

Project report

ATM S 544

Dominik Stiller

December 13, 2023

1 Introduction

For this project, I reproduced some results of Perkins and Hakim (2021)¹, who reconstructed ten climate fields over the last millennium. Their paleoclimate data assimilation (DA) method combines the ensemble forecast of a coupled ocean–atmosphere linear inverse model (LIM) with proxy observations by means of an ensemble square-root filter (EnSRF). Since I will continue this line of research, my project is less concerned with originality and more with familiarizing myself with the data, methods, and workflow.

There are some notable differences of my implementation compared to PH21:

- They use training data for the LIM from CMIP5 *past1000* simulations, while I use CMIP6 *past2k* experiment data. Therefore, my training dataset is larger and more recent.
- They use real proxy observations from the Pages2k, while I generate noisy pseudo-observations from a CMIP6 *past1000* simulation.
- They validate against instrumental datasets, while I validate against the simulation from which the pseudo-observations are drawn.

Note that using pseudoproxies and verifying against their "truth" allows me to check the correctness of my code rather than actually reconstructing the past climate.

My LIM training dataset is the *past2k* simulation (Jungclauss et al., 2017), which models 400 coupled ocean–atmosphere fields over 1–1849 CE at monthly resolution. The extended simulation compared to the *past1000* experiments allows better investigation of the medieval period; the forcings and spin-up are identical to the shorter variant. For us, the benefit are 850 additional years to learn ocean–atmosphere dynamics from. I used the *past2k* simulation of the MPI-ESM1-2-LR model, which is the only available one.

2 Forecasting and data assimilation setup

The main components of my reconstruction code are shown in Figure 1. A mapper translates between the high-dimensional physical space for DA and the low-dimensional space for forecasting (Section 2.1). The LIM provides an ensemble forecast based on initial conditions in this reduced space (Section 2.2). Observations are assimilated in the physical space by the EnSRF (Section 2.3). This section describes the technical aspects of my cycling DA setup.

¹The three Perkins & Hakim papers are hereafter referred to as PH17, PH20, and PH21.

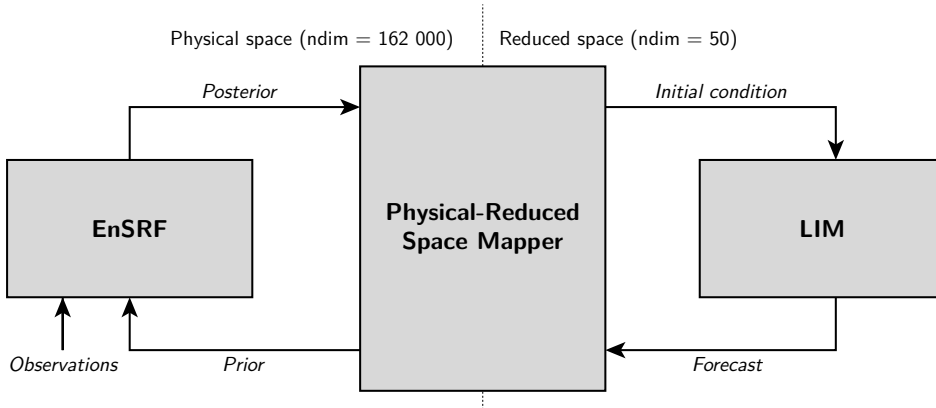


Figure 1: Setup for cycling DA.

2.1 Dimensionality reduction for LIM

The *past2k* simulation output of MPI-ESM1-2-LR contains ~ 400 fields on its native grid of $1.875^\circ \times 1.875^\circ$. Like PH21, I regridded the data to a regular $2.0^\circ \times 2.0^\circ$ latitude–longitude grid using bilinear interpolation and select ten fields: 2 m surface air temperature (tas), precipitation (pr), sea-level pressure (psl), 500 hPa geopotential height (zg500), outgoing top-of-atmosphere longwave radiation (rlut), reflected top-of-atmosphere shortwave radiation (rsut), sea-surface temperature (tos), sea-surface salinity (sos), dynamic ocean surface height (zos), and 0 m to 700 m ocean heat content (OHC700). I then formed annual averages from April to March since many real proxies do not have higher temporal resolution. This also means that the LIM forecast step is 1 year.

The original state vector for each timestep is 162 000-dimensional, even after regridding and field selection, which poses a large computational burden for the LIM. PH21 leveraged joint modes of variability to reduce the state to 50 dimensions while retaining sufficient variance. The dimensionality reduction comprises multiple steps:

1. Mask out nan values. These do not have to be included in the state vector. For example, ocean fields do not have values over land, and I needed to remove some spurious invalid values, particularly above the poles.
2. Remove trend. The LIM assumes stationary statistics (Penland & Sardeshmukh, 1995) and forecasts anomalies; removing the trend yields these anomalies.
3. Fit and project each field individually onto empirical orthogonal functions (EOFs). PH21 found that 400 components retain more than 90% of each fields variance. Each grid cell is latitude-weighted before fitting the EOFs.
4. Standardize each field individually. (The pr field is standardized before EOF projection since its small magnitude makes it prone to numerical noise.)
5. Stack all fields except ohc700, then fit and project them onto a joint EOF. I retain the leading 30 components, slightly more than PH21 since my data are different. This highlights how coupled field variability really is.
6. Truncate the ohc700 field to 20 components. This field is kept separate because its dynamics are much slower than the other fields, particularly atmospheric ones.
7. Stack the 30 joint EOF components and 20 ohc700 components into a single 50-dimensional vector.

To map back into physical space, these steps are applied in reverse order. The dimensionality reduction compresses the ten fields from 26 GB down to just 1 MB (a factor 12 of which is due to the annual averaging)/

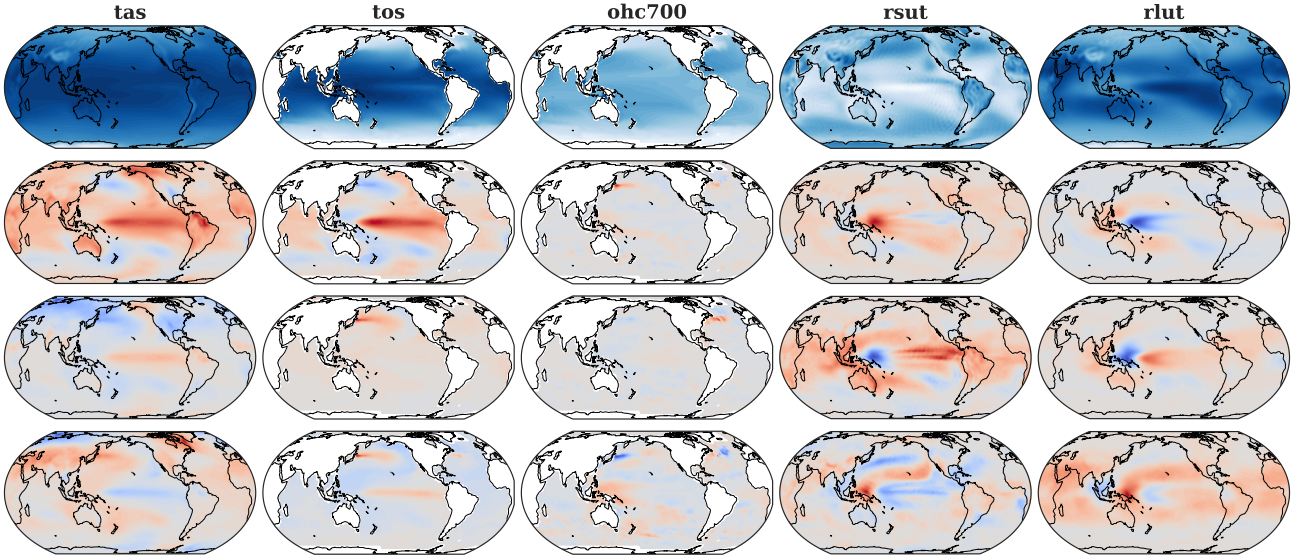


Figure 2: Mean (top) and first three EOFs (below) for a selection of fields.

Figure 2 shows the mean and leading EOFs that my code found. All of them look physical (as far as I can judge) and show well-known patterns such as ENSO and the Indonesian Throughflow. The ohc700 field is clearly distinct with less spatial variability. This supports its separate treatment in the reduction procedure.

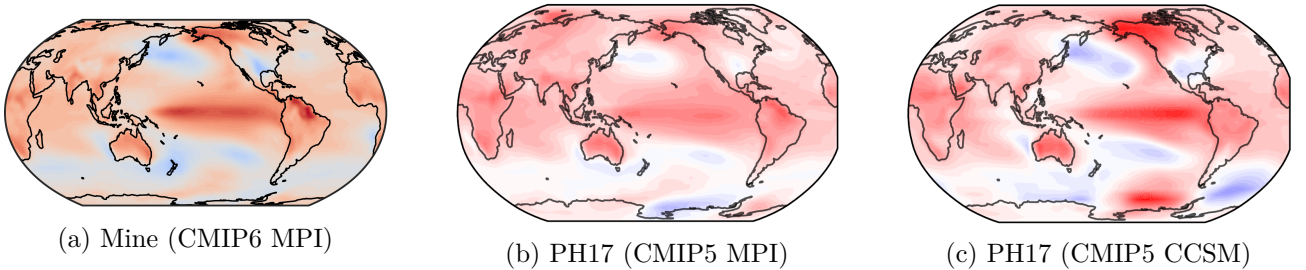


Figure 3: Leading EOFs for TAS from my and Perkins and Hakim (2017) (PH17) reconstructions.

We can compare the leading EOF of tas (first column, second row in Figure 2) to those found by Perkins and Hakim (2017) in Figure 3. Their EOFs are derived from CMIP5 *past1000* outputs of the CCSM and MPI models while mine comes from the CMIP6 version of the MPI model. The patterns are generally similar: a pronounced warm pattern in the tropical Pacific is flanked by colder regions; hot Australian and North American continents. The Southern Ocean and South America regions show more similarity between the MPI models, but the North Pacific low is smaller in the CMIP5 MPI version. Still, the EOFs are similar enough for us to conclude that the choice of training data has little effect on the TAS reconstructions.

2.2 LIM for forecasting

The LIM is an efficient forecast model that makes online DA computationally feasible. The forecast allows information to be propagated through time, particularly by inclusion of ohc700, which acts as

dynamic memory to inform low-frequency variability. In addition, the cheap forecasting enables large ensembles: forecasting my 100 ensemble members takes less than 1 s.

A LIM can represent a linear system of the form

$$\frac{d\mathbf{x}}{dt} = \mathbf{L}\mathbf{x} + \boldsymbol{\xi},$$

where $\mathbf{L} \in \mathbb{R}^{N_x \times N_x}$ is a linear operator and $\boldsymbol{\xi} \sim N(0, \mathbf{Q}/dt)$ is additive Gaussian noise. In my case, the state dimension is $N_x = 50$. The deterministic forecast equation is then (Penland & Sardeshmukh, 1995)

$$\hat{\mathbf{x}}_{k+1} = \exp(\mathbf{L} \Delta t) \mathbf{x}_k = \mathbf{G} \mathbf{x}_k,$$

where $\mathbf{G} \in \mathbb{R}^{N_x \times N_x}$ is the linear forecast operator. For probabilistic ensemble forecasts, we can use the two-step stochastic integration scheme proposed by Penland and Matrosova, 1994:

$$\begin{aligned} \mathbf{x}'(t + \delta t) &= \mathbf{x}'(t) + \mathbf{L}\mathbf{x}'(t) \delta t + \hat{\mathbf{Q}}\sqrt{\boldsymbol{\Lambda}\delta t}\mathcal{R}, \\ \mathbf{x}(t + \delta t/2) &= [\mathbf{x}'(t) + \mathbf{x}'(t + \delta t)] / 2, \end{aligned}$$

where $\hat{\mathbf{Q}}$ and $\boldsymbol{\Lambda}$ come from the eigendecomposition $\mathbf{Q} = \hat{\mathbf{Q}}\boldsymbol{\Lambda}\hat{\mathbf{Q}}^{-1}$ of the noise covariance matrix and \mathcal{R} is a vector of random numbers drawn from a standard normal distribution. The integration timestep δt must be chosen much smaller than the corresponding deterministic timestep Δt . PH21 use $\delta t \approx 6$ h, which requires 1440 integration steps over a period of $\Delta t = 1$ year.

The system dynamics \mathbf{L} and \mathbf{G} as well as the noise covariance \mathbf{Q} can be determined from data. The procedure is based on the zero-lag and τ -lag covariance matrices:

$$\mathbf{C}(0) = \langle \mathbf{x}(t)\mathbf{x}^T(t) \rangle \quad \text{and} \quad \mathbf{C}(\tau) = \langle \mathbf{x}(t+\tau)\mathbf{x}^T(t) \rangle,$$

where $\langle \cdot \rangle$ denotes the time average. In my case, the average is over 1848 timesteps (one year is removed by the split-year annual average, another by the 1-year lag). The forecast operator is then recovered as

$$\mathbf{G} = \mathbf{C}(\tau)\mathbf{C}(0)^{-1}.$$

This reveals \mathbf{G} to be the sensitivity between the state at τ -lagged timesteps, normalized by the state covariance. The linear operator \mathbf{L} required for stochastic integration is then found as

$$\mathbf{L} = \frac{\ln \mathbf{G}}{\tau}.$$

The logarithm is preferably evaluated by eigendecomposition of \mathbf{G} . Finally, we can find the noise covariance matrix as

$$\mathbf{Q} = -(\mathbf{L}\mathbf{C}(0) + \mathbf{C}(0)\mathbf{L}).$$

While covariance matrices are positive semi-definite, the \mathbf{Q} found here may have spurious negative eigenvalues as a result of a short training period or significant non-linear dynamics (Penland & Matrosova, 1994). We can remove these negative eigenvalues by eigendecomposition of \mathbf{Q} , but have to rescale the remaining ones to retain the total variance.

2.3 EnSRF for data assimilation

The ensemble forecast from the LIM requires a DA scheme that is compatible with ensemble priors. For simplicity, I use the EnSRF (Whitaker & Hamill, 2002) also employed in PH17, while PH21 use an ensemble transform Kalman filter. The results should be similar if not identical between the two. I did not have time to implement covariance localization and calibrate my filter through inflation or additive noise.

The EnSRF works on the assumption that observations are independent so that Bayes' rule can be applied sequentially. First, the ensemble is converted into perturbations

$$\mathbf{X} = [\mathbf{x}^1 - \hat{\mathbf{x}} \quad \mathbf{x}^2 - \hat{\mathbf{x}} \quad \dots \quad \mathbf{x}^{N_e} - \hat{\mathbf{x}}],$$

where \mathbf{x}^i are the ensemble priors (column vectors) and $\hat{\mathbf{x}}$ is their mean. In my case, the number of ensemble members $N_e = 100$. We then define the ensemble prior of the observations

$$\mathbf{Z} = \mathbf{H}\mathbf{X},$$

which allows us to efficiently compute the variance of the observation forecast as

$$\hat{\sigma}_p^2 = \frac{\mathbf{Z}\mathbf{Z}^\top}{N_e - 1}.$$

The analysis ensemble perturbations are then

$$\mathbf{X}^a = \mathbf{X} \left(\mathbf{I} - \alpha \frac{\mathbf{Z}\mathbf{Z}^\top}{(N_e - 1)(\hat{\sigma}_p^2 + \sigma_o^2)} \right),$$

where σ_o^2 is the observation variance. Notice that the analysis perturbations must be in the span of the prior perturbations. The scalar α such that the analysis covariance is correct is found as

$$\alpha = \left(1 + \sqrt{\frac{\sigma_o^2}{\hat{\sigma}_p^2 + \sigma_o^2}} \right)^{-1}.$$

Next, we compute the analysis ensemble mean as

$$\hat{\mathbf{x}}^a = \hat{\mathbf{x}} + \frac{\mathbf{X}\mathbf{Z}^\top}{(N_e - 1)(\hat{\sigma}_p^2 + \sigma_o^2)} (y - \mathbf{H}\hat{\mathbf{x}}).$$

Finally, we combine mean and perturbations to form the posterior ensemble

$$(\mathbf{x}^i)^a = \hat{\mathbf{x}}^a + \mathbf{X}_i^a,$$

which is the initial condition for the next forecast of the LIM.

2.4 Setup for cycling data assimilation

For this project, I only reconstruct over 850-1050 CE due to time constraints. I fitted the physical-reduced space mapper on the Casper cluster, which produced the training dataset for the LIM. The remaining steps (fitting the LIM, running the reconstruction, analyzing the results) were performed on Greg's `enkf` server.

My ensemble has 100 members, which are initialized from a random draw of 1000 years from the training dataset. My pseudo-observations come from annual averages of the CMIP6 *past1000* output

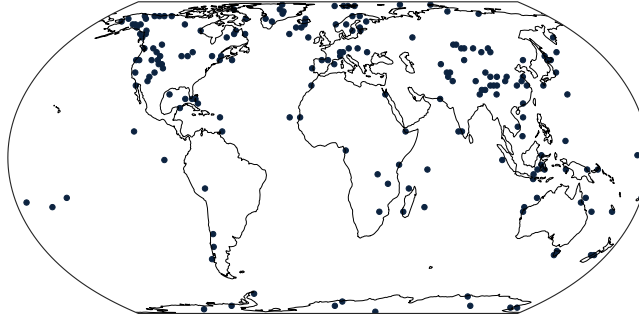


Figure 4: Pseudo-observation locations.

of the MRI-ESM2-0 model. Drawing them from a different model than the training data should lead to a more challenging testbed for my DA system. The observations are surface air temperatures from locations of the Pages2k proxies (PAGES 2k Consortium, 2017), which is the real dataset used by PH21. I draw a random sample of 200 locations (the same ones for every year of the reconstruction), which is about average of the number of proxies assimilated per year in PH21 over the whole last millenium. The 200 locations are shown in Figure 4. The pseudo-observations are perturbed by Gaussian noise with $\sigma_o^2 = 1 \text{ K}^2$, which is also the variance used in the EnSRF.

3 LIM results

also show heatmap for Q matrix

4 DA results

compare to MRI model as validation

For reconstructed GMST plot: show LIM free run, LIM + DA, training data s

use area-weighted GMST

4.1 Sanity check

variance should decrease posterior should be between prior and obs

References

- Jungclaus, J. H., Bard, E., Baroni, M., Braconnot, P., Cao, J., Chini, L. P., Egorova, T., Evans, M., González-Rouco, J. F., Goosse, H., Hurtt, G. C., Joos, F., Kaplan, J. O., Khodri, M., Klein Goldewijk, K., Krivova, N., LeGrande, A. N., Lorenz, S. J., Luterbacher, J., ... Zorita, E. (2017). The PMIP4 contribution to CMIP6 – part 3: The last millennium, scientific objective, and experimental design for the PMIP4 past1000 simulations. *Geoscientific Model Development*, 10(11), 4005–4033. <https://doi.org/10.5194/gmd-10-4005-2017>
- PAGES 2k Consortium. (2017). A global multiproxy database for temperature reconstructions of the common era. *Scientific Data*, 4(1). <https://doi.org/10.1038/sdata.2017.88>

- Penland, C., & Matrosova, L. (1994). A balance condition for stochastic numerical models with application to the el niño-southern oscillation. *Journal of Climate*, 7(9), 1352–1372. [https://doi.org/https://doi.org/10.1175/1520-0442\(1994\)007<1352:ABCFSN>2.0.CO;2](https://doi.org/https://doi.org/10.1175/1520-0442(1994)007<1352:ABCFSN>2.0.CO;2)
- Penland, C., & Sardeshmukh, P. D. (1995). The optimal growth of tropical sea surface temperature anomalies. *Journal of Climate*, 8(8), 1999–2024. [https://doi.org/10.1175/1520-0442\(1995\)008<1999:togots>2.0.co;2](https://doi.org/10.1175/1520-0442(1995)008<1999:togots>2.0.co;2)
- Perkins, W. A., & Hakim, G. J. (2021). Coupled atmosphereocean reconstruction of the last millennium using online data assimilation. *Paleoceanography and Paleoclimatology*, 36(5). <https://doi.org/10.1029/2020pa003959>
- Perkins, W. A., & Hakim, G. J. (2017). Reconstructing paleoclimate fields using online data assimilation with a linear inverse model. *Climate of the Past*, 13(5), 421–436. <https://doi.org/10.5194/cp-13-421-2017>
- Whitaker, J. S., & Hamill, T. M. (2002). Ensemble data assimilation without perturbed observations. *Monthly Weather Review*, 130(7), 1913–1924. [https://doi.org/10.1175/1520-0493\(2002\)130<1913:edawpo>2.0.co;2](https://doi.org/10.1175/1520-0493(2002)130<1913:edawpo>2.0.co;2)