



# Entropy and mutual information in models of deep neural networks

NeurIPS 2018 - Thursday Dec 06th - Spotlight  
Poster @ Room 210 & 230 AB #110

Marylou Gabrié (LPS ENS), Andre Manoel (INRIA Saclay, Owkin),  
Clément Luneau, Jean Barbier, Nicolas Macris (EPFL),  
Lenka Zdeborová (CEA Saclay), Florent Krzakala (LPS ENS)

# Motivations

## Information theoretic arguments to deep learning

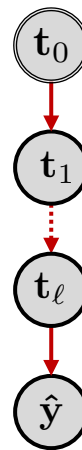
- **Theory of generalization:** Shwartz-Ziv et al. 2017, Saxe et al. 2018, Goldfeld et al. 2018 etc.
- **New regularizer:** Chalk et al. 2016, Alemi et al. 2017, Kolchinsky et al. 2017, Belghazi et al. 2017, Zhao et al. 2018, Achille et al. 2018, Hjelm et al. 2018, etc.

# Motivations

## Information theoretic arguments to deep learning

- **Theory of generalization:** Shwartz-Ziv et al. 2017, Saxe et al. 2018, Goldfeld et al. 2018 etc.
- **New regularizer:** Chalk et al. 2016, Alemi et al. 2017, Kolchinsky et al. 2017, Belghazi et al. 2017, Zhao et al. 2018, Achille et al. 2018, Hjelm et al. 2018, etc.

## Computing mutual informations in neural networks ...

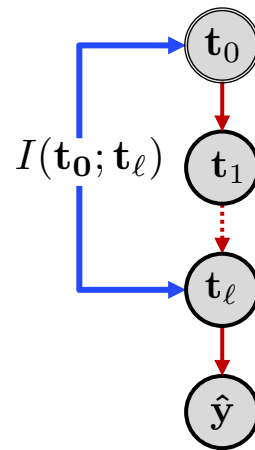


# Motivations

## Information theoretic arguments to deep learning

- **Theory of generalization:** Shwartz-Ziv et al. 2017, Saxe et al. 2018, Goldfeld et al. 2018 etc.
- **New regularizer:** Chalk et al. 2016, Alemi et al. 2017, Kolchinsky et al. 2017, Belghazi et al. 2017, Zhao et al. 2018, Achille et al. 2018, Hjelm et al. 2018, etc.

## Computing mutual informations in neural networks ...

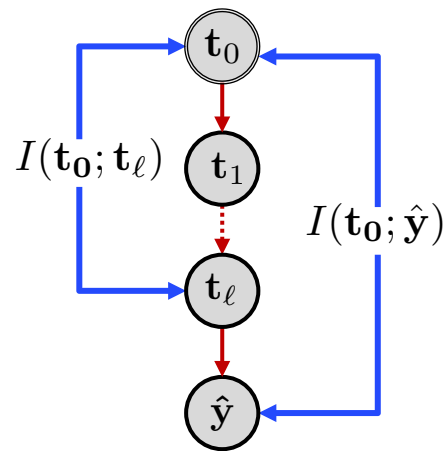


# Motivations

## Information theoretic arguments to deep learning

- **Theory of generalization:** Shwartz-Ziv et al. 2017, Saxe et al. 2018, Goldfeld et al. 2018 etc.
- **New regularizer:** Chalk et al. 2016, Alemi et al. 2017, Kolchinsky et al. 2017, Belghazi et al. 2017, Zhao et al. 2018, Achille et al. 2018, Hjelm et al. 2018, etc.

## Computing mutual informations in neural networks ...



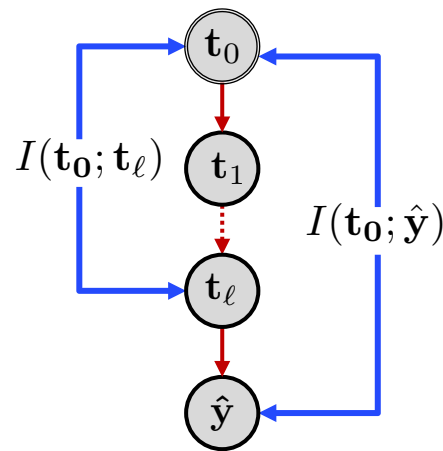
# Motivations

## Information theoretic arguments to deep learning

- **Theory of generalization:** Shwartz-Ziv et al. 2017, Saxe et al. 2018, Goldfeld et al. 2018 etc.
- **New regularizer:** Chalk et al. 2016, Alemi et al. 2017, Kolchinsky et al. 2017, Belghazi et al. 2017, Zhao et al. 2018, Achille et al. 2018, Hjelm et al. 2018, etc.

## Computing mutual informations in neural networks ... is tricky

- Intractable  $p(\mathbf{t}_\ell) = \int d\mathbf{t}_0 d\mathbf{t}_1 \cdots d\mathbf{t}_{\ell-1} p(\mathbf{t}_\ell | \mathbf{t}_{\ell-1}) \cdots p(\mathbf{t}_1 | \mathbf{t}_0) p(\mathbf{t}_0)$



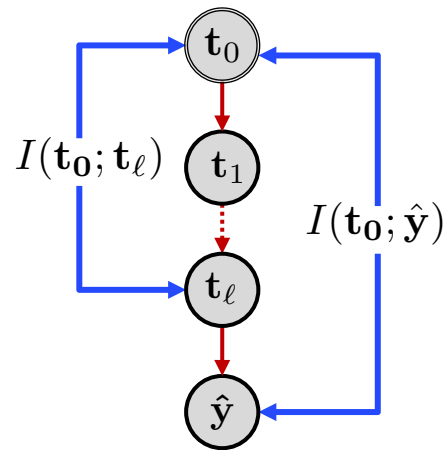
# Motivations

## Information theoretic arguments to deep learning

- **Theory of generalization:** Shwartz-Ziv et al. 2017, Saxe et al. 2018, Goldfeld et al. 2018 etc.
- **New regularizer:** Chalk et al. 2016, Alemi et al. 2017, Kolchinsky et al. 2017, Belghazi et al. 2017, Zhao et al. 2018, Achille et al. 2018, Hjelm et al. 2018, etc.

## Computing mutual informations in neural networks ... is tricky

- Intractable  $p(\mathbf{t}_\ell) = \int d\mathbf{t}_0 d\mathbf{t}_1 \cdots d\mathbf{t}_{\ell-1} p(\mathbf{t}_\ell | \mathbf{t}_{\ell-1}) \cdots p(\mathbf{t}_1 | \mathbf{t}_0) p(\mathbf{t}_0)$
- Computing  $I(\mathbf{t}_0; \mathbf{t}_\ell)$  or  $H(\mathbf{t}_\ell) = - \int d\mathbf{t}_\ell p(\mathbf{t}_\ell) \log p(\mathbf{t}_\ell)$  yet more involved



# Motivations

## Information theoretic arguments to deep learning

- Theory of generalization: Shwartz-Ziv et al. 2017, Saxe et al. 2018, Goldfeld et al. 2018 etc.
- New regularizer: Chalk et al. 2016, Alemi et al. 2017, Kolchinsky et al. 2017, Belghazi et al. 2017, Zhao et al. 2018, Achille et al. 2018, Hjelm et al. 2018, etc.

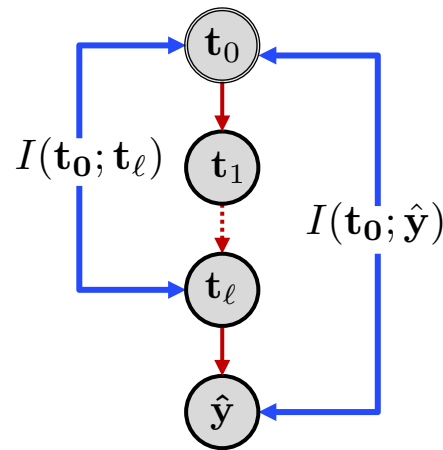
## Computing mutual informations in neural networks ... is tricky

- Intractable  $p(\mathbf{t}_\ell) = \int d\mathbf{t}_0 d\mathbf{t}_1 \cdots d\mathbf{t}_{\ell-1} p(\mathbf{t}_\ell | \mathbf{t}_{\ell-1}) \cdots p(\mathbf{t}_1 | \mathbf{t}_0) p(\mathbf{t}_0)$
- Computing  $I(\mathbf{t}_0; \mathbf{t}_\ell)$  or  $H(\mathbf{t}_\ell) = - \int d\mathbf{t}_\ell p(\mathbf{t}_\ell) \log p(\mathbf{t}_\ell)$  yet more involved

## Sampled based estimations (non-parametric, variational methods)

Kraskov et al. 2014, Kolchinsky et al. 2017, Belghazi et al. 2017, Goldfeld et al. 2018 etc.

- Less and less reliable as networks get large
- Need for benchmark controlled case





# Motivations

## Information theoretic arguments to deep learning

- Theory of generalization: Shwartz-Ziv et al. 2017, Saxe et al. 2018, Goldfeld et al. 2018 etc.
- New regularizer: Chalk et al. 2016, Alemi et al. 2017, Kolchinsky et al. 2017, Belghazi et al. 2017, Zhao et al. 2018, Achille et al. 2018, Hjelm et al. 2018, etc.

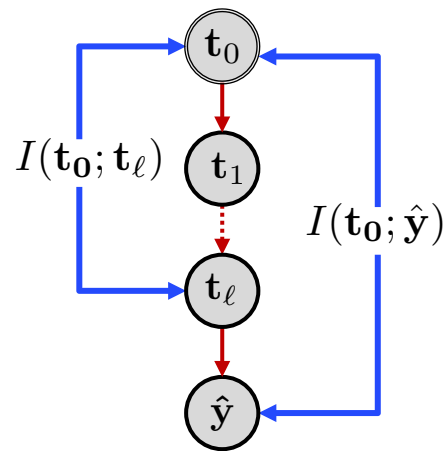
## Computing mutual informations in neural networks ... is tricky

- Intractable  $p(\mathbf{t}_\ell) = \int d\mathbf{t}_0 d\mathbf{t}_1 \cdots d\mathbf{t}_{\ell-1} p(\mathbf{t}_\ell | \mathbf{t}_{\ell-1}) \cdots p(\mathbf{t}_1 | \mathbf{t}_0) p(\mathbf{t}_0)$
- Computing  $I(\mathbf{t}_0; \mathbf{t}_\ell)$  or  $H(\mathbf{t}_\ell) = - \int d\mathbf{t}_\ell p(\mathbf{t}_\ell) \log p(\mathbf{t}_\ell)$  yet more involved

## Sampled based estimations (non-parametric, variational methods)

Kraskov et al. 2014, Kolchinsky et al. 2017, Belghazi et al. 2017, Goldfeld et al. 2018 etc.

- Less and less reliable as networks get large
- Need for benchmark controlled case



Here 1) Restrict to **models** of DNNs, 2) Leverage statistical physics **replica method**

Model: fully connected feed forward neural network  
with stochastic activations

$$\mathbf{t}_\ell \stackrel{\mathbb{R}^{N_\ell}}{=} f_\epsilon(\mathbf{W}^{(\ell)} f_\epsilon(\mathbf{W}^{(\ell-1)} \cdots f_\epsilon(\mathbf{W}^{(1)} \mathbf{t}_0 \cdots))) \stackrel{\mathbb{R}^{N_0}}{}$$

Model: fully connected feed forward neural network  
with stochastic activations

$$\mathbf{t}_\ell \in \mathbb{R}^{N_\ell} = f_\epsilon(\mathbf{W}^{(\ell)} f_\epsilon(\mathbf{W}^{(\ell-1)} \cdots f_\epsilon(\mathbf{W}^{(1)} \mathbf{t}_0 \cdots))) \in \mathbb{R}^{N_0}$$

With :

Model: fully connected feed forward neural network  
with stochastic activations

$$\mathbf{t}_\ell \in \mathbb{R}^{N_\ell} = f_\epsilon(\mathbf{W}^{(\ell)} f_\epsilon(\mathbf{W}^{(\ell-1)} \cdots f_\epsilon(\mathbf{W}^{(1)} \mathbf{t}_0) \cdots)) \in \mathbb{R}^{N_0}$$

With :

1) known factorized input distribution  $P_0(\mathbf{t}_0) = \prod_{i=1}^{N_0} P_0(t_{0i})$

Model: fully connected feed forward neural network  
with stochastic activations

$$\mathbf{t}_\ell \in \mathbb{R}^{N_\ell} = f_\epsilon(\mathbf{W}^{(\ell)} f_\epsilon(\mathbf{W}^{(\ell-1)} \cdots f_\epsilon(\mathbf{W}^{(1)} \mathbf{t}_0) \cdots)) \in \mathbb{R}^{N_0}$$

With :

1) known factorized input distribution  $P_0(\mathbf{t}_0) = \prod_{i=1}^{N_0} P_0(t_{0i})$

2) weight matrices from class of random matrix class

# Model: fully connected feed forward neural network with stochastic activations

$$\mathbb{R}^{N_\ell} \quad \mathbf{t}_\ell = f_\epsilon(\mathbf{W}^{(\ell)} f_\epsilon(\mathbf{W}^{(\ell-1)} \cdots f_\epsilon(\mathbf{W}^{(1)} \mathbf{t}_0) \cdots)) \quad \mathbb{R}^{N_0}$$

With :

1) known factorized input distribution  $P_0(\mathbf{t}_0) = \prod_{i=1}^{N_0} P_0(t_{0i})$

2) weight matrices from class of random matrix class

i.i.d Gaussian

$$W^{(\ell)}$$



# Model: fully connected feed forward neural network with stochastic activations

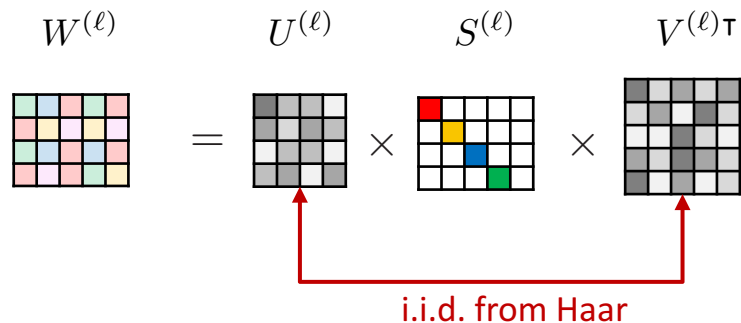
$$\mathbb{R}^{N_\ell} \mathbf{t}_\ell = f_\epsilon(\mathbf{W}^{(\ell)} f_\epsilon(\mathbf{W}^{(\ell-1)} \cdots f_\epsilon(\mathbf{W}^{(1)} \mathbf{t}_0) \cdots)) \mathbb{R}^{N_0}$$

With :

1) known factorized input distribution  $P_0(\mathbf{t}_0) = \prod_{i=1}^{N_0} P_0(t_{0i})$

2) weight matrices from class of random matrix class

i.i.d Gaussian  $\rightarrow$  Orthogonally invariant



# Model: fully connected feed forward neural network with stochastic activations

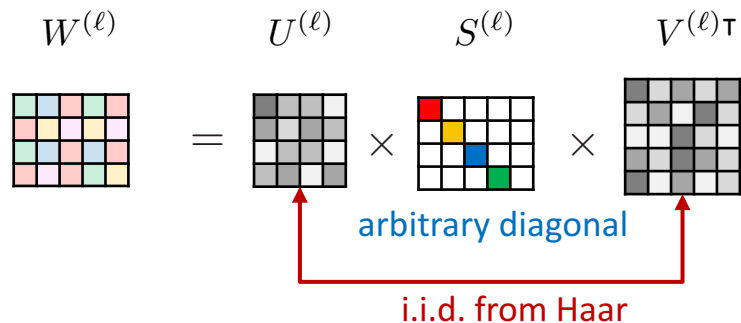
$$\mathbb{R}^{N_\ell} \mathbf{t}_\ell = f_\epsilon(\mathbf{W}^{(\ell)} f_\epsilon(\mathbf{W}^{(\ell-1)} \cdots f_\epsilon(\mathbf{W}^{(1)} \mathbf{t}_0) \cdots)) \mathbb{R}^{N_0}$$

With :

1) known factorized input distribution  $P_0(\mathbf{t}_0) = \prod_{i=1}^{N_0} P_0(t_{0i})$

2) weight matrices from class of random matrix class

i.i.d Gaussian  $\rightarrow$  Orthogonally invariant





# Main results

1) Formula:  $\lim_{N_0 \rightarrow \infty} \frac{1}{N_0} H(\mathbf{t}_\ell) = \min_{\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}} \text{extr} \phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}})$

# Main results

1) Formula:  $\lim_{N_0 \rightarrow \infty} \frac{1}{N_0} H(\mathbf{t}_\ell) = \min_{\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}} \text{extr} \phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}})$

# Main results

1) Formula:  $\lim_{N_0 \rightarrow \infty} \frac{1}{N_0} H(\mathbf{t}_\ell) = \min_{\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}} \text{extr} \phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}})$

potential: 
$$\begin{aligned} \phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}) = & I \left( t_0; t_0 + \frac{\xi_0}{\sqrt{\tilde{A}_1}} \right) - \frac{1}{2} \sum_{k=1}^{\ell} \tilde{\alpha}_{k-1} [\tilde{A}_k V_k + \alpha_k A_k \tilde{V}_k - F_{W_k}(A_k V_k)] \\ & + \sum_{k=1}^{\ell-1} \tilde{\alpha}_k \left[ H(t_k | \xi_k; \tilde{A}_{k+1}, \tilde{V}_k, \tilde{\rho}_k) - \frac{1}{2} \log(2\pi e \tilde{A}_{k+1}^{-1}) \right] + \tilde{\alpha}_\ell H(t_\ell | \xi_\ell; \tilde{V}_\ell, \tilde{\rho}_\ell) \end{aligned}$$

# Main results

1) Formula:  $\lim_{N_0 \rightarrow \infty} \frac{1}{N_0} H(\mathbf{t}_\ell) = \min_{\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}} \text{extr} \phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}})$

potential: 
$$\phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}) = I \left( t_0; t_0 + \frac{\xi_0}{\sqrt{\tilde{A}_1}} \right) - \frac{1}{2} \sum_{k=1}^{\ell} \tilde{\alpha}_{k-1} [\tilde{A}_k V_k + \alpha_k A_k \tilde{V}_k - F_{W_k}(A_k V_k)]$$
$$+ \sum_{k=1}^{\ell-1} \tilde{\alpha}_k \left[ H(t_k | \xi_k; \tilde{A}_{k+1}, \tilde{V}_k, \tilde{\rho}_k) - \frac{1}{2} \log(2\pi e \tilde{A}_{k+1}^{-1}) \right] + \tilde{\alpha}_\ell H(t_\ell | \xi_\ell; \tilde{V}_\ell, \tilde{\rho}_\ell)$$

# Main results

1) Formula:  $\lim_{N_0 \rightarrow \infty} \frac{1}{N_0} H(\mathbf{t}_\ell) = \min_{\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}} \text{extr} \phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}})$

potential: 
$$\phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}) = I\left(t_0; t_0 + \frac{\xi_0}{\sqrt{\tilde{A}_1}}\right) - \frac{1}{2} \sum_{k=1}^{\ell} \tilde{\alpha}_{k-1} [\tilde{A}_k V_k + \alpha_k A_k \tilde{V}_k - F_{W_k}(A_k V_k)]$$
$$+ \sum_{k=1}^{\ell-1} \tilde{\alpha}_k \left[ H(t_k | \xi_k; \tilde{A}_{k+1}, \tilde{V}_k, \tilde{\rho}_k) - \frac{1}{2} \log(2\pi e \tilde{A}_{k+1}^{-1}) \right] + \tilde{\alpha}_\ell H(t_\ell | \xi_\ell; \tilde{V}_\ell, \tilde{\rho}_\ell)$$

# Main results

1) Formula:  $\lim_{N_0 \rightarrow \infty} \frac{1}{N_0} H(\mathbf{t}_\ell) = \min_{\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}} \text{extr} \phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}})$

$$\text{potential: } \phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}) = I\left(t_0; t_0 + \frac{\xi_0}{\sqrt{\tilde{A}_1}}\right) - \frac{1}{2} \sum_{k=1}^{\ell} \tilde{\alpha}_{k-1} [\tilde{A}_k V_k + \alpha_k A_k \tilde{V}_k - F_{W_k}(A_k V_k)] \\ + \sum_{k=1}^{\ell-1} \tilde{\alpha}_k \left[ H(t_k | \xi_k; \tilde{A}_{k+1}, \tilde{V}_k, \tilde{\rho}_k) - \frac{1}{2} \log(2\pi e \tilde{A}_{k+1}^{-1}) \right] + \tilde{\alpha}_\ell H(t_\ell | \xi_\ell; \tilde{V}_\ell, \tilde{\rho}_\ell)$$

efficient implementation of extremization:  [sphinxtteam](#) / [dnner](#)

# Main results

1) Formula:  $\lim_{N_0 \rightarrow \infty} \frac{1}{N_0} H(\mathbf{t}_\ell) = \min_{\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}} \text{extr} \phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}})$

$$\text{potential: } \phi_\ell(\mathbf{A}, \mathbf{V}, \tilde{\mathbf{A}}, \tilde{\mathbf{V}}) = I\left(t_0; t_0 + \frac{\xi_0}{\sqrt{\tilde{A}_1}}\right) - \frac{1}{2} \sum_{k=1}^{\ell} \tilde{\alpha}_{k-1} [\tilde{A}_k V_k + \alpha_k A_k \tilde{V}_k - F_{W_k}(A_k V_k)] \\ + \sum_{k=1}^{\ell-1} \tilde{\alpha}_k \left[ H(t_k | \xi_k; \tilde{A}_{k+1}, \tilde{V}_k, \tilde{\rho}_k) - \frac{1}{2} \log(2\pi e \tilde{A}_{k+1}^{-1}) \right] + \tilde{\alpha}_\ell H(t_\ell | \xi_\ell; \tilde{V}_\ell, \tilde{\rho}_\ell)$$

efficient implementation of extremization: [sphinxtteam / dnnr](#)

## 2) Theorem:

$$\left\{ \begin{array}{l} \text{2-layer } \mathbf{t}_2 = f_\epsilon(\mathbf{W}^{(2)} f_\epsilon(\mathbf{W}^{(1)} \mathbf{t}_0)) \\ \text{Gaussian random weight matrices} \end{array} \right.$$

$\Rightarrow$

replica prediction exact  
in limit of large networks



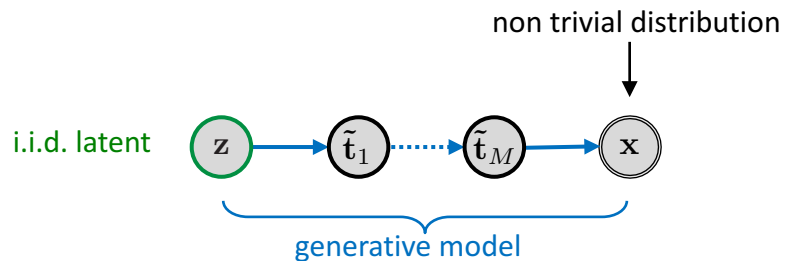
**APPLICATION:**

**Follow mutual information  
during SGD learning**



# Applying the replica formula to learning experiments

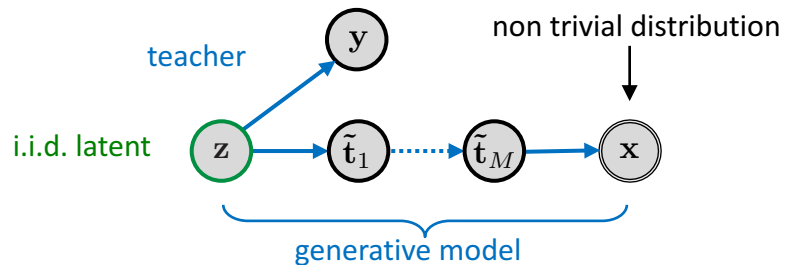
## 1) Synthetic data sets for student



Keras implementation

# Applying the replica formula to learning experiments

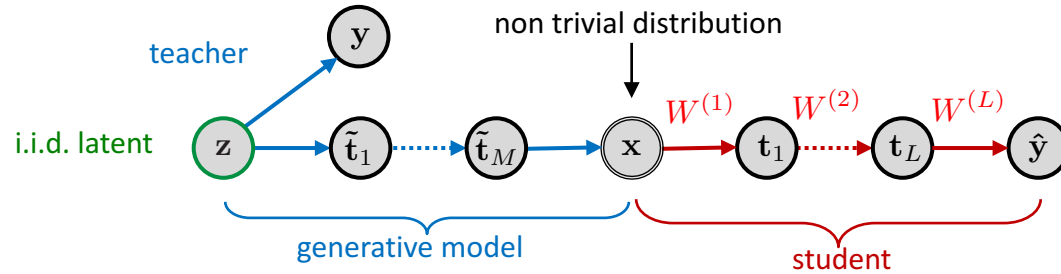
## 1) Synthetic data sets for student



Keras implementation

# Applying the replica formula to learning experiments

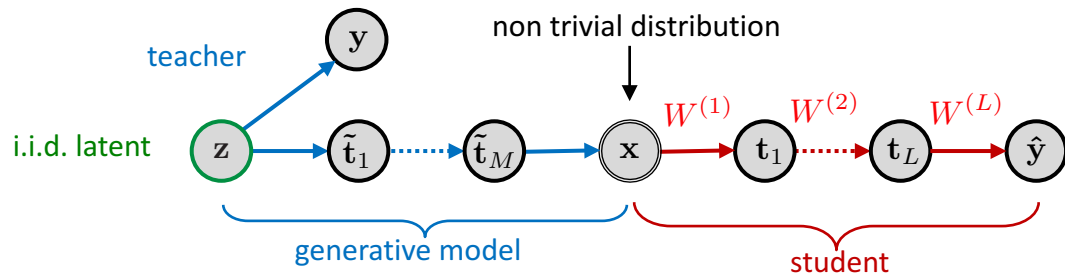
## 1) Synthetic data sets for student



Keras implementation

# Applying the replica formula to learning experiments

## 1) Synthetic data sets for student



## 2) Orthogonally invariant weight matrices during learning for the student

$$W^{(\ell)} = U^{(\ell)} \times S^{(\ell)} \times V^{(\ell)}$$

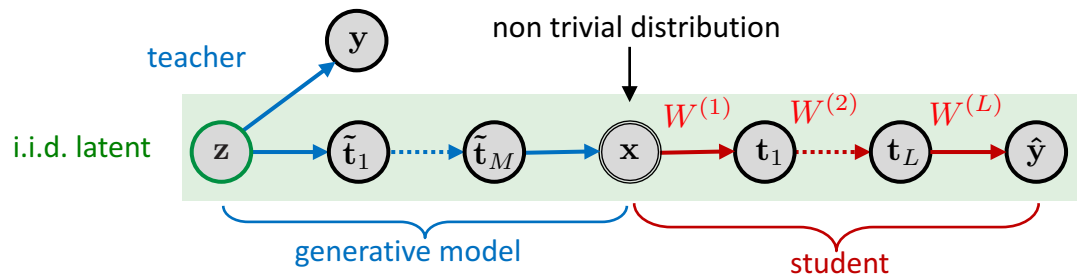
The diagram shows the decomposition of a weight matrix  $W^{(\ell)}$  into three components:  $U^{(\ell)}$  (fixed),  $S^{(\ell)}$  (updated by backprop), and  $V^{(\ell)}$  (fixed). Each component is represented by a 4x4 grid of colored squares.

- $U^{(\ell)}$ : A 4x4 grid of grey squares, labeled **fixed**.
- $S^{(\ell)}$ : A 4x4 grid with a red square at (1,1), a yellow square at (2,2), a blue square at (3,3), and a green square at (4,4), labeled **updated by backprop**.
- $V^{(\ell)}$ : A 4x4 grid of grey squares, labeled **fixed**.

Keras implementation

# Applying the replica formula to learning experiments

## 1) Synthetic data sets for student



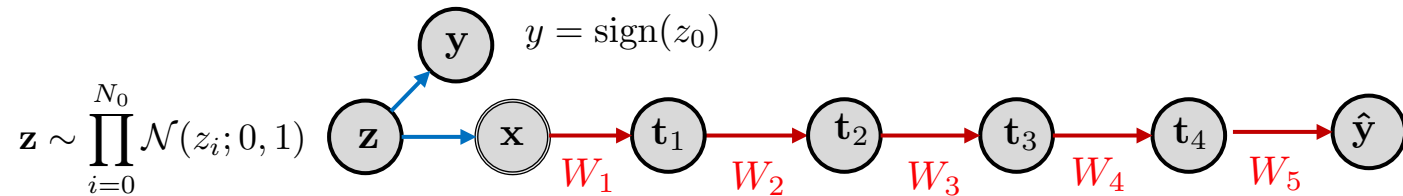
## 2) Orthogonally invariant weight matrices during learning for the student

$$W^{(\ell)} = U^{(\ell)} \times S^{(\ell)} \times V^{(\ell)}$$

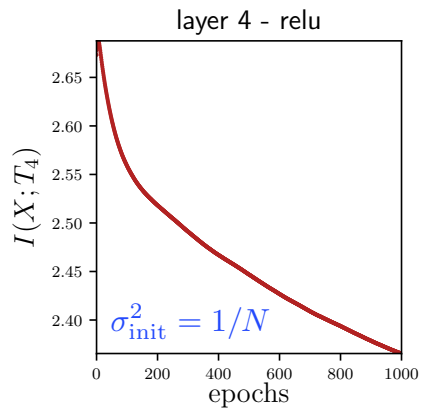
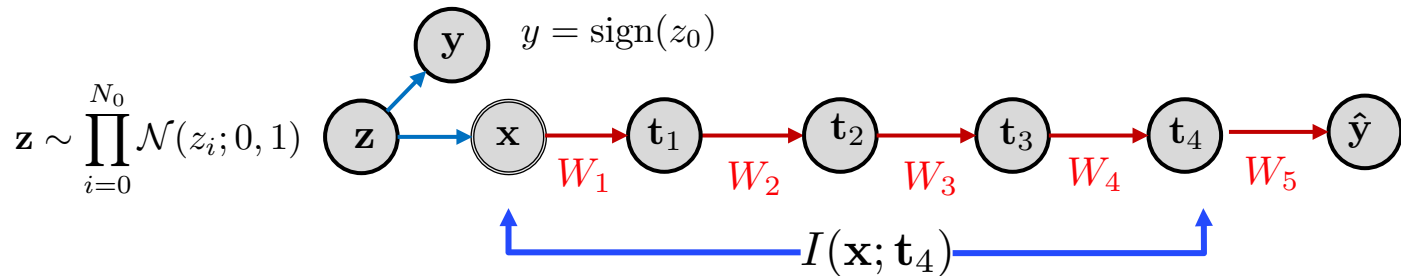
The diagram shows the decomposition of a weight matrix  $W^{(\ell)}$  into three components:  $U^{(\ell)}$  (fixed),  $S^{(\ell)}$  (updated by backprop), and  $V^{(\ell)}$  (fixed). Each component is represented by a 4x4 grid of colored squares.  $U^{(\ell)}$  is a grey grid,  $S^{(\ell)}$  is a white grid with a few colored squares (red, yellow, blue, green), and  $V^{(\ell)}$  is a grey grid.

Keras implementation

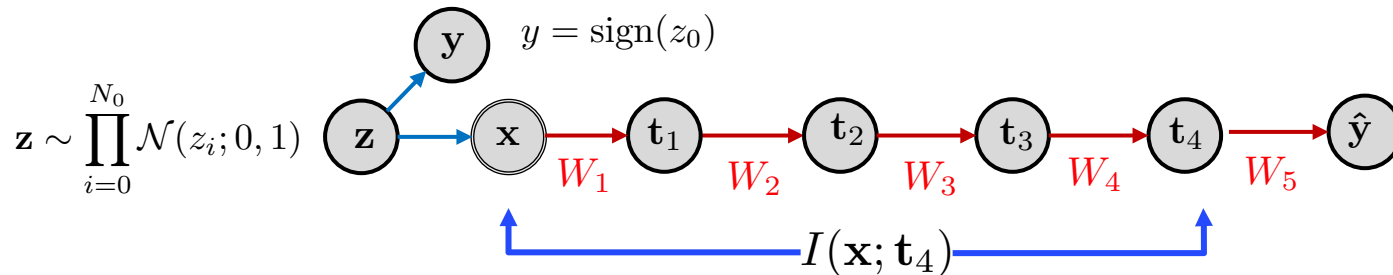
# Exploratory experiment: A binary classification example



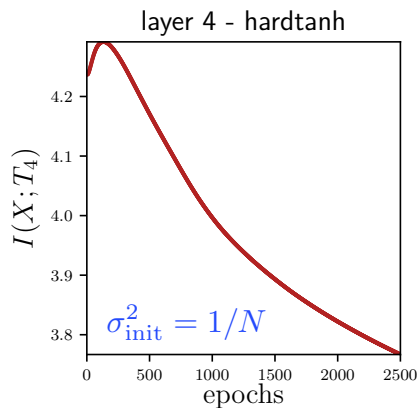
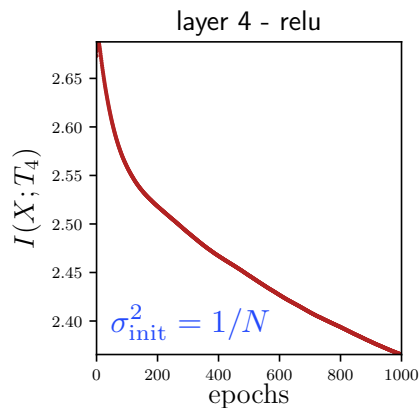
# Exploratory experiment: A binary classification example



# Exploratory experiment: A binary classification example

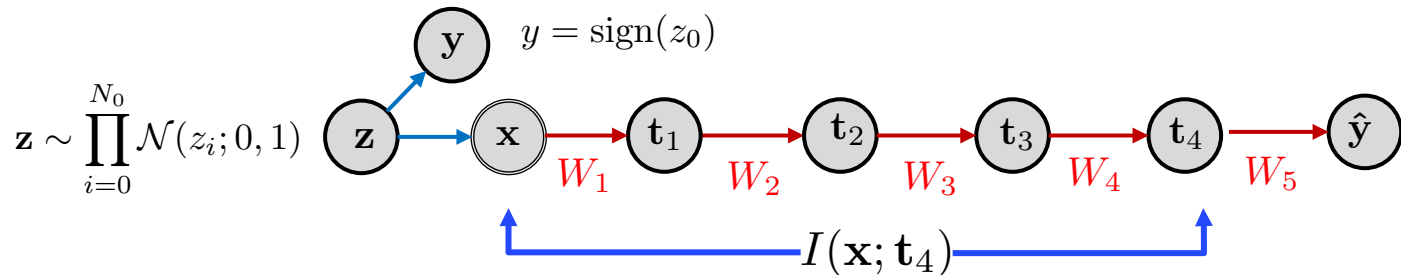


## Different non-linearities



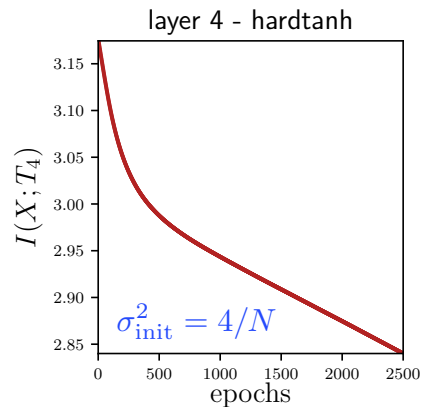
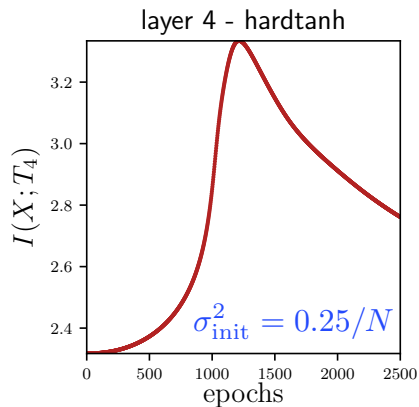
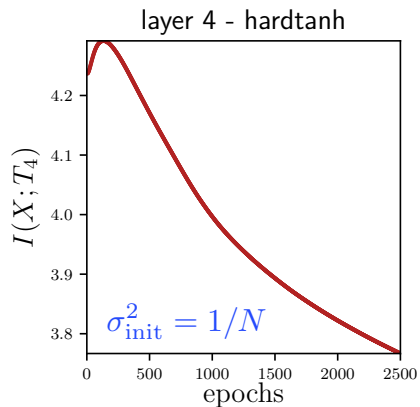
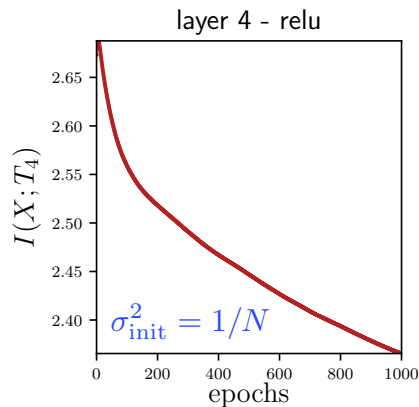


# Exploratory experiment: A binary classification example

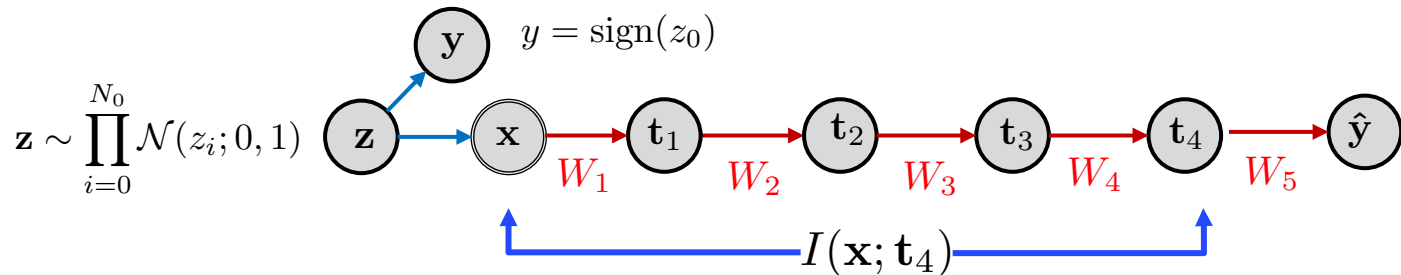


Different non-linearities

Different initializations

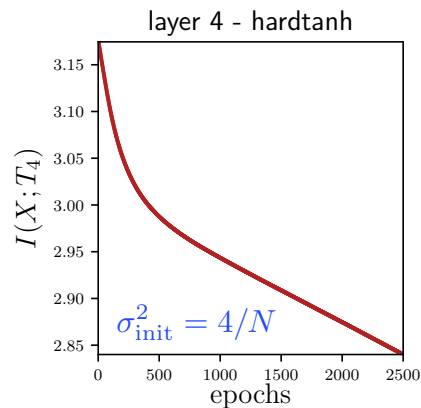
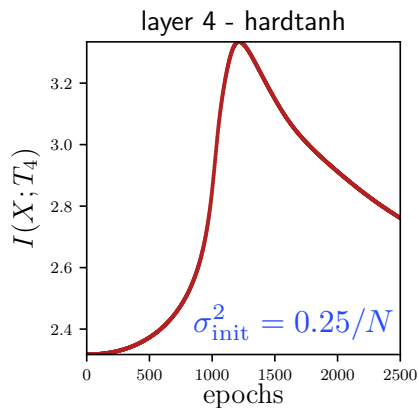
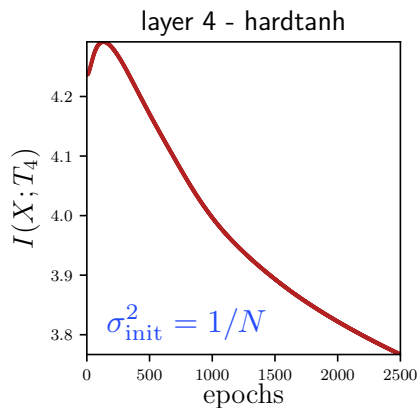
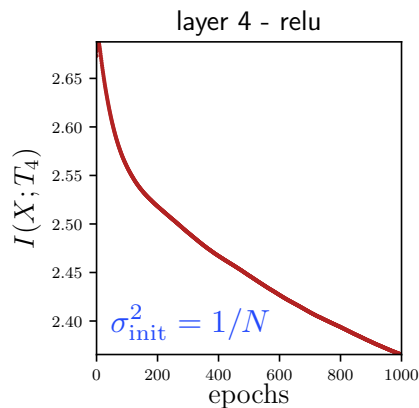


# Exploratory experiment: A binary classification example



Different non-linearities

Different initializations



Poster @ Room 210 & 230 AB #110