

# Hardware-aware Few-shot Learning on a Memristor-based Small-world Architecture

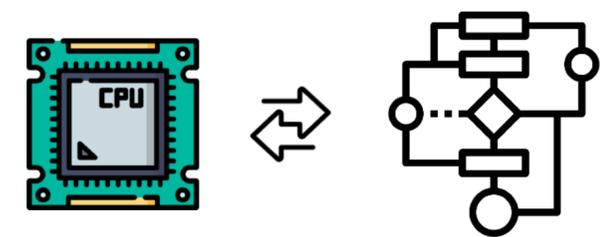
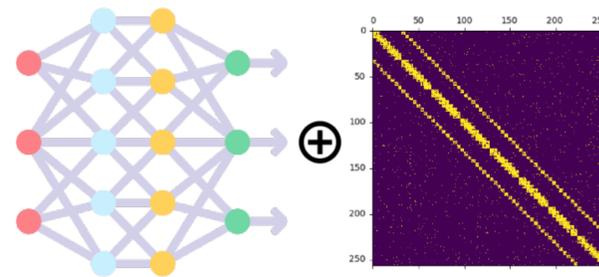
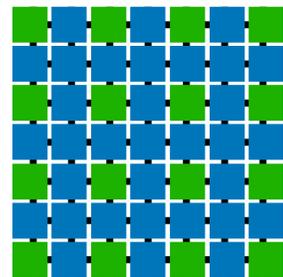
**Institute of Neuroinformatics**  
*University of Zürich and ETH Zürich*

Karthik Charan Raghunathan  
Emerging Intelligent Substrates lab  
24th April 2024

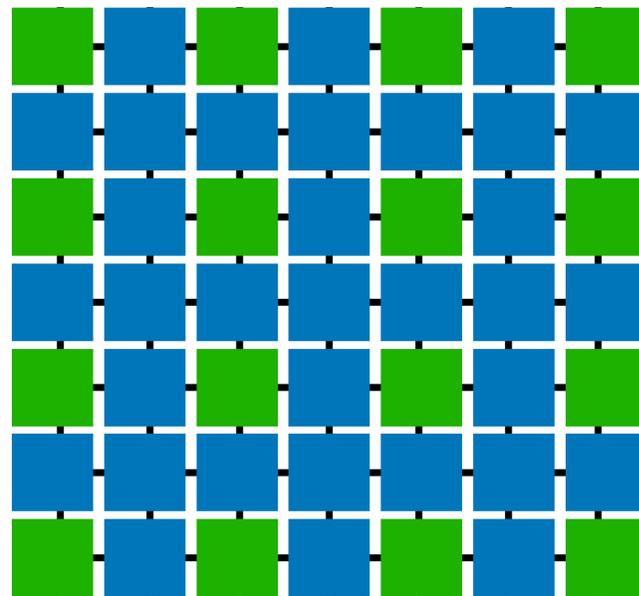


# Take home message

- Demonstrated (using simulations) Online few-shot learning on a novel, energy efficient hardware architecture
- Constraining network architecture does not terminally detriment the performance of the model
- Importance of inner loop dynamics in MAML; Highlighting the need for HW-SW co-design for efficient

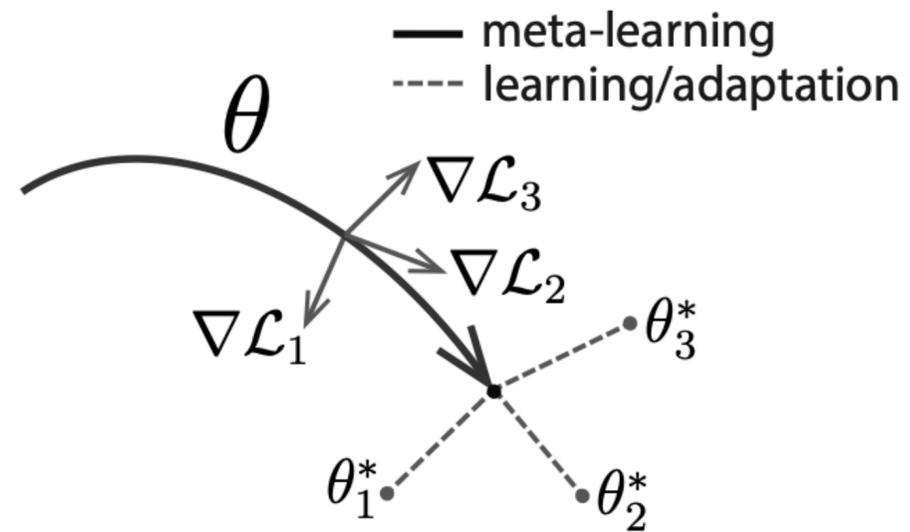


# What is this project about?



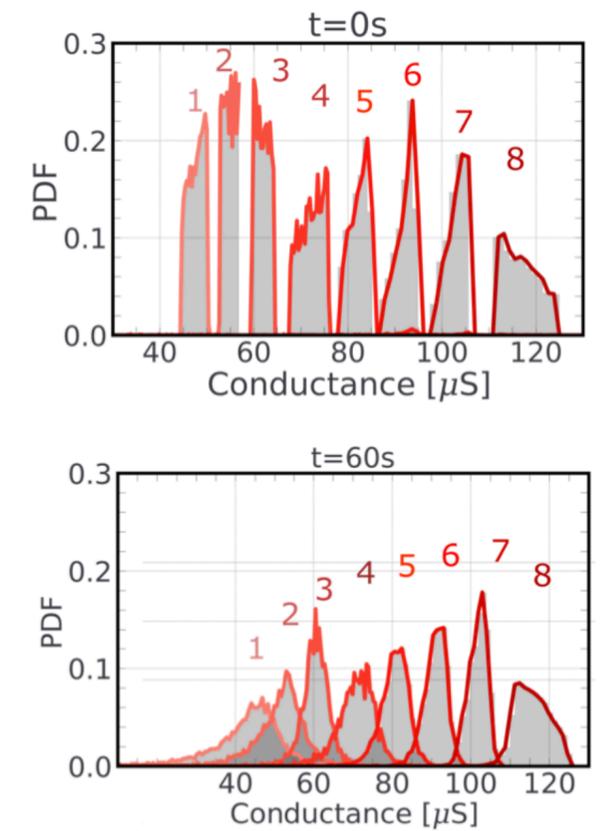
Novel energy efficient architecture

+



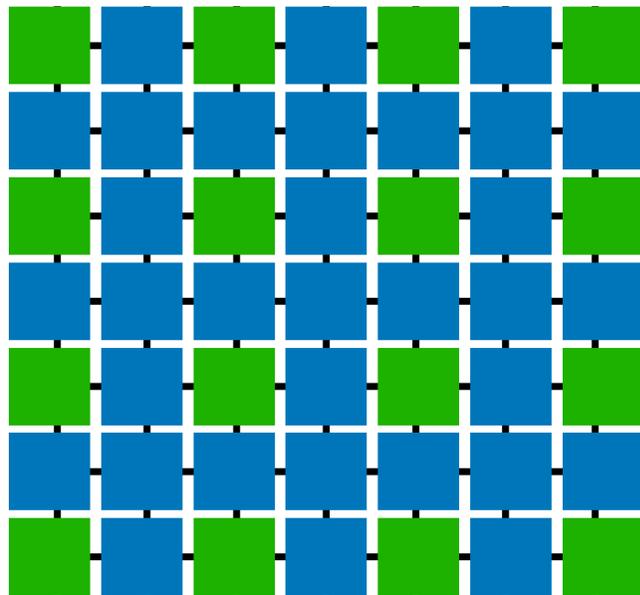
Meta-learning

+

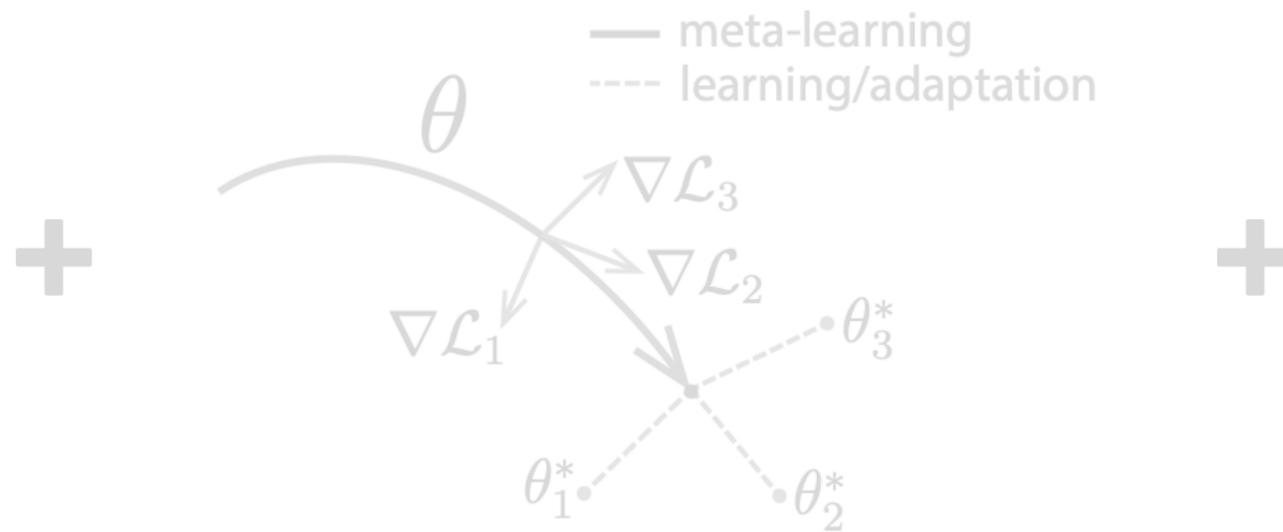


Hardware-aware training

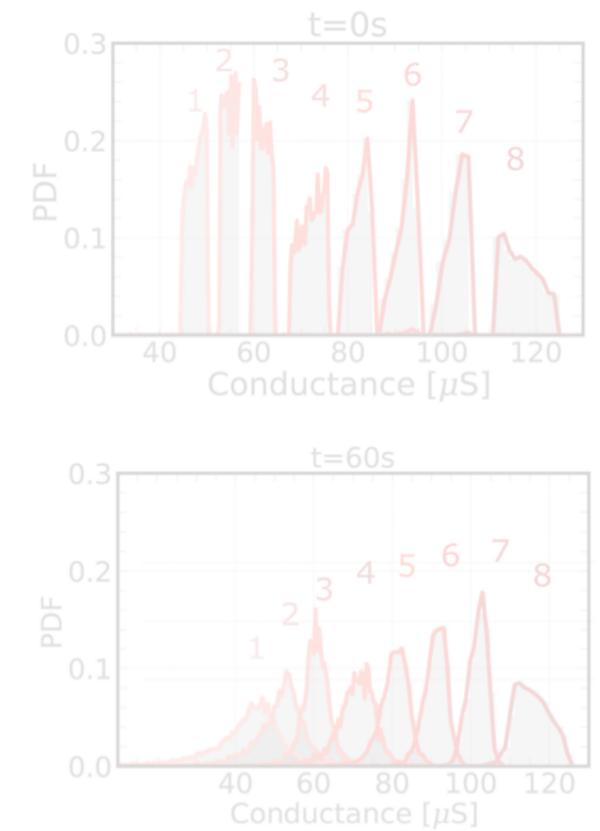
# Low-power architecture



Novel energy efficient architecture

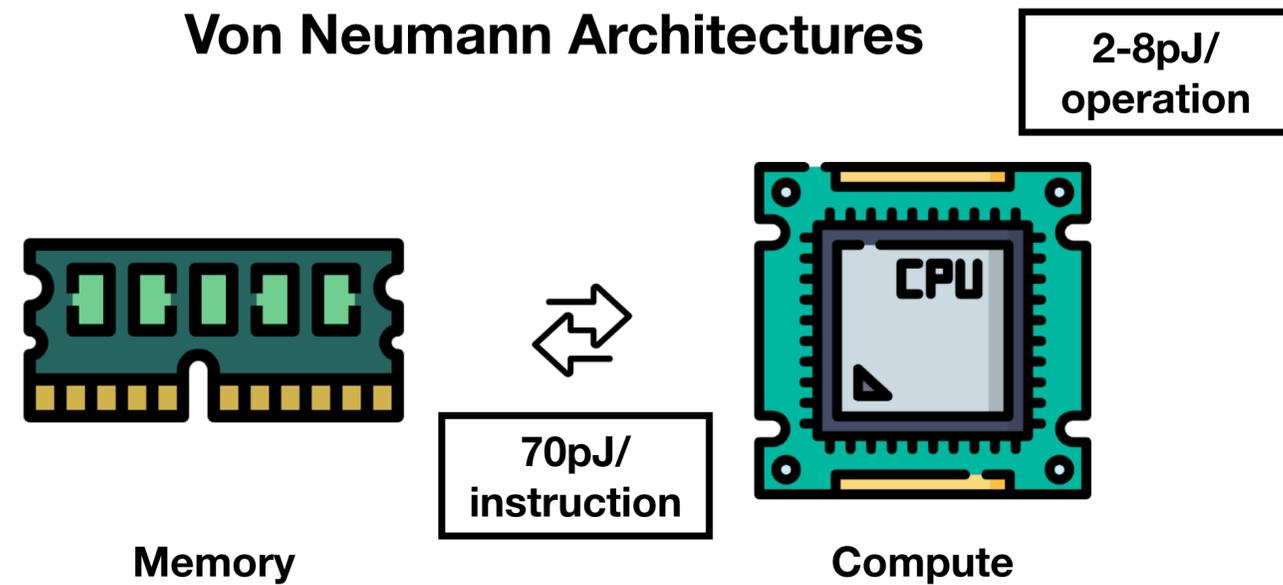


Meta-learning



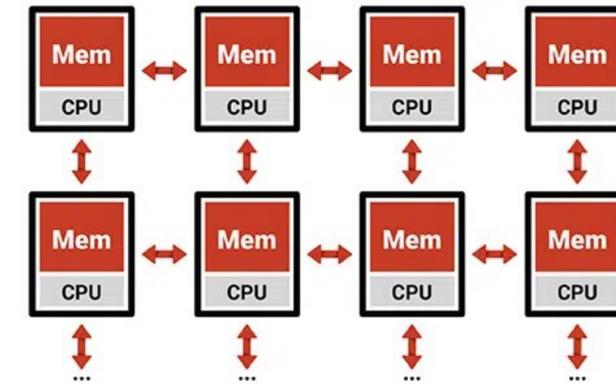
Hardware-aware training

# The Memory wall problem



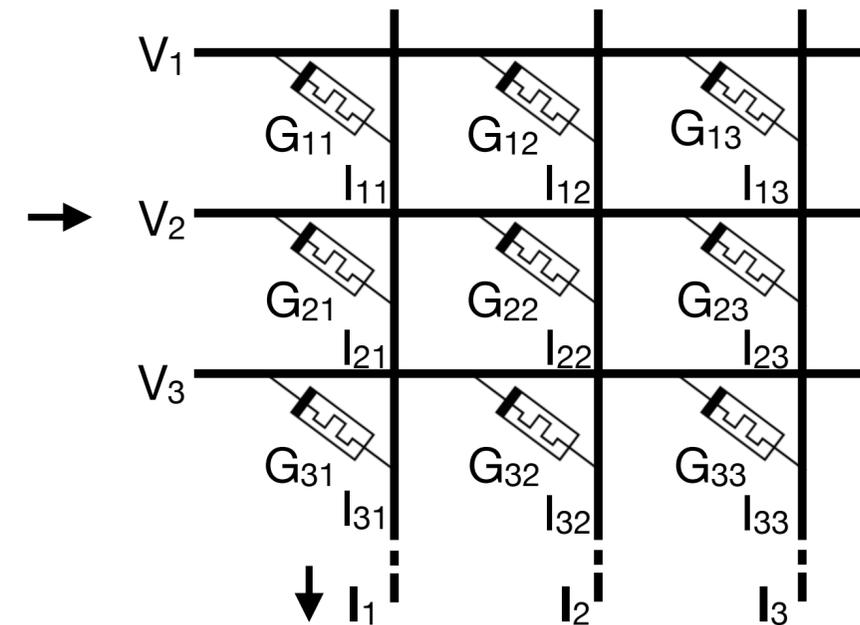
## Non-Von Neumann Architectures

### Near-memory compute



Memory & Compute

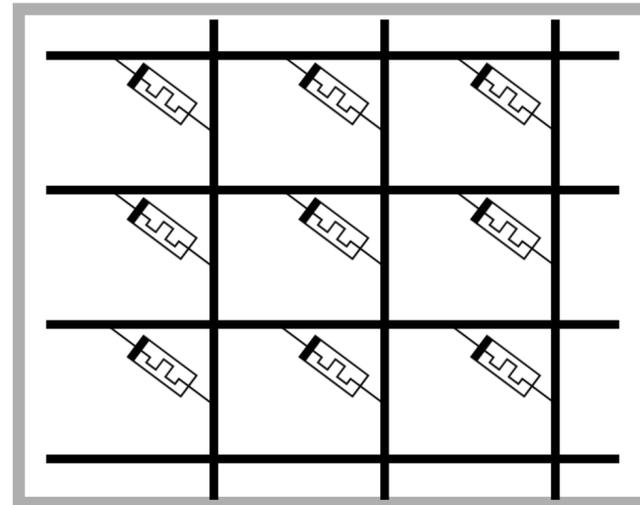
### In-memory compute



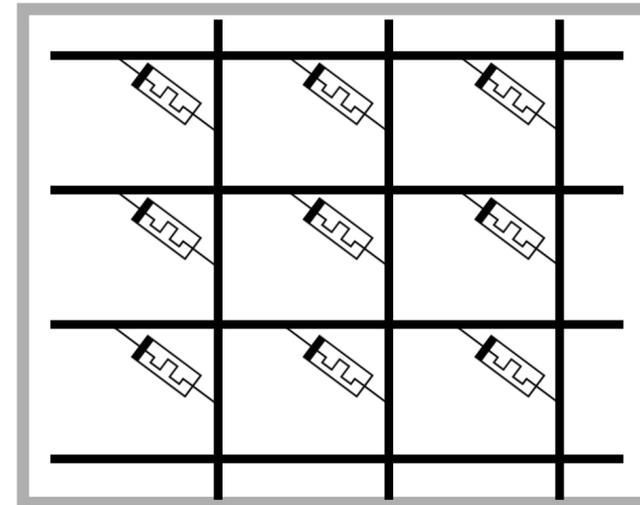
Memory + Compute

# Scaling IMC: *Possible solution*

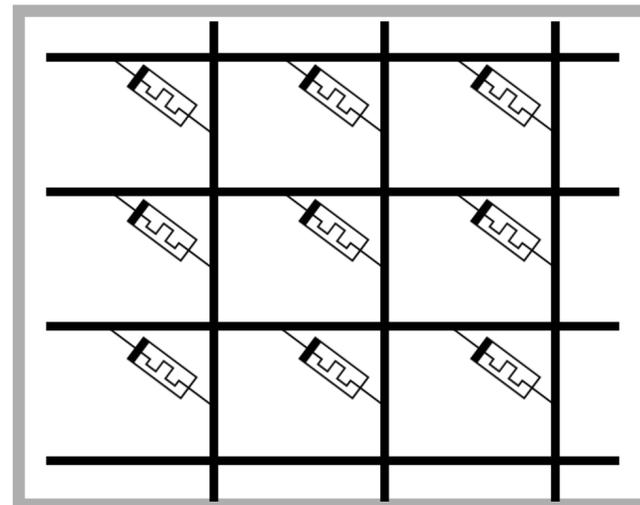
Core 1



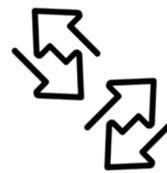
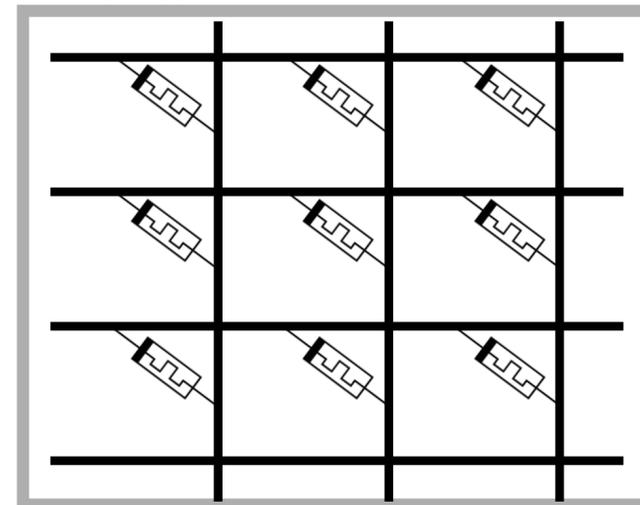
Core 2



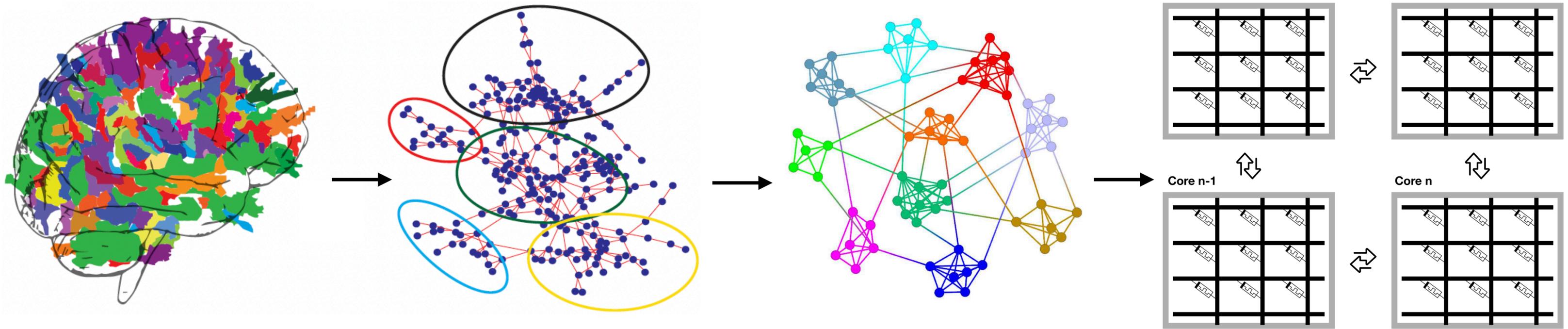
Core n-1



Core n

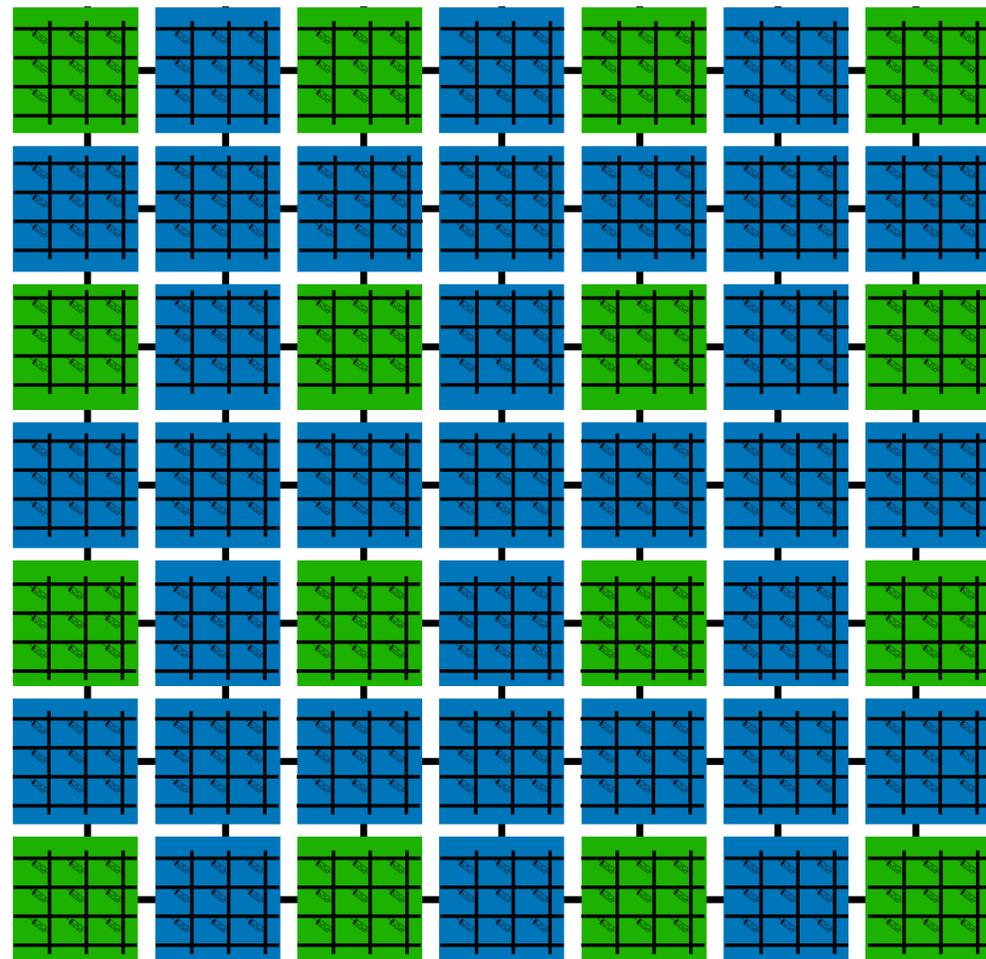


# Scaling multi-core IMC

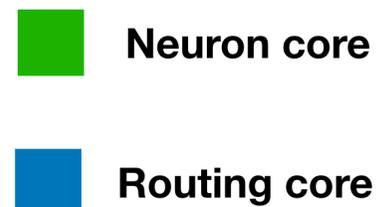


- High clustering coefficient
- Low number of total connections

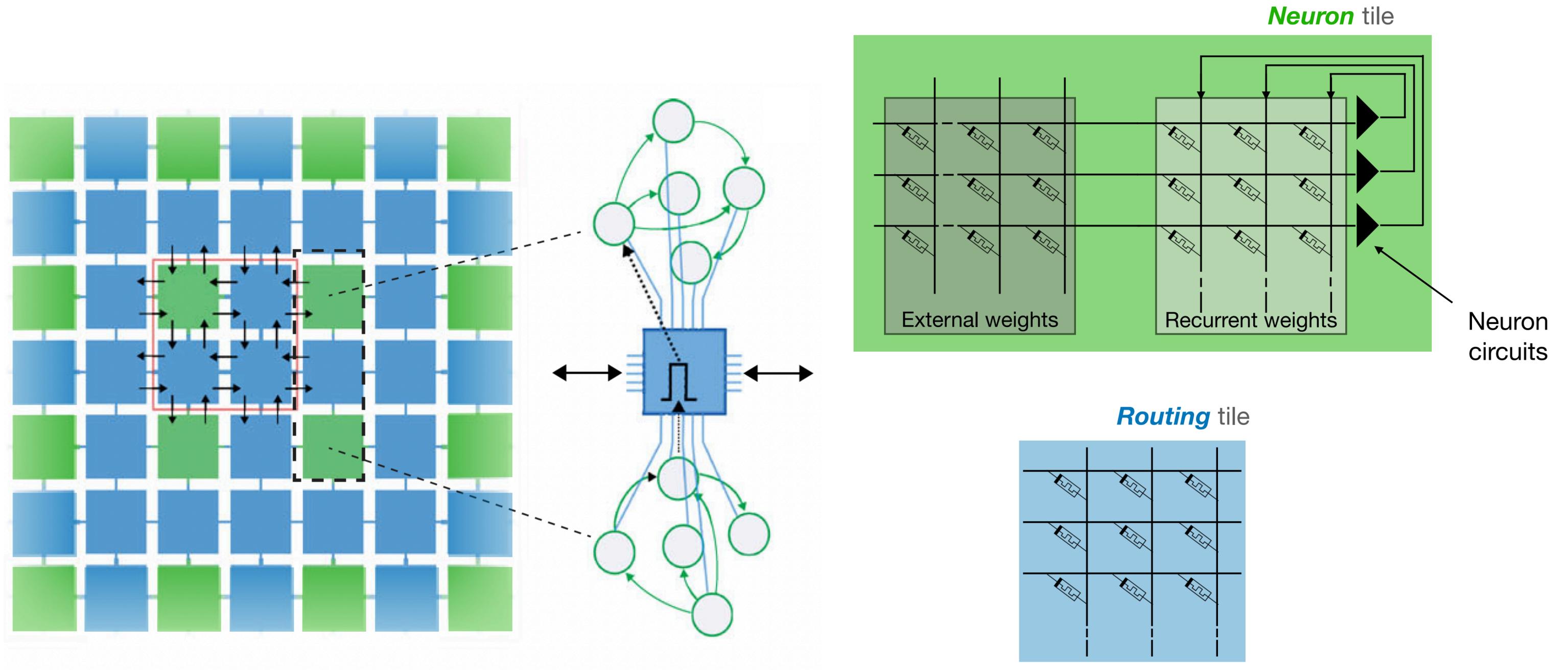
# The Neuromorphic Mosaic



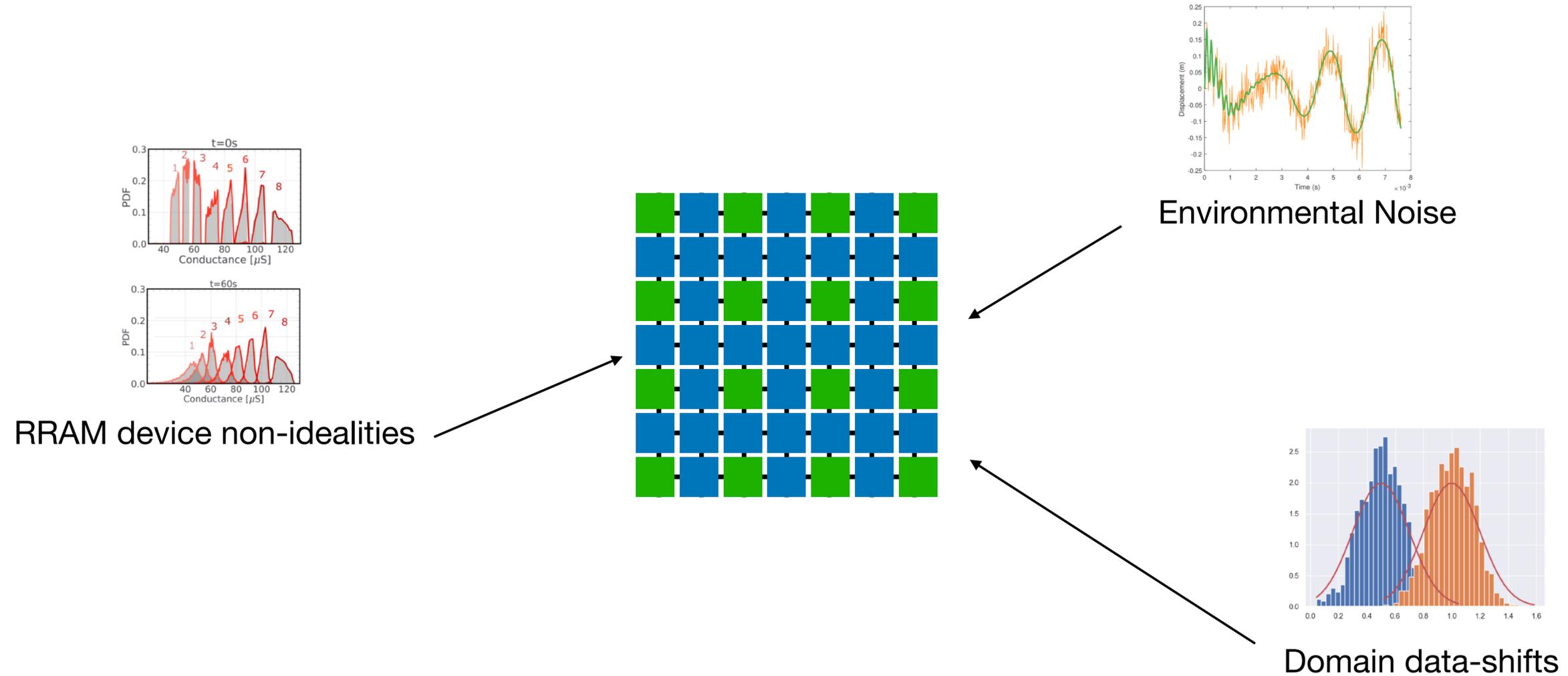
- Non-von Neumann ***systolic*** architecture
- Distributed memristors for ***in-memory computing*** and ***in-memory routing***
- Efficiently implements ***small-world*** graph topologies for Spiking Neural Networks (SNNs)



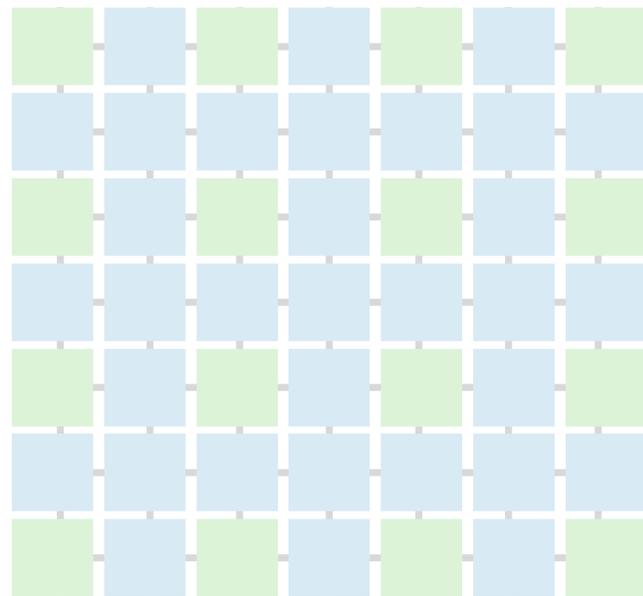
# The Neuromorphic Mosaic



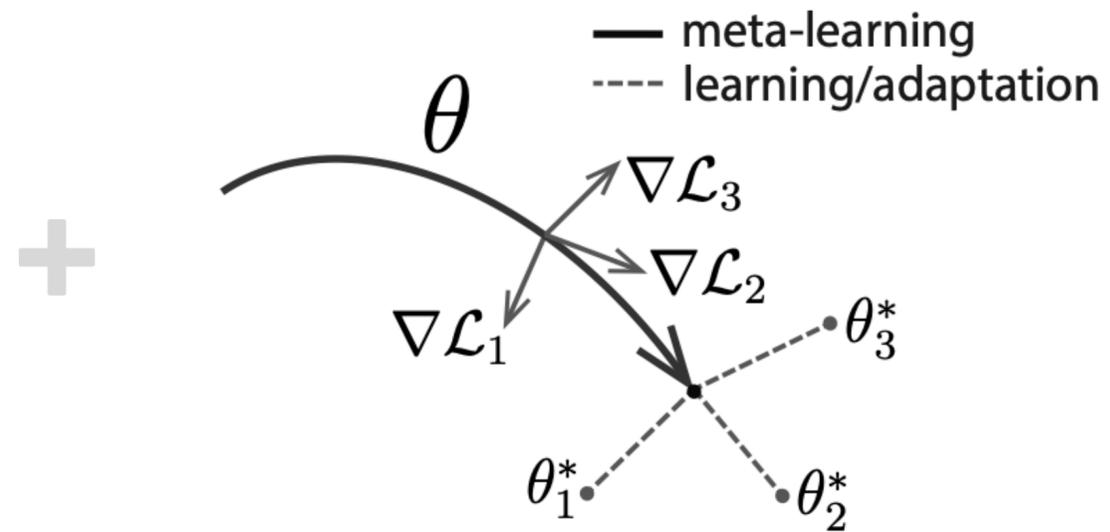
# Neuromorphic Mosaic for edge computing



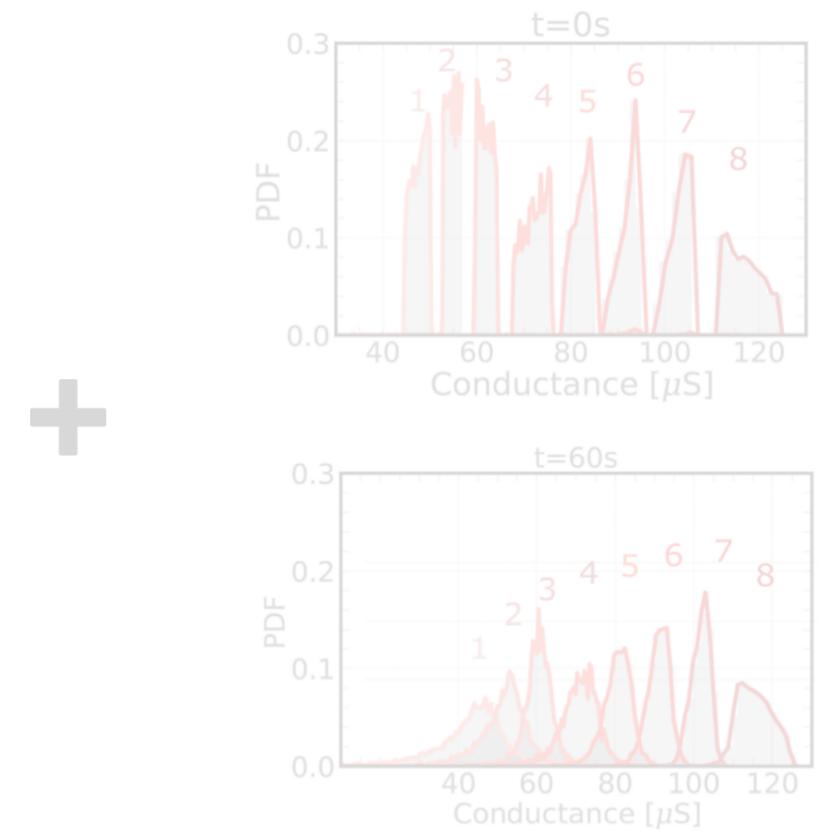
# Meta-learning algorithm



Novel energy efficient architecture



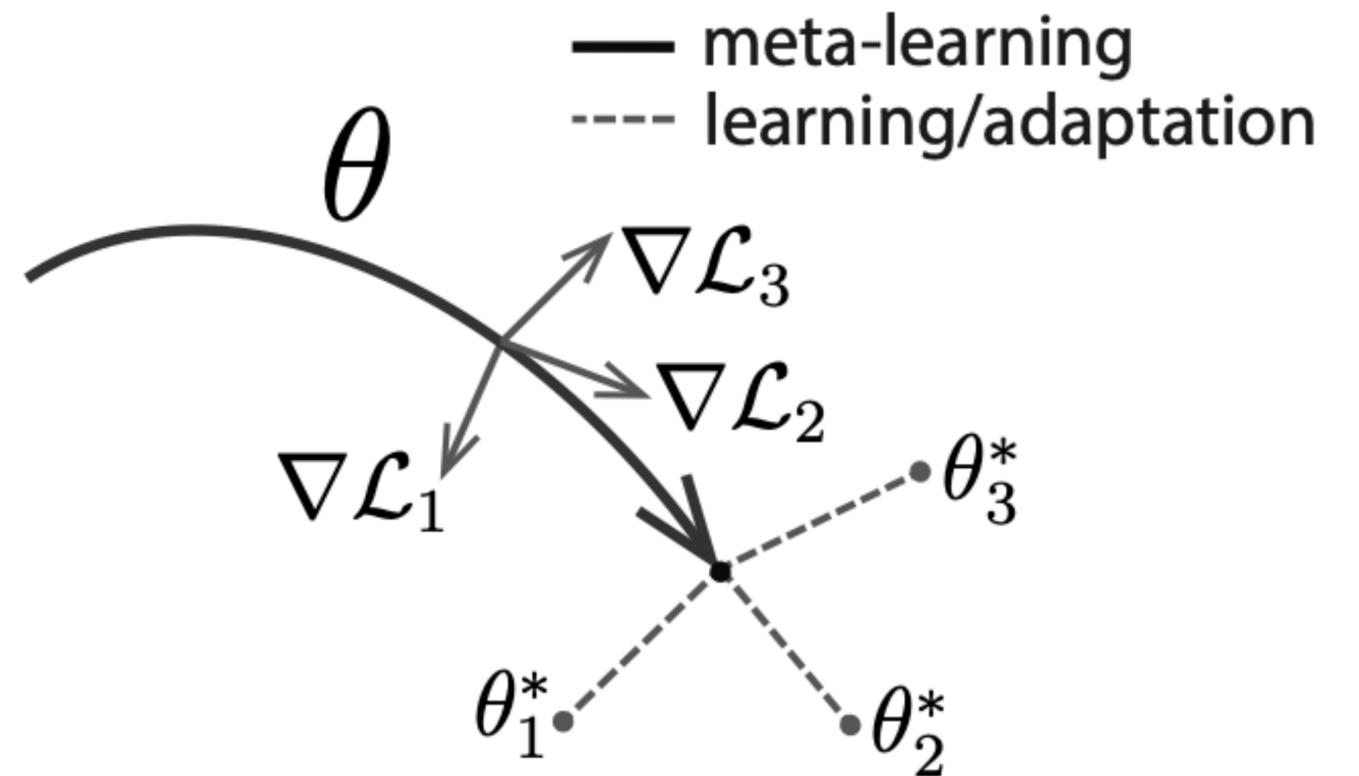
Meta-learning



Hardware-aware training

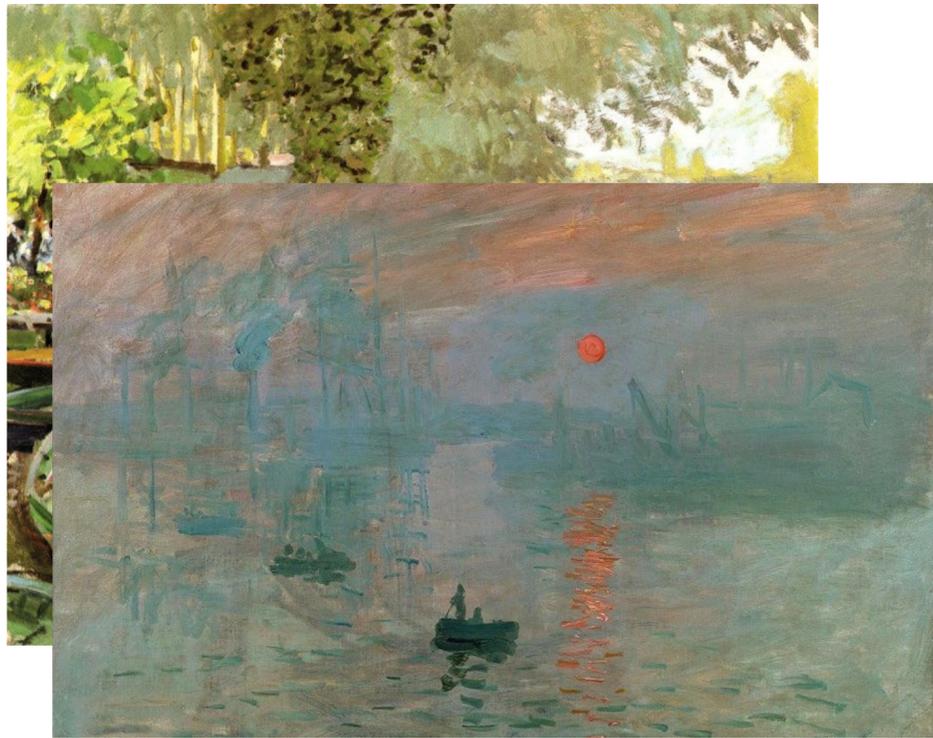
# Model Agnostic Meta-learning (MAML)

- Learning to learn approach
- Model agnostic
- Learn a good initialisation
- Fine-tune to learn new examples (few-shot)
- Gradient-based bi-level optimisation



# A few-shot task

← N-way →



↙  
**K-shot**

[“Claude Monet”, “Jackson Pollock”, “Vincent Van Gogh”]

# A few-shot data loader

*Meta-training*

$D^{train}$



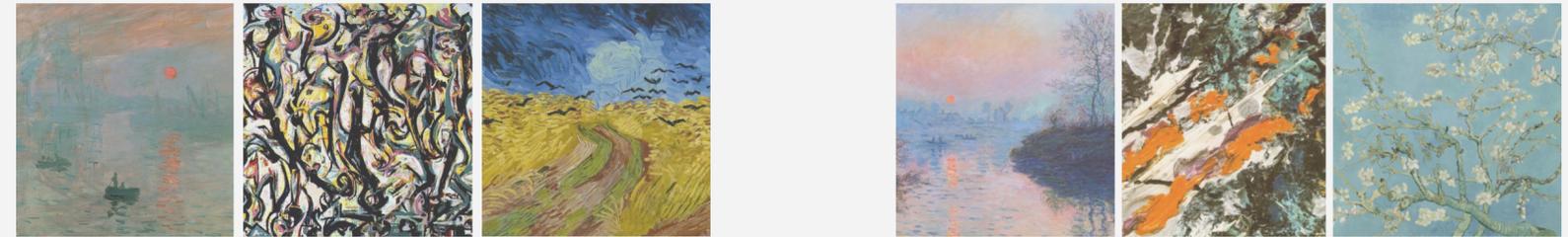
*Support set*

*Query set*

*Fine tune to this particular task*

*Meta-testing*

$D^{test}$



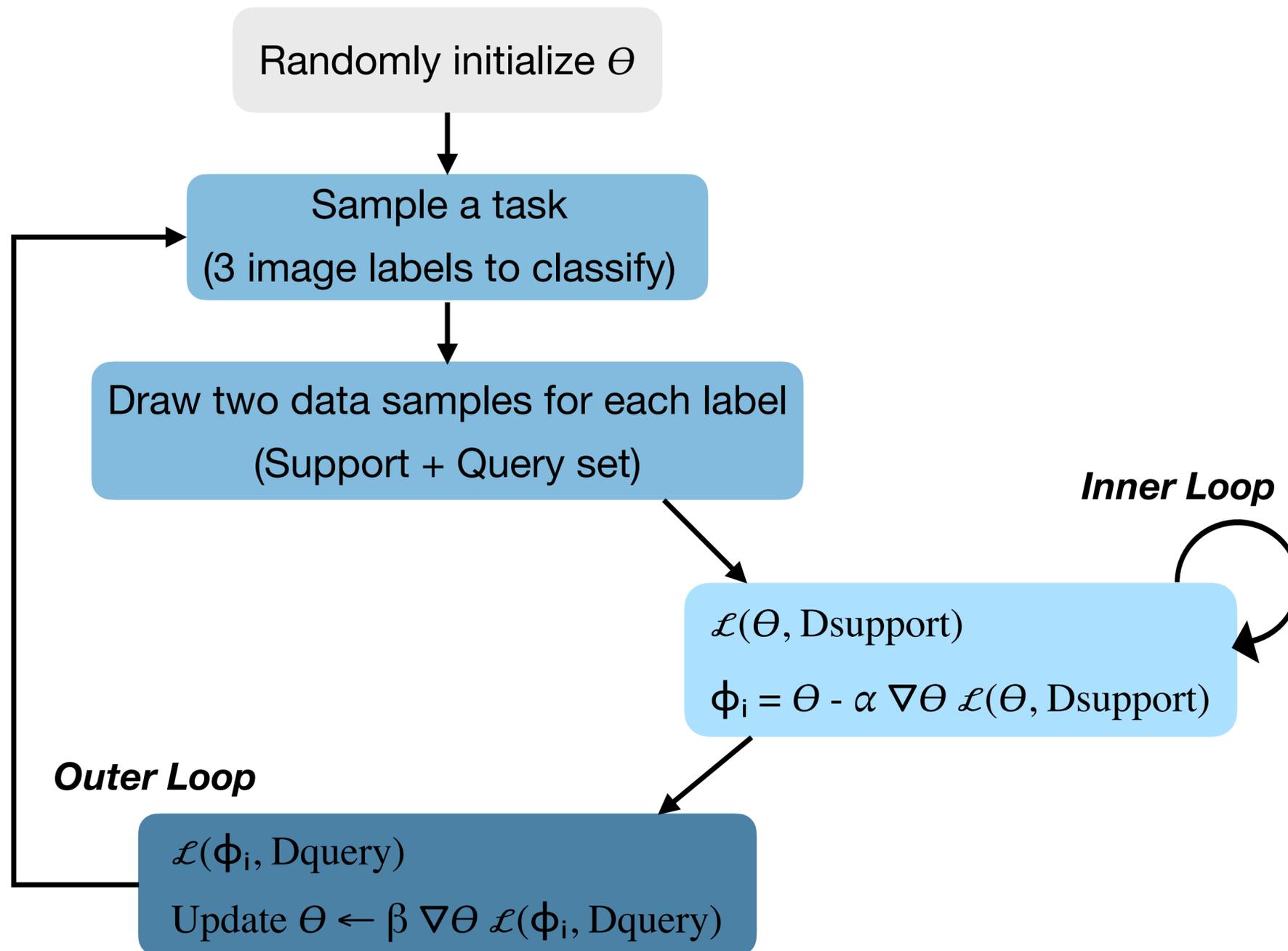
*Support set*

*Query set*

*Evaluate the FSL performance*

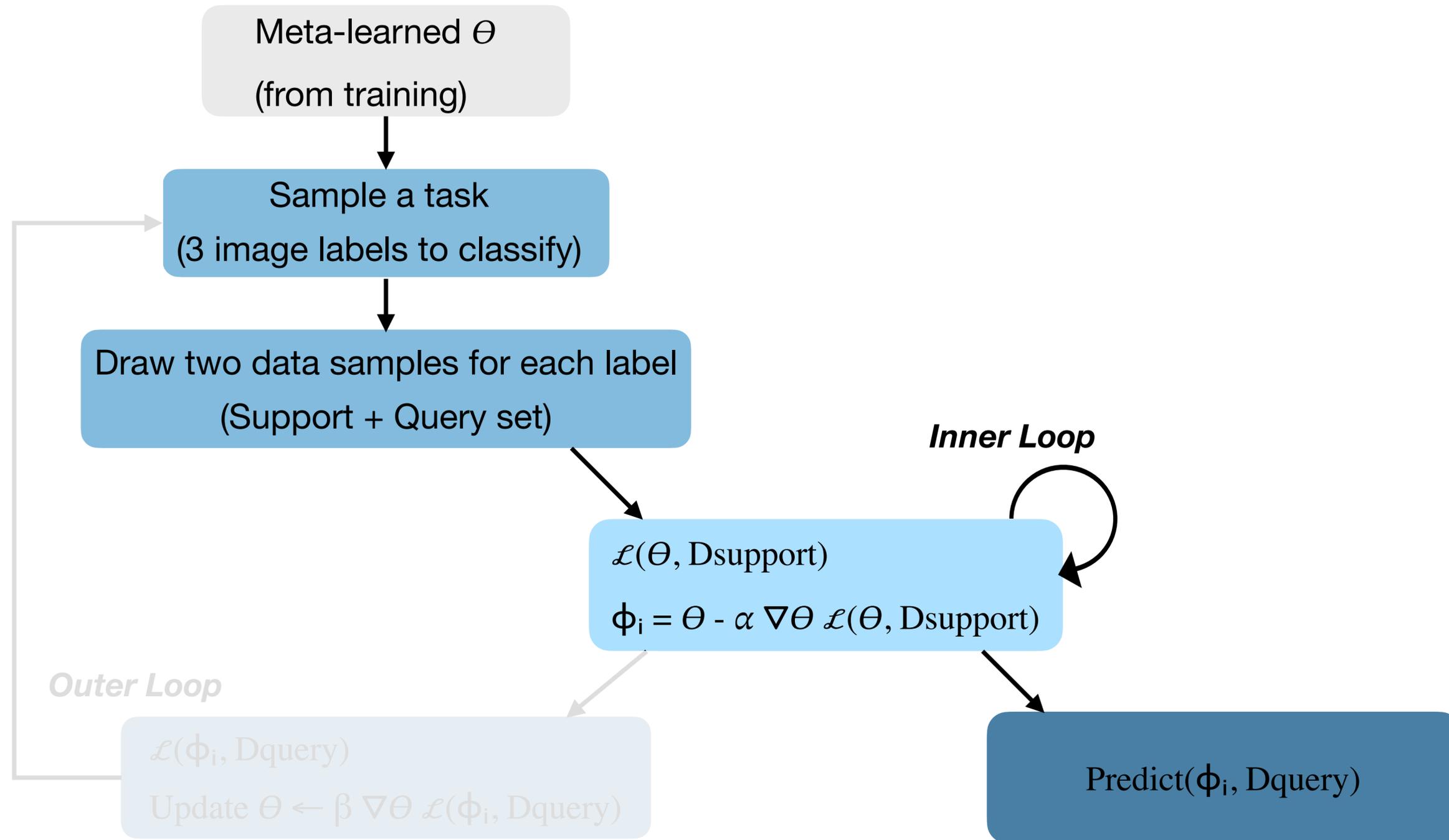
# Lifecycle of MAML

## *Meta-training phase*



# Lifecycle of MAML

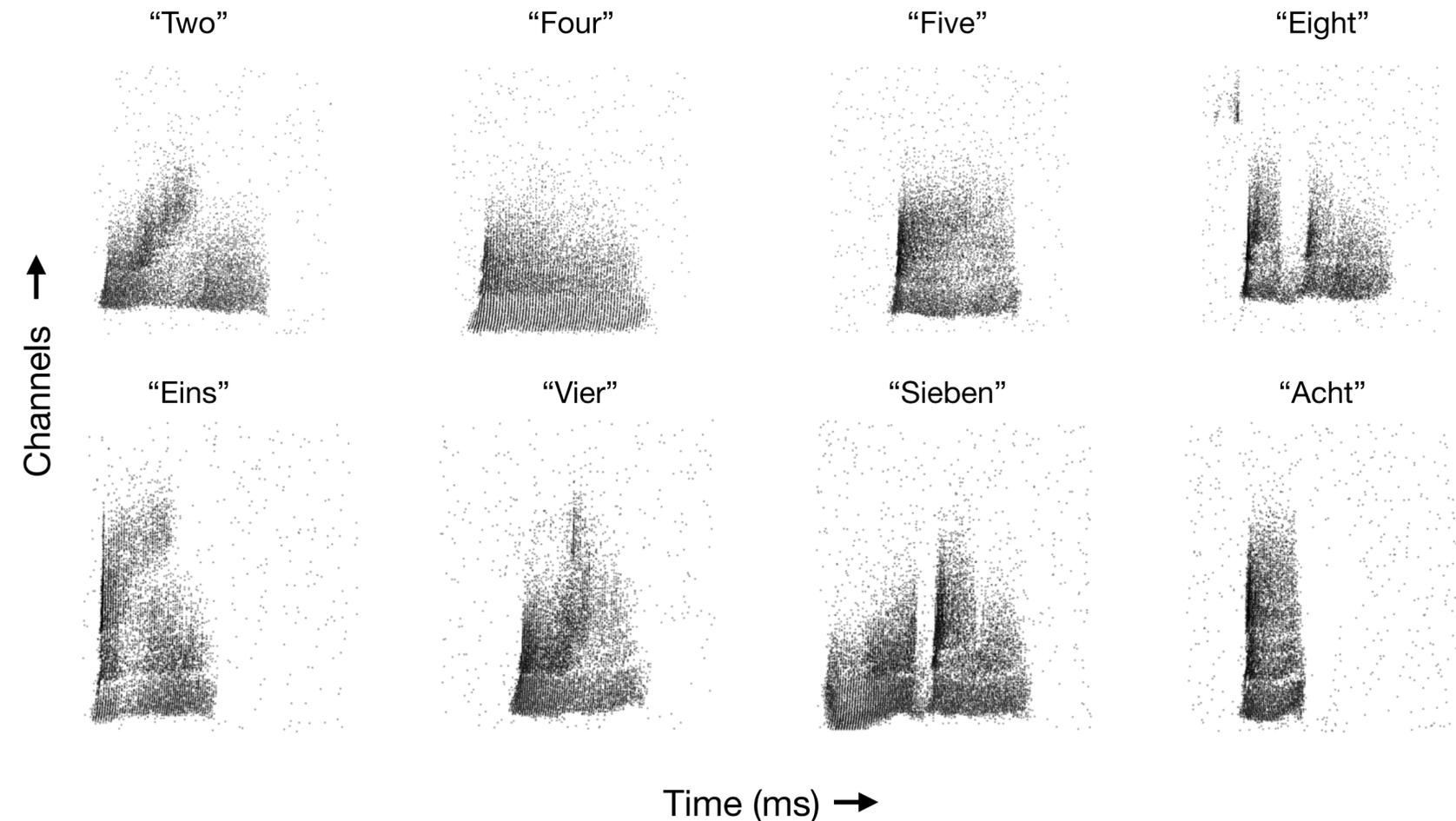
## *Meta-testing phase*



# Dataset

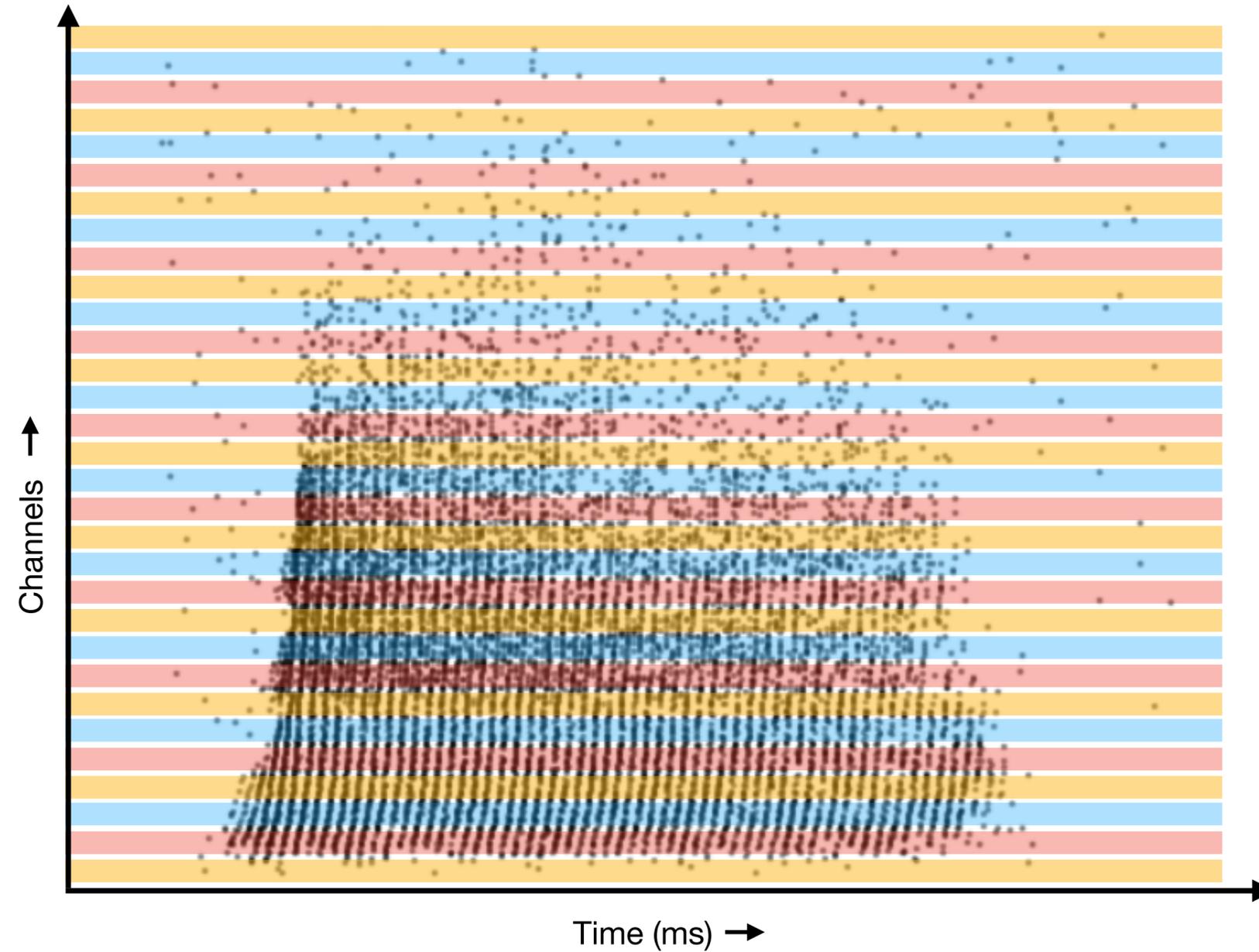
## *Spiking Heidelberg Digits*

- Audio-based classification dataset
- 700 frequency channels, 100 time steps
- “Zero”, “One”, “Two”, ..... , “Nine”  
“Null”, “Eins”, “Zwei”, ..... , “Neun”
- 12 speakers in total (*2 speakers exclusive to the testing set*)
- Class-balanced; 8k training and 2k testing samples



# Data augmentation

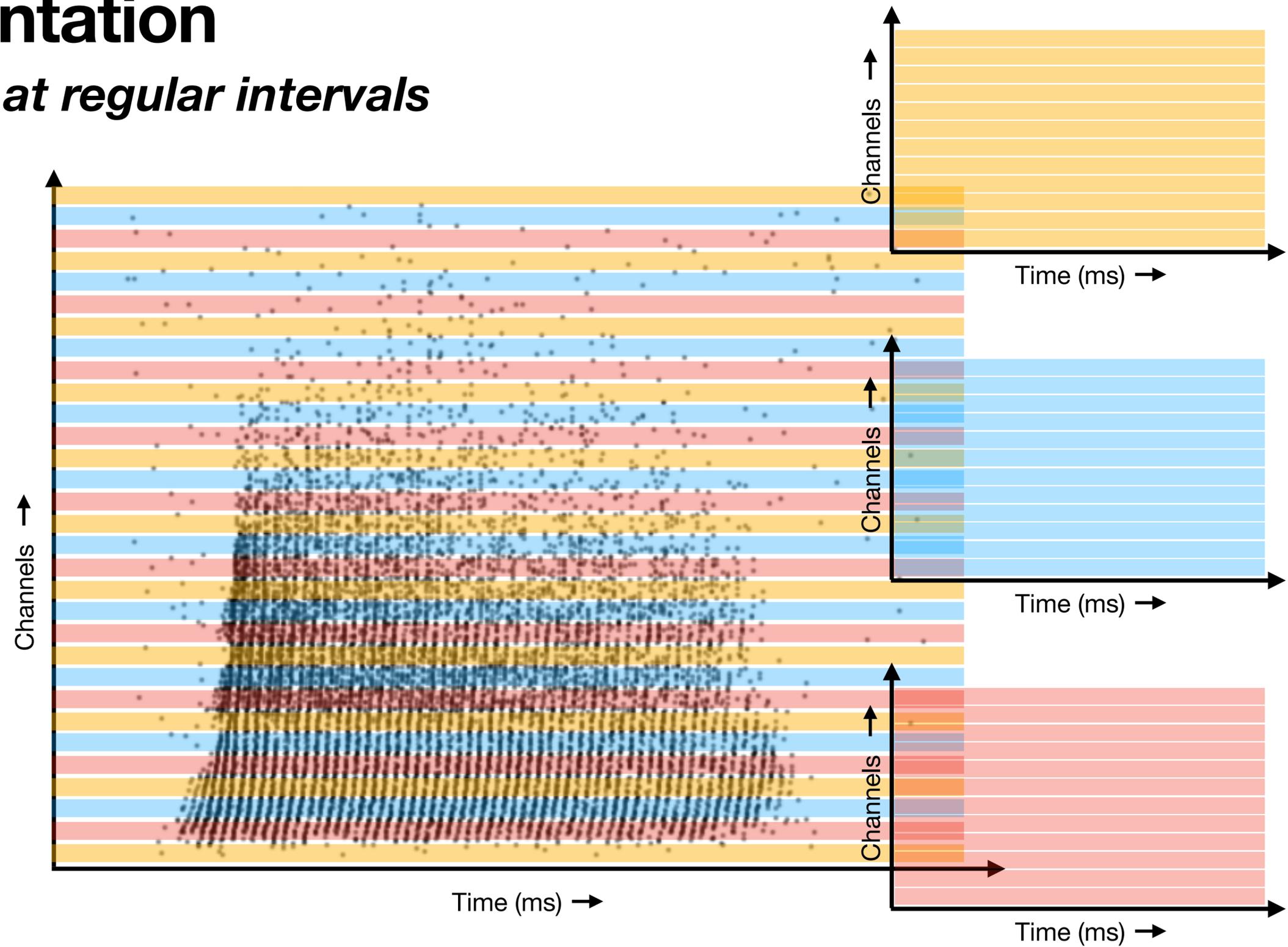
*Spatial sampling at regular intervals*



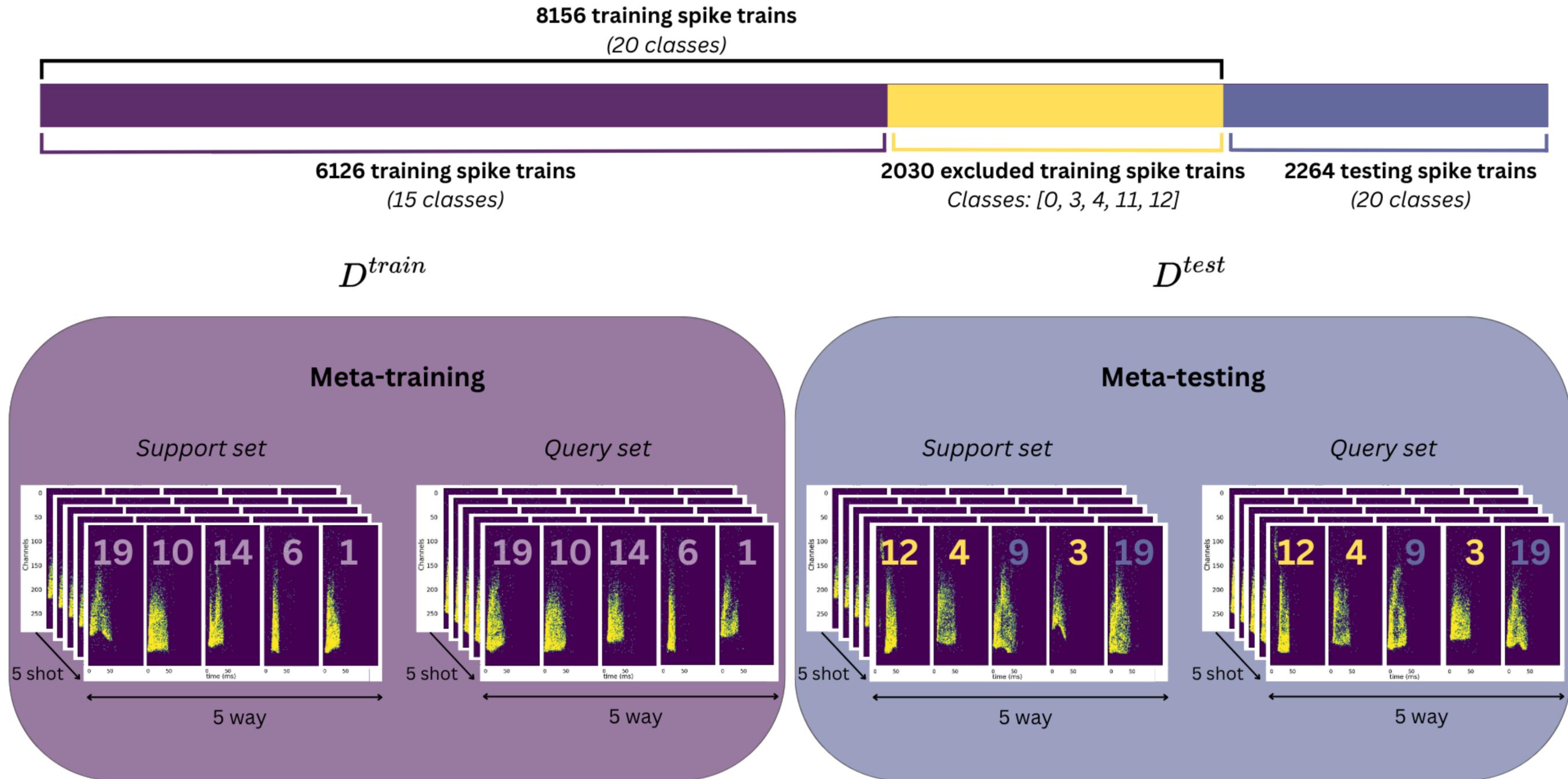
\* Not to scale

# Data augmentation

*Spatial sampling at regular intervals*

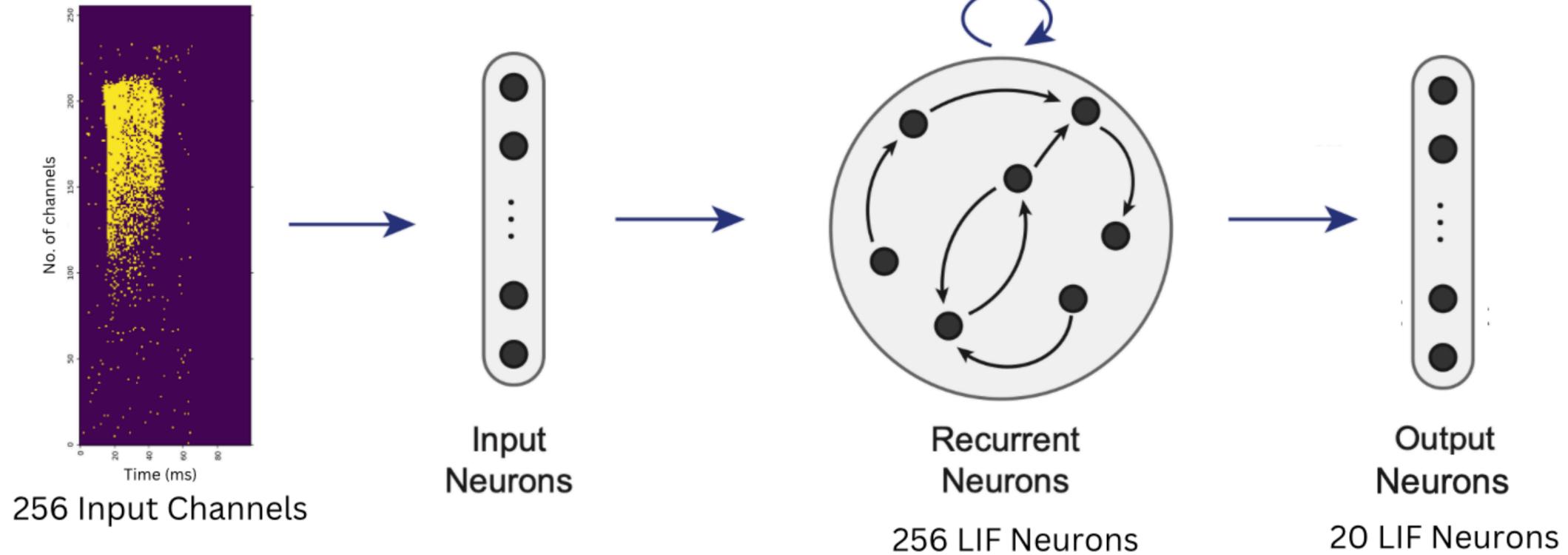


# Converting the SHD into a FSL dataset

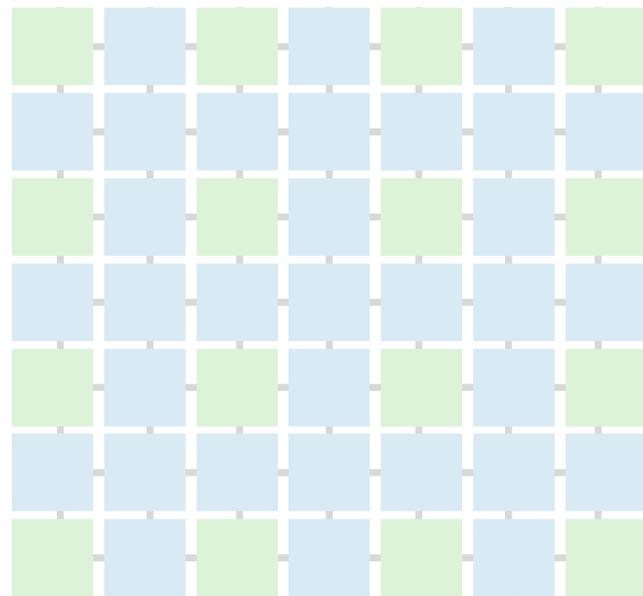


# Model Architecture

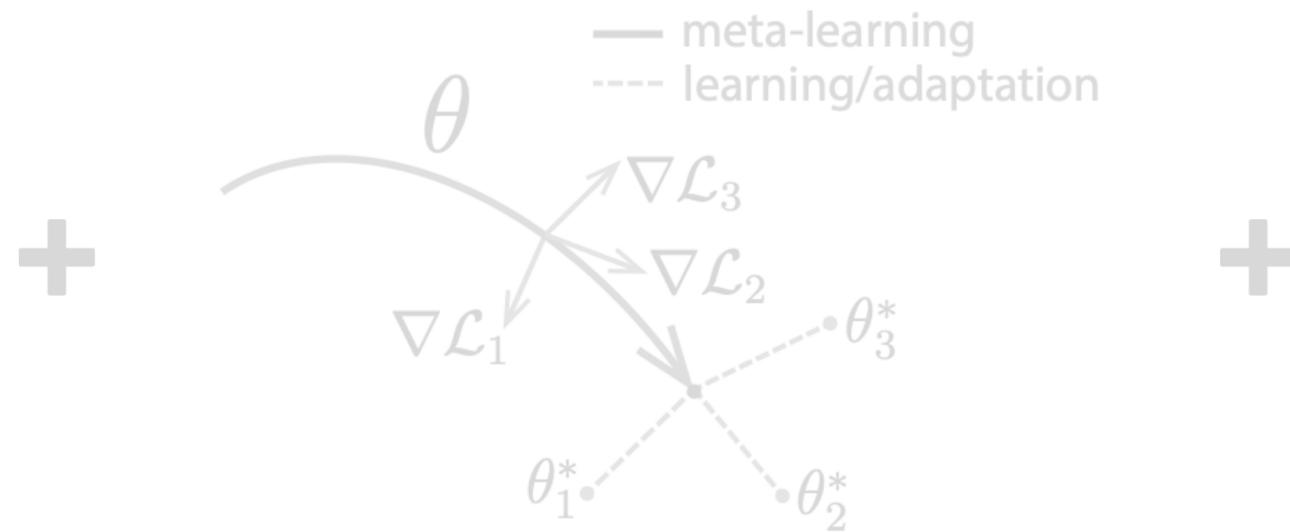
- Layer Sequen
- Inner loop: S
- Outer loop: A
- Softmax Cro
- BPTT using s  
compatible w



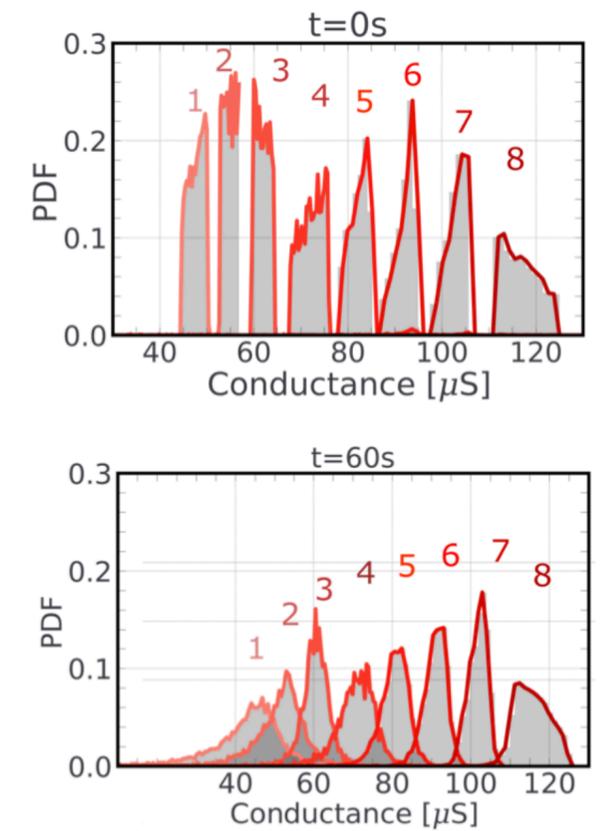
# Hardware-aware training



Novel energy efficient architecture



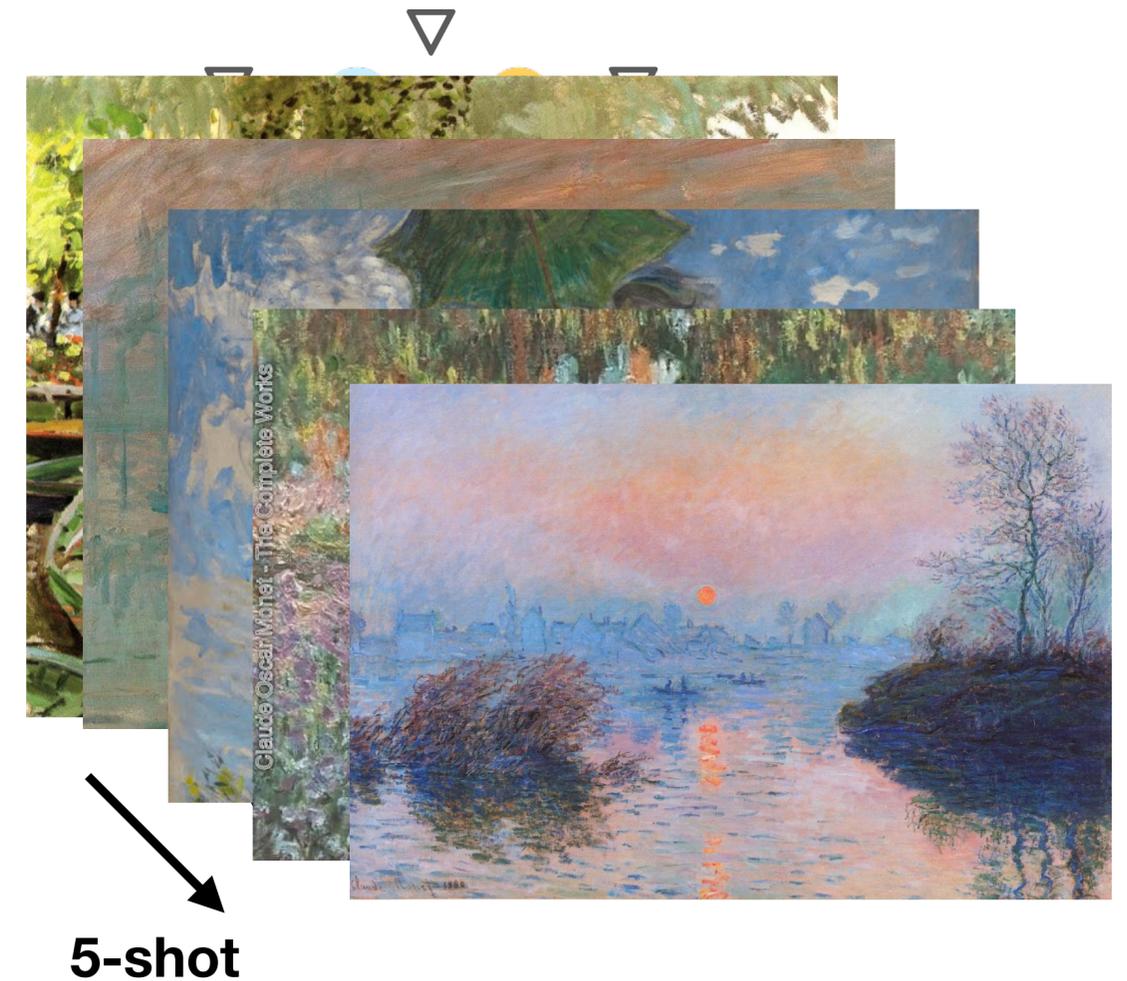
Meta-learning



Hardware-aware training

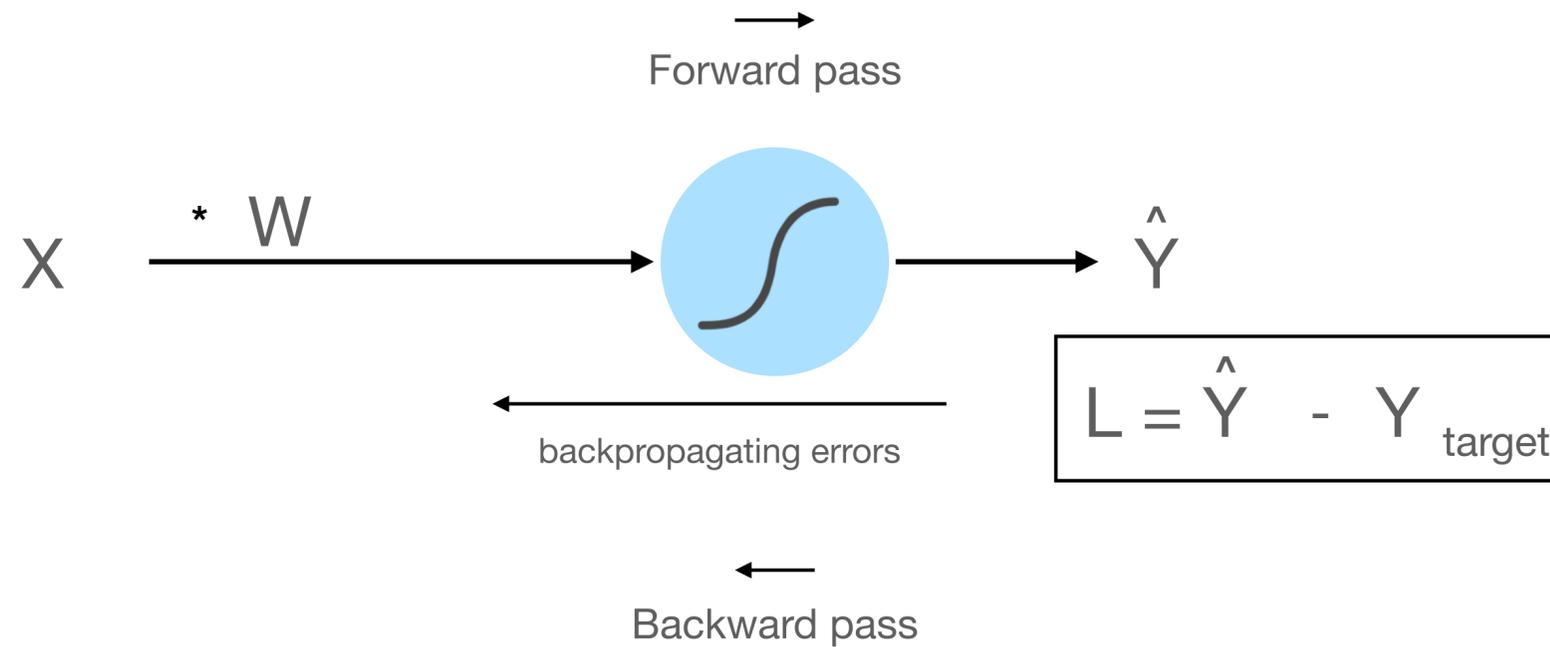
# Hardware based constraints

- Update only the output layer weights in the inner loop
- One inner loop gradient step
- 1-shot learning is ideal
- A very high inner loop learning rate

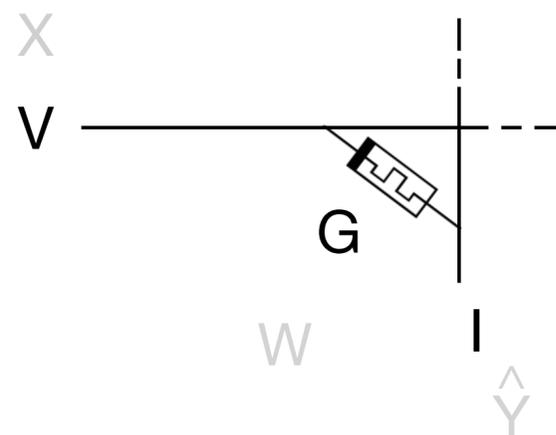


# RRAM-aware training

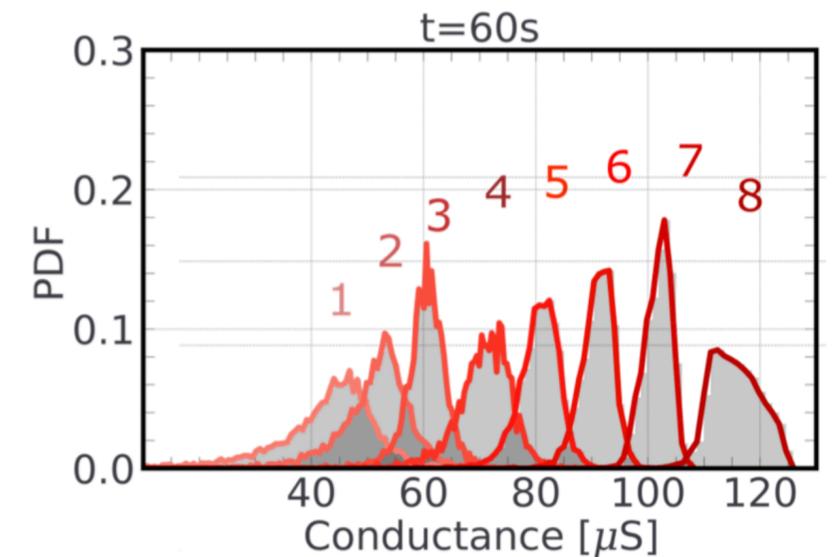
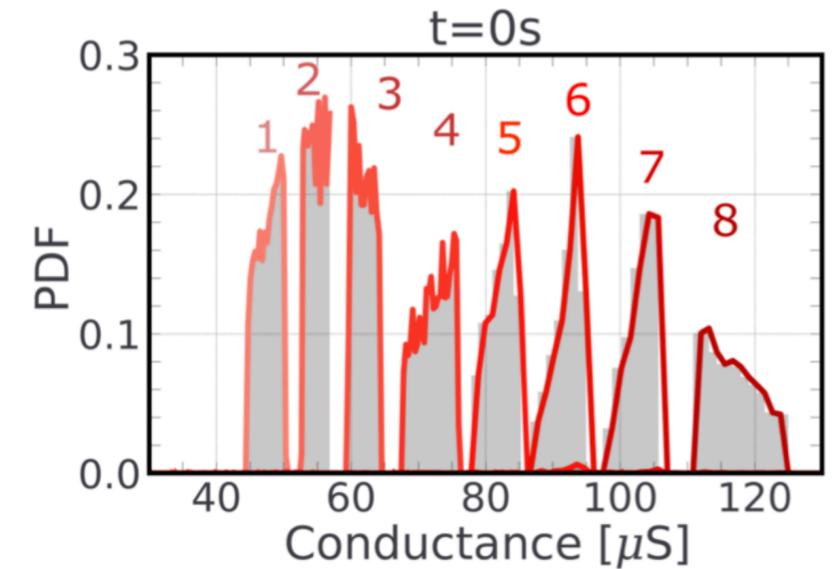
Learning the weights in a Neural network



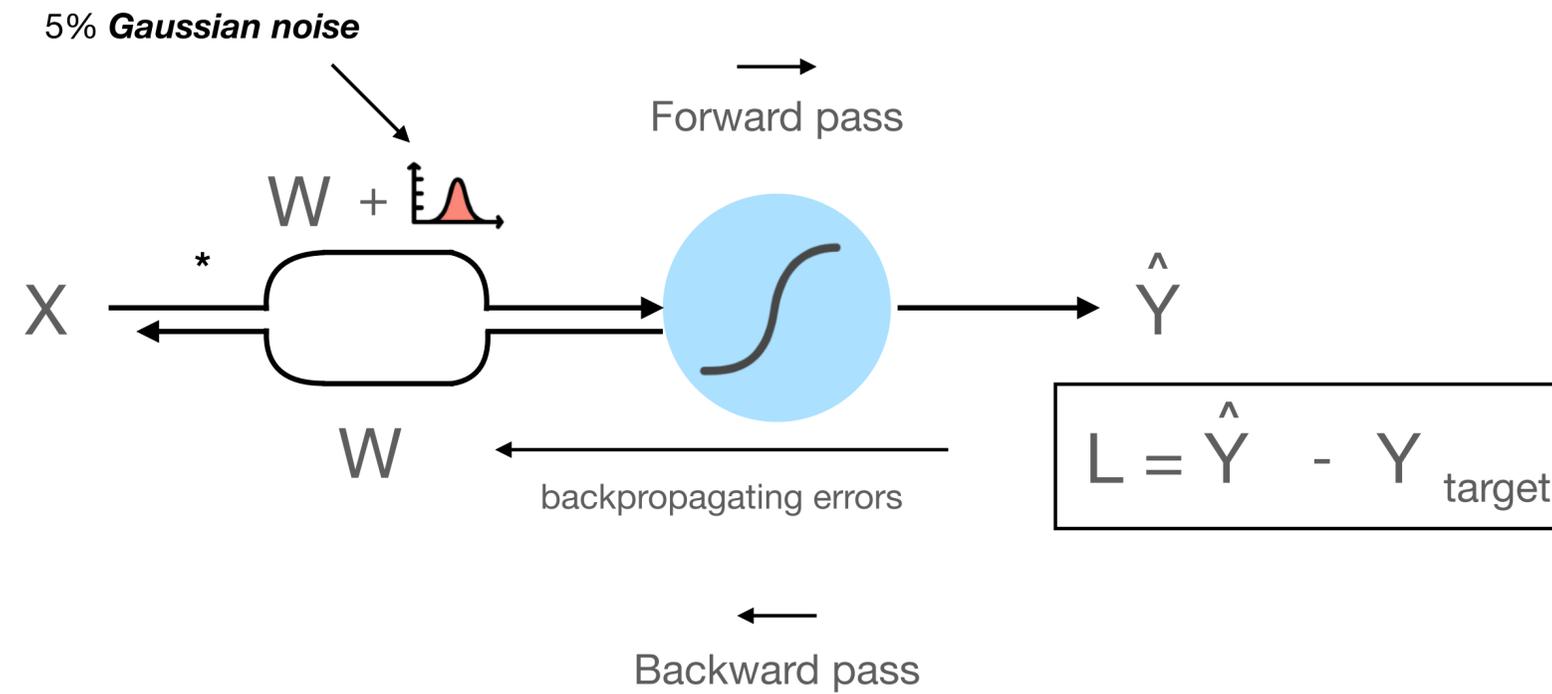
Operation on Hardware



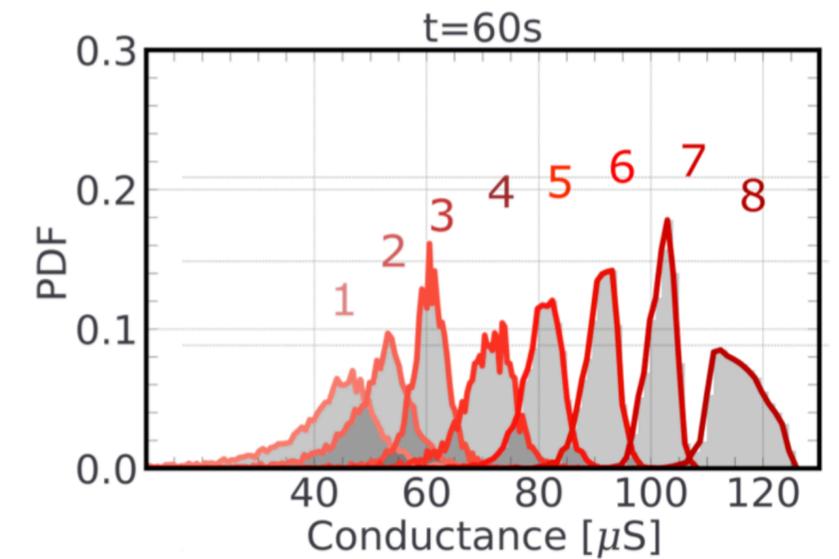
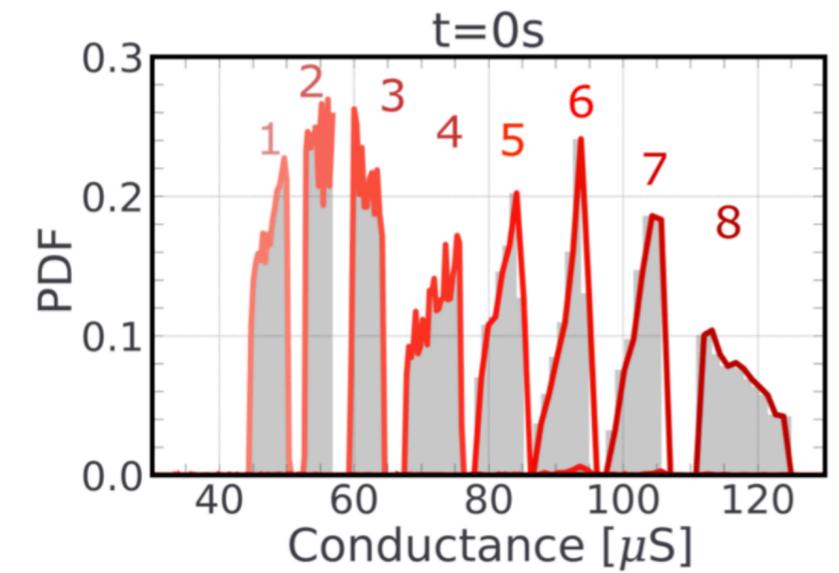
Conductance drift on RRAMs



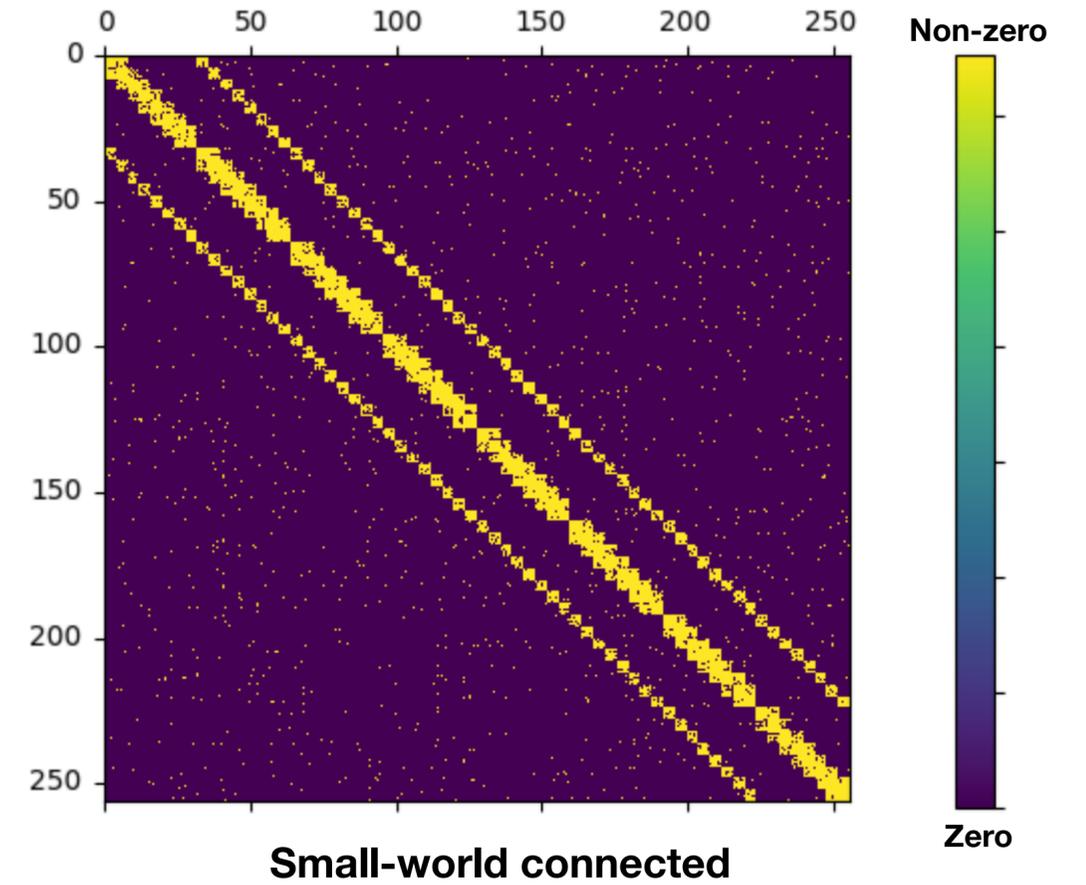
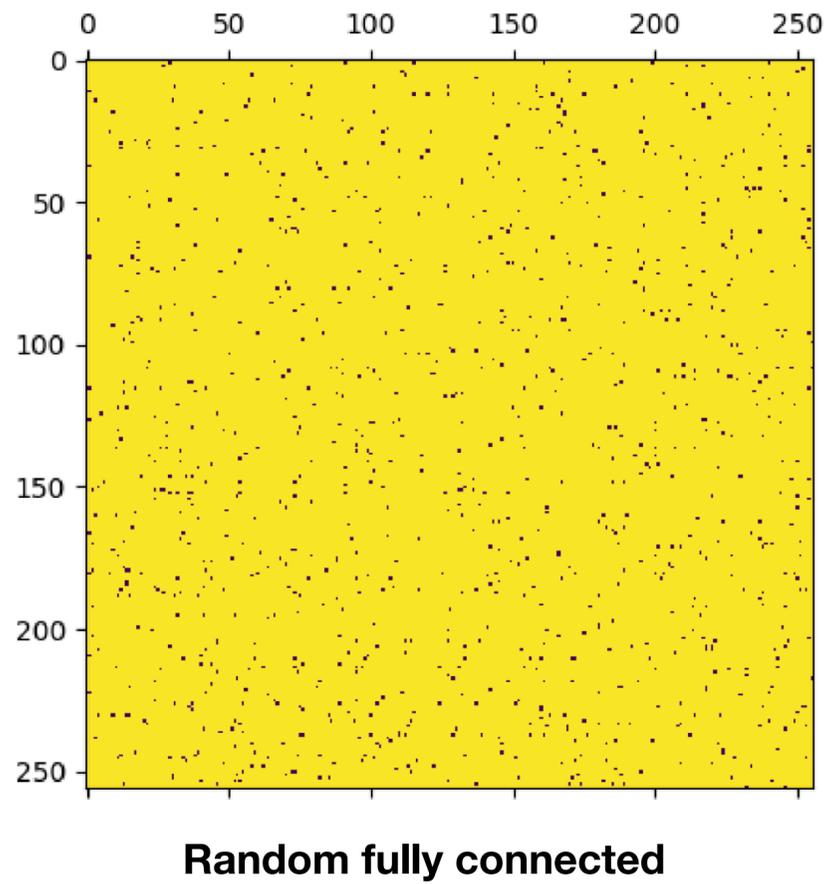
# RRAM-aware training



**Conductance drift on RRAMs**

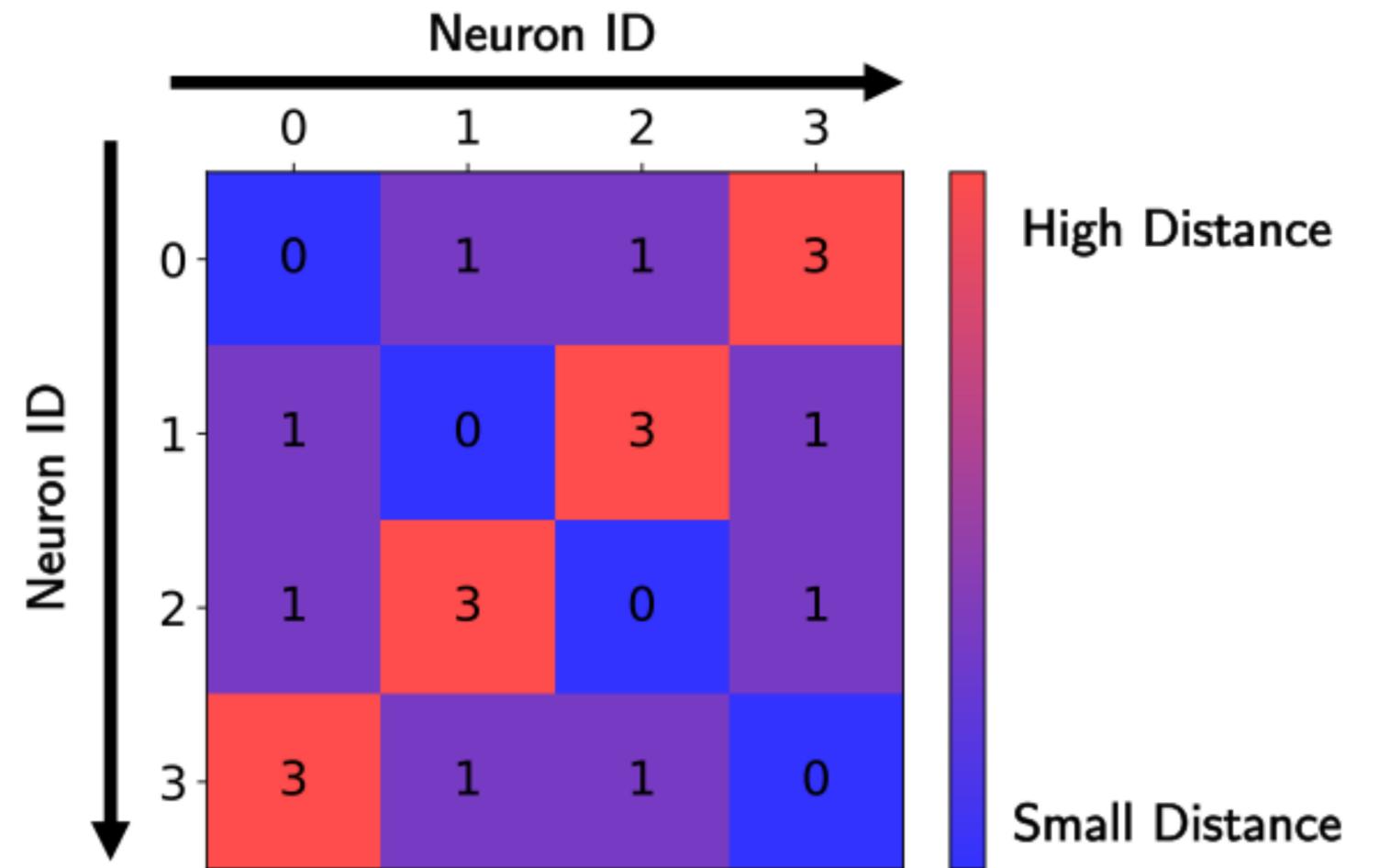
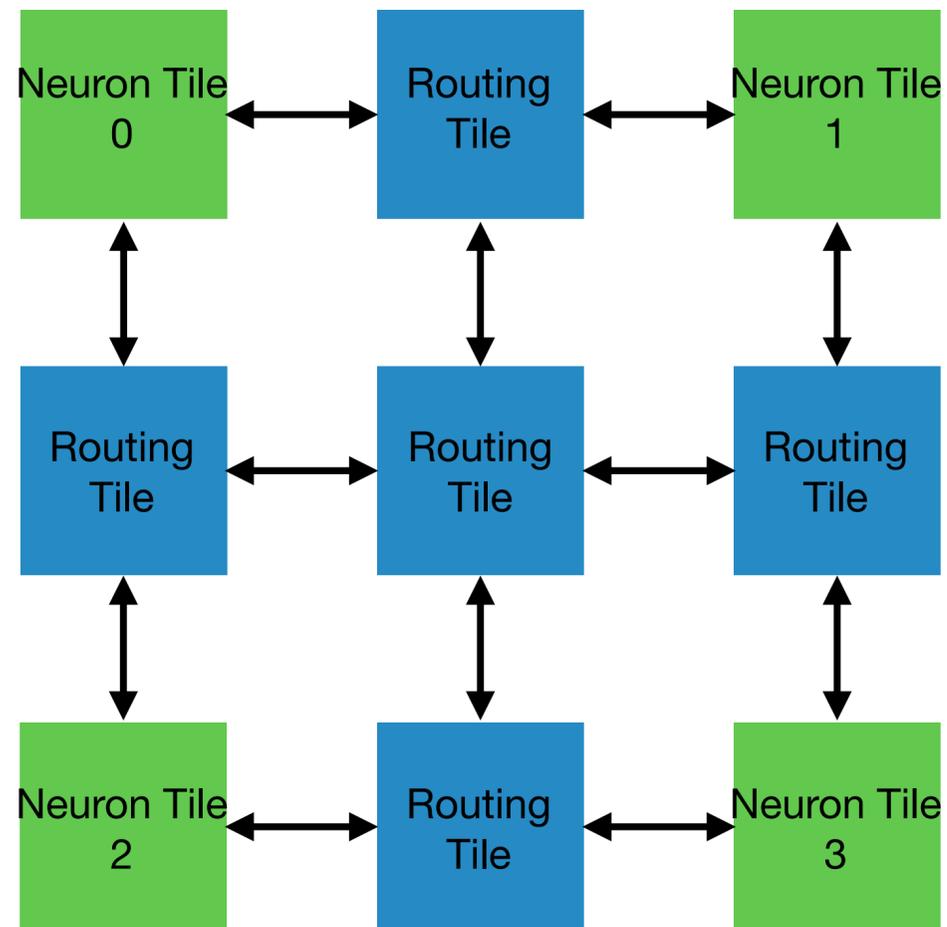


# Mosaic layout-aware training

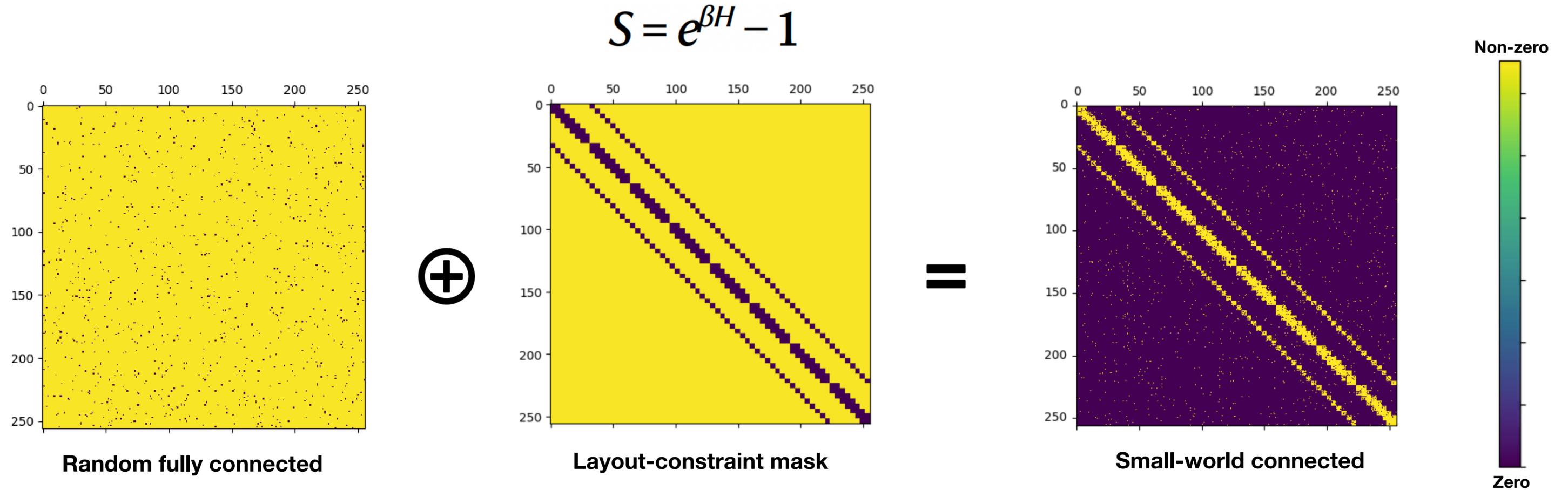


# Calculating the layout-constraint mask

## Example of matrix $H$



# Mosaic layout-aware training



$$L_M = \sum S \odot W^2$$

Total loss  $\longrightarrow$

$$L = L_{CE} + \lambda L_M$$

$\uparrow$        $\longleftarrow$

Cross entropy loss      Layout loss

# Results



# Experimental Setup

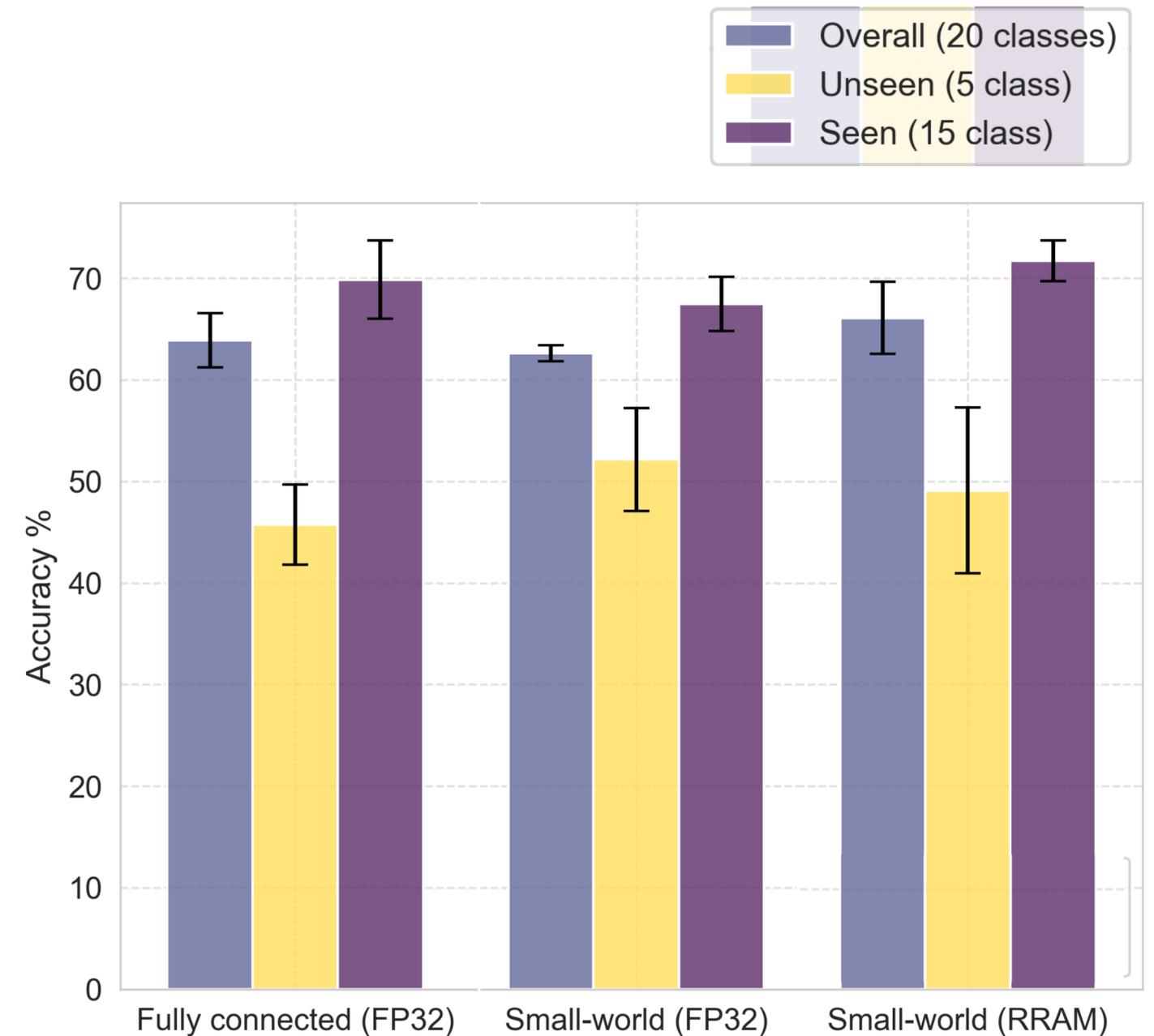
- 1 layer SRNN with 256 LIF neurons (64 neurons for fully-connected baseline)
- Trained on 15 classes (5 excluded classes); Tested on all 20 classes
- 5 way 5 shot tasks
- One inner loop update
- Report the accuracy on seen and unseen classes separately



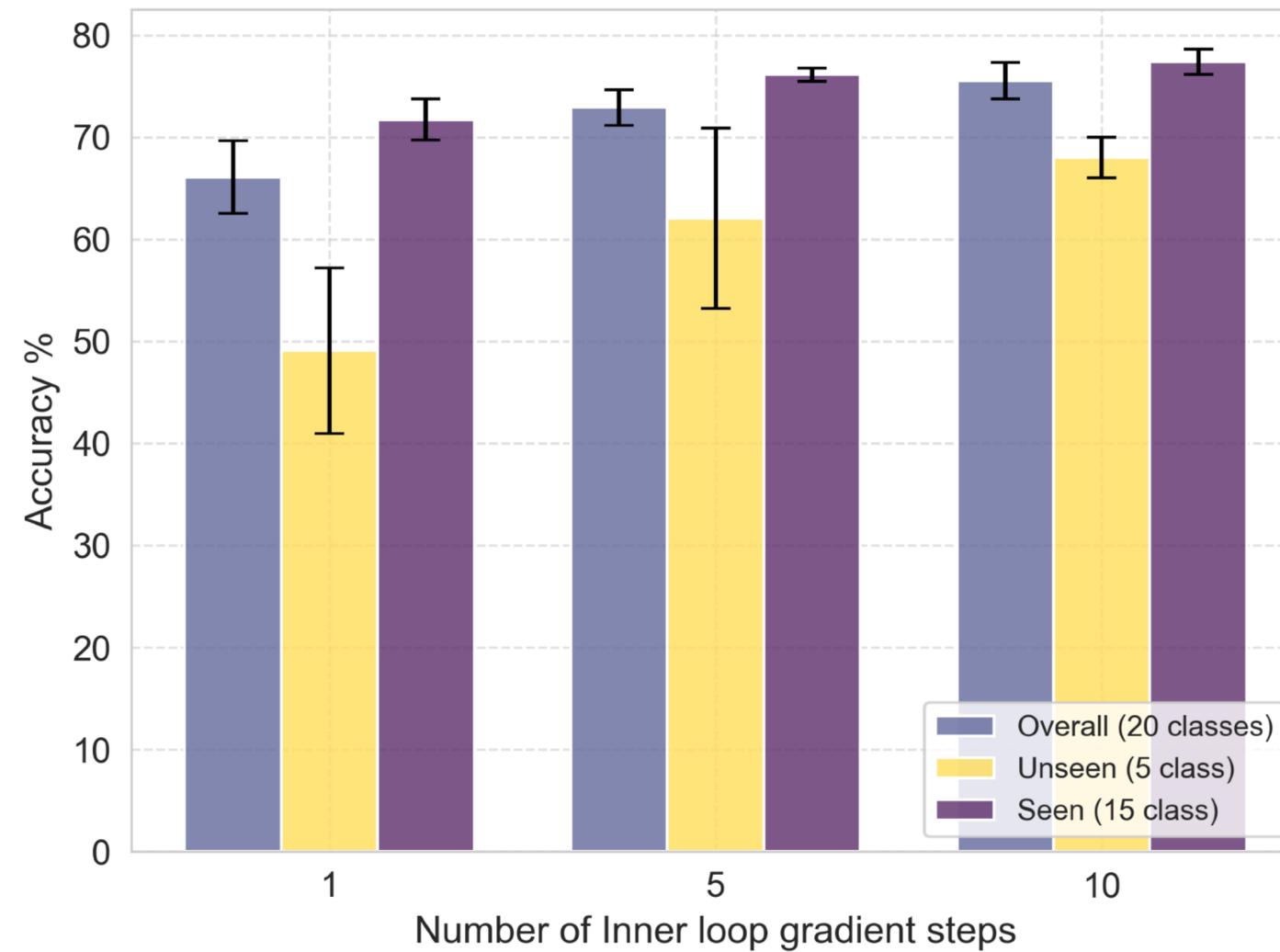
# Constrained network architecture does *not* reduce the performance of the model!

*(Compared to a fully connected network)*

- 70% accuracy, standard SHD benchmark for 1 layer SRNN
- Increase in performance on unseen classes (attributed to sparse dist. of parameters)
- Best performance on seen classes (attributed to noise aware training)



# Increasing inner loop gradient steps increases performance!



# Discussion



# Discussion

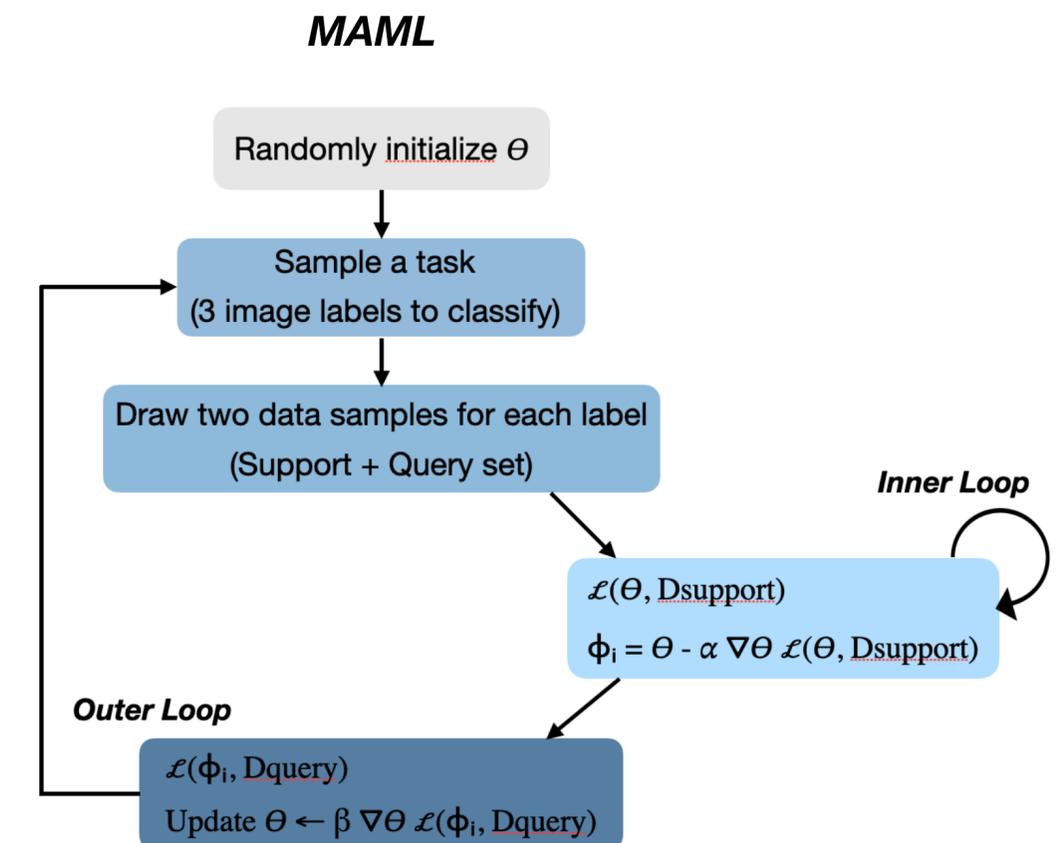
## *Interesting relationship between learning rates and no. of inner loop updates*

➔ 1 inner loop gradient step

Config parameter	Importance ⓘ ↓
lr_out	<div style="width: 80%;"></div>
Runtime	<div style="width: 20%;"></div>
lr_in	<div style="width: 5%;"></div>

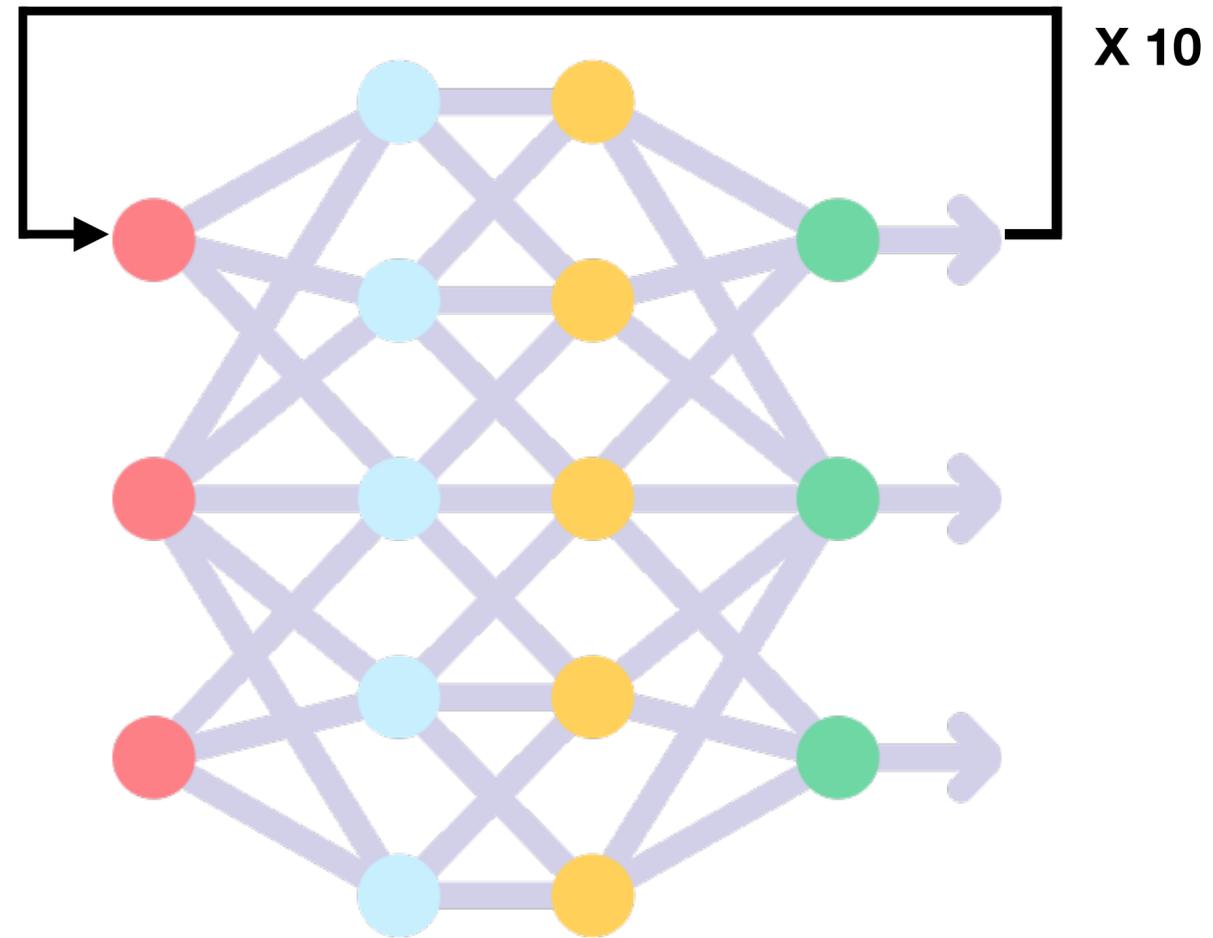
➔ 10 inner loop gradient steps

Config parameter	Importance ⓘ ↓
lr_in	<div style="width: 80%;"></div>
Runtime	<div style="width: 20%;"></div>
lr_out	<div style="width: 20%;"></div>



# Discussion

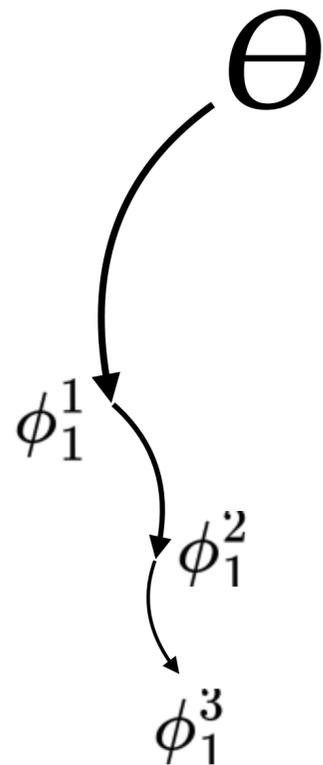
*Can re-training compensate for multiple gradient steps?*



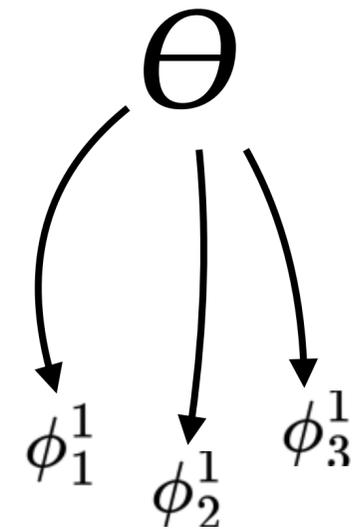
# Discussion

***Can re-training compensate for multiple gradient steps?***

*3 inner loop updates*

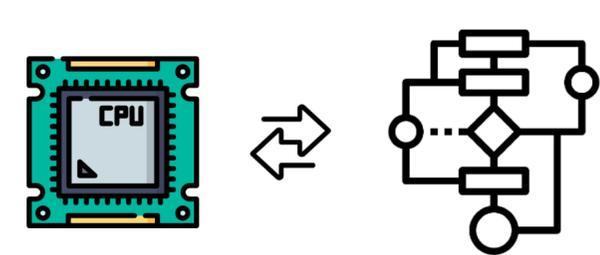
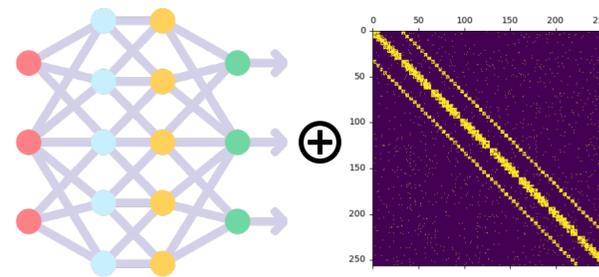
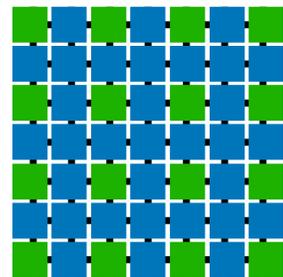


*Re-training*



# Take home message

- Demonstrated (using simulations) Online few-shot learning on the Mosaic architecture
- Constraining network architecture does not terminally detriment the performance of the model
- Importance of inner loop dynamics in MAML; Highlights the need for HW-SW co-design for efficient



# Acknowledgements



Dr. Melika Payvand  
Institute of Neuroinformatics



Yigit Demirag  
Institute of Neuroinformatics



Prof. Dr. Herbert Jaeger  
University of Groningen



Dr. Emre Nefci  
Forschungszentrum Jülich



**Emerging Intelligent Substrates lab**



This work is supported by Horizon Europe  
**METASPIN** project (grant number 101098651)  
and the SNSF Starting Grant Project UNITE  
(TMSGI2- 211461)