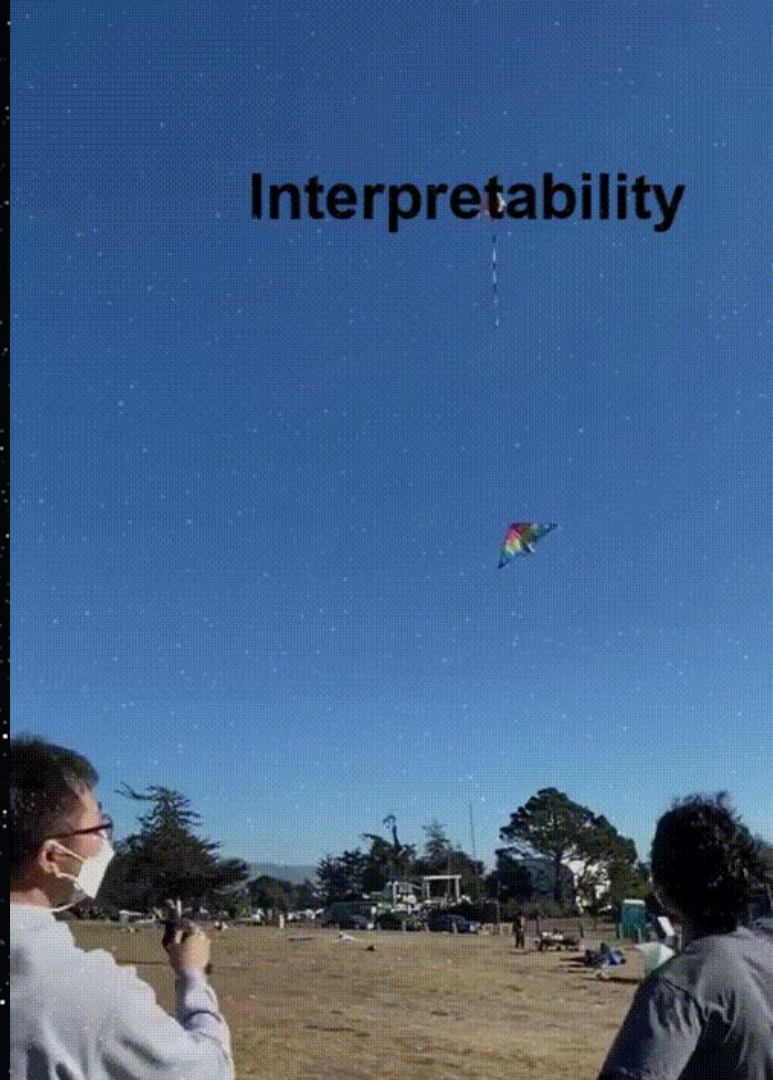


Adaptive wavelet distillation from neural networks through interpretations

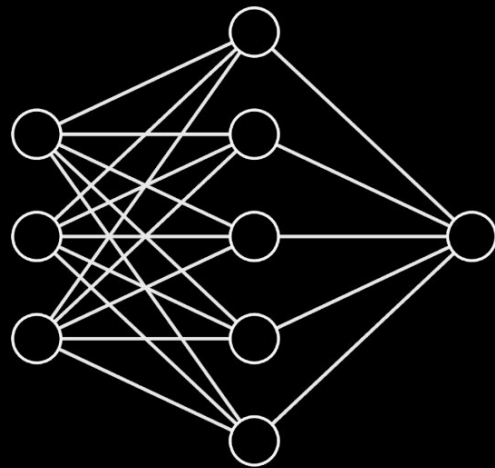
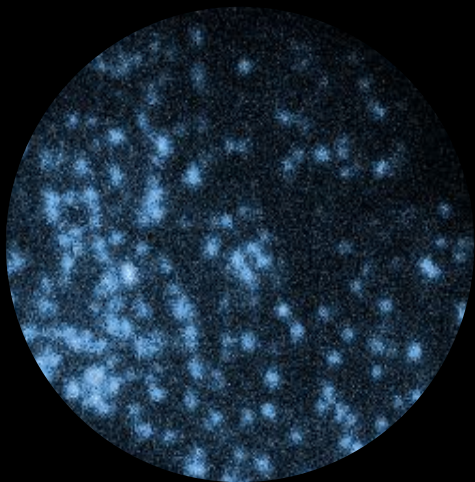
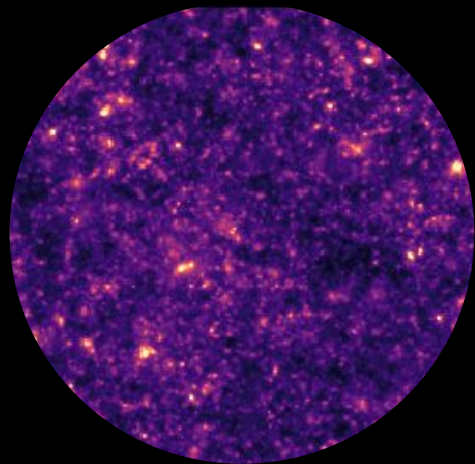


UC Berkeley

Interpretability



a deep problem: interpretability



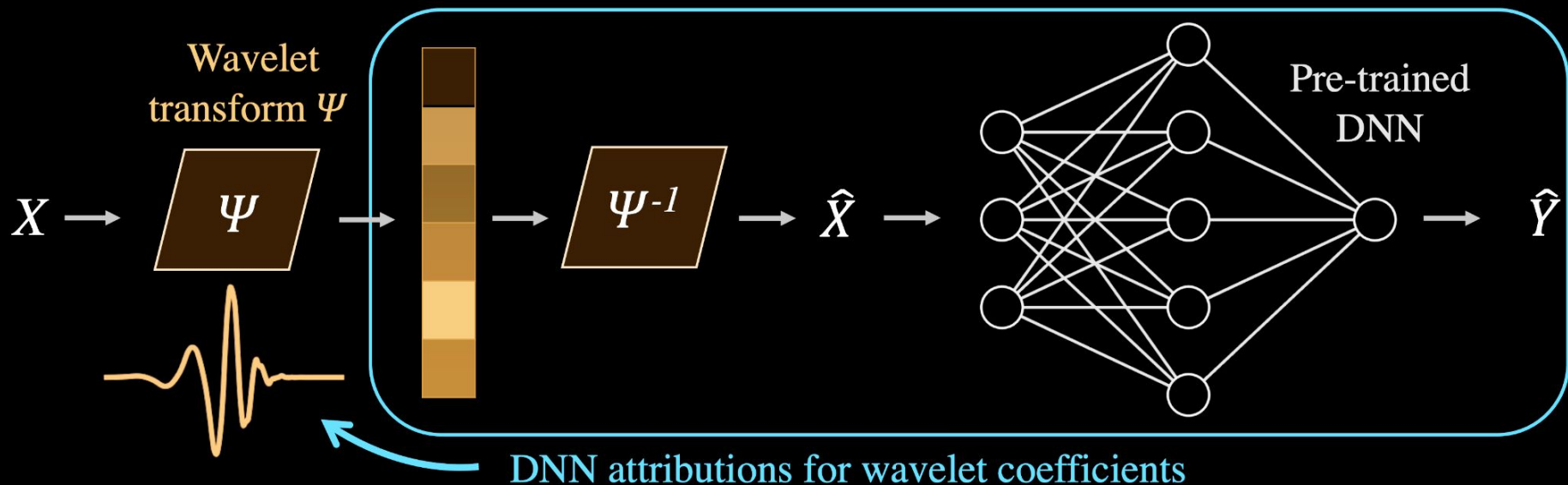
how???



prediction

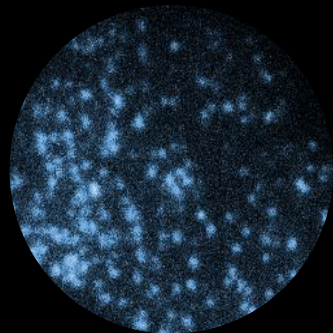
why???

a shallow solution: wavelets



$$\underset{h,g}{\text{minimize}} \mathcal{L}(h,g) = \underbrace{\frac{1}{m} \sum_i \|x_i - \hat{x}_i\|_2^2}_{\text{Reconstruction loss}} + \underbrace{\frac{1}{m} \sum_i W(h,g,x_i;\lambda)}_{\text{Wavelet loss}} + \underbrace{\gamma \sum_i \|\text{TRIM}_{\Psi,f}(\Psi x_i)\|_1}_{\text{Interpretation loss}}$$

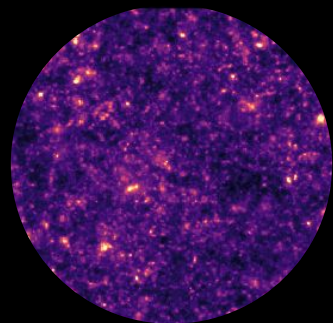
results (spoilers!)



cell-biology: CME event prediction

LSTM	SOA Baseline
0.237	0.197

R^2 (higher is better)

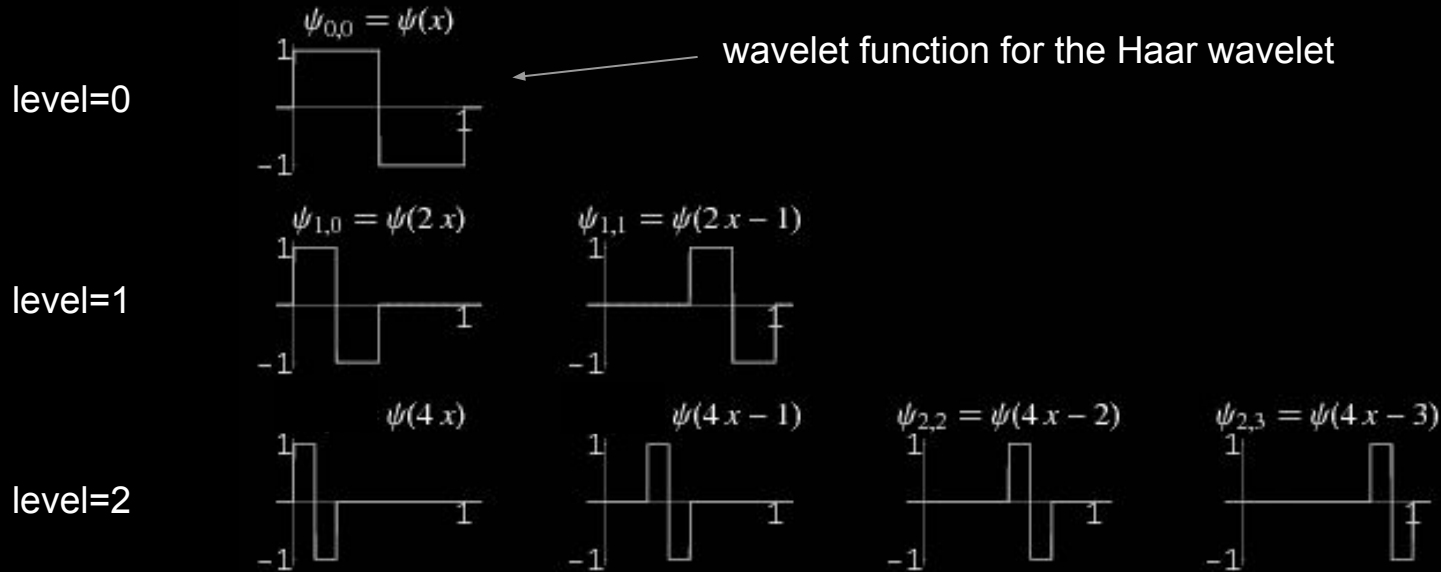


cosmological parameter prediction

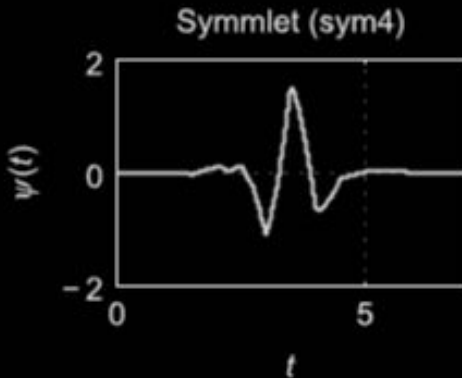
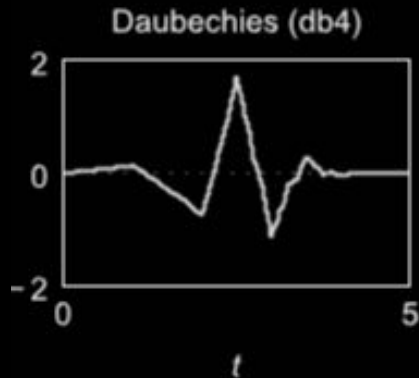
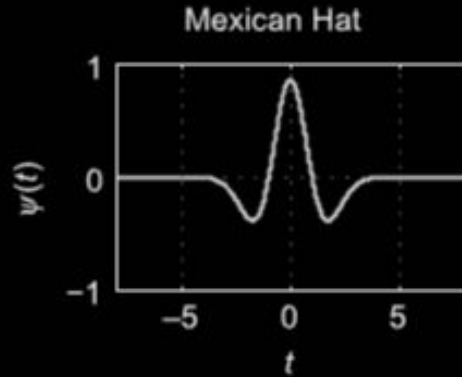
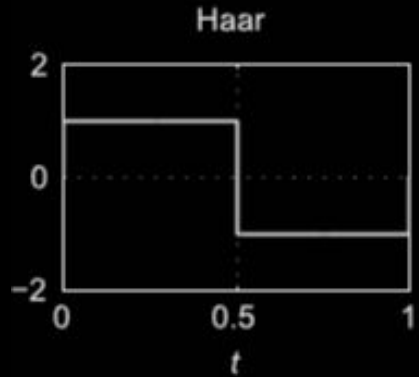
Resnet	SOA Baseline
1.156	1.259

RMSE (lower is better)

wavelet transform: multi-scale, spatially localized

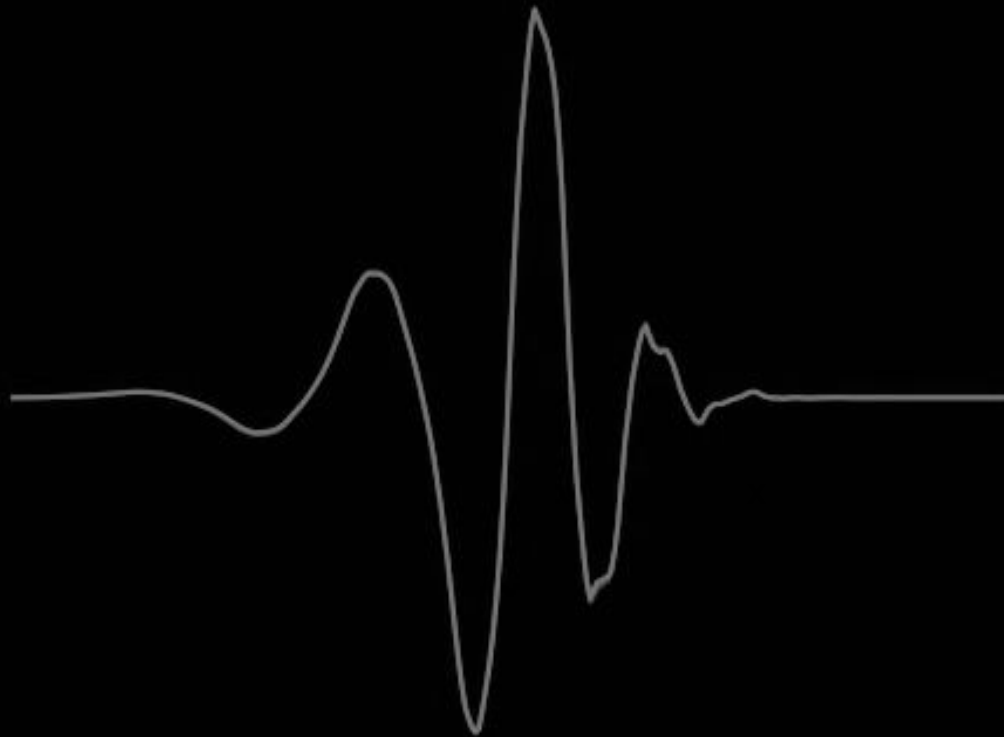


wavelet function can vary



adaptive wavelet

$$\underset{h,g}{\text{minimize}} \mathcal{L}(h, g) = \underbrace{\frac{1}{m} \sum_i \|x_i - \hat{x}_i\|_2^2}_{\text{Reconstruction loss}} + \underbrace{\frac{1}{m} \sum_i W(h, g, x_i; \lambda)}_{\text{Wavelet loss}}$$



orthonormal basis under following conditions (mallat, 1998):

$$|\hat{h}(w)|^2 + |\hat{h}(w + \pi)|^2 = 2 \quad \forall w$$

$$\sum h_i = \sqrt{2}$$

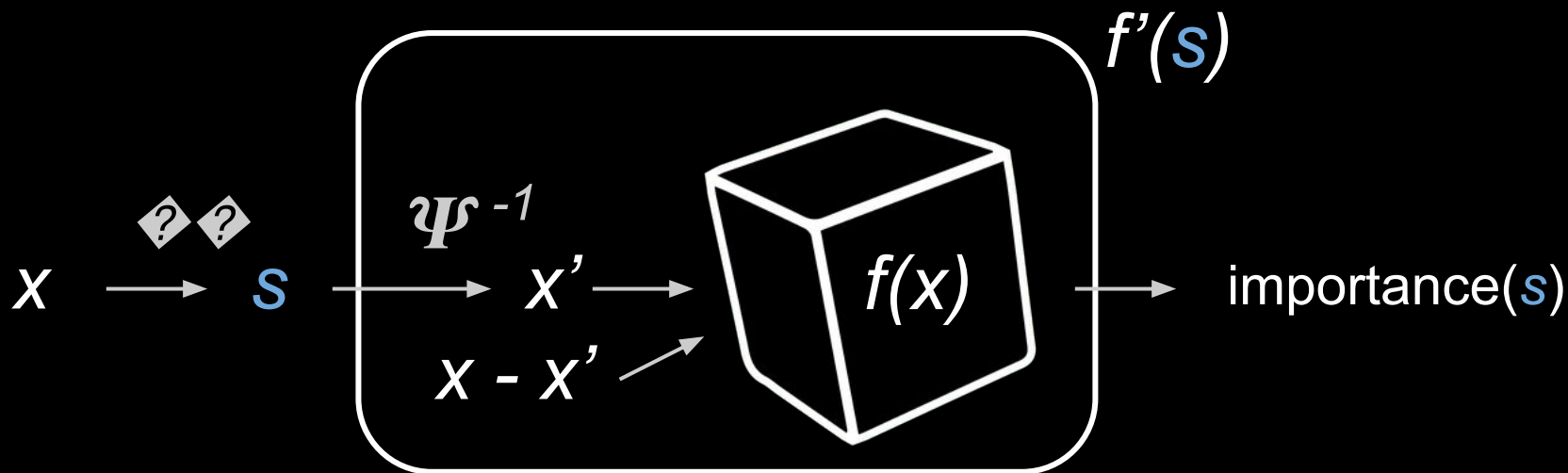
$$g(n) = (-1)^n h(N - 1 - n)$$

$$\|\Phi f\|_1 \text{ small}$$

adaptive wavelet + distillation

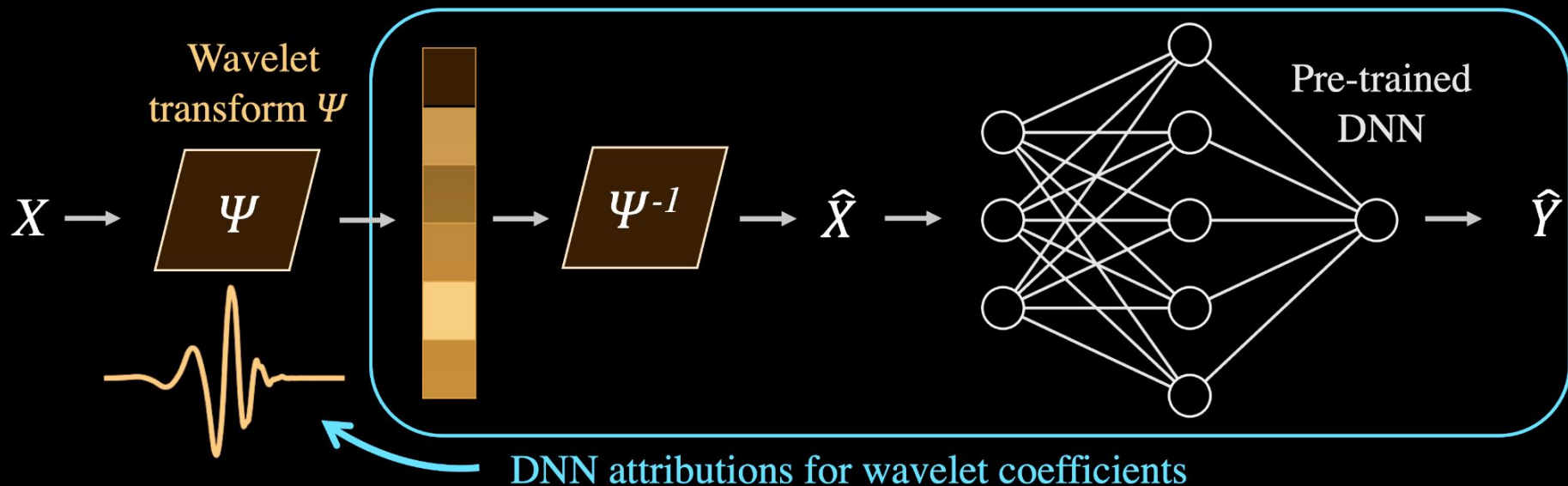
$$\underset{h,g}{\text{minimize}} \mathcal{L}(h, g) = \underbrace{\frac{1}{m} \sum_i \|x_i - \hat{x}_i\|_2^2}_{\text{Reconstruction loss}} + \underbrace{\frac{1}{m} \sum_i W(h, g, x_i; \lambda)}_{\text{Wavelet loss}} + \underbrace{\gamma \sum_i \|\text{TRIM}_{\Psi, f}(\Psi x_i)\|_1}_{\text{Interpretation loss}},$$

transformation importance



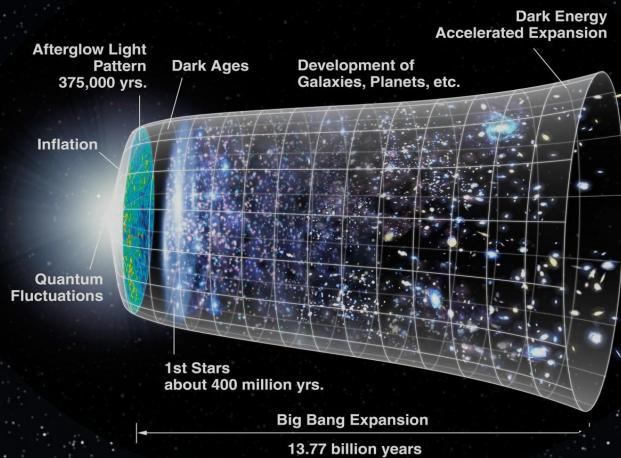
singh, ha, lanusse, boehm, liu, & yu, 2019
“transformation importance with applications to cosmology”

putting it all together

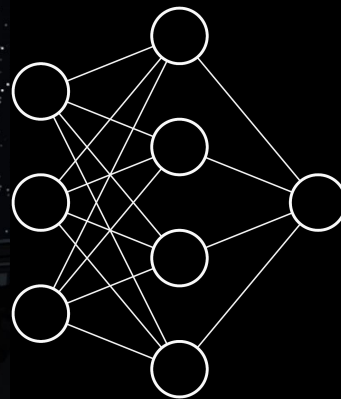
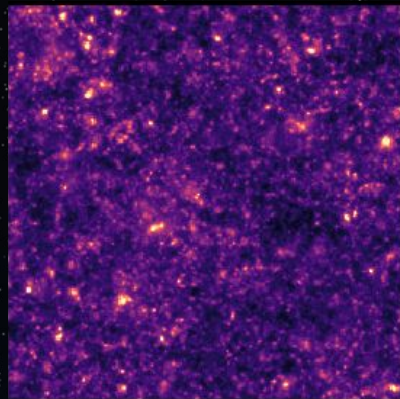


$$\underset{h,g}{\text{minimize}} \mathcal{L}(h,g) = \underbrace{\frac{1}{m} \sum_i \|x_i - \hat{x}_i\|_2^2}_{\text{Reconstruction loss}} + \underbrace{\frac{1}{m} \sum_i W(h,g,x_i;\lambda)}_{\text{Wavelet loss}} + \underbrace{\gamma \sum_i \|\text{TRIM}_{\Psi,f}(\Psi x_i)\|_1}_{\text{Interpretation loss}},$$

cosmology problem
(more in the paper)



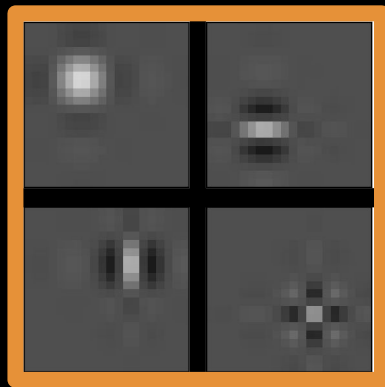
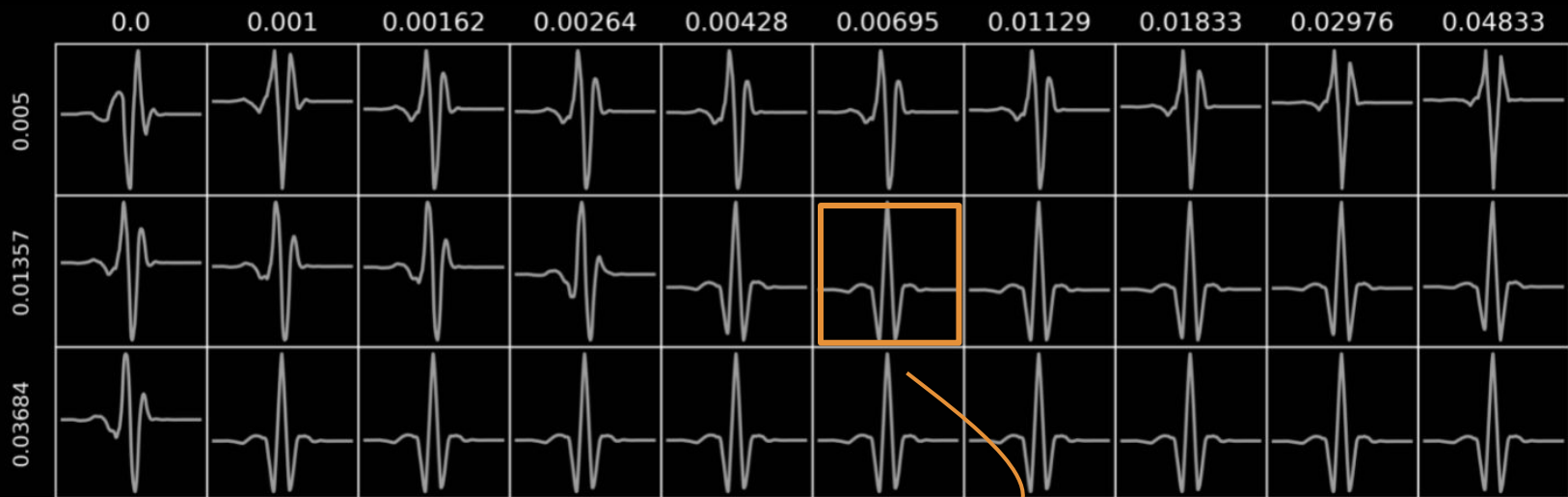
Ω_m



$\hat{\Omega}_m$

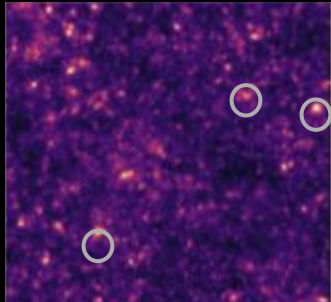
Increasing attribution penalty $\gamma \rightarrow$

\leftarrow Increasing sparsity penalty λ

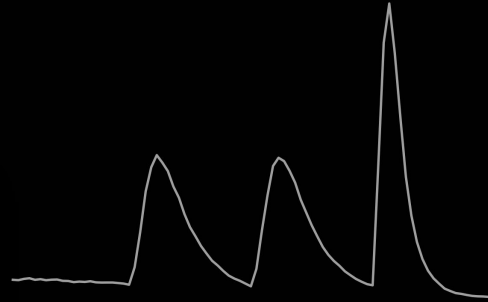
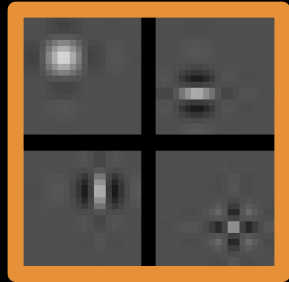


Recon loss = 0.0089
Wavelet loss = 0.0013

Predicting via peak counts



Filtered values
at the local maxima



Binning into histogram

Predict using
nearest-neighbor

$$\hat{\Omega}_m$$

Prediction error for Ω_m (RMSE)

AWD	Resnet	AWD no interp. loss
1.029	1.156	1.354

x 10⁻⁴

Peak Height	Laplace*	Roberts-Cross*	DB5
1.609	1.369	1.259	1.569

x 10⁻⁴

*Ribli et al (2019) Nature Astronomy