

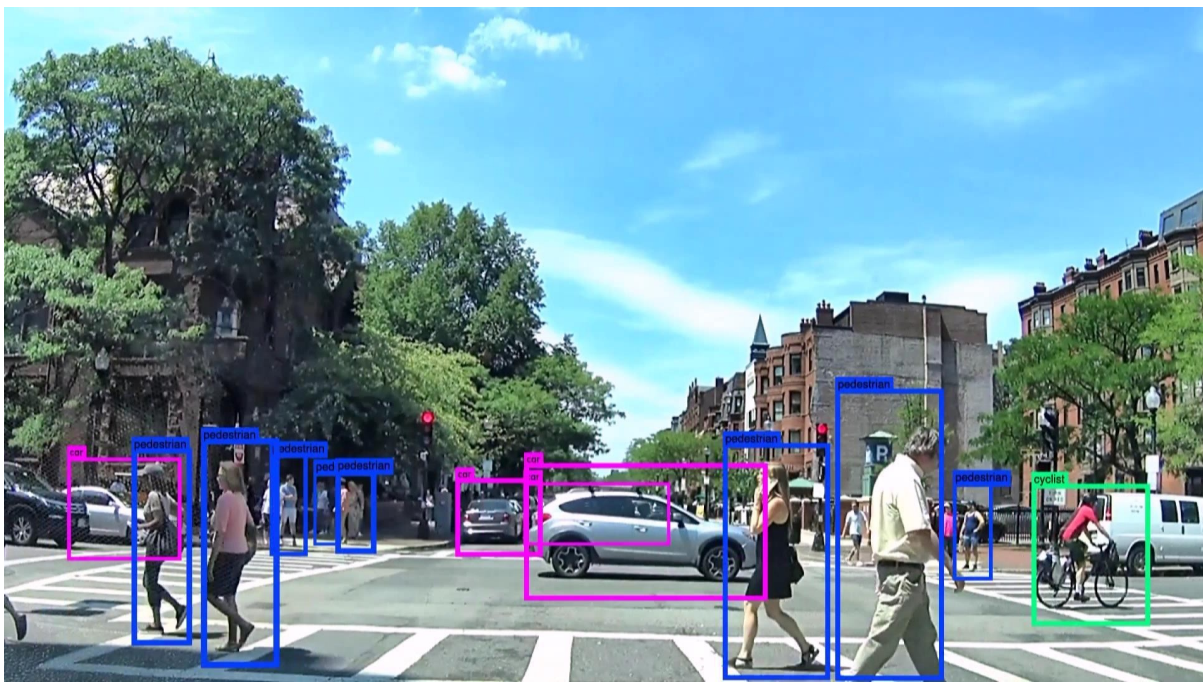
Models of representation learning dynamics

Andrew Saxe

Gatsby Unit & Sainsbury Wellcome Centre, UCL

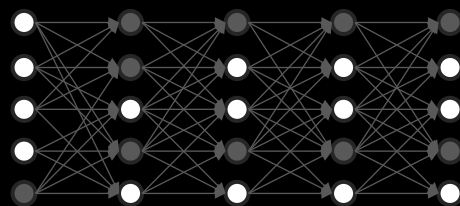








*Artificial
Intelligence*



*Neural
Networks*



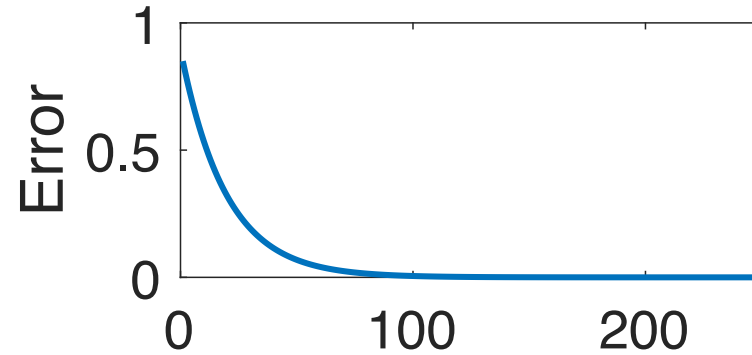
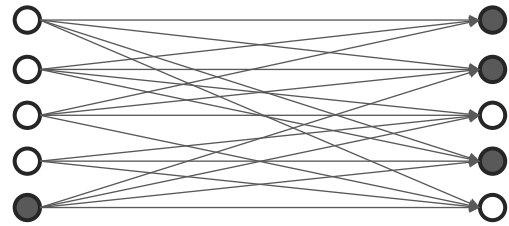
Brain & Mind

Today

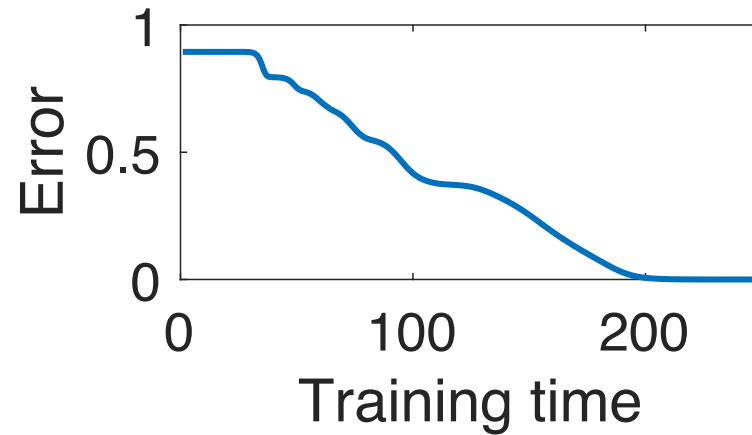
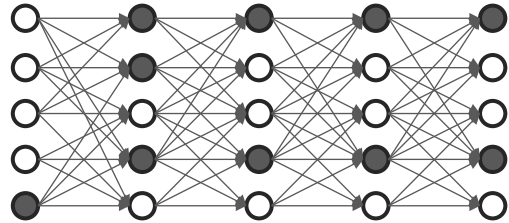
1. Deep linear network dynamics
2. Nontrivial initializations: Lazy, rich, & beyond
3. Nonlinear networks & the neural tangent reduction

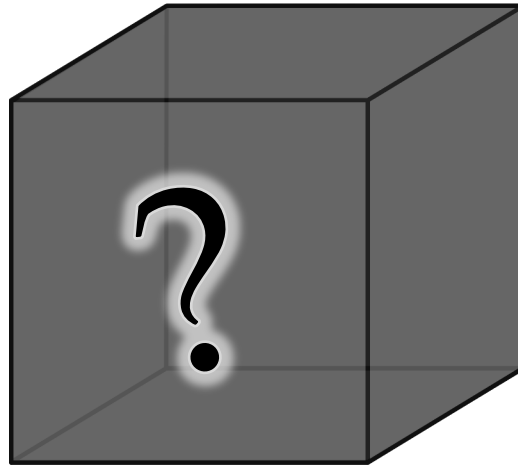
Depth complicates learning dynamics

Shallow



Deep



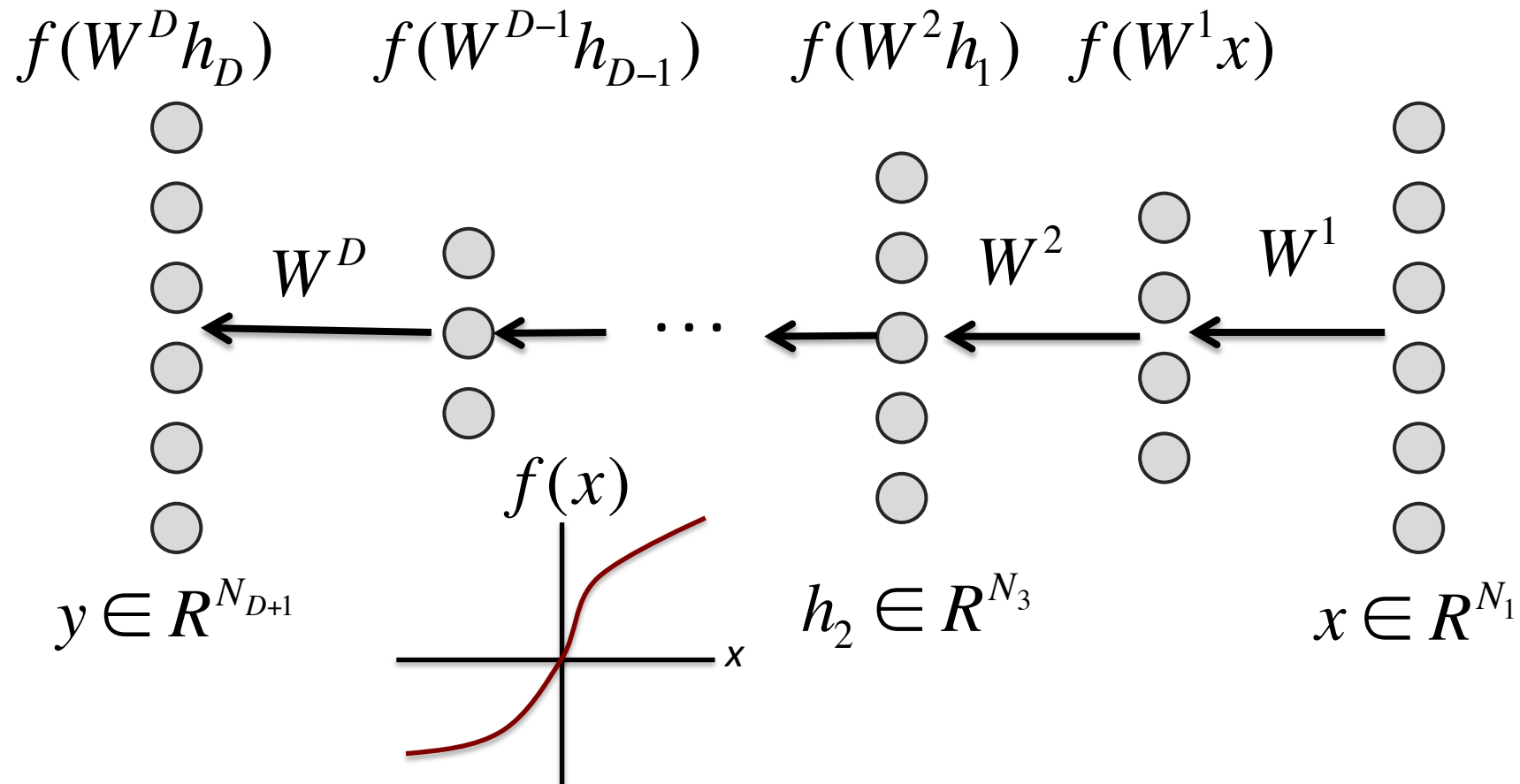


Surrogate models

- Tackling these questions in full generality is challenging
- Instead, we can analyze a surrogate model that is simpler but retains key features of the full problem
- Particularly for brain sciences, crucial to have a minimal, tractable model
 - Conceptual clarity
 - Unambiguous predictions
 - Isolate contribution of depth, data statistics, nonlinearity

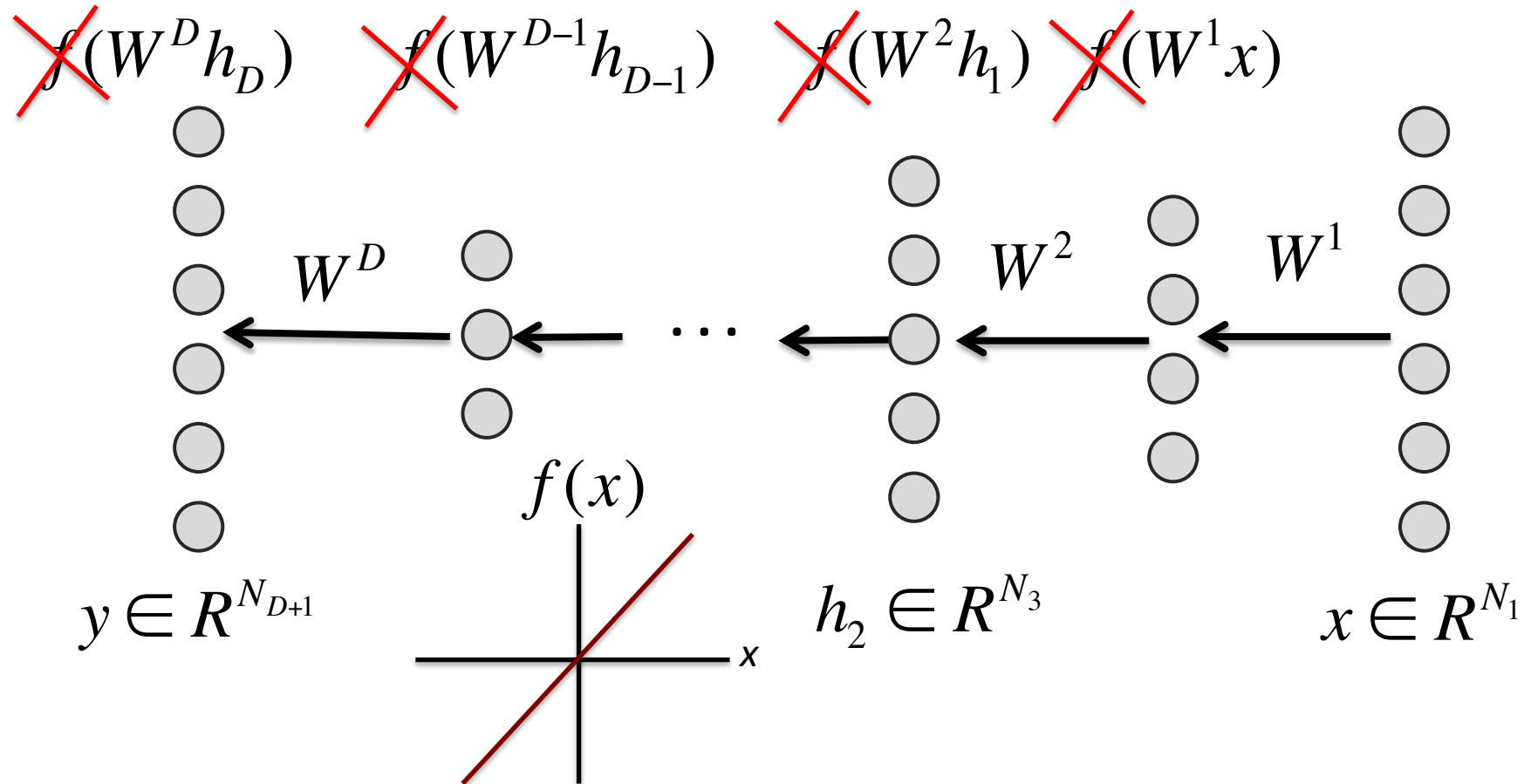
Deep network

- Little hope for a complete theory with arbitrary nonlinearities



Deep *linear* network

- Use a deep *linear* network as a starting point.



Gradient descent

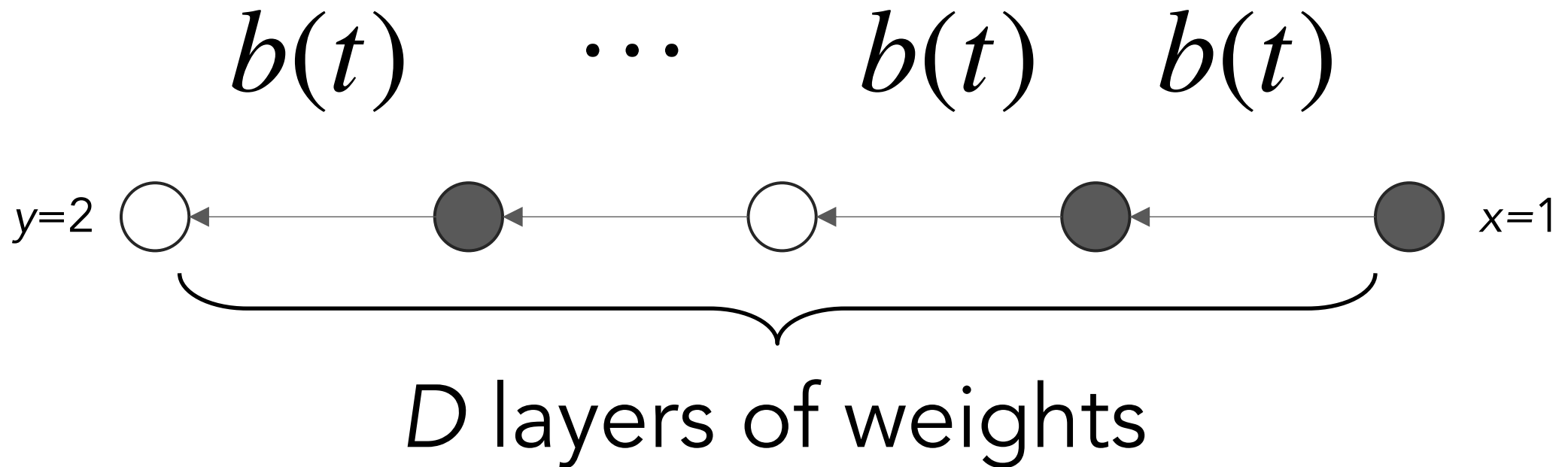
Mean squared error loss:

$$\min_{W_1, \dots, W_D} \sum_{\mu} \left\| y^{\mu} - \left(\prod_{i=1}^D W^i \right) x^{\mu} \right\|^2$$

Gradient flow dynamics:

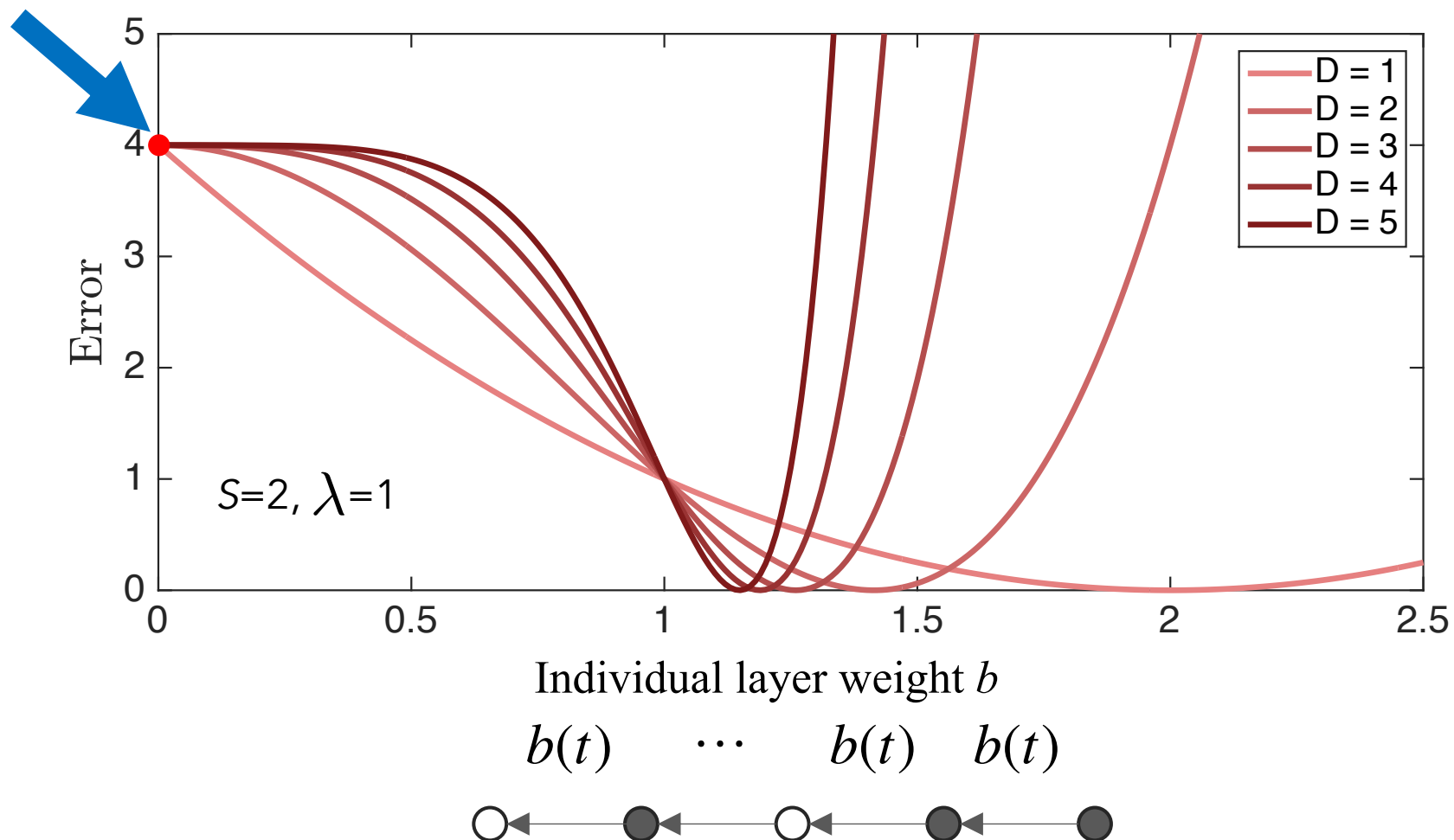
$$\tau \frac{d}{dt} W^l = \left(\prod_{i=l+1}^D W^i \right)^T \left[\Sigma^{yx} - \left(\prod_{i=1}^D W^i \right) \Sigma^{xx} \right] \left(\prod_{i=1}^{l-1} W^i \right)^T \quad l = 1, \dots, D$$

A linear chain

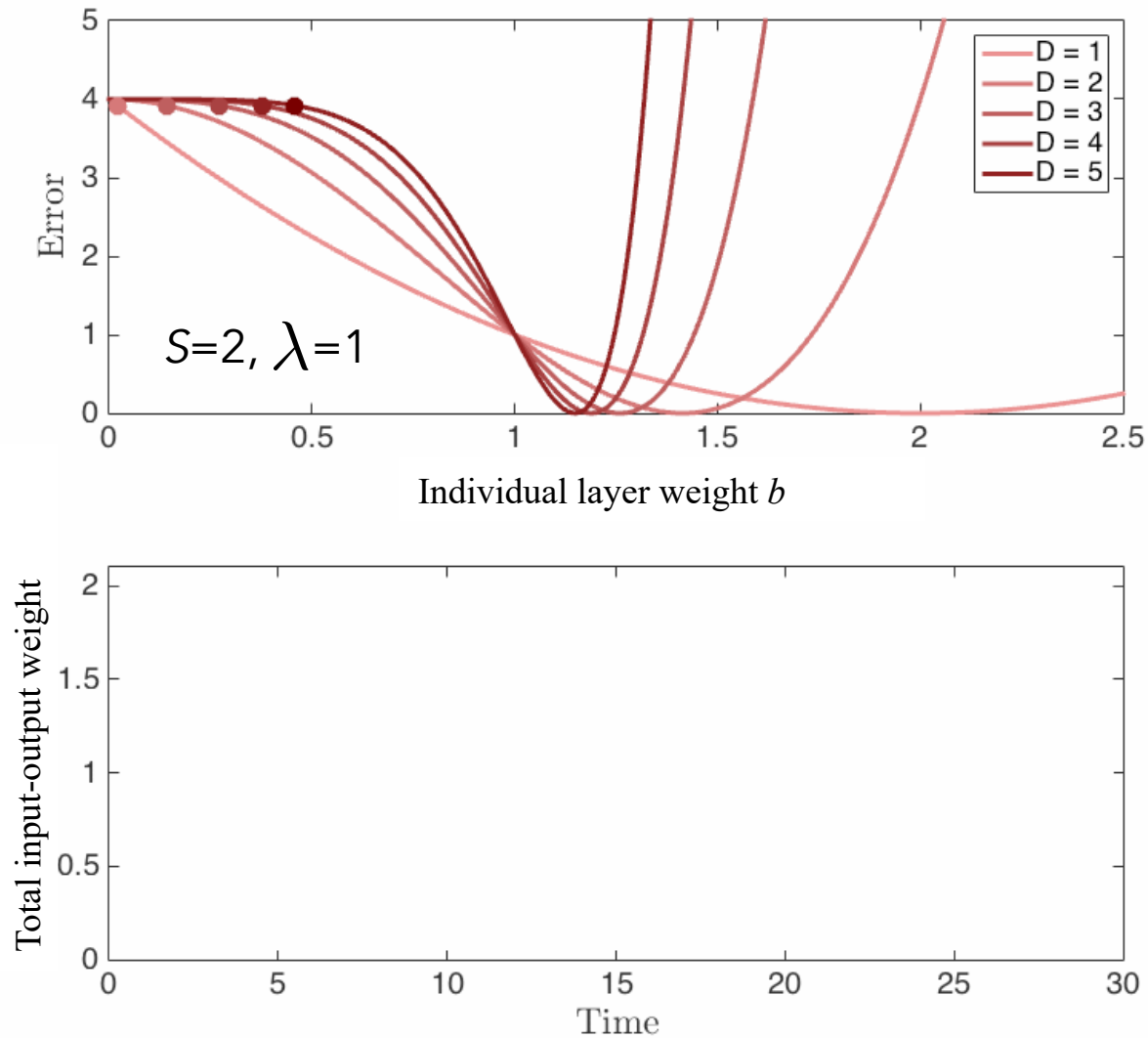


Error surface

Depth introduces a saddle point



Gradient descent dynamics



Analytic learning trajectory

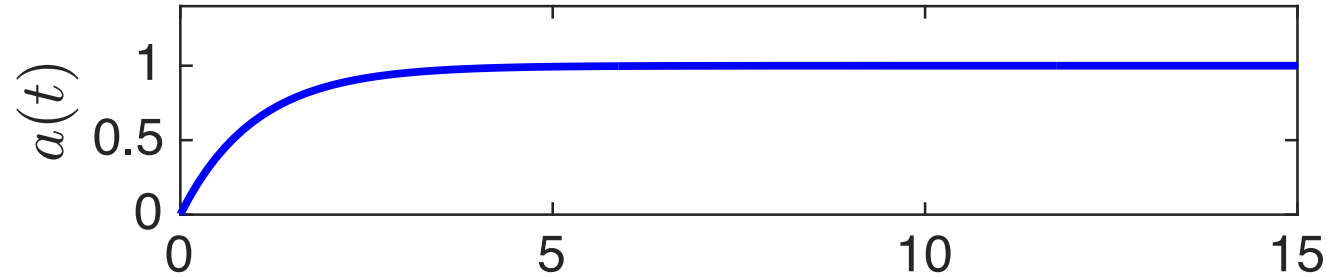
Shallow ($D=1$):
$$a(t) = \frac{s}{\lambda} \left(1 - e^{-t/\tau}\right) + a_0 e^{-t/\tau}$$

Deep ($D=2$):
$$a(t) = \frac{s/\lambda}{1 - \left(1 - \frac{s}{\lambda a_0}\right) e^{-\frac{2st}{\tau}}}$$

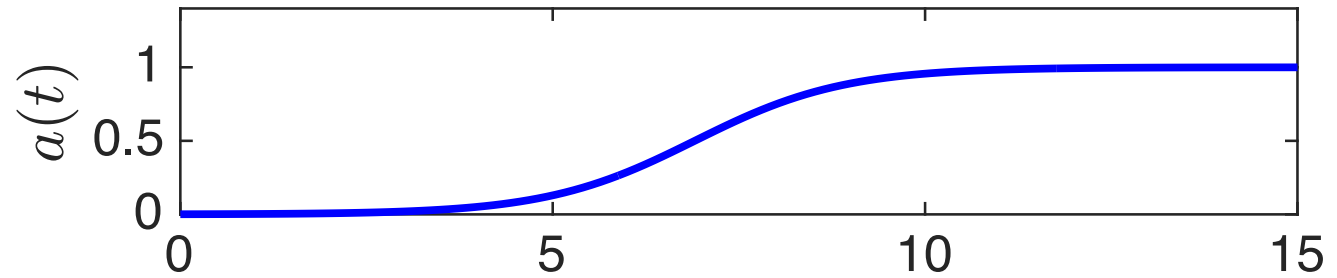
V. Deep ($D \rightarrow \infty$):
$$a(t) = \frac{s/\lambda}{1 + W \left[\left(\frac{s}{\lambda a_0} - 1 \right) e^{\frac{s}{\lambda a_0} - 1 - t/\tau} \right]}$$

Analytic learning trajectory

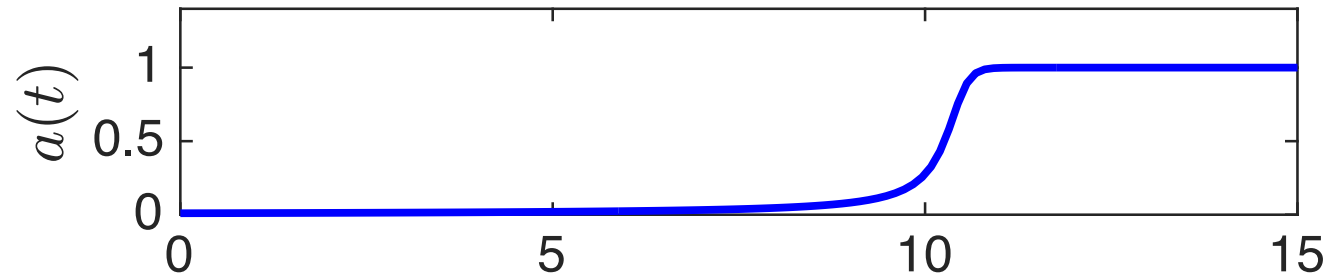
Shallow ($D=1$):



Deep ($D=2$):



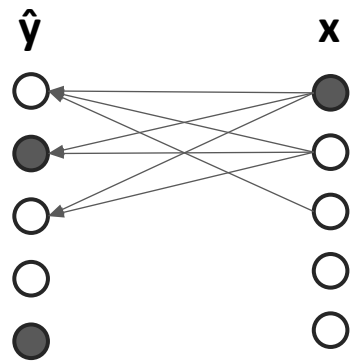
V. Deep ($D \rightarrow \infty$):



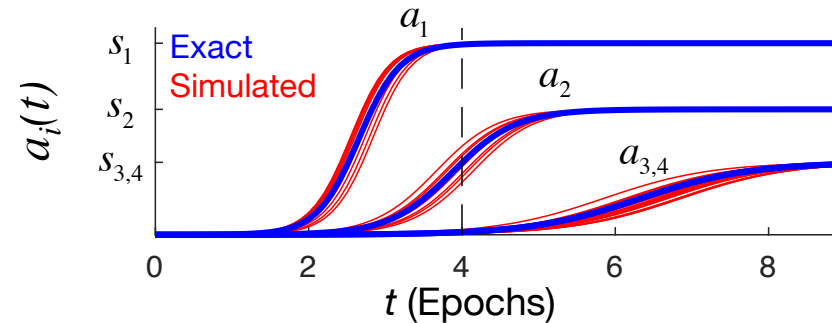
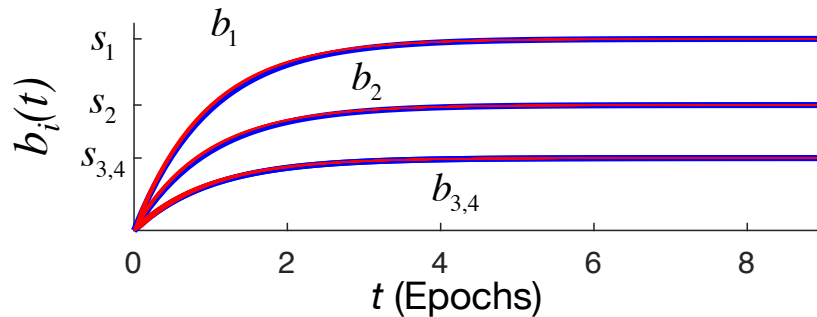
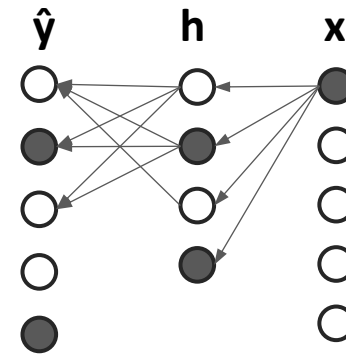
Epochs

Full networks act like several 1D chains

Shallow

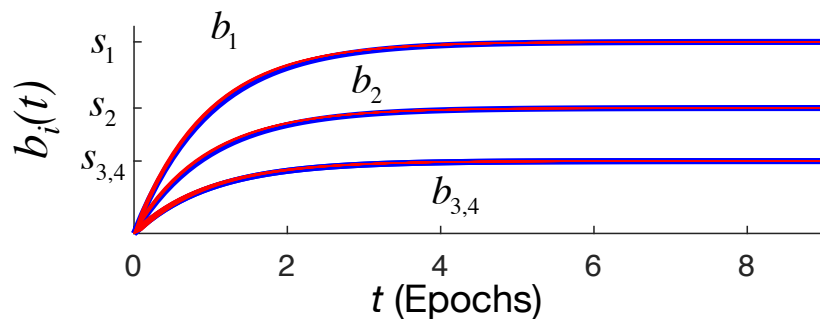
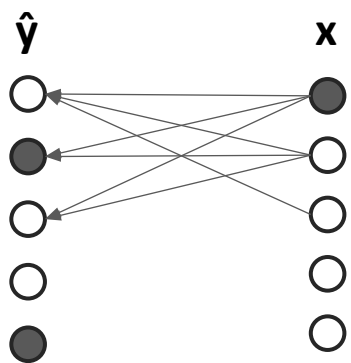


Deep

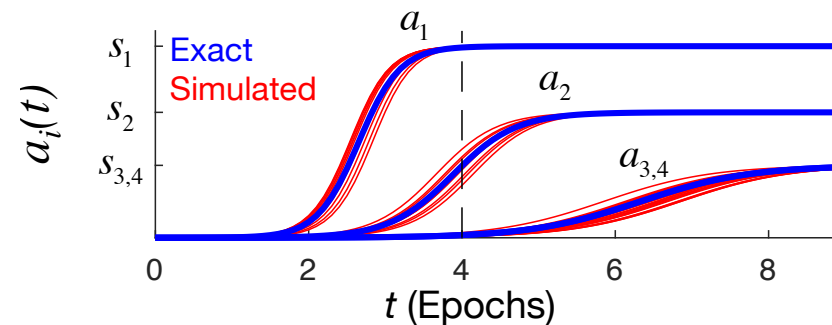
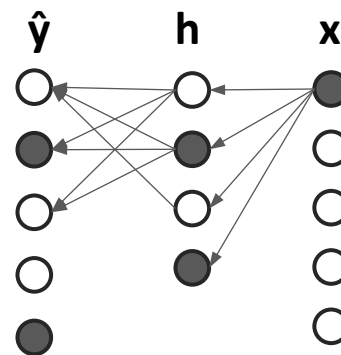


Depth introduces stage-like transitions

Shallow



Deep



Training speed

- How does training speed scale with depth?
- Time difference for deep net vs shallow net is

$$t_{\infty} - t_1 \approx O\left(\frac{1}{sb_0^D}\right)$$

t_D	epochs to train depth D network
b_0	Initial layer singular value
s	Minimum nonzero singular value
D	Depth

- Deep learning speed is highly sensitive to initial conditions

Effect of initialization

- Small random weights scale exponentially

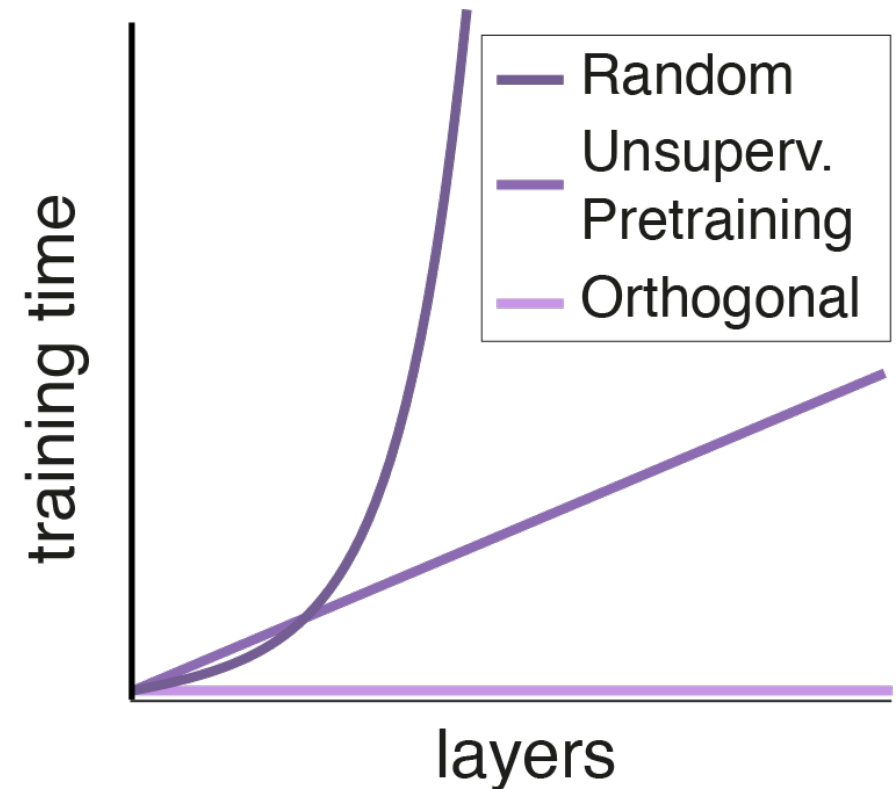
$$t_{\infty} - t_1 \approx O(1/b_0^D)$$

- Pretraining + fine-tuning scales linearly

$$t_{\infty} - t_1 \approx O(D/b_0^2)$$

- Orthogonal initialization: depth-independent

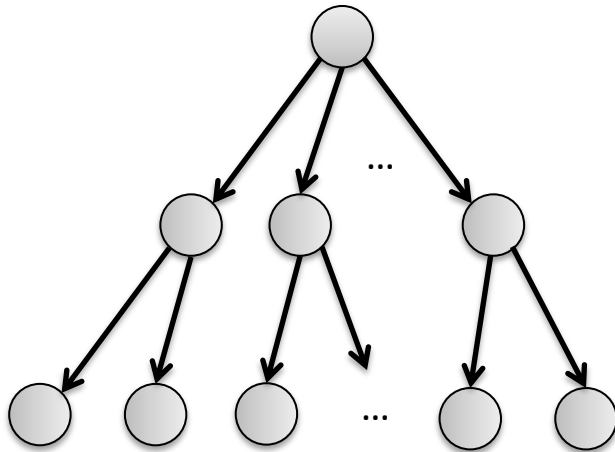
$$t_{\infty} - t_1 \approx O(1)$$



Connecting neural nets and graphical models

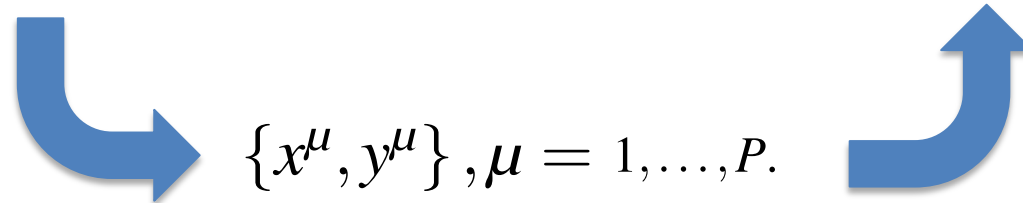
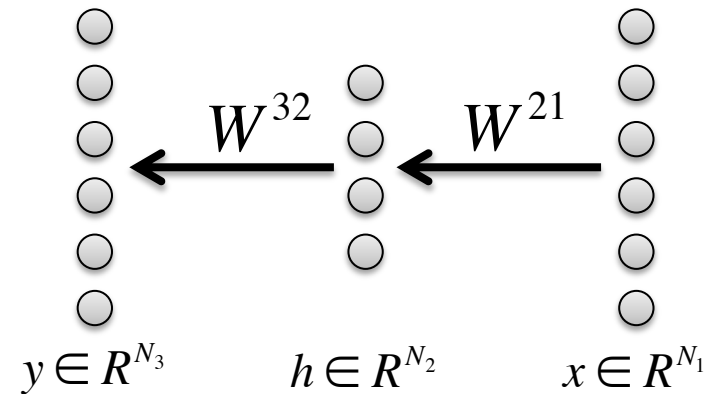
The “World”:

Structured generative model



The “Learner”:

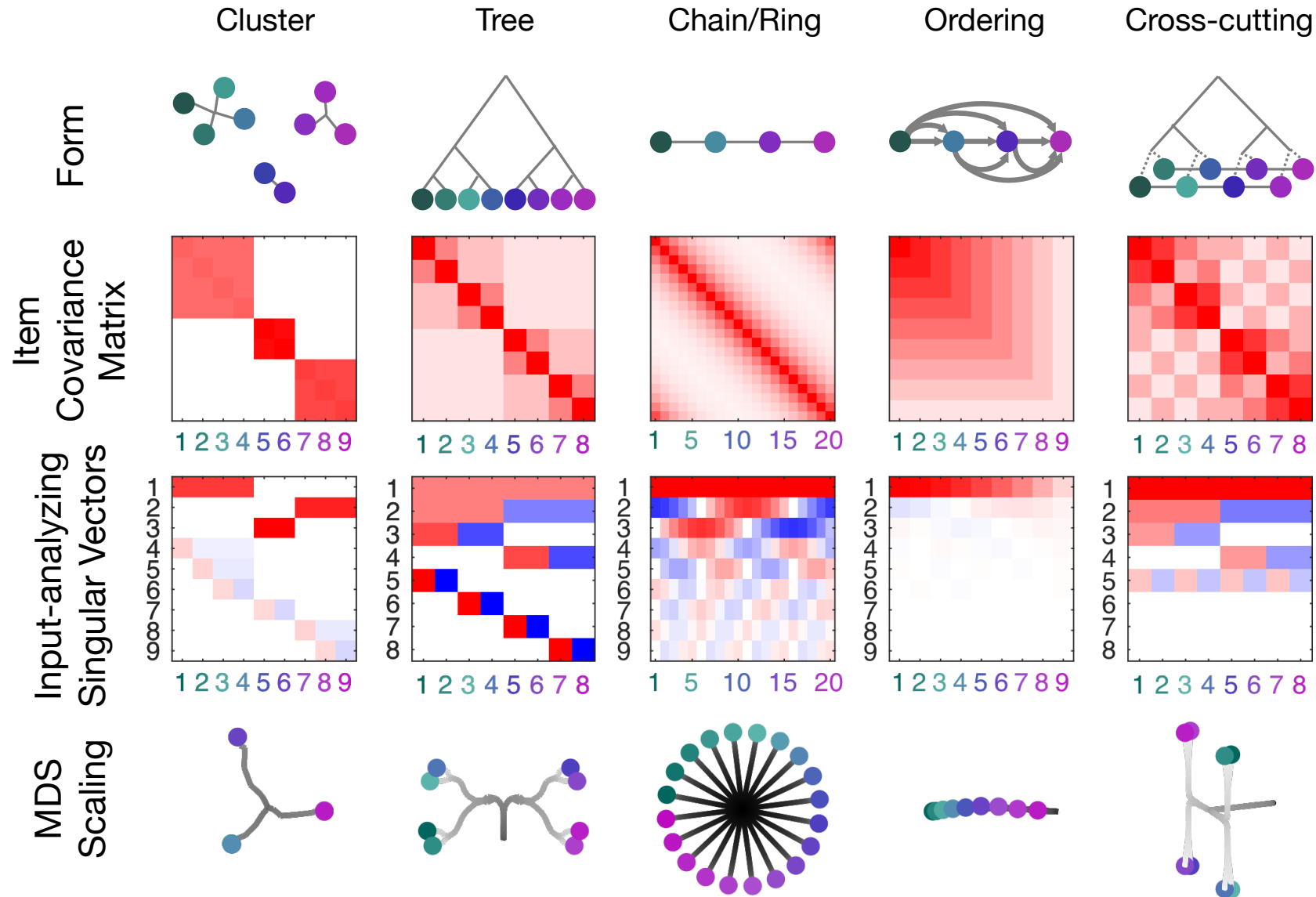
Deep linear network



Analytic link

- In the limit of many features, what matters to learning dynamics is SVD of correlation structure
- Can find this exactly for certain graphical models
 - Partitions
 - Trees
 - Grids/rings

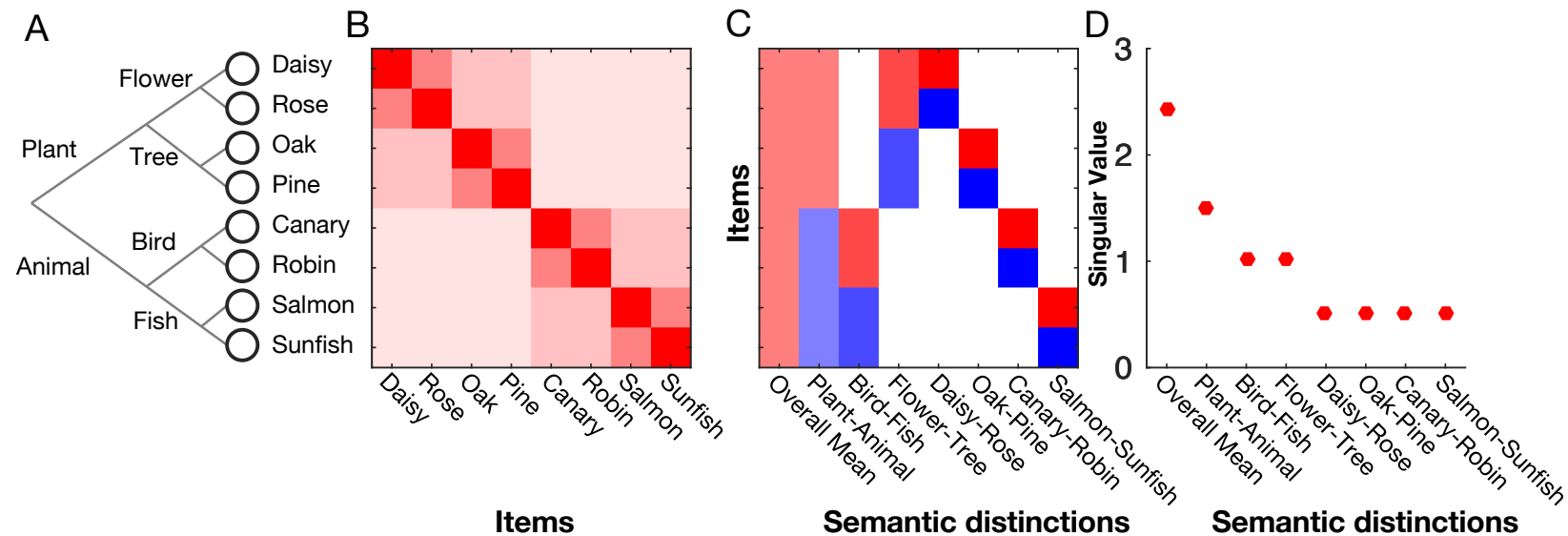
Learning diverse structures



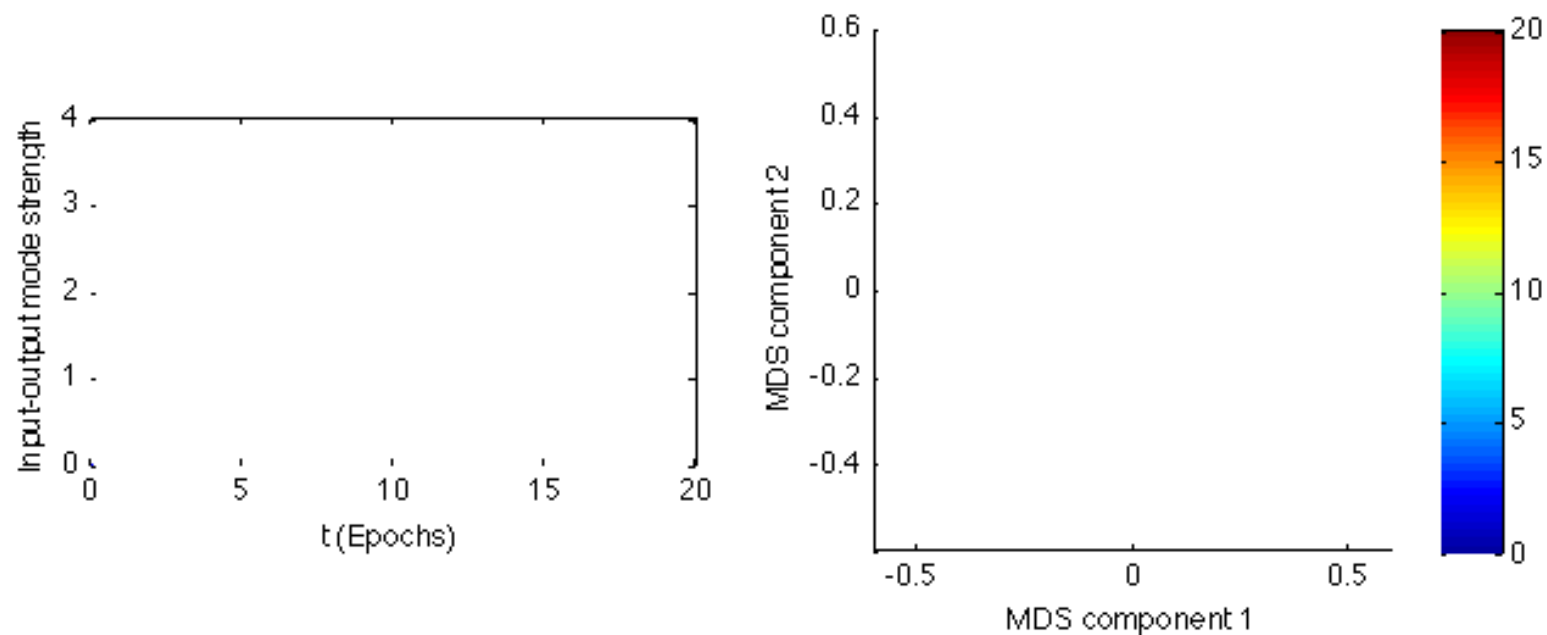
Progressive differentiation

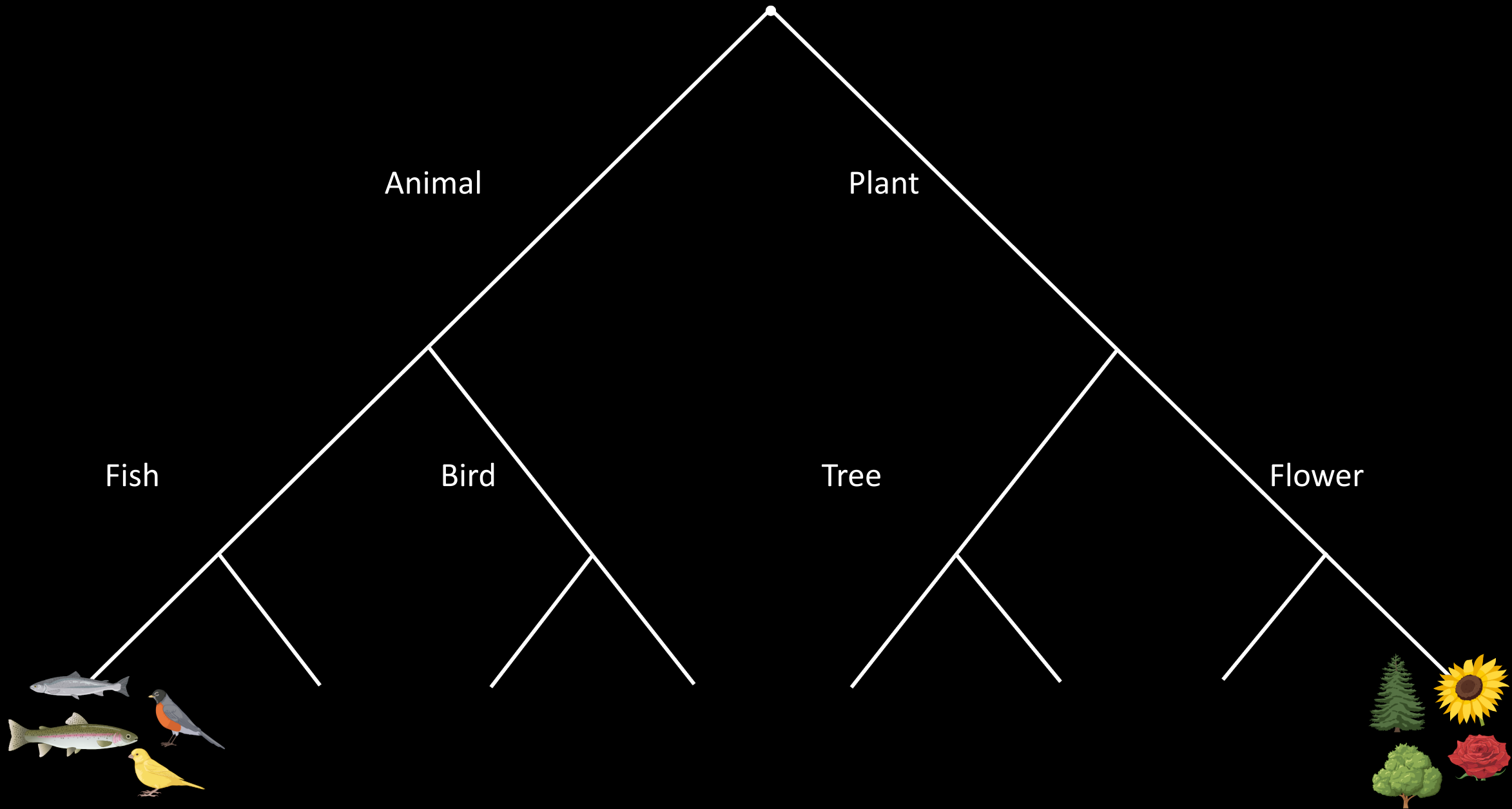
These networks **must** exhibit progressive differentiation:

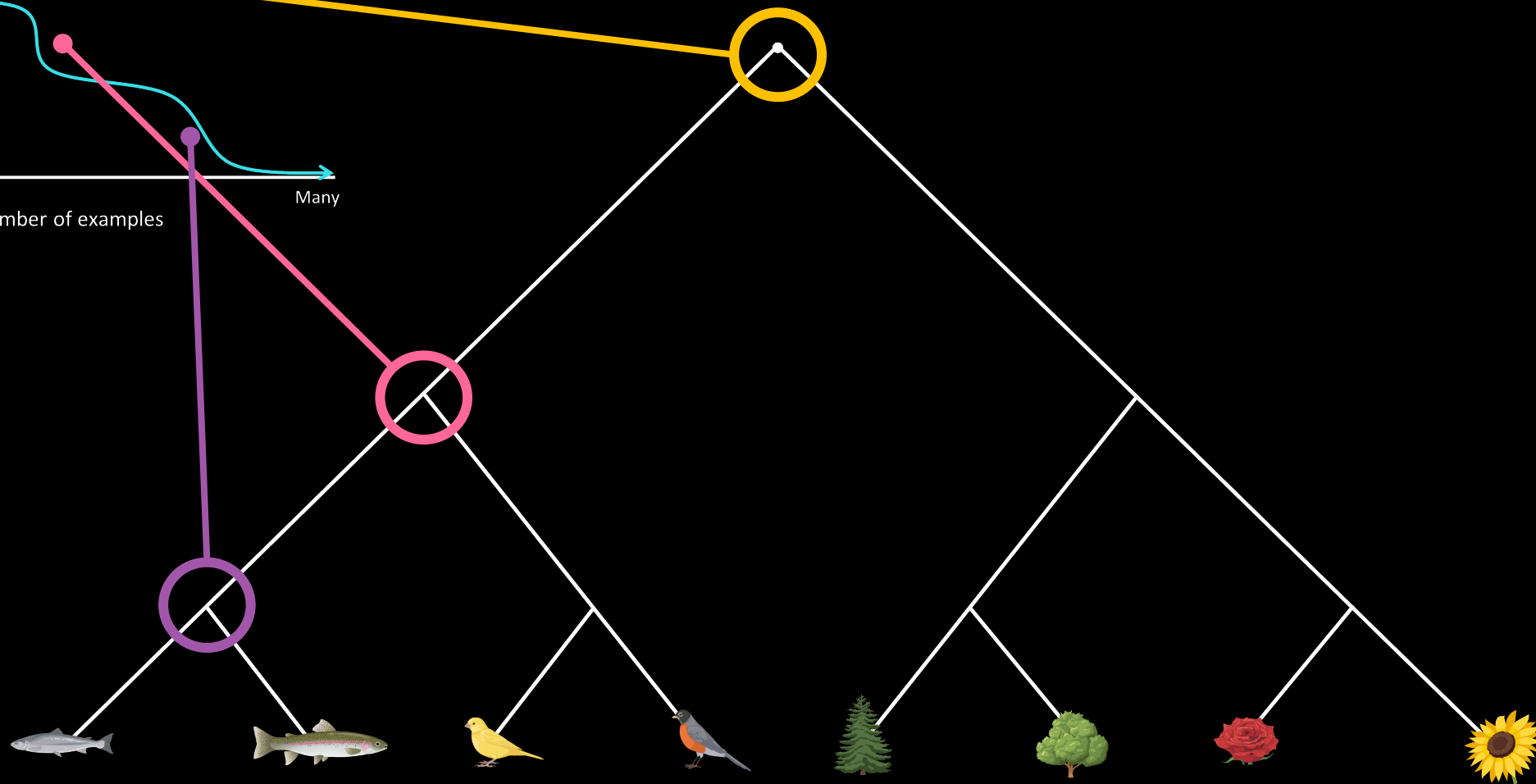
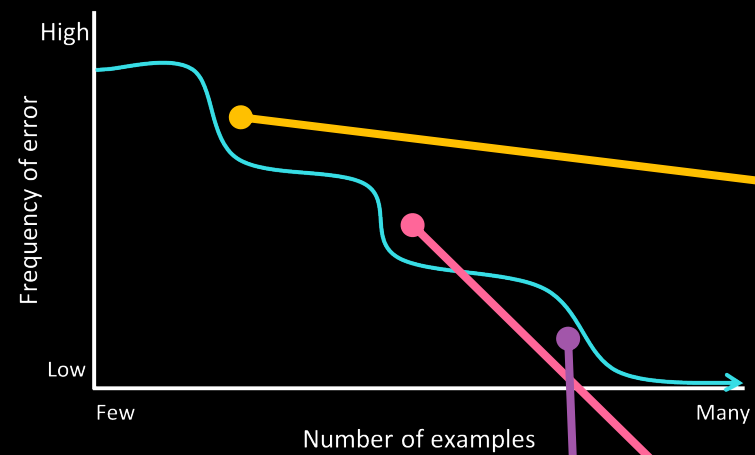
- Singular vectors mirror hierarchy
- Singular values decay with depth



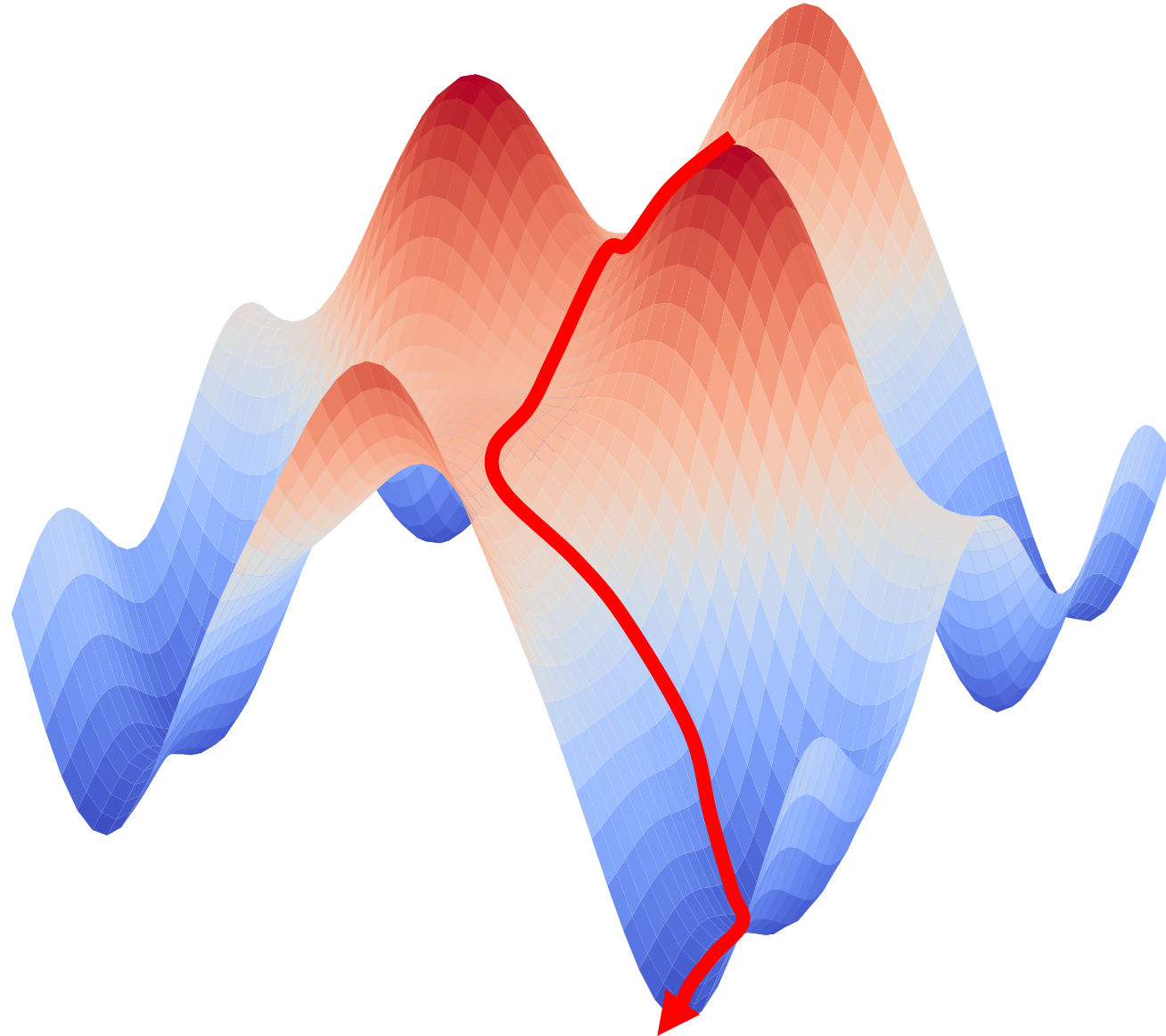
Progressive differentiation



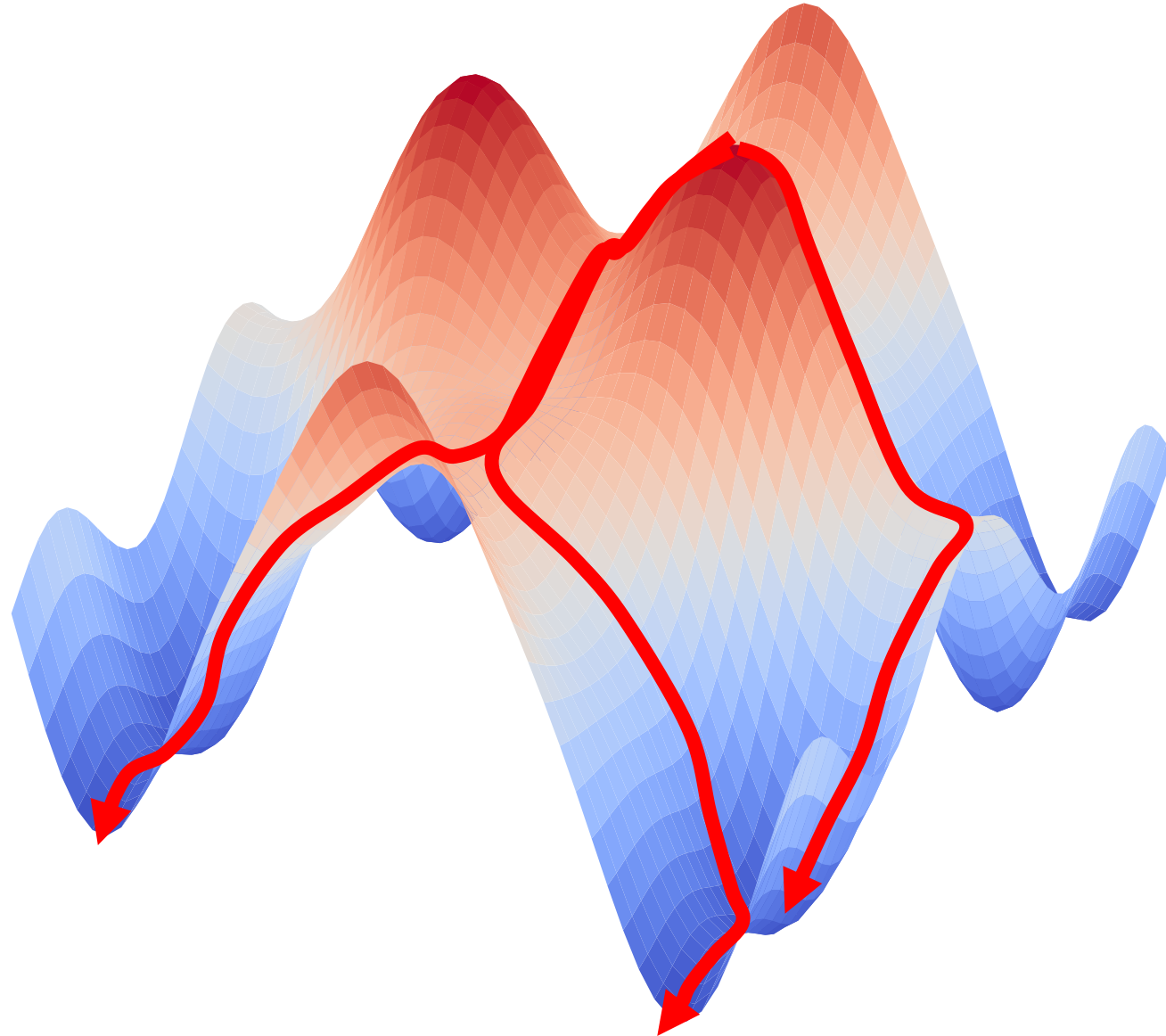




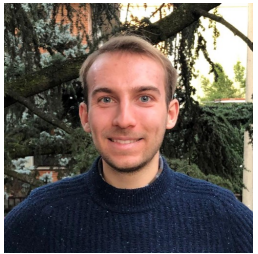
Depth introduces a hierarchy of saddle points



Individual variability amidst structure



Learning to make perceptual decisions from naïve to expert



Sam Liebana



Aeron Laffere



Chiara Toschi



Peter
Zatka-Haas



Louisa Schilling

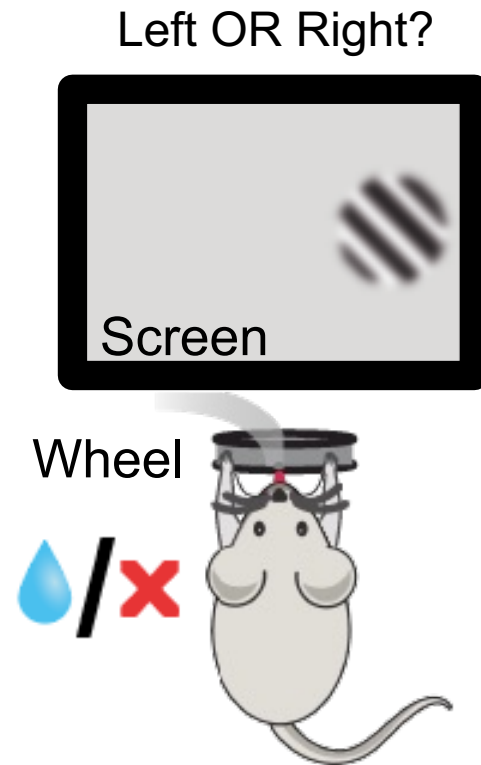







Rafal Bogacz



Armin Lak

Learning to make perceptual decisions from naïve to expert

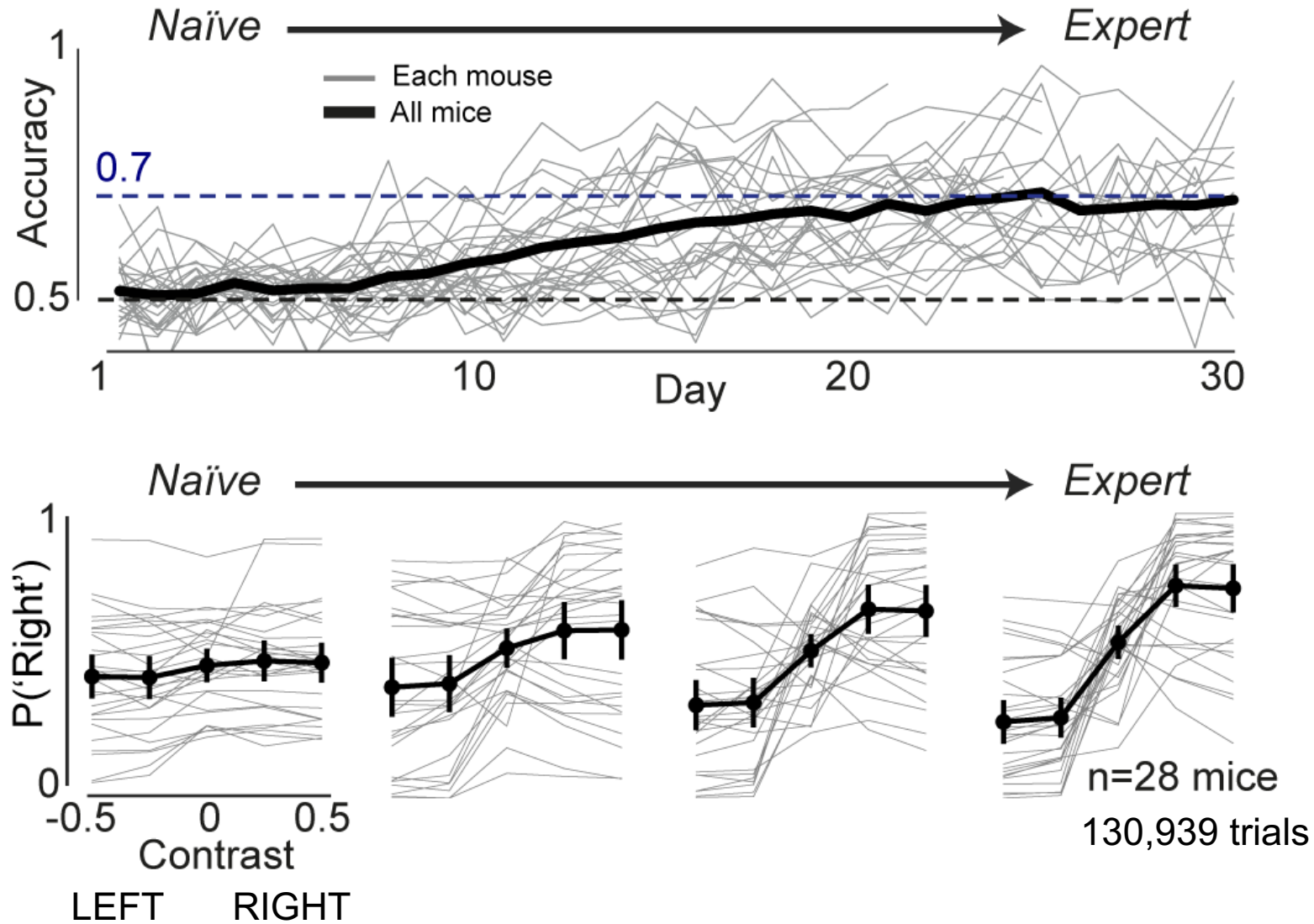


	-0.5	-0.25	0	0.25	0.5
Stimulus contrast					
Rewarded action	L	L	L/R	R	R

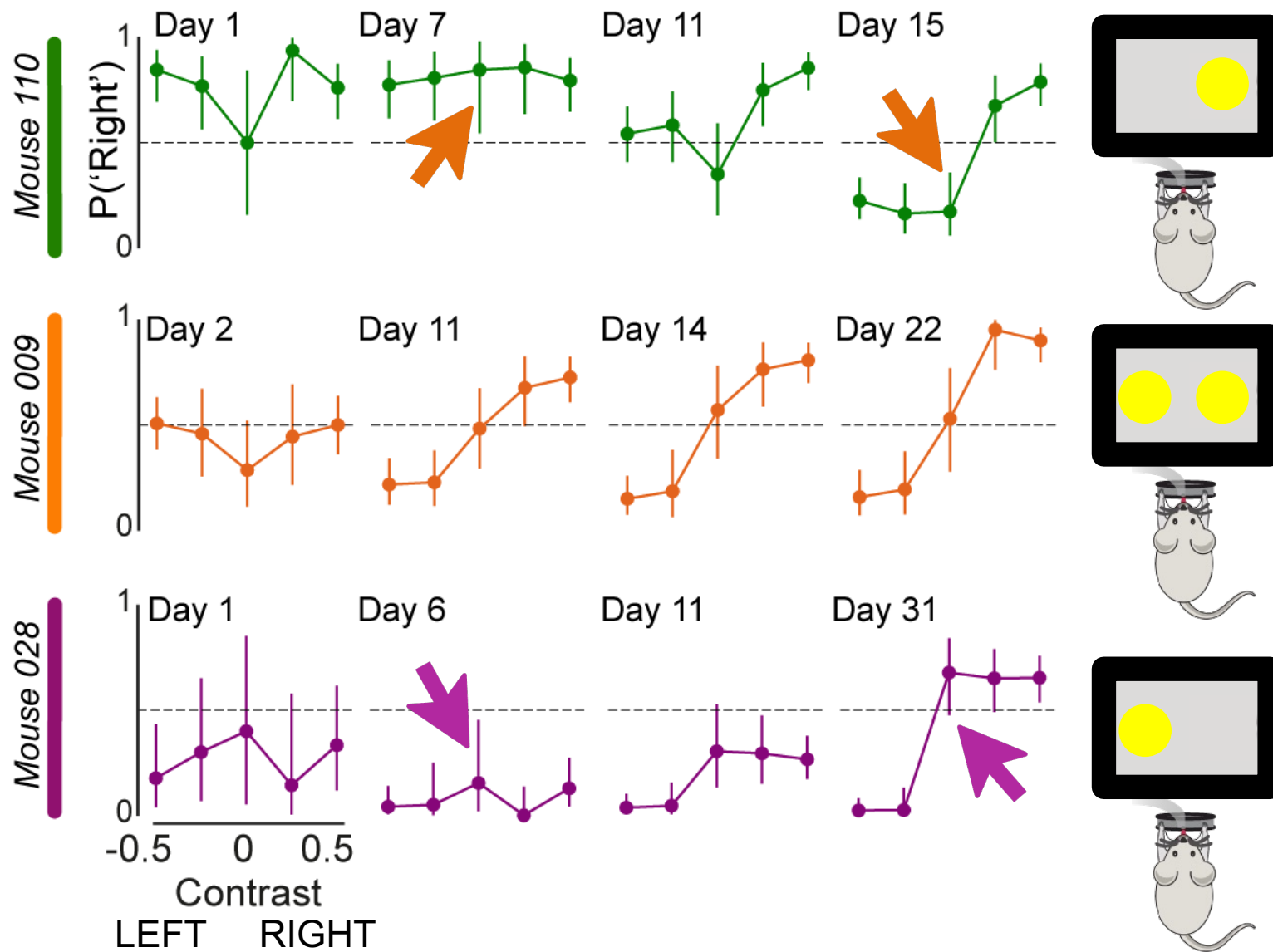
Reward  / 

Full task from day 1 without any change over learning

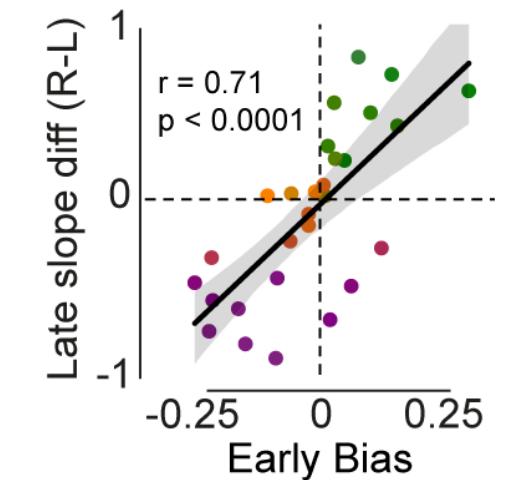
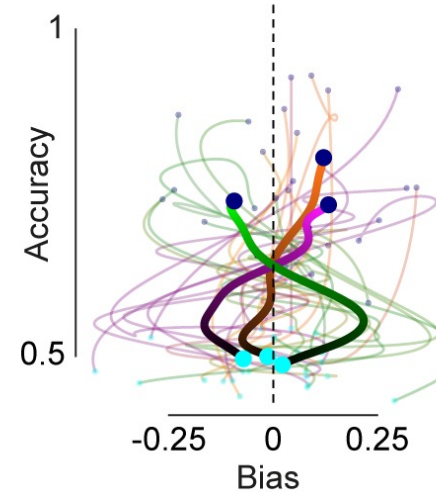
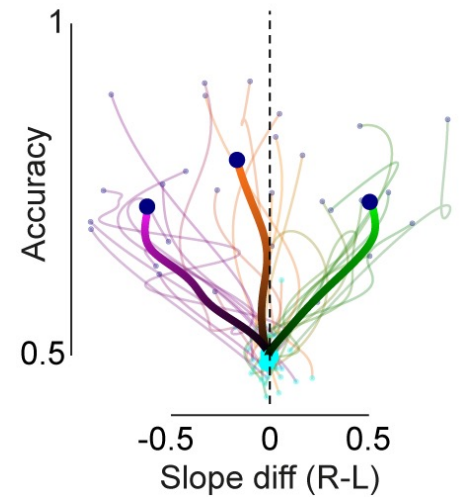
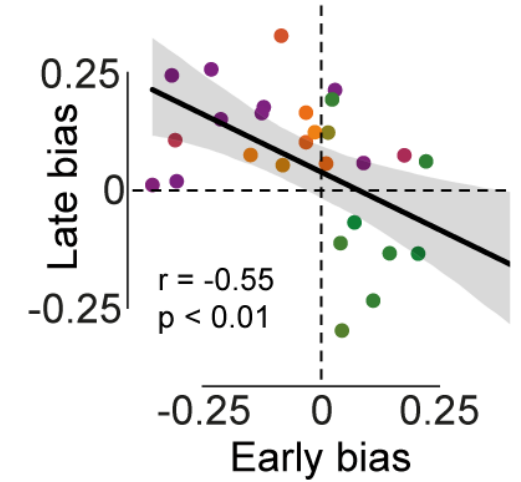
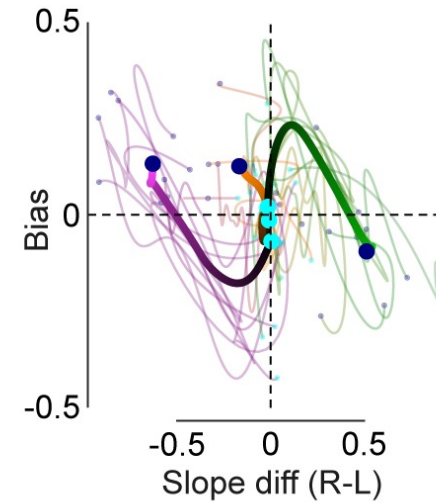
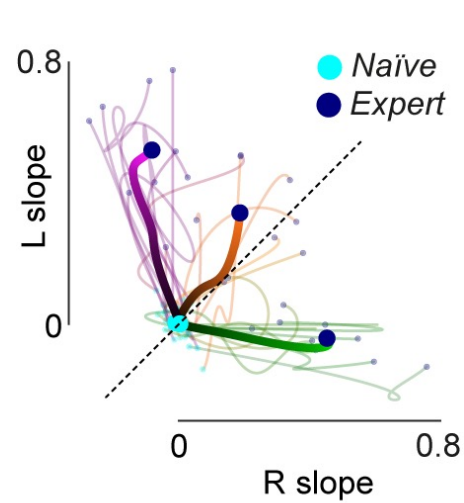
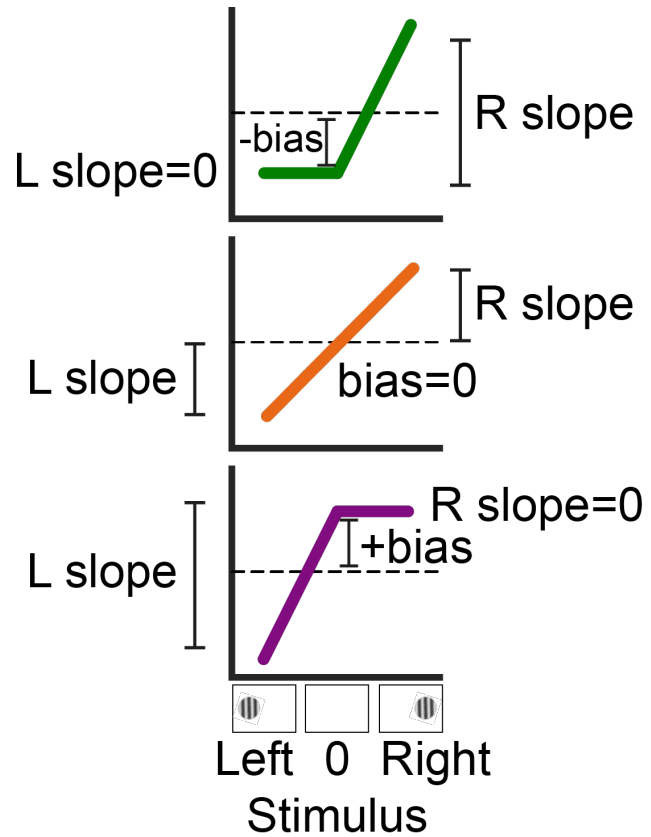
Learning to make perceptual decisions from naïve to expert



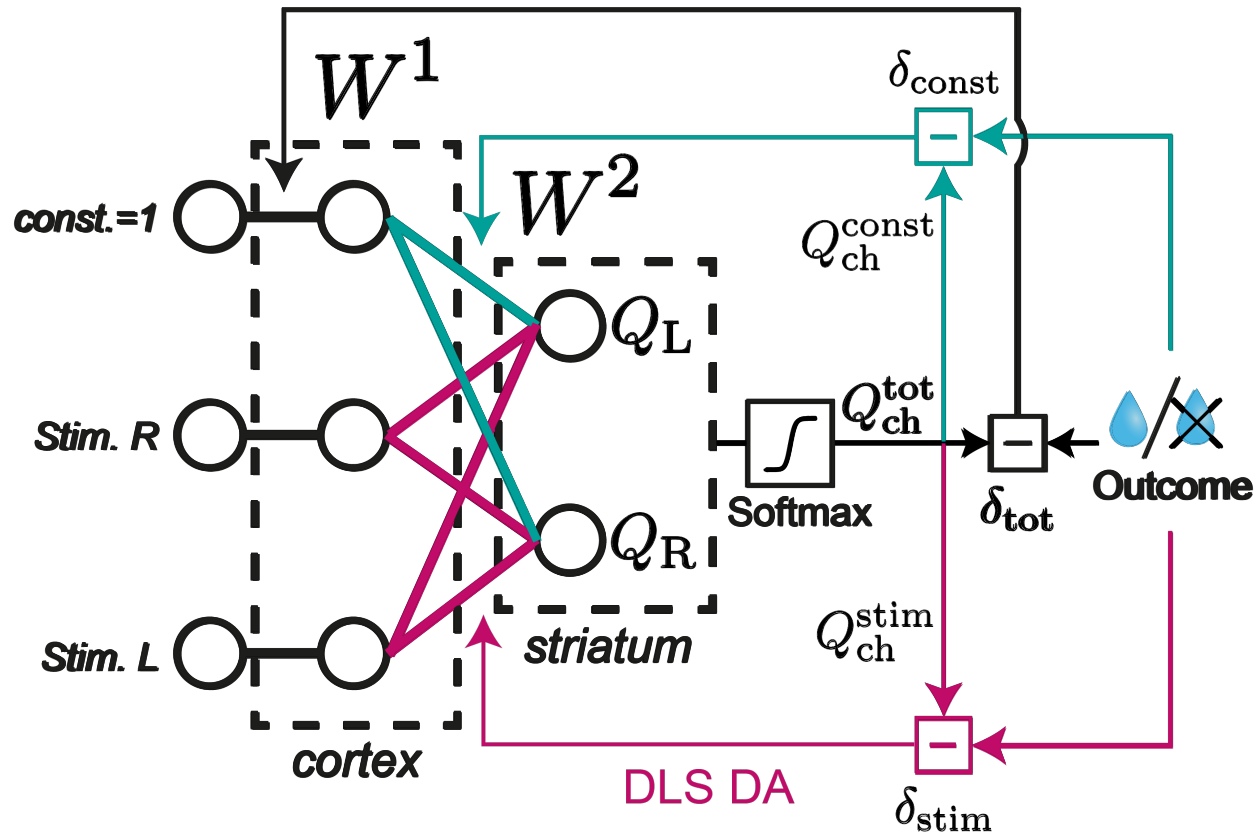
Mice exhibit diverse learning trajectories



Learning trajectories are individually diverse but systematic



A Deep RL Neural Network



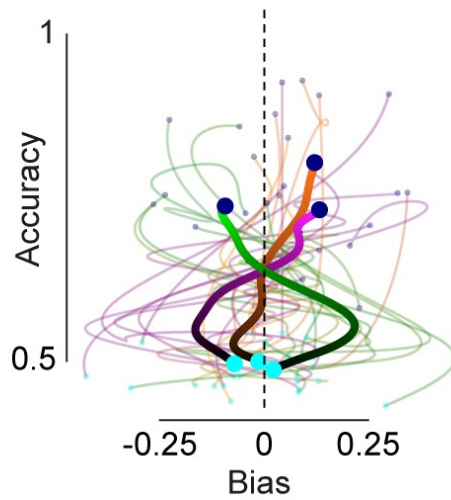
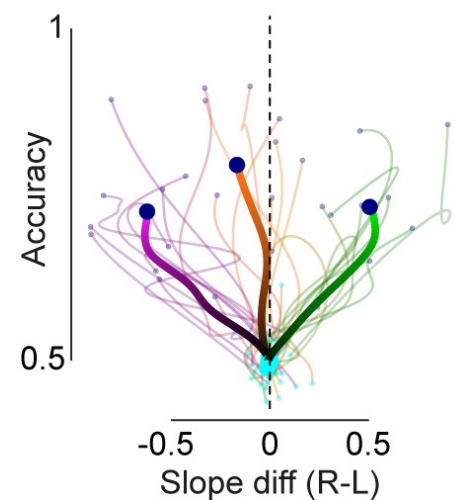
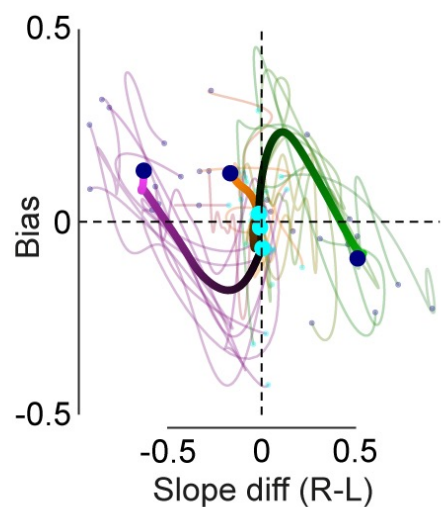
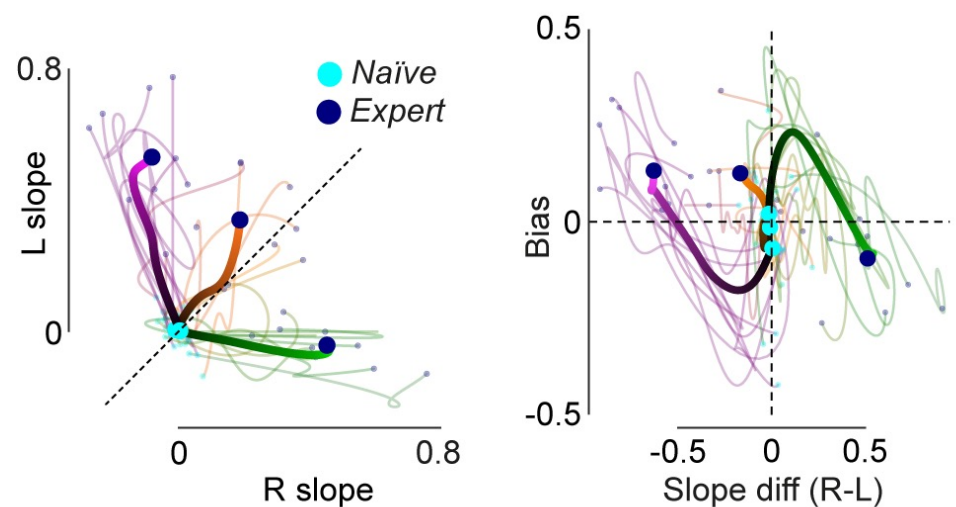
$$\mathcal{L}^{cortex} = \frac{1}{2} \delta_{tot}^2 = \frac{1}{2} (\text{Rew} - Q_{ch}^{tot})^2$$

$$\mathcal{L}^{const} = \frac{1}{2} \delta_{const}^2 = \frac{1}{2} (\text{Rew} - Q_{ch}^{const})^2$$

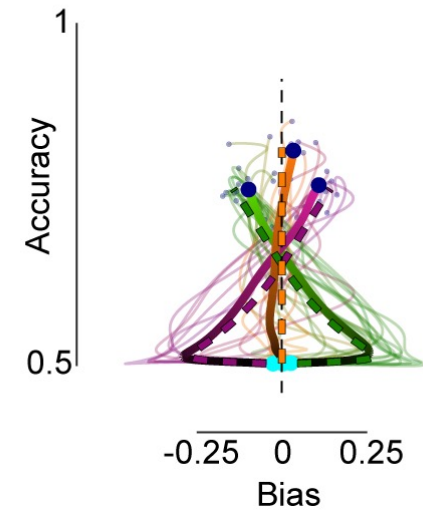
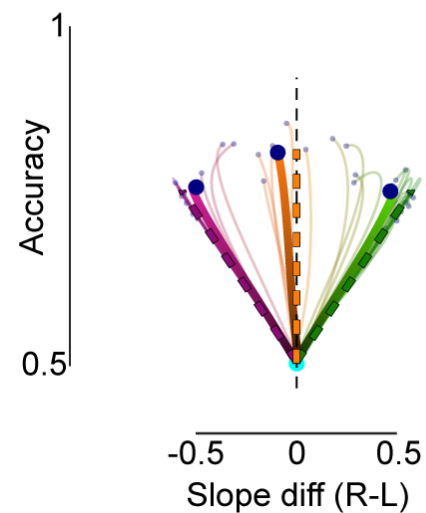
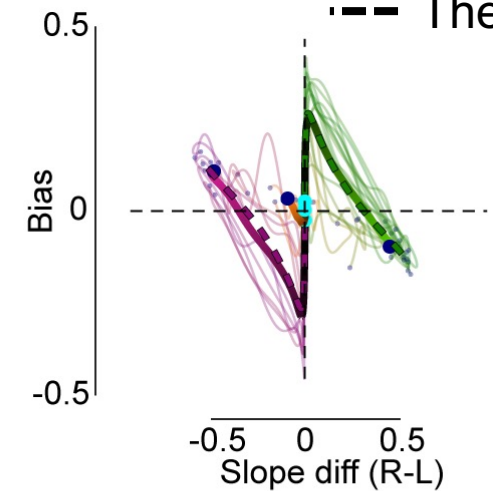
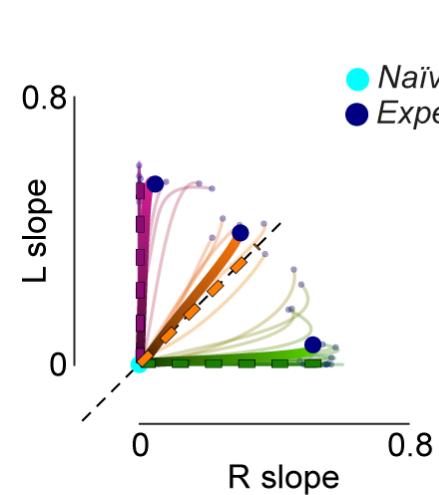
$$\mathcal{L}^{stim} = \frac{1}{2} \delta_{stim}^2 = \frac{1}{2} (\text{Rew} - Q_{ch}^{stim})^2$$

Model captures behavior

Behavior

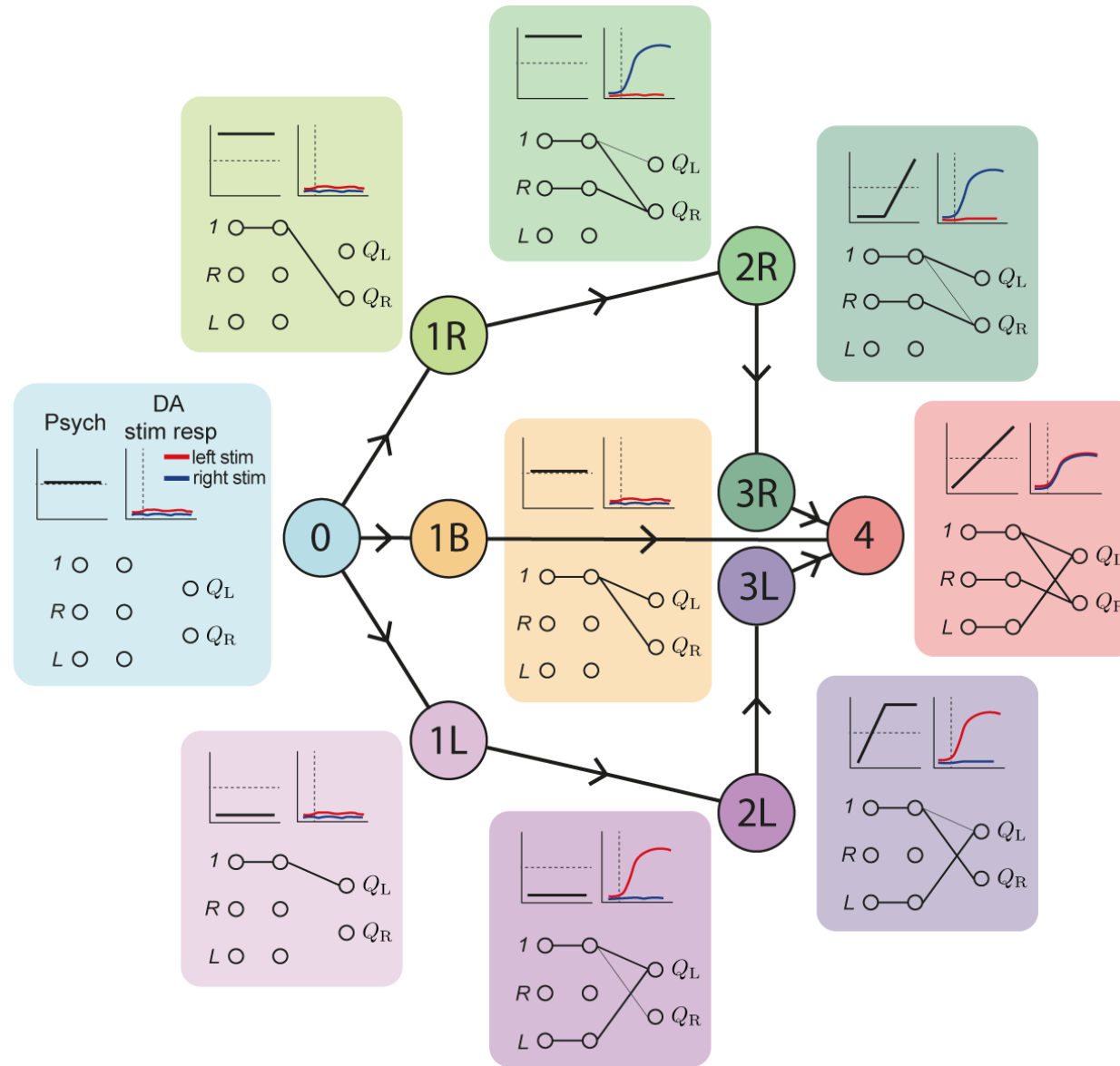


Model



— Simulation
- - - Theory

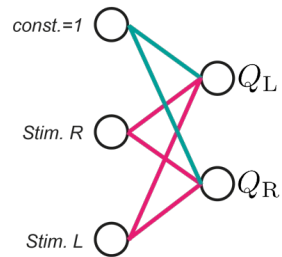
Dynamics pass near a hierarchy of saddle points



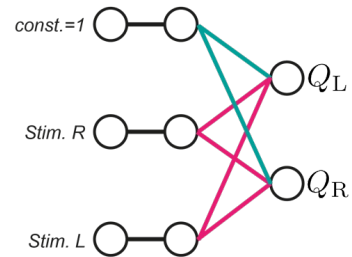
Saddle points arise through depth

Architecture

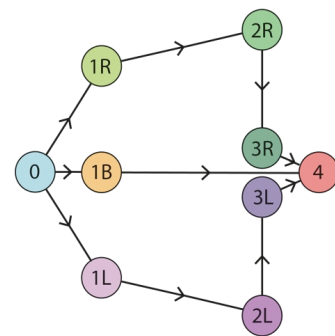
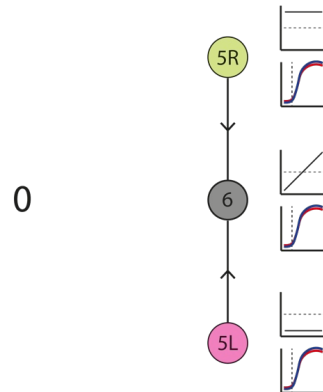
Shallow



Deep

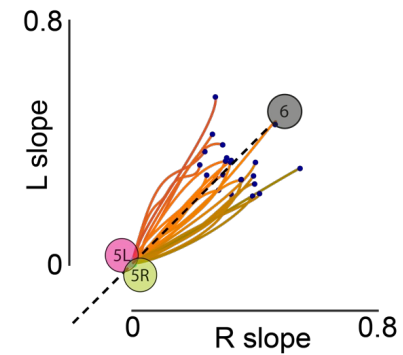


Stationary Points

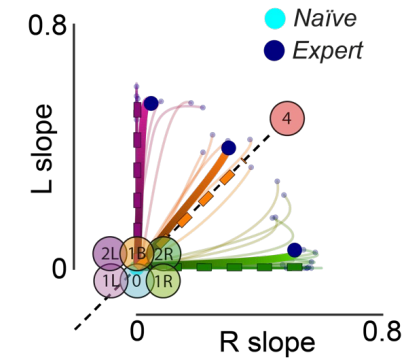


Behavioral Trajectories

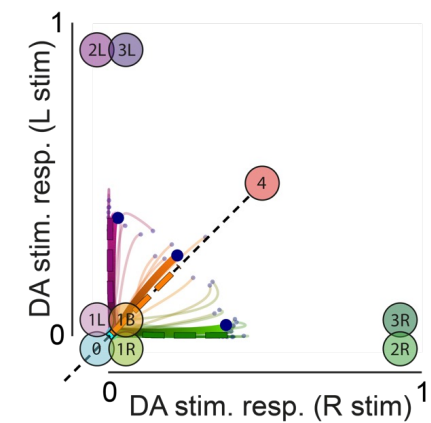
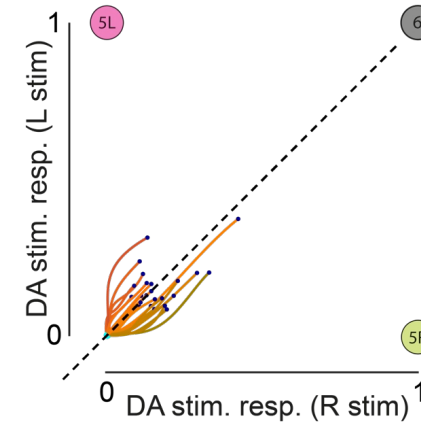
Shallow



Deep



Neural Trajectories



Today

1. Deep linear network dynamics from *tabula rasa* initialization
2. Nontrivial initializations: Lazy, rich, & beyond
3. Nonlinear networks & the neural race reduction

Partitioned solution

$$\mathbf{Q}\mathbf{Q}^T(t) = \begin{pmatrix} \mathbf{Z}_1(t)\mathbf{A}^{-1}(t)\mathbf{Z}_1^T(t) & \mathbf{Z}_1(t)\mathbf{A}^{-1}(t)\mathbf{Z}_2^T(t) \\ \mathbf{Z}_2(t)\mathbf{A}^{-1}(t)\mathbf{Z}_1^T(t) & \mathbf{Z}_2(t)\mathbf{A}^{-1}(t)\mathbf{Z}_2^T(t) \end{pmatrix},$$

with the time-dependent variables $\mathbf{Z}_1(t) \in \mathbb{R}^{N_i \times N_h}$, $\mathbf{Z}_2(t) \in \mathbb{R}^{N_o \times N_h}$, and $\mathbf{A}(t) \in \mathbb{R}^{N_h \times N_h}$:

$$\mathbf{Z}_1(t) = \frac{1}{2}\tilde{\mathbf{V}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{B}^T - \frac{1}{2}\tilde{\mathbf{V}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{V}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \mathbf{D}^T, \quad (13)$$

$$\mathbf{Z}_2(t) = \frac{1}{2}\tilde{\mathbf{U}}(\tilde{\mathbf{G}} + \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{B}^T + \frac{1}{2}\tilde{\mathbf{U}}(\tilde{\mathbf{G}} - \tilde{\mathbf{H}}\tilde{\mathbf{G}})e^{-\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} \mathbf{C}^T + \tilde{\mathbf{U}}_\perp e^{\lambda_\perp \frac{t}{\tau}} \mathbf{D}^T, \quad (14)$$

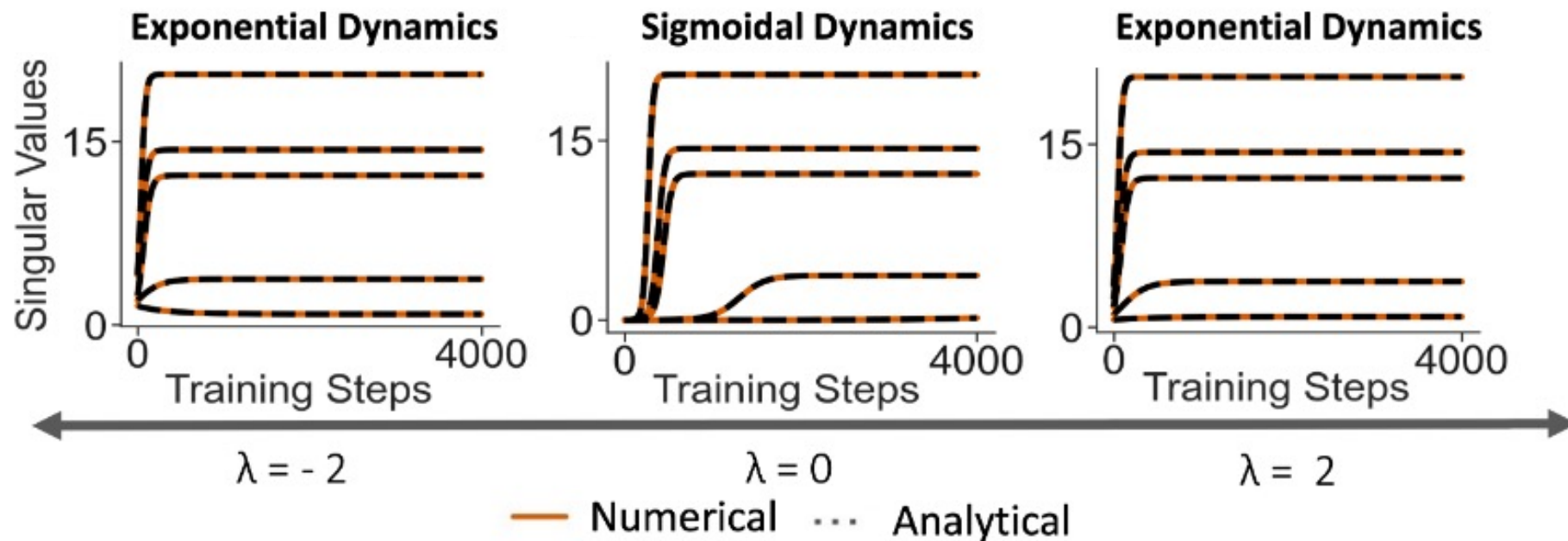
$$\mathbf{A}(t) = \mathbf{I} + \mathbf{B} \left(\frac{e^{2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I}}{4\tilde{\mathbf{S}}_\lambda} \right) \mathbf{B}^T - \mathbf{C} \left(\frac{e^{-2\tilde{\mathbf{S}}_\lambda \frac{t}{\tau}} - \mathbf{I}}{4\tilde{\mathbf{S}}_\lambda} \right) \mathbf{C}^T + \mathbf{D} \left(\frac{e^{\lambda_\perp \frac{t}{\tau}} - \mathbf{I}}{\lambda_\perp} \right) \mathbf{D}^T. \quad (15)$$

and

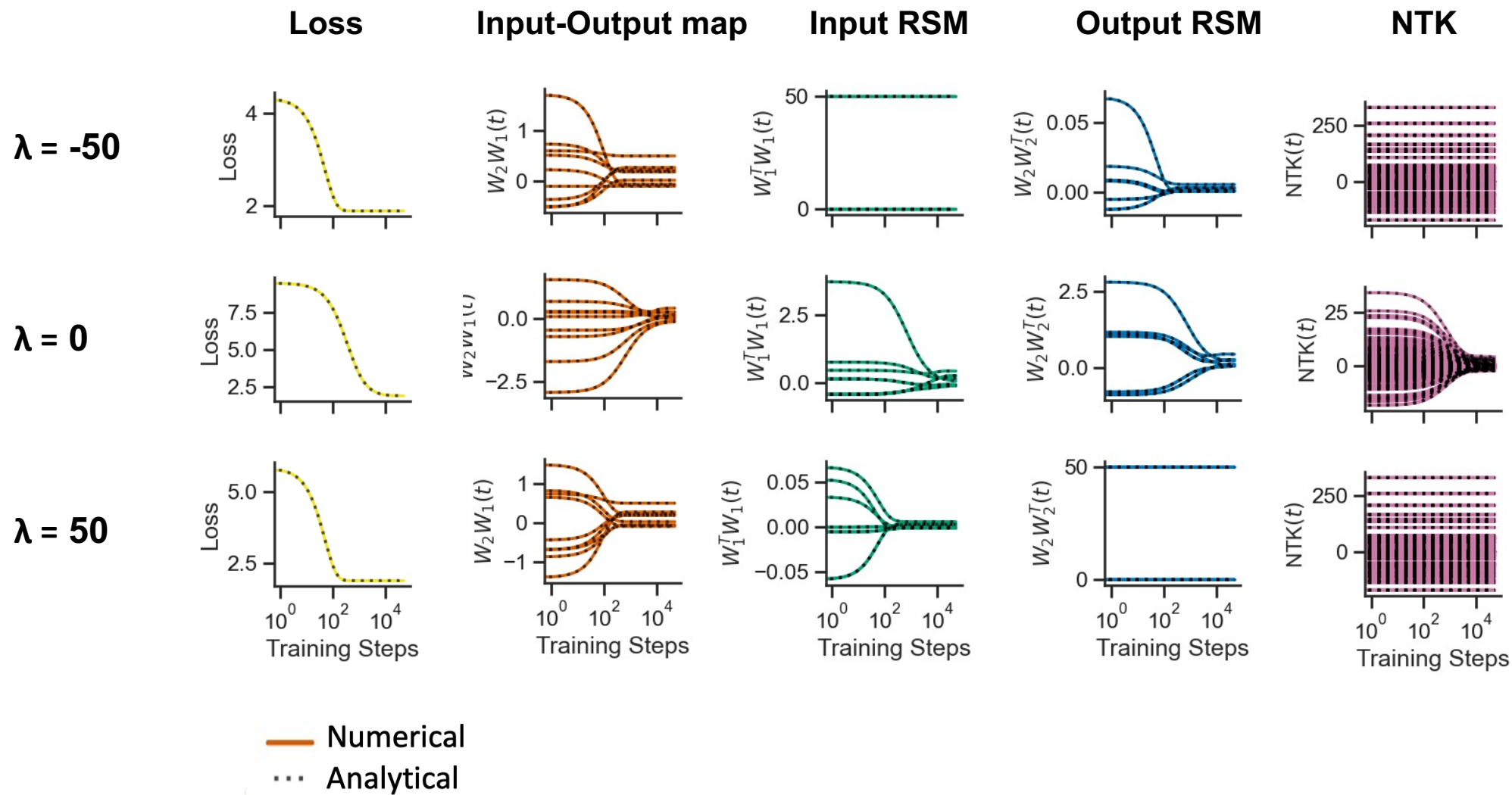
$$\tilde{\mathbf{S}}_\lambda = \sqrt{\tilde{\mathbf{S}}^2 + \frac{\lambda^2}{4}\mathbf{I}}, \quad \lambda_\perp = \text{sgn}(N_o - N_i) \frac{\lambda}{2} \mathbf{I}_{|N_o - N_i|}, \quad \tilde{\mathbf{H}} = \text{sgn}(\lambda) \sqrt{\frac{\tilde{\mathbf{S}}_\lambda - \tilde{\mathbf{S}}}{\tilde{\mathbf{S}}_\lambda + \tilde{\mathbf{S}}}}, \quad \tilde{\mathbf{G}} = \frac{1}{\sqrt{\mathbf{I} + \tilde{\mathbf{H}}^2}}.$$

where \mathbf{B} , \mathbf{C} , \mathbf{D} are initialization-dependent matrices.

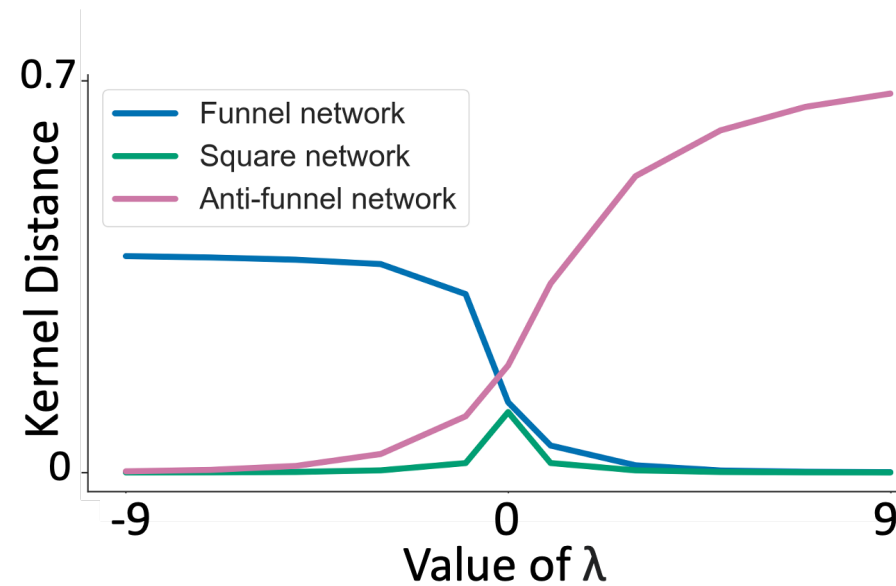
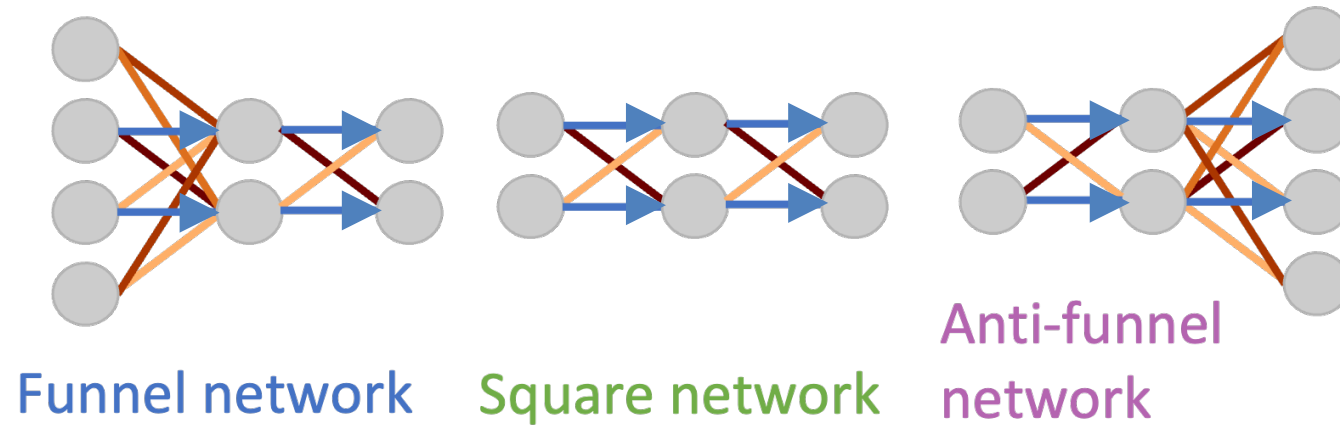
From exponential to sigmoidal dynamics



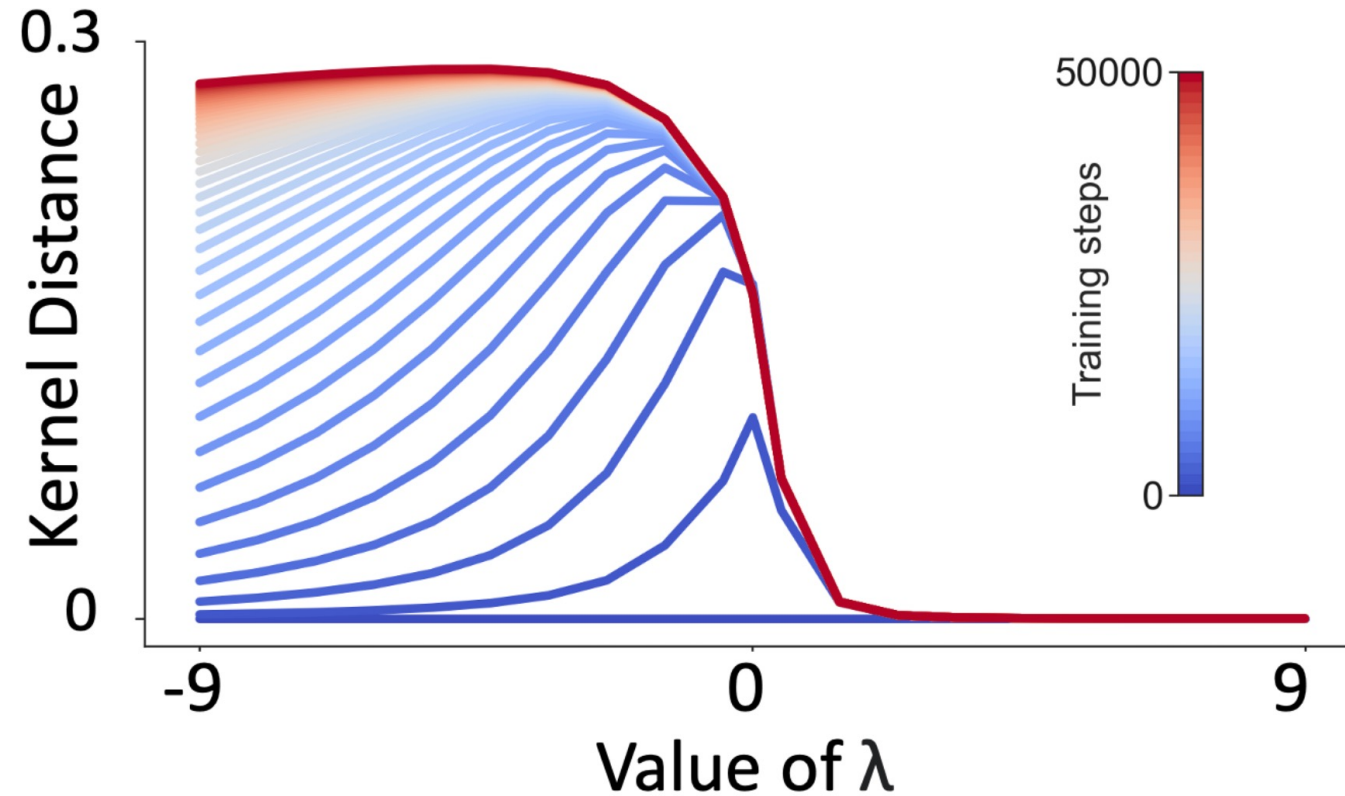
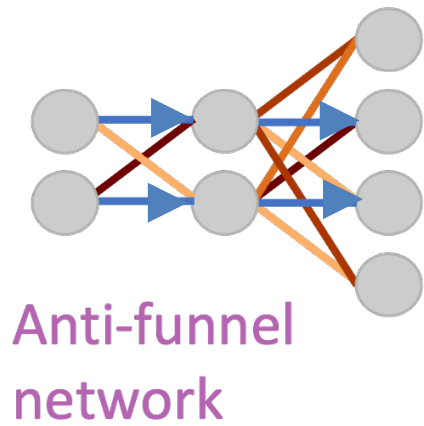
Rich and lazy learning



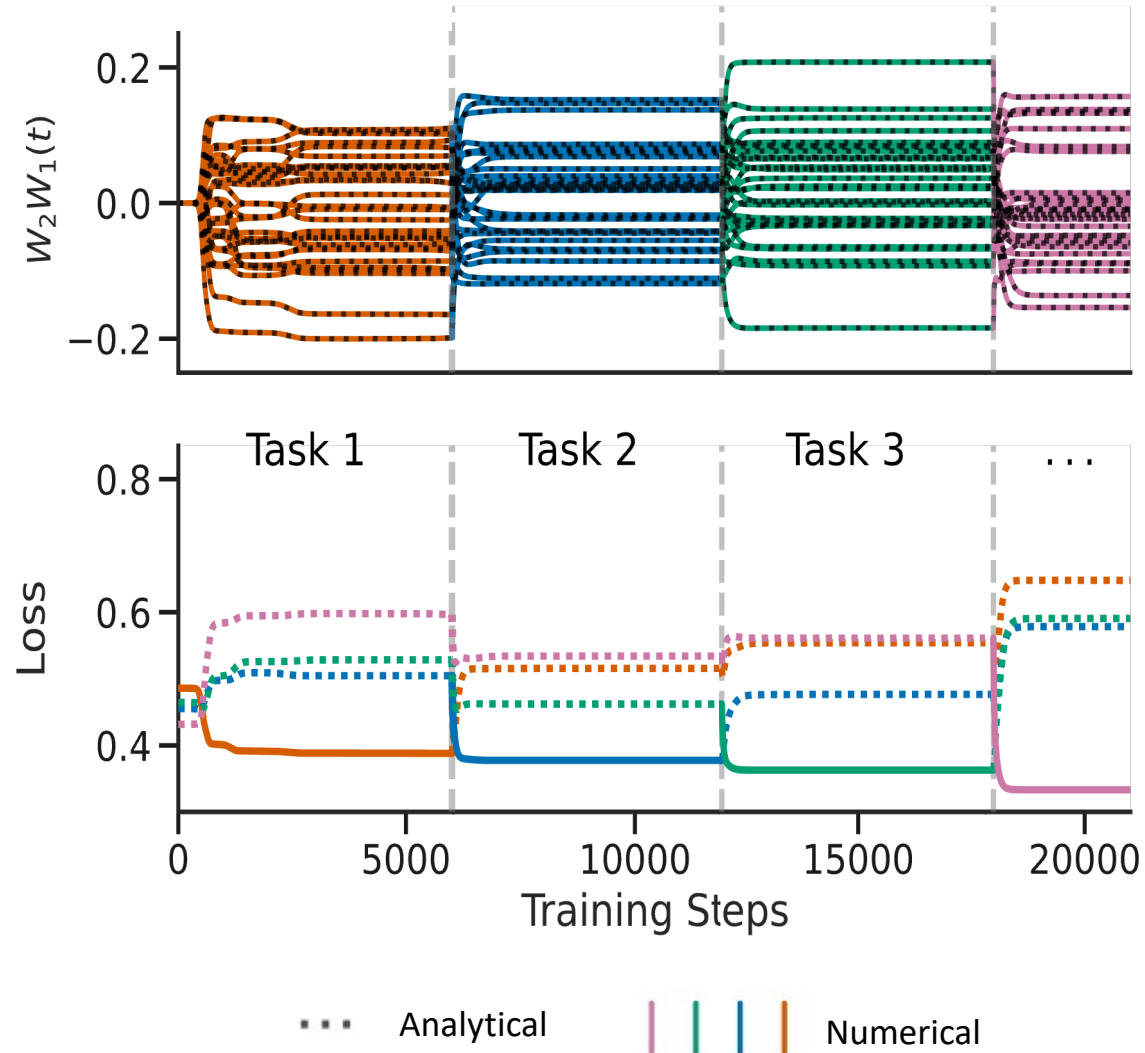
Architecture and learning regime



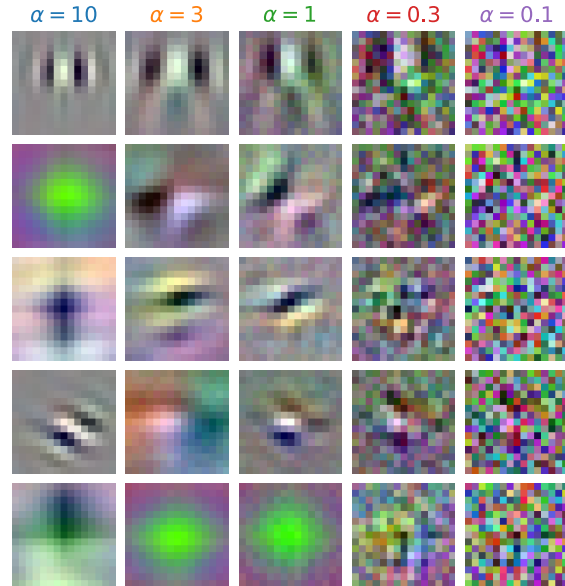
Delayed rich regime



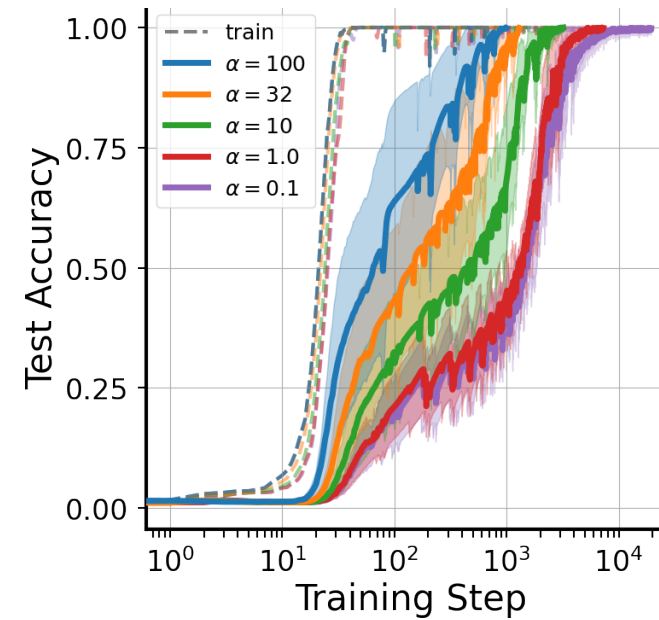
Exact continual learning dynamics



Impact of relative scale initializations in practice



Promotes interpretability of early layers in CNNs

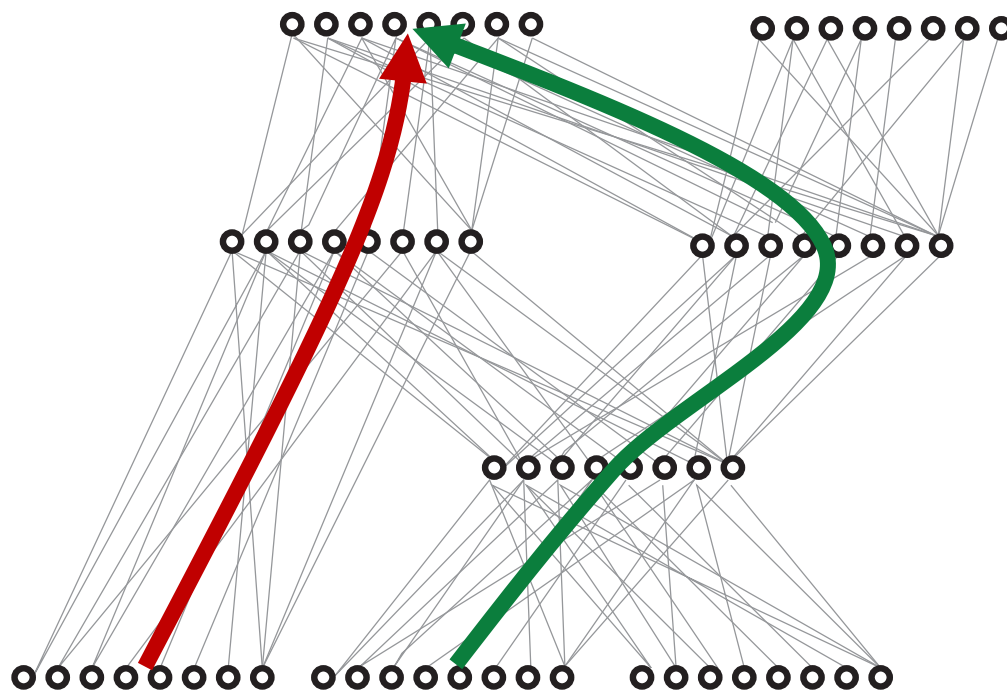


Decreases the time to grokking in modular arithmetic

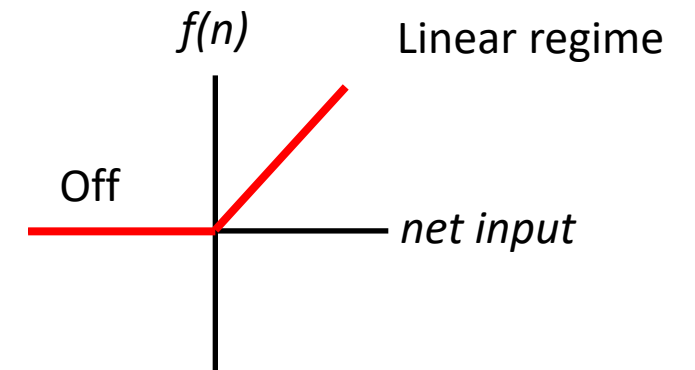
Today

1. Deep linear network dynamics from *tabula rasa* initialization
2. Nontrivial initializations: Lazy, rich, & beyond
3. Nonlinear networks & the neural race reduction

Gating: a simple view of nonlinearity



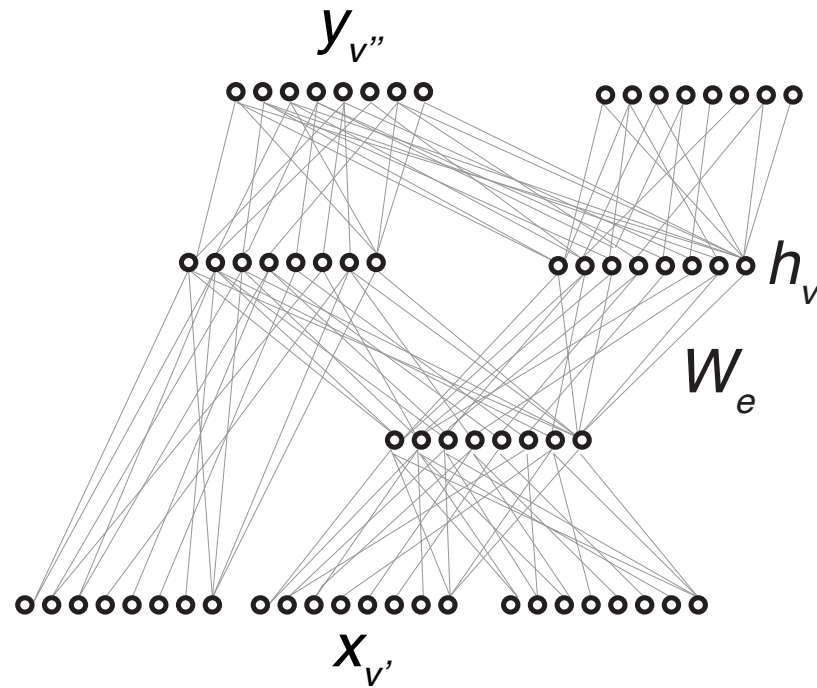
ReLU neural nonlinearity



When active, each pathway behaves like a deep linear network

Gated Deep Linear Network

Arch graph Γ : nodes V , edges E



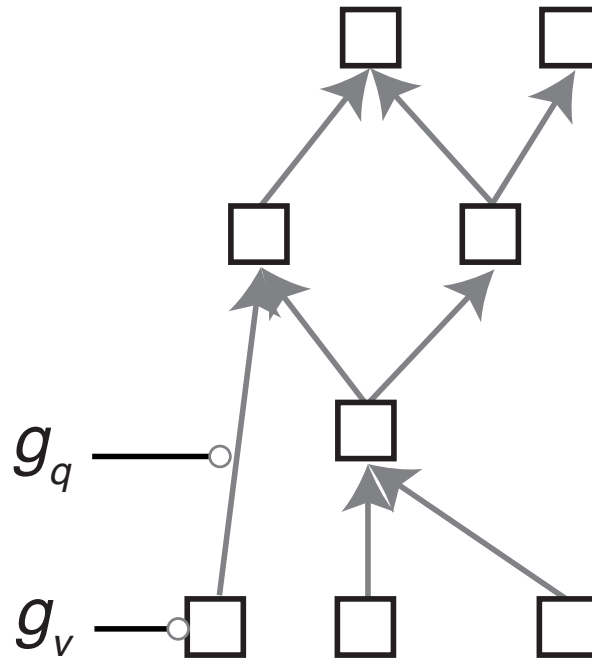
$y_{v''}$

h_v

W_e

$x_{v'}$

Gated Deep Linear Network



Forward propagation:

$$h_v = g_v \sum_{q \in E: t(q)=v} g_q W_q h_{s(q)}$$

$s(q)$: source node of edge q
 $t(q)$: target node of edge q

Gradient descent

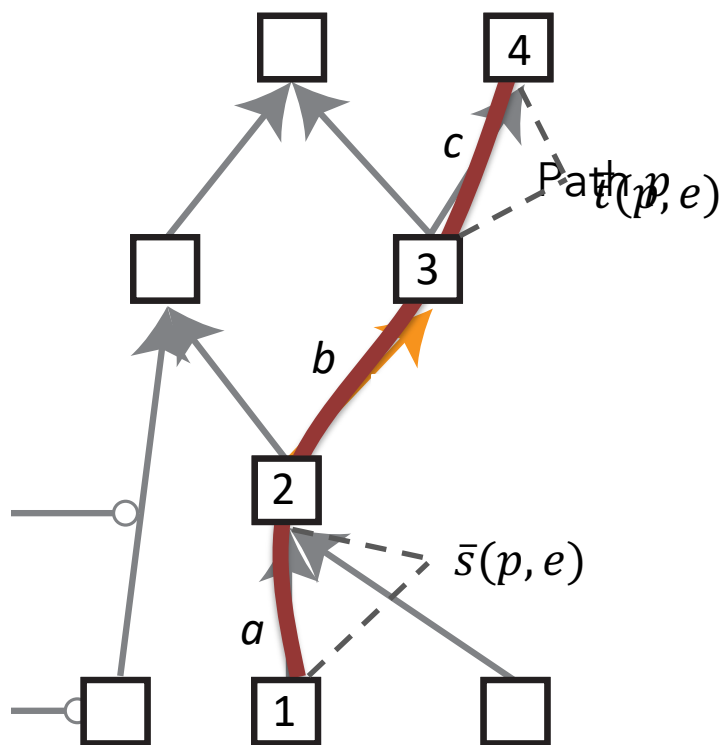
Minimize L_2 loss

$$\mathcal{L}(\{W\}) = \left\langle \frac{1}{2} \sum_{v \in \text{Out}(\Gamma)} \|y_v - h_v\|_2^2 \right\rangle_{x,y,g}$$

using gradient flow on the weights

$$\tau \frac{d}{dt} W_e = - \frac{\partial \mathcal{L}(\{W\})}{\partial W_e} \quad \forall e \in E$$

Gradient descent



Path notation

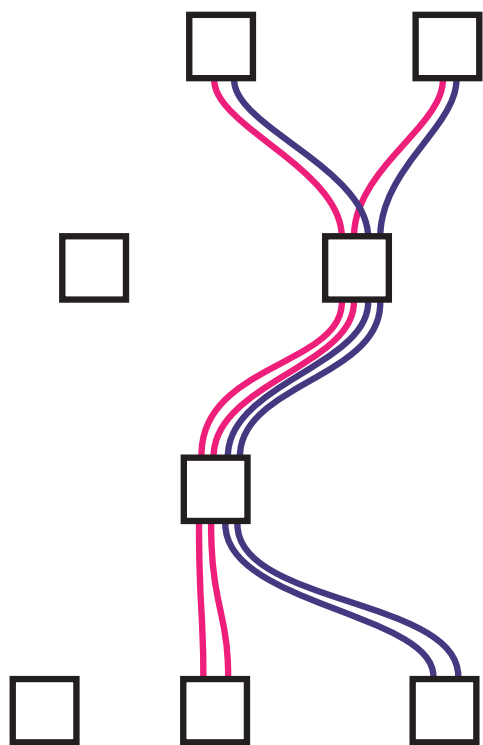
$$W_p = W_c W_b W_a$$

$$g_p = g_4 g_c g_3 g_b g_2 g_a g_1$$

$\bar{t}(p, e)$: target path of e

$\bar{s}(p, e)$: source path of e

Gradient descent



$$\tau \frac{d}{dt} W_e = \sum_{p \in \mathcal{P}(e)} \underbrace{W_{\bar{t}(p,e)}^T \mathcal{E}(p) W_{\bar{s}(p,e)}^T}_{\mathcal{P}(e): \text{All paths through } e}$$

$$\mathcal{E}(p) = \Sigma^{yx}(p) - \sum_{j \in \mathcal{T}(p)} \underbrace{W_j \Sigma^x(j, p)}_{\mathcal{T}(e): \text{All paths terminating at same node as } p}$$

Correlation matrices

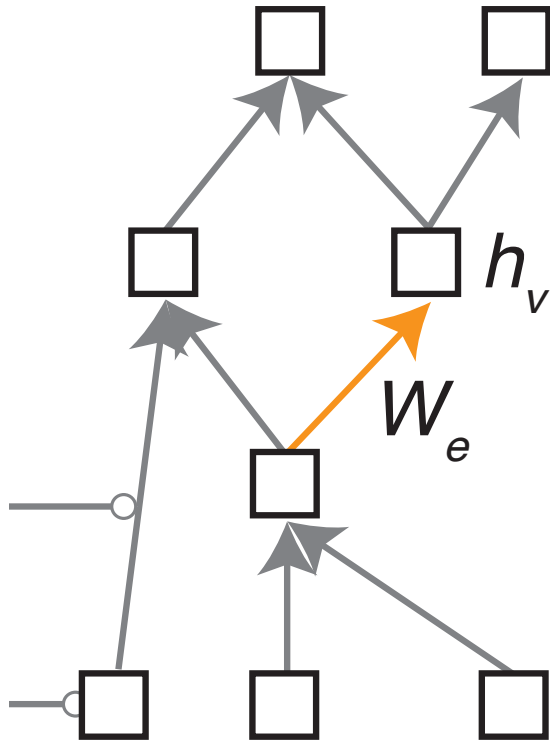
- Dynamics driven only by statistics:

$$\Sigma^{yx}(p) = \langle g_p y_{t(p)} x_{s(p)}^T \rangle_{y,x,g}$$

$$\Sigma^x(j, p) = \langle g_j x_{s(j)} x_{s(p)}^T g_p \rangle_{y,x,g}$$

- One correlation matrix per path

Intuition

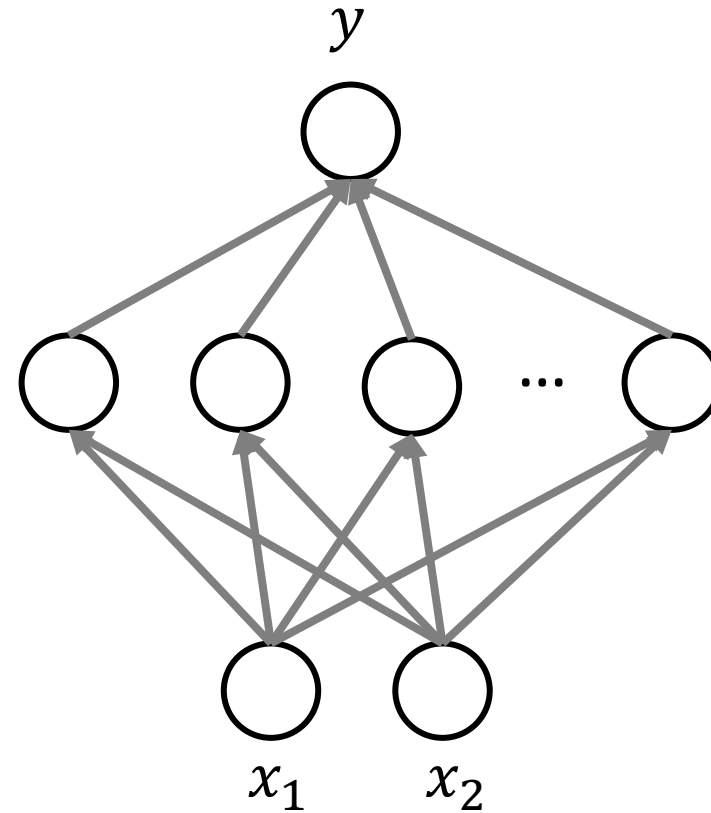
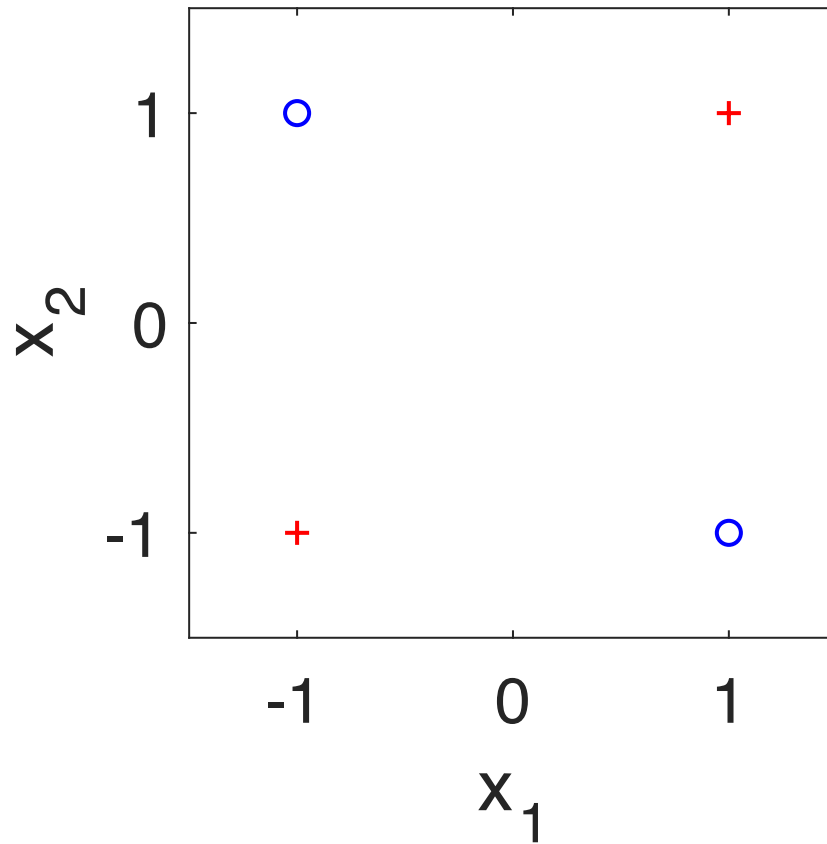


Each pathway behaves like a deep linear network

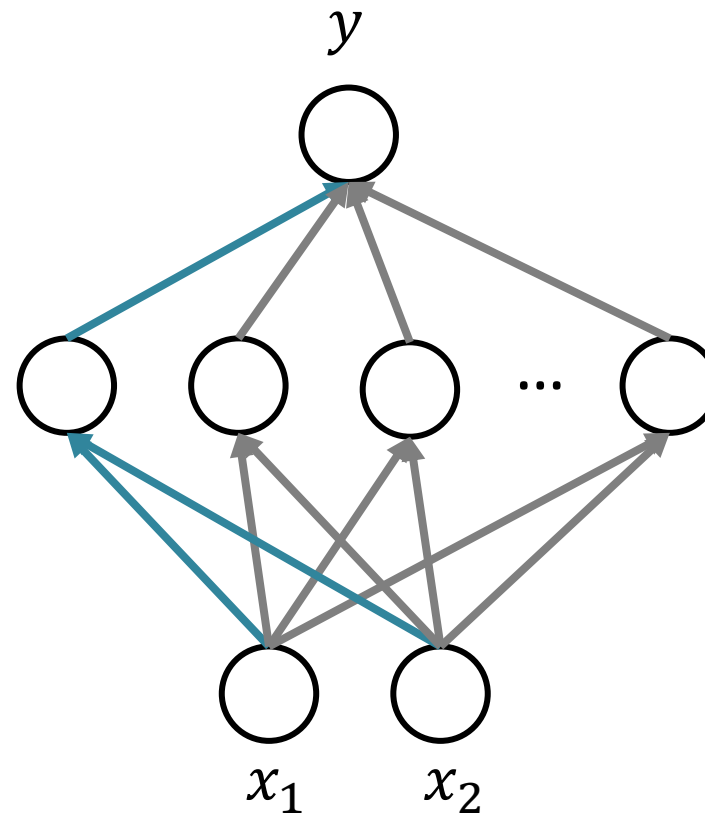
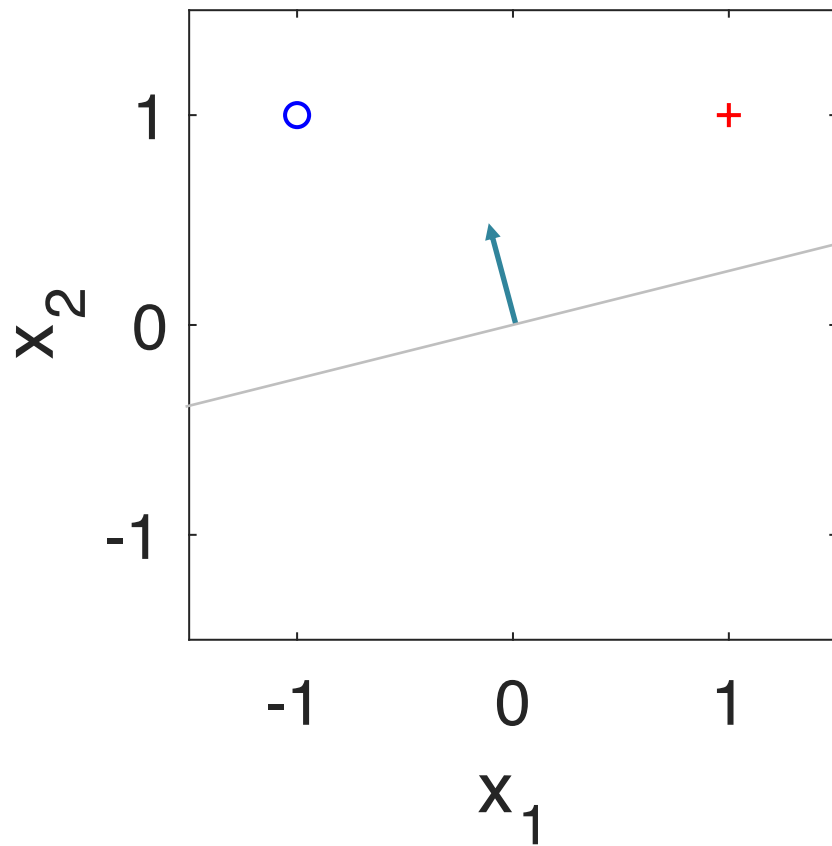
Gating controls the *effective dataset* for each pathway

All paths through an edge sum to determine dynamics

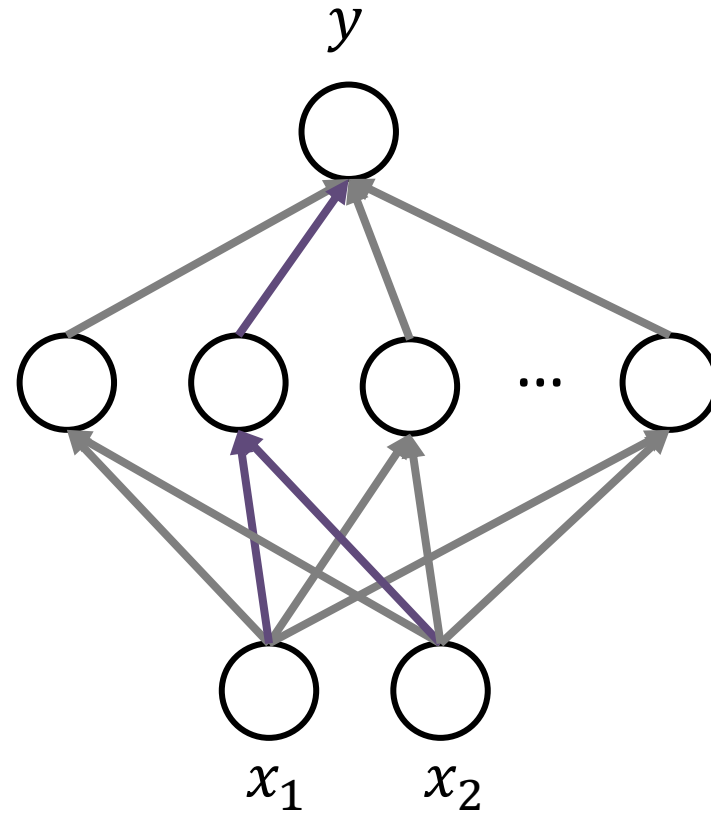
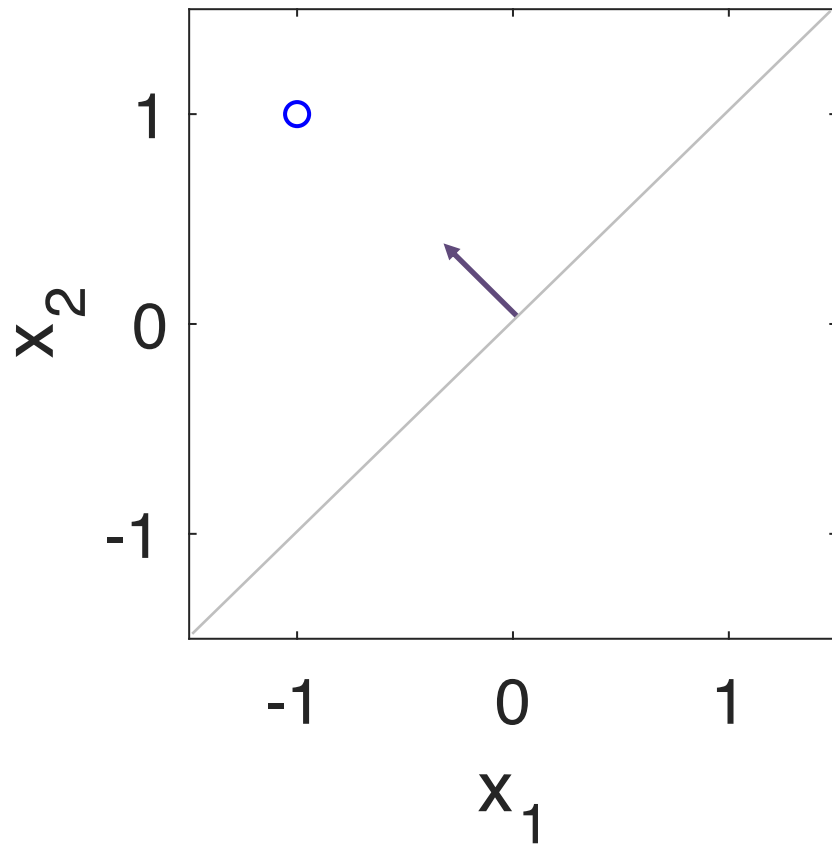
The XoR problem



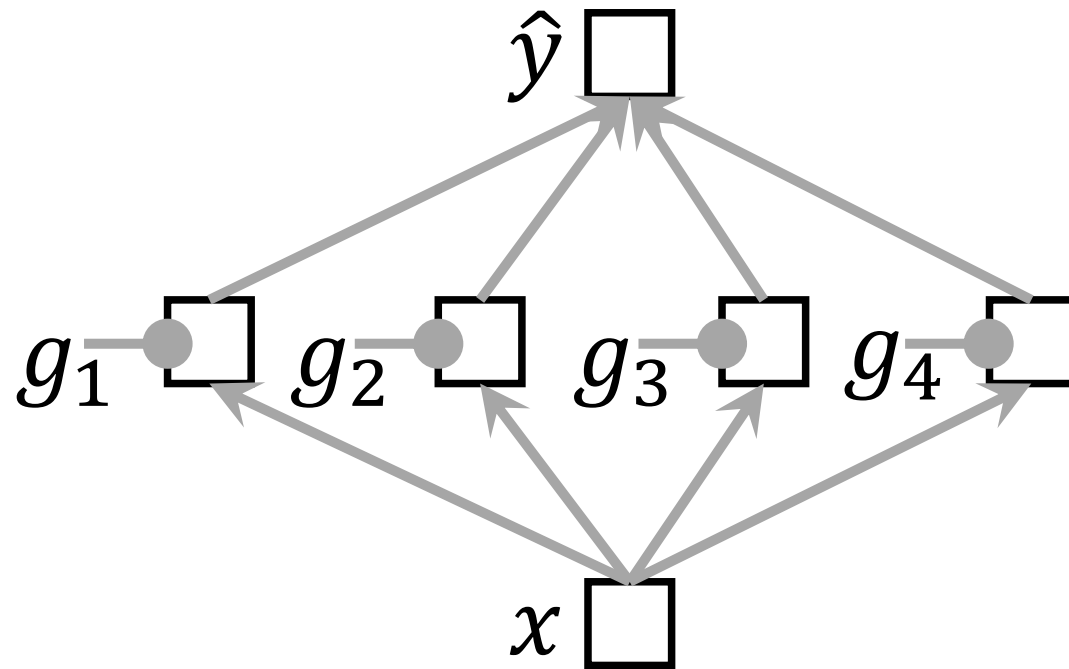
Gated dynamics



Gated dynamics

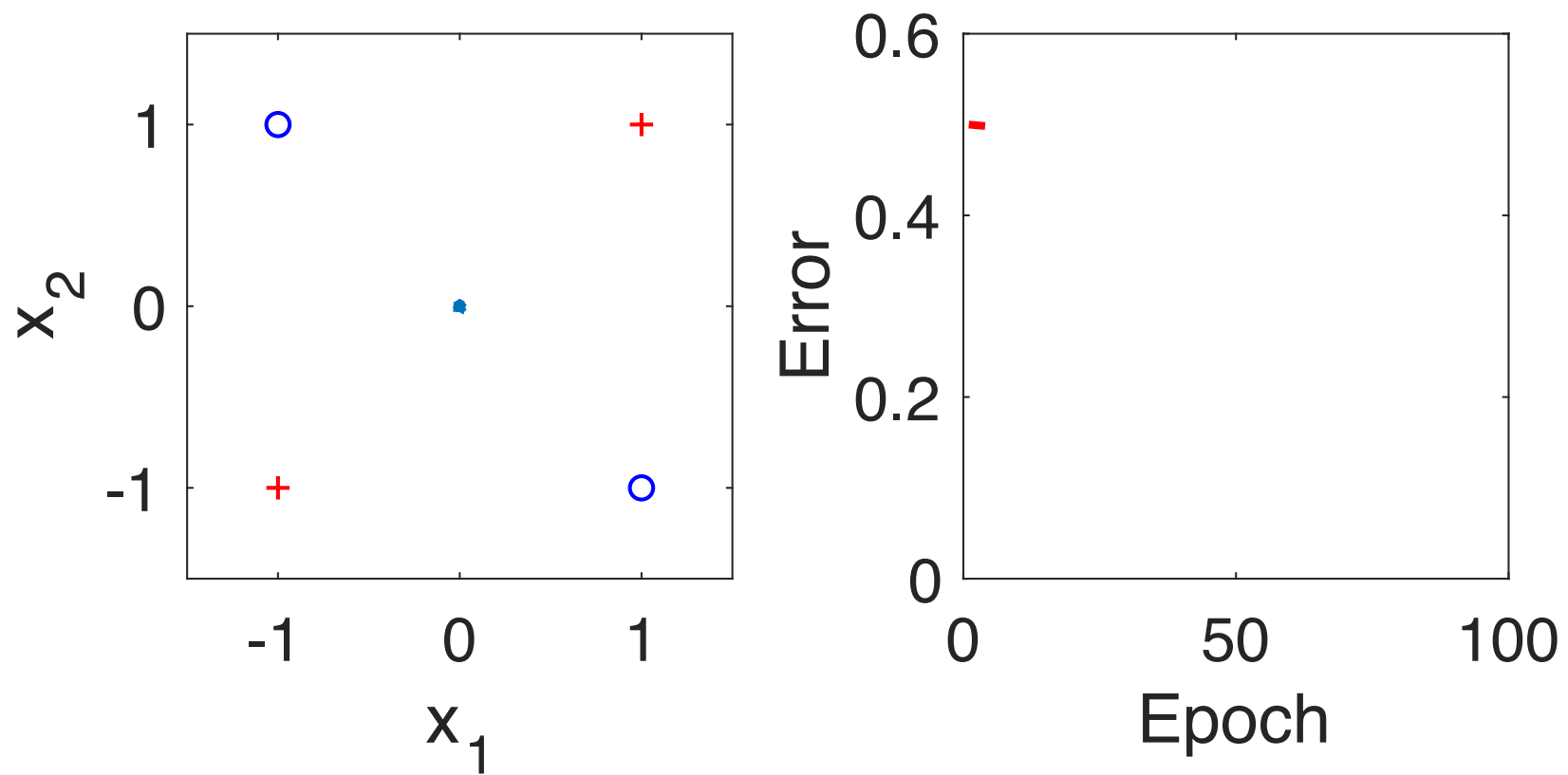


Gated DLN on XoR

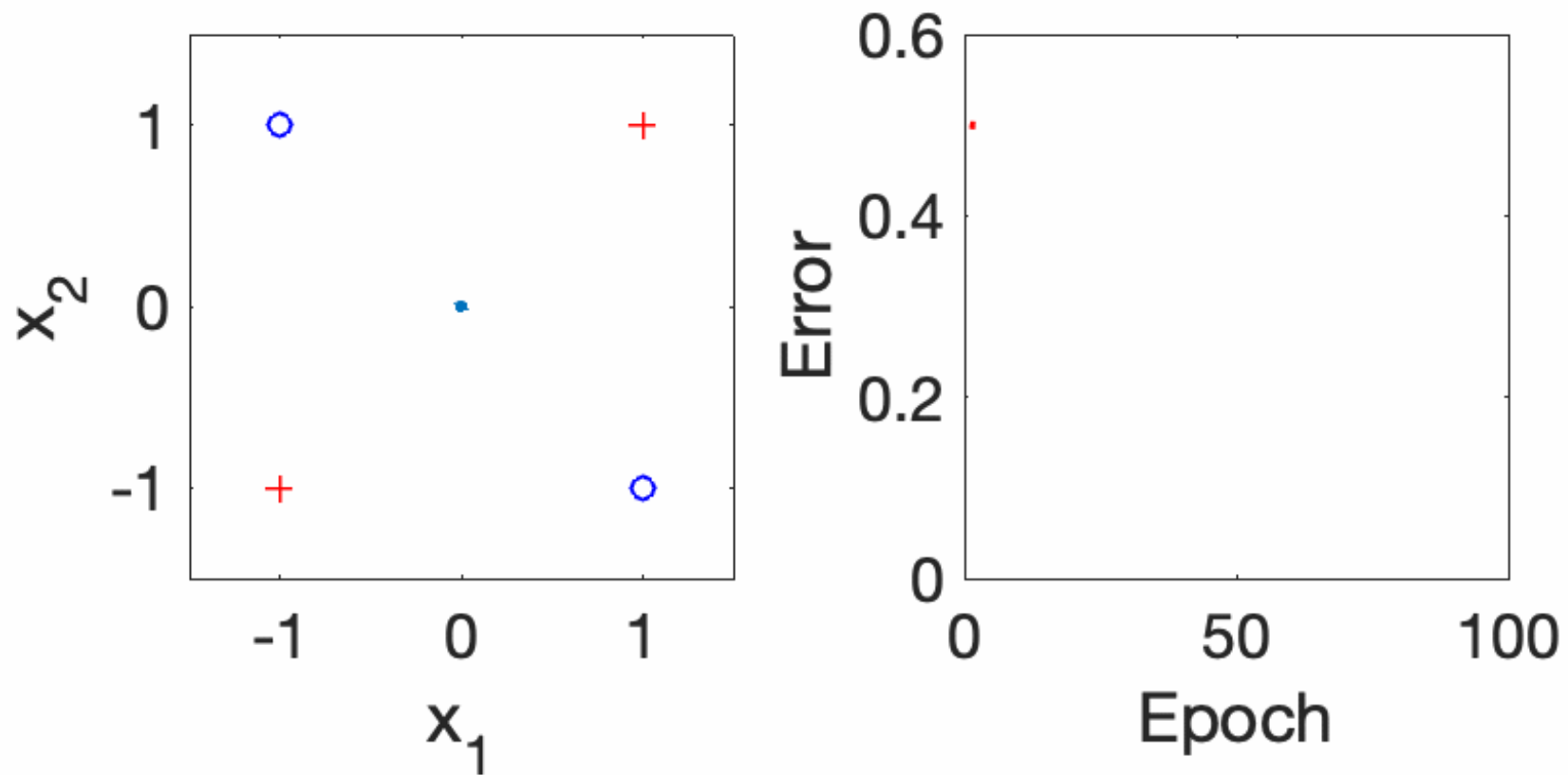


$$g_i = \begin{cases} 1 & \text{on example } i \\ 0 & \text{otherwise} \end{cases}$$

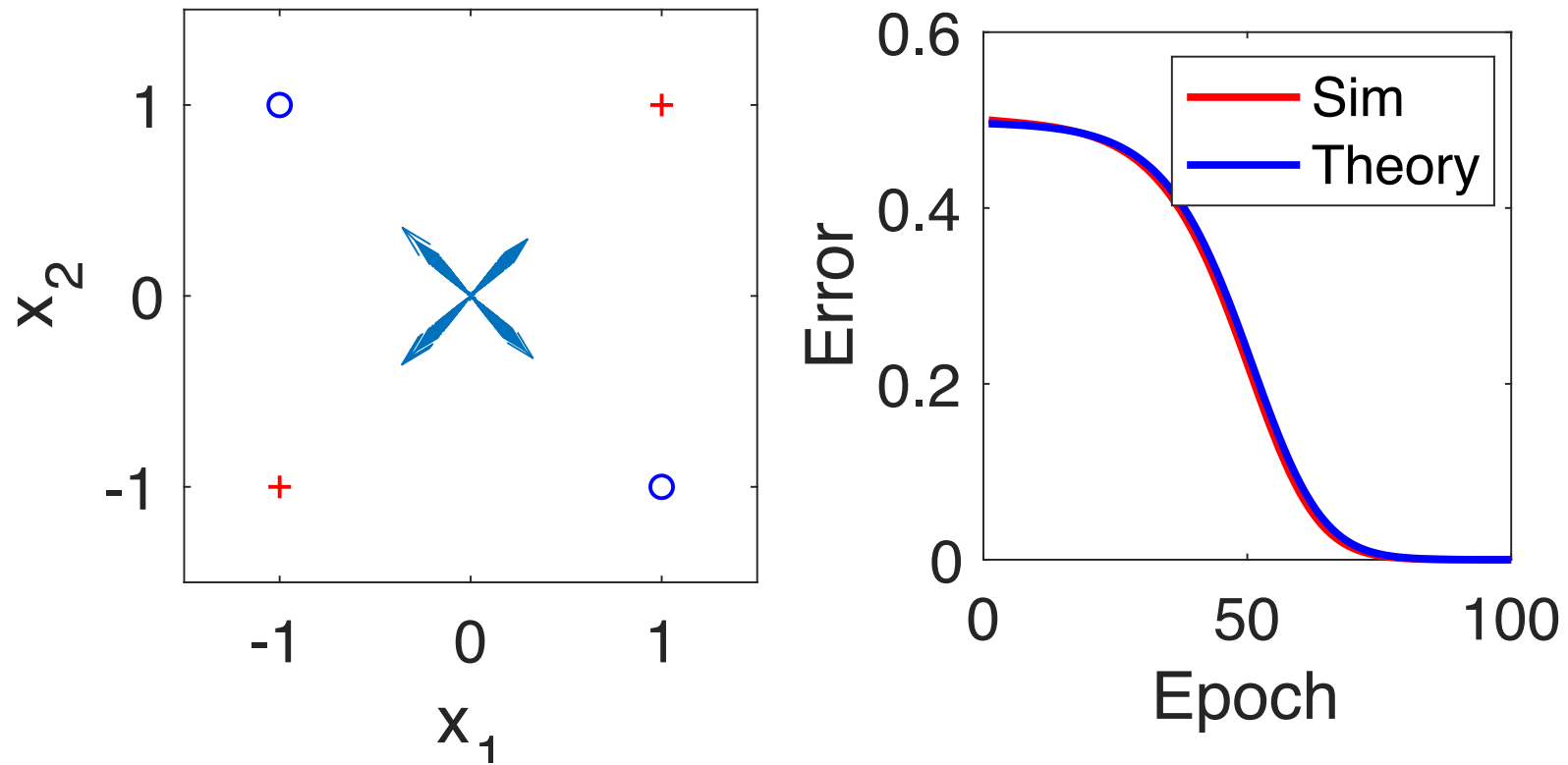
XoR Dynamics



XoR Dynamics



XoR Dynamics



Reduction and (occasionally) exact solutions

“Decoupled” initialization:

$$W_e(t) = R_{t(e)} B_e(t) R_{s(e)}^T \quad \forall e$$

Mutually diagonalizable correlations:

$$\begin{aligned} \Sigma^{yx}(p) &= U_{t(p)} S(p) V_{s(p)}^T \\ \Sigma^x(j, p) &= V_{s(j)} D(j, p) V_{s(p)}^T \end{aligned}$$

Reduction:

$$\tau \frac{d}{dt} B_e = \sum_{p \in \mathcal{P}(e)} B_{p \setminus e} \left[S(p) - \sum_{j \in \mathcal{T}(t(p))} B_j D(j, p) \right]$$

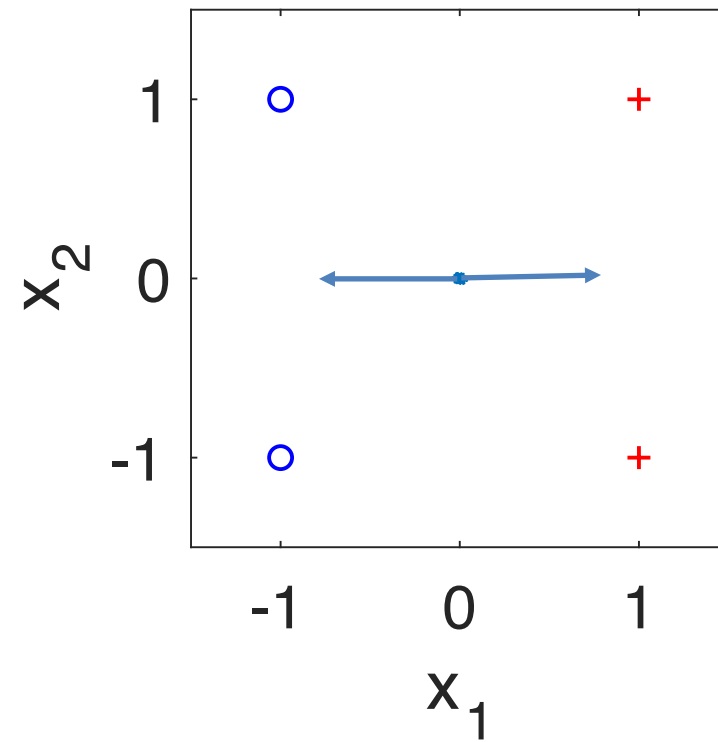
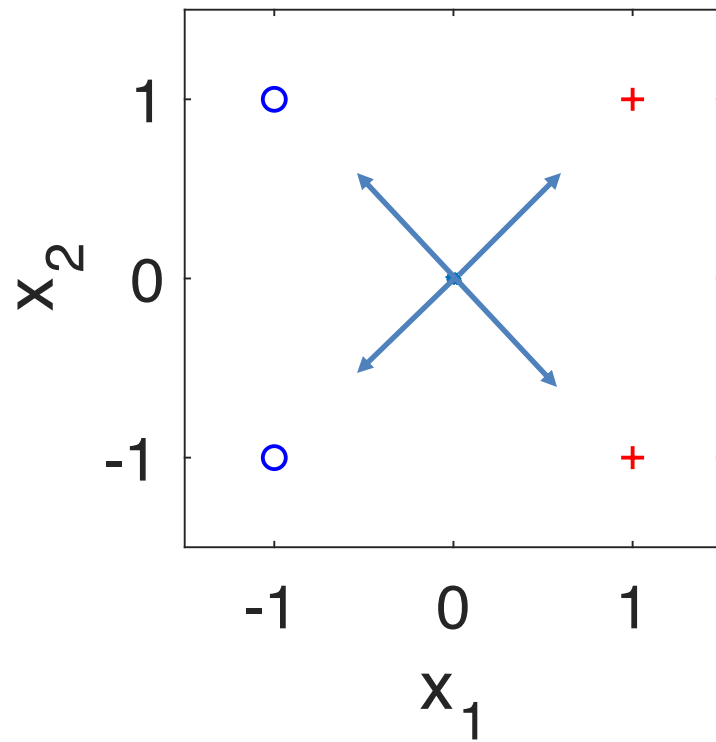
Assumptions & caveats

- Reduction exact for GDLNs
- Also exact for ReLU networks under the assumptions:
 - Gates on each example match the activity set
 - No neurons switch their activity set
 - Initial weights are decoupled
- Can approximate ReLU networks with small random weights, but not always

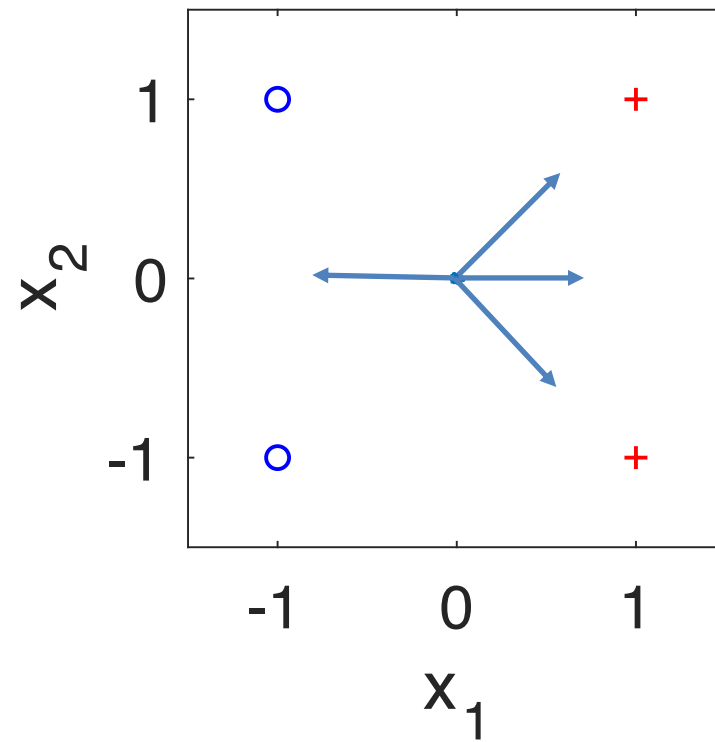
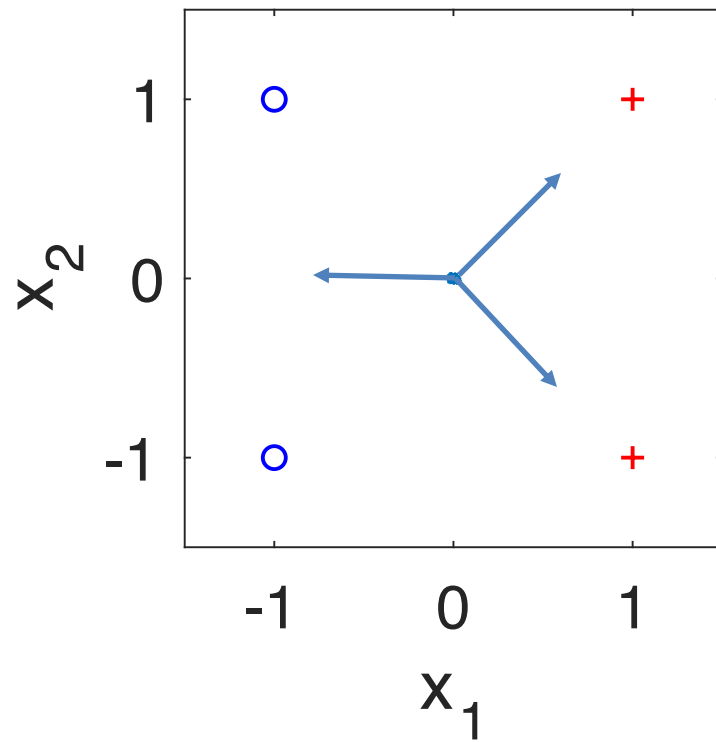
The neural race reduction

- In a large network with many pathways, these compete to reduce the global error
- A pathway's learning speed depends on:
 - Effective dataset (larger input-output correlation faster)
 - Pathway depth (deeper generally slower)
 - Initialization (larger/imbalanced generally faster)
 - Edge sharing (more pathways through edge generally faster)
- The fastest pathways can dominate the solution

Which gating structures?

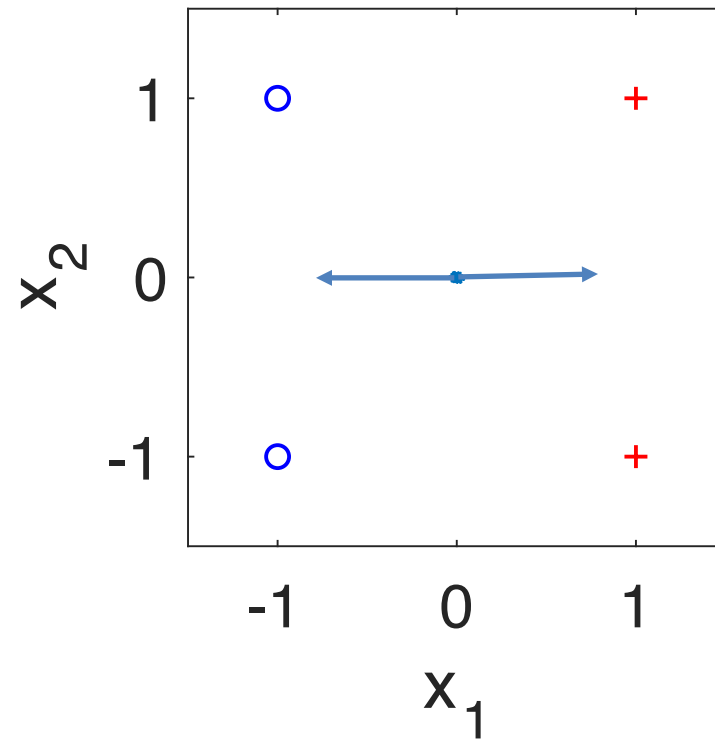
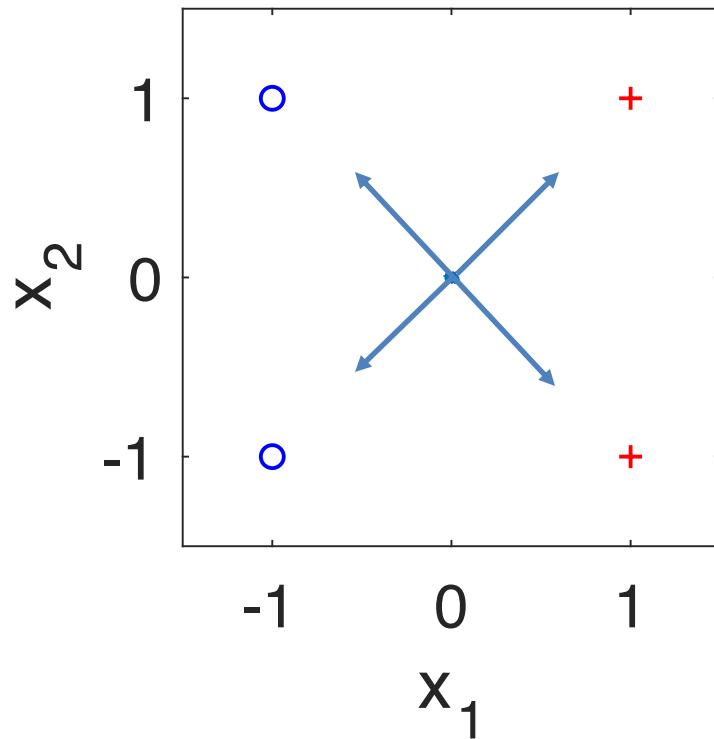


Which gating structures?

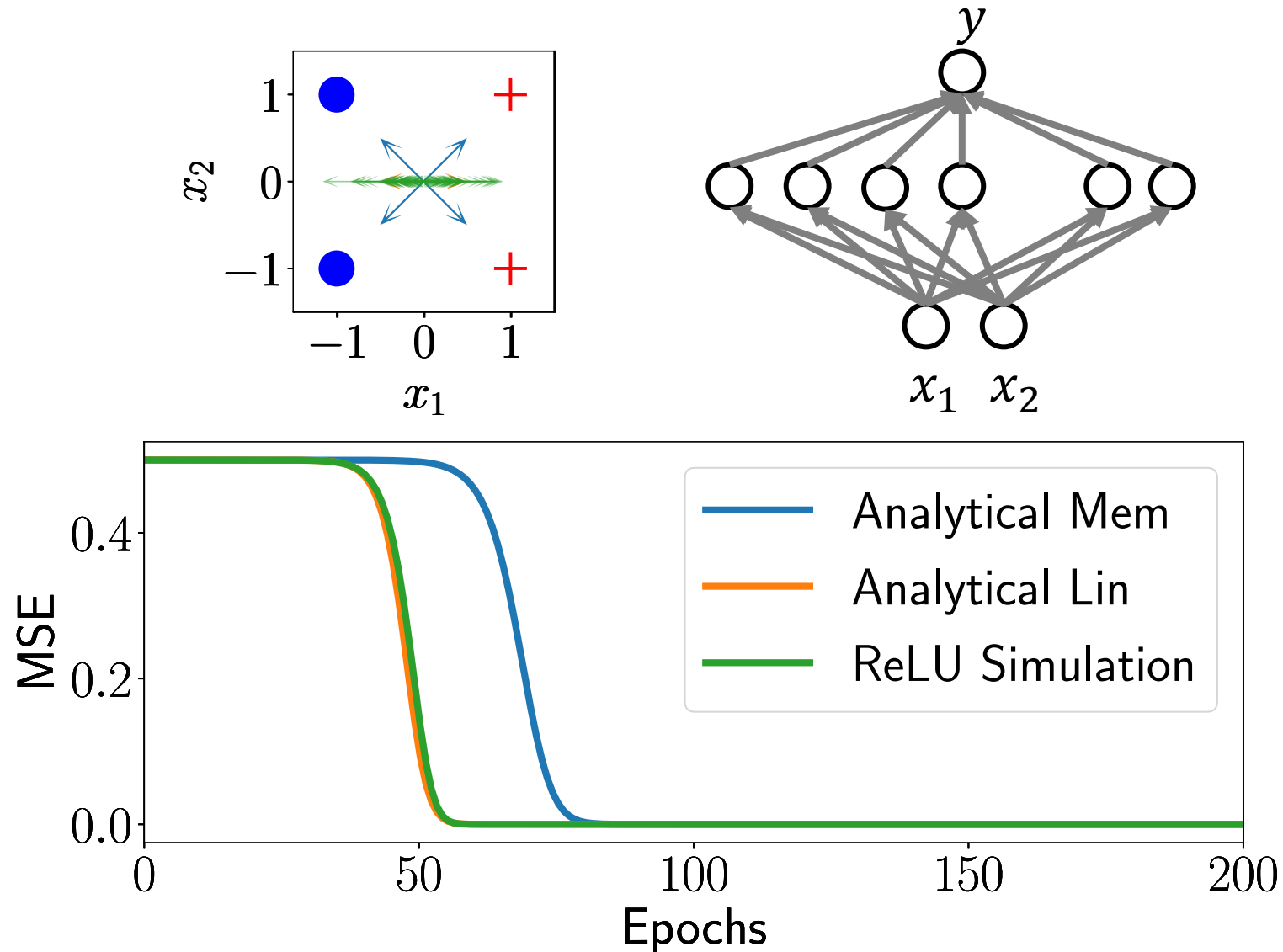


Neural Race Reduction

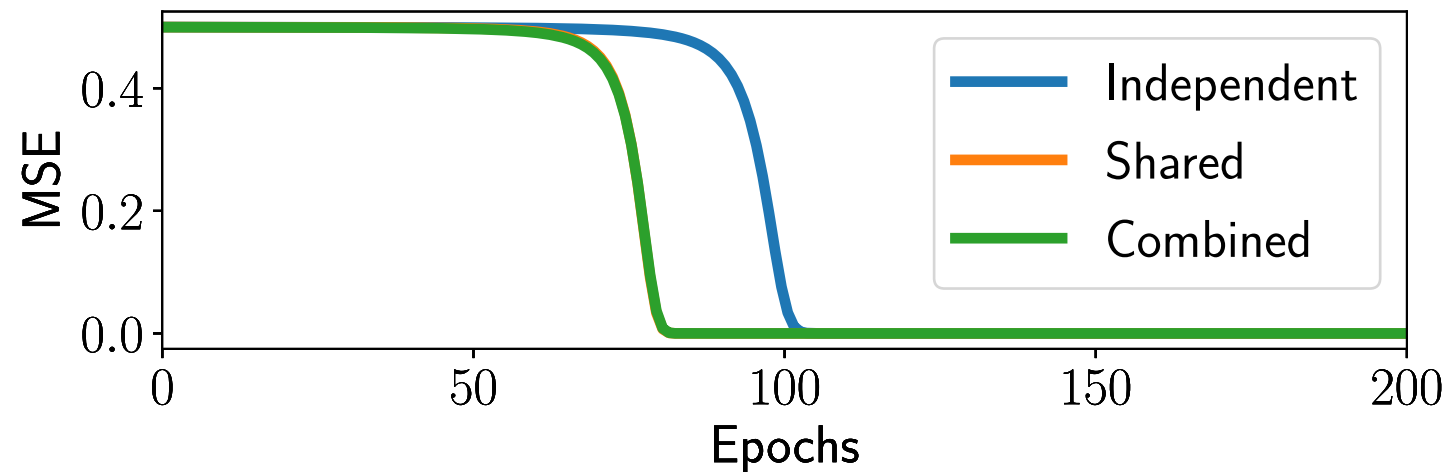
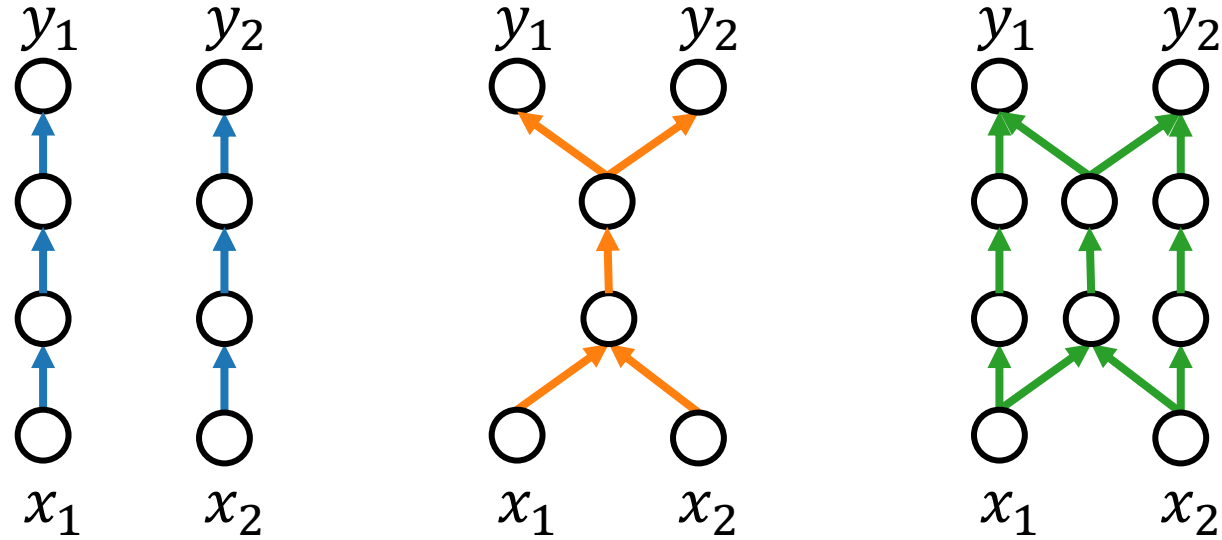
- Each gating scheme yields a distinct effective dataset and deep linear network trajectory
- The one which learns fastest dominates the solution



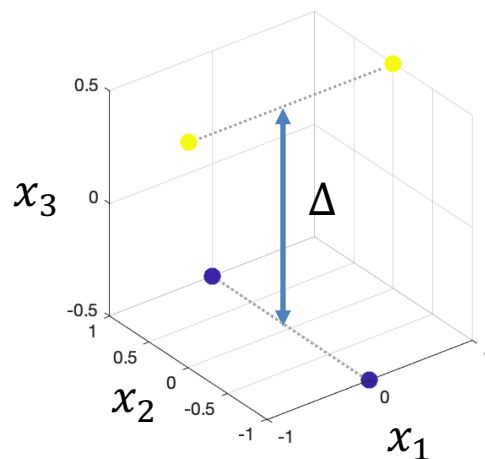
The neural race: stronger input-output correlations



The neural race: edge sharing



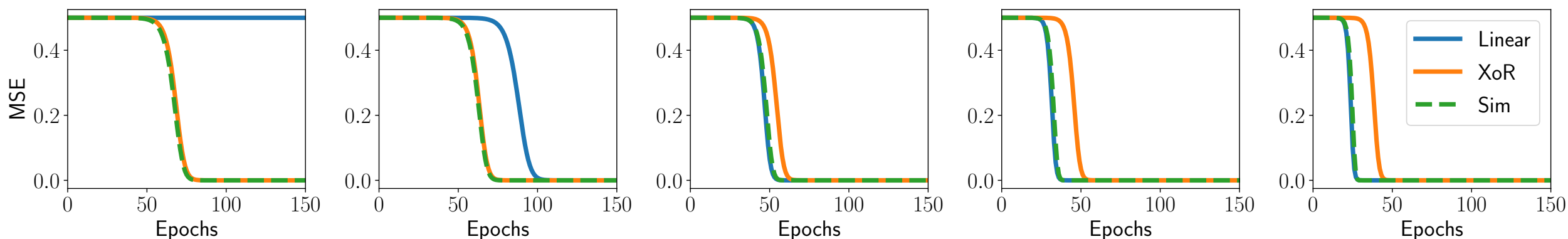
Example: transition to nonlinearity



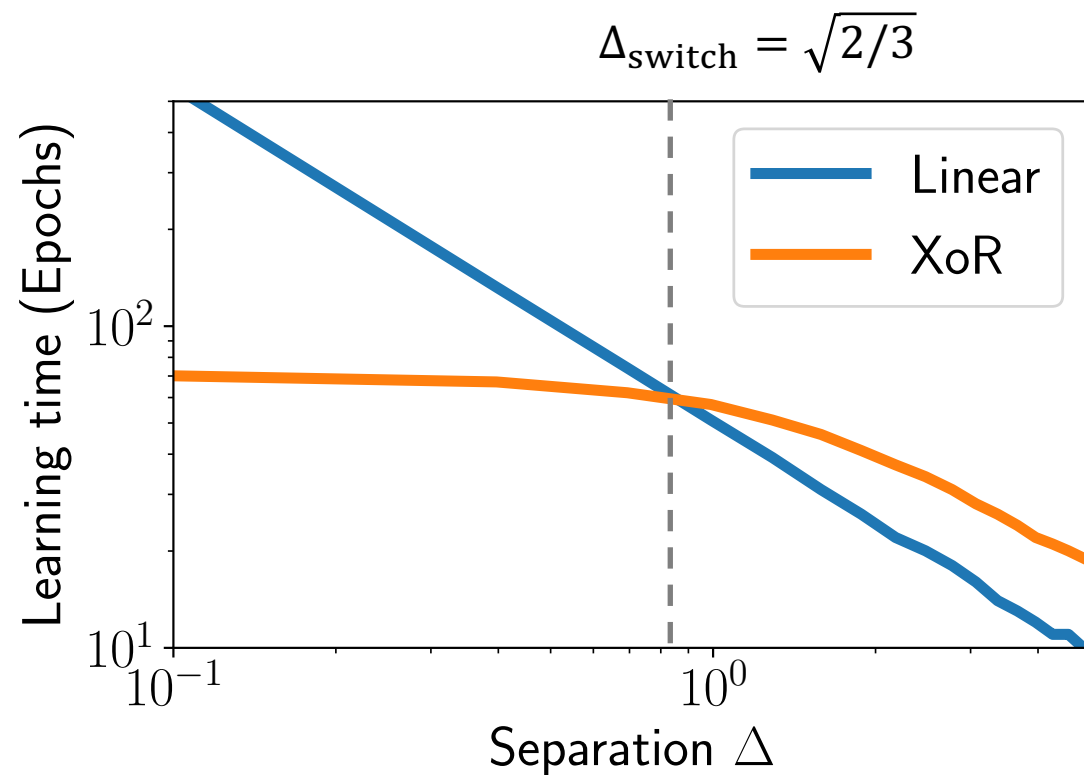
linearly separable with margin Δ for any $\Delta > 0$,
collapses to XoR at $\Delta = 0$

$\Delta = 0$

$\Delta = 5$



Example: transition to nonlinearity

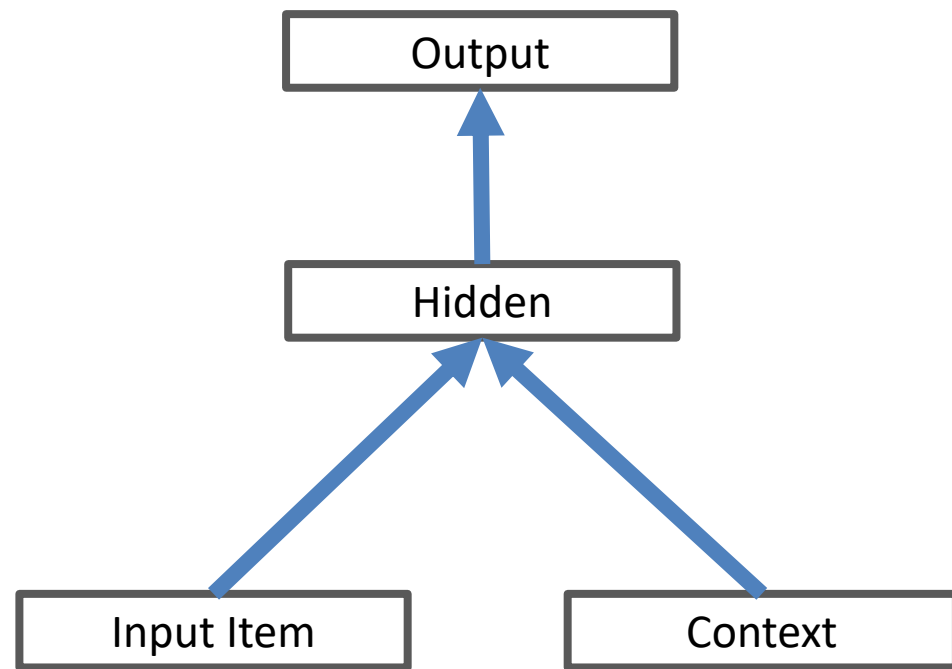
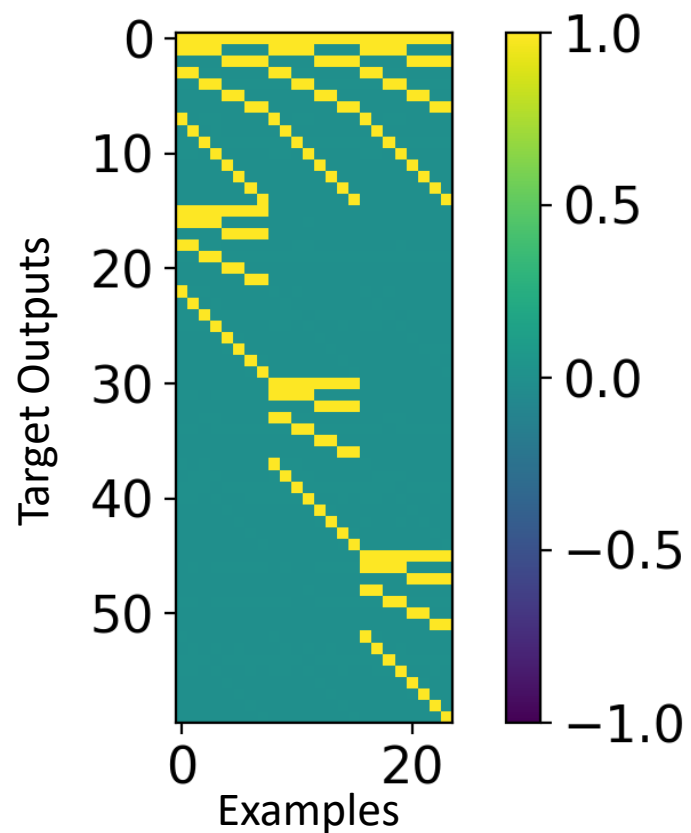


$$s_{lin} = \Delta/2$$

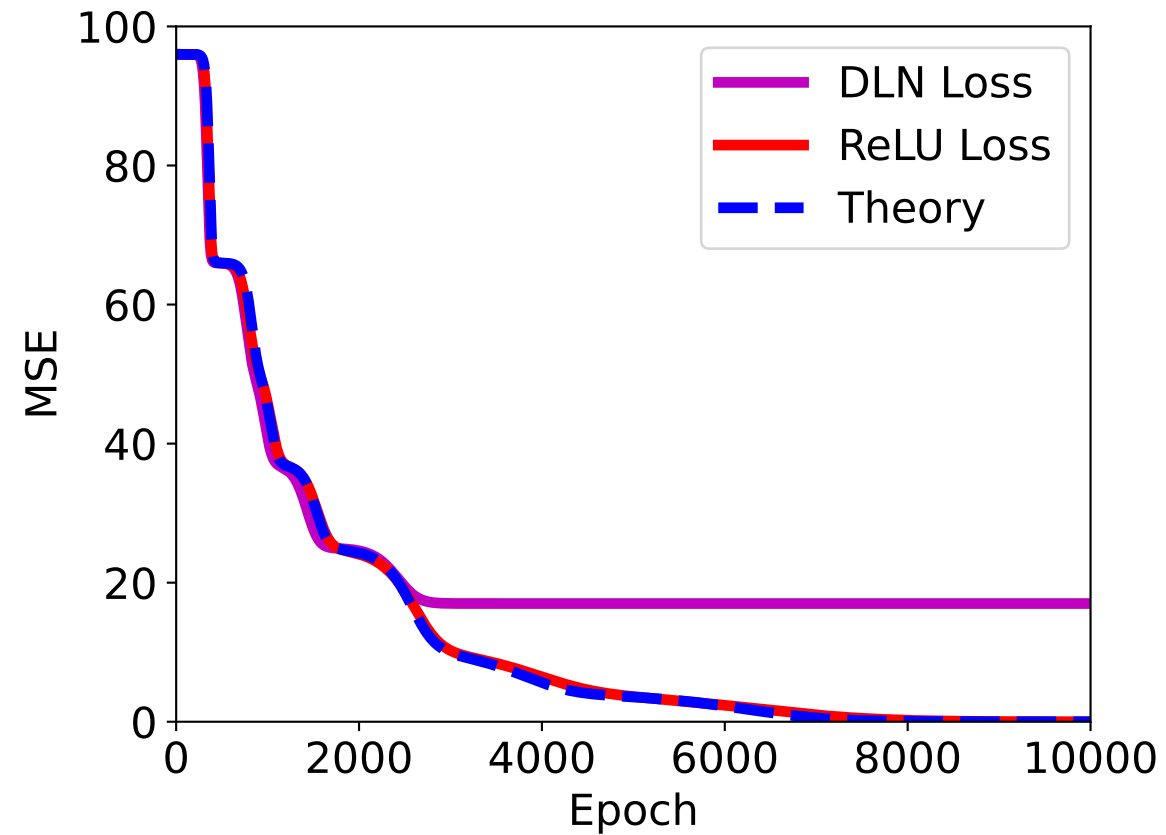
$$s_{XoR} = \frac{\sqrt{2 + \Delta^2}}{P}$$

Nonlinear representations emerge before they are strictly necessary

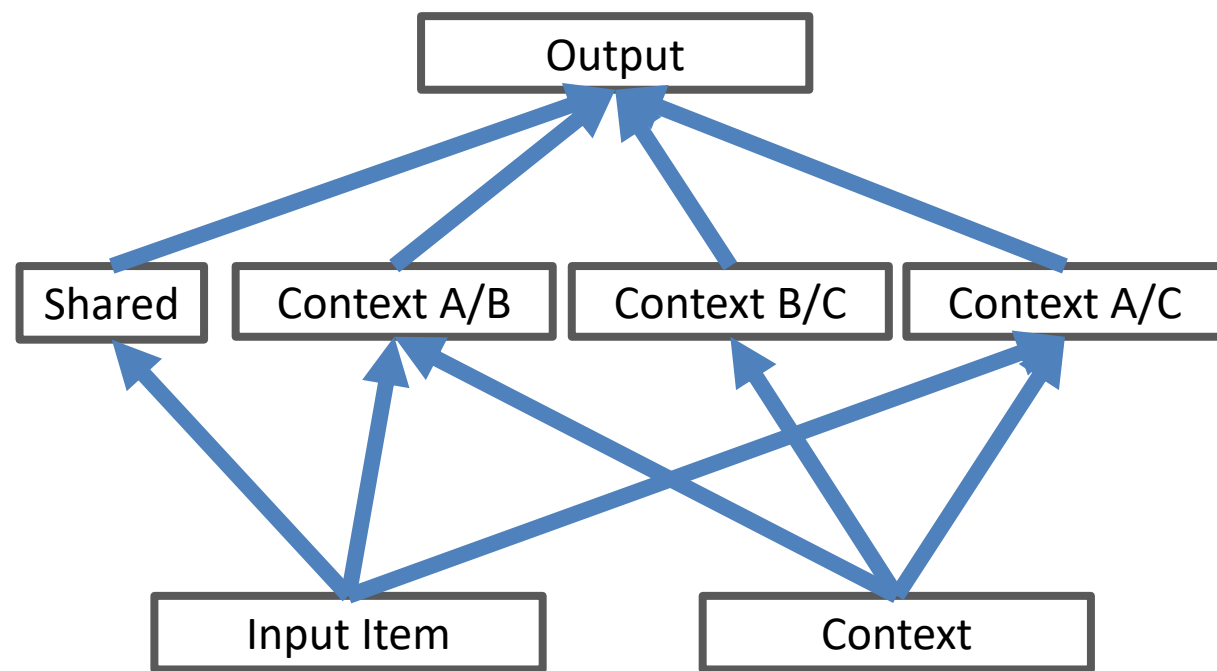
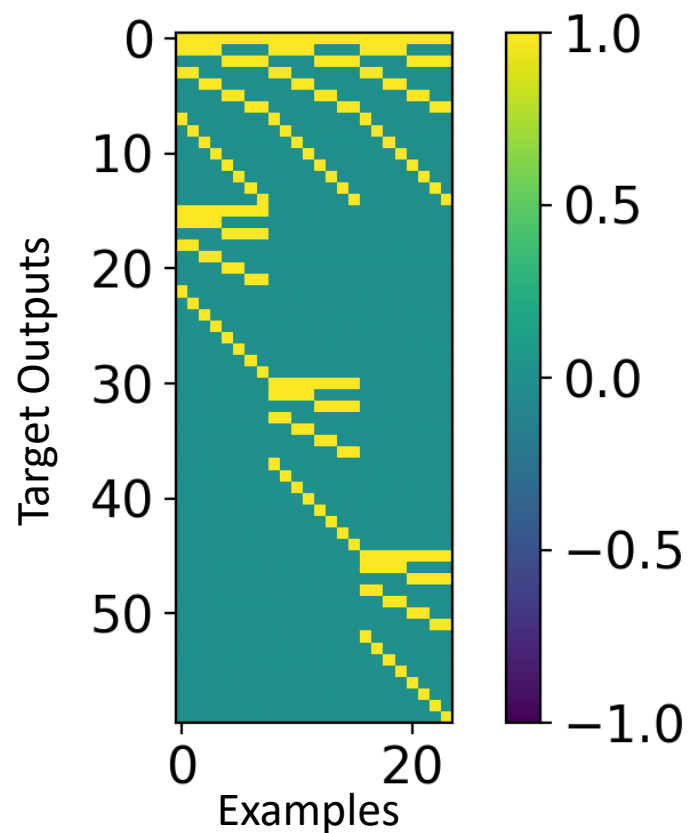
Context-dependent Processing



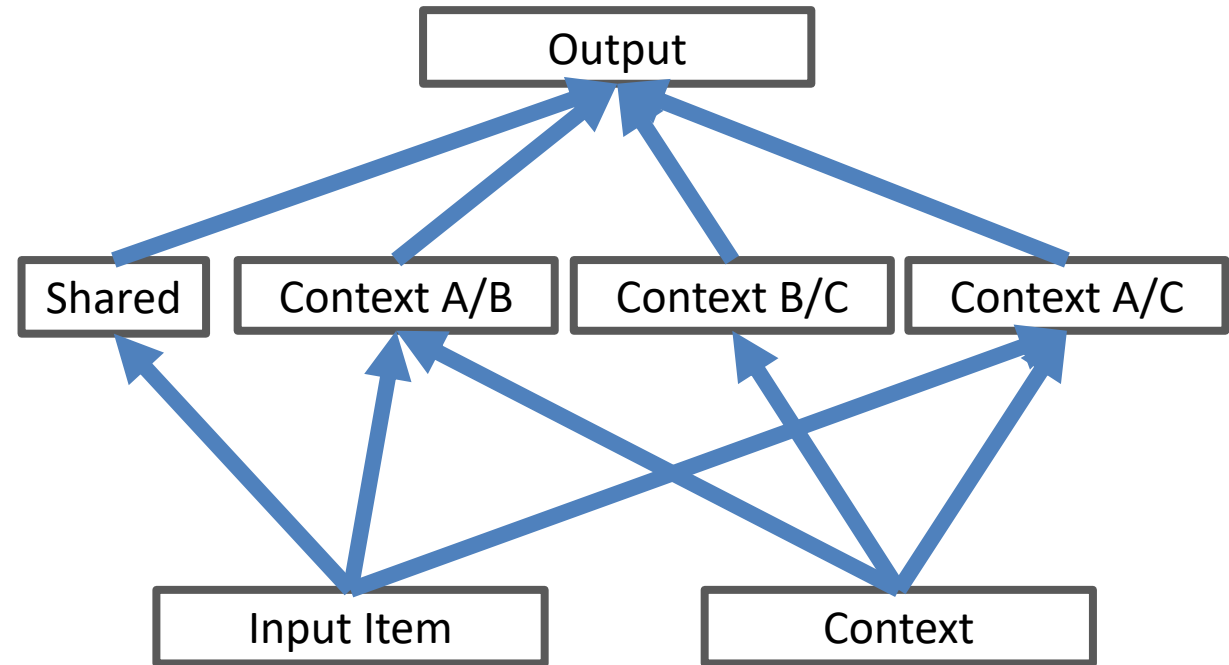
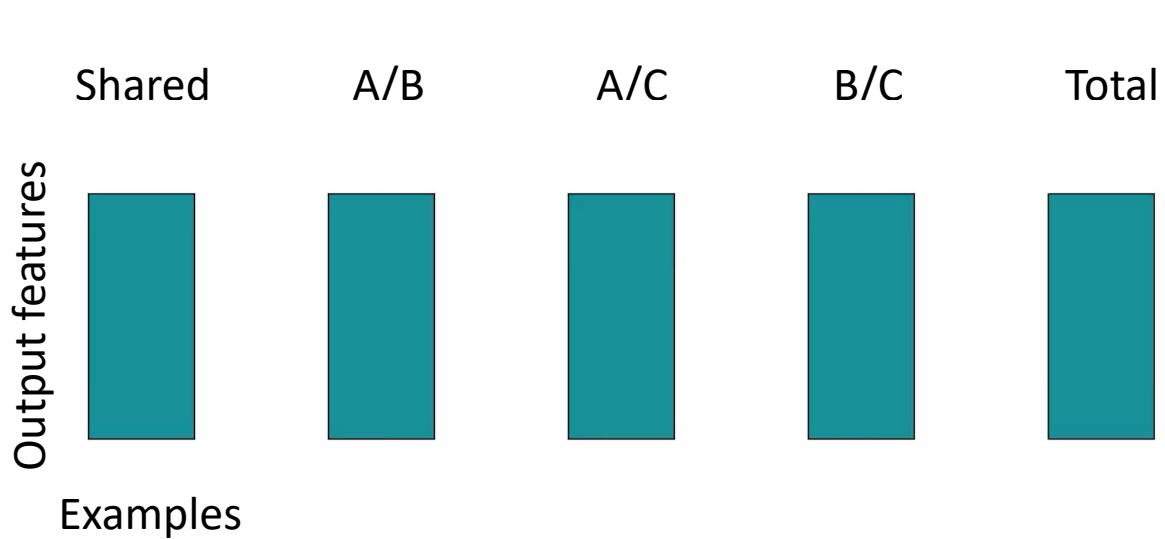
Context-dependent Processing



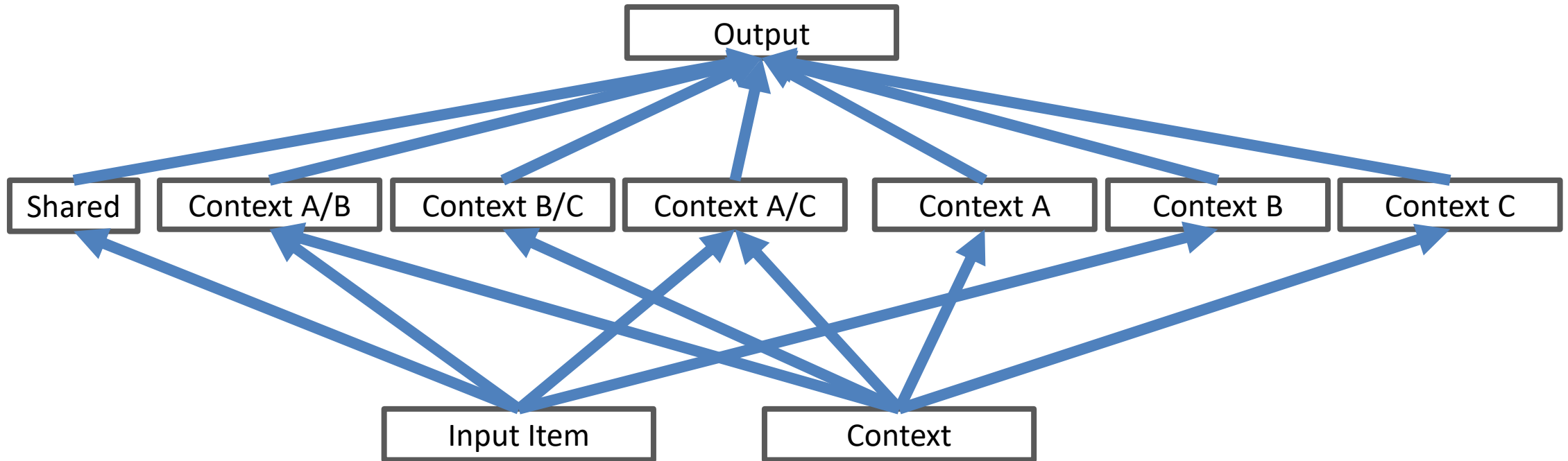
Context-dependent Processing



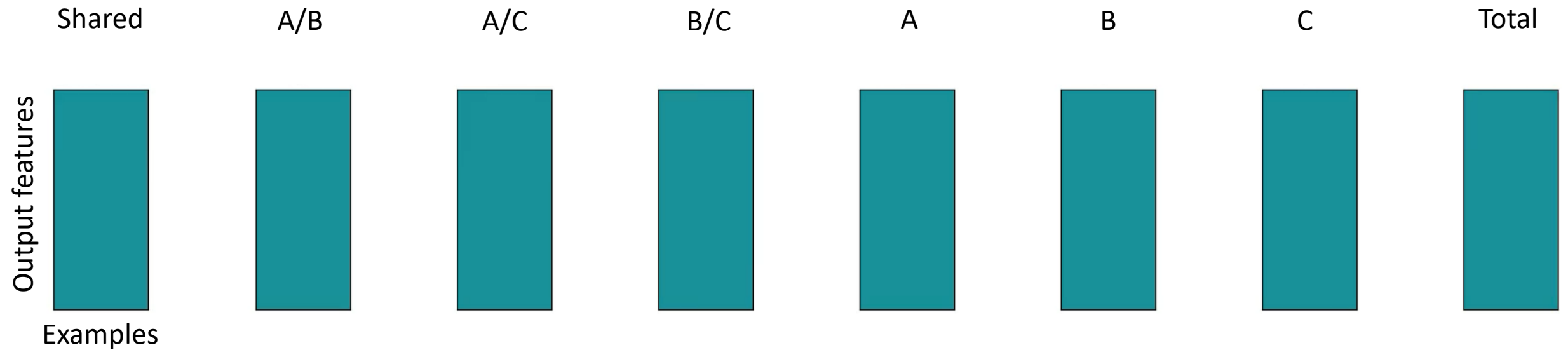
Context-dependent Processing



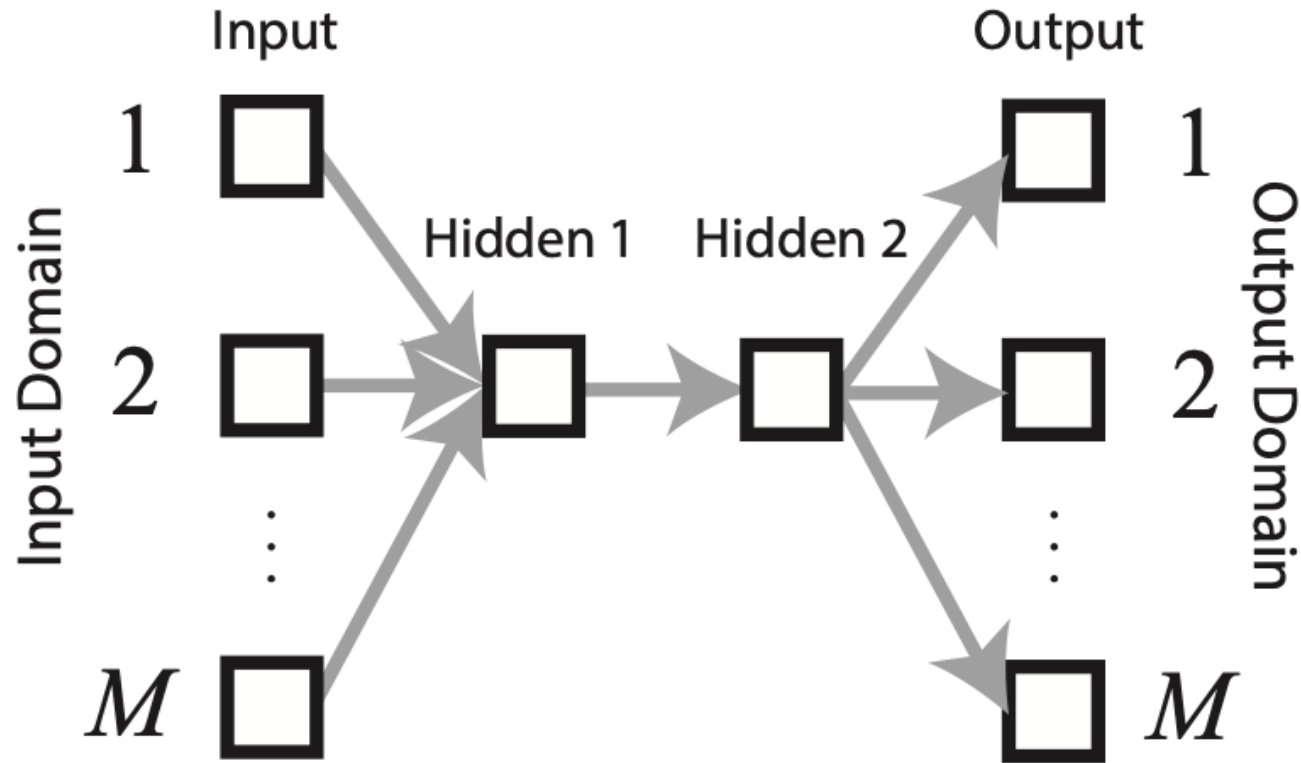
Context-dependent processing



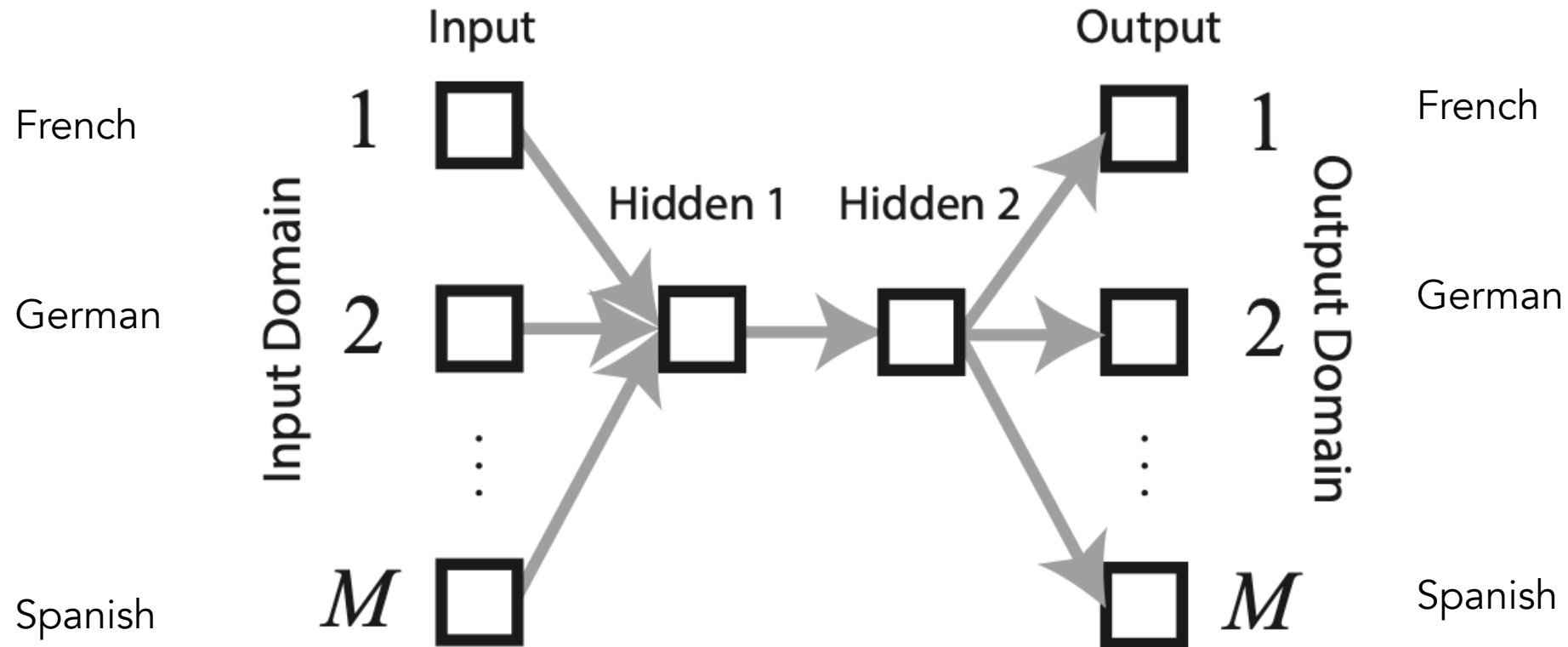
Context-dependent processing



Example: Routing network



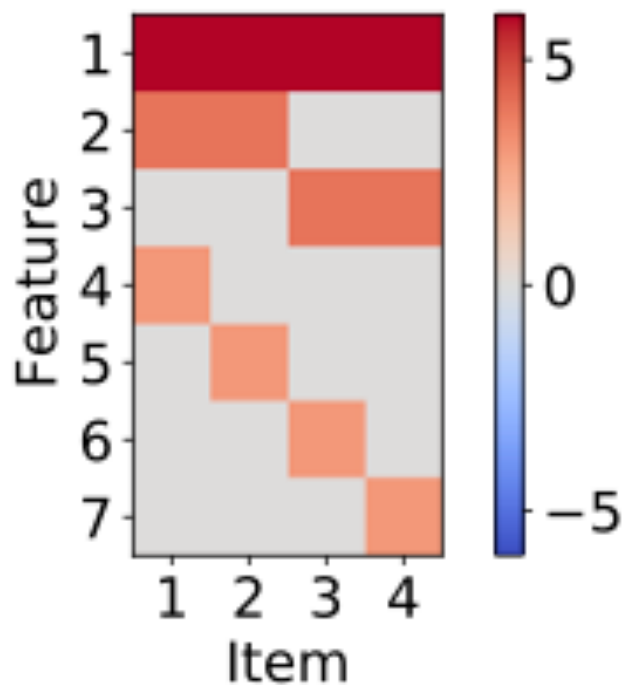
Ex: multilingual translation



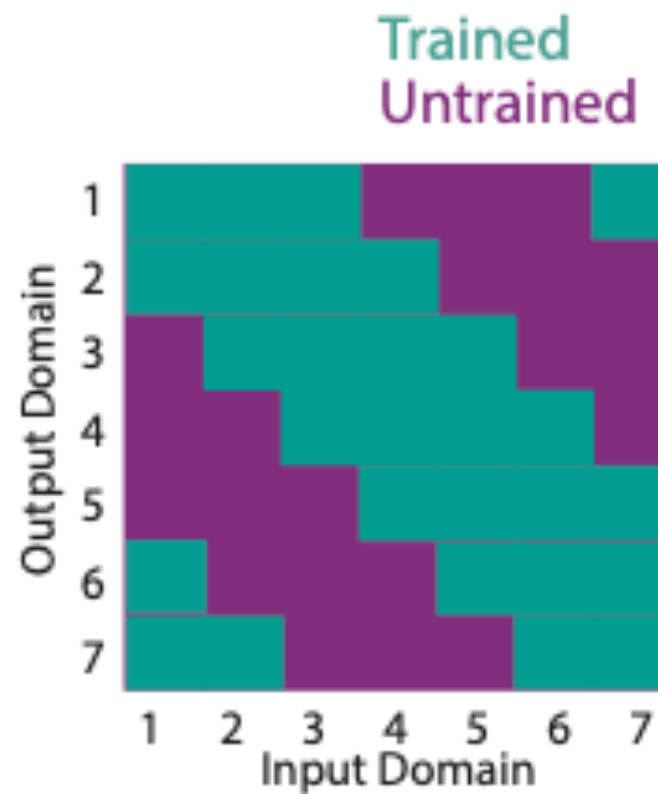
Each domain has distinctive inputs/outputs but similar underlying structural form

Dataset

Simple hierarchical dataset for each domain



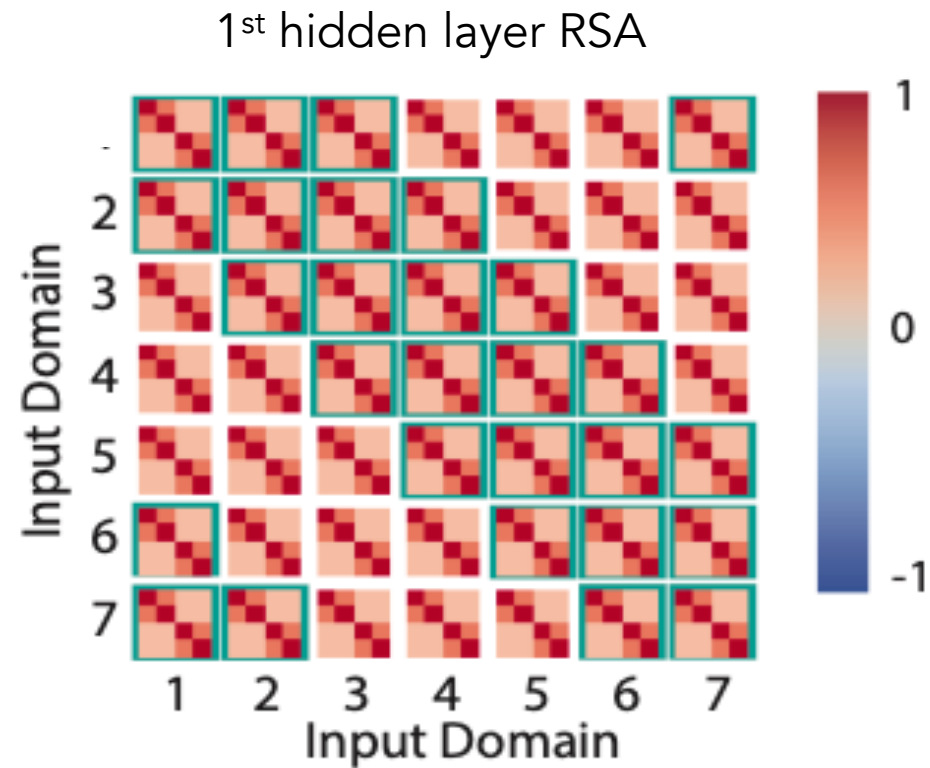
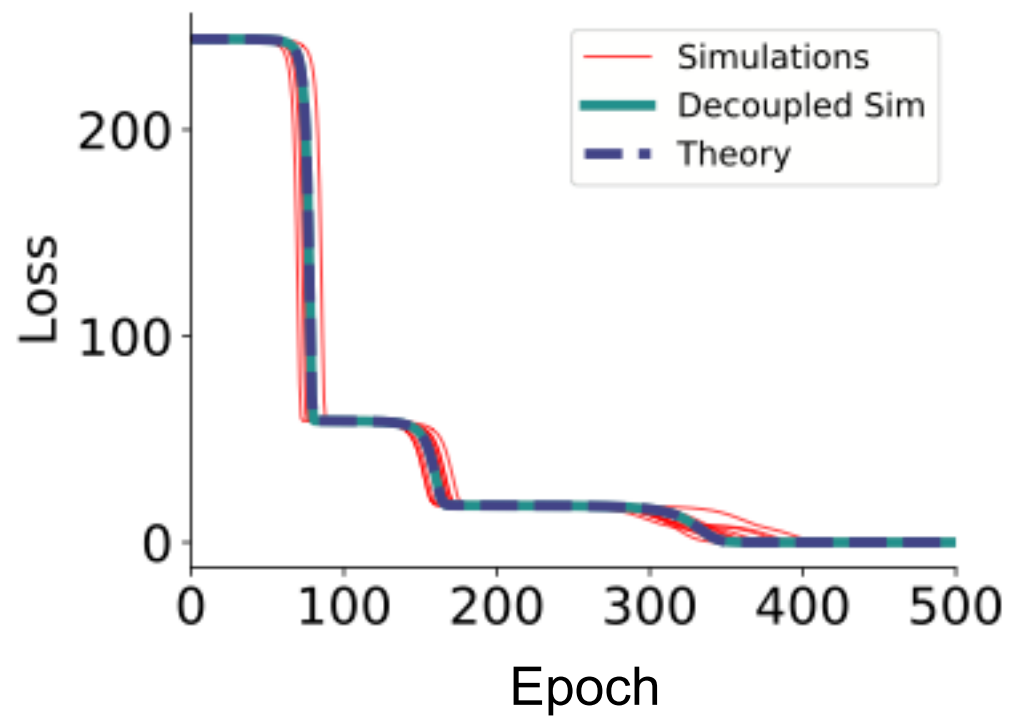
Subset of trained domain pairs



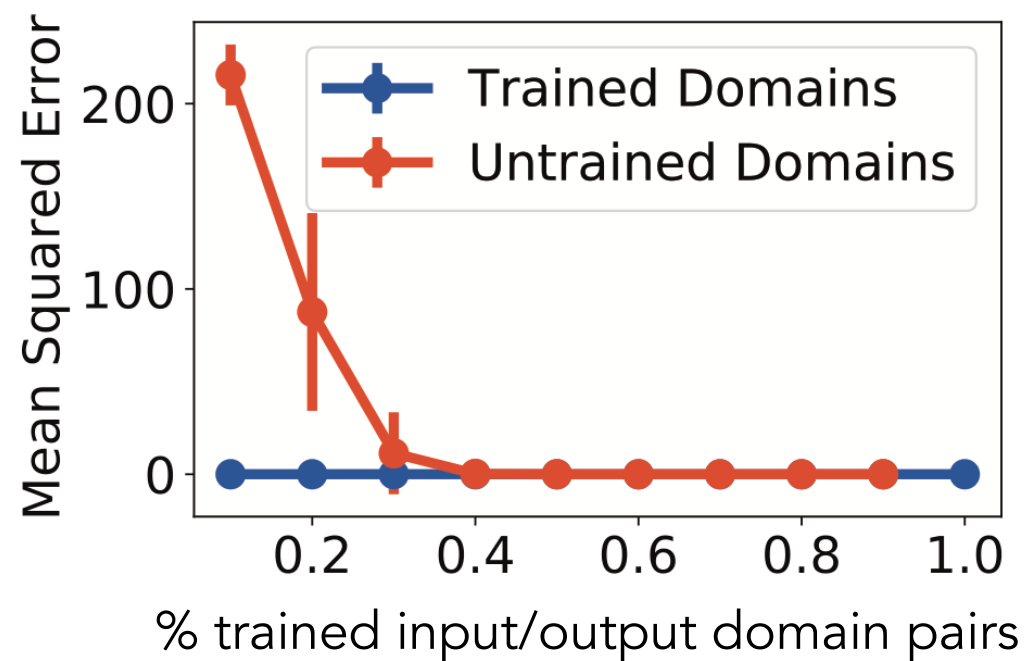
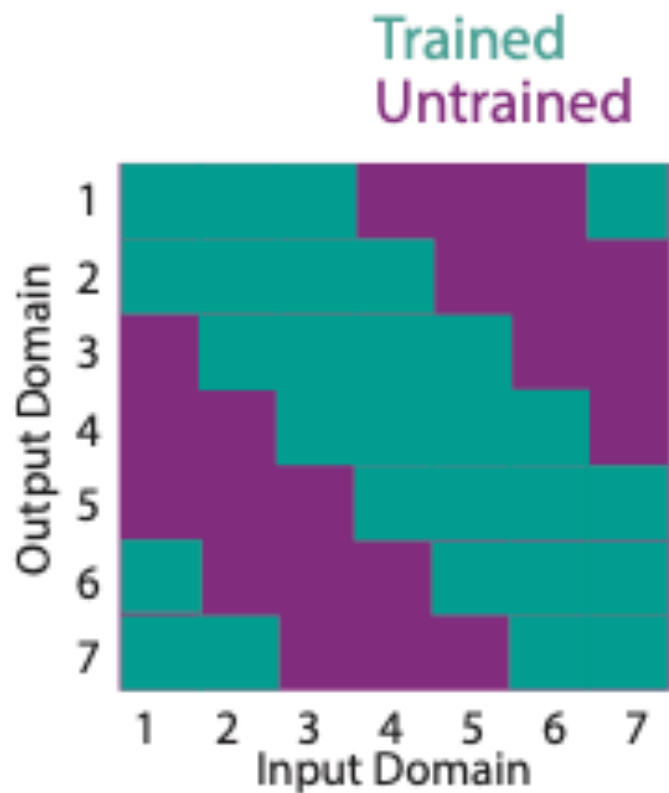
M : # domains

K : # trained output domains per input domain

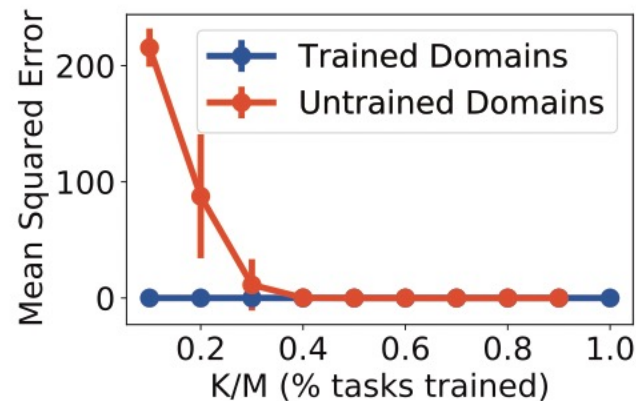
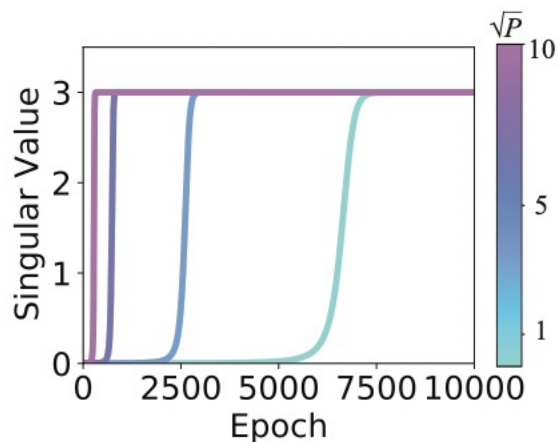
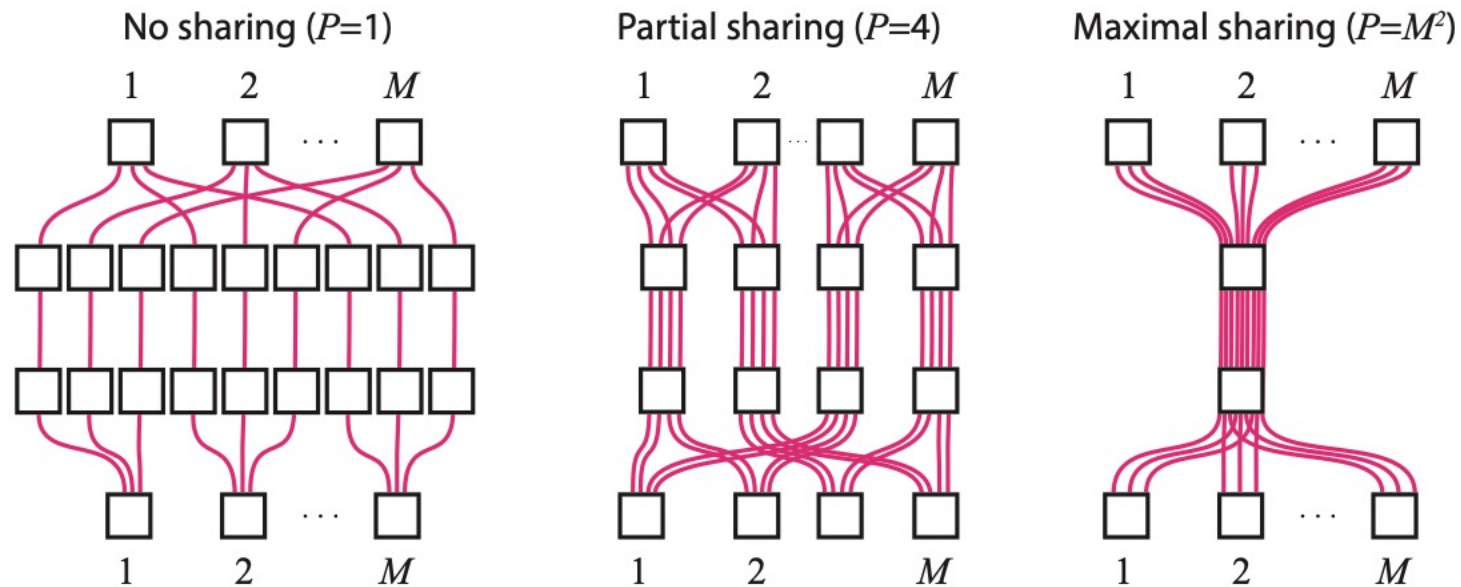
Dynamics of abstraction



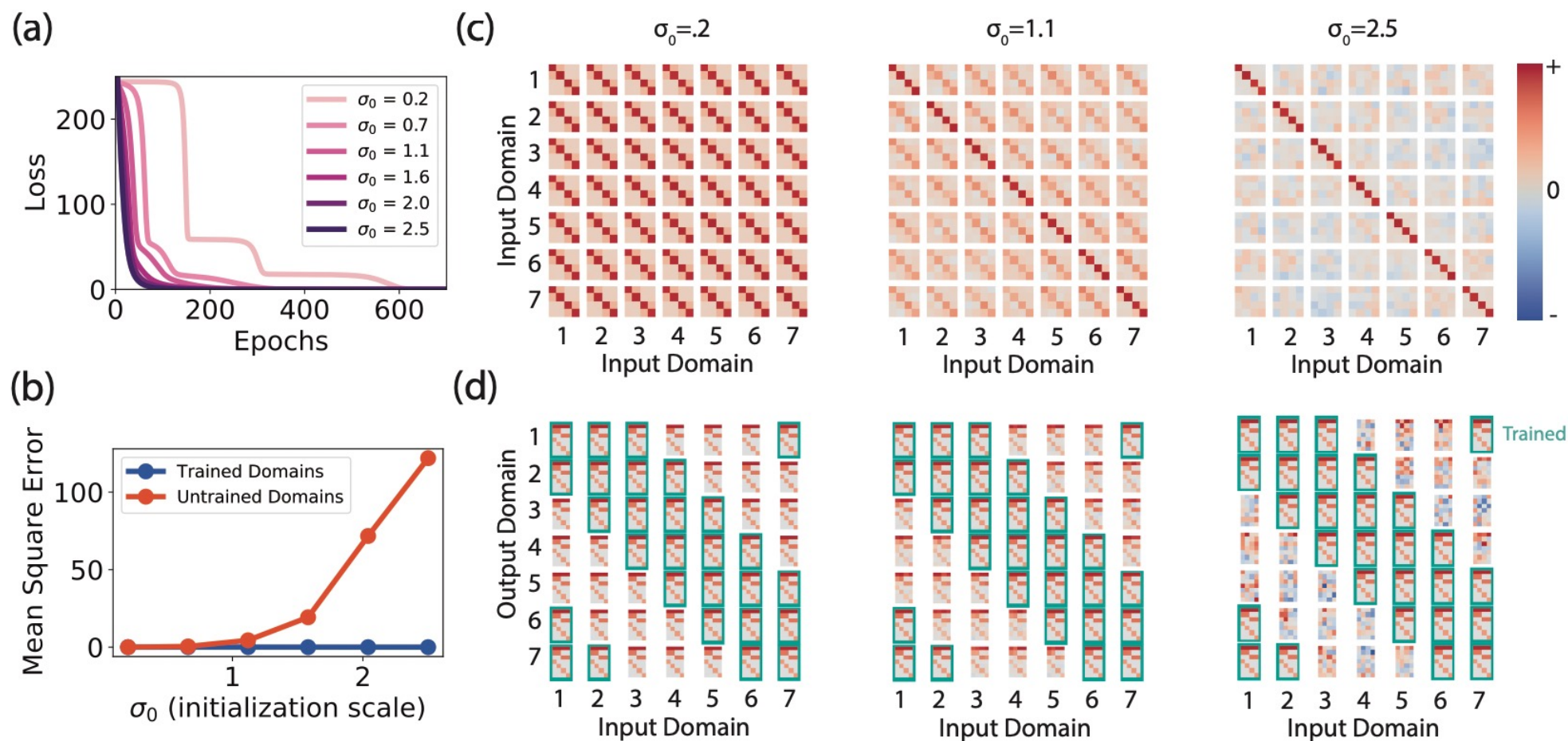
Systematic generalization



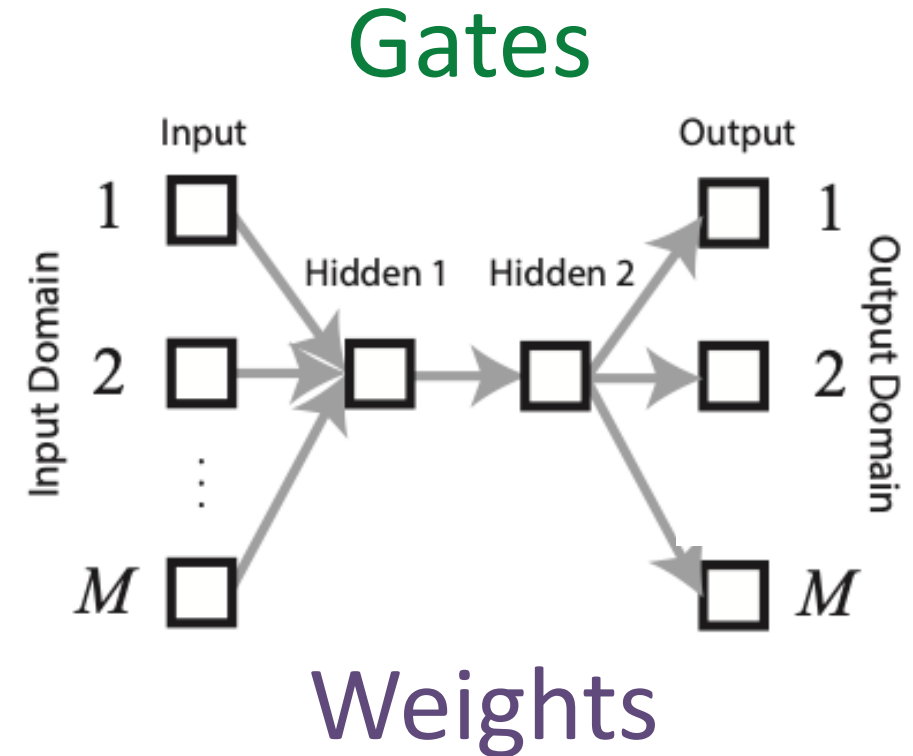
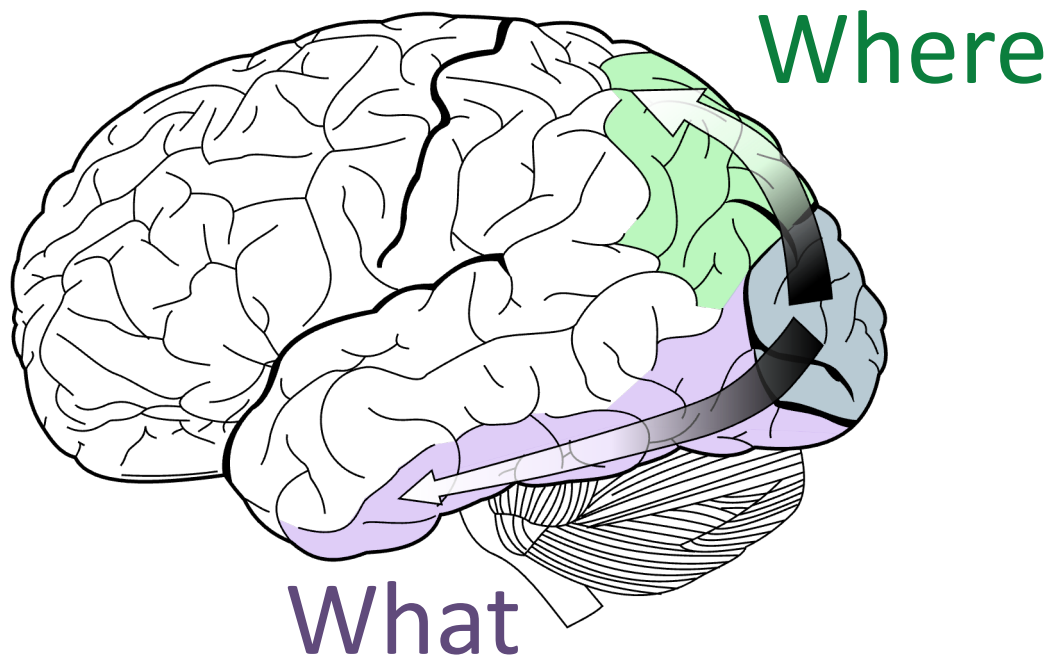
Race dynamics favor shared structure



Initialization dependence: rich vs lazy learning



Factorization & principle of convergence



Selket, <https://commons.wikimedia.org/w/index.php?curid=1679336>

Multipotential representation learning

- Animals can recombine their existing knowledge to exploit new opportunities
- In machine learning systems, this ability can emerge at scale (e.g., in context learning)
- What are the factors that give rise to *multipotential* representations?

Conclusion & outlook

- Depth introduces a hierarchy of saddle points into the loss landscape, yielding a quasi-systematic progression through stages
- Initialization determines whether these saddle points influence dynamics, yielding several learning regimes
- In nonlinear networks, pathways race to explain the dataset

References

- Saxe, A. M., McClelland, J. L., & Ganguli, S. (2019). A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23), 11537–11546.
- Liebana Garcia et al. (2025). Dopamine encodes deep network teaching signals for individual learning trajectories. *Cell*.
- Braun, L., Dominé, C.C.J., Fitzgerald, J., Saxe, A.M. (2022) Exact learning dynamics of deep linear networks with prior knowledge. In *NeurIPS*.
- Dominé, C.C.J., Anguita, N., Proca, A. M., Braun, L., Kunin, D., Mediano, P. A. M., & Saxe, A. M. (2025). From lazy to rich: Exact learning dynamics in deep linear networks. In *ICLR*.
- Saxe*, A. M., Sodhani*, S., & Lewallen, S. (2022). The Neural Race Reduction: Dynamics of Abstraction in Gated Networks. In *ICML*. *Equal contribution.
- Jarvis, D., Klein, R., Rosman, B., Saxe, A.M. (2025). Make haste slowly: A theory of emergent structured mixed selectivity in feature learning ReLU networks. In *ICLR*.

Aaditya Singh
Anika Lowe
Basile Confavreux
Clementine Domine
Cris Holobetz
Devon Jarvis
Erin Grant
Jin Lee
Jirko Rubruck
Lukas Braun
Nishil Patel
Rodrigo Carrasco Davis
Rachel Swanson
Sam Lewallen
Sam Liebana
Sarah Armstrong
Sebastian Lee
Stefano Sarao Mannelli
Tyler Boyd-Meredith
Verena Klar
Victor Pedrosa



SCHMIDT
FUTURES



THE
ROYAL
SOCIETY

CIFAR

