

The topology of representation teleportation, regularized Oja's rule, and symmetric weights

Jon Bloom

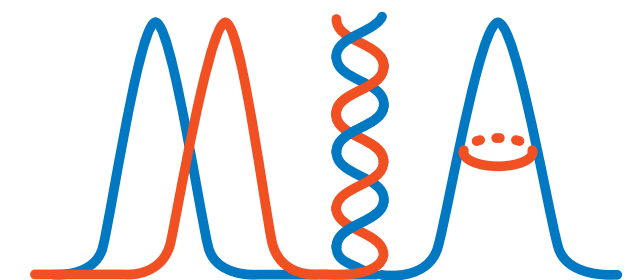
Institute Scientist, Hail Engineer

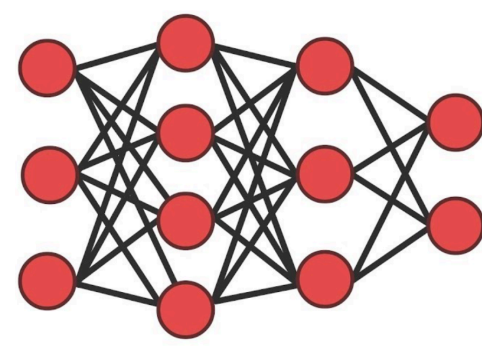
Director of Models, Inference, and Algorithms

Co-conspirators: Daniel Kunin, Aleks Goeva, Cotton Seed



MASSACHUSETTS
GENERAL HOSPITAL





Prediction in artificial neural networks is inspired by the brain.
 Is *learning* in the brain inspired by artificial neural networks?



COGNITIVE SCIENCE **11**, 23–63 (1987)

Competitive Learning: From Interactive Activation to Adaptive Resonance

STEPHEN GROSSBERG
Boston University

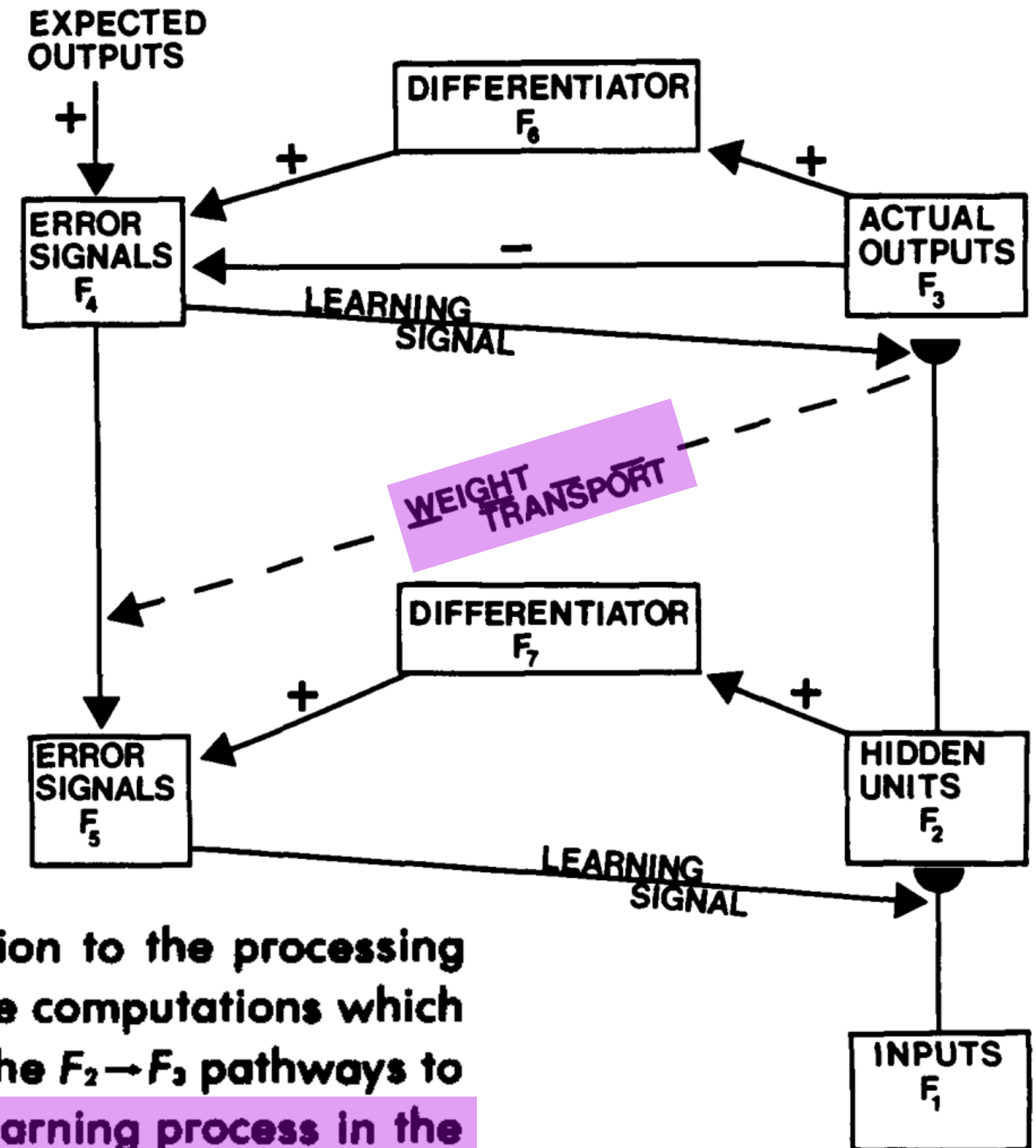
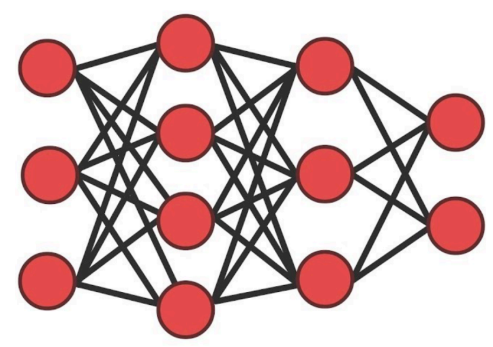


Figure 8. Circuit diagram of the back propagation model: In addition to the processing levels F_1 , F_2 , F_3 , there are also levels F_4 , F_5 , F_6 , and F_7 to carry out the computations which control the learning process. The transport of learned weights from the $F_2 \rightarrow F_3$ pathways to the $F_4 \rightarrow F_5$ pathways shows that **this algorithm cannot represent a learning process in the brain.**

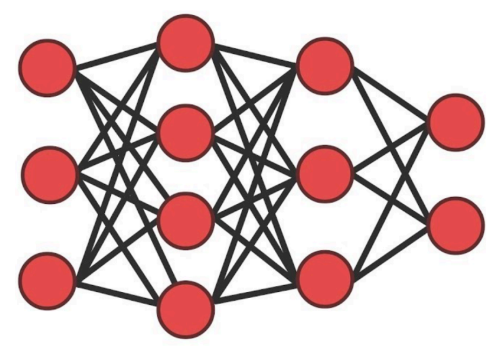


Prediction in artificial neural networks is inspired by the brain.



Is *learning* in the brain inspired by artificial neural networks?

Such a physical transport of weights has no plausible physical interpretation. The weights in the $F_2 \rightarrow F_3$ pathways must be computed *within* these pathways in order to multiply signals from F_2 to F_3 . These weights cannot also exist *within* the pathways from F_4 to F_5 in order to multiply signals from F_4 to F_5 without being physically transported from $(F_2 \rightarrow F_3)$ to $(F_4 \rightarrow F_5)$ pathways, thereby violating basic properties of locality. Moreover, the levels F_3 and F_4 cannot be lumped together, because F_3 must record actual outputs, whereas F_4 must record differences between expected and actual outputs. The *BP* model is thus not a model of a brain process.

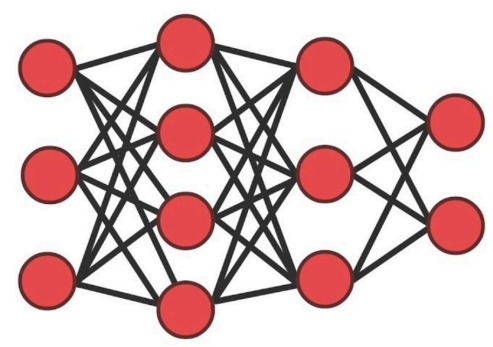


Prediction in artificial neural networks is inspired by the brain.
Is *learning* in the brain inspired by artificial neural networks?

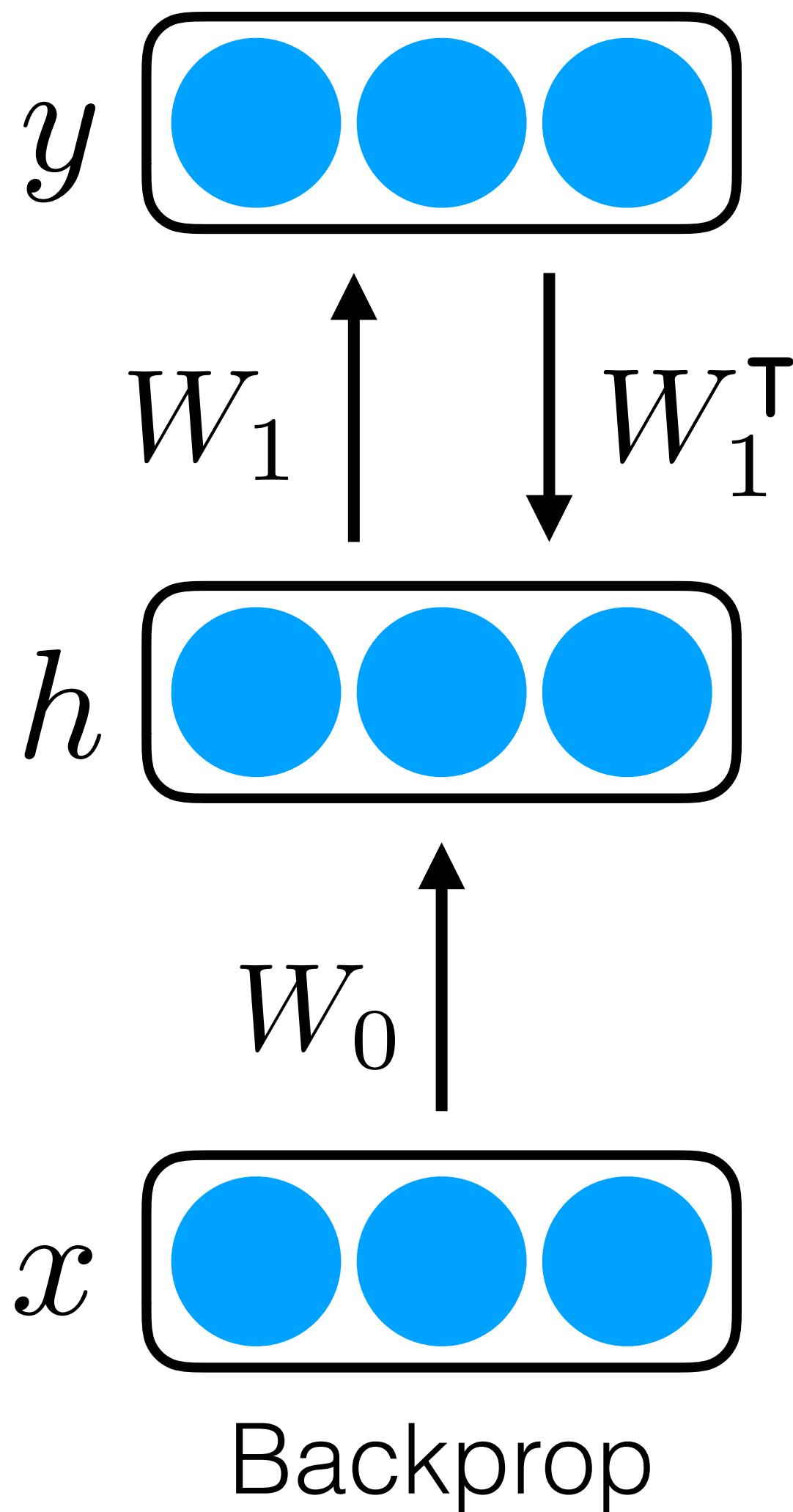


Nope.

I hope you enjoyed my talk!



Prediction in artificial neural networks is inspired by the brain.
Is *learning* in the brain inspired by artificial neural networks?



$$e = y - W_1 W_0 x$$

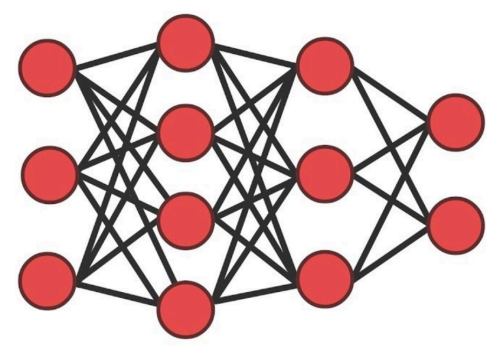
$$\Delta W_1 \propto e h^T$$

$$\Delta W_0 \propto W_1^T e x^T$$

weight
transport
problem

symmetric
weights
problem

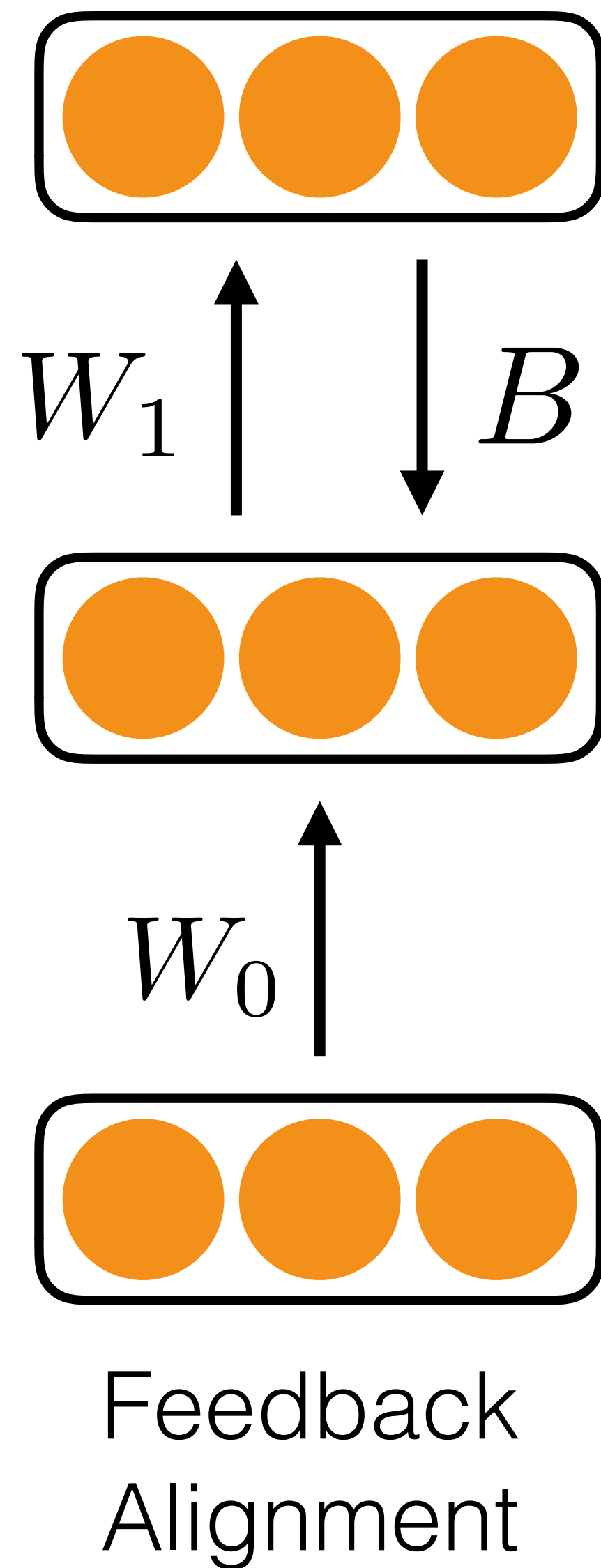
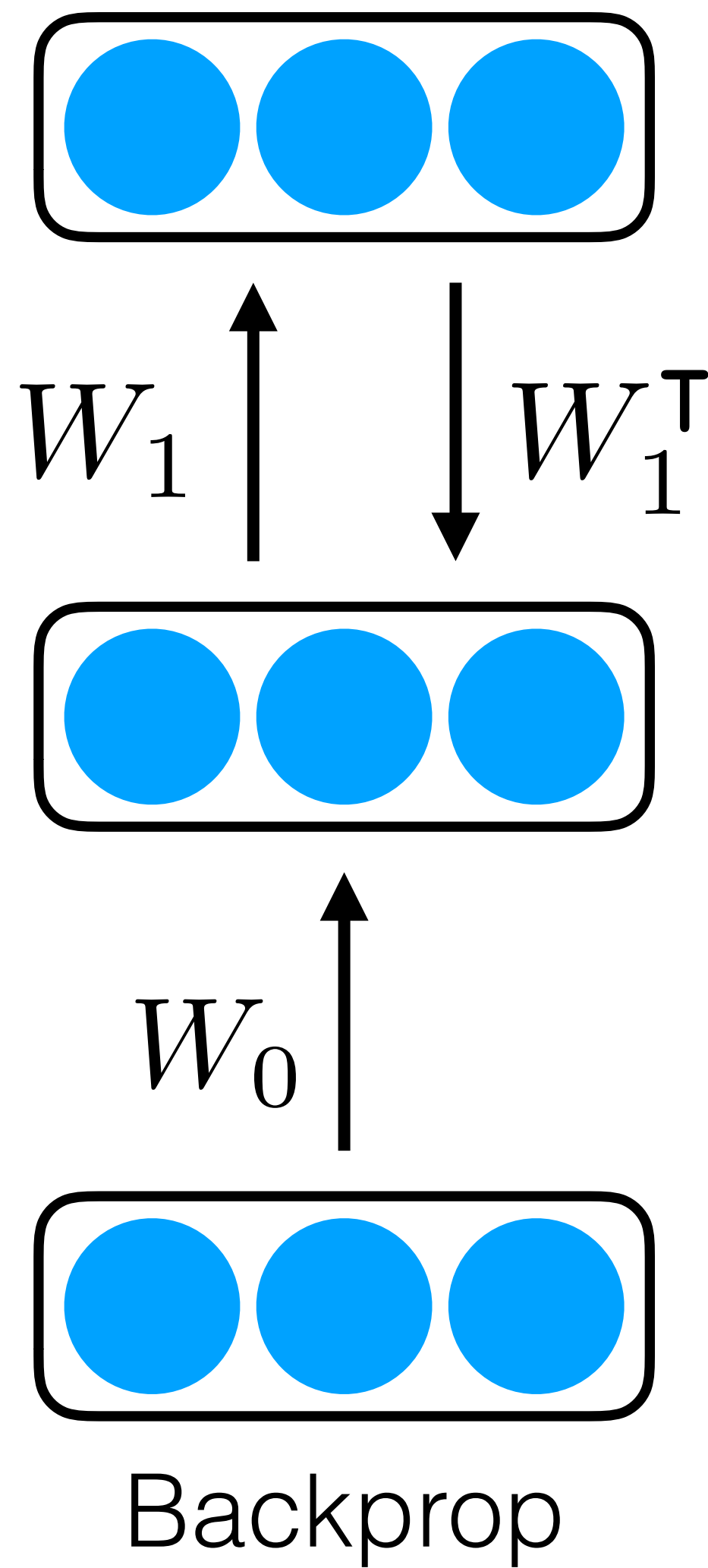
Individual neurons
are unidirectional.



Prediction in artificial neural networks is inspired by the brain.



Is *learning* in the brain inspired by artificial neural networks?



ARTICLE

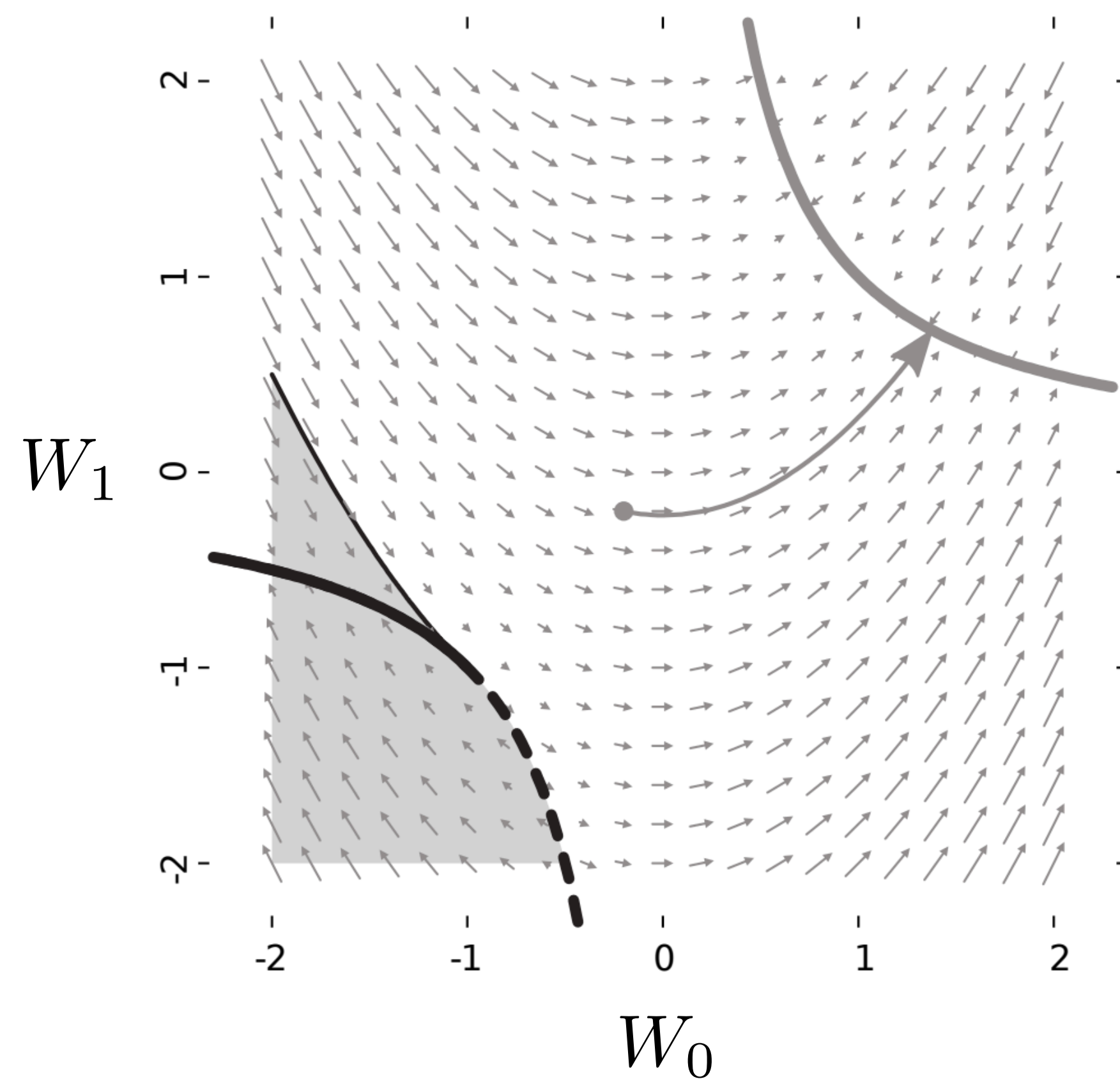
Received 7 Jan 2016 | Accepted 16 Sep 2016 | Published 8 Nov 2016

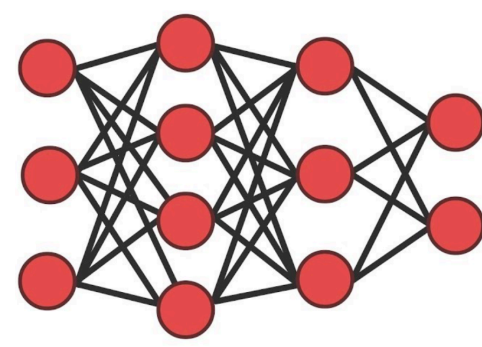
DOI: 10.1038/ncomms13276

OPEN

Random synaptic feedback weights support error backpropagation for deep learning

Timothy P. Lillicrap^{1,2}, Daniel Cownden³, Douglas B. Tweed^{4,5} & Colin J. Akerman¹





Prediction in artificial neural networks is inspired by the brain.
Is *learning* in the brain inspired by artificial neural networks?



Random feedback weights support learning in deep neural networks

Timothy P. Lillicrap, Daniel Cownden, Douglas B. Tweed, Colin J. Akerman

Difference Target Propagation

Dong-Hyun Lee, Saizheng Zhang, Asja Fischer, Yoshua Bengio

Towards Biologically Plausible Deep Learning

Yoshua Bengio, Dong-Hyun Lee, Jorg Bornschein, Thomas Mesnard, Zhouhan Lin

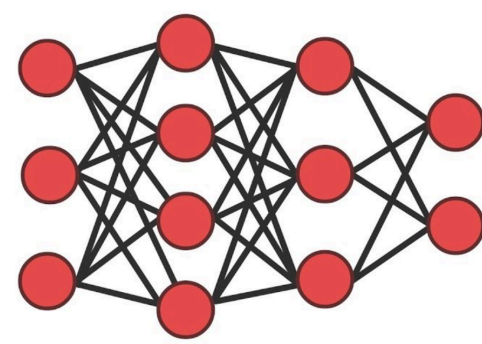
How Important is Weight Symmetry in Backpropagation?

Qianli Liao, Joel Z. Leibo, Tomaso Poggio

Equivalence of Equilibrium Propagation and Recurrent Backpropagation

Benjamin Scellier, Yoshua Bengio

The weight symmetry problem [is] arguably the crux of BP's biological implausibility.



Prediction in artificial neural networks is inspired by the brain.
Is *learning* in the brain inspired by artificial neural networks?



NeurIPS 2018

Assessing the Scalability of Biologically-Motivated Deep Learning Algorithms and Architectures

Sergey Bartunov
DeepMind

Adam Santoro
DeepMind

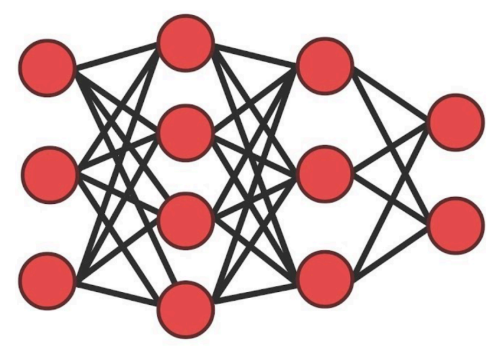
Blake A. Richards
University of Toronto

Luke Marris
DeepMind

Geoffrey E. Hinton
Google Brain

Timothy P. Lillicrap
DeepMind, University College London

Many of these algorithms perform well for MNIST, but for CIFAR and ImageNet we find that TP and FA variants perform significantly worse than BP, especially for networks composed of locally connected units, opening questions about whether new architectures and algorithms are required to scale these approaches.

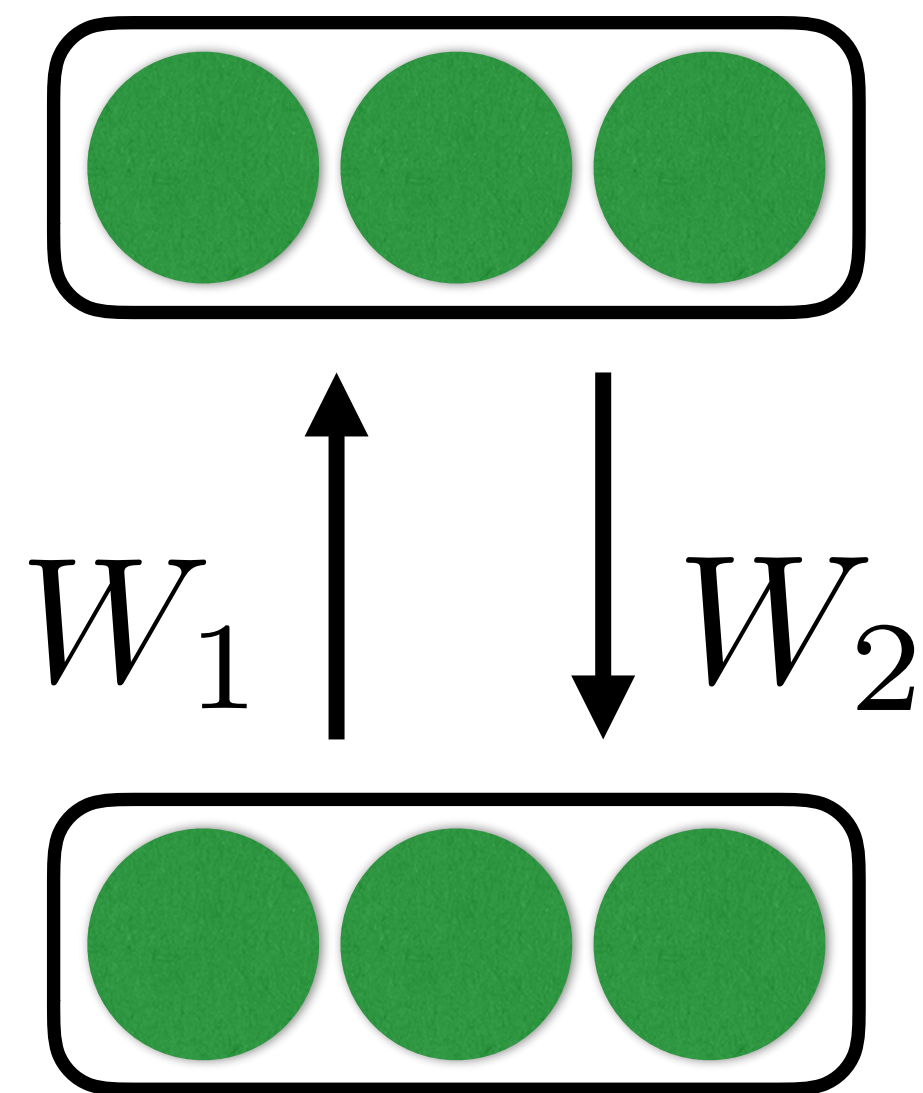


Prediction in artificial neural networks is inspired by the brain.

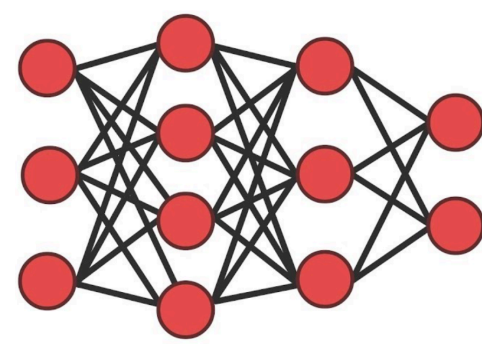


Is *learning* in the brain inspired by artificial neural networks?

Theorem: L_2 -regularized linear autoencoders are *symmetric* at all critical points.



$$W_2 = W_1^+ \xrightarrow{\text{weight decay}} W_2 = W_1^T$$



Prediction in artificial neural networks is inspired by the brain.



Is *learning* in the brain inspired by artificial neural networks?

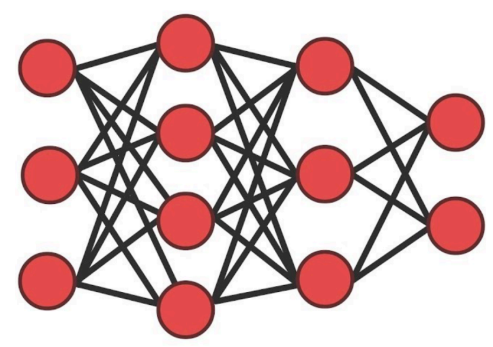
Theorem: L_2 -regularized linear autoencoders are *symmetric* at all critical points.

On Sun, Feb 10, 2019 at 2:16 PM Yoshua Bengio <yoshua.bengio@mila.quebec> wrote:

Thanks for reaching out, this is interesting.

The question of obtaining the transpose is actually pretty important for research on a biologically plausible version of backprop, because if you obtain approximate transposes, then several local learning rules give rise to gradient estimator analogues of backprop.

Cheers,
-- Yoshua

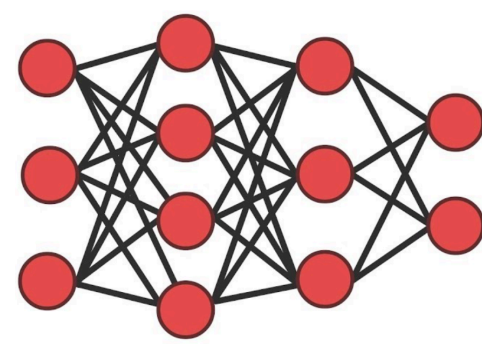


Prediction in artificial neural networks is inspired by the brain.
Is *learning* in the brain inspired by artificial neural networks?



Yes.

By pure logic.

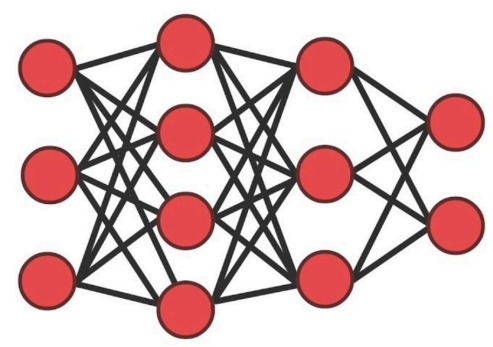


Prediction in artificial neural networks is inspired by the brain.
Is *learning* in the brain inspired by artificial neural networks?



Maybe.

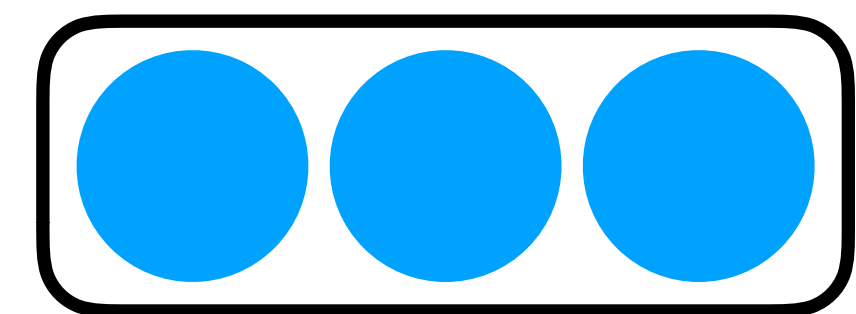
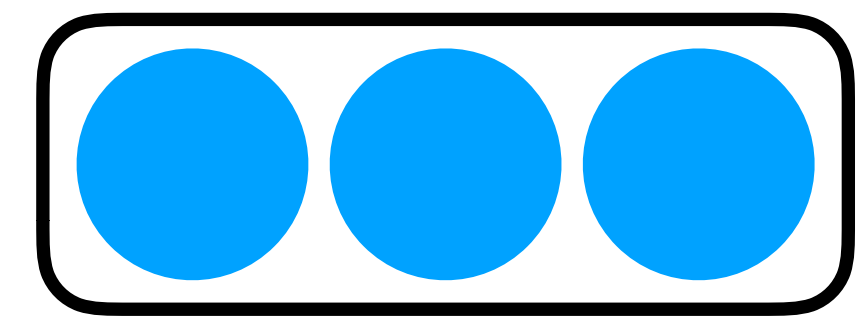
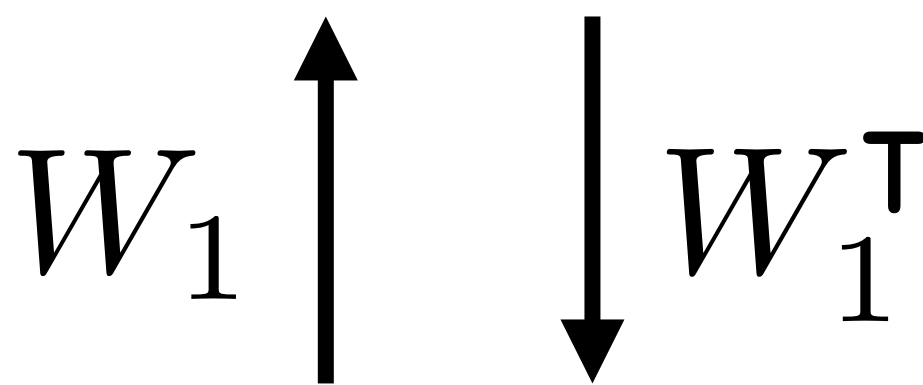
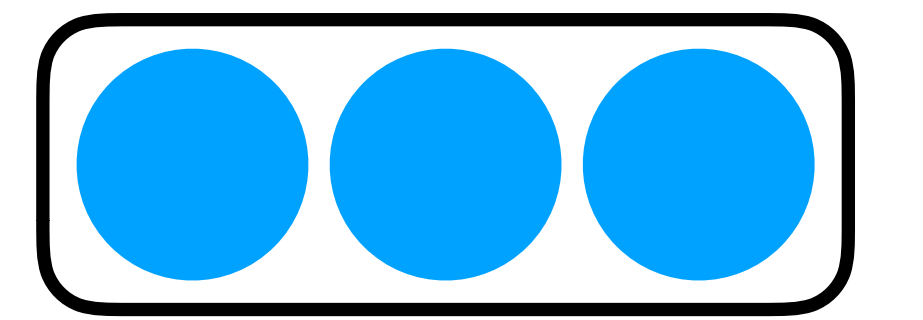
*And wouldn't it be fun to build
toward statements that could be
verified or falsified with rigor?*



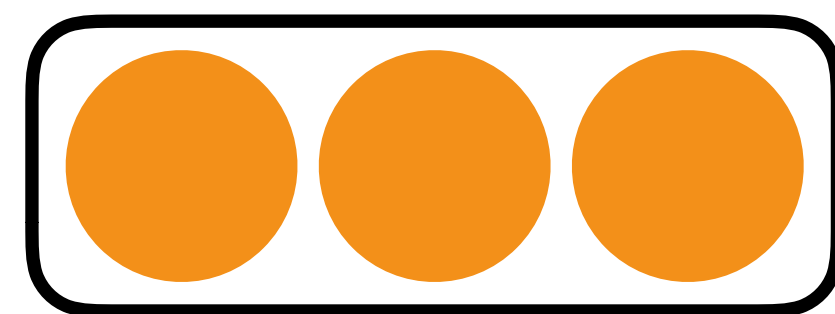
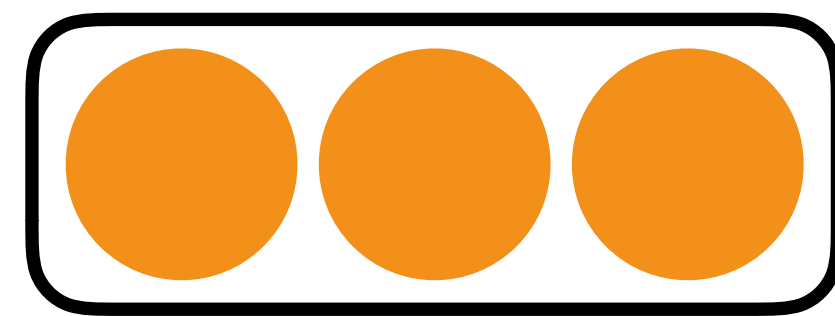
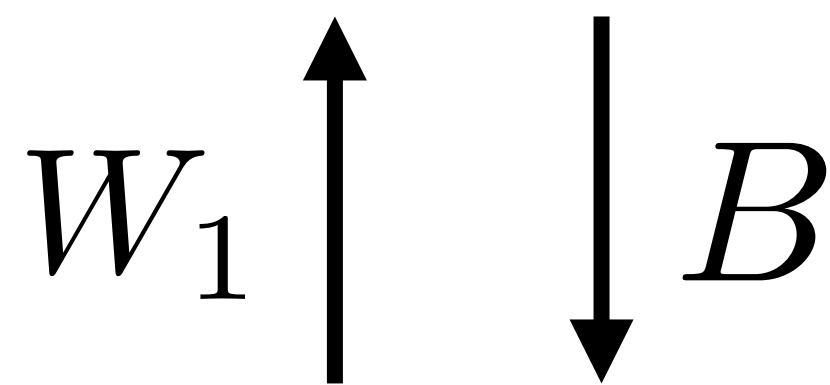
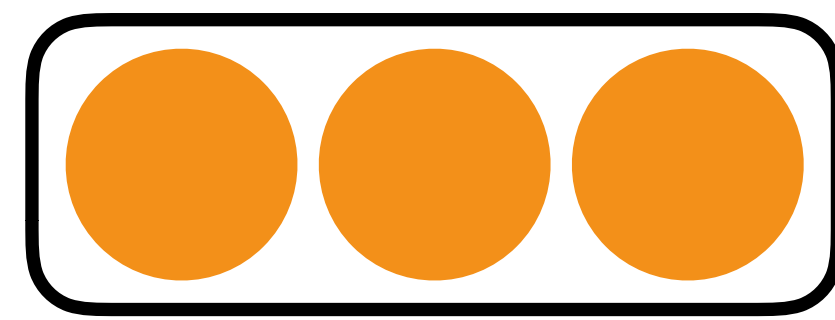
Prediction in artificial neural networks is inspired by the brain.



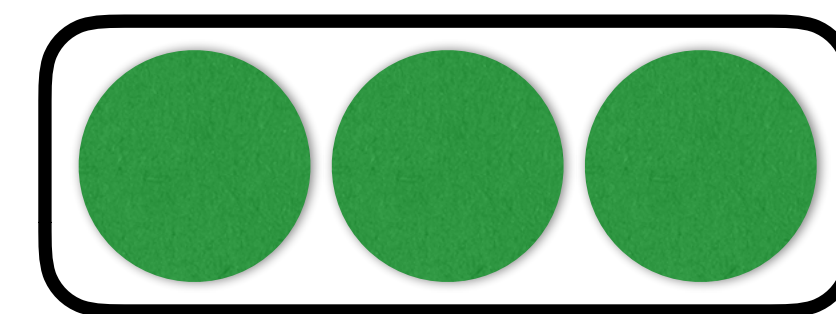
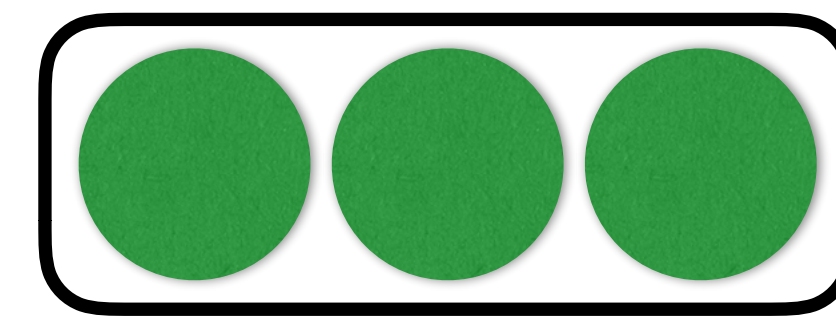
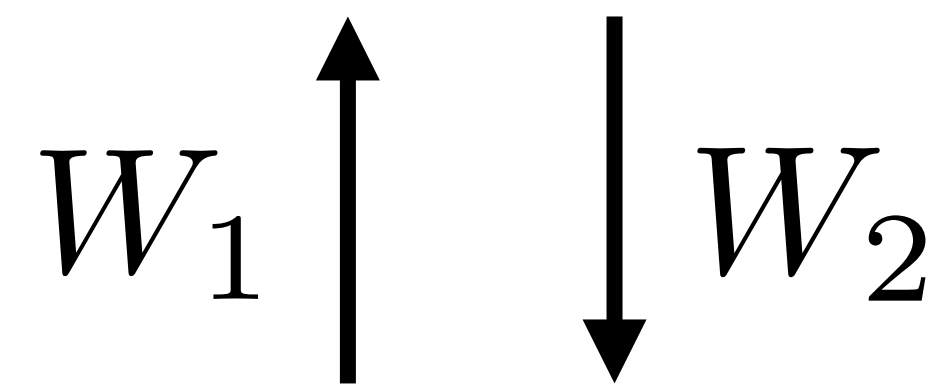
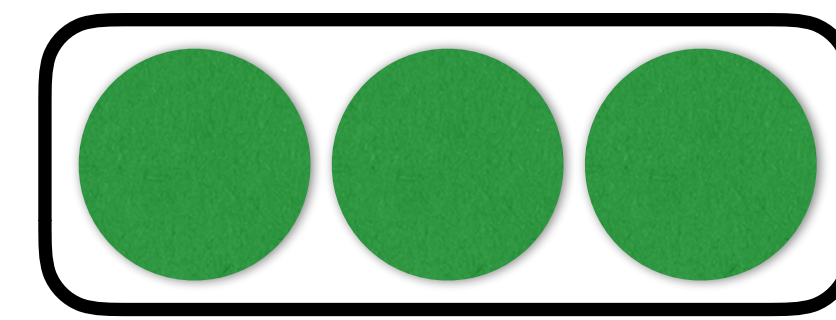
Is *learning* in the brain inspired by artificial neural networks?



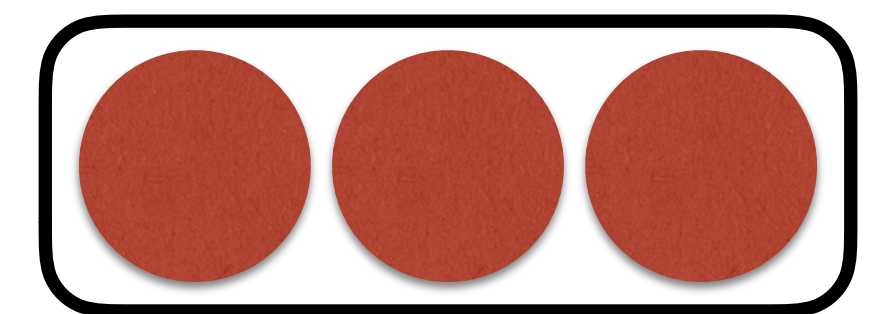
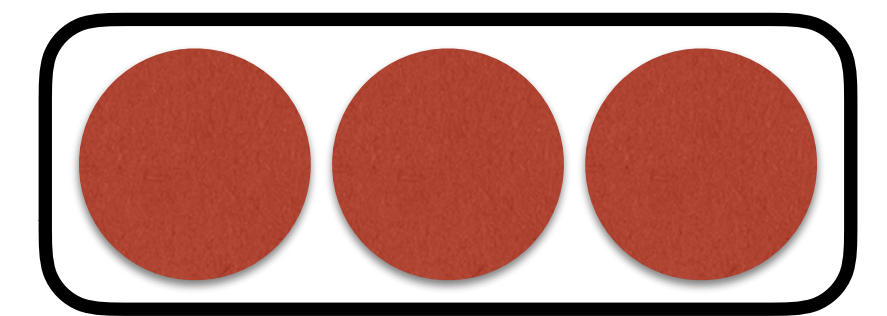
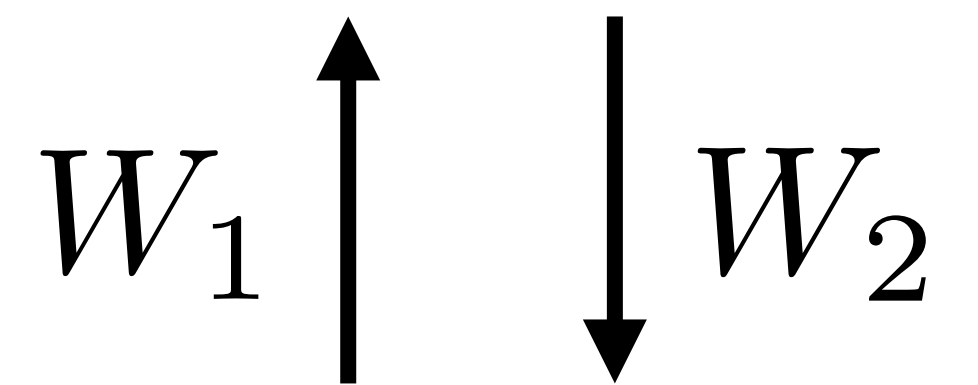
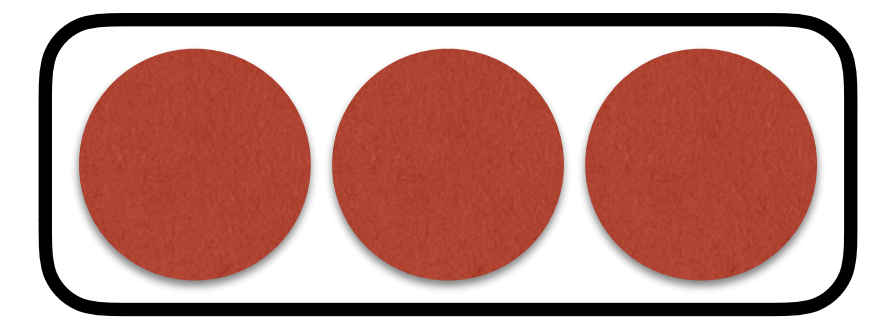
Backprop



Feedback Alignment



Information Alignment



Symmetric Alignment

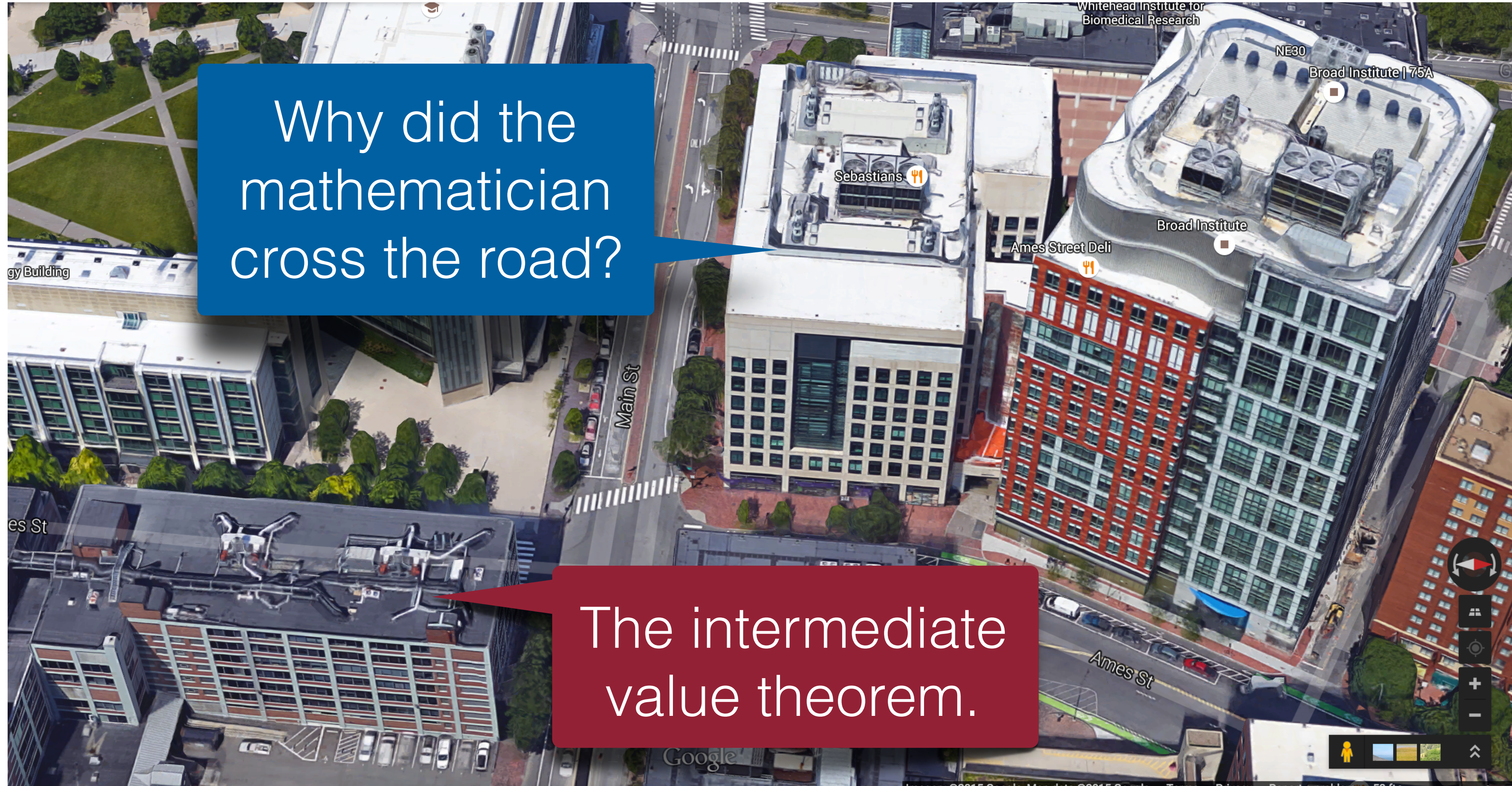


CBMM



STANLEY CENTER
FOR PSYCHIATRIC RESEARCH

AT BROAD INSTITUTE



Why did the mathematician cross the road?

The intermediate value theorem.

MIT
MATHEMATICS

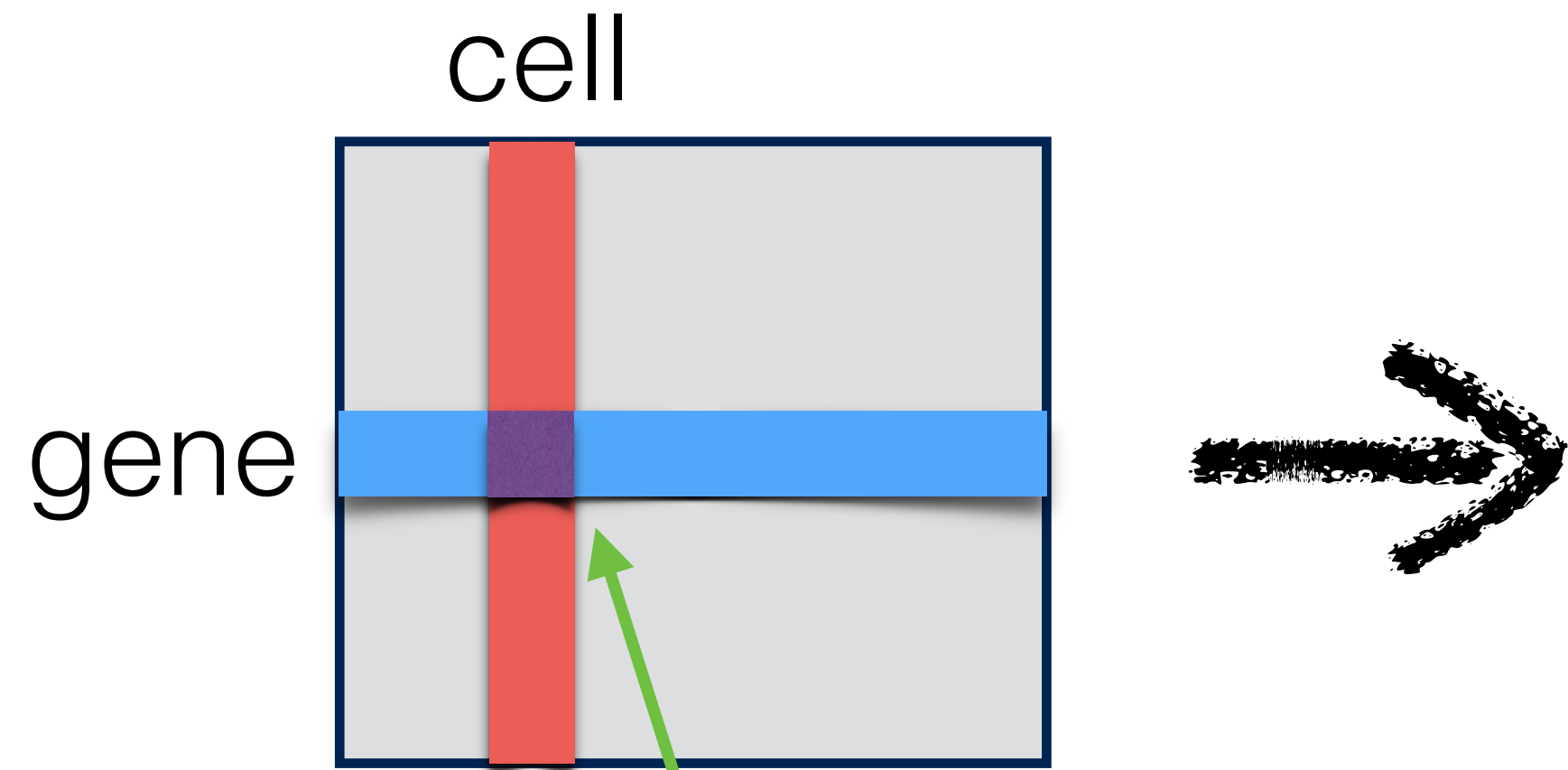




STANLEY CENTER
FOR PSYCHIATRIC RESEARCH

AT BROAD INSTITUTE



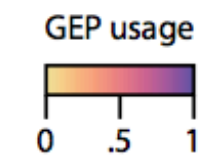


RNA molecules for that gene in that cell

Organoid cell-types and activities

Identity GEP

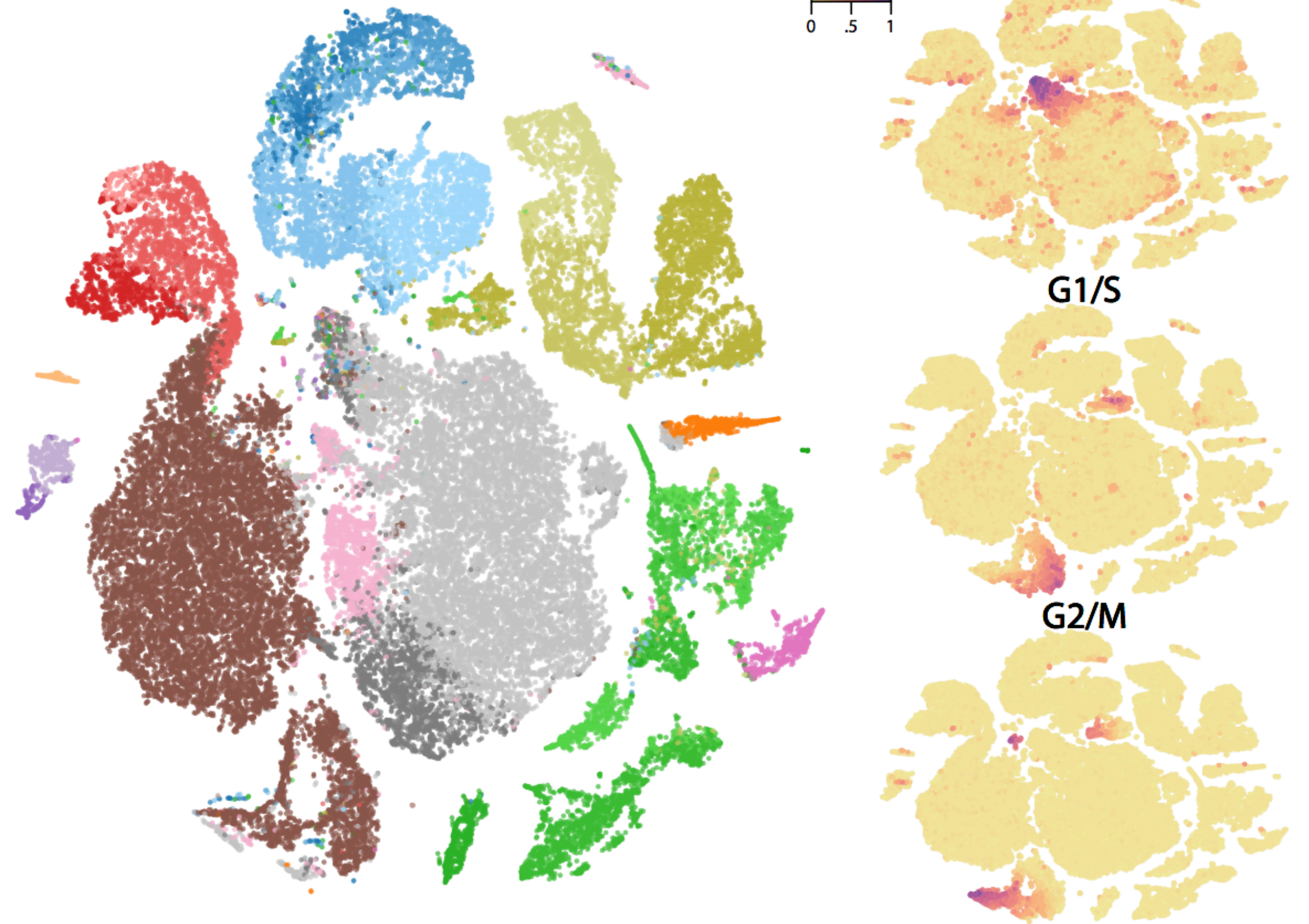
- Astro-1
- Astro-2
- Astro-3
- Astro-4
- Astro-5
- Astro-6
- Ret-1
- Ret-2
- Ret-3
- Ret-4
- Ret-5
- Ret-6
- FB-1
- FB-2
- FB-3
- Dop-1
- Dop-2
- NE-1
- NE-2
- Stem-like
- PP
- Musc-T1
- Musc-Im
- Musc-T2
- C6-1
- C6-2
- C7
- C8



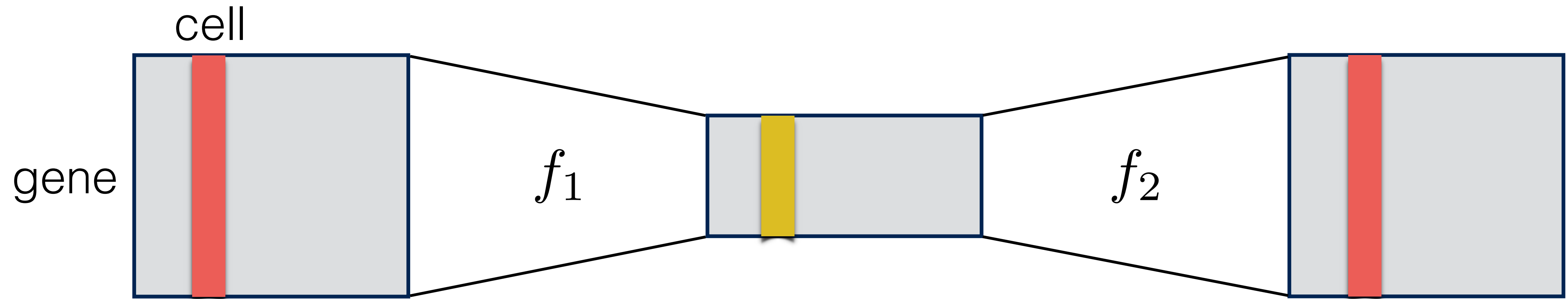
Hypoxia

G1/S

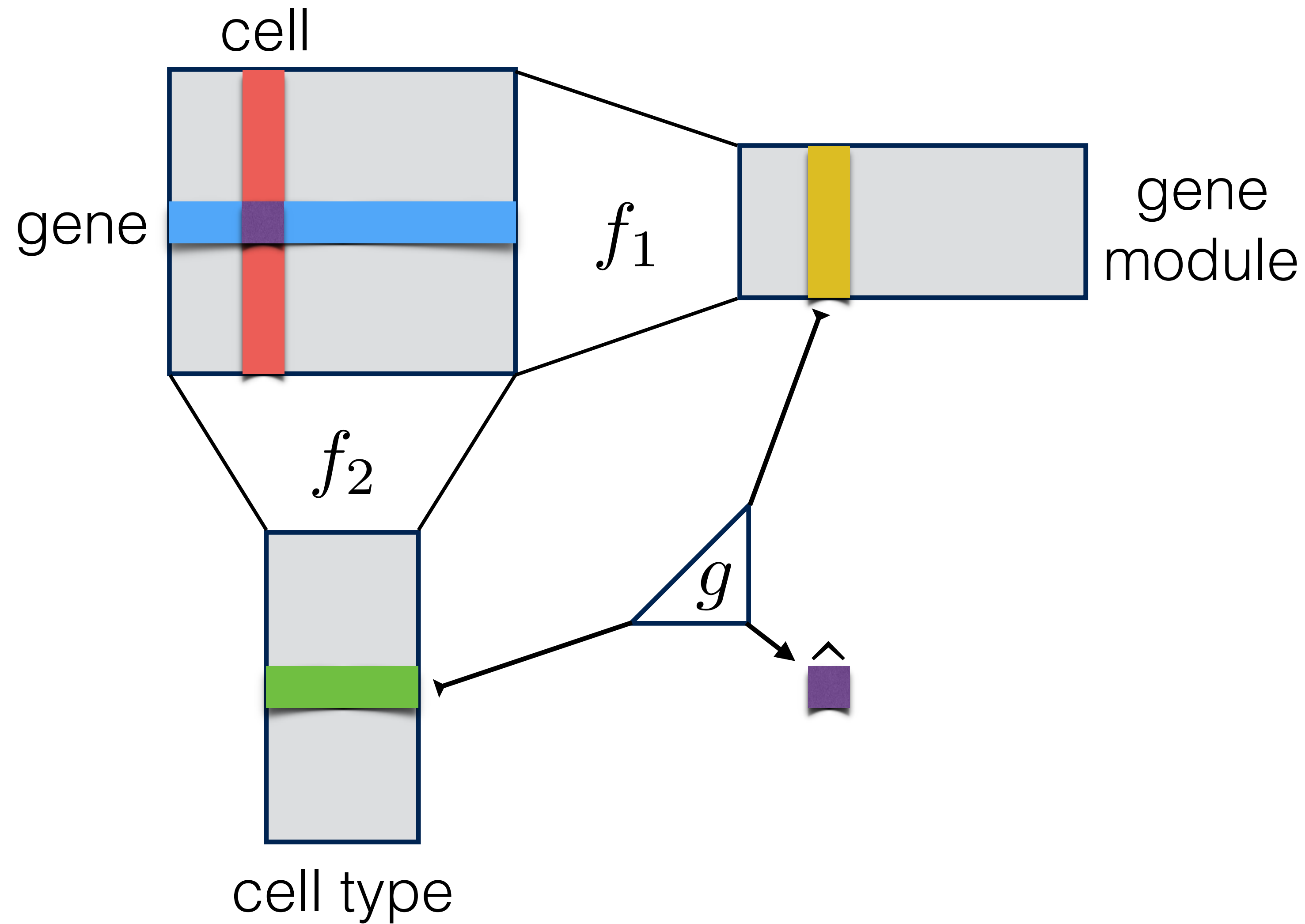
G2/M



Representation Learning of Cell



Representation Learning of Cell++

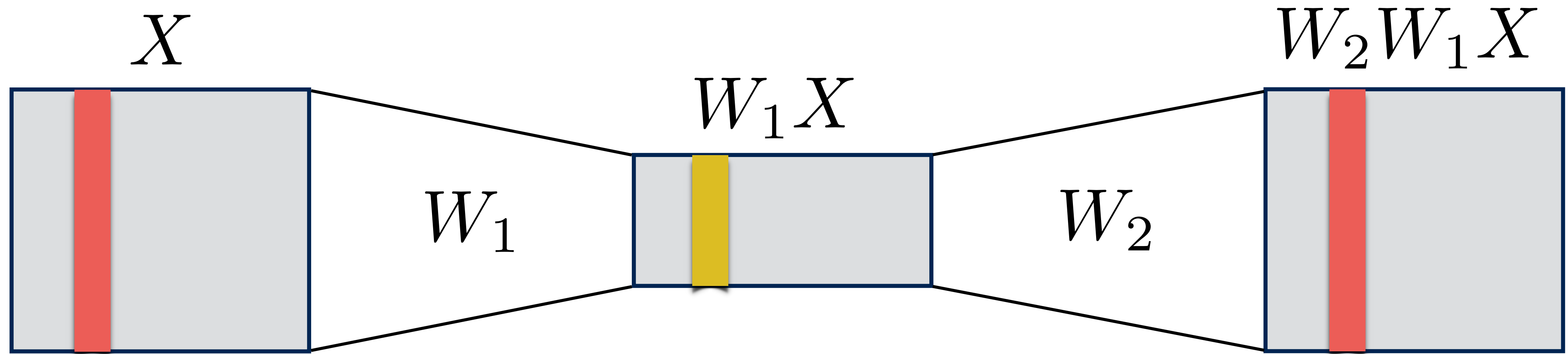


Non-linear
Singular Value
Decomposition

$$X = U\Sigma V^T$$



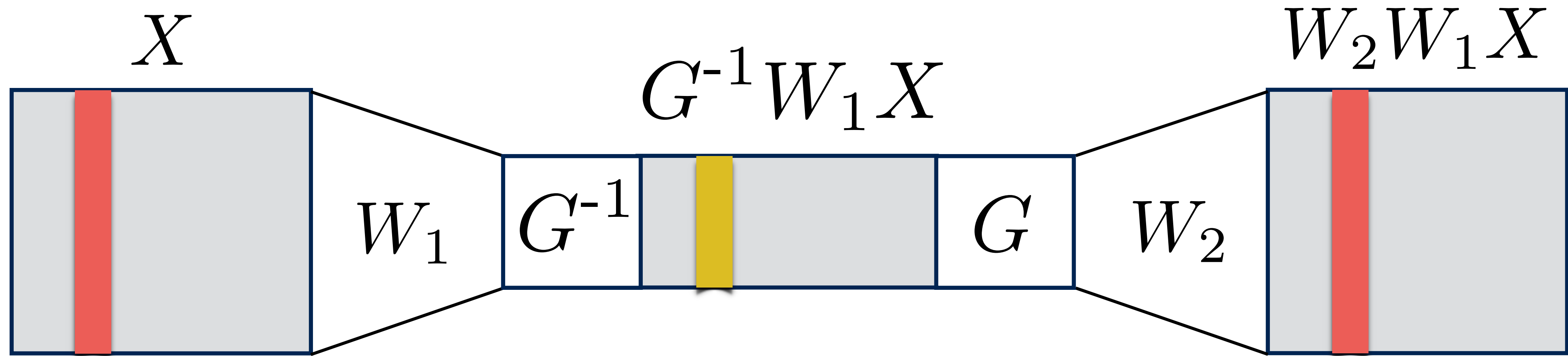
Linear Autoencoder



$$\mathcal{L}(W_1, W_2) = \|X - W_2 W_1 X\|^2$$

$$\begin{array}{lll} X = U \Sigma V^\top & W_1 = U_k^\top & W_2 W_1 = U_k U_k^\top \\ X X^\top = U \Sigma^2 U^\top & W_2 = U_k & \end{array}$$

Linear Autoencoder



$$\mathcal{L}(W_1, W_2) = \|X - W_2W_1X\|^2 = \|X - (W_2G)(G^{-1}W_1)X\|^2$$

$$X = U\Sigma V^T \quad W_1 = G^{-1}U_k^T \quad W_2W_1 = U_kU_k^T$$

$$XX^T = U\Sigma^2U^T \quad W_2 = U_kG$$

Fact: Linear autoencoders are *pseudoinverses* at all critical points: $W_2 = W_1^+$

Topology of PCA

Problem

Find the k -plane closest to a point cloud in \mathbb{R}^m .

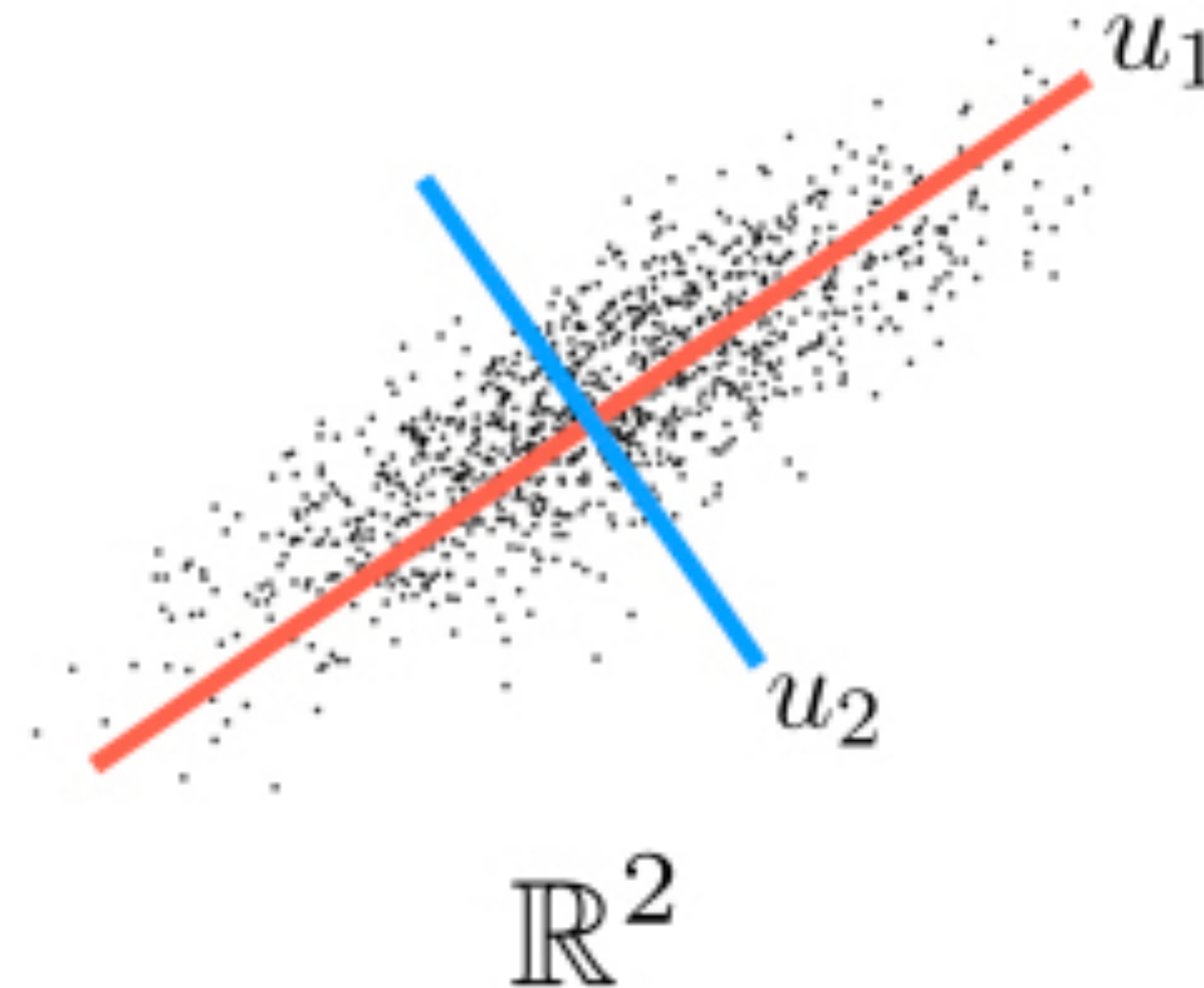
Domain

The manifold whose points are k -planes in \mathbb{R}^m .

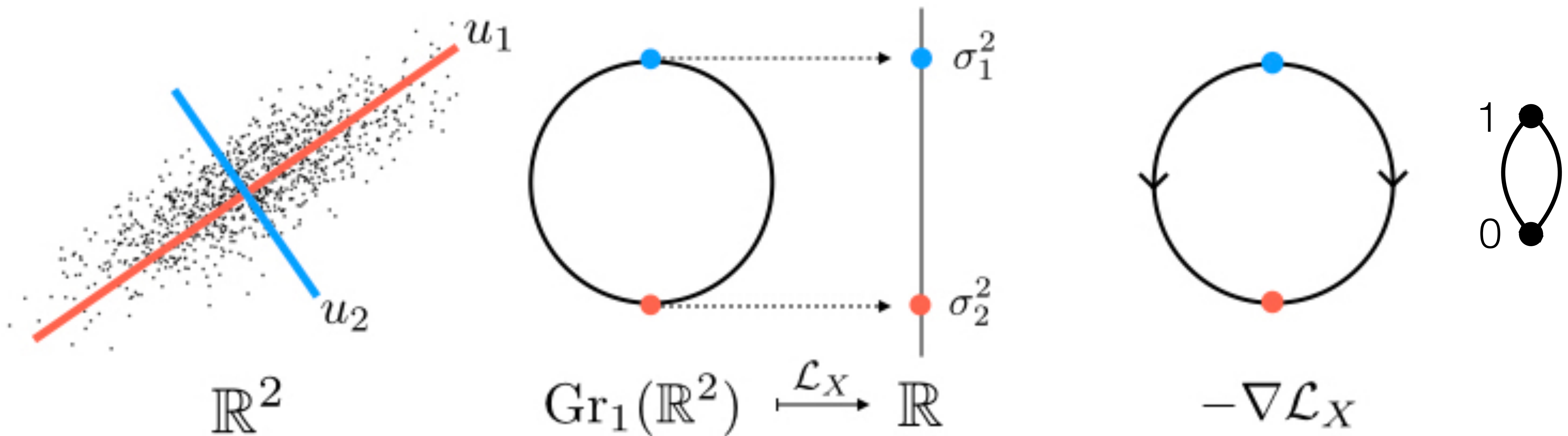
I.e., the *Grassmannian* manifold:

$$\text{Gr}_k(\mathbb{R}^m) \cong \{P = P^2, P = P^\top, \text{tr } P = k\} \subset \mathbb{R}^{m \times m}$$

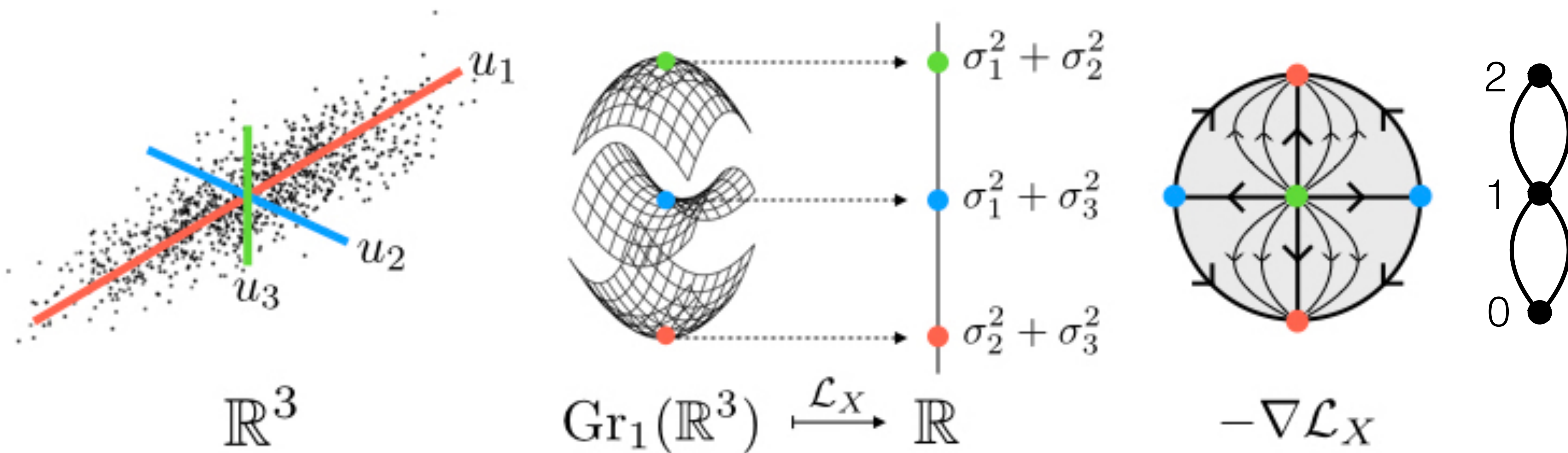
$$\mathcal{L}_X : \text{Gr}_k(\mathbb{R}^m) \rightarrow \mathbb{R} \quad \mathcal{L}_X(P) = \|X - PX\|^2$$



Topology of PCA



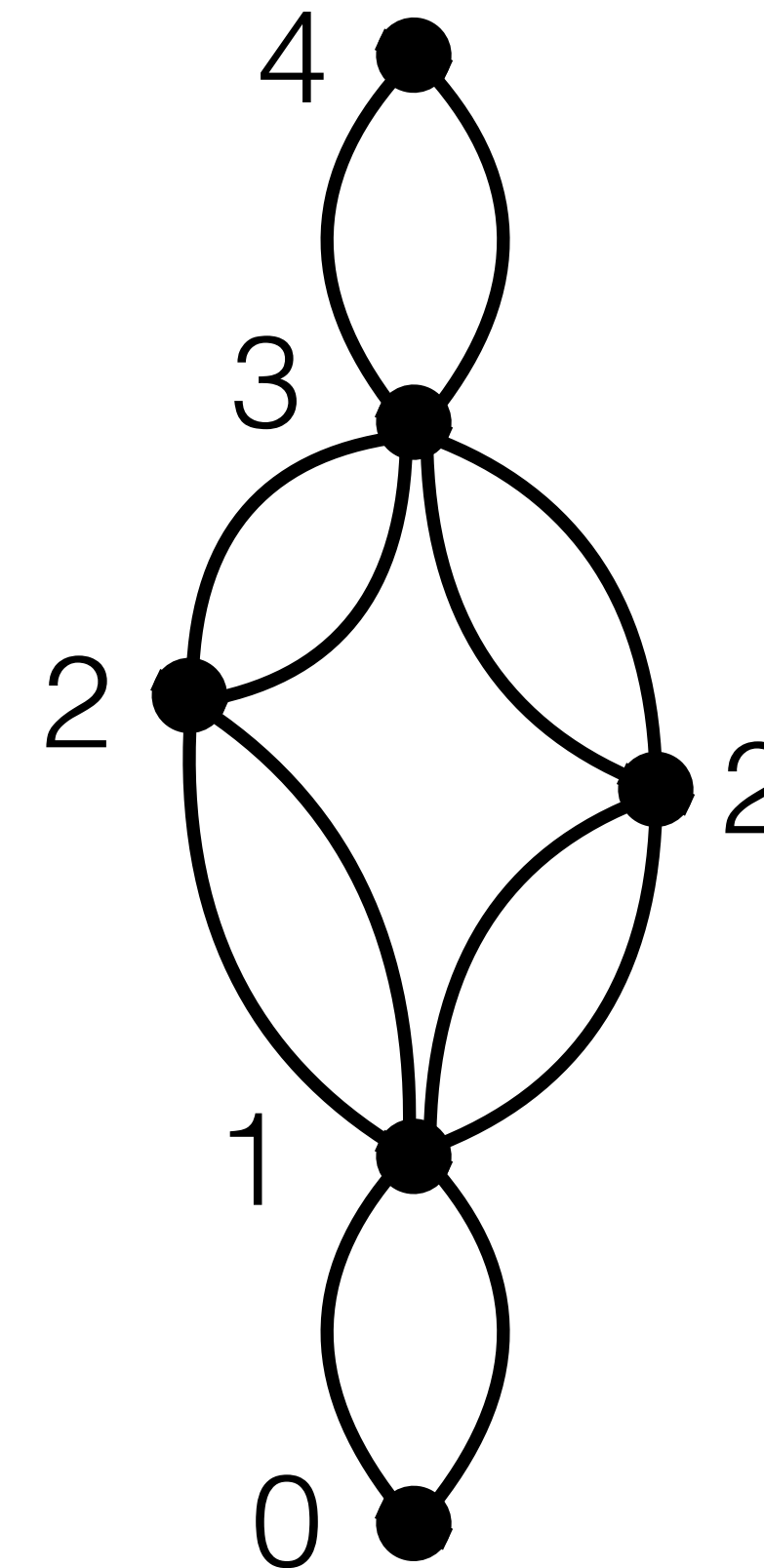
Topology of PCA



Topology of PCA

$$\mathcal{L}_X : \text{Gr}_2(\mathbb{R}^4) \rightarrow \mathbb{R}$$

d	u_1	u_2	u_3	u_4
4			•	•
3		•		•
2		•	•	
2	•			•
1	•		•	
0	•	•		

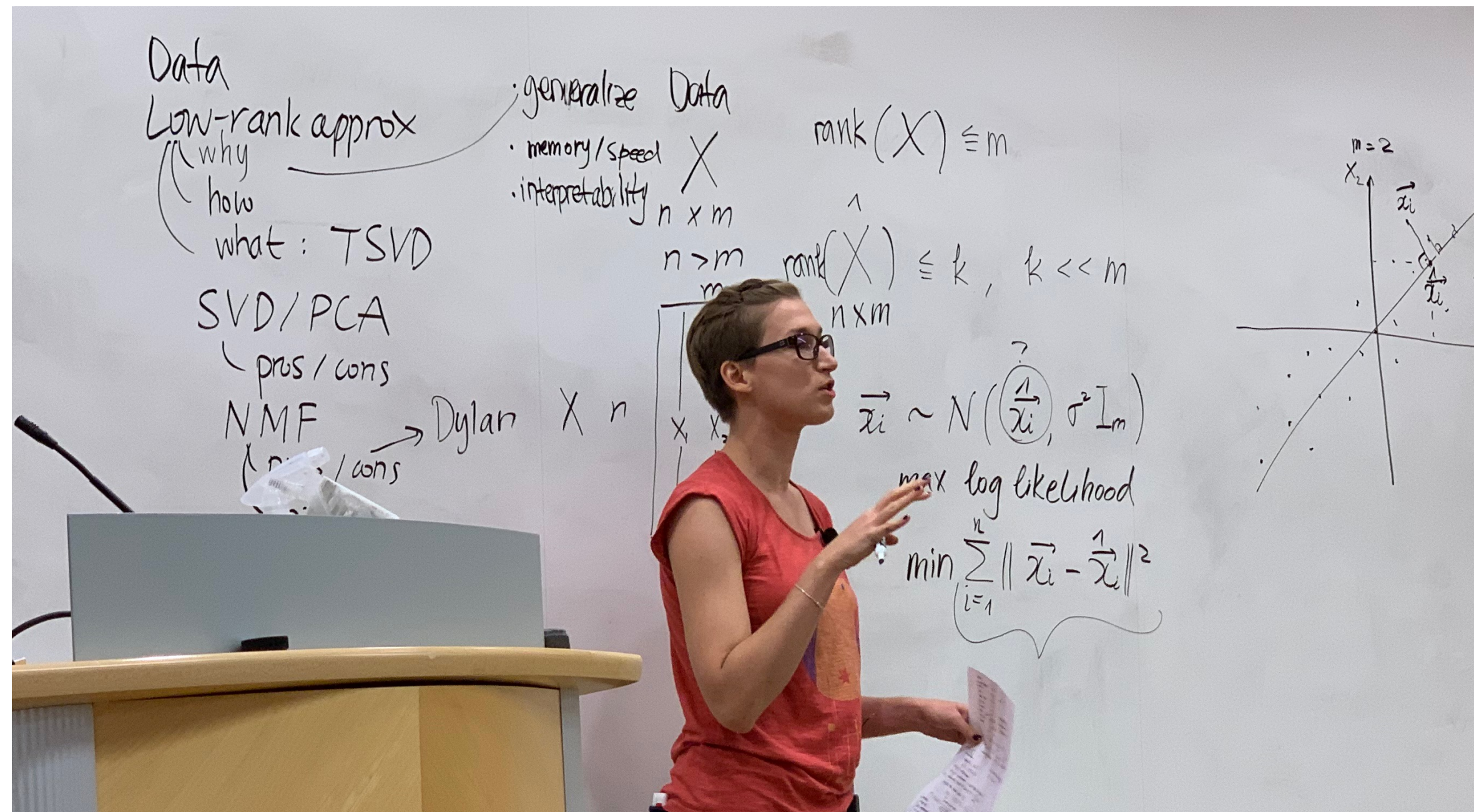
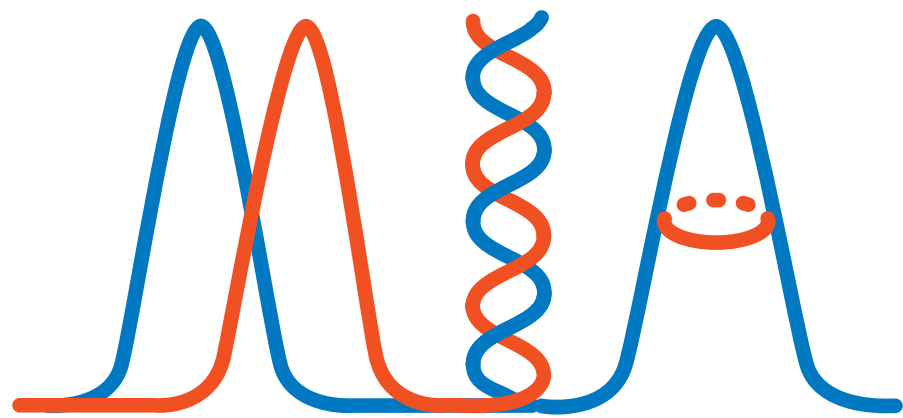


$$\dim \text{Gr}_k(\mathbb{R}^m) = k(m - k)$$

*To visualize 4 dimensions, first visualize n dimensions and then let $n = 4$.

Topology of PCA

- $\mathcal{L}_X : \text{Gr}_k(\mathbb{R}^m) \rightarrow \mathbb{R}$ is Morse iff singular values are positive and distinct.
- $\binom{m}{k}$ critical points are the principal k-planes.
- Critical values are sums of eigenvalues \Rightarrow toy model for random matrix theory.
- Gradient trajectories between principal planes in adjacent index rotate one principal direction in first plane to another in second plane fixing the rest.
- There are exactly two such trajectories \implies perfect \mathbb{F}_2 -Morse function.
- This rabbit hole goes far *deeper*...



Models, Inference and Algorithms Primer • 2019

From Morse theory to geometric ensembling via the topology of PCA

Jon Bloom & Cotton Seed

BROAD INSTITUTE

0:02 / 1:45:26

Models, Inference and Algorithms Meeting • 2019

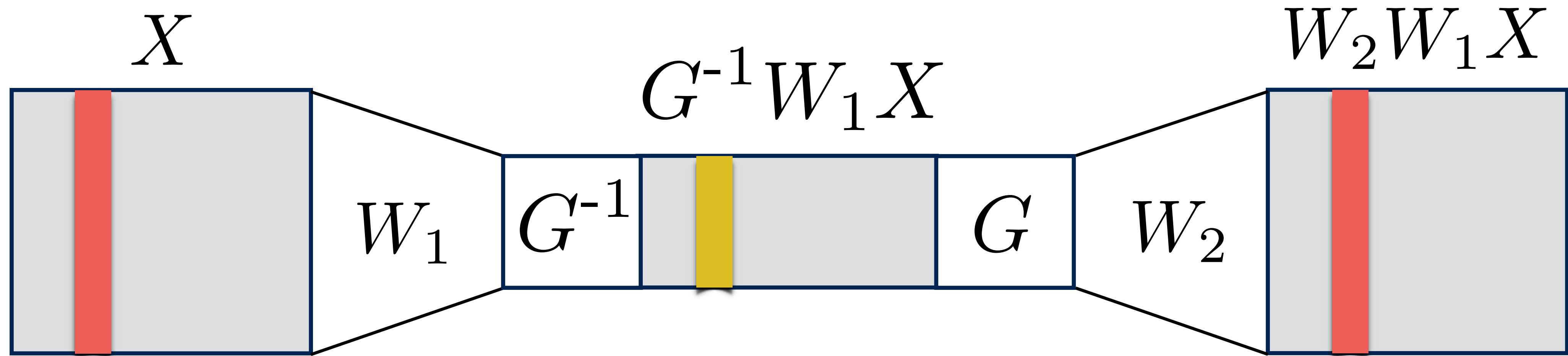
Regularized linear autoencoders, probabilistic PCA, & backpropagation in the brain

Daniel Kunin & Aleksandrina Goeva

BROAD INSTITUTE

58:00 / 1:45:26

Linear Autoencoder



$$\mathcal{L}(W_1, W_2) = \|X - W_2W_1X\|^2 = \|X - (W_2G)(G^{-1}W_1)X\|^2$$

$$X = U\Sigma V^T \quad W_1 = G^{-1}U_k^T \quad W_2W_1 = U_kU_k^T$$

$$W_2 = U_kG$$

Fact: Linear autoencoders are *pseudoinverses* at all critical points: $W_2 = W_1^+$

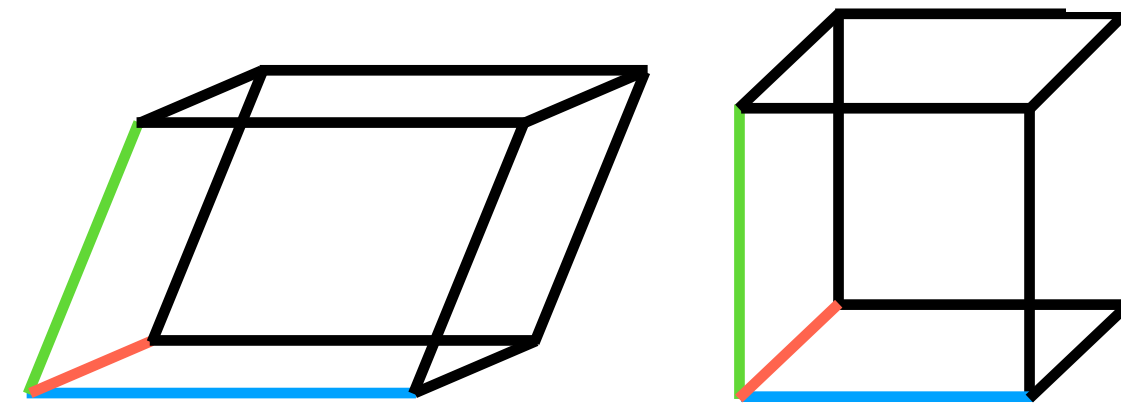
Regularization



Orthogonality

- Orthogonal matrices are the volume-preserving matrices of minimal Frobenius norm.

$$\sum \sigma_i^2 \quad \min_A \|A\|_F^2 \quad \text{s.t.} \quad \det(A)^2 = 1 \quad \prod \sigma_i^2$$

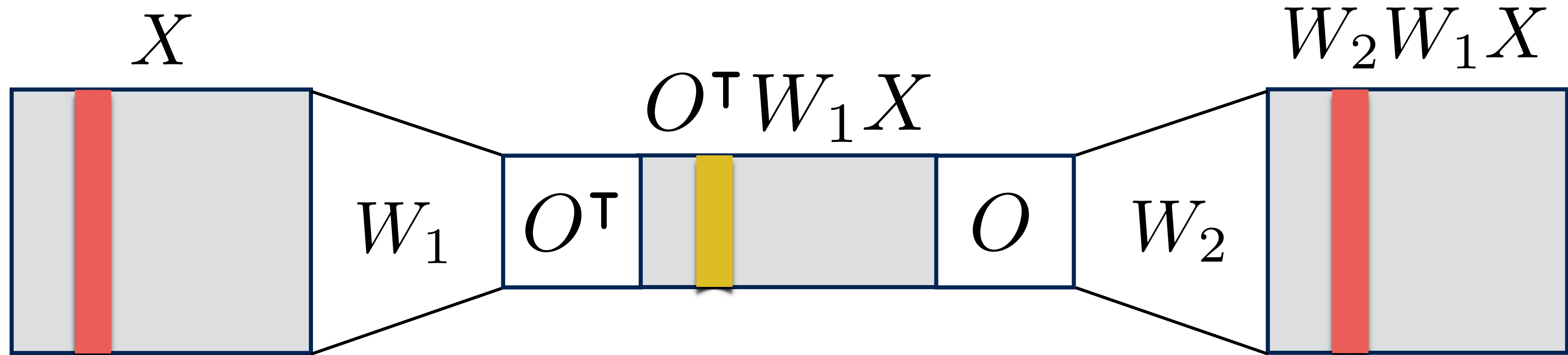


- Orthogonal matrices are the inverse matrices of minimum total Frobenius norm.

$$\sum (\sigma_i^2 + \sigma_i^{-2}) \quad \min_{A,B} \|A\|_F^2 + \|B\|_F^2 \quad \text{s.t.} \quad AB = I$$

In particular, $A = B^\top$.

Linear Autoencoder



$$\mathcal{L}_\sigma(W_1, W_2) = \|X - W_2W_1X\|^2 + \lambda(\|W_1\|^2 + \|W_2\|^2)$$

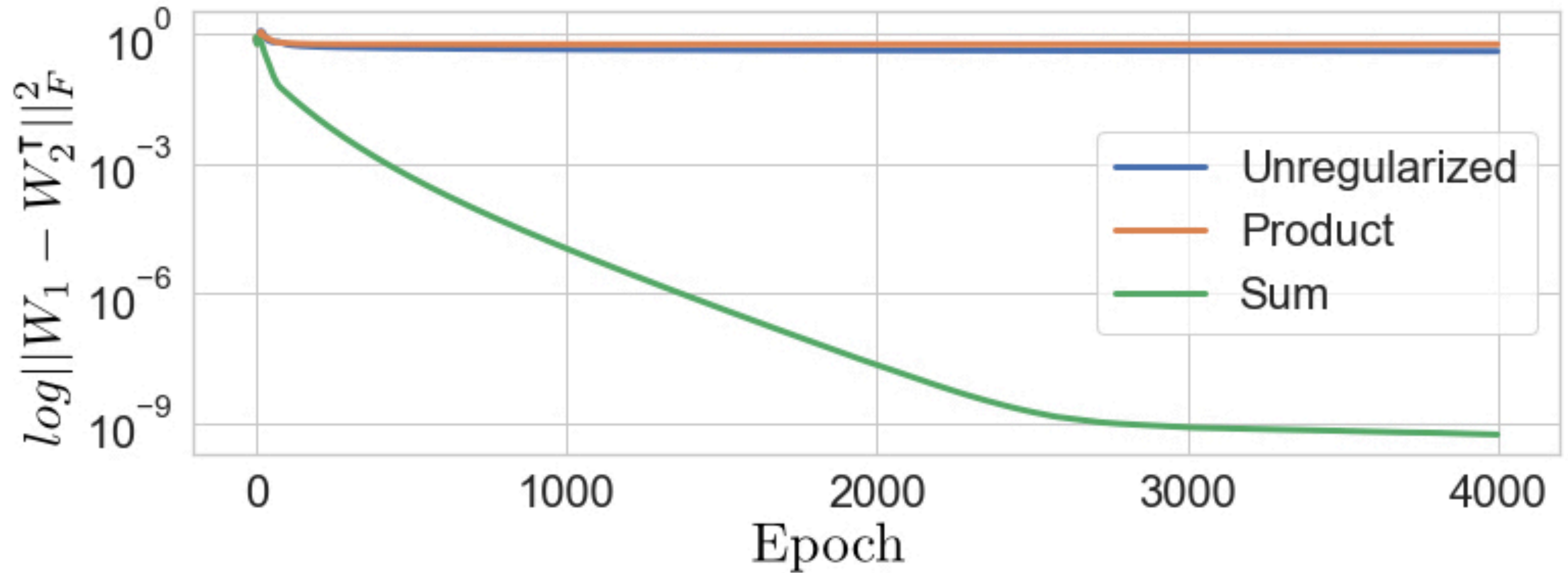
$$X = U\Sigma V^T \quad W_1 = O^T (I - \lambda\Sigma^{-2})^{-\frac{1}{2}} U_k^T$$

$$W_2W_1 = U_k (I - \lambda\Sigma^{-2})^{-1} U_k^T \quad W_2 = U_k (I - \lambda\Sigma^{-2})^{-\frac{1}{2}} O$$

Theorem: L_2 -regularized linear autoencoders are *symmetric* at all critical points.

$$W_2 = W_1^T$$

Linear Autoencoder



Linear Autoencoder

$$X = U\Sigma V^T \implies W_2 = U_k (I - \lambda \Sigma^{-2})^{-\frac{1}{2}} O$$

Theorem: L_2 -regularized linear autoencoders are *symmetric* at all critical points.

Theorem: The loss is strictly saddle (*Morse-Bott*). All minima are global.

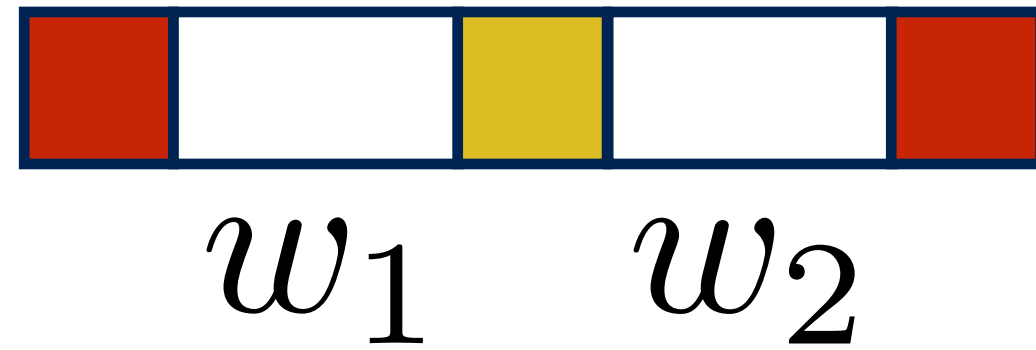
Theorem: The top principal directions of X with eigenvalue greater than λ coincide with the left singular vectors of the trained decoder W_2 . These eigenvalues are determined by the singular values of W_2 .

PCA algo: Fit a regularized LAE and SVD the decoder.



+ weight decay

Scalar Autoencoder



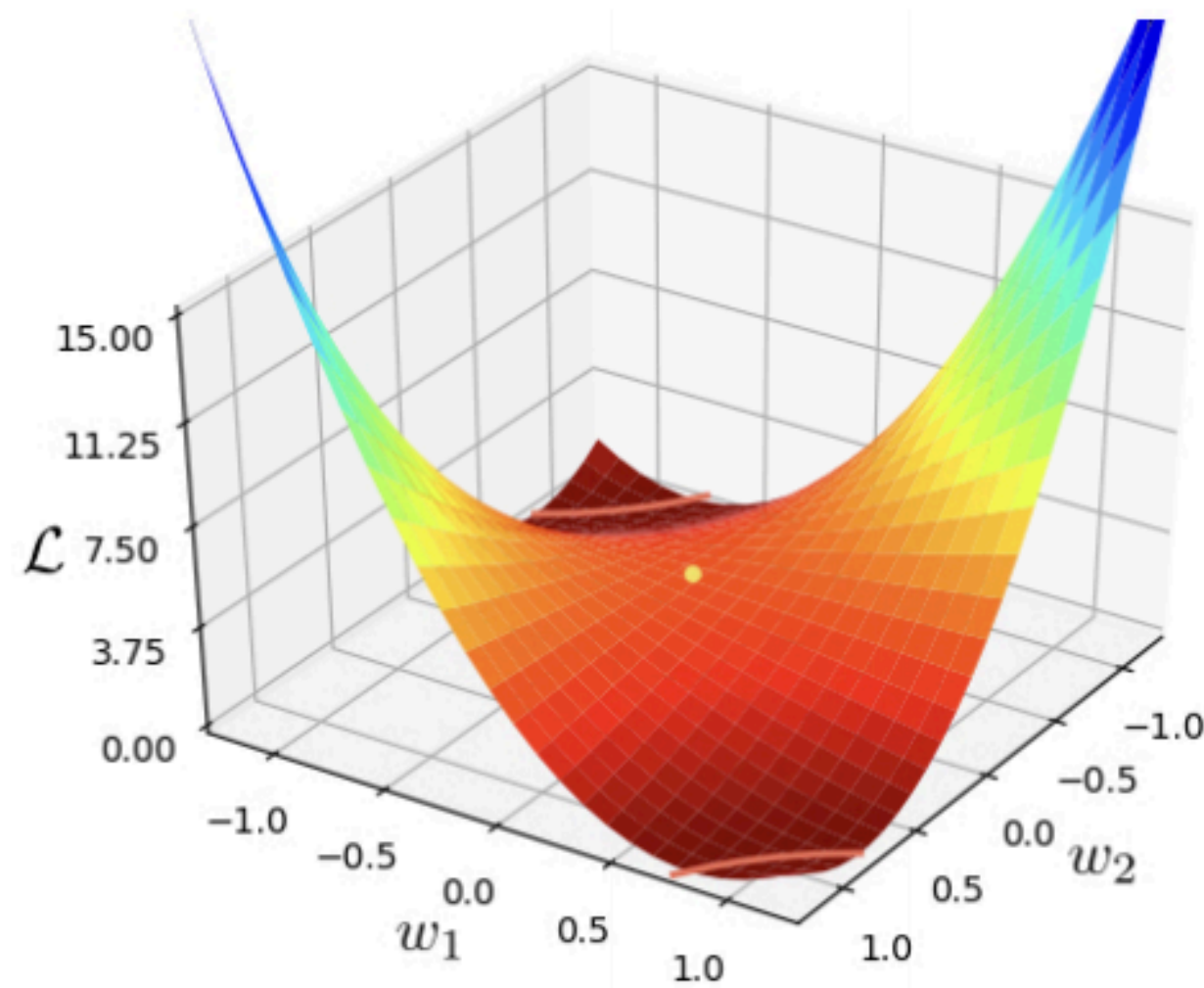
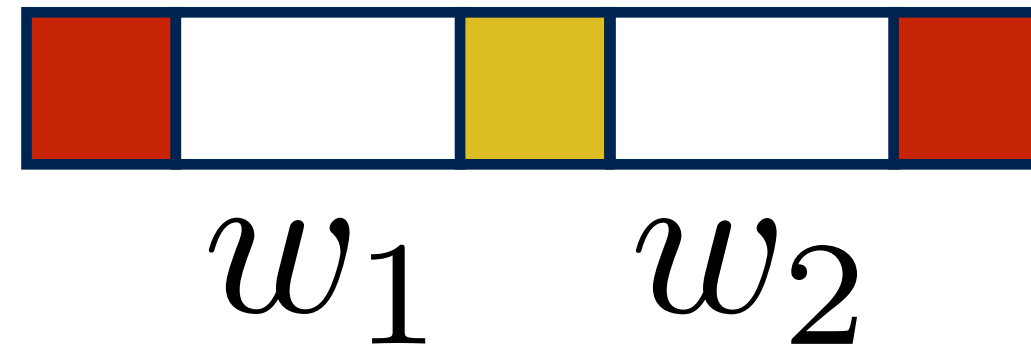
unregularized $\mathcal{L}(w_1, w_2) = (x - w_2 w_1 x)^2$

product $\mathcal{L}_\pi(w_1, w_2) = (x - w_2 w_1 x)^2 + \lambda(w_2 w_1)^2$

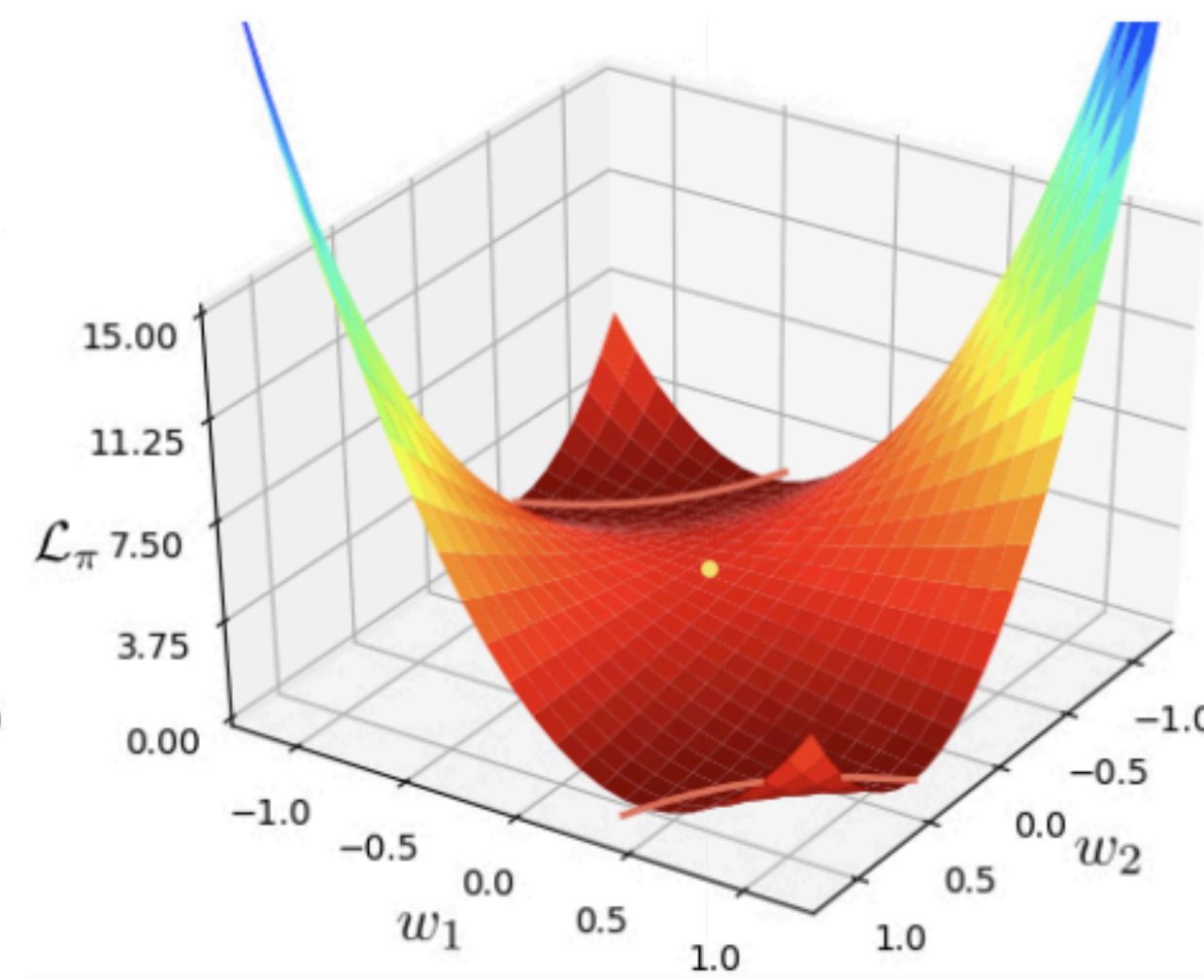
sum $\mathcal{L}_\sigma(w_1, w_2) = (x - w_2 w_1 x)^2 + \lambda(w_1^2 + w_2^2)$

[Scalar AE Visualization](#)

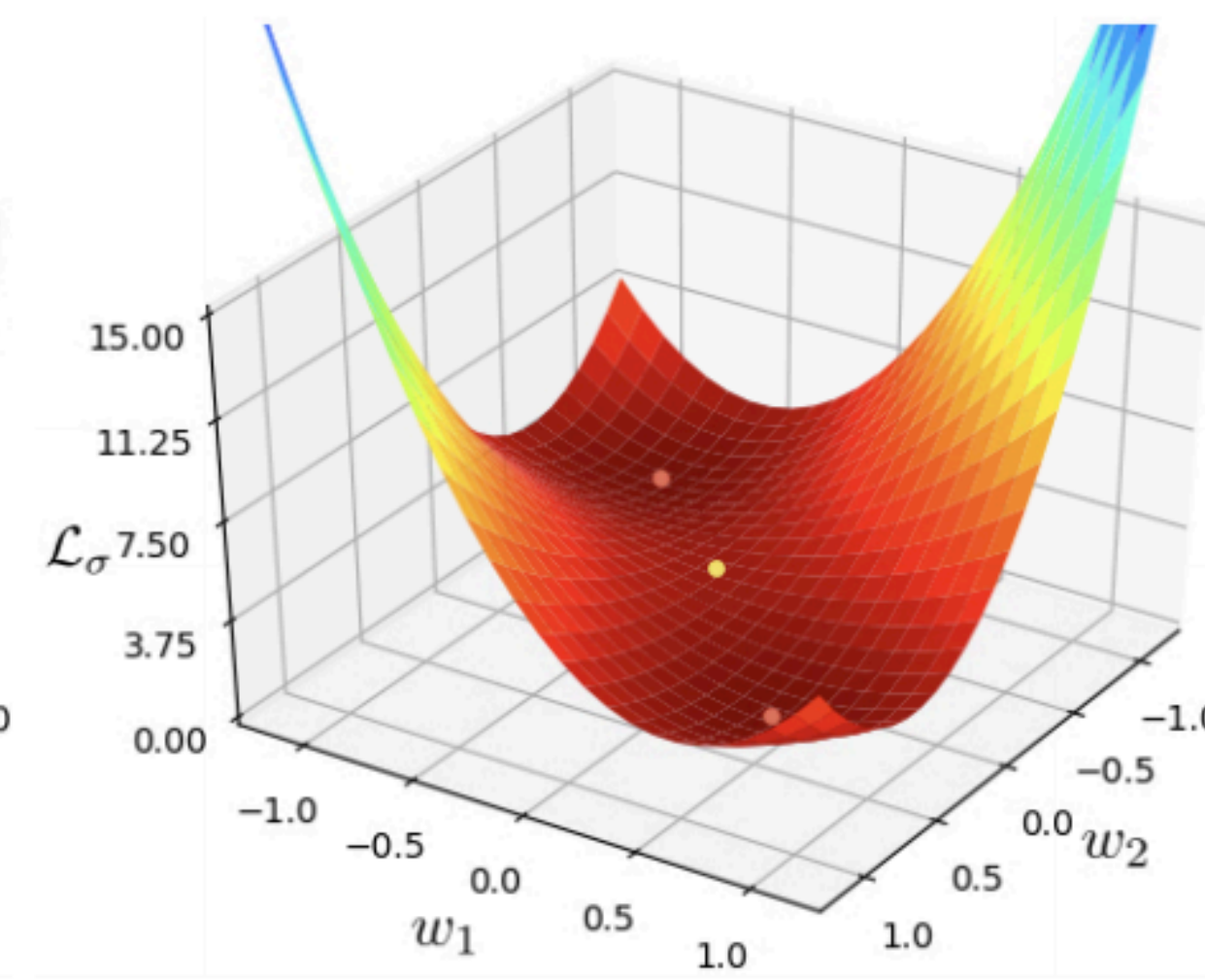
Scalar Autoencoder



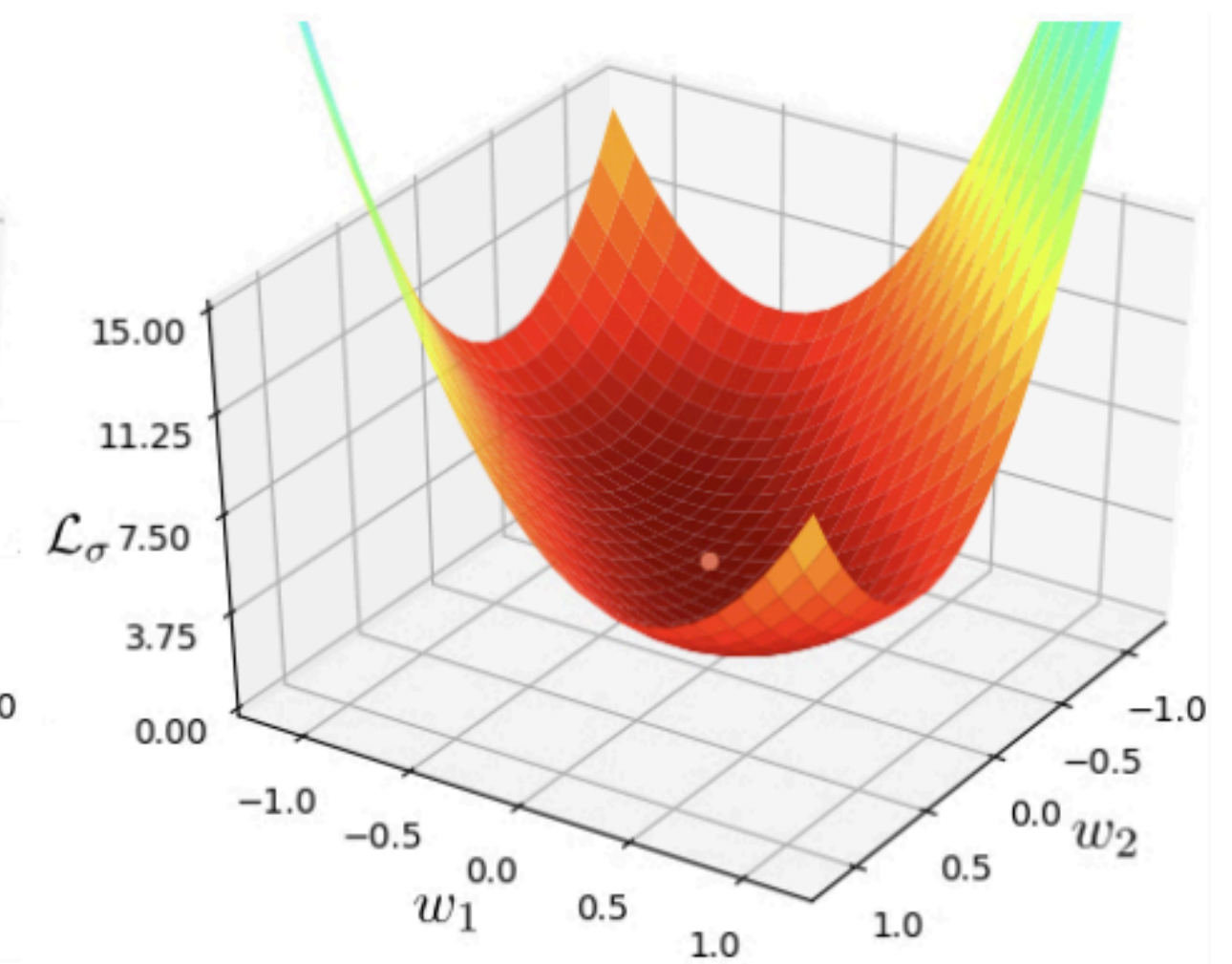
(a) Unregularized



(b) Product ($\lambda = 2$)



(c) Sum ($\lambda = 2$)



(d) Sum ($\lambda = 4$)

Figure 1. Scalar loss landscapes with $x^2 = 4$. Yellow points are saddles and red curves and points are global minima.

Theorem 4.2 (Landscape Theorem).

The critical landscape is diffeomorphic to the space of pairs (\mathcal{I}, G) or (\mathcal{I}, O) with

- $\mathcal{I} \subset \{1, \dots, m\}$ of size $0 \leq l \leq k$,
- $G \in \mathbb{R}^{k \times l}$ with independent columns,
- $O \in \mathbb{R}^{k \times l}$ with orthonormal columns.

	W_2	W_1
\mathcal{L}	$U_{\mathcal{I}} G^+$	$G U_{\mathcal{I}}^{\top}$
\mathcal{L}_{π}	$U_{\mathcal{I}} (I_{\ell} + \lambda \Sigma_{\mathcal{I}}^{-2})^{-\frac{1}{2}} G^+$	$G (I_{\ell} + \lambda \Sigma_{\mathcal{I}}^{-2})^{-\frac{1}{2}} U_{\mathcal{I}}^{\top}$
\mathcal{L}_{σ}	$U_{\mathcal{I}} (I_{\ell} - \lambda \Sigma_{\mathcal{I}}^{-2})^{\frac{1}{2}} O^{\top}$	$O (I_{\ell} - \lambda \Sigma_{\mathcal{I}}^{-2})^{\frac{1}{2}} U_{\mathcal{I}}^{\top}$

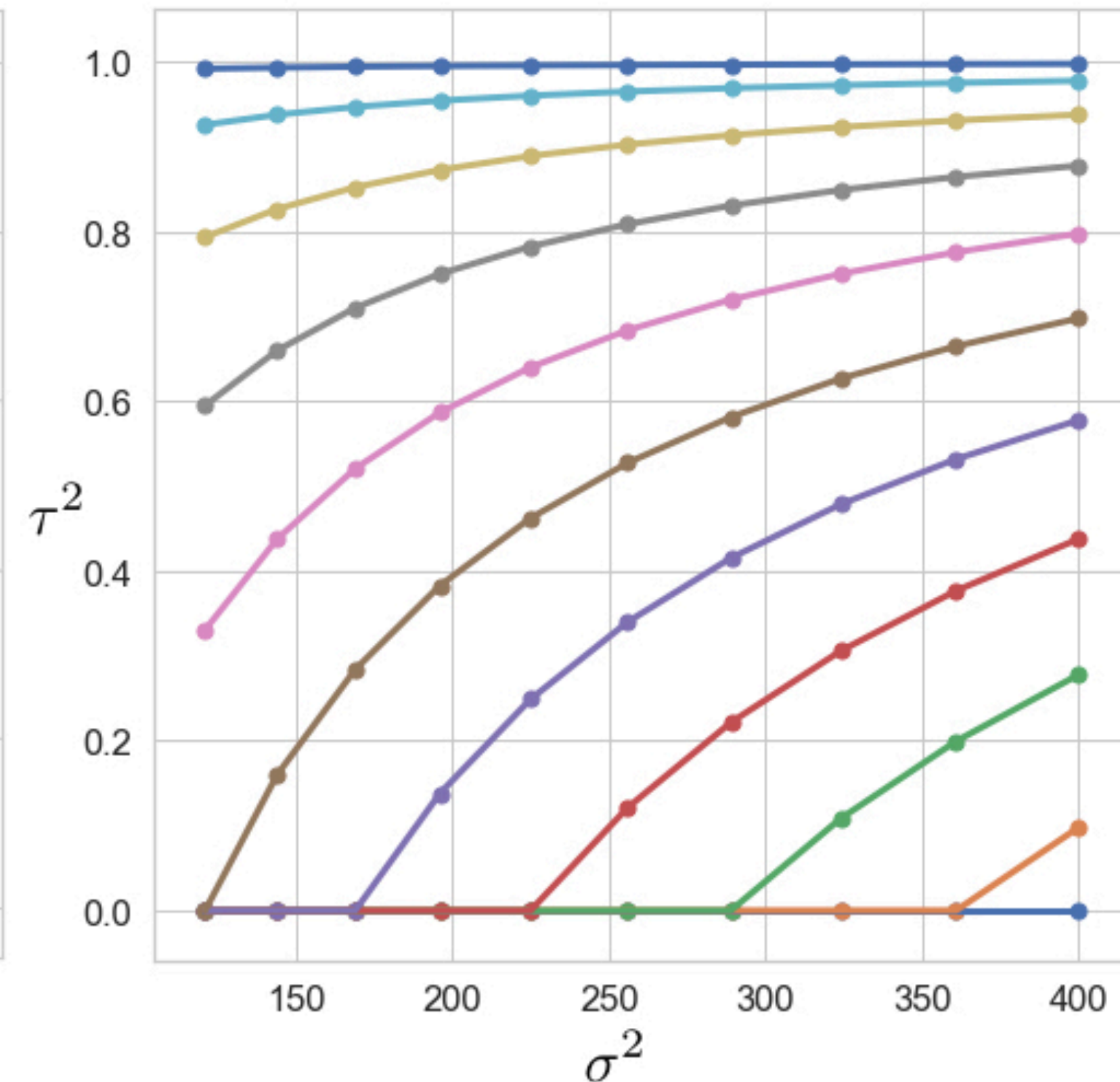
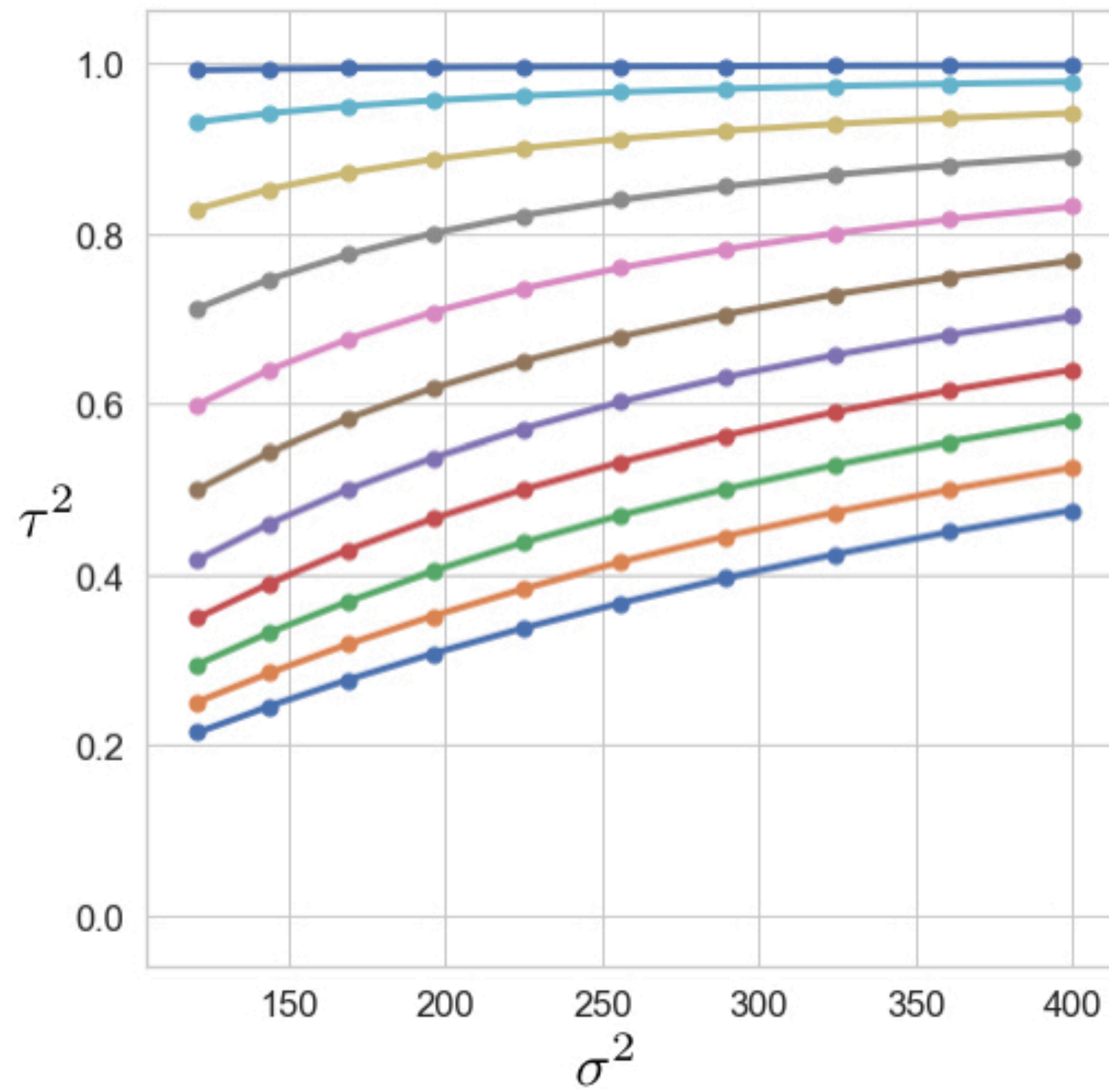
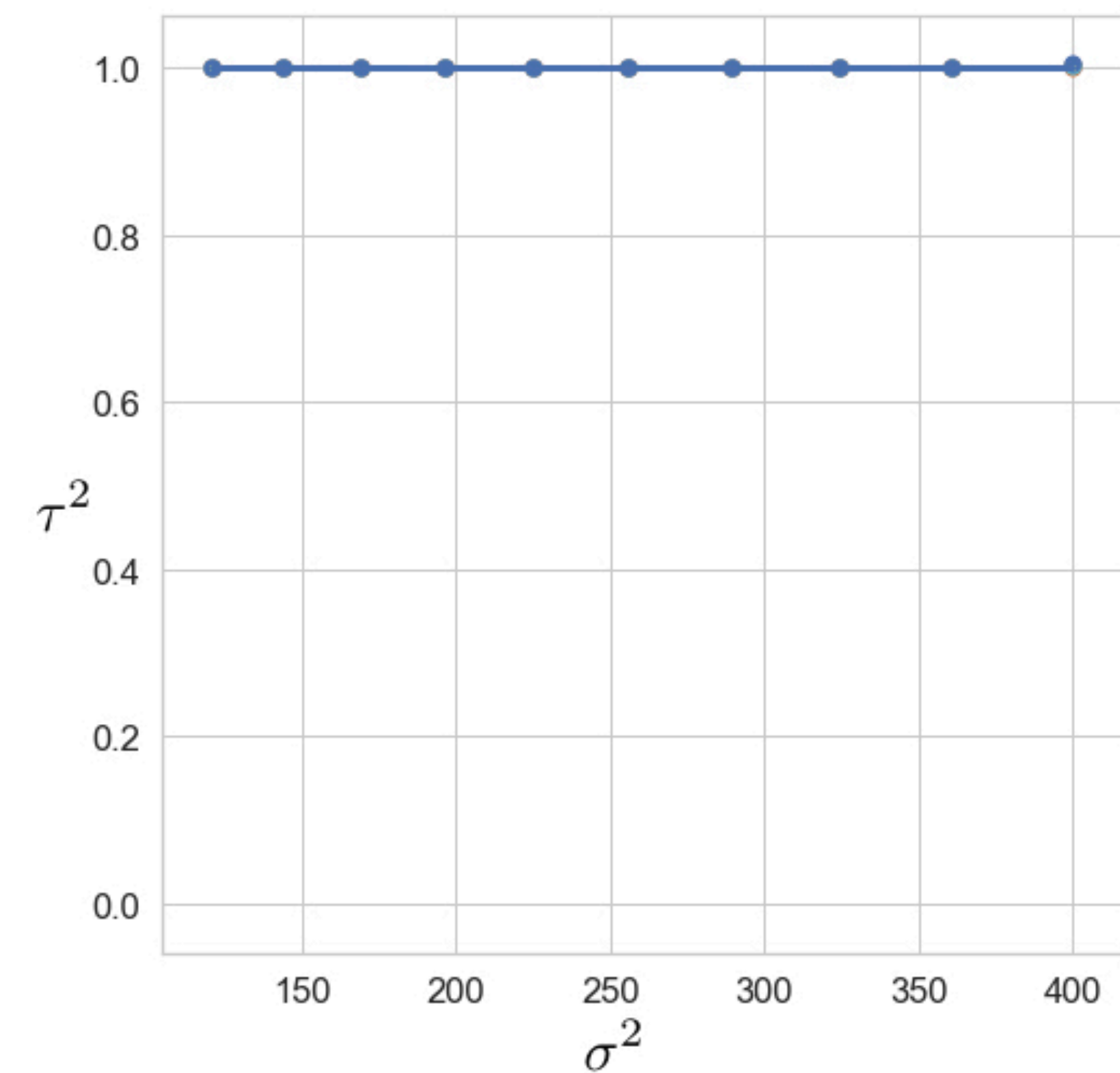
Shrinkage



None

Ridge

pPCA-like



Unregularized

Product

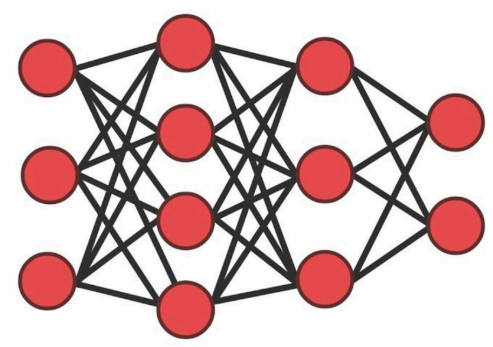
Sum

Theorem 3.1 (pPCA Theorem). *With $\sigma^2 = \lambda$, the critical points of*

$$\mathcal{L}_\sigma^0(W_0) = \mathcal{L}_\sigma(W_0^\top (XX^\top)^{-\frac{1}{2}}, (XX^\top)^{-\frac{1}{2}} W_0)$$

coincide with the critical points of pPCA.

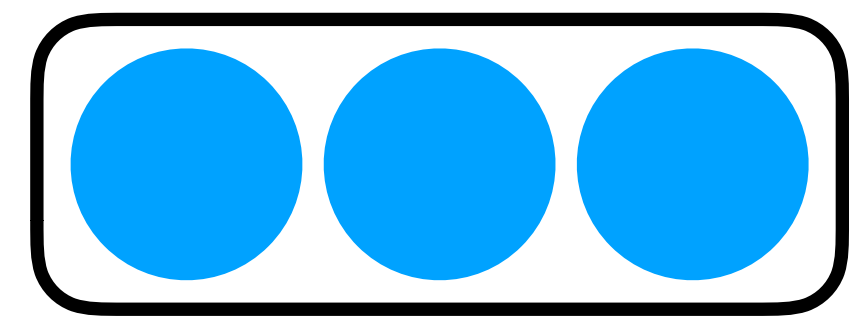
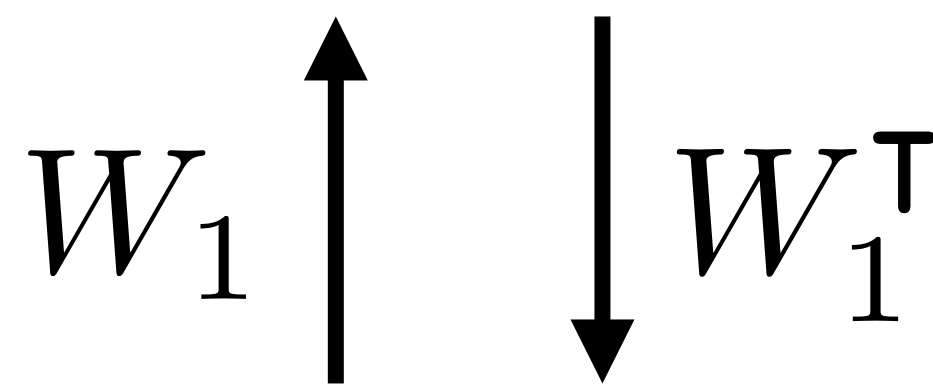
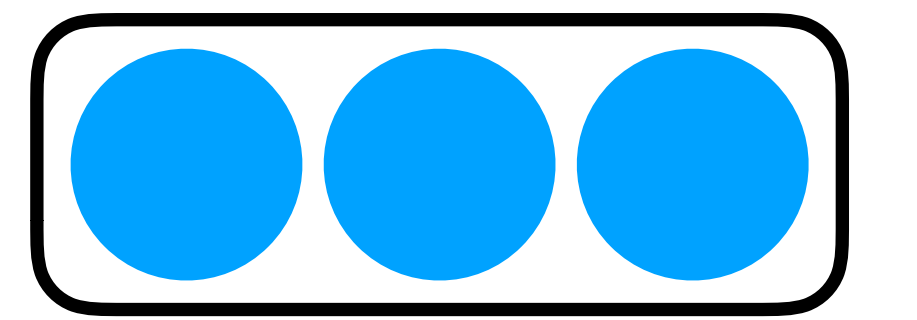
Bayesian \mathcal{L}_σ	pPCA
$W_1, W_2^\top \sim \mathcal{N}_{k \times m}(0, \lambda^{-1})$	$z_i \sim \mathcal{N}_k(0, 1)$
$\varepsilon_i \sim \mathcal{N}_m(0, 1)$	$\varepsilon_i \sim \mathcal{N}_m(0, \sigma^2)$
$x_i = W_2 W_1 x_i + \varepsilon_i$	$x_i = W_0 z_i + \varepsilon_i$



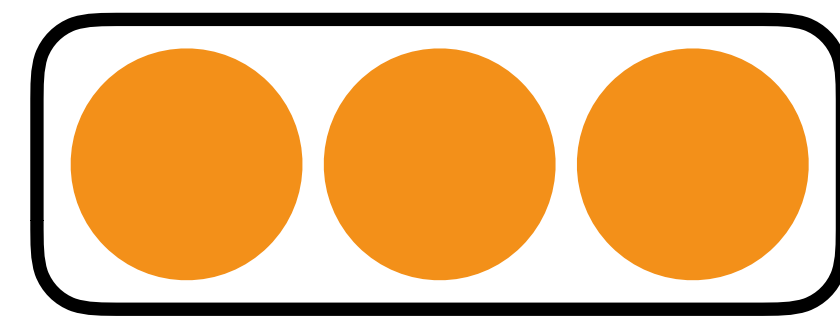
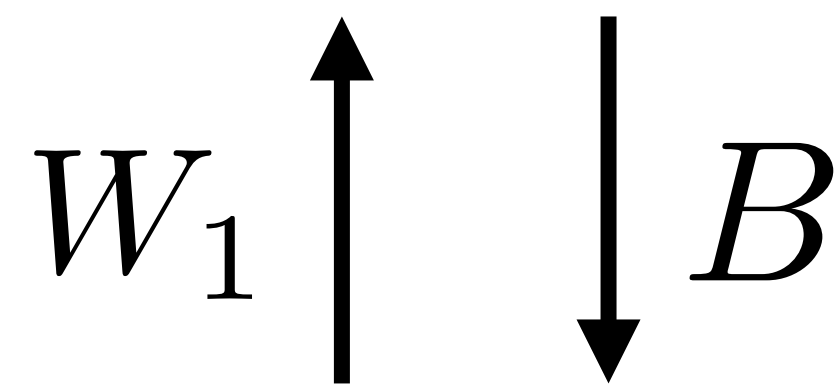
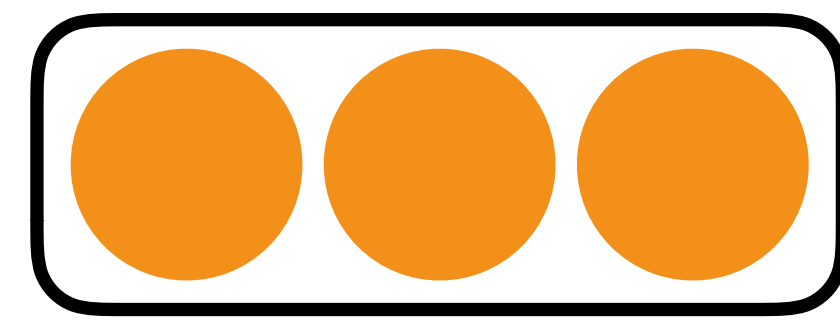
Prediction in artificial neural networks is inspired by the brain.



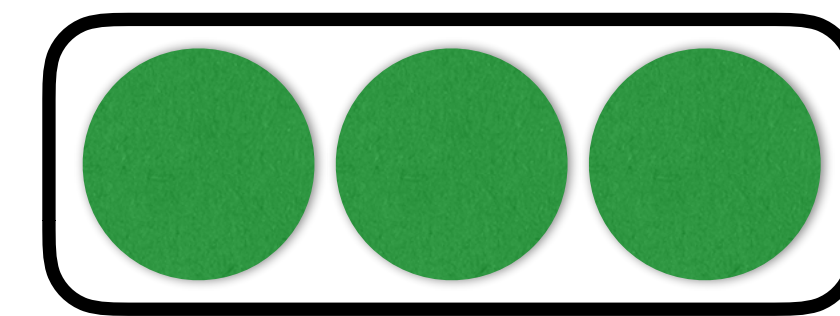
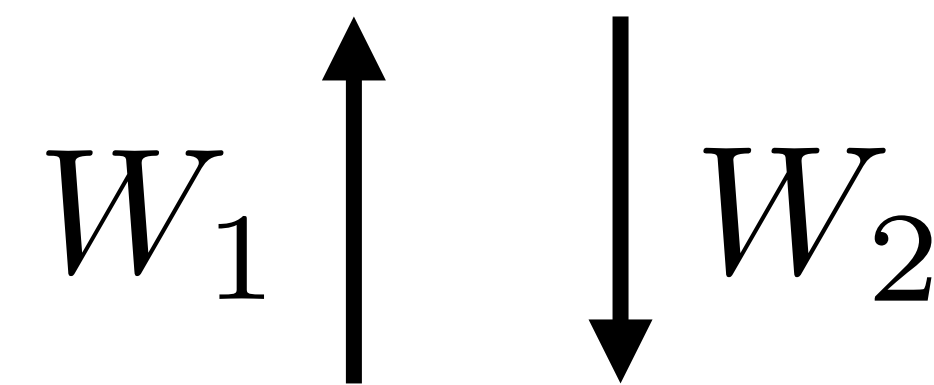
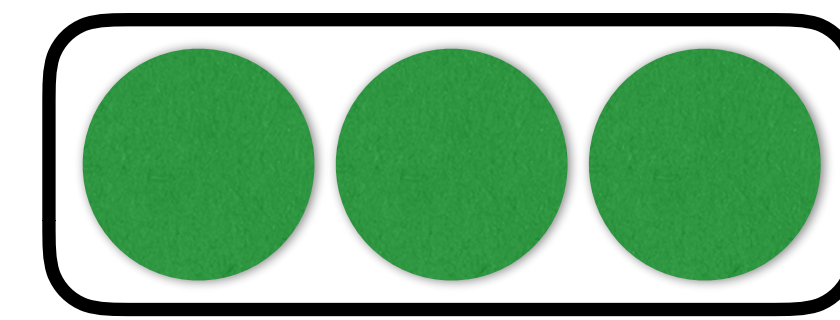
Is *learning* in the brain inspired by artificial neural networks?



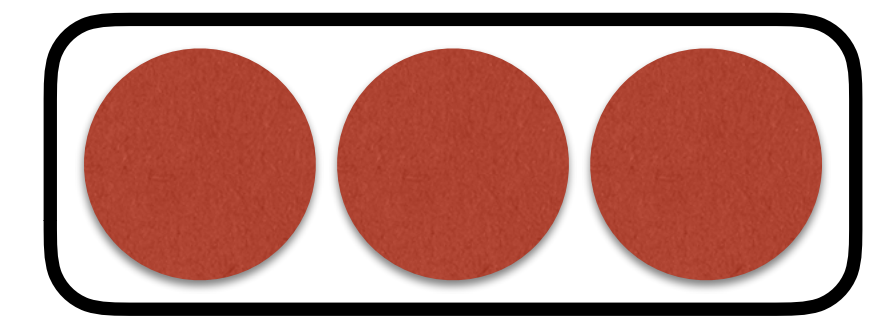
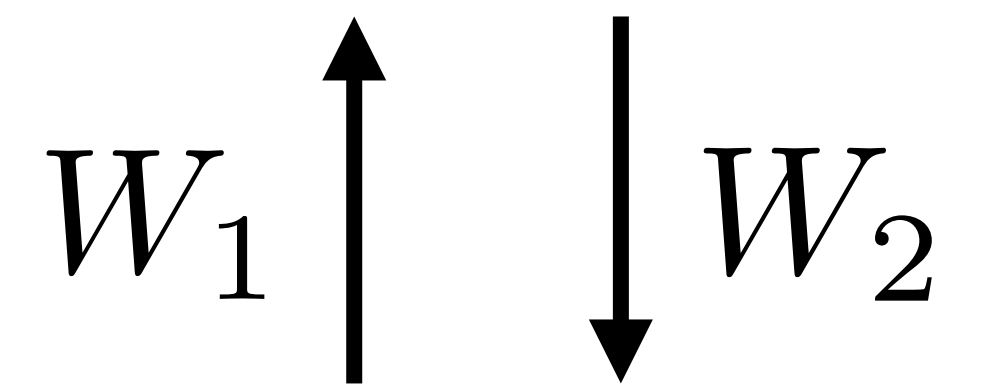
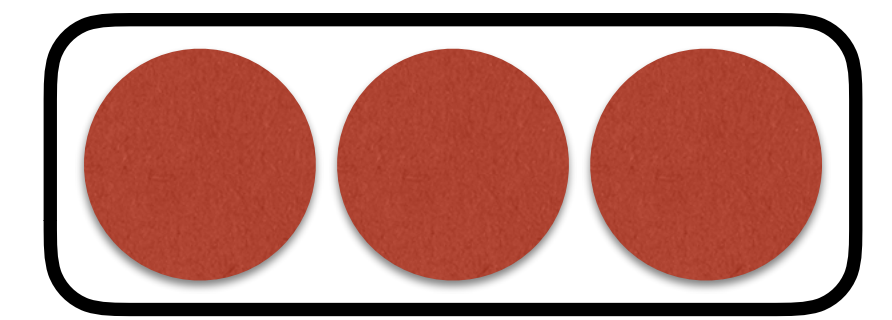
Backprop



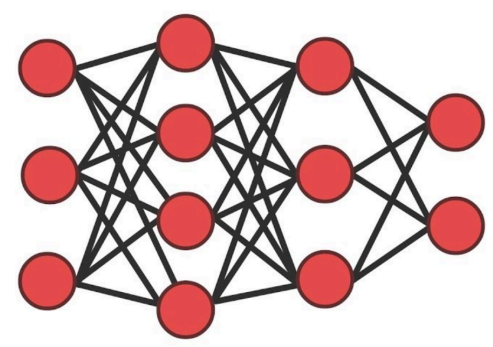
Feedback Alignment



Information Alignment



Symmetric Alignment



Prediction in artificial neural networks is inspired by the brain.



Is *learning* in the brain inspired by artificial neural networks?

Prediction
Supervised

$$\mathcal{L}_{\text{pred}} = \|y - W_1 W_0 x\|^2$$

Representation
Unsupervised

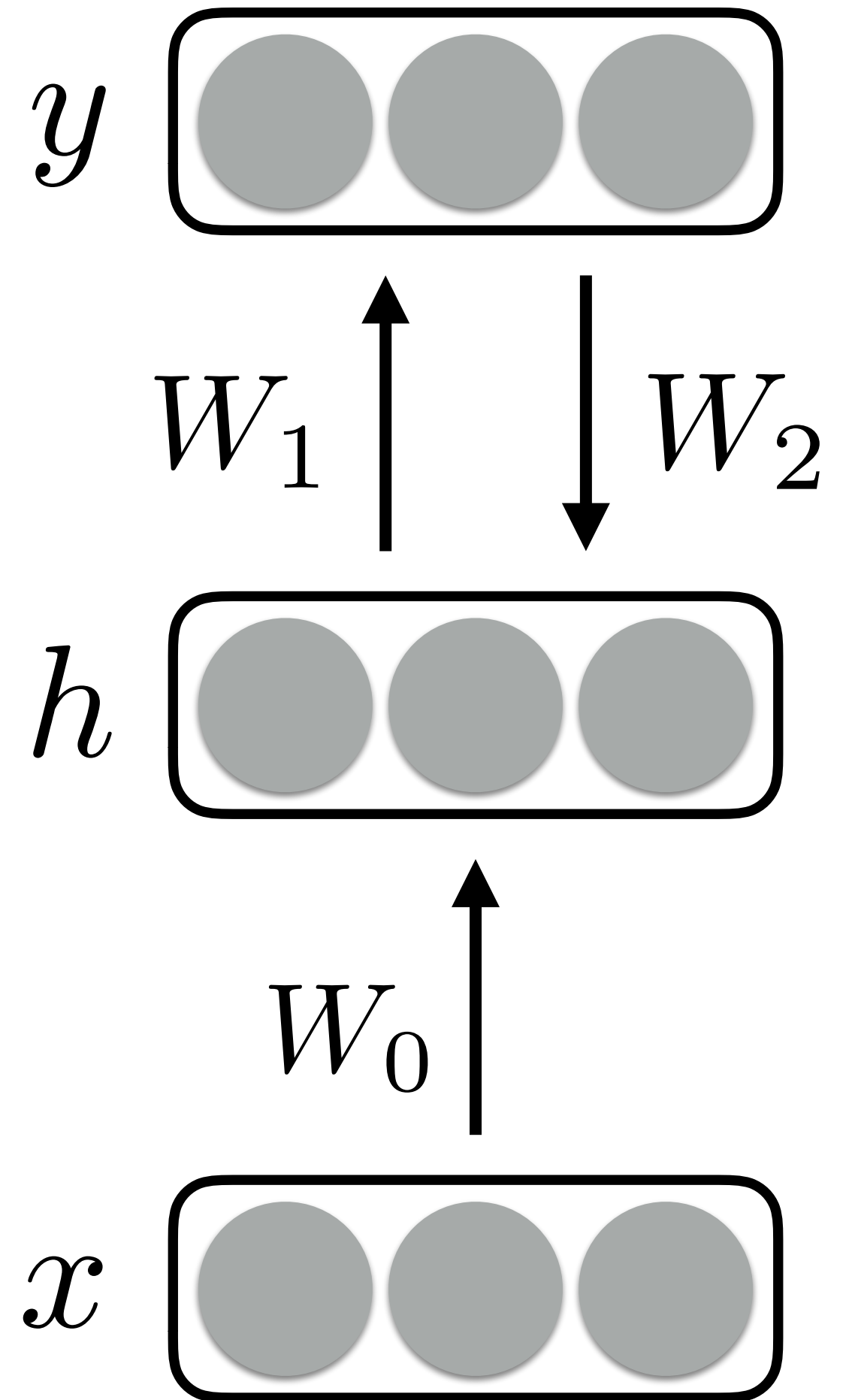
$$\mathcal{L}_{\text{info}} = \|h - W_2 W_1 h\|^2$$

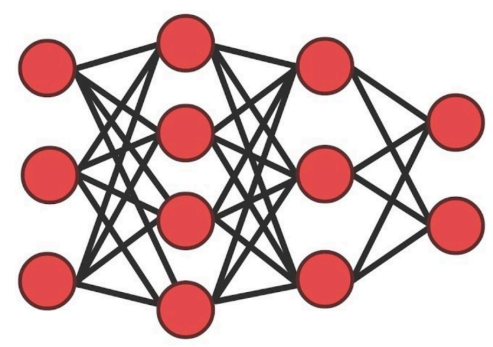
Efficiency
Sparsity

$$\mathcal{L}_{\text{reg}} = \|W_1\|^2 + \|W_2\|^2$$

Self-amplification
Feedback control

$$\mathcal{L}_{\text{self}} = -2\text{tr}(W_2 W_1)$$





Prediction in artificial neural networks is inspired by the brain.

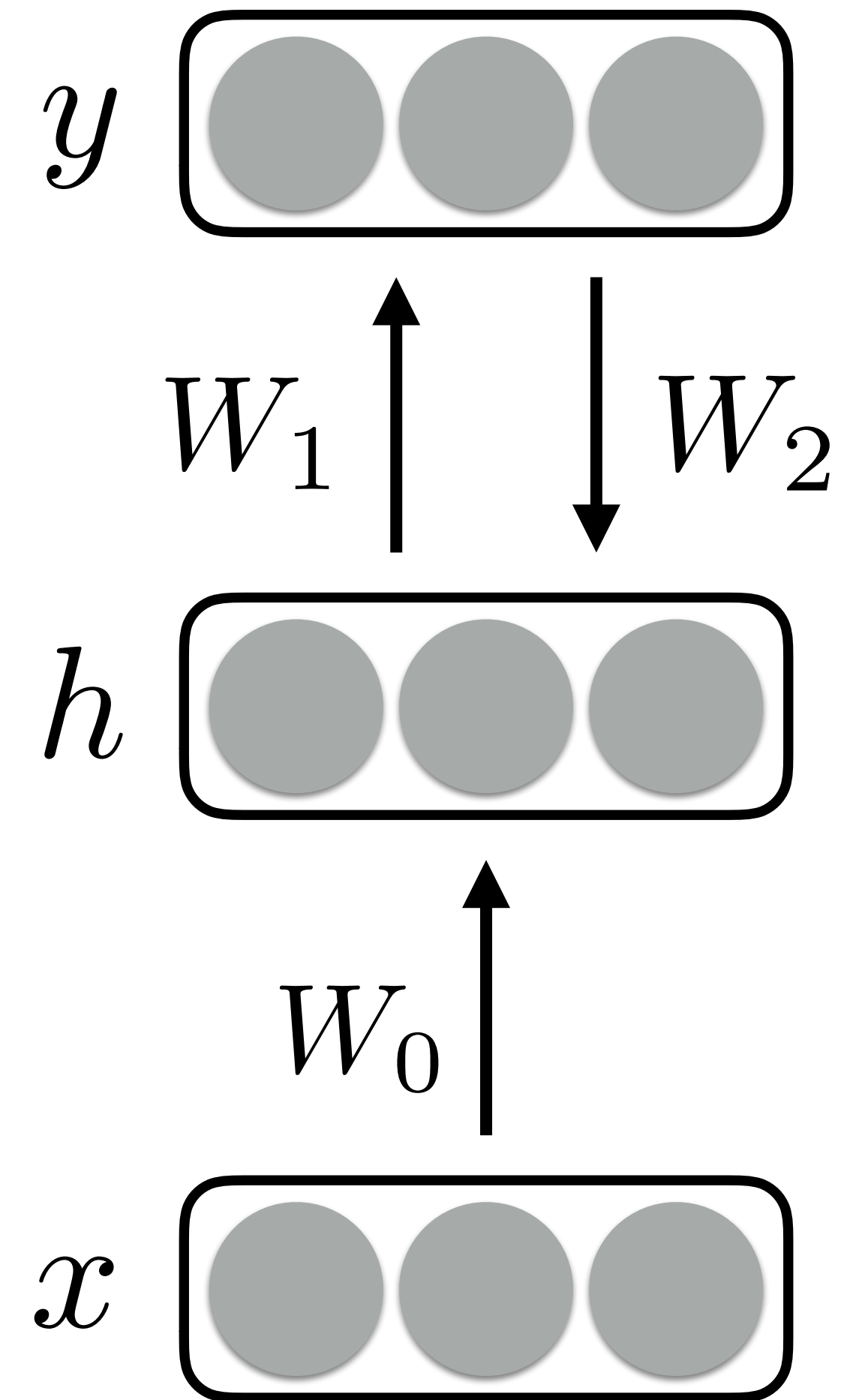


Is *learning* in the brain inspired by artificial neural networks?

$$\mathcal{L}_{\text{BP}} = \mathcal{L}_{\text{pred}} + \mathcal{L}_{\text{reg}}$$

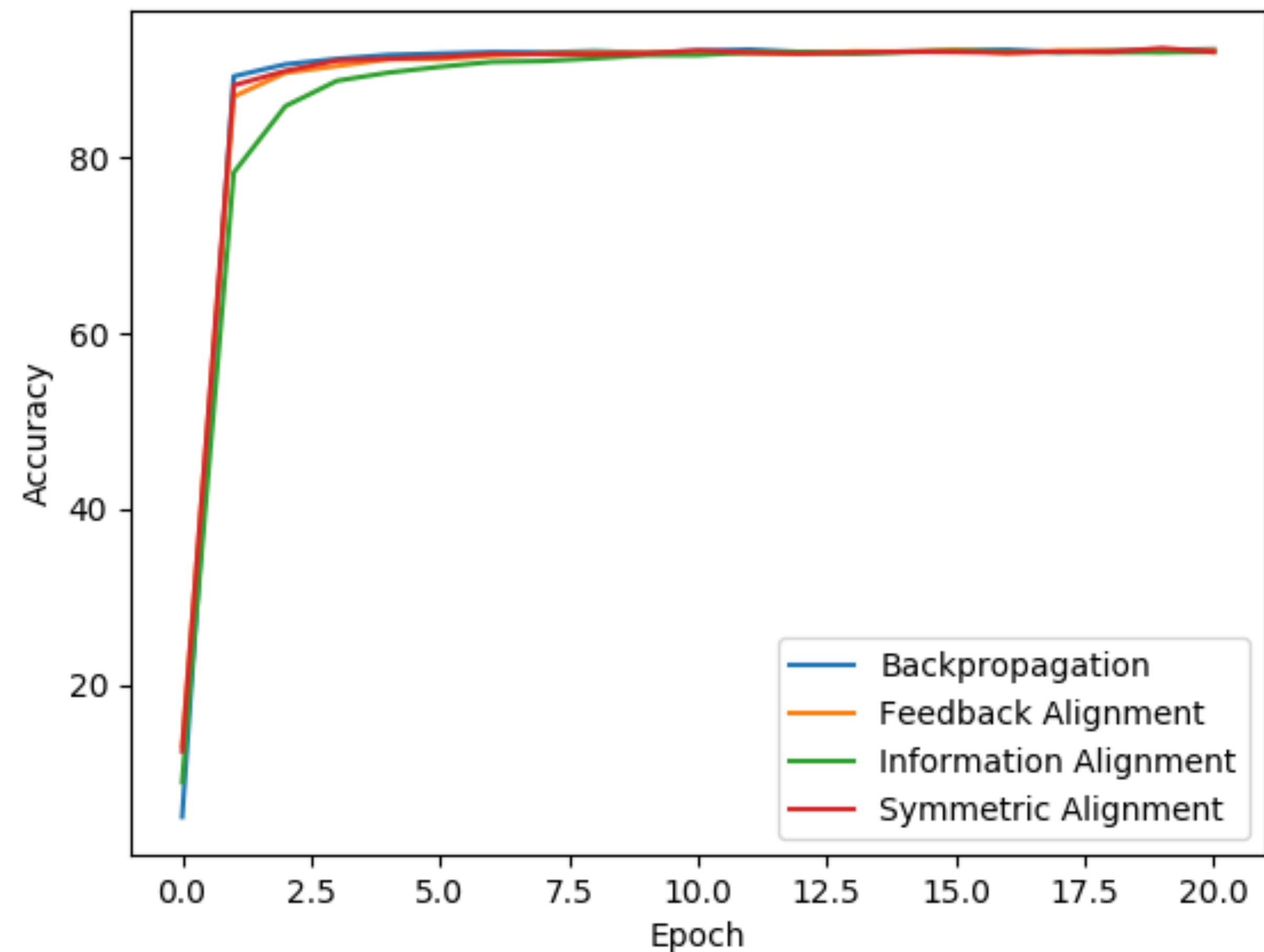
$$\mathcal{L}_{\text{IA}} = \mathcal{L}_{\text{pred}} + \mathcal{L}_{\text{info}} + \mathcal{L}_{\text{reg}}$$

$$\mathcal{L}_{\text{SA}} = \mathcal{L}_{\text{pred}} + \underbrace{\mathcal{L}_{\text{self}} + \mathcal{L}_{\text{reg}}}_{\|W_2 - W_1^T\|^2}$$

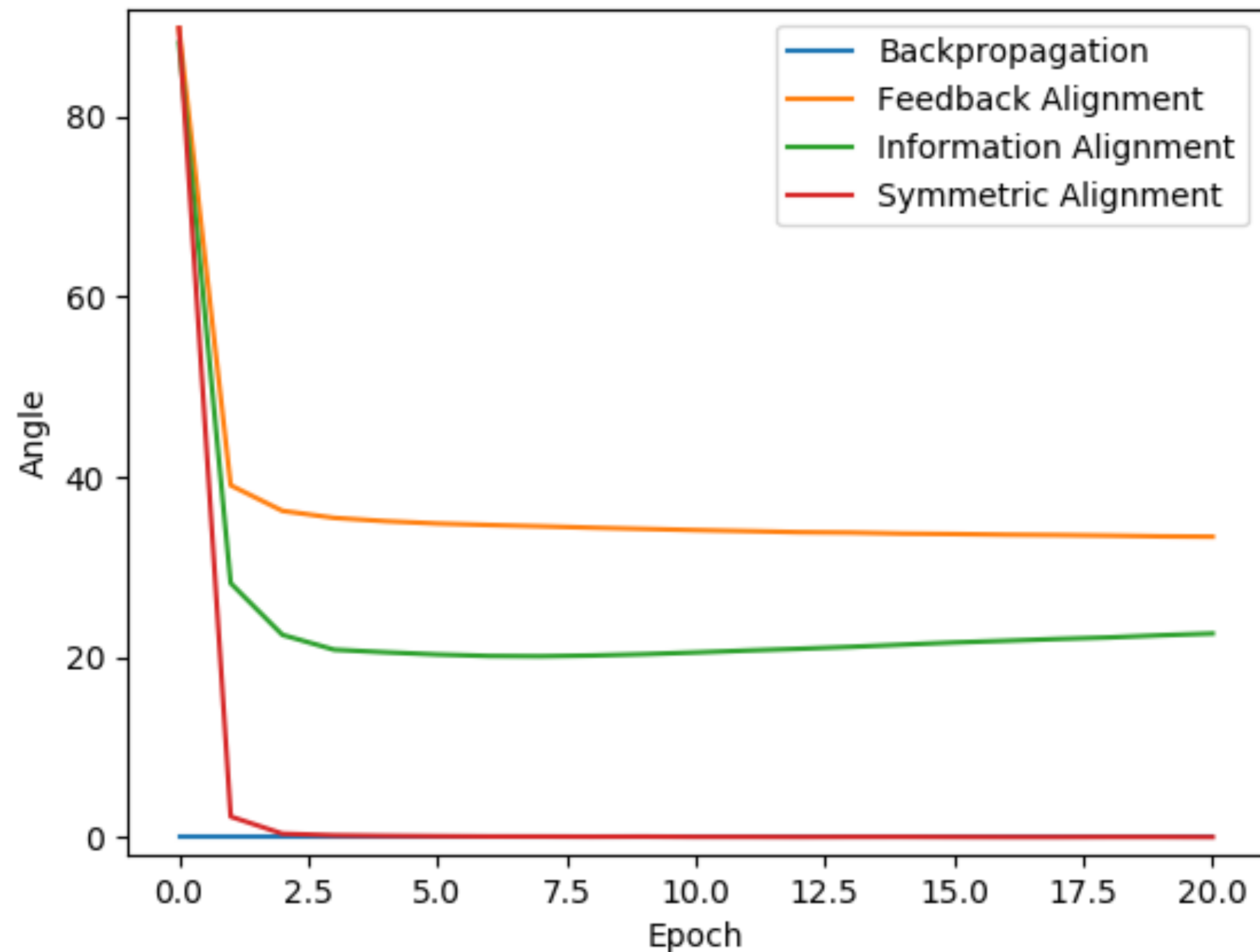


Linear, 1 hidden layer

MNIST Test Accuracy

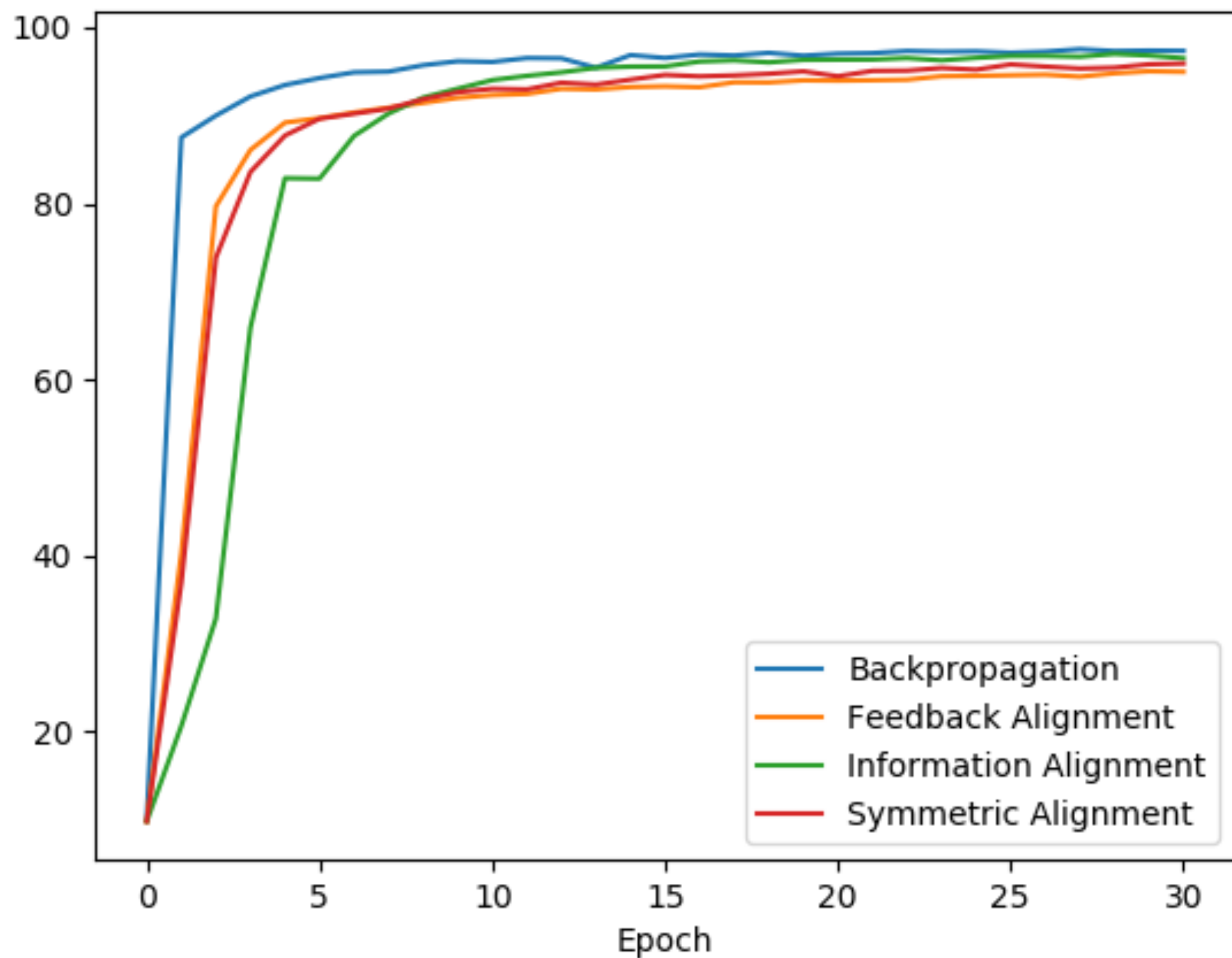


Angle Alignment



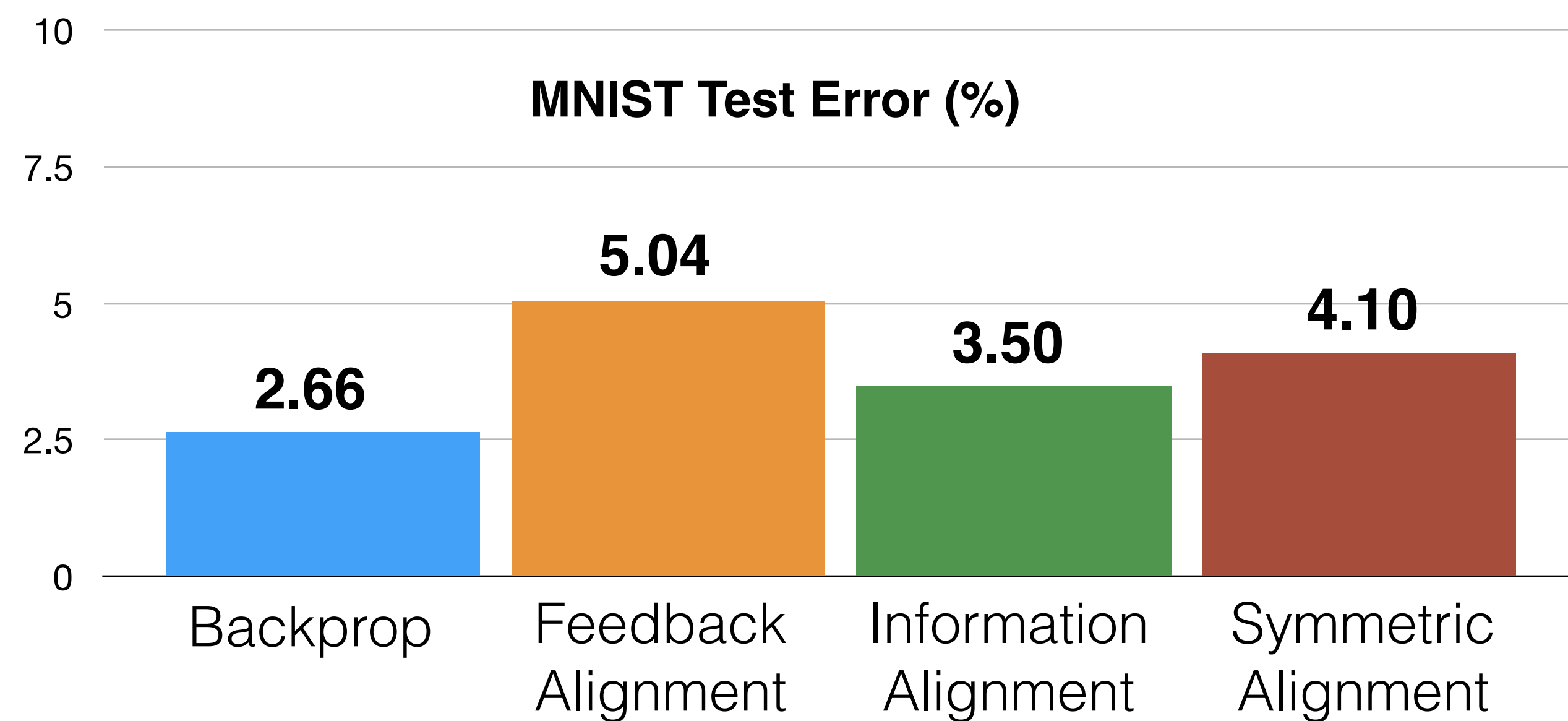
ReLU, 6 hidden layers

MNIST Test Accuracy

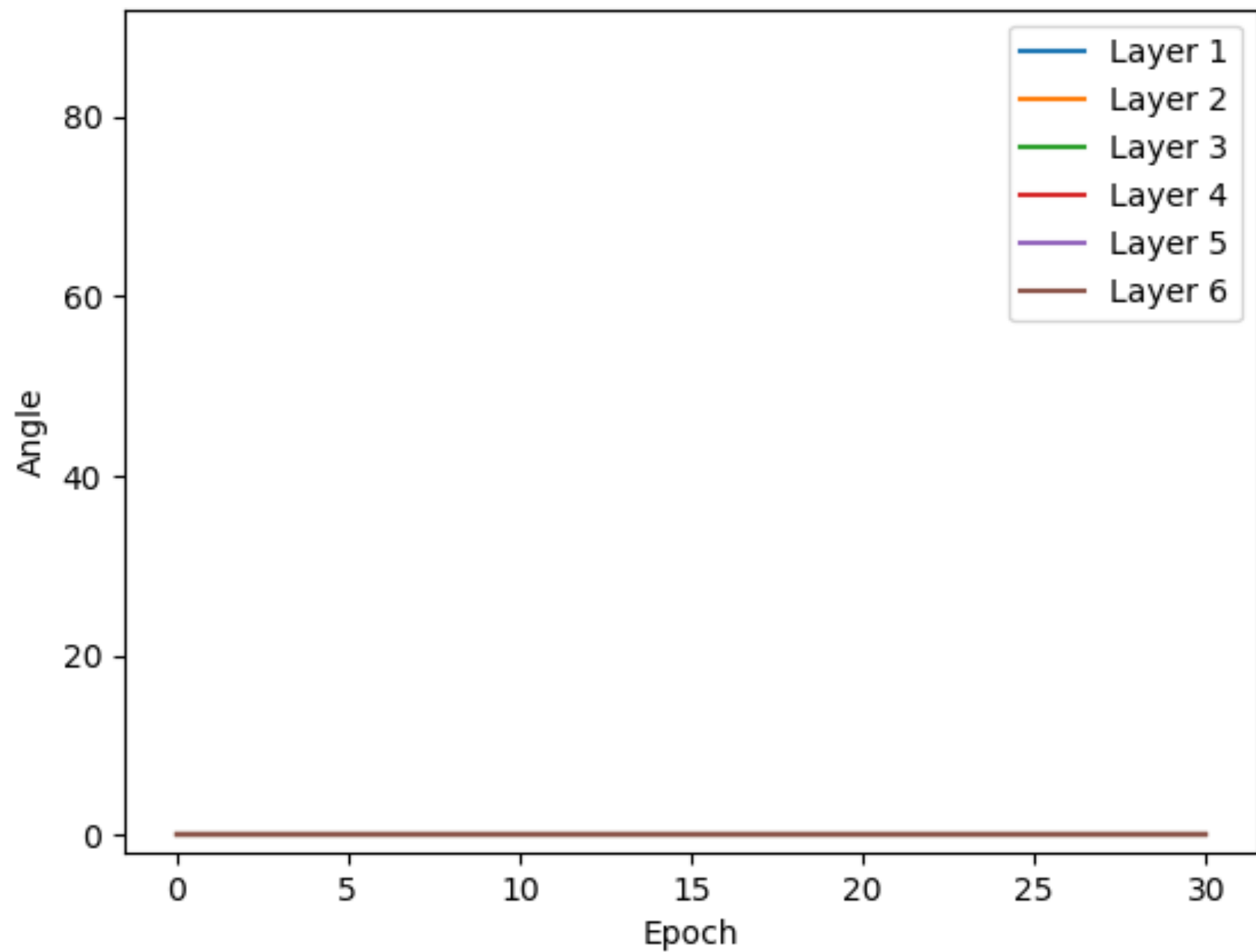


- $5e-5$ learning rate
- 1000 batch size
- 30 epochs
- ReLu activations
- 784 - 256 - 256 - 256 - 256 - 256 - 10 architecture
- Lambda = $1e-6$

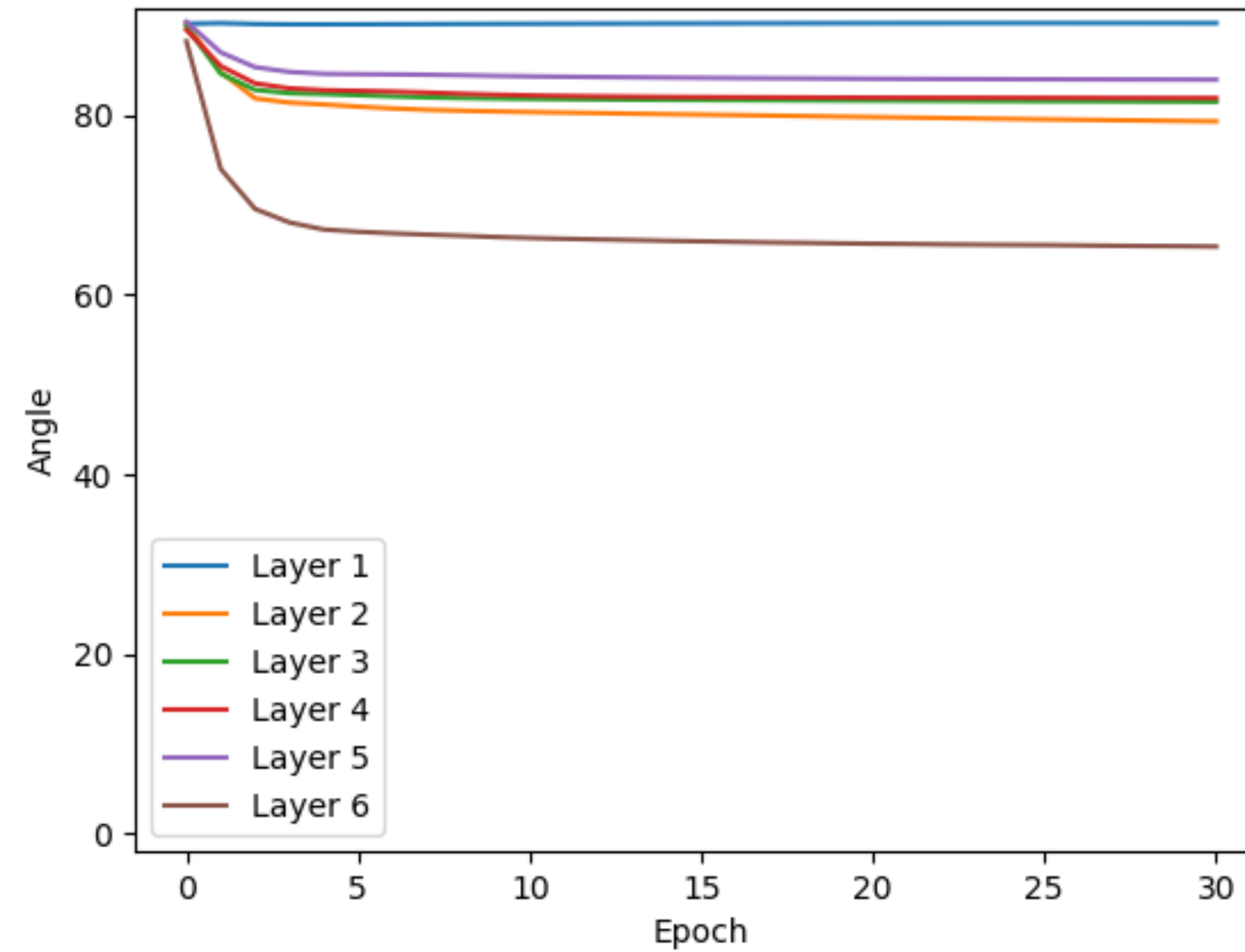
MNIST Test Error (%)



BP

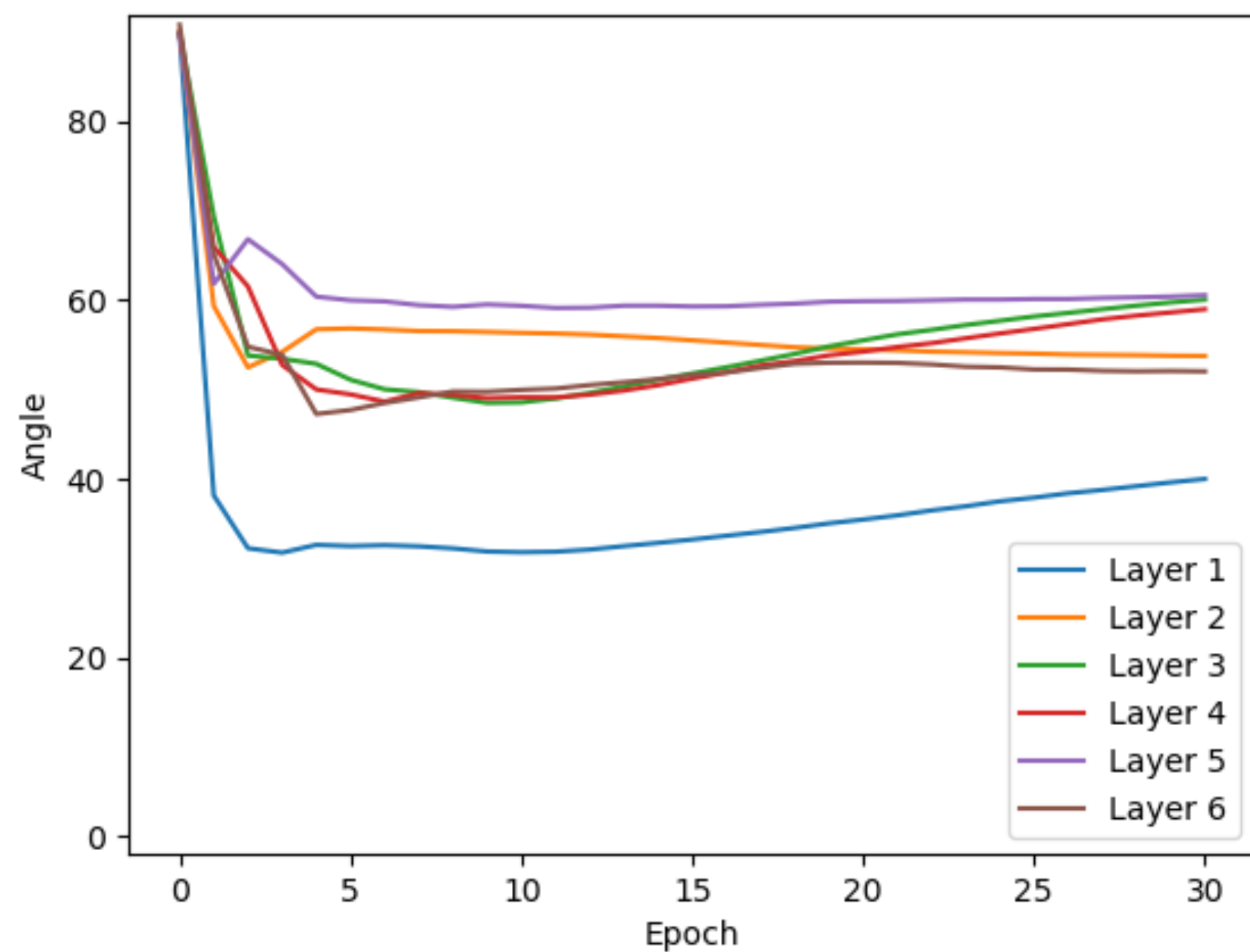


FA

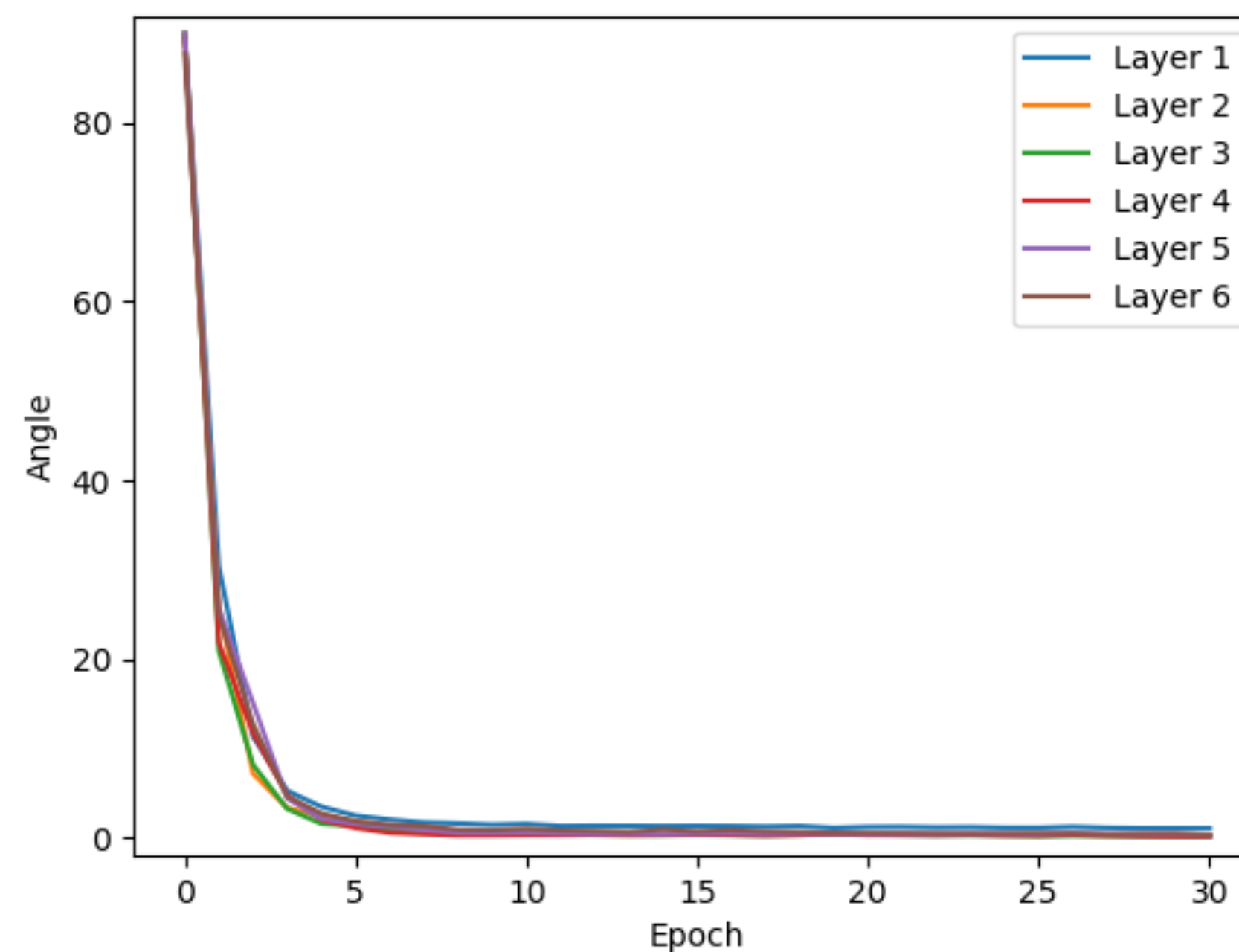


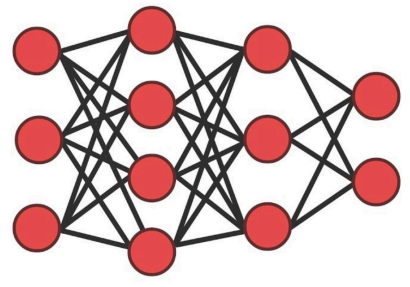
Alignment

IA



SA





in silico

- architecture search for less-biologically-implausible learning at scale
- adversarial robustness, information bottleneck, multitask learning

π

in puro

- LAE from X to Y: closely related to PLS, CCA
- derive and study continuous flows on linear alignment networks

Σ

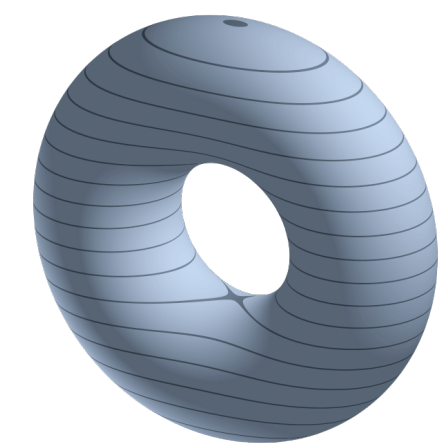
in algo

- bio-inspired take down of randomized SVD

in vivo



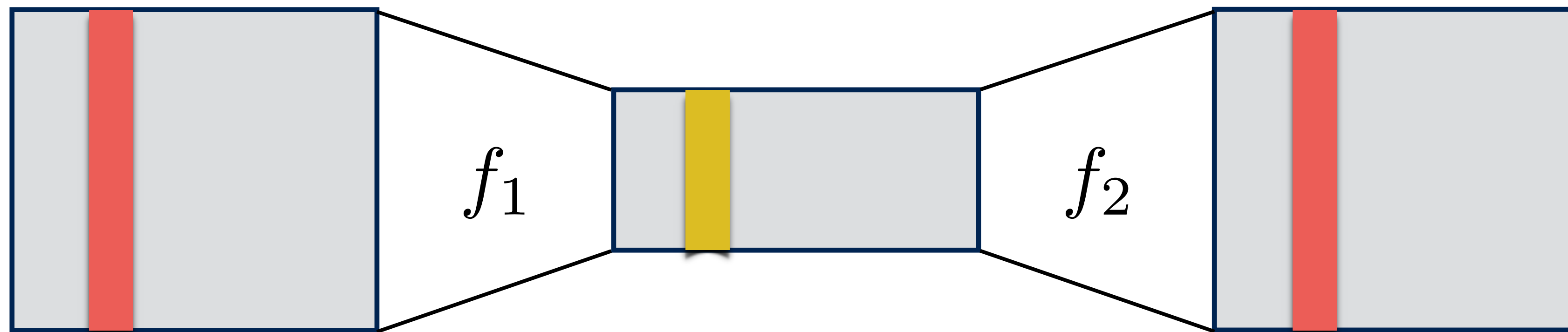
- dynamics of representation teleportation in development and learning
- neural implementation factored through genomic / molecular ontogeny in space and time



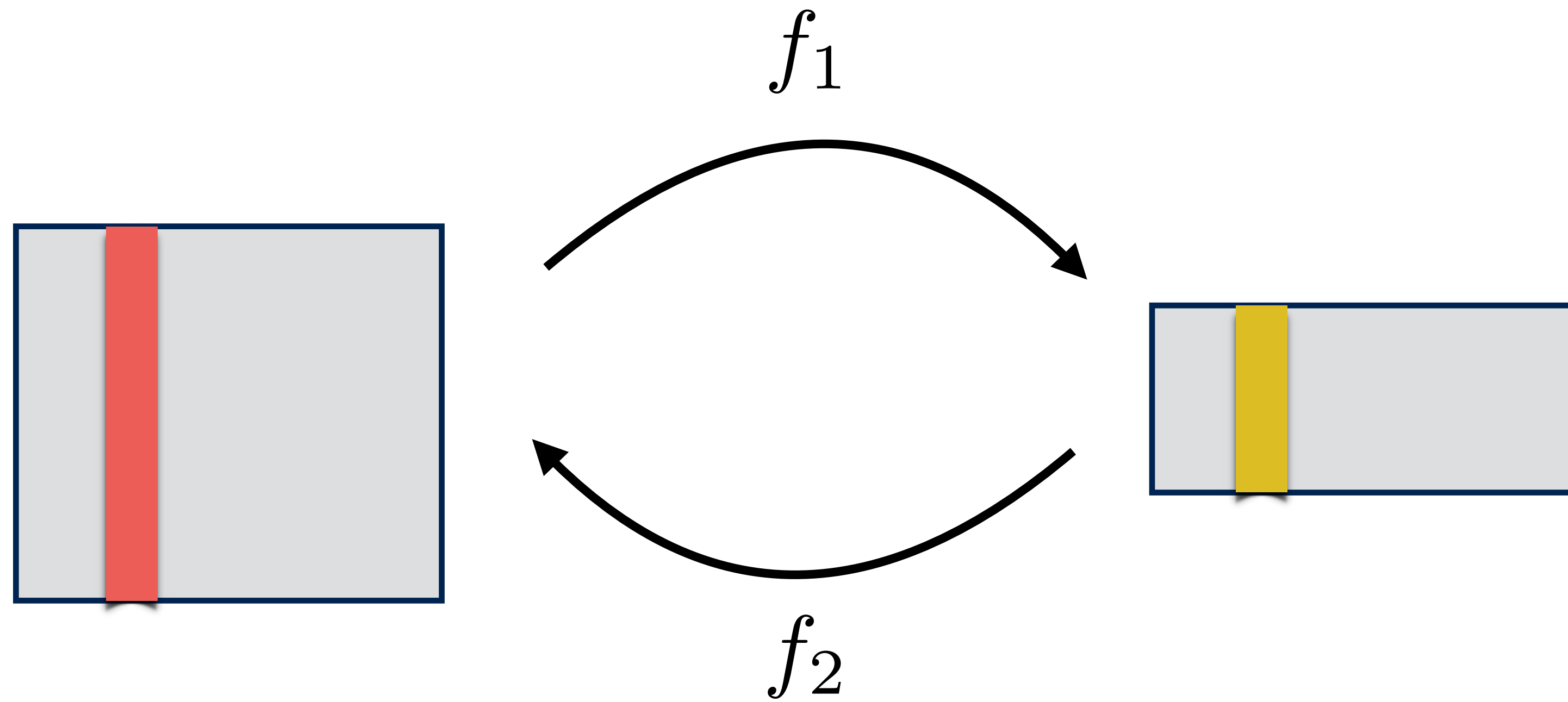
in categorico / physico

- *ground learning in functors between algebra, topology, and geometry (Morse homology, ensemble and consensus learning, TQFT)*

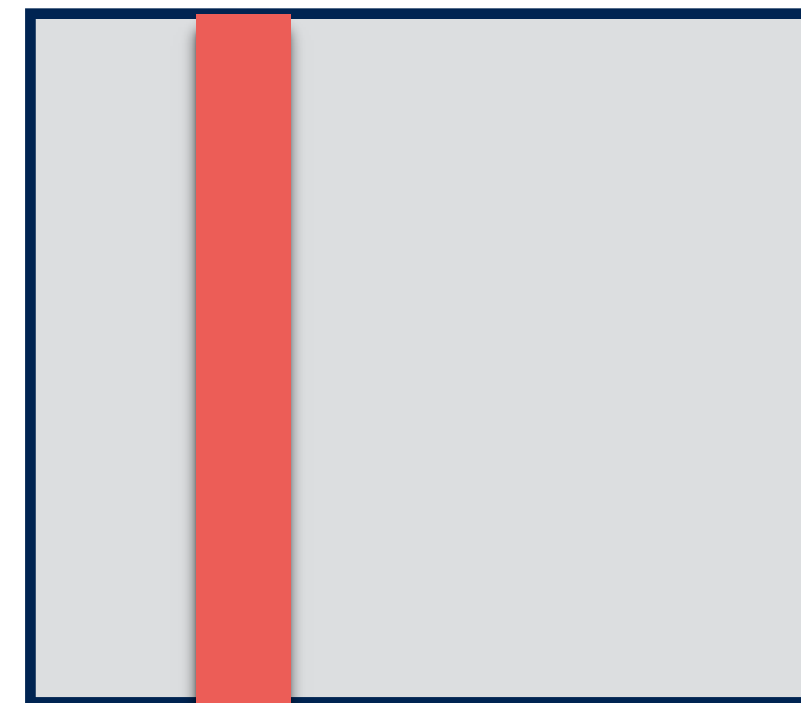
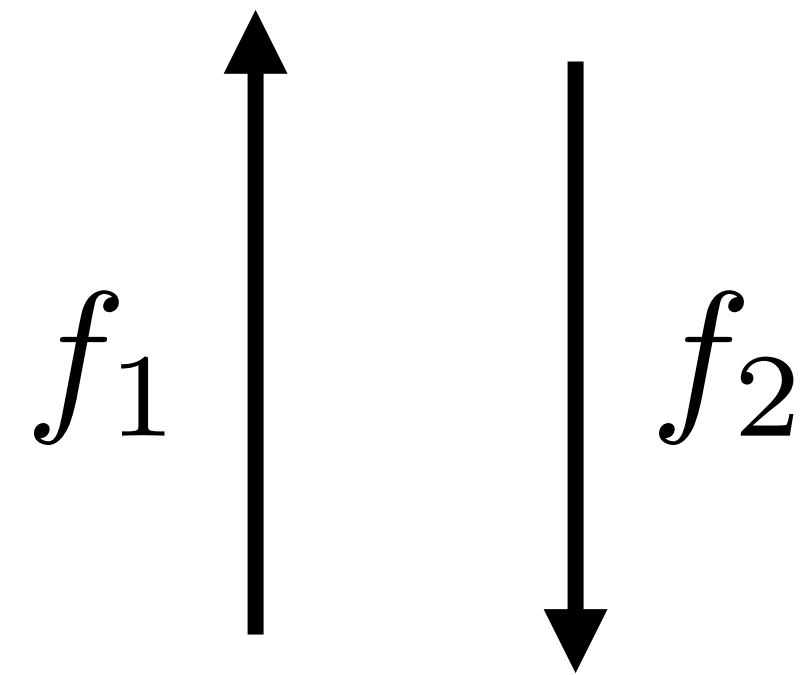
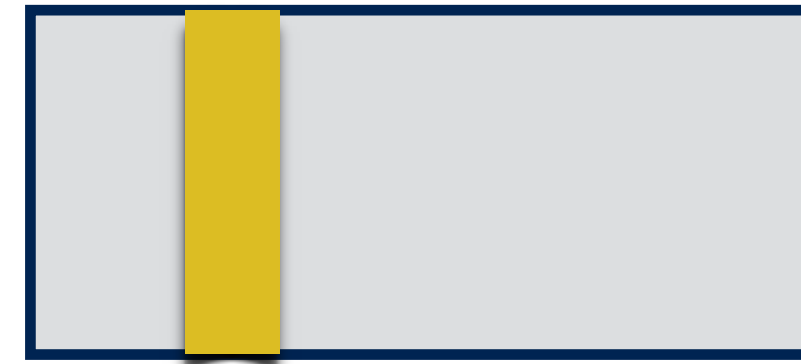
Compressor



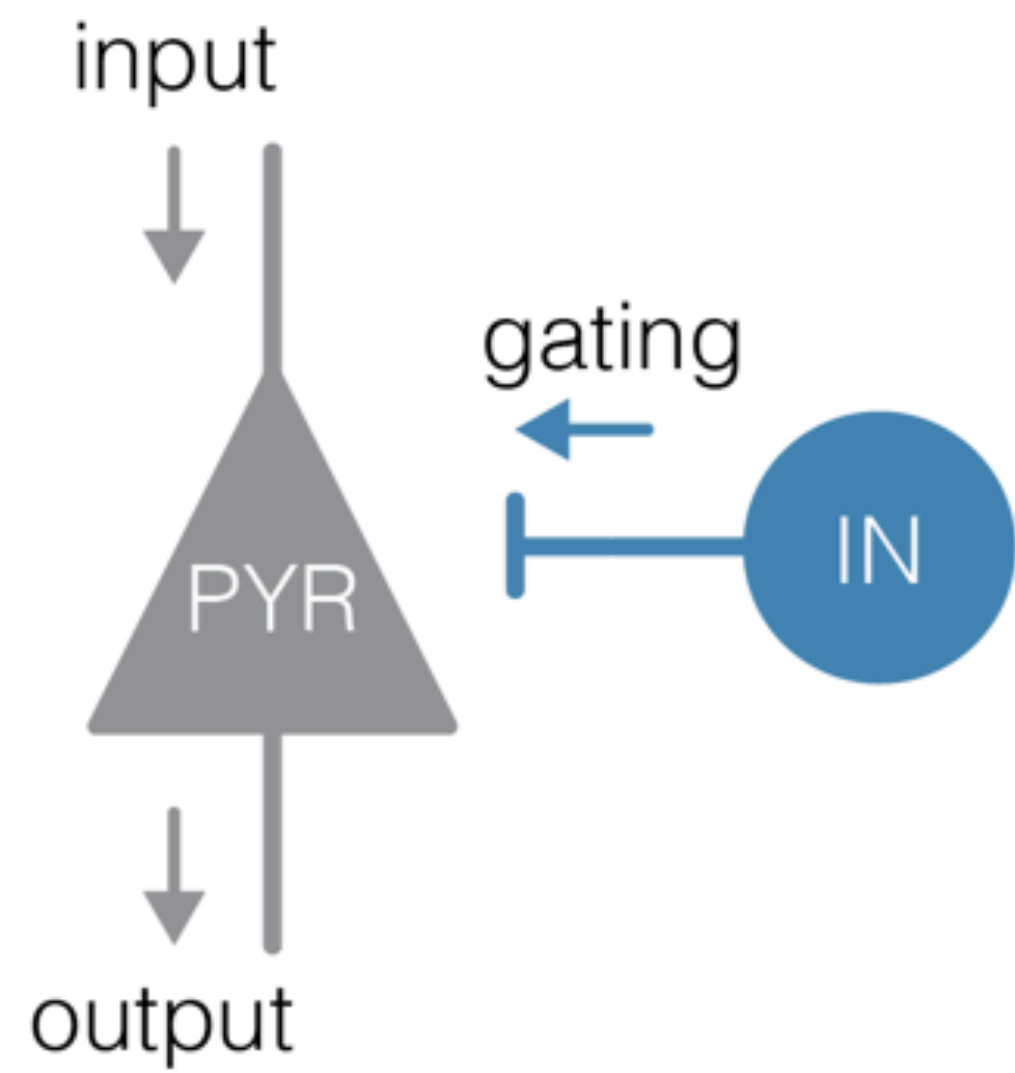
Transporter



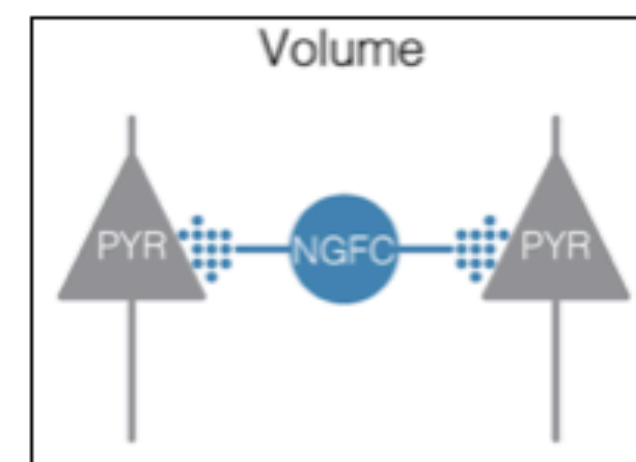
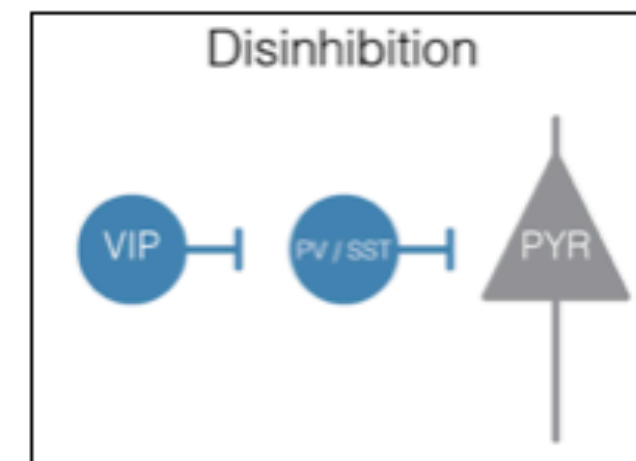
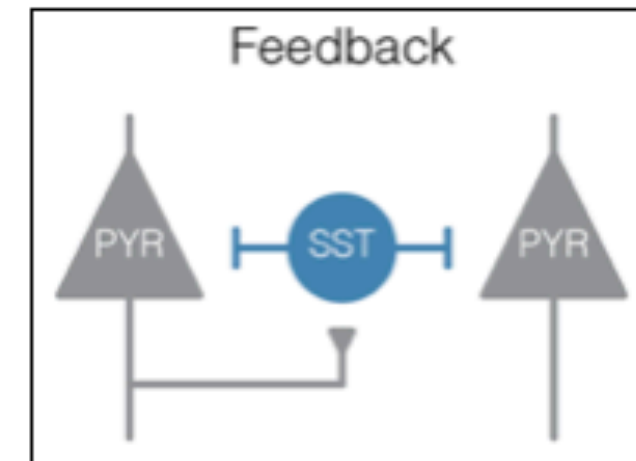
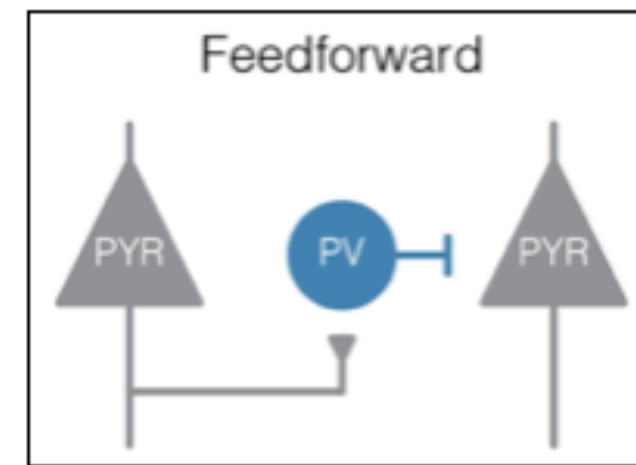
Stabilizer



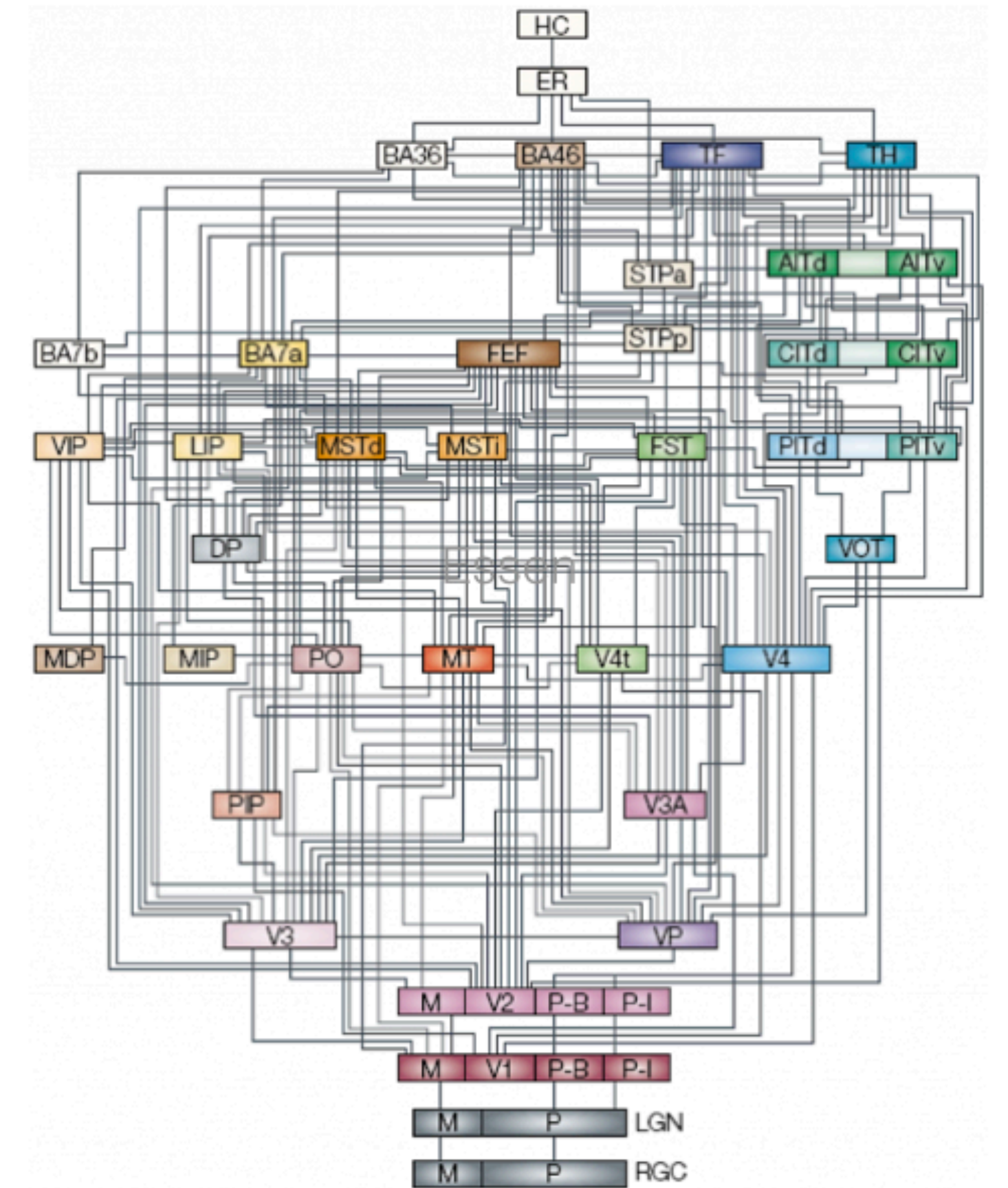
Cellular interactions



Circuit motifs



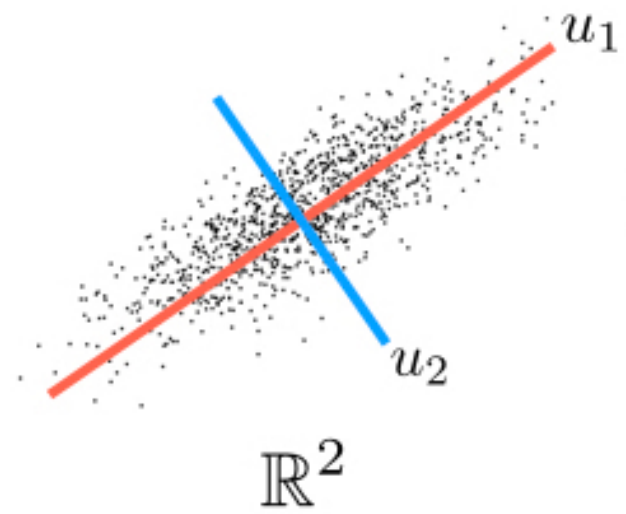
Connectivity across regions



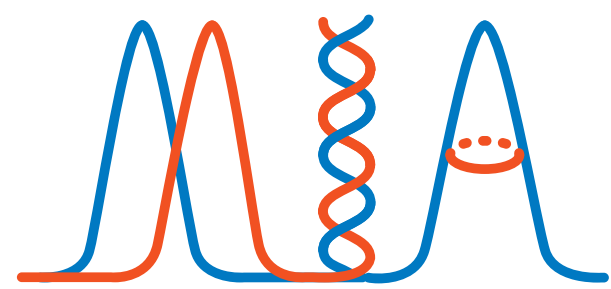
MAY THE
MORSE
BE WITH YOU



Thanks!



LAE: github.com/danielkunin/Regularized-Linear-Autoencoders



homepage: broadinstitute.org/mia

playlist: bit.ly/2l18EvO

overview: youtube.com/watch?v=gWcFJiYZNZ0



homepage: hail.is/about.html

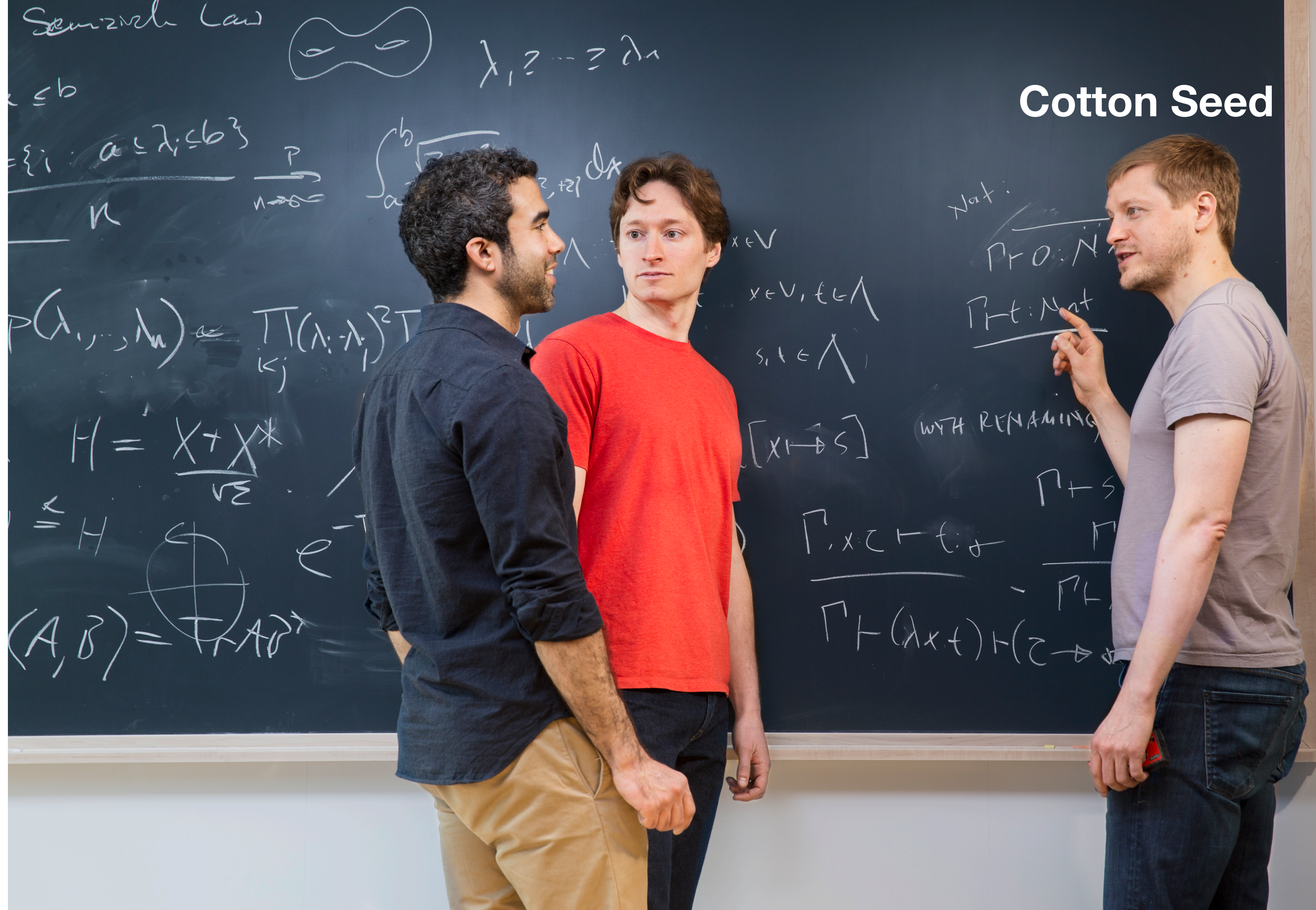
code: github.com/hail-is/hail



article: thecrimson.com/article/2019/2/28/broad-institute-scrut/

Hail Slides

Cotton Seed



Riemann zeta function $\zeta(s) = \sum_{n=1}^{\infty} \frac{1}{n^s}$



$\{i : a \leq \lambda_i \leq b\}$
 $\int_a^b \sqrt{1-x^2} dx$

$$P(\lambda_1, \dots, \lambda_n) \propto \prod_{1 \leq j < k \leq n} (\lambda_j - \lambda_k)^2$$

$$H = \frac{X + X^*}{\sqrt{2}}$$

$$(A, B) = \text{Tr}(AB^T)$$

$x \in V$
 $x \in V, t \in \Lambda$
 $s, t \in \Lambda$

$[x \mapsto s]$

Nat.
 $\Gamma \vdash 0 : \text{Nat}$
 $\Gamma \vdash t : \text{Nat}$

WITH RENAMING

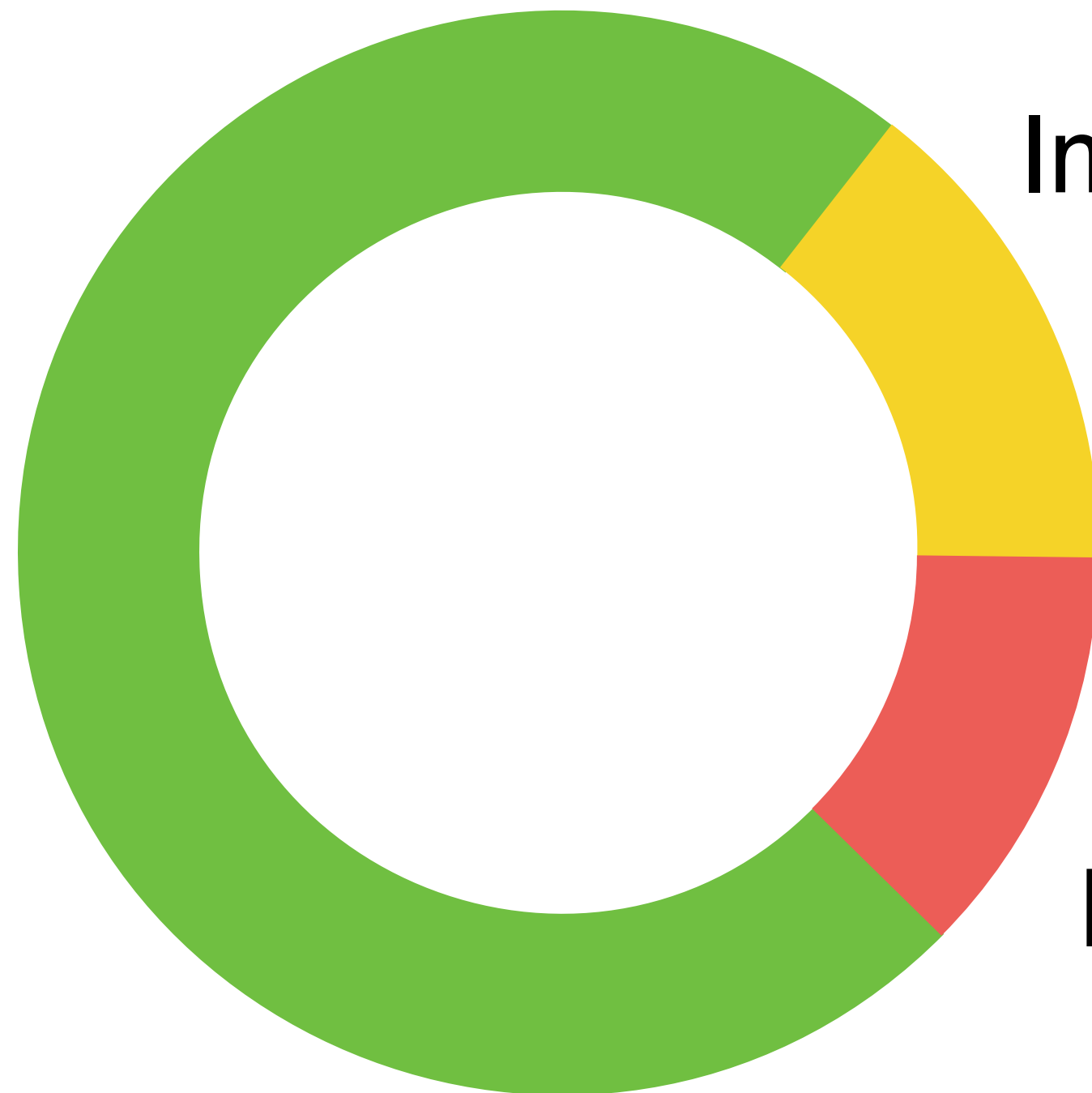
$\Gamma, x \vdash t : \sigma$
 $\Gamma \vdash (\lambda x. t) \vdash (\tau \rightarrow \sigma)$

haixl





Scientific
Reasoning



Implementation

Runtime

Statistical Genetics Tools

Custom Python/R scripts

- Filter genotypes with bad allele balance
- Call *de novo* variants
- Compute transmission disequilibrium
- Dominance-encoded GWAS
- Gene count permutation tests

Doesn't Scale

PLINK

- Detect sample duplicates or ID swaps
- Call Mendelian violations
- Relatedness
- GWAS
- ...

Doesn't Scale

SNPSift

- Genotype concordance

Doesn't Scale

EMMAX

- Sequence kernel association test
- Rare variant burden test

Scales-ish

bcftools

- Split multiallelic variants
- Filter on GQ, AD, PASS

Doesn't Scale

tabix

- Subset VCFs to intervals

Doesn't Scale

vcffilterjdk

- Filter variants

Doesn't Scale

bedtools

- Interval annotation

Doesn't Scale

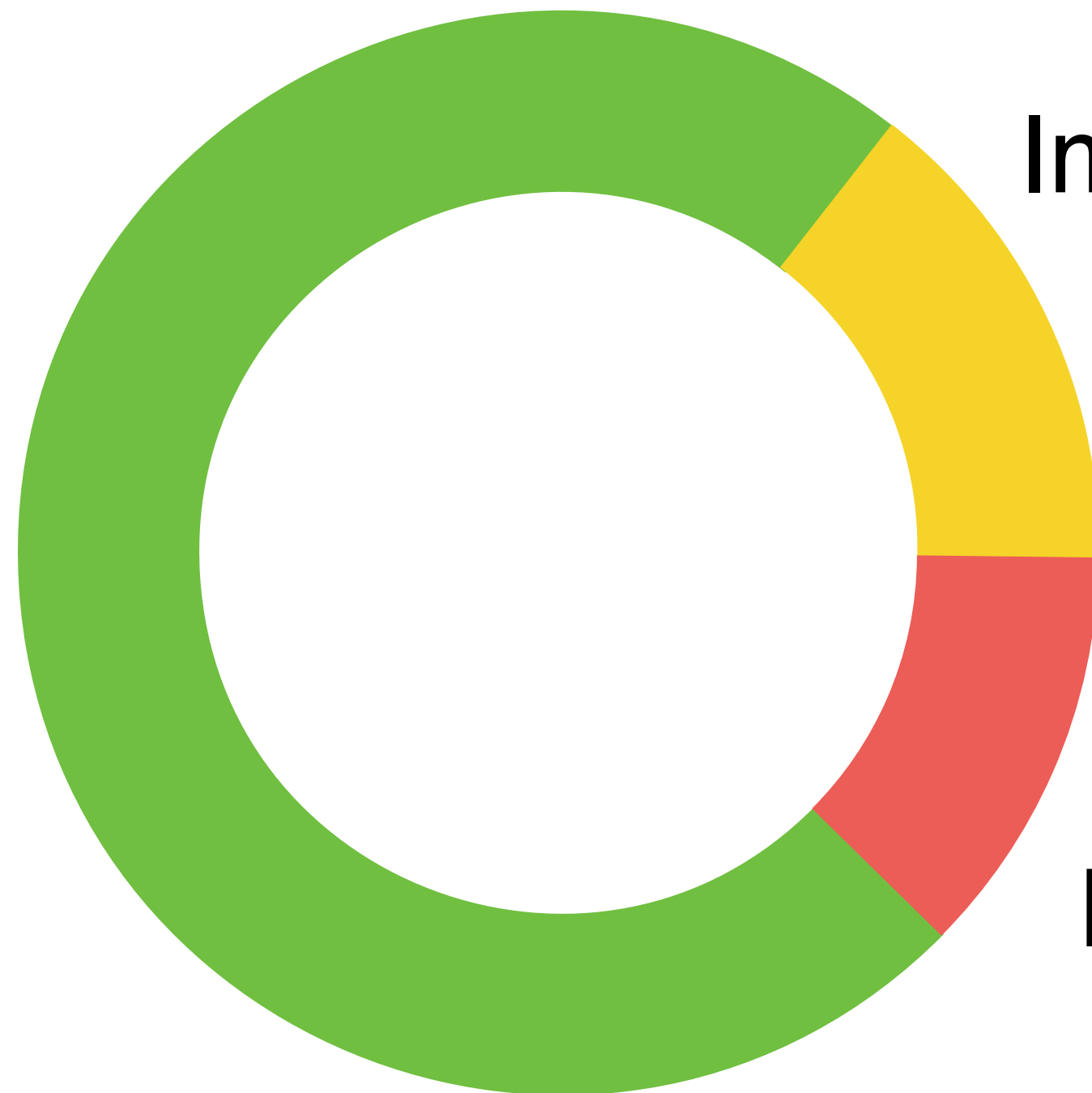
Eigenstrat

- PCA

Doesn't Scale



Scientific
Reasoning

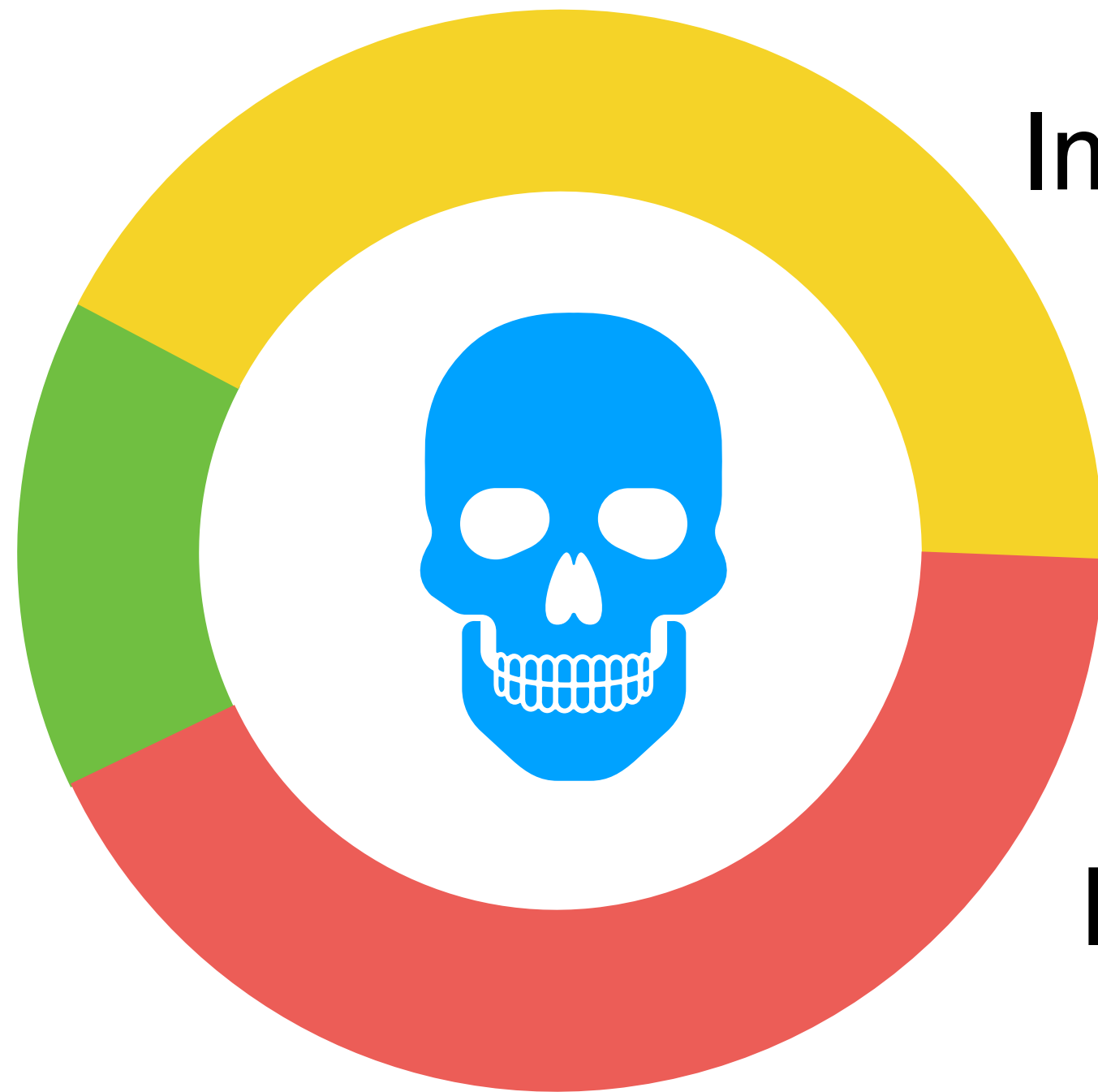


Implementation

Runtime

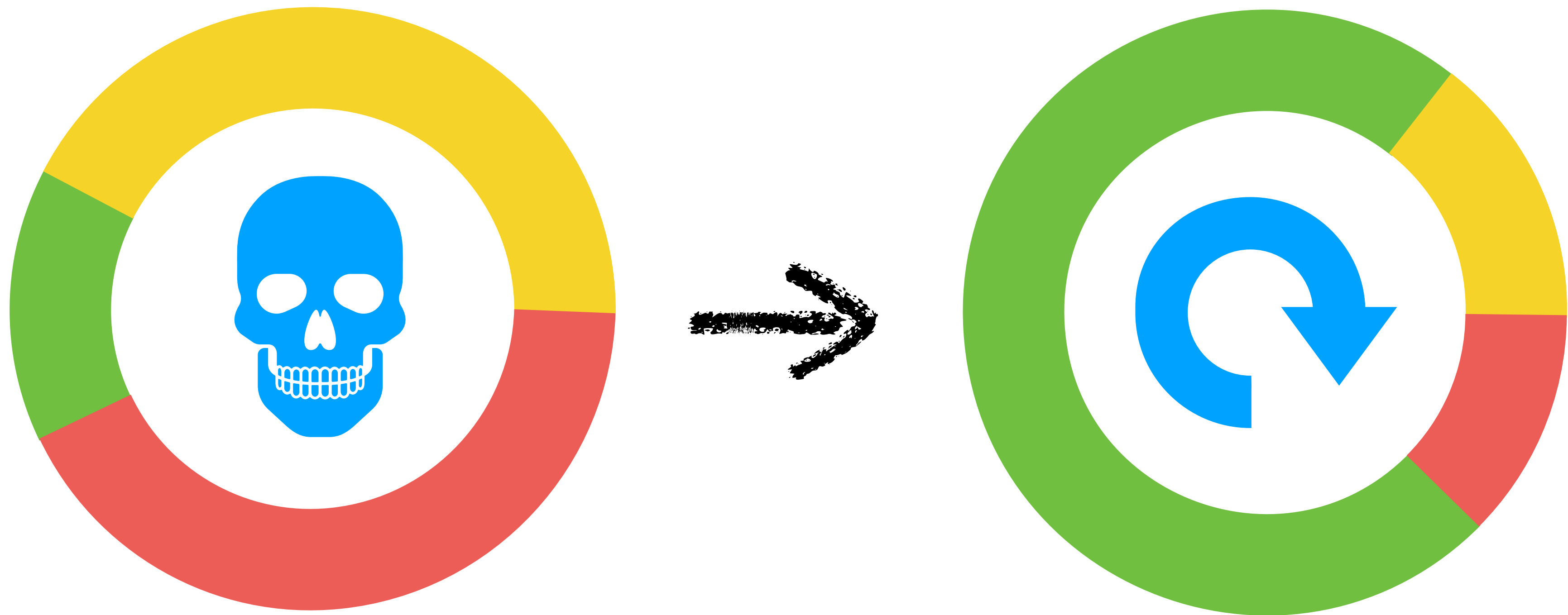


Scientific Reasoning



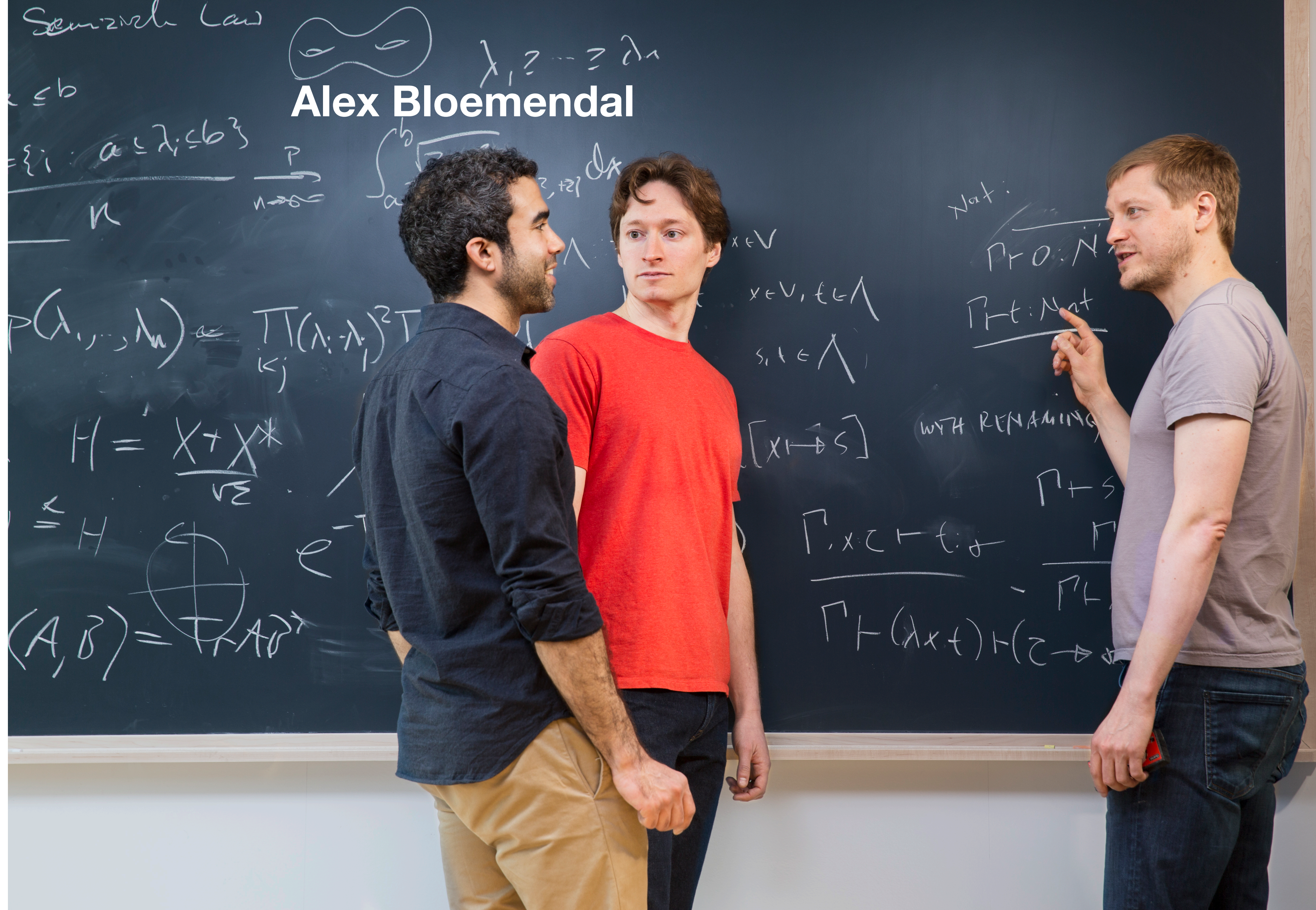
Implementation

Runtime



MIA Slides

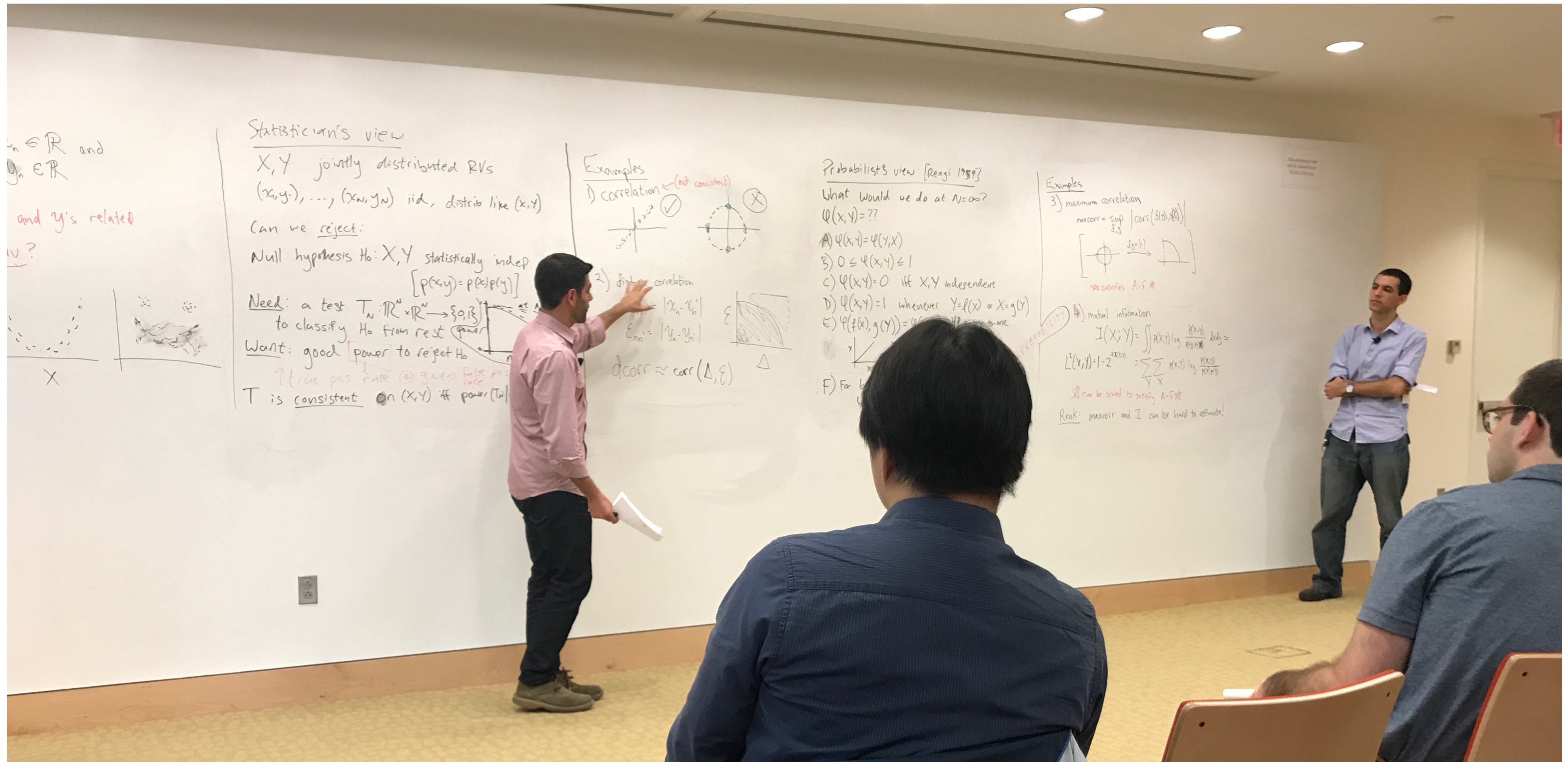
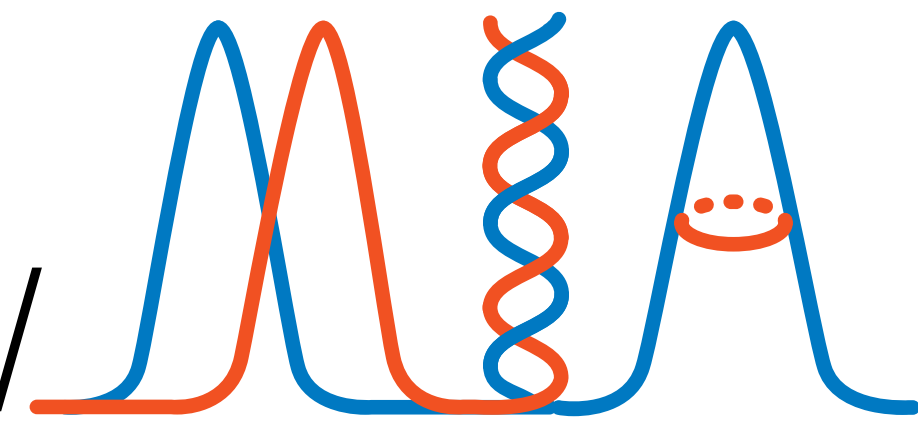
Alex Bloemendal



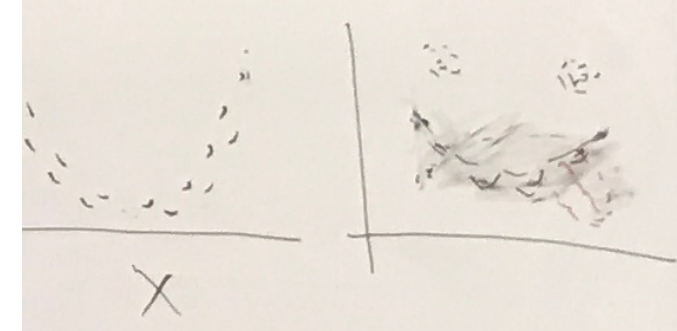
[HOME](#) » [SCIENCE](#)

MODELS, INFERENCE & ALGORITHMS





$x_n \in \mathbb{R}$ and $y_n \in \mathbb{R}$
and y 's related
iv?



Statistician's view
 X, Y jointly distributed RVs
 $(x_1, y_1), \dots, (x_n, y_n)$ iid, distrib like (x, y)
 Can we reject:
 Null hypothesis $H_0: X, Y$ statistically indep
 $[p(x, y) = p(x)p(y)]$
 Need: a test $T_n: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \{0, 1\}$
 to classify H_0 from rest (power)
 Want: good power to reject H_0
 ↑ true pos rate @ given false pos rate
 T is consistent on (X, Y) iff power (T_n)

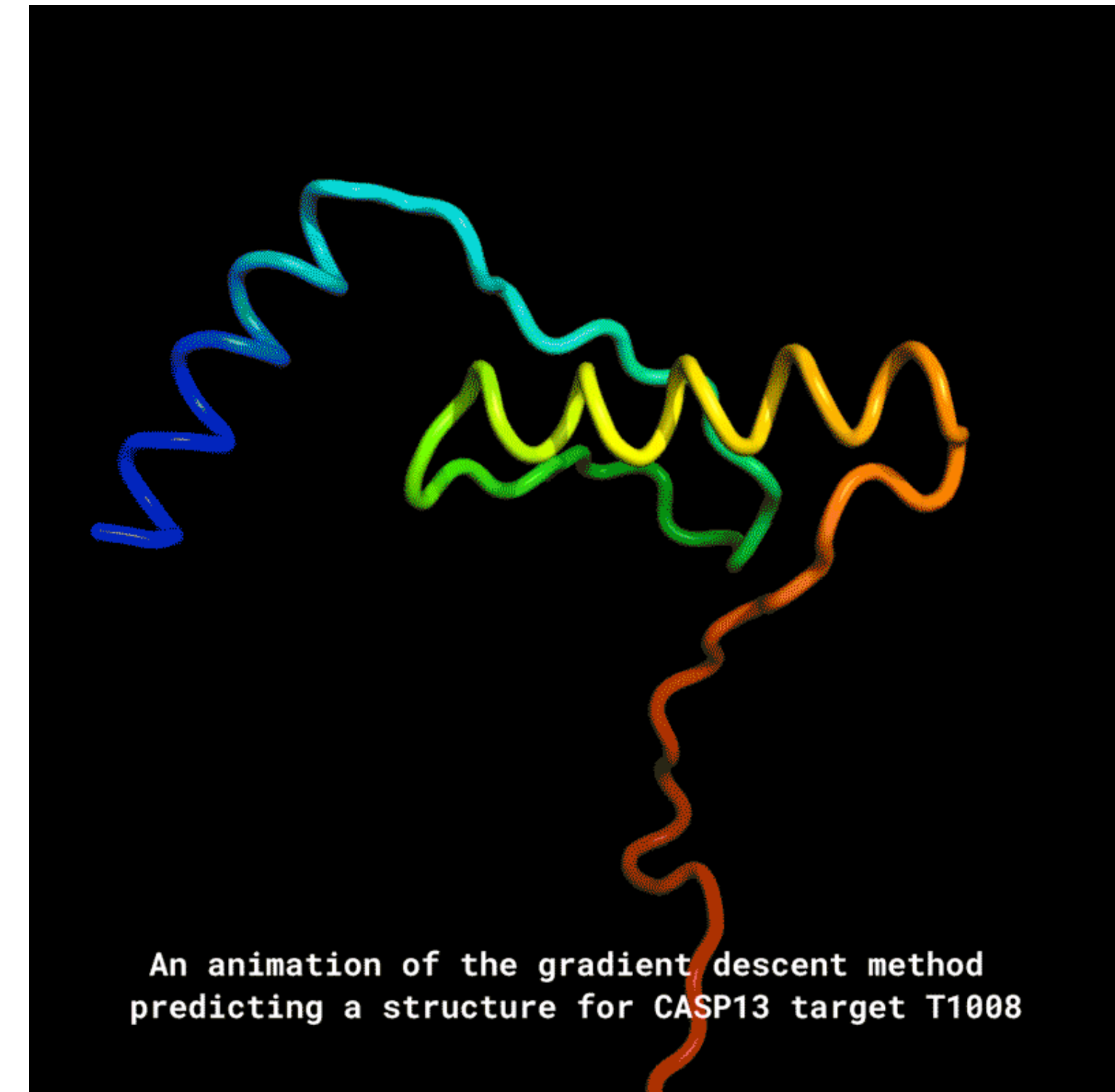
Examples
 1) correlation (not consistent)

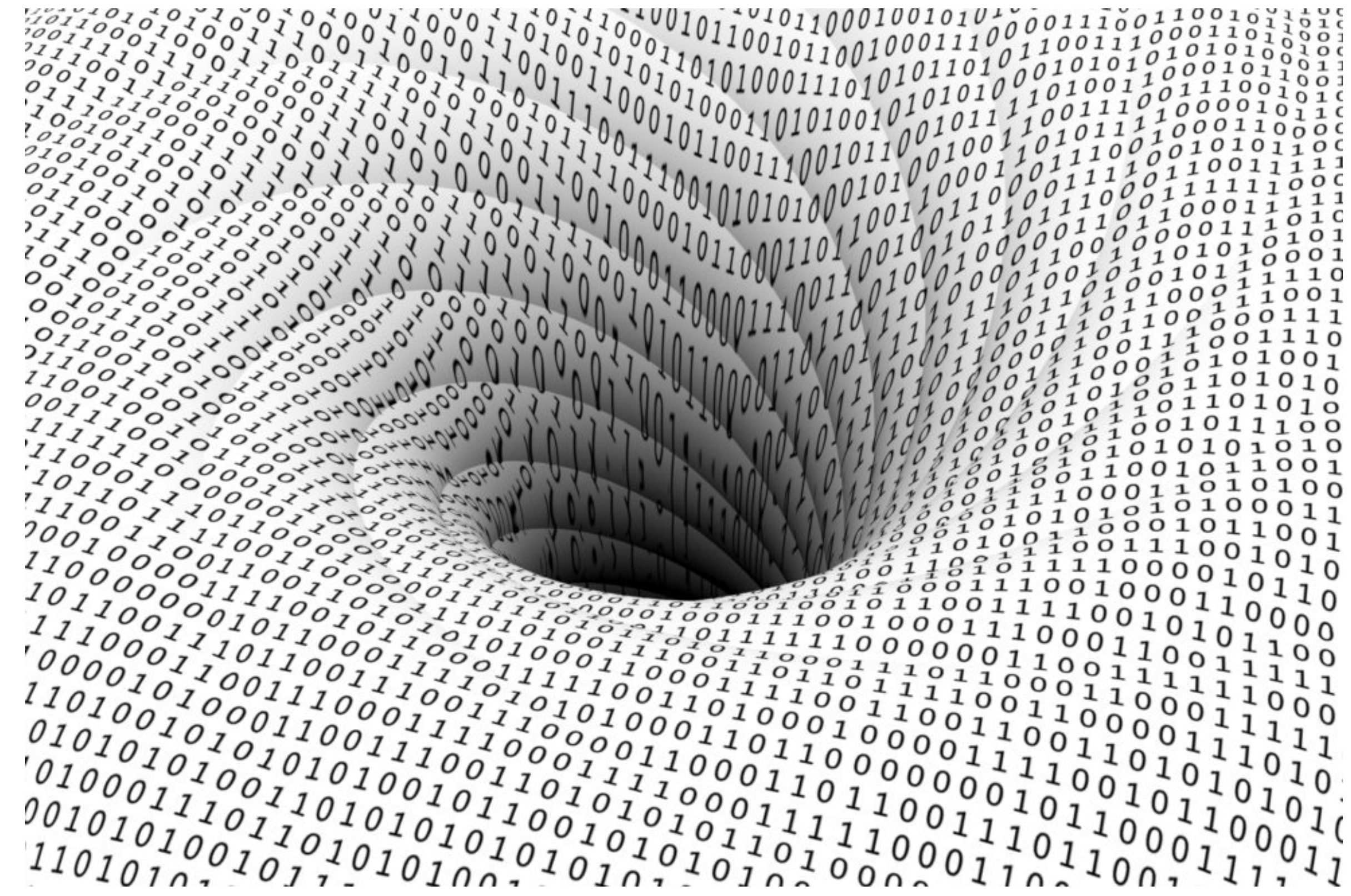
 2) distance correlation
 $\epsilon_{nn} = |x_n - x'_n|$
 $\epsilon_{nn} = |y_n - y'_n|$
 $d_{corr} \approx \text{corr}(\Delta, \epsilon)$

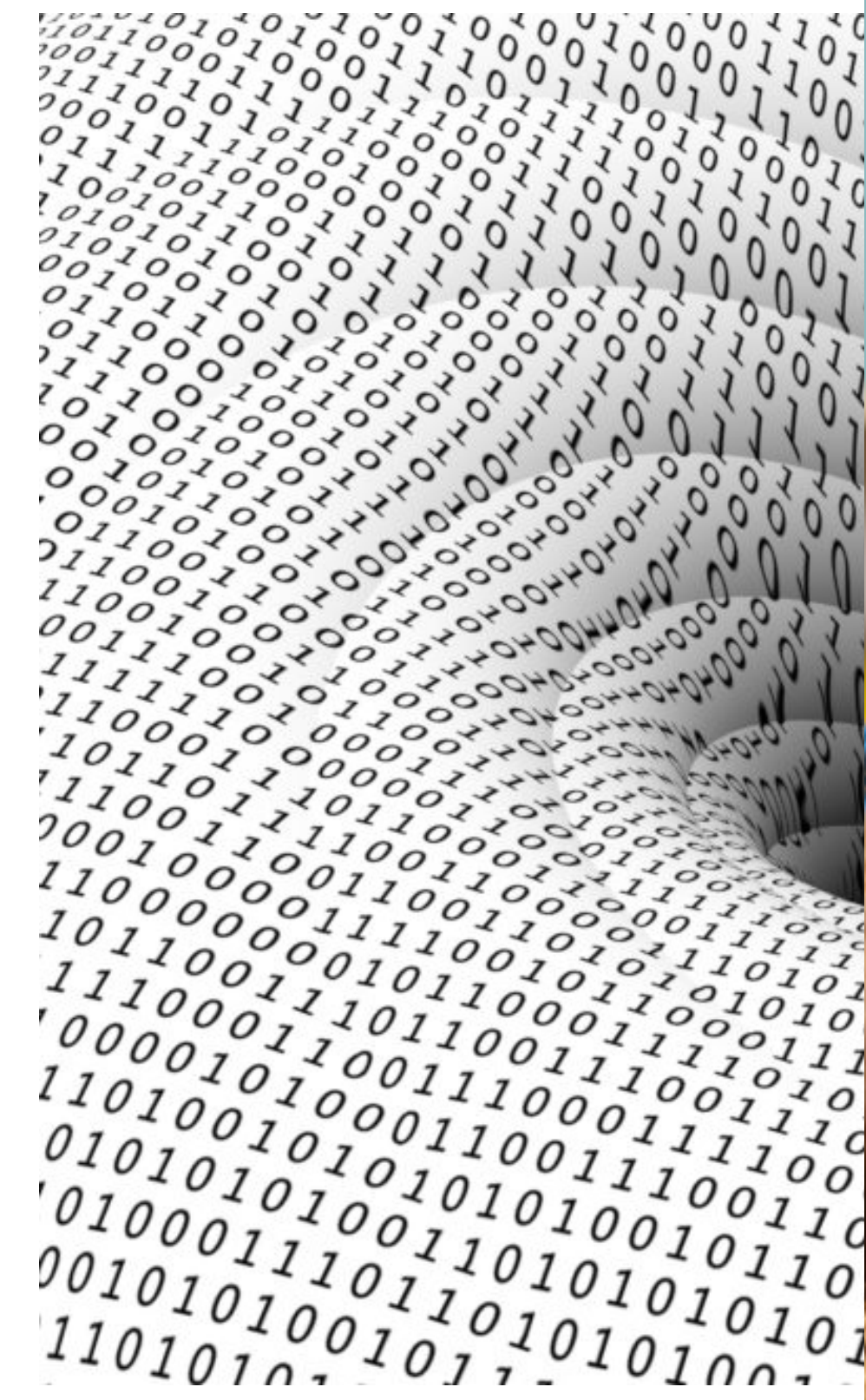
Probabilist's view [Rengji 1999]
 What would we do at $N = \infty$?
 $\varphi(x, y) = ??$
 A) $\varphi(x, y) = \varphi(y, x)$
 B) $0 \leq \varphi(x, y) \leq 1$
 C) $\varphi(x, y) = 0$ iff X, Y independent
 D) $\varphi(x, y) = 1$ whenever $Y = f(X)$ or $X = g(Y)$
 E) $\varphi(f(X), g(Y)) = \varphi(X, Y)$ iff f, g one-to-one
 F) For bivariate (X, Y)

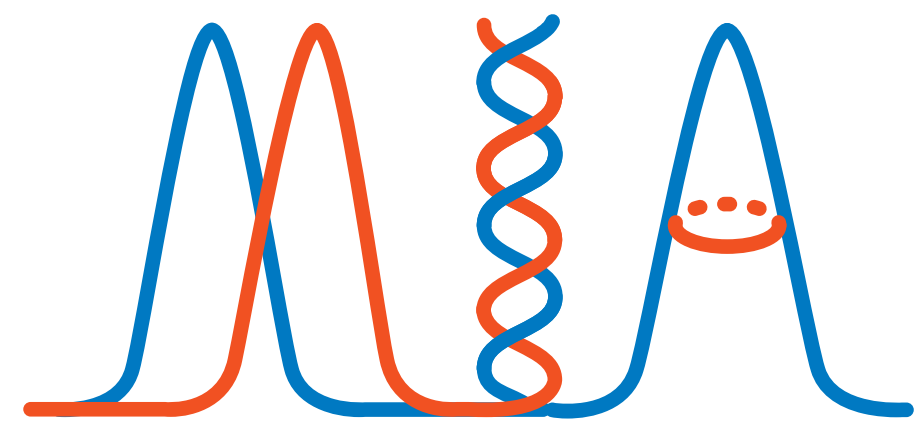
Examples
 3) maximum correlation
 $\text{maxcorr} = \sup_{f, g} |\text{corr}(f(X), g(Y))|$

 satisfies A-F
 4) mutual information
 $I(X; Y) = \int p(x, y) \log \frac{p(x, y)}{p(x)p(y)} dx dy =$
 $L^2(x, y) = 1 - 2^{-I(X; Y)} = \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$
 can be scaled to satisfy A-F
 Remark: maxcorr and I can be hard to estimate!










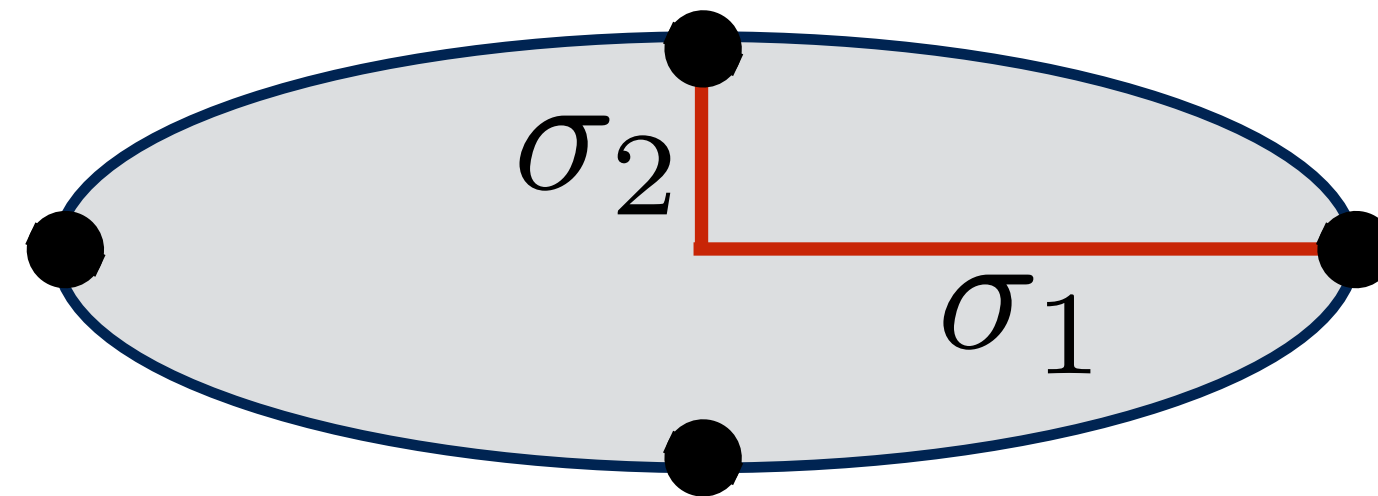
ML / AI needs biomedicine / neuroscience.



Backup Slides

Topology of PCA: Proof options

- Reduce to a simpler problem.
- Replace X with Σ for analysis or $[\Sigma \mid -\Sigma]$ for symmetry. 

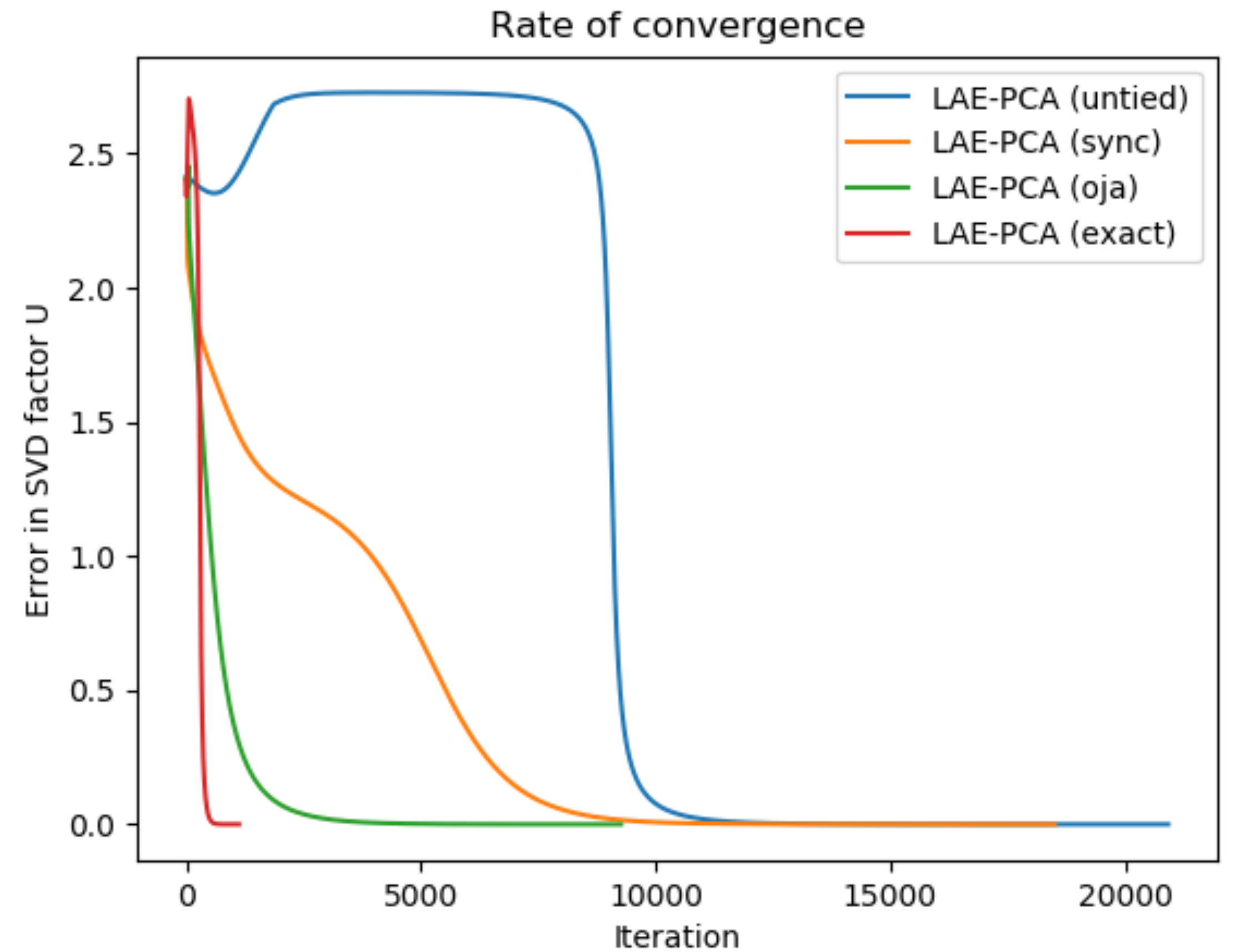
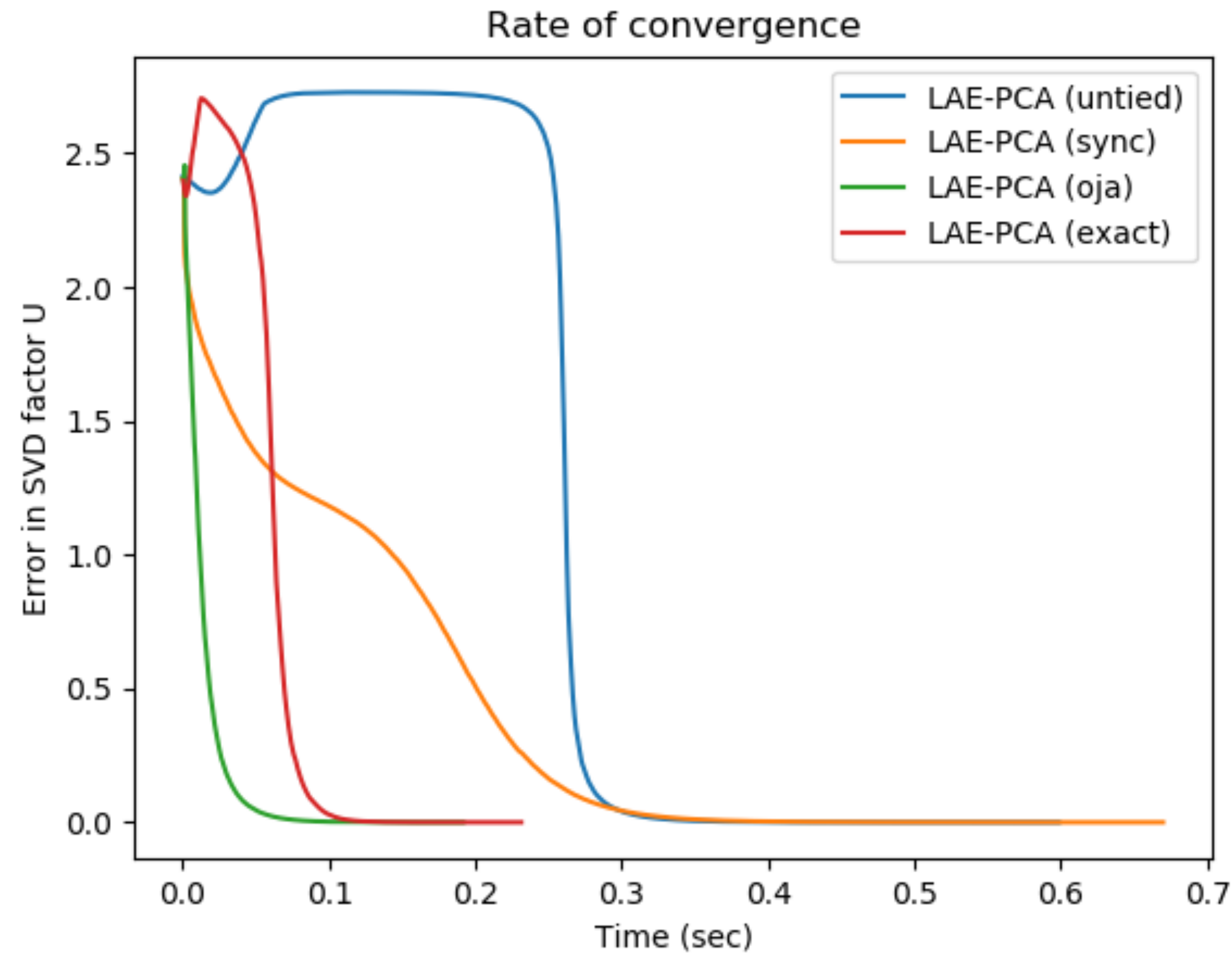


- Reduce to a harder problem that's already been solved: the LAE! 

$$\begin{array}{ccc}
 V_k(\mathbb{R}^m) & \xrightarrow{\pi: O \mapsto \text{Im}(OO^\top)} & \text{Gr}_k(\mathbb{R}^m) \\
 \downarrow \iota: O \mapsto (O^\top, O) & & \downarrow \mathcal{L}_X \\
 \mathbb{R}^{k \times m} \times \mathbb{R}^{m \times k} & \xrightarrow{\mathcal{L}} & \mathbb{R}
 \end{array}$$

LAE-SVD optimization

- Algorithm: optimize L_2 -regularized LAE and then take SVD of the decoder.



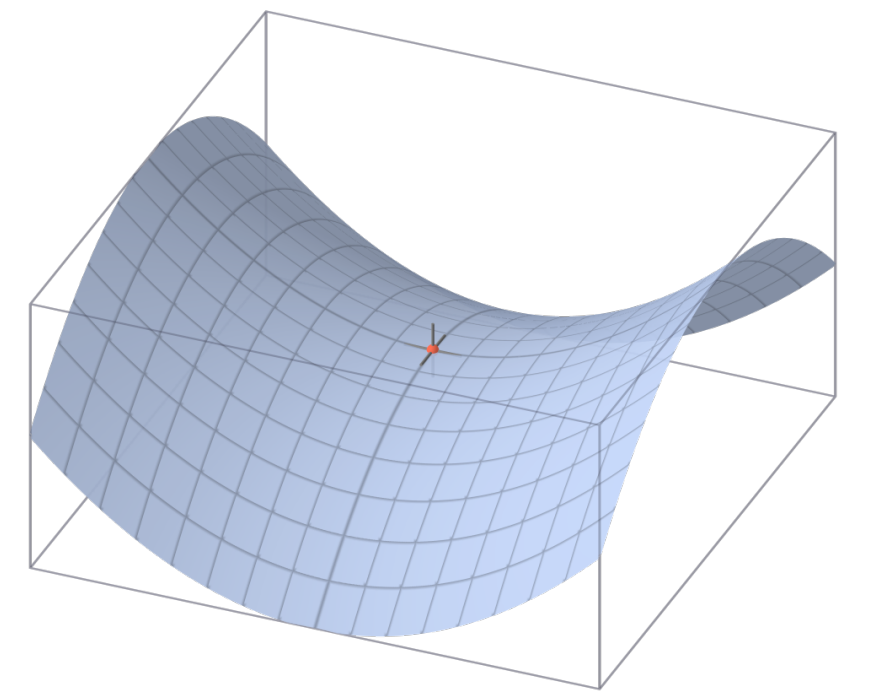
- Using SGD, algorithm resembles randomized SVD.

Morse theory

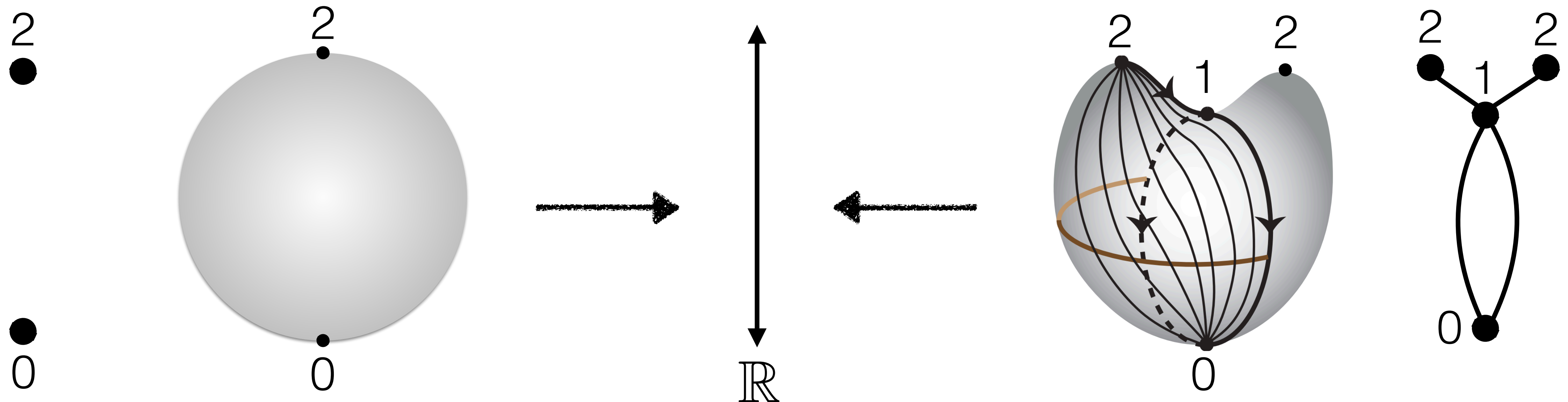
- Idea: Study the topology of a space via smooth functions on the space.

- A function is *Morse* if all critical points are non-degenerate:

$$f(x_1, \dots, x_m) = c - x_1^2 - \dots - x_d^2 + x_{d+1}^2 + \dots + x_m^2$$



- The *Morse index* d is number of negative eigenvalues of Hessian.

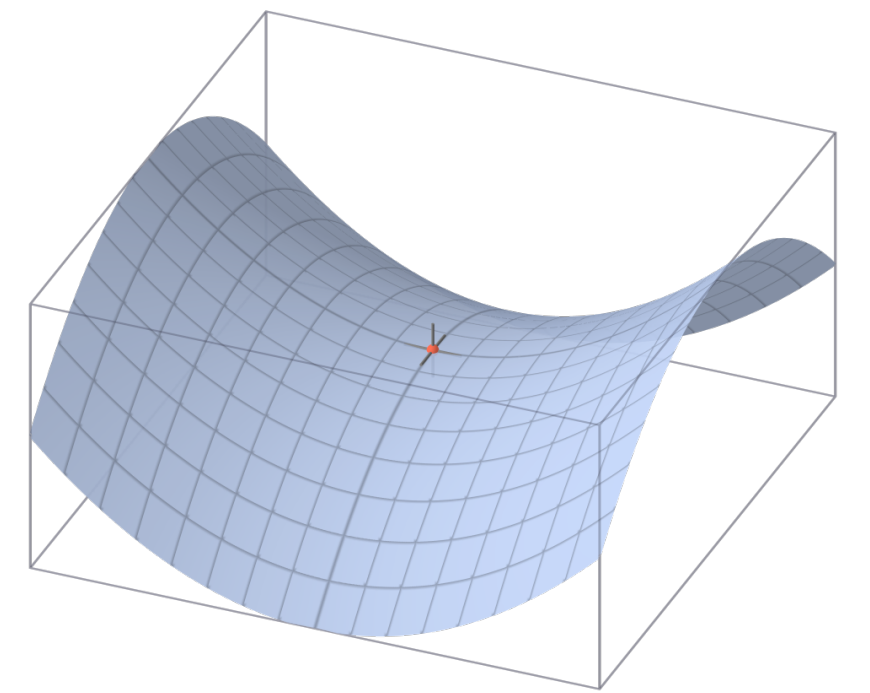


Morse theory

- Idea: Study the topology of a space via smooth functions on the space.

- A function is *Morse* if all critical points are non-degenerate:

$$f(x_1, \dots, x_m) = c - x_1^2 - \dots - x_d^2 + x_{d+1}^2 + \dots + x_m^2$$

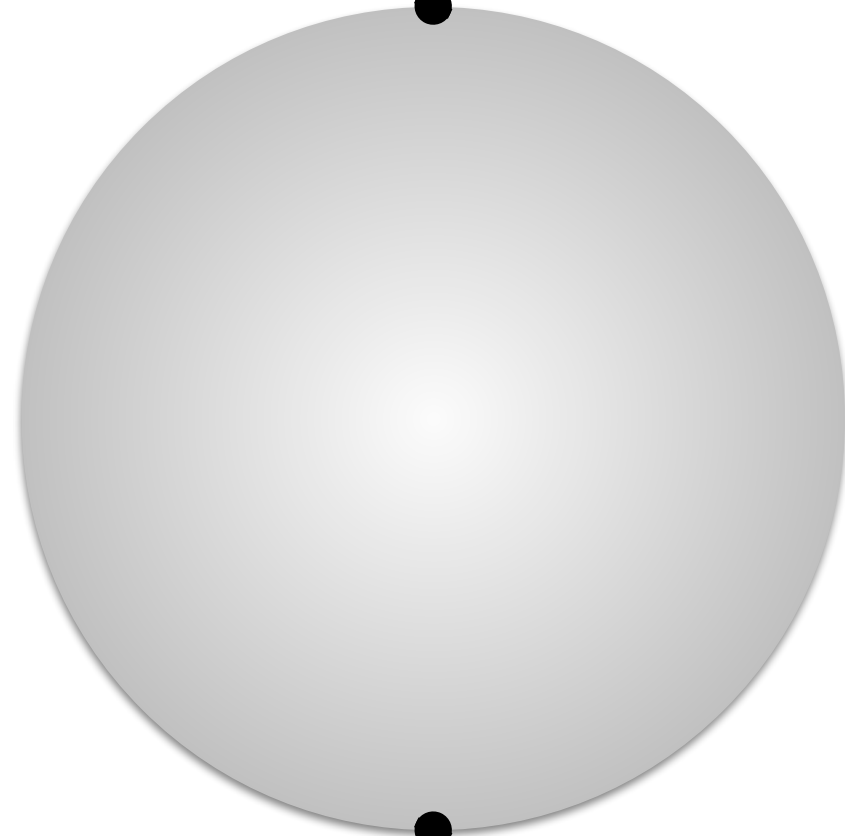


- The *Morse index* d is number of negative eigenvalues of Hessian.

2



2



0

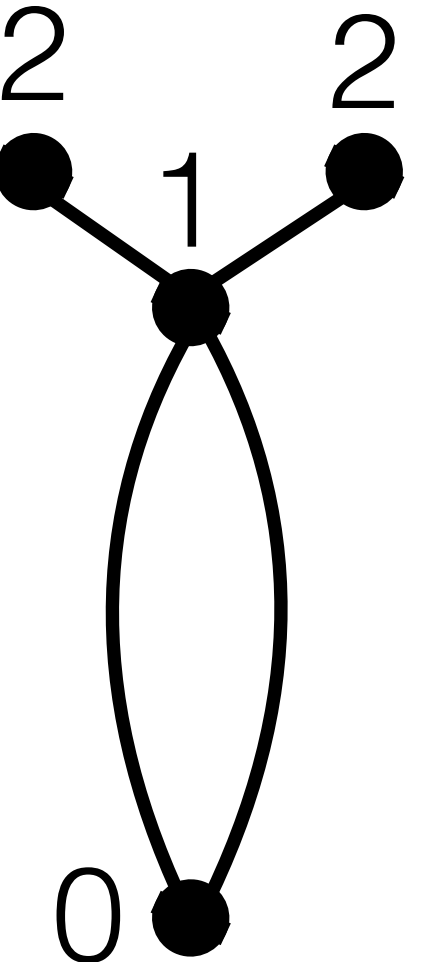
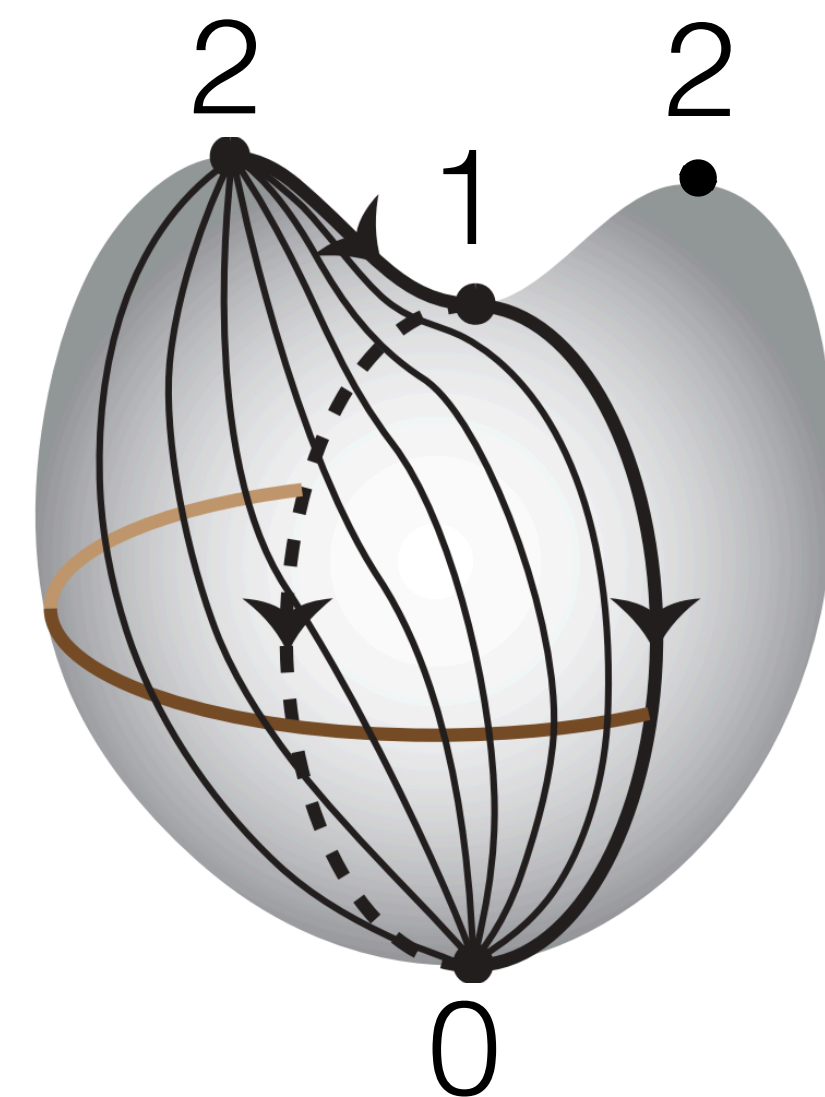
0



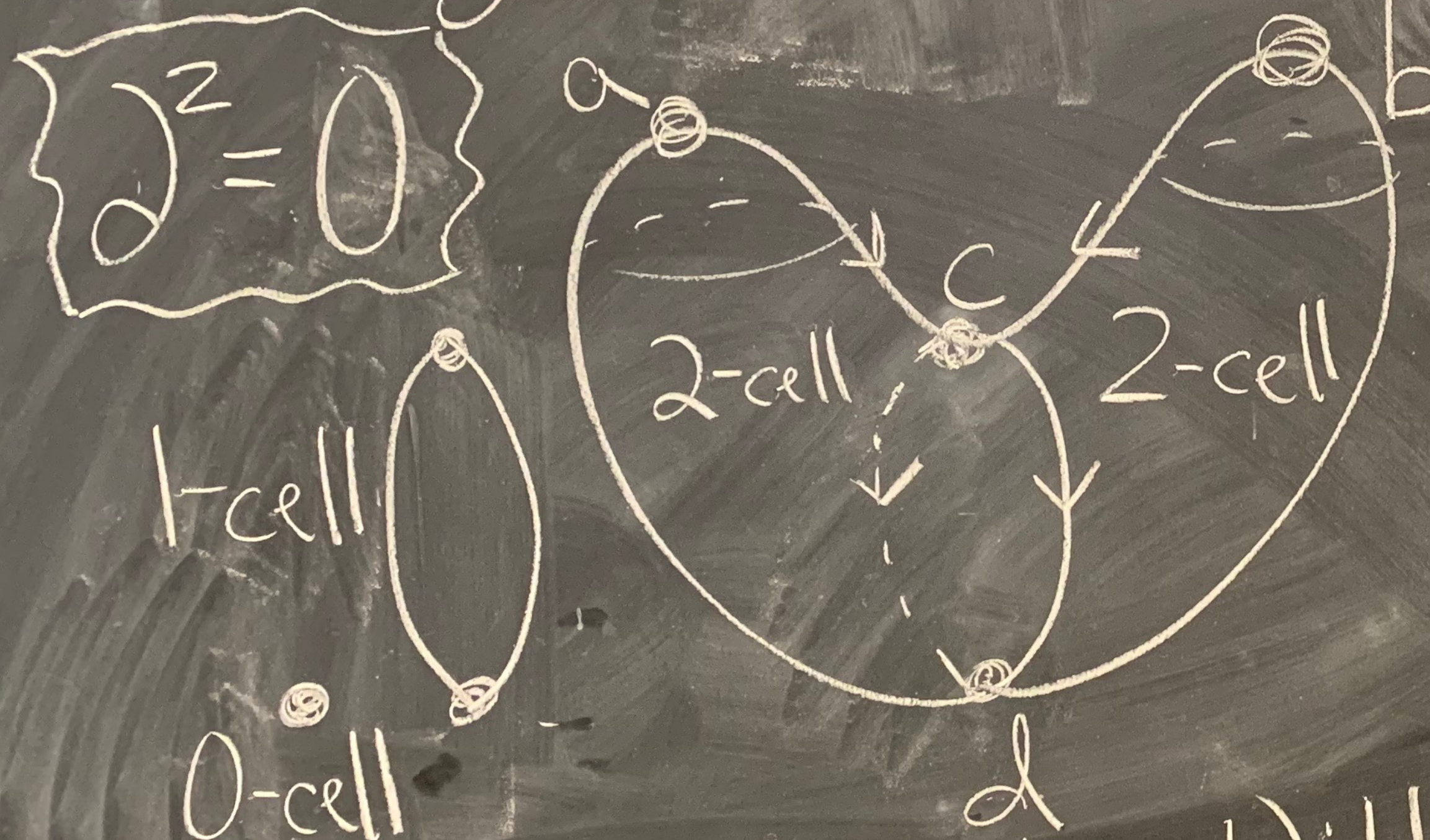
Euler characteristic

$$\chi = \sum (-1)^{d_i}$$

$$1 - 0 + 1 = 2 = 2 - 1 + 1$$



$$\langle \partial x, y \rangle = \# \left\{ \begin{array}{l} \text{gradient} \\ \text{trajectory} \\ \gamma: x \rightarrow y \end{array} \right\}$$

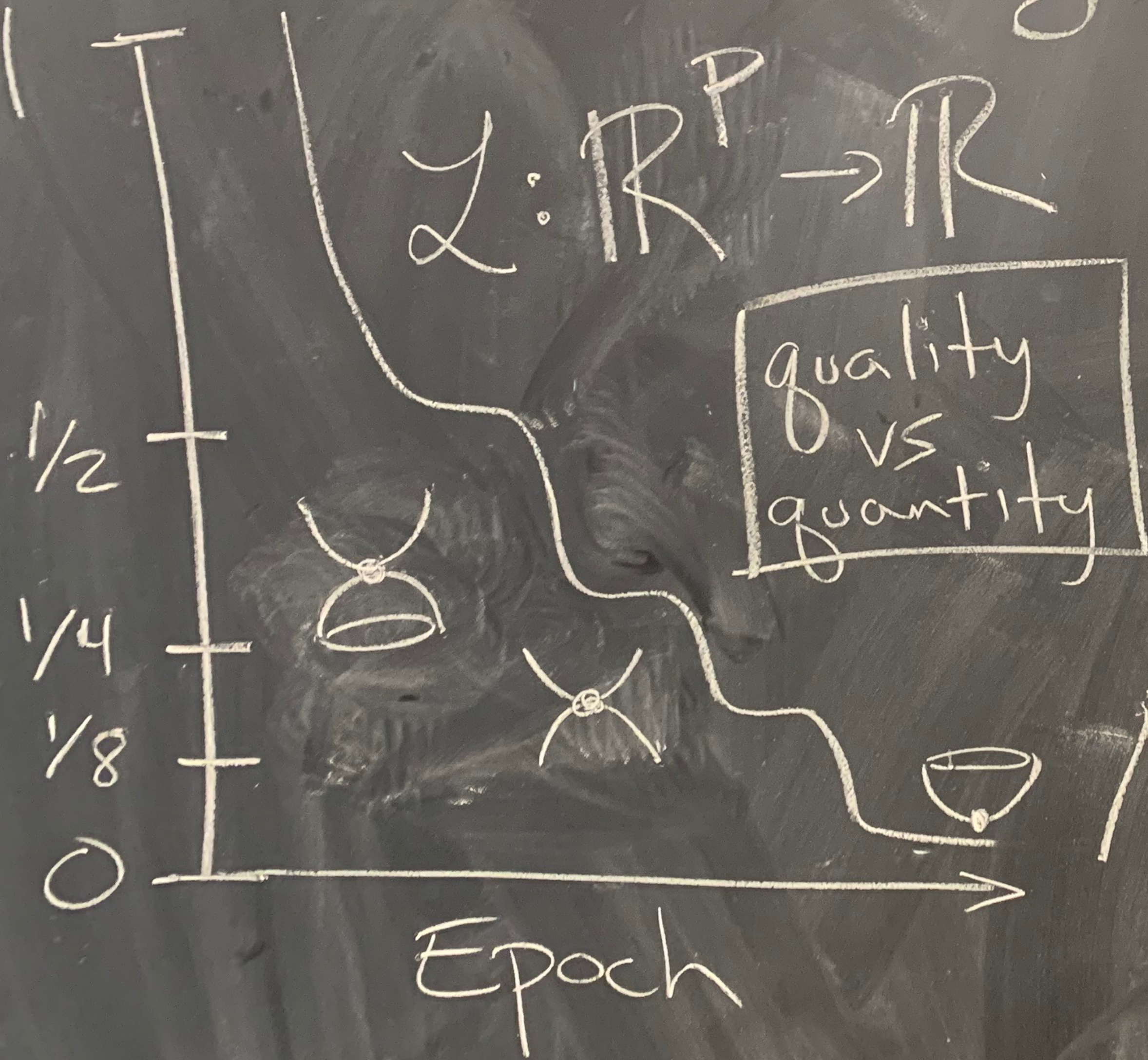


Thom, Smale, Milnor, Witten

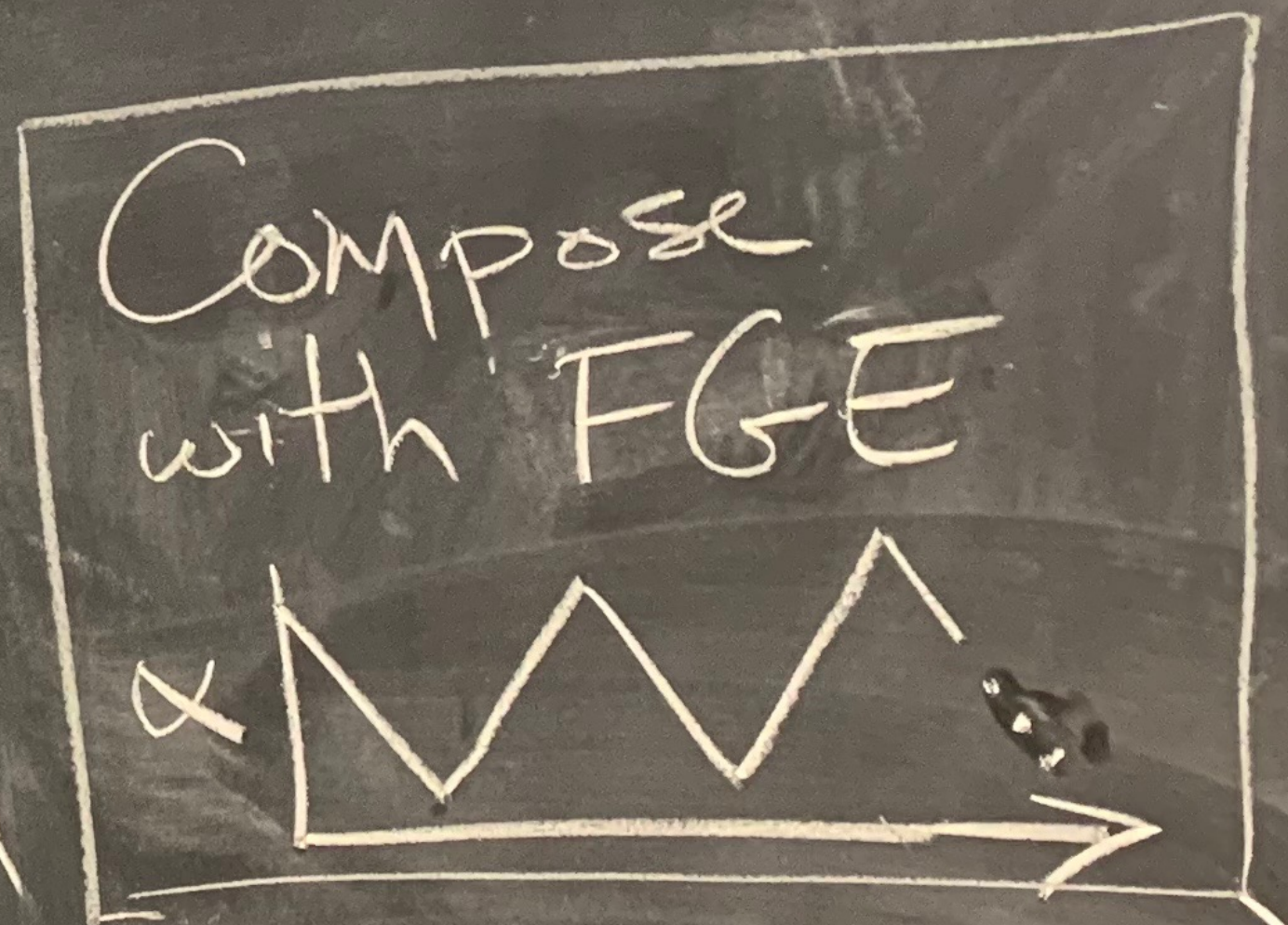
Chain Complex	(Morse) Homology
$\mathbb{F}_a \oplus \mathbb{F}_b$ $\begin{array}{c} \downarrow \\ d_2 \end{array} \quad [1 \quad 1]$ \mathbb{F}_c	$\mathbb{F}(a+b)$
$\begin{array}{c} \downarrow \\ d_1 \end{array} \quad [0]$ \mathbb{F}_d	0
	\mathbb{F}_d
$H_i(M; \mathbb{F}) = \frac{\text{Ker}(d_i)}{\text{Im}(d_{i+1})}$	

Morse ensembling = Morsembling?

$$\mathcal{L}: \mathbb{R}^P \rightarrow \mathbb{R}$$



$\log(n)$ compute



LAE-SVD optimization

- Algorithm: optimize L_2 -regularized LAE and then take SVD of the decoder.

```
XXt = X @ X.T
while np.linalg.norm(W1 - W2.T) > epsilon:
    W1 -= alpha * ((W2.T @ (W2 @ W1 - I)) @ XXt + lamb * W1)
    W2 -= alpha * (((W2 @ W1 - I) @ XXt) @ W1.T + lamb * W2)

principal_directions, s, _ = np.linalg.svd(W2, full_matrices = False)
eigenvalues = np.sqrt(lamb / (1 - s**2))
```

- This is a regularized version of Oja's rule.

LAE-SVD optimization

- Algorithm: optimize L_2 -regularized LAE and then take SVD of the decoder.

```
XXt = X @ X.T
diff = np.inf
while diff > epsilon:
    update = alpha * (((W2 @ W2.T - I) @ XXt) @ W2 + lamb * W2)
    W2 -= update
    diff = np.linalg.norm(update)

principal_directions, s, _ = np.linalg.svd(W2, full_matrices = False)
eigenvalues = np.sqrt(lamb / (1 - s**2))
```

- This is a regularized version of Oja's rule.