

Meta-Learning Representations for Continual Learning

Khurram Javed

Work done in collaboration with Martha White

Continual Learning Prediction

Continual Learning Prediction

$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$

Continual Learning Prediction

$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$

Data could be highly correlated

Continual Learning Prediction

$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$

Data could be highly correlated

AA **AA** **AA** **BB** **BB** **BB** **BB** **CC** **cc** **CC** **DD** **DD** **DD** **ZZ** **ZZ** **ZZ** **ZZ**

Continual Learning Prediction

$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$

Data could be highly correlated

A**A****A****A** **B****B****B****B****B** **C****C****C****C** **D****D****D****D** **Z****Z****Z****Z****Z**

Naive solution using a neural network:

Continual Learning Prediction

$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$

Data could be highly correlated

A**A****A****A** **B****B****B****B****B** **C****C****C****C** **D****D****D****D** **Z****Z****Z****Z****Z**

Naive solution using a neural network:

1. Get a new sample

Continual Learning Prediction

$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$

Data could be highly correlated

A**A****A****A** **B****B****B****B****B** **C****C****C** **D****D****D****D** **Z****Z****Z****Z****Z**

Naive solution using a neural network:

1. Get a new sample
2. Update parameters w.r.t latest sample

Continual Learning Prediction

$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$

Data could be highly correlated

AA **AA** **AA** **BB** **BB** **BB** **BB** **CC** **CC** **CC** **DD** **DD** **DD** **DD** **ZZ** **ZZ** **ZZ** **ZZ** **ZZ**

Naive solution using a neural network:

1. Get a new sample
2. Update parameters w.r.t latest sample
3. Repeat

Continual Learning Prediction

$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$

Data could be highly correlated

A**A****A****A** **B****B****B****B** **C****C****C** **D****D****D****D** **Z****Z****Z****Z****Z**

Naive solution using a neural network:

1. Get a new sample
2. Update parameters w.r.t latest sample
3. Repeat

Continual Learning Prediction

$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$

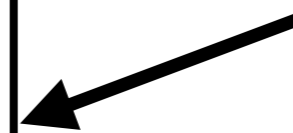
Data could be highly correlated

AA A A A BBBB B C C C C D D D D D Z Z Z Z Z Z

Naive solution using a neural network:

1. Get a new sample
2. Update parameters w.r.t latest sample
3. Repeat

**Traditional Neural
Network Stack**



Traditional Neural Network Stack

Traditional Neural Network Stack

Dataset : Omniglot

Traditional Neural Network Stack

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning

Traditional Neural Network Stack

Dataset : Omniglot

- ~~~950 characters from multiple alphabets for representation learning~~

Traditional Neural Network Stack

Dataset : Omniglot

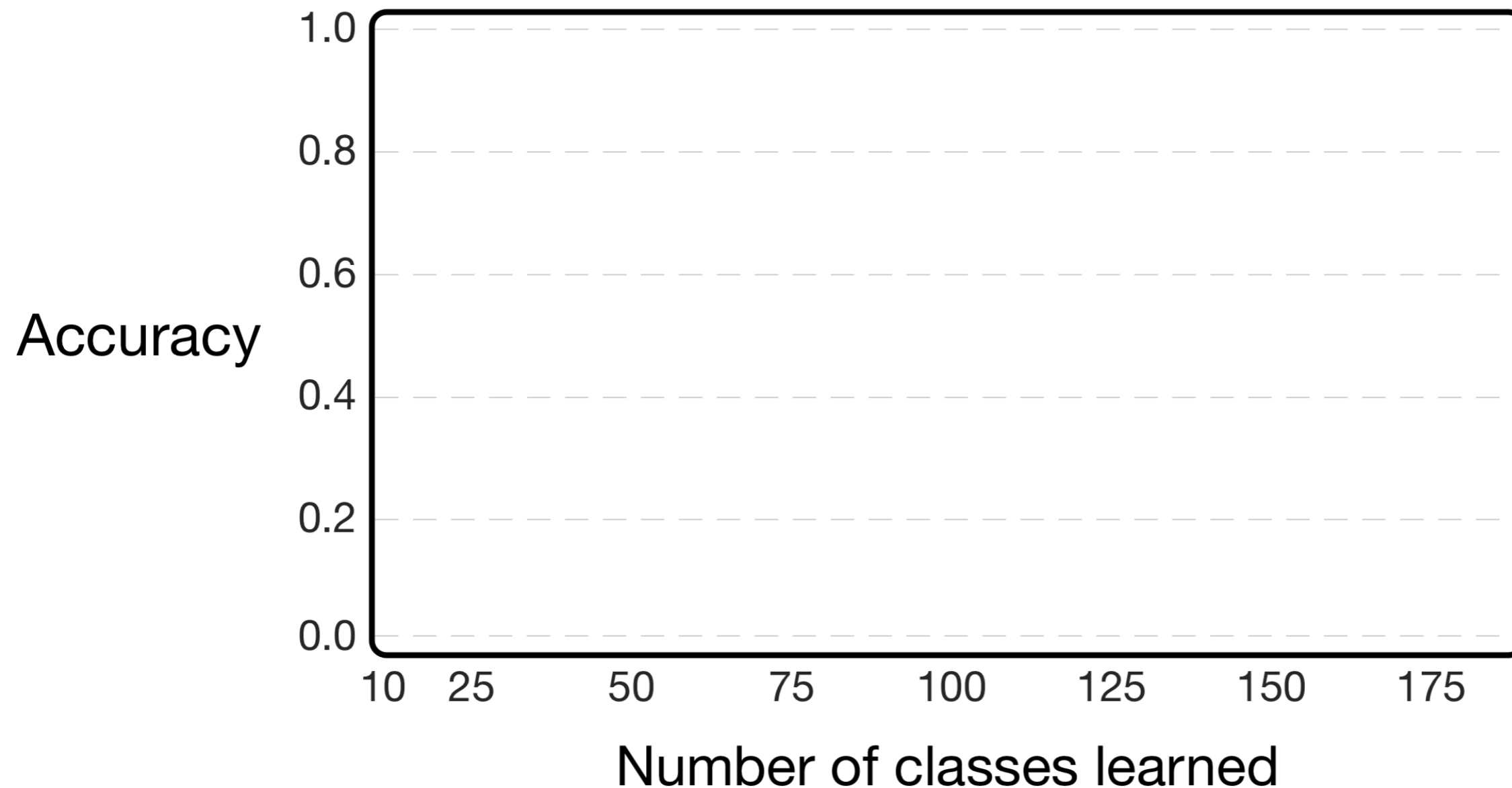
- ~~~950 characters from multiple alphabets for representation learning~~
- ~600 characters from multiple alphabets for Continual Learning Prediction

Traditional Neural Network Stack

Dataset : Omniglot

- ~~~950 characters from multiple alphabets for representation learning~~
- ~600 characters from multiple alphabets for Continual Learning Prediction

Omniglot Training Trajectory Performance

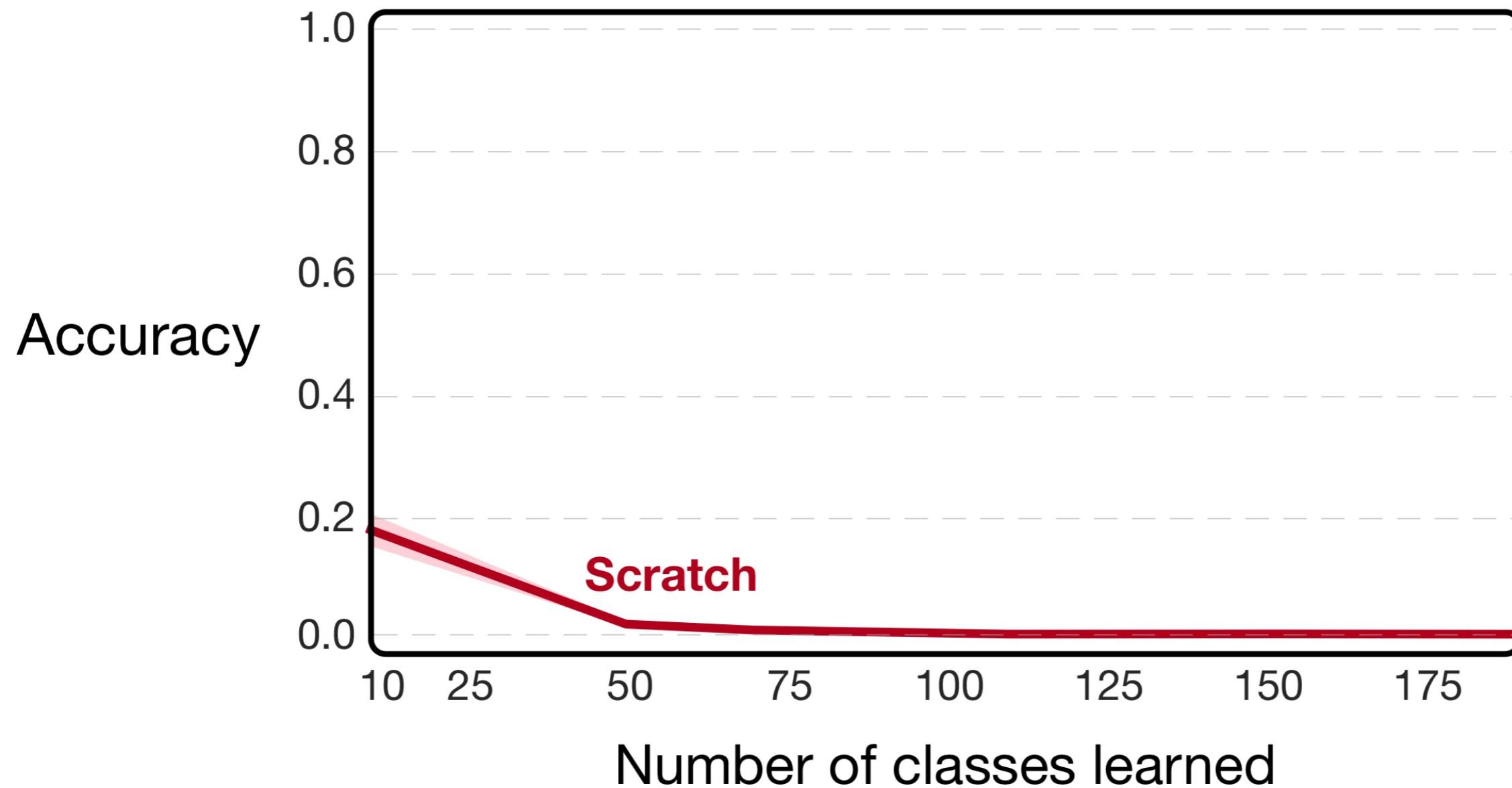


Traditional Neural Network Stack

Dataset : Omniglot

- ~~~950 characters from multiple alphabets for representation learning~~
- ~600 characters from multiple alphabets for Continual Learning Prediction

Omniglot Training Trajectory Performance



Catastrophic Interference

Catastrophic Interference

Neural networks suffer from ***catastrophic*** interference

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling
2. Dense inputs

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling
2. Dense inputs
3. Global and greedy learning update

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling

2. Dense inputs

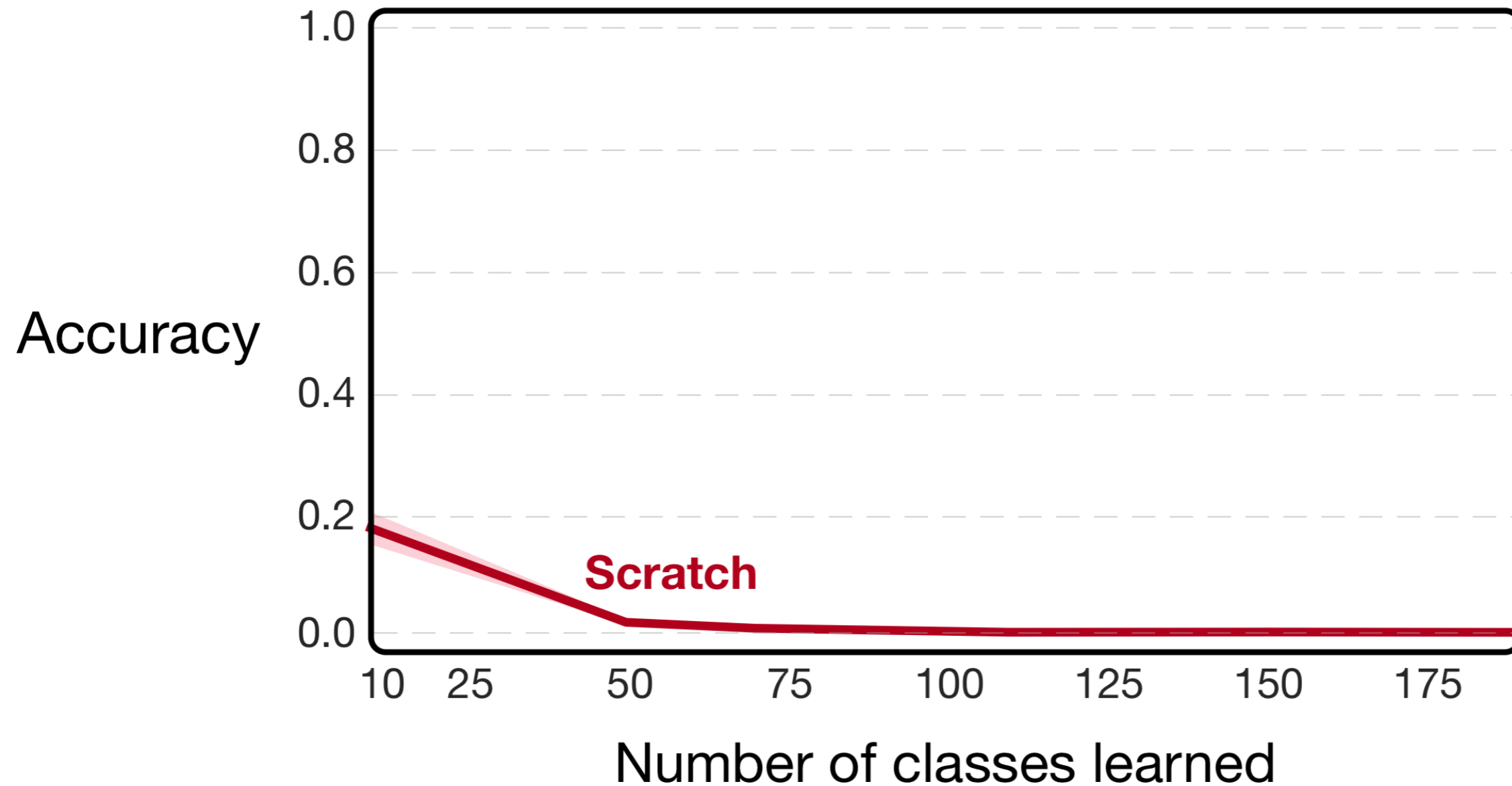
3. Global and greedy learning update

Non-IID sampling

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

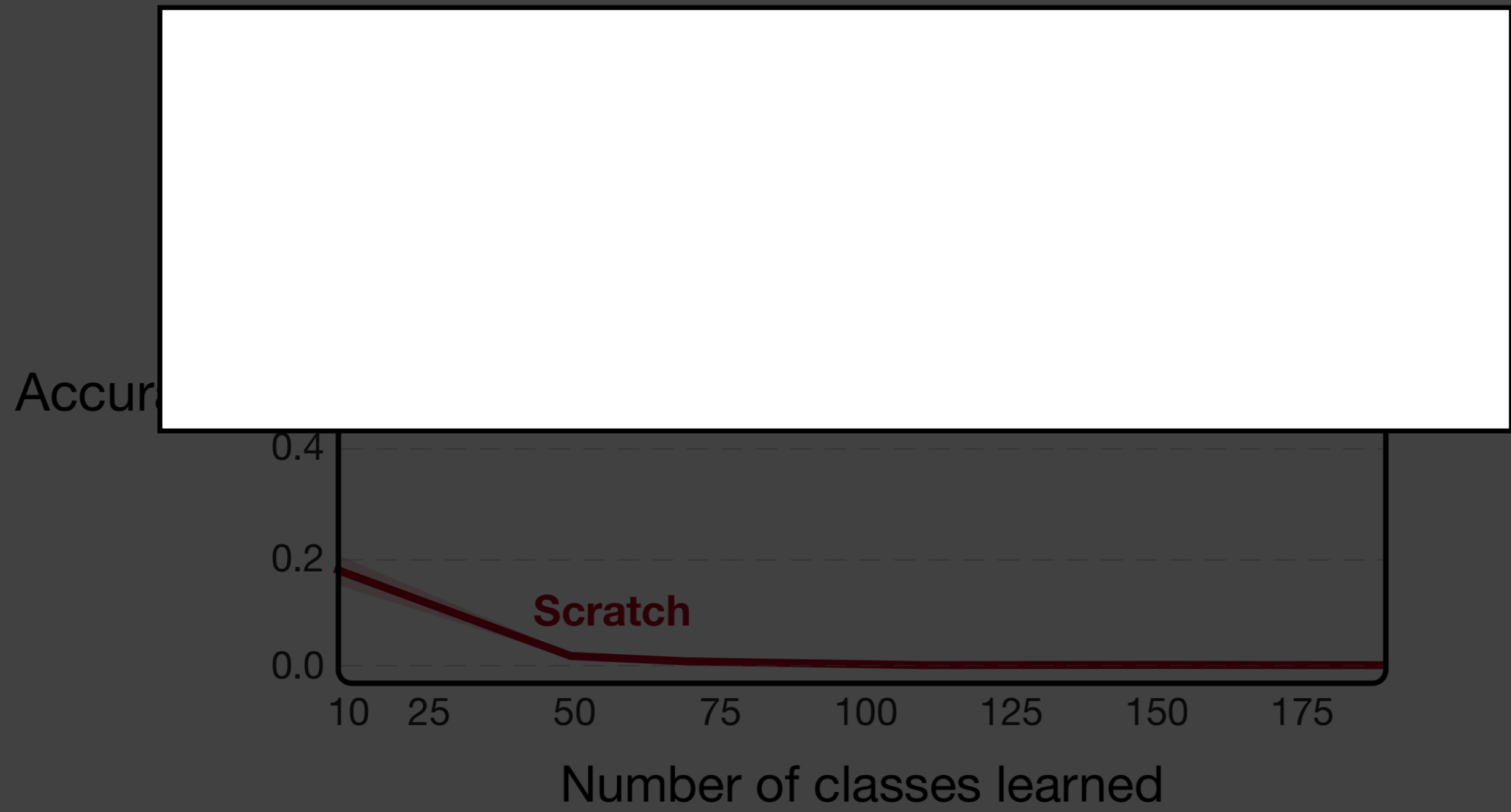
Omniglot Training Trajectory Performance



Non-IID sampling

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction



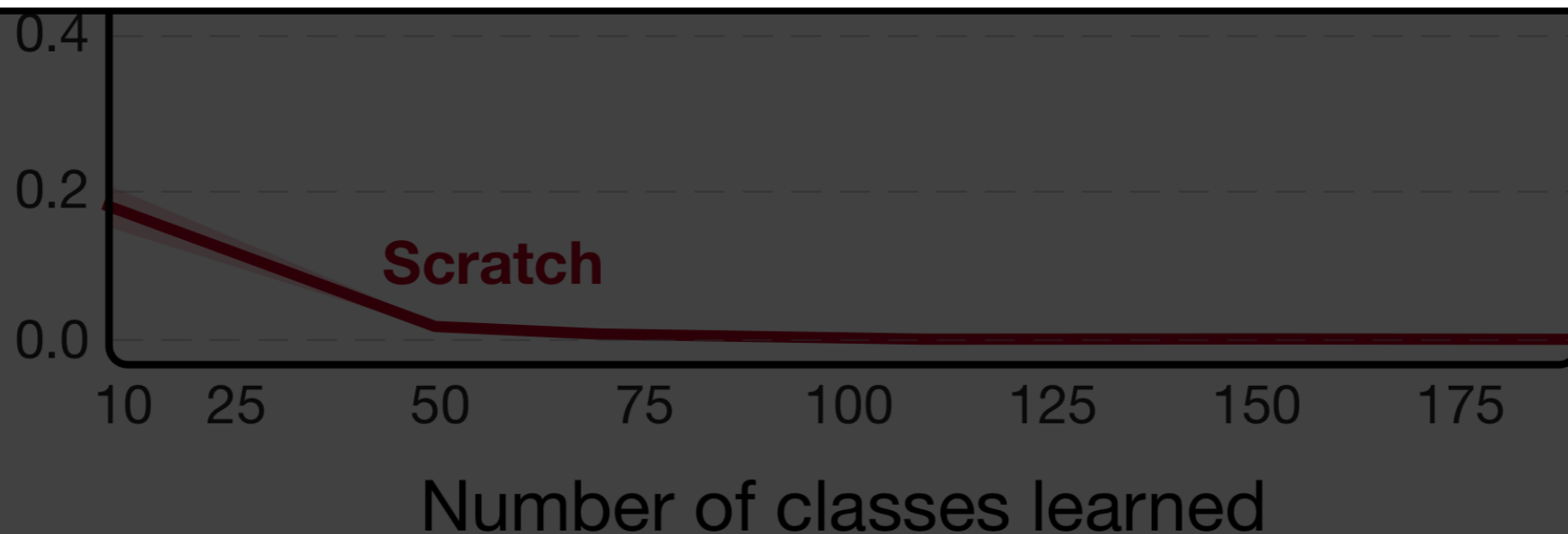
Non-IID sampling

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$

Accur

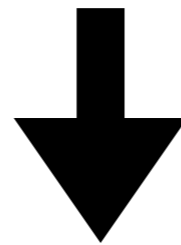


Non-IID sampling

Dataset : Omniglot

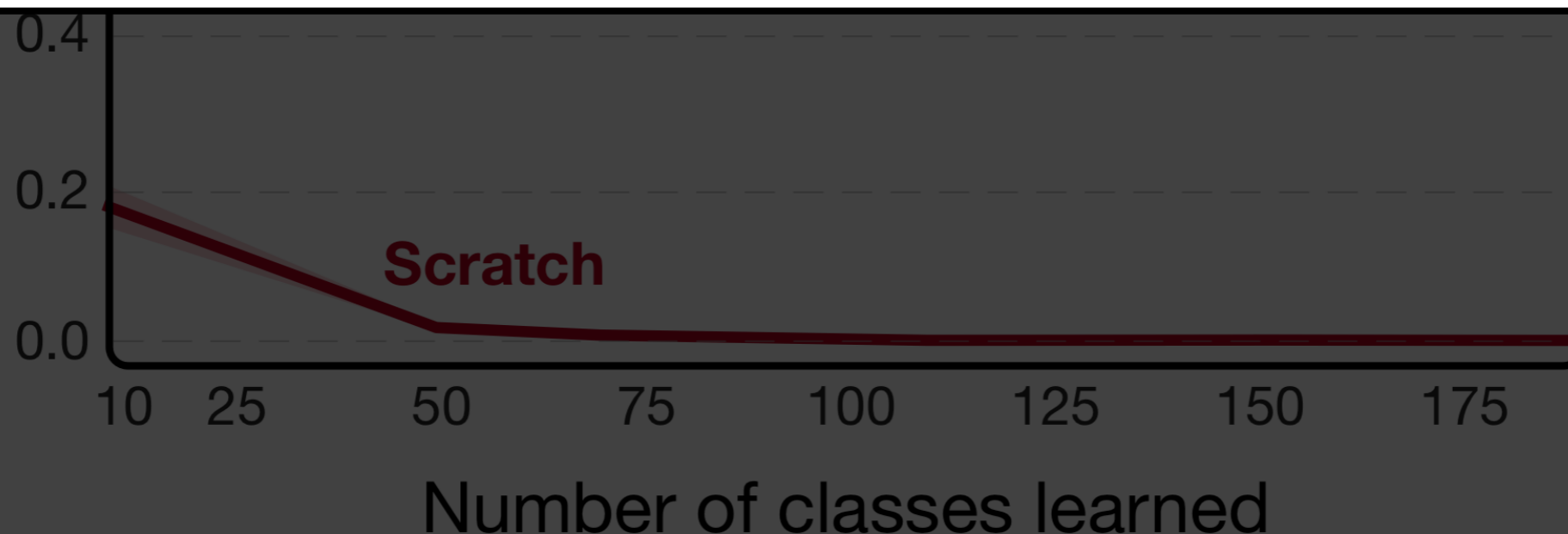
- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$



Shuffle order

Accur

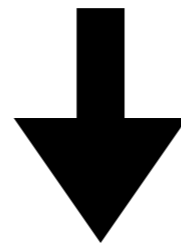


Non-IID sampling

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

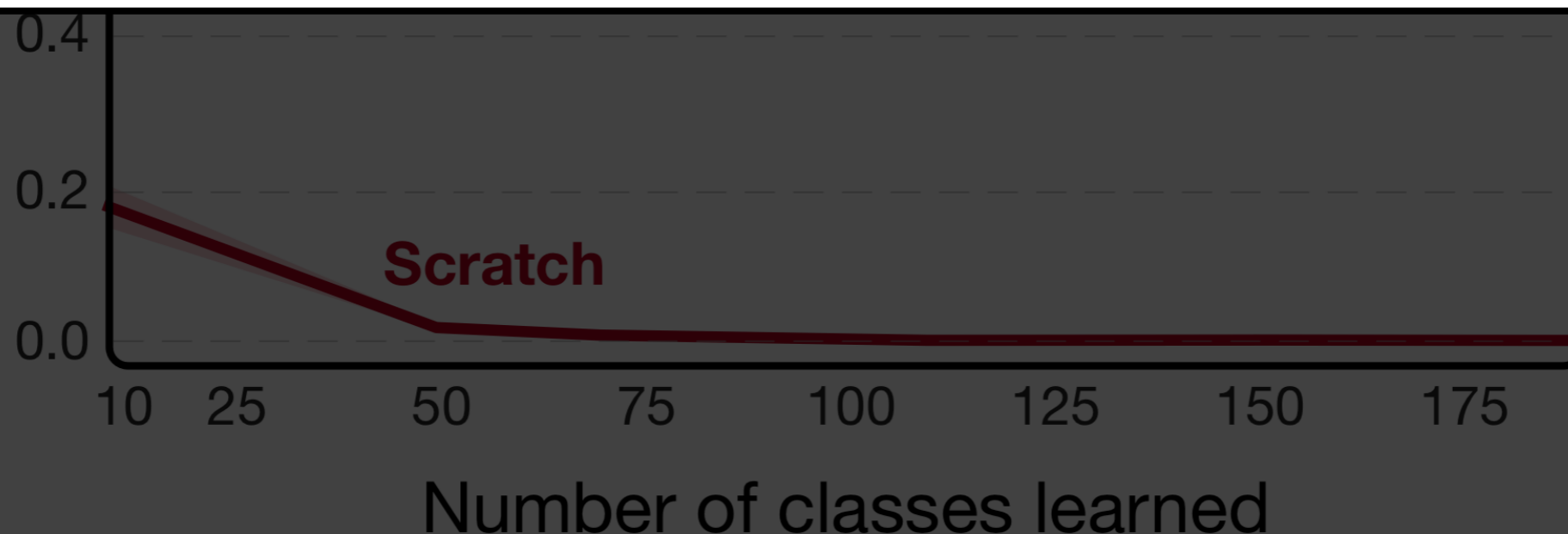
$$D_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots$$



Shuffle order

$$D_{rand} = (X_{10}, Y_{10}), (X_1, Y_1), \dots, (X_m, Y_m), \dots$$

Accur

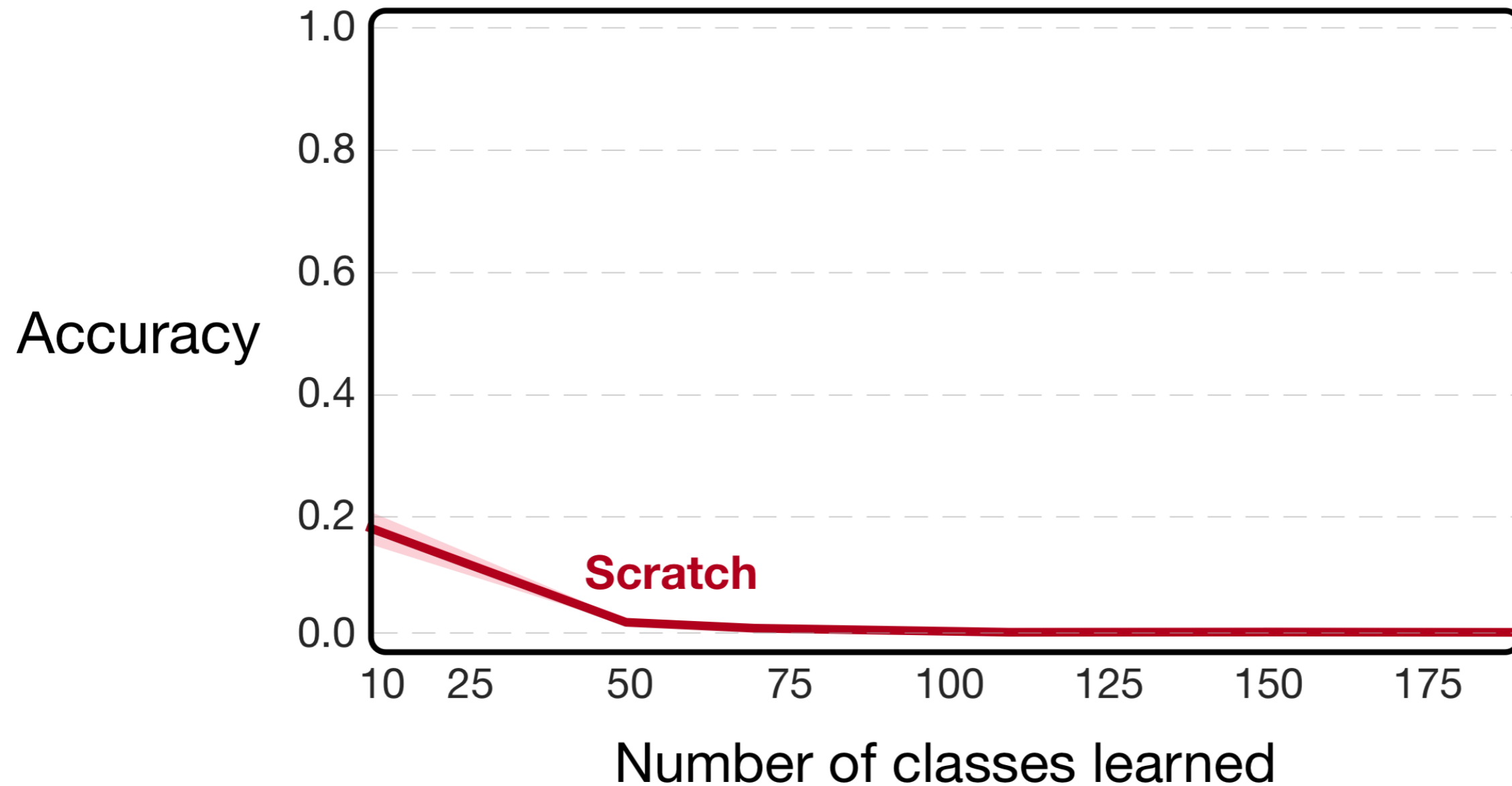


Non-IID sampling

Dataset : Omniglot

- ~~~950 characters from multiple alphabets for representation learning~~
- ~600 characters from multiple alphabets for continual learning prediction

Omniglot Training Trajectory Performance

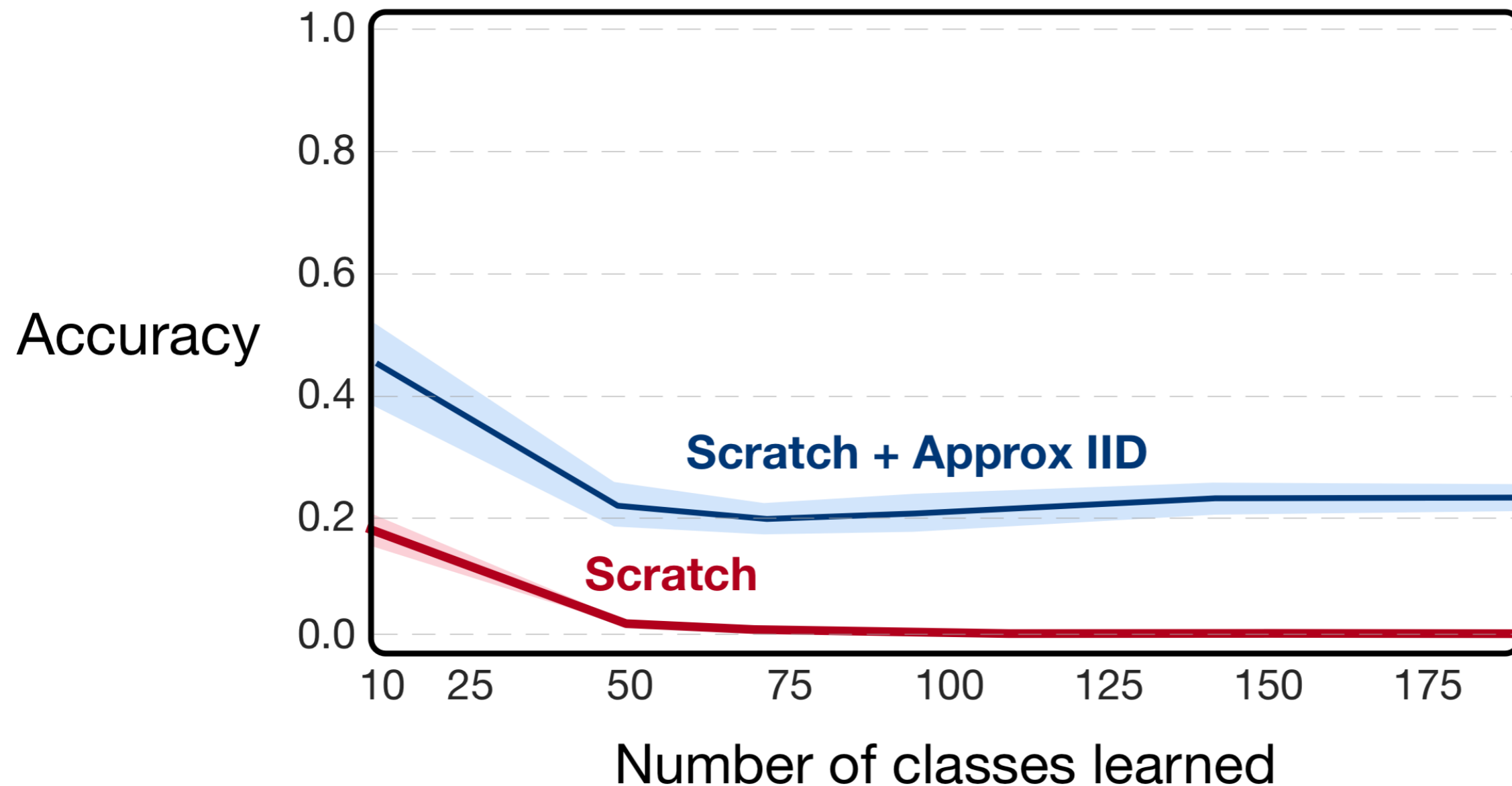


Non-IID sampling

Dataset : Omniglot

- ~~~950 characters from multiple alphabets for representation learning~~
- ~600 characters from multiple alphabets for continual learning prediction

Omniglot Training Trajectory Performance



Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling
2. Dense inputs
3. Global and greedy update

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling

2. Dense inputs

3. Global and greedy update

Dense inputs

64



64

Dense inputs

=> Input vector of length $64 \times 64 = 4096$



64

64

Dense inputs



=> Input vector of length $64 \times 64 = 4096$

Most values are non-zero for most natural images

64

64

Dense inputs



=> Input vector of length $64 \times 64 = 4096$

Most values are non-zero for most natural images

Alternative?

64

64

Dense inputs

64



64

=> Input vector of length $64 \times 64 = 4096$

Most values are non-zero for most natural images

Alternative?

Sparse input representations

Example: Tile coding

Dense Inputs

Dataset : Omniglot

- ~~~950 characters from multiple alphabets for representation learning~~
- ~600 characters from multiple alphabets for continual learning prediction

Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Plan

Learn a representation on the ~950 characters, evaluate on the same ~600

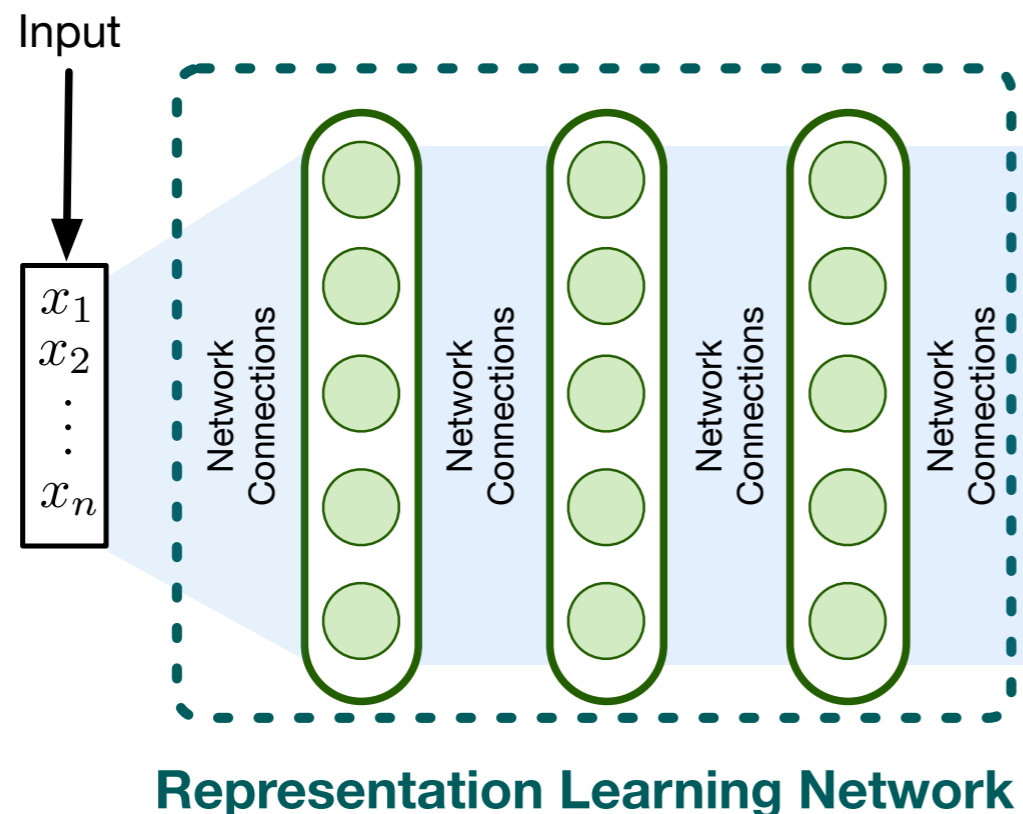
Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Plan

Learn a representation on the ~950 characters, evaluate on the same ~600



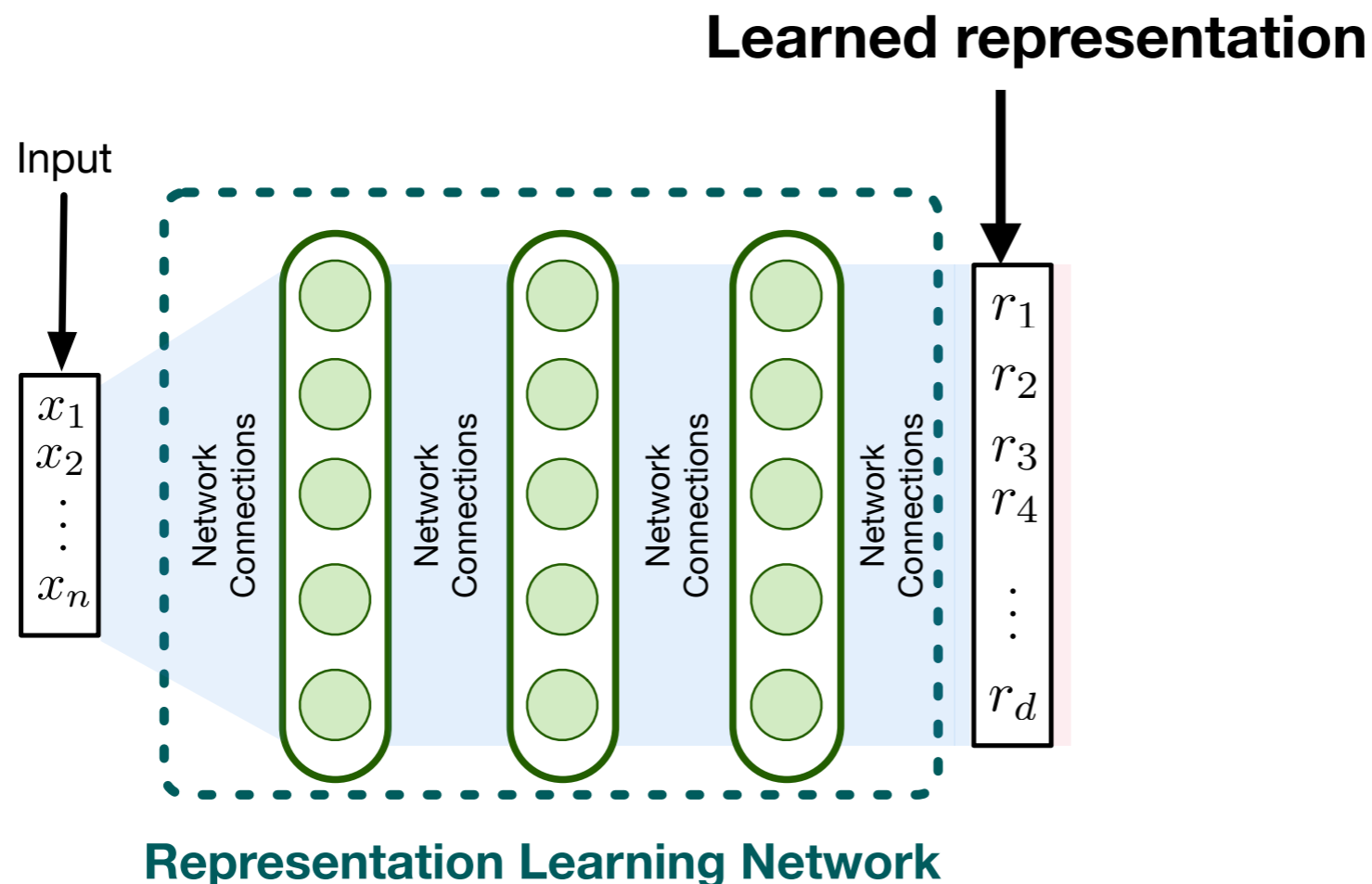
Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Plan

Learn a representation on the ~950 characters, evaluate on the same ~600



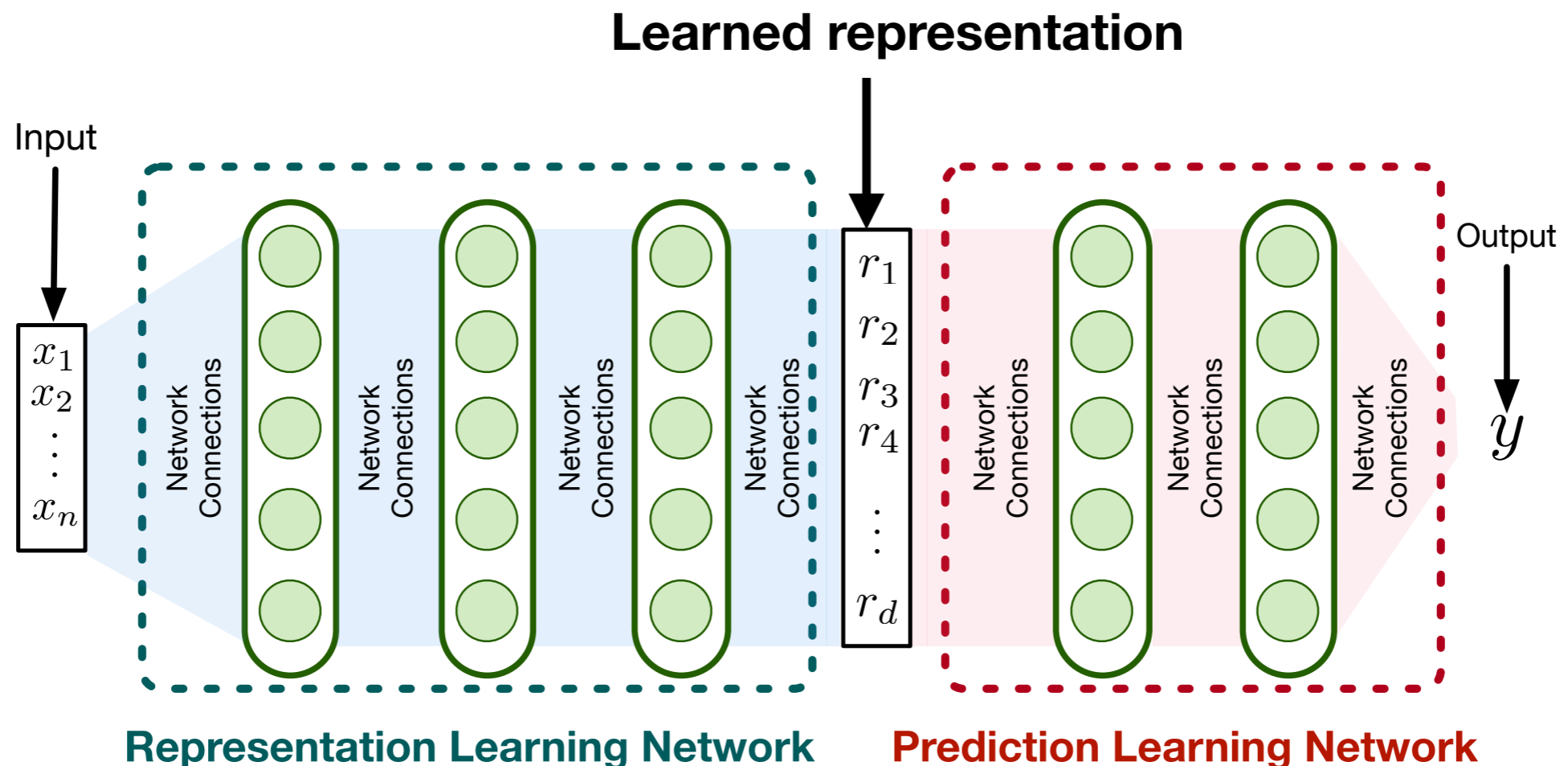
Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Plan

Learn a representation on the ~950 characters, evaluate on the same ~600



Dense Inputs

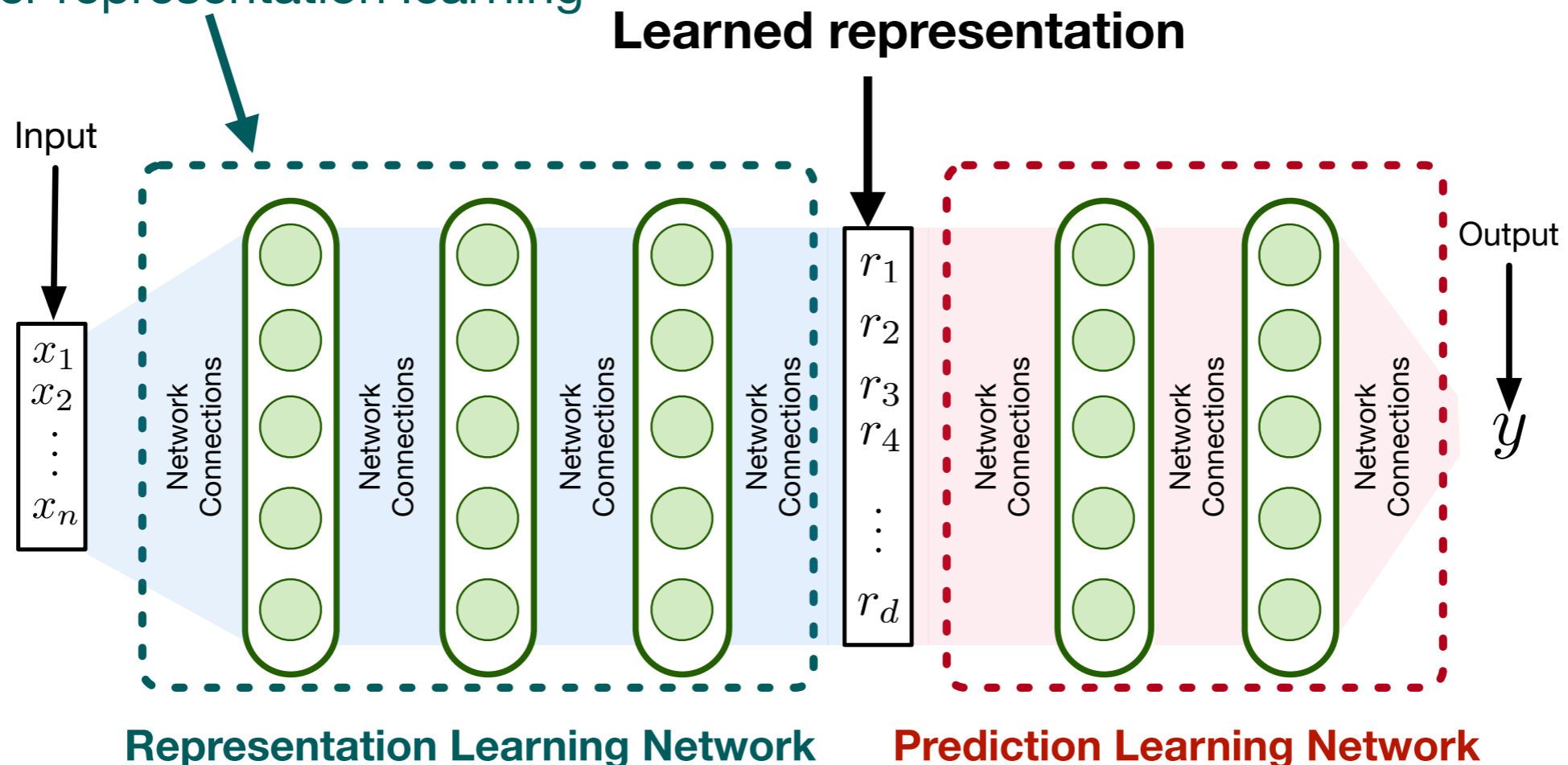
Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Plan

Learn a representation on the ~950 characters, evaluate on the same ~600

Fixed after representation learning



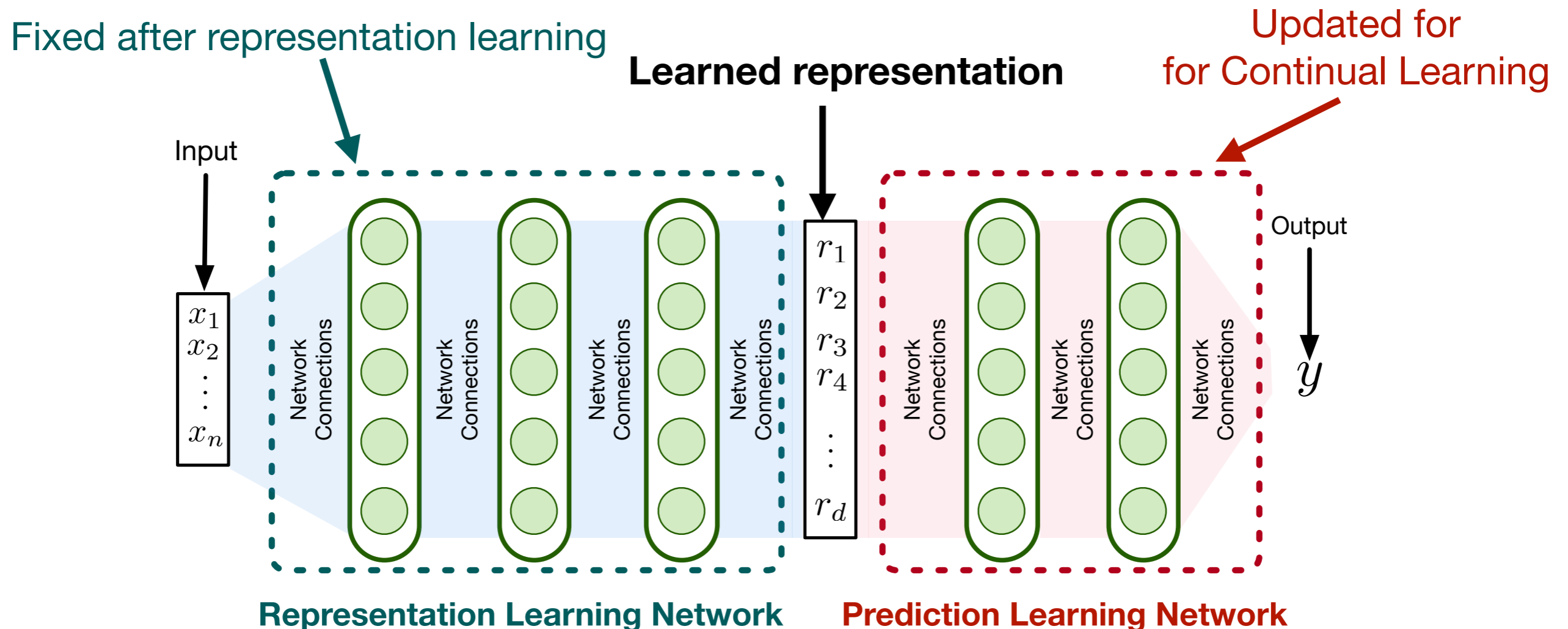
Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Plan

Learn a representation on the ~950 characters, evaluate on the same ~600



Dense Inputs

Dataset : Omniglot

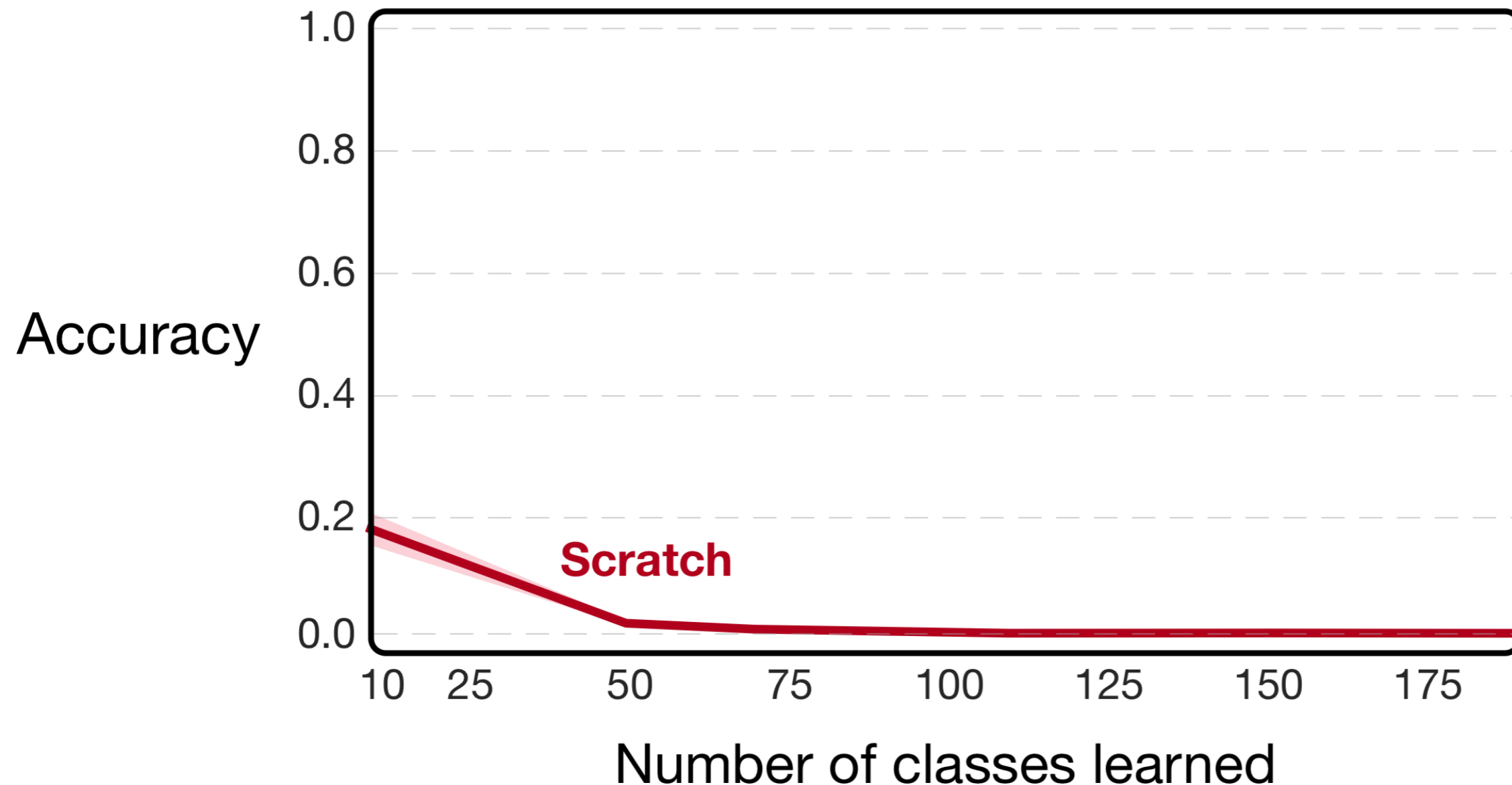
- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Omniglot Training Trajectory Performance

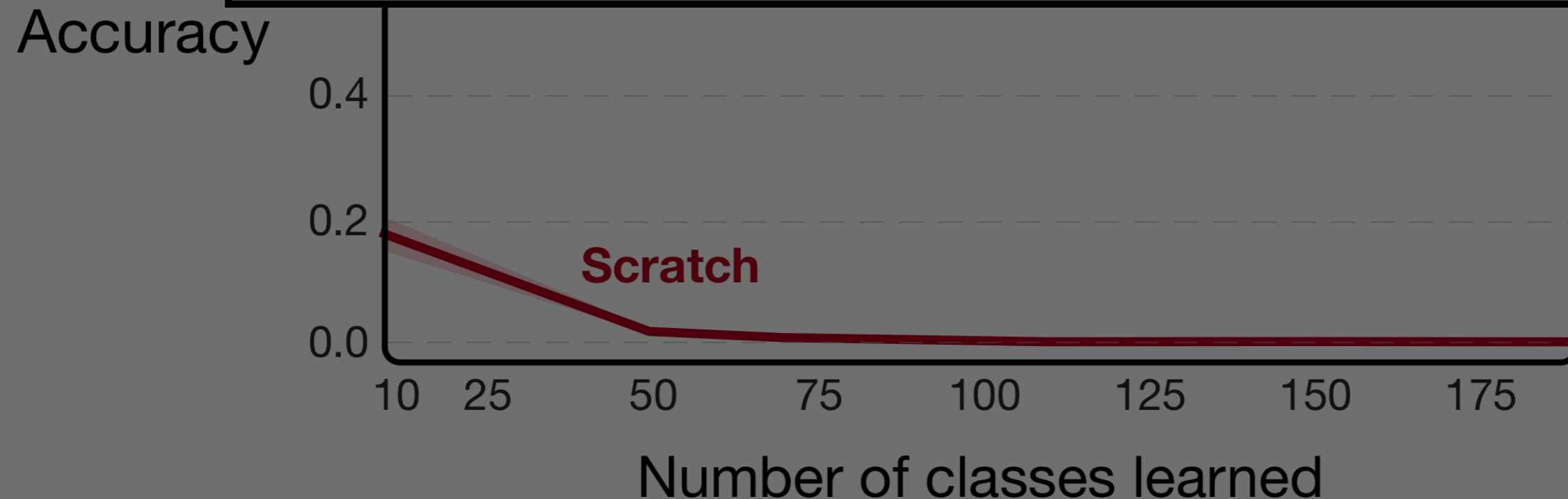


Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Is “Scratch” a fair baseline?



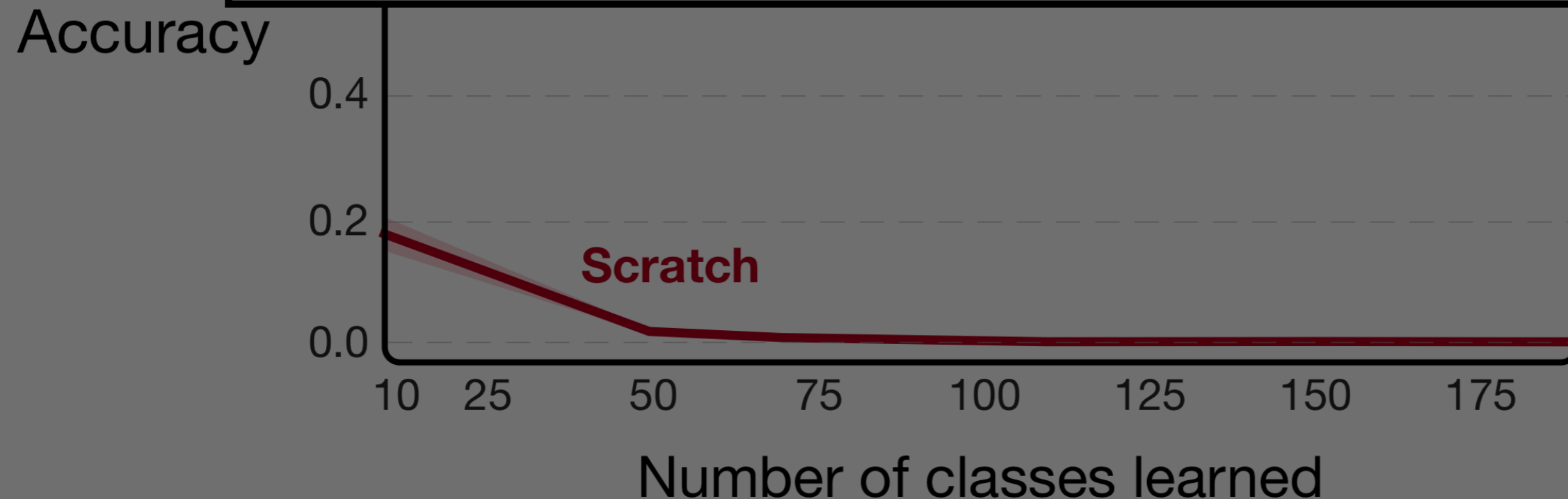
Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Is “Scratch” a fair baseline?

- Scratch is not longer a fair baseline



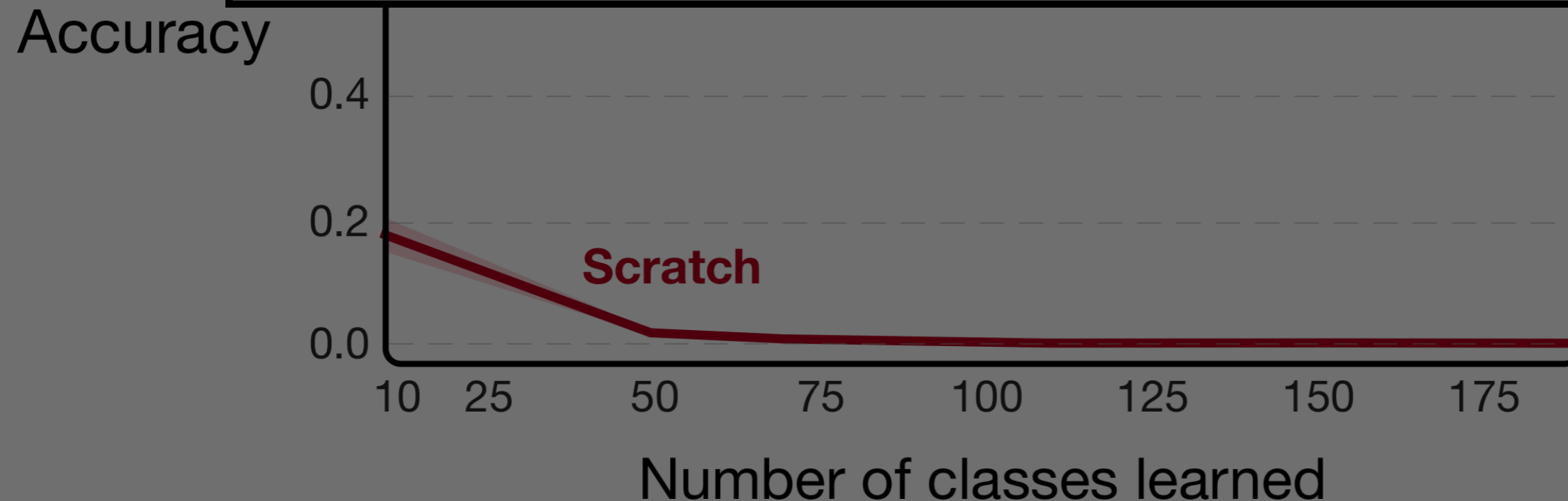
Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Is “Scratch” a fair baseline?

- Scratch is not longer a fair baseline
- Need a baseline that uses the information from the representation learning dataset



Dense Inputs

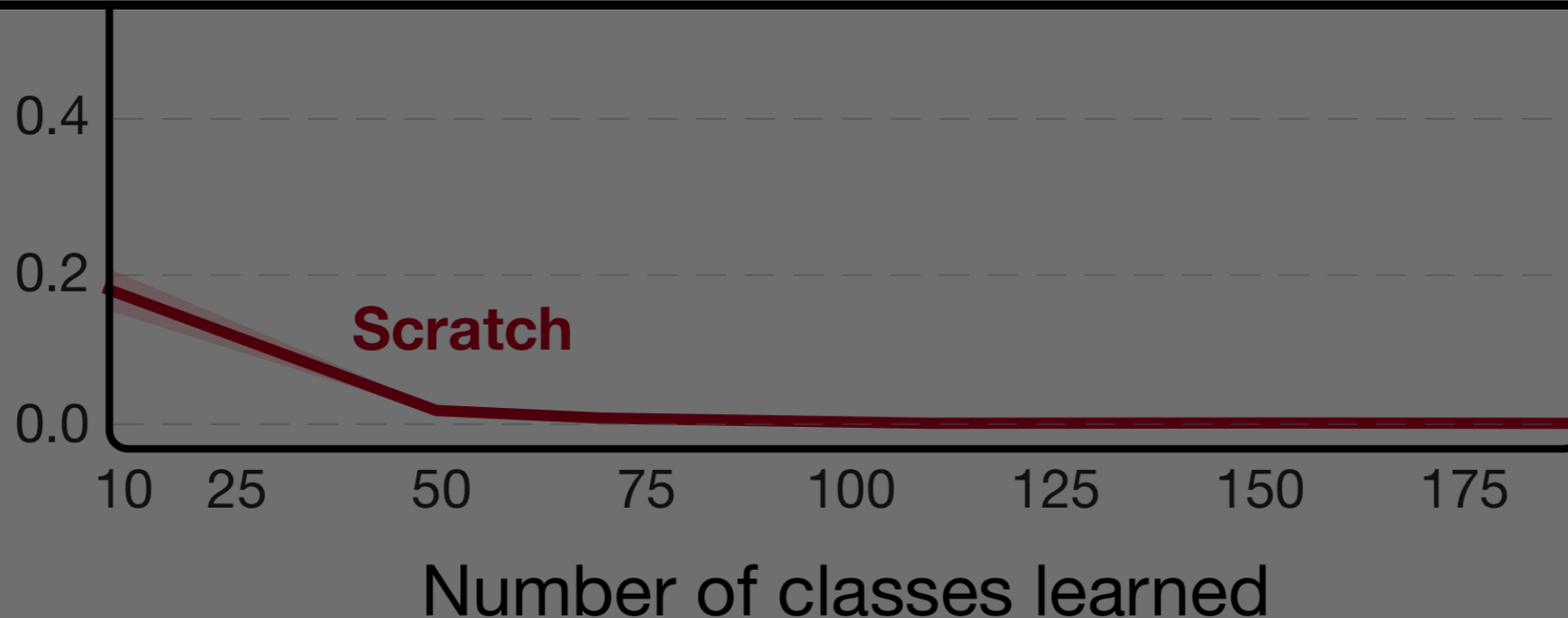
Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Is “Scratch” a fair baseline?

- Scratch is not longer a fair baseline
- Need a baseline that uses the information from the representation learning dataset
- **New baseline** : Learn a 950 way classifier using iid sampling and multiple epochs

Accuracy

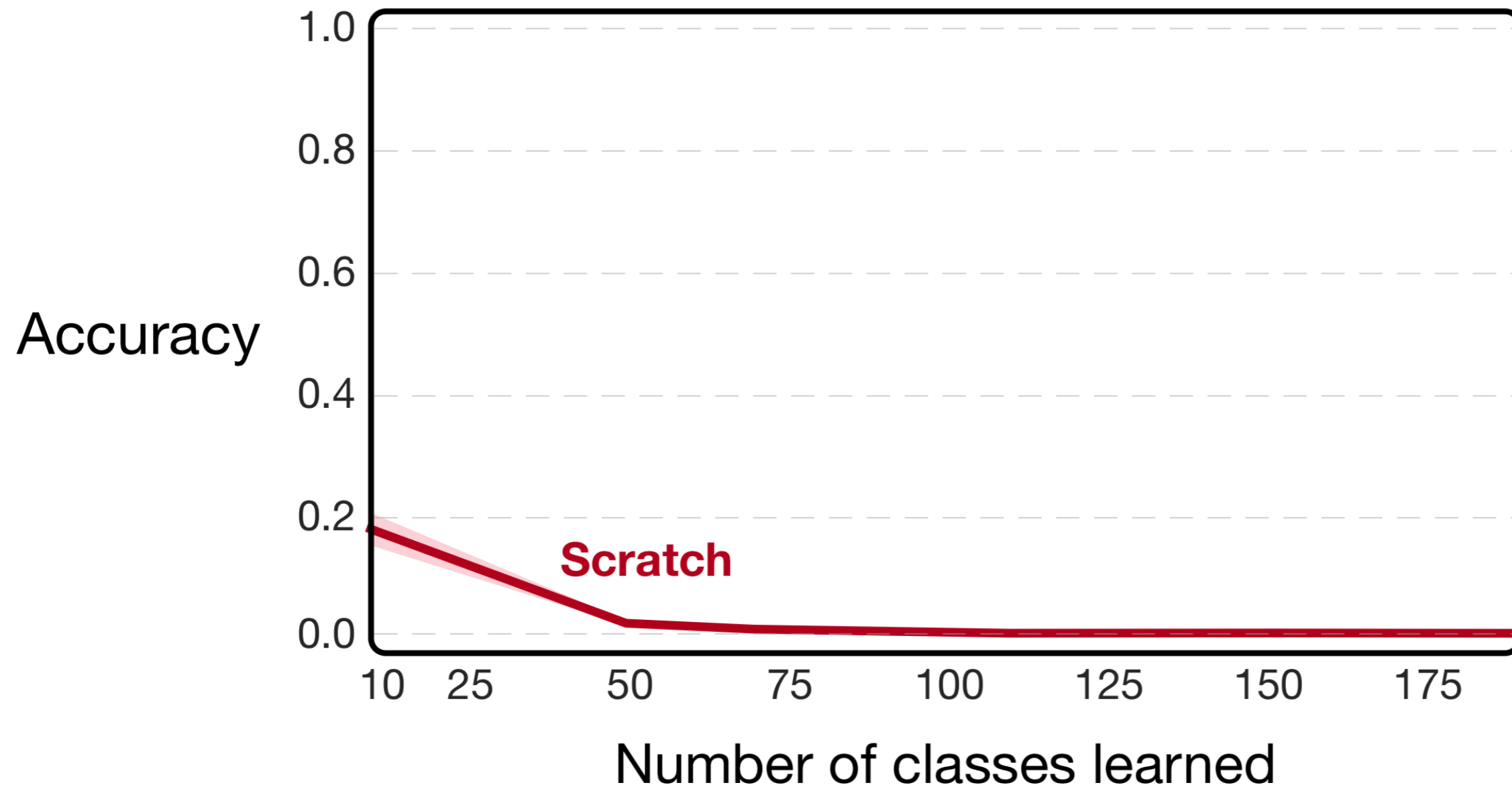


Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Omniglot Training Trajectory Performance

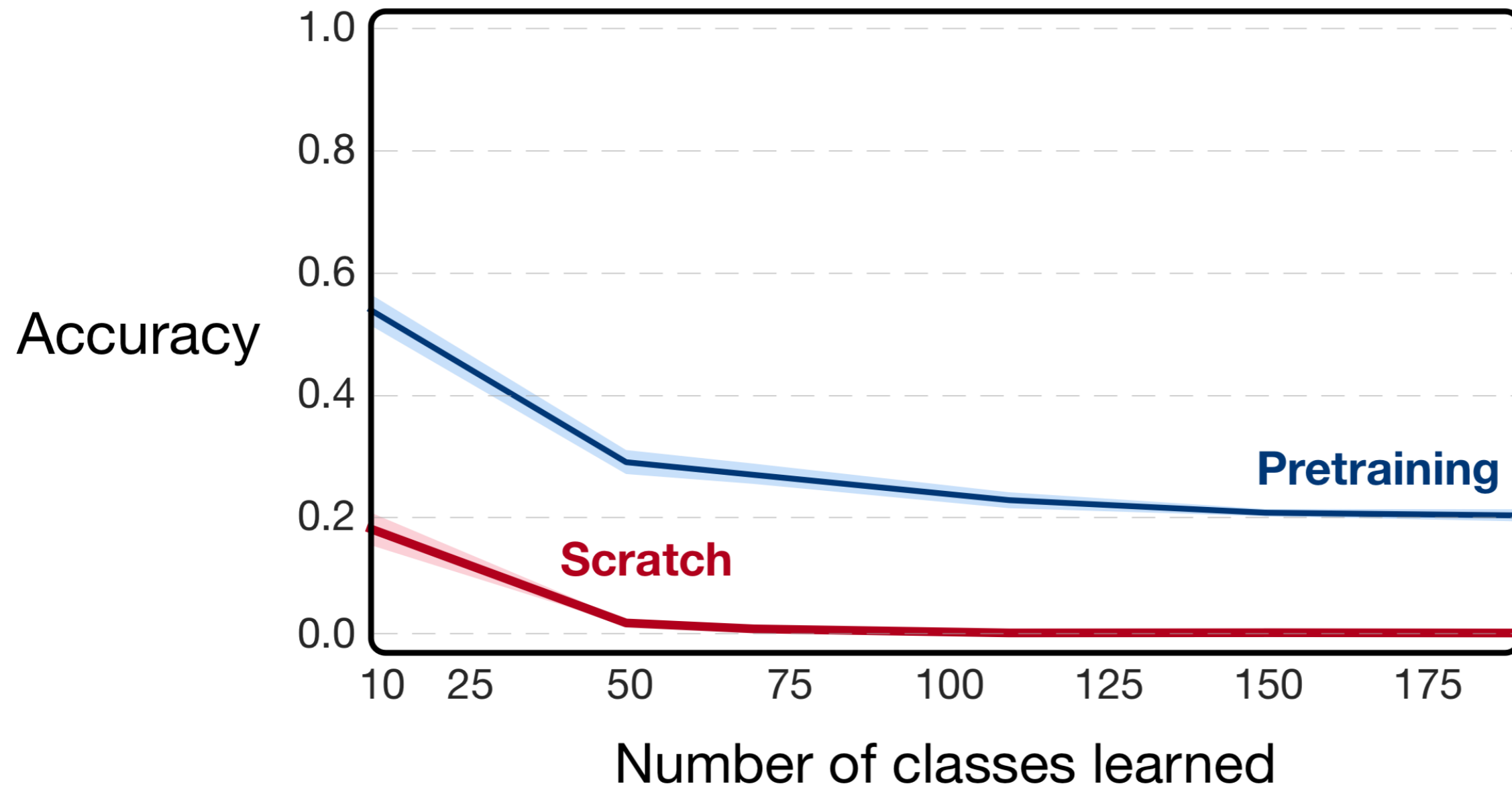


Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

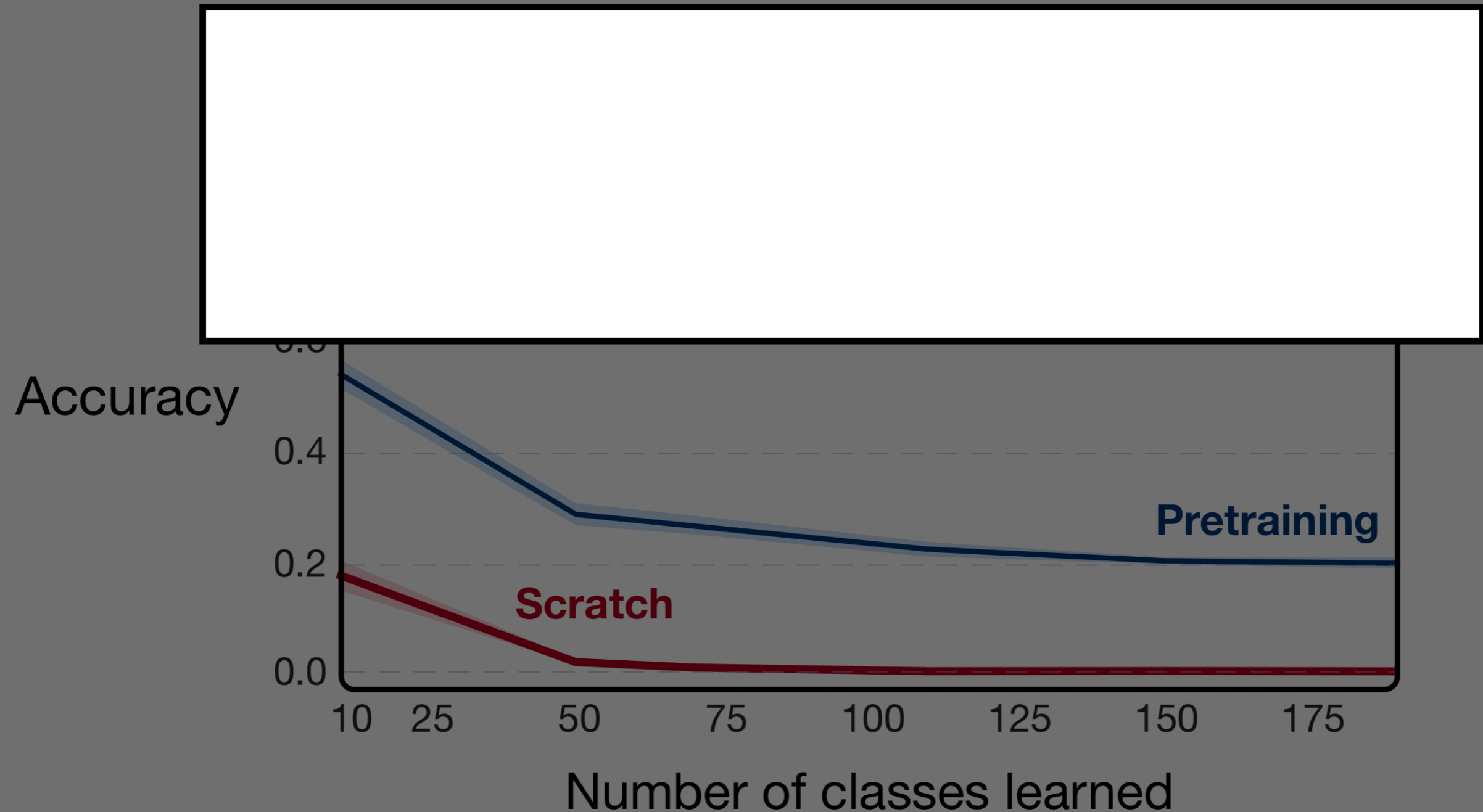
Omniglot Training Trajectory Performance



Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

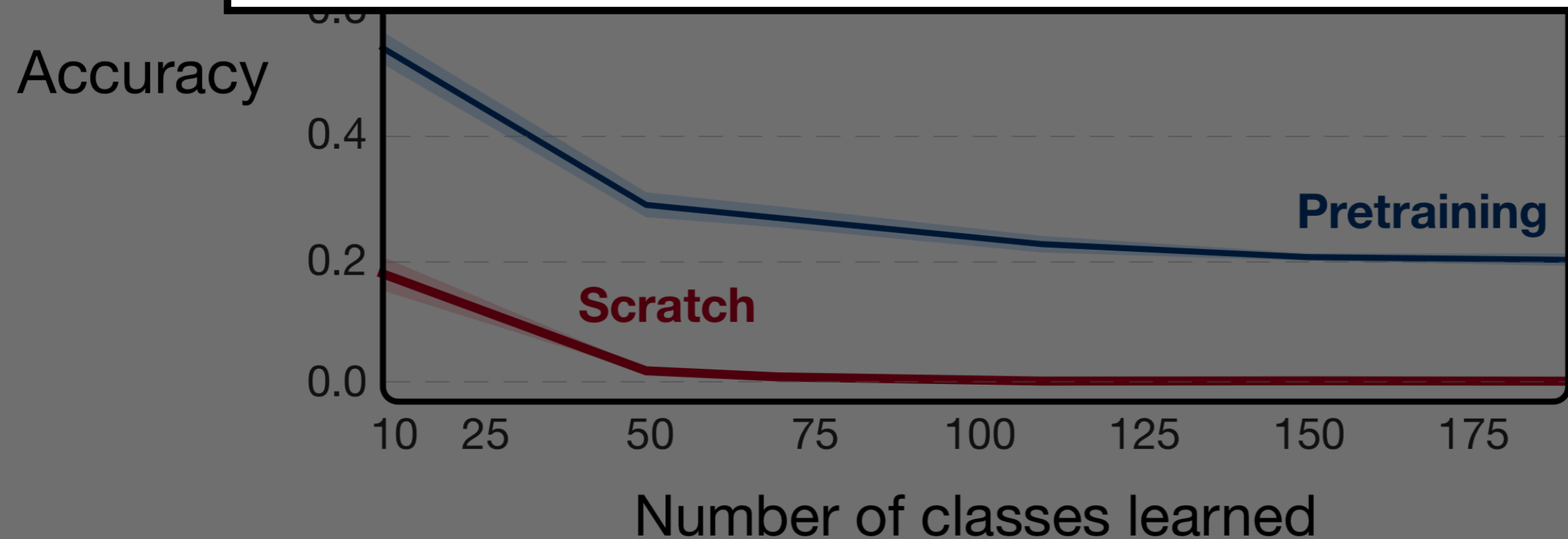


Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Sparse Representations



Dense Inputs

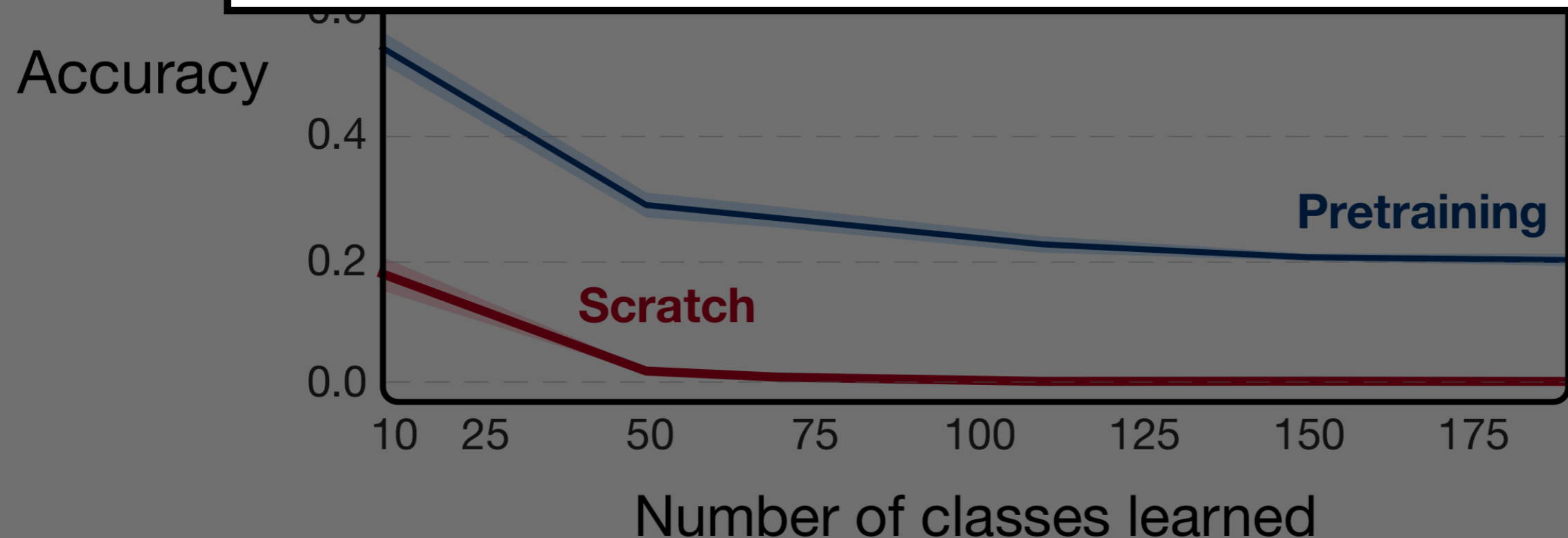
Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Sparse Representations

- We used SR-NN introduced by Vincent *et al.* [1]

[1] Liu, V., Kumaraswamy, R., Le, L., & White, M. (2018). The utility of sparse representations for control in reinforcement learning. AAI 19



Dense Inputs

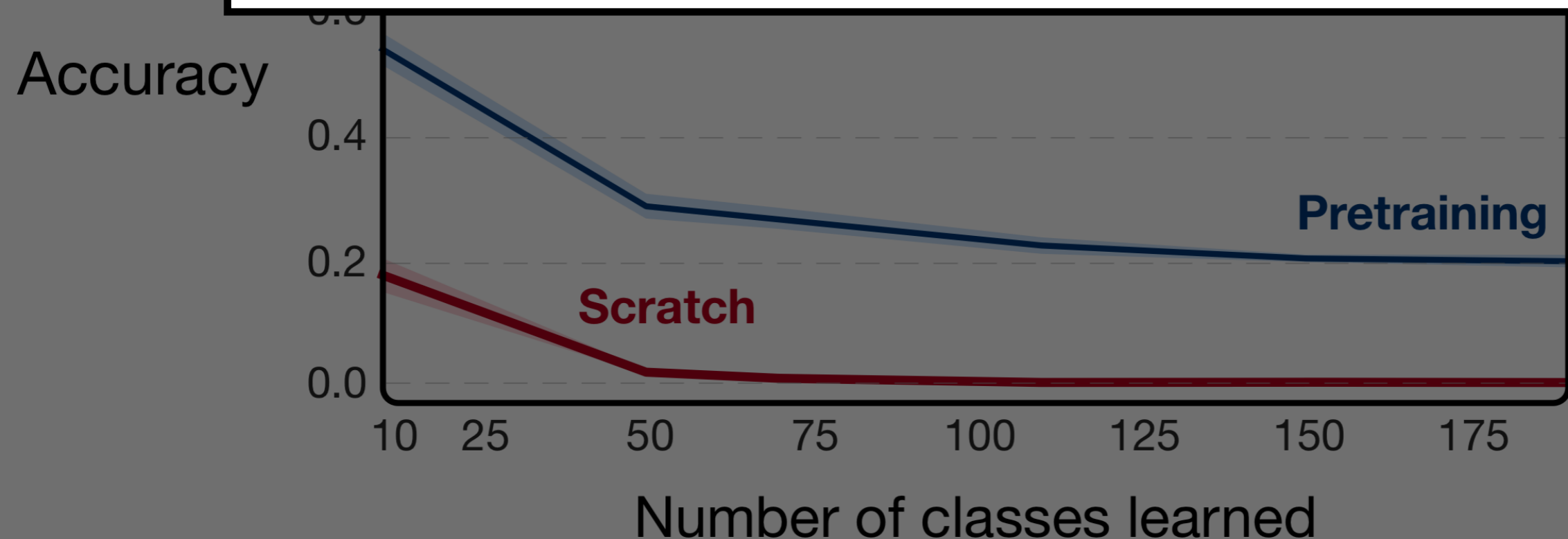
Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Sparse Representations

- We used SR-NN introduced by Vincent *et al.* [1]
- Results reported using the best performing SR-NN

[1] Liu, V., Kumaraswamy, R., Le, L., & White, M. (2018). The utility of sparse representations for control in reinforcement learning. AAI 19

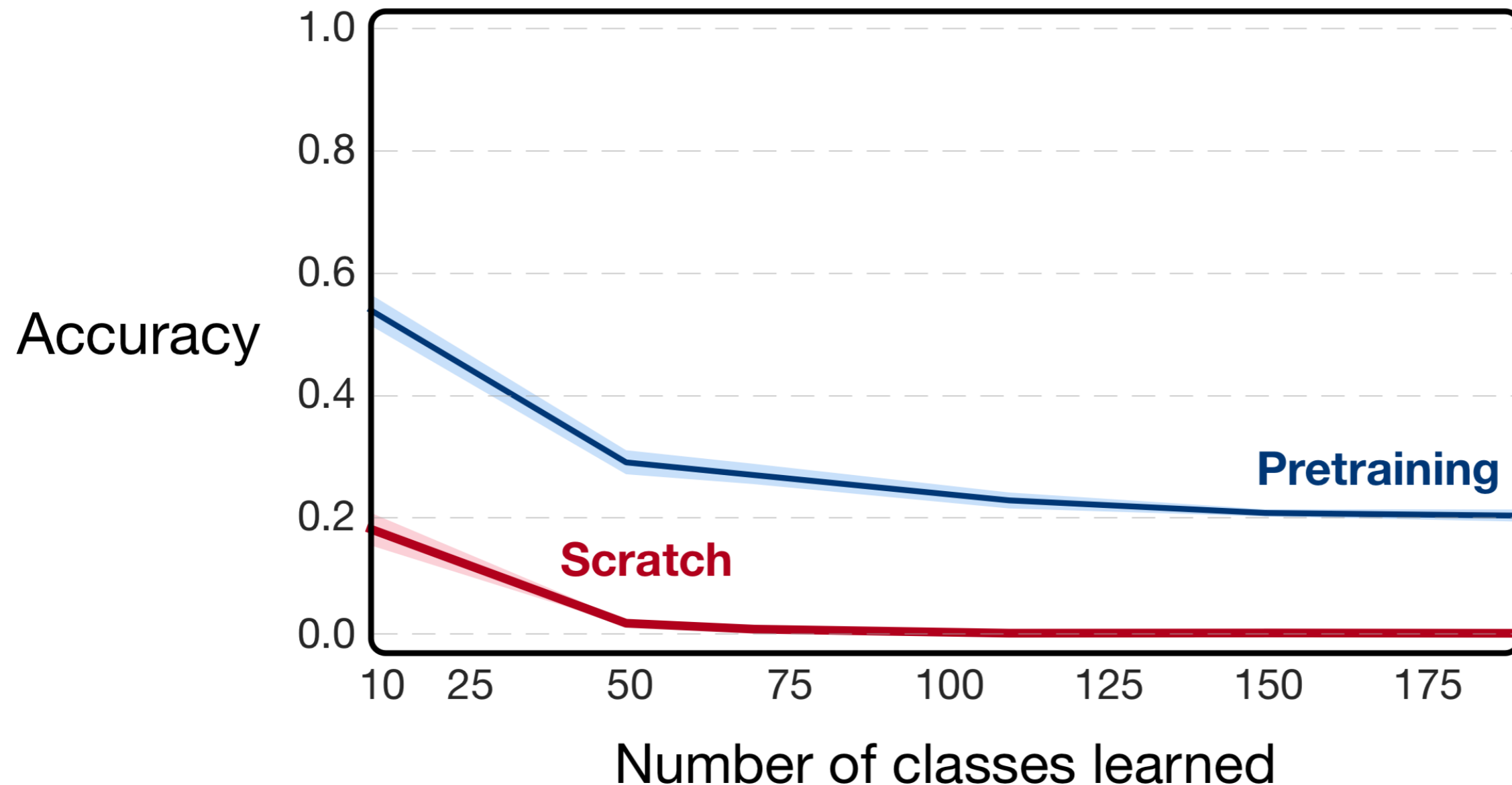


Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Omniglot Training Trajectory Performance

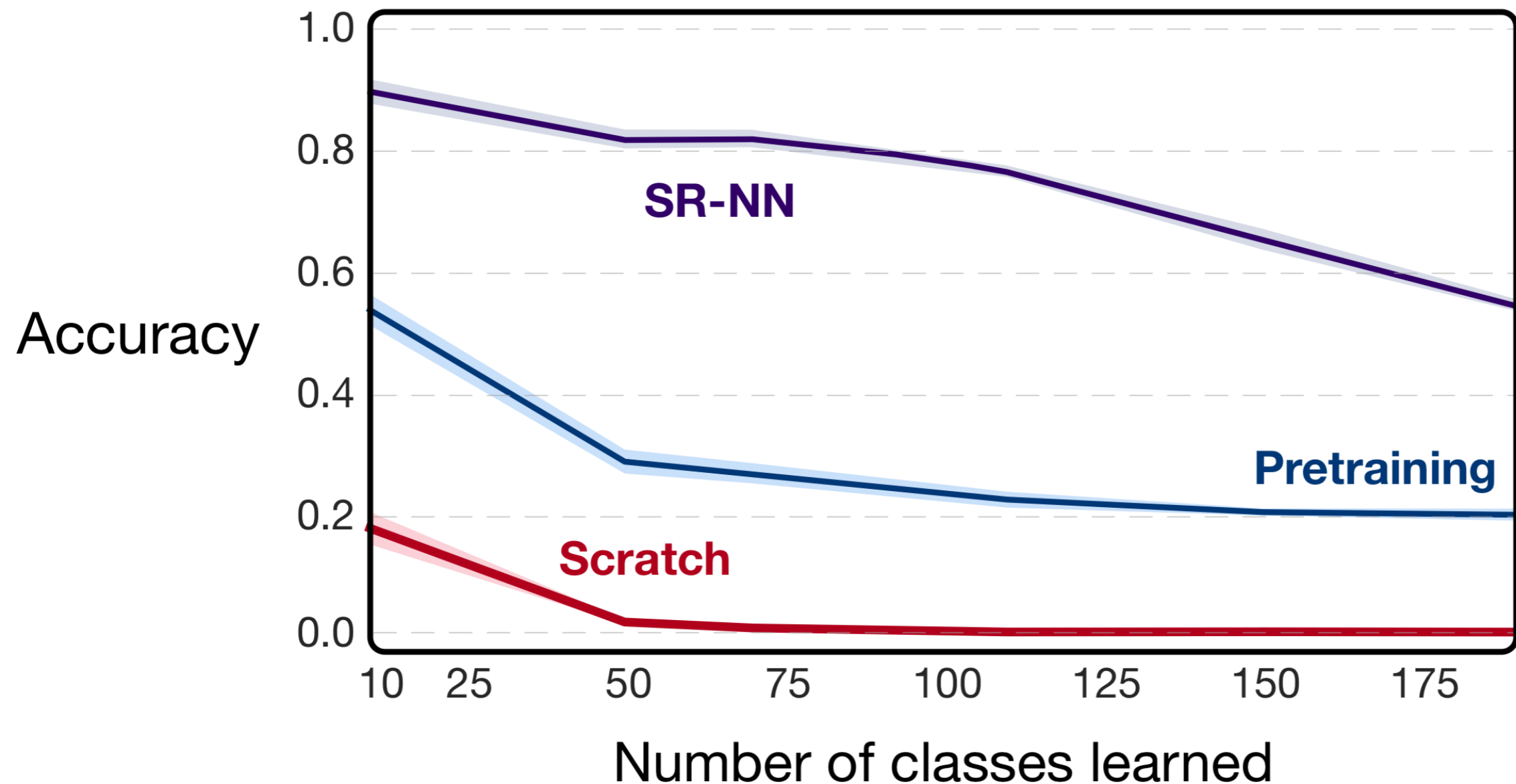


Dense Inputs

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Omniglot Training Trajectory Performance



Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling

2. Dense inputs

3. Global and greedy update

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling
2. Dense inputs
3. Global and greedy update

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling
2. Dense inputs
3. Global and greedy update

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling
2. Dense inputs
3. Global and greedy update

Experience Replay

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling
2. Dense inputs
3. Global and greedy update

Experience Replay

Tile coding

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling

Experience Replay

2. Dense inputs

Tile coding

3. Global and greedy update

Adam optimizer?

Catastrophic Interference

Neural networks suffer from *catastrophic* interference

Consequences : *Forgetting* and *slow learning*

When and why?

1. Non-IID sampling

2. Dense inputs

3. Global and greedy update

Our work



Experience Replay

Tile coding

Adam optimizer?

Meta-Learning Representations for Continual Learning

Meta-Learning Representations for Continual Learning

Core idea

Meta-Learning Representations for Continual Learning

Core idea

- Don't use sparsity as a proxy to a good representation

Meta-Learning Representations for Continual Learning

Core idea

- Don't use sparsity as a proxy to a good representation
- Directly optimize for representations which allow for continual learning

Meta-Learning Representations for Continual Learning

Meta-Learning Representations for Continual Learning

Data stream

$$\mathcal{D} = (X_0, Y_0), (X_1, Y_1), \dots, (X_k, Y_k), \dots, (X_n, Y_n), \dots,$$

Meta-Learning Representations for Continual Learning

Data stream

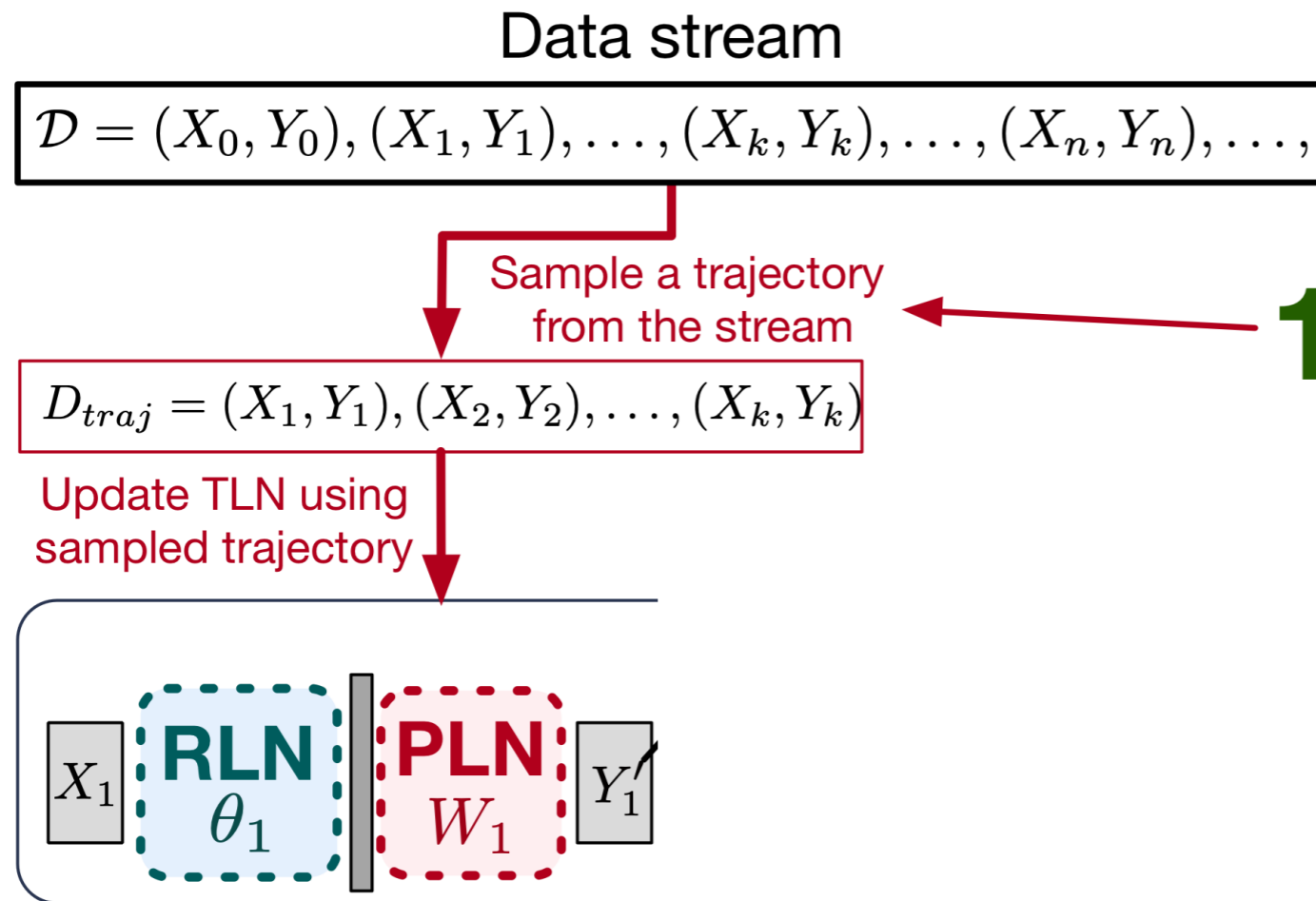
$$\mathcal{D} = (X_0, Y_0), (X_1, Y_1), \dots, (X_k, Y_k), \dots, (X_n, Y_n), \dots,$$

Sample a trajectory
from the stream

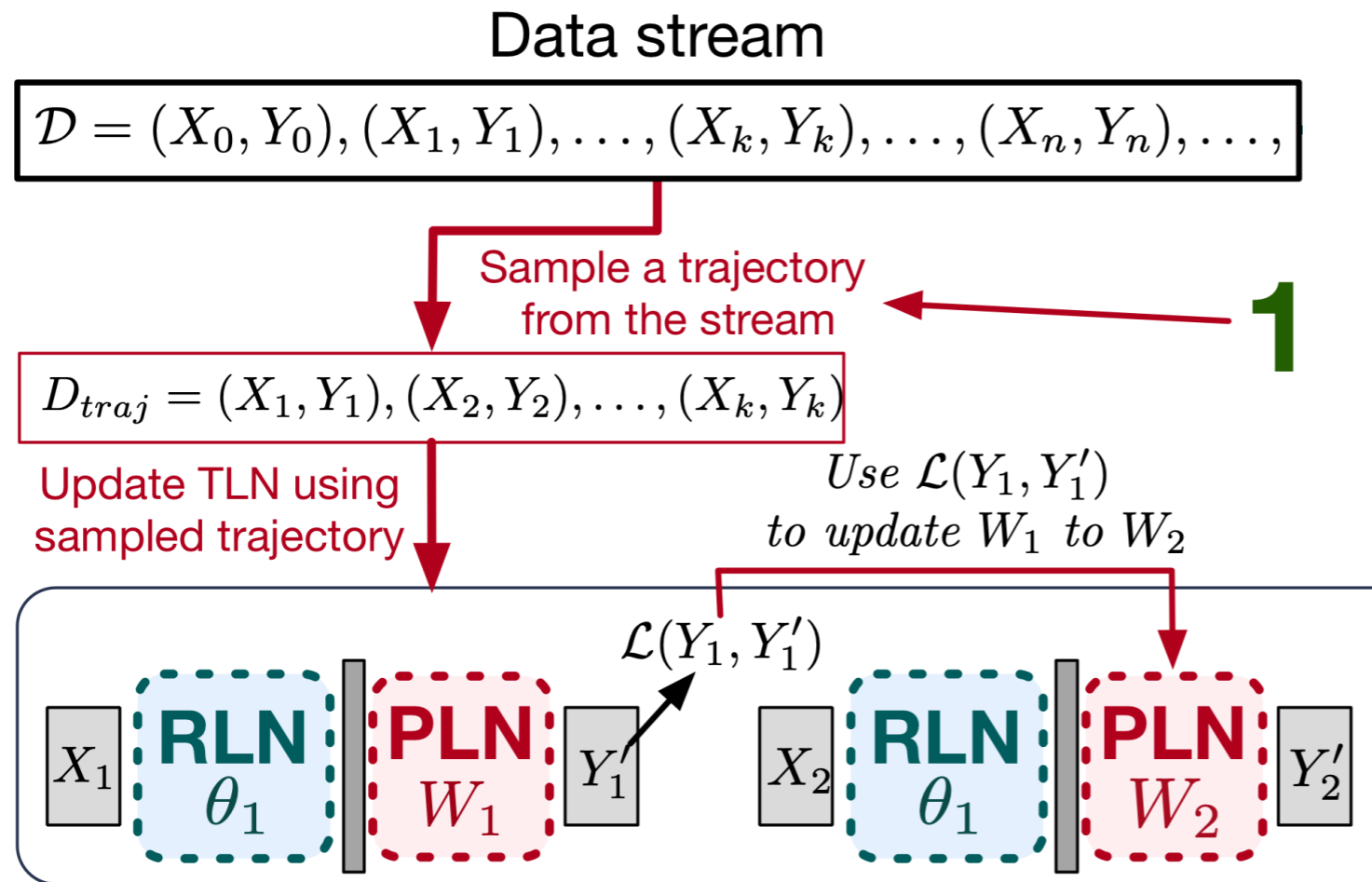
$$D_{traj} = (X_1, Y_1), (X_2, Y_2), \dots, (X_k, Y_k)$$

1

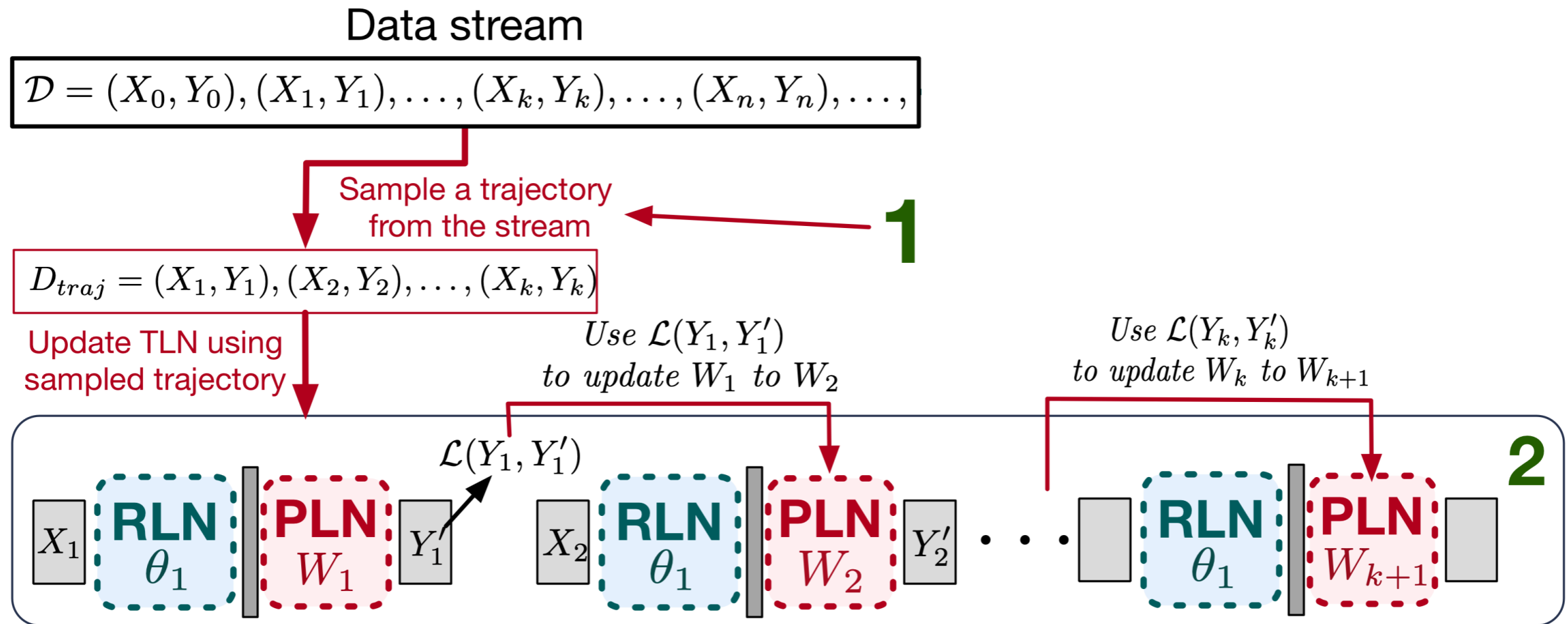
Meta-Learning Representations for Continual Learning



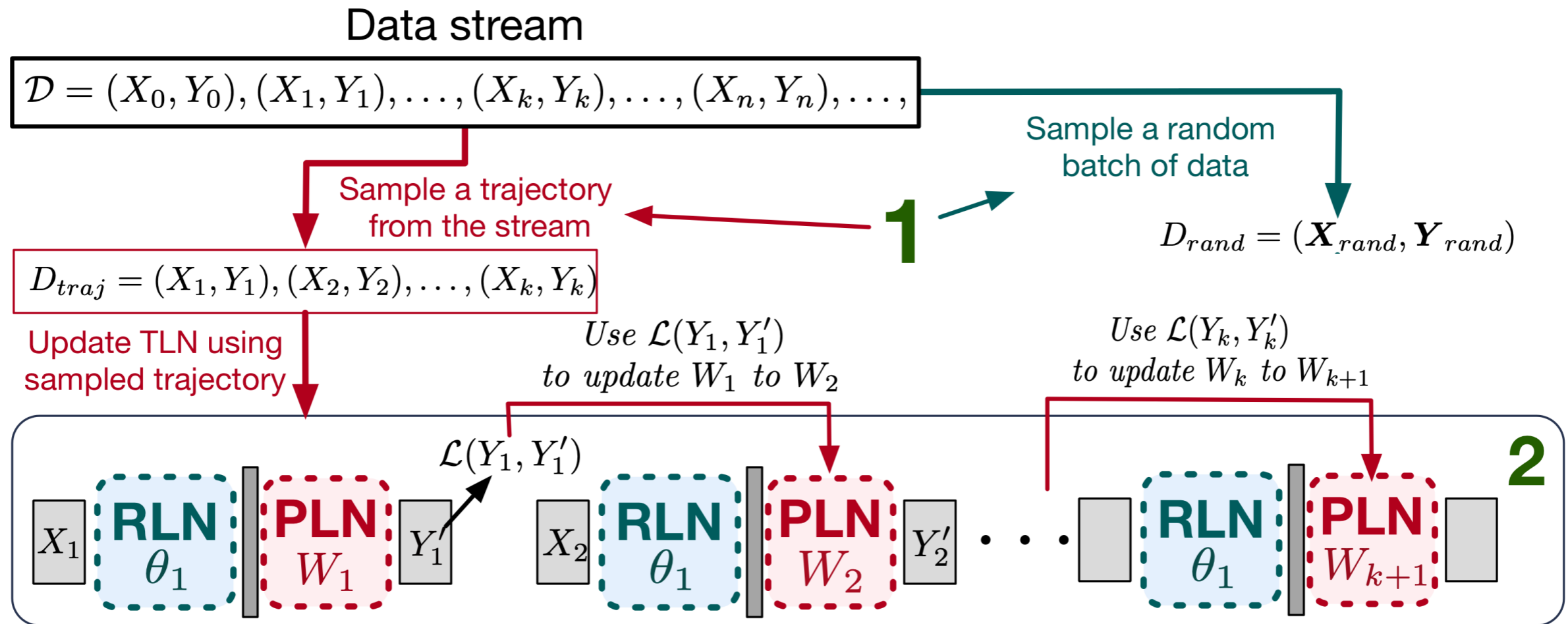
Meta-Learning Representations for Continual Learning



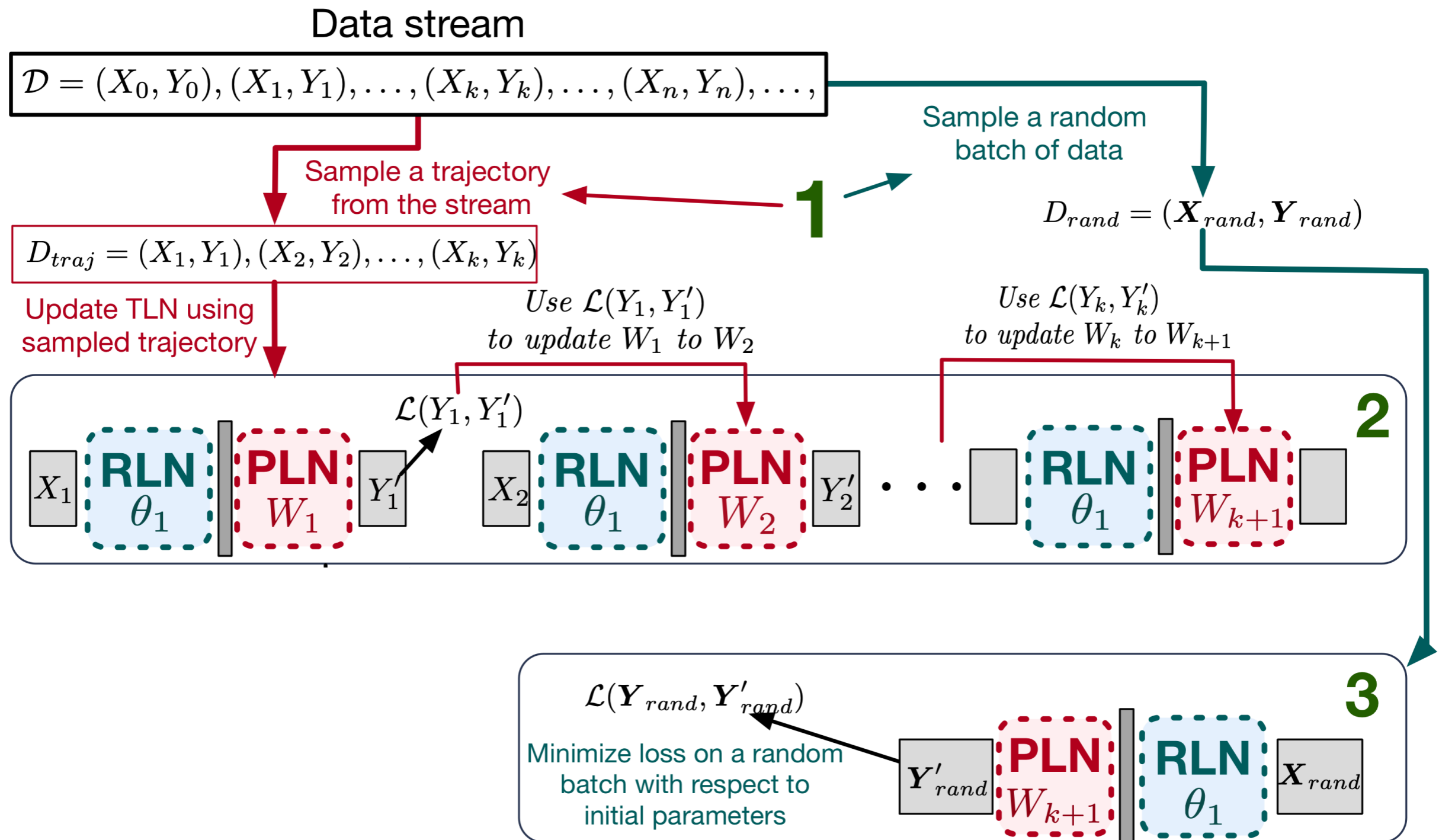
Meta-Learning Representations for Continual Learning



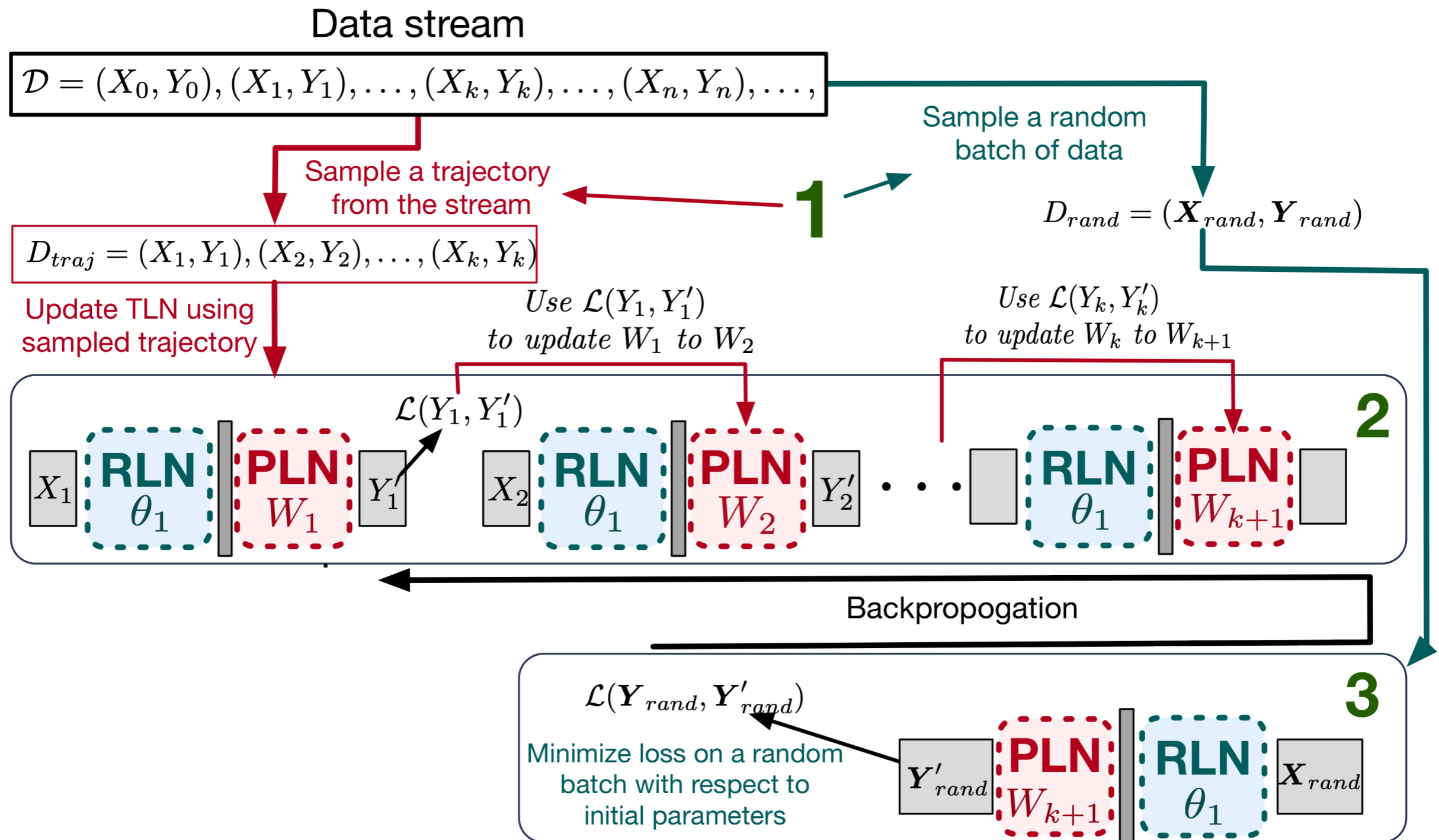
Meta-Learning Representations for Continual Learning



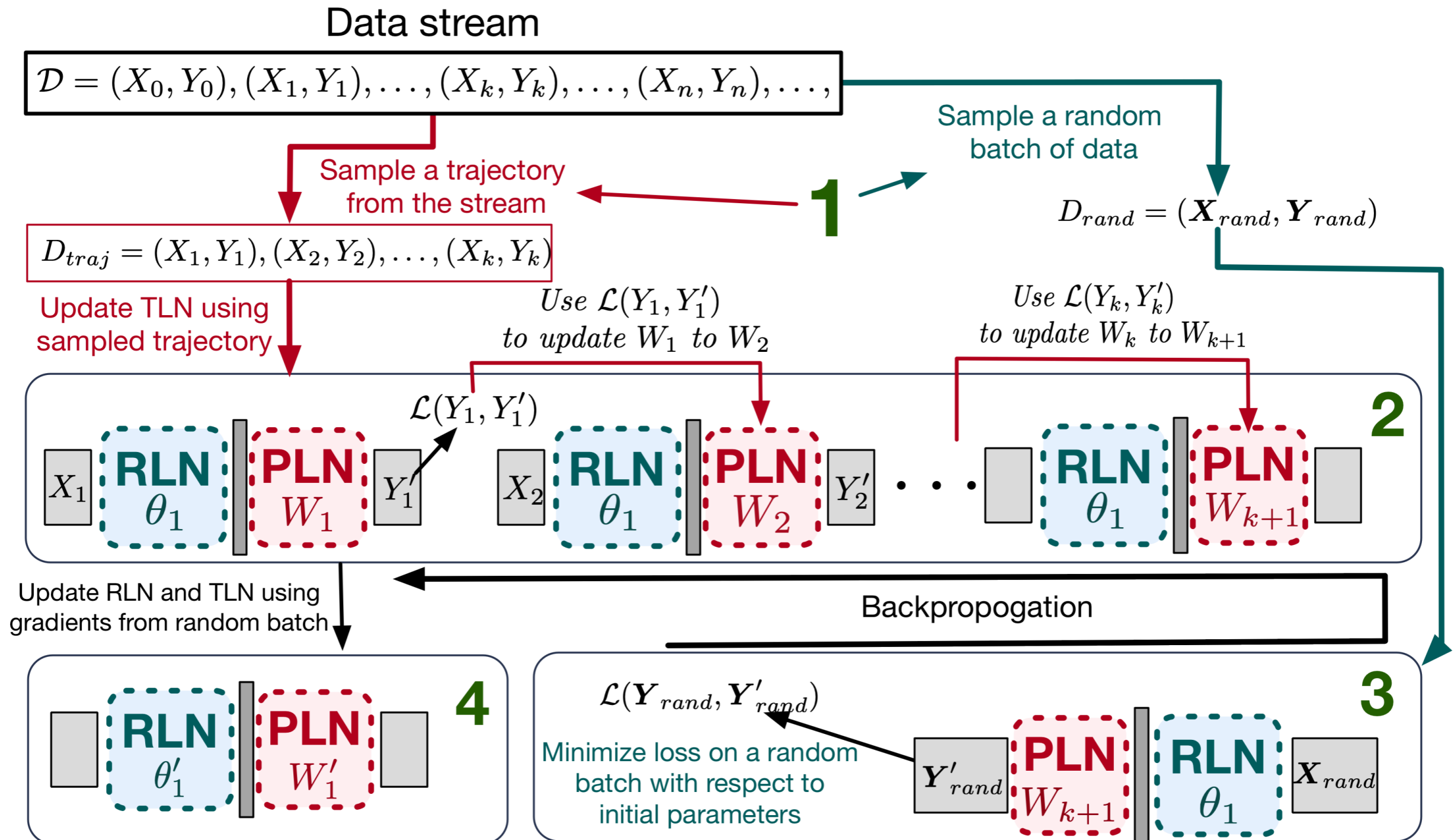
Meta-Learning Representations for Continual Learning



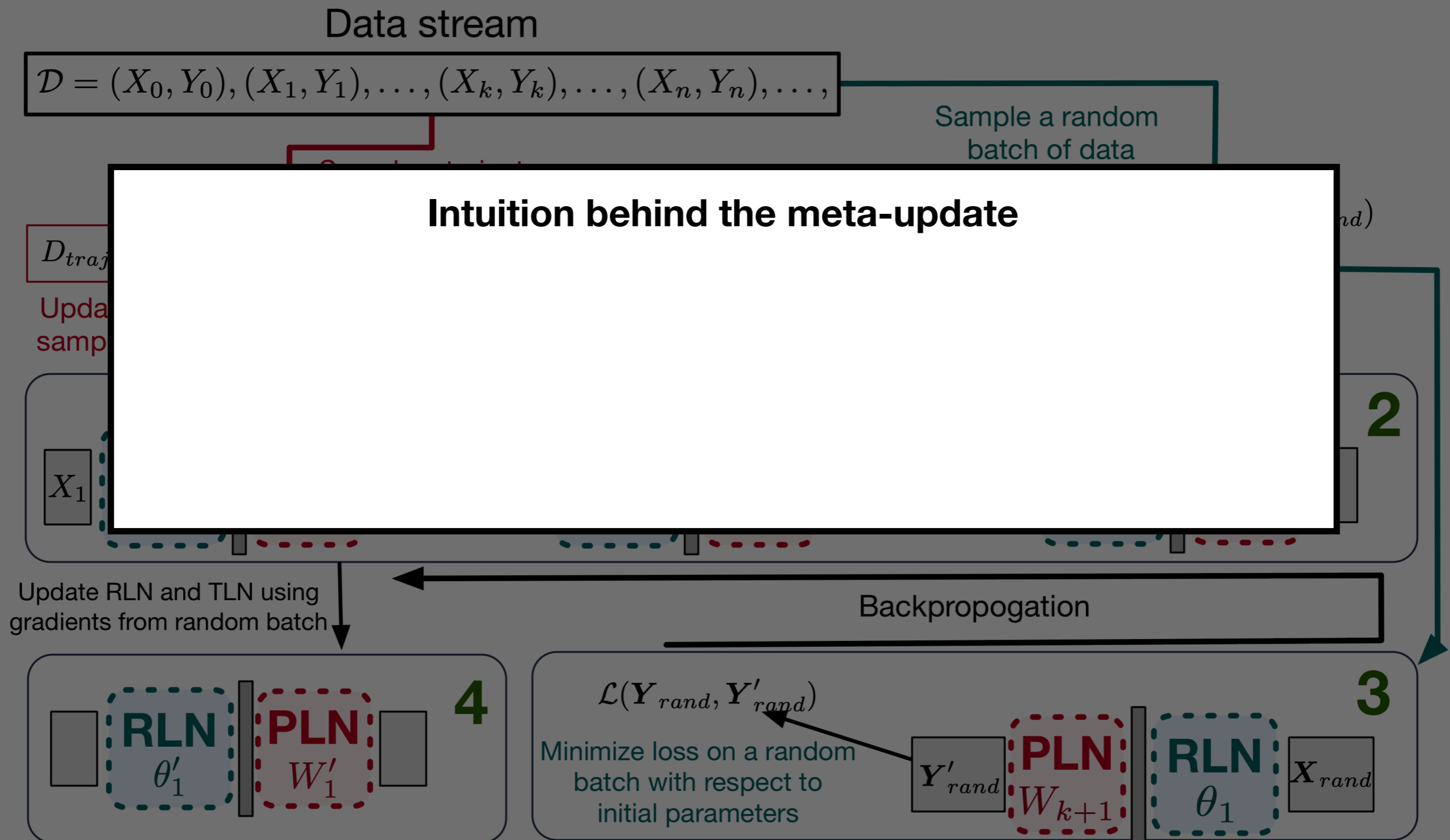
Meta-Learning Representations for Continual Learning



Meta-Learning Representations for Continual Learning



Meta-Learning Representations for Continual Learning



Meta-Learning Representations for Continual Learning

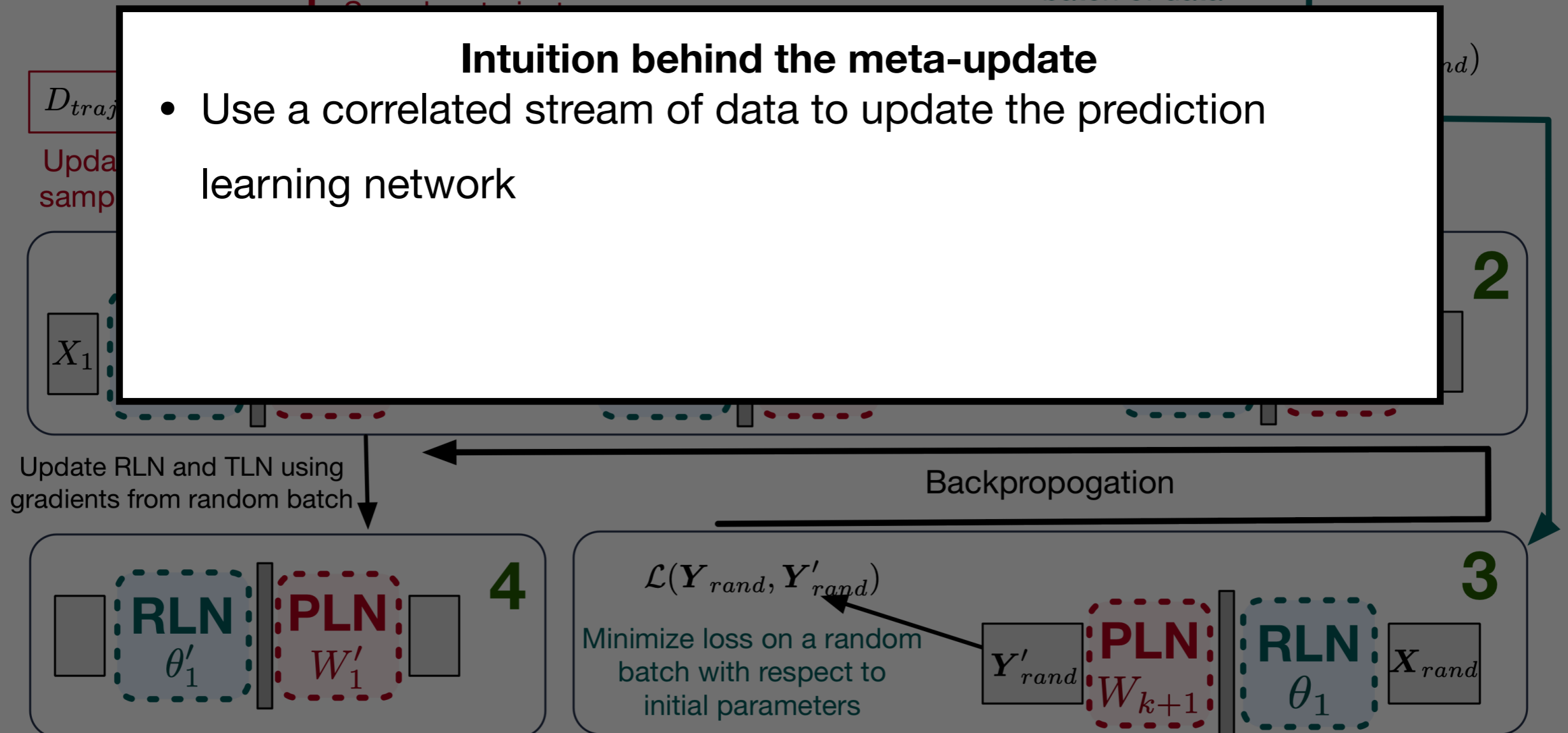
Data stream

$$\mathcal{D} = (X_0, Y_0), (X_1, Y_1), \dots, (X_k, Y_k), \dots, (X_n, Y_n), \dots,$$

Sample a random batch of data

Intuition behind the meta-update

- Use a correlated stream of data to update the prediction learning network



Meta-Learning Representations for Continual Learning

Data stream

$$\mathcal{D} = (X_0, Y_0), (X_1, Y_1), \dots, (X_k, Y_k), \dots, (X_n, Y_n), \dots,$$

Sample a random batch of data

Intuition behind the meta-update

- Use a correlated stream of data to update the prediction learning network
- Compute the degree of interference/forgetting

D_{traj}

Update samp

X_1

Update RLN and TLN using gradients from random batch

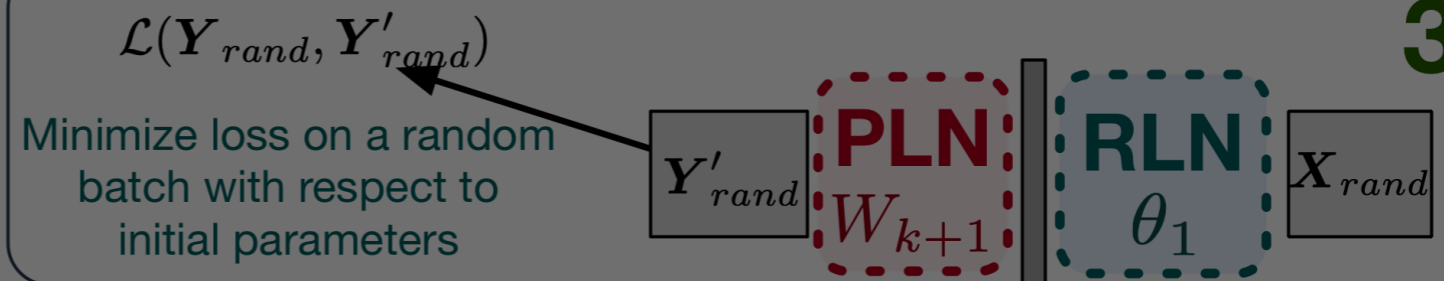
Backpropogation

2

4



3



Meta-Learning Representations for Continual Learning

Data stream

$$\mathcal{D} = (X_0, Y_0), (X_1, Y_1), \dots, (X_k, Y_k), \dots, (X_n, Y_n), \dots,$$

Sample a random batch of data

Intuition behind the meta-update

- Use a correlated stream of data to update the prediction learning network
- Compute the degree of interference/forgetting
- Use interference as a training signal for updating representation

D_{traj}

Update samp

X_1

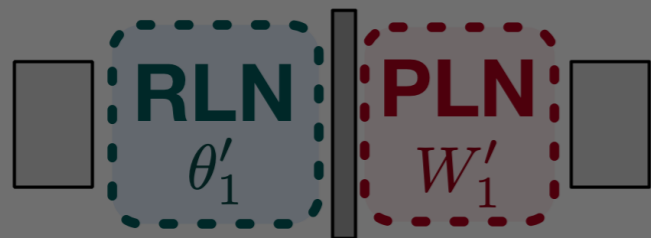
Update RLN and TLN using gradients from random batch

Backpropogation

2

3

4



$$\mathcal{L}(Y_{rand}, Y'_{rand})$$

Minimize loss on a random batch with respect to initial parameters

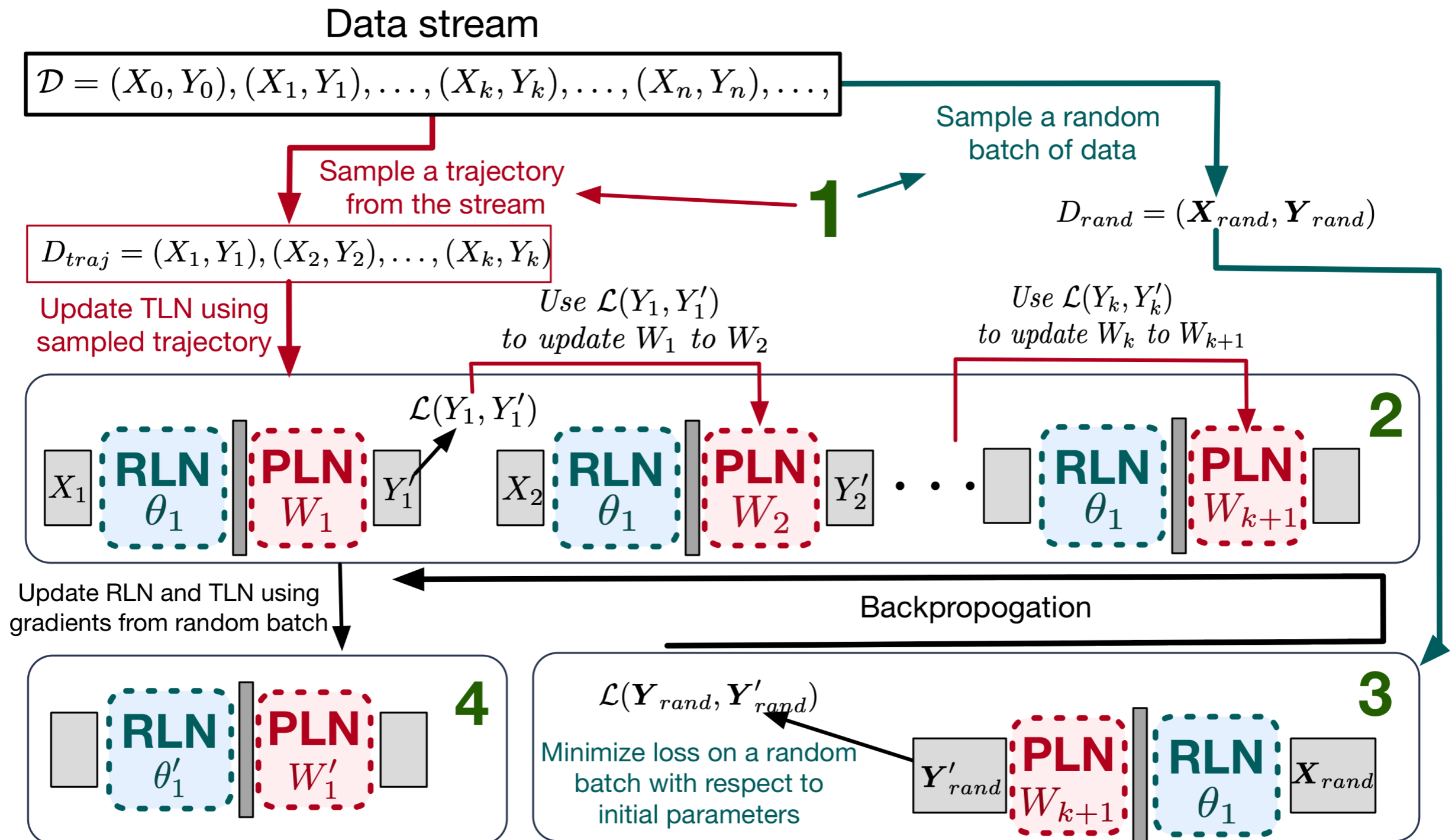
Y'_{rand}

PLN
 W_{k+1}

RLN
 θ_1

X_{rand}

Meta-Learning Representations for Continual Learning

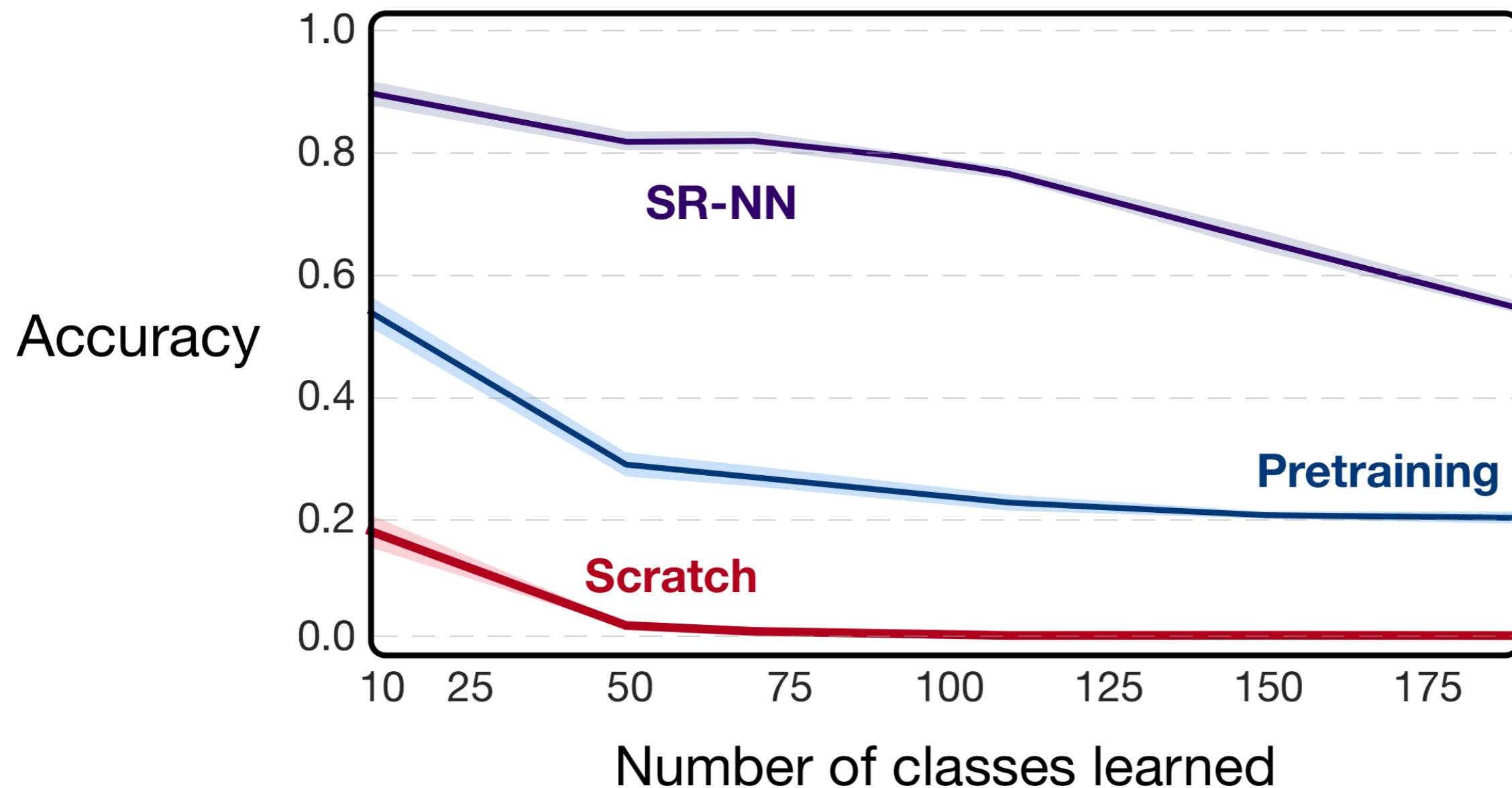


Meta-Learning Representations for Continual Learning

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Omniglot Training Trajectory Performance

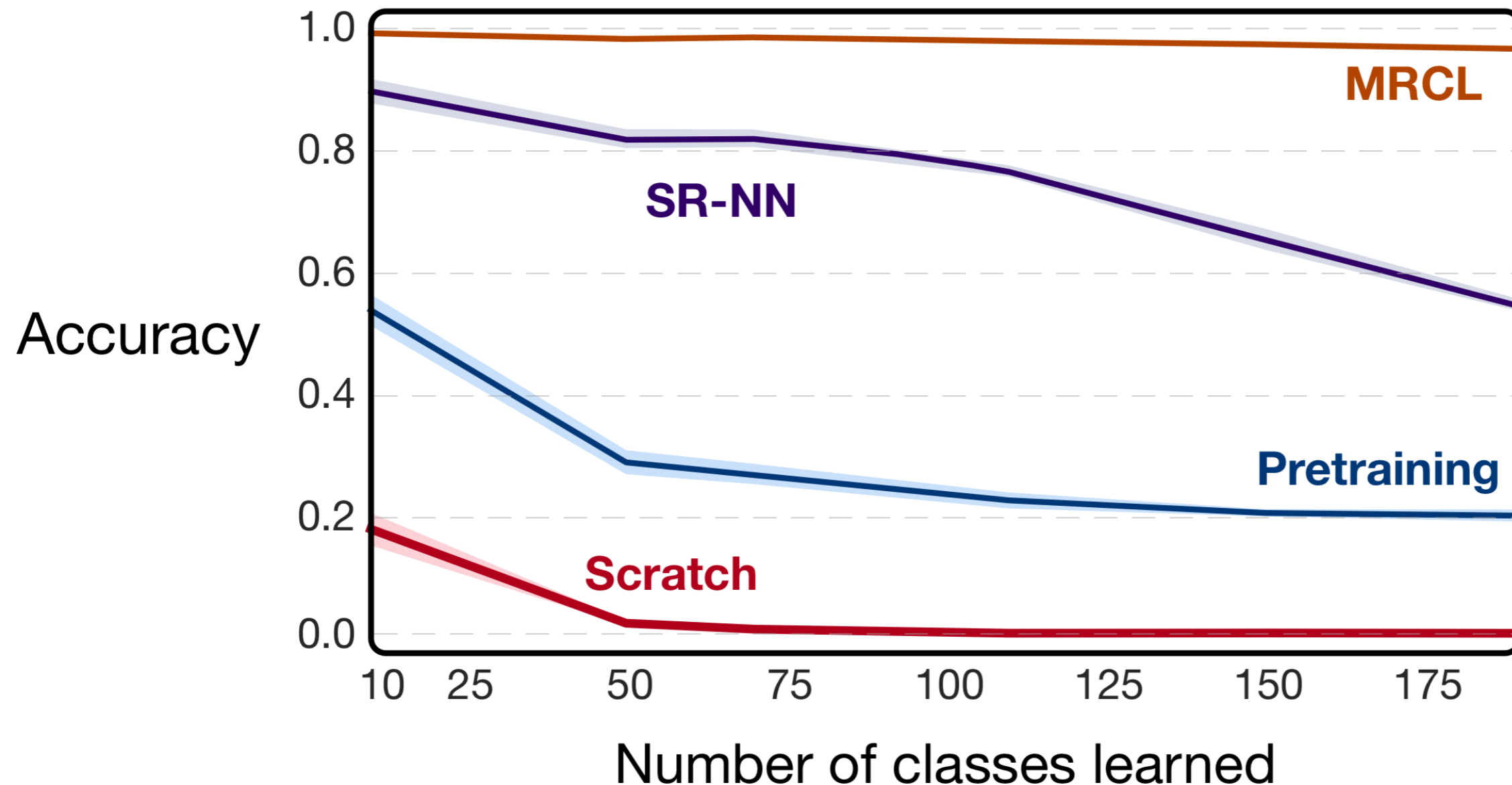


Meta-Learning Representations for Continual Learning

Dataset : Omniglot

- ~950 characters from multiple alphabets for representation learning
- ~600 characters from multiple alphabets for continual learning prediction

Omniglot Training Trajectory Performance



Analyzing representations

Analyzing representations

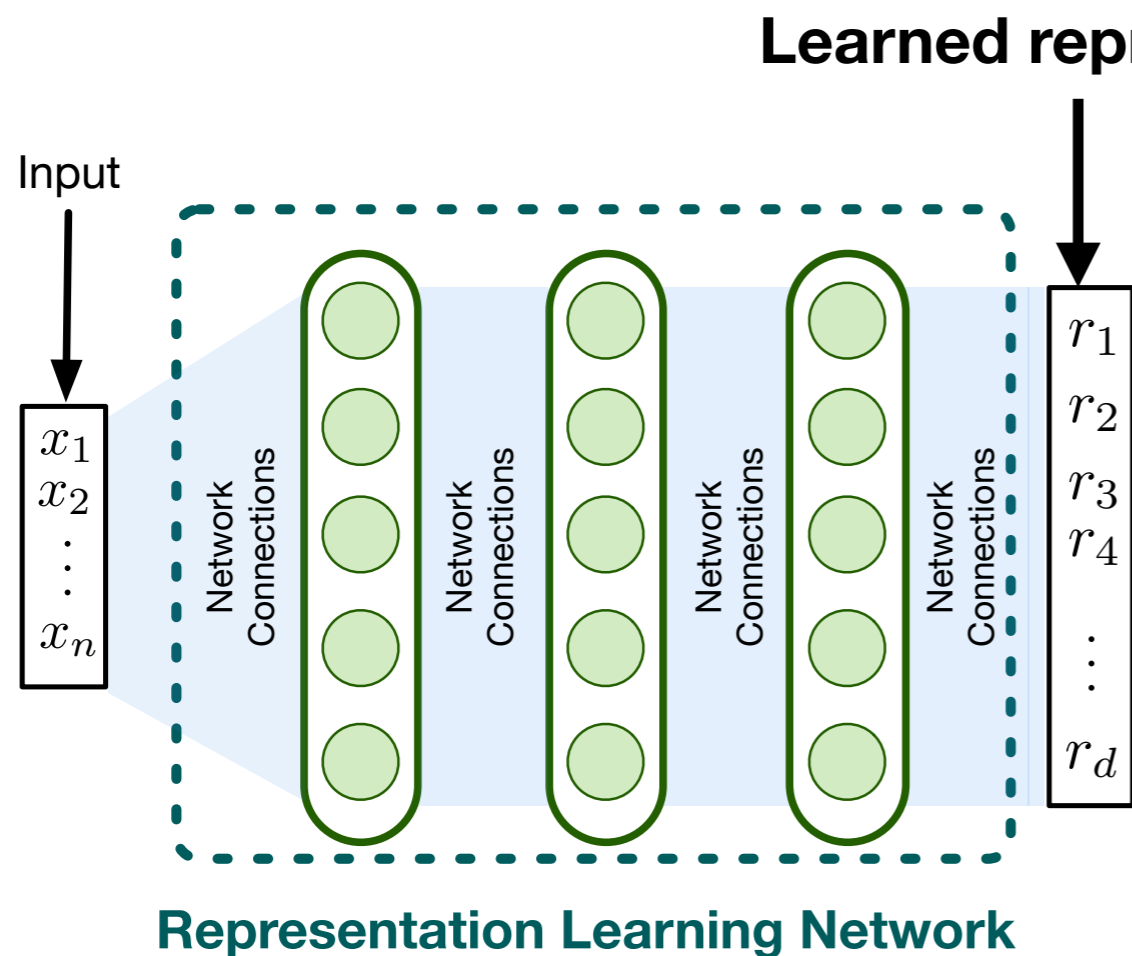
Instance sparsity

Percentage of non-zero entries used to represent an input on average

Analyzing representations

Instance sparsity

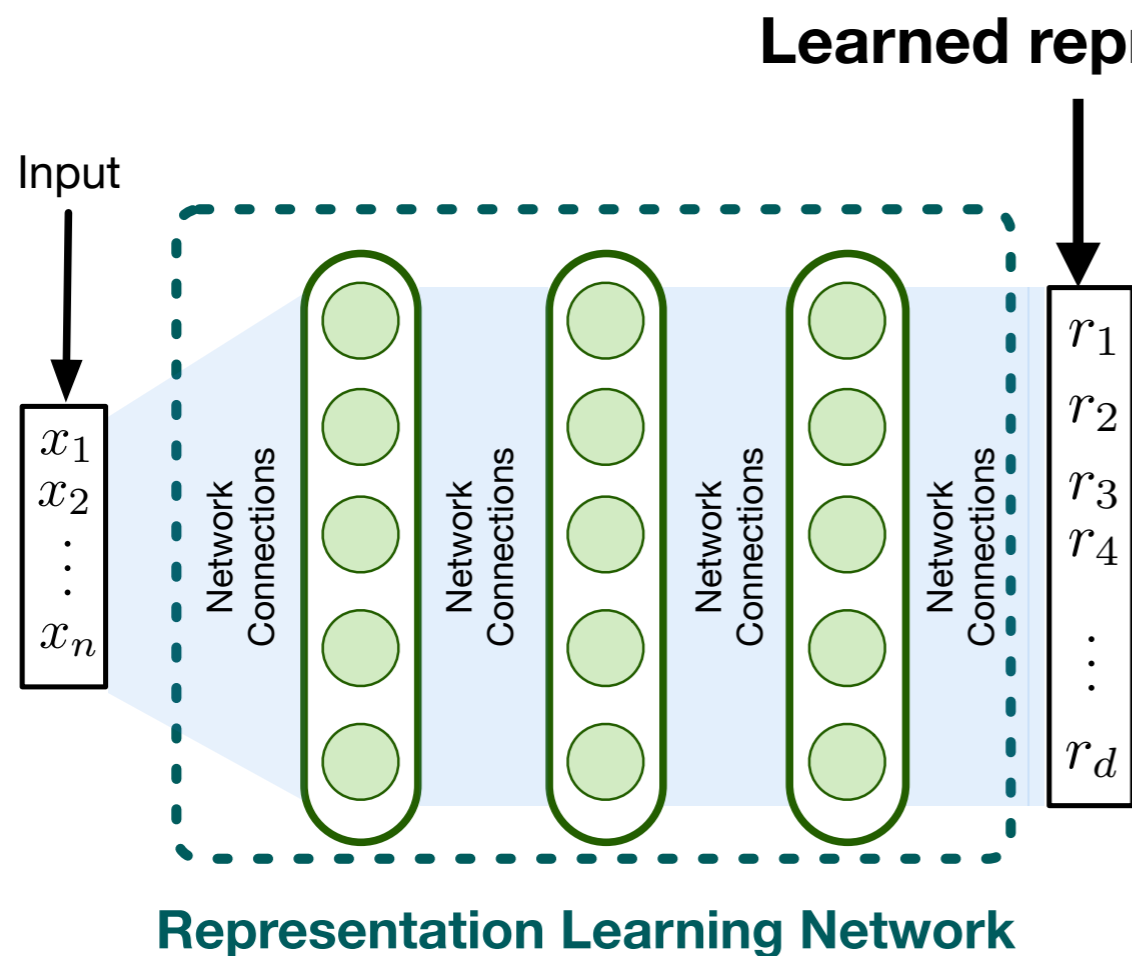
Percentage of non-zero entries used to represent an input on average



Analyzing representations

Instance sparsity

Percentage of non-zero entries used to represent an input on average

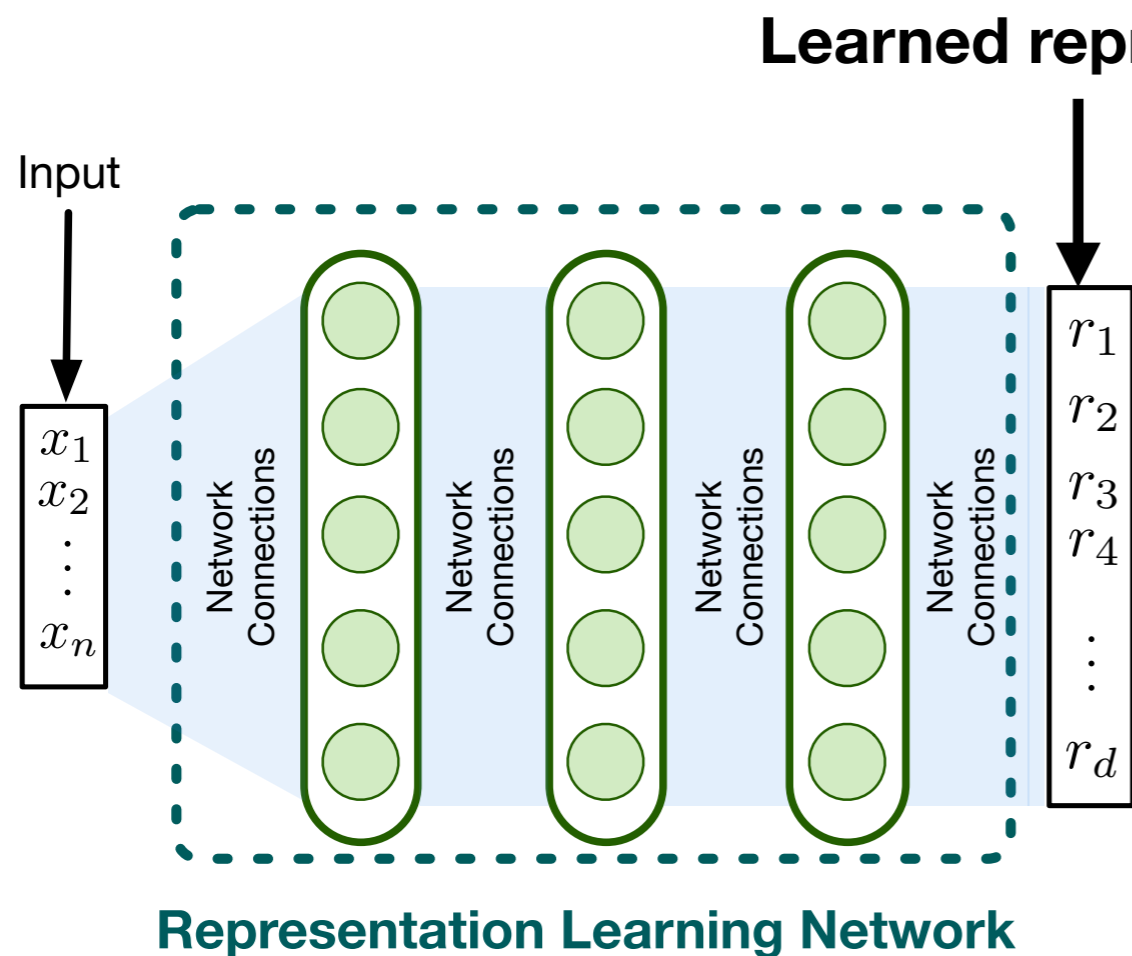


- **Pretraining: 38%**

Analyzing representations

Instance sparsity

Percentage of non-zero entries used to represent an input on average

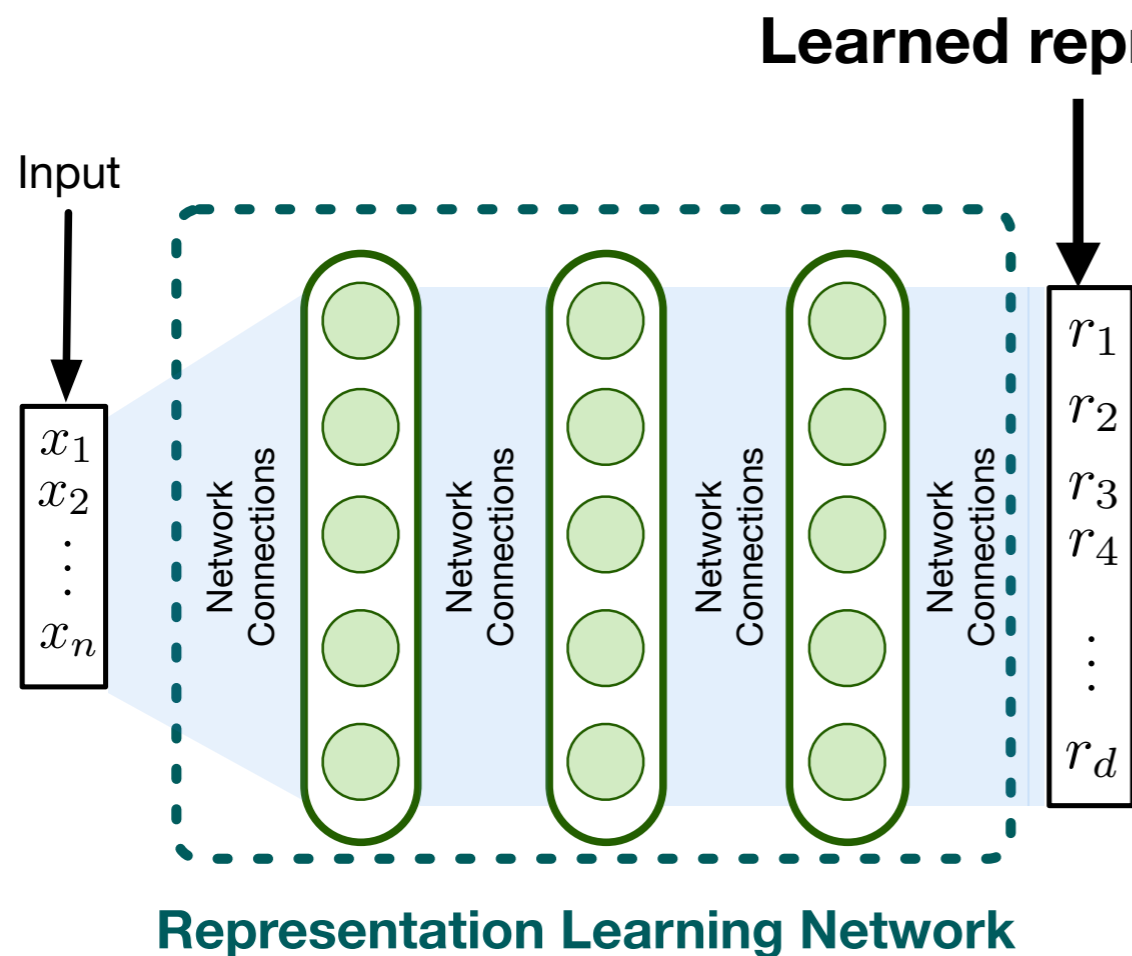


- **Pretraining: 38%**
- **SR-NN: 15%**

Analyzing representations

Instance sparsity

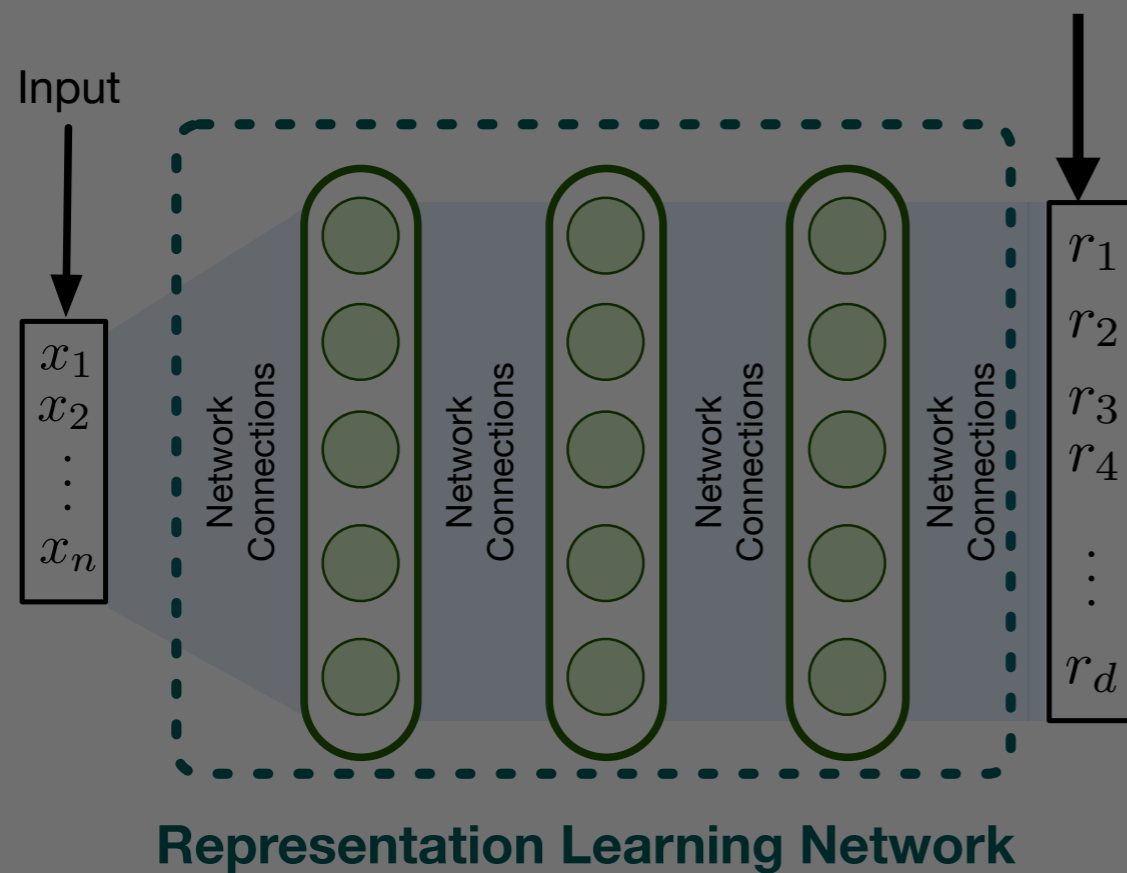
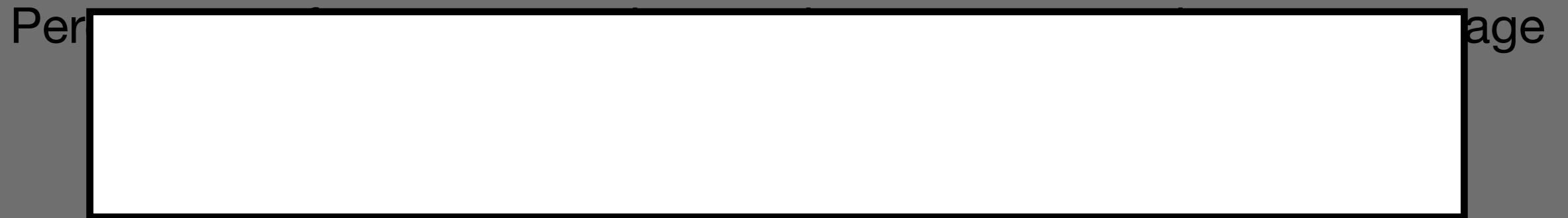
Percentage of non-zero entries used to represent an input on average



- **Pretraining:** 38%
- **SR-NN:** 15%
- **MRCL:** 3.8%

Analyzing representations

Instance sparsity



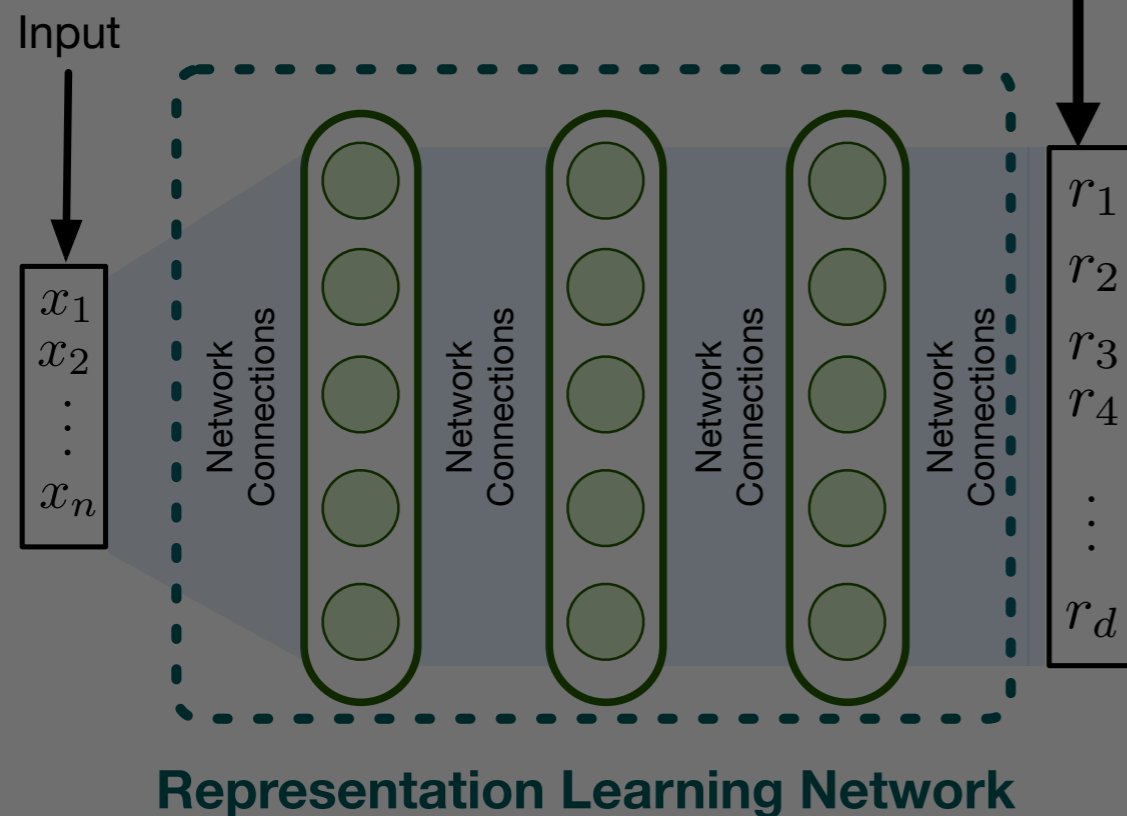
- **Pretraining: 38%**
- **SR-NN: 15%**
- **MRCL: 3.8%**

Analyzing representations

Instance sparsity

Per Page

We can train SR-NN with different hyper-parameter to achieve same degree of sparsity

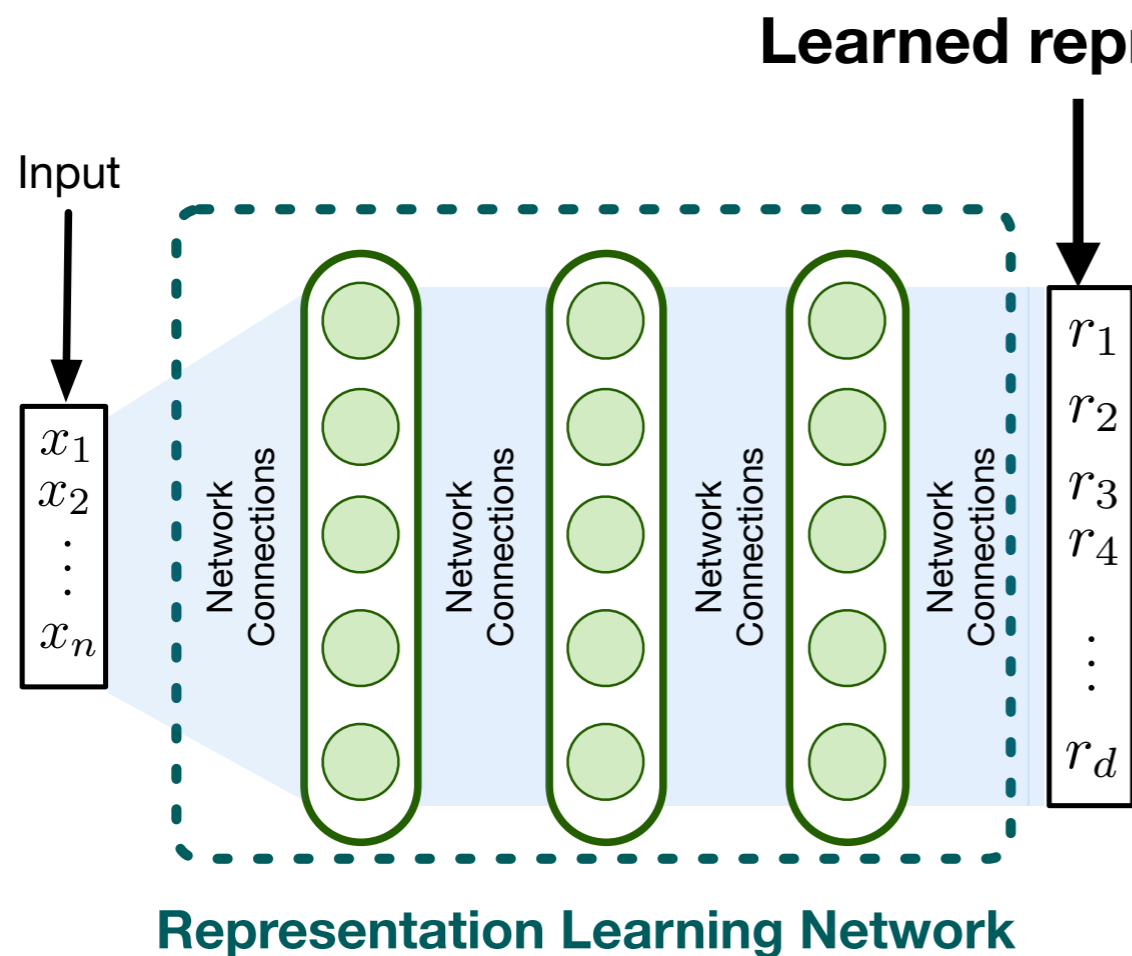


- **Pretraining:** 38%
- **SR-NN:** 15%
- **MRCL:** 3.8%

Analyzing representations

Instance sparsity

Percentage of non-zero entries used to represent an input on average

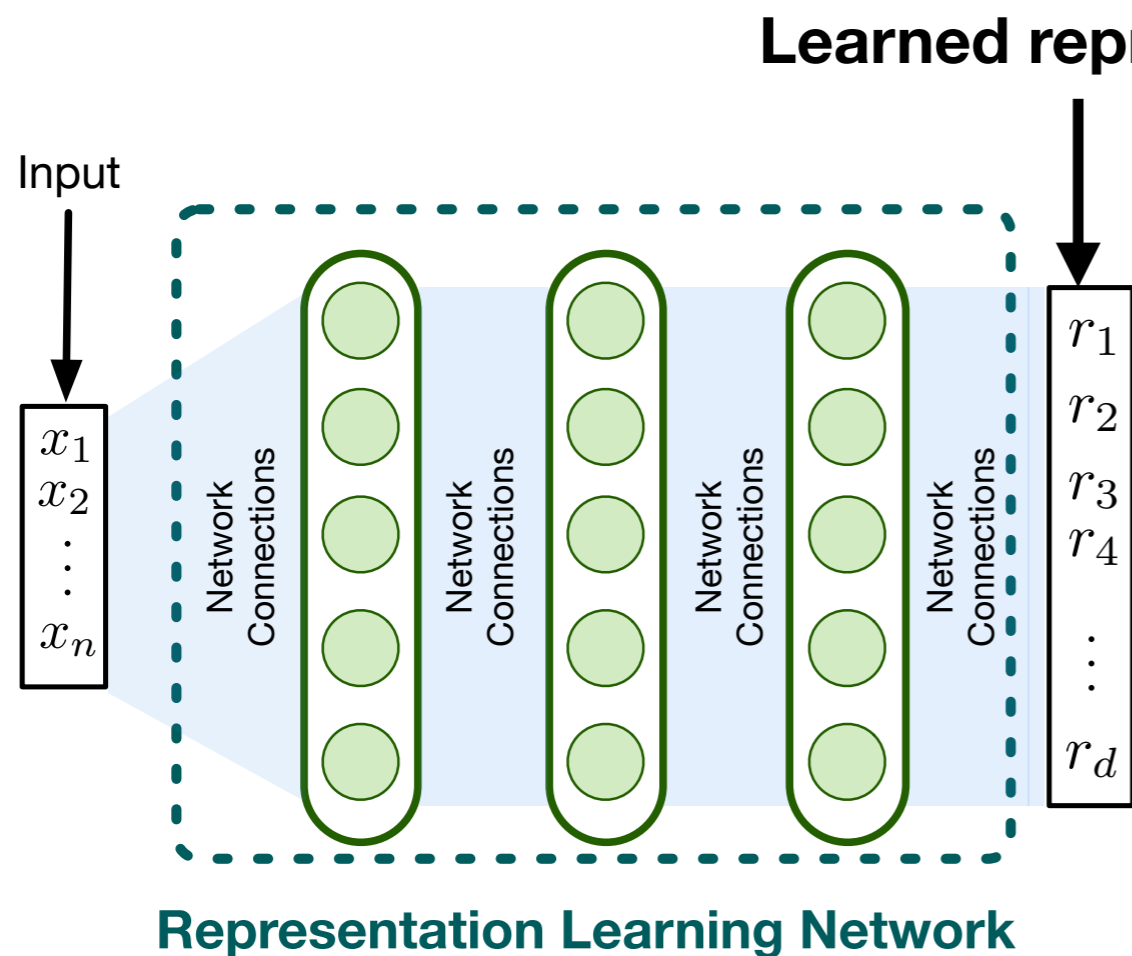


- **Pretraining:** 38%
- **SR-NN:** 15%
- **MRCL:** 3.8%

Analyzing representations

Instance sparsity

Percentage of non-zero entries used to represent an input on average



- **Pretraining:** 38%
- **SR-NN:** 15%
- **MRCL:** 3.8%
- **SR-NN-2:** 4.9%

Visualizing representations

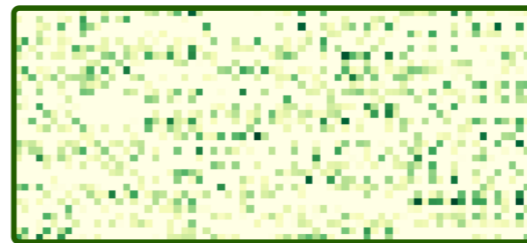
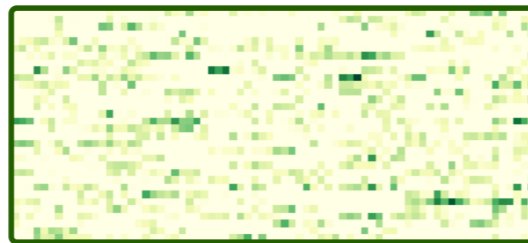
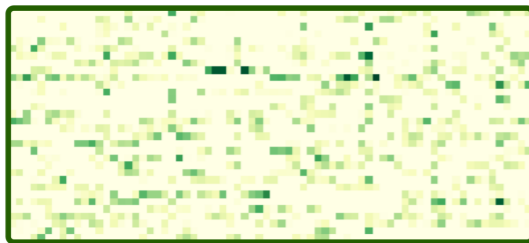
Visualizing representations

Random Instance 1

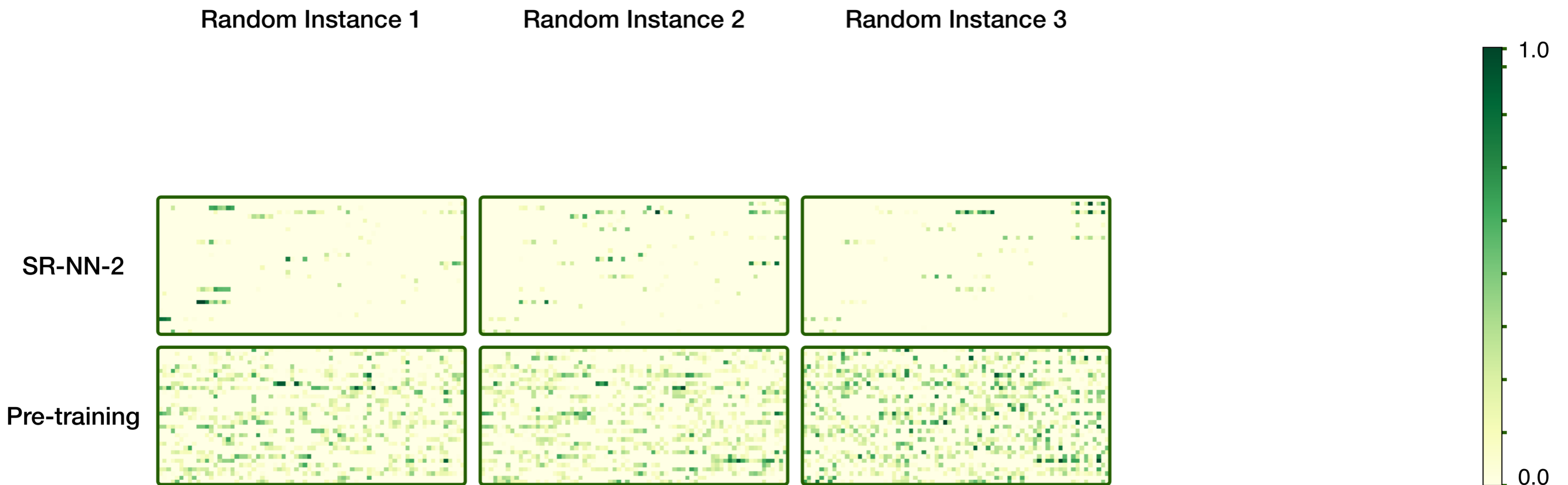
Random Instance 2

Random Instance 3

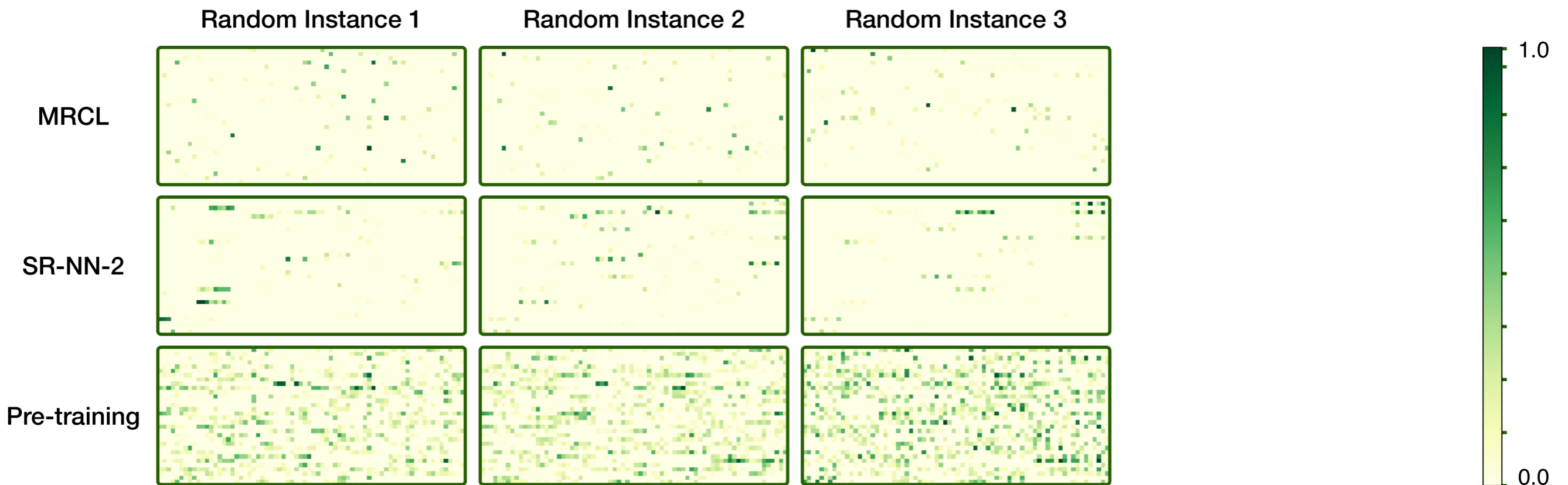
Pre-training



Visualizing representations

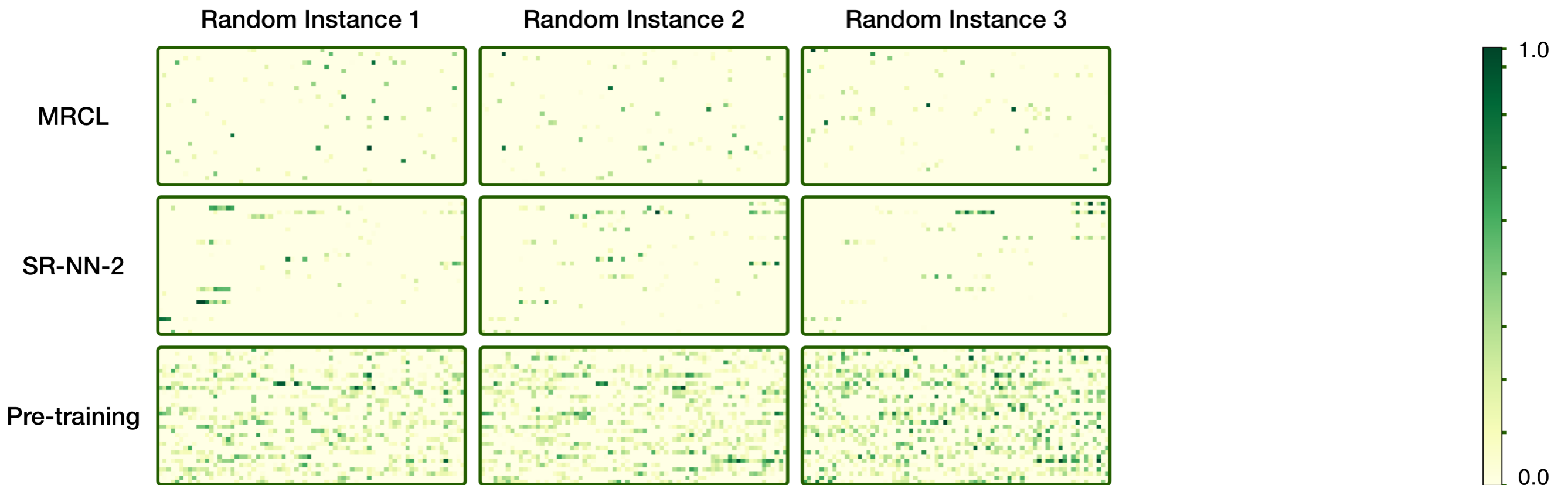


Visualizing representations



Visualizing representations

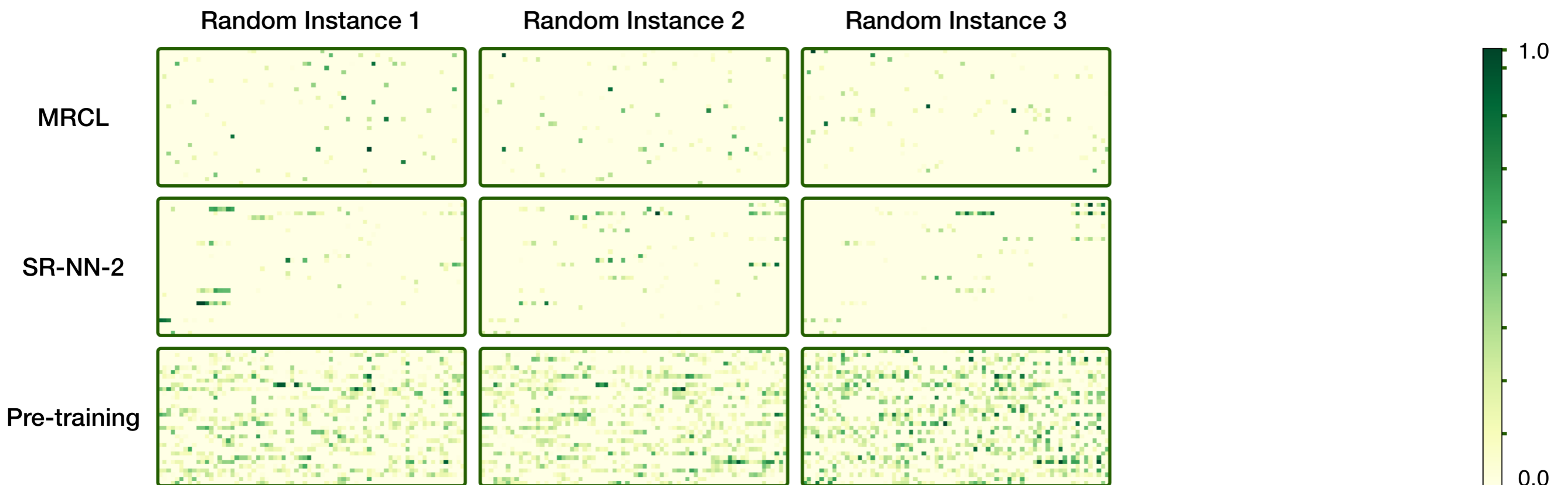
Average representation



Visualizing representations

Average representation

Compute representation of the complete representation learning dataset

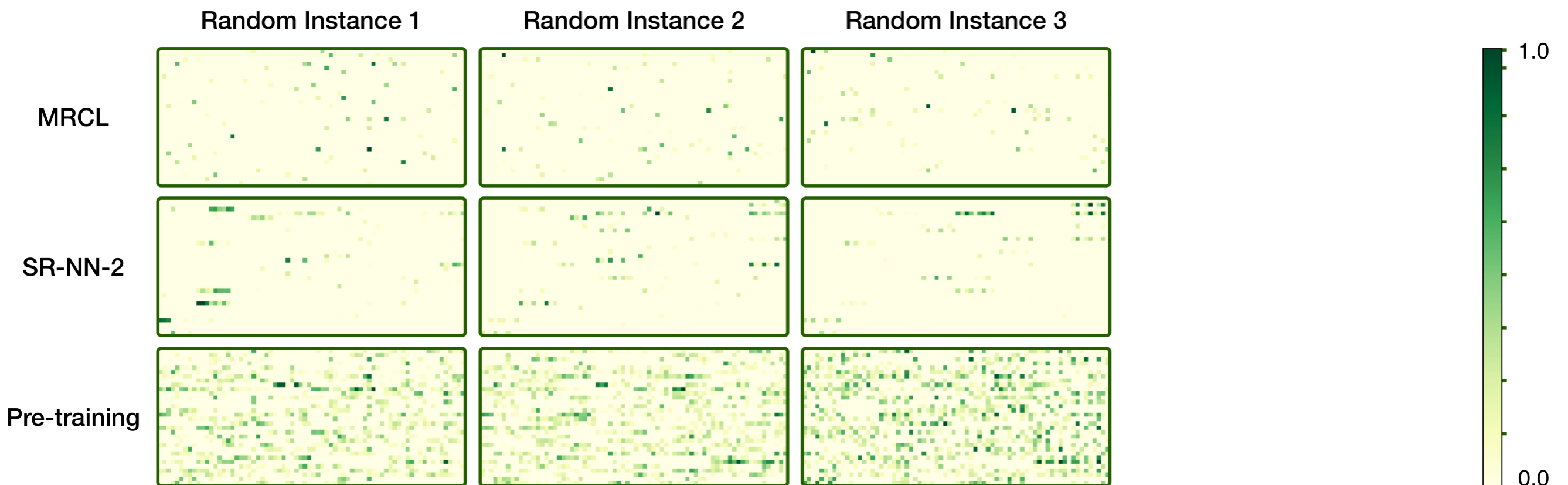


Visualizing representations

Average representation

Compute representation of the complete representation learning dataset

Compute the average representation vector

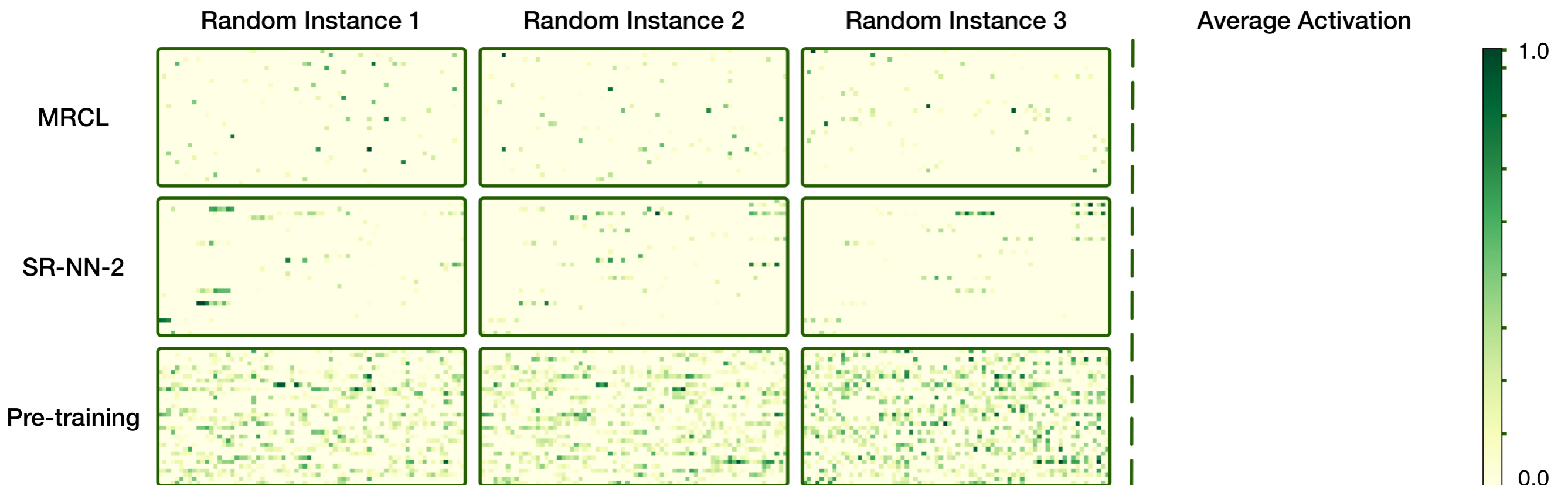


Visualizing representations

Average representation

Compute representation of the complete representation learning dataset

Compute the average representation vector

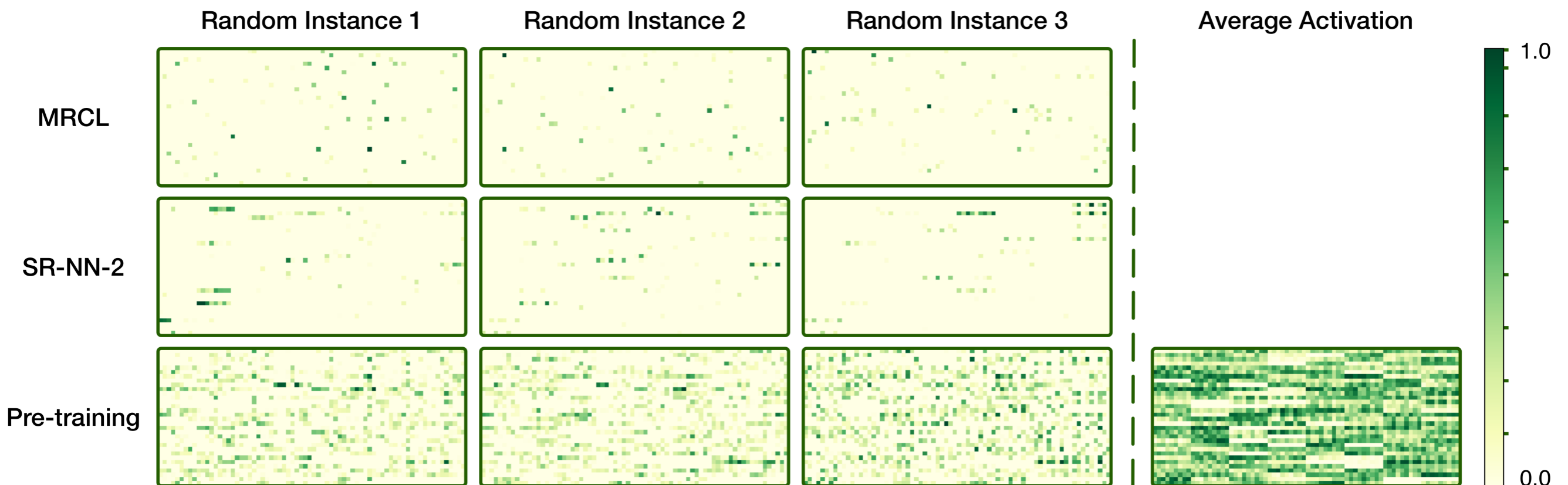


Visualizing representations

Average representation

Compute representation of the complete representation learning dataset

Compute the average representation vector

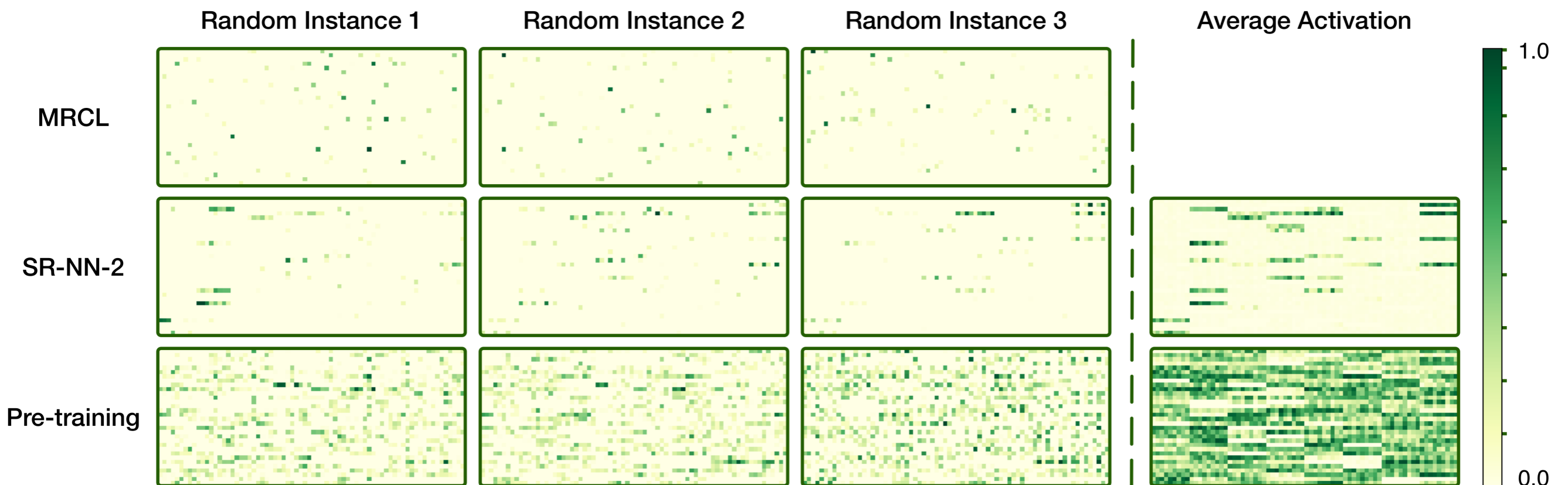


Visualizing representations

Average representation

Compute representation of the complete representation learning dataset

Compute the average representation vector

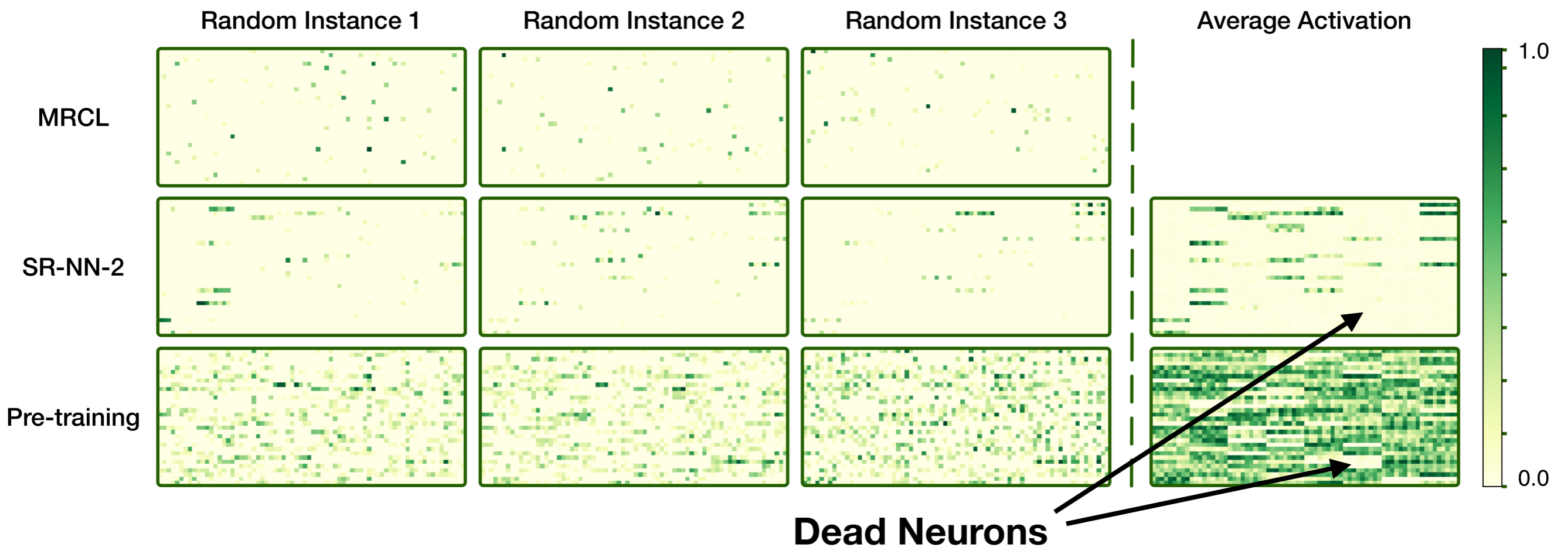


Visualizing representations

Average representation

Compute representation of the complete representation learning dataset

Compute the average representation vector

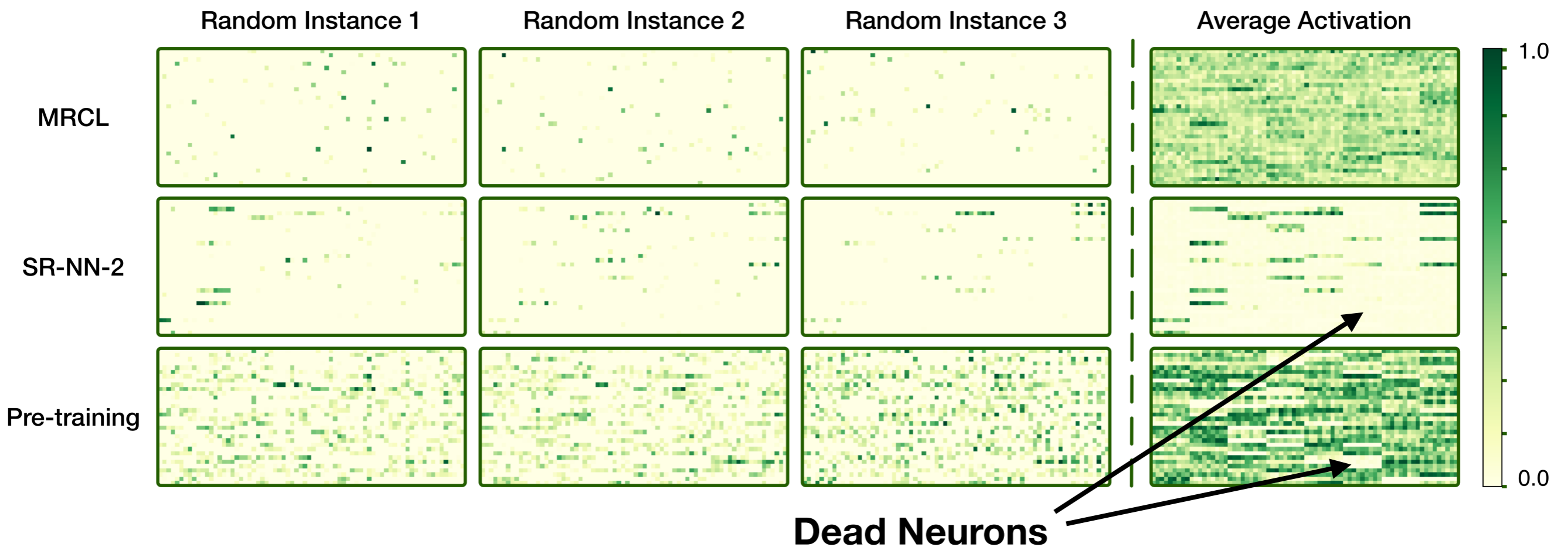


Visualizing representations

Average representation

Compute representation of the complete representation learning dataset

Compute the average representation vector



Visualizing representations

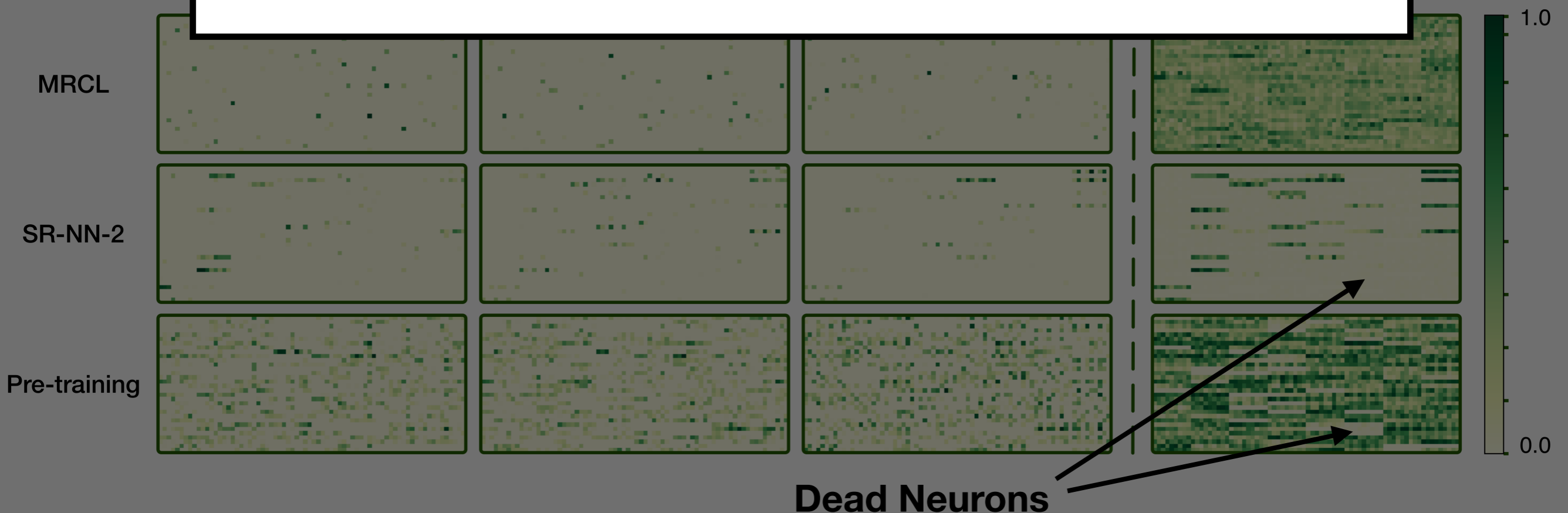
Average representation

Compute representation of the complete representation learning dataset

Compute the average representation vector

Possible explanation

MRCL does better because it is learning the right kind of sparsity!



Is catastrophic interference solved?

Is catastrophic interference solved?

- No! We assumed we have access to a dataset for learning representations

Is catastrophic interference solved?

- No! We assumed we have access to a dataset for learning representations
- Representation learning is not online

Is catastrophic interference solved?

- No! We assumed we have access to a dataset for learning representations
- Representation learning is not online
- However, positive preliminary results for online representation learning

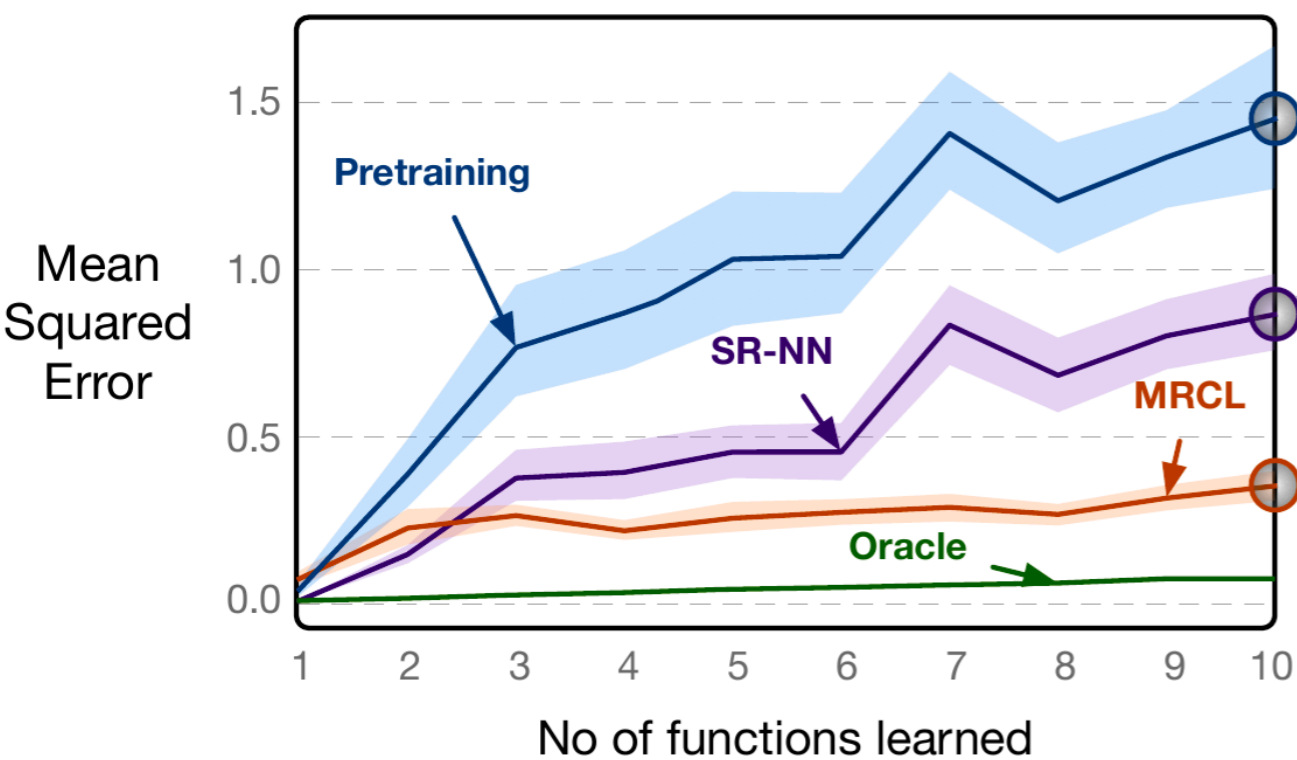
Questions?

Paper: <https://arxiv.org/pdf/1905.12588.pdf>

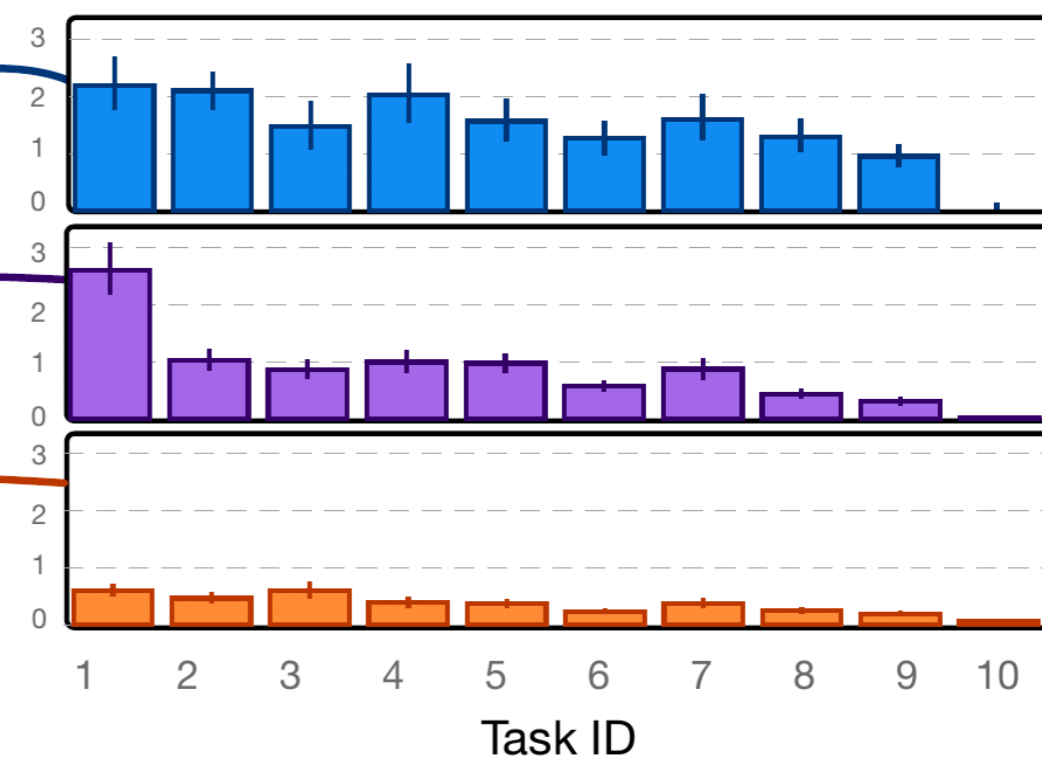
Code: <https://github.com/khurramjaved96/mrcl>

Continual Regression Tasks

Continual Regression Experiment



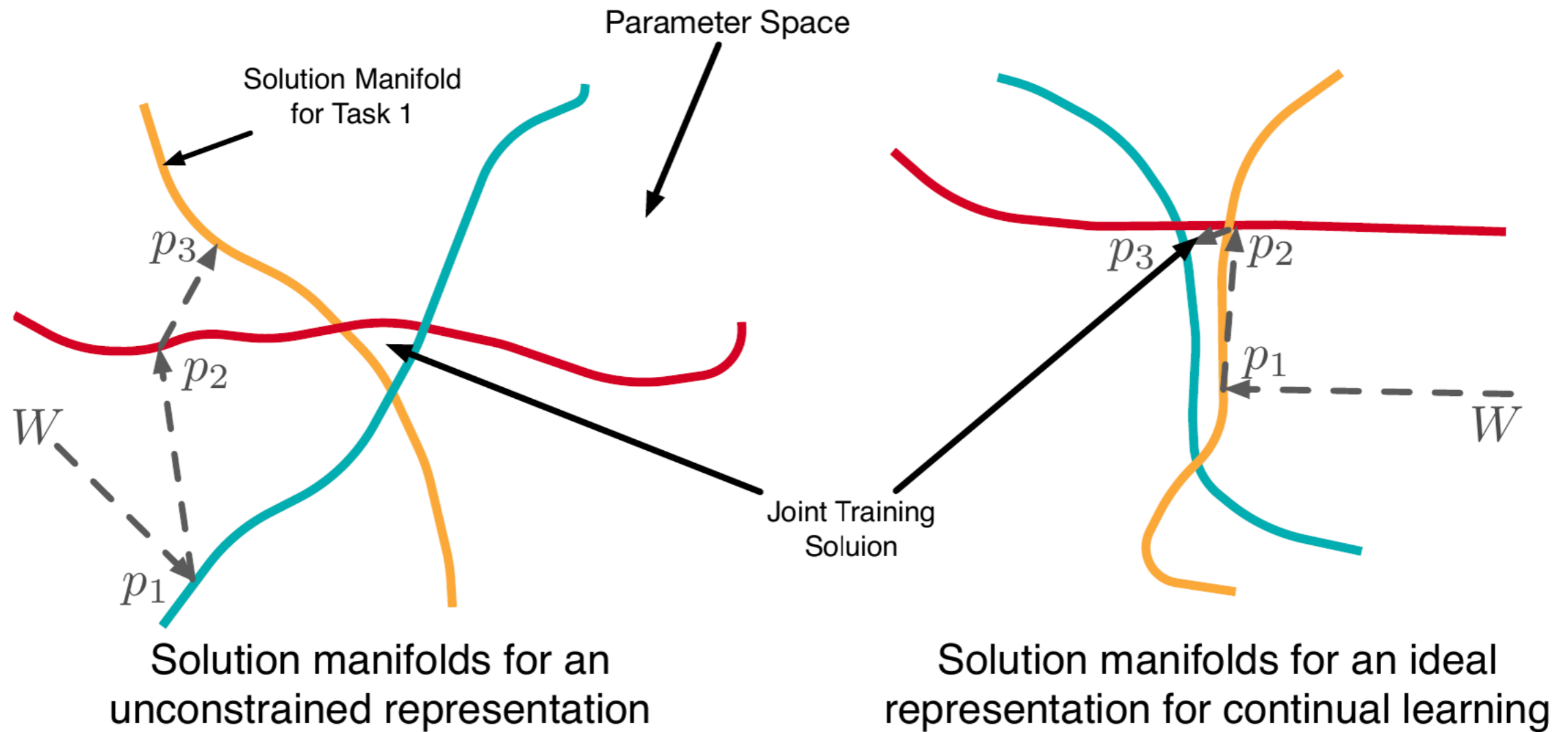
Error Distribution



Combining with existing methods

Method	Split-Omniglot					
	One class per task, 50 tasks			Five classes per task, 20 tasks		
	Standard	MRCL	Pretraining	Standard	MRCL	Pretraining
Online	04.64 \pm 2.61	64.72 \pm 2.57	21.16 \pm 2.71	01.40 \pm 0.43	55.32 \pm 2.25	11.80 \pm 1.92
Approx IID	53.95 \pm 5.50	75.12 \pm 3.24	54.29 \pm 3.48	48.02 \pm 5.67	67.03 \pm 2.10	46.02 \pm 2.83
ER-Reservoir	52.56 \pm 2.12	68.16 \pm 3.12	36.72 \pm 3.06	24.32 \pm 5.37	60.92 \pm 2.41	37.44 \pm 1.67
MER	54.88 \pm 4.12	76.00 \pm 2.07	62.76 \pm 2.16	29.02 \pm 4.01	62.05 \pm 2.19	42.05 \pm 3.71
EWC	05.08 \pm 2.47	64.44 \pm 3.13	18.72 \pm 3.97	02.04 \pm 0.35	56.03 \pm 3.20	10.03 \pm 1.53

Solution Manifolds



Pseudo-code

Algorithm 1: MRCL for optimizing objective in (3)

Require: $\mathcal{D}_{stream} = (X_1, Y_1), (X_2, Y_2), \dots, (X_t, Y_t), \dots;$

Require: $g_W(\phi_\theta(x))$ as a parametrized function;

Require: $\alpha, \beta, \mathcal{L}, n$ as meta learning rate, inner learning rate, loss metric over (X_i, Y_i) and total gradient updates;

Initialize RLN and TLN to θ and W ;

for *iterations* $1, 2, 3, \dots, n$ **do**

 Sample trajectory $(\mathbf{X}_{traj}, \mathbf{Y}_{traj}) = (X_{i+1}, Y_{i+1}) \dots (X_{i+k}, Y_{i+k}) \sim \mathcal{D}_{stream};$

$W_0 = W;$

for j in $1, 2, 3, \dots, k$ **do**

$W_j = W_{j-1} - \beta \nabla_{W_{j-1}} (\mathcal{L}(g_{W_{j-1}}(\phi_\theta(X_{i+j})), Y_{i+j});$

end

$(\mathbf{X}_{rand}, \mathbf{Y}_{rand}) \sim \mathcal{D}_{stream};$

Sample a random batch of data

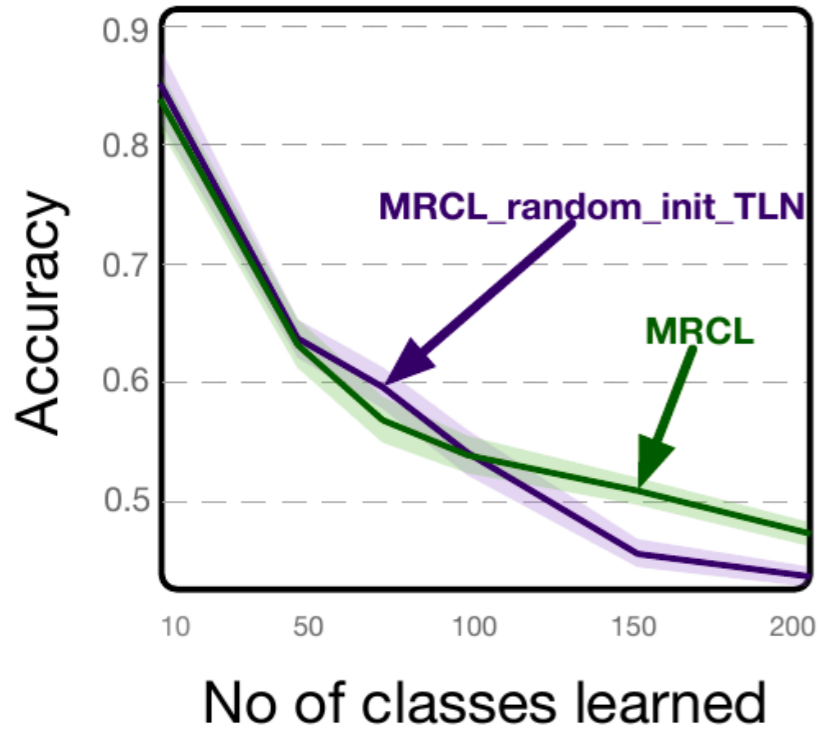
$(\mathbf{X}_{meta}, \mathbf{Y}_{meta}) = (\mathbf{X}_{rand} + \mathbf{X}_{traj}, \mathbf{Y}_{rand} + \mathbf{Y}_{traj});$

$W, \theta = (W, \theta) - \alpha \nabla_{W, \theta} \mathcal{L}(g_{W_k}(\phi_\theta(\mathbf{X}_{meta})), \mathbf{Y}_{meta});$

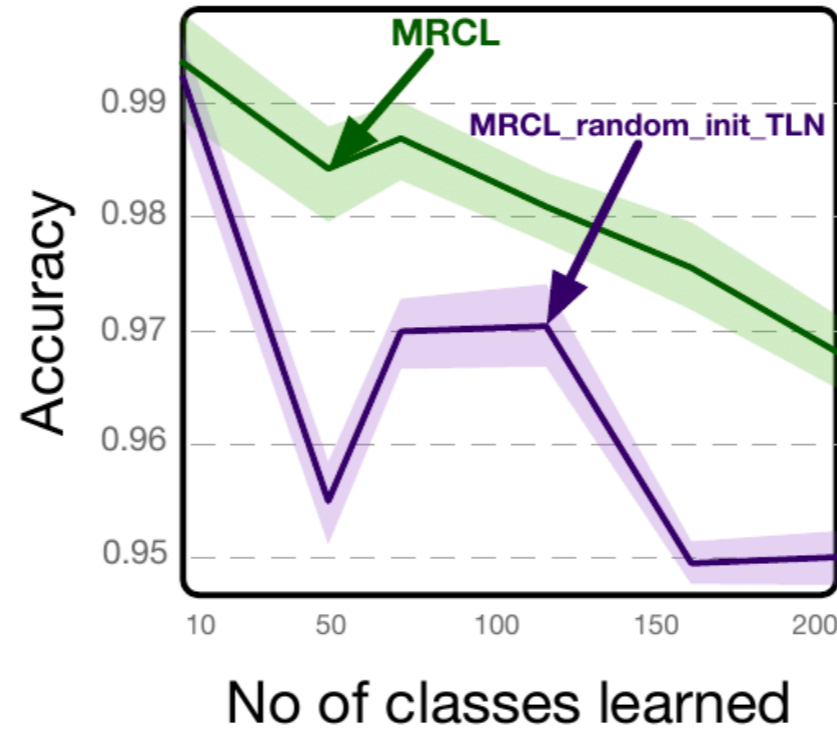
end

PLN Initialization

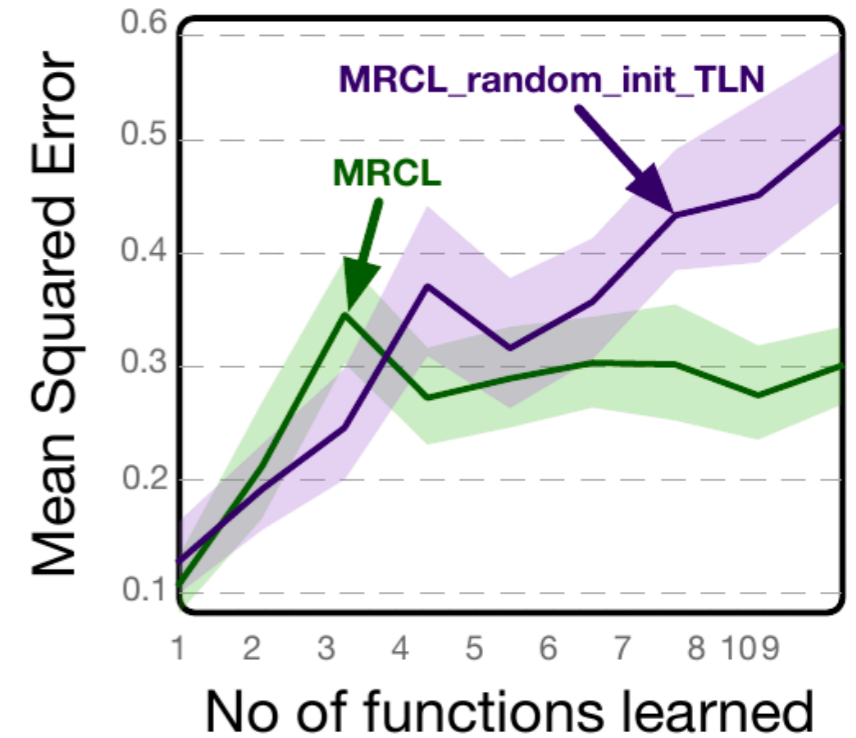
Omniglot Test



Omniglot Training Trajectory

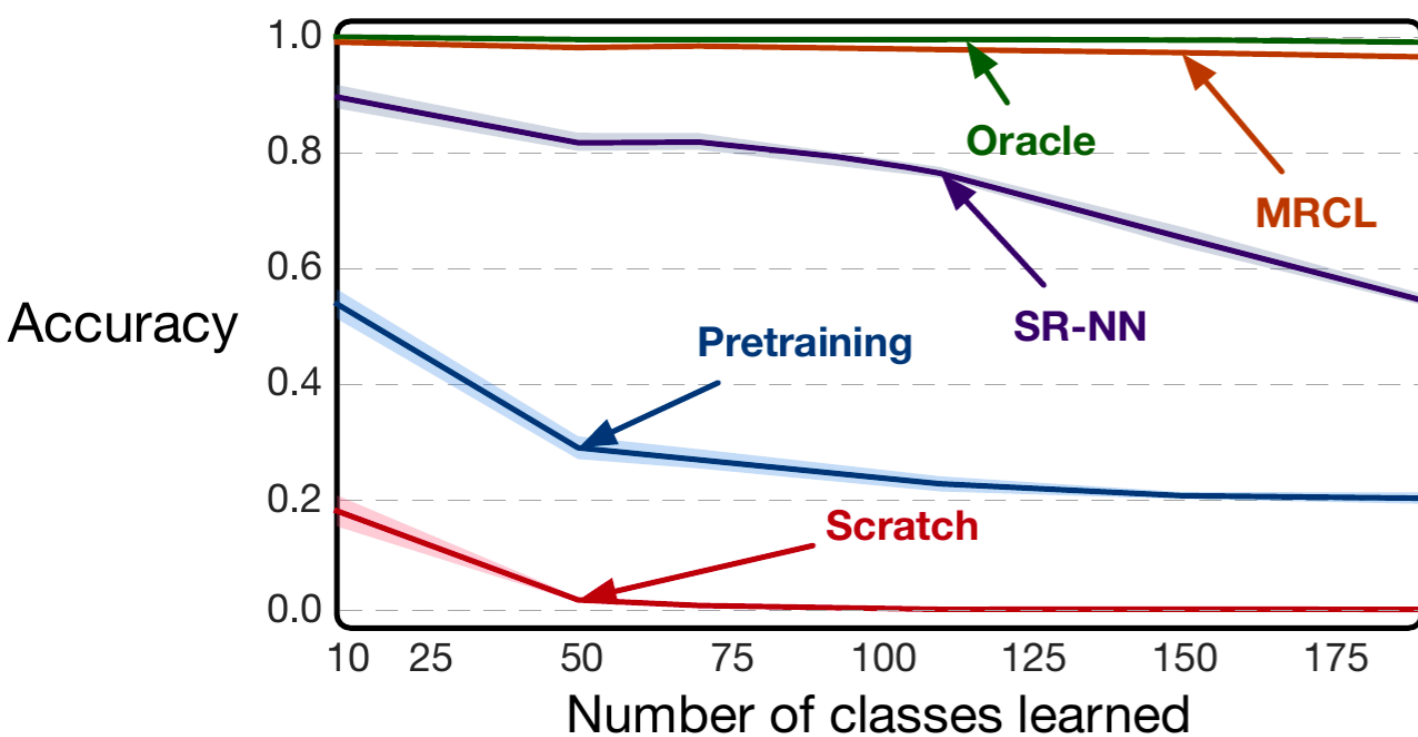


Regression



Generalization Error

Omniglot Training Trajectory Performance



Omniglot Test Set Performance

