

Interpretability of AI weather models via intermediate decoding

Matthias Beylich, Kirsten Tempest and George Craig

Meteorological Institute Munich, Ludwig Maximilian University of Munich

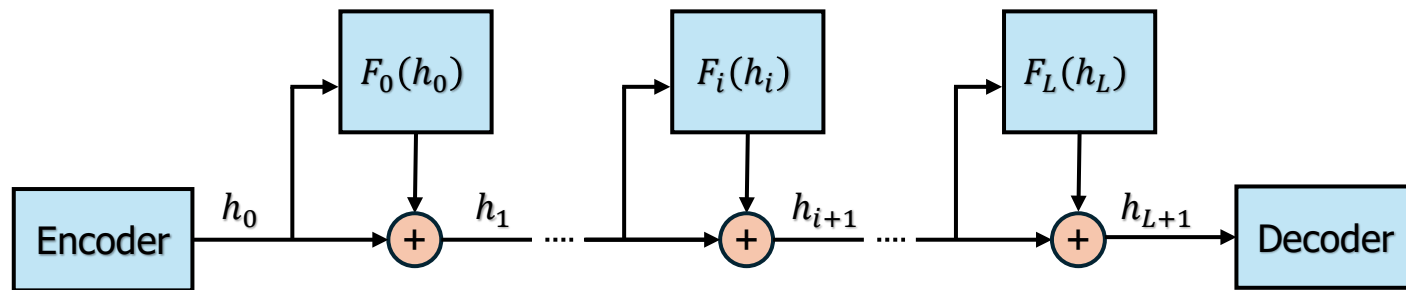
Introduction

- AI weather models are still black boxes
- Apply method of Intermediate Decoding to three different models:
 - Graphcast (Graph Neural Network)
 - Aurora (Transformer model)
 - ArchesWeather (Transformer model)
- How does the forecast evolve with increasing model depth?
 - Model loss
 - Spatial scale progression
 - Physical balances

Methods

- Models employ a residual network architecture:

$$\mathbf{h}_{i+1} = \mathbf{h}_i + F_i(\mathbf{h}_i)$$



- Send intermediate model states to the decoder to create a prediction \hat{X}^{t+1} :

$$\hat{X}^{t+1}(i) = \text{Decoder}(\mathbf{h}_i)$$

- Apply linear transformation (translator) before decoding:

$$\hat{X}^{t+1}(i) = \text{Decoder}(\mathbf{h}_i W_i + \mathbf{b}_i)$$

Methods

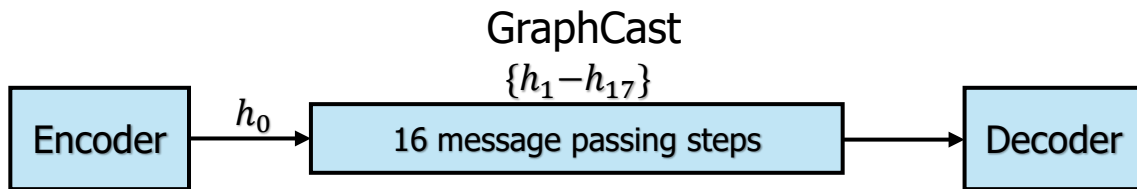
- Models generate a global forecast:

$$\hat{X}^{t+1} = Model(X^t, X^{t-1}), \quad \Delta t_{GC, Au} = 6h, \Delta t_{Ar} = 24h$$

Methods

- Models generate a global forecast:

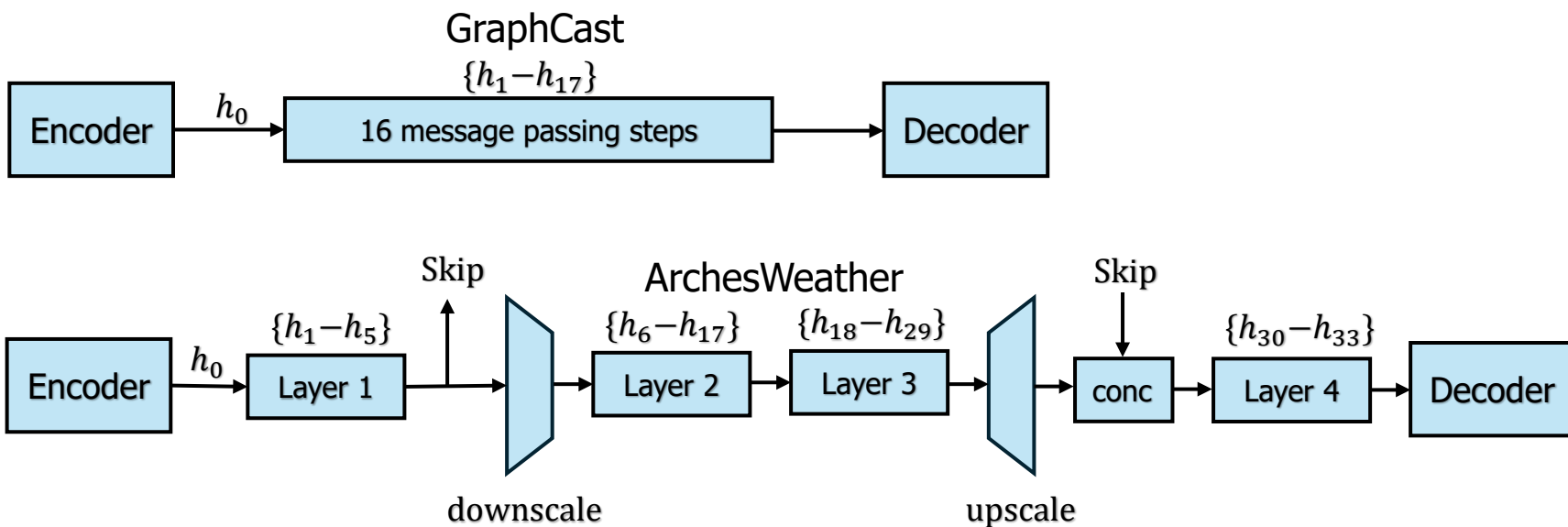
$$\hat{X}^{t+1} = Model(X^t, X^{t-1}), \quad \Delta t_{GC, Au} = 6h, \Delta t_{Ar} = 24h$$



Methods

- Models generate a global forecast:

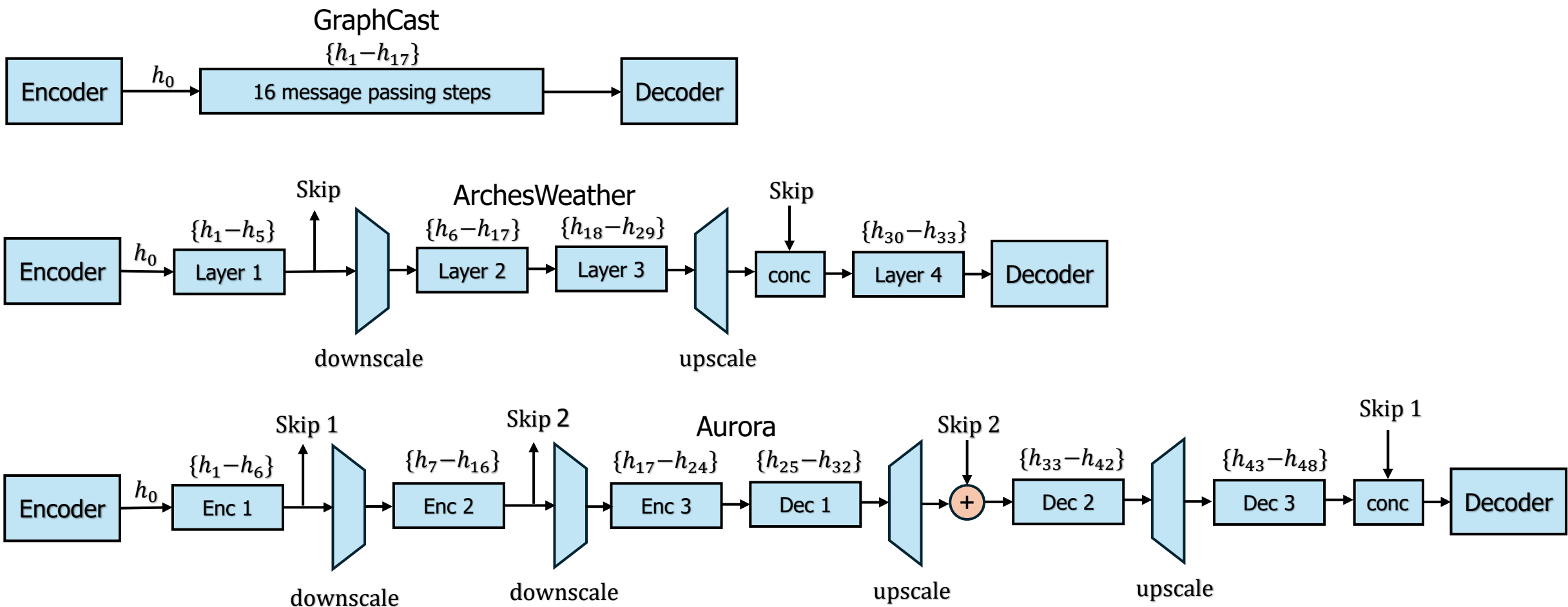
$$\hat{X}^{t+1} = Model(X^t, X^{t-1}), \quad \Delta t_{GC, Au} = 6h, \Delta t_{Ar} = 24h$$



Methods

- Models generate a global forecast:

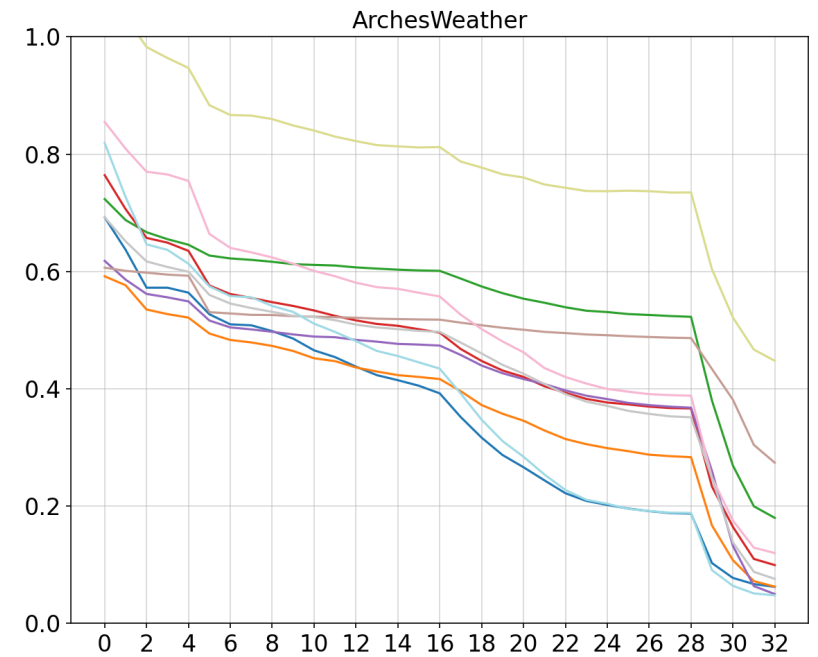
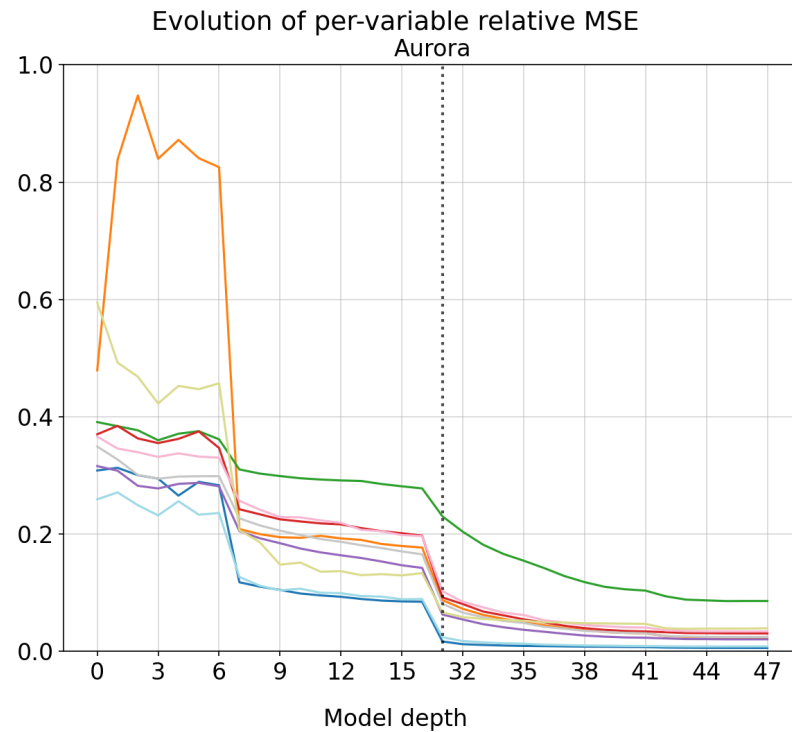
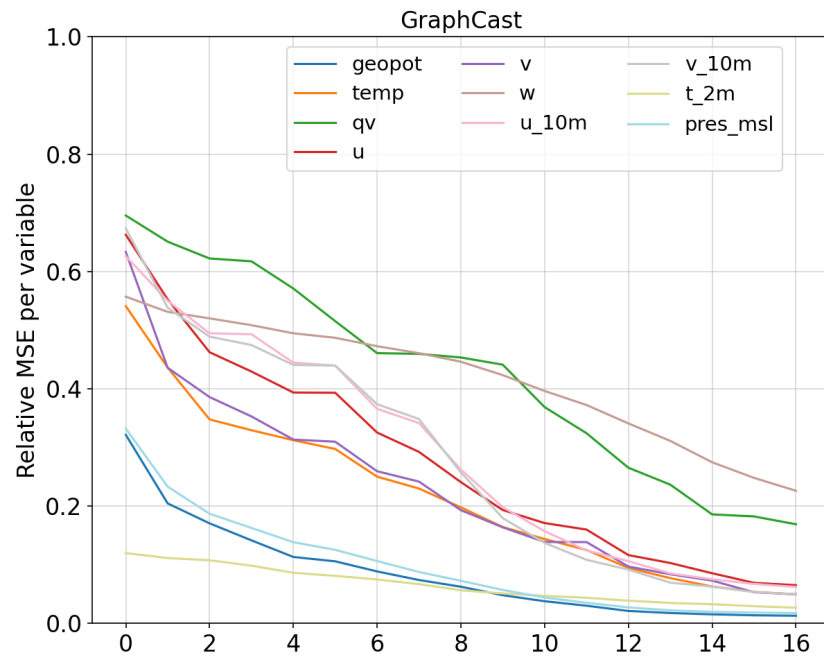
$$\hat{X}^{t+1} = Model(X^t, X^{t-1}), \quad \Delta t_{GC, Au} = 6h, \Delta t_{Ar} = 24h$$



Forecast performance

- Calculate MSE for intermediate decodings:

$$MSE_{rel}(i) = \frac{MSE(\hat{X}_i^{t+1}, X^{t+1})}{MSE(X^t, X^{t+1})}$$

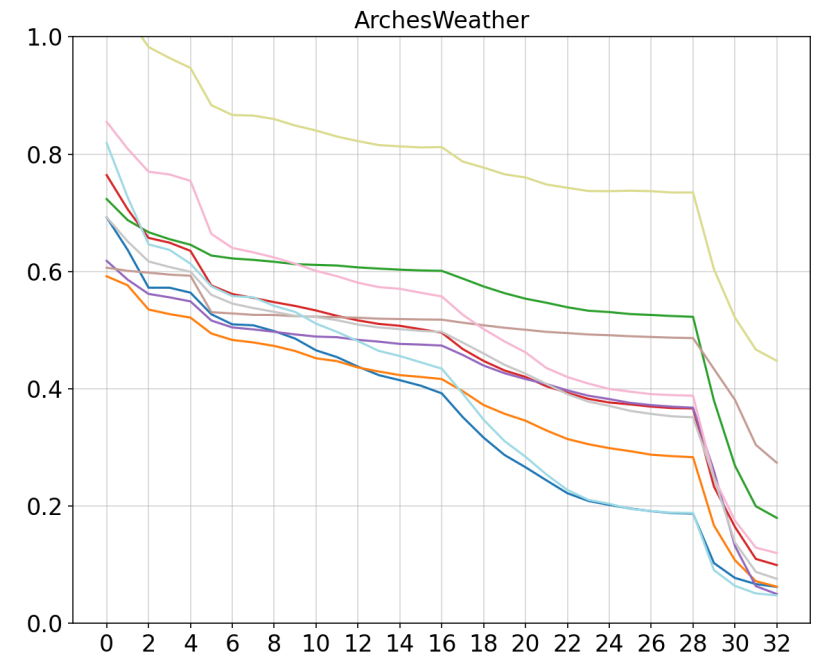
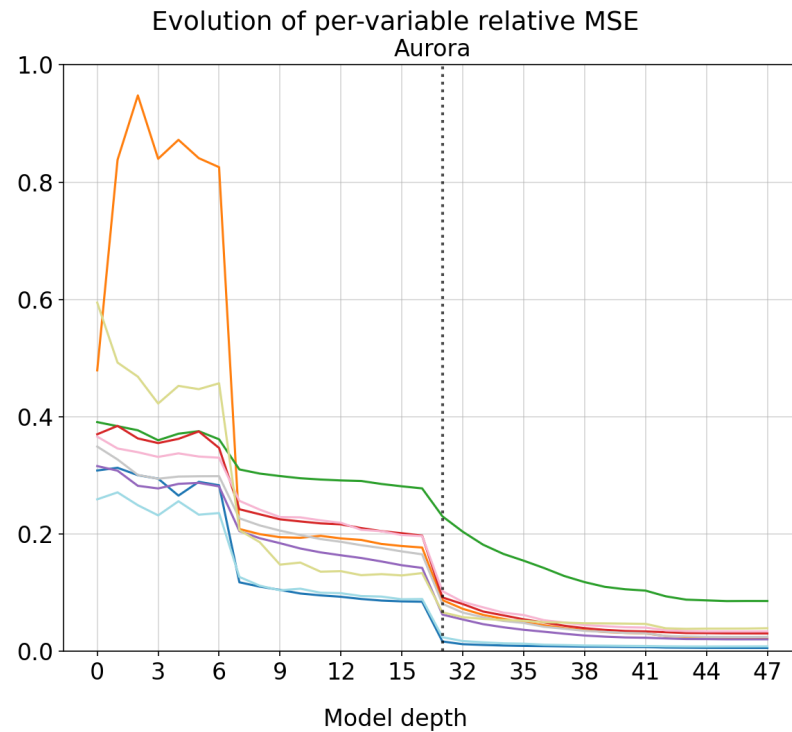
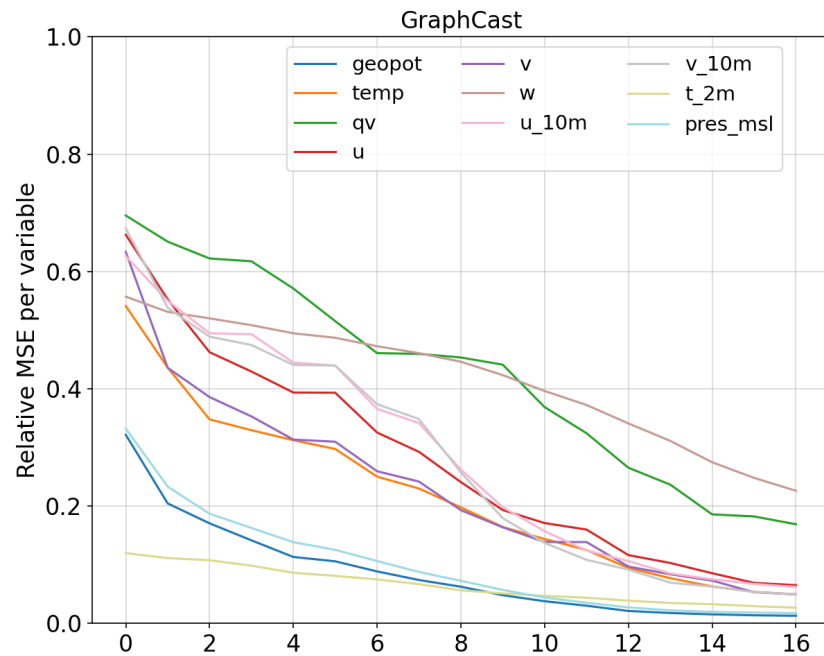


Forecast performance

- Calculate MSE for intermediate decodings:

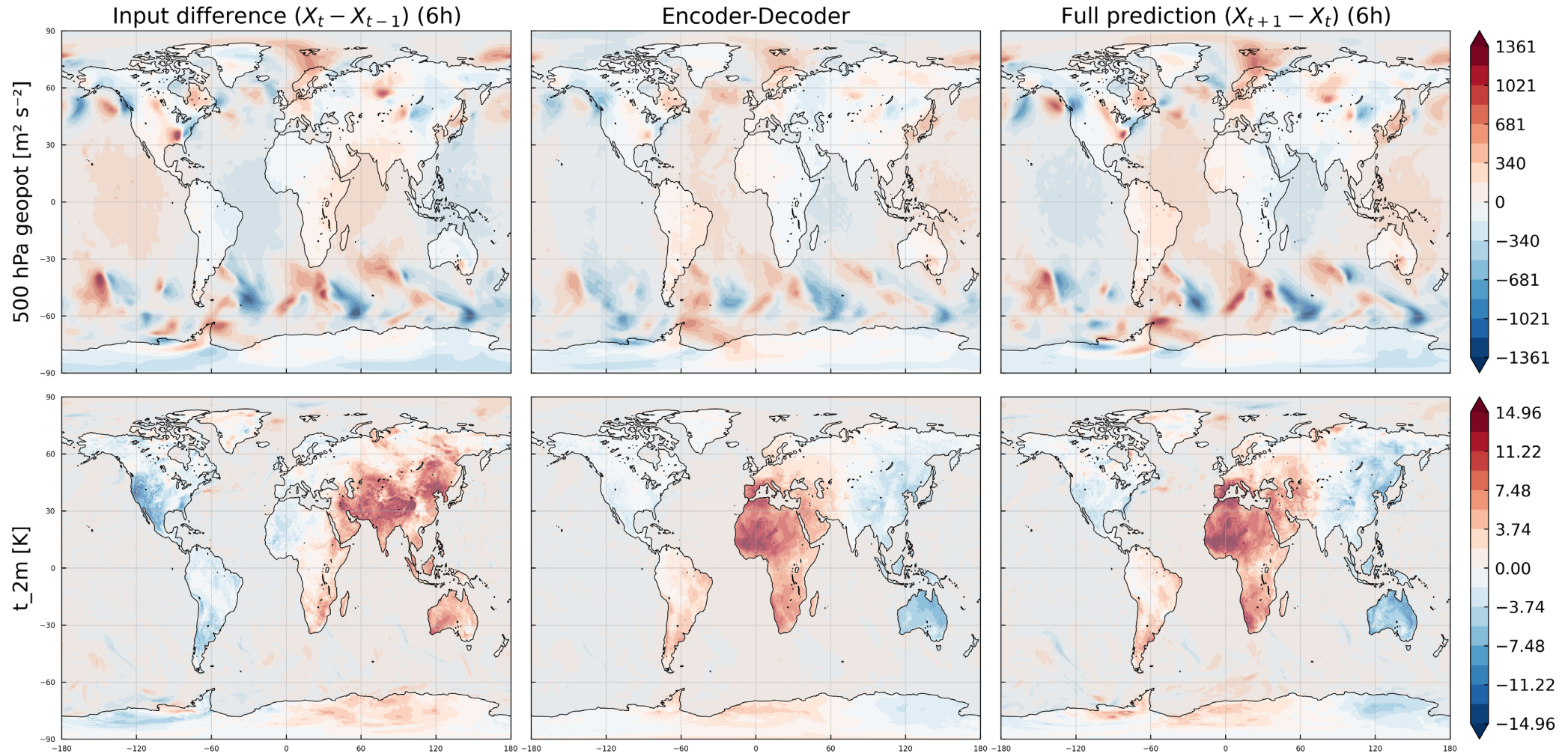
$$MSE_{rel}(i) = \frac{MSE(\hat{X}_i^{t+1}, X^{t+1})}{MSE(X^t, X^{t+1})}$$

- Encoder immediately creates prediction
- Stage structure evident
- Gradual convergence towards final prediction within consistent stages



Encoder-Decoder only forecast

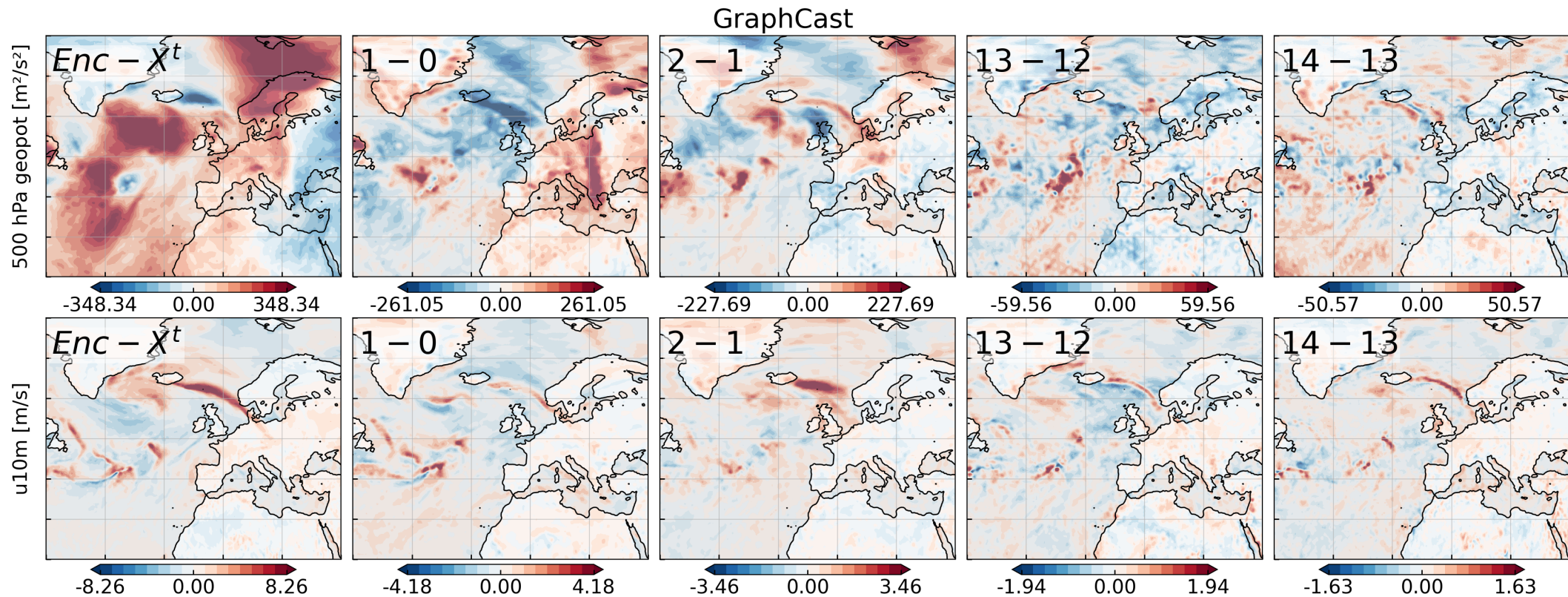
GraphCast



Forecast progression

- Calculate increment between intermediate decodings:

$$\Delta \hat{X}^{t+1}(i) = \hat{X}^{t+1}(i) - \hat{X}^{t+1}(i-1)$$

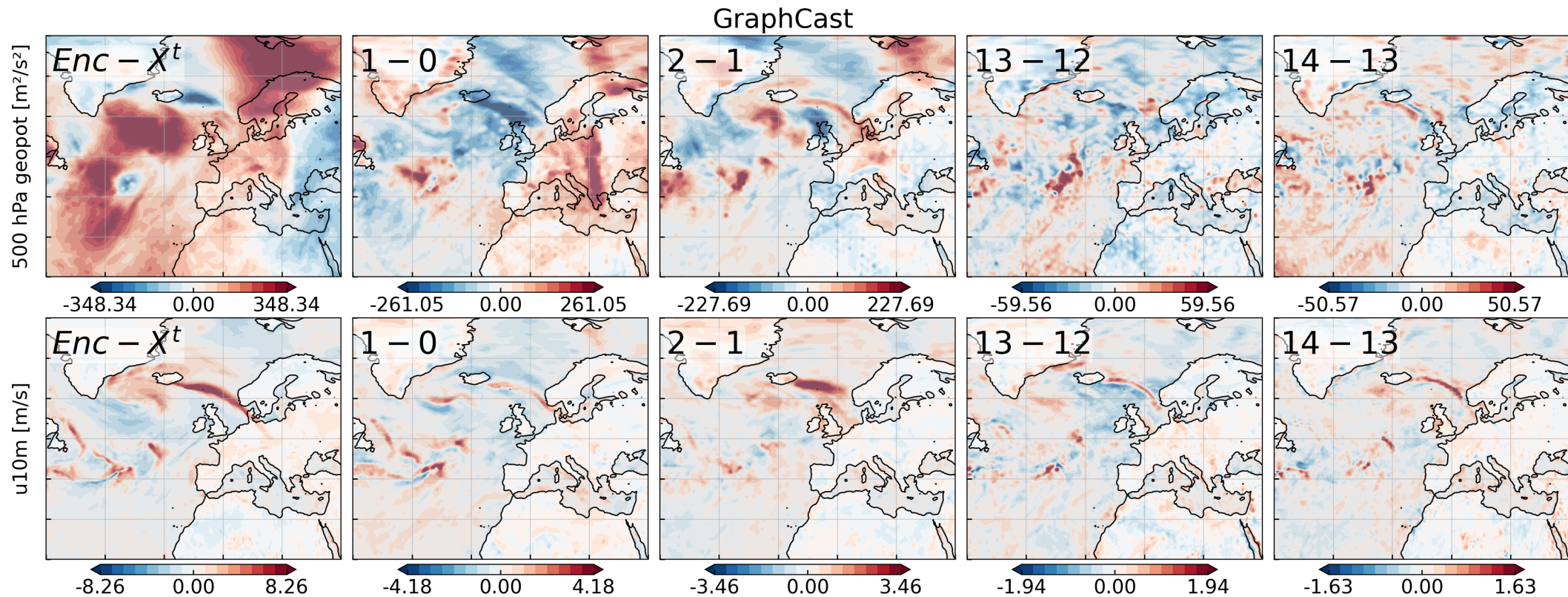


Forecast progression

- Calculate increment between intermediate decodings:

$$\Delta \hat{X}^{t+1}(i) = \hat{X}^{t+1}(i) - \hat{X}^{t+1}(i-1)$$

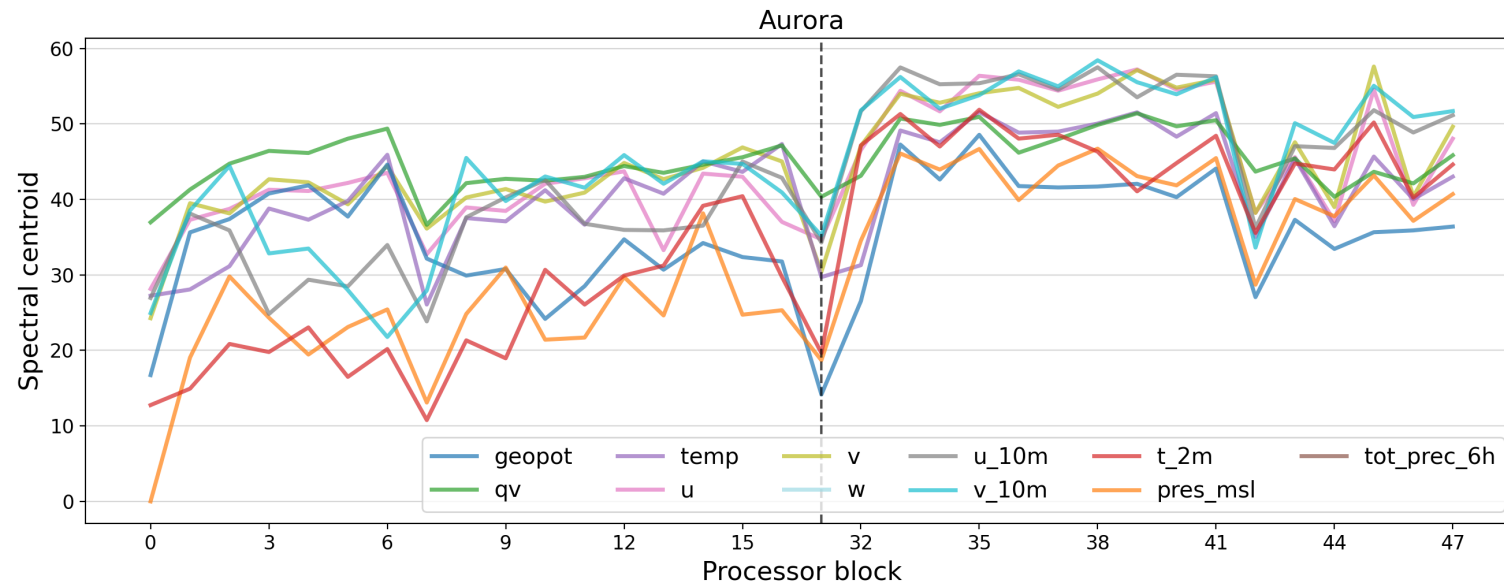
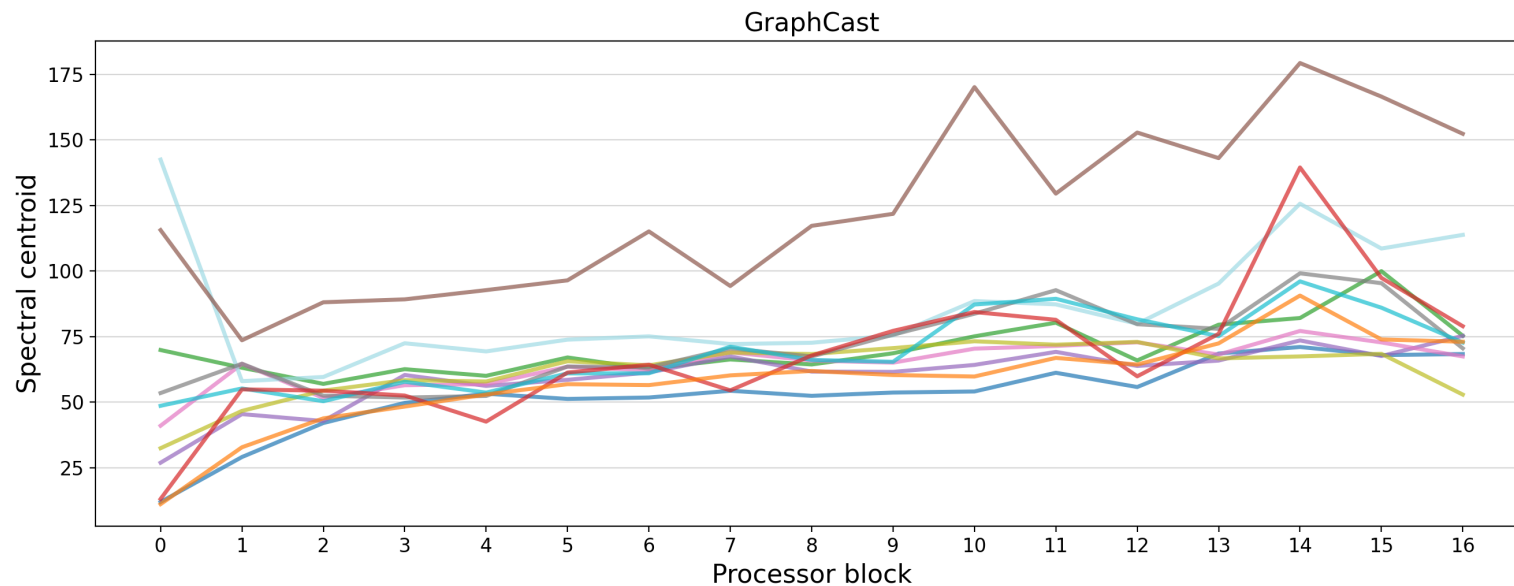
- Progression toward smaller spatial scales and reduced amplitudes with increasing model depth
- Features are iteratively refined



Forecast progression: Scales

- Calculate spectral centroids of increment fields:

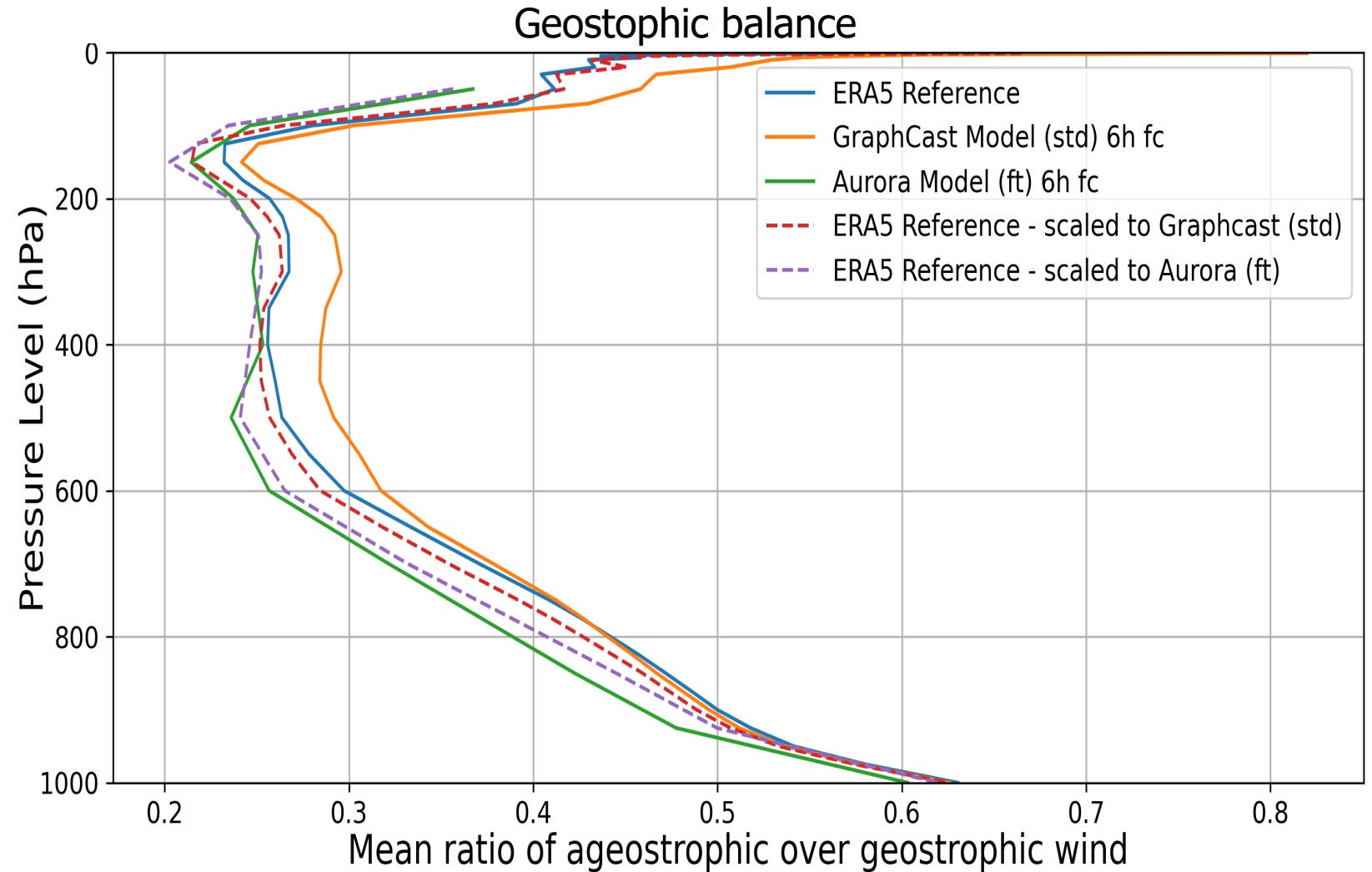
$$l_{centroid} = \frac{\sum l P_l}{\sum P_l}$$



Geostrophic balance

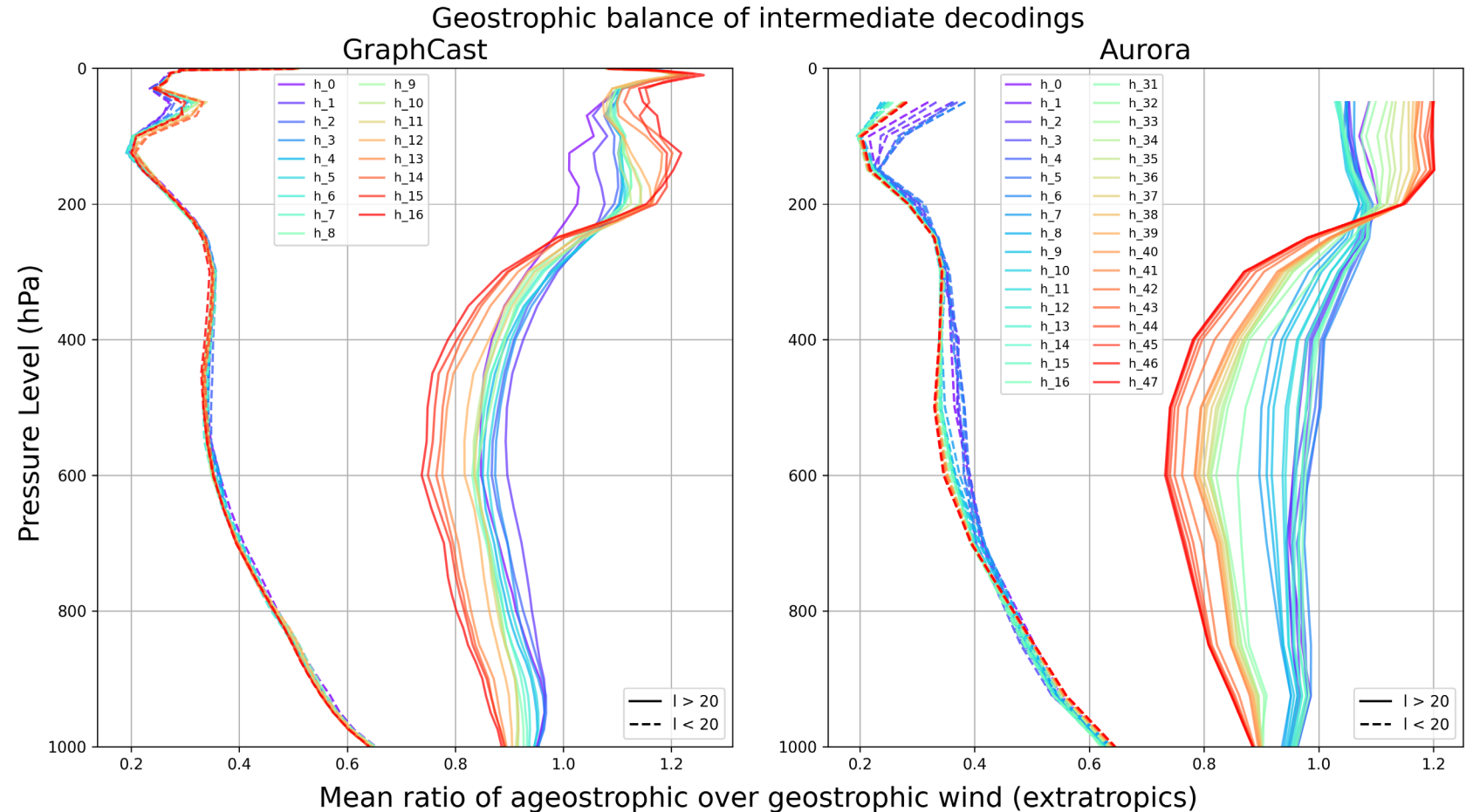
- Calculate geostrophic balance of model forecasts and compare to ERA5

$$r = \frac{\|\mathbf{u}_{\text{ageo}}\|}{\|\mathbf{u}_{\text{geo}}\|}, \quad \text{lon} > 20^\circ$$



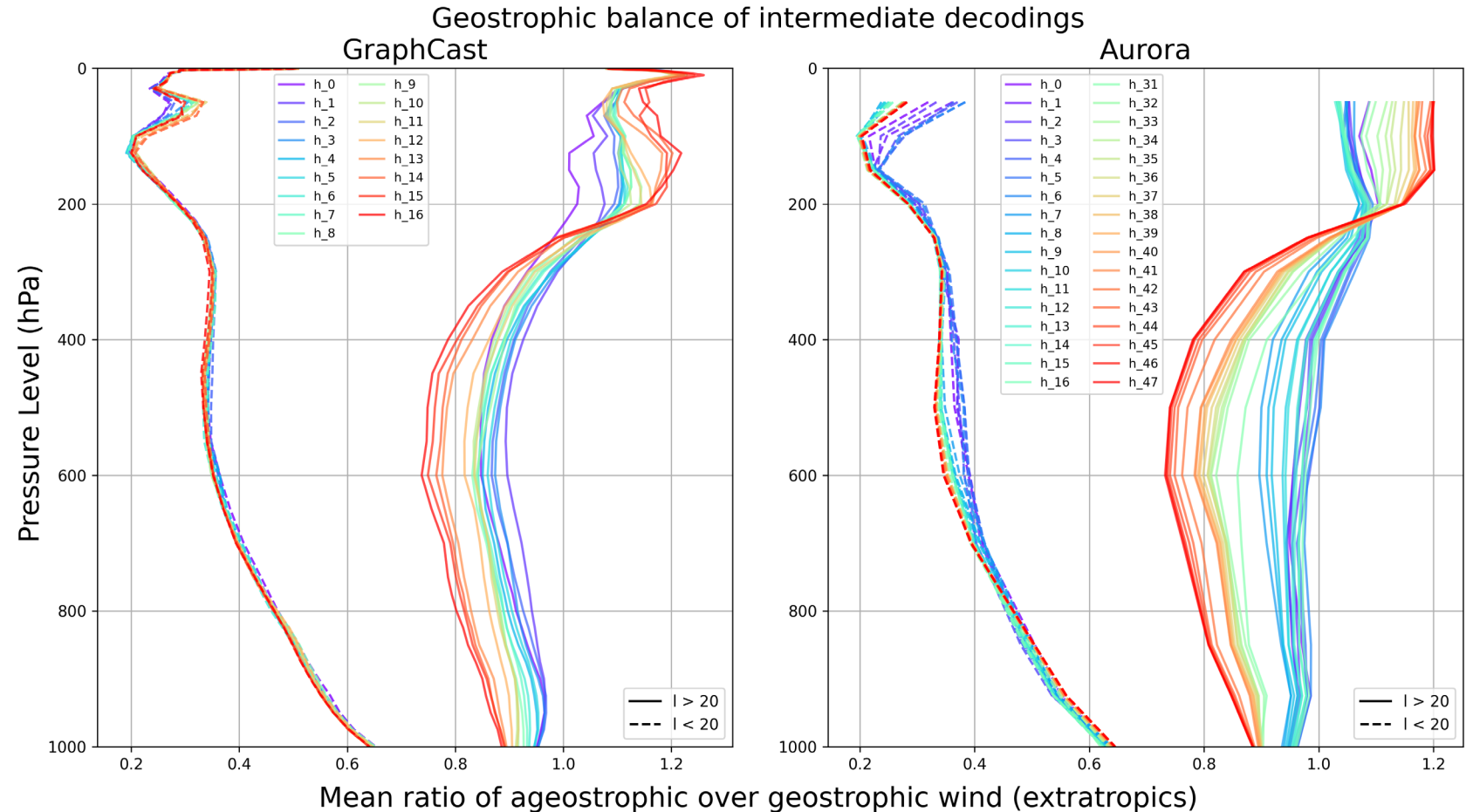
Geostrophic balance: intermediate decodings

- Create intermediate decoding forecasts and filter scales:
 - Planetary/synoptic ($l \leq 20$)
 - Meso ($l > 20$)
- Calculate geostrophic balance on scale separated intermediate decoding forecasts



Geostrophic balance: intermediate decodings

- Create intermediate decoding forecasts and filter scales:
 - Planetary/synoptic ($l \leq 20$)
 - Meso ($l > 20$)
- Calculate geostrophic balance on scale separated intermediate decoding forecasts



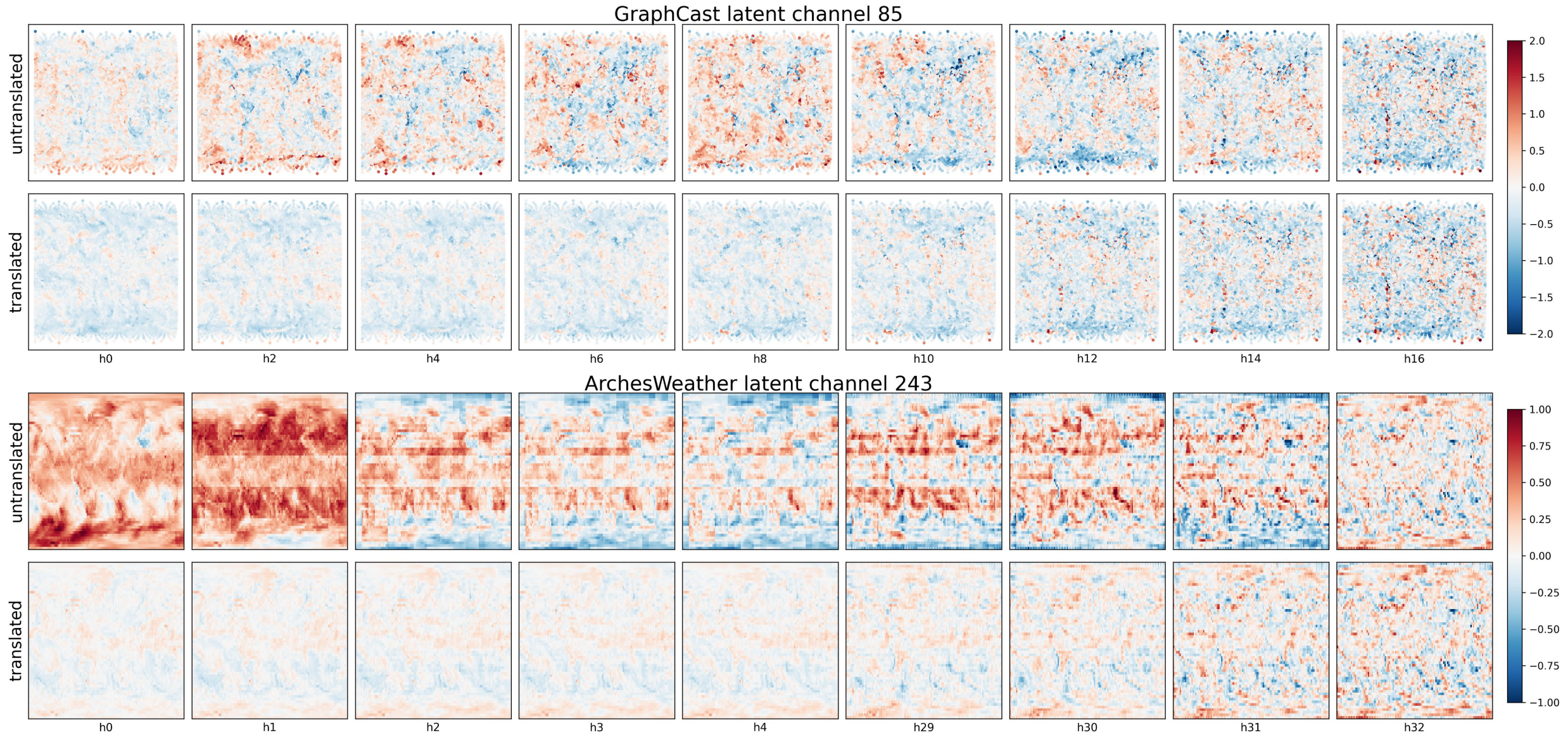
- Geostrophic balance increases with model depth, driven by small scales

Conclusion

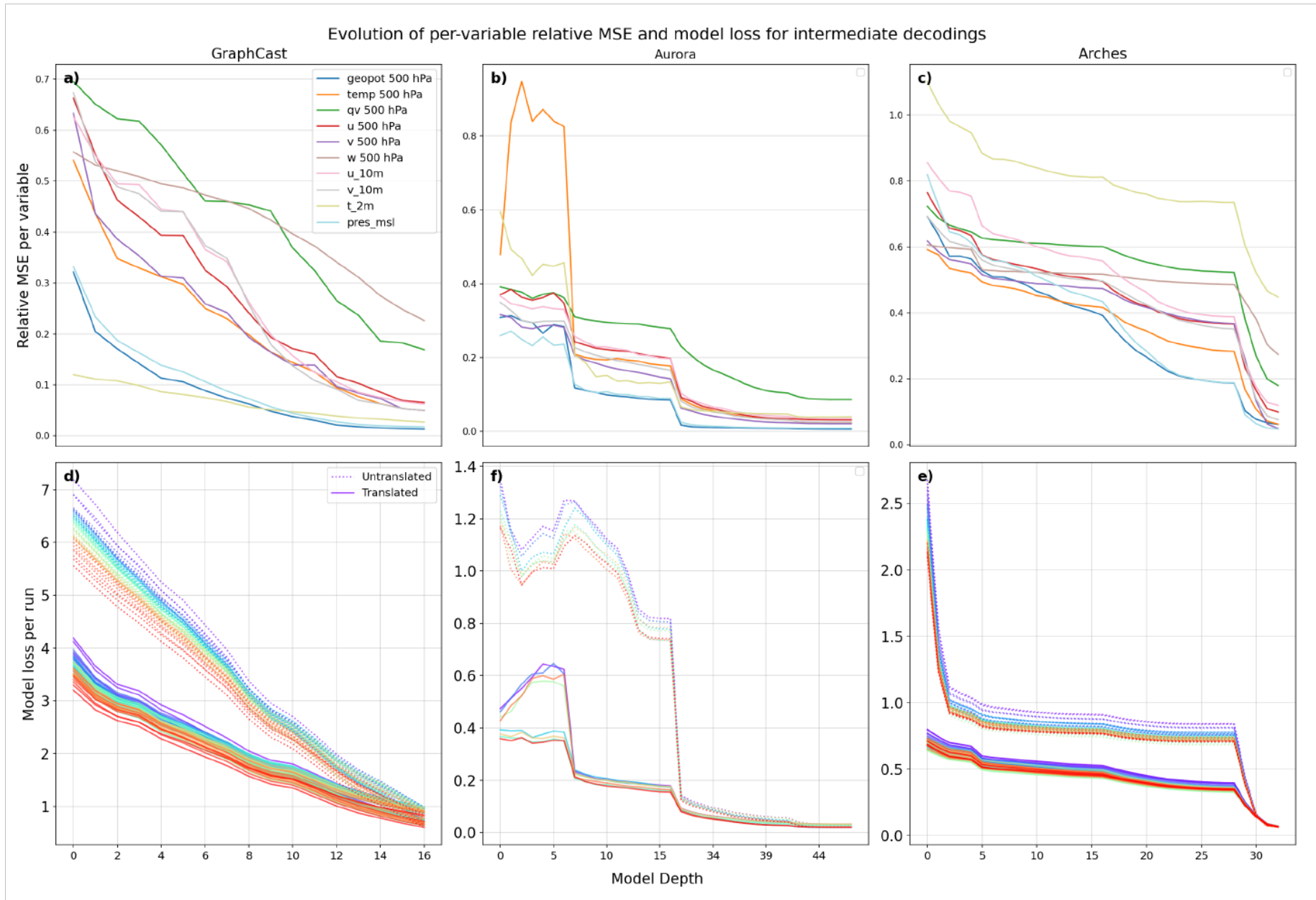
- We apply a method to track how forecasts evolve in physical space:
 - Stage structure relevant for forecast construction
 - Encoder-Decoder create sensible first guess
 - Progression to smaller scales
 - Geostrophic balance with later layers
- However: No direct access to interpretable latent features
- Next steps: Examine latent space directly, PCA, sparse autoencoders, Interactive tool (https://github.com/ktempestuous/latent_space_visualiser_weather_models) (k.tempest@physik.uni-muenchen.de)

Appendix

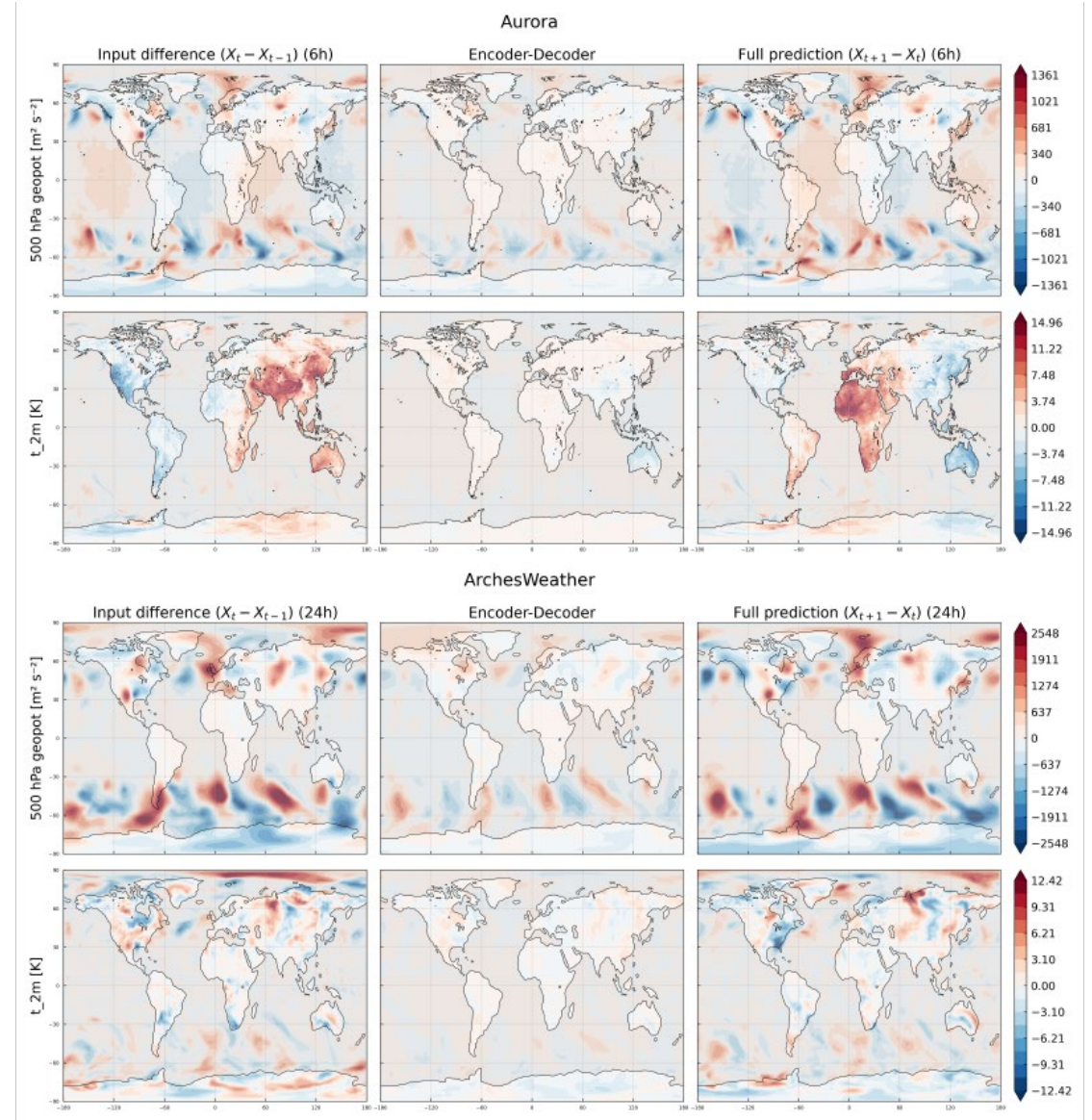
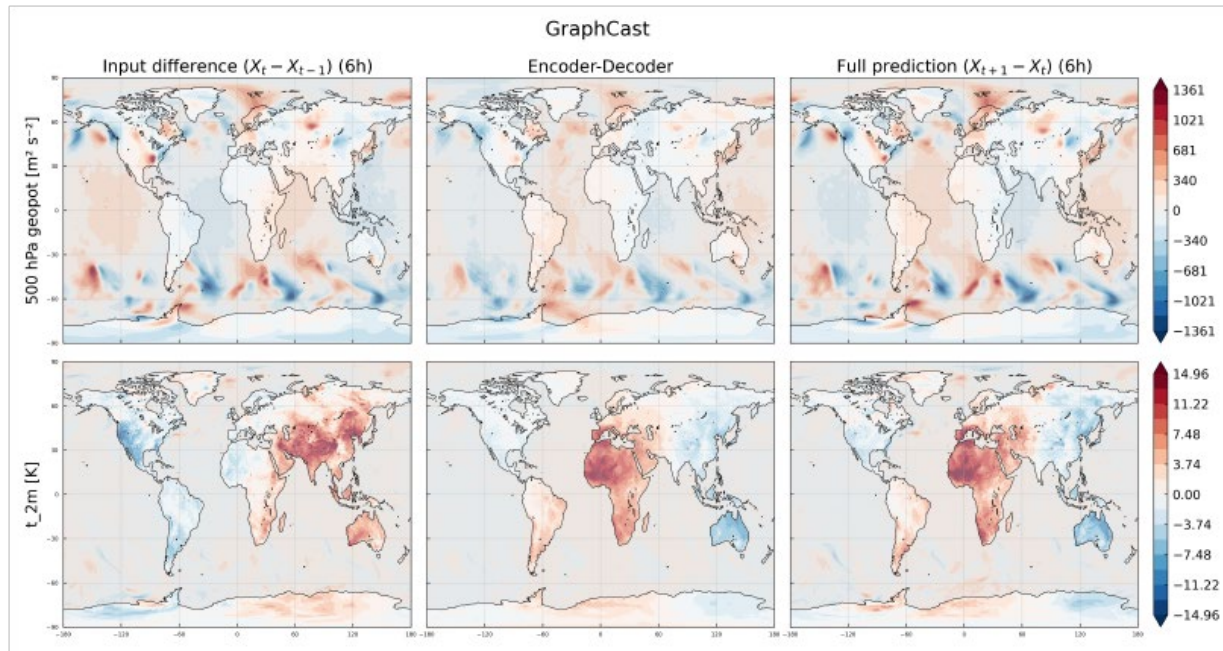
Latent channel development



Forecast progression



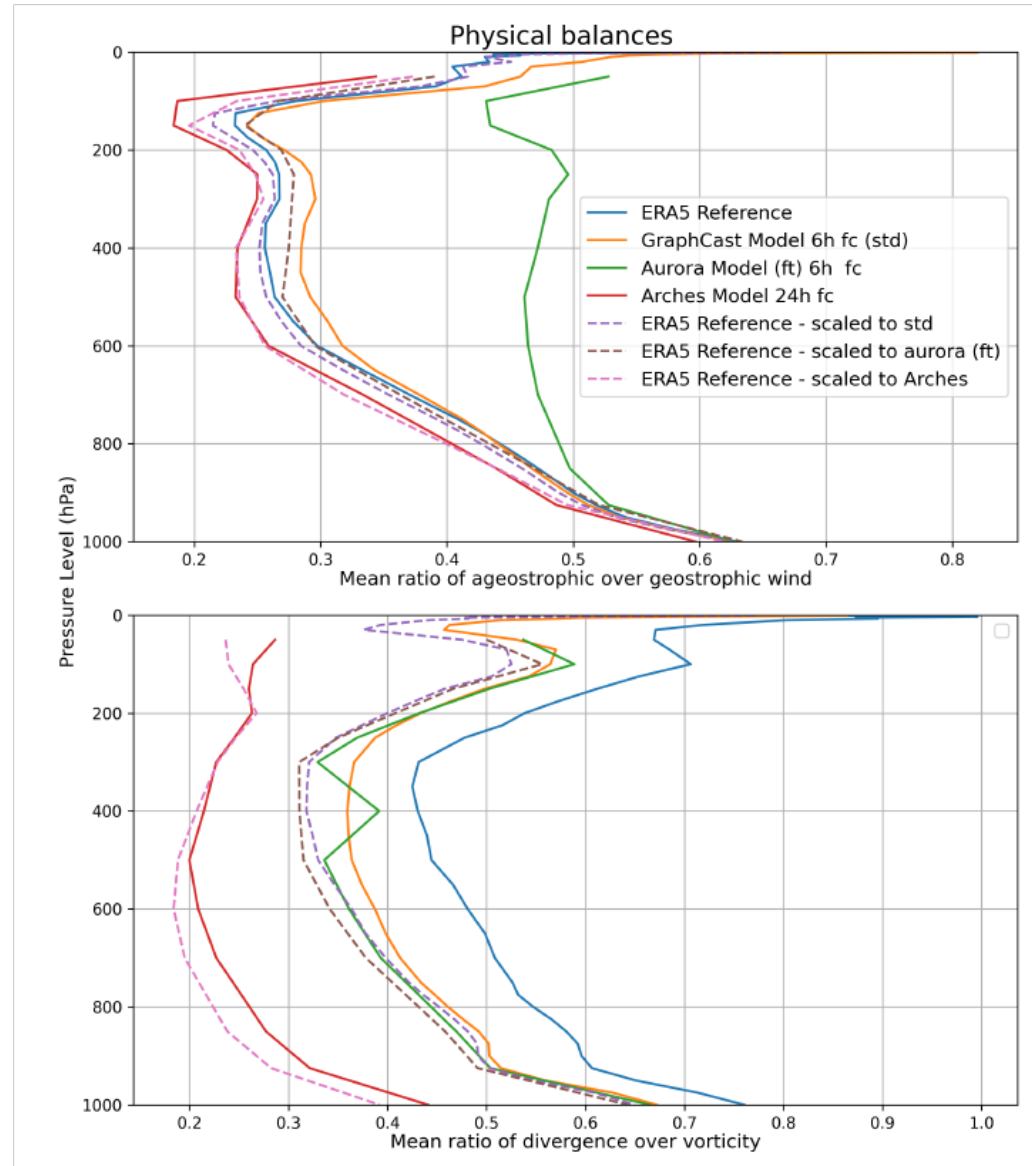
Encoder-Decoder



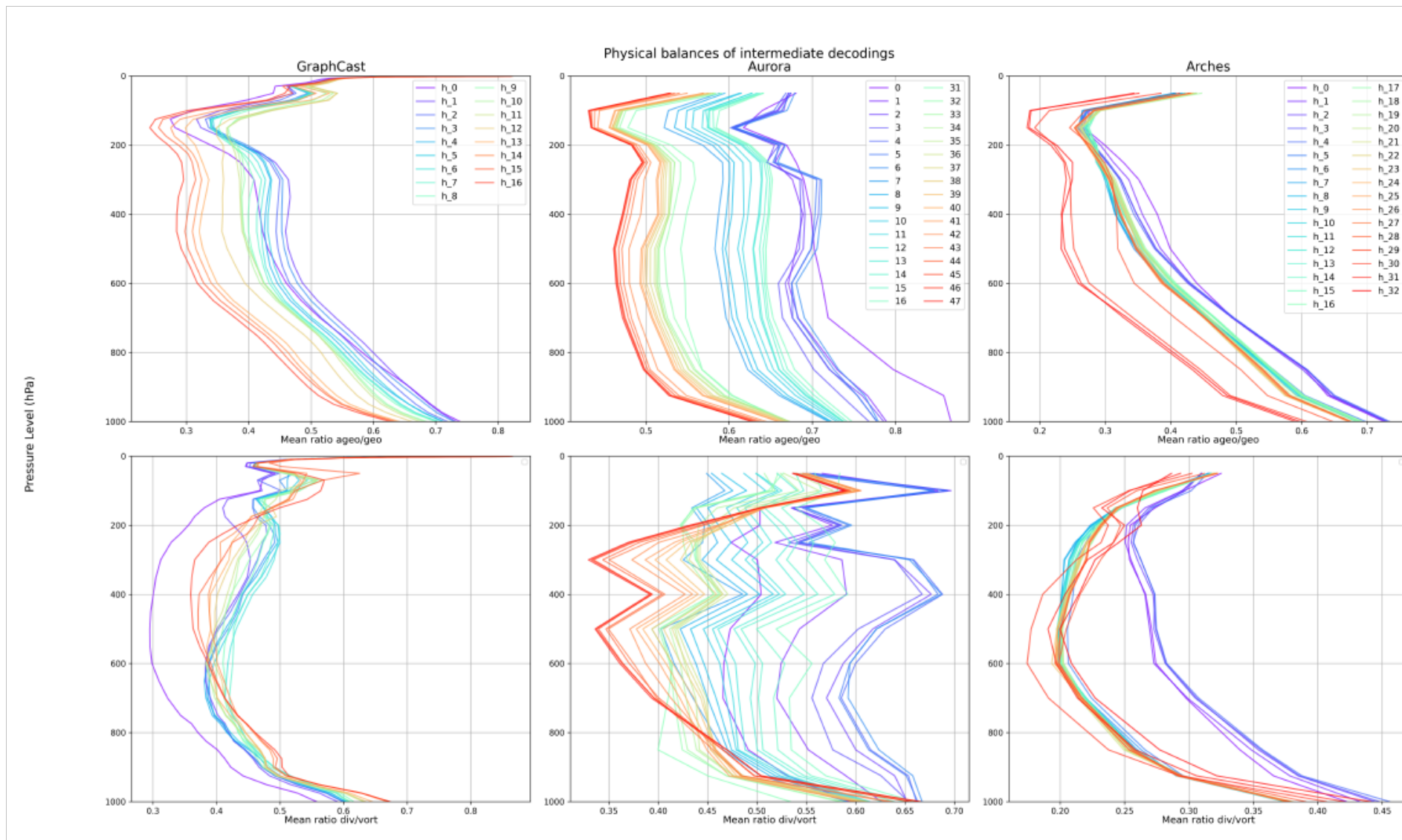
Scales



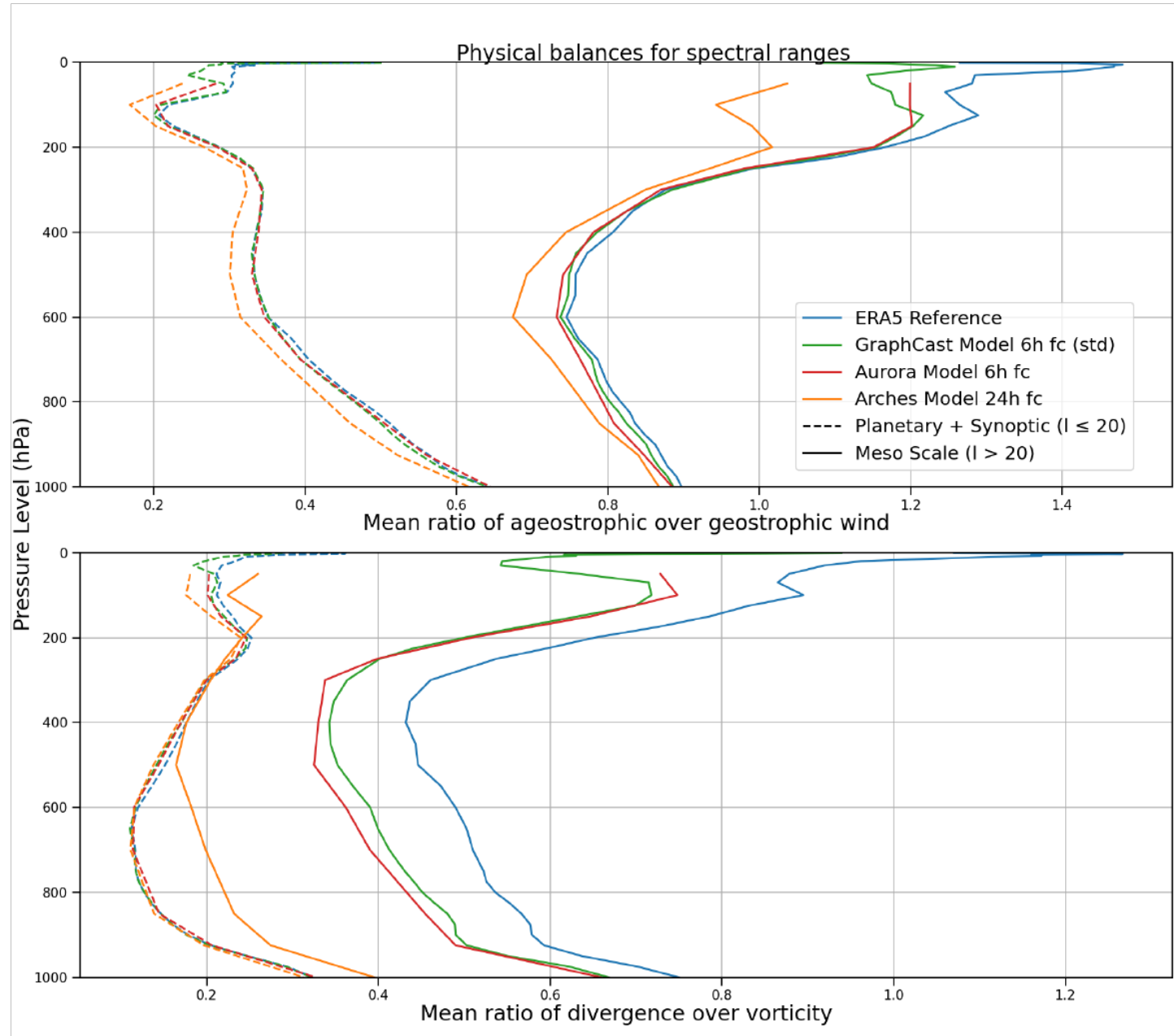
Physical balance



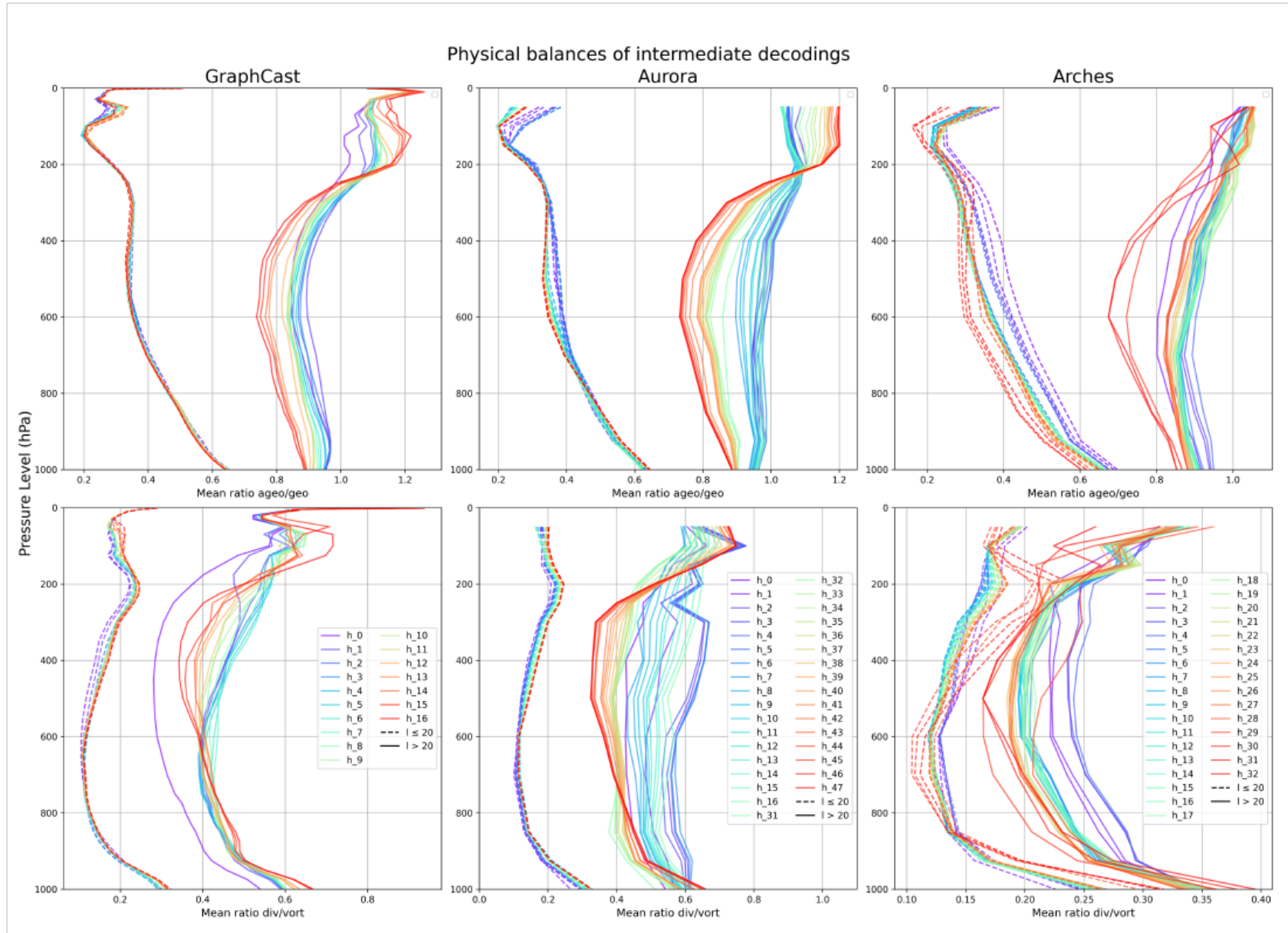
Physical balance



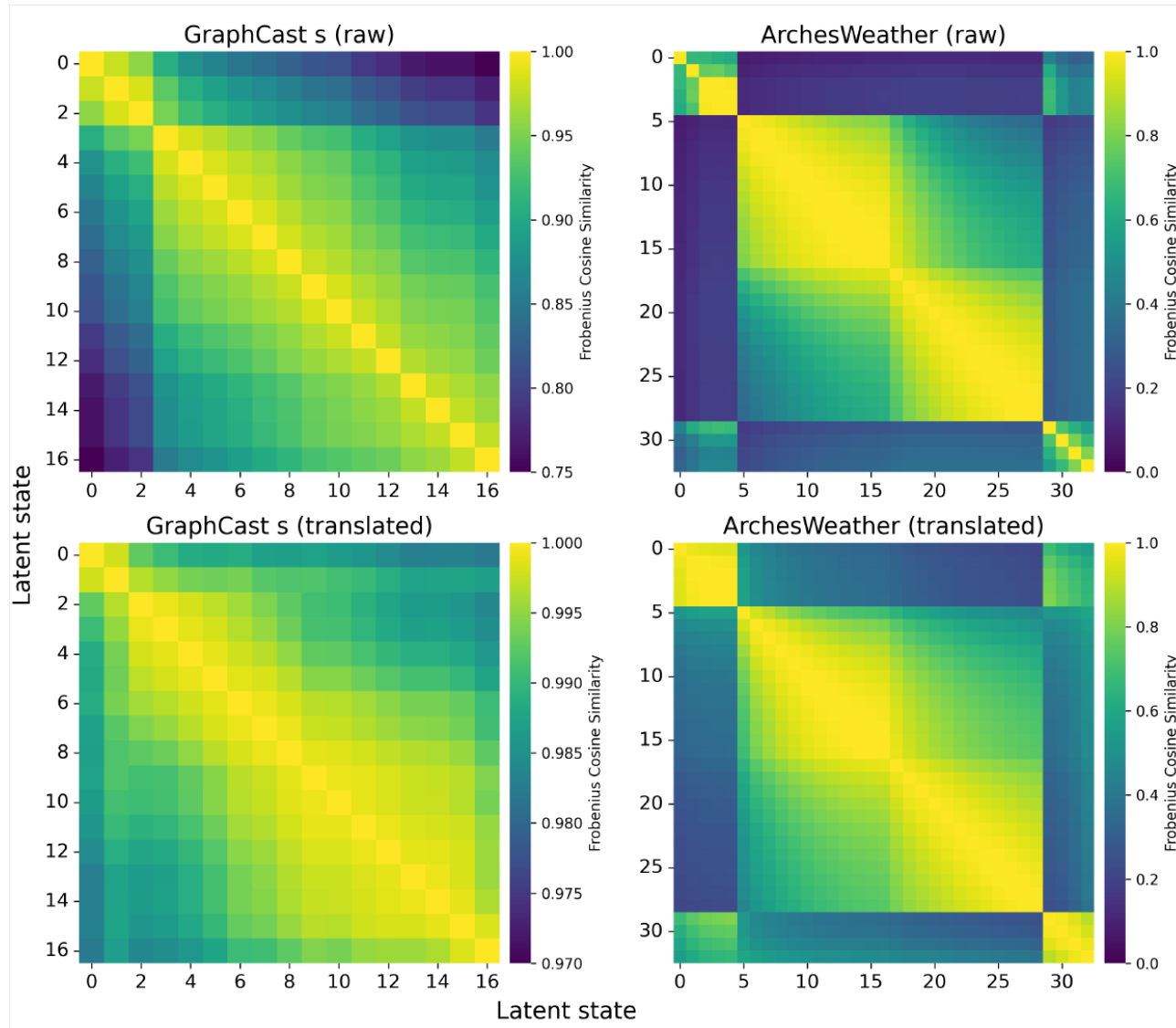
Physical balance



Physical balance

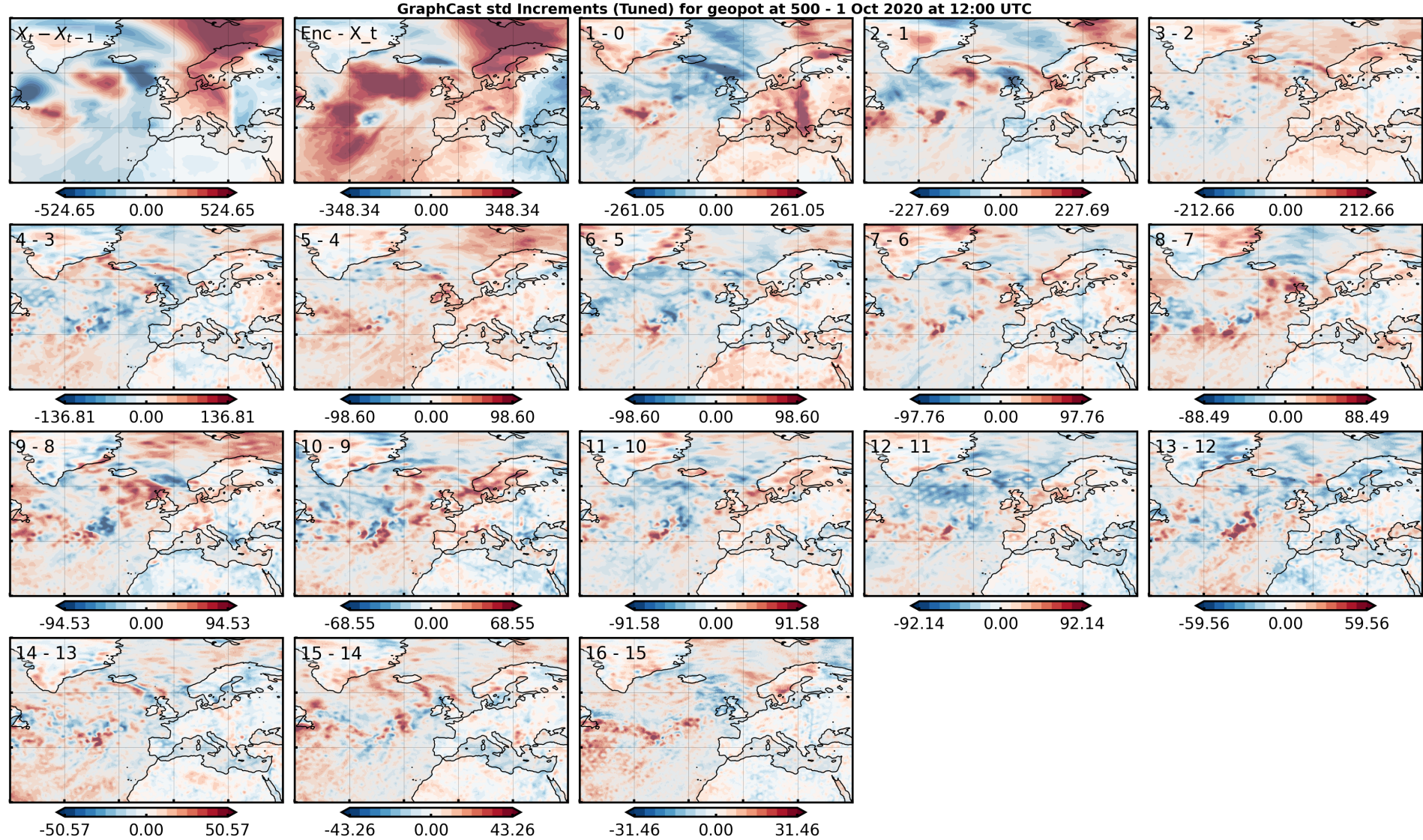


Representational stability

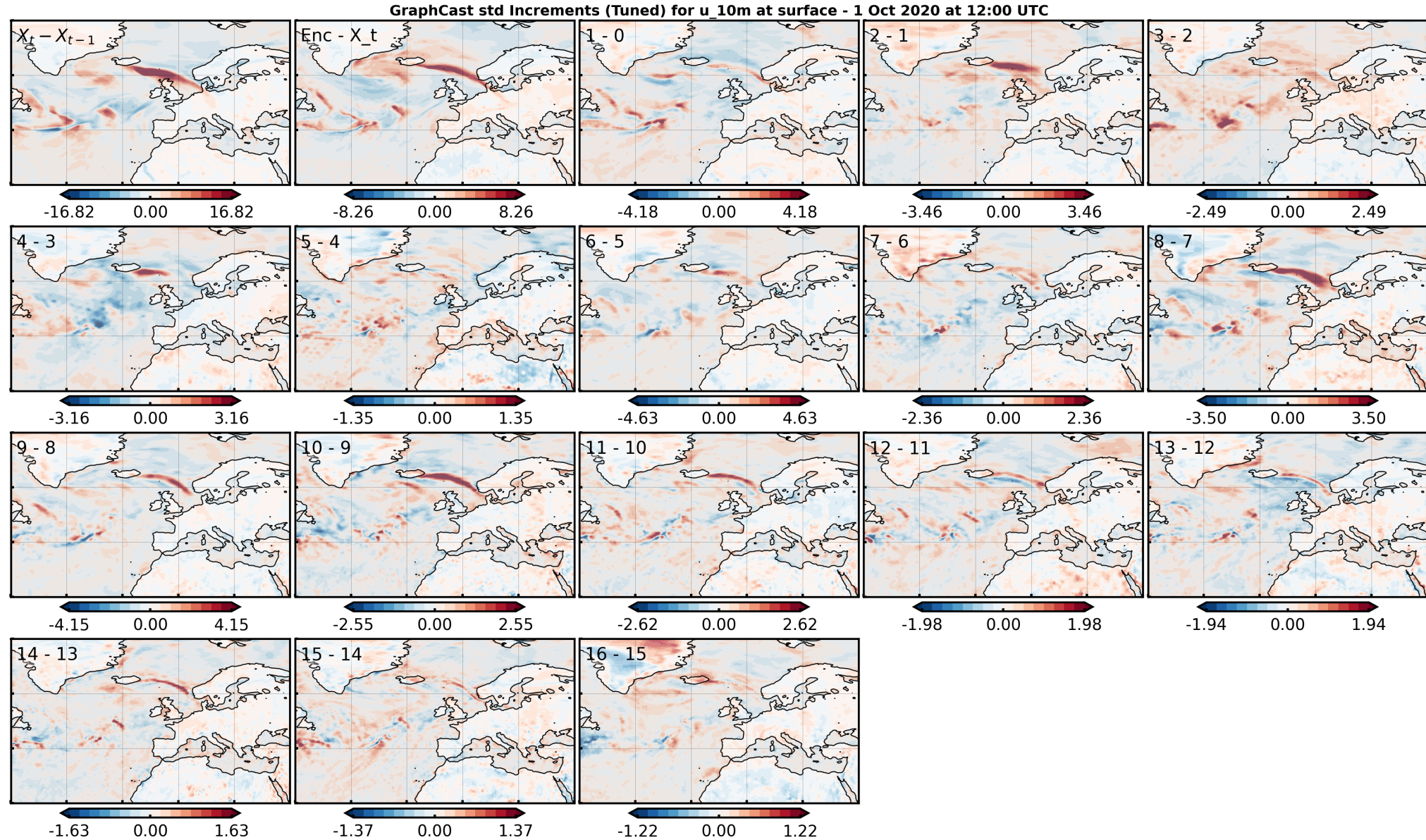


$$\cos_F(A, B) = \frac{\langle A, B \rangle_F}{\|A\|_F \|B\|_F},$$

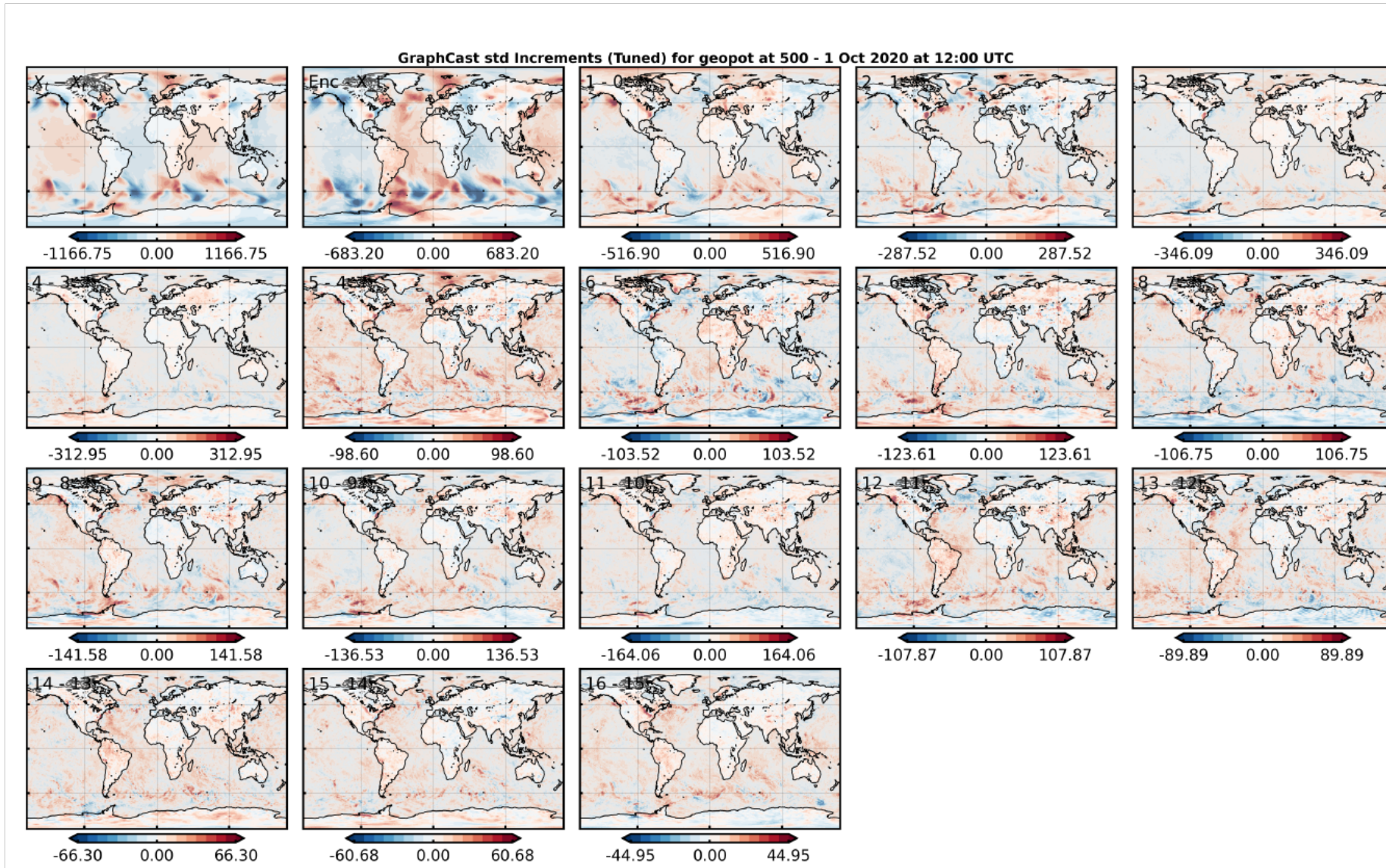
GraphCast all increments



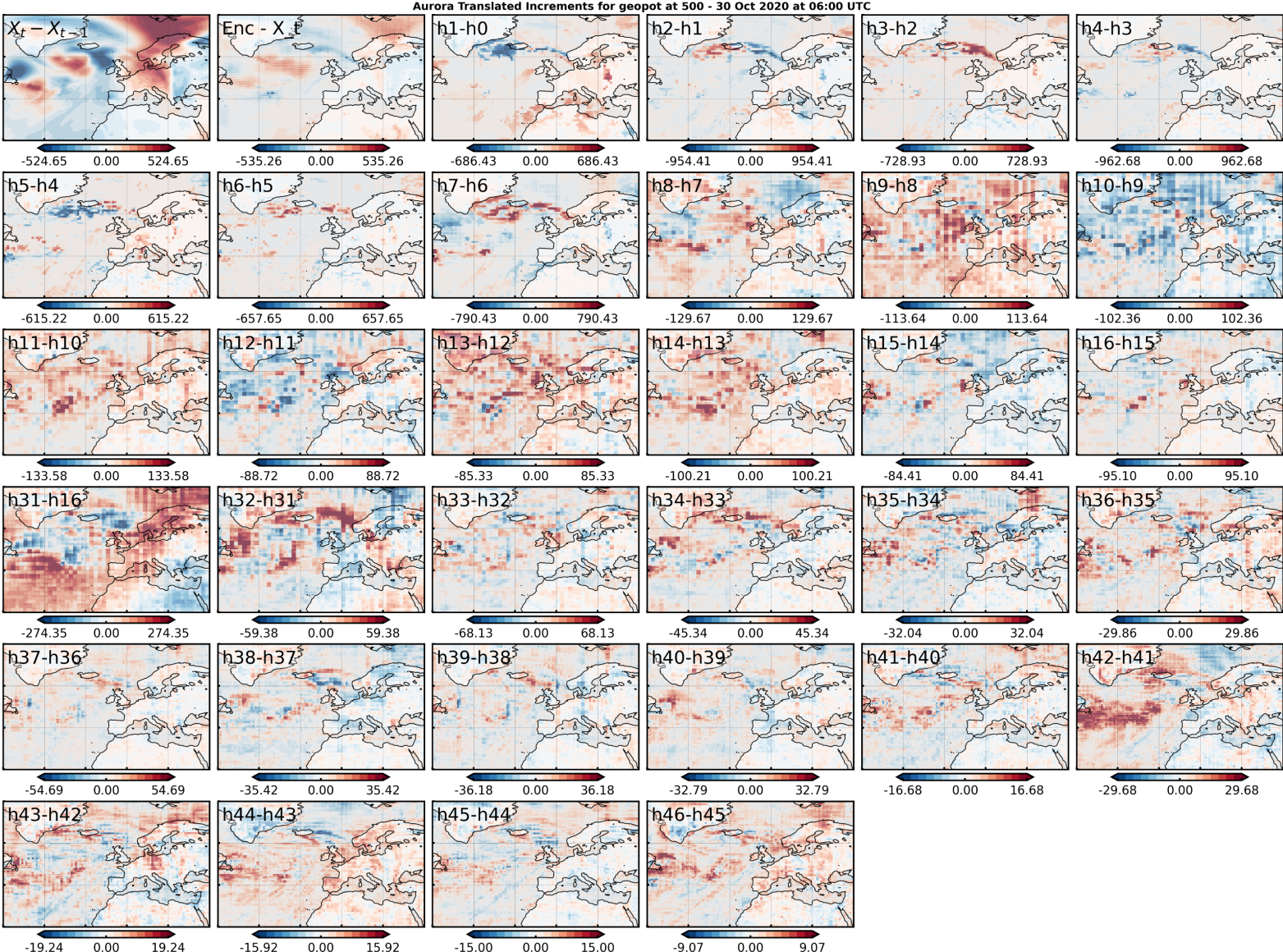
GraphCast all increments



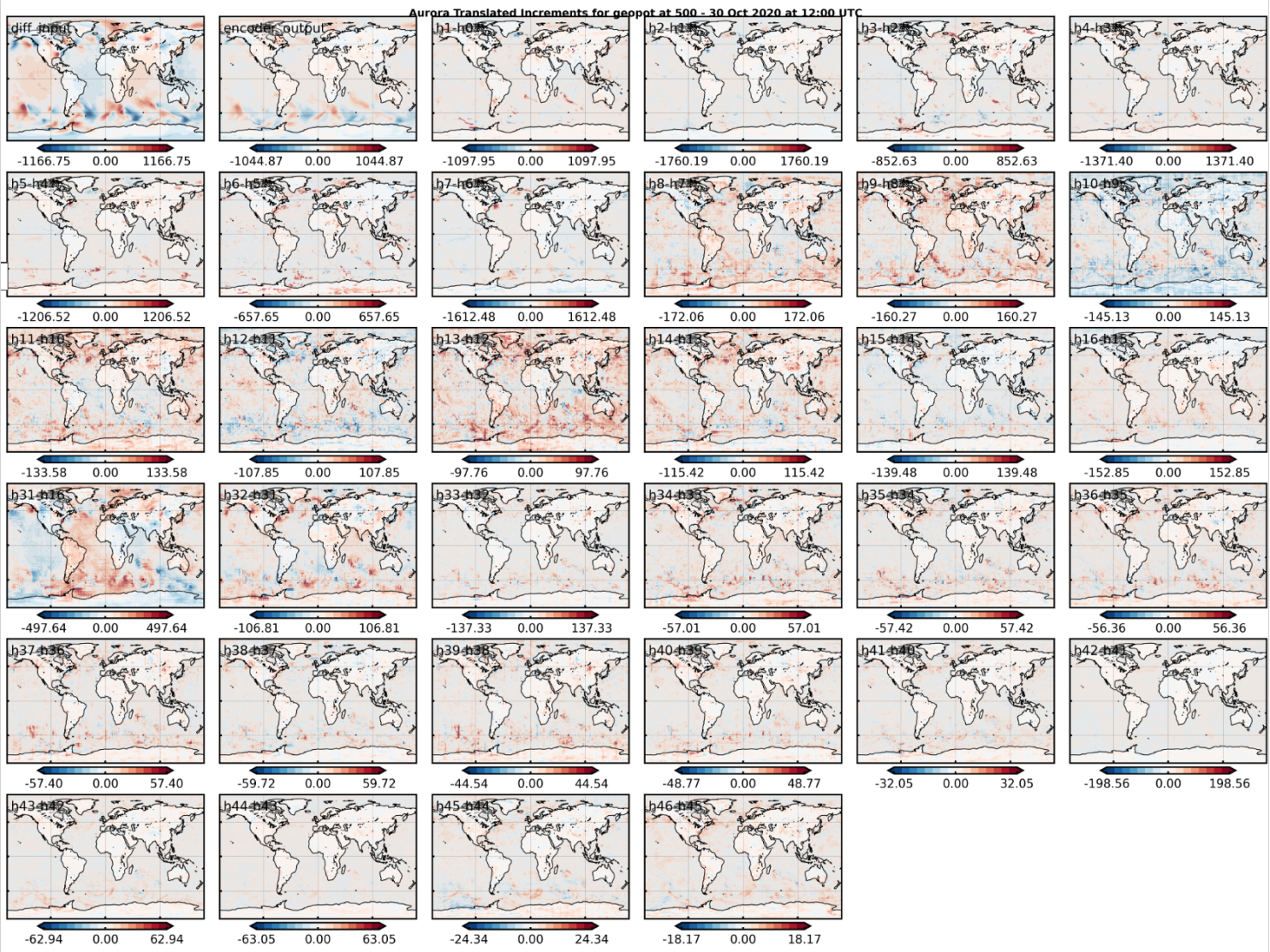
GraphCast all increments



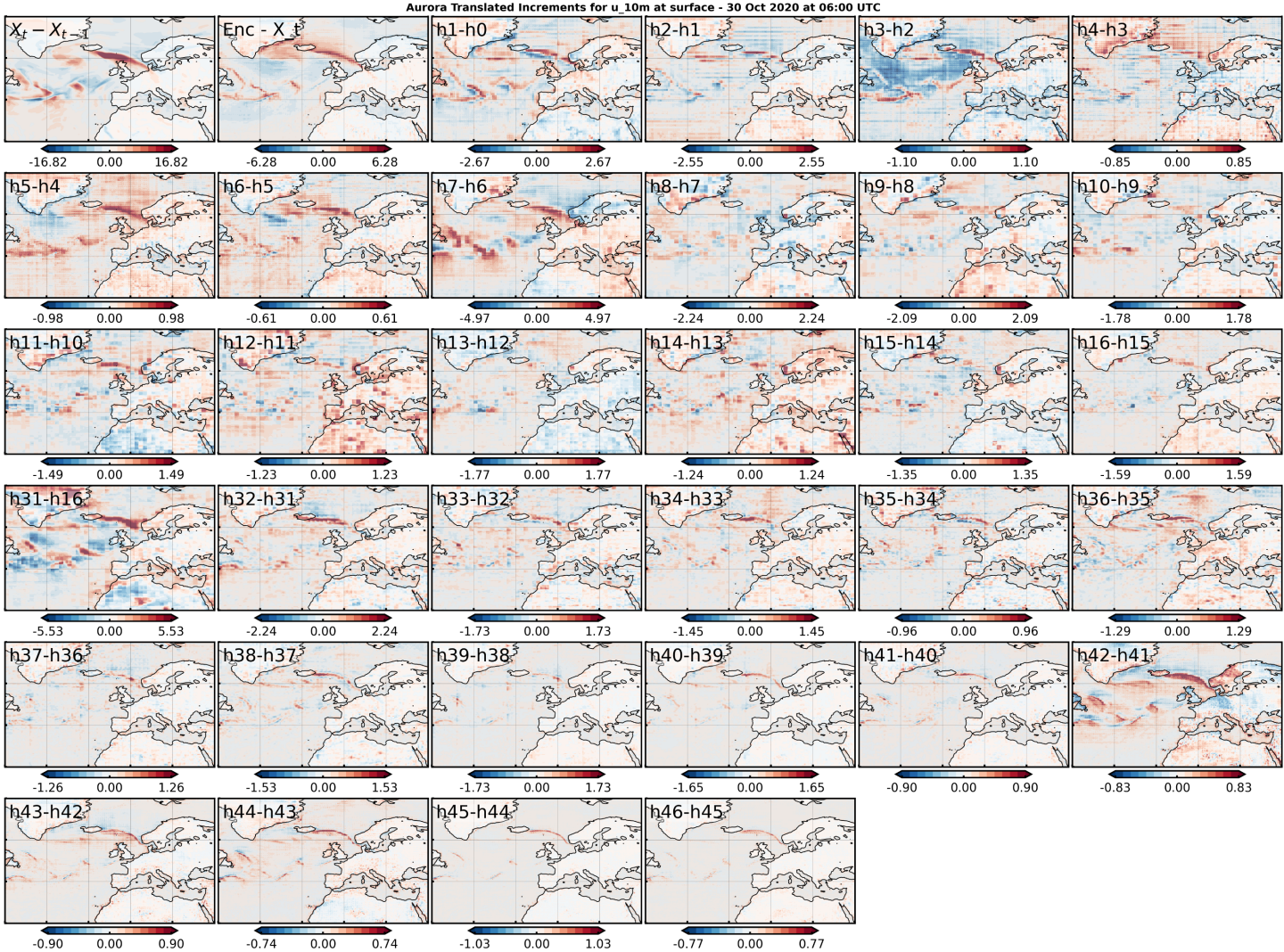
Aurora all increments



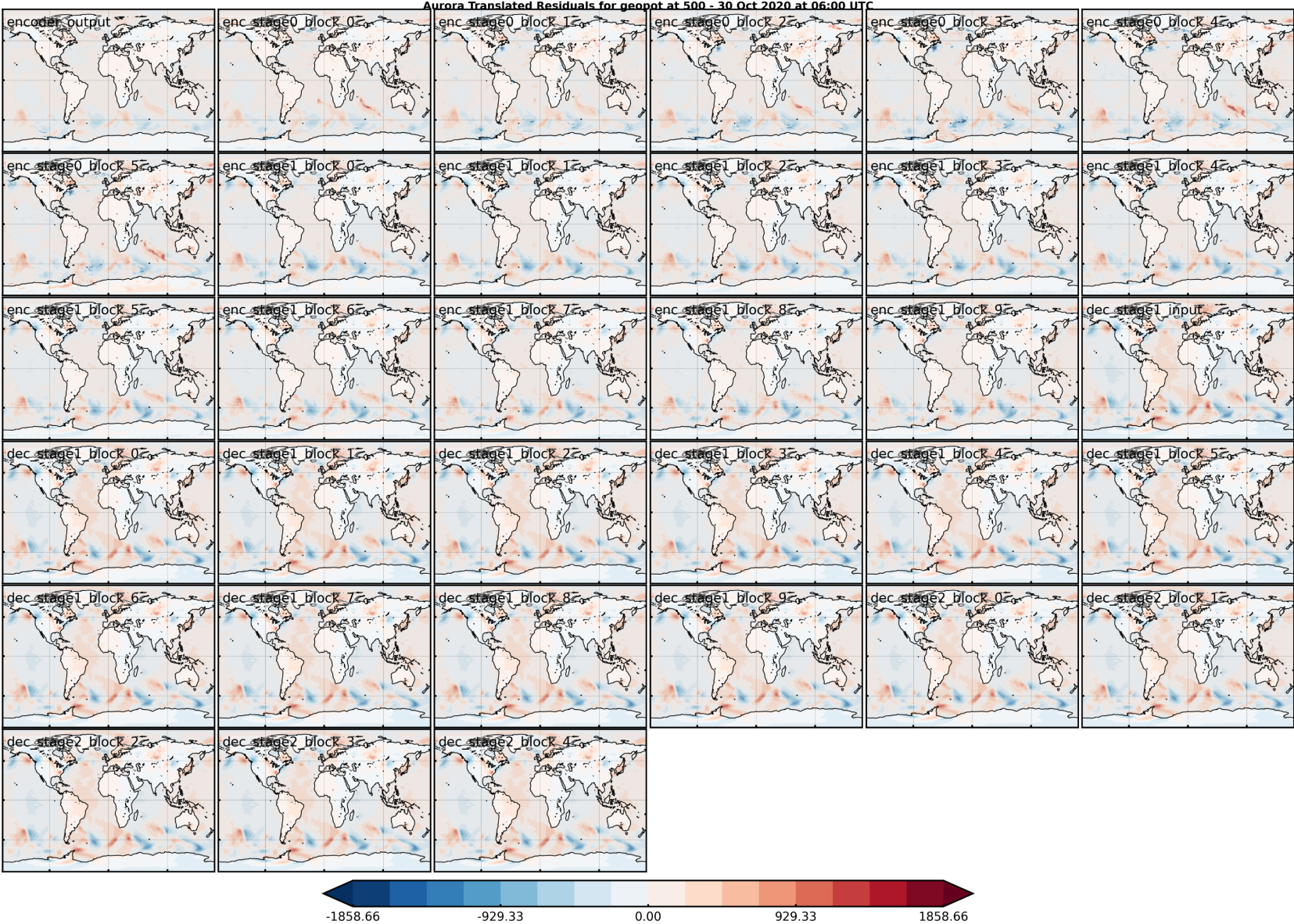
Aurora all increments



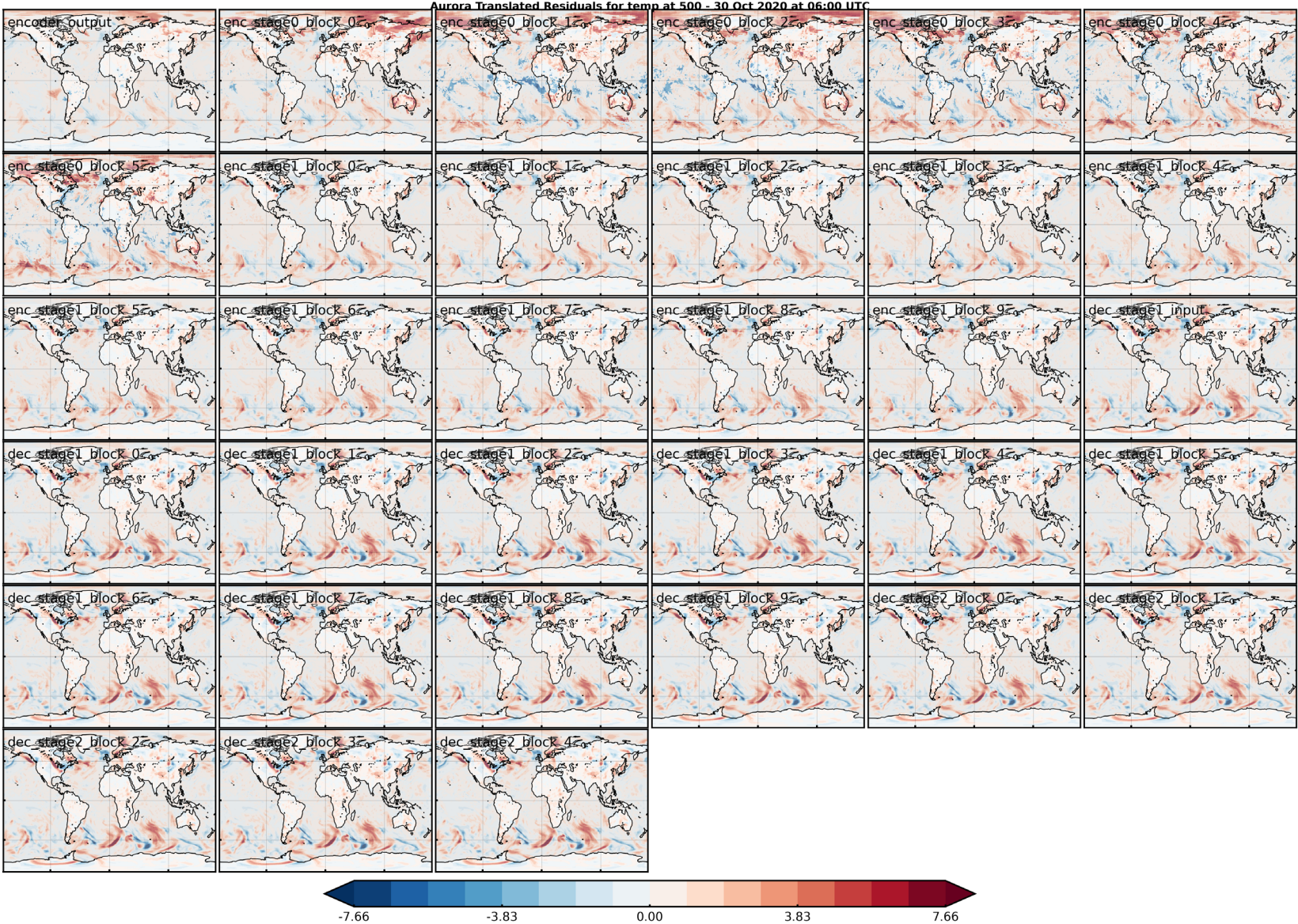
Aurora all increments



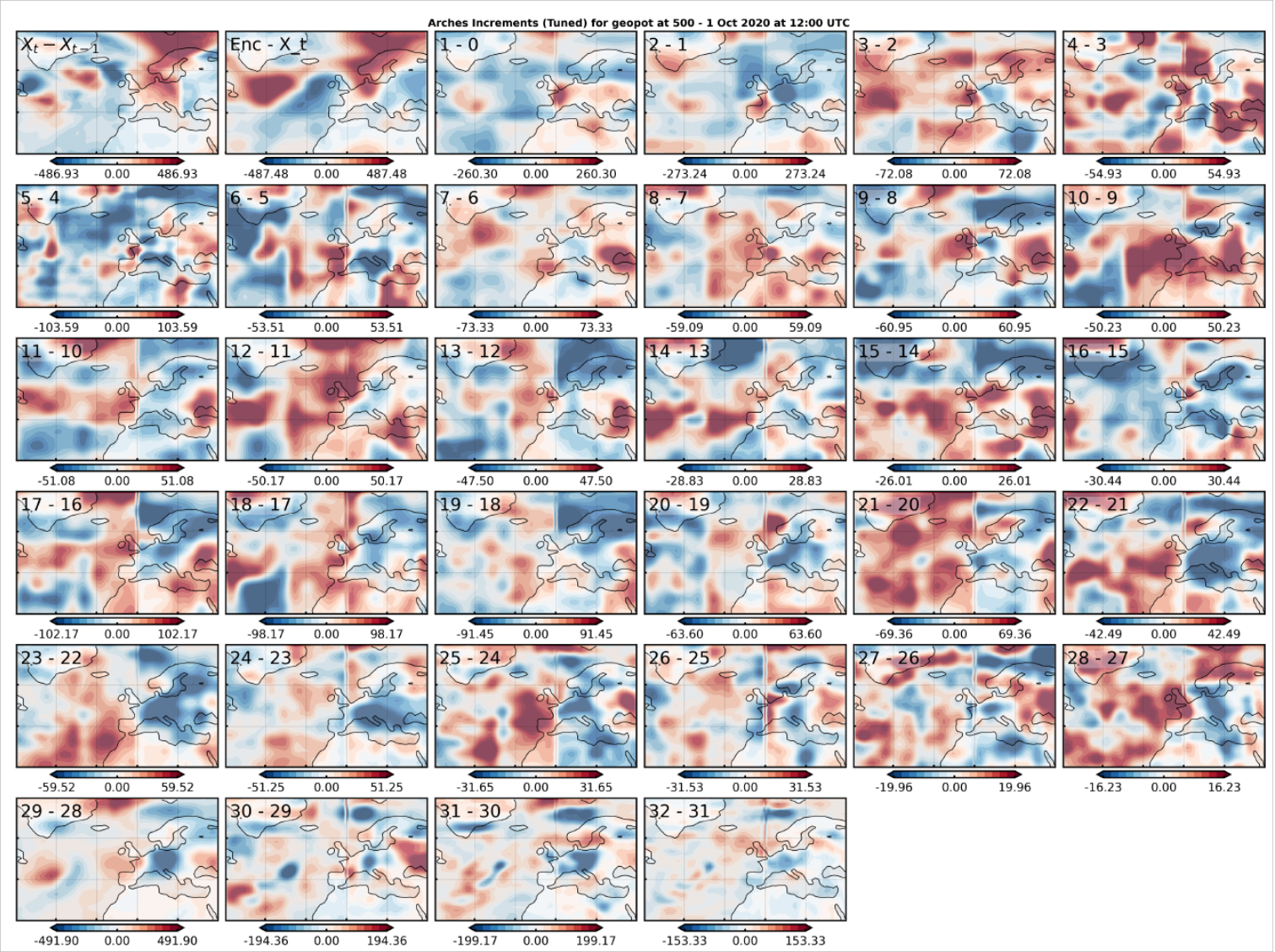
Aurora all increments



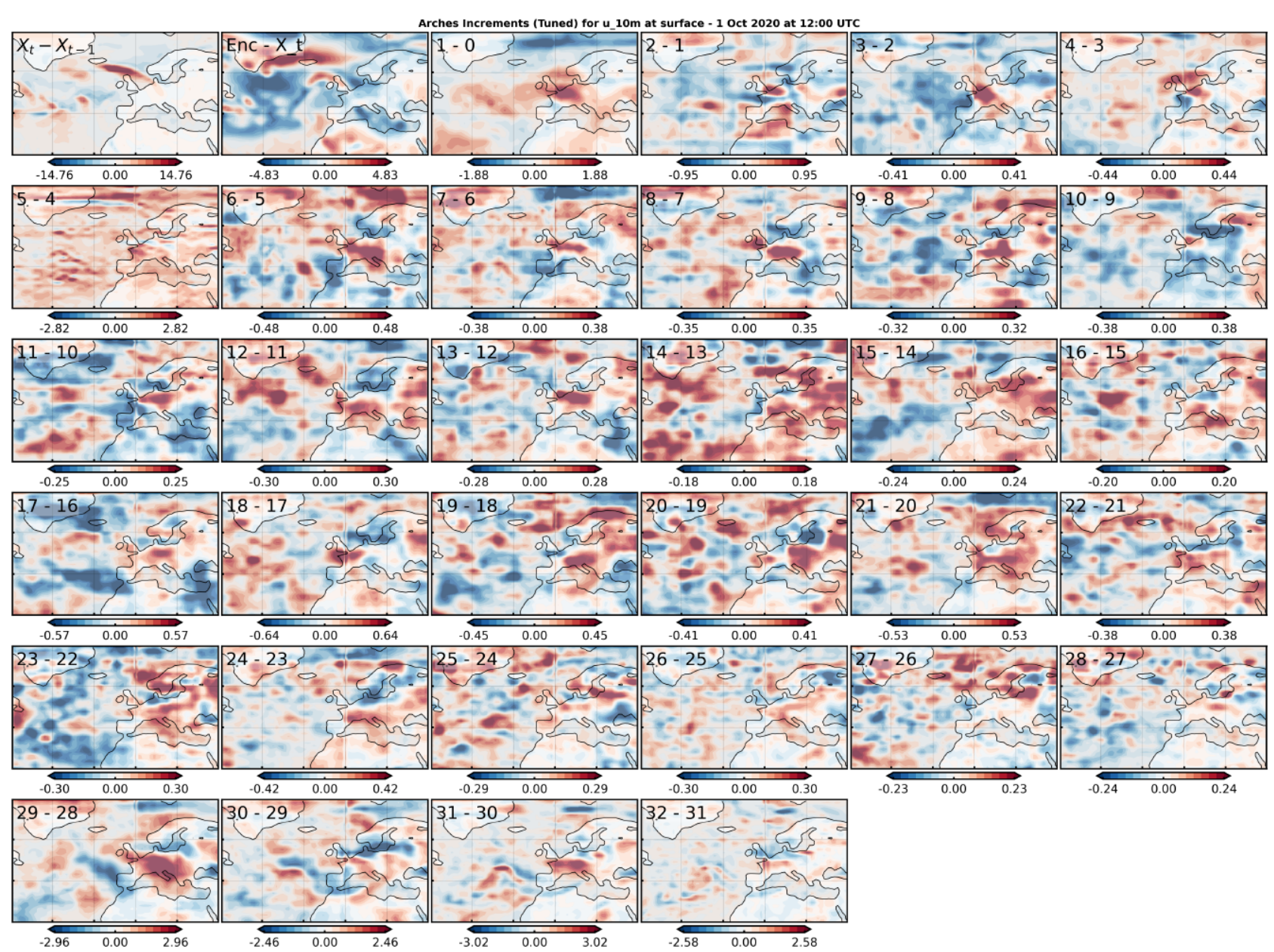
Aurora all increments



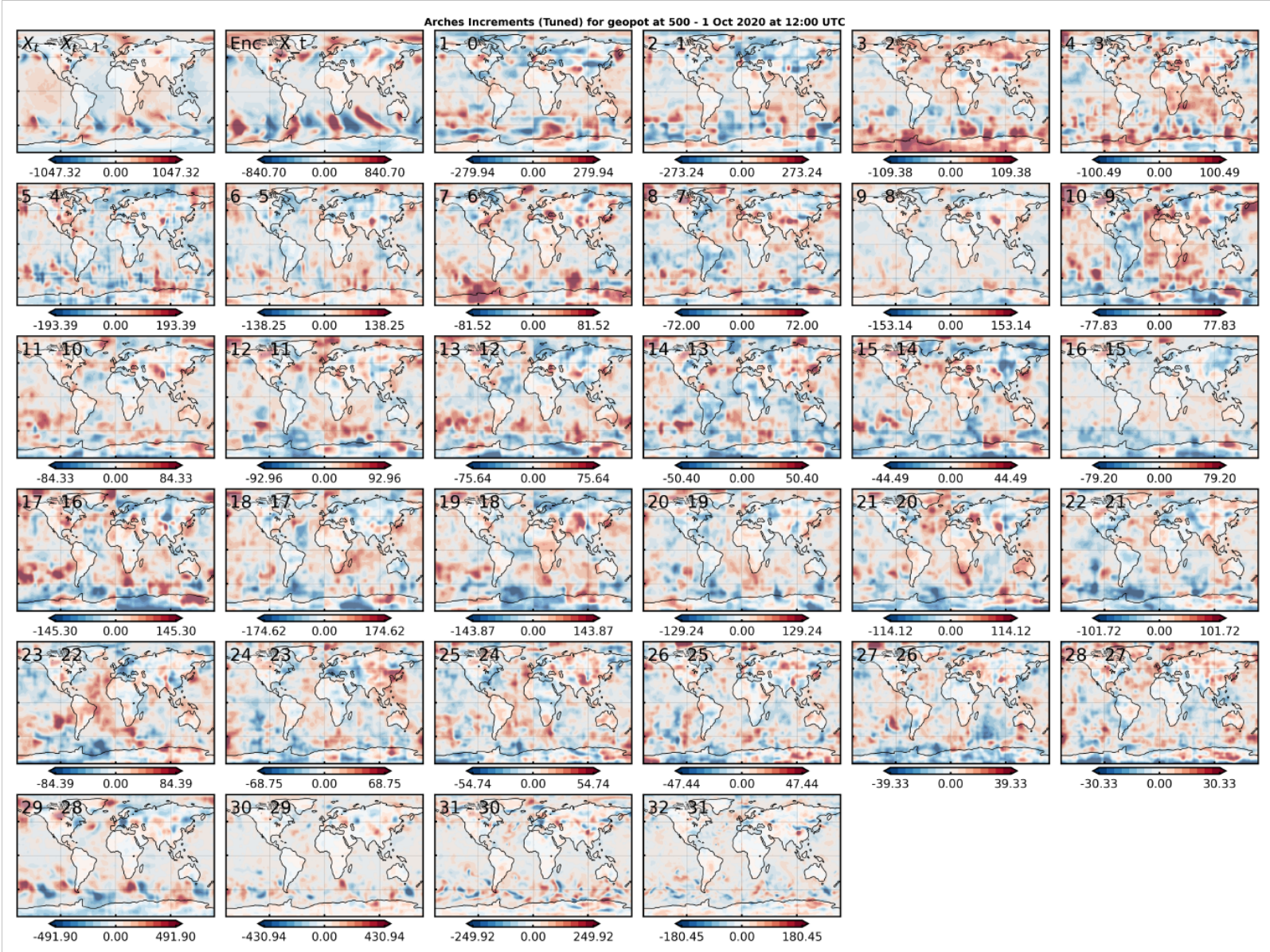
Arches all increments



Arches all increments



Arches all increments



Arches all increments

