# Sparse Coding, Artificial Neural Networks, and the Brain: Toward "Substantive Intelligence"

Demba Ba[1,2]

Assoc. Prof. of Electrical Engineering and Bioengineering

[1]Harvard University
School of Engineering and Applied Sciences (SEAS)
[2]Institute for AI and Fundamental Interactions (IAIFI)

MIT, April 15[th] 2021

# Outline
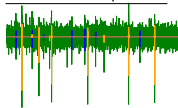
The many faces of sparse coding

Interpretable neural networks for source separation/sparse
dictionary learning

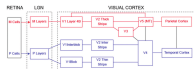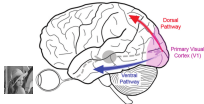Deep sparse coding and hierarchical sensory processing

# Neuroscience, ReLU nets and sparse coding

Unsupervised learning
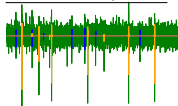to decipher neural code

Blind source separation



Hierarchical sensory
processing principles

# Neuroscience, ReLU nets and sparse coding

## Unsupervised learning to decipher neural code

### Blind source separation



### Hierarchical sensory processing principles



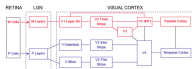## How to design interpretable deep nets?

# Neuroscience, ReLU nets and sparse coding



Unsupervised learning to decipher neural code

Blind source separation

Hierarchical sensory processing principles

How to design interpretable deep nets?

$\log p(\mathbf{y}|\mathbf{Hx}) \propto f(\mathbf{Hx})^{\mathsf{T}}\mathbf{y}$

$\mathbf{x} \sim$ sparse

Encoder

Decoder

Repeat $T$ times

(Deep) Sparse Coding

$\mathbf{y} = \mathbf{A}^{(2)}$ $\mathbf{A}^{(1)}$ $\mathbf{x}$
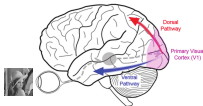
$\mathbf{x} \sim$ sparse

# Neuroscience, ReLU nets and sparse coding
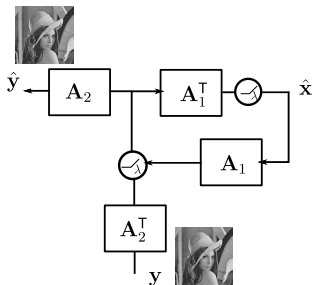


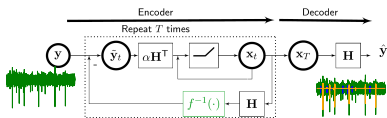Unsupervised learning to decipher neural code

Blind source separation
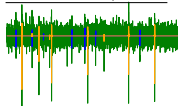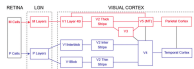
Hierarchical sensory processing principles

How to design interpretable deep nets?

$\log p(\mathbf{y}|\mathbf{H}\mathbf{x}) \propto f(\mathbf{H}\mathbf{x})^{\mathsf{T}}\mathbf{y}$

$\mathbf{x} \sim$ sparse

(Deep) Sparse Coding

$\mathbf{y} = \mathbf{A}^{(2)}$ $\mathbf{A}^{(1)}$ $\mathbf{x}$

$\mathbf{x} \sim$ sparse

Unique opportunity to use deep learning, *in a principled fashion*.

# Outline

# Electrophysiological recordings of neural spiking
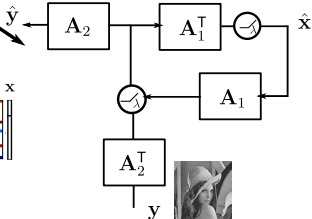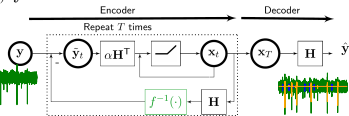
Tetrodes: Four electrodes "listen" to $\approx 10$ neighboring neurons (circa $\approx 2000$)



Data rate

$32$-bit time series

$10$ kHz per second

$\approx 1$ Gb per 6 hours!

Neuropixels: Hundreds of electrodes "listen" to $\approx 100$ neighboring neurons (circa $\approx 2020$)



Data rate

$14$-bit time series

$10$ kHz per second

$\approx 400$ Gb per hour!

Goal: identify neural sources (spike sorting).

# Auto-encoders for spike sorting

### Electrophysiology



### Sparse Coding

$$\mathbf{y} = [\mathbf{H}_1 | \cdots | \mathbf{H}_C] \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_C \end{bmatrix} + \mathbf{v} = \mathbf{H}\mathbf{x} + \mathbf{v}$$

$$\mathbf{x} \sim \text{sparse}$$

# Auto-encoders for spike sorting

# Auto-encoders for spike sorting

# Blind source separation by dictionary learning

Sparse Codes $x_c[n]$

Filters/Dictionary elements $h_c[n]$

# Convolutional generative model

The generative model for shift-invariant sparse representation:

$$y_n = \sum_{c=1}^{C} h_c[n] * x_c[n] + v_n$$

Sources                          Data

# Convolutional generative model

The generative model for shift-invariant sparse representation:

$$y_n = \sum_{c=1}^{C} h_c[n] * x_c[n] + v_n$$

(Linear-algebraic form) $\quad \mathbf{y} = \begin{bmatrix} \mathbf{H}_1 | \cdots | \mathbf{H}_C \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_C \end{bmatrix} + \mathbf{v} = \mathbf{H}\mathbf{x} + \mathbf{v}$

Given only the data $\mathbf{y}$, how to solve for $\mathbf{H}$ and $\mathbf{x}$?

# Spatio-temporal generalizations of convolutional model

$$y_{n,k_1,k_2} = \sum_{c=1}^{C} \sum_{i=1}^{N_c} x_c[n,k_1,k_2] * h_c[n] g_c[k_1,k_2] + v_{n,k_1,k_2}$$

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{v}$$

# Pervasiveness of $\mathbf{y} = \mathbf{Hx} + \mathbf{v}$ model

$$\mathbf{y} = \mathbf{Hx} + \mathbf{v}$$

| Discipline | y | H | x | Conv? |
|---|---|---|---|---|
| Neuro | Electrode array | Neurons | event times/amplitudes | Y |
| Acoustics | Mic. array | Sound source | event times/amplitudes | Y |
| Astro | Tel. array | Celest. objects | location/intensity | Y |
| Part. physics | Detector observ. (mixture) | Particle "topics" | Proportions | N |
| Optics | Intensity | Scattering matrix | Object | N/Y |
| Genomics | Cell gene expr. | Gene expr. modules | Gene expr. levels | N |
| Text | Word hist. | Topic hist. | Topic proportions | N |
| Radar | Antenna array | Sources | Locations | Y |
| ... | | | | |

# Optimization perspective

Given $J$ examples of data $\{\mathbf{y}^j\}_{j=1}^J$,

CDL solves:

$$\min_{(\mathbf{x}^j)_{j=1}^J, (\mathbf{h}_c)_{c=1}^C} \sum_{j=1}^J \frac{1}{2} \left|\left|\mathbf{y}^j - \mathbf{H}\mathbf{x}^j\right|\right|_2^2 + \lambda \left|\left|\mathbf{x}^j\right|\right|_1$$

$$\text{s.t. } ||\mathbf{h}_c||_2 \leq 1 \quad \text{for } c = 1, \cdots, C,$$

where $\lambda > 0$ (regularization parameter enforcing sparsity).

# Convolutional sparse coding step

Given the filters,

CSC is separable over $J$ examples

$$\min_{\mathbf{x}^j} \frac{1}{2} \left|\left|\mathbf{y}^j - \mathbf{H}\mathbf{x}^j\right|\right|_2^2 + \lambda \left|\left|\mathbf{x}^j\right|\right|_1$$

This step is

- Embarrassingly parallelizable.
- Amenable to GPU processing (long recordings).

# ReLU nonlinearity from sparse approximation in 1D

$$\begin{aligned}
\hat{x} &= \underset{x \in \mathbb{R}^+}{\arg\min} \; \frac{1}{2}(y-x)^2 + \lambda|x| \\
&= \mathsf{max}(y - \lambda, 0) \\
&= \mathsf{ReLU}(y - \lambda).
\end{aligned}$$



$\mathsf{ReLU}(y - \lambda)$

# Neural net with ReLU from sparse approximation

<u>Case 1</u>: $\mathbf{H}^\mathsf{T}\mathbf{H} = \mathbf{I}$

$$
\begin{aligned}
\hat{\mathbf{x}} &= \operatorname*{arg\,min}_{\mathbf{x}\in\mathbb{R}^r} \frac{1}{2}\|\mathbf{y}-\mathbf{Hx}\|_2^2 + \lambda\|\mathbf{x}\|_1 \\
&= \mathsf{ReLU}\left(\mathbf{H}^\mathsf{T}\mathbf{y} - \lambda\right).
\end{aligned}
$$



<u>Case 2</u>: $\mathbf{H}^\mathsf{T}\mathbf{H} \neq \mathbf{I}$, $\sigma_{\mathsf{max}}(\mathbf{H}^\mathsf{T}\mathbf{H}) \leq L$

$$
\mathbf{x}_t = \mathsf{ReLU}\left(\mathbf{x}_{t-1} + \alpha\mathbf{H}^\mathsf{T}(\mathbf{y}-\mathbf{Hx}_{t-1}) - \frac{\lambda}{L}\right)
$$

## Other nonlinearities?

$$\log p_X(x) \propto$$
$$-\frac{x^2}{2\sigma_-^2}\mathbf{1}_{\{x\leq 0\}} - \frac{x^2}{2\sigma_+^2}\mathbf{1}_{\{x\geq 0\}}$$



$$\begin{aligned}
\hat{x} &= \underset{x\in\mathbb{R}}{\arg\min}\,\frac{1}{2}(y-x)^2 - \log p_X(x) = \\
&= \frac{\sigma_-^2}{1+\sigma_-^2}y\mathbf{1}_{\{y\leq 0\}} + \frac{\sigma_+^2}{1+\sigma_+^2}y\mathbf{1}_{\{y\geq 0\}} \\
&= \mathsf{LeakyReLu}\,(y;\sigma_+,\sigma_-)\,.
\end{aligned}$$



LeakyReLu$(y;\sigma_+,\sigma_-)$

# Convolutional dictionary update Step

$$\min_{(\mathbf{h}_c)_{c=1}^{C}} \sum_{j=1}^{J} \frac{1}{2} \left\| \mathbf{y}^j - \mathbf{H}\mathbf{x}^j \right\|_2^2 \text{ s.t. } \|\mathbf{h}_c\|_2 \leq 1 \quad \text{for } c = 1, \cdots, C.$$

This step is

- Computationally Expensive.
- Not parallelizable over $J$ examples.

# Auto-encoder for CDL by deep unfolding and weight tying

1. <u>Nonlinear encoder</u>:
   $\mathbf{x}_t = \mathsf{ReLU}\left(\mathbf{x}_{t-1} + \alpha\mathbf{H}^\mathsf{T}(\mathbf{y} - \mathbf{H}\mathbf{x}_{t-1}) - \frac{\lambda}{L}\right).$
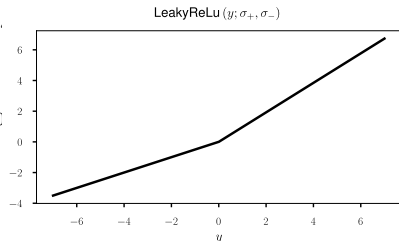2. <u>Linear decoder</u>: $\hat{\mathbf{y}} = \mathbf{H}\mathbf{x}_T.$
3. <u>Training</u>: backprop with MSE loss.

# Simulated data



Sources          Data

Data: 17 minutes of electrical activity from 3 neurons!
Sampling rate: $f_s = 10$ kHz.
Firing rate: 30 Hz.

# AE performs dictionary learning

$$\text{err}(\mathbf{h}_c, \hat{\mathbf{h}}_c) = \sqrt{1 - \frac{\langle \mathbf{h}_c, \hat{\mathbf{h}}_c \rangle^2}{\|\mathbf{h}_c\|_2^2 \|\hat{\mathbf{h}}_c\|_2^2}}$$

# Training of AE is fast

|                      |            | CRsAE      | Sporco     | CBP      |
|----------------------|------------|------------|------------|----------|
| Learning Spike Shapes | runtime    | **69.27 s** | $319.52$ s |          |
|                      | iterations | **10**     | $89$       |          |
| Spike Sorting        | runtime    | **0.93 s** |            | 17 hours |

# Harris dataset: spike sorting and denoising

# Unfolded auto-encoders for exponential family data



Components for different distributions
$(\mathcal{S}_b(\cdot) = \mathsf{ReLU}(\cdot - b), \ b \text{ depends on } \lambda)$

|          | $\mathbf{y}$ | $f^{-1}(\cdot)$ | Encoder Unfolding $(\mathbf{x}_t)$ | Decoder $(f(\hat{\boldsymbol{\mu}}))$ |
|----------|:------------:|:---------------:|:-----------------------------------:|:-------------------------------------:|
| Gaussian | $\mathbb{R}$ | $I(\cdot)$ | $\mathcal{S}_b\left(\mathbf{x}_{t-1} + \alpha\mathbf{H}^{\mathsf{T}}\widetilde{\mathbf{y}}_t\right)$ | $\mathbf{H}\mathbf{x}_T$ |
| Binomial | $[0..M]$ | $\mathsf{sigmoid}(\cdot)$ | $\mathcal{S}_b\left(\mathbf{x}_{t-1} + \alpha\mathbf{H}^{\mathsf{T}}(\frac{1}{M}\widetilde{\mathbf{y}}_t)\right)$ | $\mathbf{H}\mathbf{x}_T$ |
| Poisson  | $[0..\infty)$ | $\exp(\cdot)$ | $\mathcal{S}_b\left(\mathbf{x}_{t-1} + \alpha\mathbf{H}^{\mathsf{T}}\left(\mathsf{Elu}(\widetilde{\mathbf{y}}_t)\right)\right)$ | $\mathbf{H}\mathbf{x}_T$ |

# Auto-encoders for denoising Poisson images

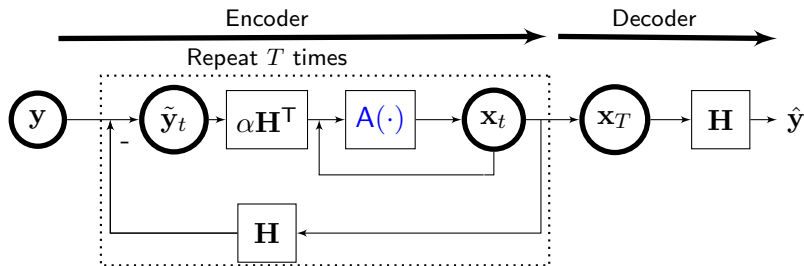Low photon count image



Denoised image



Poisson Convolutional Model

$$\mathbf{y} \sim \text{Poisson}(\boldsymbol{\lambda})$$

$$\log \boldsymbol{\lambda} \;=\; [\mathbf{H}_1 | \cdots | \mathbf{H}_C] \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_C \end{bmatrix} = \mathbf{Hx}$$

$$\mathbf{x} \sim \text{sparse}$$

Expl'ble AI

Fast

Neural Network



|  |  | SPDA | CA | DCEA-C | DCEA-UC |
|---|---|---|---|---|---|
| Peak 1 | Set12 | 20.39 | **21.51** | 20.72 | 21.37 |
|  | BSD68 | . | 21.78 | 21.27 | **21.84** |
| Peak 2 | Set12 | 21.70 | **22.97** | 22.02 | 22.79 |
|  | BSD68 | . | 22.90 | 22.31 | **22.92** |
| Peak 4 | Set12 | 22.56 | **24.40** | 23.51 | 24.37 |
|  | BSD68 | . | 23.98 | 23.54 | **24.10** |
| # of Params |  | 160K | 655K | 20K | 61K |

Tolooshams et al., ICML 2020

# Unfolded auto-encoders with different priors

Encoder — Repeat $T$ times — Decoder

$\mathbf{y} \rightarrow \tilde{\mathbf{y}}_t \rightarrow \alpha\mathbf{H}^{\mathsf{T}} \rightarrow A(\cdot) \rightarrow \mathbf{x}_t \rightarrow \mathbf{x}_T \rightarrow \mathbf{H} \rightarrow \hat{\mathbf{y}}$
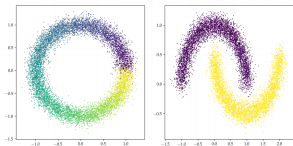
Activation functions for different priors, for (input $\tilde{\mathbf{x}}_t = \mathbf{x}_{t-1} + \alpha\mathbf{H}^{\mathsf{T}}\tilde{\mathbf{y}}_t$ )
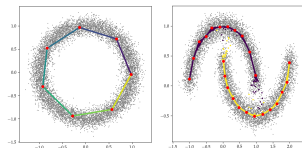
| Prior | Constraints | Expression | Activation |
|---|---|---|---|
| Gaussian/Ridge | None | $\frac{1}{\sigma^2}\sum x_j^2$ | $\frac{\sigma^2}{1/L+\sigma^2}\tilde{\mathbf{x}}_t = c(\sigma)\tilde{\mathbf{x}}_t$ |
| Sparsity | $\mathbf{x} \geq 0$ | $\lambda \sum x_j$ | $\mathrm{ReLU}(\tilde{\mathbf{x}}_t - \frac{\lambda}{L})$ |
| 2-sided Gaussian | None | $\sum \frac{x_j^2}{\sigma_+^2}\mathbf{1}_{\{x_j\geq 0\}} + \frac{x_j^2}{\sigma_-^2}\mathbf{1}_{\{x_j\leq 0\}}$ | $c(\sigma_+)\tilde{\mathbf{x}}_t\mathbf{1}_{\{\tilde{\mathbf{x}}_t\geq 0\}} + c(\sigma_-)\tilde{\mathbf{x}}_t\mathbf{1}_{\{\tilde{\mathbf{x}}_t\leq 0\}}$ |
| Group sparsity | $\mathbf{x} \geq 0$ | $\sum \|x_g\|_2$ | $\left(\tilde{\mathbf{x}}_{t,g} - \lambda\frac{\tilde{\mathbf{x}}_{t,g}}{\|\tilde{\mathbf{x}}_{t,g}\|_2}\right)\mathbf{1}_{\{\|\tilde{\mathbf{x}}_{t,g}\|_2\geq\lambda\}}$ |
| Simplex | $\mathbf{x} \geq 0, \mathbf{x}^{\mathsf{T}}\mathbf{1} = 1$ | None | $\mathrm{ReLU}(\tilde{\mathbf{x}}_t - b(\tilde{\mathbf{x}}_t))$ |
| ... | | | |

# Auto-encoders for manifold learning

### Nonlinear manifold learning



### Sparse Coding on *simplices*
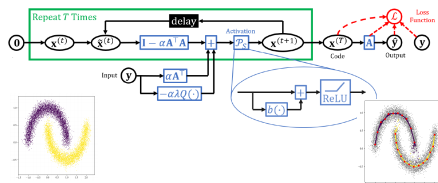
$$\min_{\mathbf{x}\in\mathbb{R}^{m\times n}} \mathcal{L}(\mathbf{A},\mathbf{y},\mathbf{x}) = \frac{1}{2}\|\mathbf{y}-\mathbf{Ax}\|^2 + \lambda\sum_{j=1}^{m} x_j\|\mathbf{y}-\mathbf{a}_j\|^2$$

$$\mathbf{x}^\top\mathbf{1} = \mathbf{1},$$
$$x_j \geq 0, \quad \text{for all } j.$$

Expl'ble AI

Fast

### Piecewise-Linear Manifold *Approximation*



### Neural Network

# Outline

# Deep sparse coding, hierarchical "virtual" brains

Natural images/Hierar. rep.



Deep Sparse Coding
Generative Model

$\mathbf{y} = \quad \mathbf{A}^{(2)} \quad \mathbf{A}^{(1)} \quad \mathbf{x}$



$\mathbf{x} \sim$ sparse
$\mathbf{x}_2 = \mathbf{A}^{(1)} \mathbf{x} \sim$ sparse

# Deep sparse coding, hierarchical "virtual" brains



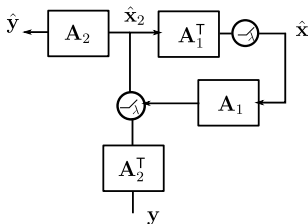Natural images/Hierar. rep.

Deep Sparse Coding
Generative Model

$\mathbf{y} = \mathbf{A}^{(2)} \quad \mathbf{A}^{(1)} \quad \mathbf{x}$

$\mathbf{x} \sim$ sparse
$\mathbf{x}_2 = \mathbf{A}^{(1)}\mathbf{x} \sim$ sparse
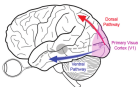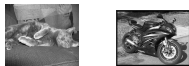
Expl'ble AI

Fast

"virtual" Two-layer Brain
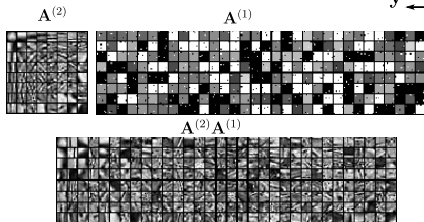
Ba, IEEE TSP 2020

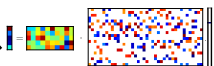# Deep sparse coding, hierarchical "virtual" brains

Natural images/Hierar. rep.



Deep Sparse Coding
Generative Model

$$\mathbf{y} = \quad \mathbf{A}^{(2)} \qquad \mathbf{A}^{(1)} \qquad \mathbf{x}$$



$$\mathbf{x} \sim \text{sparse}$$
$$\mathbf{x}_2 = \mathbf{A}^{(1)}\mathbf{x} \sim \text{sparse}$$

Expl'ble AI

Fast

"virtual" Two-layer Brain

Learned filters

$\mathbf{A}^{(2)}$ $\qquad\qquad \mathbf{A}^{(1)}$



$$\mathbf{A}^{(2)}\mathbf{A}^{(1)}$$



$\hat{\mathbf{y}} \leftarrow \boxed{\mathbf{A}_2} \leftarrow \hat{\mathbf{x}}_2 \rightarrow \boxed{\mathbf{A}_1^\mathsf{T}} \leftarrow \hat{\mathbf{x}}$

$\boxed{\mathbf{A}_1}$

$\boxed{\mathbf{A}_2^\mathsf{T}}$

$\mathbf{y}$

Ba, IEEE TSP 2020

# Deep sparse coding theory guides architecture design

Deeply-sparse coding model



$\mathbf{y} = \quad \mathbf{A}^{(2)} \quad\quad \mathbf{A}^{(1)} \quad\quad \mathbf{x}$

# Deep sparse coding theory guides architecture design

Deeply-sparse coding model

Associated AE

$\mathbf{y} = \quad \mathbf{A}^{(2)} \qquad \mathbf{A}^{(1)} \qquad \mathbf{x}$



1. Pick $r_2$, # units in 1st hidden layer.

$\mathbf{A}^{(2)\mathsf{T}} \sim$ Dense
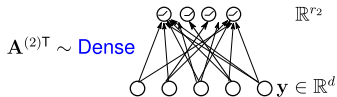
$\mathbb{R}^{r_2}$

$\mathbf{y} \in \mathbb{R}^d$

# Deep sparse coding theory guides architecture design



Deeply-sparse coding model

$\mathbf{y} = \quad \mathbf{A}^{(2)} \qquad \mathbf{A}^{(1)} \qquad \mathbf{x}$

1. Pick $r_2$, # units in 1st hidden layer.

2. $r_1 = \mathcal{O}\left(\max(r_2^2, r_2 s_{\mathbf{A}^{(1)}})\right)$, i.e. expansion.

Associated AE

$\mathbf{A}^{(1)\mathsf{T}} \sim$ Sparse

$\mathbf{A}^{(2)\mathsf{T}} \sim$ Dense

$\mathbb{R}^{r_1}$

$\mathbb{R}^{r_2}$

$\mathbf{y} \in \mathbb{R}^d$

# Deep sparse coding theory guides architecture design



Deeply-sparse coding model

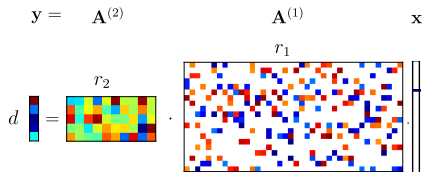$\mathbf{y} = \quad \mathbf{A}^{(2)} \quad \mathbf{A}^{(1)} \quad \mathbf{x}$
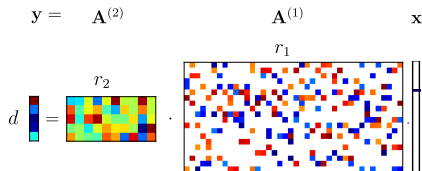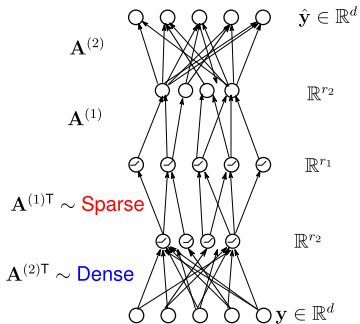
$d \;\; = \quad \overbrace{\quad}^{r_2} \quad \cdot \quad \overbrace{\quad}^{r_1} \quad \cdot$

1. Pick $r_2$, # units in 1st hidden layer.

2. $r_1 = \mathcal{O}\left(\max(r_2^2, r_2 s_{\mathbf{A}^{(1)}})\right)$, i.e. expansion.

3. Reconstruct.

Associated AE

$\hat{\mathbf{y}} \in \mathbb{R}^d$

$\mathbf{A}^{(2)}$

$\mathbb{R}^{r_2}$

$\mathbf{A}^{(1)}$

$\mathbb{R}^{r_1}$

$\mathbf{A}^{(1)\mathsf{T}} \sim$ Sparse

$\mathbb{R}^{r_2}$

$\mathbf{A}^{(2)\mathsf{T}} \sim$ Dense

$\mathbf{y} \in \mathbb{R}^d$

Theorem (Ba 2020): need $n = \mathcal{O}\left(\max(r_1^2, r_1 s_{\mathbf{X}})\right)$ examples to learn $\mathbf{A}^{(1)}$ and $\mathbf{A}^{(2)}$.

# Predictions of deep sparse coding theory



- ▶ Predictions for hierarchical sensory processing
  - ▶ Sparse activation of neurons at every level
  - ▶ *Increasingly sparse* activations as one goes up the hierarchy: Barlow 1961, (cat neuron, dog neuron?).
  - ▶ Sparse receptive fields at higher levels of hierarchy.
- ▶ Can we design experiments to test/verify/modify?

# Implications/Concluding thoughts

1. A principled way to design neural networks. Choices of
   - Nonlinearities (e.g. ReLU) $\iff$ *prior* on hidden activations.
   - Training loss $\iff$ assumptions on data *likelihood*.
2. Fast learning and inference using GPUs.
3. Interpretable architectures; orders of mag. fewer params. than standard neural nets.
4. A call for utilizing neural networks in a principled way to
   4.1 Solve pattern discovery problems in neuroscience.
   4.2 Elucidate principles of hierarchical sensory processing, in conjunction with experiments.

# Thank you

Demba Ba
*demba@seas.harvard.edu*
https://crip.seas.harvard.edu/

https://github.com/demba/
https://bitbucket.org/demba/

Tankala, Tassissa, Murphy and **Ba**. "K-deep simplex: deep manifold learning via local dictionaries". Submitted to ICML 2021.

**Ba**. "Deeply-sparse signal representations". IEEE Transactions on Signal Processing, 2020.

Tolooshams, Dey and **Ba**. "Constrained recurrent sparse auto-encoders: Expectation-maximization based architectures for convolutional dictionary learning". IEEE Transcations on Neural Networks, 2020.

Tolooshams, Song, Temereanca and **Ba**. "Convolutional dictionary learning based auto-encoders for natural exponential-family distributions". ICML 2020.

Zazo, Tolooshams and **Ba**. "Convolutional dictionary learning in hierarchical networks". in Proceedings of IEEE CAMSAP 2019.