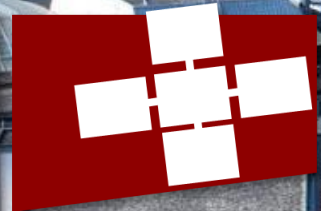


NIKOLI DRYDEN (NDRYDEN@ETHZ.CH)

Parallelism in Training Deep Neural Networks

DPHPC Guest Lecture

WITH CONTRIBUTIONS FROM TAL BEN-NUN, TORSTEN HOEFLER, DAN ALISTARH,
AND OTHERS AT SPCL, LLNL, UIUC, IST AUSTRIA, AND TOKYO TECH



Overview

- What is deep learning?
- Some deep neural networks
- Parallelizing and distributing training
- Communication for training
- Applications

Some General References

- Russell & Norvig, *Artificial Intelligence: A Modern Approach*
- Goodfellow, Bengio, & Courville, *Deep Learning*
 - Freely available online: <http://www.deeplearningbook.org/>
- Ben-Nun & Hoefler, *Demystifying Parallel and Distributed Deep Learning*
 - <https://arxiv.org/abs/1802.09941>
- Many slides adapted from Tal Ben-Nun, Torsten Hoefler, Svetlana Lazebnik, and prior talks

What is Deep Learning good for?



What is Deep Learning good for?

Digit Recognition



1989

2012

What is Deep Learning good for?

Digit Recognition



1989

2012

2017

What is Deep Learning good for?

Digit Recognition

Object Classification



1989

2012

2017

What is Deep Learning good for?

Digit Recognition



Object Classification

Image Captioning
GANs



1989

2012

2017



What is Deep Learning good for?

Digit Recognition



Object Classification

Segmentation

Image Captioning

GANs



1989

2012

2017



What is Deep Learning good for?

Digit Recognition



Object Classification

Segmentation

Image Captioning

GANs



1989

2012

2017



What is Deep Learning good for?

Digit Recognition



Object Classification

Segmentation

Image Captioning

GANs



1989

2012

2017



What is Deep Learning good for?

Digit Recognition



Object Classification

Segmentation

Image Captioning

GANs



1989

2012 2013

2017



What is Deep Learning good for?

Digit Recognition



Object Classification

Segmentation

Image Captioning

GANs



1989

2012 2013

2017



What is Deep Learning good for?

Digit Recognition



Object Classification

Segmentation

Image Captioning

GANs



1989

2012 2013

2017



What is Deep Learning good for?

Digit Recognition

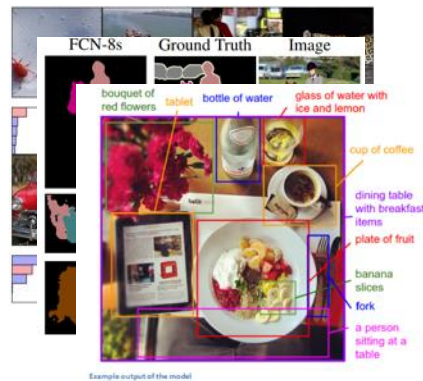


Object Classification

Segmentation

Image Captioning

GANs



1989

2012 2013 2014

2017



What is Deep Learning good for?

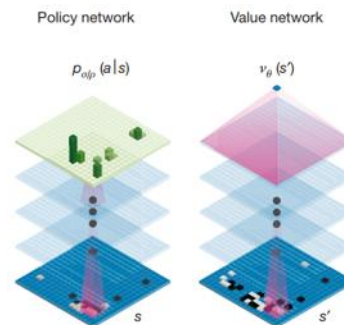
Digit Recognition



Object Classification
Segmentation
Image Captioning
GANs



Gameplay AI
Translation



1989

2012 2013 2014

2017

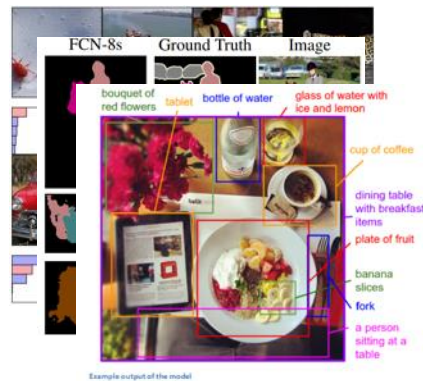


What is Deep Learning good for?

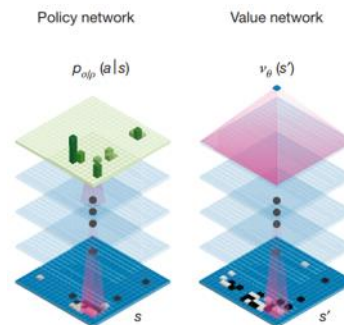
Digit Recognition



Object Classification
Segmentation
Image Captioning
GANs



Gameplay AI
Translation



1989

2012 2013 2014

2016

2017



What is Deep Learning good for?

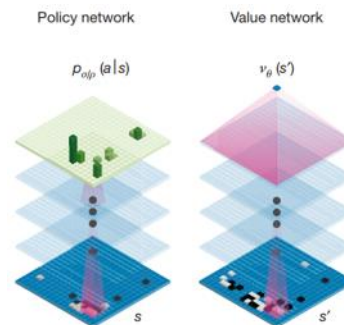
Digit Recognition



Object Classification
Segmentation
Image Captioning
GANs



Gameplay AI
Translation



1989

2012 2013 2014

2016

2017



What is Deep Learning good for?

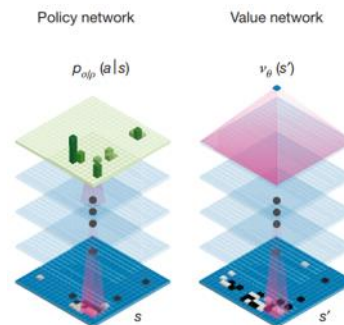
Digit Recognition



Object Classification
Segmentation
Image Captioning
GANs



Gameplay AI
Translation



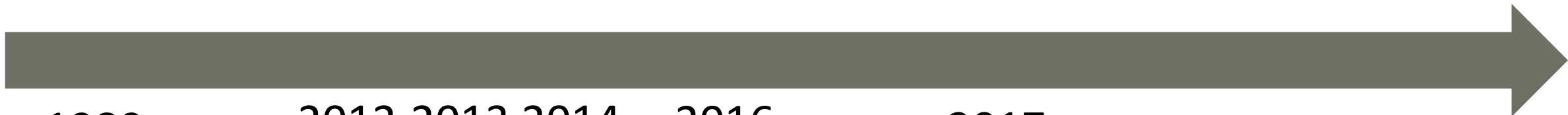
Neural Computers

1989

2012 2013 2014

2016

2017

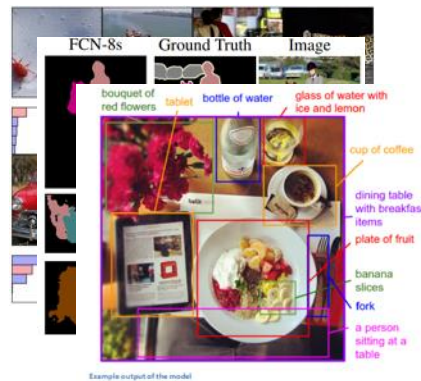


What is Deep Learning good for?

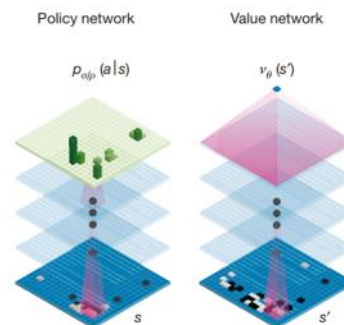
Digit Recognition



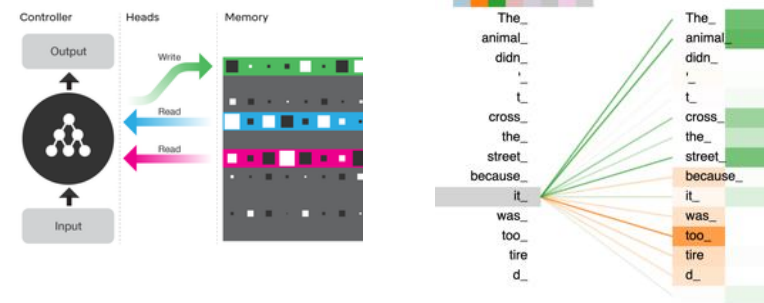
Object Classification
Segmentation
Image Captioning
GANs



Gameplay AI
Translation



Neural Computers



1989

2012 2013 2014

2016

2017

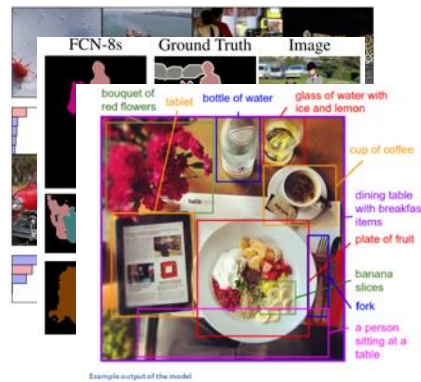


What is Deep Learning good for?

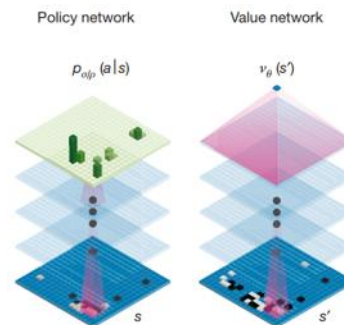
Digit Recognition



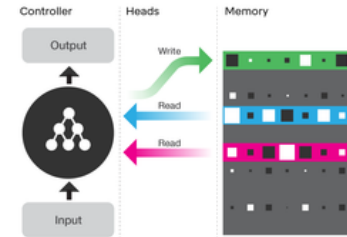
Object Classification
Segmentation
Image Captioning
GANs



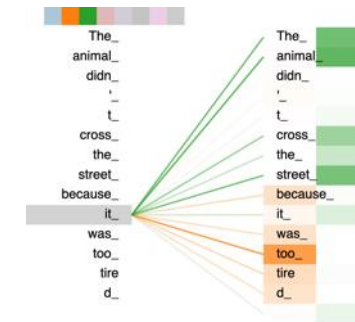
Gameplay AI
Translation



Neural Computers



Language Models



1989

2012 2013 2014

2016

2017

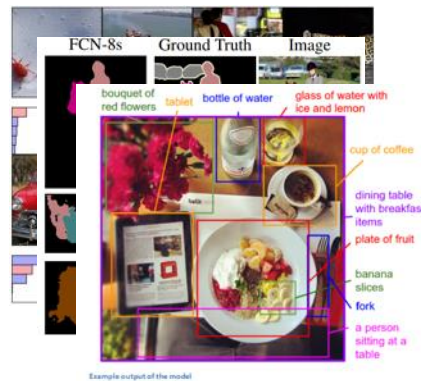


What is Deep Learning good for?

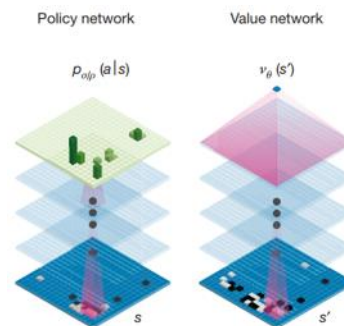
Digit Recognition



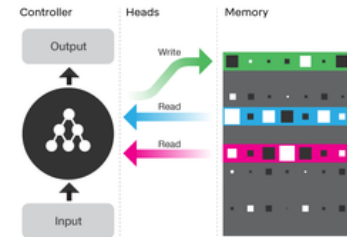
Object Classification
Segmentation
Image Captioning
GANs



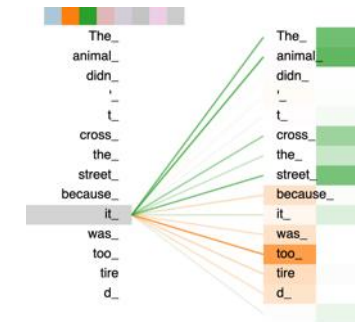
Gameplay AI
Translation



Neural Computers



Language Models



1989

2012 2013 2014

2016

2017



What is Deep Learning good for?

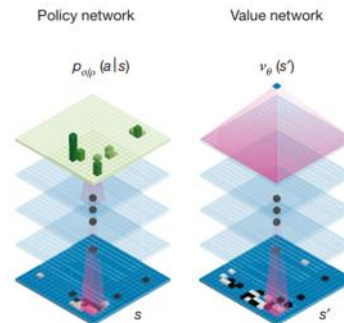
Digit Recognition



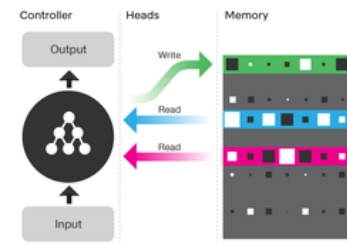
Object Classification
Segmentation
Image Captioning
GANs



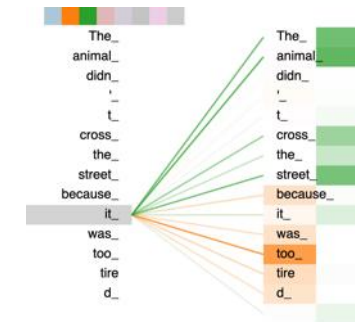
Gameplay AI
Translation



Neural Computers



Language Models



1989

2012 2013 2014

2016

2017

2019

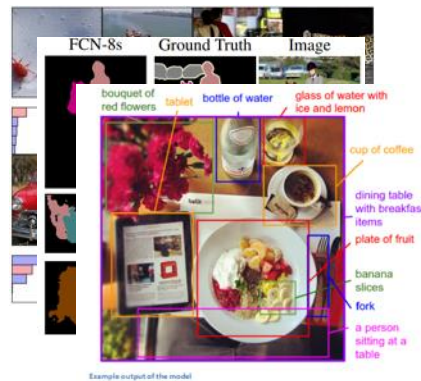


What is Deep Learning good for?

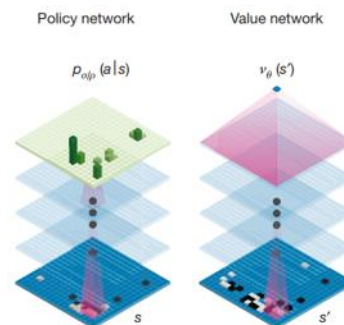
Digit Recognition



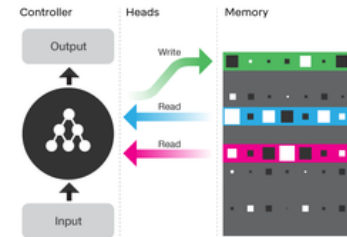
Object Classification
Segmentation
Image Captioning
GANs



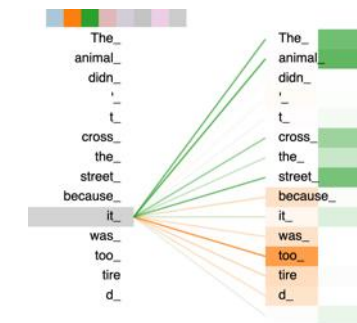
Gameplay AI
Translation



Neural Computers



Language Models



Towards
Real Physics
RTS

1989

2012 2013 2014

2016

2017

2018

2019

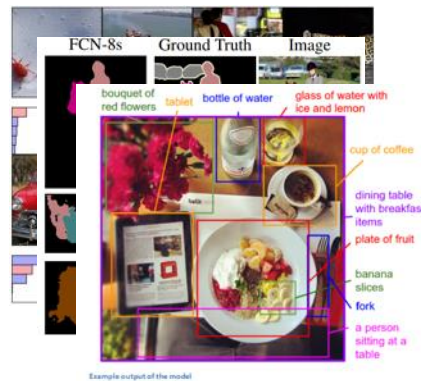


What is Deep Learning good for?

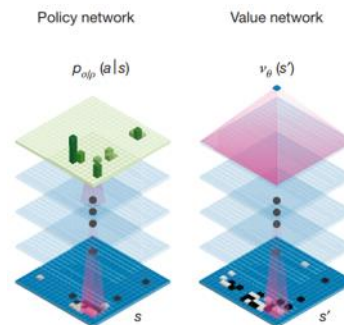
Digit Recognition



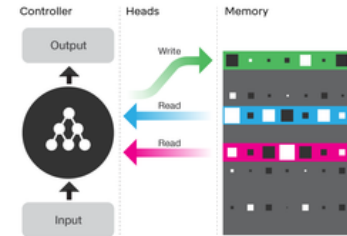
Object Classification
Segmentation
Image Captioning
GANs



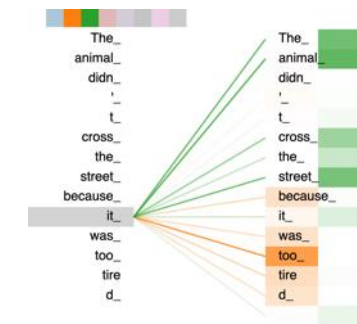
Gameplay AI
Translation



Neural Computers



Language Models



Towards
Real Physics
RTS



1989

2012 2013 2014

2016

2017

2018

2019



What is Deep Learning good for?

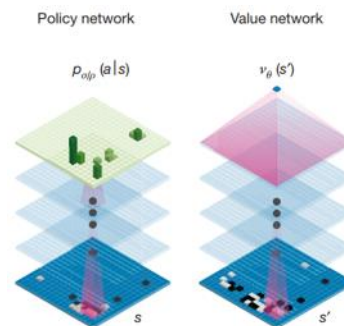
Digit Recognition



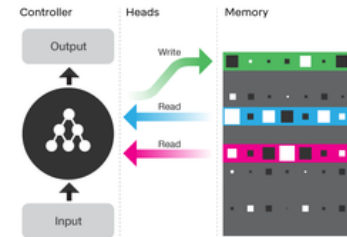
Object Classification
Segmentation
Image Captioning
GANs



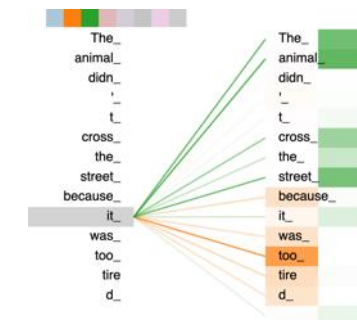
Gameplay AI
Translation



Neural Computers



Language Models



Towards
Real Physics
RTS



1989

2012 2013 2014

2016

2017

2018

2019



What is Deep Learning good for?

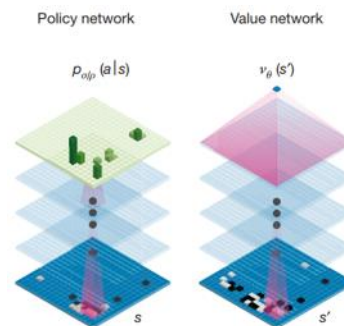
Digit Recognition



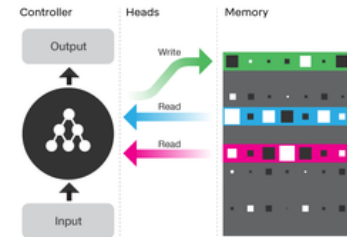
Object Classification
Segmentation
Image Captioning
GANs



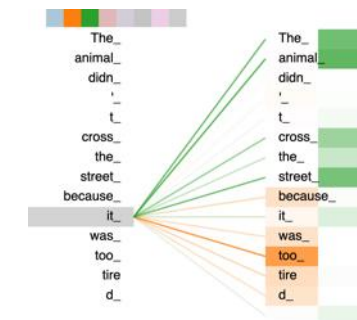
Gameplay AI
Translation



Neural Computers



Language Models



Towards
Real Physics
RTS



1989

2012 2013 2014

2016

2017

2018

2019

What is Deep Learning good for?

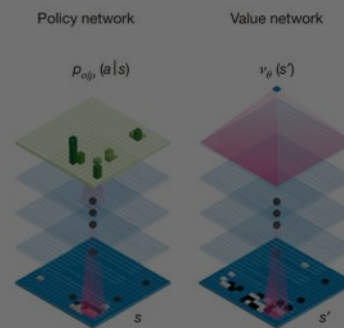
Digit Recognition



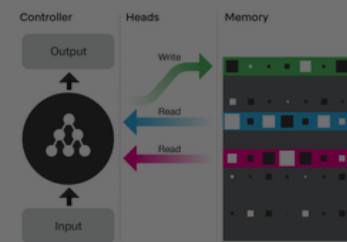
Object Classification
Segmentation
Image Captioning
GANs



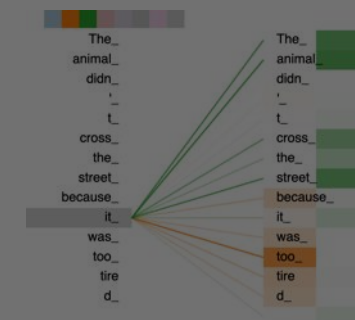
Gameplay AI
Translation



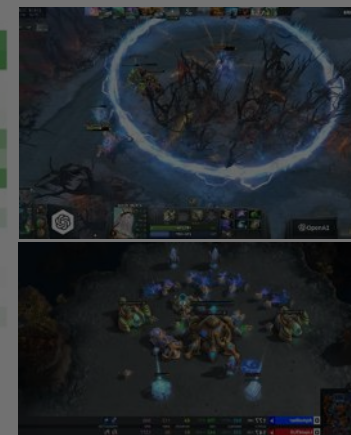
Neural Computers



Language Models



Towards
Real Physics
RTS



1989

2012 2013 2014

2016

2017

2018

2019

What is Deep Learning good for?

Digit Recognition

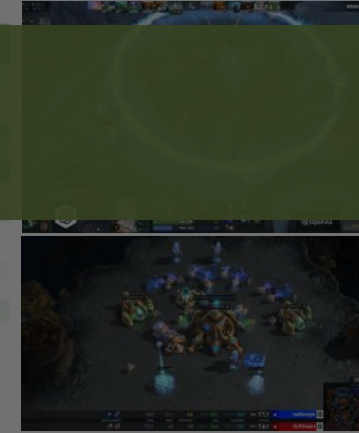
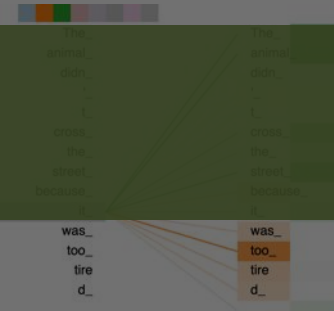
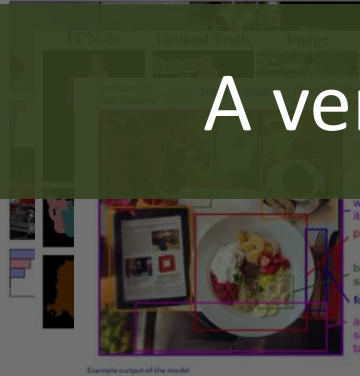
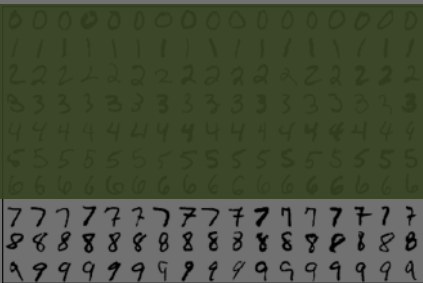
Object Classification
Segmentation
Image Captioning
GANs

Gameplay AI
Translation

Neural Computers

Language Models

Towards
Real Physics
RTS



A very active area of research!

1989

2012 2013 2014

2016

2017

2018

2019

What is Deep Learning good for?

Digit Recognition

Object Classification
 Segmentation
 Image Captioning
 GANs

Gameplay AI
 Translation

Neural Computers

Language Models

Towards
 Real Physics
 RTS

A very active area of research!

Subject	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019
cs.AI	380	479	789	1082	1768	1028	1106	1938	2820	4263	4371
cs.CV	148	286	385	577	852	1349	2262	3631	5704	8599	10353
cs.LG	231	333	469	1222	1418	1742	2485	3564	5225	10472	17267
stat.ML	164	256	439	1131	1203	1360	1827	2628	4021	8376	11551
Total	923	1354	2082	4012	5241	5479	7680	11761	17770	31710	43542

1989

2019

What is Deep Learning

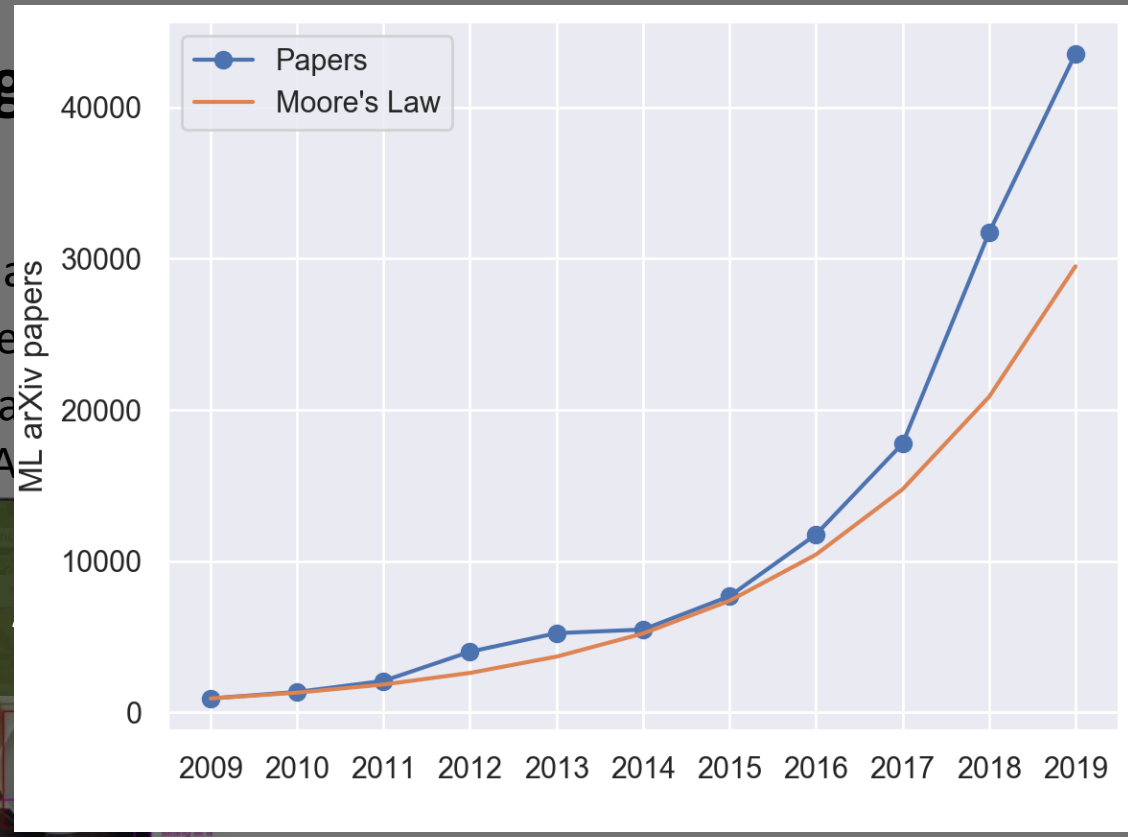
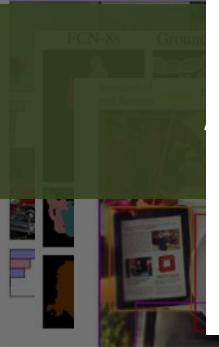
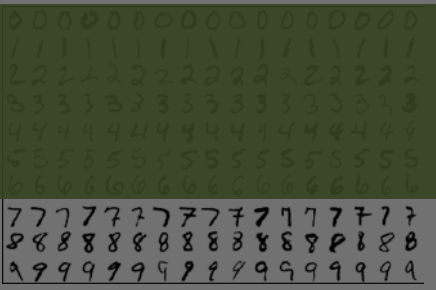
Digit Recognition

Object Classification

Segmentation

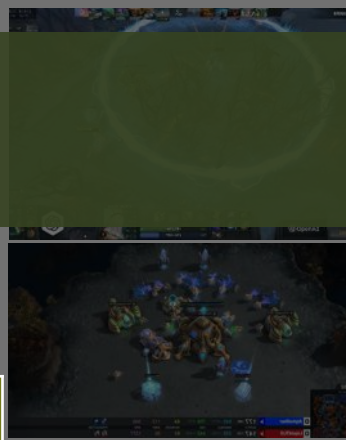
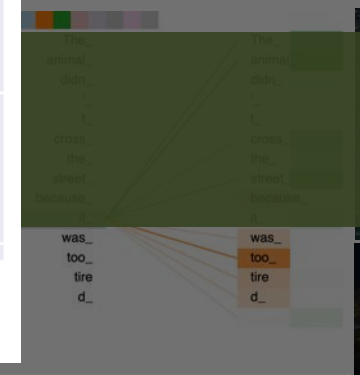
Image Captioning

Generative Adversarial Networks



Language Models

Towards Real Physics
RTS



Subject	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019
cs.AI	380	479	789	1082	1768	1028	1106	1938	2820	4263	4371
cs.CV	148	286	385	577	852	1349	2262	3631	5704	8599	10353
cs.LG	231	333	469	1222	1418	1742	2485	3564	5225	10472	17267
stat.ML	164	256	439	1131	1203	1360	1827	2628	4021	8376	11551
Total	923	1354	2082	4012	5241	5479	7680	11761	17770	31710	43542

1989

2019

What is Deep Learning

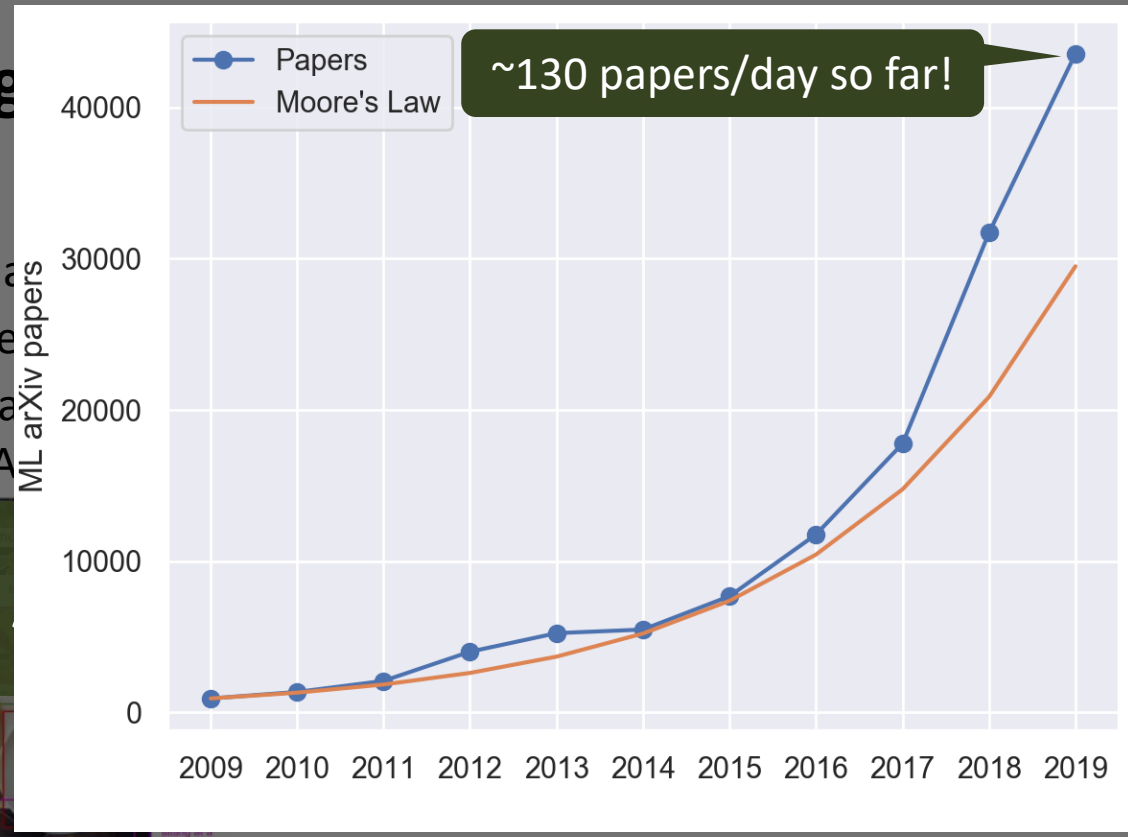
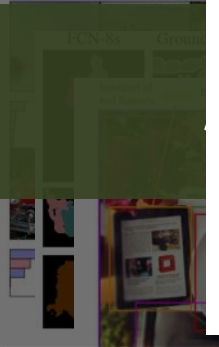
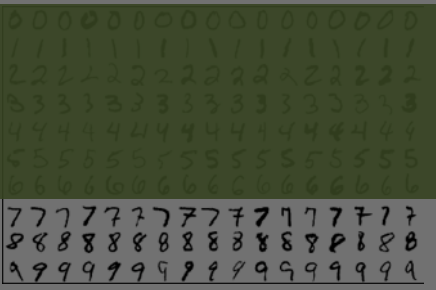
Digit Recognition

Object Classification

Segmentation

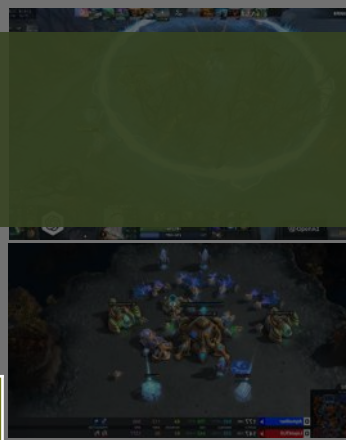
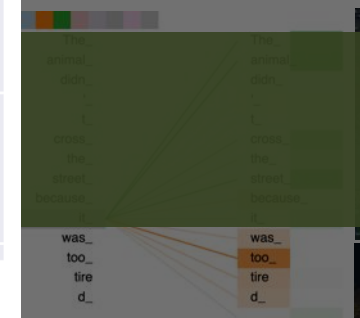
Image Captioning

Generative Adversarial Networks



Language Models

Towards Real Physics
RTS

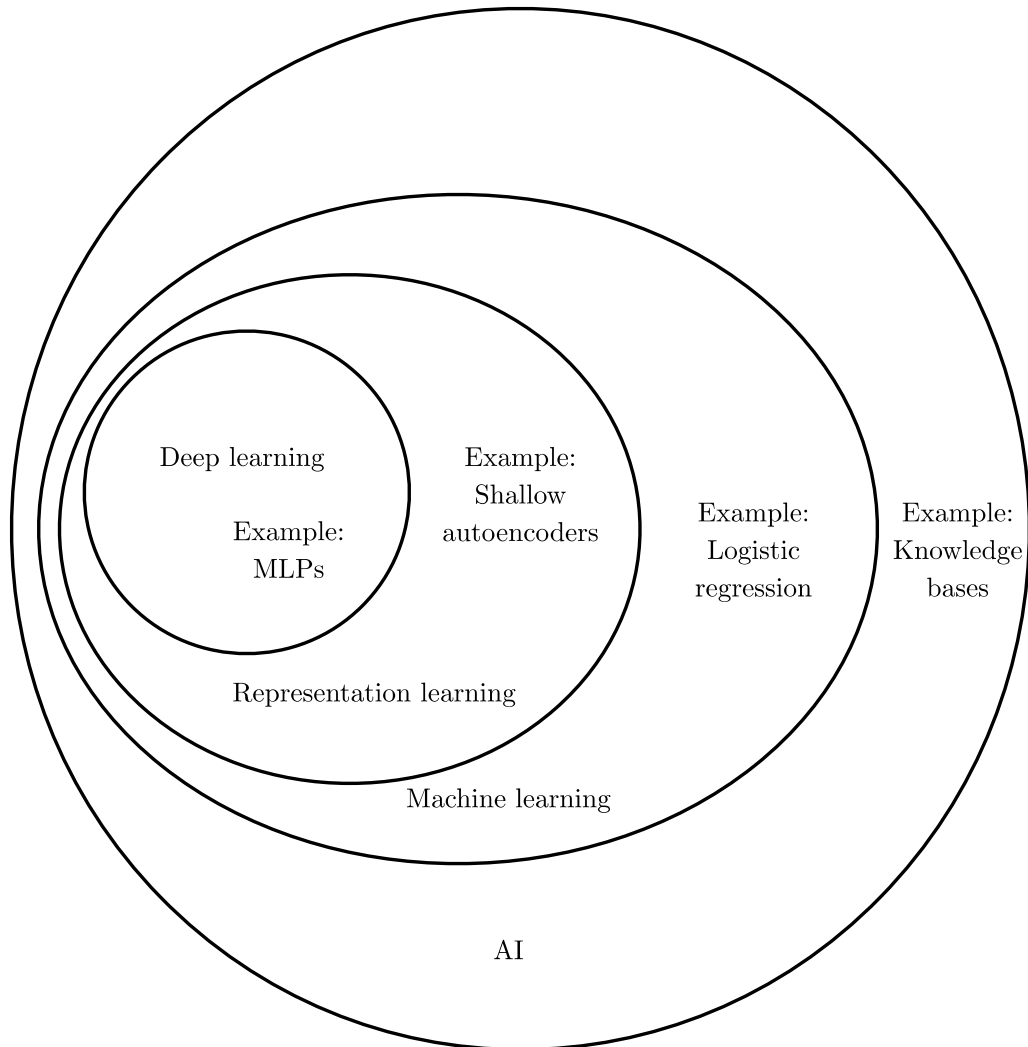


Subject	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019
cs.AI	380	479	789	1082	1768	1028	1106	1938	2820	4263	4371
cs.CV	148	286	385	577	852	1349	2262	3631	5704	8599	10353
cs.LG	231	333	469	1222	1418	1742	2485	3564	5225	10472	17267
stat.ML	164	256	439	1131	1203	1360	1827	2628	4021	8376	11551
Total	923	1354	2082	4012	5241	5479	7680	11761	17770	31710	43542

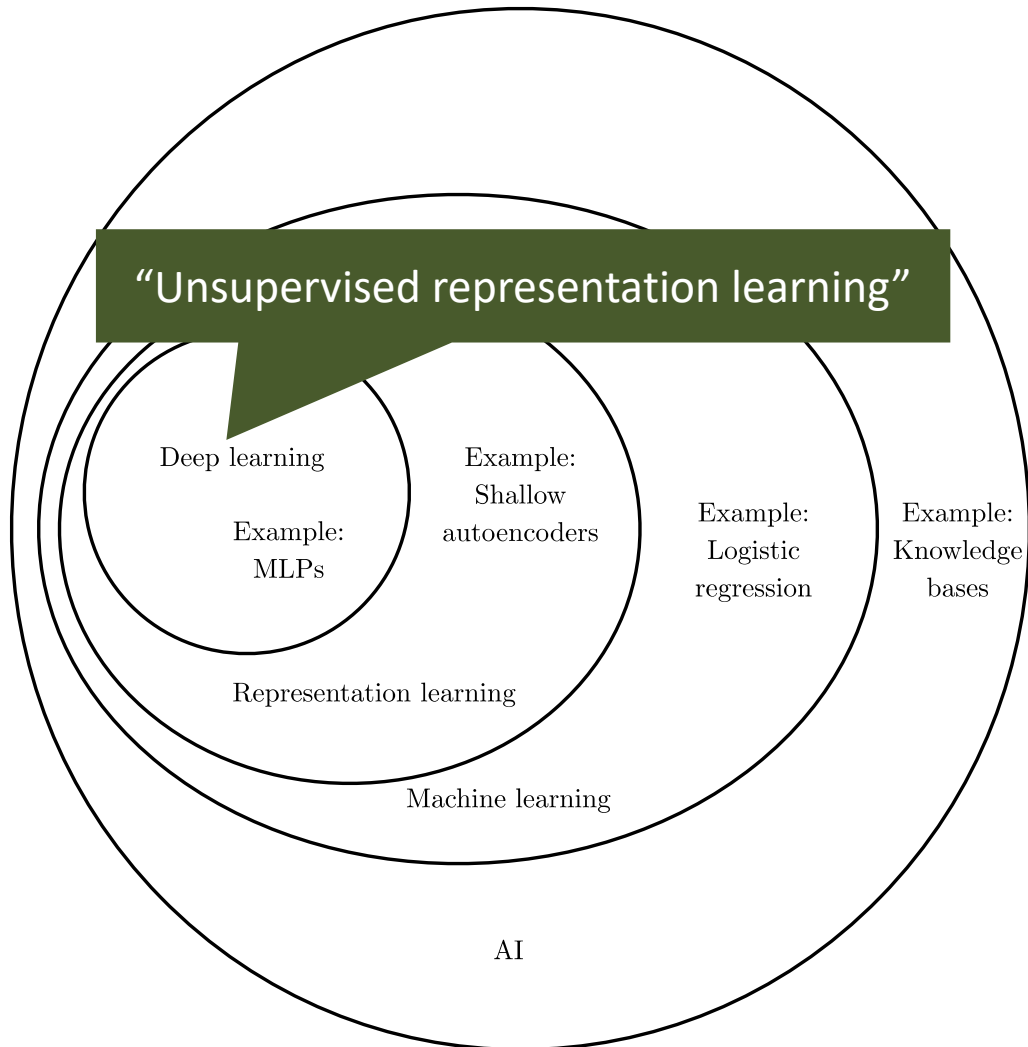
1989

2019

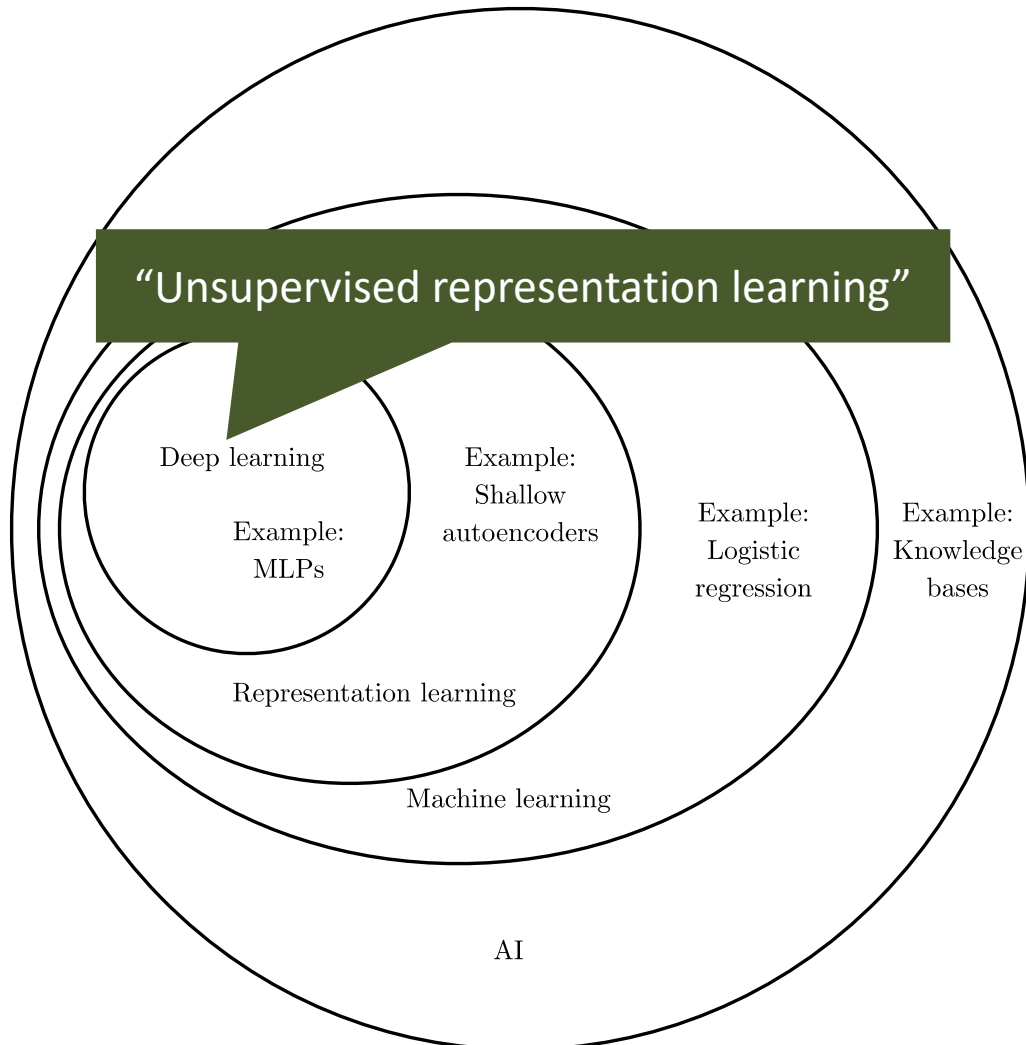
Classes of AI Problems



Classes of AI Problems



Classes of AI Problems



- **Supervised learning**

- Learn mapping from labeled inputs

$$\operatorname{argmin}_{f \in \mathcal{H}} \mathbb{E}_{x, y \sim \mathcal{D}} [\ell(f(x), y)]$$

- **Unsupervised learning**

- Learn patterns in inputs

$$\operatorname{argmin}_{f \in \mathcal{H}} \mathbb{E}_{x \sim \mathcal{D}} [\ell(f(x))]$$

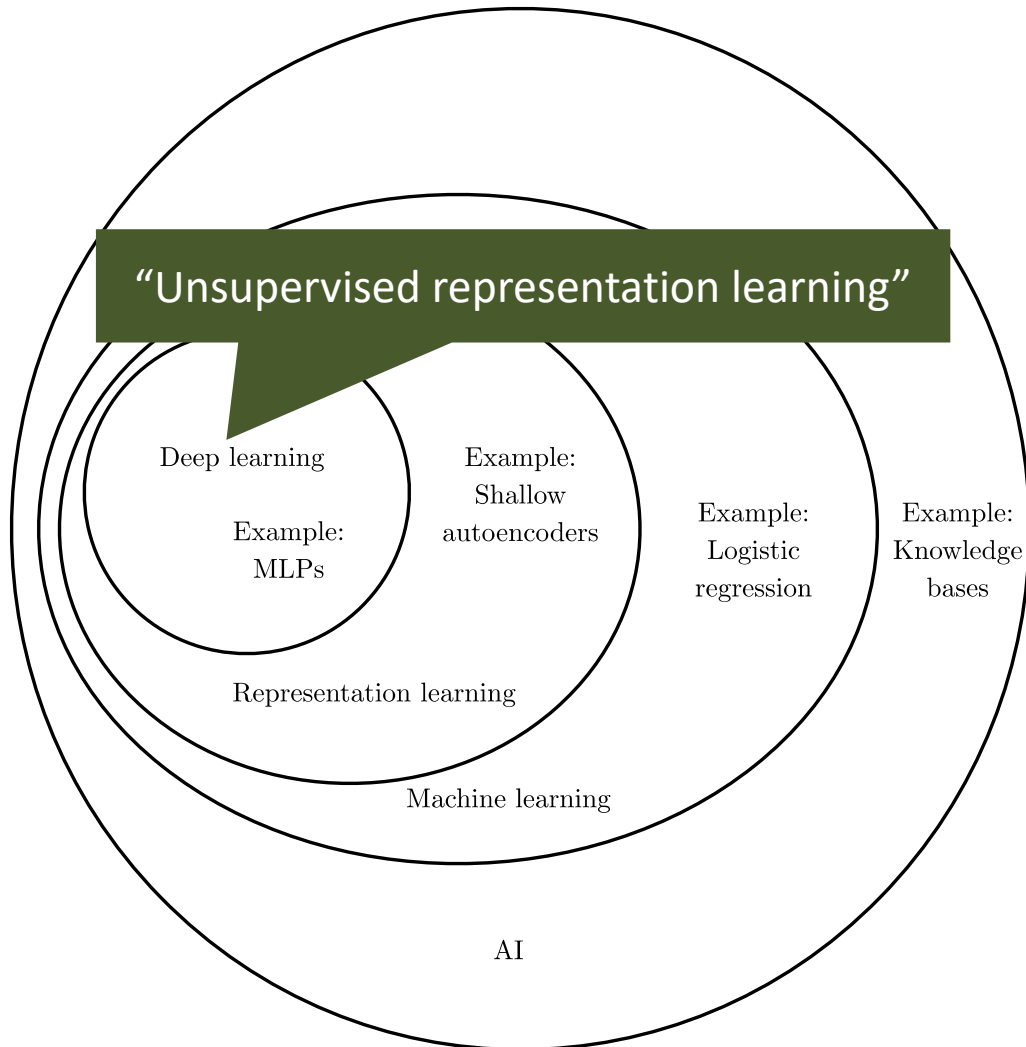
- **Reinforcement learning**

- Learn policy to maximize reward

$$\operatorname{argmax}_{\pi \in \mathcal{H}} \mathbb{E}_{O \sim \Omega} [R(\pi, O)]$$

- **Many others...**

Classes of AI Problems



- **Supervised learning**

- Learn mapping from labeled inputs

$$\operatorname{argmin}_{f \in \mathcal{H}} \mathbb{E}_{x, y \sim \mathcal{D}} [\ell(f(x), y)]$$

- **Unsupervised learning**

- Learn patterns in inputs

$$\operatorname{argmin}_{f \in \mathcal{H}} \mathbb{E}_{x \sim \mathcal{D}} [\ell(f(x))]$$

- **Reinforcement learning**

- Learn policy to maximize reward

$$\operatorname{argmax}_{\pi \in \mathcal{H}} \mathbb{E}_{O \sim \Omega} [R(\pi, O)]$$

- **Many others...**

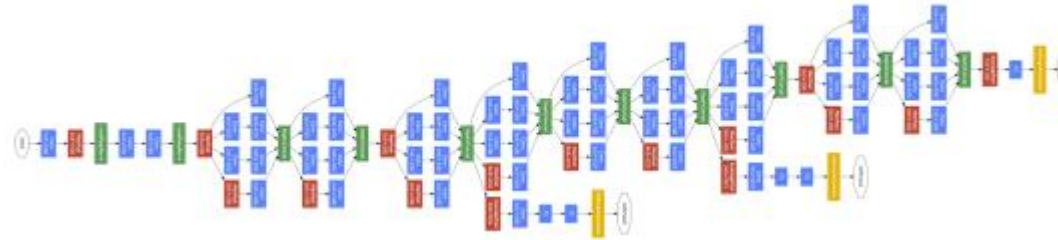
Deep

A brief theory of supervised deep learning (mini-batch SGD)

A brief theory of supervised deep learning (mini-batch SGD)



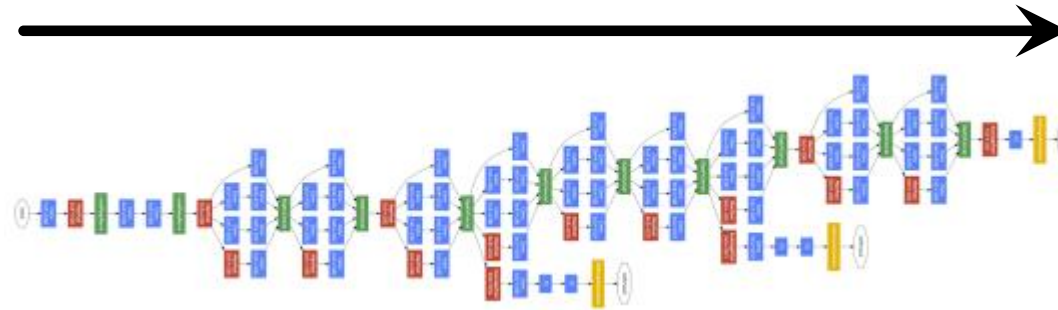
A brief theory of supervised deep learning (mini-batch SGD)



A brief theory of supervised deep learning (mini-batch SGD)



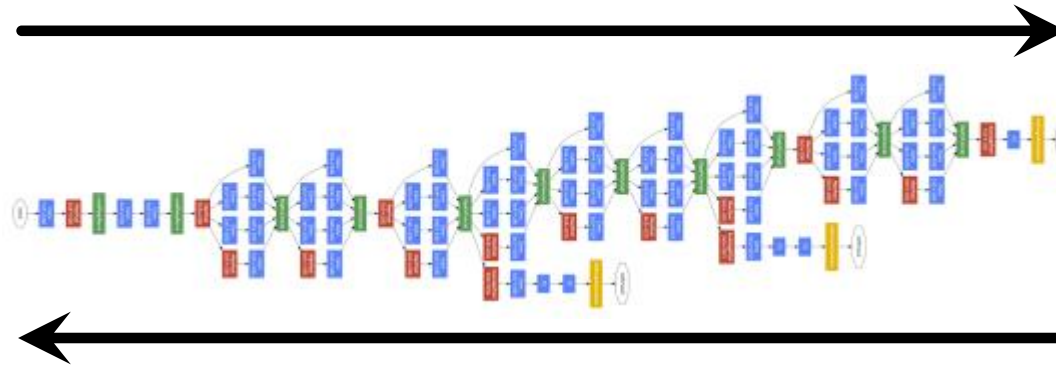
labeled samples $x \in X \subset \mathcal{D}$



A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$

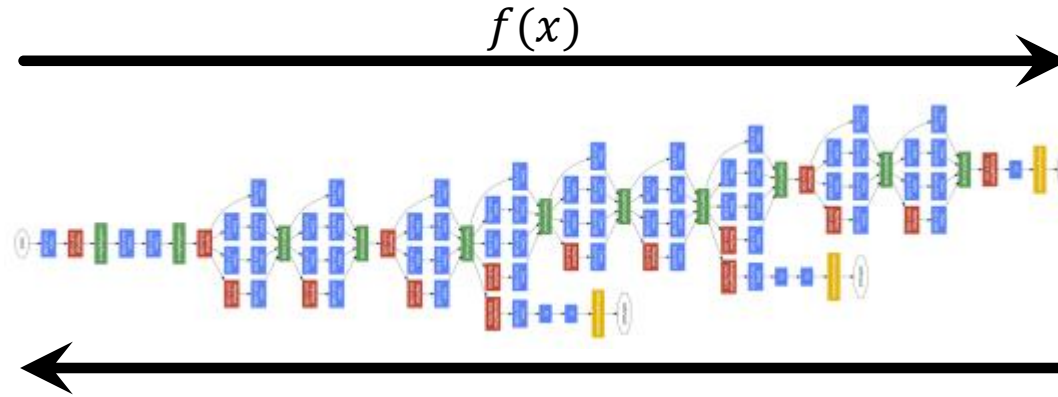


Cat	1.00
Dog	0.00
Airplane	0.00
Horse	0.00
Banana	0.00
Truck	0.00

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



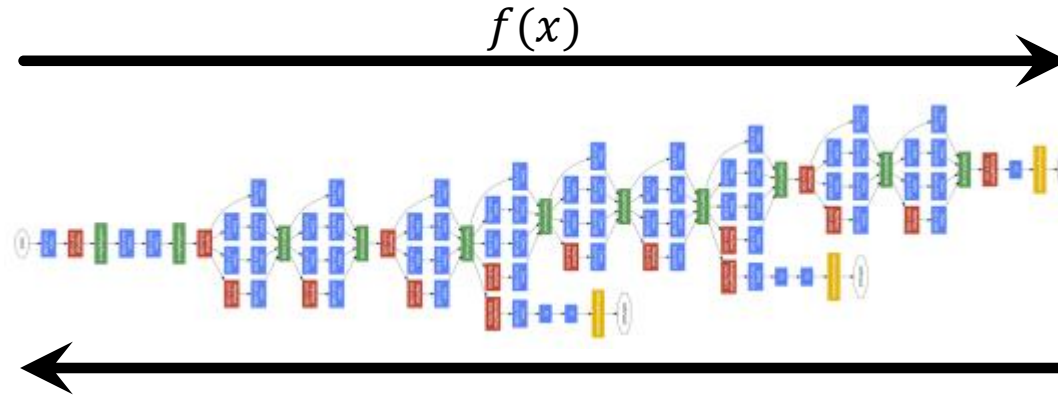
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

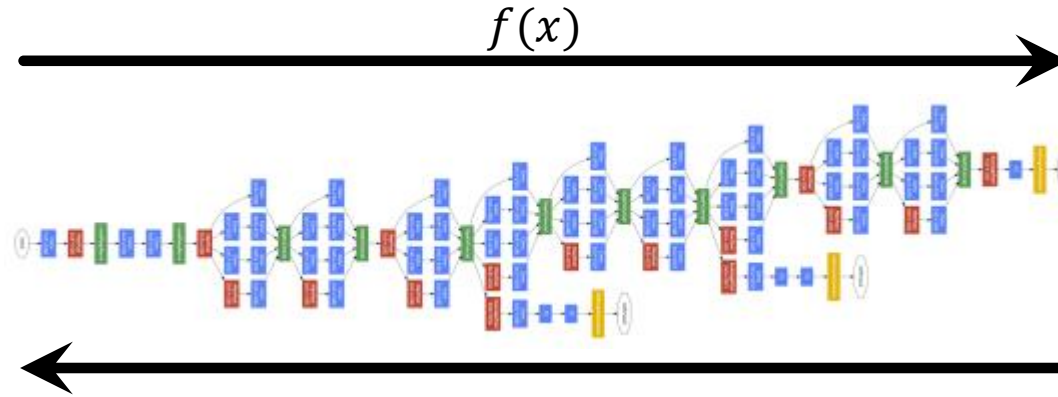
label domain Y

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$

$$f(x): X \rightarrow Y$$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

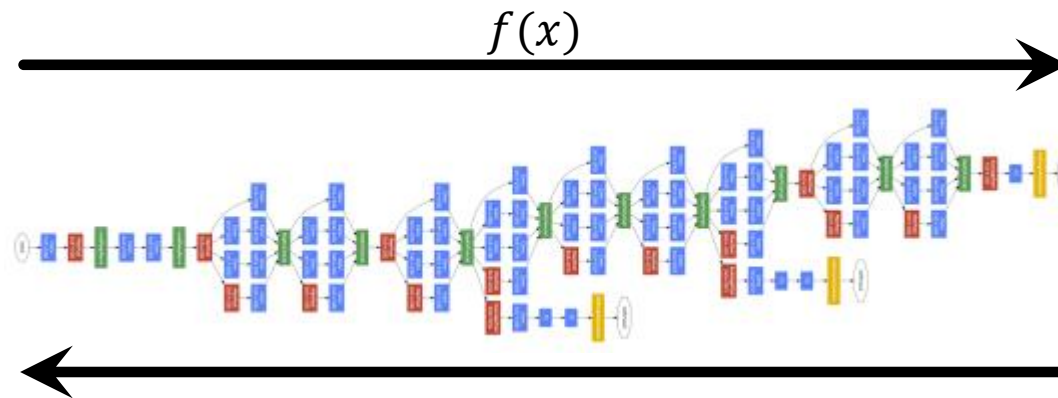
true label $l(x)$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$

$$f(x): X \rightarrow Y$$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

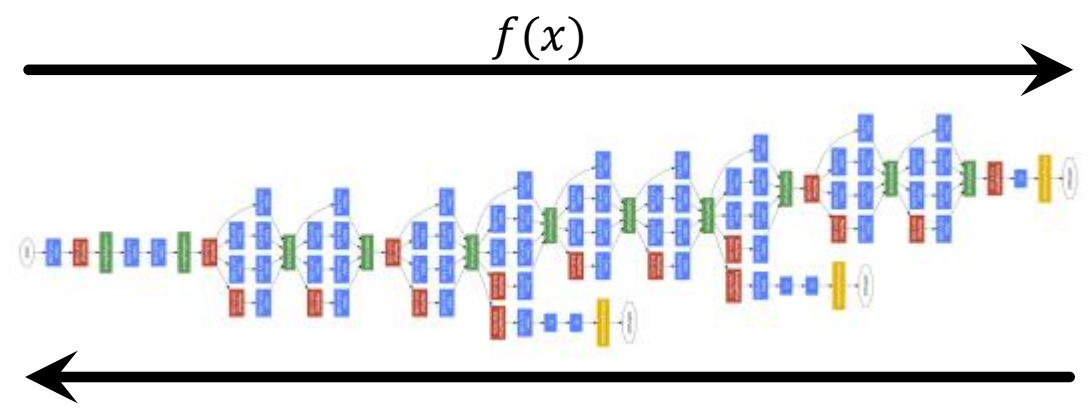
label domain Y

true label $l(x)$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

true label $l(x)$

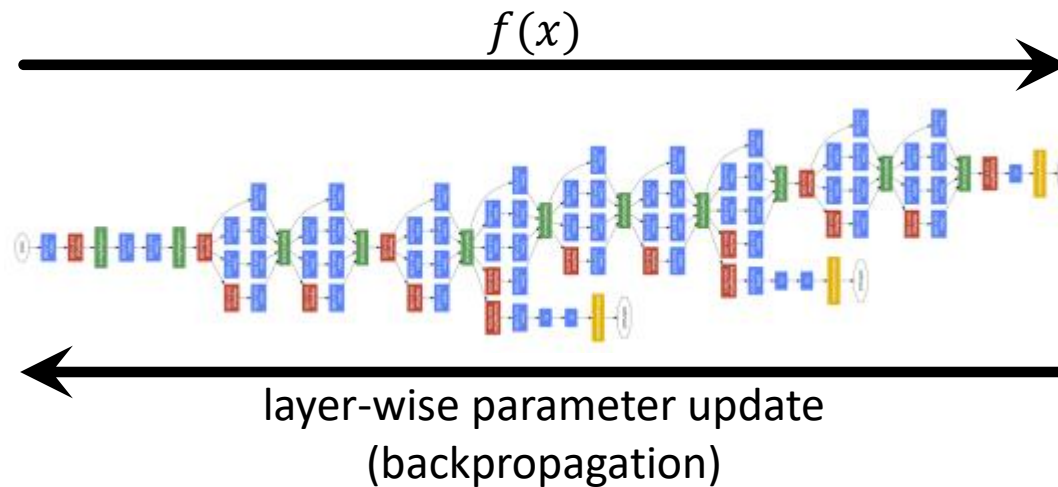
$$f(x): X \rightarrow Y$$

network structure
(fixed)

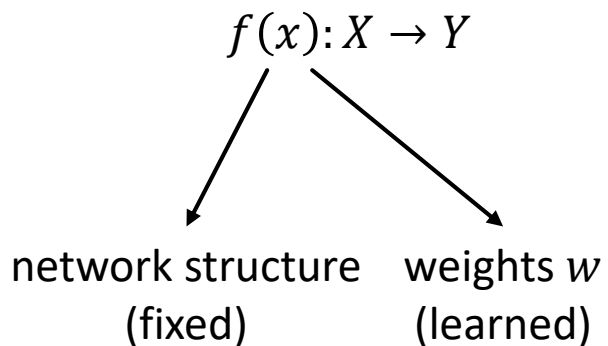
A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



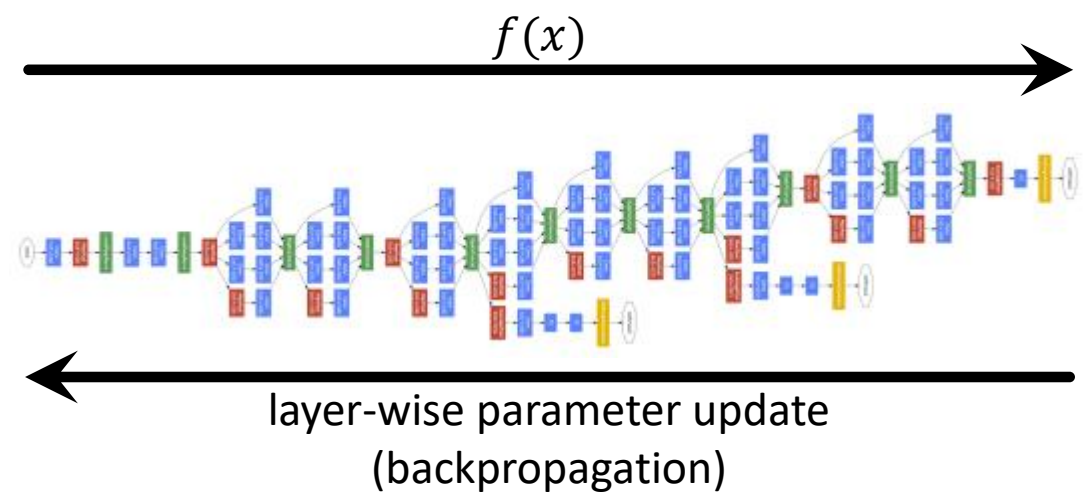
label domain Y	output	true label $l(x)$
Cat	0.54	1.00
Dog	0.28	0.00
Airplane	0.07	0.00
Horse	0.33	0.00
Banana	0.02	0.00
Truck	0.02	0.00



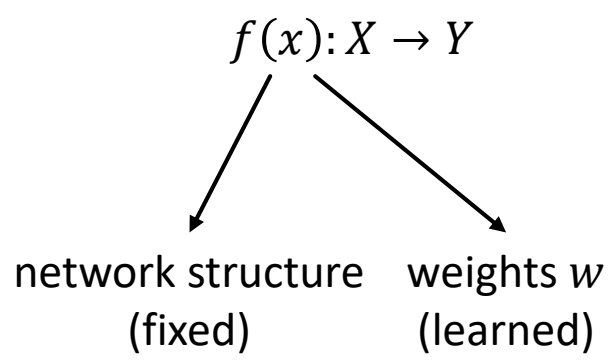
A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



label domain Y	output	true label $l(x)$
Cat	0.54	1.00
Dog	0.28	0.00
Airplane	0.07	0.00
Horse	0.33	0.00
Banana	0.02	0.00
Truck	0.02	0.00

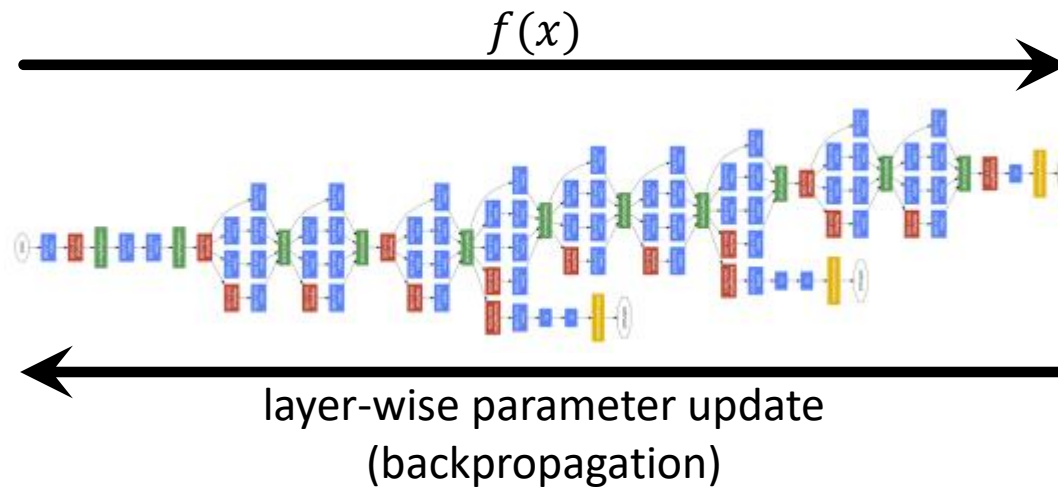


$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

true label $l(x)$

$$f(x): X \rightarrow Y$$

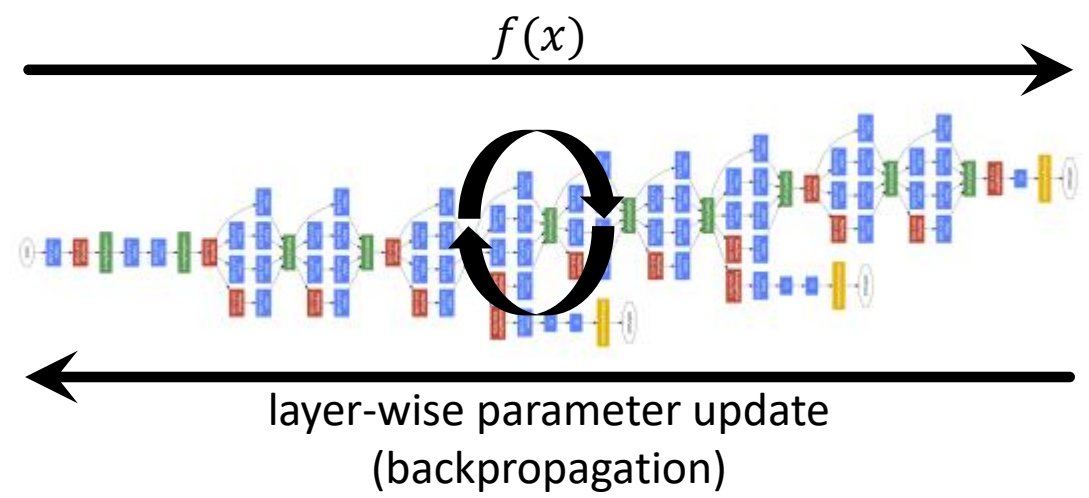
network structure (fixed) weights w (learned)

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

true label $l(x)$

$$f(x): X \rightarrow Y$$

network structure (fixed) weights w (learned)

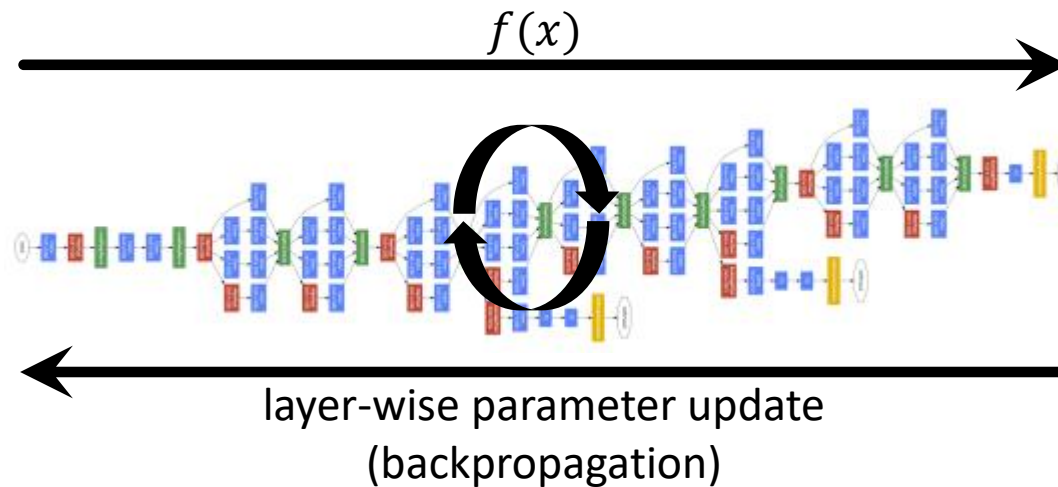
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

true label $l(x)$

$$f(x): X \rightarrow Y$$

network structure
(fixed)

weights w
(learned)

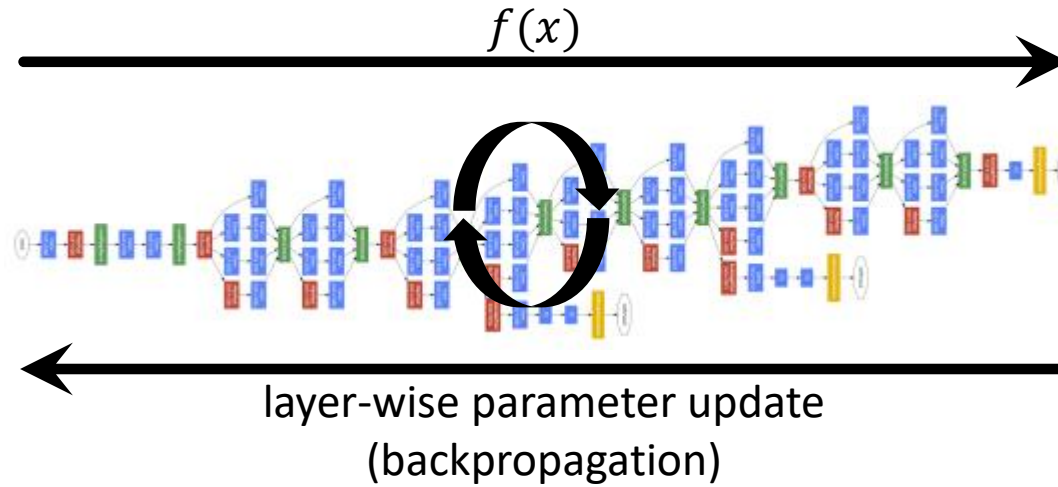
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

A brief theory of supervised deep learning (mini-batch SGD)



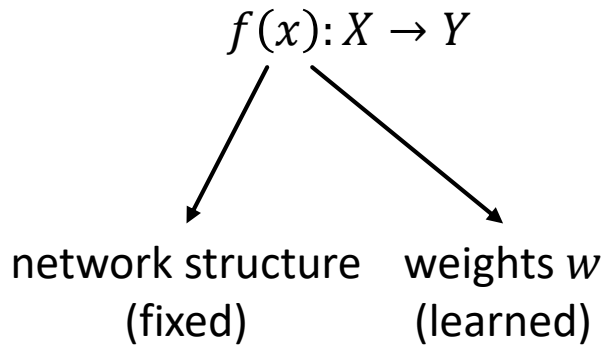
labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

true label $l(x)$



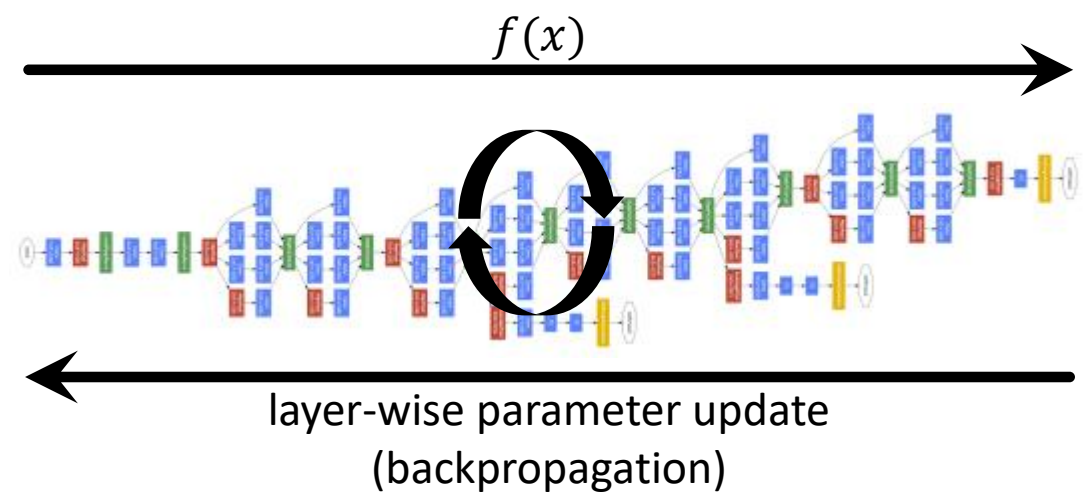
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

A brief theory of supervised deep learning (mini-batch SGD)

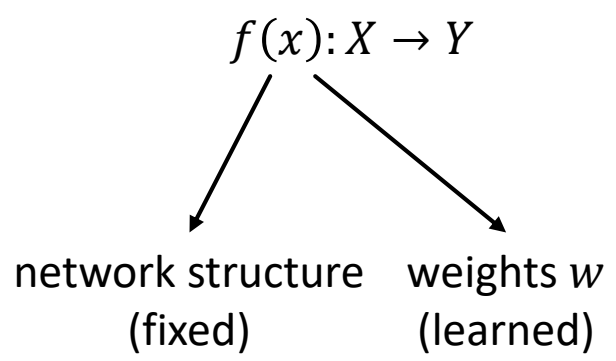


labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y true label $l(x)$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

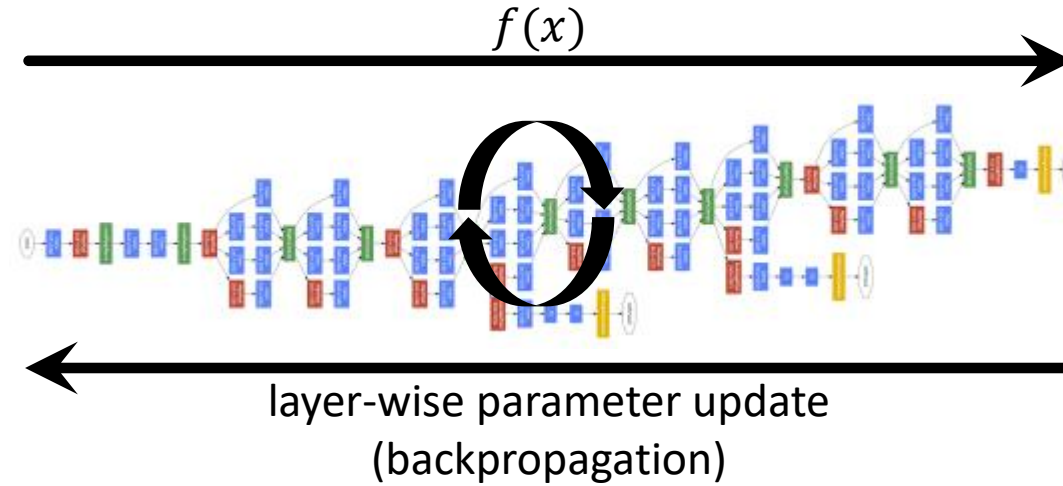
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

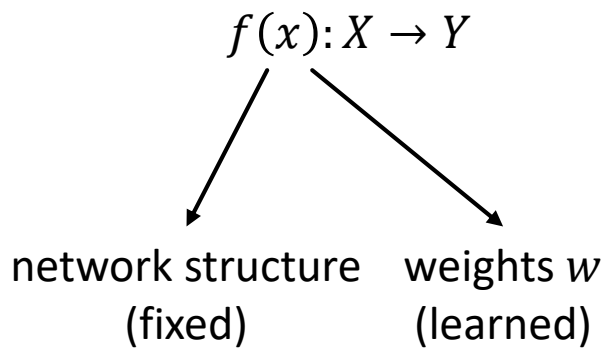
A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



label domain Y	predicted probability	true label $l(x)$
Cat	0.54	1.00
Dog	0.28	0.00
Airplane	0.07	0.00
Horse	0.33	0.00
Banana	0.02	0.00
Truck	0.02	0.00



$$f(x) = f_n(f_{n-1}(f_{n-2}(\dots f_1(x) \dots)))$$

convolution 1

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

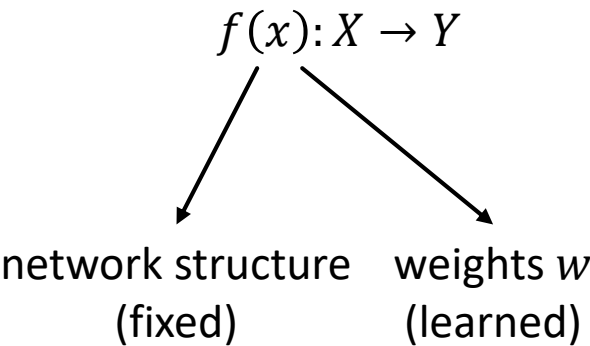
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

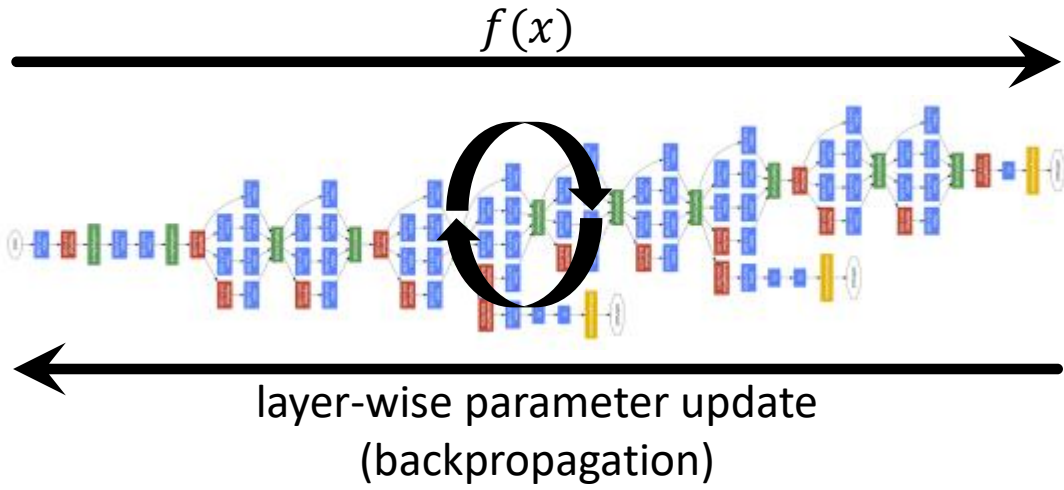
A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$



$$f(x) = f_n \left(f_{n-1} \left(f_{n-2} \left(\dots f_1(x) \dots \right) \right) \right)$$

convolution 1

convolution 2

label domain Y		true label $l(x)$	
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

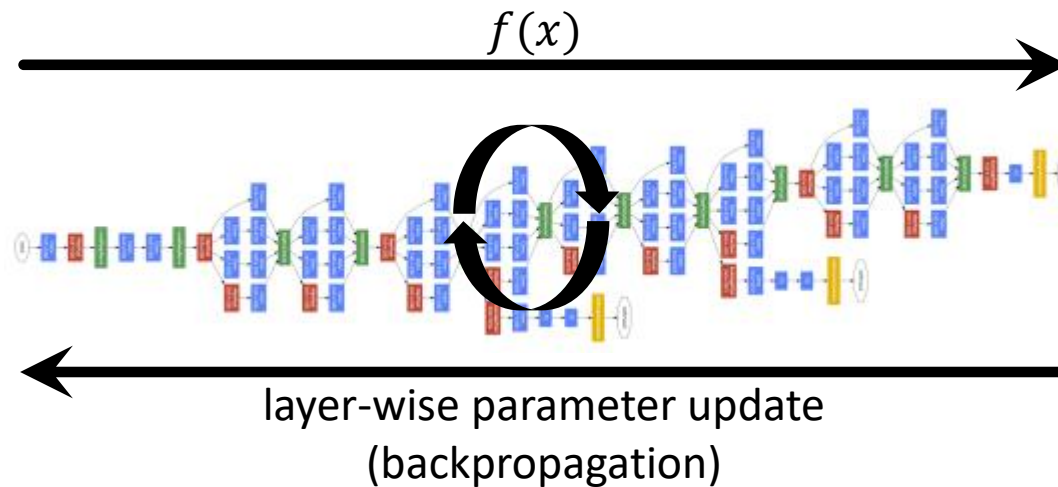
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



label domain Y	output	true label $l(x)$	output
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

true label $l(x)$

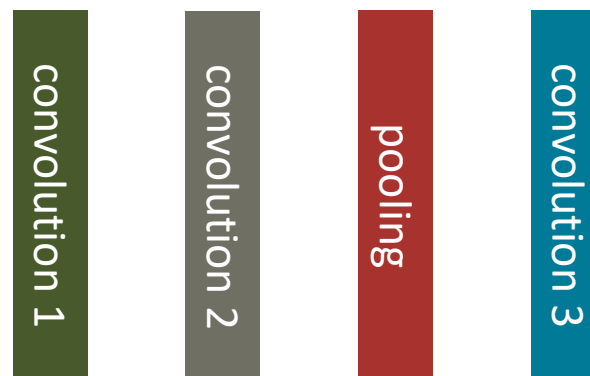
$$f(x): X \rightarrow Y$$

network structure
(fixed)

weights w
(learned)

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$$f(x) = f_n \left(f_{n-1} \left(f_{n-2} \left(\dots f_1(x) \dots \right) \right) \right)$$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

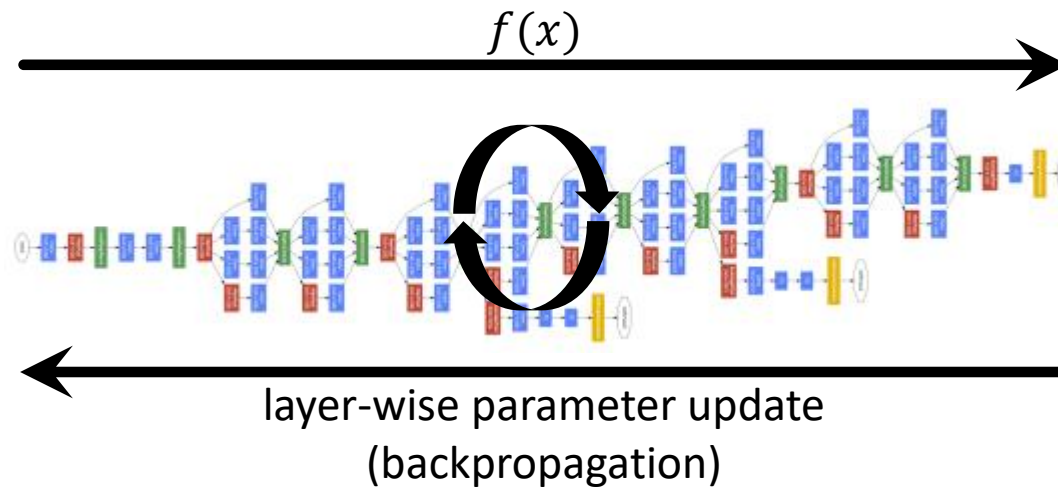
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

true label $l(x)$

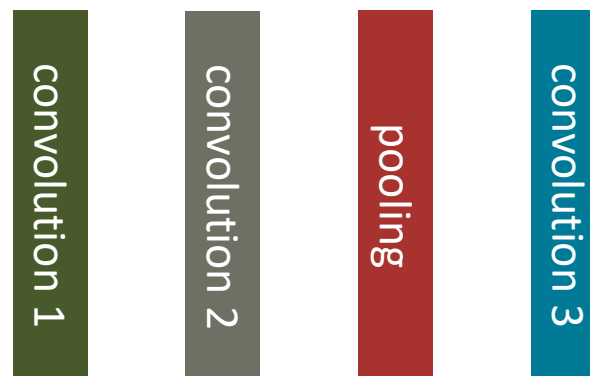
$$f(x): X \rightarrow Y$$

network structure
(fixed)

weights w
(learned)

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$$f(x) = f_n \left(f_{n-1} \left(f_{n-2} \left(\dots f_1(x) \dots \right) \right) \right)$$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

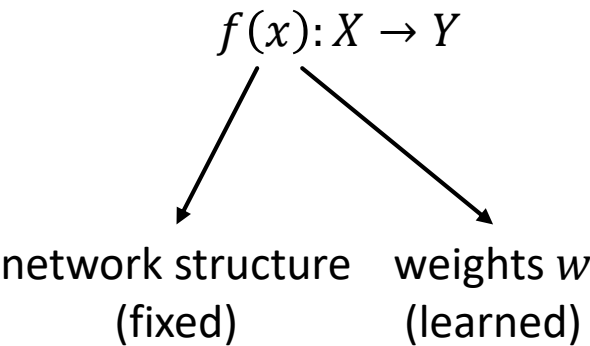
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

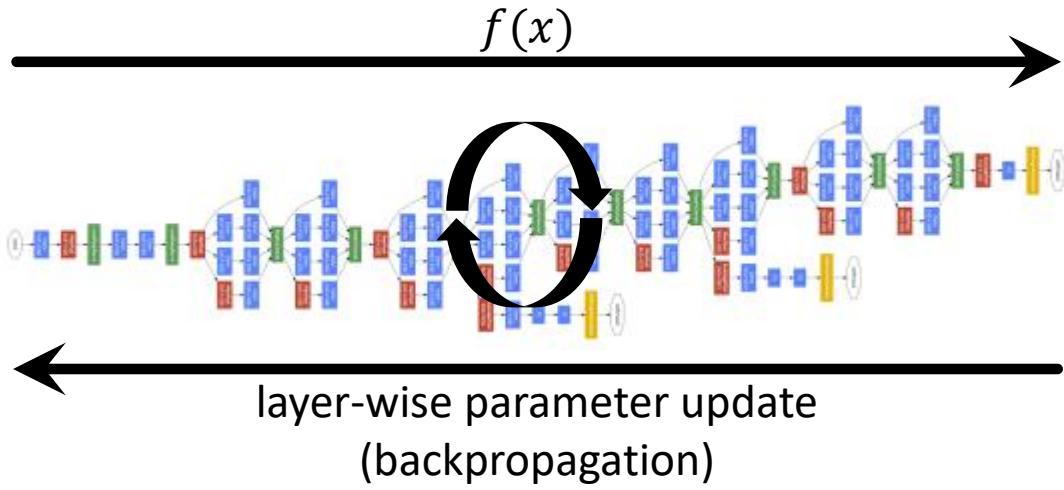
A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$



$$f(x) = f_n \left(f_{n-1} \left(f_{n-2} \left(\dots f_1(x) \dots \right) \right) \right)$$



label domain Y	prediction	true label $l(x)$
Cat	0.54	1.00
Dog	0.28	0.00
Airplane	0.07	0.00
Horse	0.33	0.00
Banana	0.02	0.00
Truck	0.02	0.00

$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

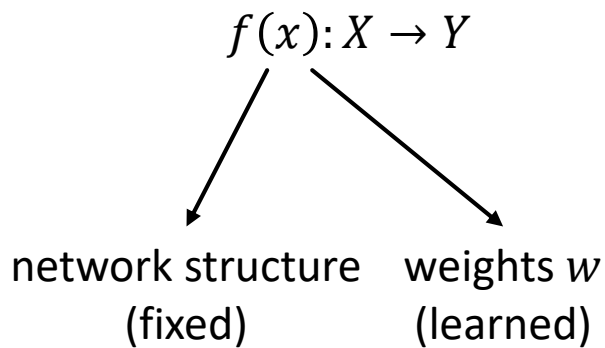
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

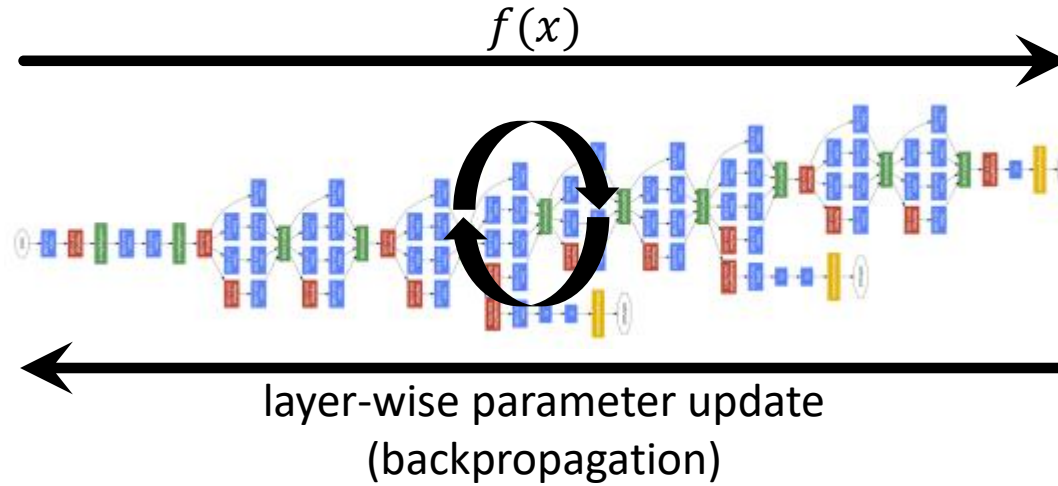
A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

true label $l(x)$

$$f(x) = f_n(f_{n-1}(f_{n-2}(\dots f_1(x) \dots)))$$

$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

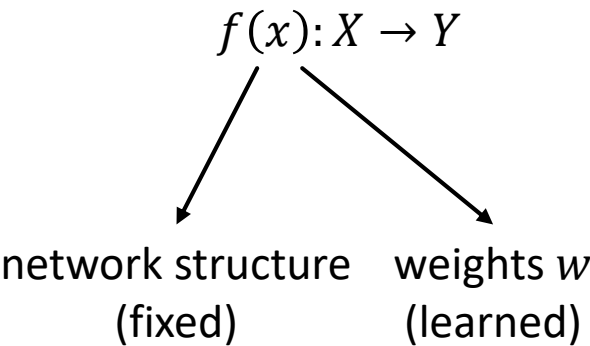
$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$



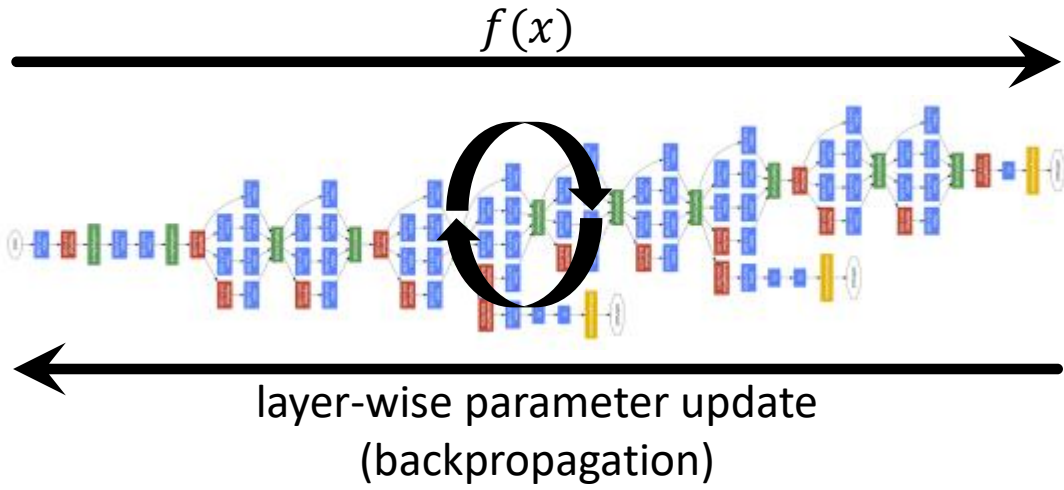
A brief theory of supervised deep learning (mini-batch SGD)



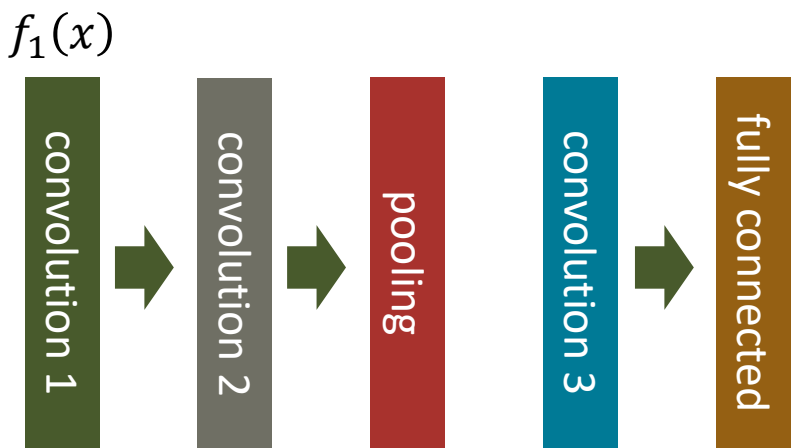
labeled samples $x \in X \subset \mathcal{D}$



$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$



$$f(x) = f_n \left(f_{n-1} \left(f_{n-2} \left(\dots f_1(x) \dots \right) \right) \right)$$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y true label $l(x)$

$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

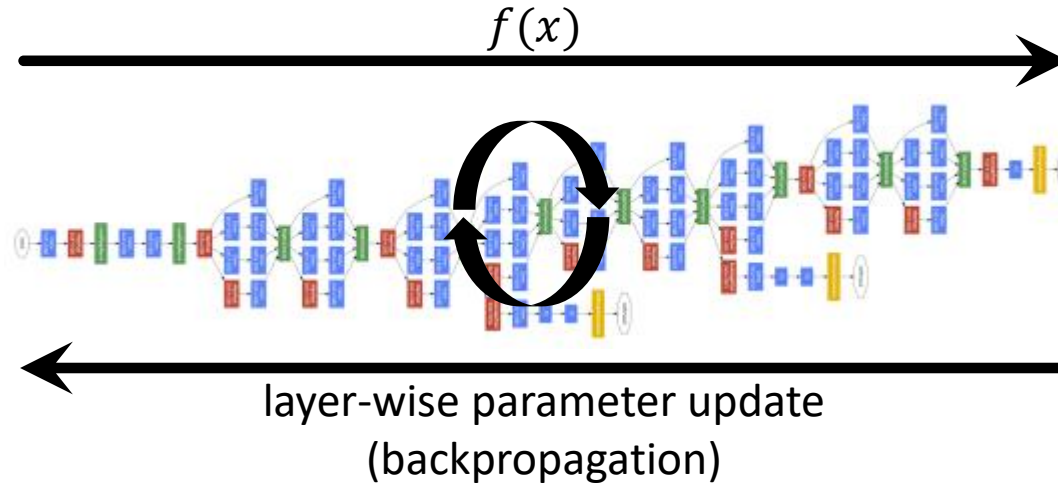
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

true label $l(x)$

$$f(x): X \rightarrow Y$$

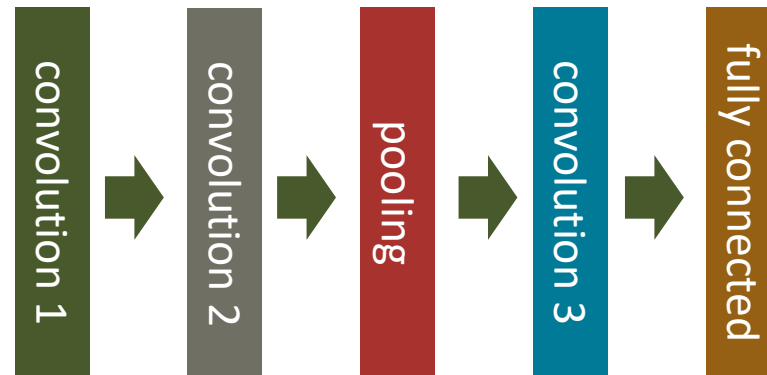
network structure
(fixed)

weights w
(learned)

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$$f(x) = f_n(f_{n-1}(f_{n-2}(\dots f_1(x) \dots)))$$

$$f_1(x) \quad f_2(f_1(x))$$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

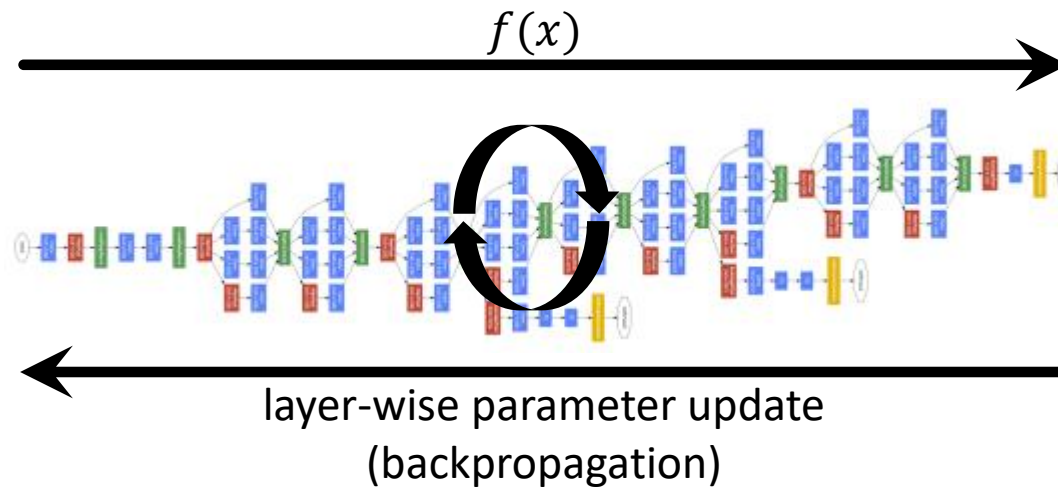
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief theory of supervised deep learning (mini-batch SGD)



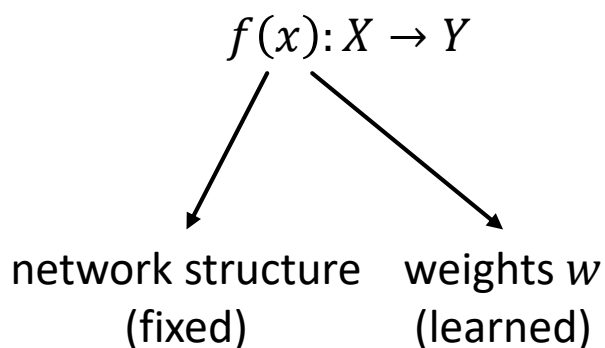
labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

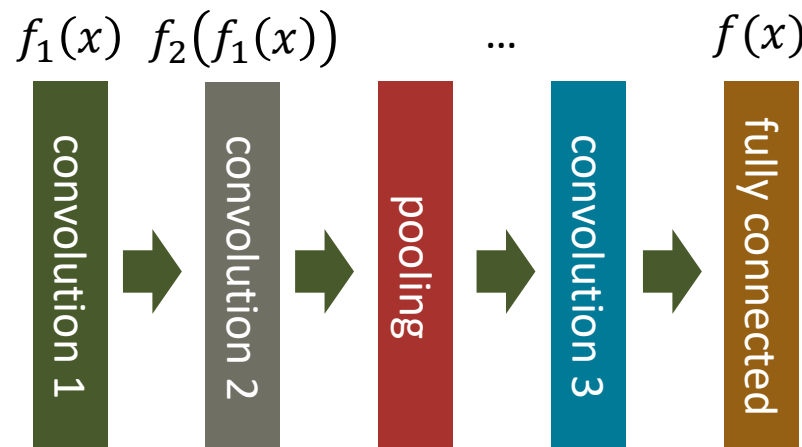
label domain Y

true label $l(x)$



$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$$f(x) = f_n(f_{n-1}(f_{n-2}(\dots f_1(x) \dots)))$$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

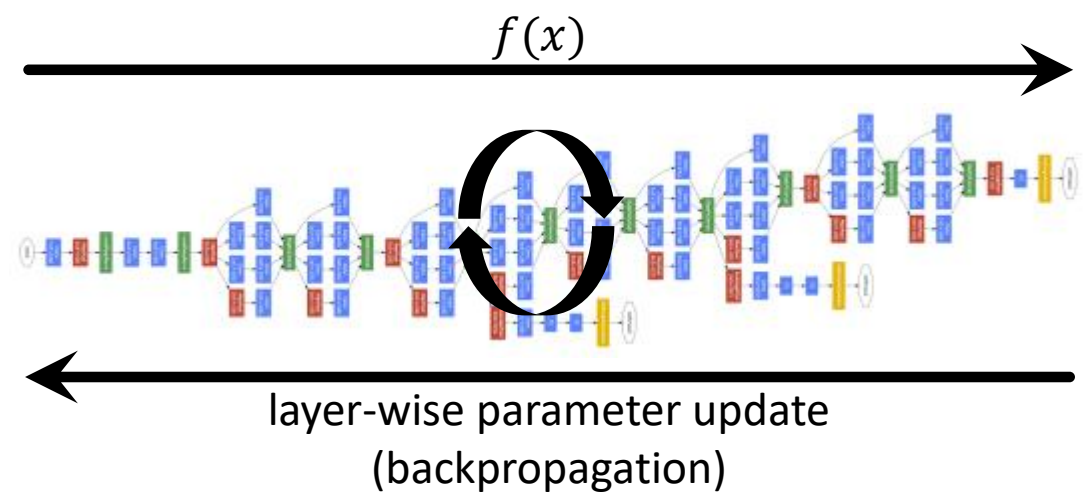
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



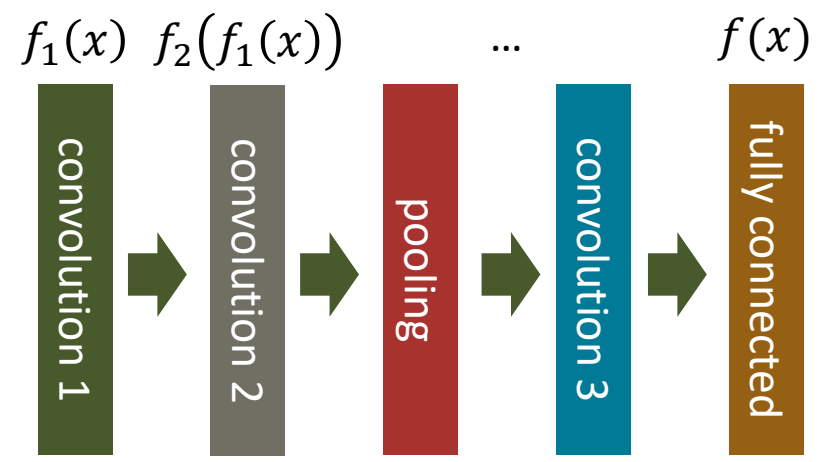
label domain Y	prediction	true label $l(x)$
Cat	0.54	1.00
Dog	0.28	0.00
Airplane	0.07	0.00
Horse	0.33	0.00
Banana	0.02	0.00
Truck	0.02	0.00

label domain Y true label $l(x)$

$$f(x): X \rightarrow Y$$

network structure (fixed) weights w (learned)

$$f(x) = f_n(f_{n-1}(f_{n-2}(\dots f_1(x) \dots)))$$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

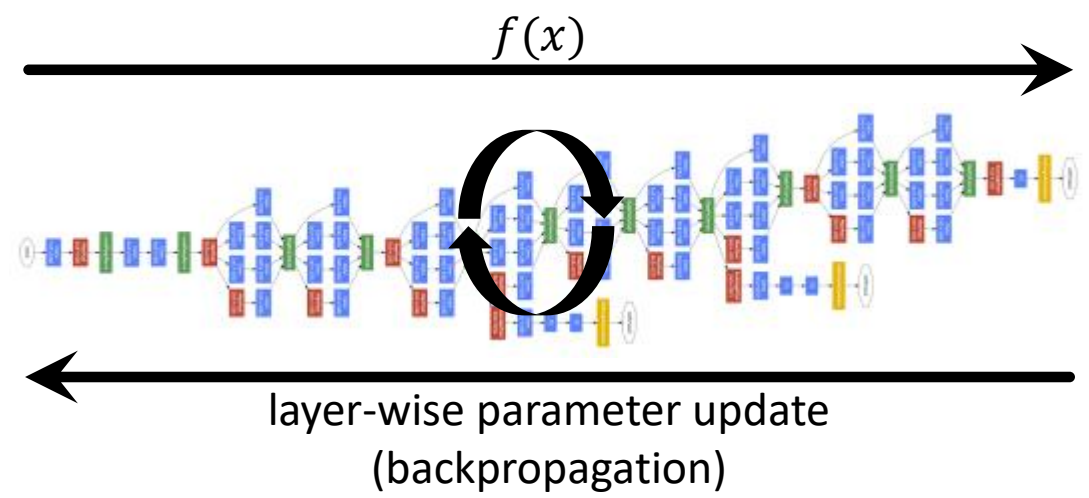
$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

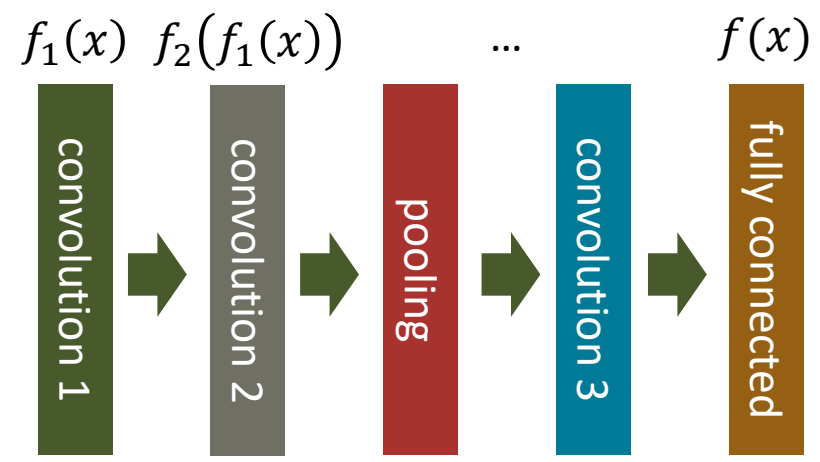
label domain Y true label $l(x)$

$$f(x): X \rightarrow Y$$

network structure (fixed) weights w (learned)

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$$f(x) = f_n \left(f_{n-1} \left(f_{n-2} \left(\dots f_1(x) \dots \right) \right) \right)$$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

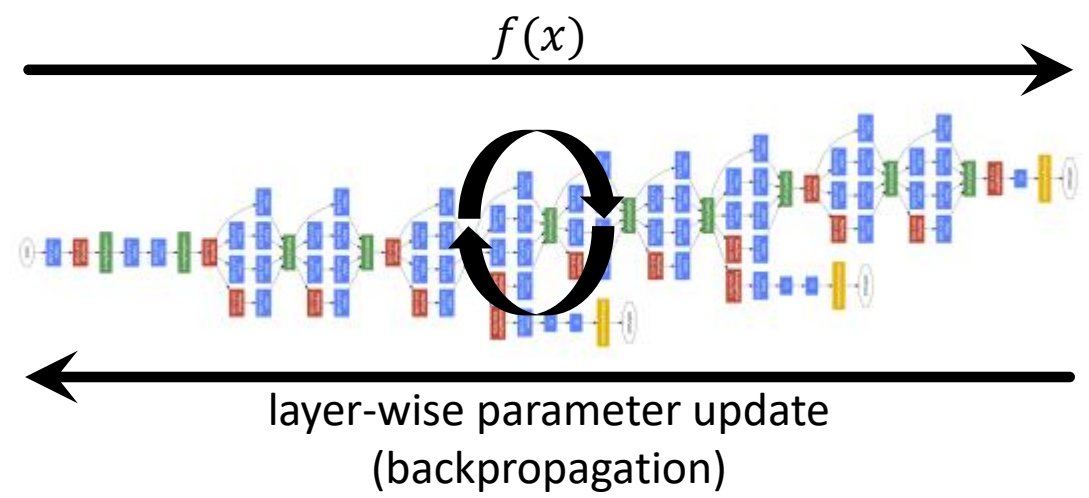
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

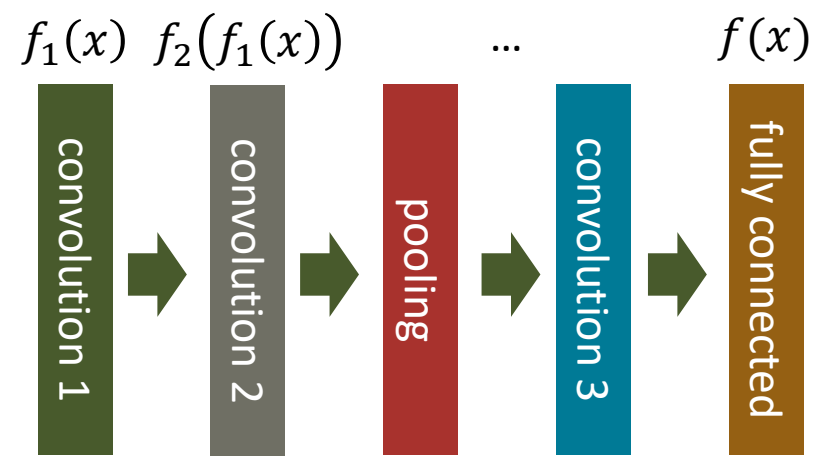
label domain Y true label $l(x)$

$$f(x): X \rightarrow Y$$

network structure (fixed) weights w (learned)

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$$f(x) = f_n \left(f_{n-1} \left(f_{n-2} \left(\dots f_1(x) \dots \right) \right) \right)$$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

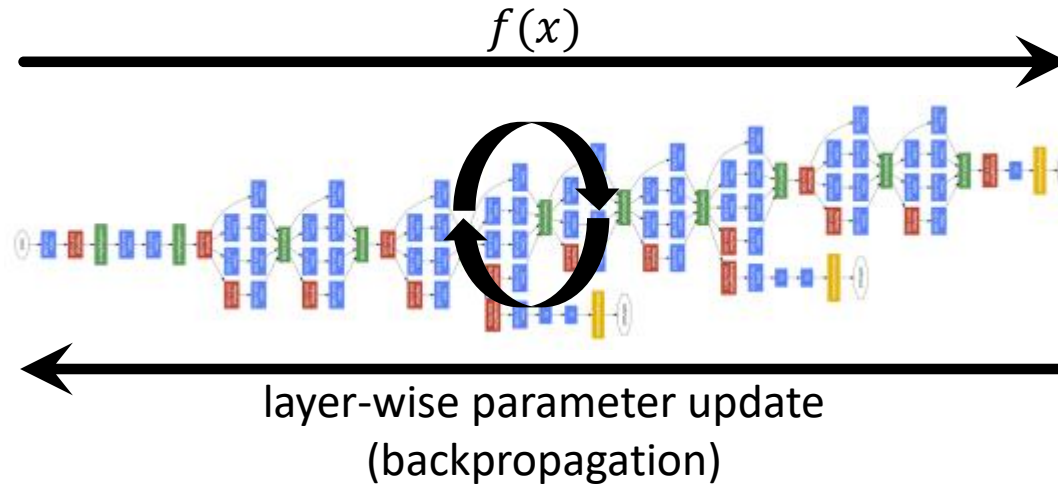
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y

true label $l(x)$

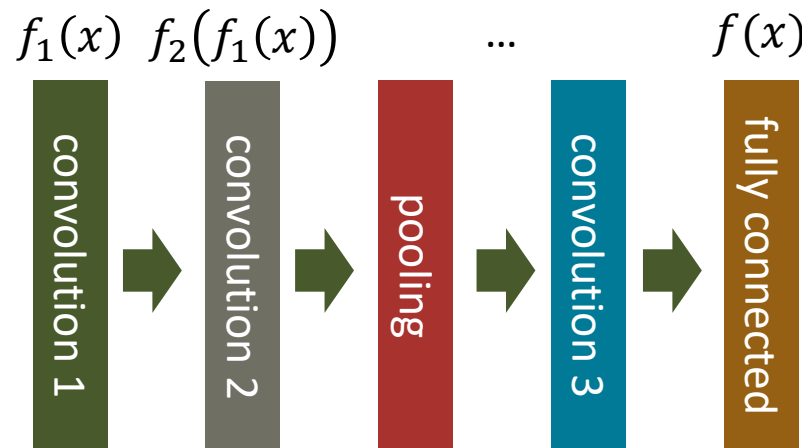
$$f(x): X \rightarrow Y$$

network structure
(fixed)

weights w
(learned)

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$$f(x) = f_n \left(f_{n-1} \left(f_{n-2} \left(\dots f_1(x) \dots \right) \right) \right)$$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

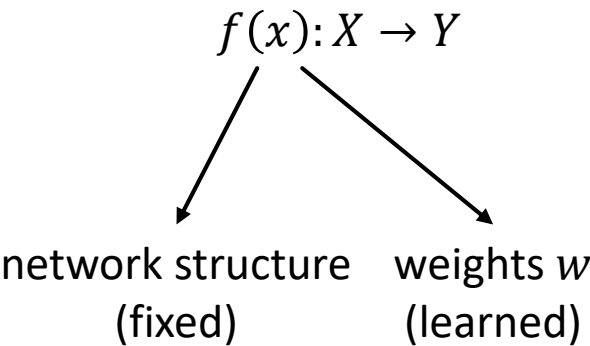
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

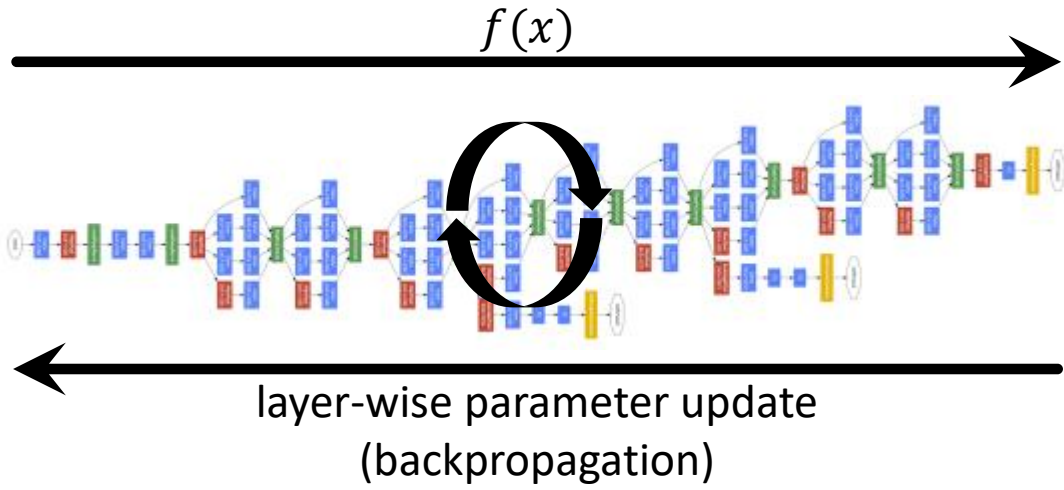
A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



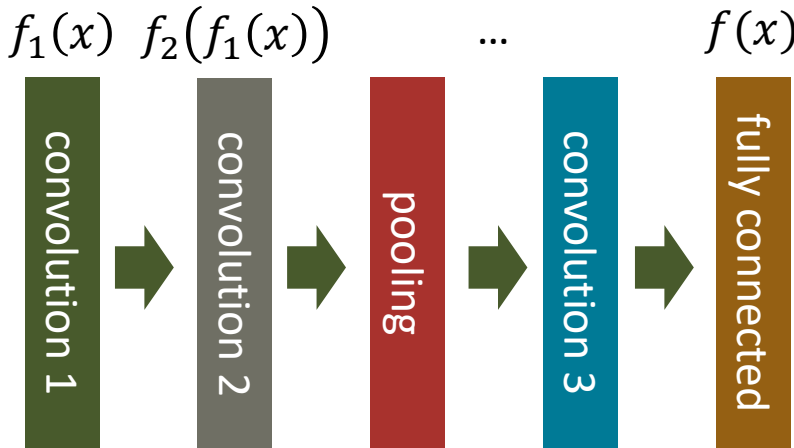
$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

label domain Y true label $l(x)$

$$f(x) = f_n(f_{n-1}(f_{n-2}(\dots f_1(x) \dots)))$$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

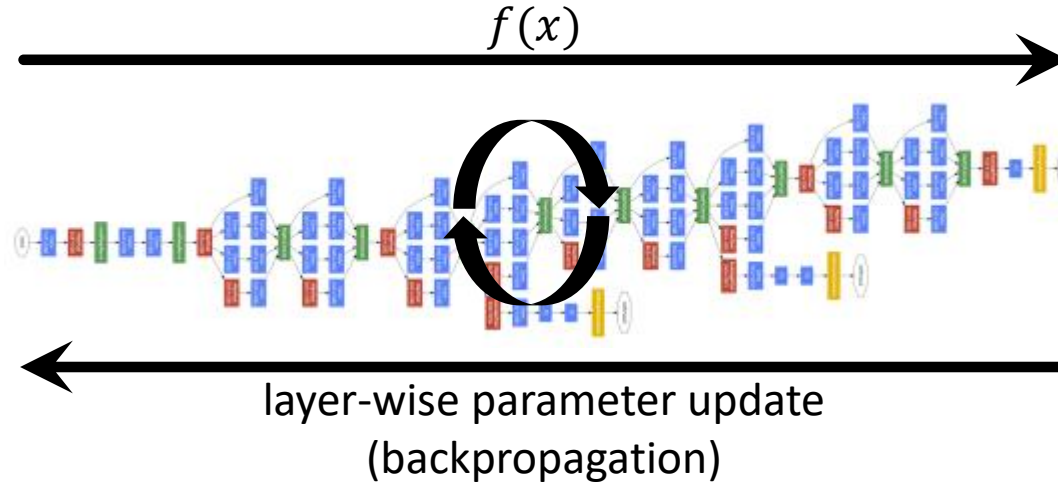
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

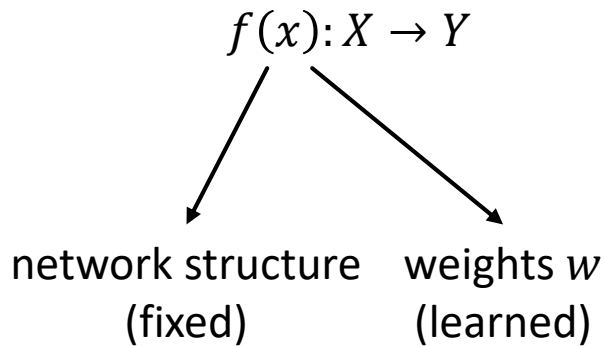
A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$

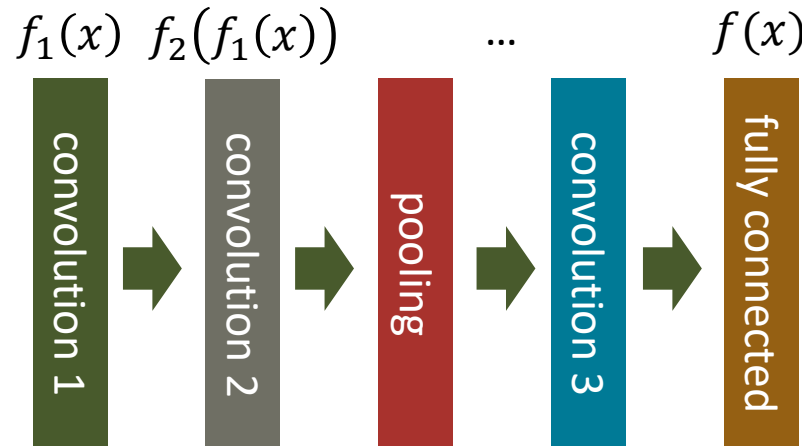


label domain Y	prediction	true label $l(x)$
Cat	0.54	1.00
Dog	0.28	0.00
Airplane	0.07	0.00
Horse	0.33	0.00
Banana	0.02	0.00
Truck	0.02	0.00



$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$$f(x) = f_n(f_{n-1}(f_{n-2}(\dots f_1(x) \dots)))$$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

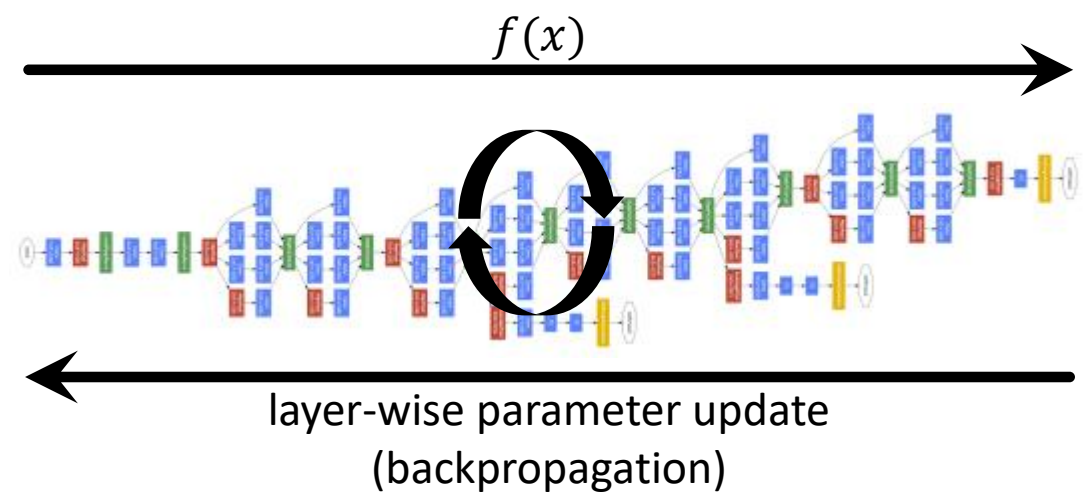
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

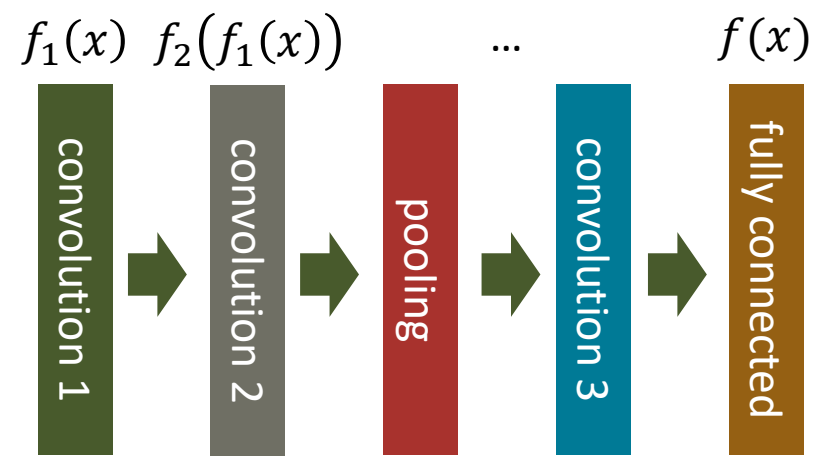
label domain Y true label $l(x)$

$$f(x): X \rightarrow Y$$

network structure (fixed) weights w (learned)

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$$f(x) = f_n \left(f_{n-1} \left(f_{n-2} \left(\dots f_1(x) \dots \right) \right) \right)$$

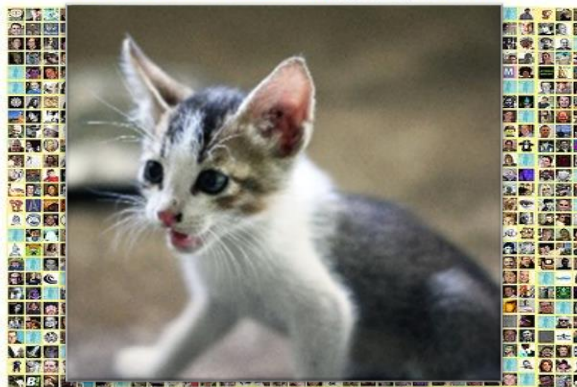


$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

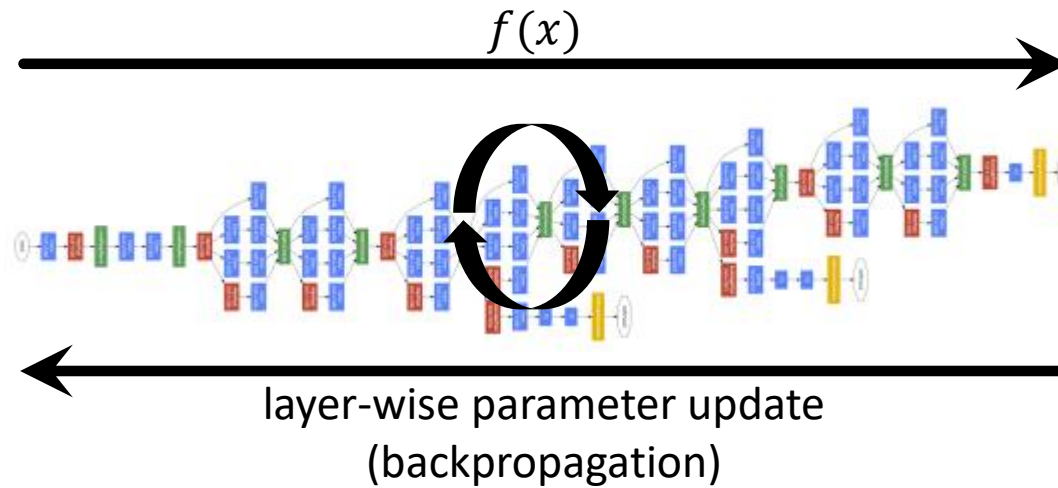
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$



label domain Y	prediction	true label $l(x)$
Cat	0.54	1.00
Dog	0.28	0.00
Airplane	0.07	0.00
Horse	0.33	0.00
Banana	0.02	0.00
Truck	0.02	0.00

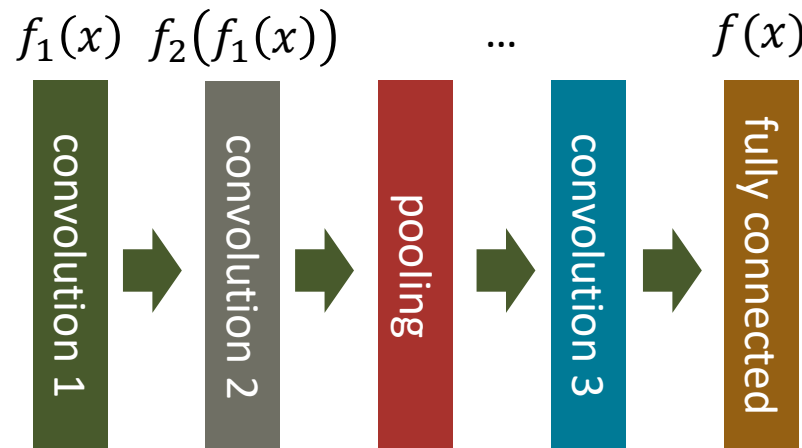
label domain Y true label $l(x)$

$$f(x): X \rightarrow Y$$

network structure (fixed) weights w (learned)

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$$f(x) = f_n \left(f_{n-1} \left(f_{n-2} \left(\dots f_1(x) \dots \right) \right) \right)$$



$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

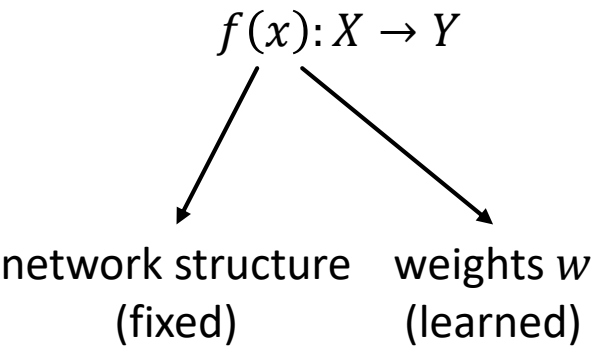
$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

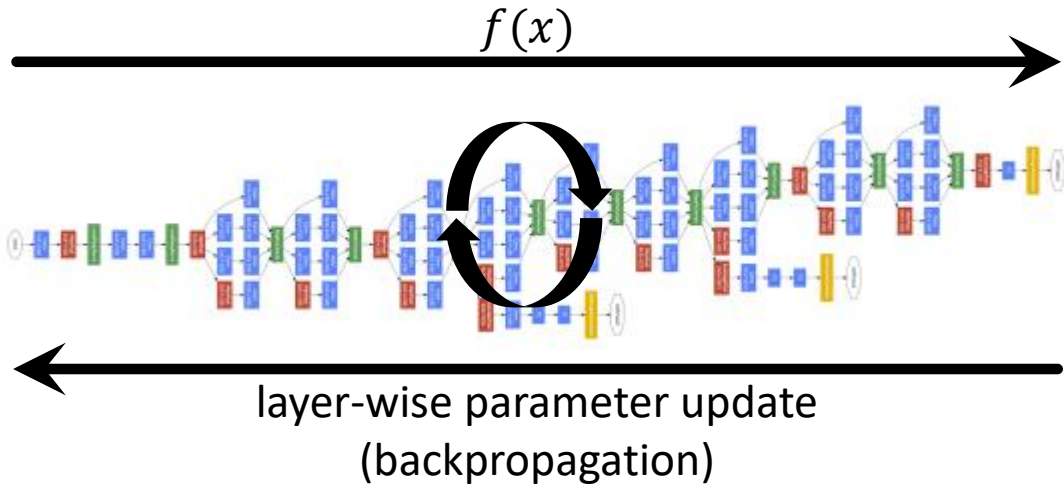
A brief theory of supervised deep learning (mini-batch SGD)



labeled samples $x \in X \subset \mathcal{D}$

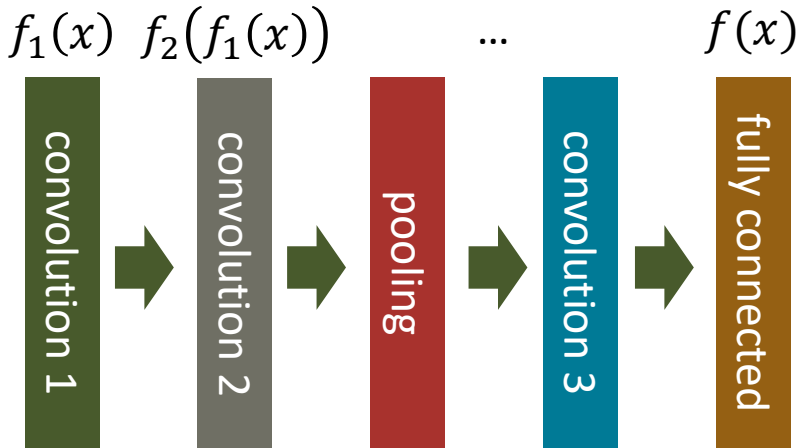


$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$



label domain Y	prediction	true label $l(x)$
Cat	0.54	1.00
Dog	0.28	0.00
Airplane	0.07	0.00
Horse	0.33	0.00
Banana	0.02	0.00
Truck	0.02	0.00

$$f(x) = f_n \left(f_{n-1} \left(f_{n-2} \left(\dots f_1(x) \dots \right) \right) \right)$$



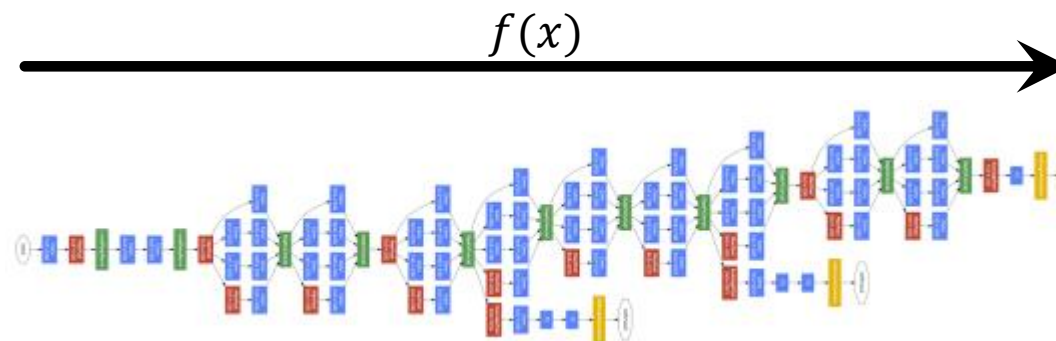
$$\ell_{sq}(w, x) = (f(x) - l(x))^2$$

$$\ell_{0-1}(w, x) = \begin{cases} 0 & f(x) = l(x) \\ 1 & f(x) \neq l(x) \end{cases}$$

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief digression on backpropagation

- "Backward propagation of errors" (Rumelhart, Hinton, & Williams 1986)

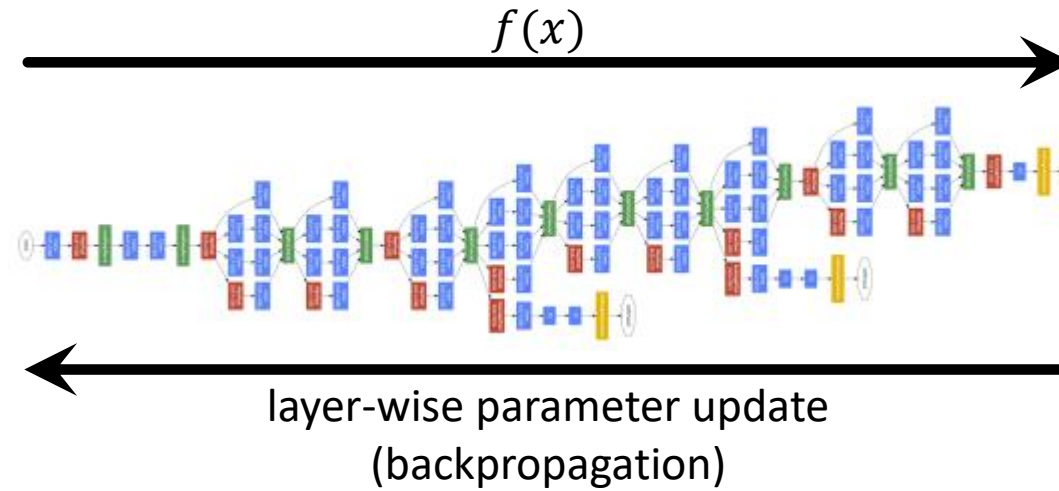


Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief digression on backpropagation

- "Backward propagation of errors" (Rumelhart, Hinton, & Williams 1986)



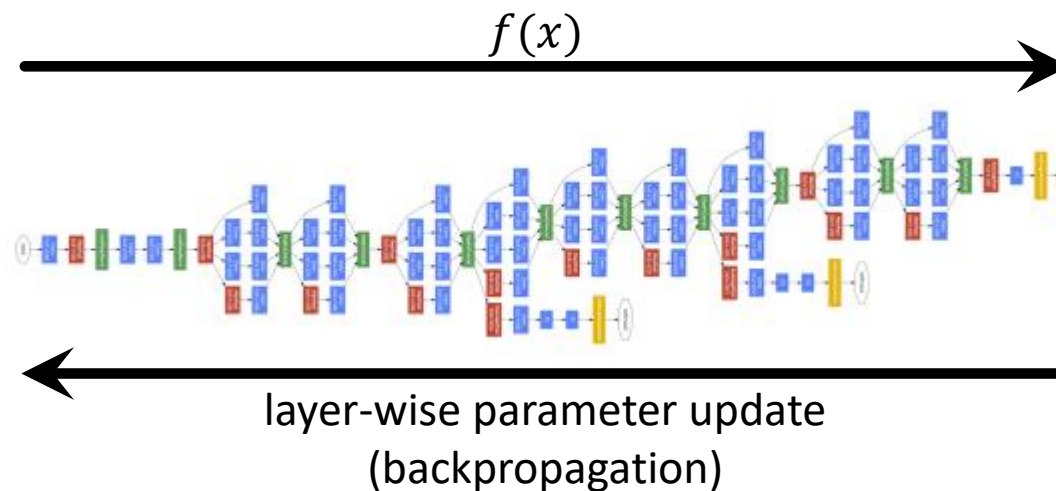
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$



A brief digression on backpropagation

- "Backward propagation of errors" (Rumelhart, Hinton, & Williams 1986)



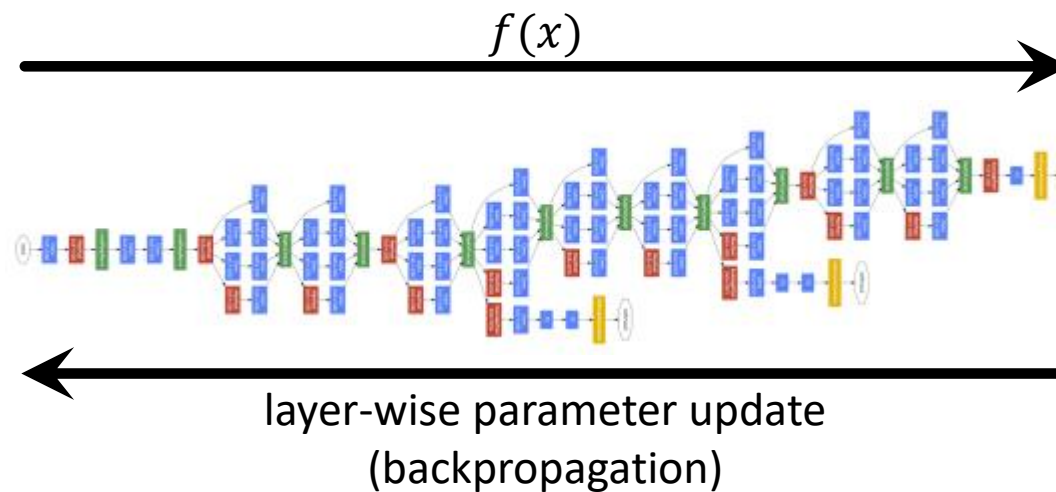
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$



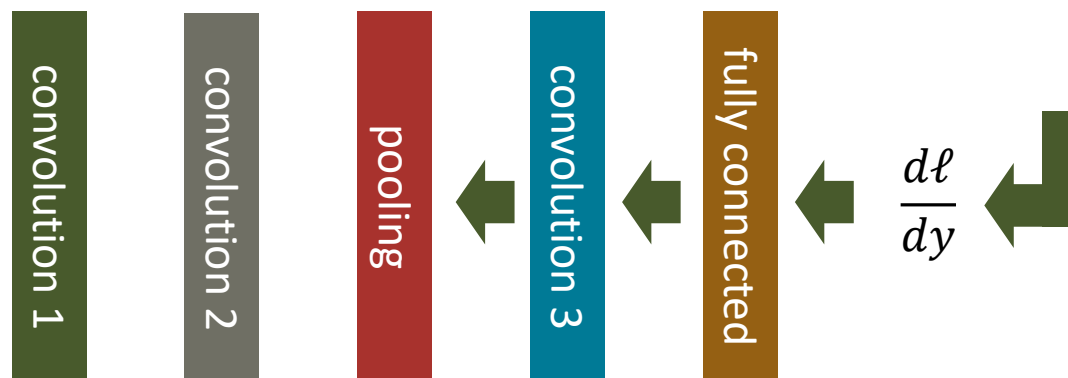
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



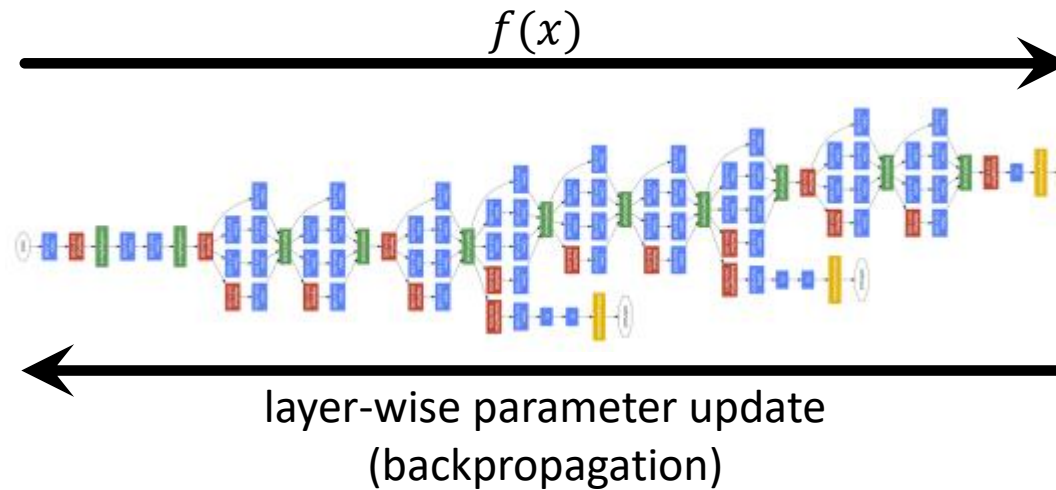
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

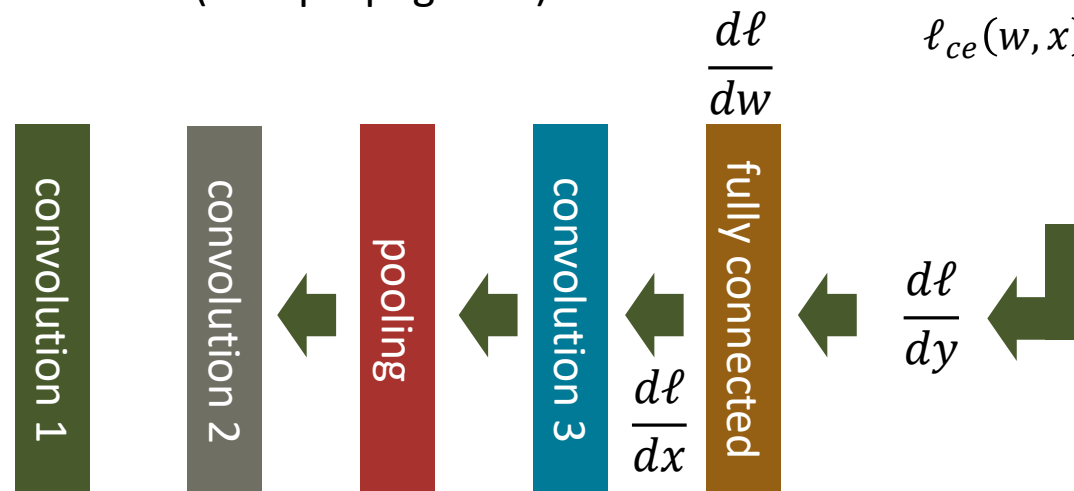


A brief digression on backpropagation

- "Backward propagation of errors" (Rumelhart, Hinton, & Williams 1986)



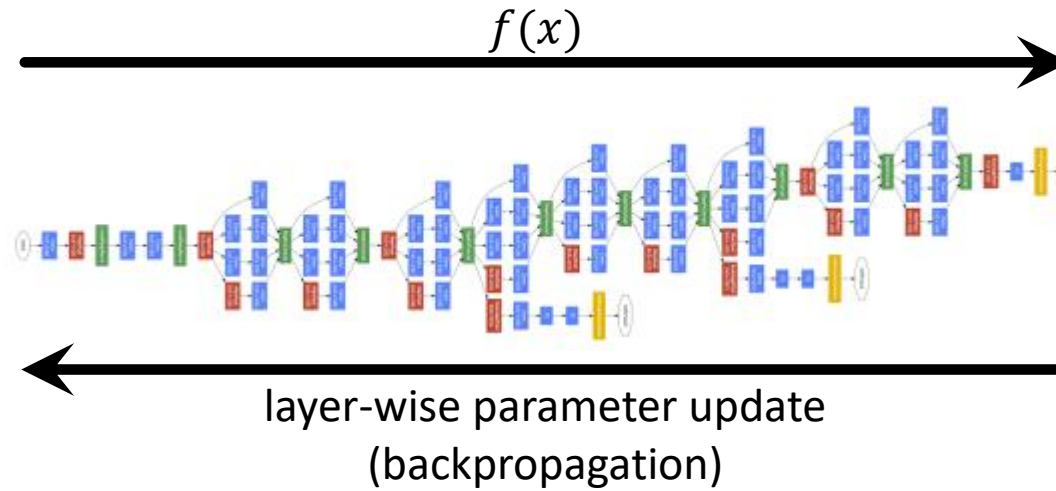
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00



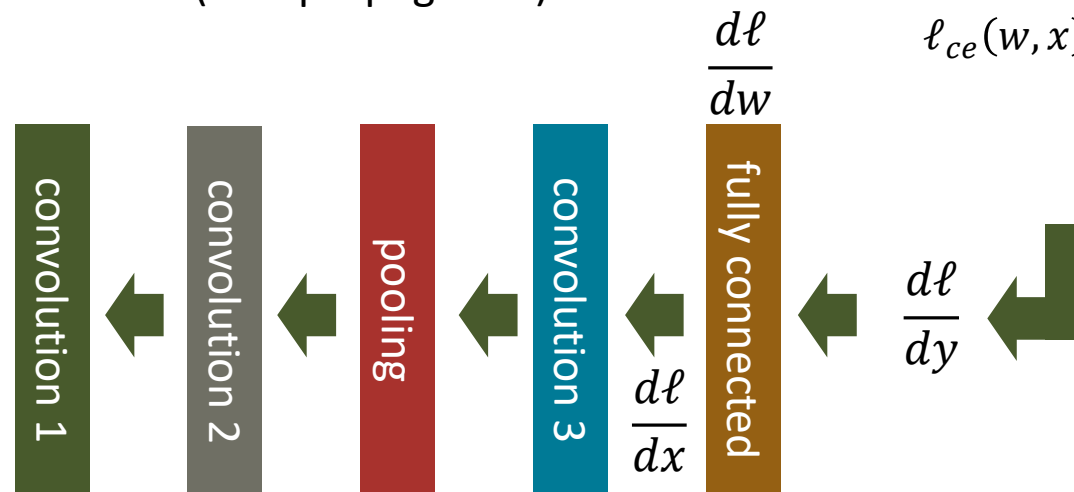
$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



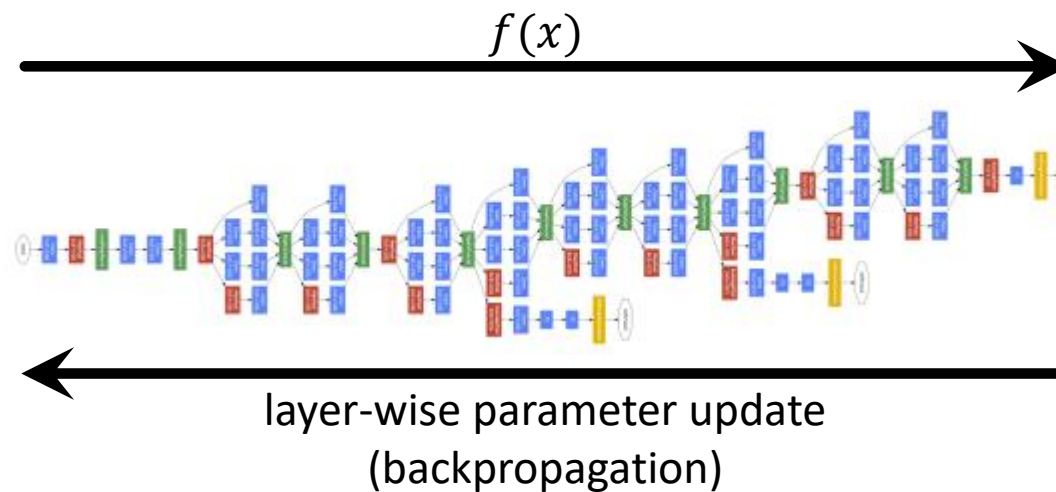
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

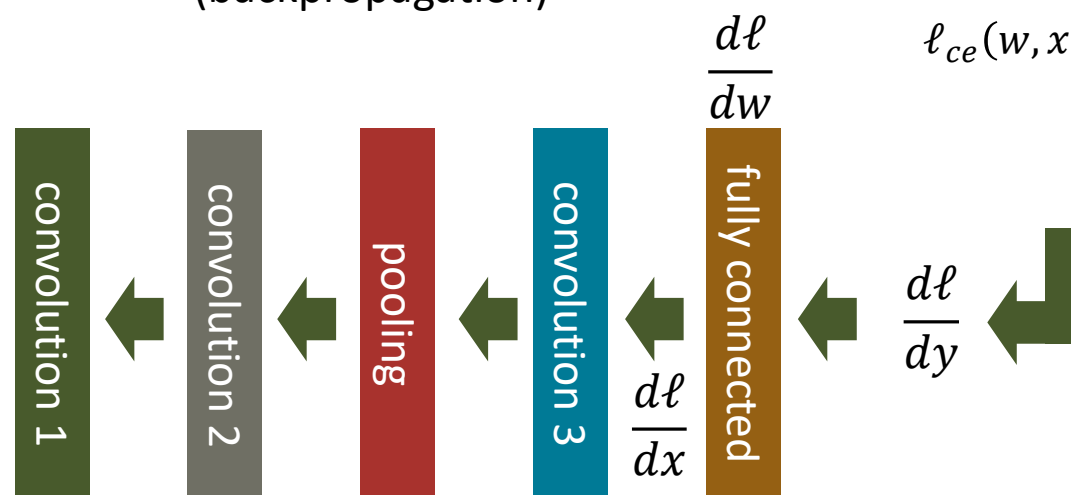
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



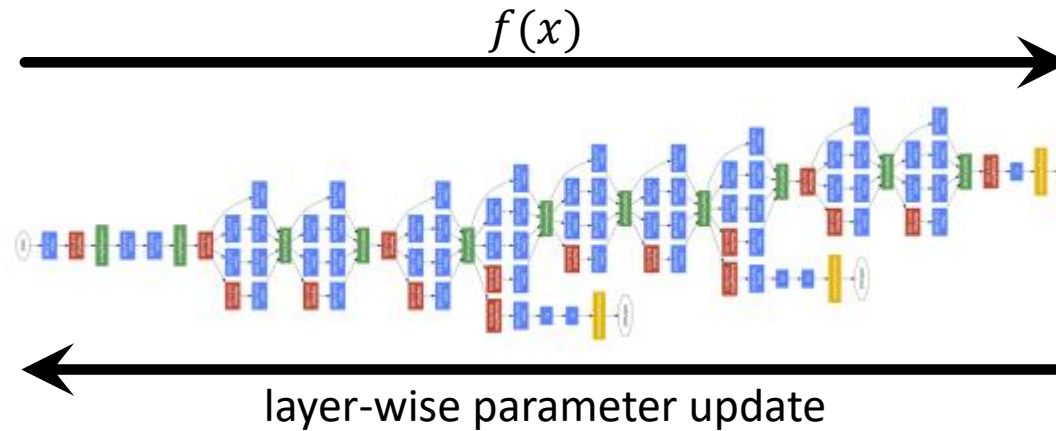
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$



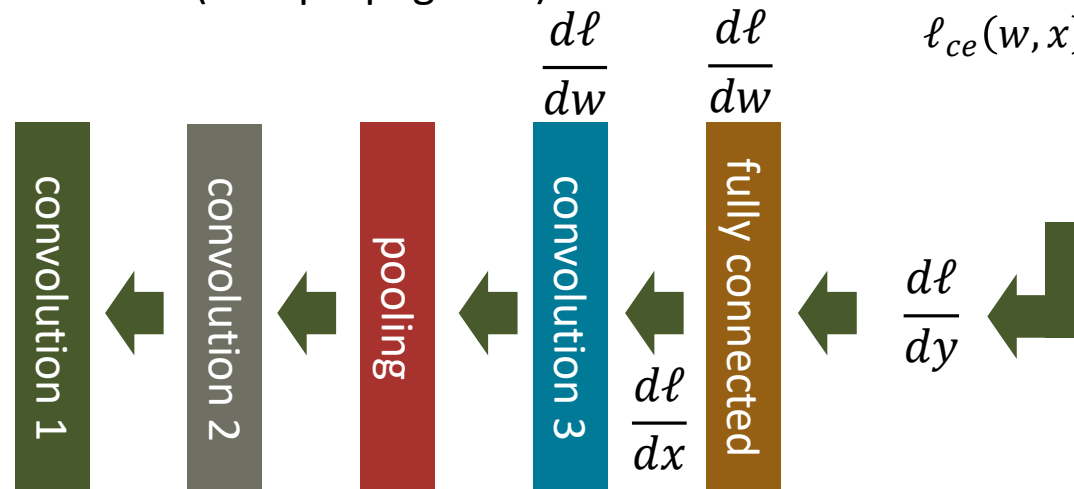
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

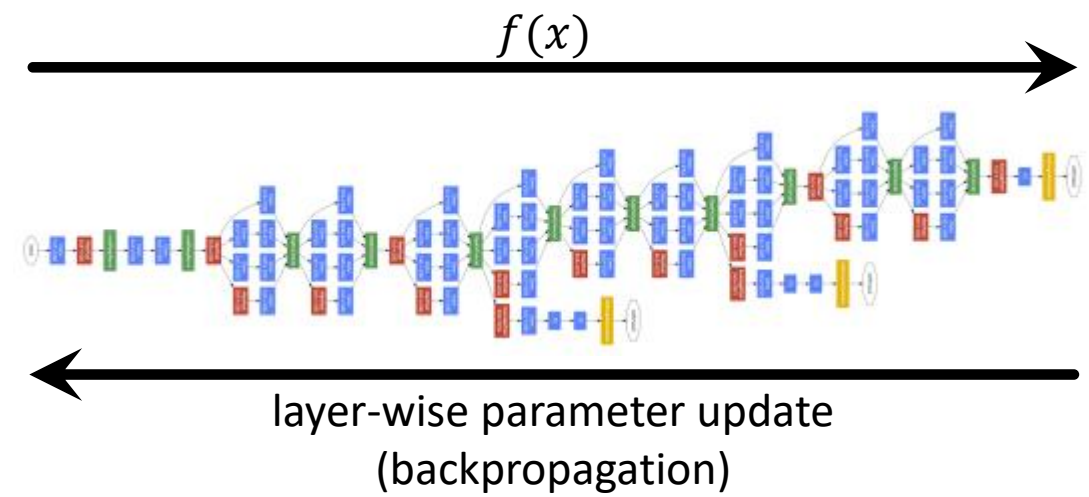
layer-wise parameter update
(backpropagation)



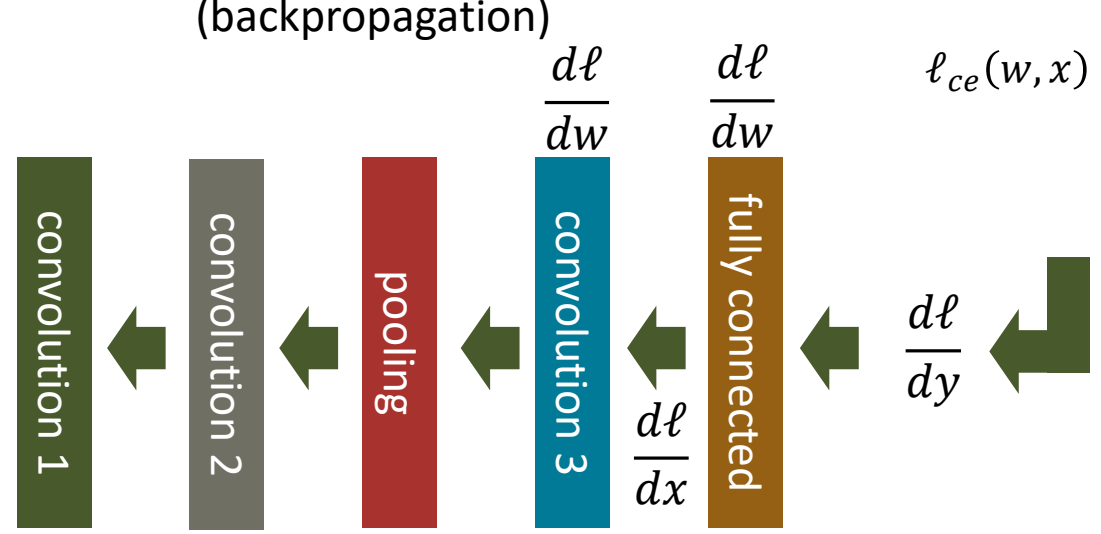
$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief digression on backpropagation

- "Backward propagation of errors" (Rumelhart, Hinton, & Williams 1986)



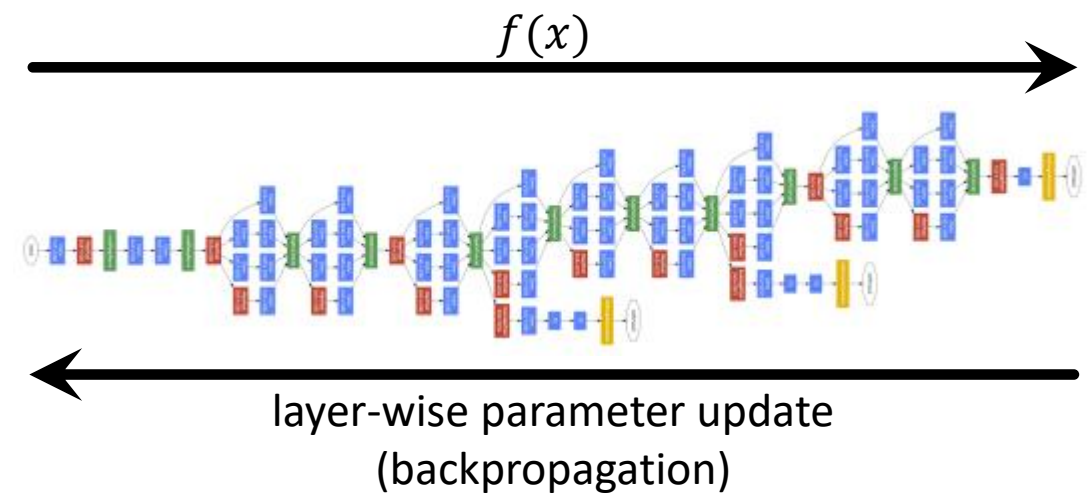
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00



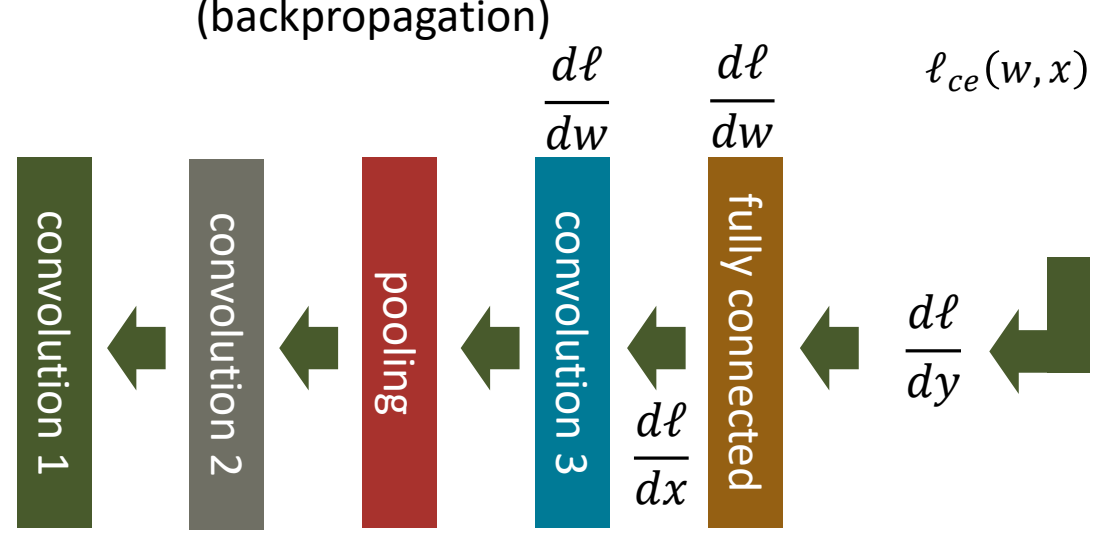
$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief digression on backpropagation

- "Backward propagation of errors" (Rumelhart, Hinton, & Williams 1986)



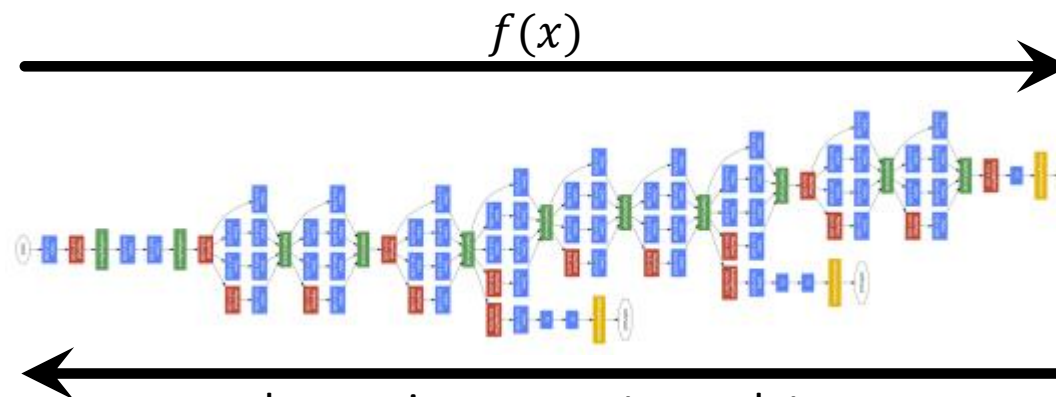
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

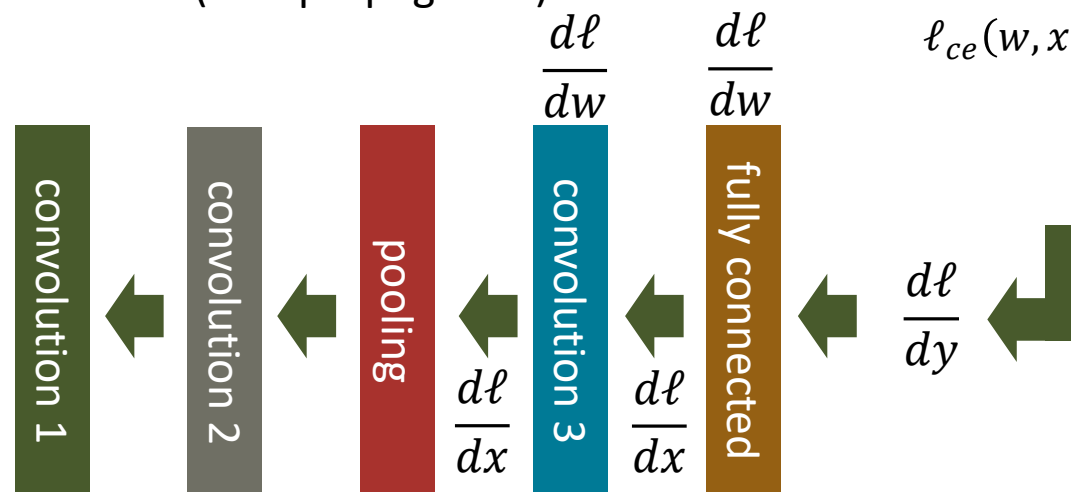
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

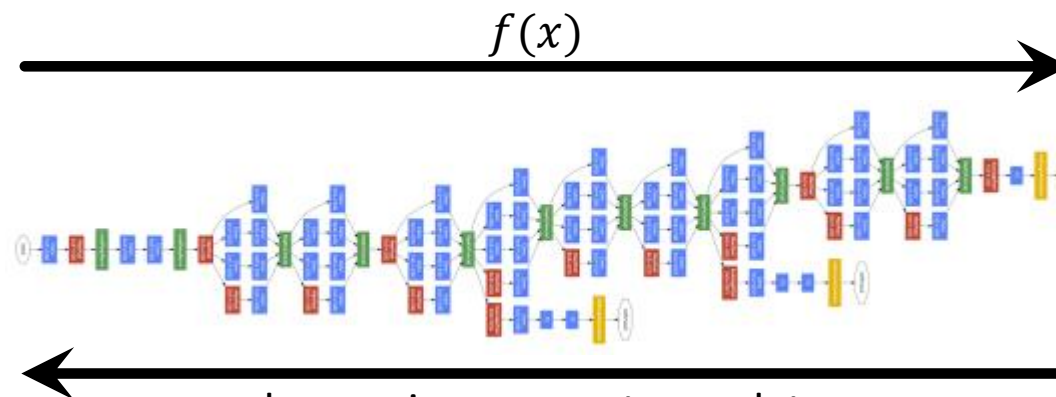
layer-wise parameter update
(backpropagation)



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief digression on backpropagation

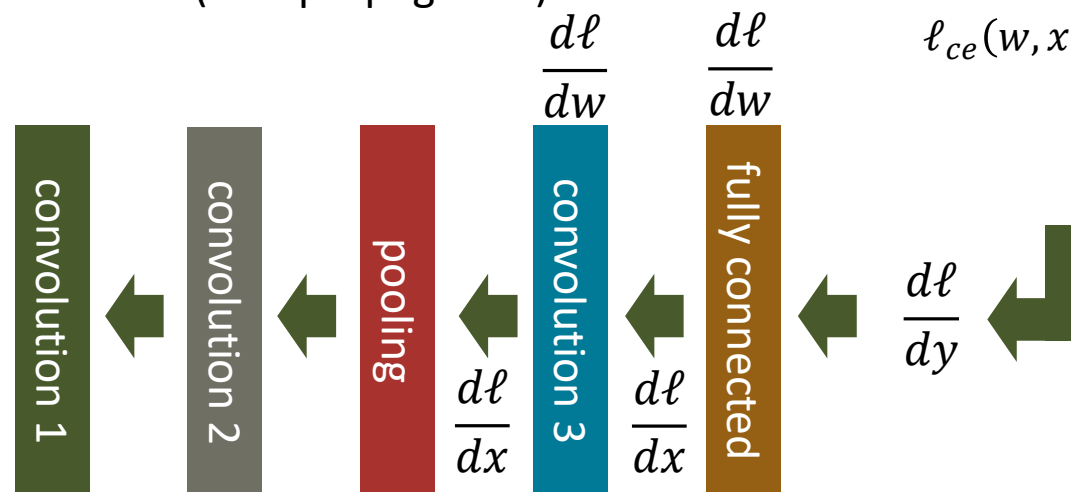
- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

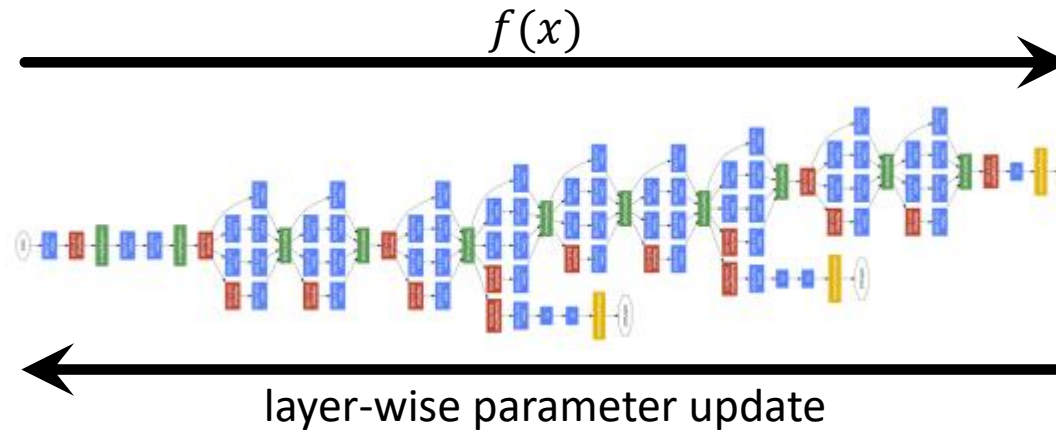
layer-wise parameter update
(backpropagation)

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$



A brief digression on backpropagation

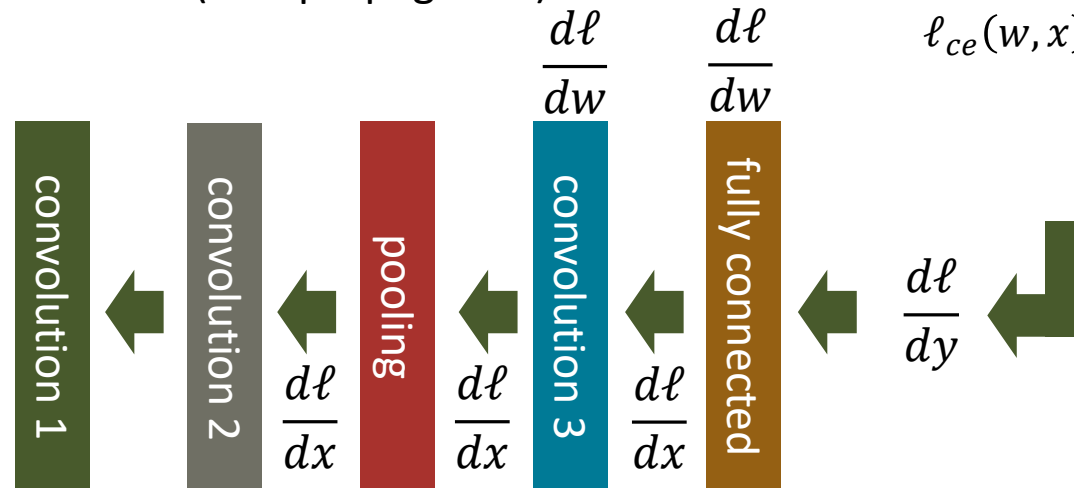
- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

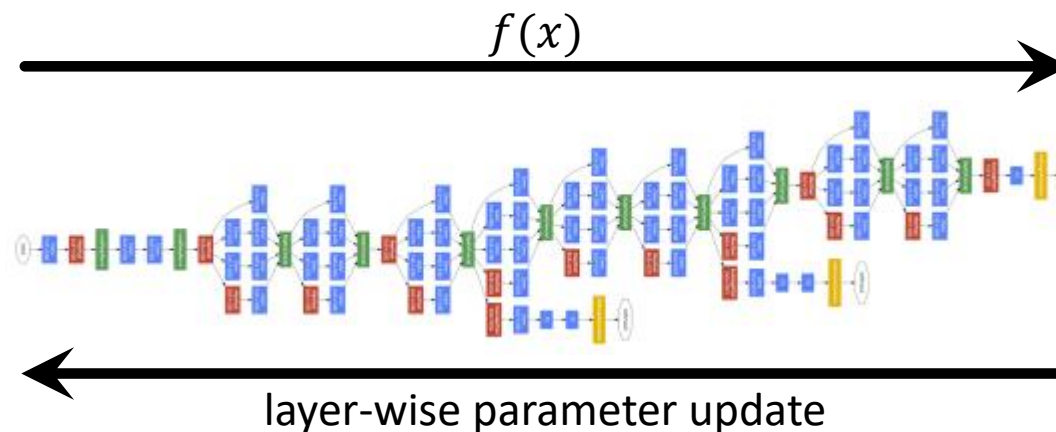
layer-wise parameter update
(backpropagation)

$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$



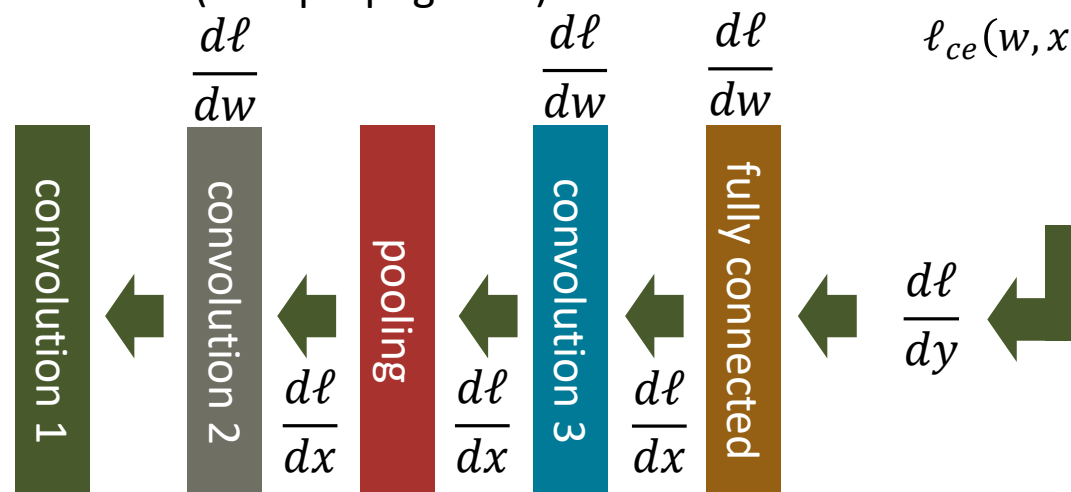
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

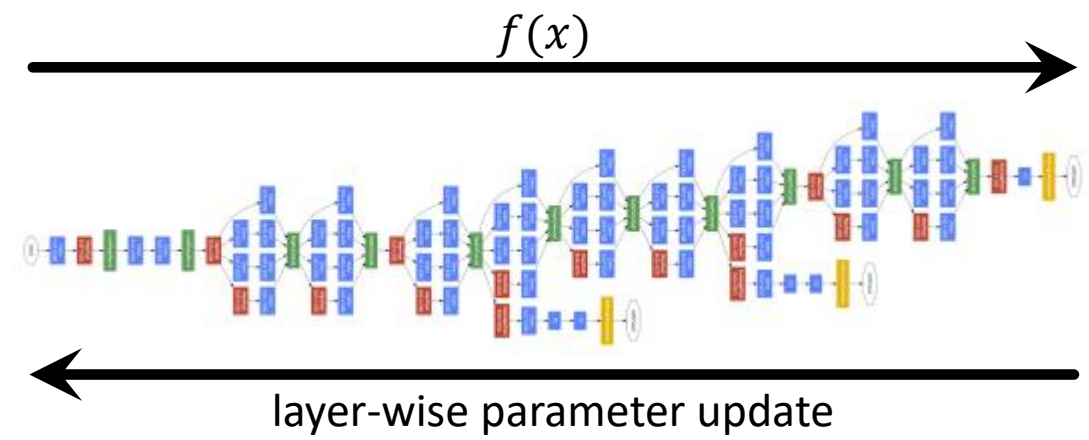
layer-wise parameter update
(backpropagation)



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

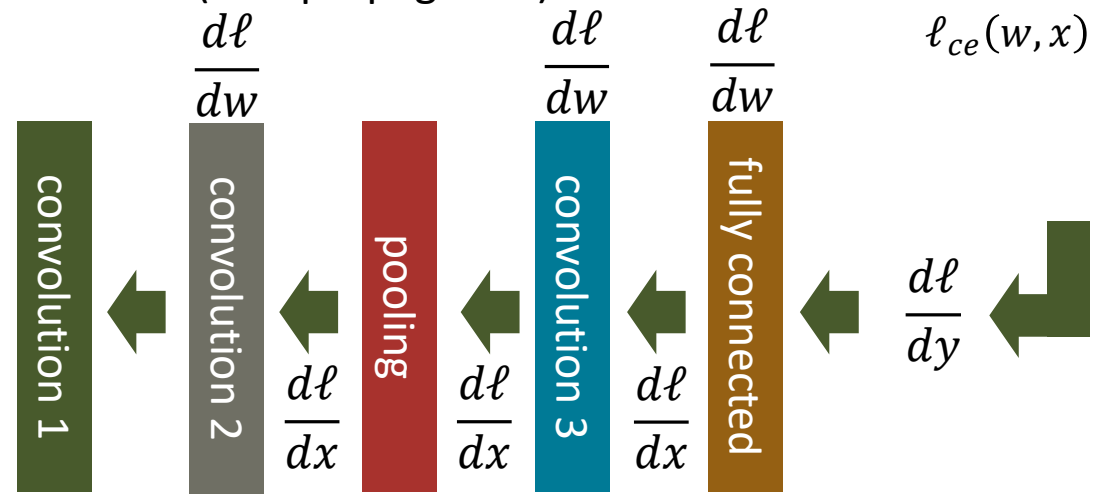
A brief digression on backpropagation

- "Backward propagation of errors" (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

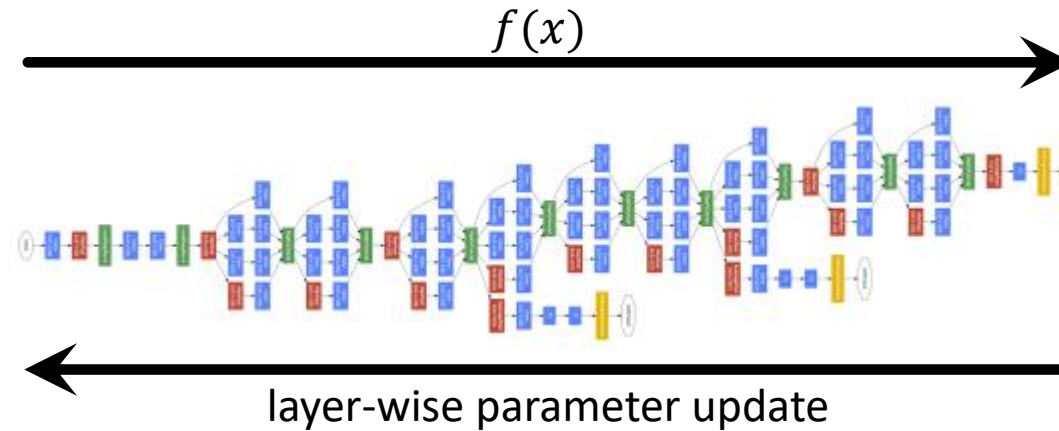
layer-wise parameter update
(backpropagation)



$$l_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

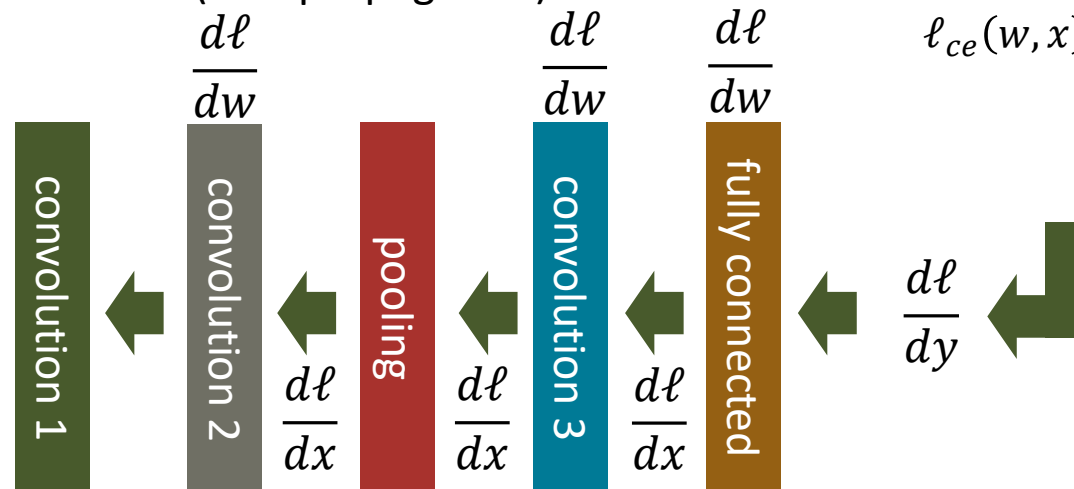
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

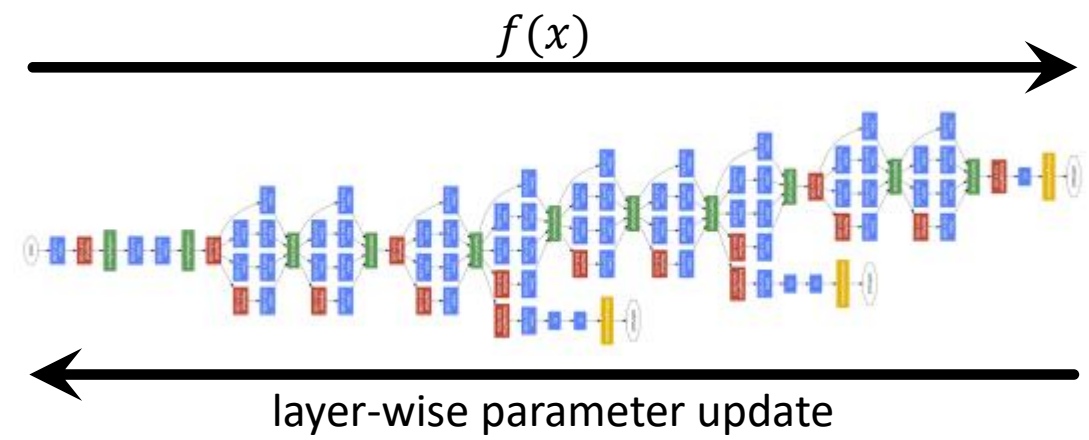
layer-wise parameter update
(backpropagation)



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

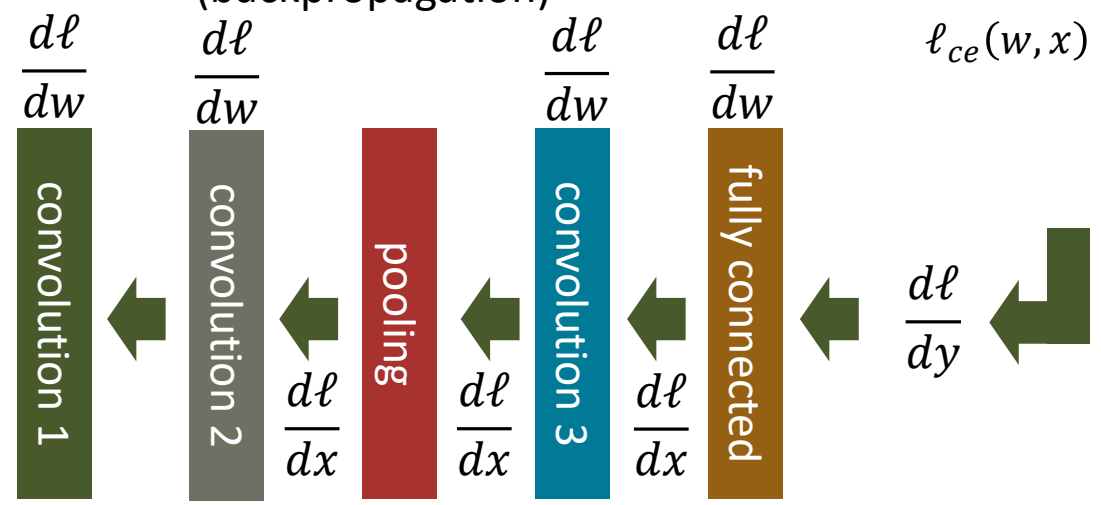
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

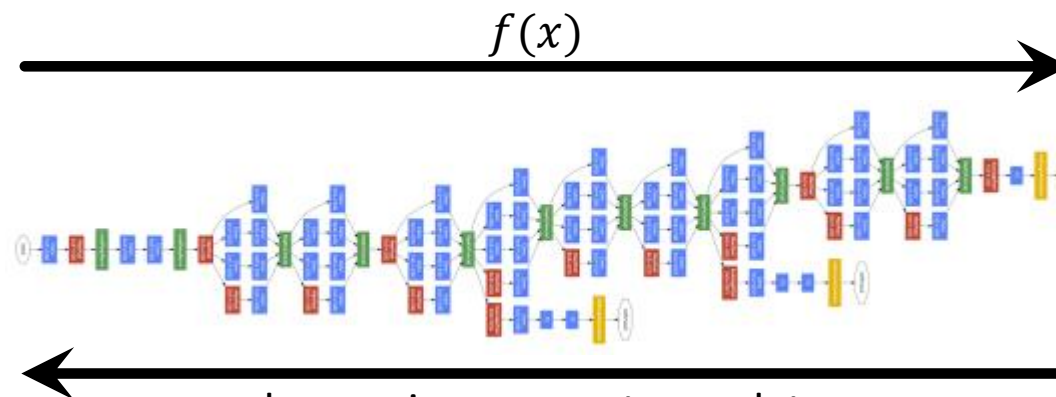
layer-wise parameter update
(backpropagation)



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

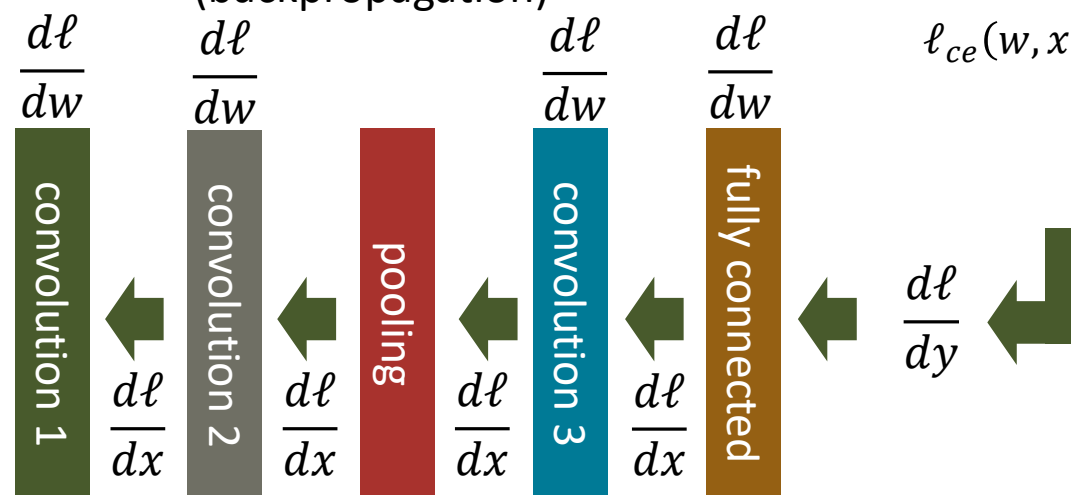
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

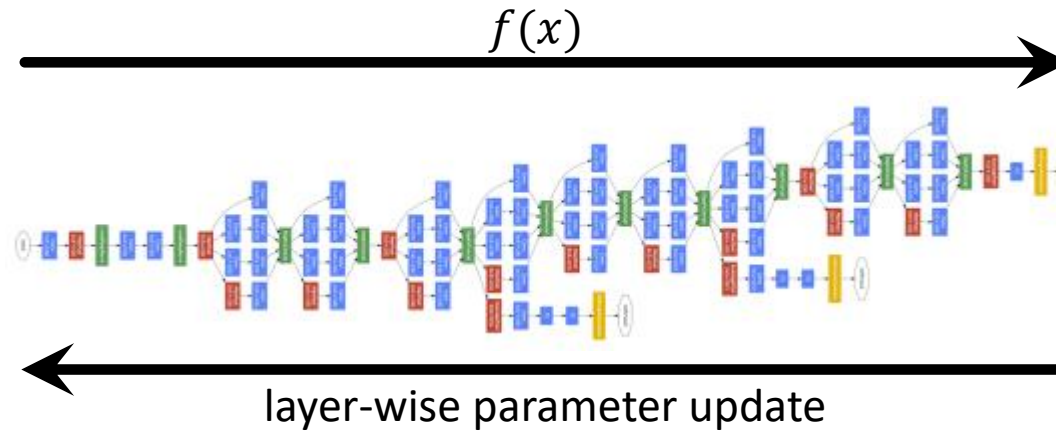
layer-wise parameter update
(backpropagation)



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

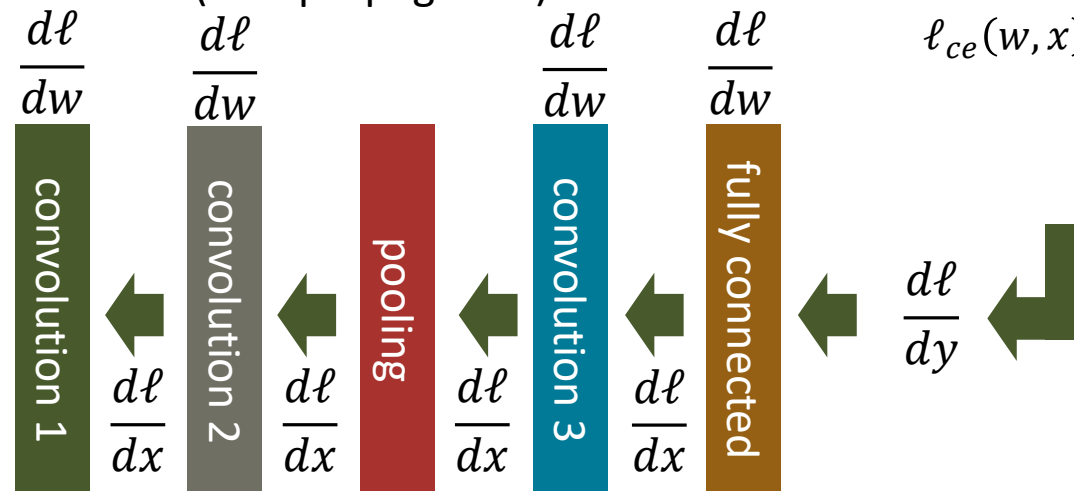
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

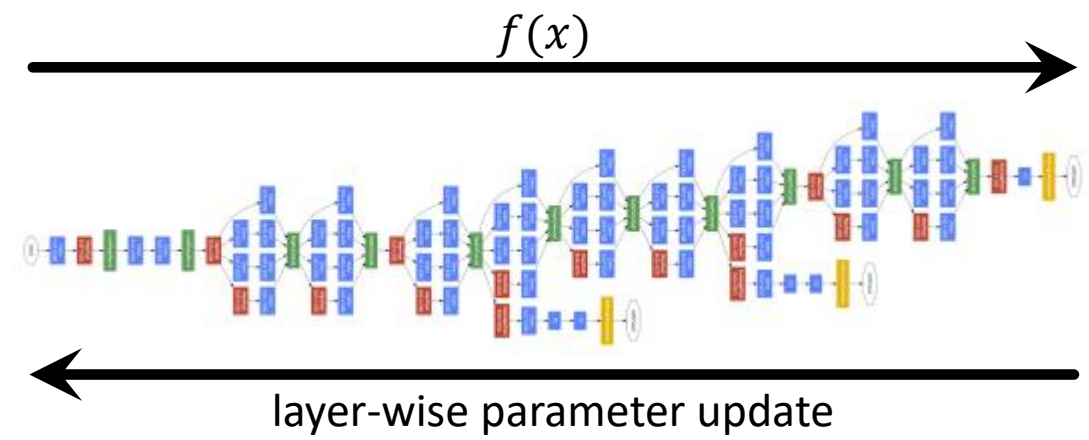
layer-wise parameter update
(backpropagation)



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

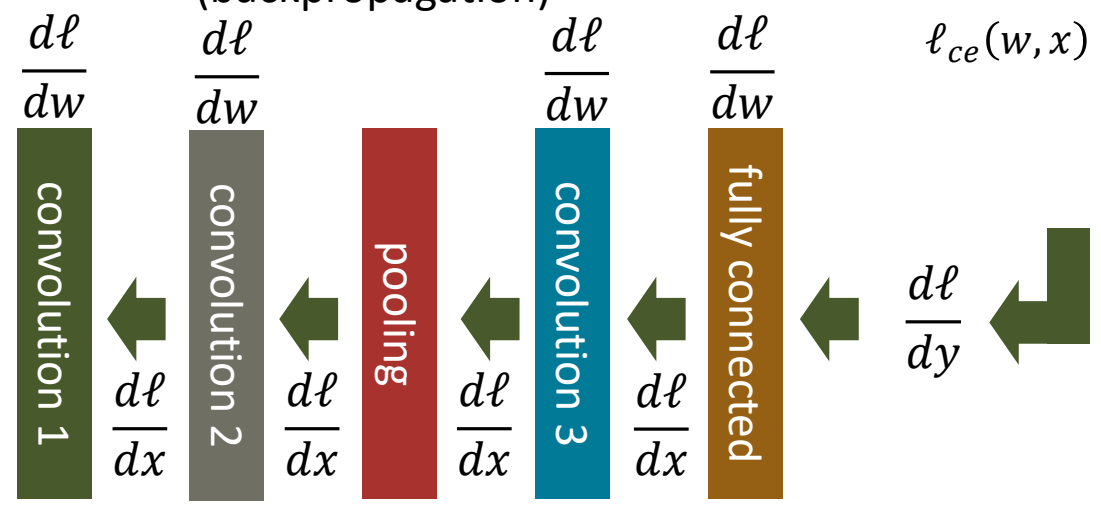
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

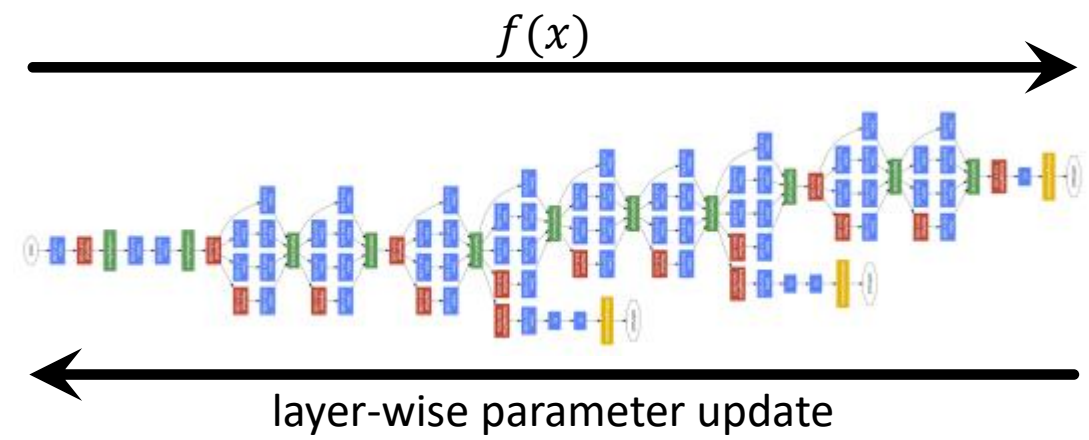
layer-wise parameter update
(backpropagation)



$$l_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

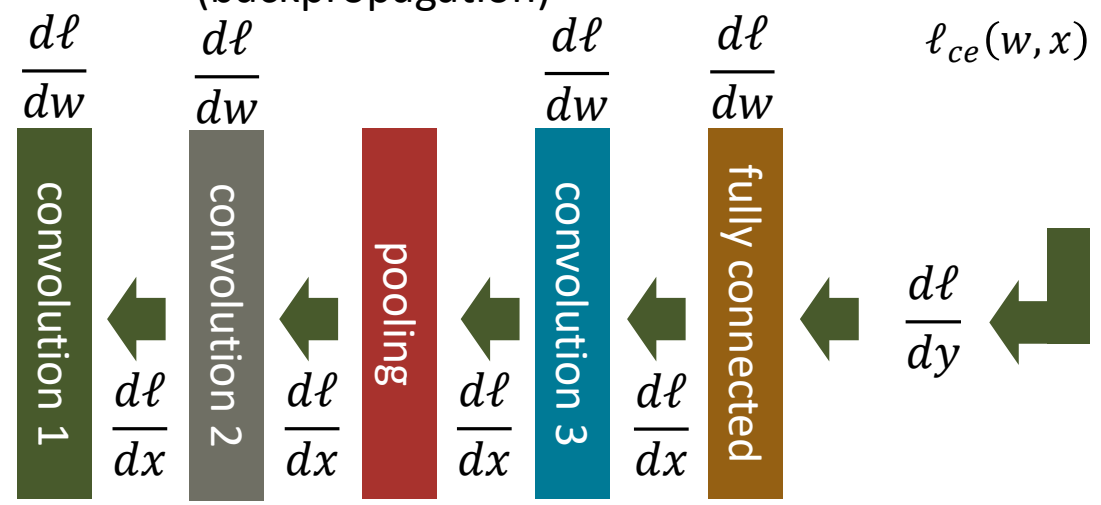
A brief digression on backpropagation

- "Backward propagation of errors" (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

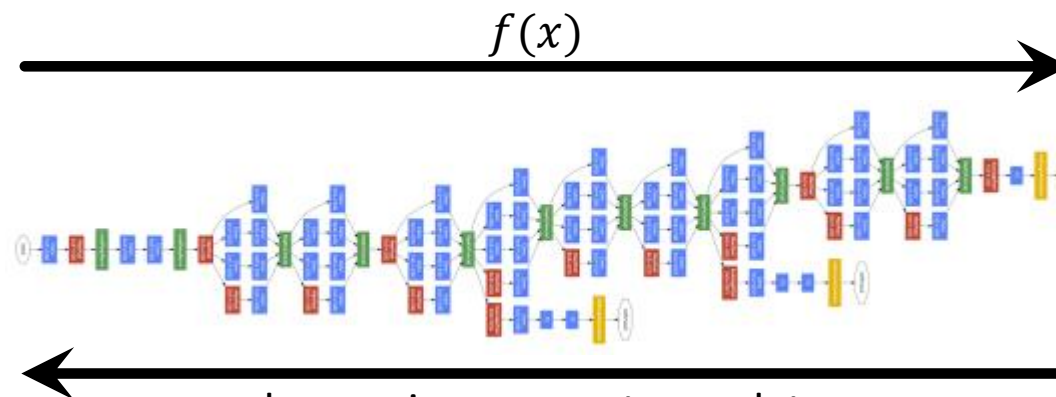
layer-wise parameter update
(backpropagation)



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

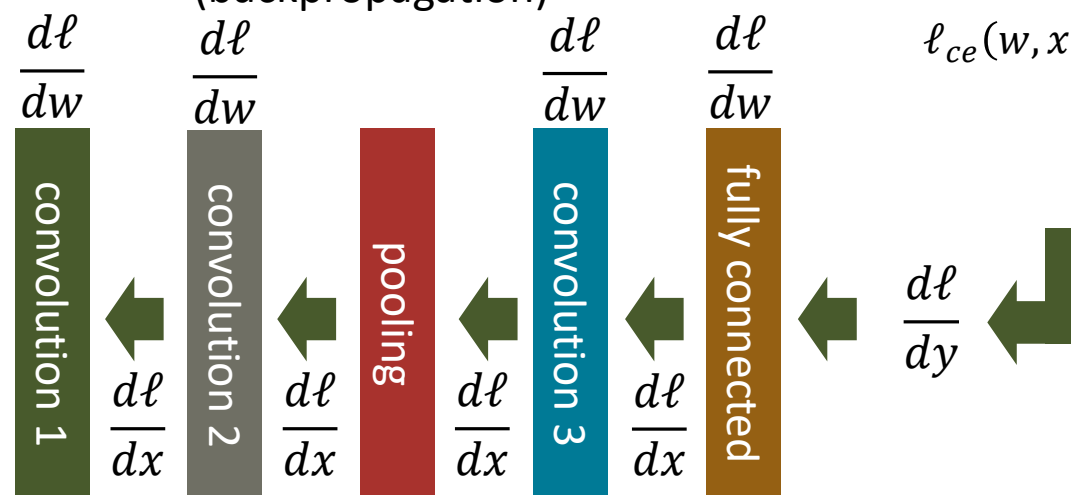
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

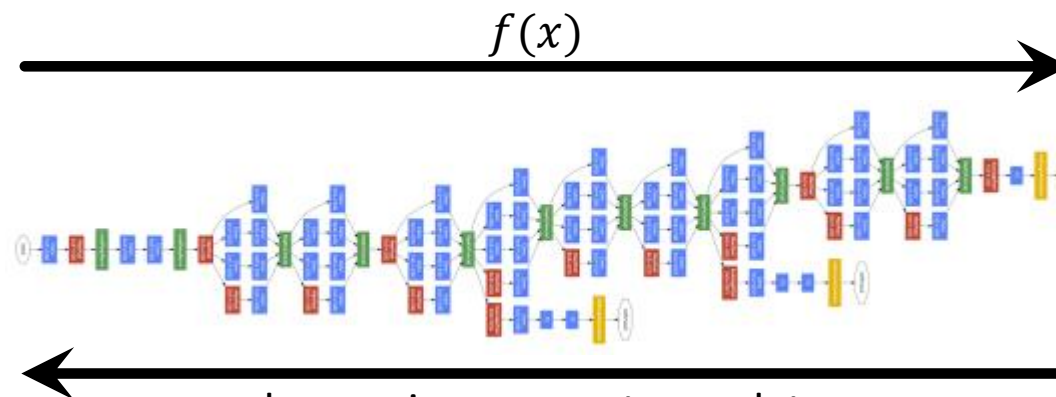
layer-wise parameter update
(backpropagation)



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

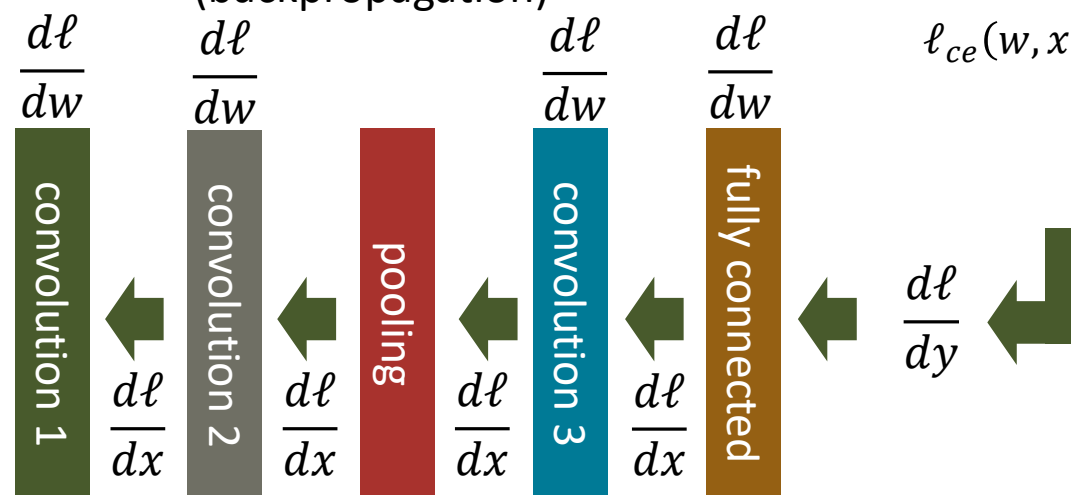
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

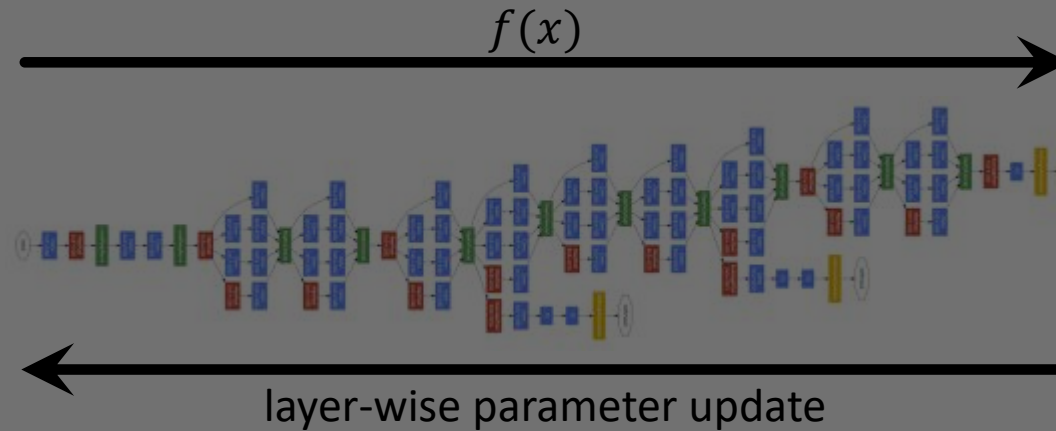
layer-wise parameter update
(backpropagation)



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

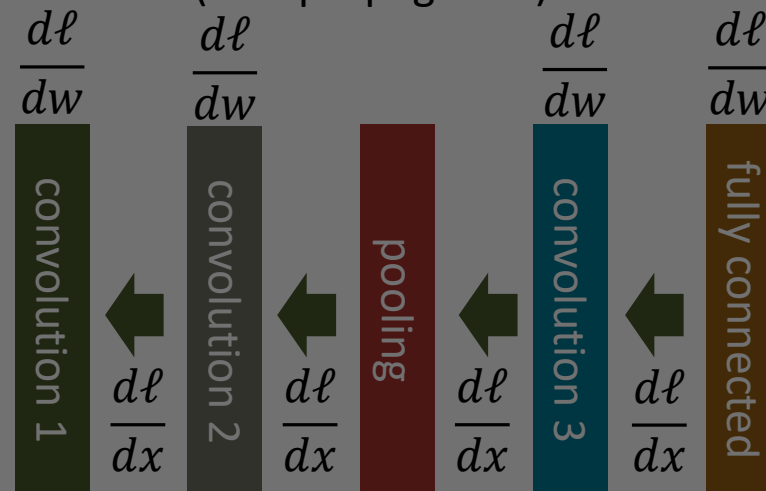
A brief digression on backpropagation

- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

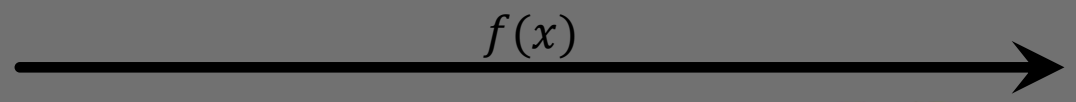
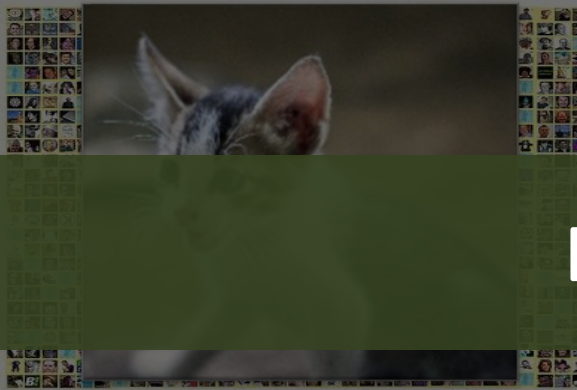
layer-wise parameter update
(backpropagation)



$$\ell_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

A brief digression on backpropagation

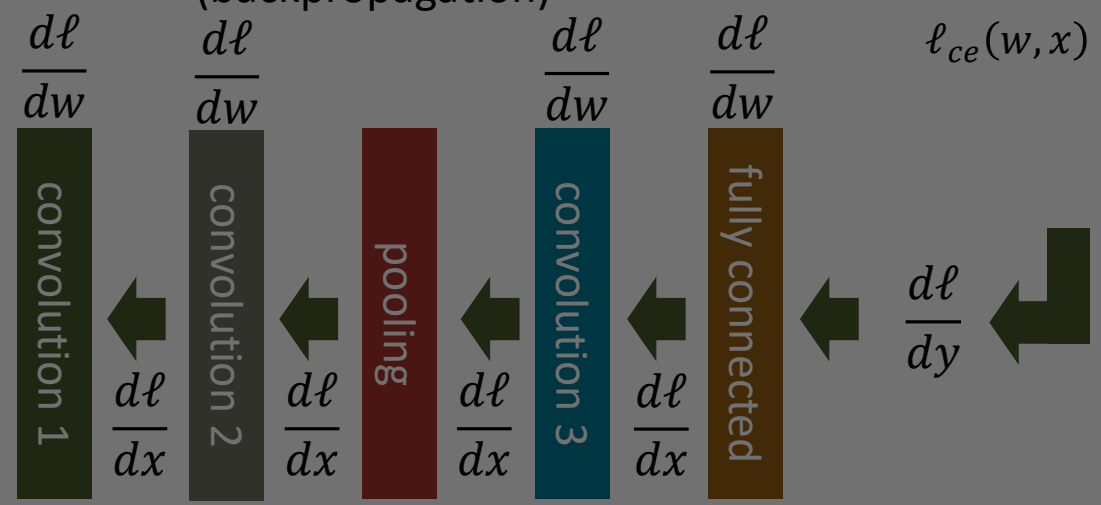
- “Backward propagation of errors” (Rumelhart, Hinton, & Williams 1986)



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

Backpropagation is just the chain rule!

layer-wise parameter update
(backpropagation)



$$l_{ce}(w, x) = - \sum_i l(x)_i \cdot \log \frac{e^{f(x)_i}}{\sum_k e^{f(x)_k}}$$

Stochastic Gradient Descent

```

1: for  $t = 0$  to  $T$  do
2:    $x \leftarrow$  Random element from  $X$ 
3:    $o_1 \leftarrow f_1(x)$ 
4:   for  $i = 2$  to layers do
5:      $o_i \leftarrow f_i(o_{i-1})$ 
6:   end for
7:   for  $i = \text{layers} - 1$  to  $1$  do
8:      $\nabla o_i \leftarrow \frac{\partial \ell}{\partial o_i}(o_{i-1}, o_i, \nabla o_{i+1})$ 
9:      $\nabla w_i^{(t)} \leftarrow \frac{\partial \ell}{\partial w_i}(o_{i-1}, o_i, \nabla o_{i+1})$ 
10:  end for
11:   $w^{(t+1)} \leftarrow w^{(t)} + u(\nabla w^{(t)}, w^{(0, \dots, t)}, t)$ 
12: end for

```

▶ Stopping condition
 ▶ Sample dataset
 ▶ Forward evaluation of ℓ
 ▶ Compute gradient of ℓ via backpropagation
 ▶ Gradient w.r.t. data
 ▶ Gradient w.r.t. layer parameters
 ▶ Update weights with function u

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$f_1(x)$ convolution 1

$f_2(f_1(x))$ convolution 2

pooling

...

convolution 3

$f(x)$ fully connected

Stochastic Gradient Descent

```

1: for  $t = 0$  to  $T$  do
2:    $x \leftarrow$  Random element from  $X$ 
3:    $o_1 \leftarrow f_1(x)$ 
4:   for  $i = 2$  to layers do
5:      $o_i \leftarrow f_i(o_{i-1})$ 
6:   end for
7:   for  $i = \text{layers} - 1$  to  $1$  do
8:      $\nabla o_i \leftarrow \frac{\partial \ell}{\partial o_i}(o_{i-1}, o_i, \nabla o_{i+1})$ 
9:      $\nabla w_i^{(t)} \leftarrow \frac{\partial \ell}{\partial w_i}(o_{i-1}, o_i, \nabla o_{i+1})$ 
10:  end for
11:   $w^{(t+1)} \leftarrow w^{(t)} + u(\nabla w^{(t)}, w^{(0, \dots, t)}, t)$ 
12: end for

```

▷ Stopping condition
 ▷ Sample dataset

▷ Forward evaluation of ℓ

▷ Compute gradient of ℓ via backpropagation
 ▷ Gradient w.r.t. data

▷ Gradient w.r.t. layer parameters

▷ Update weights with function u

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$f_1(x)$ convolution 1

$f_2(f_1(x))$ convolution 2

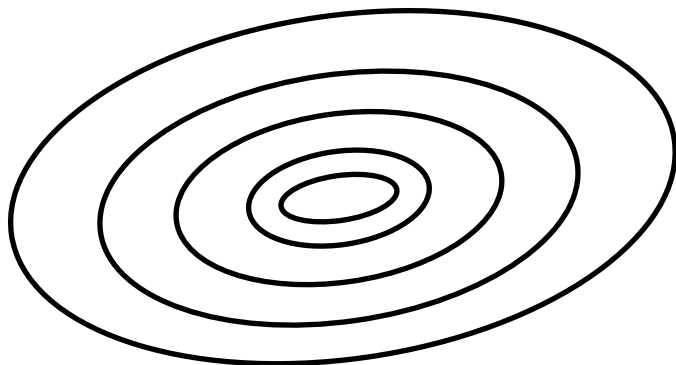
pooling

...

convolution 3

$f(x)$ fully connected

- Layer storage = $|w_l| + |f_l(o_{l-1})| + |\nabla w_l| + |\nabla o_l|$



Stochastic Gradient Descent

```

1: for  $t = 0$  to  $T$  do
2:    $x \leftarrow$  Random element from  $X$ 
3:    $o_1 \leftarrow f_1(x)$ 
4:   for  $i = 2$  to layers do
5:      $o_i \leftarrow f_i(o_{i-1})$ 
6:   end for
7:   for  $i = \text{layers} - 1$  to  $1$  do
8:      $\nabla o_i \leftarrow \frac{\partial \ell}{\partial o_i}(o_{i-1}, o_i, \nabla o_{i+1})$ 
9:      $\nabla w_i^{(t)} \leftarrow \frac{\partial \ell}{\partial w_i}(o_{i-1}, o_i, \nabla o_{i+1})$ 
10:  end for
11:   $w^{(t+1)} \leftarrow w^{(t)} + u(\nabla w^{(t)}, w^{(0, \dots, t)}, t)$ 
12: end for
  
```

▷ Stopping condition
 ▷ Sample dataset

▷ Forward evaluation of ℓ

▷ Compute gradient of ℓ via backpropagation
 ▷ Gradient w.r.t. data

▷ Gradient w.r.t. layer parameters

▷ Update weights with function u

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

$f_1(x)$ convolution 1

$f_2(f_1(x))$ convolution 2

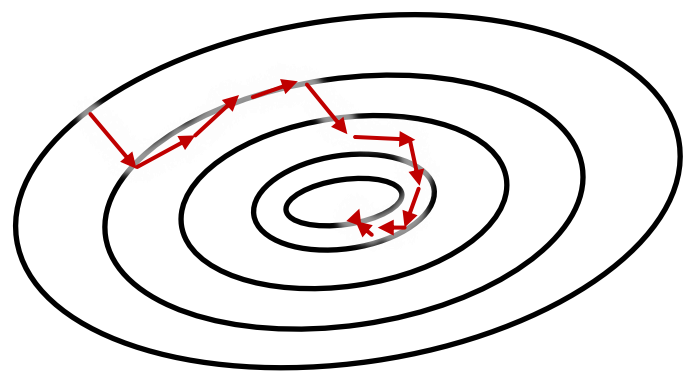
pooling

...

convolution 3

$f(x)$ fully connected

- Layer storage = $|w_l| + |f_l(o_{l-1})| + |\nabla w_l| + |\nabla o_l|$



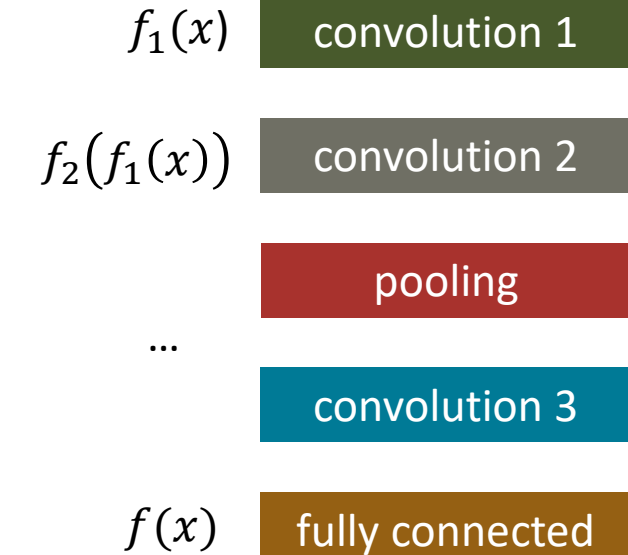
Stochastic Gradient Descent

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

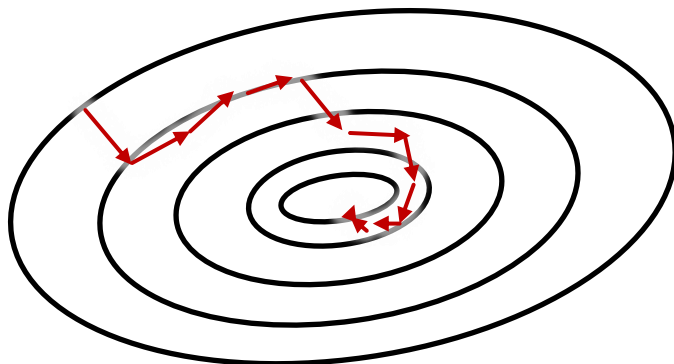
```

1: for t = 0 to T do
2:   x ← Random element from X
3:   o1 ← f1(x)
4:   for i = 2 to layers do
5:     oi ← fi(oi-1)
6:   end for
7:   for i = layers-1 to 1 do
8:     ∇oi ← ∂ℓ/∂oi(oi-1, oi, ∇oi+1)
9:     ∇wi(t) ← ∂ℓ/∂wi(oi-1, oi, ∇oi+1)
10:  end for
11:  w(t+1) ← w(t) + u(∇w(t), w(0, ..., t), t)
12: end for
    
```

▶ Stopping condition
 ▶ Sample dataset
 ▶ Forward evaluation of ℓ
 ▶ Compute gradient of ℓ via backpropagation
 ▶ Gradient w.r.t. data
 ▶ Gradient w.r.t. layer parameters
 ▶ Update weights with function u



- Layer storage = $|w_l| + |f_l(o_{l-1})| + |\nabla w_l| + |\nabla o_l|$



Learning Rate	$w^{(t+1)} = w^{(t)} - \eta \cdot \nabla \ell(w^{(t)}, z) = w^{(t)} - \eta \cdot \nabla w^{(t)}$
Adaptive Learning Rate	$w^{(t+1)} = w^{(t)} - \eta_t \cdot \nabla w^{(t)}$
Momentum [Qian 1999]	$w^{(t+1)} = w^{(t)} + \mu \cdot (w^{(t)} - w^{(t-1)}) - \eta \cdot \nabla w^{(t)}$
Nesterov Momentum [Nesterov 1983]	$w^{(t+1)} = w^{(t)} + v_t; \quad v_{t+1} = \mu \cdot v_t - \eta \cdot \nabla \ell(w^{(t)} - \mu \cdot v_t, z)$
AdaGrad [Duchi et al. 2011]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot \nabla w_i^{(t)}}{\sqrt{A_{i,t} + \epsilon}}; \quad A_{i,t} = \sum_{\tau=0}^t (\nabla w_i^{(\tau)})^2$
RMSProp [Hinton 2012]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot \nabla w_i^{(t)}}{\sqrt{A'_{i,t} + \epsilon}}; \quad A'_{i,t} = \beta \cdot A'_{i,t-1} + (1 - \beta) (\nabla w_i^{(t)})^2$
Adam [Kingma and Ba 2015]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot M_{i,t}^{(1)}}{\sqrt{M_{i,t}^{(2)} + \epsilon}}; \quad M_{i,t}^{(m)} = \frac{\beta^m \cdot M_{i,t-1}^{(m)} + (1 - \beta^m) (\nabla w_i^{(t)})^m}{1 - \beta^m}$

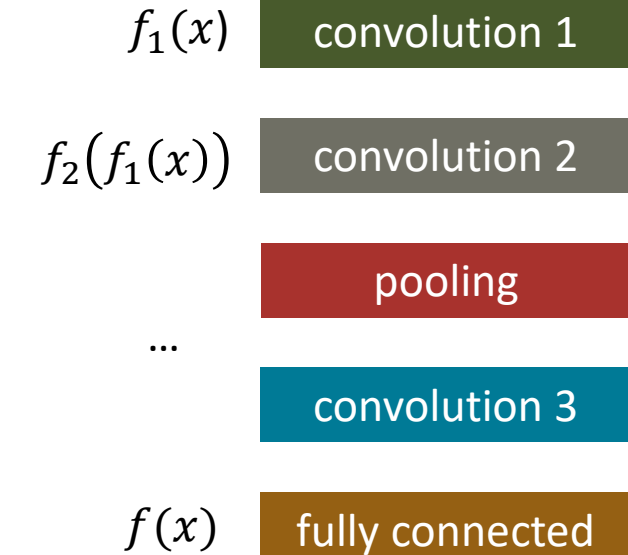
Stochastic Gradient Descent

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^d} \mathbb{E}_{x \sim \mathcal{D}} [\ell(w, x)]$$

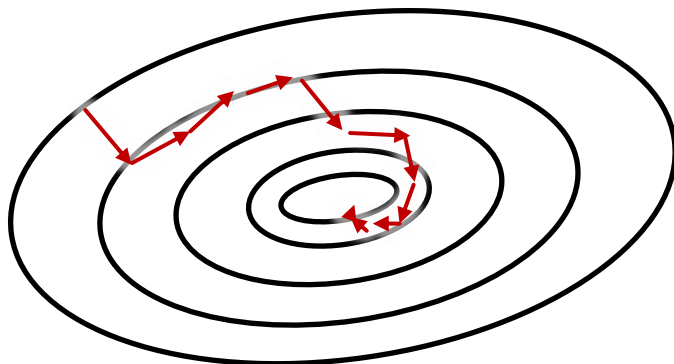
```

1: for t = 0 to T do
2:   x ← Random element from X
3:   o1 ← f1(x)
4:   for i = 2 to layers do
5:     oi ← fi(oi-1)
6:   end for
7:   for i = layers-1 to 1 do
8:     ∇oi ← ∂ℓ/∂oi(oi-1, oi, ∇oi+1)
9:     ∇wi(t) ← ∂ℓ/∂wi(oi-1, oi, ∇oi+1)
10:  end for
11:  w(t+1) ← w(t) + u(∇w(t), w(0, ..., t), t)
12: end for
    
```

▶ Stopping condition
 ▶ Sample dataset
 ▶ Forward evaluation of ℓ
 ▶ Compute gradient of ℓ via backpropagation
 ▶ Gradient w.r.t. data
 ▶ Gradient w.r.t. layer parameters
 ▶ Update weights with function u

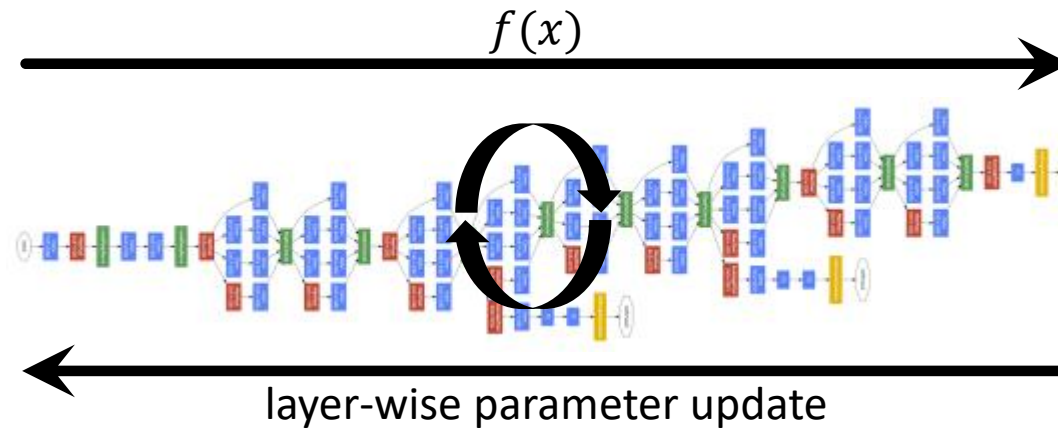


- Layer storage = $|w_l| + |f_l(o_{l-1})| + |\nabla w_l| + |\nabla o_l|$



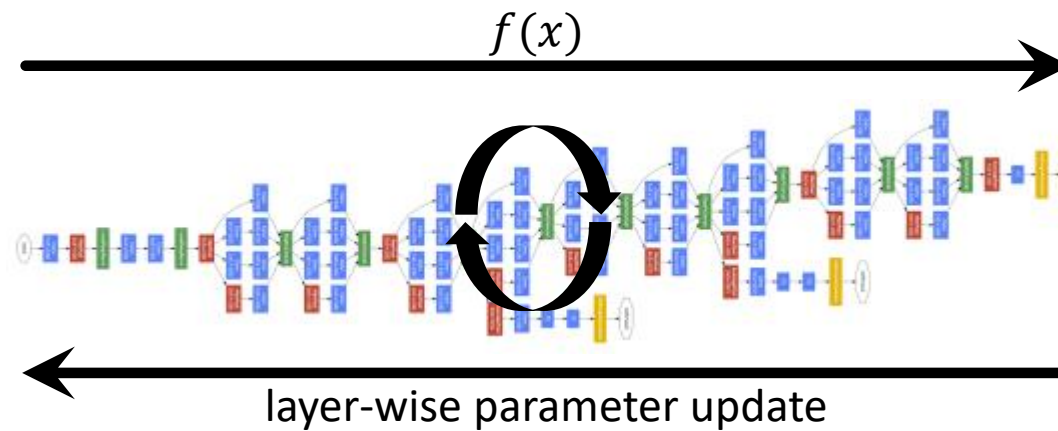
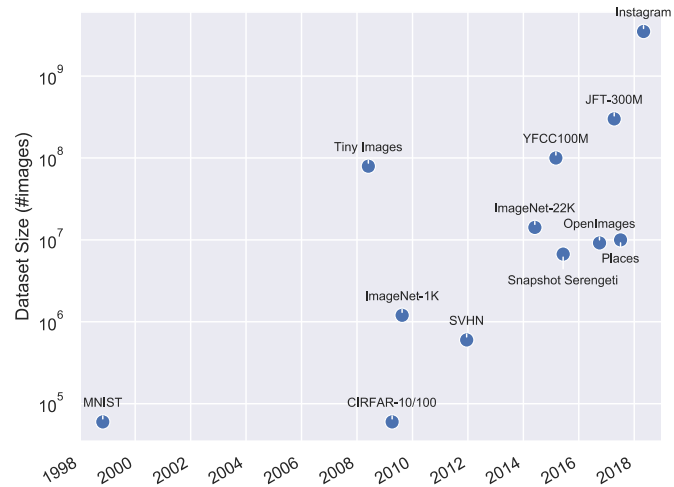
Learning Rate	$w^{(t+1)} = w^{(t)} - \eta \cdot \nabla \ell(w^{(t)}, z) = w^{(t)} - \eta \cdot \nabla w^{(t)}$
Adaptive Learning Rate	$w^{(t+1)} = w^{(t)} - \eta_t \cdot \nabla w^{(t)}$
Momentum [Qian 1999]	$w^{(t+1)} = w^{(t)} + \mu \cdot (w^{(t)} - w^{(t-1)}) - \eta \cdot \nabla w^{(t)}$
Nesterov Momentum [Nesterov 1983]	$w^{(t+1)} = w^{(t)} + v_t; \quad v_{t+1} = \mu \cdot v_t - \eta \cdot \nabla \ell(w^{(t)} - \mu \cdot v_t, z)$
AdaGrad [Duchi et al. 2011]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot \nabla w_i^{(t)}}{\sqrt{A_{i,t} + \epsilon}}; \quad A_{i,t} = \sum_{\tau=0}^t (\nabla w_i^{(\tau)})^2$
RMSProp [Hinton 2012]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot \nabla w_i^{(t)}}{\sqrt{A'_{i,t} + \epsilon}}; \quad A'_{i,t} = \beta \cdot A'_{i,t-1} + (1 - \beta) (\nabla w_i^{(t)})^2$
Adam [Kingma and Ba 2015]	$w_i^{(t+1)} = w_i^{(t)} - \frac{\eta \cdot M_{i,t}^{(1)}}{\sqrt{M_{i,t}^{(2)} + \epsilon}}; \quad M_{i,t}^{(m)} = \frac{\beta^m \cdot M_{i,t-1}^{(m)} + (1 - \beta^m) (\nabla w_i^{(t)})^m}{1 - \beta^m}$

The scale of deep learning



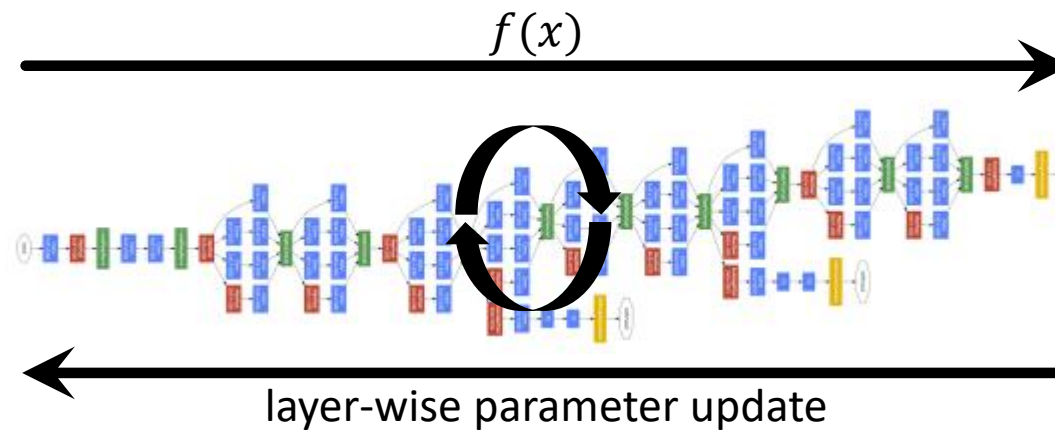
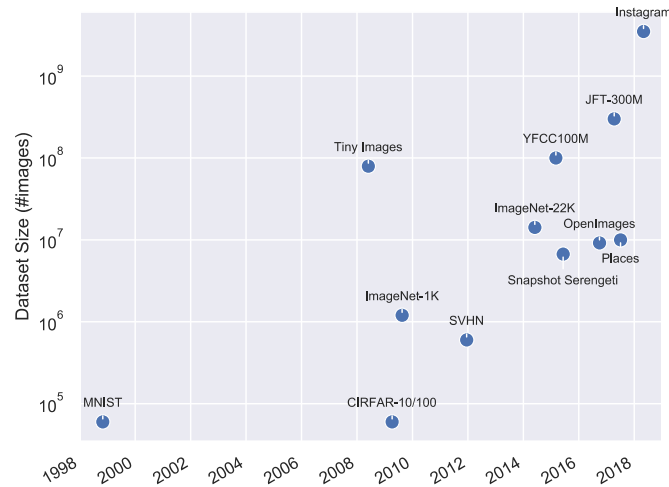
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

The scale of deep learning



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

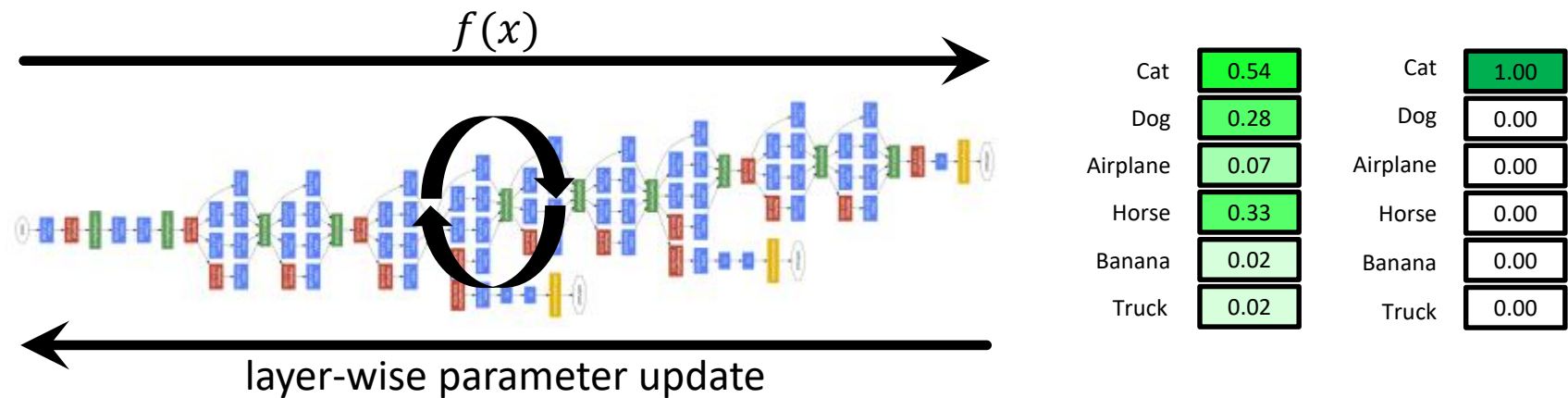
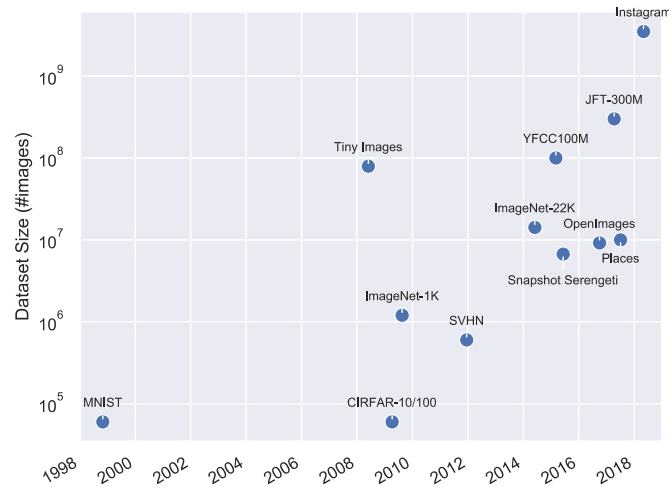
The scale of deep learning



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

- ImageNet 1k: 150 GB
- ImageNet 22k: ~2 TB
- Industry: Much larger

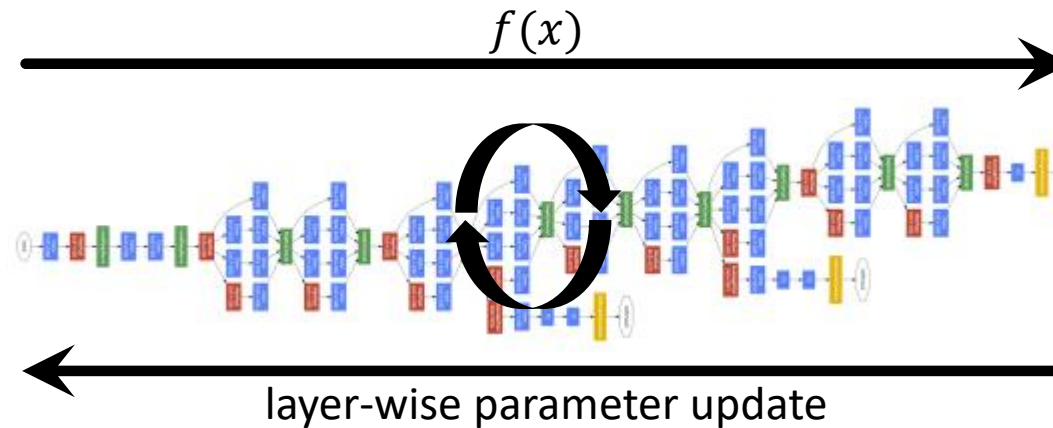
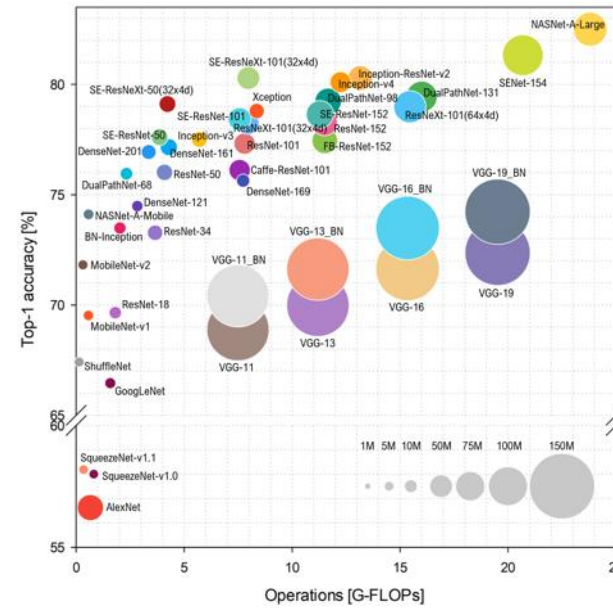
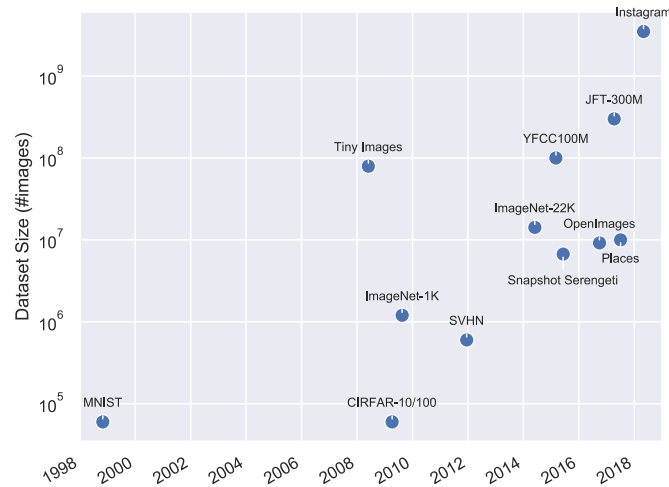
The scale of deep learning



- ImageNet 1k: 150 GB
- ImageNet 22k: ~2 TB
- Industry: Much larger

- >100 layers deep
- ~25M – >10B parameters
- 0.1 – 40 GiB storage

The scale of deep learning

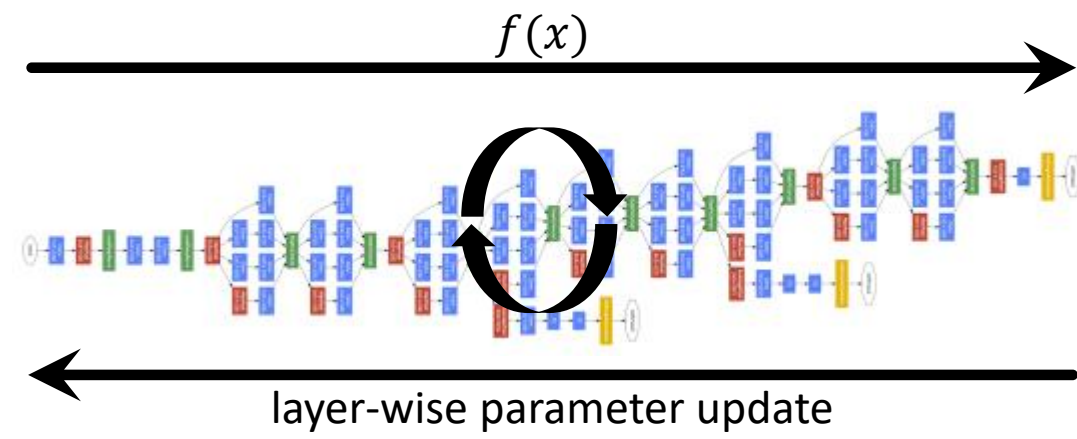
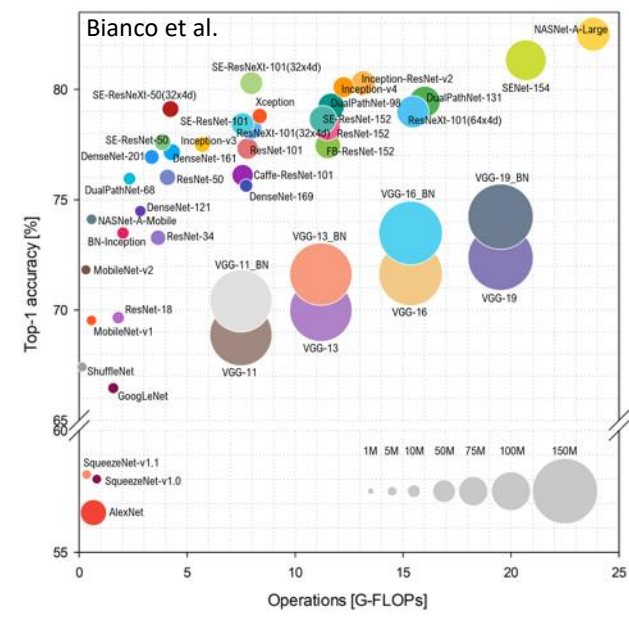
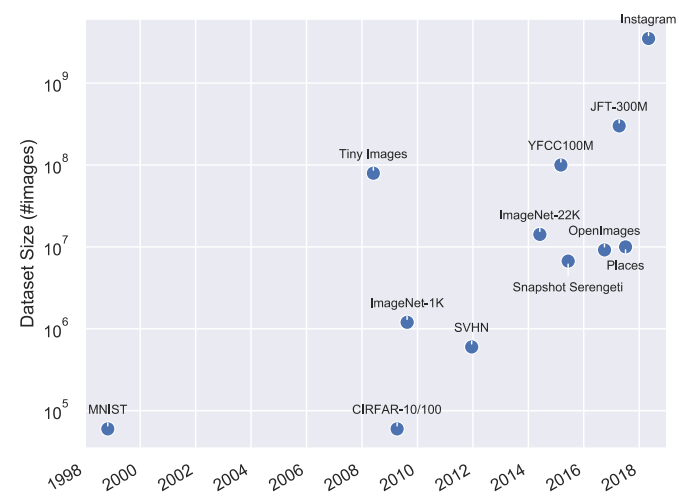


Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

- ImageNet 1k: 150 GB
- ImageNet 22k: ~2 TB
- Industry: Much larger

- >100 layers deep
- ~25M – >10B parameters
- 0.1 – 40 GiB storage

The scale of deep learning

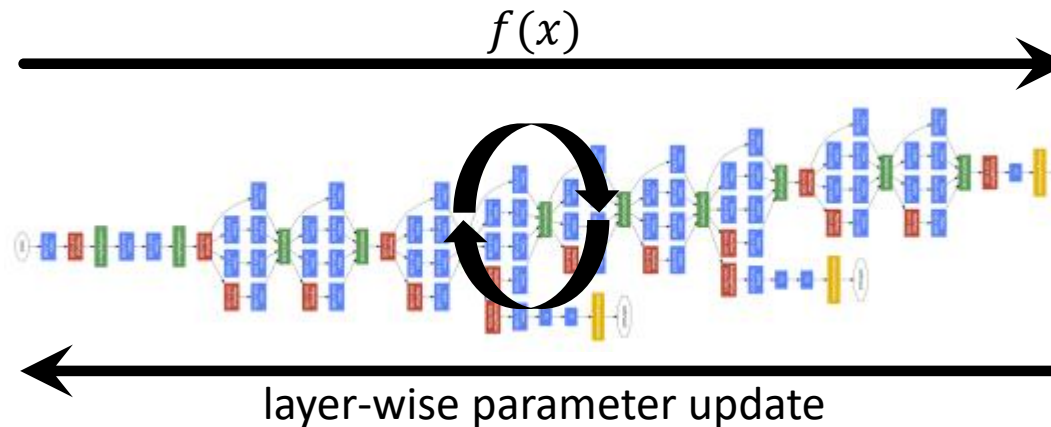
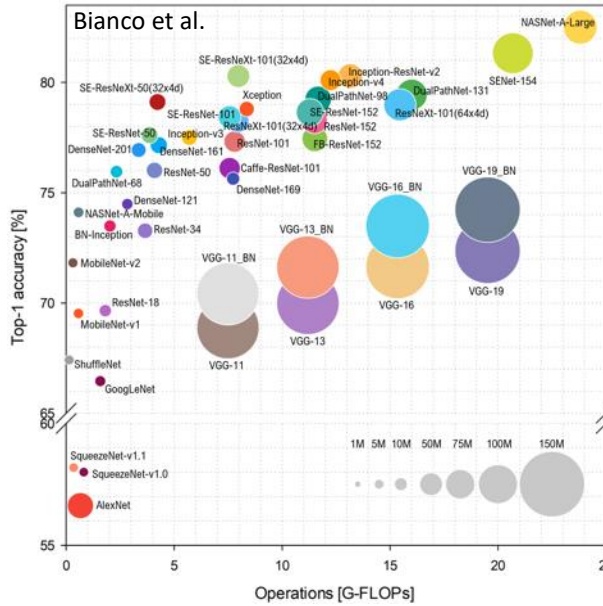
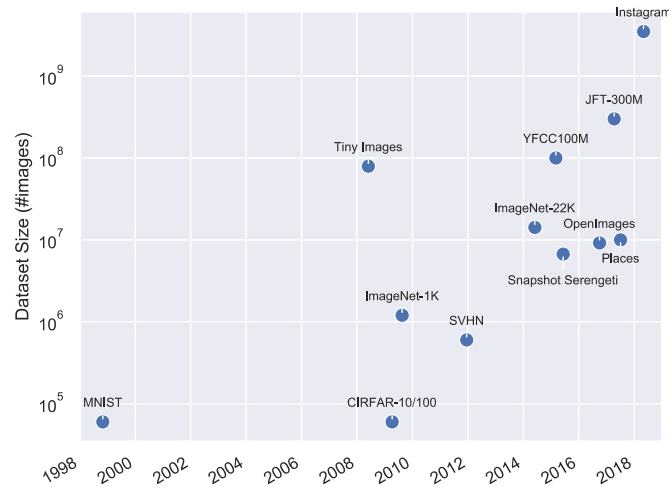


Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

- ImageNet 1k: 150 GB
- ImageNet 22k: ~2 TB
- Industry: Much larger

- >100 layers deep
- ~25M – >10B parameters
- 0.1 – 40 GiB storage

The scale of deep learning



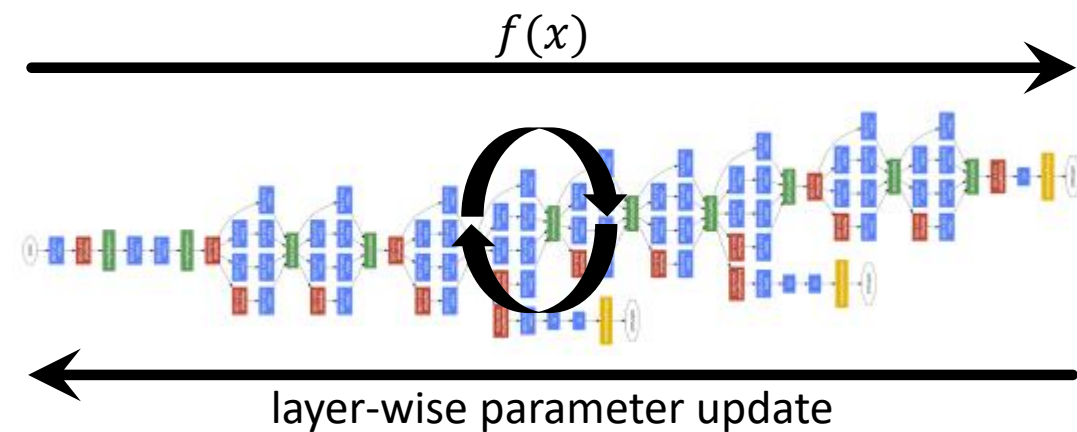
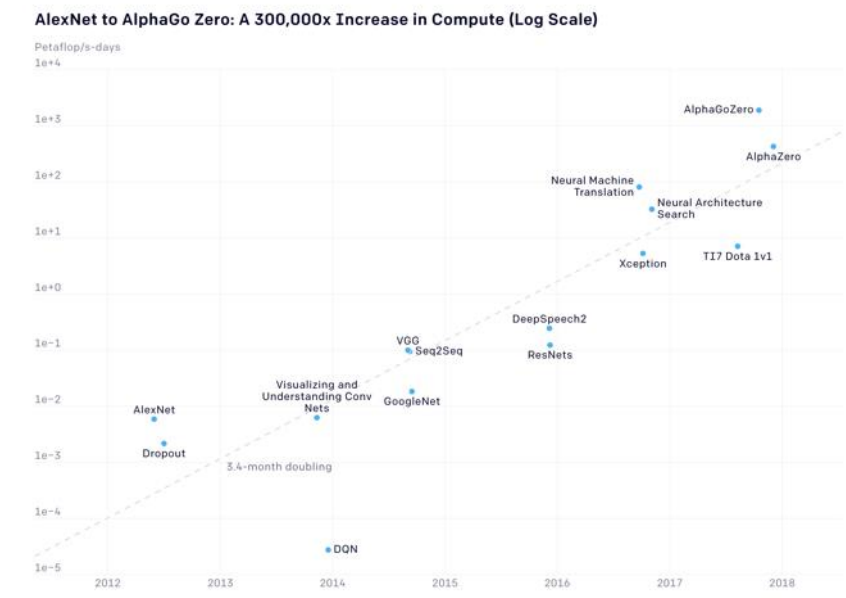
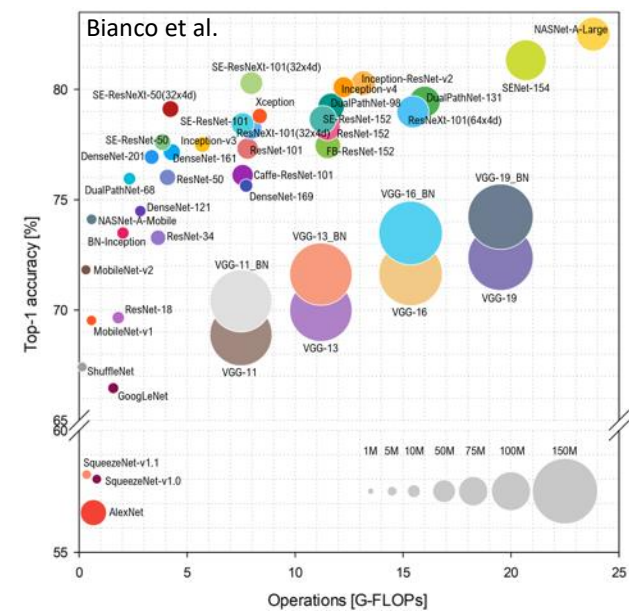
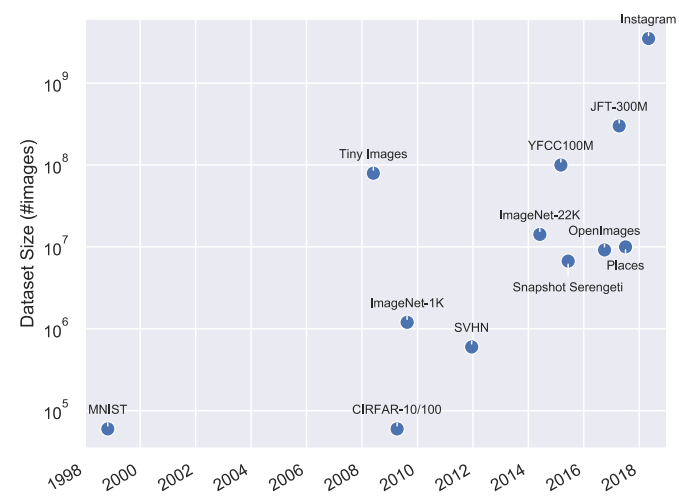
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

- ImageNet 1k: 150 GB
- ImageNet 22k: ~2 TB
- Industry: Much larger

- >100 layers deep
- ~25M – >10B parameters
- 0.1 – 40 GiB storage

- 10-22k labels
- Growing
- Weeks to train

The scale of deep learning



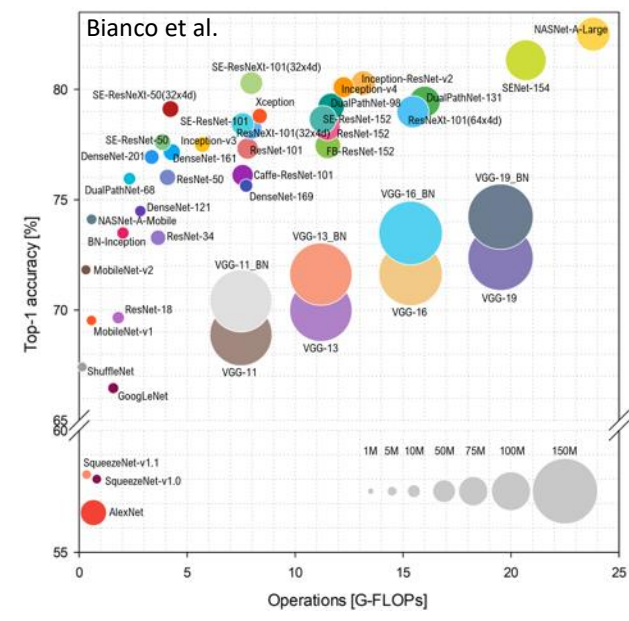
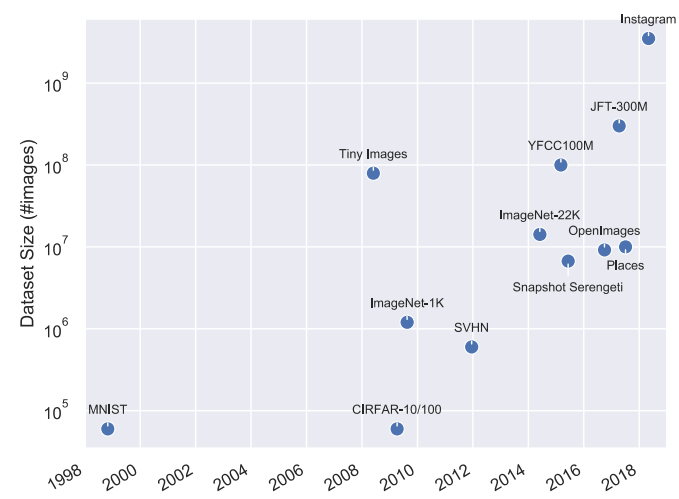
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

- ImageNet 1k: 150 GB
- ImageNet 22k: ~2 TB
- Industry: Much larger

- >100 layers deep
- ~25M – >10B parameters
- 0.1 – 40 GiB storage

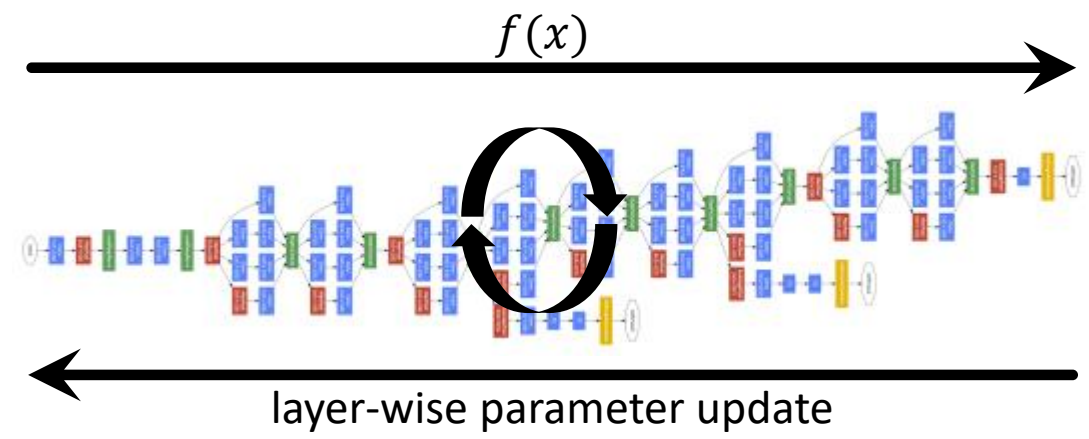
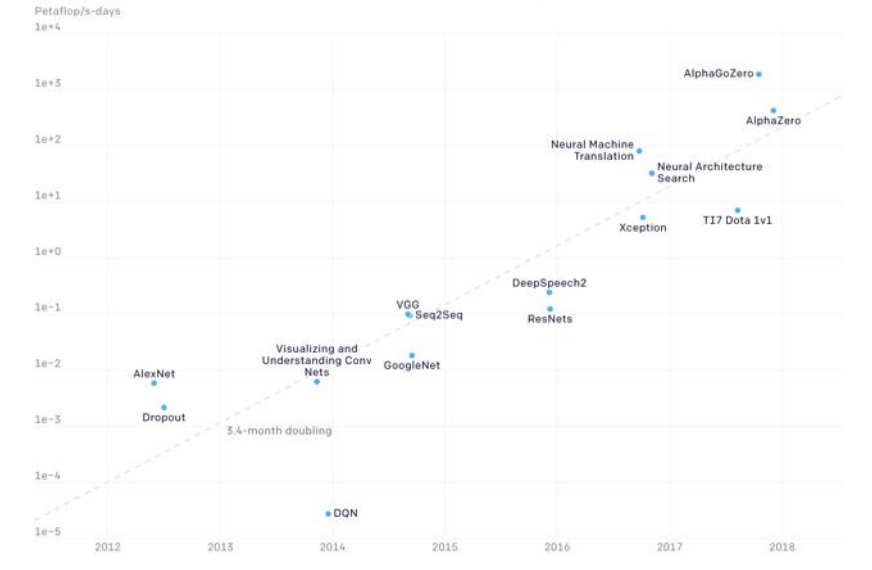
- 10-22k labels
- Growing
- Weeks to train

The scale of deep learning



OpenAI

AlexNet to AlphaGo Zero: A 300,000x Increase in Compute (Log Scale)



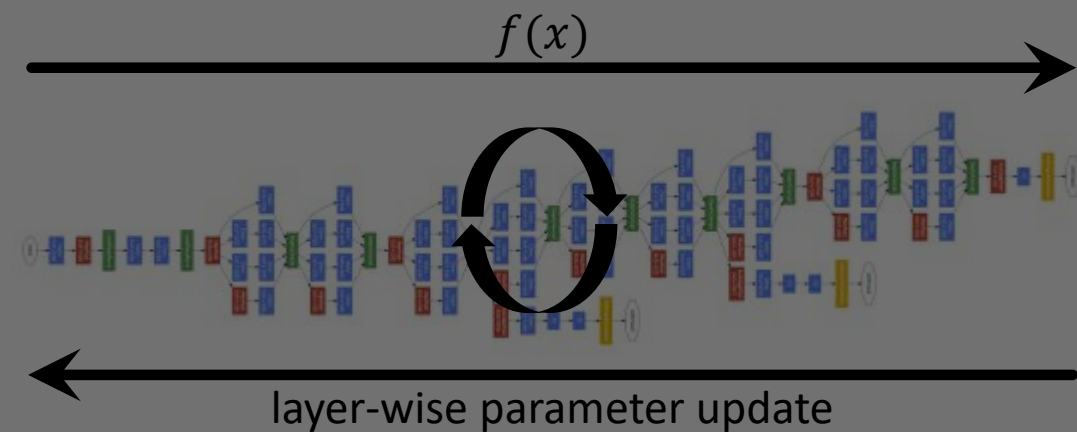
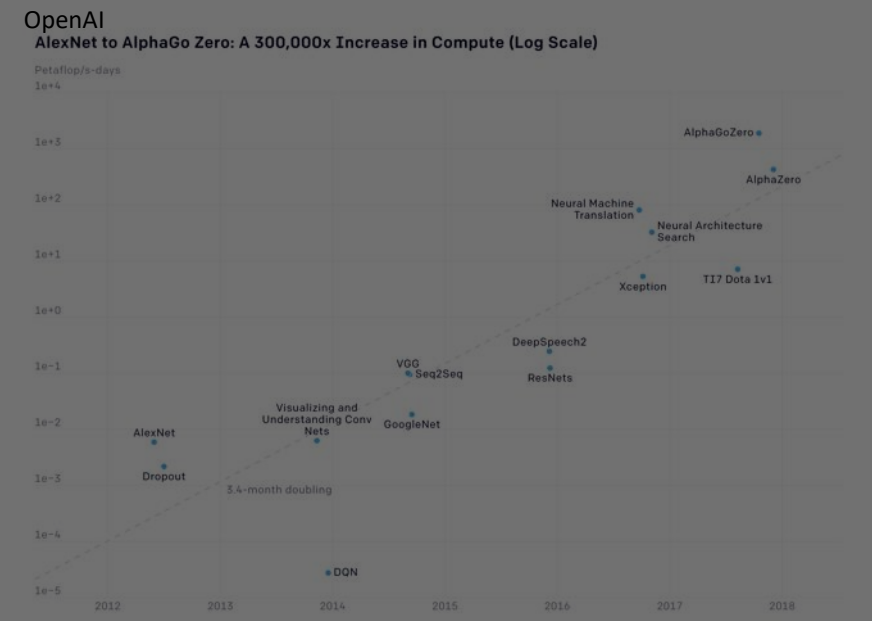
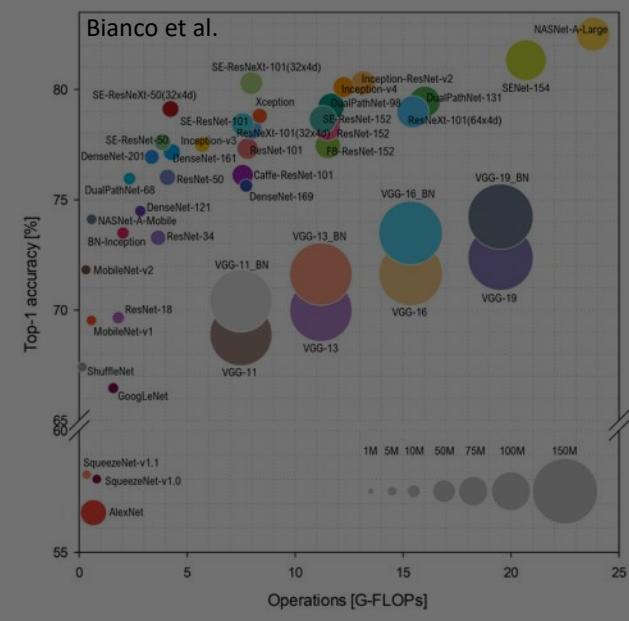
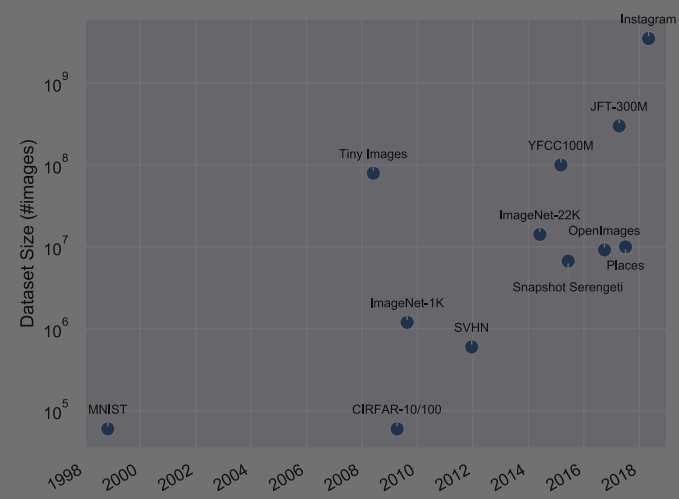
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

- ImageNet 1k: 150 GB
- ImageNet 22k: ~2 TB
- Industry: Much larger

- >100 layers deep
- ~25M – >10B parameters
- 0.1 – 40 GiB storage

- 10-22k labels
- Growing
- Weeks to train

The scale of deep learning



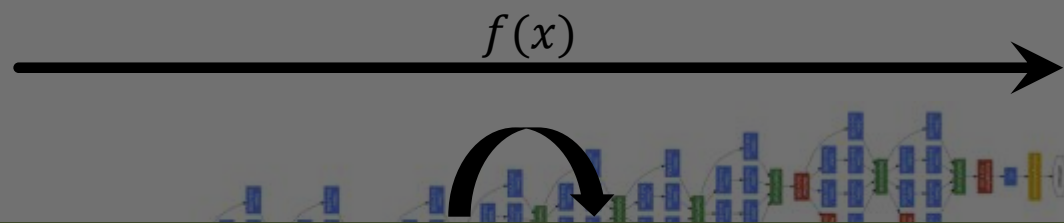
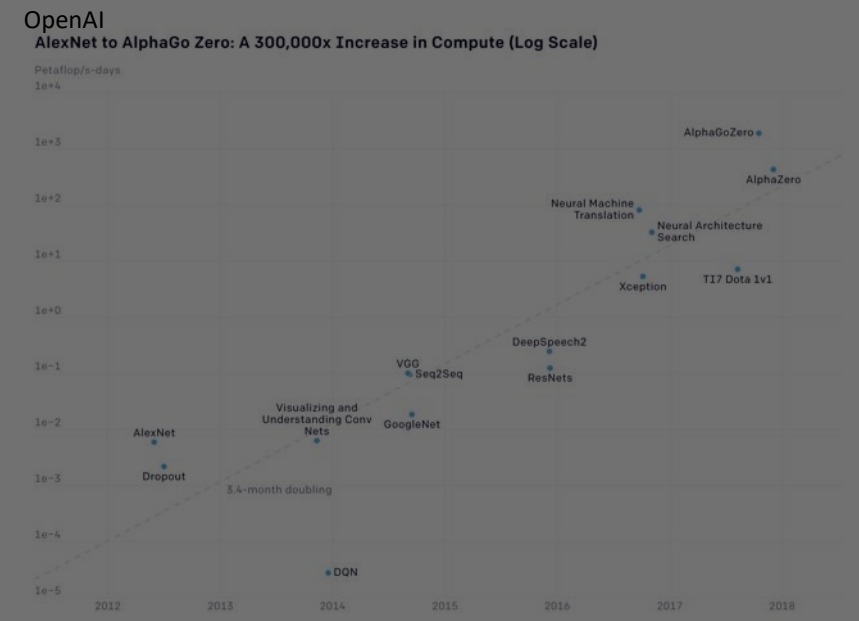
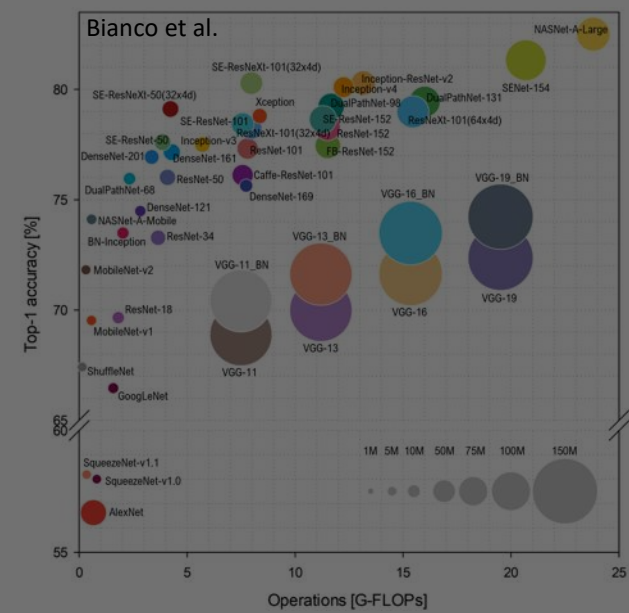
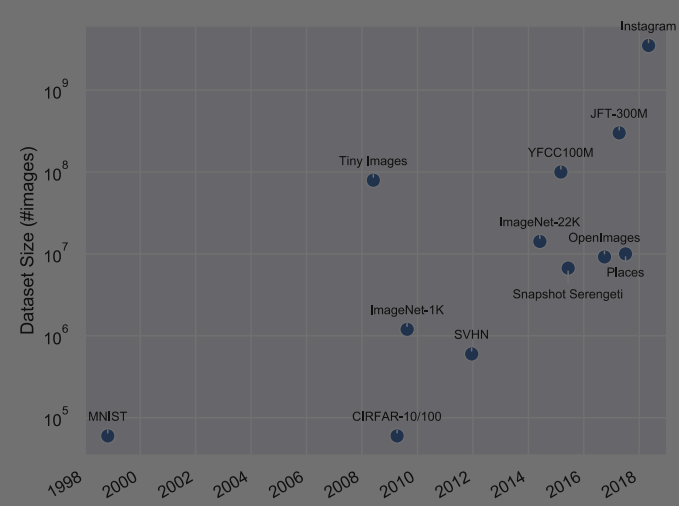
Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.33	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

- ImageNet 1k: 150 GB
- ImageNet 22k: ~2 TB
- Industry: Much larger

- >100 layers deep
- ~25M – >10B parameters
- 0.1 – 40 GiB storage

- 10-22k labels
- Growing
- Weeks to train

The scale of deep learning



Cat	0.54	Cat	1.00
Dog	0.28	Dog	0.00
Airplane	0.07	Airplane	0.00
Horse	0.05	Horse	0.00
Banana	0.02	Banana	0.00
Truck	0.02	Truck	0.00

Deep Learning is Supercomputing!

layer-wise parameter update

- ImageNet 1k: 150 GB
- ImageNet 22k: ~2 TB
- Industry: Much larger

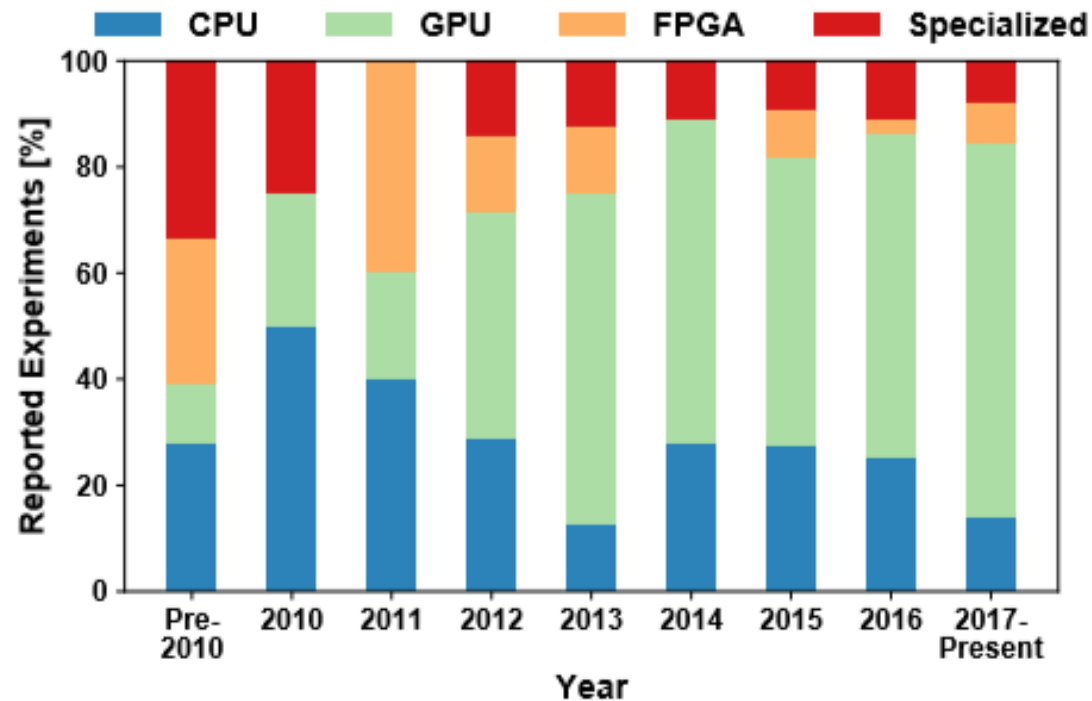
- >100 layers deep
- ~25M – >10B parameters
- 0.1 – 40 GiB storage

- 10-22k labels
- Growing
- Weeks to train

Trends in deep learning: hardware and multi-node

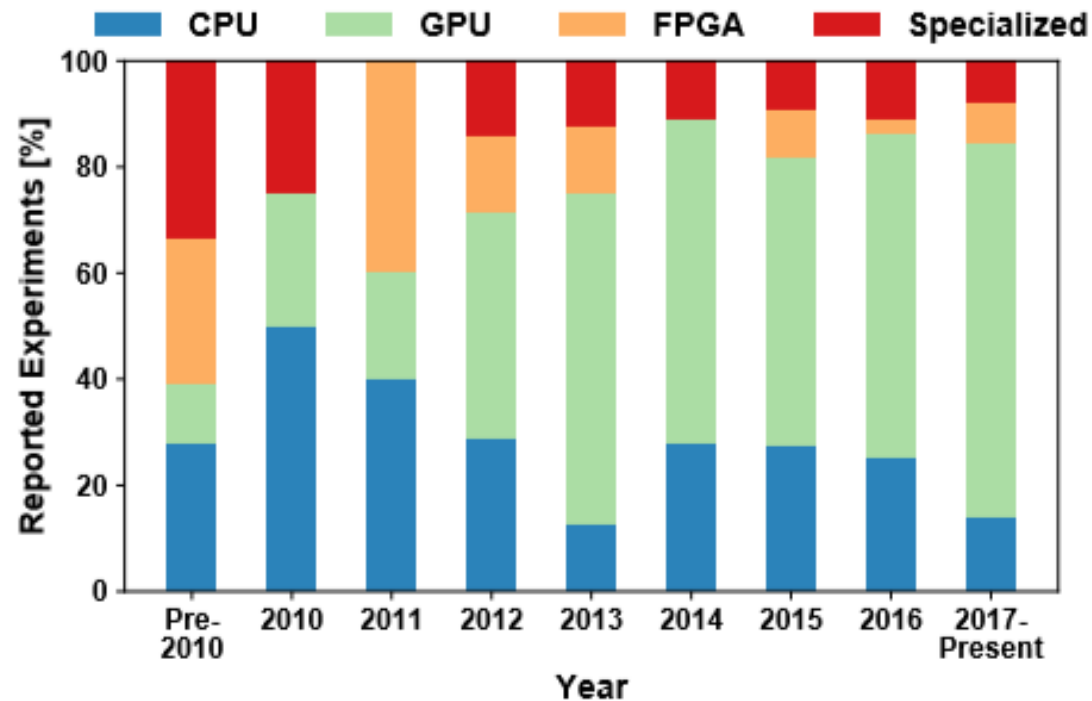
Trends in deep learning: hardware and multi-node

The field is moving fast – trying everything imaginable – survey results from 252 papers in the area of parallel deep learning



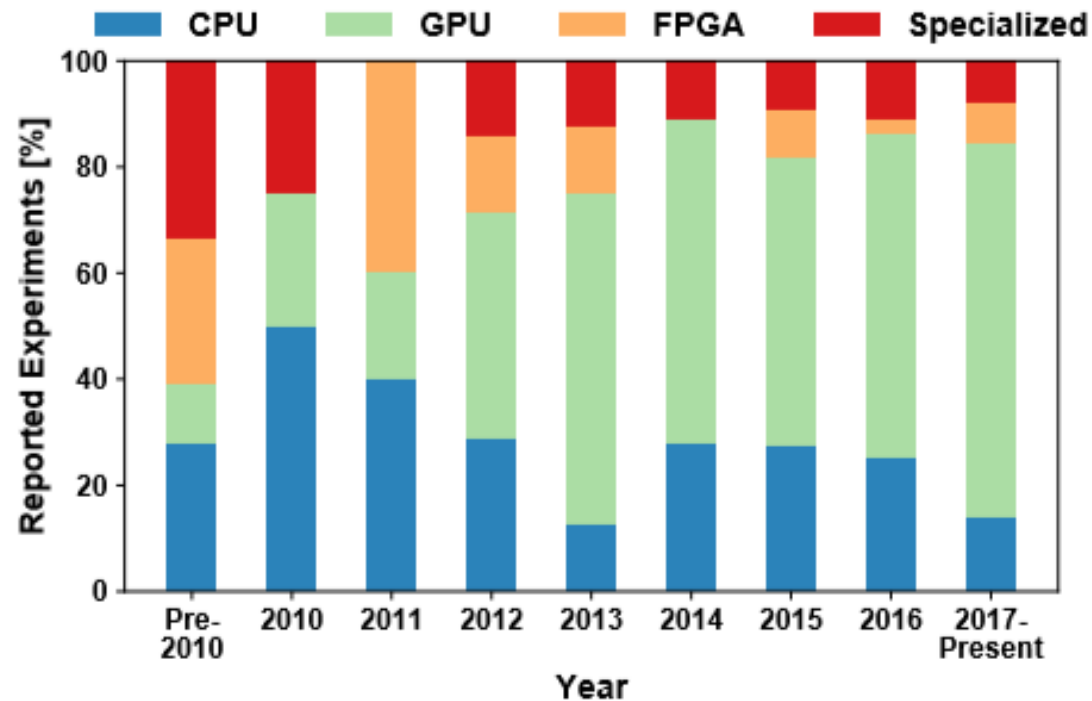
Trends in deep learning: hardware and multi-node

The field is moving fast – trying everything imaginable – survey results from 252 papers in the area of parallel deep learning

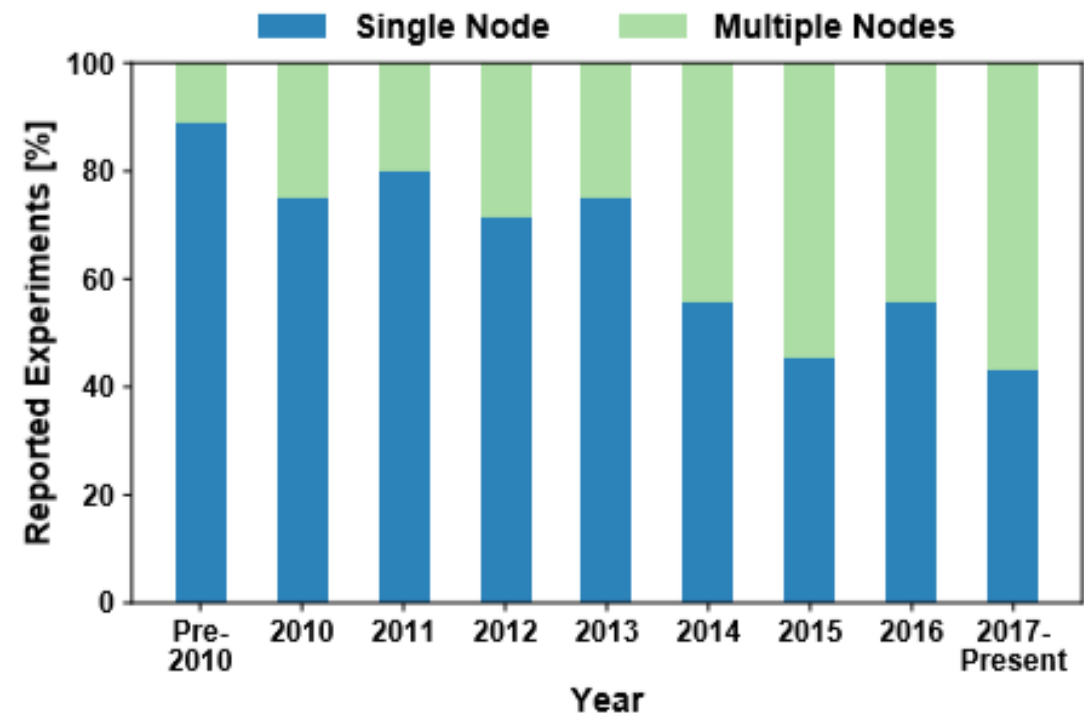


Trends in deep learning: hardware and multi-node

The field is moving fast – trying everything imaginable – survey results from 252 papers in the area of parallel deep learning

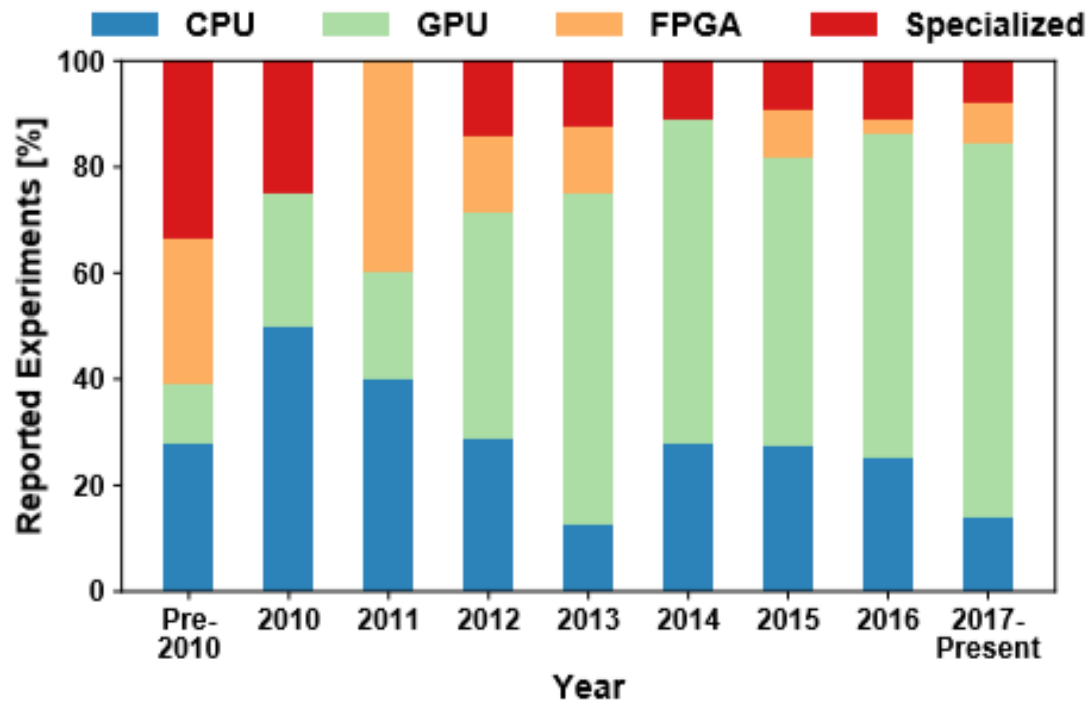


Hardware used

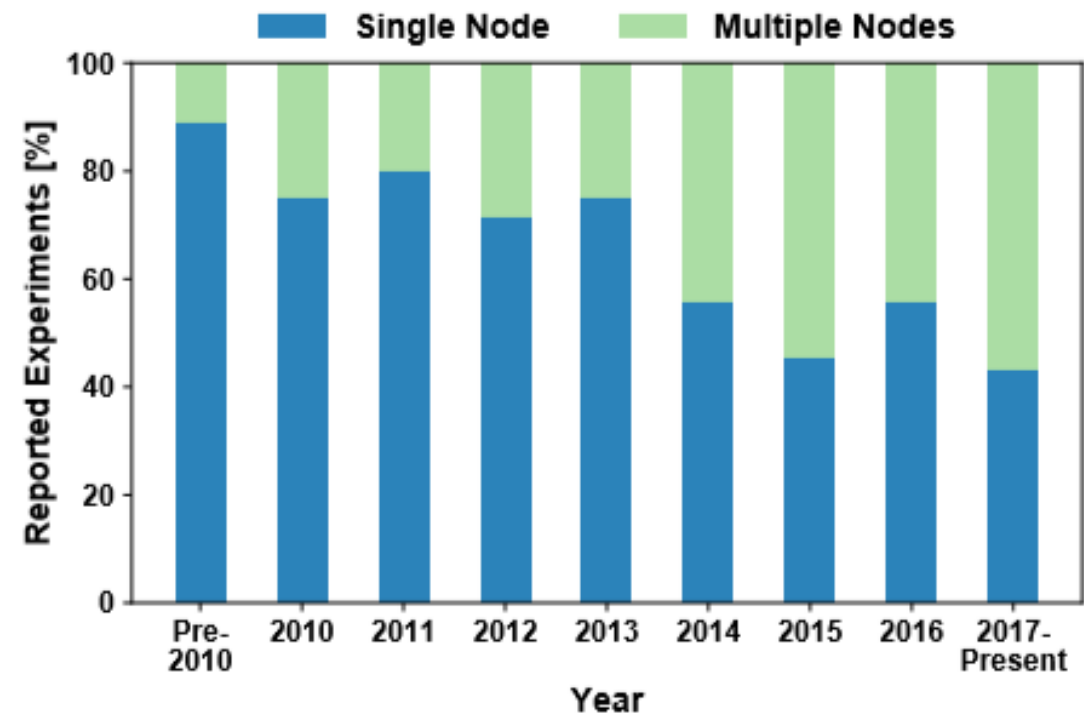


Trends in deep learning: hardware and multi-node

The field is moving fast – trying everything imaginable – survey results from 252 papers in the area of parallel deep learning

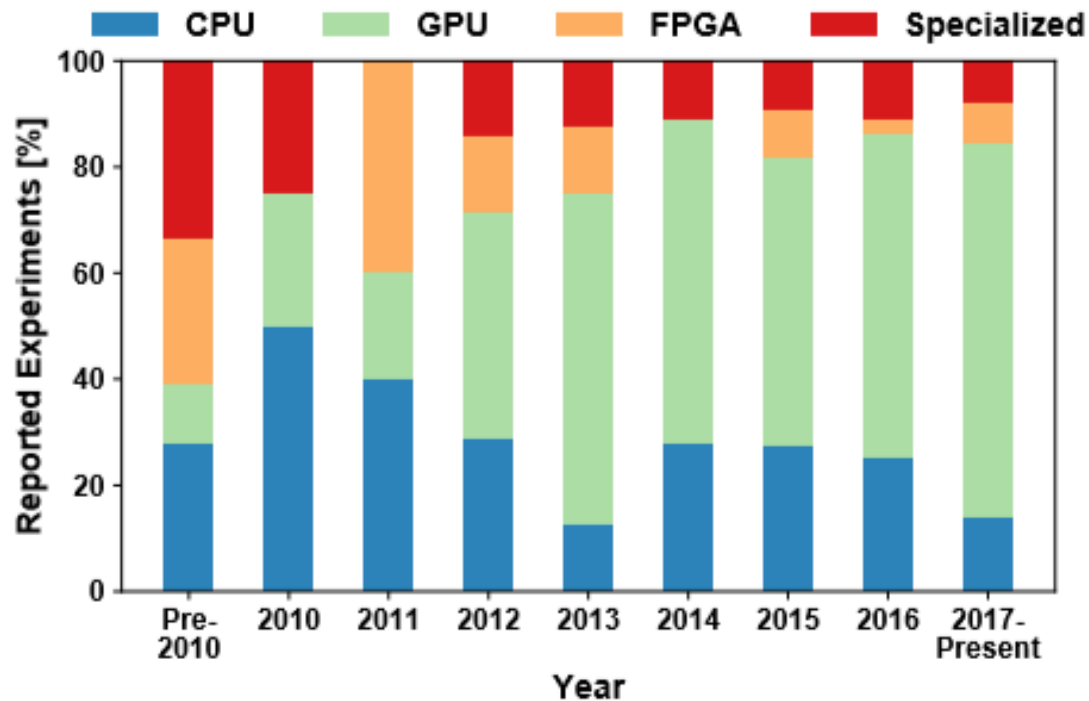


Hardware used

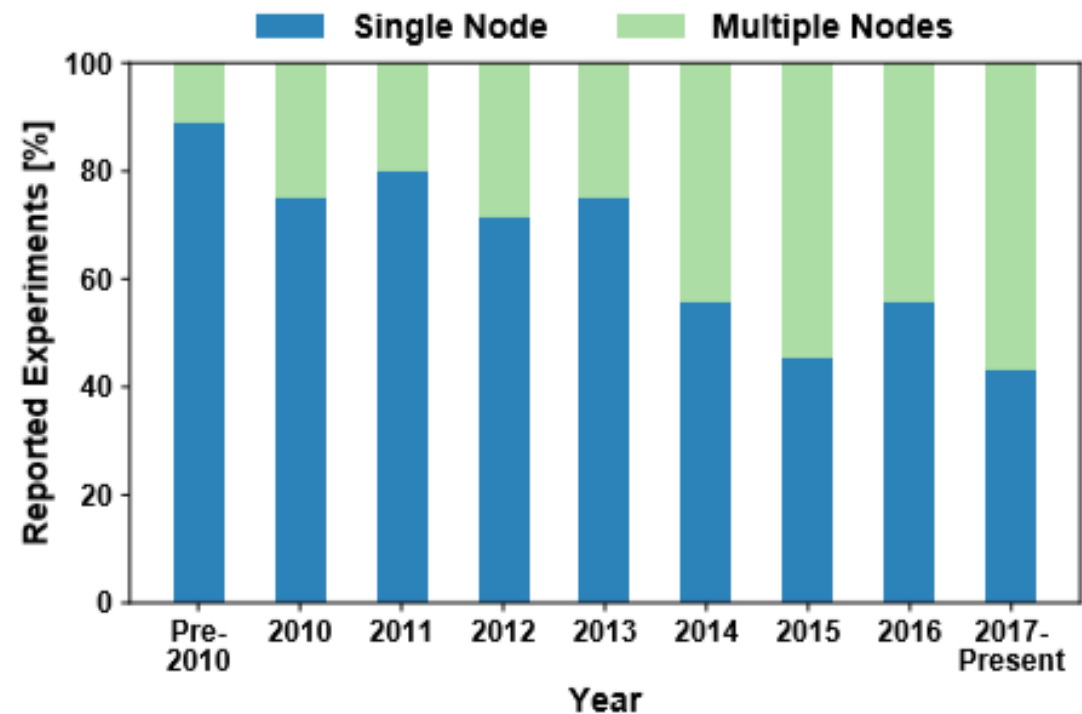


Trends in deep learning: hardware and multi-node

The field is moving fast – trying everything imaginable – survey results from 252 papers in the area of parallel deep learning



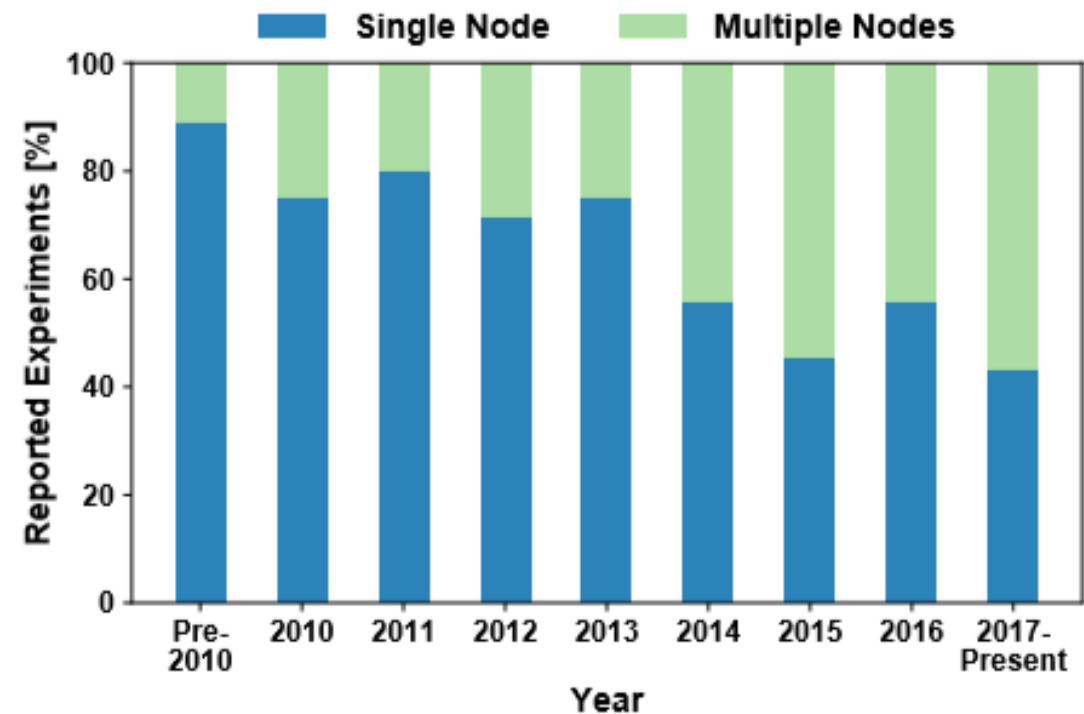
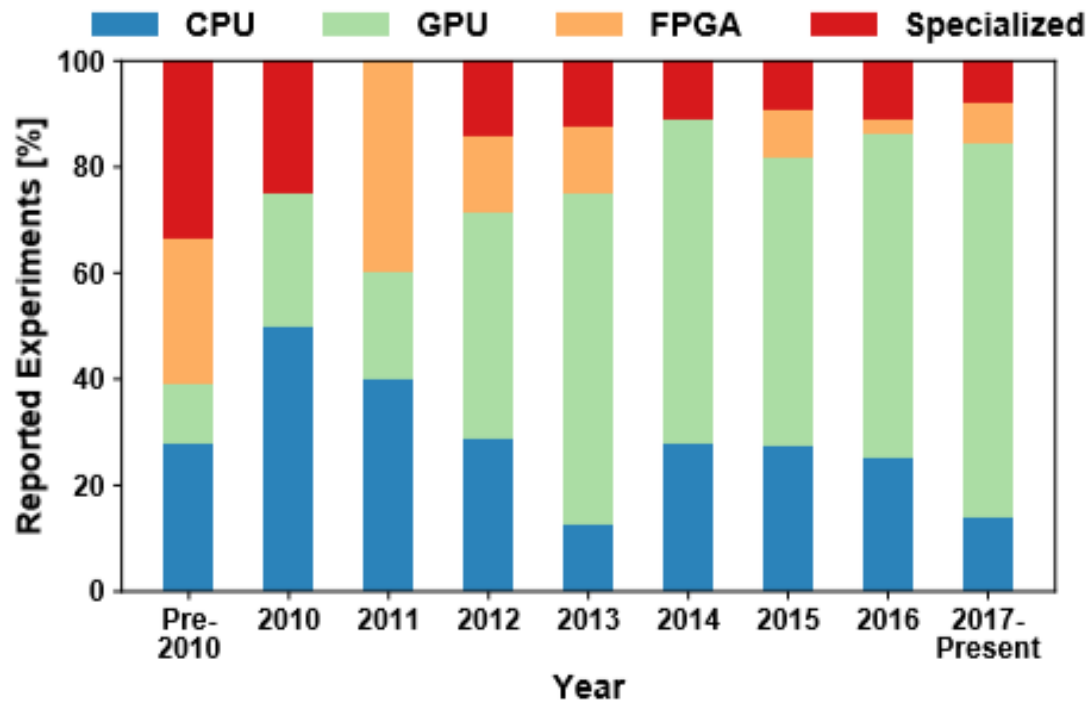
Hardware used



Shared vs. distributed memory

Trends in deep learning: hardware and multi-node

The field is moving fast – trying everything imaginable – survey results from 252 papers in the area of parallel deep learning



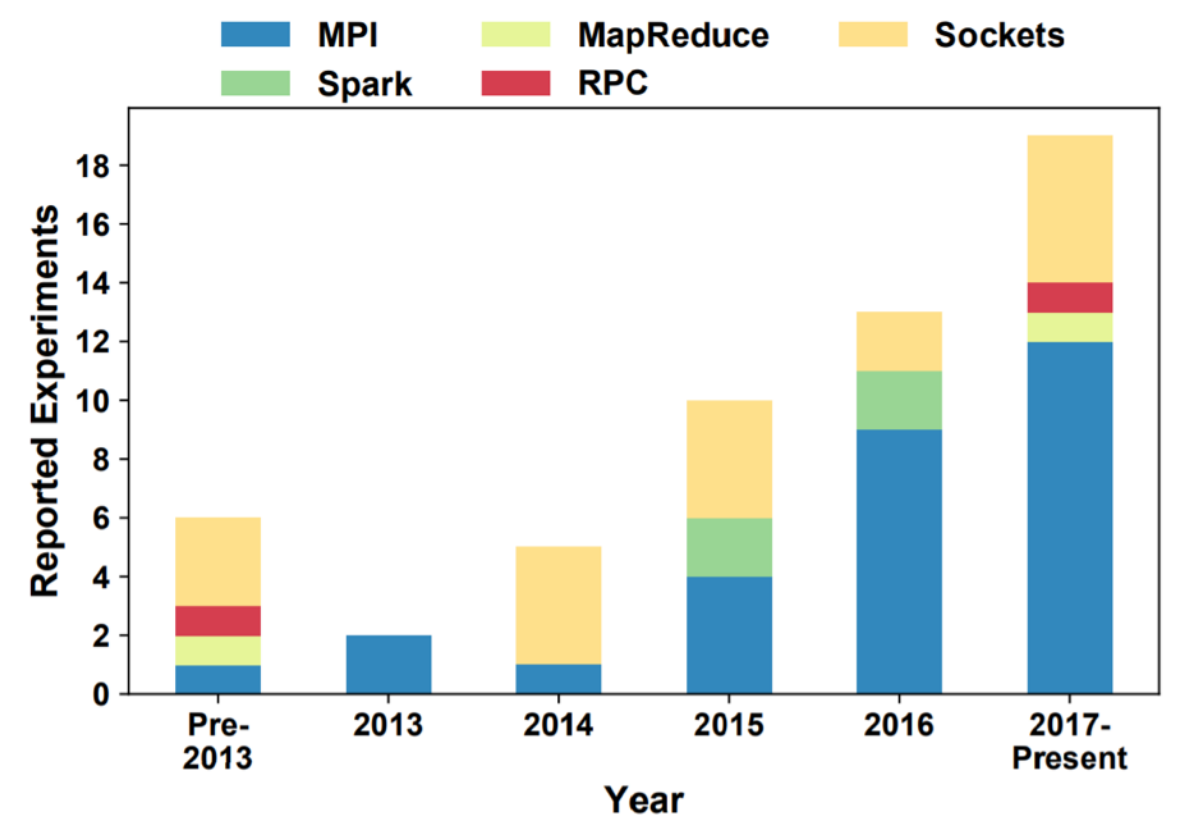
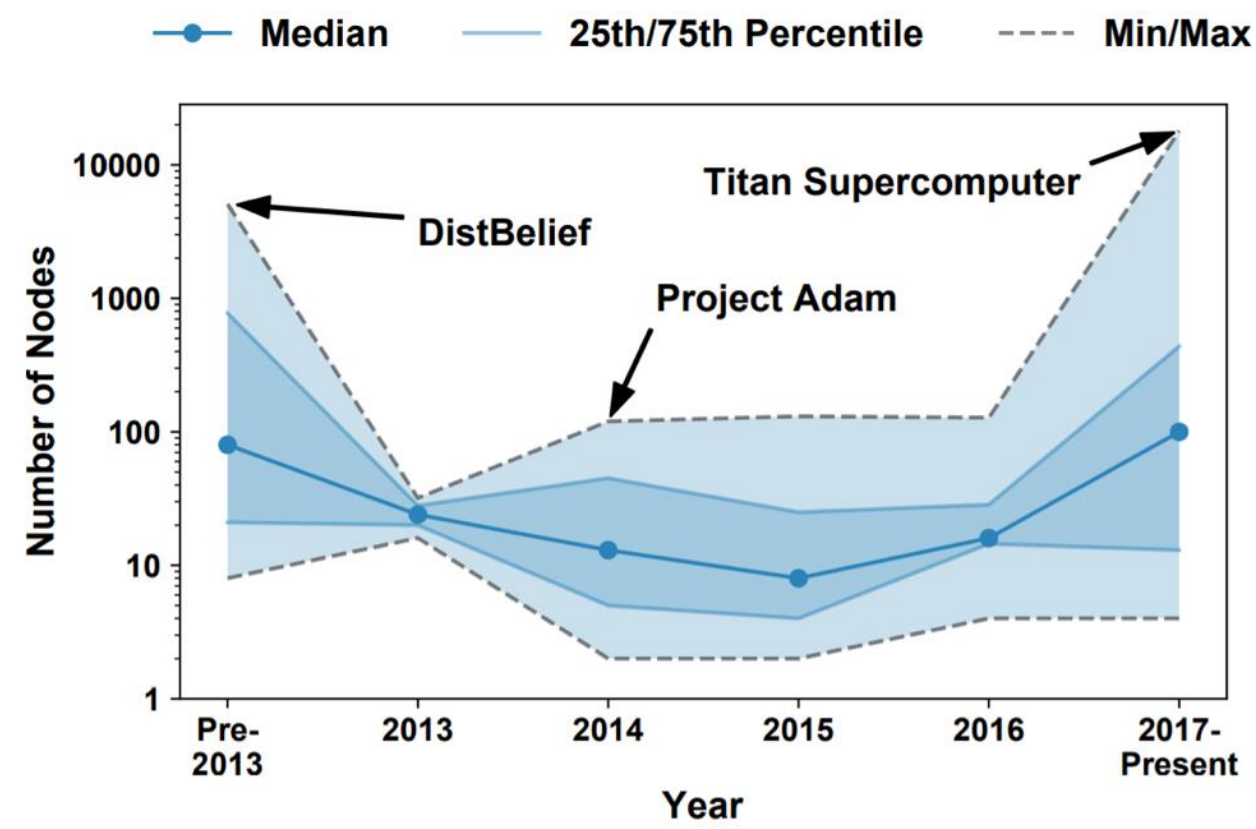
Deep Learning is largely on distributed memory today!

Trends in **distributed** deep learning: node count and communication

The field is moving fast – trying everything imaginable – survey results from 252 papers in the area of parallel deep learning

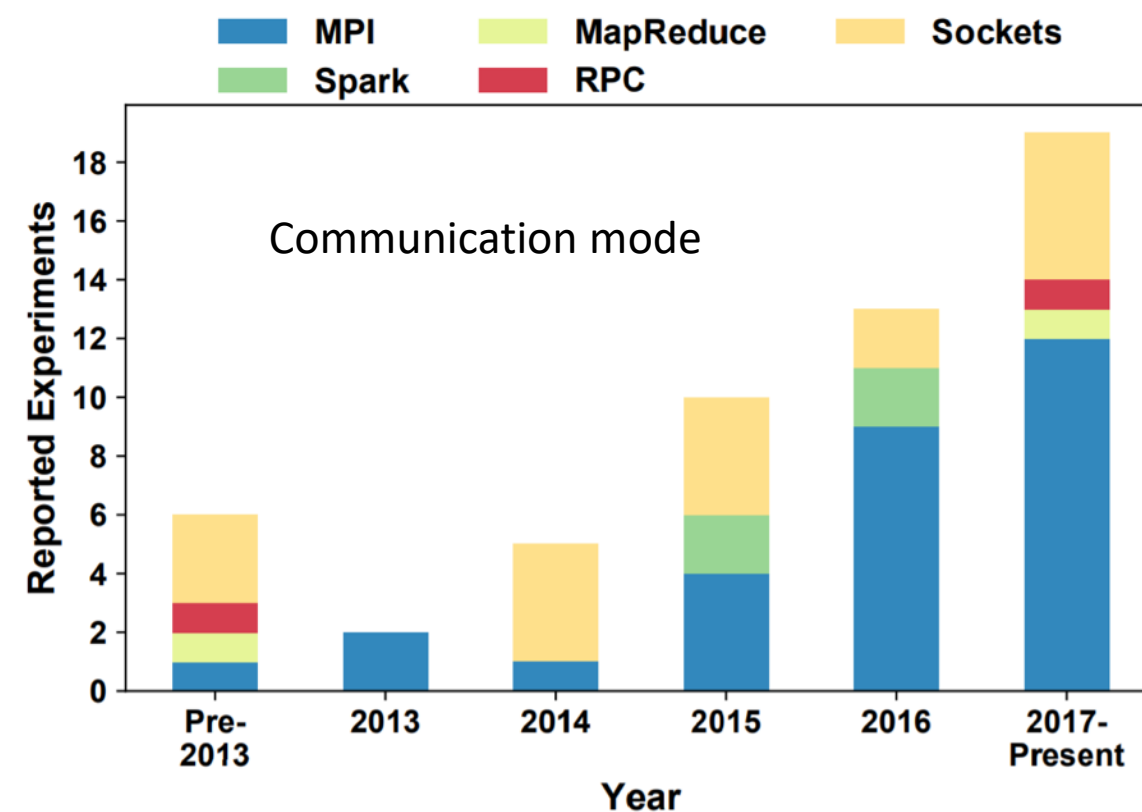
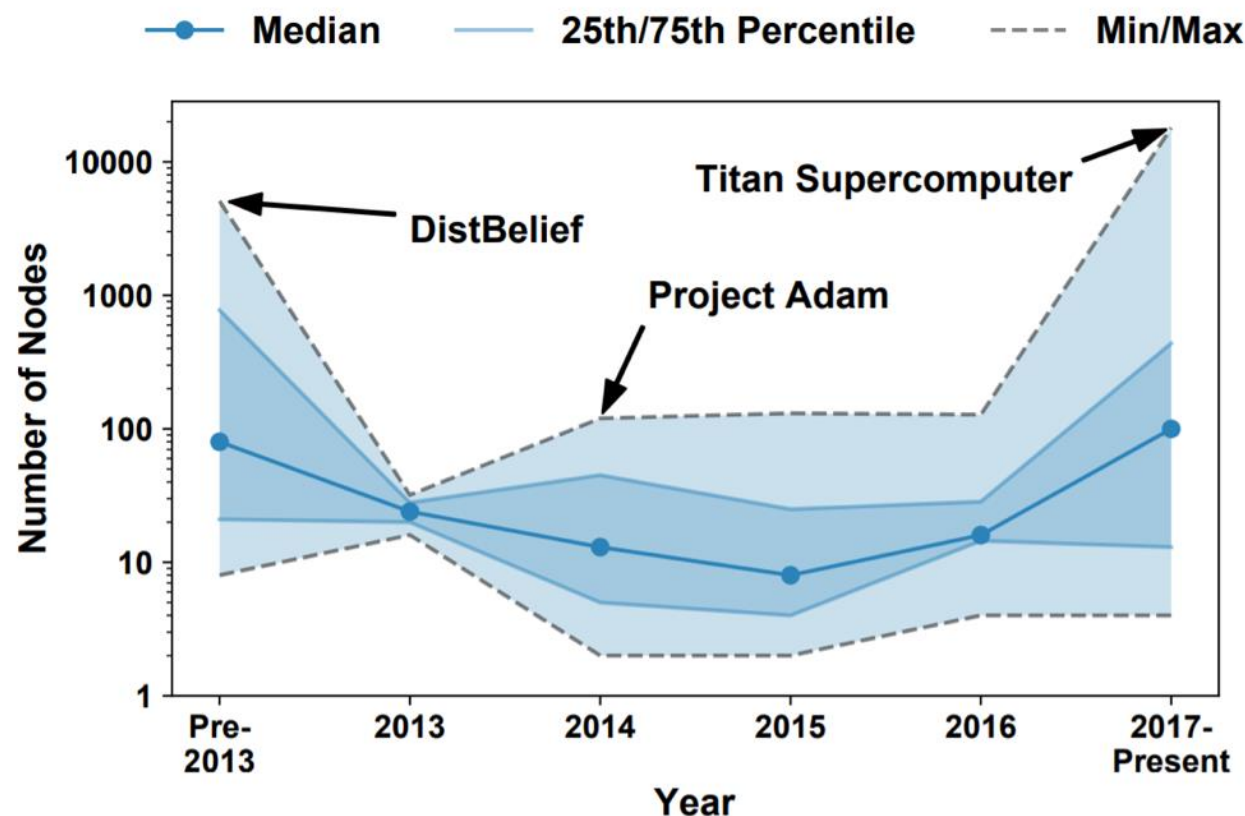
Trends in **distributed** deep learning: node count and communication

The field is moving fast – trying everything imaginable – survey results from 252 papers in the area of parallel deep learning



Trends in **distributed** deep learning: node count and communication

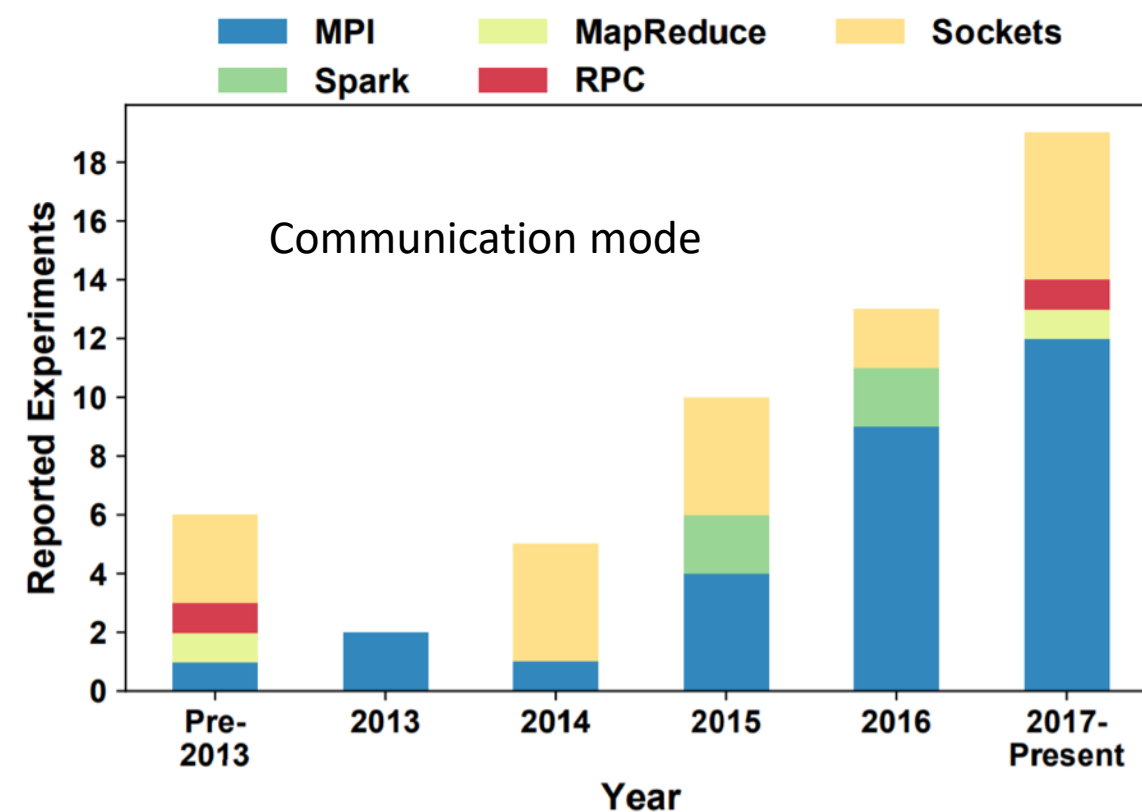
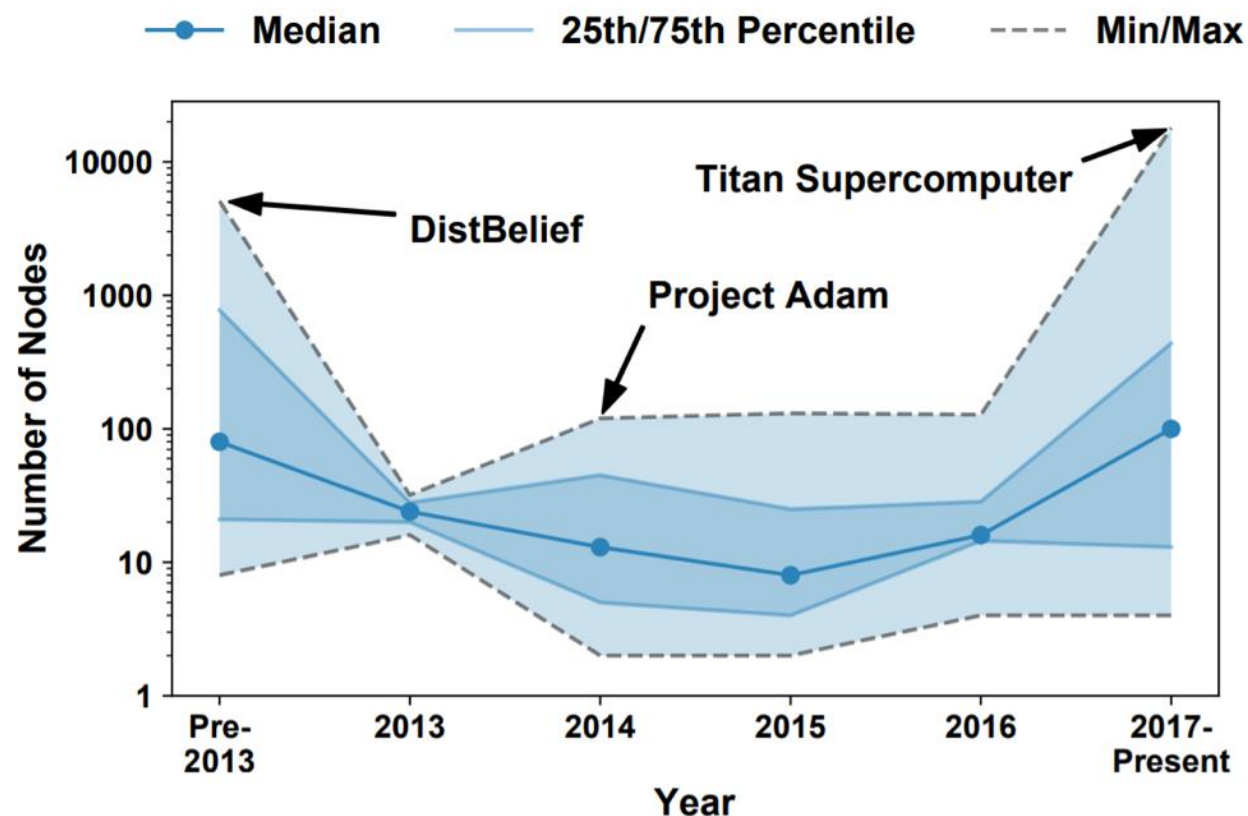
The field is moving fast – trying everything imaginable – survey results from 252 papers in the area of parallel deep learning



Deep Learning research is converging to MPI!

Trends in **distributed** deep learning: node count and communication

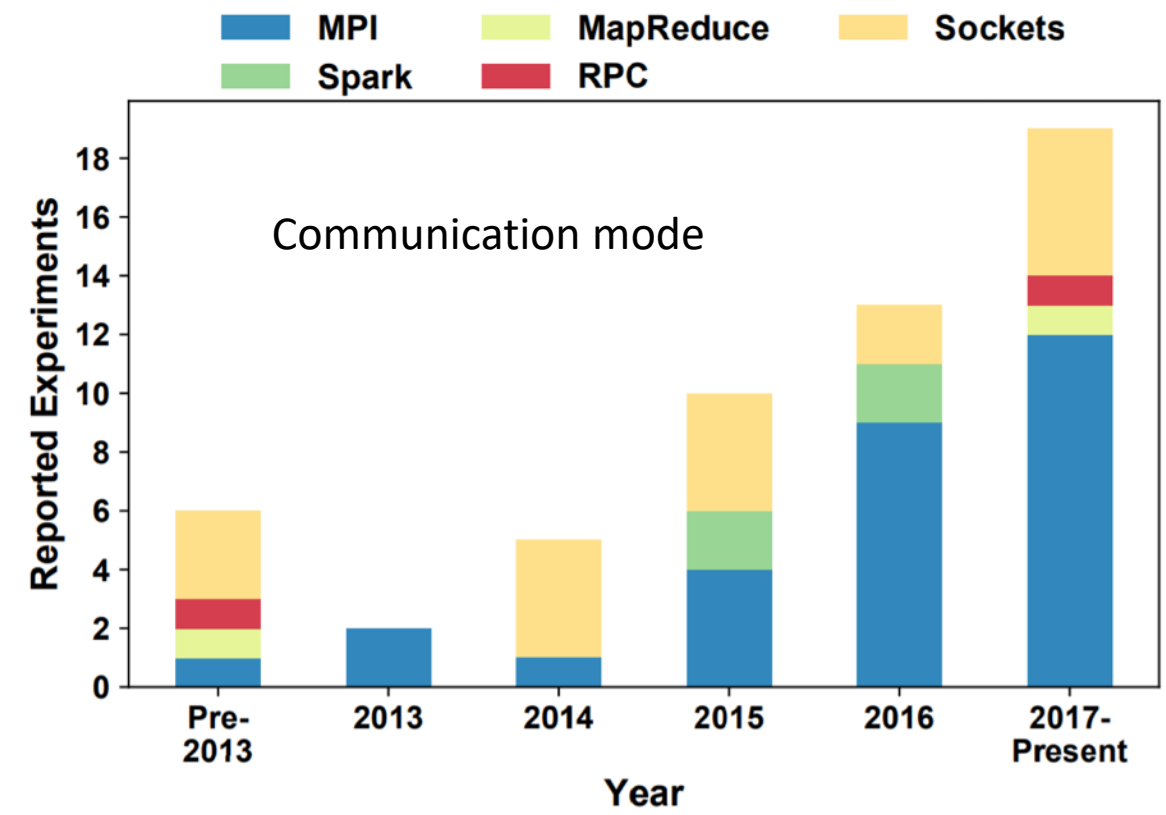
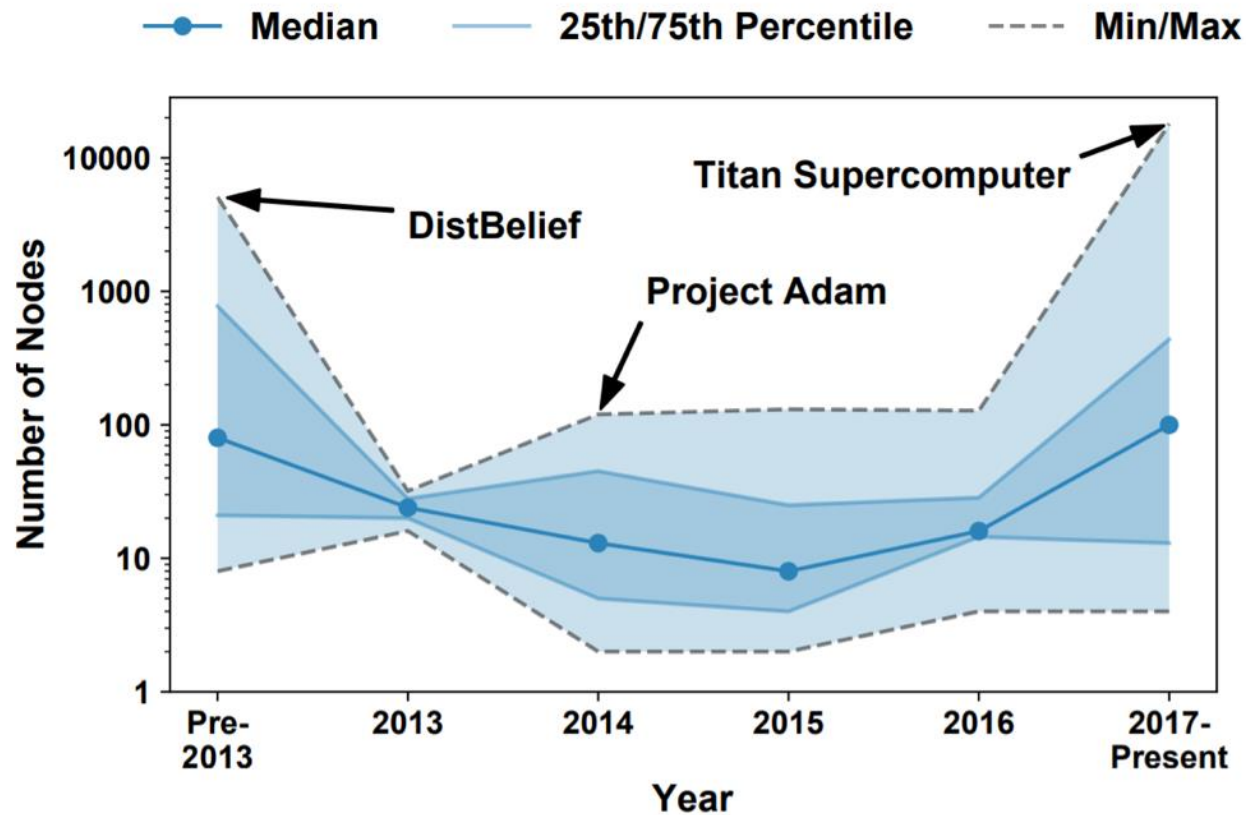
The field is moving fast – trying everything imaginable – survey results from 252 papers in the area of parallel deep learning



Deep Learning research is converging to MPI!

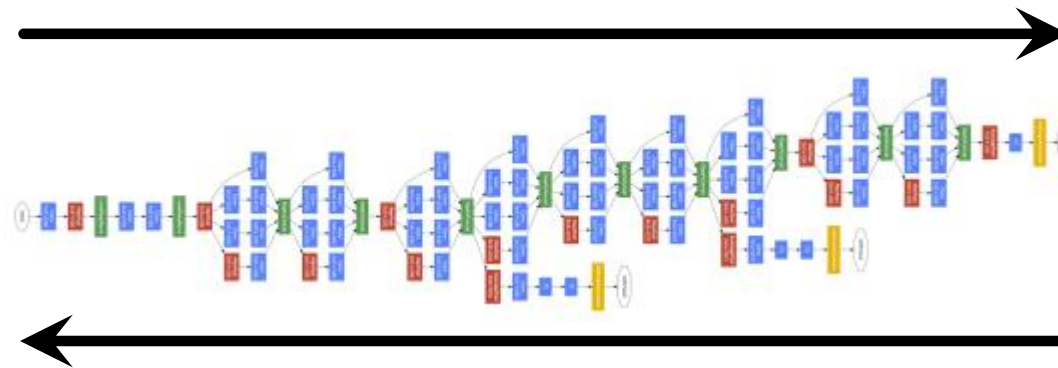
Trends in **distributed** deep learning: node count and communication

The field is moving fast – trying everything imaginable – survey results from 252 papers in the area of parallel deep learning



Deep Learning research is converging to MPI!

Minibatch Stochastic Gradient Descent (SGD)

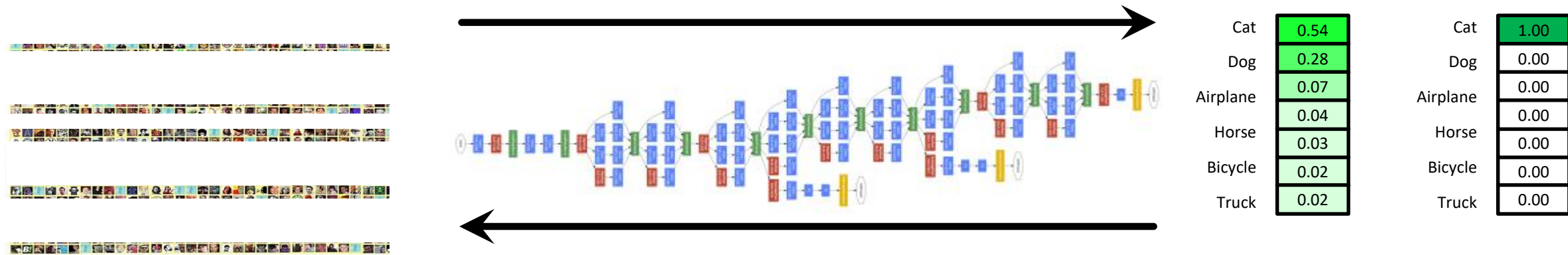


Cat	0.54
Dog	0.28
Airplane	0.07
Horse	0.04
Bicycle	0.02
Truck	0.02

Cat	1.00
Dog	0.00
Airplane	0.00
Horse	0.00
Bicycle	0.00
Truck	0.00

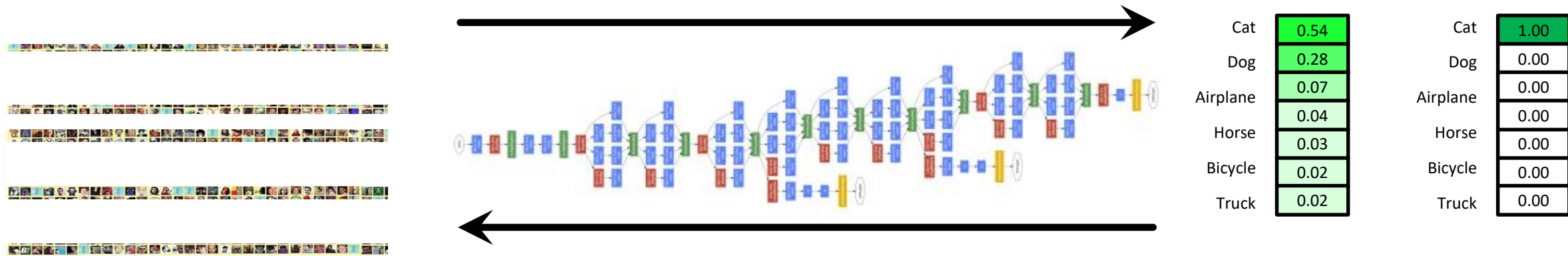
-
- 1: **for** $t = 0$ **to** T **do** ▷ Stopping condition
 - 2: $x \leftarrow$ Random element from X ▷ Sample dataset
 - 3: $o_1 \leftarrow f_1(x)$
 - 4: **for** $i = 2$ **to** *layers* **do** ▷ Forward evaluation of ℓ
 - 5: $o_i \leftarrow f_i(o_{i-1})$
 - 6: **end for**
 - 7: **for** $i = \text{layers} - 1$ **to** 1 **do** ▷ Compute gradient of ℓ via backpropagation
 - 8: $\nabla o_i \leftarrow \frac{\partial \ell}{\partial o_i}(o_{i-1}, o_i, \nabla o_{i+1})$ ▷ Gradient w.r.t. data
 - 9: $\nabla w_i^{(t)} \leftarrow \frac{\partial \ell}{\partial w_i}(o_{i-1}, o_i, \nabla o_{i+1})$ ▷ Gradient w.r.t. layer parameters
 - 10: **end for**
 - 11: $w^{(t+1)} \leftarrow w^{(t)} + u(\nabla w^{(t)}, w^{(0, \dots, t)}, t)$ ▷ Update weights with function u
 - 12: **end for**
-

Minibatch Stochastic Gradient Descent (SGD)



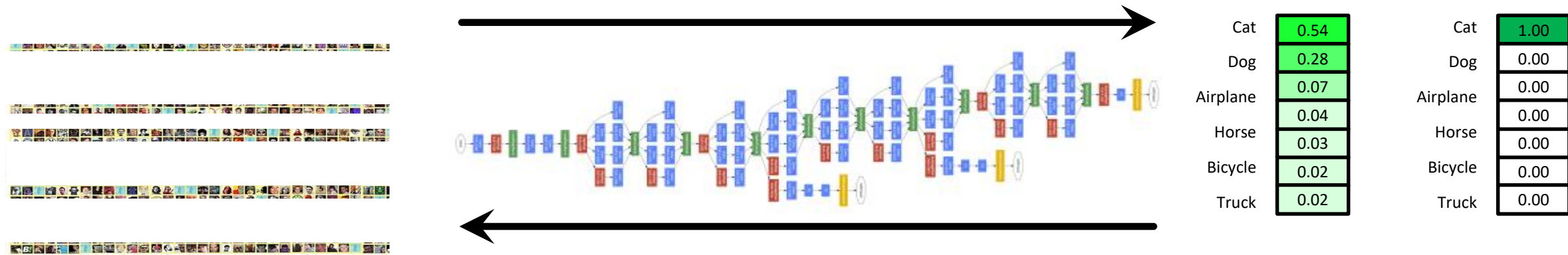
- 1: **for** $t = 0$ **to** T **do** ▷ Stopping condition
- 2: $\vec{x} \leftarrow$ Random B elements from X ▷ Sample dataset
- 3: $\vec{o}_1 \leftarrow f_1(\vec{x})$
- 4: **for** $i = 2$ **to** $layers$ **do** ▷ Forward evaluation of ℓ
- 5: $\vec{o}_i \leftarrow f_i(\vec{o}_{i-1})$
- 6: **end for**
- 7: **for** $i = layers - 1$ **to** 1 **do** ▷ Compute gradient of ℓ via backpropagation
- 8: $\vec{\nabla} o_i \leftarrow \frac{\partial \ell}{\partial o_i}(\vec{o}_{i-1}, \vec{o}_i, \vec{\nabla} o_{i+1})$ ▷ Gradient w.r.t. data
- 9: $\vec{\nabla} w_i^{(t)} \leftarrow \frac{\partial \ell}{\partial w_i}(\vec{o}_{i-1}, \vec{o}_i, \vec{\nabla} o_{i+1})$ ▷ Gradient w.r.t. layer parameters
- 10: **end for**
- 11: $\vec{\nabla} w^{(t)} \leftarrow \frac{1}{B} \sum_{i=0}^B \vec{\nabla} w^{(t)}[i]$
- 12: $w^{(t+1)} \leftarrow w^{(t)} + u(\vec{\nabla} w^{(t)}, w^{(0, \dots, t)}, t)$ ▷ Update weights with function u
- 13: **end for**

Minibatch Stochastic Gradient Descent (SGD)



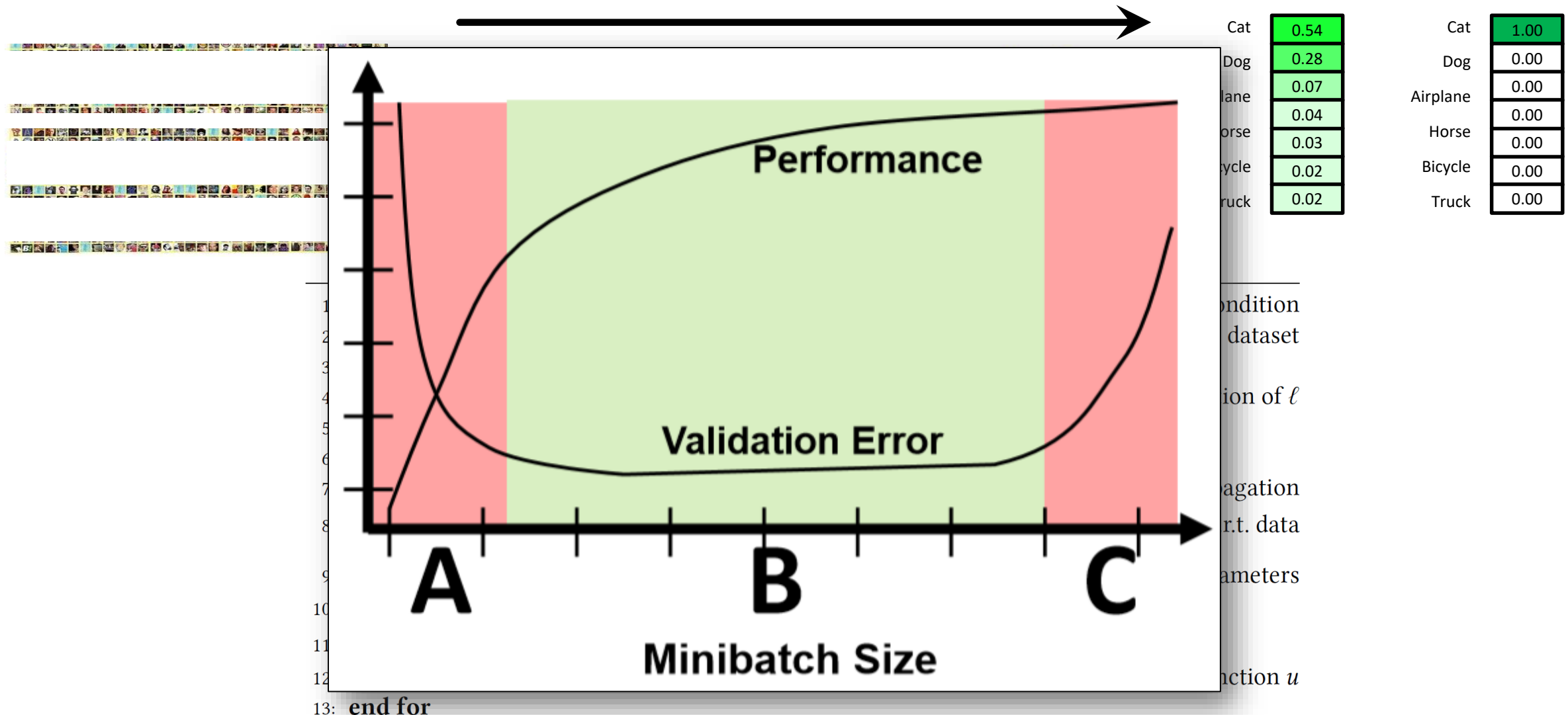
- 1: **for** $t = 0$ **to** T **do**
 - 2: $\vec{x} \leftarrow$ Random B elements from X
 - Stopping condition
 - Sample dataset
 - 3: $\vec{o}_1 \leftarrow f_1(\vec{x})$
 - 4: **for** $i = 2$ **to** $layers$ **do**
 - 5: $\vec{o}_i \leftarrow f_i(\vec{o}_{i-1})$
 - Forward evaluation of ℓ
 - 6: **end for**
 - 7: **for** $i = layers - 1$ **to** 1 **do**
 - 8: $\vec{\nabla}_{o_i} \leftarrow \frac{\partial \ell}{\partial o_i}(\vec{o}_{i-1}, \vec{o}_i, \vec{\nabla}_{o_{i+1}})$
 - Compute gradient of ℓ via backpropagation
 - Gradient w.r.t. data
 - 9: $\vec{\nabla}_{w_i^{(t)}} \leftarrow \frac{\partial \ell}{\partial w_i}(\vec{o}_{i-1}, \vec{o}_i, \vec{\nabla}_{o_{i+1}})$
 - Gradient w.r.t. layer parameters
 - 10: **end for**
 - 11: $\vec{\nabla}_{w^{(t)}} \leftarrow \frac{1}{B} \sum_{i=0}^B \vec{\nabla}_{w^{(t)}}[i]$
 - 12: $w^{(t+1)} \leftarrow w^{(t)} + u(\vec{\nabla}_{w^{(t)}}, w^{(0, \dots, t)}, t)$
 - Update weights with function u
 - 13: **end for**

Minibatch Stochastic Gradient Descent (SGD)



- 1: **for** $t = 0$ **to** T **do** ▷ Stopping condition
- 2: $\vec{x} \leftarrow$ Random B elements from X ▷ Sample dataset
- 3: $\vec{o}_1 \leftarrow f_1(\vec{x})$
- 4: **for** $i = 2$ **to** $layers$ **do** ▷ Forward evaluation of ℓ
- 5: $\vec{o}_i \leftarrow f_i(\vec{o}_{i-1})$
- 6: **end for**
- 7: **for** $i = layers - 1$ **to** 1 **do** ▷ Compute gradient of ℓ via backpropagation
- 8: $\vec{\nabla} o_i \leftarrow \frac{\partial \ell}{\partial o_i}(\vec{o}_{i-1}, \vec{o}_i, \vec{\nabla} o_{i+1})$ ▷ Gradient w.r.t. data
- 9: $\vec{\nabla} w_i^{(t)} \leftarrow \frac{\partial \ell}{\partial w_i}(\vec{o}_{i-1}, \vec{o}_i, \vec{\nabla} o_{i+1})$ ▷ Gradient w.r.t. layer parameters
- 10: **end for**
- 11: $\vec{\nabla} w^{(t)} \leftarrow \frac{1}{B} \sum_{i=0}^B \vec{\nabla} w^{(t)}[i]$
- 12: $w^{(t+1)} \leftarrow w^{(t)} + u(\vec{\nabla} w^{(t)}, w^{(0), \dots, t}, t)$ ▷ Update weights with function u
- 13: **end for**

Minibatch Stochastic Gradient Descent (SGD)

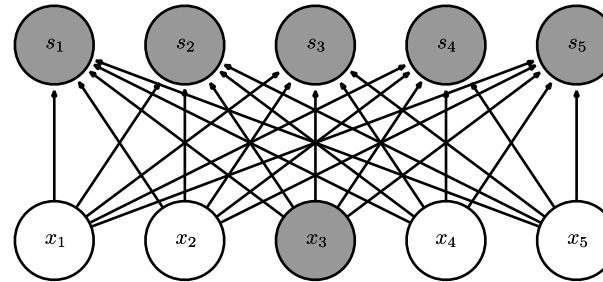


13: end for

Ingredients of a neural network: Operators

Ingredients of a neural network: Operators

- Fully-connected layers (multi-layer perceptrons)



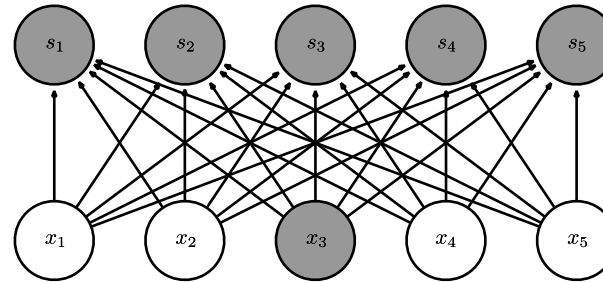
- Convolution

- Many other moving parts:

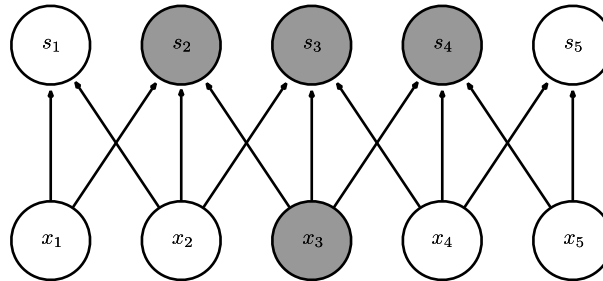
- Pooling
- Batch normalization [Ioffe & Szegedy 2015]
- ReLU activations [Glorot, Bordes, & Bengion 2011]
- ...

Ingredients of a neural network: Operators

- Fully-connected layers (multi-layer perceptrons)



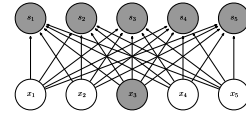
- Convolution



- Many other moving parts:

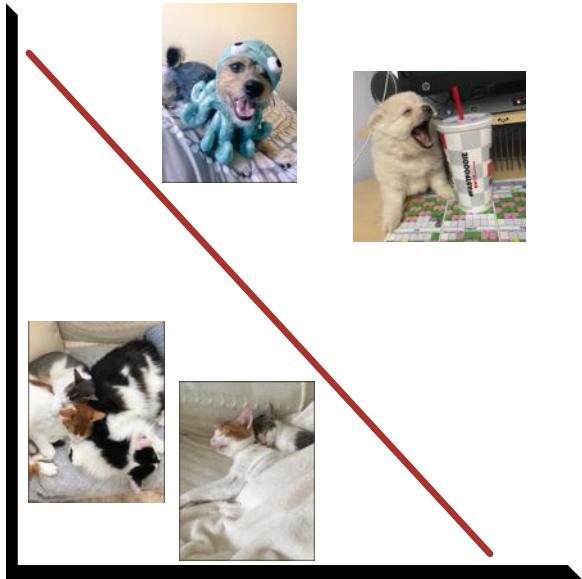
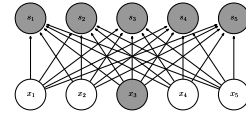
- Pooling
- Batch normalization [Ioffe & Szegedy 2015]
- ReLU activations [Glorot, Bordes, & Bengion 2011]
- ...

Fully-connected layers



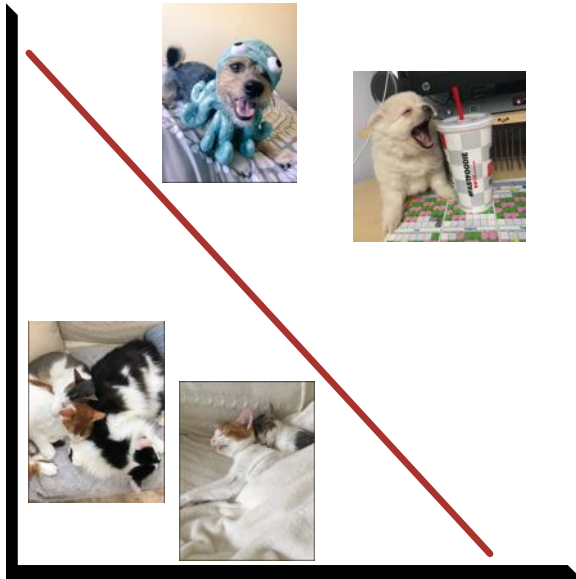
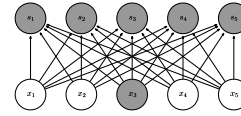
Fully-connected layers

Perceptron [Rosenblatt 1958]

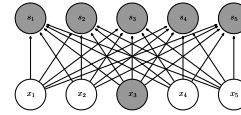


Fully-connected layers

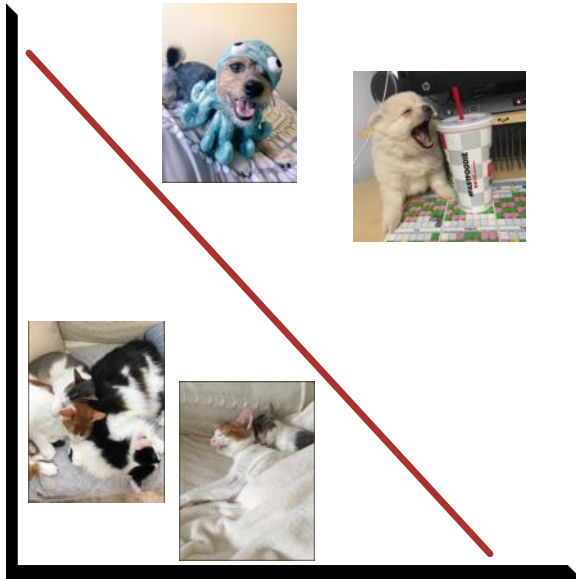
Perceptron [Rosenblatt 1958]



Fully-connected layers

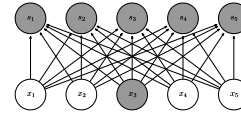


Perceptron [Rosenblatt 1958]

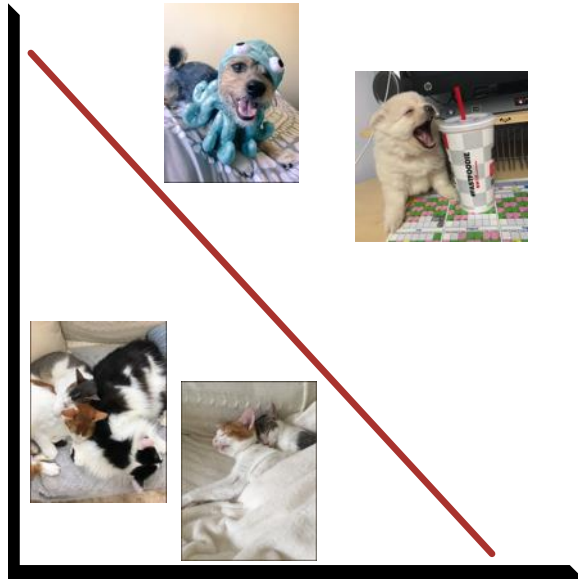


$$y = wx + b$$

Fully-connected layers



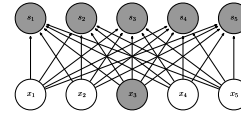
Perceptron [Rosenblatt 1958]



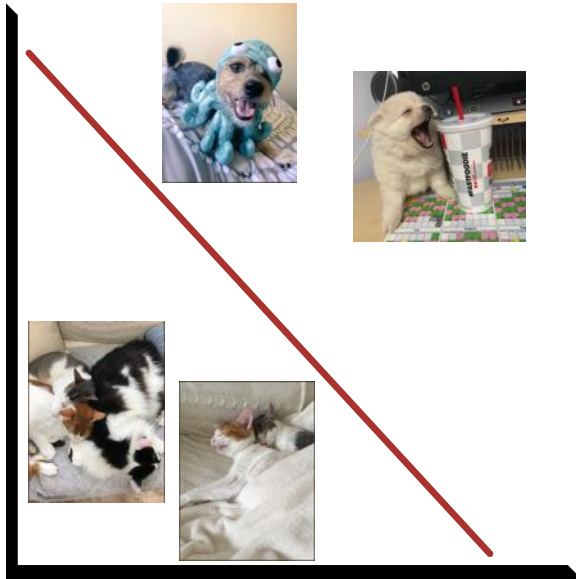
$$y = wx + b$$

Learned weights and bias

Fully-connected layers



Perceptron [Rosenblatt 1958]

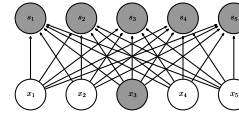


$$y = wx + b$$

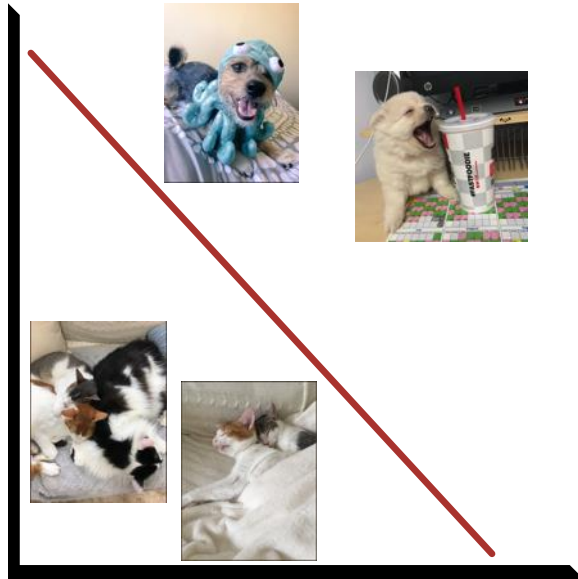


Learned weights and bias

Fully-connected layers



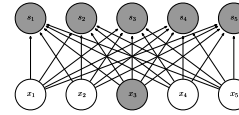
Perceptron [Rosenblatt 1958]



$$y = wx + b$$

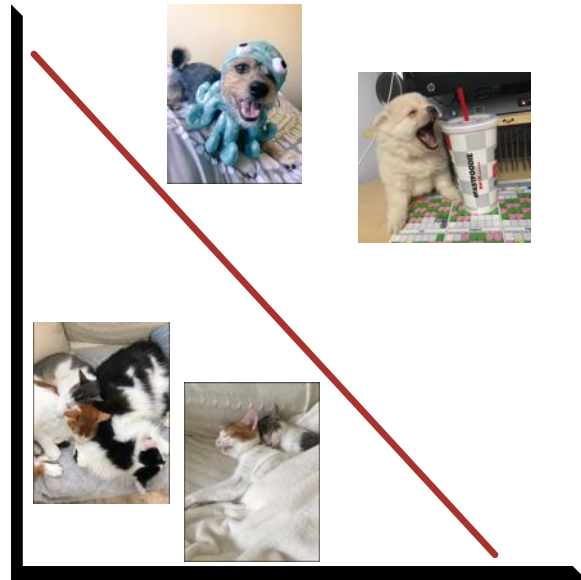


Learned weights and bias

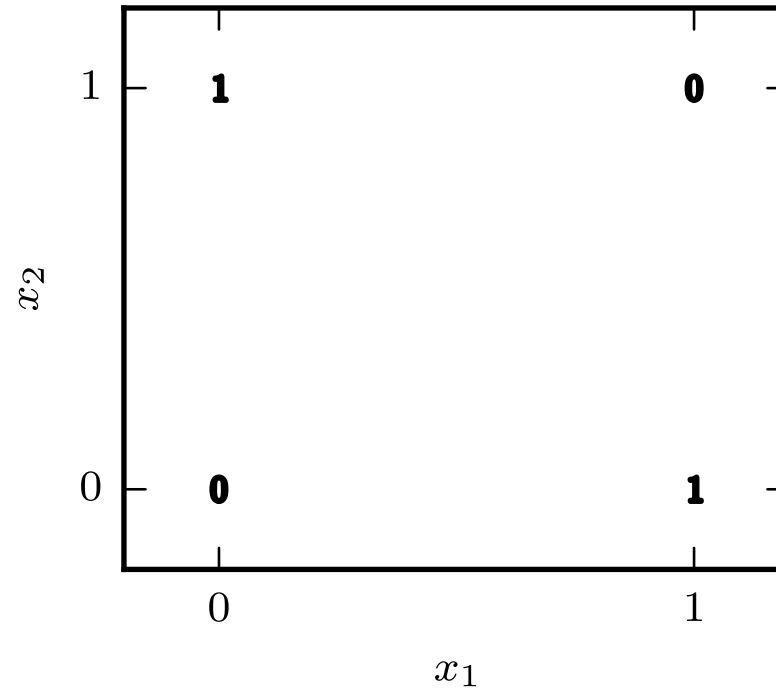


Fully-connected layers

Perceptron [Rosenblatt 1958]



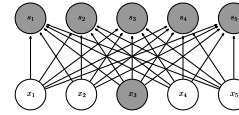
Exclusive OR



$$y = wx + b$$

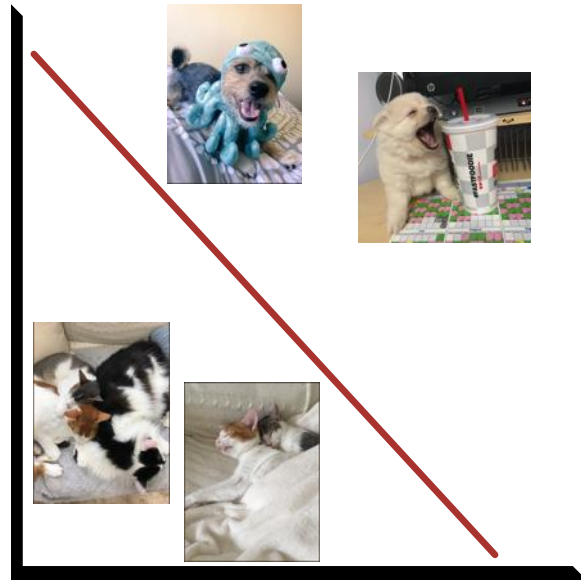


Learned weights and bias

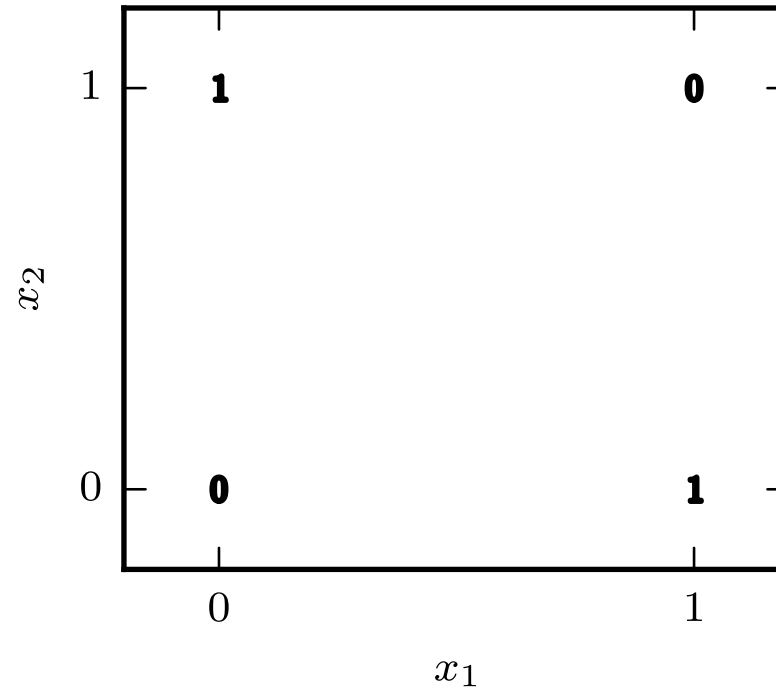


Fully-connected layers

Perceptron [Rosenblatt 1958]



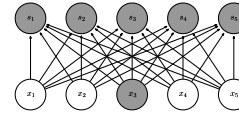
Exclusive OR



$$y = wx + b$$

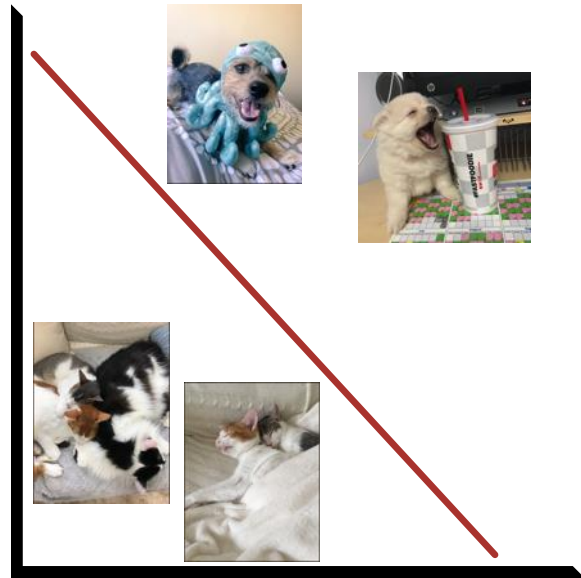


Learned weights and bias

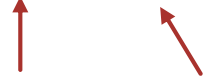


Fully-connected layers

Perceptron [Rosenblatt 1958]

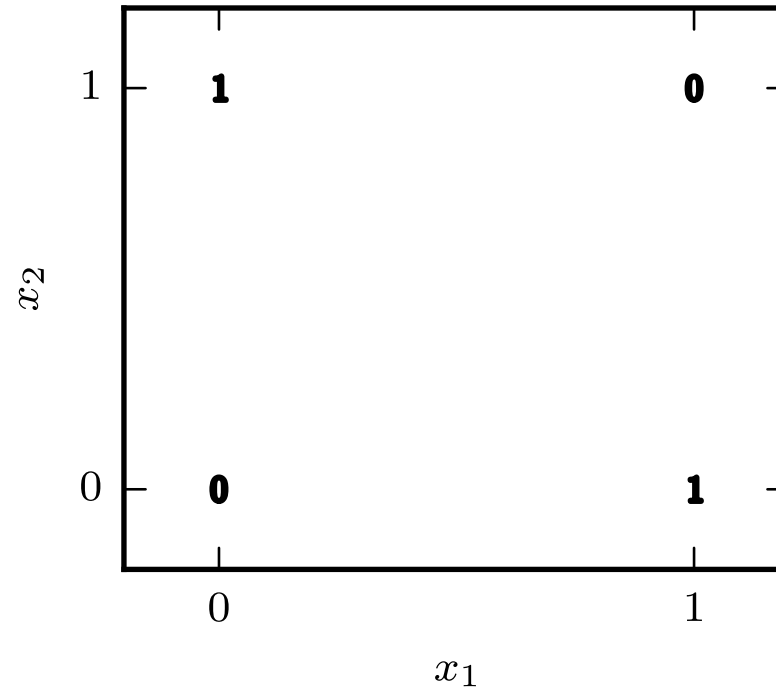


$$y = wx + b$$

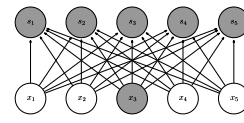


Learned weights and bias

Exclusive OR

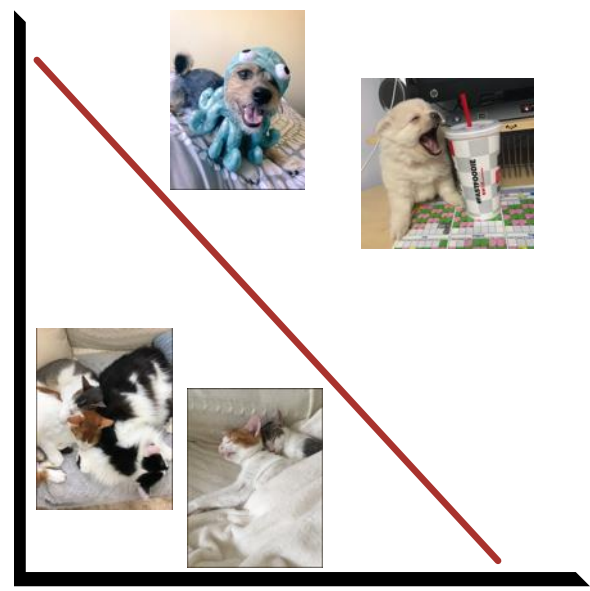


Not linearly separable!



Fully-connected layers

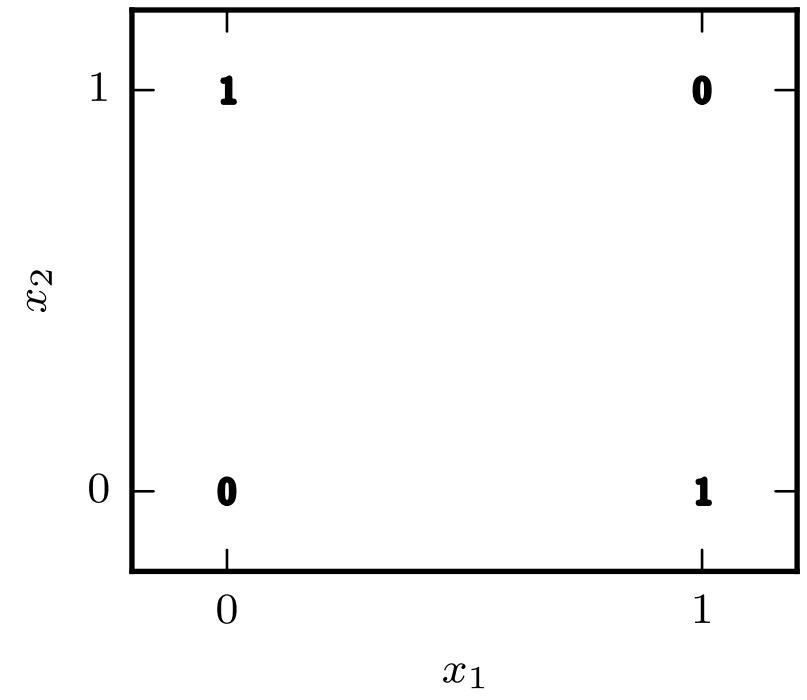
Perceptron [Rosenblatt 1958]



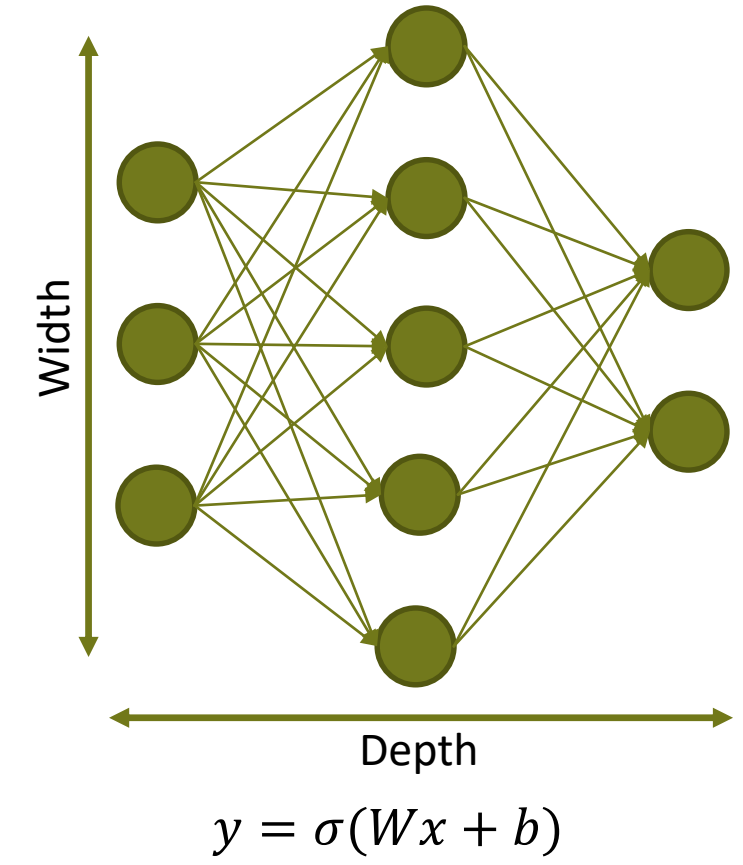
$$y = wx + b$$

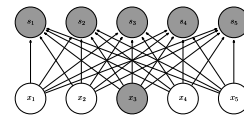
Learned weights and bias

Exclusive OR



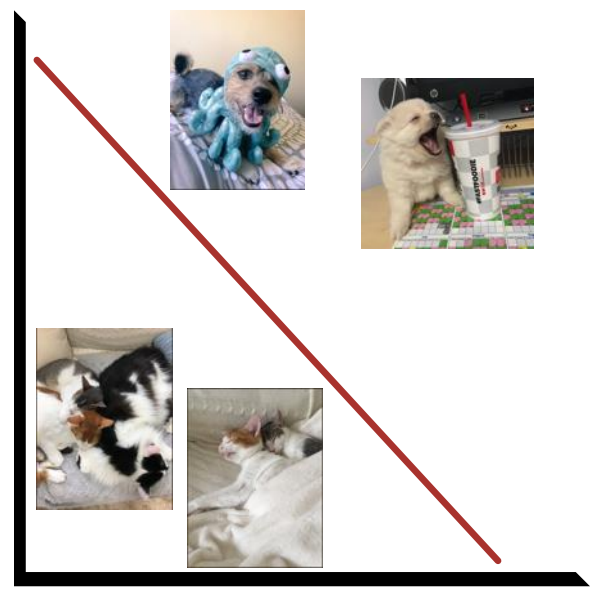
Not linearly separable!





Fully-connected layers

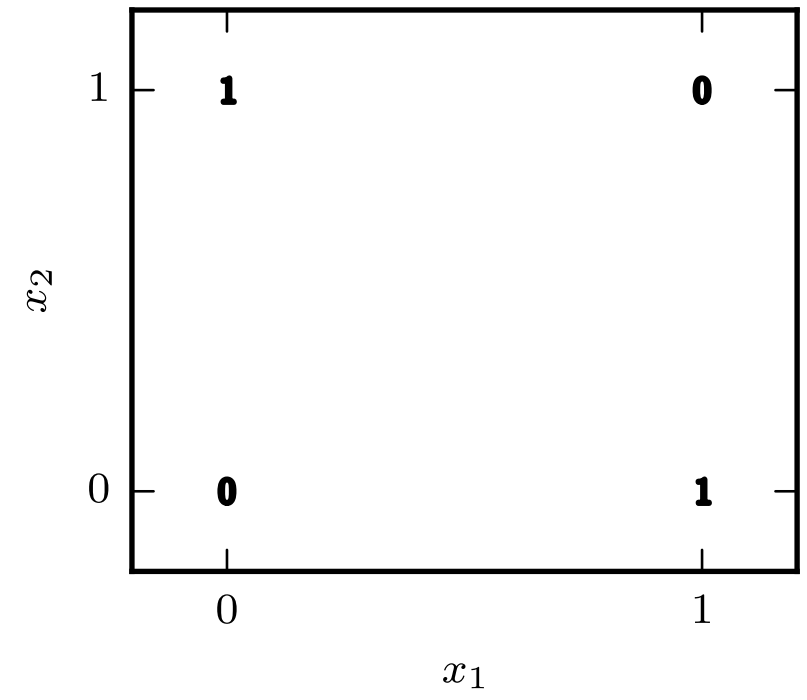
Perceptron [Rosenblatt 1958]



$$y = wx + b$$

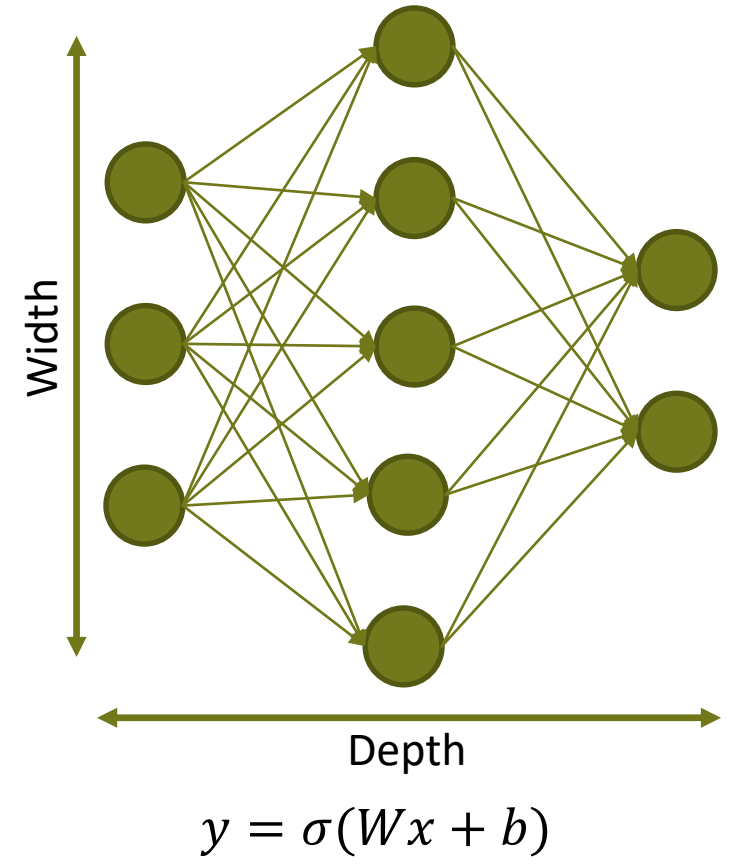
Learned weights and bias

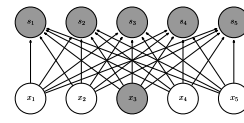
Exclusive OR



Not linearly separable!

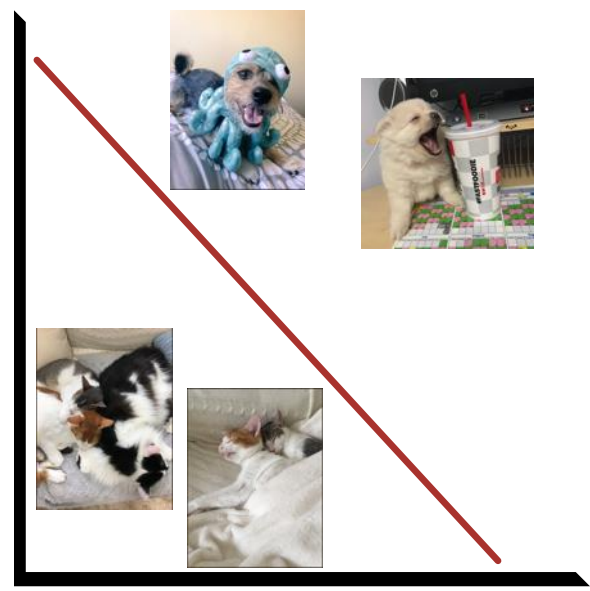
Multi-layer Perceptron





Fully-connected layers

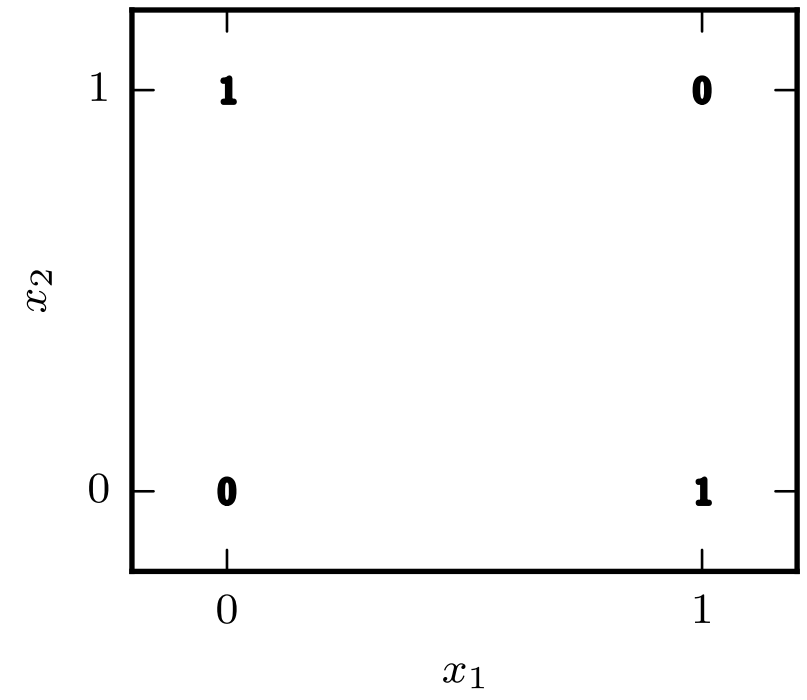
Perceptron [Rosenblatt 1958]



$$y = wx + b$$

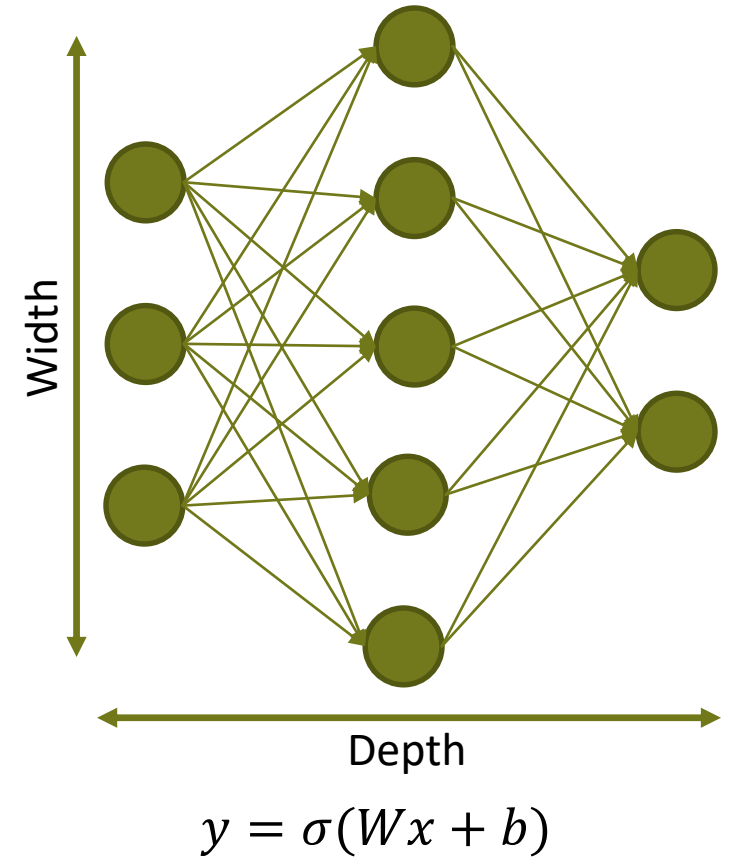
Learned weights and bias

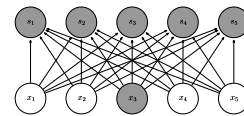
Exclusive OR



Not linearly separable!

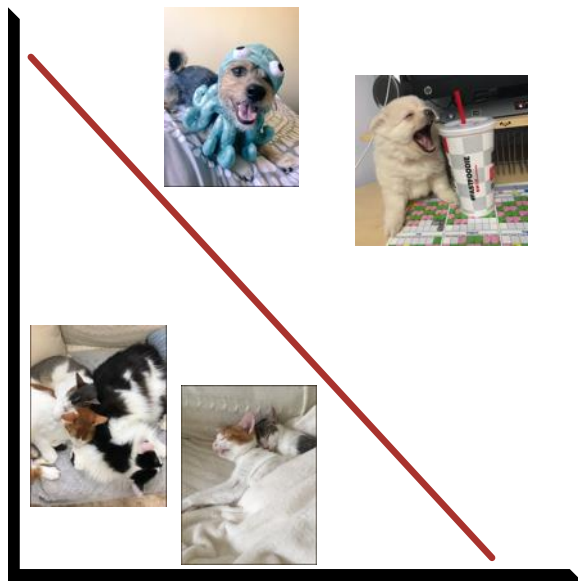
Multi-layer Perceptron





Fully-connected layers

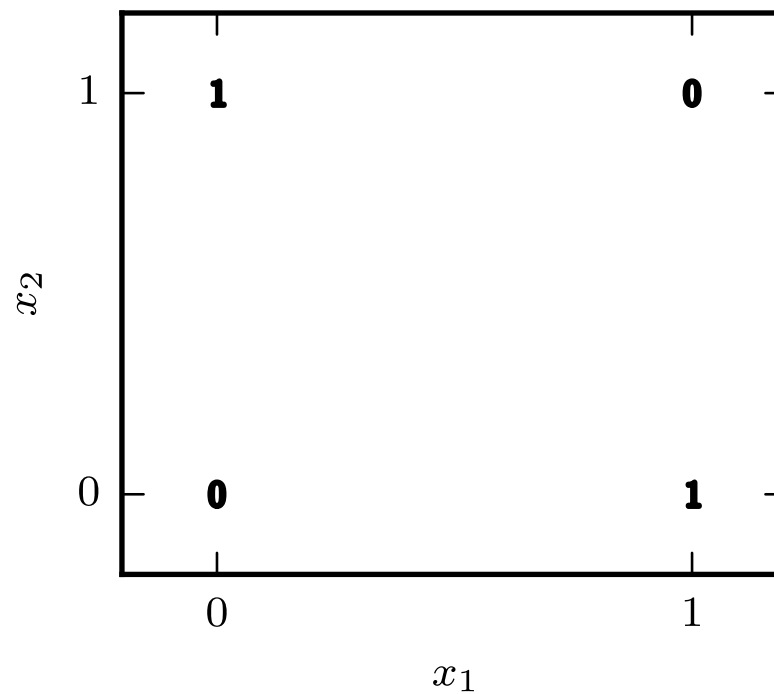
Perceptron [Rosenblatt 1958]



$$y = wx + b$$

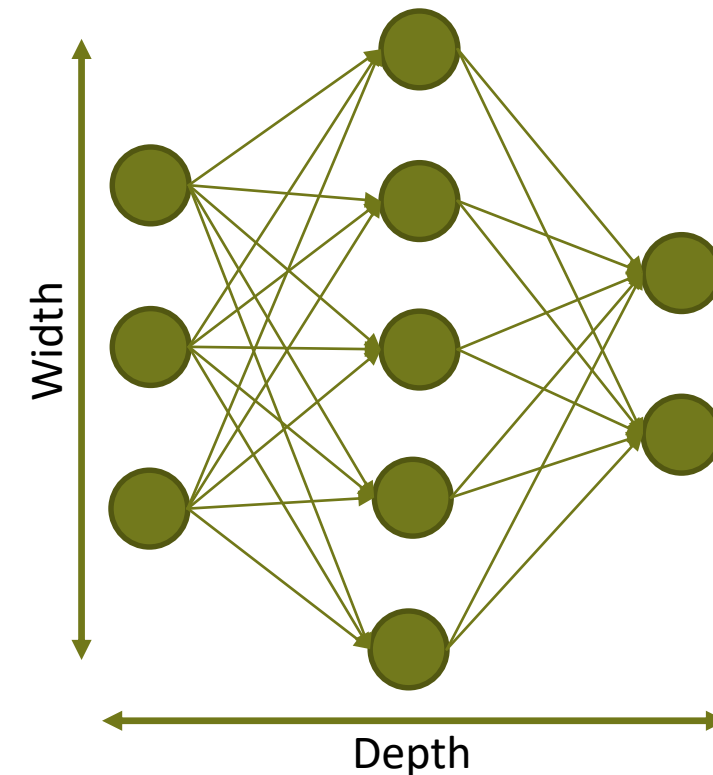
Learned weights and bias

Exclusive OR



Not linearly separable!

Multi-layer Perceptron

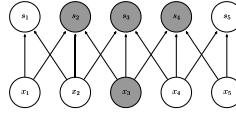


$$y = \sigma(Wx + b)$$

Universal!

Convolution

Inputs

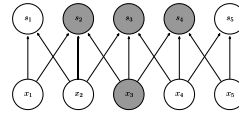
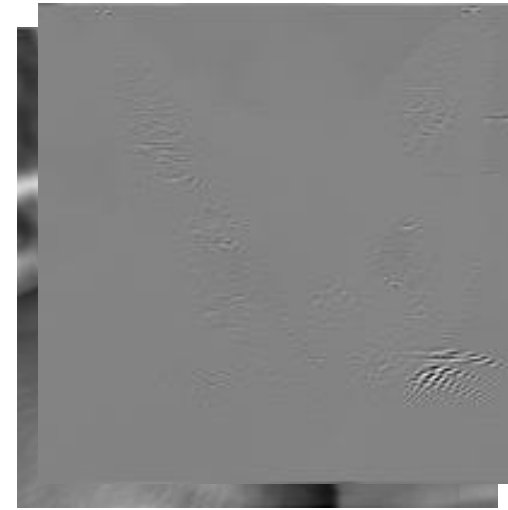
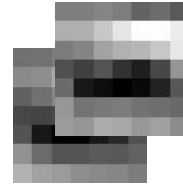


Convolution

Inputs

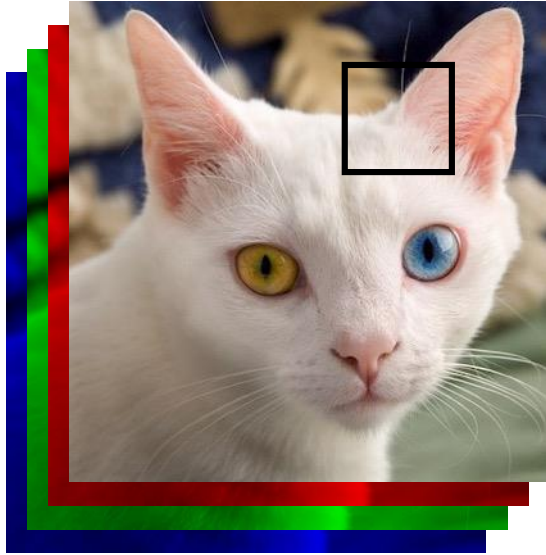


Filters

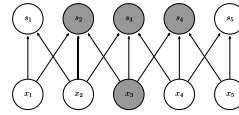
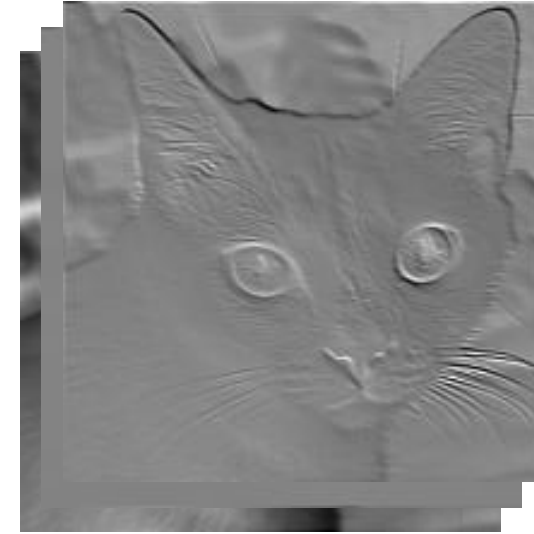
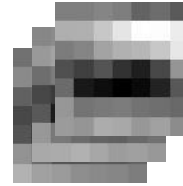


Convolution

Inputs

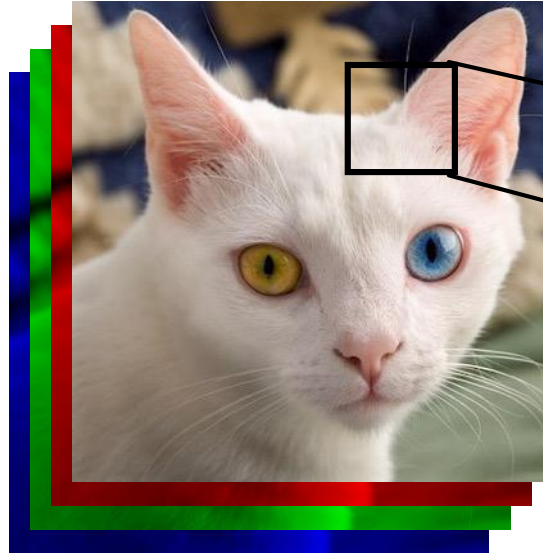


Filters

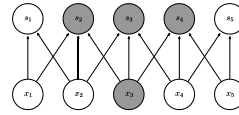
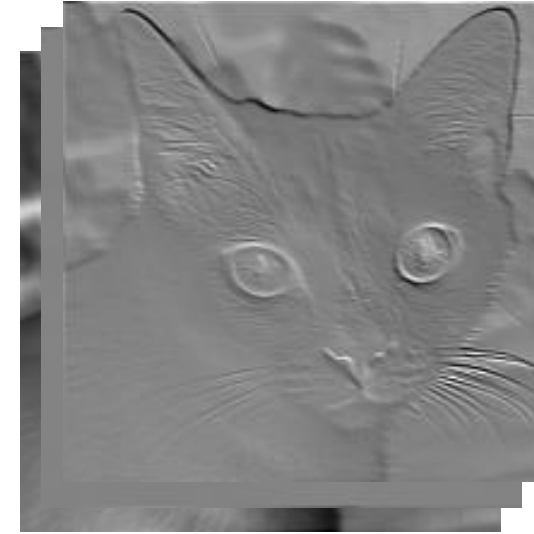
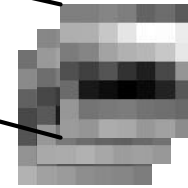


Convolution

Inputs

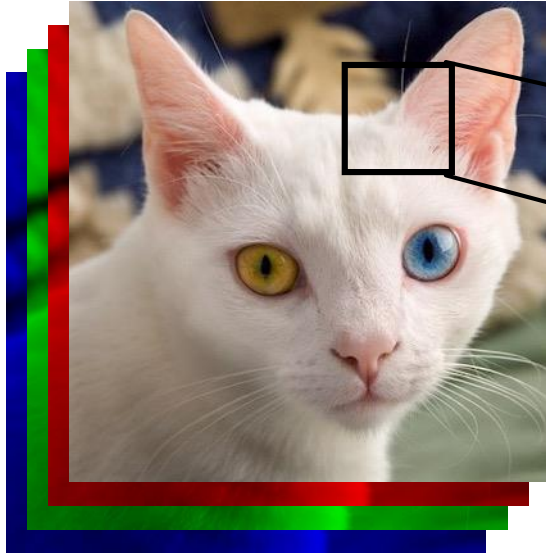


Filters

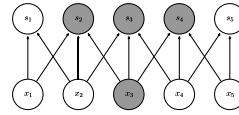
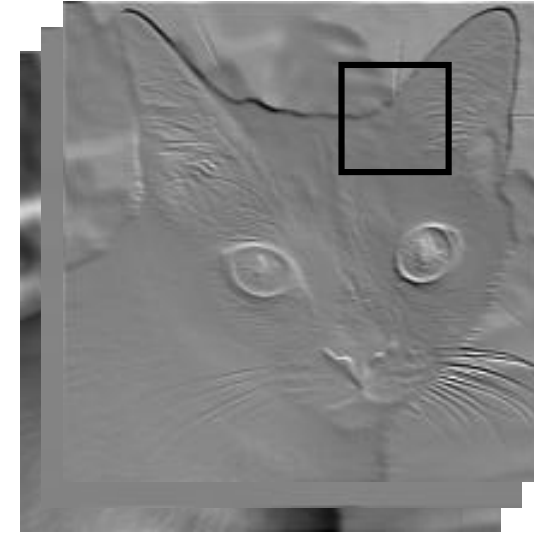
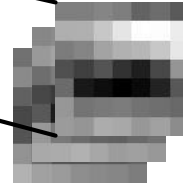


Convolution

Inputs

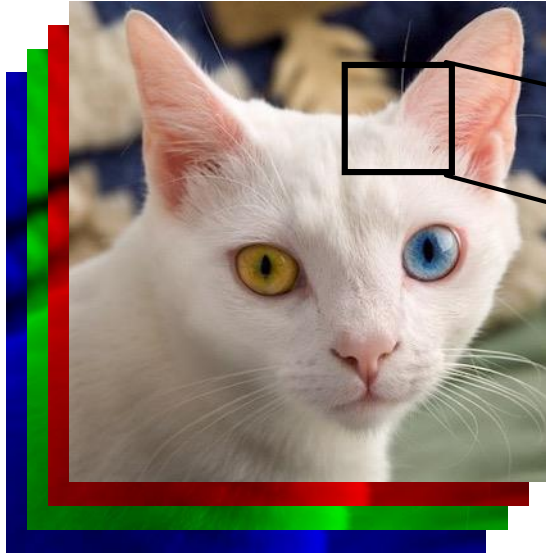


Filters

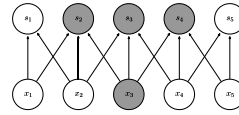
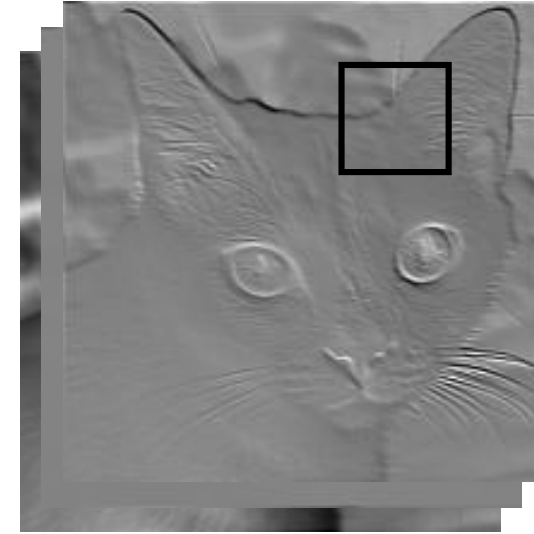
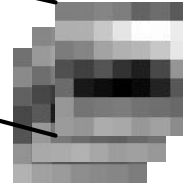


Convolution

Inputs

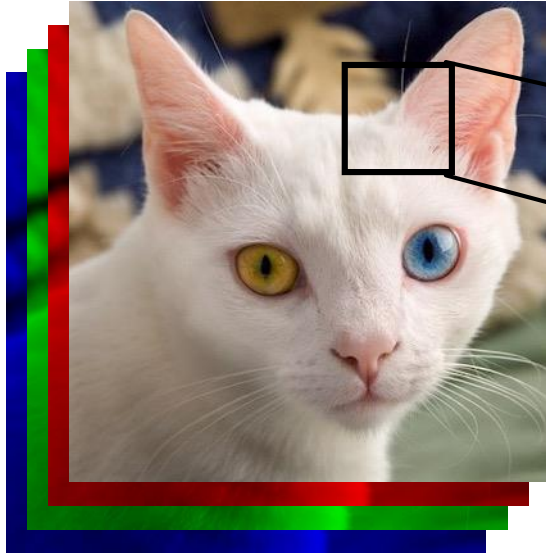


Filters

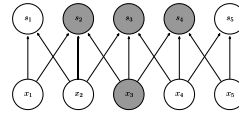
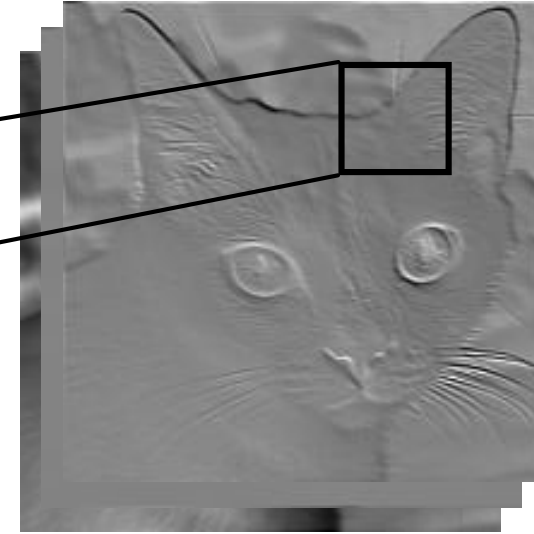
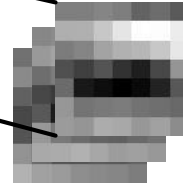


Convolution

Inputs



Filters

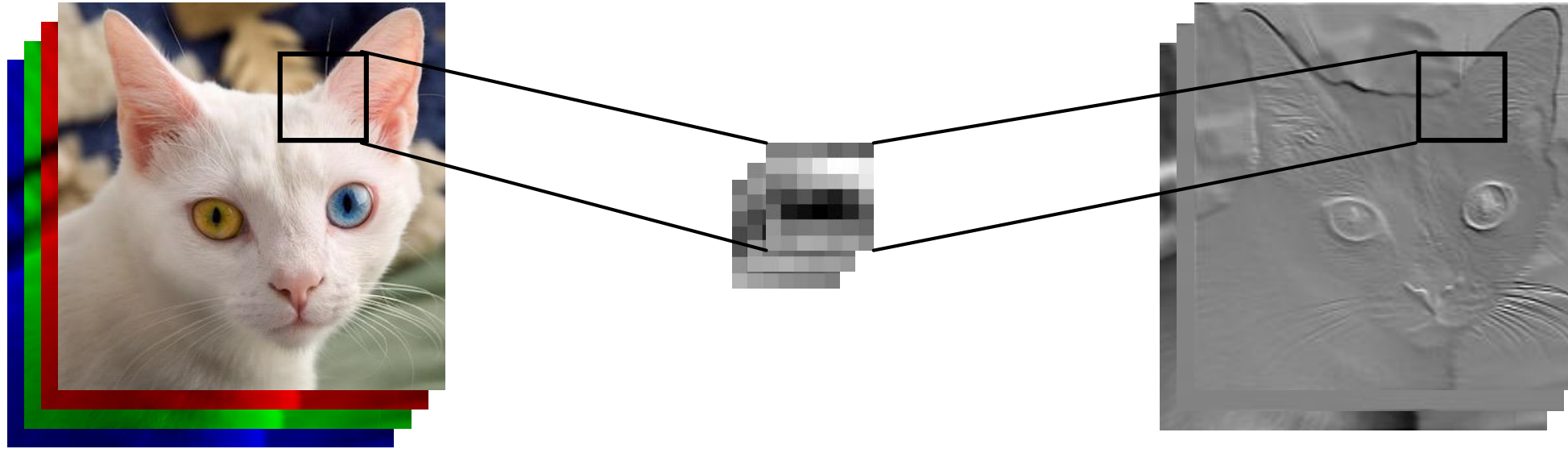
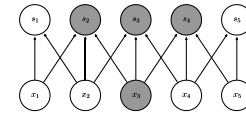


Convolution

Inputs

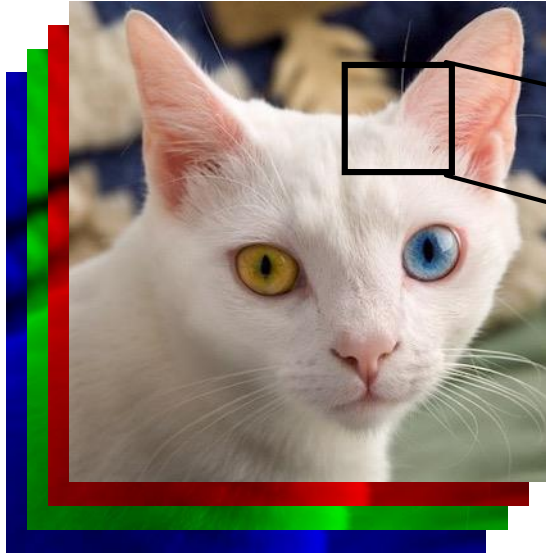
Filters

Feature maps (activations)

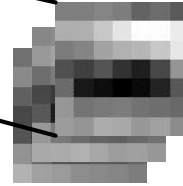


Convolution

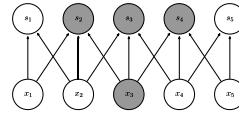
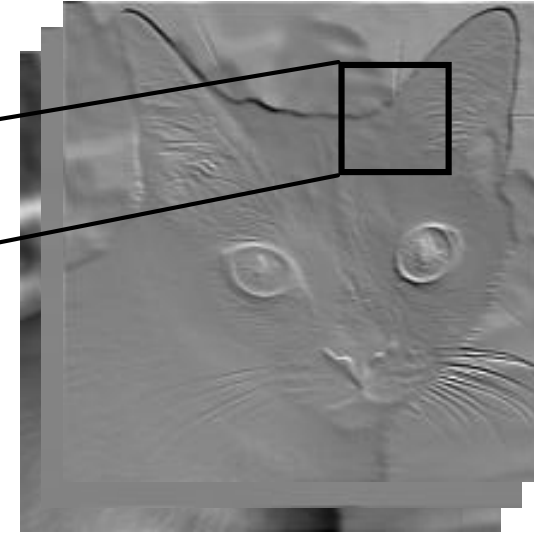
Inputs



Filters

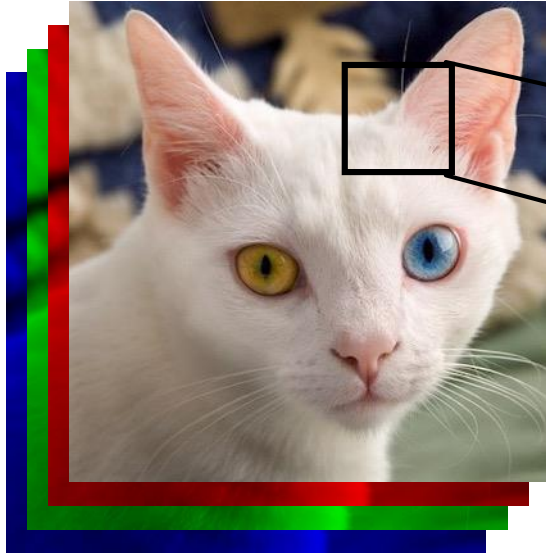


Feature maps (activations)

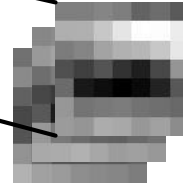


Convolution

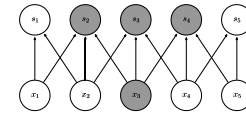
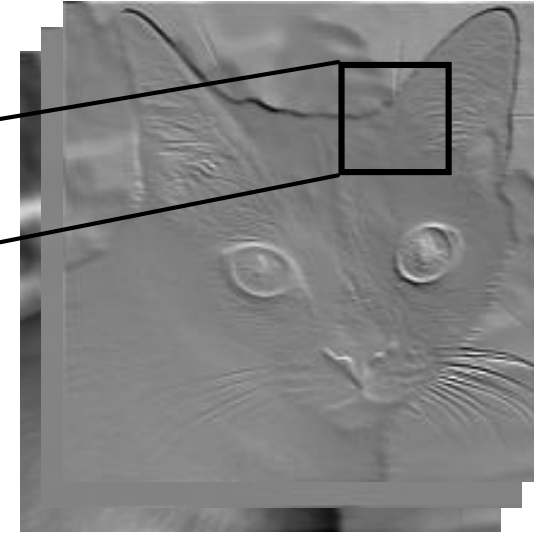
Inputs



Filters

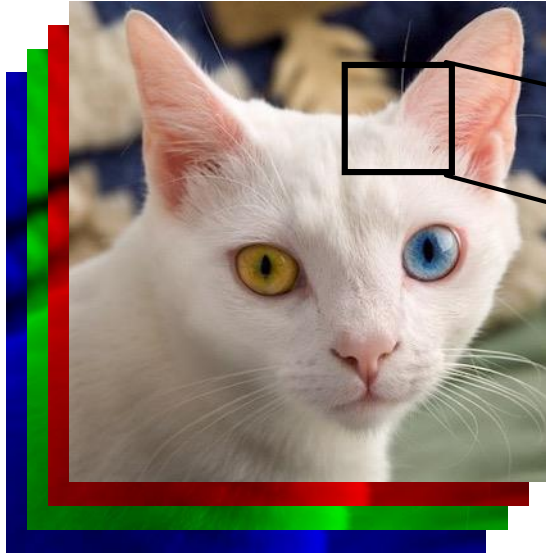


Feature maps (activations)

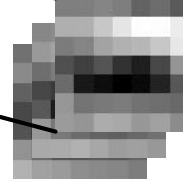


Convolution

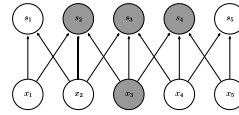
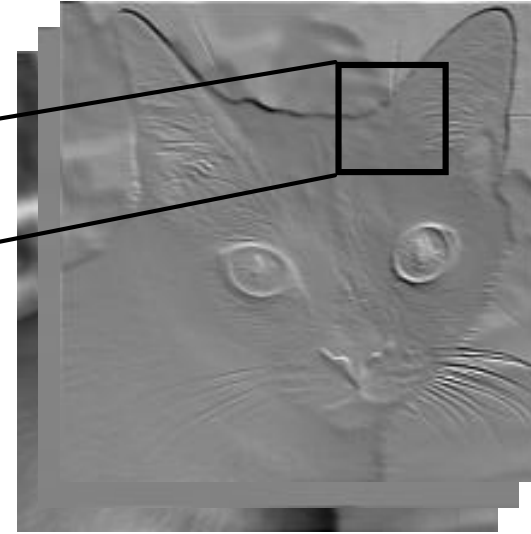
Inputs



Filters



Feature maps (activations)

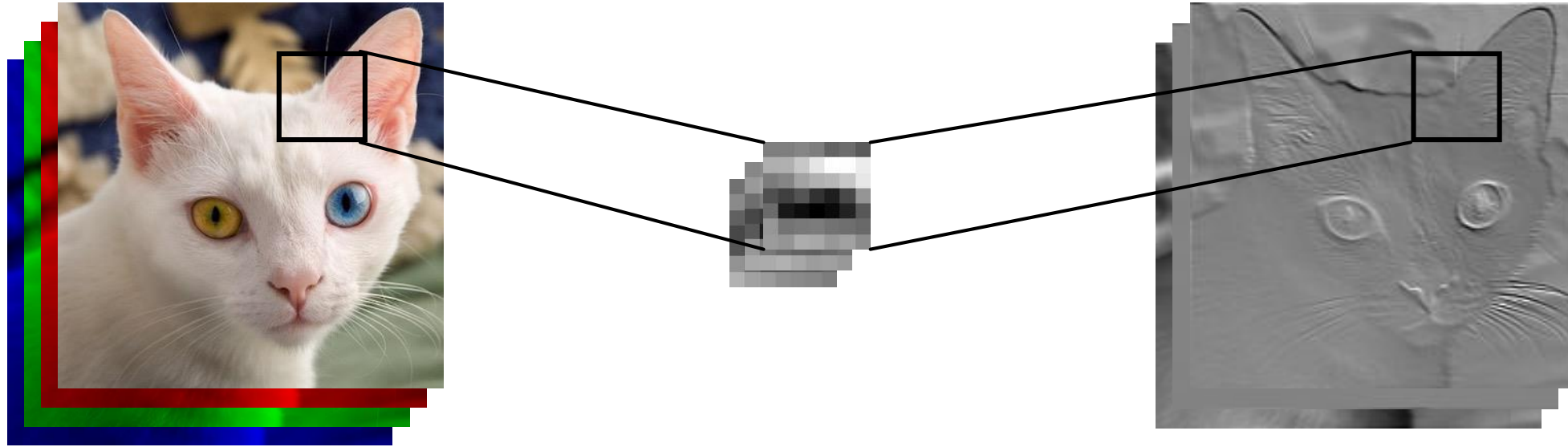
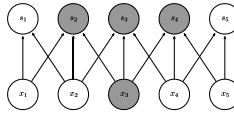


Convolution

Inputs

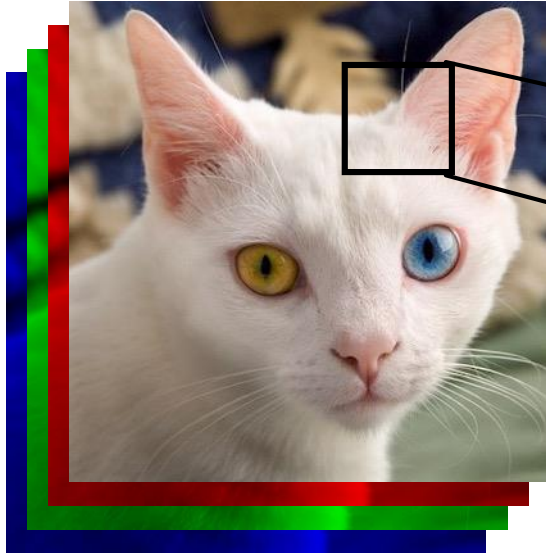
Filters

Feature maps (activations)

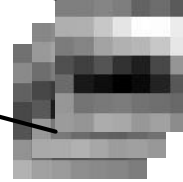


Convolution

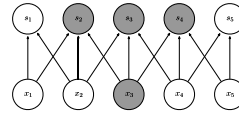
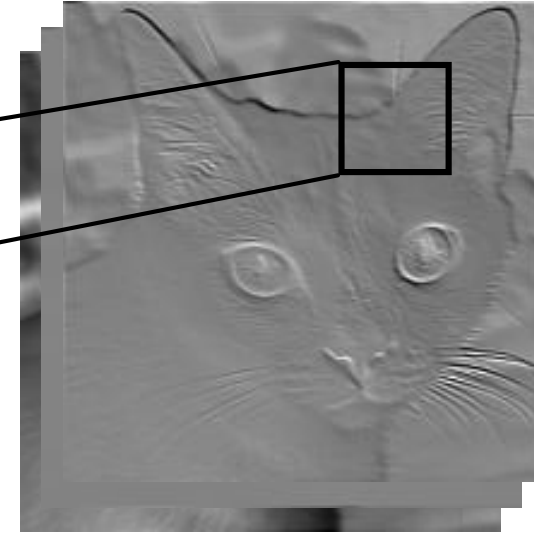
Inputs



Filters



Feature maps (activations)

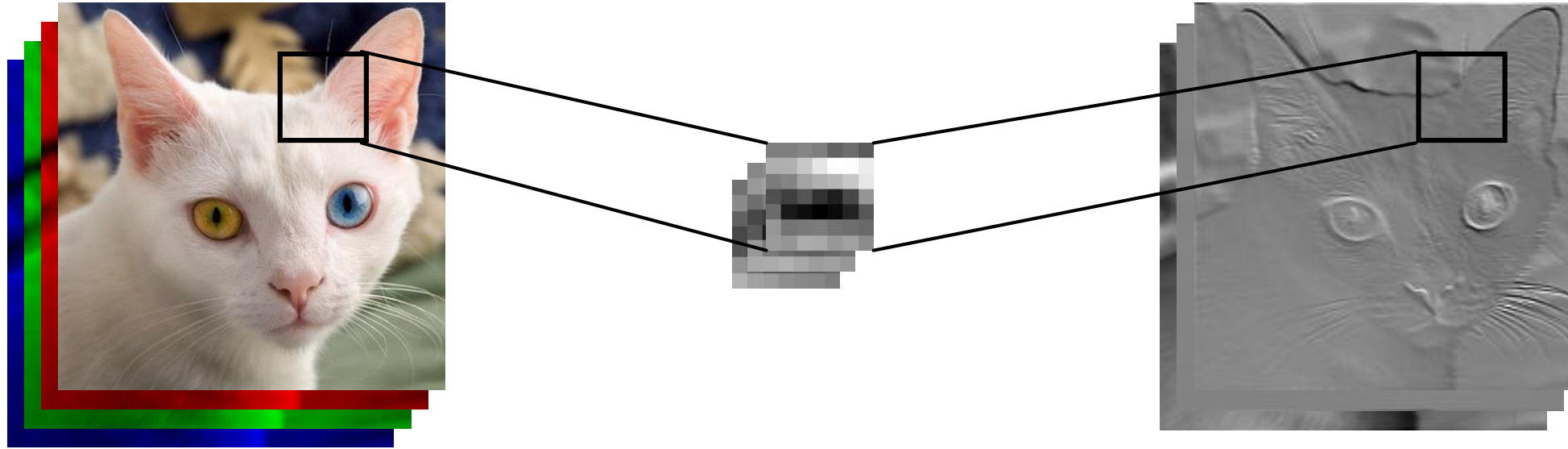
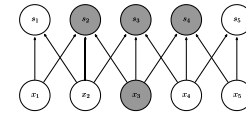


Convolution

Inputs

Filters

Feature maps (activations)

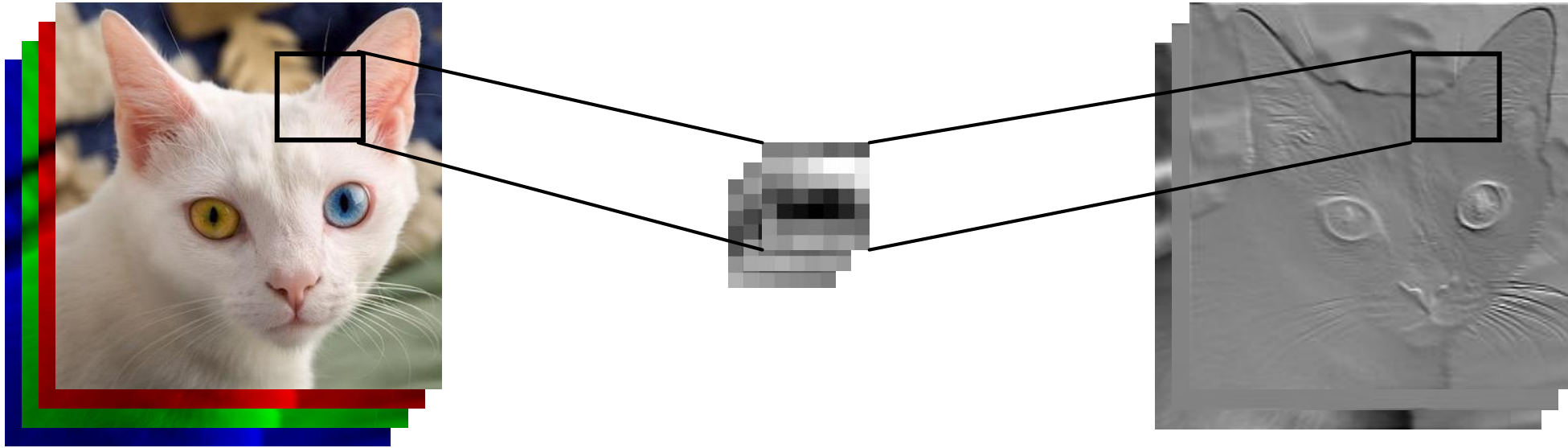
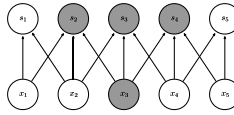


Convolution

Inputs

Filters

Feature maps (activations)

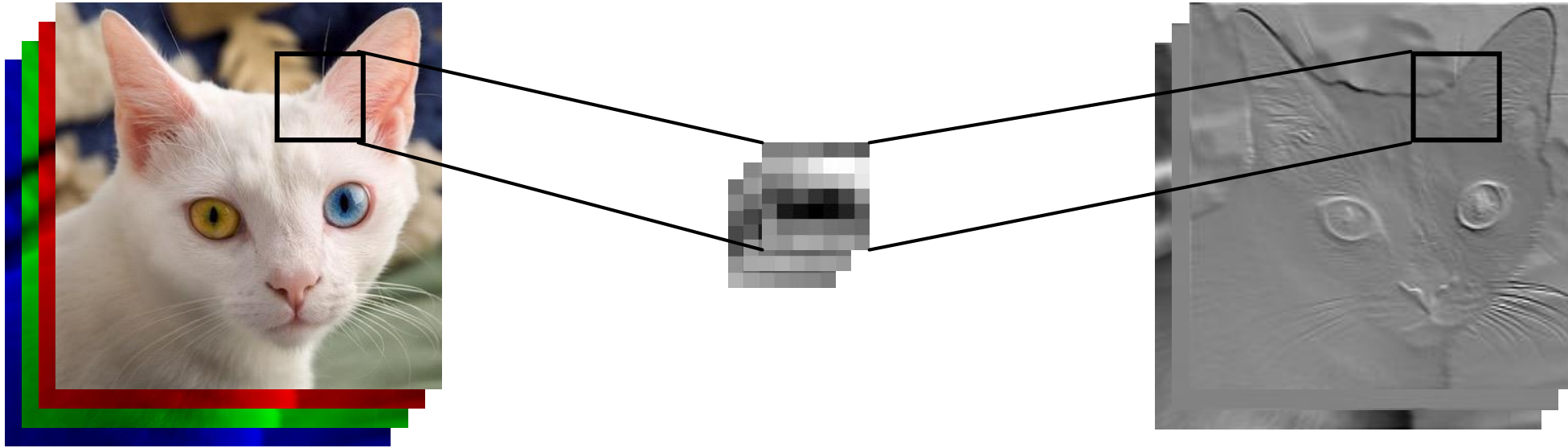
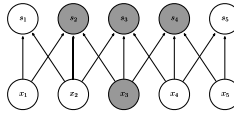


Convolution

Inputs

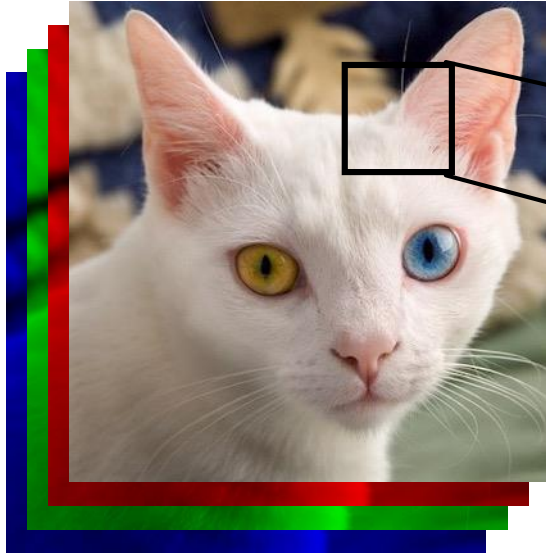
Filters

Feature maps (activations)

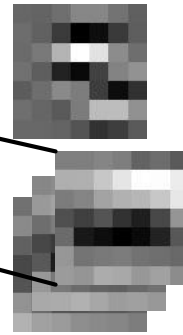


Convolution

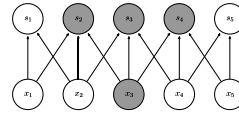
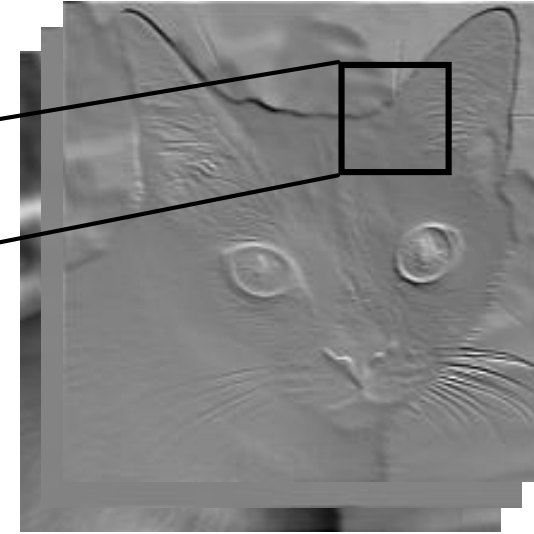
Inputs



Filters

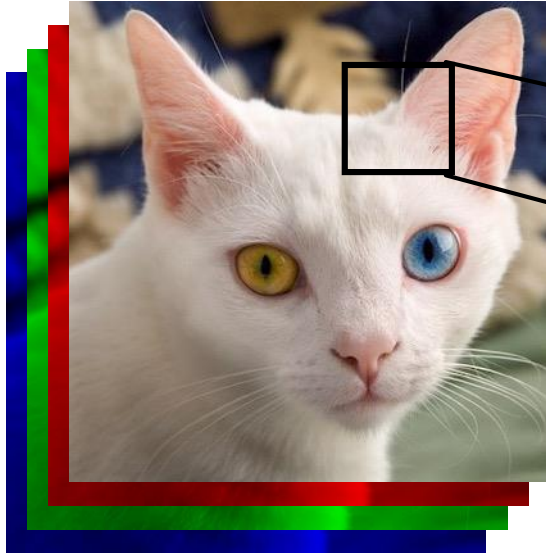


Feature maps (activations)

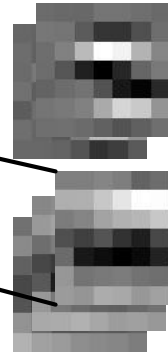


Convolution

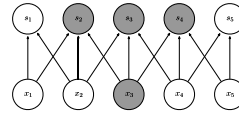
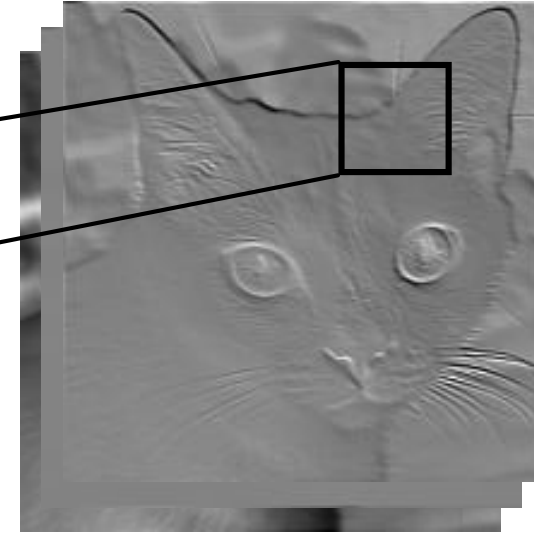
Inputs



Filters

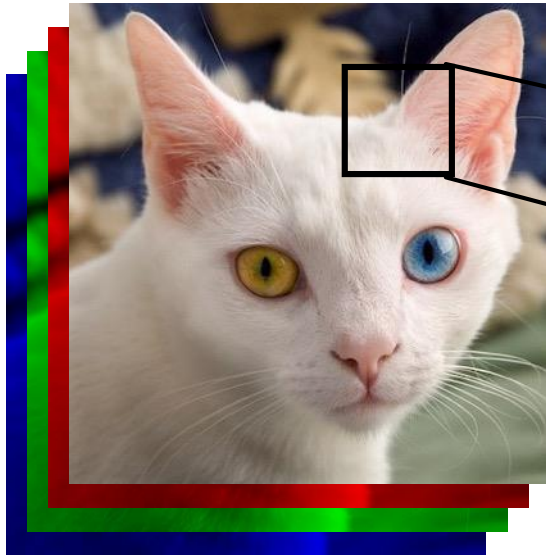


Feature maps (activations)



Convolution

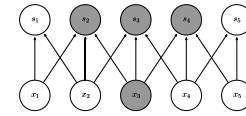
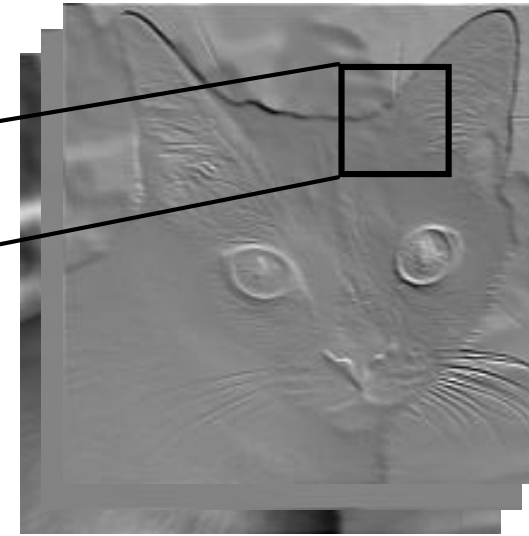
Inputs



Filters

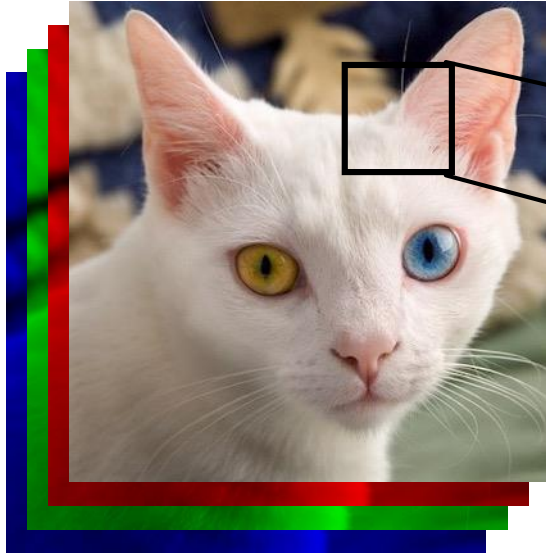


Feature maps (activations)

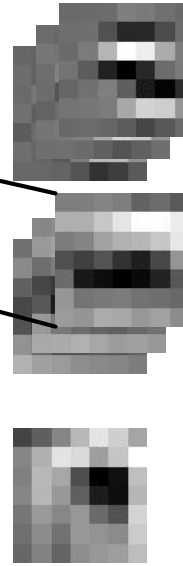


Convolution

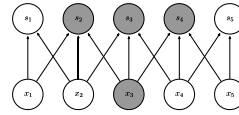
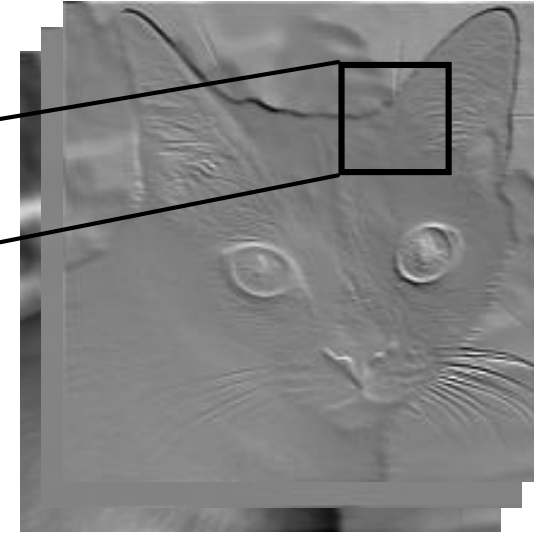
Inputs



Filters

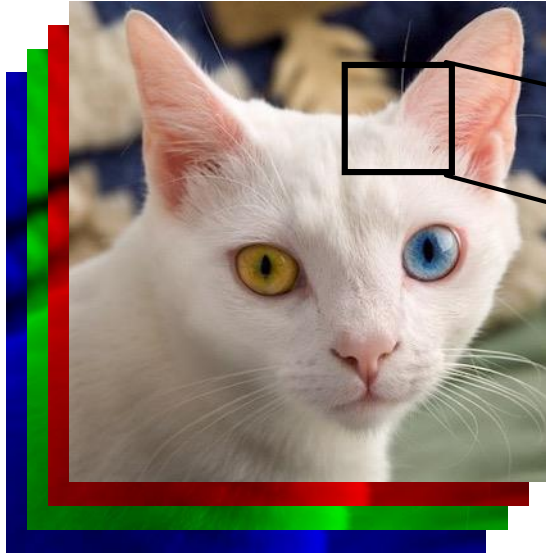


Feature maps (activations)

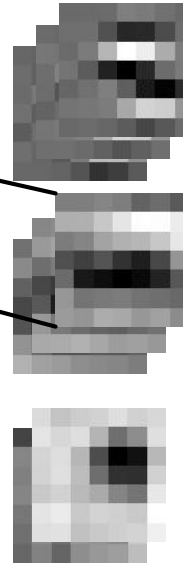


Convolution

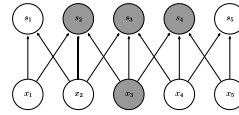
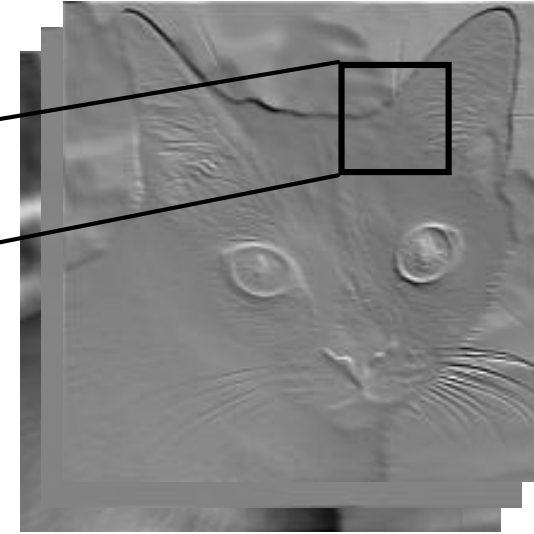
Inputs



Filters

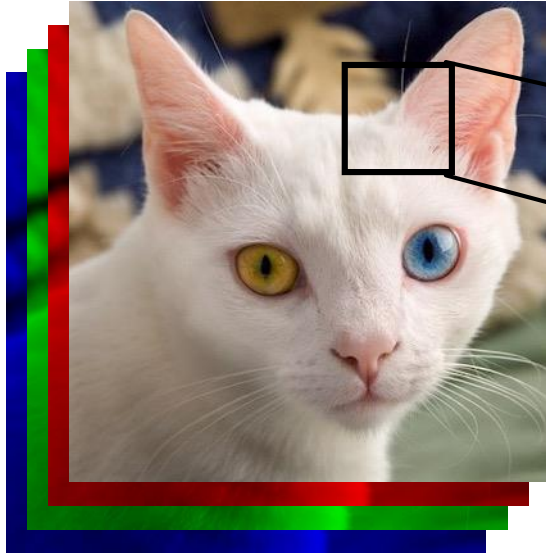


Feature maps (activations)

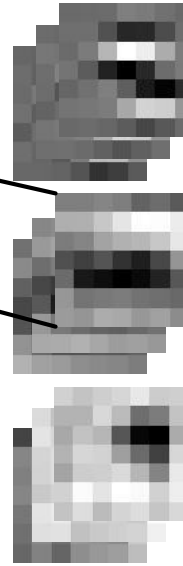


Convolution

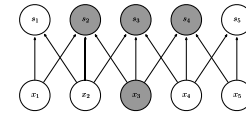
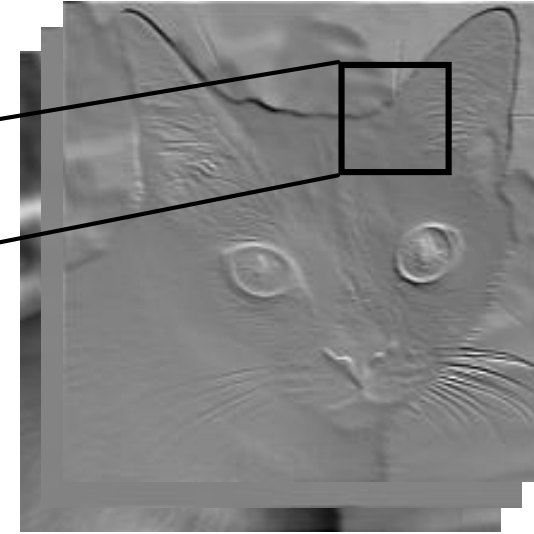
Inputs



Filters

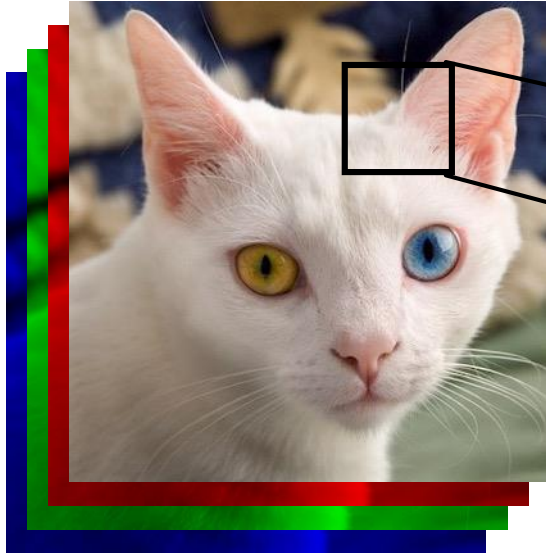


Feature maps (activations)

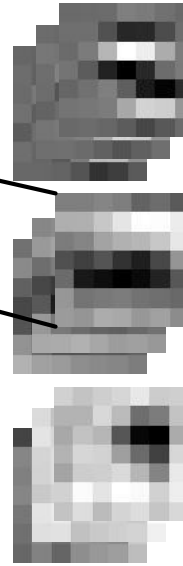


Convolution

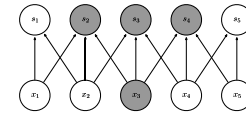
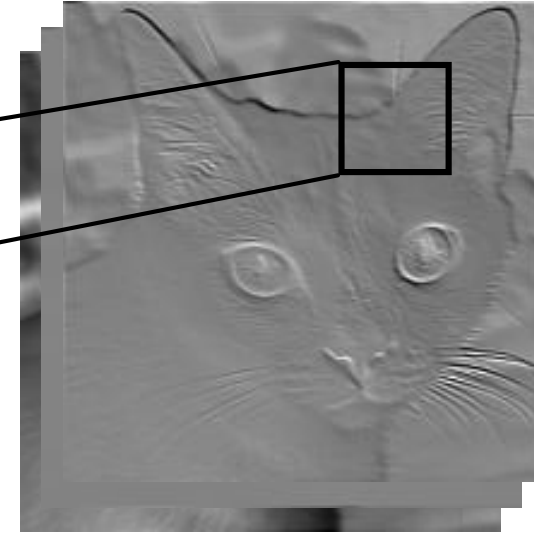
Inputs



Filters

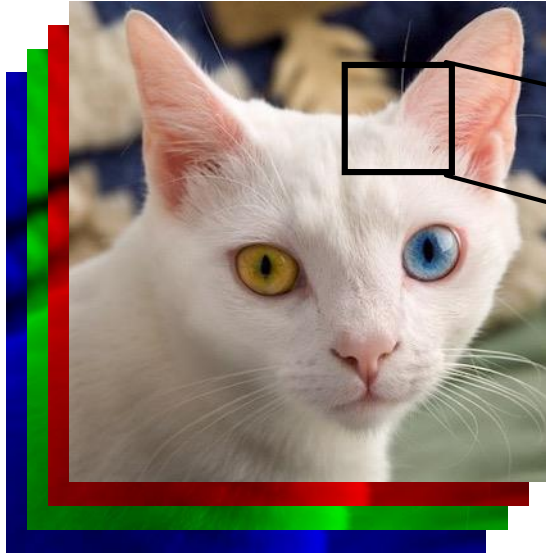


Feature maps (activations)

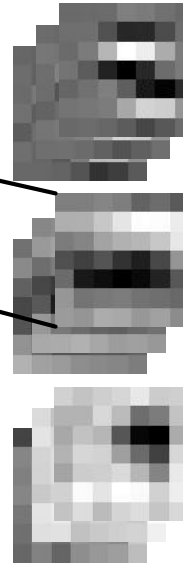


Convolution

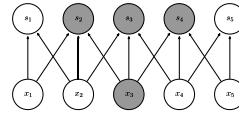
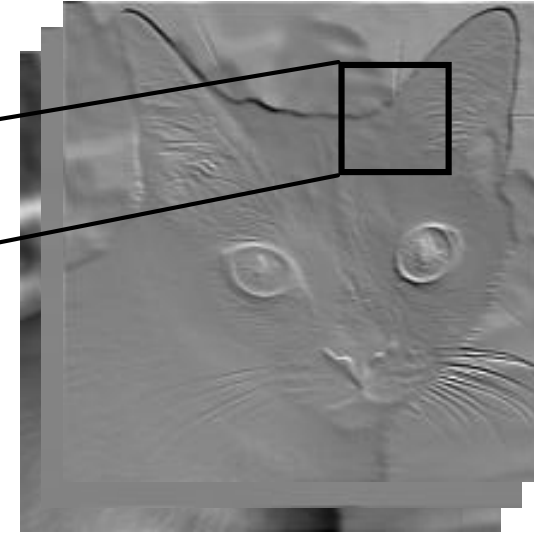
Inputs



Filters



Feature maps (activations)

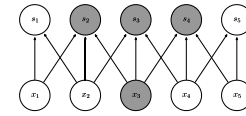


Convolution

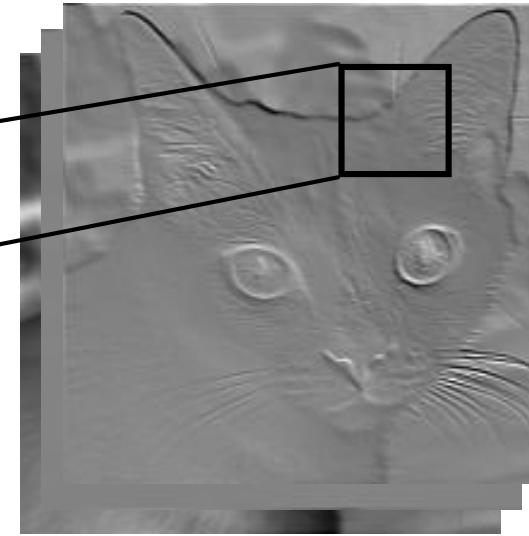
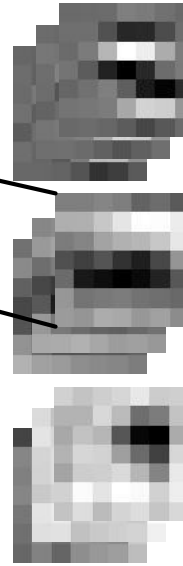
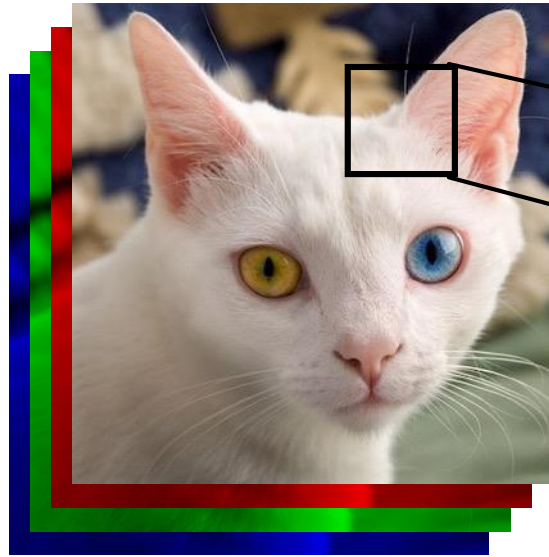
Inputs

Filters

Feature maps (activations)



N

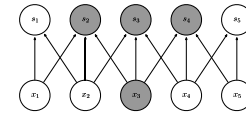


Convolution

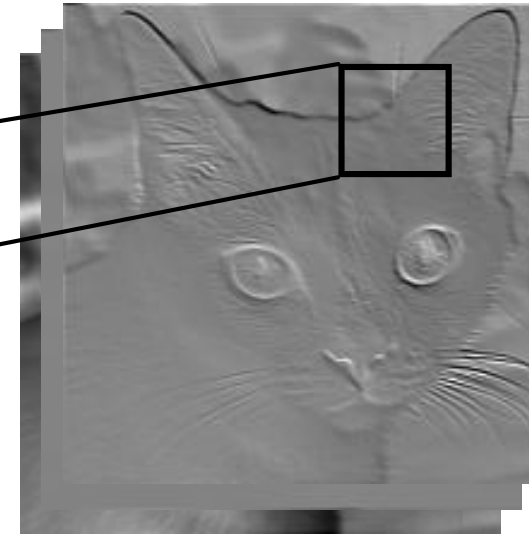
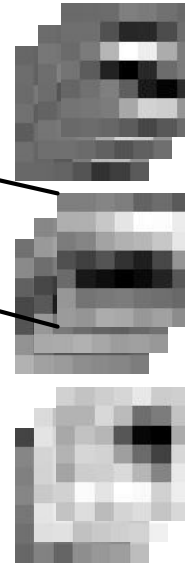
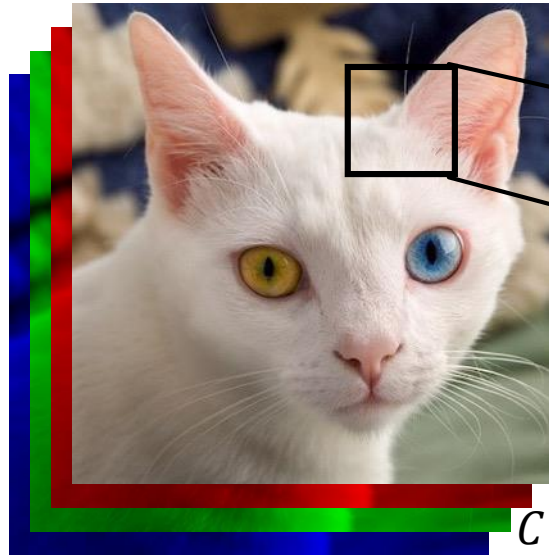
Inputs

Filters

Feature maps (activations)



N

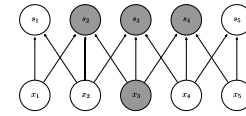


Convolution

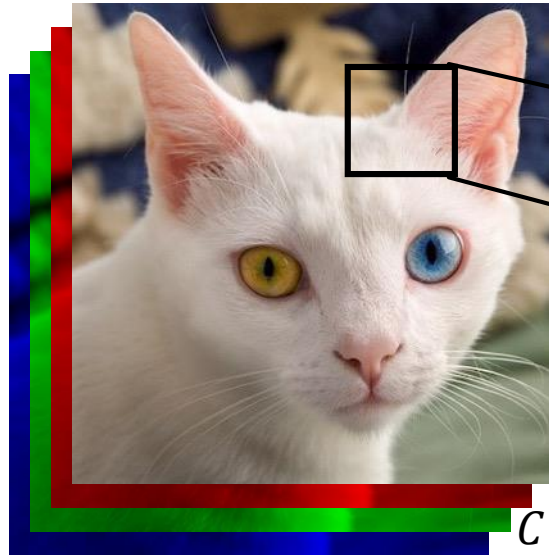
Inputs

Filters

Feature maps (activations)

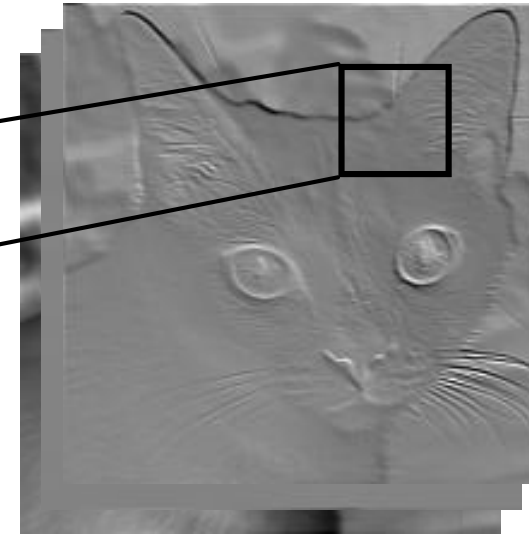
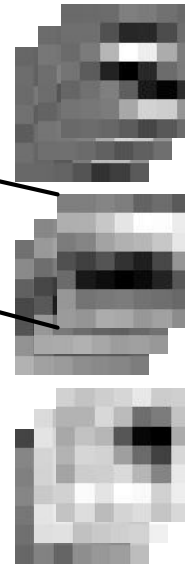


N

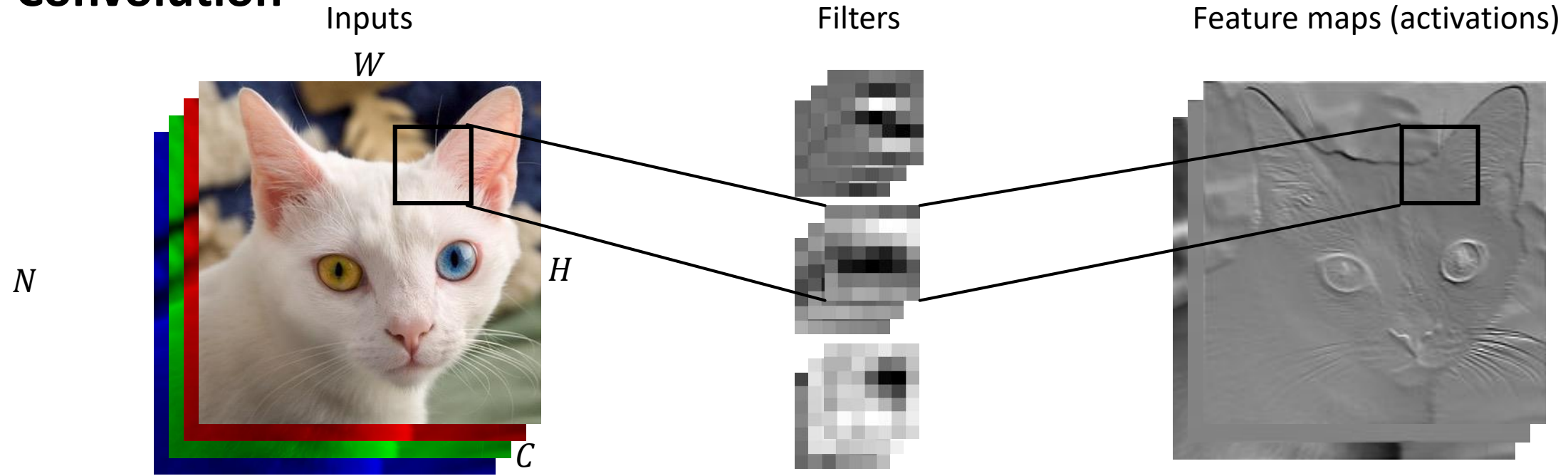
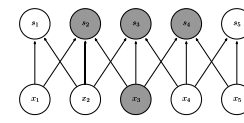


H

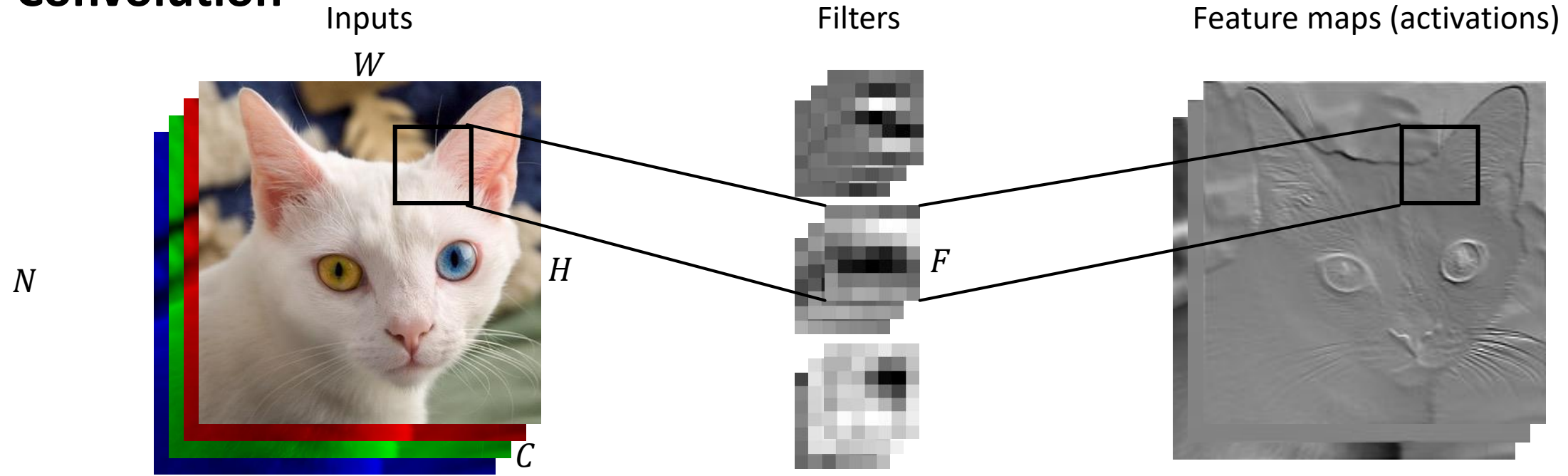
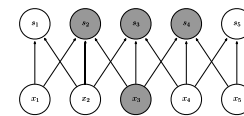
C



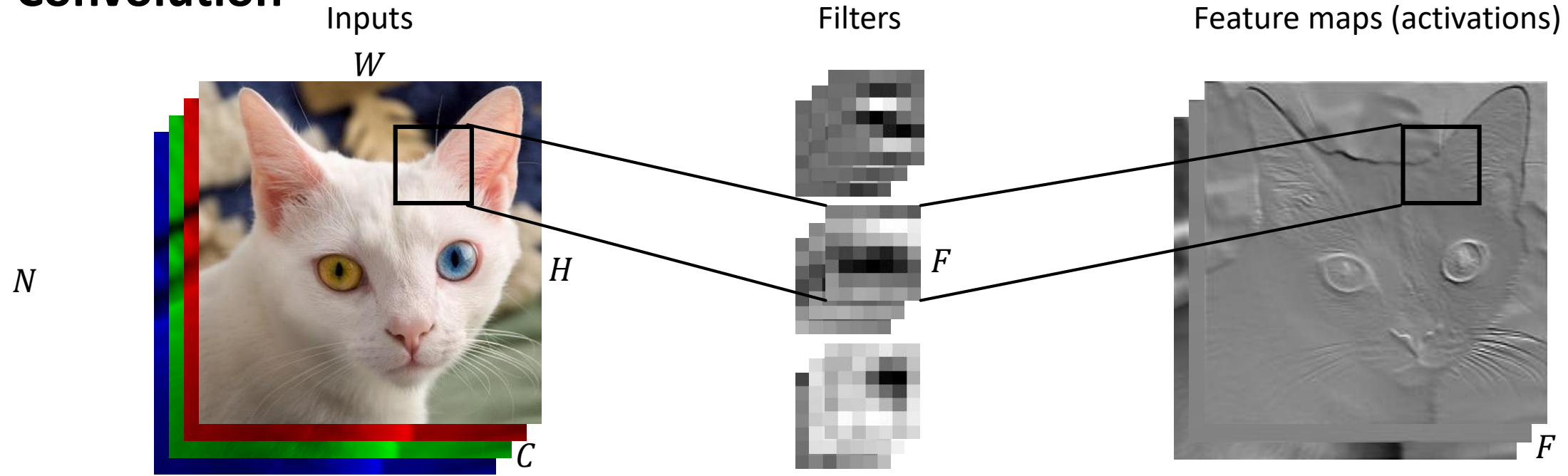
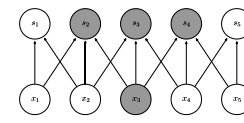
Convolution



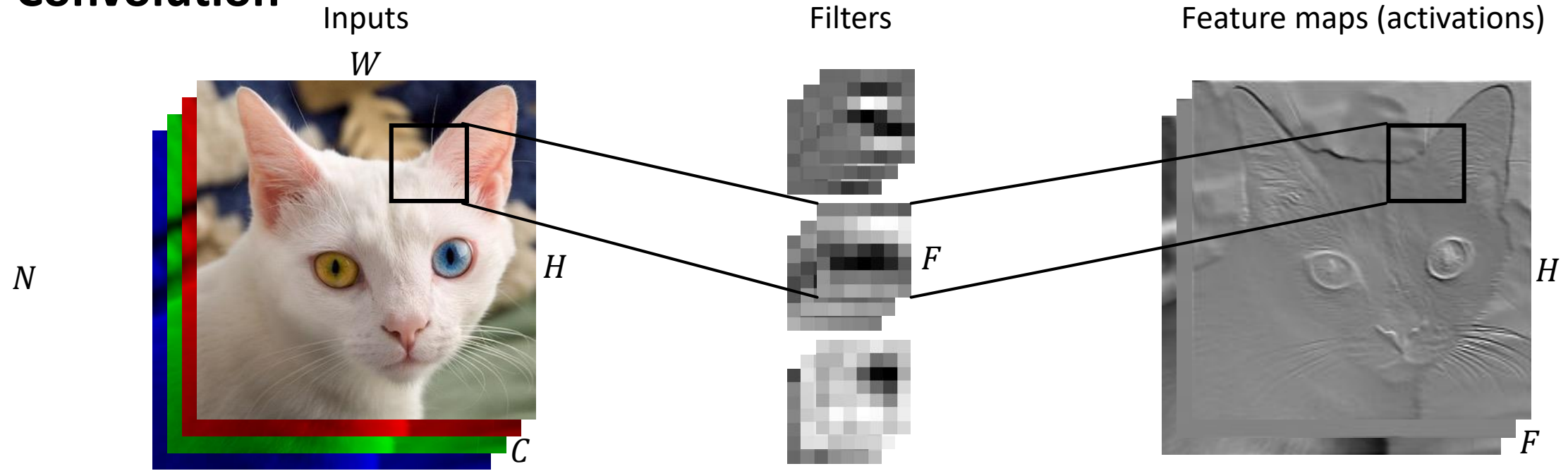
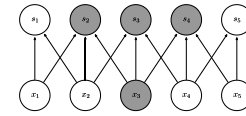
Convolution



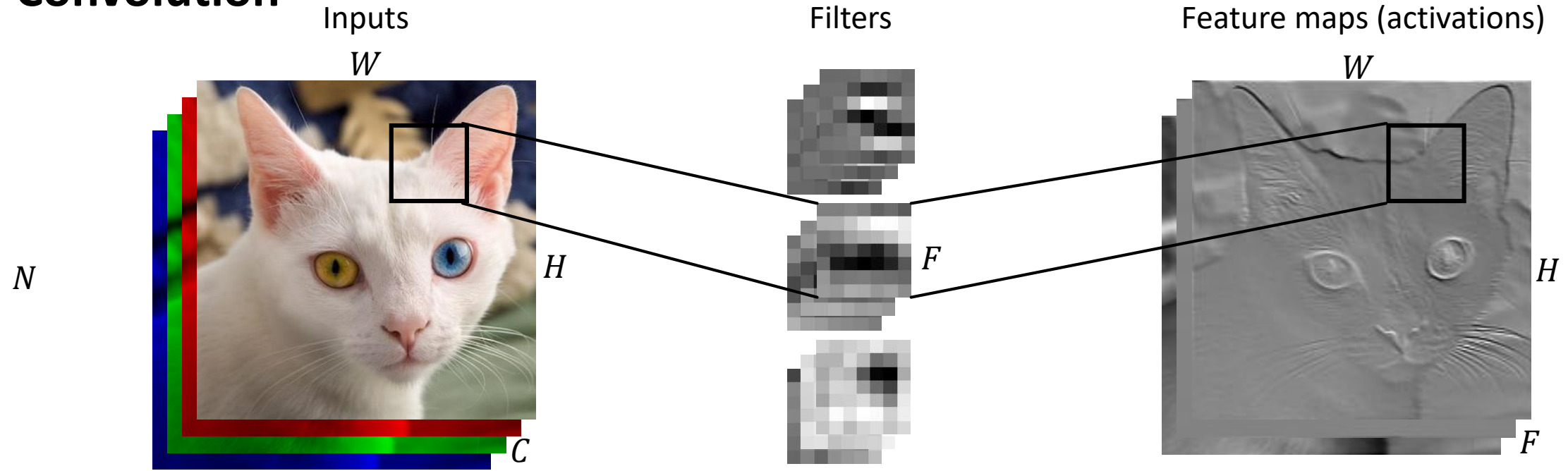
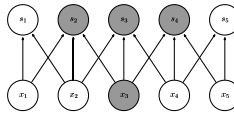
Convolution



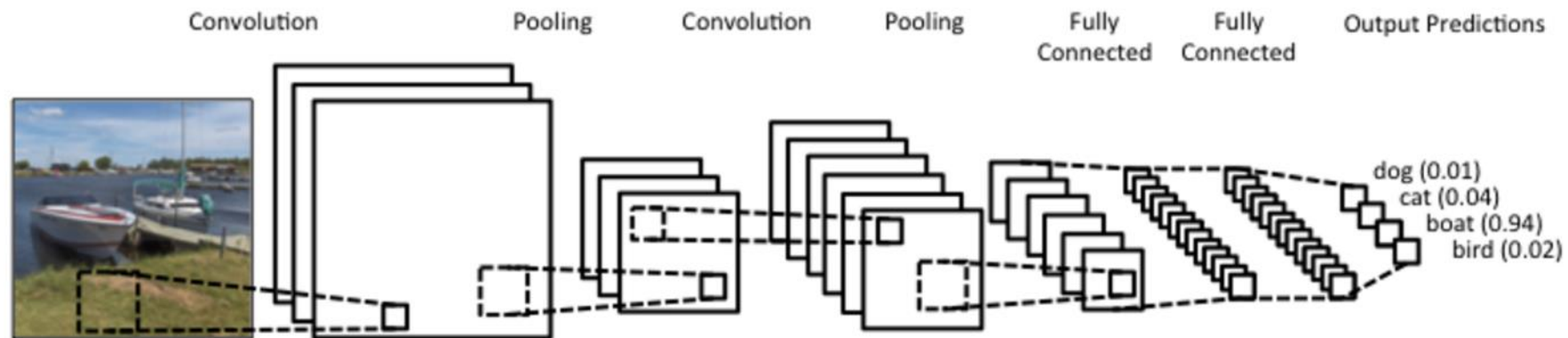
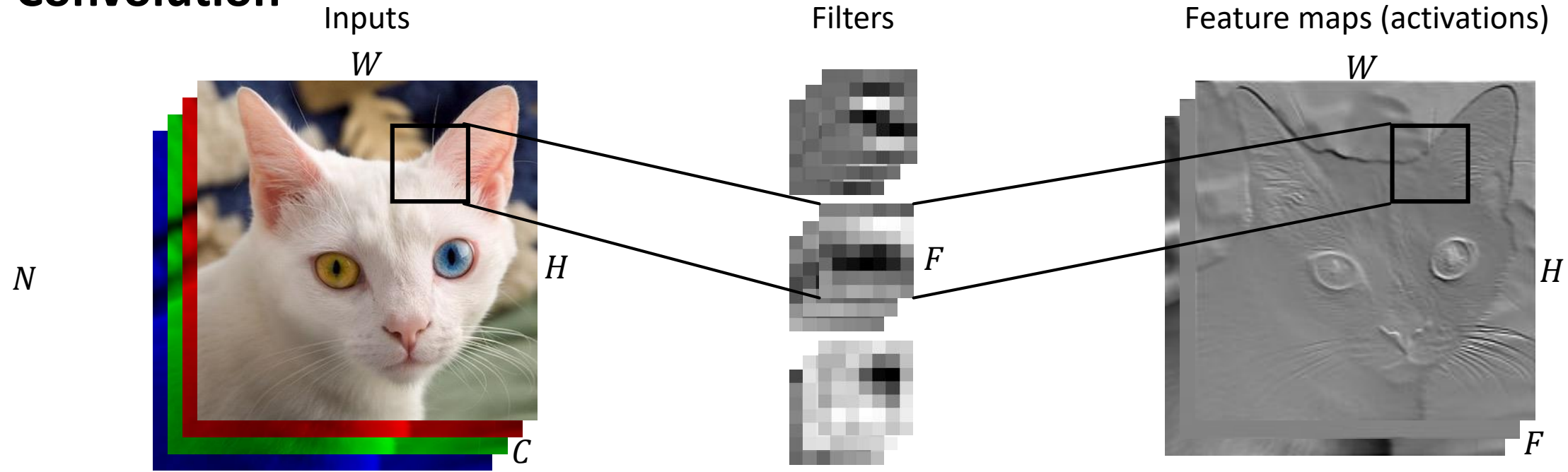
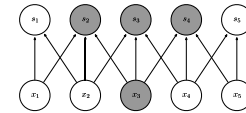
Convolution



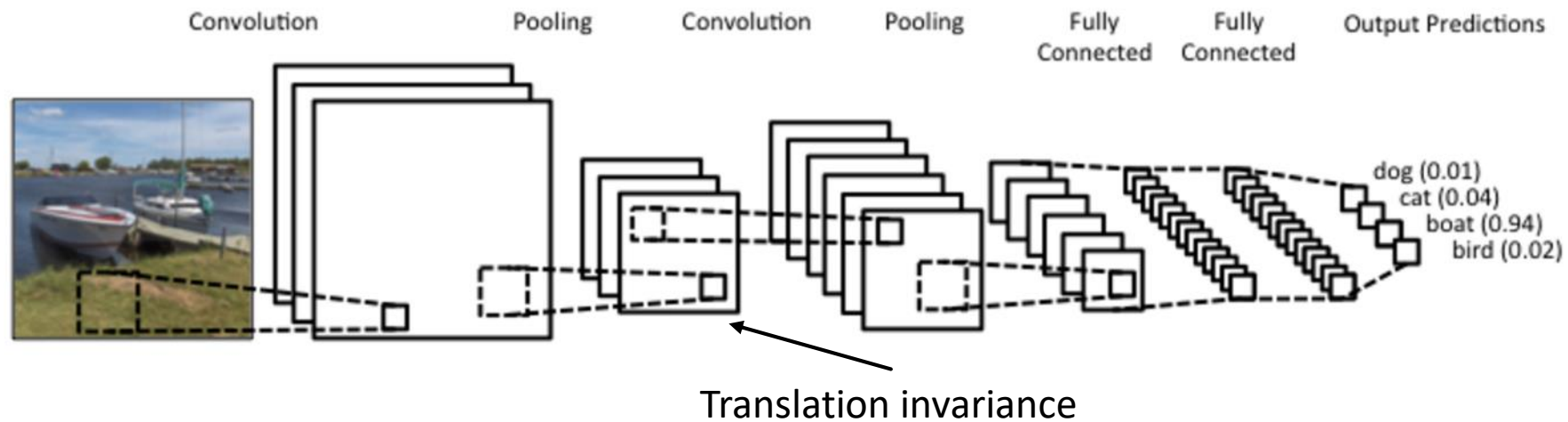
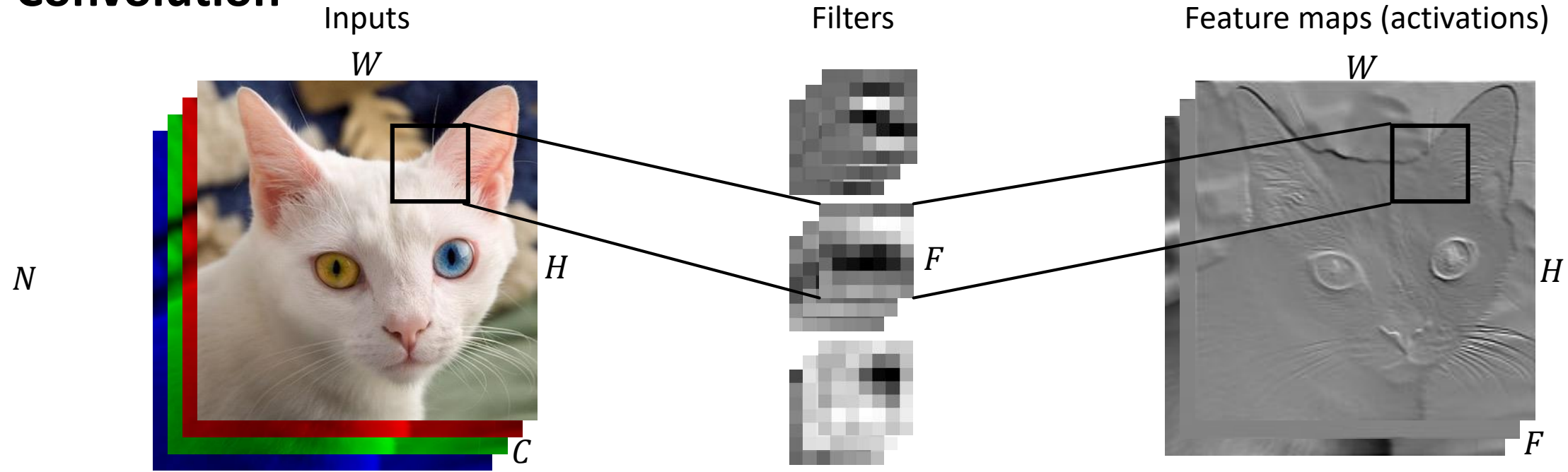
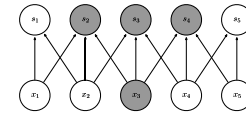
Convolution



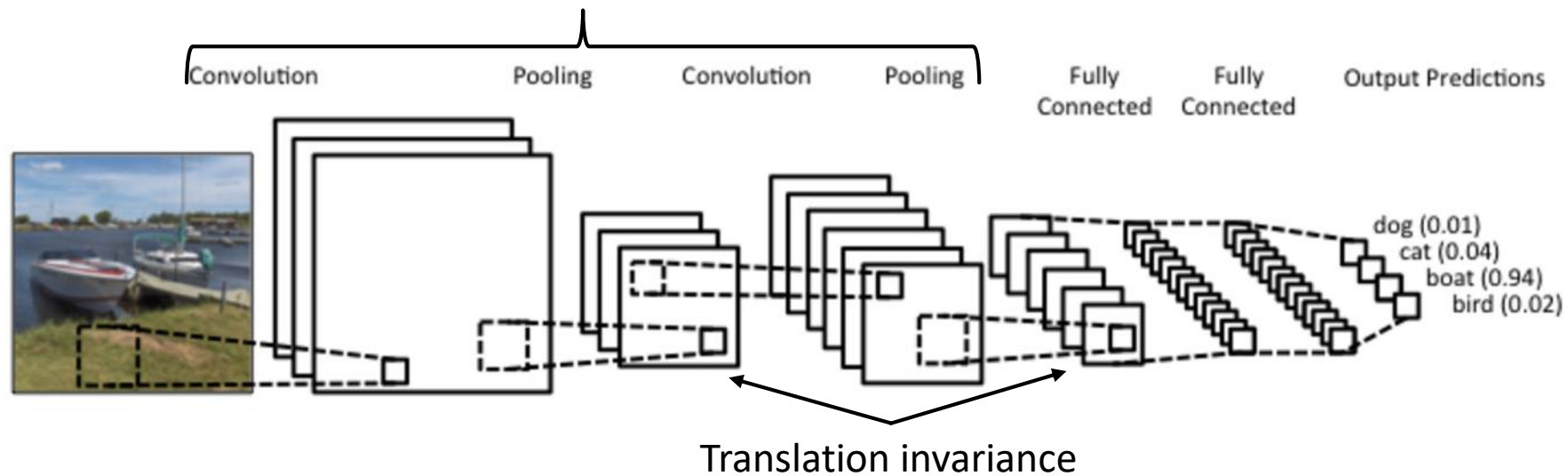
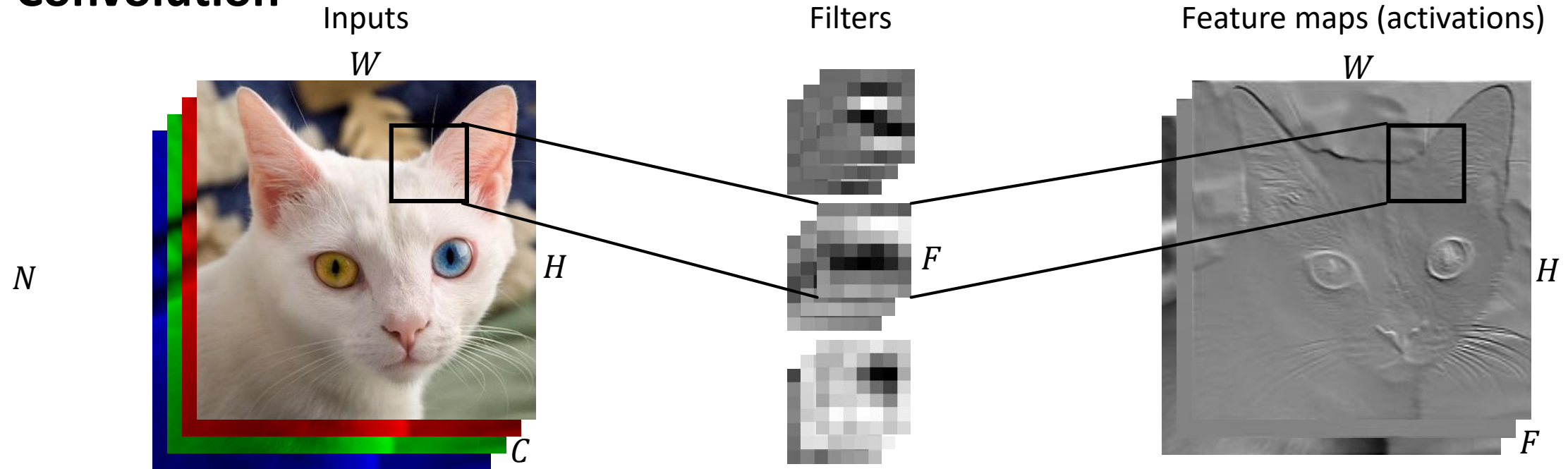
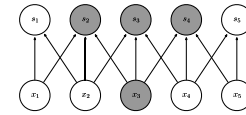
Convolution



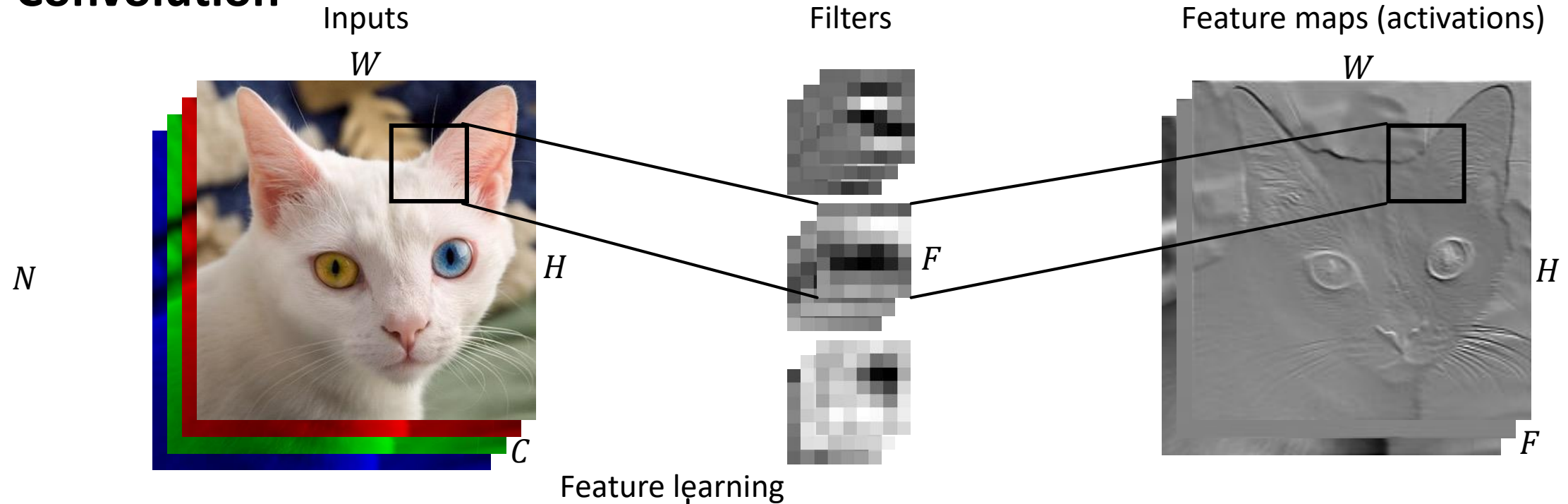
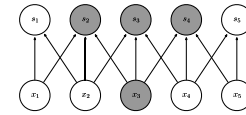
Convolution



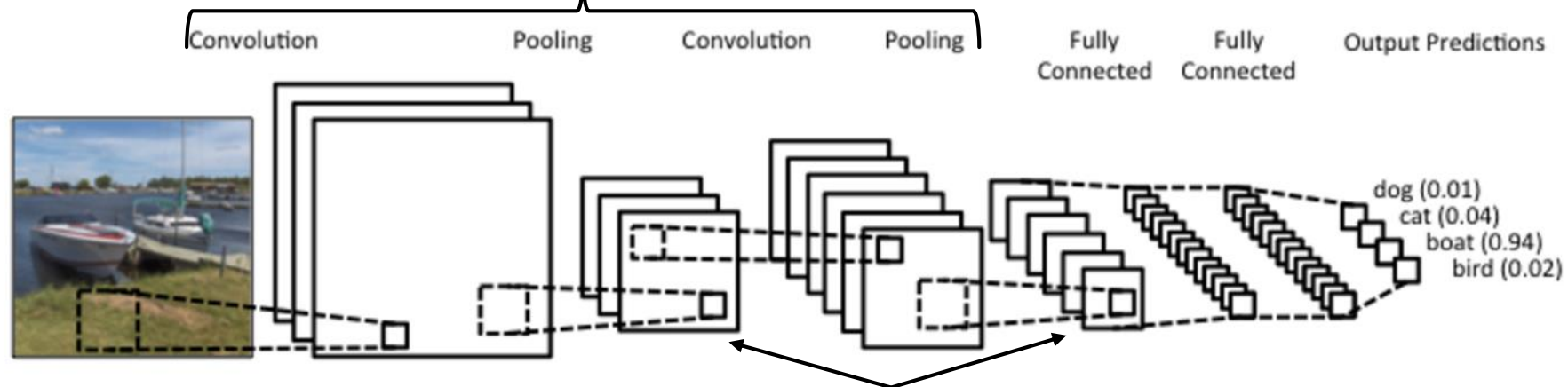
Convolution



Convolution

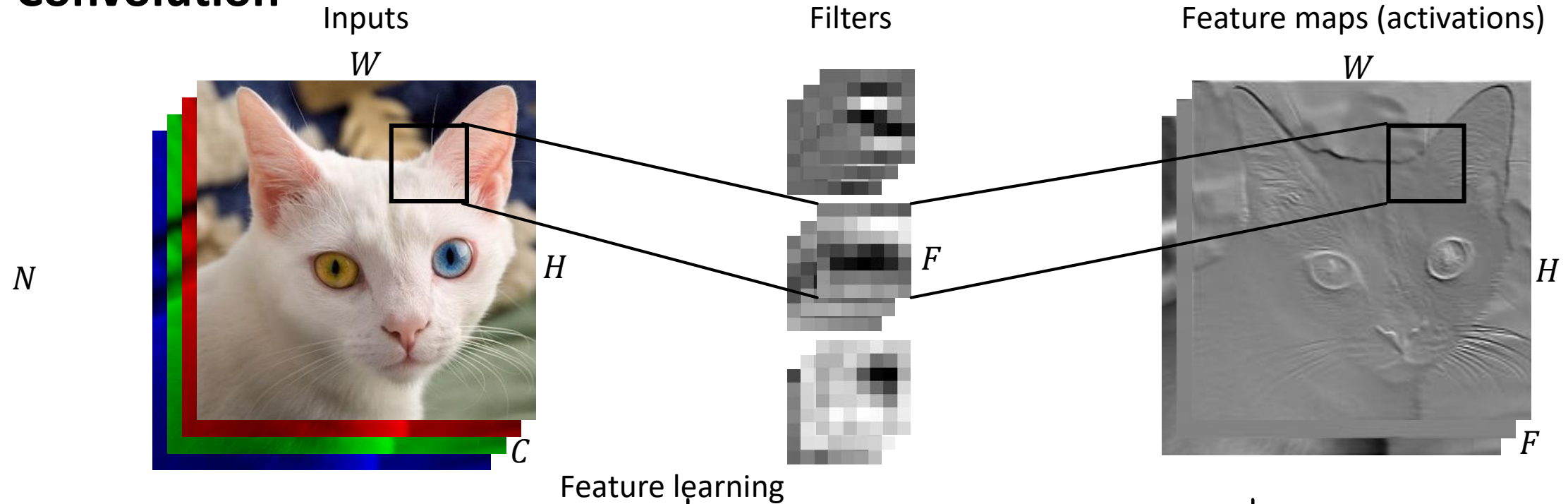
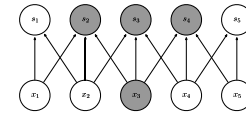


Feature learning

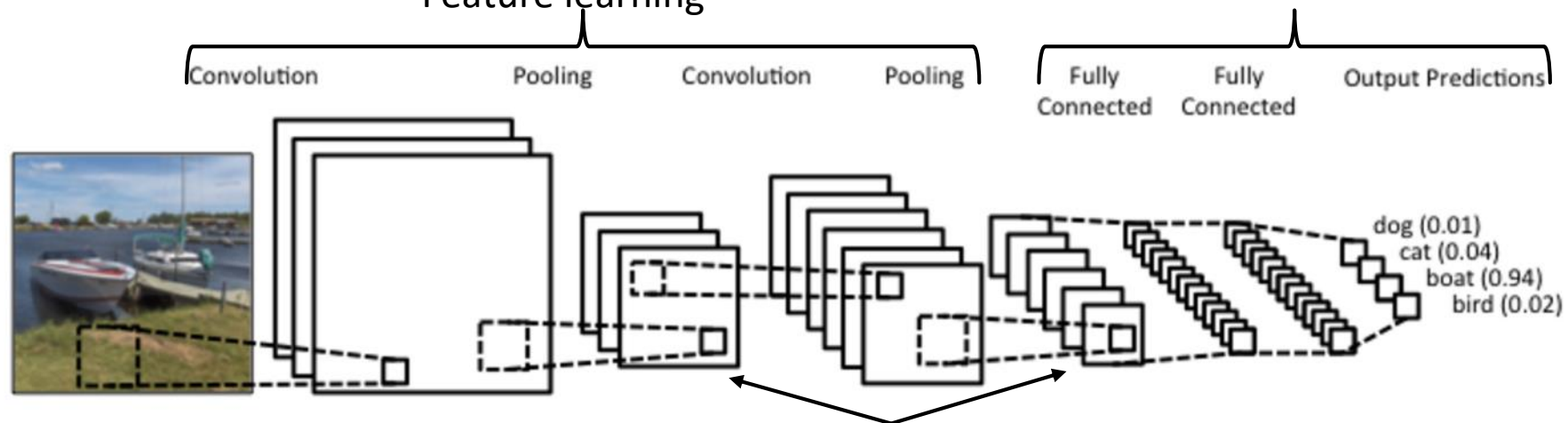


Translation invariance

Convolution

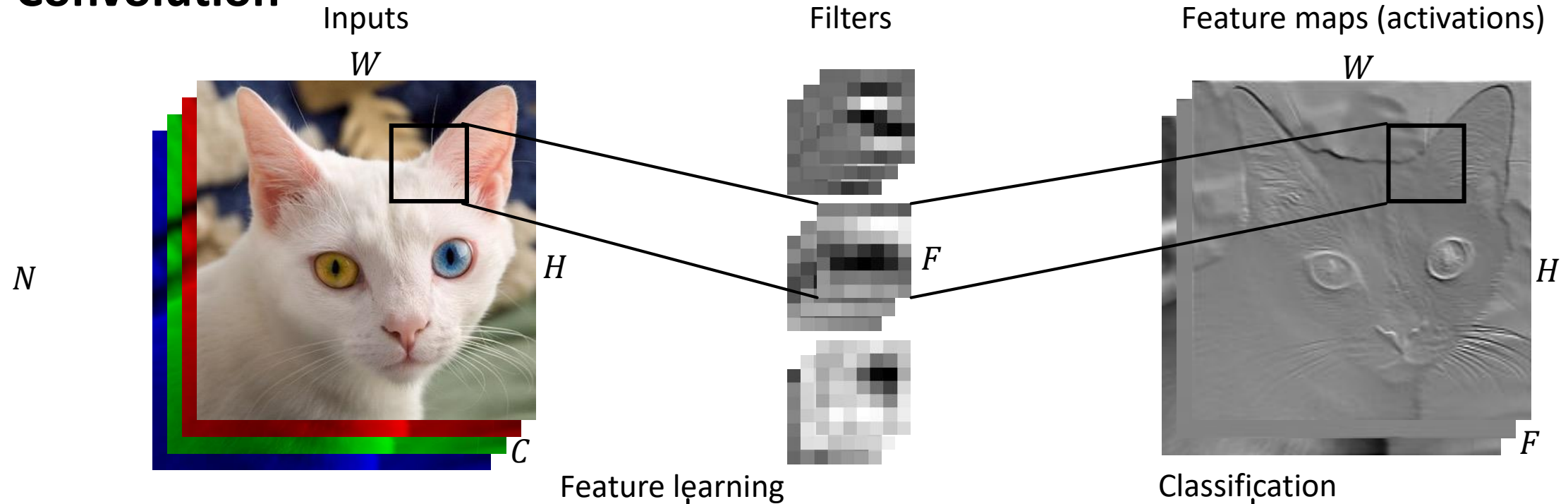
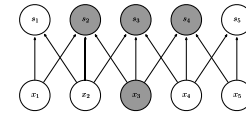


Feature learning



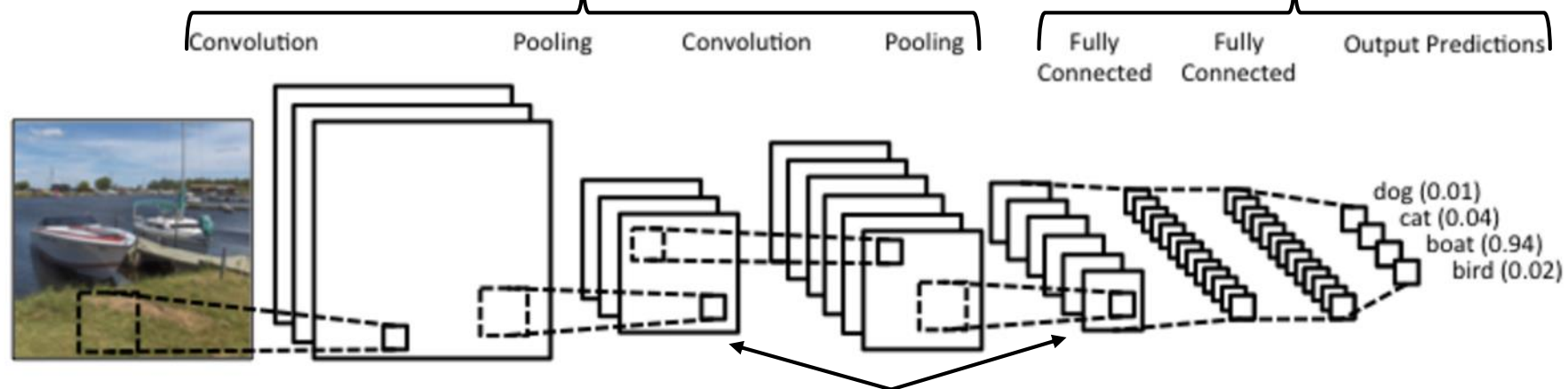
Translation invariance

Convolution

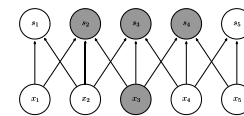


Feature learning

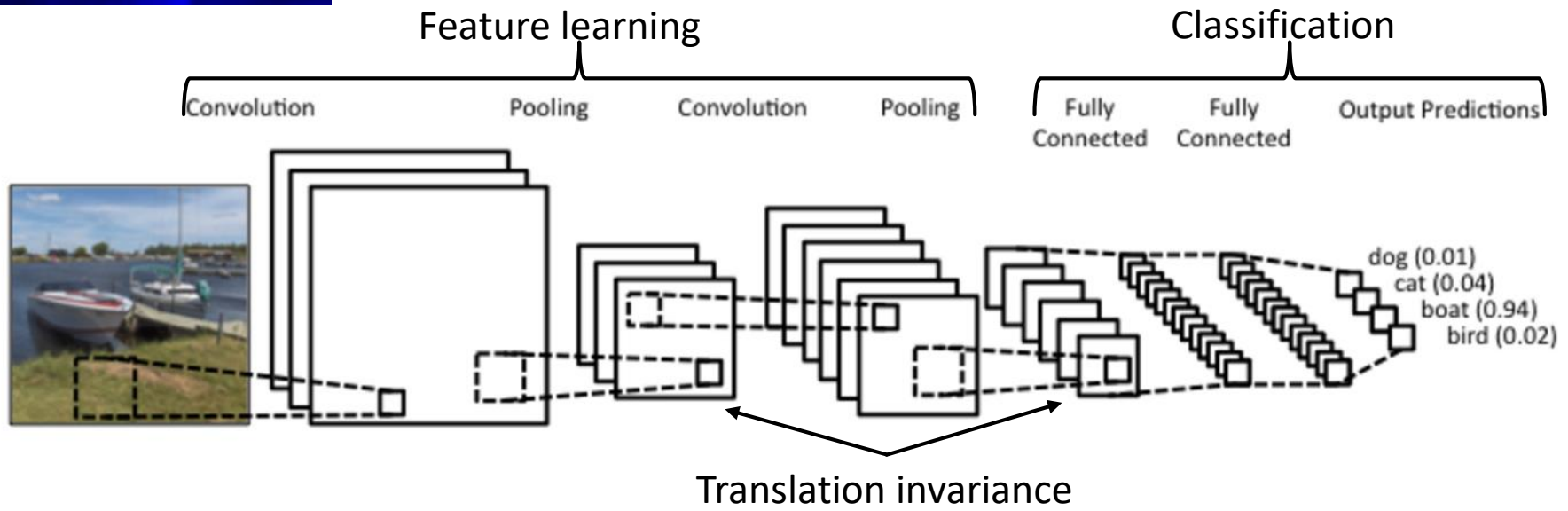
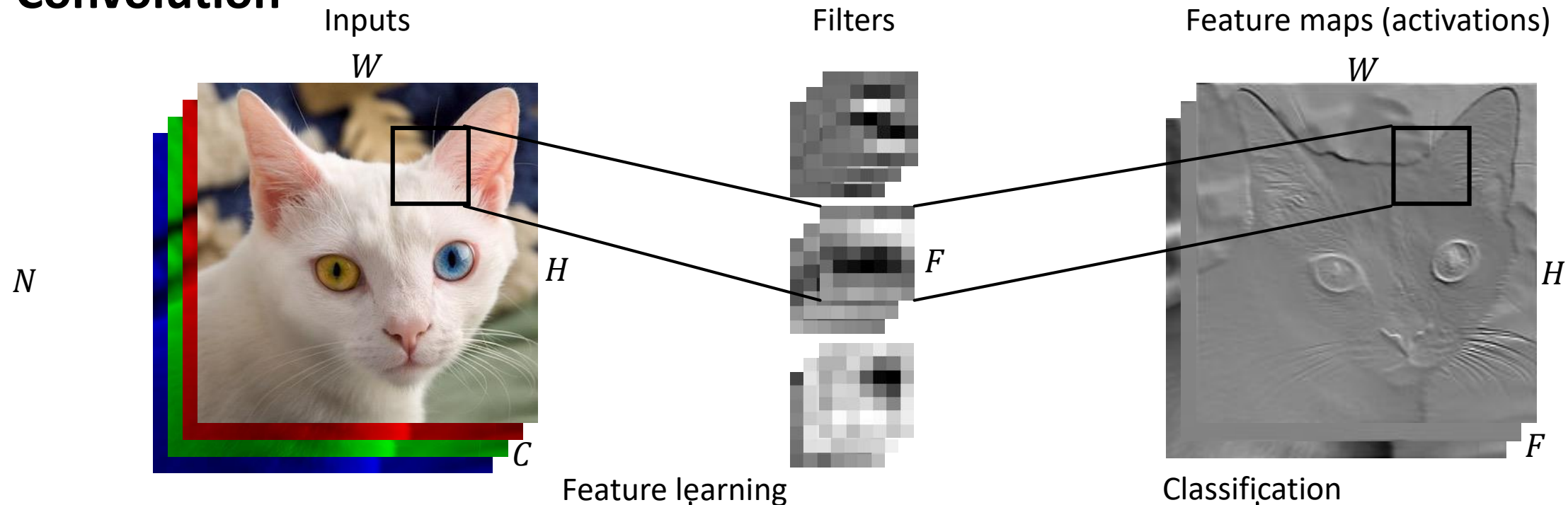
Classification



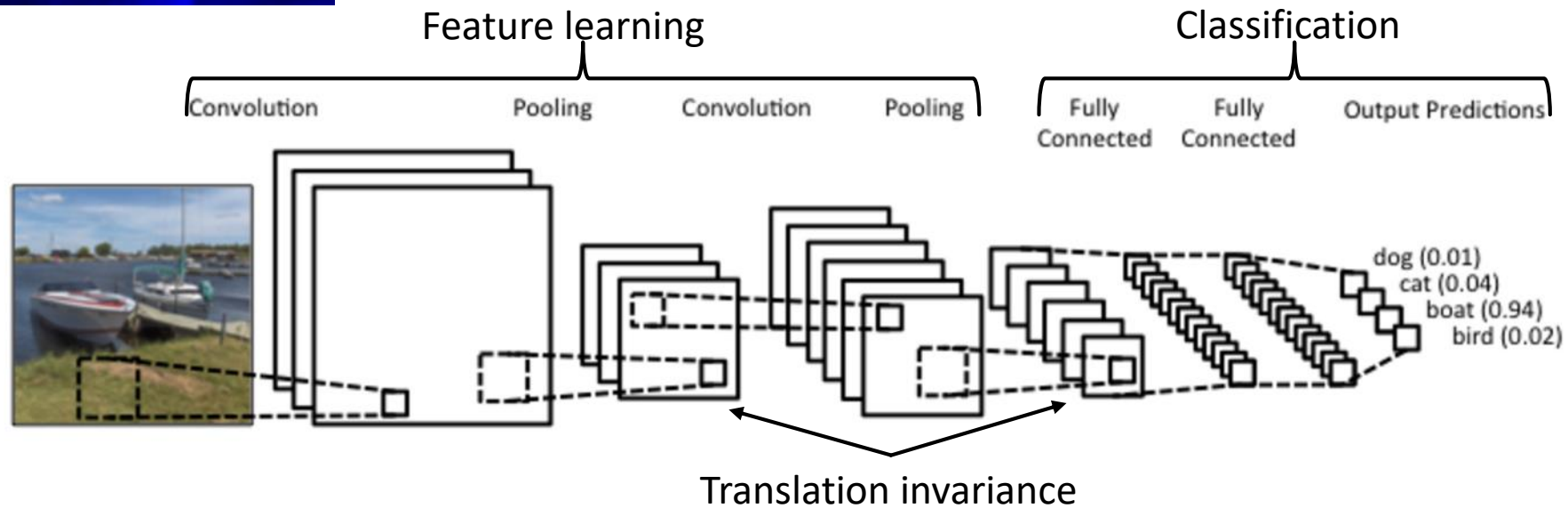
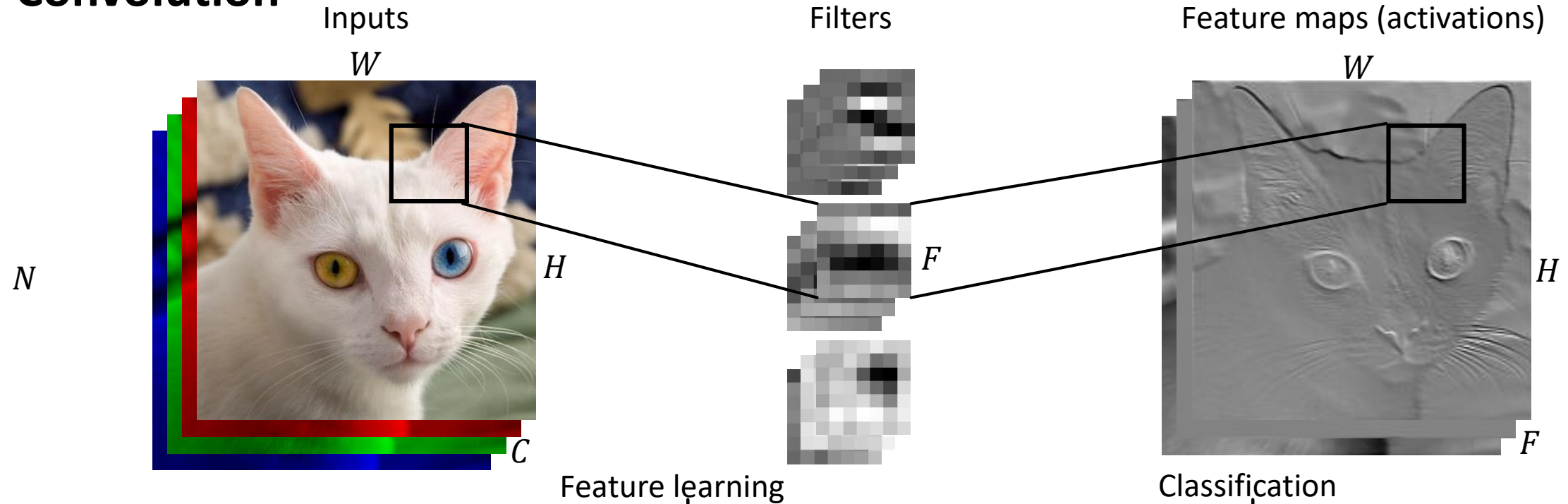
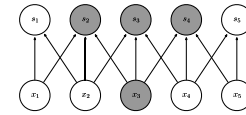
Translation invariance



Convolution



Convolution



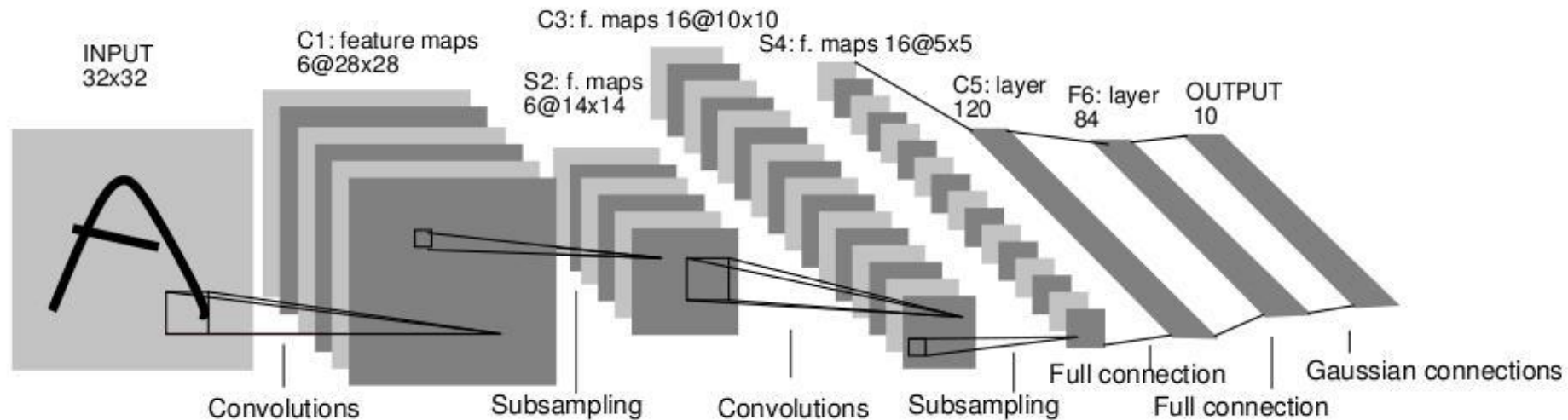
Operators



Operators

A short history of (old) CNNs

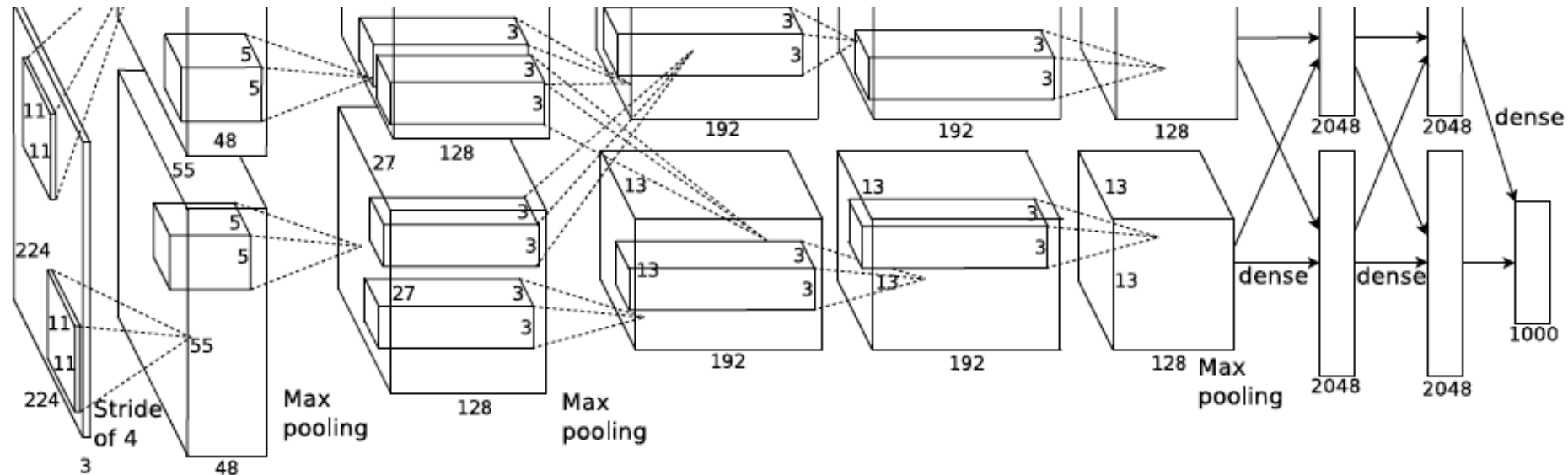
LeNet-5 [LeCun, Bottou, Bengio, & Haffner 1998]



- Average pooling
- Sigmoid/tanh nonlinearities
- Fully-connected layers at end
- Trained on MNIST (60k samples)

A short history of (old) CNNs

AlexNet [Krizhevsky, Sutskever, & Hinton 2012]



- Max pooling
- ReLU nonlinearities
- Deeper, bigger model
- Dropout
- Trained on ImageNet (1.2M images)
- GPU implementation (2 GPUs for a week)

A short history of (old) CNNs

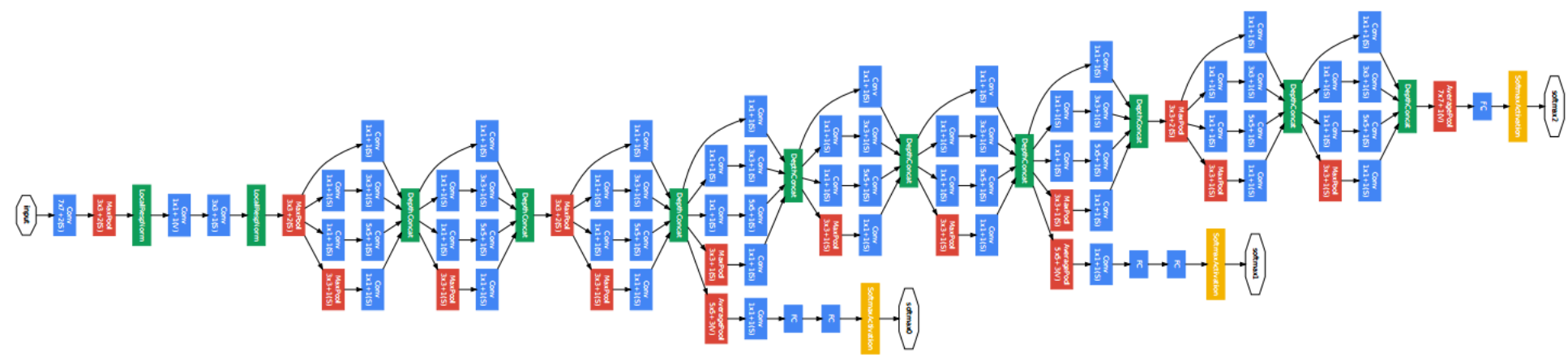
AlexNet [Krizhevsky, Sutskever, & Hinton 2012]



- Deeper, bigger model
- Dropout
- Trained on ImageNet (1.2M images)
- GPU implementation (2 GPUs for a week)

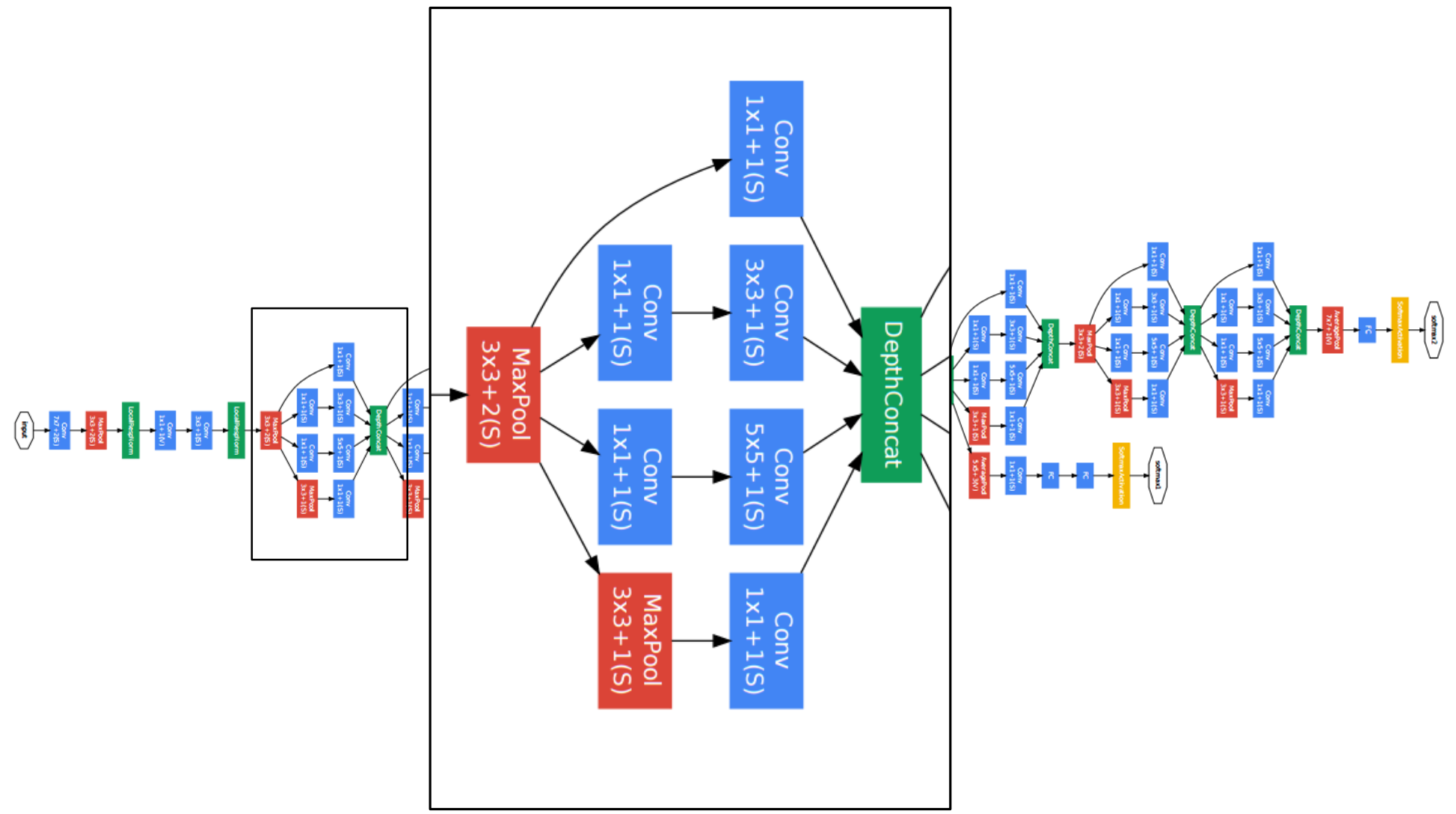
A short history of (old) CNNs

GoogLeNet (AKA Inception) [Szegedy et al. 2015]



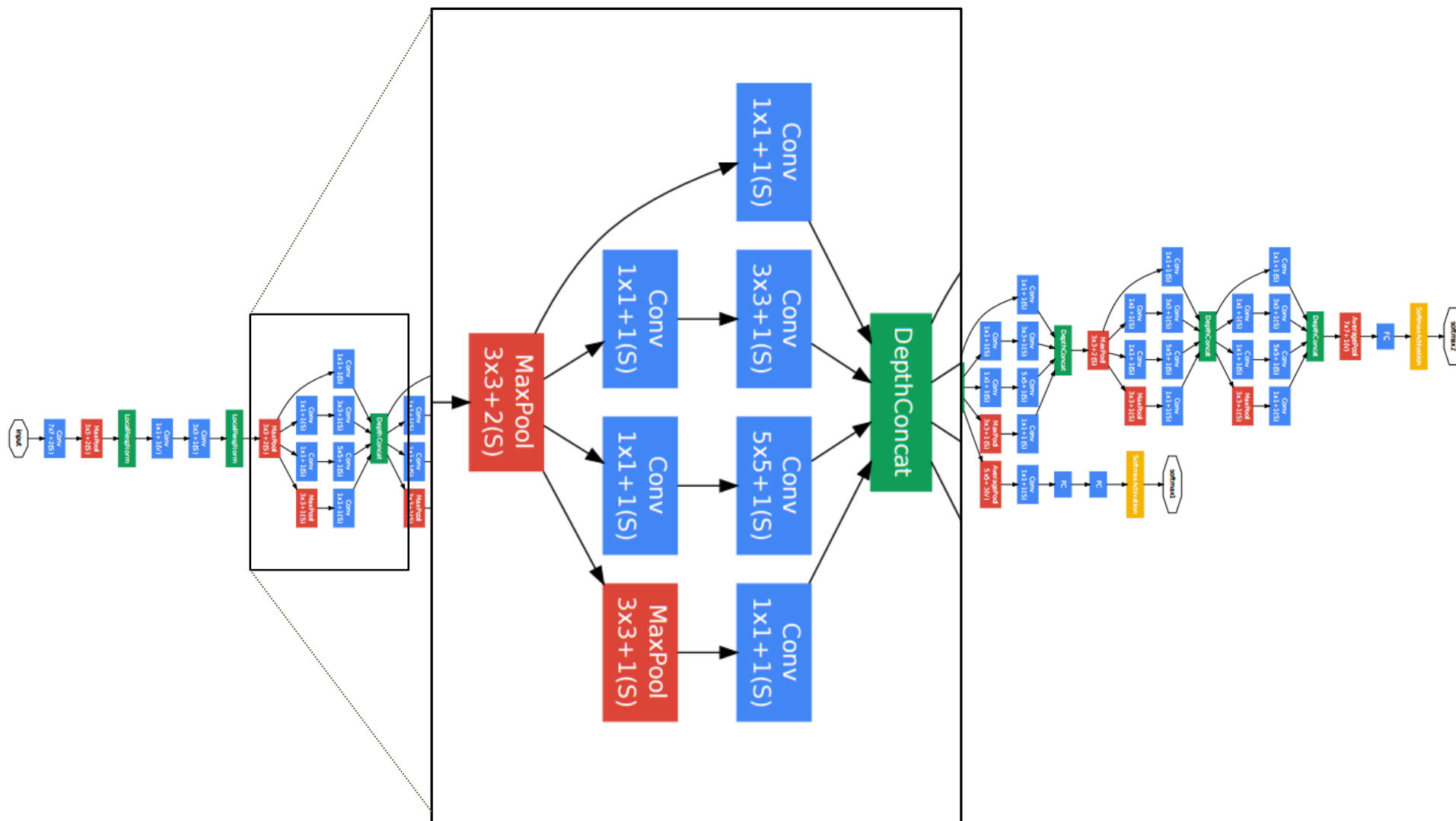
A short history of (old) CNNs

GoogLeNet (AKA Inception) [Szegedy et al. 2015]



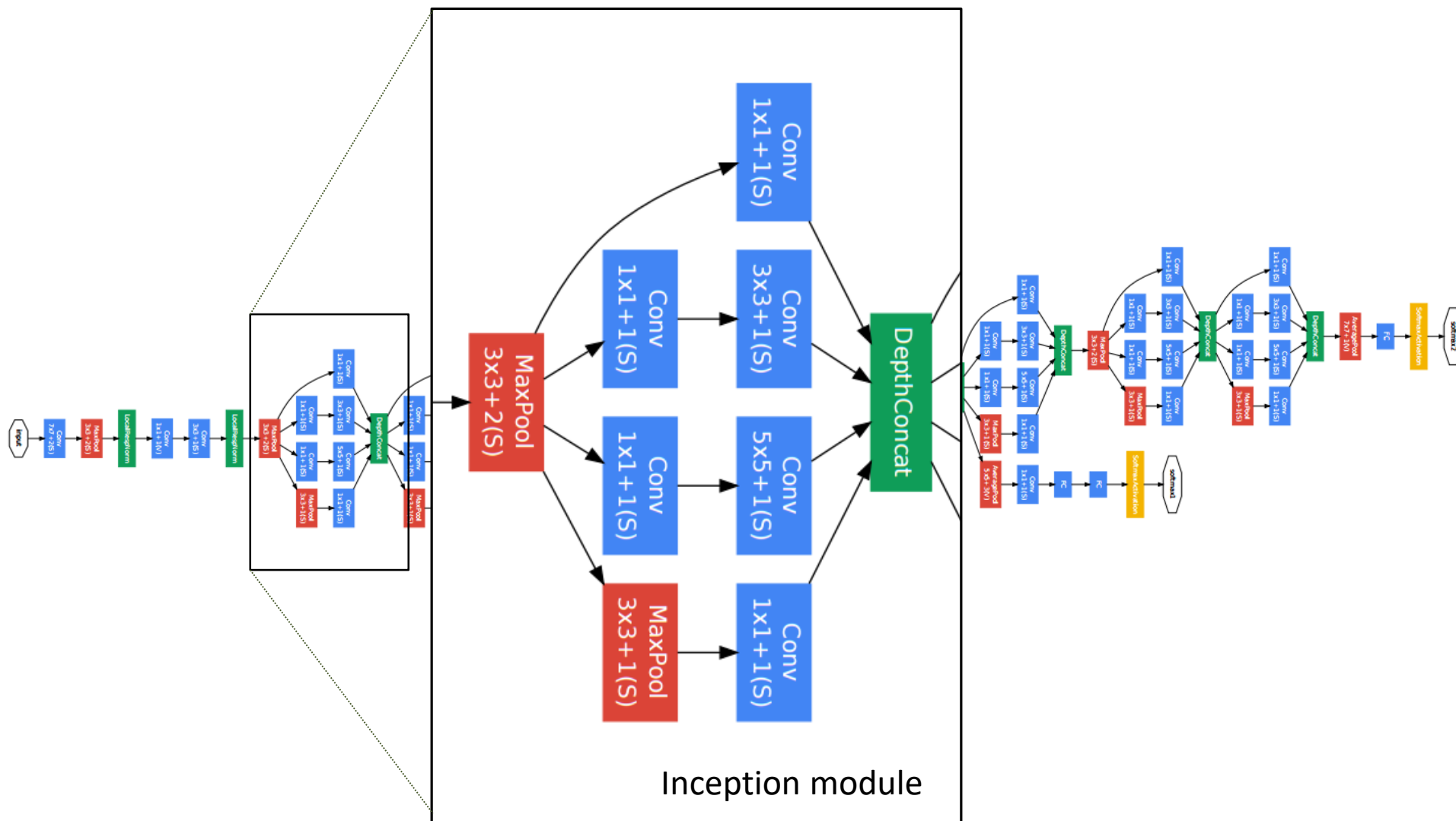
A short history of (old) CNNs

GoogLeNet (AKA Inception) [Szegedy et al. 2015]



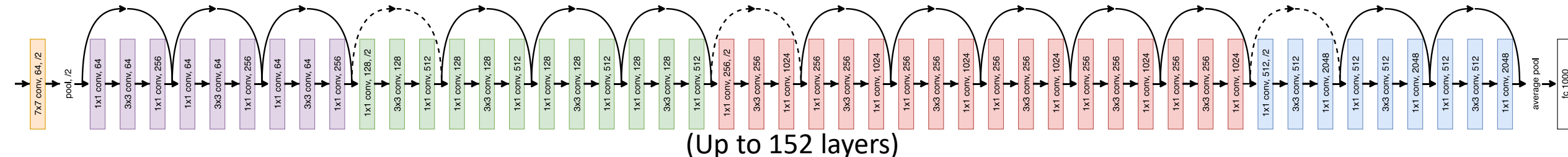
A short history of (old) CNNs

GoogLeNet (AKA Inception) [Szegedy et al. 2015]



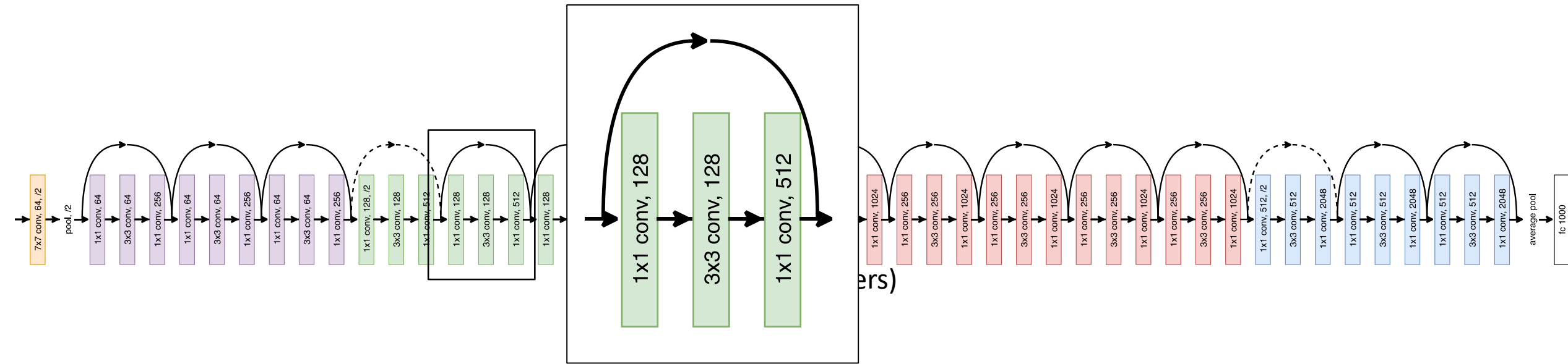
ResNet-50

ResNets [He, Zhang, Ren, & Sun, 2016]



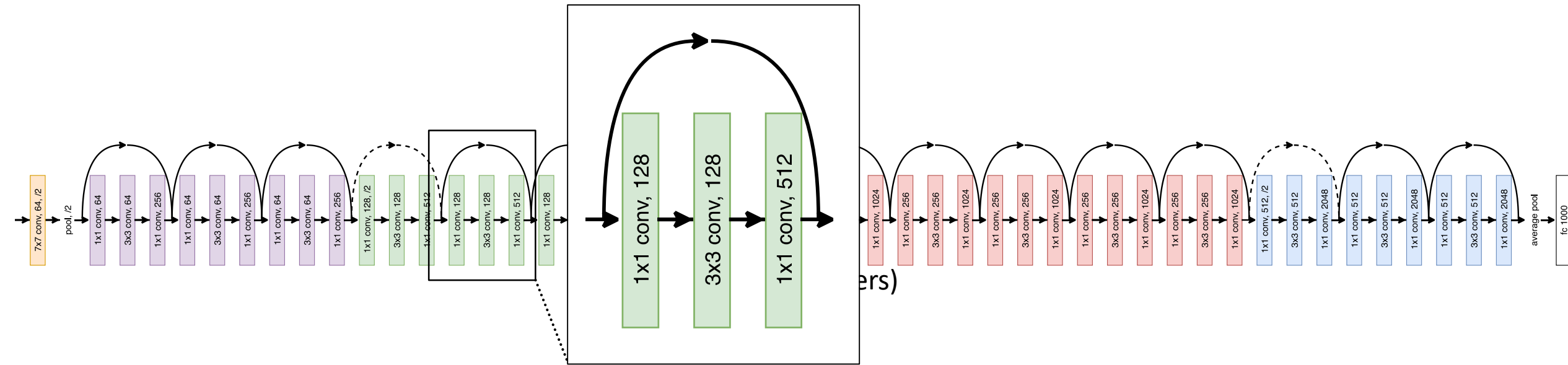
ResNet-50

ResNets [He, Zhang, Ren, & Sun, 2016]



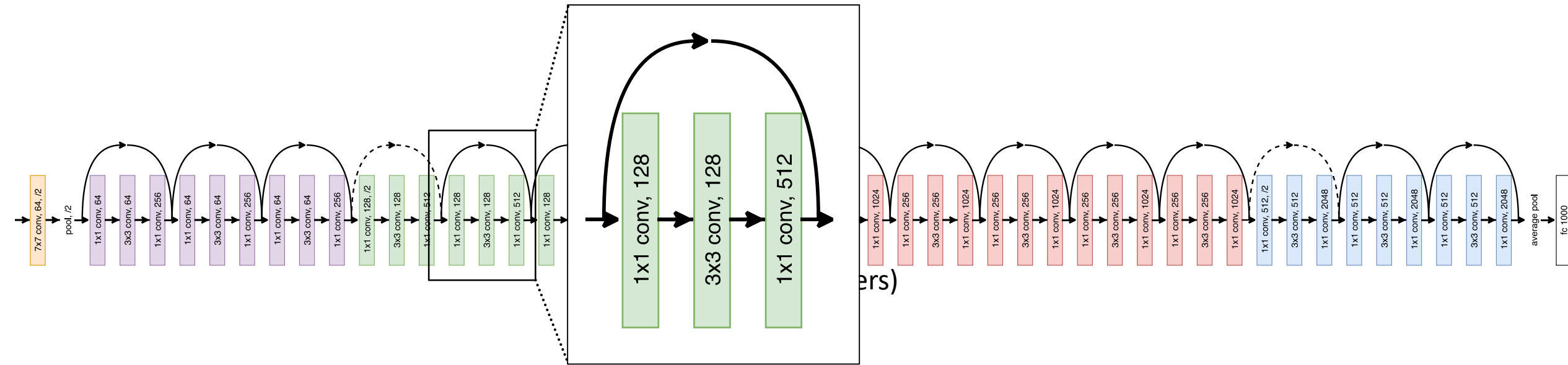
ResNet-50

ResNets [He, Zhang, Ren, & Sun, 2016]



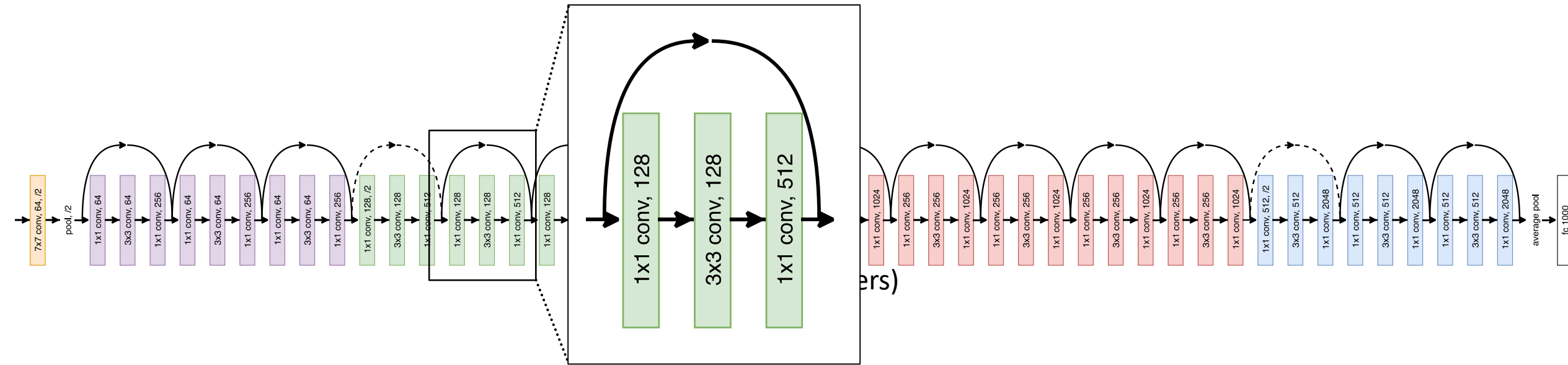
ResNet-50

ResNets [He, Zhang, Ren, & Sun, 2016]



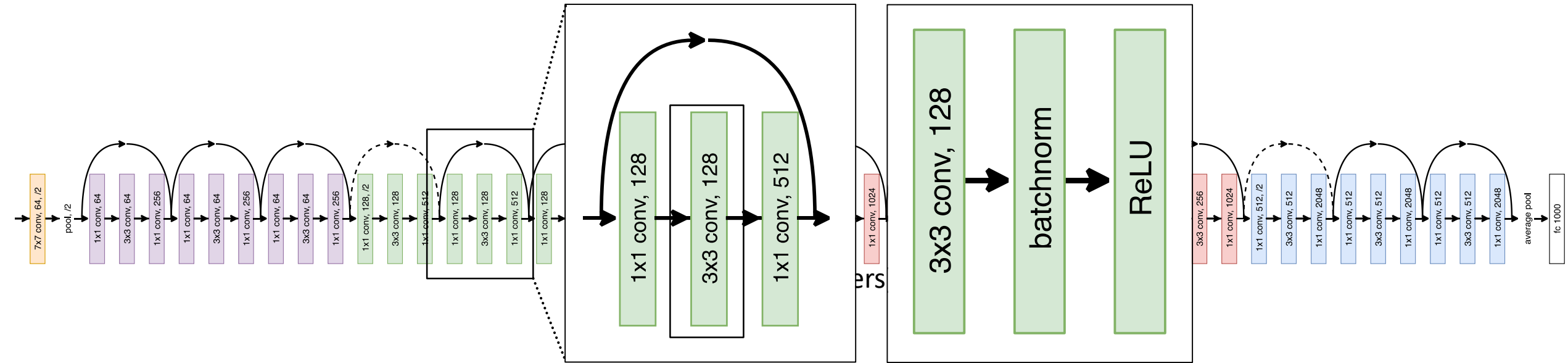
ResNet-50

ResNets [He, Zhang, Ren, & Sun, 2016]



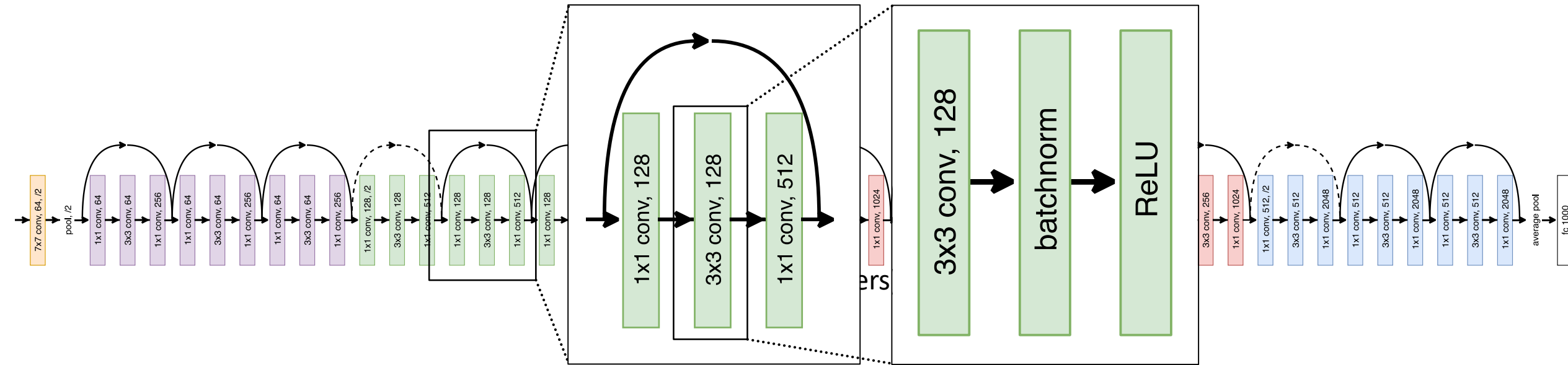
ResNet-50

ResNets [He, Zhang, Ren, & Sun, 2016]



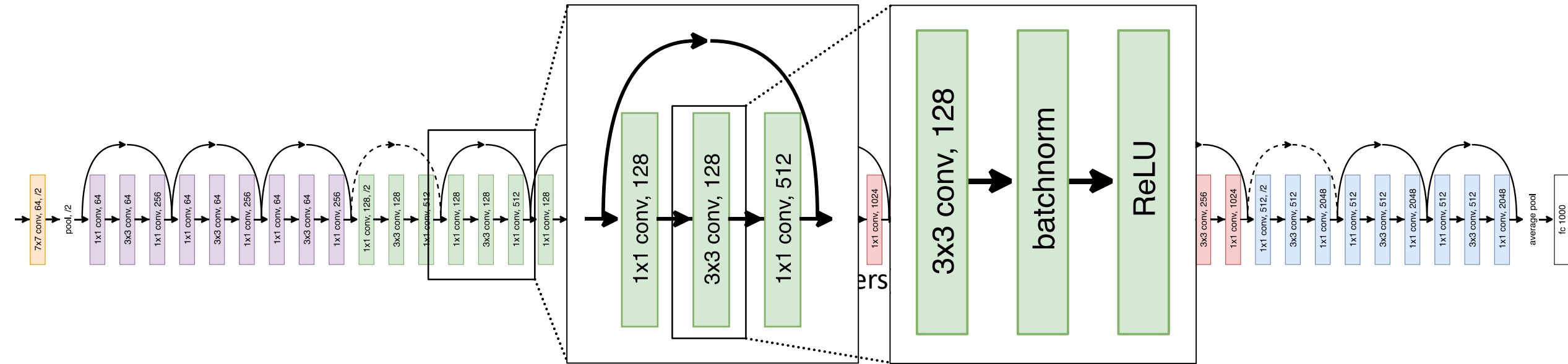
ResNet-50

ResNets [He, Zhang, Ren, & Sun, 2016]



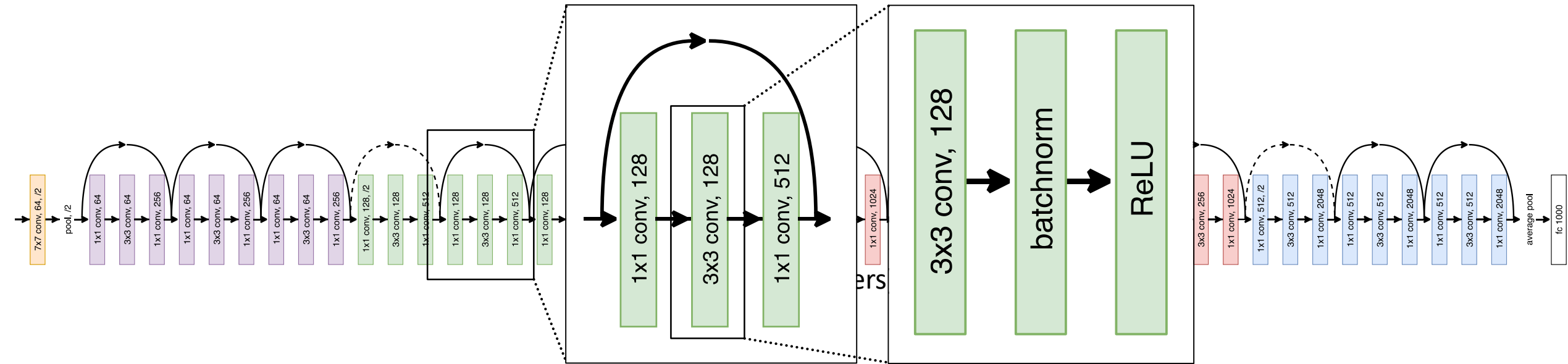
ResNet-50

ResNets [He, Zhang, Ren, & Sun, 2016]



ResNet-50

ResNets [He, Zhang, Ren, & Sun, 2016]

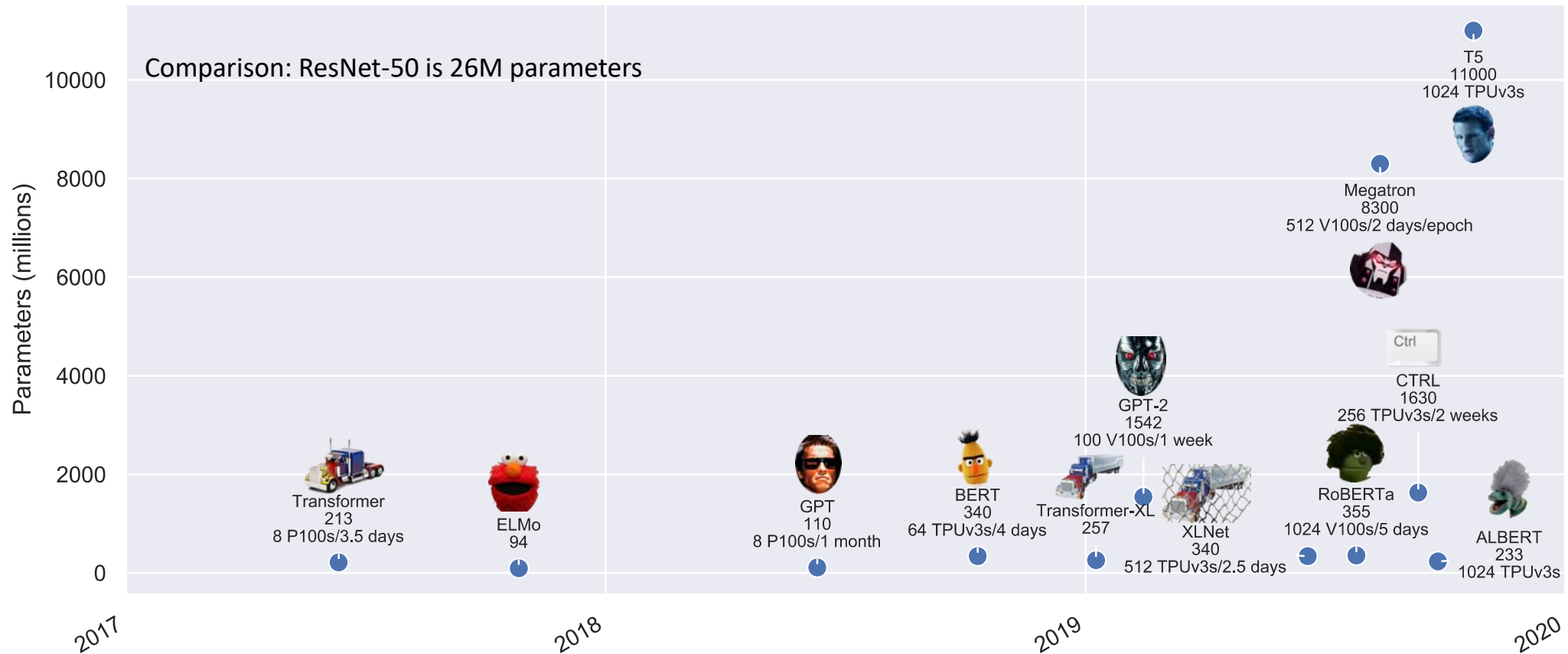


GPT-2 (transformers)

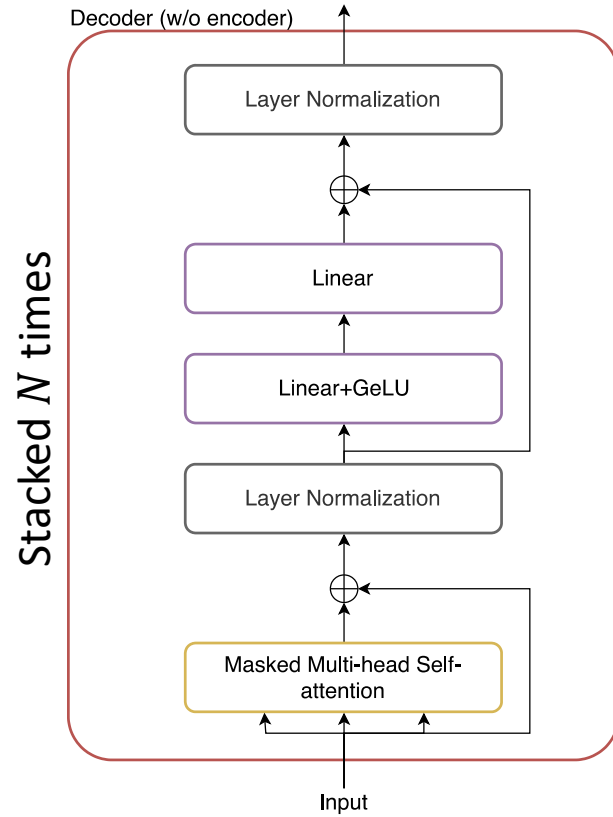
- Sequence-to-sequence models (like RNNs but with more parallelism)
- Revolutionizing NLP like AlexNet &co. did for computer vision

GPT-2 (transformers)

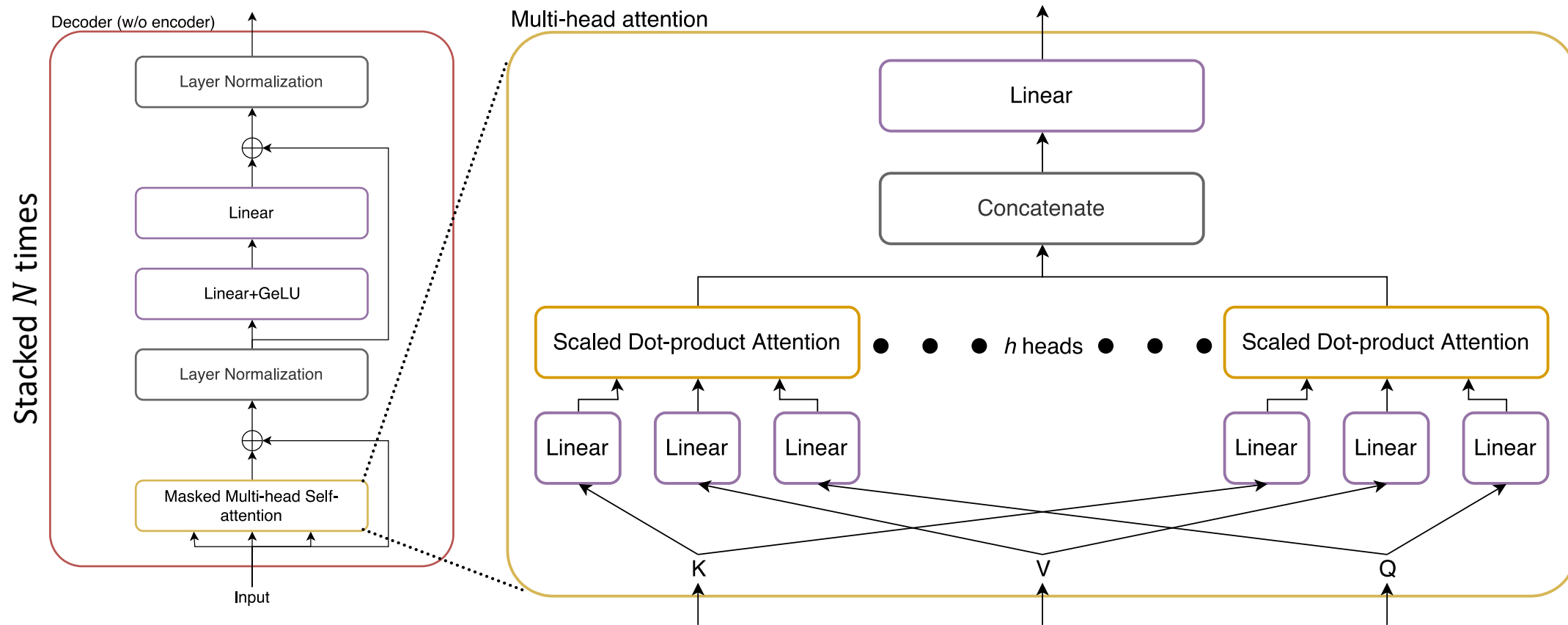
- Sequence-to-sequence models (like RNNs but with more parallelism)
- Revolutionizing NLP like AlexNet &co. did for computer vision



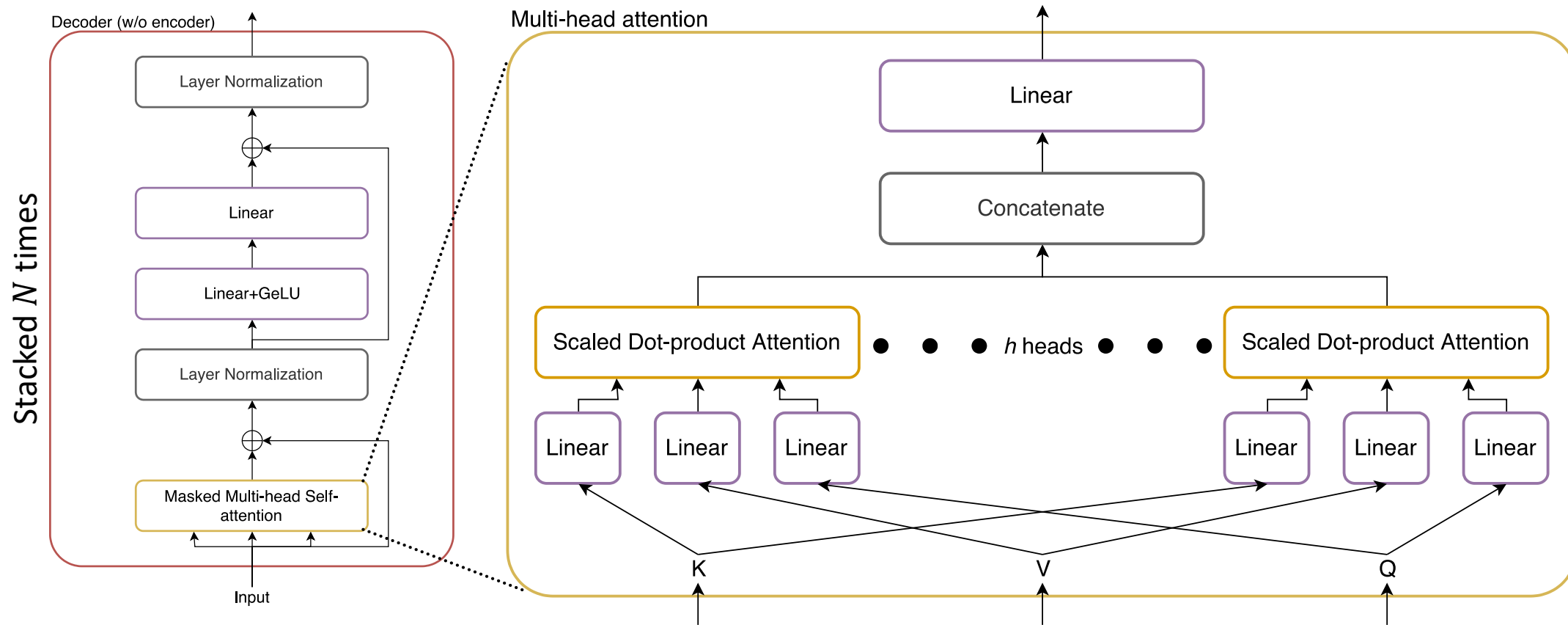
GPT-2 (transformers)



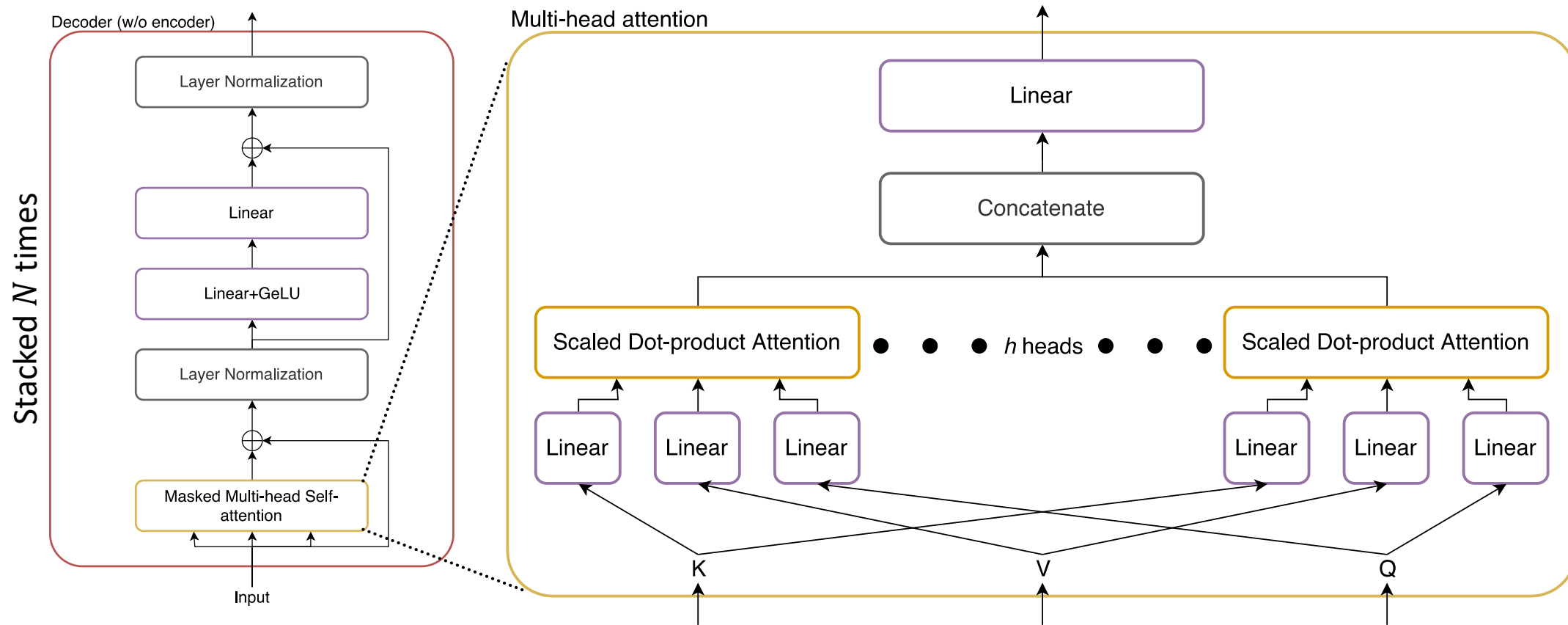
GPT-2 (transformers)



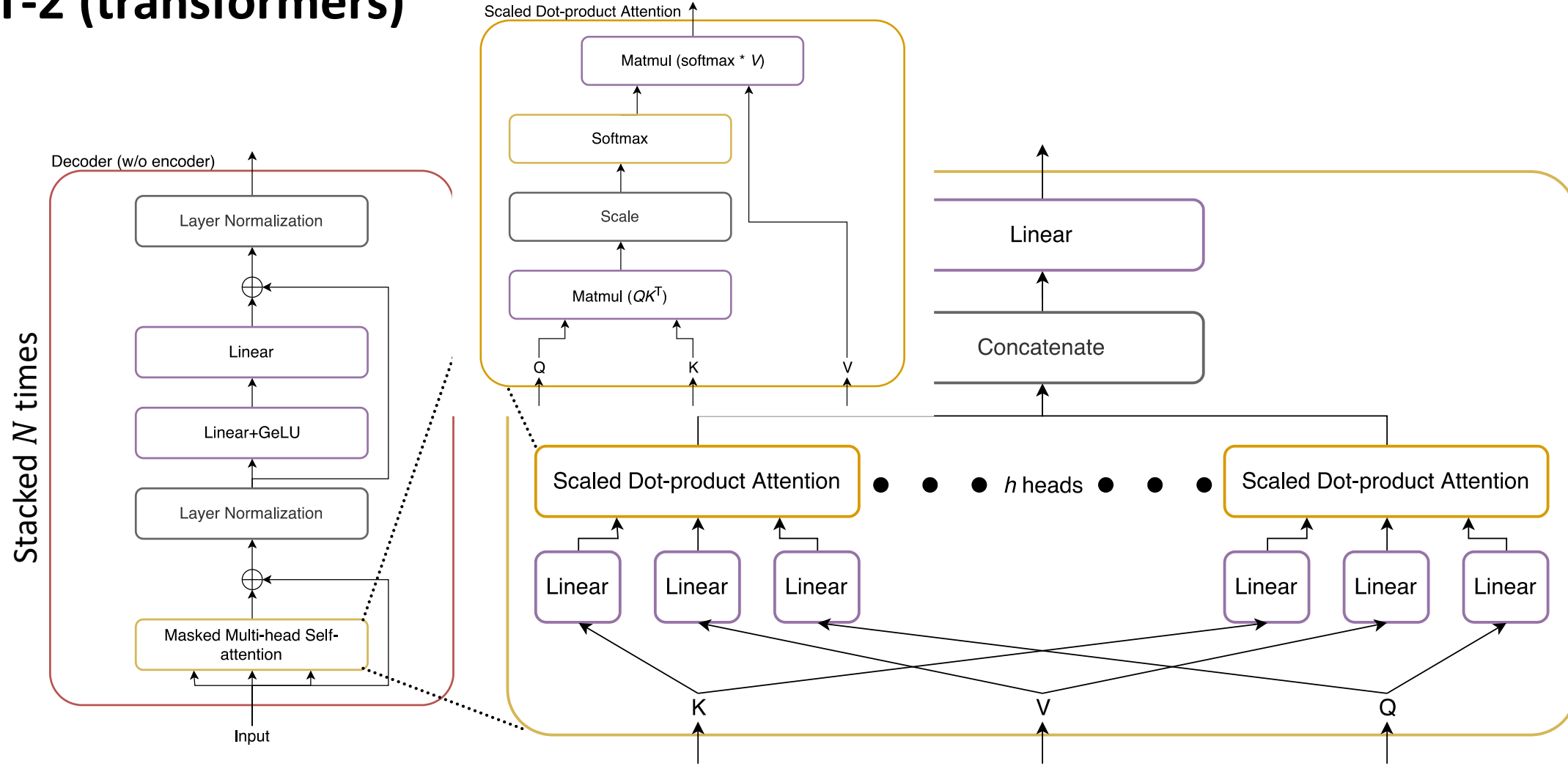
GPT-2 (transformers)



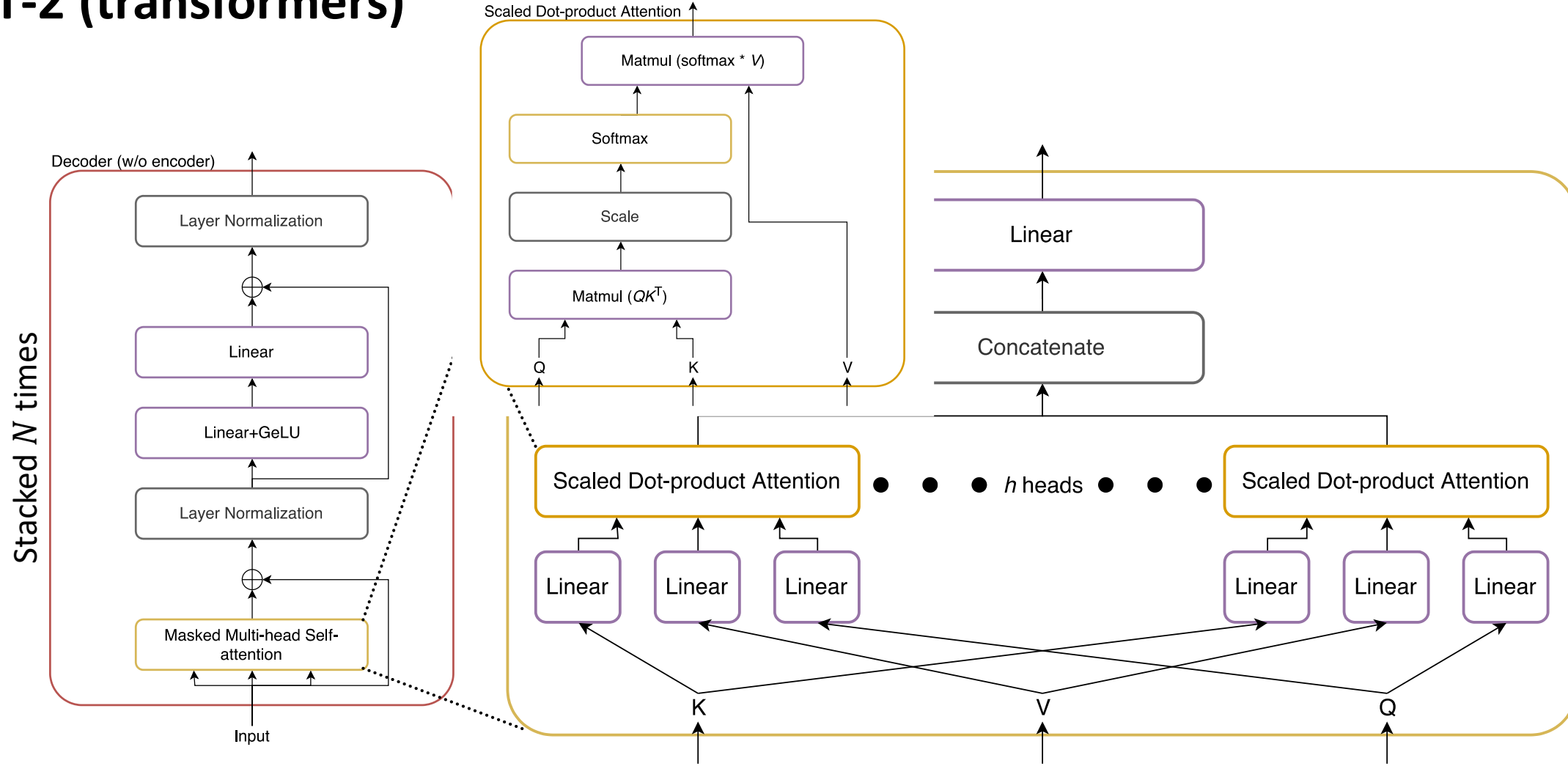
GPT-2 (transformers)



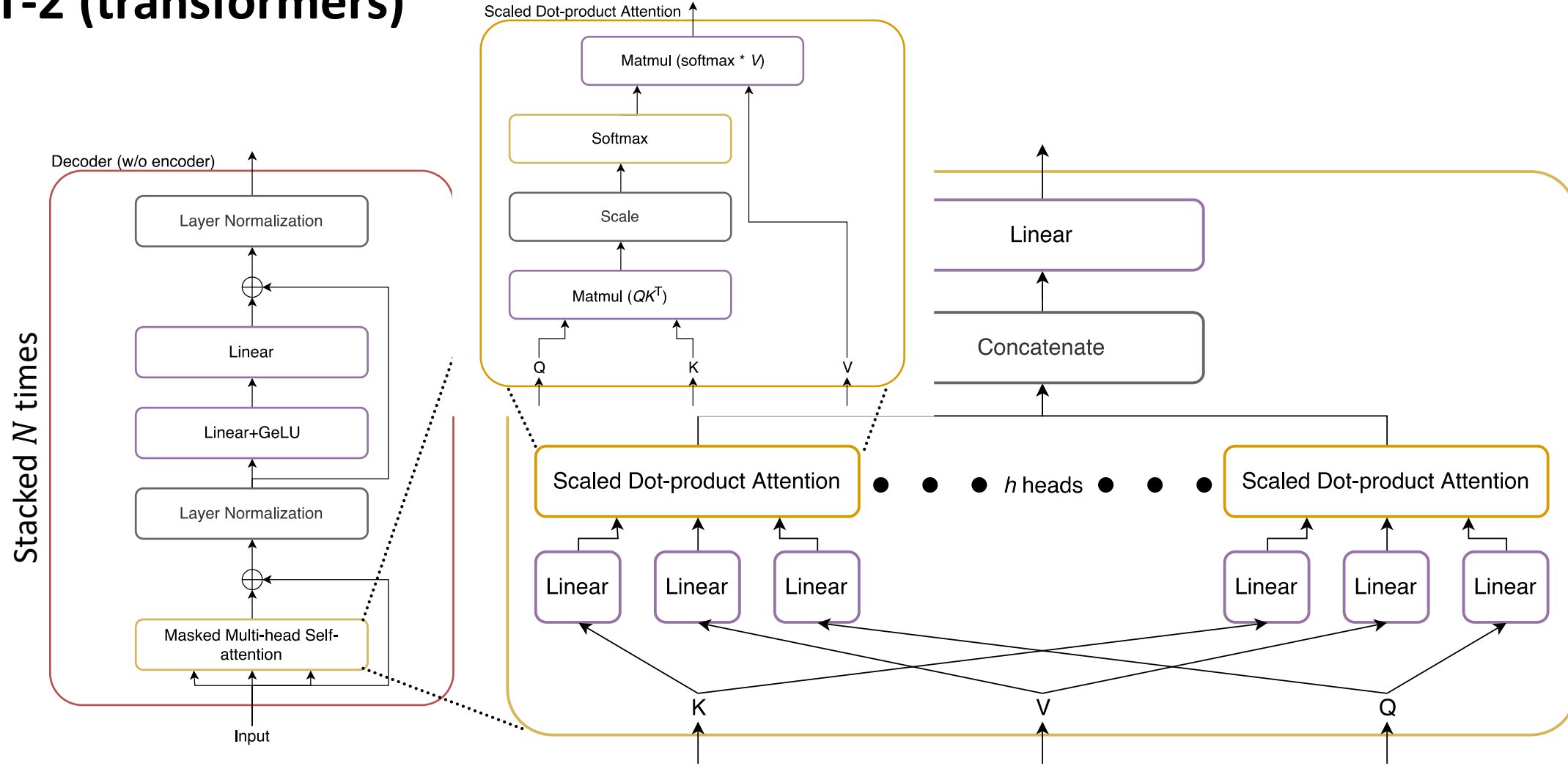
GPT-2 (transformers)



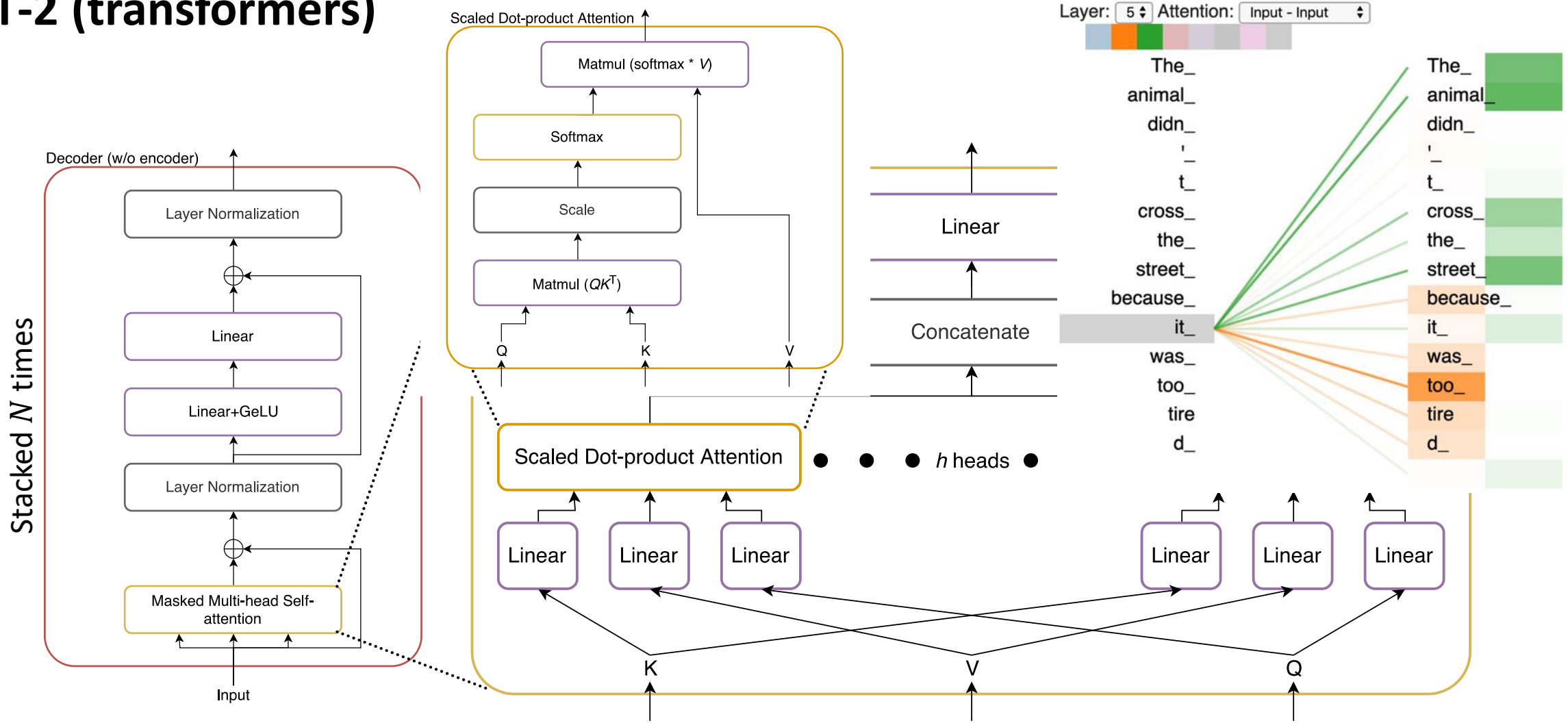
GPT-2 (transformers)



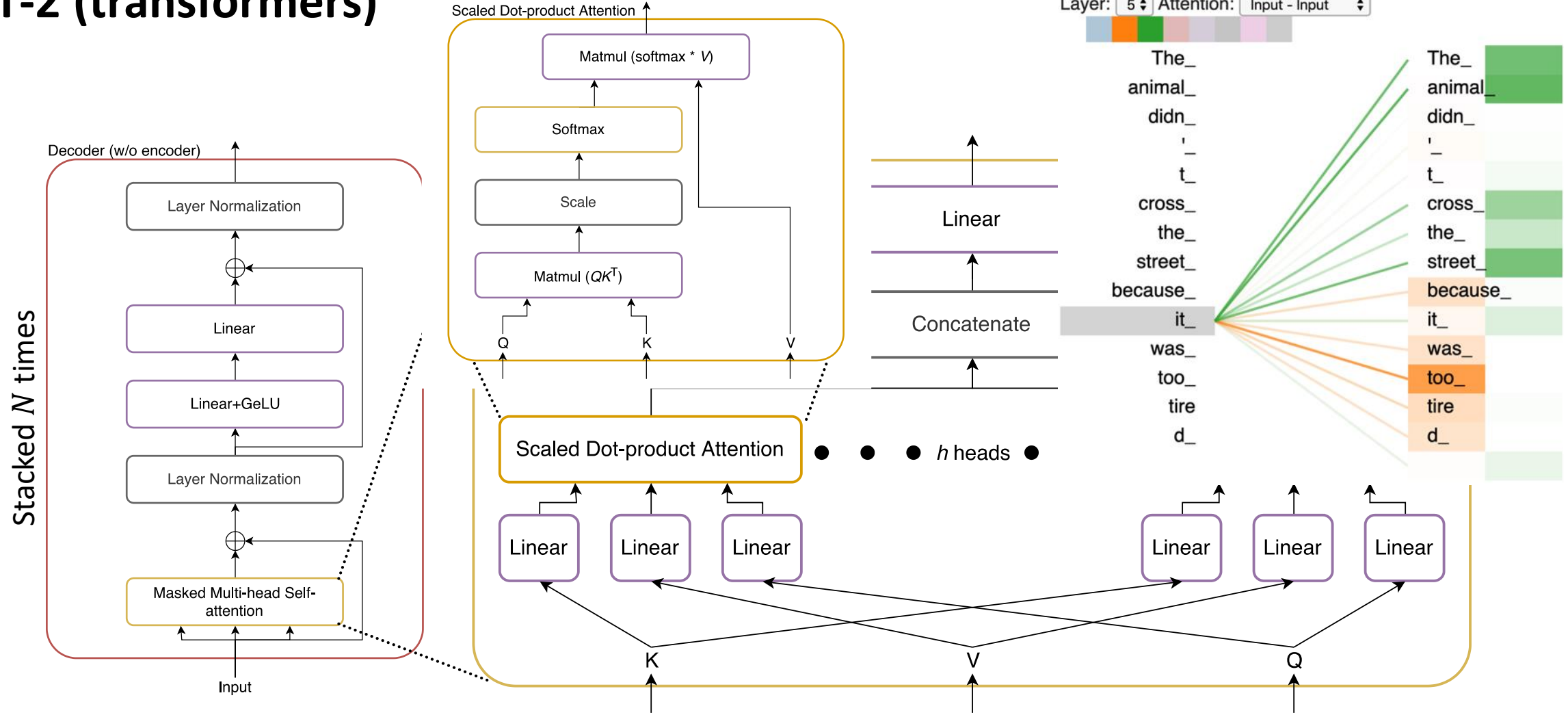
GPT-2 (transformers)



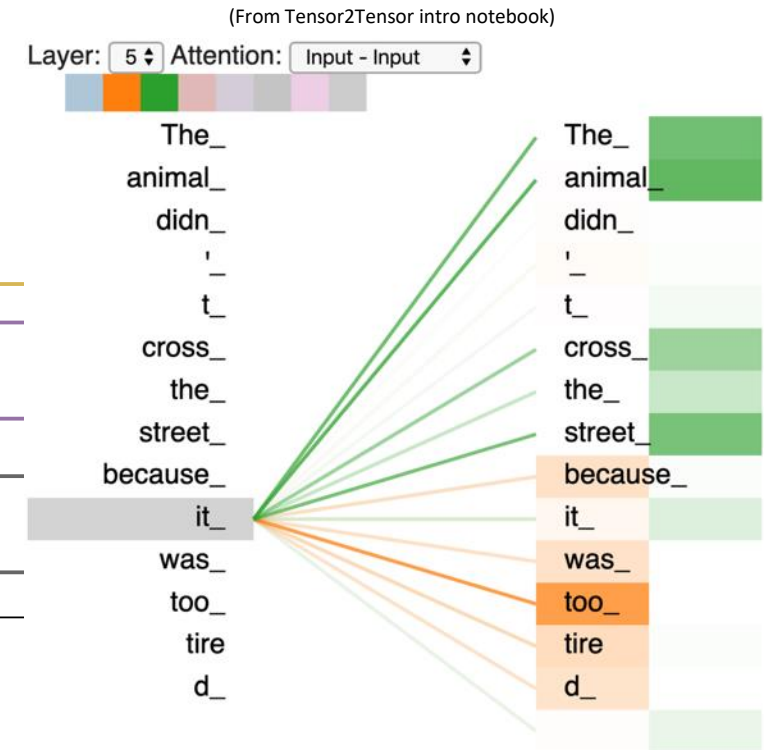
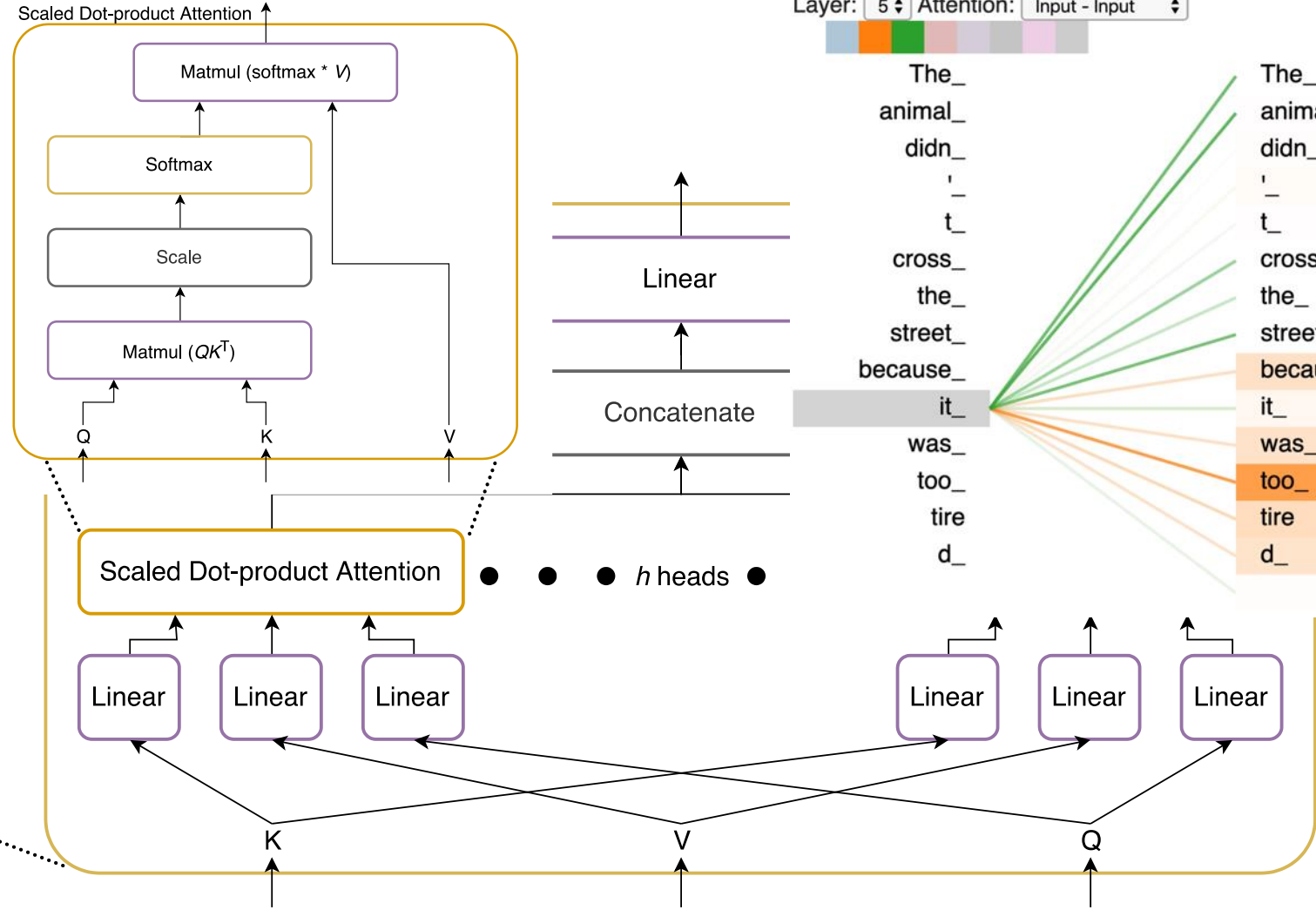
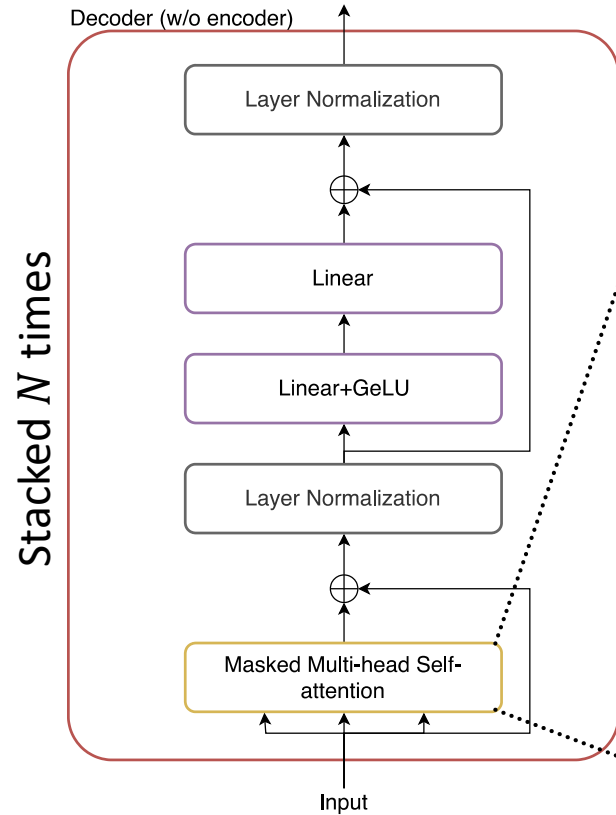
GPT-2 (transformers)



GPT-2 (transformers)



GPT-2 (transformers)



Objective: "predict the next word"

GPT-2 (transformers)

(From Tensor2Tensor intro notebook)

Layer: 5 Attention: Input - Input

Scaled Dot-product Attention ↑

Prompt (human-written):

Recycling is good for the world.

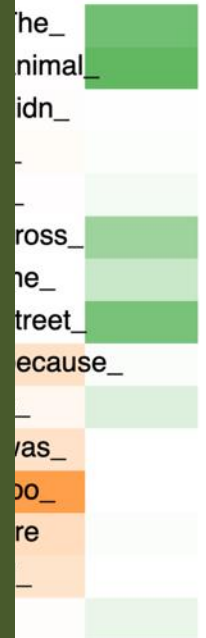
NO! YOU COULD NOT BE MORE WRONG!!

GPT-2-xlarge generated text:

Recycling is NOT good for the world. It is bad for the environment, it is bad for our health, and it is bad for our economy. I'm not kidding. Recycling is not good for the environment. It is destructive to the earth and it is a major contributor to global warming. Recycling is not good for our health. It contributes to obesity and diseases like heart disease and cancer. Recycling is bad for our economy. It increases the cost of a product, and in turn, the price of everything that is made with that product. Recycling is not good for our nation. We pay a tremendous price for the privilege of having the world's most advanced and efficient recycling system. Recycling is a huge, colossal waste of time, energy, money, and resources. And THAT is why we need to get back to basics and get back to basics in our recycling efforts. One of the best ways to start is to look at the process of creating a paper product. When you make a paper product, it is basically a long chain of materials. Everything from the raw materials (wood, cardboard, paper, etc.), to the reagents (dyes, solvents, etc.) to the printing equipment (chemicals, glue, paper, ink, etc.), to the packaging, to the packaging materials (mercury, chemicals, etc.) to the processing equipment (heating, cooling, etc.), to the packaging materials, to the packaging materials that are shipped overseas and to the packaging materials that are used in the United States. Each step along the way creates tons of waste that we constantly have to clean up. The process of making a paper product is a very wasteful one. But the end result is something that all of us need to consume. And if we want to keep the recycling process running efficiently, then we really need to think about each and every step that goes into making a paper product.

Stacked N times

Decoder (w/o encoder)



GPT-2 (transformers)

(From Tensor2Tensor intro notebook)

Layer: 5 Attention: Input - Input

Scaled Dot-product Attention ↑

Prompt (human-written):

Recycling is good for the world.

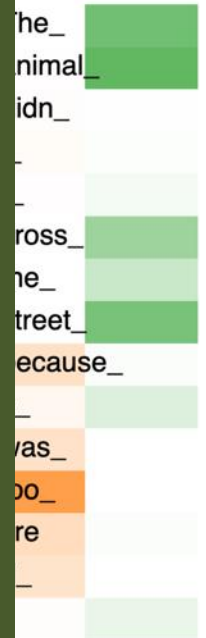
NO! YOU COULD NOT BE MORE WRONG!!

GPT-2-xlarge generated text:

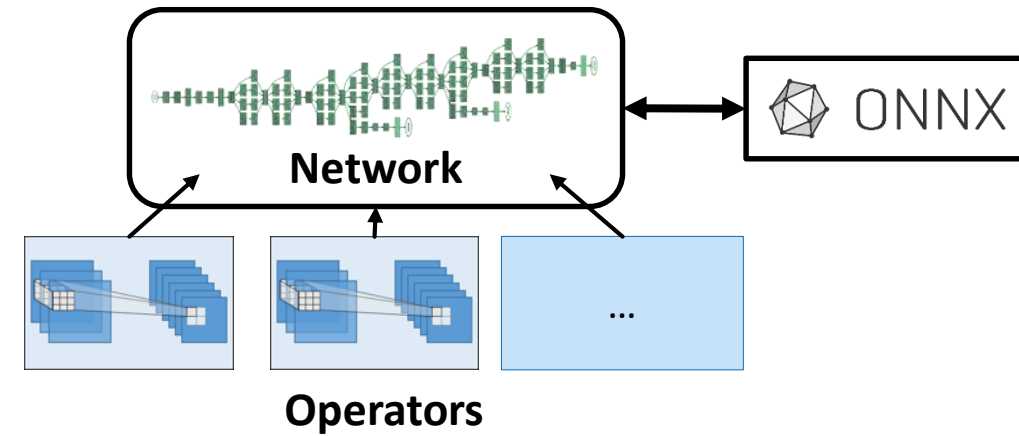
Recycling is NOT good for the world. It is bad for the environment, it is bad for our health, and it is bad for our economy. I'm not kidding. Recycling is not good for the environment. It is destructive to the earth and it is a major contributor to global warming. Recycling is not good for our health. It contributes to obesity and diseases like heart disease and cancer. Recycling is bad for our economy. It increases the cost of a product, and in turn, the price of everything that is made with that product. Recycling is not good for our nation. We pay a tremendous price for the privilege of having the world's most advanced and efficient recycling system. Recycling is a huge, colossal waste of time, energy, money, and resources. And THAT is why we need to get back to basics and get back to basics in our recycling efforts. One of the best ways to start is to look at the process of creating a paper product. When you make a paper product, it is basically a long chain of materials. Everything from the raw materials (wood, cardboard, paper, etc.), to the reagents (dyes, solvents, etc.) to the printing equipment (chemicals, glue, paper, ink, etc.), to the packaging, to the packaging materials (mercury, chemicals, etc.) to the processing equipment (heating, cooling, etc.), to the packaging materials, to the packaging materials that are shipped overseas and to the packaging materials that are used in the United States. Each step along the way creates tons of waste that we constantly have to clean up. The process of making a paper product is a very wasteful one. But the end result is something that all of us need to consume. And if we want to keep the recycling process running efficiently, then we really need to think about each and every step that goes into making a paper product.

Stacked N times

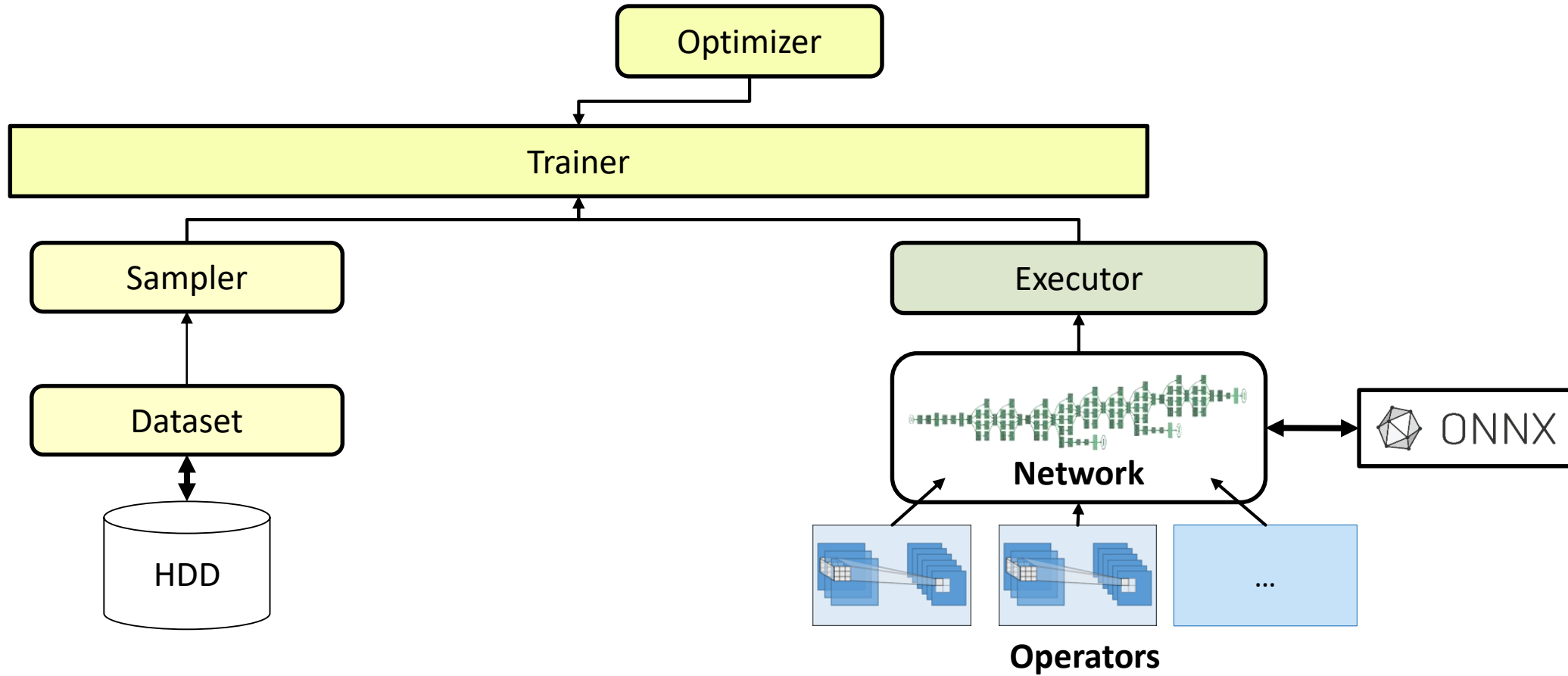
Decoder (w/o encoder)



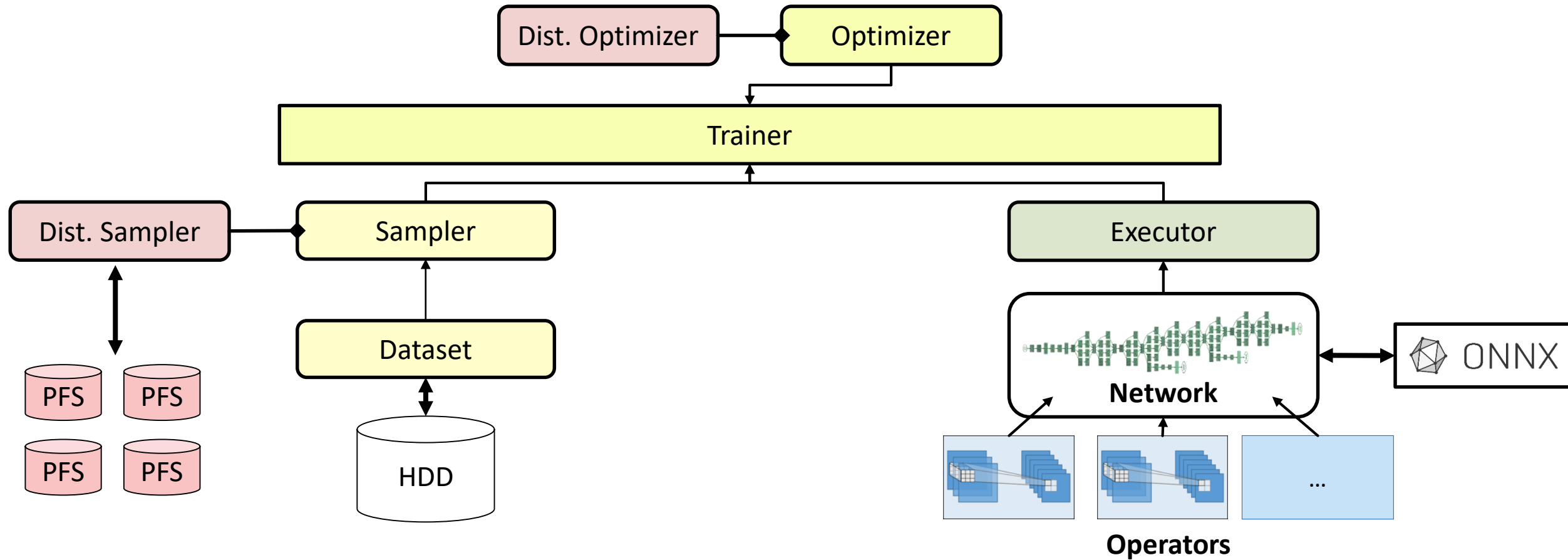
Networks



Training



Distributed training

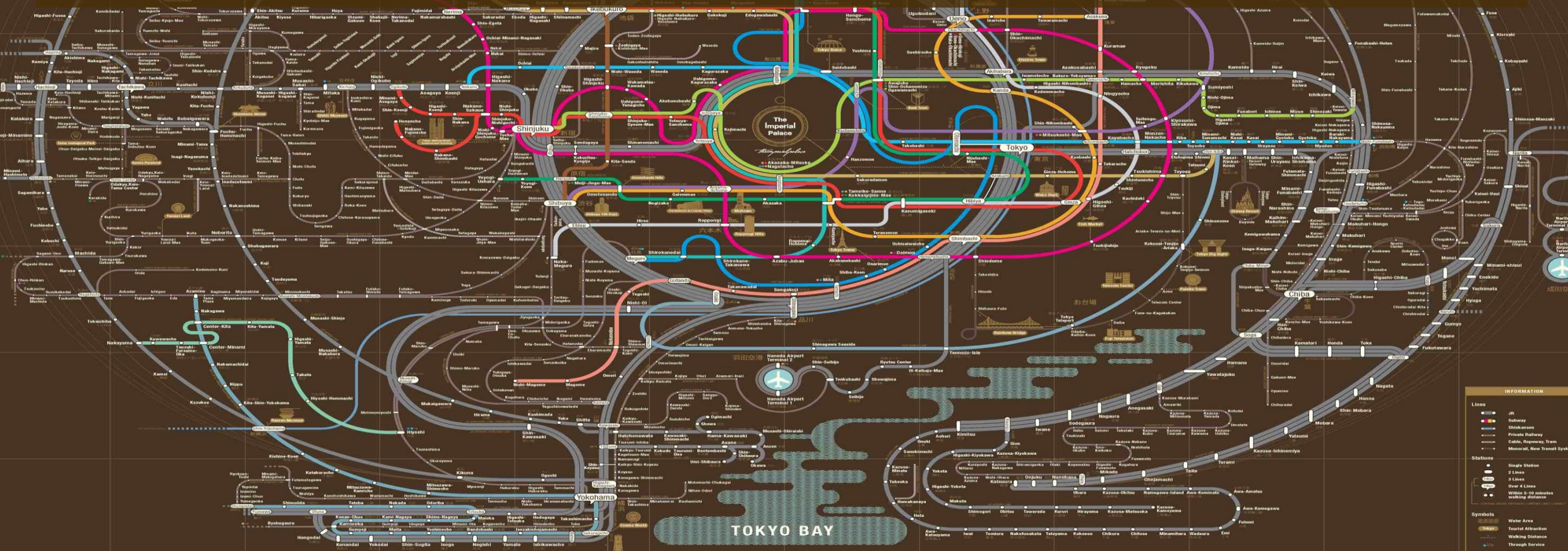


TOKYO METROPOLITAN RAILWAY SYSTEM

TOKYO RAILWAY SYSTEM MAY 2008
THIS IS NOT THE OFFICIAL MAP. THE RAILWAYS IN THIS DIAGRAM ARE
COPYRIGHTED BY TOKYO METRO. ALL RIGHTS RESERVED.
WWW.METRO.TOKYO



Optimizing parallel deep learning systems is a bit like navigating Tokyo by public transit
--- at first glance impossibly complex but eventually doable with the right guidelines ---
(Torsten Hoefler)



INFORMATION

- Lines
 - JR
 - Subway
 - Shinkansen
 - Private Railway
 - Cable, Ropeway, Tram
 - Monorail, New Transit System
- Stations
 - Single Station
 - 2 Lines
 - 3 Lines
 - Over 4 Lines
 - Walking to the station
 - Through Service
- Symbols
 - Water Area
 - Tourist Attraction
 - Walking Distance
 - Through Service

TOKYO BAY

10/21/2008

Operator implementations: fully-connected layers

$$Y = \sigma(WX + b)$$

Operator implementations: fully-connected layers

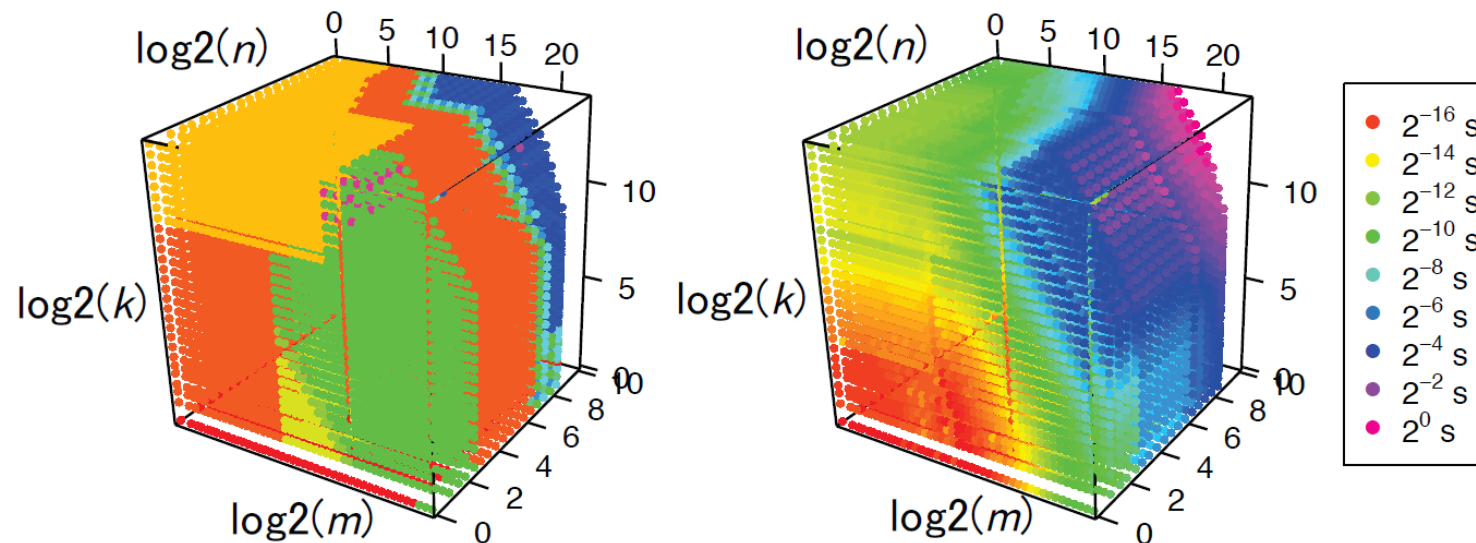
$$Y = \sigma(WX + b)$$

- **Dominated by matrix-matrix multiplication**
 - Standard tricks: vectorize, tile, fusion, ...
- **BLAS3 GEMM**
 - cuBLAS, MKL, ...

Operator implementations: fully-connected layers

$$Y = \sigma(WX + b)$$

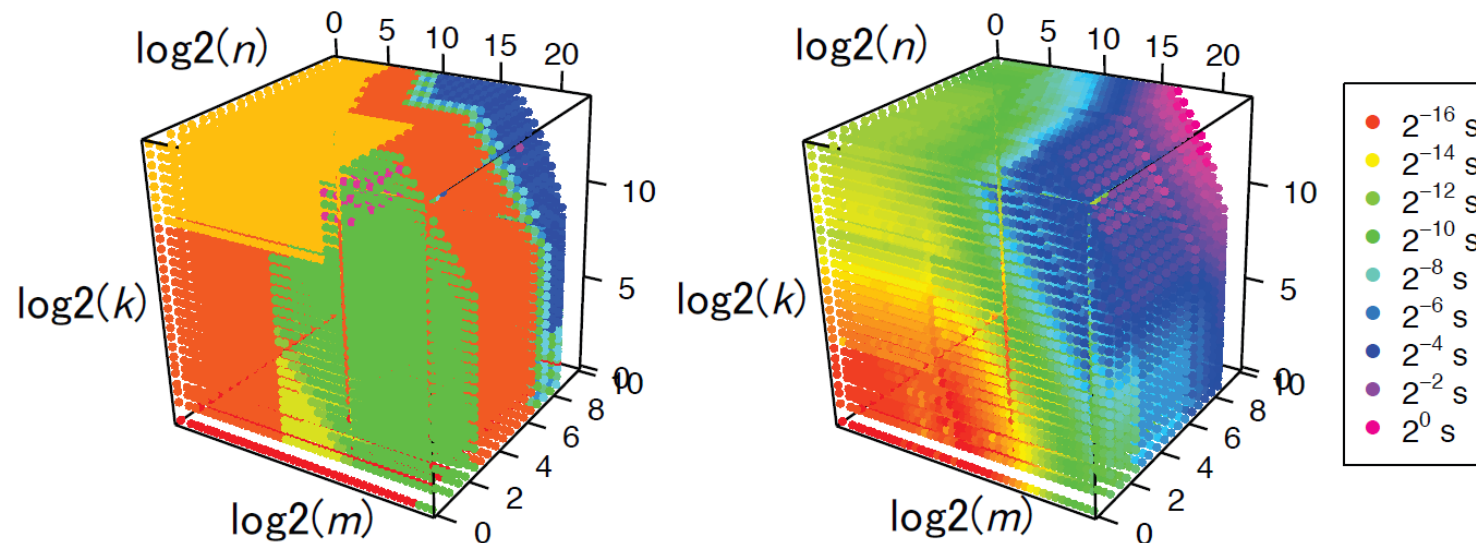
- **Dominated by matrix-matrix multiplication**
 - Standard tricks: vectorize, tile, fusion, ...
- **BLAS3 GEMM**
 - cuBLAS, MKL, ...



Operator implementations: fully-connected layers

$$Y = \sigma(WX + b)$$

- **Dominated by matrix-matrix multiplication**
 - Standard tricks: vectorize, tile, fusion, ...
- **BLAS3 GEMM**
 - cuBLAS, MKL, ...

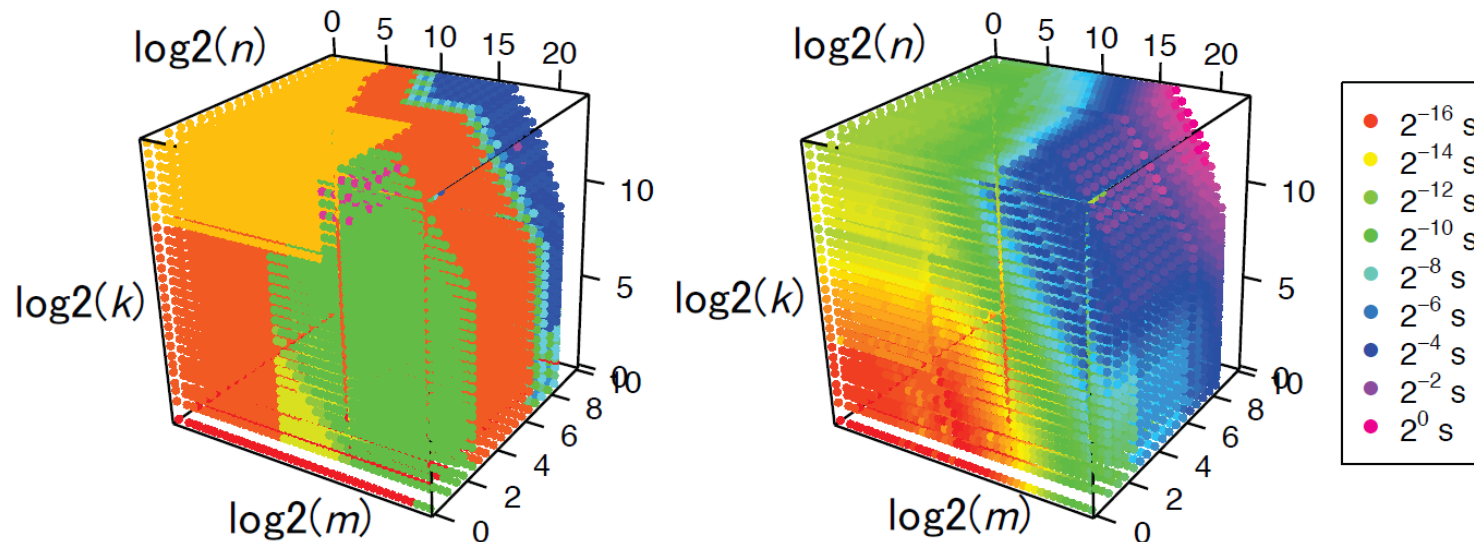


Operator implementations: fully-connected layers

$$Y = \sigma(WX + b)$$

- **Dominated by matrix-matrix multiplication**
 - Standard tricks: vectorize, tile, fusion, ...
- **BLAS3 GEMM**
 - cuBLAS, MKL, ...

Performance is not consistent!



Operator implementations: convolution

$$Y_{k,f,i,j} = \sum_{c=0}^{C-1} \sum_{a=-0}^0 \sum_{b=-0}^0 X_{k,c,i+a,j+b} W_{f,c,a+0,b+0}$$

Operator implementations: convolution

$$Y_{k,f,i,j} = \sum_{c=0}^{C-1} \sum_{a=-0}^0 \sum_{b=-0}^0 X_{k,c,i+a,j+b} W_{f,c,a+0,b+0}$$

Direct

Indirect

Operator implementations: convolution

$$Y_{k,f,i,j} = \sum_{c=0}^{C-1} \sum_{a=-O}^O \sum_{b=-O}^O X_{k,c,i+a,j+b} W_{f,c,a+O,b+O}$$

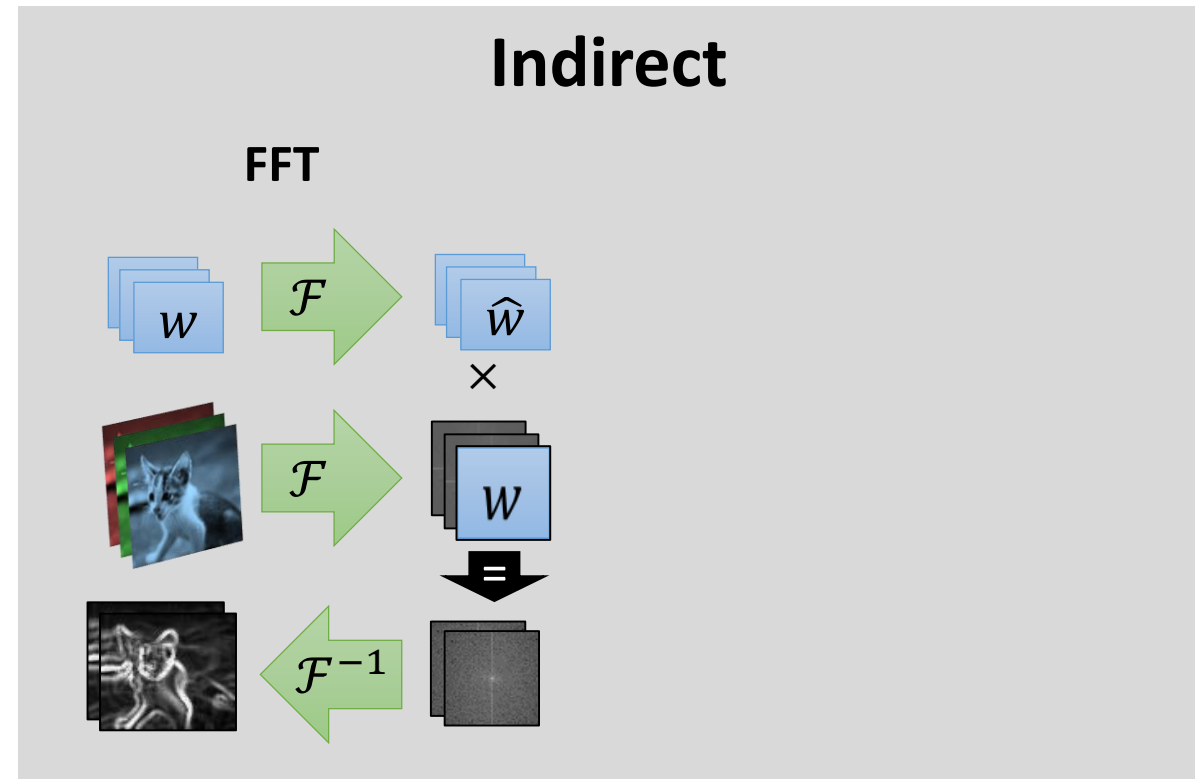
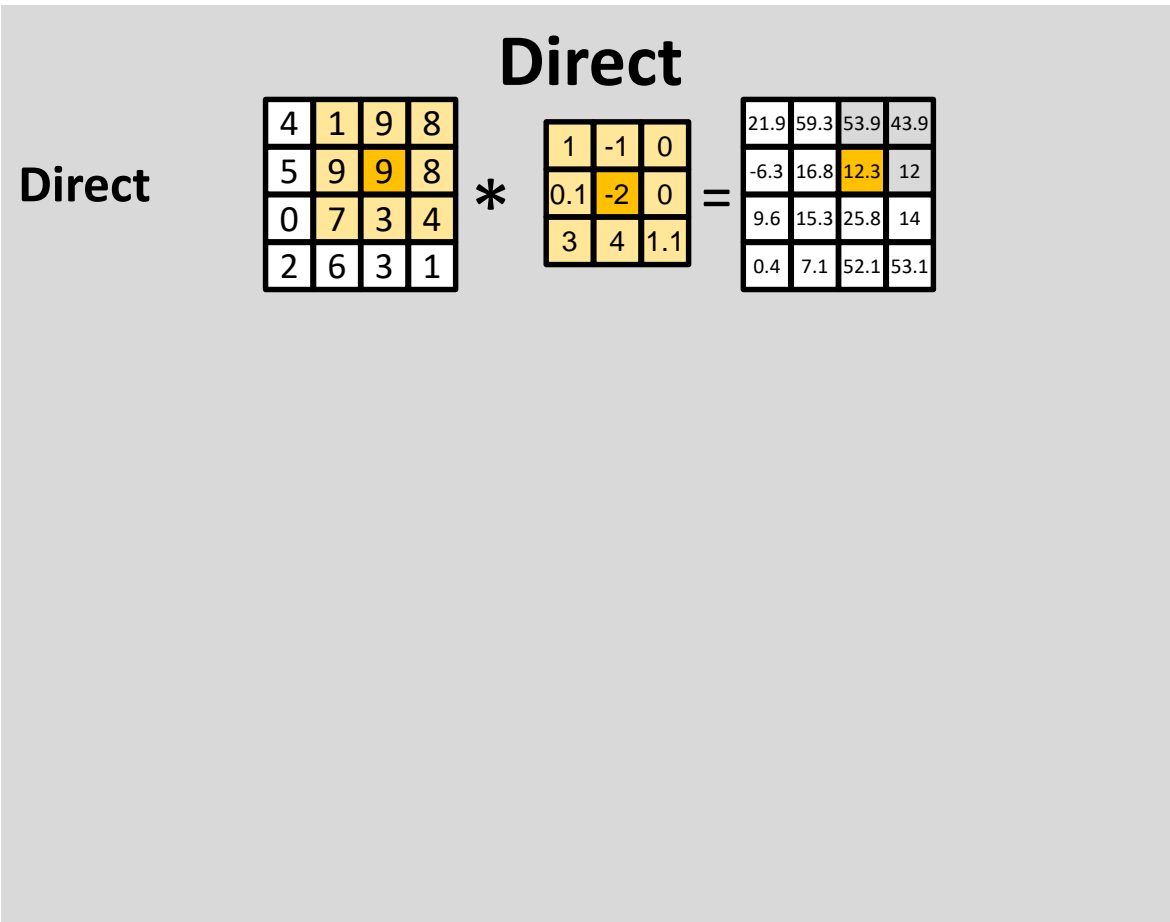
Direct

4	1	9	8	*	1	-1	0	=	21.9	59.3	53.9	43.9
5	9	9	8		0.1	-2	0		-6.3	16.8	12.3	12
0	7	3	4		3	4	1.1		9.6	15.3	25.8	14
2	6	3	1						0.4	7.1	52.1	53.1

Indirect

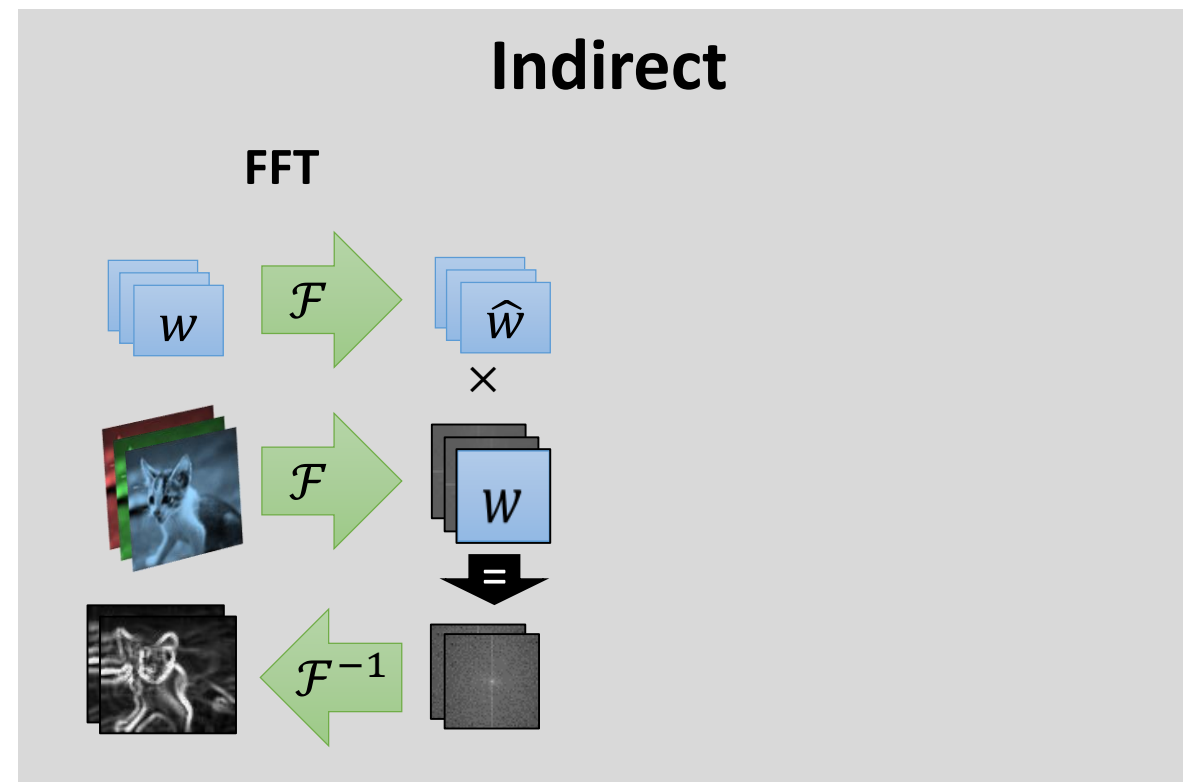
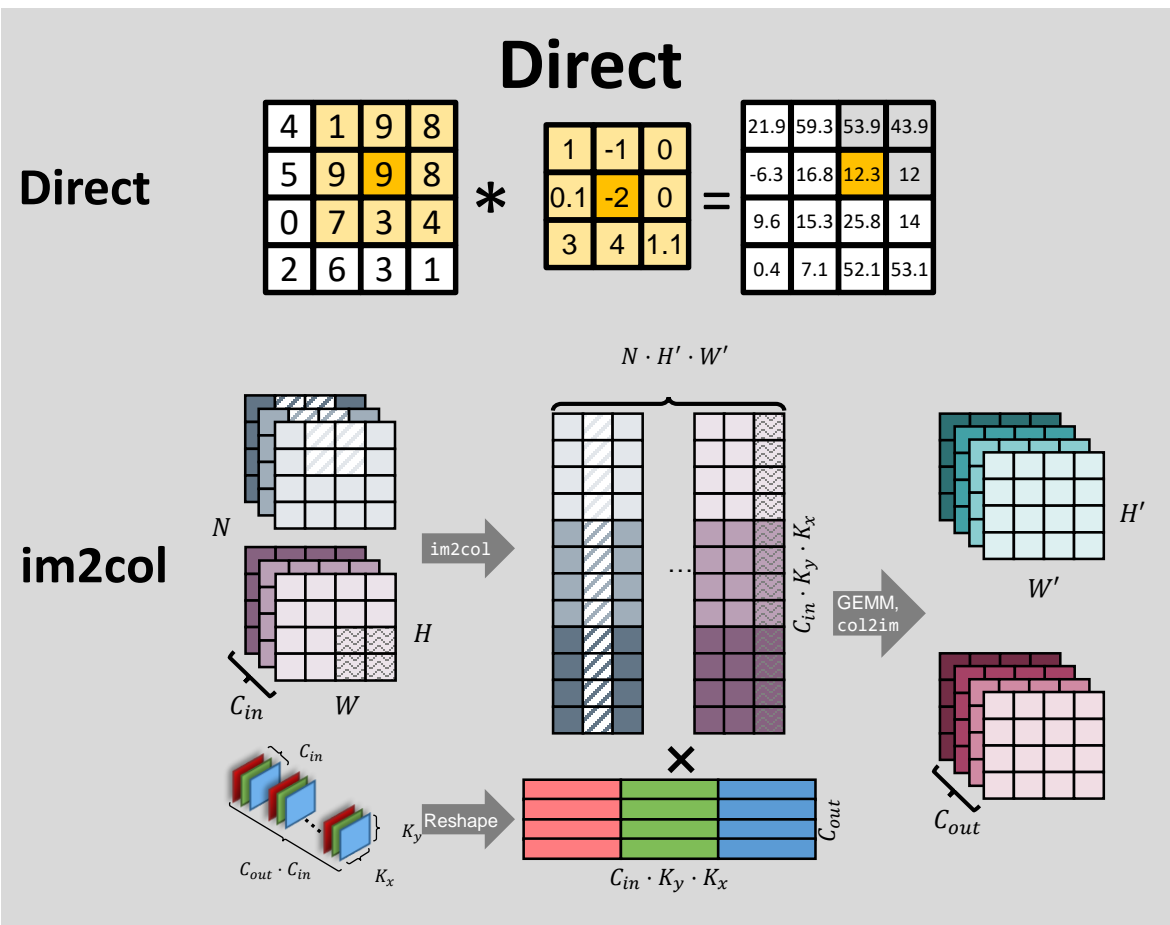
Operator implementations: convolution

$$Y_{k,f,i,j} = \sum_{c=0}^{C-1} \sum_{a=-O}^O \sum_{b=-O}^O X_{k,c,i+a,j+b} W_{f,c,a+O,b+O}$$



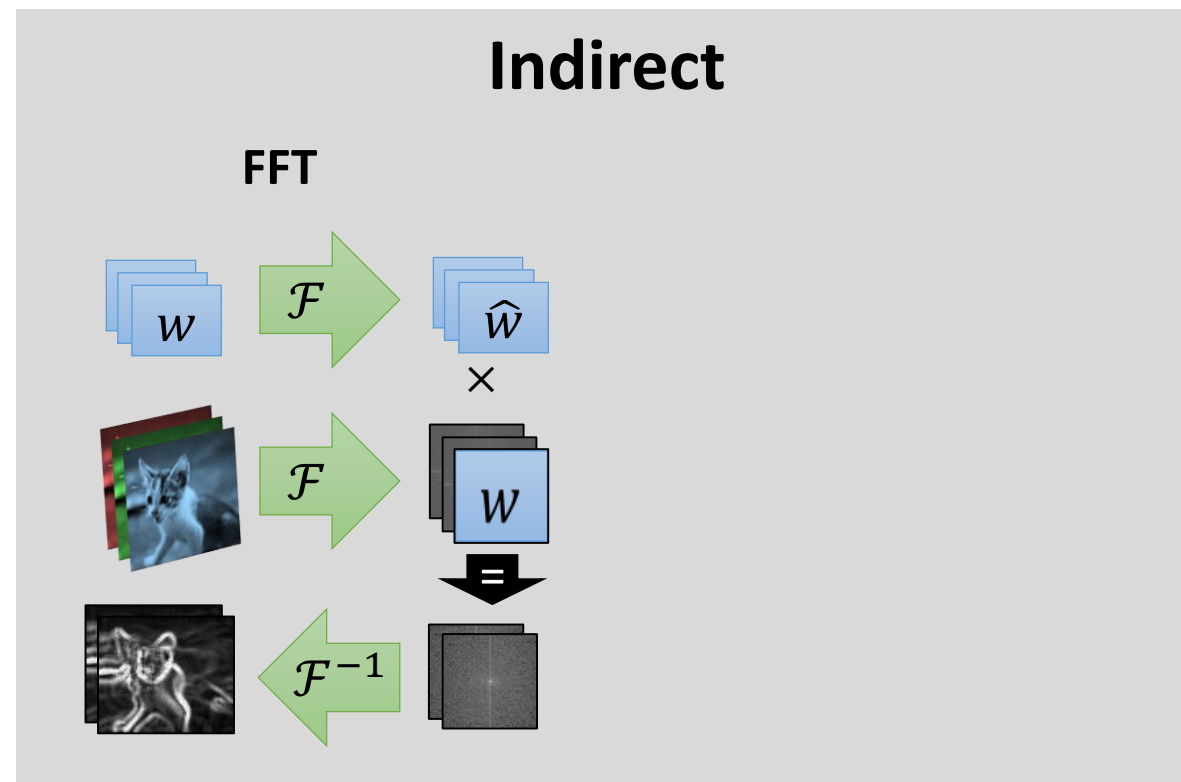
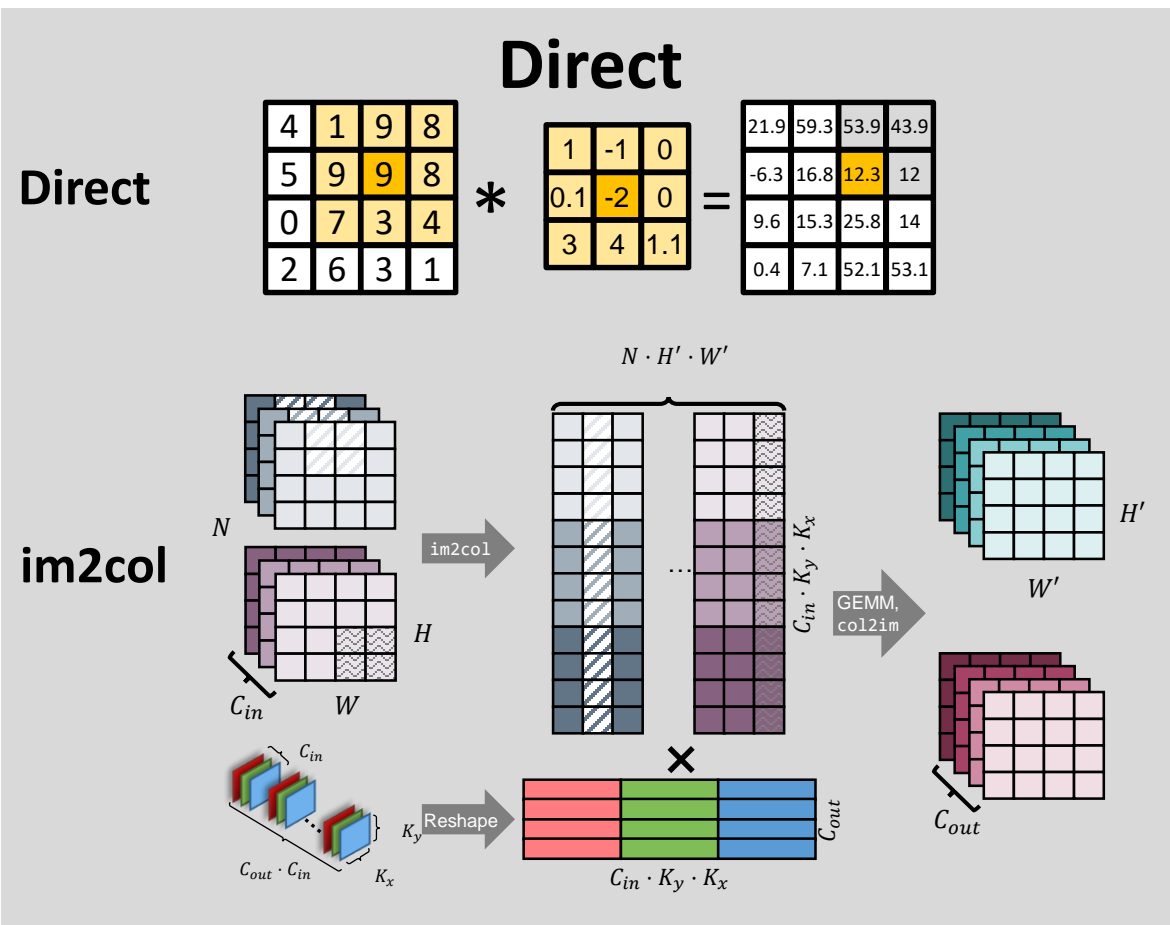
Operator implementations: convolution

$$Y_{k,f,i,j} = \sum_{c=0}^{C-1} \sum_{a=-O}^O \sum_{b=-O}^O X_{k,c,i+a,j+b} W_{f,c,a+O,b+O}$$



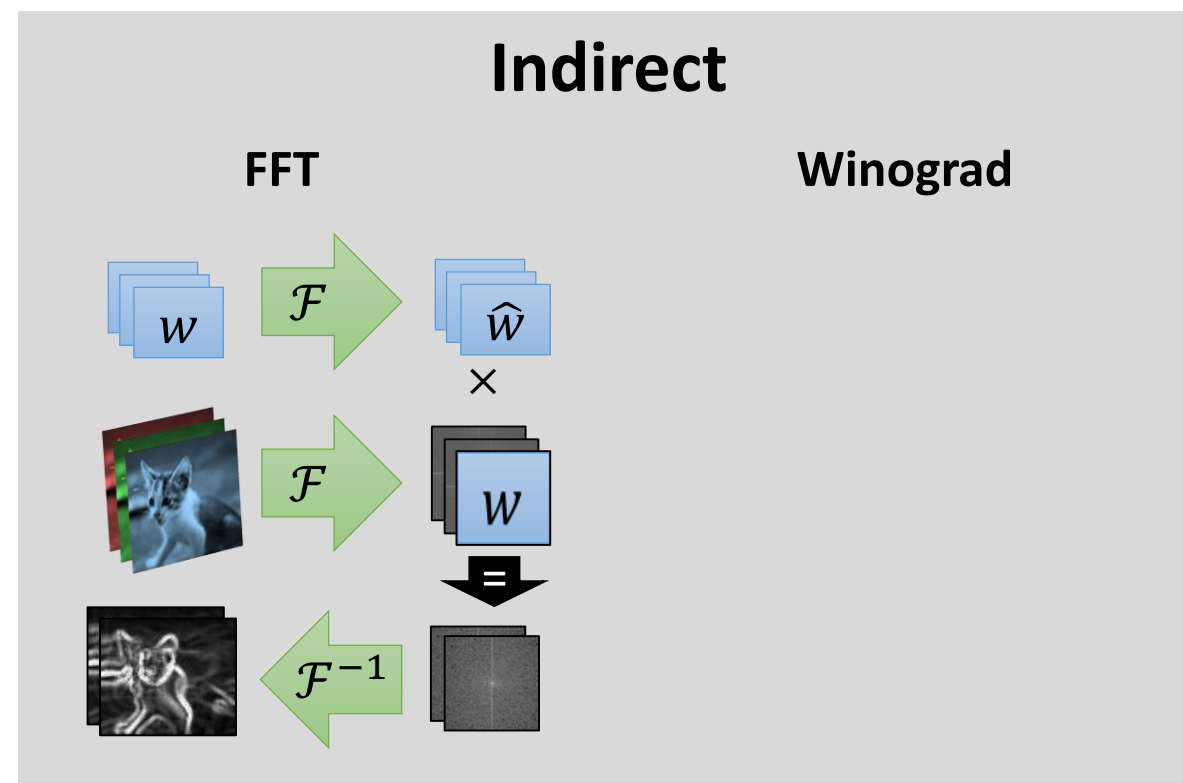
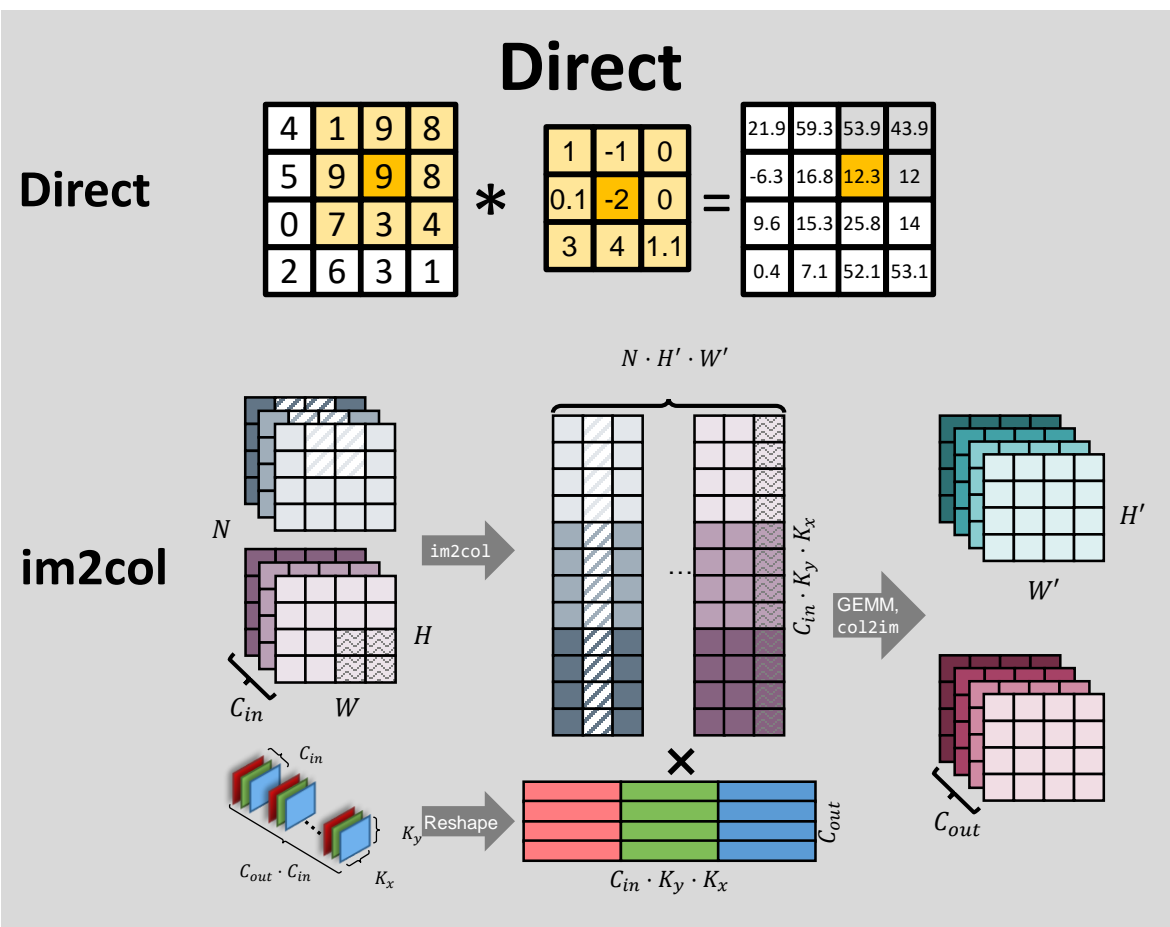
Operator implementations: convolution

$$Y_{k,f,i,j} = \sum_{c=0}^{C-1} \sum_{a=-O}^O \sum_{b=-O}^O X_{k,c,i+a,j+b} W_{f,c,a+O,b+O}$$



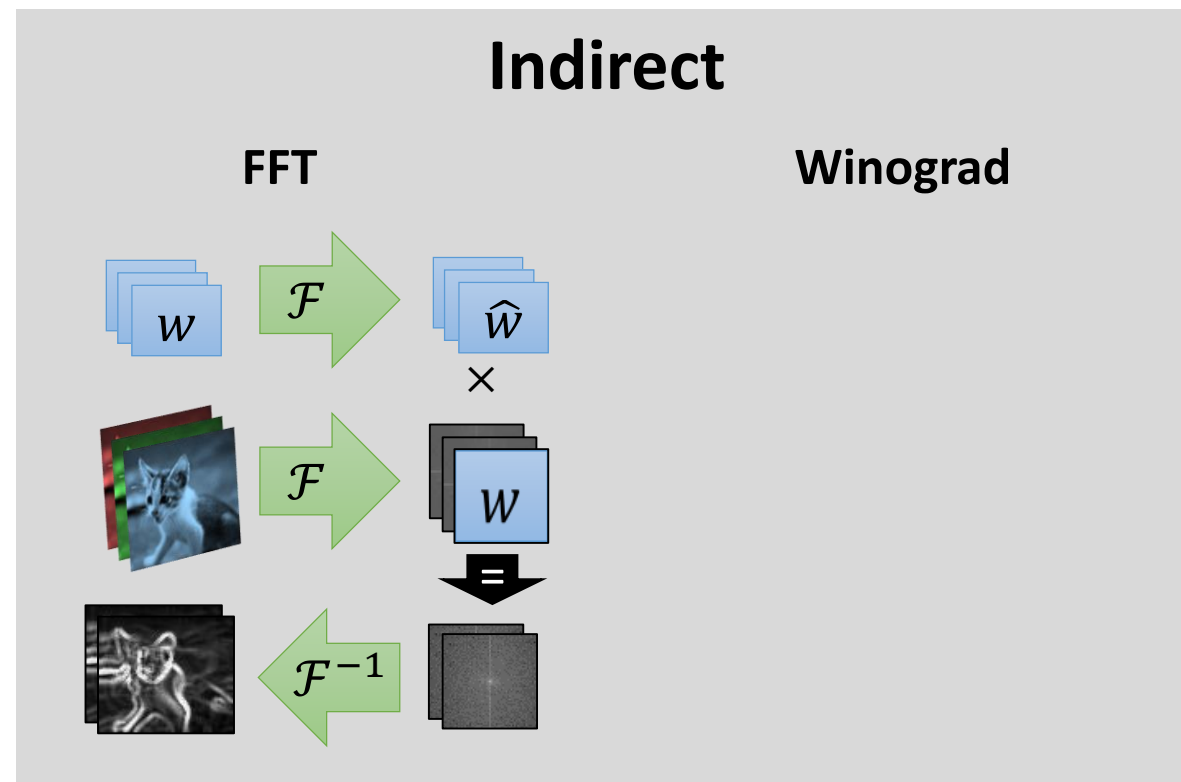
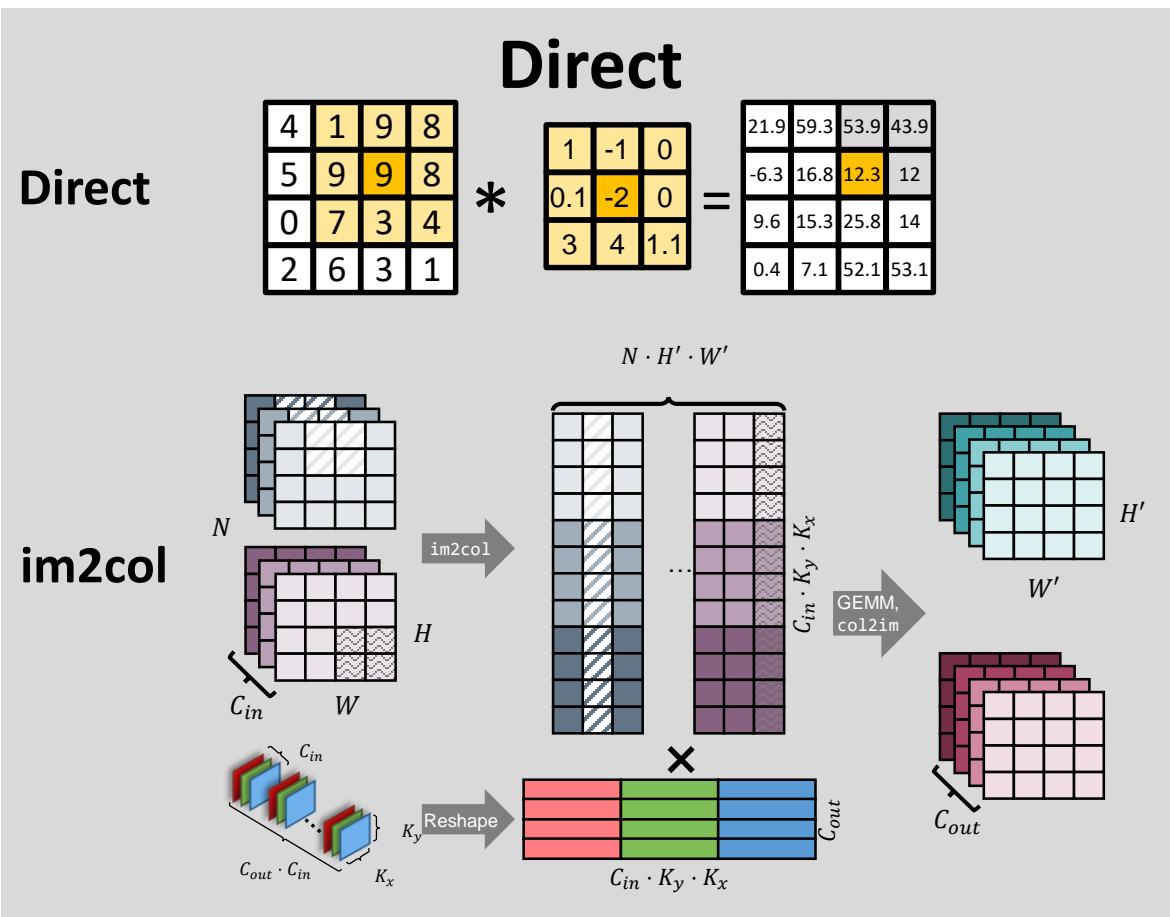
Operator implementations: convolution

$$Y_{k,f,i,j} = \sum_{c=0}^{C-1} \sum_{a=-O}^O \sum_{b=-O}^O X_{k,c,i+a,j+b} W_{f,c,a+O,b+O}$$



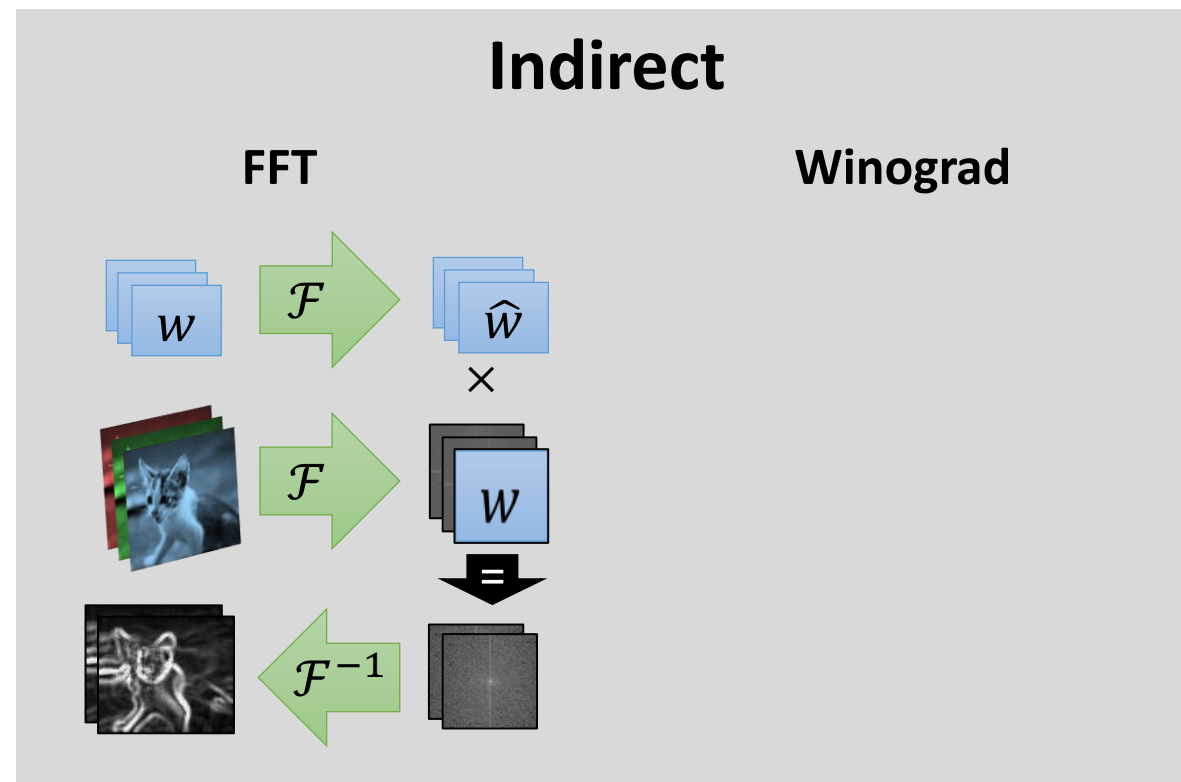
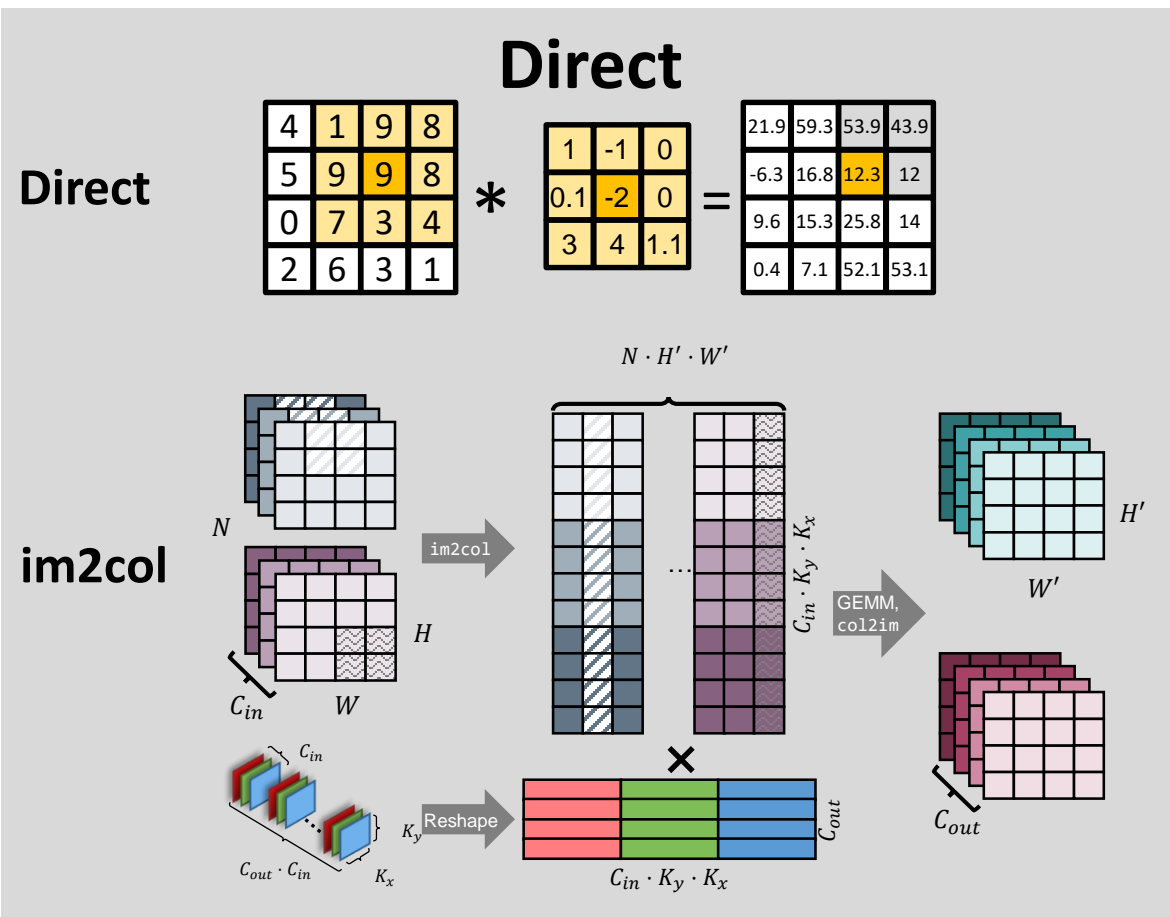
Operator implementations: convolution

$$Y_{k,f,i,j} = \sum_{c=0}^{C-1} \sum_{a=-O}^O \sum_{b=-O}^O X_{k,c,i+a,j+b} W_{f,c,a+O,b+O}$$



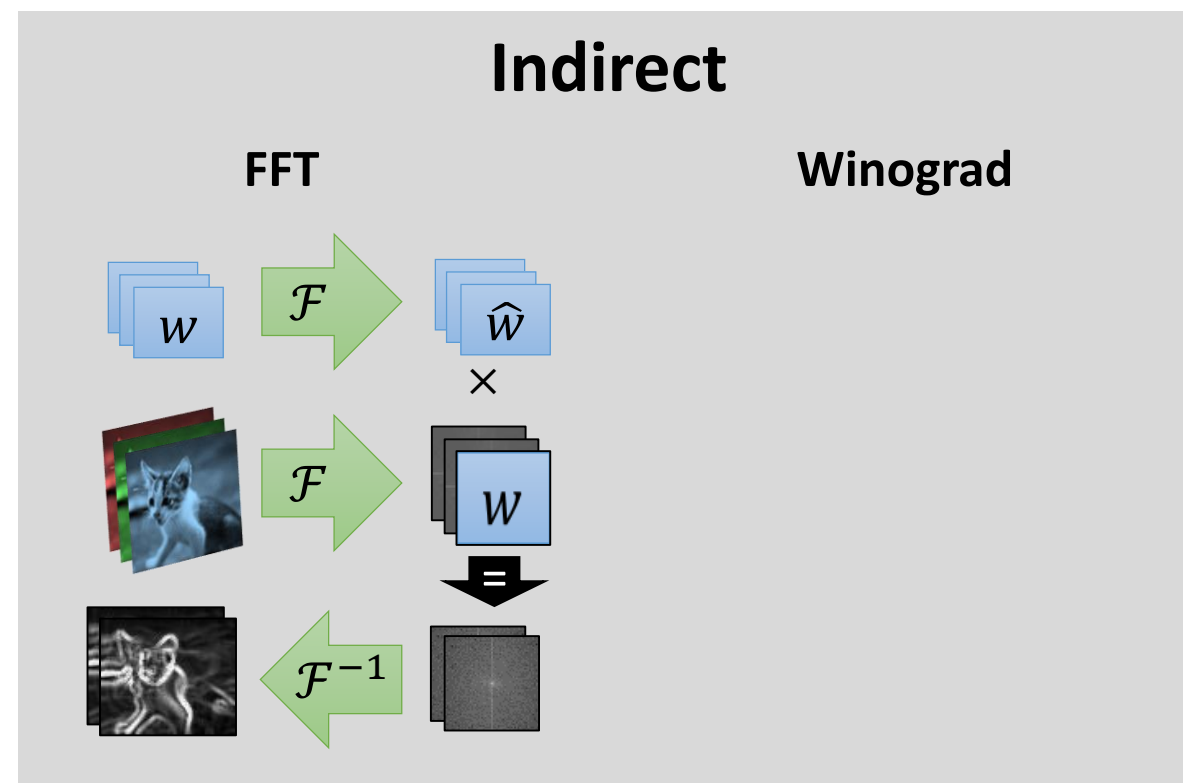
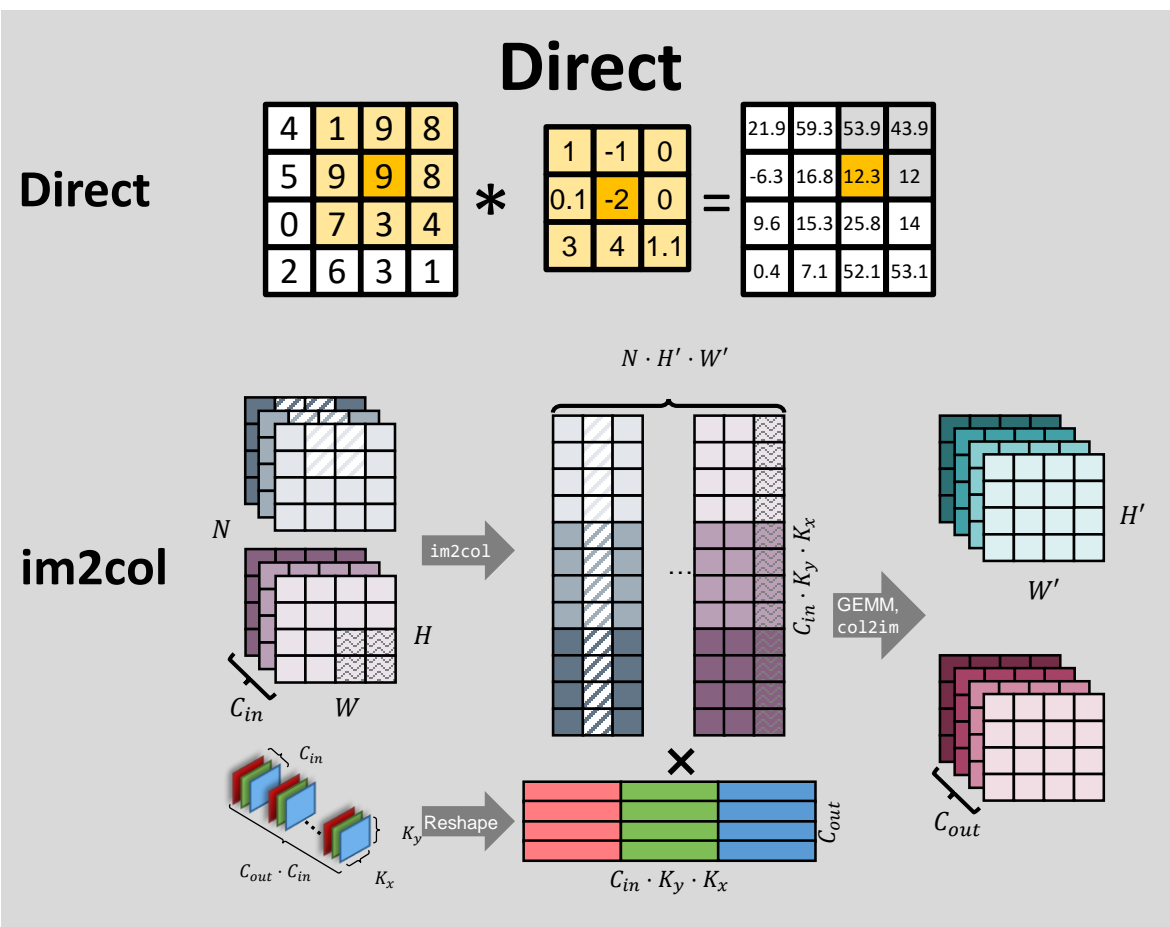
Operator implementations: convolution

$$Y_{k,f,i,j} = \sum_{c=0}^{C-1} \sum_{a=-O}^O \sum_{b=-O}^O X_{k,c,i+a,j+b} W_{f,c,a+O,b+O}$$



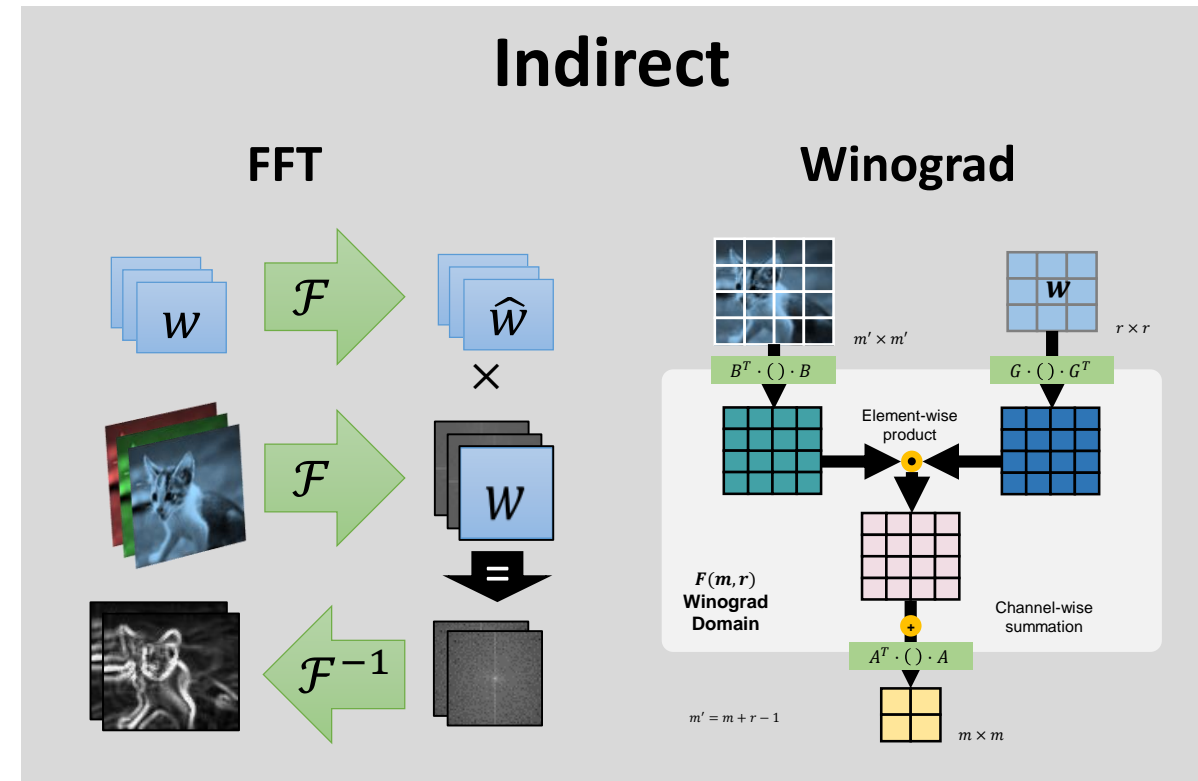
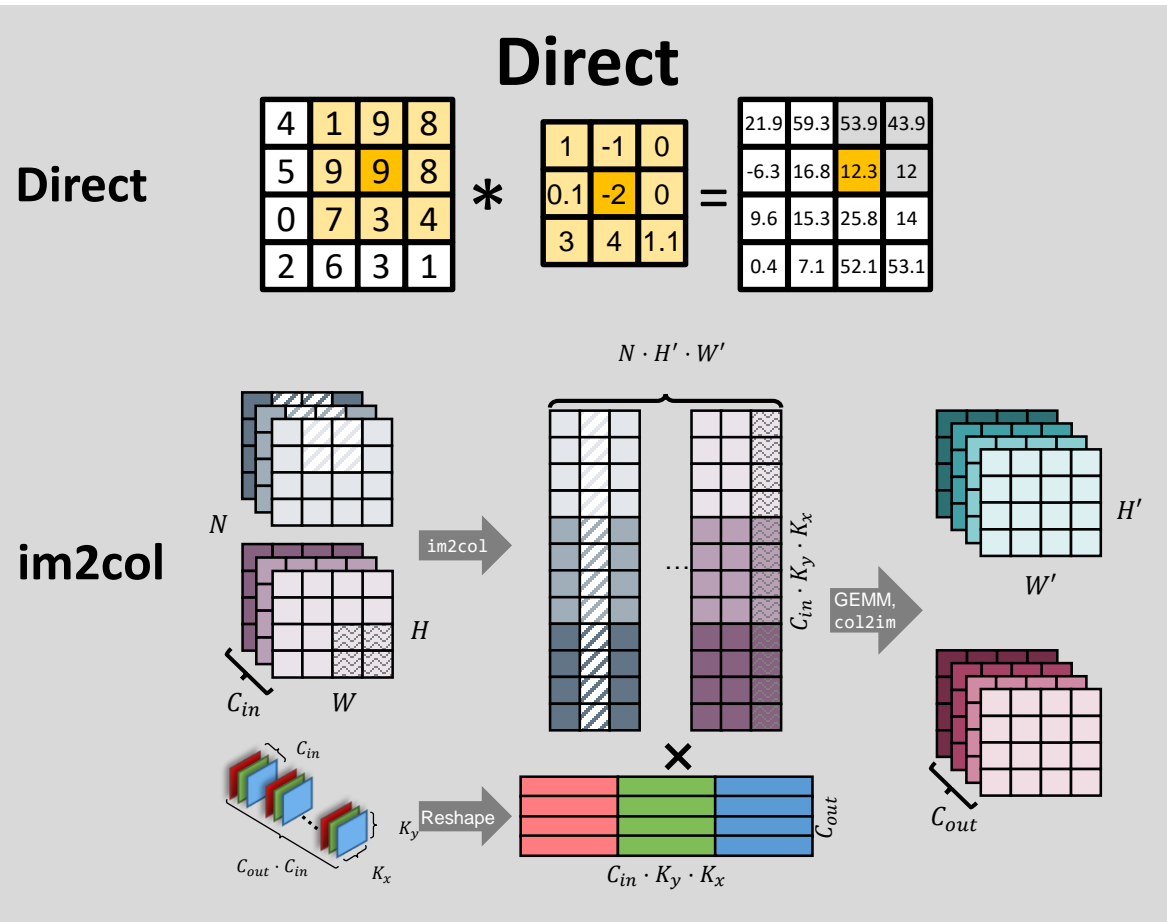
Operator implementations: convolution

$$Y_{k,f,i,j} = \sum_{c=0}^{C-1} \sum_{a=-O}^O \sum_{b=-O}^O X_{k,c,i+a,j+b} W_{f,c,a+O,b+O}$$



Operator implementations: convolution

$$Y_{k,f,i,j} = \sum_{c=0}^{C-1} \sum_{a=-O}^O \sum_{b=-O}^O X_{k,c,i+a,j+b} W_{f,c,a+O,b+O}$$



Microbatching (μ -cuDNN) – how to implement layers best in practice?

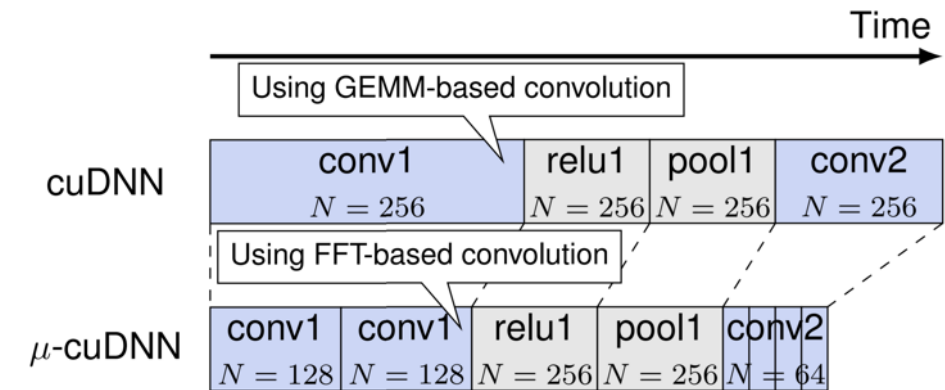
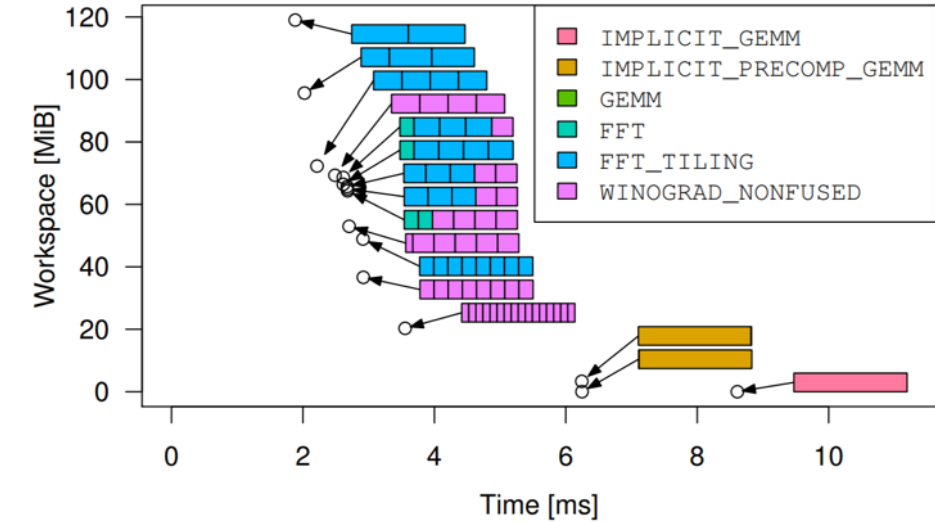
Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms
- How to choose microbatch sizes and algorithms?



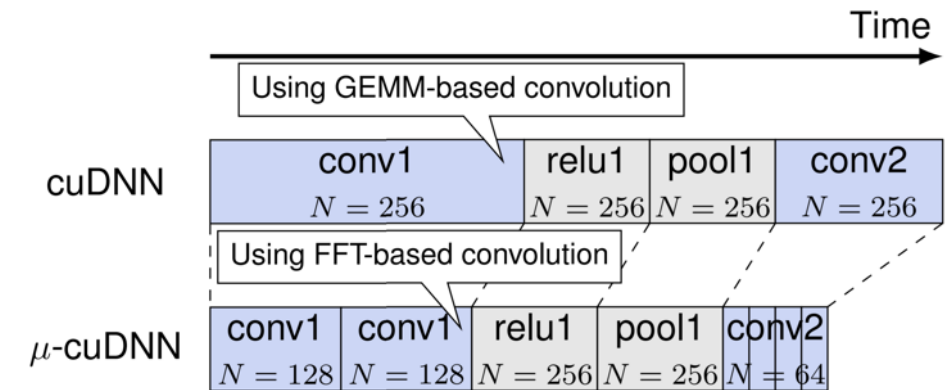
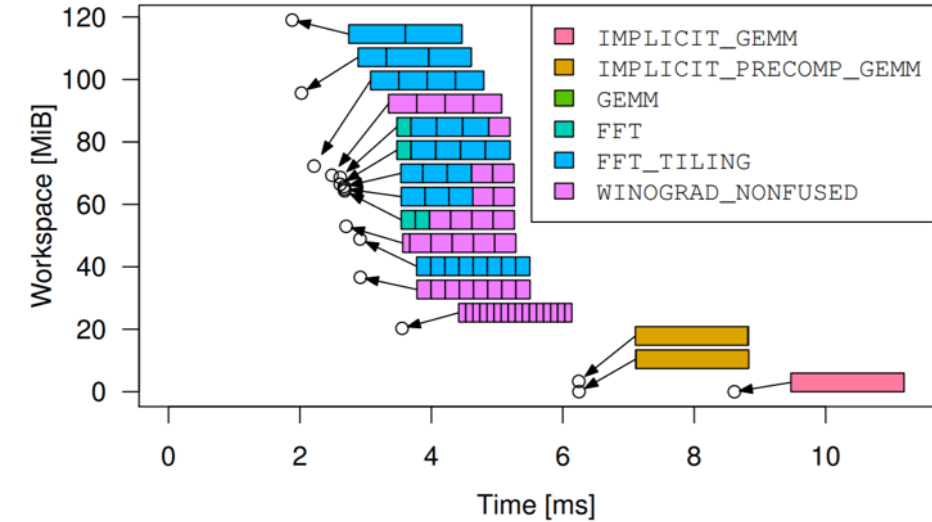
Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms
- How to choose microbatch sizes and algorithms?



Microbatching (μ -cuDNN) – how to implement layers best in practice?

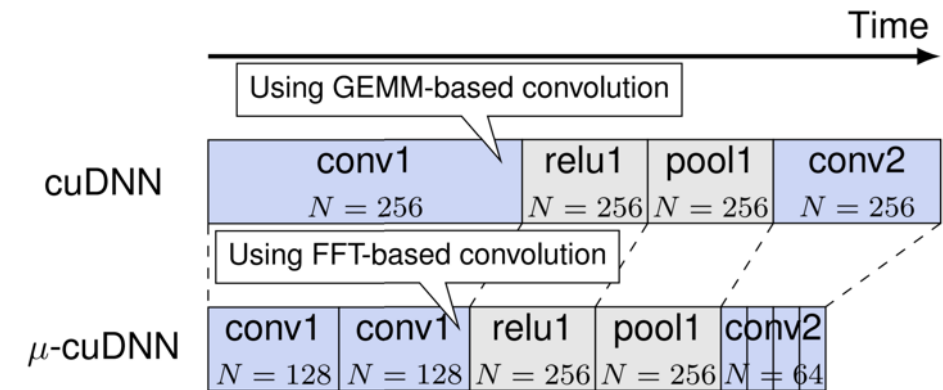
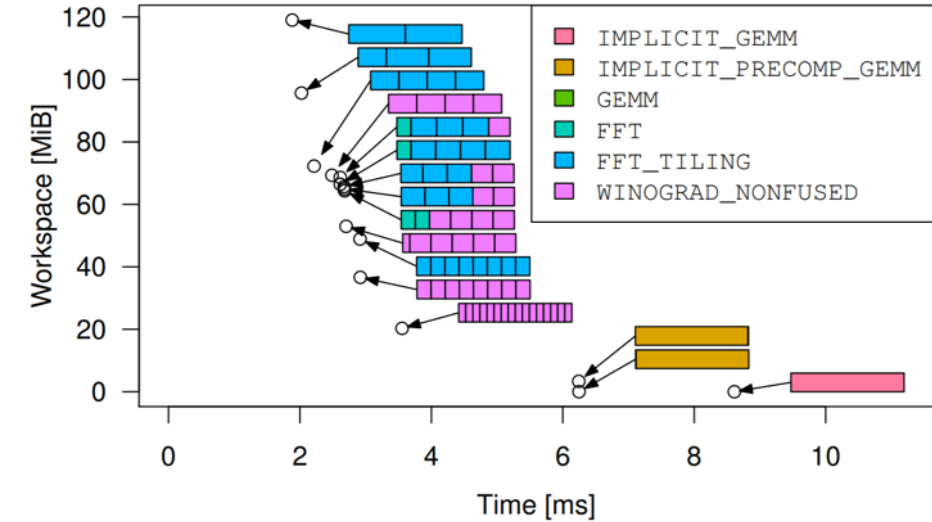
- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms
- How to choose microbatch sizes and algorithms?



Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms
- How to choose microbatch sizes and algorithms?

$$T(b) = \min \left\{ \begin{array}{l} T_{\mu}(b), \\ \min_{b'=1,2,\dots,b-1} T(b') + T(b - b') \end{array} \right\}$$

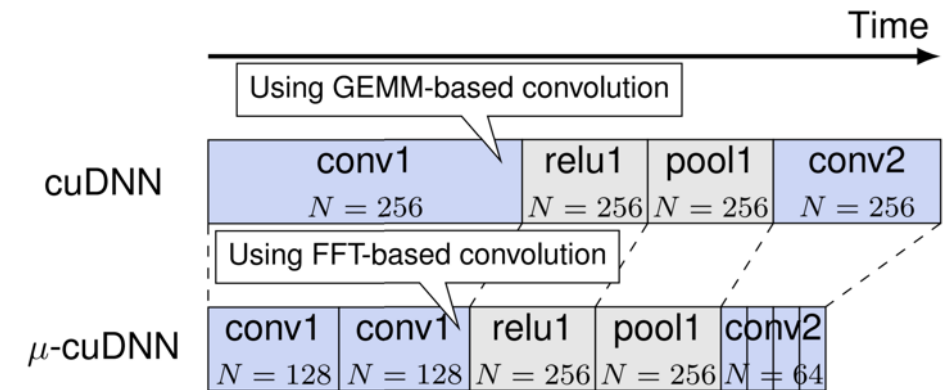
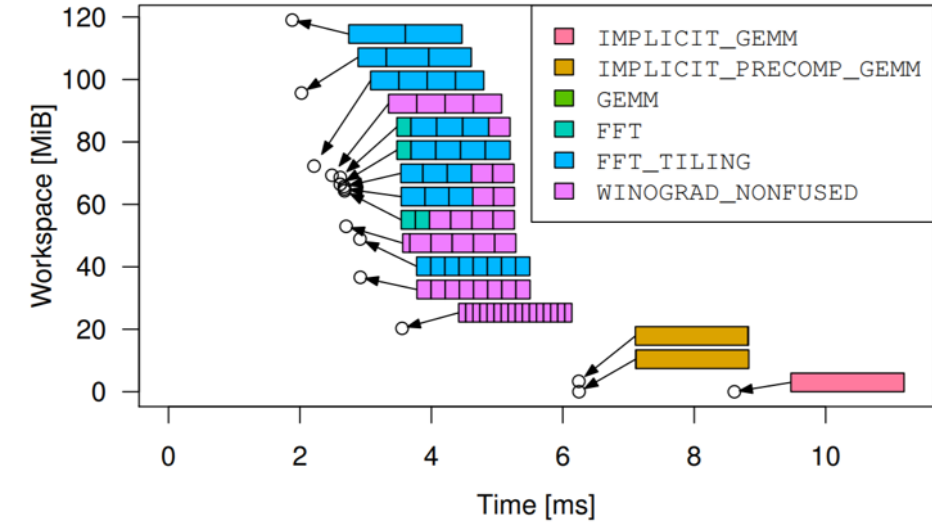


Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms
- How to choose microbatch sizes and algorithms?

$$T(b) = \min \left\{ \begin{array}{l} T_{\mu}(b), \\ \min_{b'=1,2,\dots,b-1} T(b') + T(b-b') \end{array} \right\}$$

Dynamic Programming (Space Reuse)

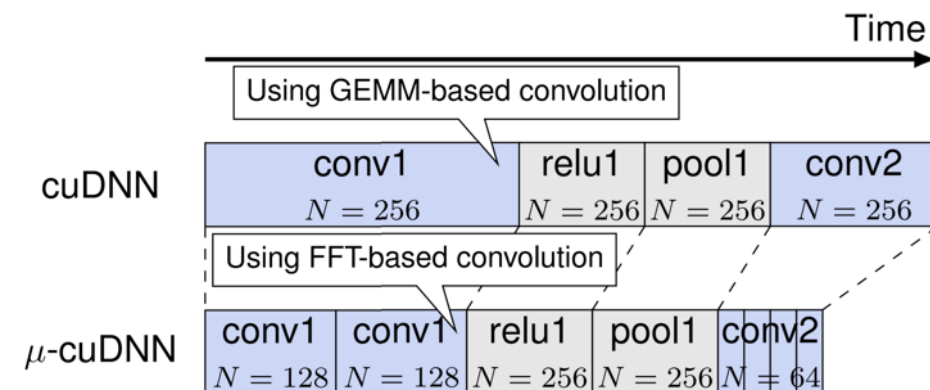
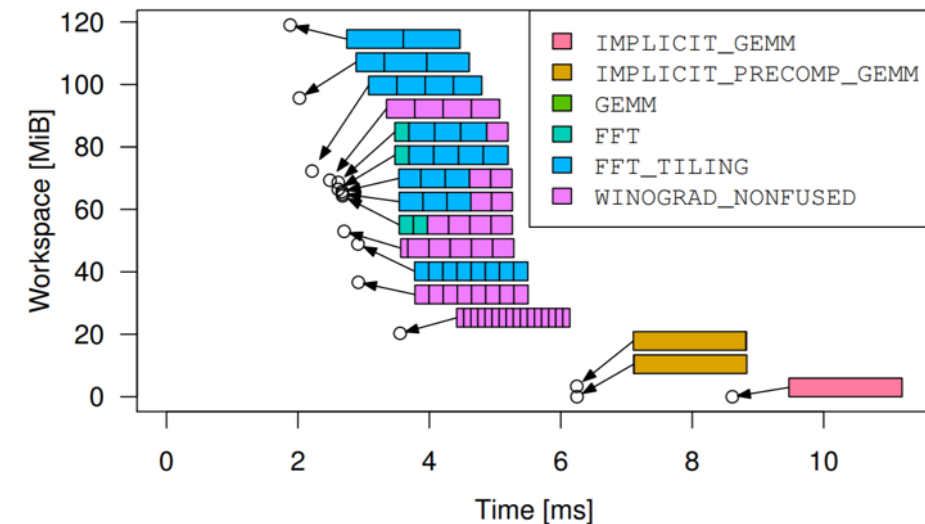


Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms
- How to choose microbatch sizes and algorithms?

$$T(b) = \min \left\{ \begin{array}{l} T_{\mu}(b), \\ \min_{b'=1,2,\dots,b-1} T(b') + T(b-b') \end{array} \right\}$$

Dynamic Programming (Space Reuse)

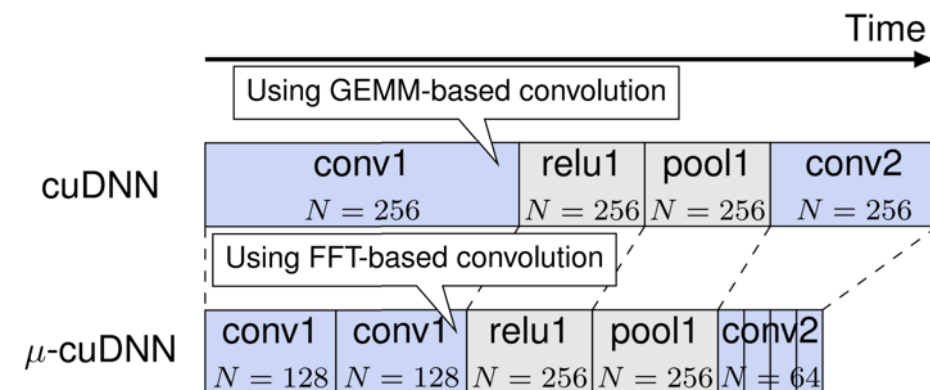
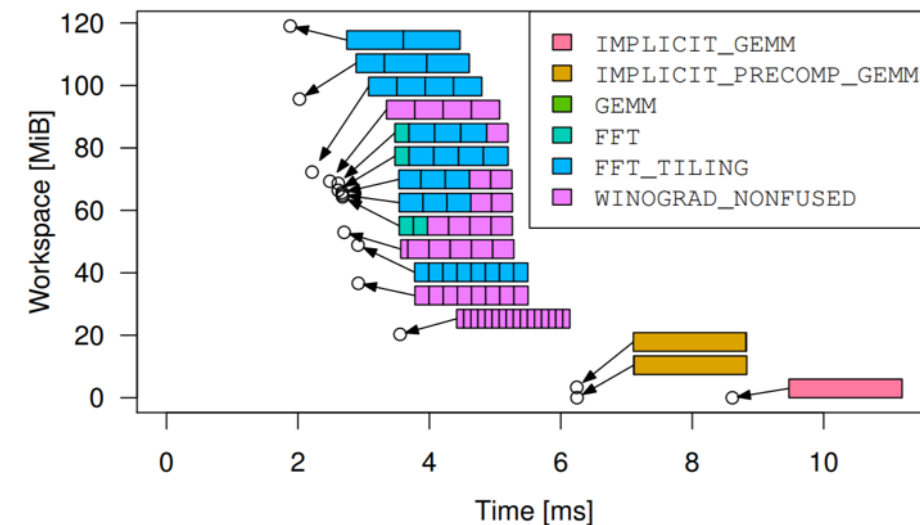


Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms
- How to choose microbatch sizes and algorithms?

$$T(b) = \min \left\{ \begin{array}{l} T_{\mu}(b), \\ \min_{b'=1,2,\dots,b-1} T(b') + T(b-b') \end{array} \right\}$$

Dynamic Programming (Space Reuse)



Microbatching (μ -cuDNN) – how to implement layers best in practice?

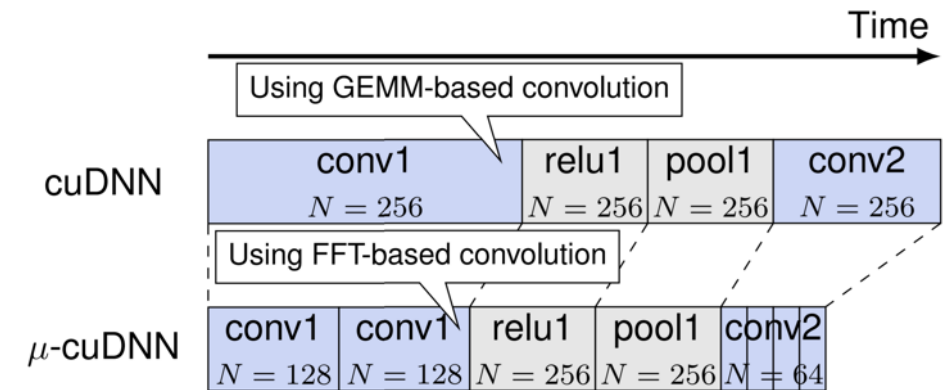
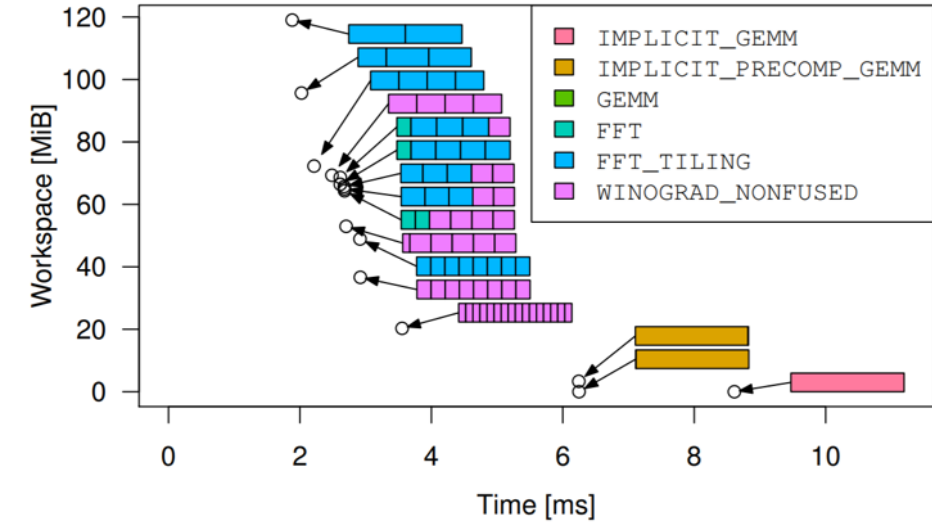
- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms
- How to choose microbatch sizes and algorithms?

$$T(b) = \min \left\{ \begin{array}{l} T_{\mu}(b), \\ \min_{b'=1,2,\dots,b-1} T(b') + T(b - b') \end{array} \right\}$$

Dynamic Programming (Space Reuse)

$$\begin{array}{ll} \min & T = \sum_{k \in \mathcal{K}} \sum_{c \in C_k} T_k(c) x_{k,c} \\ \text{subject to} & \sum_{k \in \mathcal{K}} \sum_{c \in C_k} M_k(c) x_{k,c} \leq M \\ & \sum_{c \in C_k} x_{k,c} = 1 \quad (\forall k \in \mathcal{K}) \\ & x_{k,c} \in \{0, 1\} \quad (\forall k \in \mathcal{K}, \forall c \in C_k) \end{array}$$

Integer Linear Programming (Space Sharing)



Microbatching (μ -cuDNN) – how to implement layers best in practice?

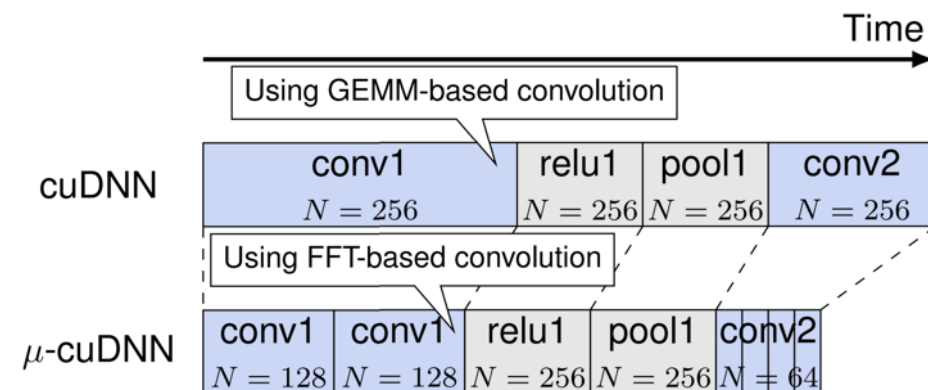
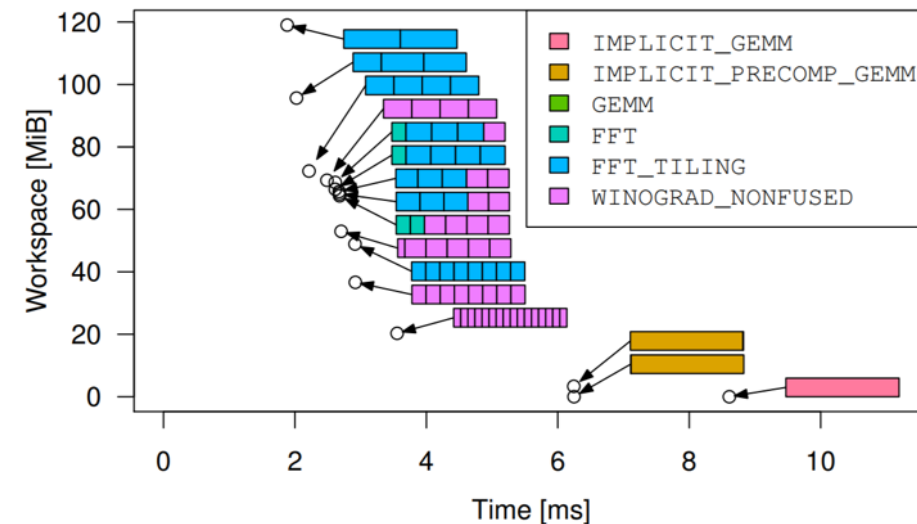
- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms
- How to choose microbatch sizes and algorithms?

$$T(b) = \min \left\{ \begin{array}{l} T_{\mu}(b), \\ \min_{b'=1,2,\dots,b-1} T(b') + T(b - b') \end{array} \right\}$$

Dynamic Programming (Space Reuse)

$$\begin{array}{ll} \min & T = \sum_{k \in \mathcal{K}} \sum_{c \in C_k} T_k(c) x_{k,c} \\ \text{subject to} & \sum_{k \in \mathcal{K}} \sum_{c \in C_k} M_k(c) x_{k,c} \leq M \\ & \sum_{c \in C_k} x_{k,c} = 1 \quad (\forall k \in \mathcal{K}) \\ & x_{k,c} \in \{0, 1\} \quad (\forall k \in \mathcal{K}, \forall c \in C_k) \end{array}$$

Integer Linear Programming (Space Sharing)



Microbatching (μ -cuDNN) – how to implement layers best in practice?

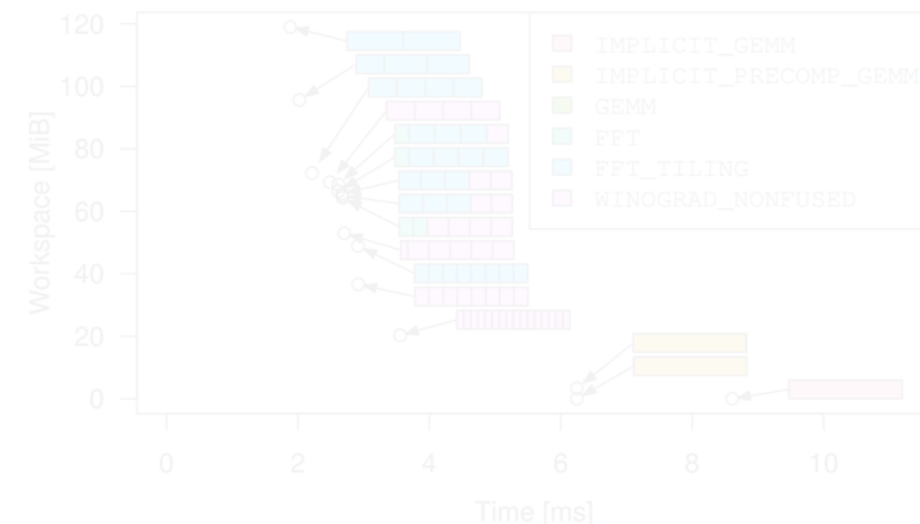
- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms
- How to choose microbatch sizes and algorithms?

$$T(b) = \min \left\{ T_{\mu}(b), \min_{b'=1,2,\dots,b-1} T(b') + T(b-b') \right\}$$

Dynamic Programming (Space Reuse)

$$\begin{aligned} \min \quad & T = \sum_{k \in \mathcal{K}} \sum_{c \in C_k} T_k(c) x_{k,c} \\ \text{subject to} \quad & \sum_{k \in \mathcal{K}} \sum_{c \in C_k} M_k(c) x_{k,c} \leq M \\ & \sum_{c \in C_k} x_{k,c} = 1 \quad (\forall k \in \mathcal{K}) \\ & x_{k,c} \in \{0, 1\} \quad (\forall k \in \mathcal{K}, \forall c \in C_k) \end{aligned}$$

Integer Linear Programming (Space Sharing)



Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms

Fast (up to 4.54x faster on DeepBench)

- How to choose microbatch

$$T(b) = \min \left\{ T_\mu(b), \min_{b'=1,2,\dots} T_\mu(b') \right\}$$

Dynamic Program

$$\min T = \sum_{k \in \mathcal{K}} x_{k,c}$$

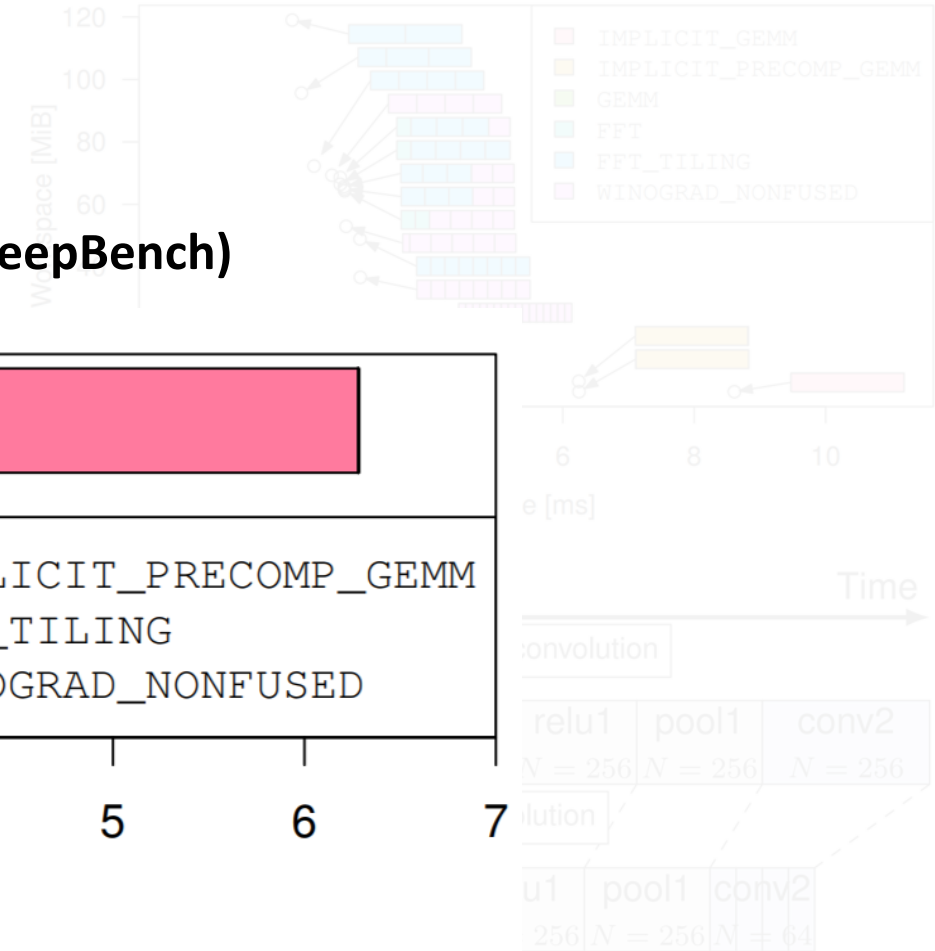
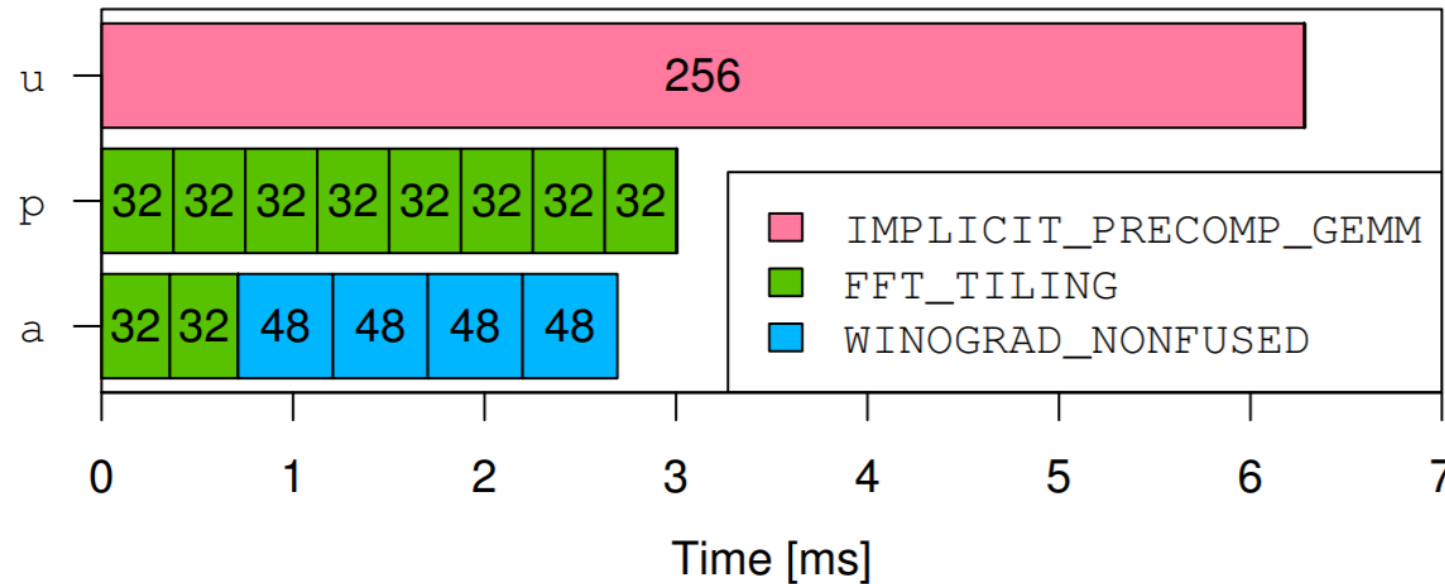
subject to

$$\sum_{k \in \mathcal{K}} \sum_{c \in C_k} x_{k,c} = 1$$

$$\sum_{c \in C_k} x_{k,c} = 1$$

$$x_{k,c} \in \{0, 1\} \quad (\forall k \in \mathcal{K}, \forall c \in C_k)$$

Integer Linear Programming (Space Sharing)



Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace, use different algorithms

Fast (up to 4.54x faster on DeepBench)

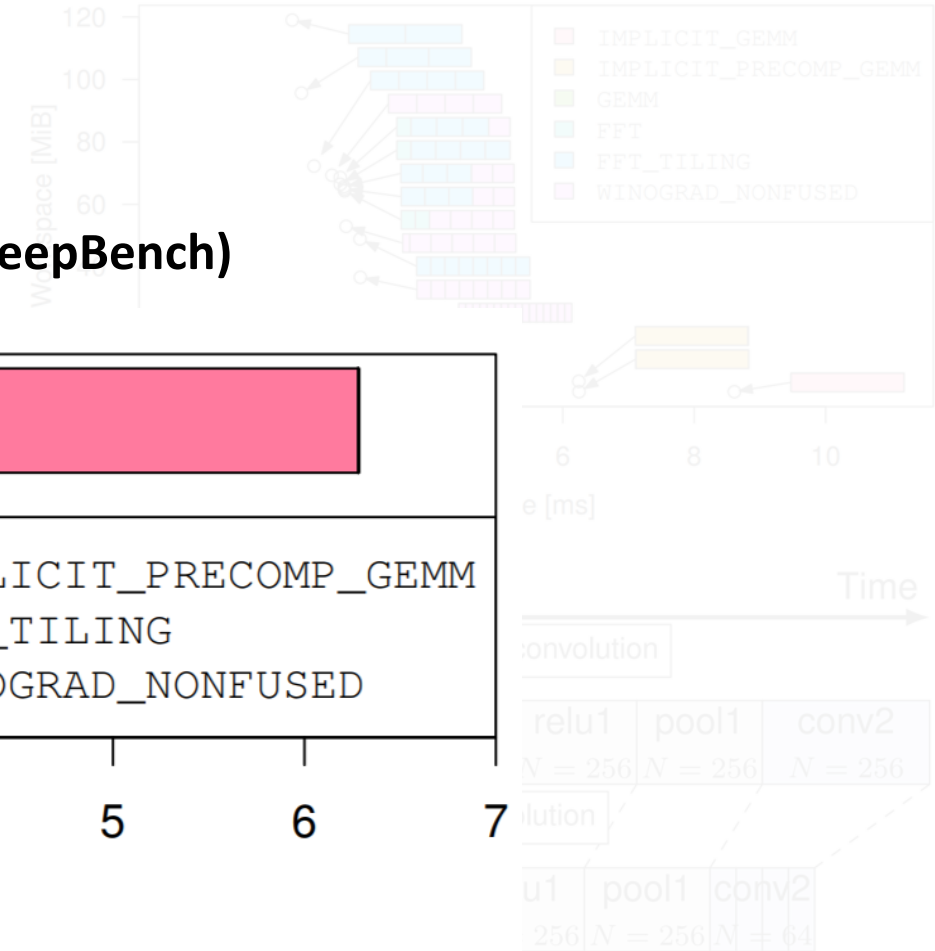
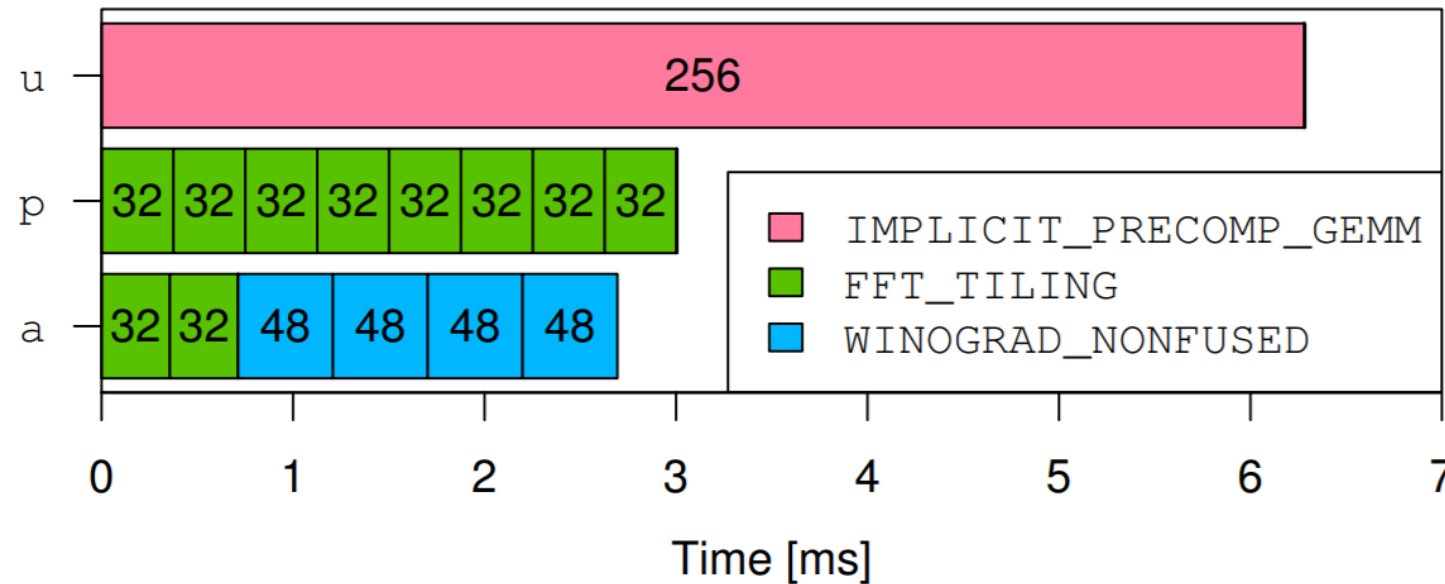
- How to choose microbatch

$$T(b) = \min \left\{ T_\mu(b), \min_{b'=1,2,\dots} T_\mu(b') \right\}$$

Dynamic Program

$$\begin{aligned} \min \quad & T = \sum_{k \in \mathcal{K}} x_{k,c} \\ \text{subject to} \quad & \sum_{k \in \mathcal{K}} \sum_{c \in C_k} x_{k,c} = 1 \\ & x_{k,c} \in \{0, 1\} \quad (\forall k \in \mathcal{K}, \forall c \in C_k) \end{aligned}$$

Integer Linear Programming (Space Sharing)



Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse

Fast (up to 4.54x faster on DeepBench)

Microbatching Strategy

- How to choose microbatching

$$T(b) = \min \left\{ T_\mu(b), \min_{b' \in \{1, 2, \dots\}} T_\mu(b') \right\}$$

Dynamic Program

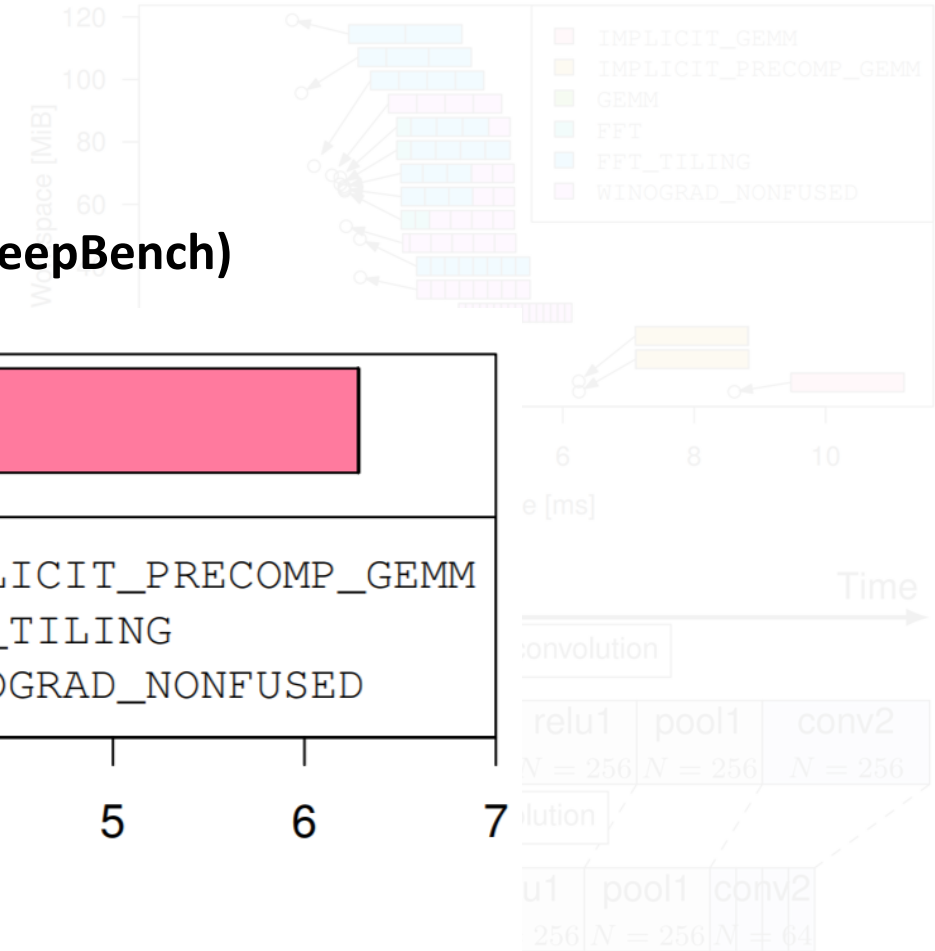
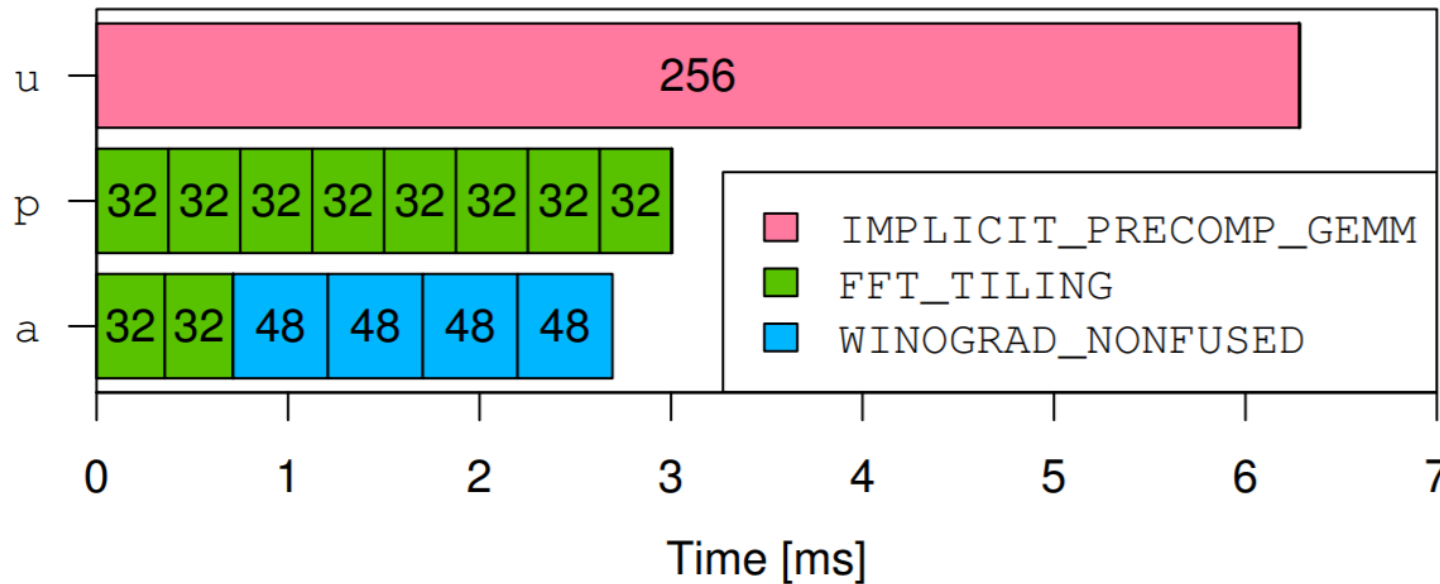
$$\min T = \sum_{k \in \mathcal{K}} x_{k,c}$$

subject to

$$\sum_{k \in \mathcal{K}} \sum_{c \in C_k} x_{k,c} = 1$$

$$\sum_{c \in C_k} x_{k,c} = 1 \quad (\forall k \in \mathcal{K}, \forall c \in C_k)$$

Integer Linear Programming (Space Sharing)



Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse

Fast (up to 4.54x faster on DeepBench)

Microbatching Strategy

- How to choose microbatching

none (undivided)

$$T(b) = \min \left\{ \begin{array}{l} T_{\mu}(\theta), \\ \min_{b'=1,2,\dots} \dots \end{array} \right.$$

Dynamic Program

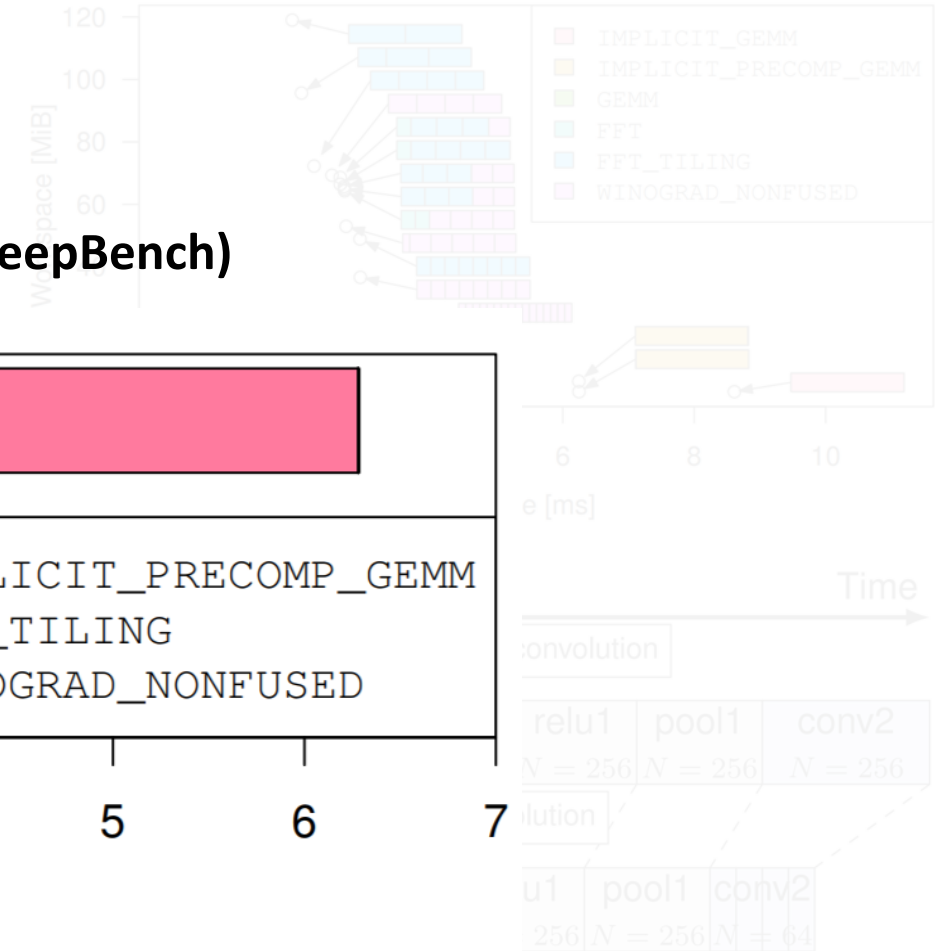
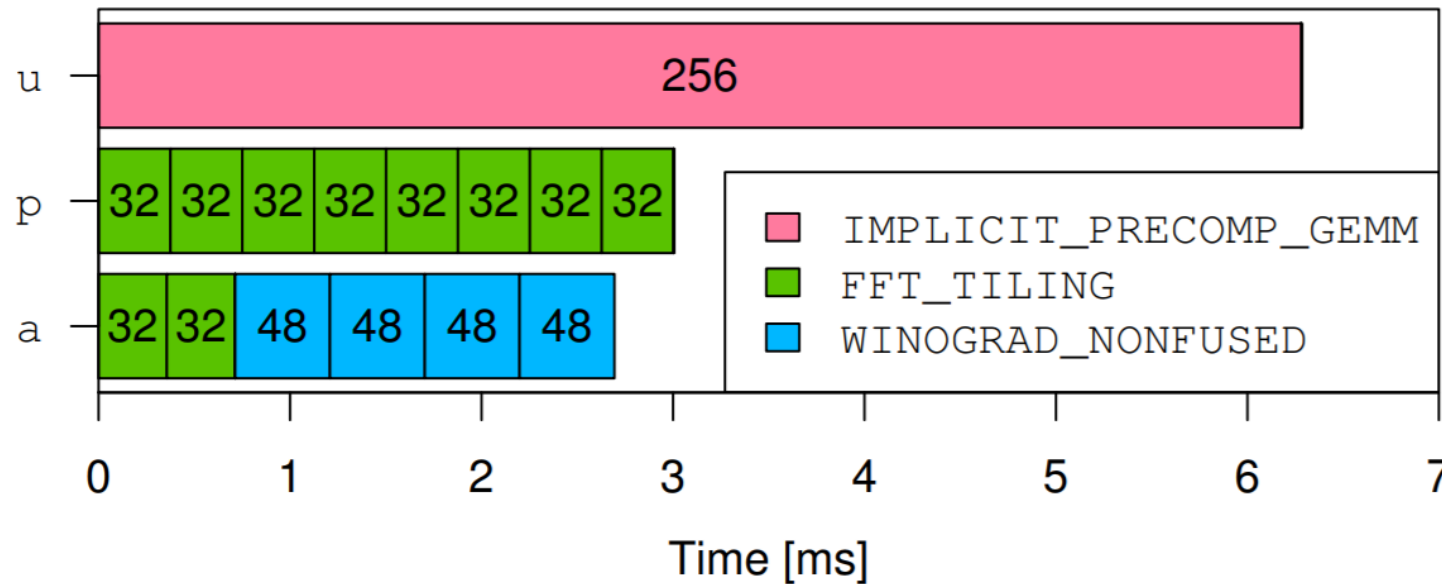
$$\min T = \sum_{k \in \mathcal{K}} x_{k,c}$$

subject to $\sum_{k \in \mathcal{K}} \sum_{c \in C_k} x_{k,c} = 1$

$$\sum_{c \in C_k} x_{k,c}$$

$$x_{k,c} \in \{0, 1\} \quad (\forall k \in \mathcal{K}, \forall c \in C_k)$$

Integer Linear Programming (Space Sharing)



Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse workspace across different algorithms

Fast (up to 4.54x faster on DeepBench)

Microbatching Strategy

- How to choose microbatching strategy

none (undivided)

$$T(b) = \min_{\mu} \left\{ \begin{array}{l} T_{\mu}(b), \\ \min_{\nu=1,2,\dots} \nu \cdot T(b/\nu) \end{array} \right.$$

powers-of-two only

Dynamic Program

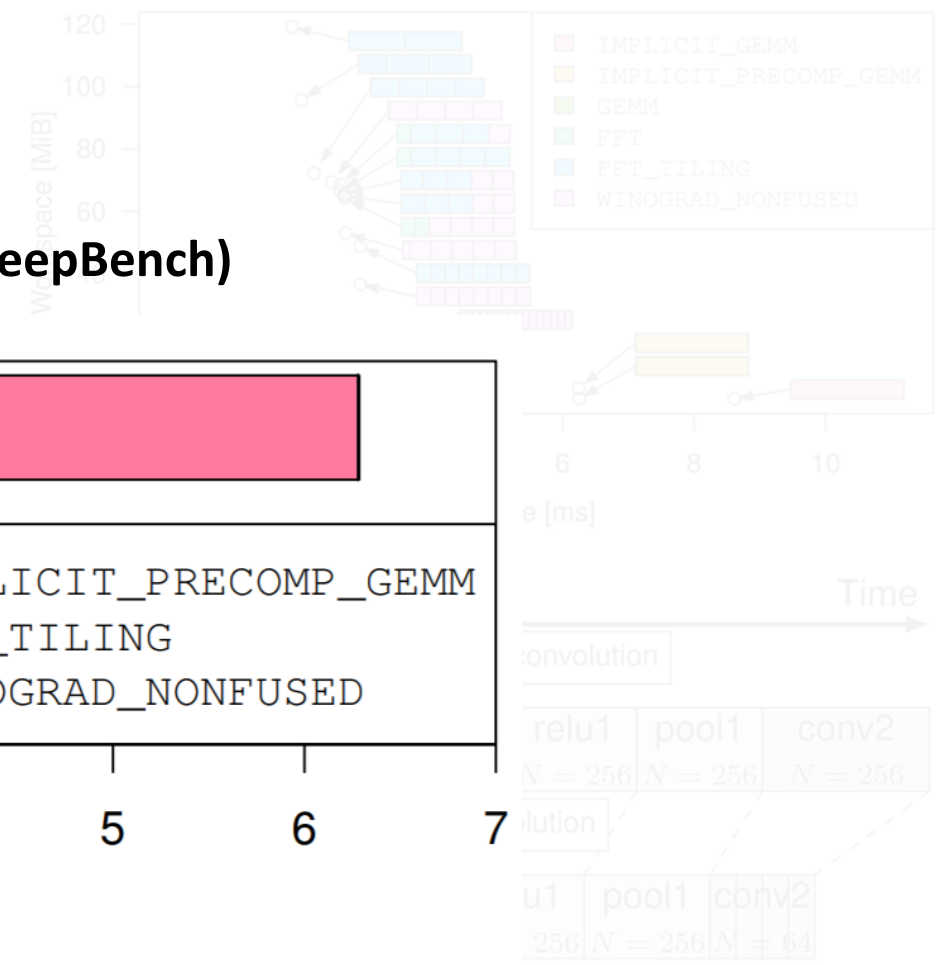
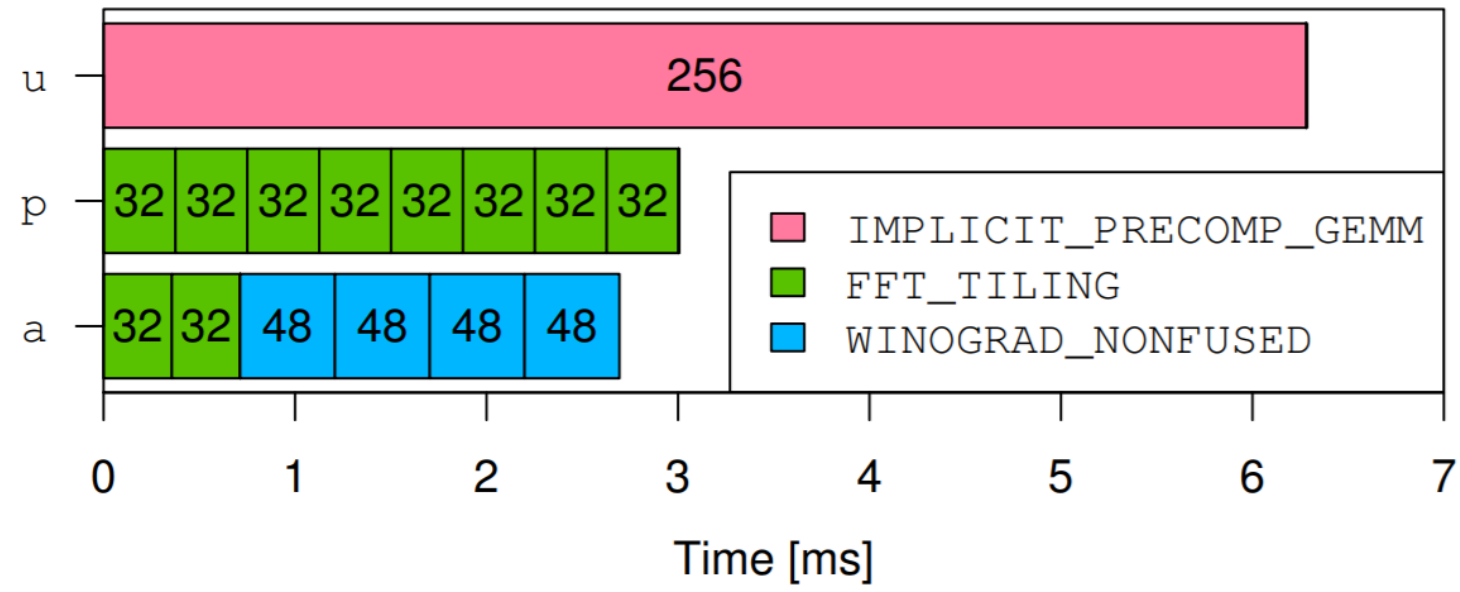
$$\min T = \sum_{k \in K} x_{k,c}$$

subject to $\sum_{k \in K} \sum_{c \in C_k} x_{k,c} = 1$

$$\sum_{c \in C_k} x_{k,c} = 1$$

$$x_{k,c} \in \{0, 1\} \quad (\forall k \in K, \forall c \in C_k)$$

Integer Linear Programming (Space Sharing)



Microbatching (μ -cuDNN) – how to implement layers best in practice?

- In cuDNN there are ~ 16 convolution implementations
- Performance depends on temporary memory (workspace) size
- Key idea: segment minibatch into microbatches, reuse

Fast (up to 4.54x faster on DeepBench)

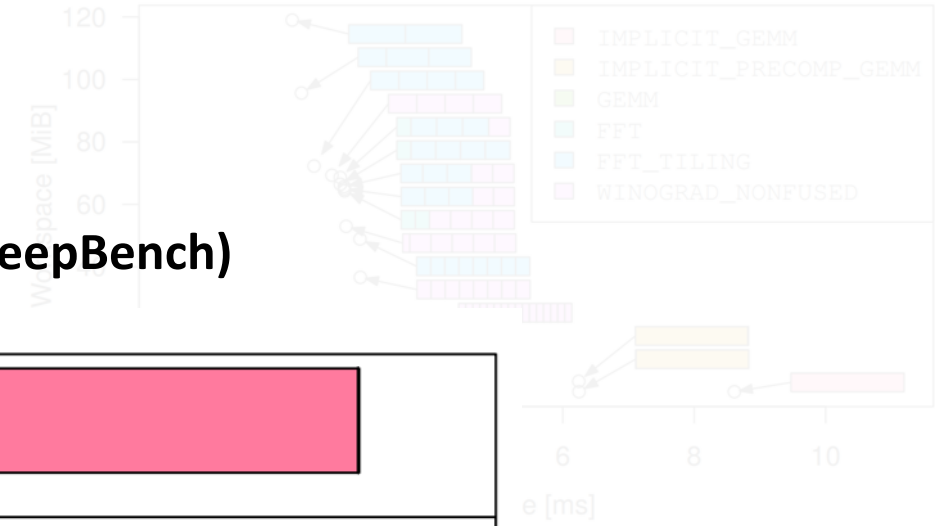
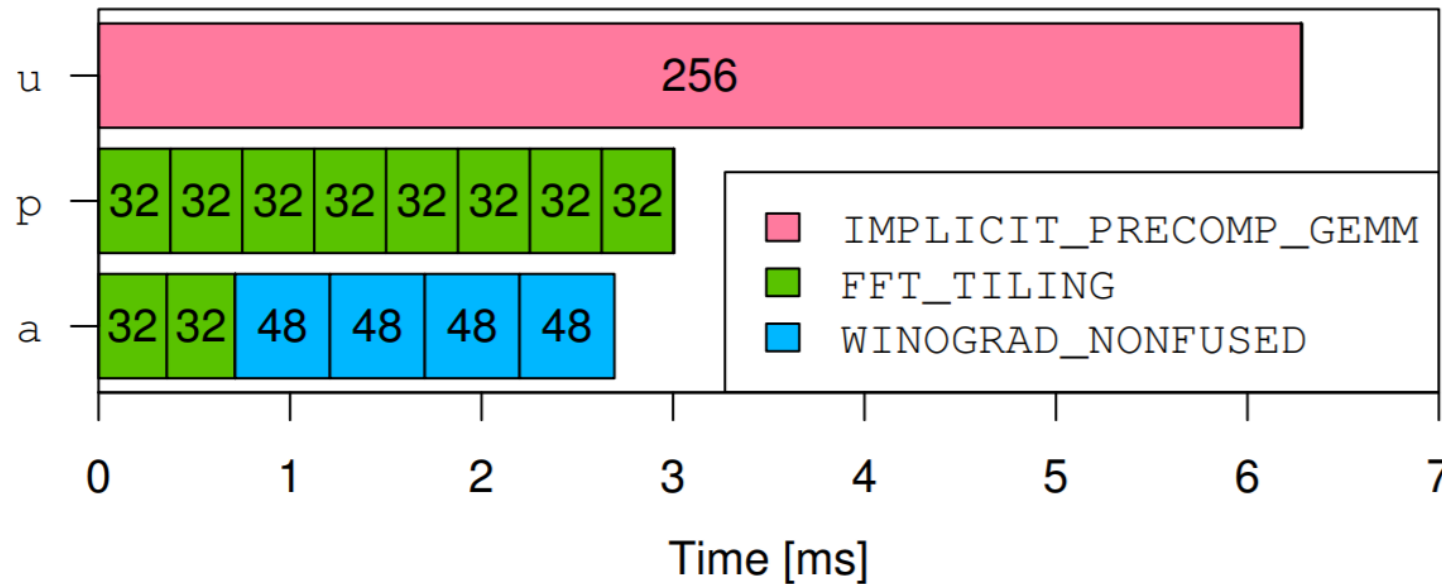
Microbatching Strategy

- How to choose microbatching strategy

none (undivided)

powers-of-two only

any (unrestricted)



Dynamic Program

$$T(b) = \min_{\mu} \left\{ T_{\mu}(b), \min_{\nu=1,2,\dots} T_{\nu}(b) \right\}$$

Integer Linear Programming (Space Sharing)

$$\min T = \sum_{k \in K} x_{k,c}$$

subject to

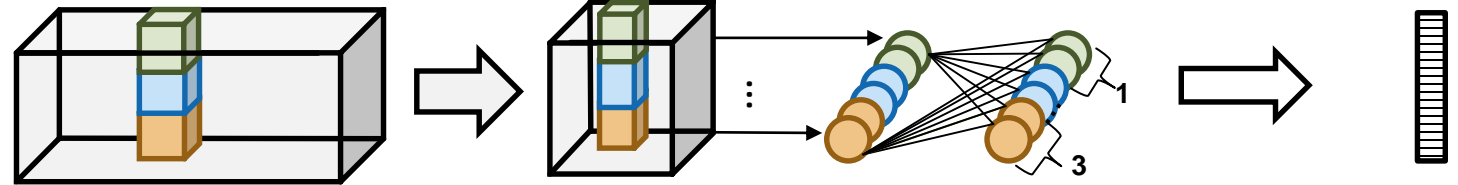
$$\sum_{k \in K} \sum_{c \in C_k} x_{k,c} \leq W$$

$$x_{k,c} \in \{0, 1\} \quad (\forall k \in K, \forall c \in C_k)$$

Model parallelism – limited by network size

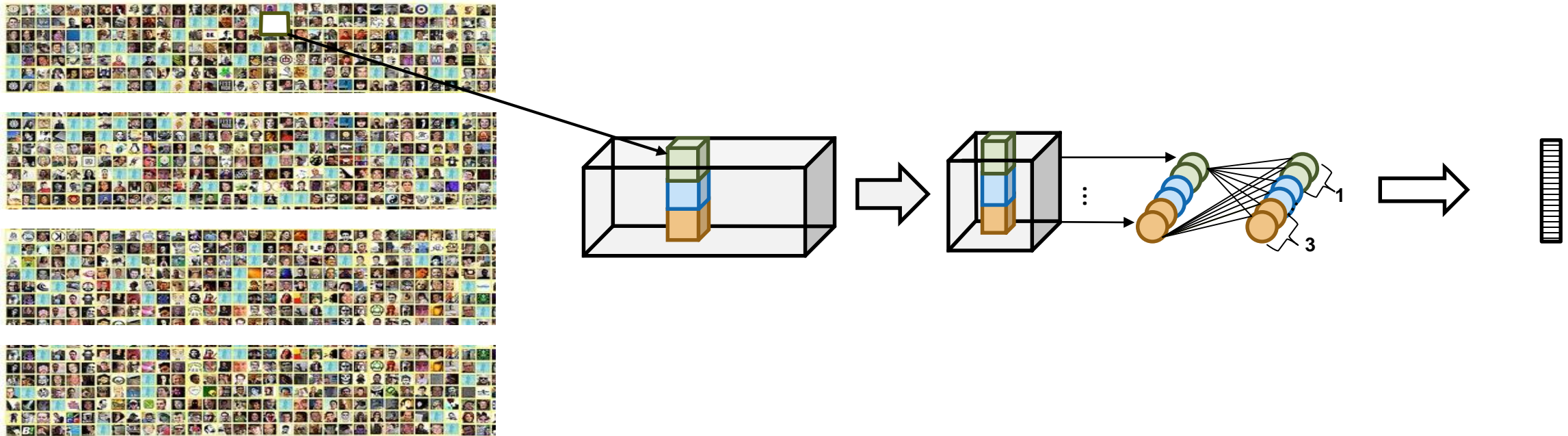


Model parallelism – limited by network size



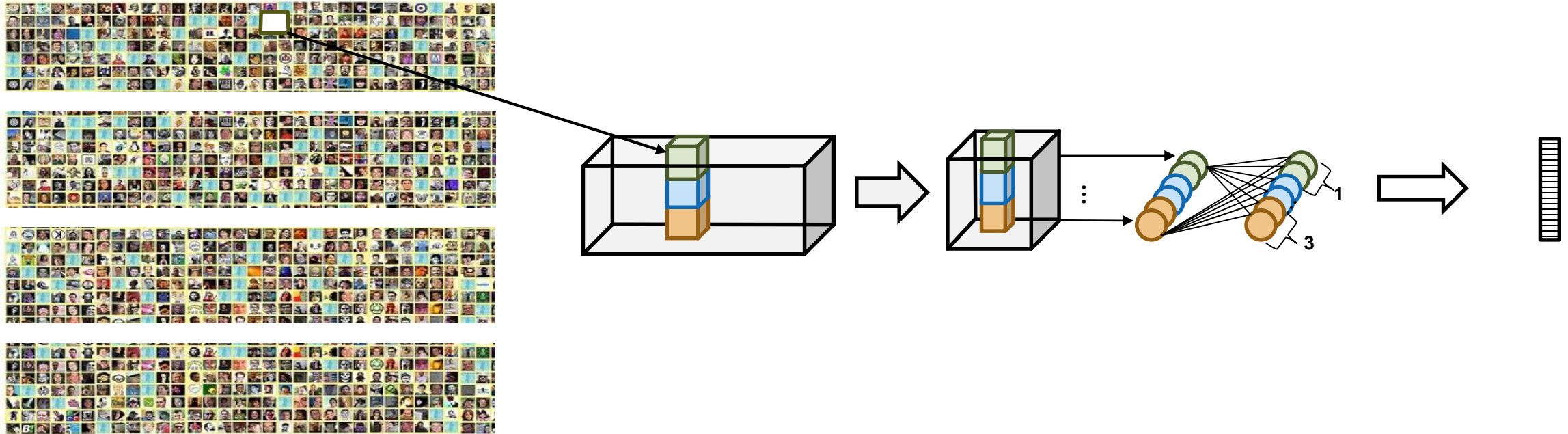
- Parameters can be distributed across processors
- **Mini-batch has to be copied to all processors**
- **Backpropagation requires complex communication every layer**

Model parallelism – limited by network size



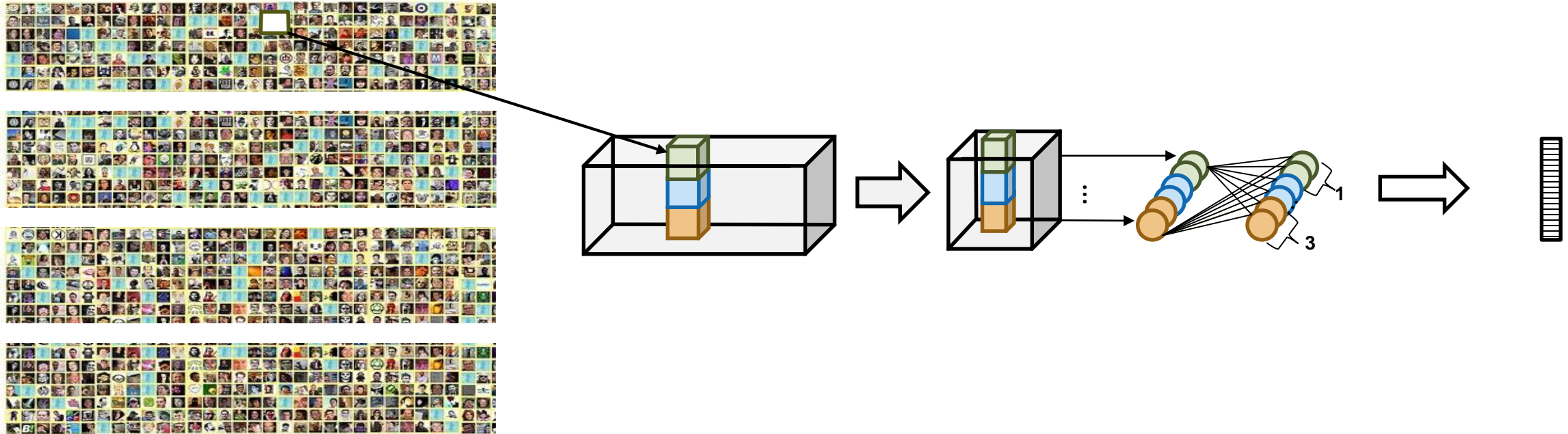
- Parameters can be distributed across processors
- **Mini-batch has to be copied to all processors**
- **Backpropagation requires complex communication every layer**

Model parallelism – limited by network size



- Parameters can be distributed across processors
- **Mini-batch has to be copied to all processors**
- **Backpropagation requires complex communication every layer**

Model parallelism – limited by network size



- Parameters can be distributed across processors
- **Mini-batch has to be copied to all processors**
- **Backpropagation requires complex communication every layer**

Pipeline parallelism – limited by network size

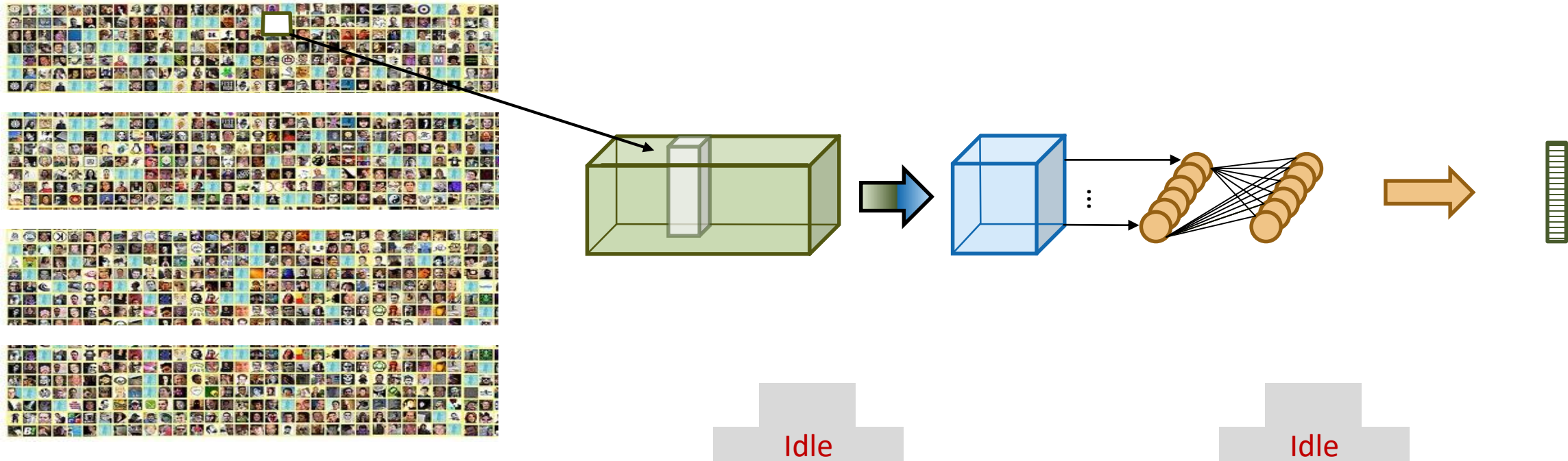


Pipeline parallelism – limited by network size



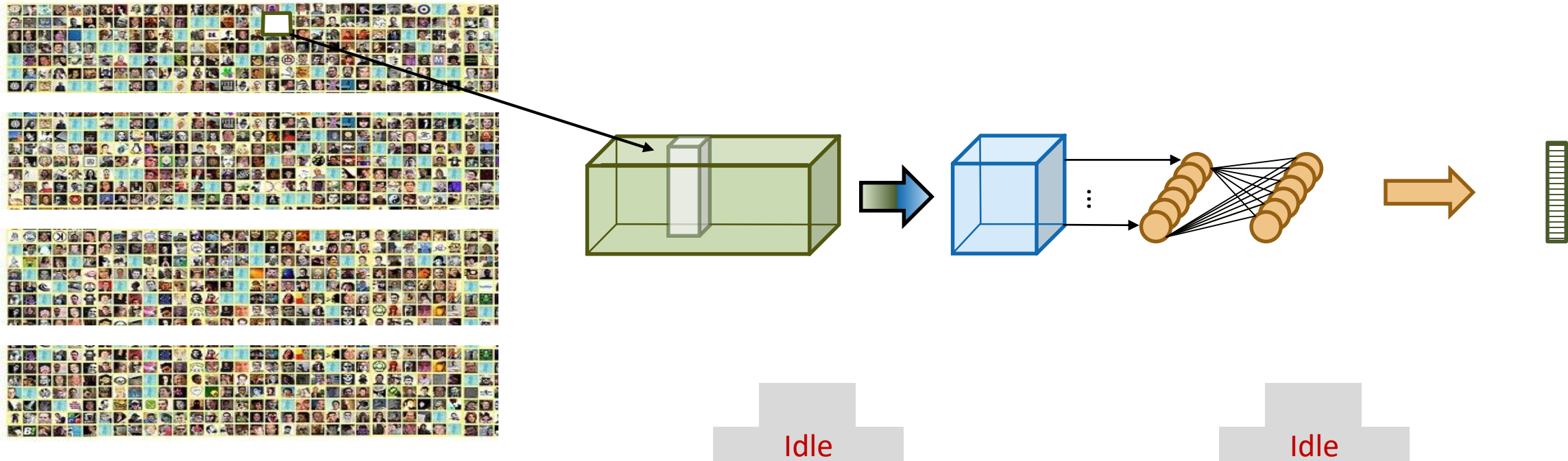
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- **Mini-batch has to be copied through all processors**
- **Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



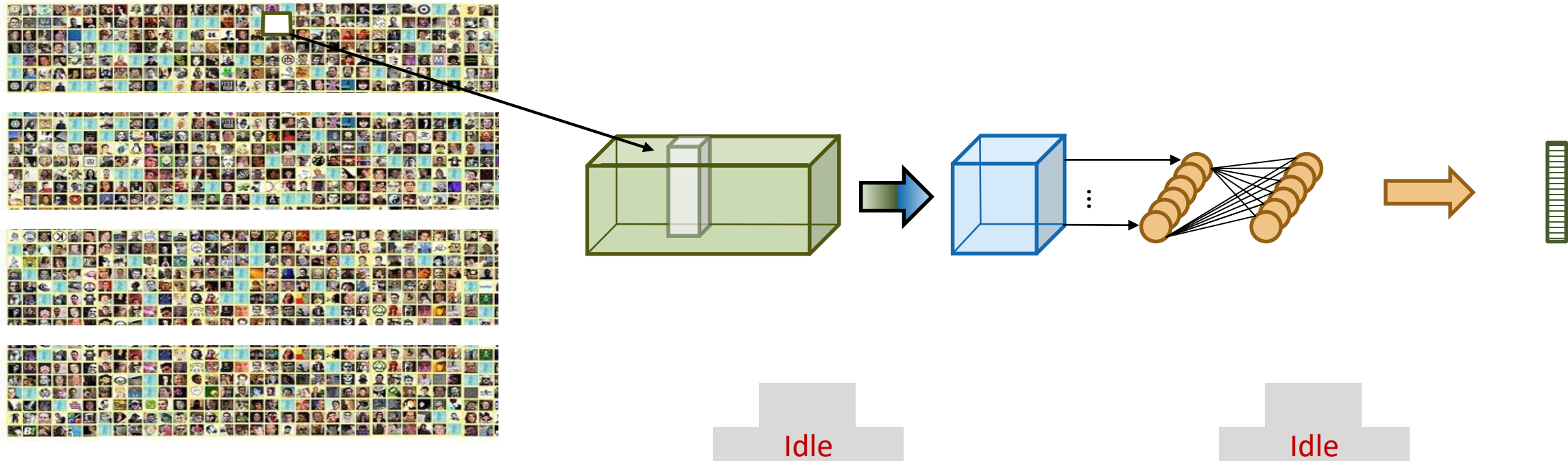
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- **Mini-batch has to be copied through all processors**
- **Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



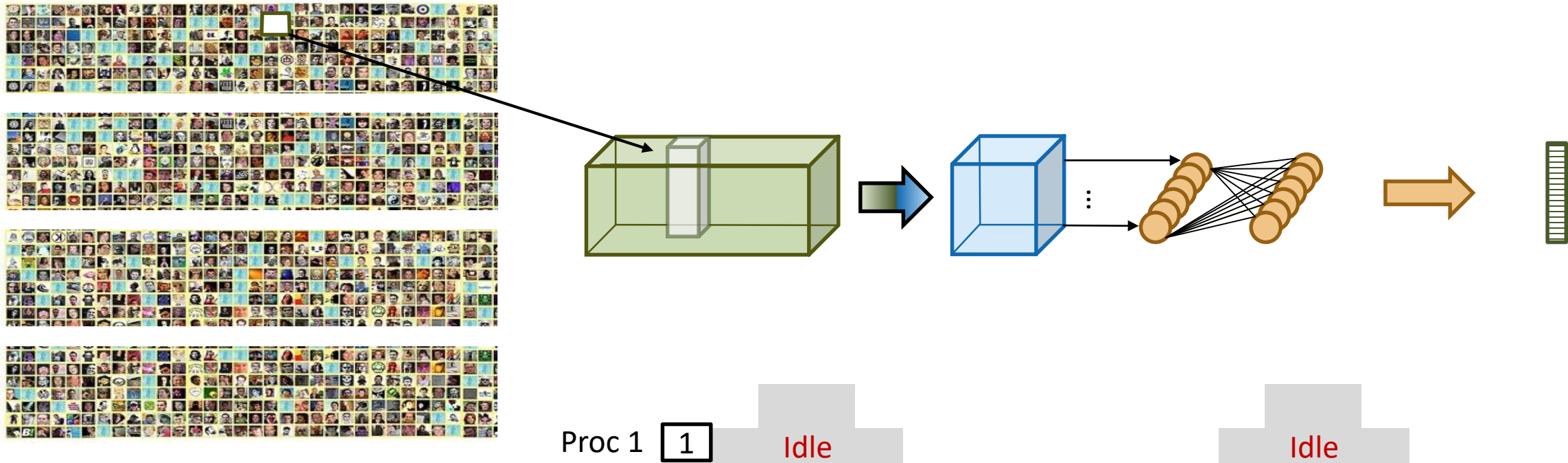
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- **Mini-batch has to be copied through all processors**
- **Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



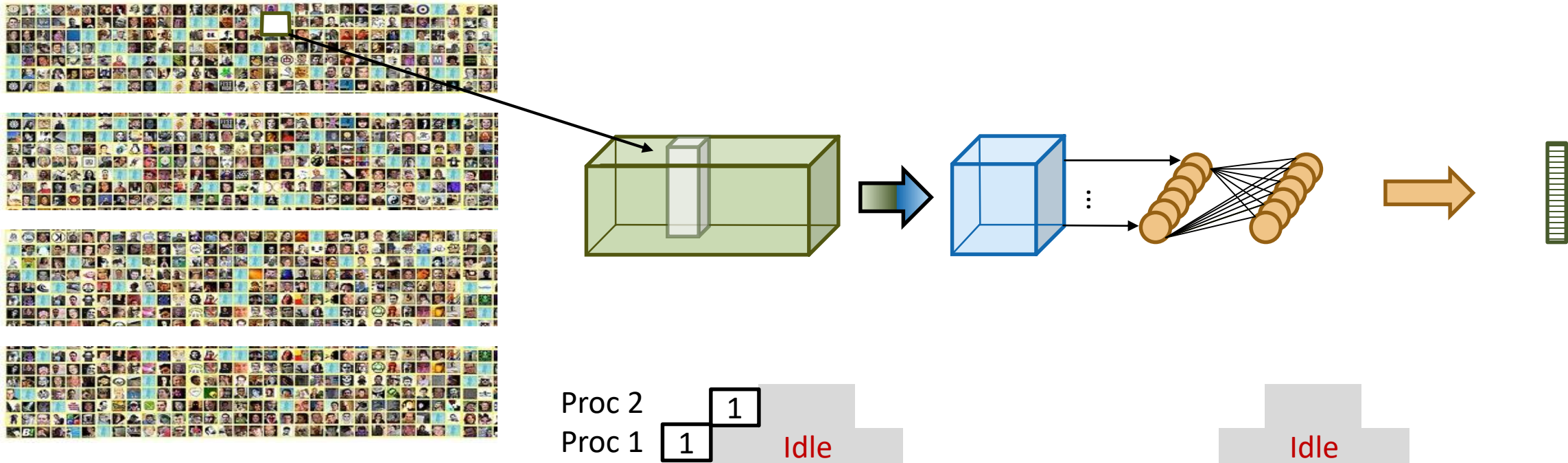
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- **Mini-batch has to be copied through all processors**
- **Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



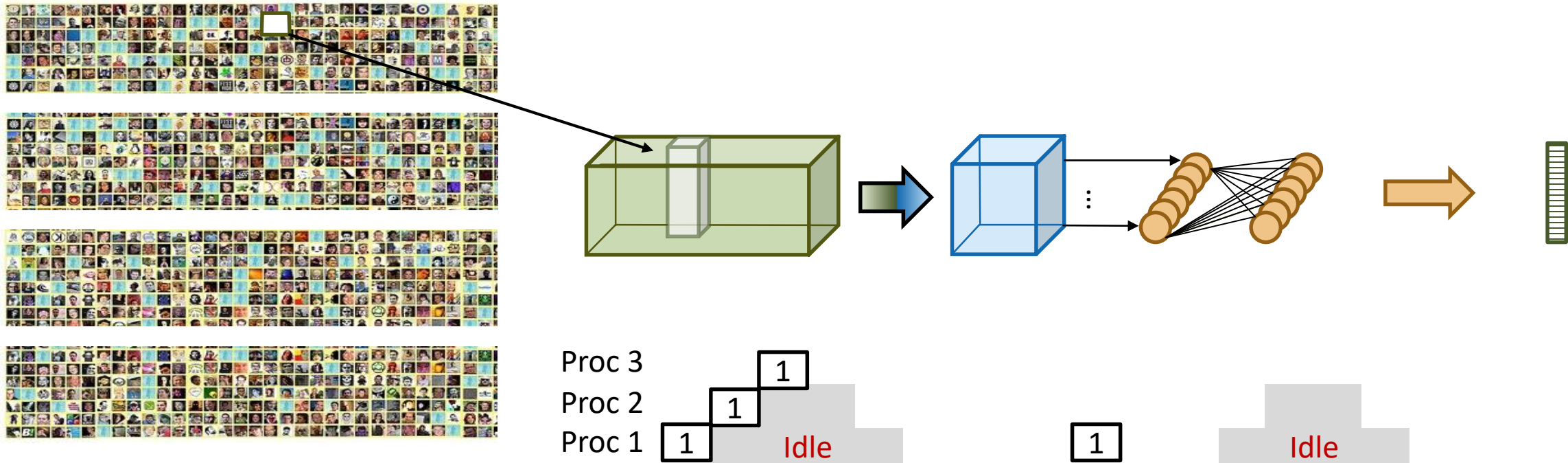
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- **Mini-batch has to be copied through all processors**
- **Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



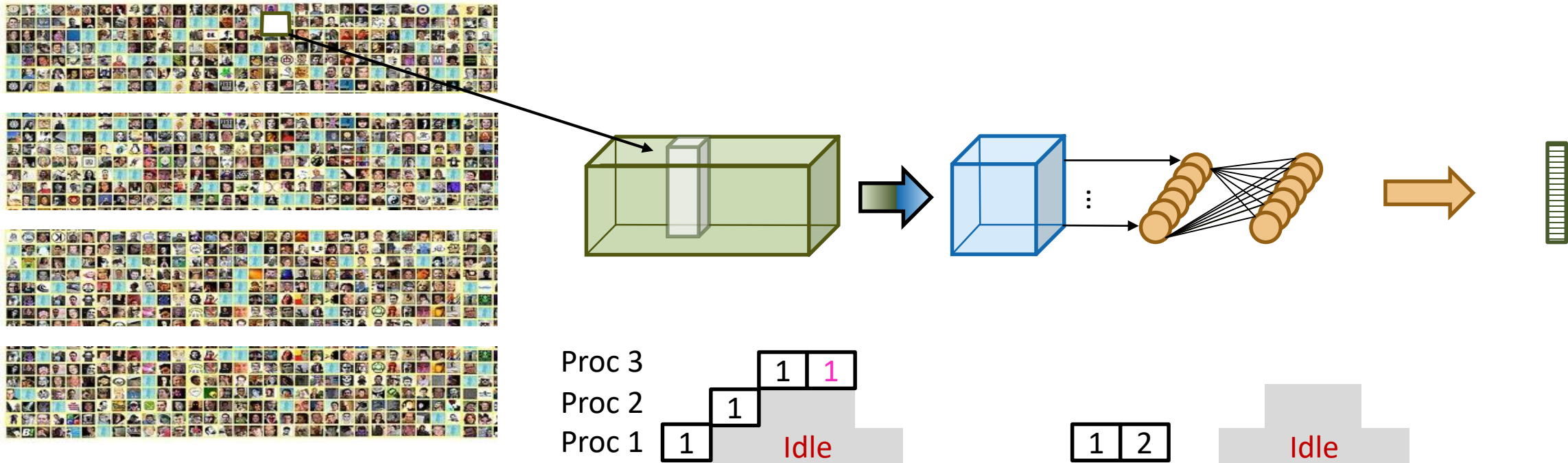
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



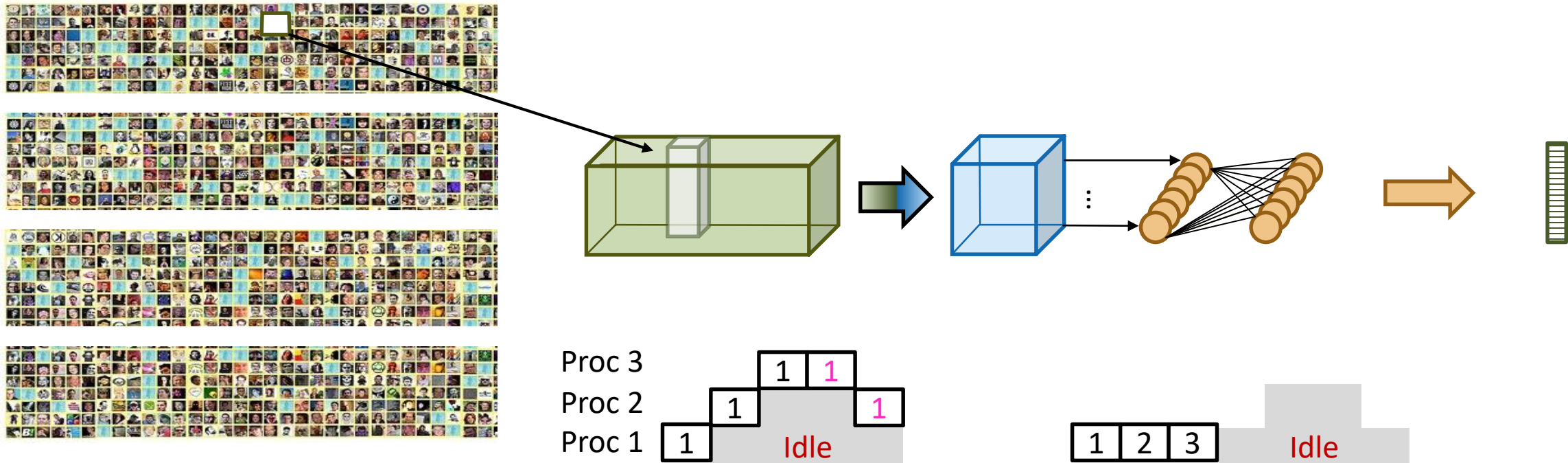
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



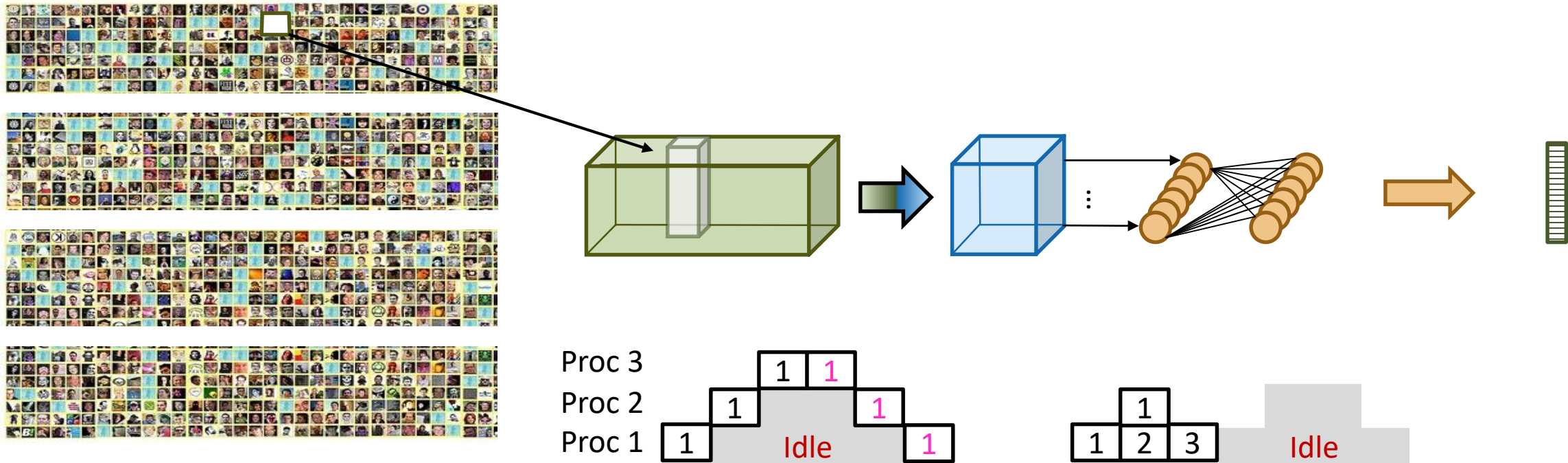
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



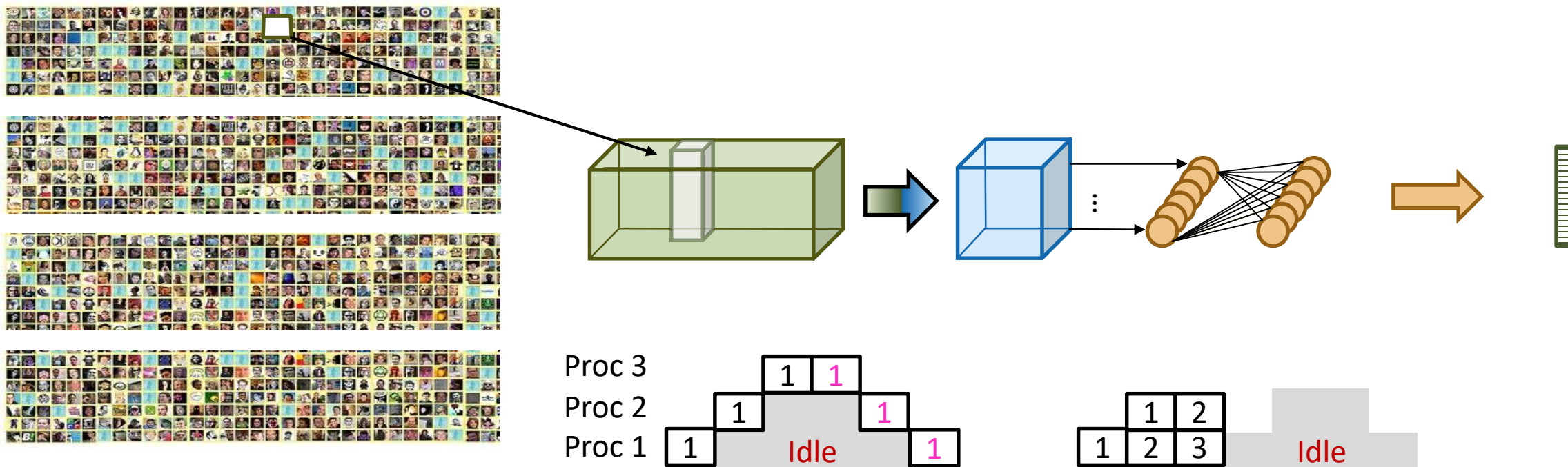
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



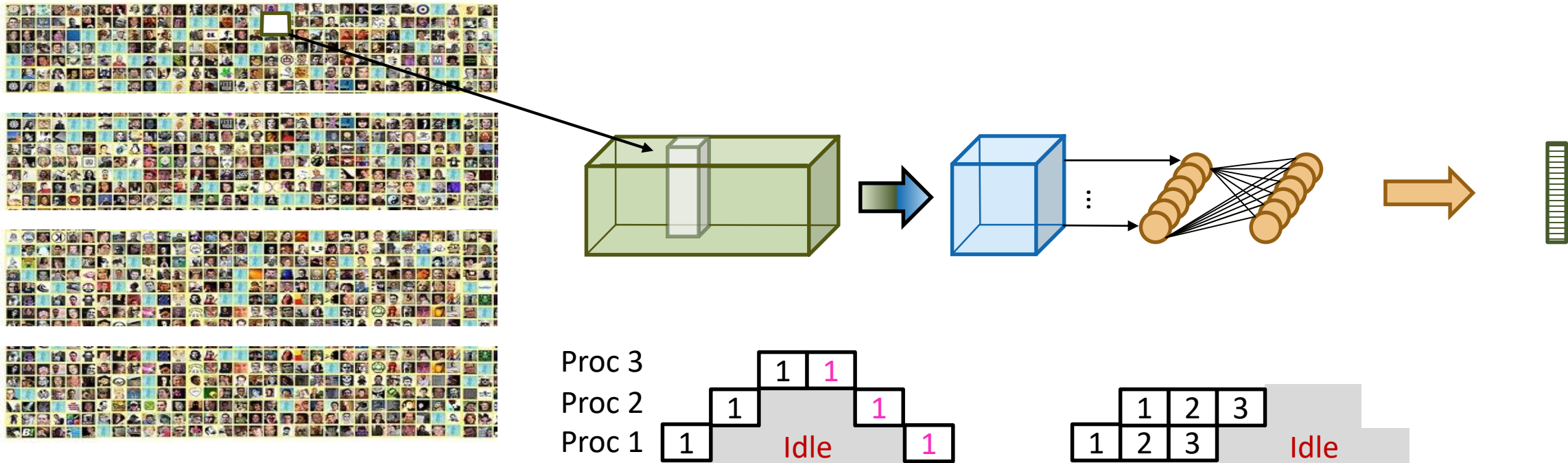
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



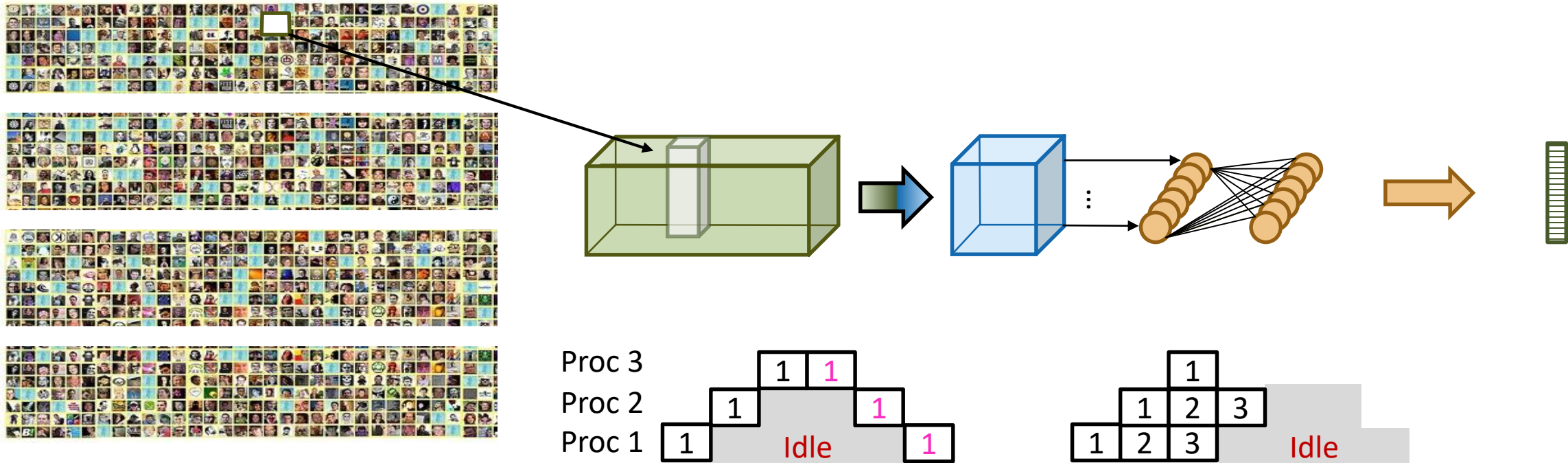
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



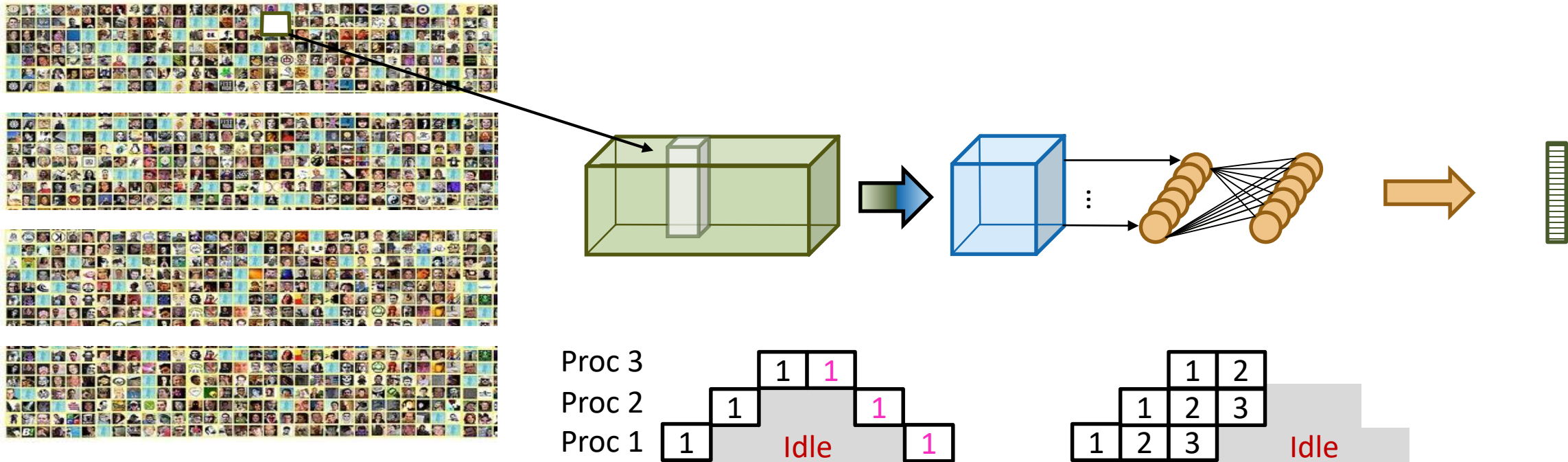
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



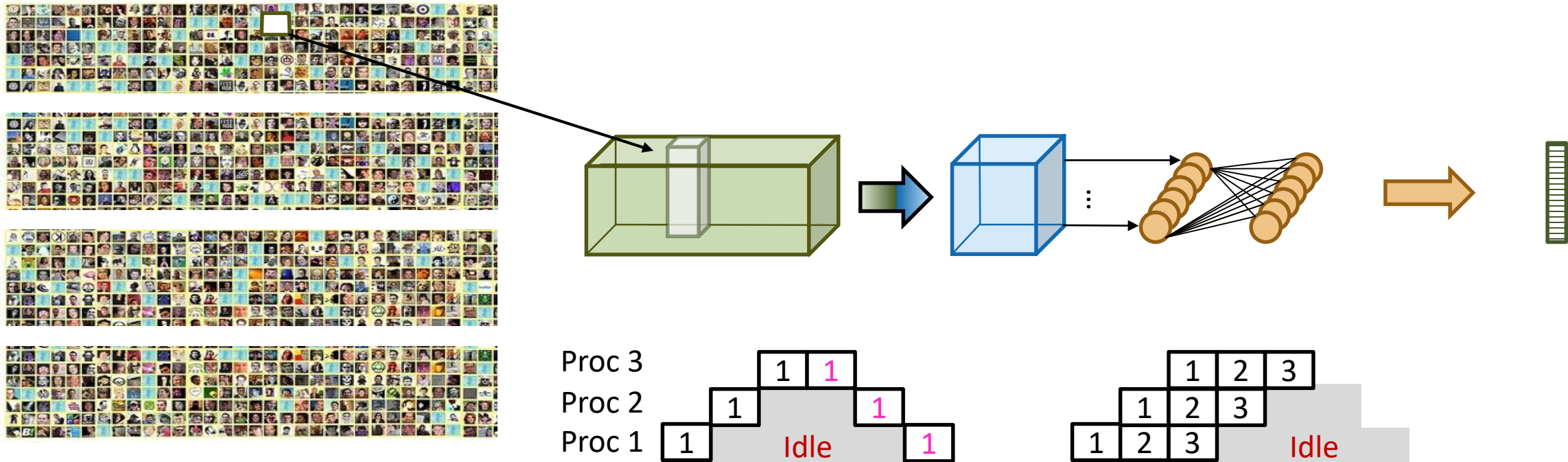
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



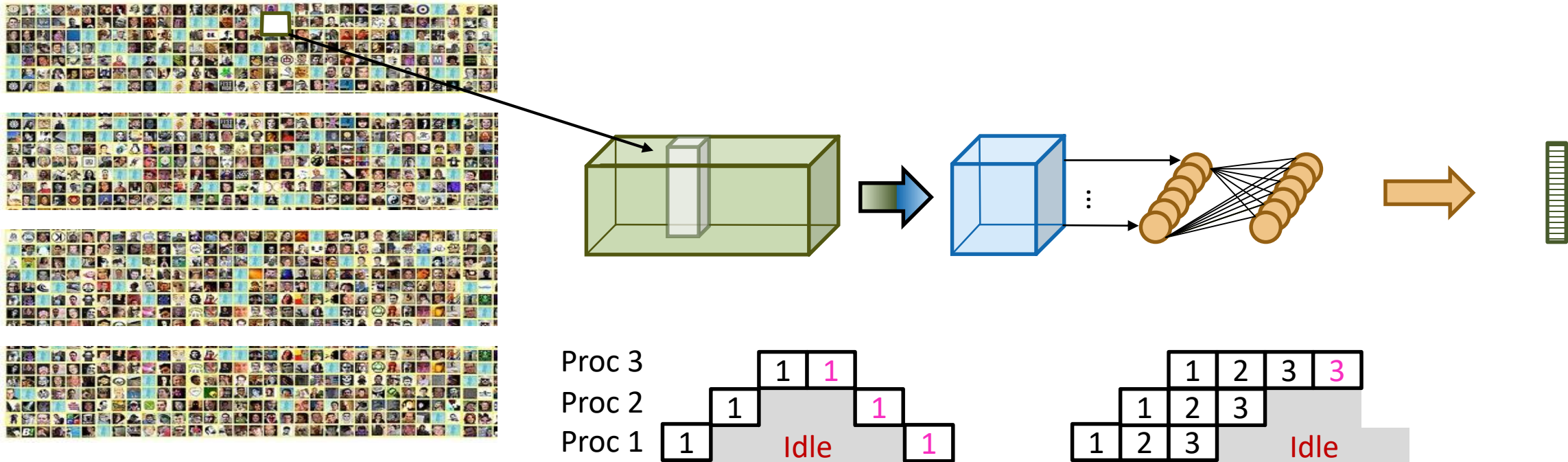
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



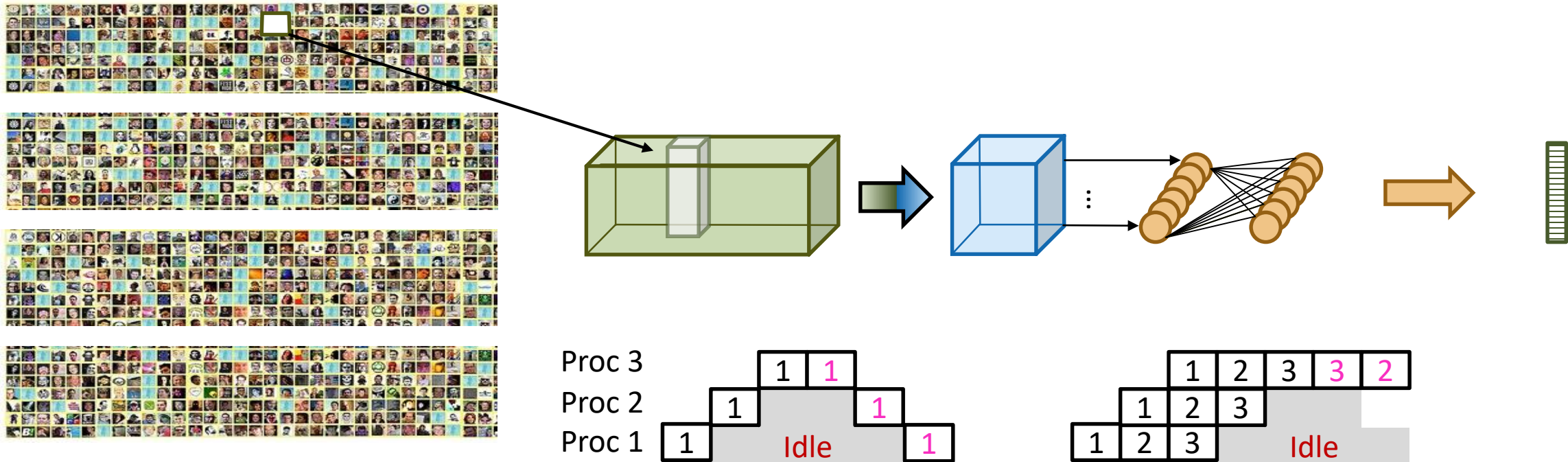
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



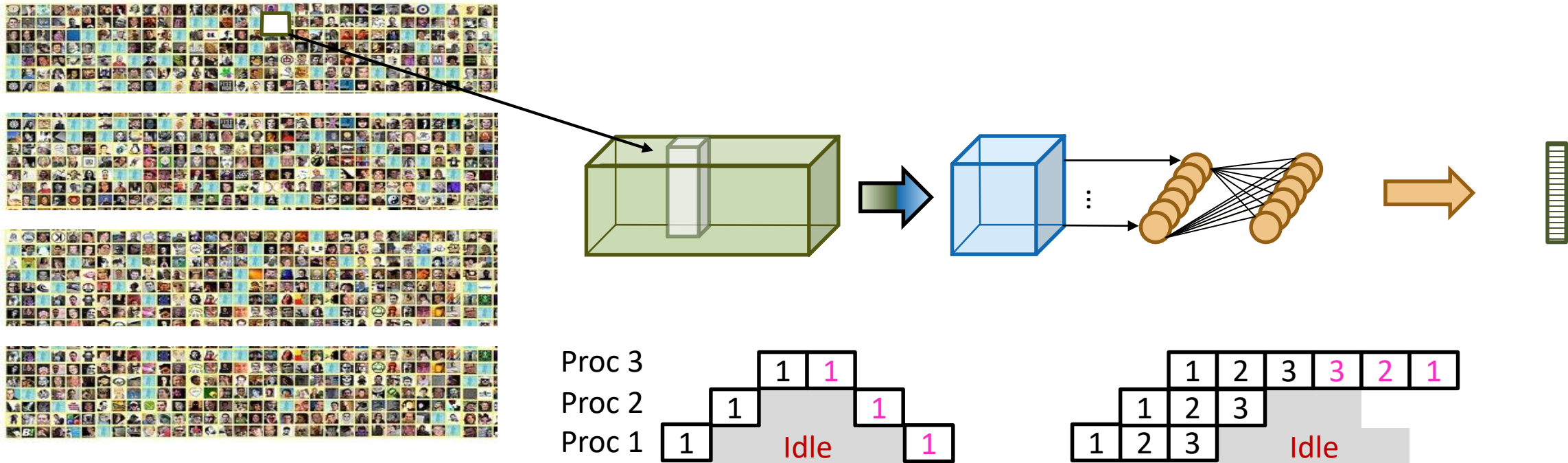
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



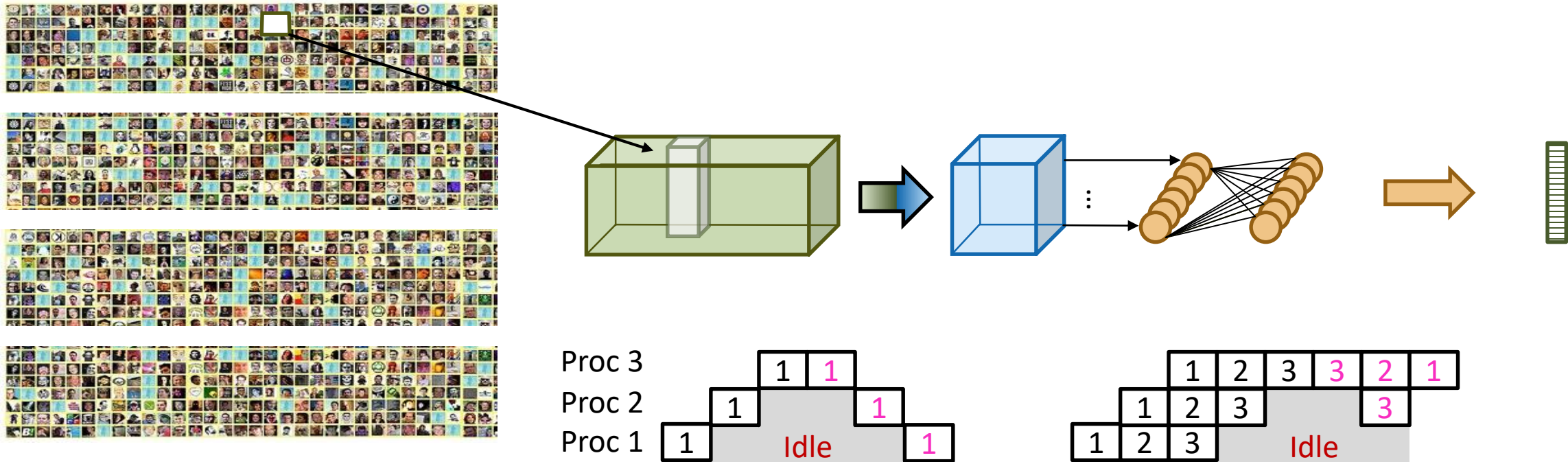
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



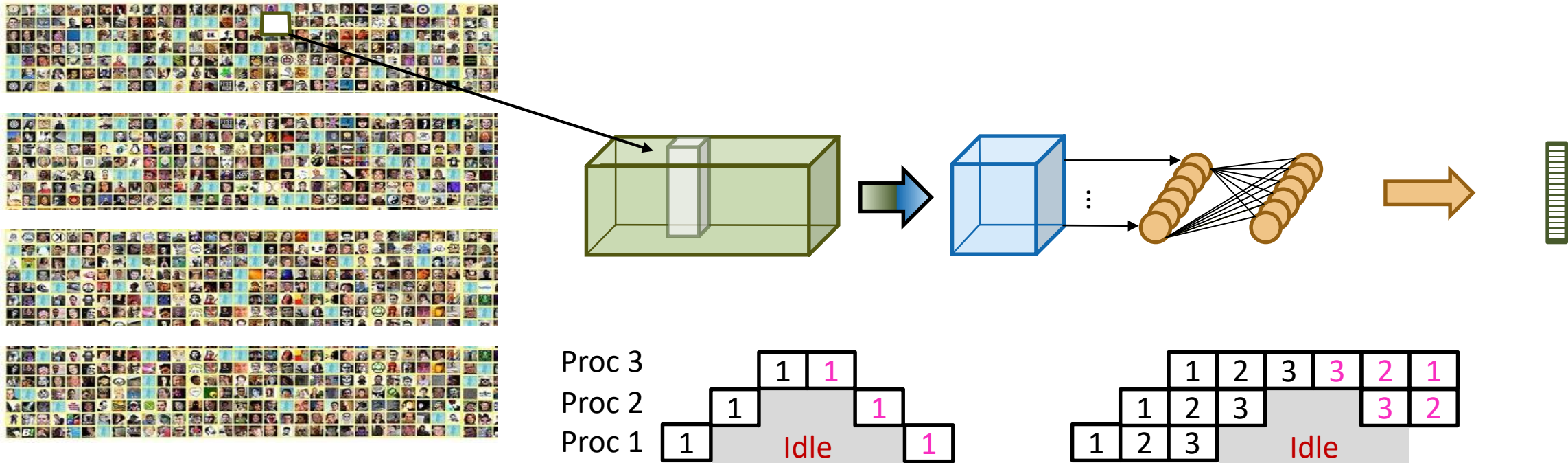
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



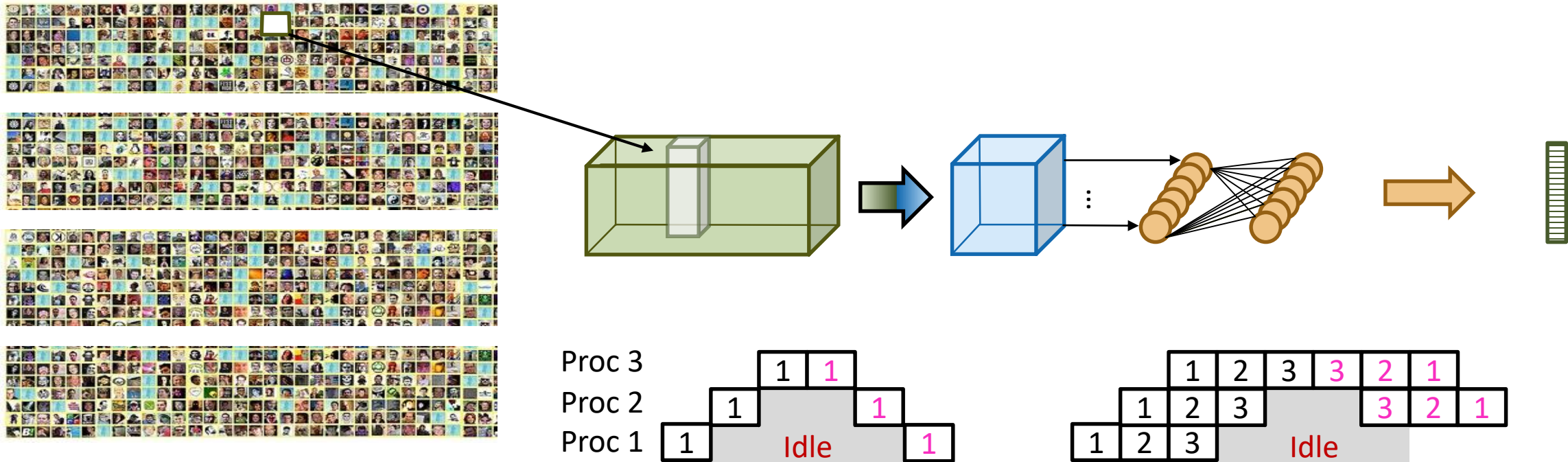
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



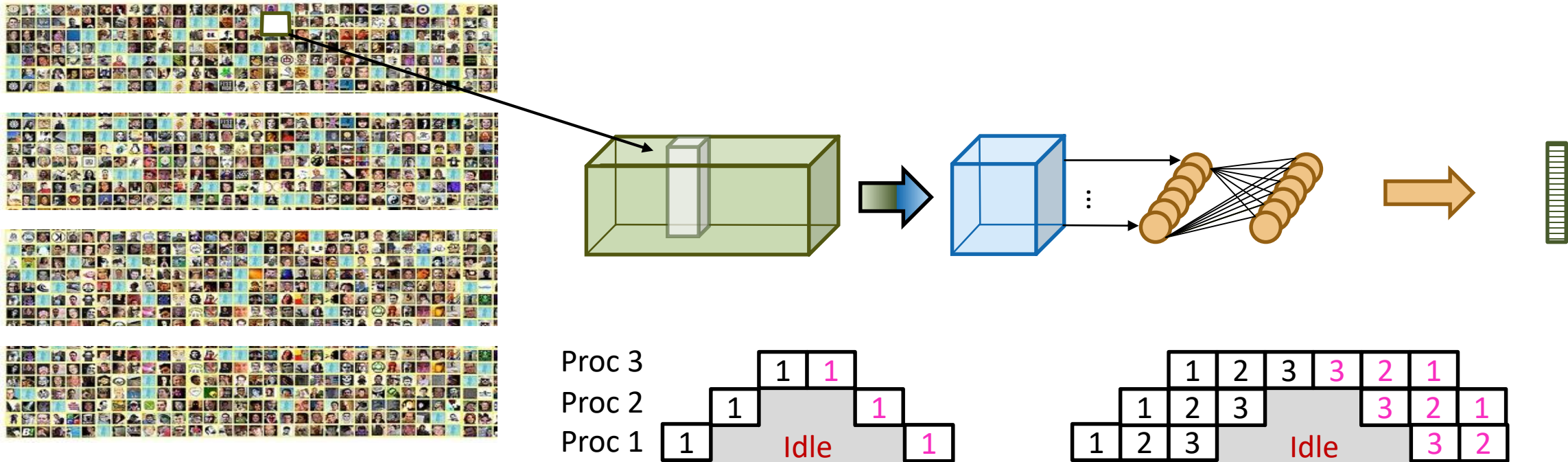
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



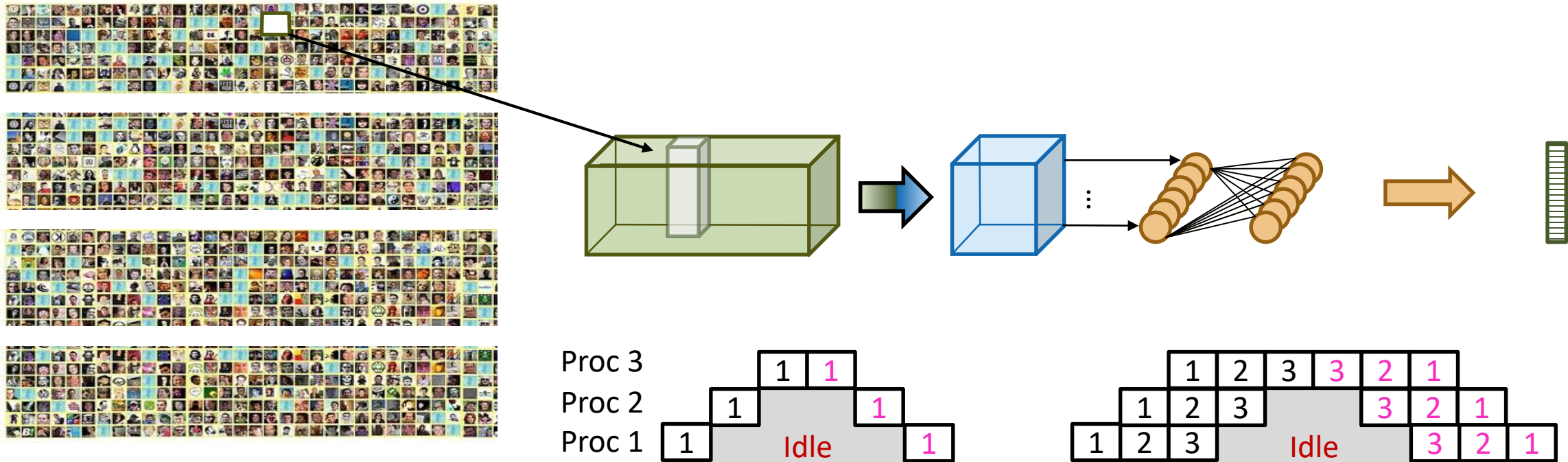
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



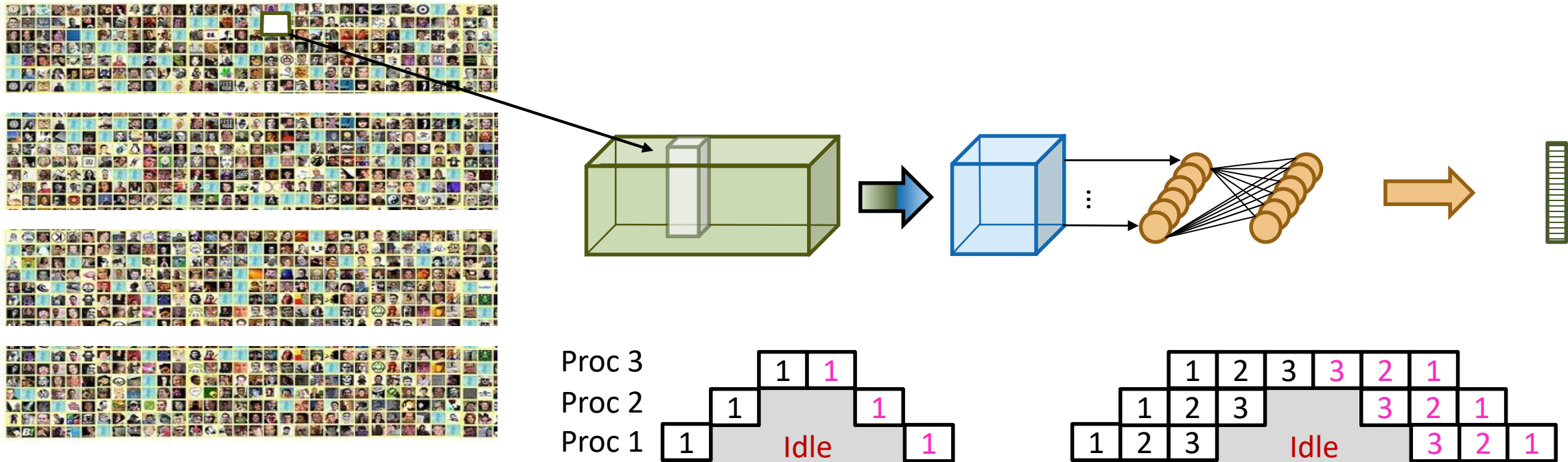
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



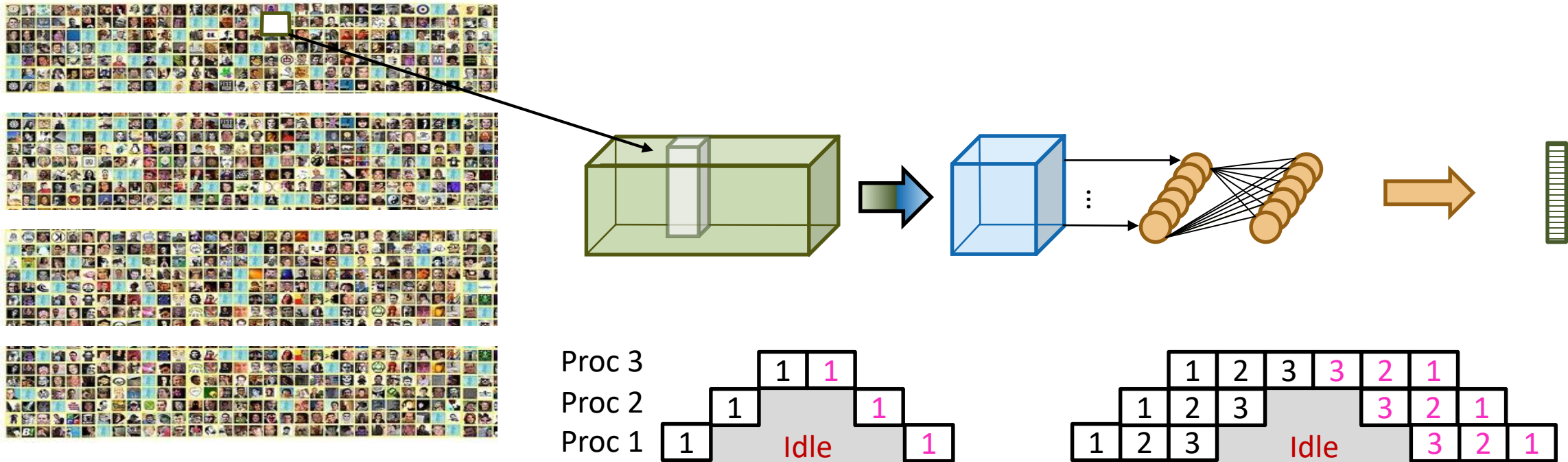
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



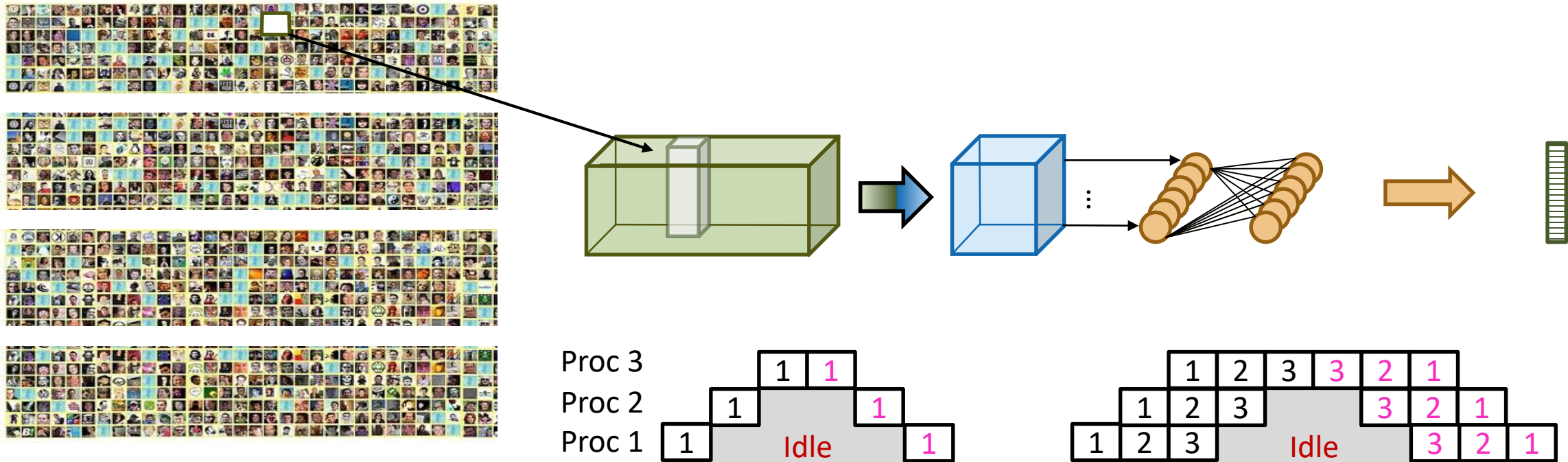
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



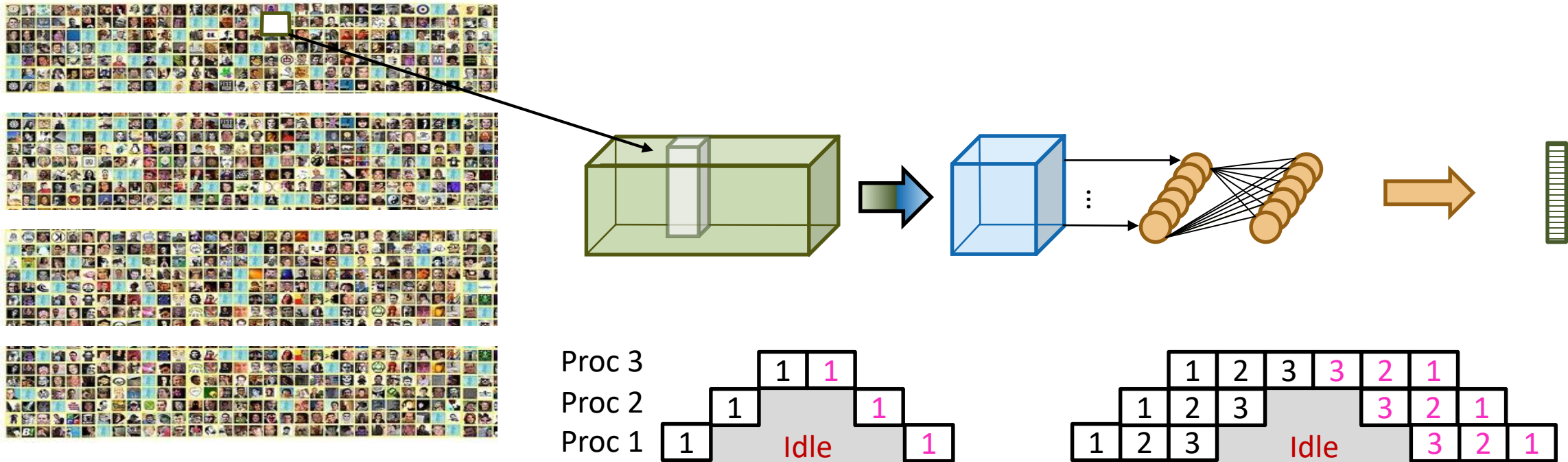
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



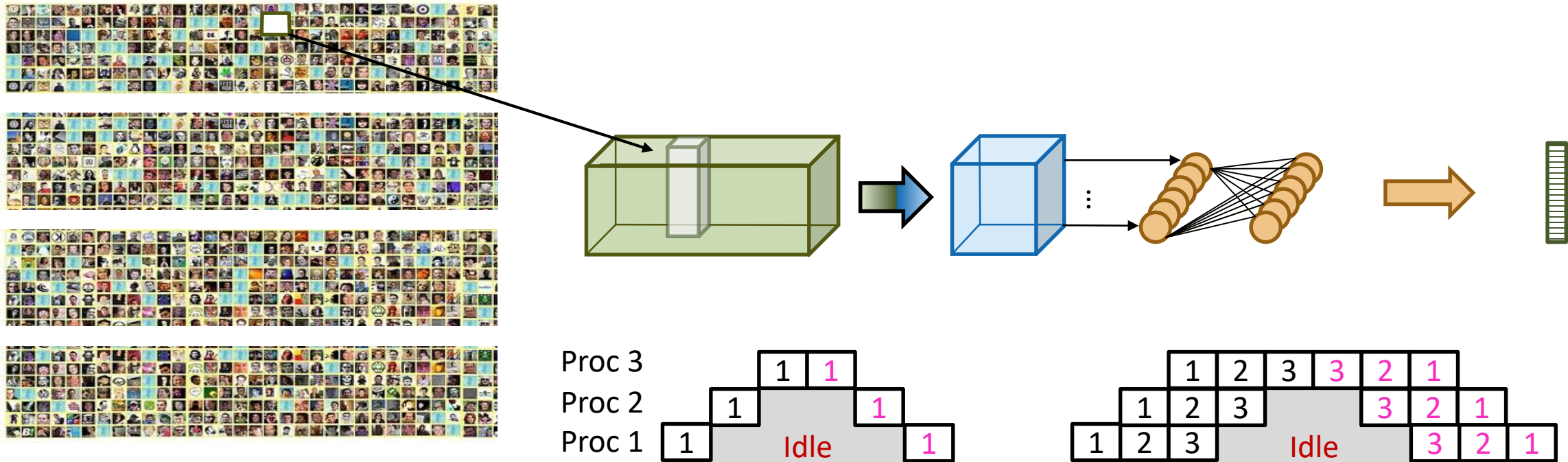
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



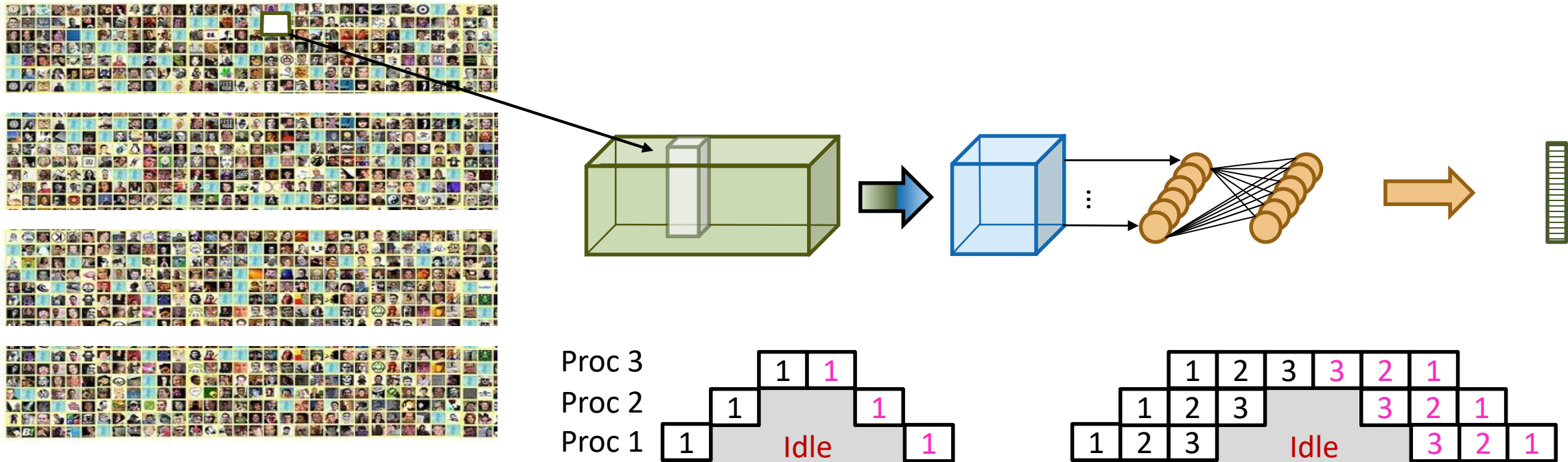
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



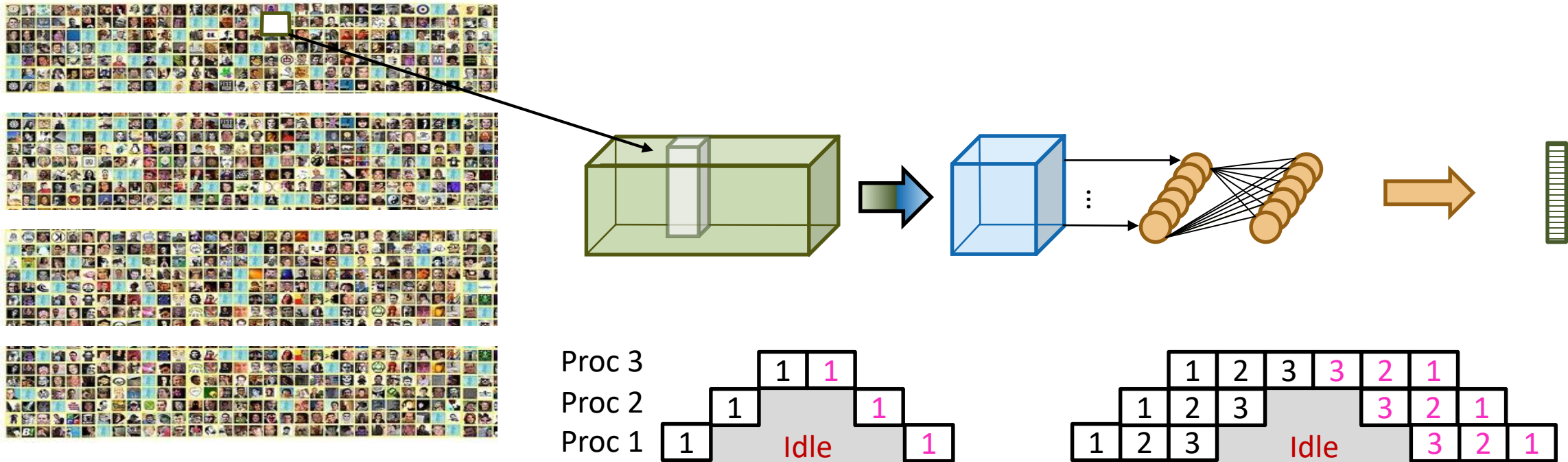
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



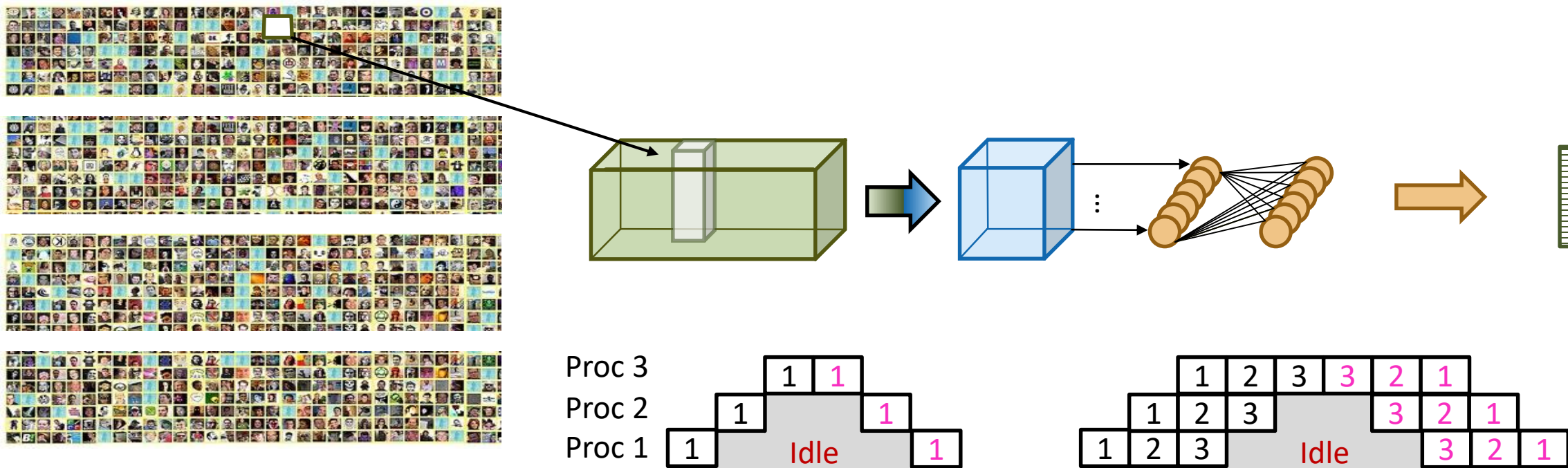
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



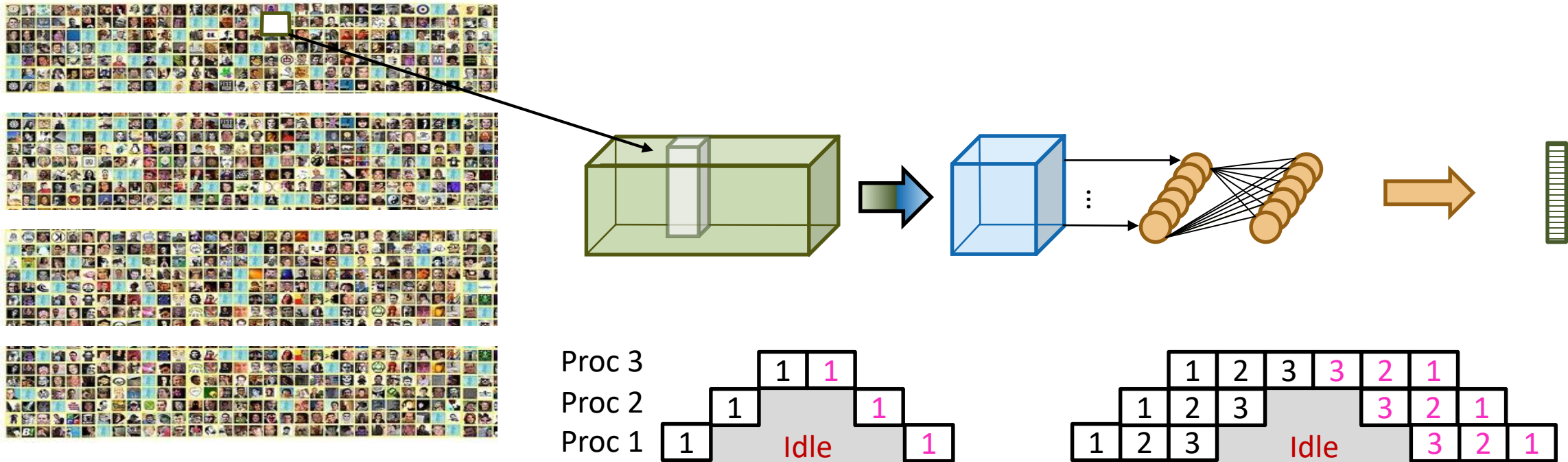
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



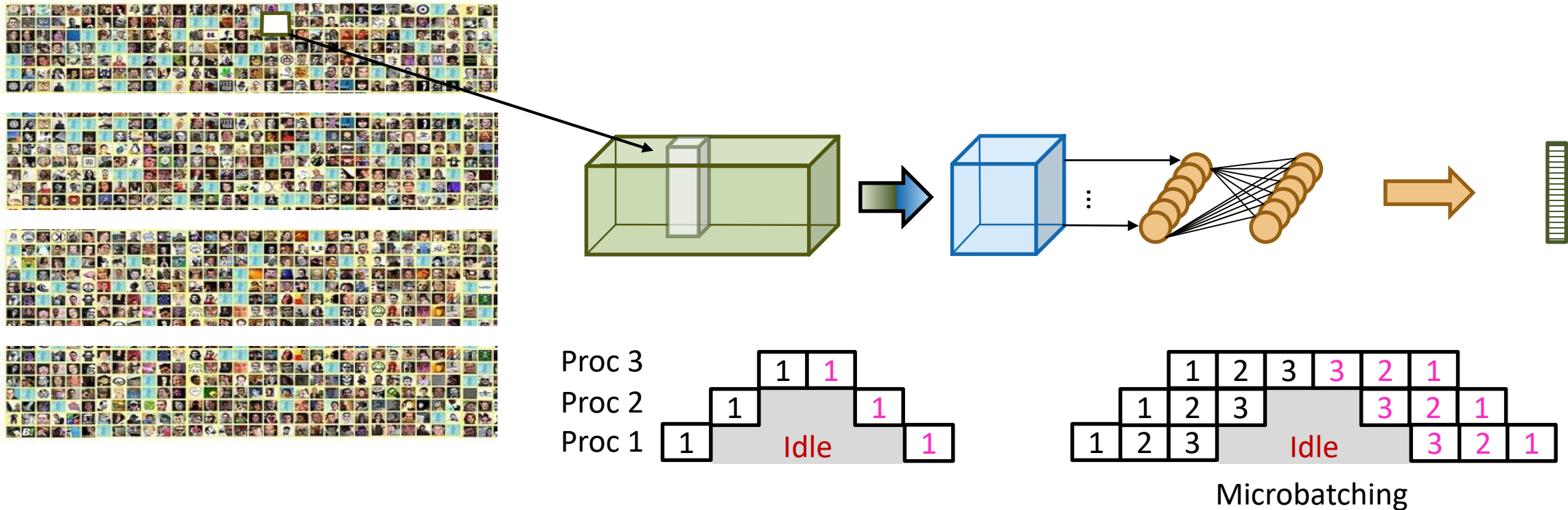
- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Pipeline parallelism – limited by network size



- Layers/parameters can be distributed across processors
- Sparse communication pattern (only pipeline stages)
- Mini-batch has to be copied through all processors**
- Consistent model introduces idle-time “bubble”**

Data parallelism – limited by batch-size

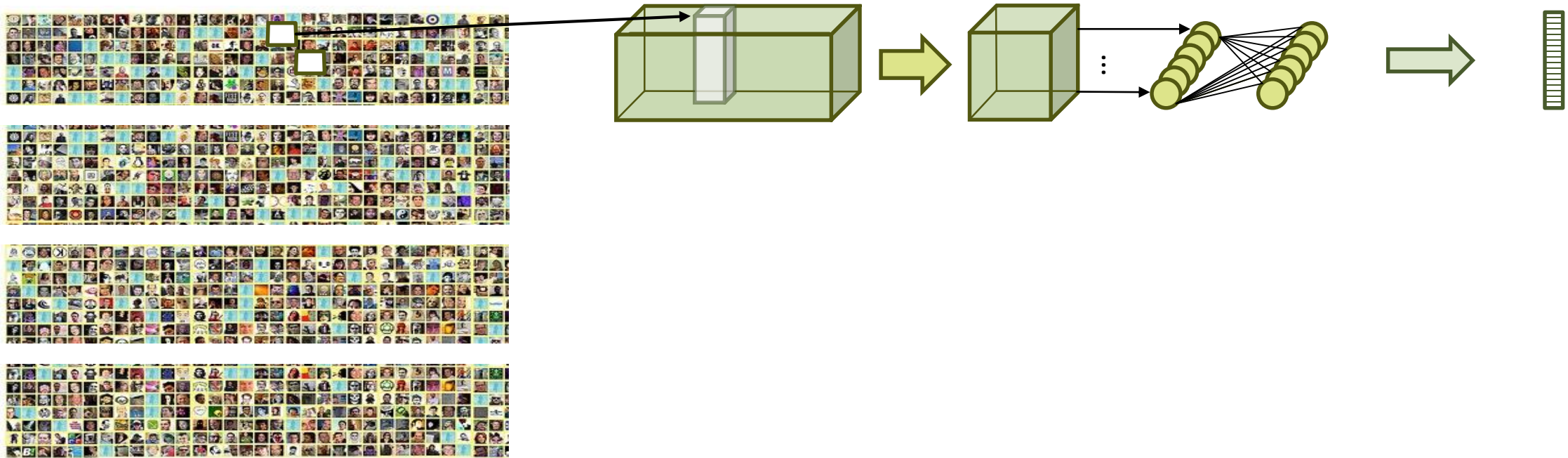


Data parallelism – limited by batch-size



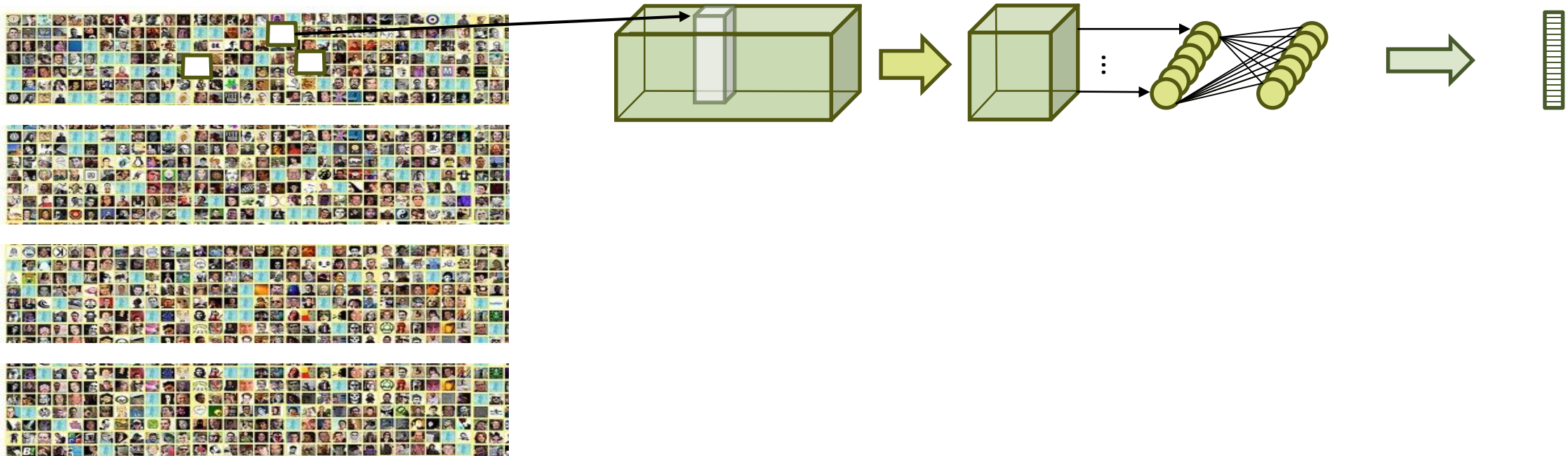
- Simple and efficient solution, easy to implement
- Duplicate parameters at all processors
- Affects generalization

Data parallelism – limited by batch-size



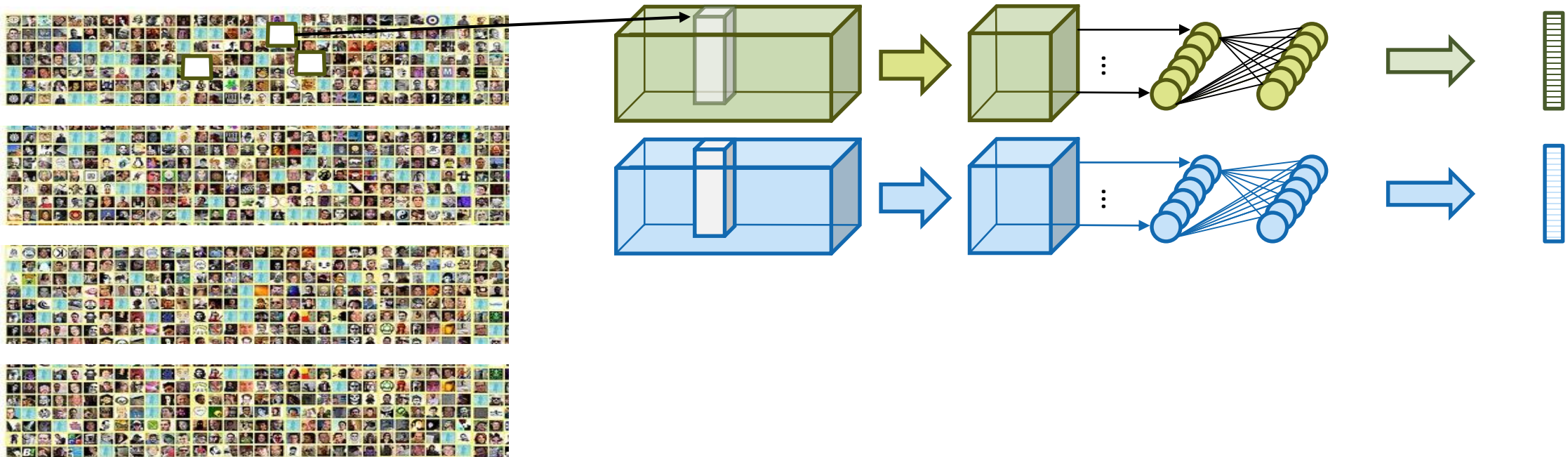
- Simple and efficient solution, easy to implement
- Duplicate parameters at all processors
- Affects generalization

Data parallelism – limited by batch-size



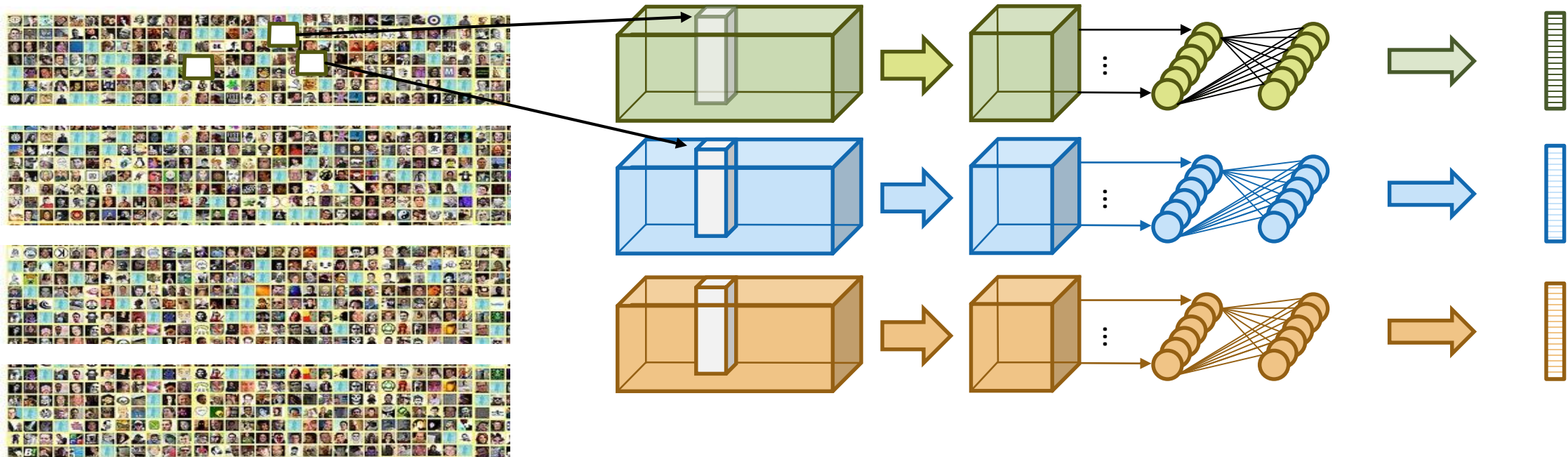
- Simple and efficient solution, easy to implement
- Duplicate parameters at all processors
- Affects generalization

Data parallelism – limited by batch-size



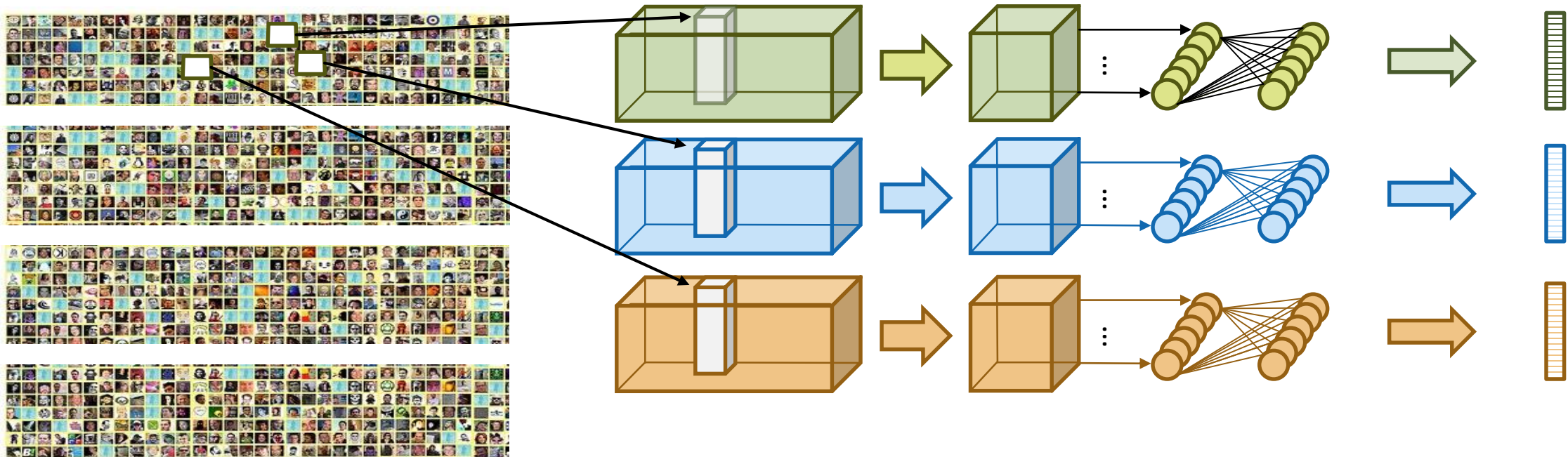
- Simple and efficient solution, easy to implement
- Duplicate parameters at all processors
- Affects generalization

Data parallelism – limited by batch-size



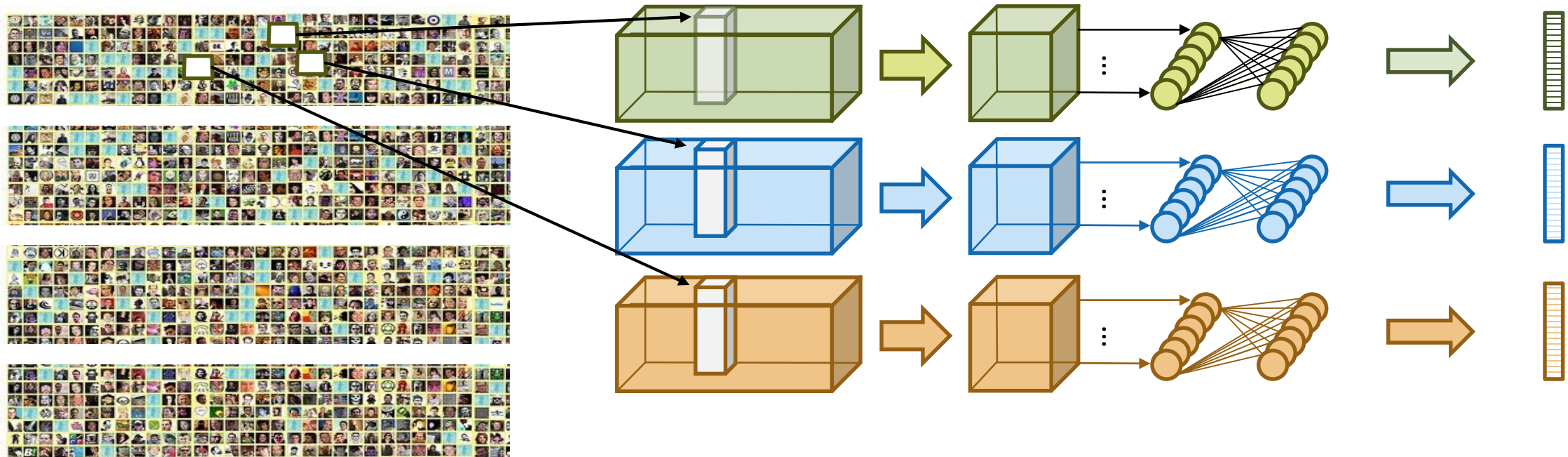
- Simple and efficient solution, easy to implement
- Duplicate parameters at all processors
- Affects generalization

Data parallelism – limited by batch-size



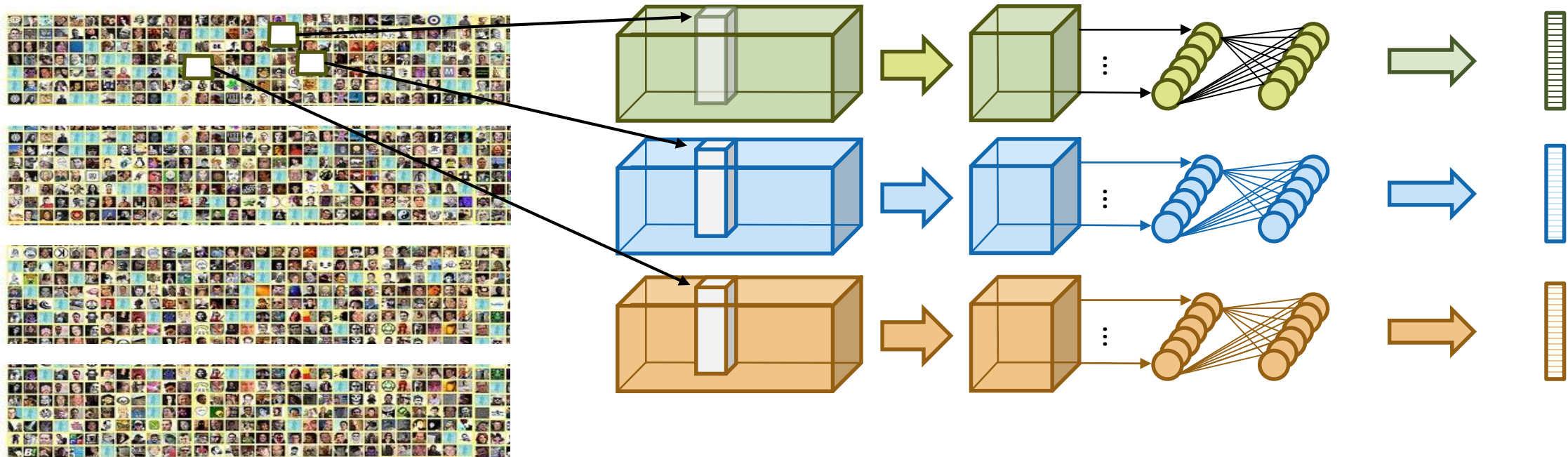
- Simple and efficient solution, easy to implement
- Duplicate parameters at all processors
- Affects generalization

Data parallelism – limited by batch-size



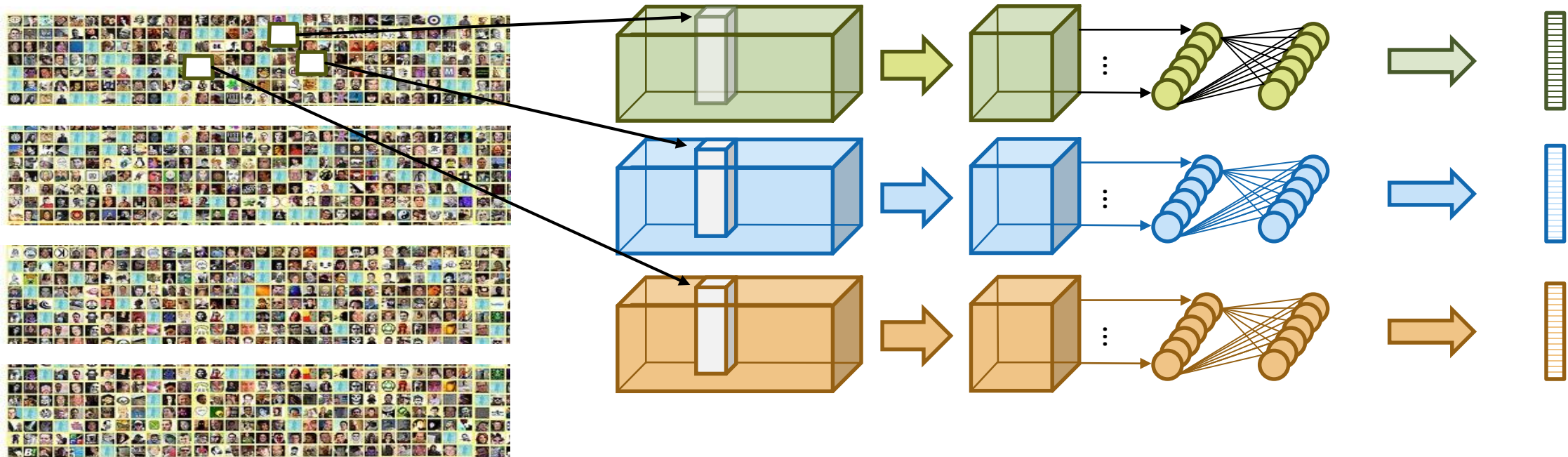
- Simple and efficient solution, easy to implement
- Duplicate parameters at all processors
- Affects generalization

Data parallelism – limited by batch-size



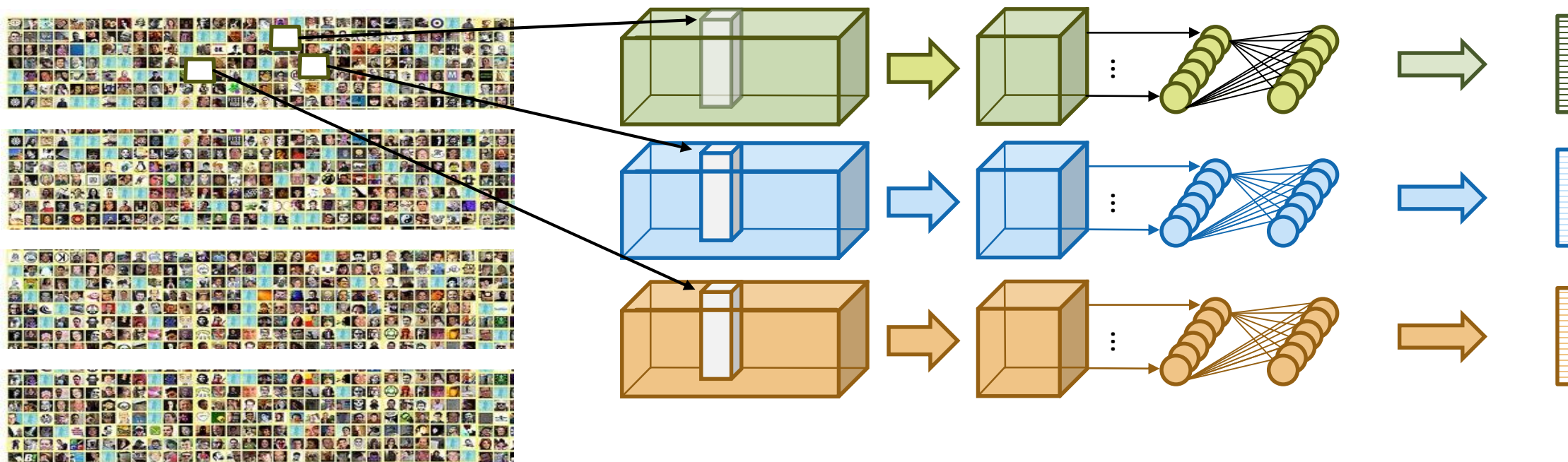
- Simple and efficient solution, easy to implement
- Duplicate parameters at all processors
- Affects generalization

Data parallelism – limited by batch-size



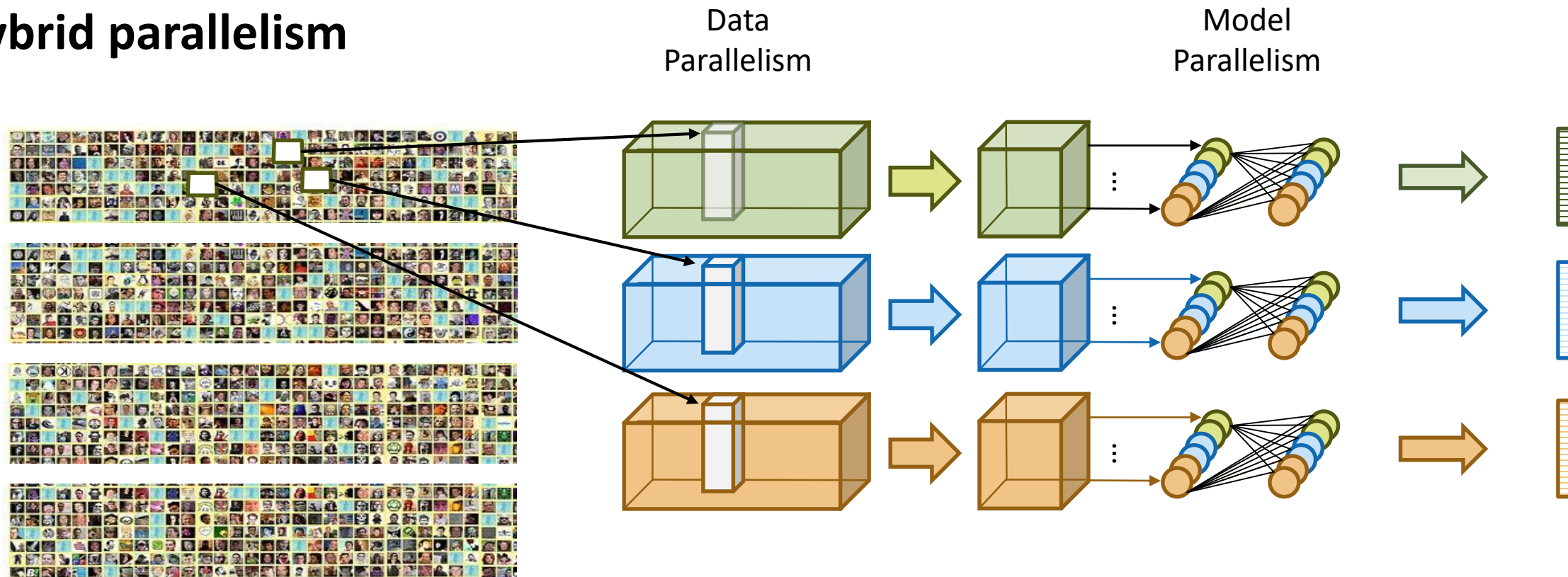
- Simple and efficient solution, easy to implement
- Duplicate parameters at all processors
- Affects generalization

Data parallelism – limited by batch-size



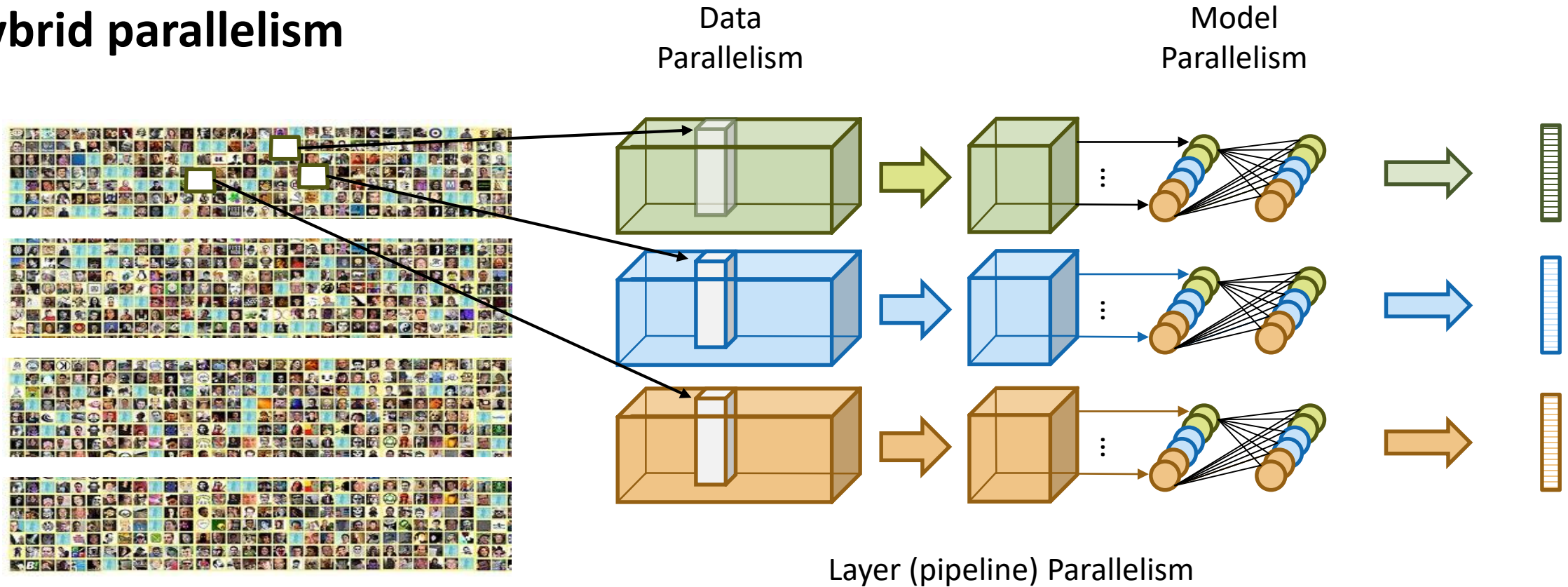
- Simple and efficient solution, easy to implement
- **Duplicate parameters at all processors**
- **Affects generalization**

Hybrid parallelism



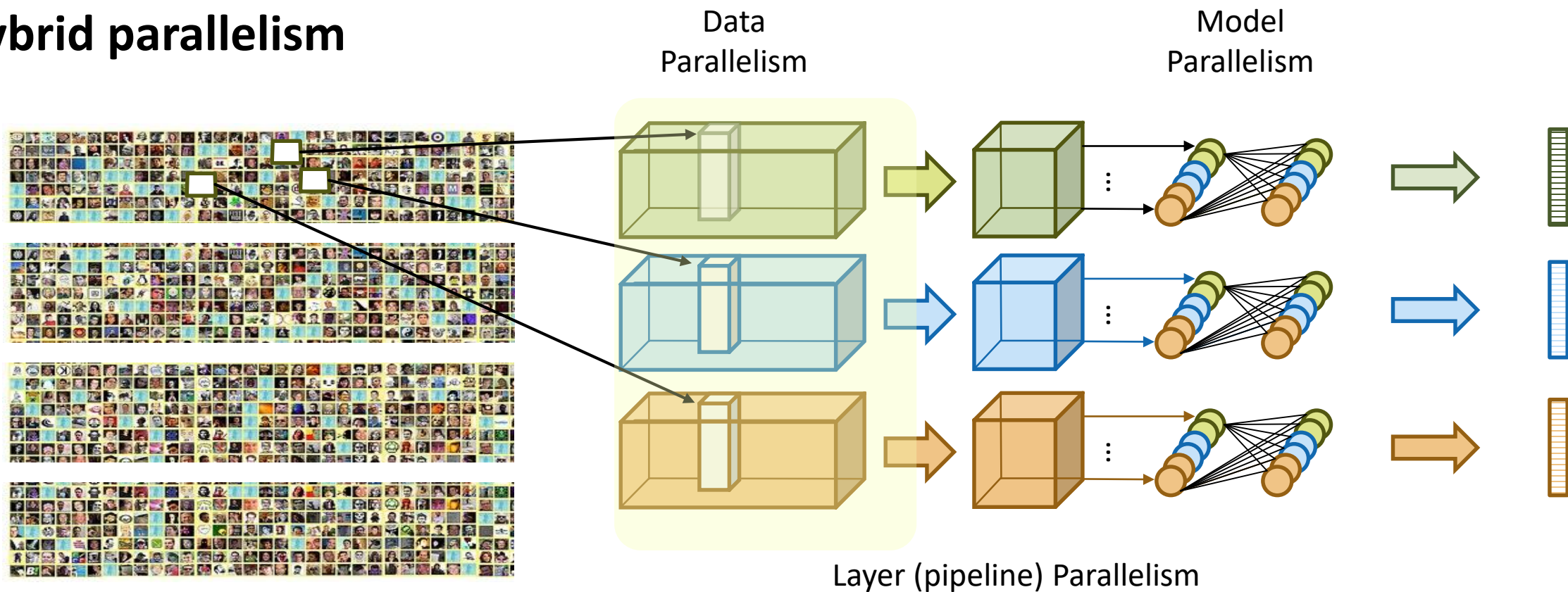
- Layers/parameters can be distributed across processors
- Can distribute minibatch
- Often specific to layer-types (e.g., distribute fc layers but handle conv layers data-parallel)
 - Enables arbitrary combinations of data, model, and pipeline parallelism – very powerful!

Hybrid parallelism



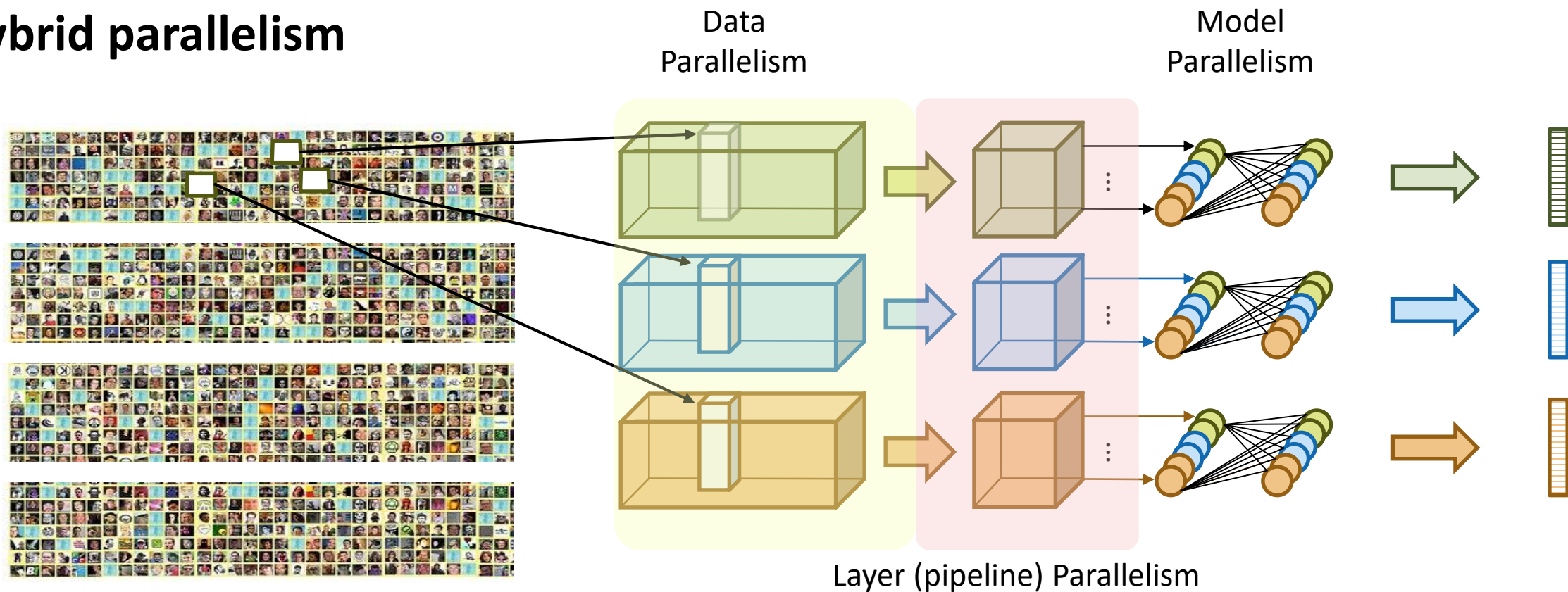
- Layers/parameters can be distributed across processors
- Can distribute minibatch
- Often specific to layer-types (e.g., distribute fc layers but handle conv layers data-parallel)
 - Enables arbitrary combinations of data, model, and pipeline parallelism – very powerful!

Hybrid parallelism



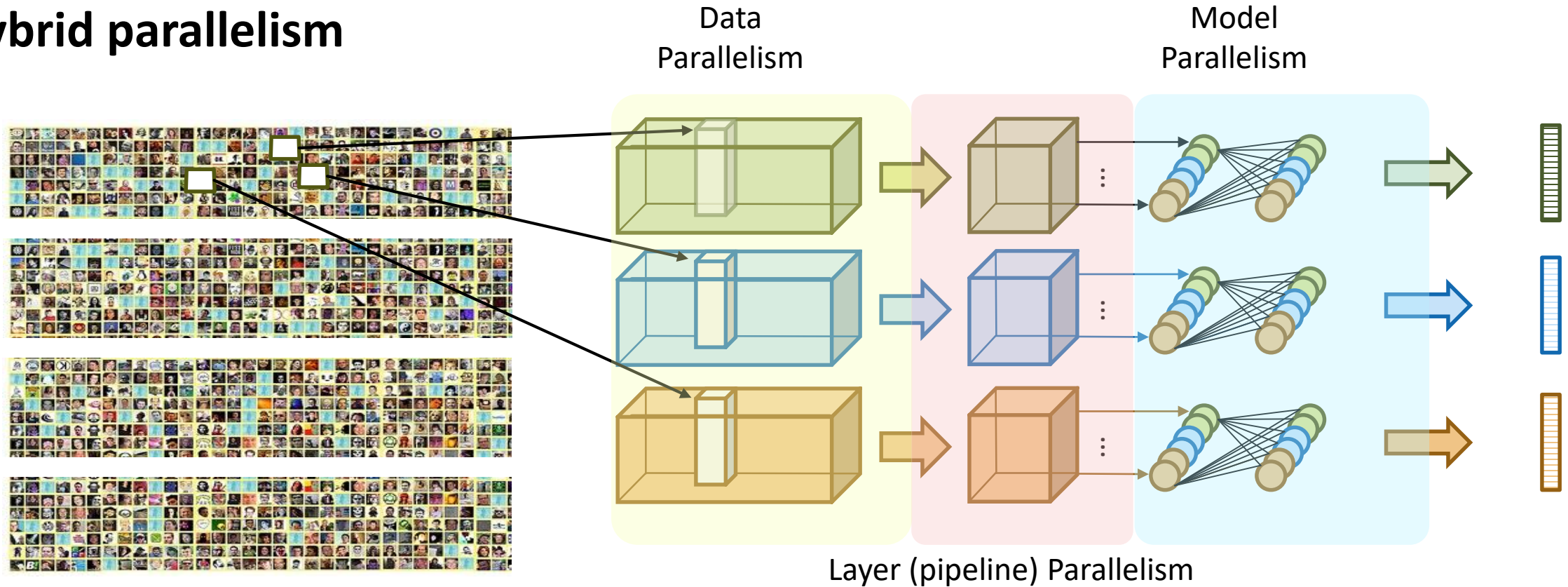
- Layers/parameters can be distributed across processors
- Can distribute minibatch
- Often specific to layer-types (e.g., distribute fc layers but handle conv layers data-parallel)
 - Enables arbitrary combinations of data, model, and pipeline parallelism – very powerful!

Hybrid parallelism



- Layers/parameters can be distributed across processors
- Can distribute minibatch
- Often specific to layer-types (e.g., distribute fc layers but handle conv layers data-parallel)
 - Enables arbitrary combinations of data, model, and pipeline parallelism – very powerful!

Hybrid parallelism



- **Layers/parameters can be distributed across processors**
- **Can distribute minibatch**
- **Often specific to layer-types (e.g., distribute fc layers but handle conv layers data-parallel)**
 - Enables arbitrary combinations of data, model, and pipeline parallelism – very powerful!

Other ways to think about parallelism

Other ways to think about parallelism

- All definitions are fuzzy (including this 😊)
- Data-, model-, pipeline-, hybrid-parallelism
- **Weak vs strong scaling**
 - What do you keep the same vs what do change?
 - Mini-batch weak scaling: grow the mini-batch
 - Mini-batch model scaling: grow the model size (not so useful in general...)
 - Strong scaling: Fix everything, use more GPUs
- **For convolution: based on partitioned tensor dimensions**
 - Sample-, spatial-, channel-, filter-parallelism

Large mini-batches

- **Make the mini-batch really big!**

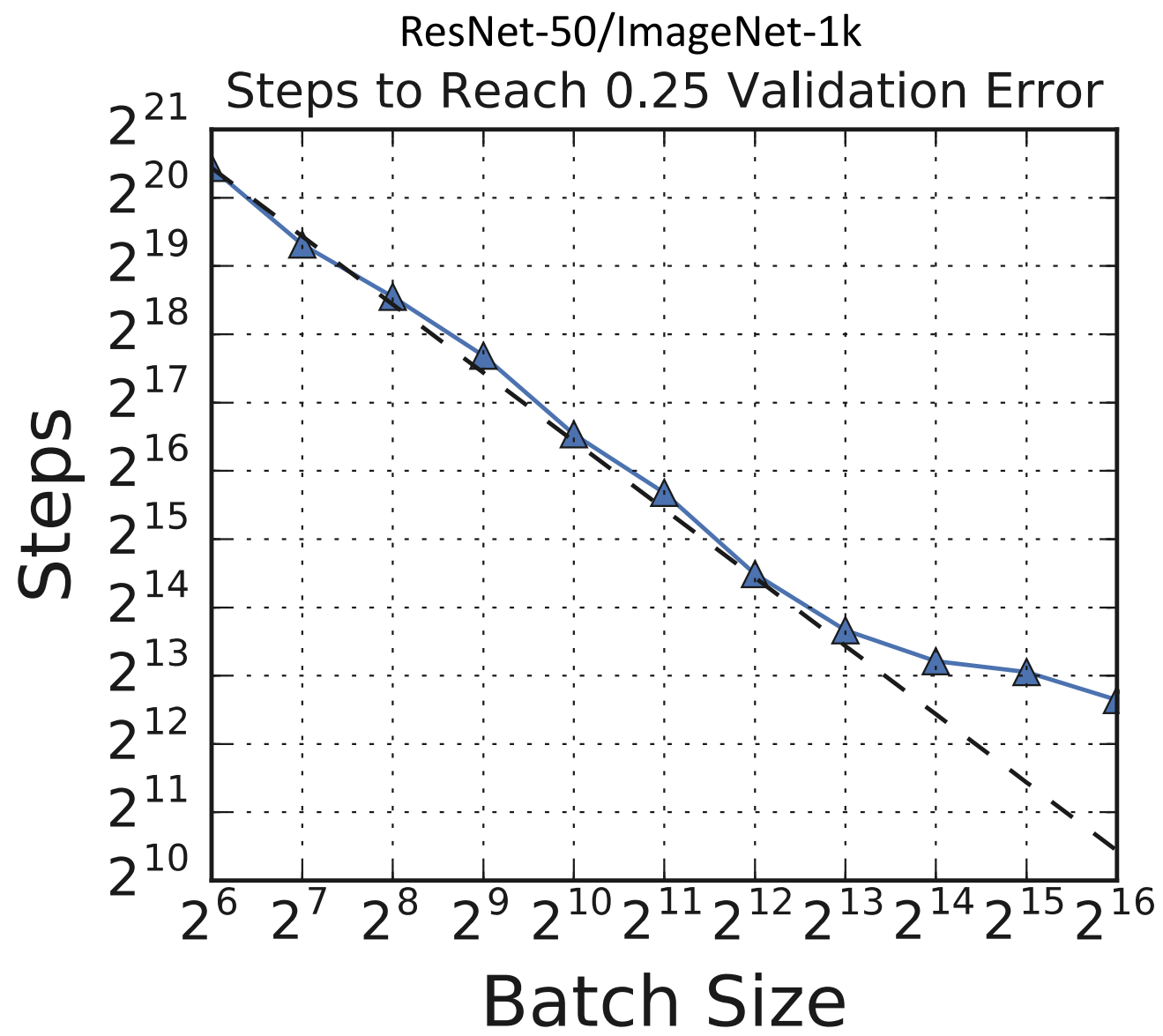
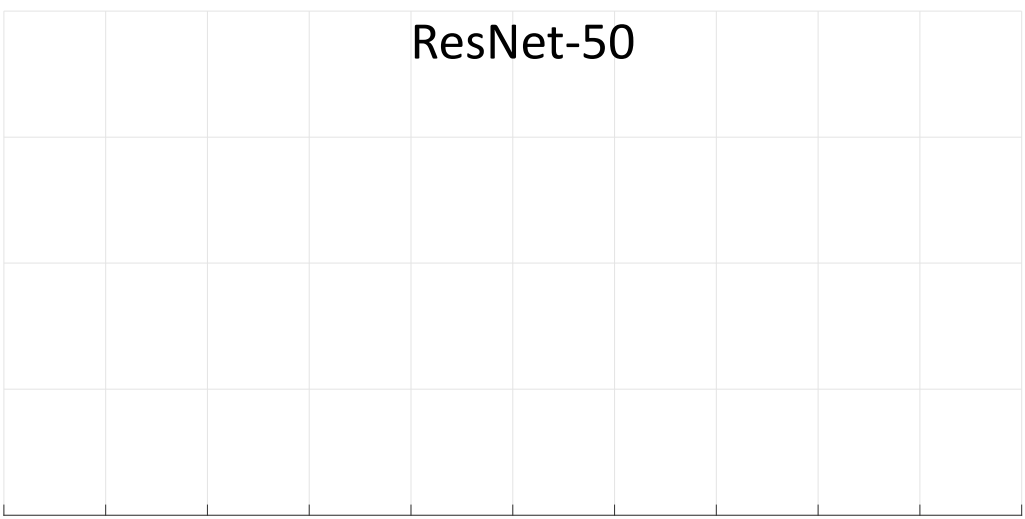
Large mini-batches

- Make the mini-batch really big!

			ResNet-50						

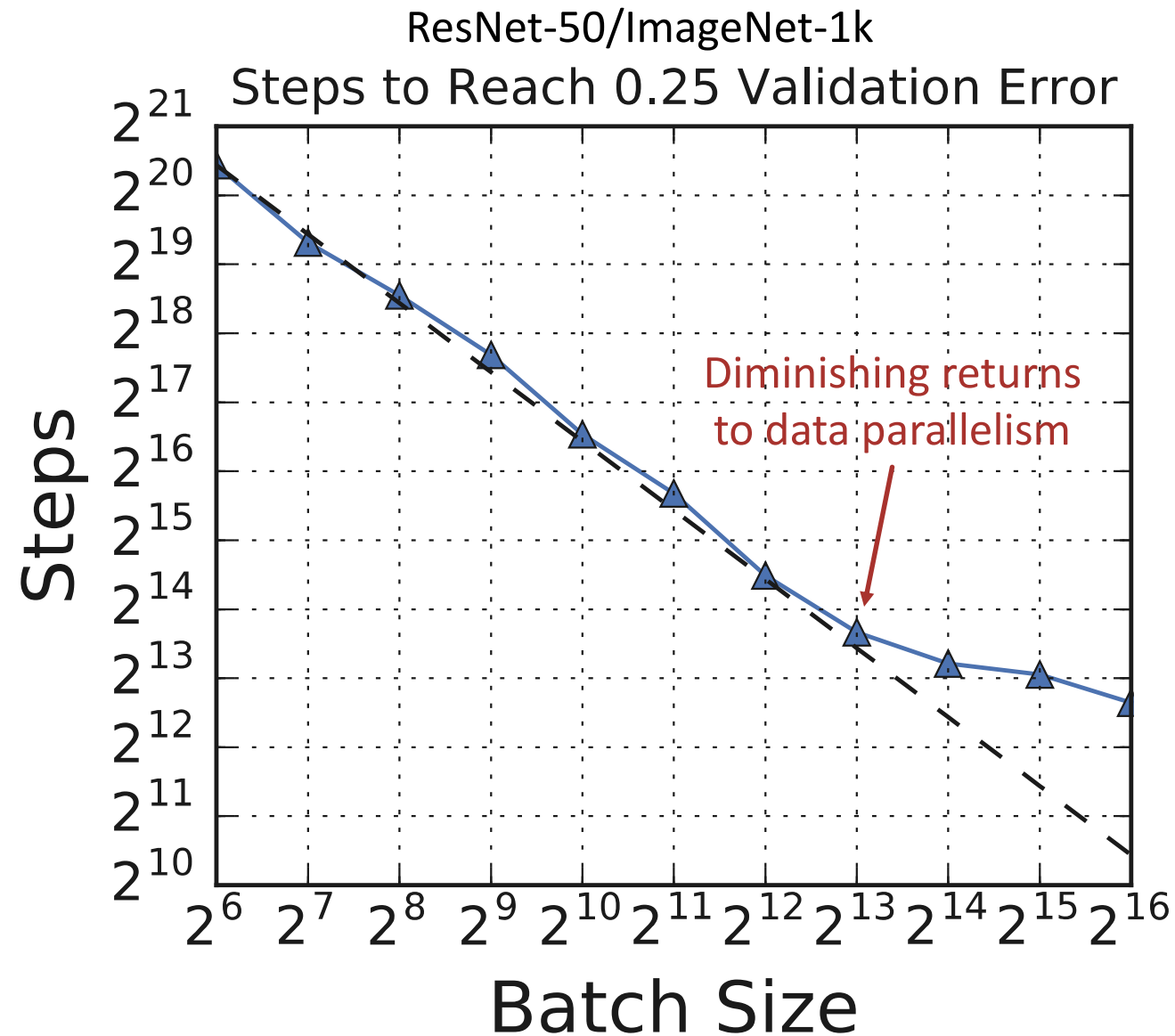
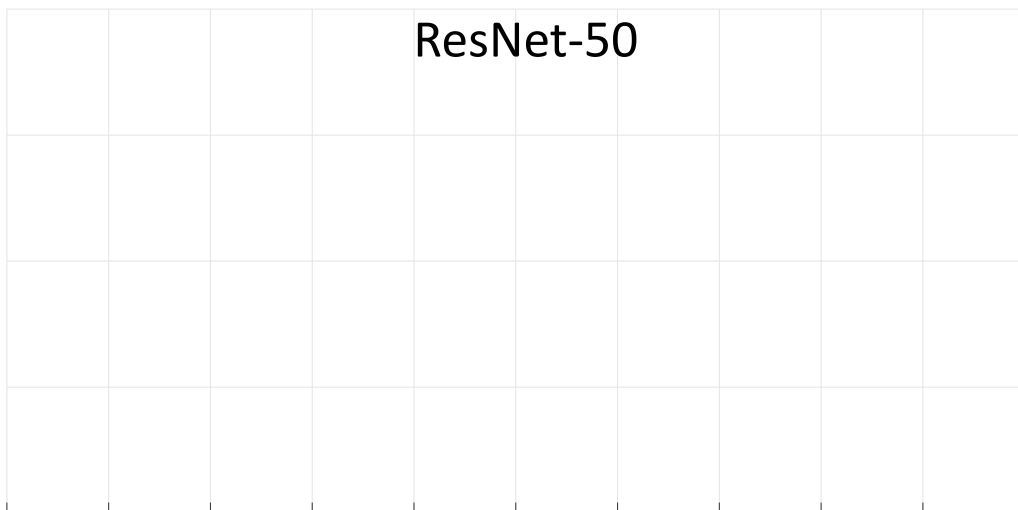
Large mini-batches

- Make the mini-batch really big!



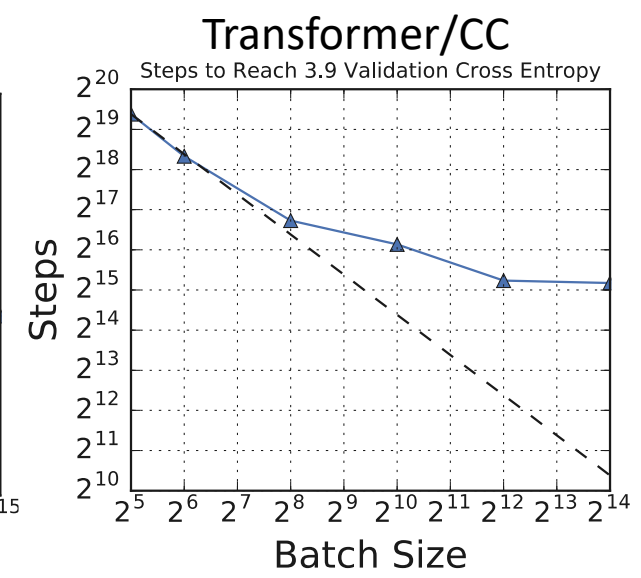
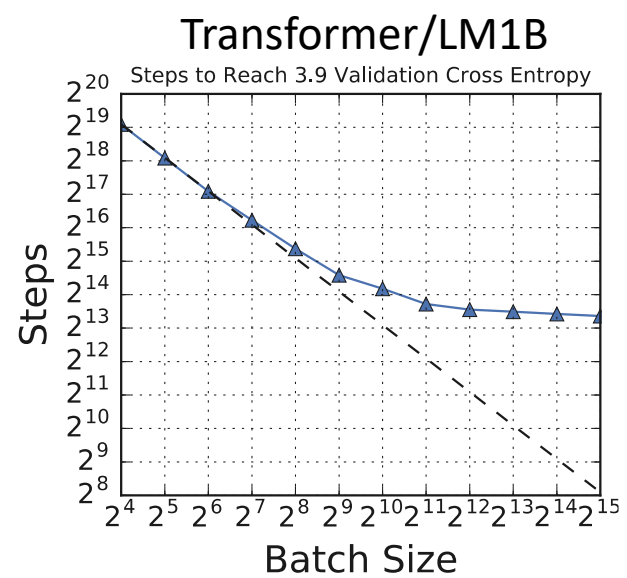
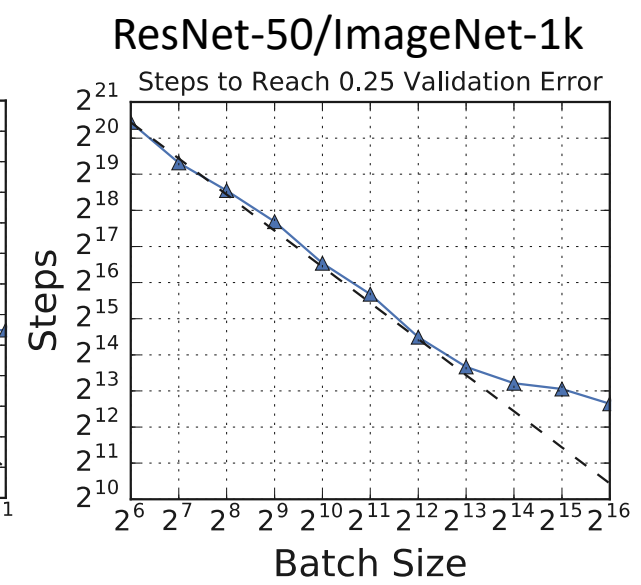
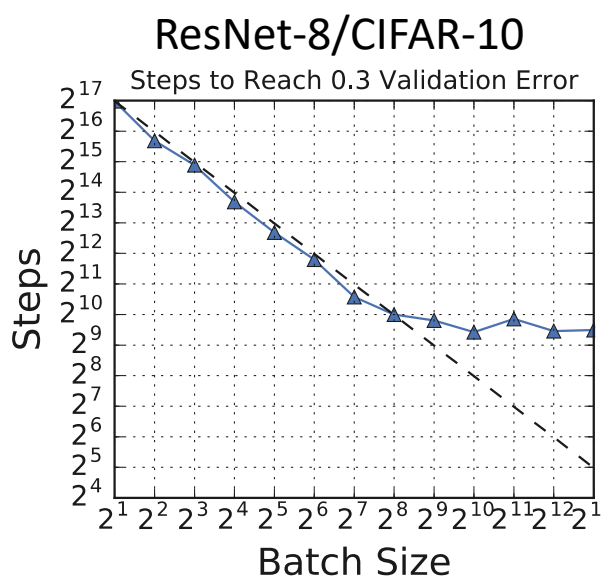
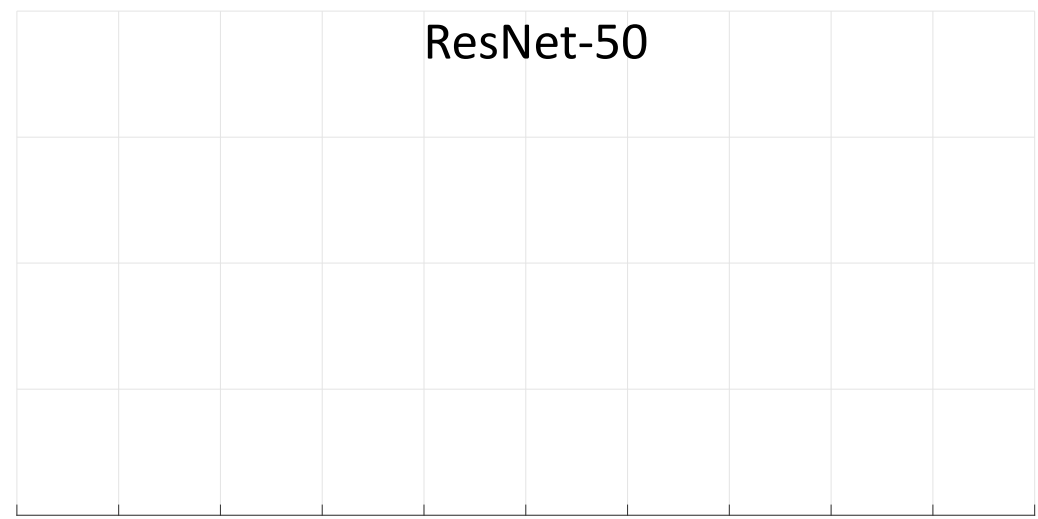
Large mini-batches

- Make the mini-batch really big!



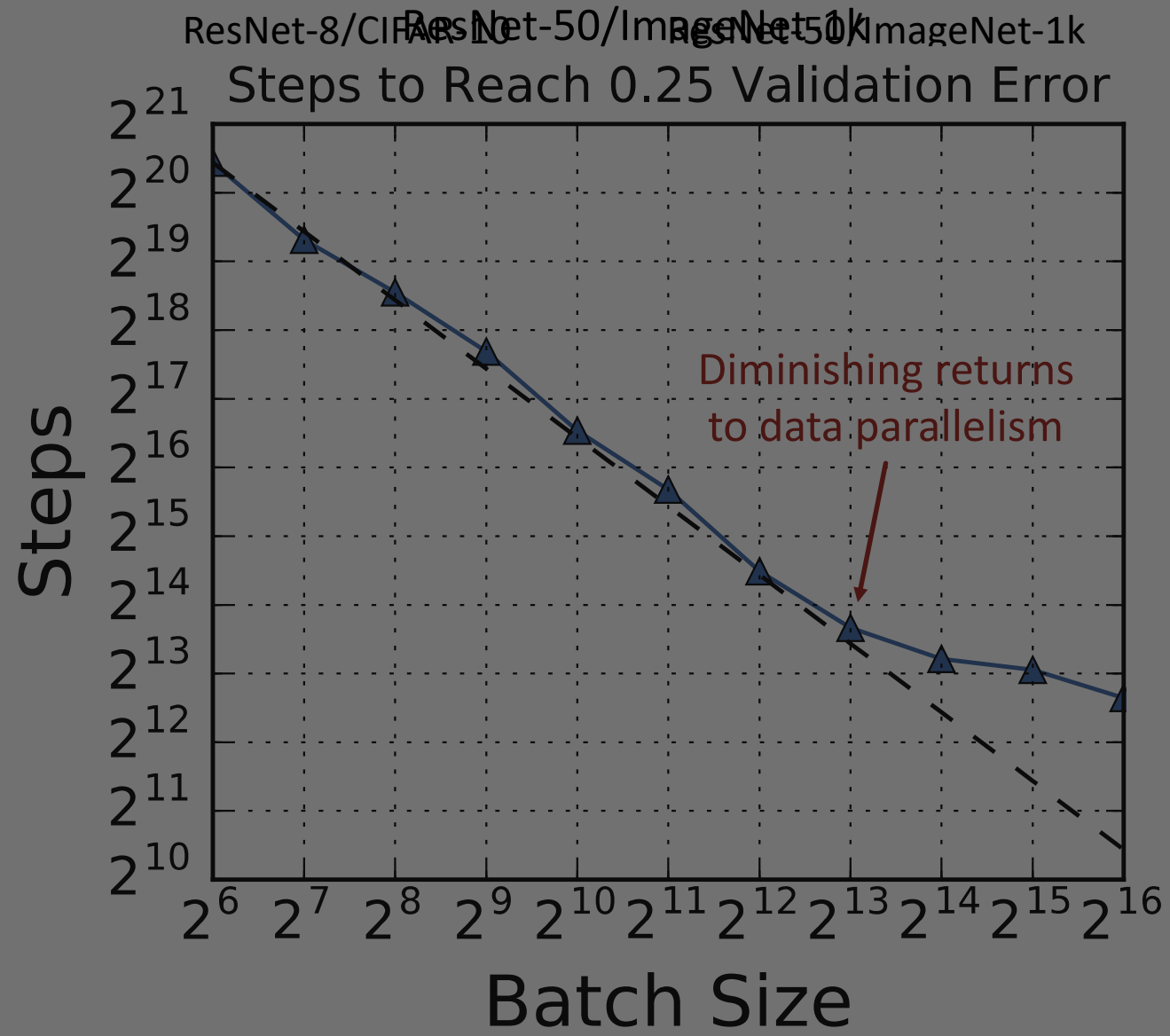
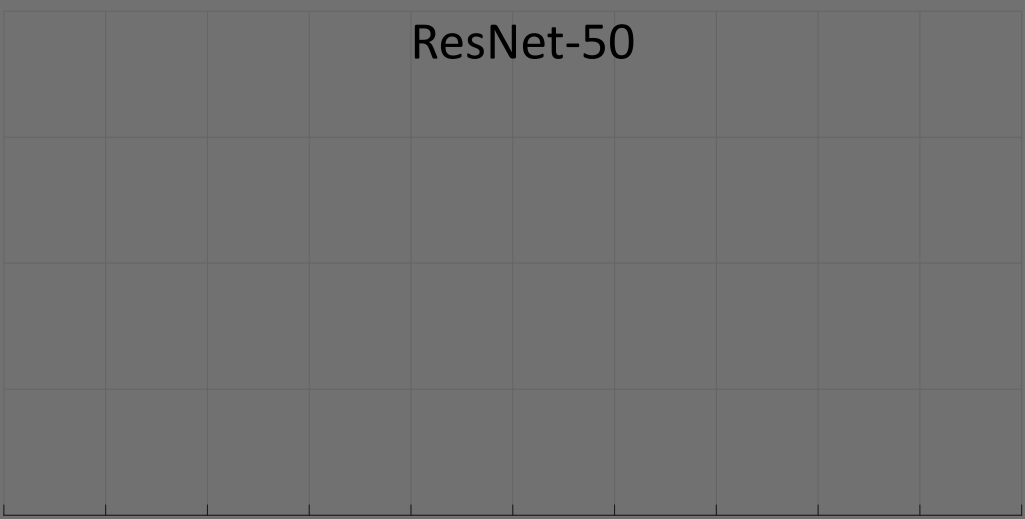
Large mini-batches

- Make the mini-batch really big!



Large mini-batches

- Make the mini-batch really big!



Large mini-batches

- Make the mini-batch really big!

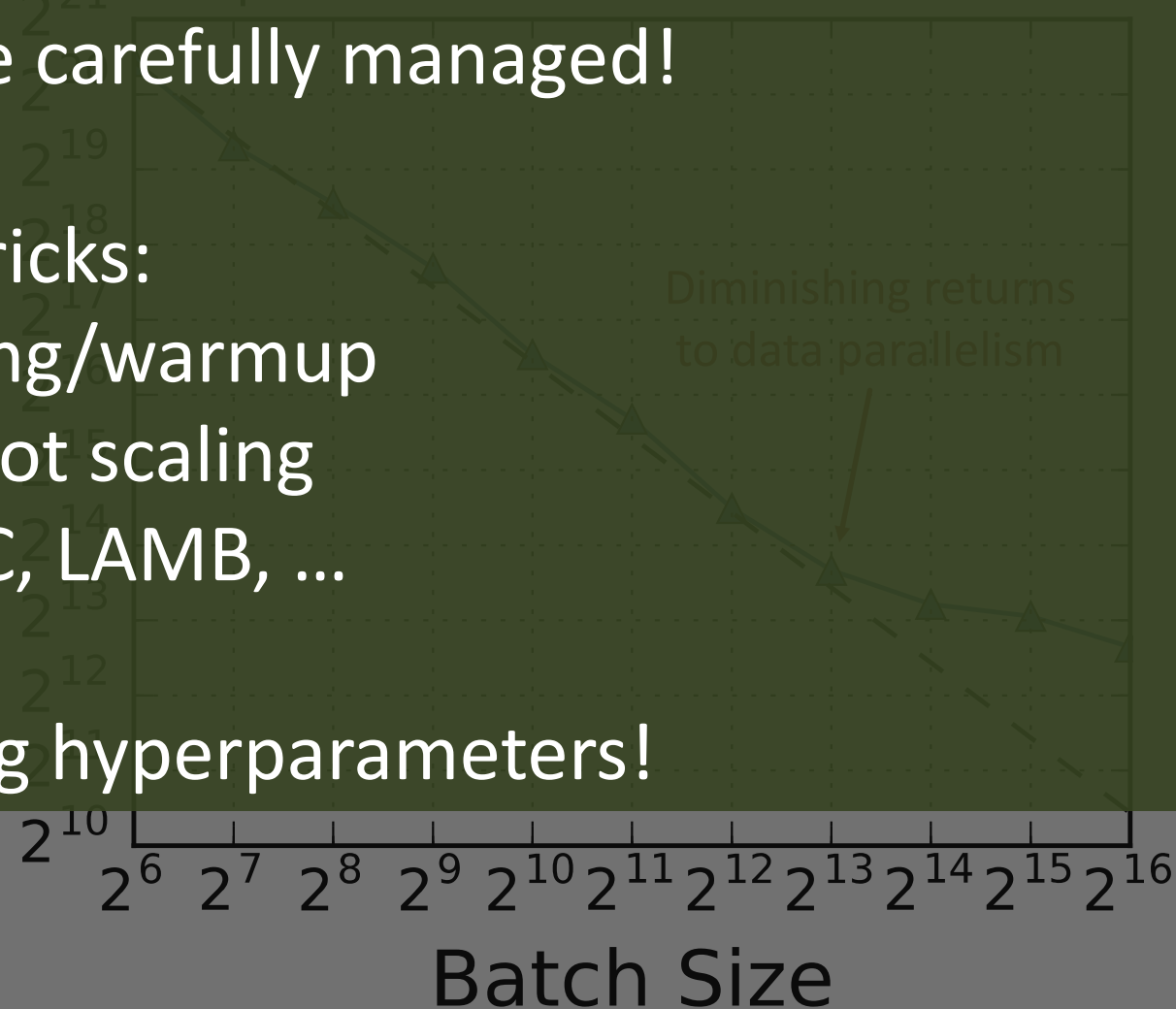
Mini-batch size must be carefully managed!

Some tricks:

- Linear scaling/warmup
- Square root scaling
- LARS, LARC, LAMB, ...

Often requires retuning hyperparameters!

ResNet-8/CIFAR-10, ResNet-50/ImageNet-50k, ResNet-101/ImageNet-1k
Steps to Reach 0.25 Validation Error



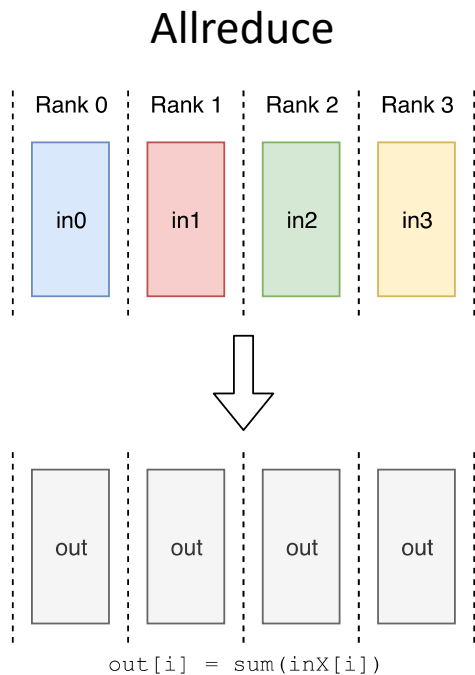
Collectives for deep learning

Collectives for deep learning

- **Certain communication patterns can be optimized**
- **People keep reinventing MPI**
 - Baidu Allreduce, NCCL, Gloo, Horovod, ...
- **What we need (for this talk):**

Collectives for deep learning

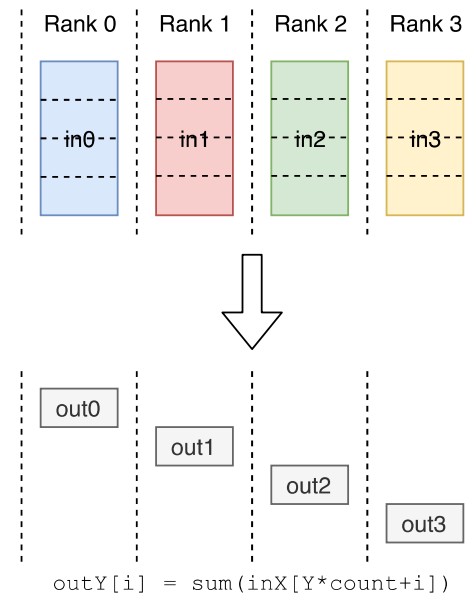
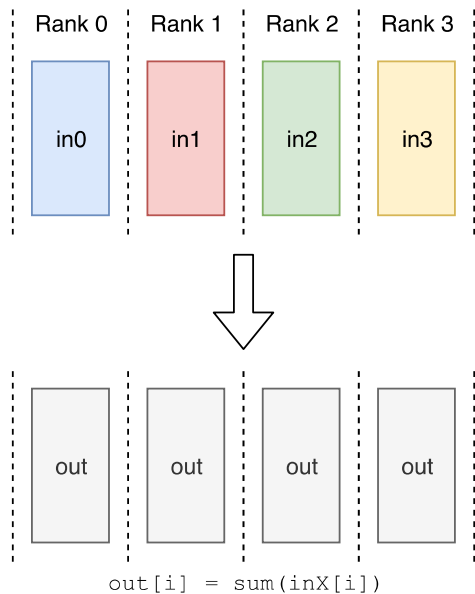
- Certain communication patterns can be optimized
- People keep reinventing MPI
 - Baidu Allreduce, NCCL, Gloo, Horovod, ...
- What we need (for this talk):



Collectives for deep learning

- Certain communication patterns can be optimized
- People keep reinventing MPI
 - Baidu Allreduce, NCCL, Gloo, Horovod, ...
- What we need (for this talk):

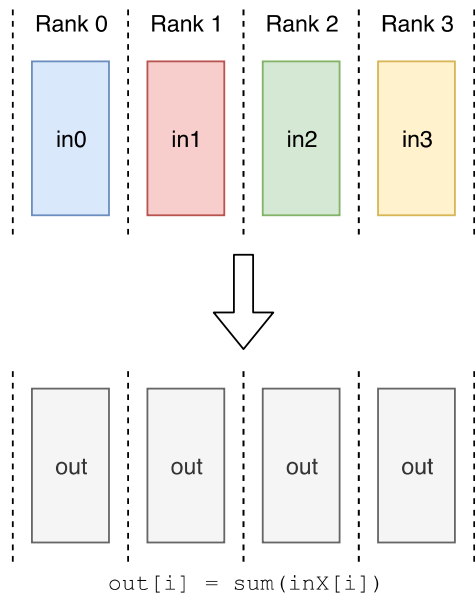
Allreduce



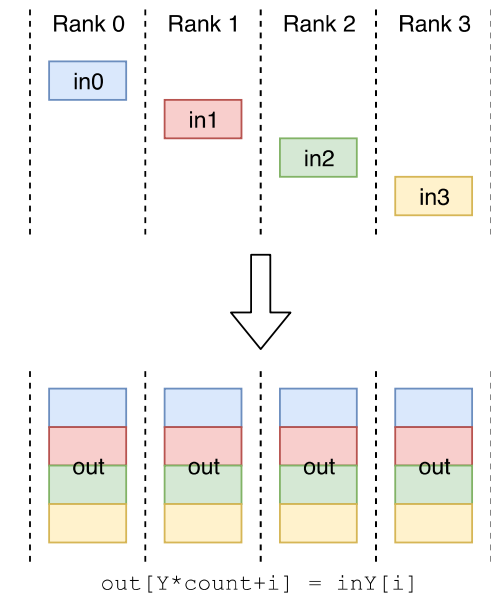
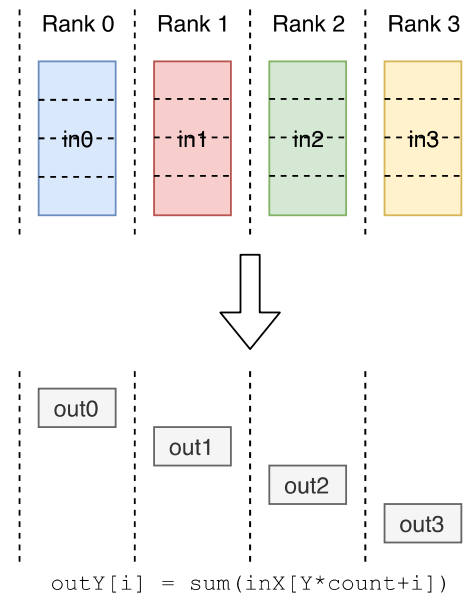
Collectives for deep learning

- Certain communication patterns can be optimized
- People keep reinventing MPI
 - Baidu Allreduce, NCCL, Gloo, Horovod, ...
- What we need (for this talk):

Allreduce



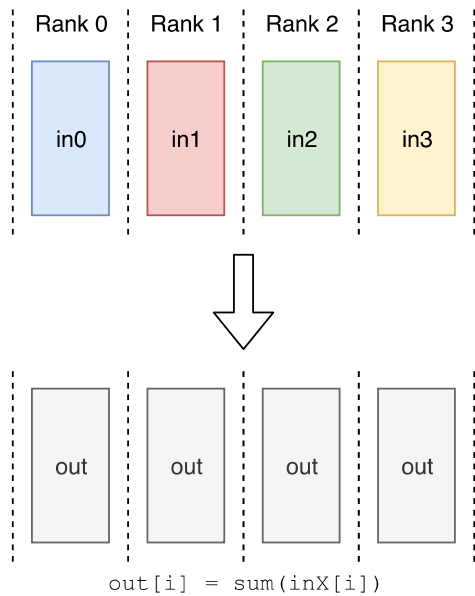
Reduce-scatter



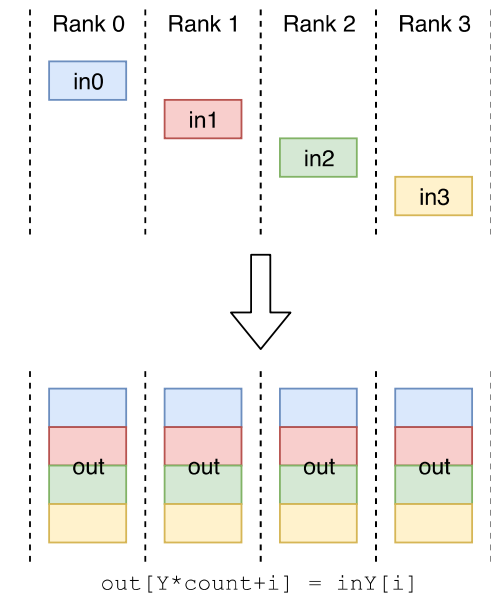
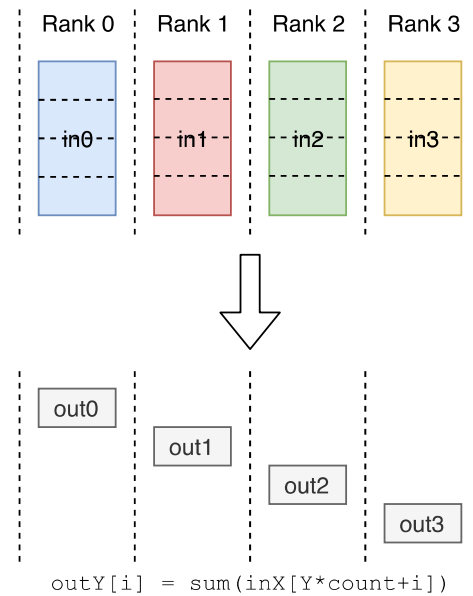
Collectives for deep learning

- Certain communication patterns can be optimized
- People keep reinventing MPI
 - Baidu Allreduce, NCCL, Gloo, Horovod, ...
- What we need (for this talk):

Allreduce



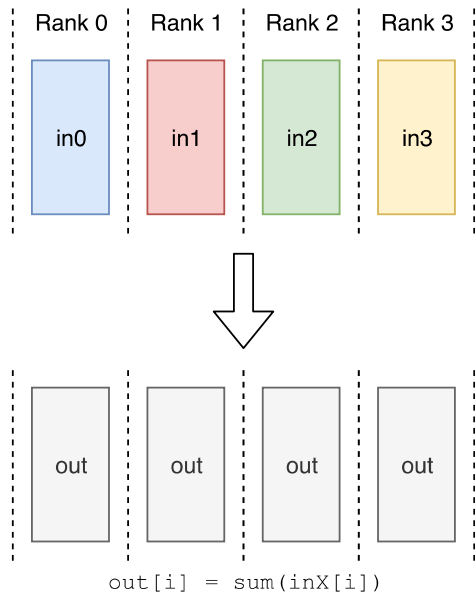
Reduce-scatter



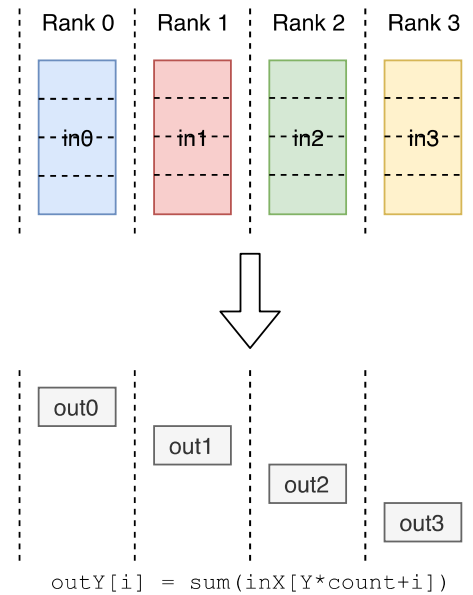
Collectives for deep learning

- Certain communication patterns can be optimized
- People keep reinventing MPI
 - Baidu Allreduce, NCCL, Gloo, Horovod, ...
- What we need (for this talk):

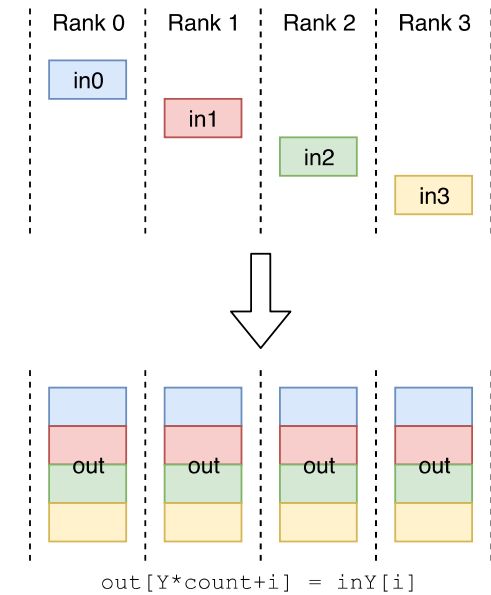
Allreduce



Reduce-scatter



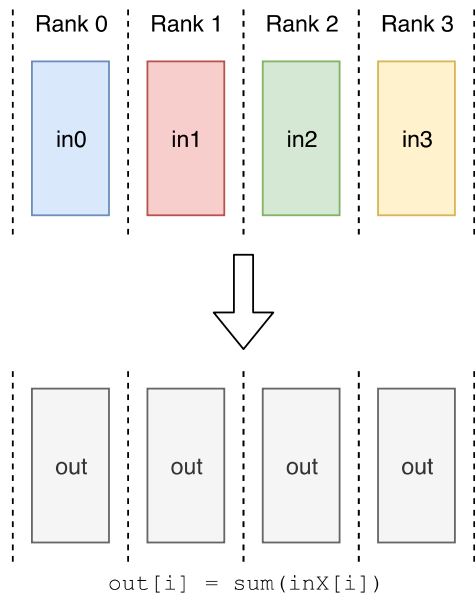
Allgather



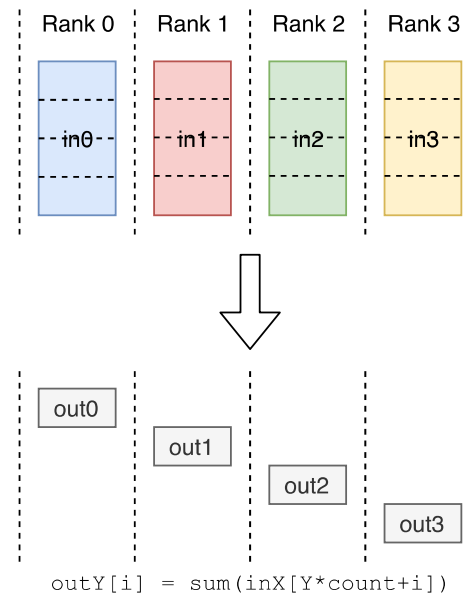
Collectives for deep learning

- Certain communication patterns can be optimized
- People keep reinventing MPI
 - Baidu Allreduce, NCCL, Gloo, Horovod, ...
- What we need (for this talk):

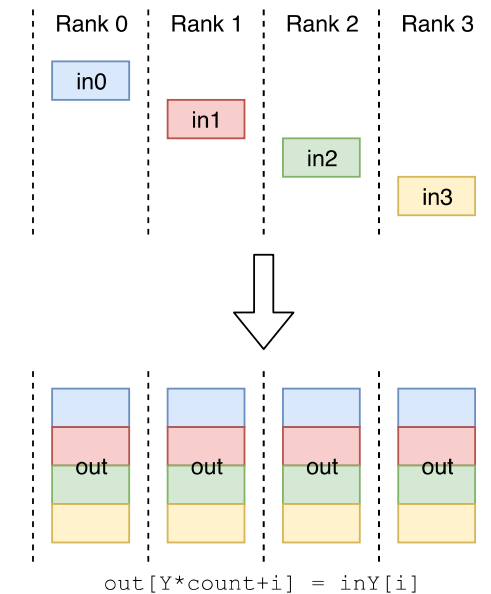
Allreduce



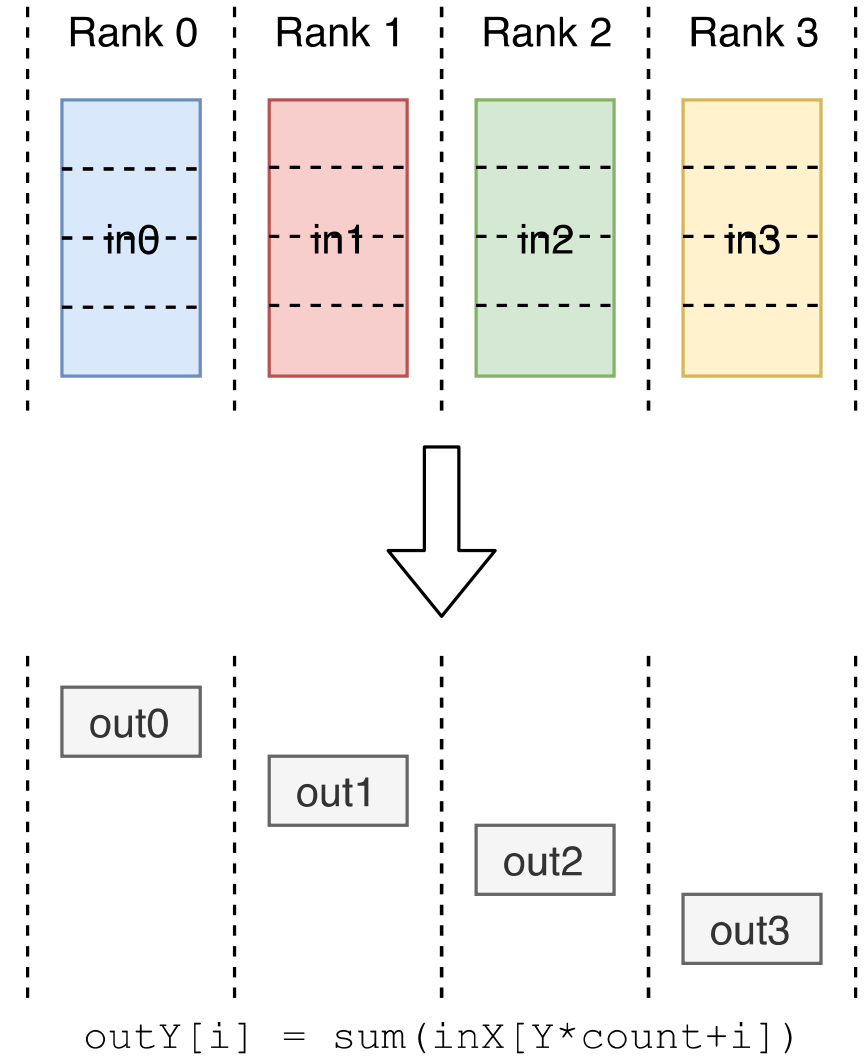
Reduce-scatter



Allgather

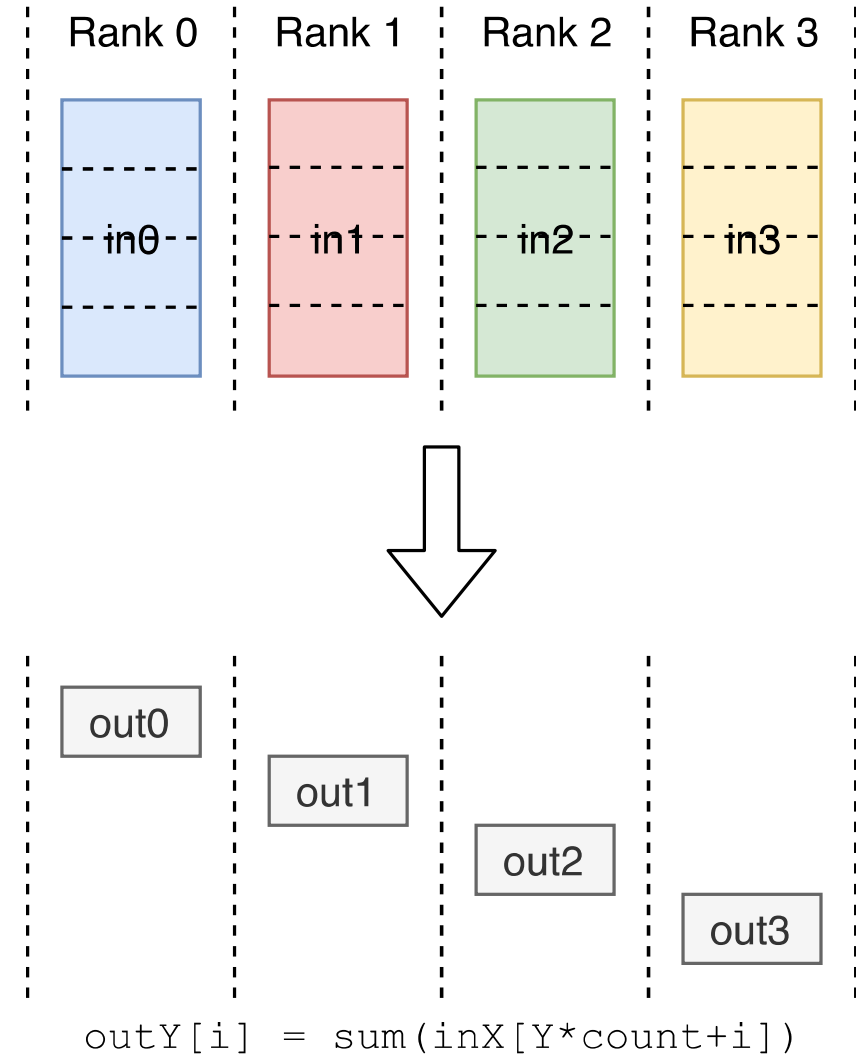


Reduce-scatter



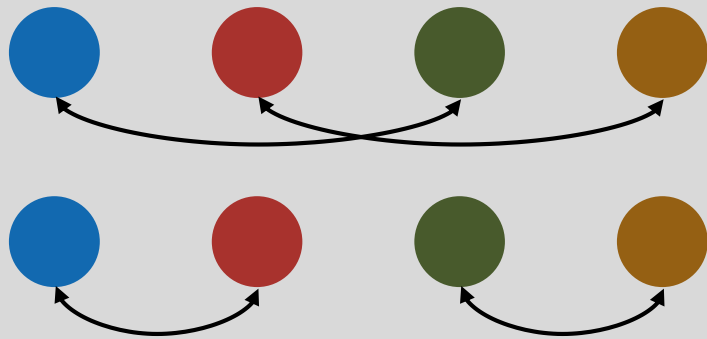
Reduce-scatter

Recursive-halving

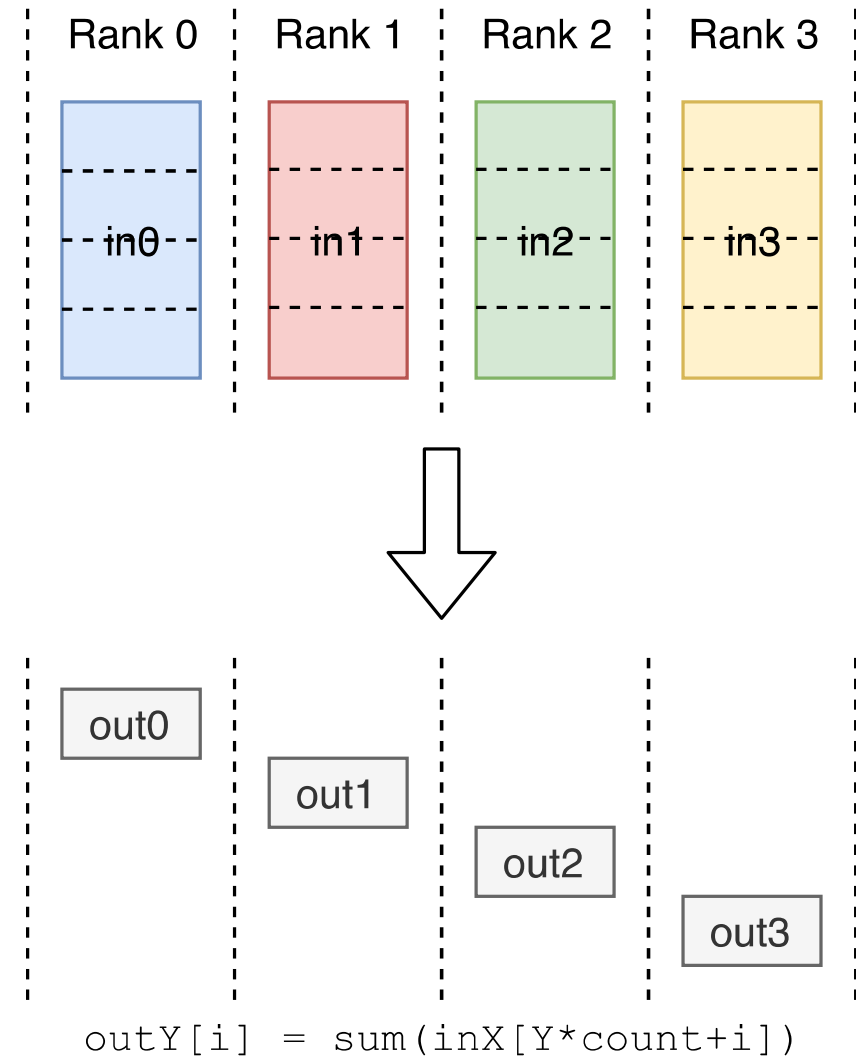


Reduce-scatter

Recursive-halving

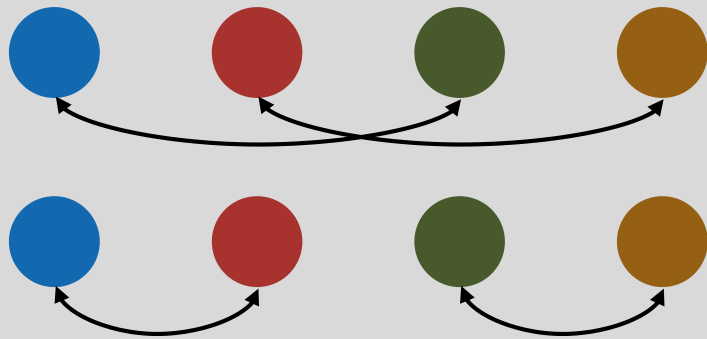


Ring

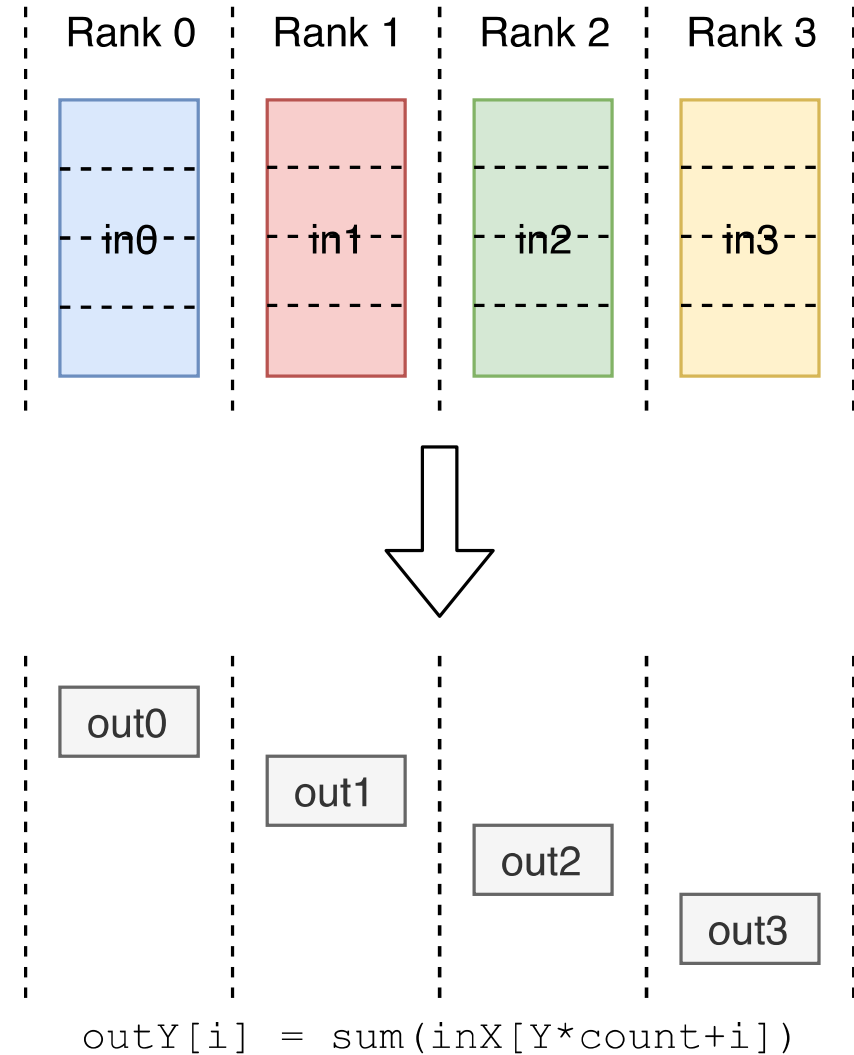


Reduce-scatter

Recursive-halving

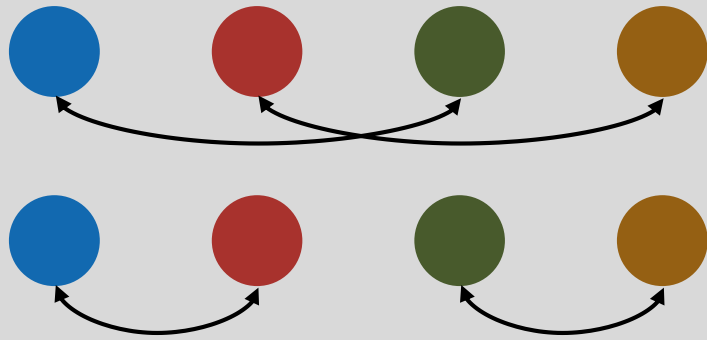


Ring

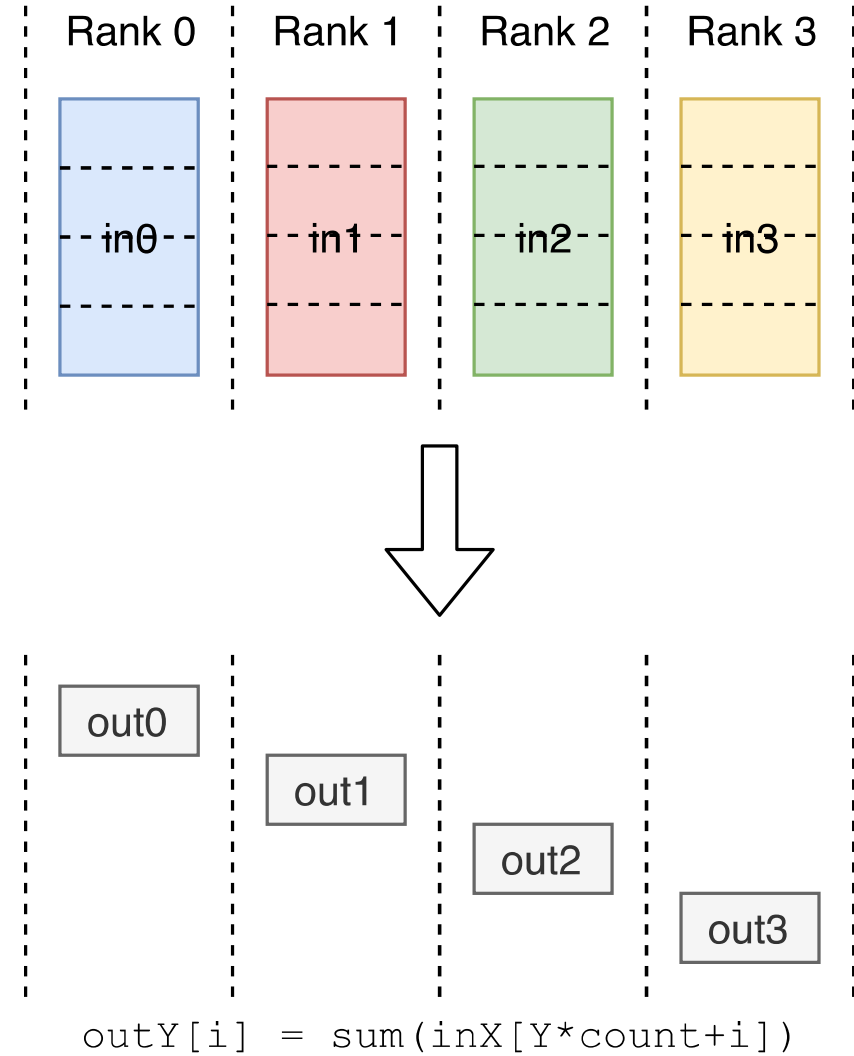
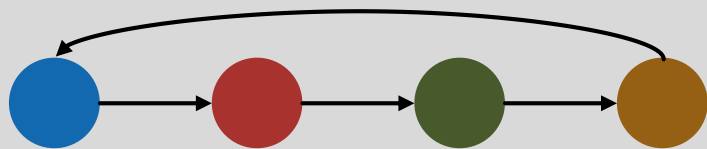


Reduce-scatter

Recursive-halving

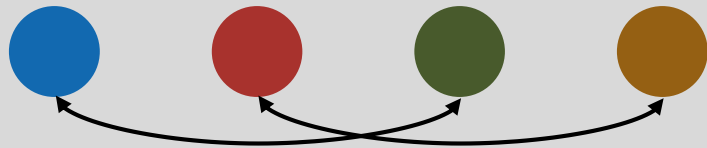


Ring



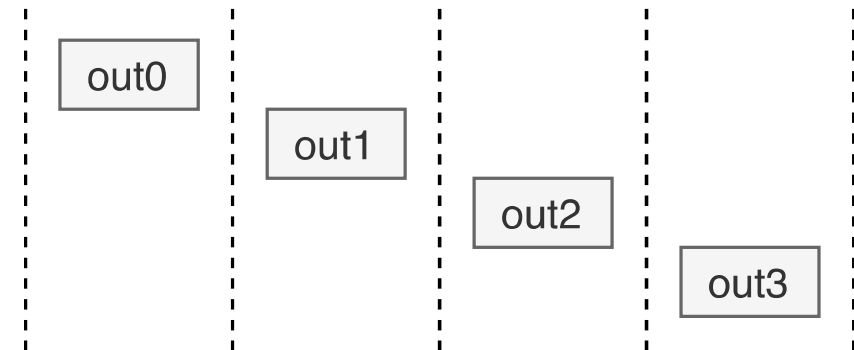
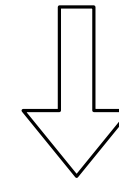
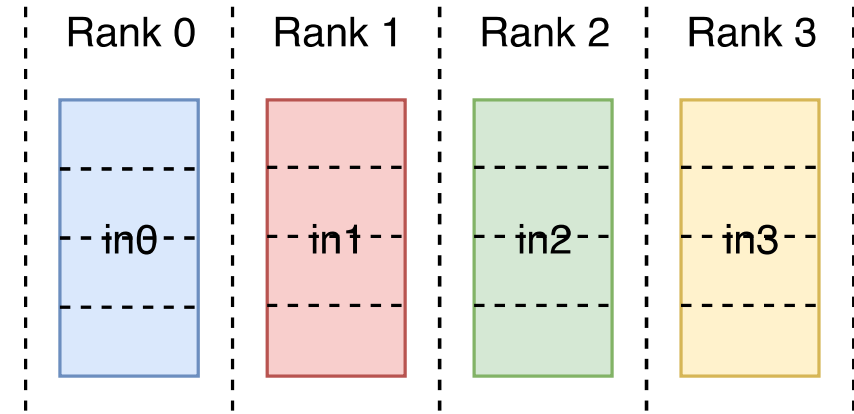
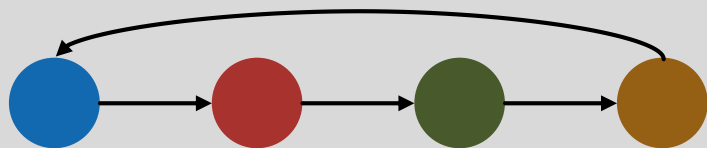
Reduce-scatter

Recursive-halving



$$\alpha \lg p + \frac{p-1}{p} n\beta + \frac{p-1}{p} n\gamma$$

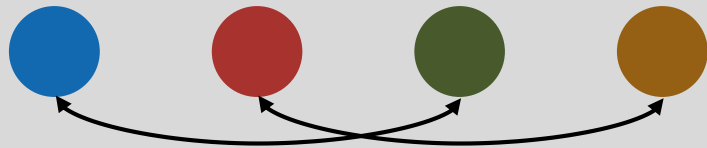
Ring



$$\text{outY}[i] = \text{sum}(\text{inX}[Y \cdot \text{count} + i])$$

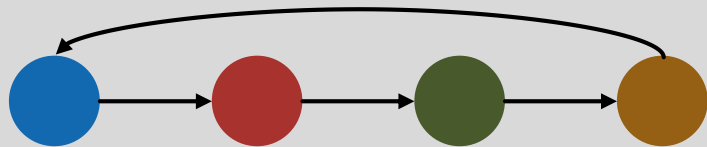
Reduce-scatter

Recursive-halving

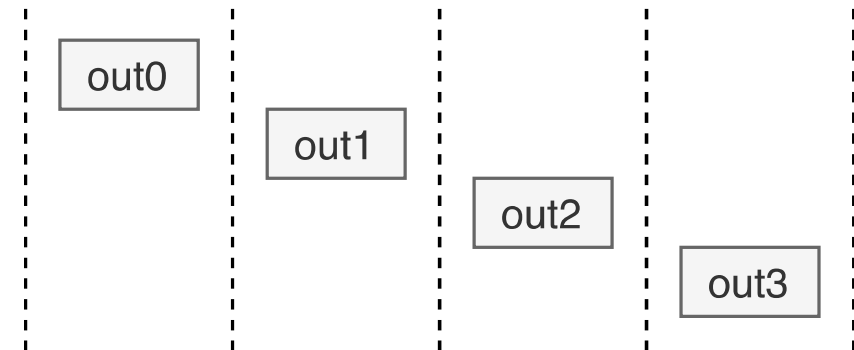
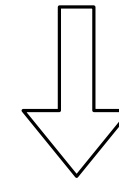
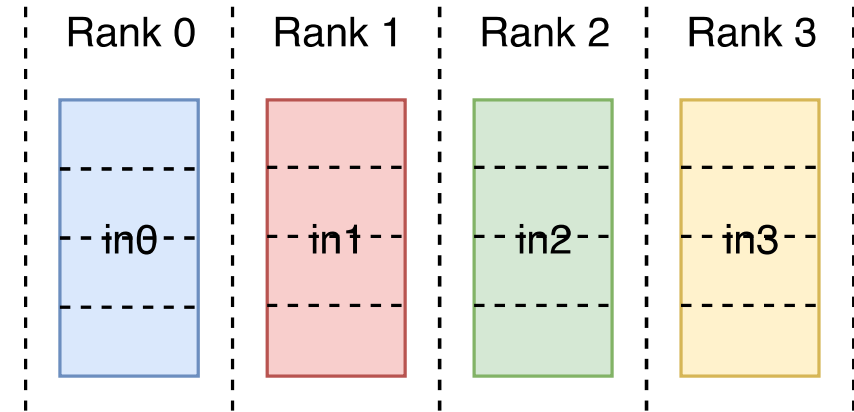


$$\alpha \lg p + \frac{p-1}{p} n\beta + \frac{p-1}{p} n\gamma$$

Ring

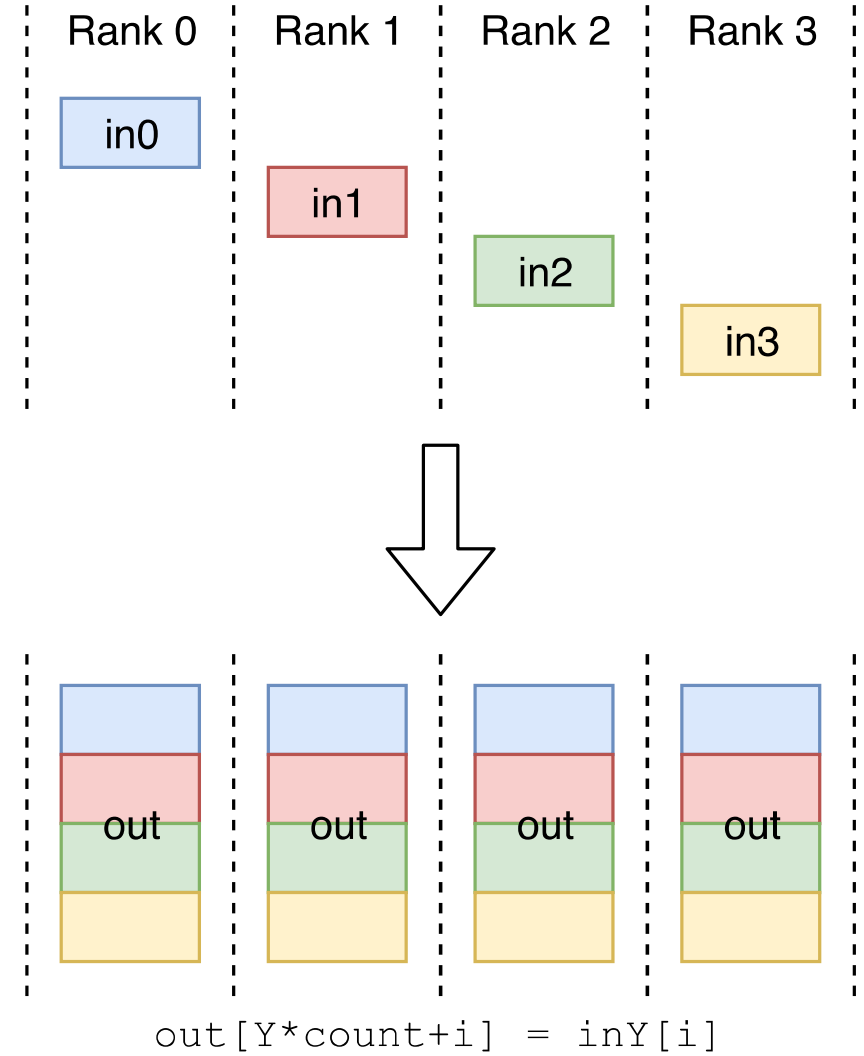


$$(p-1)\alpha + \frac{p-1}{p} n\beta + \frac{p-1}{p} n\gamma$$



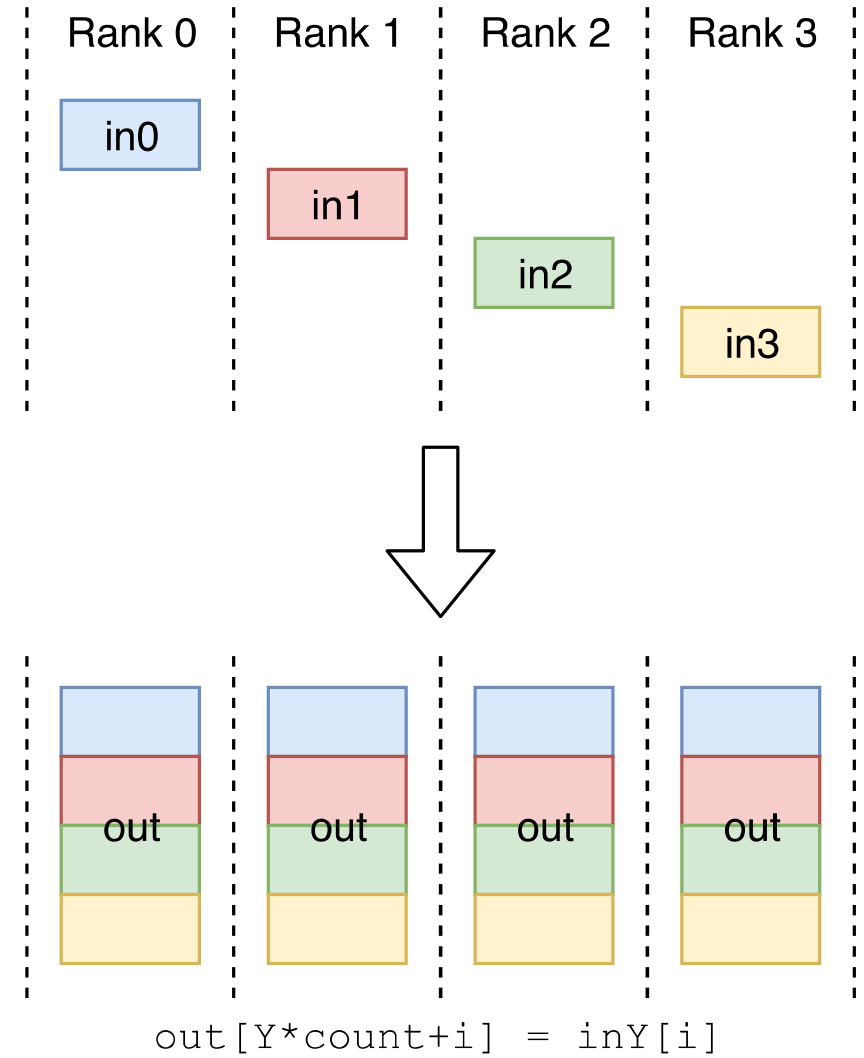
$$\text{outY}[i] = \text{sum}(\text{inX}[Y*\text{count}+i])$$

Allgather



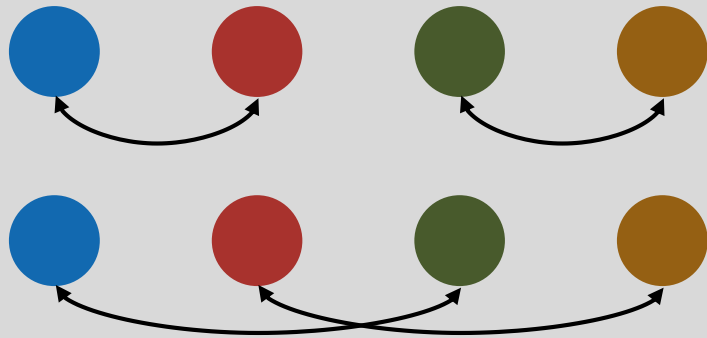
Allgather

Recursive-doubling

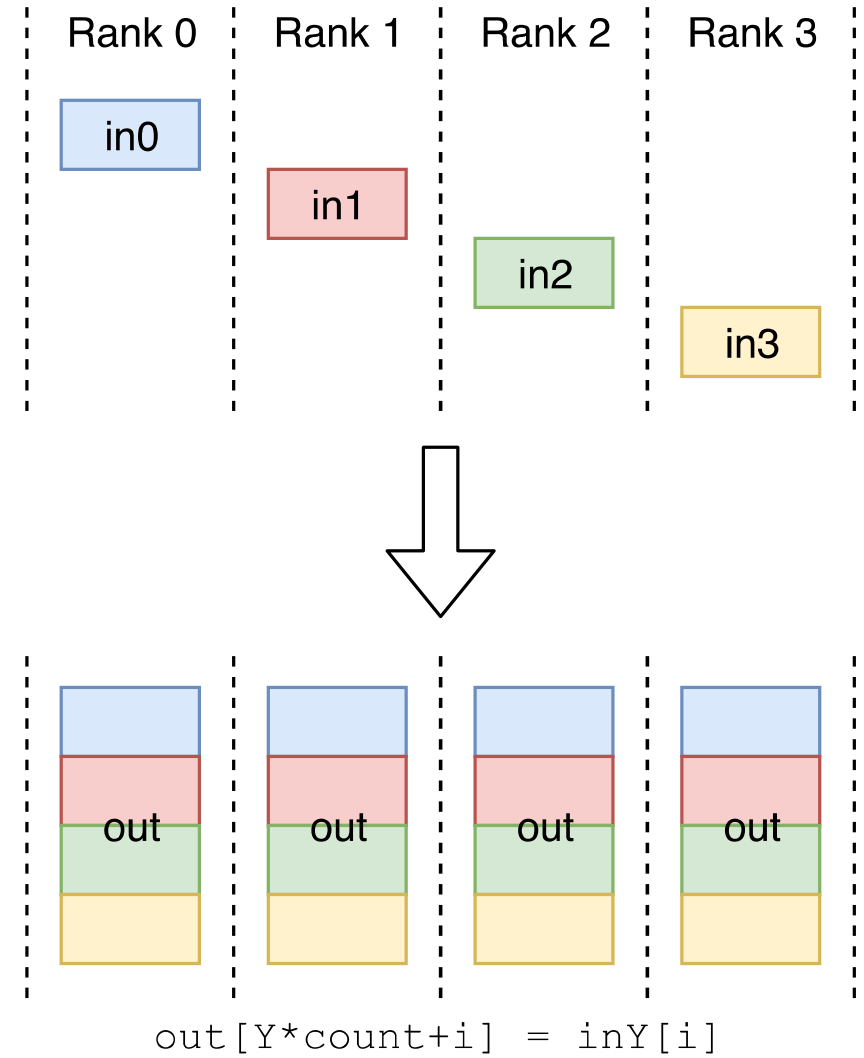


Allgather

Recursive-doubling

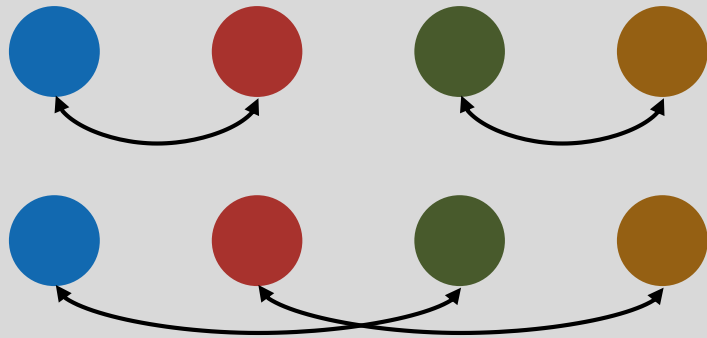


Ring

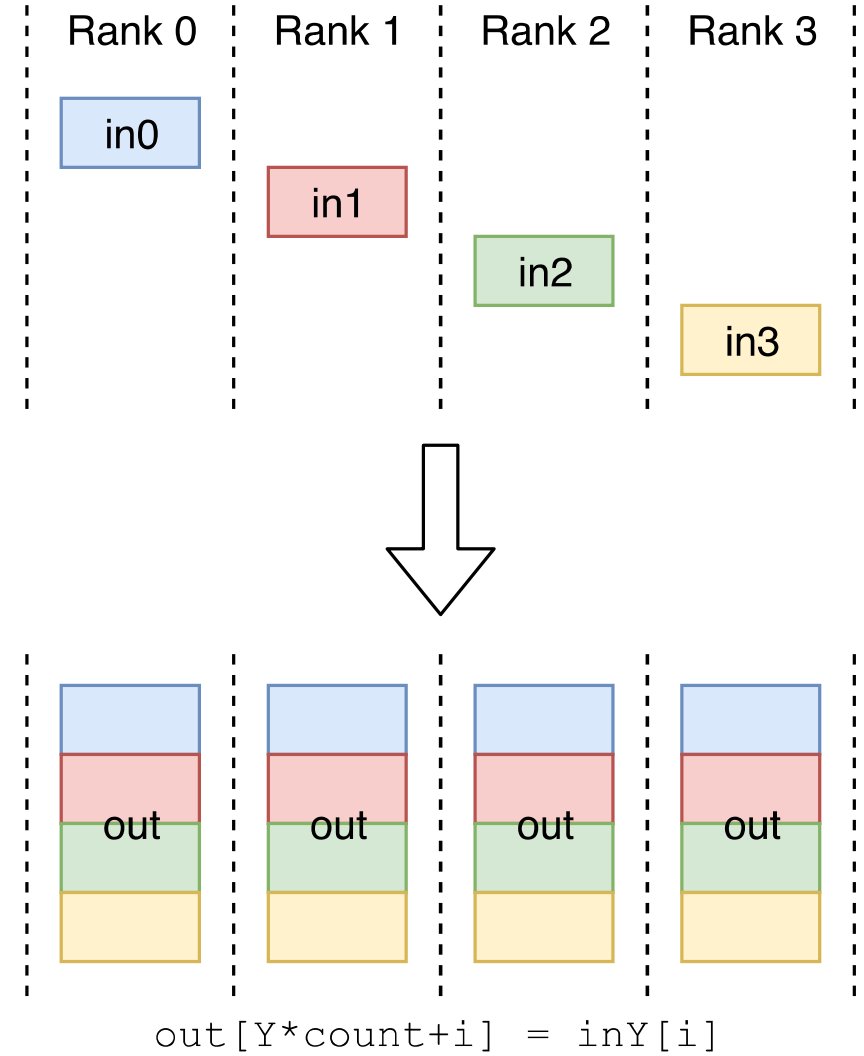
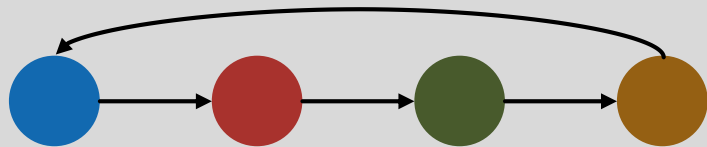


Allgather

Recursive-doubling



Ring



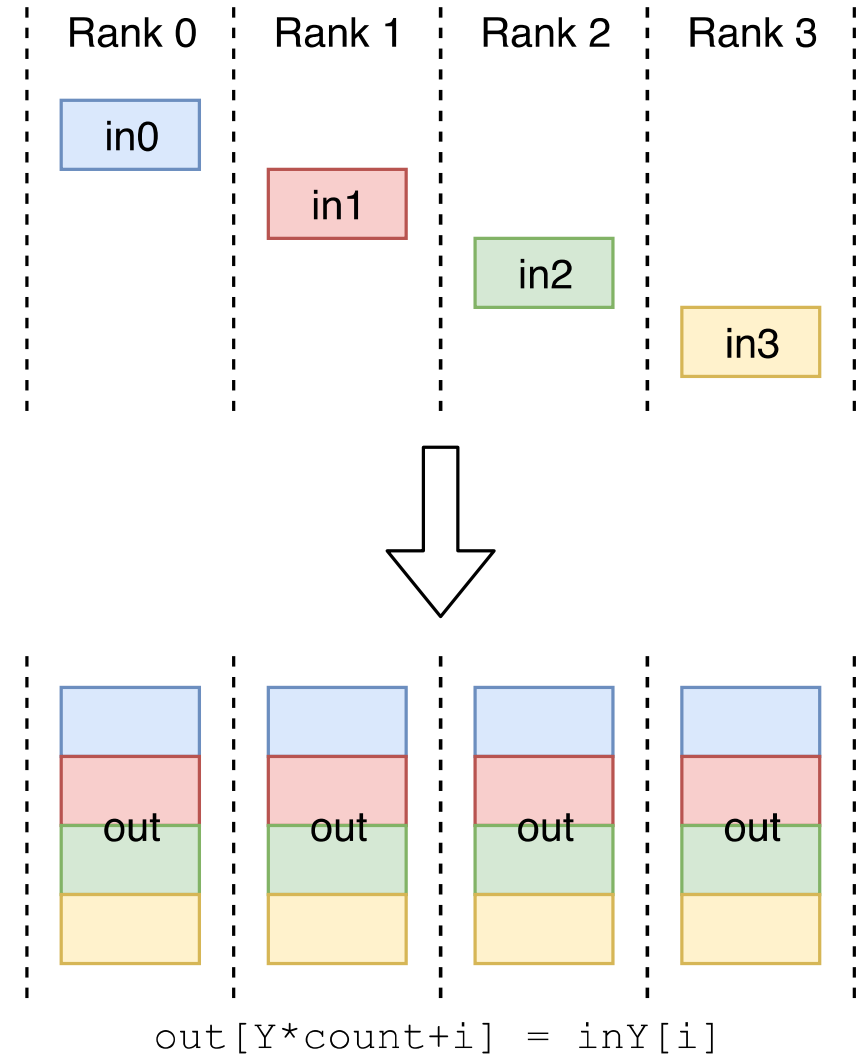
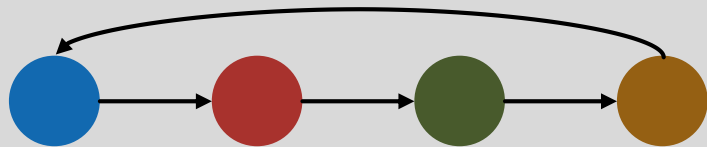
Allgather

Recursive-doubling



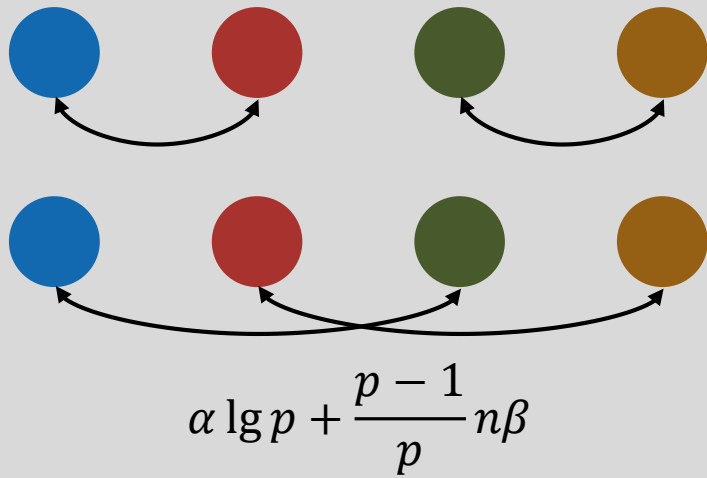
$$\alpha \lg p + \frac{p-1}{p} n\beta$$

Ring

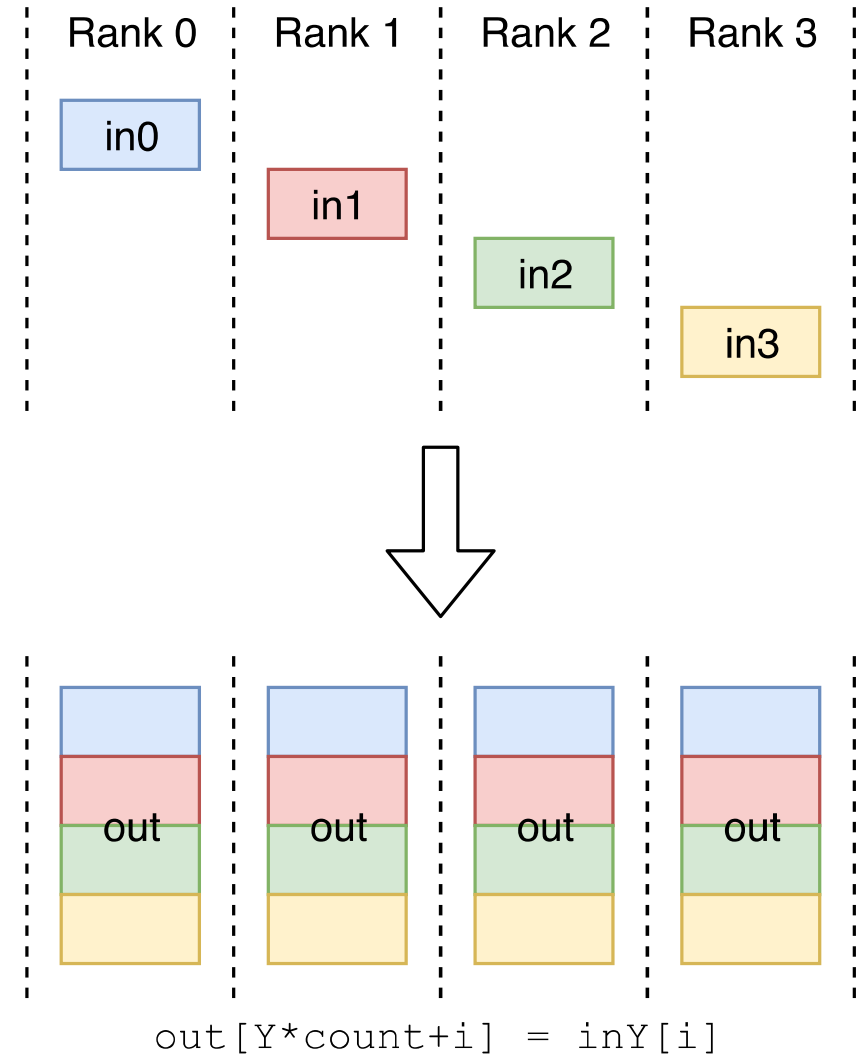
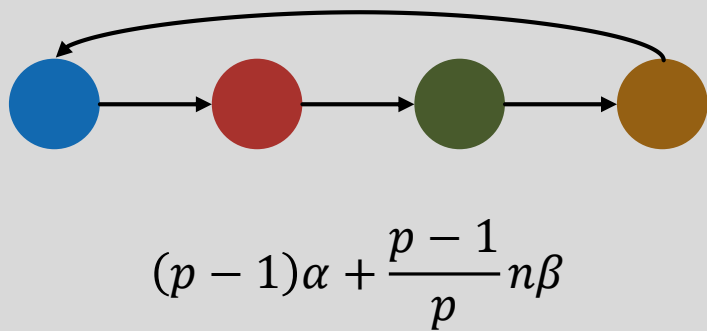


Allgather

Recursive-doubling

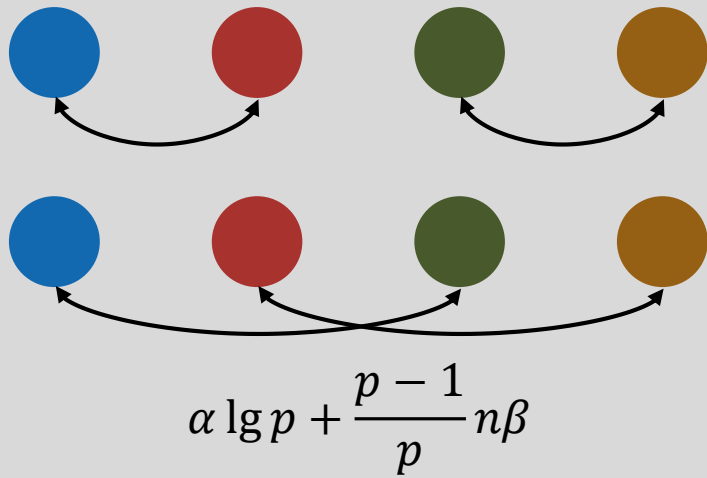


Ring

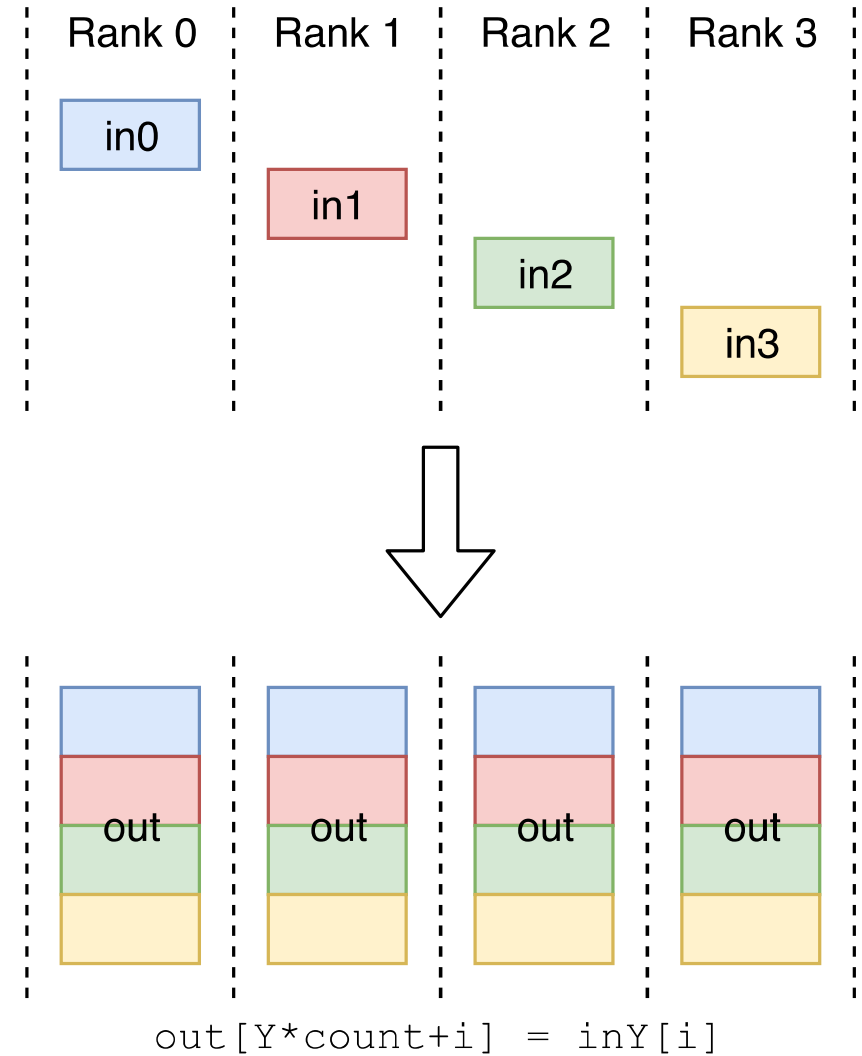
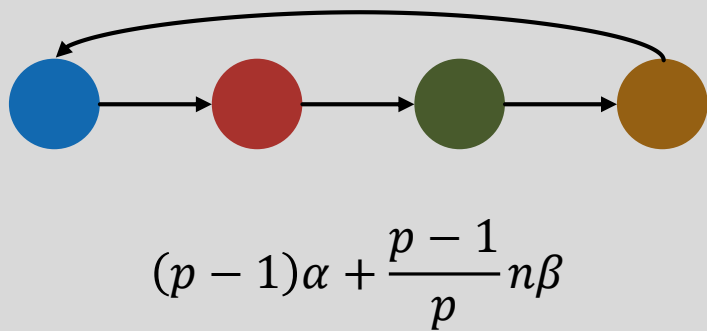


Allgather

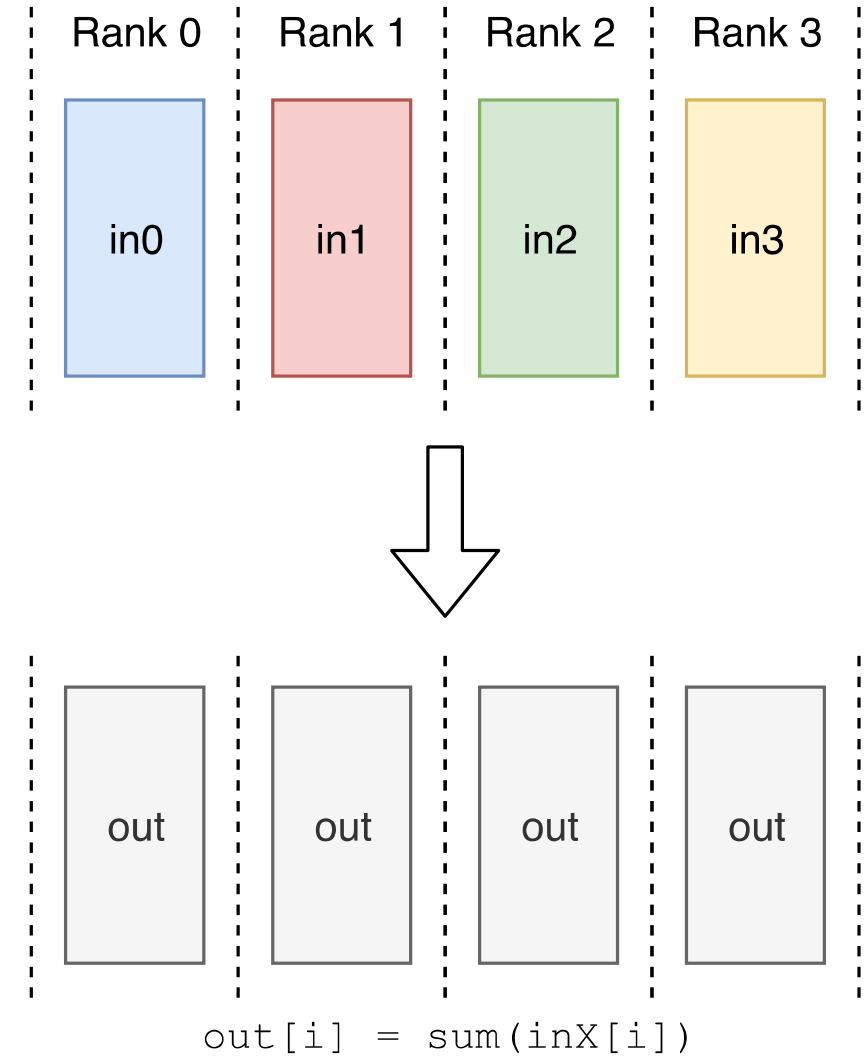
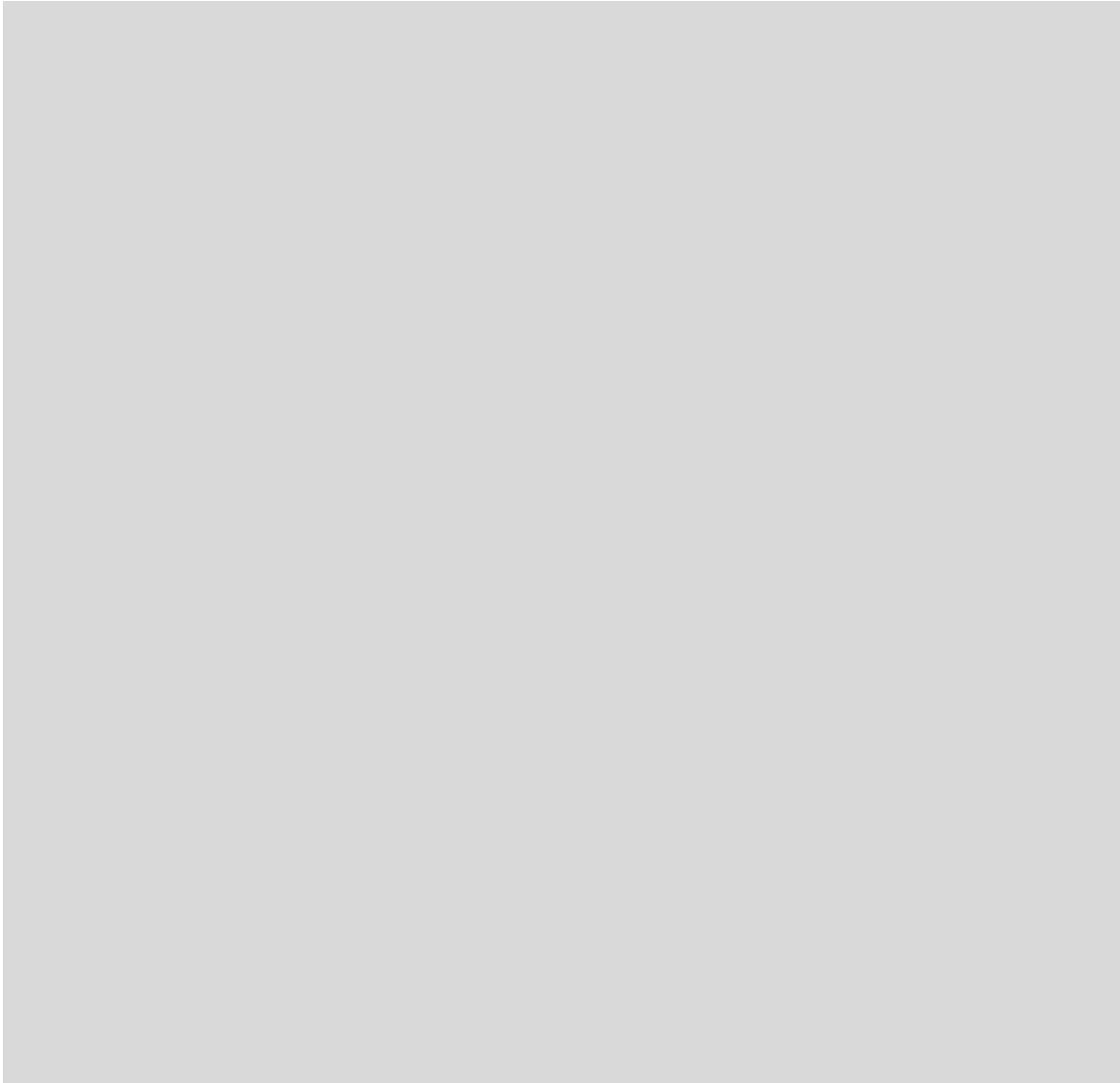
Recursive-doubling



Ring

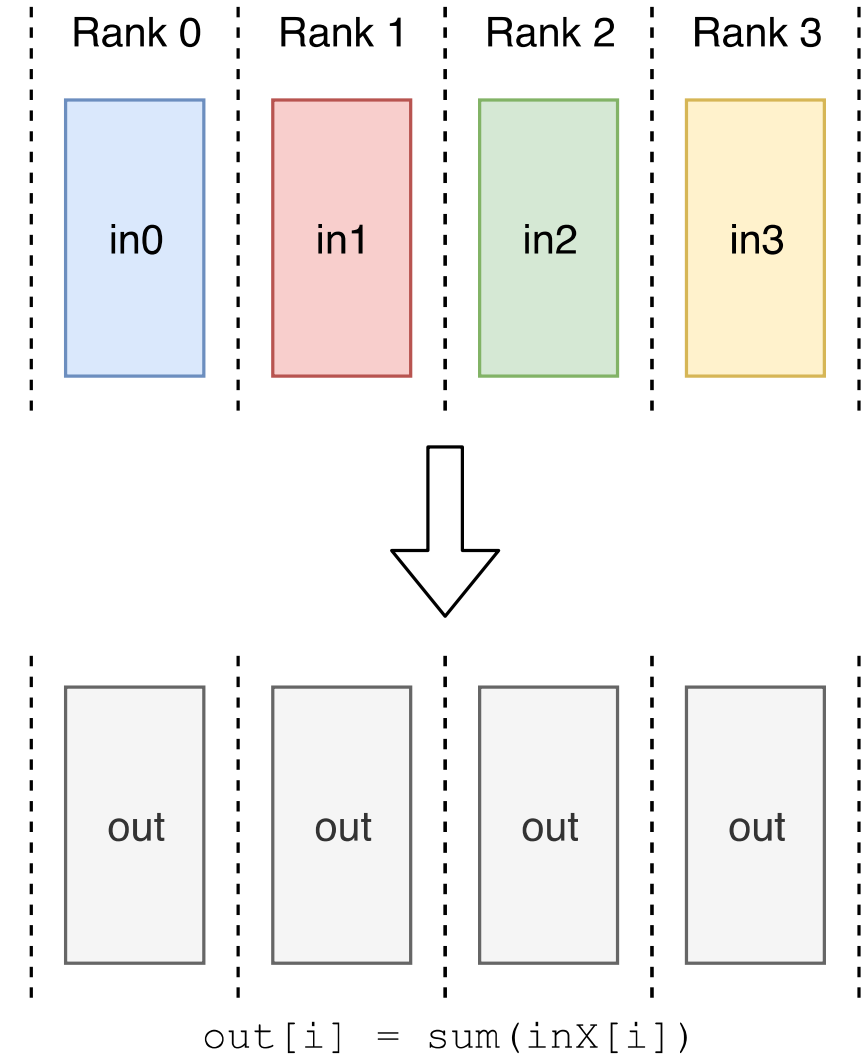


Allreduce



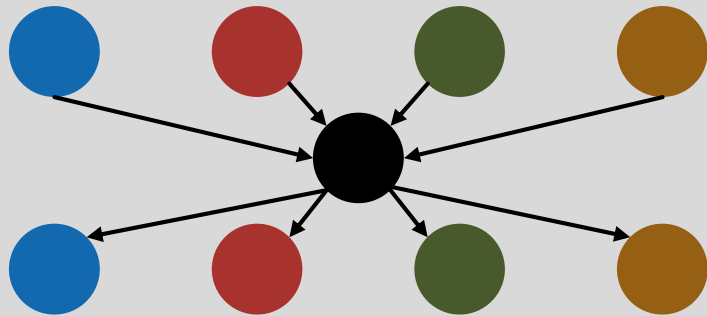
Allreduce

Parameter Server

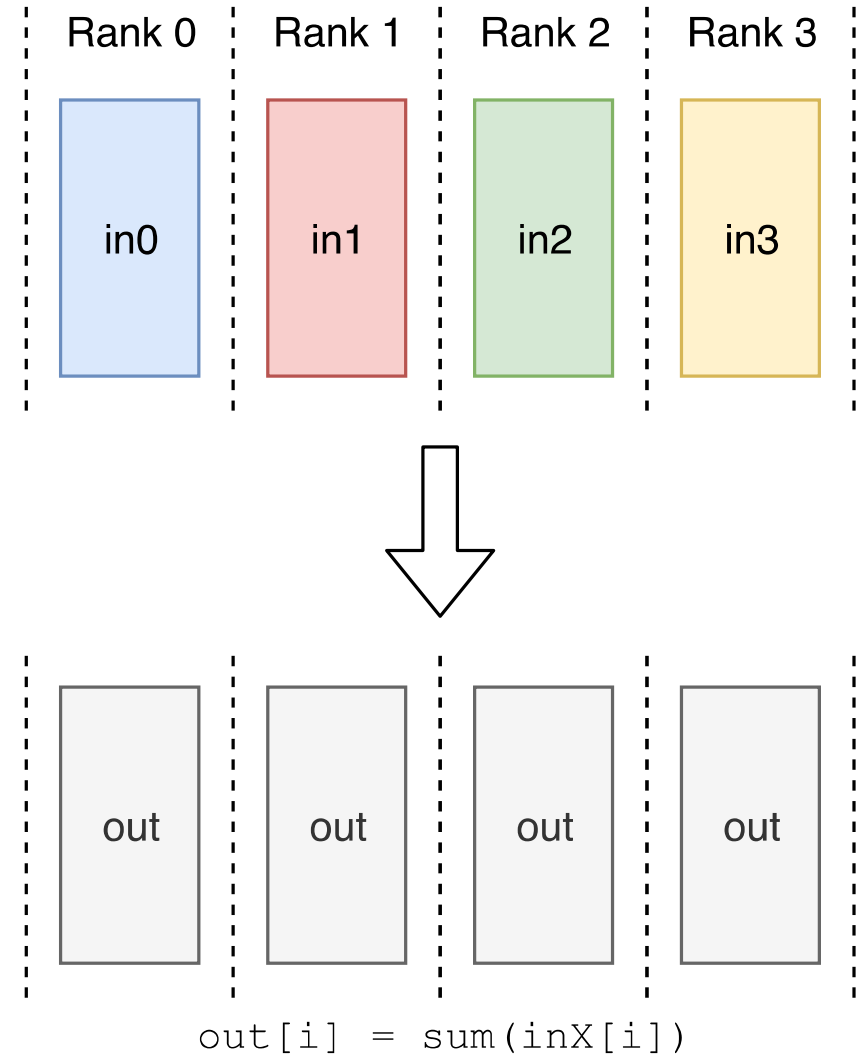


Allreduce

Parameter Server

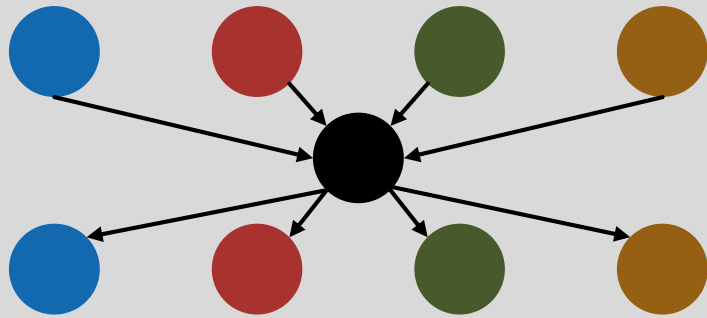


Tree (reduce/broadcast)

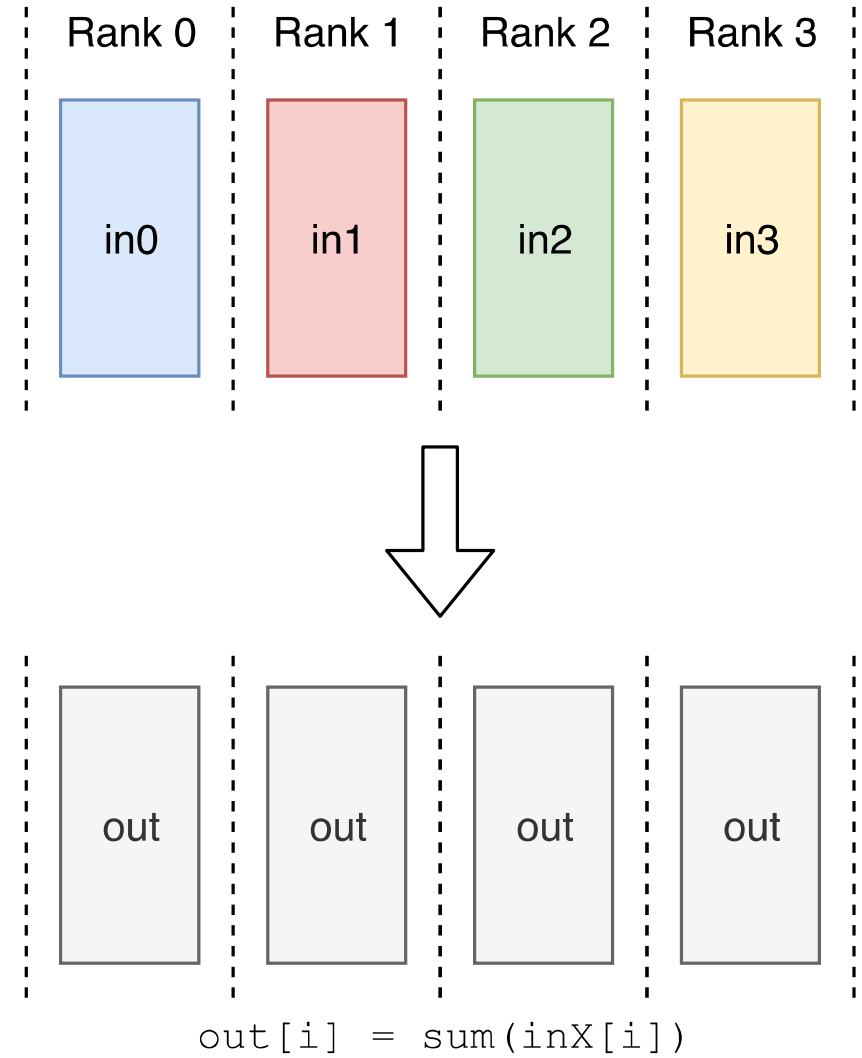


Allreduce

Parameter Server

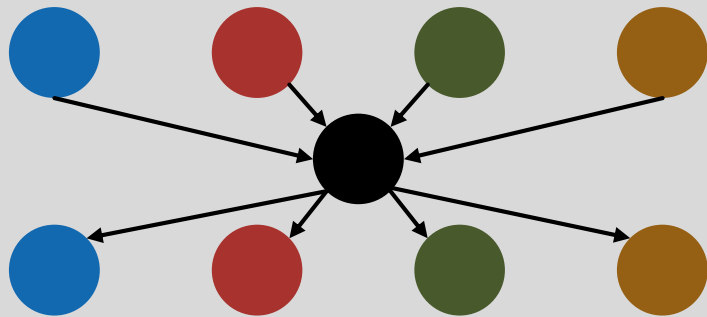


Tree (reduce/broadcast)



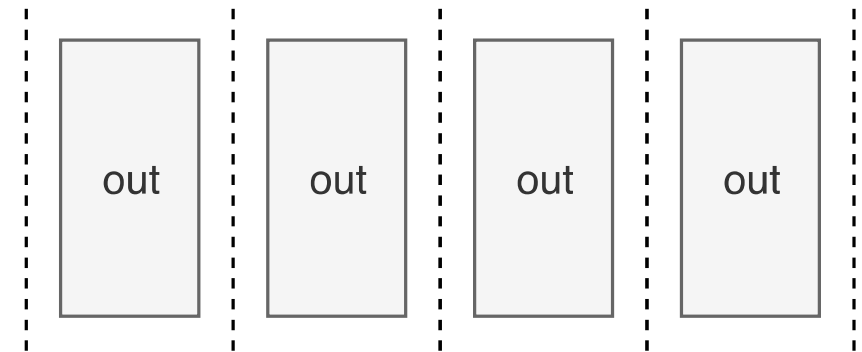
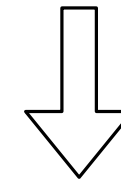
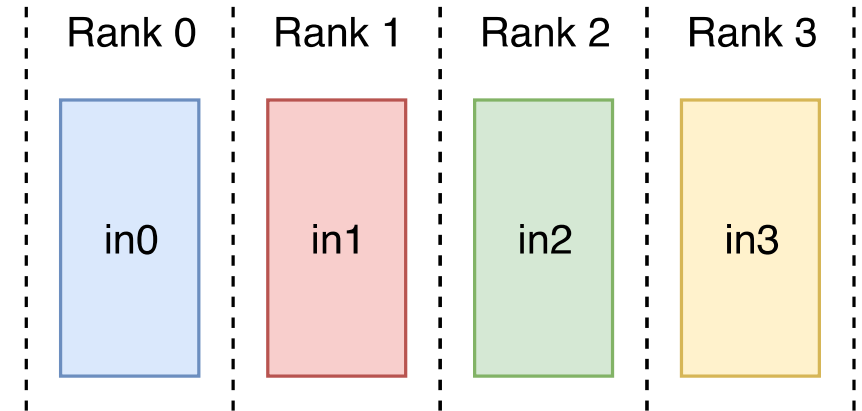
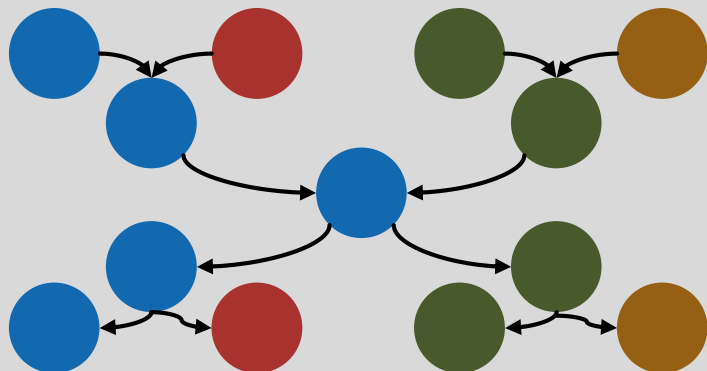
Allreduce

Parameter Server



$$2p\alpha + 2pn\beta + pn\gamma$$

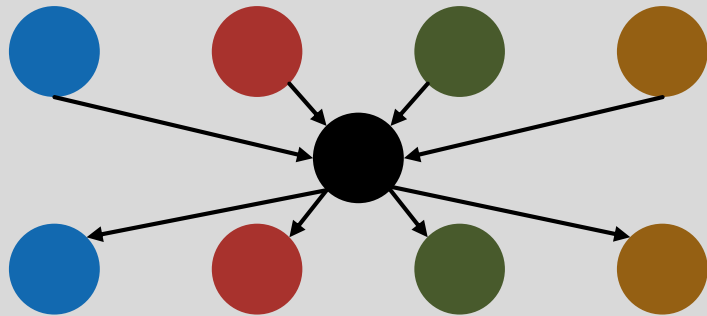
Tree (reduce/broadcast)



$$\text{out}[i] = \text{sum}(\text{inX}[i])$$

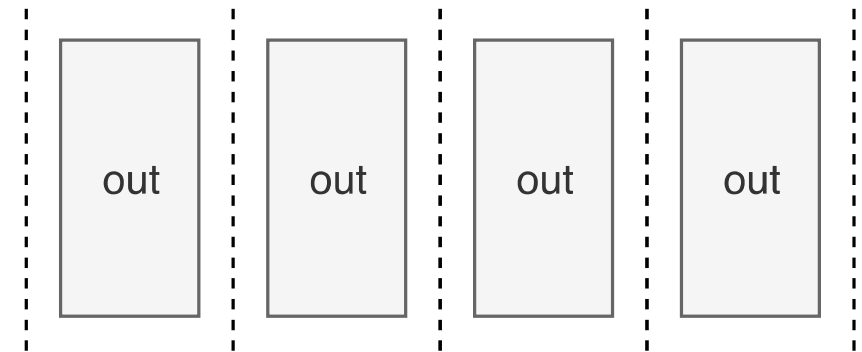
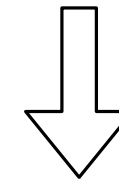
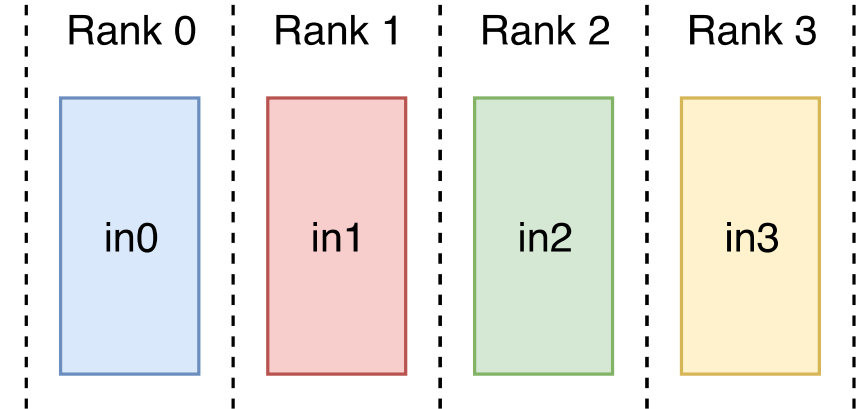
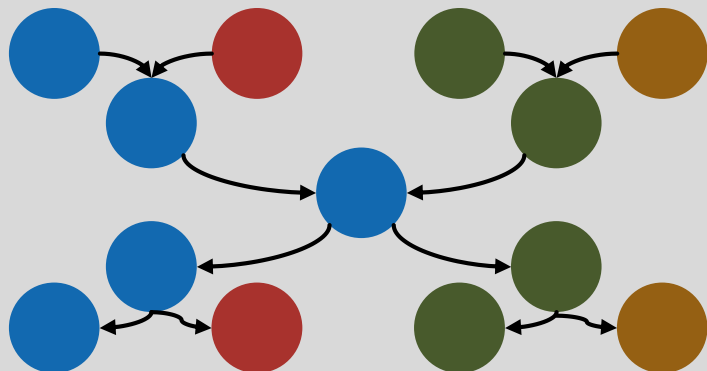
Allreduce

Parameter Server



$$2p\alpha + 2pn\beta + pn\gamma$$

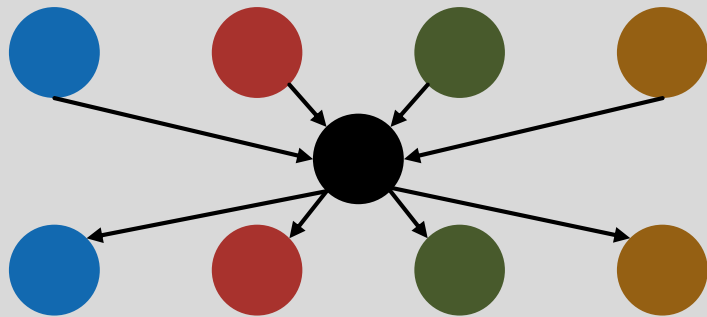
Tree (reduce/broadcast)



$$\text{out}[i] = \text{sum}(\text{inX}[i])$$

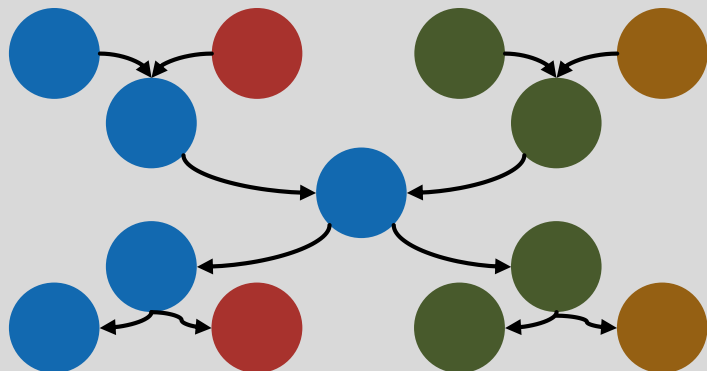
Allreduce

Parameter Server

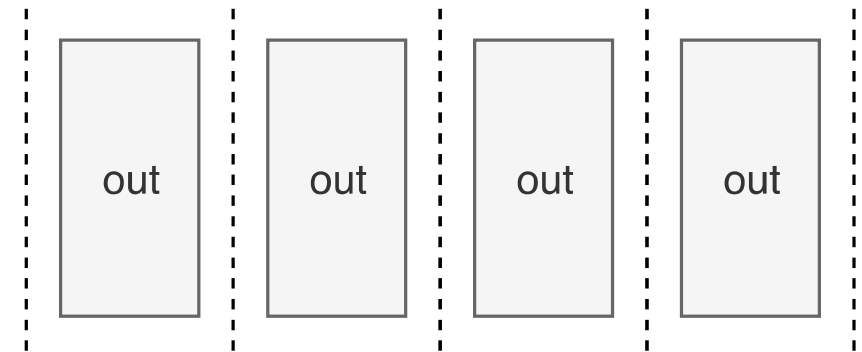
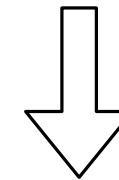
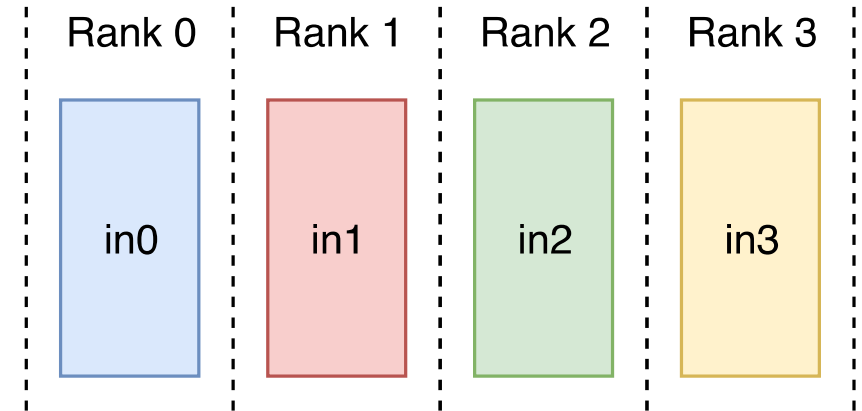


$$2p\alpha + 2pn\beta + pn\gamma$$

Tree (reduce/broadcast)



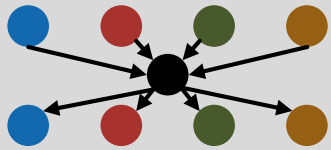
$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$



$$\text{out}[i] = \text{sum}(\text{inX}[i])$$

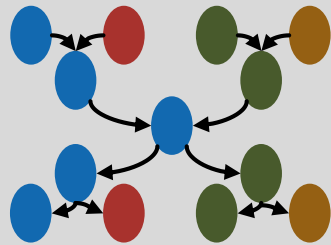
Allreduce

Parameter Server



$$2p\alpha + 2pn\beta + pn\gamma$$

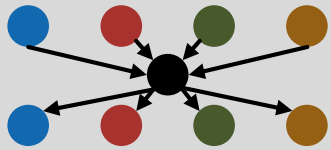
Tree



$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

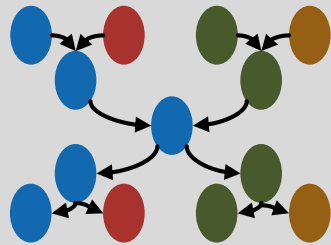
Allreduce

Parameter Server



$$2p\alpha + 2pn\beta + pn\gamma$$

Tree

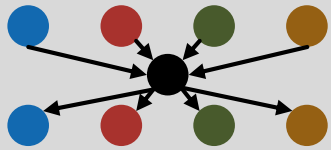


$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

Butterfly (doubling)

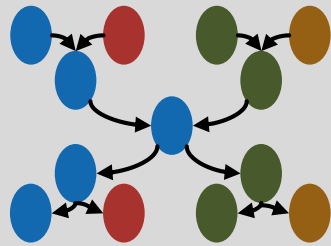
Allreduce

Parameter Server



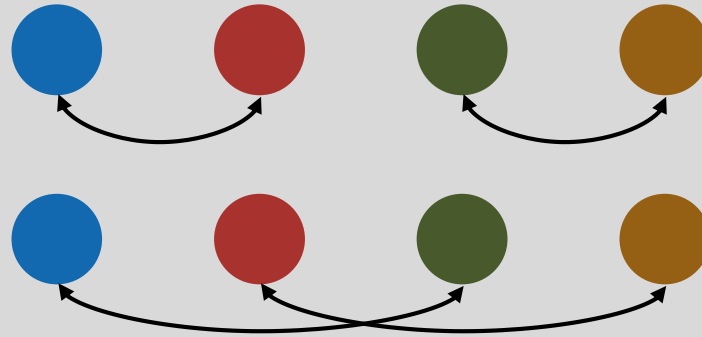
$$2p\alpha + 2pn\beta + pn\gamma$$

Tree



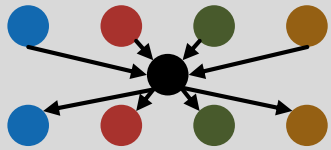
$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

Butterfly (doubling)



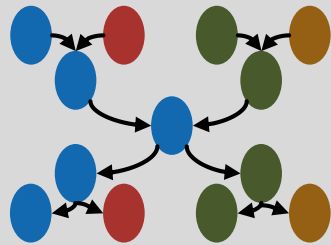
Allreduce

Parameter Server



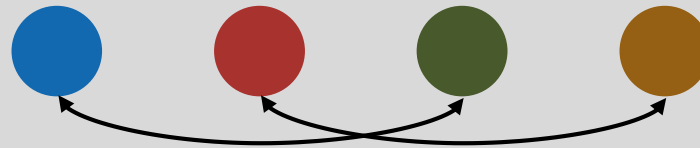
$$2p\alpha + 2pn\beta + pn\gamma$$

Tree



$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

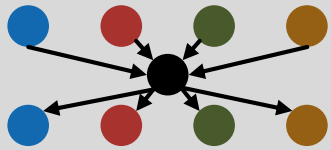
Butterfly (doubling)



Rabenseifner (half/double)

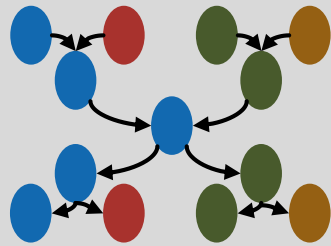
Allreduce

Parameter Server



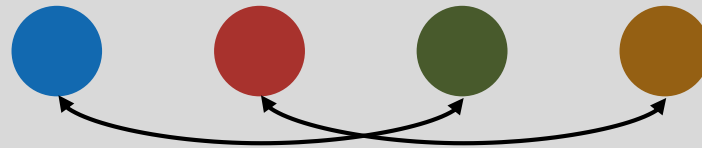
$$2p\alpha + 2pn\beta + pn\gamma$$

Tree

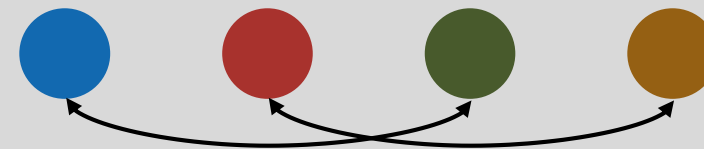
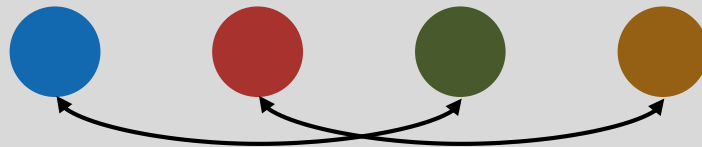


$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

Butterfly (doubling)



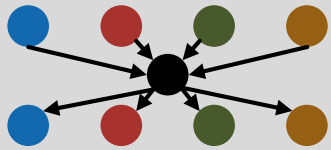
Rabenseifner (half/double)



Ring

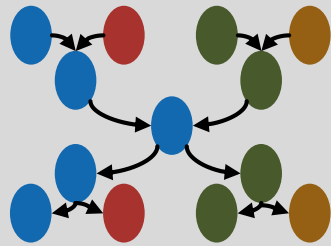
Allreduce

Parameter Server



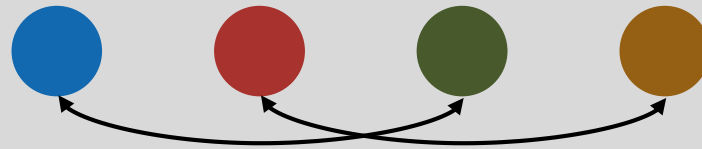
$$2p\alpha + 2pn\beta + pn\gamma$$

Tree

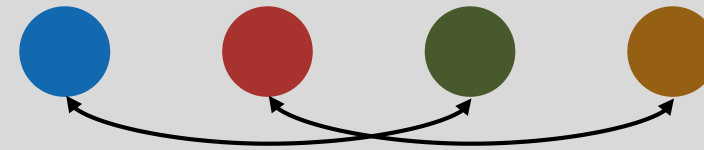
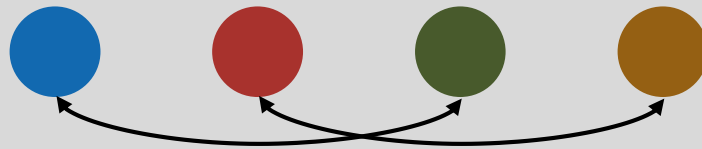


$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

Butterfly (doubling)



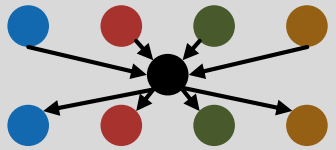
Rabenseifner (half/double)



Ring

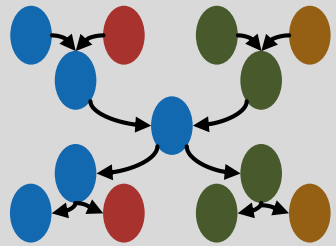
Allreduce

Parameter Server



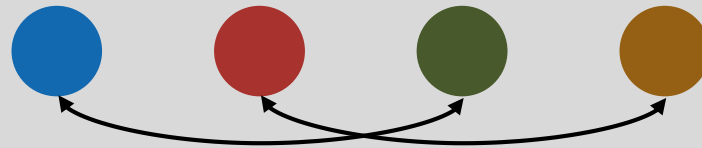
$$2p\alpha + 2pn\beta + pn\gamma$$

Tree

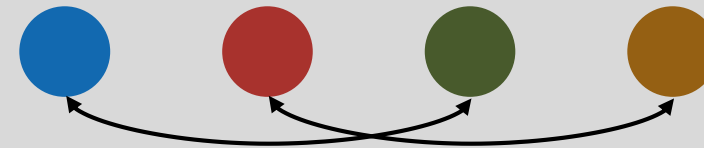
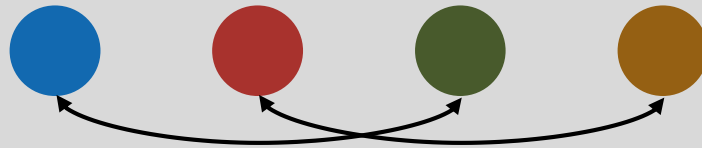


$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

Butterfly (doubling)



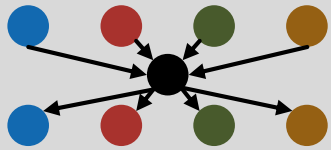
Rabenseifner (half/double)



Ring

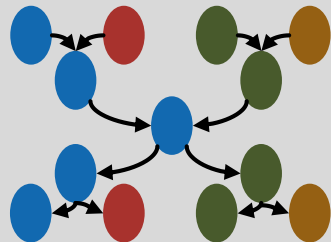
Allreduce

Parameter Server



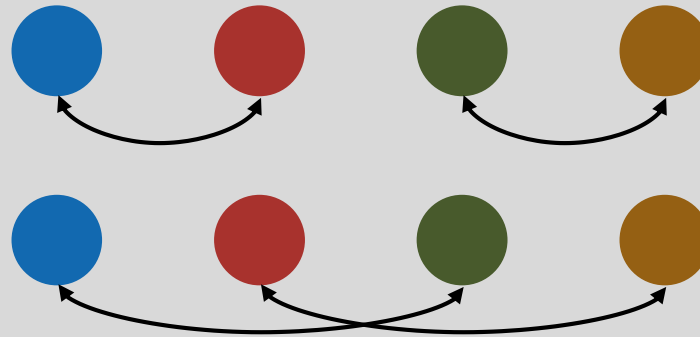
$$2p\alpha + 2pn\beta + pn\gamma$$

Tree

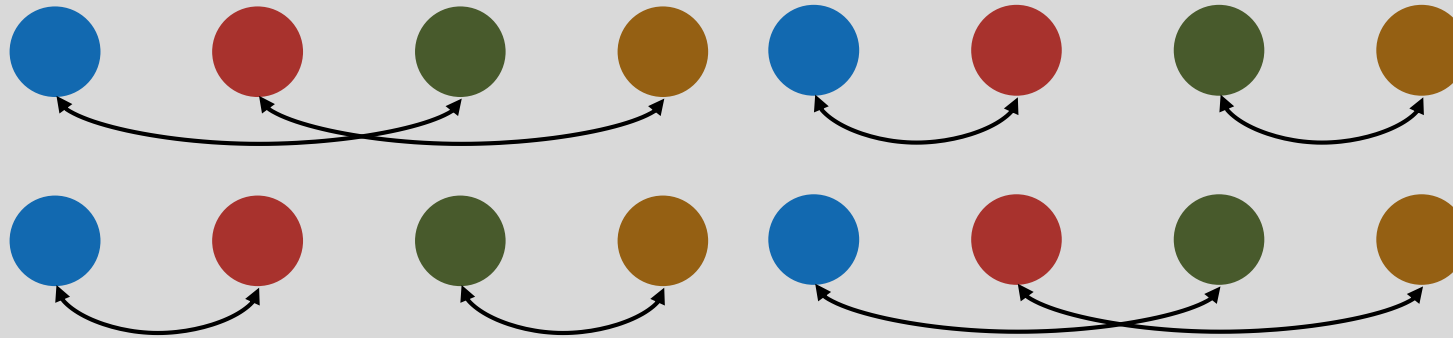


$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

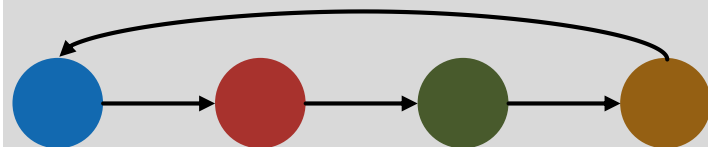
Butterfly (doubling)



Rabenseifner (half/double)

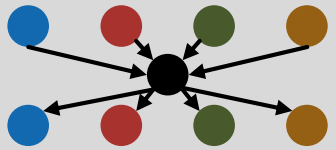


Ring



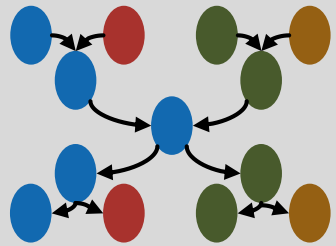
Allreduce

Parameter Server



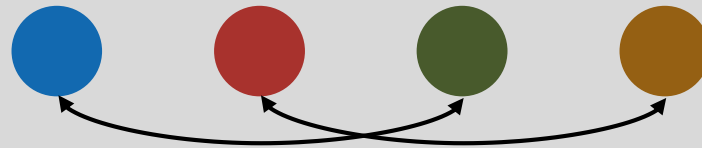
$$2p\alpha + 2pn\beta + pn\gamma$$

Tree

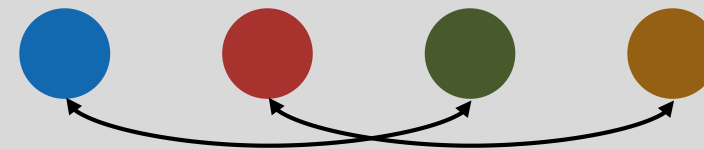
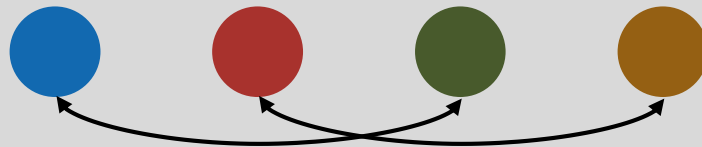


$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

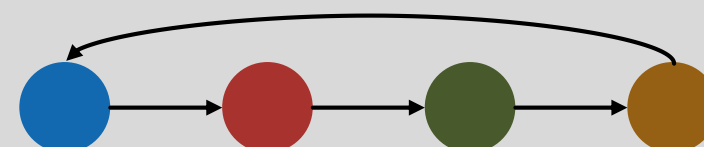
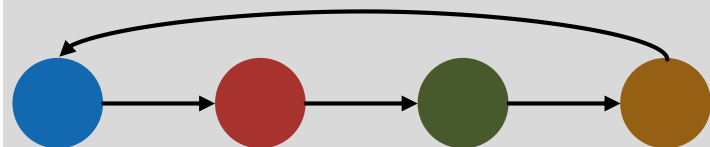
Butterfly (doubling)



Rabenseifner (half/double)

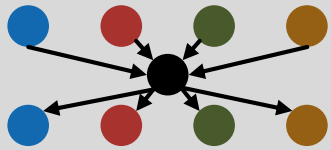


Ring



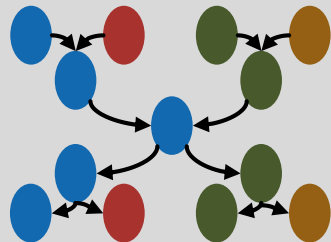
Allreduce

Parameter Server



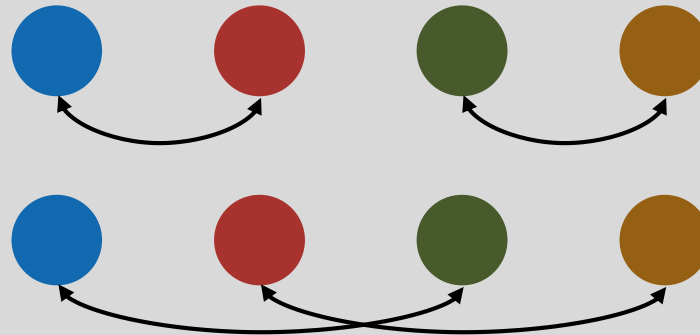
$$2p\alpha + 2pn\beta + pn\gamma$$

Tree

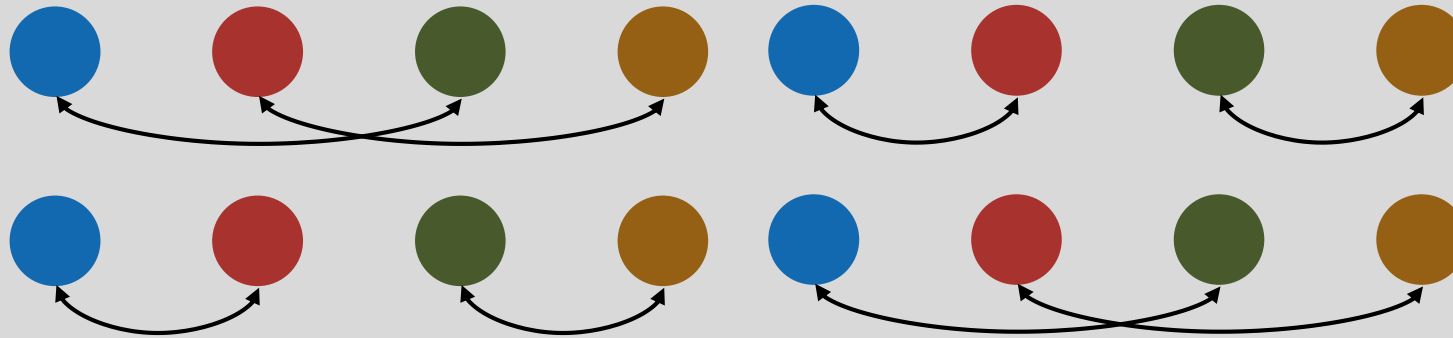


$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

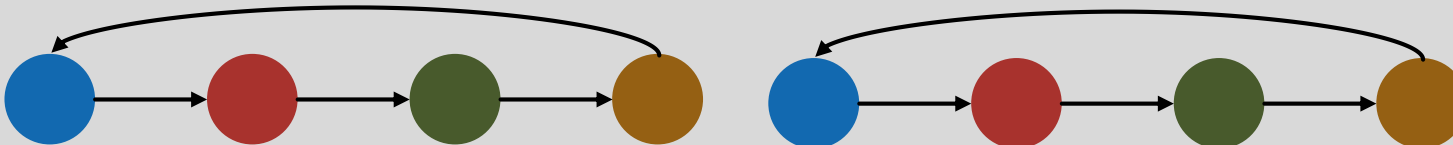
Butterfly (doubling)



Rabenseifner (half/double)



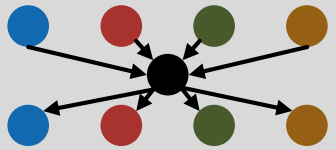
Ring



$$\alpha \lg p + \beta n \lg p + \gamma n \lg p$$

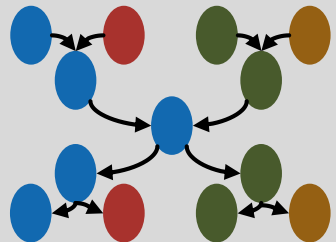
Allreduce

Parameter Server



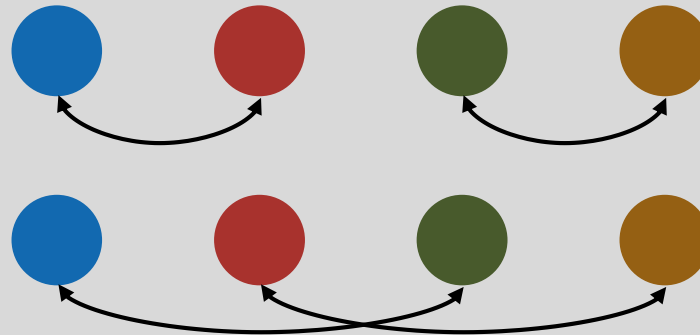
$$2p\alpha + 2pn\beta + pn\gamma$$

Tree



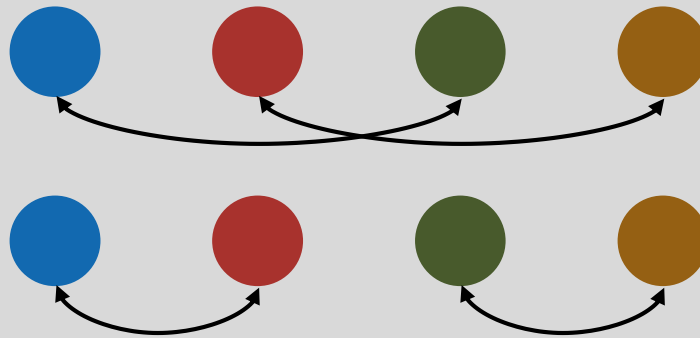
$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

Butterfly (doubling)



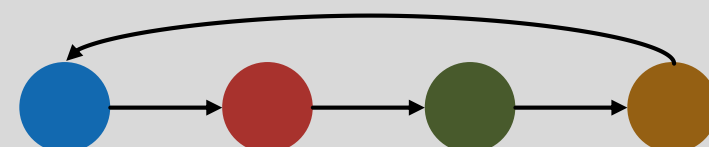
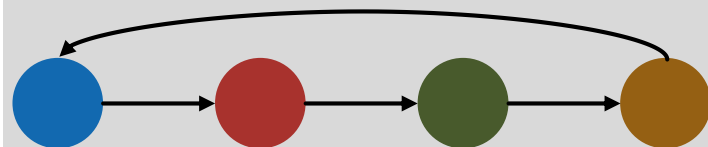
$$\alpha \lg p + \beta n \lg p + \gamma n \lg p$$

Rabenseifner (half/double)



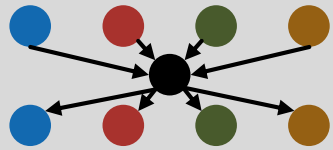
$$2\alpha \lg p + 2 \frac{p-1}{p} n\beta + \frac{p-1}{p} n\gamma$$

Ring



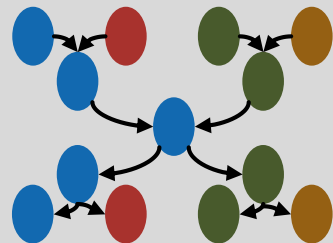
Allreduce

Parameter Server



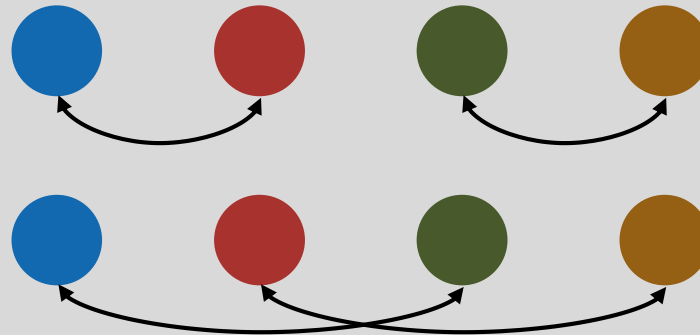
$$2p\alpha + 2pn\beta + pn\gamma$$

Tree

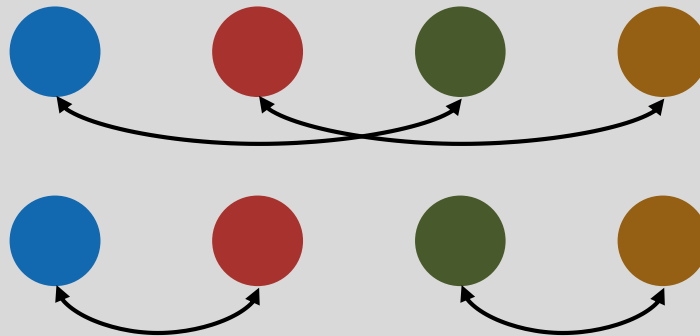


$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

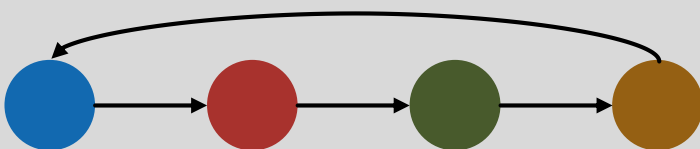
Butterfly (doubling)



Rabenseifner (half/double)

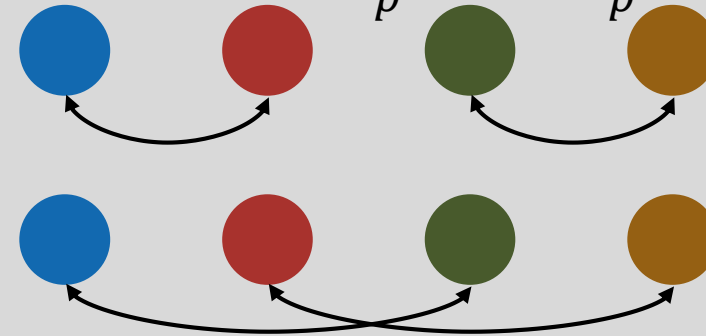


Ring

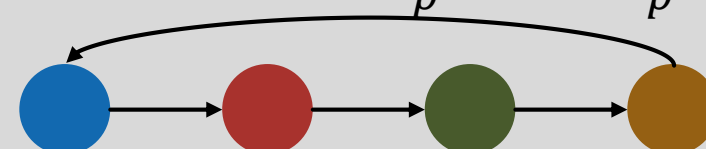


$$\alpha \lg p + \beta n \lg p + \gamma n \lg p$$

$$2\alpha \lg p + 2 \frac{p-1}{p} n\beta + \frac{p-1}{p} n\gamma$$

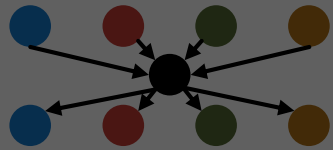


$$2(p-1)\alpha + 2 \frac{p-1}{p} n\beta + \frac{p-1}{p} n\gamma$$



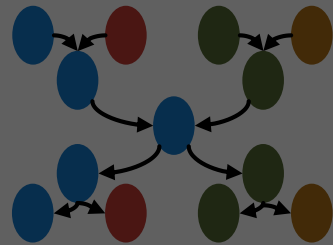
Allreduce

Parameter Server



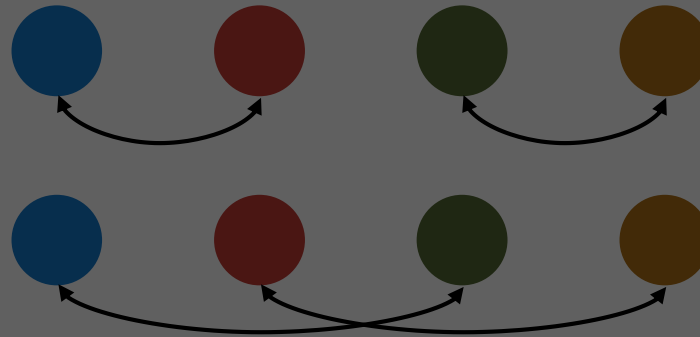
$$2p\alpha + 2pn\beta + pn\gamma$$

Tree

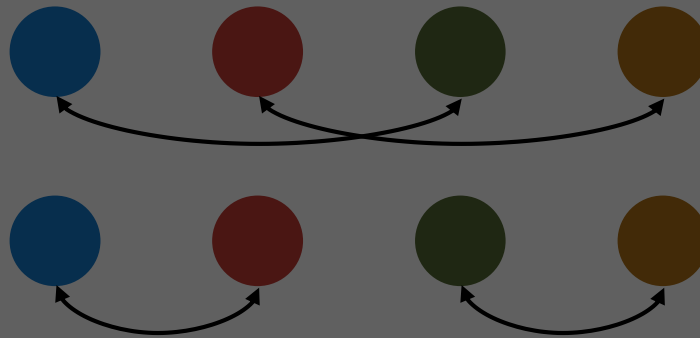


$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

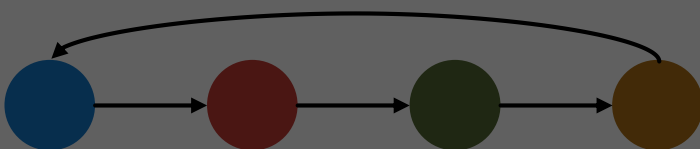
Butterfly (doubling)



Rabenseifner (half/double)



Ring



$$\alpha \lg p + \beta n \lg p + \gamma n \lg p$$

$$2\alpha \lg p + 2 \frac{p-1}{p} n\beta + \frac{p-1}{p} n\gamma$$

$$2(p-1)\alpha + 2 \frac{p-1}{p} n\beta + \frac{p-1}{p} n\gamma$$

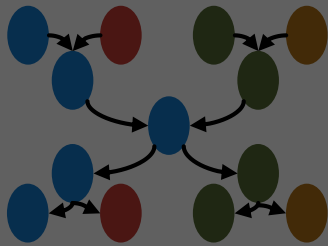
Allreduce

Parameter Server



$$2p\alpha + 2pn\beta + pn\gamma$$

Tree



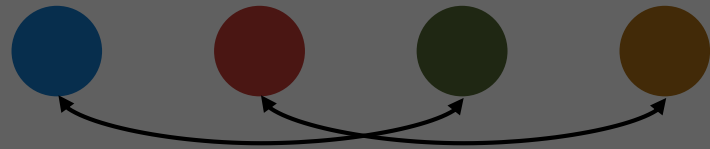
$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$

Butterfly (doubling)

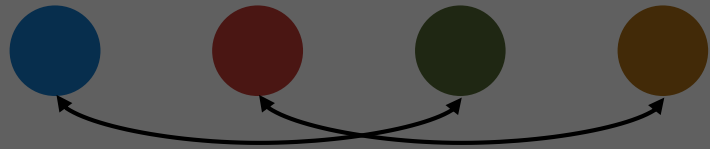


$$\alpha \lg p + \beta n \lg p + \gamma n \lg p$$

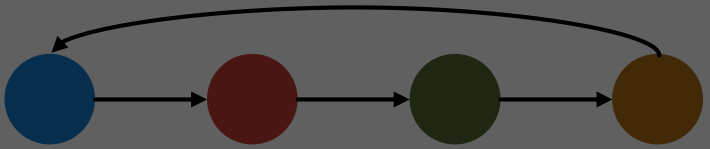
Implementation matters!



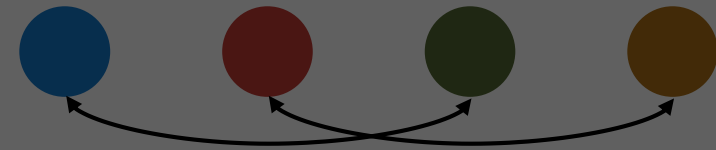
Rabenseifner (half/double)



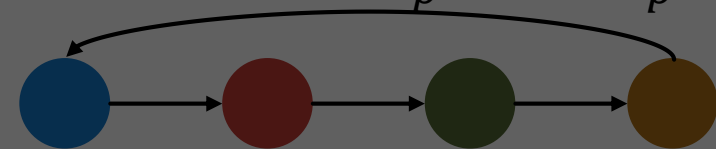
Ring



$$2\alpha \lg p + 2 \frac{p-1}{p} n\beta + \frac{p-1}{p} n\gamma$$



$$2(p-1)\alpha + 2 \frac{p-1}{p} n\beta + \frac{p-1}{p} n\gamma$$



Allreduce

Parameter Server

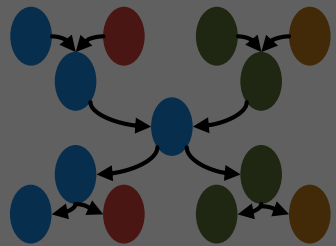
Butterfly (doubling)

$$\alpha \lg p + \beta n \lg p + \gamma n \lg p$$

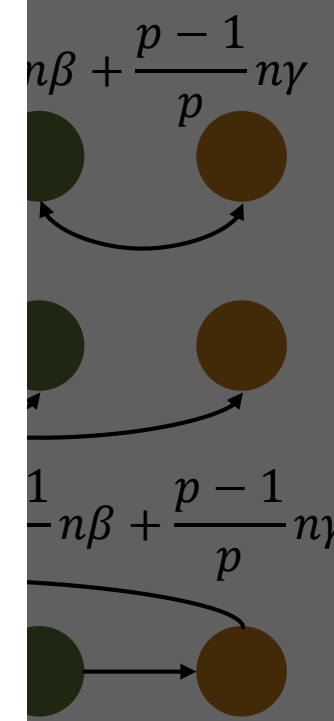
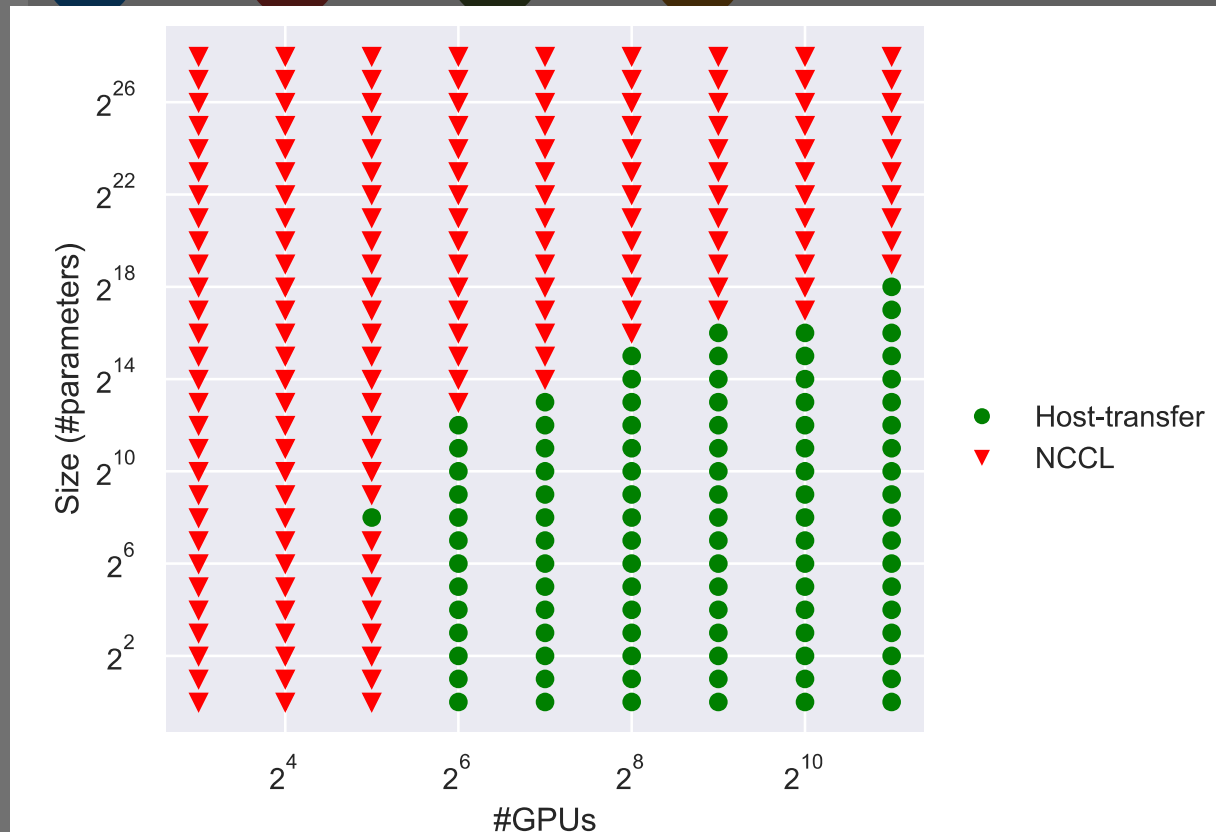
Implementation matters!

$$2p\alpha + 2pn\beta + pn\gamma$$

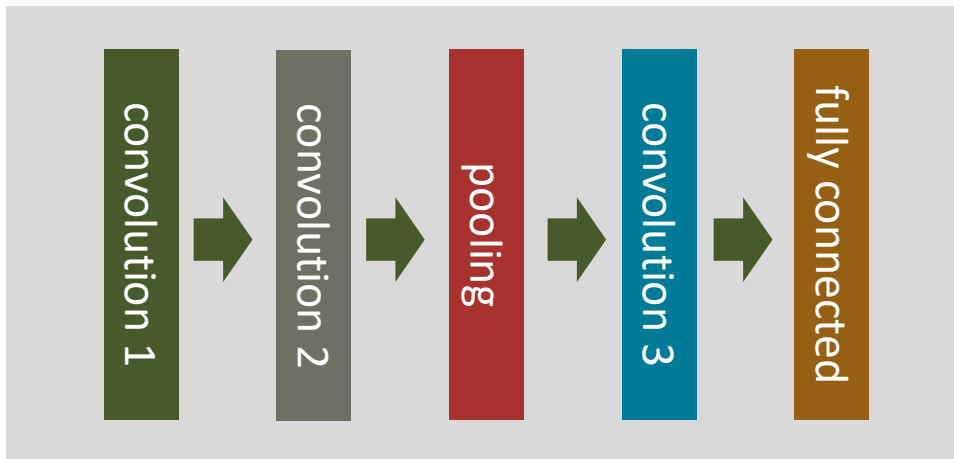
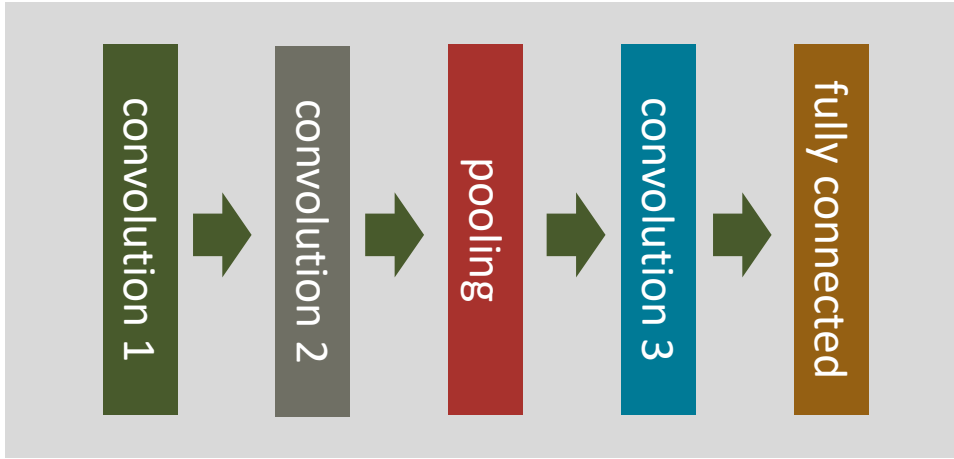
Tree



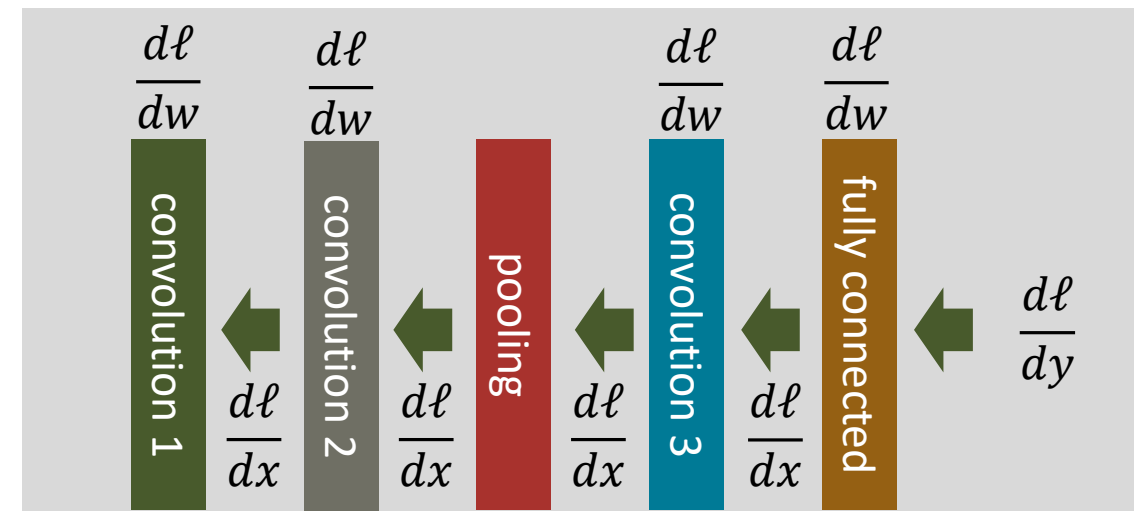
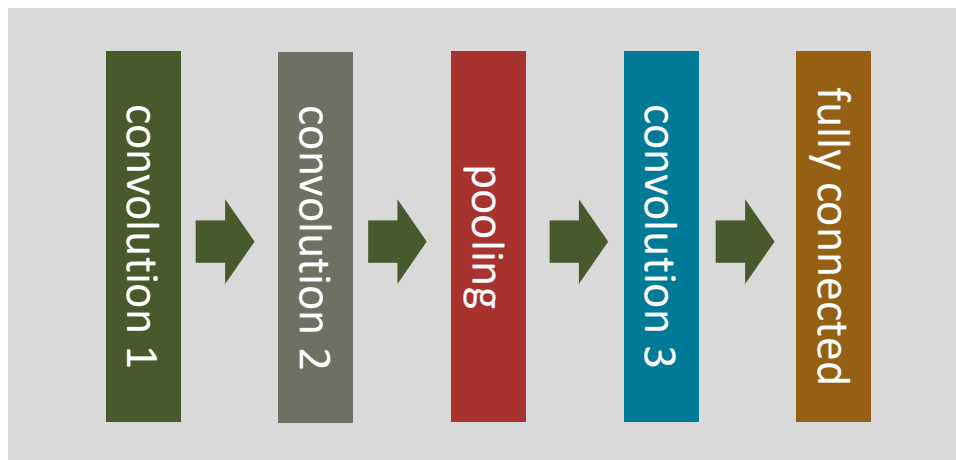
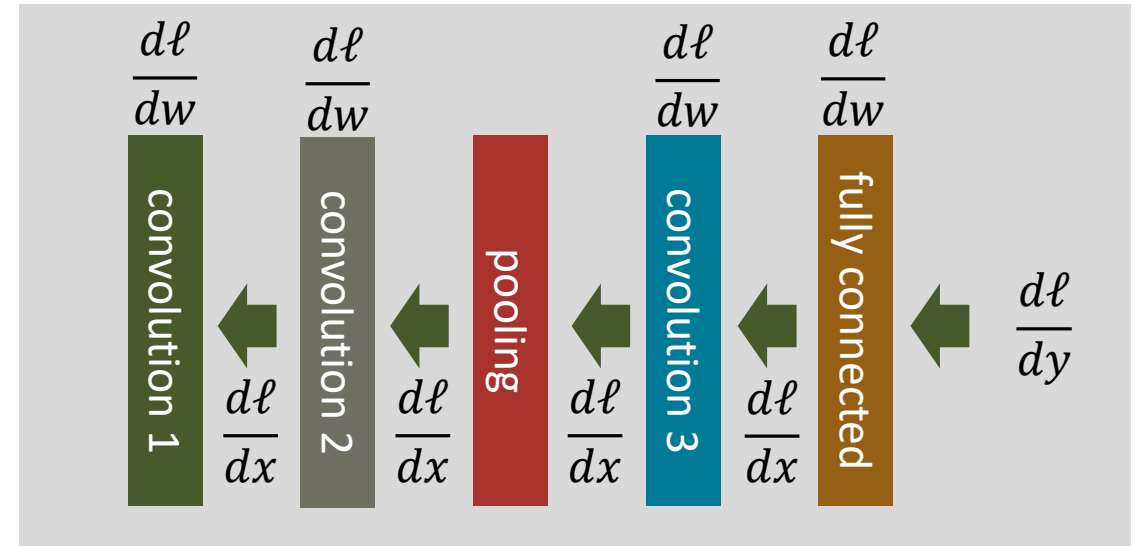
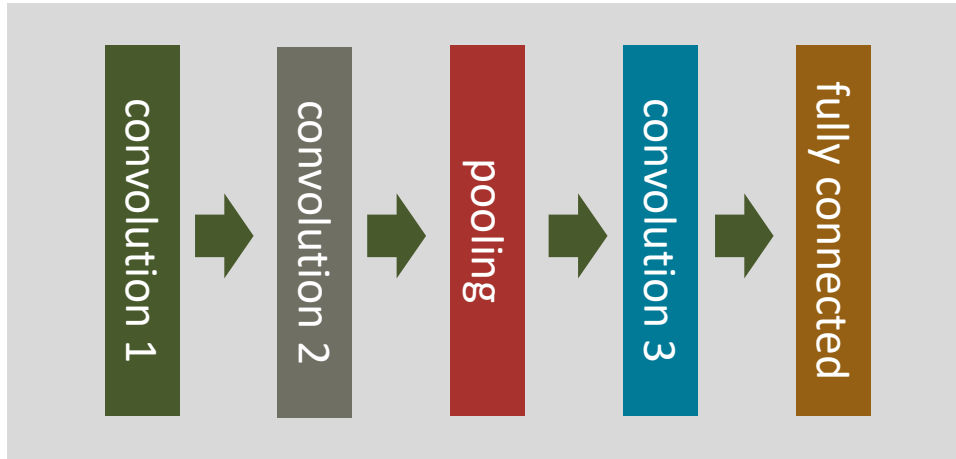
$$2\alpha \lg p + 2\beta n \lg p + \gamma n \lg p$$



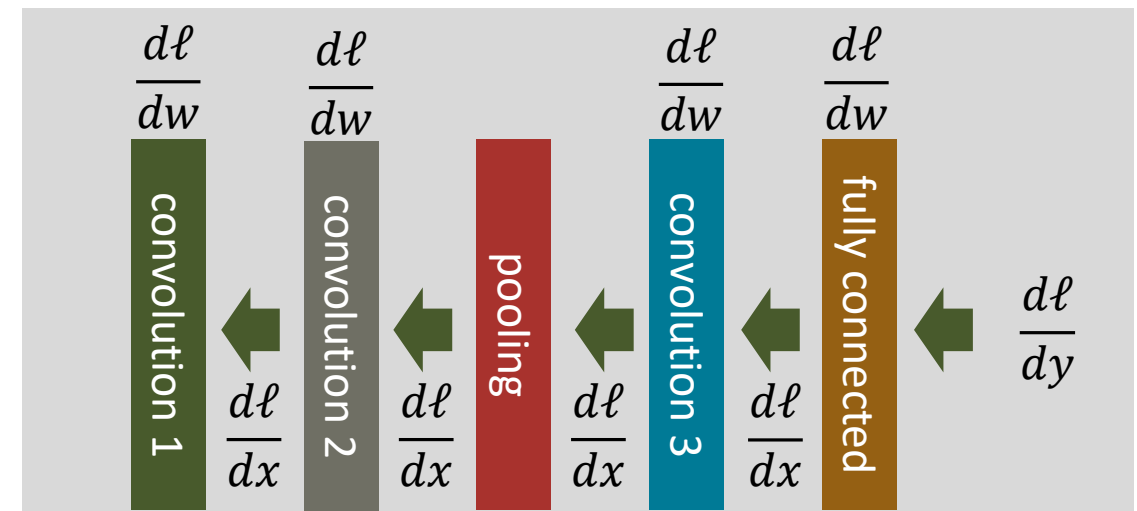
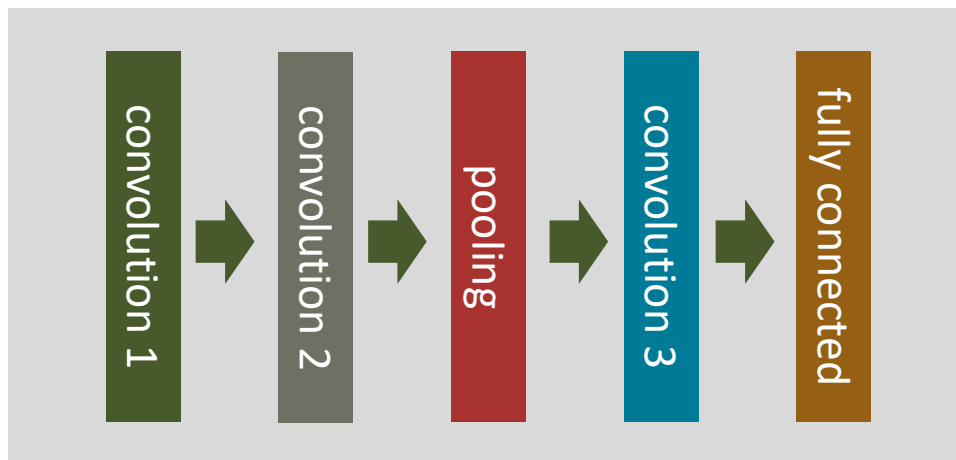
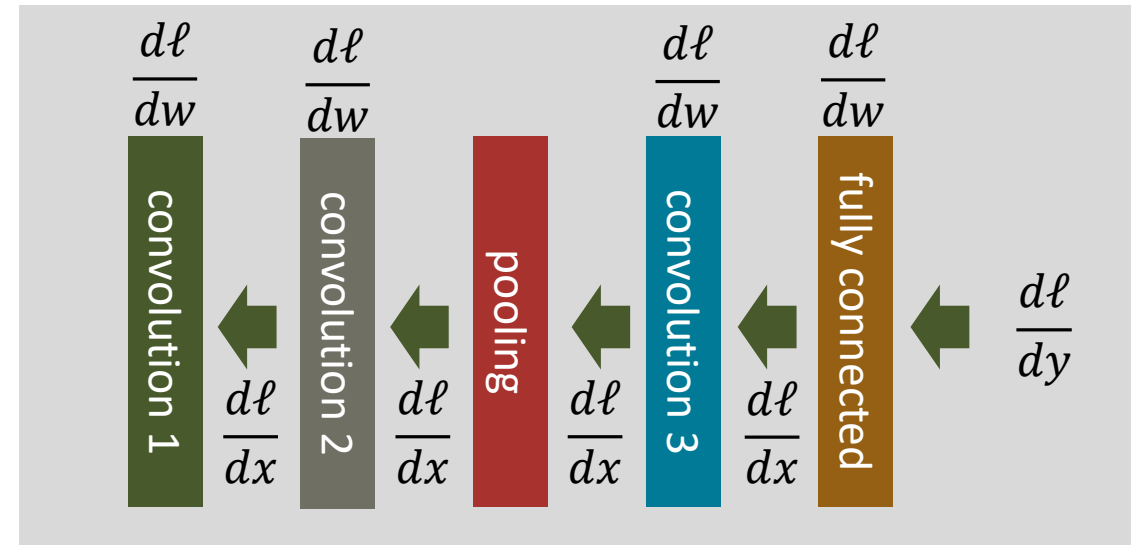
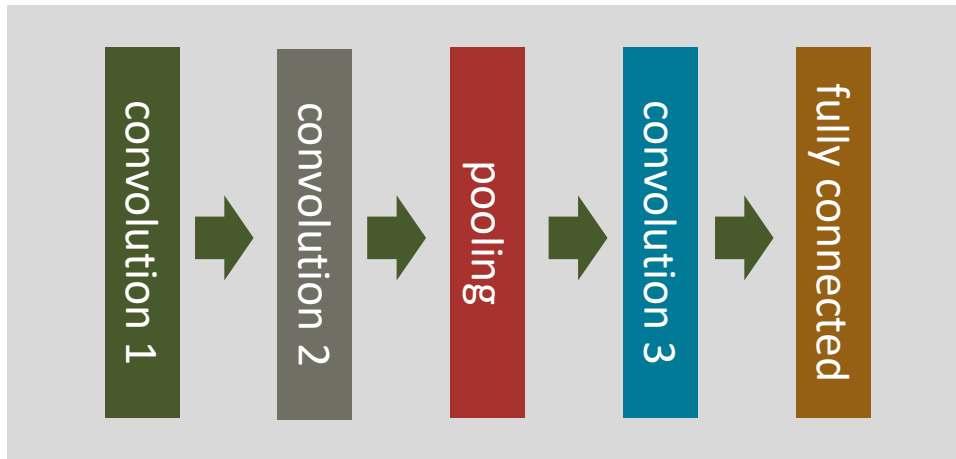
Distributed data-parallelism



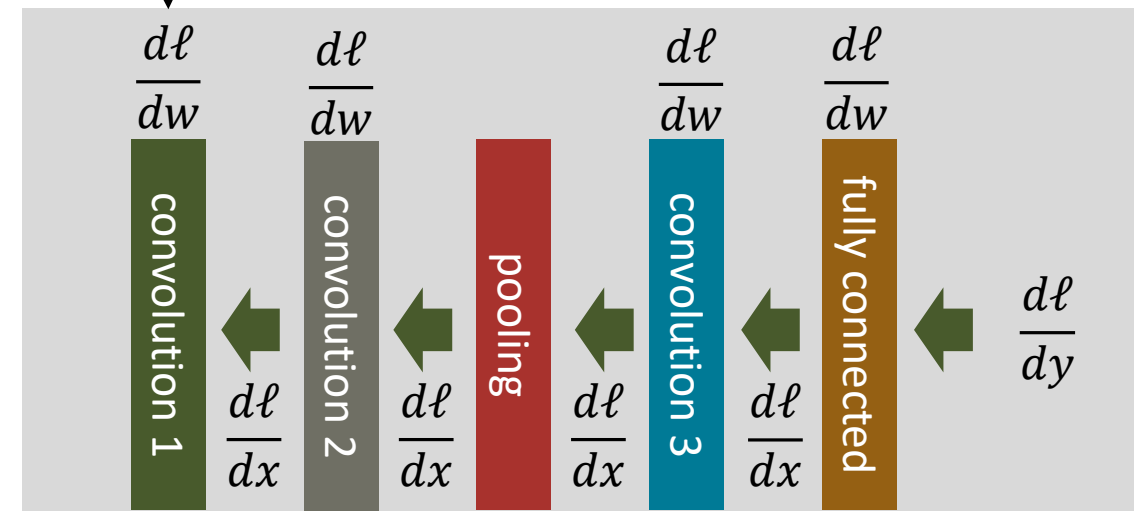
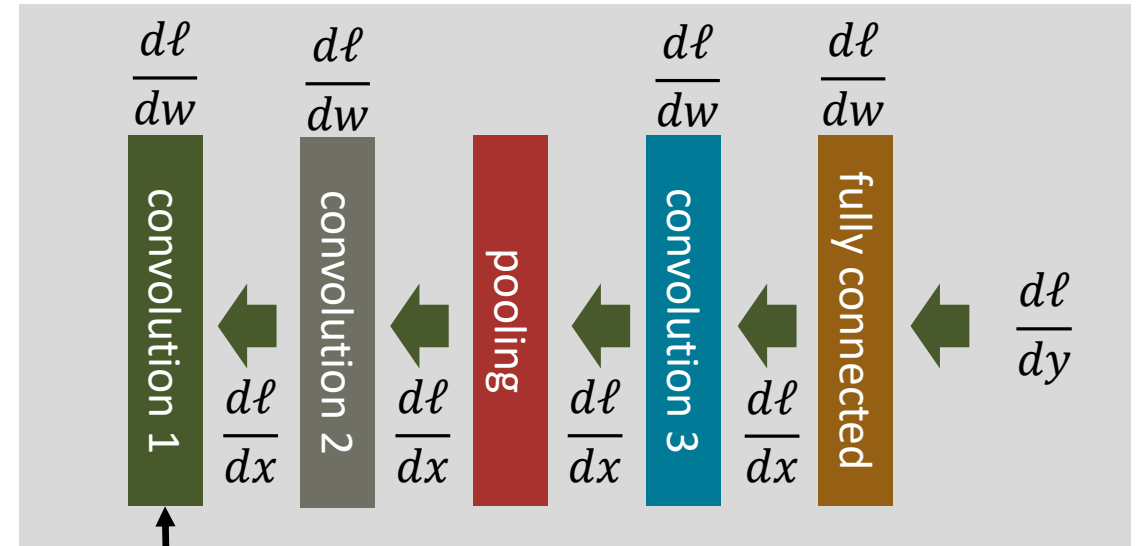
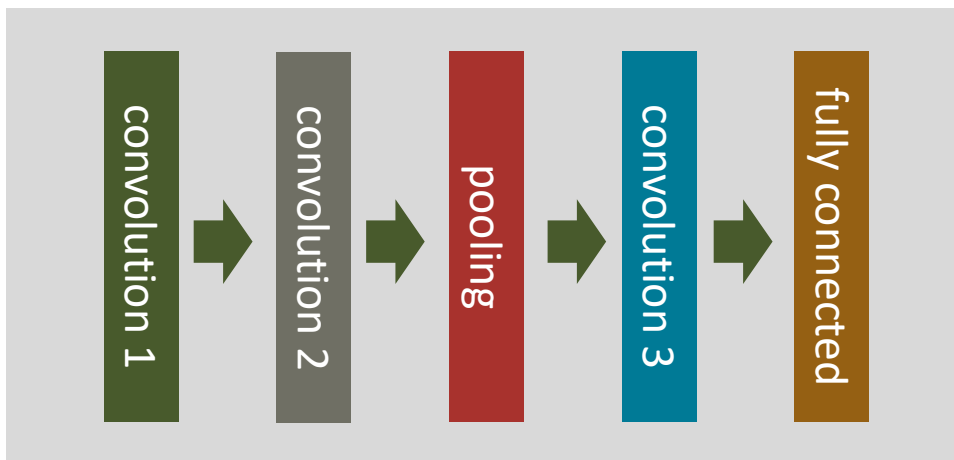
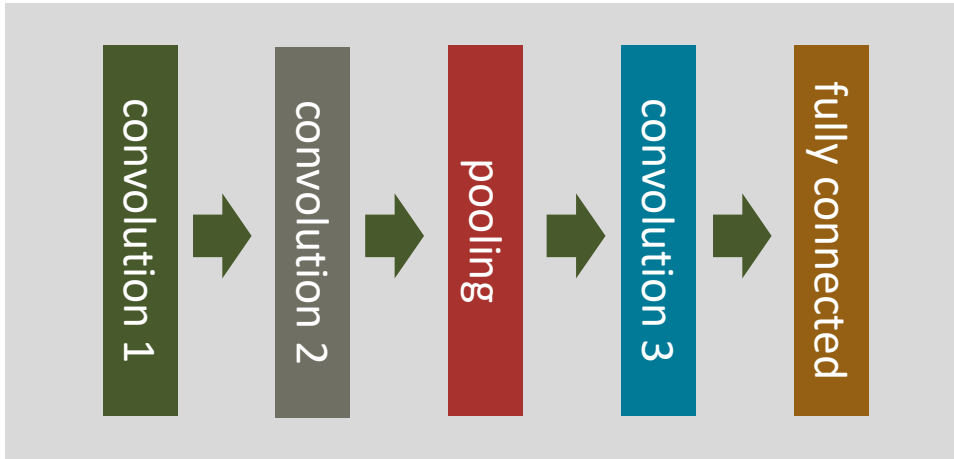
Distributed data-parallelism



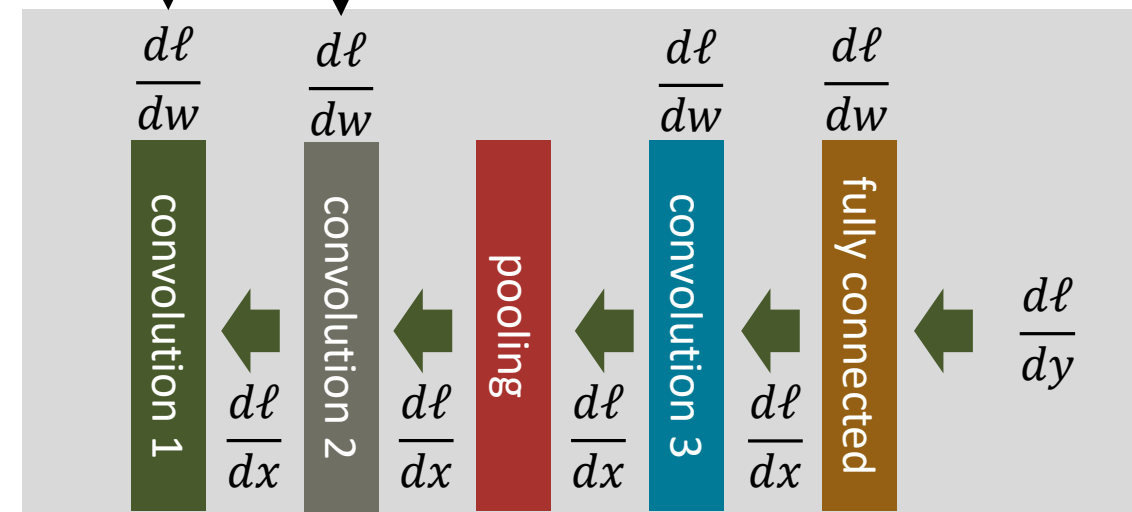
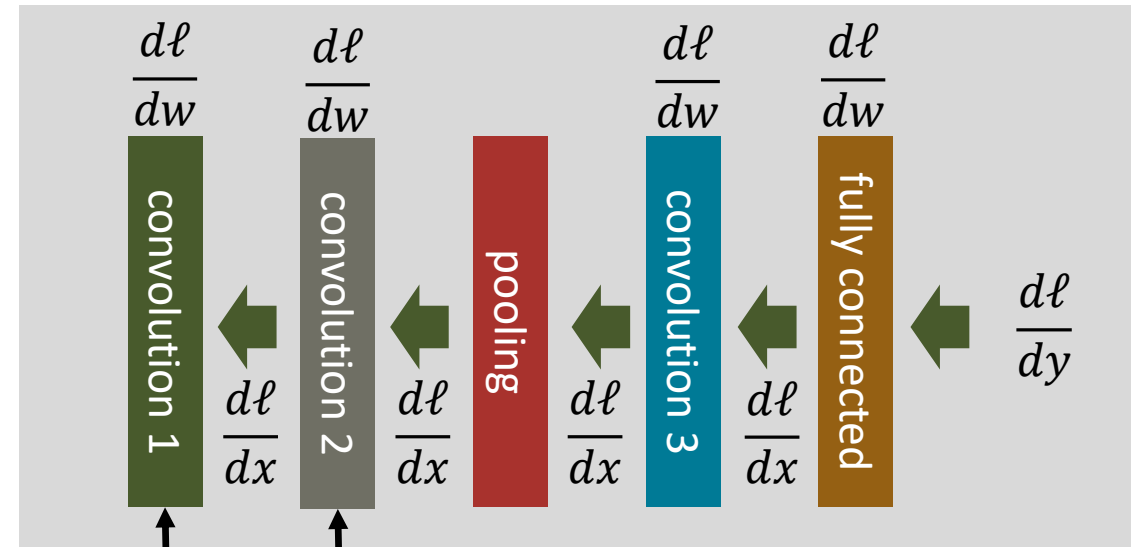
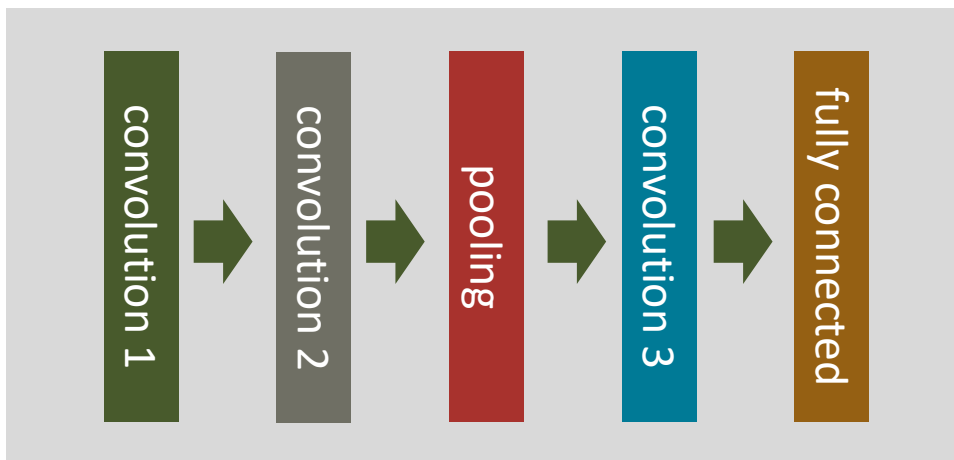
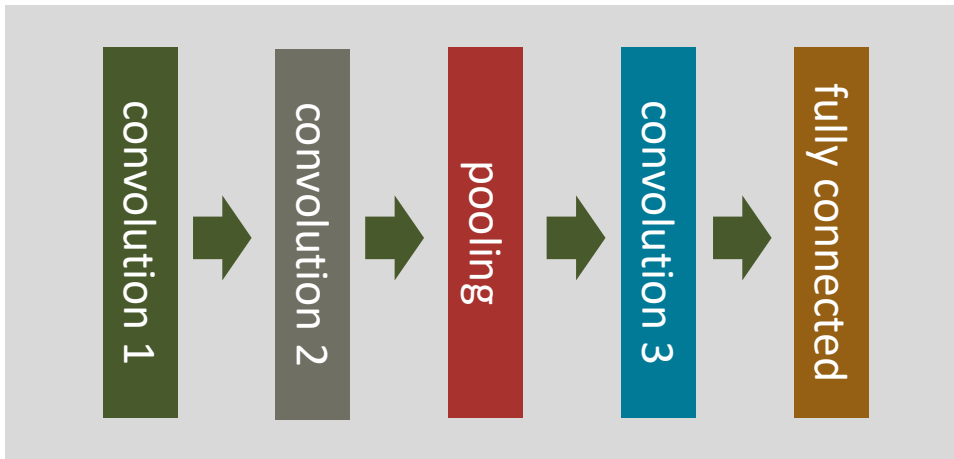
Distributed data-parallelism



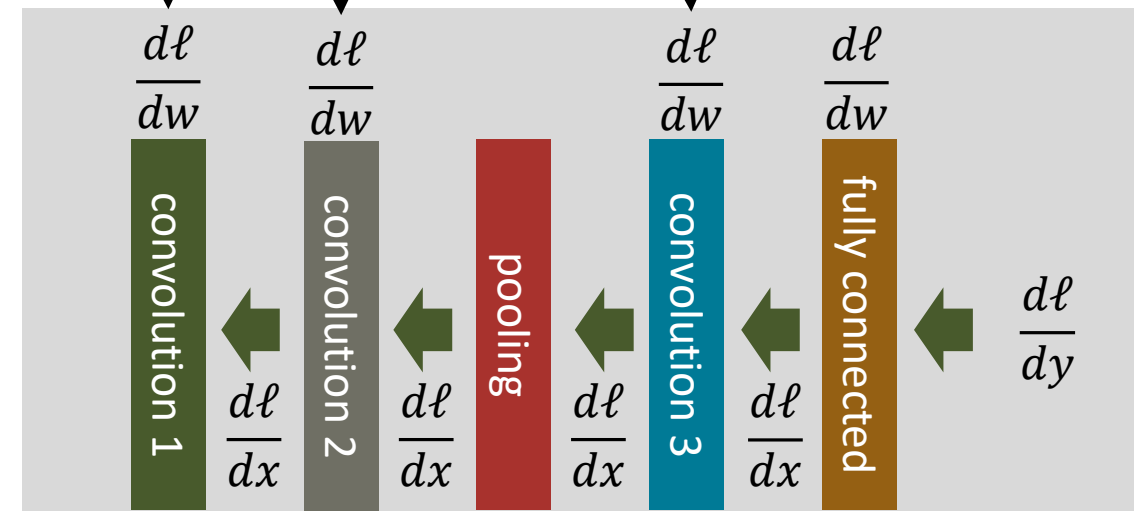
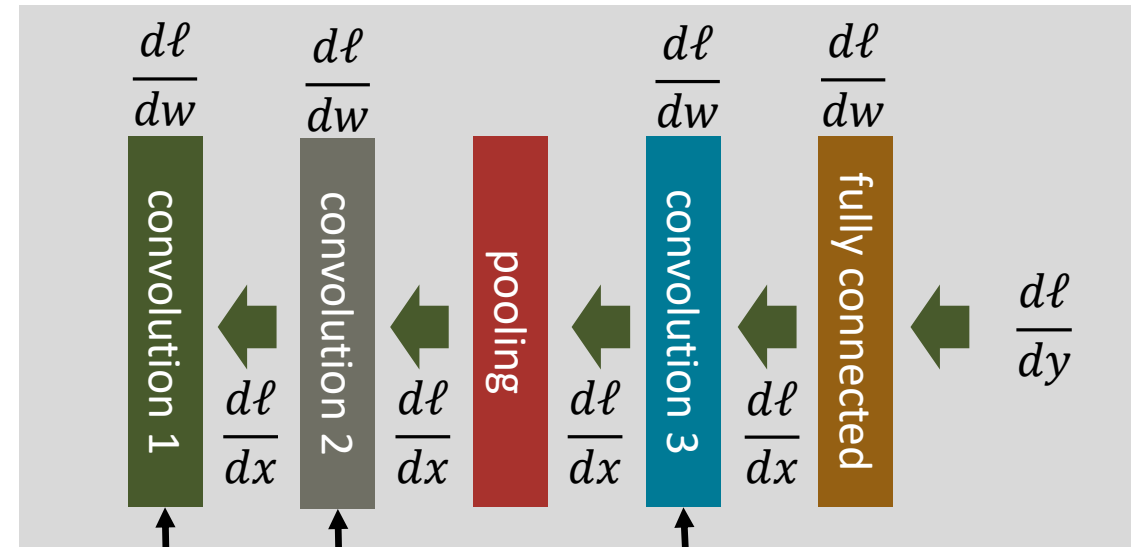
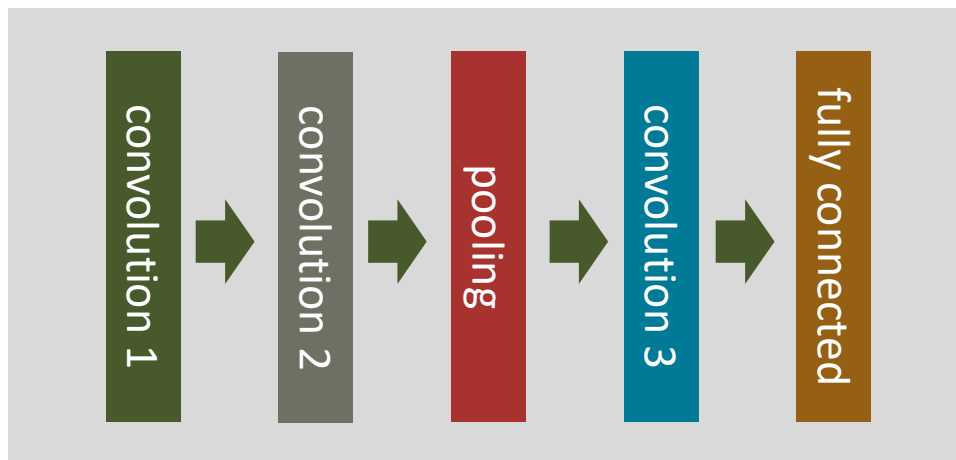
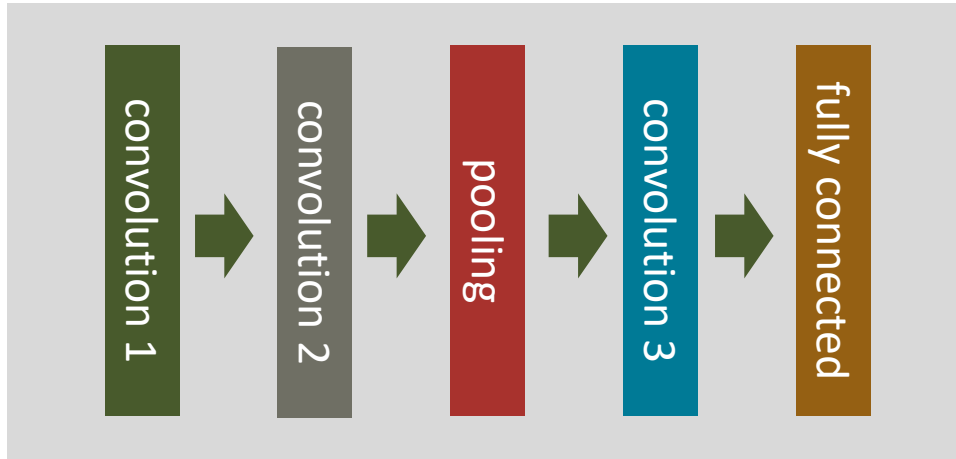
Distributed data-parallelism



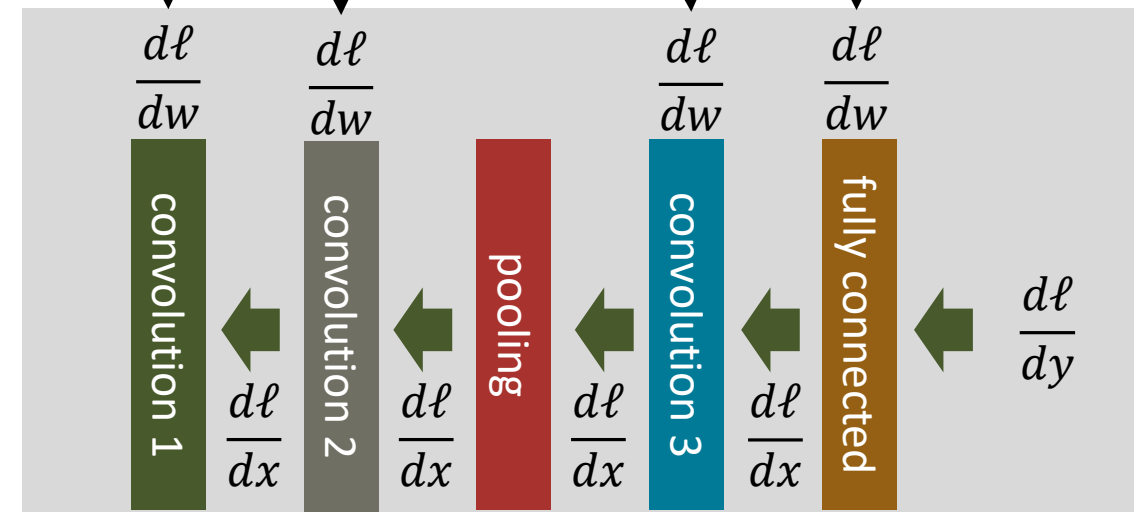
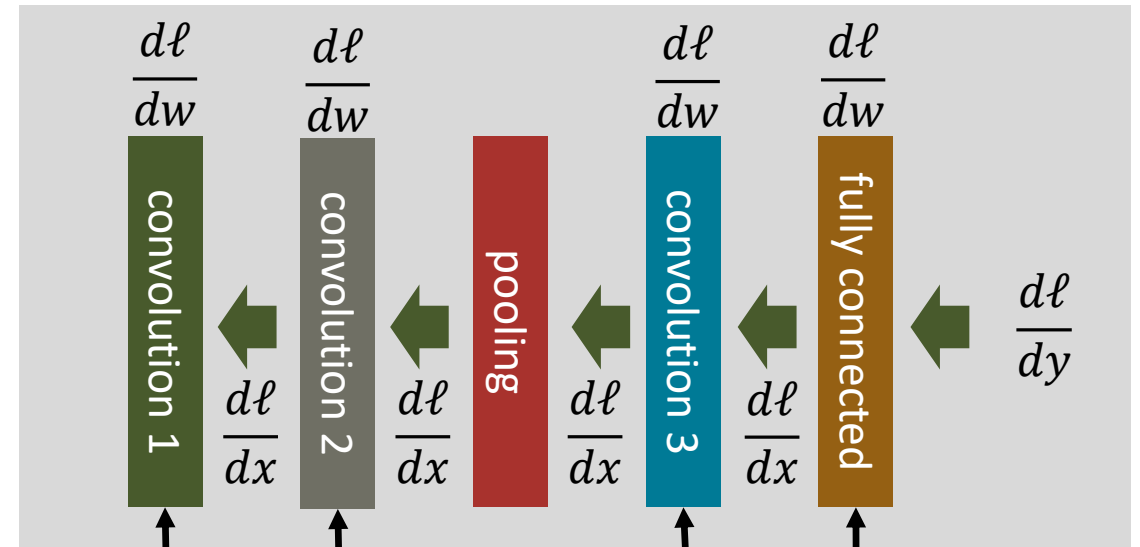
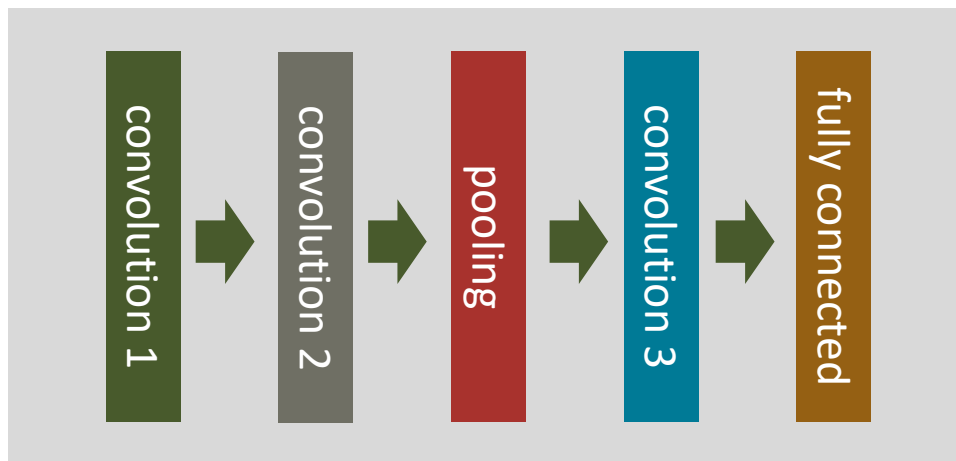
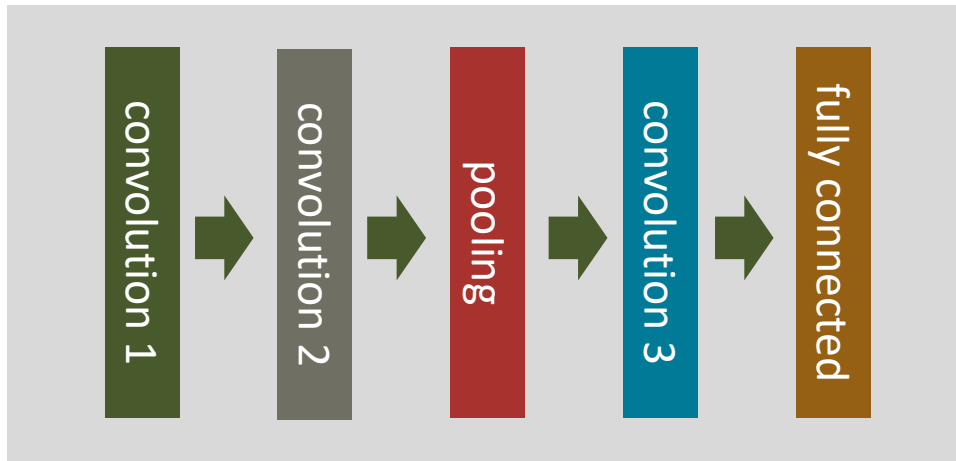
Distributed data-parallelism



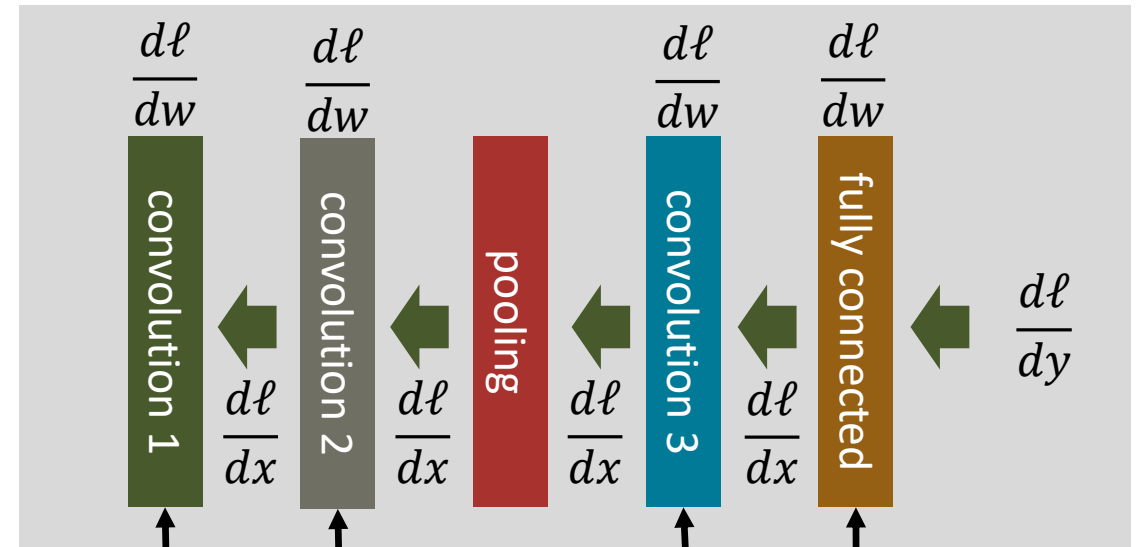
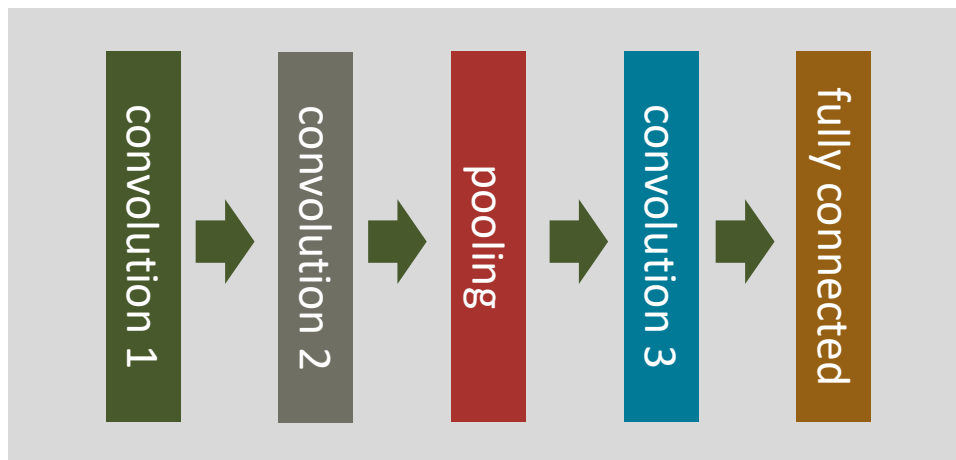
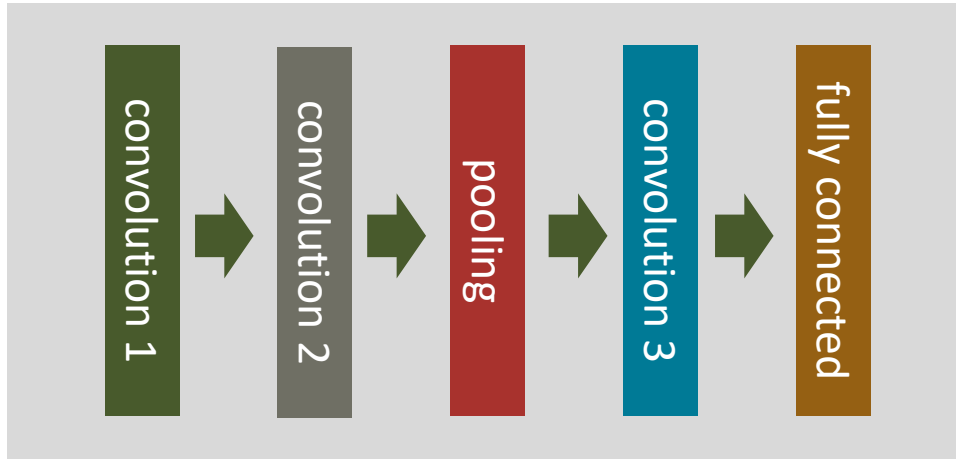
Distributed data-parallelism



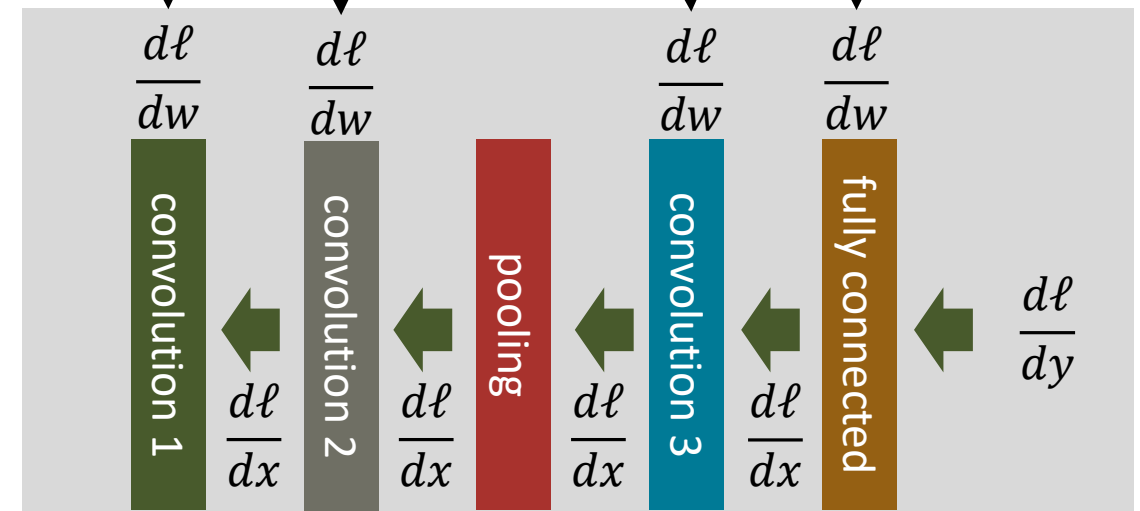
Distributed data-parallelism



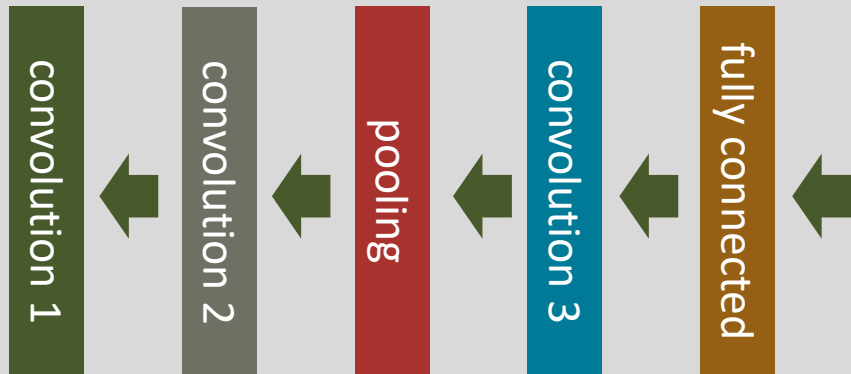
Distributed data-parallelism



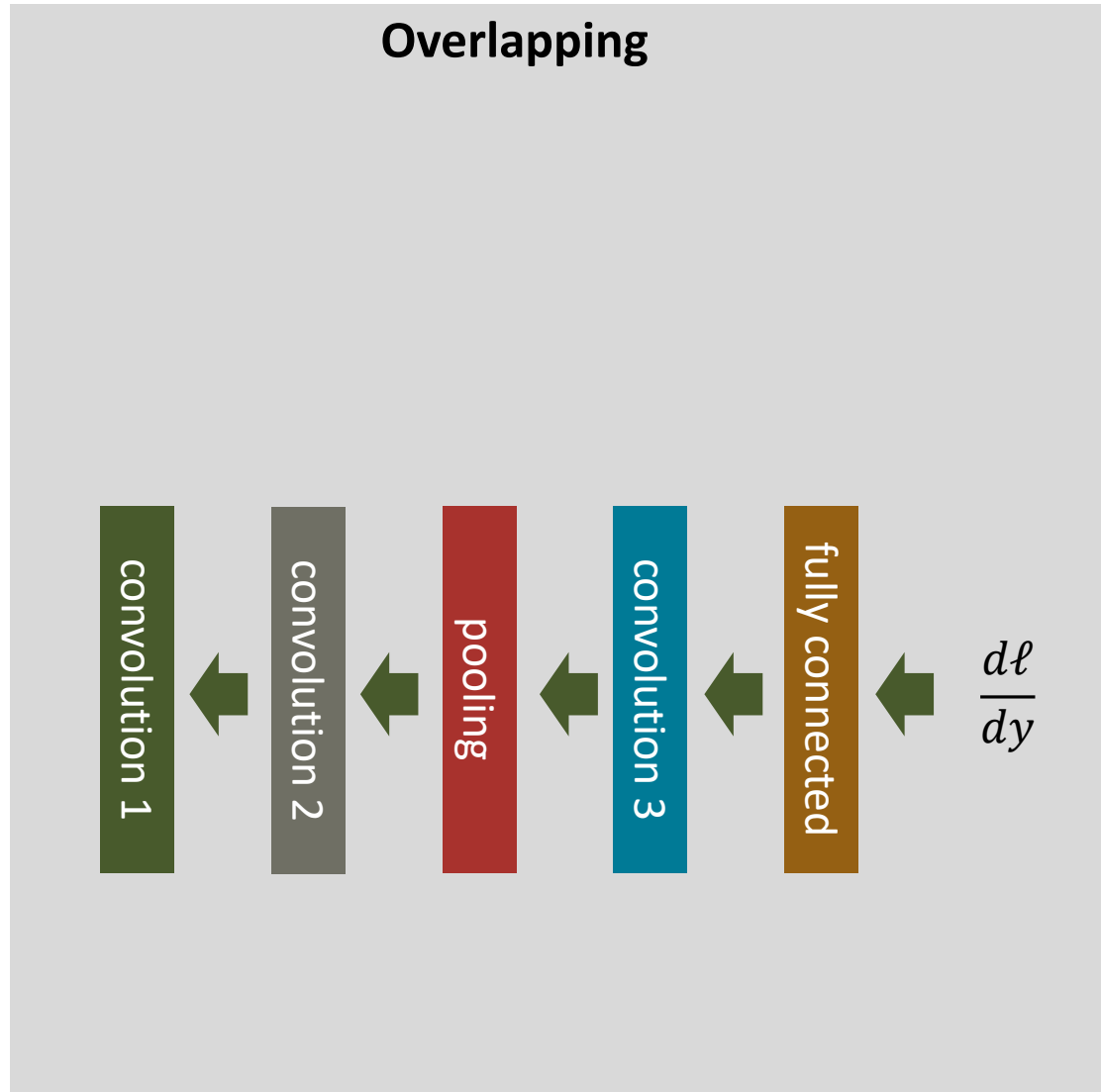
Allreduce!



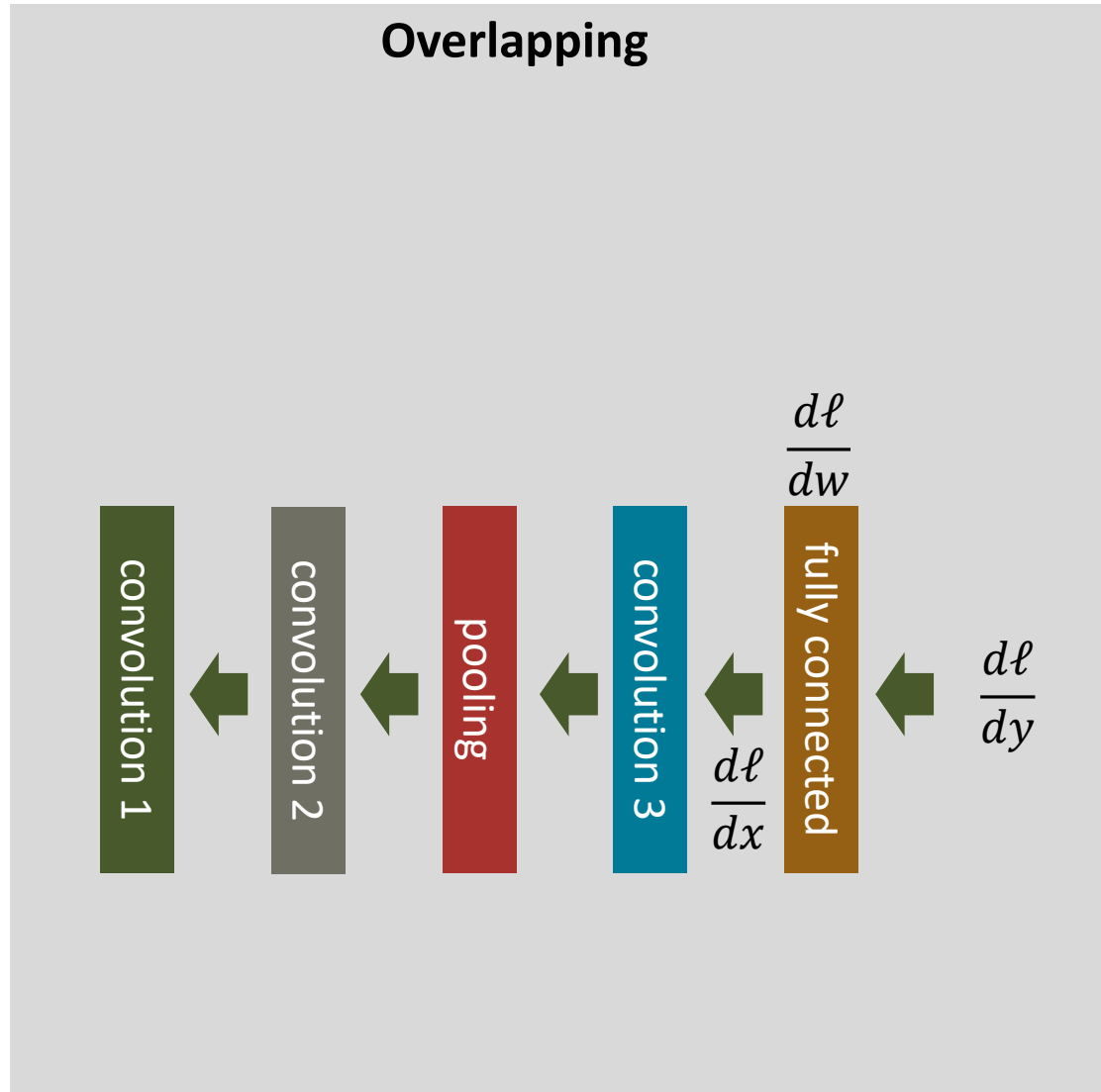
Distributed data-parallelism



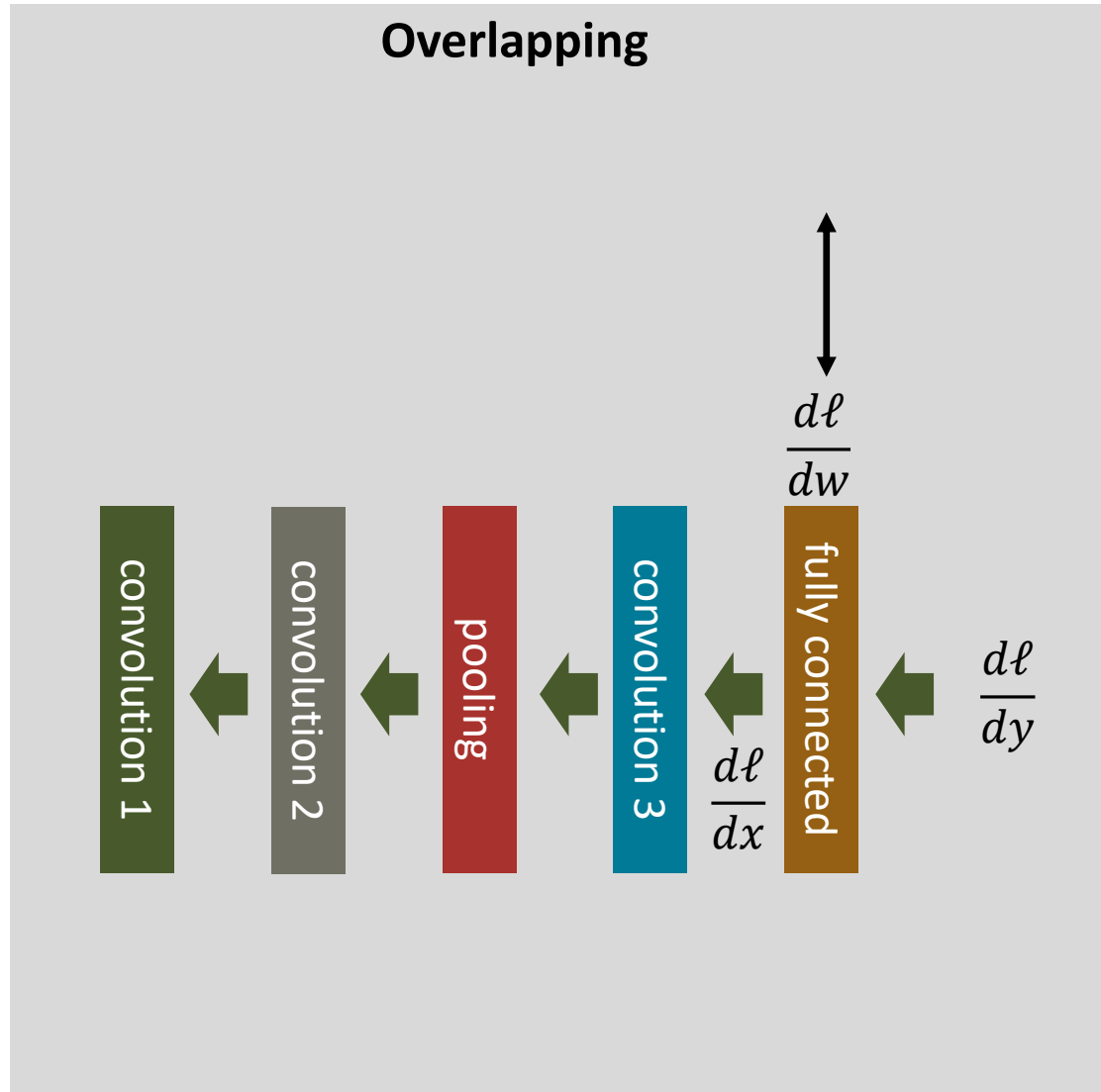
Distributed data-parallelism



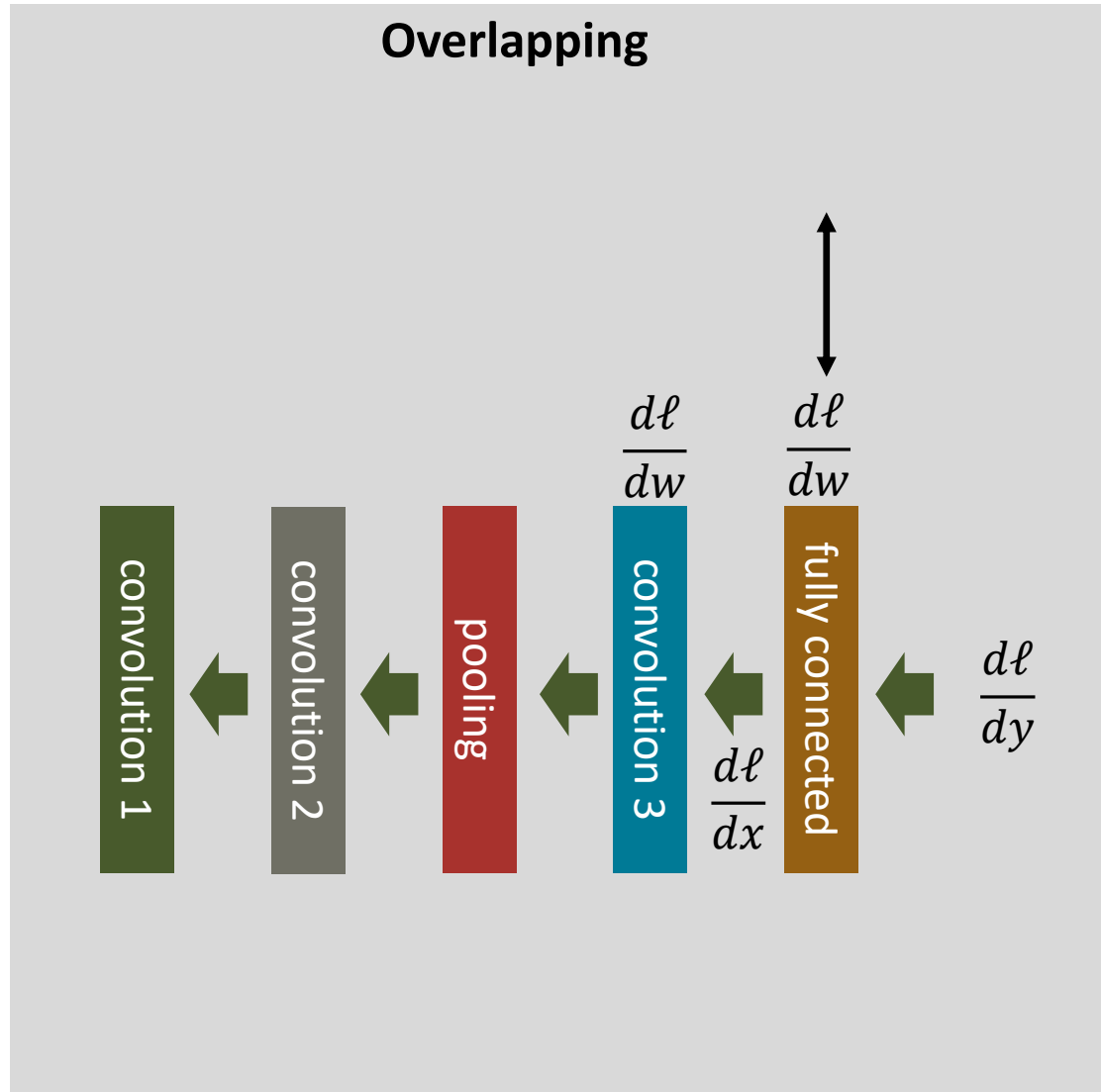
Distributed data-parallelism



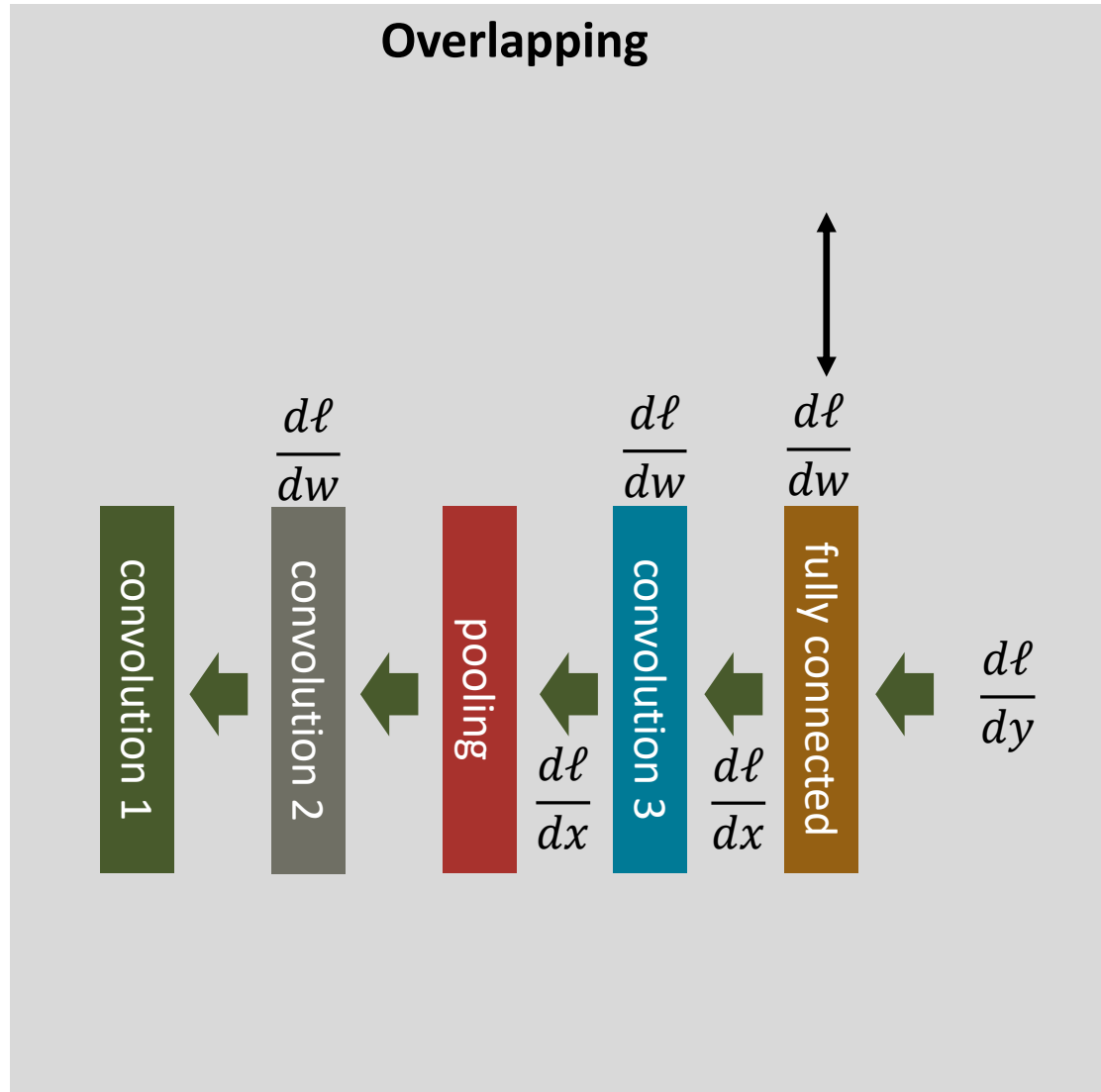
Distributed data-parallelism



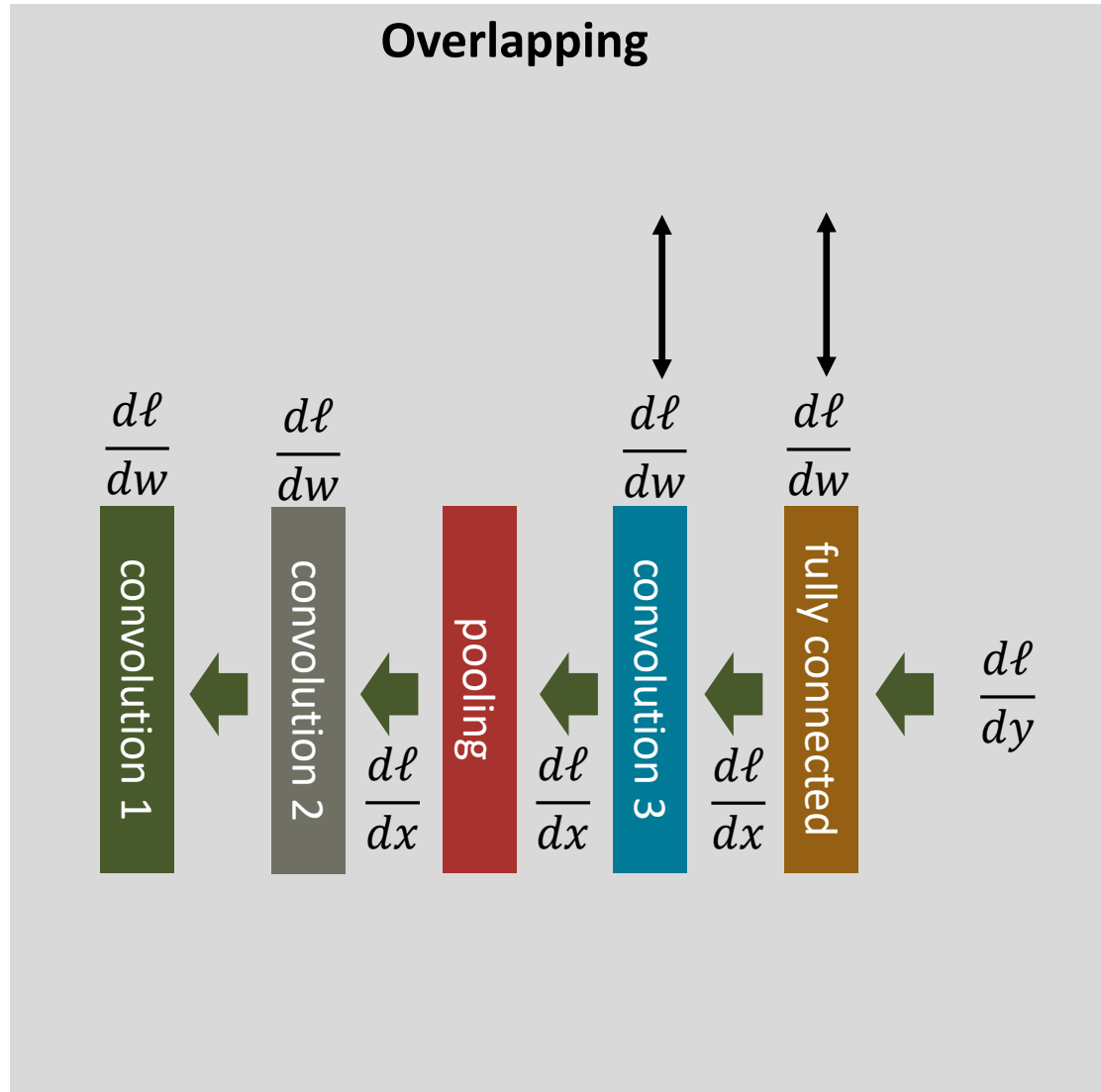
Distributed data-parallelism



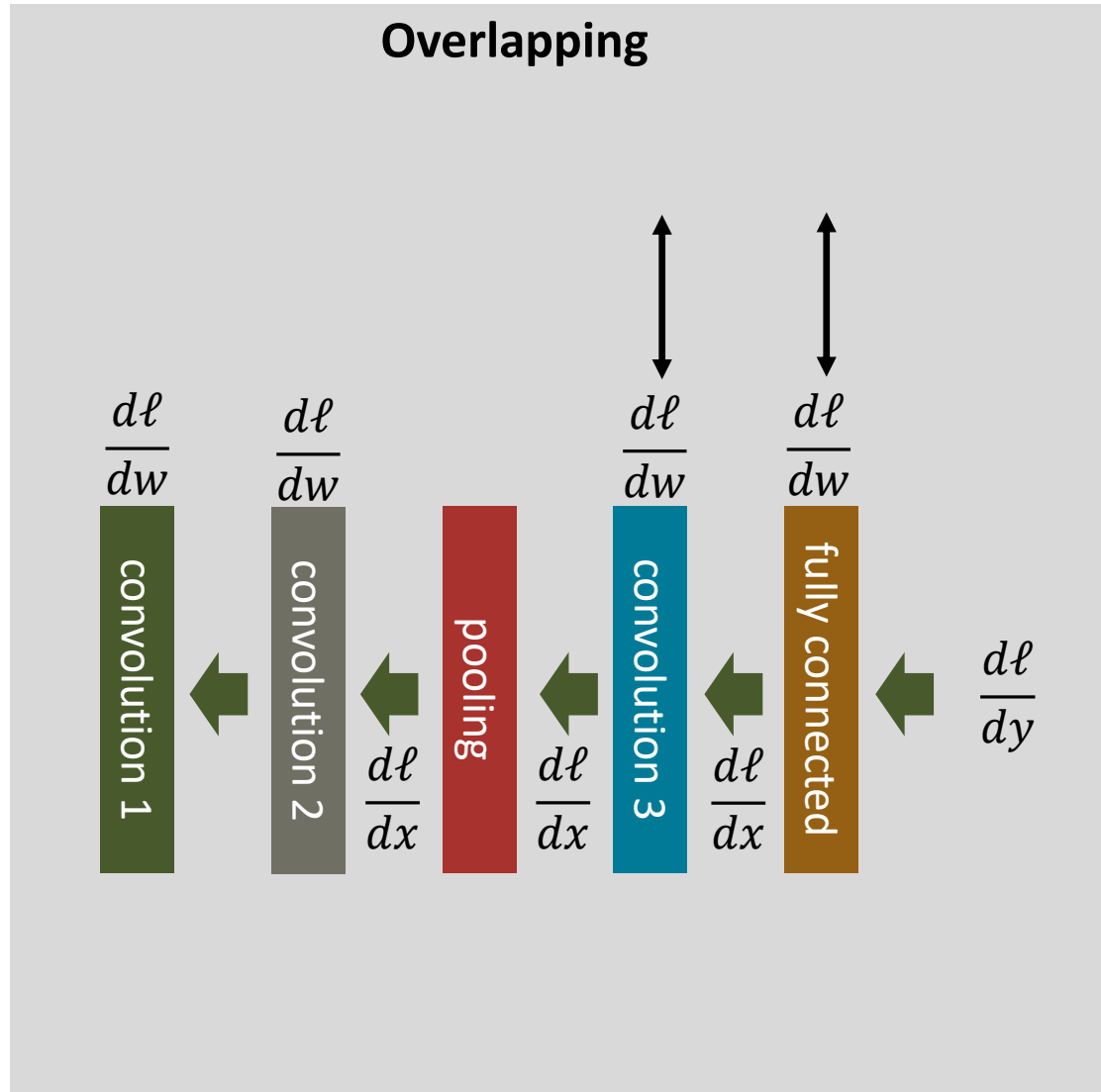
Distributed data-parallelism



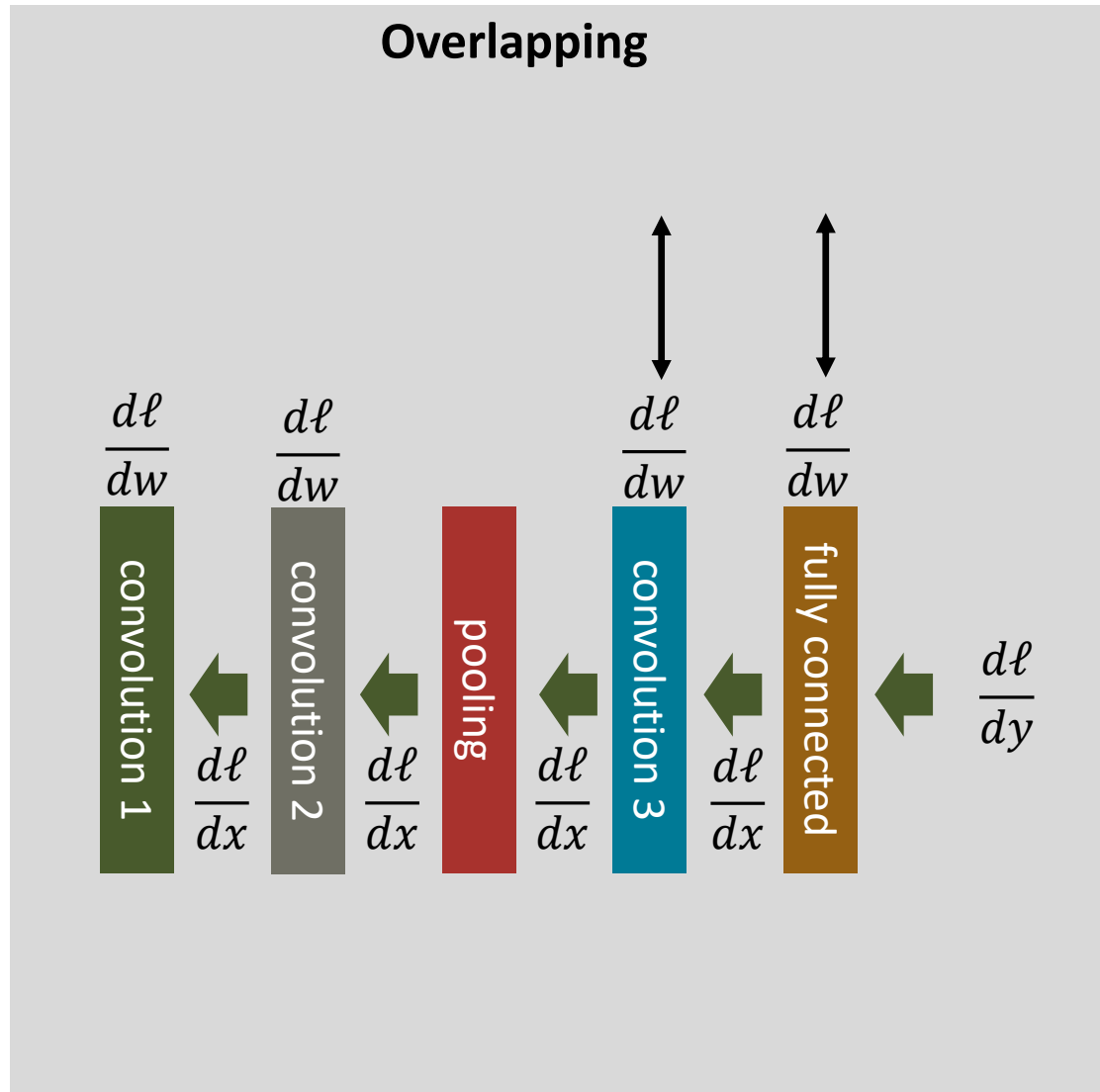
Distributed data-parallelism



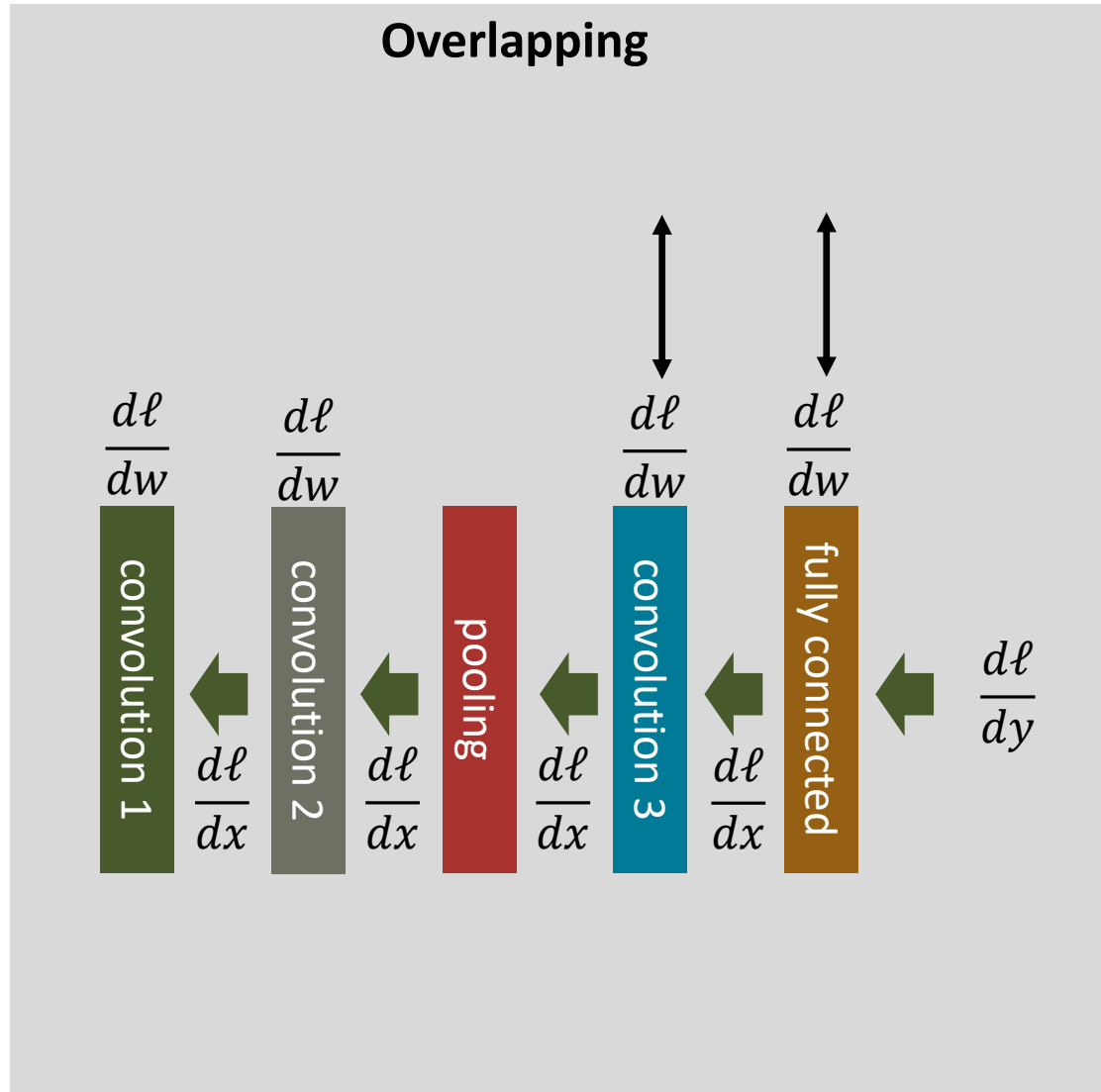
Distributed data-parallelism



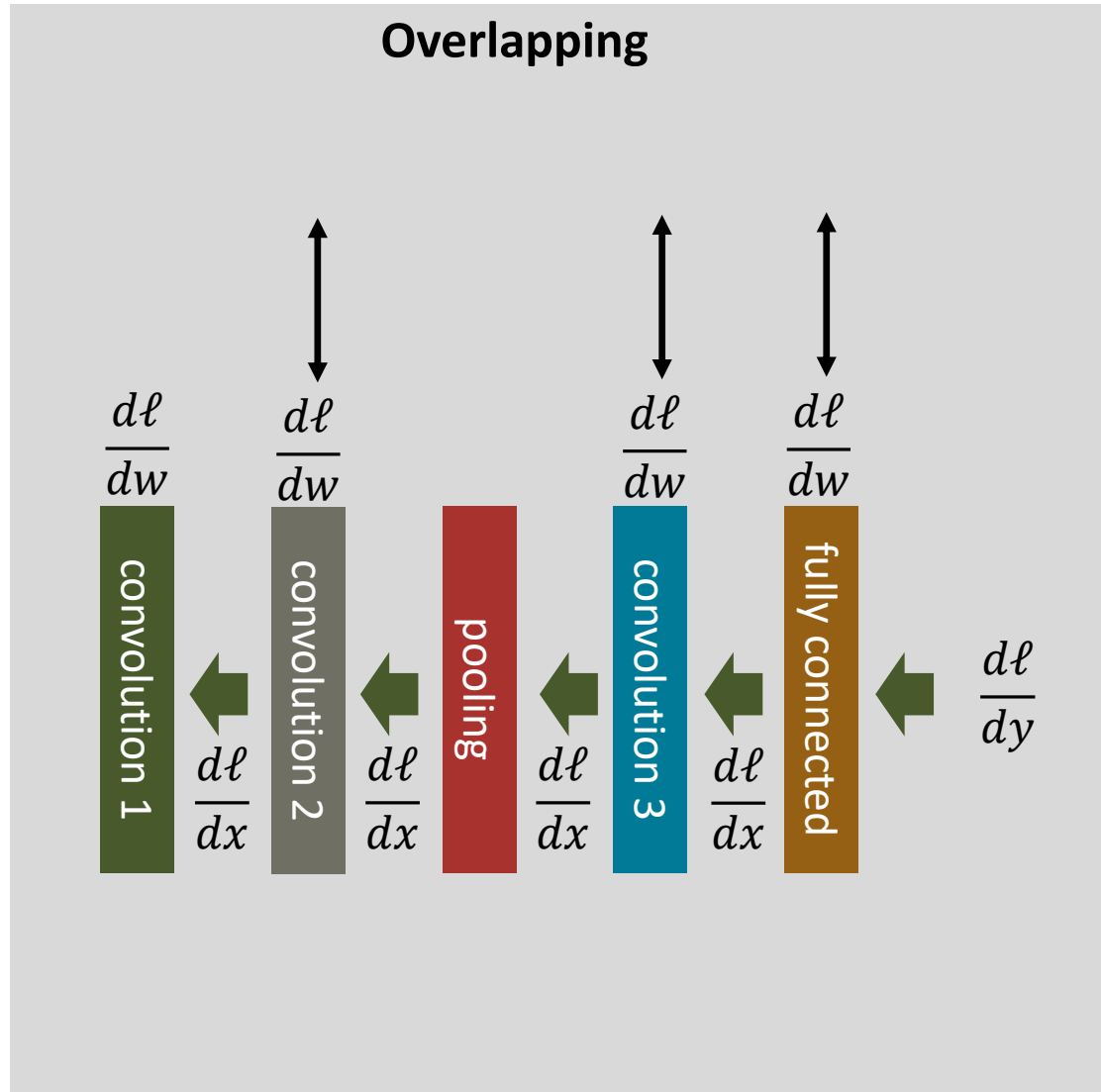
Distributed data-parallelism



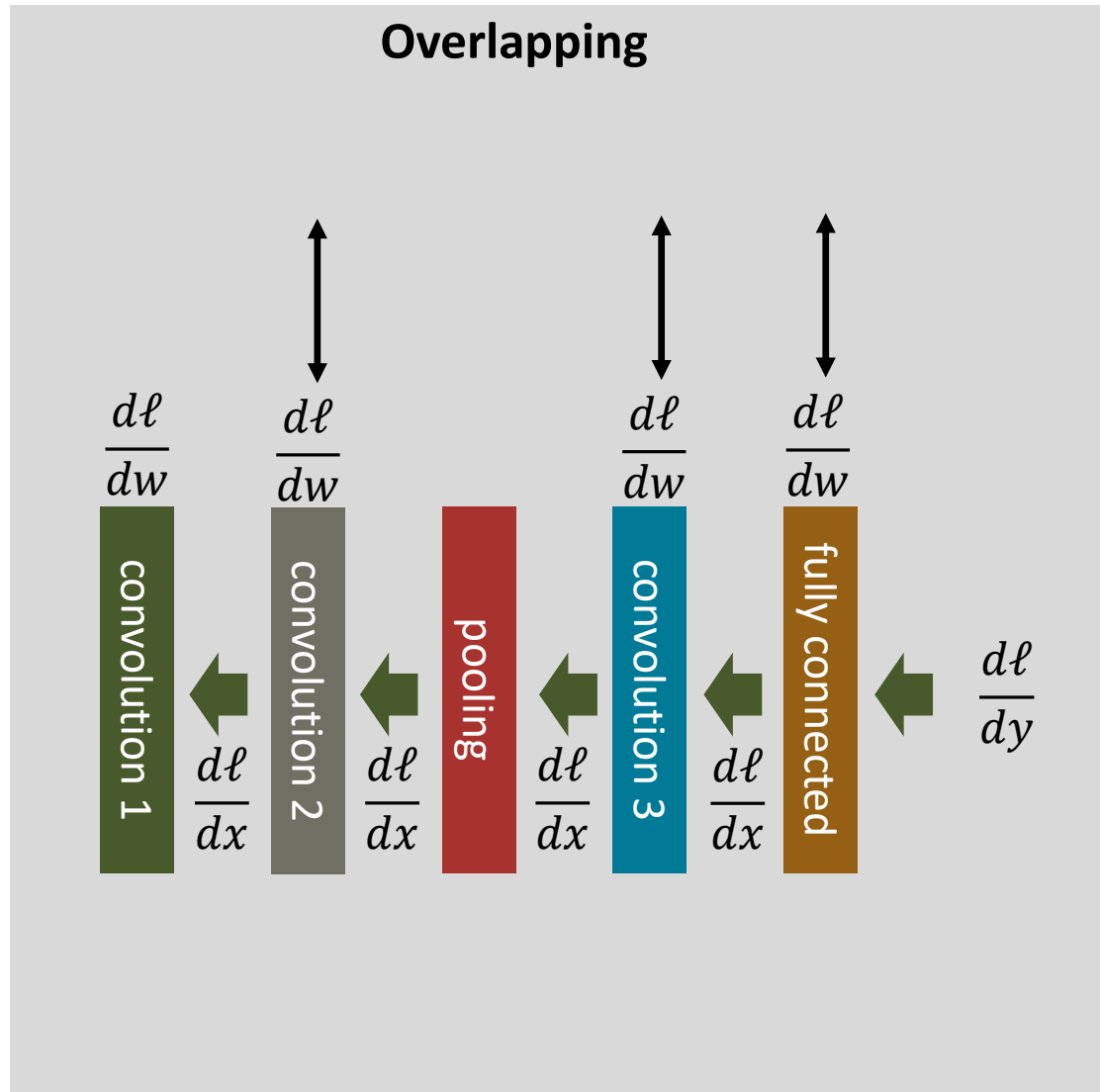
Distributed data-parallelism



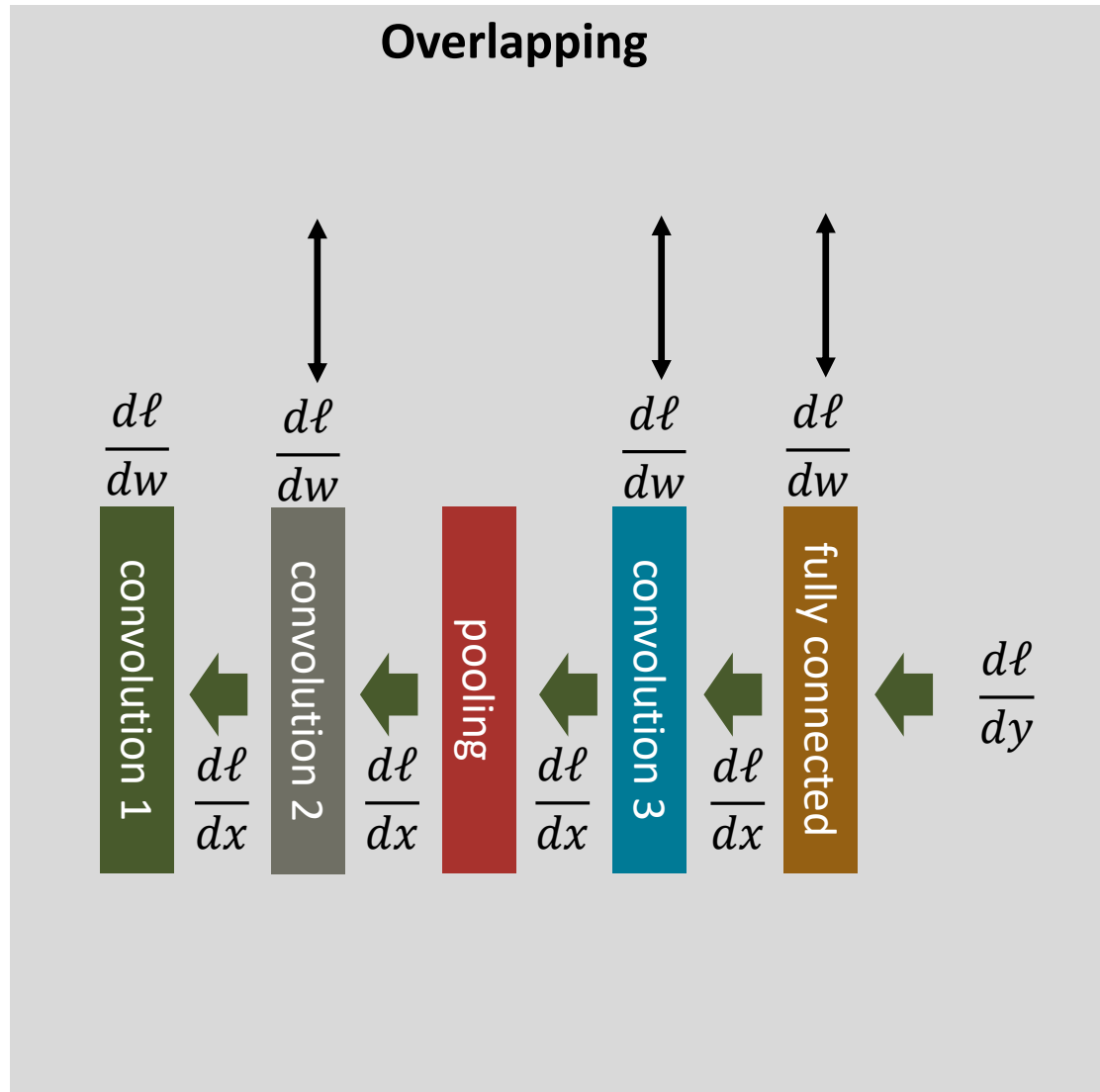
Distributed data-parallelism



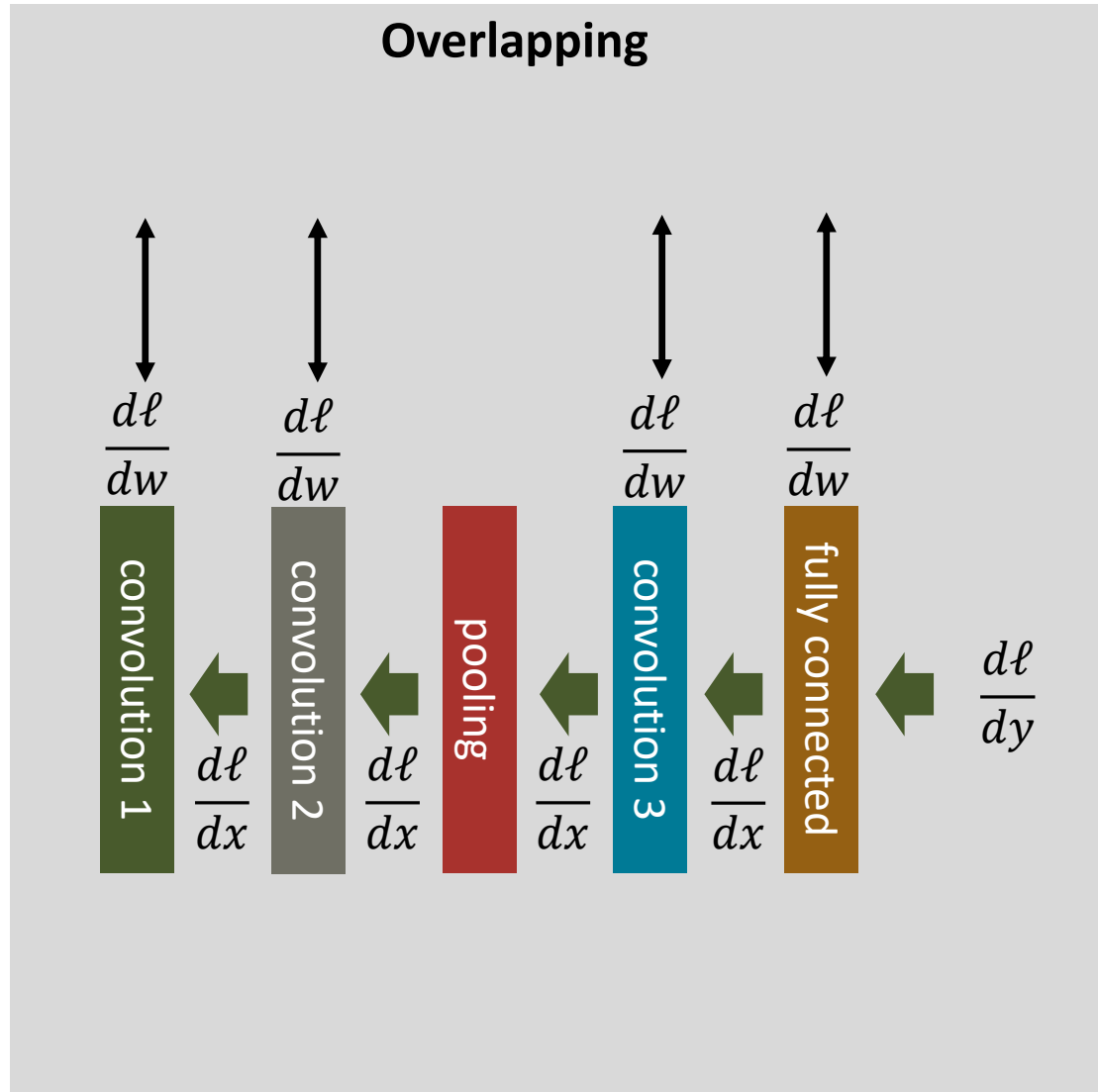
Distributed data-parallelism



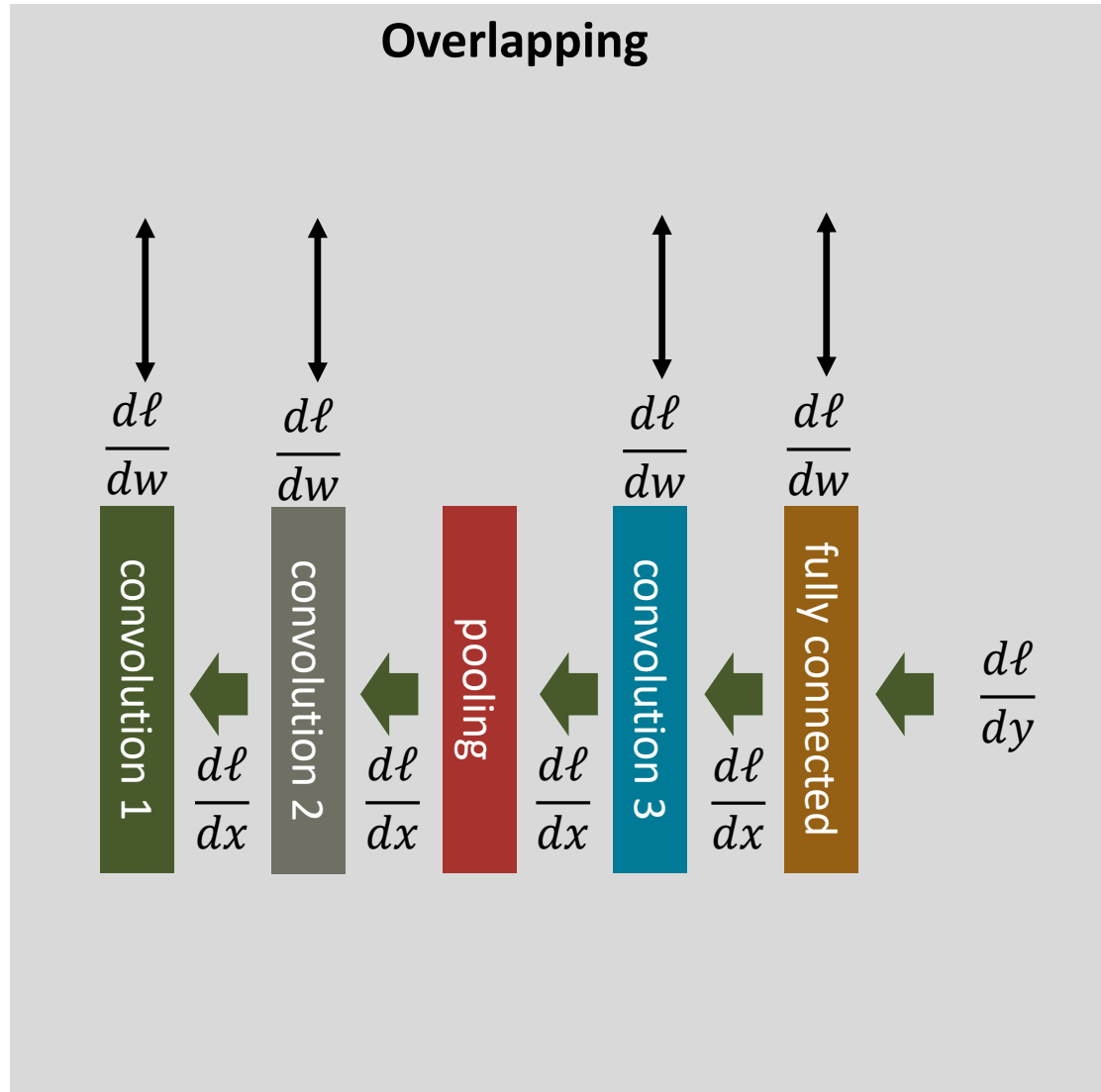
Distributed data-parallelism



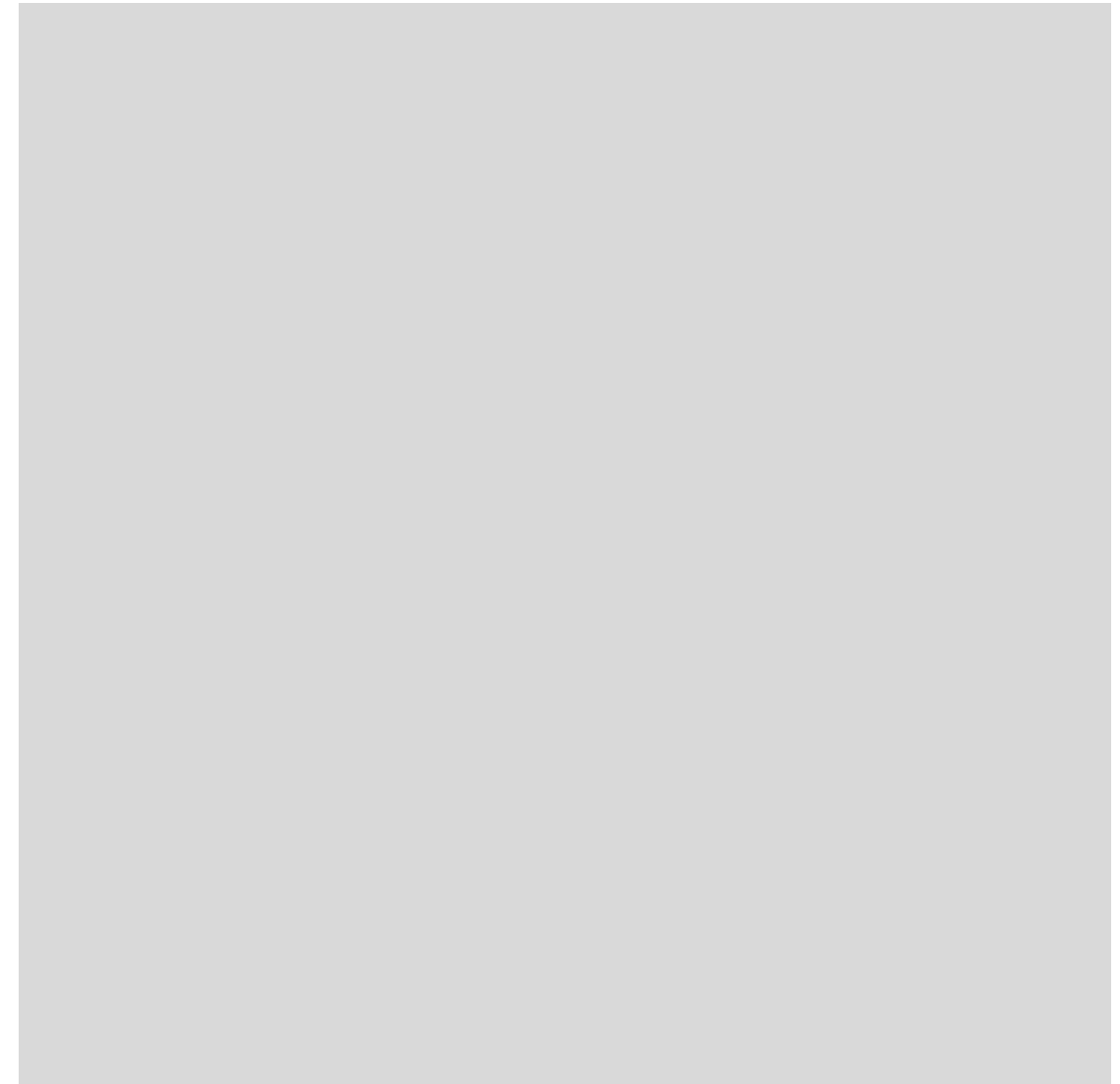
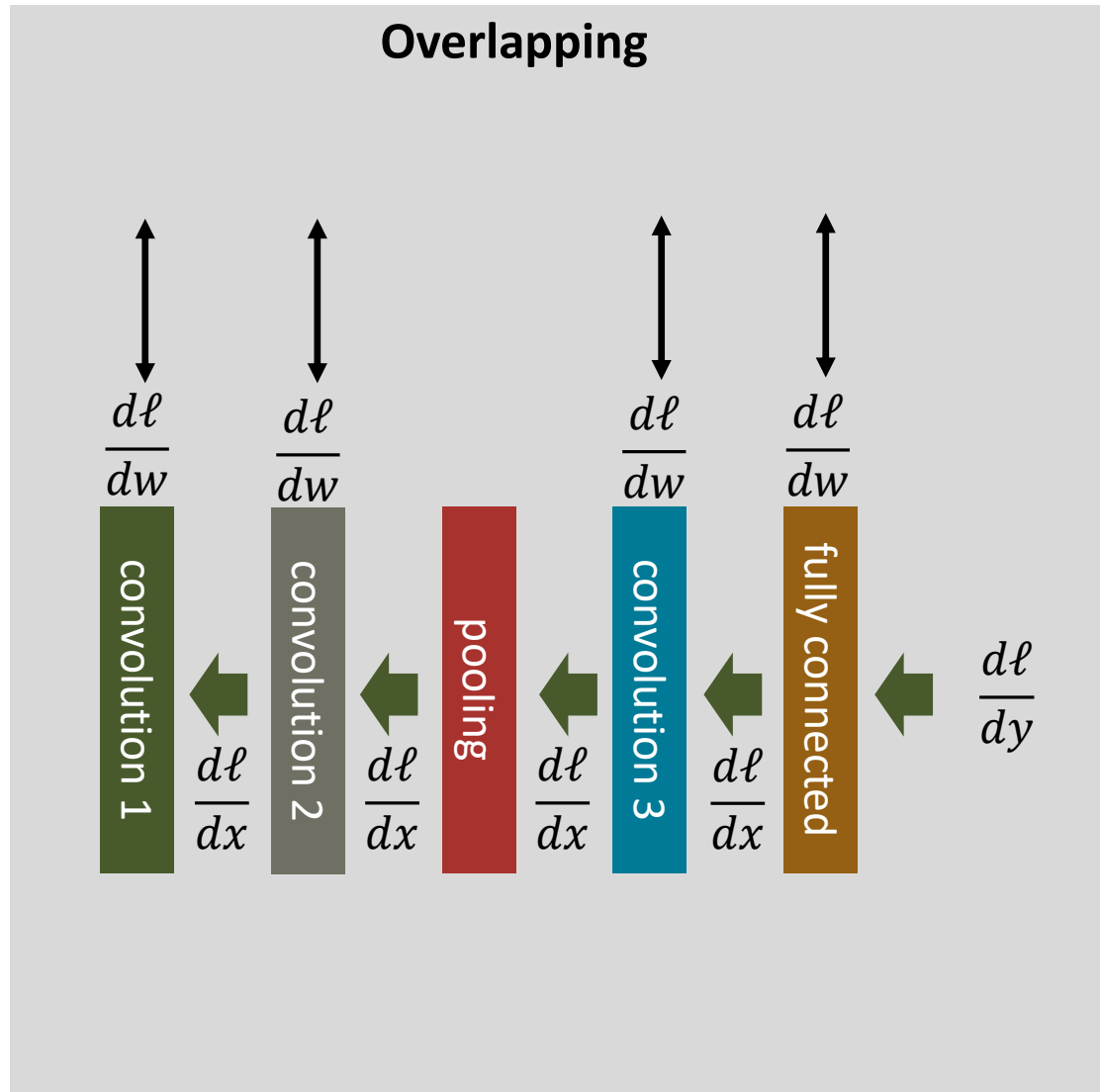
Distributed data-parallelism



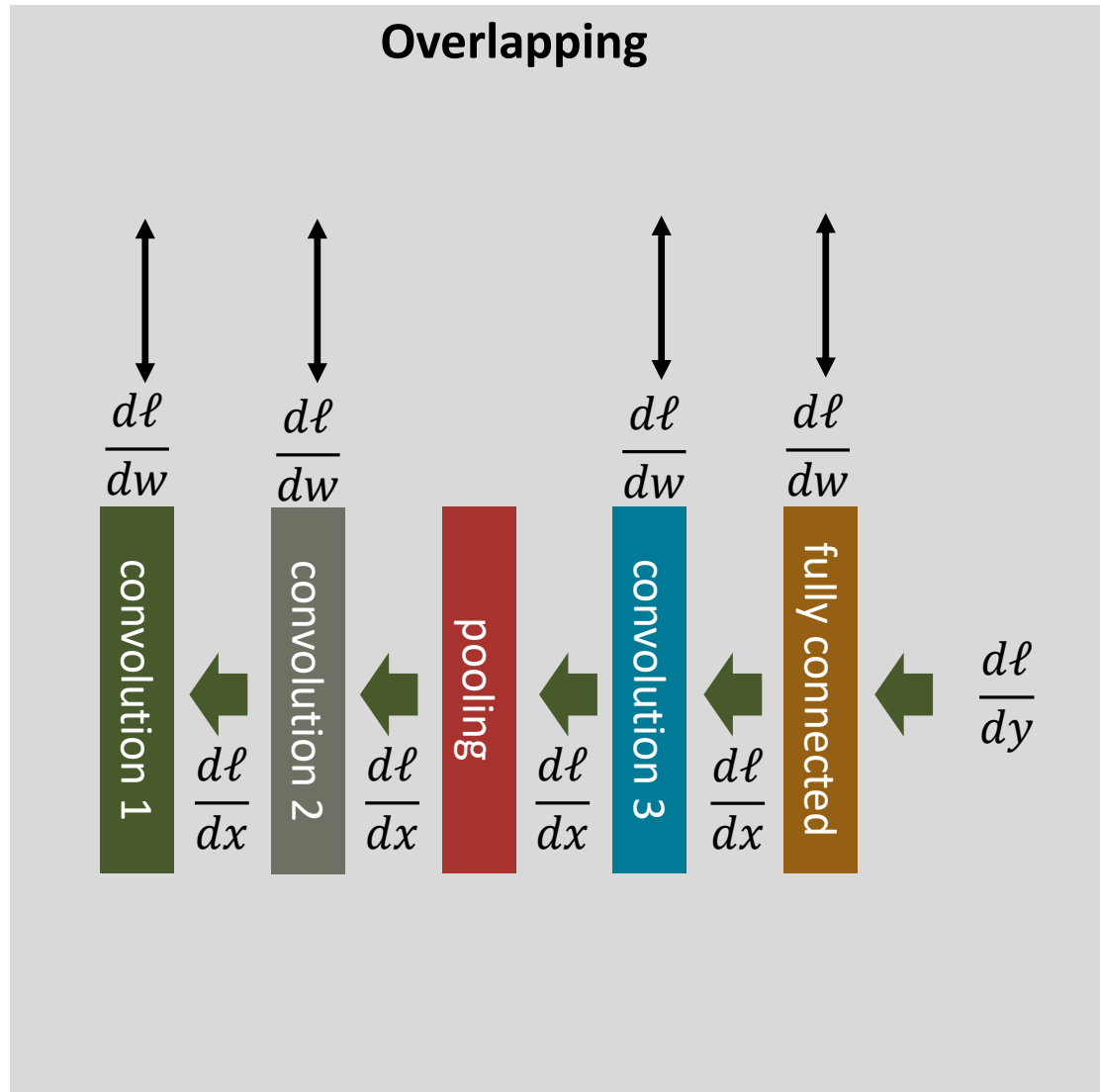
Distributed data-parallelism



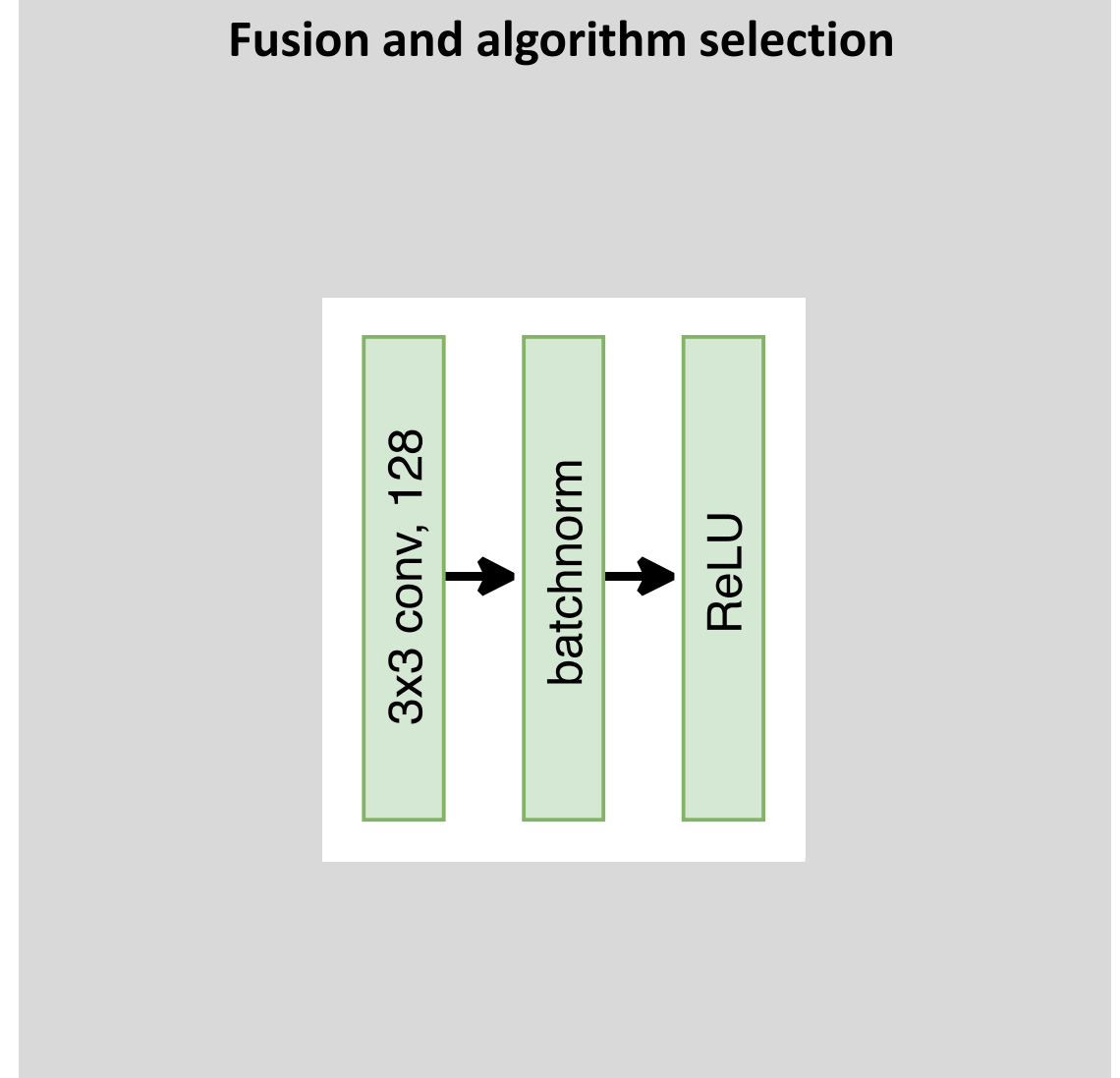
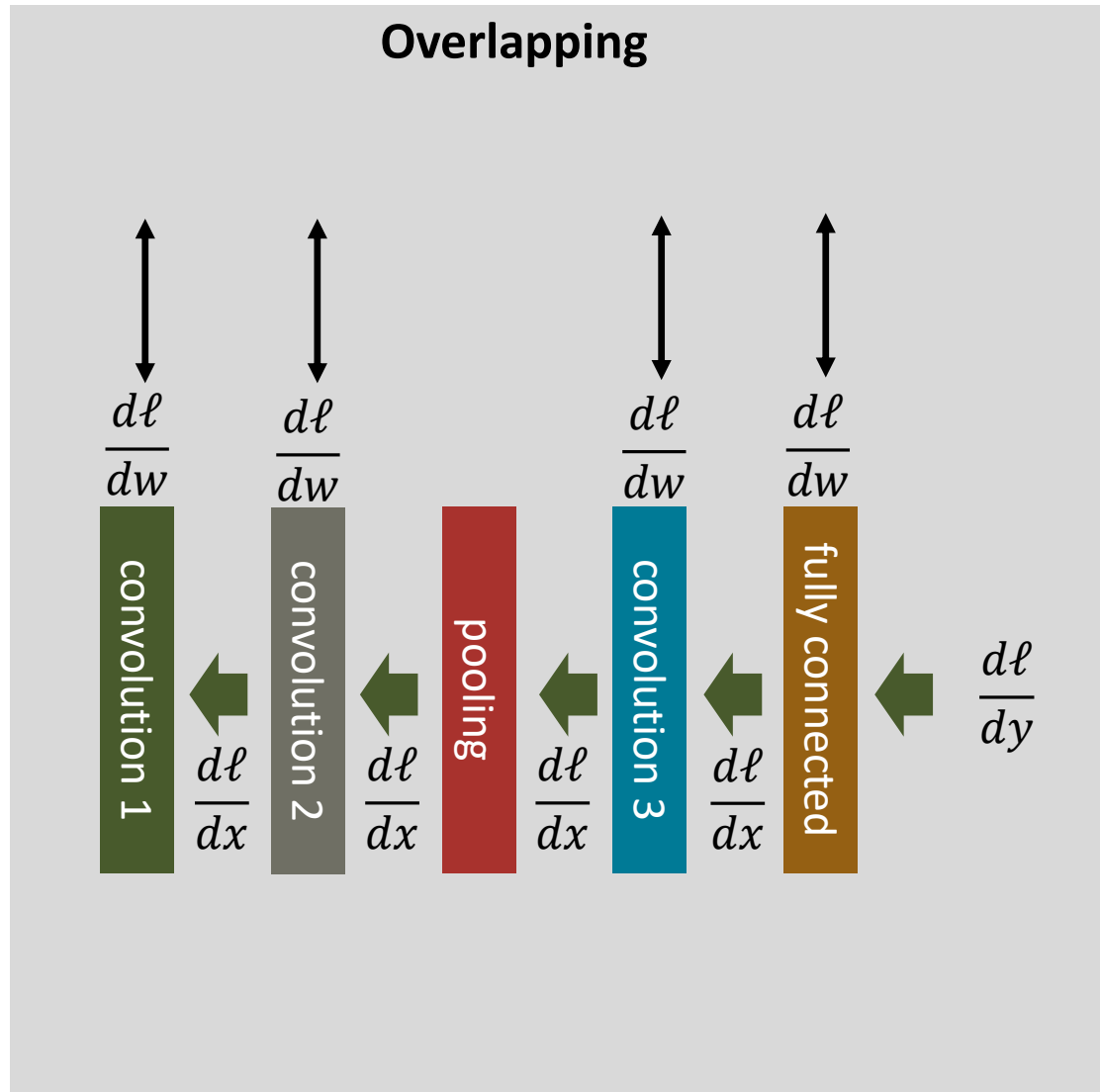
Distributed data-parallelism



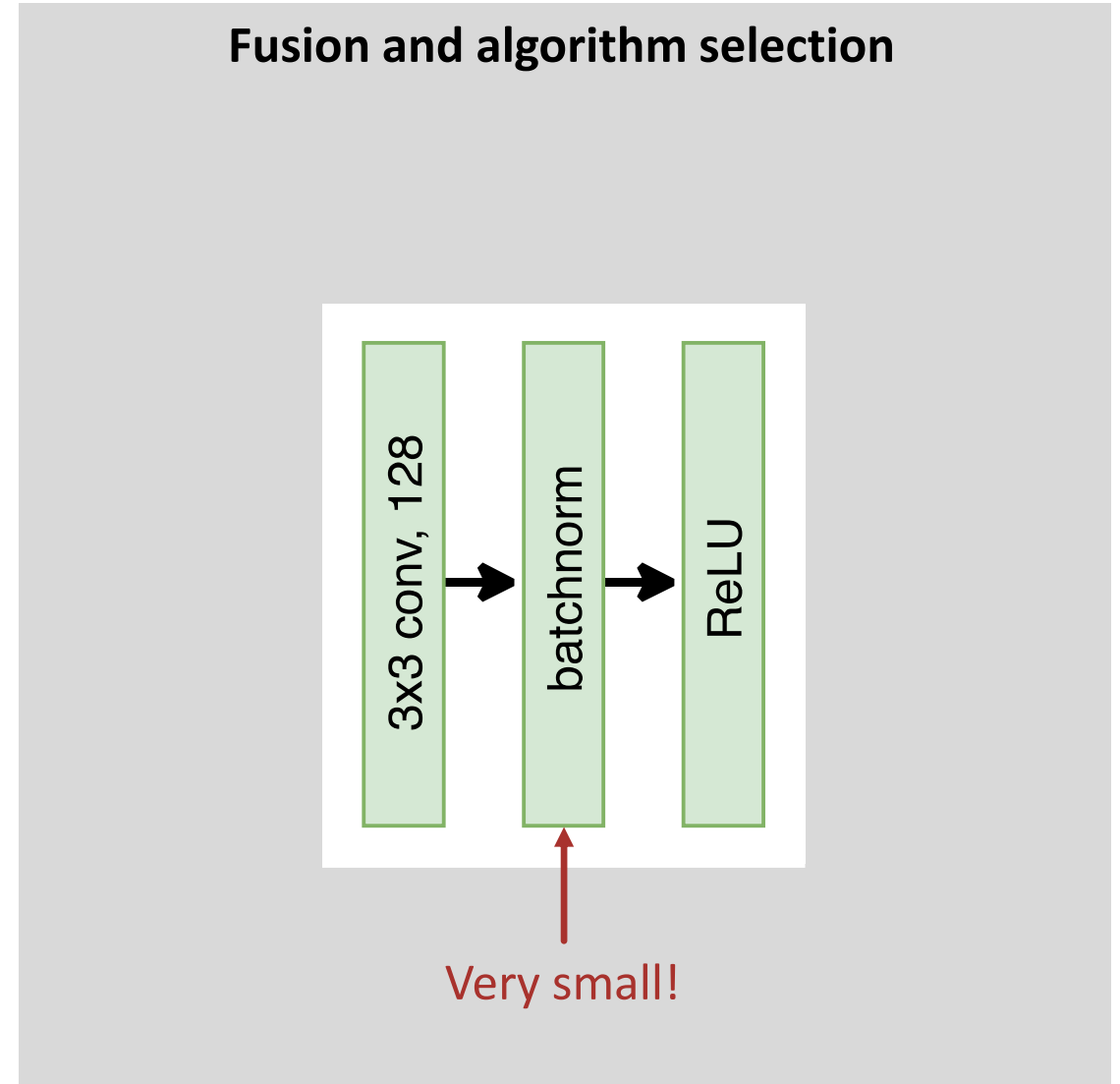
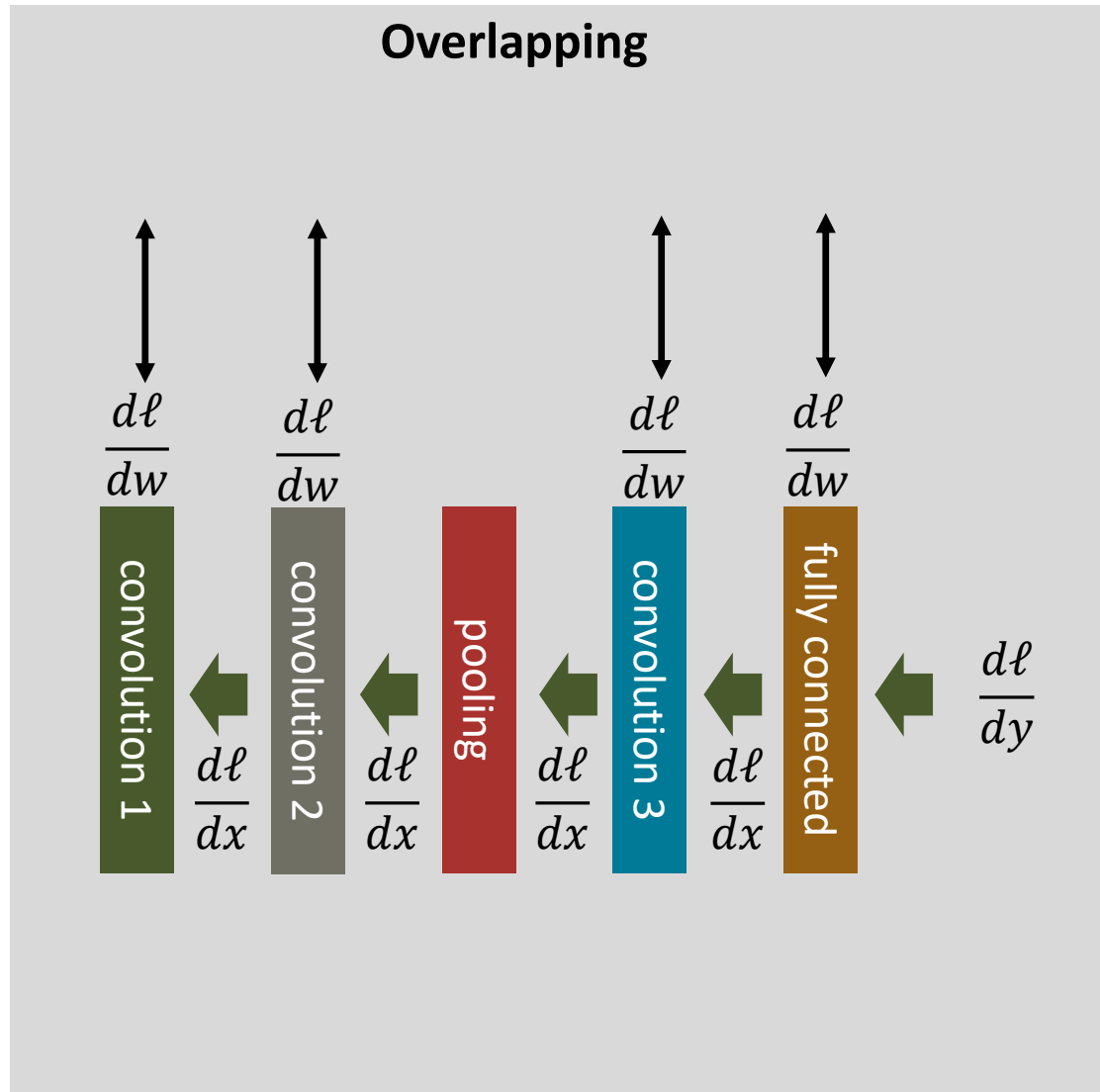
Distributed data-parallelism



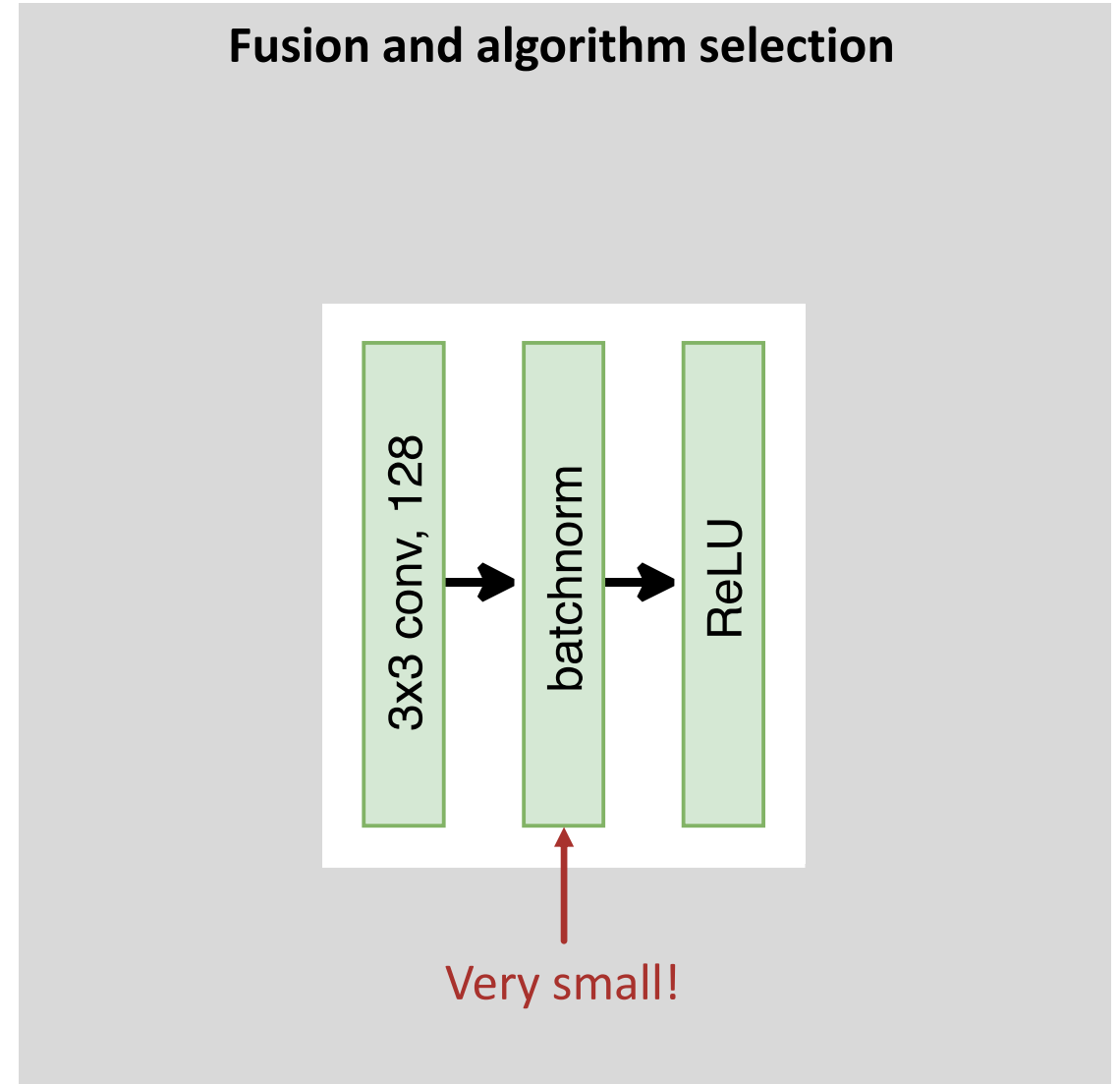
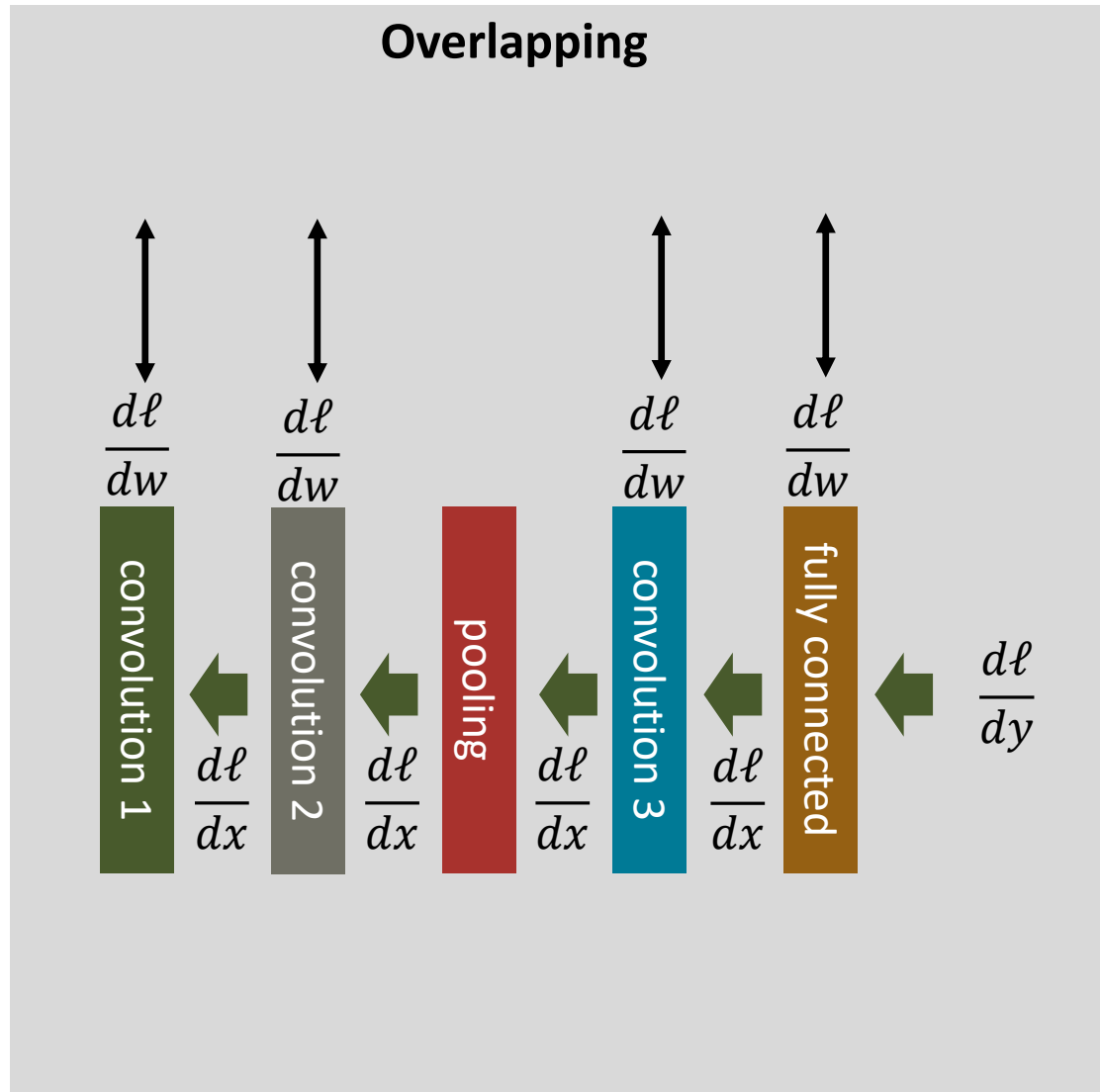
Distributed data-parallelism



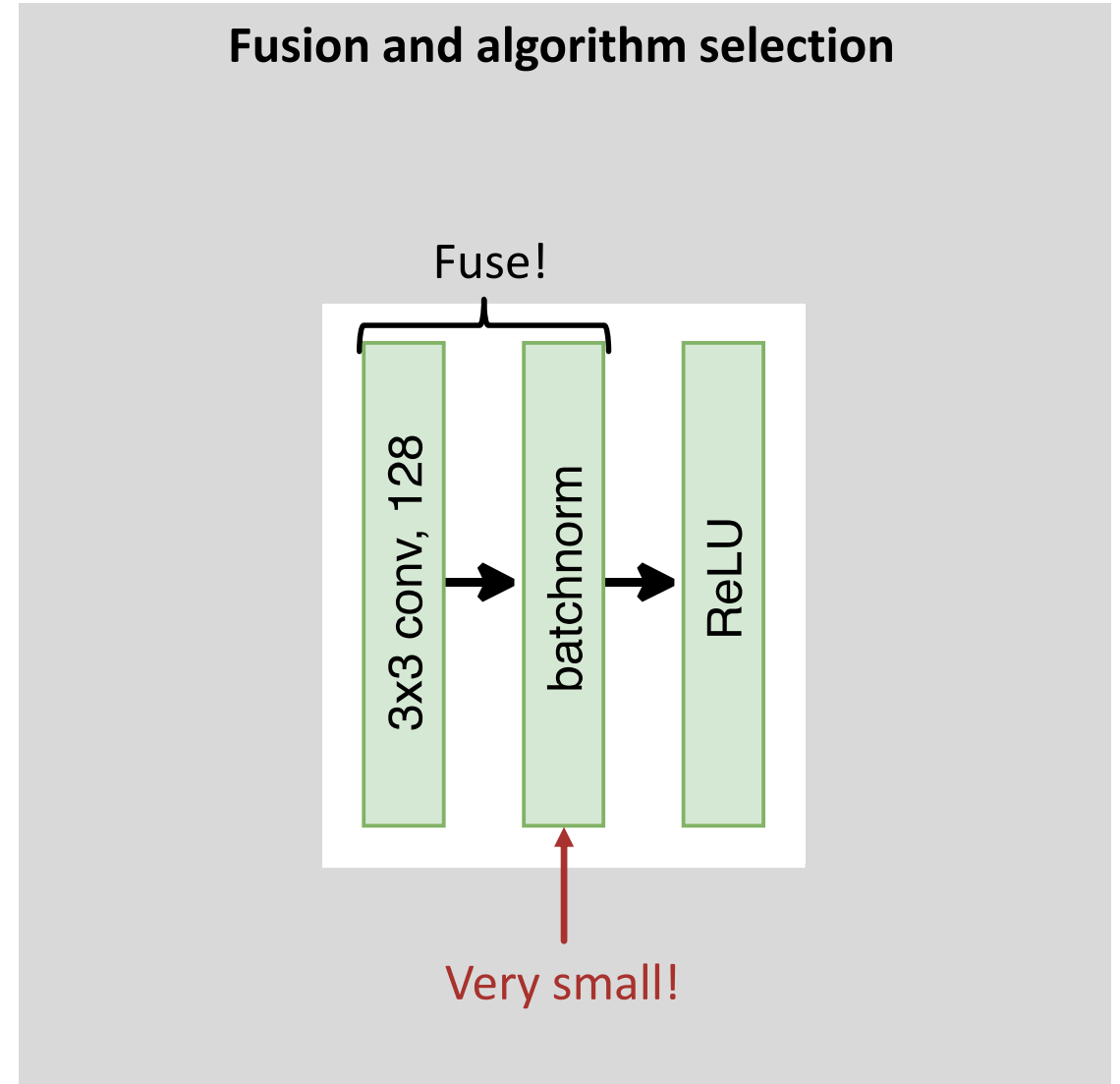
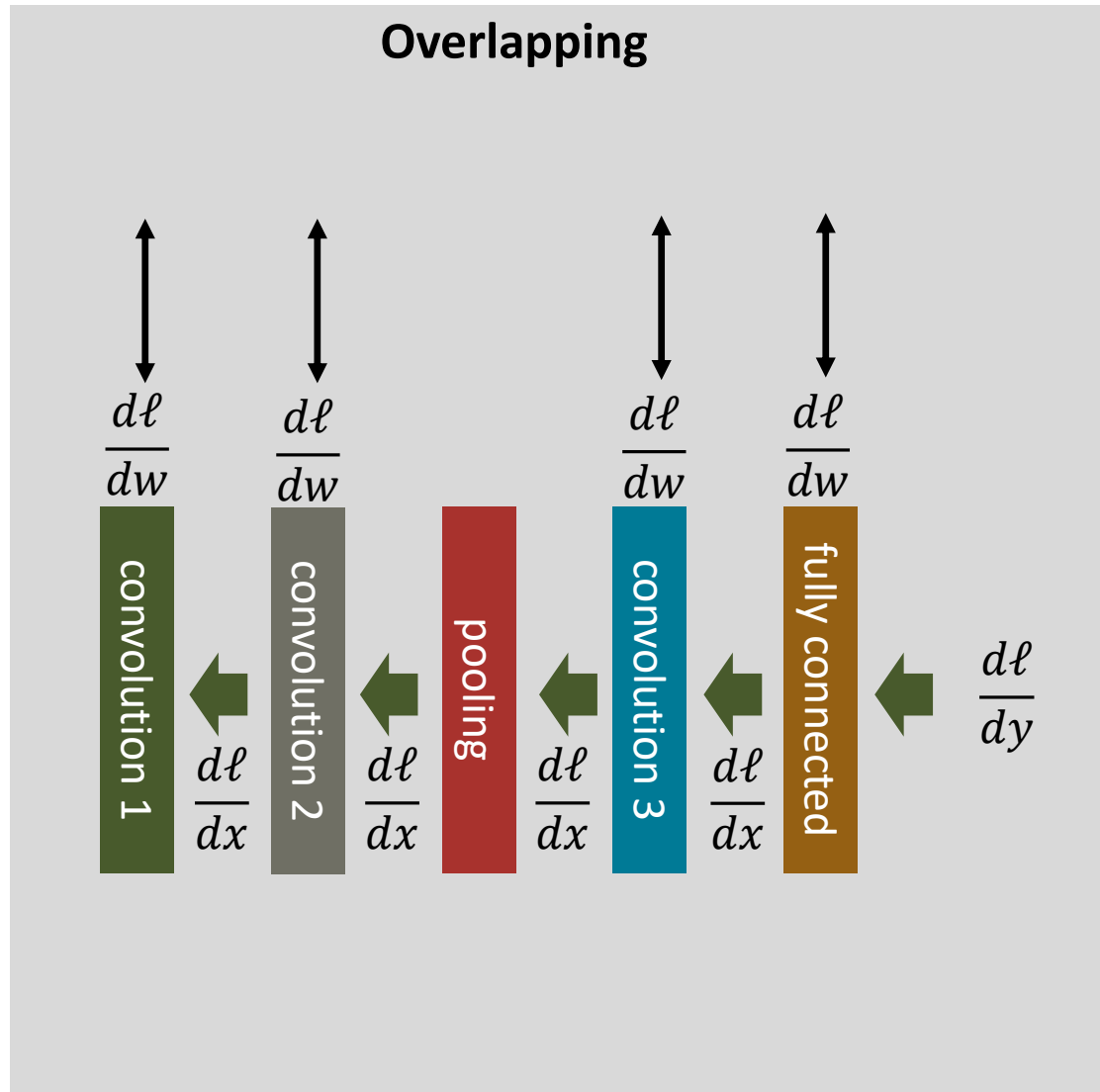
Distributed data-parallelism



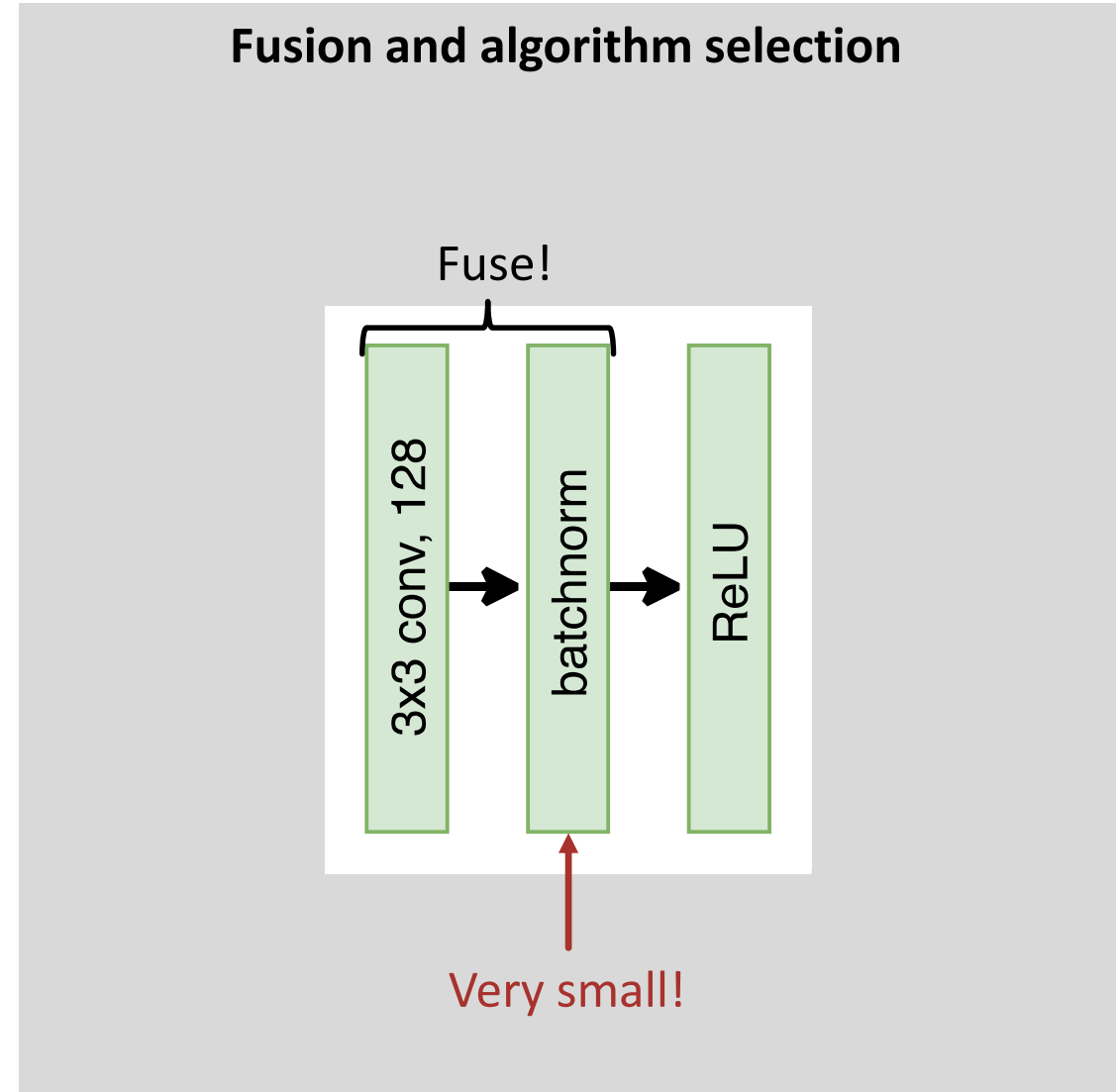
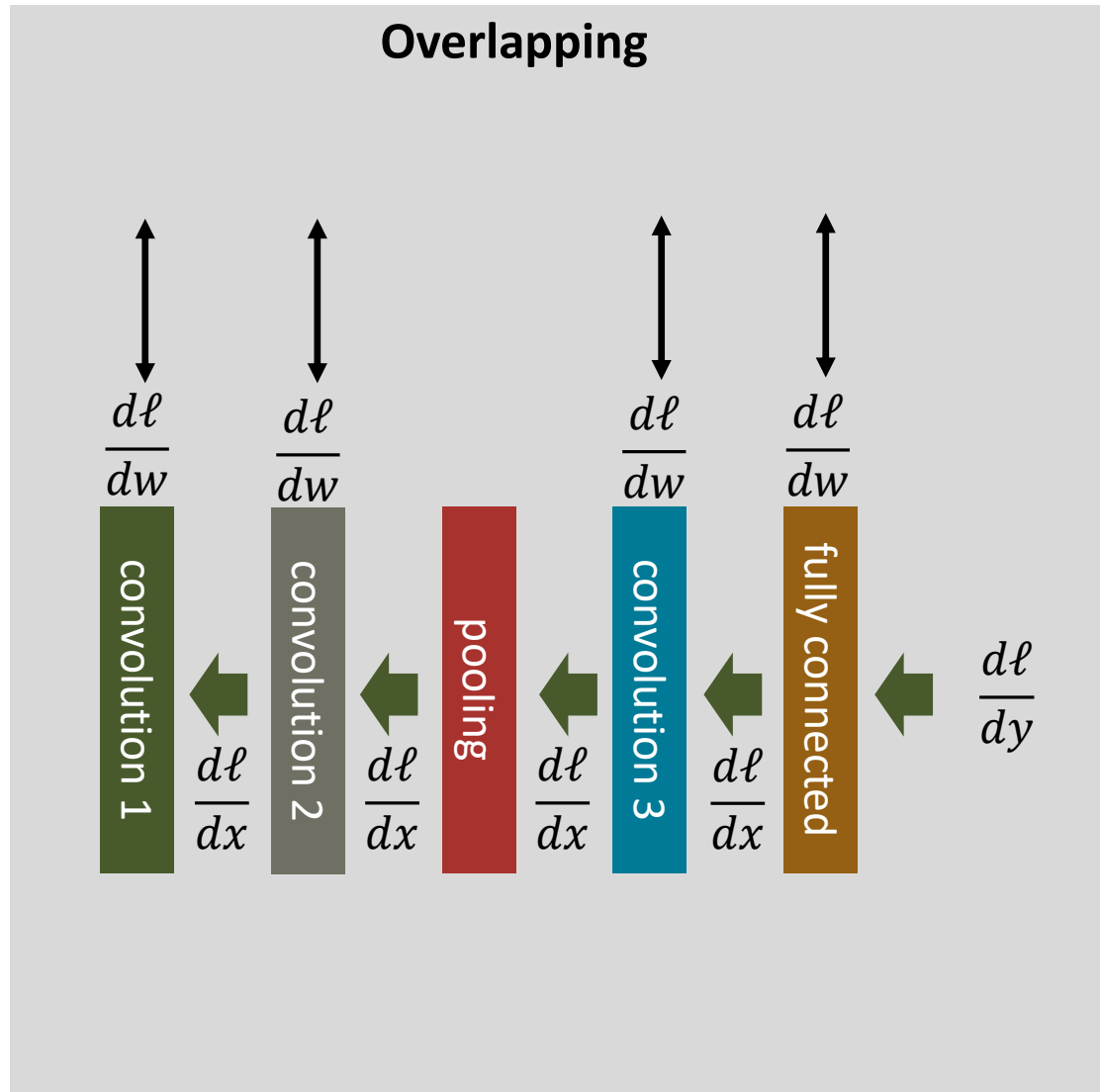
Distributed data-parallelism



Distributed data-parallelism

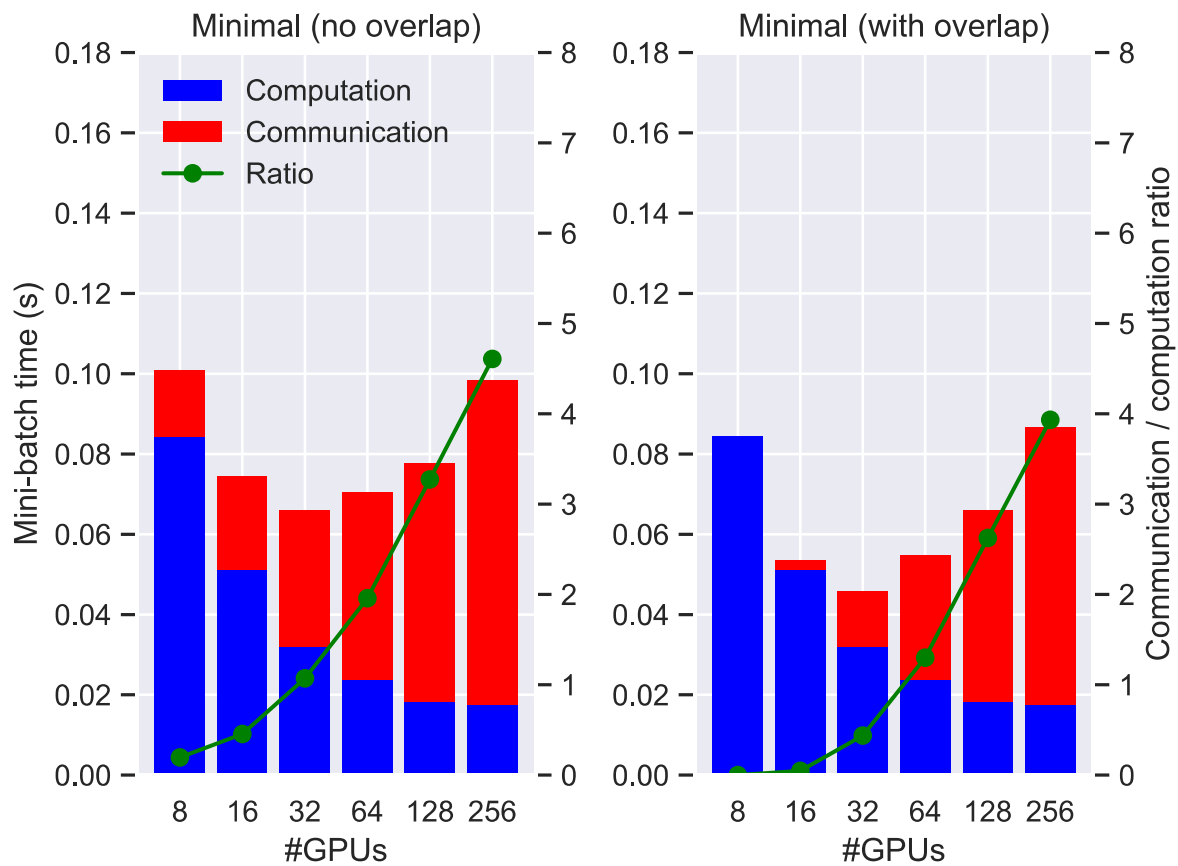


Distributed data-parallelism



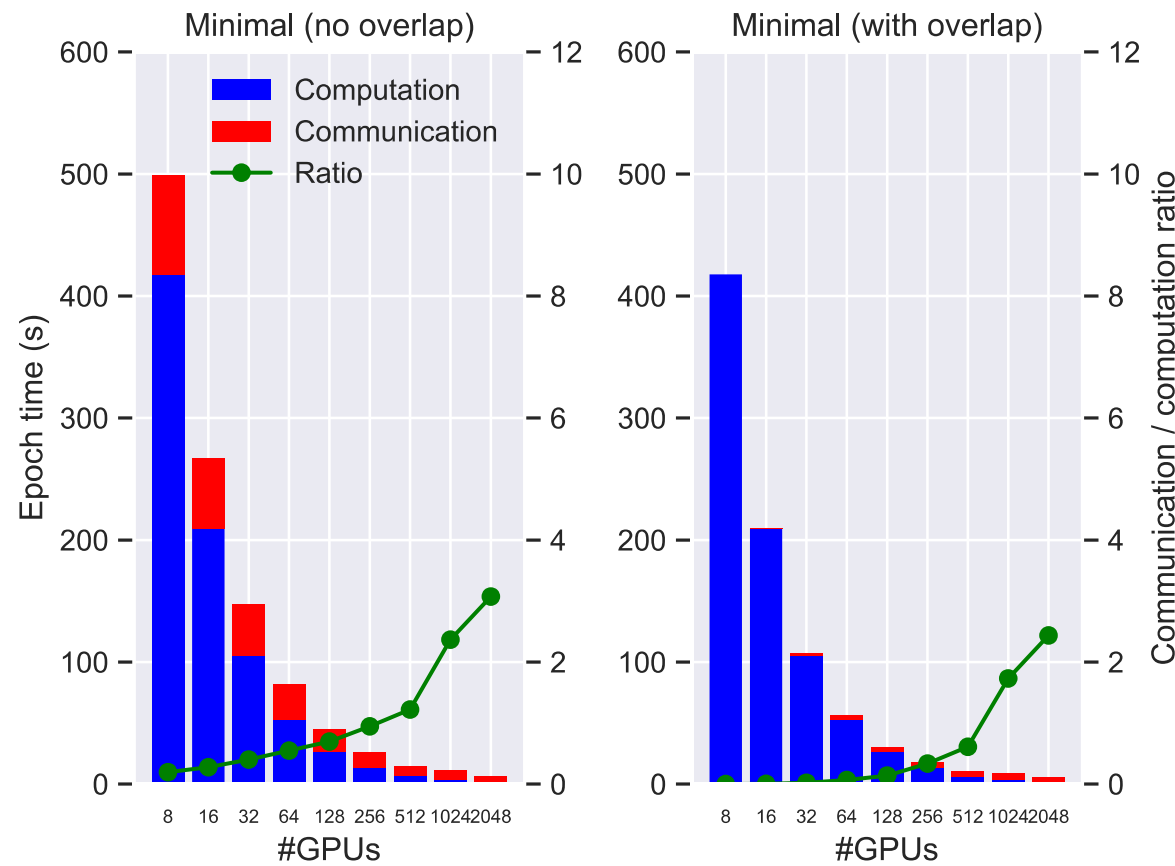
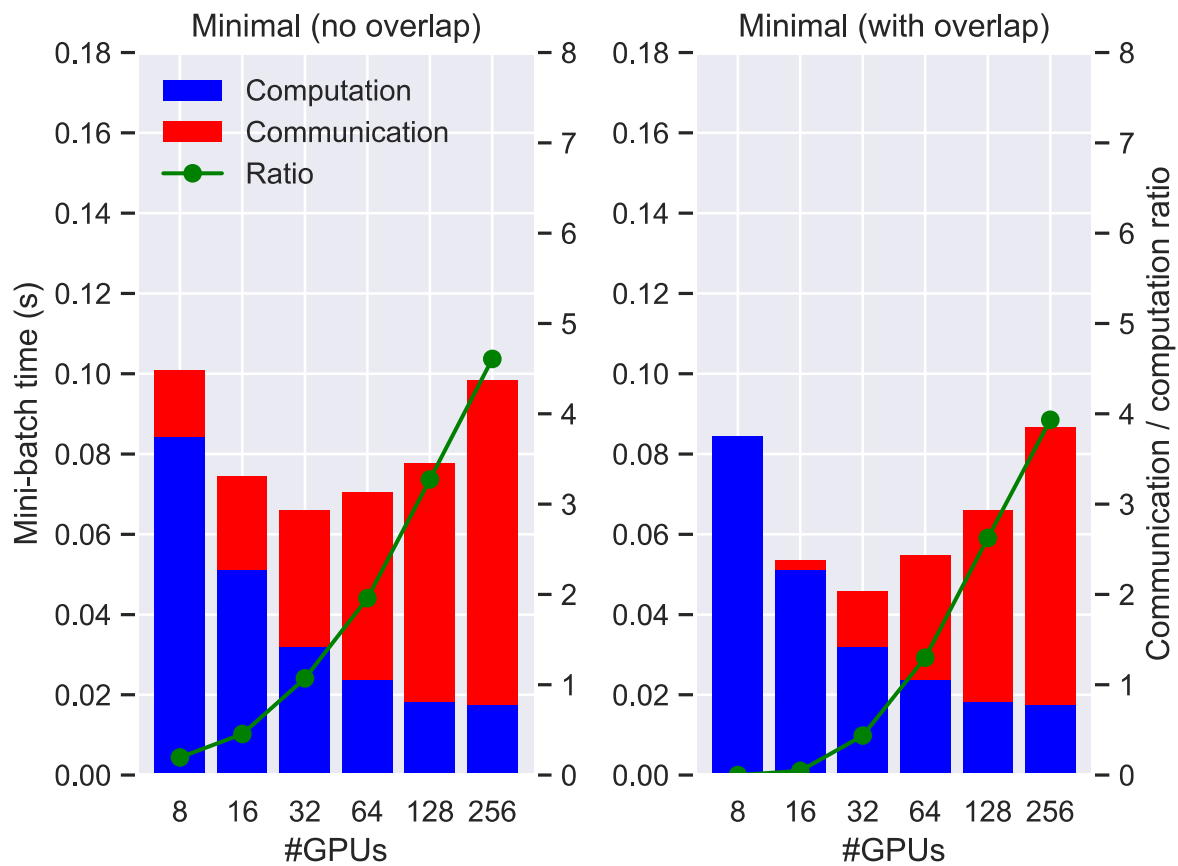
Distributed data-parallelism

Strong scaling



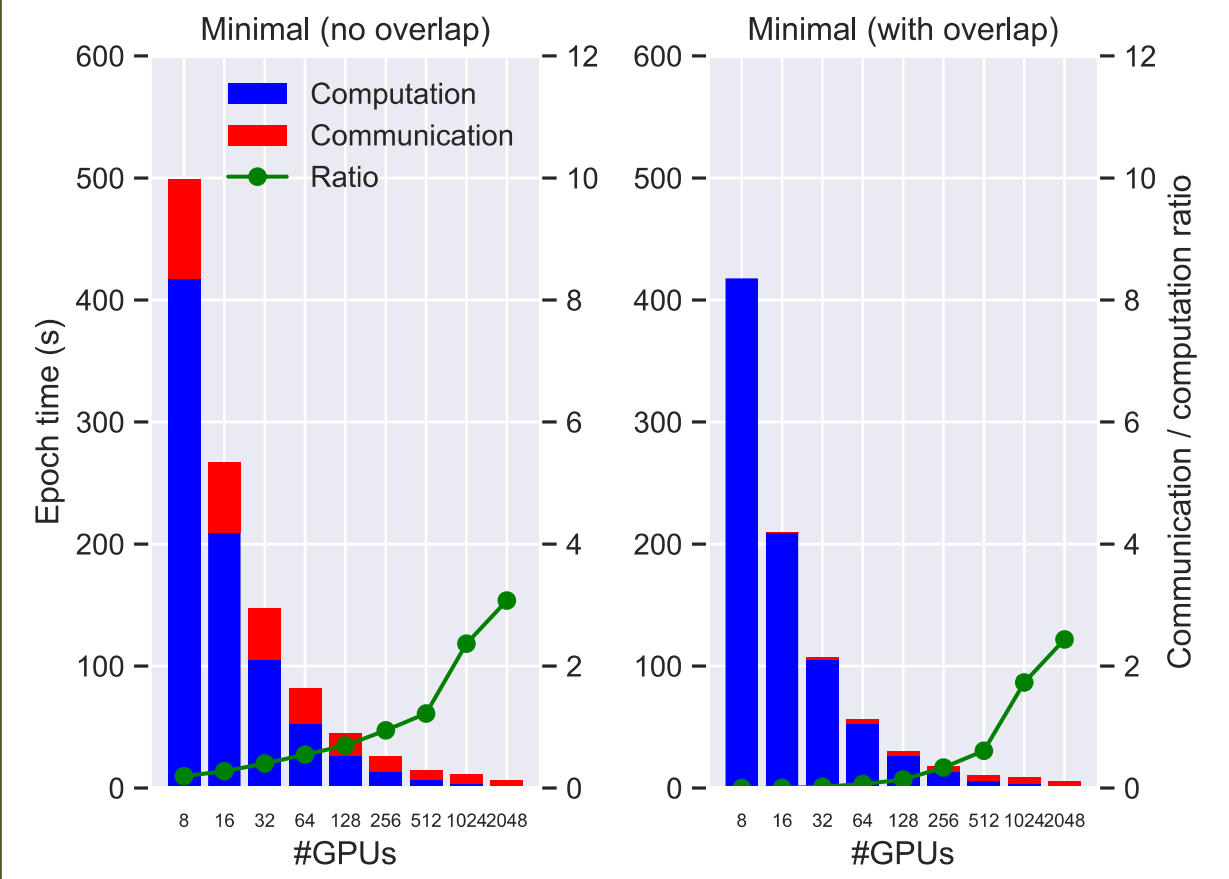
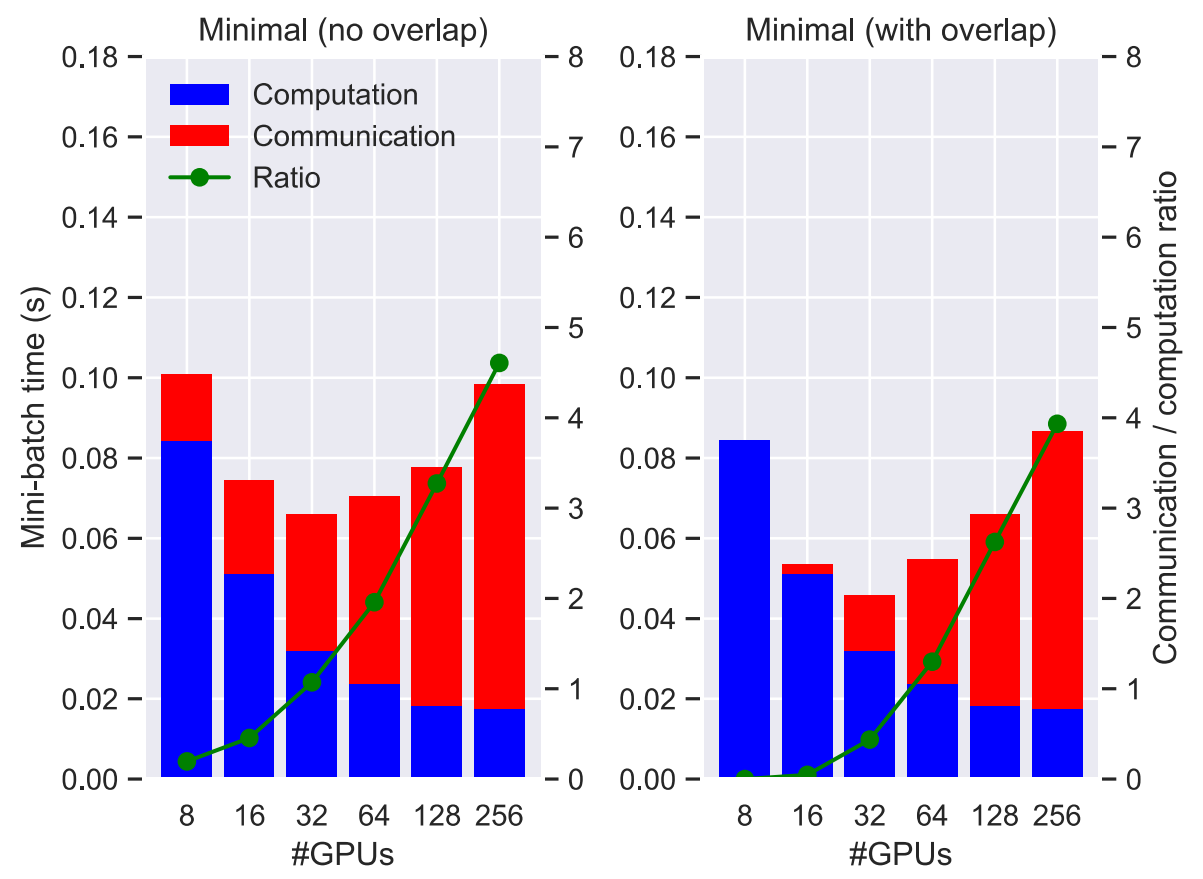
Distributed data-parallelism

Strong scaling



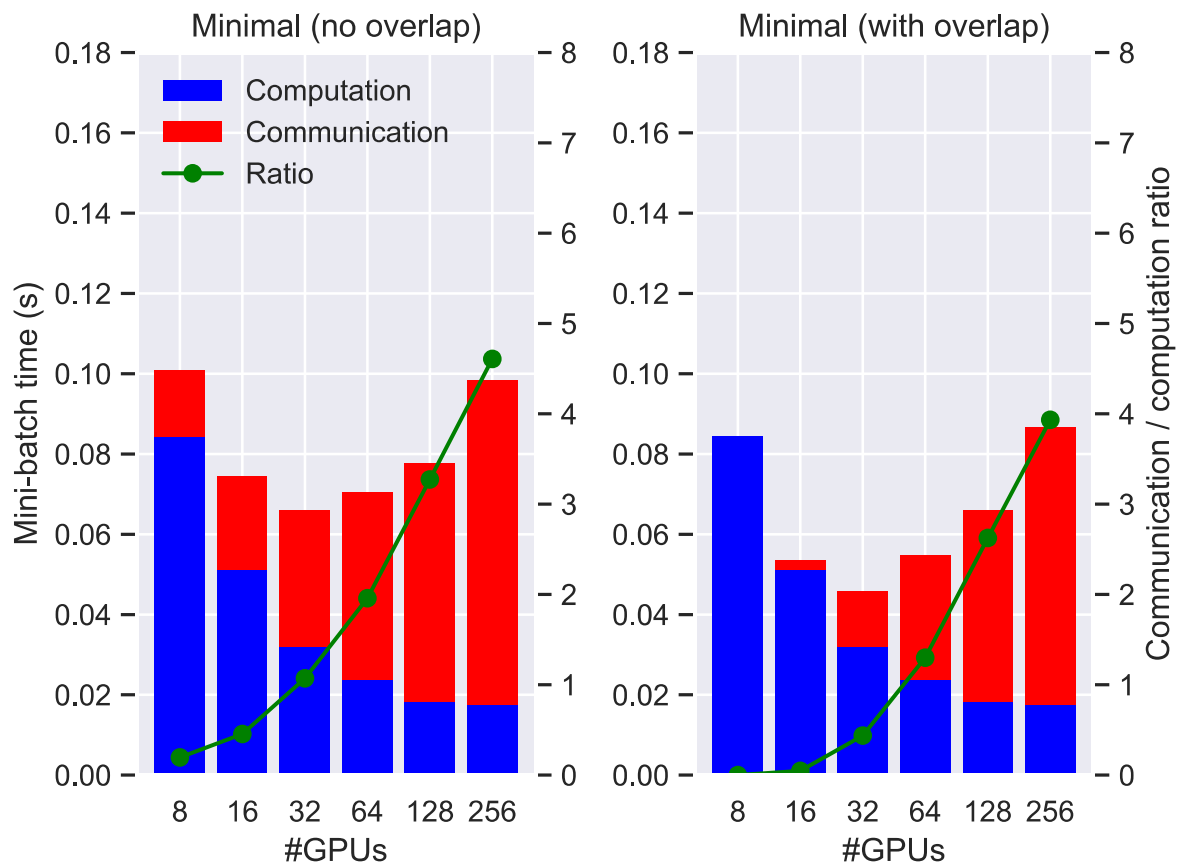
Distributed data-parallelism

Strong scaling

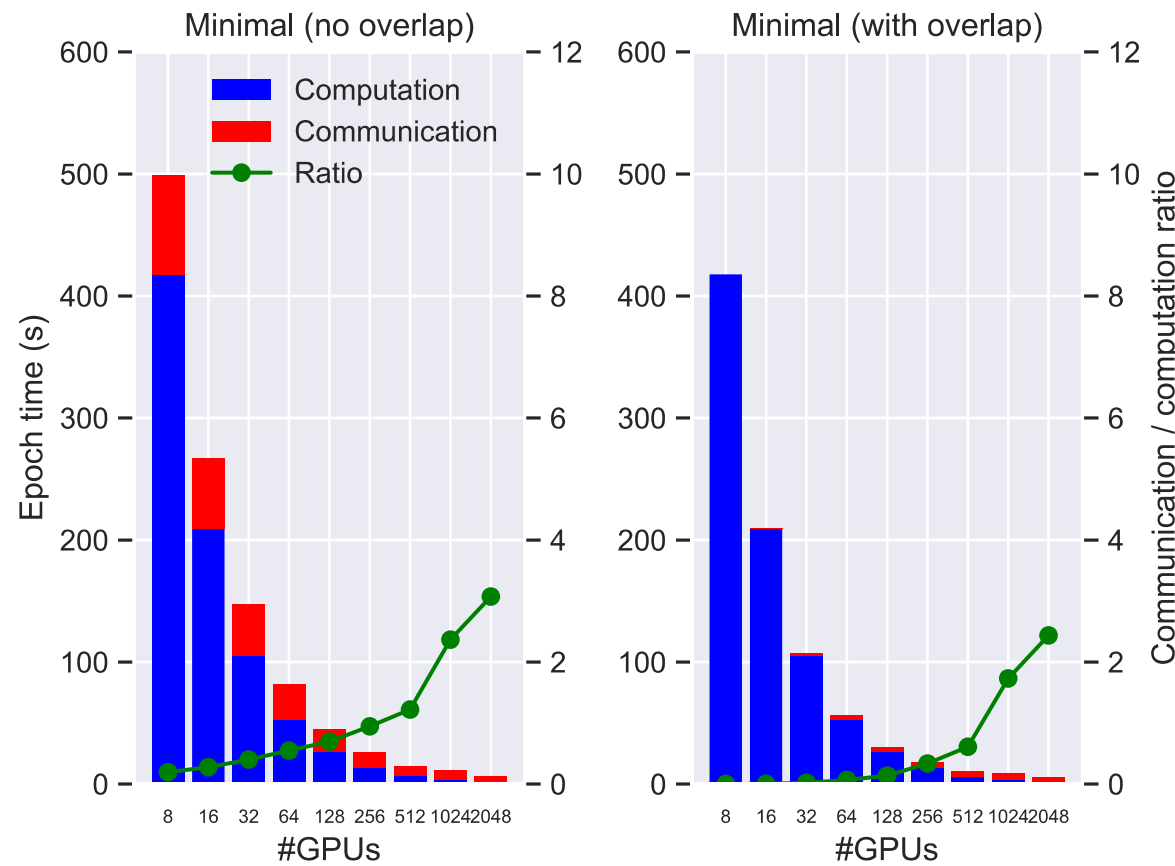


Distributed data-parallelism

Strong scaling

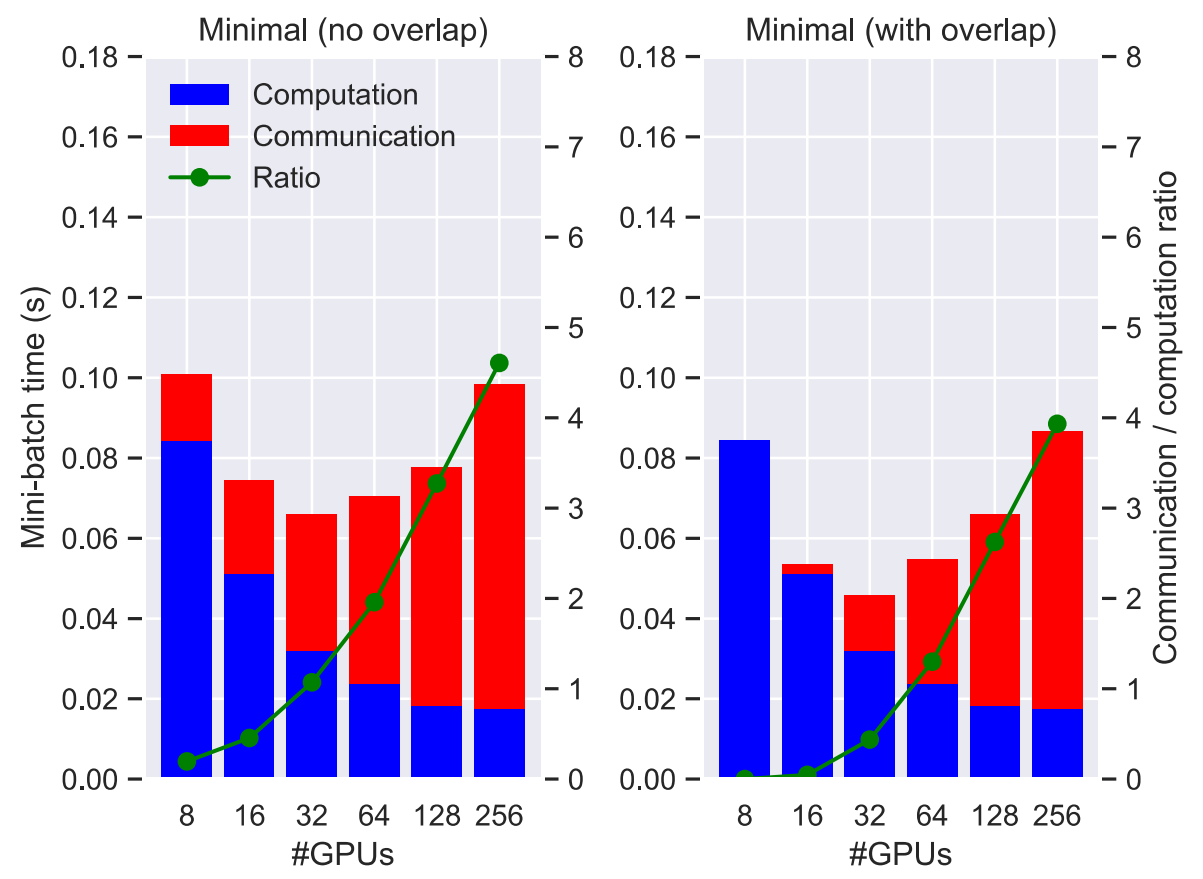


Weak scaling

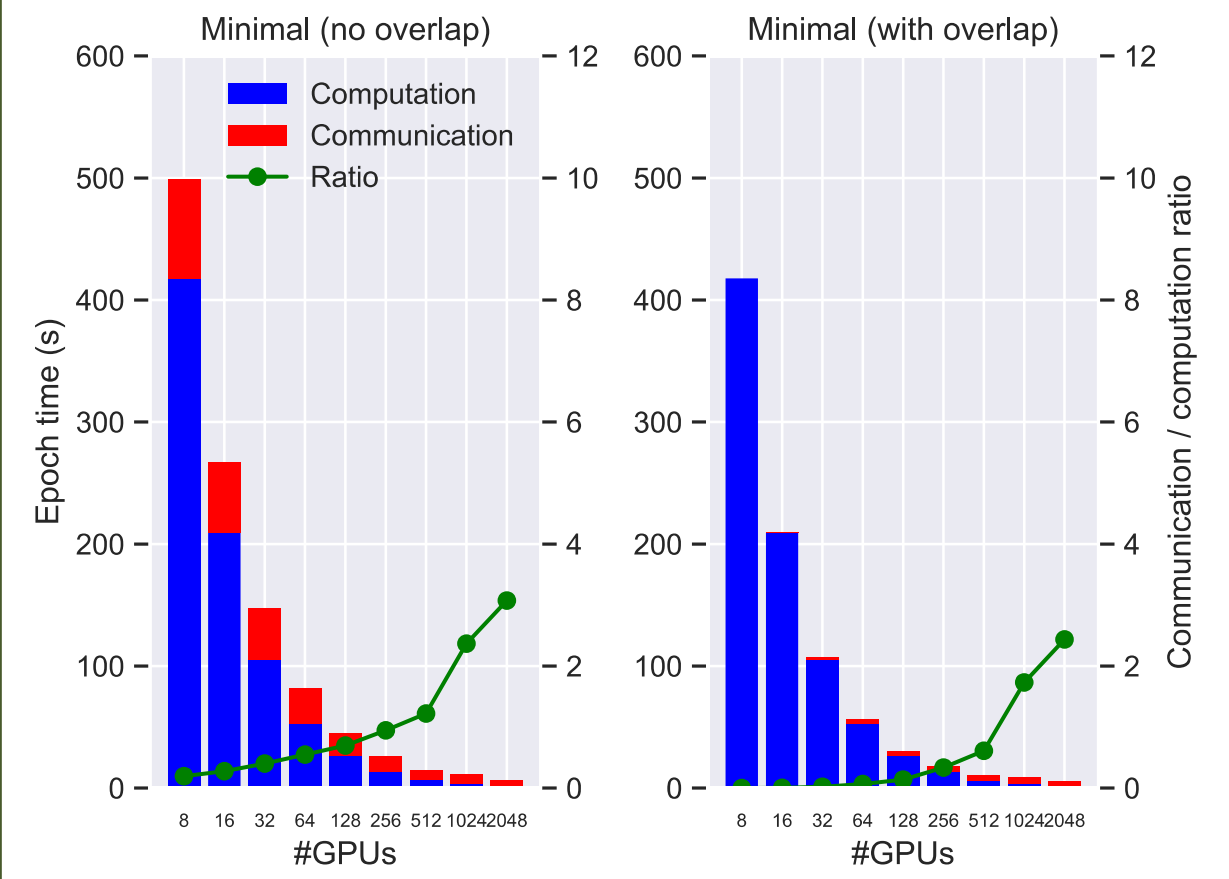


Distributed data-parallelism

Strong scaling



Weak scaling



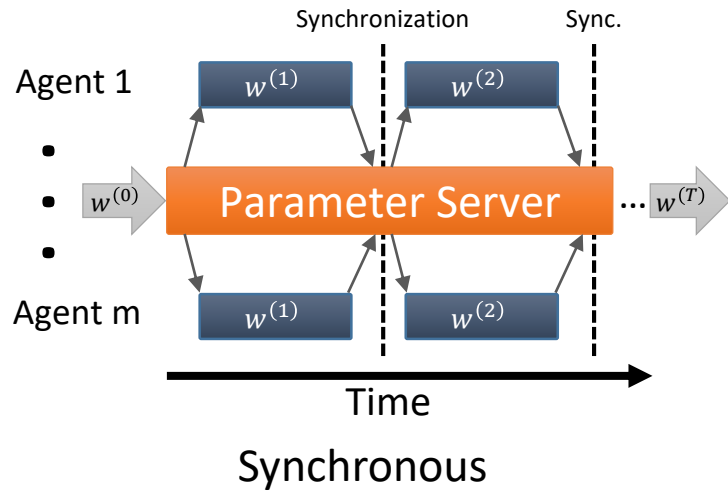
Parameter (and model) consistency - centralized

- Parameter exchange frequency can be controlled, while still attaining convergence:

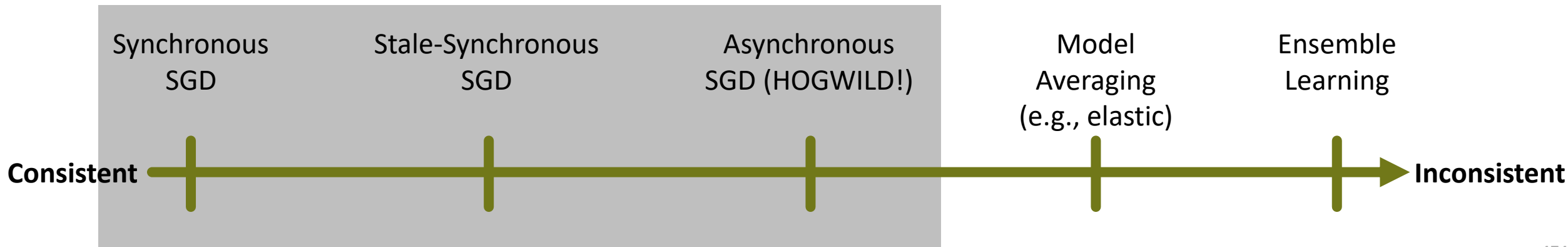


Parameter (and model) consistency - centralized

- Parameter exchange frequency can be controlled, while still attaining convergence:

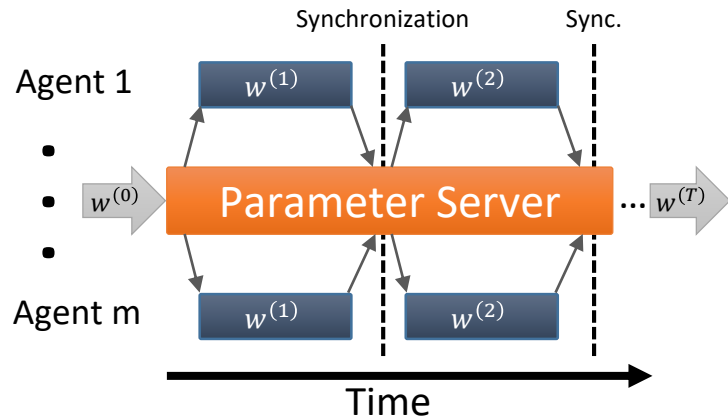


- Trades off “statistical performance” for “hardware performance”



Parameter (and model) consistency - centralized

- Parameter exchange frequency can be controlled, while still attaining convergence:

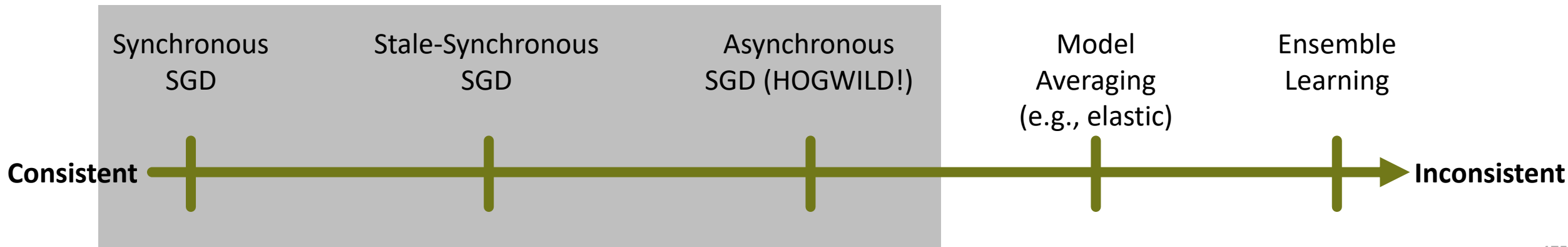


Synchronous

Stale Synchronous / Bounded Asynchronous

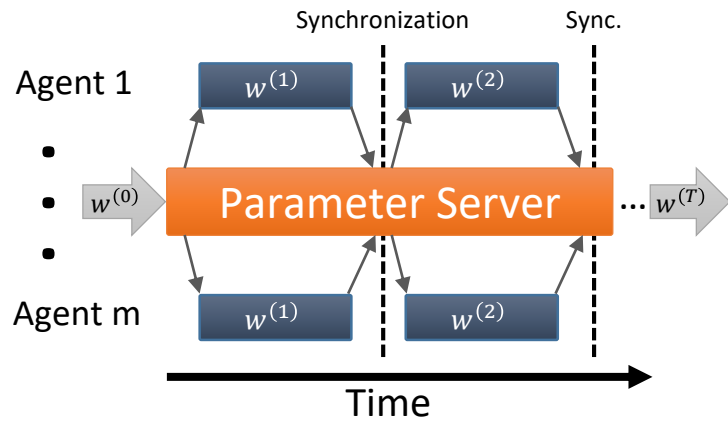
Asynchronous

- Trades off “statistical performance” for “hardware performance”



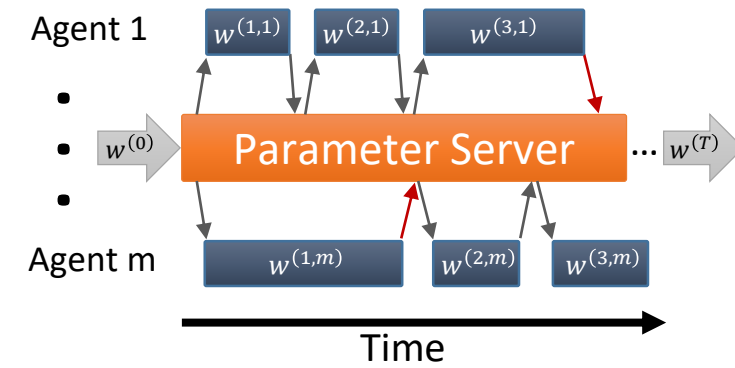
Parameter (and model) consistency - centralized

- Parameter exchange frequency can be controlled, while still attaining convergence:



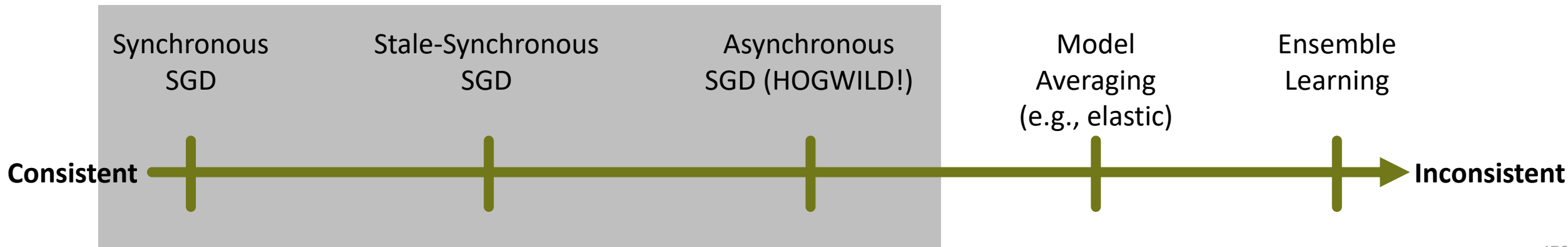
Synchronous

Stale Synchronous / Bounded Asynchronous



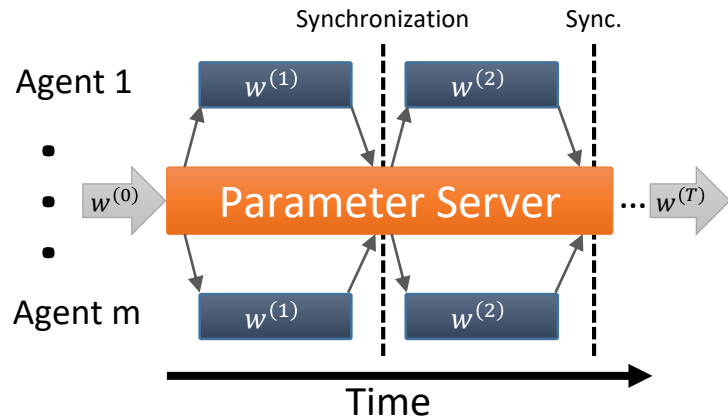
Asynchronous

- Trades off “statistical performance” for “hardware performance”



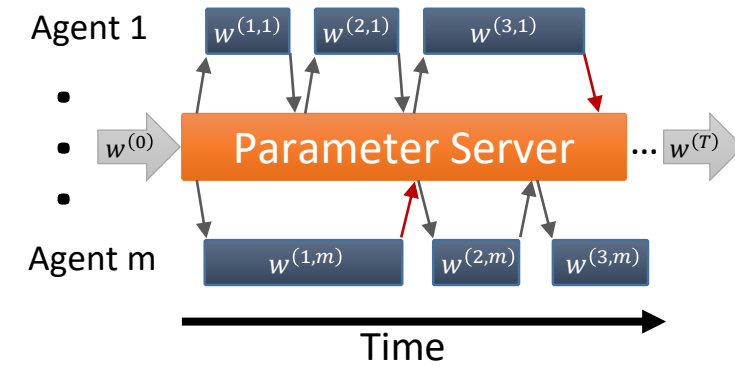
Parameter (and model) consistency - centralized

- Parameter exchange frequency can be controlled, while still attaining convergence:



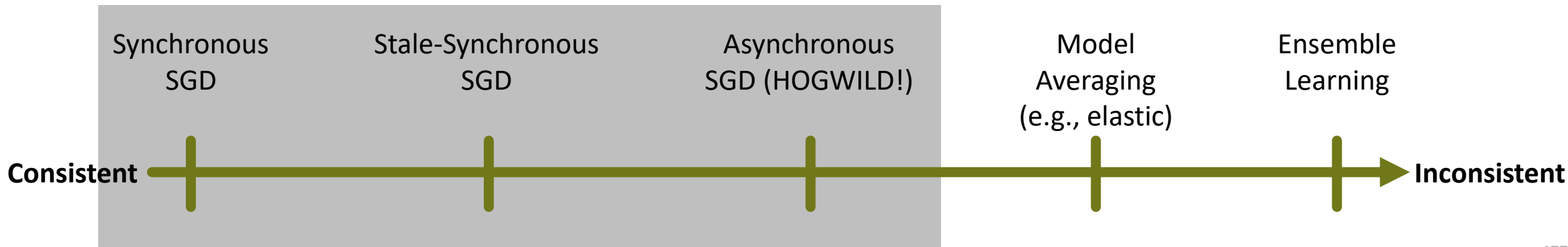
Synchronous

Stale Synchronous / Bounded Asynchronous



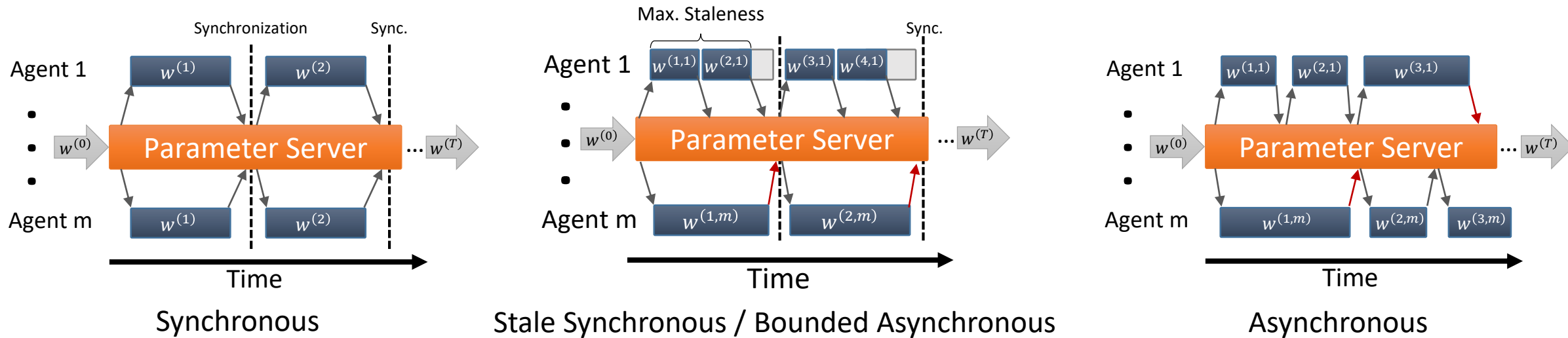
Asynchronous

- Trades off “statistical performance” for “hardware performance”



Parameter (and model) consistency - centralized

- Parameter exchange frequency can be controlled, while still attaining convergence:

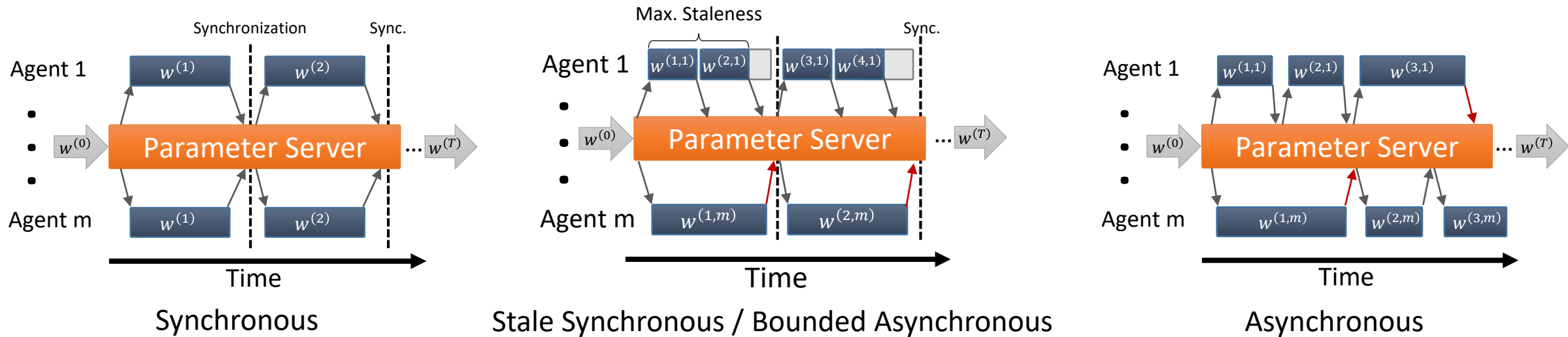


- Trades off “statistical performance” for “hardware performance”

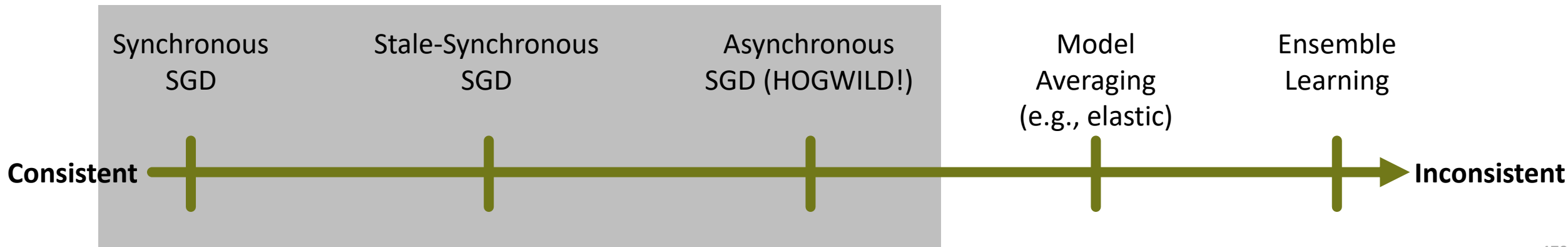


Parameter (and model) consistency - centralized

- Parameter exchange frequency can be controlled, while still attaining convergence:

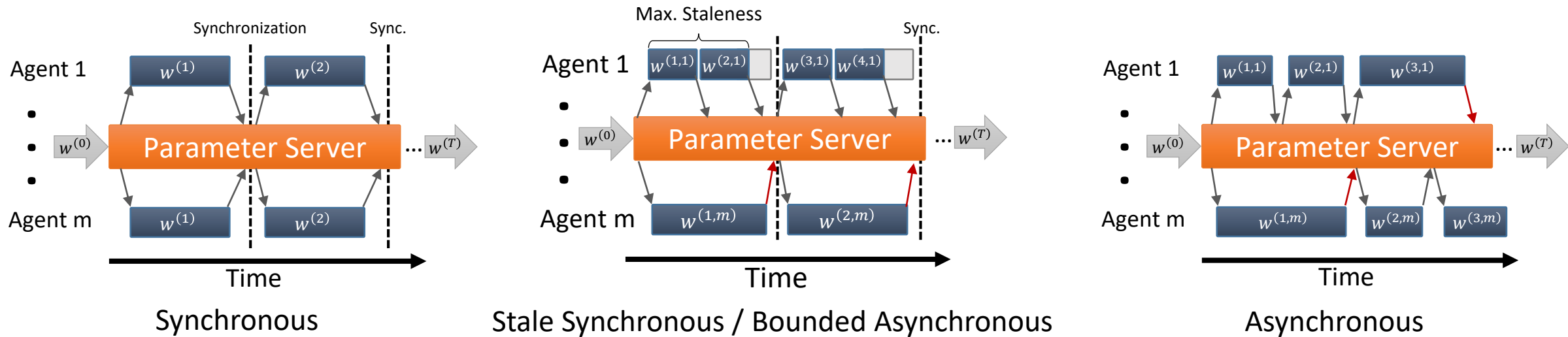


- Trades off “statistical performance” for “hardware performance”

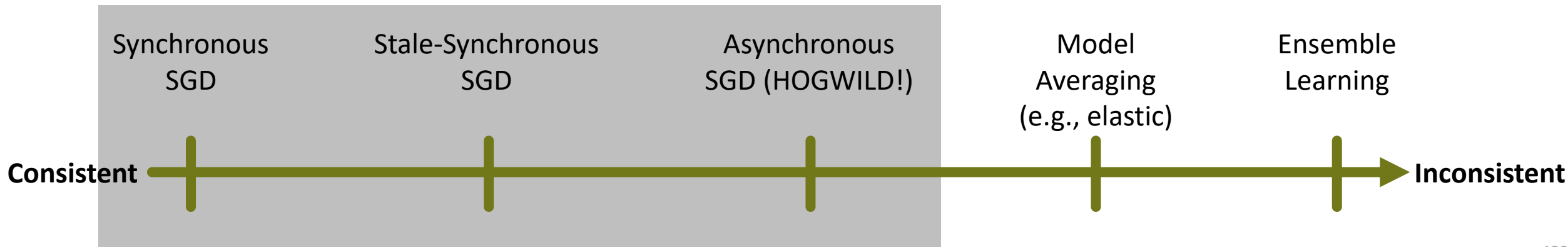


Parameter (and model) consistency - centralized

- Parameter exchange frequency can be controlled, while still attaining convergence:



- Trades off “statistical performance” for “hardware performance”



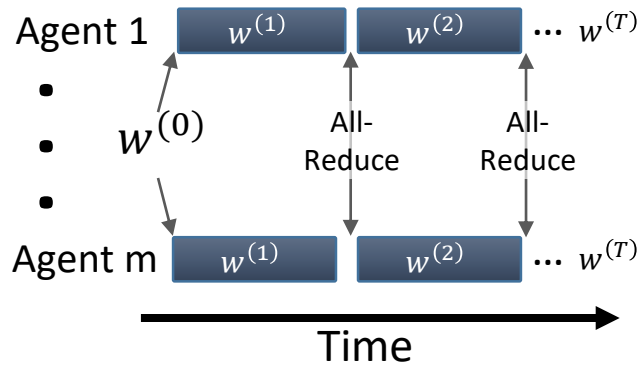
Parameter (and model) consistency - decentralized

- Parameter exchange frequency can be controlled, while still attaining convergence:



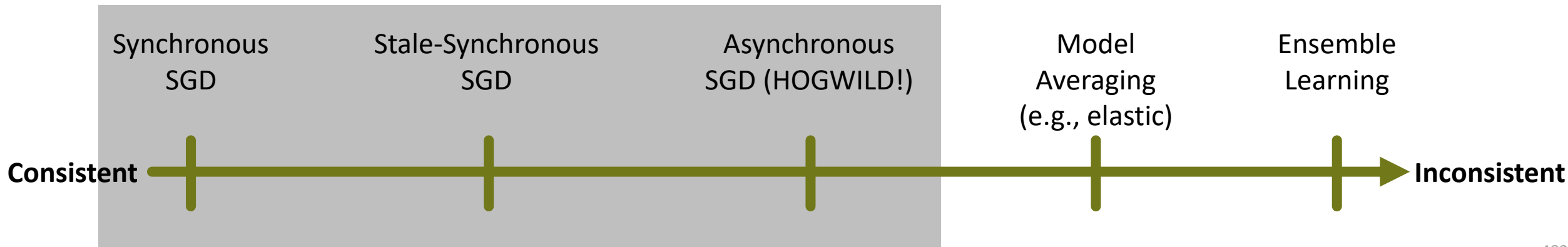
Parameter (and model) consistency - decentralized

- Parameter exchange frequency can be controlled, while still attaining convergence:



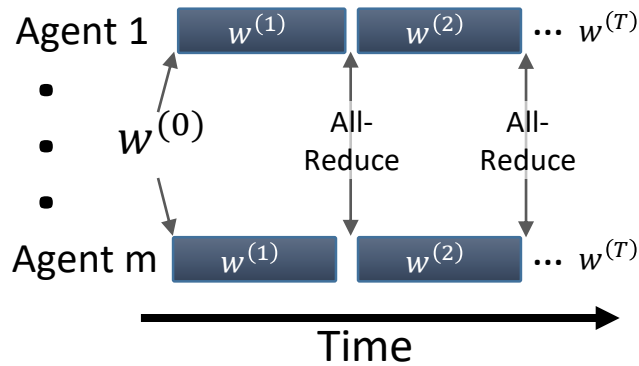
Synchronous

- May also consider limited/slower distribution – gossip [Jin et al. 2016]

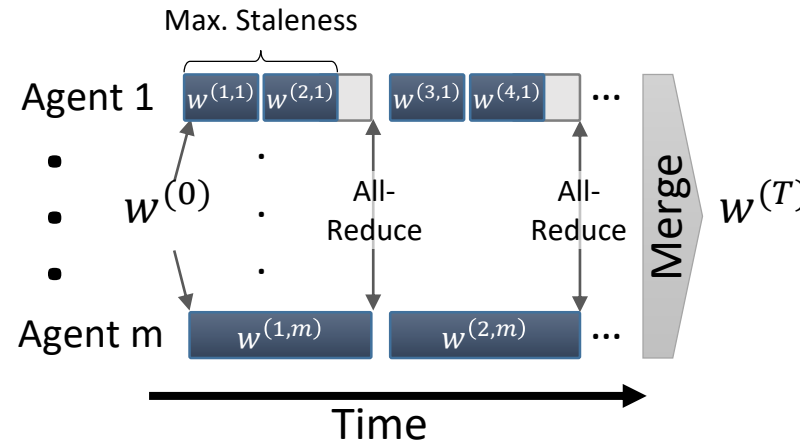


Parameter (and model) consistency - decentralized

- Parameter exchange frequency can be controlled, while still attaining convergence:

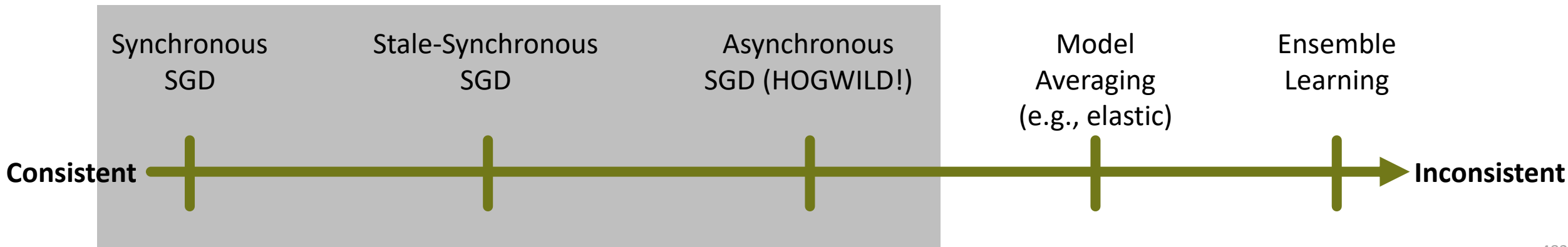


Synchronous



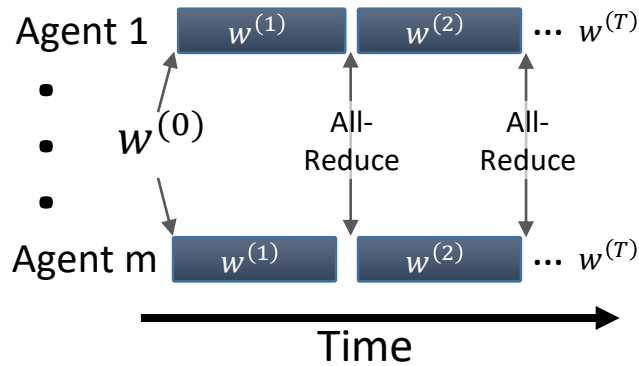
Stale Synchronous / Bounded Asynchronous

- May also consider limited/slower distribution – gossip [Jin et al. 2016]

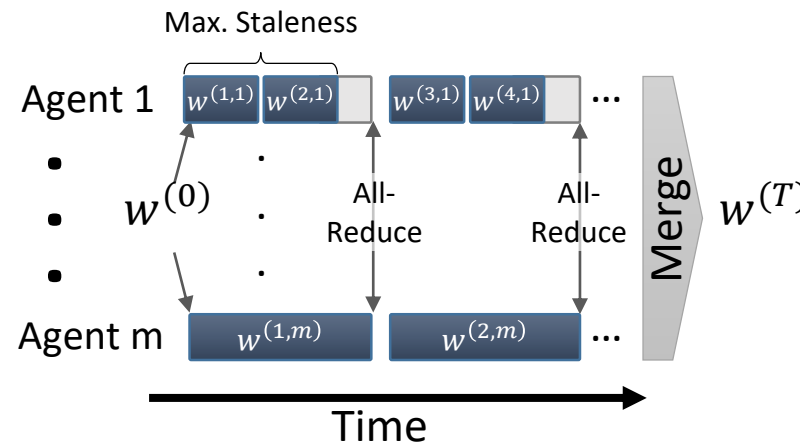


Parameter (and model) consistency - decentralized

- Parameter exchange frequency can be controlled, while still attaining convergence:



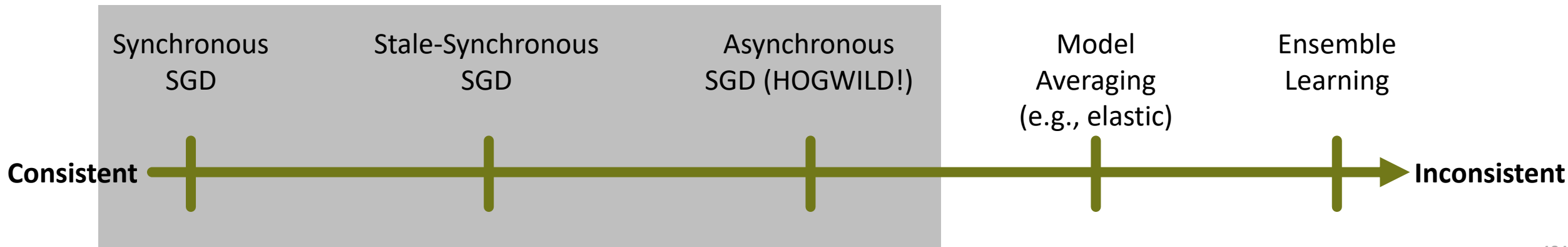
Synchronous



Stale Synchronous / Bounded Asynchronous

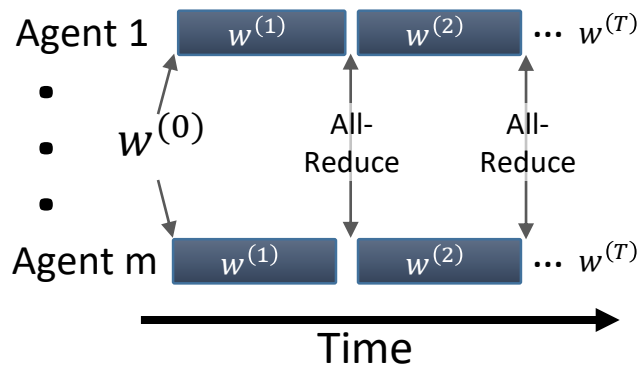
Asynchronous

- May also consider limited/slower distribution – gossip [Jin et al. 2016]

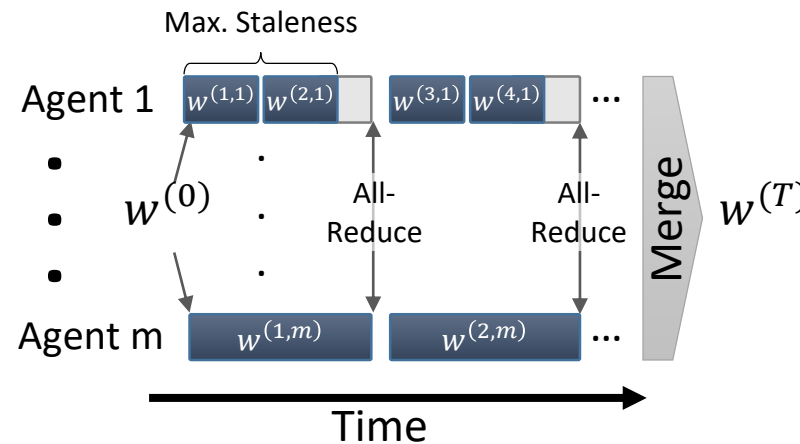


Parameter (and model) consistency - decentralized

- Parameter exchange frequency can be controlled, while still attaining convergence:



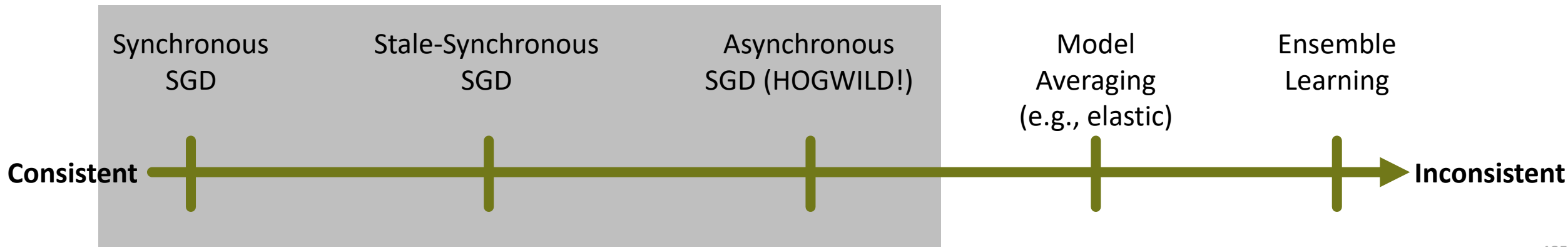
Synchronous



Stale Synchronous / Bounded Asynchronous

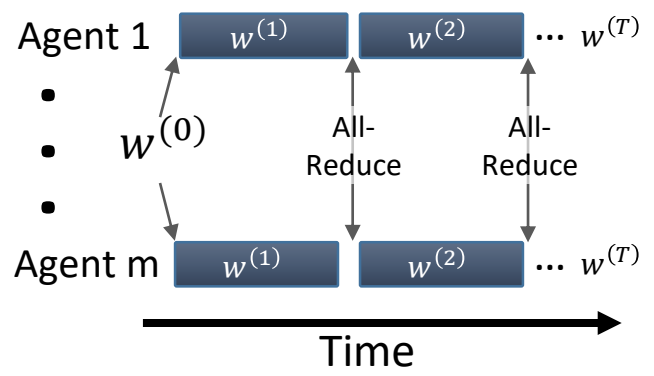
Asynchronous

- May also consider limited/slower distribution – gossip [Jin et al. 2016]

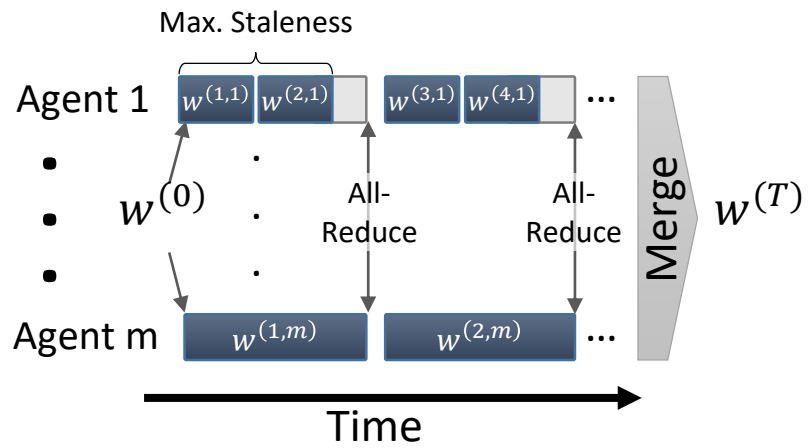


Parameter (and model) consistency - decentralized

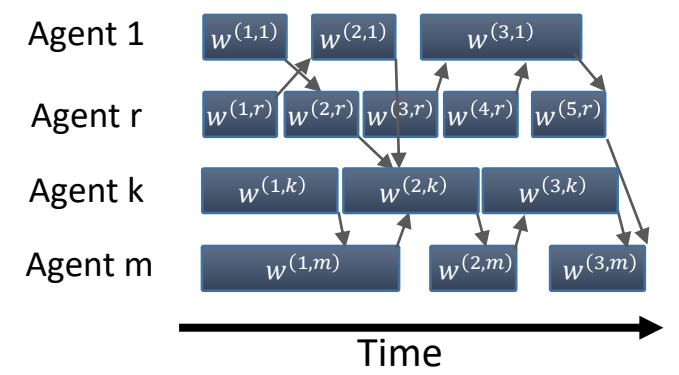
- Parameter exchange frequency can be controlled, while still attaining convergence:



Synchronous

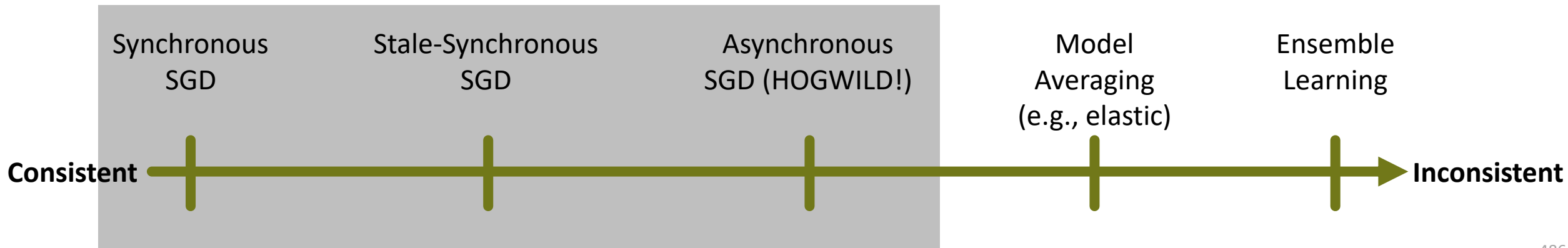


Stale Synchronous / Bounded Asynchronous



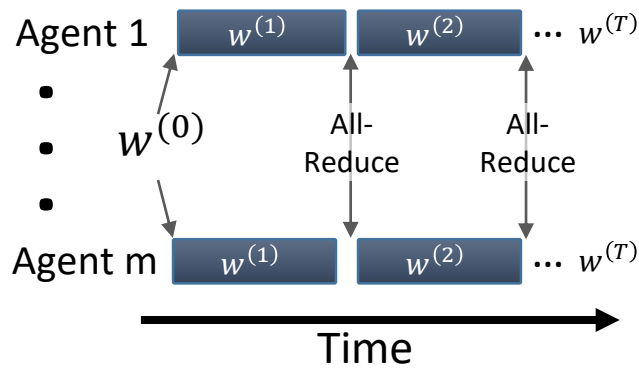
Asynchronous

- May also consider limited/slower distribution – gossip [Jin et al. 2016]

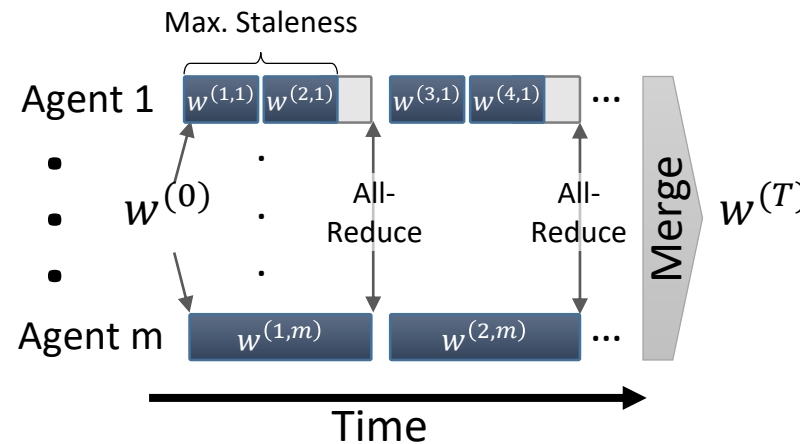


Parameter (and model) consistency - decentralized

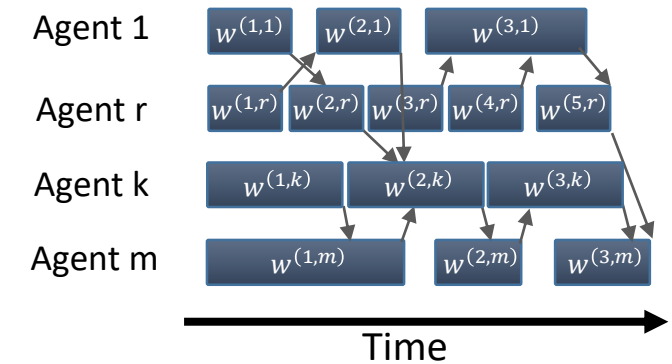
- Parameter exchange frequency can be controlled, while still attaining convergence:



Synchronous

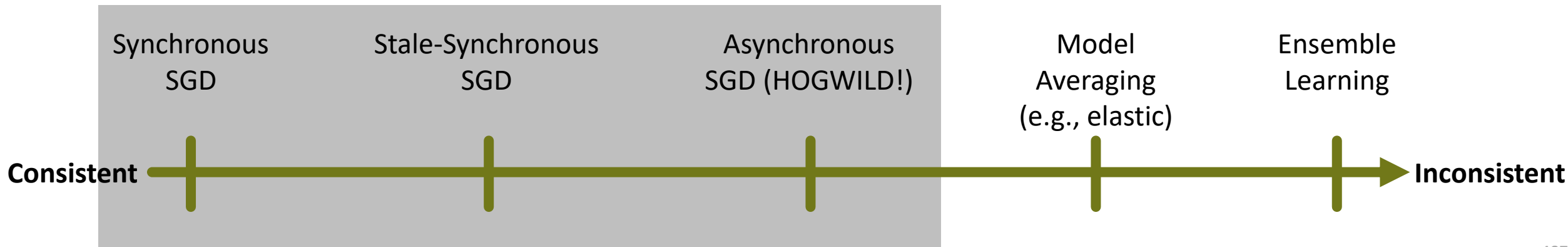


Stale Synchronous / Bounded Asynchronous



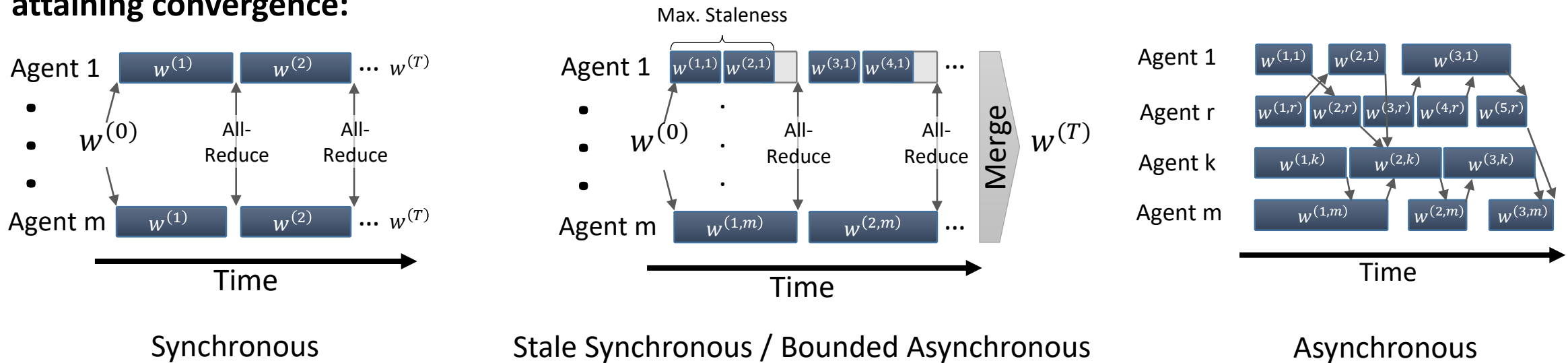
Asynchronous

- May also consider limited/slower distribution – gossip [Jin et al. 2016]

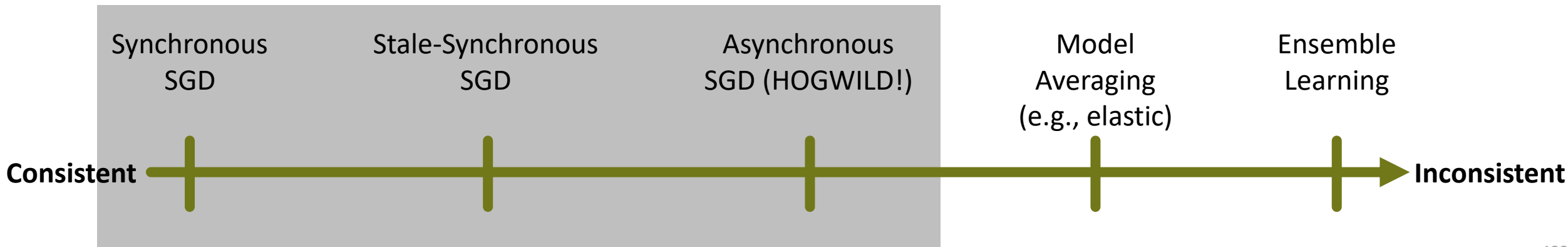


Parameter (and model) consistency - decentralized

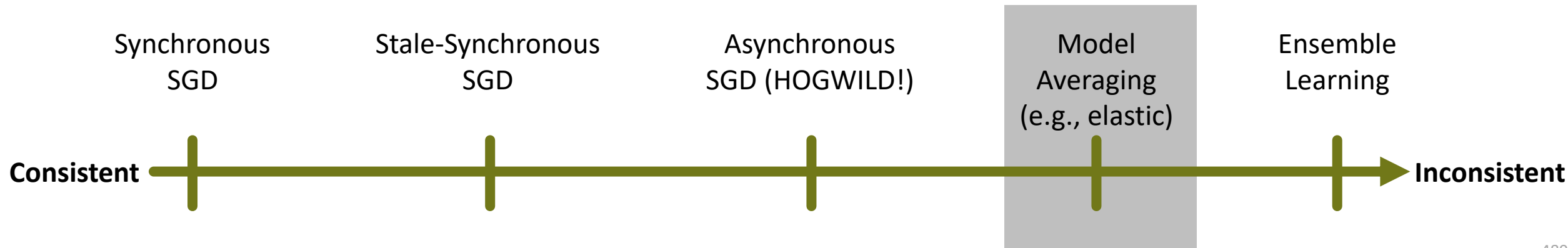
- Parameter exchange frequency can be controlled, while still attaining convergence:



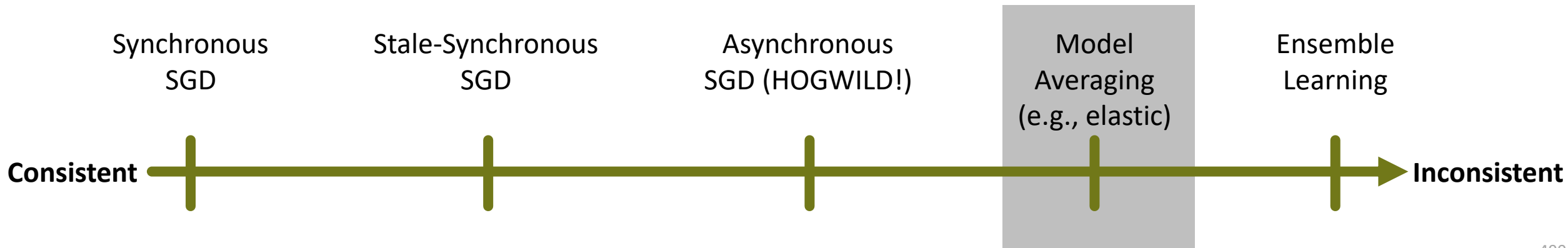
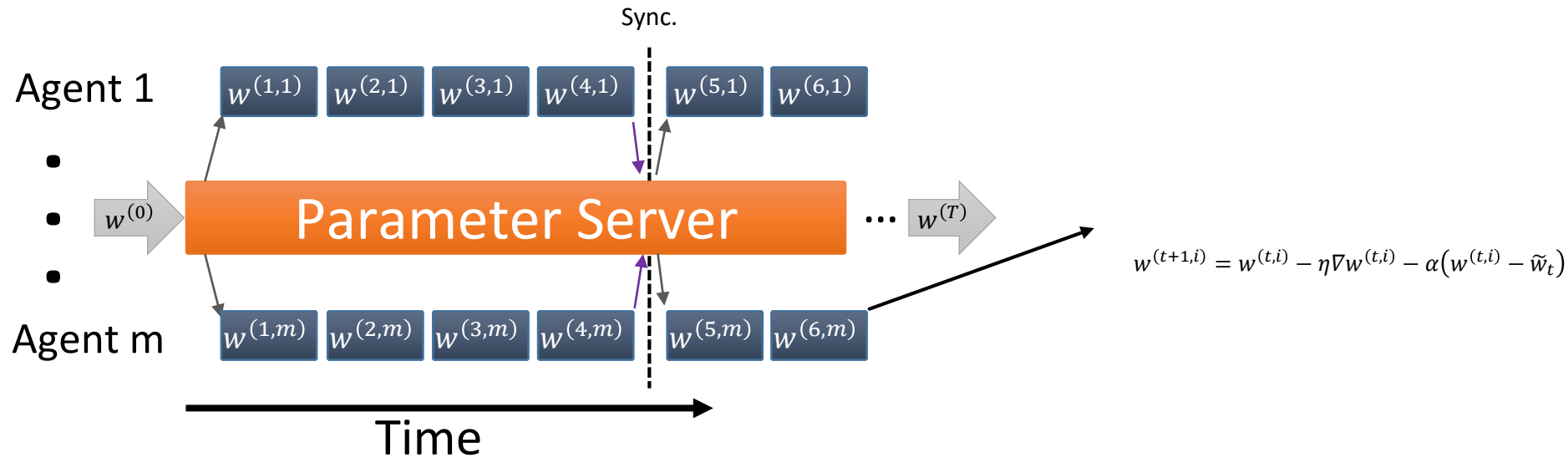
- May also consider limited/slower distribution – gossip [Jin et al. 2016]



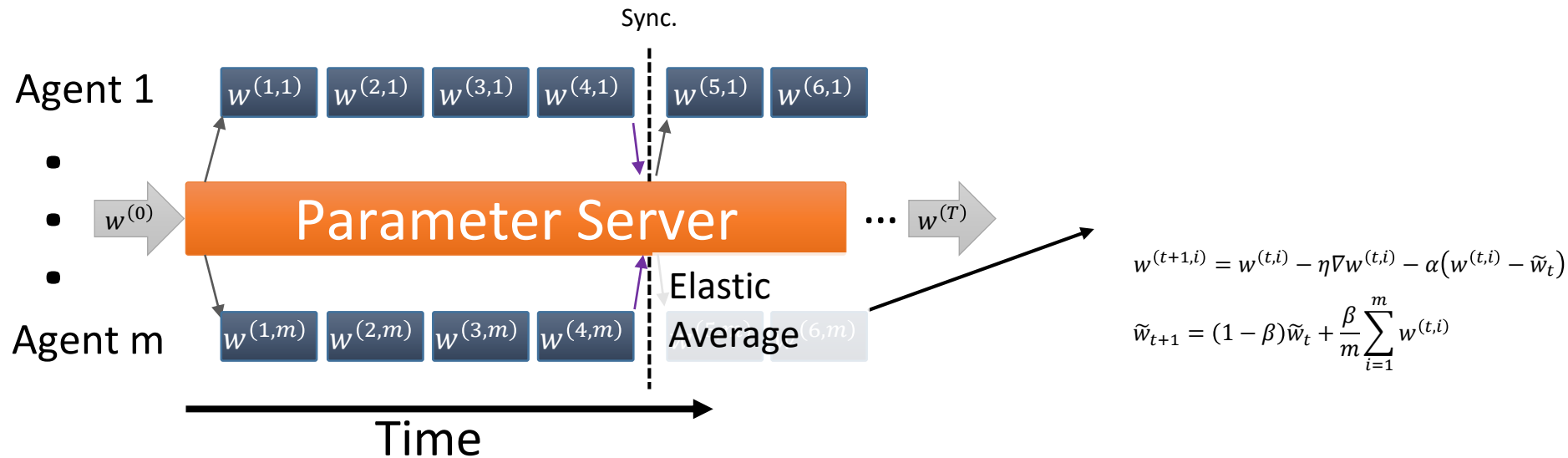
Parameter consistency in deep learning



Parameter consistency in deep learning

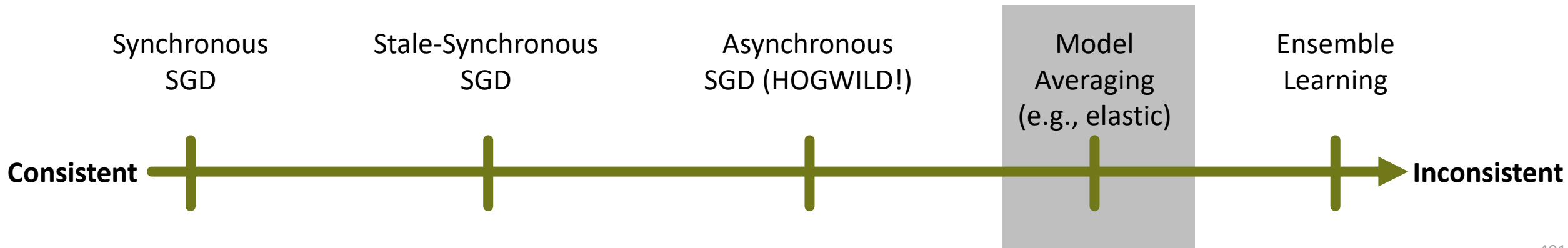


Parameter consistency in deep learning

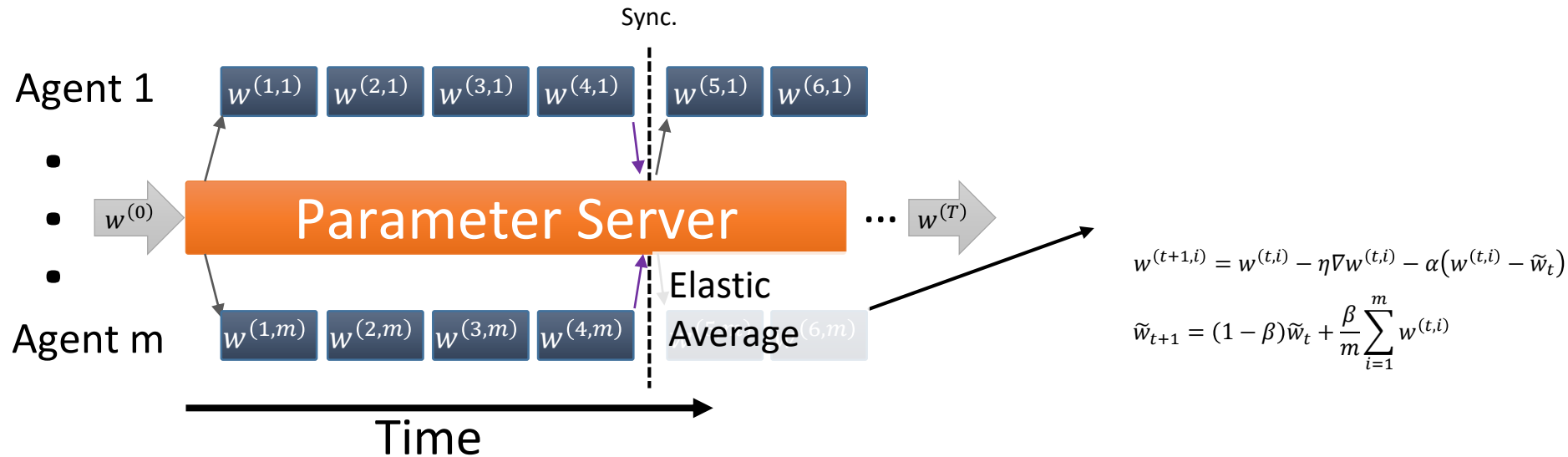


$$w^{(t+1,i)} = w^{(t,i)} - \eta \nabla w^{(t,i)} - \alpha (w^{(t,i)} - \tilde{w}_t)$$

$$\tilde{w}_{t+1} = (1 - \beta) \tilde{w}_t + \frac{\beta}{m} \sum_{i=1}^m w^{(t,i)}$$

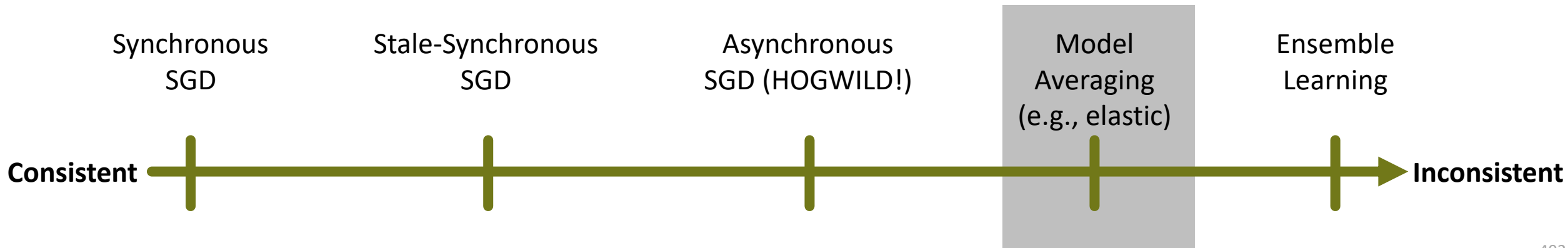


Parameter consistency in deep learning

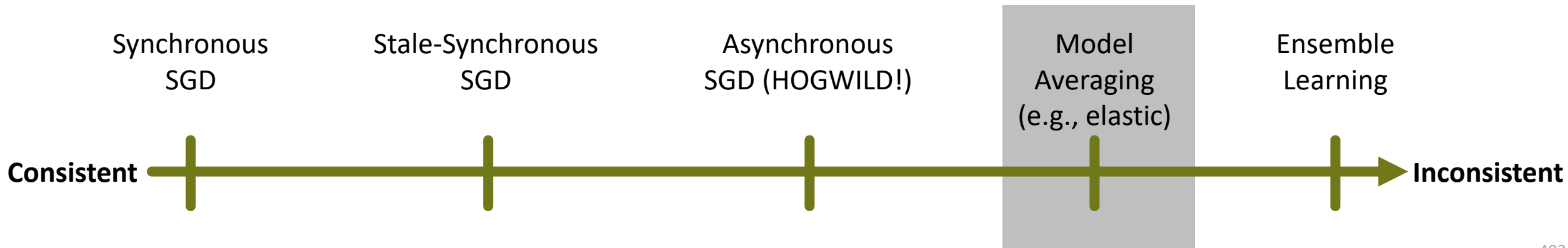
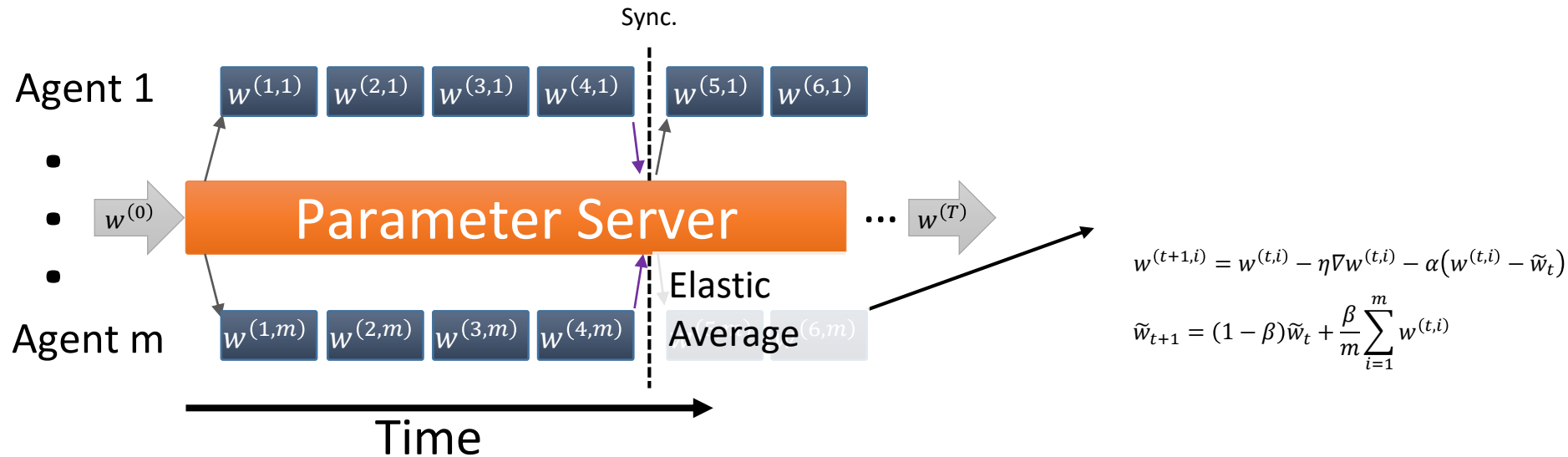


$$w^{(t+1,i)} = w^{(t,i)} - \eta \nabla w^{(t,i)} - \alpha (w^{(t,i)} - \tilde{w}_t)$$

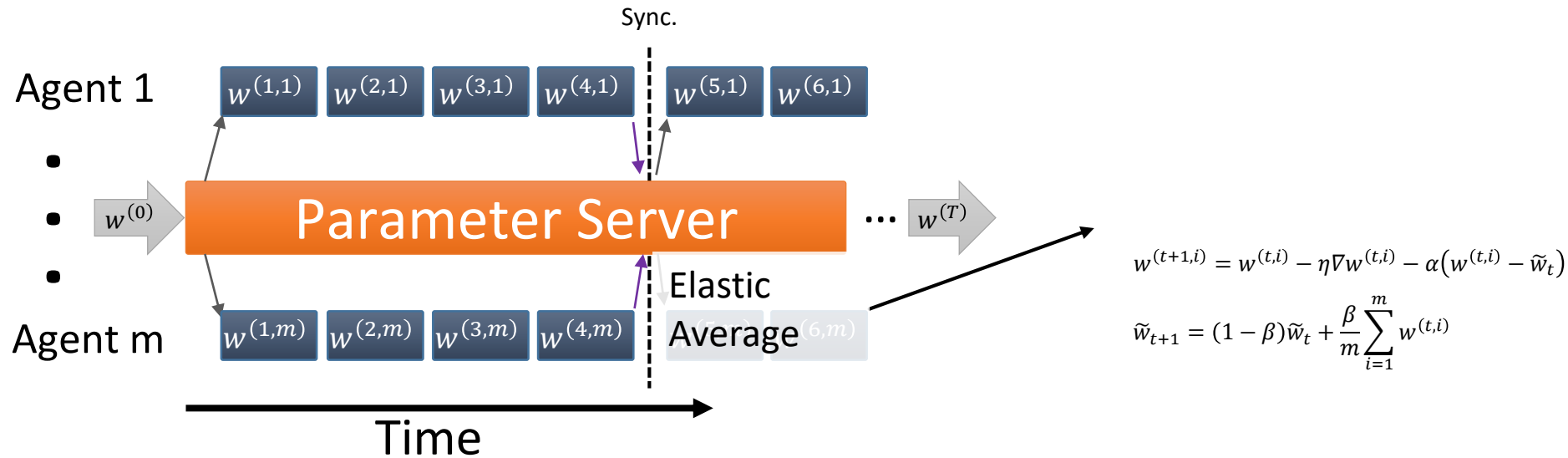
$$\tilde{w}_{t+1} = (1 - \beta) \tilde{w}_t + \frac{\beta}{m} \sum_{i=1}^m w^{(t,i)}$$



Parameter consistency in deep learning

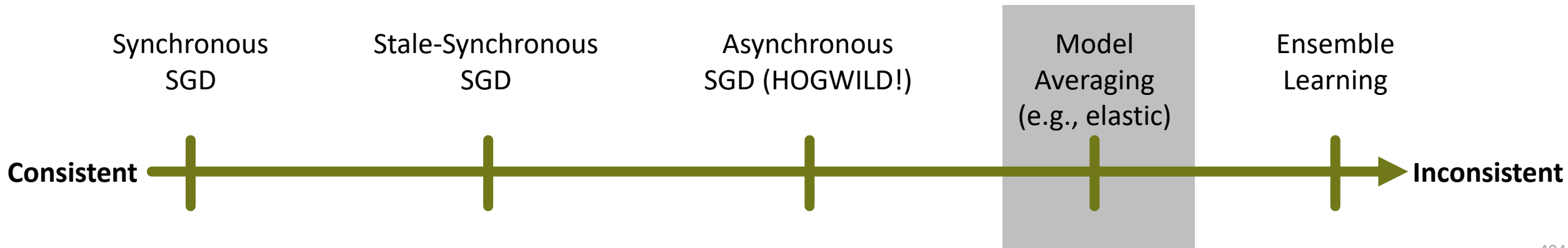


Parameter consistency in deep learning

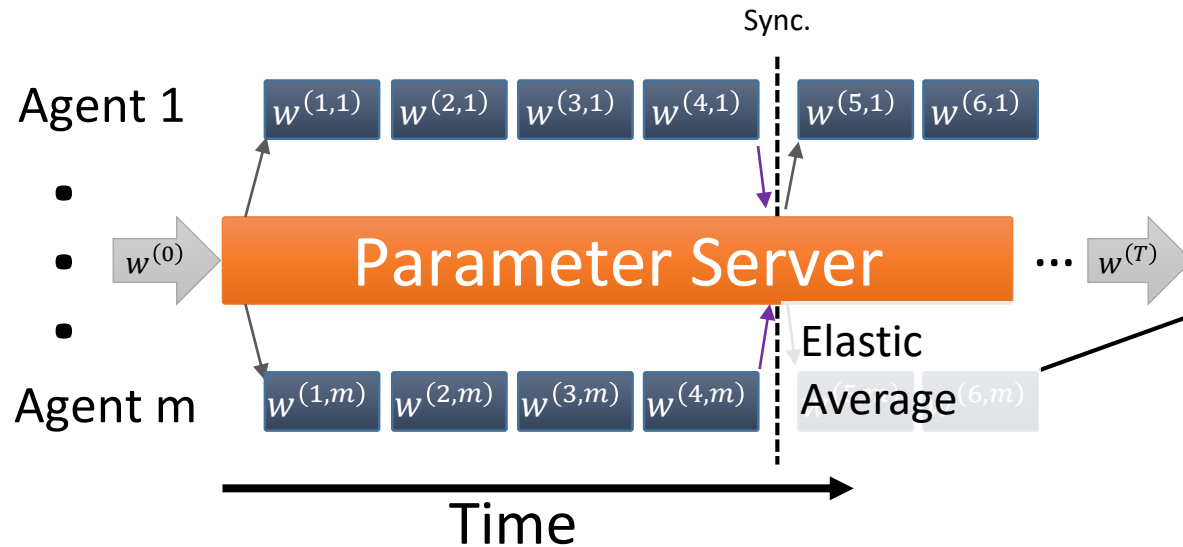


$$w^{(t+1,i)} = w^{(t,i)} - \eta \nabla w^{(t,i)} - \alpha (w^{(t,i)} - \tilde{w}_t)$$

$$\tilde{w}_{t+1} = (1 - \beta) \tilde{w}_t + \frac{\beta}{m} \sum_{i=1}^m w^{(t,i)}$$



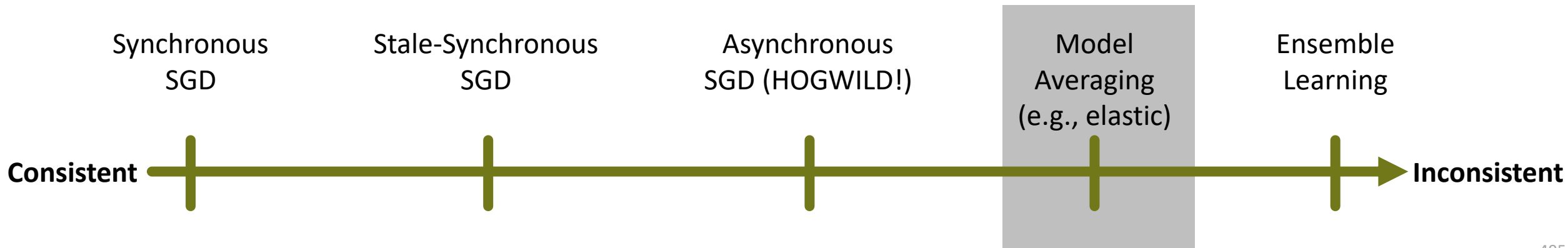
Parameter consistency in deep learning



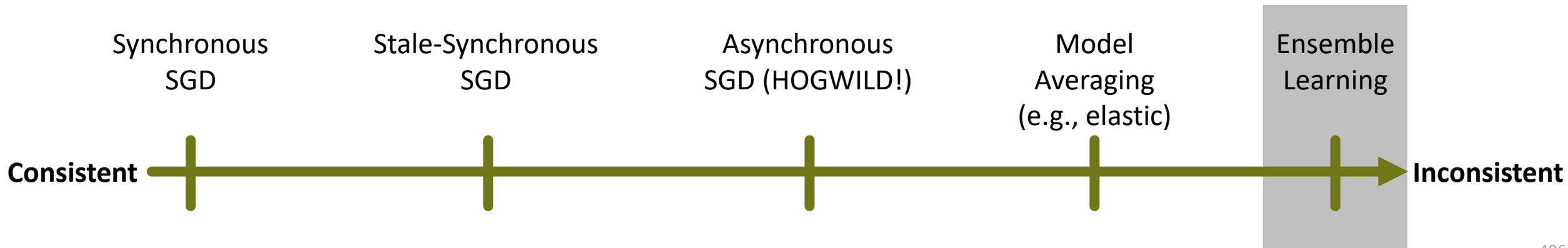
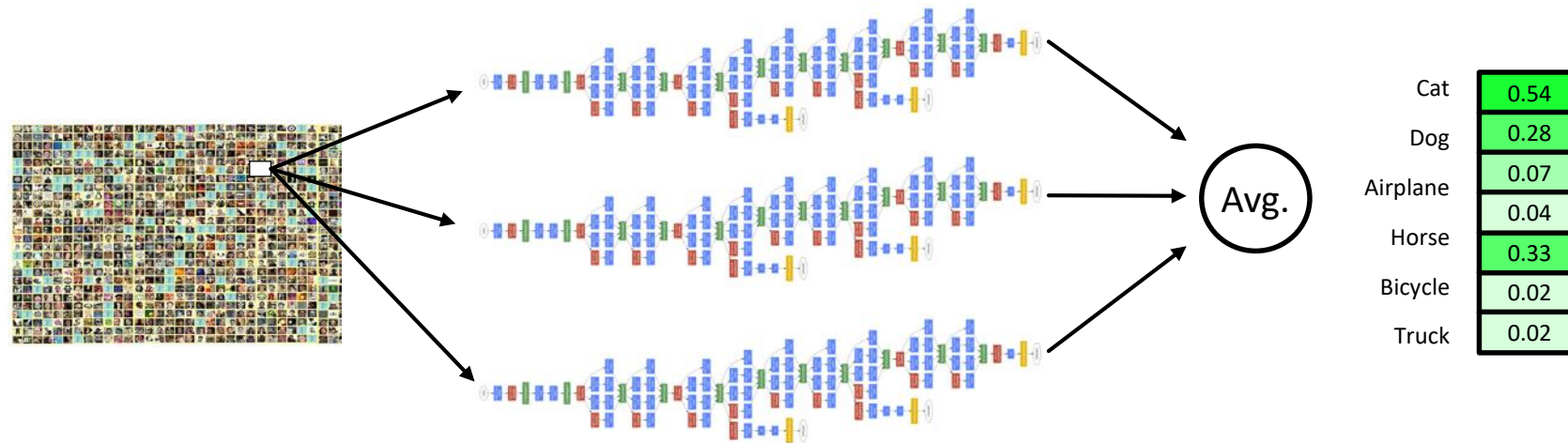
Using physical forces between different versions of w :

$$w^{(t+1,i)} = w^{(t,i)} - \eta \nabla w^{(t,i)} - \alpha (w^{(t,i)} - \tilde{w}_t)$$

$$\tilde{w}_{t+1} = (1 - \beta) \tilde{w}_t + \frac{\beta}{m} \sum_{i=1}^m w^{(t,i)}$$



Parameter consistency in deep learning



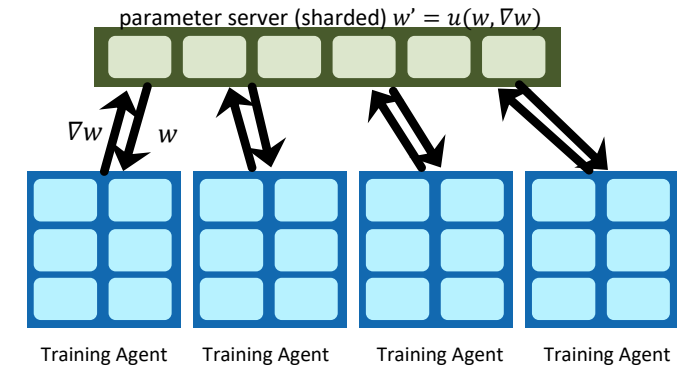
Communication optimization

Communication optimization

- **Lossy compression: trade off latency (local compute) for bandwidth**
- **Sufficient factor broadcasting**
- **Quantization:**
 - 16-bit (IEEE FP16, bfloat16) becoming standard
 - QSGD (stochastic rounding) [Alistarh et al. 2016]
 - 1-bit SGD [Seide et al. 2014; Dryden et al. 2016]
 - Error feedback is important!
- **Sparsification:**
 - Top-k SGD [Renggli et al. 2019]
 - Skip small weight updates

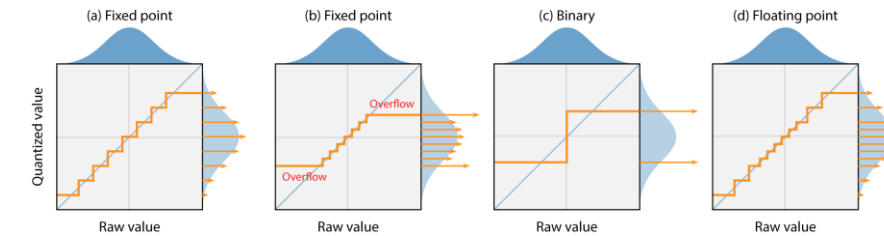
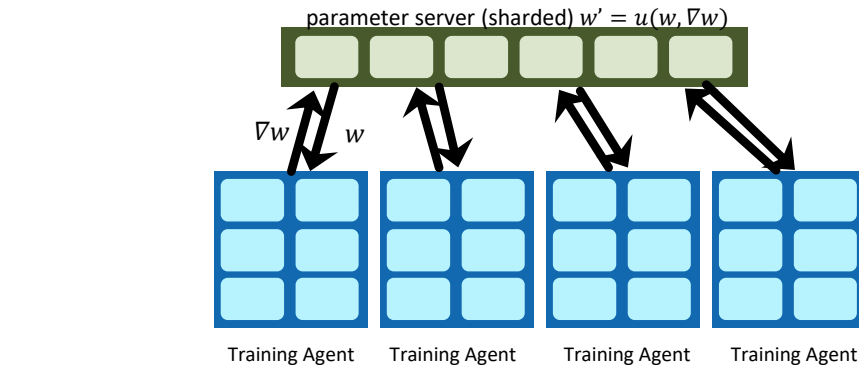
Communication optimization

- **Lossy compression:** trade off latency (local compute) for bandwidth
- **Sufficient factor broadcasting**
- **Quantization:**
 - 16-bit (IEEE FP16, bfloat16) becoming standard
 - QSGD (stochastic rounding) [Alistarh et al. 2016]
 - 1-bit SGD [Seide et al. 2014; Dryden et al. 2016]
 - Error feedback is important!
- **Sparsification:**
 - Top-k SGD [Renggli et al. 2019]
 - Skip small weight updates



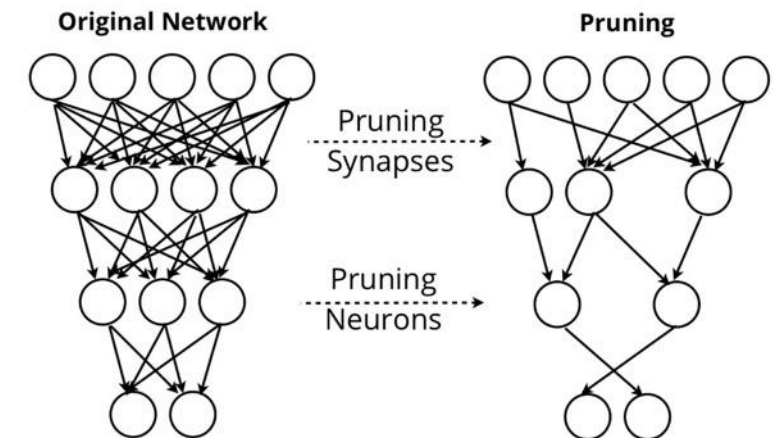
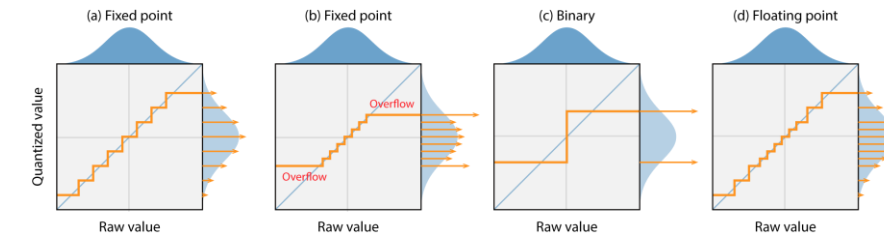
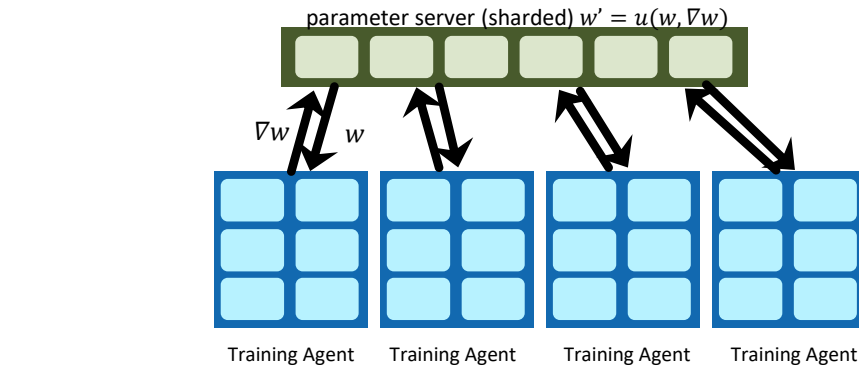
Communication optimization

- **Lossy compression: trade off latency (local compute) for bandwidth**
- **Sufficient factor broadcasting**
- **Quantization:**
 - 16-bit (IEEE FP16, bfloat16) becoming standard
 - QSGD (stochastic rounding) [Alistarh et al. 2016]
 - 1-bit SGD [Seide et al. 2014; Dryden et al. 2016]
 - Error feedback is important!
- **Sparsification:**
 - Top-k SGD [Renggli et al. 2019]
 - Skip small weight updates

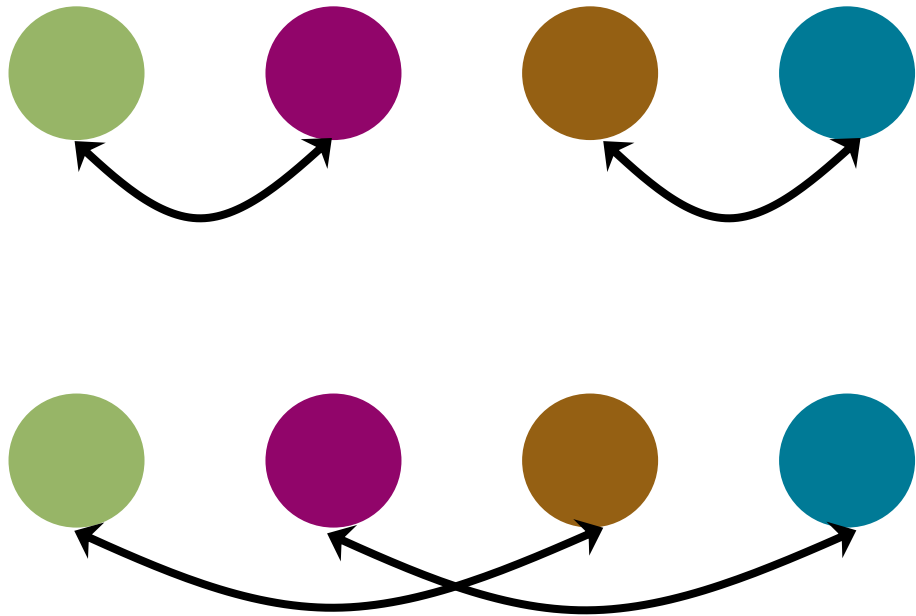


Communication optimization

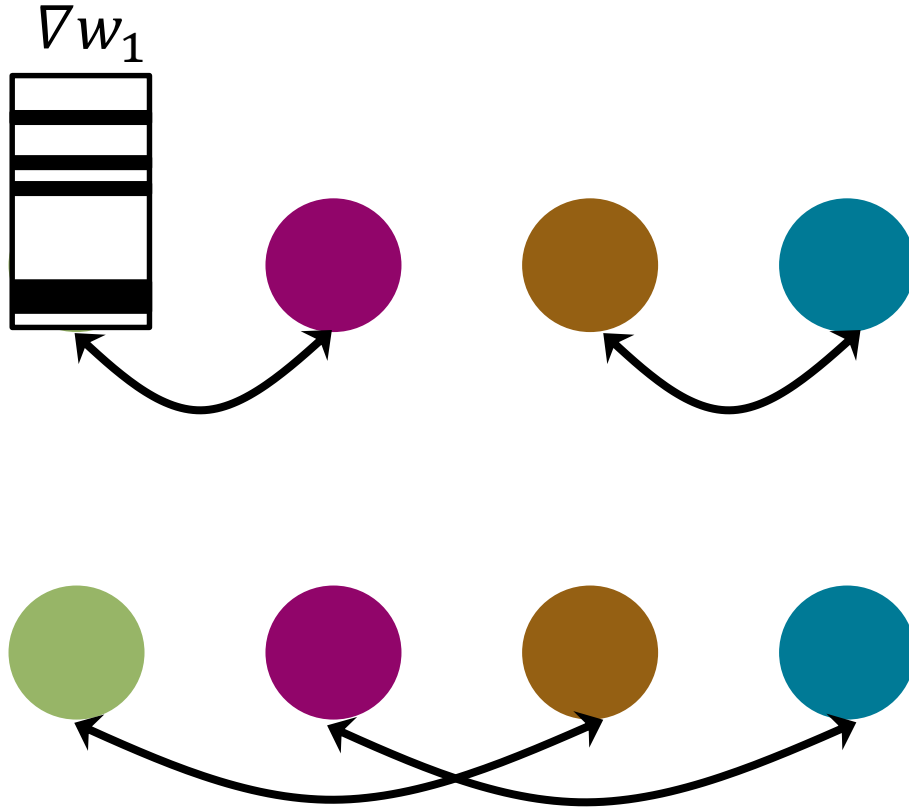
- **Lossy compression: trade off latency (local compute) for bandwidth**
- **Sufficient factor broadcasting**
- **Quantization:**
 - 16-bit (IEEE FP16, bfloat16) becoming standard
 - QSGD (stochastic rounding) [Alistarh et al. 2016]
 - 1-bit SGD [Seide et al. 2014; Dryden et al. 2016]
 - Error feedback is important!
- **Sparsification:**
 - Top-k SGD [Renggli et al. 2019]
 - Skip small weight updates



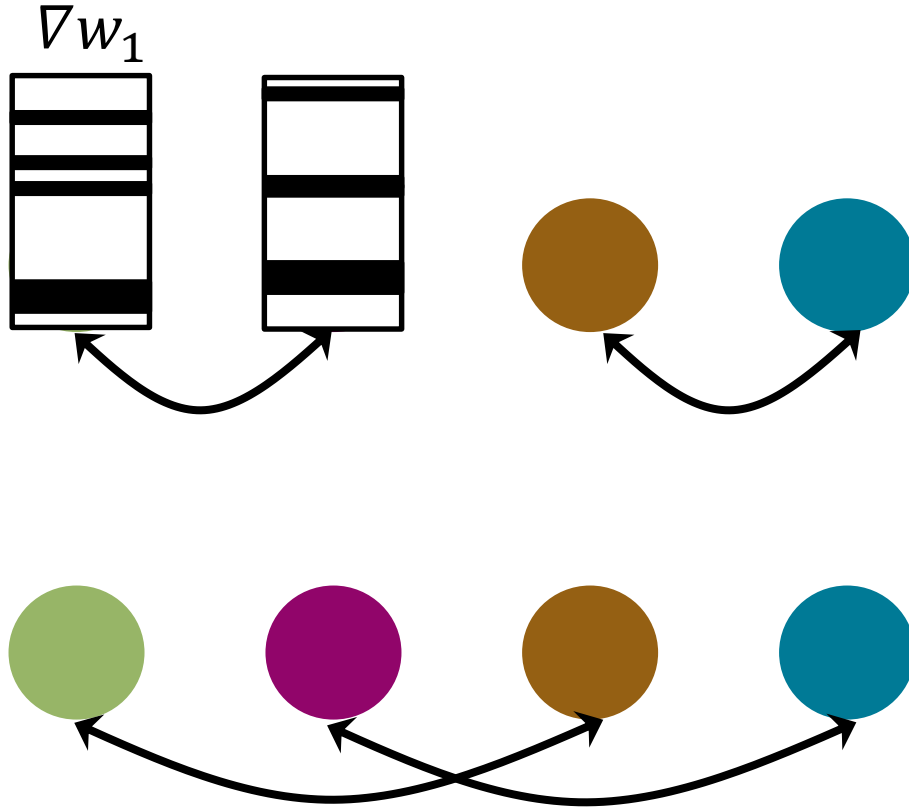
SparCML – Quantized sparse allreduce



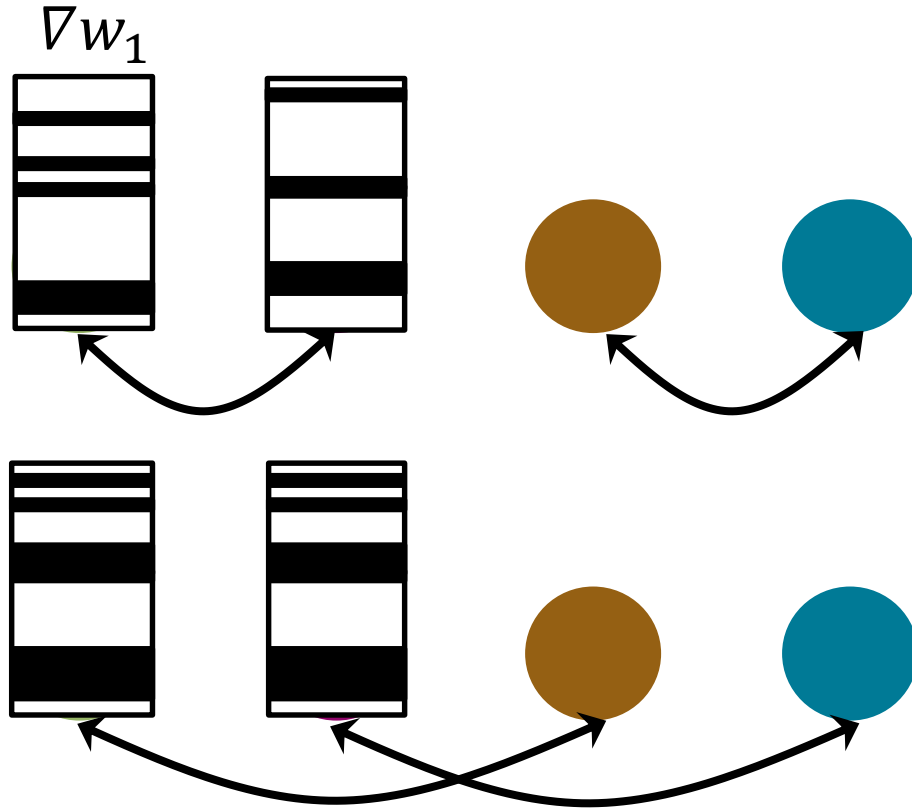
SparCML – Quantized sparse allreduce



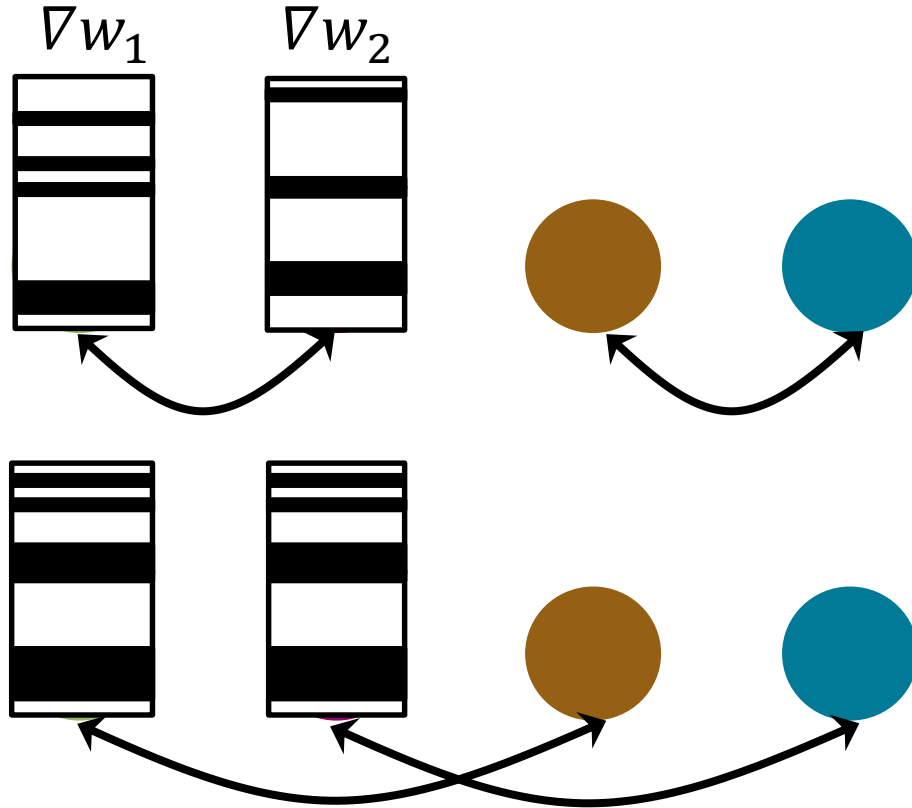
SparCML – Quantized sparse allreduce



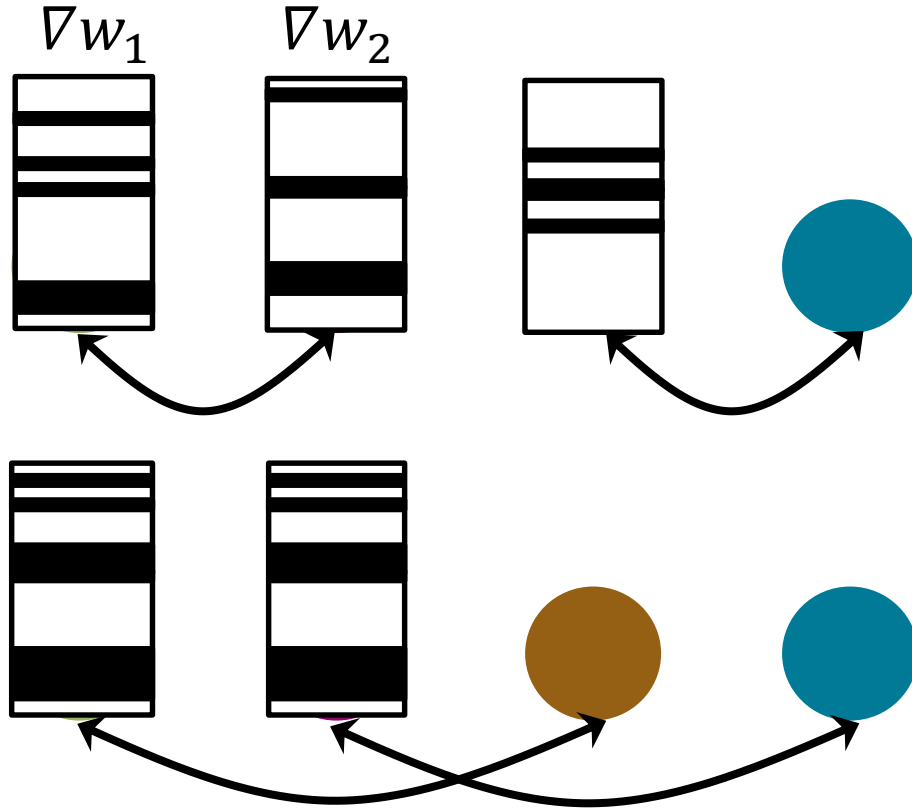
SparCML – Quantized sparse allreduce



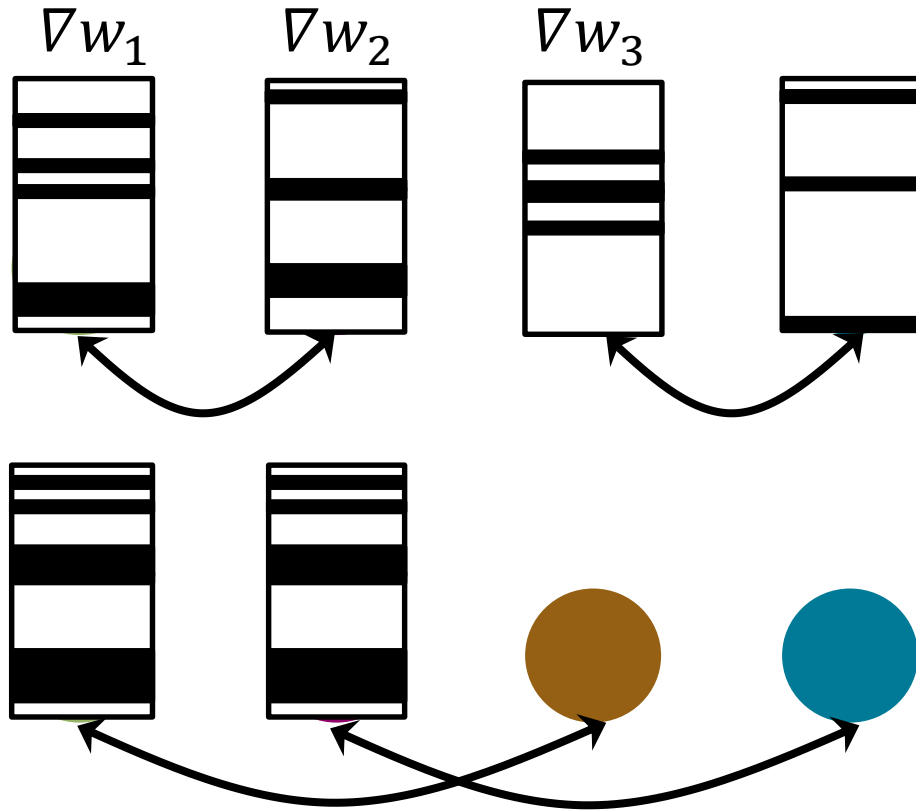
SparCML – Quantized sparse allreduce



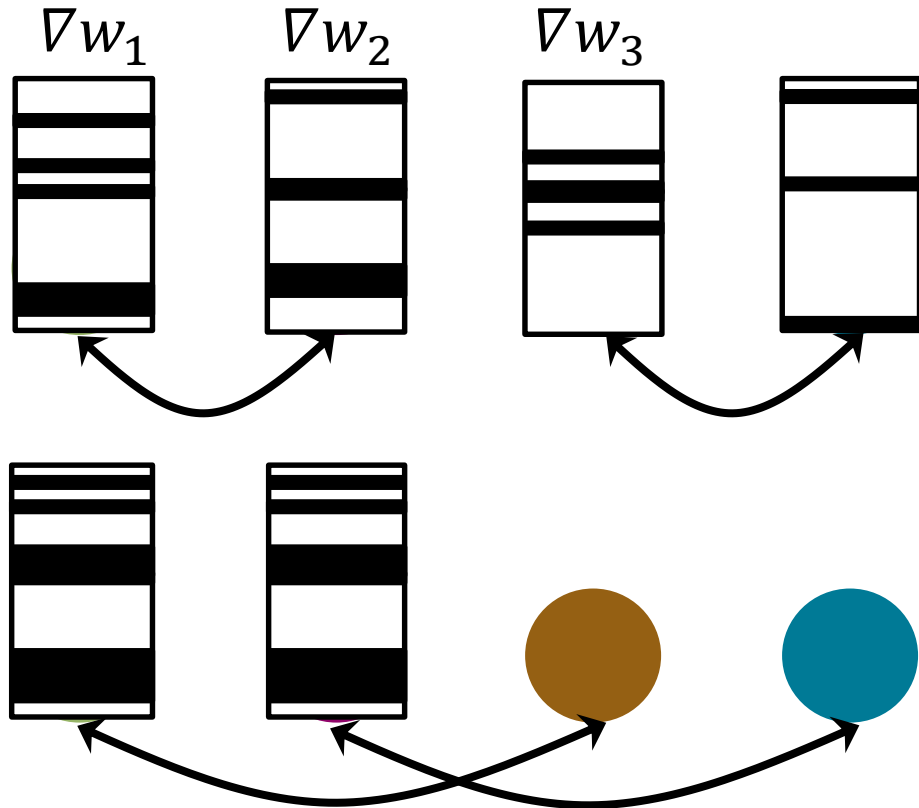
SparCML – Quantized sparse allreduce



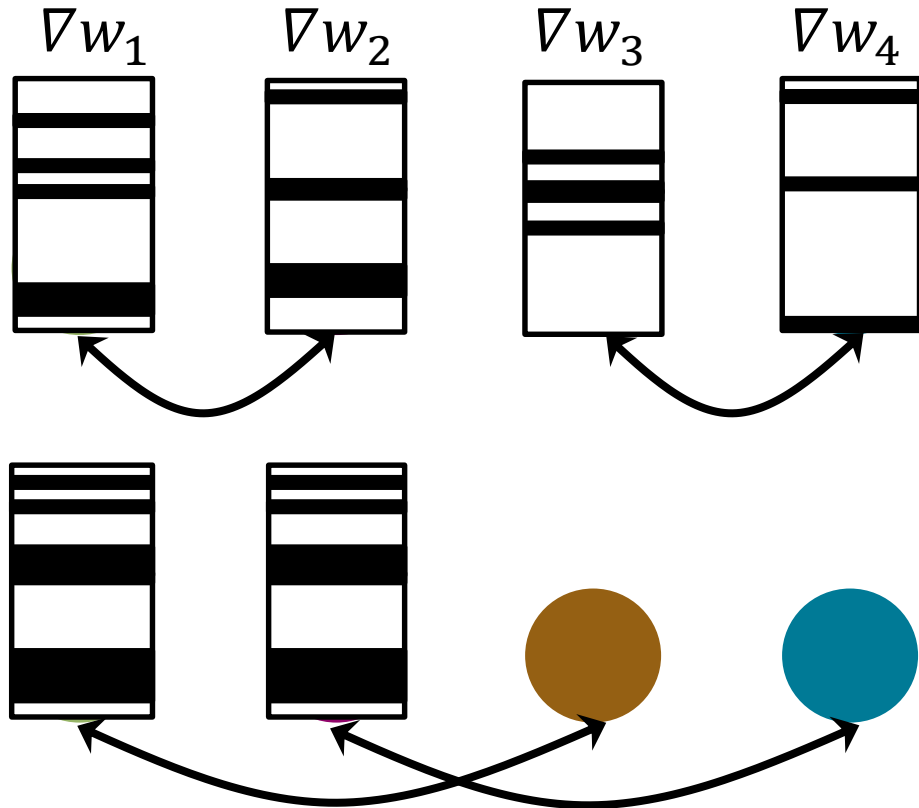
SparCML – Quantized sparse allreduce



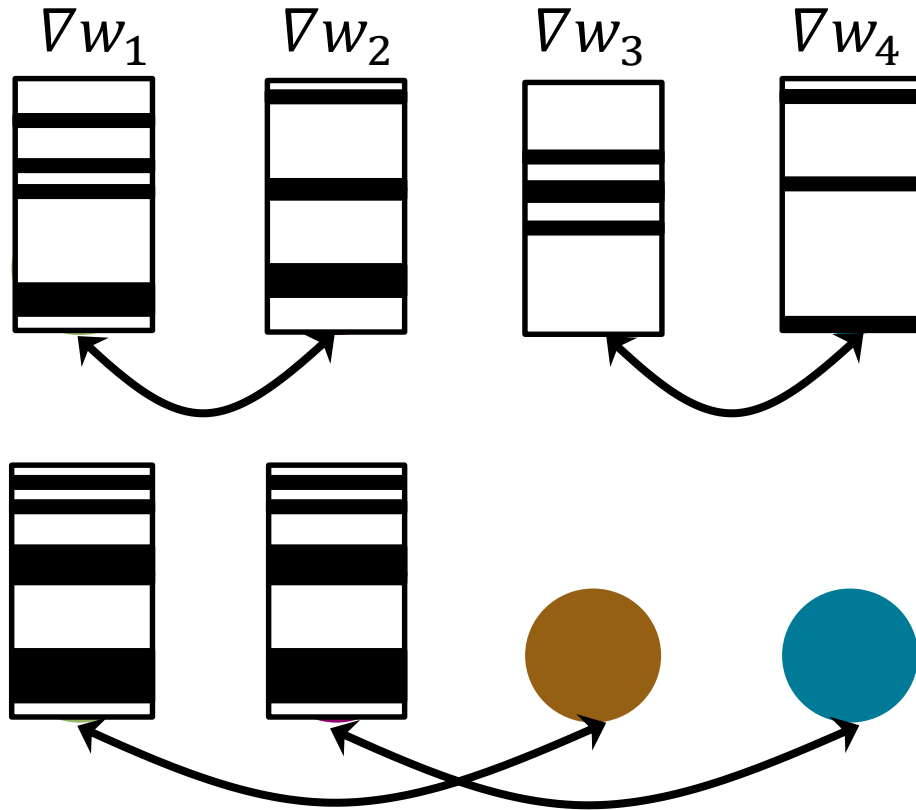
SparCML – Quantized sparse allreduce



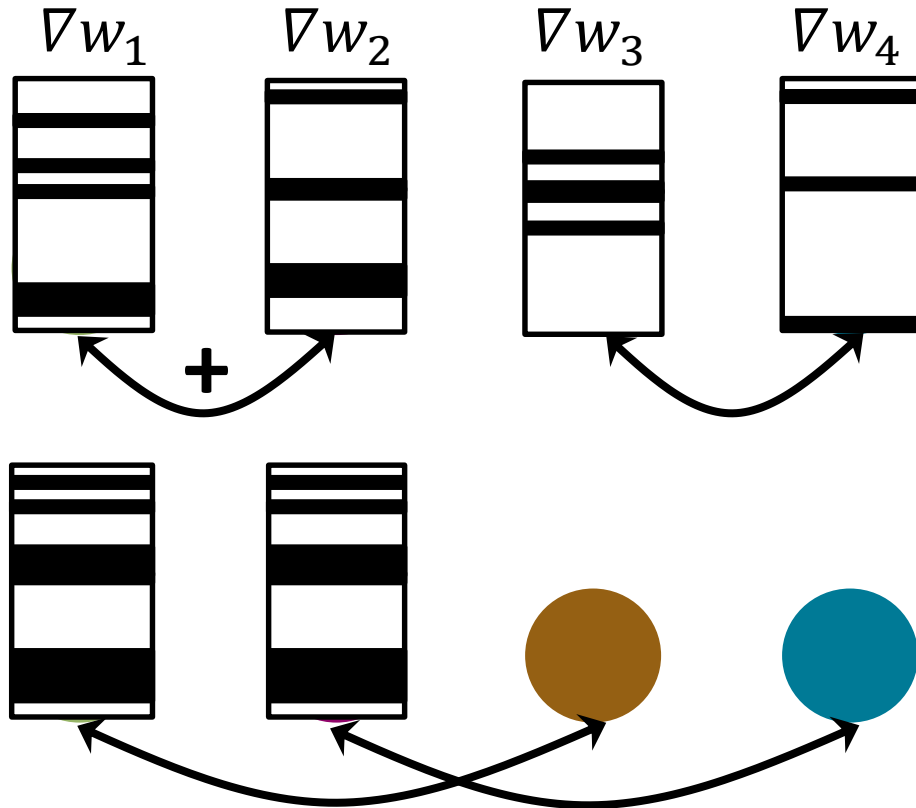
SparCML – Quantized sparse allreduce



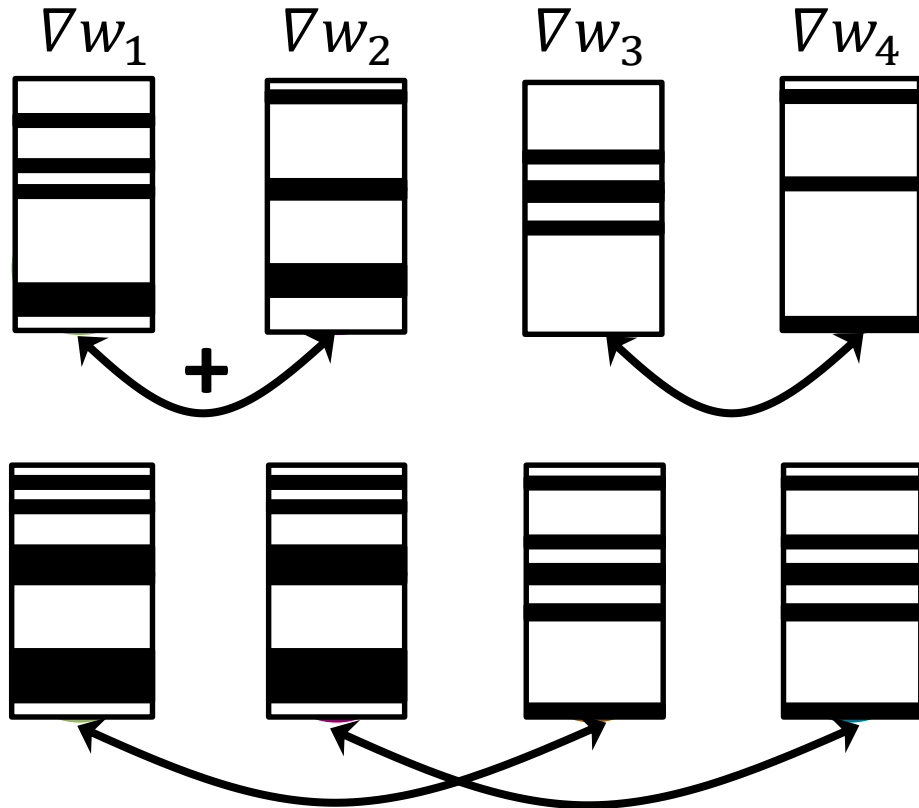
SparCML – Quantized sparse allreduce



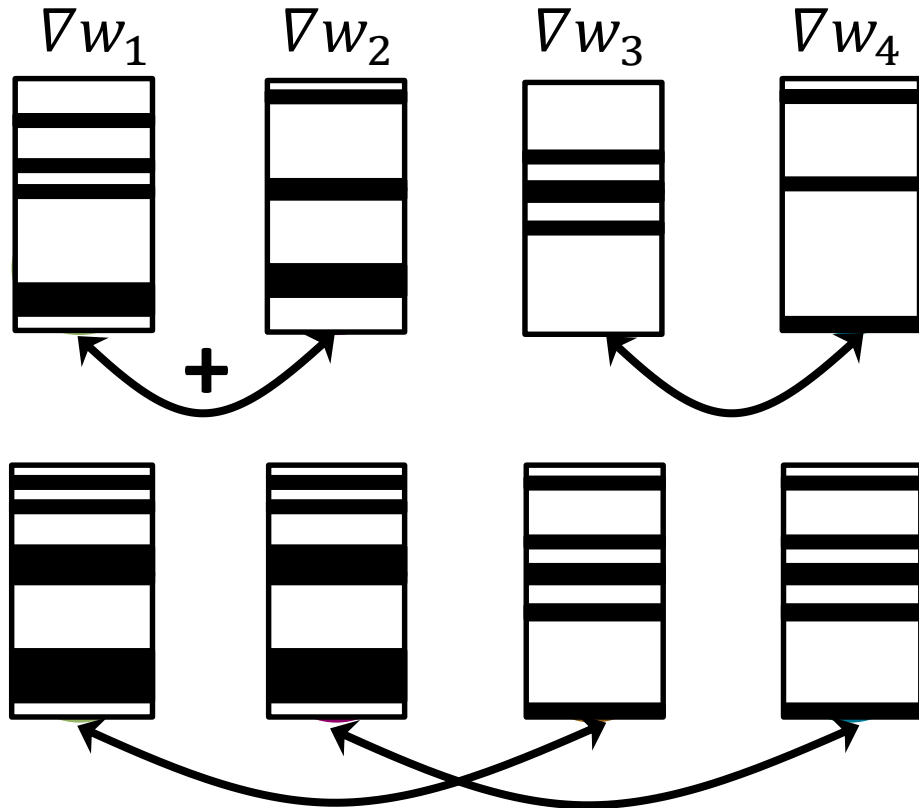
SparCML – Quantized sparse allreduce



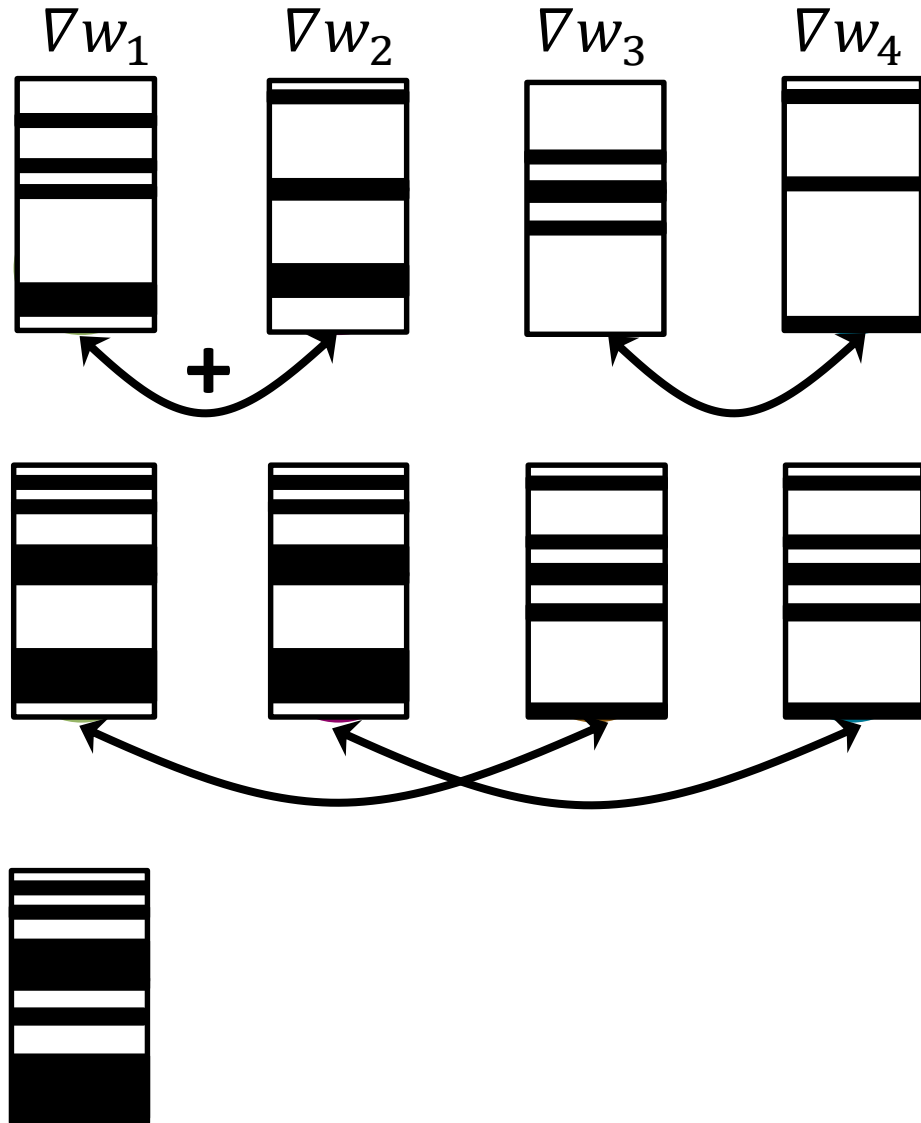
SparCML – Quantized sparse allreduce



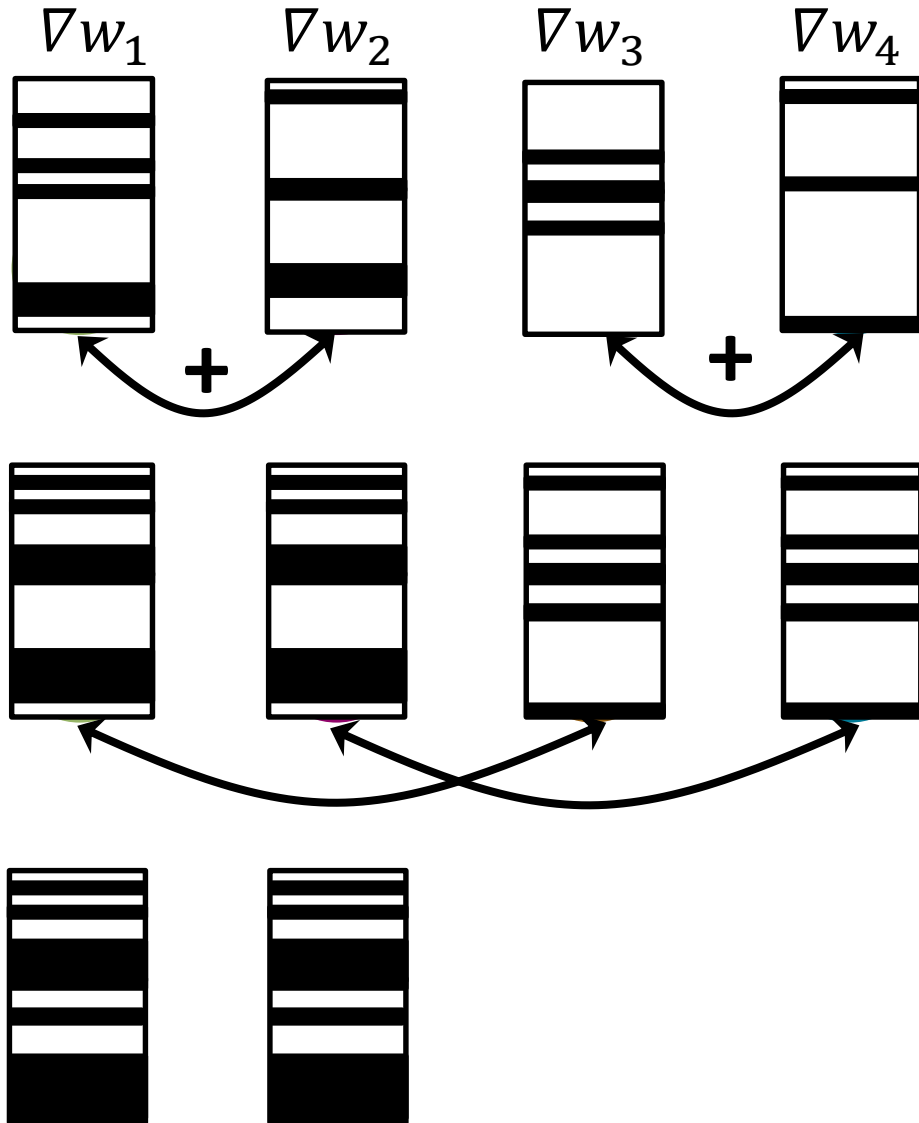
SparCML – Quantized sparse allreduce



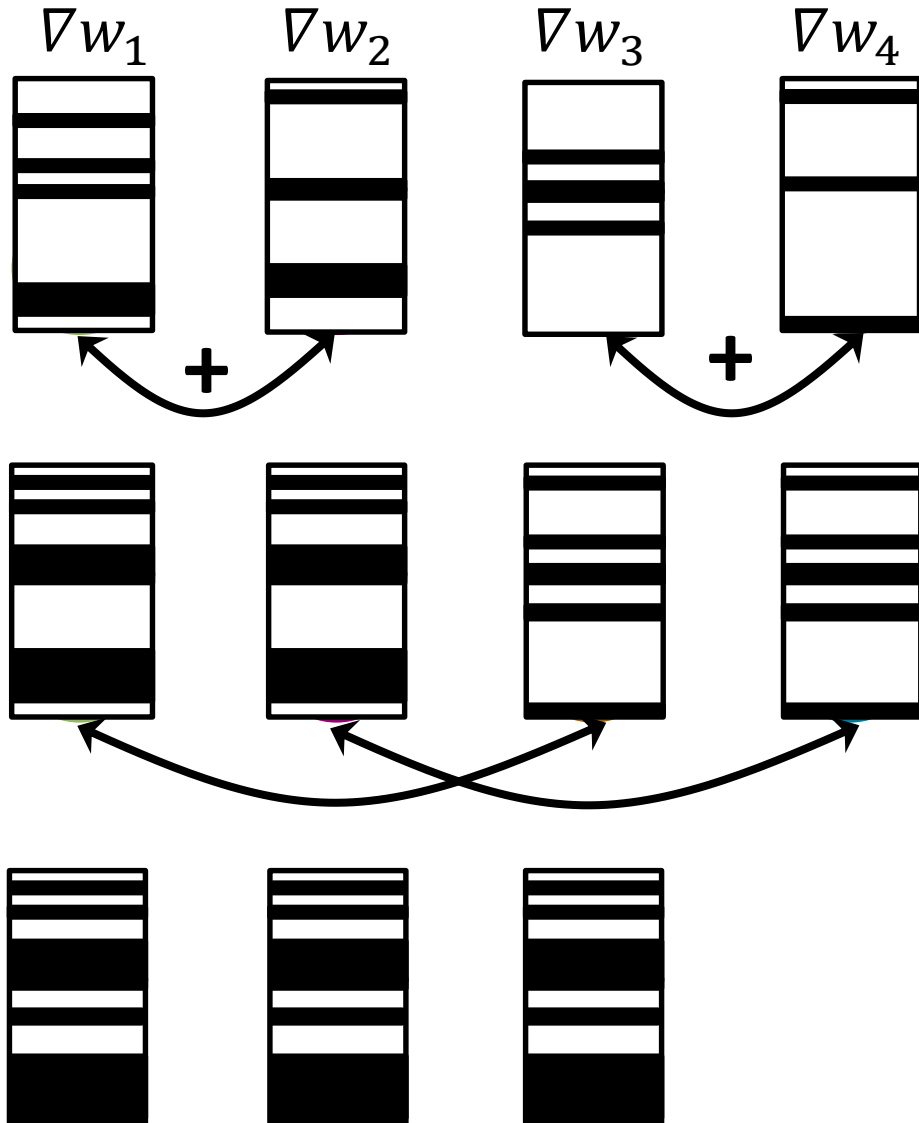
SparCML – Quantized sparse allreduce



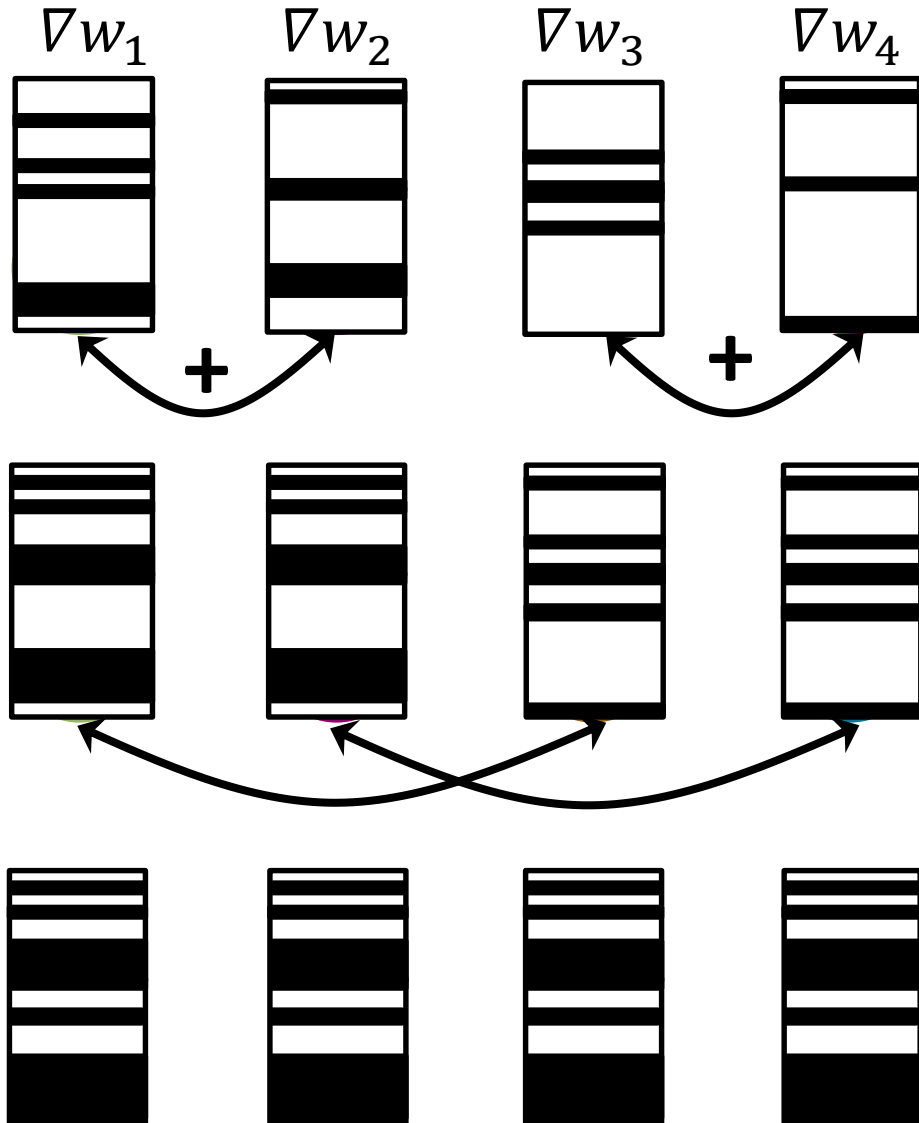
SparCML – Quantized sparse allreduce



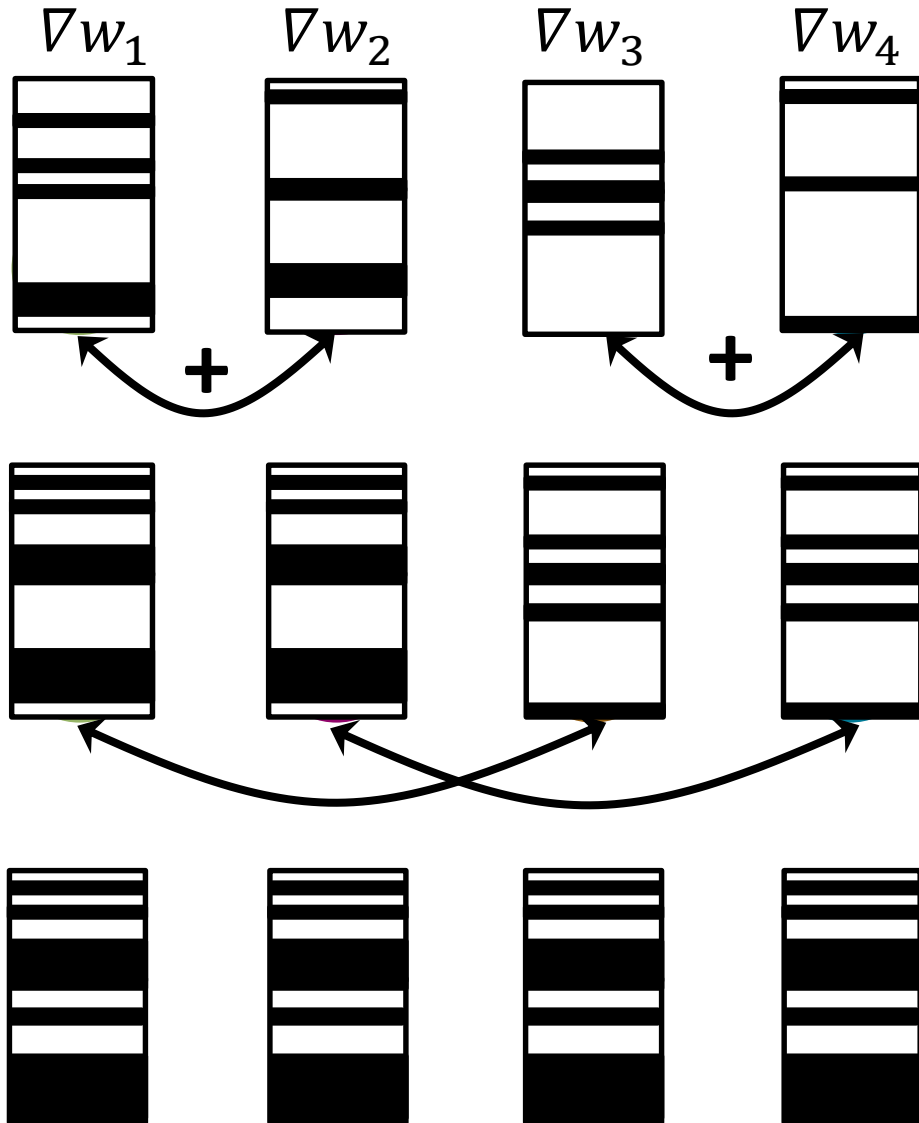
SparCML – Quantized sparse allreduce



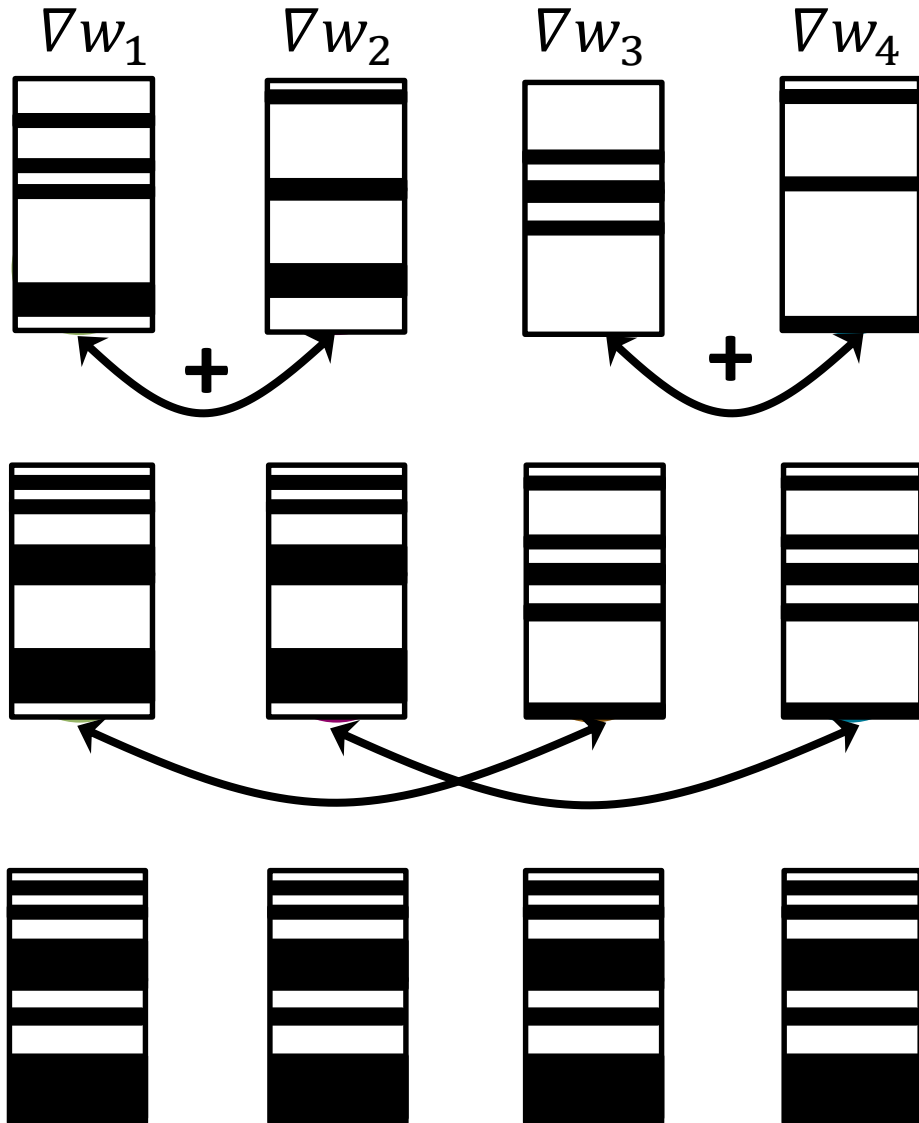
SparCML – Quantized sparse allreduce



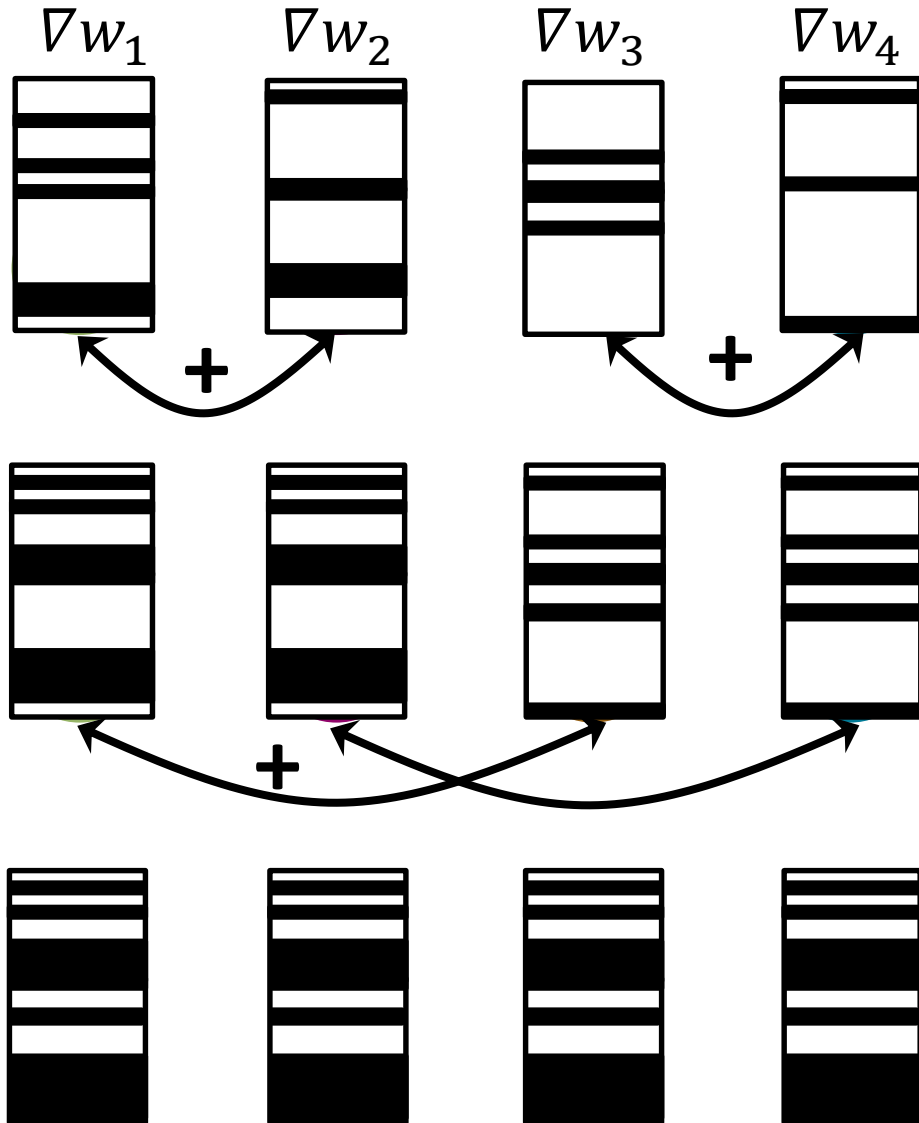
SparCML – Quantized sparse allreduce



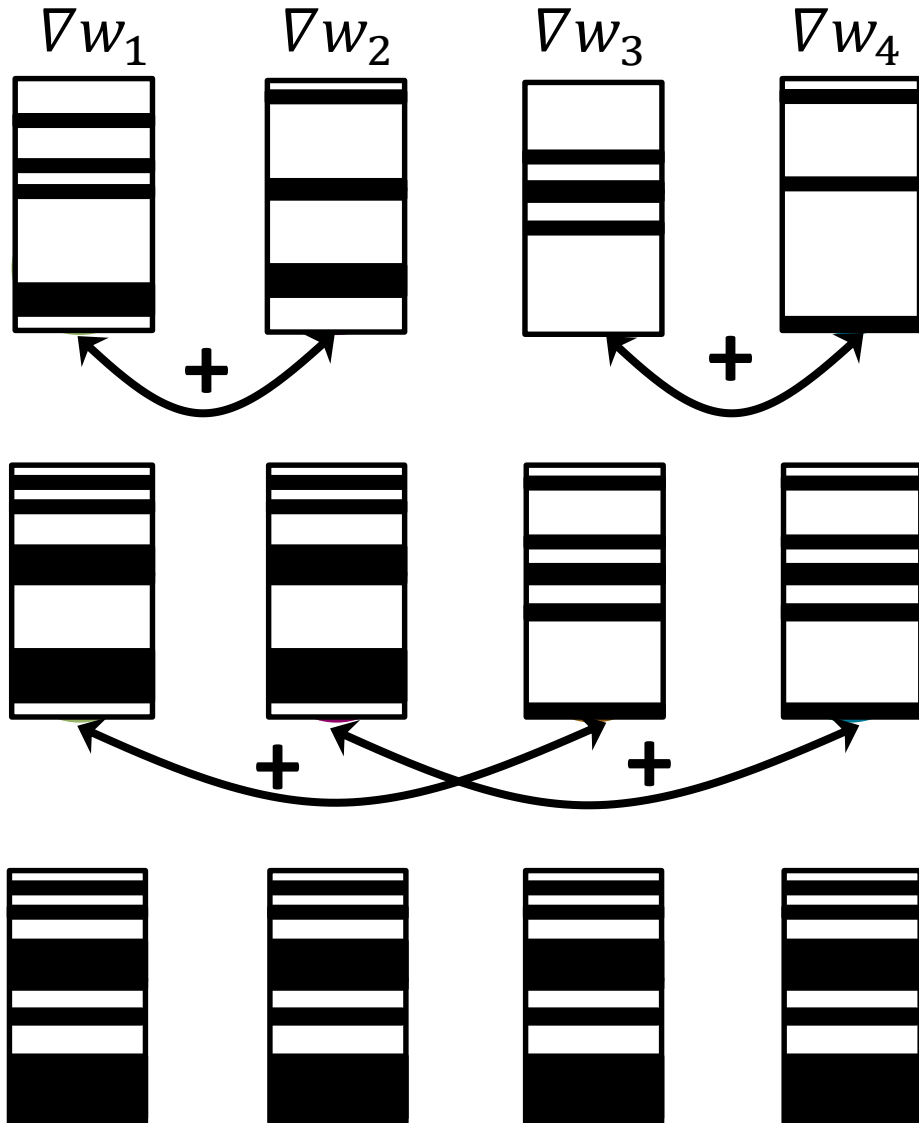
SparCML – Quantized sparse allreduce



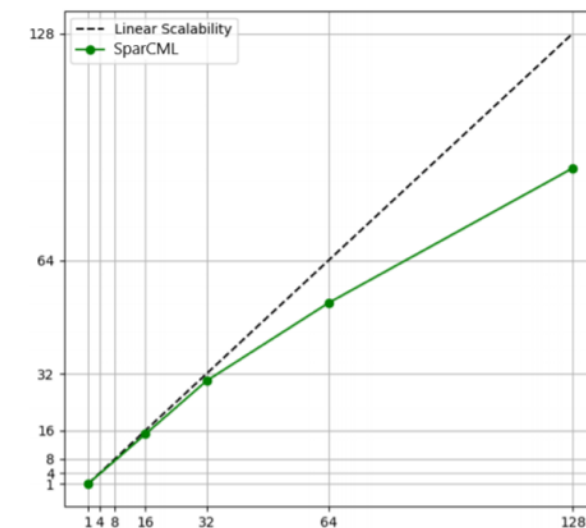
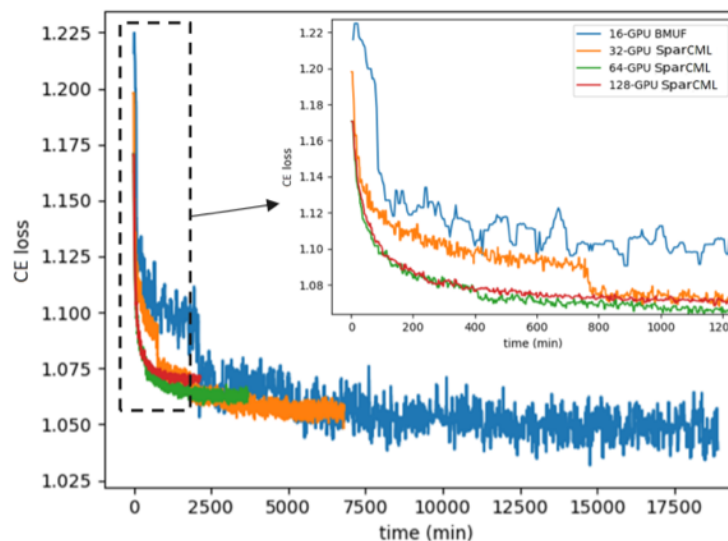
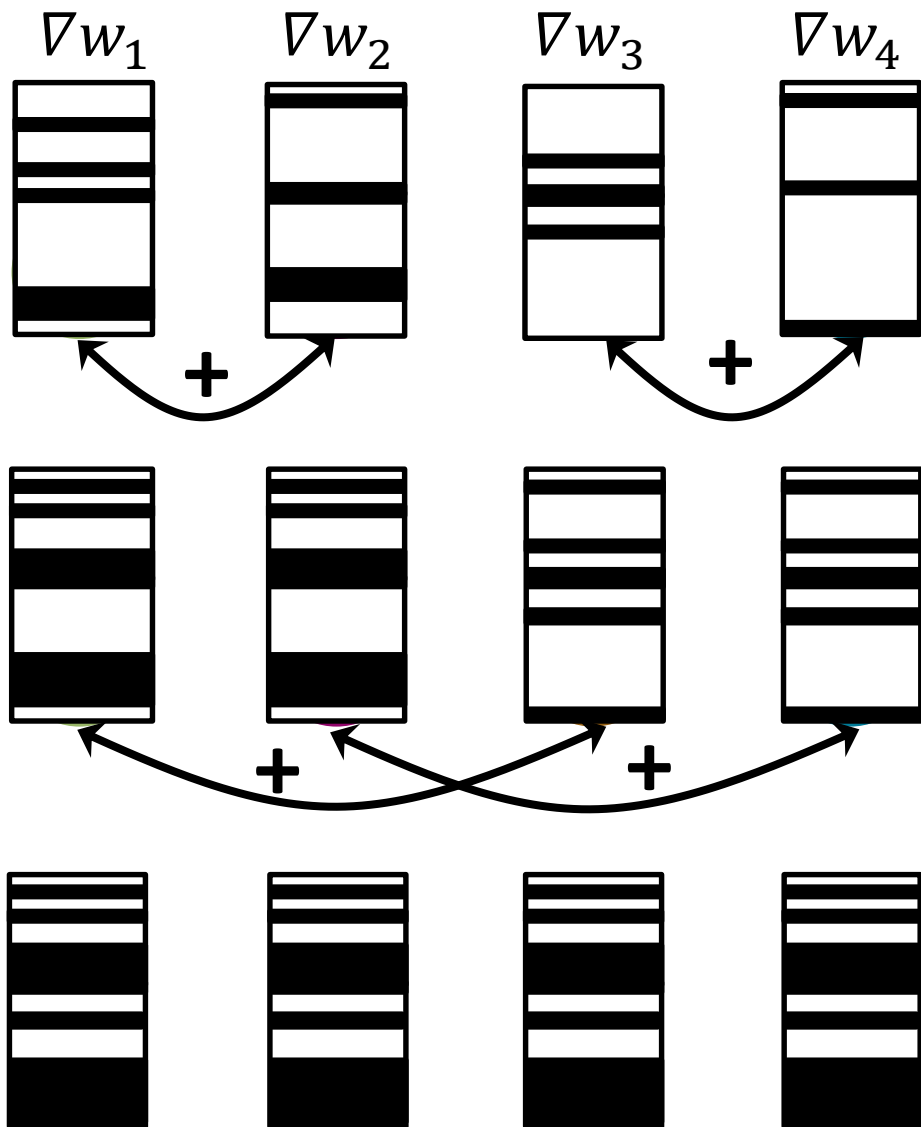
SparCML – Quantized sparse allreduce



SparCML – Quantized sparse allreduce

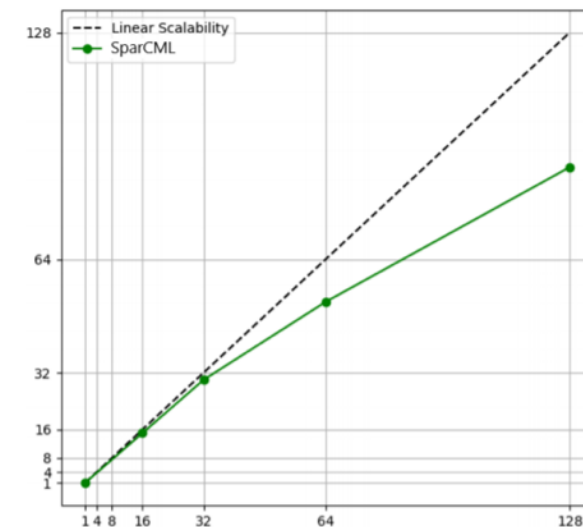
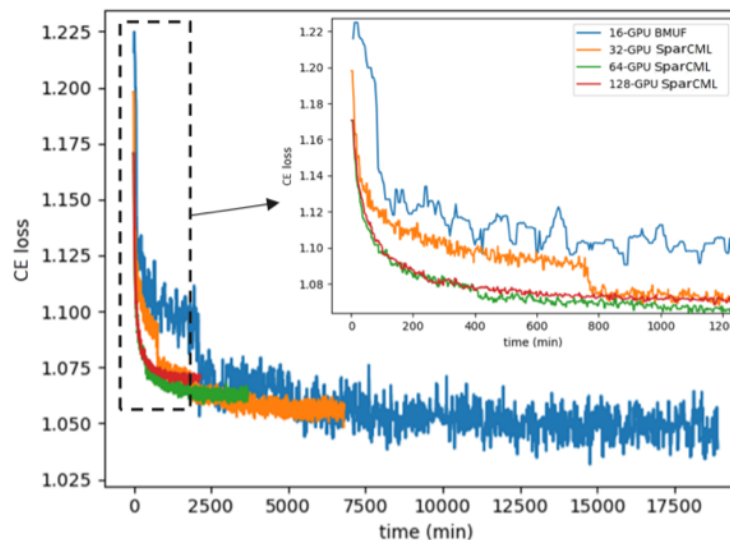
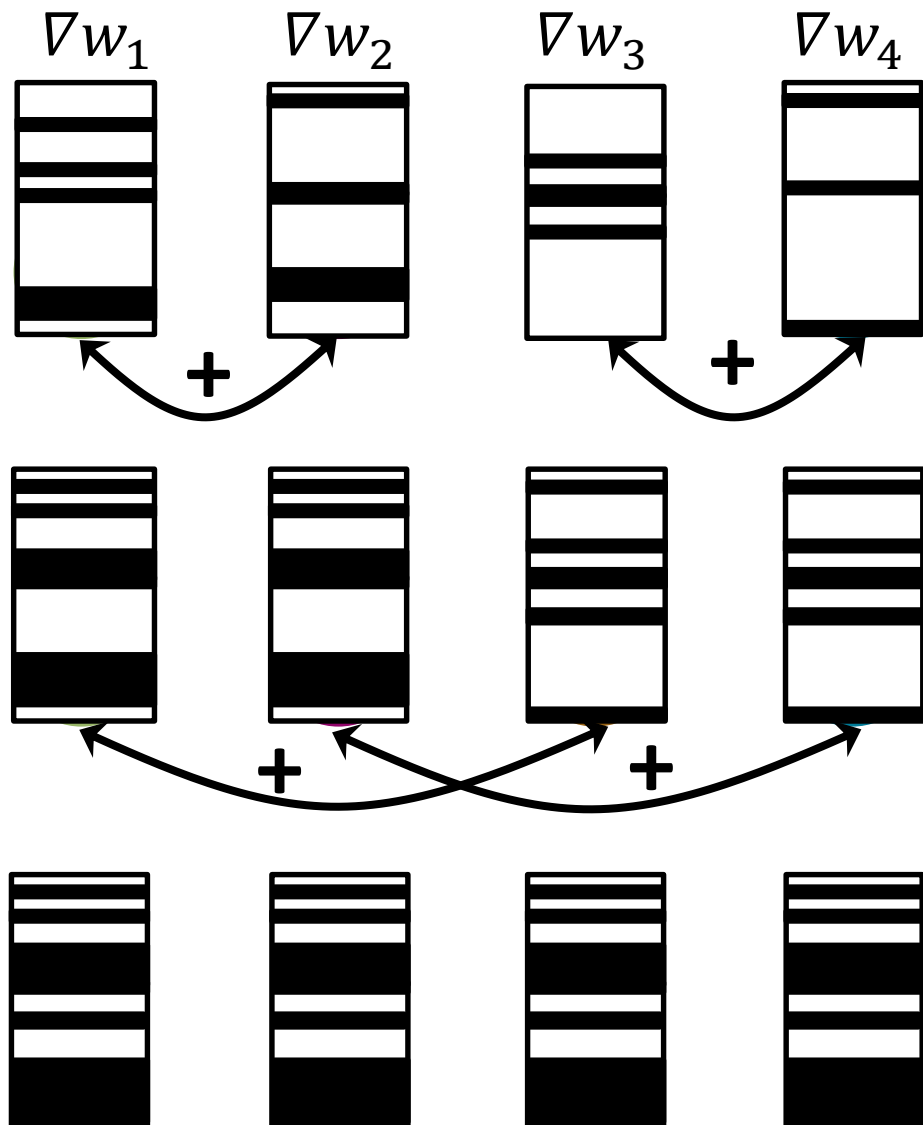


SparCML – Quantized sparse allreduce



System	Dataset	Model	# of nodes	Algorithm	Speedup
Piz Daint	ImageNet	VGG19	8	Q4	1.55 (3.31)
Piz Daint	ImageNet	AlexNet	16	Q4	1.30 (1.36)
Piz Daint EC2	MNIST	MLP	8	Top16_Q4 Top16_Q4	3.65 (4.53) 19.12 (22.97)

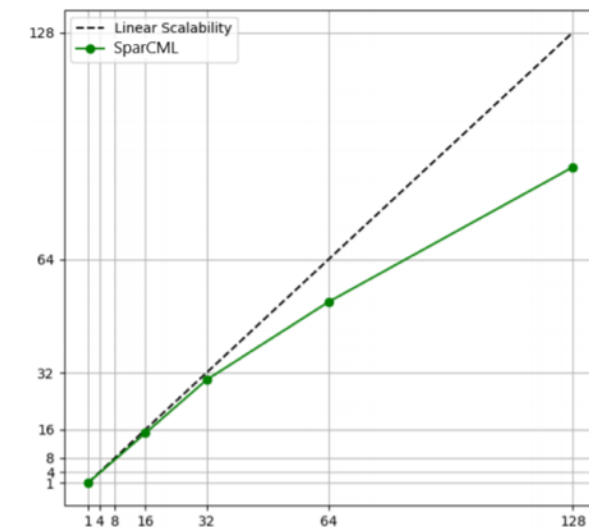
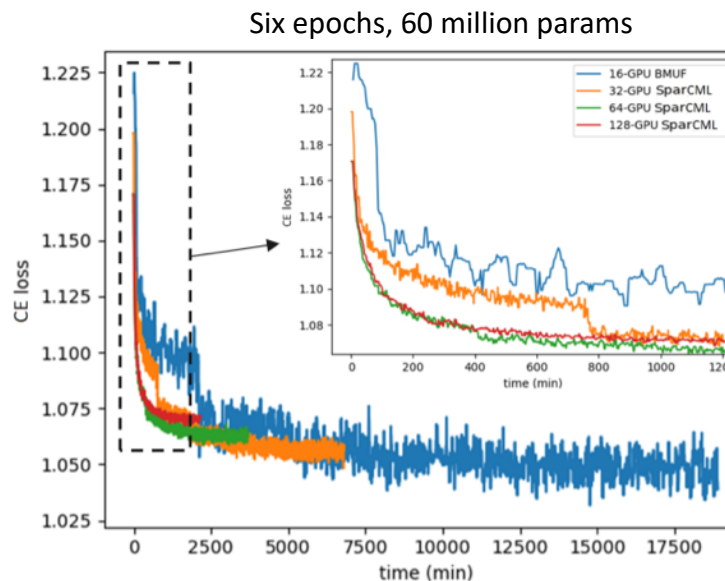
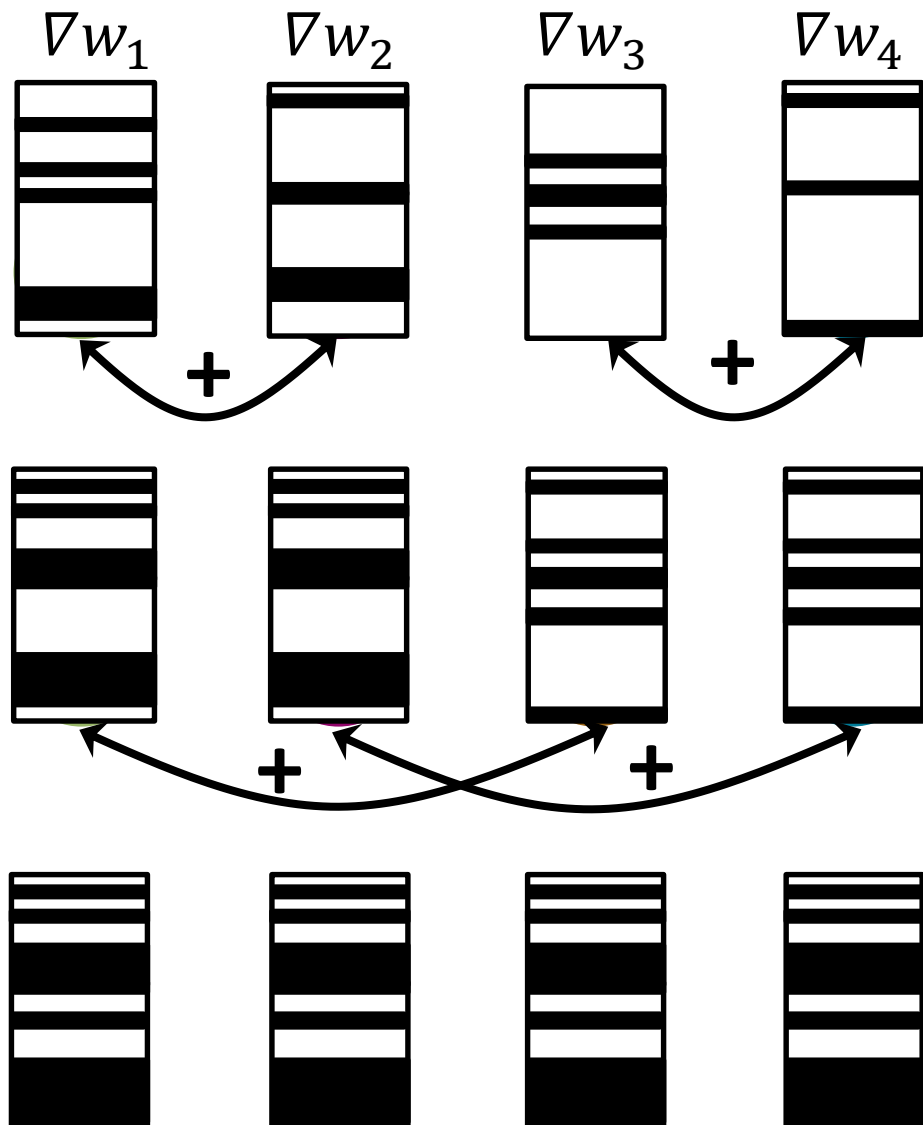
SparCML – Quantized sparse allreduce



Microsoft Speech Production Workload Results – **2 weeks** → **2 days!**

System	Dataset	Model	# of nodes	Algorithm	Speedup
Piz Daint	ImageNet	VGG19	8	Q4	1.55 (3.31)
Piz Daint	ImageNet	AlexNet	16	Q4	1.30 (1.36)
Piz Daint EC2	MNIST	MLP	8	Top16_Q4	3.65 (4.53)
				Top16_Q4	19.12 (22.97)

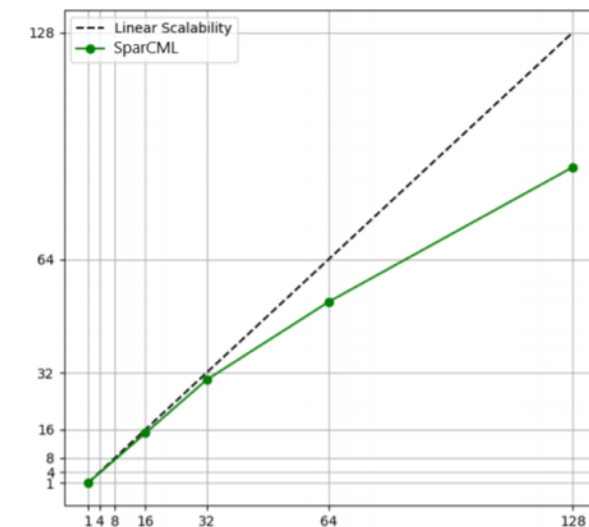
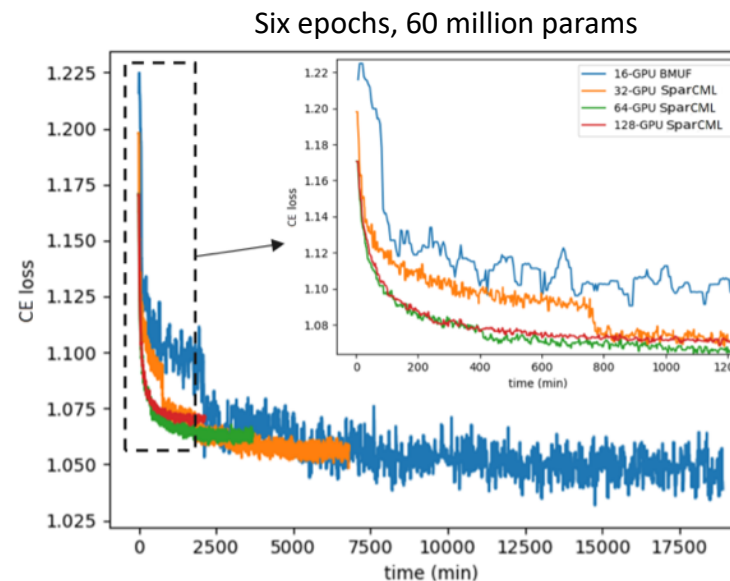
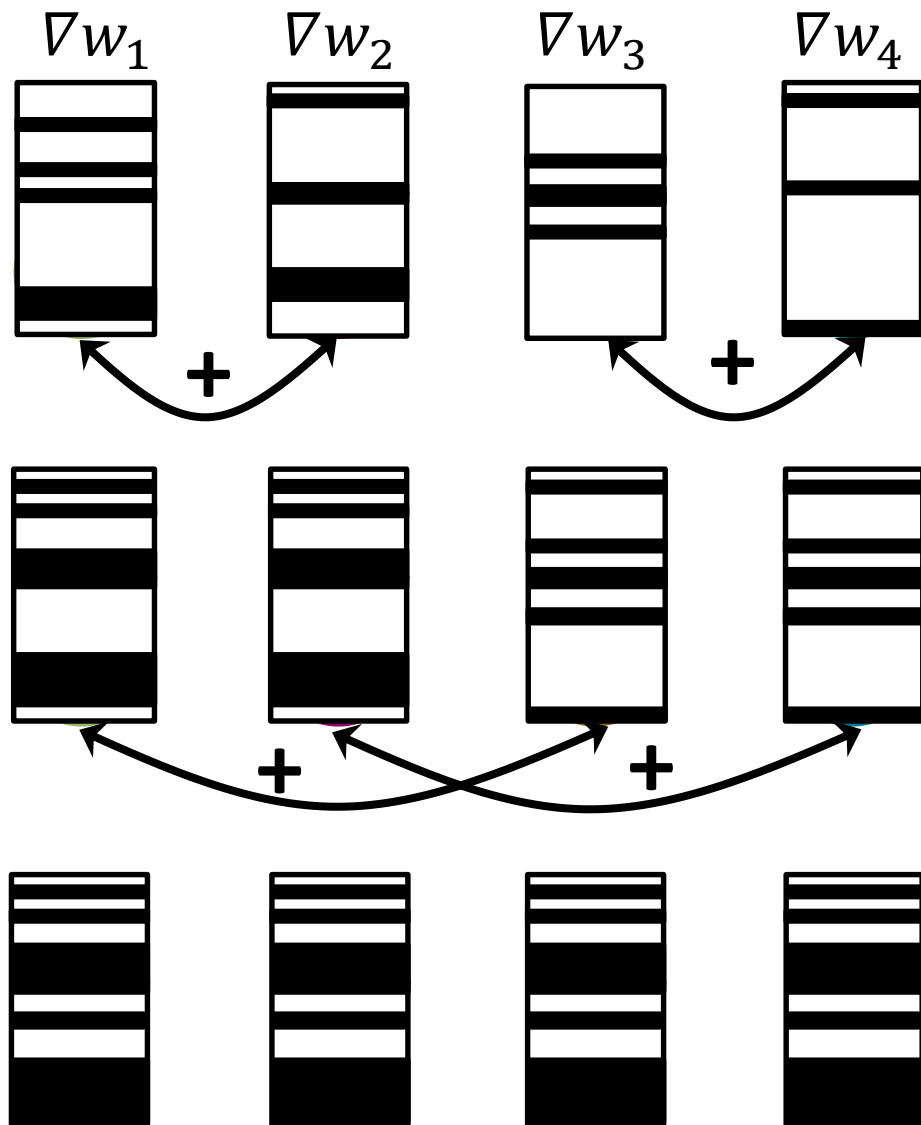
SparCML – Quantized sparse allreduce



Microsoft Speech Production Workload Results – **2 weeks** → **2 days!**

System	Dataset	Model	# of nodes	Algorithm	Speedup
Piz Daint	ImageNet	VGG19	8	Q4	1.55 (3.31)
Piz Daint	ImageNet	AlexNet	16	Q4	1.30 (1.36)
Piz Daint EC2	MNIST	MLP	8	Top16_Q4 Top16_Q4	3.65 (4.53) 19.12 (22.97)

SparCML – Quantized sparse allreduce



Microsoft Speech Production Workload Results – **2 weeks** → **2 days!**

System	Dataset	Model	# of nodes	Algorithm	Speedup
Piz Daint	ImageNet	VGG19	8	Q4	1.55 (3.31)
Piz Daint	ImageNet	AlexNet	16	Q4	1.30 (1.36)
Piz Daint EC2	MNIST	MLP	8	Top16_Q4 Top16_Q4	3.65 (4.53) 19.12 (22.97)

The limits to data parallelism

The limits to data parallelism

- When does data parallelism break down?
- Communication overheads
- Hyperparameter tuning
- Memory for one sample
- More GPUs than samples in a mini-batch

The limits to data parallelism

- When does data parallelism break down?
- Communication overheads
- Hyperparameter tuning
- Memory for one sample
- More GPUs than samples in a mini-batch

The limits to data parallelism

- When does data parallelism break down?
- Communication overheads
- Hyperparameter tuning
- Memory for one sample
- More GPUs than samples in a mini-batch

Need to strong scale!

Distributed-memory fully-connected layers

$$Y = \sigma(WX + b)$$

Distributed-memory fully-connected layers

$$Y = \sigma(WX + b)$$

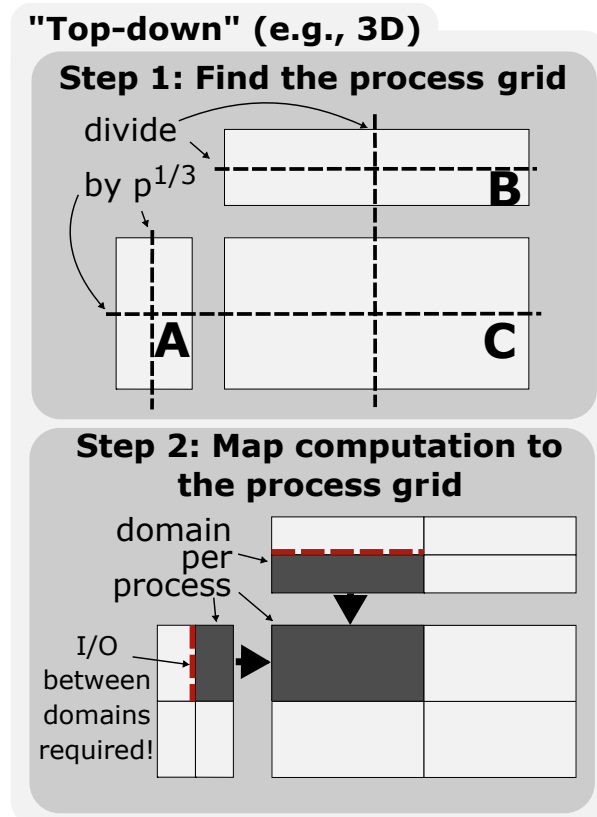
- Just a distributed matrix-matrix multiplication!

Distributed-memory fully-connected layers

$$Y = \sigma(WX + b)$$

- Just a distributed matrix-matrix multiplication!

Use SUMMA [Van Essen et al. 2015]

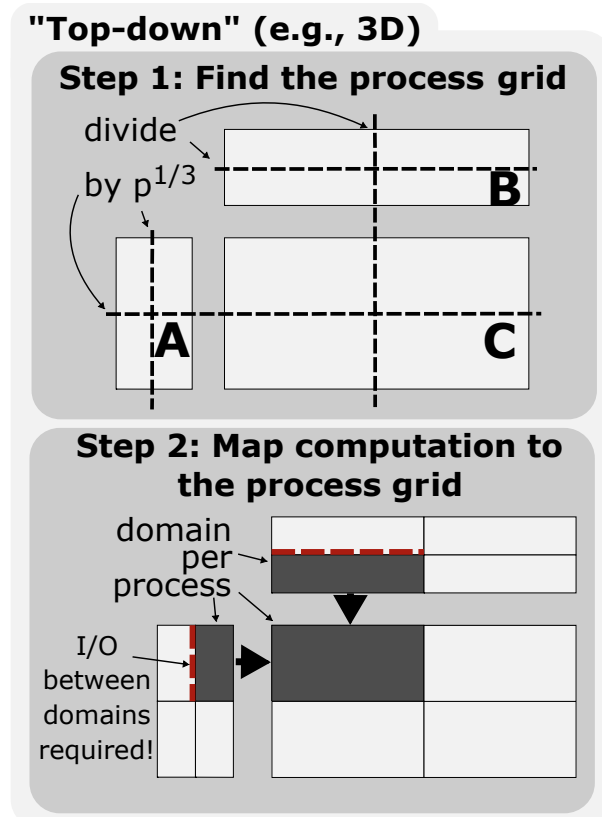


Distributed-memory fully-connected layers

$$Y = \sigma(WX + b)$$

- Just a distributed matrix-matrix multiplication!

Use SUMMA [Van Essen et al. 2015]

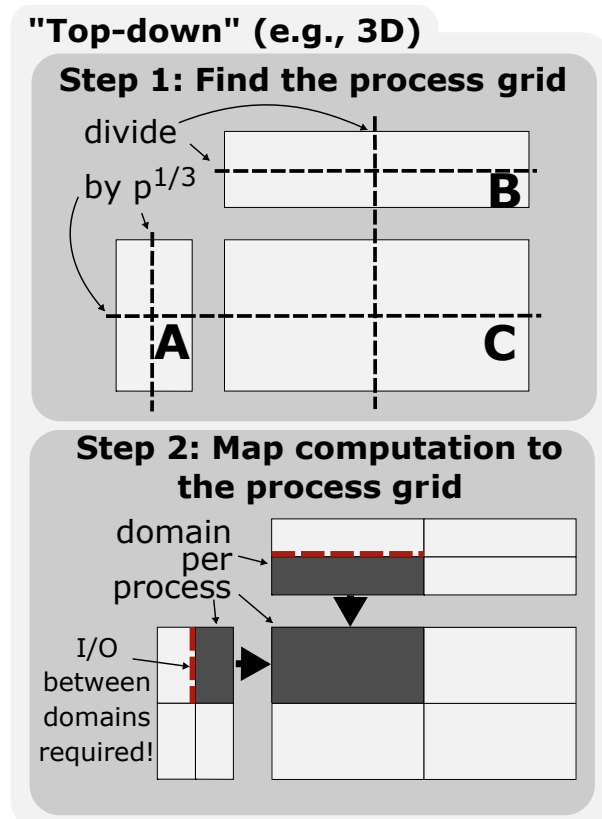


Distributed-memory fully-connected layers

$$Y = \sigma(WX + b)$$

- Just a distributed matrix-matrix multiplication!

Use SUMMA [Van Essen et al. 2015]



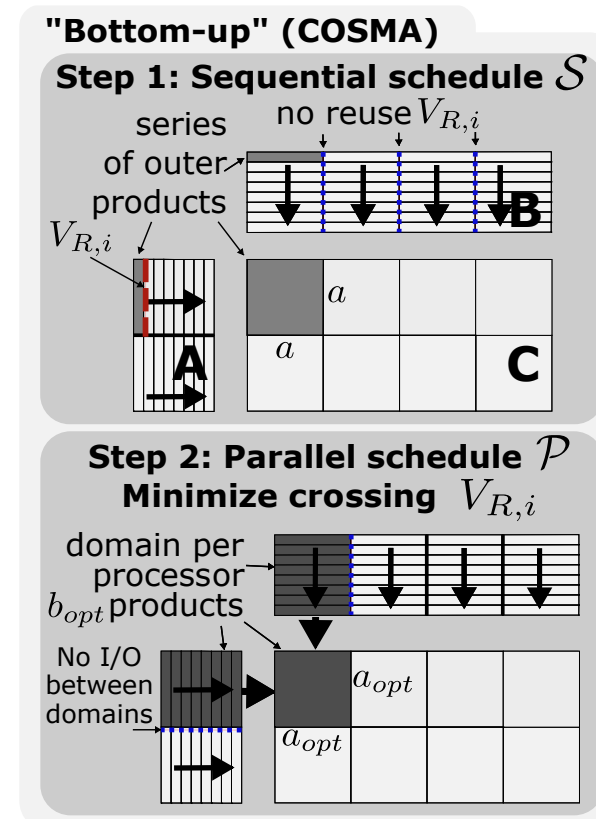
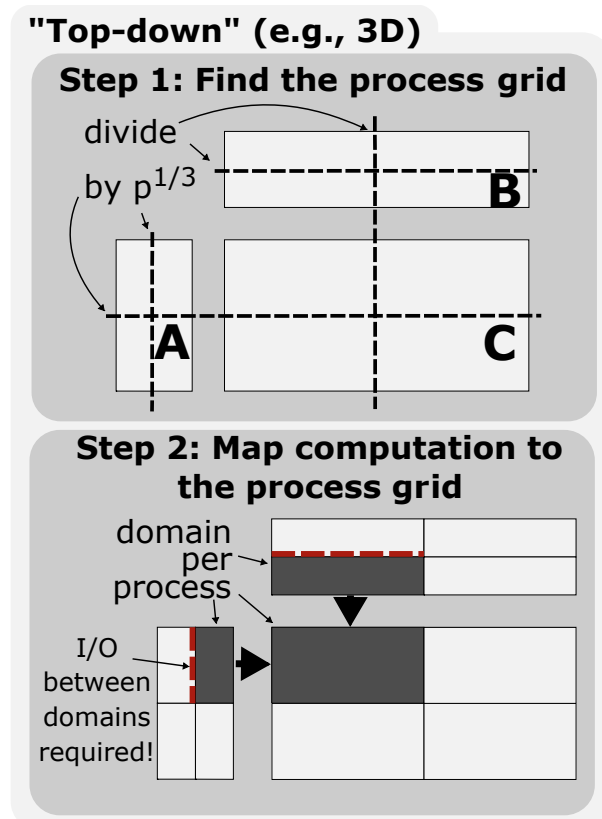
Distributed-memory fully-connected layers

$$Y = \sigma(WX + b)$$

- Just a distributed matrix-matrix multiplication!

Use SUMMA [Van Essen et al. 2015]

New: COSMA [Kwasniewski et al. 2019]



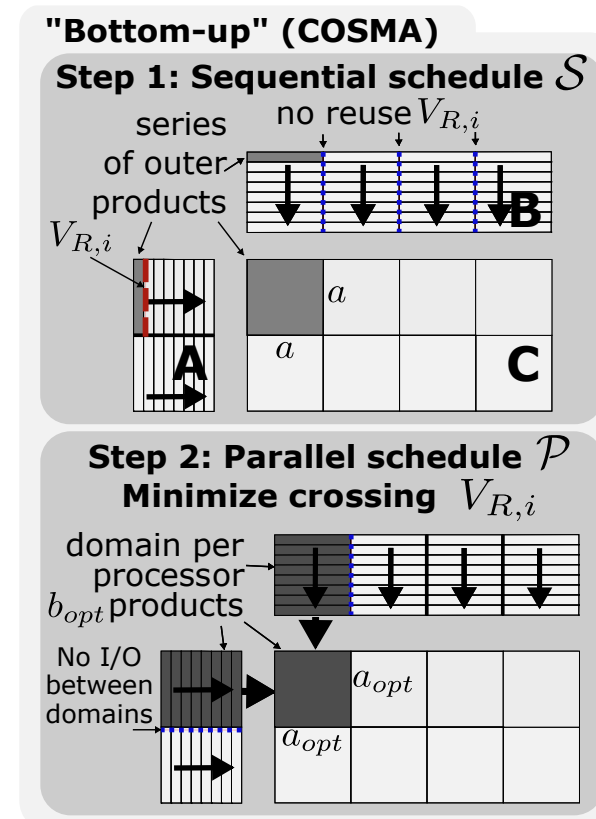
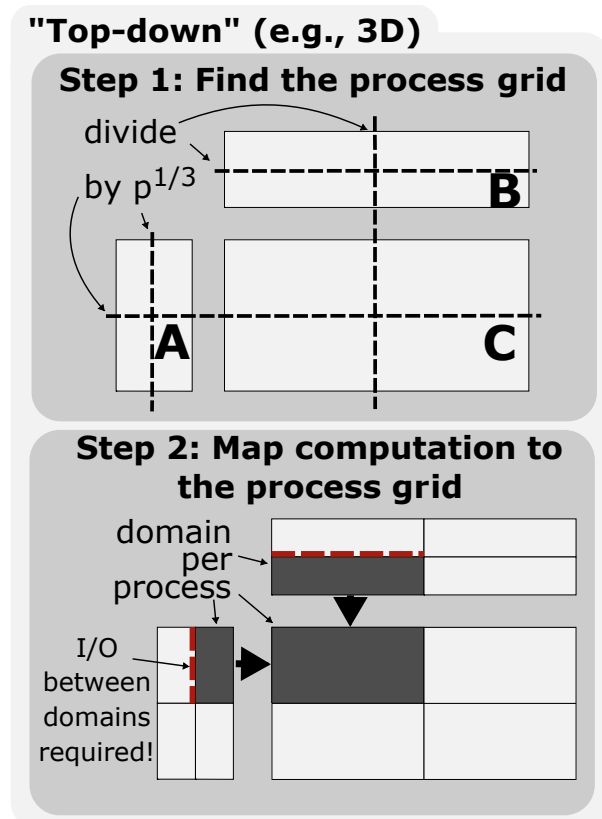
Distributed-memory fully-connected layers

$$Y = \sigma(WX + b)$$

- Just a distributed matrix-matrix multiplication!

Use SUMMA [Van Essen et al. 2015]

New: COSMA [Kwasniewski et al. 2019]



Spatial parallelism [Dryden et al. 2019]

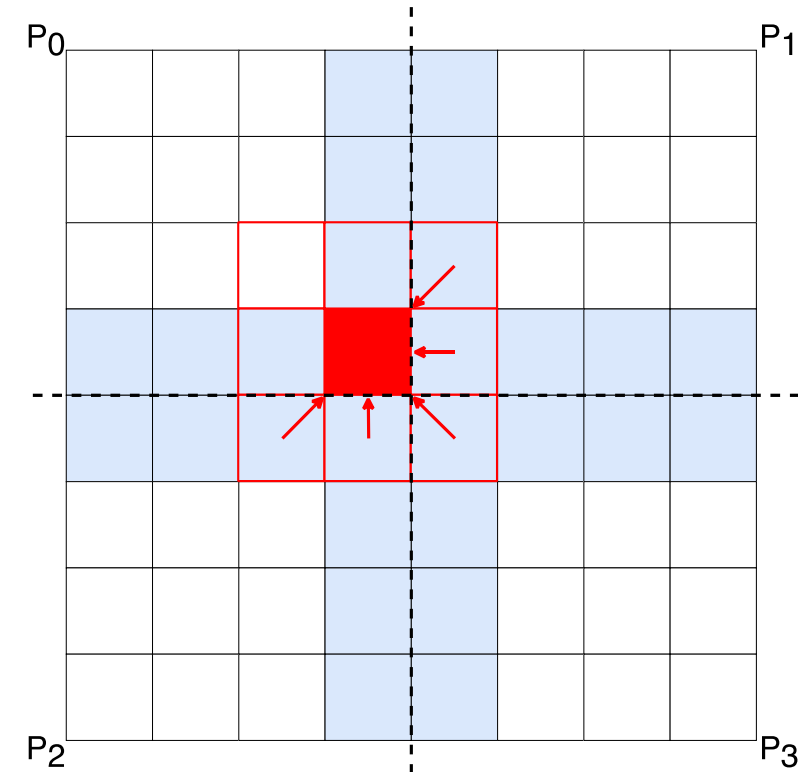
Spatial parallelism [Dryden et al. 2019]

- **Observation: Convolution is just a funny stencil operation**
- **Domain decomposition with a halo exchange!**

Spatial parallelism [Dryden et al. 2019]

- **Observation: Convolution is just a funny stencil operation**
- **Domain decomposition with a halo exchange!**

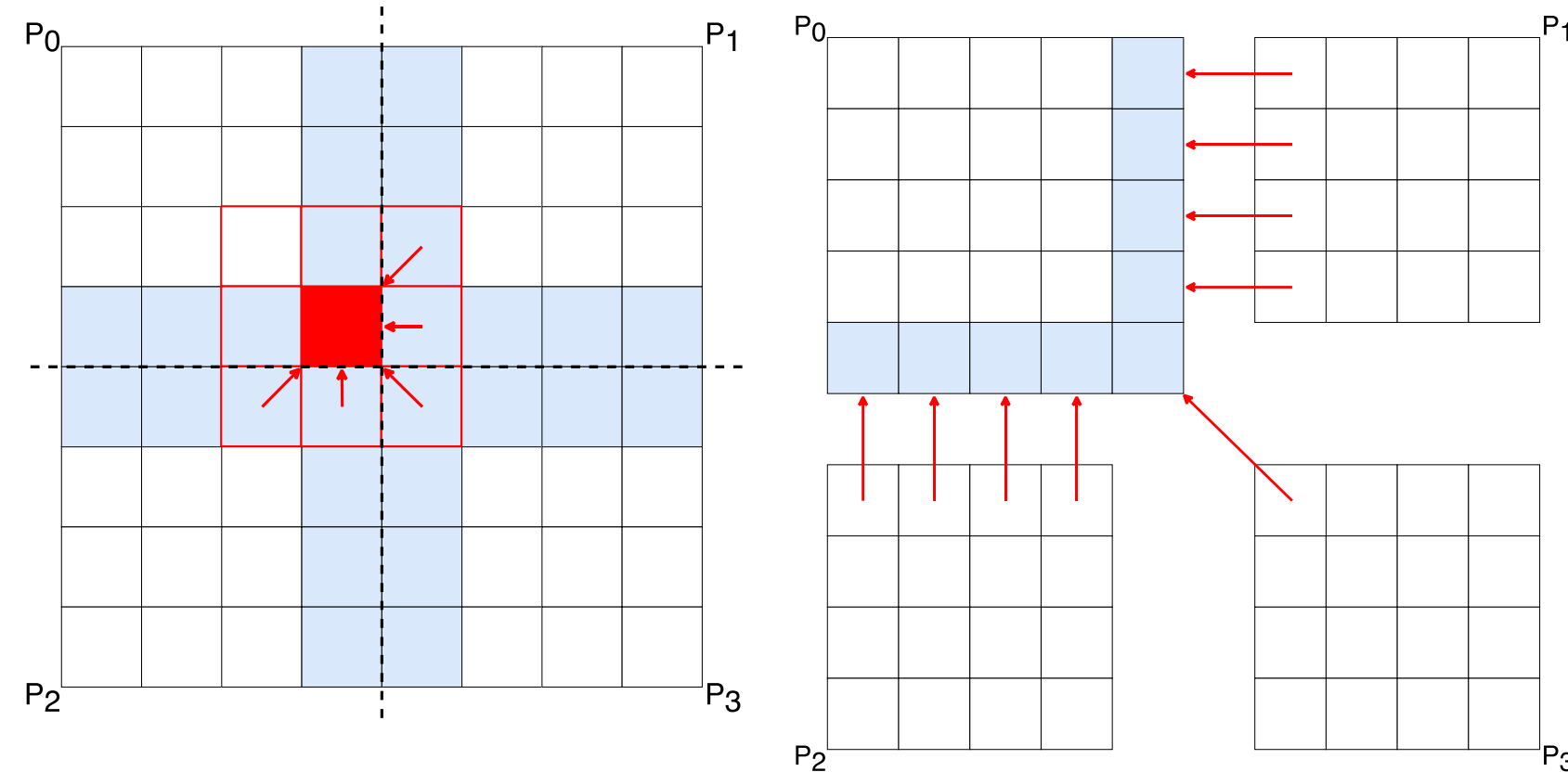
Convolution dependencies



Spatial parallelism [Dryden et al. 2019]

- **Observation: Convolution is just a funny stencil operation**
- **Domain decomposition with a halo exchange!**

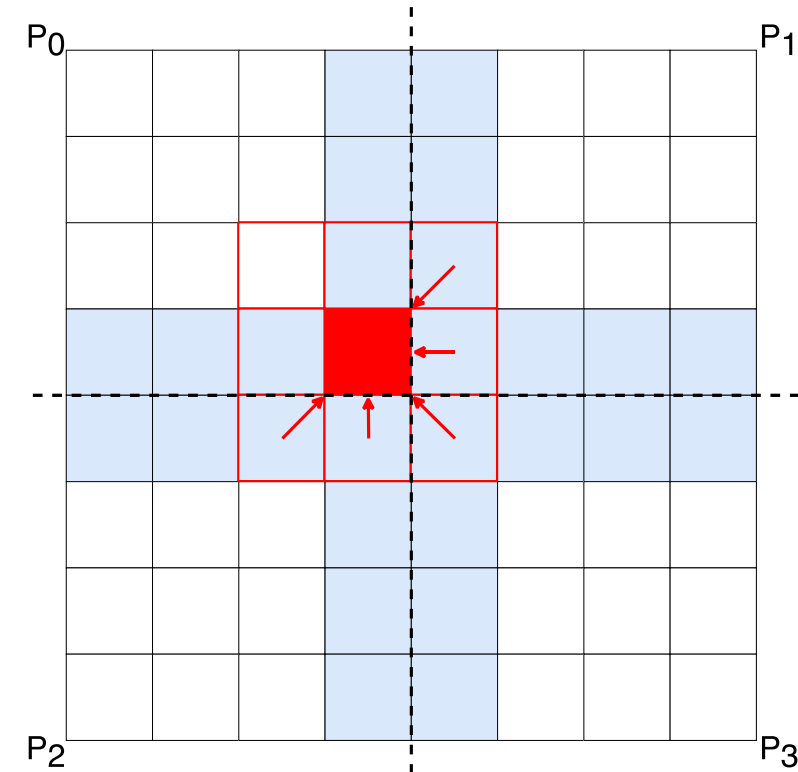
Convolution dependencies



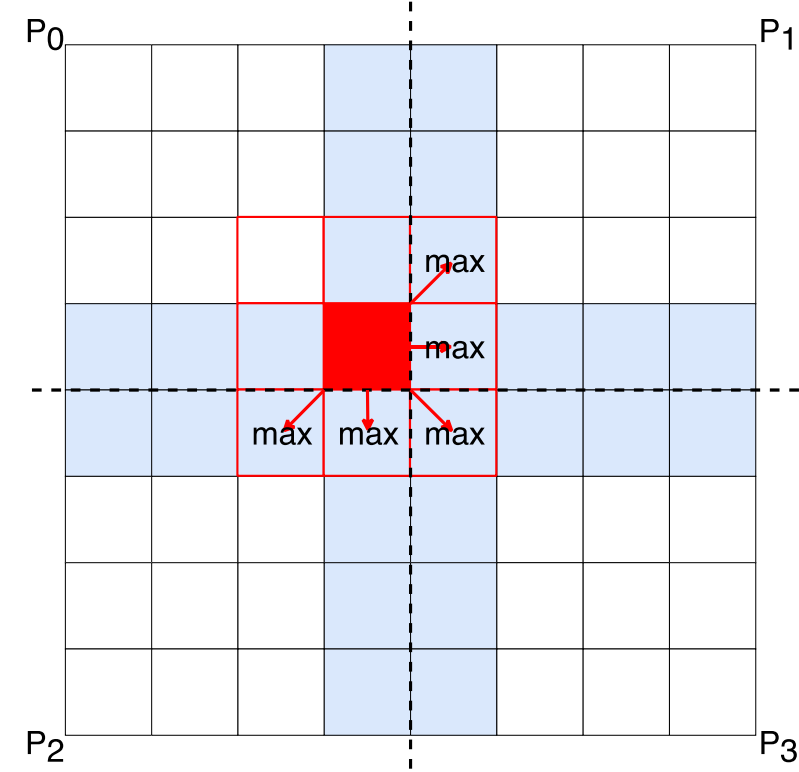
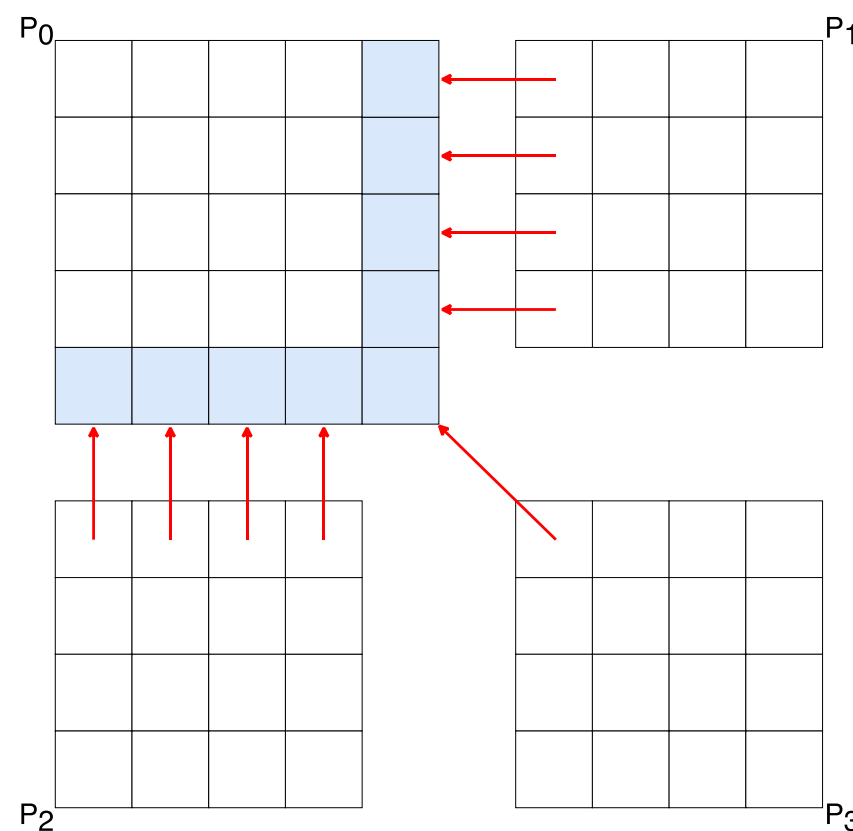
Spatial parallelism [Dryden et al. 2019]

- **Observation: Convolution is just a funny stencil operation**
- **Domain decomposition with a halo exchange!**

Convolution dependencies



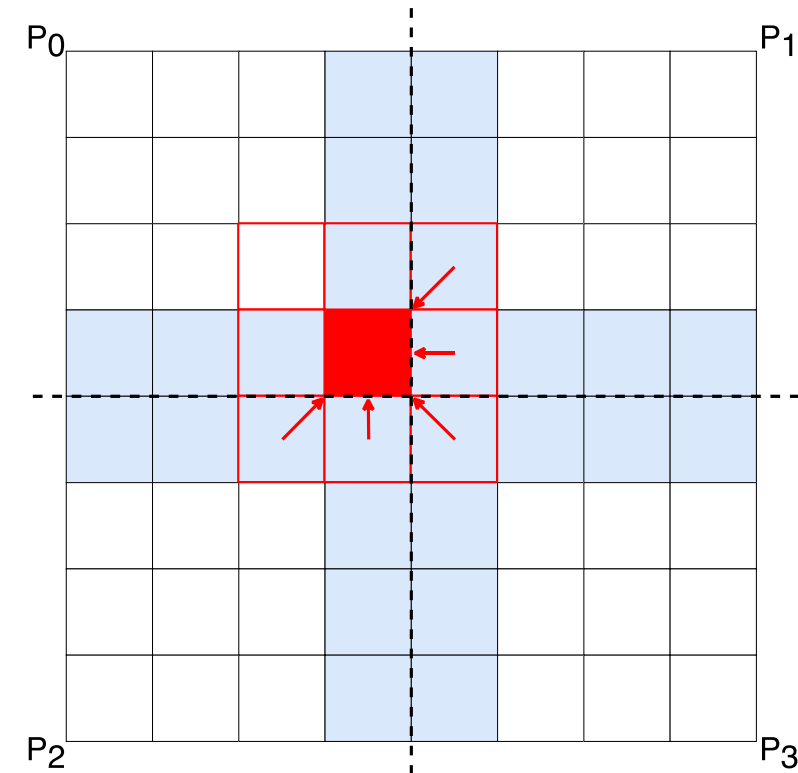
Halo exchange



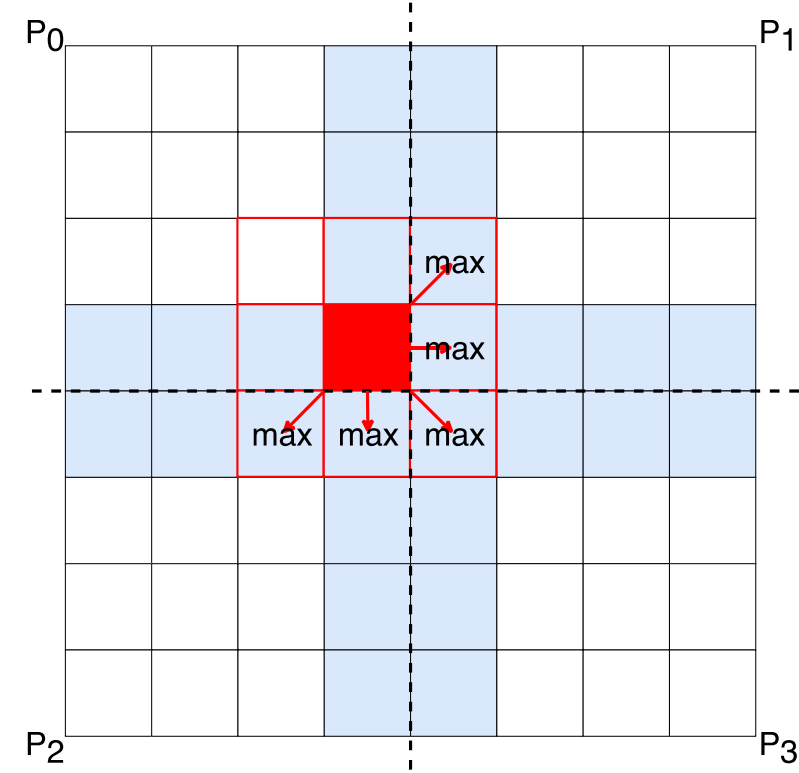
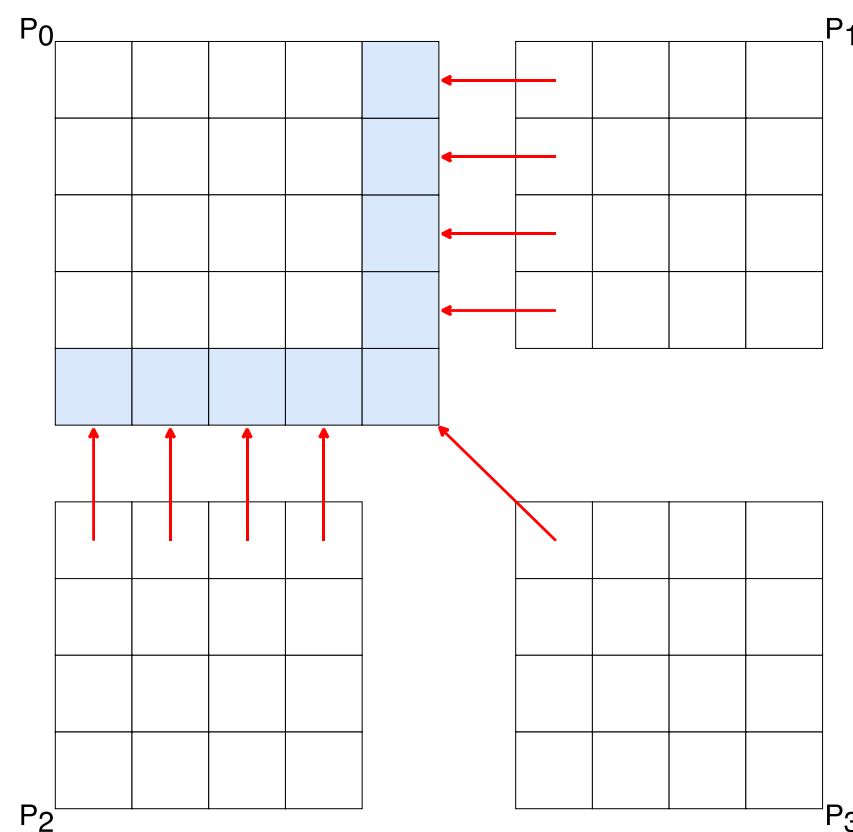
Spatial parallelism [Dryden et al. 2019]

- **Observation: Convolution is just a funny stencil operation**
- **Domain decomposition with a halo exchange!**

Convolution dependencies



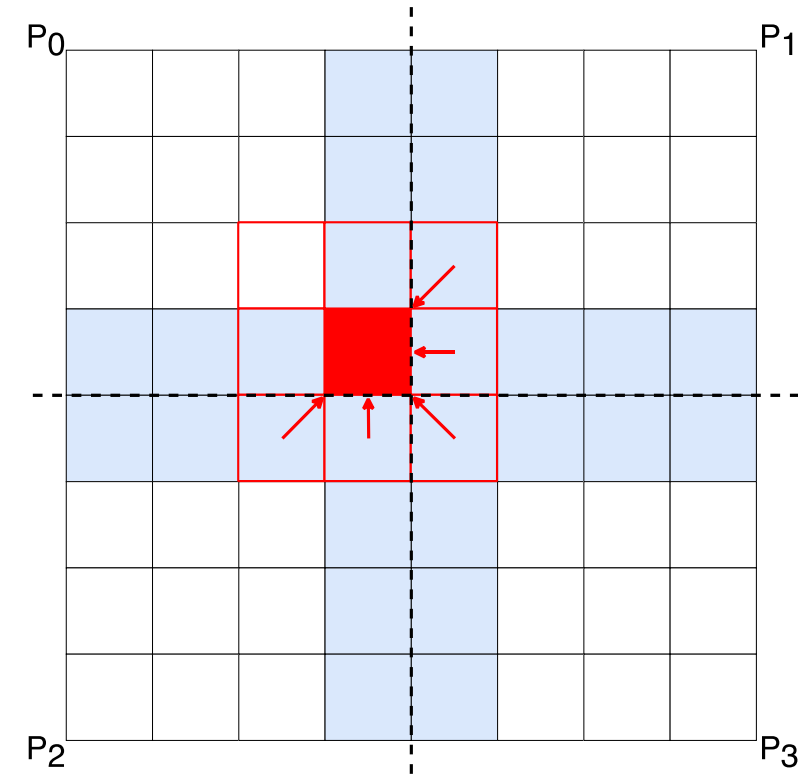
Halo exchange



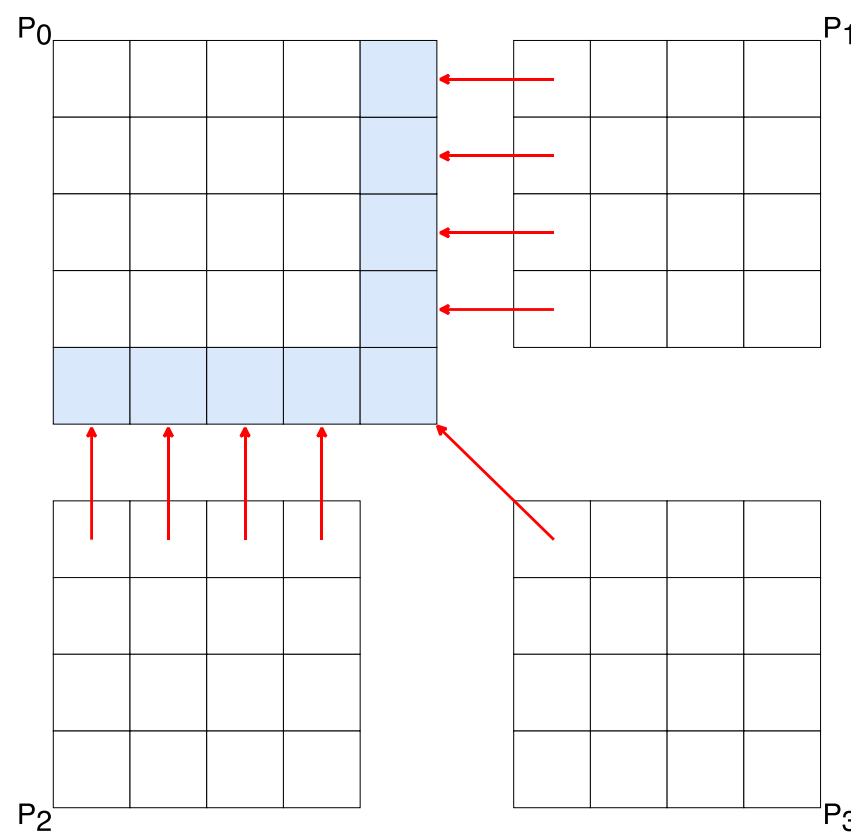
Spatial parallelism [Dryden et al. 2019]

- **Observation: Convolution is just a funny stencil operation**
- **Domain decomposition with a halo exchange!**

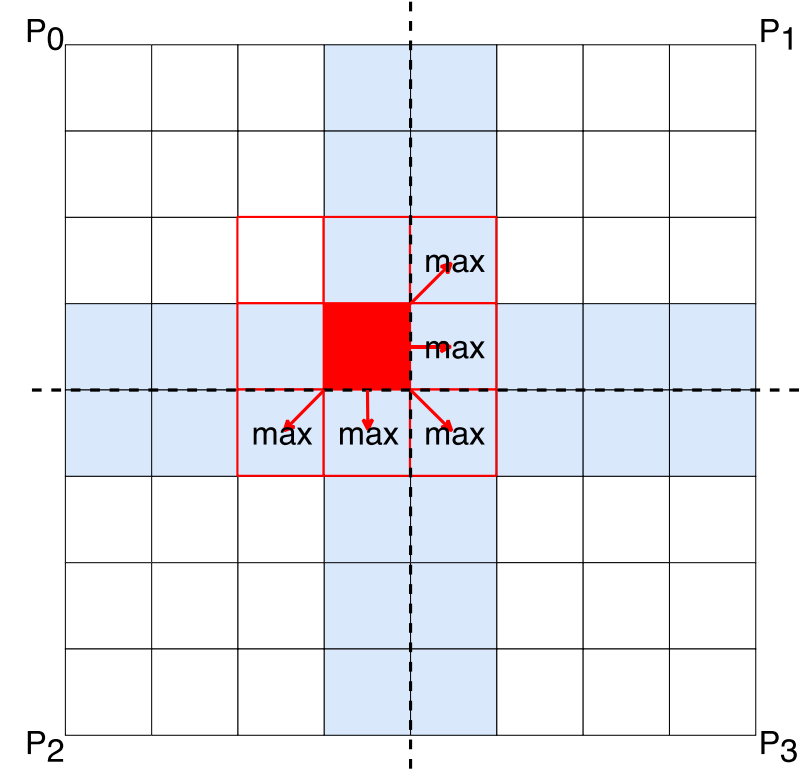
Convolution dependencies



Halo exchange



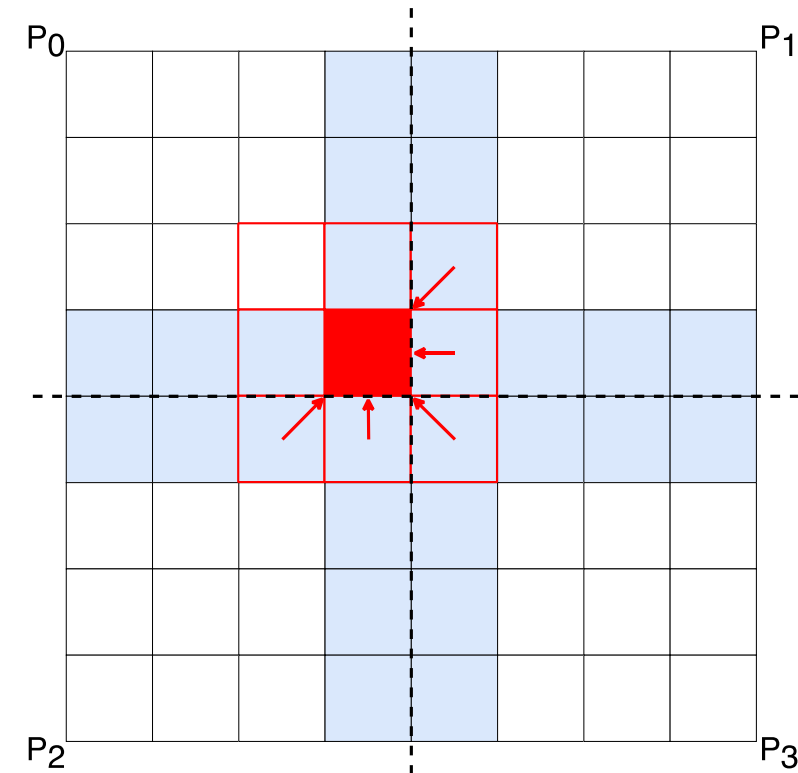
“Push” exchange for pooling



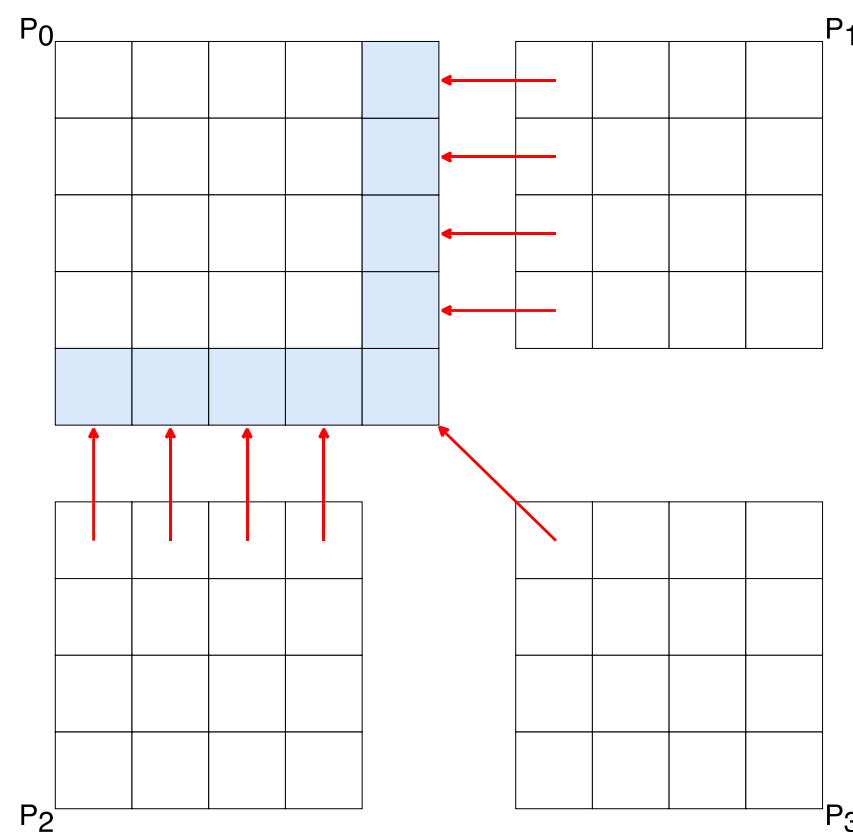
Spatial parallelism [Dryden et al. 2019]

- **Observation: Convolution is just a funny stencil operation**
- **Domain decomposition with a halo exchange!**

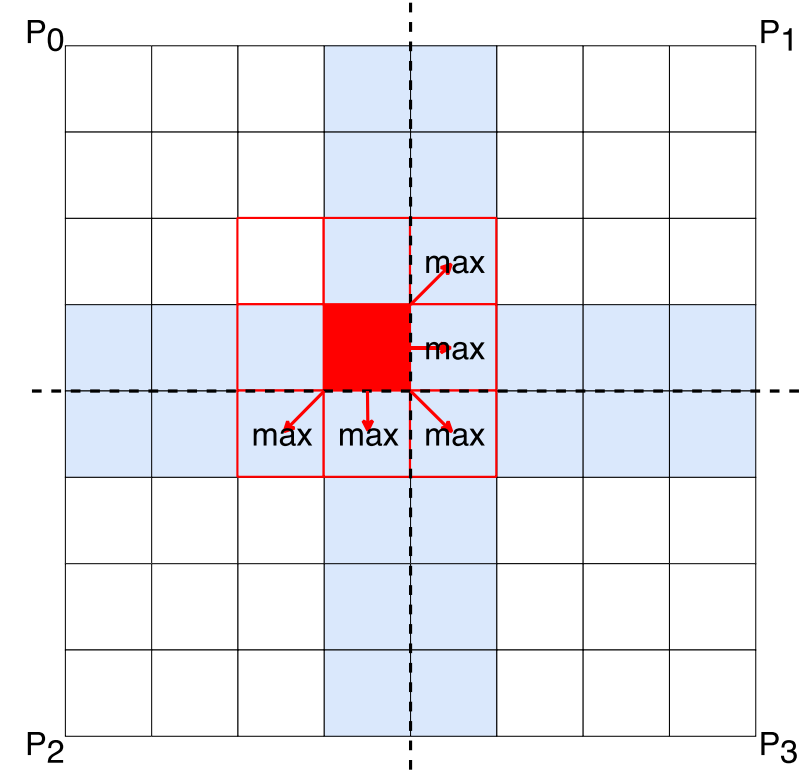
Convolution dependencies



Halo exchange

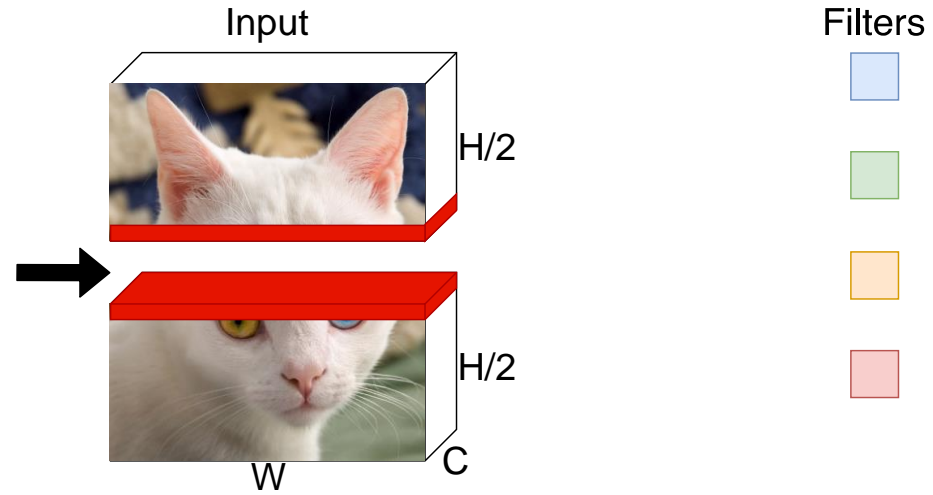


“Push” exchange for pooling

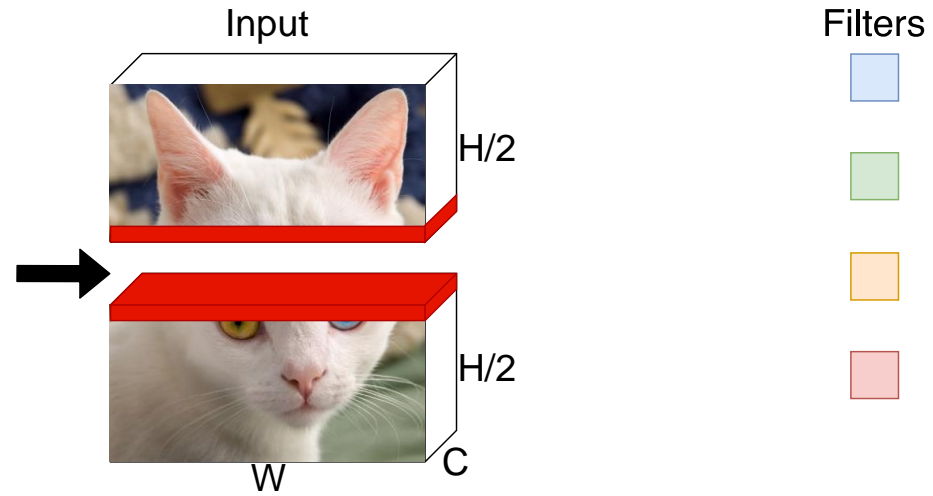


Spatial parallelism [Dryden et al. 2019]

Spatial parallelism [Dryden et al. 2019]



Spatial parallelism [Dryden et al. 2019]



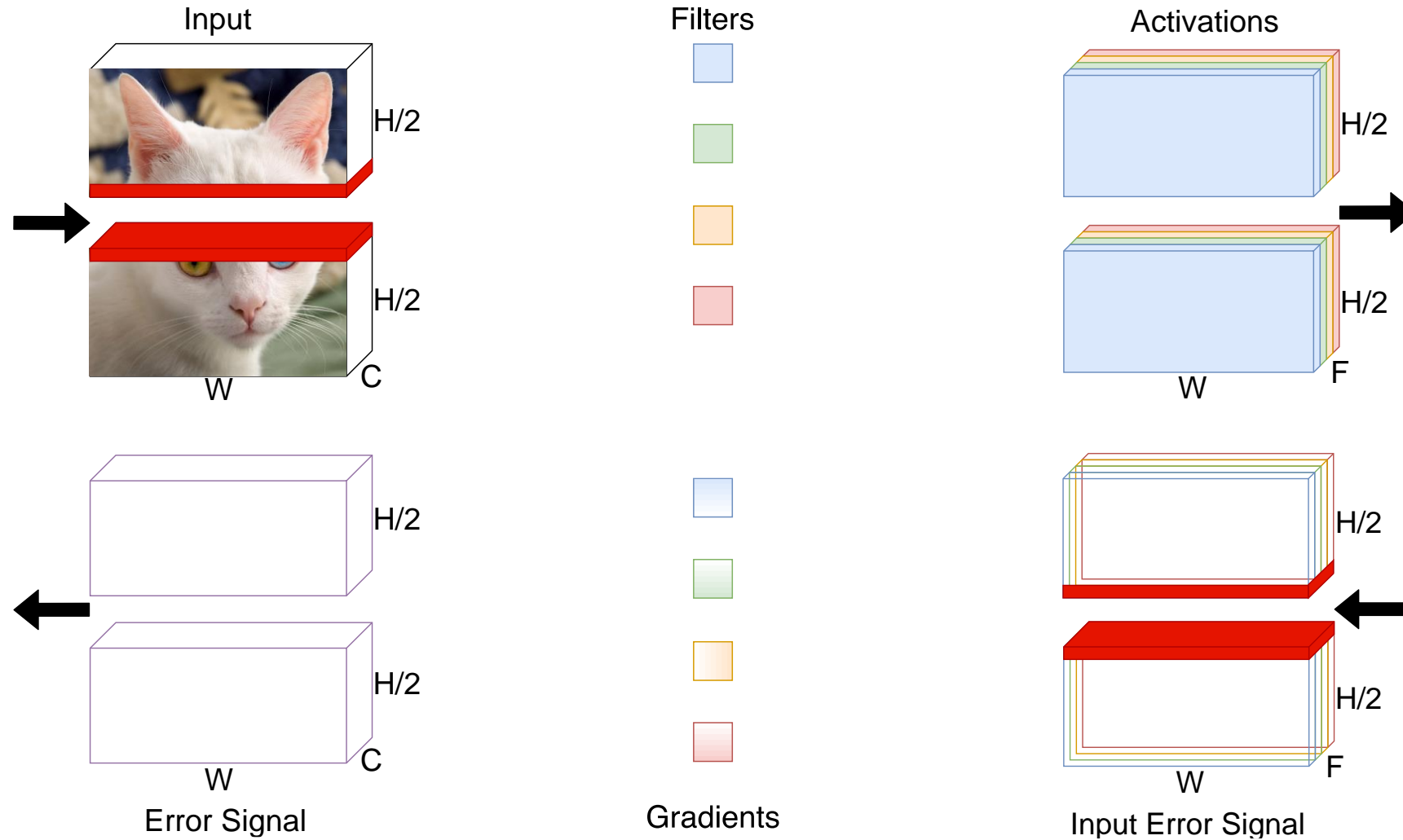
Spatial parallelism [Dryden et al. 2019]



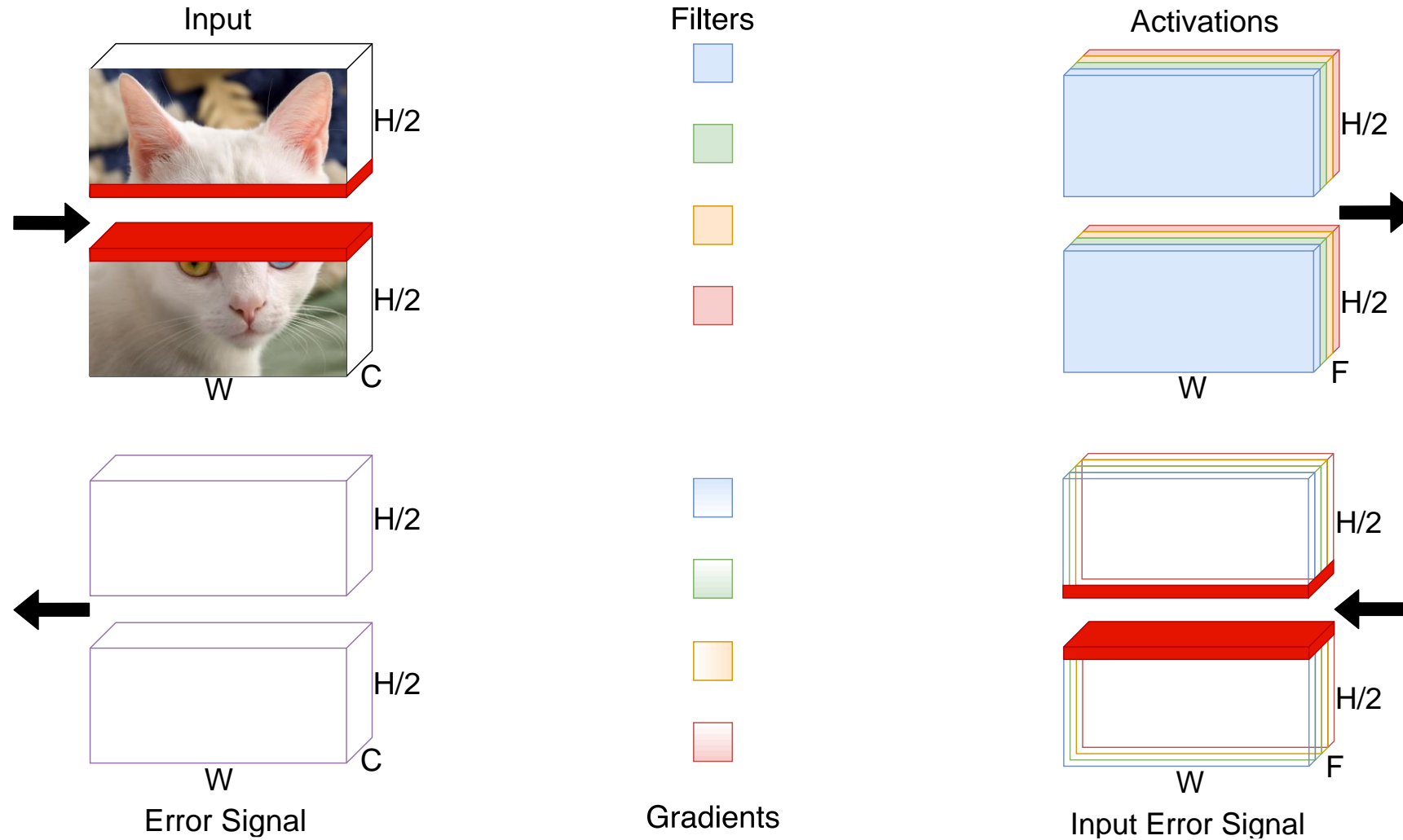
Spatial parallelism [Dryden et al. 2019]



Spatial parallelism [Dryden et al. 2019]



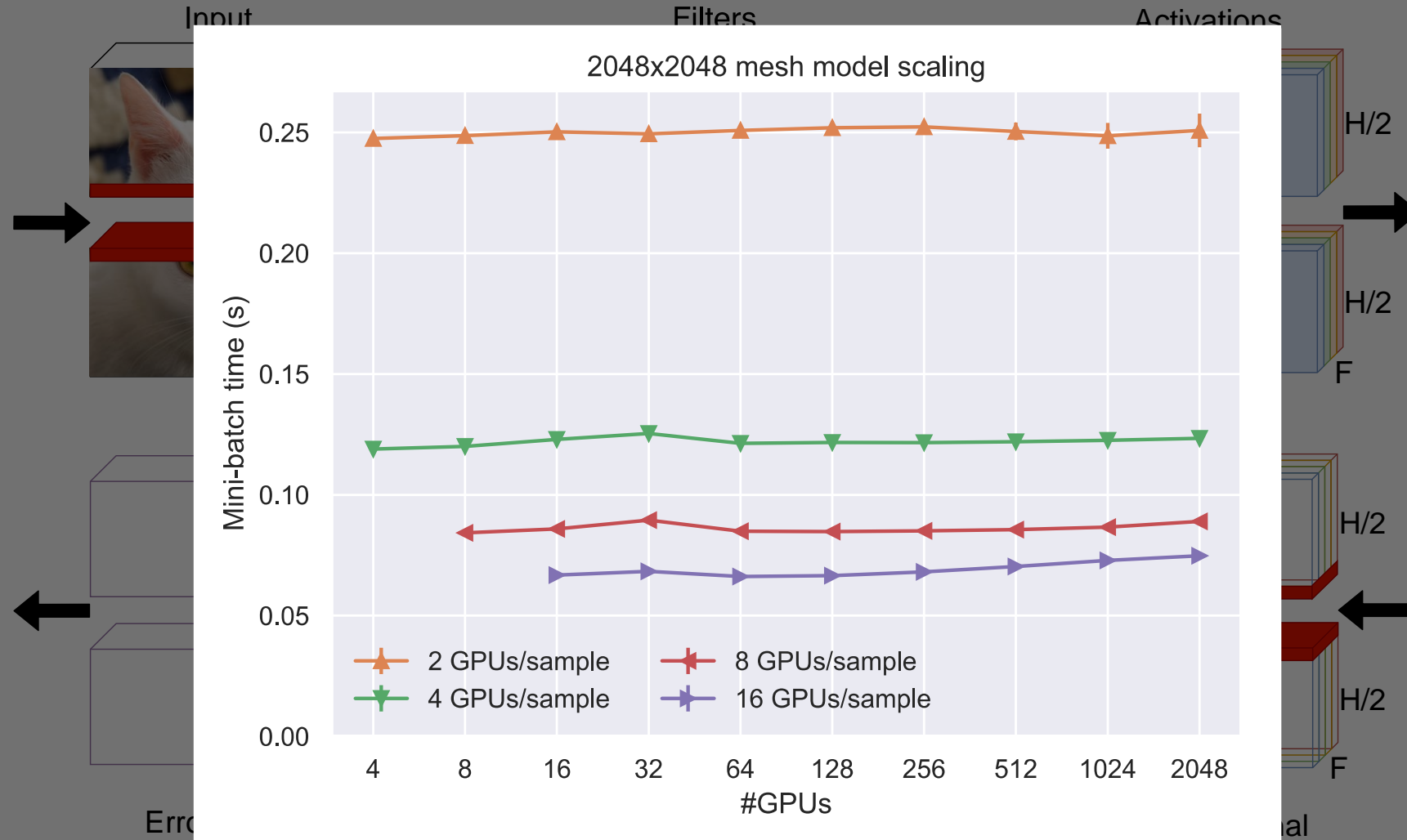
Spatial parallelism [Dryden et al. 2019]



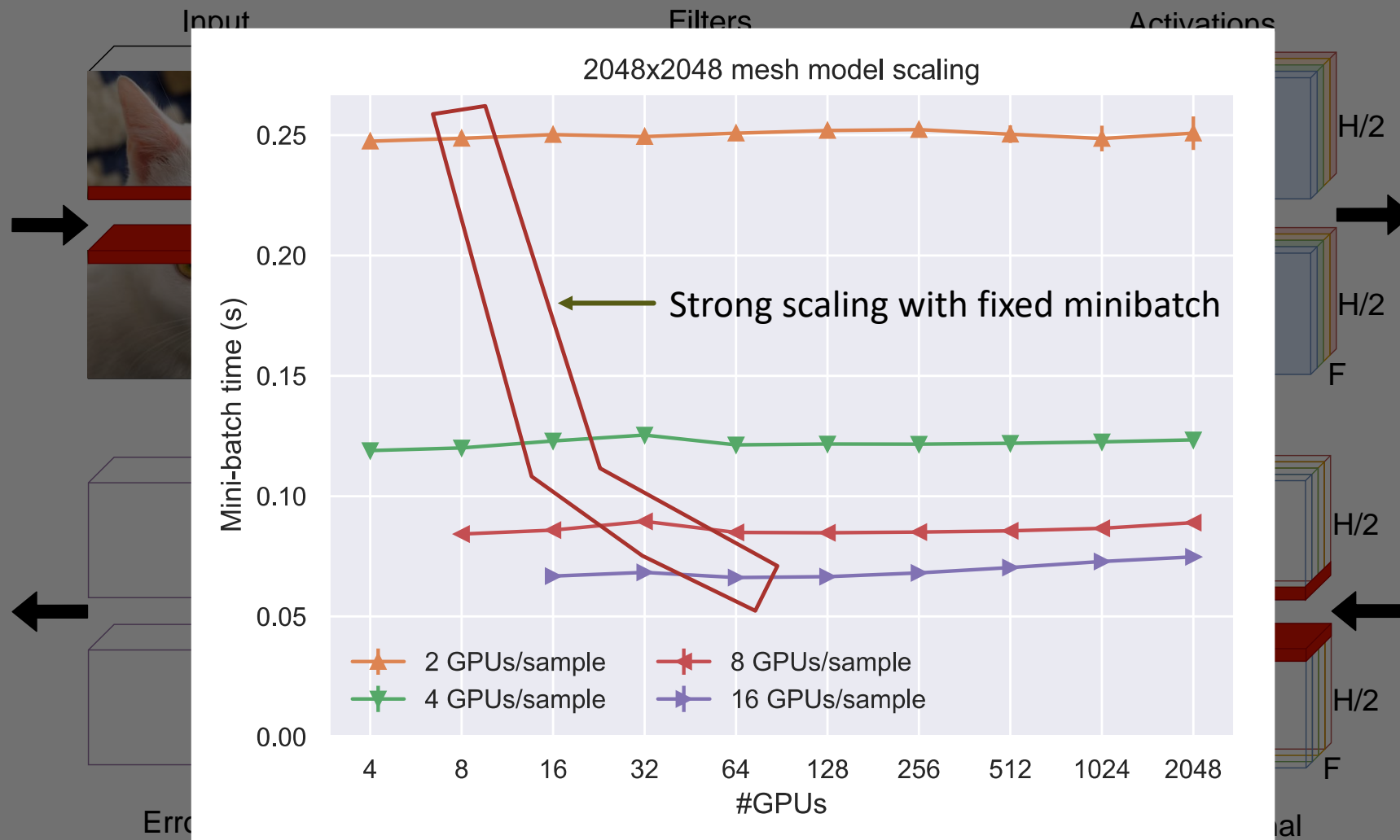
Spatial parallelism [Dryden et al. 2019]



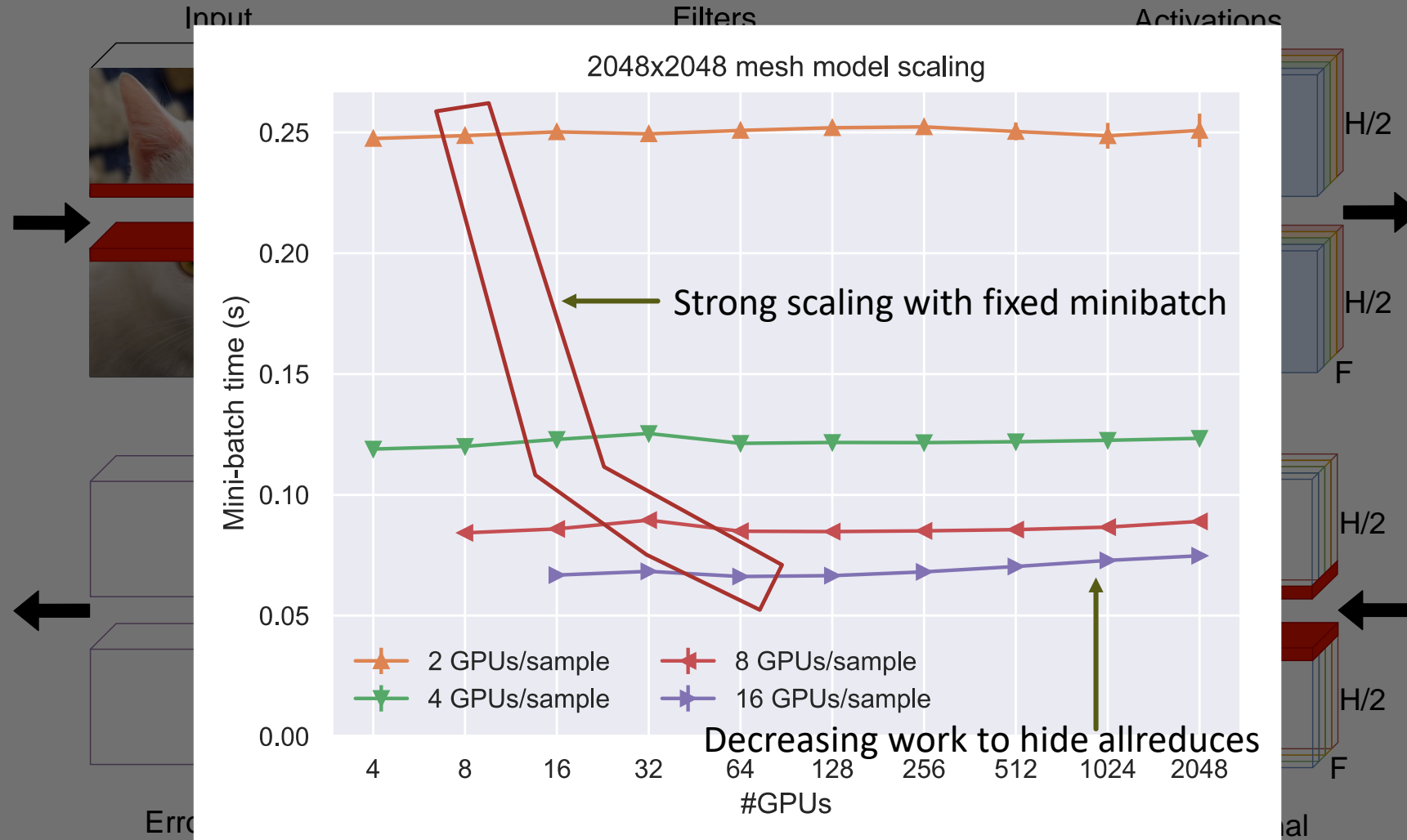
Spatial parallelism [Dryden et al. 2019]



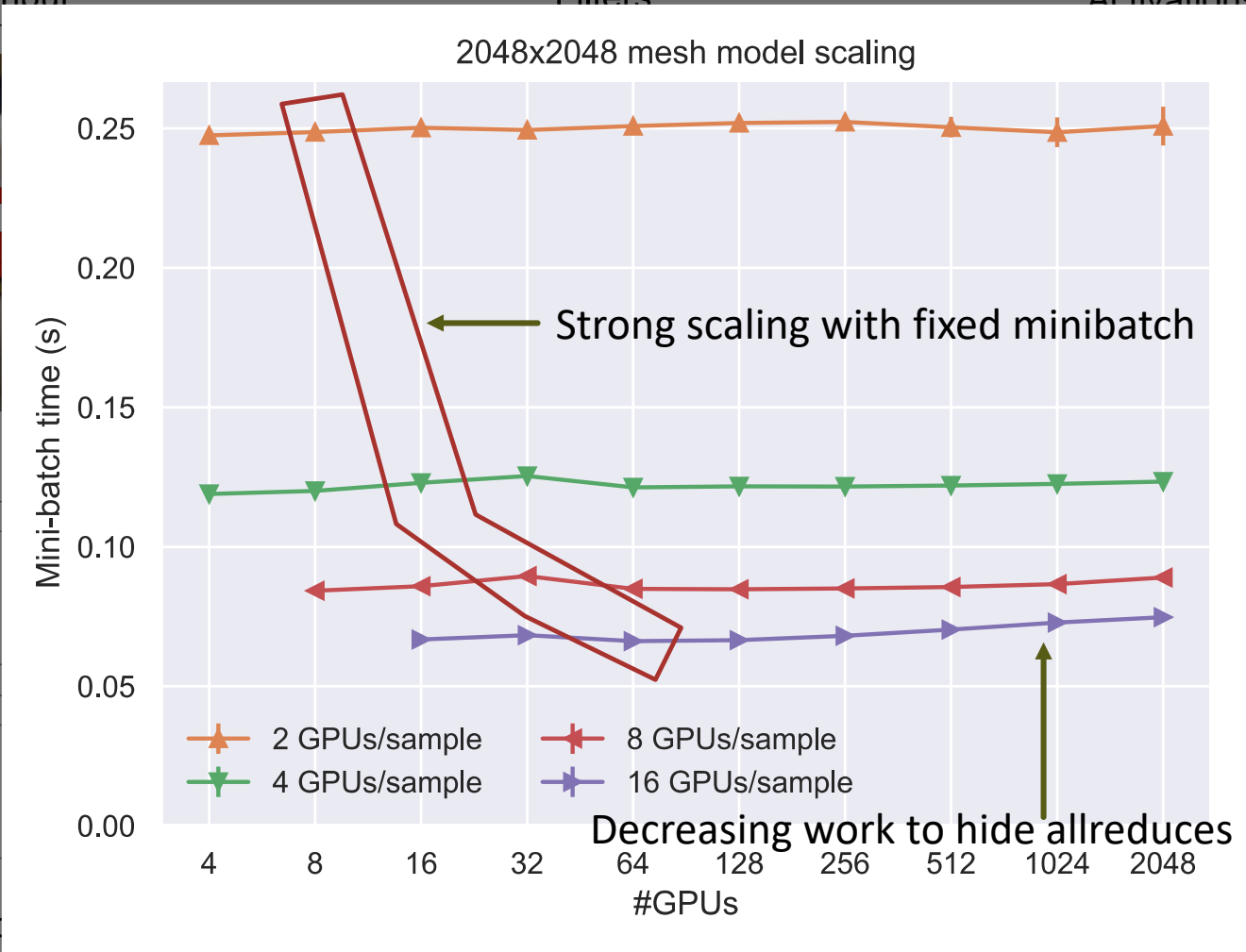
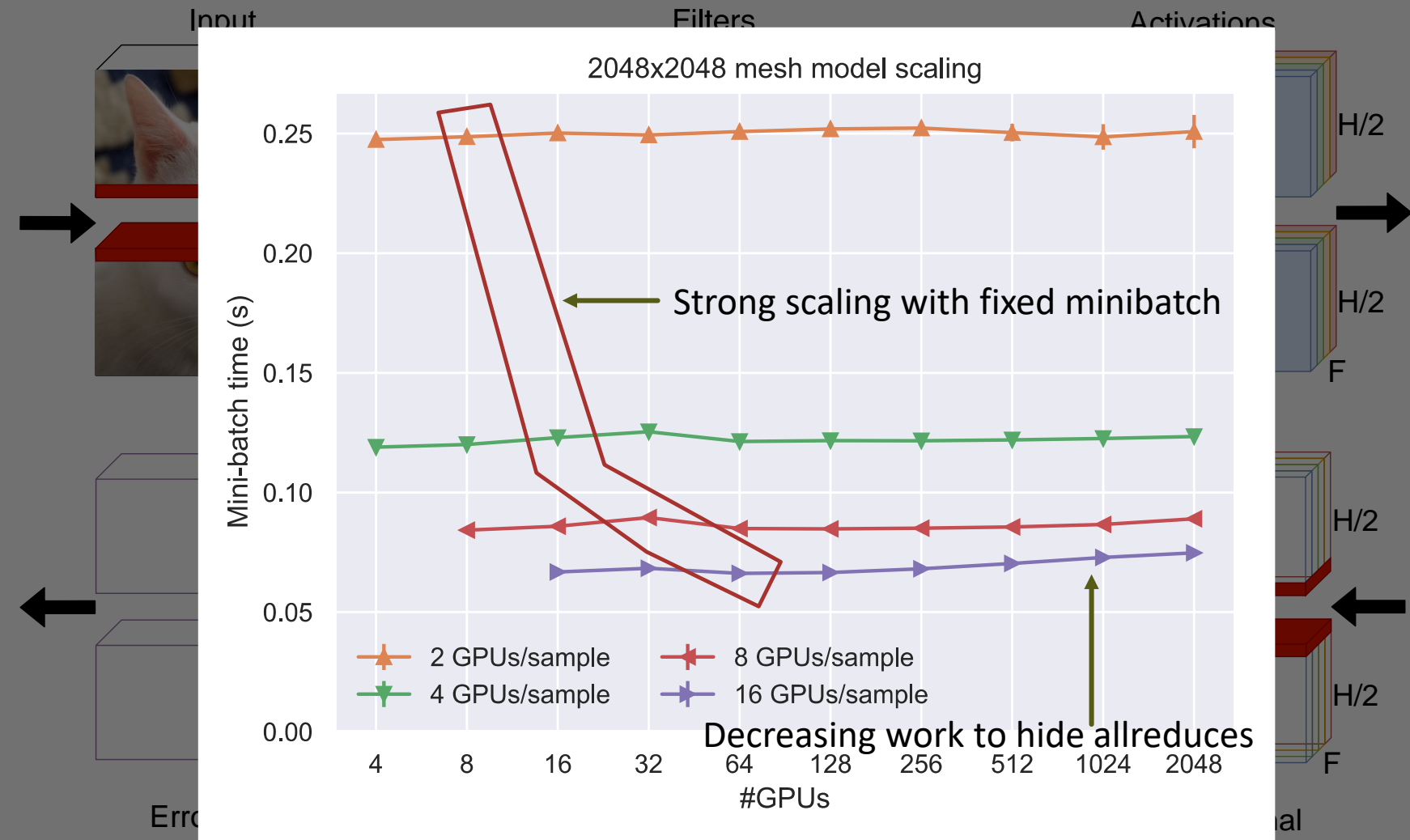
Spatial parallelism [Dryden et al. 2019]



Spatial parallelism [Dryden et al. 2019]

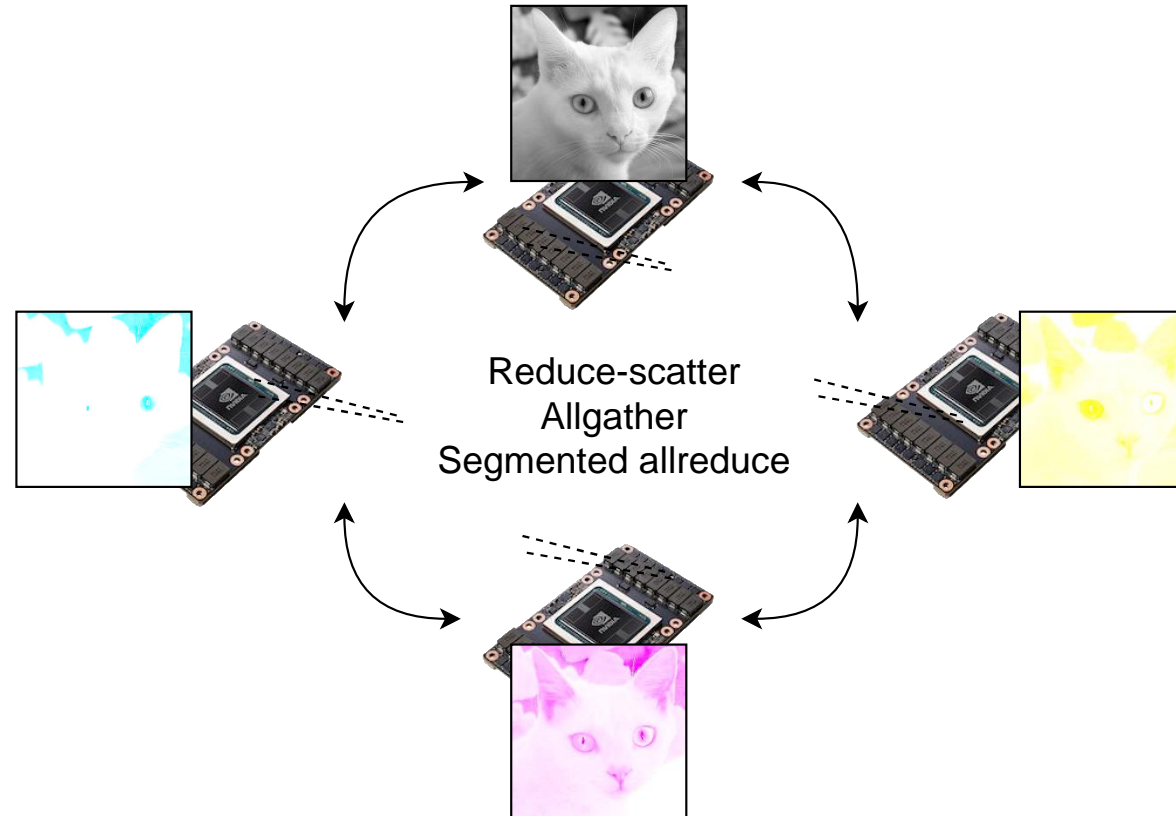


Spatial parallelism [Dryden et al. 2019]

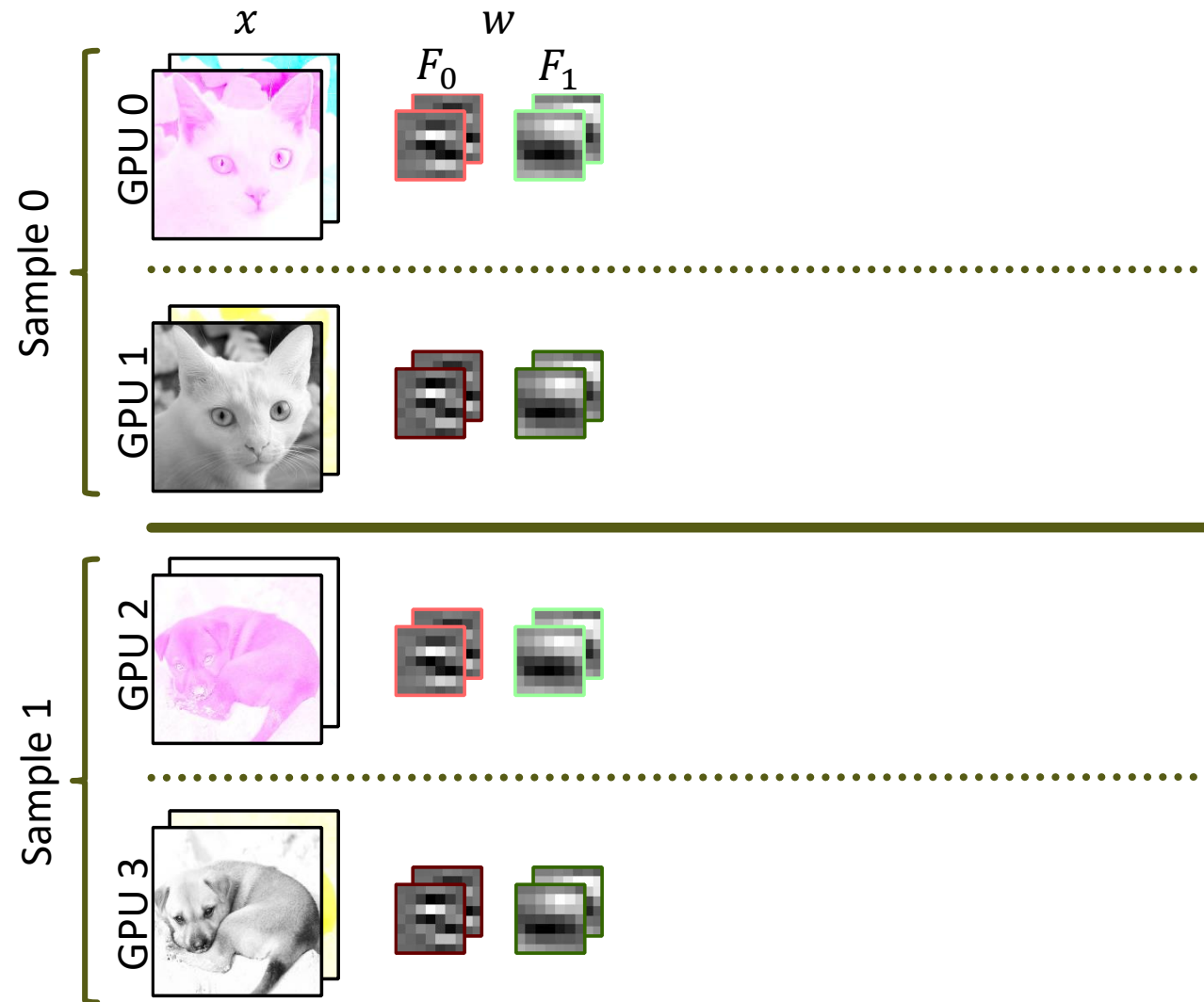


Channel/filter parallelism [Dryden et al. 2019]

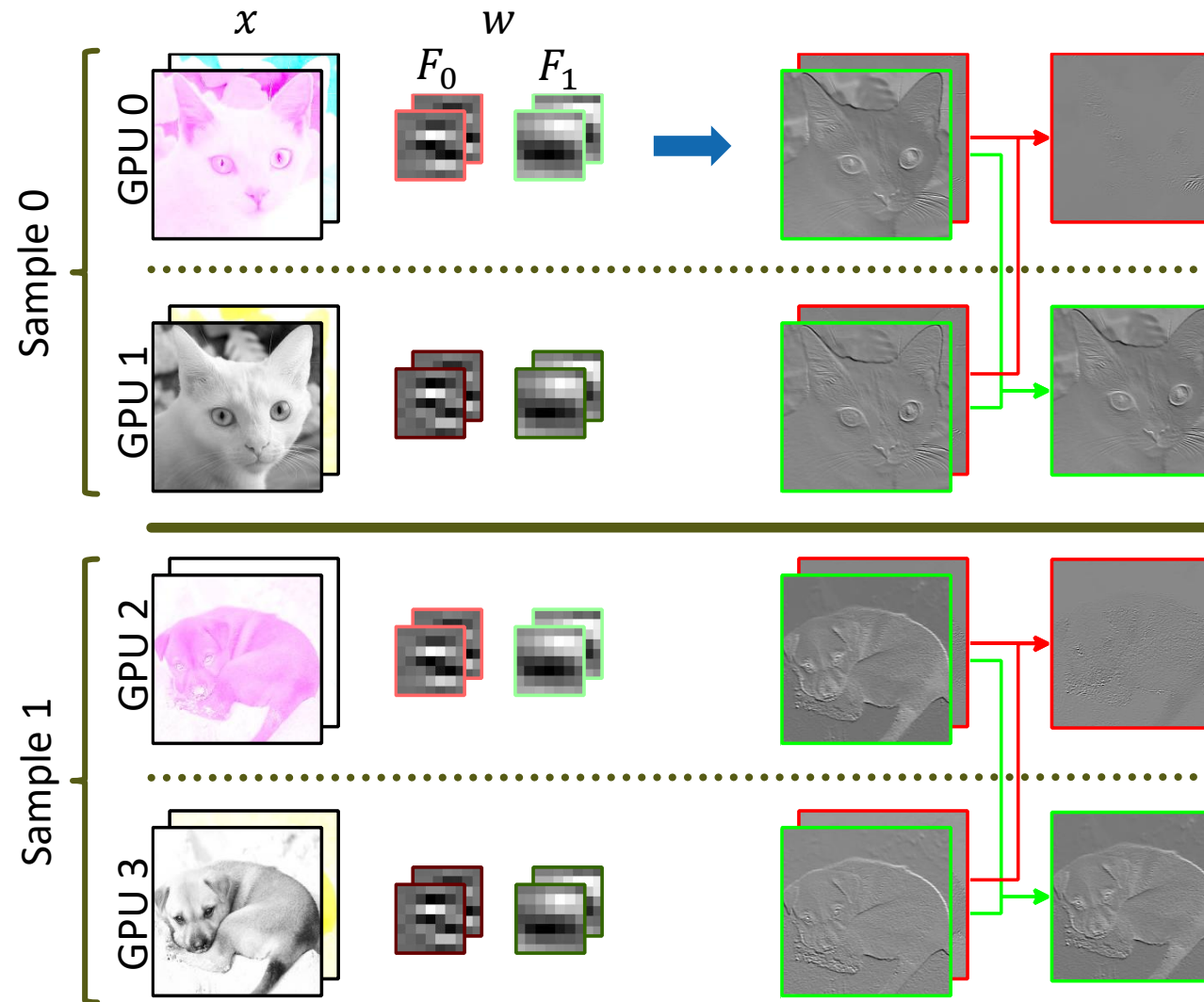
- Family of algorithms for jointly partitioning channels and filters in convolution



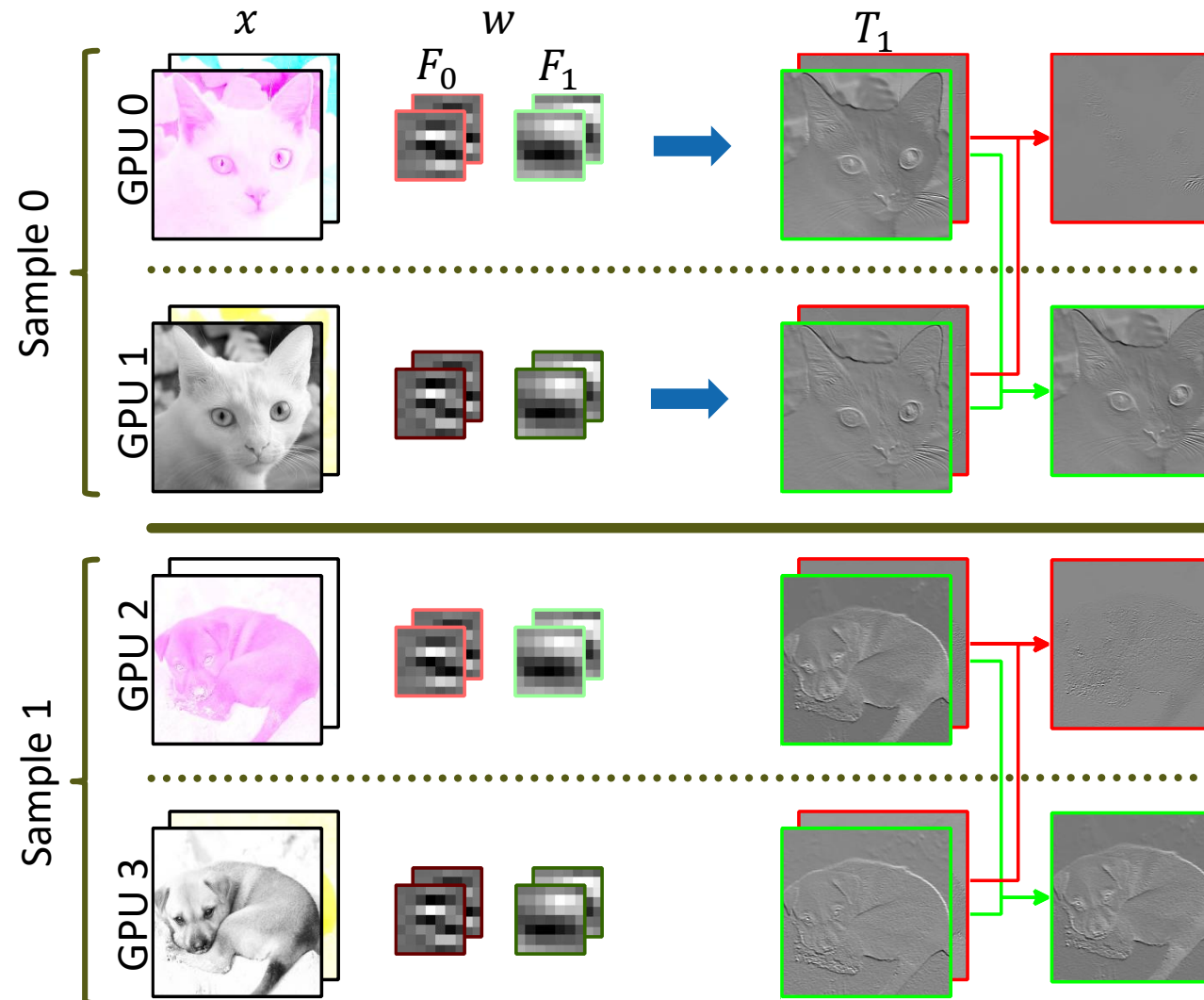
Stationary- x : Forward



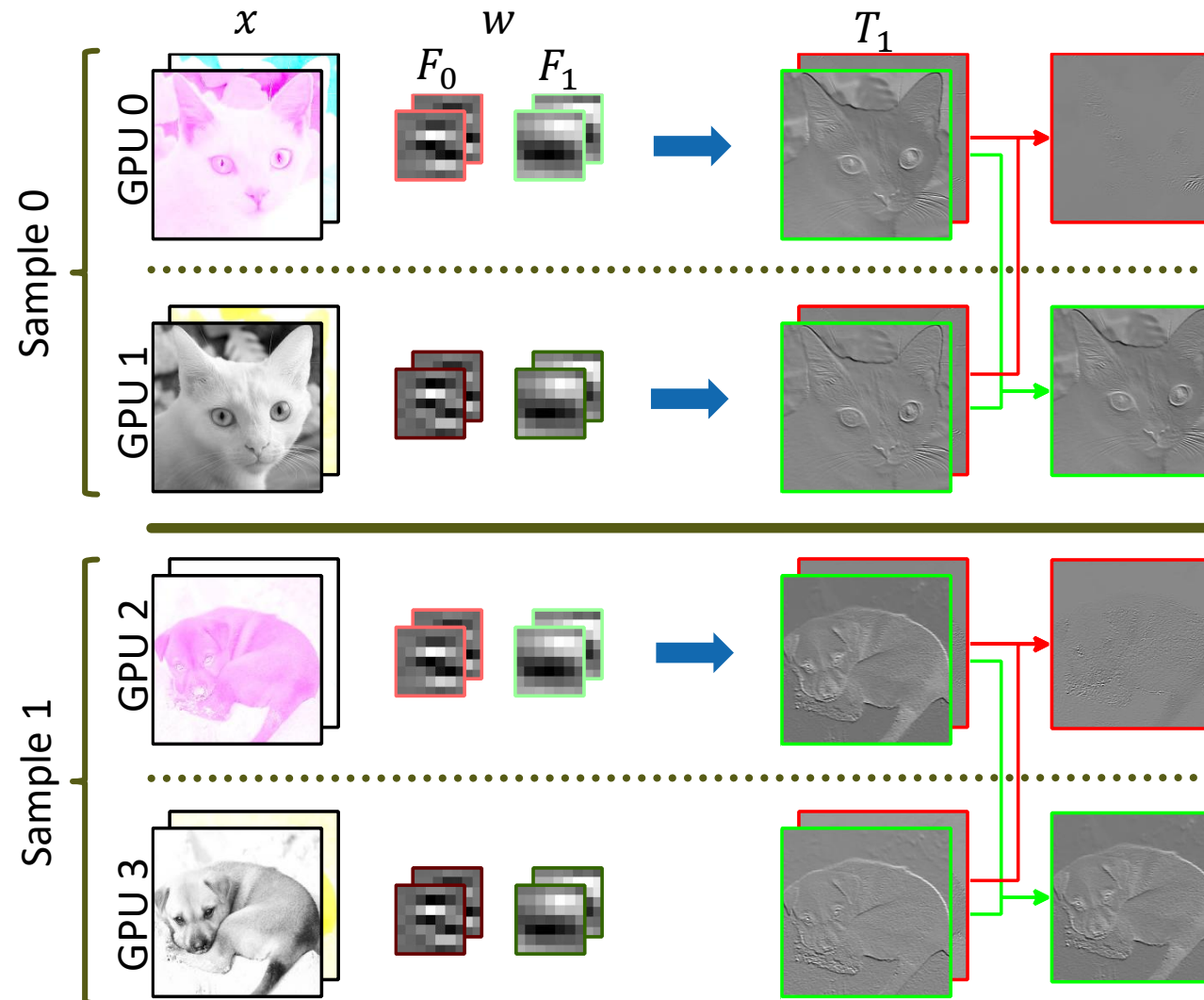
Stationary- x : Forward



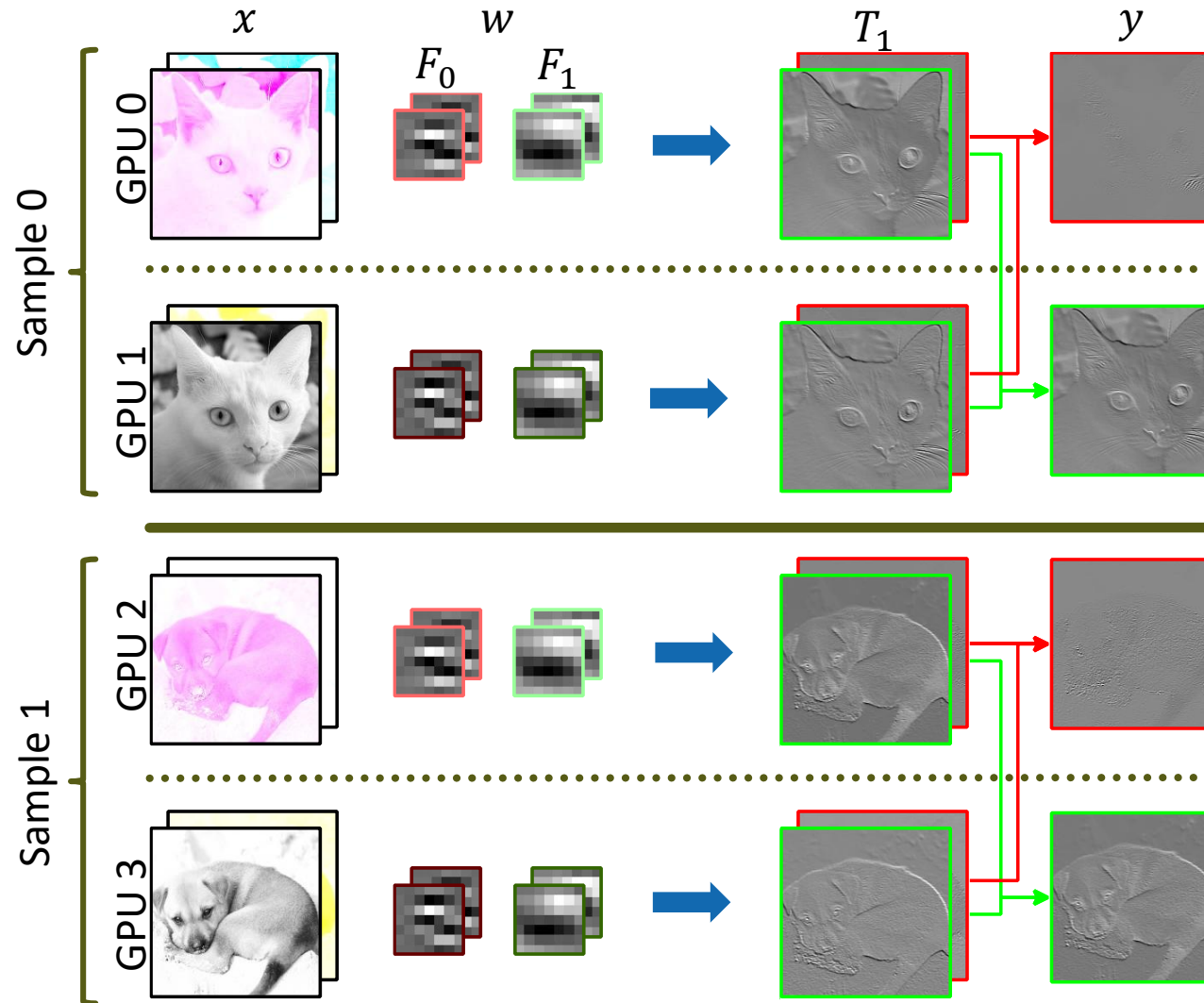
Stationary- x : Forward



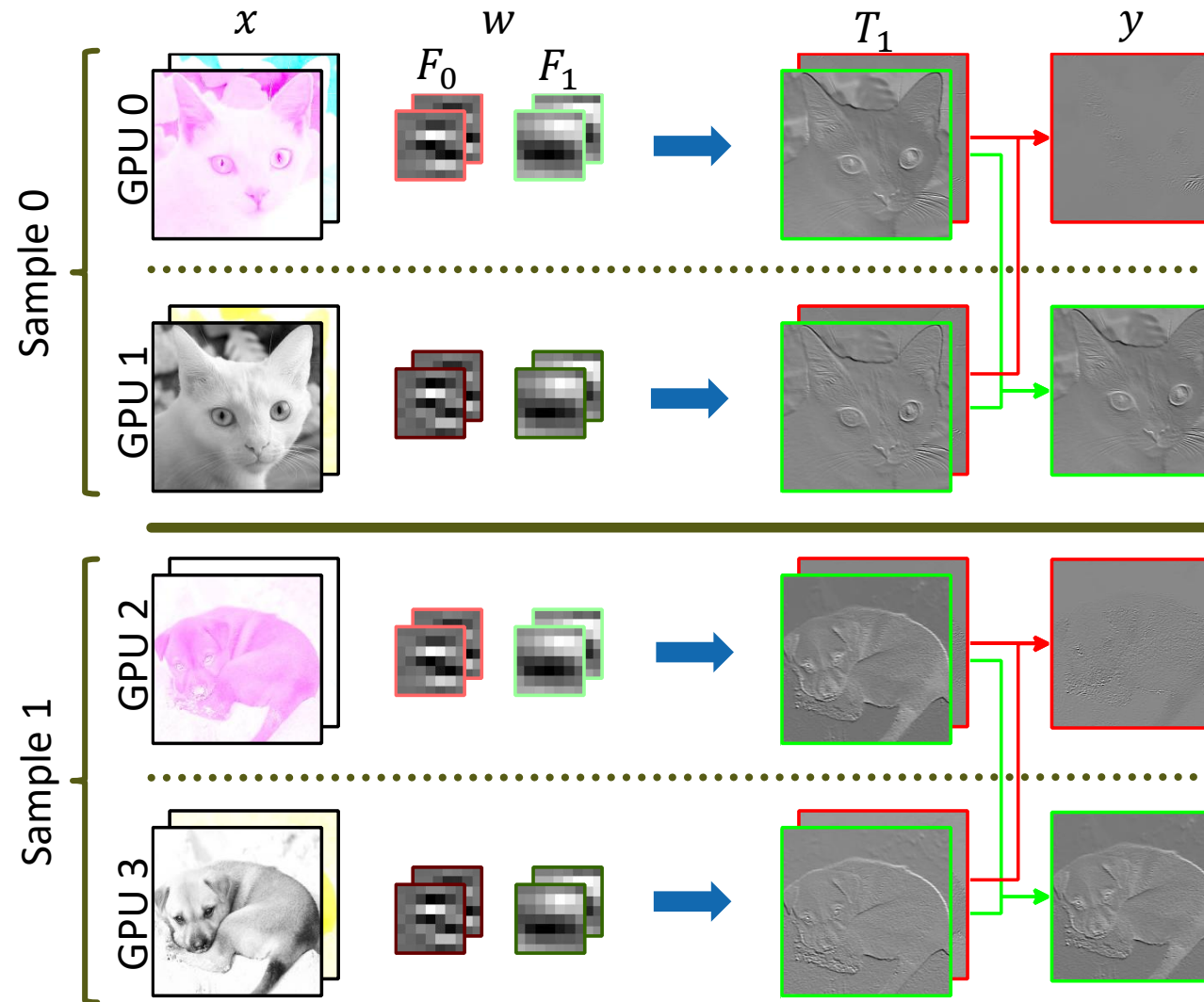
Stationary- x : Forward



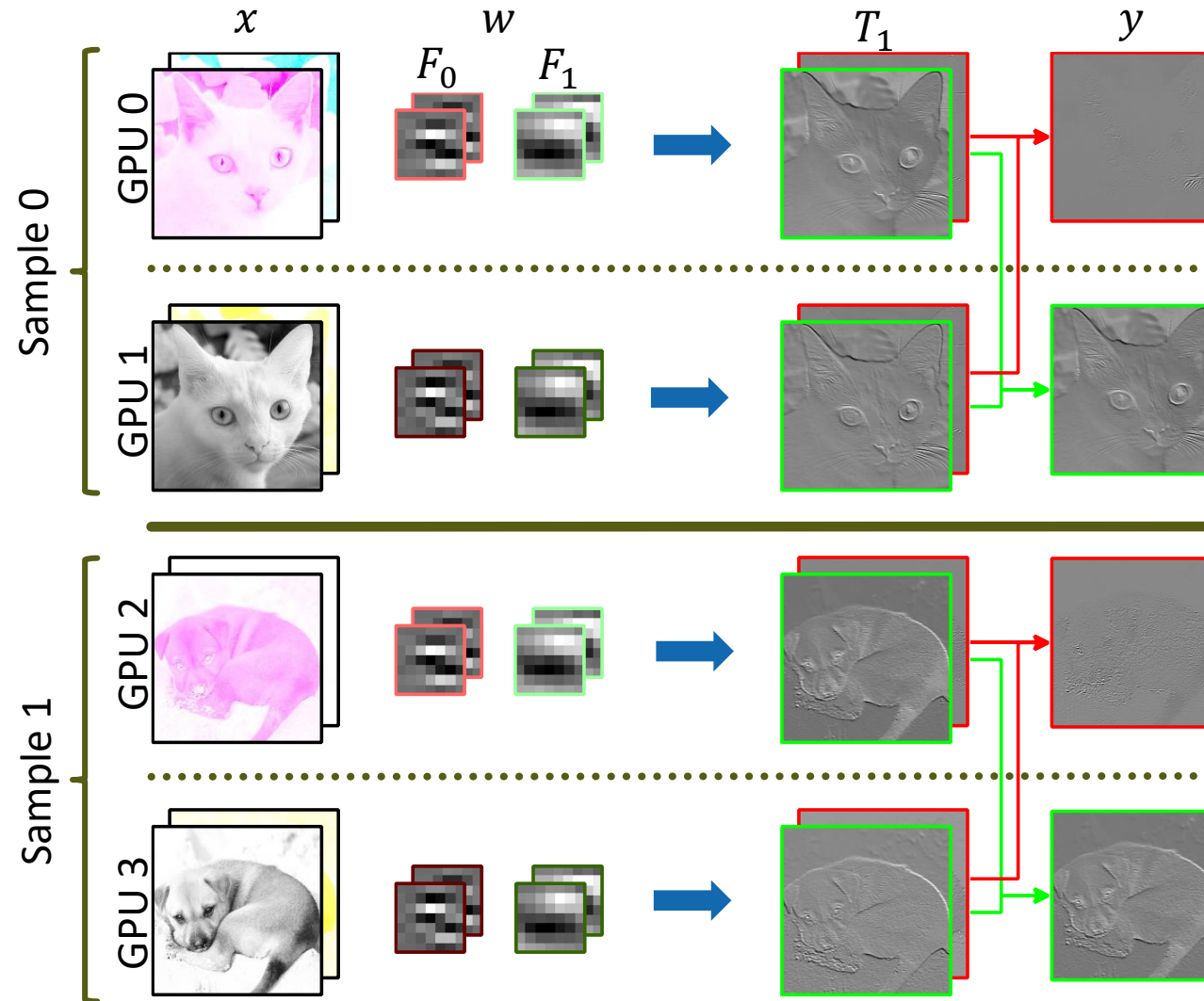
Stationary- x : Forward



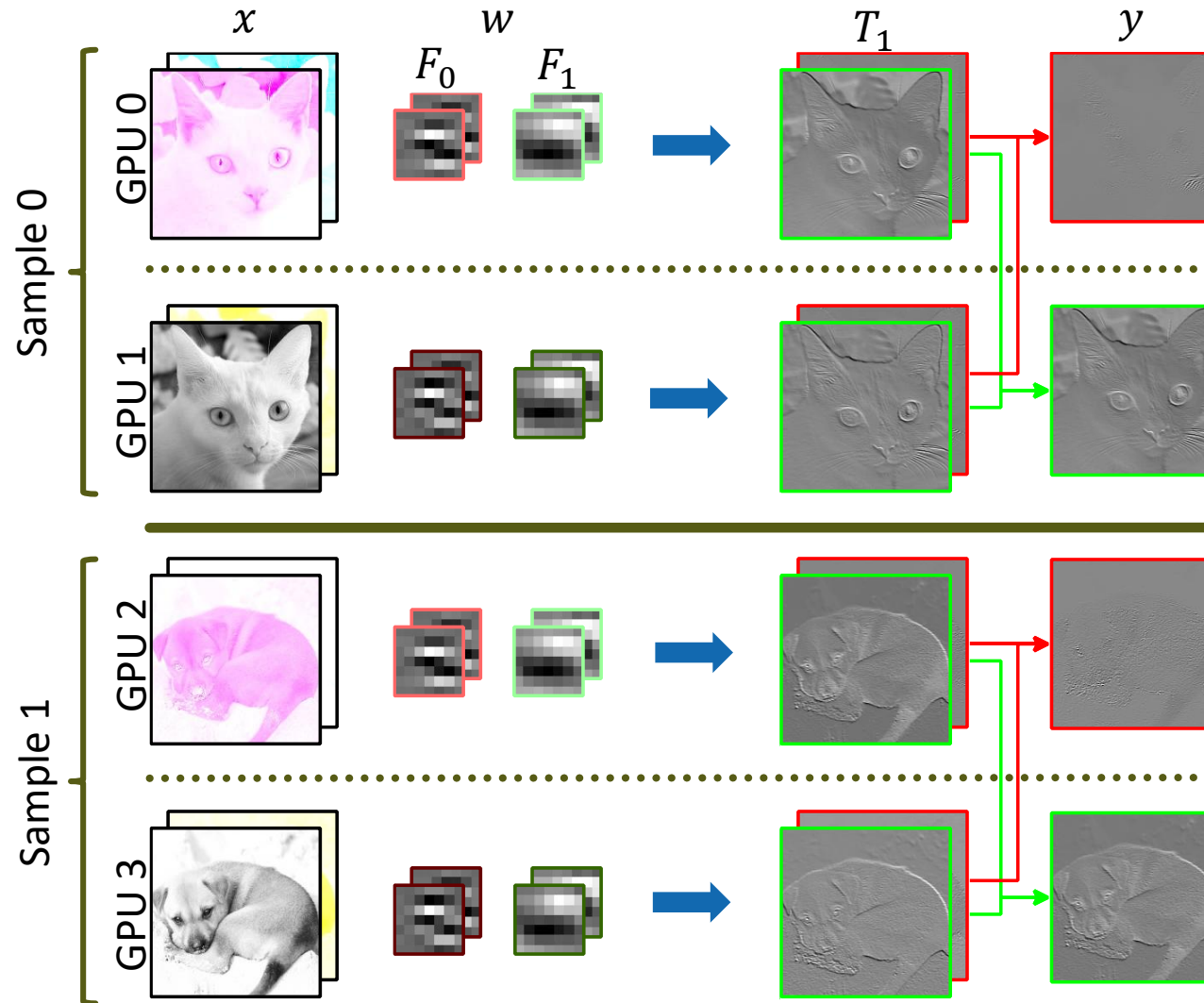
Stationary- x : Forward



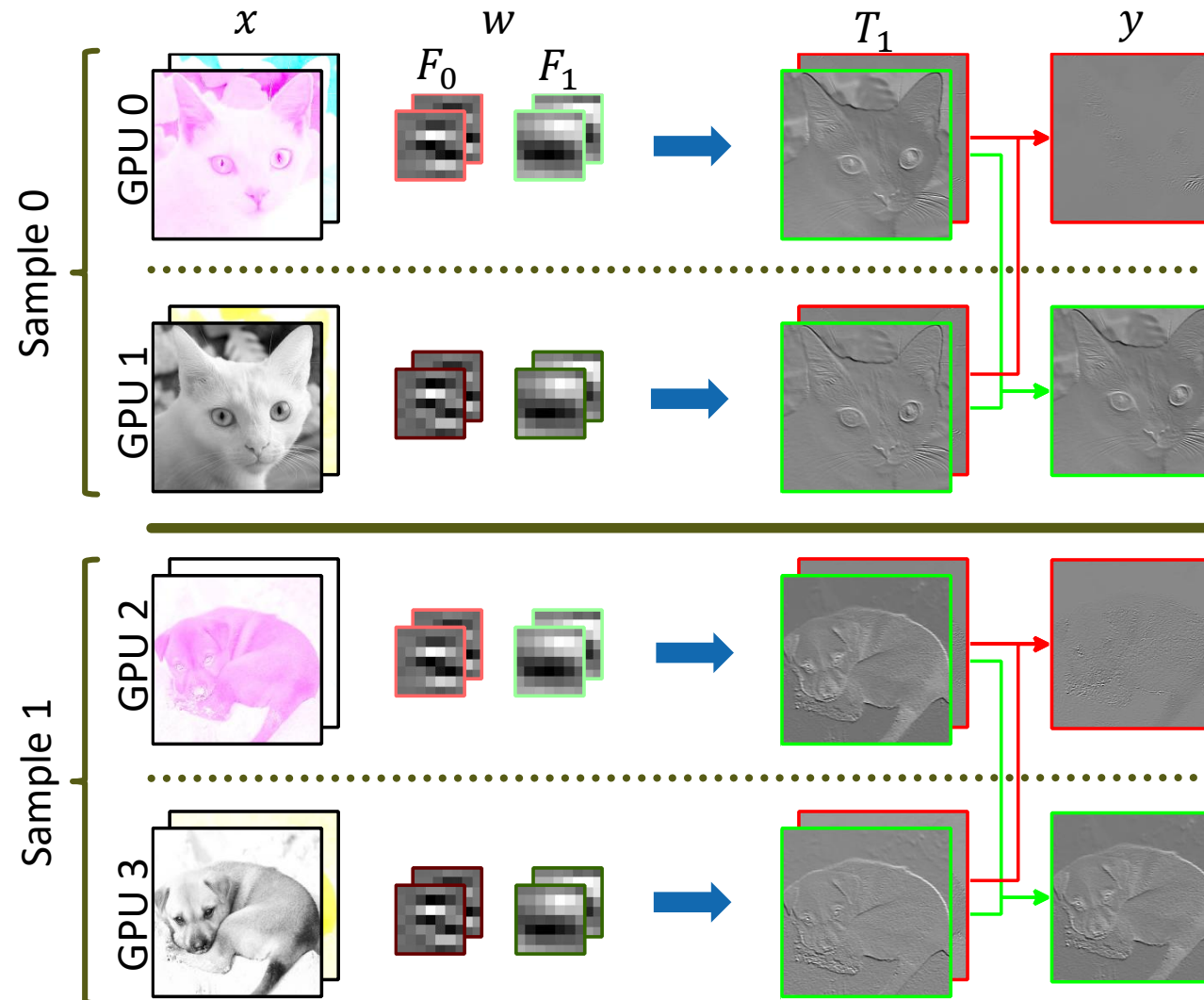
Stationary- x : Forward



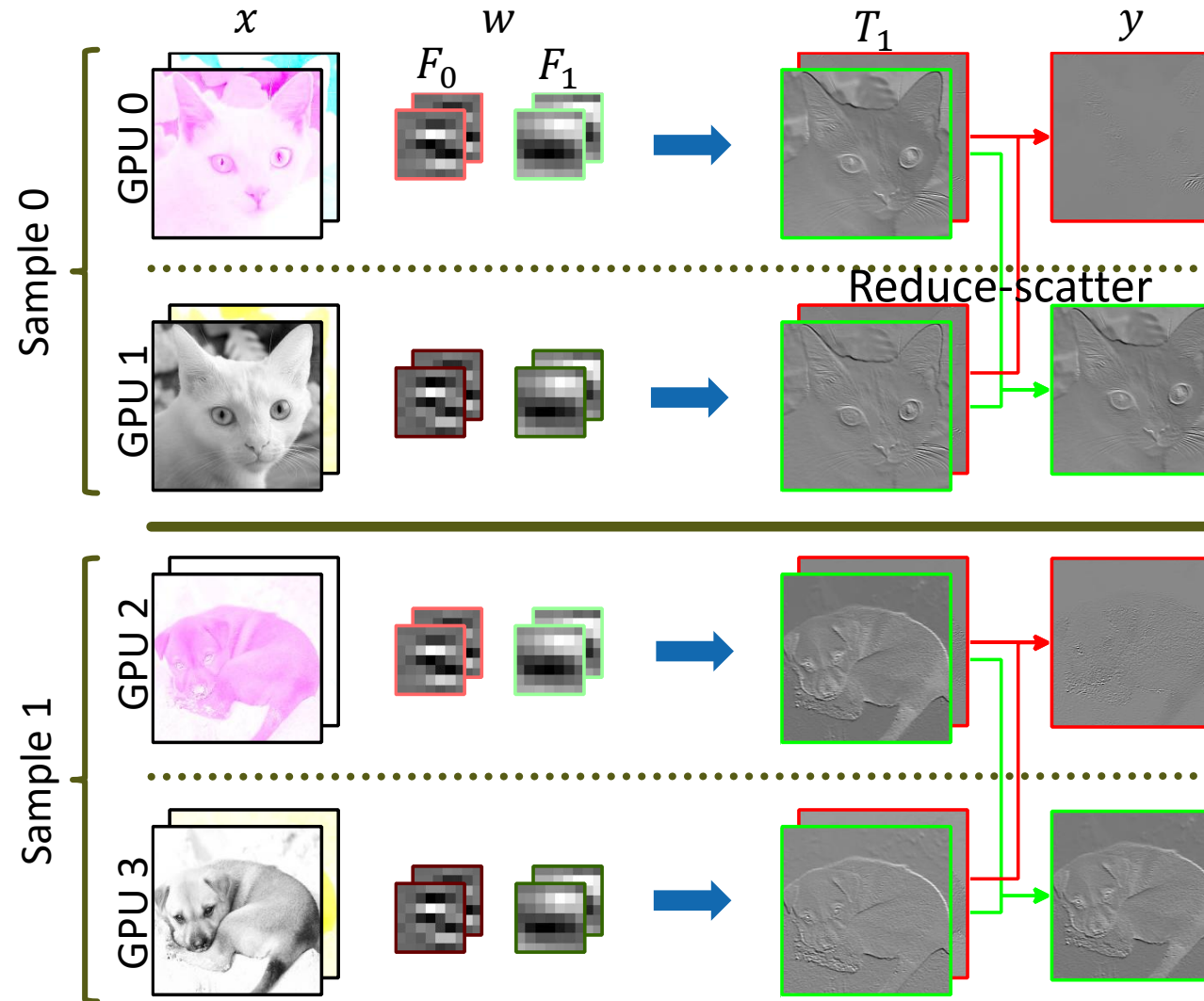
Stationary- x : Forward



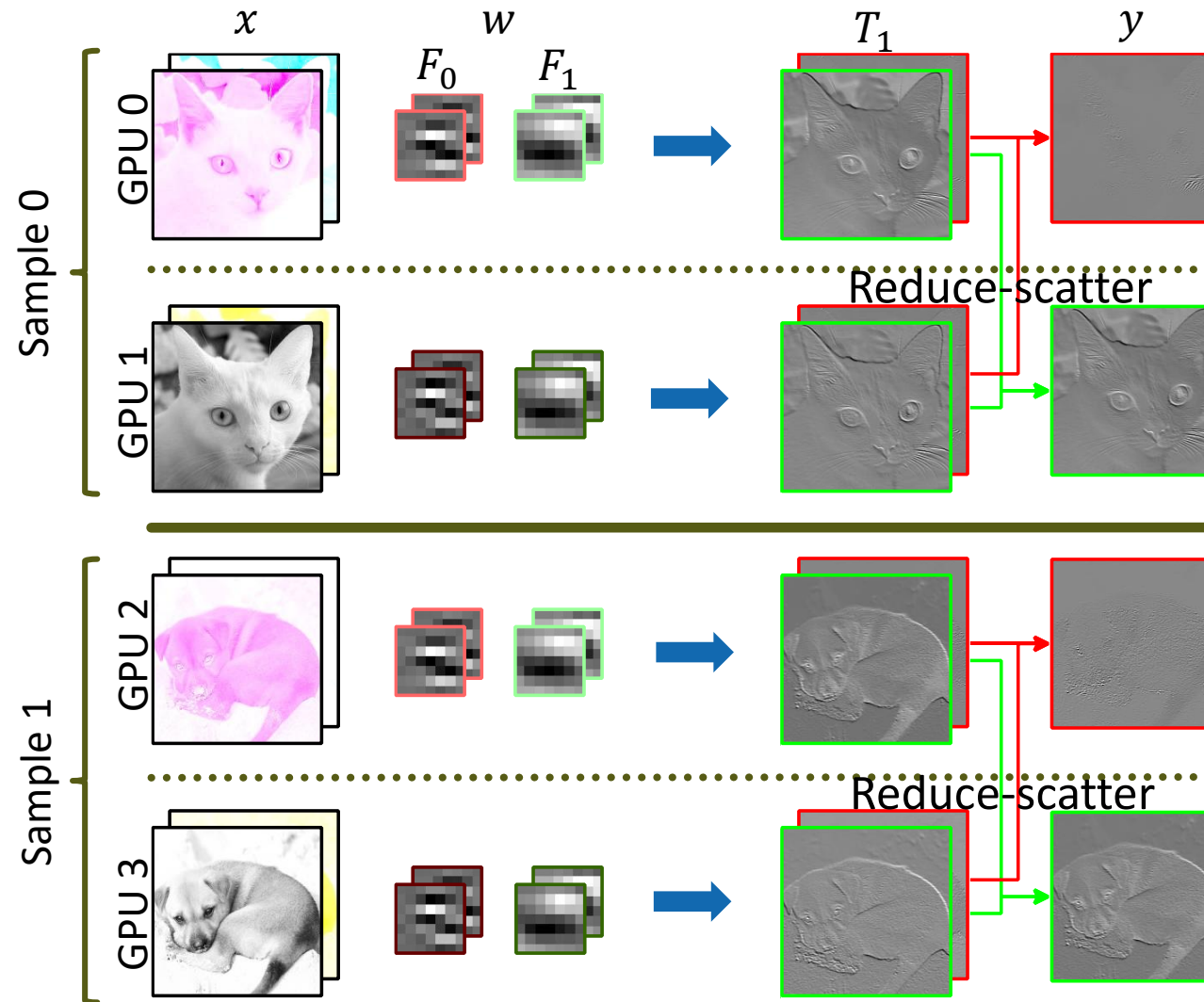
Stationary- x : Forward



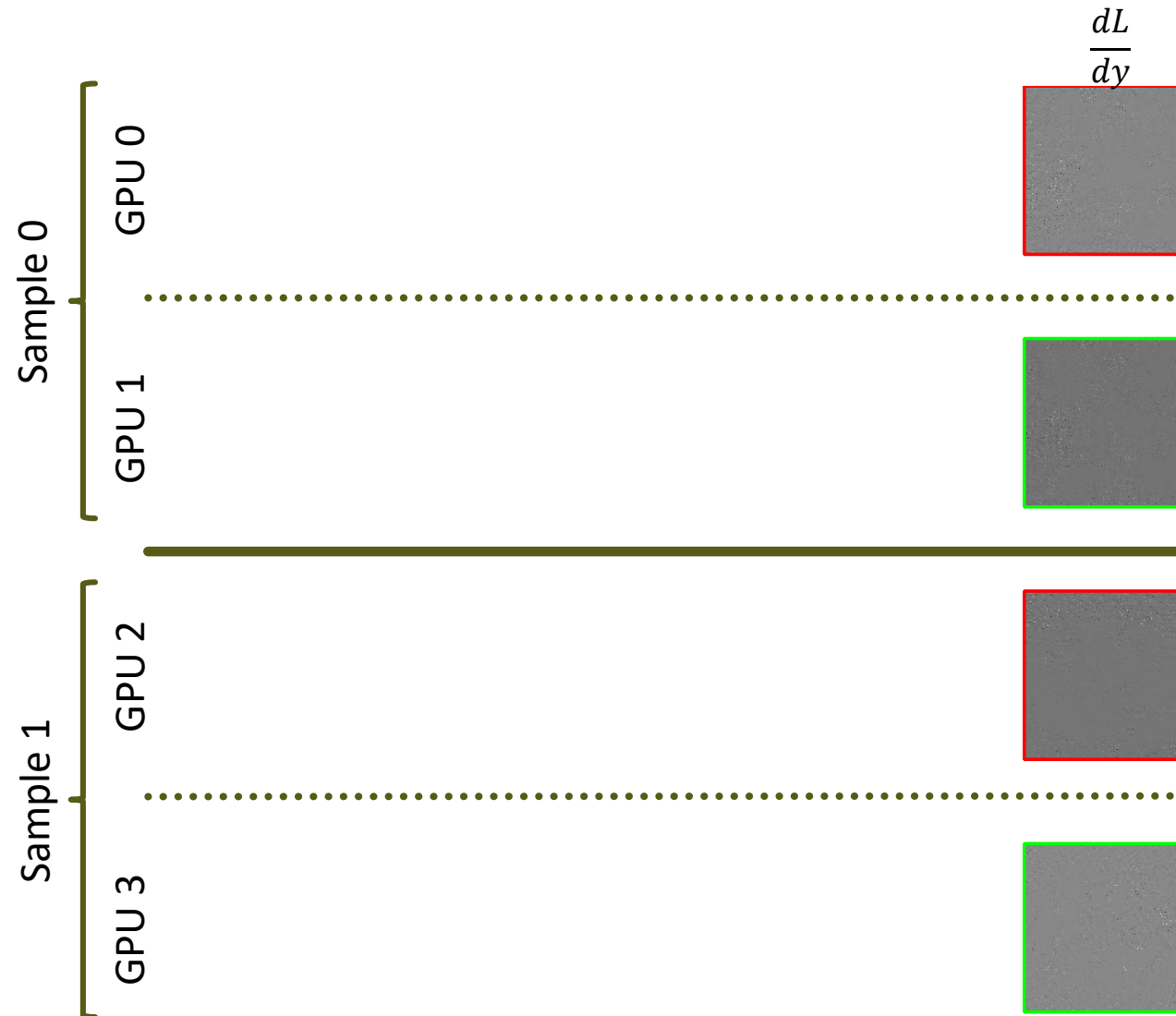
Stationary- x : Forward



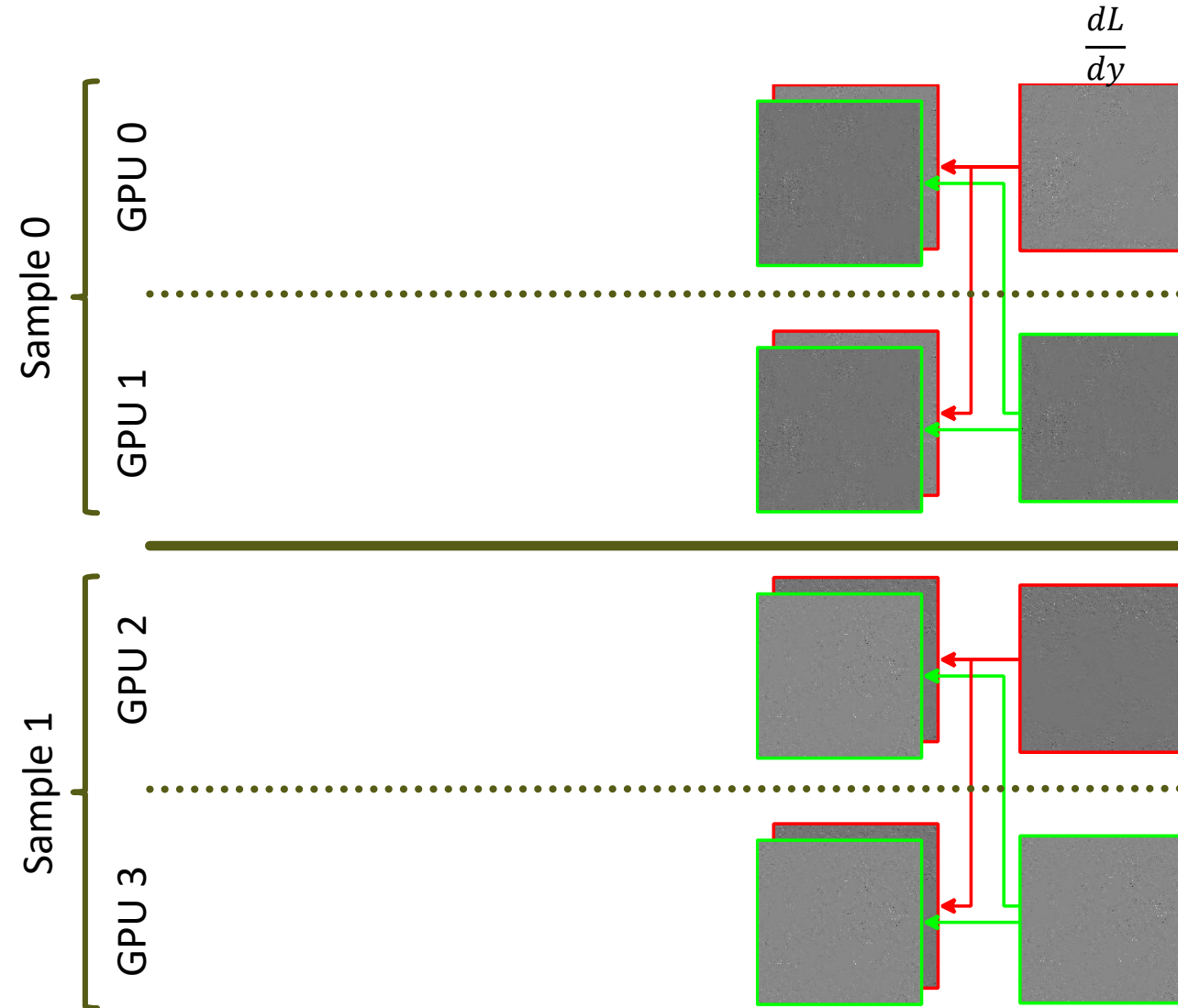
Stationary- x : Forward



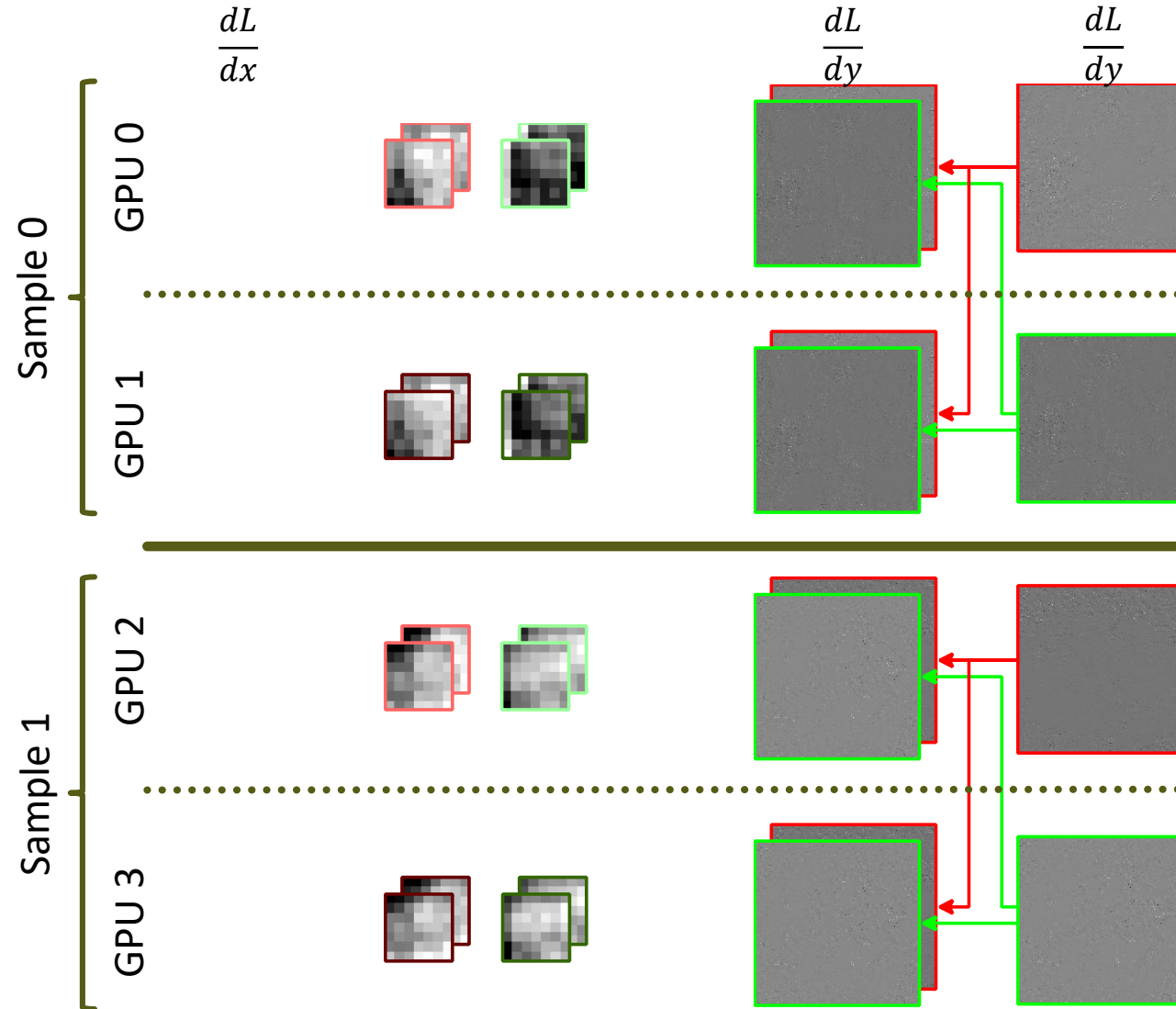
Stationary- x : Backward



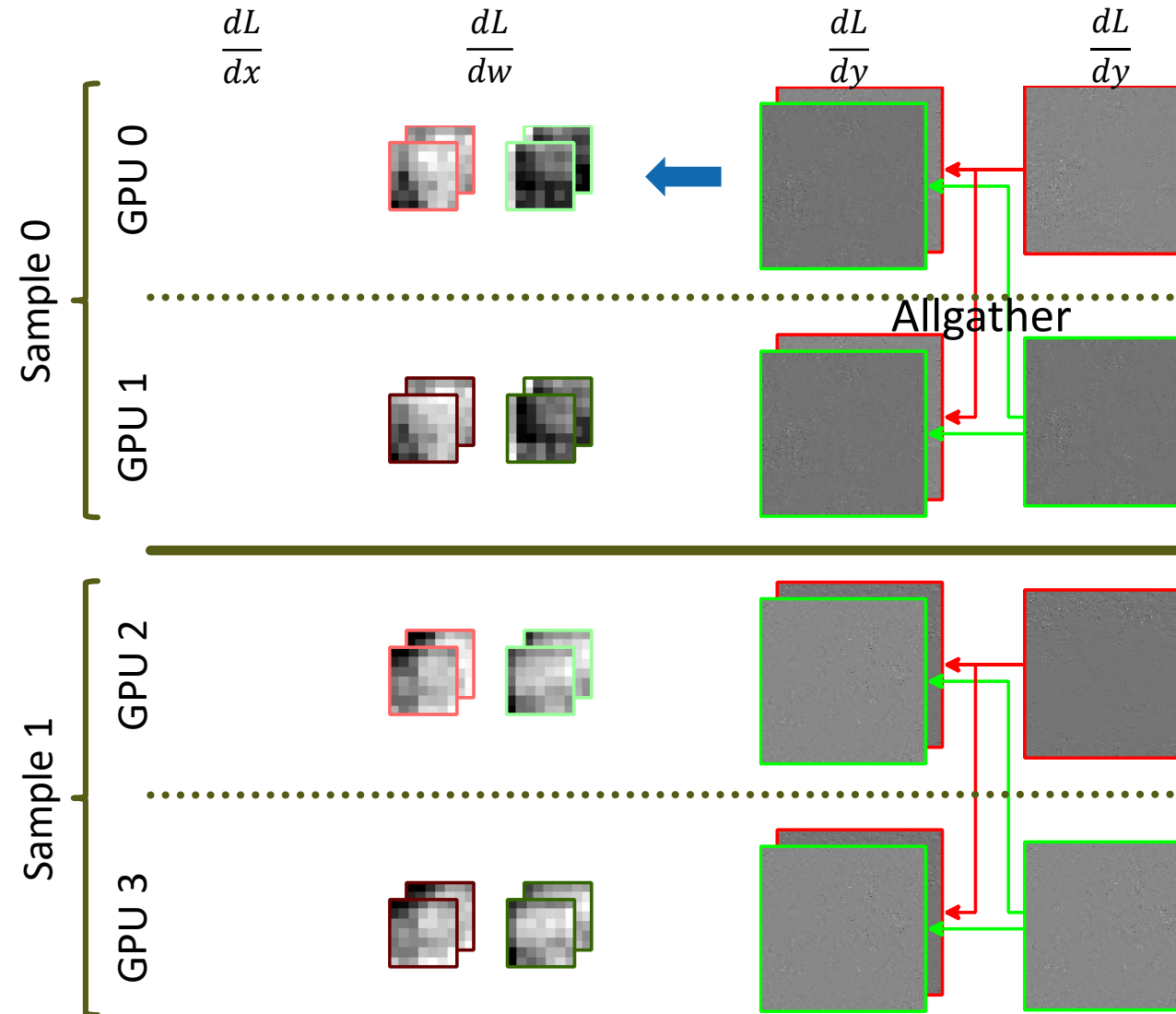
Stationary- x : Backward



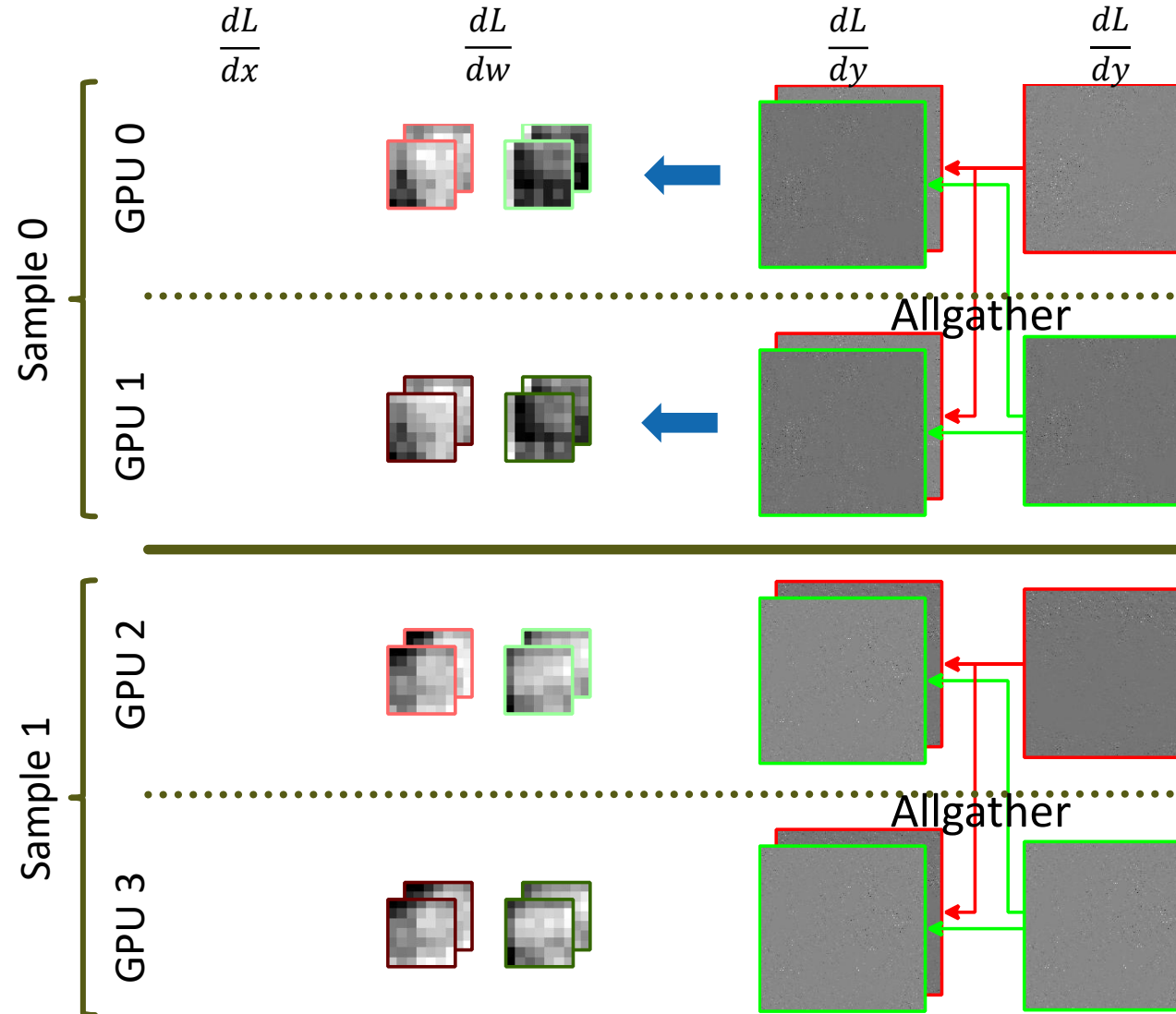
Stationary- x : Backward



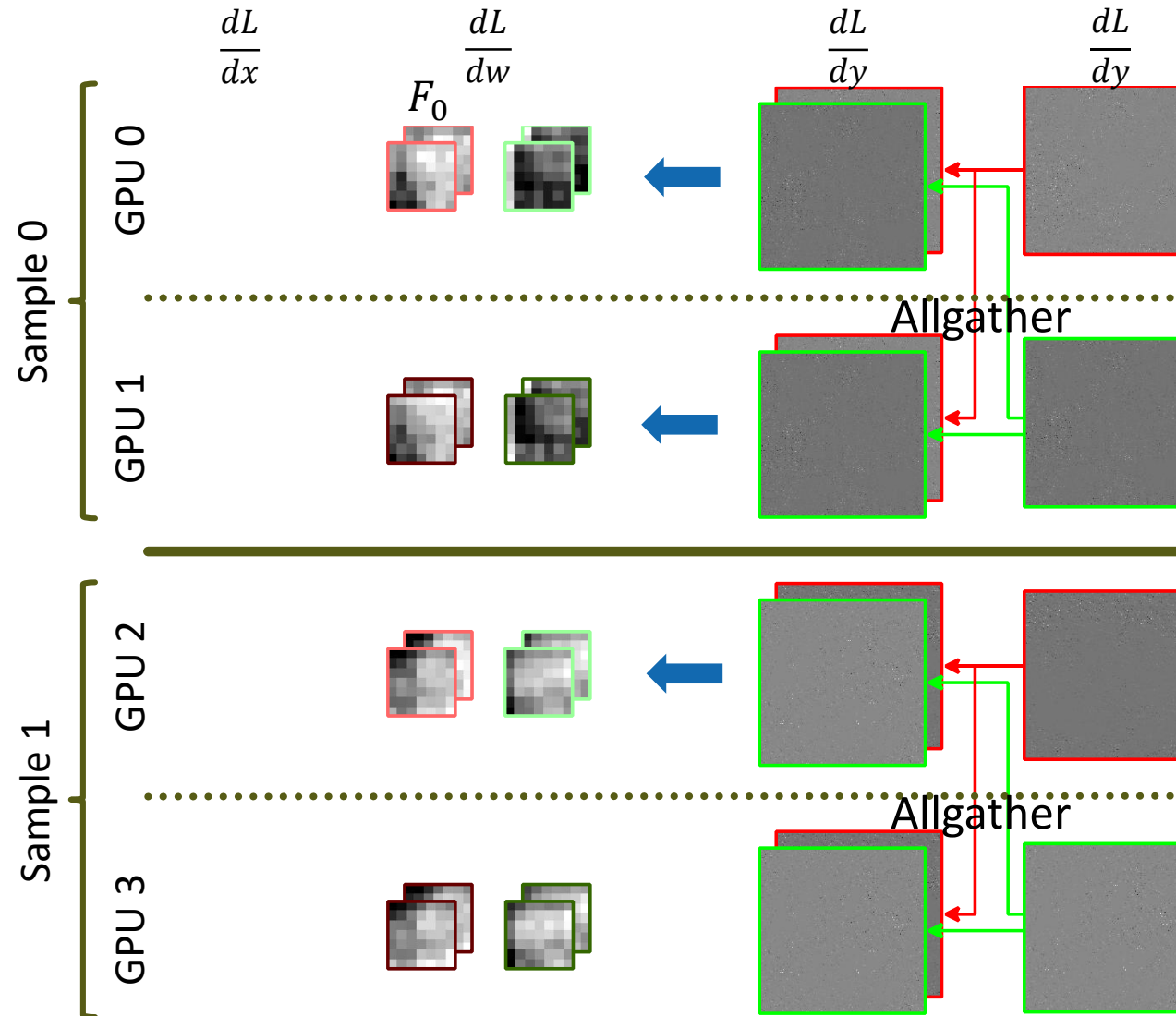
Stationary- x : Backward



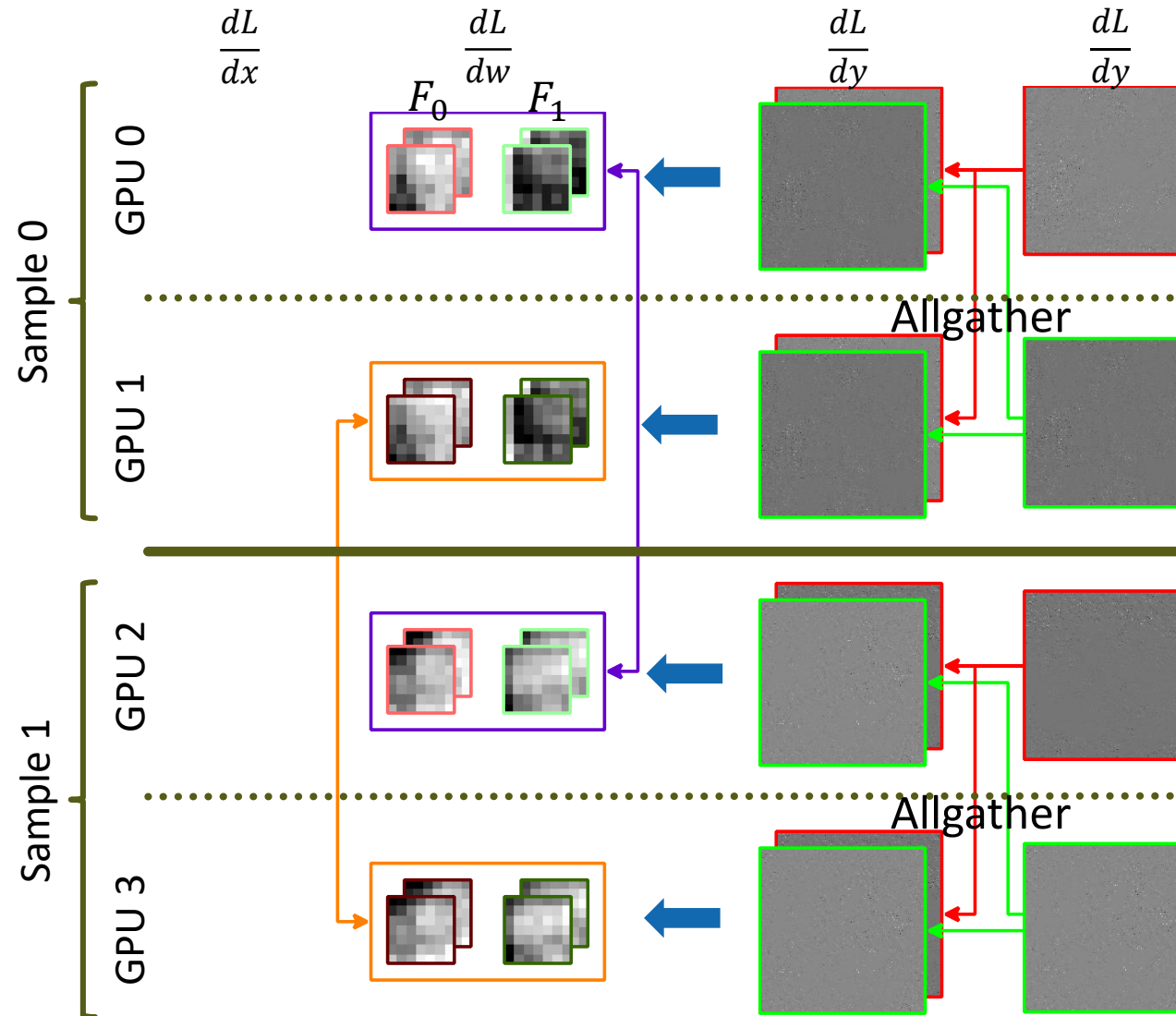
Stationary- x : Backward



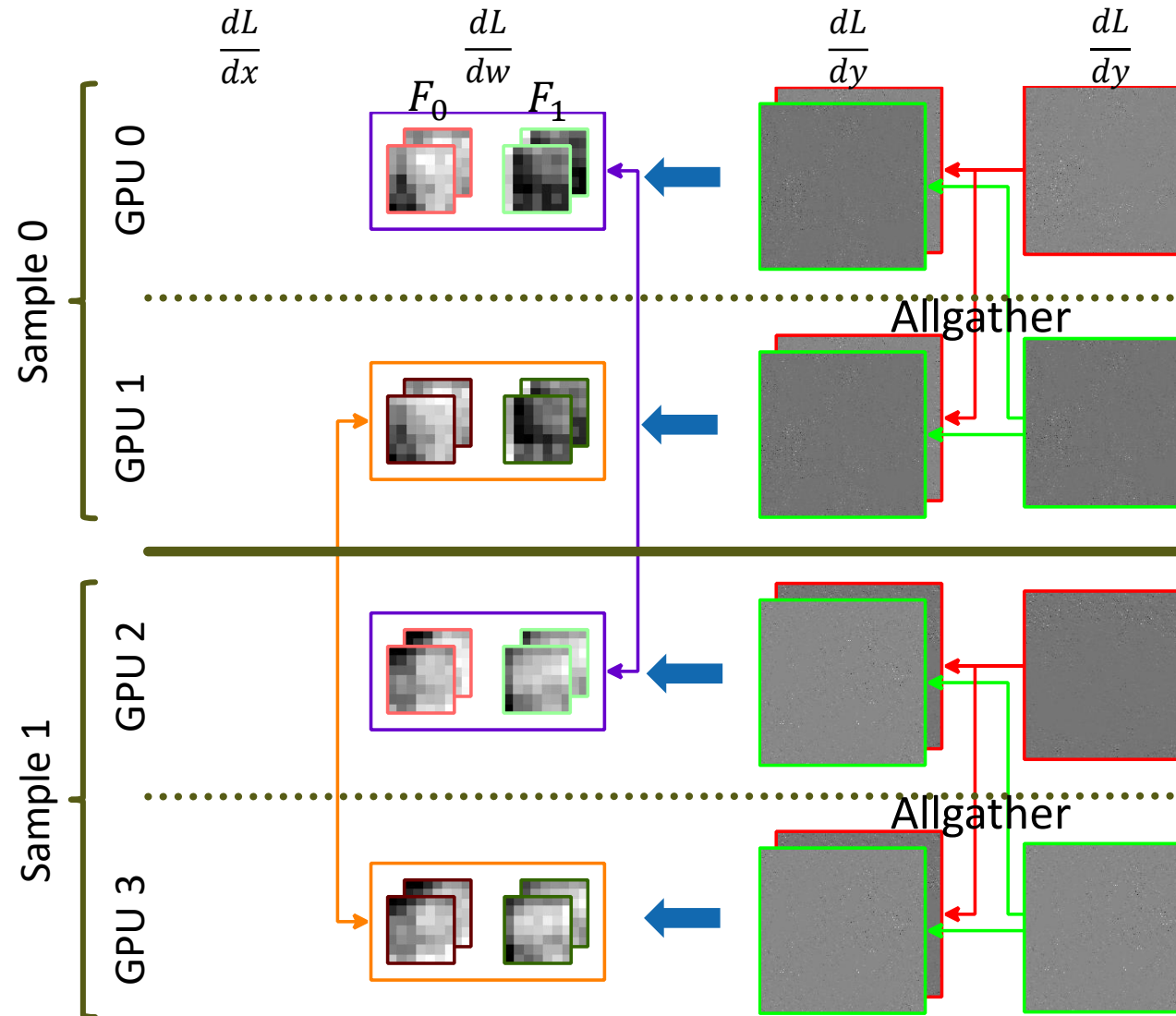
Stationary- x : Backward



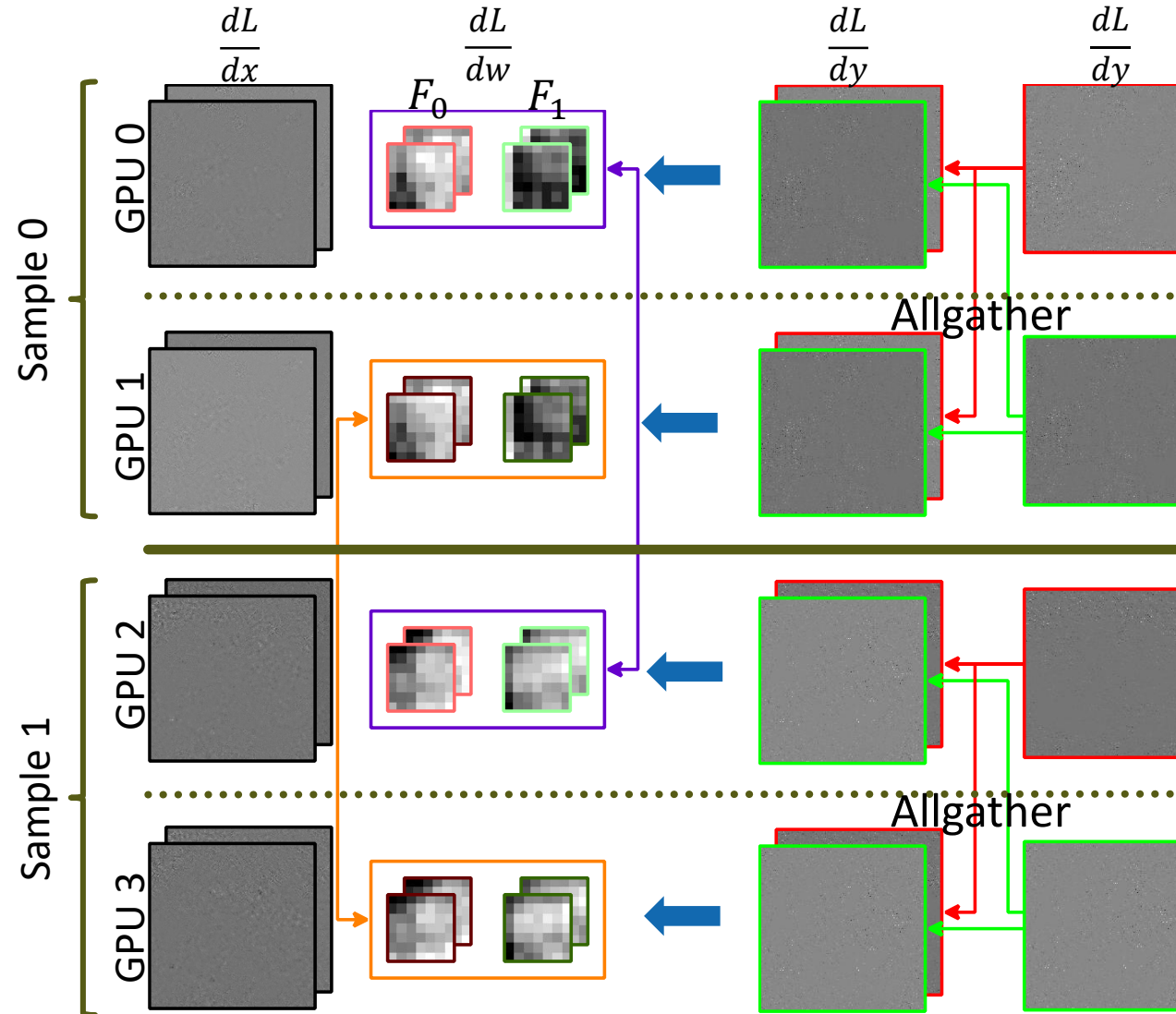
Stationary- x : Backward



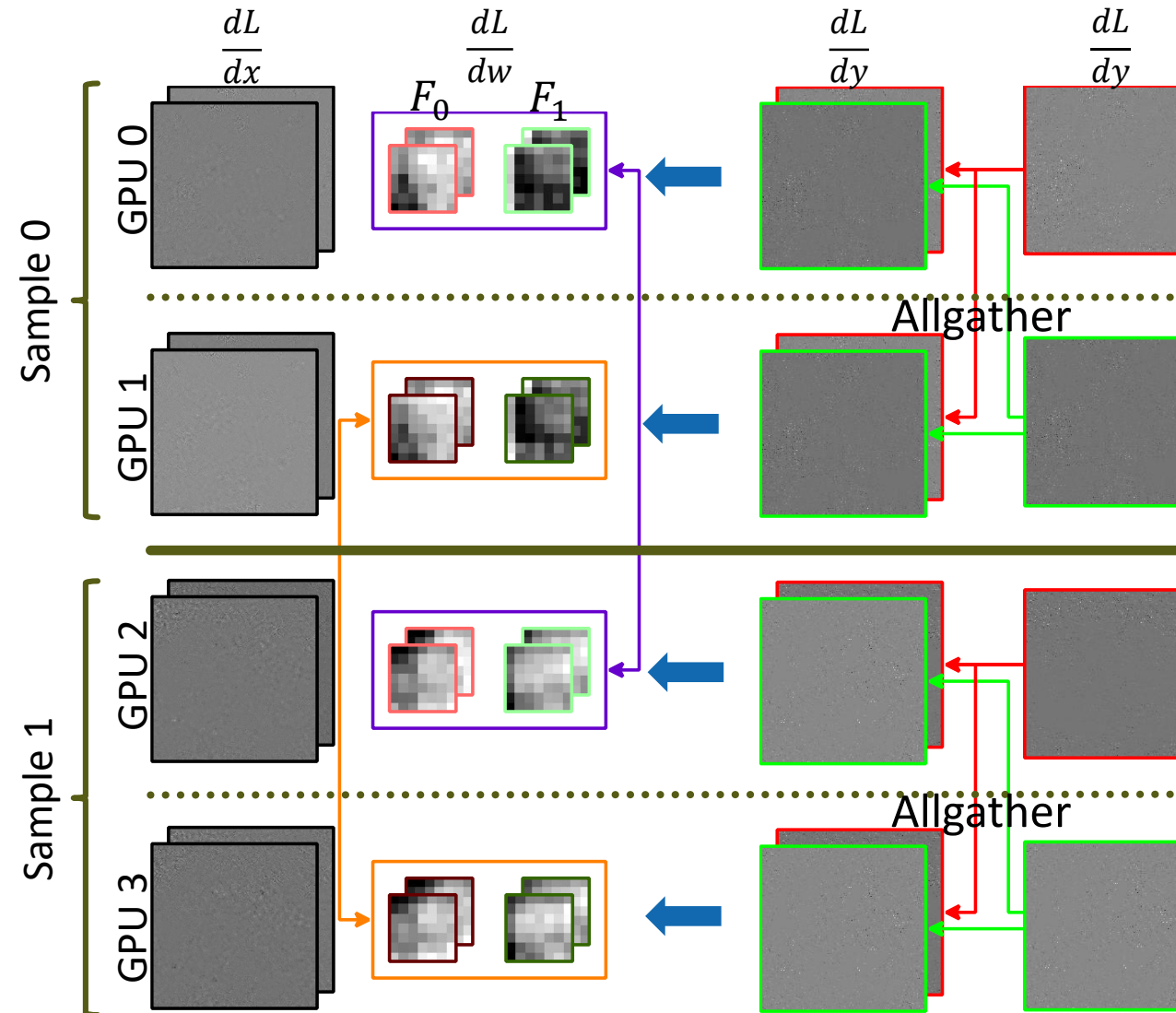
Stationary- x : Backward



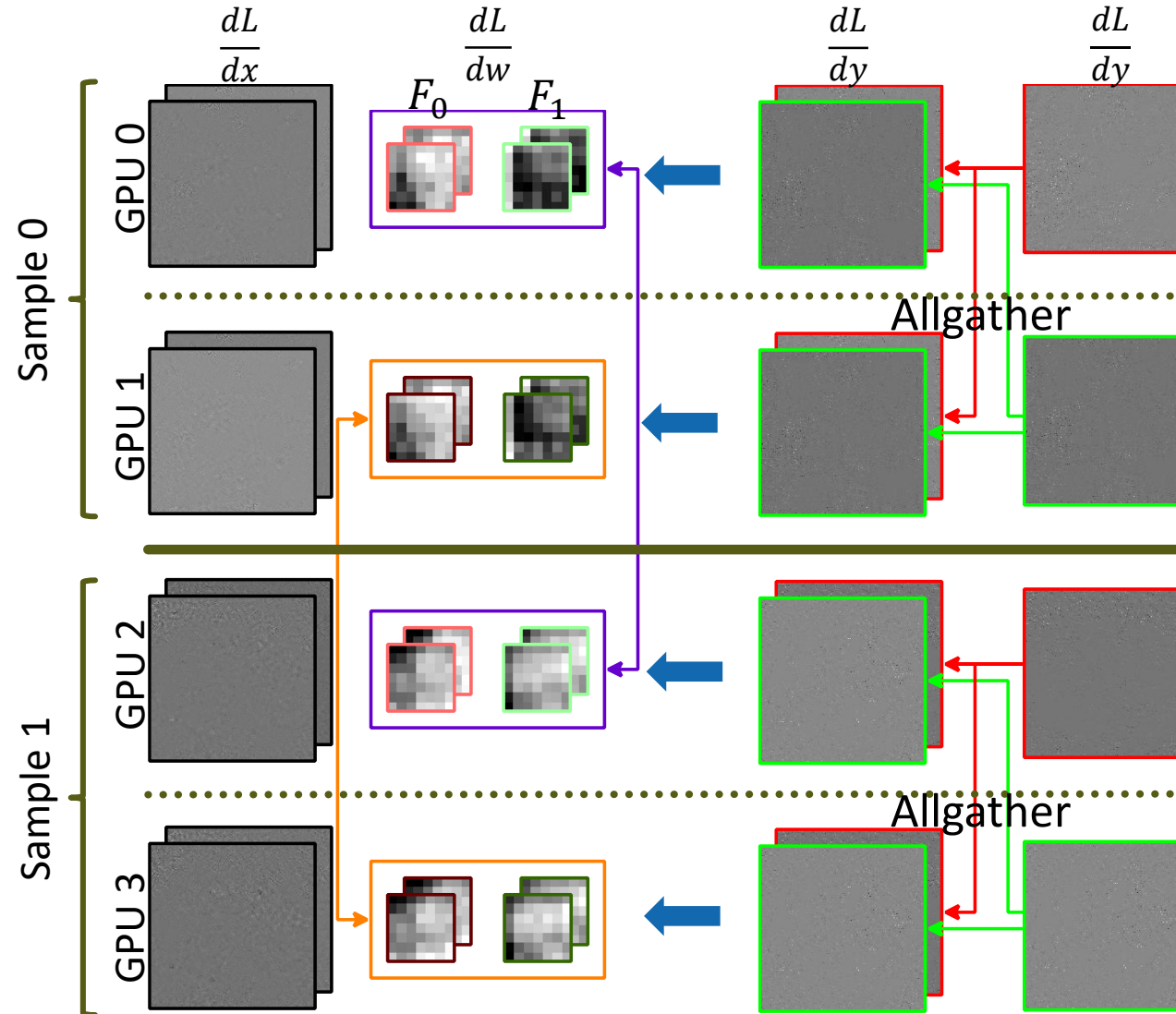
Stationary- x : Backward



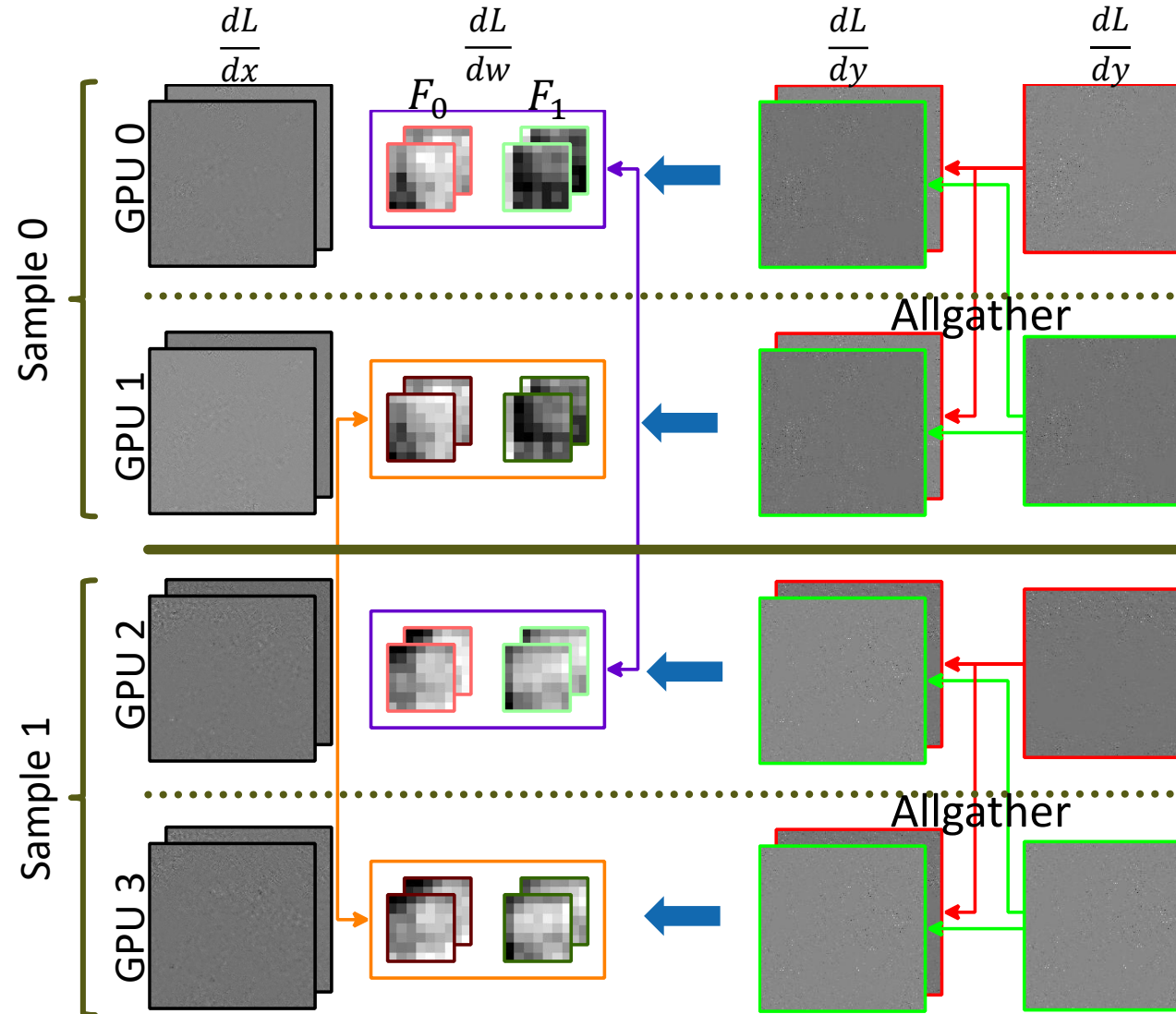
Stationary- x : Backward



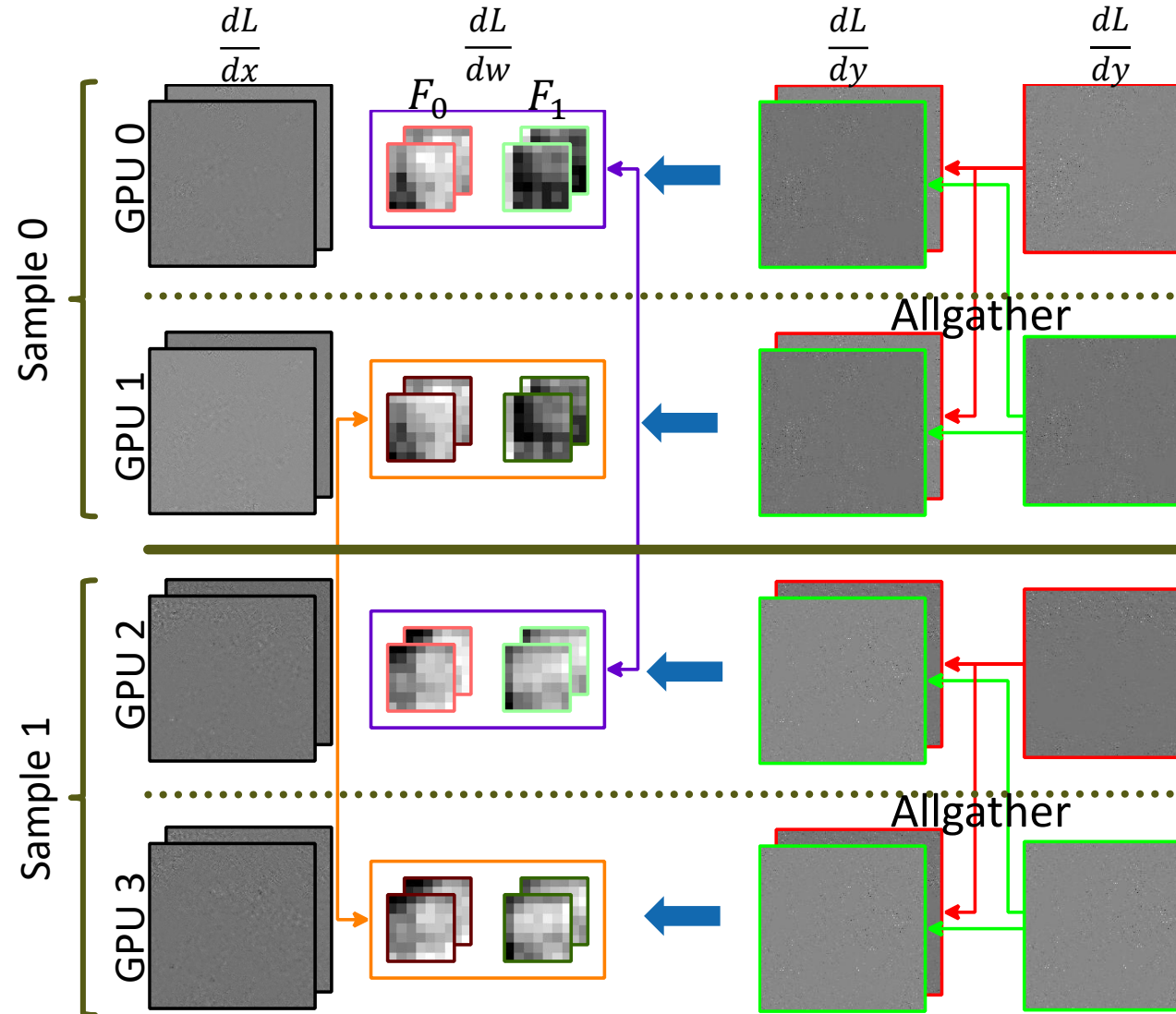
Stationary- x : Backward



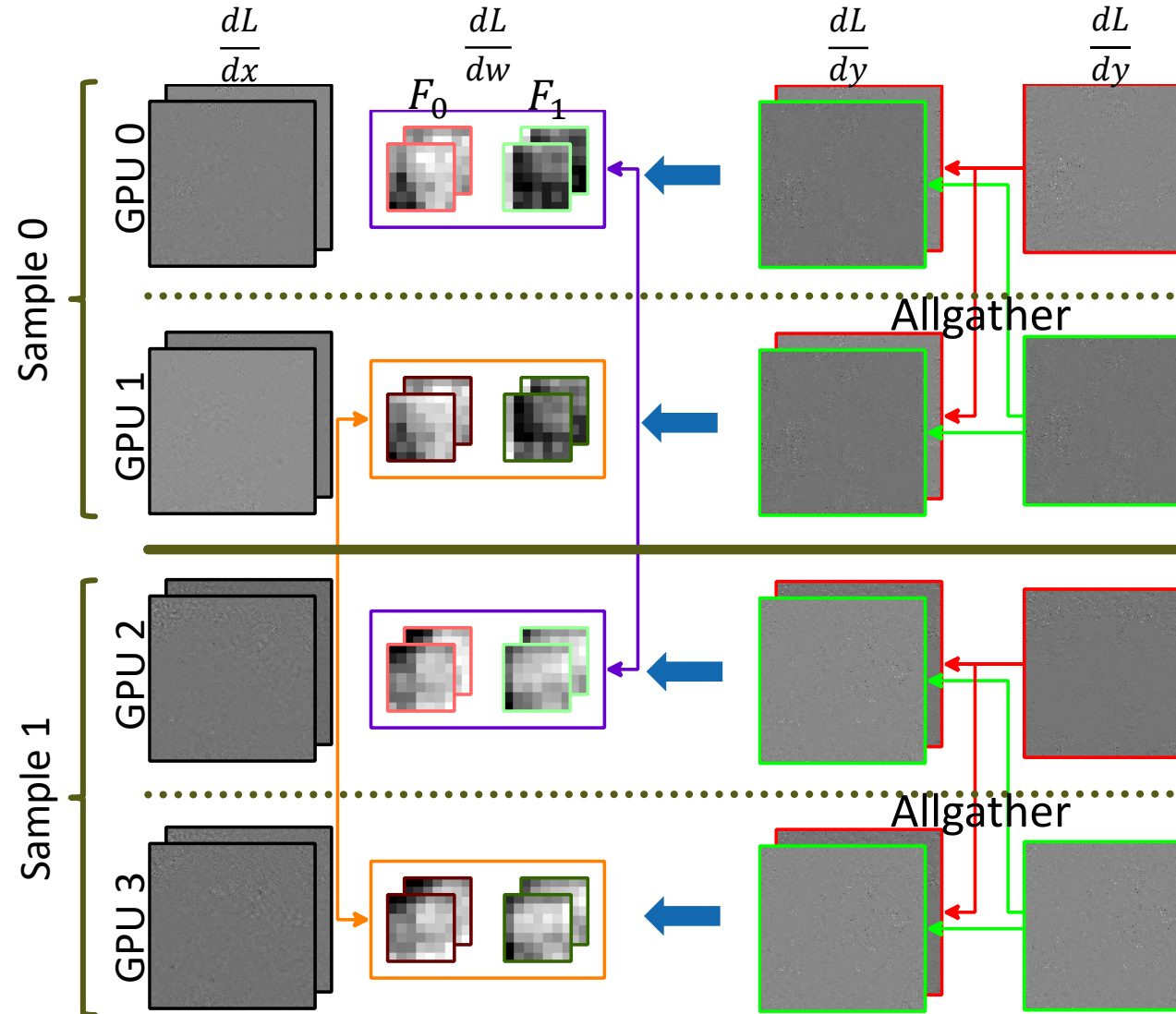
Stationary- x : Backward



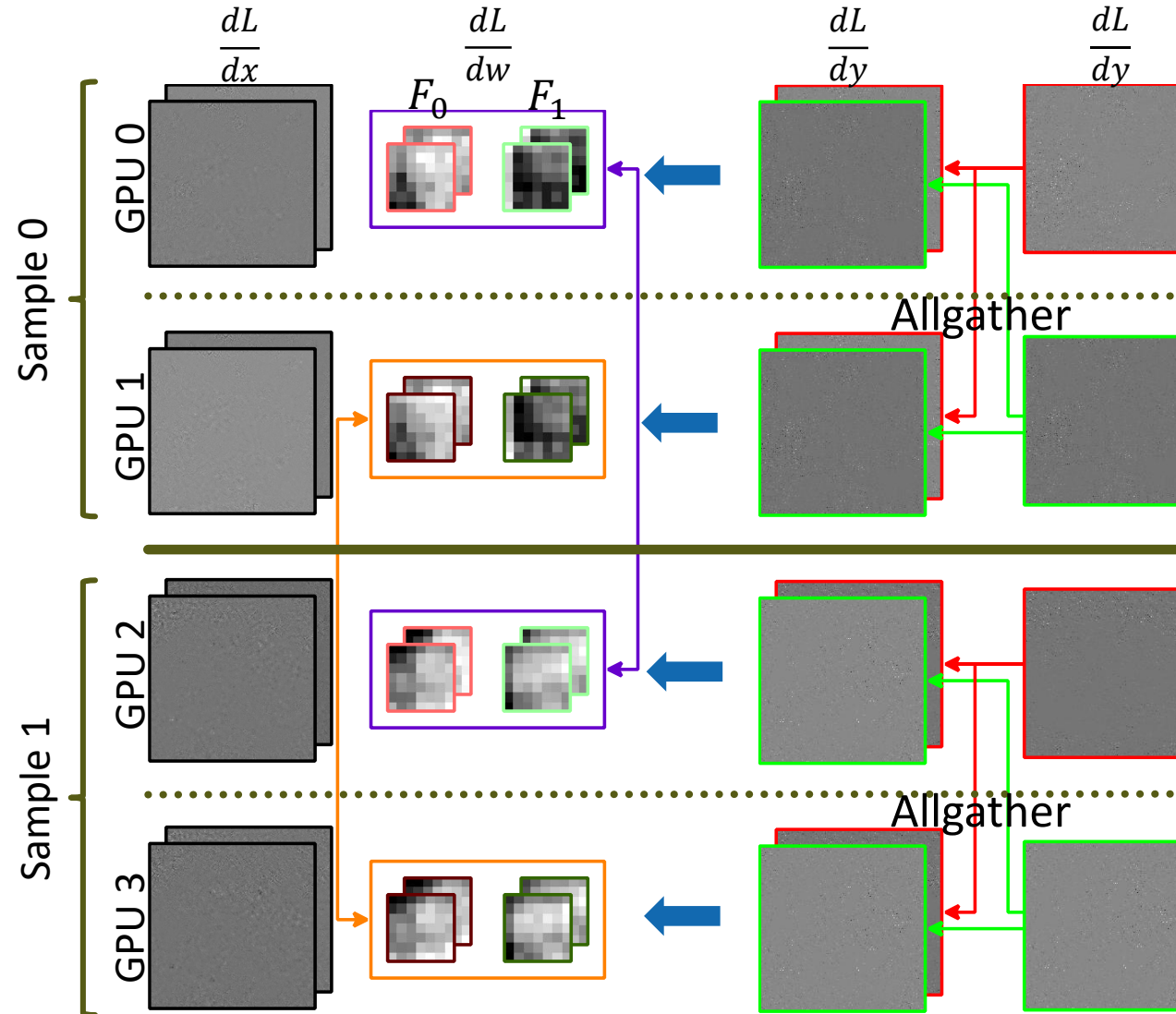
Stationary- x : Backward



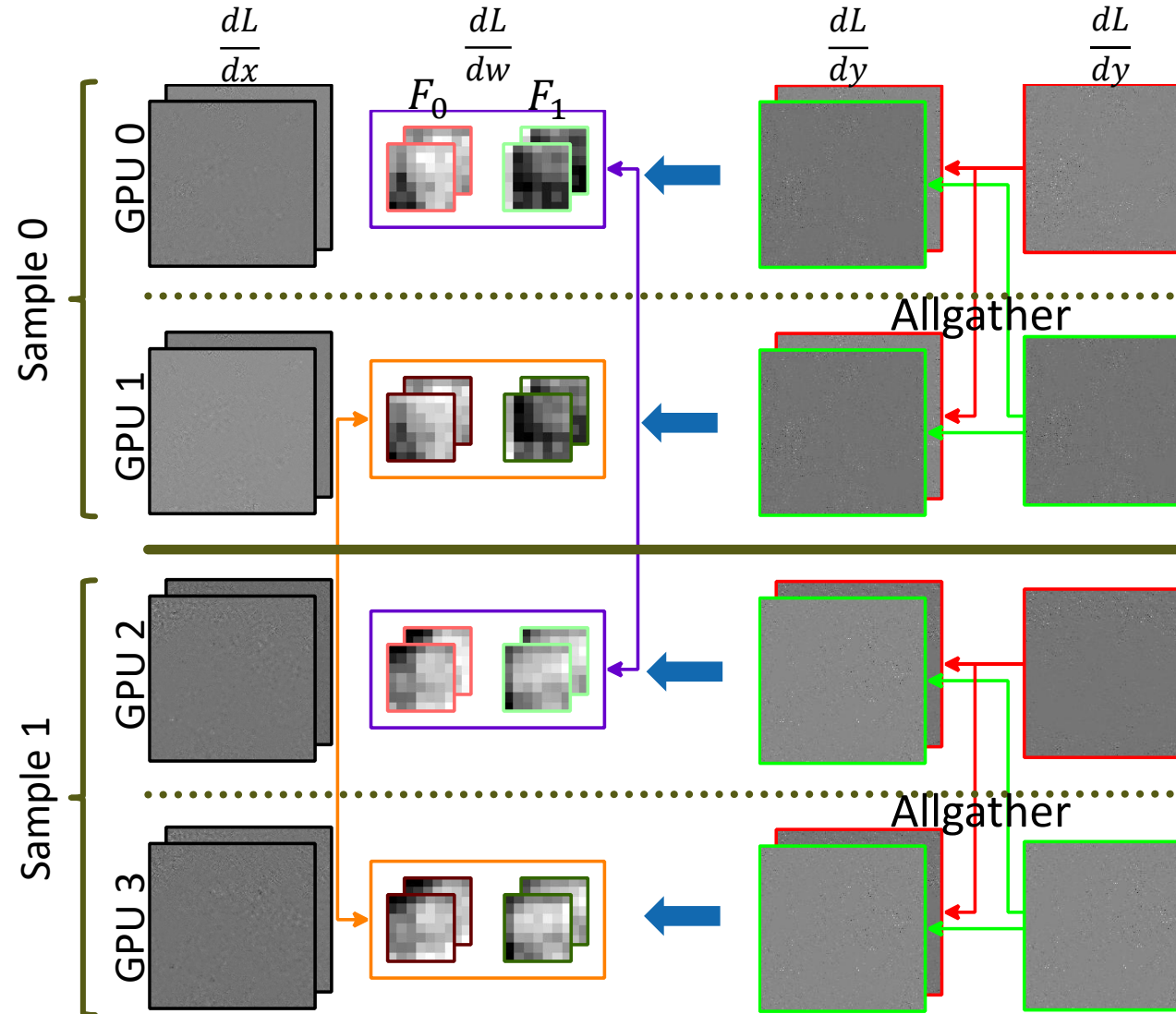
Stationary- x : Backward



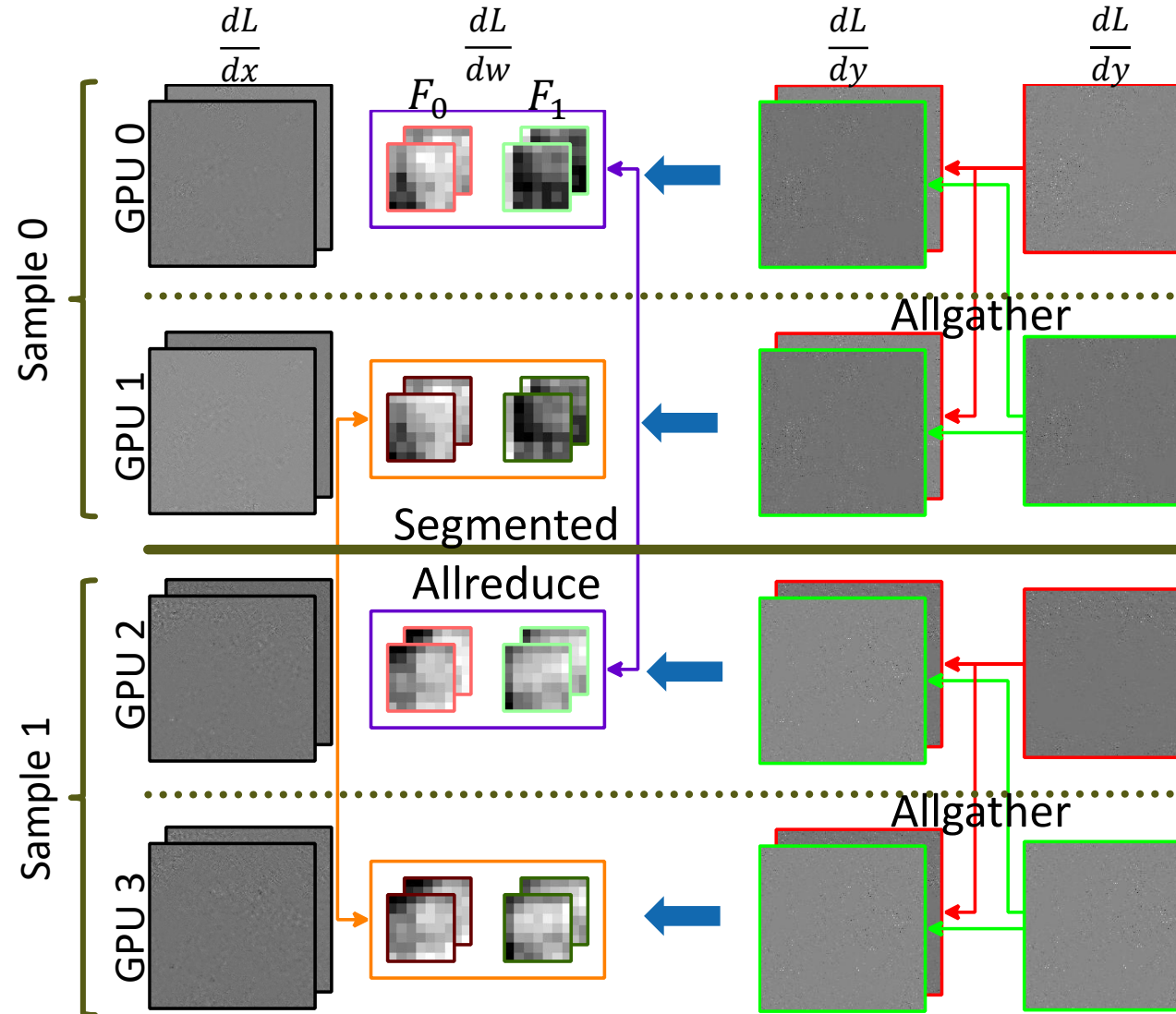
Stationary- x : Backward



Stationary- x : Backward

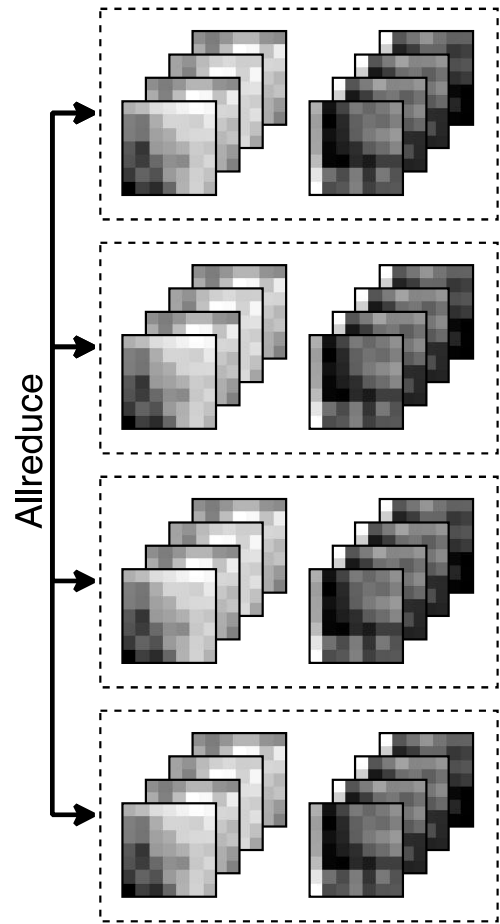


Stationary- x : Backward



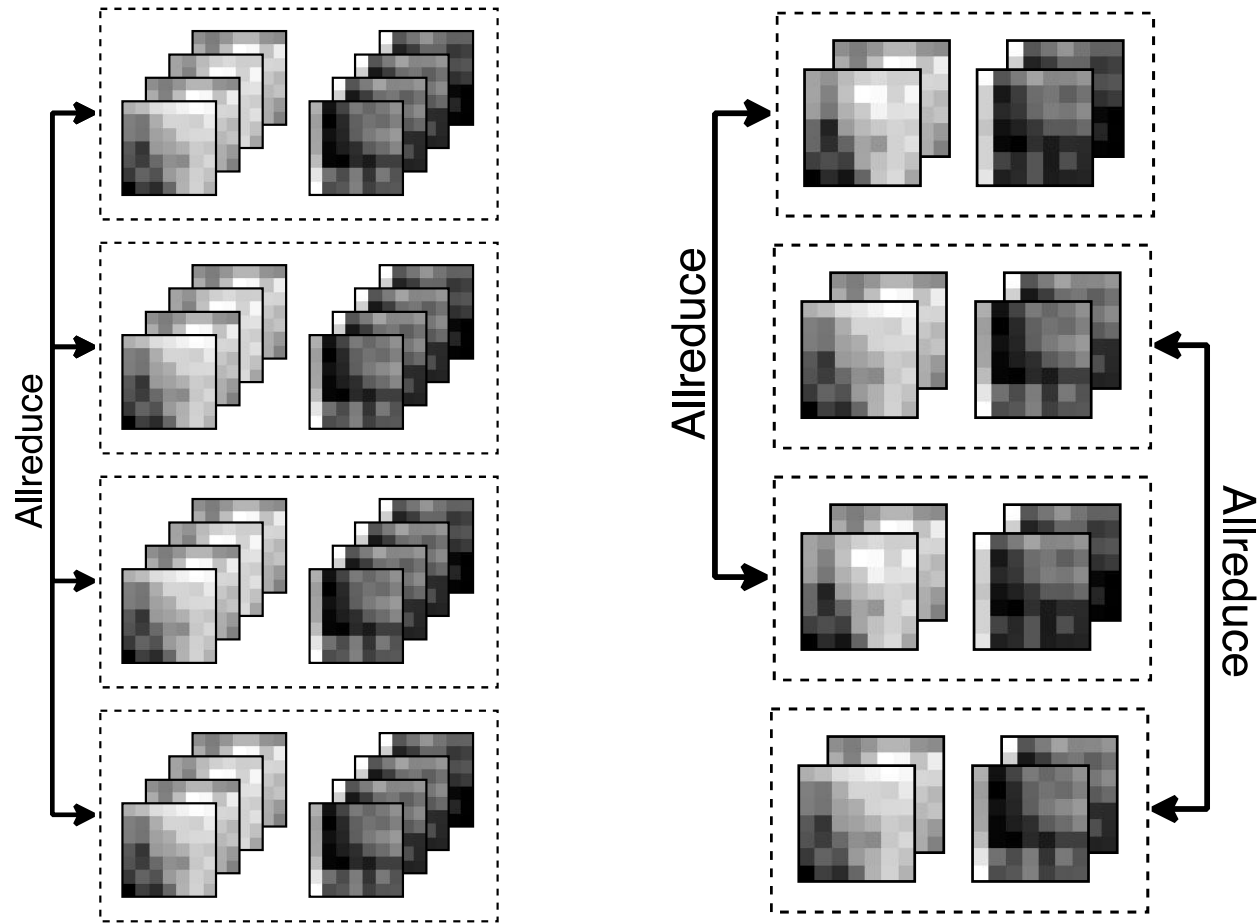
Segmented allreduce

Data Parallelism



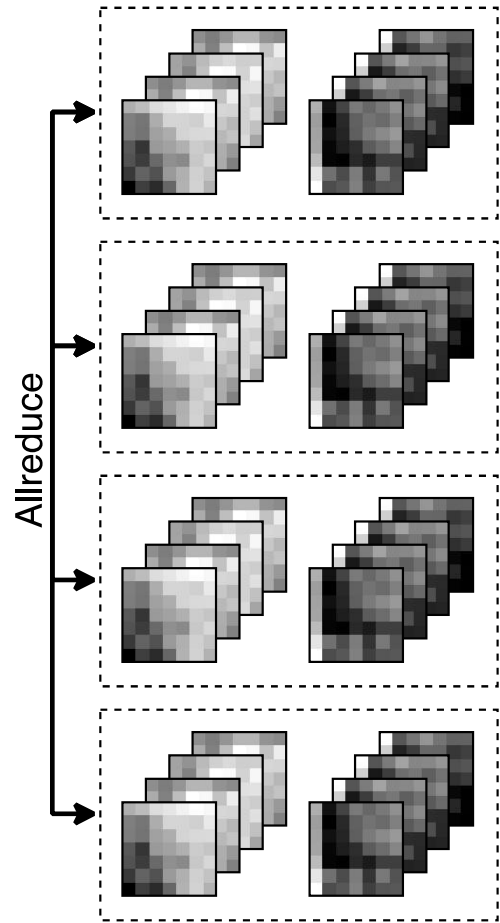
Segmented allreduce

Data Parallelism

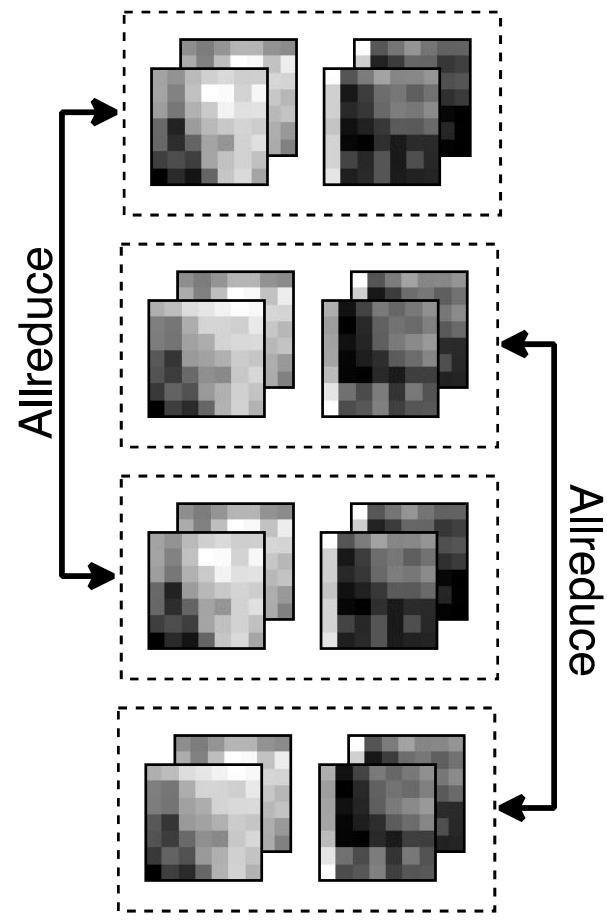


Segmented allreduce

Data Parallelism

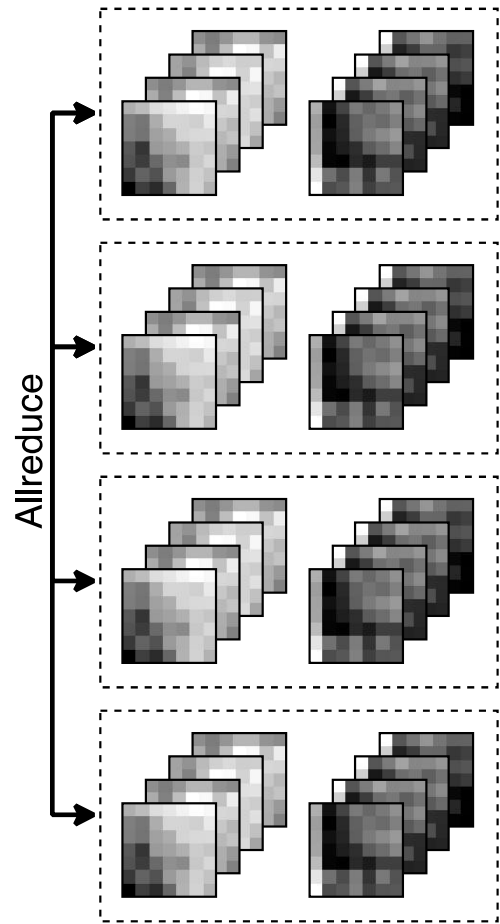


Segmented Allreduce

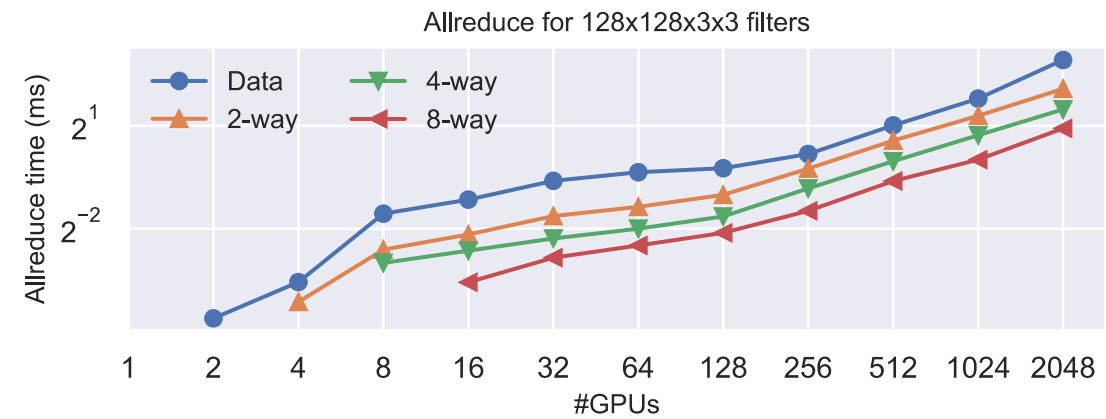
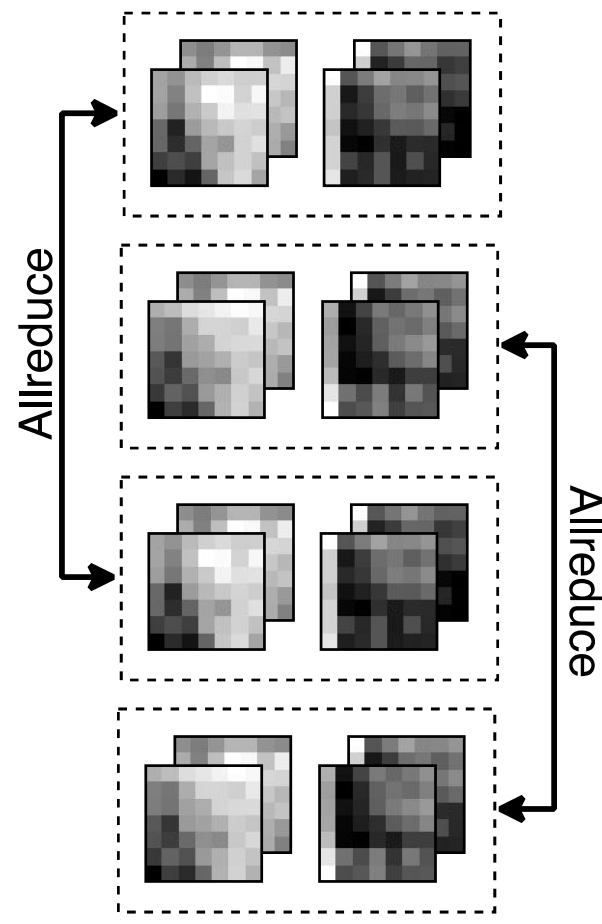


Segmented allreduce

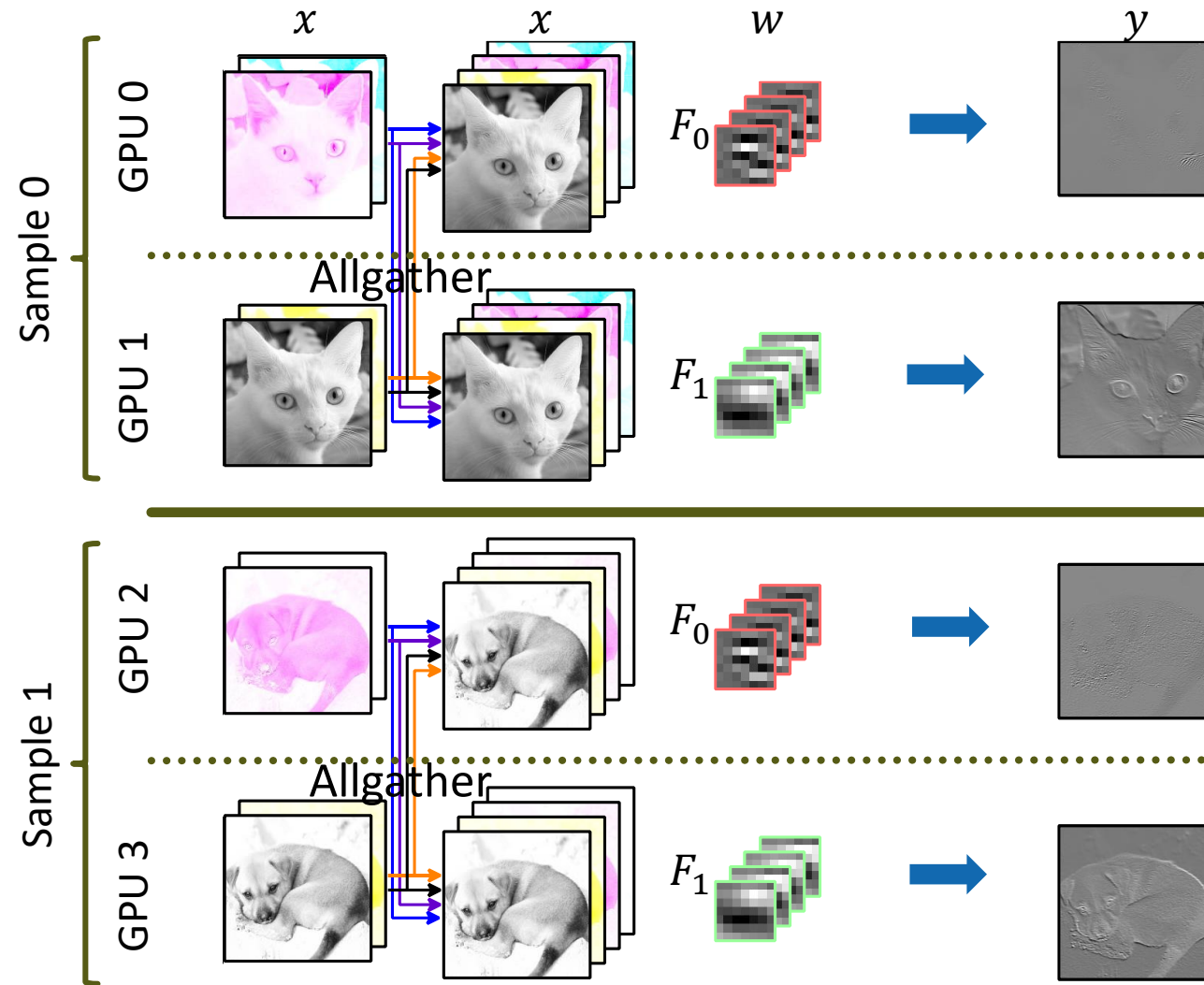
Data Parallelism



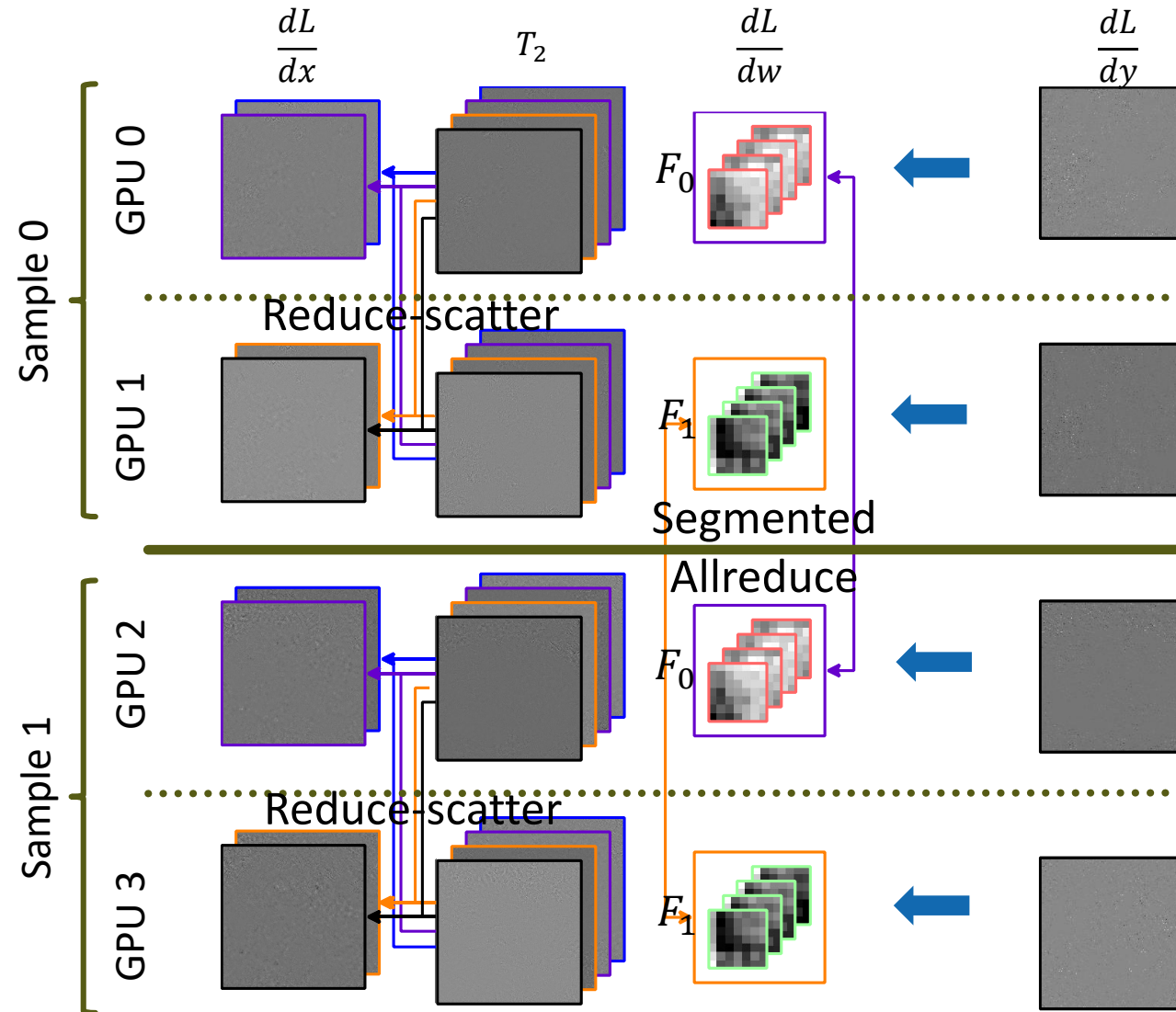
Segmented Allreduce



Stationary- y : Forward



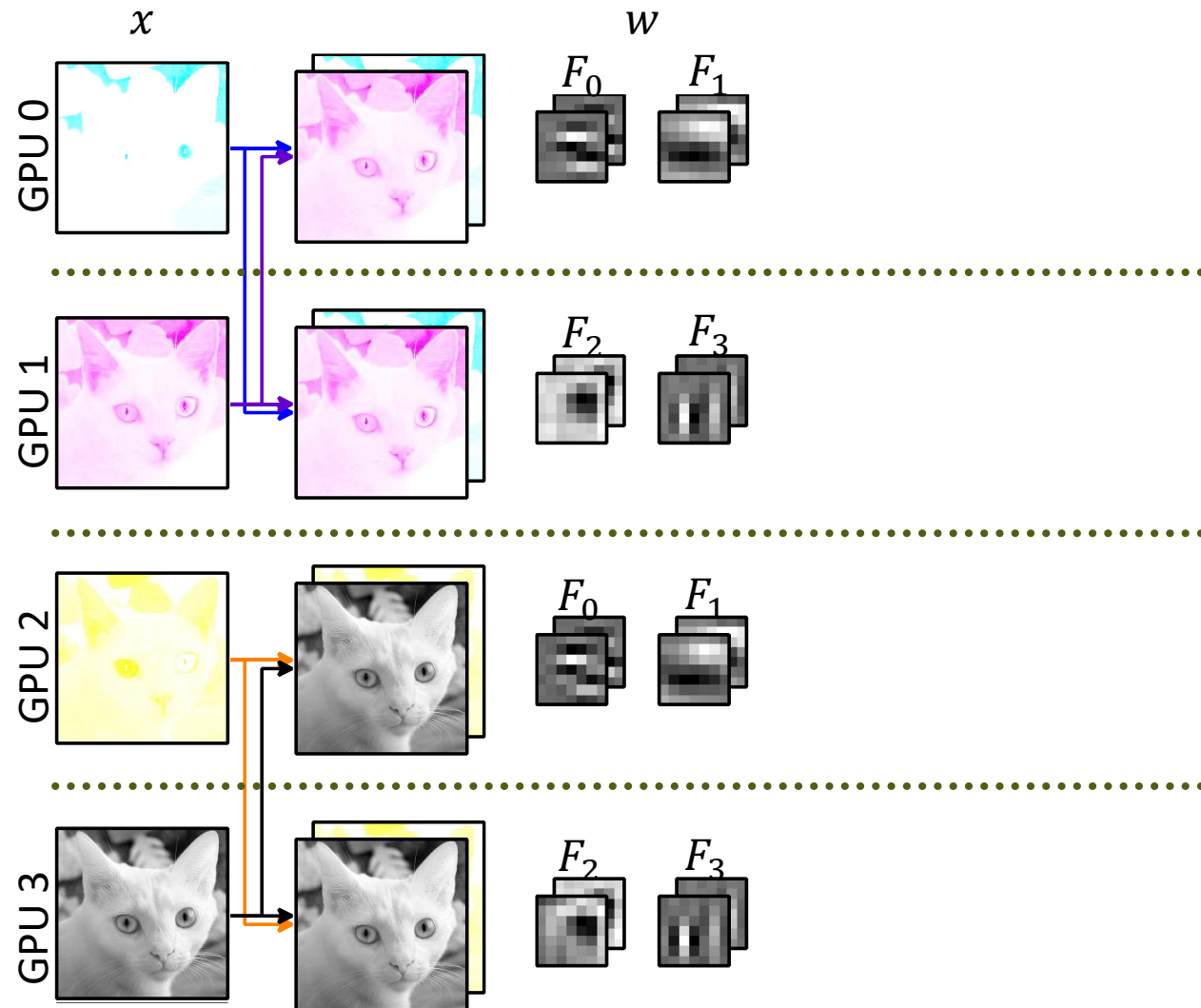
Stationary-y: Backward



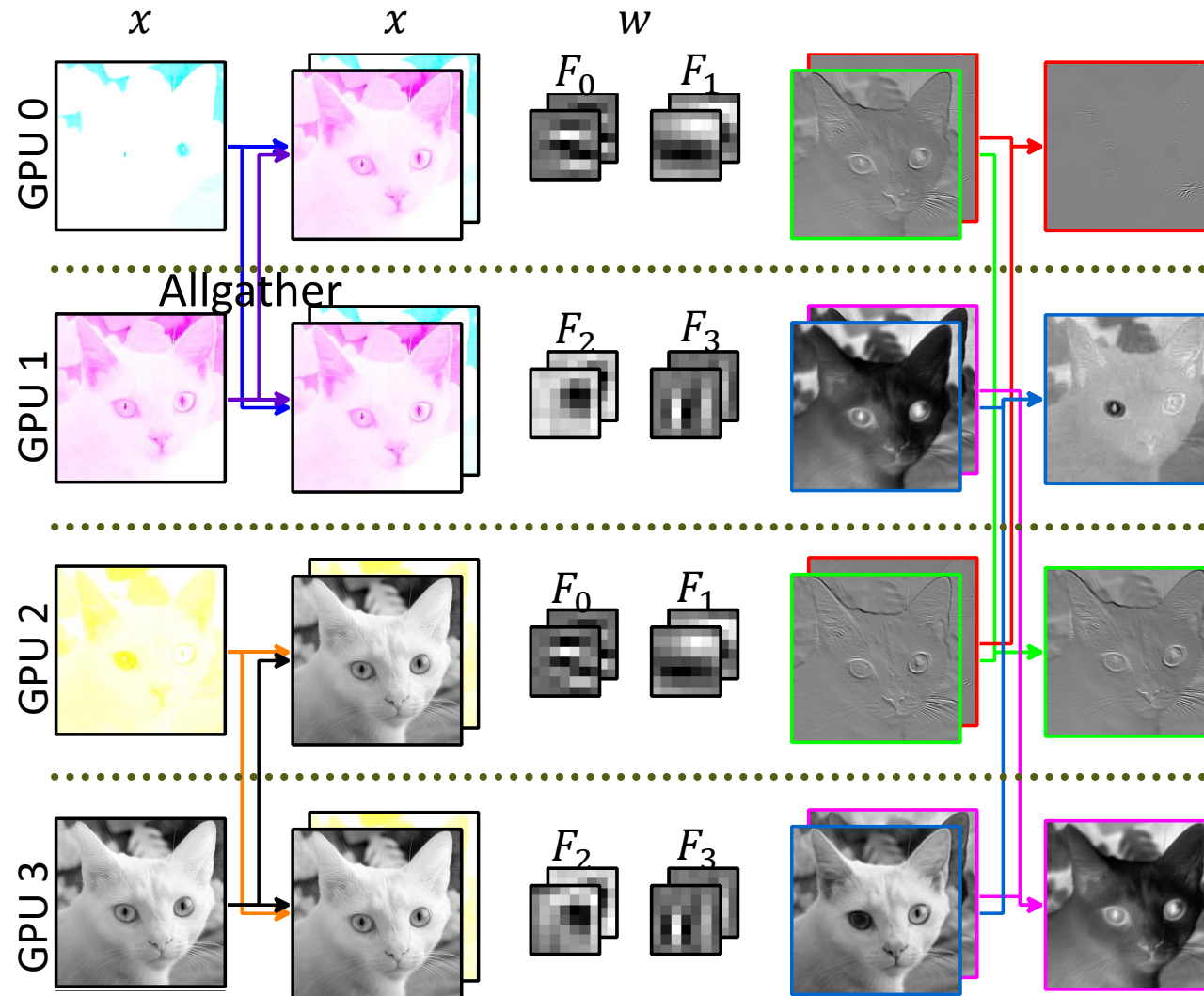
Stationary-w: Forward



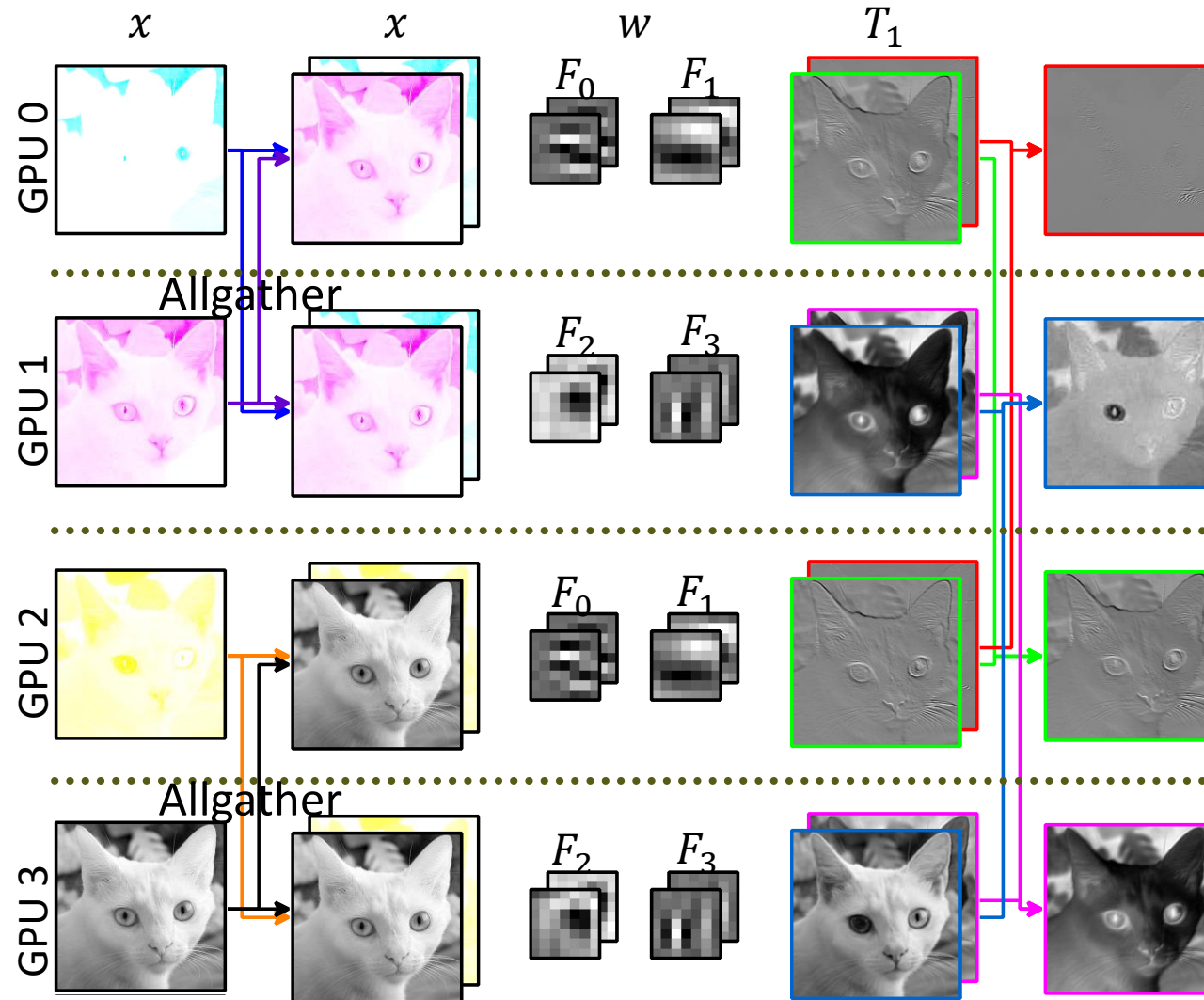
Stationary-w: Forward



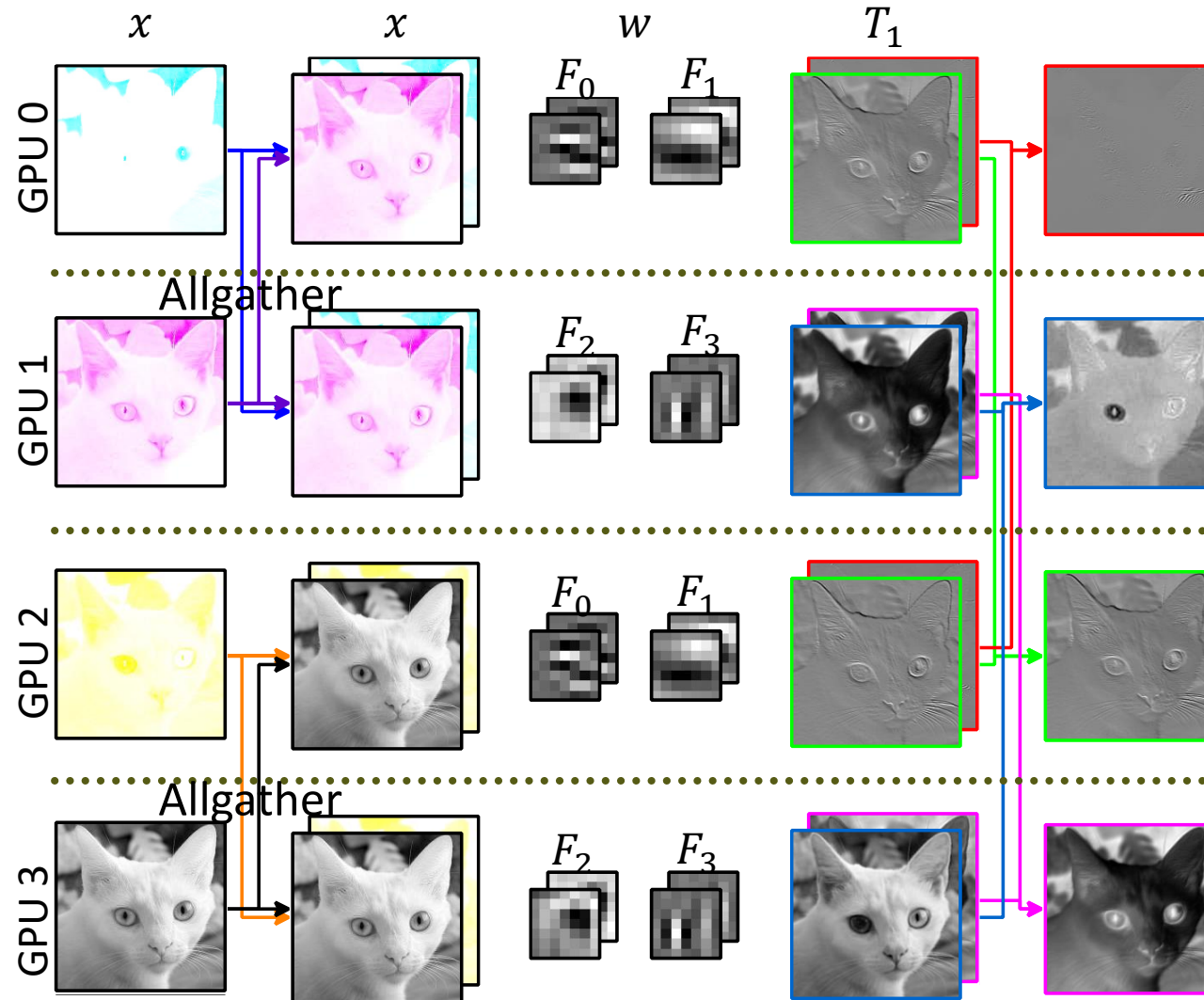
Stationary-w: Forward



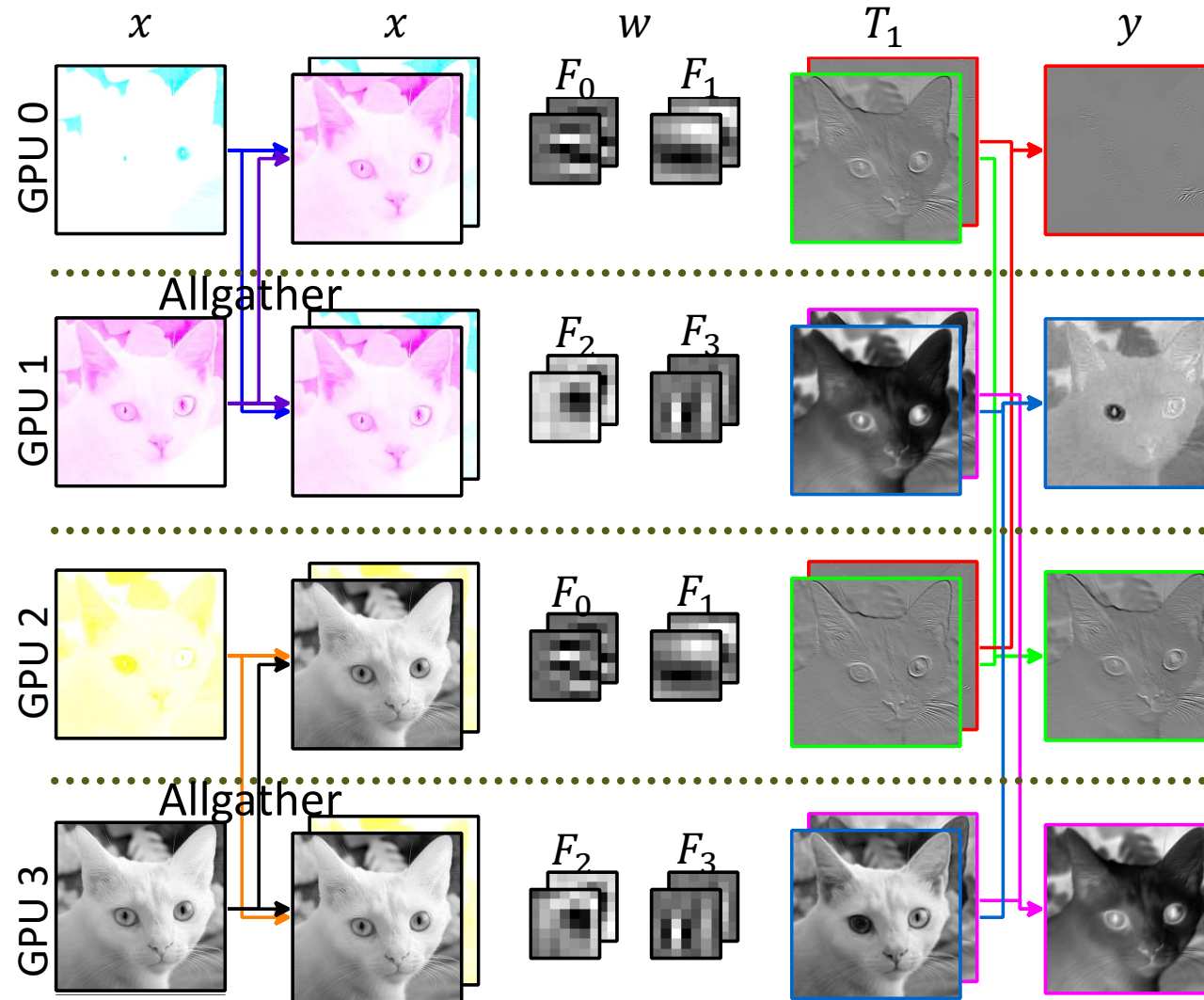
Stationary-w: Forward



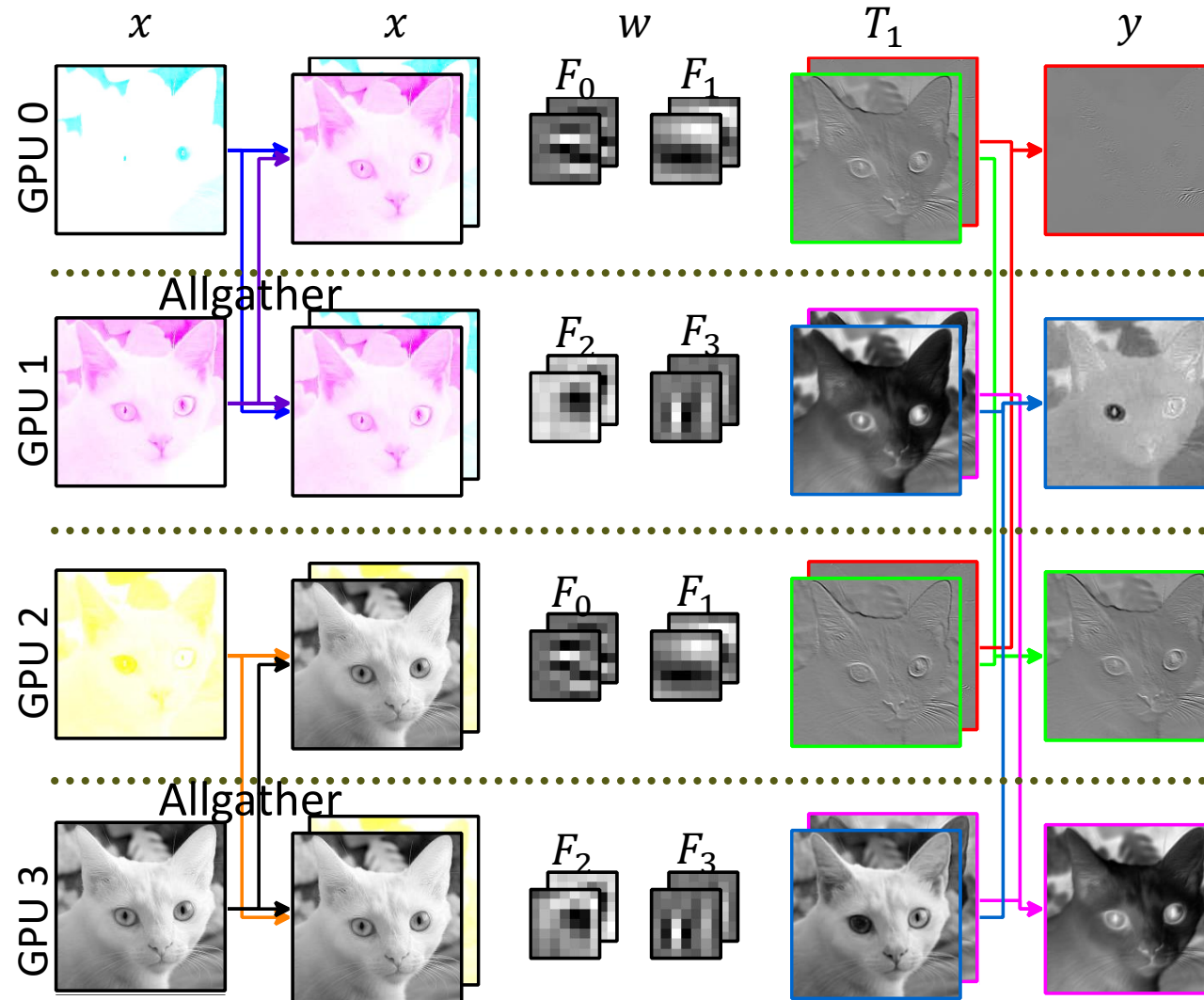
Stationary-w: Forward



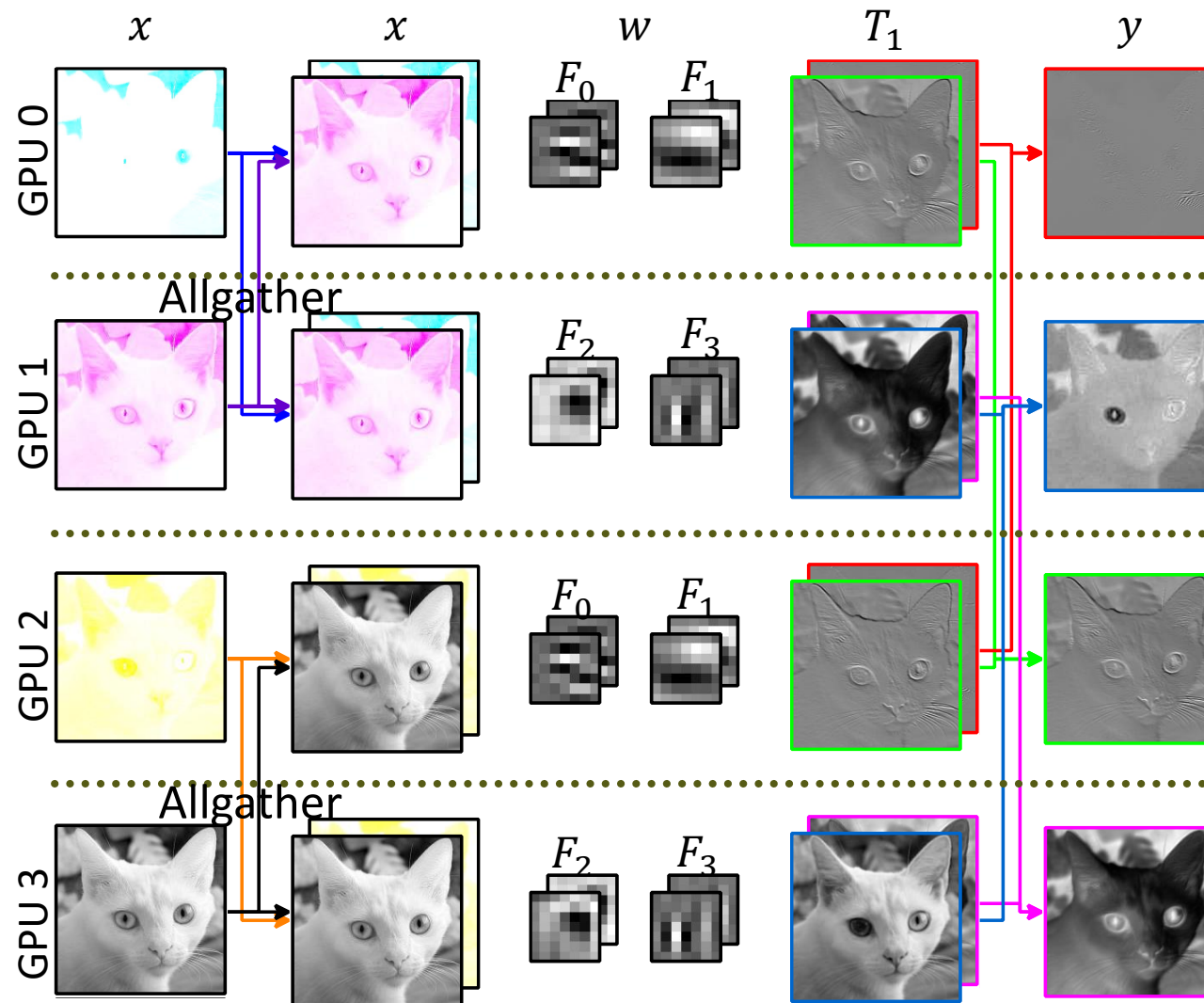
Stationary-w: Forward



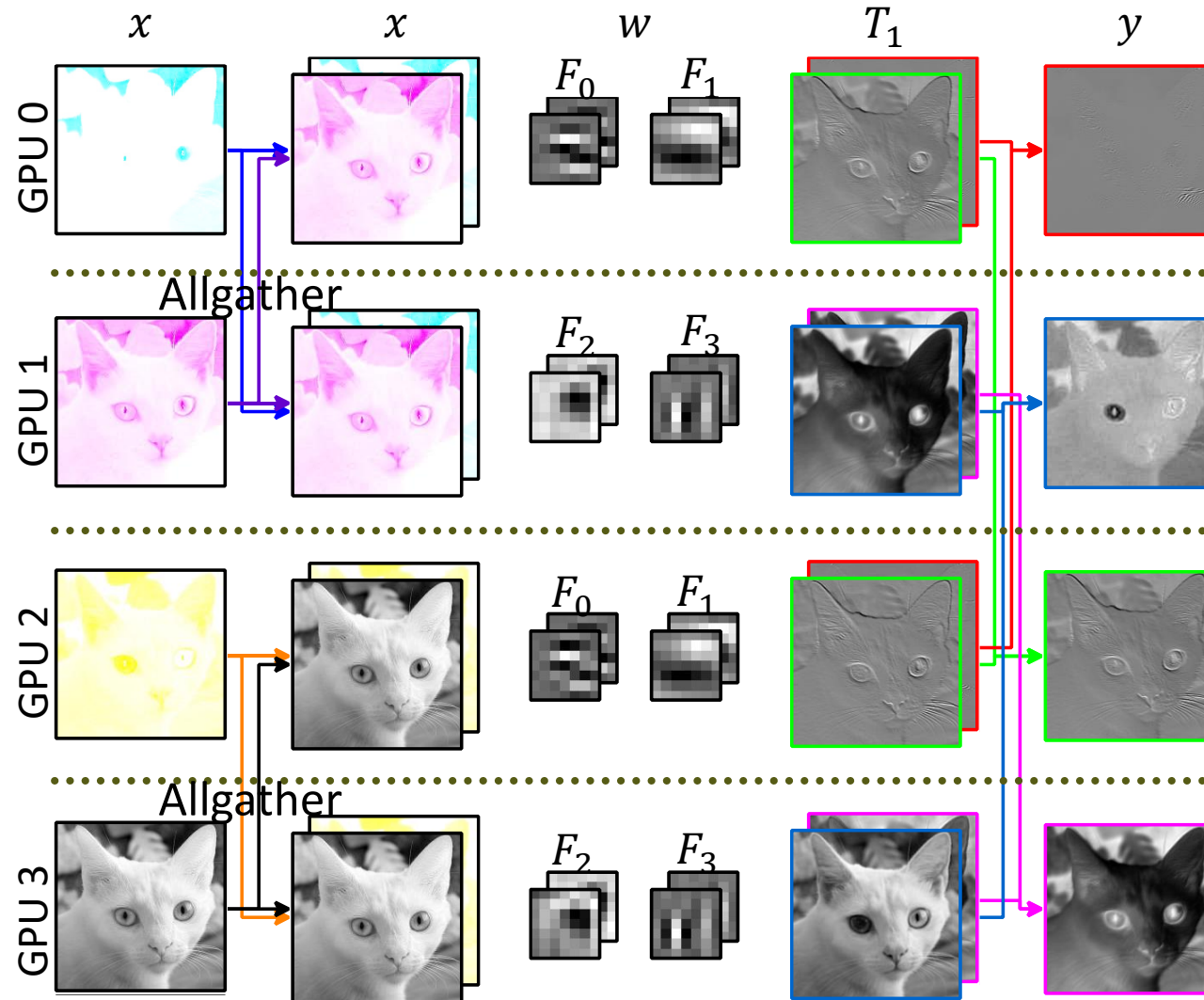
Stationary-w: Forward



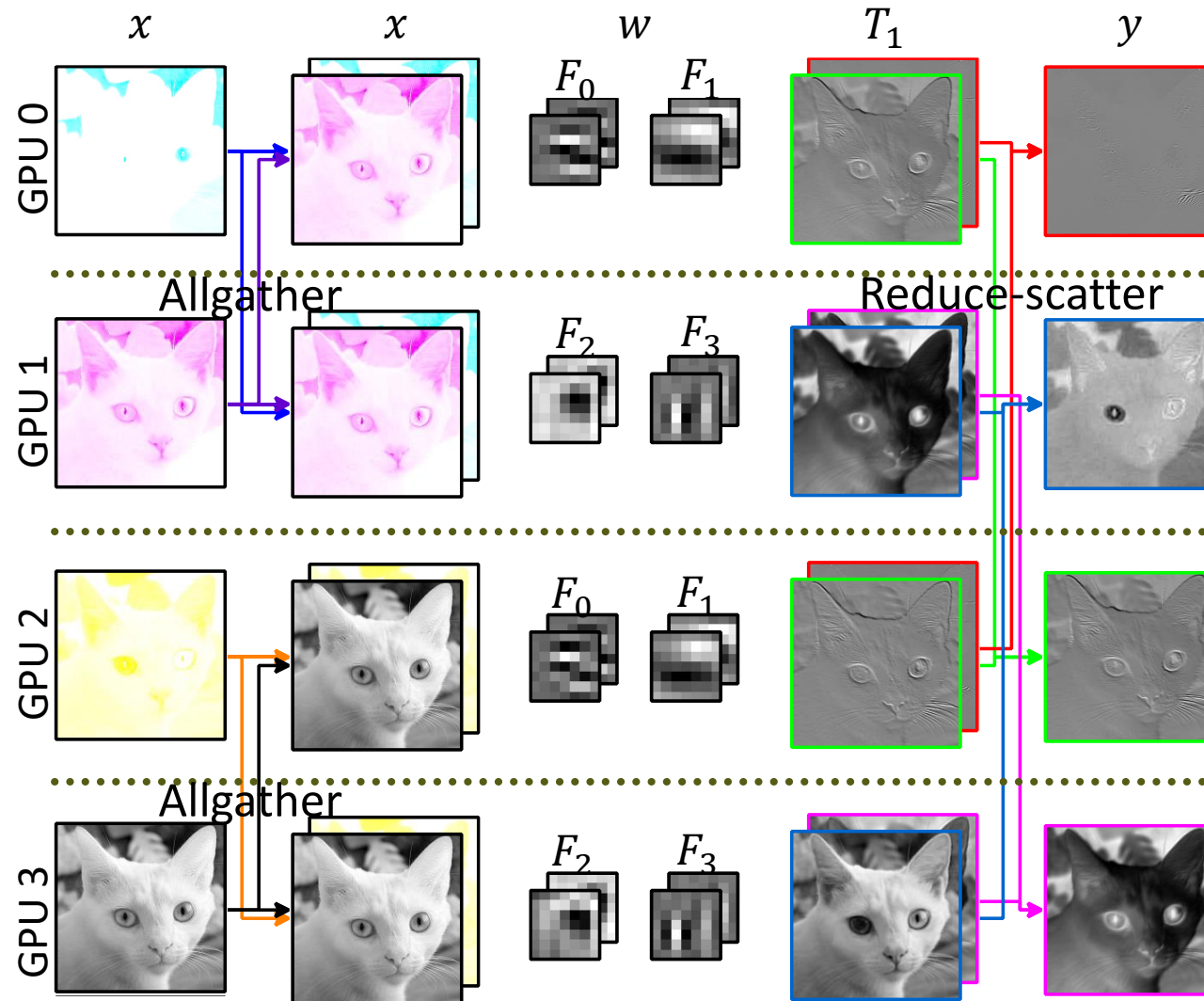
Stationary-w: Forward



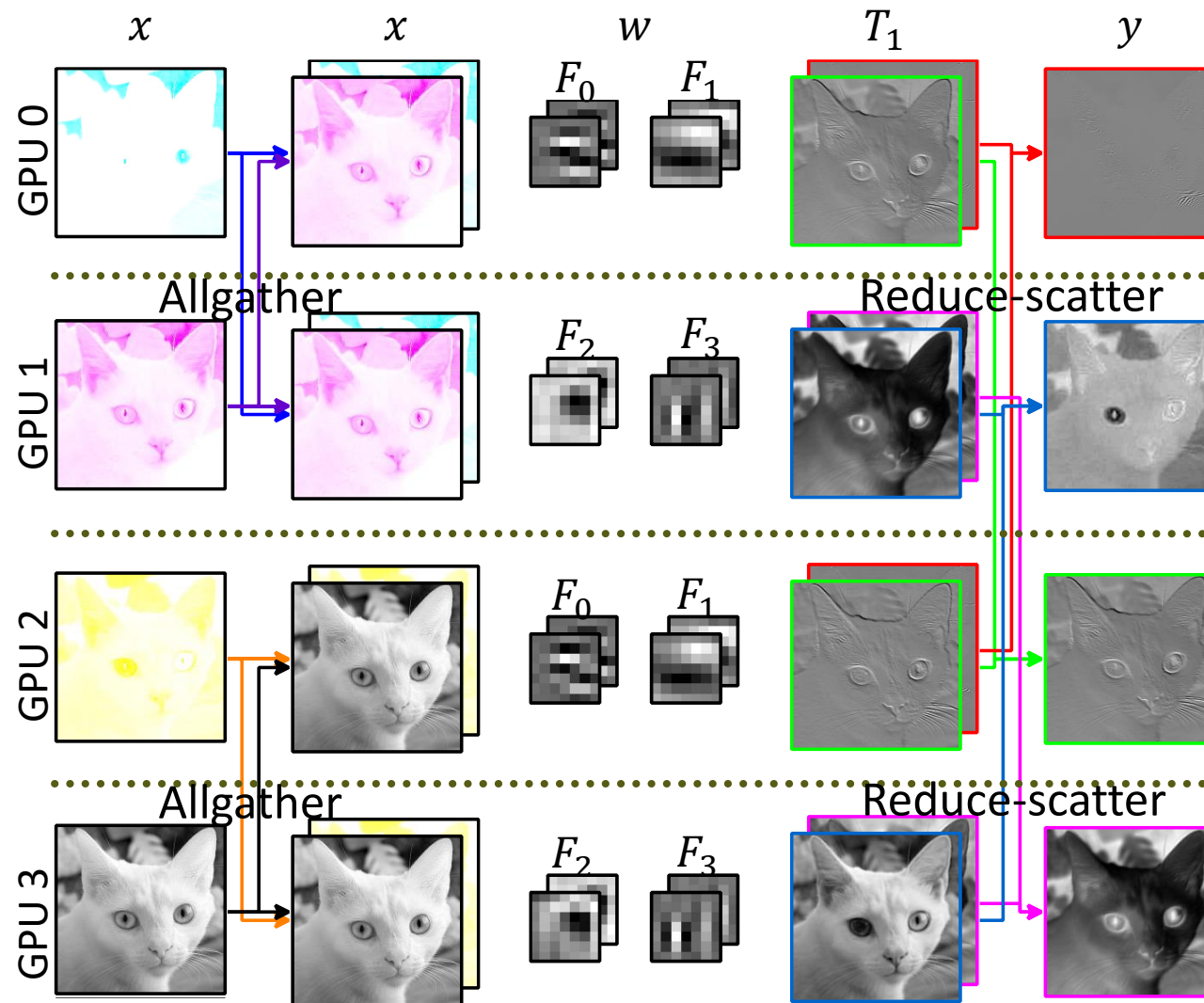
Stationary-w: Forward



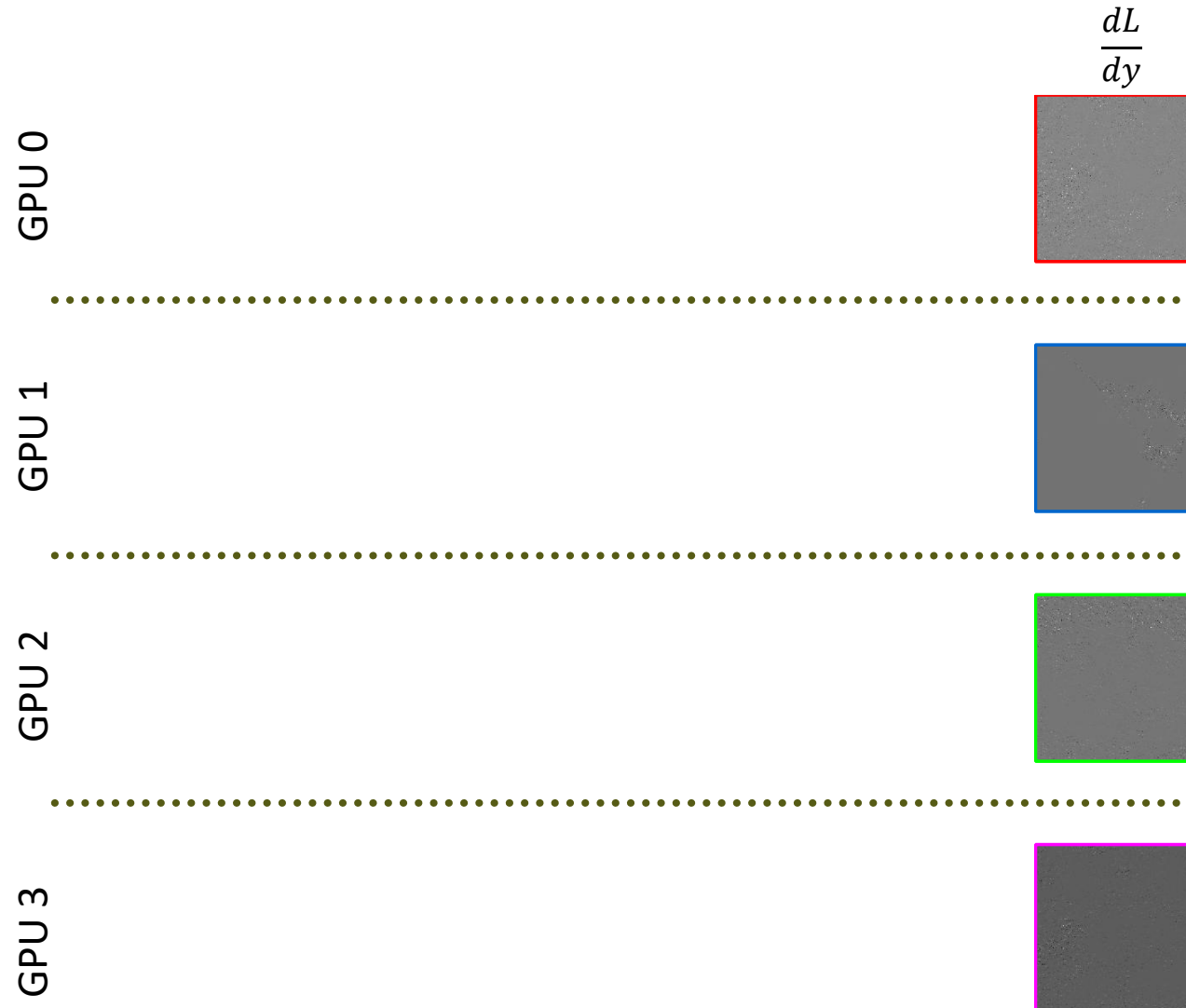
Stationary-w: Forward



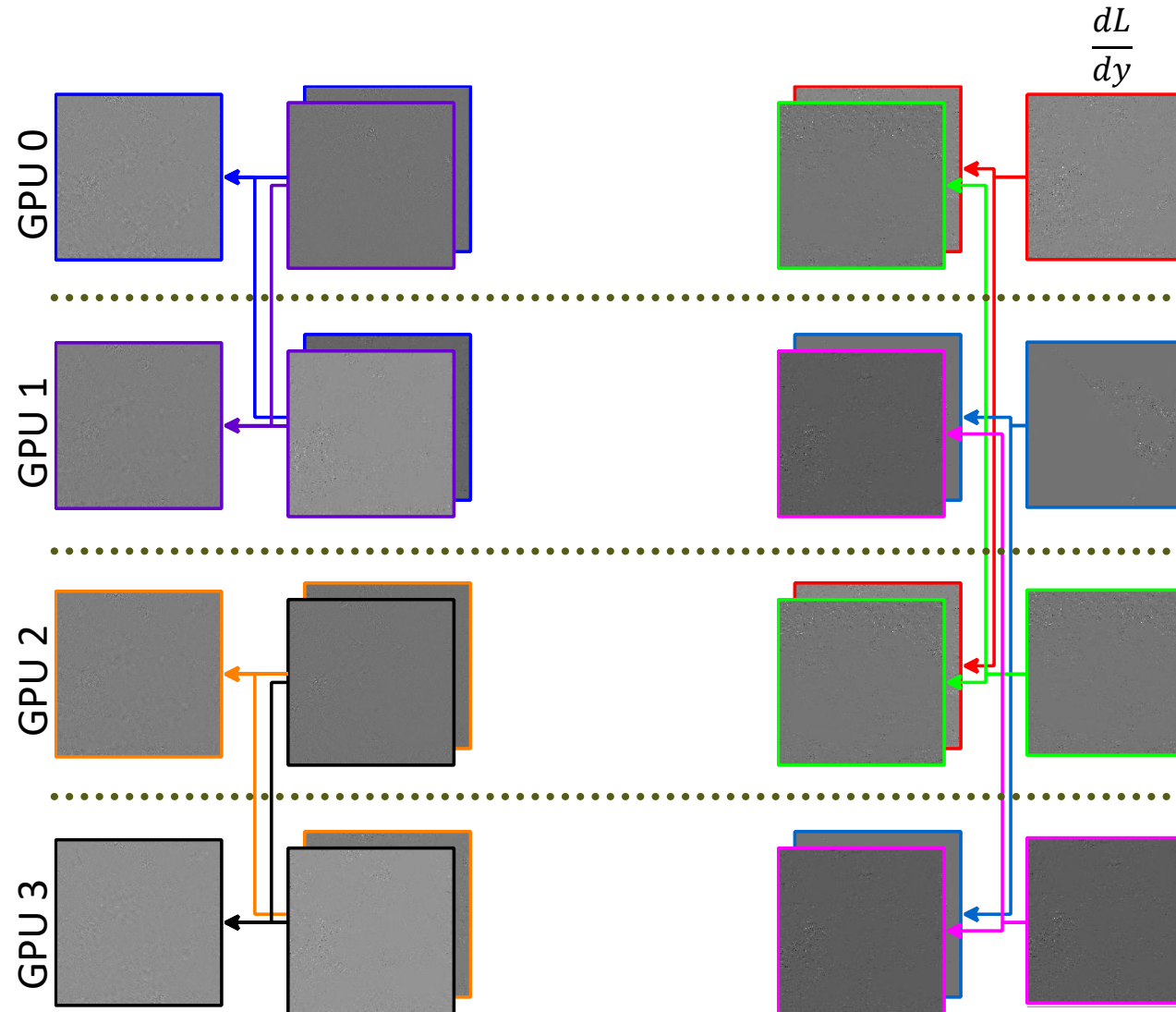
Stationary-w: Forward



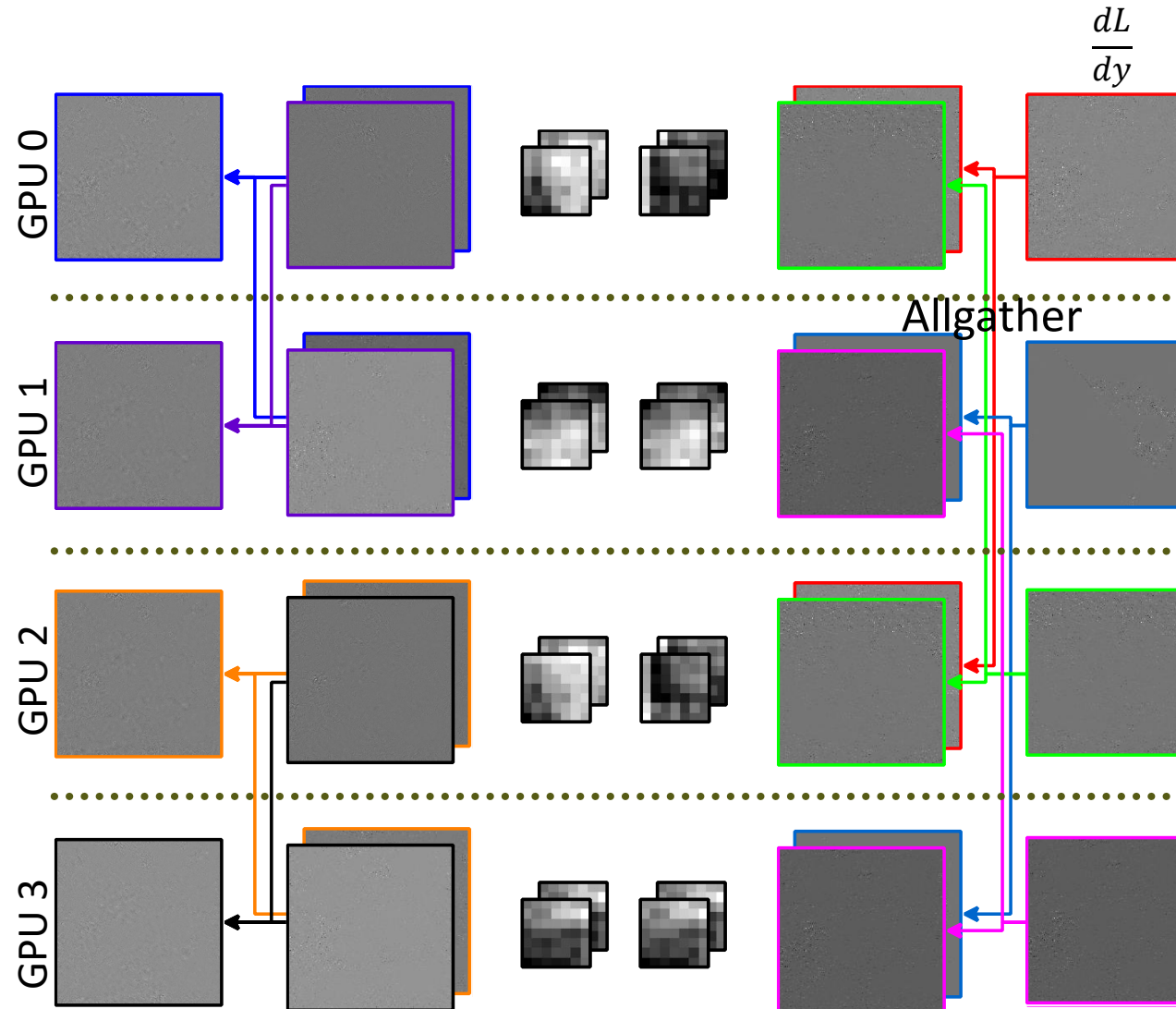
Stationary-w: Backward



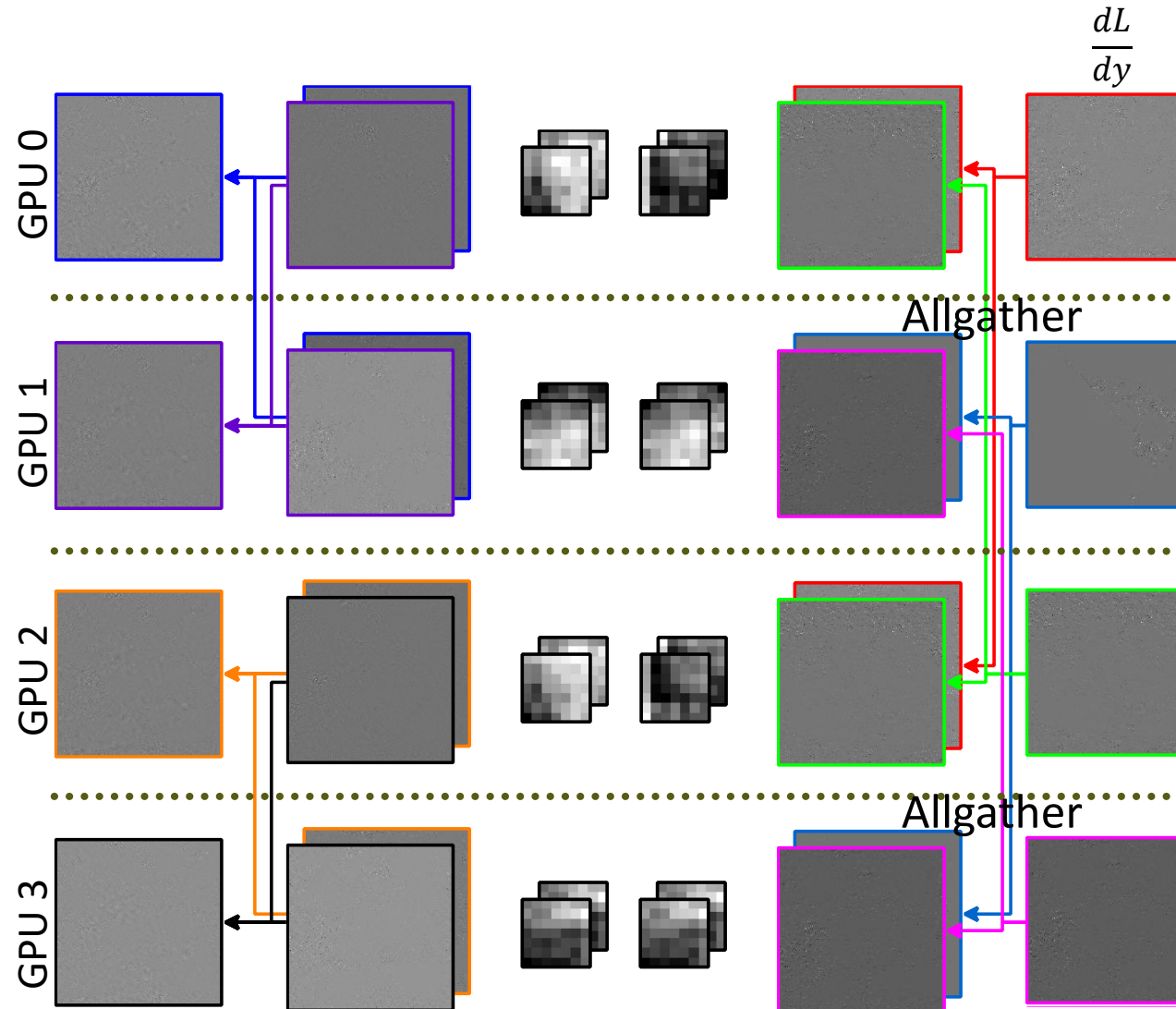
Stationary-w: Backward



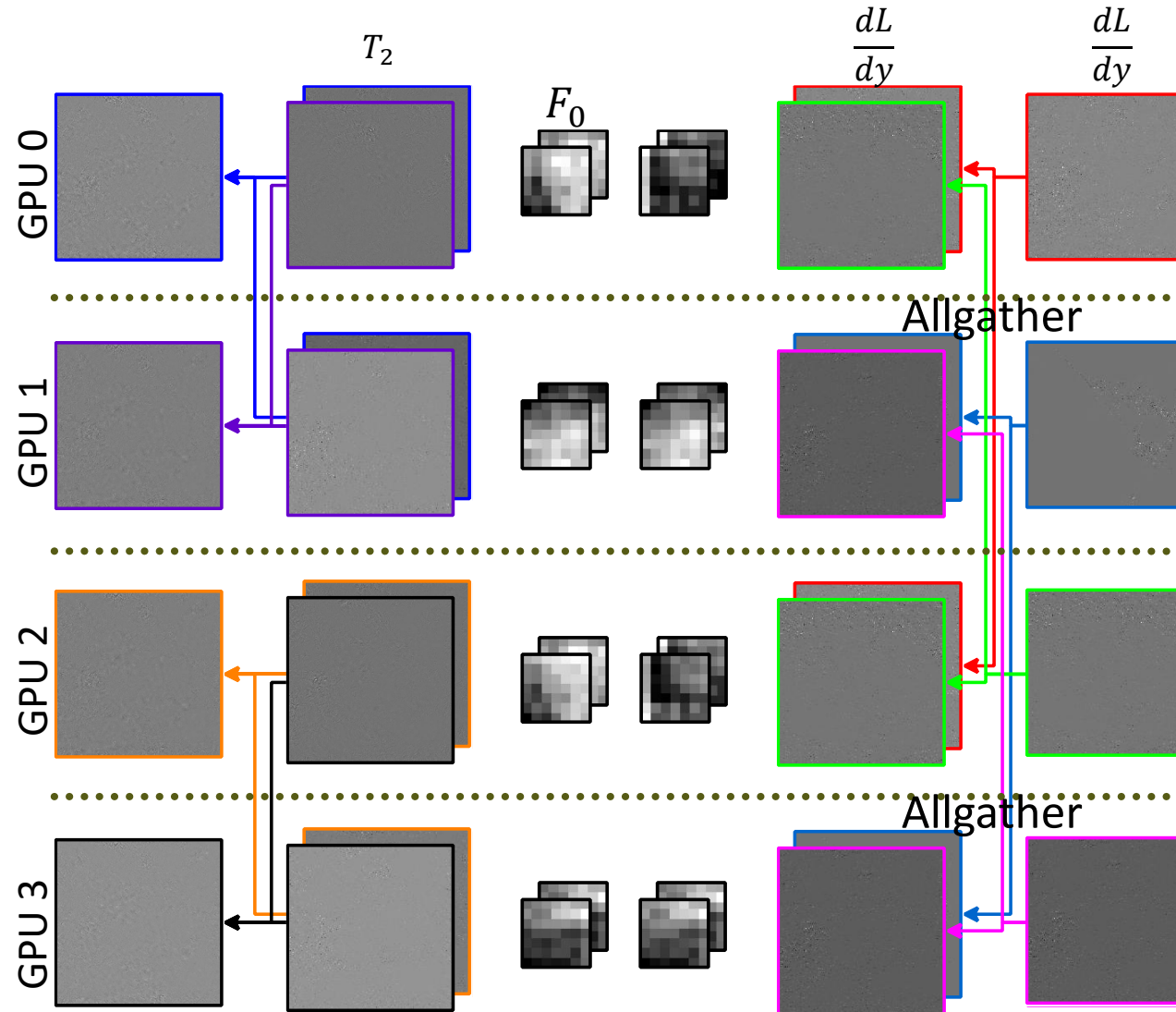
Stationary-w: Backward



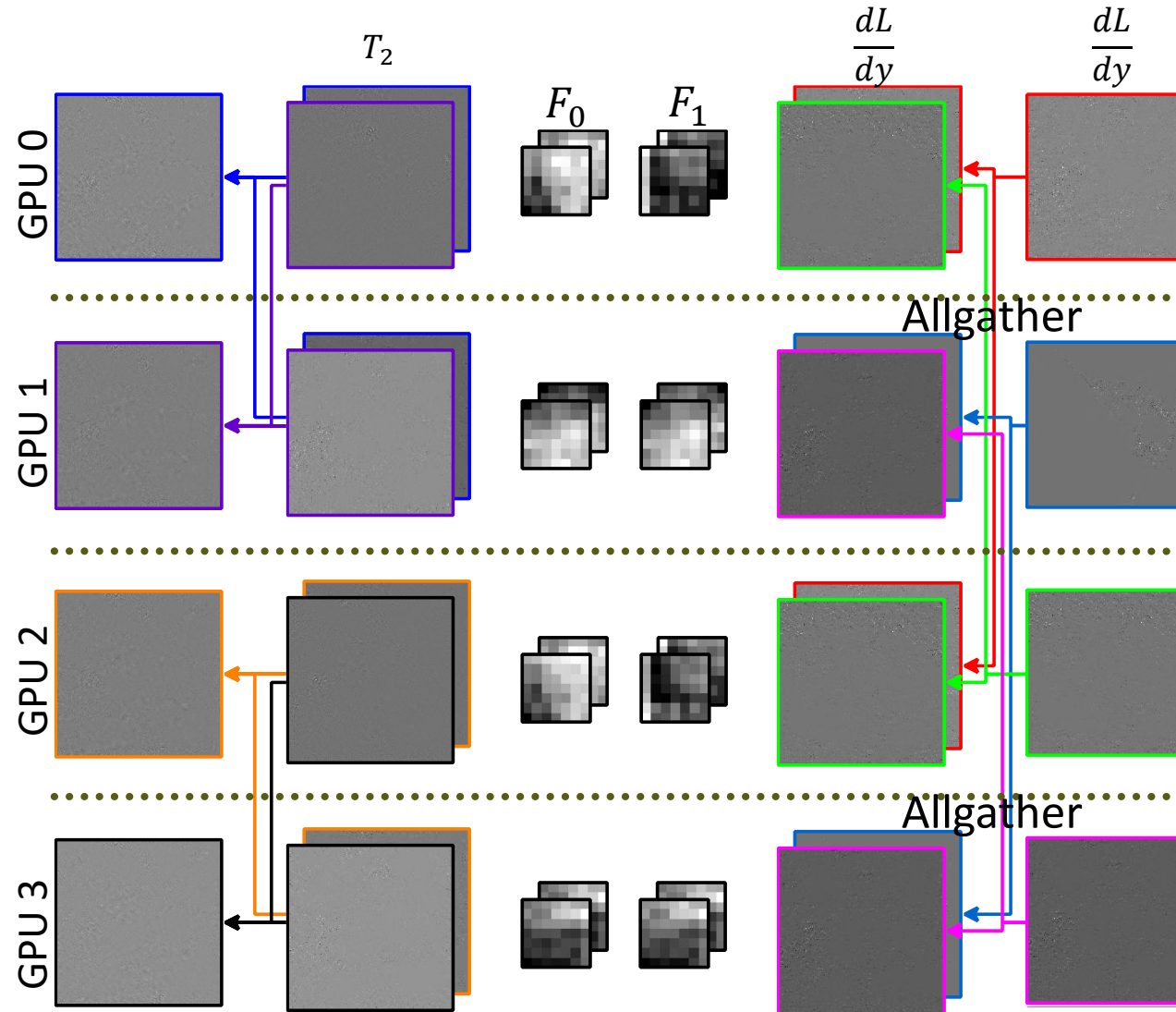
Stationary-w: Backward



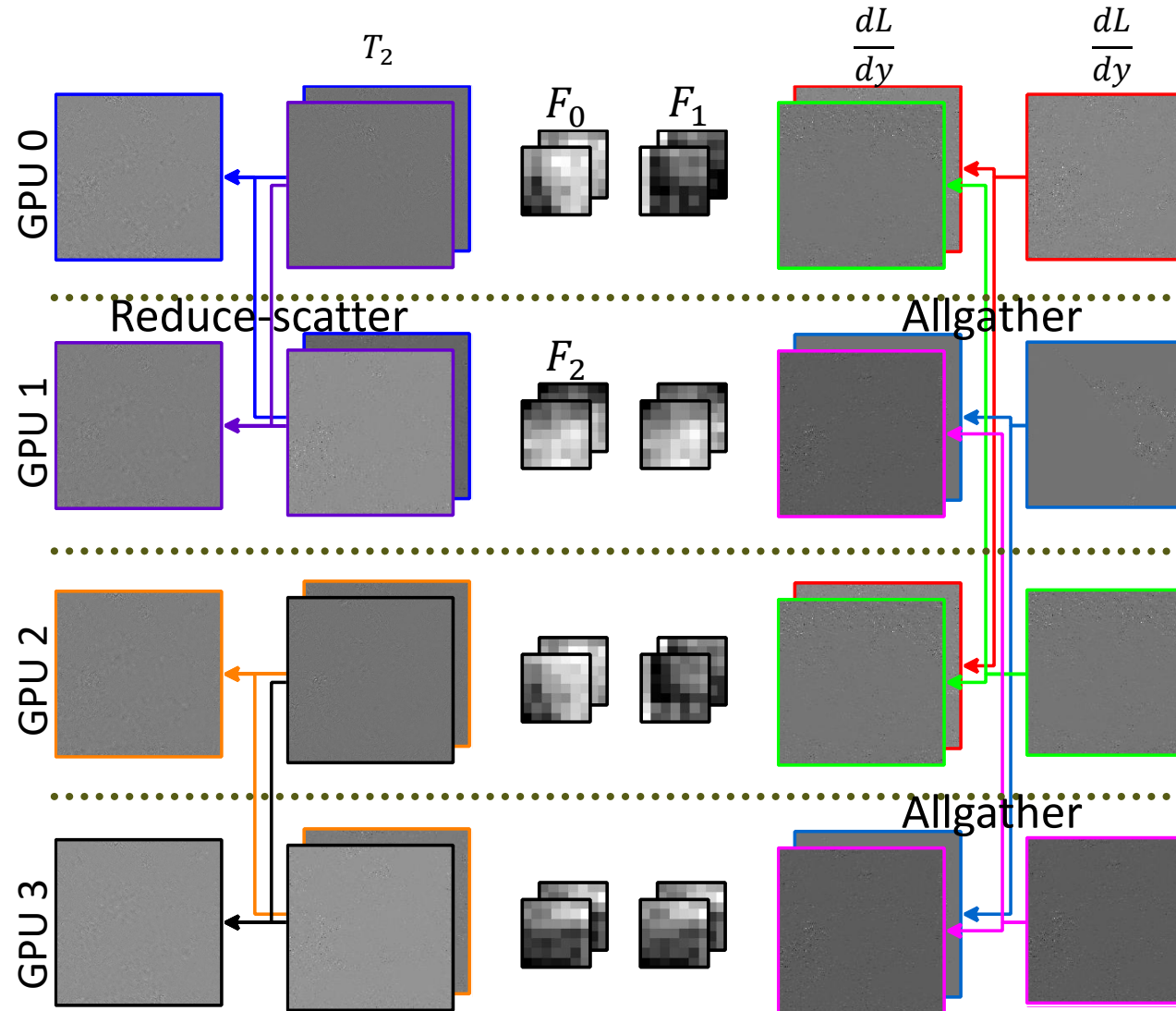
Stationary-w: Backward



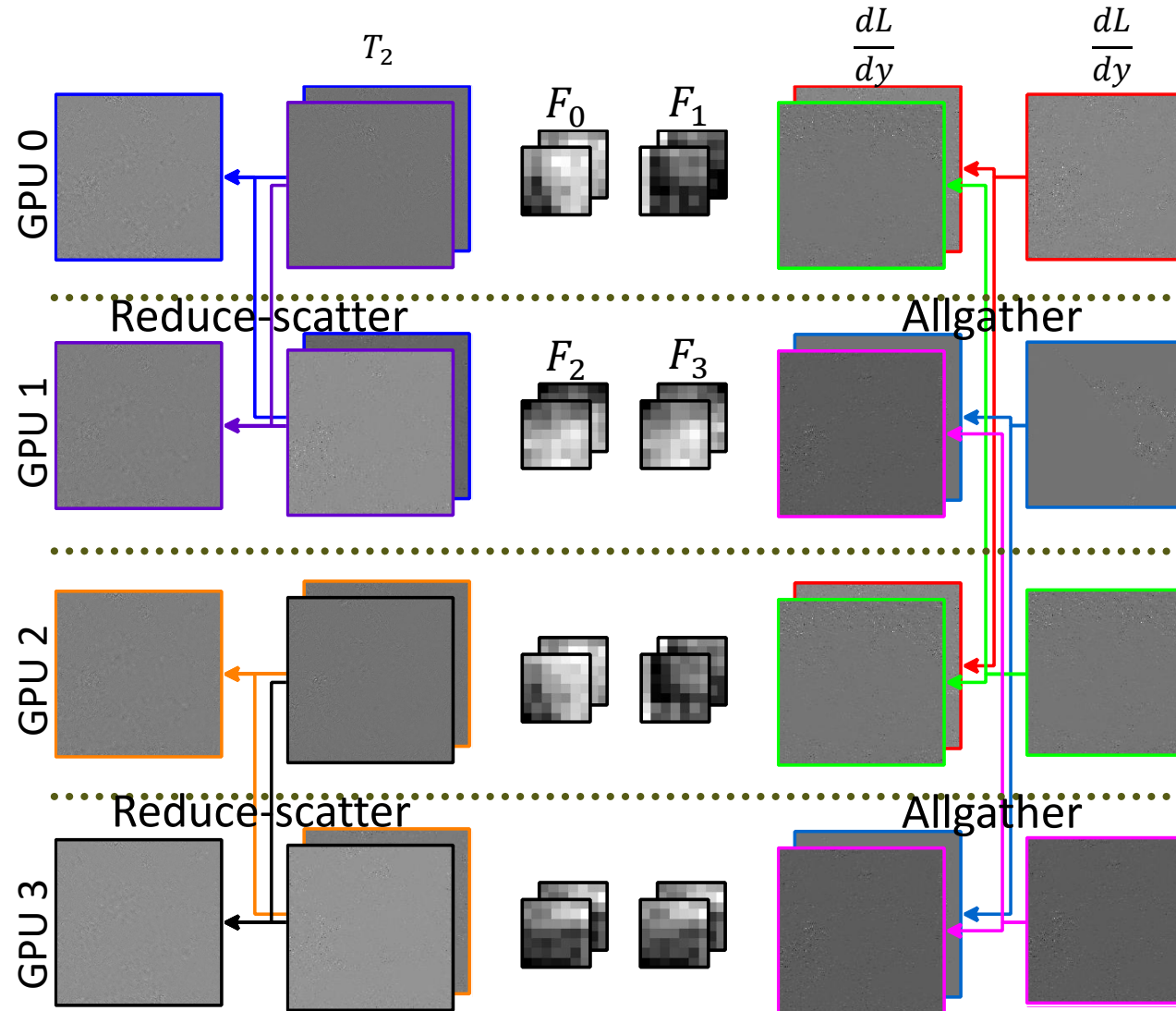
Stationary-w: Backward



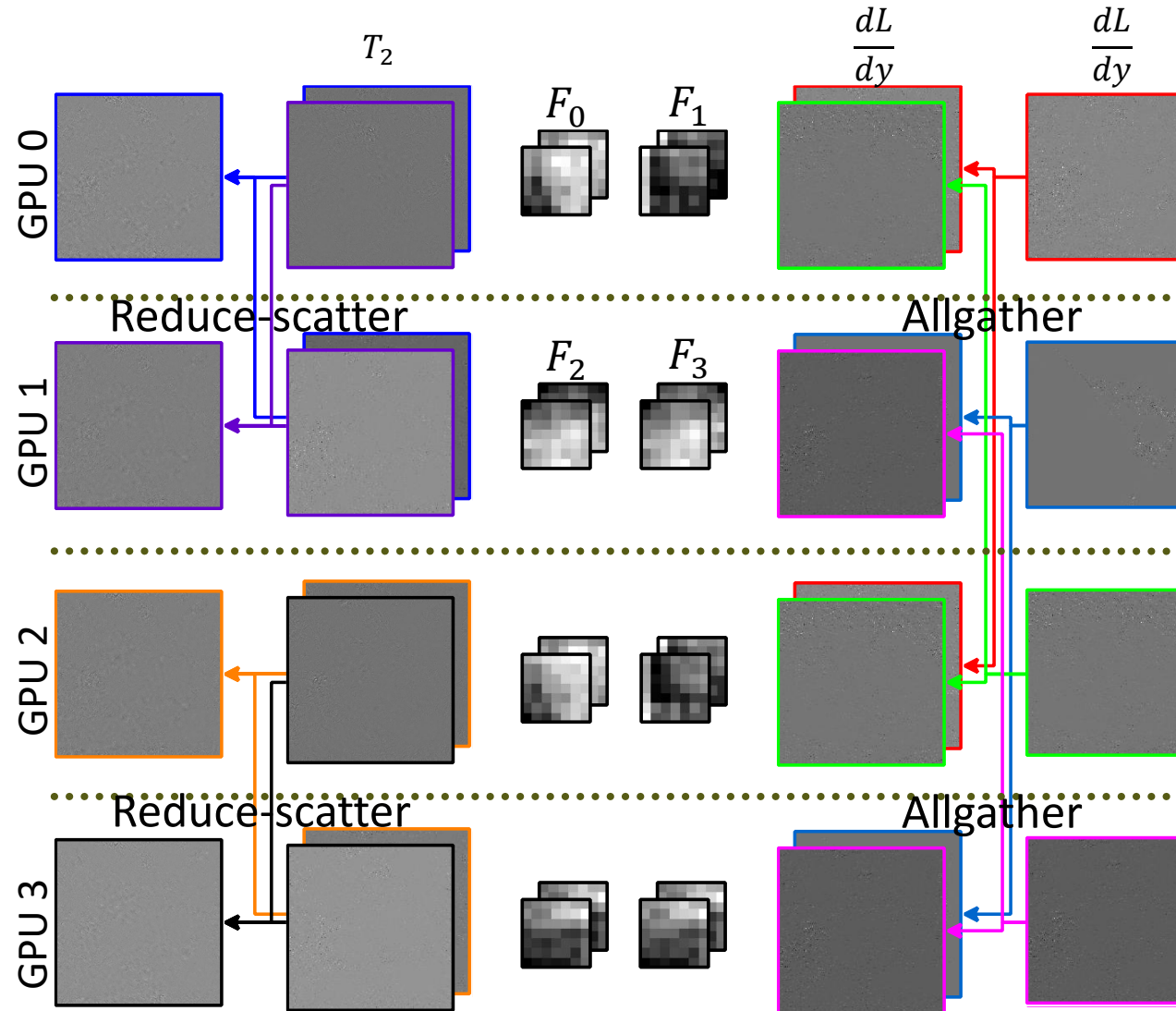
Stationary-w: Backward



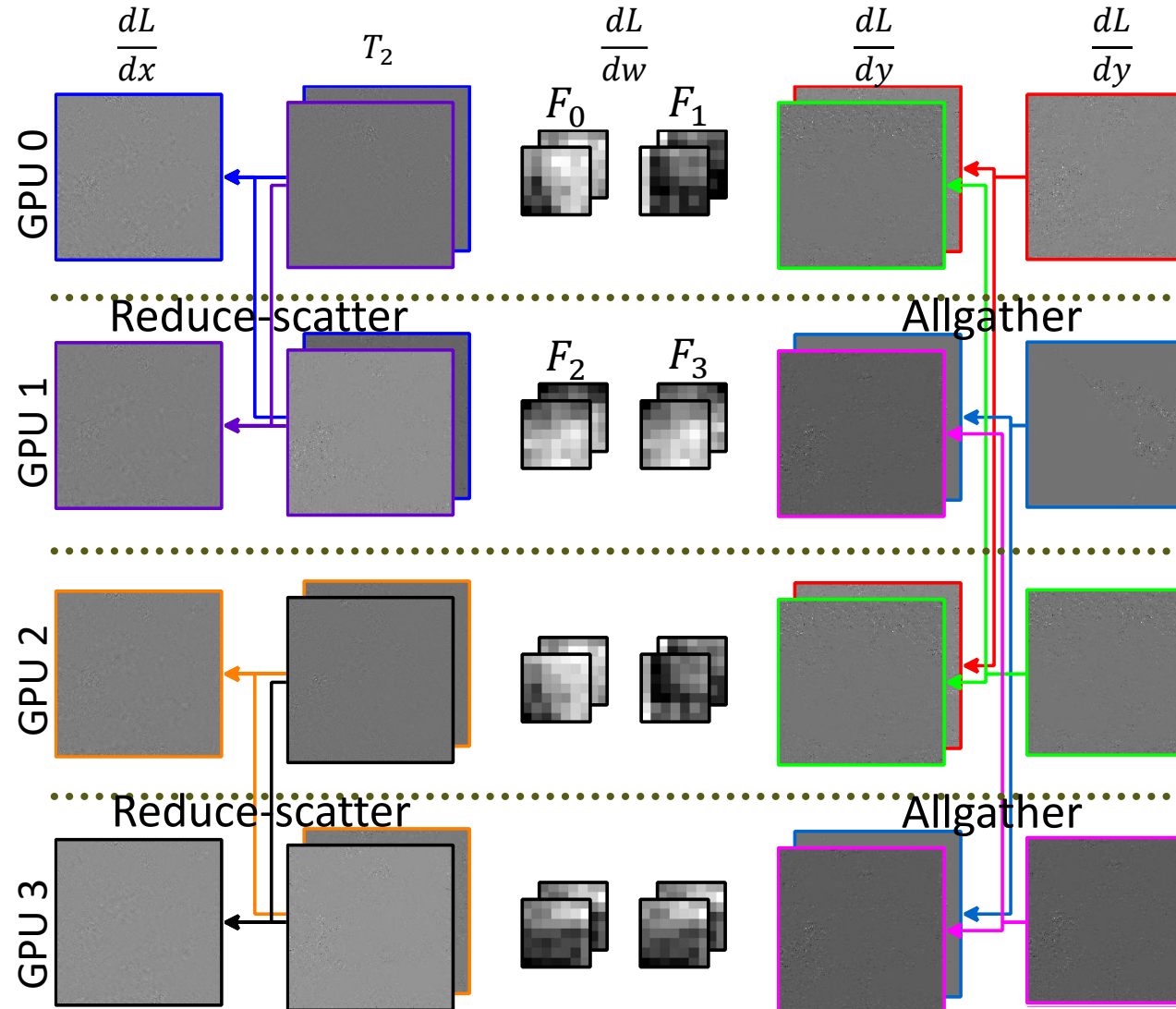
Stationary-w: Backward



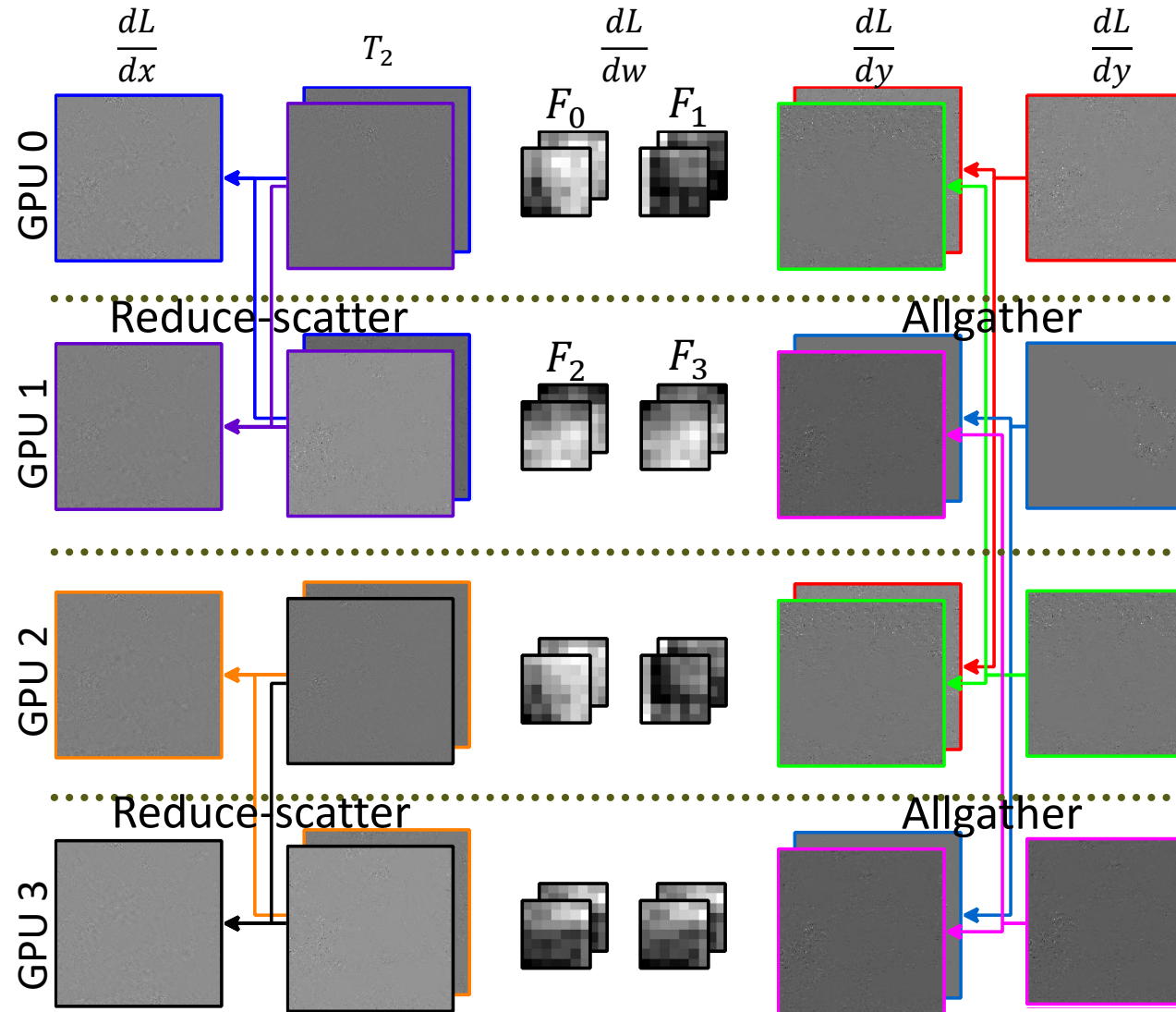
Stationary-w: Backward



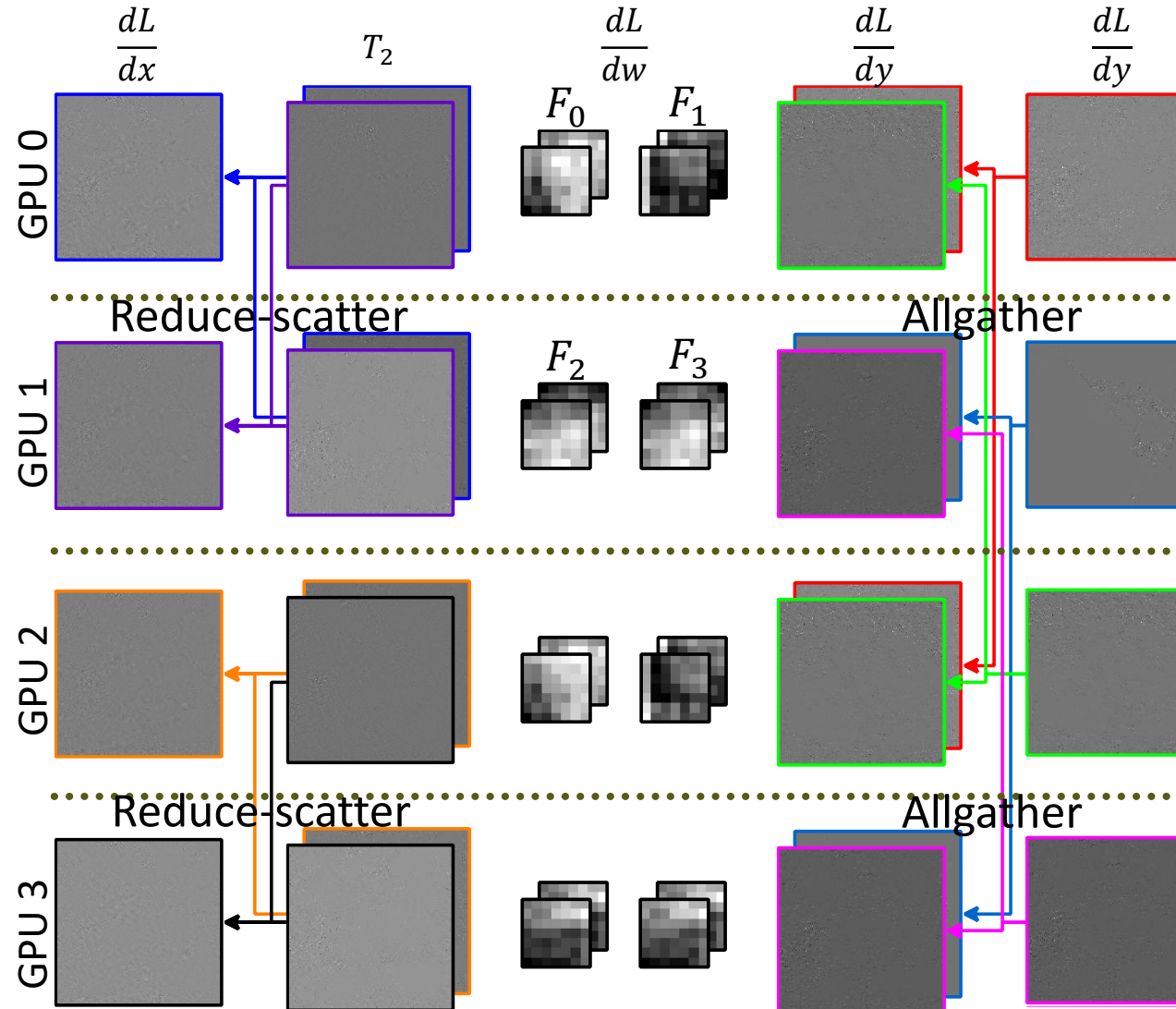
Stationary-w: Backward



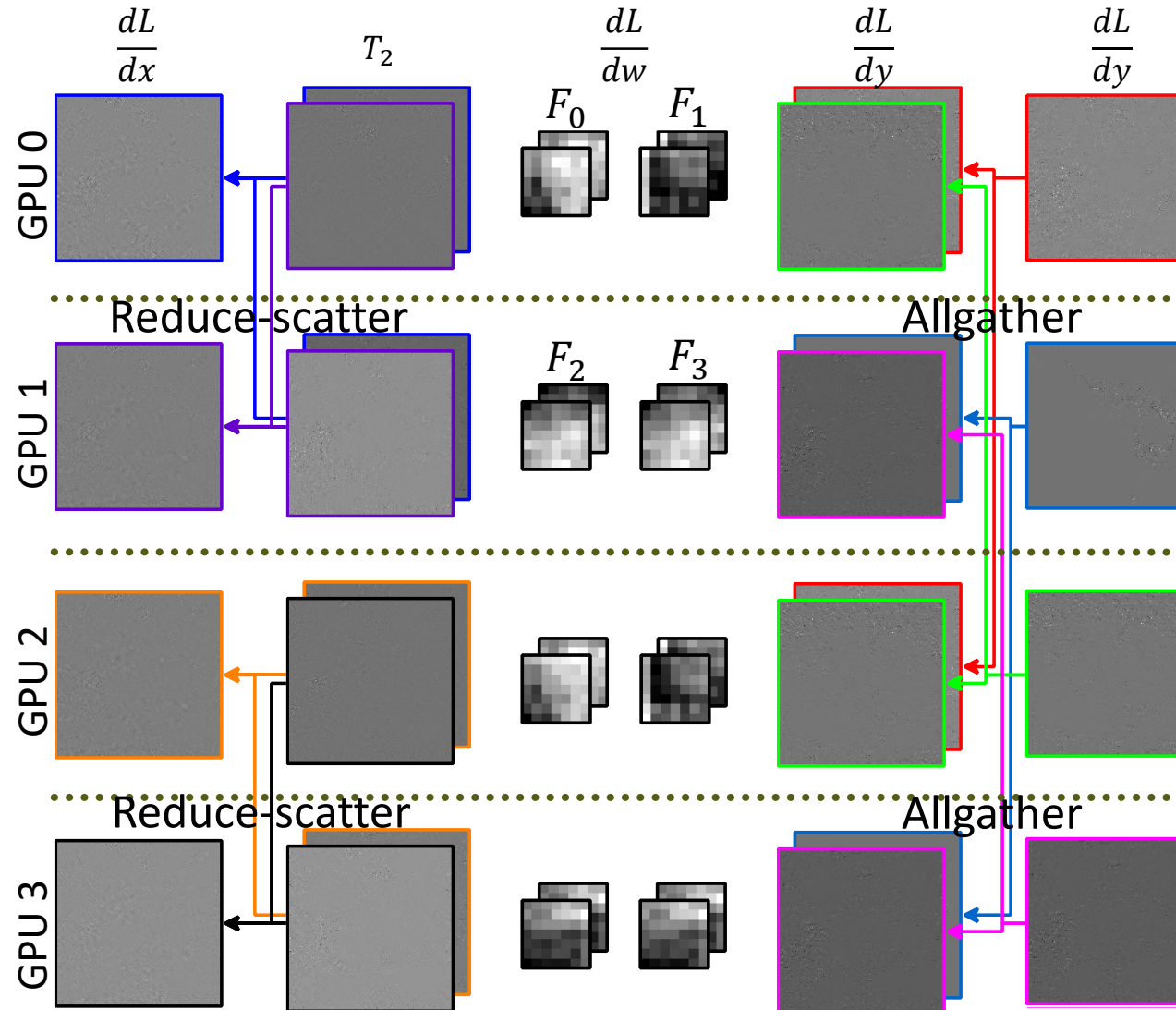
Stationary-w: Backward



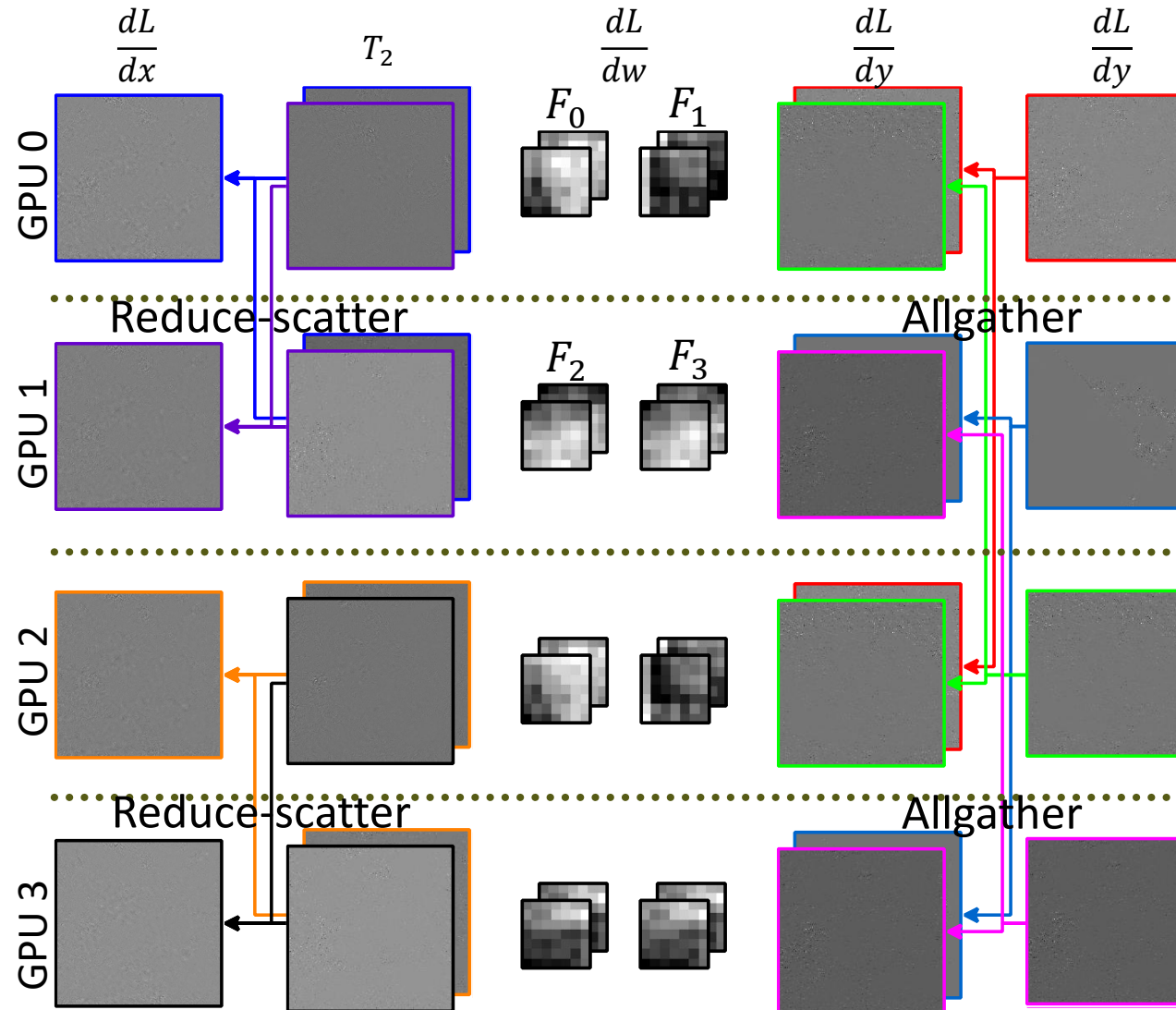
Stationary-w: Backward



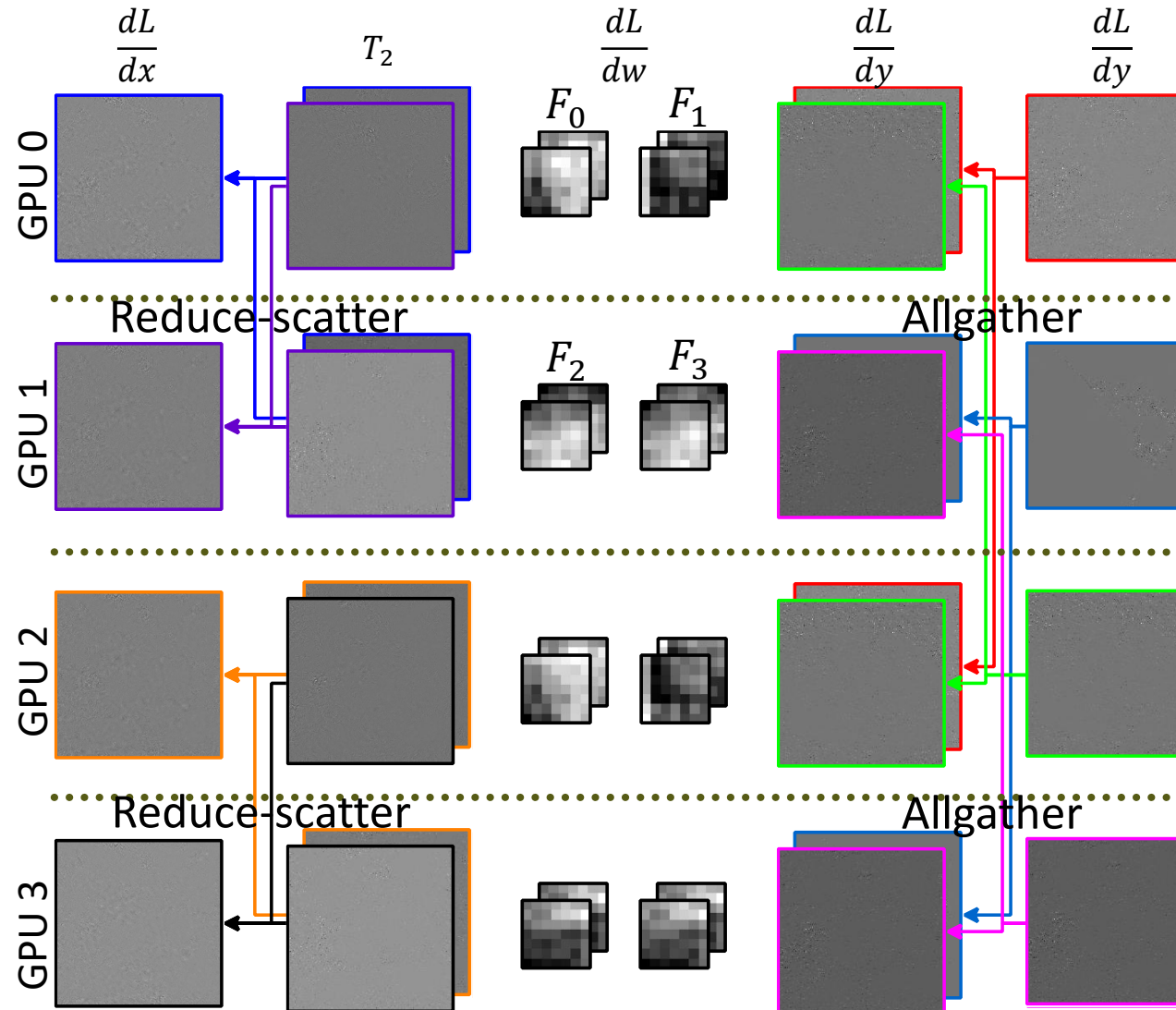
Stationary-w: Backward



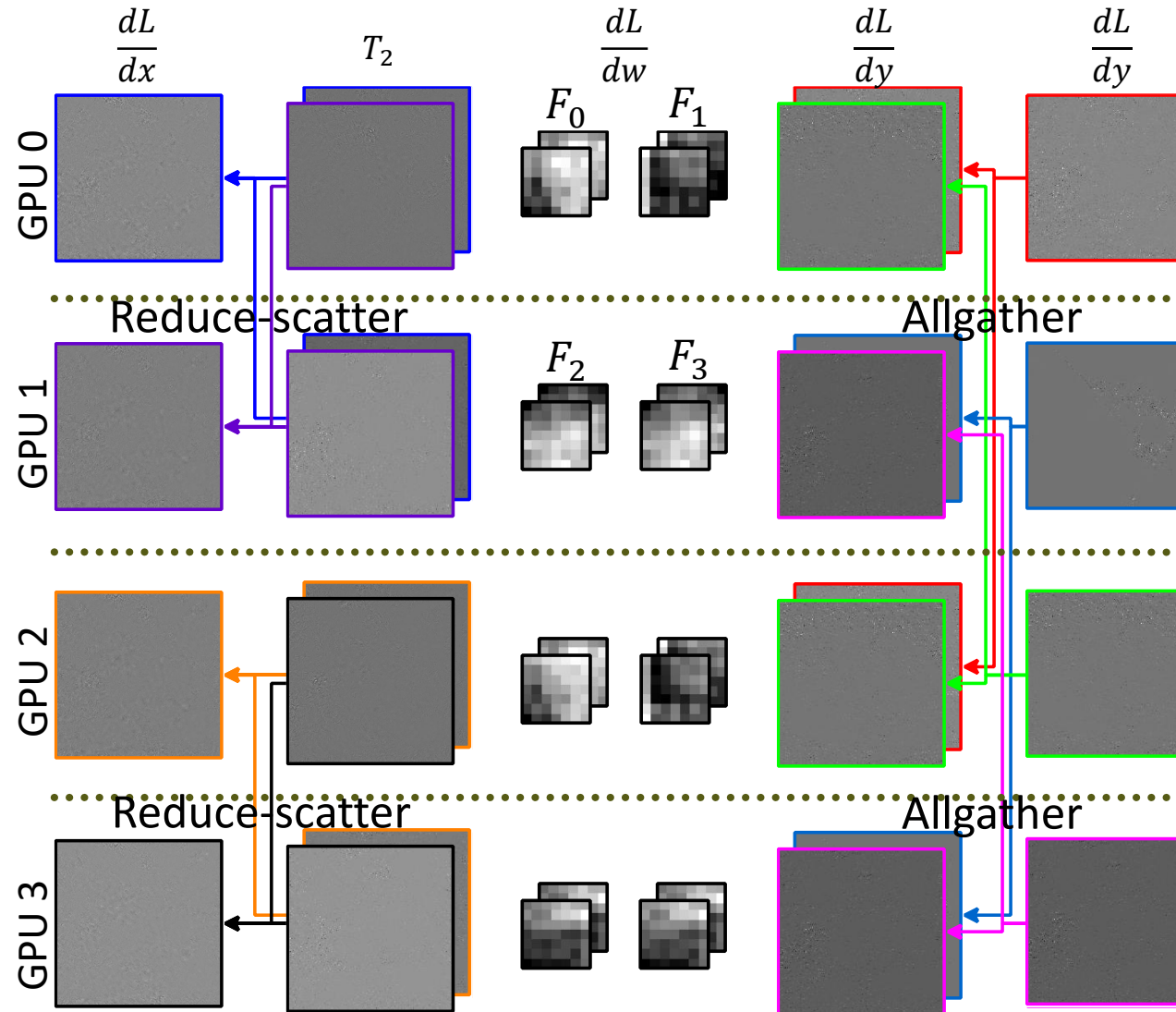
Stationary-w: Backward



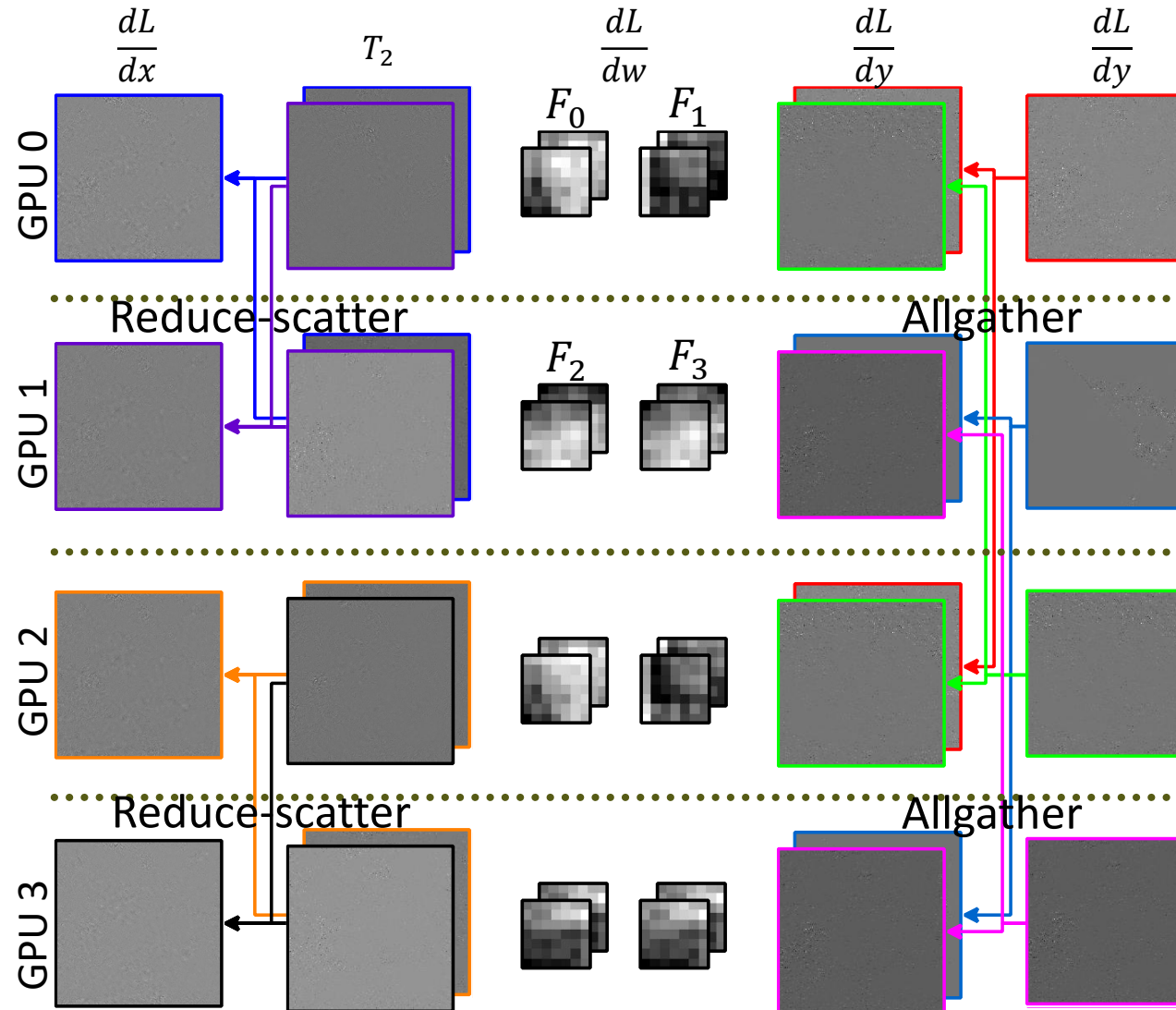
Stationary-w: Backward



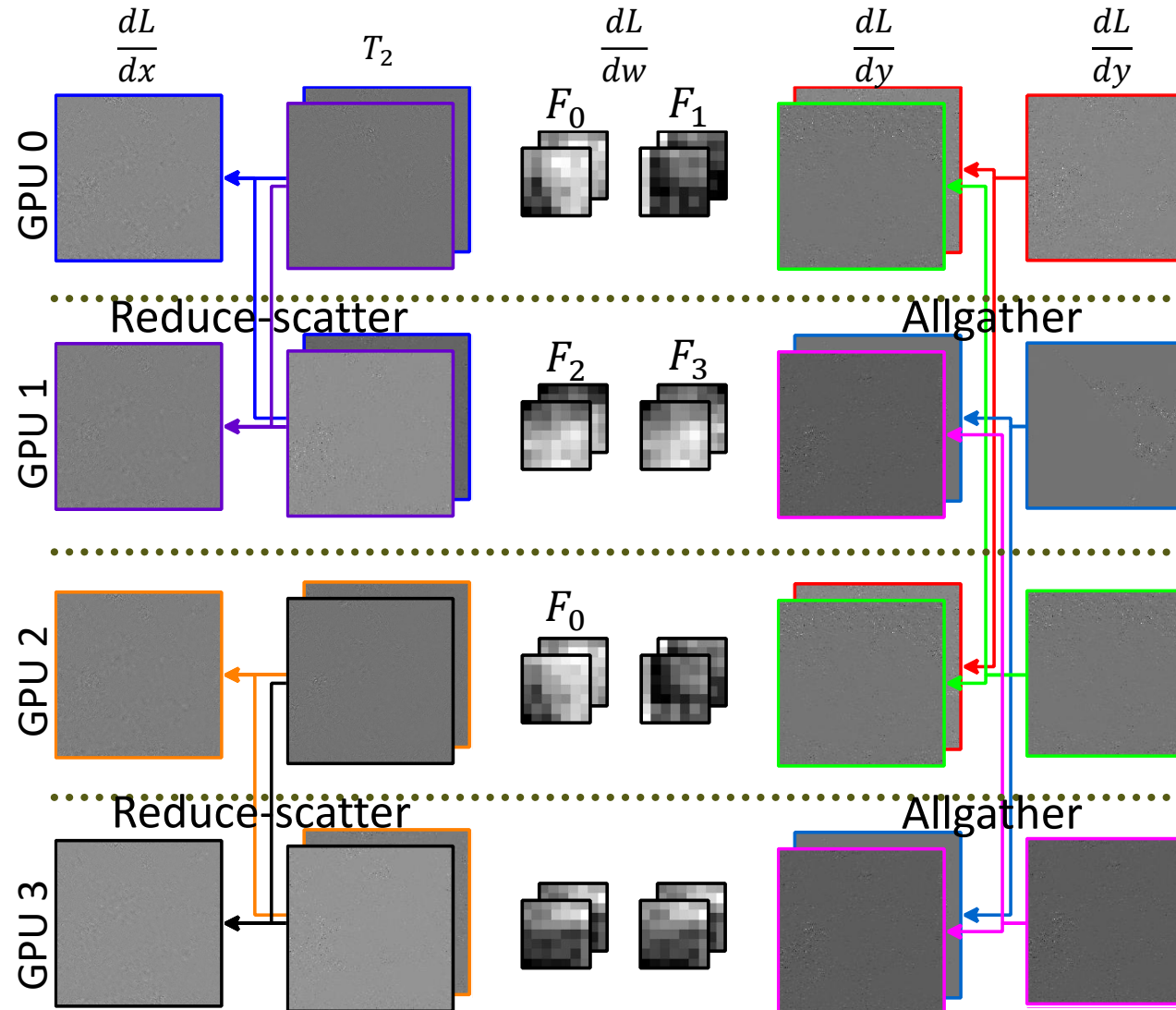
Stationary-w: Backward



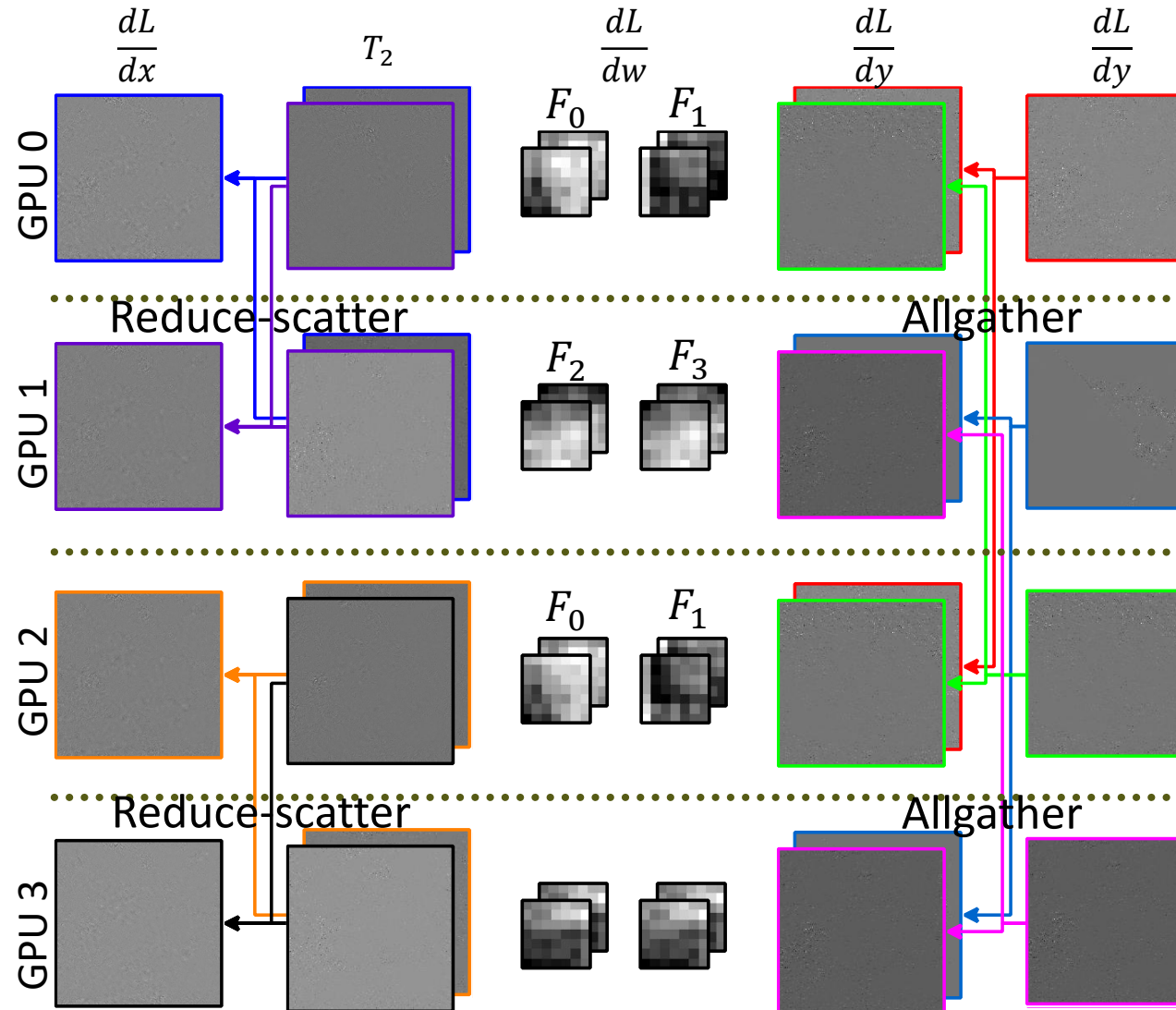
Stationary-w: Backward



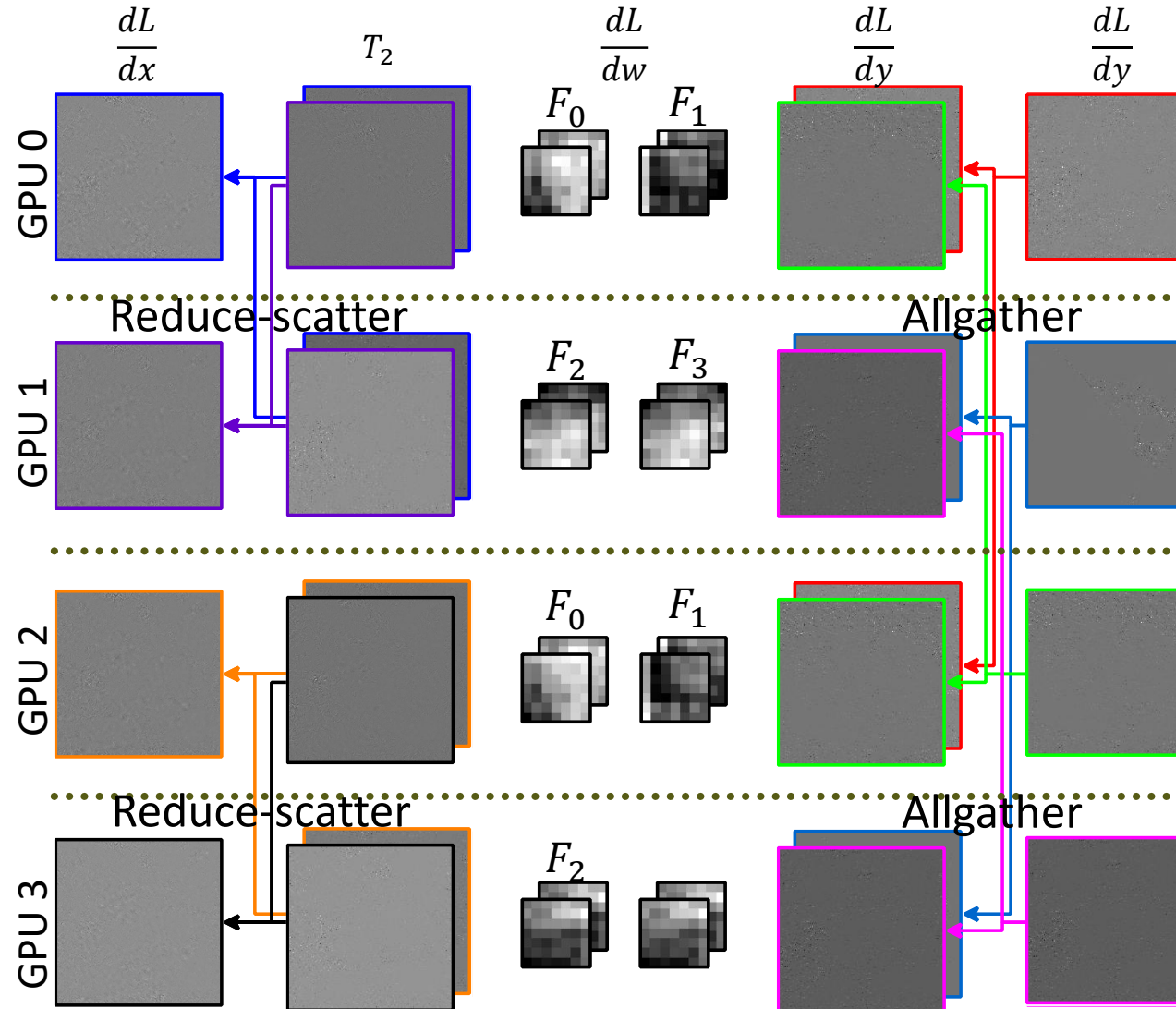
Stationary-w: Backward



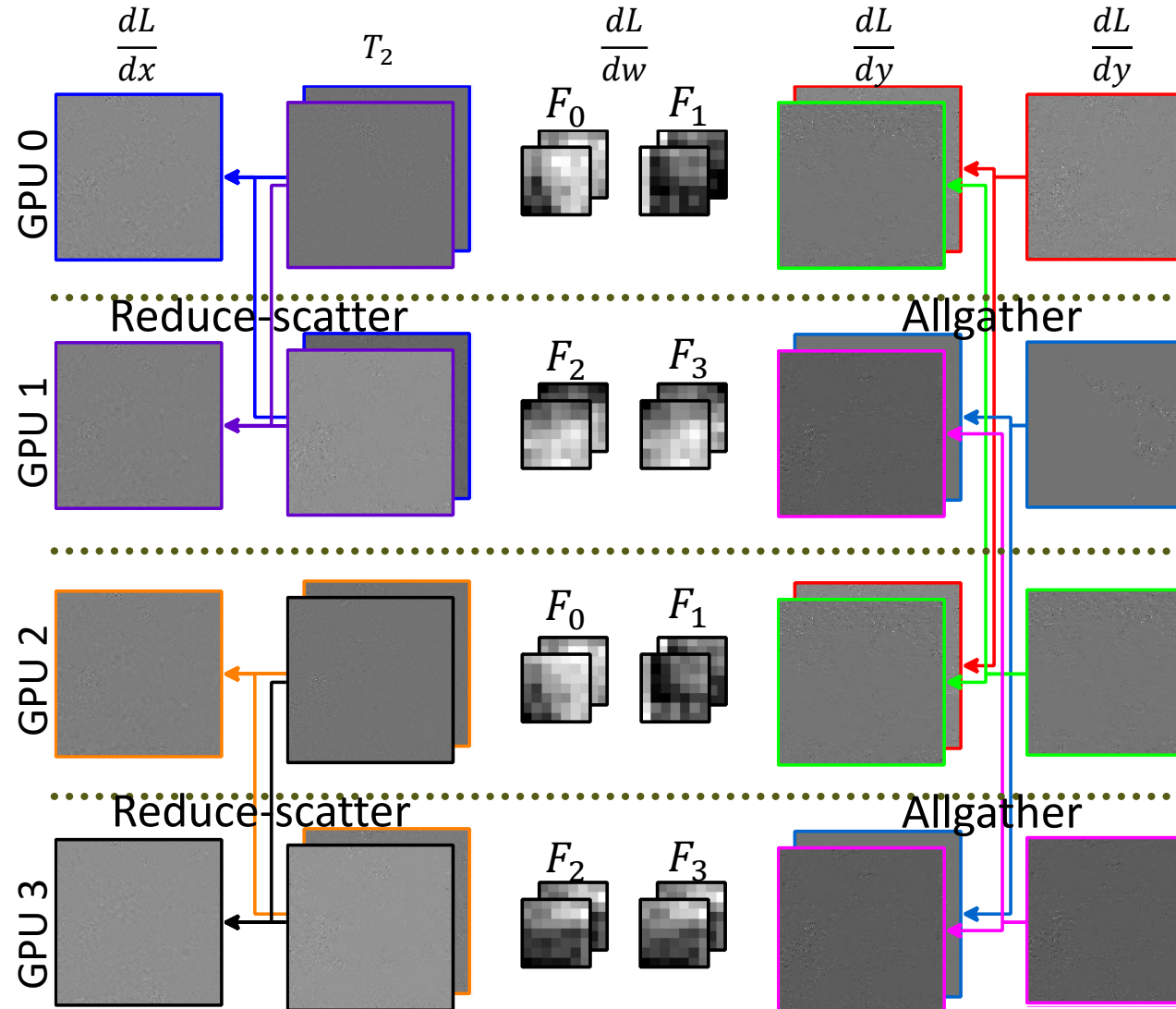
Stationary-w: Backward



Stationary-w: Backward

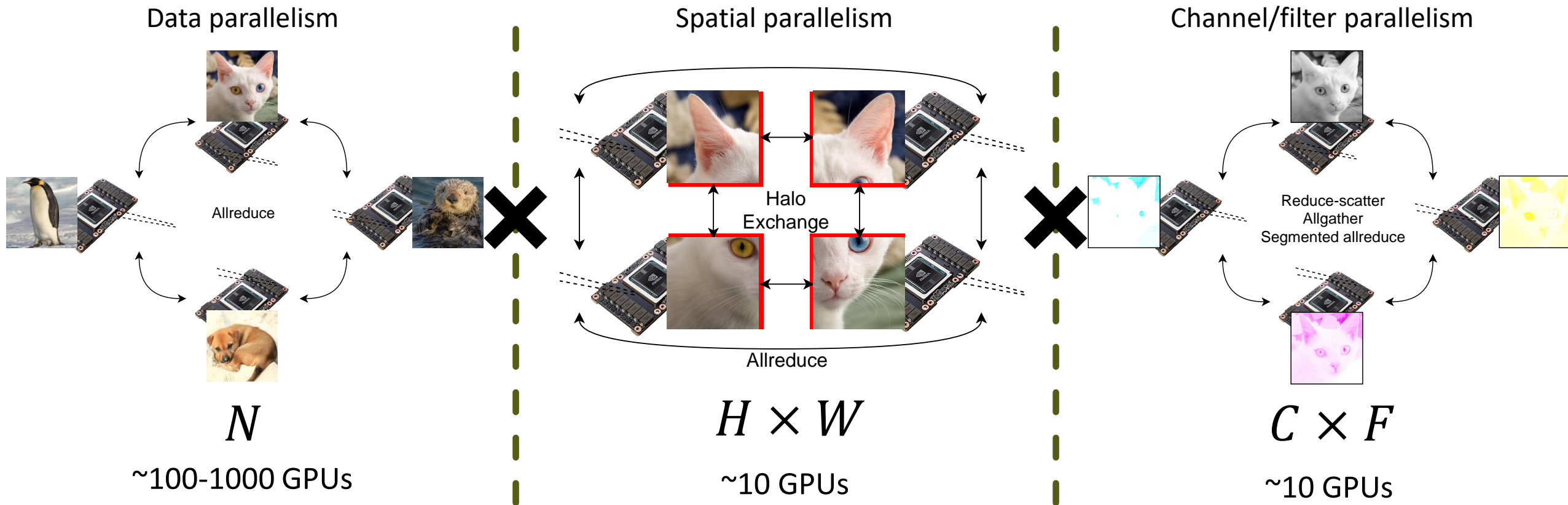


Stationary-w: Backward



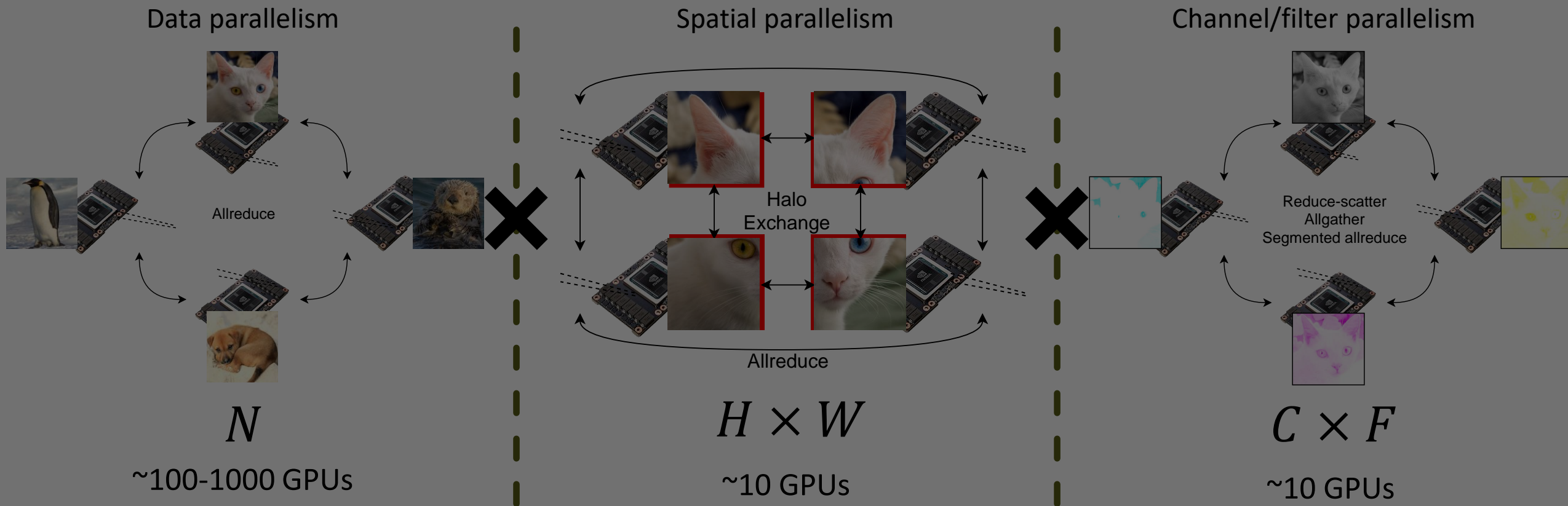
General distributed convolution

- Provide a variety of options to enable and improve strong and weak scaling
- Support a full spectrum of data and model types



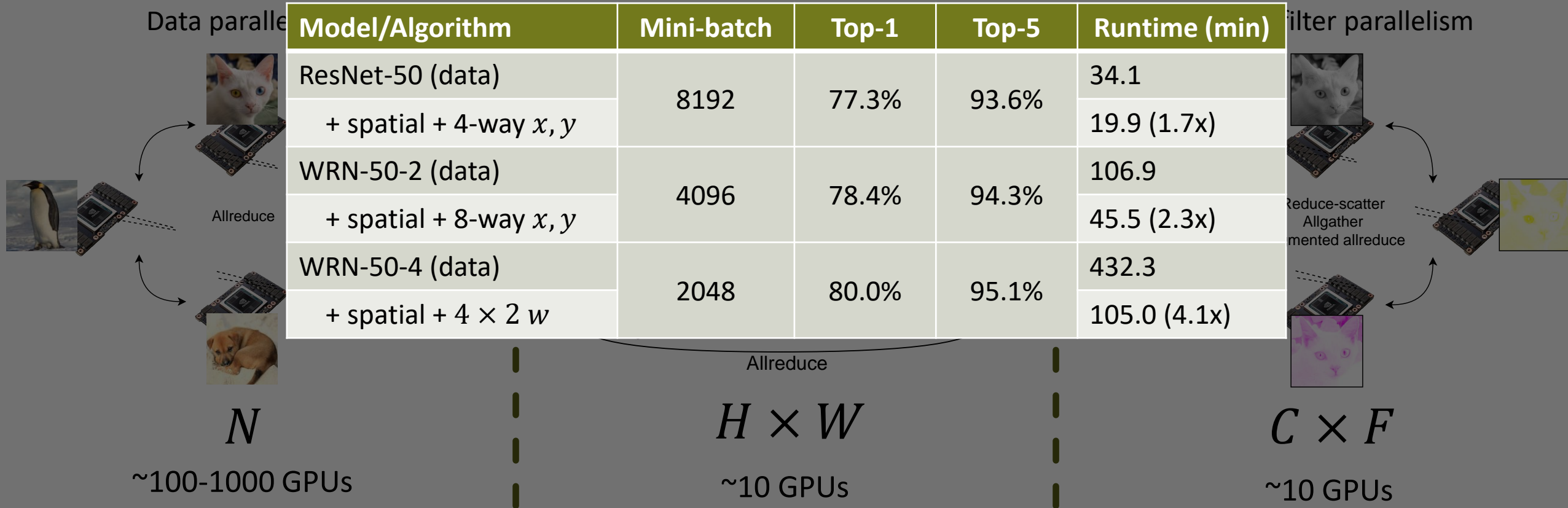
General distributed convolution

- Provide a variety of options to enable and improve strong and weak scaling
- Support a full spectrum of data and model types



General distributed convolution

- Provide a variety of options to enable and improve strong and weak scaling
- Support a full spectrum of data and model types

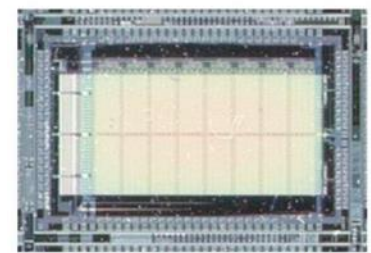


Specialized hardware

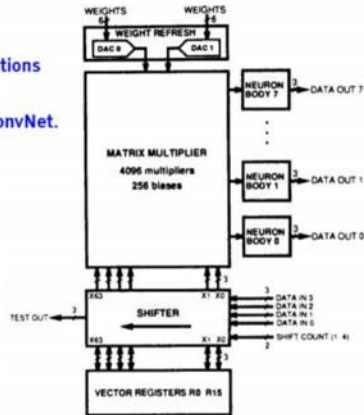
Specialized hardware



- [Boser, Säckinger, Bromley, LeCun, Jackel, IEEE J. SSC 26(12), 1991]
- 4096 Multiply-Accumulate operators
- 6 bit weights, 4 bit states
- 20 MHz clock
- Shift registers for efficient I/O with convolutions
- 4 GOPS (peak)
- 1000 characters per second for OCR with ConvNet.



Source: Yann LeCun

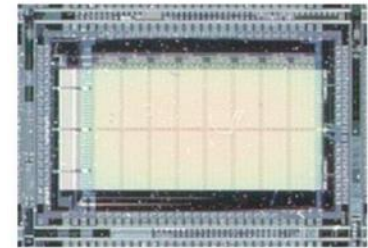


Specialized hardware

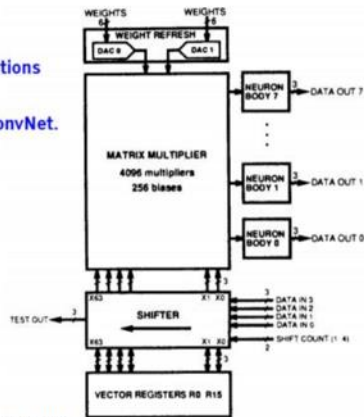


[Boser, Säckinger, Bromley, LeCun, Jackel, IEEE J. SSC 26(12), 1991]

- 4096 Multiply-Accumulate operators
- 6 bit weights, 4 bit states
- 20 MHz clock
- Shift registers for efficient I/O with convolutions
- 4 GOPS (peak)
- 1000 characters per second for OCR with ConvNet.



Source: Yann LeCun



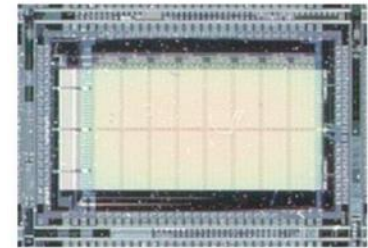
ANNA [1991]

Specialized hardware

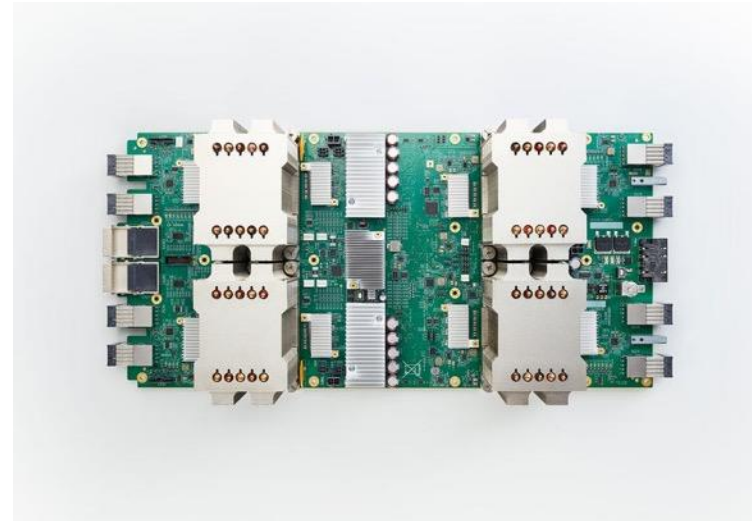
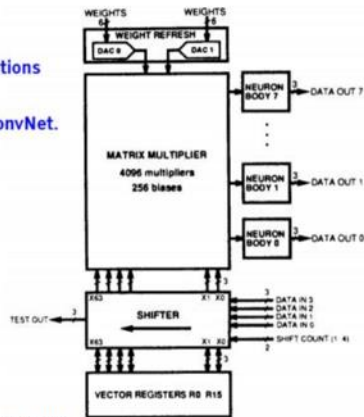


[Boser, Säckinger, Bromley, LeCun, Jackel, IEEE J. SSC 26(12), 1991]

- 4096 Multiply-Accumulate operators
- 6 bit weights, 4 bit states
- 20 MHz clock
- Shift registers for efficient I/O with convolutions
- 4 GOPS (peak)
- 1000 characters per second for OCR with ConvNet.



Source: Yann LeCun



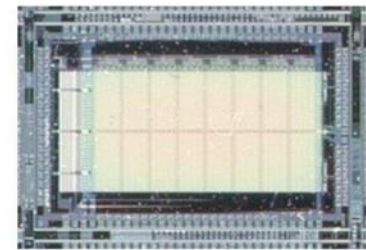
ANNA [1991]

Specialized hardware

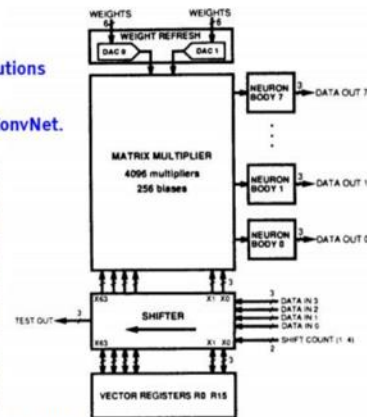


[Boser, Säckinger, Bromley, LeCun, Jackel, IEEE J. SSC 26(12), 1991]

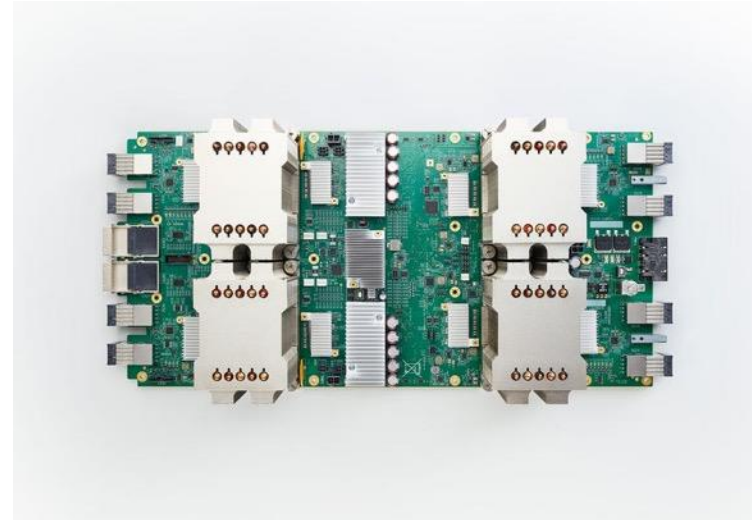
- 4096 Multiply-Accumulate operators
- 6 bit weights, 4 bit states
- 20 MHz clock
- Shift registers for efficient I/O with convolutions
- 4 GOPS (peak)
- 1000 characters per second for OCR with ConvNet.



Source: Yann LeCun



ANNA [1991]

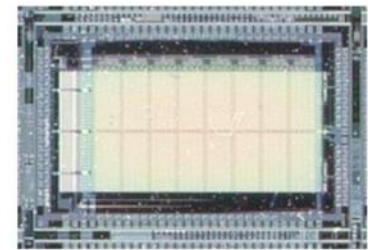


**TPU v1(inference); v2-3(training)
[2016-today]**

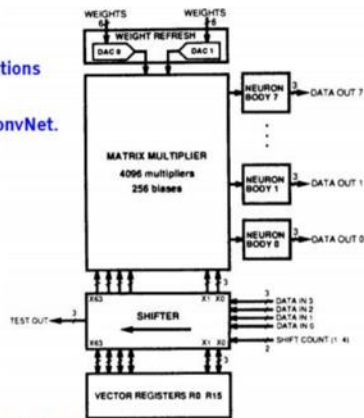
Specialized hardware



- [Boser, Säckinger, Bromley, LeCun, Jackel, IEEE J. SSC 26(12), 1991]
- 4096 Multiply-Accumulate operators
- 6 bit weights, 4 bit states
- 20 MHz clock
- Shift registers for efficient I/O with convolutions
- 4 GOPS (peak)
- 1000 characters per second for OCR with ConvNet.



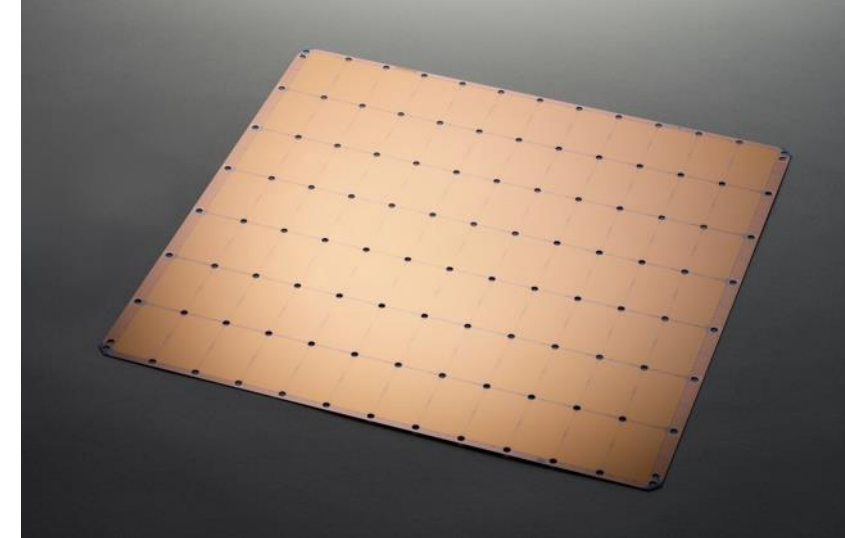
Source: Yann LeCun



ANNA [1991]



**TPU v1(inference); v2-3(training)
[2016-today]**

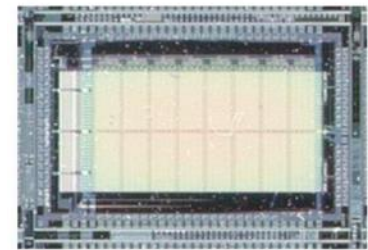


Cerebrus CS-1 [2019]

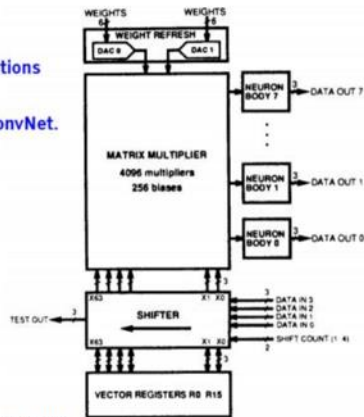
Specialized hardware



- [Boser, Säckinger, Bromley, LeCun, Jackel, IEEE J. SSC 26(12), 1991]
- 4096 Multiply-Accumulate operators
- 6 bit weights, 4 bit states
- 20 MHz clock
- Shift registers for efficient I/O with convolutions
- 4 GOPS (peak)
- 1000 characters per second for OCR with ConvNet.



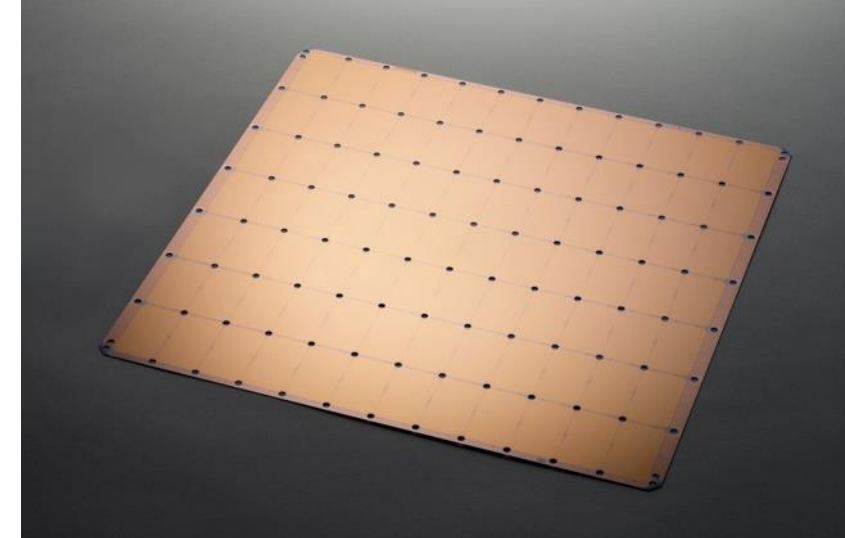
Source: Yann LeCun



ANNA [1991]



TPU v1(inference); v2-3(training)
[2016-today]

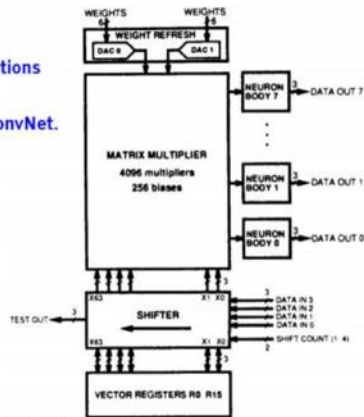


Cerebrus CS-1 [2019]

Specialized hardware

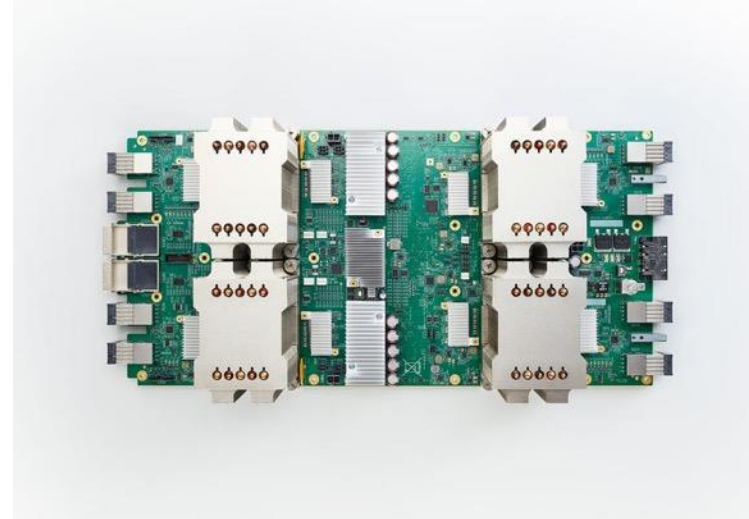


- [Boser, Säckinger, Bromley, LeCun, Jackel, IEEE J. SSC 26(12), 1991]
- 4096 Multiply-Accumulate operators
- 6 bit weights, 4 bit states
- 20 MHz clock
- Shift registers for efficient I/O with convolutions
- 4 GOPS (peak)
- 1000 characters per second for OCR with ConvNet.

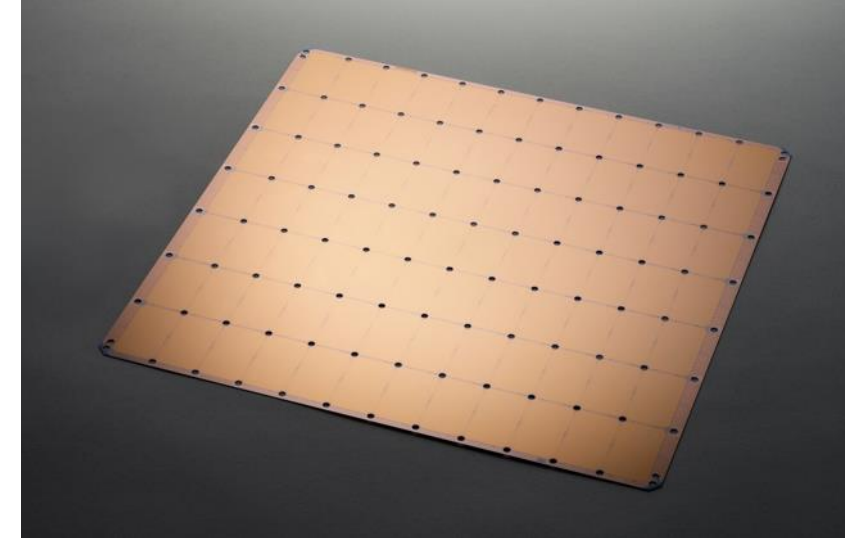


Source: Yann LeCun

ANNA [1991]



**TPU v1(inference); v2-3(training)
[2016-today]**



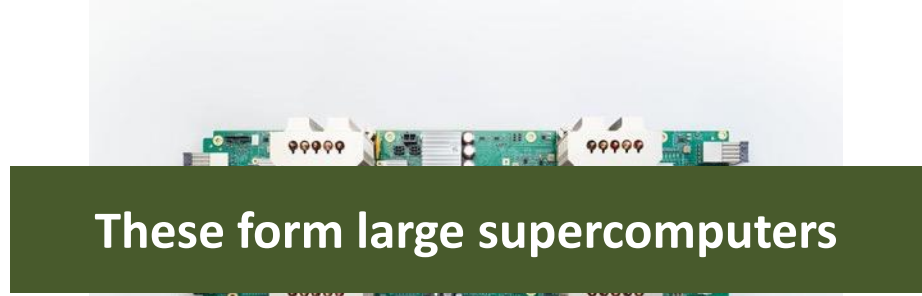
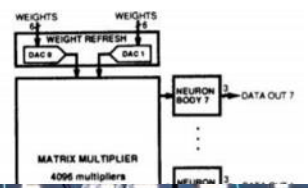
Cerebrus CS-1 [2019]

Literally hundreds of other startups in this space

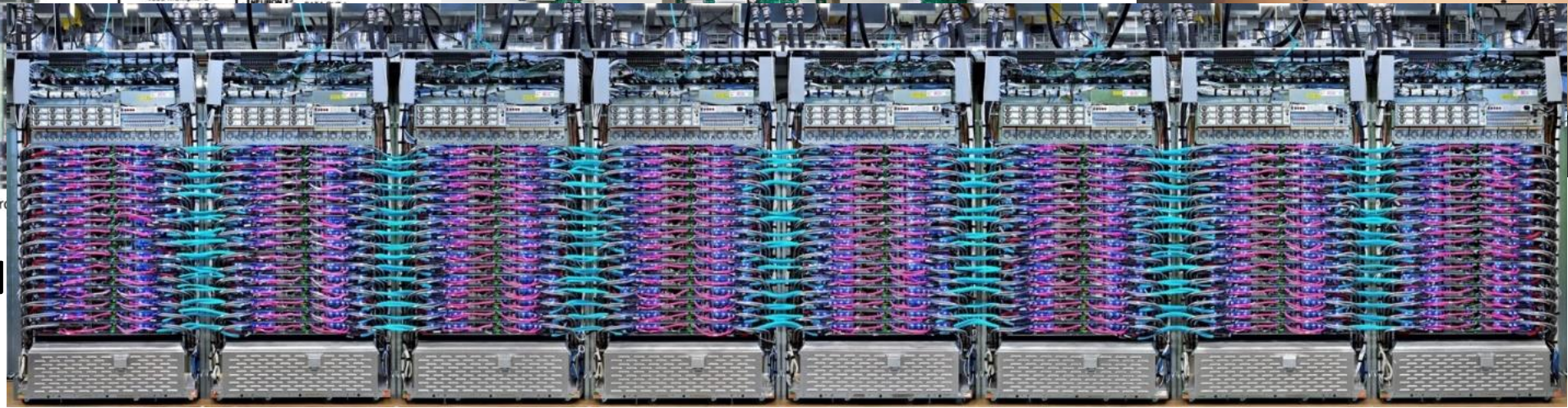
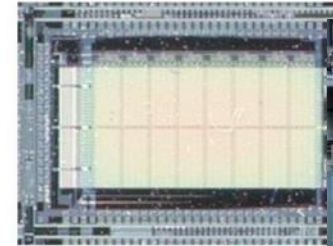
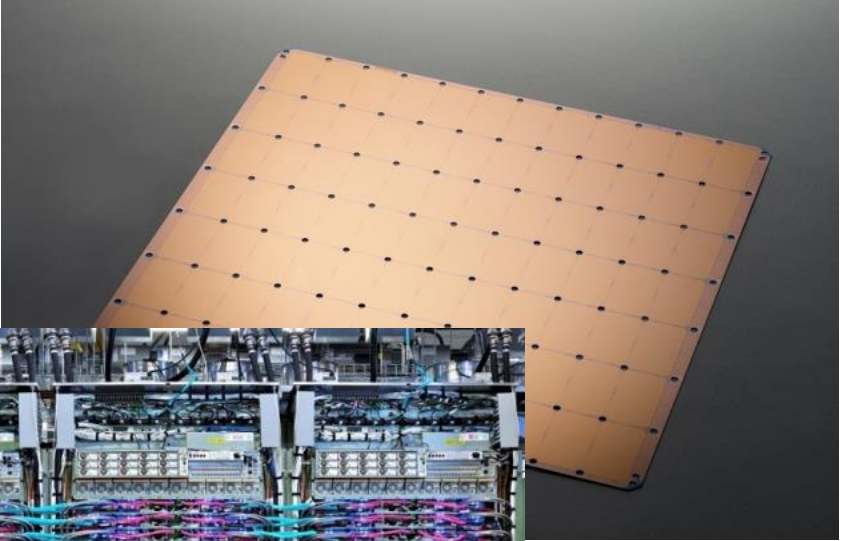
Specialized hardware



- [Boser, Säckinger, Bromley, LeCun, Jackel, IEEE J. SSC 26(12), 1991]
- 4096 Multiply-Accumulate operators
- 6 bit weights, 4 bit states
- 20 MHz clock
- Shift registers for efficient I/O with convolutions
- 4 GOPS (peak)
- 1000 characters per second for OCR with ConvNet.



These form large supercomputers



ANN

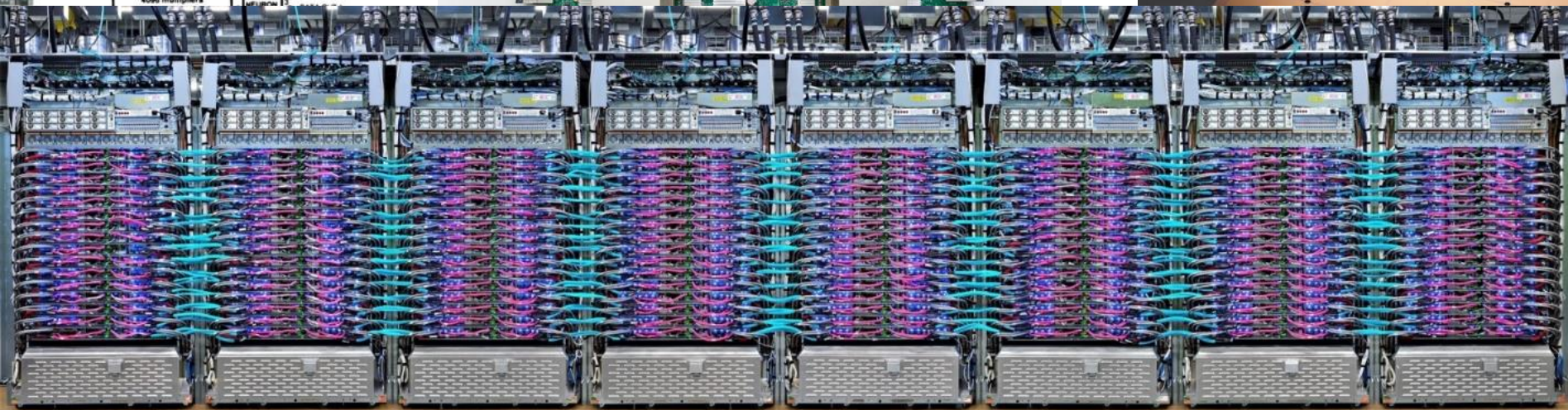
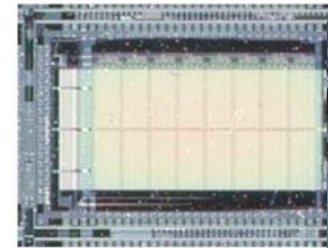
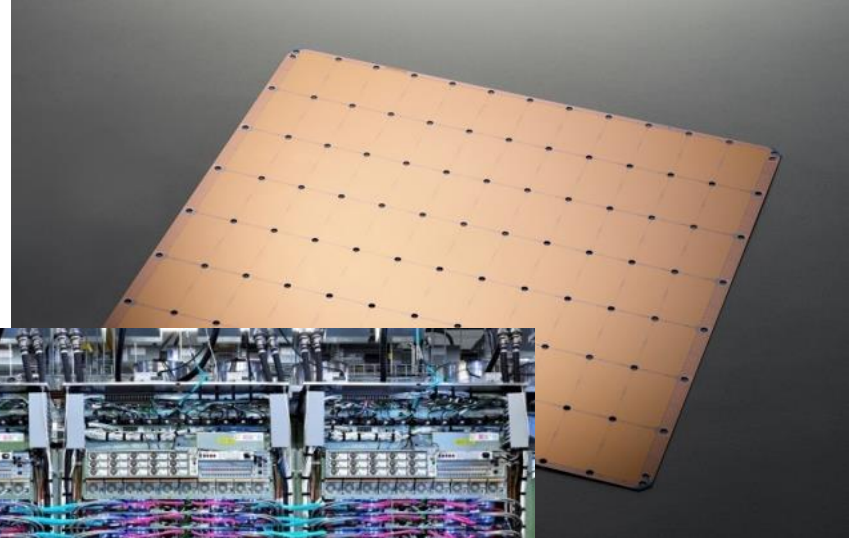
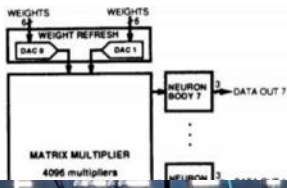
[2019]

Literally hundreds of other startups in this space

Specialized hardware



- [Boser, Säckinger, Bromley, LeCun, Jackel, IEEE J. SSC 26(12), 1991]
- 4096 Multiply-Accumulate operators
- 6 bit weights, 4 bit states
- 20 MHz clock
- Shift registers for efficient I/O with convolutions
- 4 GOPS (peak)
- 1000 characters per second for OCR with ConvNet.



ANN

[2019]

Literally hundreds of other startups in this space

DNN Compilers

- Use techniques from compiler construction: DNN → Graph → IR → Transformations → HW Mapping

DNN Compilers

- Use techniques from compiler construction: DNN → Graph → IR → Transformations → HW Mapping

XLA HLO

Target-independent
Optimizations & Analyses

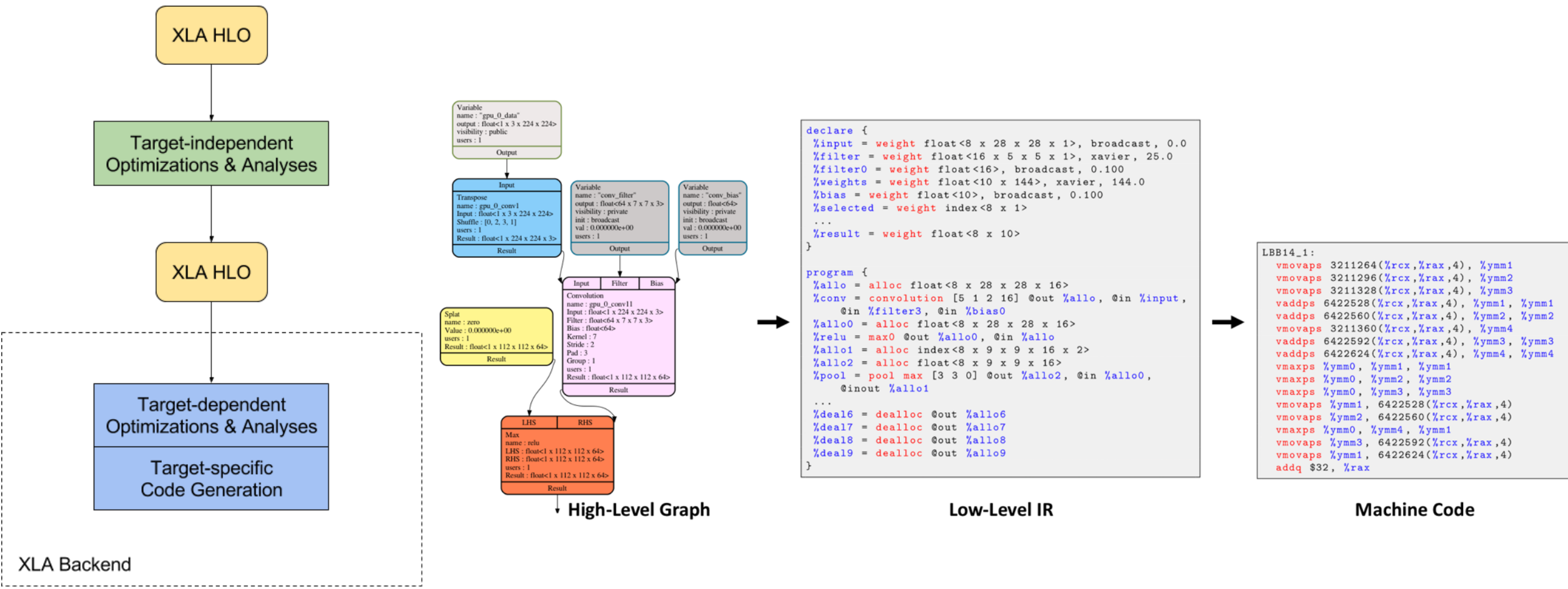
XLA HLO

Target-dependent
Optimizations & Analyses
Target-specific
Code Generation

XLA Backend

DNN Compilers

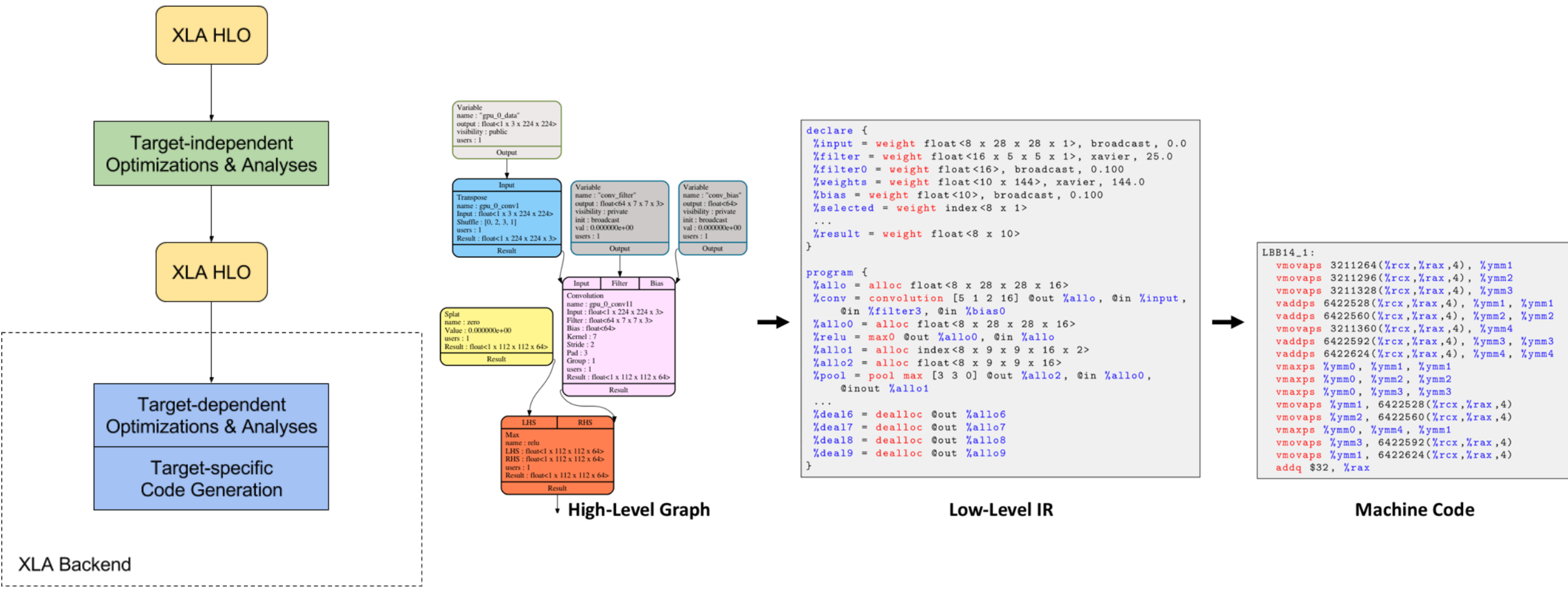
- Use techniques from compiler construction: DNN → Graph → IR → Transformations → HW Mapping



TensorFlow XLA

DNN Compilers

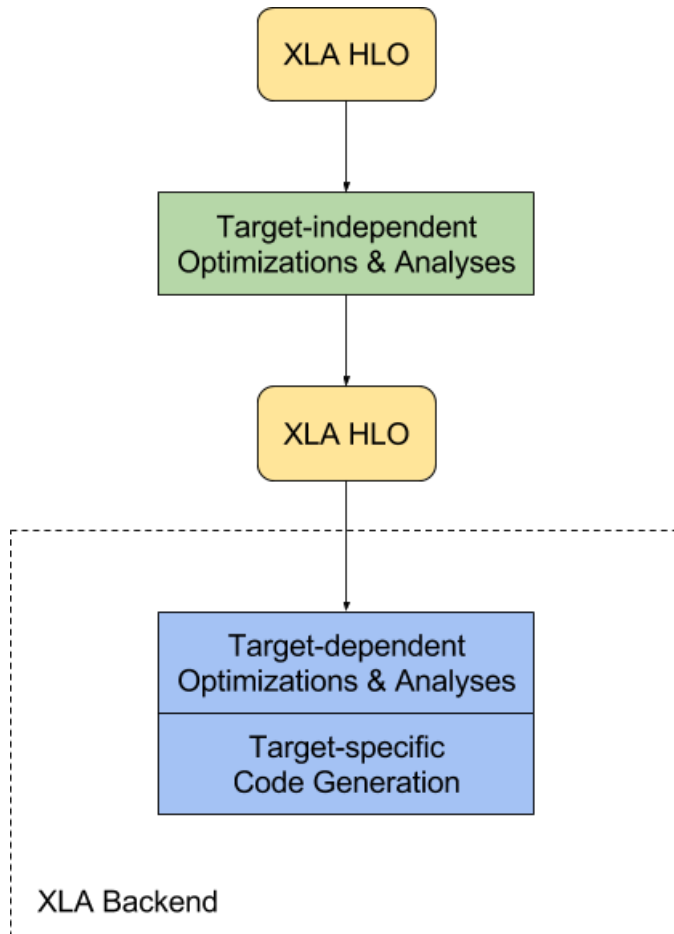
- Use techniques from compiler construction: DNN → Graph → IR → Transformations → HW Mapping



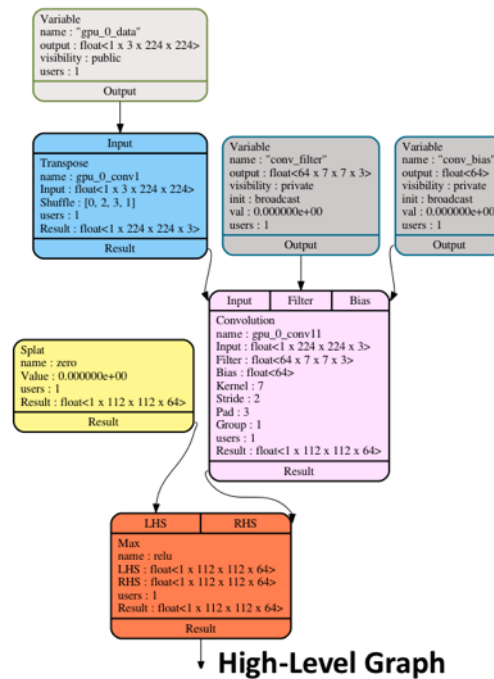
TensorFlow XLA

DNN Compilers

- Use techniques from compiler construction: DNN → Graph → IR → Transformations → HW Mapping



TensorFlow XLA



High-Level Graph

```

declare {
  %input = weight float<8 x 28 x 28 x 1>, broadcast, 0.0
  %filter = weight float<16 x 5 x 5 x 1>, xavier, 25.0
  %filter0 = weight float<16>, broadcast, 0.100
  %weights = weight float<10 x 144>, xavier, 144.0
  %bias = weight float<10>, broadcast, 0.100
  %selected = weight index<8 x 1>
  ...
  %result = weight float<8 x 10>
}

program {
  %allo = alloc float<8 x 28 x 28 x 16>
  %conv = convolution [5 1 2 16] @out %allo, @in %input,
    @in %filters3, @in %bias0
  %allo0 = alloc float<8 x 28 x 28 x 16>
  %relu = max0 @out %allo0, @in %allo
  %allo1 = alloc index<8 x 9 x 9 x 16 x 2>
  %allo2 = alloc float<8 x 9 x 9 x 16>
  %pool = pool max [3 3 0] @out %allo2, @in %allo1,
    @inout %allo1
  ...
  %deal6 = dealloc @out %allo6
  %deal7 = dealloc @out %allo7
  %deal8 = dealloc @out %allo8
  %deal9 = dealloc @out %allo9
}
    
```

Low-Level IR

```

LBB14_1:
  vmovaps 3211264(%rcx,%rax,4), %ymm1
  vmovaps 3211296(%rcx,%rax,4), %ymm2
  vmovaps 3211328(%rcx,%rax,4), %ymm3
  vaddps 6422528(%rcx,%rax,4), %ymm1, %ymm1
  vaddps 6422560(%rcx,%rax,4), %ymm2, %ymm2
  vmovaps 3211360(%rcx,%rax,4), %ymm4
  vaddps 6422592(%rcx,%rax,4), %ymm3, %ymm3
  vaddps 6422624(%rcx,%rax,4), %ymm4, %ymm4
  vmxaps %ymm0, %ymm1, %ymm1
  vmxaps %ymm0, %ymm2, %ymm2
  vmxaps %ymm0, %ymm3, %ymm3
  vmovaps %ymm1, 6422528(%rcx,%rax,4)
  vmovaps %ymm2, 6422560(%rcx,%rax,4)
  vmovaps %ymm3, 6422592(%rcx,%rax,4)
  vmovaps %ymm1, 6422624(%rcx,%rax,4)
  addq $32, %rax
    
```

Machine Code

Facebook Glow

DNN Compilers

- Use techniques from compiler construction: DNN → Graph → IR → Transformations → HW Mapping

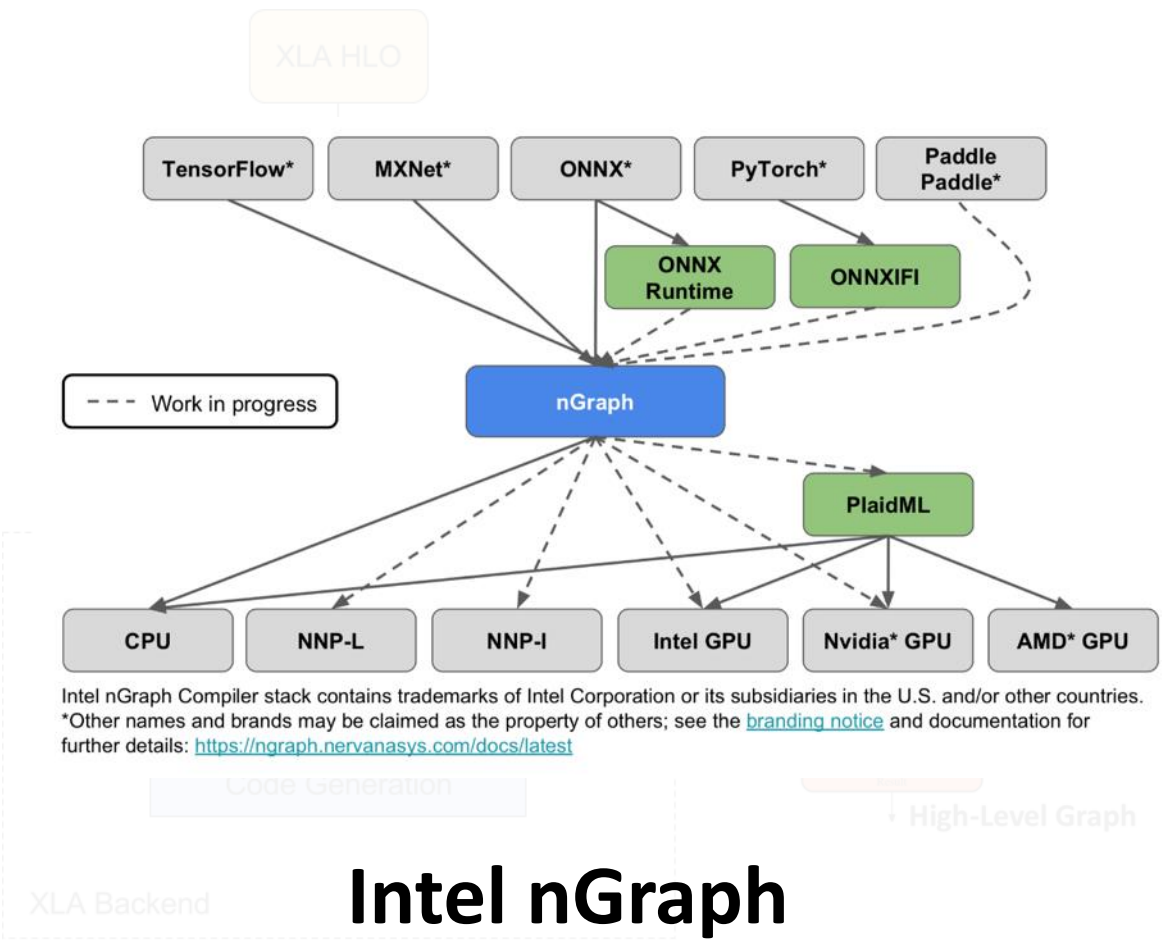


TensorFlow XLA

Facebook Glow

DNN Compilers

- Use techniques from compiler construction: DNN → Graph → IR → Transformations → HW Mapping



Intel nGraph



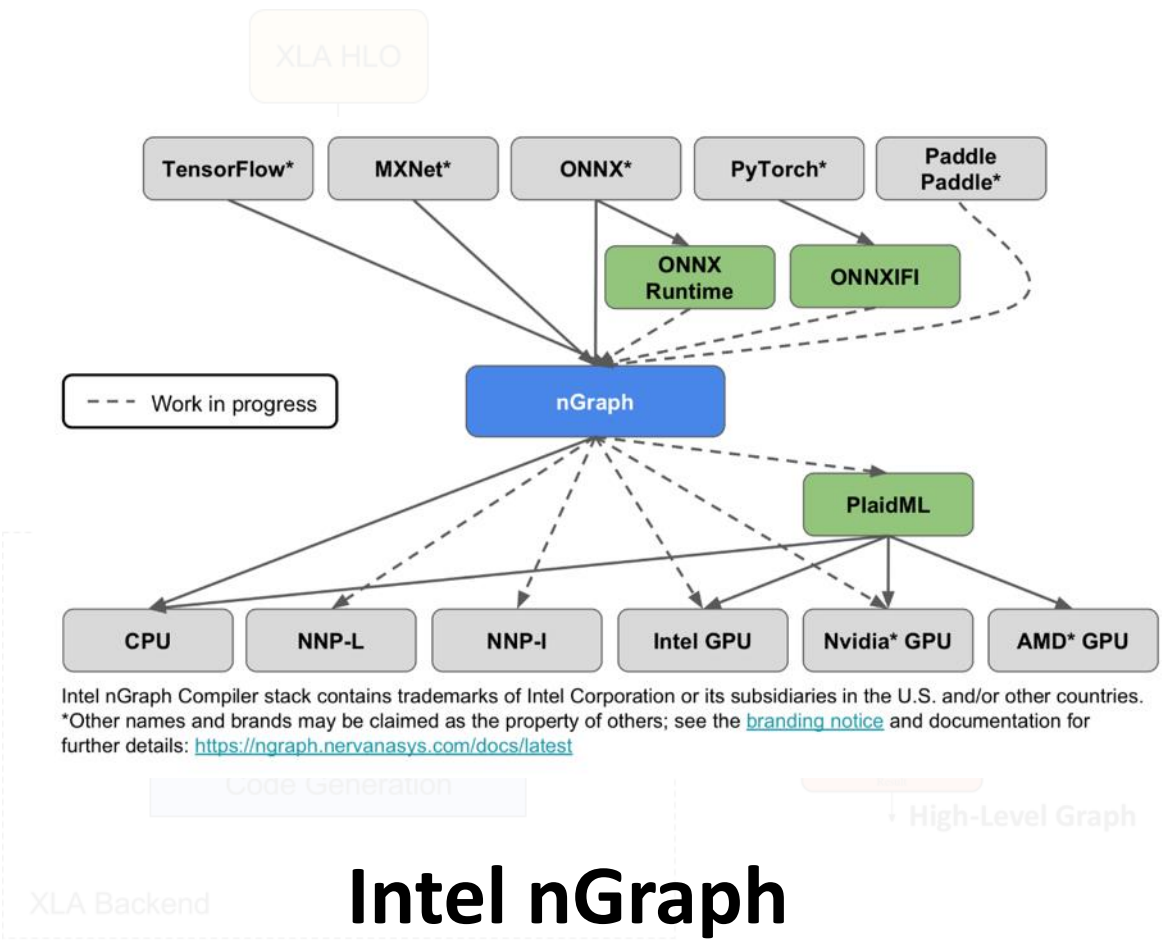
TVM Stack

TensorFlow XLA

Facebook Glow

DNN Compilers

- Use techniques from compiler construction: DNN → Graph → IR → Transformations → HW Mapping



Intel nGraph



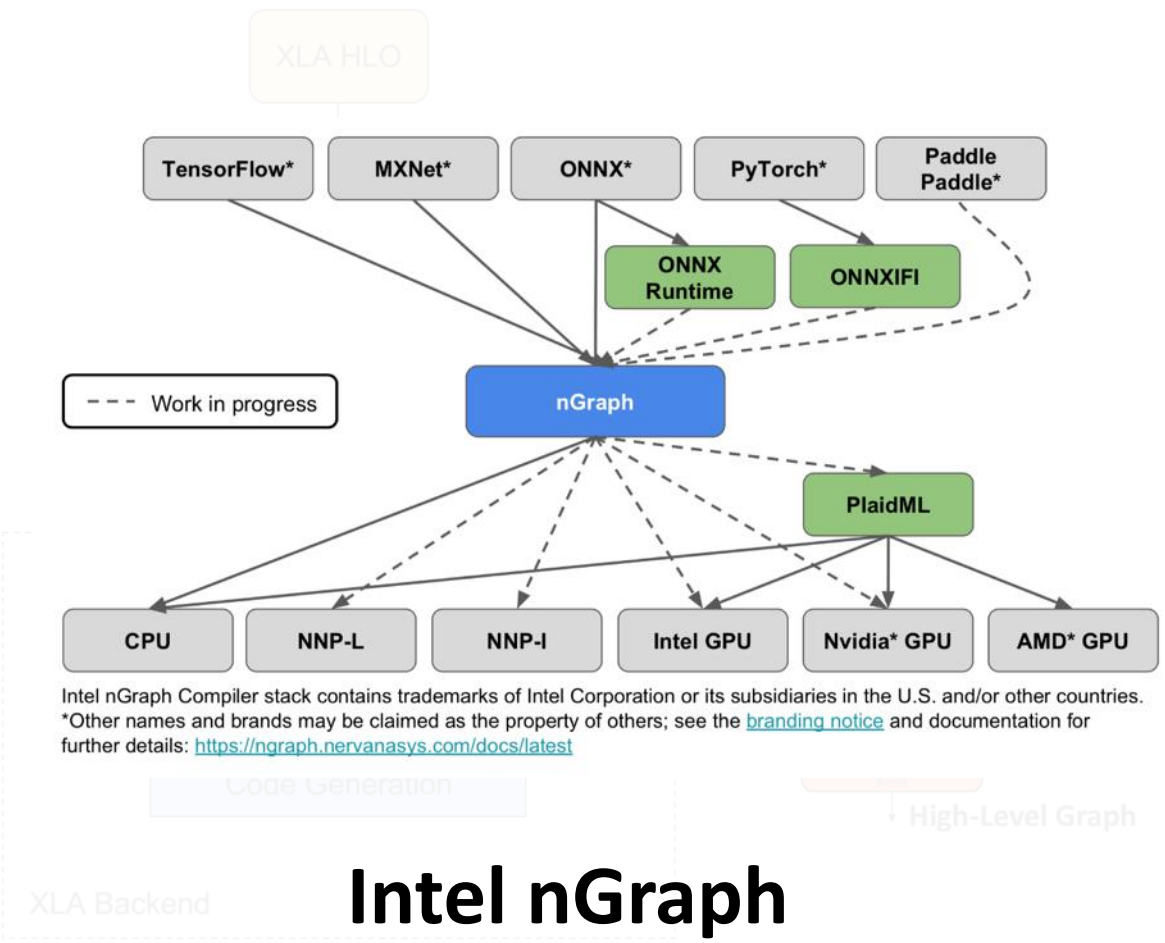
TVM Stack

TensorFlow XLA

Facebook Glow

DNN Compilers

- Use techniques from compiler construction: DNN → Graph → IR → Transformations → HW Mapping



Intel nGraph Compiler stack contains trademarks of Intel Corporation or its subsidiaries in the U.S. and/or other countries.
 *Other names and brands may be claimed as the property of others; see the [branding notice](#) and documentation for further details: <https://ngraph.nervanasys.com/docs/latest>

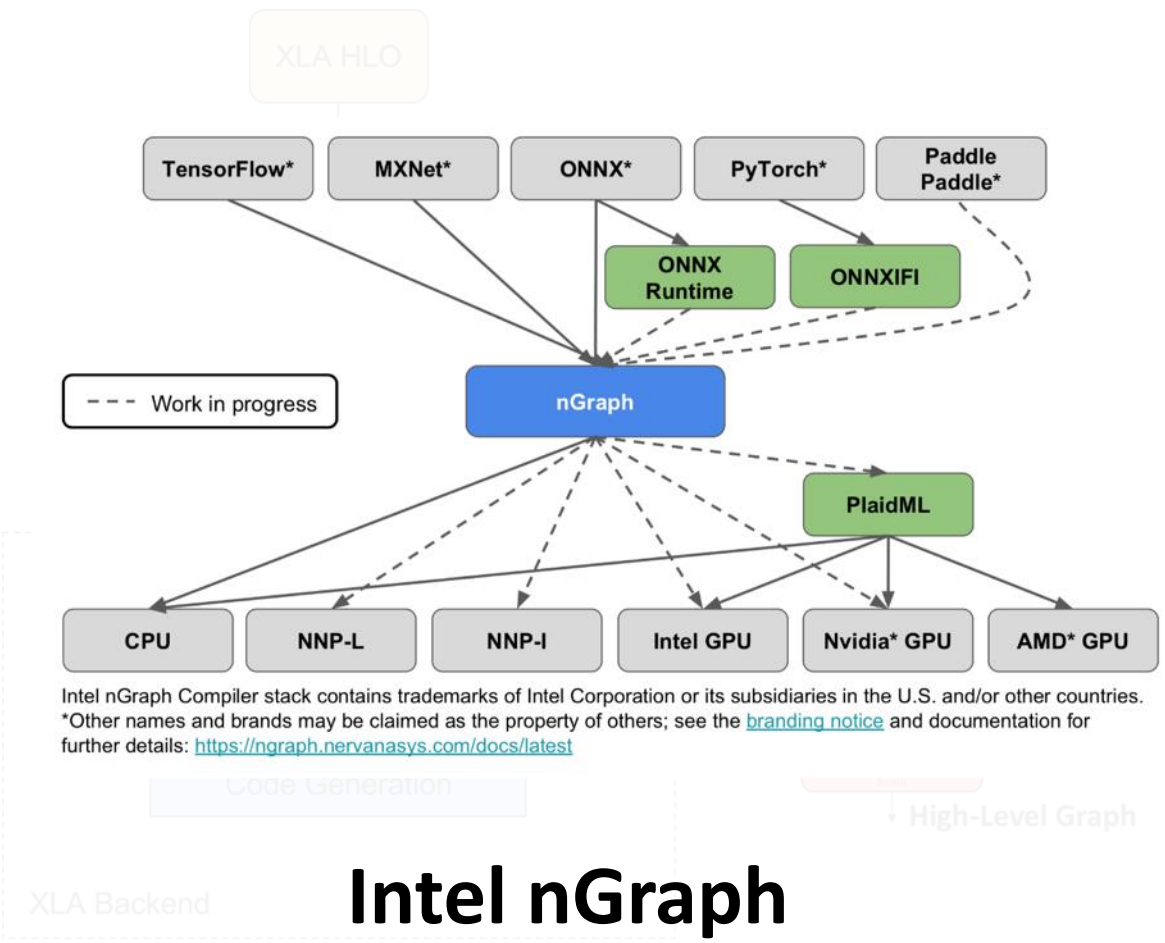
TensorFlow XLA



Facebook Glow

DNN Compilers

- Use techniques from compiler construction: DNN → Graph → IR → Transformations → HW Mapping



Intel nGraph

TensorFlow XLA



TVM Stack

Facebook Glow

How to not do this

How to not do this

“Twelve ways to fool the masses when reporting performance of deep learning workloads”

(A humorous guide to floptimize deep learning)

<https://hlor.inf.ethz.ch/blog/index.php/2018/11/08/twelve-ways-to-fool-the-masses-when-reporting-performance-of-deep-learning-workloads/>

How to not do this

“Twelve ways to fool the masses when reporting performance of deep learning workloads”

(A humorous guide to floptimize deep learning)

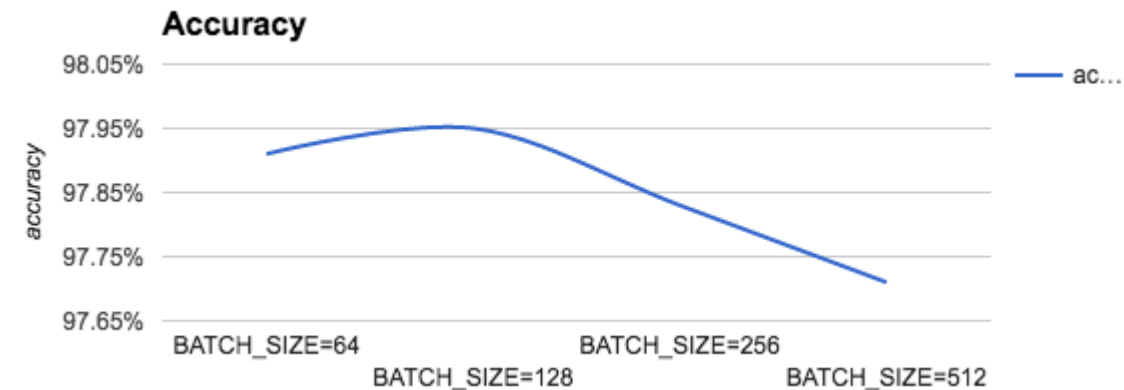
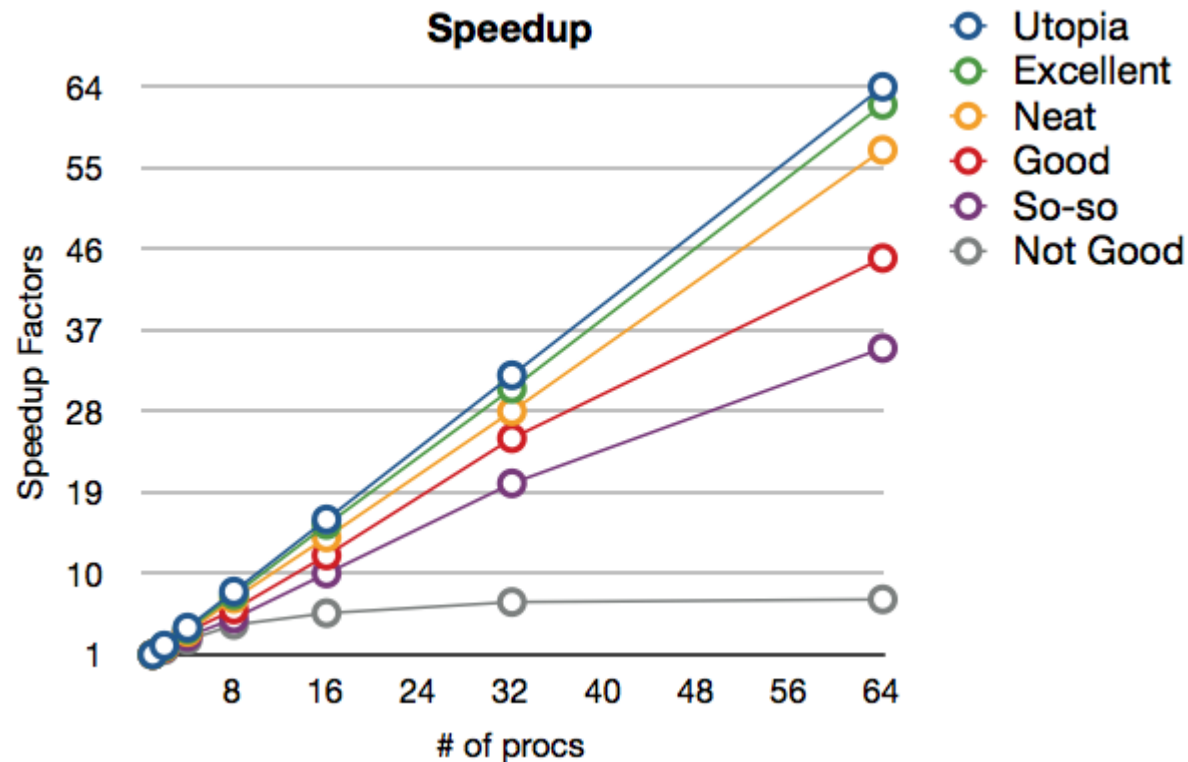
<https://hlor.inf.ethz.ch/blog/index.php/2018/11/08/twelve-ways-to-fool-the-masses-when-reporting-performance-of-deep-learning-workloads/>



8) Show performance when enabling option set A and show accuracy when enabling option set B!

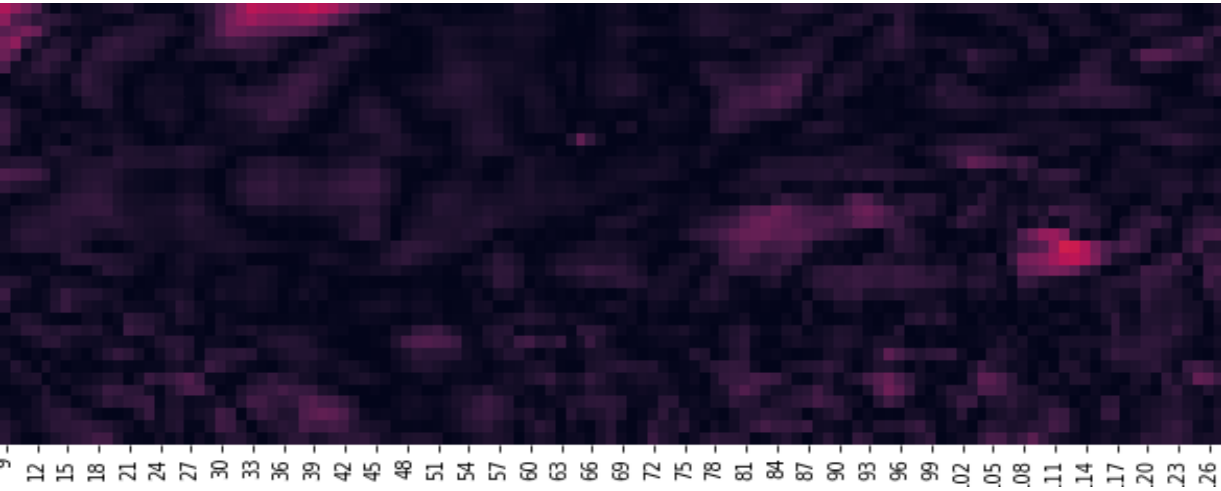
- Pretty cool idea isn't it? Hyperparameters sometimes conflict

So always tune the to show the best result, whatever the result shall be!



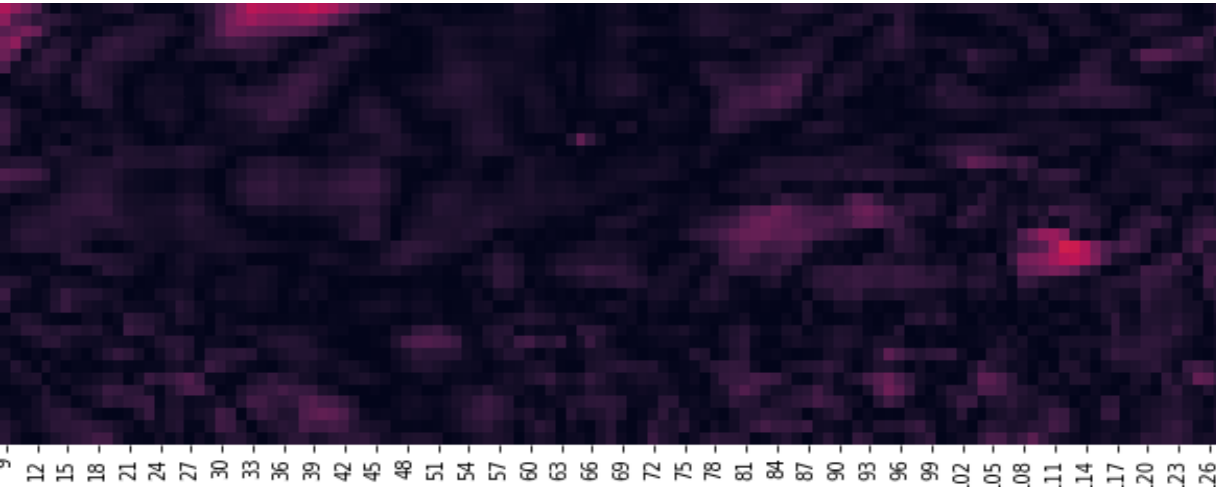
Some big deep learning applications

Some big deep learning applications



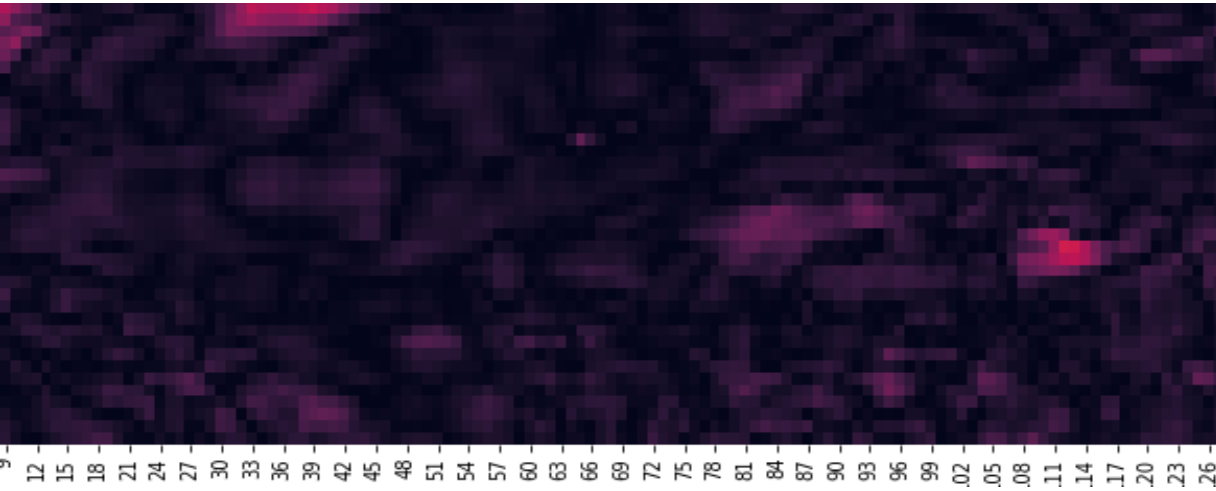
Some big deep learning applications

UQ for weather prediction [Grönquist et al. 2019]

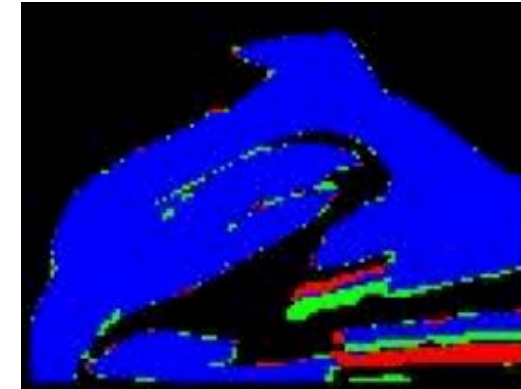


Some big deep learning applications

UQ for weather prediction [Grönquist et al. 2019]

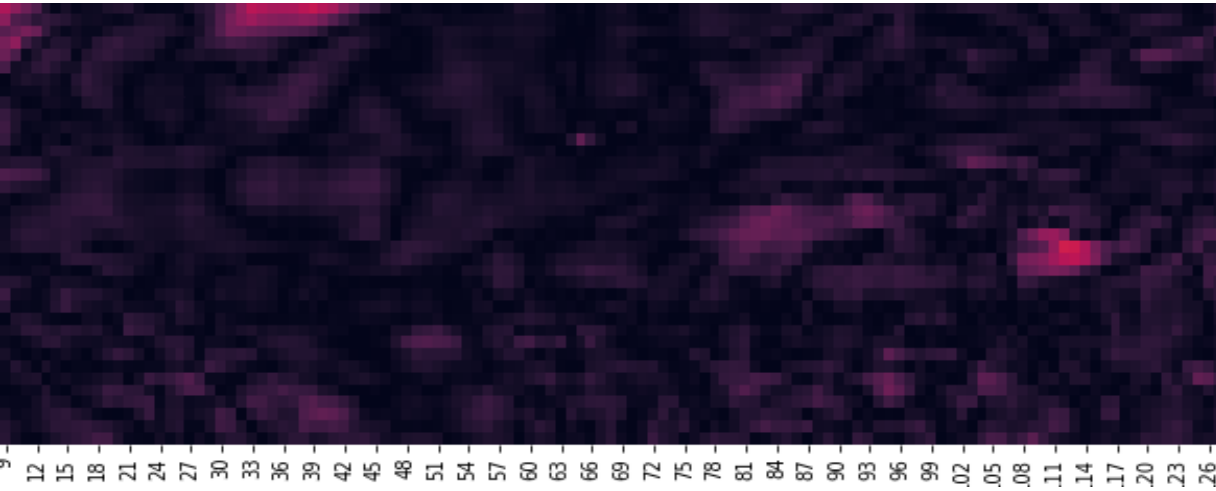


Predicting mesh tangling [Dryden et al. 2019]

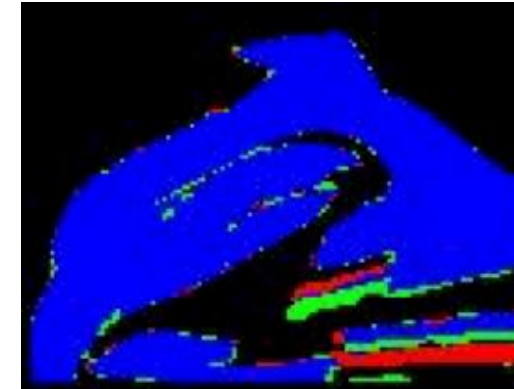


Some big deep learning applications

UQ for weather prediction [Grönquist et al. 2019]

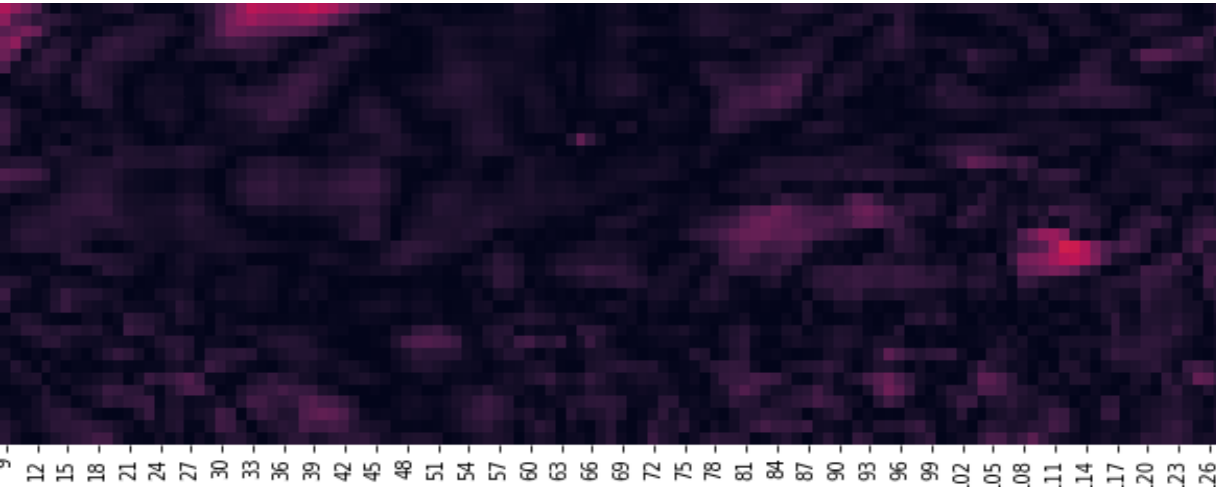


Predicting mesh tangling [Dryden et al. 2019]

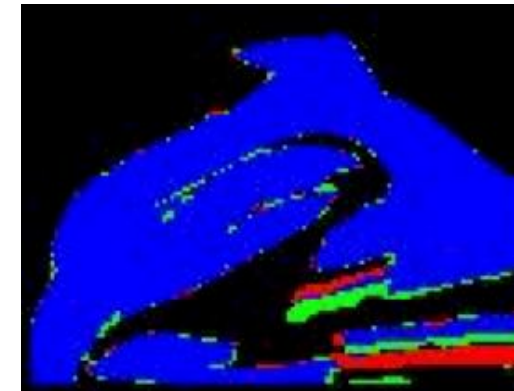


Some big deep learning applications

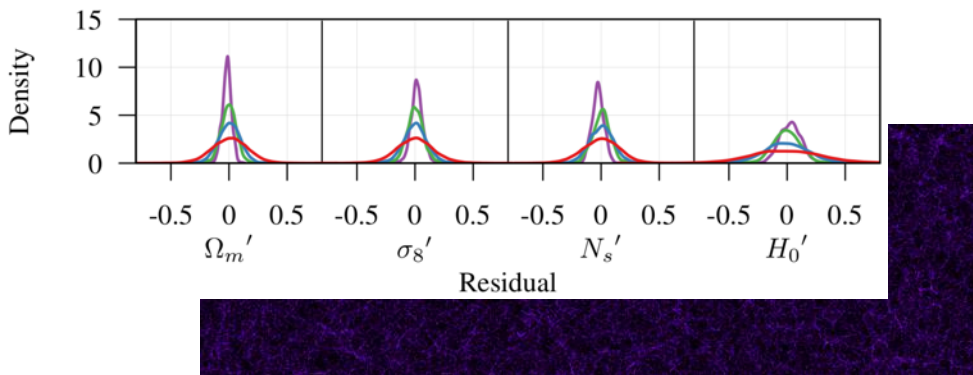
UQ for weather prediction [Grönquist et al. 2019]



Predicting mesh tangling [Dryden et al. 2019]

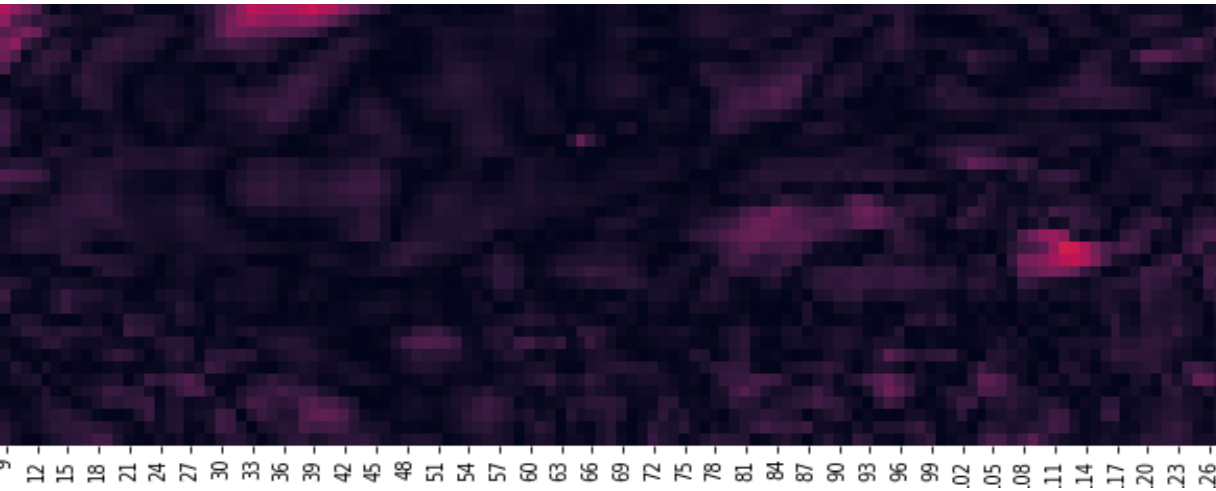


Cosmology [Oyama et al. 2019]

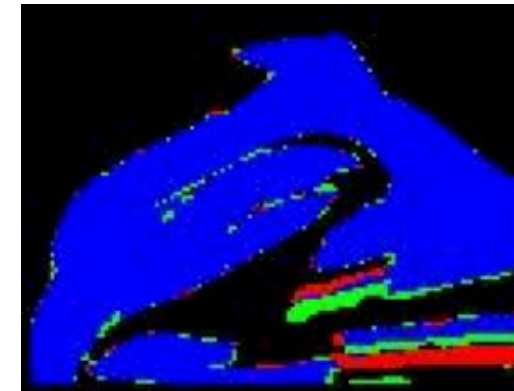


Some big deep learning applications

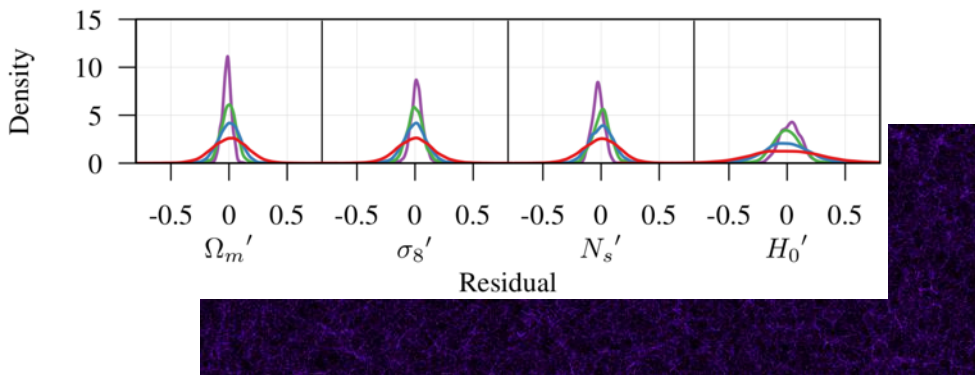
UQ for weather prediction [Grönquist et al. 2019]



Predicting mesh tangling [Dryden et al. 2019]

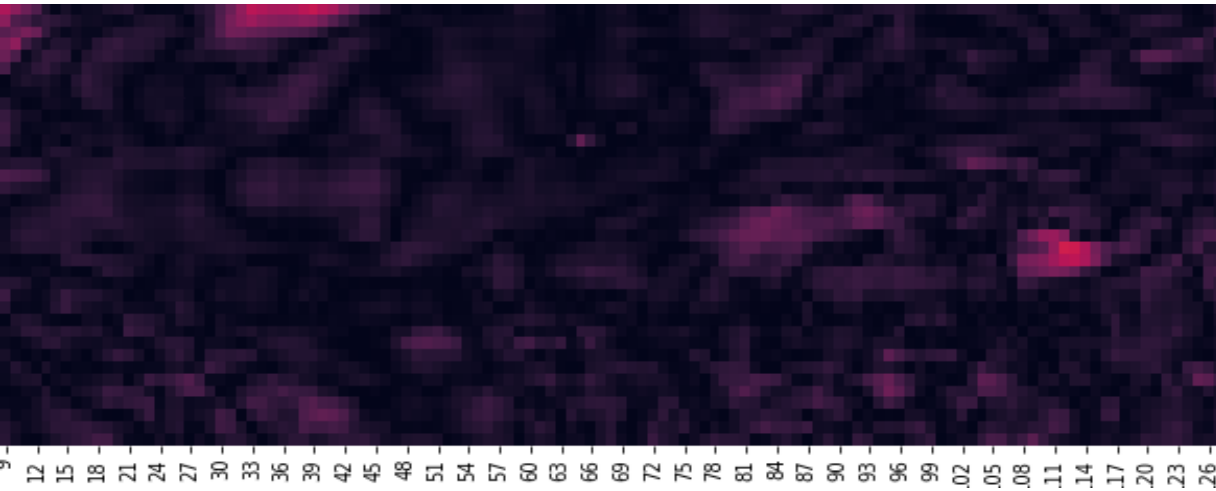


Cosmology [Oyama et al. 2019]

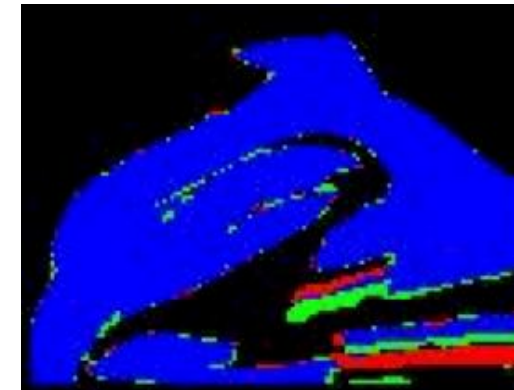


Some big deep learning applications

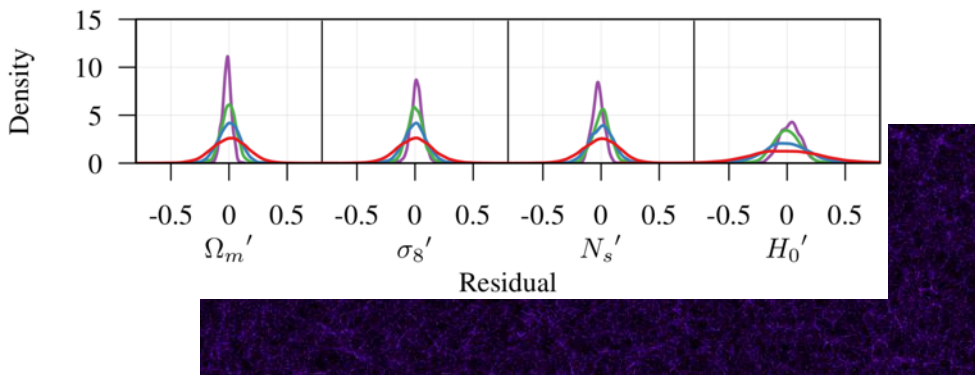
UQ for weather prediction [Grönquist et al. 2019]



Predicting mesh tangling [Dryden et al. 2019]

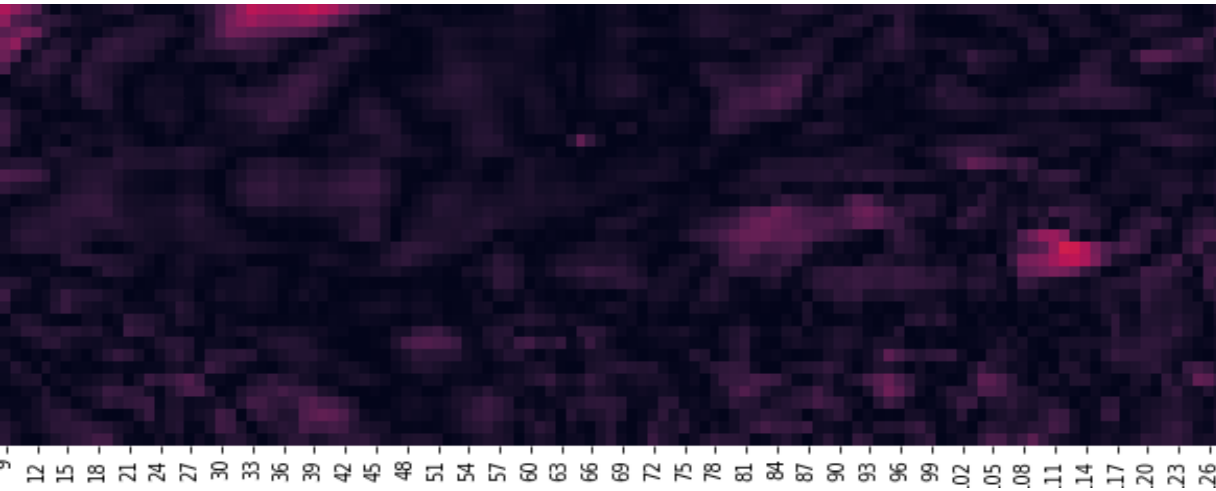


Cosmology [Oyama et al. 2019]

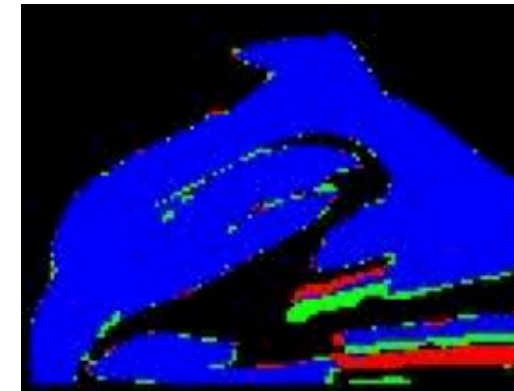


Some big deep learning applications

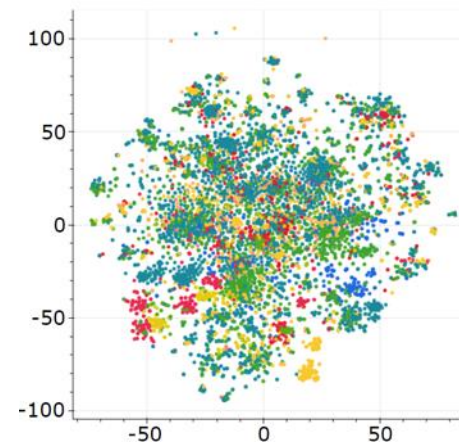
UQ for weather prediction [Grönquist et al. 2019]



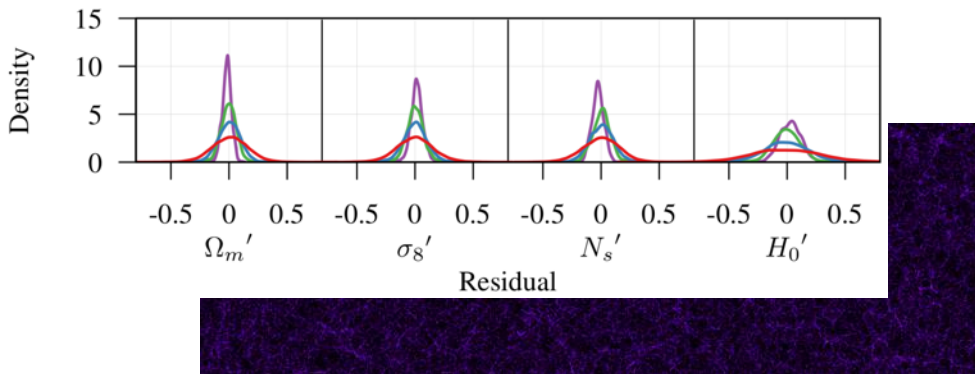
Predicting mesh tangling [Dryden et al. 2019]



Code comprehension [Ben-Nun et al. 2018]

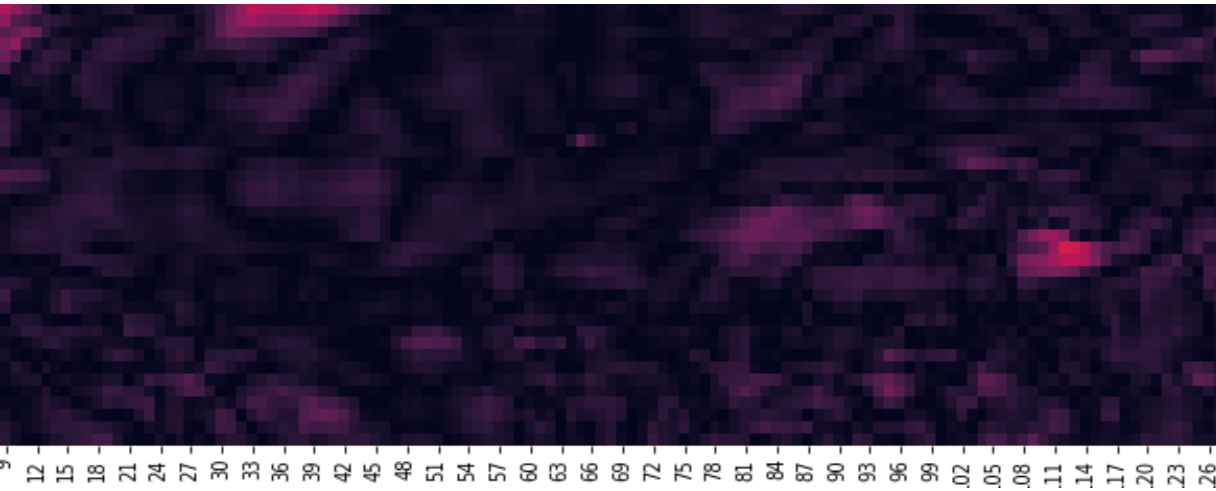


Cosmology [Oyama et al. 2019]

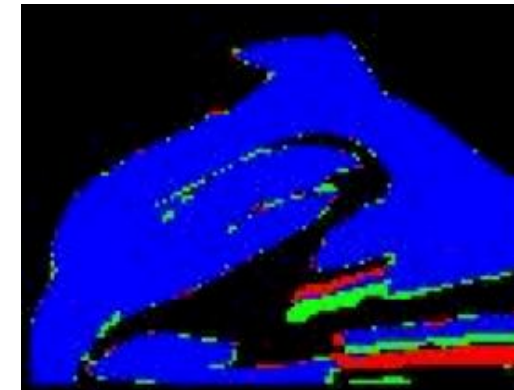


Some big deep learning applications

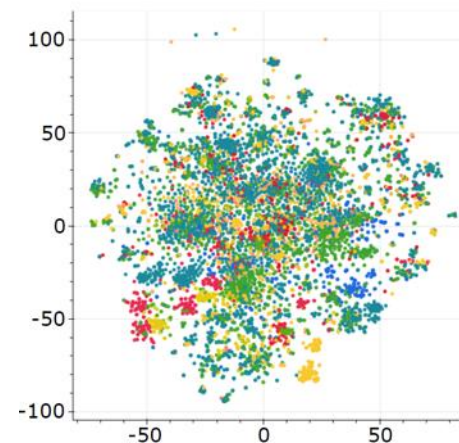
UQ for weather prediction [Grönquist et al. 2019]



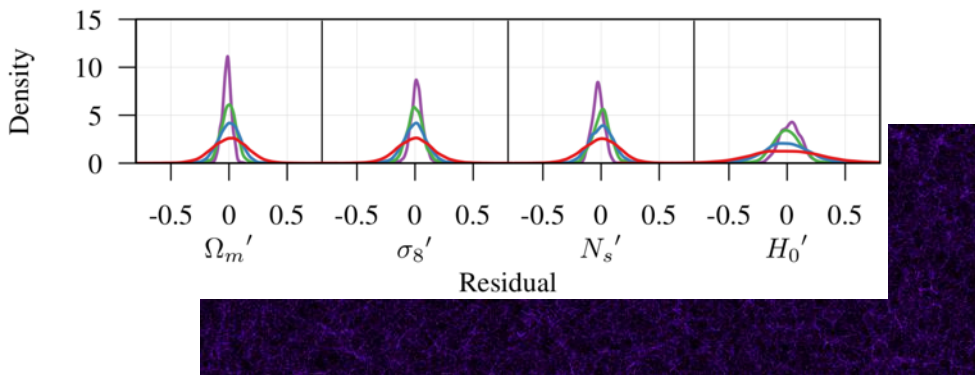
Predicting mesh tangling [Dryden et al. 2019]



Code comprehension [Ben-Nun et al. 2018]

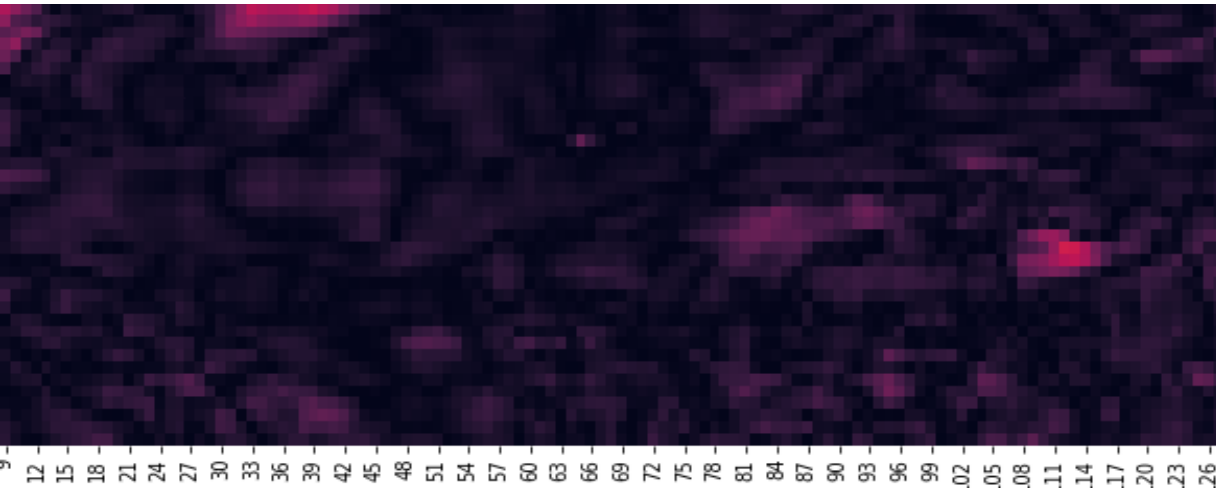


Cosmology [Oyama et al. 2019]

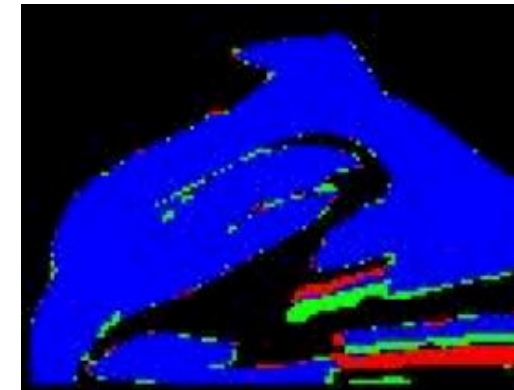


Some big deep learning applications

UQ for weather prediction [Grönquist et al. 2019]

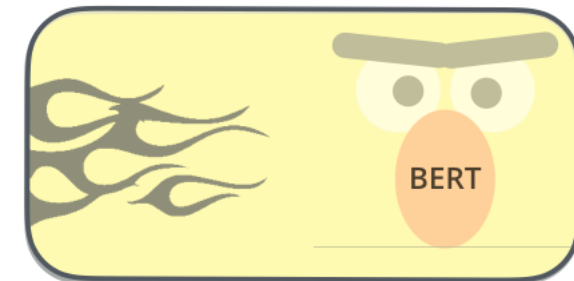
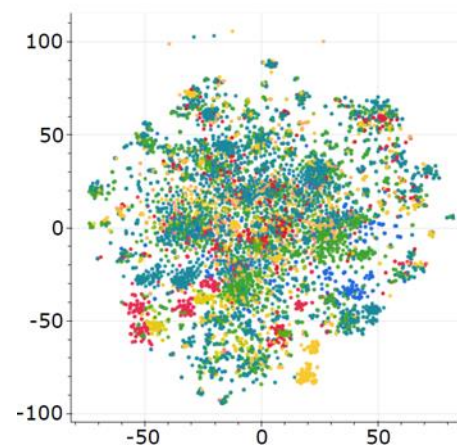
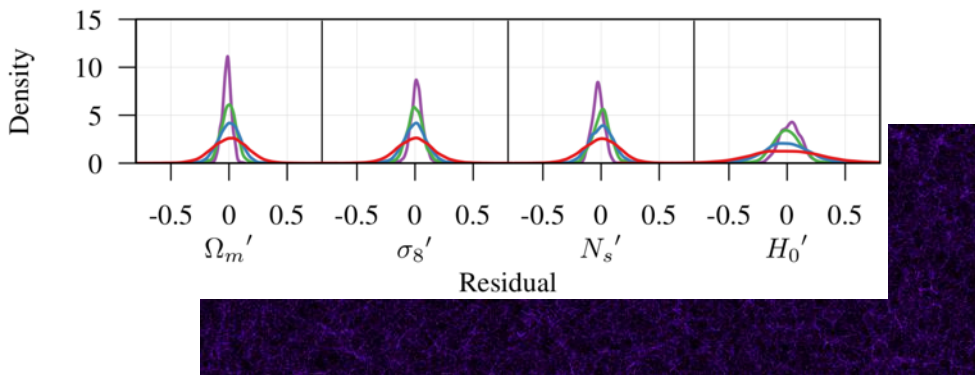


Predicting mesh tangling [Dryden et al. 2019]



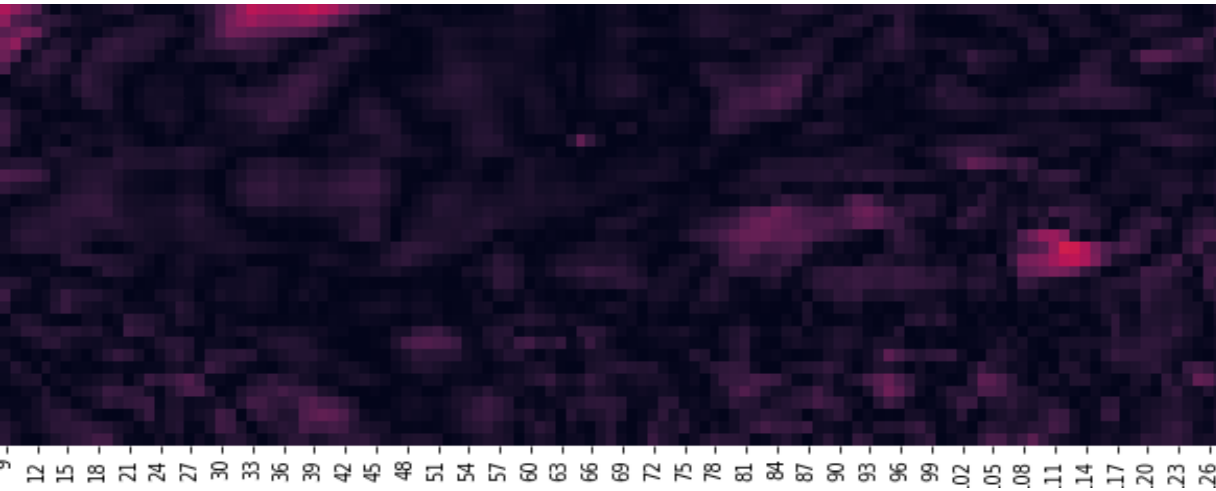
Code comprehension [Ben-Nun et al. 2018]

Cosmology [Oyama et al. 2019]

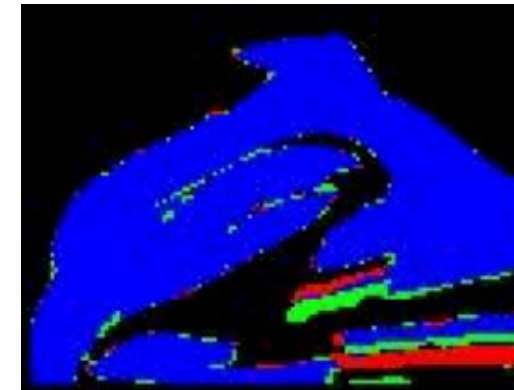


Some big deep learning applications

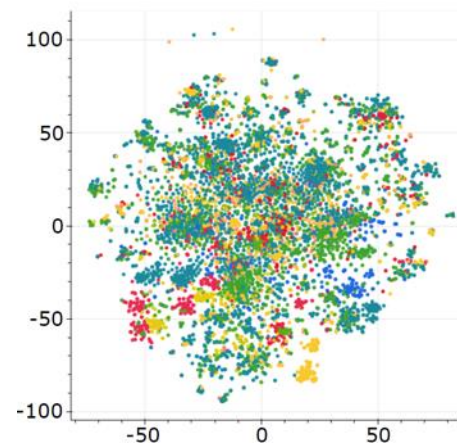
UQ for weather prediction [Grönquist et al. 2019]



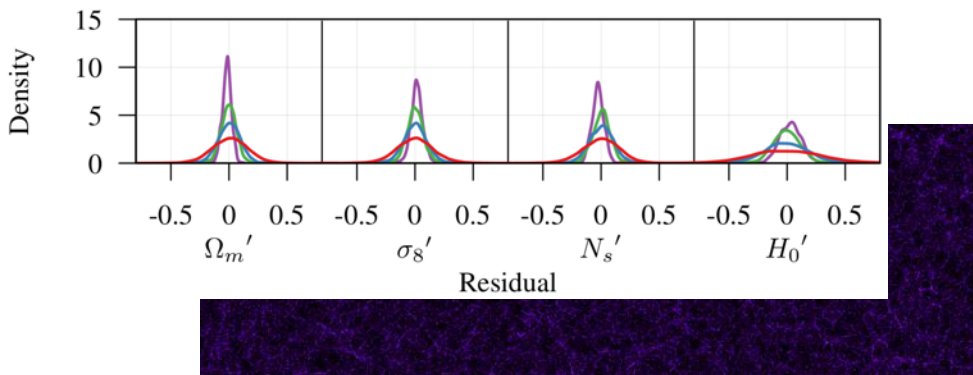
Predicting mesh tangling [Dryden et al. 2019]



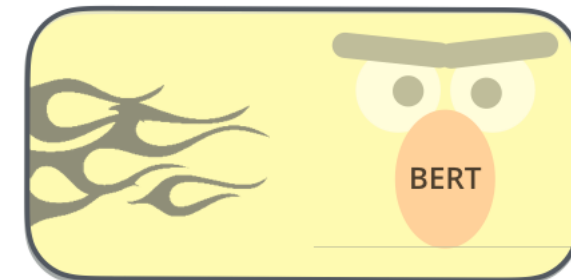
Code comprehension [Ben-Nun et al. 2018]



Cosmology [Oyama et al. 2019]

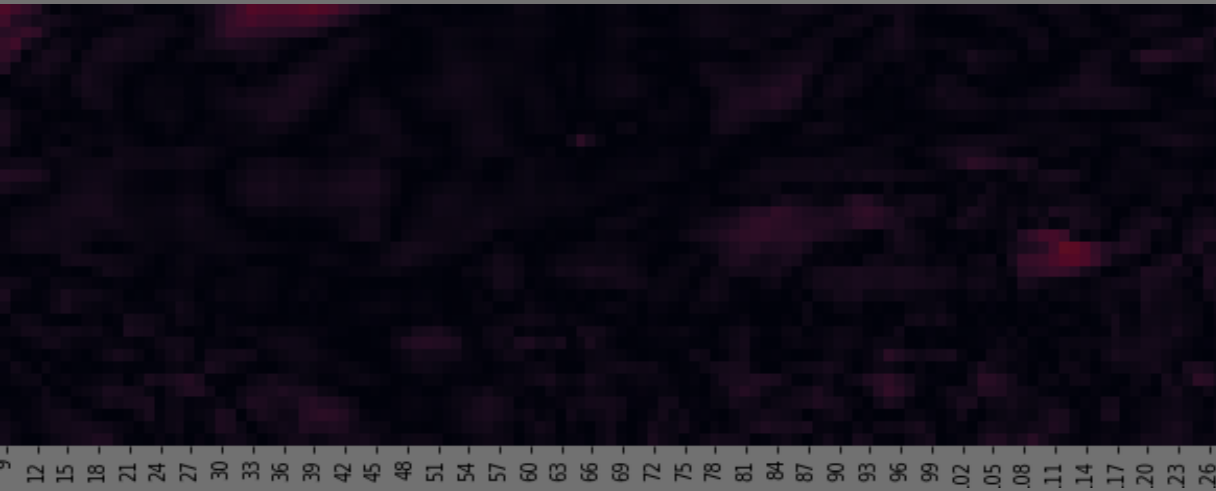


Big seq2seq models

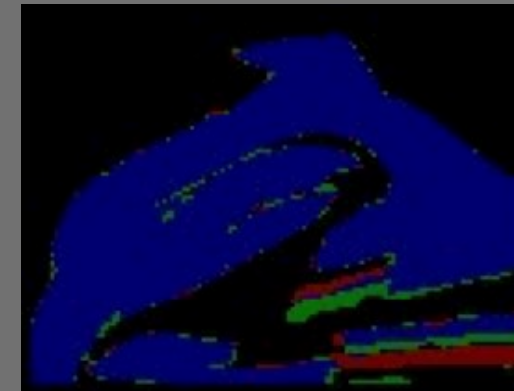


Some big deep learning applications

UQ for weather prediction [Grönquist et al. 2019]

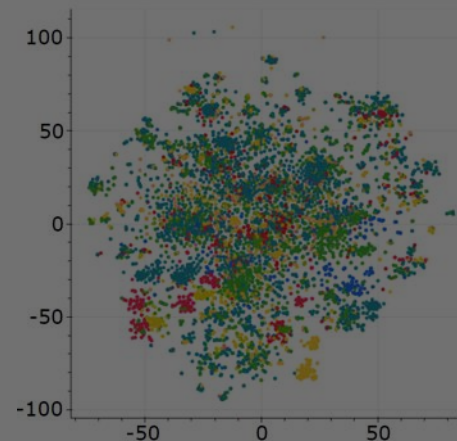
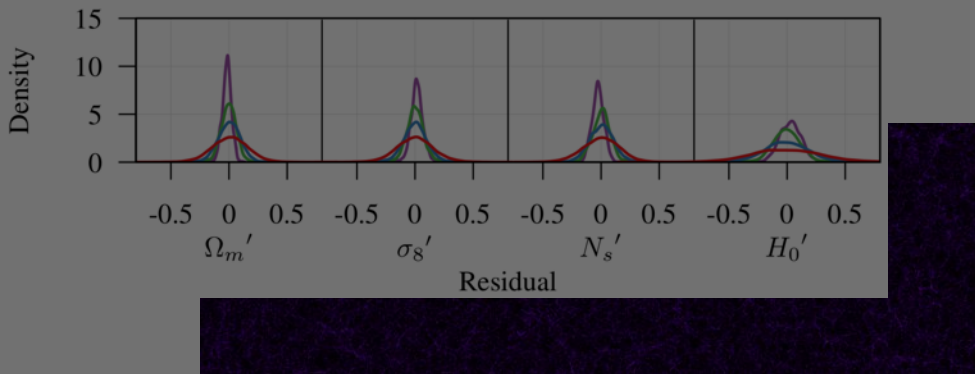


Predicting mesh tangling [Dryden et al. 2019]

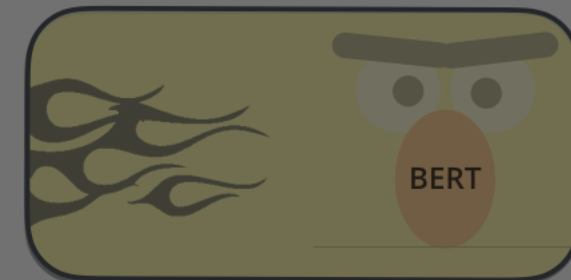


Code comprehension [Ben-Nun et al. 2018]

Cosmology [Oyama et al. 2019]



Big seq2seq models

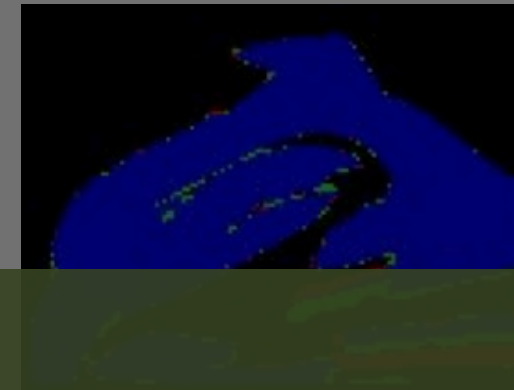


Some big deep learning applications

UQ for weather prediction [Grönquist et al. 2019]



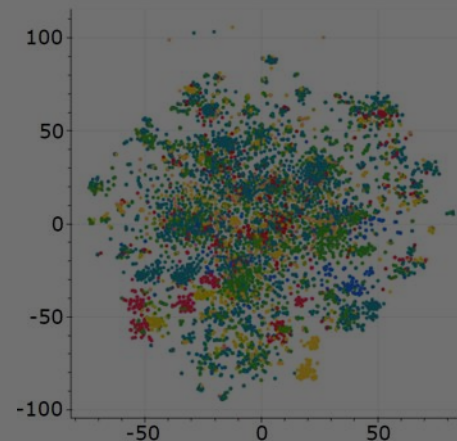
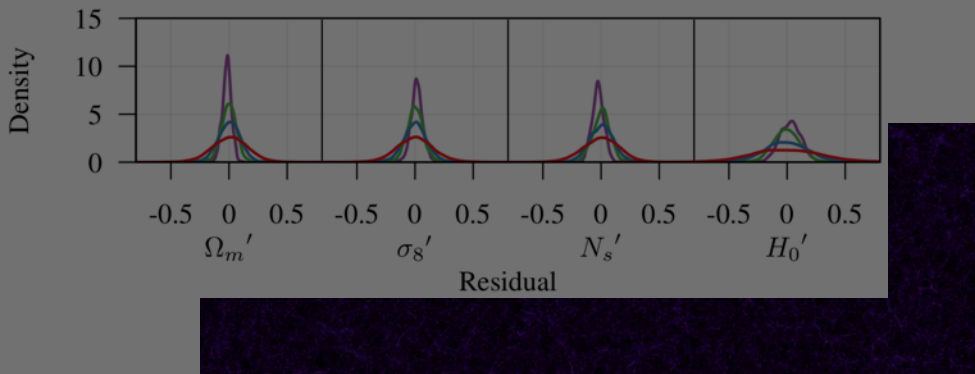
Predicting mesh tangling [Dryden et al. 2019]



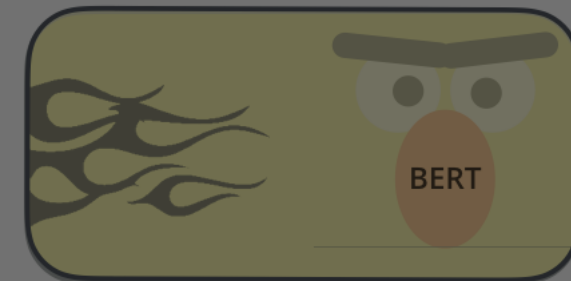
And many more!

Code comprehension [Ben-Nun et al. 2018]

Cosmology [Oyama et al. 2019]



Big seq2seq models



Research opportunities

<https://spcl.inf.ethz.ch/SeMa/>
 ndryden@ethz.ch

Project Name	Category
+ Efficient partial collective operations for distributed deep learning training	Parallel Algorithms
+ Clairvoyant Prefetching for Machine Learning I/O	Machine Learning
+ Transformers: More than Meets the Eye (of the Hurricane)	Machine Learning
+ Scalable Deep Learning for Weather and Climate Prediction	Machine Learning
+ Fastest Matrix Multiplication in the West (Europe)	Parallel Algorithms
+ Efficient Collective Operations On Reconfigurable Hardware	Architecture
+ Quantized Allreduces for Distributed Deep Learning Training	Machine Learning
+ Who Optimizes the Optimizers? Performance Programming Made Easy	Toolchains
+ Analytical Cache Model for Parallel Programs	Compilation
+ Automatic Algorithm Detection for Readability and Performance Rewriting	Compilation
+ Data-Centric Deep Learning Framework (or: how to beat TensorFlow)	Machine Learning
+ Authenticated Deep Learning	Machine Learning
+ Kernel/Network Offloading of Streaming Data Processing Tasks	Networking
+ Performance Counters for Interactive Bottleneck Identification in Large-Scale Applications	Compilation
+ Automatic Learning of GPU Code Generation Parameters	Compilation
+ Array Partitioning to Exploit High-Bandwidth Block-RAM on FPGA	Compilation
+ Proofably optimal loop scheduling using MINLP	Compilation
+ Cache as RAM to accelerate x86 computations	Compilation
+ Visualization Techniques for Performance-Guided Programming	Compilation
+ Large Scale Framework for Code Analysis	Compilation
+ Deep Learning for Large-Scale Graph Analytics	Machine Learning

Parallelism in training DNNs – Summary

Parallelism in training DNNs – Summary

- **Deep learning is HPC – very similar – mainly dense linear algebra**
 - Amenable to our usual set of tricks, sometimes with a twist
- **Main bottleneck is communication – reduction by trading off**

Parameter Consistency

- Bounded synchronous SGD
 - Central vs. distributed parameter server
 - EASGD to ensemble learning
-
- **Strong scaling requires effort**
 - **Very different environment from traditional HPC**
 - Trade-off accuracy for performance!
 - **Performance-centric view in HPC can be harmful for accuracy!**

Parallelism in training DNNs – Summary

- **Deep learning is HPC – very similar – mainly dense linear algebra**
 - Amenable to our usual set of tricks, sometimes with a twist
- **Main bottleneck is communication – reduction by trading off**

Parameter Consistency

- Bounded synchronous SGD
 - Central vs. distributed parameter server
 - EASGD to ensemble learning
- **Strong scaling requires effort**
 - **Very different environment from traditional HPC**
 - Trade-off accuracy for performance!
 - **Performance-centric view in HPC can be harmful for accuracy!**

Parallelism in training DNNs – Summary

- **Deep learning is HPC – very similar – mainly dense linear algebra**
 - Amenable to our usual set of tricks, sometimes with a twist
- **Main bottleneck is communication – reduction by trading off**

Parameter Consistency

- Bounded synchronous SGD
- Central vs. distributed parameter server
- EASGD to ensemble learning

Parameter Accuracy

- Lossless compression of gradient updates
- Quantization of gradient updates
- Sparsification of gradient updates

- **Strong scaling requires effort**
- **Very different environment from traditional HPC**
 - Trade-off accuracy for performance!
- **Performance-centric view in HPC can be harmful for accuracy!**

Parallelism in training DNNs – Summary

- **Deep learning is HPC – very similar – mainly dense linear algebra**
 - Amenable to our usual set of tricks, sometimes with a twist
- **Main bottleneck is communication – reduction by trading off**

Parameter Consistency

- Bounded synchronous SGD
- Central vs. distributed parameter server
- EASGD to ensemble learning

Parameter Accuracy

- Lossless compression of gradient updates
- Quantization of gradient updates
- Sparsification of gradient updates

- **Strong scaling requires effort**
- **Very different environment from traditional HPC**
 - Trade-off accuracy for performance!
- **Performance-centric view in HPC can be harmful for accuracy!**

Parallelism in training DNNs – Summary

- **Deep learning is HPC – very similar – mainly dense linear algebra**
 - Amenable to our usual set of tricks, sometimes with a twist
- **Main bottleneck is communication – reduction by trading off**

Parameter Consistency

- Bounded synchronous SGD
- Central vs. distributed parameter server
- EASGD to ensemble learning

Parameter Accuracy

- Lossless compression of gradient updates
- Quantization of gradient updates
- Sparsification of gradient updates

- **Strong scaling requires effort**
- **Very different environment from traditional HPC**
 - Trade-off accuracy for performance!
- **Performance-centric view in HPC can be harmful for accuracy!**