### Advancing the TOPAZ ocean and sea ice data

### assimilation system: Algorithmic innovations





Yue (Michael) Ying



EnKF Workshop 2025, Ullensvang, Norway

## **Data assimilation software at NERSC**



Previously, each model has a specific DA software separately maintained:



NEDAS now provides support to all NERSC models

- Rapid research development
- Access to rich Python ecosystem for scientific computation, machine learning etc.



NorESM, HYCOM-CICE, neXtSIM, ECOSMO, WRF, ...



EnKF (different flavors: ETKF or EAKF), EnKS, non-Gaussian filters, alignment, emulators...



## Python tools for high-performance computation

213%



- drives most of the linear algebra NumPy
- mpi4py provides support for MPI parallelism
- JIT technique (numba, JAX, taichi) to convert python functions to machine code



Versatile Ocean Simulation in Pure Python

#### Taichi vs. CUDA **High-performance** Taichi's JIT compiler automatically compiles Python functions into fast GPU 121% 112% or CPU machine code for parallel execution. While Taichi lives in Python, it can approach or even outrun the speed of C++ or CUDA. Stencil N-Body

### @numba.njit

def obs\_increment\_eakf(obs\_prior, obs, obs\_err) -> np.ndarray: nens = obs\_prior.size

##obs error variance obs\_var = obs\_err\*\*2

##obs\_prior separate into mean+perturbation obs\_prior\_mean = np.mean(obs\_prior) obs\_prior\_pert = obs\_prior - obs\_prior\_mean

##compute prior error variance obs\_prior\_var = np.sum(obs\_prior\_pert\*\*2) / (nens-1)

"""ensemble adjustment Kalman filter (Anderson 2003)""" var\_ratio = obs\_var / (obs\_prior\_var + obs\_var)

##new mean is weighted average between obs\_prior\_mean and obs obs\_post\_mean = var\_ratio \* obs\_prior\_mean + (1 - var\_ratio) \* obs

##new pert is adjusted by sqrt(var\_ratio), a deterministic square-r obs\_post\_pert = np.sqrt(var\_ratio) \* obs\_prior\_pert

##assemble the increments obs\_incr = obs\_post\_mean + obs\_post\_pert - obs\_prior

return obs\_incr



### Python as control scripts (thanks to "Code Clinic" by ICCS)

• YAML configuration files

nens: 30
run\_analysis: True
run\_diagnose: True
debug: True

time\_start: 2021-07-05T00:00:00Z time\_end: 2021-08-09T00:00:00Z time\_analysis\_start: 2021-07-12T00:00:00Z time\_analysis\_end: 2021-08-09T00:00:00Z cycle\_period: 168

model\_def: topaz.v5: config\_file: '{nedas\_root}/models/topaz/v5/default.yml' model\_env: '{nedas\_root}/models/topaz/env/setup.src' reanalysis\_code: '/cluster/home/vingvue/code/ReanalysisTP5' basedir: '/cluster/work/users/yingyue/TP5a0.06' nproc\_per\_run: 512 nproc\_per\_util: 50 walltime: 3600 restart\_dt: 168 forcing\_dt: 6 ens\_run\_type: scheduler use\_job\_array: False stagnant\_log\_timeout: 100 ens\_init\_dir: '/cluster/work/users/yingyue/26118/FORECAST' forcing\_file: '/cluster/work/users/yingyue/TP5a0.06/force/syr truth dir: ''

analysis\_scheme: offline\_filter
assimilator\_def:
 type: TopazDEnKF

- Better error handling
- Multiplatform support: HPC, jupyter-lab, containers

```
##build the shell command line
model_exe = os.path.join(self.basedir, f'expt_{self.X}', 'build', f'src_{self.V}
shell_cmd = ""
if self.model_env:
    shell_cmd = ". "+self.model_env+"; " ##enter topaz5 env
shell_cmd += "cd "+run_dir+"; "
                                            ##enter run directory
shell_cmd += 'JOB_EXECUTE '+model_exe+" >& run.log"
##run the model, give it 3 attempts
for i in range(3):
    try:
        run_job(shell_cmd, job_name='topaz5', run_dir=run_dir,
                nproc=self.nproc, offset=task_id*self.nproc_per_run,
                walltime=self.walltime, log_file=log_file, **kwargs)
    except RuntimeError as e:
        print(f"{e}, retrying ({2-i} attempts remain)")
        run_command(f"cp {log_file} {log_file}.attempt{i}")
        continue
    ##check output
    if find_keyword_in_file(log_file, 'Exiting hycom_cice'):
        run_success = True
        break
```



### The offline filter (using model restart files)

#### class OfflineFilterAnalysisScheme:

def filter(self, c):
 for c.iter in range(c.niter):
 c.transform\_funcs = assim\_tools.transforms.get\_transform\_funcs(c)
 state = assim\_tools.state.get\_state(c)
 timer(c)(state.prepare\_state)(c)
 obs = assim\_tools.obs.get\_obs(c, state)
 timer(c)(obs.prepare\_obs)(c, state)
 timer(c)(obs.prepare\_obs\_from\_state)(c, state, 'prior')
 assim\_lator = assim\_tools.assimilators.get\_assimilator(c)
 timer(c)(assimilator.assimilate)(c, state, obs)

updator) = assim\_tools.updators.get\_updator(c)
timer(c)(updator.update)(c, state)

for *s* in 1, ..., *N<sub>s</sub>*:

 $\mathbf{x}_{s}^{b} = F_{s}(\mathbf{x}^{b})$   $\mathbf{y}_{s}^{o} = F_{s}^{o}(\mathbf{y}^{o})$   $\mathbf{y}_{s}^{b} = F_{s}^{o}[\mathcal{H}(\mathbf{x}^{b})]$   $\delta \mathbf{x}_{s} = \mathcal{A}(\mathbf{x}_{s}^{b}, \mathbf{y}_{s}^{b}, \mathbf{y}_{s}^{o}, \mathbf{R}_{s}, \boldsymbol{\rho}_{s})$   $\mathbf{x}^{b} = \mathcal{U}(\mathbf{x}^{b}, \delta \mathbf{x}_{s})$   $\mathbf{x}^{a} \leftarrow \mathbf{x}^{b}$ 



### **Analysis schemes**





### **Assimilator Classes**





# **TopazDEnKFAssimilator**

DEnKF (Sakov et al. 2012): 100-member ensemble without vertical localization: Compute ensemble transform weights once, applies it to 50 levels

In NEDAS (ETKF available): Ifactor = hlfactor \* vlfactor \* tlfactor \* impact\_on\_state

Exact replica of TOPAZ with similar runtime efficiency:

- Cache the transform weights for same localization factor
- Early exit when localization factor is zero

Computing the distance itself can be costly:

- Use L1 distance for subsetting, then refine with L2 distance
- How to keep distance-free types of algorithm efficient?





### Uncertainty estimates versus model resolution – the dilemma





image credit: Laurent Bertino

30-100 ensemble members = factor 3-4 in resolution

### **Do we need vertical localization in TOPAZ?**



- Even with 100 members, there are still sampling noises
- Smaller ensemble size
   + vertical localization: cost effective?



### **QCEFAssimilator as an alternative**



Non-Gaussian error distribution:

- Sea ice variables
- BGC tracers

Current solution: postprocessing (fixhycom)

Quantile Conserving Ensemble Filter (Anderson 2023) expected to better handle bounded error distribution and nonlinearity

Porting QCEF subroutines to NEDAS using **f2py** 



### **TOPAZ perturbation scheme**

Evensen 1994 approach to draw random fields: specified hcorr and variance

Real forcast errors grow in time, but our perturbations are fixed amplitude



### **Additional perturbation schemes**





### Al enhancements (low-hanging fruits compared to end-to-end)



### AlignmentUpdator



Applied to assimilation of deformation in neXtSIM (to be coupled to HYCOM in next version TOPAZ)



### AlignmentUpdator in use for s = 1:



### AlignmentUpdator





### **Alternative optical flow algorithms**

In Ying 2019: HornSchunk with pyramid method

from OpenCV: Farneback (parametric dense flow) and DIS (based on deep learning)

```
import cv2
frame1 = ((img1+4) / 8 * 256).astype(np.uint8)
frame2 = ((img2+4) / 8 * 256).astype(np.uint8)
flow = cv2.calcOpticalFlowFarneback(frame1, frame2, None, 0.5, 3, 15, 3, 5, 1.2, 0)
```

```
dis = cv2.DISOpticalFlow_create(cv2.DISOPTICAL_FLOW_PRESET_FAST)
```

```
flow = dis.calc(frame1, frame2, None)
```



### **Alternative optical flow algorithms**



NERSC



Advancing TOPAZ assimilation algorithm:

- Vertical localization (finally!) and cross-variable impact factor
- Trying new perturbation schemes
- QCEF for sea ice and BGC variables (via f2py)
- Alignment technique (better alternatives from OpenCV)
- ...Any suggestions?

