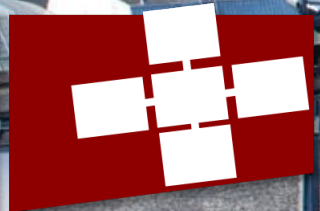


Oliver Rausch\*, Tal Ben-Nun\*, Nikoli Dryden, Andrei Ivanov, Shigang Li, Torsten Hoefler

# DaCeML: A Data-Centric Optimization Framework for Machine Learning



This project received funding from the European Research Council under the European Union's Horizon 2020 programme (Project PSAP, No. 101002047); and receives EuroHPC-JU funding under grants MAELSTROM, No. 955513 and DEEP-SEA, No. 955606, with support from the Horizon 2020 programme







# Contemporary ML Systems are Inflexible

- Researchers demand efficient, large-scale compute more than ever
- Almost all optimizations relate to the data movement
- But current compilers are highly specialized towards
  - Operators
  - Transformations
  - Models
- Tuning, or even just visualizing the output of a compiler is difficult

# DaCeML

- Data-centric lowering and multi-level optimization of DNNs

Usability

Generality

Interactivity

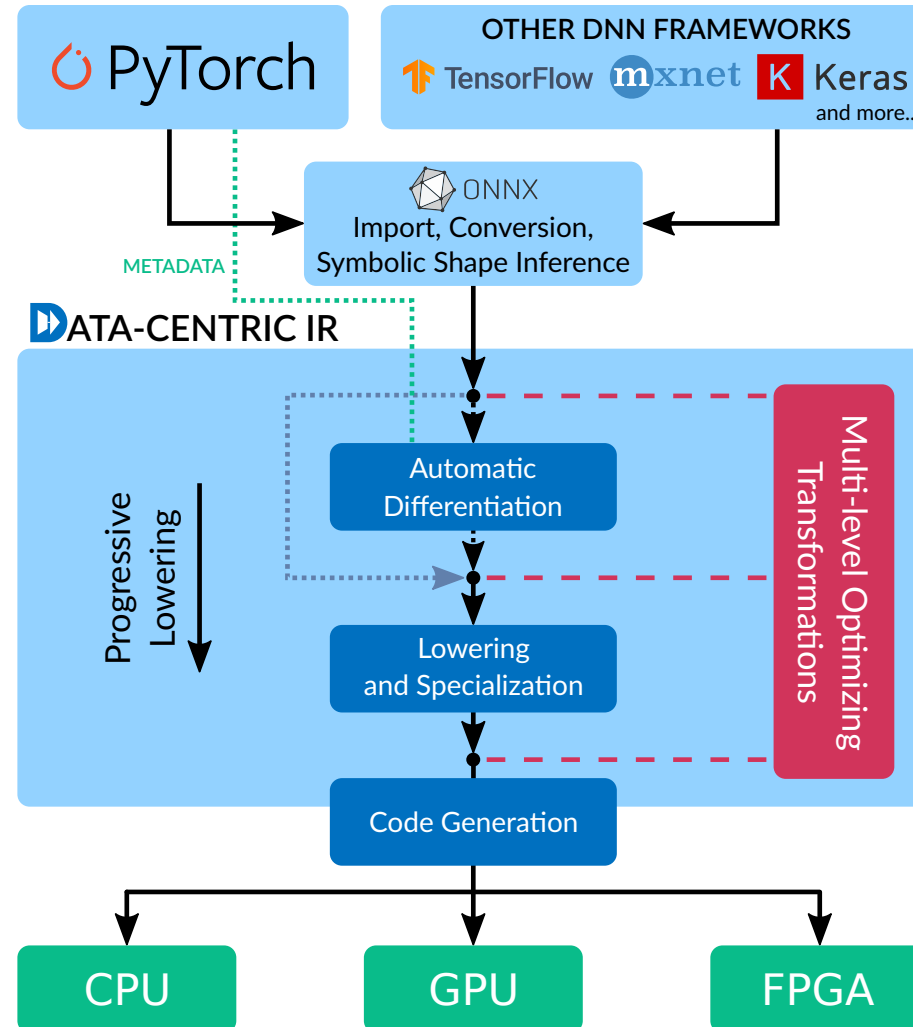
# Usage

```
from torch import nn
from daceml.pytorch import dace_module

@dace_module
class MyModule(nn.Module):
    def __init__(self, n_in, n_out):
        super().__init__()
        self.linear = nn.Linear(n_in, n_out)
        self.fanout = n_out

    def forward(self, x):
        return self.linear(x) / self.fanout
```

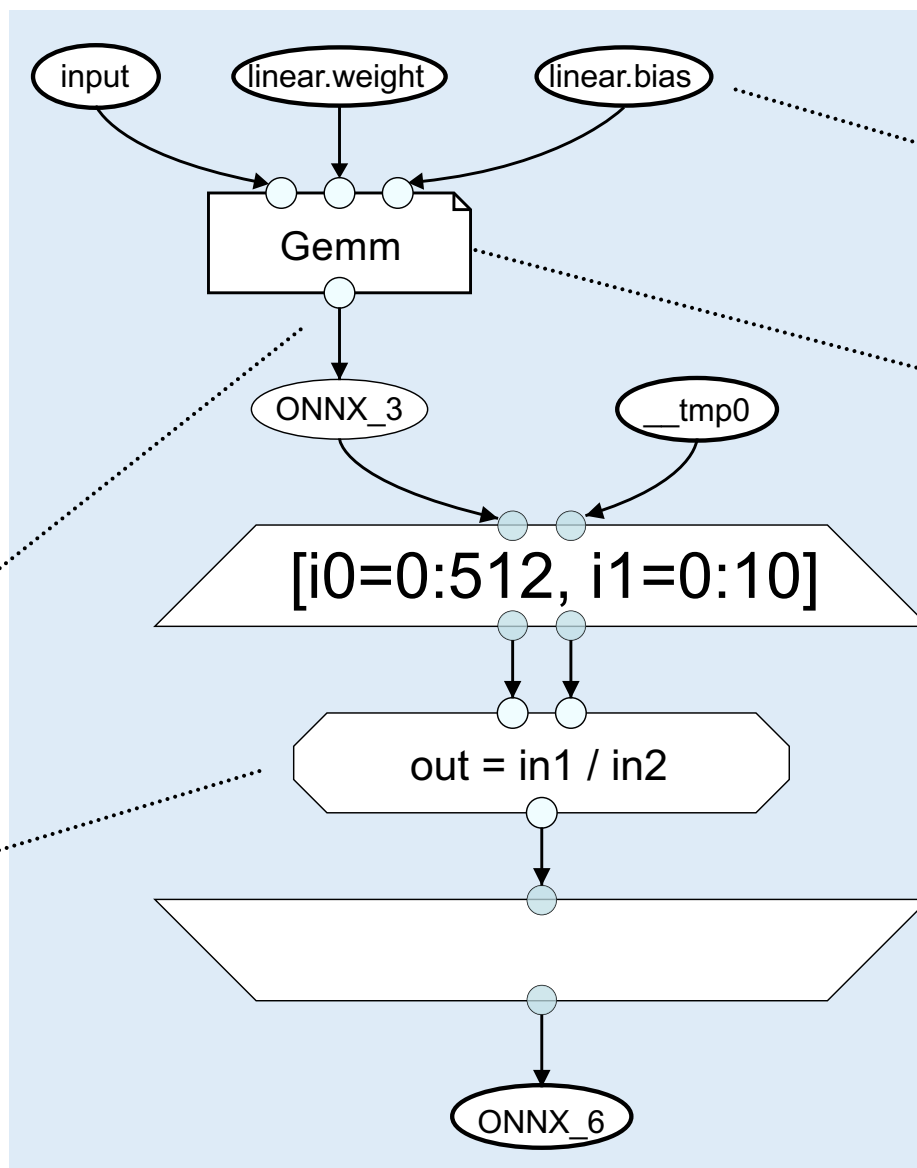
# System Overview



```
from torch import nn
from daceml.pytorch import dace_module
```

```
@dace_module
class MyModule(nn.Module):
    def __init__(self, n_in, n_out):
        super().__init__()
        self.linear = nn.Linear(n_in, n_out)
        self.fanout = n_out

    def forward(self, x):
        return self.linear(x) / self.fanout
```



**Access nodes** to data containers

**Library nodes** with domain-specific semantics

**Maps:** Parametric parallelism scopes

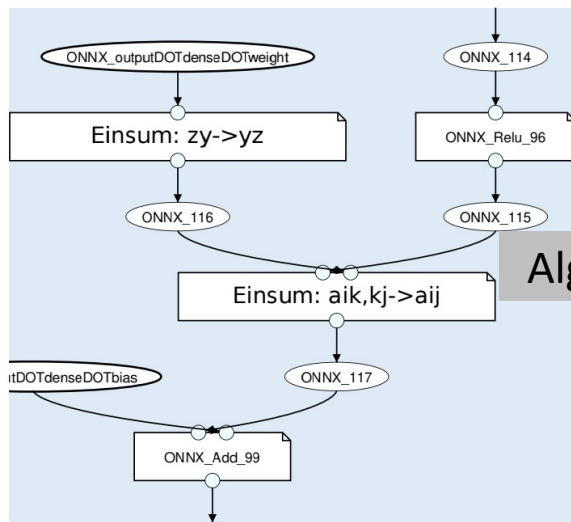
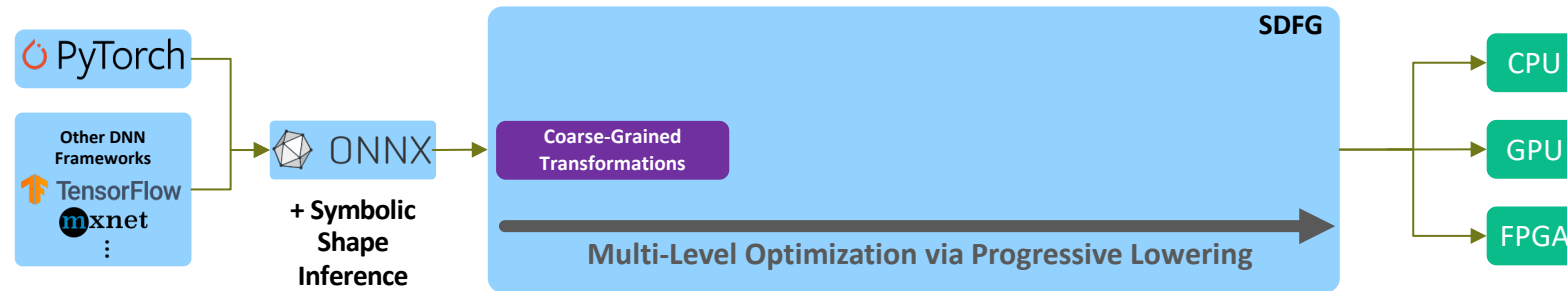
**Memlets:** explicit data movement at all granularities

**Tasklets:** stateless computations

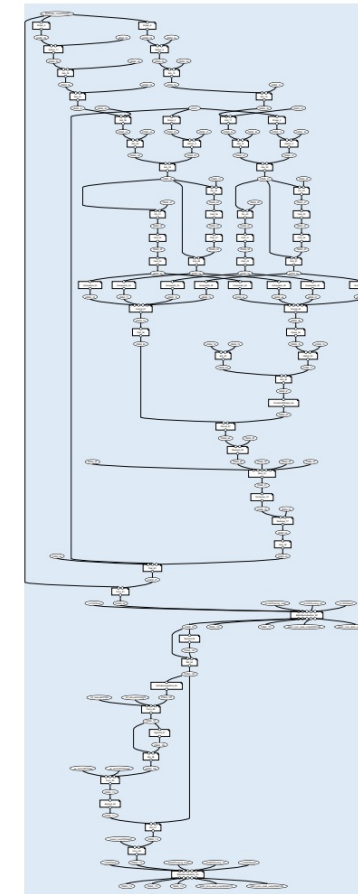
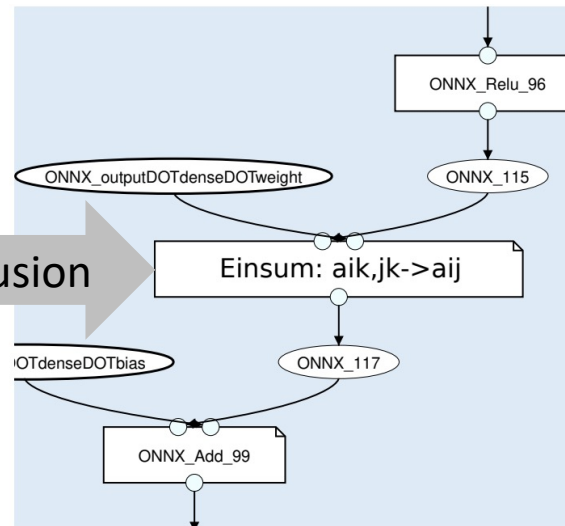
# Optimization Pipeline



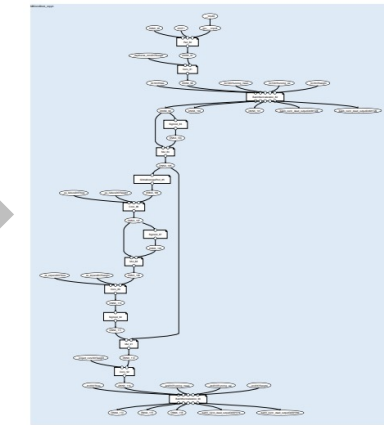


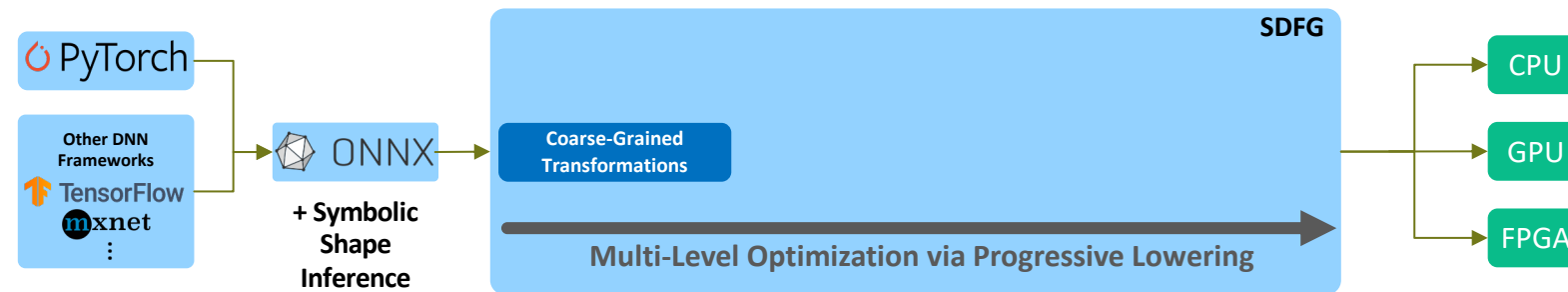


Algebraic Fusion



Folding

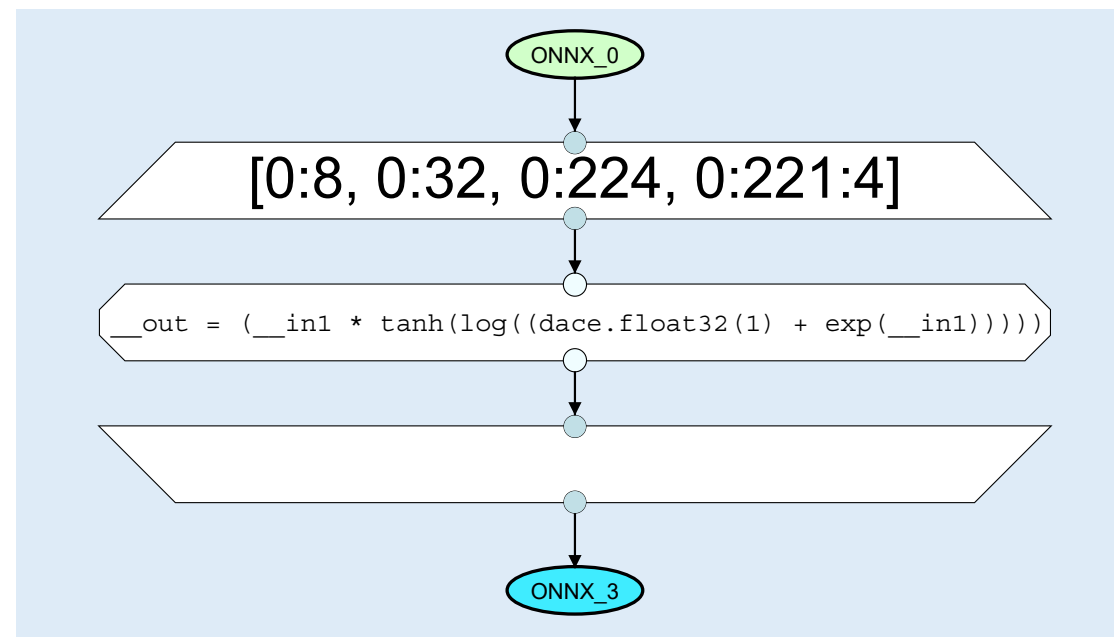
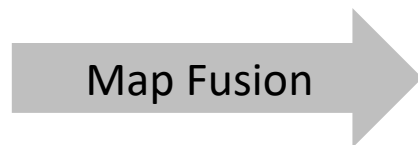
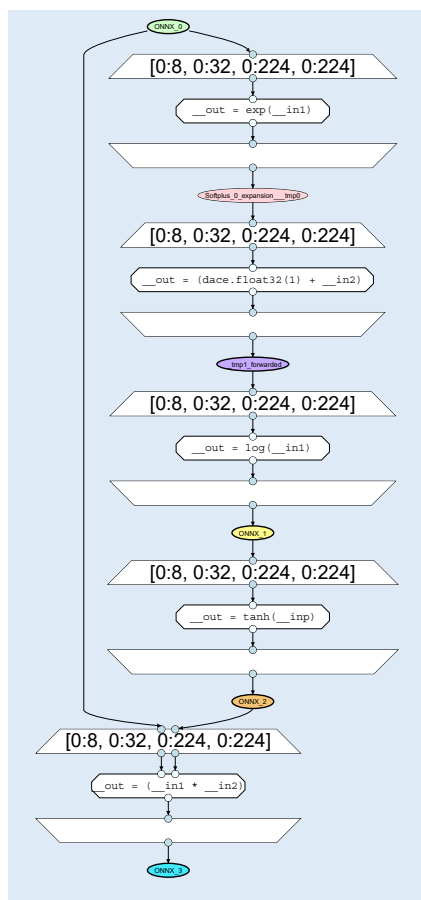
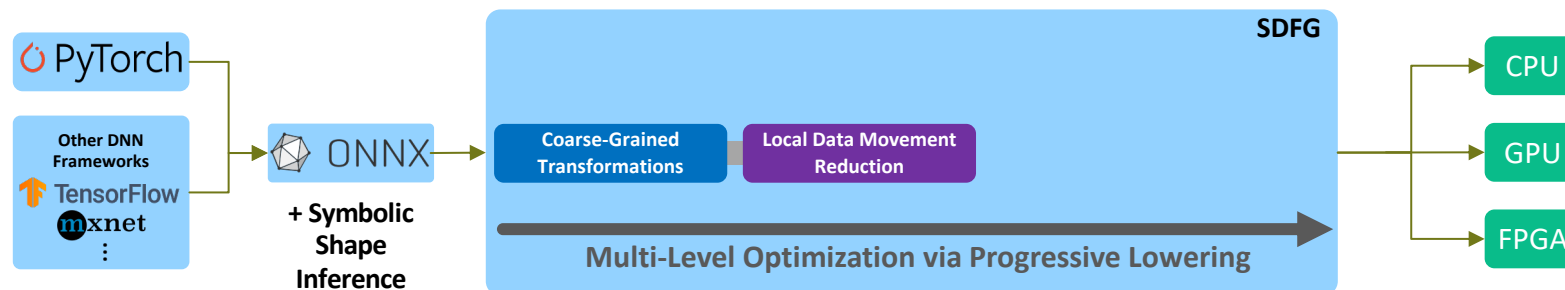


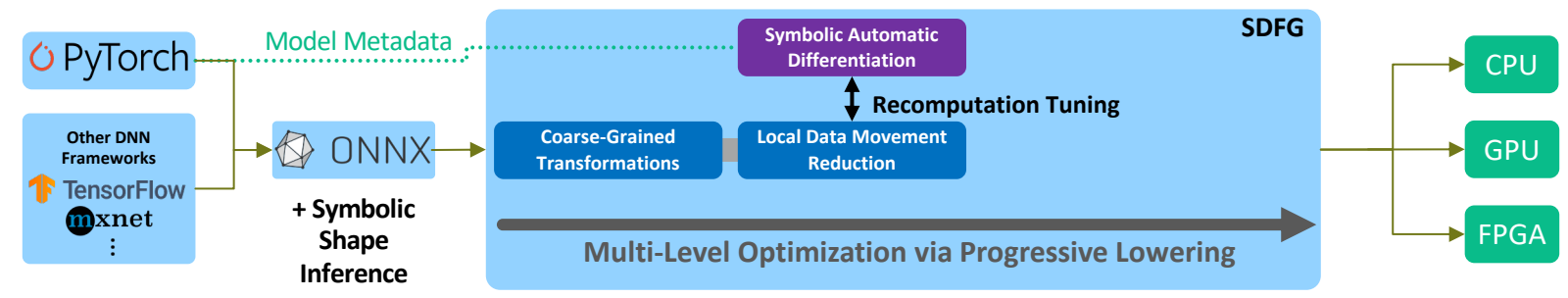


## Lowering ONNX nodes

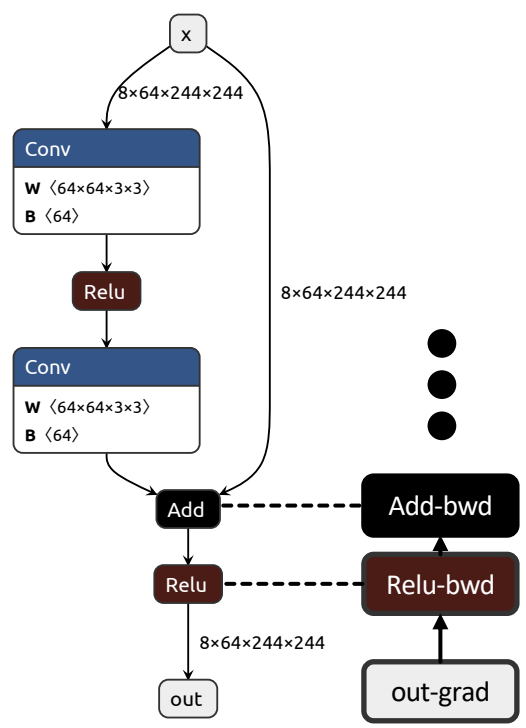
```
@python_pure_op_implementation  
def Softplus(X, Y):  
    Y[:] = numpy.log(1 + numpy.exp(X))
```

```
def forward(self, x):
    y = F.softplus(x)
    y = torch.tanh(y)
    return x * y
```

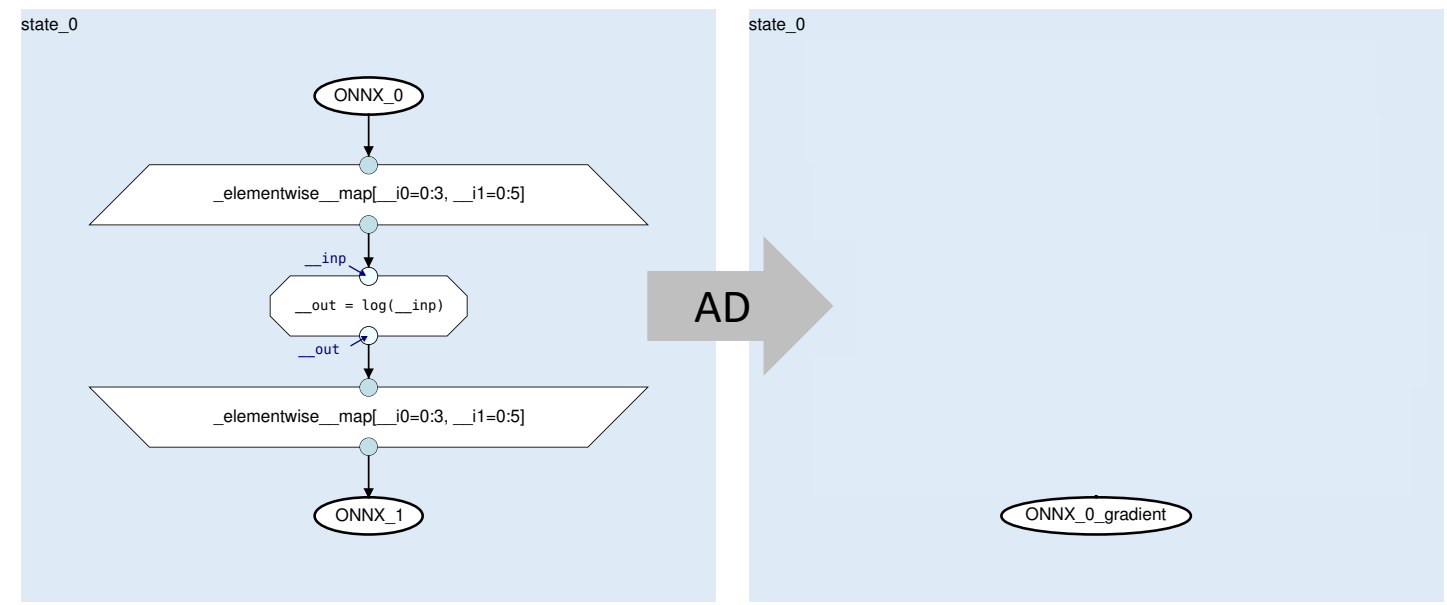




### Operator-Level AD

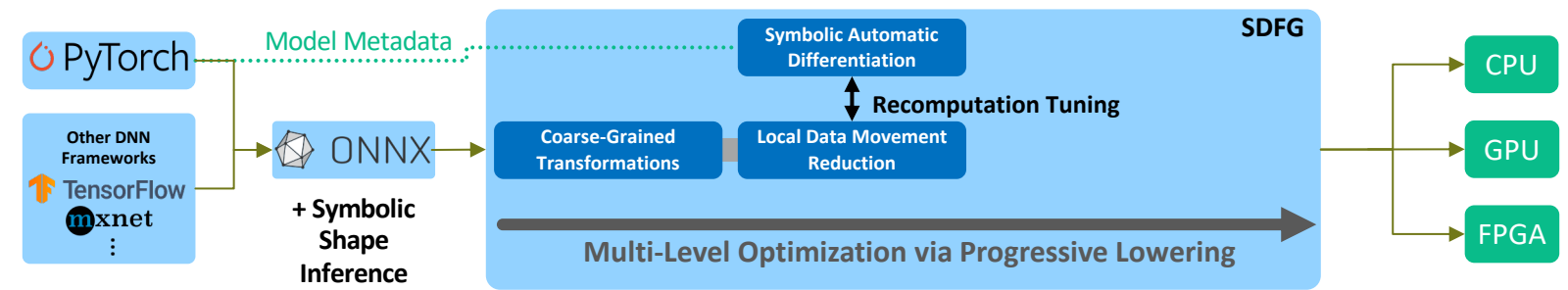


### DaCeML: data-centric, Symbolic AD

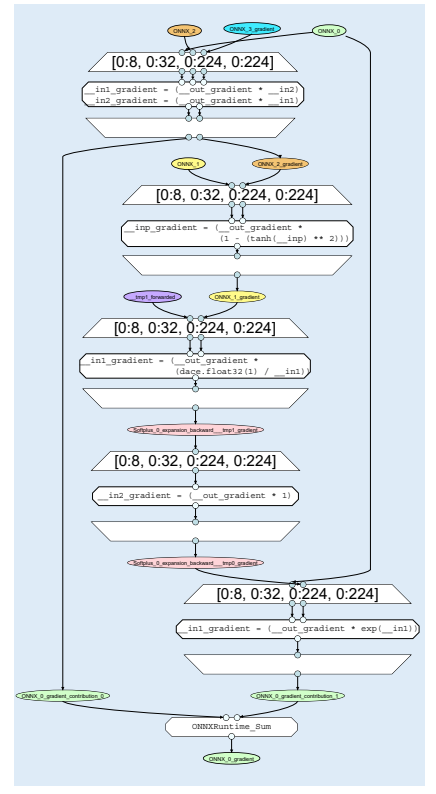


Data-movement information enables automatic backward-kernel synthesis!



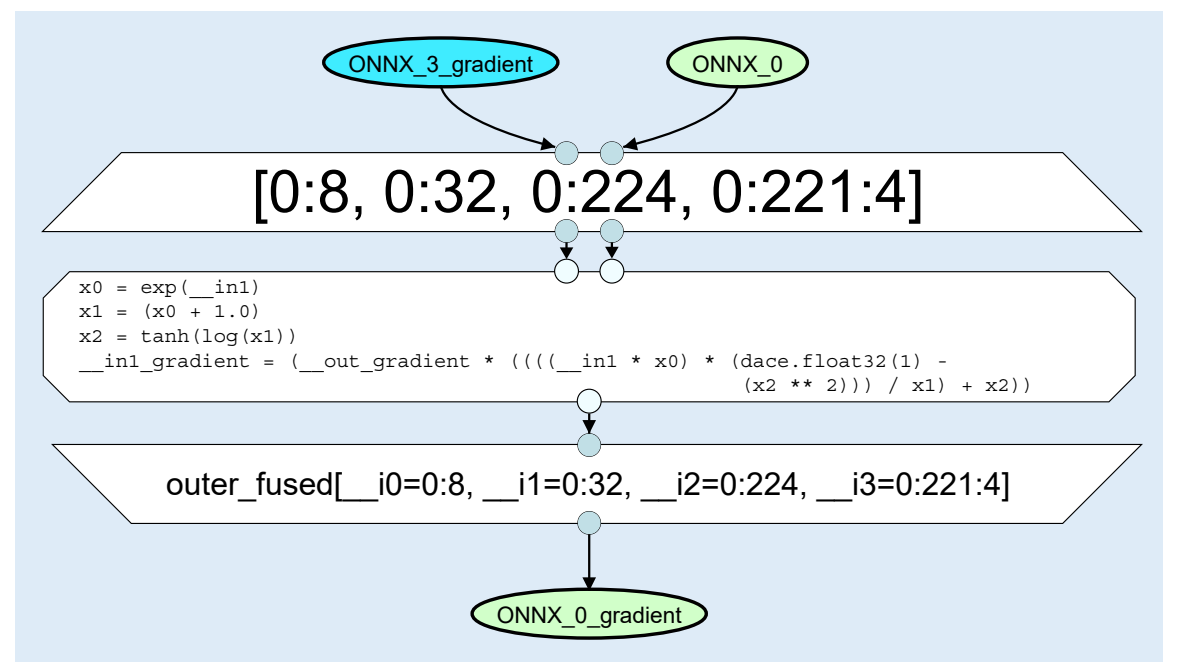


### Post-AD Optimization



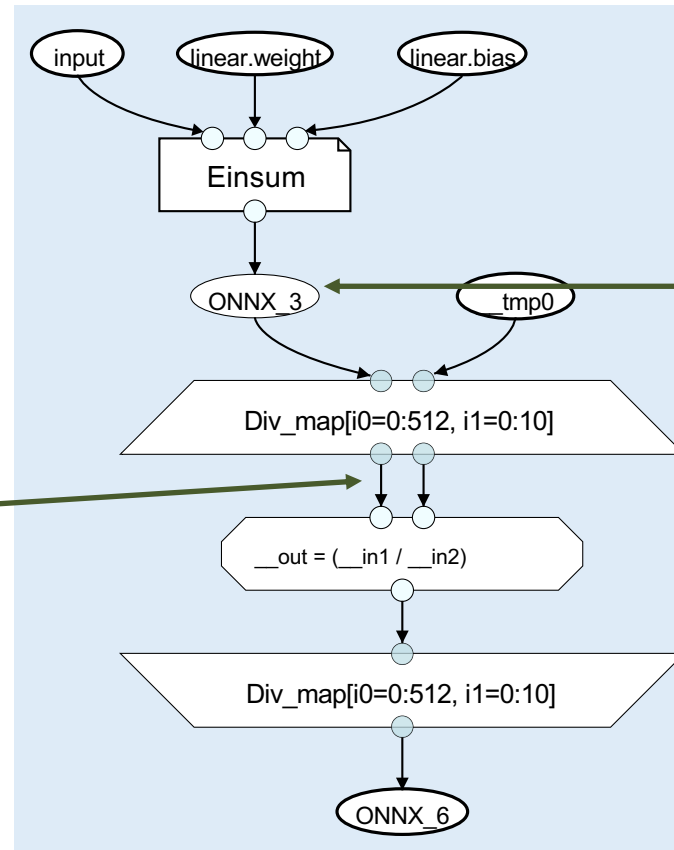
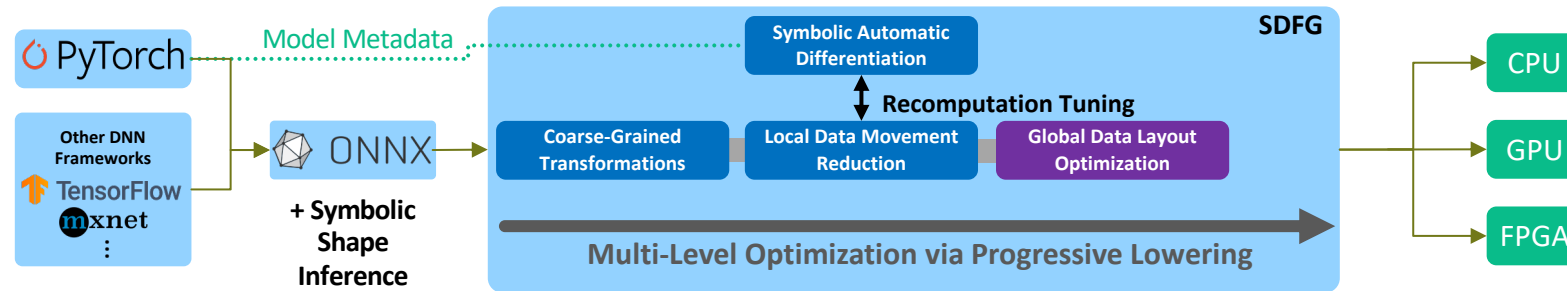
Intermediate Memory: 134.75MiB

### Pre-AD Optimization



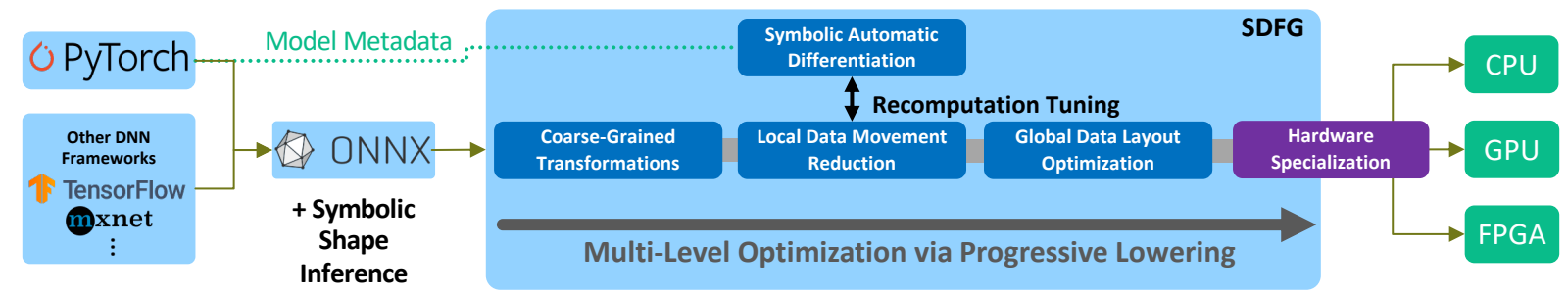
Intermediate Memory: 0MiB

No other framework can perform this type of low-level pre-AD optimization



Prune layouts  
That don't permit  
lowering to BLAS

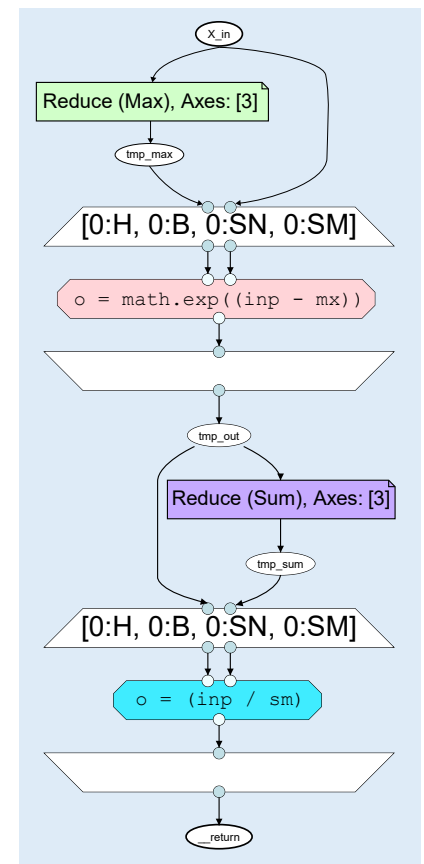
Memlet propagation  
& codegen is data-  
layout aware



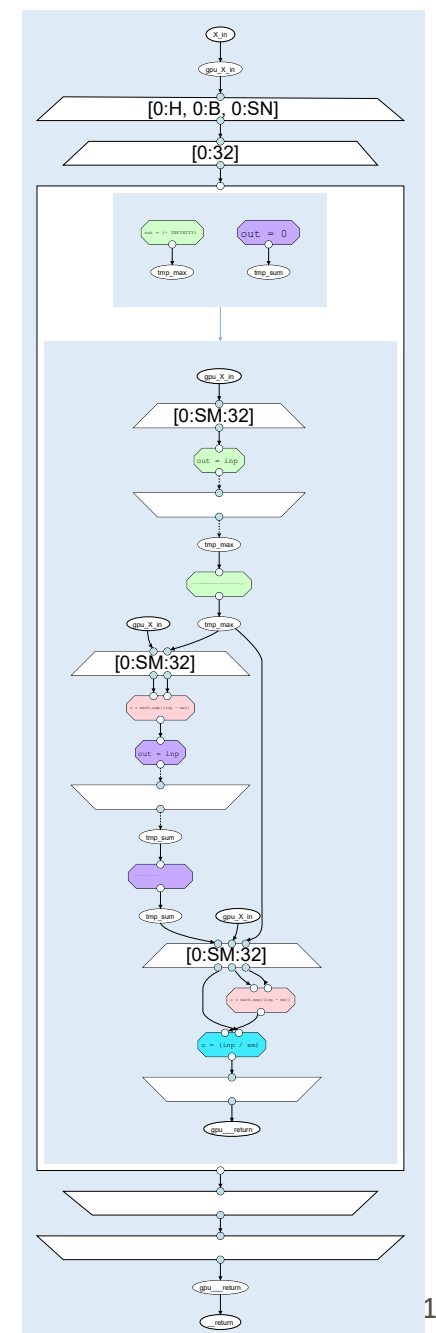
# Operator implementation in NumPy

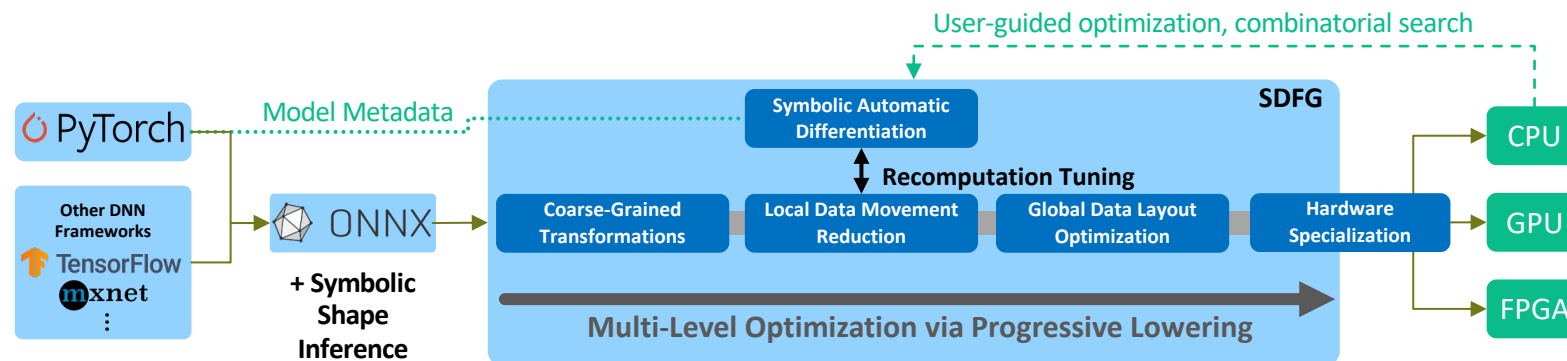
```
def Softmax(input, output):
    max = input.max(axis=axis, keepdims=True)
    exp = np.exp(input - max)
    sum = exp.sum(axis=axis, keepdims=True)
    output[:] = exp / sum
```

Naive Lowering



WarpTiling





**TRANSFORMATIONS**

- Selection
- Viewport
- FPGATransformState
- MapTiling
- StripMining
- MapDimShuffle
- ReduceExpansion
- MapReduceFusion
- GPUTransformMap
- GPUTransformMap
- MapTilingWithOverlap
- MapExpansion
- OuterProductOperation
- GPUTransformLocalStorage
- GPUTransformLocalStorage
- Global
- FPGATransformSDFG
- NestSDFG
- GPUTransformSDFG
- Uncategorized

```

@gemm.py M X
3 import dace
4 import numpy as np
5
6 M, N, K = (dace.symbol(s) for s in 'MNK')
7
8 @dace.program
9 def gemm(A: dace.float64[M, K],
10         B: dace.float64[K, N],
11         C: dace.float64[M, N]):
12     # Transient variable
13     tmp = dace.define_local([M, N, K], dtype=A.dtype)
14
15     for i, j, k in dace.map[0:M, 0:N, 0:K]:
16         tmp[i, j, k] = A[i, k] * B[k, j]
17
18 C[:] = np.sum(tmp, axis=2)
19
20 N = M = K = 128
21
22 A = np.random.rand(M, K)
23 B = np.random.rand(K, N)
24 C = np.random.rand(M, N)
25
26 gemm(A, B, C)
27

```

**program.sdfg M X**

Search the graph

Display Breakpoints Refresh SDFG

SDFG gemm

Go to Generated Code Clear Info x

General

arg\_names [A, B, C]

constants\_prop

exit\_code

global\_code

init\_code

instrument No\_Instrumentation

openmp\_sections

symbols

```

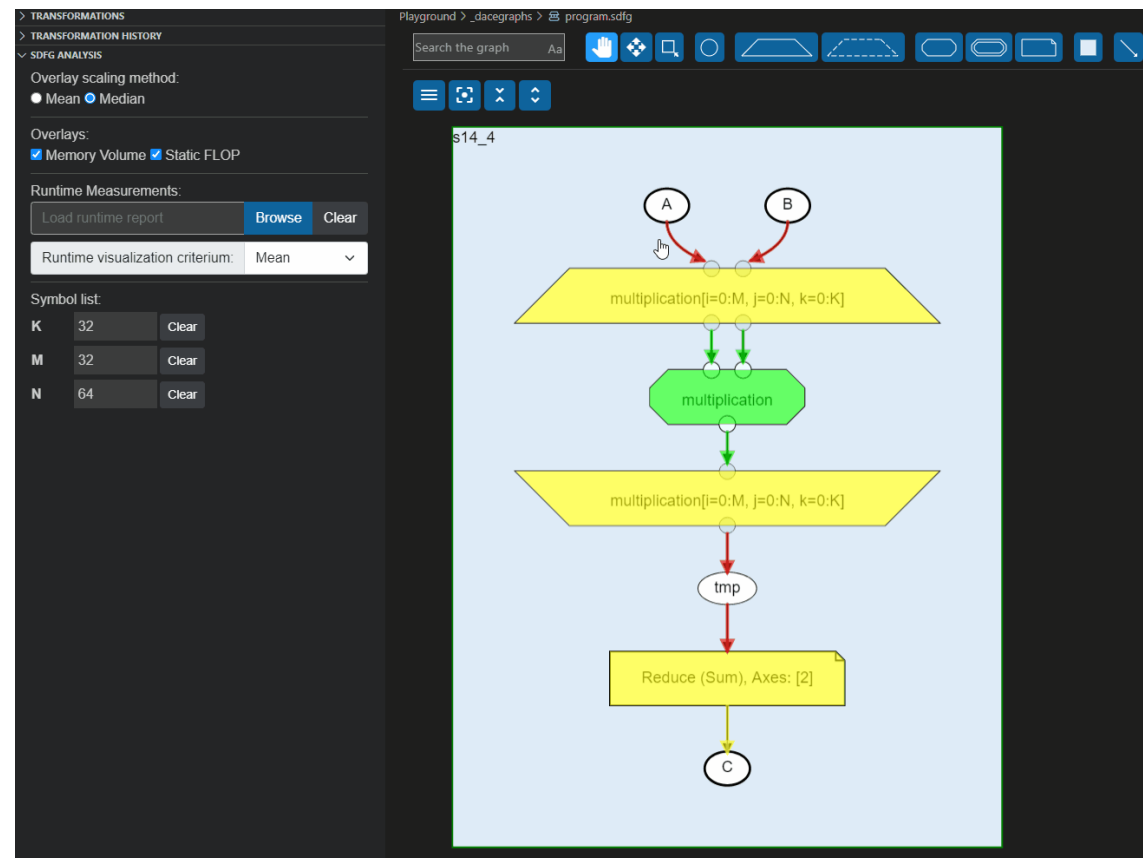
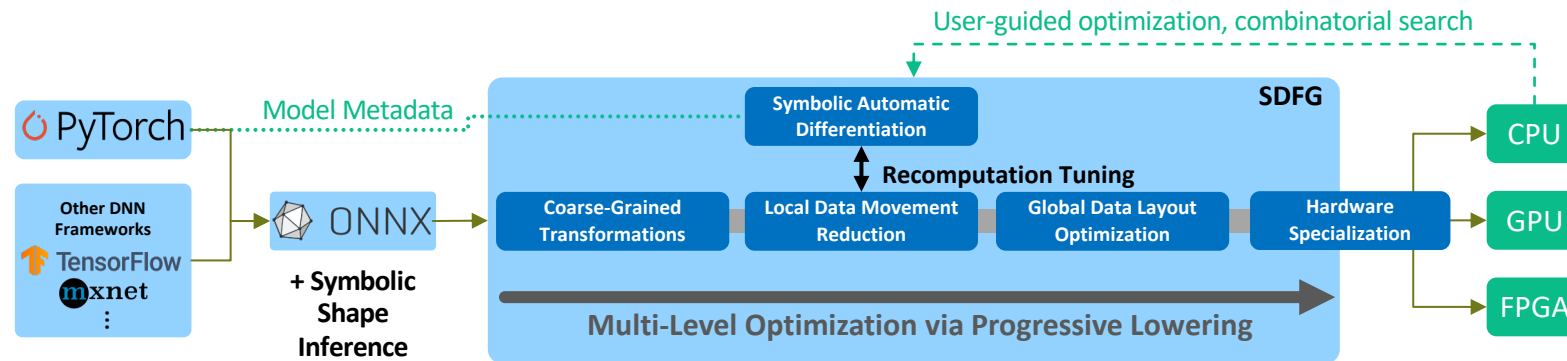
{
  "K": "int32",
  "M": "int32",
  "N": "int32"
}

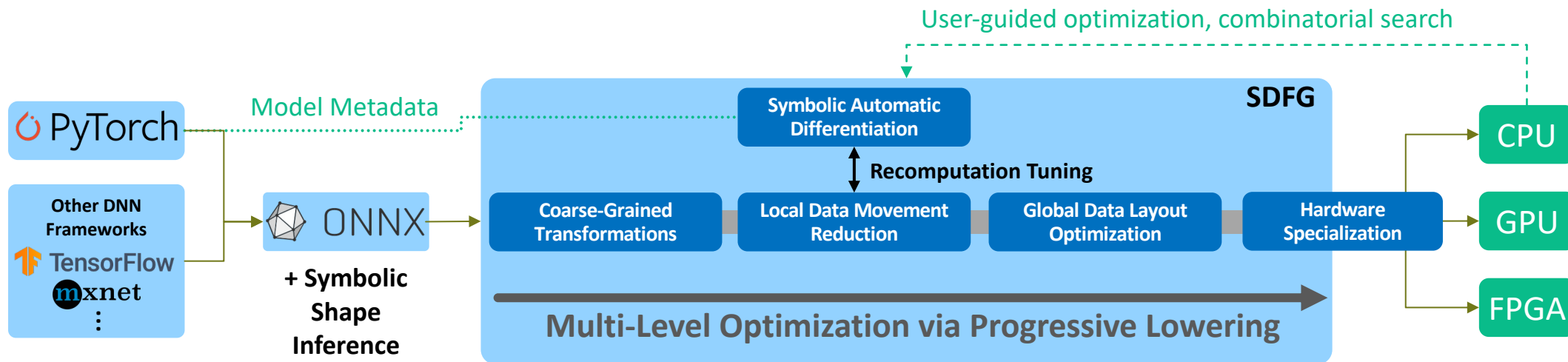
```

Uncategorized

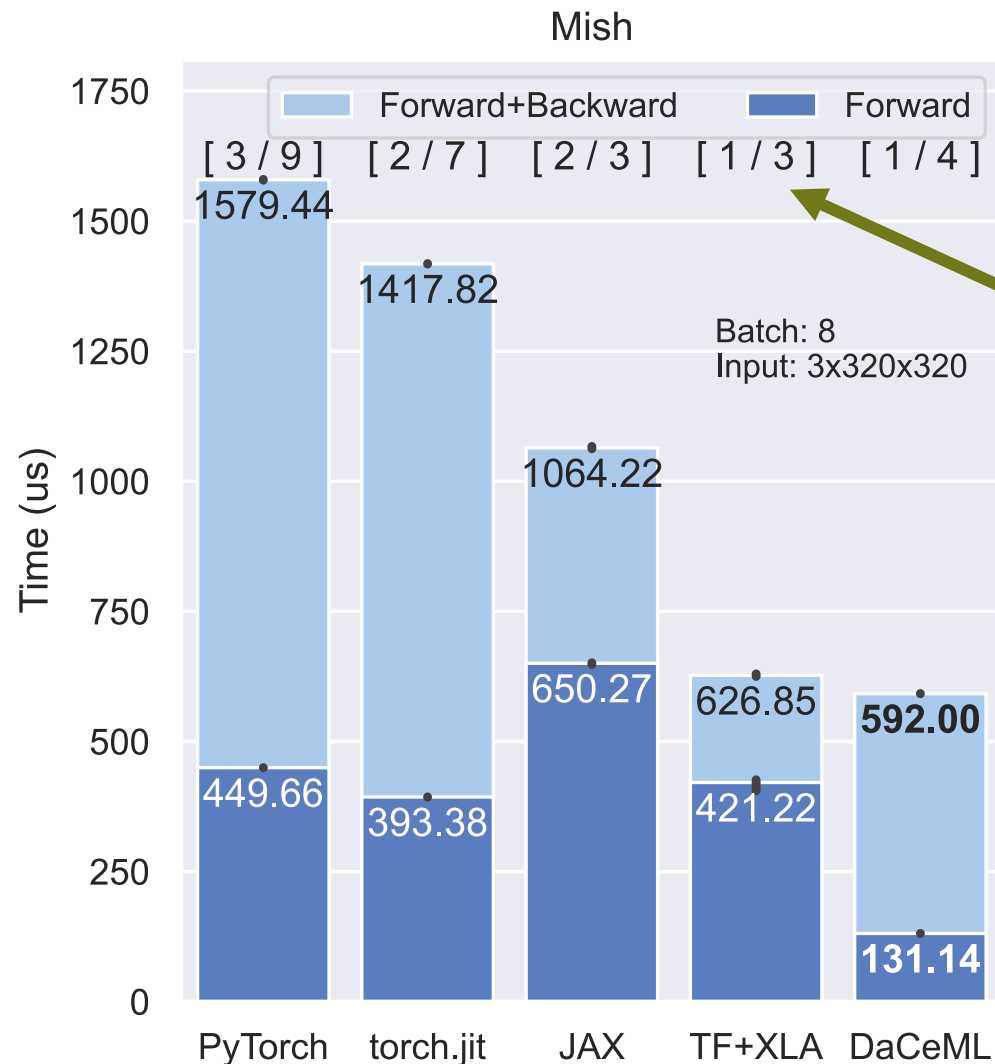
name gemm





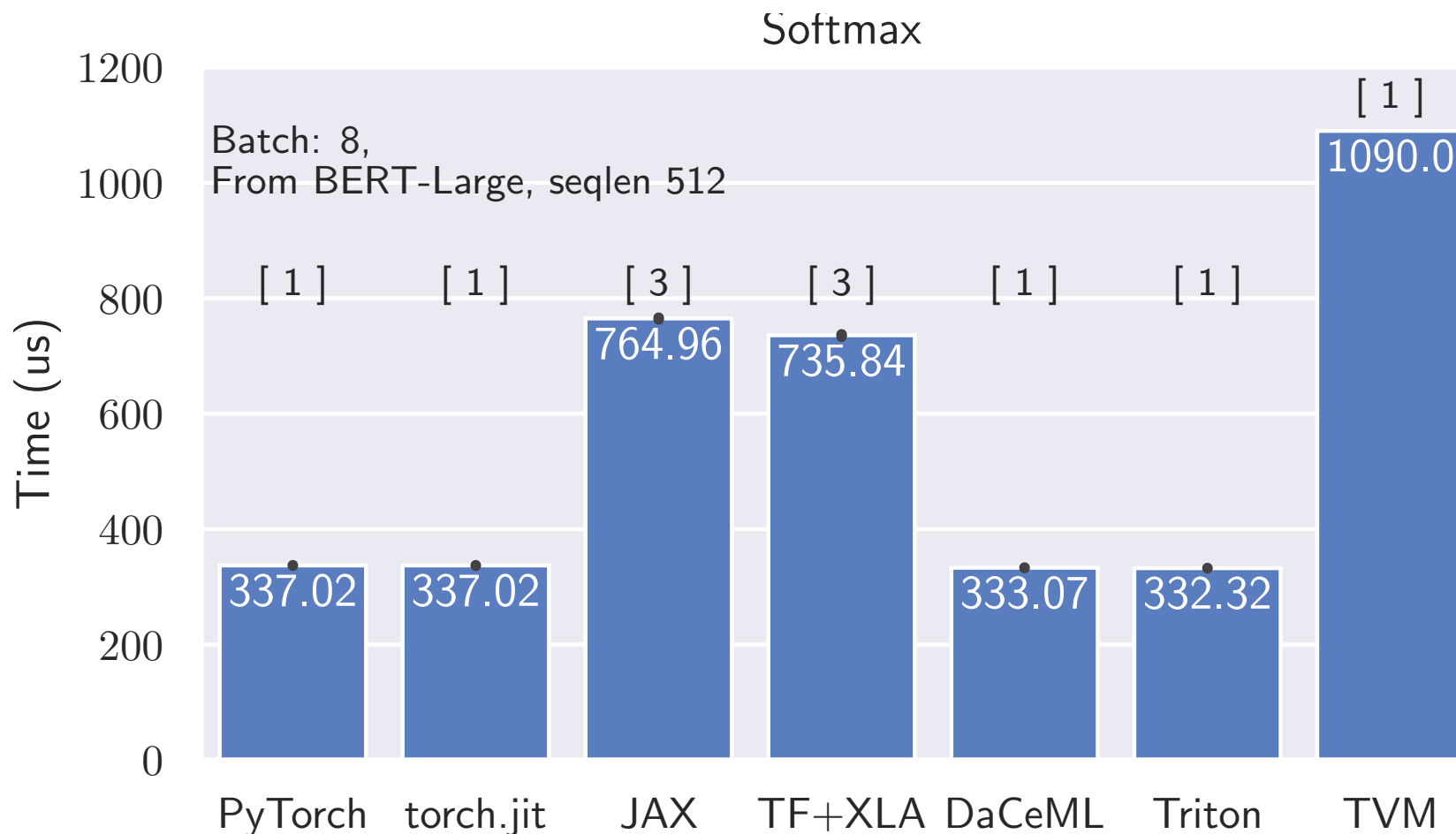


# Results - Pre-AD Fusion (Mish)



- XLA manages to produce a single fused kernel
- But fails to eliminate global loads/stores introduced to stash intermediate values

# Results – Softmax



Transformation recipe **matches hand written kernels** and generalizes to other operators (e.g. Layer Normalization)

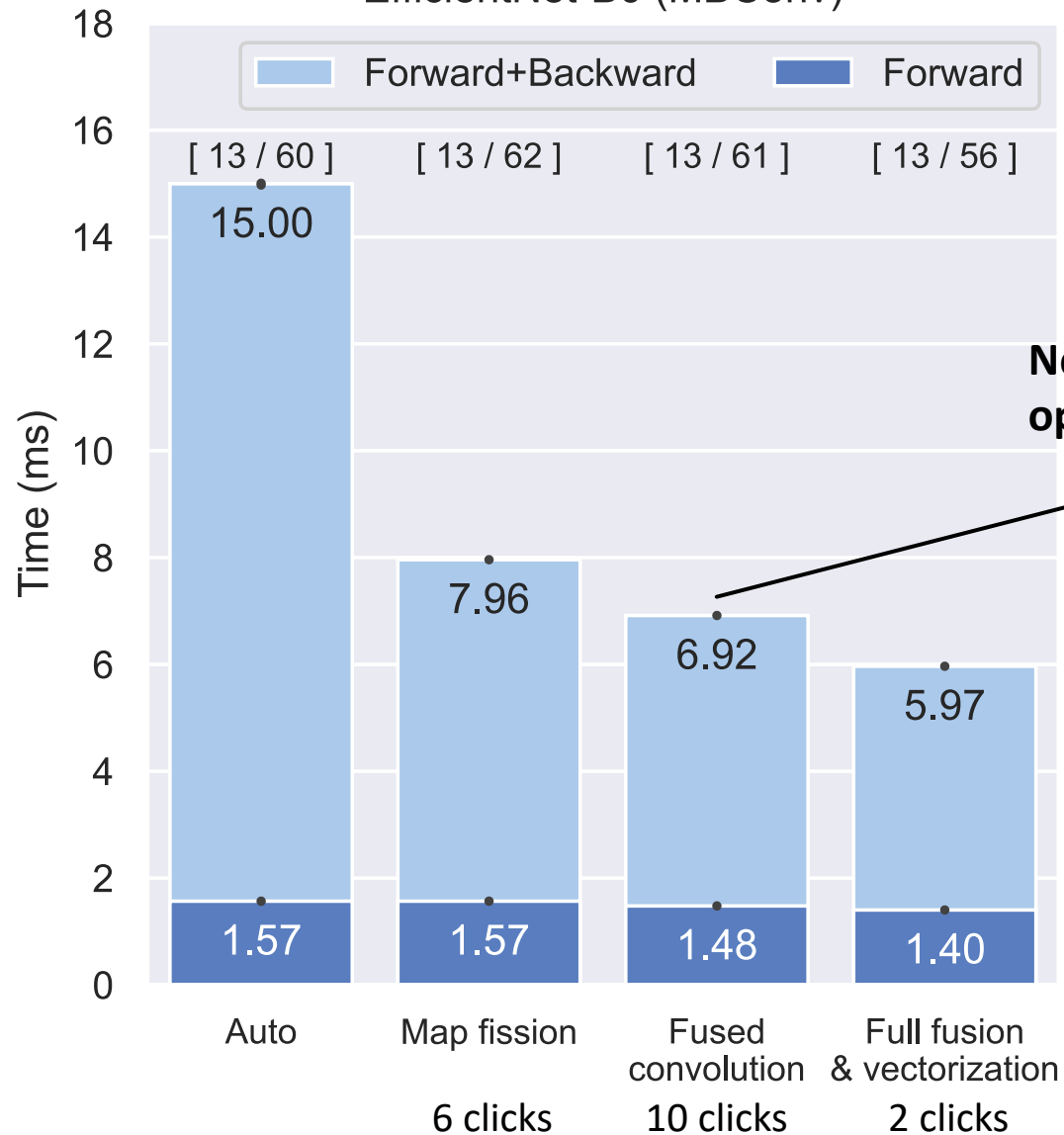


# Results – Automatic Optimization

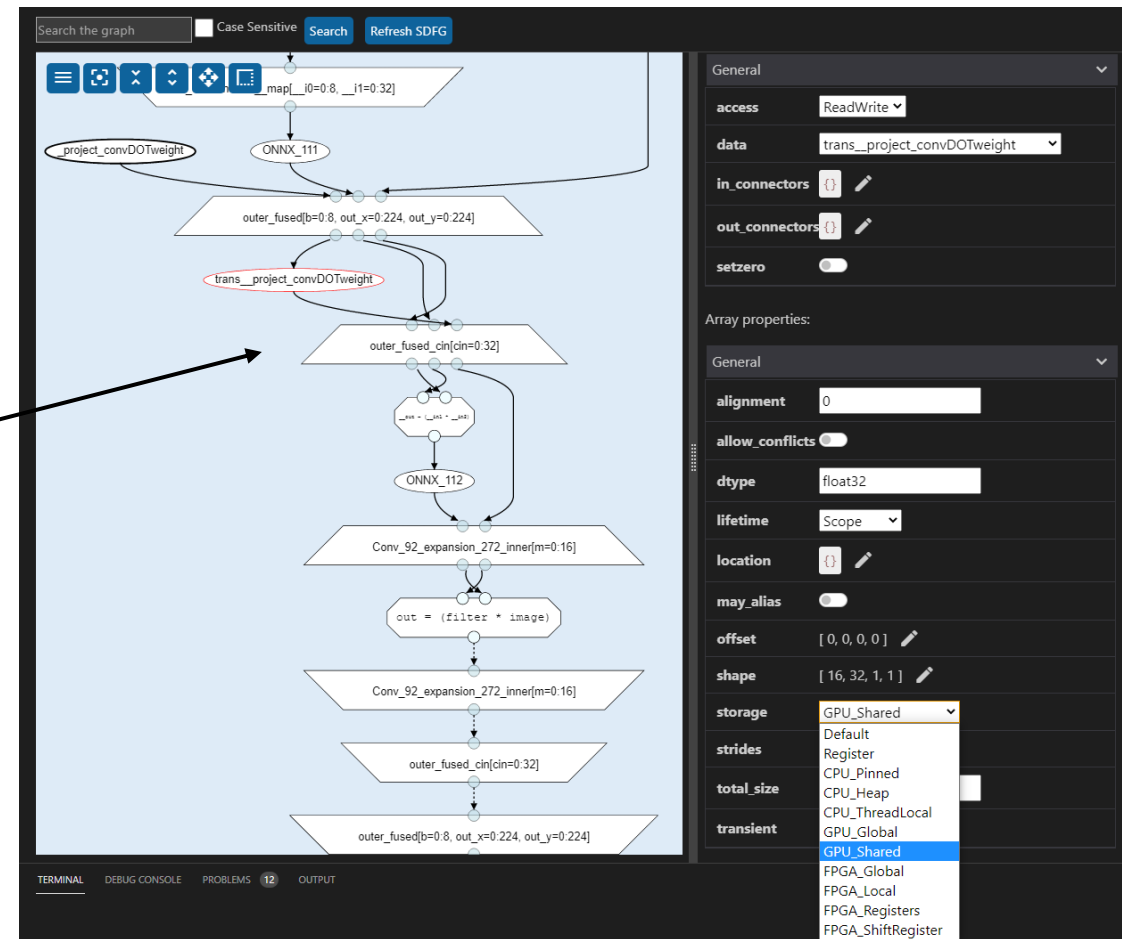
	PyTorch		torch.jit		JAX		TF+XLA		DaCeML	
	→	↔	→	↔	→	↔	→	↔	→	↔
Automatic										
ResNet-50 (👁️)	14.55	32.04	<b>9.98</b>	<b>31.94</b>	14.17	33.93	12.33	35.57	10.03	32.45
Wide ResNet-50-2 (👁️)	22.50	70.94	22.45	70.83	40.49	98.13	32.79	99.06	<b>20.62</b>	<b>67.99</b>
MobileNet V2 (👁️)	9.98	18.45	6.22	15.53	—	—	7.42	20.29	<b>4.74</b>	<b>14.77</b>
EfficientNet (👁️)	2.05	6.90	2.04	6.94	2.39	7.40	<b>1.54</b>	<b>6.37</b>	1.57	15.00
MLP Mixer (👁️)	1.63	<b>3.65</b>	<b>1.36</b>	3.66	1.77	4.01	—	—	1.48	4.25
FCN8s (🧩)	46.85	158.42	46.82	<b>158.40</b>	—	—	—	—	<b>45.97</b>	166.30
WaveNet (🔊)	23.21	46.39	<b>18.67</b>	41.49	—	—	—	—	26.16	<b>41.07</b>
BERT <sub>LARGE</sub> (single) (🏠)	11.05	31.76	11.05	31.82	<b>10.93</b>	<b>29.94</b>	11.14	38.73	11.44	32.98
BERT <sub>LARGE</sub> (mixed) (🏠)	2.94	8.18	<b>2.92</b>	8.20	3.19	<b>8.11</b>	3.80	10.76	3.34	9.25
DLRM (📈)	118.07	126.55	<b>117.38</b>	126.83	—	—	—	—	117.69	<b>126.42</b>

# Results – Guided Optimization

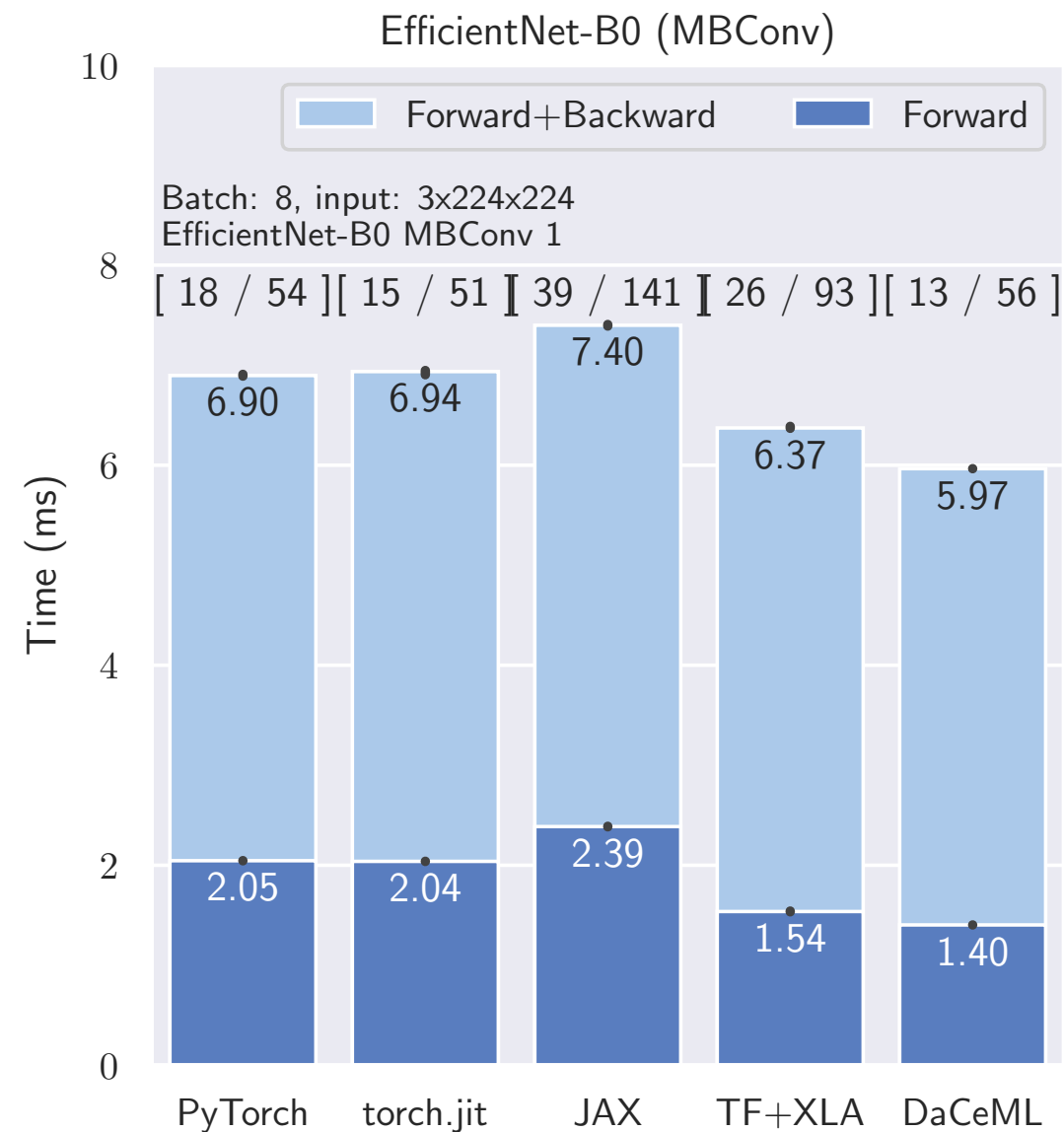
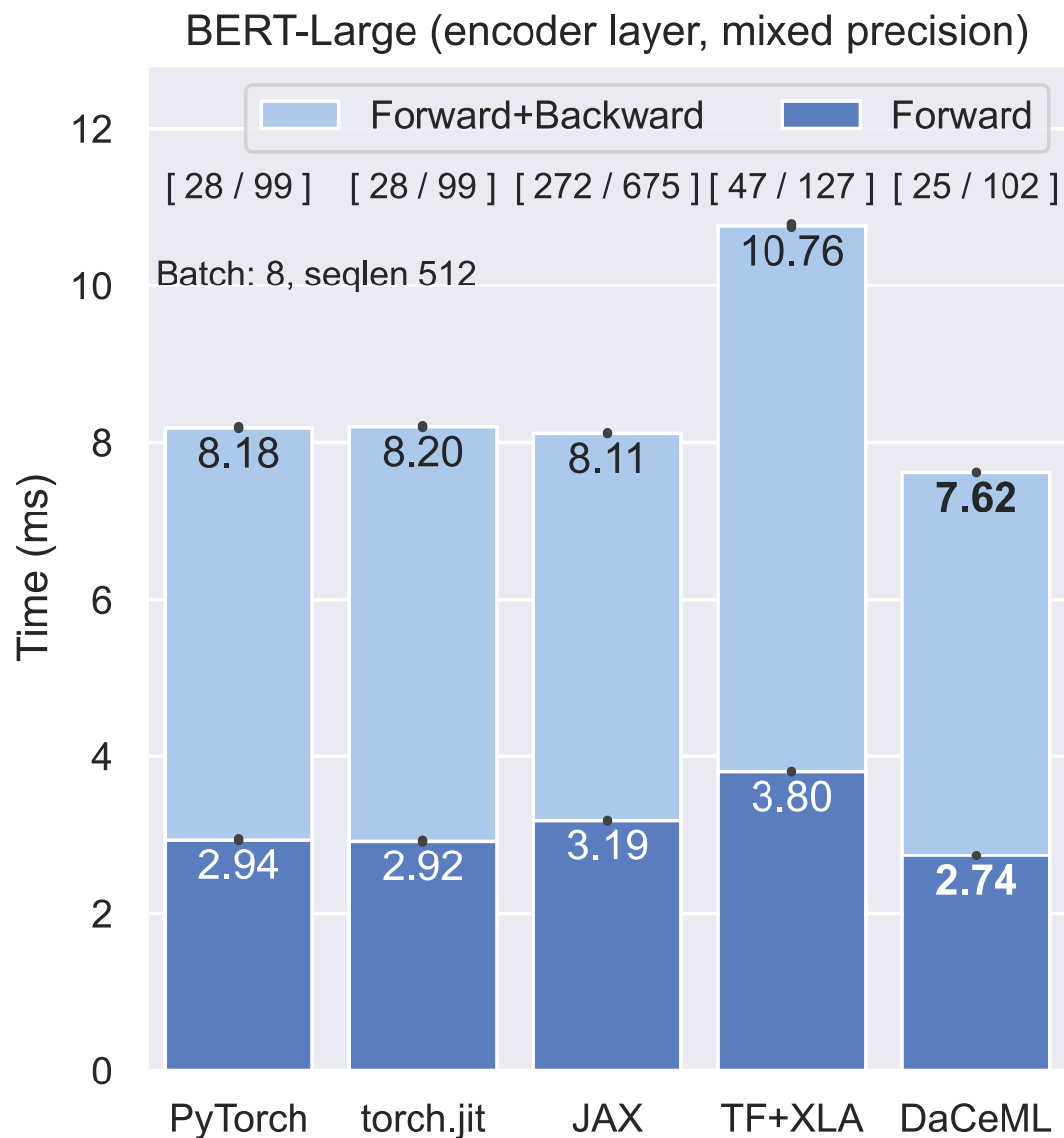
EfficientNet-B0 (MBCConv)



**New fusion opportunity**



# Results – Guided Optimization



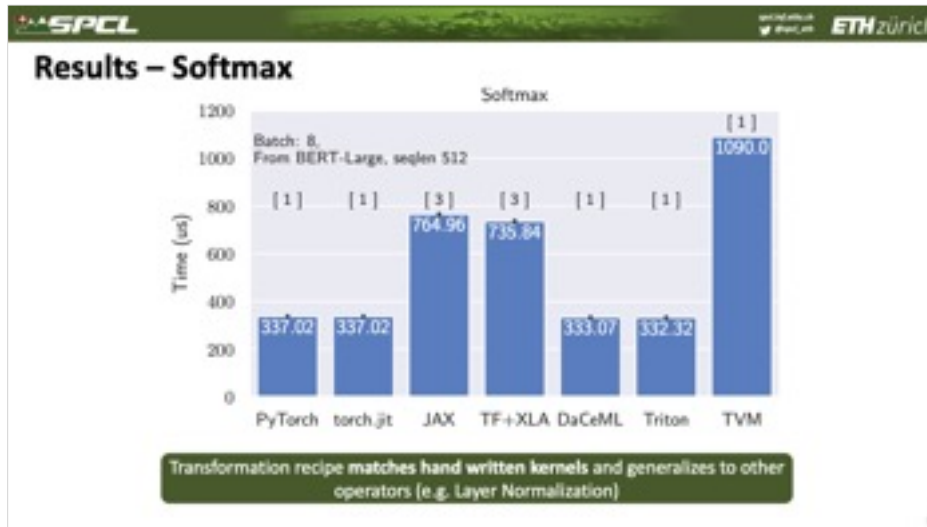
# Summary

```

from torch import nn
from daceml.pytorch import daceml_module

@daceml_module
class MyModule(nn.Module):
    def __init__(self, n_in, n_out):
        super().__init__()
        self.linear = nn.Linear(n_in, n_out)
        self.fanout = n_out

    def forward(self, x):
        return self.linear(x) / self.fanout
    
```



### Results – Models

	PyTorch		torch.jit		JAX		TF+XLA		DaCeML	
	→	⇄	→	⇄	→	⇄	→	⇄	→	⇄
ResNet-50 (🏆)	14.55	32.04	<b>9.98</b>	<b>31.94</b>	14.17	33.93	12.33	35.57	10.03	32.45
Wide ResNet-50-2 (🏆)	22.50	70.94	22.45	70.83	40.49	98.13	32.79	99.06	<b>20.62</b>	<b>67.99</b>
MobileNet V2 (🏆)	9.98	18.45	6.22	15.53	—	—	7.42	20.29	<b>4.74</b>	<b>14.77</b>
EfficientNet (🏆)	2.05	6.90	2.04	6.94	2.39	7.40	<b>1.54</b>	<b>6.37</b>	1.57	15.00
MLP Mixer (🏆)	1.63	<b>3.65</b>	<b>1.36</b>	3.66	1.77	4.01	—	—	1.48	4.25
FCNs (🏆)	46.85	158.42	46.82	<b>158.40</b>	—	—	—	—	<b>45.97</b>	166.30
WaveNet (🏆)	23.21	46.39	<b>18.67</b>	41.49	—	—	—	—	26.16	<b>41.07</b>
BERT <sub>LARGE</sub> (🏆) (🏆)	11.05	31.76	11.05	31.82	<b>10.93</b>	<b>29.94</b>	11.14	38.73	11.44	32.98
BERT <sub>LARGE</sub> (🏆) (🏆)	2.94	8.18	<b>2.92</b>	8.20	3.19	<b>8.11</b>	3.80	10.76	3.34	9.25
DLRM (🏆)	118.07	126.55	<b>117.38</b>	126.83	—	—	—	—	117.69	<b>126.42</b>
Guided										
EfficientNet (🏆)	2.05	6.90	2.04	6.94	2.39	7.40	1.54	6.37	<b>1.40</b>	<b>5.97</b>
BERT <sub>LARGE</sub> (🏆) (🏆)	2.94	8.18	2.92	8.20	3.19	8.11	3.80	10.76	<b>2.74</b>	<b>7.62</b>

