# Measures of Information Reflect Memorization Patterns

Rachit Bansal, Danish Pruthi, Yonatan Belinkov

https://information-measures.cs.technion.ac.il
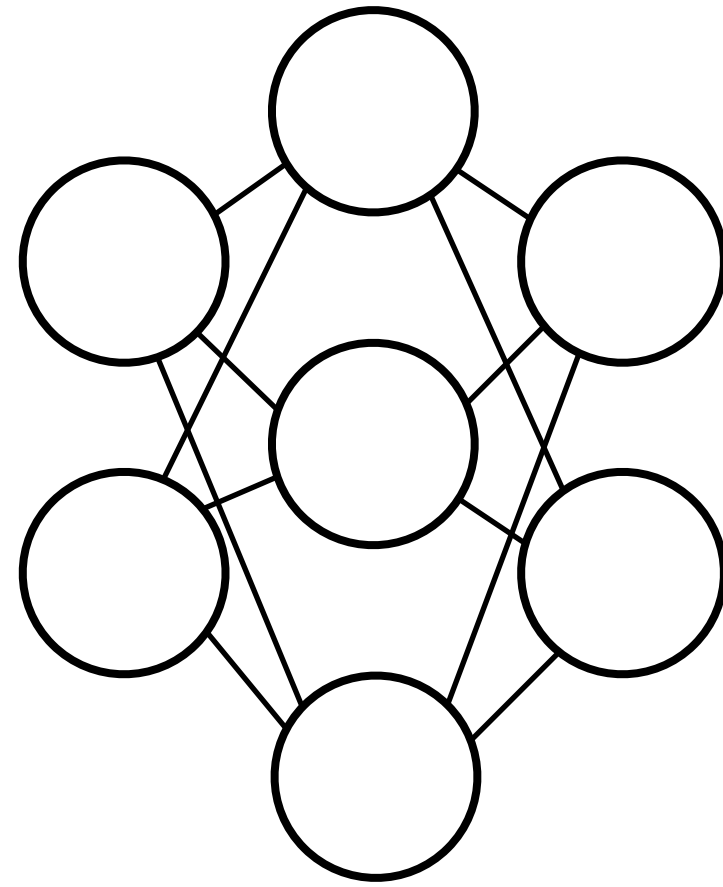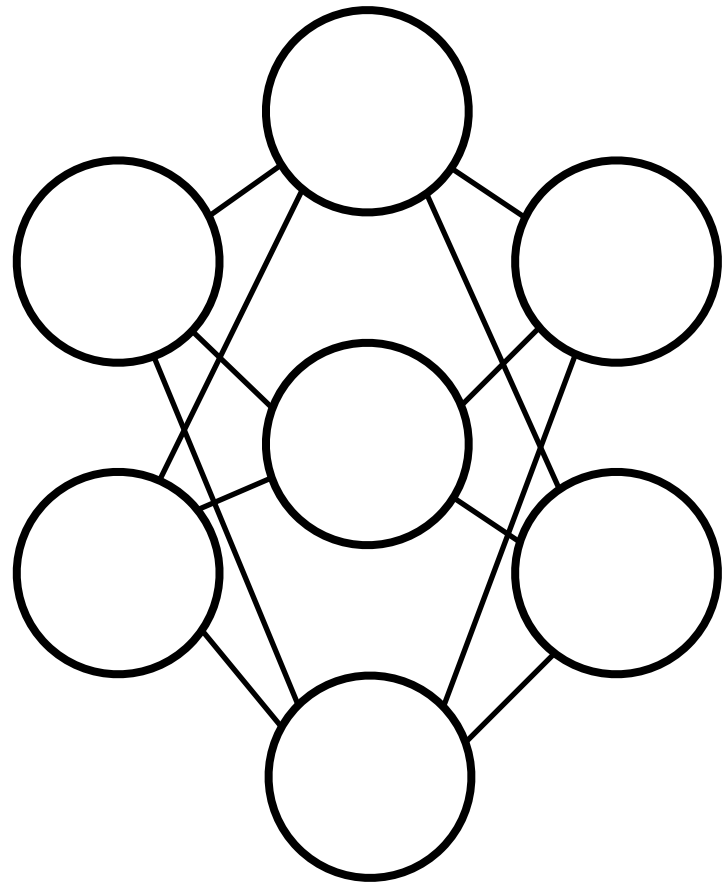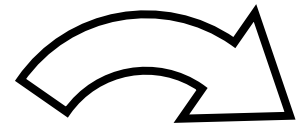Correspondence: racbansa@gmail.com

"Evaluating model generalization using intrinsic, information-theoretic metrics."

"*Evaluating* model *generalization* using *intrinsic,* *information-theoretic* metrics."
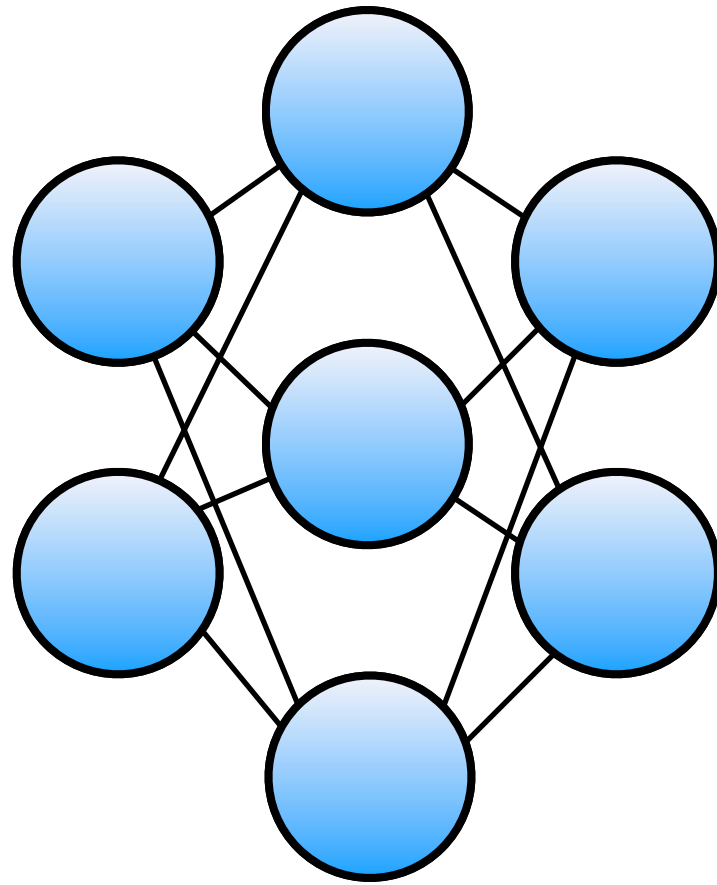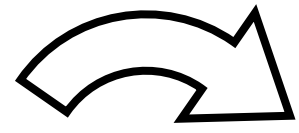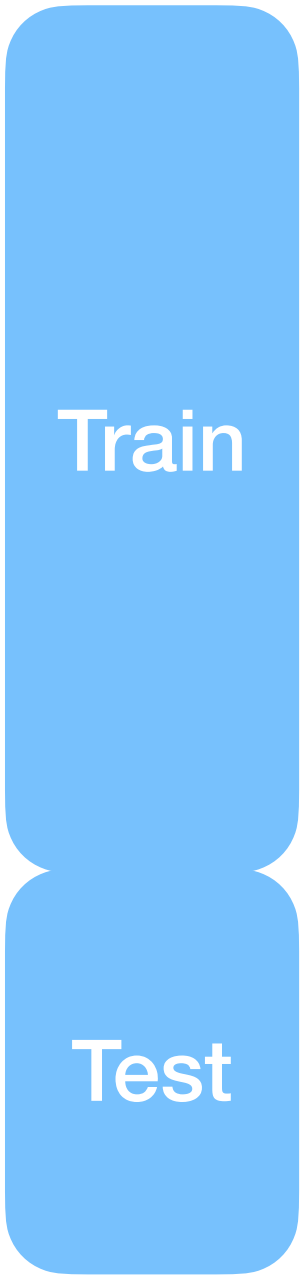
"*Evaluating* model *generalization* using *intrinsic, information-theoretic* metrics."

Data

Specialised sets are required
to identify these cases of memorization.

"*Evaluating* model *generalization* using *intrinsic, information-theoretic* metrics."

"*Evaluating* model *generalization* using *intrinsic, information-theoretic* metrics."
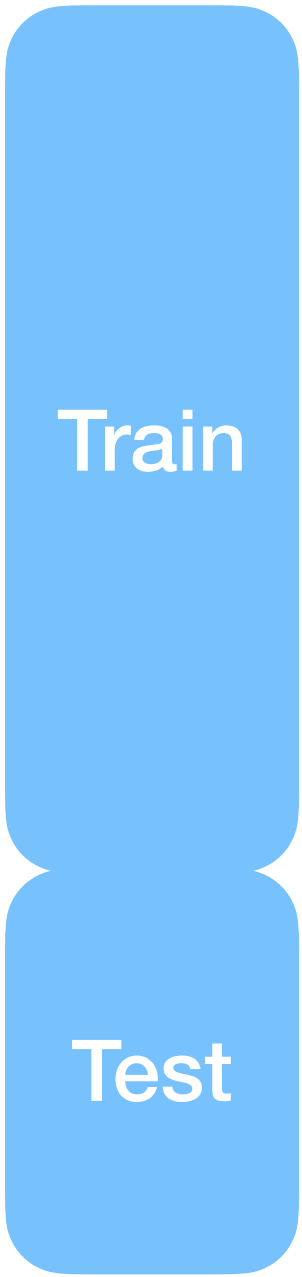
Example-level
Memorization

Example-level
Memorization

Heuristic
Memorization

"*Evaluating* model *generalization* using *intrinsic, information-theoretic* metrics."

"*Evaluating* model *generalization* using *intrinsic*, *information-theoretic* metrics."

Example-level
Memorization

Heuristic
Memorization

Example-level
Memorization

Heuristic
Memorization

Example-level
Memorization

Heuristic
Memorization

Example-level Memorization

Heuristic Memorization

Hypothesis: Neuron diversity is reflective of model generalization

Hypothesis: Neuron diversity is reflective of model generalization

"*Evaluating* model *generalization* using *intrinsic*, *information-theoretic* metrics."

"*Evaluating* model *generalization* using *intrinsic*, *information-theoretic* metrics."

Hypothesis: Neuron diversity is reflective of model generalization

$$A_1, \ldots, A_N \leftarrow \{f(x_i)\}_{i=1}^{S}$$

inter-neuron diversity

intra-neuron diversity

Examples

Neuron Index

inter-neuron diversity

intra-neuron diversity

Examples

Neuron Index

inter-neuron diversity

intra-neuron diversity

Examples

Neuron Index

$$A_1, \ldots, A_N \leftarrow \{f(x_i)\}_{i=1}^S$$

Intra-neuron Diversity
$\sim$
***Entropy***

inter-neuron diversity

intra-neuron diversity

Examples

Neuron Index

Intra-neuron Diversity

~

*Entropy*

$A_i$

$$A_1, \ldots, A_N \leftarrow \{f(x_i)\}_{i=1}^S$$

inter-neuron diversity

intra-neuron diversity

Examples

Neuron Index

$$A_1, \ldots, A_N \leftarrow \{f(x_i)\}_{i=1}^S$$

Intra-neuron Diversity
~
*Entropy*

$A_i$

inter-neuron diversity

intra-neuron diversity

Examples

Neuron Index

inter-neuron diversity

intra-neuron diversity

Examples

Neuron Index

inter-neuron diversity

intra-neuron diversity

Examples

Neuron Index

$$A_1, \ldots, A_N \leftarrow \{f(x_i)\}_{i=1}^S$$

Intra-neuron Diversity

$\sim$

**Entropy**

$A_i$

$$H(A_i) = \mathop{\mathbb{E}}_{\hat{a}_i^s \in A_i}[h(\hat{a}_i^s)] = \sum_{j=1}^{N_{\text{bins}}} p(\hat{a}_i^j) \log(\frac{1}{p(\hat{a}_i^j)})$$

inter-neuron diversity

intra-neuron diversity

Examples

Neuron Index

Intra-neuron Diversity
~
*Entropy*

$$H(A_i) = \mathop{\mathbb{E}}_{\hat{a}_i^s \in A_i}[h(\hat{a}_i^s)] = \sum_{j=1}^{N_{\text{bins}}} p(\hat{a}_i^j) \log\left(\frac{1}{p(\hat{a}_i^j)}\right)$$

$$A_1, \ldots, A_N \leftarrow \{f(x_i)\}_{i=1}^{S}$$

Inter-neuron Diversity
~
*Mutual Information*

Intra-neuron Diversity
~
**Entropy**

$$H(A_i) = \mathop{\mathbb{E}}_{\hat{a}_i^s \in A_i}[h(\hat{a}_i^s)] = \sum_{j=1}^{N_{\text{bins}}} p(\hat{a}_i^j) \log(\frac{1}{p(\hat{a}_i^j)})$$

$$A_1, \ldots, A_N \leftarrow \{f(x_i)\}_{i=1}^S$$

Inter-neuron Diversity
~
**Mutual Information**

inter-neuron diversity

intra-neuron diversity

Examples

Neuron Index

inter-neuron diversity

intra-neuron diversity

Examples

Neuron Index

$$A_1, \ldots, A_N \leftarrow \{f(x_i)\}_{i=1}^{S}$$

Intra-neuron Diversity

$\sim$

***Entropy***

$$H(A_i) = \mathop{\mathbb{E}}_{\hat{a}_i^s \in A_i}[h(\hat{a}_i^s)] = \sum_{j=1}^{N_{\text{bins}}} p(\hat{a}_i^j) \log\left(\frac{1}{p(\hat{a}_i^j)}\right)$$

Inter-neuron Diversity

$\sim$

***Mutual Information***

$$I(A_i) = \{I(A_i; A_1), \ldots, I(A_i; A_N)\}$$

$$A_1, \ldots, A_N \leftarrow \{f(x_i)\}_{i=1}^S$$

Intra-neuron Diversity
$\widetilde{\phantom{xx}}$
***Entropy***

$$H(A_i) = \mathop{\mathbb{E}}_{\hat{a}_i^s \in A_i} [h(\hat{a}_i^s)] = \sum_{j=1}^{N_{\text{bins}}} p(\hat{a}_i^j) \log(\frac{1}{p(\hat{a}_i^j)})$$
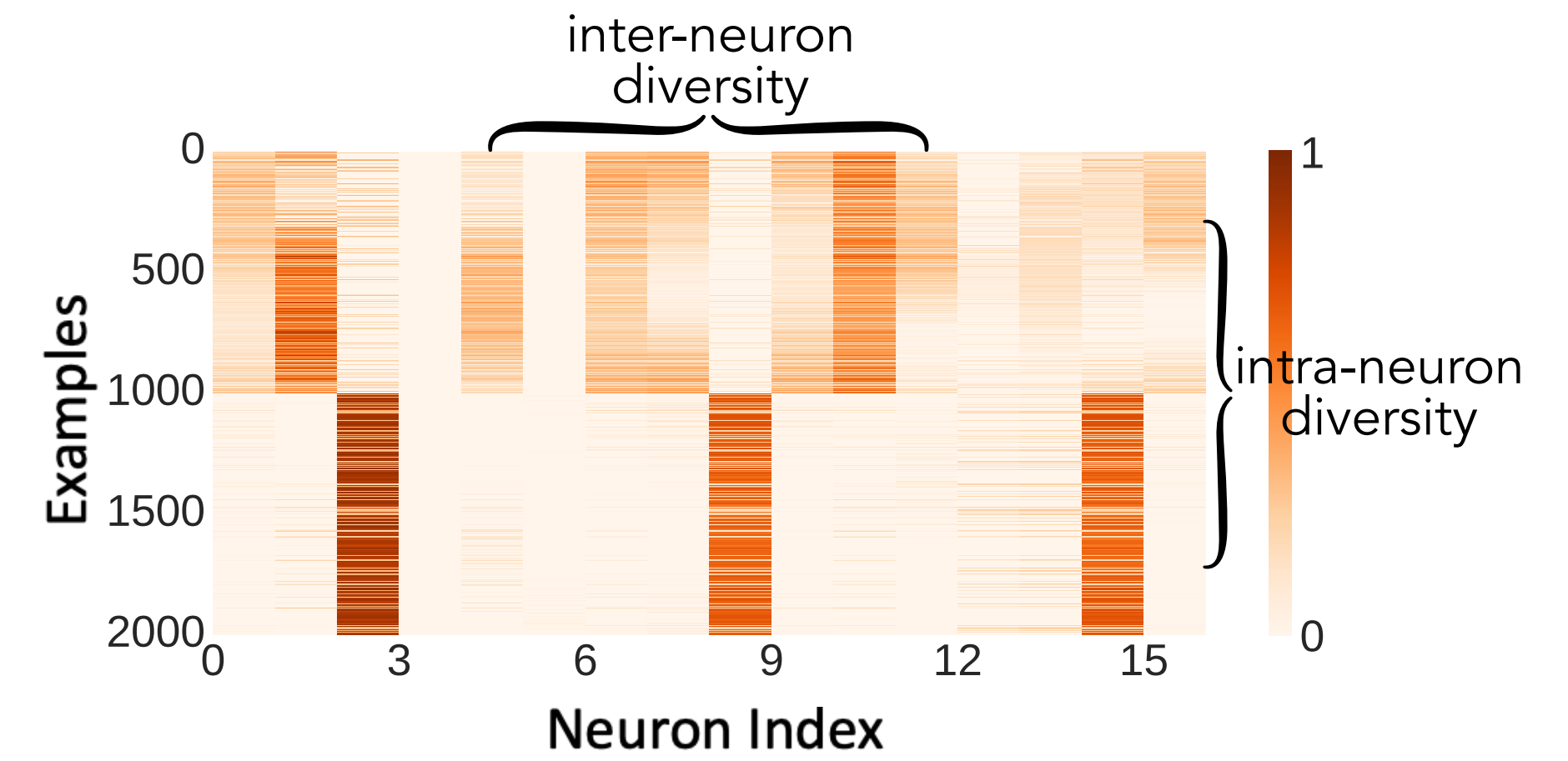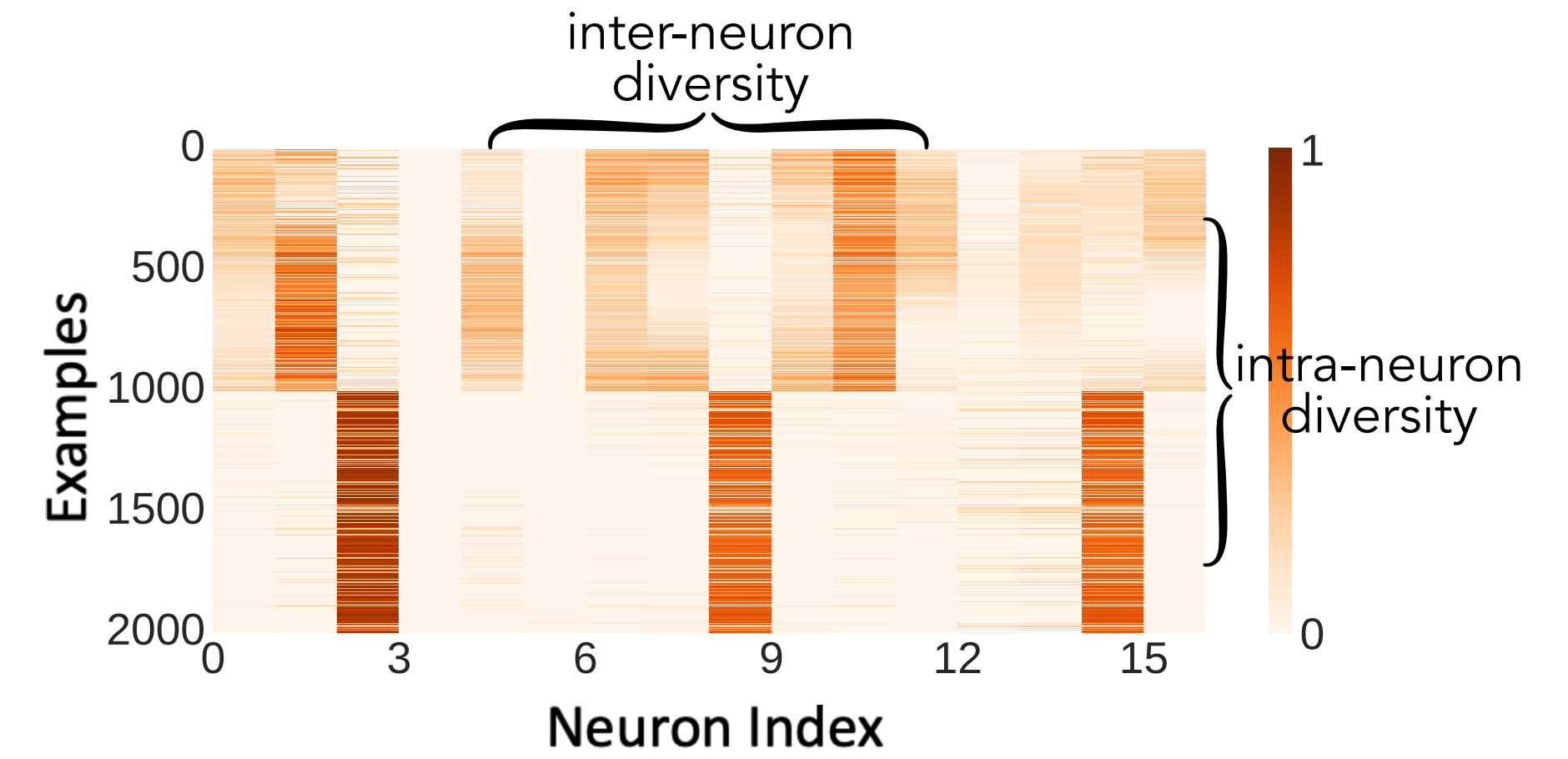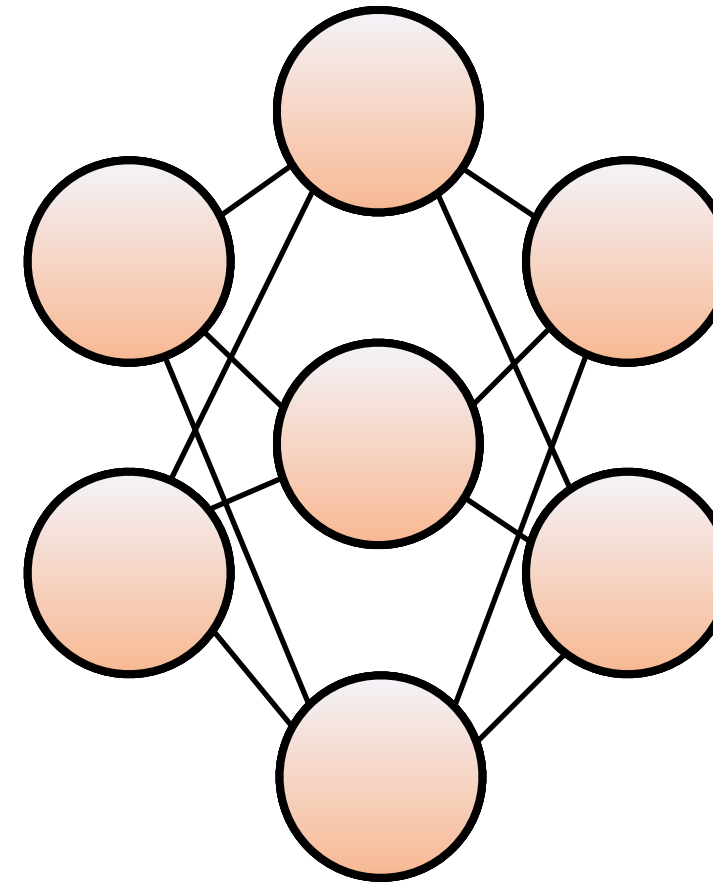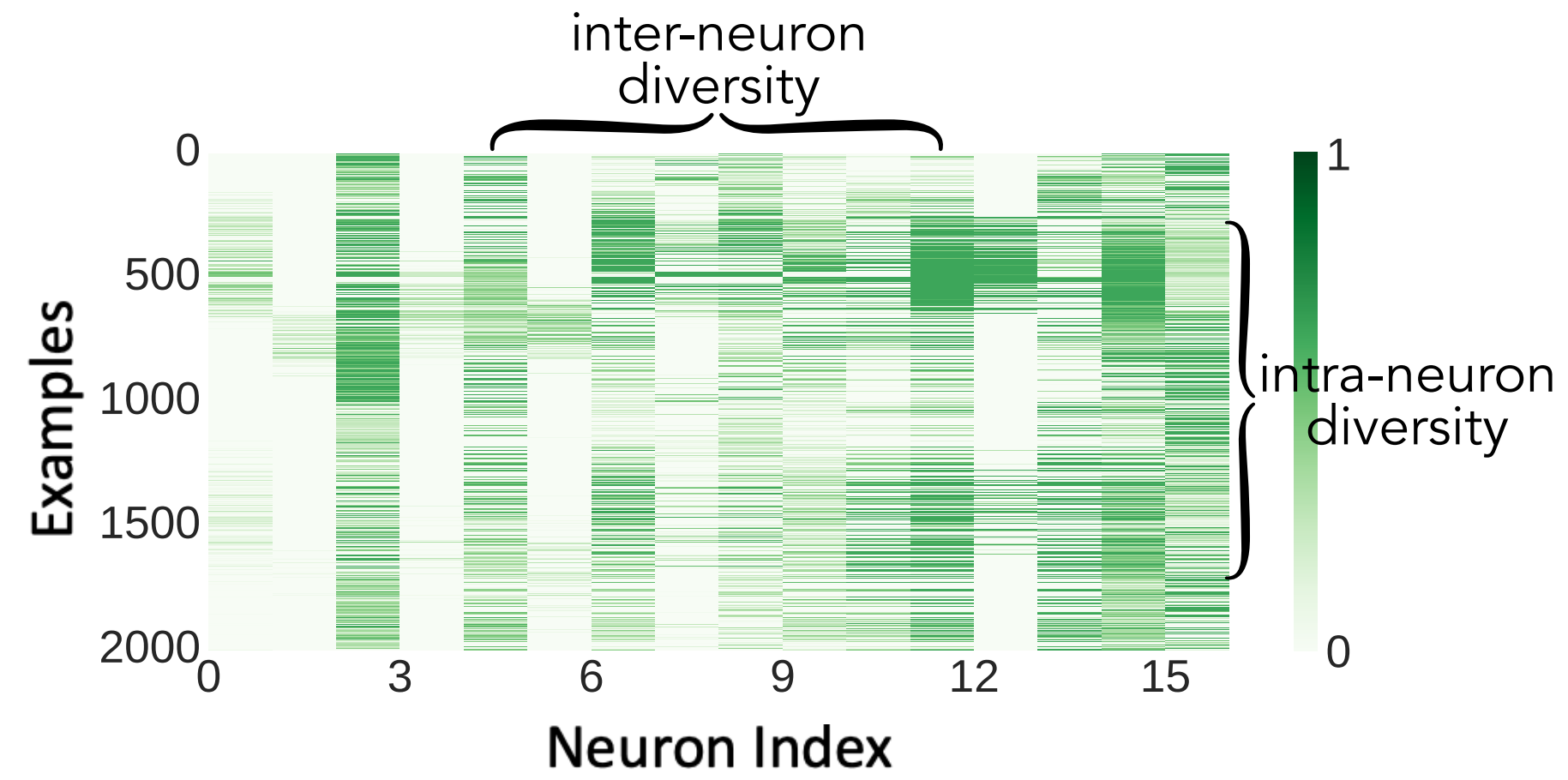
| Memorization | Diversity | |
|---|---|---|
| | Intra-neuron ($\propto$ Entropy) | Inter-neuron ($\propto$ MI$^{-1}$) |
| Heuristic Example-level | $\downarrow$ $\uparrow$ | $\downarrow$ $\uparrow$ |

Inter-neuron Diversity
$\widetilde{\phantom{xx}}$
***Mutual Information***

$$I(A_i) = \{I(A_i; A_1), \ldots, I(A_i; A_N)\}$$

*Some experimental results*

3-layered MLP on Colored MNIST

# 3-layered MLP on Colored MNIST

3-layered MLP on Colored MNIST

Hypothesis:
↓ Entropy
↑ MI

$\alpha$

"proportion of the labels synthetically
changed to follow the spurious correlation"

3-layered MLP on Colored MNIST

Hypothesis:
↓ Entropy
↑ MI

Neuron Entropy

$\alpha$

"proportion of the labels synthetically
changed to follow the spurious correlation"

3-layered MLP on Colored MNIST

Hypothesis:
↓ Entropy
↑ MI

α

"proportion of the labels synthetically
changed to follow the spurious correlation"

*3-layered MLP on Colored MNIST*

Hypothesis:
↓ Entropy
↑ MI

Neuron Entropy

Validation Accuracy

Layer-1
Layer-2
Layer-3

$\alpha$

"proportion of the labels synthetically changed to follow the spurious correlation"

*3-layered MLP on Colored MNIST*

Hypothesis:
↓ Entropy
↑ MI

$\alpha$

"proportion of the labels synthetically
changed to follow the spurious correlation"

Layer-1
Layer-2
Layer-3

*3-layered MLP on Colored MNIST*

Hypothesis:
↓ Entropy
↑ MI

Layer-1
Layer-2
Layer-3

$\alpha$

"proportion of the labels synthetically
changed to follow the spurious correlation"

3-layered MLP on Shuffled MNIST

Hypothesis:
↑ Entropy
↓ MI

*3-layered MLP on Shuffled MNIST*

Neuron Entropy

Validation Accuracy

$\beta$

"proportion of the label noise
added to the training dataset"

Hypothesis:
↑ Entropy
↓ MI

3-layered MLP on Shuffled MNIST

Hypothesis:
↑ Entropy
↓ MI

"proportion of the label noise added to the training dataset"

# Take aways

# Take aways

- Intrinsic activation patterns in neural networks are indicative of memorization.



RoBERTa-base (Seed=0) on BiasInBios

# Take aways

- Intrinsic activation patterns in neural networks are indicative of memorization.

- Neuron diversity (as a measure of information) correlates to generalization.



RoBERTa-base (Seed=0) on BiasInBios

# Take aways

- Intrinsic activation patterns in neural networks are indicative of memorization.

- Neuron diversity (as a measure of information) correlates to generalization.

- Useful for model selection. Potentially extensible for OOD detection, model regularization, and understanding training dynamics.



RoBERTa-base (Seed=0) on BiasInBios

# Take aways

- Intrinsic activation patterns in neural networks are indicative of memorization.

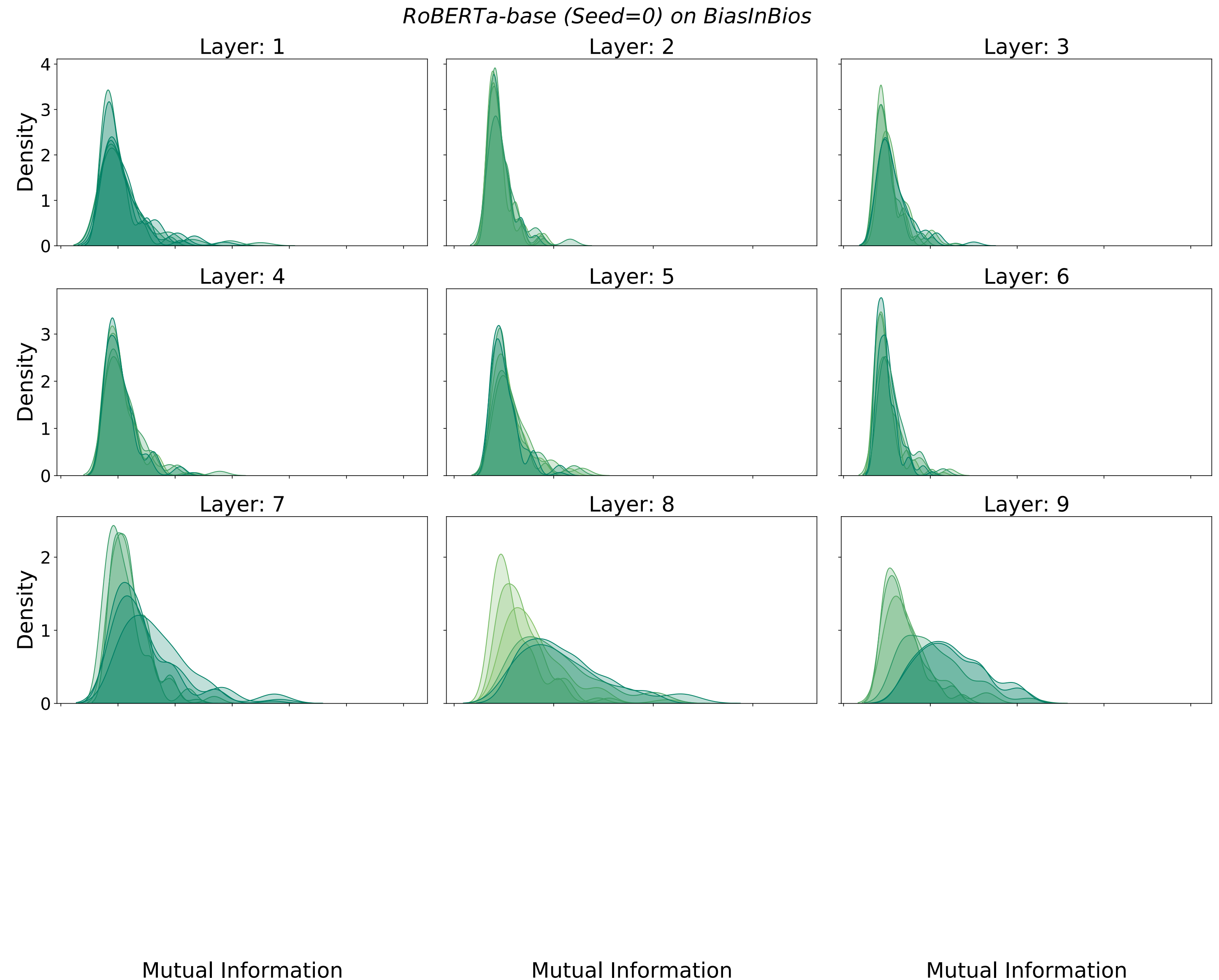- Neuron diversity (as a measure of information) correlates to generalization.

- Useful for model selection. Potentially extensible for OOD detection, model regularization, and understanding training dynamics.



RoBERTa-base (Seed=0) on BiasInBios