

# Estimating the Information Flow in Deep Neural Networks

Ziv Goldfeld

MIT

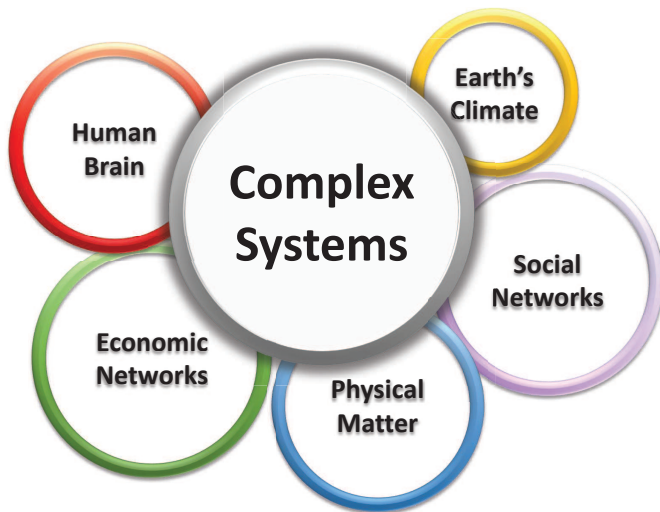
# Information in the Modern Age & Complex Systems

# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions

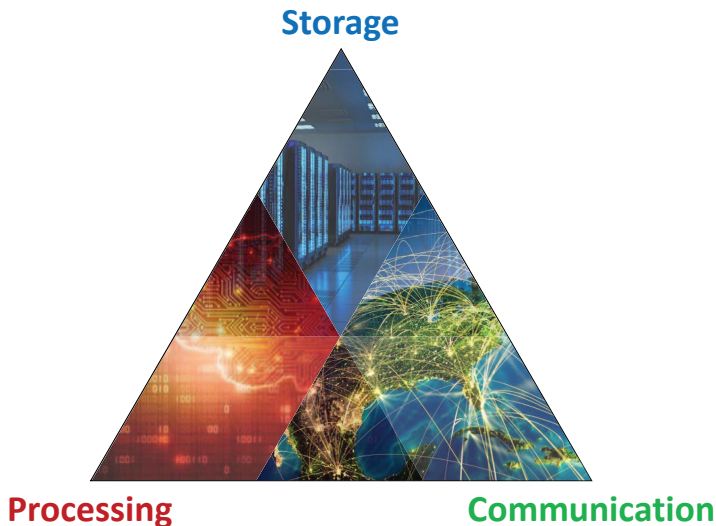
# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions



# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions



# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions

**Storage**



**Processing**

**Communication**

# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions

**Storage**

Emerging Technologies:

Shrink magnetic region per bit



**Processing**

**Communication**

# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions

## Storage

### Emerging Technologies:

Shrink magnetic region per bit

### Challenges:

Stabilization of written data



**Processing**

**Communication**

# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions

## Storage

### Emerging Technologies:

Shrink magnetic region per bit

### Challenges:

Stabilization of written data

### Model & Study:

Interacting particle sys.



Storage capacity &

HDD designs

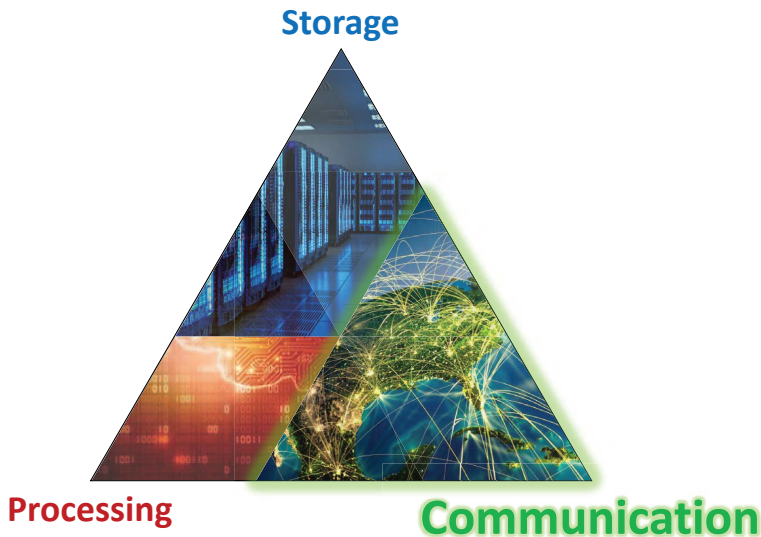


**Processing**

**Communication**

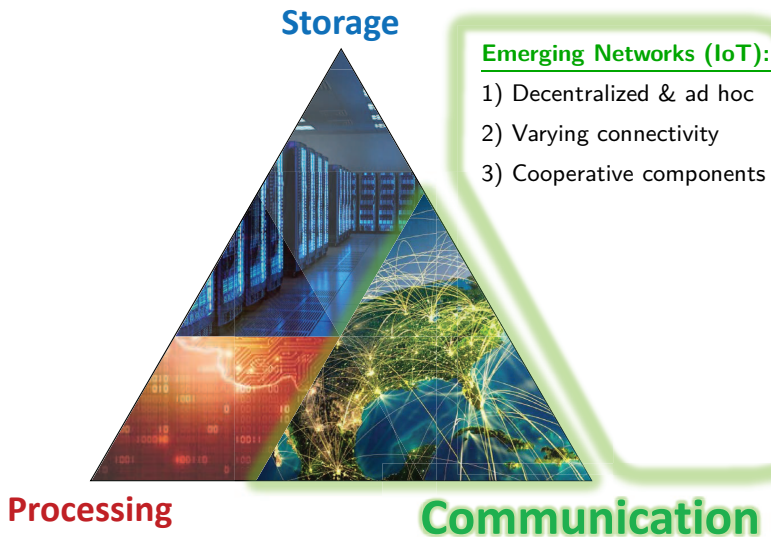
# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions



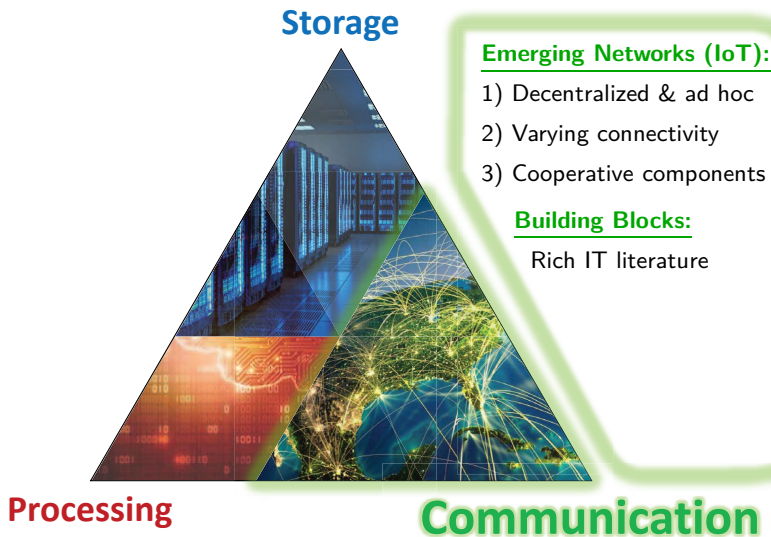
# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions



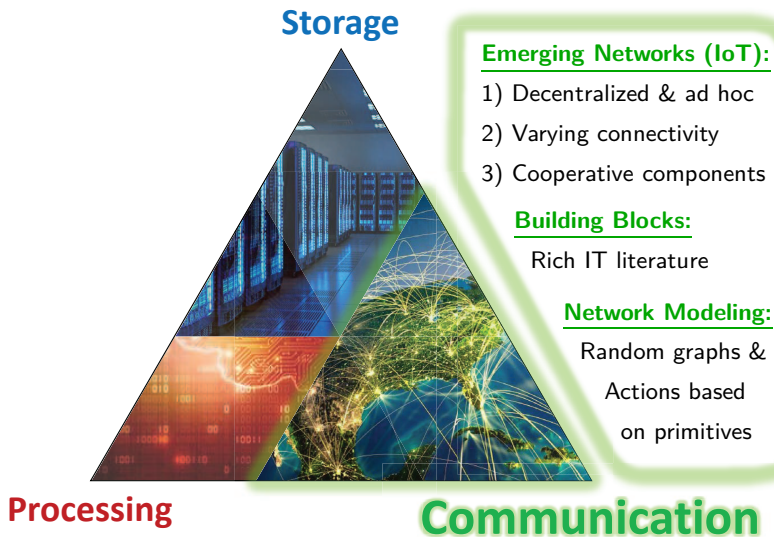
# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions



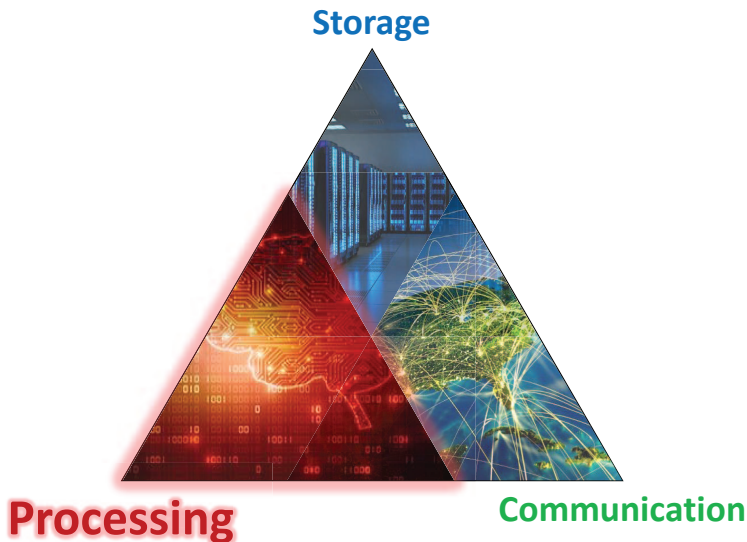
# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions



# Information in the Modern Age & Complex Systems

**Complex System:** Multi-component system driven by local interactions

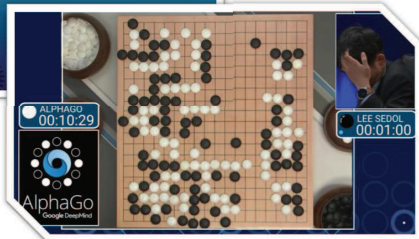
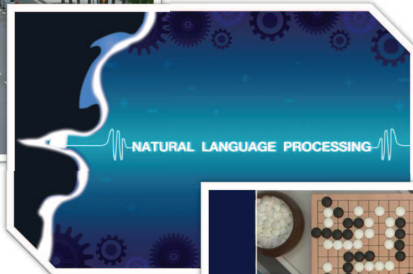


# Deep Neural Networks

- Unprecedented practical success in hosts of tasks

# Deep Neural Networks

- Unprecedented practical success in hosts of tasks



# Deep Neural Networks

- Unprecedented practical success in hosts of tasks
- Lacking theory:

# Deep Neural Networks

- Unprecedented practical success in hosts of tasks
- Lacking theory:
  - ▶ What drives the evolution of hidden representations?

# Deep Neural Networks

- Unprecedented practical success in hosts of tasks
- Lacking theory:
  - ▶ What drives the evolution of hidden representations?
  - ▶ What are properties of learned representations?

# Deep Neural Networks

- Unprecedented practical success in hosts of tasks
- Lacking theory:
  - ▶ What drives the evolution of hidden representations?
  - ▶ What are properties of learned representations?
  - ▶ How fully trained networks process information?

# Deep Neural Networks

- Unprecedented practical success in hosts of tasks
- Lacking theory:
  - ▶ What drives the evolution of hidden representations?
  - ▶ What are properties of learned representations?
  - ▶ How fully trained networks process information?
  - ▶

# Deep Neural Networks

- Unprecedented practical success in hosts of tasks
- Lacking theory:
  - ▶ What drives the evolution of hidden representations?
  - ▶ What are properties of learned representations?
  - ▶ How fully trained networks process information?
  - ▶
- Some past attempts to understand effectiveness of deep learning

# Deep Neural Networks

- Unprecedented practical success in hosts of tasks
- Lacking theory:
  - ▶ What drives the evolution of hidden representations?
  - ▶ What are properties of learned representations?
  - ▶ How fully trained networks process information?
  - ▶  $\vdots$
- Some past attempts to understand effectiveness of deep learning
  - ▶ Optimization dynamics in parameter space  
[Saxe-McClelland-Ganguli'14, Choromanska-et al'15, Wei-Lee-Ma'18]

# Deep Neural Networks

- Unprecedented practical success in hosts of tasks
- Lacking theory:
  - ▶ What drives the evolution of hidden representations?
  - ▶ What are properties of learned representations?
  - ▶ How fully trained networks process information?
  - ▶  $\vdots$
- Some past attempts to understand effectiveness of deep learning
  - ▶ Optimization dynamics in parameter space  
[Saxe-McClelland-Ganguli'14, Choromanska-et al'15, Wei-Lee-Ma'18]
  - ▶ Classes of efficiently representable functions  
[Montufar-Pascanu-Cho-Bengio'14, Poggio-Mhaskar-Rosasco-Miranda-Liao'17]

# Deep Neural Networks

- Unprecedented practical success in hosts of tasks
- Lacking theory:
  - ▶ What drives the evolution of hidden representations?
  - ▶ What are properties of learned representations?
  - ▶ How fully trained networks process information?
  - ▶  $\vdots$
- Some past attempts to understand effectiveness of deep learning
  - ▶ Optimization dynamics in parameter space  
[Saxe-McClelland-Ganguli'14, Choromanska-et al'15, Wei-Lee-Ma'18]
  - ▶ Classes of efficiently representable functions  
[Montufar-Pascanu-Cho-Bengio'14, Poggio-Mhaskar-Rosasco-Miranda-Liao'17]
  - ▶ Information Bottleneck Theory  
[Tishby-Zaslavsky'15, Shwartz-Tishby'17, Saxe et al.'18]

# Deep Neural Networks

- Unprecedented practical success in hosts of tasks
- Lacking theory:
  - ▶ What drives the evolution of hidden representations?
  - ▶ What are properties of learned representations?
  - ▶ How fully trained networks process information?
  - ▶  $\vdots$
- Some past attempts to understand effectiveness of deep learning
  - ▶ Optimization dynamics in parameter space  
[Saxe-McClelland-Ganguli'14, Choromanska-et al'15, Wei-Lee-Ma'18]
  - ▶ Classes of efficiently representable functions  
[Montufar-Pascanu-Cho-Bengio'14, Poggio-Mhaskar-Rosasco-Miranda-Liao'17]
  - ▶ **Information Bottleneck Theory**  
[Tishby-Zaslavsky'15, Shwartz-Tishby'17, Saxe et al.'18]

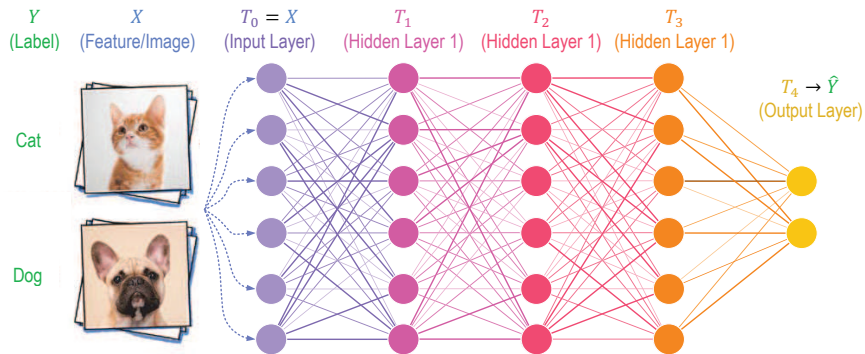
# Deep Neural Networks

- Unprecedented practical success in hosts of tasks
- Lacking theory:
  - ▶ What drives the evolution of hidden representations?
  - ▶ What are properties of learned representations?
  - ▶ How fully trained networks process information?
  - ▶  $\vdots$
- Some past attempts to understand effectiveness of deep learning
  - ▶ Optimization dynamics in parameter space  
[Saxe-McClelland-Ganguli'14, Choromanska-et al'15, Wei-Lee-Ma'18]
  - ▶ Classes of efficiently representable functions  
[Montufar-Pascanu-Cho-Bengio'14, Poggio-Mhaskar-Rosasco-Miranda-Liao'17]
  - ▶ **Information Bottleneck Theory**  
[Tishby-Zaslavsky'15, Shwartz-Tishby'17, Saxe et al.'18]

★ Goal: Explain 'compression' in Information Bottleneck framework

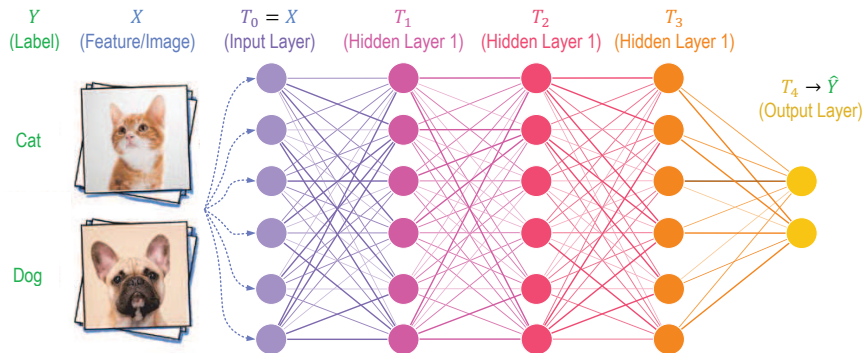
# Setup and Preliminaries

## Feedforward DNN for Classification:



# Setup and Preliminaries

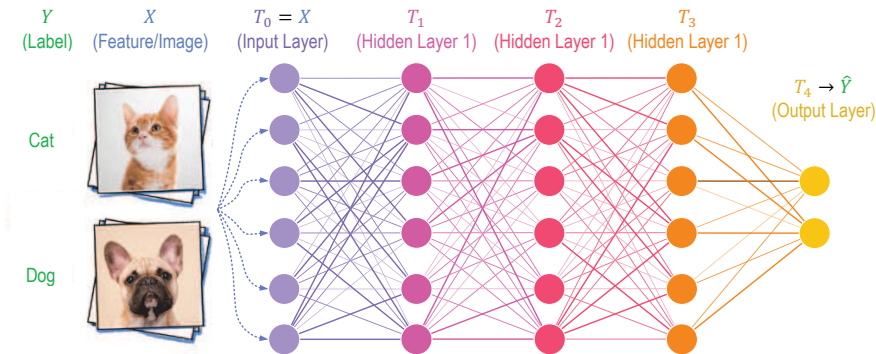
## Feedforward DNN for Classification:



- **Deterministic DNN:**  $T_\ell = f_\ell(T_{\ell-1})$  (MLP:  $T_\ell = \sigma(W_\ell T_{\ell-1} + b_\ell)$ )

# Setup and Preliminaries

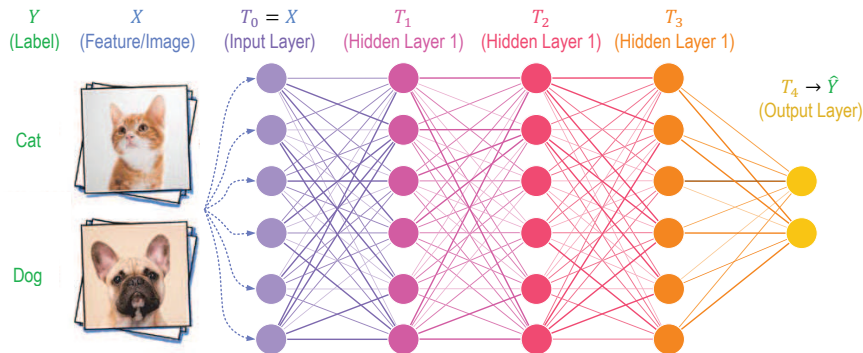
## Feedforward DNN for Classification:



- **Deterministic DNN:**  $T_\ell = f_\ell(T_{\ell-1})$  (MLP:  $T_\ell = \sigma(W_\ell T_{\ell-1} + b_\ell)$ )
- **Joint Distribution:**  $P_{X,Y}$

# Setup and Preliminaries

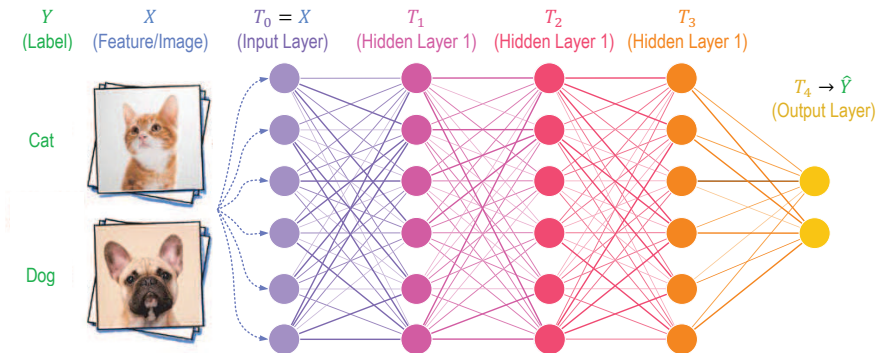
## Feedforward DNN for Classification:



- **Deterministic DNN:**  $T_\ell = f_\ell(T_{\ell-1})$  (MLP:  $T_\ell = \sigma(W_\ell T_{\ell-1} + b_\ell)$ )
- **Joint Distribution:**  $P_{X,Y} \implies P_{X,Y} \cdot P_{T_1, \dots, T_L | X}$

# Setup and Preliminaries

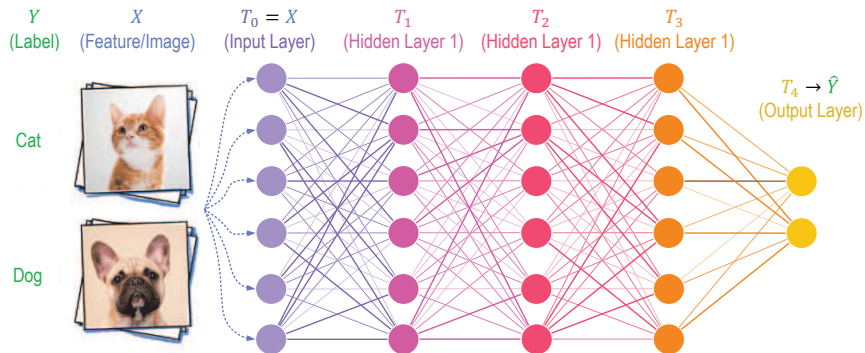
## Feedforward DNN for Classification:



- **Deterministic DNN:**  $T_\ell = f_\ell(T_{\ell-1})$  (**MLP:**  $T_\ell = \sigma(W_\ell T_{\ell-1} + b_\ell)$ )
- **Joint Distribution:**  $P_{X,Y} \implies P_{X,Y} \cdot P_{T_1, \dots, T_L | X}$
- **IB Theory:** Track MI pairs  $(I(X; T_\ell), I(Y; T_\ell))$  (information plane)

# Setup and Preliminaries

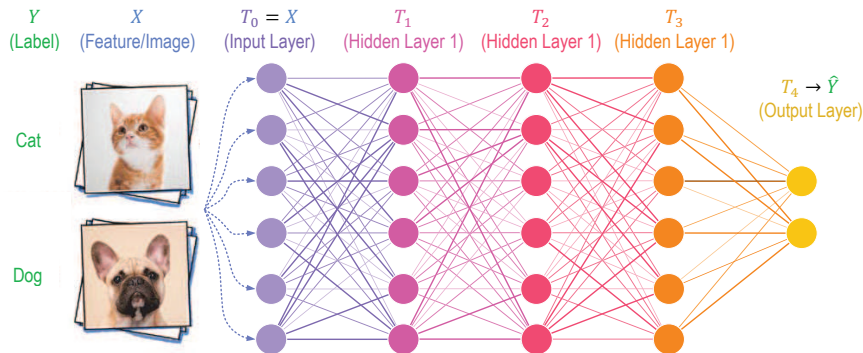
## Feedforward DNN for Classification:



IB Theory Claim: Training comprises 2 phases

# Setup and Preliminaries

## Feedforward DNN for Classification:

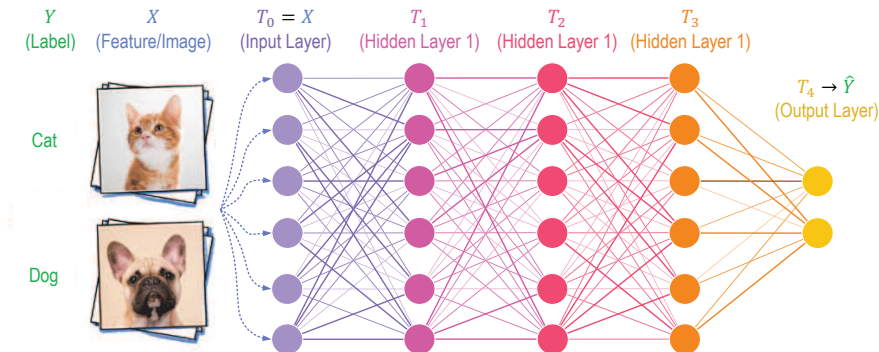


**IB Theory Claim:** Training comprises 2 phases

- **Fitting:**  $I(Y; T_\ell)$  &  $I(X; T_\ell)$  rise (short)

# Setup and Preliminaries

## Feedforward DNN for Classification:

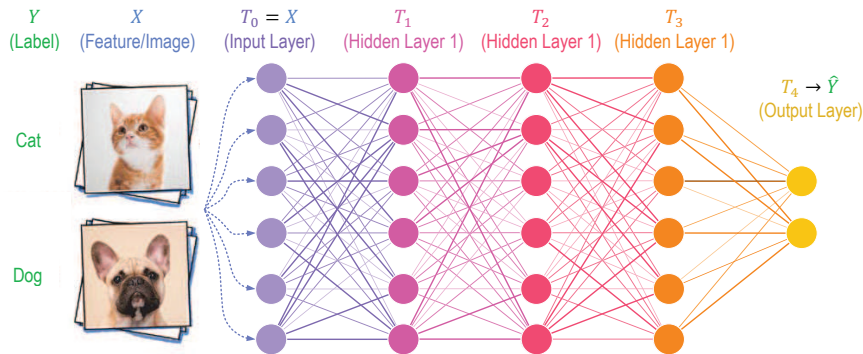


**IB Theory Claim:** Training comprises 2 phases

- **Fitting:**  $I(Y; T_\ell)$  &  $I(X; T_\ell)$  rise (short)
- **Compression:**  $I(X; T_\ell)$  slowly drops (long)

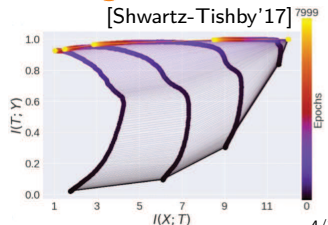
# Setup and Preliminaries

## Feedforward DNN for Classification:



IB Theory Claim: Training comprises 2 phases

- **Fitting:**  $I(Y; T_\ell)$  &  $I(X; T_\ell)$  rise (short)
- **Compression:**  $I(X; T_\ell)$  slowly drops (long)



# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is independent of the DNN parameters

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is independent of the DNN parameters

Why?

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is independent of the DNN parameters

## Why?

- **Continuous  $X$ :**

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is independent of the DNN parameters

## Why?

- **Continuous  $X$ :**  $I(X; T_\ell) = h(T_\ell) - h(T_\ell|X)$

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is independent of the DNN parameters

## Why?

- **Continuous  $X$ :**  $I(X; T_\ell) = h(T_\ell) - \textcolor{red}{h}(T_\ell | \textcolor{red}{X})$

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is independent of the DNN parameters

## Why?

• **Continuous  $X$ :** 
$$I(X; T_\ell) = h(T_\ell) - \textcolor{red}{h}(\textcolor{red}{\tilde{f}_\ell(X)}|\textcolor{red}{X})$$

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is independent of the DNN parameters

## Why?

- **Continuous  $X$ :**

$$I(X; T_\ell) = h(T_\ell) - \underbrace{h(\tilde{f}_\ell(X)|X)}_{=-\infty}$$

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is independent of the DNN parameters

## Why?

- **Continuous  $X$ :**

$$I(X; T_\ell) = h(T_\ell) - h(\tilde{f}_\ell(X)|X) = \infty$$

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is independent of the DNN parameters

## Why?

- **Continuous  $X$ :**  $I(X; T_\ell) = h(T_\ell) - h(\tilde{f}_\ell(X)|X) = \infty$
- **Discrete  $X$ :**

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is independent of the DNN parameters

## Why?

- **Continuous  $X$ :**  $I(X; T_\ell) = h(T_\ell) - h(\tilde{f}_\ell(X)|X) = \infty$
- **Discrete  $X$ :** The map  $X \mapsto T_\ell$  is injective<sup>★</sup>

★ For almost all weight matrices and bias vectors

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is independent of the DNN parameters

## Why?

- **Continuous  $X$ :**  $I(X; T_\ell) = h(T_\ell) - h(\tilde{f}_\ell(X)|X) = \infty$
- **Discrete  $X$ :** The map  $X \mapsto T_\ell$  is injective<sup>★</sup>  $\implies I(X; T_\ell) = H(X)$

★ For almost all weight matrices and bias vectors

# Vacuous Mutual Information

## Observation

*Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)*  
 $\implies I(X; T_\ell)$  is **independent of the DNN parameters**

## Why?

- **Continuous  $X$ :**  $I(X; T_\ell) = h(T_\ell) - h(\tilde{f}_\ell(X)|X) = \infty$
- **Discrete  $X$ :** The map  $X \mapsto T_\ell$  is injective\*  $\implies I(X; T_\ell) = \mathbf{H}(X)$

# Vacuous Mutual Information

## Observation

Det. DNNs with strictly monotone nonlinearities (e.g., tanh or sigmoid)

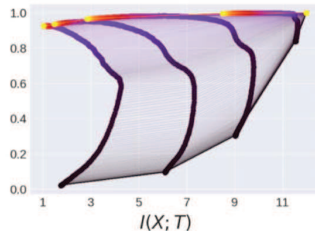
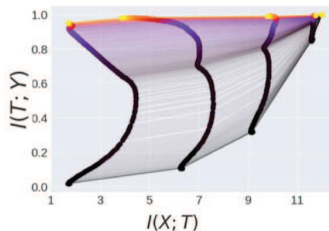
$\implies I(X; T_\ell)$  is **independent of the DNN parameters**

## Why?

- **Continuous  $X$ :**  $I(X; T_\ell) = h(T_\ell) - h(\tilde{f}_\ell(X)|X) = \infty$
- **Discrete  $X$ :** The map  $X \mapsto T_\ell$  is injective\*  $\implies I(X; T_\ell) = \mathbf{H}(X)$

## Past Works:

[Shwartz-Tishby'17,  
Saxe et al.'18]



# What is going on here?

- Plots via binning-based estimator of  $I(X; T_\ell)$ , for  $X \sim \text{Unif}(\text{dataset})$

# What is going on here?

- Plots via binning-based estimator of  $I(X; T_\ell)$ , for  $X \sim \text{Unif}(\text{dataset})$   
 $\implies$  Plotted values are  $I(X; \text{Bin}(T_\ell))$

# What is going on here?

- Plots via binning-based estimator of  $I(X; T_\ell)$ , for  $X \sim \text{Unif}(\text{dataset})$   
 $\implies$  Plotted values are  $I(X; \text{Bin}(T_\ell)) \overset{??}{\approx} I(X; T_\ell)$

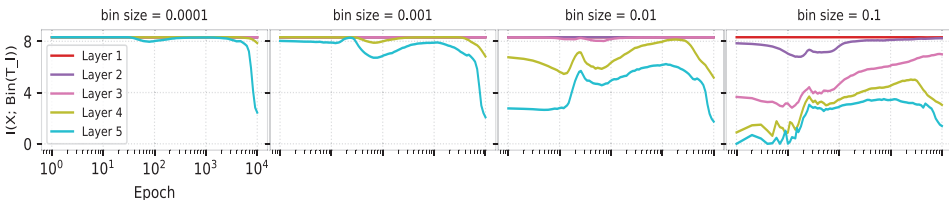
# What is going on here?

- Plots via binning-based estimator of  $I(X; T_\ell)$ , for  $X \sim \text{Unif}(\text{dataset})$   
 $\implies$  Plotted values are  $I(X; \text{Bin}(T_\ell)) \overset{??}{\approx} I(X; T_\ell)$  **No!**

# What is going on here?

- Plots via binning-based estimator of  $I(X; T_\ell)$ , for  $X \sim \text{Unif}(\text{dataset})$

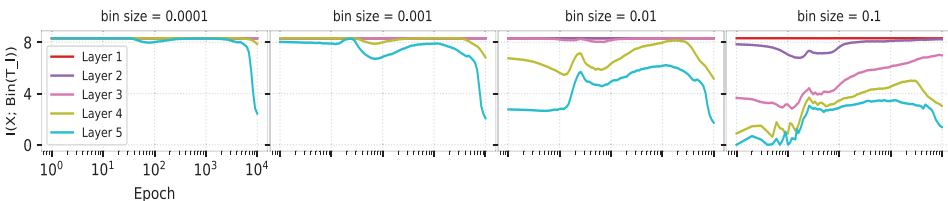
$\Rightarrow$  Plotted values are  $I(X; \text{Bin}(T_\ell)) \stackrel{??}{\approx} I(X; T_\ell)$  **No!**



# What is going on here?

- Plots via binning-based estimator of  $I(X; T_\ell)$ , for  $X \sim \text{Unif}(\text{dataset})$

$\Rightarrow$  Plotted values are  $I(X; \text{Bin}(T_\ell)) \stackrel{??}{\approx} I(X; T_\ell)$  **No!**

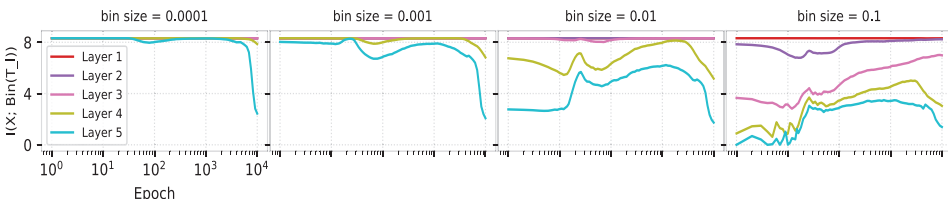


- Smaller bins  $\Rightarrow$  Closer to truth:  $I(X; T_\ell) = \ln(2^{12}) \approx 8.31$

# What is going on here?

- Plots via binning-based estimator of  $I(X; T_\ell)$ , for  $X \sim \text{Unif}(\text{dataset})$

$\implies$  Plotted values are  $I(X; \text{Bin}(T_\ell)) \stackrel{??}{\approx} I(X; T_\ell)$  **No!**

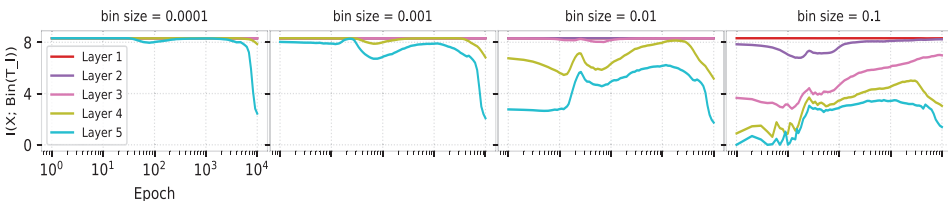


- Smaller bins  $\implies$  Closer to truth:  $I(X; T_\ell) = \ln(2^{12}) \approx 8.31$
- Binning introduces “noise” into estimator (not present in the DNN)

# What is going on here?

- Plots via binning-based estimator of  $I(X; T_\ell)$ , for  $X \sim \text{Unif}(\text{dataset})$

$\Rightarrow$  Plotted values are  $I(X; \text{Bin}(T_\ell)) \stackrel{??}{\approx} I(X; T_\ell)$  **No!**

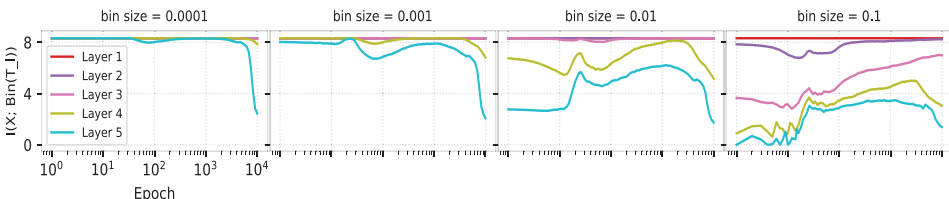


- Smaller bins  $\Rightarrow$  Closer to truth:  $I(X; T_\ell) = \ln(2^{12}) \approx 8.31$
- Binning introduces “noise” into estimator (not present in the DNN)
- Plots showing estimation errors

# What is going on here?

- Plots via binning-based estimator of  $I(X; T_\ell)$ , for  $X \sim \text{Unif}(\text{dataset})$

$\Rightarrow$  Plotted values are  $I(X; \text{Bin}(T_\ell)) \stackrel{??}{\approx} I(X; T_\ell)$  **No!**



- Smaller bins  $\Rightarrow$  Closer to truth:  $I(X; T_\ell) = \ln(2^{12}) \approx 8.31$
- Binning introduces “noise” into estimator (not present in the DNN)
- Plots showing estimation errors

**\* Real Problem:**  $I(X; T_\ell)$  is meaningless in det. DNNs

# Auxiliary Framework - Noisy Deep Neural Networks

**Modification:** Inject (small) Gaussian noise to neurons' output

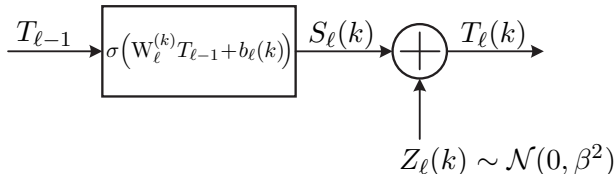
[G.-Berg-Greenewald-Melnyk-Nguyen-Kingsbury-Polyanskiy'18]

# Auxiliary Framework - Noisy Deep Neural Networks

**Modification:** Inject (small) Gaussian noise to neurons' output

[G.-Berg-Greenewald-Melnyk-Nguyen-Kingsbury-Polyanskiy'18]

- **Formally:**  $T_\ell = f_\ell(T_{\ell-1}) + Z_\ell$ , where  $Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

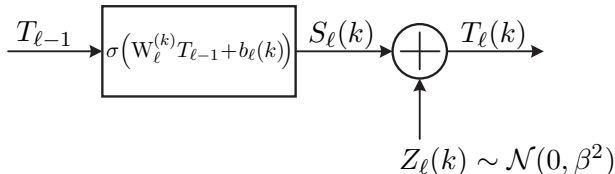


# Auxiliary Framework - Noisy Deep Neural Networks

**Modification:** Inject (small) Gaussian noise to neurons' output

[G.-Berg-Greenewald-Melnyk-Nguyen-Kingsbury-Polyanskiy'18]

- **Formally:**  $T_\ell = f_\ell(T_{\ell-1}) + Z_\ell$ , where  $Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



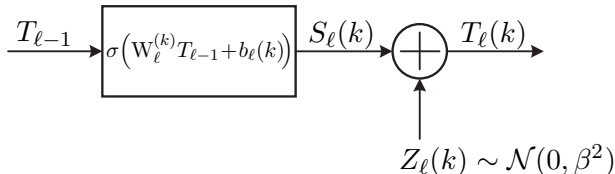
$\implies X \mapsto T_\ell$  is a **parametrized channel** that depends on DNN param.!

# Auxiliary Framework - Noisy Deep Neural Networks

**Modification:** Inject (small) Gaussian noise to neurons' output

[G.-Berg-Greenewald-Melnyk-Nguyen-Kingsbury-Polyanskiy'18]

- **Formally:**  $T_\ell = f_\ell(T_{\ell-1}) + Z_\ell$ , where  $Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



$\implies X \mapsto T_\ell$  is a **parametrized channel** that depends on DNN param.!

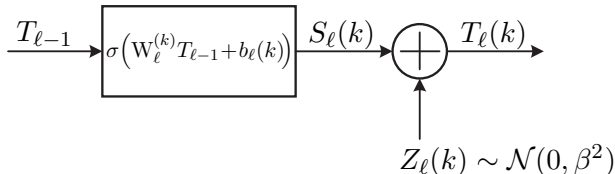
$\implies I(X; T_\ell)$  is a **function** of weights and biases!

# Auxiliary Framework - Noisy Deep Neural Networks

**Modification:** Inject (small) Gaussian noise to neurons' output

[G.-Berg-Greenewald-Melnyk-Nguyen-Kingsbury-Polyanskiy'18]

- **Formally:**  $T_\ell = f_\ell(T_{\ell-1}) + Z_\ell$ , where  $Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



$\implies X \mapsto T_\ell$  is a **parametrized channel** that depends on DNN param.!

$\implies I(X; T_\ell)$  is a **function** of weights and biases!

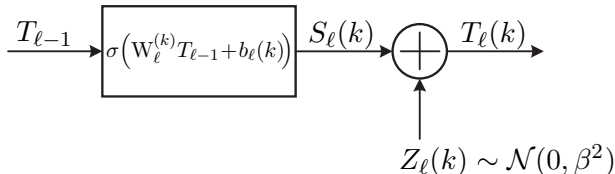
- **Operational Perspective:**

# Auxiliary Framework - Noisy Deep Neural Networks

**Modification:** Inject (small) Gaussian noise to neurons' output

[G.-Berg-Greenewald-Melnyk-Nguyen-Kingsbury-Polyanskiy'18]

- **Formally:**  $T_\ell = f_\ell(T_{\ell-1}) + Z_\ell$ , where  $Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



$\implies X \mapsto T_\ell$  is a **parametrized channel** that depends on DNN param.!

$\implies I(X; T_\ell)$  is a **function** of weights and biases!

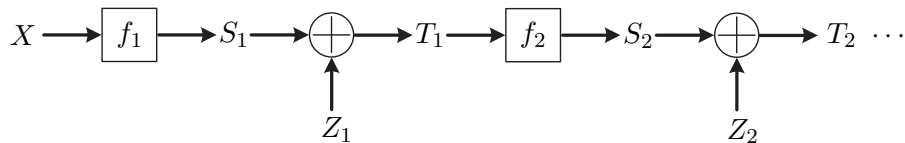
- **Operational Perspective:**

Performance & learned representations similar to det. DNNs ( $\beta \approx 10^{-1}$ )

# Mutual Information in Noisy DNNs

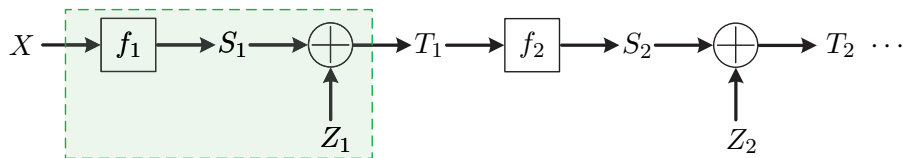
# Mutual Information Estimation in Noisy DNNs

## Noisy DNN:



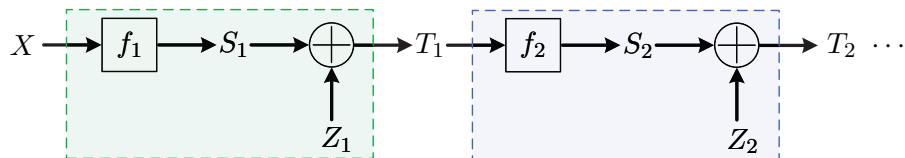
# Mutual Information Estimation in Noisy DNNs

## Noisy DNN:



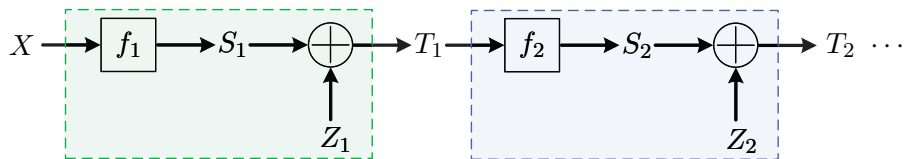
# Mutual Information Estimation in Noisy DNNs

## Noisy DNN:



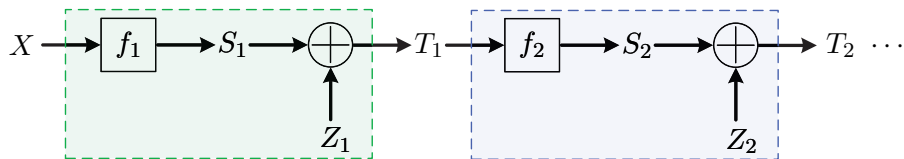
# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1})$



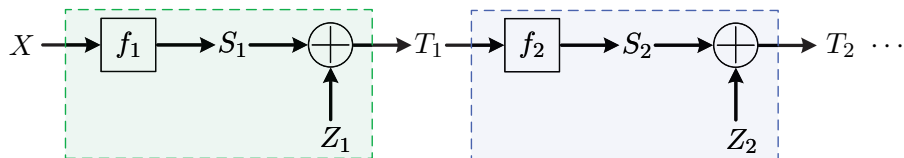
# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



# Mutual Information Estimation in Noisy DNNs

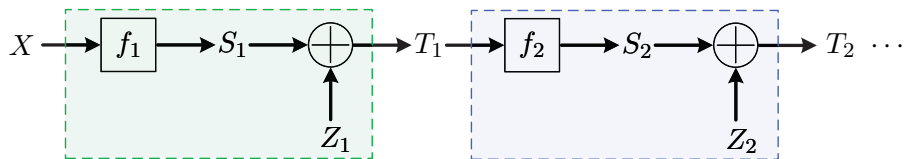
Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



- **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

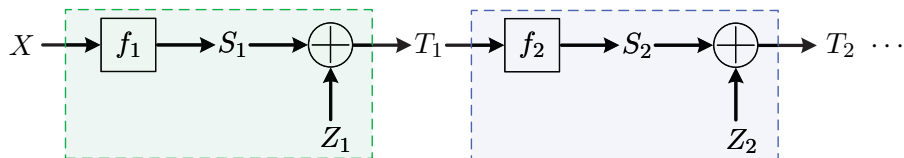


● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



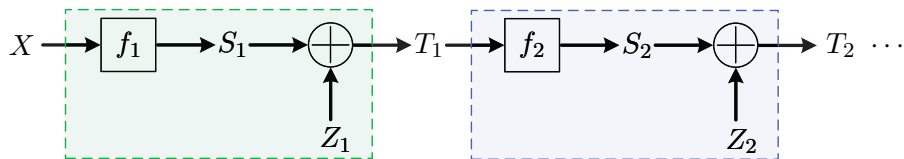
- Assume:  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

- Structure:  $S_\ell \perp Z_\ell \implies T_\ell = S_\ell + Z_\ell \sim P * \varphi$

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



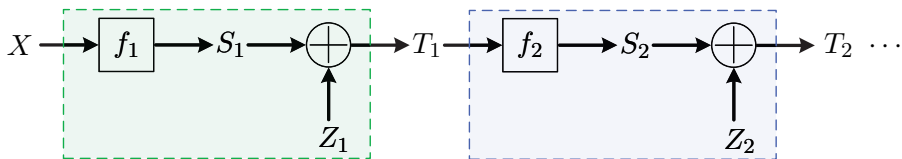
- Assume:  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

- Structure:**  $S_\ell \perp Z_\ell \implies T_\ell = \textcolor{red}{S}_\ell + Z_\ell \sim \textcolor{red}{P} * \varphi$

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



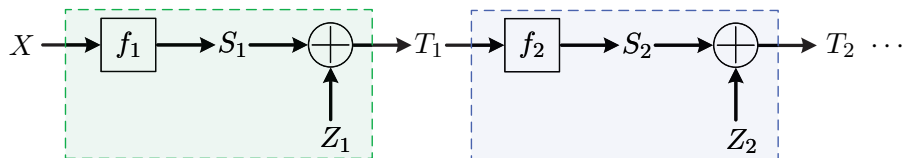
- Assume:  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

- Structure:  $S_\ell \perp Z_\ell \implies T_\ell = S_\ell + \mathbf{Z}_\ell \sim P * \varphi$

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

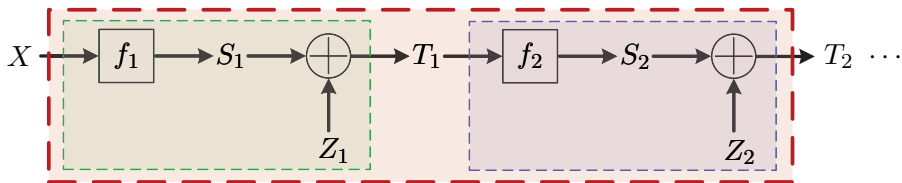
$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

● **Structure:**  $S_\ell \perp Z_\ell \implies T_\ell = S_\ell + Z_\ell \sim P * \varphi$

⊛ **Know** the distribution  $\varphi$  of  $Z_\ell$  (noise injected by design)

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

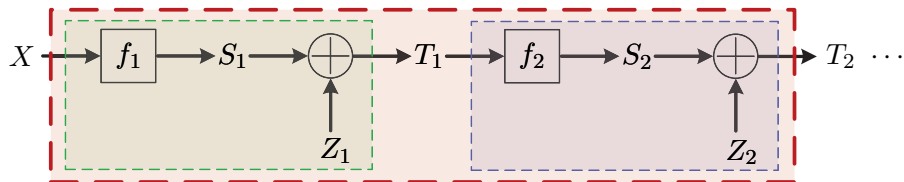
$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

● **Structure:**  $S_\ell \perp Z_\ell \implies T_\ell = S_\ell + Z_\ell \sim P * \varphi$

⊛ **Know** the distribution  $\varphi$  of  $Z_\ell$  (noise injected by design)

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



• **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

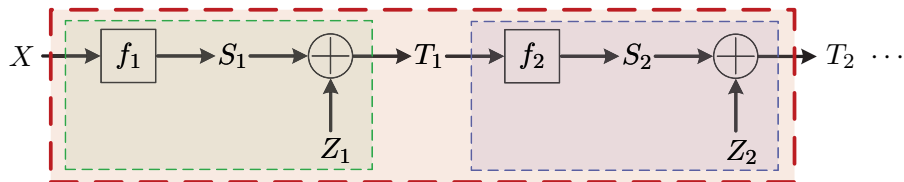
• **Structure:**  $S_\ell \perp Z_\ell \implies T_\ell = S_\ell + Z_\ell \sim P * \varphi$

⊛ **Know** the distribution  $\varphi$  of  $Z_\ell$  (noise injected by design)

⊛  $P_{S_\ell}$  and  $P_{S_\ell | X}$  are **extremely complicated** to compute/evaluate

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



• **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

• **Structure:**  $S_\ell \perp Z_\ell \implies T_\ell = S_\ell + Z_\ell \sim P * \varphi$

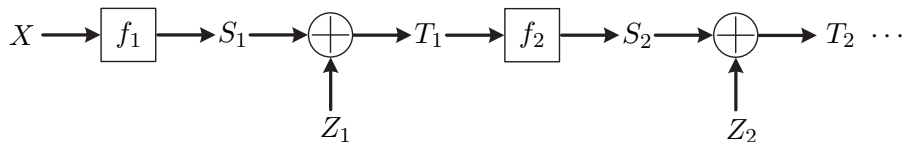
⊛ **Know** the distribution  $\varphi$  of  $Z_\ell$  (noise injected by design)

⊛  $P_{S_\ell}$  and  $P_{S_\ell | X}$  are **extremely complicated** to compute/evaluate

⊛ But both are **easily sampled** via the DNN forward pass

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

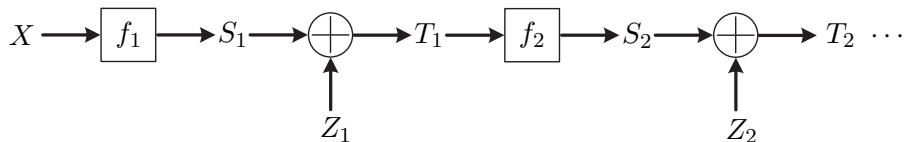


## Differential Entropy Estimation under Gaussian Convolutions

Estimate  $h(P * \varphi)$  based on  $n$  i.i.d. samples from  $P \in \mathcal{F}_d$  (nonparametric class) and knowledge of  $\varphi$  (PDF of  $\mathcal{N}(0, \beta^2 \mathbf{I}_d)$ ).

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



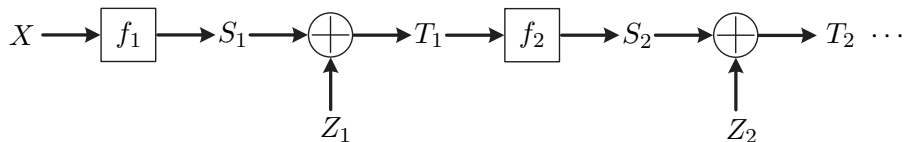
## Differential Entropy Estimation under Gaussian Convolutions

Estimate  $h(P * \varphi)$  based on  $n$  i.i.d. samples from  $P \in \mathcal{F}_d$  (nonparametric class) and knowledge of  $\varphi$  (PDF of  $\mathcal{N}(0, \beta^2 \mathbf{I}_d)$ ).

Estimation Results [G.-Greenewald-Polyanskiy'18]:

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



## Differential Entropy Estimation under Gaussian Convolutions

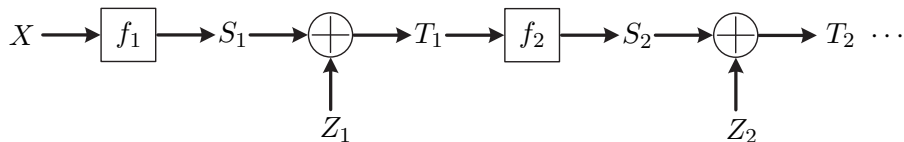
Estimate  $h(P * \varphi)$  based on  $n$  i.i.d. samples from  $P \in \mathcal{F}_d$  (nonparametric class) and knowledge of  $\varphi$  (PDF of  $\mathcal{N}(0, \beta^2 \mathbf{I}_d)$ ).

Estimation Results [G.-Greenewald-Polyanskiy'18]:

- Efficient & parallelizable estimator  $h(\hat{P}_n * \varphi) \approx h(P * \varphi)$

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



## Differential Entropy Estimation under Gaussian Convolutions

Estimate  $h(P * \varphi)$  based on  $n$  i.i.d. samples from  $P \in \mathcal{F}_d$  (nonparametric class) and knowledge of  $\varphi$  (PDF of  $\mathcal{N}(0, \beta^2 \mathbf{I}_d)$ ).

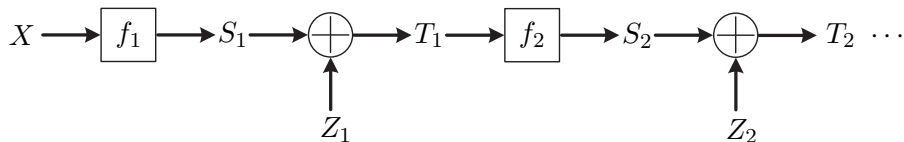
### Estimation Results [G.-Greenewald-Polyanskiy'18]:

- ▶ Efficient & parallelizable estimator  $h(\hat{P}_n * \varphi) \approx h(P * \varphi)$
- ▶ **Guarantees:** Estimation risk is  $O(1/\sqrt{n})$  (all constants explicit)\*

\* Exponentially large in  $d$  though constants, which is provably necessary.

# Mutual Information Estimation in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



## Differential Entropy Estimation under Gaussian Convolutions

Estimate  $h(P * \varphi)$  based on  $n$  i.i.d. samples from  $P \in \mathcal{F}_d$  (nonparametric class) and knowledge of  $\varphi$  (PDF of  $\mathcal{N}(0, \beta^2 \mathbf{I}_d)$ ).

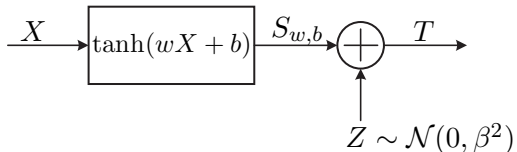
### Estimation Results [G.-Greenewald-Polyanskiy'18]:

- ▶ Efficient & parallelizable estimator  $h(\hat{P}_n * \varphi) \approx h(P * \varphi)$
- ▶ **Guarantees:** Estimation risk is  $O(1/\sqrt{n})$  (all constants explicit)\*
- ▶ **Faster Rate:** kNN/KDE est. via 'noisy' samples attain  $O\left(n^{-\frac{a}{b+d}}\right)$

Back to Noisy DNNs

# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

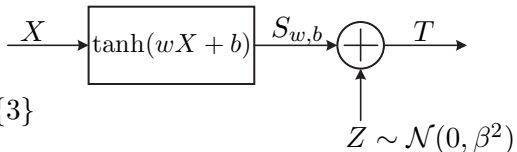


# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

- **Input:**  $X \sim \text{Unif}\{\pm 1, \pm 3\}$

$$\mathcal{X}_{y=-1} \triangleq \{-3, -1, 1\}, \quad \mathcal{X}_{y=1} \triangleq \{3\}$$

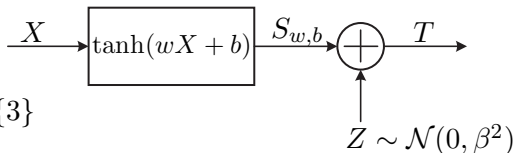
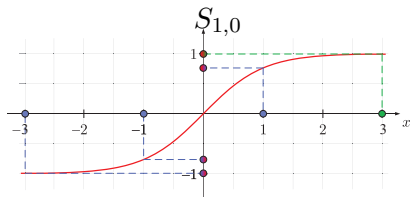


# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

● **Input:**  $X \sim \text{Unif}\{\pm 1, \pm 3\}$

$$\mathcal{X}_{y=-1} \triangleq \{-3, -1, 1\}, \quad \mathcal{X}_{y=1} \triangleq \{3\}$$

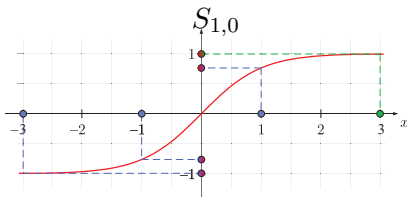
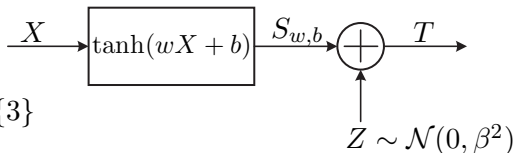


# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

- **Input:**  $X \sim \text{Unif}\{\pm 1, \pm 3\}$

$$\mathcal{X}_{y=-1} \triangleq \{-3, -1, 1\}, \quad \mathcal{X}_{y=1} \triangleq \{3\}$$



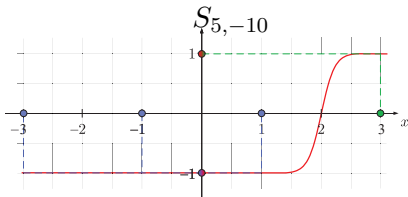
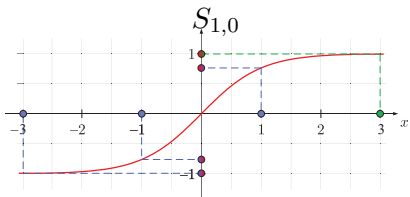
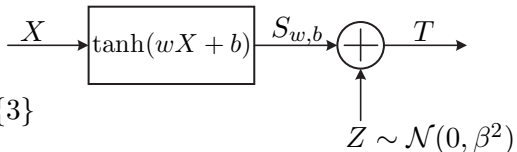
⊛ Center & sharpen transition (  $\iff$  increase  $w$  and keep  $b = -2w$  )

# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

- **Input:**  $X \sim \text{Unif}\{\pm 1, \pm 3\}$

$$\mathcal{X}_{y=-1} \triangleq \{-3, -1, 1\}, \quad \mathcal{X}_{y=1} \triangleq \{3\}$$

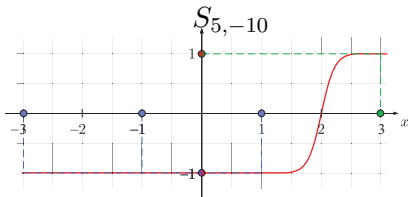
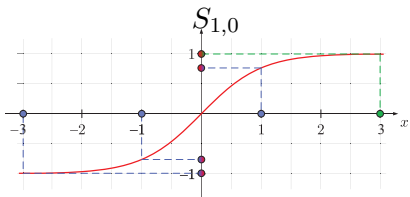
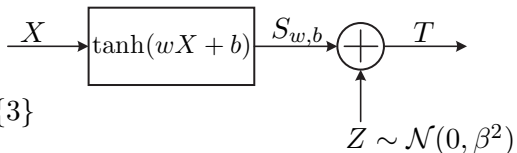


# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

- **Input:**  $X \sim \text{Unif}\{\pm 1, \pm 3\}$

$$\mathcal{X}_{y=-1} \triangleq \{-3, -1, 1\}, \quad \mathcal{X}_{y=1} \triangleq \{3\}$$



✓ Correct classification performance

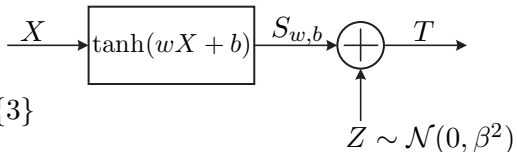
# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

- **Input:**  $X \sim \text{Unif}\{\pm 1, \pm 3\}$

$$\mathcal{X}_{y=-1} \triangleq \{-3, -1, 1\}, \quad \mathcal{X}_{y=1} \triangleq \{3\}$$

- **Mutual Information:**



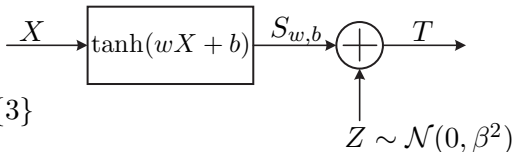
# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

- **Input:**  $X \sim \text{Unif}\{\pm 1, \pm 3\}$

$$\mathcal{X}_{y=-1} \triangleq \{-3, -1, 1\}, \quad \mathcal{X}_{y=1} \triangleq \{3\}$$

- **Mutual Information:**  $I(X; T) = I(S_{w,b}; S_{w,b} + Z)$



# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

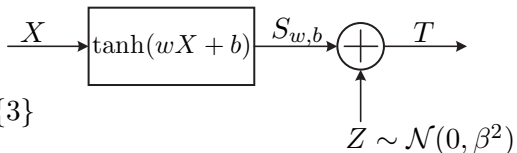
- **Input:**  $X \sim \text{Unif}\{\pm 1, \pm 3\}$

$$\mathcal{X}_{y=-1} \triangleq \{-3, -1, 1\}, \quad \mathcal{X}_{y=1} \triangleq \{3\}$$

- **Mutual Information:**  $I(X; T) = I(S_{w,b}; S_{w,b} + Z)$

$\implies I(X; T)$  is # bits (nats) transmittable over AWGN with symbols

$$\mathcal{S}_{w,b} \triangleq \{\tanh(-3w+b), \tanh(-w+b), \tanh(w+b), \tanh(3w+b)\}$$



# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

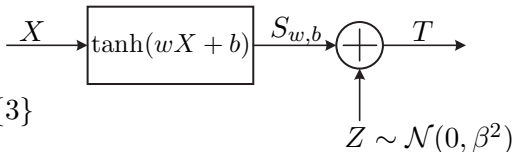
- **Input:**  $X \sim \text{Unif}\{\pm 1, \pm 3\}$

$$\mathcal{X}_{y=-1} \triangleq \{-3, -1, 1\}, \mathcal{X}_{y=1} \triangleq \{3\}$$

- **Mutual Information:**  $I(X; T) = I(S_{w,b}; S_{w,b} + Z)$

$\implies I(X; T)$  is # bits (nats) transmittable over AWGN with symbols

$$\mathcal{S}_{w,b} \triangleq \{\tanh(-3w+b), \tanh(-w+b), \tanh(w+b), \tanh(3w+b)\} \longrightarrow \{\pm 1\}$$



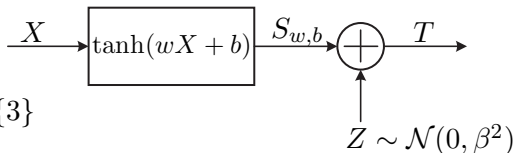
# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

- **Input:**  $X \sim \text{Unif}\{\pm 1, \pm 3\}$

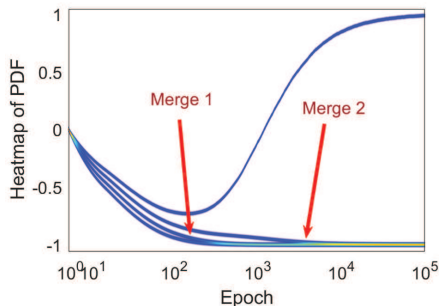
$$\mathcal{X}_{y=-1} \triangleq \{-3, -1, 1\}, \mathcal{X}_{y=1} \triangleq \{3\}$$

- **Mutual Information:**  $I(X; T) = I(S_{w,b}; S_{w,b} + Z)$



$\implies I(X; T)$  is # bits (nats) transmittable over AWGN with symbols

$$S_{w,b} \triangleq \{\tanh(-3w+b), \tanh(-w+b), \tanh(w+b), \tanh(3w+b)\} \longrightarrow \{\pm 1\}$$



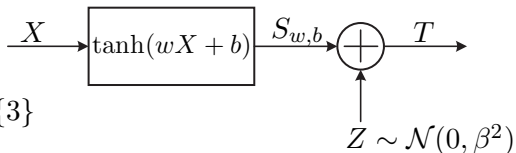
# $I(X; T_\ell)$ Dynamics - Illustrative Minimal Example

## Single Neuron Classification:

- **Input:**  $X \sim \text{Unif}\{\pm 1, \pm 3\}$

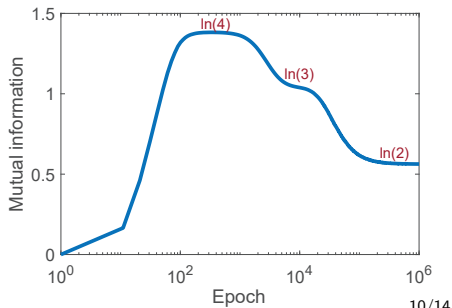
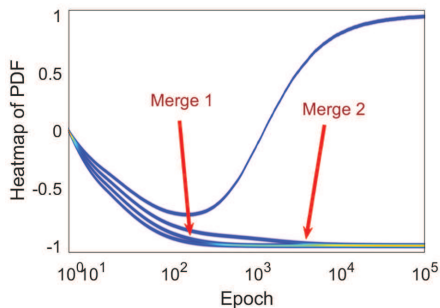
$$\mathcal{X}_{y=-1} \triangleq \{-3, -1, 1\}, \mathcal{X}_{y=1} \triangleq \{3\}$$

- **Mutual Information:**  $I(X; T) = I(S_{w,b}; S_{w,b} + Z)$



$\implies I(X; T)$  is # bits (nats) transmittable over AWGN with symbols

$$S_{w,b} \triangleq \{\tanh(-3w+b), \tanh(-w+b), \tanh(w+b), \tanh(3w+b)\} \longrightarrow \{\pm 1\}$$



# Clustering of Representations - Larger Networks

Noisy version of DNN from [Shwartz-Tishby'17]:

# Clustering of Representations - Larger Networks

## Noisy version of DNN from [Shwartz-Tishby'17]:

- **Binary Classification:** 12-bit input & 12–10–7–5–4–3–2 MLP arch.

# Clustering of Representations - Larger Networks

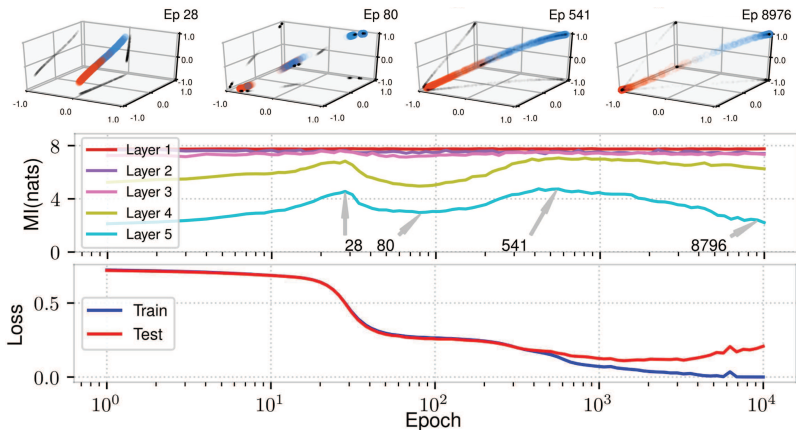
## Noisy version of DNN from [Shwartz-Tishby'17]:

- **Binary Classification:** 12-bit input & 12-10-7-5-4-3-2 MLP arch.
- **Noise std.:** Set to  $\beta = 0.01$

# Clustering of Representations - Larger Networks

## Noisy version of DNN from [Shwartz-Tishby'17]:

- **Binary Classification:** 12-bit input & 12-10-7-5-4-3-2 MLP arch.
- **Noise std.:** Set to  $\beta = 0.01$



# Clustering of Representations - Larger Networks

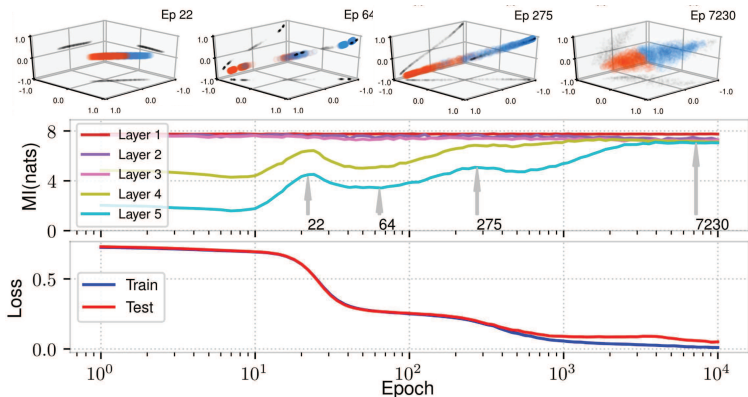
## Noisy version of DNN from [Shwartz-Tishby'17]:

- **Binary Classification:** 12-bit input & 12-10-7-5-4-3-2 MLP arch.
- **Noise std.:** Set to  $\beta = 0.01$

# Clustering of Representations - Larger Networks

## Noisy version of DNN from [Shwartz-Tishby'17]:

- **Binary Classification:** 12-bit input & 12-10-7-5-4-3-2 MLP arch.
- **Noise std.:** Set to  $\beta = 0.01$



⊛ weight orthonormality regularization

# Clustering of Representations - Larger Networks

## Noisy version of DNN from [Shwartz-Tishby'17]:

- **Binary Classification:** 12-bit input & 12-10-7-5-4-3-2 MLP arch.
- **Noise std.:** Set to  $\beta = 0.01$
- Verified in multiple additional experiments

# Clustering of Representations - Larger Networks

## Noisy version of DNN from [Shwartz-Tishby'17]:

- **Binary Classification:** 12-bit input & 12–10–7–5–4–3–2 MLP arch.
  - **Noise std.:** Set to  $\beta = 0.01$
  - Verified in multiple additional experiments
- ⇒ Compression of  $I(X; T_\ell)$  driven by clustering of representations

# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant

## Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering

## Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):

# Circling back to Deterministic DNNs

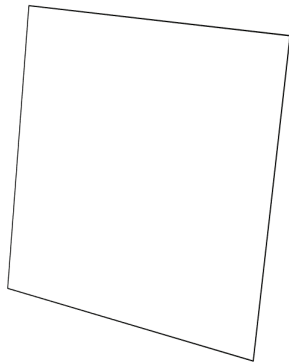
- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):
  - ▶ Scatter plots (up to 3D layers)

# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):
  - ▶ Scatter plots (up to 3D layers)
  - ▶ Binned entropy  $H(\text{Bin}(T_\ell))$

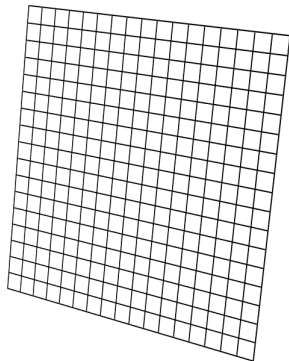
# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):
  - ▶ Scatter plots (up to 3D layers)
  - ▶ Binned entropy  $H(\text{Bin}(T_\ell))$



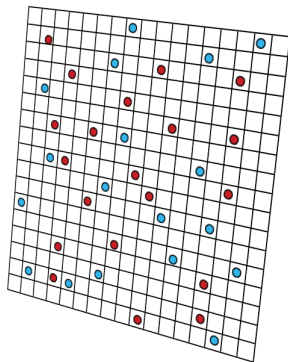
# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):
  - ▶ Scatter plots (up to 3D layers)
  - ▶ Binned entropy  $H(\text{Bin}(T_\ell))$



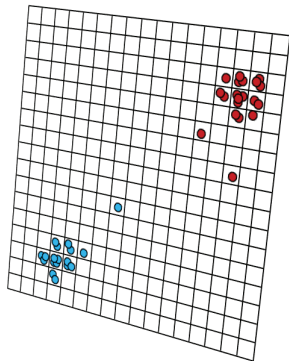
# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):
  - ▶ Scatter plots (up to 3D layers)
  - ▶ Binned entropy  $H(\text{Bin}(T_\ell)) \uparrow$



# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):
  - ▶ Scatter plots (up to 3D layers)
  - ▶ Binned entropy  $H(\text{Bin}(T_\ell)) \downarrow$



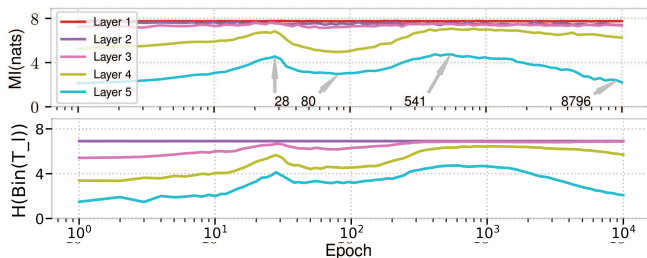
# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):
  - ▶ Scatter plots (up to 3D layers)
  - ▶ Binned entropy  $H(\text{Bin}(T_\ell))$
- **Noisy DNNs:**  $I(X; T_\ell)$  and  $H(\text{Bin}(T_\ell))$  highly correlated!\*

\* When bin size chosen  $\propto$  noise std.

# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):
  - ▶ Scatter plots (up to 3D layers)
  - ▶ Binned entropy  $H(\text{Bin}(T_\ell))$
- **Noisy DNNs:**  $I(X; T_\ell)$  and  $H(\text{Bin}(T_\ell))$  highly correlated!\*



★ When bin size chosen  $\propto$  noise std.

# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):
  - ▶ Scatter plots (up to 3D layers)
  - ▶ Binned entropy  $H(\text{Bin}(T_\ell))$
- **Noisy DNNs:**  $I(X; T_\ell)$  and  $H(\text{Bin}(T_\ell))$  highly correlated!\*
- **Det. DNNs:**  $H(\text{Bin}(T_\ell)) = I(X; \text{Bin}(T_\ell))$  compresses

# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):
  - ▶ Scatter plots (up to 3D layers)
  - ▶ Binned entropy  $H(\text{Bin}(T_\ell))$
- **Noisy DNNs:**  $I(X; T_\ell)$  and  $H(\text{Bin}(T_\ell))$  highly correlated!\*
- **Det. DNNs:**  $H(\text{Bin}(T_\ell)) = I(X; \text{Bin}(T_\ell))$  compresses
  - ✗ Incapable of accurately estimating MI values

# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
- Alternative measures for clustering (det. and noisy DNNs):
  - ▶ Scatter plots (up to 3D layers)
  - ▶ Binned entropy  $H(\text{Bin}(T_\ell))$
- **Noisy DNNs:**  $I(X; T_\ell)$  and  $H(\text{Bin}(T_\ell))$  highly correlated!\*
- **Det. DNNs:**  $H(\text{Bin}(T_\ell)) = I(X; \text{Bin}(T_\ell))$  compresses
  - ✗ Incapable of accurately estimating MI values
  - ✓ Does track clustering!

# Circling back to Deterministic DNNs

- $I(X; T_\ell)$  is constant  $\implies$  Doesn't measure clustering
  - Alternative measures for clustering (det. and noisy DNNs):
    - ▶ Scatter plots (up to 3D layers)
    - ▶ Binned entropy  $H(\text{Bin}(T_\ell))$
  - **Noisy DNNs:**  $I(X; T_\ell)$  and  $H(\text{Bin}(T_\ell))$  highly correlated!\*
  - **Det. DNNs:**  $H(\text{Bin}(T_\ell)) = I(X; \text{Bin}(T_\ell))$  compresses
    - ✗ Incapable of accurately estimating MI values
    - ✓ Does track clustering!
- $\implies$  Past works were not showing MI but clustering (via binned-MI)!

- **Reexamined Information Bottleneck Compression:**

- **Reexamined Information Bottleneck Compression:**
  - ▶  $I(X;T)$  fluctuations in det. DNNs are theoretically impossible

- **Reexamined Information Bottleneck Compression:**

- ▶  $I(X; T)$  fluctuations in det. DNNs are theoretically impossible
- ▶ Yet, past works presented (binned)  $I(X; T)$  dynamics during training

- **Reexamined Information Bottleneck Compression:**
  - ▶  $I(X; T)$  fluctuations in det. DNNs are theoretically impossible
  - ▶ Yet, past works presented (binned)  $I(X; T)$  dynamics during training
- **Noisy DNN Framework:** Studying IT quantities over DNNs

- **Reexamined Information Bottleneck Compression:**
  - ▶  $I(X; T)$  fluctuations in det. DNNs are theoretically impossible
  - ▶ Yet, past works presented (binned)  $I(X; T)$  dynamics during training
- **Noisy DNN Framework:** Studying IT quantities over DNNs
  - ▶ Developed estimator for accurate MI estimation over this framework

- **Reexamined Information Bottleneck Compression:**
  - ▶  $I(X;T)$  fluctuations in det. DNNs are theoretically impossible
  - ▶ Yet, past works presented (binned)  $I(X;T)$  dynamics during training
- **Noisy DNN Framework:** Studying IT quantities over DNNs
  - ▶ Developed estimator for accurate MI estimation over this framework
  - ▶ Clustering of the learned representations is the source of compression

- **Reexamined Information Bottleneck Compression:**
  - ▶  $I(X;T)$  fluctuations in det. DNNs are theoretically impossible
  - ▶ Yet, past works presented (binned)  $I(X;T)$  dynamics during training
- **Noisy DNN Framework:** Studying IT quantities over DNNs
  - ▶ Developed estimator for accurate MI estimation over this framework
  - ▶ Clustering of the learned representations is the source of compression
- **Clarify Past Observations of Compression:** in fact show clustering

- **Reexamined Information Bottleneck Compression:**
  - ▶  $I(X;T)$  fluctuations in det. DNNs are theoretically impossible
  - ▶ Yet, past works presented (binned)  $I(X;T)$  dynamics during training
- **Noisy DNN Framework:** Studying IT quantities over DNNs
  - ▶ Developed estimator for accurate MI estimation over this framework
  - ▶ Clustering of the learned representations is the source of compression
- **Clarify Past Observations of Compression:** in fact show clustering
  - ⇒ **Clustering** is the common phenomenon of interest!

# Future/Ongoing Clustering Inspired Research

- **Track Clustering in High-Dimensions:**

# Future/Ongoing Clustering Inspired Research

- **Track Clustering in High-Dimensions:**

- ▶ Lower-dimensional embedding

# Future/Ongoing Clustering Inspired Research

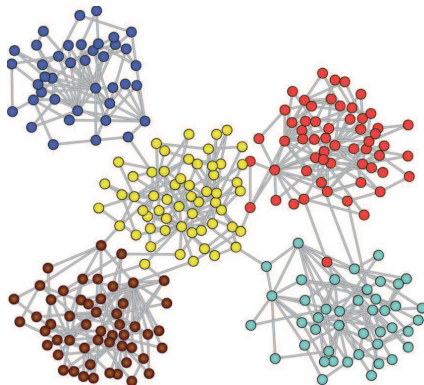
- **Track Clustering in High-Dimensions:**

- ▶ Lower-dimensional embedding
- ▶ Summarizing statistics

# Future/Ongoing Clustering Inspired Research

## ● Track Clustering in High-Dimensions:

- ▶ Lower-dimensional embedding
- ▶ Summarizing statistics
- ▶ Graph clusterability [Czumaj-Peng-Sohler'15]



# Future/Ongoing Clustering Inspired Research

- **Track Clustering in High-Dimensions:**

- ▶ Lower-dimensional embedding
- ▶ Summarizing statistics
- ▶ Graph clusterability [Czumaj-Peng-Sohler'15]

- **Role of Compression/Clustering:**

# Future/Ongoing Clustering Inspired Research

- **Track Clustering in High-Dimensions:**

- ▶ Lower-dimensional embedding
- ▶ Summarizing statistics
- ▶ Graph clusterability [Czumaj-Peng-Sohler'15]

- **Role of Compression/Clustering:**

- ▶ Is it necessary? Desirable?

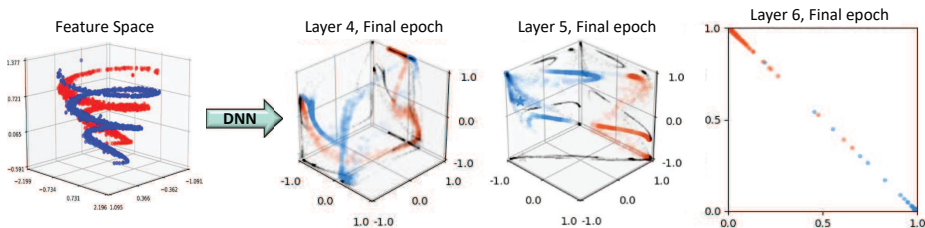
# Future/Ongoing Clustering Inspired Research

- **Track Clustering in High-Dimensions:**

- ▶ Lower-dimensional embedding
- ▶ Summarizing statistics
- ▶ Graph clusterability [Czumaj-Peng-Sohler'15]

- **Role of Compression/Clustering:**

- ▶ Is it necessary? Desirable?



# Future/Ongoing Clustering Inspired Research

- **Track Clustering in High-Dimensions:**

- ▶ Lower-dimensional embedding
- ▶ Summarizing statistics
- ▶ Graph clusterability [Czumaj-Peng-Sohler'15]

- **Role of Compression/Clustering:**

- ▶ Is it necessary? Desirable?
- ▶ Design tool for DNN architectures

# Future/Ongoing Clustering Inspired Research

- **Track Clustering in High-Dimensions:**

- ▶ Lower-dimensional embedding
- ▶ Summarizing statistics
- ▶ Graph clusterability [Czumaj-Peng-Sohler'15]

- **Role of Compression/Clustering:**

- ▶ Is it necessary? Desirable?
- ▶ Design tool for DNN architectures

- **Algorithmic Perspective:**

# Future/Ongoing Clustering Inspired Research

- **Track Clustering in High-Dimensions:**

- ▶ Lower-dimensional embedding
- ▶ Summarizing statistics
- ▶ Graph clusterability [Czumaj-Peng-Sohler'15]

- **Role of Compression/Clustering:**

- ▶ Is it necessary? Desirable?
- ▶ Design tool for DNN architectures

- **Algorithmic Perspective:**

- ▶ Better understanding of internal representation evolution & final state

# Future/Ongoing Clustering Inspired Research

## ● **Track Clustering in High-Dimensions:**

- ▶ Lower-dimensional embedding
- ▶ Summarizing statistics
- ▶ Graph clusterability [Czumaj-Peng-Sohler'15]

## ● **Role of Compression/Clustering:**

- ▶ Is it necessary? Desirable?
- ▶ Design tool for DNN architectures

## ● **Algorithmic Perspective:**

- ▶ Better understanding of internal representation evolution & final state
- ▶ Enhanced DNN training alg. (regularize intermediate layers wrt  $I(Y;T)$ )

# Future/Ongoing Clustering Inspired Research

## ● Track Clustering in High-Dimensions:

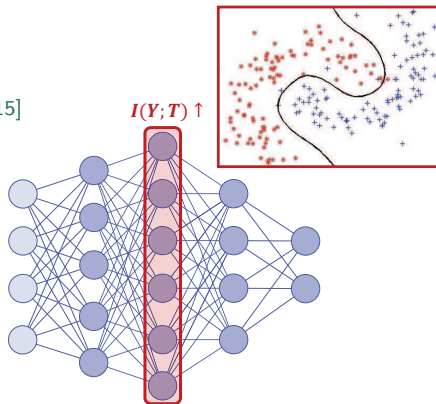
- ▶ Lower-dimensional embedding
- ▶ Summarizing statistics
- ▶ Graph clusterability [Czumaj-Peng-Sohler'15]

## ● Role of Compression/Clustering:

- ▶ Is it necessary? Desirable?
- ▶ Design tool for DNN architectures

## ● Algorithmic Perspective:

- ▶ Better understanding of internal representation evolution & final state
- ▶ Enhanced DNN training alg. (regularize intermediate layers wrt  $I(Y;T)$ )



# References

- [1] Z. Goldfeld, E. van den Berg, K. Greenewald, I. Melnyk, N. Nguyen, B. Kingsbury and Y. Polyanskiy, "Estimating information flow in neural networks," Arxiv preprint <https://arxiv.org/abs/1810.05728>, October 2018.
  
- [2] Z. Goldfeld, K. Greenewald and Y. Polyanskiy, "Estimating differential entropy under Gaussian convolutions," Submitted to the *IEEE Transactions on Information Theory*, October 2018.  
Arxiv: <https://arxiv.org/abs/1810.11589>
  
- [3] Z. Goldfeld, G. Bresler and Y. Polyanskiy, "Information storage in the stochastic Ising model," Submitted to the *IEEE Transactions on Information Theory*, May 2018.  
Arxiv: <https://arxiv.org/abs/1805.03027>

# Information Storage in Interacting Particle Systems

Motivation: Demand for high-capacity data storage devices

# Information Storage in Interacting Particle Systems

Motivation: Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit

# Information Storage in Interacting Particle Systems

Motivation: Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit
- **Challenges** in stabilizing stored data (interparticle interactions)

# Information Storage in Interacting Particle Systems

Motivation: Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit
- **Challenges** in stabilizing stored data (interparticle interactions)

Approach:

# Information Storage in Interacting Particle Systems

Motivation: Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit
- **Challenges** in stabilizing stored data (interparticle interactions)

Approach: View storage medium as **Interacting Particle System:**

# Information Storage in Interacting Particle Systems

Motivation: Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit
- **Challenges** in stabilizing stored data (interparticle interactions)

Approach: View storage medium as **Interacting Particle System:**

- Hard-drive topology

# Information Storage in Interacting Particle Systems

Motivation: Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit
- **Challenges** in stabilizing stored data (interparticle interactions)

Approach: View storage medium as **Interacting Particle System:**

- Hard-drive topology  $\implies$  Graph

# Information Storage in Interacting Particle Systems

Motivation: Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit
- **Challenges** in stabilizing stored data (interparticle interactions)

Approach: View storage medium as **Interacting Particle System:**

- Hard-drive topology  $\implies$  Graph
- Interactions

# Information Storage in Interacting Particle Systems

Motivation: Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit
- **Challenges** in stabilizing stored data (interparticle interactions)

Approach: View storage medium as **Interacting Particle System:**

- Hard-drive topology  $\implies$  Graph
- Interactions  $\implies$  Stochastic dynamics

# Information Storage in Interacting Particle Systems

Motivation: Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit
- **Challenges** in stabilizing stored data (interparticle interactions)

Approach: View storage medium as **Interacting Particle System:**

- Hard-drive topology  $\implies$  Graph
- Interactions  $\implies$  Stochastic dynamics

# Information Storage in Interacting Particle Systems

Motivation: Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit
- **Challenges** in stabilizing stored data (interparticle interactions)

Approach: View storage medium as **Interacting Particle System:**

- Hard-drive topology  $\implies$  Graph
- Interactions  $\implies$  Stochastic dynamics

Setup: **Design** initial configuration and **recover** data after  $t$  time

# Information Storage in Interacting Particle Systems

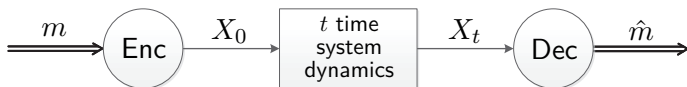
Motivation: Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit
- **Challenges** in stabilizing stored data (interparticle interactions)

Approach: View storage medium as **Interacting Particle System**:

- Hard-drive topology  $\implies$  Graph
- Interactions  $\implies$  Stochastic dynamics

Setup: **Design** initial configuration and **recover** data after  $t$  time



# Information Storage in Interacting Particle Systems

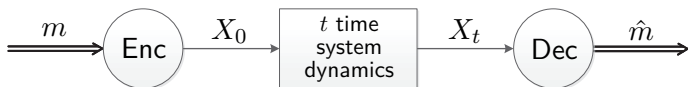
**Motivation:** Demand for high-capacity data storage devices

- **Emerging technologies** drastically shrink magnetic region per bit
- **Challenges** in stabilizing stored data (interparticle interactions)

**Approach:** View storage medium as **Interacting Particle System:**

- Hard-drive topology  $\implies$  Graph
- Interactions  $\implies$  Stochastic dynamics

**Setup:** **Design** initial configuration and **recover** data after  $t$  time



[G.-Bresler-Polyanskiy'18] Performance benchmarks & hard-drive designs

# Cooperative Communication Networks

Motivation: Modern networks are large, decentralized and ad hoc

# Cooperative Communication Networks

**Motivation:** Modern networks are large, decentralized and ad hoc



# Cooperative Communication Networks

Motivation: Modern networks are large, decentralized and ad hoc

- Most devices are **simple & low-complexity**

# Cooperative Communication Networks

Motivation: Modern networks are large, decentralized and ad hoc

- Most devices are **simple & low-complexity**
- Nodes constantly **join** and **leave** the network

# Cooperative Communication Networks

Motivation: Modern networks are large, decentralized and ad hoc

- Most devices are **simple & low-complexity**
- Nodes constantly **join** and **leave** the network
- **Cooperation** for long-distance communication (mitigate interference)

# Cooperative Communication Networks

Motivation: Modern networks are large, decentralized and ad hoc

- Most devices are **simple & low-complexity**
- Nodes constantly **join** and **leave** the network
- **Cooperation** for long-distance communication (mitigate interference)

Modus Operandi: Primitive building blocks  $\implies$  Entire network

# Cooperative Communication Networks

Motivation: Modern networks are large, decentralized and ad hoc

- Most devices are **simple & low-complexity**
- Nodes constantly **join** and **leave** the network
- **Cooperation** for long-distance communication (mitigate interference)

Modus Operandi: Primitive building blocks  $\implies$  Entire network

- P2P, feedback, relay, uplink, downlink channels

# Cooperative Communication Networks

Motivation: Modern networks are large, decentralized and ad hoc

- Most devices are **simple & low-complexity**
- Nodes constantly **join** and **leave** the network
- **Cooperation** for long-distance communication (mitigate interference)

Modus Operandi: Primitive building blocks  $\implies$  Entire network

- P2P, feedback, relay, uplink, downlink channels
- Public vs. private vs. confidential transmissions

# Cooperative Communication Networks

Motivation: Modern networks are large, decentralized and ad hoc

- Most devices are **simple & low-complexity**
- Nodes constantly **join** and **leave** the network
- **Cooperation** for long-distance communication (mitigate interference)

Modus Operandi: Primitive building blocks  $\implies$  Entire network

- P2P, feedback, relay, uplink, downlink channels
- Public vs. private vs. confidential transmissions

Next Steps: **Model & study** large-scale networks

# Cooperative Communication Networks

Motivation: Modern networks are large, decentralized and ad hoc

- Most devices are **simple & low-complexity**
- Nodes constantly **join** and **leave** the network
- **Cooperation** for long-distance communication (mitigate interference)

Modus Operandi: Primitive building blocks  $\implies$  Entire network

- P2P, feedback, relay, uplink, downlink channels
- Public vs. private vs. confidential transmissions

Next Steps: **Model & study** large-scale networks

- **Stochastic** connectivity patterns (random / small-world / free-scale)

# Cooperative Communication Networks

Motivation: Modern networks are large, decentralized and ad hoc

- Most devices are **simple & low-complexity**
- Nodes constantly **join** and **leave** the network
- **Cooperation** for long-distance communication (mitigate interference)

Modus Operandi: Primitive building blocks  $\implies$  Entire network

- P2P, feedback, relay, uplink, downlink channels
- Public vs. private vs. confidential transmissions

Next Steps: **Model & study** large-scale networks

- **Stochastic** connectivity patterns (random / small-world / free-scale)
- Local interactions **distilled** from primitive results

# Cooperative Communication Networks

Motivation: Modern networks are large, decentralized and ad hoc

- Most devices are **simple & low-complexity**
- Nodes constantly **join** and **leave** the network
- **Cooperation** for long-distance communication (mitigate interference)

Modus Operandi: Primitive building blocks  $\implies$  Entire network

- P2P, feedback, relay, uplink, downlink channels
- Public vs. private vs. confidential transmissions

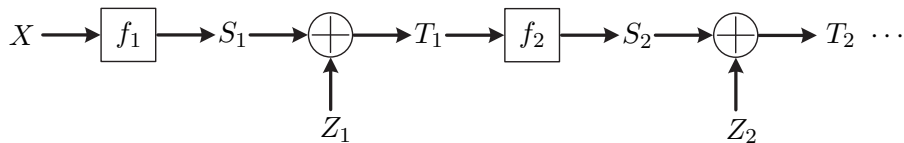
Next Steps: **Model & study** large-scale networks

- **Stochastic** connectivity patterns (random / small-world / free-scale)
- Local interactions **distilled** from primitive results

Q: Reliable (& secure) information passing protocols? Fundamental limits?

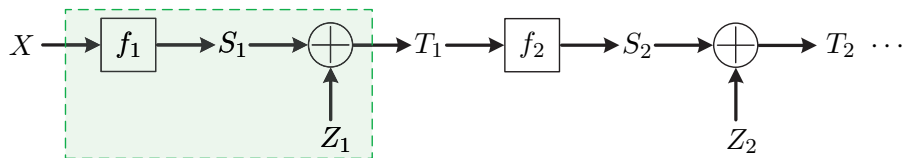
# Mutual Information in Noisy DNNs

## Noisy DNN:



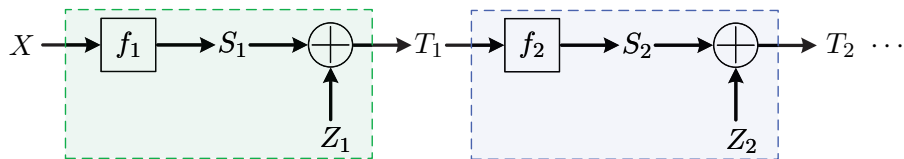
# Mutual Information in Noisy DNNs

## Noisy DNN:



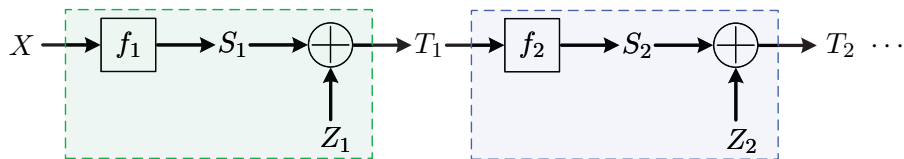
# Mutual Information in Noisy DNNs

## Noisy DNN:



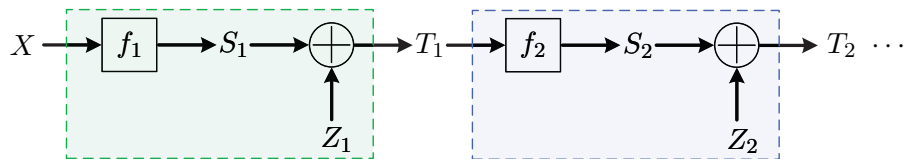
# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1})$



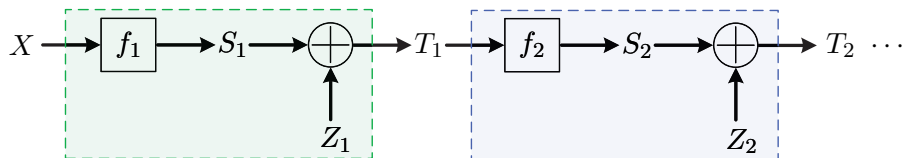
# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



# Mutual Information in Noisy DNNs

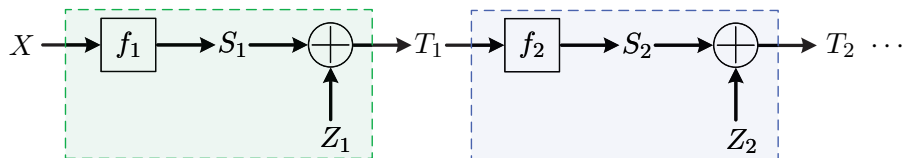
Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



- Assume:  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

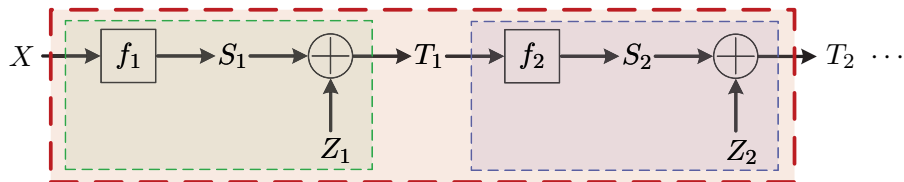


● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

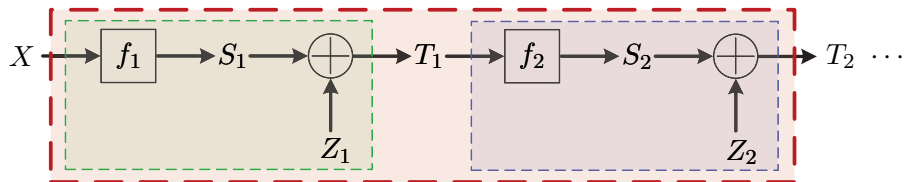


● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



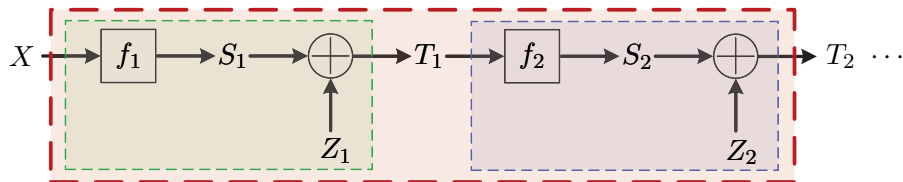
● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

⊛  $P_{T_\ell}$  and  $P_{T_\ell|X}$  are **extremely** complicated to compute/evaluate

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

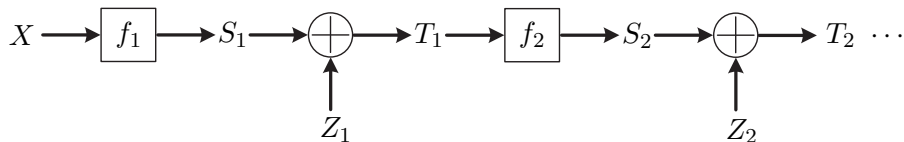
$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

⊗  $P_{T_\ell}$  and  $P_{T_\ell|X}$  are **extremely** complicated to compute/evaluate

⊗ But both are **easily** sampled via the DNN forward pass

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

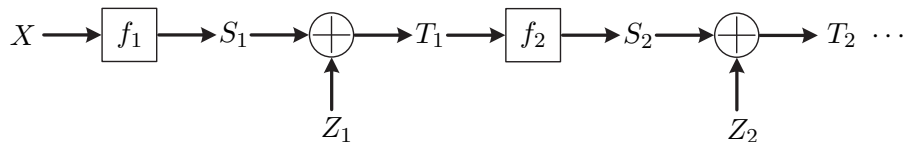
$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

⊗  $P_{T_\ell}$  and  $P_{T_\ell|X}$  are **extremely** complicated to compute/evaluate

⊗ But both are **easily** sampled via the DNN forward pass

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

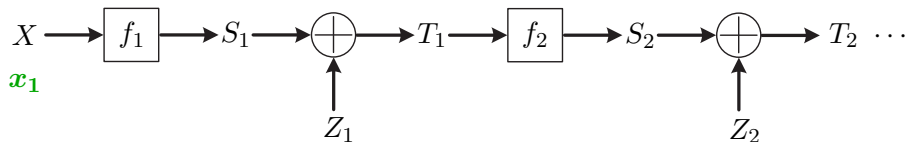
⊗  $P_{T_\ell}$  and  $P_{T_\ell | X}$  are **extremely** complicated to compute/evaluate

⊛ But both are **easily** sampled via the DNN forward pass

► **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

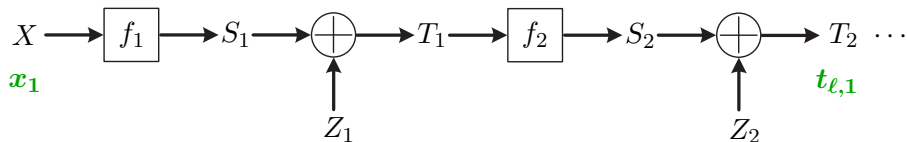
⊗  $P_{T_\ell}$  and  $P_{T_\ell | X}$  are **extremely** complicated to compute/evaluate

⊗ But both are **easily** sampled via the DNN forward pass

► **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

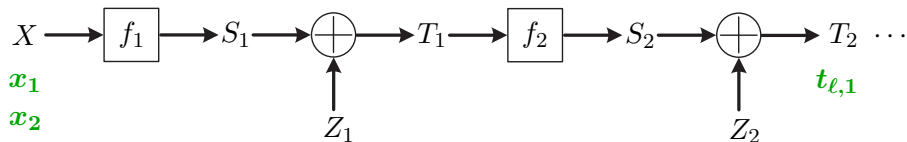
⊗  $P_{T_\ell}$  and  $P_{T_\ell | X}$  are **extremely** complicated to compute/evaluate

⊛ But both are **easily** sampled via the DNN forward pass

► **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

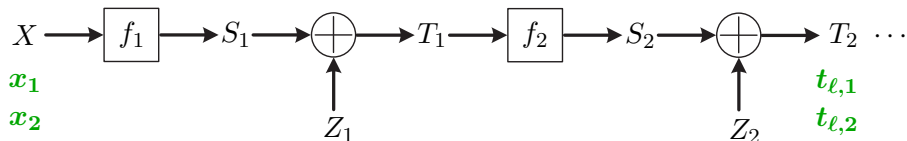
⊛  $P_{T_\ell}$  and  $P_{T_\ell|X}$  are **extremely** complicated to compute/evaluate

⊛ But both are **easily** sampled via the DNN forward pass

► **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

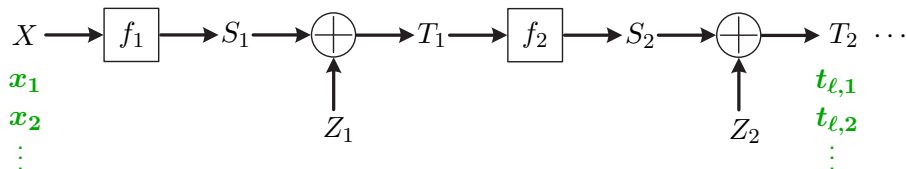
⊛  $P_{T_\ell}$  and  $P_{T_\ell|X}$  are **extremely** complicated to compute/evaluate

⊛ But both are **easily** sampled via the DNN forward pass

► **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

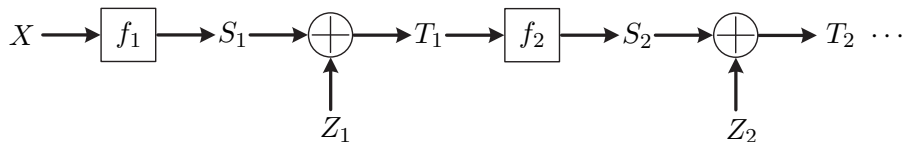
⊛  $P_{T_\ell}$  and  $P_{T_\ell | X}$  are **extremely** complicated to compute/evaluate

⊛ But both are **easily** sampled via the DNN forward pass

► **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

⊛  $P_{T_\ell}$  and  $P_{T_\ell|X}$  are **extremely** complicated to compute/evaluate

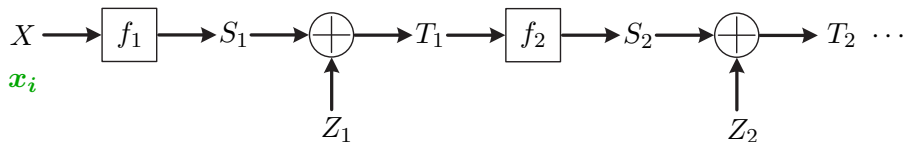
⊛ But both are **easily** sampled via the DNN forward pass

▶ **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

▶ **Sampling  $P_{T_\ell|X=x_i}$ :** Feed  $x_i$  multiples times & read  $T_\ell$  values

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

⊛  $P_{T_\ell}$  and  $P_{T_\ell | X}$  are **extremely** complicated to compute/evaluate

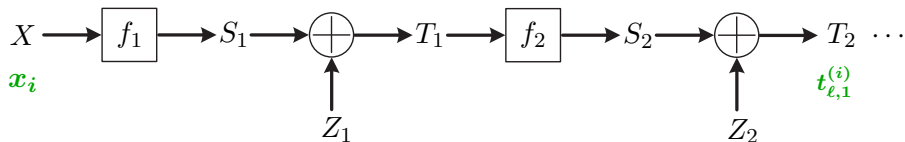
⊛ But both are **easily** sampled via the DNN forward pass

▶ **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

▶ **Sampling  $P_{T_\ell | X=x_i}$ :** Feed  $x_i$  multiples times & read  $T_\ell$  values

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

⊛  $P_{T_\ell}$  and  $P_{T_\ell | X}$  are **extremely** complicated to compute/evaluate

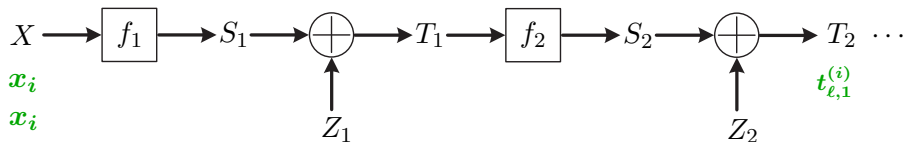
⊛ But both are **easily** sampled via the DNN forward pass

▶ **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

▶ **Sampling  $P_{T_\ell | X=x_i}$ :** Feed  $x_i$  multiples times & read  $T_\ell$  values

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

⊛  $P_{T_\ell}$  and  $P_{T_\ell | X}$  are **extremely** complicated to compute/evaluate

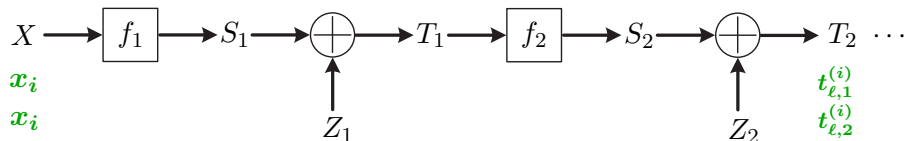
⊛ But both are **easily** sampled via the DNN forward pass

▶ **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

▶ **Sampling  $P_{T_\ell | X=x_i}$ :** Feed  $x_i$  multiples times & read  $T_\ell$  values

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

⊛  $P_{T_\ell}$  and  $P_{T_\ell | X}$  are **extremely** complicated to compute/evaluate

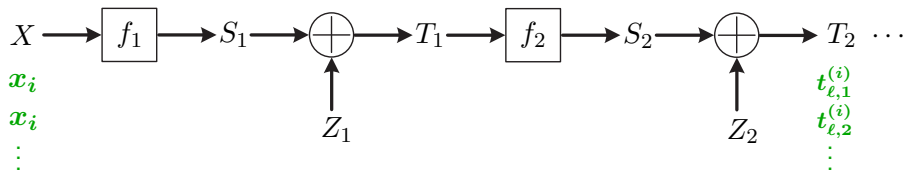
⊛ But both are **easily** sampled via the DNN forward pass

▶ **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

▶ **Sampling  $P_{T_\ell | X=x_i}$ :** Feed  $x_i$  multiples times & read  $T_\ell$  values

# Mutual Information in Noisy DNNs

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



● **Assume:**  $X \sim \text{Unif}(\mathcal{X})$ , where  $\mathcal{X} \triangleq \{x_i\}_{i=1}^m$  is empirical dataset

$\implies$  **Mutual Information:**  $I(X; T_\ell) = h(T_\ell) - \frac{1}{m} \sum_{i=1}^m h(T_\ell | X = x_i)$

⊛  $P_{T_\ell}$  and  $P_{T_\ell|X}$  are **extremely** complicated to compute/evaluate

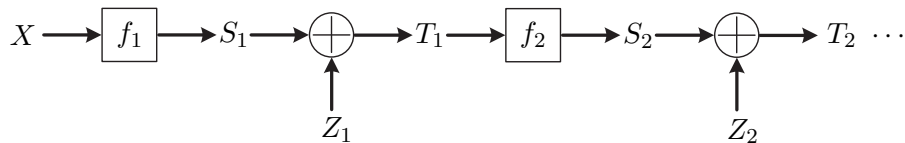
⊛ But both are **easily** sampled via the DNN forward pass

▶ **Sampling  $P_{T_\ell}$ :** Feed randomly chosen  $x_i$ 's & read  $T_\ell$  values

▶ **Sampling  $P_{T_\ell|X=x_i}$ :** Feed  $x_i$  multiples times & read  $T_\ell$  values

# General-Purpose Differential Entropy Estimators

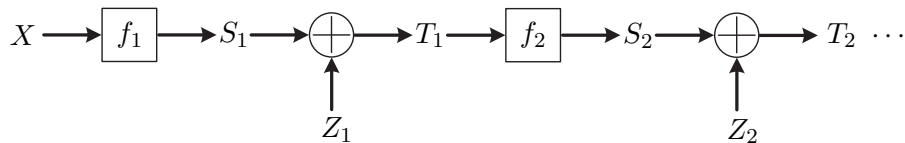
Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



$\implies$  Estimate  $I(X; T_\ell)$  from samples via **general-purpose  $h(P)$  est.:**

# General-Purpose Differential Entropy Estimators

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

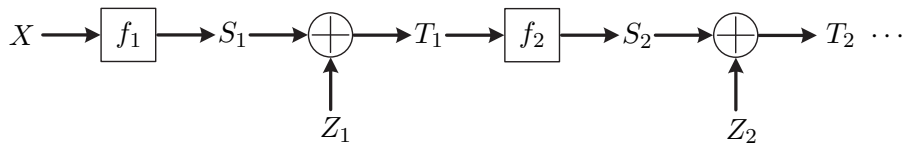


$\implies$  Estimate  $I(X; T_\ell)$  from samples via **general-purpose  $h(P)$  est.:**

- Most results assume lower bounded density

# General-Purpose Differential Entropy Estimators

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

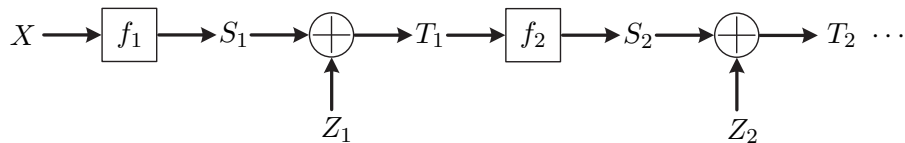


$\implies$  Estimate  $I(X; T_\ell)$  from samples via **general-purpose  $h(P)$  est.:**

- Most results assume lower bounded density  $\implies$  **Inapplicable**

# General-Purpose Differential Entropy Estimators

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

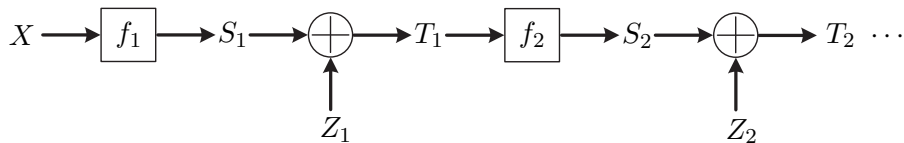


$\implies$  Estimate  $I(X; T_\ell)$  from samples via **general-purpose  $h(P)$  est.:**

- Most results assume lower bounded density  $\implies$  **Inapplicable**
- 2 Works Drop Assumption:**

# General-Purpose Differential Entropy Estimators

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



$\implies$  Estimate  $I(X; T_\ell)$  from samples via **general-purpose  $h(P)$  est.:**

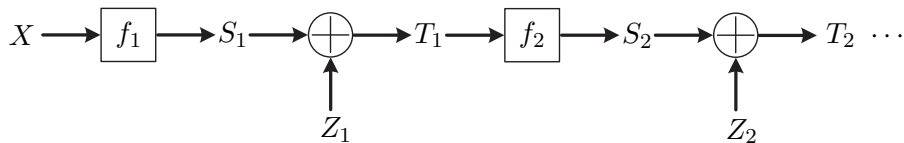
- Most results assume lower bounded density  $\implies$  **Inapplicable**

- 2 Works Drop Assumption:**

  - 1 KDE + Best poly. approximation [Han-Jiao-Weissman-Wu'17]

# General-Purpose Differential Entropy Estimators

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



$\implies$  Estimate  $I(X; T_\ell)$  from samples via **general-purpose  $h(P)$  est.:**

- Most results assume lower bounded density  $\implies$  **Inapplicable**

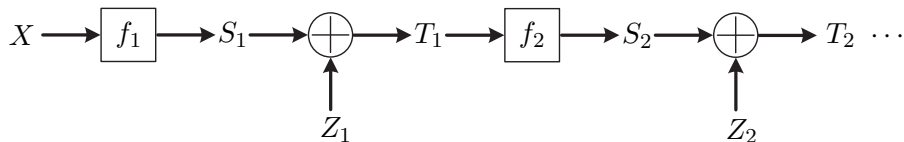
- 2 Works Drop Assumption:**

- 1 KDE + Best poly. approximation [Han-Jiao-Weissman-Wu'17]

- 2 Kozachenko-Leonenko (kNN) estimator [Jiao-Gao-Han'17]

# General-Purpose Differential Entropy Estimators

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

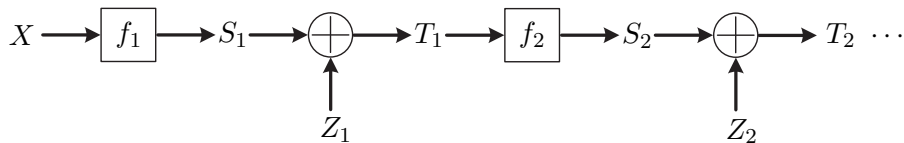


$\implies$  Estimate  $I(X; T_\ell)$  from samples via **general-purpose  $h(P)$  est.:**

- Most results assume lower bounded density  $\implies$  **Inapplicable**
- **2 Works Drop Assumption:**
  - 1 KDE + Best poly. approximation [Han-Jiao-Weissman-Wu'17]
  - 2 Kozachenko-Leonenko (kNN) estimator [Jiao-Gao-Han'17]
- **Assume:**  $\text{supp} = [0, 1]^d$  & Periodic BC &  $s \in (0, 2]$

# General-Purpose Differential Entropy Estimators

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



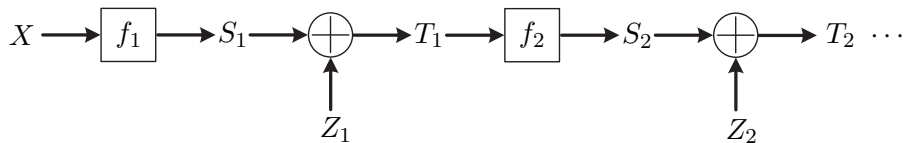
$\implies$  Estimate  $I(X; T_\ell)$  from samples via **general-purpose  $h(P)$  est.:**

- Most results assume lower bounded density  $\implies$  **Inapplicable**
- **2 Works Drop Assumption:**
  - ① KDE + Best poly. approximation [Han-Jiao-Weissman-Wu'17]
  - ② Kozachenko-Leonenko (kNN) estimator [Jiao-Gao-Han'17]
- **Assume:**  $\text{supp} = [0, 1]^d$  & Periodic BC &  $s \in (0, 2] \implies$  **Inapplicable\***

\* Except sub-Gaussian result from [Han-Jiao-Weissman-Wu'17]

# General-Purpose Differential Entropy Estimators

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

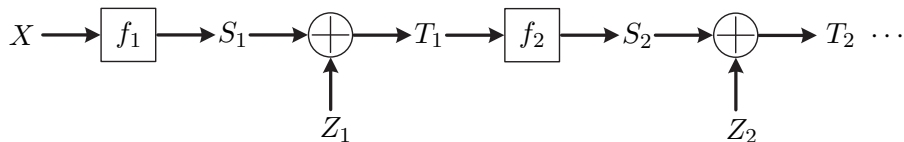


$\implies$  Estimate  $I(X; T_\ell)$  from samples via **general-purpose  $h(P)$  est.:**

- Most results assume lower bounded density  $\implies$  **Inapplicable**
- **2 Works Drop Assumption:**
  - 1 KDE + Best poly. approximation [Han-Jiao-Weissman-Wu'17]
  - 2 Kozachenko-Leonenko (kNN) estimator [Jiao-Gao-Han'17]
- **Assume:**  $\text{supp} = [0, 1]^d$  & Periodic BC &  $s \in (0, 2] \implies$  **Inapplicable\***
- **Rate:**  $\text{Risk} \leq O\left(n^{-\frac{\alpha s}{\beta s + d}}\right), \quad \text{w/ } \alpha, \beta \in \mathbb{N}, s \text{ smoothness, } d \text{ dimension}$

# Exploit Structure - Ad Hoc Estimation

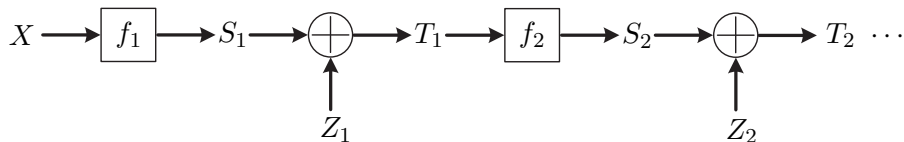
Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



⊛ **Exploit structure:** We know  $T_\ell = S_\ell + Z_\ell \sim P * \varphi$  and:

# Exploit Structure - Ad Hoc Estimation

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

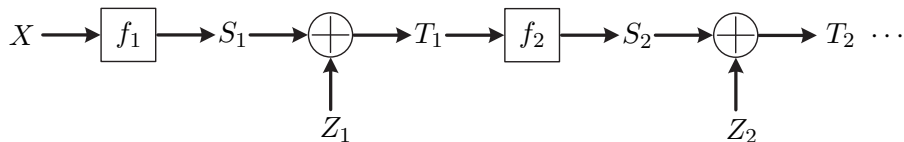


⊛ **Exploit structure:** We know  $T_\ell = S_\ell + Z_\ell \sim P * \varphi$  and:

- **Genie1:** Sample  $P = P_{S_\ell}$  and  $P = P_{S_\ell|X=x_i}$  (sample  $T_{\ell-1}$  & apply  $f_\ell$ )

# Exploit Structure - Ad Hoc Estimation

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$

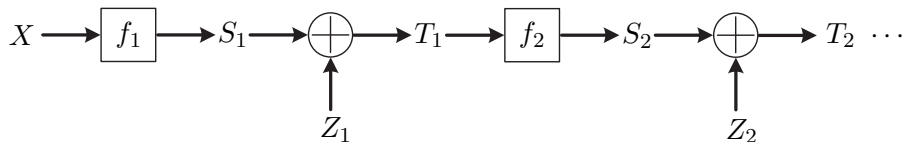


⊛ **Exploit structure:** We know  $T_\ell = S_\ell + Z_\ell \sim P * \varphi$  and:

- **Genie1:** Sample  $P = P_{S_\ell}$  and  $P = P_{S_\ell|X=x_i}$  (sample  $T_{\ell-1}$  & apply  $f_\ell$ )
- **Genie2:** Know the distribution  $\varphi$  of  $Z_\ell$  (noise injected by design)

# Exploit Structure - Ad Hoc Estimation

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



⊛ **Exploit structure:** We know  $T_\ell = S_\ell + Z_\ell \sim P * \varphi$  and:

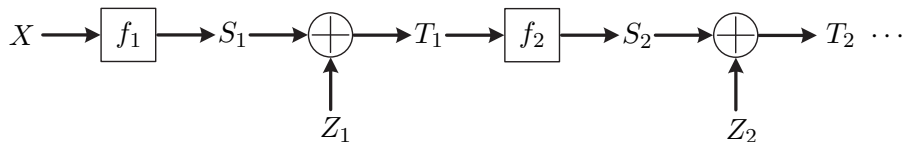
- **Genie1:** Sample  $P = P_{S_\ell}$  and  $P = P_{S_\ell|X=x_i}$  (sample  $T_{\ell-1}$  & apply  $f_\ell$ )
- **Genie2:** Know the distribution  $\varphi$  of  $Z_\ell$  (noise injected by design)

## Differential Entropy Estimation under Gaussian Convolutions

Estimate  $h(P * \varphi)$  based on  $n$  i.i.d. samples from  $P \in \mathcal{F}_d$  (nonparametric class) and knowledge of  $\varphi$  (PDF of  $\mathcal{N}(0, \beta^2 \mathbf{I}_d)$ ).

# Exploit Structure - Ad Hoc Estimation

Noisy DNN:  $S_\ell \triangleq f_\ell(T_{\ell-1}) \implies T_\ell = S_\ell + Z_\ell, \quad Z_\ell \sim \mathcal{N}(0, \beta^2 \mathbf{I})$



⊛ **Exploit structure:** We know  $T_\ell = S_\ell + Z_\ell \sim P * \varphi$  and:

- **Genie1:** Sample  $P = P_{S_\ell}$  and  $P = P_{S_\ell|X=x_i}$  (sample  $T_{\ell-1}$  & apply  $f_\ell$ )
- **Genie2:** Know the distribution  $\varphi$  of  $Z_\ell$  (noise injected by design)

## Differential Entropy Estimation under Gaussian Convolutions

Estimate  $h(P * \varphi)$  based on  $n$  i.i.d. samples from  $P \in \mathcal{F}_d$  (nonparametric class) and knowledge of  $\varphi$  (PDF of  $\mathcal{N}(0, \beta^2 \mathbf{I}_d)$ ).

Nonparametric Class: Depends on DNN architecture (nonlinearities)

# The Sample Propagation Estimator

**Abs. Error Minimax Risk:**  $S^n$  are  $n$  i.i.d. samples from  $P$ , define

$$\mathcal{R}_d^*(n, \beta) \triangleq \inf_{\hat{h}} \sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}(S^n, \beta) \right|$$

# The Sample Propagation Estimator

**Abs. Error Minimax Risk:**  $S^n$  are  $n$  i.i.d. samples from  $P$ , define

$$\mathcal{R}_d^*(n, \beta) \triangleq \inf_{\hat{h}} \sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}(S^n, \beta) \right|$$

⊛ **Curse of Dimensionality:** Sample complexity exponential in  $d$

# The Sample Propagation Estimator

**Abs. Error Minimax Risk:**  $S^n$  are  $n$  i.i.d. samples from  $P$ , define

$$\mathcal{R}_d^*(n, \beta) \triangleq \inf_{\hat{h}} \sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}(S^n, \beta) \right|$$

**\* Curse of Dimensionality:** Sample complexity exponential in  $d$

**'Sample Propagation' Estimator:** Empirical distribution  $\hat{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{S_i}$

$$\hat{h}_{\text{SP}}(S^n, \beta) \triangleq h(\hat{P}_n * \varphi)$$

# The Sample Propagation Estimator

**Abs. Error Minimax Risk:**  $S^n$  are  $n$  i.i.d. samples from  $P$ , define

$$\mathcal{R}_d^*(n, \beta) \triangleq \inf_{\hat{h}} \sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}(S^n, \beta) \right|$$

**⊛ Curse of Dimensionality:** Sample complexity exponential in  $d$

**'Sample Propagation' Estimator:** Empirical distribution  $\hat{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{S_i}$

$$\hat{h}_{\text{SP}}(S^n, \beta) \triangleq h(\hat{P}_n * \varphi)$$

**Comments:**

# The Sample Propagation Estimator

**Abs. Error Minimax Risk:**  $S^n$  are  $n$  i.i.d. samples from  $P$ , define

$$\mathcal{R}_d^*(n, \beta) \triangleq \inf_{\hat{h}} \sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}(S^n, \beta) \right|$$

**⊛ Curse of Dimensionality:** Sample complexity exponential in  $d$

**'Sample Propagation' Estimator:** Empirical distribution  $\hat{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{S_i}$

$$\hat{h}_{\text{SP}}(S^n, \beta) \triangleq h(\hat{P}_n * \varphi)$$

**Comments:**

- **Plug-in:**  $\hat{h}_{\text{SP}}$  is just plug-in est. for the functional  $T_\varphi(P) \triangleq h(P * \varphi)$

# The Sample Propagation Estimator

**Abs. Error Minimax Risk:**  $S^n$  are  $n$  i.i.d. samples from  $P$ , define

$$\mathcal{R}_d^*(n, \beta) \triangleq \inf_{\hat{h}} \sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}(S^n, \beta) \right|$$

⊛ **Curse of Dimensionality:** Sample complexity exponential in  $d$

**'Sample Propagation' Estimator:** Empirical distribution  $\hat{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{S_i}$

$$\hat{h}_{\text{SP}}(S^n, \beta) \triangleq h(\hat{P}_n * \varphi)$$

## Comments:

- **Plug-in:**  $\hat{h}_{\text{SP}}$  is just plug-in est. for the functional  $T_\varphi(P) \triangleq h(P * \varphi)$
- **Mixture:**  $\hat{h}_{\text{SP}}$  is the diff. entropy of a **known** Gaussian mixture

# The Sample Propagation Estimator

**Abs. Error Minimax Risk:**  $S^n$  are  $n$  i.i.d. samples from  $P$ , define

$$\mathcal{R}_d^*(n, \beta) \triangleq \inf_{\hat{h}} \sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}(S^n, \beta) \right|$$

**\* Curse of Dimensionality:** Sample complexity exponential in  $d$

**'Sample Propagation' Estimator:** Empirical distribution  $\hat{P}_n = \frac{1}{n} \sum_{i=1}^n \delta_{S_i}$

$$\hat{h}_{\text{SP}}(S^n, \beta) \triangleq h(\hat{P}_n * \varphi)$$

## Comments:

- **Plug-in:**  $\hat{h}_{\text{SP}}$  is just plug-in est. for the functional  $T_\varphi(P) \triangleq h(P * \varphi)$
- **Mixture:**  $\hat{h}_{\text{SP}}$  is the diff. entropy of a **known** Gaussian mixture
- **Computing:** Can be efficiently computed via MC integration

# The Sample Propagation Estimator - Convergence

## Theorem (ZG-Greenewald-Polyanskiy '18)

For  $\mathcal{F}_d \triangleq \{P | \text{supp}(P) \subseteq [-1, 1]^d\}$  and any  $\beta > 0$  and  $d \geq 1$ , we have

$$\sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}_{\text{SP}}(S^n, \beta) \right| \leq O_\beta \left( \frac{(\log n)^{d/4}}{\sqrt{n}} \right).$$

# The Sample Propagation Estimator - Convergence

## Theorem (ZG-Greenewald-Polyanskiy '18)

For  $\mathcal{F}_d \triangleq \{P | \text{supp}(P) \subseteq [-1, 1]^d\}$  and any  $\beta > 0$  and  $d \geq 1$ , we have

$$\sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}_{\text{SP}}(S^n, \beta) \right| \leq O_\beta \left( \frac{(\log n)^{d/4}}{\sqrt{n}} \right).$$

$$\begin{aligned} & \sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}_{\text{SP}}(S^n, \beta) \right| \\ & \leq \frac{1}{2(4\pi\beta^2)^{\frac{d}{4}}} \log \left( \frac{n \left( 2 + 2\beta \sqrt{(2 + \epsilon) \log n} \right)^d}{(\pi\beta^2)^{\frac{d}{2}}} \right) \left( 2 + 2\beta \sqrt{(2 + \epsilon) \log n} \right)^{\frac{d}{2}} \frac{1}{\sqrt{n}} \\ & \quad + \left( c_{\beta,d}^2 + \frac{2c_{\beta,d}d(1 + \beta^2)}{\beta^2} + \frac{8d(d + 2\beta^4 + d\beta^4)}{\beta^4} \right) \frac{2}{n} \end{aligned}$$

where  $c_{\beta,d} \triangleq \frac{d}{2} \log(2\pi\beta^2) + \frac{d}{\beta^2}$ .

# The Sample Propagation Estimator - Convergence

## Theorem (ZG-Greenewald-Polyanskiy '18)

For  $\mathcal{F}_d \triangleq \{P | \text{supp}(P) \subseteq [-1, 1]^d\}$  and any  $\beta > 0$  and  $d \geq 1$ , we have

$$\sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}_{\text{SP}}(S^n, \beta) \right| \leq O_\beta \left( \frac{(\log n)^{d/4}}{\sqrt{n}} \right).$$

Pf. Technique:

# The Sample Propagation Estimator - Convergence

## Theorem (ZG-Greenewald-Polyanskiy '18)

For  $\mathcal{F}_d \triangleq \{P | \text{supp}(P) \subseteq [-1, 1]^d\}$  and any  $\beta > 0$  and  $d \geq 1$ , we have

$$\sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}_{\text{SP}}(S^n, \beta) \right| \leq O_\beta \left( \frac{(\log n)^{d/4}}{\sqrt{n}} \right).$$

**Pf. Technique:** Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

# The Sample Propagation Estimator - Convergence

## Theorem (ZG-Greenewald-Polyanskiy '18)

For  $\mathcal{F}_d \triangleq \{P | \text{supp}(P) \subseteq [-1, 1]^d\}$  and any  $\beta > 0$  and  $d \geq 1$ , we have

$$\sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}_{\text{SP}}(S^n, \beta) \right| \leq O_\beta \left( \frac{(\log n)^{d/4}}{\sqrt{n}} \right).$$

**Pf. Technique:** Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

- **Inside  $\mathcal{R}$ :** Modulus of cont. & Convex analysis & Functional opt.

# The Sample Propagation Estimator - Convergence

## Theorem (ZG-Greenewald-Polyanskiy '18)

For  $\mathcal{F}_d \triangleq \{P | \text{supp}(P) \subseteq [-1, 1]^d\}$  and any  $\beta > 0$  and  $d \geq 1$ , we have

$$\sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}_{\text{SP}}(S^n, \beta) \right| \leq O_\beta \left( \frac{(\log n)^{d/4}}{\sqrt{n}} \right).$$

Pf. Technique: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

- **Inside  $\mathcal{R}$ :** Modulus of cont. & Convex analysis & Functional opt.
- **Outside  $\mathcal{R}$ :** Chi-squared distribution tail bounds

# The Sample Propagation Estimator - Convergence

## Theorem (ZG-Greenewald-Polyanskiy '18)

For  $\mathcal{F}_d \triangleq \{P | \text{supp}(P) \subseteq [-1, 1]^d\}$  and any  $\beta > 0$  and  $d \geq 1$ , we have

$$\sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}_{\text{SP}}(S^n, \beta) \right| \leq O_\beta \left( \frac{(\log n)^{d/4}}{\sqrt{n}} \right).$$

**Pf. Technique:** Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

- **Inside  $\mathcal{R}$ :** Modulus of cont. & Convex analysis & Functional opt.
- **Outside  $\mathcal{R}$ :** Chi-squared distribution tail bounds

**Comments:**

# The Sample Propagation Estimator - Convergence

## Theorem (ZG-Greenewald-Polyanskiy '18)

For  $\mathcal{F}_d \triangleq \{P | \text{supp}(P) \subseteq [-1, 1]^d\}$  and any  $\beta > 0$  and  $d \geq 1$ , we have

$$\sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}_{\text{SP}}(S^n, \beta) \right| \leq O_\beta \left( \frac{(\log n)^{d/4}}{\sqrt{n}} \right).$$

**Pf. Technique:** Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

- **Inside  $\mathcal{R}$ :** Modulus of cont. & Convex analysis & Functional opt.
- **Outside  $\mathcal{R}$ :** Chi-squared distribution tail bounds

## Comments:

- **Faster rate** than  $O\left(n^{-\frac{\alpha s}{\beta s + d}}\right)$  for kNN/KDE est. via 'noisy' samples

# The Sample Propagation Estimator - Convergence

## Theorem (ZG-Greenewald-Polyanskiy '18)

For  $\mathcal{F}_d \triangleq \{P | \text{supp}(P) \subseteq [-1, 1]^d\}$  and any  $\beta > 0$  and  $d \geq 1$ , we have

$$\sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}_{\text{SP}}(S^n, \beta) \right| \leq O_\beta \left( \frac{(\log n)^{d/4}}{\sqrt{n}} \right).$$

Pf. Technique: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

- **Inside  $\mathcal{R}$ :** Modulus of cont. & Convex analysis & Functional opt.
- **Outside  $\mathcal{R}$ :** Chi-squared distribution tail bounds

## Comments:

- **Faster rate** than  $O\left(n^{-\frac{\alpha s}{\beta s + d}}\right)$  for kNN/KDE est. via 'noisy' samples
- **Explicit expression** enables **concrete error bounds** in simulations

# The Sample Propagation Estimator - Convergence

## Theorem (ZG-Greenewald-Polyanskiy '18)

For  $\mathcal{F}_d \triangleq \{P | \text{supp}(P) \subseteq [-1, 1]^d\}$  and any  $\beta > 0$  and  $d \geq 1$ , we have

$$\sup_{P \in \mathcal{F}_d} \mathbb{E}_{S^n} \left| h(P * \varphi) - \hat{h}_{\text{SP}}(S^n, \beta) \right| \leq O_\beta \left( \frac{(\log n)^{d/4}}{\sqrt{n}} \right).$$

Pf. Technique: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

- **Inside  $\mathcal{R}$ :** Modulus of cont. & Convex analysis & Functional opt.
- **Outside  $\mathcal{R}$ :** Chi-squared distribution tail bounds

## Comments:

- **Faster rate** than  $O\left(n^{-\frac{\alpha s}{\beta s + d}}\right)$  for kNN/KDE est. via 'noisy' samples
- **Explicit expression** enables **concrete error bounds** in simulations
- **Extension:**  $P$  with sub-Gaussian marginals (ReLU + Weight regular.)

# The Sample Propagation Estimator - Proof Ideas

Strategy: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

# The Sample Propagation Estimator - Proof Ideas

Strategy: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

- **Restricted Entropy:**  $h_{\mathcal{R}}(p) \triangleq \mathbb{E}[-\log p(X) \mathbb{1}_{\{X \in \mathcal{R}\}}]$

# The Sample Propagation Estimator - Proof Ideas

Strategy: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

• **Restricted Entropy:**  $h_{\mathcal{R}}(p) \triangleq \mathbb{E}[-\log p(X) \mathbb{1}_{\{X \in \mathcal{R}\}}]$

$$\sup \mathbb{E} |h(P * \varphi) - h(\hat{P}_n * \varphi)| \leq \sup \mathbb{E} |h_{\mathcal{R}}(P * \varphi) - h_{\mathcal{R}}(\hat{P}_n * \varphi)| + 2 \sup |h_{\mathcal{R}^c}(P * \varphi)|$$

# The Sample Propagation Estimator - Proof Ideas

Strategy: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

• **Restricted Entropy:**  $h_{\mathcal{R}}(p) \triangleq \mathbb{E}[-\log p(X) \mathbb{1}_{\{X \in \mathcal{R}\}}]$

$$\sup \mathbb{E} |h(P * \varphi) - h(\hat{P}_n * \varphi)| \leq \sup \mathbb{E} |h_{\mathcal{R}}(P * \varphi) - h_{\mathcal{R}}(\hat{P}_n * \varphi)| + 2 \sup |h_{\mathcal{R}^c}(P * \varphi)|$$

• **Inside  $\mathcal{R}$ :**  $\blacktriangleright -t \log t$  modulus of cont. for  $x \mapsto x \log x$  & Jensen's ineq.

# The Sample Propagation Estimator - Proof Ideas

Strategy: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

• **Restricted Entropy:**  $h_{\mathcal{R}}(p) \triangleq \mathbb{E}[-\log p(X) \mathbb{1}_{\{X \in \mathcal{R}\}}]$

$$\sup \mathbb{E} |h(P * \varphi) - h(\hat{P}_n * \varphi)| \leq \sup \mathbb{E} |h_{\mathcal{R}}(P * \varphi) - h_{\mathcal{R}}(\hat{P}_n * \varphi)| + 2 \sup |h_{\mathcal{R}^c}(P * \varphi)|$$

• **Inside  $\mathcal{R}$ :**  $\blacktriangleright -t \log t$  modulus of cont. for  $x \mapsto x \log x$  & Jensen's ineq.

$$\implies \text{Focus on analyzing } \mathbb{E} \left| (P * \varphi)(x) - (\hat{P}_n * \varphi)(x) \right|$$

# The Sample Propagation Estimator - Proof Ideas

Strategy: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

• **Restricted Entropy:**  $h_{\mathcal{R}}(p) \triangleq \mathbb{E}[-\log p(X) \mathbb{1}_{\{X \in \mathcal{R}\}}]$

$$\sup \mathbb{E} |h(P * \varphi) - h(\hat{P}_n * \varphi)| \leq \sup \mathbb{E} |h_{\mathcal{R}}(P * \varphi) - h_{\mathcal{R}}(\hat{P}_n * \varphi)| + 2 \sup |h_{\mathcal{R}^c}(P * \varphi)|$$

• **Inside  $\mathcal{R}$ :** ▶  $-t \log t$  modulus of cont. for  $x \mapsto x \log x$  & Jensen's ineq.

⇒ Focus on analyzing  $\mathbb{E} |(P * \varphi)(x) - (\hat{P}_n * \varphi)(x)|$

▶ Bias & variance analysis

# The Sample Propagation Estimator - Proof Ideas

Strategy: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

• **Restricted Entropy:**  $h_{\mathcal{R}}(p) \triangleq \mathbb{E}[-\log p(X) \mathbb{1}_{\{X \in \mathcal{R}\}}]$

$$\sup \mathbb{E} |h(P * \varphi) - h(\hat{P}_n * \varphi)| \leq \sup \mathbb{E} |h_{\mathcal{R}}(P * \varphi) - h_{\mathcal{R}}(\hat{P}_n * \varphi)| + 2 \sup |h_{\mathcal{R}^c}(P * \varphi)|$$

• **Inside  $\mathcal{R}$ :** ►  $-t \log t$  modulus of cont. for  $x \mapsto x \log x$  & Jensen's ineq.

$$\implies \text{Focus on analyzing } \mathbb{E} \left| (P * \varphi)(x) - (\hat{P}_n * \varphi)(x) \right|$$

► Bias & variance analysis

$$\implies \mathbb{E} \left| (P * \varphi)(x) - (\hat{P}_n * \varphi)(x) \right| \leq c_1 \sqrt{\frac{(P * \tilde{\varphi})(x)}{n}}, \quad \tilde{\varphi} = \mathcal{N}\left(0, \frac{\beta^2}{2} \mathbf{I}\right)$$

# The Sample Propagation Estimator - Proof Ideas

Strategy: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

• **Restricted Entropy:**  $h_{\mathcal{R}}(p) \triangleq \mathbb{E}[-\log p(X) \mathbb{1}_{\{X \in \mathcal{R}\}}]$

$$\sup \mathbb{E} |h(P * \varphi) - h(\hat{P}_n * \varphi)| \leq \sup \mathbb{E} |h_{\mathcal{R}}(P * \varphi) - h_{\mathcal{R}}(\hat{P}_n * \varphi)| + 2 \sup |h_{\mathcal{R}^c}(P * \varphi)|$$

• **Inside  $\mathcal{R}$ :** ▶  $-t \log t$  modulus of cont. for  $x \mapsto x \log x$  & Jensen's ineq.

$$\implies \text{Focus on analyzing } \mathbb{E} \left| (P * \varphi)(x) - (\hat{P}_n * \varphi)(x) \right|$$

▶ Bias & variance analysis

$$\implies \mathbb{E} \left| (P * \varphi)(x) - (\hat{P}_n * \varphi)(x) \right| \leq c_1 \sqrt{\frac{(P * \tilde{\varphi})(x)}{n}}, \quad \tilde{\varphi} = \mathcal{N}\left(0, \frac{\beta^2}{2} \mathbf{I}\right)$$

▶ Plug back in & Convex analysis

# The Sample Propagation Estimator - Proof Ideas

Strategy: Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

• **Restricted Entropy:**  $h_{\mathcal{R}}(p) \triangleq \mathbb{E}[-\log p(X) \mathbb{1}_{\{X \in \mathcal{R}\}}]$

$$\sup \mathbb{E} |h(P * \varphi) - h(\hat{P}_n * \varphi)| \leq \sup \mathbb{E} |h_{\mathcal{R}}(P * \varphi) - h_{\mathcal{R}}(\hat{P}_n * \varphi)| + 2 \sup |h_{\mathcal{R}^c}(P * \varphi)|$$

• **Inside  $\mathcal{R}$ :** ▶  $-t \log t$  modulus of cont. for  $x \mapsto x \log x$  & Jensen's ineq.

$$\implies \text{Focus on analyzing } \mathbb{E} \left| (P * \varphi)(x) - (\hat{P}_n * \varphi)(x) \right|$$

▶ Bias & variance analysis

$$\implies \mathbb{E} \left| (P * \varphi)(x) - (\hat{P}_n * \varphi)(x) \right| \leq c_1 \sqrt{\frac{(P * \tilde{\varphi})(x)}{n}}, \quad \tilde{\varphi} = \mathcal{N}\left(0, \frac{\beta^2}{2} \mathbf{I}\right)$$

▶ Plug back in & Convex analysis

$$\implies \sup \mathbb{E} |h_{\mathcal{R}}(P * \varphi) - h_{\mathcal{R}}(\hat{P}_n * \varphi)| \leq c_2 \log \left( \frac{n \lambda(\mathcal{R})}{c_3} \right) \sqrt{\frac{\lambda(\mathcal{R})}{n}}$$

# The Sample Propagation Estimator - Proof Ideas

**Strategy:** Split analysis to  $\mathcal{R} \triangleq [-1, 1]^d + \mathcal{B}(0, \sqrt{c \log n})$  and  $\mathcal{R}^c$

• **Restricted Entropy:**  $h_{\mathcal{R}}(p) \triangleq \mathbb{E}[-\log p(X) \mathbb{1}_{\{X \in \mathcal{R}\}}]$

$$\sup \mathbb{E} |h(P * \varphi) - h(\hat{P}_n * \varphi)| \leq \sup \mathbb{E} |h_{\mathcal{R}}(P * \varphi) - h_{\mathcal{R}}(\hat{P}_n * \varphi)| + 2 \sup |h_{\mathcal{R}^c}(P * \varphi)|$$

• **Inside  $\mathcal{R}$ :** ▶  $-t \log t$  modulus of cont. for  $x \mapsto x \log x$  & Jensen's ineq.

$$\implies \text{Focus on analyzing } \mathbb{E} \left| (P * \varphi)(x) - (\hat{P}_n * \varphi)(x) \right|$$

▶ Bias & variance analysis

$$\implies \mathbb{E} \left| (P * \varphi)(x) - (\hat{P}_n * \varphi)(x) \right| \leq c_1 \sqrt{\frac{(P * \tilde{\varphi})(x)}{n}}, \quad \tilde{\varphi} = \mathcal{N}\left(0, \frac{\beta^2}{2} \mathbf{I}\right)$$

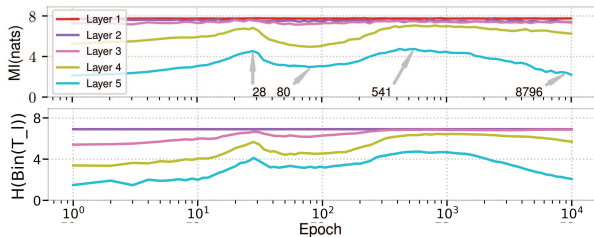
▶ Plug back in & Convex analysis

$$\implies \sup \mathbb{E} |h_{\mathcal{R}}(P * \varphi) - h_{\mathcal{R}}(\hat{P}_n * \varphi)| \leq c_2 \log \left( \frac{n \lambda(\mathcal{R})}{c_3} \right) \sqrt{\frac{\lambda(\mathcal{R})}{n}}$$

• **Outside  $\mathcal{R}$ :**  $O\left(\frac{1}{n}\right)$  decay via Chi-squared distribution tail bounds

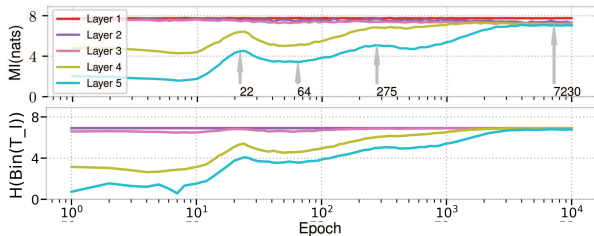
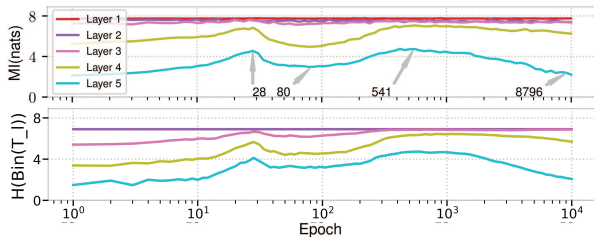
# Binning vs True Mutual Information

## Comparing to Previously Shown MI Plots:



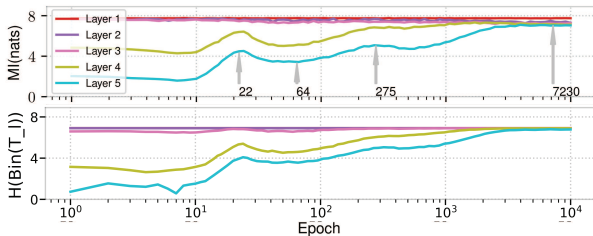
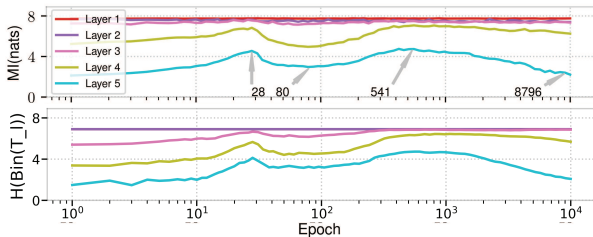
# Binning vs True Mutual Information

## Comparing to Previously Shown MI Plots:



# Binning vs True Mutual Information

## Comparing to Previously Shown MI Plots:



⇒ Past works were not showing MI but clustering (via binned-MI)!