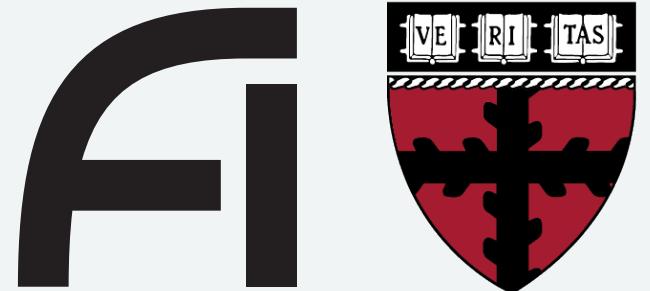


# Constraining neural networks to craft representations

“How can we structurally enforce desirable properties on neural networks?”

Friday, March 17, 2022

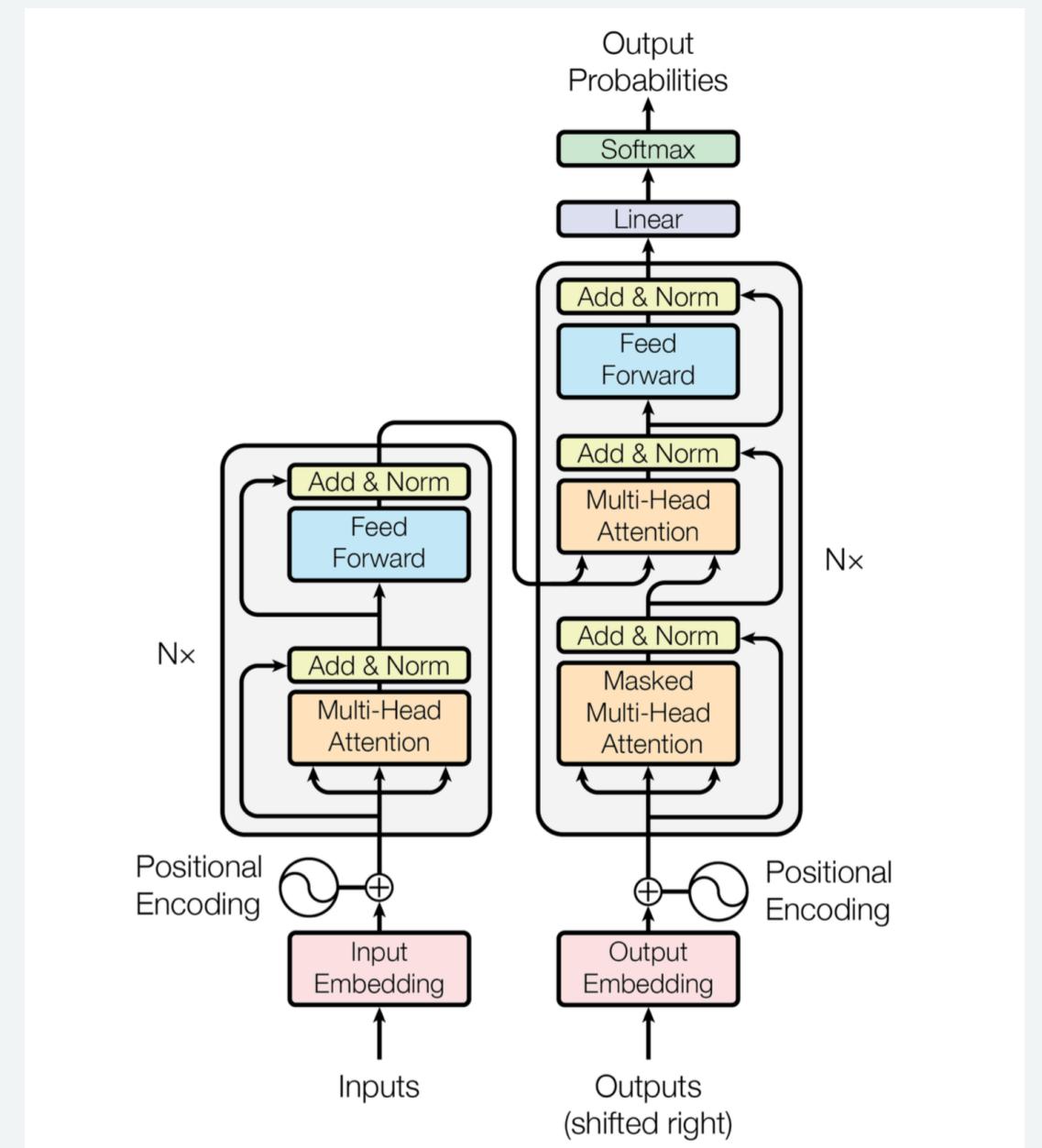
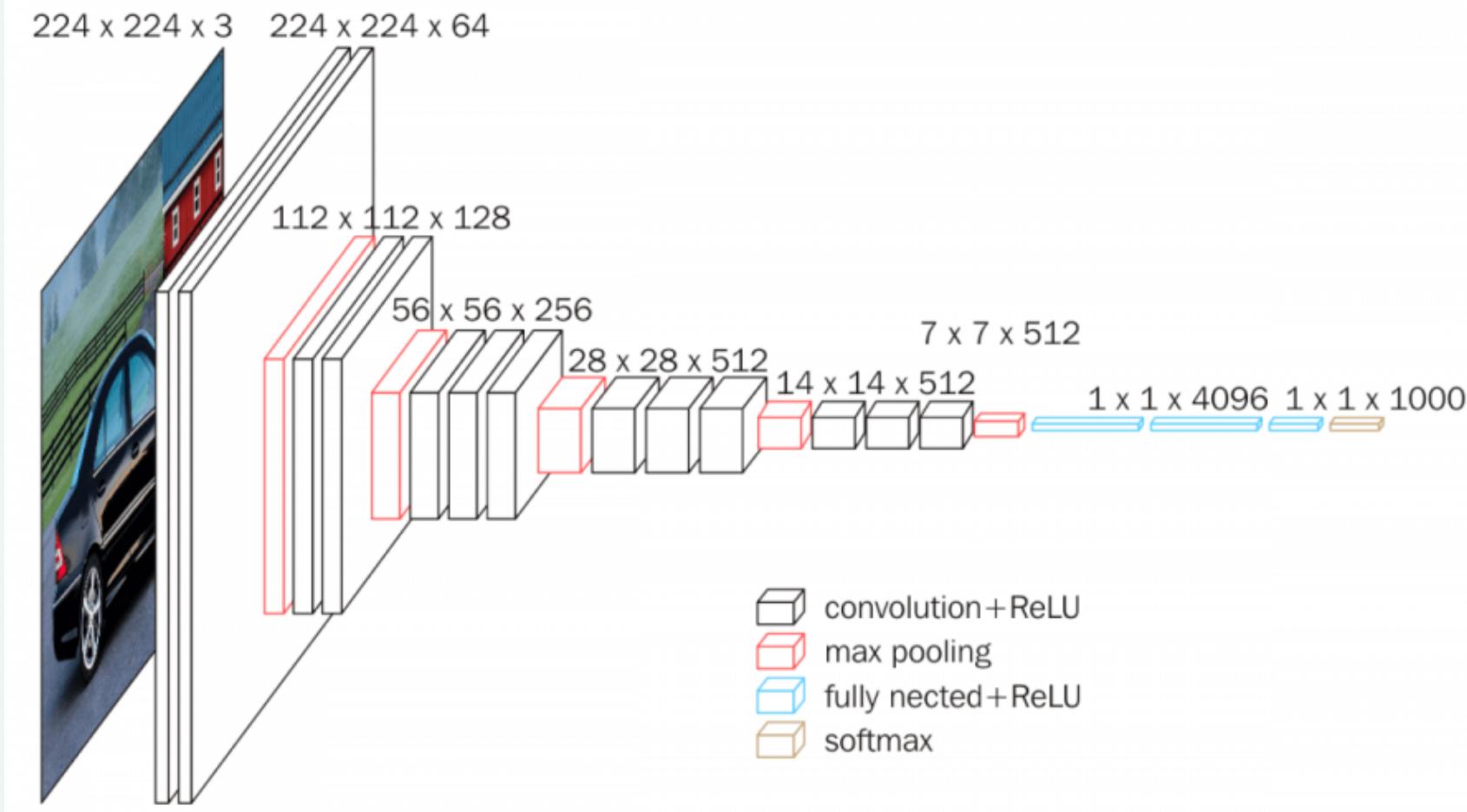
Emmanouil Theodosis  
[etheodosis@g.harvard.edu](mailto:etheodosis@g.harvard.edu)  
[manosth.github.io/](https://manosth.github.io/)  
[github.com/manosth/](https://github.com/manosth/)



# Motivation

## Interpretability

Deep learning is empirical and hard to reason



## Crafting representations

How do we instill domain knowledge?

ReLU vs sigmoid vs tanh vs ....

Analysis is usually post-hoc.

# Optimization-informed deep learning

Linear generative model

$$\mathbf{y} = \mathbf{A}\mathbf{x}$$

Inverse problem

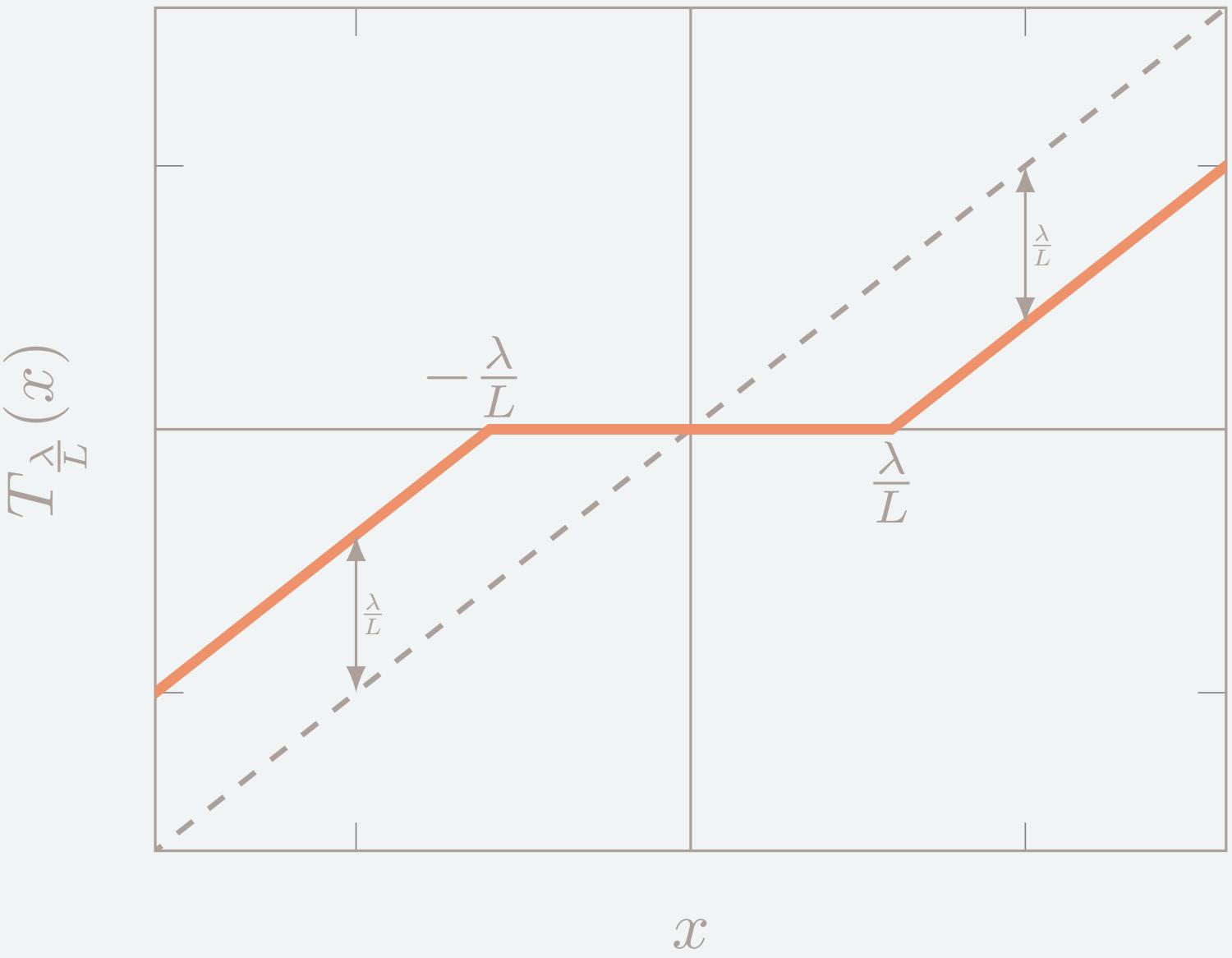
$$\mathbf{x}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$$

In deep learning  $\mathbf{A}$  is overcomplete

$$\mathbf{y} = \mathbf{A}\mathbf{x} \quad \text{s.t. } \mathbf{x} \text{ sparse}$$

Iterative soft thresholding

$$\mathbf{x}^{(l+1)} \leftarrow T_{\frac{\lambda}{L}}(\mathbf{x}^{(l)} + \frac{1}{L} \mathbf{A}^T (\mathbf{y} - \mathbf{A}\mathbf{x}^{(l)}))$$

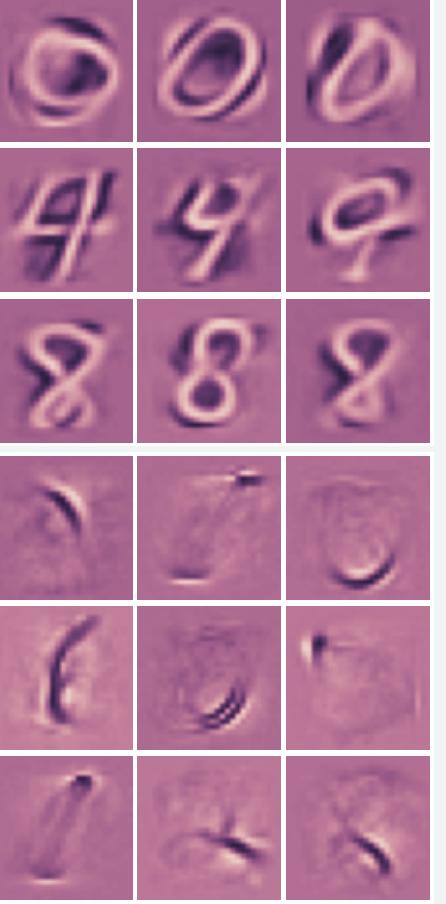
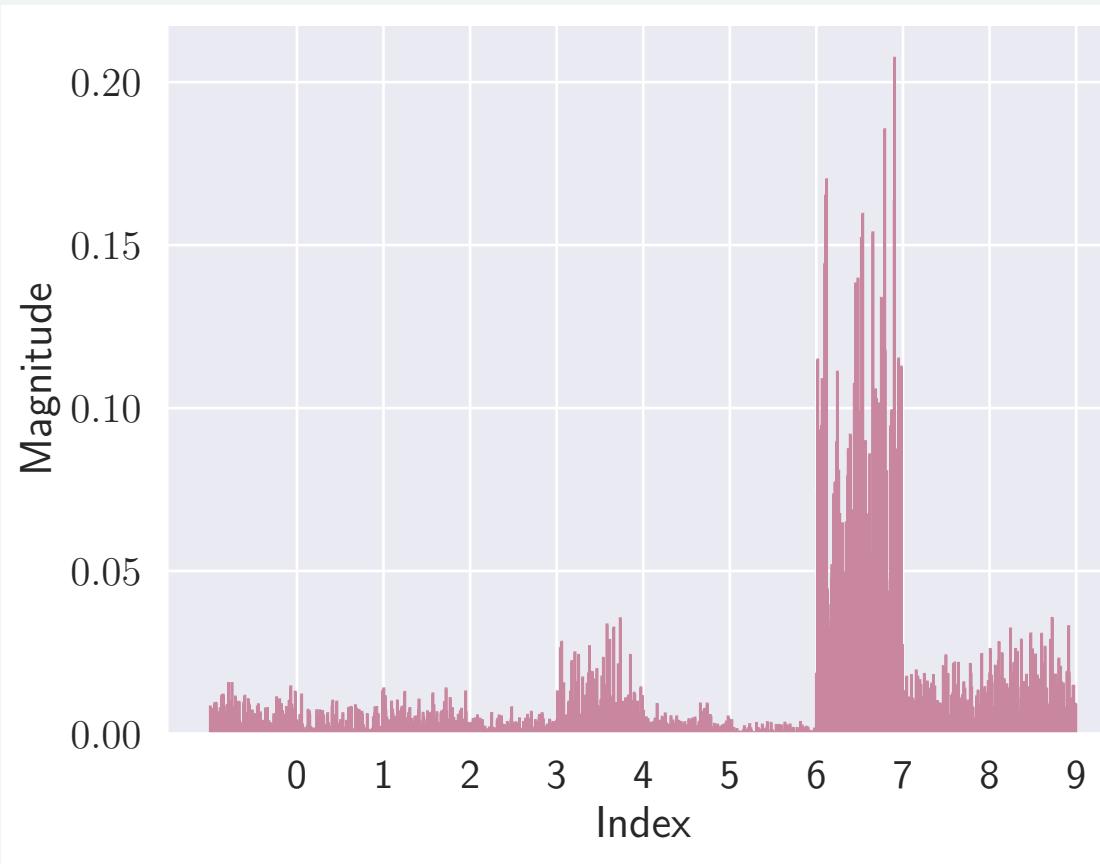


# Group sparse priors

## Model

Assume a group sparse prior

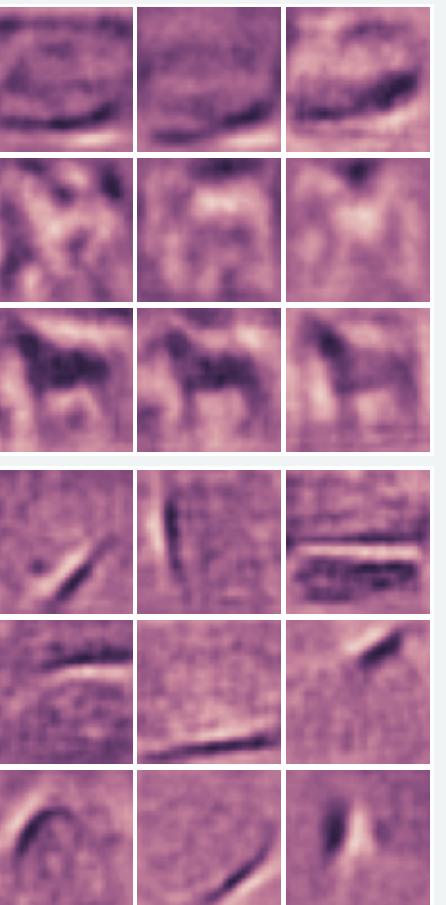
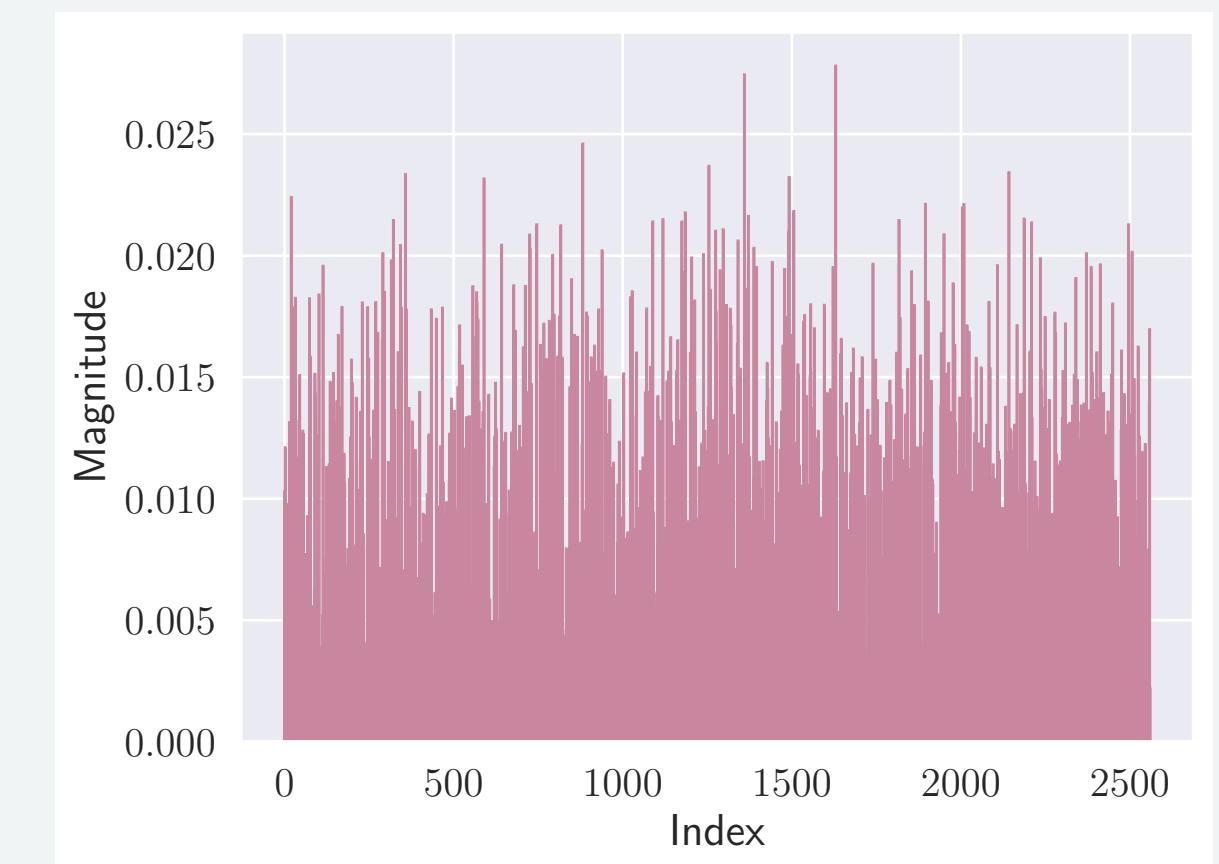
$$\min_{\boldsymbol{x}} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2 + \lambda \sum_{g \in S} \|\boldsymbol{x}_g\|_2$$



## Proximal operator

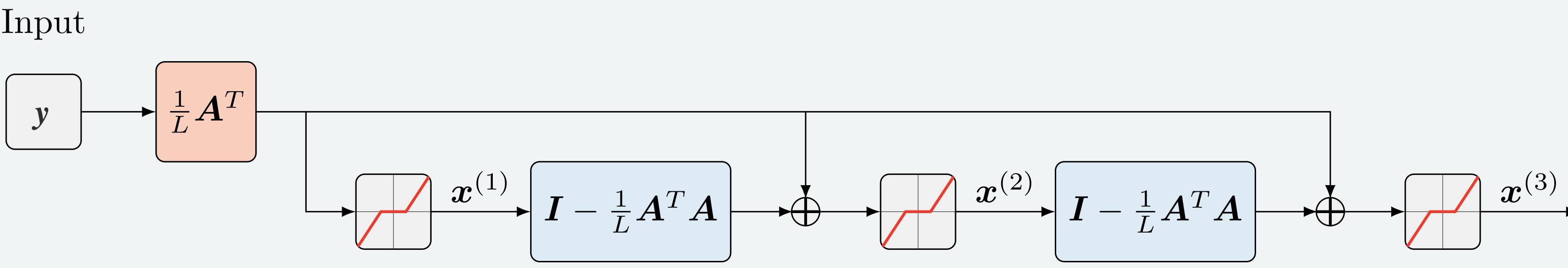
Projection to the solution set via

$$\sigma_\lambda(\boldsymbol{x}_g) = \left( 1 - \frac{\lambda}{\|\boldsymbol{x}_g\|_2} \right)_+ \boldsymbol{x}_g$$



## Architecture

Residual connections to the input

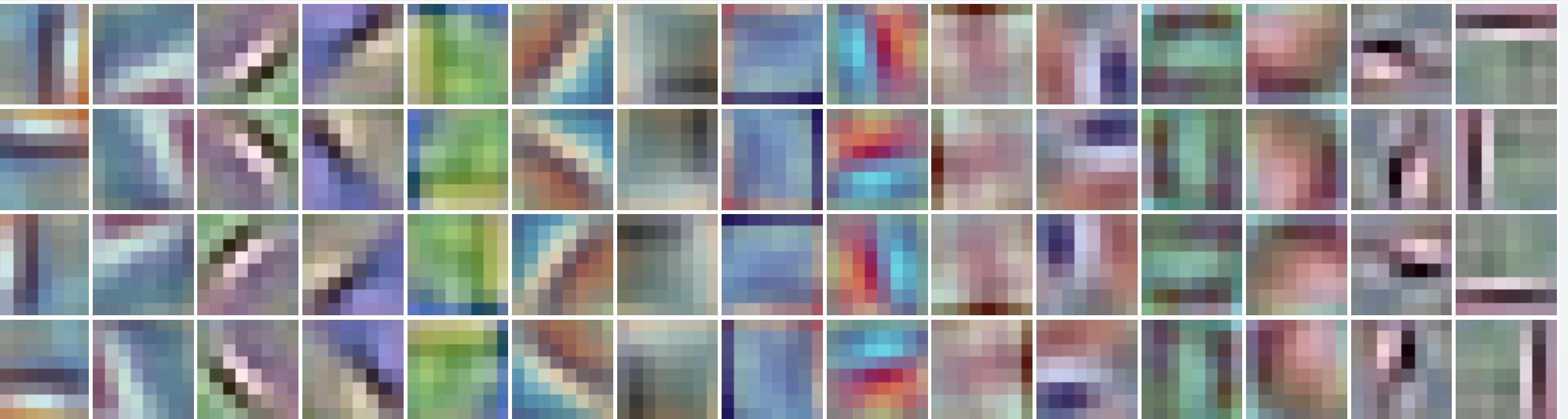


# Unrolling equivariances

## Equivariance

“Consistent” behavior w.r.t. an operator

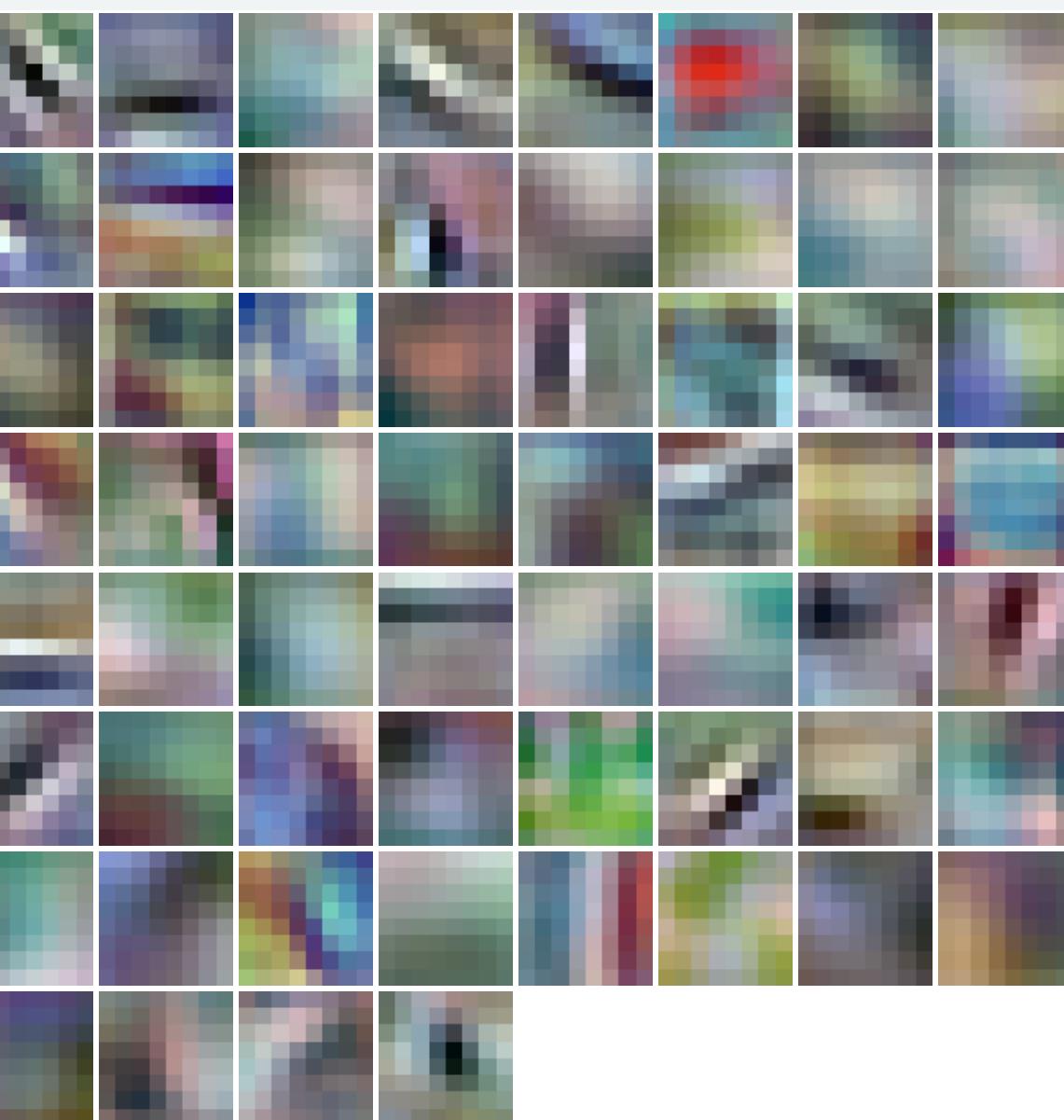
$$f(T(\mathbf{x})) = T'(f(\mathbf{x}))$$



## Symmetry

Model fixed rotations

$$G = \{e, R_\theta, R_\theta^2, \dots, R_\theta^{k-1}\}$$



## Unrolled network

Layer weights

$$\mathbf{W}_l = [\mathbf{w}_l \quad R_\theta(\mathbf{w}_l) \quad \dots R_\theta^{k-1}(\mathbf{w}_l)]$$

# Dense and sparse decomposition

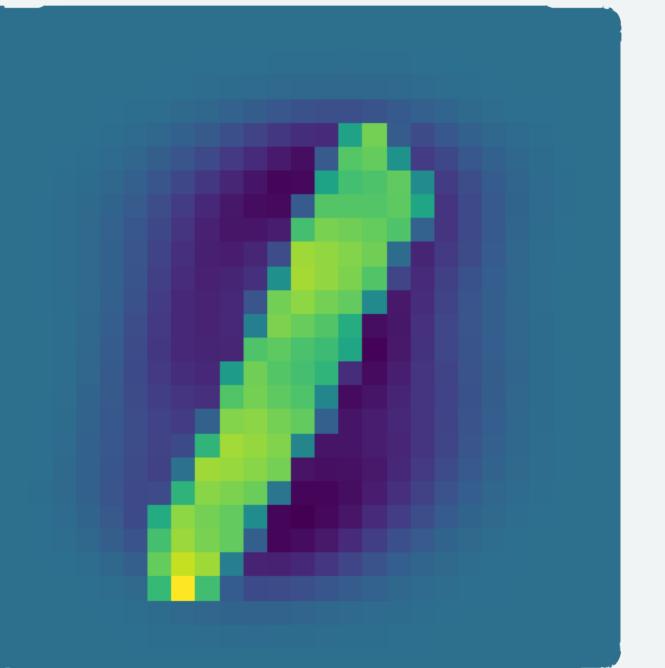
## Model

Decompose observations

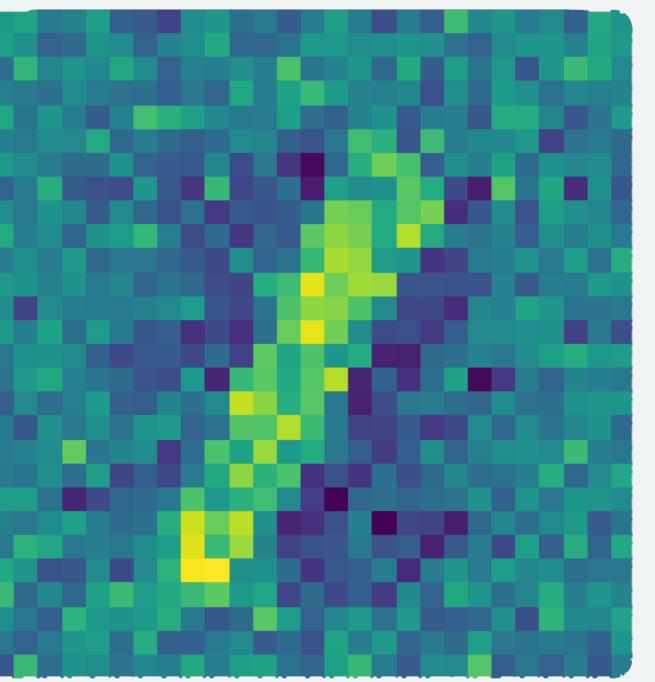
$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$$

where

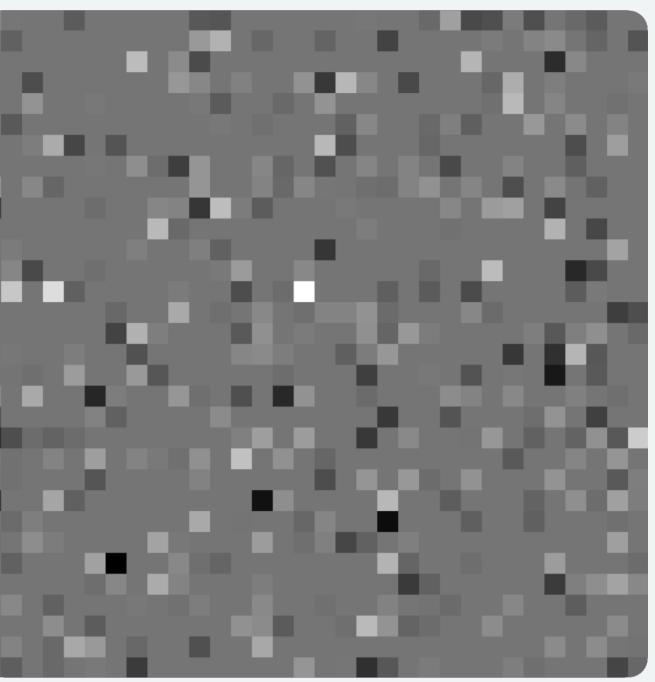
- $\mathbf{u}$  is sparse
- $\mathbf{A}\mathbf{x}$  is dense



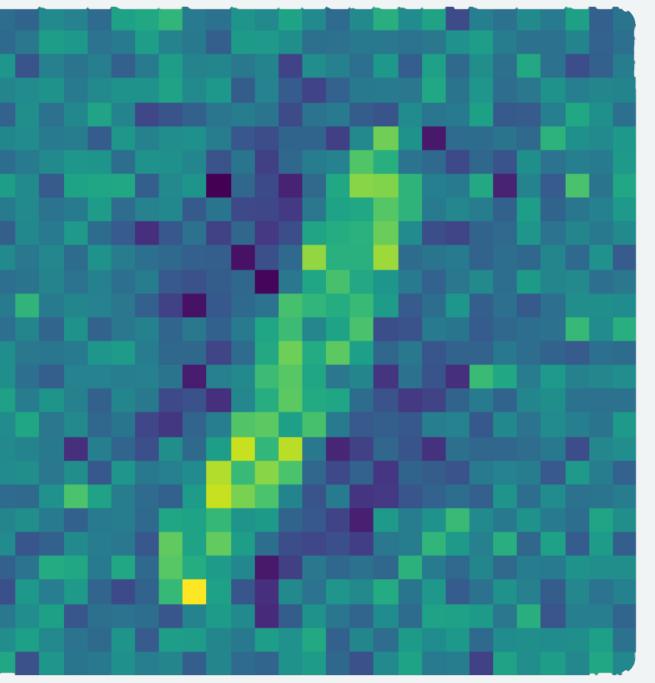
$\mathbf{y}$



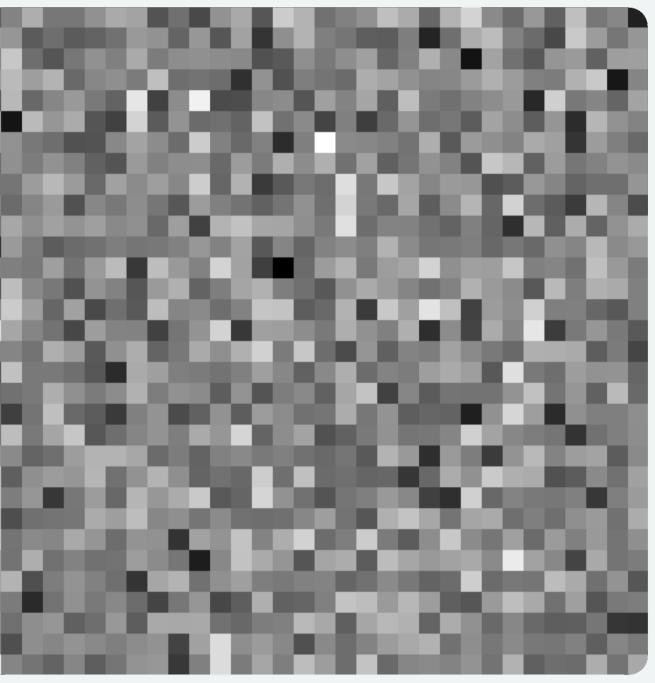
$\mathbf{B}\mathbf{u}$



$\mathbf{u}$



$\mathbf{A}\mathbf{x}$



$\mathbf{x}$

## Optimization

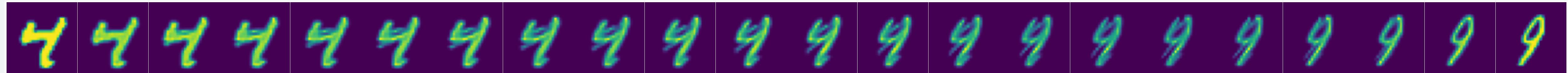
Unrolled network for

$$\min_{\mathbf{x}, \mathbf{u}} \|\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{B}\mathbf{u}\|_2^2 + \mu \|\mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{u}\|_1$$



# Current and future ideas

## Robustness



- Group sparse supports are harder to change
- Robustness to noise (natural or adversarial)

## Generalization

$$\text{sample complexity}_{\text{GS}} < \text{sample complexity}_{\text{S}}$$

- Able to reach similar performance with fewer samples
- Generalization error is smaller

## (Partial) Symmetry learning

$$W_l = [w_l \quad Aw_l \quad \dots A^{K-1}w_l]$$

- Learn residual symmetries in convolutional settings
- General framework for implicit group learning

# Collaborators and publications

## Publications

E. Theodosis and D. Ba, “[Learning silhouettes with group sparse autoencoders](#)”, in *International Conference on Acoustics, Speech, and Signal Processing*, 2023

A. Tasissa, E. Theodosis, B. Tolooshams, and D. Ba, “[Discriminative reconstruction via simultaneous dense and sparse coding](#)”, *Under review*, 2022

E. Theodosis and D. Ba, “[Learning unfolded networks with a cyclic group structure](#)”, in *NeurIPS Workshop on Symmetry and Geometry in Neural Representations*, 2022

E. Theodosis, B. Tolooshams, P. Tankala, A. Tasissa, and D. Ba, “[On the convergence of group sparse autoencoders](#)”, in *arXiv*, 2020



Bahareh Tolooshams  
[btolooshams.github.io/](https://btolooshams.github.io/)



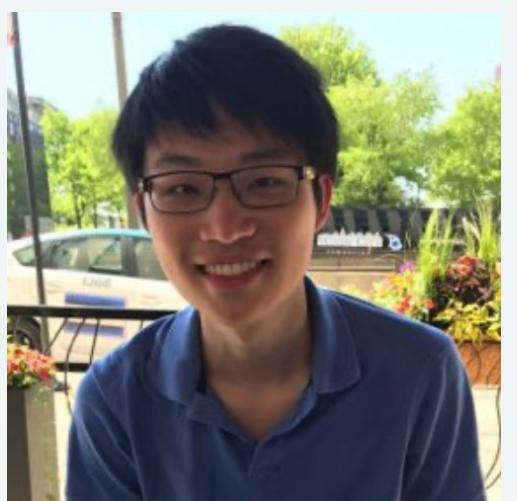
Demba Ba  
[demba-ba.org/](https://demba-ba.org/)



Abiy Tasissa  
[sites.tufts.edu/atasissa/](https://sites.tufts.edu/atasissa/)

## Acknowledgements

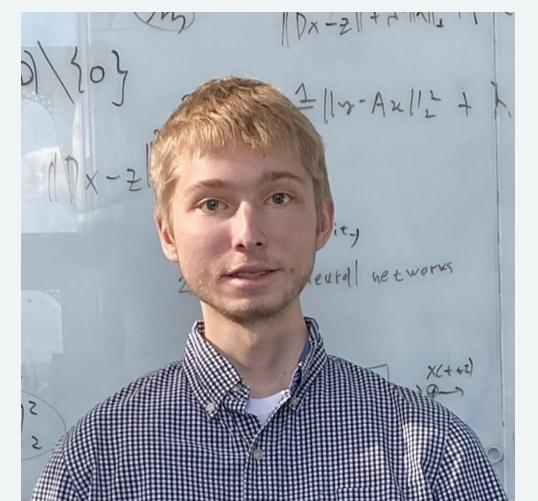
Members of CRISP



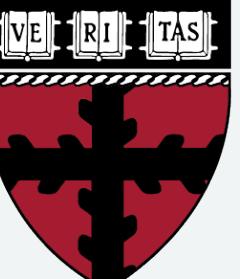
Alex



Shubham



Jon



# THANK YOU

Questions?

