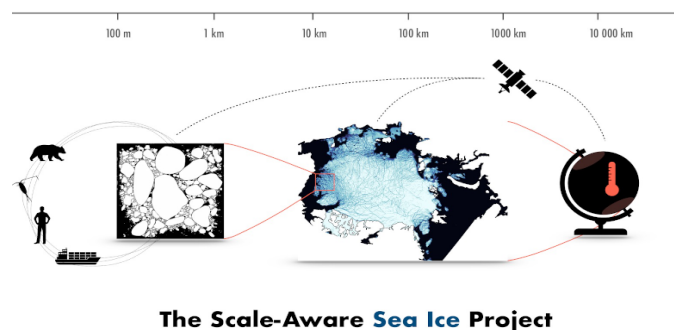


Introducing NEDAS: A Light-weight and Scalable Python Solution for Ensemble Data Assimilation

Yue (Michael) Ying

NERSC

Supported by: *SASIP, ACCIBERG*



Ensemble Data Assimilation

Model state ψ updated by observation φ to find best estimate

$$p(\psi|\varphi) = \frac{p(\varphi|\psi)p(\psi)}{p(\varphi)}$$

Use an ensemble of state $\Psi = (\psi_1, \dots, \psi_{N_e}) \in \mathbb{R}^{N_{state} \times N_e}$
as samples of $p(\psi)$

and “observation priors” $\Phi = (\varphi_1, \dots, \varphi_{N_e}) \in \mathbb{R}^{N_{obs} \times N_e}$
comparing with actual observation φ^o to give the likelihood $p(\varphi|\psi)$

How to update Ψ so that it characterizes $p(\psi|\varphi)$?

Algorithm: $\Psi \leftarrow \mathcal{A}(\Psi, \Phi, \varphi^o)$

How much effort is needed for testing novel algorithms in real models?

Simple method: just implement in the model code (WRF nudging/fdda)

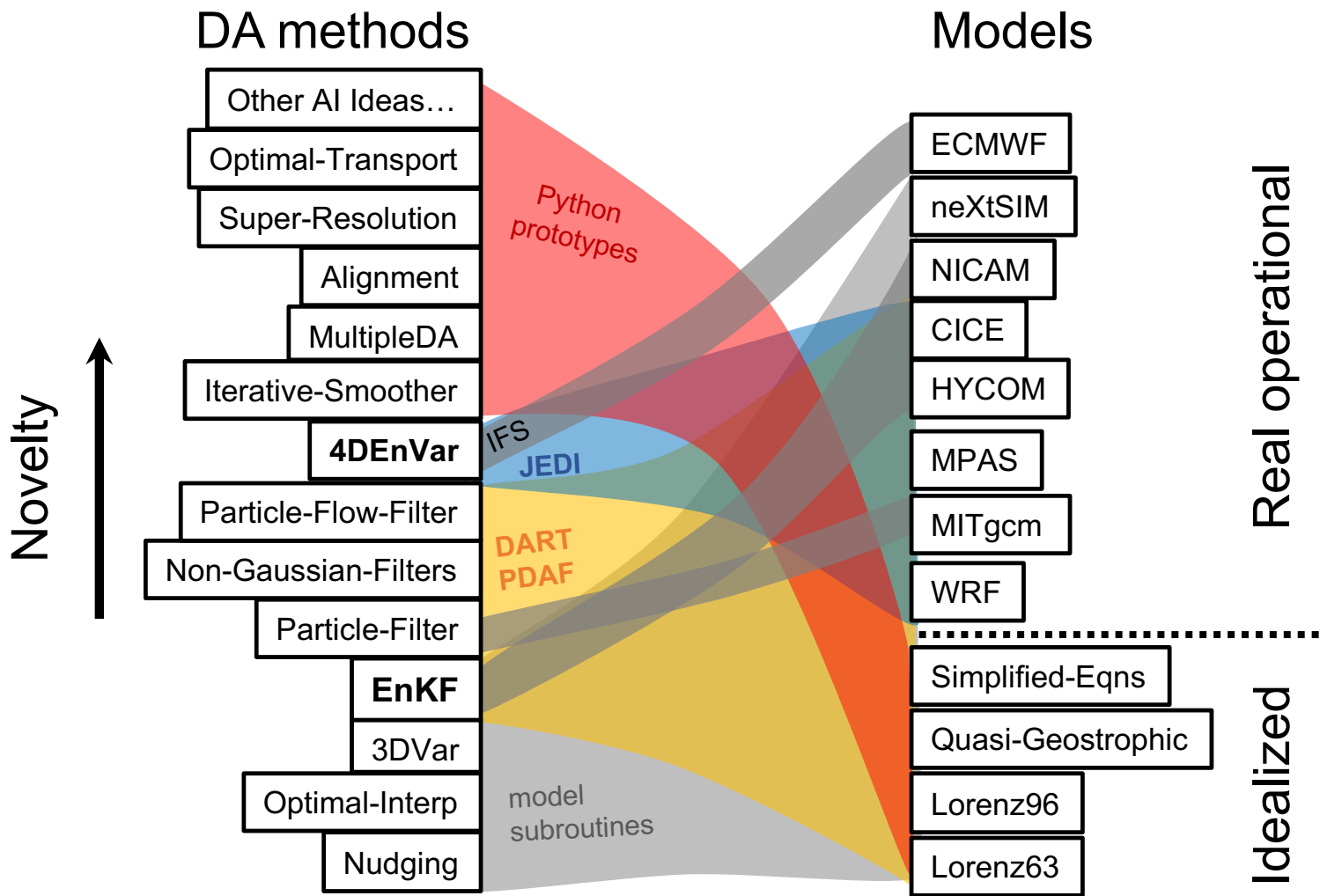
Complex method: dedicated DA software:

Data Assimilation Research Testbed (DART; Anderson et al. 2009)

Parallel Data Assimilation Framework (PDAF; Nerger & Hiller 2013)

Joint Effort in DA Integration (JEDI; JCSDA)

Conception → Python prototype → implement in *DART / PDAF / JEDI*
→ test in real model → operational use



New ideas for nonlinear filtering for large dimensional systems, but a lot of them stuck at Python prototype phase...

Next-generation Ensemble Data Assimilation System **NEDAS** enters the market

Conception → Python prototype → test in real models →
implement in DART / PDAF / JEDI → operational use

Python code is light-weight and easier to maintain

mpi4py, numpy, numba.jit allow scalability and efficiency

operating system integration, error handling and testing

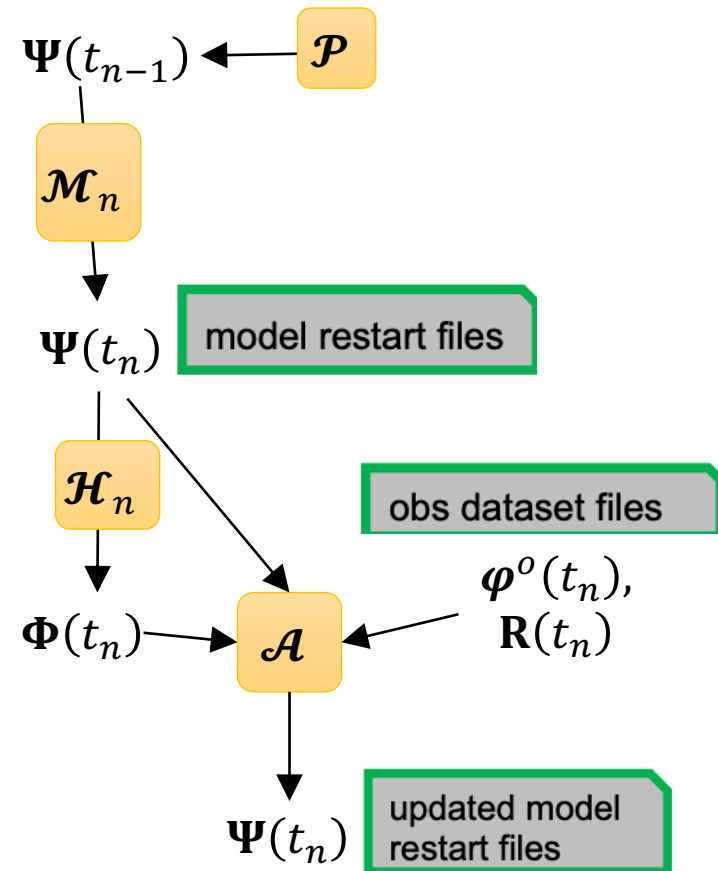
plenty of open-source libraries: machine learning, optimization...

NEDAS implementation

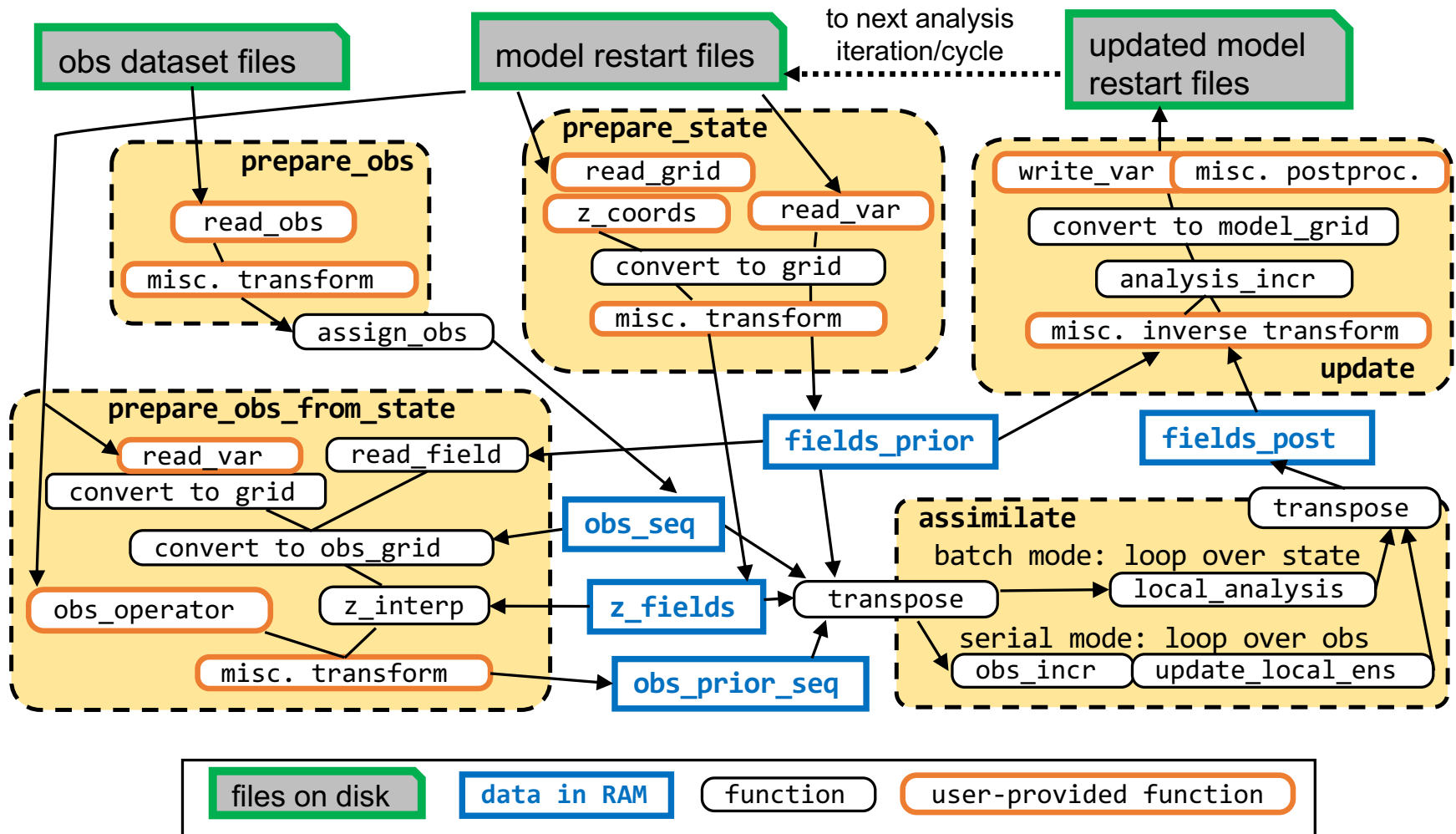
Sequential DA with pause-restart strategy

Require: $\Psi(t_0), \varphi^o(t_{1:N_t}), \mathbf{R}(t_{1:N_t})$

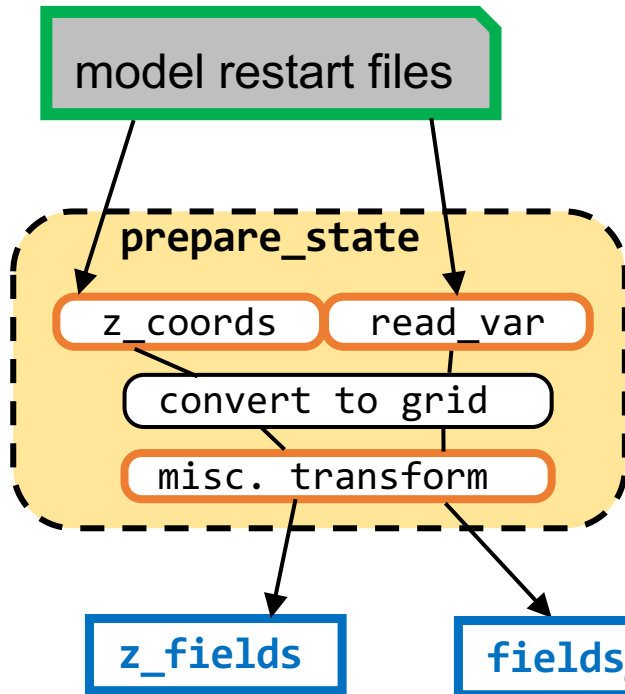
- 1: **for** $n = 1, \dots, N_t$ **do**
 - 2: $\Psi(t_{n-1}) \leftarrow \mathcal{P}[\Psi(t_{n-1})]$
 - 3: $\Psi(t_n) = \mathcal{M}_n[\Psi(t_{n-1})]$
 - 4: $\Phi(t_n) = \mathcal{H}_n[\Psi(t_n)]$
 - 5: $\Psi(t_n) \leftarrow \mathcal{A}[\Psi(t_n), \Phi(t_n), \varphi^o(t_n), \mathbf{R}(t_n)]$
 - 6: **end for**
 - 7: **return** $\Psi(t_{1:N_t})$
-



NEDAS implementation: more details



model.<model_name> modules



read_var: reading a variable **field** from restart

read_grid: getting the 2D grid

Grid.convert: supports many grid types; caches interpolation coeffs for efficiency

z_coords: getting the z coordinate fields

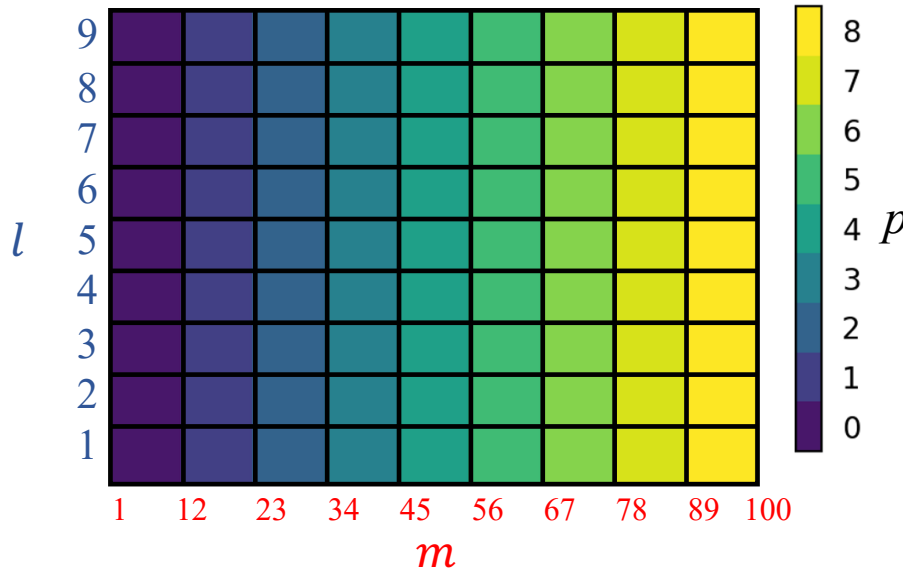
dimensions: member, location, field record
(x,y) (z,time,varname)

$$\Psi \left[\begin{array}{c} \mathbf{m}, \\ N_e \end{array}, \underbrace{\begin{array}{c} \mathbf{l}, \\ \mathbf{n} \end{array}}_{N_{\text{state}}} \right]$$

The smallest I/O task for a processor is to read a field with record index n and member index m

Parallel processing of state data

$$\Psi[m, l, n]$$



In this example:

$m = 1 \dots 100$ members

$l = 1 \dots 9$ spatial locations

$n = 1 \dots 4$ field records

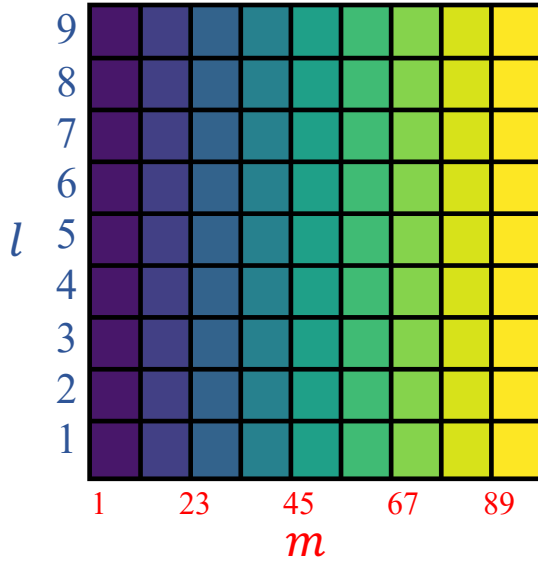
$Nq = 2$ groups of $Np = 9$ processors,
total 18 processors

The data is “**field-complete**”
(also called “state-complete” in
Anderson & Collins 2007)

$q =$	$p =$	1	2	...	9
1		1 data[1:11, 1:9, 1:2]	2 data[12:23, 1:9, 1:2]		9 data[89:100, 1:9, 1:2]
2		10 data[1:11, 1:9, 3:4]	11 data[12:23, 1:9, 3:4]		18 data[89:100, 1:9, 3:4]

field-complete

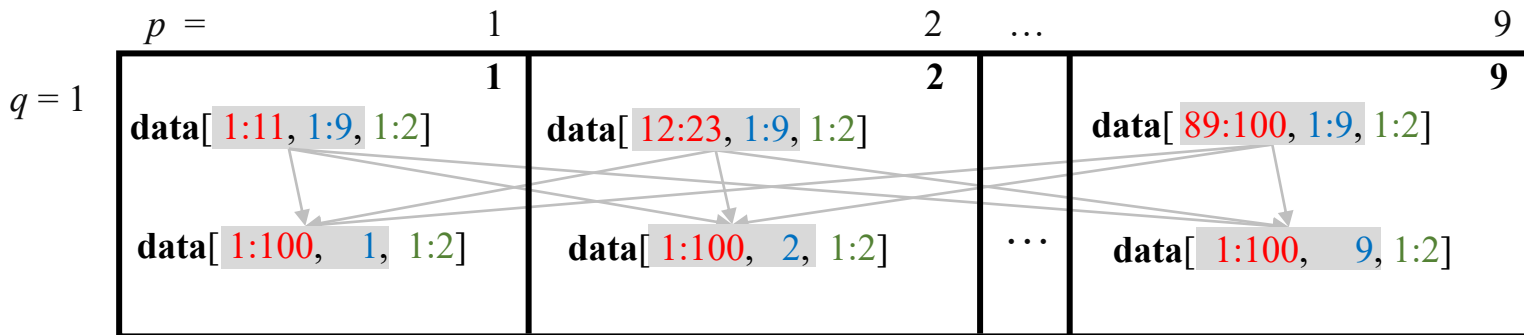
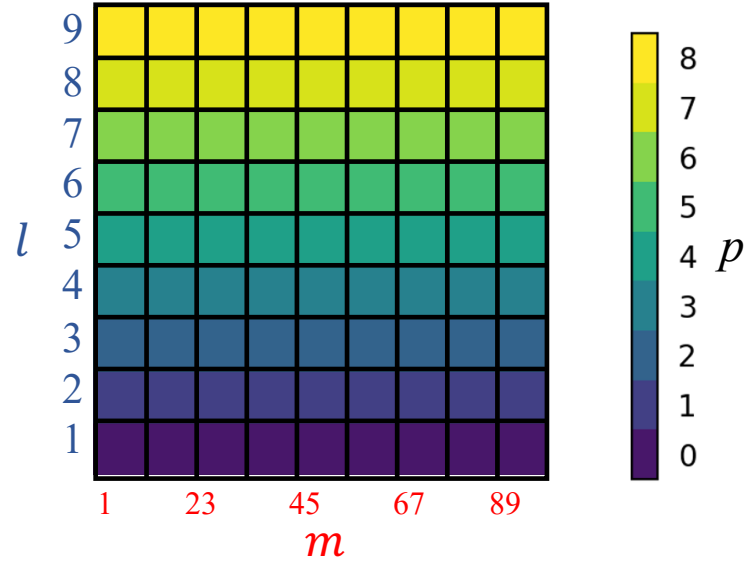
$$\Psi = (\psi_1, \dots, \psi_{N_e})$$



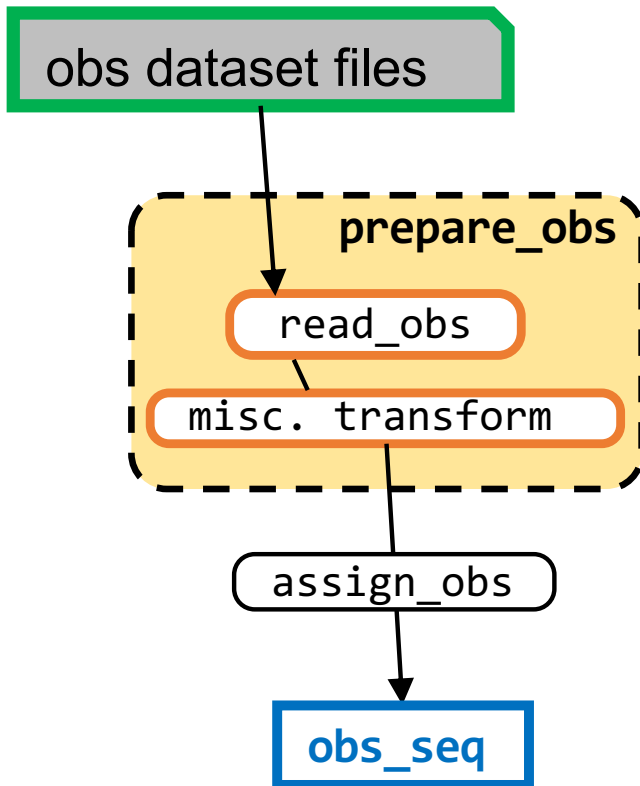
transpose
→

ensemble-complete

$$\Psi = (\psi_1^e, \dots, \psi_{N_{\text{state}}}^e)^T$$



dataset.<dataset_name> modules



read_obs reads obs dataset, preprocessing

random_network generates synthetic obs network

assign_obs finds indices of “local obs”

misc. transform: placeholder for algorithmic flexibility

dimensions

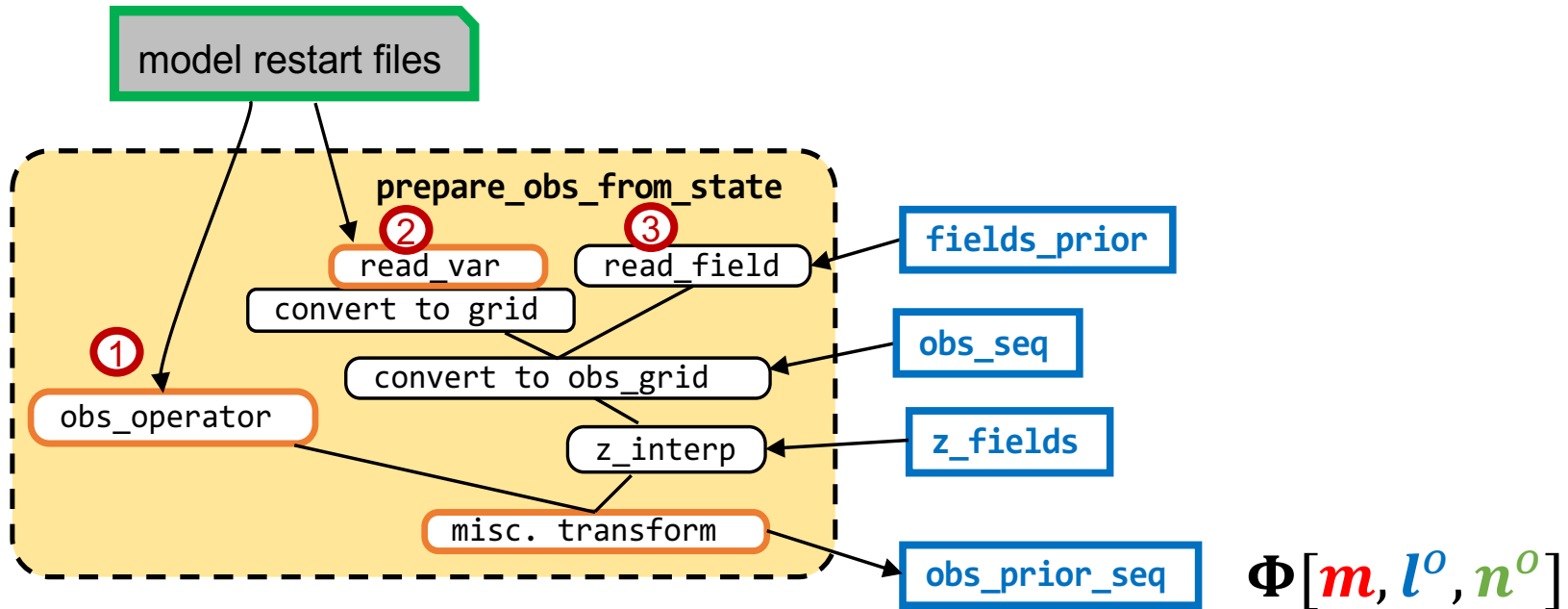
location
(x,y,z)

record
(time,varname)

$$\varphi^o [\quad l^o, \quad n^o \quad]$$

processor $p = 0$ is in charge of reading the obs

Observation operator $\Phi = \mathcal{H}(\Psi)$

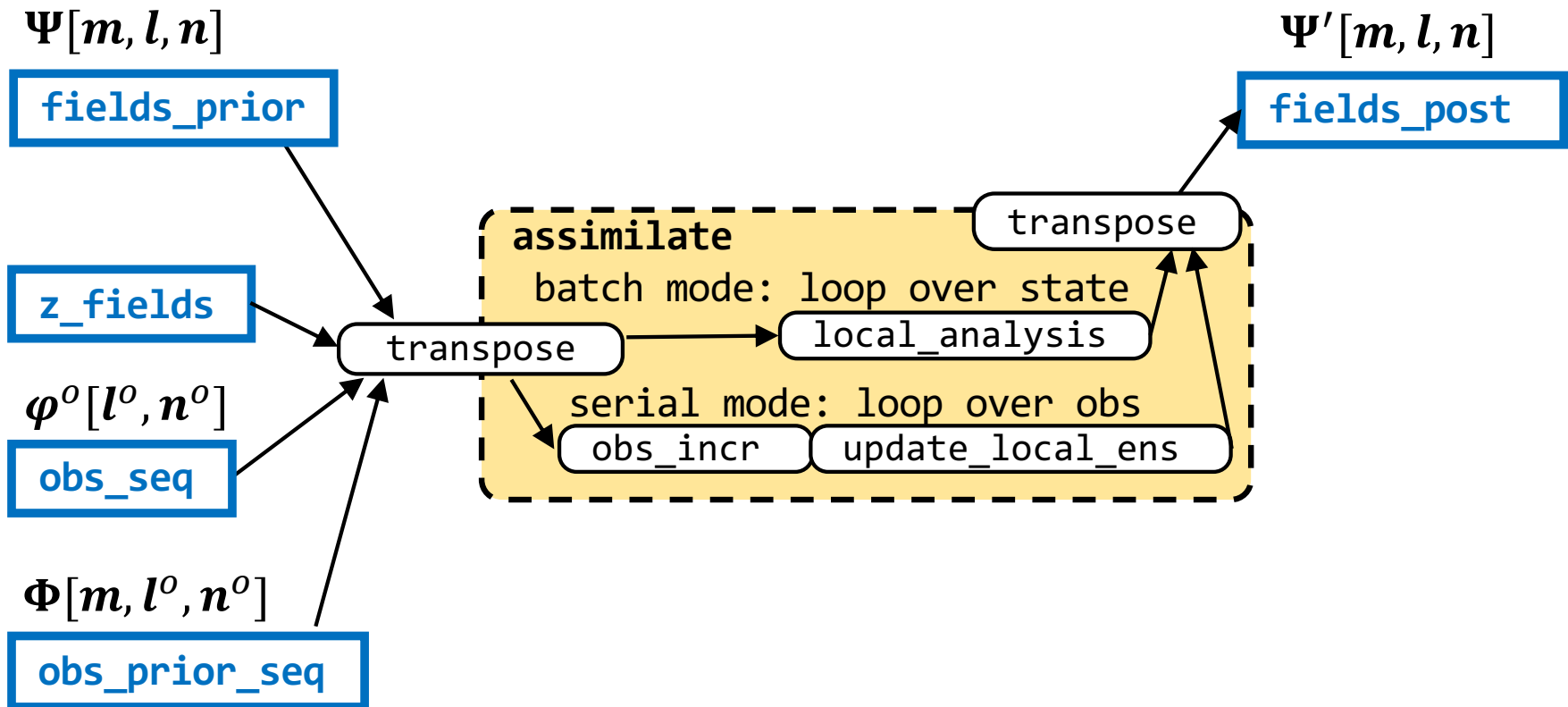


- ①: obs_operator provided by the dataset module (space-time integral, nonlinear function, etc.)
- ②: obs fields from read_var provided by the model module
- ③: obs is one of the state variable: just read from fields_prior

Same parallel processing as the state:

each processor reads an obs network with record index n^o and member index m

Assimilation algorithm $\mathcal{A}(\Psi, \Phi, \varphi^0)$



EnKF implementation

Make Gauss-linear assumption $p(\boldsymbol{\psi}) = \mathcal{N}(\bar{\boldsymbol{\psi}}, \mathbf{C}_{\boldsymbol{\psi}\boldsymbol{\psi}})$

where $\bar{\boldsymbol{\psi}} = \boldsymbol{\Psi}\mathbf{1}^T/N_e$ and $\mathbf{C}_{\boldsymbol{\psi}\boldsymbol{\psi}} = \boldsymbol{\Psi}(\boldsymbol{\Psi} - \bar{\boldsymbol{\psi}}\mathbf{1}^T)^T/(N_e - 1)$

and observation $\boldsymbol{\varphi}^o = \mathcal{H}(\boldsymbol{\psi}^{\text{tr}}) + \boldsymbol{\varepsilon}^o$, $\boldsymbol{\varepsilon}^o \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$

Algorithm:

$$\boldsymbol{\Phi} = \mathcal{H}(\boldsymbol{\Psi})$$

$$\boldsymbol{\Psi} \leftarrow \boldsymbol{\Psi} + \mathbf{C}_{\boldsymbol{\psi},\boldsymbol{\varphi}}(\mathbf{C}_{\boldsymbol{\varphi},\boldsymbol{\varphi}} + \mathbf{R})^{-1}(\boldsymbol{\Phi}^o - \boldsymbol{\Phi}) \quad (\text{perturbed-obs EnKF})$$

There are two different strategies in parallelization:

- batch assimilation
- serial assimilation

Parallelization strategy

Batch assimilation (e.g. PDAF)

for $i = 1, \dots, N_{\text{state}}$:

$$\mathbf{S} = \mathbf{R}^{-1/2}(\mathbf{\Phi} - \bar{\boldsymbol{\varphi}}\mathbf{1}^T) \circ (\boldsymbol{\rho}\mathbf{1}^T) / \sqrt{N_e - 1}$$

$$\mathbf{d} = \mathbf{R}^{-1/2}(\boldsymbol{\varphi}^o - \bar{\boldsymbol{\varphi}}) \circ \boldsymbol{\rho} / \sqrt{N_e - 1}$$

$$\mathbf{\Xi} = (\mathbf{I} + \mathbf{S}^T\mathbf{S})^{-1}$$

$$\mathbf{T} = \mathbf{\Xi}\mathbf{S}^T\mathbf{d}\mathbf{1}^T + \mathbf{\Xi}^{1/2}$$

$$\boldsymbol{\psi}_i^{eT} \leftarrow \boldsymbol{\psi}_i^{eT}\mathbf{T}$$

return $\boldsymbol{\Psi} = (\boldsymbol{\psi}_1^e, \dots, \boldsymbol{\psi}_{N_{\text{state}}}^e)^T$

cost: $\mathcal{O}(N_{\text{lobs}}N_e^2 + N_e^3) \times N_{\text{state}}$

“local analysis”

Serial assimilation (e.g. DART)

for $j = 1, \dots, N_{\text{obs}}$:

$$\xi = \sigma_{o,j}^2 / (\sigma_j^2 + \sigma_{o,j}^2)$$

$$\boldsymbol{\delta}_j^e = \xi\bar{\boldsymbol{\varphi}}_j + (1 - \xi)\boldsymbol{\varphi}_j^o + \sqrt{\xi}(\boldsymbol{\varphi}_j^e - \bar{\boldsymbol{\varphi}}_j) - \boldsymbol{\varphi}_j^e$$

broadcast $\boldsymbol{\delta}_j^e$

$$\boldsymbol{\Psi} \leftarrow \boldsymbol{\Psi} + (\boldsymbol{\rho}^\psi \circ \mathbf{c}_{\psi,\varphi_j} / \sigma_j^2) \boldsymbol{\delta}_j^{eT}$$

$$\boldsymbol{\Phi} \leftarrow \boldsymbol{\Phi} + (\boldsymbol{\rho}^\varphi \circ \mathbf{c}_{\varphi,\varphi_j} / \sigma_j^2) \boldsymbol{\delta}_j^{eT}$$

return $\boldsymbol{\Psi}$

cost:

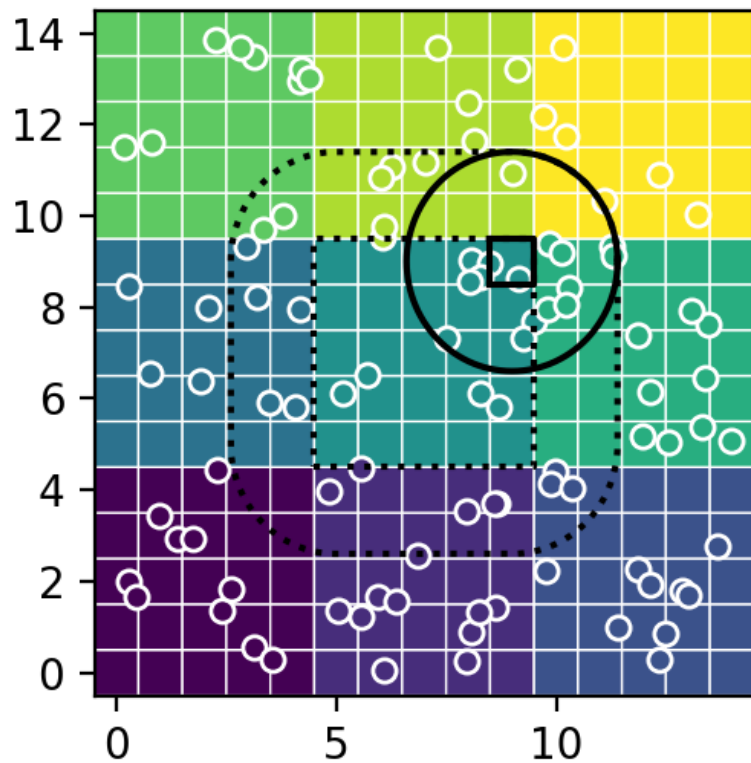
$\mathcal{O}(N_e \log N_p + N_e N_{\text{state}} + N_e N_{\text{lobs}}) \times N_{\text{obs}}$

“obs_incr”: nonlinear filters possible

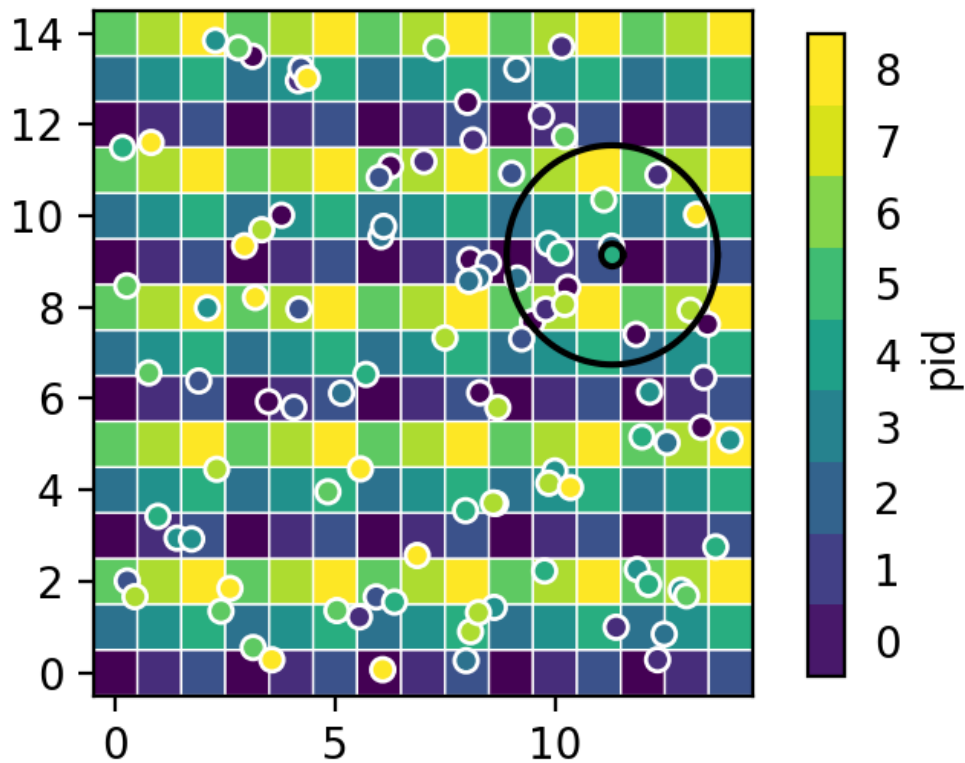
“obs_updates_ens”: linear, probit-space

Memory layout for state/obs

(a) batch assimilation



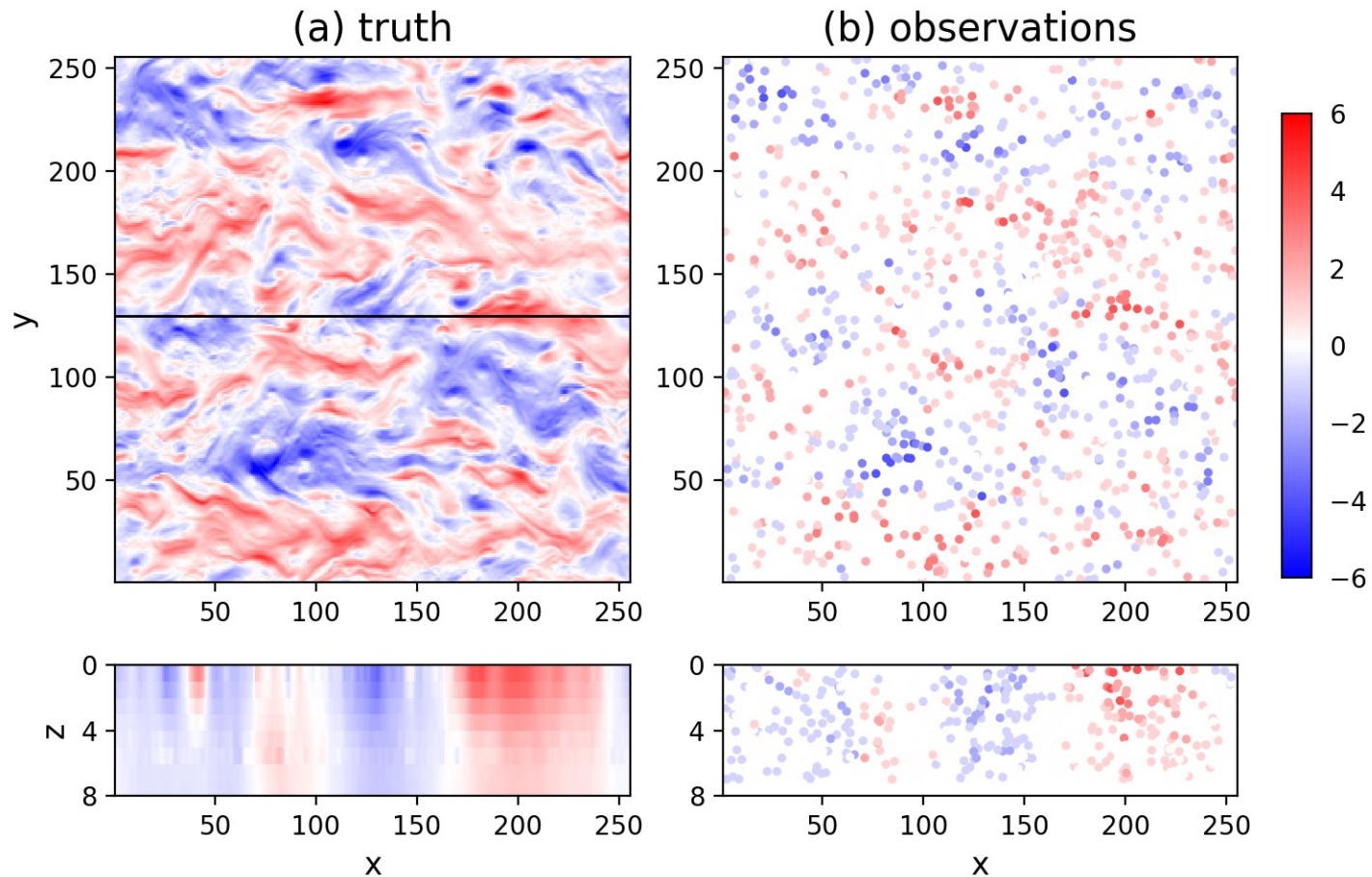
(b) serial assimilation



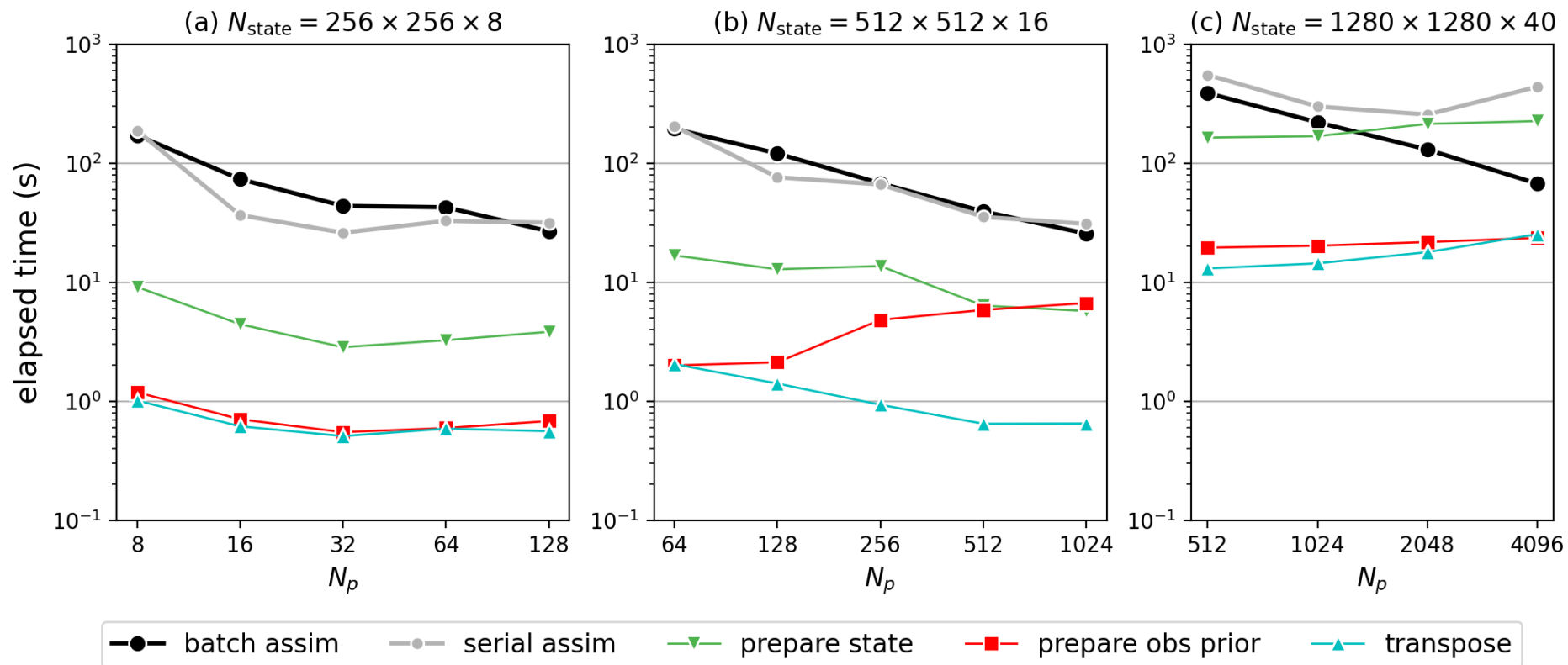
Benchmarking: QG model example

state and observations: velocity fields

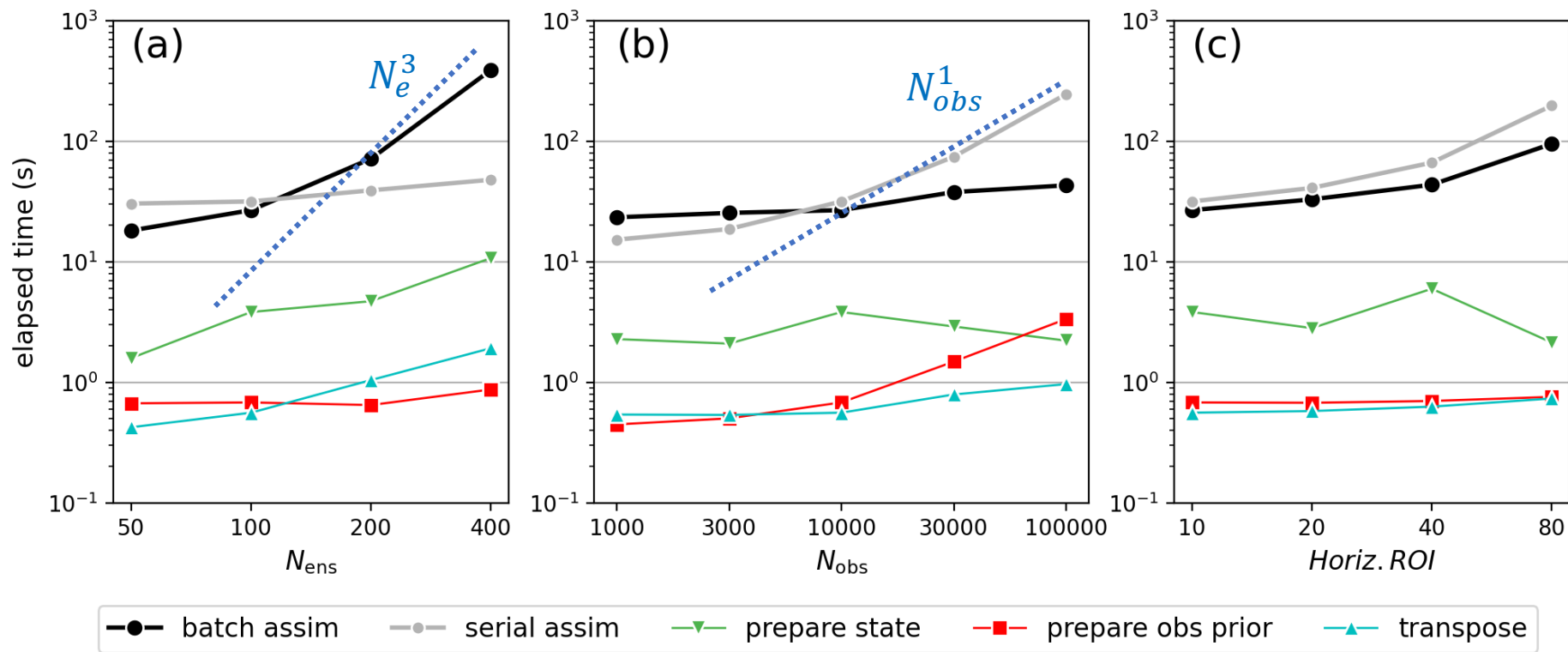
Nstate = 256x256x8, Nobs = 10000, obs_err = 0.5, Ne = 100, hroi=10, vroi=5



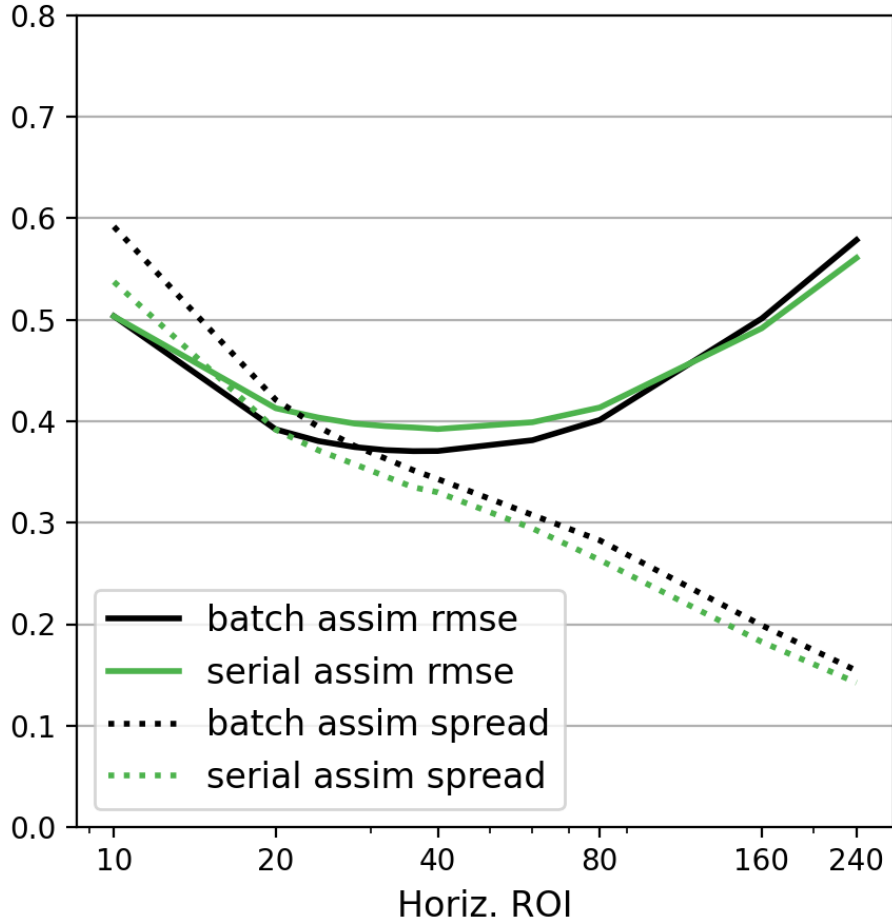
Scalability of \mathcal{A} as N_p increases



How \mathcal{A} scales as dimensionality increases



Analysis error/spread comparison



Both strategies produced comparable results

Serial assimilation fits more to observations (lower posterior spread); slightly less accurate (higher rmse)

Consistent with previous findings (Holland & Wang 2012; Nerger 2015).

Algorithmic flexibility: Miscellaneous transform functions \mathcal{T}

1: for $s = 1, \dots, N_s$ do

2: $\tilde{\varphi}^o = \mathcal{T}_s^o(\varphi^o)$

3: for $m = 1, \dots, N_e$ do

4: $\tilde{\varphi}_m = \mathcal{T}_s^o(\varphi_m)$

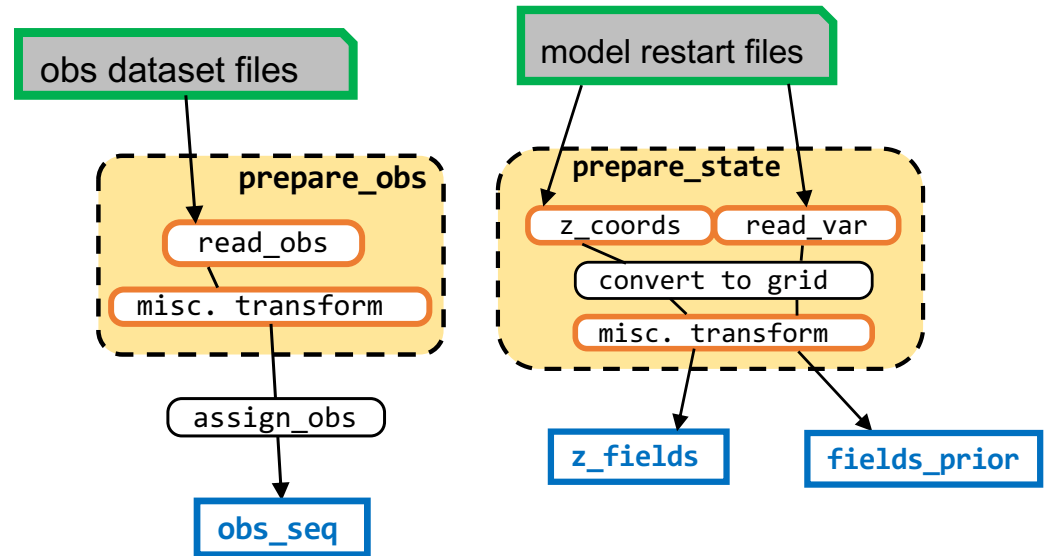
5: $\tilde{\psi}_m = \mathcal{T}_s(\psi_m)$

6: end for

7: $\tilde{\Psi} = (\tilde{\psi}_1, \dots, \tilde{\psi}_{N_e})$

8: $\tilde{\Phi} = (\tilde{\varphi}_1, \dots, \tilde{\varphi}_{N_e})$

9: $(\tilde{\psi}'_1, \dots, \tilde{\psi}'_{N_e}) = \tilde{\Psi}' = \tilde{\mathcal{A}}(\tilde{\Psi}, \tilde{\Phi}, \tilde{\varphi}^o, \tilde{\mathbf{R}}_s, \mathbf{r}_s^\psi, \mathbf{r}_s^\varphi, \mathbf{L}_s)$



- Gaussian anamorphosis (Simon & Bertino 2009)
- Multiscale decomposition (Ying 2019, 2020), gcm-filters (Grooms et al. 2021)
- Super-resolution (Barthelemy et al 2022)
- Mapping to latent space (Chipilski 2023)

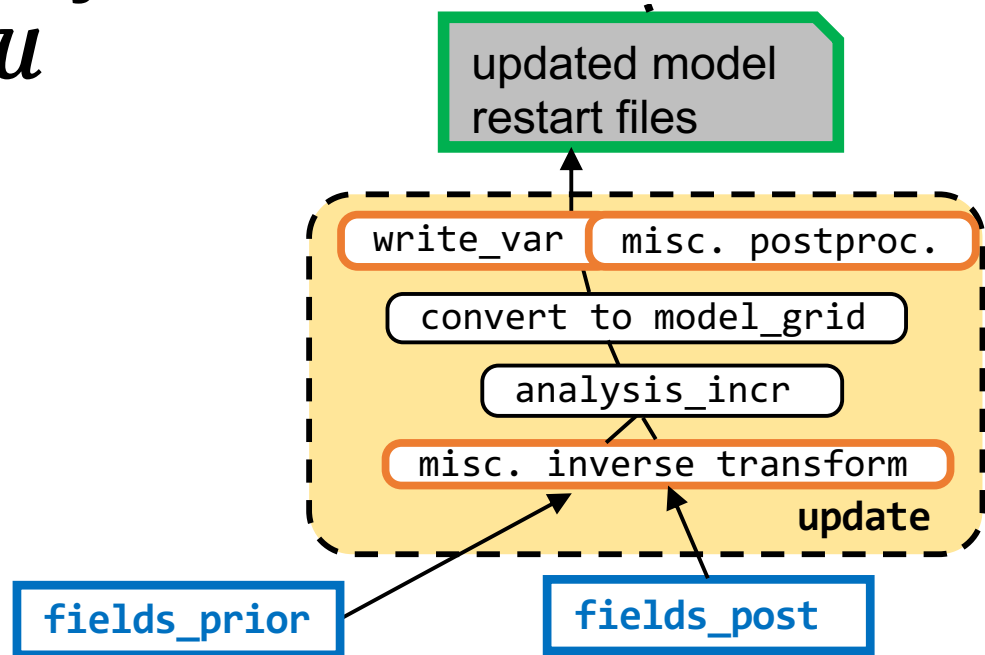
...

Algorithmic flexibility: Update functions \mathcal{U}

```

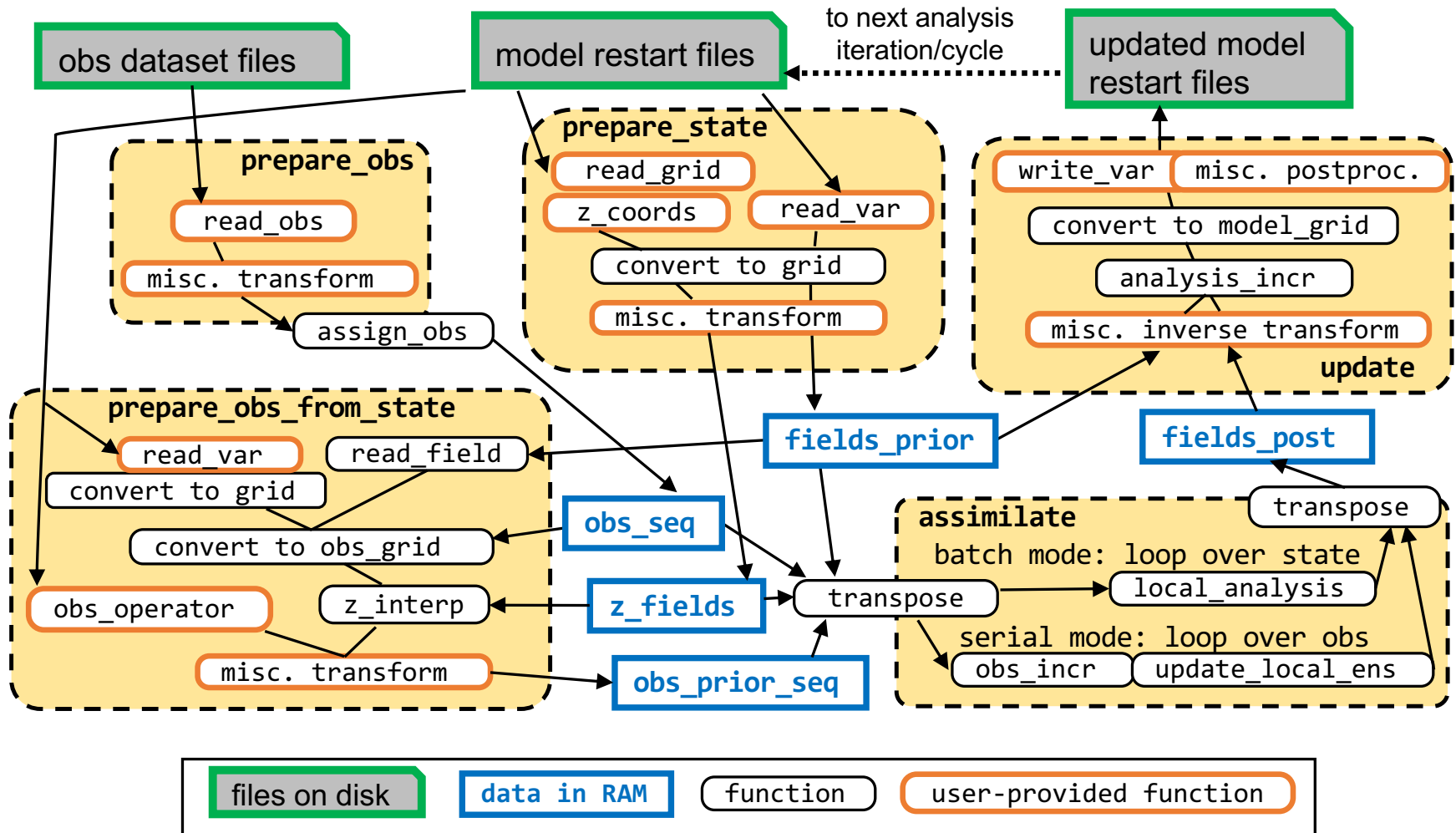
1: for  $s = 1, \dots, N_s$  do
2:    $\tilde{\varphi}^o = \mathcal{T}_s^o(\varphi^o)$ 
3:   for  $m = 1, \dots, N_e$  do
4:      $\tilde{\varphi}_m = \mathcal{T}_s^o(\varphi_m)$ 
5:      $\tilde{\psi}_m = \mathcal{T}_s(\psi_m)$ 
6:   end for
7:    $\tilde{\Psi} = (\tilde{\psi}_1, \dots, \tilde{\psi}_{N_e})$ 
8:    $\tilde{\Phi} = (\tilde{\varphi}_1, \dots, \tilde{\varphi}_{N_e})$ 
9:    $(\tilde{\psi}'_1, \dots, \tilde{\psi}'_{N_e}) = \tilde{\Psi}' = \tilde{\mathcal{A}}(\tilde{\Psi}, \tilde{\Phi}, \tilde{\varphi}^o, \tilde{\mathbf{R}}_s, \mathbf{r}_s^\psi, \mathbf{r}_s^\varphi, \mathbf{L}_s)$ 
10:  for  $m = 1, \dots, N_e$  do
11:     $\psi_m \leftarrow \mathcal{U}_s(\psi_m, \tilde{\psi}_m, \tilde{\psi}'_m)$ 
12:  end for
13: end for

```



- Inverse transform, additive increments
- Alignment techniques (Ying 2019)
- Update not only model states, but also hyperparameters.

NEDAS analysis workflow

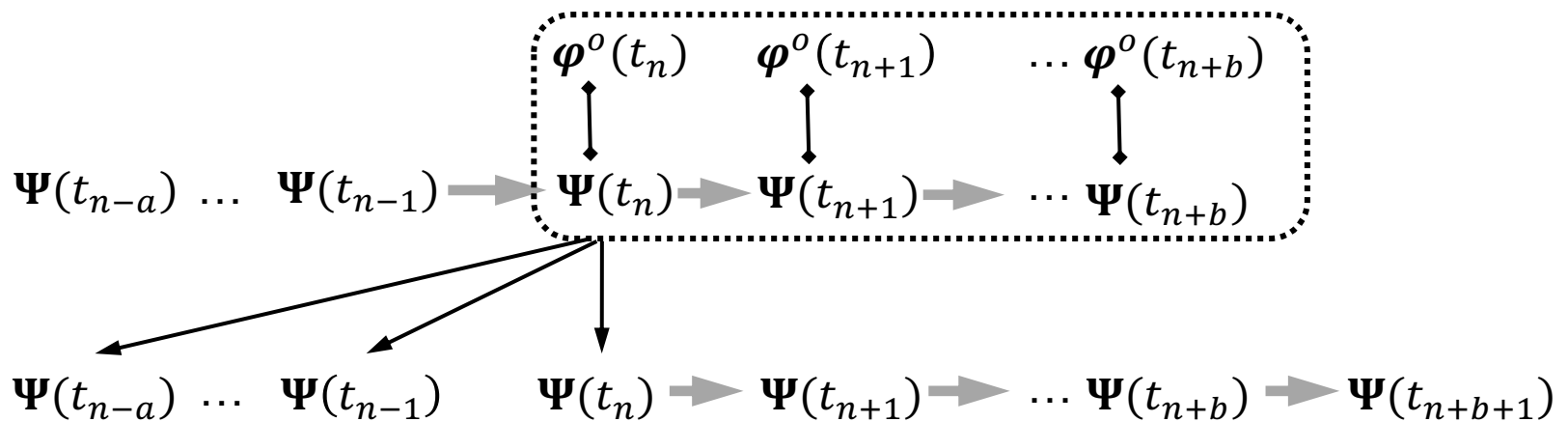


NEDAS sequential DA includes 4D model states

```

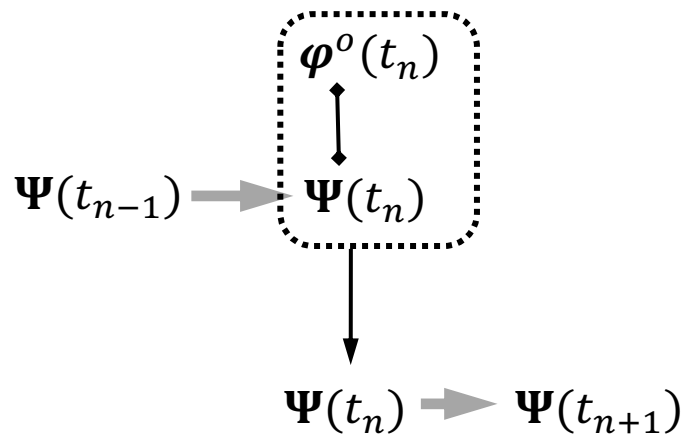
1: for  $n = 1, \dots, N_t$  do
2:    $\Psi(t_{n-1}) \leftarrow \mathcal{P}[\Psi(t_{n-1})]$ 
3:   for  $k = 0, \dots, b$  do
4:      $\Psi(t_{n+k}) = \mathcal{M}_{n+k} [\Psi(t_{n+k-1})]$ 
5:      $\Phi(t_{n+k}) = \mathcal{H}_{n+k} [\Psi(t_{n+k})]$ 
6:   end for
7:    $\Psi(t_{n-a:n}) \leftarrow \mathcal{A}[\Psi(t_{n-a:n}), \Phi(t_{n:n+b}), \varphi^o(t_{n:n+b}), (b+1)\mathbf{R}(t_{n:n+b})]$ 
8: end for
9: return  $\Psi(t_{1:N_t})$ 

```



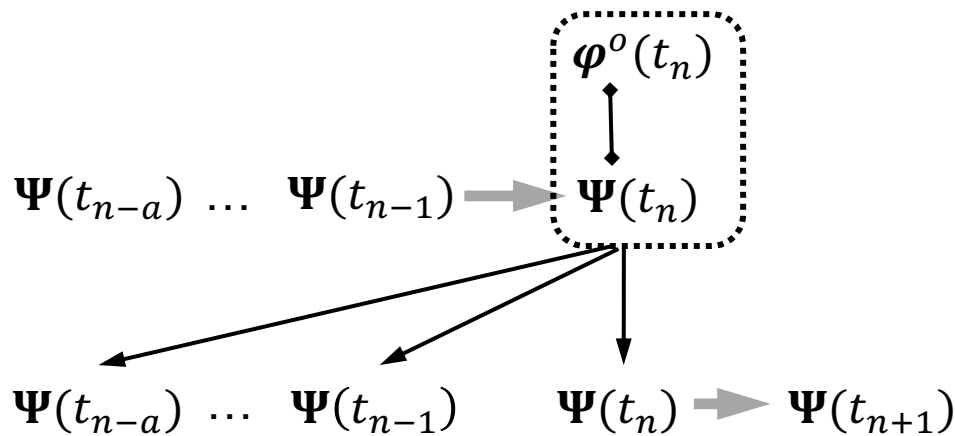
$a = b = 0$: filter

- 1: **for** $n = 1, \dots, N_t$ **do**
 - 2: $\Psi(t_{n-1}) \leftarrow \mathcal{P}[\Psi(t_{n-1})]$
 - 3: **for** $k = 0, \dots, b$ **do**
 - 4: $\Psi(t_{n+k}) = \mathcal{M}_{n+k} [\Psi(t_{n+k-1})]$
 - 5: $\Phi(t_{n+k}) = \mathcal{H}_{n+k} [\Psi(t_{n+k})]$
 - 6: **end for**
 - 7: $\Psi(t_{n-a:n}) \leftarrow \mathcal{A}[\Psi(t_{n-a:n}), \Phi(t_{n:n+b}), \varphi^o(t_{n:n+b}), (b+1)\mathbf{R}(t_{n:n+b})]$
 - 8: **end for**
 - 9: **return** $\Psi(t_{1:N_t})$
-



$b = 1, a > 0$: recursive smoother

1: **for** $n = 1, \dots, N_t$ **do** Evensen & van Leeuwen 2000
2: $\Psi(t_{n-1}) \leftarrow \mathcal{P}[\Psi(t_{n-1})]$
3: **for** $k = 0, \dots, b$ **do**
4: $\Psi(t_{n+k}) = \mathcal{M}_{n+k} [\Psi(t_{n+k-1})]$
5: $\Phi(t_{n+k}) = \mathcal{H}_{n+k} [\Psi(t_{n+k})]$
6: **end for**
7: $\Psi(t_{n-a:n}) \leftarrow \mathcal{A}[\Psi(t_{n-a:n}), \Phi(t_{n:n+b}), \varphi^o(t_{n:n+b}), (b+1)\mathbf{R}(t_{n:n+b})]$
8: **end for**
9: **return** $\Psi(t_{1:N_t})$



$b = 1, a = 0$: one-step-ahead smoother

1: **for** $n = 1, \dots, N_t$ **do**

2: $\Psi(t_{n-1}) \leftarrow \mathcal{P}[\Psi(t_{n-1})]$

3: **for** $k = 0, \dots, b$ **do**

4: $\Psi(t_{n+k}) = \mathcal{M}_{n+k} [\Psi(t_{n+k-1})]$

5: $\Phi(t_{n+k}) = \mathcal{H}_{n+k} [\Psi(t_{n+k})]$

6: **end for**

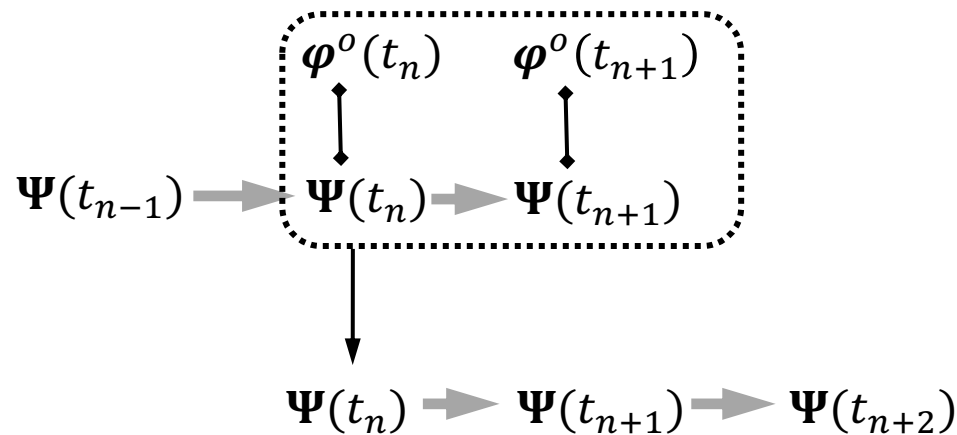
7: $\Psi(t_{n-a:n}) \leftarrow \mathcal{A}[\Psi(t_{n-a:n}), \Phi(t_{n:n+b}), \varphi^o(t_{n:n+b}), (b+1)\mathbf{R}(t_{n:n+b})]$

8: **end for**

9: **return** $\Psi(t_{1:N_t})$

Desbouvries et al. 2011; Gharamti et al. 2015;
Ait-El-Fquih & Hoteit 2022

More general smoother formulation:
Khare et al., 2008; Bocquet & Sakov, 2014;
Grudzien & Bocquet, 2022



Summary

- NEDAS provides a light-weight Python solution for testing ensemble DA algorithms in real model settings.
- Both batch and serial assimilation strategies are efficient and scalable for large models, but they are favored in different scenarios.
- NEDAS has a modular design, which gives flexibility for integrating new algorithmic ideas in its workflow.

Code publicly available:

<https://github.com/nansencenter/NEDAS>

Manuscript submitted to JAMES:

Ying: “Introducing NEDAS: a light-weight and scalable Python solution for ensemble data assimilation”

Contact me: yue.ying@nersc.no