



An Online Paleoclimate Data Assimilation with a Deep Learning-based Network

Haohao Sun¹ Lili Lei¹ Zhengyu Liu² Liang Ning³ Zhe-Min Tan¹

¹School of Atmospheric Sciences, Nanjing University, Nanjing
²Department of Geography, The Ohio State University, Columbus, Ohio
³School of Geography, Nanjing Normal University, Nanjing



Introduction

Paleoclimate data assimilation (PDA) combines information from climate simulations and proxy data to provide an optimal estimate of the past climate. Due to the tremendous computational cost required by lone-time simulations of an Earth system model, PDA has been traditionally performed in an offline scenario, in which priors are selected from existing paleoclimate simulations.

Currently, deep learning-based networks have been promptly developed, which are data-driven and computationally efficient after trained, and thus provide an alternative for online PDA. The deep learning-based network along with an integrated hybrid ensemble Kalman filter is proposed as an online PDA here. This online PDA can efficiently provide climate forecasts with predictive skills and effectively assimilate sparse proxy data, leading to well reconstructed surface air temperature.

Surrogate models

- **Linear regression model:** The linear inverse model (LIM, Penland and Sardeshmukh 1995) is a linear Markov process that represents the evolution of a dynamic system. The model state at t mapped from the model state at $t-1$ is given by

$$\mathbf{x}_t = \mathbf{M}_{LIM} \mathbf{x}_{t-1} + \sigma_{t-1} \quad (1)$$

where the linear operator matrix \mathbf{M}_{LIM} is given by the regression of covariance matrices of the state vectors with time lags of 0 and 1 in the empirical orthogonal function (EOF) space.

- **Deep learning-based network:** A deep learning-based network (NET) is built to map the model state at $t-1$ to that at t , which is a convolutional decoder-encoder (fig. 1). The input and output are three-dimensional state variables, which contain augmented state variables of annual mean and six seasonal averages that are required by the proxy system models (PSMs)

$$\mathbf{x}_t = \mathbf{M}_{NET}(\mathbf{x}_{t-1}) \quad (2)$$

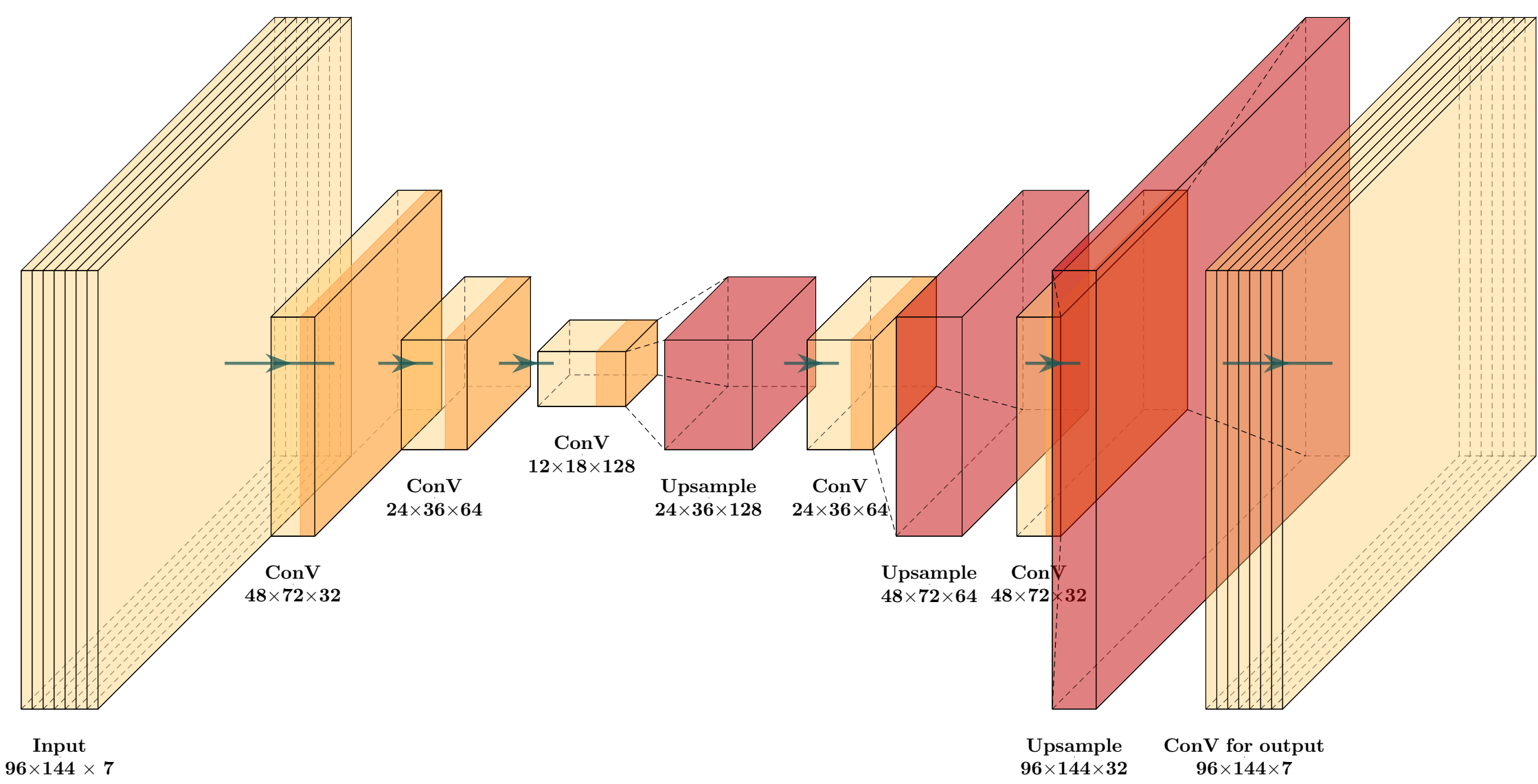


Figure 1. Architecture of the convolutional decoder-encoder network. The encoder comprises one input layer, three convolutional layers (ConV), followed by three upsample-convolutional layers for the decoder. The activation layer is denoted by a darker orange shading.

Data assimilation methods

- **Offline data assimilation:** The offline EnKF [Hakim et al. 2016] is used as a benchmark. Ensemble priors \mathbf{x}^f are randomly sampled from a climatological simulation. The ensemble-square root filter (EnSRF), as a deterministic filter, updates the ensemble mean and ensemble perturbations separately.

$$\bar{\mathbf{x}}^a = \bar{\mathbf{x}}^f + \mathbf{P}^f \mathbf{H}^T [\mathbf{H} \mathbf{P}^f \mathbf{H}^T + \mathbf{R}]^{-1} [\mathbf{y} - H(\bar{\mathbf{x}}^f)] \quad (3)$$

$$\mathbf{x}_i^{ra} = \mathbf{x}_i^f - \mathbf{P}^f \mathbf{H}^T \left[\left(\sqrt{\mathbf{H} \mathbf{P}^f \mathbf{H}^T + \mathbf{R}} \right)^{-1} \right]^T \left[\sqrt{\mathbf{H} \mathbf{P}^f \mathbf{H}^T + \mathbf{R} + \sqrt{\mathbf{R}}} \right]^{-1} (\mathbf{H} \mathbf{x}_i^f) \quad (4)$$

- **Online data assimilation:** Ensemble priors \mathbf{x}_{cyc}^f of online DA are short-term forecasts based on the surrogate models. Ensemble mean and ensemble perturbations are updated separately by the hybrid background error covariances through integrated hybrid EnKF (IHenKF, Lei et al. 2021). The hybrid background error covariances \mathbf{P}_{hyb}^f are a combination of the flow-dependent background error covariances \mathbf{P}_{cyc}^f and the static background error covariances \mathbf{P}_{climo}^f , with weights α and $1-\alpha$ respectively.

$$\bar{\mathbf{x}}_{cyc}^a = \bar{\mathbf{x}}_{cyc}^f + \mathbf{P}_{hyb}^f \mathbf{H}^T [\mathbf{H} \mathbf{P}_{hyb}^f \mathbf{H}^T + \mathbf{R}]^{-1} [\mathbf{y} - H(\bar{\mathbf{x}}_{cyc}^f)] \quad (5)$$

$$\mathbf{x}_{cyc,i}^{ra} = \mathbf{x}_{cyc,i}^f - \mathbf{P}_{hyb}^f \mathbf{H}^T \left[\left(\sqrt{\mathbf{H} \mathbf{P}_{hyb}^f \mathbf{H}^T + \mathbf{R}} \right)^{-1} \right]^T \left[\sqrt{\mathbf{H} \mathbf{P}_{hyb}^f \mathbf{H}^T + \mathbf{R} + \sqrt{\mathbf{R}}} \right]^{-1} (\mathbf{H} \mathbf{x}_{cyc,i}^f) \quad (6)$$

Since the surrogate model tends to lose ensemble variances along with lead times, a blended prior is adopted [Perkins and Hakim 2017]. The selected climatological \mathbf{x}_{climo}^f or analog priors \mathbf{x}_{analog}^f can be added to the ensemble priors \mathbf{x}_{cyc}^f with coefficient β :

$$\mathbf{x}_{cyc,i}^f = \beta \mathbf{M}(\mathbf{x}_{cyc,i}^a) + (1-\beta) \mathbf{x}_i^f \quad (7)$$

Experimental design

- **Proxy and simulation data:** The proxy records from the PAGES 2k Consortium [PAGES2k Consortium et al. 2017] that are temperature sensitive are assimilated, with temporal coverage of the Common Era (0-2000CE). The maximum number of proxies in 1945 CE is 554 (fig. 2).
- **Pseudoproxy and real proxy experiments:** Pseudoproxy data is generated by adding random error from a normal distribution to the true observed value; Real proxies from PAGES 2k are further utilized; Six different proxy networks with 50, 150, 250, 350, 450, and the full proxies are tested.
- **Evaluation metrics:** Root-mean square error (RMSE) and coefficient of efficiency (CE) are computed.

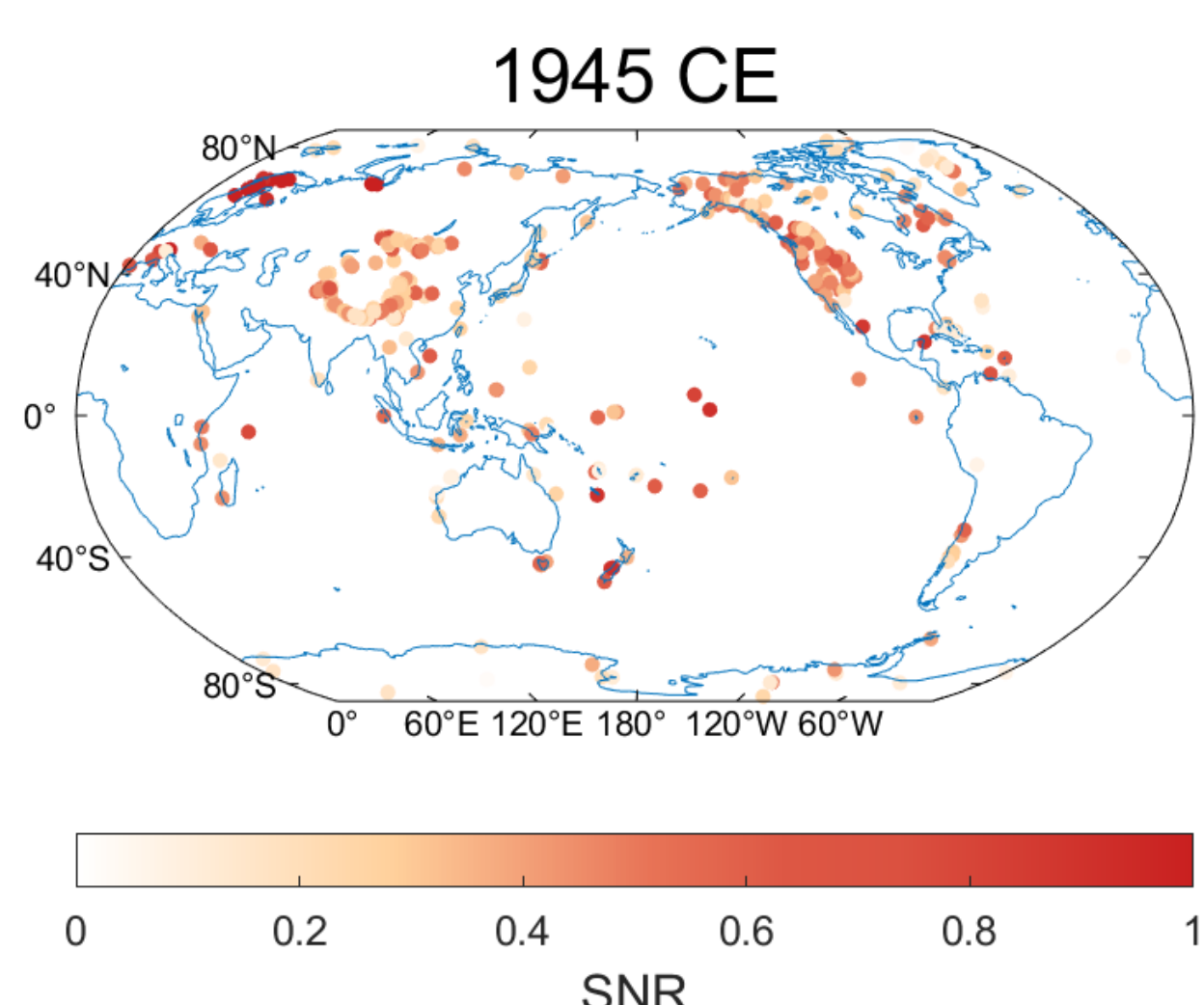


Figure 2. Spatial distributions of proxy records in 1945 CE with colors indicating the signal-to-noise ratio

Predictive skills of surrogate models

- Compared to the offline sampling, LIM has improved predictive skills within five years, and the shorter the lead time the better the predictive skill.
- Compared to the LIM, the NET generally further improves the predictive skills.

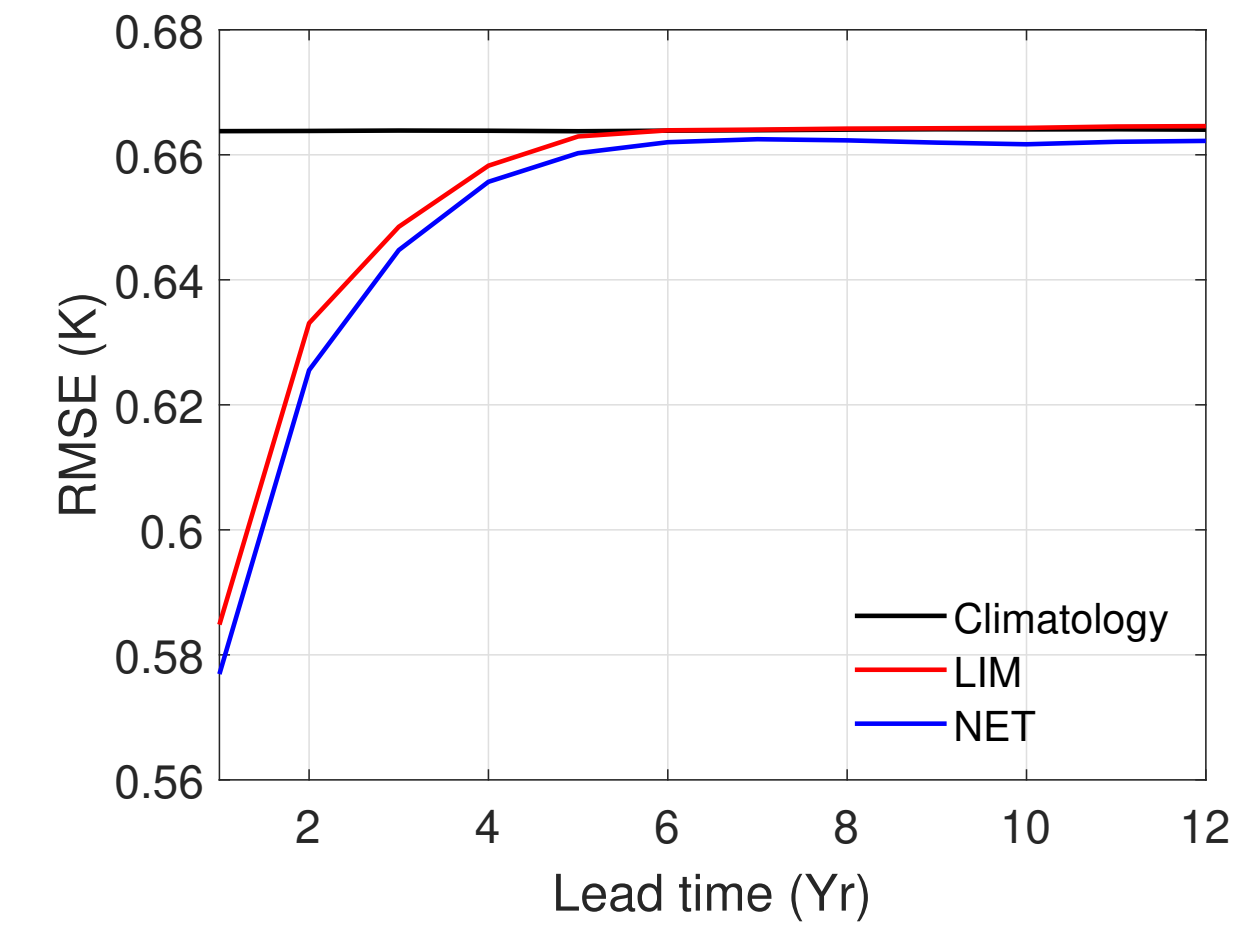


Figure 3. Spatially and temporally averaged RMSEs for climatology, LIM, and NET at various forecast lead times (years)

Results from pseudoproxy experiments

- Exp.LIM and Exp.NET have larger CEs than Exp.OFF, which indicates improved reconstruction skills from the online DA with surrogate models than the offline DA.
- Compared to Exp.LIM, Exp.NET obtains larger CEs of priors and posteriors, the improved predictive skills contributed by the NET than the LIM are beneficial for online DA.
- The differences between Exp.LIM and Exp.NET are enlarged with sparse proxy data.

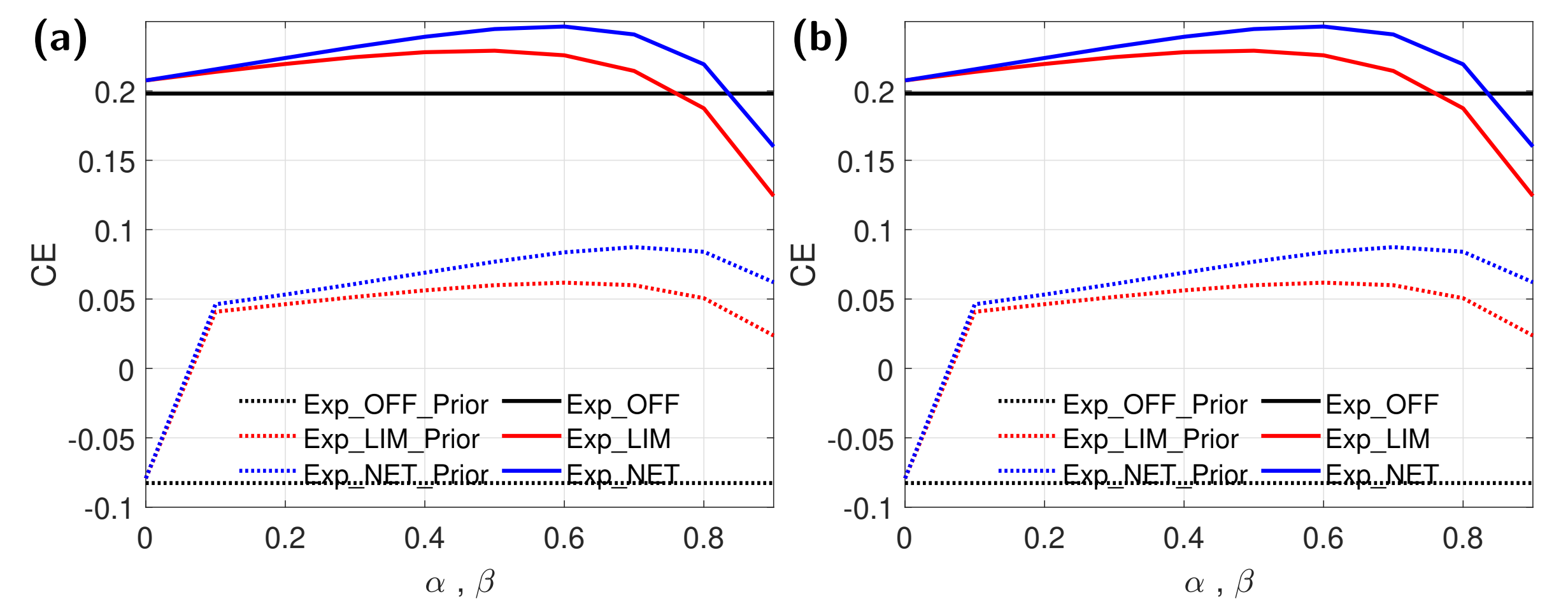


Figure 4. Area-weighted global averaged CEs for Exp.OFF, Exp.LIM, and Exp.NET with varying hybrid weights / blending weights, assimilating a proxy network with (a) 150 pseudoproxies and (b) full pseudoproxies.

Results from real proxy experiments

- When a full set of proxy data is used, offline DA and online DA with either LIM or NET have similar GMT reconstructions, which are mainly forced by the proxy information.
- When a limited set of proxies is used, Exp.OFF has reduced GMT variability than Exp.LIM and Exp.NETANA.
- Compared to Exp.LIM with linear predictive skills obtained by the LIM, Exp.NETANA benefits from retaining nonlinear dynamics provided by the NET and further incorporating analog climatological ensembles, which leads to improved GMT reconstruction even with limited proxies.

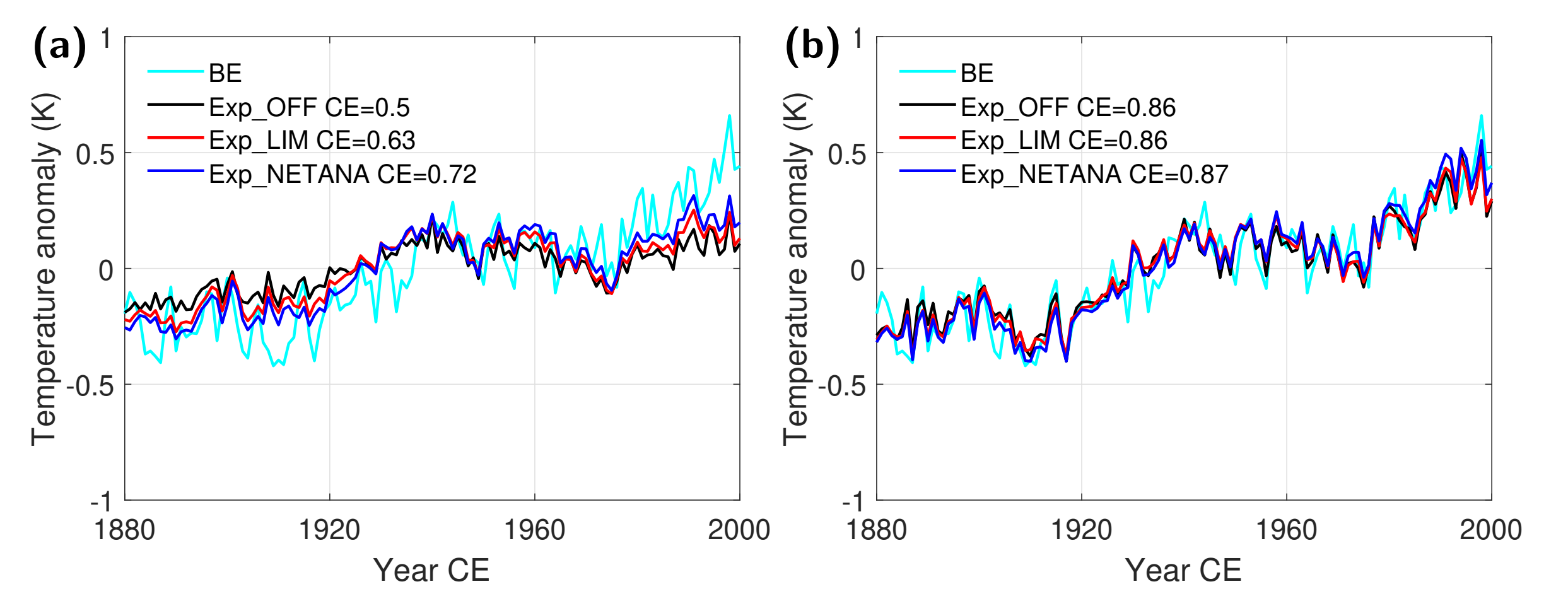


Figure 5. The reconstruction of GMTs over the instrumental period for Exp.OFF, Exp.LIM, and Exp.NETANA given a hybrid weight / blending weight of 0.7, assimilating a proxy network with (a) 150 proxies and (b) full proxies.

- The GMT reconstruction provided by the online PDA has smaller spread but more consistent error statistics than those provided by the online PDA with LIM and offline PDA, during the early period with limited proxy data.

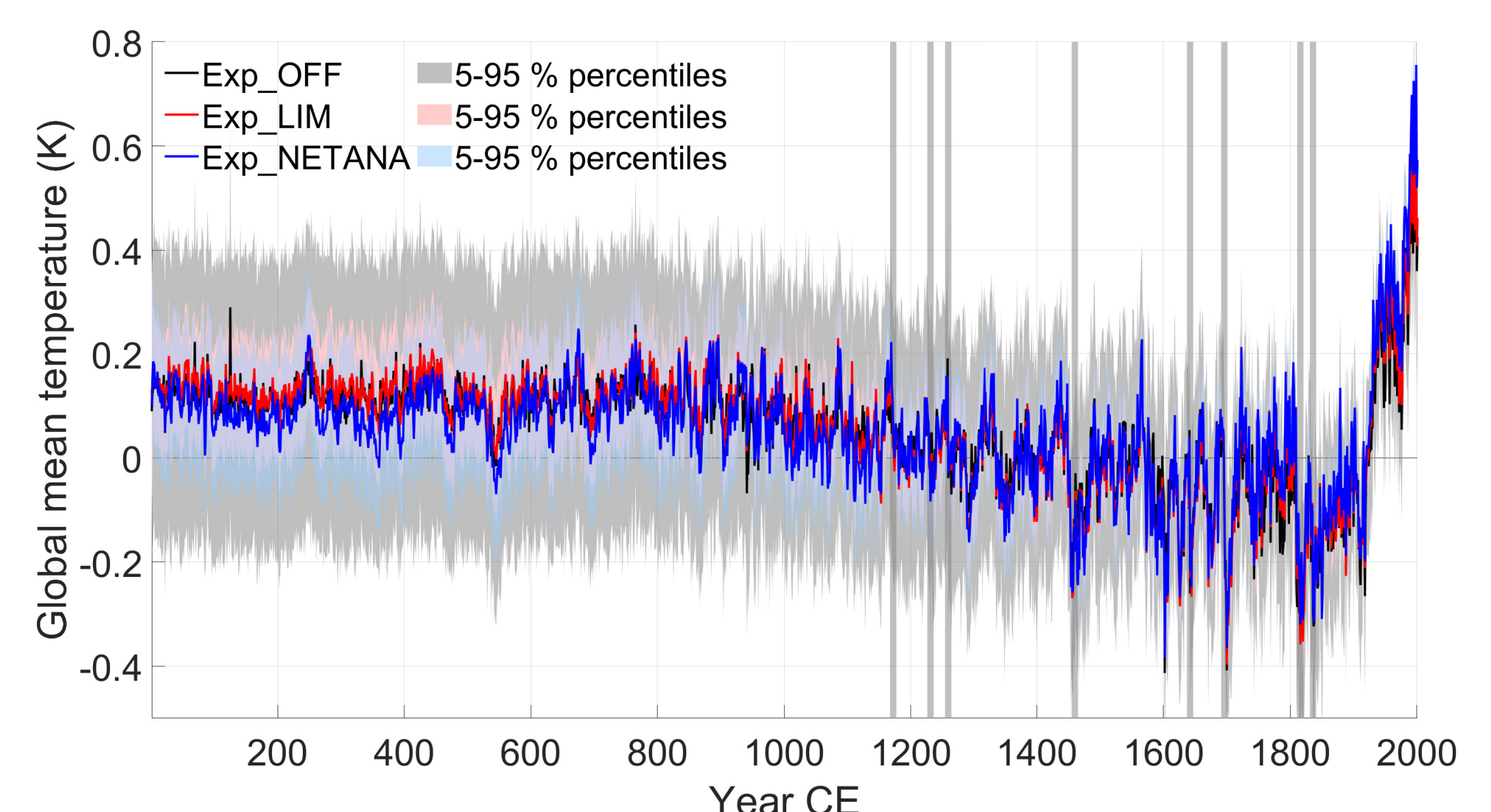


Figure 6. The reconstruction of GMTs over the Common Era for Exp.OFF, Exp.LIM, and Exp.NETANA, given a hybrid weight / blending weight of 0.7, assimilating a proxy network with full proxies. The vertical grey lines denote the volcanic eruption years.