

[embodied-learning-vision-course.github.io](https://embodied-learning-vision-course.github.io)



# Lecture Slides for Note Taking



Module 5:  
Continual Learning,  
Few-Shot Learning, Meta-Learning

# Why Continual Learning?

- The world is not a dataset that allows you to get IID samples.



# Why Continual Learning?

- The world is not a dataset that allows you to get IID samples.
- The world keeps changing and evolving.

# Why Continual Learning?

- The world is not a dataset that allows you to get IID samples.
- The world keeps changing and evolving.
- Online vs. Continual
  - Online means that samples arrive in a streaming / temporal partial order, but they may still come from a static distribution.

$$\theta_t = f(x_t, \theta_{t-1}) \quad x_{1:T} \sim \mathcal{X}$$

- Example: Online reinforcement learning, trajectory roll out is online, but the environment is the same.
- Continual learning means that there will be distribution shift.

# What is Continual Learning?

- Distribution shift: Forgetting
  - Learning on A and then B, results in worse performance on A.

# What is Continual Learning?

- Distribution shift: Forgetting
  - Learning on A and then B, results in worse performance on A.
- Multi-task learning: Forward transfer
  - Learning Task A + B results in better learning in Task C compared to learning C alone.
  - Leverage the similarity between tasks.

# What is Continual Learning?

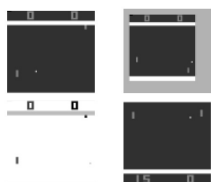
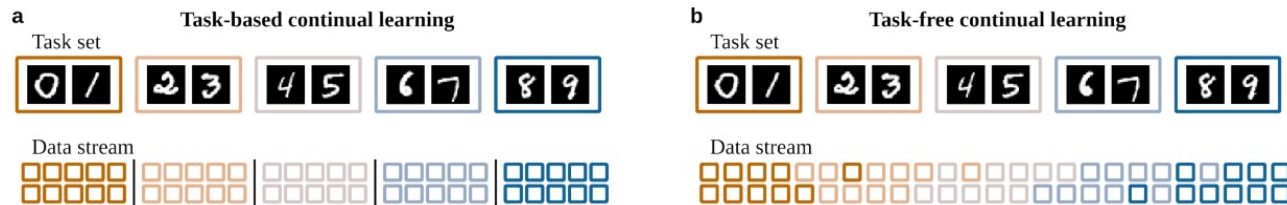
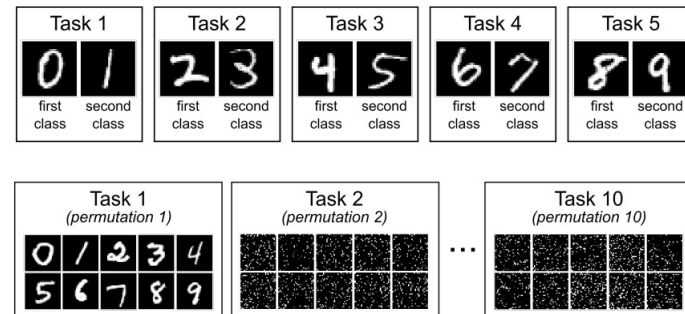
- Distribution shift: Forgetting
  - Learning on A and then B, results in worse performance on A.
- Multi-task learning: Forward transfer
  - Learning Task A + B results in better learning in Task C compared to learning C alone.
  - Leverage the similarity between tasks.
- Compositionality
  - Learning A and B first, and then learning tasks with composed A+B.

# What is Continual Learning?

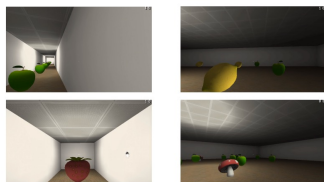
- Distribution shift: Forgetting
  - Learning on A and then B, results in worse performance on A.
- Multi-task learning: Forward transfer
  - Learning Task A + B results in better learning in Task C compared to learning C alone.
  - Leverage the similarity between tasks.
- Compositionality
  - Learning A and B first, and then learning tasks with composed A+B.
- Incremental/curriculum Learning
  - Learning A->B->C is easier than at random order.

# Continual Learning

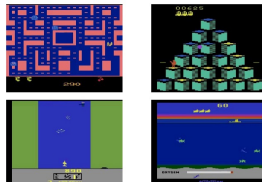
- Learning a sequence of tasks without looking back.



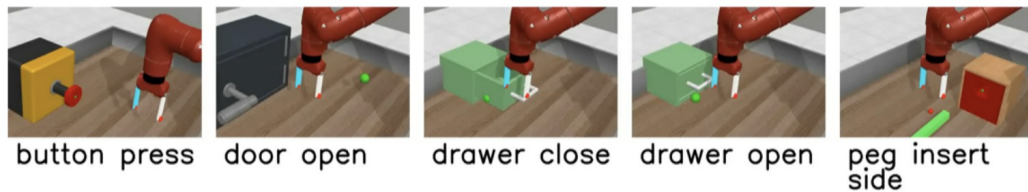
(a) Pong variants



(b) Labyrinth games

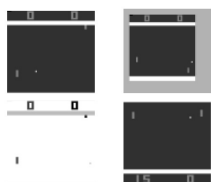
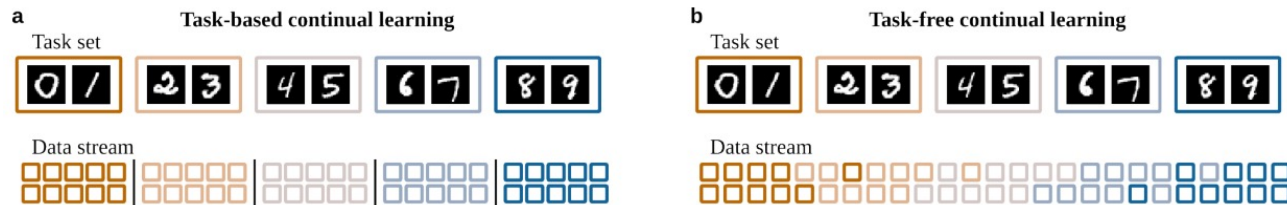
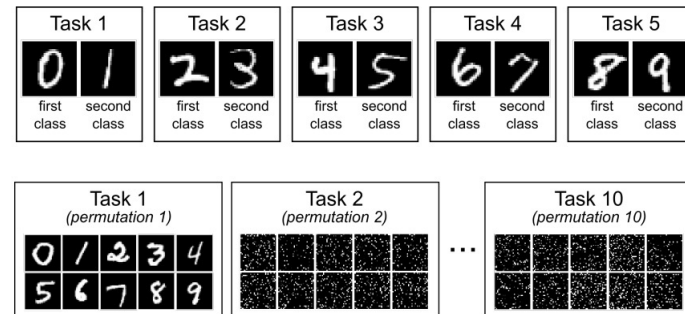


(c) Atari games

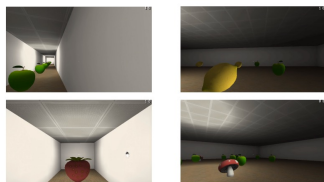


# Continual Learning

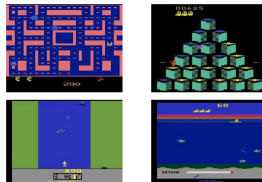
- Learning a sequence of tasks without looking back.
- Goal is to do well on all of the tasks at the end.



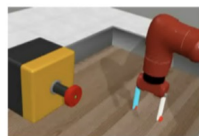
(a) Pong variants



(b) Labyrinth games



(c) Atari games



button press



door open



drawer close



drawer open

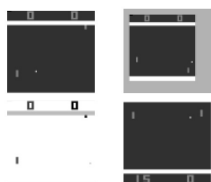
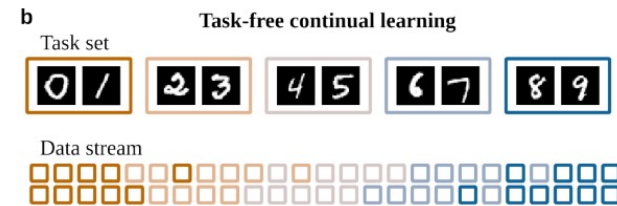
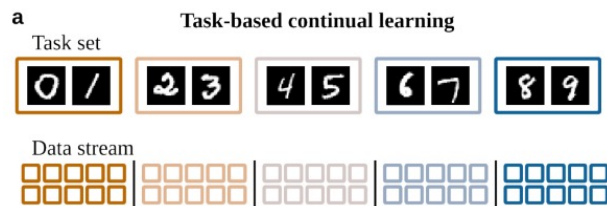
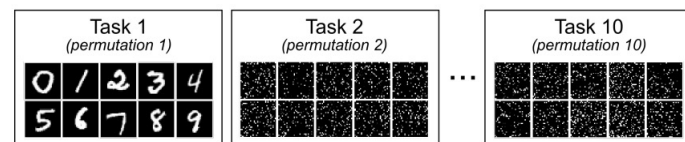
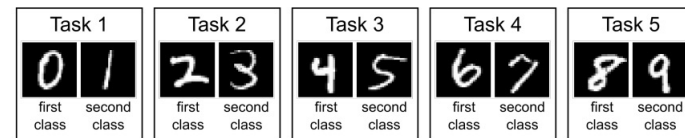


peg insert side

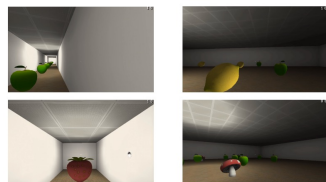


# Continual Learning

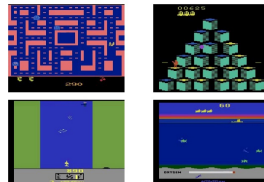
- Learning a sequence of tasks without looking back.
- Goal is to do well on all of the tasks at the end.
- Task boundary



(a) Pong variants



(b) Labyrinth games



(c) Atari games



button press



door open



drawer close



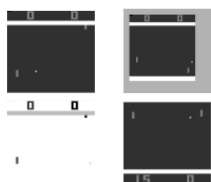
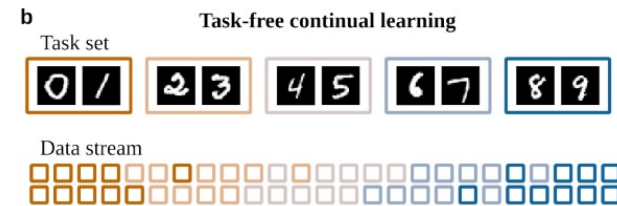
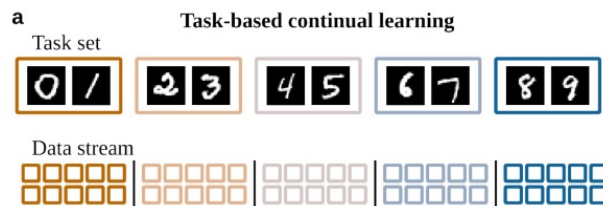
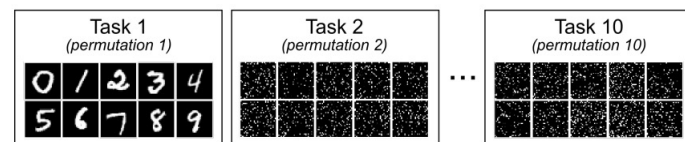
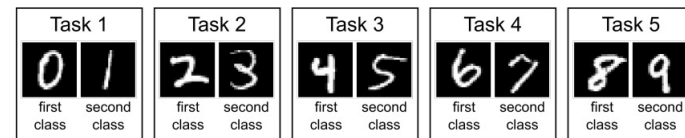
drawer open



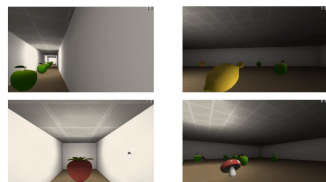
peg insert side

# Continual Learning

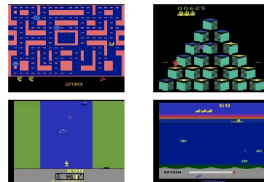
- Learning a sequence of tasks without looking back.
- Goal is to do well on all of the tasks at the end.
- Task boundary
- Memory constraints



(a) Pong variants



(b) Labyrinth games



(c) Atari games



button press



door open



drawer close



drawer open



peg insert side

# Parameter Regularization

- Over-completeness Assumption. A multitude of models can reach equivalent performance.

$$\mathcal{S}_A = \{\theta \mid \ell_A(\theta) < \epsilon\}$$

$$\mathcal{S}_A \cap \mathcal{S}_B \neq \emptyset$$

# Parameter Regularization

- Over-completeness Assumption. A multitude of models can reach equivalent performance.
- What is left is to efficiently find the intersection between A and B.

$$\mathcal{S}_A = \{\theta \mid \ell_A(\theta) < \epsilon\}$$

$$\mathcal{S}_A \cap \mathcal{S}_B \neq \emptyset$$

$$p(\theta \mid \mathcal{D}_A) = \mathcal{N}(\theta; \theta^*, \Sigma)$$

# Parameter Regularization

- Over-completeness Assumption. A multitude of models can reach equivalent performance.
- What is left is to efficiently find the intersection between A and B.

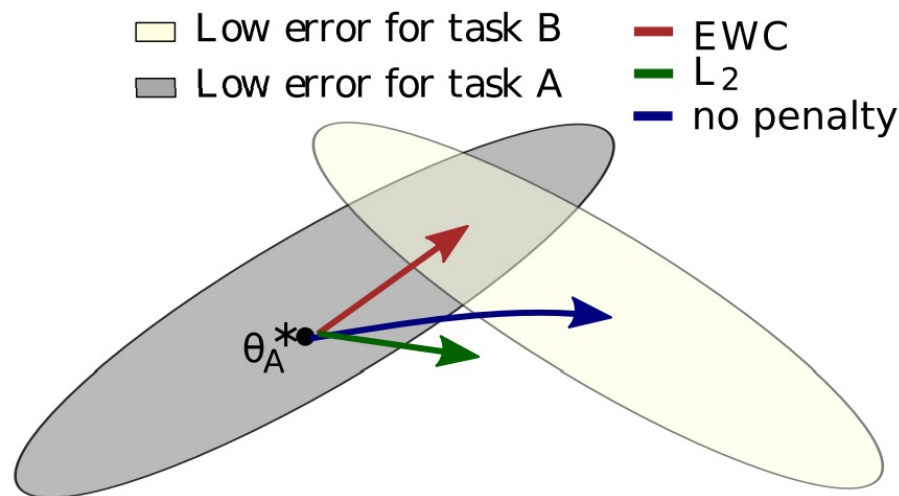
$$p(\theta \mid \mathcal{D}_A) = \mathcal{N}(\theta; \theta^*, \Sigma)$$

- Elastic Weight Consolidation (EWC):

$$\mathcal{L}(\theta) = \mathcal{L}_B(\theta) + \sum_i \frac{\lambda}{2} \underbrace{F_i}_{\text{EWC}} (\theta_i - \theta_{A,i}^*)^2$$

$$\mathcal{S}_A = \{\theta \mid \ell_A(\theta) < \epsilon\}$$

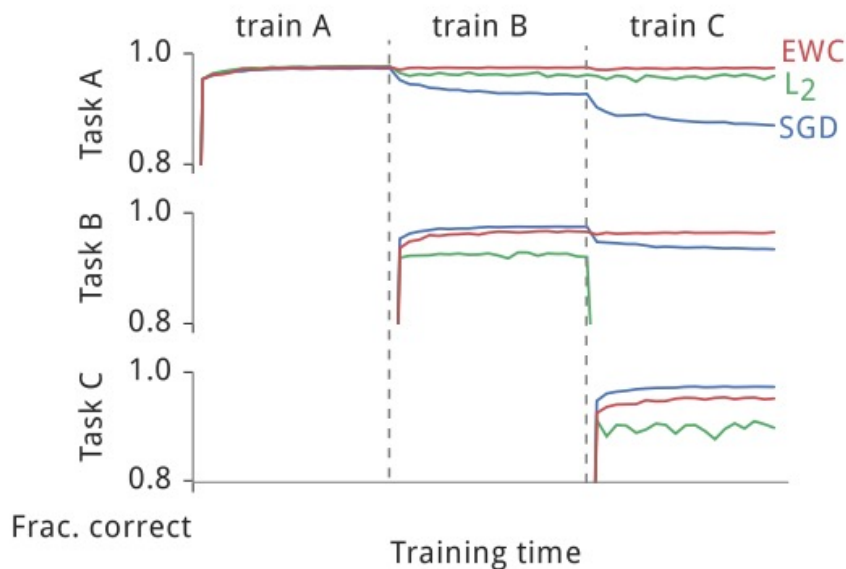
$$\mathcal{S}_A \cap \mathcal{S}_B \neq \emptyset$$



# Computing Fisher

- At the end of each epoch, compute the gradient squared:

$$F_i = \left( \frac{d\mathcal{L}}{d\theta_i} \right)^2$$

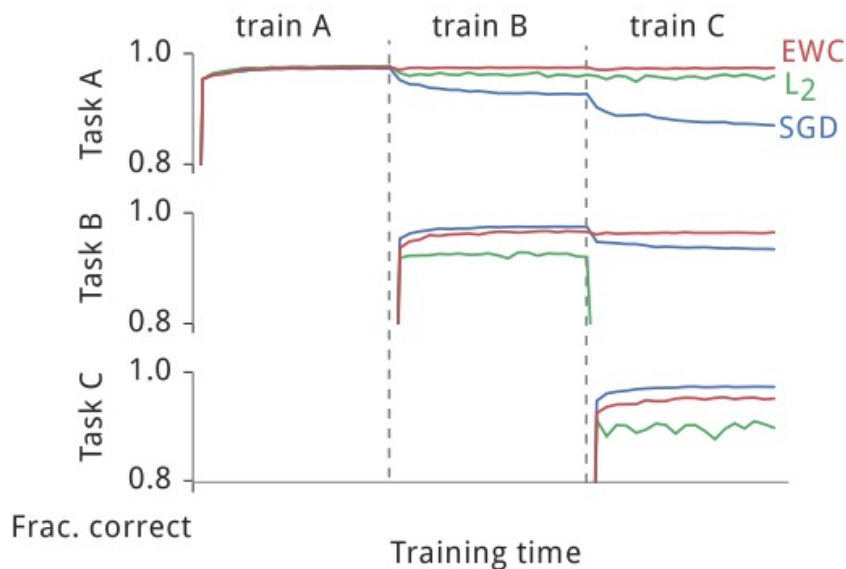


# Computing Fisher

- At the end of each epoch, compute the gradient squared:

$$F_i = \left( \frac{d\mathcal{L}}{d\theta_i} \right)^2$$

- Measures the sensitivity on each parameter dimension.

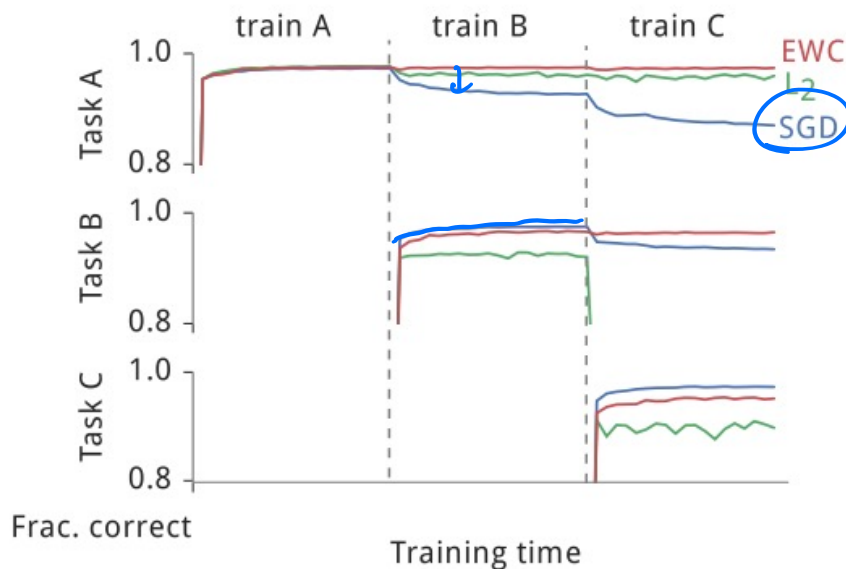


# Computing Fisher

- At the end of each epoch, compute the gradient squared:

$$F_i = \left( \frac{d\mathcal{L}}{d\theta_i} \right)^2$$

- Measures the sensitivity on each parameter dimension.
- You can also accumulate an online estimate.





# Variational Continual Learning (VCL)

- Bayesian formulation:

$$p(\underline{\theta} \mid \underline{\mathcal{D}_{1:T}}) \propto p(\theta) \prod_{t=1}^T p(\mathcal{D}_t \mid \theta) \propto p(\theta \mid \mathcal{D}_{1:T-1}) p(\mathcal{D}_T \mid \theta).$$

# Variational Continual Learning (VCL)

- Bayesian formulation:

$$p(\boldsymbol{\theta} \mid \mathcal{D}_{1:T}) \propto p(\boldsymbol{\theta}) \prod_{t=1}^T p(\mathcal{D}_t \mid \boldsymbol{\theta}) \propto p(\boldsymbol{\theta} \mid \mathcal{D}_{1:T-1}) p(\mathcal{D}_T \mid \boldsymbol{\theta}).$$

- Variational approach:

$$\boxed{q_t(\boldsymbol{\theta})} = \operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL} \left( q(\boldsymbol{\theta}) \parallel \frac{1}{Z_t} \underbrace{q_{t-1}(\boldsymbol{\theta})}_{\text{prior}} \underbrace{p(\mathcal{D}_t \mid \boldsymbol{\theta})}_{\text{likelihood for the current dataset}} \right).$$

# Variational Continual Learning (VCL)

- Bayesian formulation:

$$p(\boldsymbol{\theta} \mid \mathcal{D}_{1:T}) \propto p(\boldsymbol{\theta}) \prod_{t=1}^T p(\mathcal{D}_t \mid \boldsymbol{\theta}) \propto p(\boldsymbol{\theta} \mid \mathcal{D}_{1:T-1}) p(\mathcal{D}_T \mid \boldsymbol{\theta}).$$

- Variational approach:

$$q_t(\boldsymbol{\theta}) = \operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL} \left( q(\boldsymbol{\theta}) \parallel \frac{1}{Z_t} q_{t-1}(\boldsymbol{\theta}) p(\mathcal{D}_t \mid \boldsymbol{\theta}) \right).$$

- Loss:  $\mathcal{L}(q_t(\boldsymbol{\theta})) = \underbrace{\mathbb{E}_{\boldsymbol{\theta} \sim q_t(\boldsymbol{\theta})} [-\log p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})]}_{\text{data loss}} + \underbrace{\operatorname{KL}(q_t(\boldsymbol{\theta}) \parallel q_{t-1}(\boldsymbol{\theta}))}_{\text{KL divergence}}.$

# Variational Continual Learning (VCL)

- Bayesian formulation:

$$p(\boldsymbol{\theta} \mid \mathcal{D}_{1:T}) \propto p(\boldsymbol{\theta}) \prod_{t=1}^T p(\mathcal{D}_t \mid \boldsymbol{\theta}) \propto p(\boldsymbol{\theta} \mid \mathcal{D}_{1:T-1}) p(\mathcal{D}_T \mid \boldsymbol{\theta}).$$

- Variational approach:

$$q_t(\boldsymbol{\theta}) = \operatorname{argmin}_{q \in \mathcal{Q}} \operatorname{KL} \left( q(\boldsymbol{\theta}) \parallel \frac{1}{Z_t} q_{t-1}(\boldsymbol{\theta}) p(\mathcal{D}_t \mid \boldsymbol{\theta}) \right).$$

- Loss:  $\mathcal{L}(q_t(\boldsymbol{\theta})) = \mathbb{E}_{\boldsymbol{\theta} \sim q_t(\boldsymbol{\theta})} [-\log p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})] + \operatorname{KL}(q_t(\boldsymbol{\theta}) \parallel q_{t-1}(\boldsymbol{\theta})).$   
 $q_t(\boldsymbol{\theta}) = \prod_{d=1}^D \mathcal{N}(\theta_{t,d}; \mu_{t,d}, \sigma_{t,d}^2).$

# Variational Continual Learning (VCL)

- Bayesian formulation:

$$p(\boldsymbol{\theta} \mid \mathcal{D}_{1:T}) \propto p(\boldsymbol{\theta}) \prod_{t=1}^T p(\mathcal{D}_t \mid \boldsymbol{\theta}) \propto p(\boldsymbol{\theta} \mid \mathcal{D}_{1:T-1}) p(\mathcal{D}_T \mid \boldsymbol{\theta}).$$

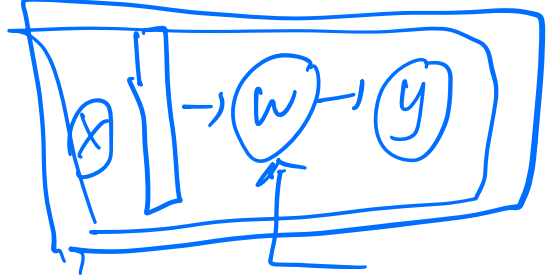
- Variational approach:

$$q_t(\boldsymbol{\theta}) = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \operatorname{KL} \left( q(\boldsymbol{\theta}) \parallel \frac{1}{Z_t} q_{t-1}(\boldsymbol{\theta}) p(\mathcal{D}_t \mid \boldsymbol{\theta}) \right).$$

- Loss:  $\mathcal{L}(q_t(\boldsymbol{\theta})) = \mathbb{E}_{\boldsymbol{\theta} \sim q_t(\boldsymbol{\theta})} [-\log p(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta})] + \operatorname{KL}(q_t(\boldsymbol{\theta}) \parallel q_{t-1}(\boldsymbol{\theta})).$

$$q_t(\boldsymbol{\theta}) = \prod_{d=1}^D \mathcal{N}(\theta_{t,d}; \mu_{t,d}, \sigma_{t,d}^2).$$

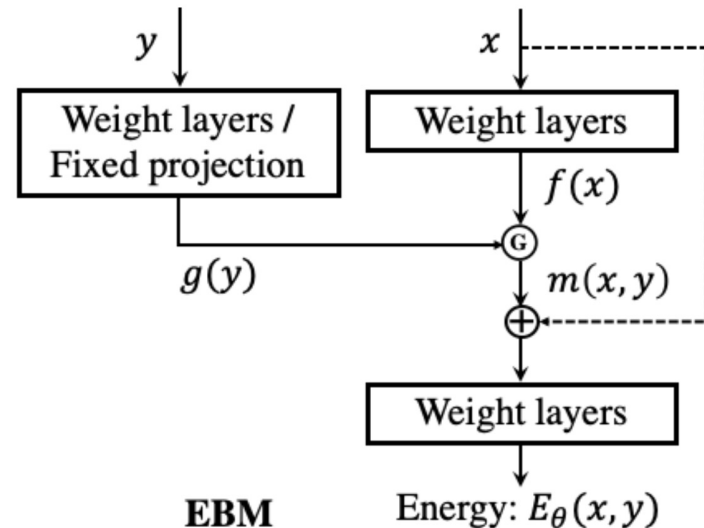
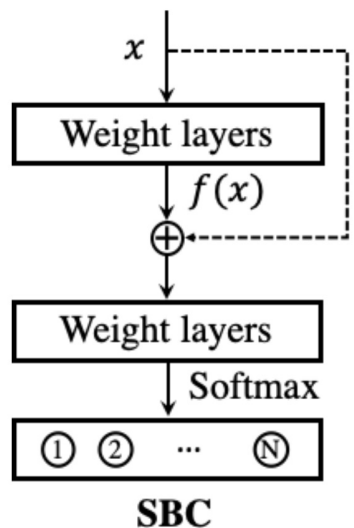
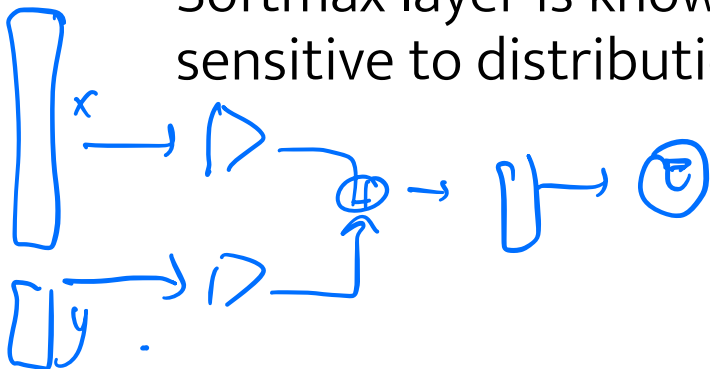
- Compare to EWC: Maintains uncertainty throughout training.



# EBM for Continual Learning

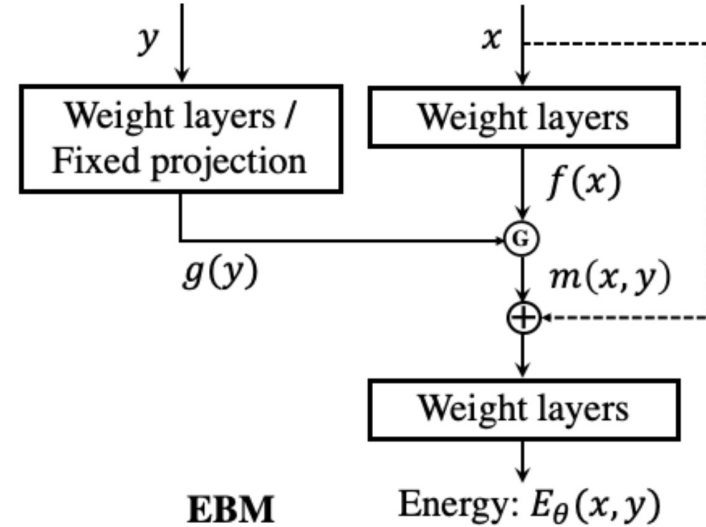
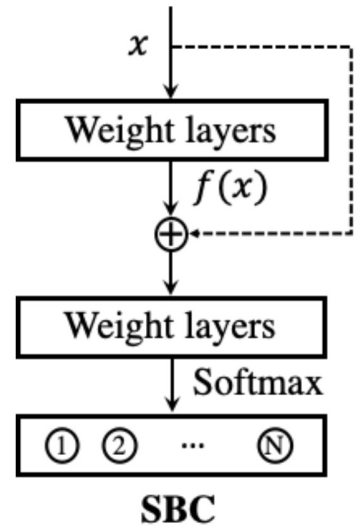
$$E(x, y)$$

- Softmax layer is known to be sensitive to distribution shift.



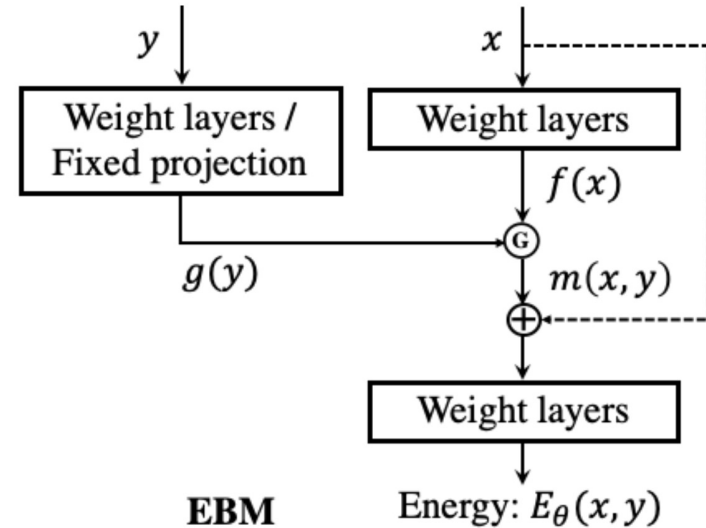
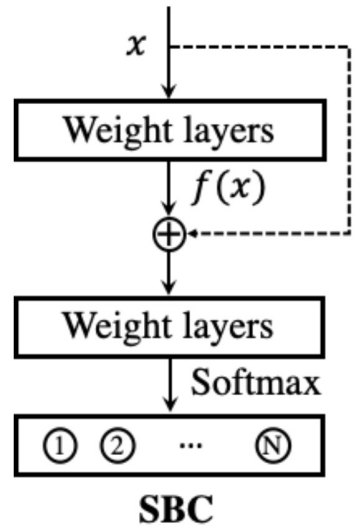
# EBM for Continual Learning

- Softmax layer is known to be sensitive to distribution shift.
- A common approach is to use nearest mean classifier.



# EBM for Continual Learning

- Softmax layer is known to be sensitive to distribution shift.
- A common approach is to use nearest mean classifier.
- Can be generalized to EBMs



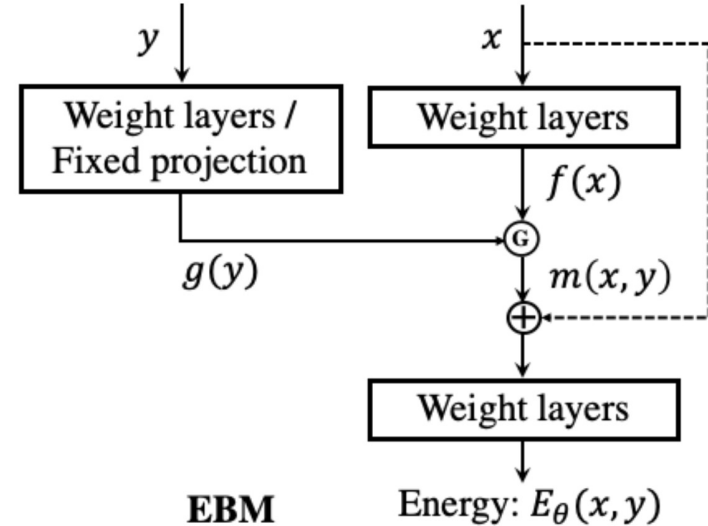
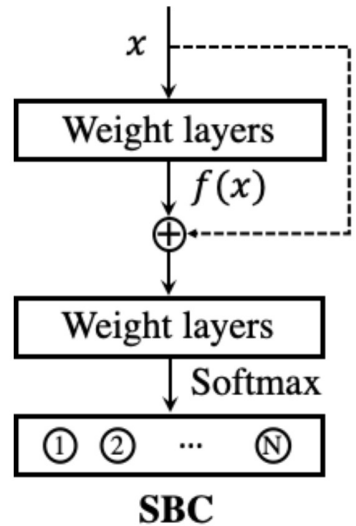


# EBM for Continual Learning

- Softmax layer is known to be sensitive to distribution shift.
- A common approach is to use nearest mean classifier.
- Can be generalized to EBMs
- Energy between inputs and labels:

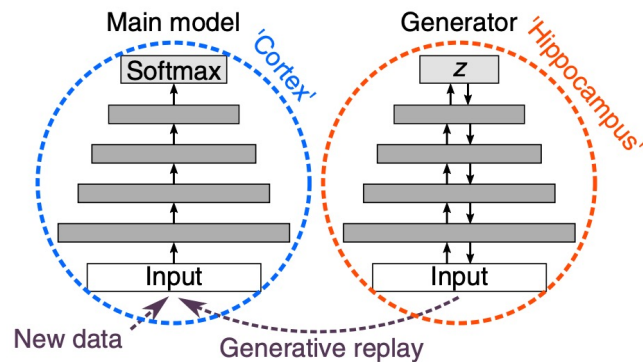
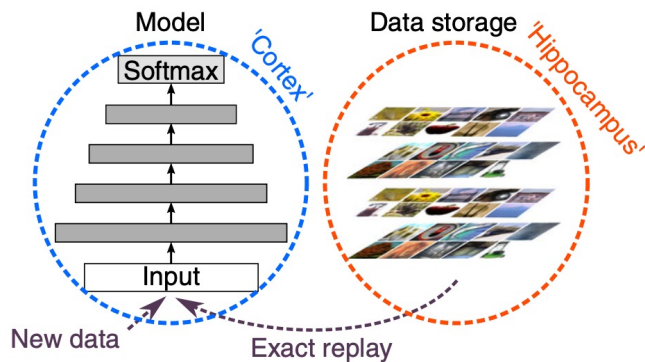
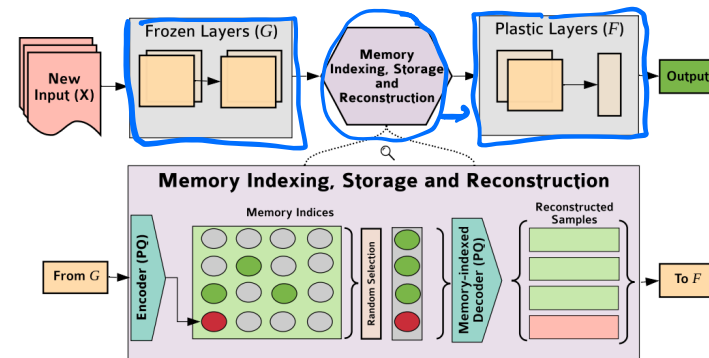
$$m(\mathbf{x}, y) = G(f(\mathbf{x}), g(y))$$

$$L_{CD}(\theta) = \mathbb{E}_{(\mathbf{x}, y, y^-) \sim p_D} [E_{\theta}(\mathbf{x}, y) - E_{\theta}(\mathbf{x}, y^-)].$$

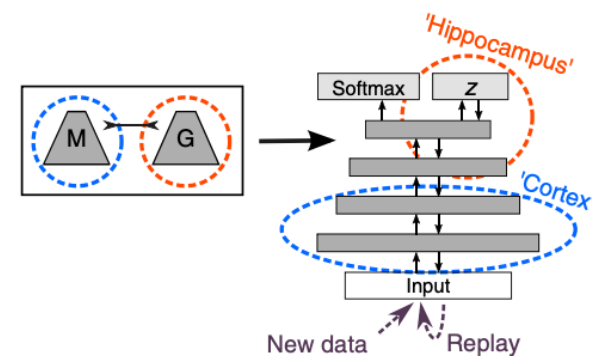


# Replay

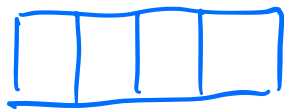
- Store raw data, representations, or train a generative model *generative replay*.



## Feedback Connections



reservoir sampling.



beginning high prob. string samples.

Replay

- Store raw data, representations, or train a generative model
- Coreset selection

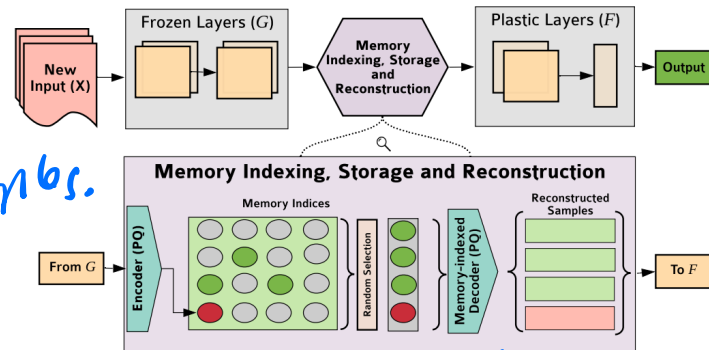
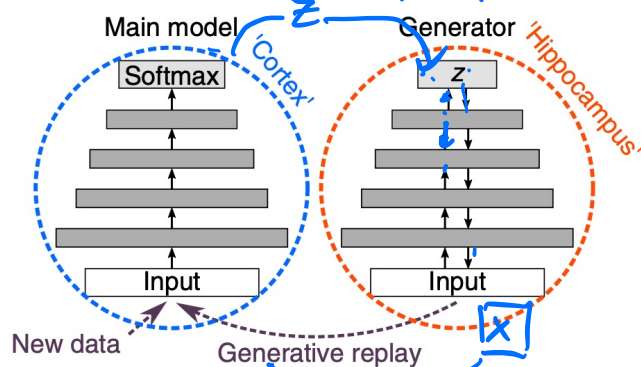
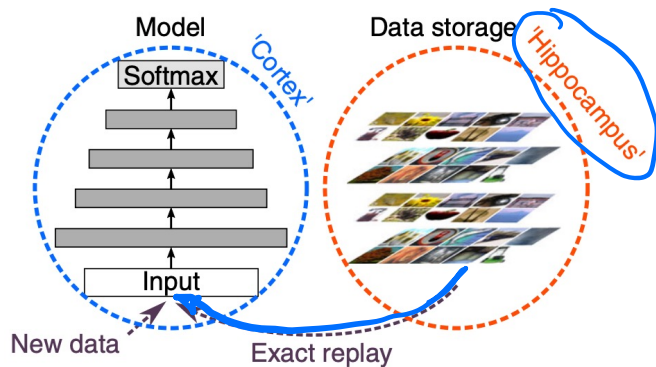
$$\mu \leftarrow \frac{1}{n} \sum_{x \in X} \varphi(x)$$

$$p_k \leftarrow \underset{x \in X}{\operatorname{argmin}} \left\| \mu - \frac{1}{k} \left[ \varphi(x) + \sum_{j=1}^{k-1} \varphi(p_j) \right] \right\|.$$

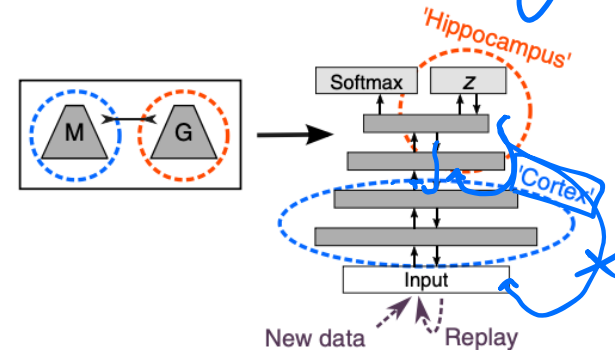
average rep

classification self-sup.

closer to avg representation.



Feedback Connections



# Knowledge Distillation

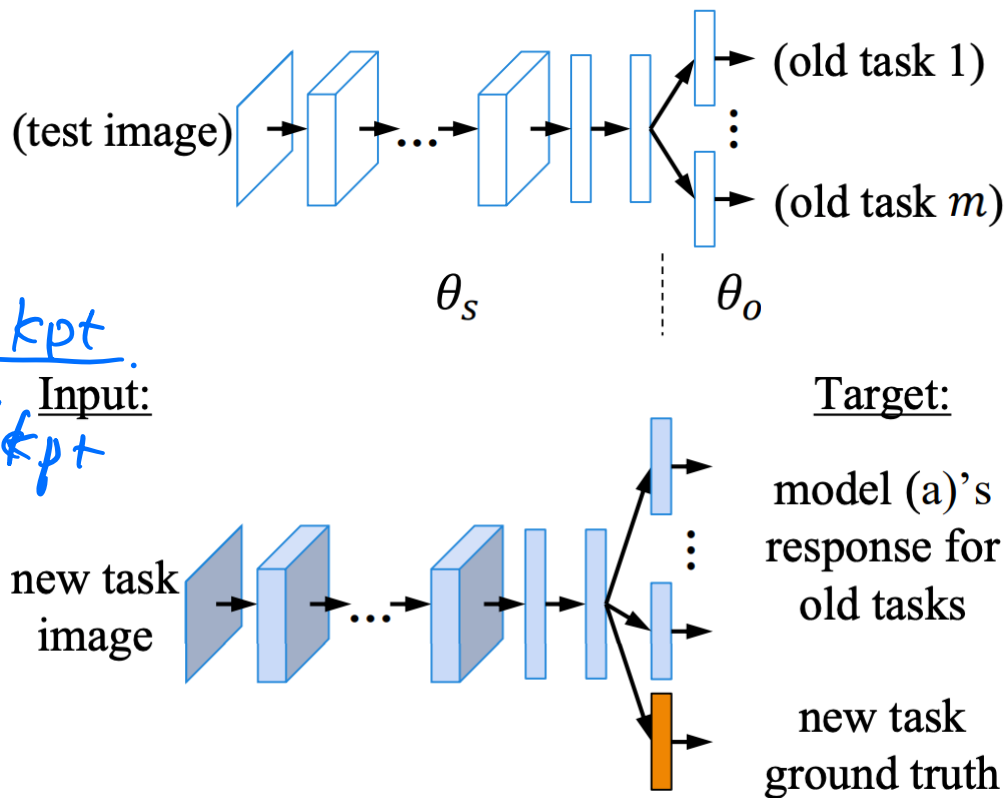
- Instead of saving the data points, we can also save the previous model checkpoint.

$$y_o = f(x_n; \theta_o) \quad \text{old ckpt.}$$

$$\hat{y}_o = f(x_n; \theta_n) \quad \text{new ckpt.}$$

$$\mathcal{L}(y_o, \hat{y}_o)$$

$$y_n = f(x_n, \theta_n) \quad \text{new data}$$



# Knowledge Distillation

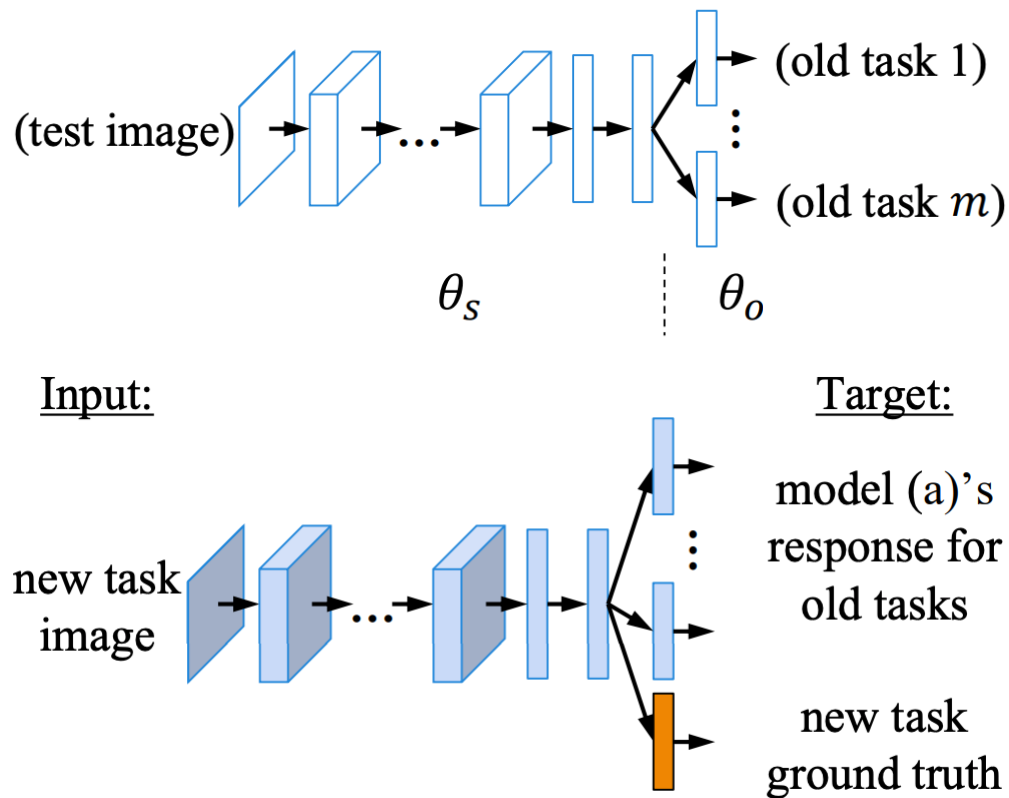
- Instead of saving the data points, we can also save the previous model checkpoint.

$$y_o = f(x_n; \theta_o).$$

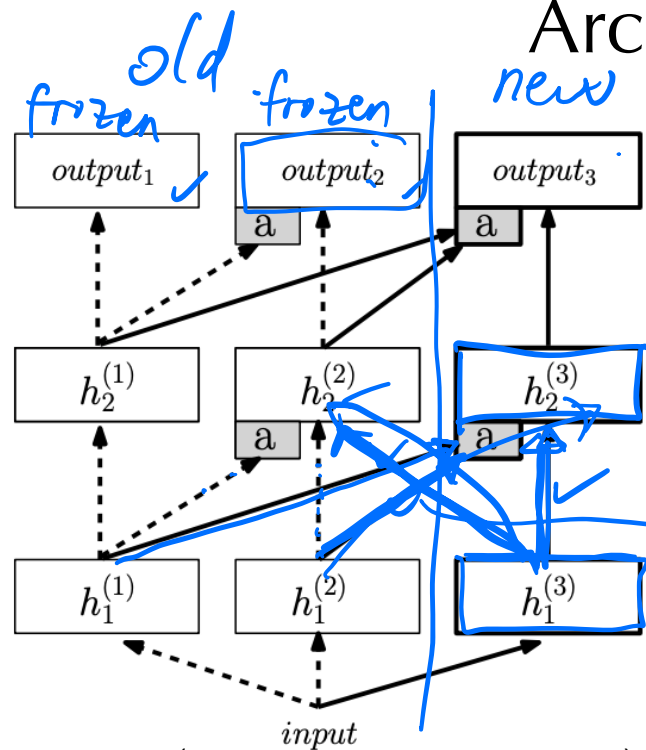
$$\hat{y}_o = f(x_n; \theta_n).$$

$$\mathcal{L}(y_o, \hat{y}_o)$$

- Use new data points and old weights to “distill”



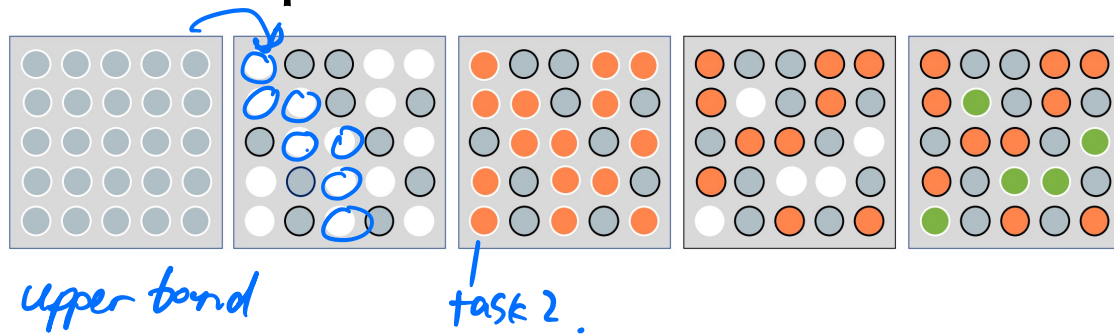
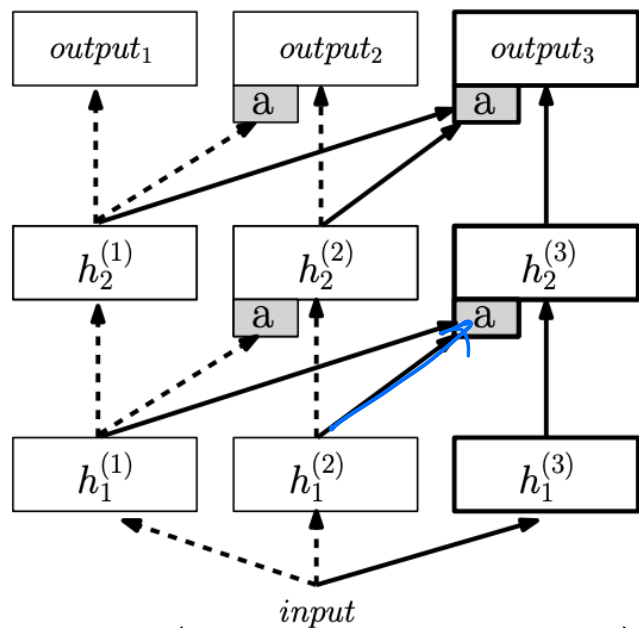
# Architecture Expansion



leverage knowledge  
from previous tasks.

$$h_i^{(k)} = f \left( W_i^{(k)} h_{i-1}^{(k)} + \sum_{j < k} U_i^{(k:j)} h_{i-1}^{(j)} \right).$$

# Architecture Expansion



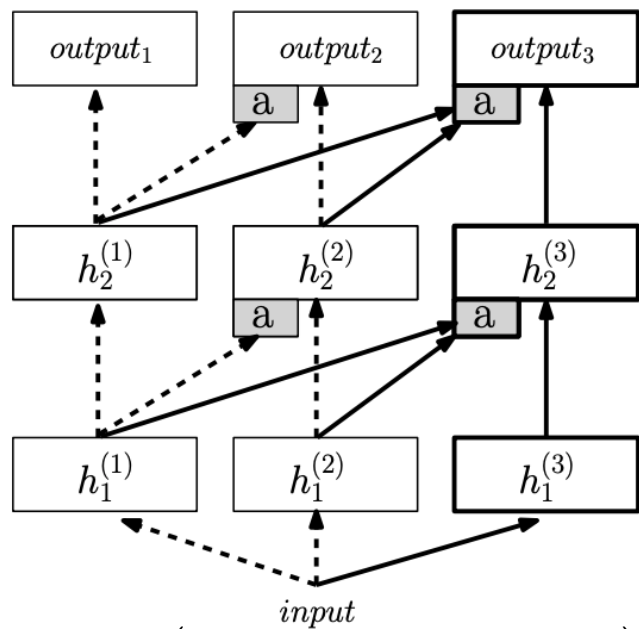
$$h_i^{(k)} = f \left( W_i^{(k)} h_{i-1}^{(k)} + \sum_{j < k} U_i^{(k:j)} h_{i-1}^{(j)} \right).$$

Rusu et al. *Progressive Neural Networks*. NIPS 2016 Deep Learning Symposium.

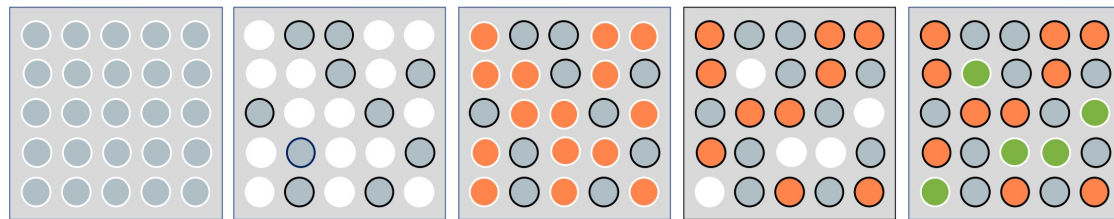
PackNet: Adding Multiple Tasks to a Single Network by Iterative Pruning. CVPR 2018.

Yoon et al. *Lifelong Learning with Dynamically Expandable Networks*. ICLR 2018.

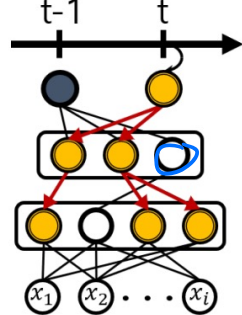
# Architecture Expansion



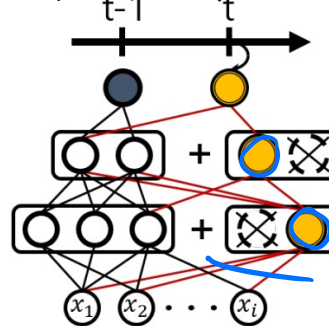
$$h_i^{(k)} = f \left( W_i^{(k)} h_{i-1}^{(k)} + \sum_{j < k} U_i^{(k:j)} h_{i-1}^{(j)} \right).$$



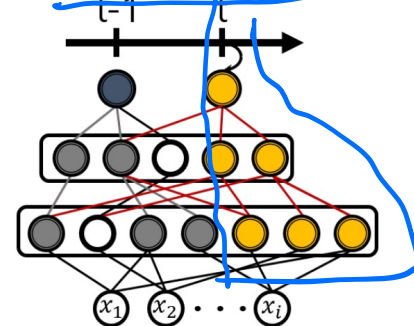
Selective Training



Dynamic Expansion



Network Split



Selective Training

$$\mathcal{L}(\mathbf{W}_L^t, \mathbf{W}_{1:L-1}^{t-1}, \mathcal{D}_t) + \mu \|\mathbf{W}_L^t\|_1$$

Dynamic Expansion

$$\mathcal{L}(\mathbf{W}_L^N, \mathbf{W}_L^{t-1}, \mathcal{D}_t) + \mu \|\mathbf{W}_L^N\|_1 + \gamma \sum_g \|\mathbf{W}_{L,g}^N\|_2$$

Group sparsity

Rusu et al. Progressive Neural Networks. NIPS 2016 Deep Learning Symposium.

PackNet: Adding Multiple Tasks to a Single Network by Iterative Pruning. CVPR 2018.

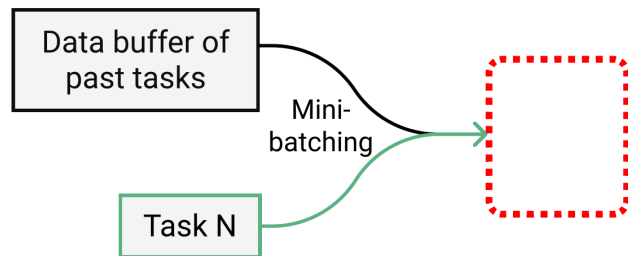
Yoon et al. Lifelong Learning with Dynamically Expandable Networks. ICLR 2018.




# Adapting Pretrained Models

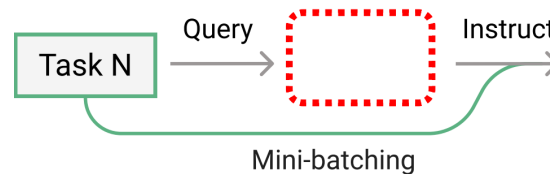
- Pretrained models have general knowledge that can be adapted to a continual stream of tasks.

Rehearsal-based methods:  
Fine-tuning



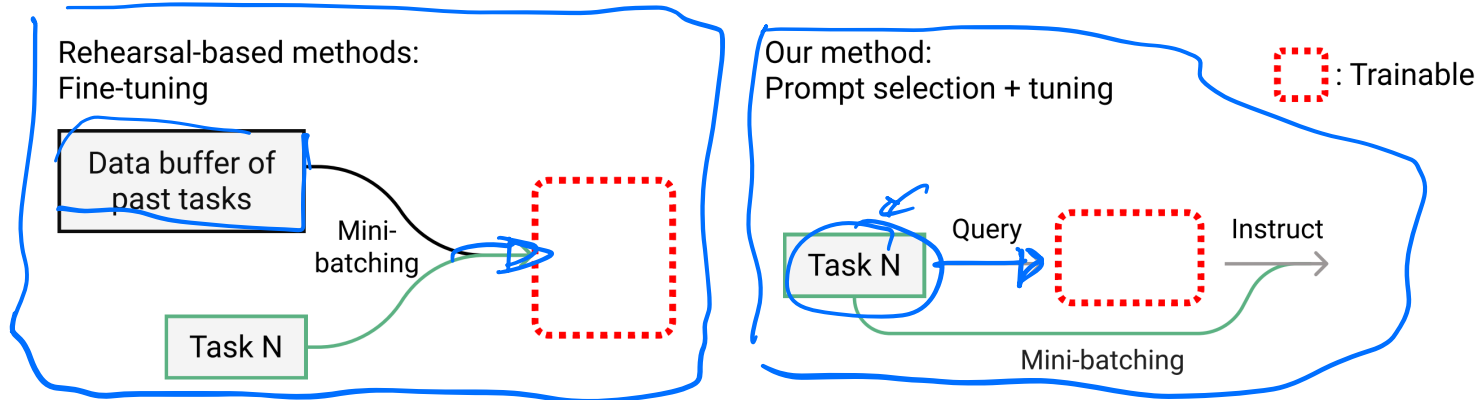
Our method:  
Prompt selection + tuning

 : Trainable



# Adapting Pretrained Models

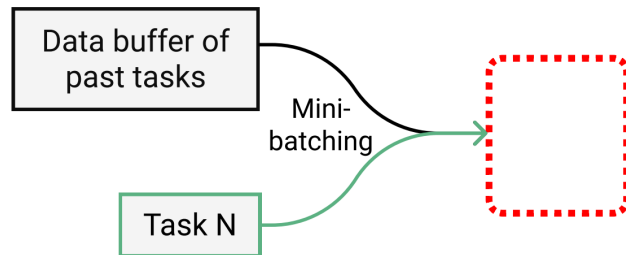
- Pretrained models have general knowledge that can be adapted to a continual stream of tasks.
- Learn adaptation parameters for each task and store these as “task embeddings.”




# Adapting Pretrained Models

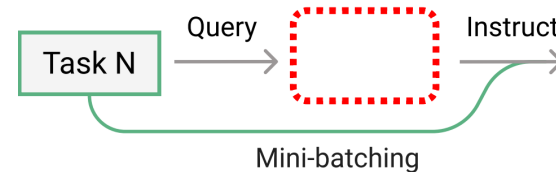
- Pretrained models have general knowledge that can be adapted to a continual stream of tasks.
- Learn adaptation parameters for each task and store these as “task embeddings.”
- Main model is frozen.

Rehearsal-based methods:  
Fine-tuning



Our method:  
Prompt selection + tuning

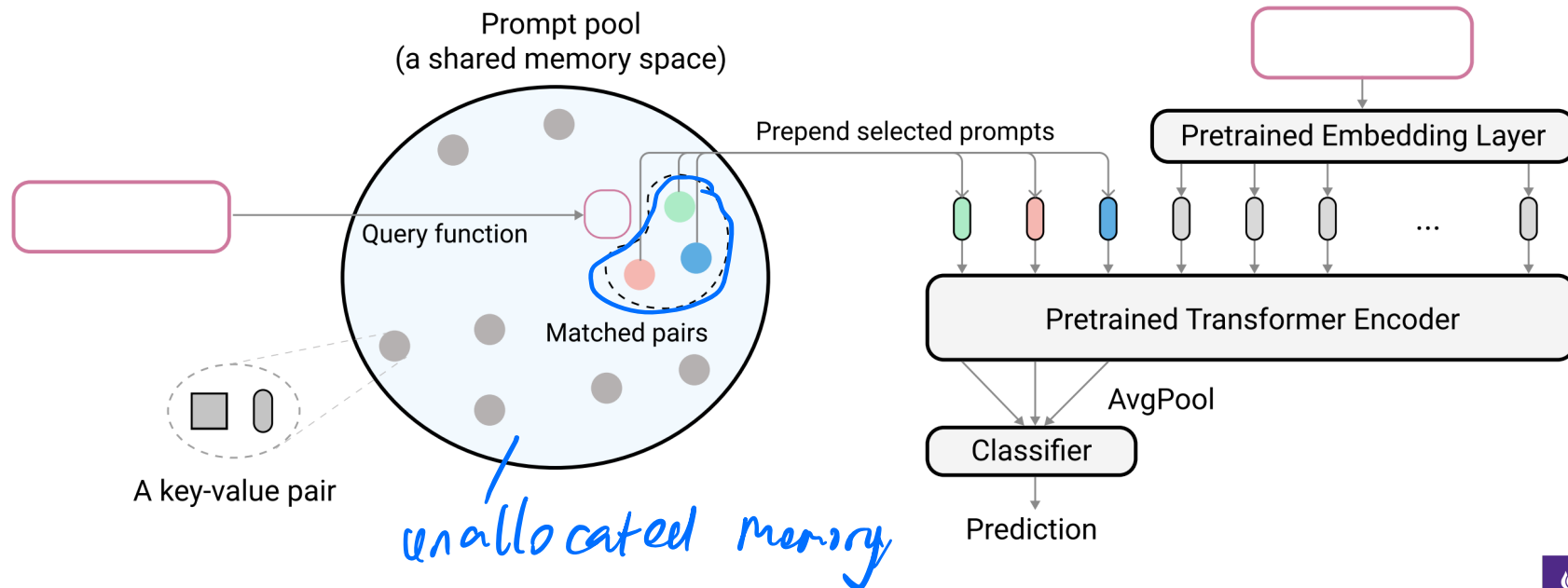
 : Trainable



# Learning to Prompt

Prompt pool (slot memory)

$$\{(k_1, p_1), \dots, (k_M, p_M)\}$$

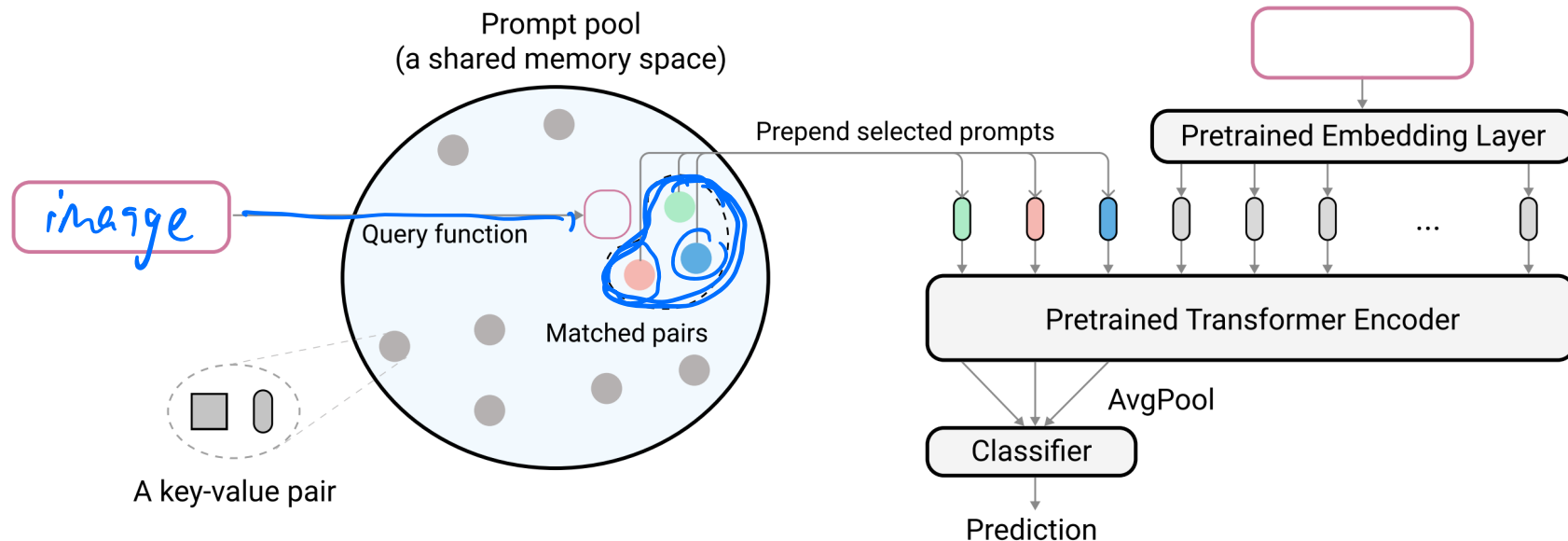


# Learning to Prompt

Prompt pool (slot memory)

$\{(k_1, p_1), \dots, (k_M, p_M)\}$

$$K_s = \underline{\text{argmin}}_S \sum_{i \in S} \underline{\gamma}(q(x), \underline{k_i})$$



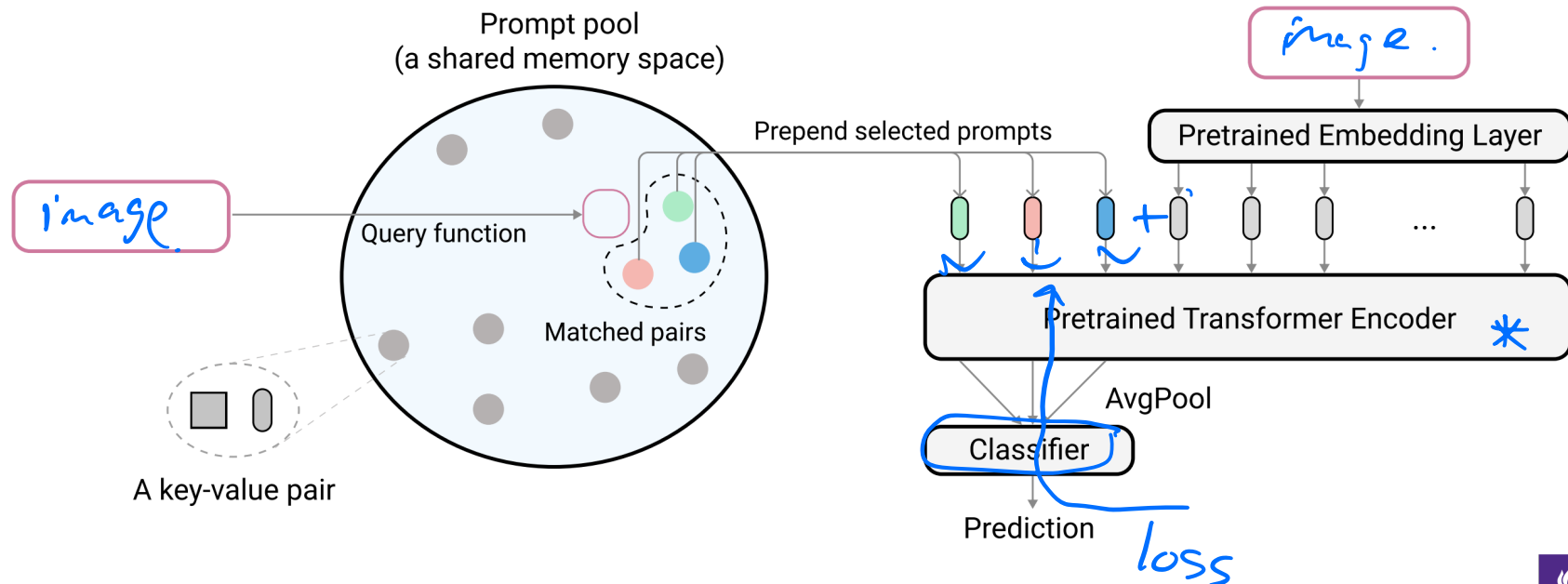
# Learning to Prompt

Prompt pool (slot memory)

$\{(k_1, p_1), \dots, (k_M, p_M)\}$

$$K_s = \operatorname{argmin}_{\mathcal{S}} \sum_{i \in \mathcal{S}} \gamma(q(x), k_i)$$

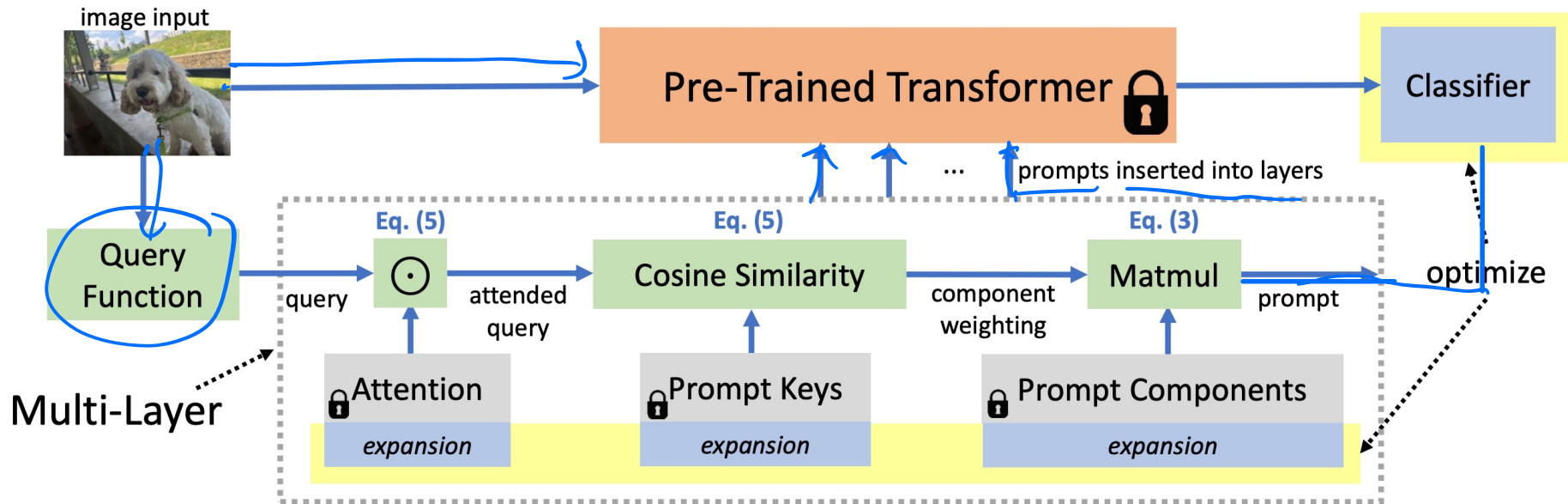
$$\min_{P, K, \phi} \mathcal{L}(g_{\phi}(x), y) + \lambda \sum_{i \in K_s} \gamma(q(x), k_i)$$



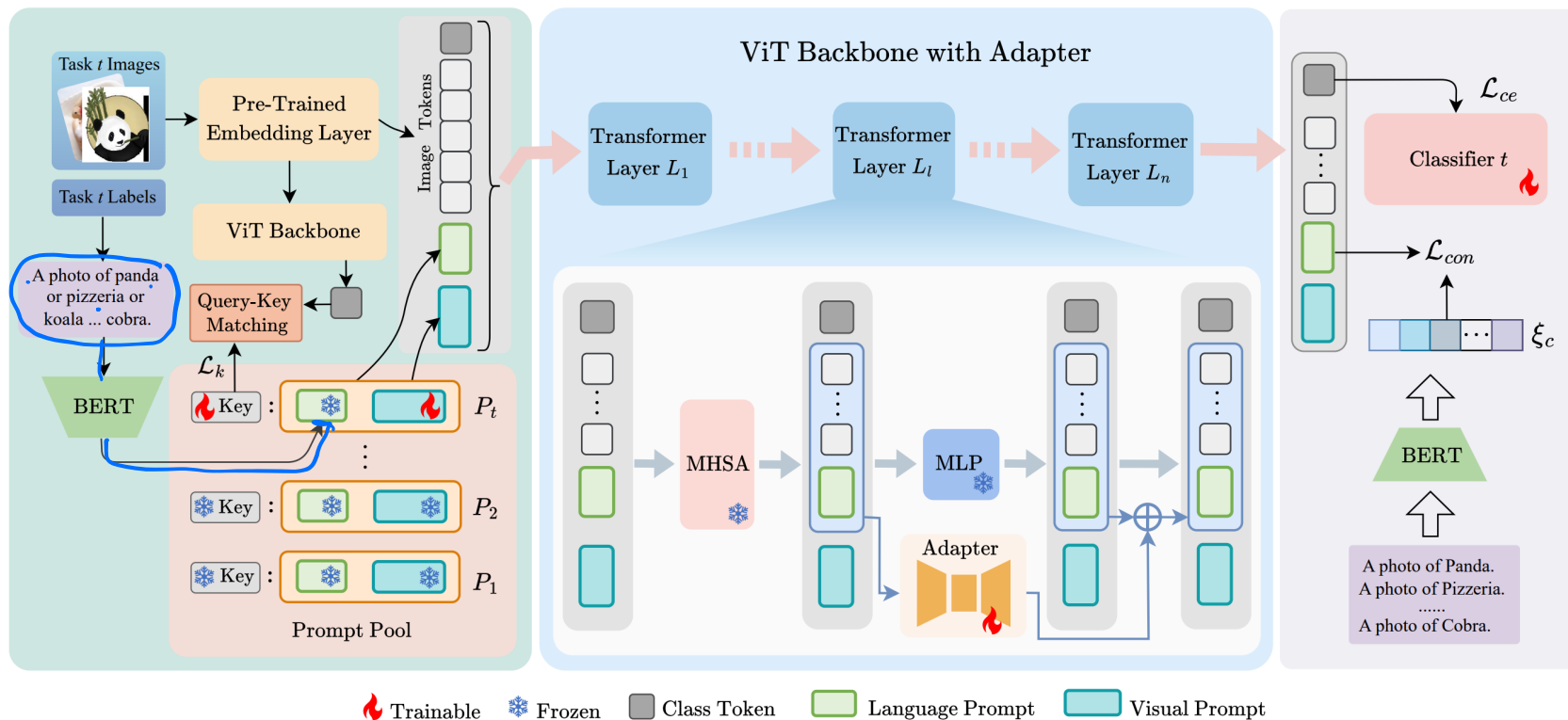
# Learned Prompt Query

- The query function can be end-to-end learned.

$$\alpha = \{\gamma(q(x) \odot A_1, K_1), \dots, \gamma(q(x) \odot A_M, K_M)\}$$



# Multimodal Semantic Prompts





# Continual Learning and Memory

forgetting.

- Fragility of feedforward gradient descent of the entire networks

# Continual Learning and Memory

- Fragility of feedforward gradient descent of the entire networks
- If we have representations ready, continual learning is just memorizing a sequence of new tasks.

# Continual Learning and Memory

- Fragility of feedforward gradient descent of the entire networks
- If we have representations ready, continual learning is just memorizing a sequence of new tasks.
- In prompting approaches:
  - prompt pool = memory
  - pretrained network = representations

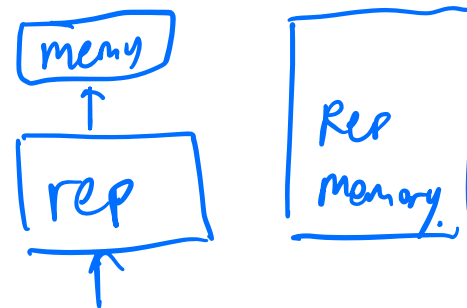
*memory.*

# Continual Learning and Memory

- Fragility of feedforward gradient descent of the entire networks
- If we have representations ready, continual learning is just memorizing a sequence of new tasks.
- In prompting approaches:
  - prompt pool = memory
  - pretrained network = representations
- But what if representations also need to be built sequentially?

# Continual Learning and Memory

- Fragility of feedforward gradient descent of the entire networks
- If we have representations ready, continual learning is just memorizing a sequence of new tasks.
- In prompting approaches:
  - prompt pool = memory
  - pretrained network = representations
- But what if representations also need to be built sequentially?
- It's also plausible that representations are just “deeper memory.”



# Associative Memory

- Memory aims to store content for easy retrieval



# Associative Memory

- Memory aims to store content for easy retrieval
  - Associative memories (Hopfield Networks) can be viewed as energy-based models



# Associative Memory

- Memory aims to store content for easy retrieval
  - Associative memories (Hopfield Networks) can be viewed as energy-based models

$$E = -\frac{1}{2} \sum_{i,j=1}^N s_i W_{ij} s_j$$

$s = \{ \dots \}$  Stored patterns. retrieval.

$$W_{ij} = \sum_{k=1}^K \xi_i^k \xi_j^k$$

"Superposition" of k slots  
Hebbian learning

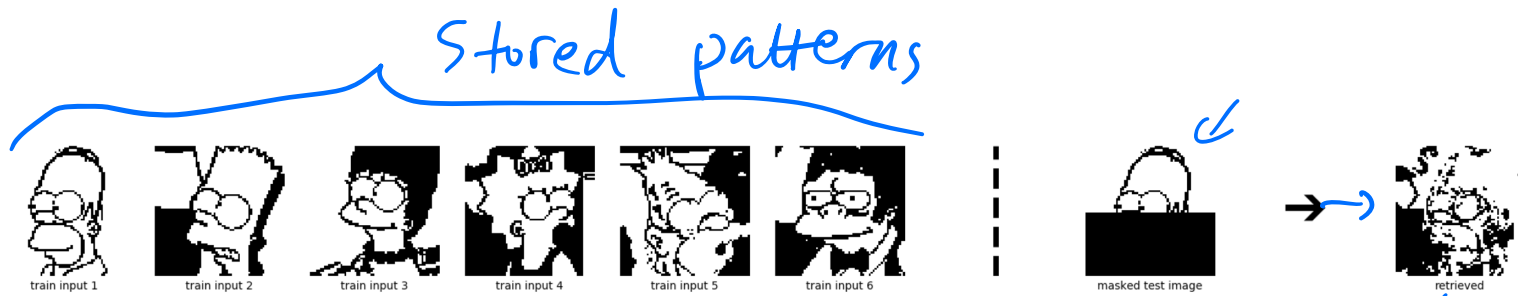
$\xi$  Stored pattern.





# Associative Memory

- Memory aims to store content for easy retrieval
    - Associative memories (Hopfield Networks) can be viewed as energy-based models
- $$E = -\frac{1}{2} \sum_{i,j=1}^N s_i W_{ij} s_j, \quad W_{ij} = \sum_{k=1}^K \xi_i^k \xi_j^k.$$
- "Superposition" of k slots  
Hebbian learning
- When presented with a new pattern the network should respond with a stored memory which most closely resembles the input.



# Associative Memory

- Memory aims to store content for easy retrieval

- Associative memories (Hopfield Networks) can be viewed as energy-based models

$$E = -\frac{1}{2} \sum_{i,j=1}^N s_i W_{ij} s_j, \quad W_{ij} = \sum_{k=1}^K \xi_i^k \xi_j^k.$$

"Superposition" of  $k$  slots  
Hebbian learning

*binary.*

- When presented with a new pattern the network should respond with a stored memory which most closely resembles the input.

- Retrieval:  $s_i = \text{sign}(\sum_j W_{ij} s_j)$

Storage:  $C \approx \frac{d}{2 \log(d)} = 0.14d.$

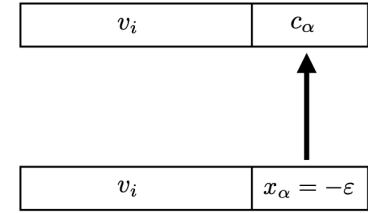
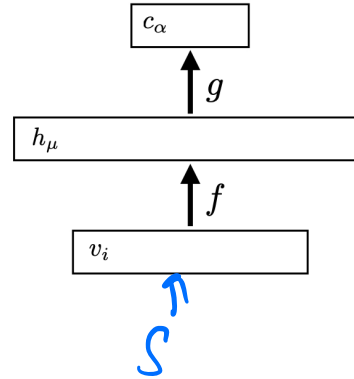
*100.*



# Memory Duality

- Duality with a feedforward network.

$$\underline{E} = - \sum_k \underline{F}(\sum_i \xi_i^k s_i)$$

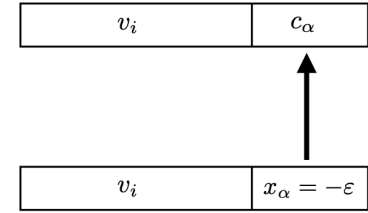
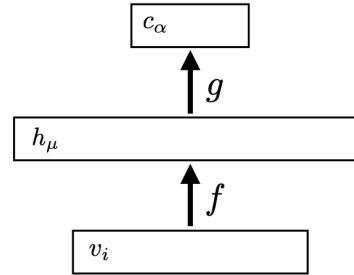


# Memory Duality

- Duality with a feedforward network.

$$E = - \sum_k F(\sum_i \xi_i^k s_i)$$

- Non-linearity allows us to store more patterns.

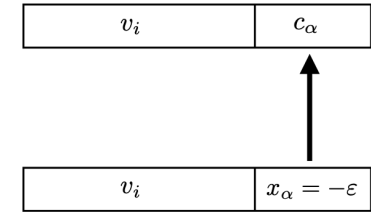
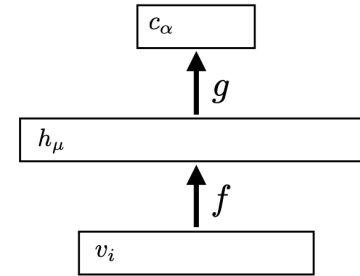


# Memory Duality

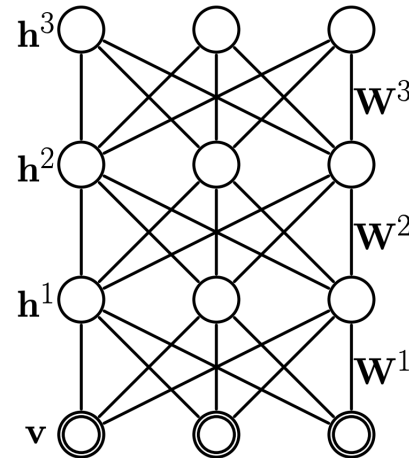
- Duality with a feedforward network.

$$E = - \sum_k F(\sum_i \xi_i^k s_i)$$

- Non-linearity allows us to store more patterns.
- Deep Boltzmann machines



**Deep Boltzmann Machine**

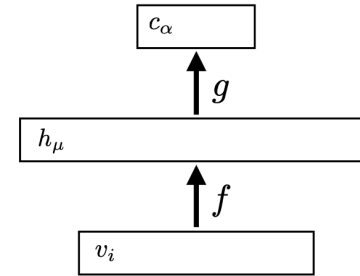


# Memory Duality

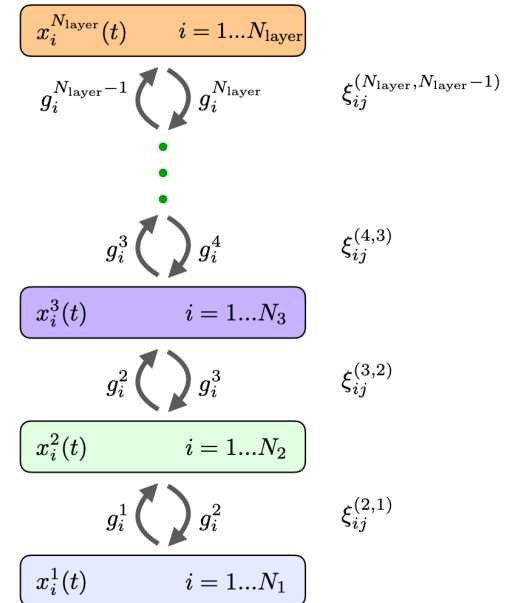
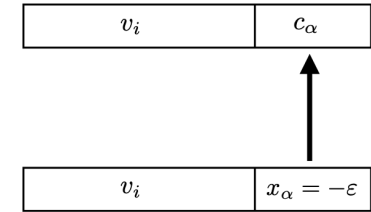
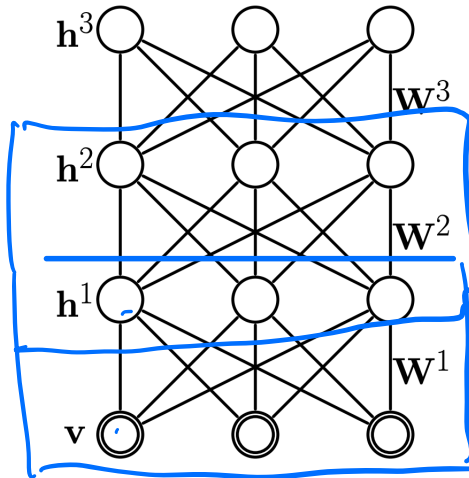
- Duality with a feedforward network.

$$E = - \sum_k F(\sum_i \xi_i^k s_i)$$

- Non-linearity allows us to store more patterns.
- Deep Boltzmann machines
- Hierarchical associative memory



**Deep Boltzmann Machine**



# Relation to Transformers

- General form:

$$\underbrace{E} = - \underbrace{\sum_k}_{\text{Stored patterns}} \underbrace{F}_{\text{retrieved vector}} \left( \underbrace{\sum_i}_{\text{retrieved vector}} \xi_i^k s_i \right).$$

Diagram illustrating the general form of the energy function  $E$  in Hopfield Networks, with handwritten annotations:

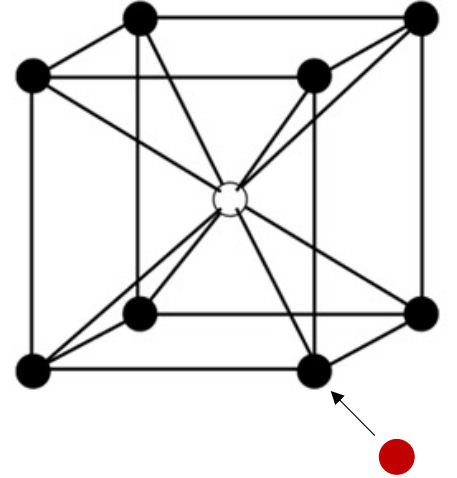
- $E$  is underlined with a blue bracket.
- The summation over  $k$  is underlined with a blue bracket and labeled "Stored patterns".
- The function  $F$  is circled in blue and labeled "retrieved vector".
- The summation over  $i$  is underlined with a blue bracket and labeled "retrieved vector".
- The term  $\xi_i^k$  is underlined with a blue bracket.
- The term  $s_i$  is underlined with a blue bracket.

# Relation to Transformers

- General form:

$$E = - \sum_k F\left(\sum_i \xi_i^k s_i\right).$$

- When  $F(z) = z^2$  it gives the classic HN.



$$\nabla_{s_i} E = - \sum_j W_{ij} s_j$$

$s_i$   $\leftarrow \text{sign}\left(\sum_j W_{ij} s_j\right)$

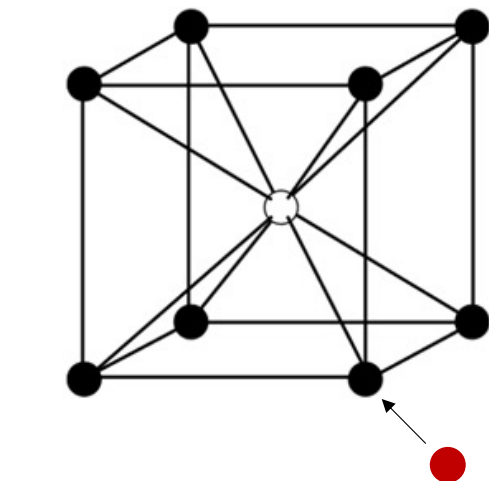


# Relation to Transformers

- General form:

$$E = - \sum_k F\left(\sum_i \xi_i^k s_i\right).$$

- When  $F(z) = z^2$  it gives the classic HN.
- Transformer-like attention operation:



$$\nabla_{s_i} E = - \sum_j W_{ij} s_j$$
$$s_i \leftarrow \text{sign}\left(\sum_j W_{ij} s_j\right)$$

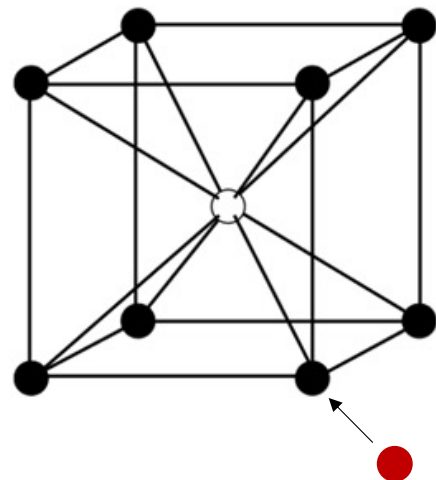
# Relation to Transformers

- General form:

$$E = - \sum_k F\left(\sum_i \xi_i^k s_i\right).$$

- When  $F(z) = z^2$  it gives the classic HN.
- Transformer-like attention operation:

$$\mathbf{Z} \leftarrow \text{softmax}\left(\beta \underbrace{\mathbf{X} \mathbf{W}_q}_{\text{blue}} \underbrace{\mathbf{W}_k^\top}_{\text{blue}} \underbrace{(\mathbf{Y})^\top}_{\text{blue}}\right) \underbrace{\mathbf{Y}_i \mathbf{W}_v}_{\text{blue}}.$$



$$\nabla_{s_i} E = - \sum_j W_{ij} s_j$$

$$s_i \leftarrow \text{sign}\left(\sum_j W_{ij} s_j\right)$$

# Relation to Transformers

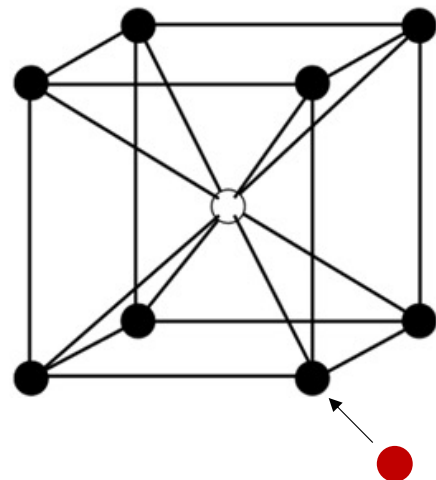
- General form:

$$E = - \sum_k F\left(\sum_i \xi_i^k s_i\right).$$

- When  $F(z) = z^2$  it gives the classic HN.
- Transformer-like attention operation:

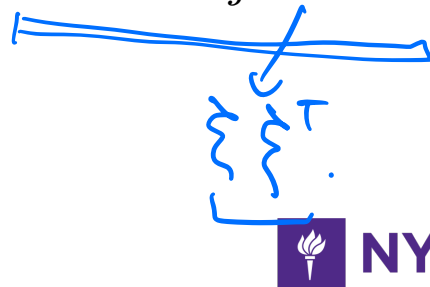
$$\mathbf{Z} \leftarrow \text{softmax}(\beta \mathbf{X} \mathbf{W}_q \mathbf{W}_k^\top \mathbf{Y}^\top) \mathbf{Y}_i \mathbf{W}_v.$$

$$\mathbf{S} \leftarrow \text{softmax}(\beta \mathbf{S} \mathbf{\Xi}^\top) \mathbf{\Xi}. \quad \rightarrow \nabla E.$$



$$\nabla_{s_i} E = - \sum_j W_{ij} s_j$$

$$s_i \leftarrow \text{sign}\left(\sum_j W_{ij} s_j\right)$$

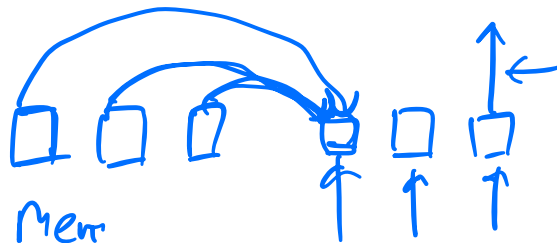


# Relation to Transformers

- General form:

$$E = - \sum_k F\left(\sum_i \xi_i^k s_i\right).$$

*Mem*

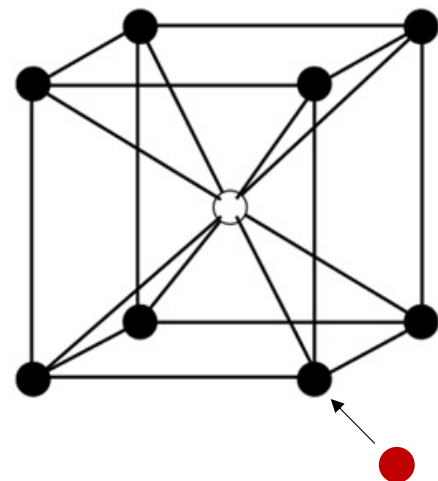


- When  $F(z) = z^2$  it gives the classic HN.
- Transformer-like attention operation:

$$\mathbf{Z} \leftarrow \text{softmax}(\beta \mathbf{X} \mathbf{W}_q \mathbf{W}_k^\top \mathbf{Y}^\top) \mathbf{Y}_i \mathbf{W}_v.$$

$$\mathbf{S} \leftarrow \text{softmax}(\beta \mathbf{S} \mathbf{\Xi}^\top) \mathbf{\Xi}. \quad \text{1 step.}$$

$$E = -\text{logsumexp}(\beta, \mathbf{\Xi}^\top \mathbf{s}) + \frac{1}{2} \mathbf{s}^\top \mathbf{s} + \beta^{-1} \log N + \frac{1}{2} M^2.$$

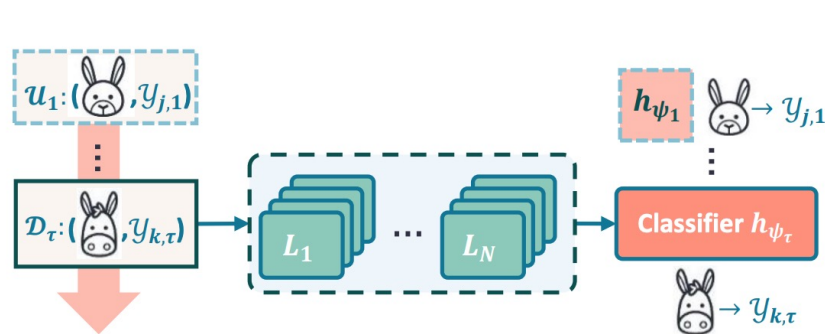


$$\nabla_{s_i} E = - \sum_j W_{ij} s_j$$

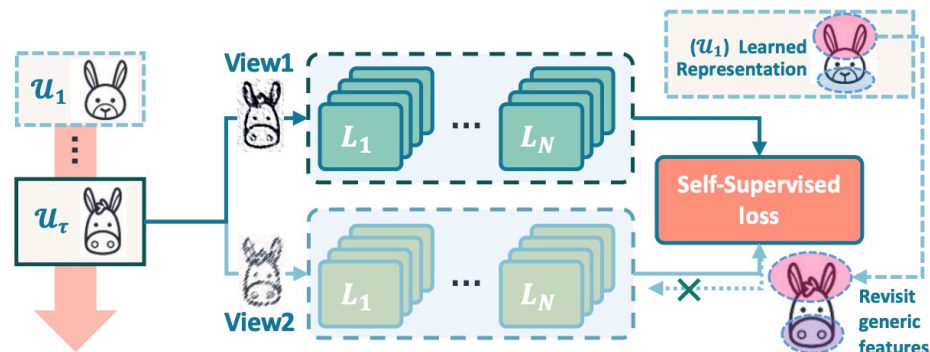
$$s_i \leftarrow \text{sign}\left(\sum_j W_{ij} s_j\right)$$

# Continual Self-Supervised Learning

- Learning from a stream of unlabeled inputs.



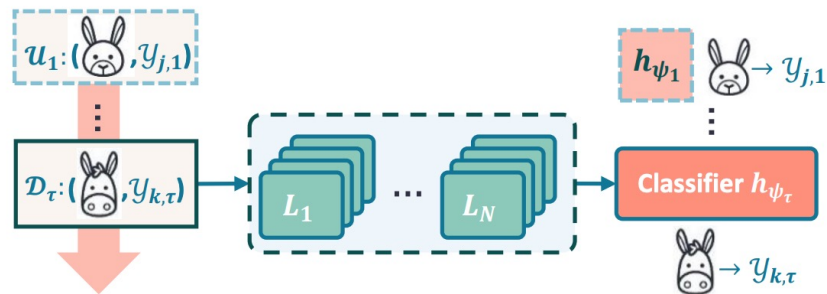
SUPERVISED CONTINUAL LEARNING (SCL)



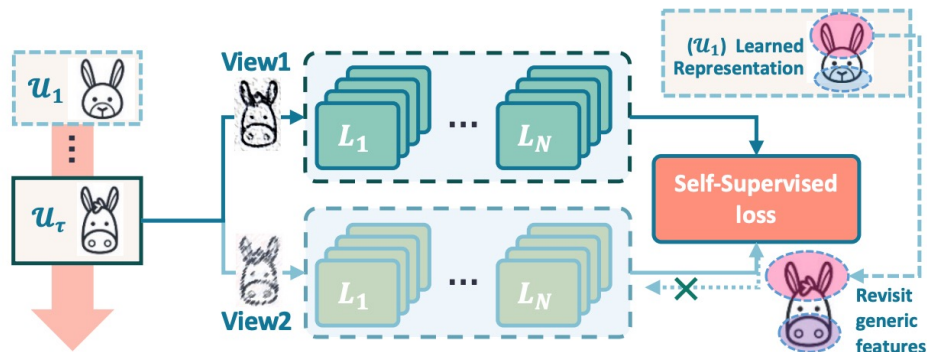
UNSUPERVISED CONTINUAL LEARNING (UCL)

# Continual Self-Supervised Learning

- Learning from a stream of unlabeled inputs.
- Bring SSL to the dynamic world.



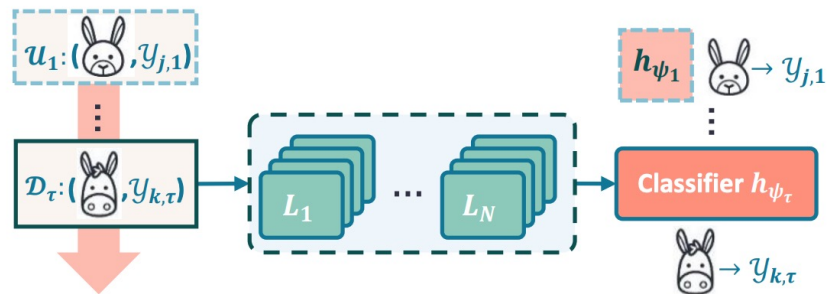
SUPERVISED CONTINUAL LEARNING (SCL)



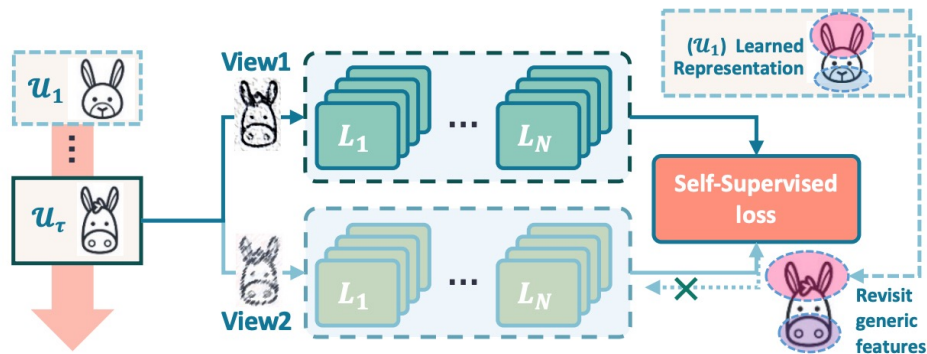
UNSUPERVISED CONTINUAL LEARNING (UCL)

# Continual Self-Supervised Learning

- Learning from a stream of unlabeled inputs.
- Bring SSL to the dynamic world.
- SSL can still suffer from distributional shifts.



SUPERVISED CONTINUAL LEARNING (SCL)

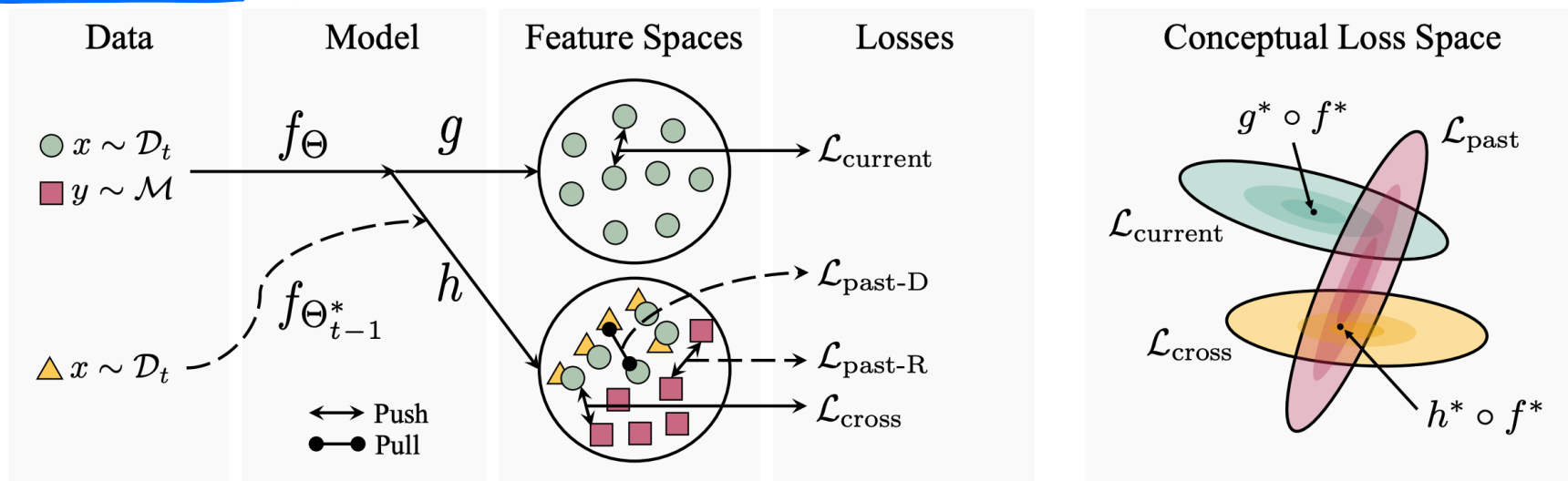


UNSUPERVISED CONTINUAL LEARNING (UCL)

# Continual Self-Supervised Learning

- Integrating learning objectives of both past and present.

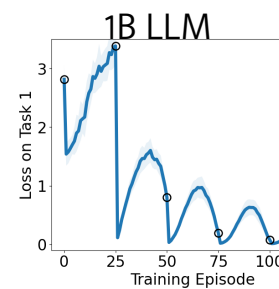
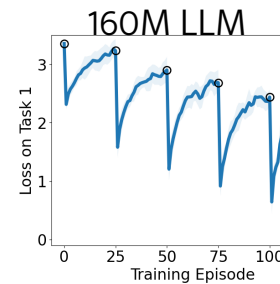
$$\underbrace{\mathcal{L}(X; g \circ f_t)}_{\text{Current}} + \underbrace{\mathcal{L}(X, X; h \circ f_{t-1}, h \circ f_t)}_{\text{Past (Distillation)}} + \underbrace{\mathcal{L}(Y; h \circ f_t)}_{\text{Past (Replay)}} + \underbrace{\mathcal{L}(X, Y; h \circ f_t)}_{\text{Cross-Consolidation}}$$





# Outlooks

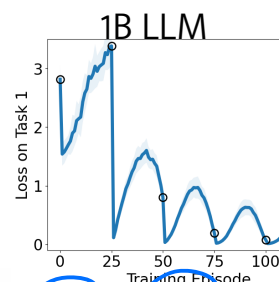
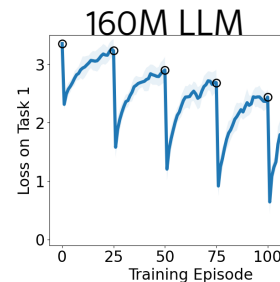
- Understand continual learning at scale



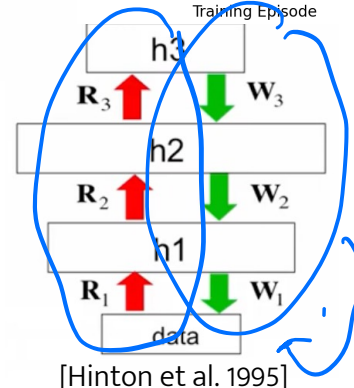
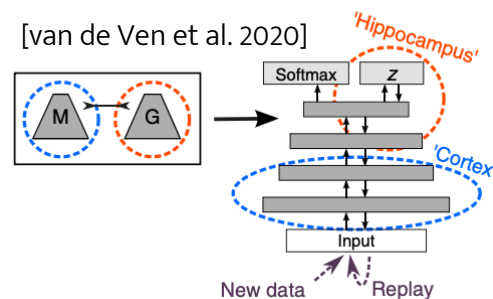
[Mayo et al. 2023]

# Outlooks

- Understand continual learning at scale
- Unified learning architecture, objective and replay, role of sleep



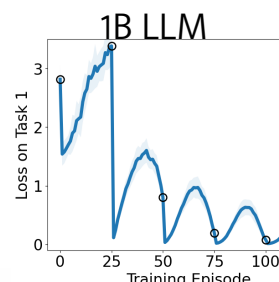
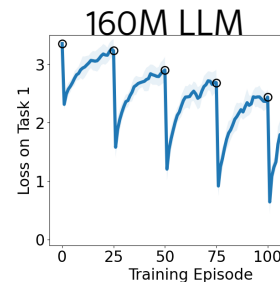
[van de Ven et al. 2020]



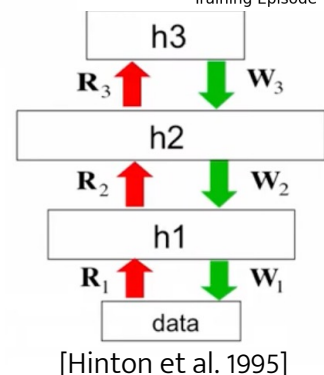
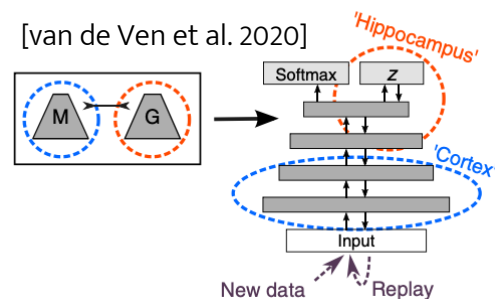
[Mayo et al. 2023]

# Outlooks

- Understand continual learning at scale
- Unified learning architecture, objective and replay, role of sleep
- Continual learning with real world structure



[van de Ven et al. 2020]



Typical setting in continual, few-shot, transfer, and meta-learning



Typical setting of human learning [Mayo et al. 2023]



Tasks are interspersed and recur  
No opportunity to master one before confronted with another

# Summary: Continual Learning

- Regularization, Distillation, Architecture Expansion/Isolation

# Summary: Continual Learning

- Regularization, Distillation, Architecture Expansion/Isolation
- Frozen representation: prompt learning

# Summary: Continual Learning

- Regularization, Distillation, Architecture Expansion/Isolation
- Frozen representation: prompt learning
- Integration of memory and representations

# Summary: Continual Learning

- Regularization, Distillation, Architecture Expansion/Isolation
- Frozen representation: prompt learning
- Integration of memory and representations
- Combination with self-supervised learning

# Summary: Continual Learning

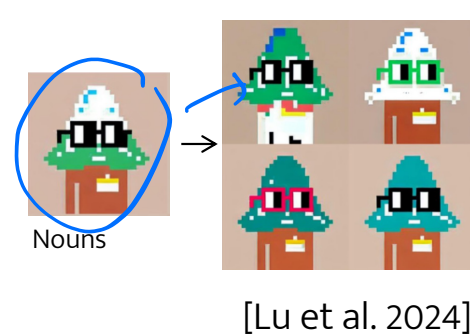
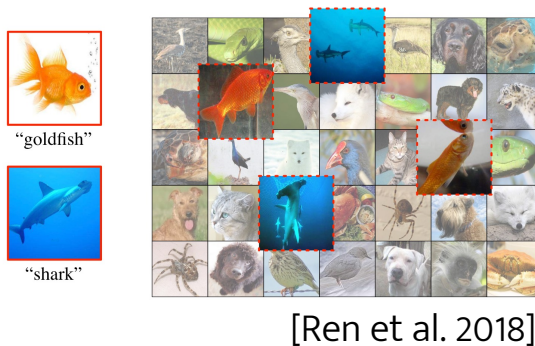
- Regularization, Distillation, Architecture Expansion/Isolation
- Frozen representation: prompt learning
- Integration of memory and representations
- Combination with self-supervised learning
- Exploration of multimodal continual learning from embodied environments



# Few-Shot Learning and Meta-Learning

# Few-Shot Learning (FSL)

- Humans can quickly learn new concepts with a few examples.
- Learning in embodied agents also needs to be adaptive and swift.
- Examples: Recognize new objects, perform new skills, map new areas, etc.

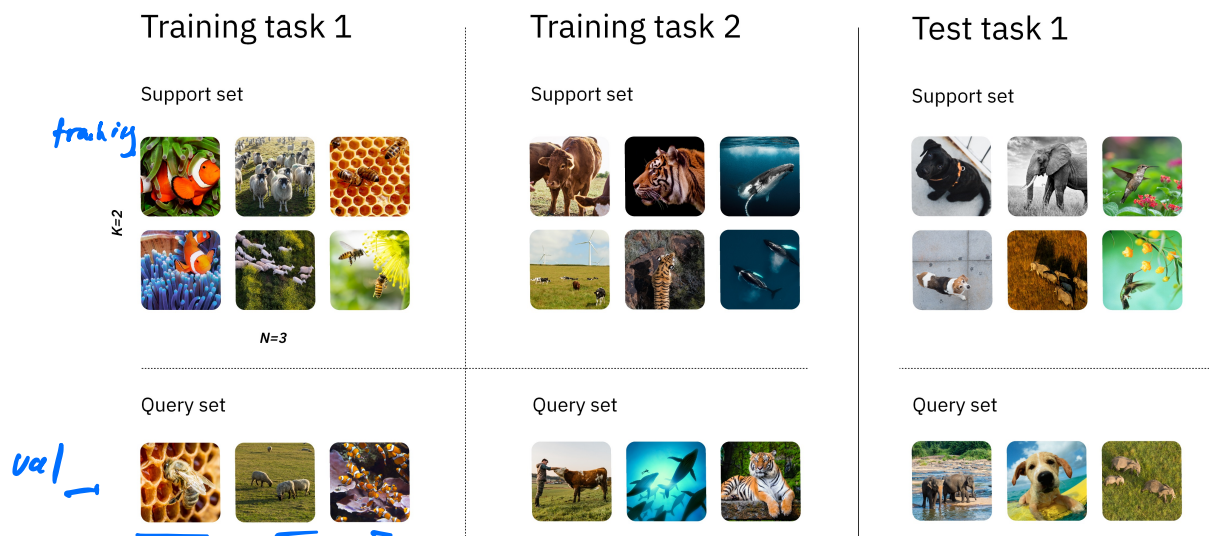


# FSL: General Setup

- Quickly learn a task (learning episode) with very few number of training examples and get evaluated on a set of test examples.

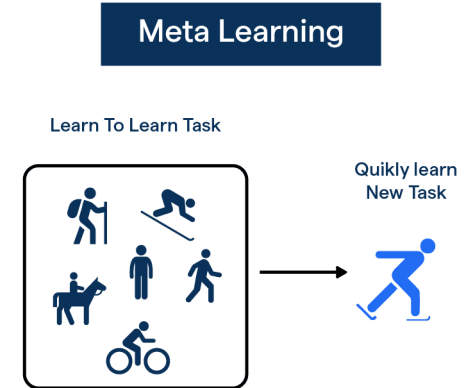
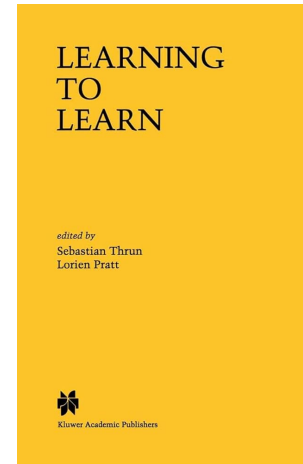
# FSL: General Setup

- Quickly learn a task (learning episode) with very few number of training examples and get evaluated on a set of test examples.
- During training, going through many episodes of the same structure.



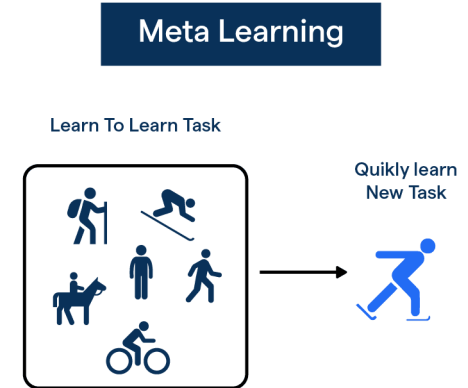
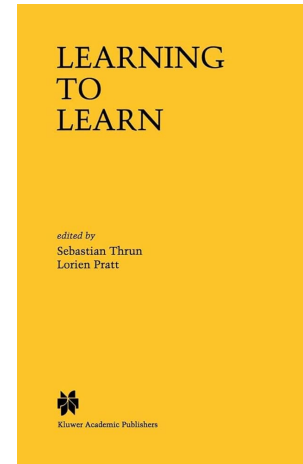
# Learning to Learn

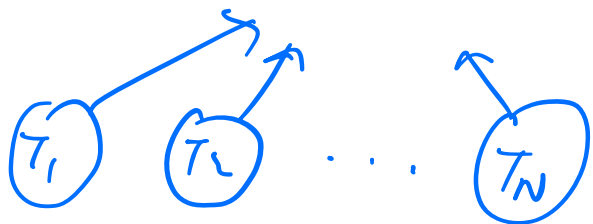
- Conceptually, we'd like to generalize new learning experiences.



# Learning to Learn

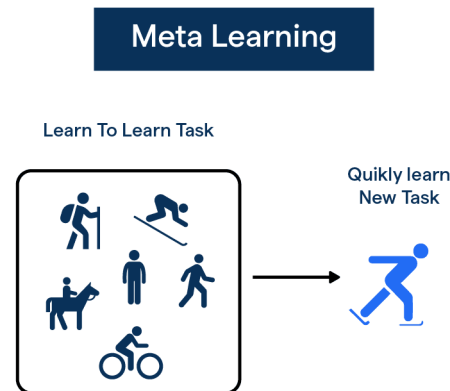
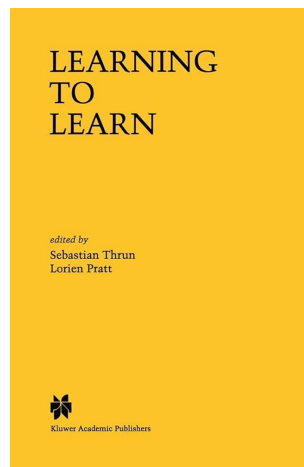
- Conceptually, we'd like to generalize new learning experiences.
- 1 learning experience = 1 training example





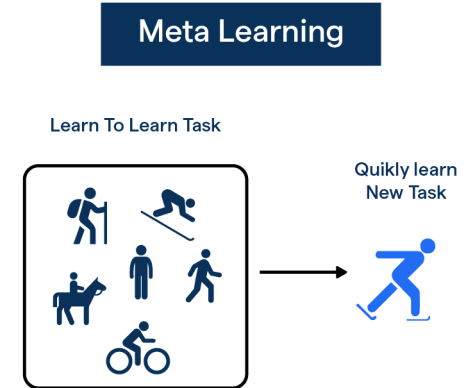
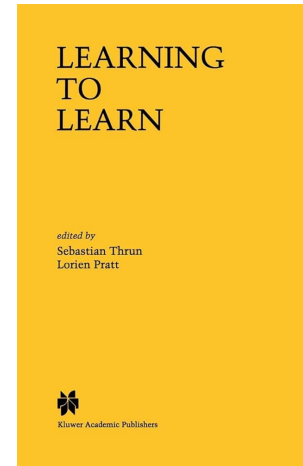
# Learning to Learn

- Conceptually, we'd like to generalize new learning experiences.
- 1 learning experience = 1 training example
- Related to multi-task learning



# Learning to Learn

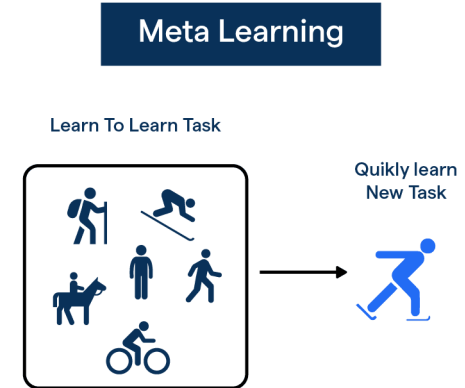
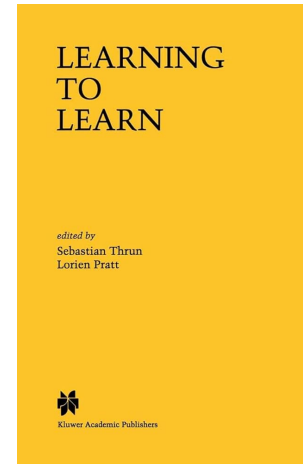
- Conceptually, we'd like to generalize new learning experiences.
- 1 learning experience = 1 training example
- Related to multi-task learning
- What can be meta-learned?
  - Optimizer
  - Initialization
  - Architecture
  - Representations
  - Abstraction of tasks





# Learning to Learn

- Conceptually, we'd like to generalize new learning experiences.
- 1 learning experience = 1 training example
- Related to multi-task learning
- What can be meta-learned?
  - Optimizer
  - Initialization
  - Architecture
  - Representations
  - Abstraction of tasks
- Discover learning algorithms that support FSL.



# Meta-Optimization (Bi-level Optimization)

- One can formulate meta-learning as a meta-optimization problem.

$$\min_{\lambda} \mathbb{E}_{D \sim \mathcal{D}} \left( \min_{\theta} \mathbb{E}_{x \sim D} \mathcal{L}(x; \theta, \lambda) \right).$$

hyperparameter

# Meta-Optimization (Bi-level Optimization)

- One can formulate meta-learning as a meta-optimization problem.

$$\min_{\lambda} \mathbb{E}_{D \sim \mathcal{D}} \min_{\theta} \mathbb{E}_{x \sim D} \mathcal{L}(x; \theta, \lambda).$$

- Need to optimize through the inner optimization.
  - BPTT
  - Fixed point (implicit differentiation)
  - Zeroth order optimization

# Meta-Optimization (Bi-level Optimization)

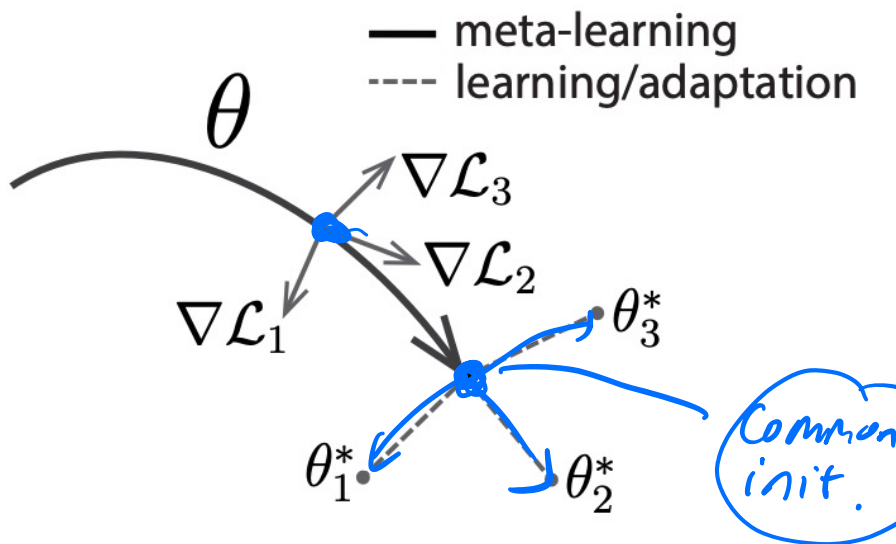
- One can formulate meta-learning as a meta-optimization problem.

$$\min_{\lambda} \mathbb{E}_{D \sim \mathcal{D}} \min_{\theta} \mathbb{E}_{x \sim D} \mathcal{L}(x; \theta, \lambda).$$

- Need to optimize through the inner optimization.
  - BPTT
  - Fixed point (implicit differentiation)
  - Zeroth order optimization
- Short-horizon bias: Optimal actions in the next few steps may not be optimal in the long run.
  - Example: Lowering learning rate will always result in short-term gains.

# MAML (Truncated Optimization)

- For few-shot learning, short horizon is actually needed.
- Unroll the gradient graph for a few iterations.
- MAML (Model-agnostic meta-learning)



---

**Algorithm 1** Model-Agnostic Meta-Learning

---

**Require:**  $p(\mathcal{T})$ : distribution over tasks

**Require:**  $\alpha, \beta$ : step size hyperparameters

- 1: randomly initialize  $\theta$
  - 2: **while** not done **do**
  - 3:   Sample batch of tasks  $\mathcal{T}_i \sim p(\mathcal{T})$
  - 4:   **for all**  $\mathcal{T}_i$  **do**
  - 5:     Evaluate  $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$  with respect to  $K$  examples
  - 6:     Compute adapted parameters with gradient descent:  $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
  - 7:   **end for**
  - 8:   Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$
  - 9: **end while**
-

# Hypernetworks, Best Response Functions

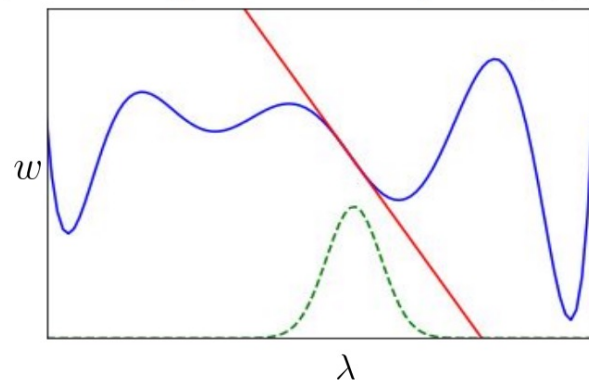
- Another way to “amortize” the optimization is to use a network to predict the optimal inner parameters.

# Hypernetworks, Best Response Functions

- Another way to “amortize” the optimization is to use a network to predict the optimal inner parameters.
- Best response function:  $w^* = \operatorname{argmin}_w \mathcal{L}_T(\lambda, w)$ .  $\lambda^* = \operatorname{argmin}_\lambda \mathcal{L}(\lambda, w^*)$ .

# Hypernetworks, Best Response Functions

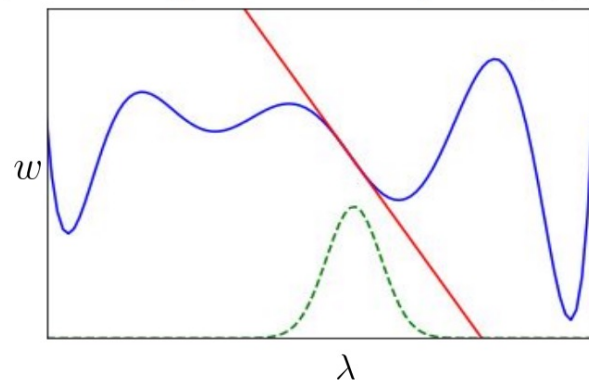
- Another way to “amortize” the optimization is to use a network to predict the optimal inner parameters.
- Best response function:  $w^* = \operatorname{argmin}_w \mathcal{L}_T(\lambda, w)$ .  $\lambda^* = \operatorname{argmin}_\lambda \mathcal{L}(\lambda, w^*)$ .
- Approximation:  $\lambda^* \approx \operatorname{argmin}_\lambda \mathcal{L}(\lambda, \hat{w}_\phi(\lambda))$ .  $\min_\phi \mathbb{E}_\epsilon [f(\lambda + \epsilon, \hat{w}(\lambda + \epsilon))]$ .





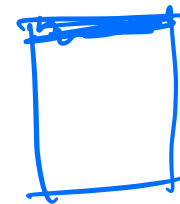
# Hypernetworks, Best Response Functions

- Another way to “amortize” the optimization is to use a network to predict the optimal inner parameters.
- Best response function:  $w^* = \operatorname{argmin}_w \mathcal{L}_T(\lambda, w)$ .  $\lambda^* = \operatorname{argmin}_\lambda \mathcal{L}(\lambda, w^*)$ .
- Approximation:  $\lambda^* \approx \operatorname{argmin}_\lambda \mathcal{L}(\lambda, \hat{w}_\phi(\lambda))$ .  $\min_\phi \mathbb{E}_\epsilon[f(\lambda + \epsilon, \hat{w}(\lambda + \epsilon))]$ .
- Scalable version? Instead of predicting full parameters?

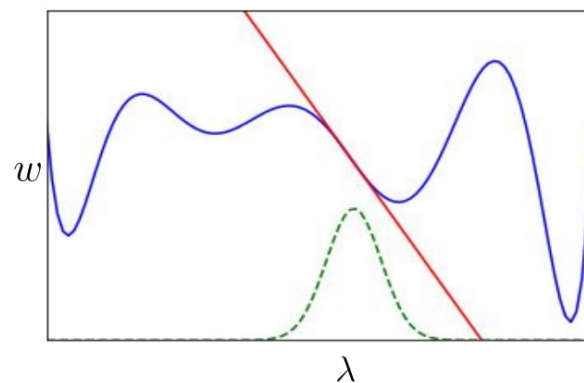
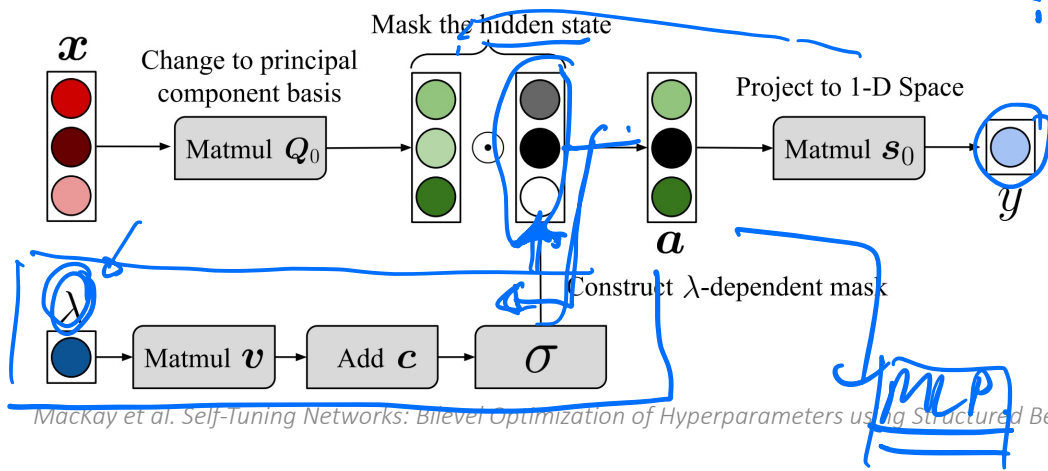


approx. inner loop.

# Hypernetworks, Best Response Functions



- Another way to “amortize” the optimization is to use a network to predict the optimal inner parameters.
- Best response function:  $w^* = \operatorname{argmin}_w \mathcal{L}_T(\lambda, w)$ .  $\lambda^* = \operatorname{argmin}_\lambda \mathcal{L}(\lambda, w^*)$ .
- Approximation:  $\lambda^* \approx \operatorname{argmin}_\lambda \mathcal{L}(\lambda, \hat{w}_\phi(\lambda))$ .  $\min_\phi \mathbb{E}_\epsilon [f(\lambda + \epsilon, \hat{w}(\lambda + \epsilon))]$ .
- Scalable version? Instead of predicting full parameters?



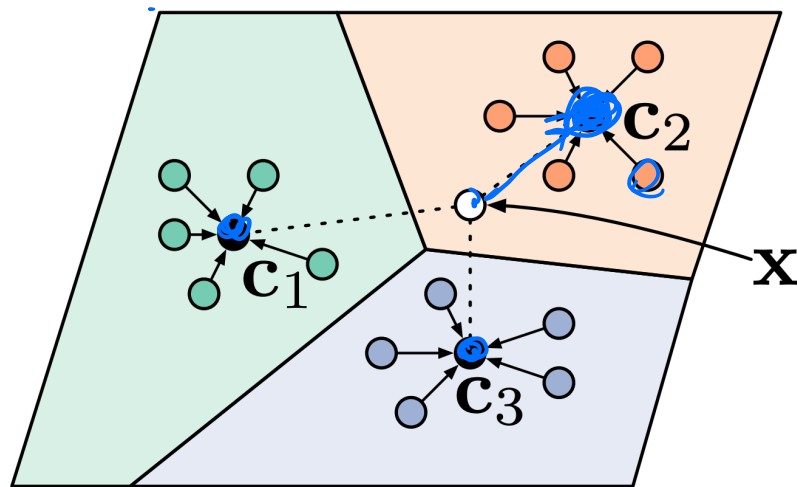
# Representation and Memory

- Prototypical Network: Few-shot Classification

# Representation and Memory

- Prototypical Network: Few-shot Classification
- Prototype = Avg. representation of a class

$$\mathbf{p}_k = \frac{1}{|S_k|} \sum_{(\mathbf{x}, y) \in S_k} f_\phi(\mathbf{x}_i).$$

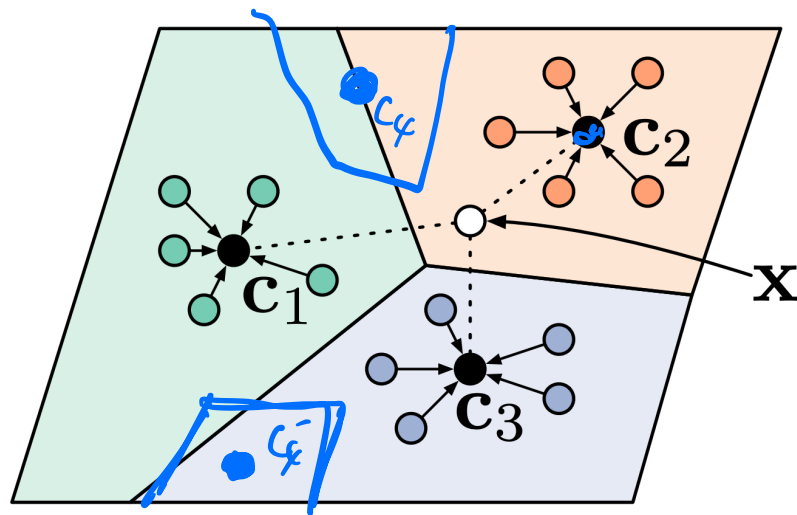


$$p(y = k \mid \mathbf{x}) = \text{softmax}(-d(f_\phi(\mathbf{x}), \mathbf{p}_k)).$$

# Representation and Memory

- Prototypical Network: Few-shot Classification
- Prototype = Avg. representation of a class
- 1 example = exemplar-based. Can be in between.

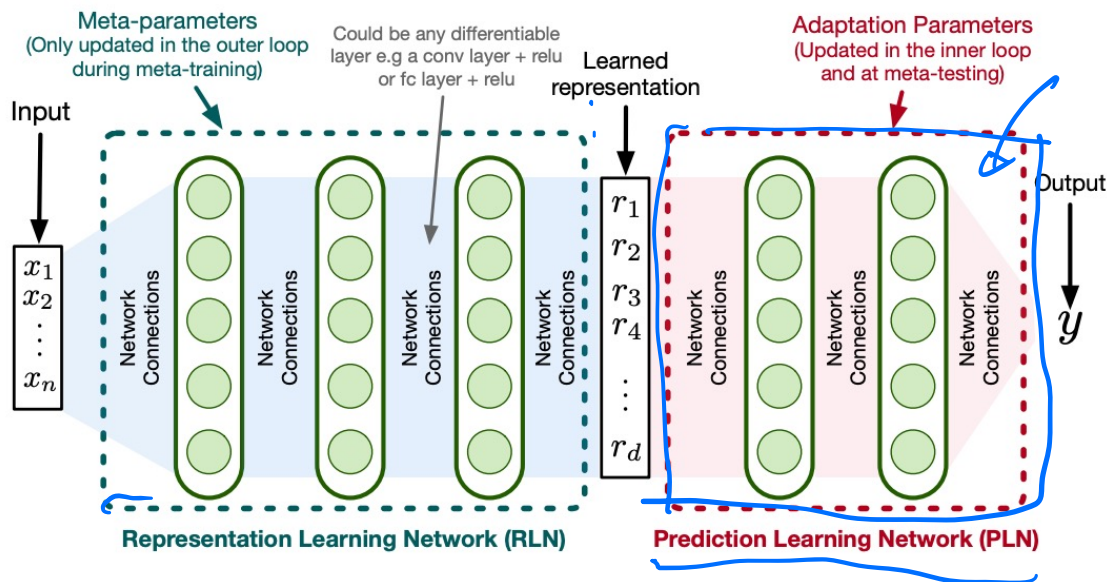
$$\mathbf{p}_k = \frac{1}{|S_k|} \sum_{(\mathbf{x}, y) \in S_k} f_\phi(\mathbf{x}_i).$$



$$p(y = k \mid \mathbf{x}) = \text{softmax}(-d(f_\phi(\mathbf{x}), \mathbf{p}_k)).$$

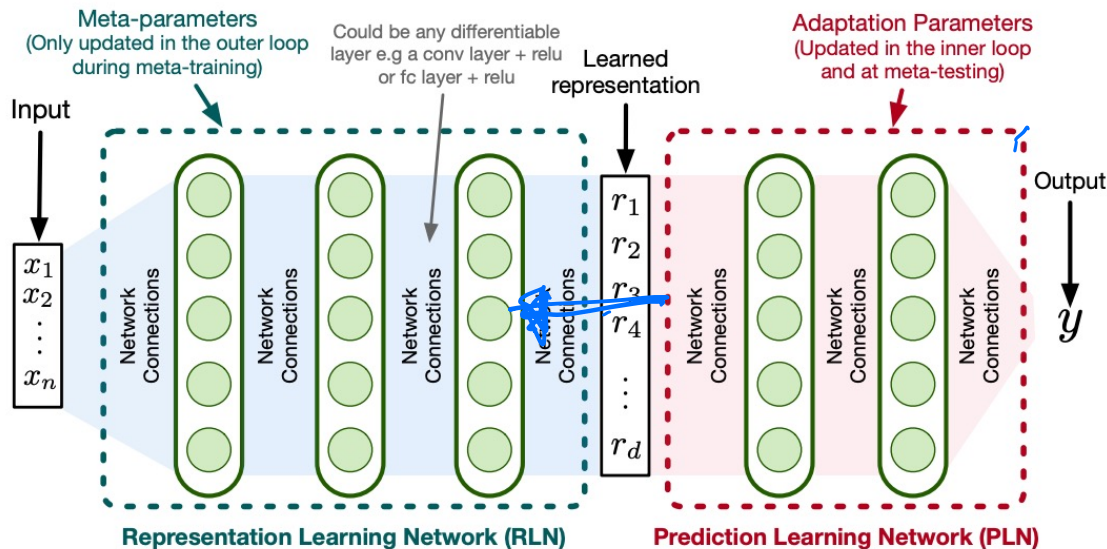
# Representation and Memory

- Representation vs. Memory Layers



# Representation and Memory

- Representation vs. Memory Layers
- Learning to continually learn



## Algorithm 2: Meta-Training : OML

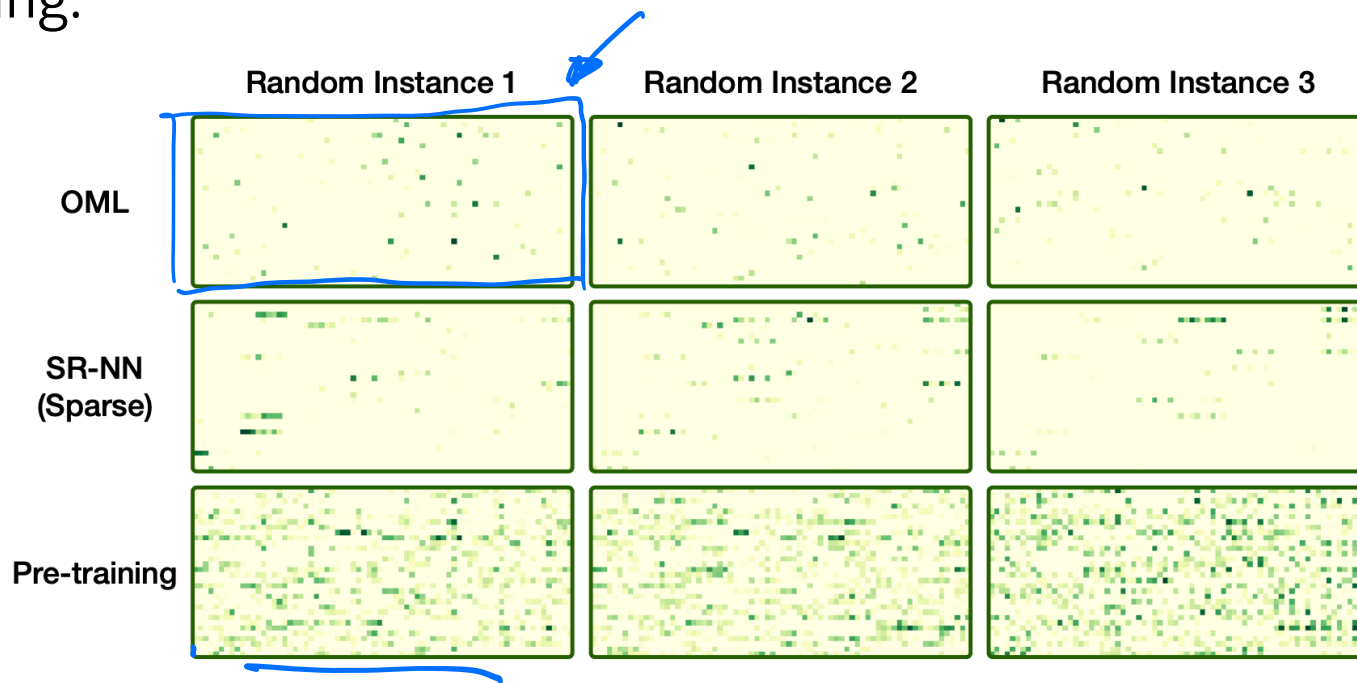
**Require:**  $p(\mathcal{T})$ : distribution over CLP problems

**Require:**  $\alpha, \beta$ : step size hyperparameters

- 1: randomly initialize  $\theta$
- 2: **while** not done **do**
- 3:   randomly initialize  $W$
- 4:   Sample CLP problem  $\mathcal{T}_i \sim p(\mathcal{T})$
- 5:   Sample  $\mathcal{S}_{train}$  from  $p(\mathcal{S}_k | \mathcal{T}_i)$
- 6:    $W_0 = W$
- 7:   **for**  $j = 1, 2, \dots, k$  **do**
- 8:      $(X_j, Y_j) = \mathcal{S}_{train}[j]$
- 9:      $W_j = W_{j-1} - \alpha \nabla_{W_{j-1}} \ell_i(f_{\theta, W_{j-1}}(X_j), Y_j)$
- 10:   **end for**
- 11:   Sample  $\mathcal{S}_{test}$  from  $p(\mathcal{S}_k | \mathcal{T}_i)$
- 12:   Update  $\theta \leftarrow \theta - \beta \nabla_{\theta} \ell_i(f_{\theta, W_k}(\mathcal{S}_{test}[:, 0]), \mathcal{S}_{test}[:, 1])$
- 13: **end while**

# Representation and Memory

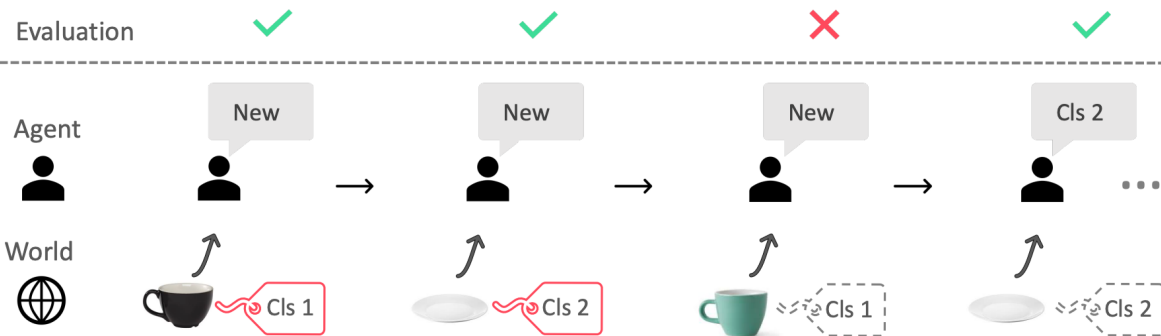
- Meta-learning leads to sparse representation suitable for continual learning.



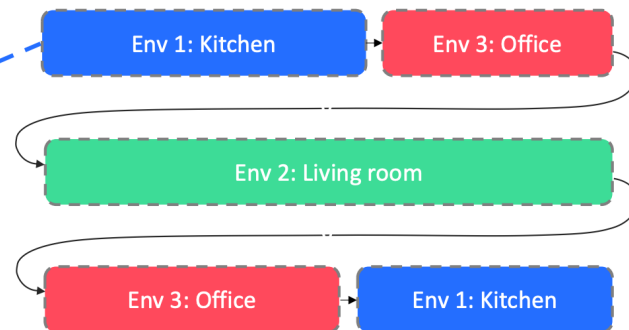


# Online Continual Few-Shot Learning

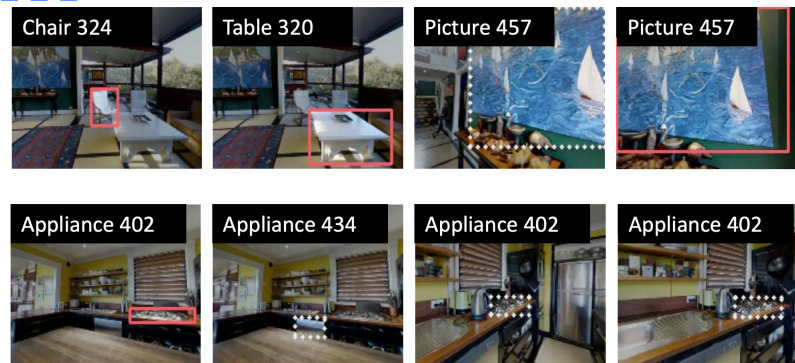
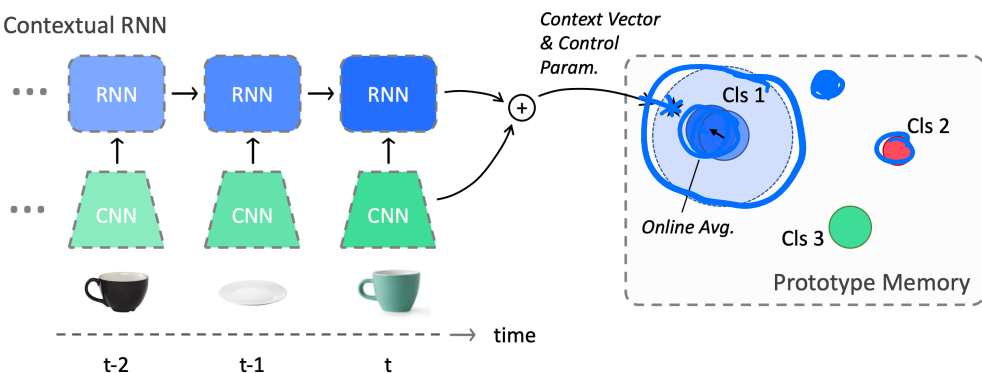
## A) Online evaluation with old and new classes



## B) Context switching



### Contextual RNN



# Online Continual Few-Shot Learning

- Compared to OML: Using prototype memory vs. generic MLP.

# Online Continual Few-Shot Learning

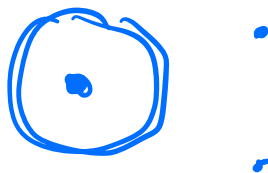
- Compared to OML: Using prototype memory vs. generic MLP.
- Both learning representations through online learning episodes.

# Online Continual Few-Shot Learning

- Compared to OML: Using prototype memory vs. generic MLP.
- Both learning representations through online learning episodes.
- Learning contextual representations (for context shifts).

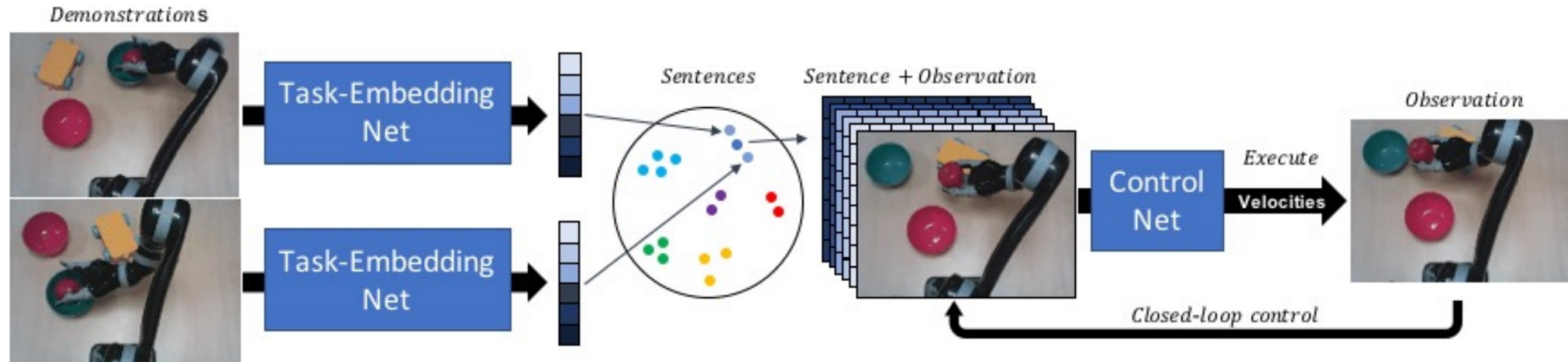
# Online Continual Few-Shot Learning

- Compared to OML: Using prototype memory vs. generic MLP.
- Both learning representations through online learning episodes.
- Learning contextual representations (for context shifts).
- Learning to output unknowns.



# Few-Shot Imitation Learning

- The idea of prototype learning can also be applied to skill learning.



Finn et al. One-Shot Visual Imitation Learning via Meta-Learning. CoRL 2017.

James et al. Task-Embedded Control Networks for Few-Shot Imitation Learning. CoRL 2018.

# Few-Shot Imitation Learning

- The idea of prototype learning can also be applied to skill learning.
- Embedding learning:

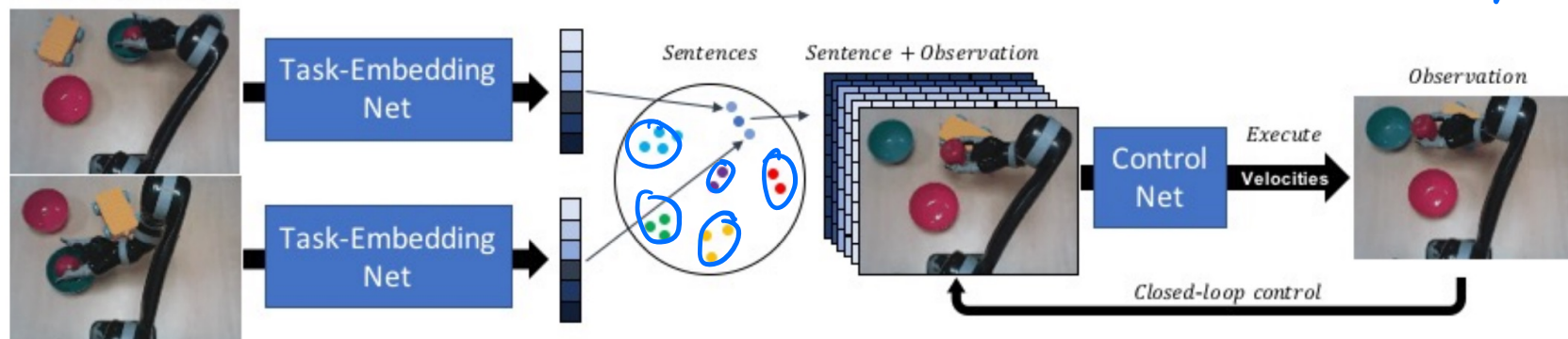
*traj.*

$$s^j = \frac{1}{K} \sum_{\tau_k^j \in \mathcal{T}^j} f_{\theta}(\tau_k^j).$$

*prototype.*

$$\mathcal{L}_{emb} = \sum_k \sum_{i \neq j} \max(0, \Delta - \tilde{s}_k^j \cdot \tilde{s}^j + \tilde{s}_k^j \cdot \tilde{s}^i).$$

*max-margin*

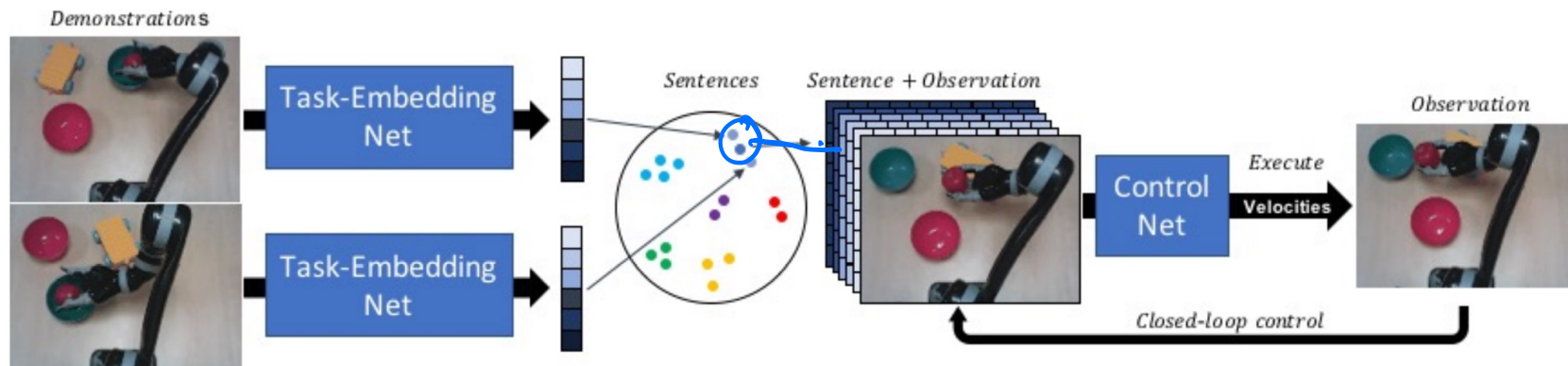


# Few-Shot Imitation Learning



- Control learning:

$$\mathcal{L}_{ctr} = \sum_j \sum_{(o,a) \in \tau^j} \left\| \underbrace{\pi(o, s^j)}_{\text{Output}} - \underbrace{a}_{\text{action}} \right\|_2^2$$

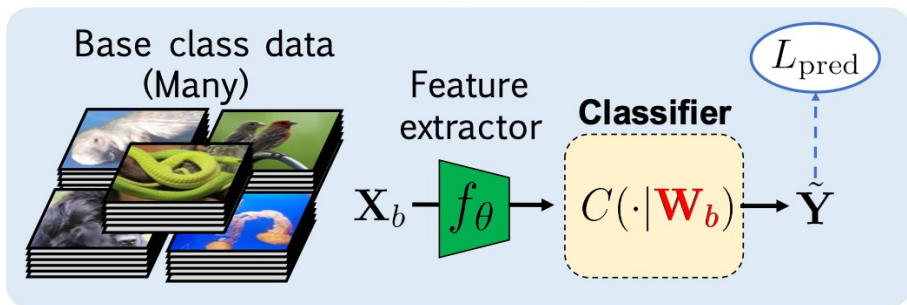




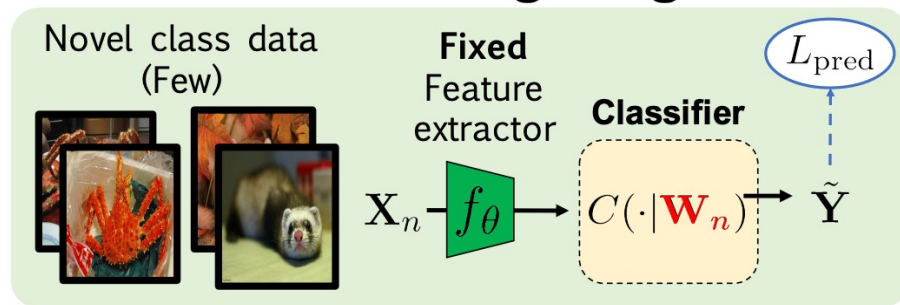
# Leveraging Pretrained Representations

- Need a strong pretrained network.

## Training stage



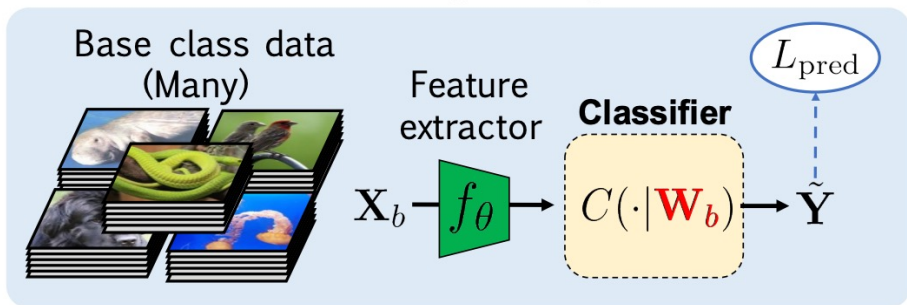
## Fine-tuning stage



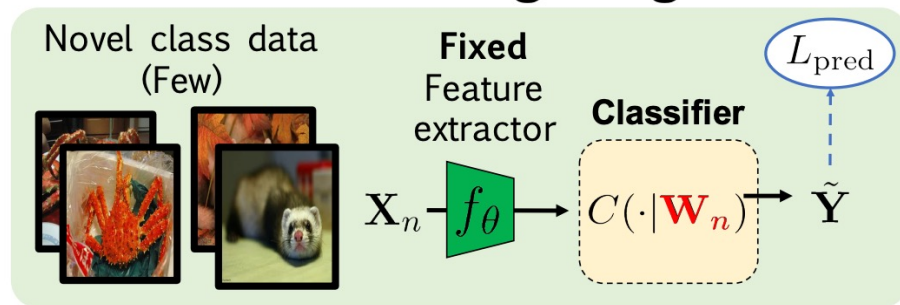
# Leveraging Pretrained Representations

- Need a strong pretrained network.
- Works well for few-shot classification.

## Training stage



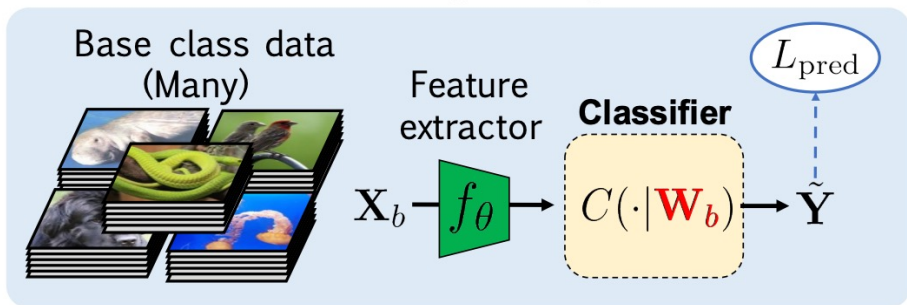
## Fine-tuning stage



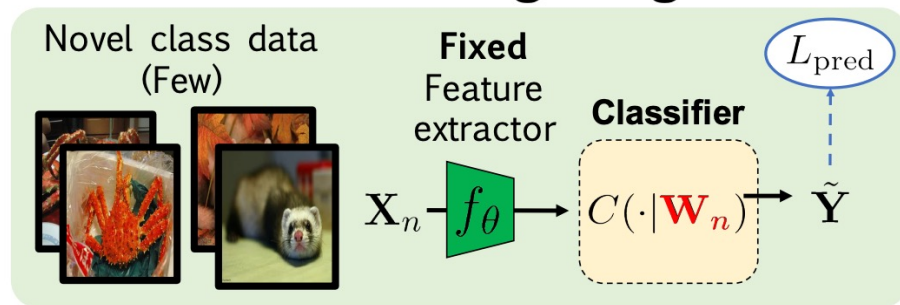
# Leveraging Pretrained Representations

- Need a strong pretrained network.
- Works well for few-shot classification.
- Once again proves that representation is crucial.

## Training stage



## Fine-tuning stage



# In-Context Learning (ICL)

- Traditionally, learning through parameters  $\theta$ .

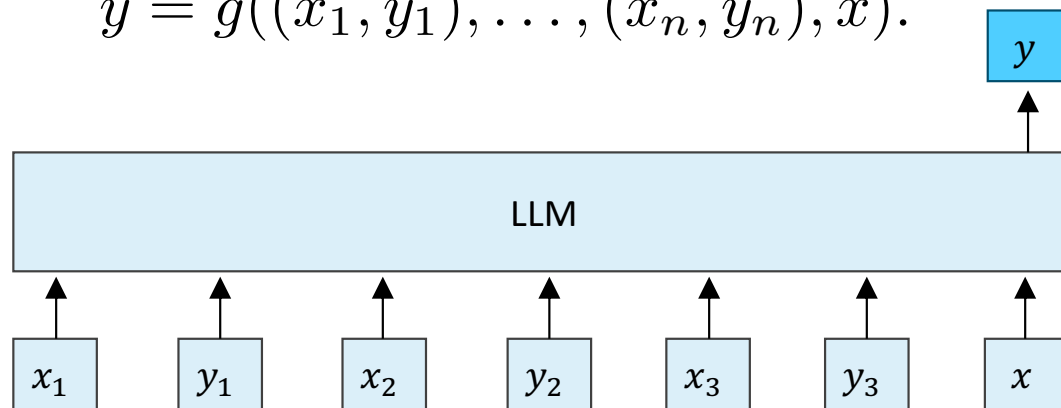
# In-Context Learning (ICL)

- Traditionally, learning through parameters  $\theta$ .
- ICL does not optimize any parameters, just put the training data in the context.

# In-Context Learning (ICL)

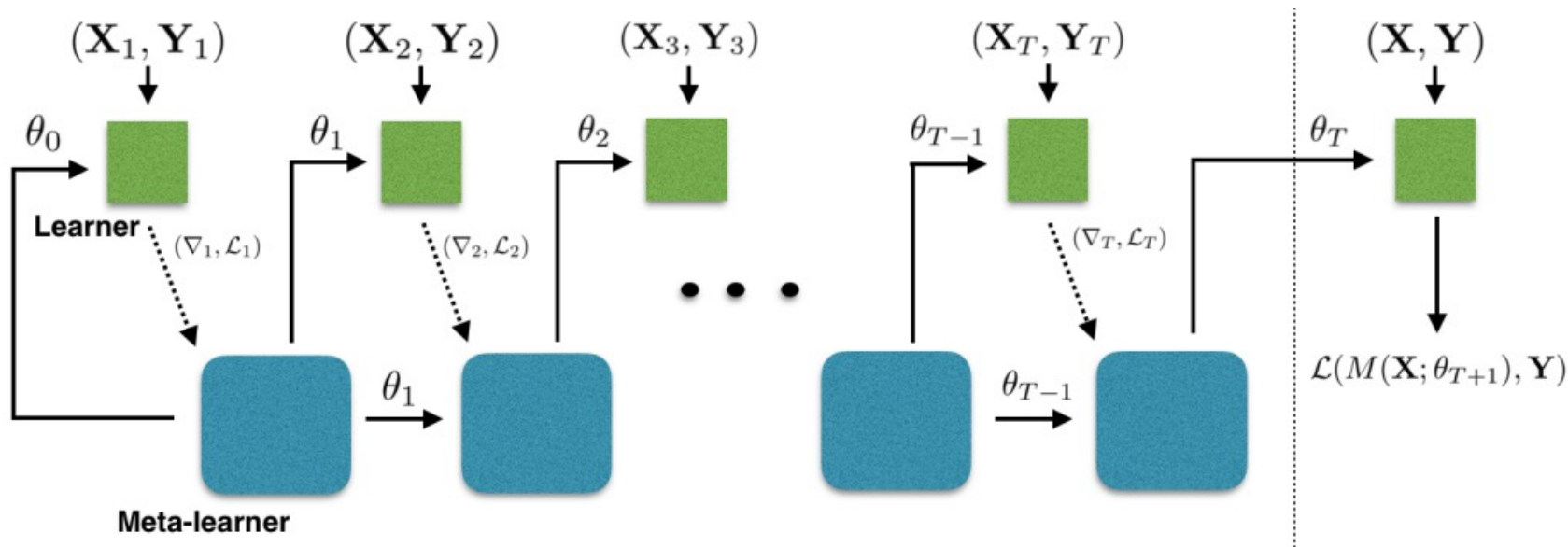
- Traditionally, learning through parameters  $\theta$ .
- ICL does not optimize any parameters, just put the training data in the context.
- Inner optimization loop done in a sequence model (RNN, Transformer, etc.)

$$y = g((x_1, y_1), \dots, (x_n, y_n), x).$$



# Meta-Learning LSTM

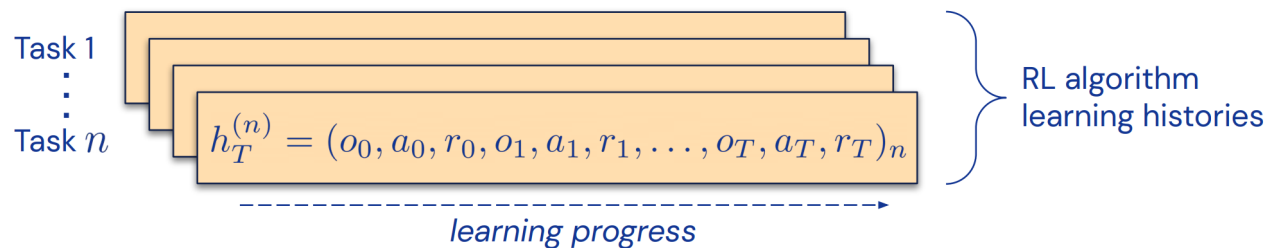
- An earlier sequential meta-learning paradigm before ICL.
- Using hidden states as “parameters”



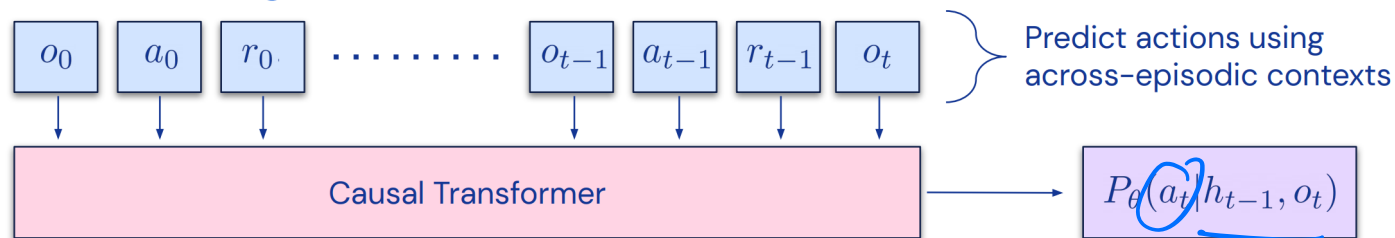
# In-Context RL

- Use sequences of Observation ( $o_t$ ), Action ( $a_t$ ) and Reward ( $r_t$ ) generated by standard RL algorithms.

## Data Generation



## Model Training





# Test-Time Tuning/Adaptation

- Context Retrieval, e.g. nearest neighbors
- Finetuning on examples
  - Full finetuning
  - Low-rank adaptation
  - Prompt tuning
- Self-supervision objectives

# Summary: Towards Embodied Learning Agents

- A modularized but differentiable end-to-end architecture
  - Perception
  - Prediction
  - Planning
  - Mapping
  - Memory

# Summary: Towards Embodied Learning Agents

- A modularized but differentiable end-to-end architecture
  - Perception
  - Prediction
  - Planning
  - Mapping
  - Memory
- A combination of learning signals

# Summary: Towards Embodied Learning Agents

- A modularized but differentiable end-to-end architecture
  - Perception
  - Prediction
  - Planning
  - Mapping
  - Memory
- A combination of learning signals
  - Direct supervision, energy-based (MSE, Contrastive, Max-Margin, Denoising)

# Summary: Towards Embodied Learning Agents

- A modularized but differentiable end-to-end architecture
  - Perception
  - Prediction
  - Planning
  - Mapping
  - Memory
- A combination of learning signals
  - Direct supervision, energy-based (MSE, Contrastive, Max-Margin, Denoising)
  - Self-supervision (InfoNCE, MSE, Cross-entropy)

# Summary: Towards Embodied Learning Agents

- A modularized but differentiable end-to-end architecture
  - Perception
  - Prediction
  - Planning
  - Mapping
  - Memory
- A combination of learning signals
  - Direct supervision, energy-based (MSE, Contrastive, Max-Margin, Denoising)
  - Self-supervision (InfoNCE, MSE, Cross-entropy)
  - Reconstruction

# Summary: Towards Embodied Learning Agents

- A modularized but differentiable end-to-end architecture
  - Perception
  - Prediction
  - Planning
  - Mapping
  - Memory
- A combination of learning signals
  - Direct supervision, energy-based (MSE, Contrastive, Max-Margin, Denoising)
  - Self-supervision (InfoNCE, MSE, Cross-entropy)
  - Reconstruction
  - Future prediction

# Summary: Towards Embodied Learning Agents

- A modularized but differentiable end-to-end architecture
  - Perception
  - Prediction
  - Planning
  - Mapping
  - Memory
- A combination of learning signals
  - Direct supervision, energy-based (MSE, Contrastive, Max-Margin, Denoising)
  - Self-supervision (InfoNCE, MSE, Cross-entropy)
  - Reconstruction
  - Future prediction ✓
  - Reinforcement learning



# Summary: Towards Embodied Learning Agents

- Useful inductive biases

# Summary: Towards Embodied Learning Agents

- Useful inductive biases
  - Spatial grounding

# Summary: Towards Embodied Learning Agents

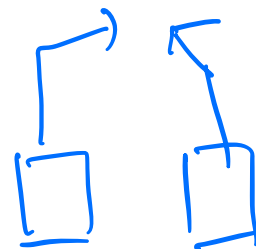
- Useful inductive biases
  - Spatial grounding
  - Geometric projections / transformations

# Summary: Towards Embodied Learning Agents

- Useful inductive biases
  - Spatial grounding
  - Geometric projections / transformations
  - Permutation invariance / equivariance

# Summary: Towards Embodied Learning Agents

- Useful inductive biases
  - Spatial grounding
  - Geometric projections / transformations
  - Permutation invariance / equivariance
  - Representation invariance



# Summary: Towards Embodied Learning Agents

- Useful inductive biases
  - Spatial grounding
  - Geometric projections / transformations
  - Permutation invariance / equivariance
  - Representation invariance
  - Motion grounding

# Summary: Towards Embodied Learning Agents

- Useful inductive biases
  - Spatial grounding
  - Geometric projections / transformations
  - Permutation invariance / equivariance
  - Representation invariance
  - Motion grounding
  - Disentanglement, object-centric latents

# Summary: Towards Embodied Learning Agents

- Useful inductive biases
  - Spatial grounding
  - Geometric projections / transformations
  - Permutation invariance / equivariance
  - Representation invariance
  - Motion grounding
  - Disentanglement, object-centric latents
  - Cost volumes



# Summary: Towards Embodied Learning Agents

- Useful inductive biases
  - Spatial grounding
  - Geometric projections / transformations
  - Permutation invariance / equivariance
  - Representation invariance
  - Motion grounding
  - Disentanglement, object-centric latents
  - Cost volumes
  - Recursive iteration

# Summary: Towards Embodied Learning Agents

- Useful inductive biases
  - Spatial grounding
  - Geometric projections / transformations
  - Permutation invariance / equivariance
  - Representation invariance
  - Motion grounding
  - Disentanglement, object-centric latents
  - Cost volumes
  - Recursive iteration
  - Optimization / fixed point iteration

# Summary: Towards Embodied Learning Agents

- Useful inductive biases
  - Spatial grounding
  - Geometric projections / transformations
  - Permutation invariance / equivariance
  - Representation invariance
  - Motion grounding
  - Disentanglement, object-centric latents
  - Cost volumes
  - Recursive iteration
  - Optimization / fixed point iteration
  - Memory, replay, sparsity

# Summary: Towards Embodied Learning Agents

- Useful inductive biases
  - Spatial grounding
  - Geometric projections / transformations
  - Permutation invariance / equivariance
  - Representation invariance
  - Motion grounding
  - Disentanglement, object-centric latents
  - Cost volumes
  - Recursive iteration
  - Optimization / fixed point iteration
  - Memory, replay, sparsity
  - Learning inductive biases, learning to learn

# Topic Presentations

	Component 1	Component 2
Week 7 (Mar 6)	Continual Learning, Few-shot Learning (1 hr)	<b>Deep Learning for Structured Prediction</b> Tanishq Sardana, Qing Mu, Owais Shuja
Week 8 (Mar 13)	Guest Lecture – Prof. Wei-Chiu Ma (1 hr)	<b>3D Vision and Mapping</b> Sihang Li, Kanishkha Jaisankar, Denis Mbey Akola, Zijin Hu
Week 9 (Mar 20)	<b>SSL and Object Discovery</b> Anurup Naskar, Dahye Kim, Sal Yeung, Surbhi (1.5 hr)	<b>World Model 1</b> Sidhartha Reddy Potu, Andrew Deur
Week 11 (Apr 3)	<b>World Model 2</b> Pratyaksh Prabhav Rao, Sergey Sedov, Rooholla Khorrambakht	<b>End-to-End Planning</b> Raman Kumar Jha, Jovita Gandhi, Sushma Mareddy, Mrunal Sarvaiya
Week 12 (Apr 10)	<b>Continual Learning</b> Akshay Raman, Amey Joshi, Zifan Zhao	<b>Few-Shot Learning</b> Ellen Su, Xu Zhang, Swarali Borde
Week 13 (Apr 17)	Guest Lecture – Dr. Andrei Barsan (1hr)	<b>LLM Agents</b> Solim LeGris, Ravan Budda, Dan Zhao, Sunidhi Tandel

# Project Presentations

Week 14 (Apr 24)      Project Presentations (7 teams)

Week 15 (May 1)      Project Presentations (7 teams)

- Apr 10: Sign up for a presentation slot. Week 14 Presenters get 2% bonus. First come first serve.
- Project topics and proposals to be shared in the class.

# Today

- Tanishq Sardana: Segment Anything
- Qing Mu: DETR: End-to-End Object Detection
- Owais Saad Shuja: Latent Diffusion Models
- Discussion