# DiffDA: a Diffusion Model for Weather-scale Data Assimilation
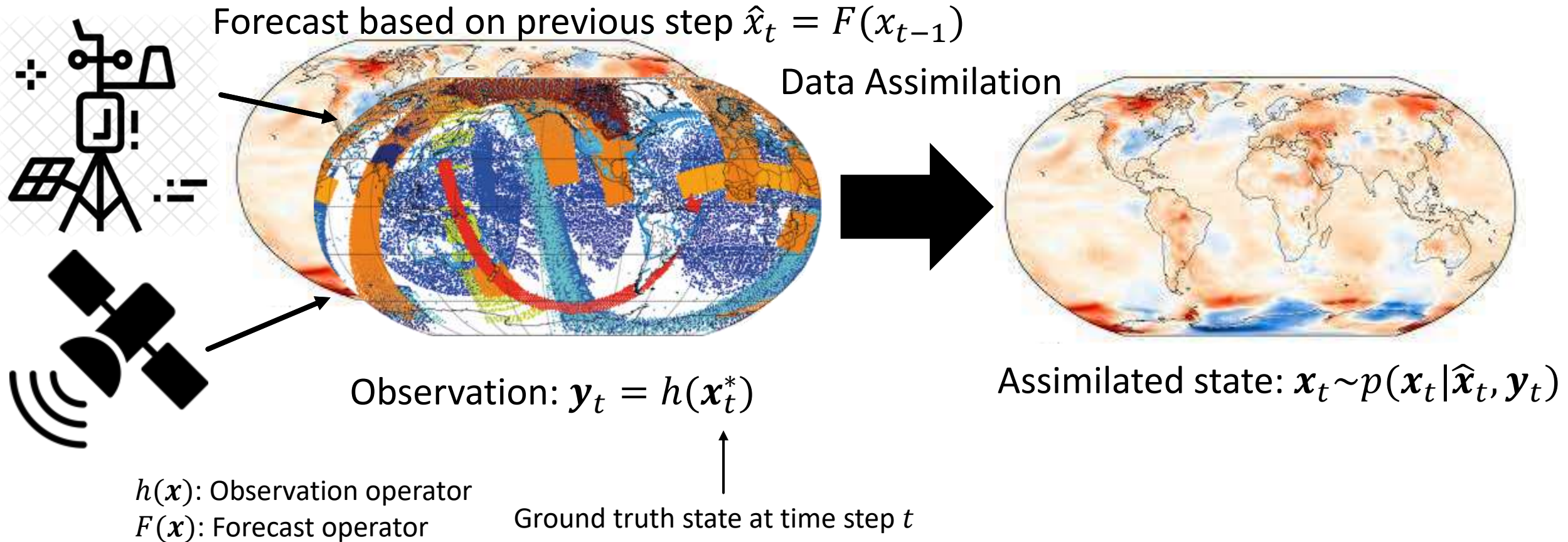
LANGWEN HUANG, Lukas Gianinazzi, Yuejiang Yu, Peter Dominik Dueben, TORSTEN HOEFLER
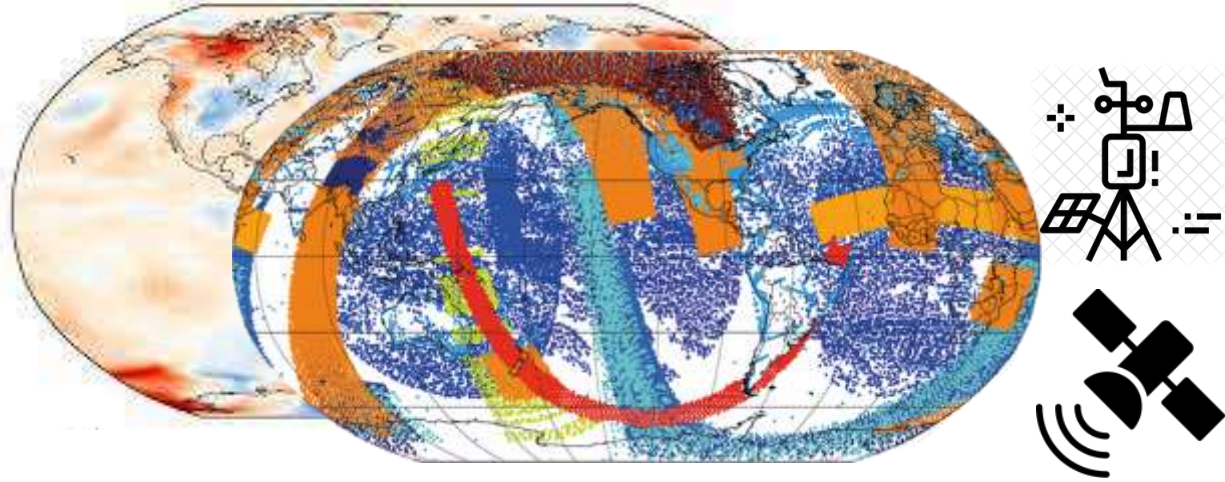
**SPCL**
spcl.ethz.ch
@spcl
@spcl_eth
CSCS
**ETH** zürich

# Motivation

- **Traditional DA methods are slow.**

- **Traditional DA methods make point estimatimation for posterior distribution.**

- **DA tools are not easily available**

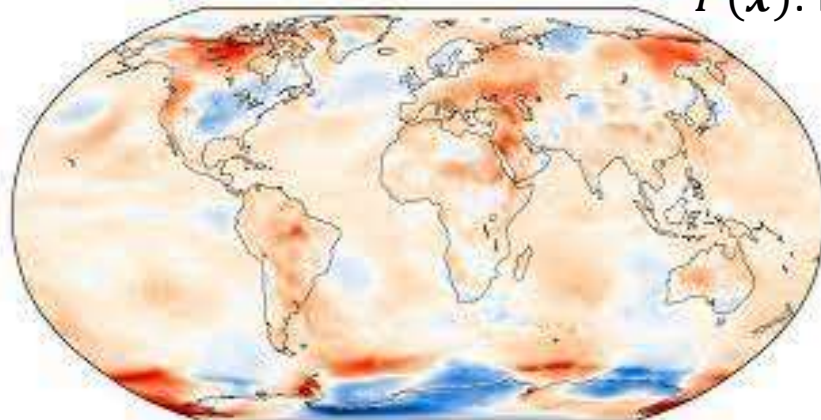- **AI weather models rely on reanalysis datasets.**

Forecast based on previous step $\hat{x}_t = F(x_{t-1})$

Data Assimilation



Assimilated state: $\boldsymbol{x}_t \sim p(\boldsymbol{x}_t | \hat{\boldsymbol{x}}_t, \boldsymbol{y}_t)$

Observation: $\boldsymbol{y}_t = h(\boldsymbol{x}_t^*)$

$h(\boldsymbol{x})$: Observation operator

$F(\boldsymbol{x})$: Forecast operator

Ground truth state at time step $t$

https://www.ecmwf.int/en/about/media-centre/news/2019/forecasting-system-upgrade-set-improve-global-weather-forecasts

2

# Motivation

Predicted state $\widehat{\boldsymbol{x}}_t = F(\boldsymbol{x}_{t-1})$  Observation: $\boldsymbol{y}_t = h(\boldsymbol{x}_t^*)$



**Data Assimilation**

$h(\boldsymbol{x})$: Observation operator
$F(\boldsymbol{x})$: Forecast operator

Assimilated state: $\boldsymbol{x}_t \sim p(\boldsymbol{x}_t | \widehat{\boldsymbol{x}}_t, \boldsymbol{y}_t)$
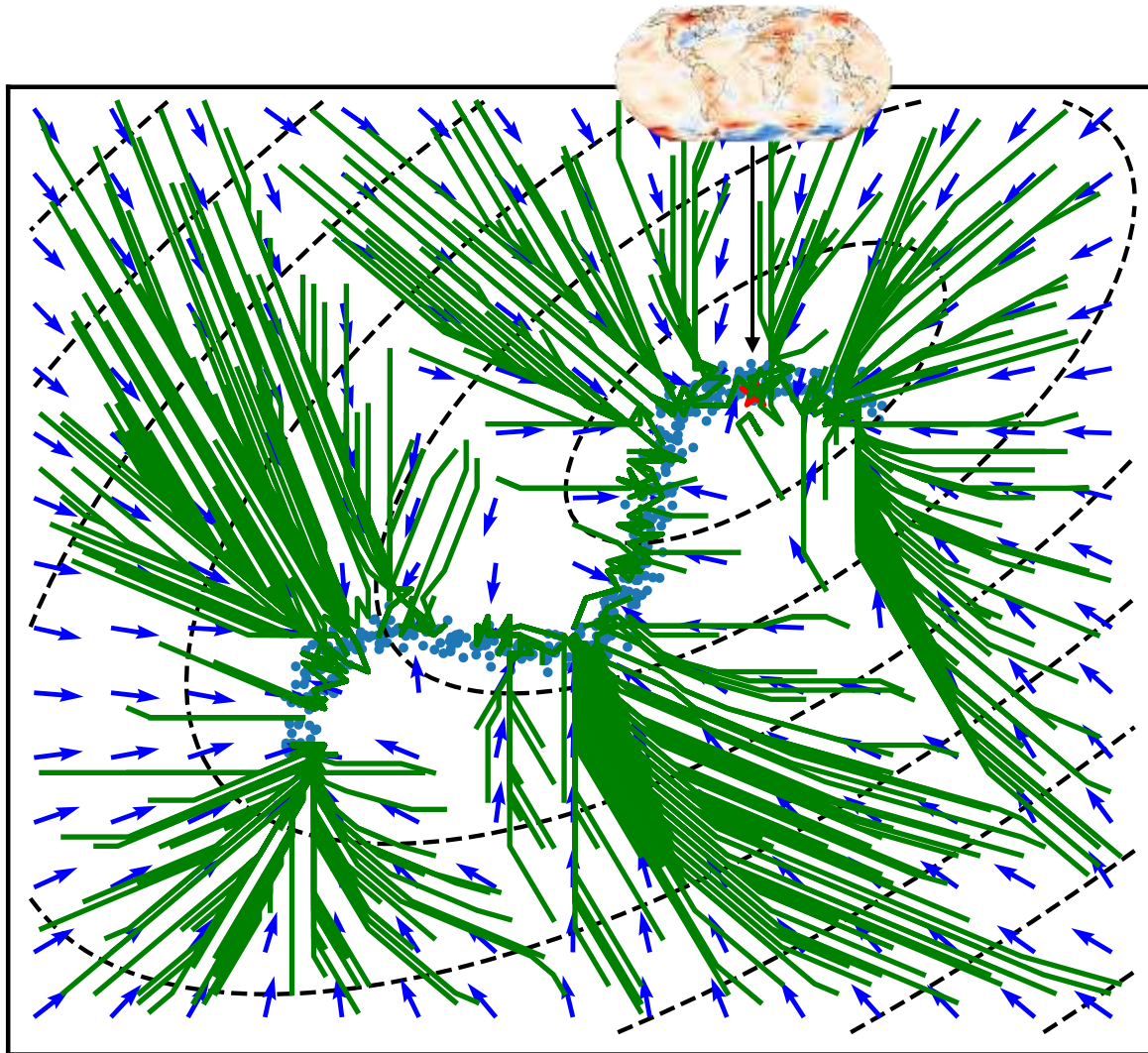
3

# General Idea



- Valid atmosphere states forms a (high-dimensional) manifold
- 3DVar performs maximum likelihood estimation of posterior distribution through minimizing a quadrate loss function
  - Assume Gaussian Process
  - Need to design covariance matrix
  - Need to "invert" covariance matrix when minimizing
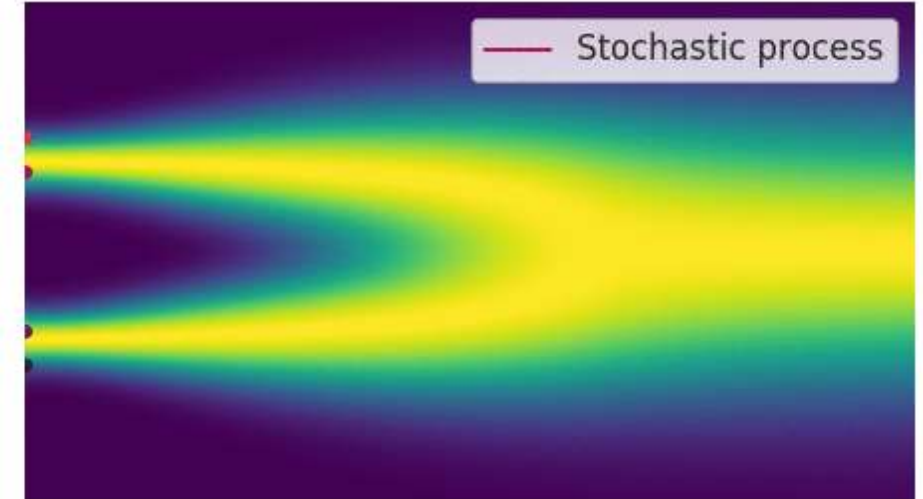  - Perform gradient descent / Newton's method to numerically find minima

# General Idea



- Generalize gradient of loss function into a **learned** vector field
- Denoising diffusion model defines the vector field via the reverse of adding gaussian white noise

# General Idea



- Generalize gradient of loss function into a **learned** vector field
- Denoising diffusion model defines the vector field via the reverse of adding gaussian white noise
- A diffusion model generates (unconditional) samples of possible atmosphere state from randomly generated states

# General Idea

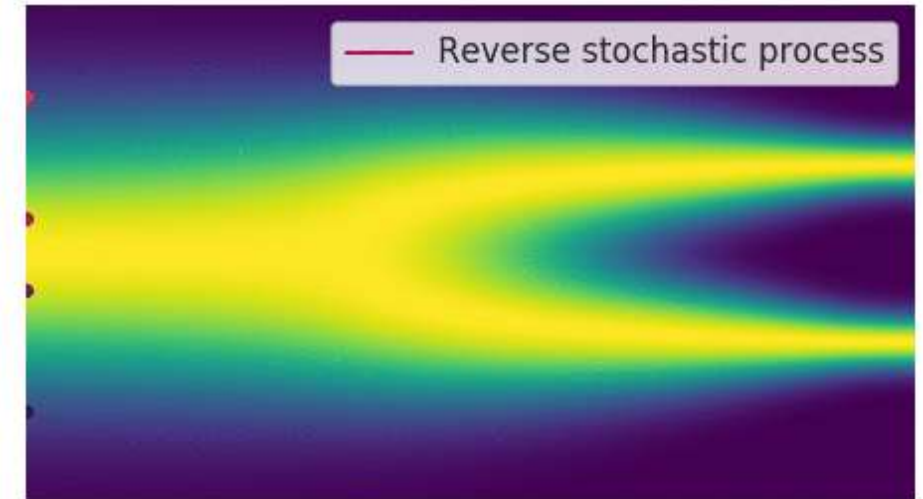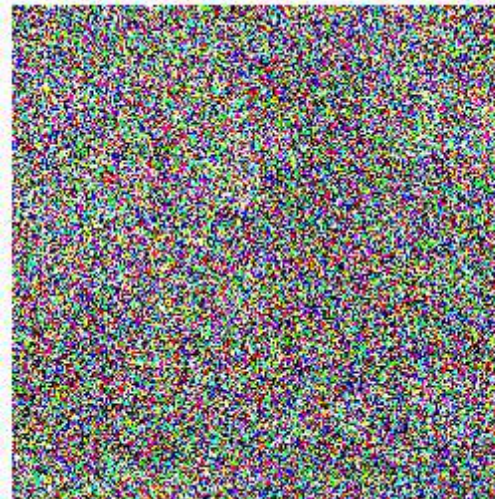- **Sampling from probability distribution with denoising diffusion model**

**Forward Process:**
  complex distribution ->
  simple realizable distribution
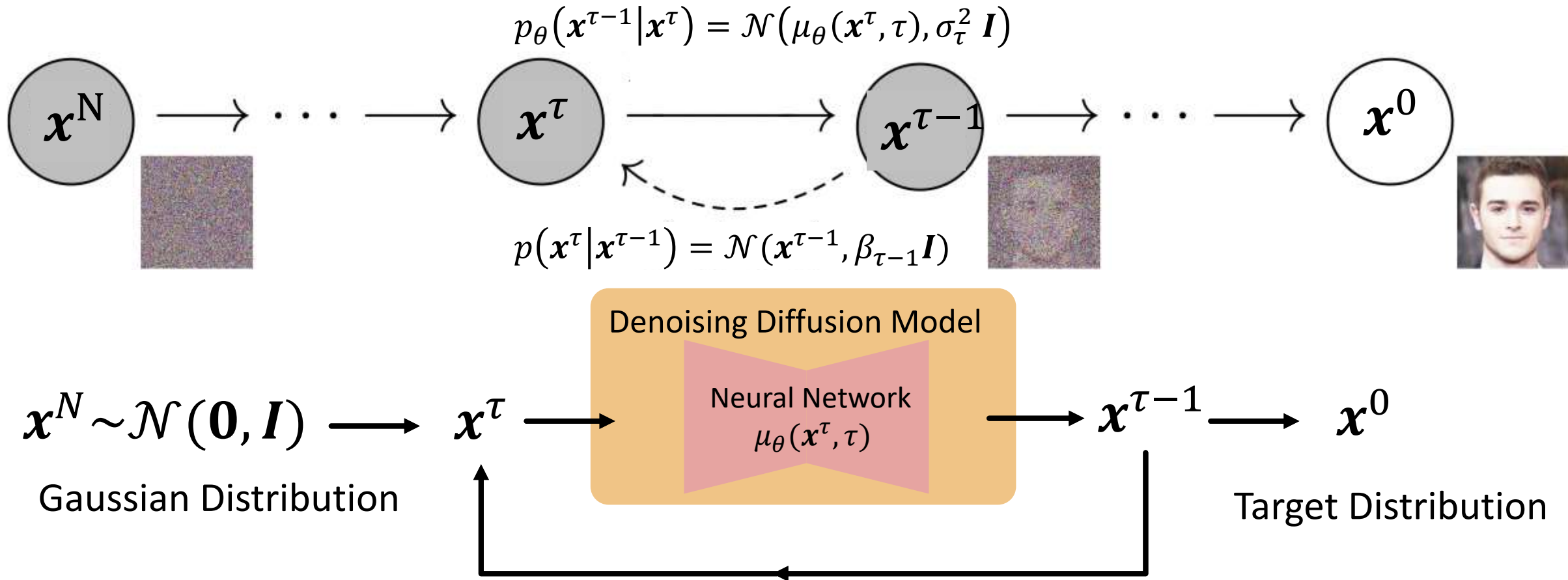
**Backward Process:**
  simple realizable distribution
  -> complex distribution



https://yang-song.net/blog/2021/score/

# General Idea

- **Sampling from probability distribution with denoising diffusion model**

$$p_\theta(x^{\tau-1}|x^\tau) = \mathcal{N}\big(\mu_\theta(x^\tau, \tau), \sigma_\tau^2 I\big)$$

$$x^N \longrightarrow \cdots \longrightarrow x^\tau \longrightarrow x^{\tau-1} \longrightarrow \cdots \longrightarrow x^0$$

$$p(x^\tau|x^{\tau-1}) = \mathcal{N}(x^{\tau-1}, \beta_{\tau-1}I)$$

**Denoising Diffusion Model**

**Neural Network**
$\mu_\theta(x^\tau, \tau)$

$$x^N \sim \mathcal{N}(0, I) \longrightarrow x^\tau \longrightarrow x^{\tau-1} \longrightarrow x^0$$

Gaussian Distribution

Target Distribution

Ho, J., Jain, A. and Abbeel, P., 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems, 33*, pp.6840-6851.
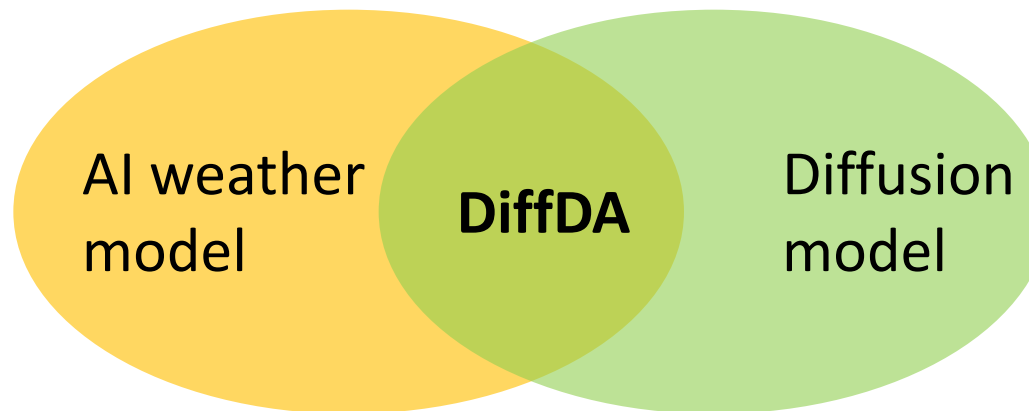
# Challenge 1: How to diffuse on high dimensional fields?

- **Normal input shape for diffusion model: (3, 512, 512)**

- **Shape of atmosphere state: (6x13+5, 721, 1440)**

  - High spatial resolution                  => Need special treatment!

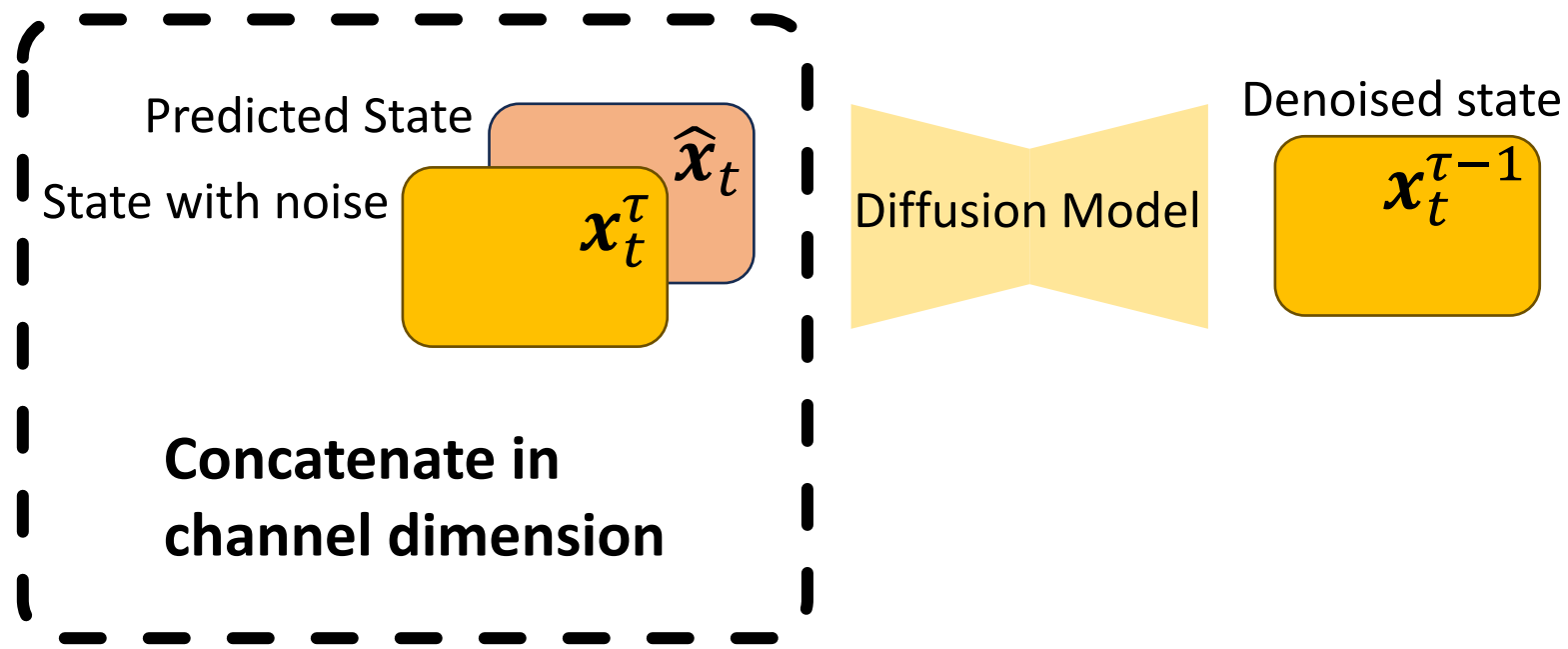  - Dimension size is not power of 2

**Options:**

- Develop a new dedicated network structure

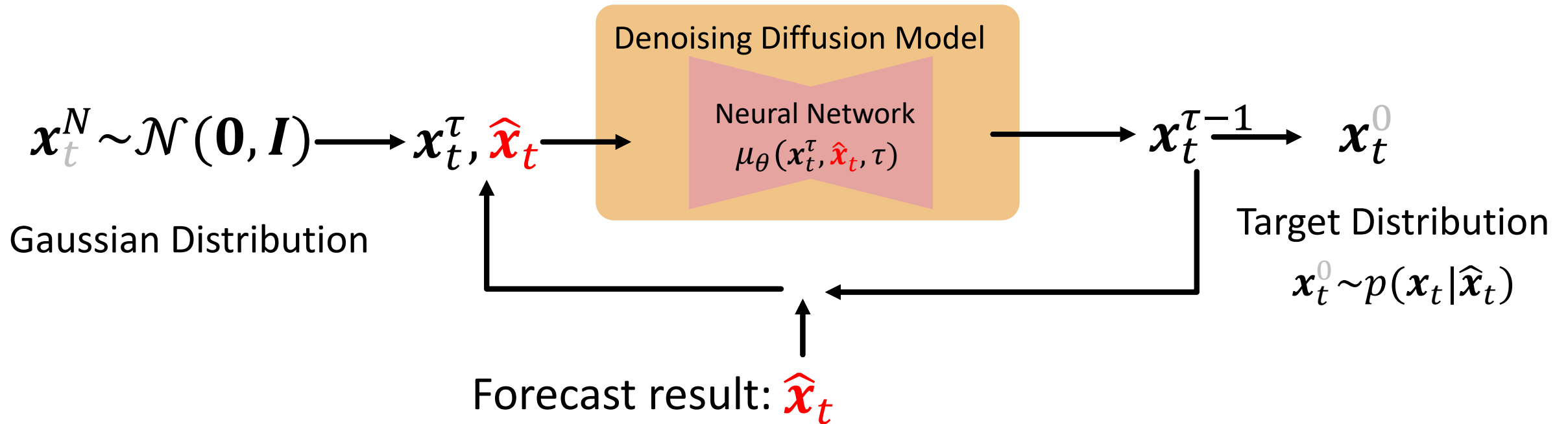- **Use the structure from AI weather model having similar input/output shape!**



AI weather model        **DiffDA**        Diffusion model

# Challenge 2: How to add conditioning?

- **Conditioning for predicted state $\widehat{x}_t$**
  - $\widehat{x}_t$ has the same shape as the assimilated state $x_t^\tau$
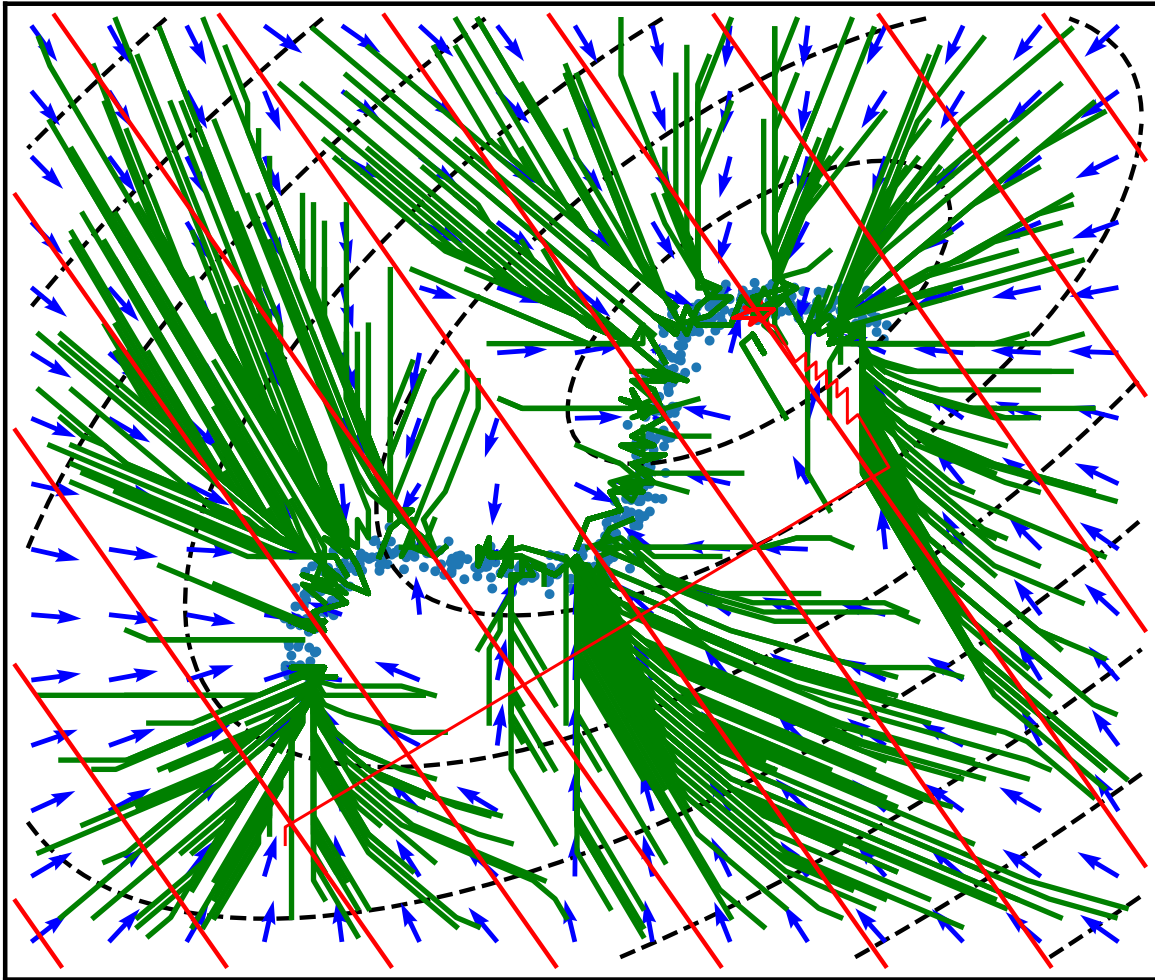  - Replace $\mu_\theta(x_t^\tau, \tau)$ (unconditional) with $\mu_\theta(x_t^\tau, \widehat{x}_t, \tau)$ (conditional)



Predicted State $\widehat{x}_t$

State with noise $x_t^\tau$

**Concatenate in channel dimension**

Diffusion Model

Denoised state $x_t^{\tau-1}$

11

# Overall Process



$$x_t^N \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \longrightarrow x_t^\tau, \hat{x}_t \longrightarrow$$

Denoising Diffusion Model

Neural Network
$\mu_\theta(x_t^\tau, \hat{x}_t, \tau)$

$$\longrightarrow x_t^{\tau-1} \rightarrow x_t^0$$

Gaussian Distribution

Target Distribution
$$x_t^0 \sim p(x_t | \hat{x}_t)$$

Forecast result: $\hat{x}_t$

# Challenge 2: How to add conditioning?

- **Conditioning for observations $y_t$: An inpainting approach**



- Add additional pentalty to guide generation
  - Simpler than 3Dvar loss function
  - $\|\boldsymbol{y} - h(\boldsymbol{x})\|^2$
- Operator splitting
- One step solution to penalty term

# Challenge 2: How to add conditioning?

- **Conditioning for observations $y_t$: An inpainting approach**
  - Assuming $y_t$ is sparse measurement of $x_t^*$: $y_t = H x_t^*$, where 0,1 matrix $H$ has only one nonzero value in each row
  - proof: similar to "classifier" guidance $\nabla_{x_t^\tau} \log p(x_t^\tau | \hat{x}_t, y_t) = \nabla_{x_t^\tau} \log p(x_t^\tau | \hat{x}_t) + \nabla_{x_t^\tau} \log p(y_t | x_t^\tau)$

$$p(y_t | x_t^\tau) \approx \mathcal{N}\left(y_t | H\mathbb{E}[x_t^0 | \hat{x}_t], \Sigma_y\right) \Rightarrow \nabla_{x_t^\tau} \log p(y_t | x_t^\tau) \approx \nabla_{x_t^\tau} \left\| y_t - H\mathbb{E}[x_t^0 | \hat{x}_t] \right\|_{\Sigma_y}^2$$

**Inpainting Pipeline**



**Treatment of Sparse Mask**



Problem: sparse signal often suppressed in the downsampling layer

Solution: enlarge the mask with interpolated data (assuming data is smooth)

Lugmayr, A., Danelljan, M., Romero, A., Yu, F., Timofte, R. and Van Gool, L., 2022. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 11461-11471).

14

# Overall Process

$$x_t^N \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \longrightarrow x_t^\tau, \widehat{x}_t \longrightarrow$$

Denoising Diffusion Model

Neural Network
$\mu_\theta(x_t^\tau, \widehat{x}_t, \tau)$

$\odot (\mathbf{1} - \boldsymbol{m})$

$$x_t^{\tau-1} \rightarrow x_t^0$$

+

Gaussian Distribution

Target Distribution
$$x_t^0 \sim p(x_t | \widehat{x}_t, \boldsymbol{y}_t)$$

$\odot \boldsymbol{m}$

$\widehat{x}_t$

$$\widetilde{x}_t + \boldsymbol{\epsilon}^\tau$$

Interpolated State     $\widetilde{x}_t = \text{Interpolate}(\boldsymbol{m}' \odot x_t^*)$

Soft Mask     $\boldsymbol{m} = \text{Softbleed}(\boldsymbol{m}', \sigma_G)$

Hard Mask     $m_i' = \sum_j H_{ji}, \; \boldsymbol{y} = H x_t^*$

Noise     $\boldsymbol{\epsilon}^\tau = \left(1 - \prod_{s=1}^{j-1} \beta_s\right) \boldsymbol{\epsilon}', \; \boldsymbol{\epsilon}' \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

Kernel

$\sigma_G$

Hard mask

(max, ×)-convolution

scale by 1/cos(lat)

Soft mask

16

# Overall Process

# Experiment Settings

- Backbone model: GraphCast operational (0.25deg 721x1440, 13 levels)

- Training data: ERA5 1979 - 2016 6hour resolution

- Emulate observations from ERA5: randomly sample horizontal coordinates + take all vertical levels & variables

- Batch size: 48 (global), 1 (local)

- Num epochs: 20

- Optimizer: Adam, LR scheduler: warmup_cos_decay

- LR: 1e-5 (0%) -> 1e-4 (12%) -> 3e-6 (100%)

- $\sigma_G = 1.5$

- Compute resources: 48 A100 80G, 4 GPUs per node, 2 days

# Experiments Overview

# Experiment 1 : Single step data assimilation

# Experiment 1: Result

**(non-weighted)-RMSE**

Assimilated Data based on 48h Forecast | GraphCast Forecast on ERA5

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Geopotential at 500hPa | 86.574 | 76.831 | 63.671 | 47.150 | 42.142 | 27.809 | 18.628 | 24.006 | 32.917 | 37.897 | 48.467 | 93.158 |
| Temperature at 850hPa | 1.113 | 1.070 | 0.982 | 0.861 | 0.811 | 0.641 | 0.496 | 0.404 | 0.618 | 0.671 | 0.763 | 0.998 |

- DA results converge to ERA5 data with increasing number of observations
- With 3.8% gird points "observed", the DA result is comparable to 12h forecast error (reduce lead time by 36 hours)

Number of observed columns

Lead time

3.8% total columns

Relative RMSE (w.r.t. 48h GraphCast forecast RMSE, lower is better)

0.2   0.4   0.6   0.8   1.0   1.2

# Case study: Assimilated 2m temperature at 2022-01-03 06z

48h GraphCast Forecast

Assimilated Data

DiffDA

Raw Observations   Interpolated Observations
with softmask

# Case study: Assimilated 2m temperature at 2022-01-03 06z

48h GraphCast Forecast Error



Interpolation Error with 8000 observation columns



DiffDA

Error of Assimilated Data



Assimilated data is better than both inputs
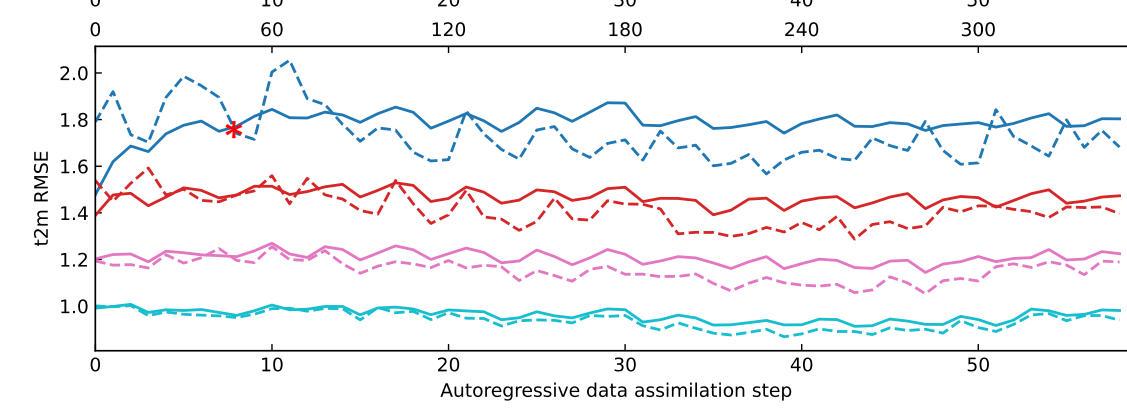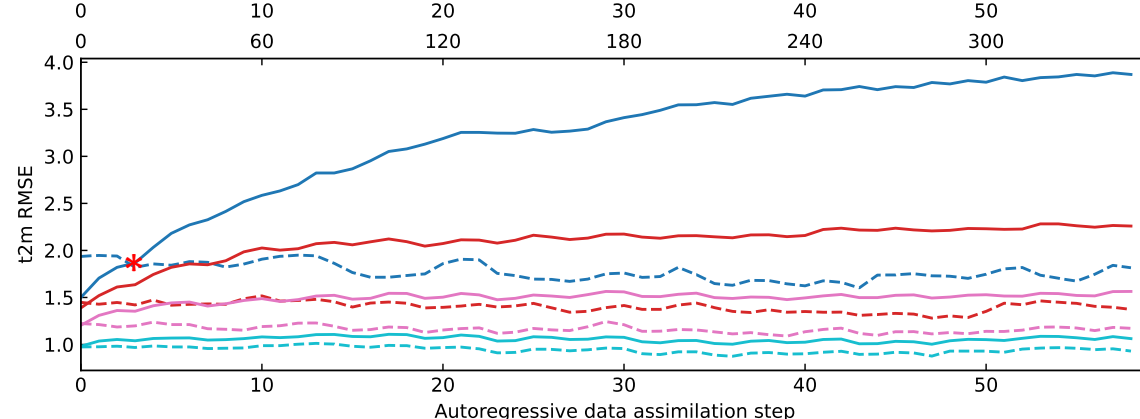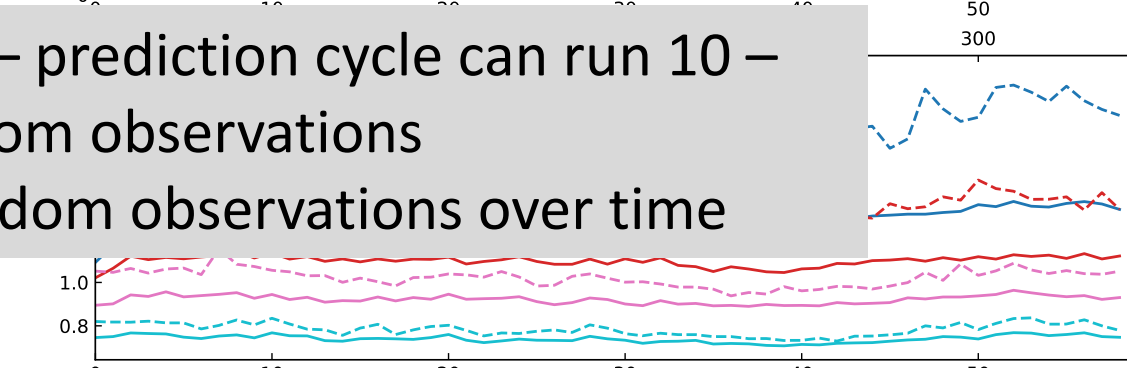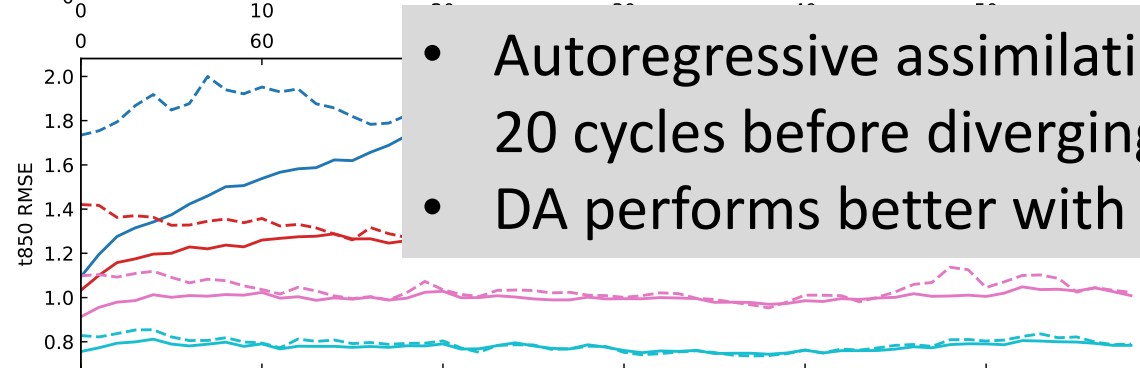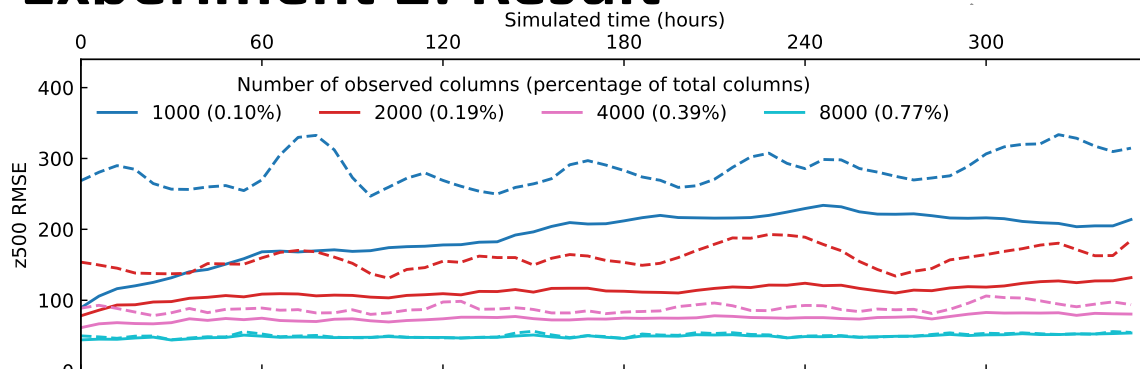(48h forecast and interpolated observations)

# Experiment 2: Autoregressive data assimilation

# Experiment 2: Result



RMSE of Autoregressive assimilated data
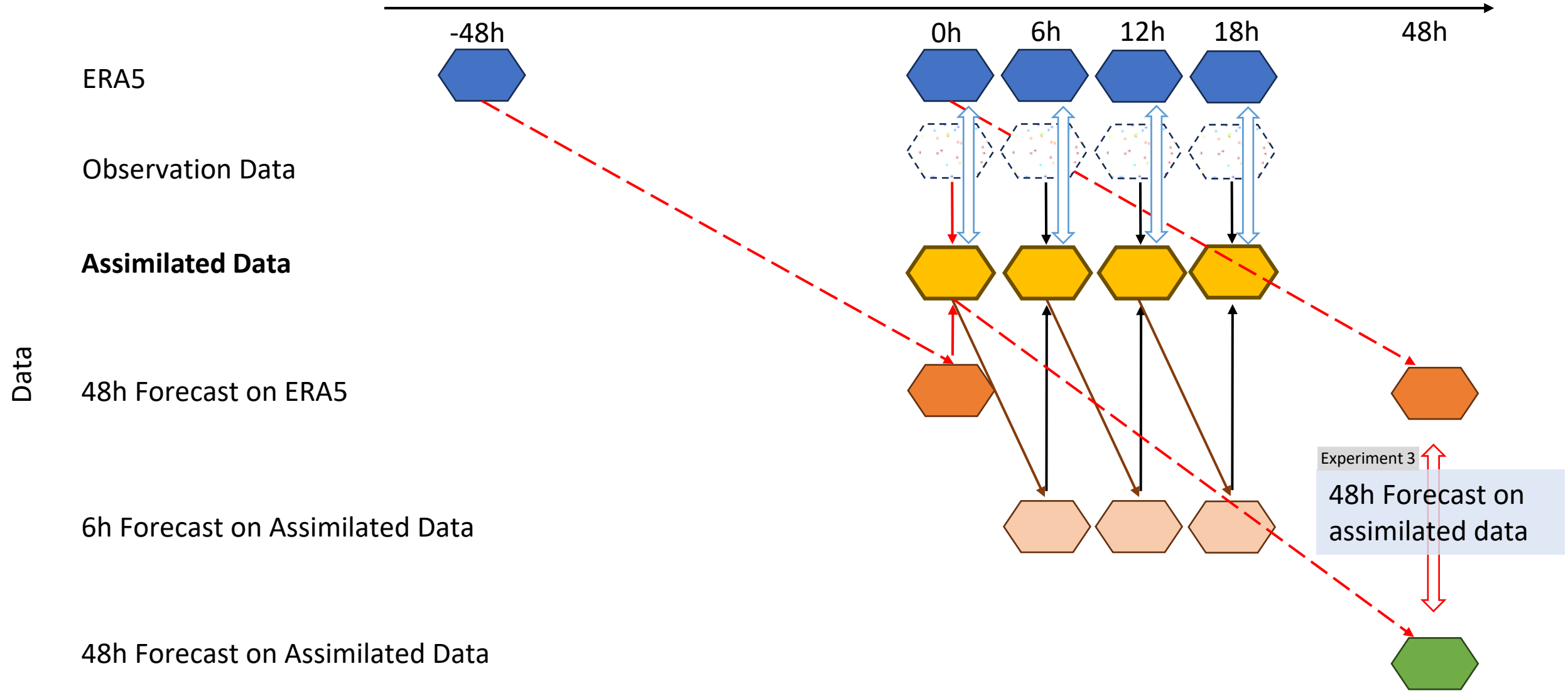
RMSE of interpolated observations

- Autoregressive assimilation – prediction cycle can run 10 – 20 cycles before diverging from observations
- DA performs better with random observations over time
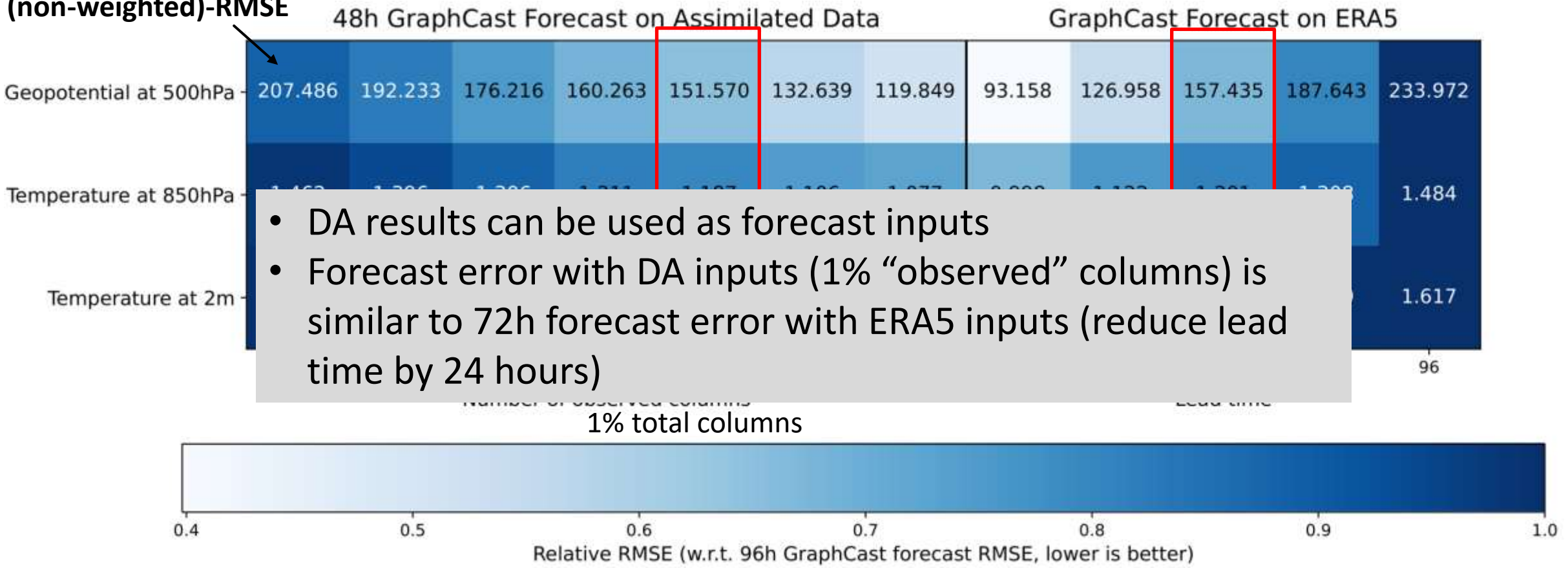
Fixed observations over time

Random observations over time

26

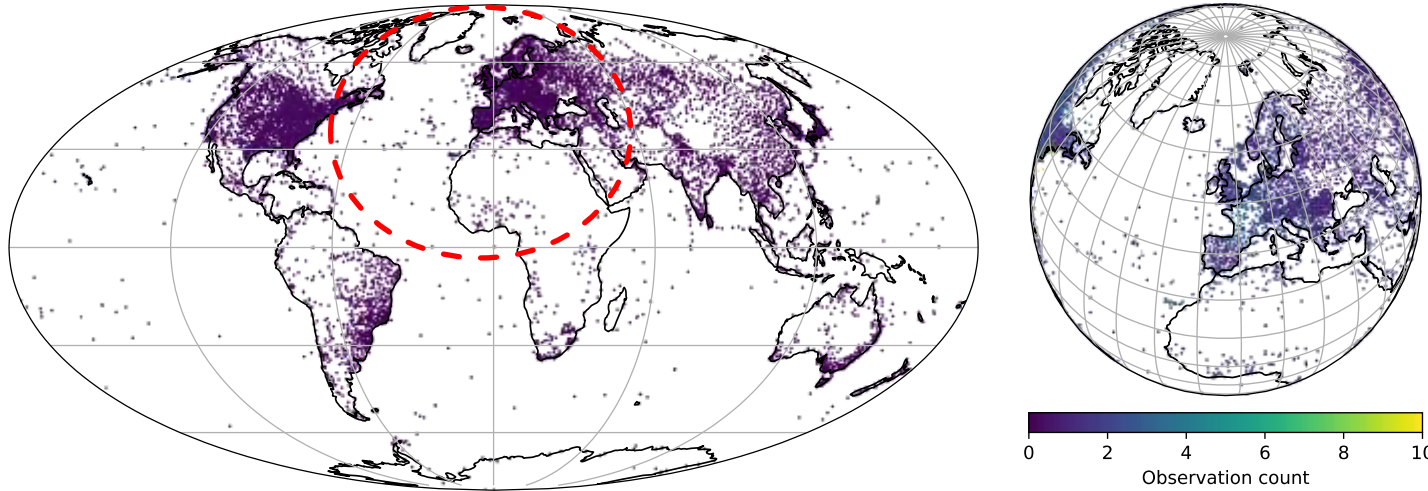# Experiment 3 : 48h forecast on single step assimilated data

# Experiment 3: Result

**(non-weighted)-RMSE**

48h GraphCast Forecast on Assimilated Data | GraphCast Forecast on ERA5

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Geopotential at 500hPa | 207.486 | 192.233 | 176.216 | 160.263 | 151.570 | 132.639 | 119.849 | 93.158 | 126.958 | 157.435 | 187.643 | 233.972 |
| Temperature at 850hPa | 1.463 | 1.306 | 1.306 | 1.211 | 1.187 | 1.106 | 1.077 | 0.908 | 1.122 | 1.281 | 1.398 | 1.484 |
| Temperature at 2m | | | | | | | | | | | | 1.617 |

1% total columns

- DA results can be used as forecast inputs
- Forecast error with DA inputs (1% "observed" columns) is similar to 72h forecast error with ERA5 inputs (reduce lead time by 24 hours)

Relative RMSE (w.r.t. 96h GraphCast forecast RMSE, lower is better)

0.4    0.5    0.6    0.7    0.8    0.9    1.0

# Next Step: Towards Assimilating Real-World Observations



**Challenges:**
- Non-uniform distribution
- Only a subset of variables are measured
- Less observations at higher levels
- Observations are collected in a time window, e.g. : (-3h, 3h)
- Need quality control

Observations of 2m temperature at 30.12.2022 00z from GDAS
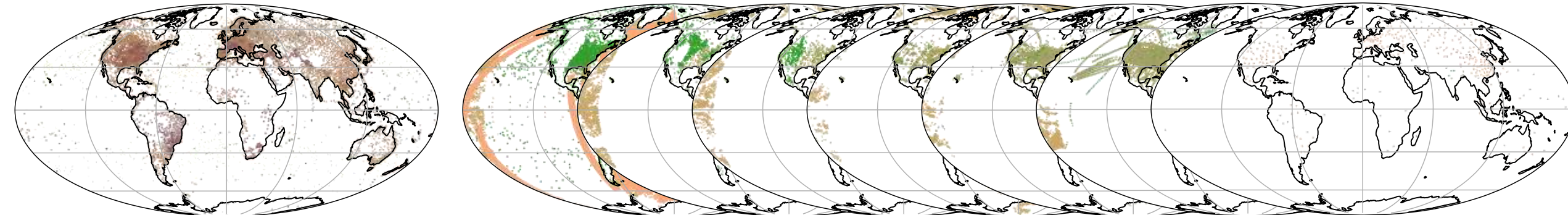Measurements/Total Grid Points: 10054/1036800
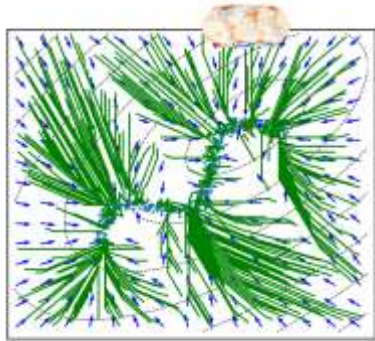Fraction: **0.97 %**

Fraction: **16%**
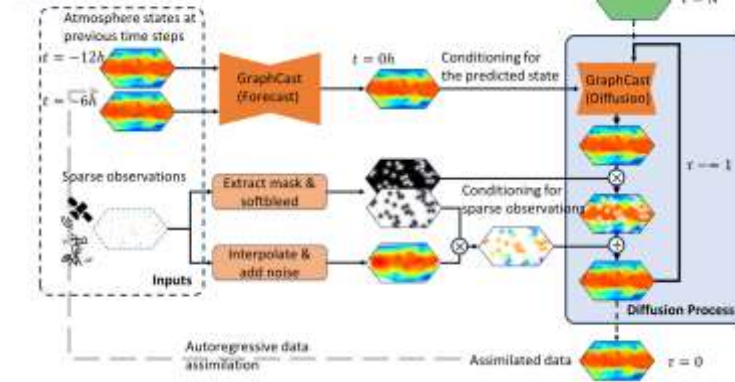
Fraction: **0.05%**



Surface

1000 hPa   925 hPa   850 hPa   500 hPa   300 hPa   200 hPa   100 hPa

TMP    RH    PRMSL    PS    TP06

SPFH    TMP    HGT    UGRD    VGRD    RH

# Conclusions

## More of SPCL's research:

**... or spcl.ethz.ch**









**Next:**
- assimilate real-world observations
- 4D Assimilation
- Incorporate errors in observations
- Incorporate non sparse observations