

PDAF – Features and Recent Developments

Lars Nerger

Alfred Wegener Institute, Helmholtz Center for Polar and Marine Research
Bremerhaven, Germany

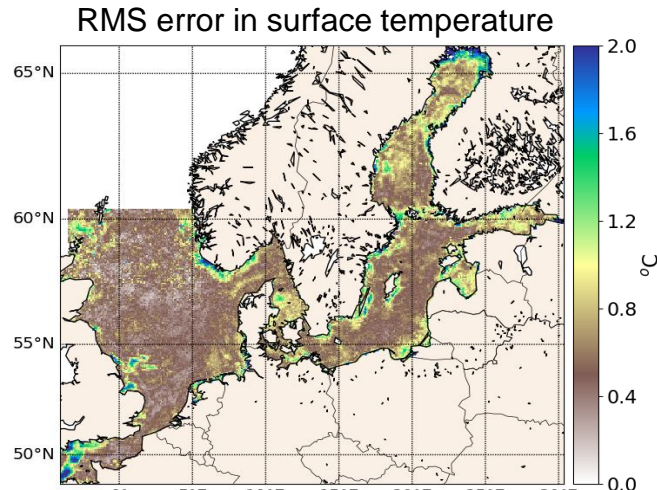
A tool for data assimilation ...

- program library for ensemble data assimilation
- provide support for parallel ensemble forecasts
- provide filters and smoothers - fully-implemented & parallelized (EnKF, LETKF, LESTKF, NETF, PF ... easy to add more)
- easily useable with (probably) any numerical model (coupled to e.g. NEMO, MITgcm, FESOM/AWI-CM, ICON, SCHISM/ESMF)
- run from notebooks to supercomputers (Fortran, MPI & OpenMP)
- usable for real assimilation applications and to study assimilation methods
- ~500 registered users; community contributions

Open source:
Code, documentation, and tutorial available at
<http://pdaf.awi.de>

HBM-ERGOM:

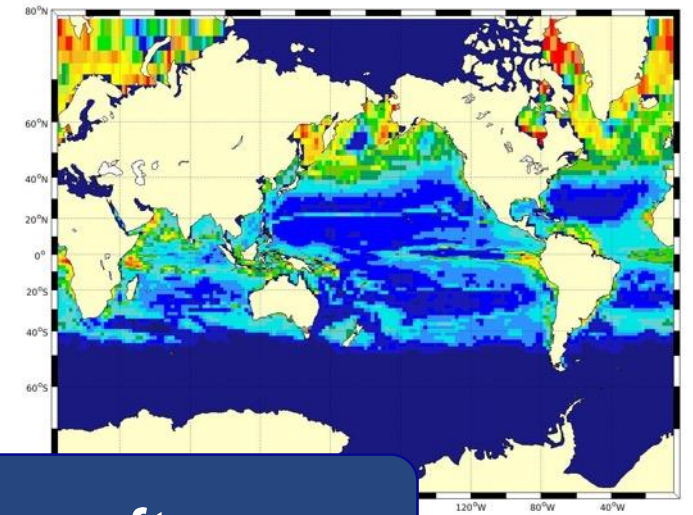
coupled physics/
biogeochemistry
coastal assimilation
(Goodliff et al., 2019)



MITgcm-REcoM:

global ocean color
assimilation into
biogeochemical
model
(Pradhan et al., 2019/20)

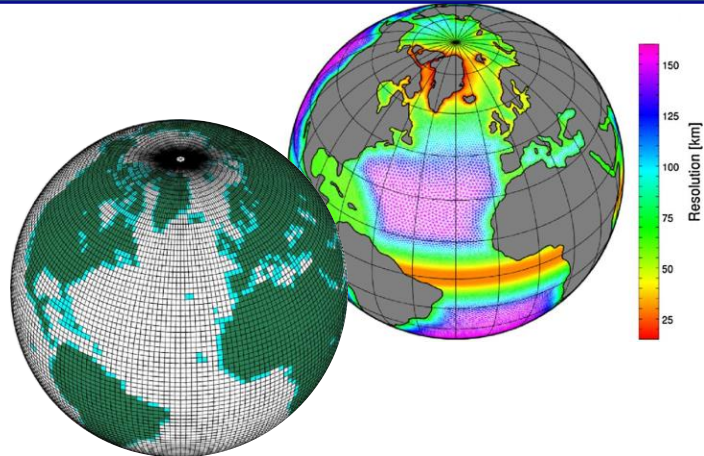
Total chlorophyll concentration June 30, 2012



Different models – same assimilation software

AWI-CM:

coupled atmos.-
ocean assimilation
(Tang et al., 2020
Mu et al., 2020
Nerger et al., 2020)



- MITgcm sea-ice assim (*operational*, NMEFC Beijing)
- CMEMS Baltic-MFC (*operational*, DMI/BSH/SMHI)
- NEMO (U Reading, P. J. van Leeuwen)
- SCHISM/ESMF (VIMS, J. Zhang)
- TerrSysMP-PDAF (hydrology, FZ Juelich, U Bonn)
- TIE-GCM (U Bonn, J. Kutsche)
- VILMA (GFZ Potsdam)
- Parody geodynamo (IPGP Paris, A. Fournier)

Goal: Enable easy and fast setup of a DA system,
and allow for extension to fully featured system

Assumption: Users know their model

→ let users implement DA system in model context

For users, model is not just a forward operator

→ let users extend their model for data assimilation

Keep code simple for the user side:

- Define subroutine interfaces to DA code based on arrays
(also simplifies interaction with languages like C/C++/Python)
- No object-oriented programming
(most models don't use it; most model developers don't know it;
many objects we would only have for observations – see later)
- Users directly implement case-specific routines
(no indirect description (XML, YAML, ...) of e.g. observation layout)

operational centers
might have other
priorities – but the
concept is still correct

1. Focus on ensemble methods

2. Efficiency:

- Direct (online/in-memory) coupling of model and data assimilation method (file-based offline coupling also supported)
- complete parallelism in model, DA method, and ensemble integrations

3. Ease of use:

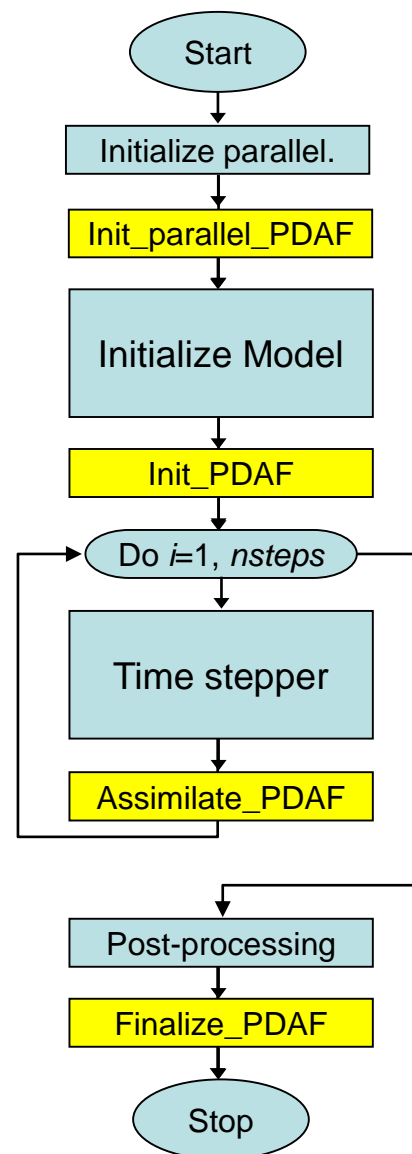
- require just standard compilers and libraries, no containers, etc.
- just add subroutine calls into model code when combining with PDAF
- model time stepper not required to be a subroutine
- model controls the assimilation program
- case-specific routines implemented like model code
- simple switching between different filters and data sets
- Separation of concerns: model, DA methods, observations

Online coupling - Augmenting a Model for Data Assimilation

revised parallelization enables
ensemble forecast

Data assimilation: run model with
additional options

PDAF also supports file-
based (offline) coupling
of separate programs for
model and DA
(but it is less efficient)



Model

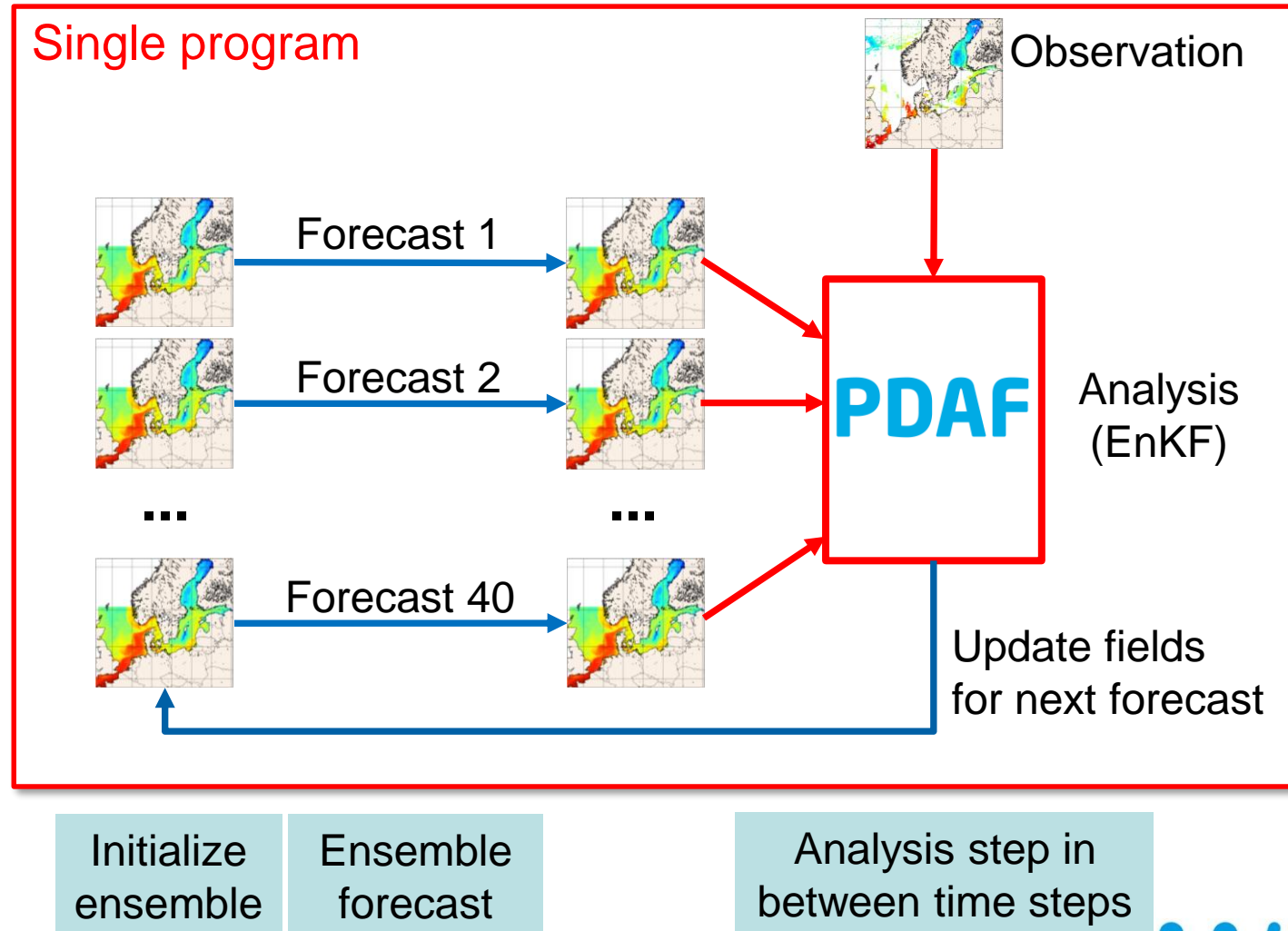
Extension for
data assimilation:
4 subroutine calls

plus:
Possible
model-specific
adaption

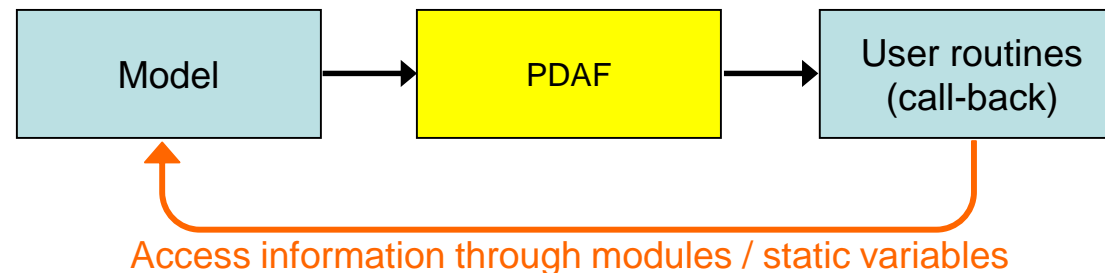
e.g. in NEMO
or ECHAM:
treat leap-frog
time stepping

Couple a model with PDAF

- Modify model to simulate ensemble of model states
- Insert analysis step/solver to be executed at prescribed interval
- Run model as usual, but with more processors and additional options



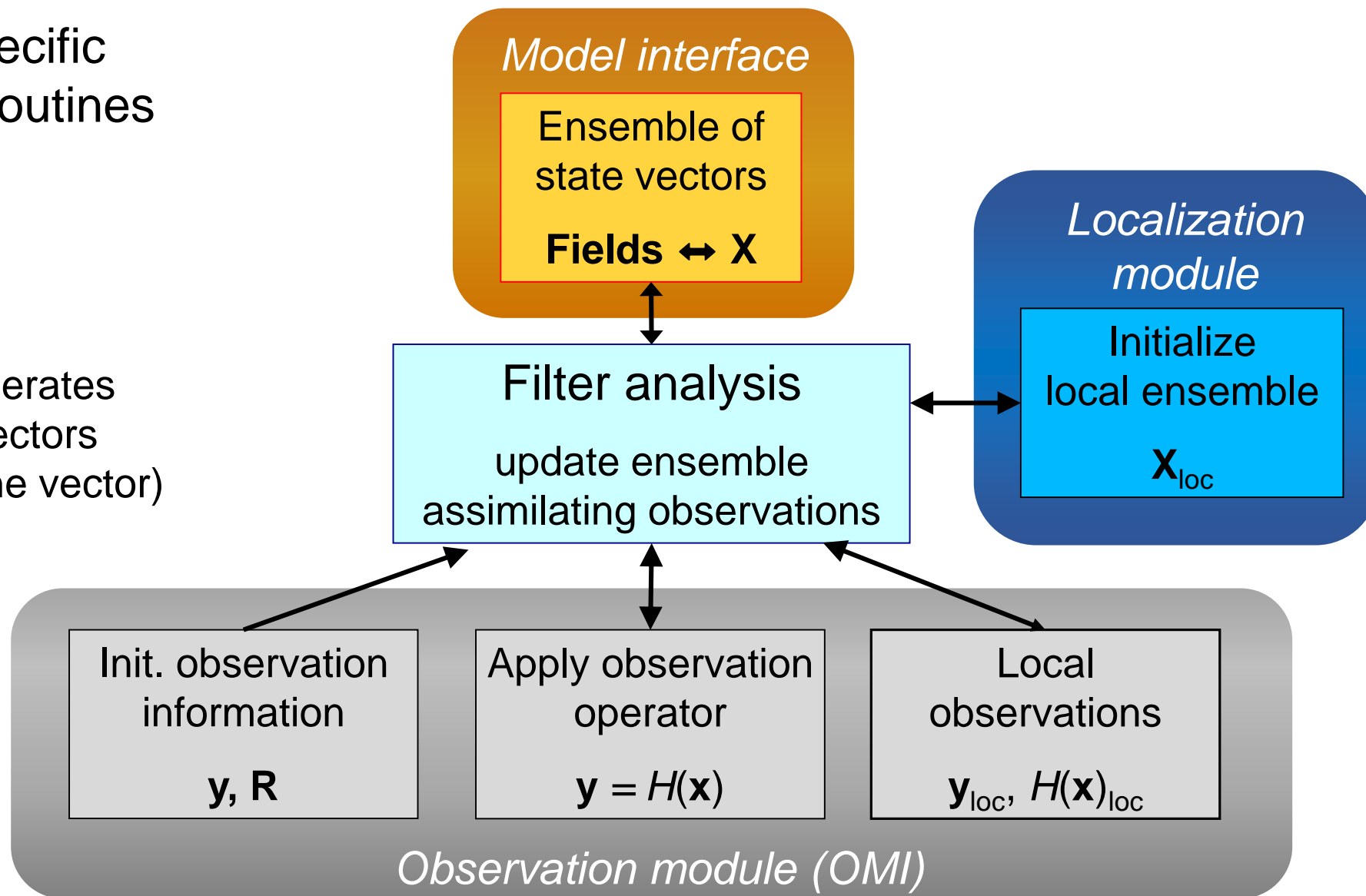
- **Model-sided API:** Defined calls to PDAF routines
- **Case-related API:** User-supplied call-back routines for elementary operations:
 - transfers between model fields and ensemble of state vectors
 - observation-related operations
- User supplied routines can be implemented as routines of the model and can share data with it (low abstraction level)



Implementing the Ensemble Analysis Step (Solver)

case-specific
call-back routines

Analysis operates
on state vectors
(all fields in one vector)



PDAF originated from comparison studies of different filters

Filters and smoothers - *global and localized versions*

- EnKF (Evensen, 1994 + perturbed obs.)
- (L)ETKF (Bishop et al., 2001)
- ESTKF (Nerger et al., 2012)
- NETF (Toedter & Ahrens, 2015)
- Particle filter
- *EnOI mode*

Model bindings

- MITgcm
- AWI-CM / FESOM

Toy models

- Lorenz-96 / Lorenz-63

Community provided:

SCHISM/ESMF

TerrSysMP-PDAF

Upcoming:

- Ensemble 3D-Var
- Hybrid 3D-Var
- Hybrid NETF/LETKF
(see my poster)

Upcoming:

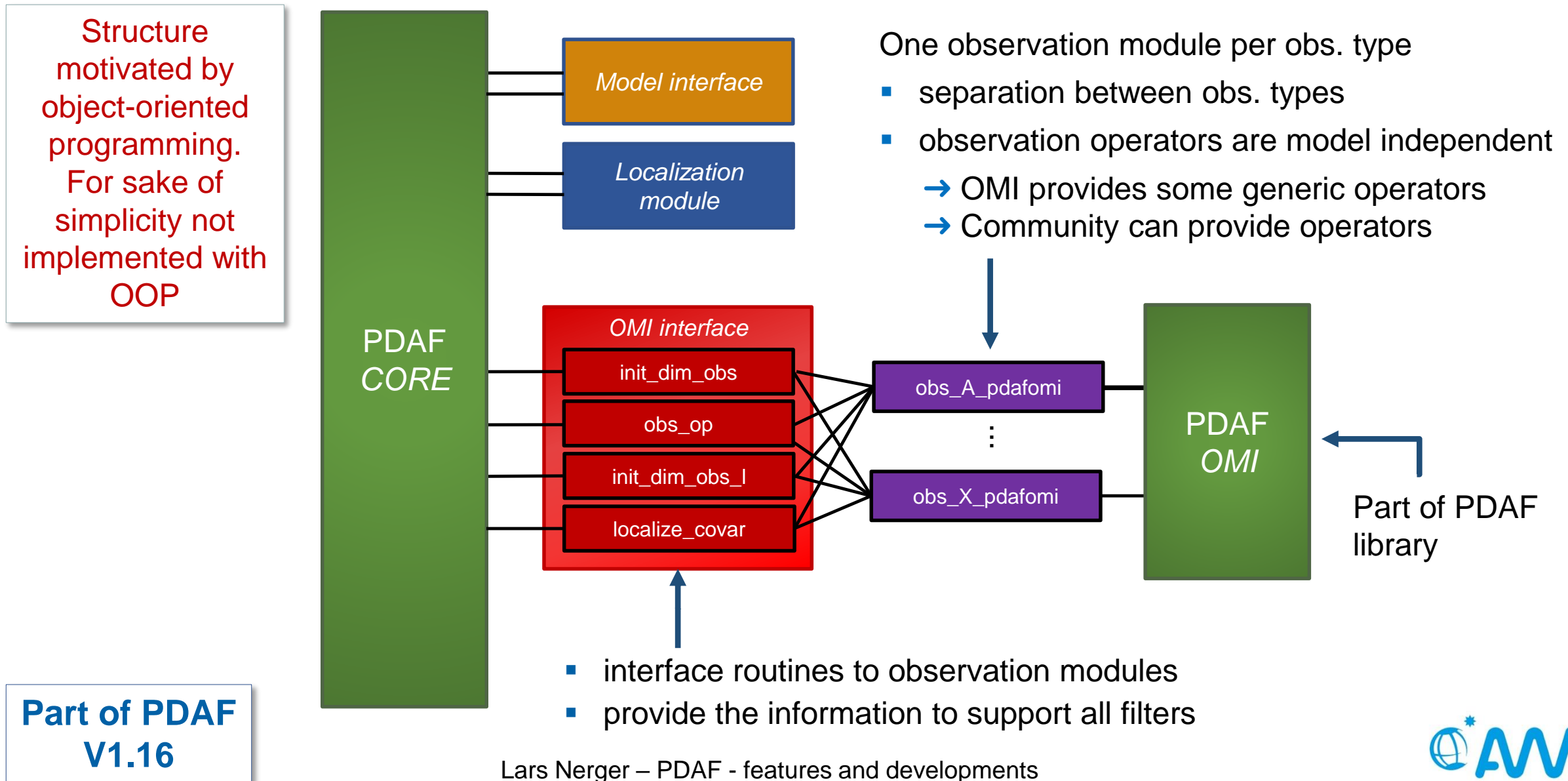
- NEMO 4 (U Reading)
- GOTM/FABM (BB ApS)

Upcoming:

- Lorenz-2005 II/III

Recent and current developments

OMI: Code structure (Observation Module Infrastructure)



See Poster 18 by
Qi Tang et al.

Strongly coupled DA:

Assimilate observation of component A into component B

PDAF supports strongly coupled DA:

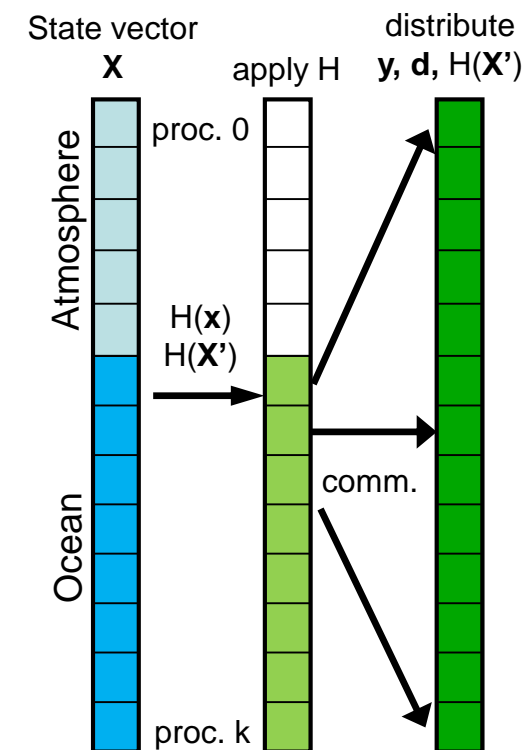
achieved by adapting MPI communicator for the filter processes

- joint state vector decomposed over the processes
- Provide observation operator that only performs MPI communication

need innovation $\mathbf{d} = \mathbf{H}(\mathbf{x}) - \mathbf{y}$ and observed ensemble perturbations $\mathbf{H}(\mathbf{X}')$

Observation operator \mathbf{H} links different compartments

1. Compute part of \mathbf{d} and $\mathbf{H}(\mathbf{X}')$ on process 'owning' the observation
2. Communicate \mathbf{d} and $\mathbf{H}(\mathbf{X}')$ to processes for which observation is within localization radius



Observation handling in
strongly coupled DA

Part of PDAF
V1.16

Activity in EU-project SEAMLESS

1D Prototype (in development):

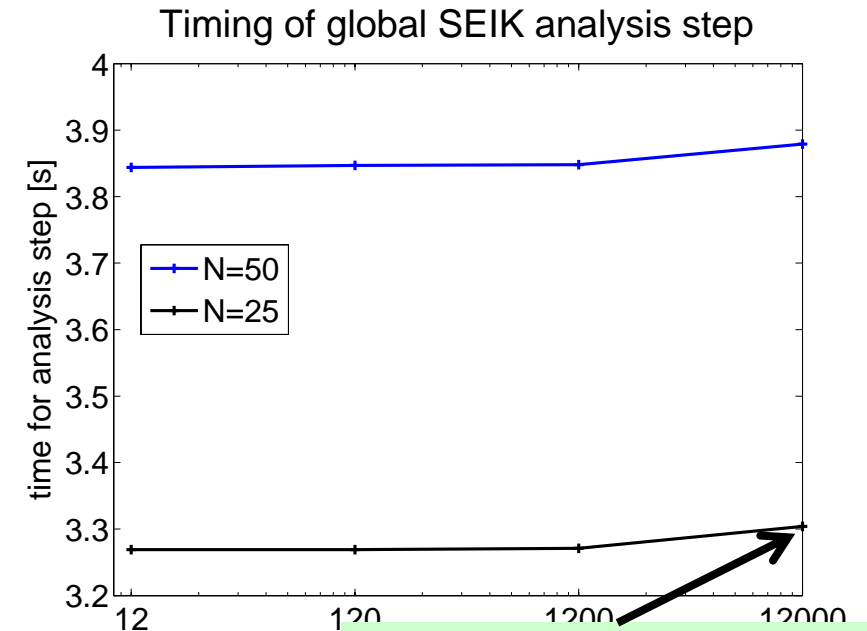
- GOTM/FABM + ecosystem models
- DA functionality provided by PDAF
- **Ensemble/Hybrid 3D-Var**
 - Some partners (PML, OGS) use 3D-Var
 - Integrate in PDAF analogous to EnKFs/PFs
 - Focus on infrastructure with optimizers as core
 - (Future PDAF release)



www.seamlessproject.org

Services based on Ecosystem
data AssiMiLation: Essential
Science and Solutions

- Simulate a “model”
- Choose an ensemble
 - state vector per processor: 10^7
 - observations per processor: $2 \cdot 10^5$
 - Ensemble size: 25
 - 2GB memory per processor
- Apply analysis step for different processor numbers
 - 12 – 120 – 1200 – 12000
- Very small increase in analysis time (~1%)
(Ideal would be constant time)
- Didn't try to run a real ensemble of largest state size (no model yet)
- Latest test: analysis step using 57600 processor cores; state dimension $8.6e11$



State dimension:
 $1.2e11$

Observation
dimension: $2.4e9$

- a program library for ensemble modeling and data assimilation
- provides support for ensemble forecasts, DA diagnostics, and fully-implemented filter and smoother algorithms
- makes good use of supercomputers
- separation of concerns: model, DA methods, observations
- easy to couple to models and to code case-specific routines
- easy to add new DA methods
- efficient for research and operational use

Open source:
Code, documentation, and tutorial available at
<http://pdaf.awi.de>

PDAF adds DA
functionality to
models

Couple model and
PDAF within days

Get DA capability
in a month

Run DA in known
environment

Access new DA
methods by
updating PDAF

- <http://pdaf.awi.de>
- Nerger, L., Hiller, W. (2013). Software for Ensemble-based Data Assimilation Systems - Implementation Strategies and Scalability. Computers and Geosciences, 55, 110-118. [doi:10.1016/j.cageo.2012.03.026](https://doi.org/10.1016/j.cageo.2012.03.026)
- Nerger, L., Tang, Q., Mu, L. (2020). Efficient ensemble data assimilation for coupled models with the Parallel Data Assimilation Framework: Example of AWI-CM. Geoscientific Model Development, 13, 4305–4321, [doi:10.5194/gmd-13-4305-2020](https://doi.org/10.5194/gmd-13-4305-2020)
- Tang, Q., Mu, L., Sidorenko, D., Goessling, H., Semmler, T., Nerger, L. (2020) Improving the ocean and atmosphere in a coupled ocean-atmosphere model by assimilating satellite sea surface temperature and subsurface profile data. Q. J. Royal Meteorol. Soc., in press [doi:10.1002/qj.3885](https://doi.org/10.1002/qj.3885)
- Mu, L., Nerger, L., Tang, Q., Losa, S. N., Sidorenko, D., Wang, Q., Semmler, T., Zampieri, L., Losch, M., Goessling, H. F. (2020) Towards a data assimilation system for seamless sea ice prediction based on the AWI climate model. Journal of Advances in Modeling Earth Systems, 12, e2019MS001937 [doi:10.1029/2019MS001937](https://doi.org/10.1029/2019MS001937)

- Fortran compiler
- MPI library
- BLAS & LAPACK
- make

- PDAF is at least tested (often used) on various computers:
 - Notebook & Workstation: MacOS, Linux (gfortran)
 - Cray XC30/40 & CS400 (Cray ftn and ifort)
 - NEC SX-8R / SX-ACE / SX-Aurora TSUBASA
 - ATOS Bull Sequana X (ifort)
 - HPE Cray Apollo (ARM)
 - Legacy:
 - SGI Altix & UltraViolet (ifort) / IBM Power (xlf) / IBM Blue Gene/Q