

Continual Machine Learning

Summer 2024

Teacher

Dr. Martin Mundt,
Research Group on Open World Lifelong Learning

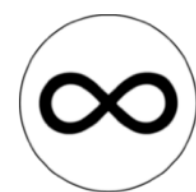
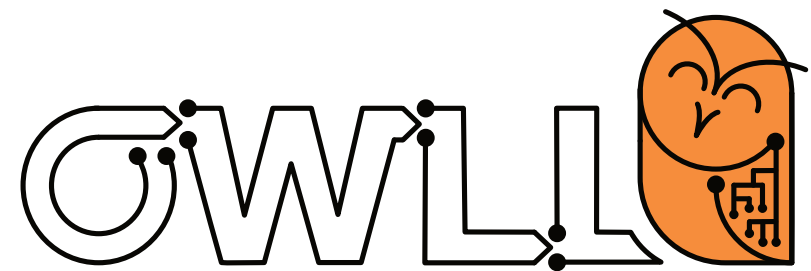
Time

Every Friday 14:25 - 16:05 CEST

Course Homepage

http://owll-lab.com/teaching/cl_lecture_24

<https://www.youtube.com/playlist?list=PLm6QXeaB-XkA5-IVBB-h7XeYzFzgSh6sk>



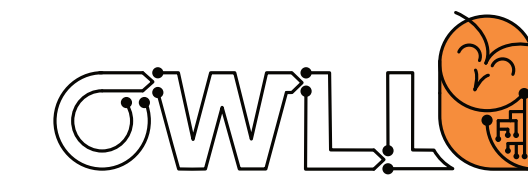
Continual **AI**



hessian.AI



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Week 4: Rehearsal - Knowledge Retention Part 2

Recall: How to avoid forgetting?



Paradigms for Continual Learning

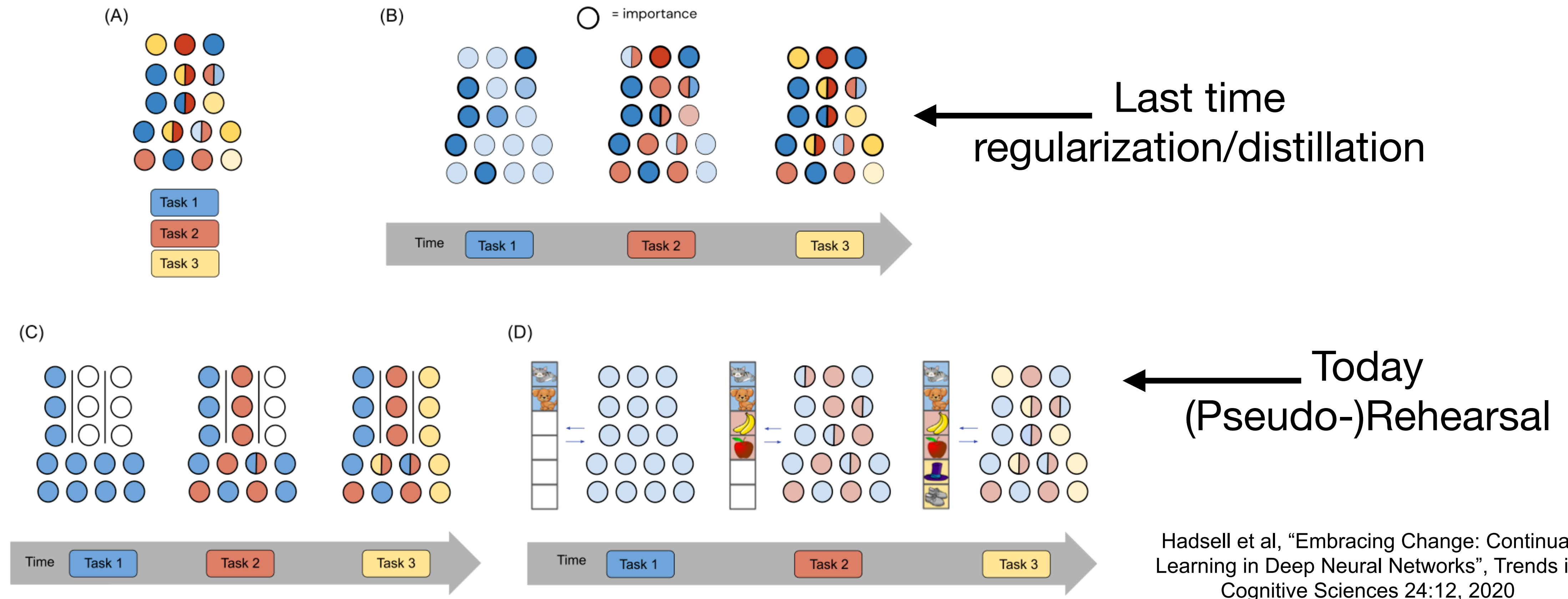
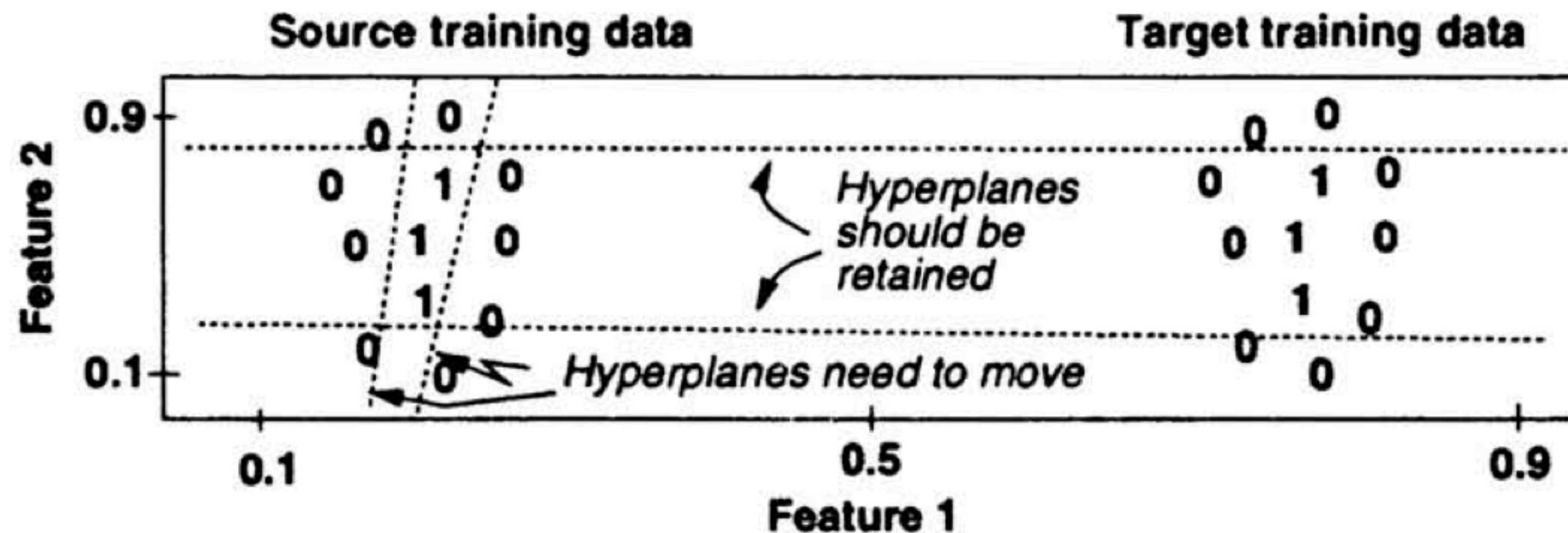


Figure 1. (A) Independent and identically distributed learning methods are standard for nonsequential, multitask learning. In this regime, tasks are learned simultaneously to avoid forgetting and instability. (B) Gradient-based approaches preserve parameters based on their importance to previously learned tasks. (C) Modularity-based methods define hard boundaries to separate task-specific parameters (often accompanied by shared parameters to allow transfer). (D) Memory-based methods write experience to memory to avoid forgetting.

Rehearsal: intuition

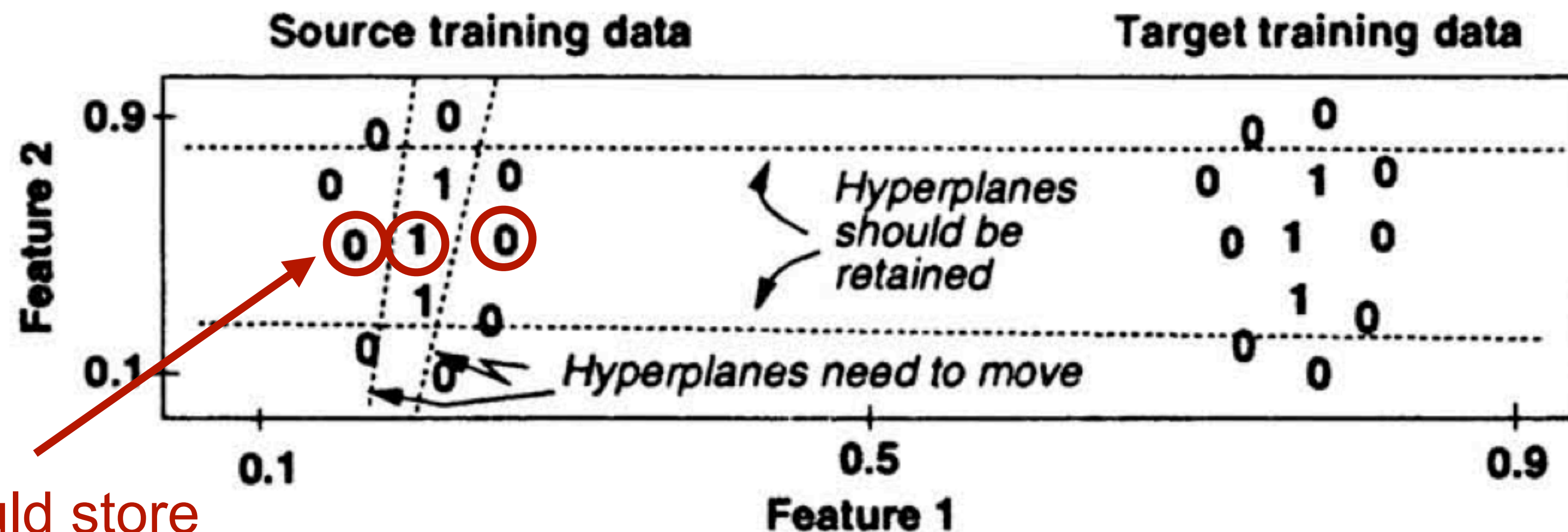


Assuming privacy is not a concern & that we have auxiliary memory:
some data is more relevant than other, can we retain a subset?



Rehearsal: intuition

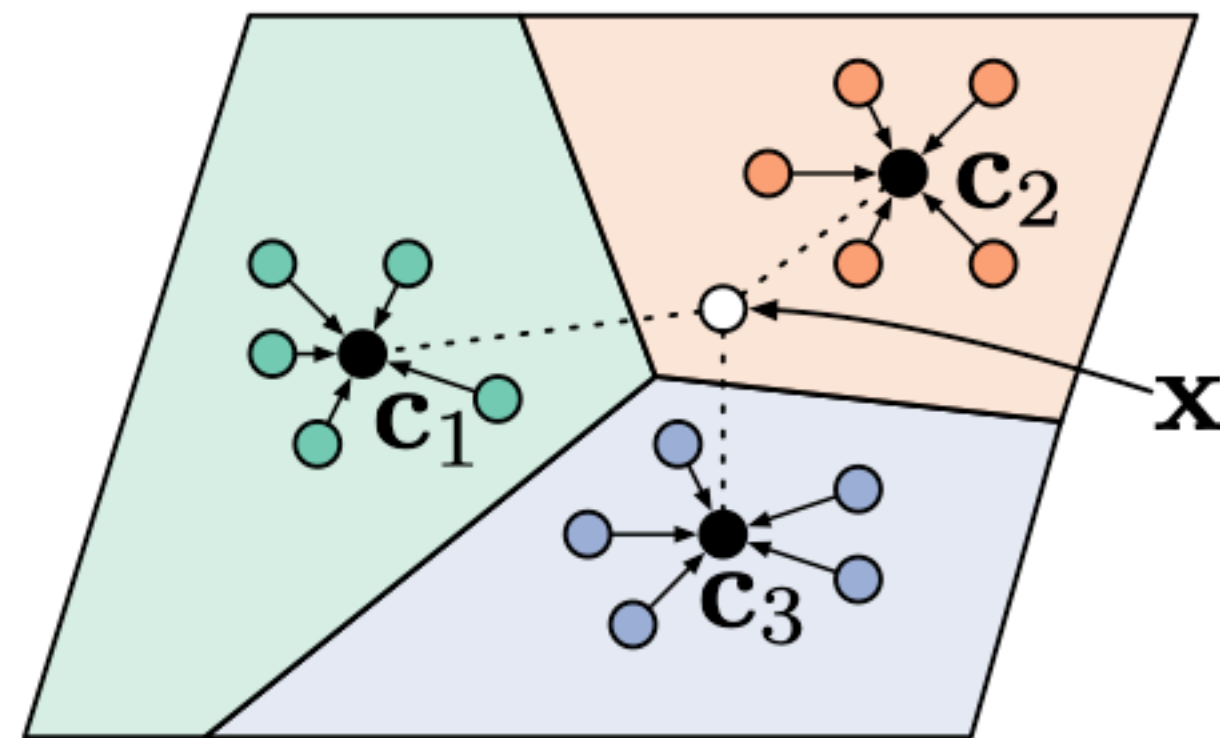
Assuming privacy is not a concern & that we have auxiliary memory:
some data is more relevant than other, can we retain a subset?



Maybe we could store
these few examples?

Nearest means classifiers

Let's start with an example this time and then discuss desiderata



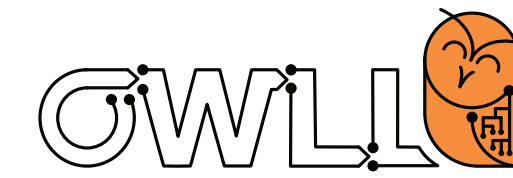
(a) Few-shot

Recall lecture 2: few shot learning

“Prototypical Networks for Few-shot Learning”,
Snell et al, NeurIPS 2017

- Get prototype as mean vector of each class
- Given a distance function d , classify according to distances to the prototypes in embedding space
- Note, nomenclature is inconsistent: prototypes, exemplars, core sets (will be defined later)

iCaRL: classification



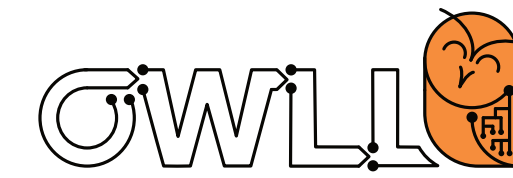
iCaRL: incremental classifier and representation learning

- Stores a subset of data in a fixed size memory buffer
- Classifies based on nearest class means
- Consecutively replaces parts of memory buffer with new examples

Algorithm 1 iCaRL CLASSIFY

```
input  $x$  // image to be classified
require  $\mathcal{P} = (P_1, \dots, P_t)$  // class exemplar sets
require  $\varphi : \mathcal{X} \rightarrow \mathbb{R}^d$  // feature map
  for  $y = 1, \dots, t$  do
     $\mu_y \leftarrow \frac{1}{|P_y|} \sum_{p \in P_y} \varphi(p)$  // mean-of-exemplars
  end for
   $y^* \leftarrow \operatorname{argmin}_{y=1, \dots, t} \|\varphi(x) - \mu_y\|$  // nearest prototype
output class label  $y^*$ 
```

iCaRL: picking exemplars



How is our memory buffer filled with specific examples?

- Iteratively: one by one, based on “herding” (Welling ICML 2009)
- Pick exemplars to best approximate the overall mean
- For a size of k exemplars: loop k times

Algorithm 4 iCaRL CONSTRUCTEXEMPLARSET

input image set $X = \{x_1, \dots, x_n\}$ of class y

input m target number of exemplars

require current feature function $\varphi : \mathcal{X} \rightarrow \mathbb{R}^d$

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

for $k = 1, \dots, m$ **do**

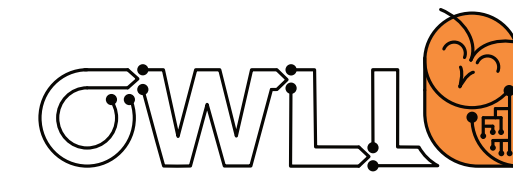
$p_k \leftarrow \operatorname{argmin}_{x \in X} \left\| \mu - \frac{1}{k} [\varphi(x) + \sum_{j=1}^{k-1} \varphi(p_j)] \right\|$

end for

$P \leftarrow (p_1, \dots, p_m)$

output exemplar set P

iCaRL: replacing exemplars



Our memory buffer is limited, how do we later remove samples?

- Memory buffer is like a prioritized list
- Later picked exemplars for a task “weigh” less
- Simply cut and repopulate

Algorithm 5 iCaRL REDUCEEXEMPLARSET

input m // target number of exemplars
input $P = (p_1, \dots, p_{|P|})$ // current exemplar set
 $P \leftarrow (p_1, \dots, p_m)$ // *i.e.* keep only first m
output exemplar set P

iCaRL: training



How do we train incrementally?

- Concatenate dataset with exemplars/
interleave exemplars into training
- Pick new exemplars (not shown on
the right) + replace existing
- Additionally use knowledge distillation

Algorithm 3 iCaRL UPDATE REPRESENTATION

input X^s, \dots, X^t // training images of classes s, \dots, t
require $\mathcal{P} = (P_1, \dots, P_{s-1})$ // exemplar sets
require Θ // current model parameters
// form combined training set:

$$\mathcal{D} \leftarrow \bigcup_{y=s, \dots, t} \{(x, y) : x \in X^y\} \cup \bigcup_{y=1, \dots, s-1} \{(x, y) : x \in P^y\}$$

// store network outputs with pre-update parameters:

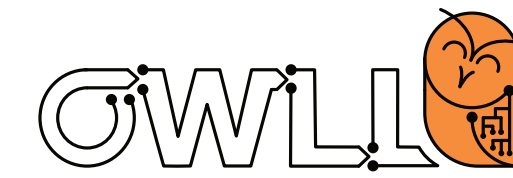
for $y = 1, \dots, s - 1$ **do**
 $q_i^y \leftarrow g_y(x_i)$ for all $(x_i, \cdot) \in \mathcal{D}$
end for

run network training (e.g. BackProp) with loss function

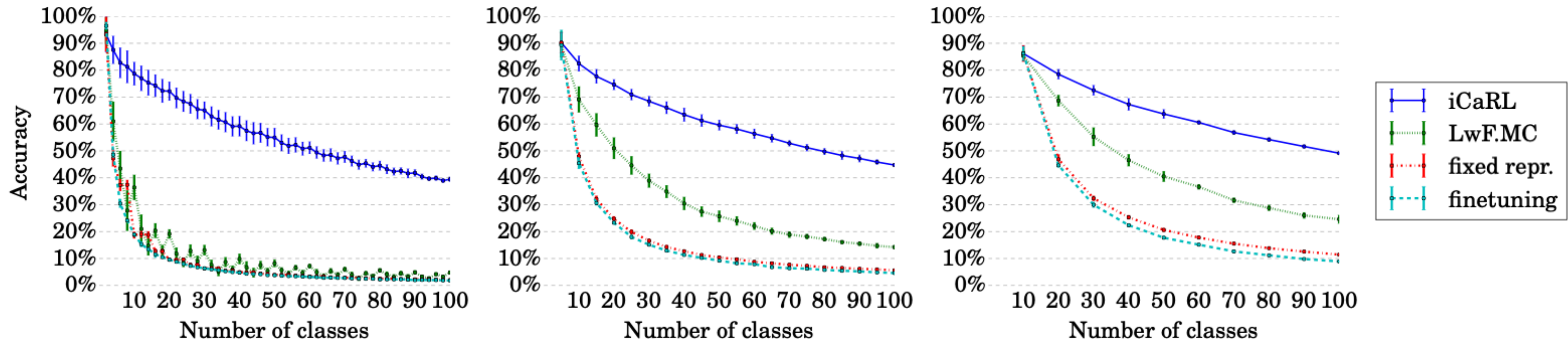
$$\ell(\Theta) = - \sum_{(x_i, y_i) \in \mathcal{D}} \left[\sum_{y=s}^t \delta_{y=y_i} \log g_y(x_i) + \delta_{y \neq y_i} \log(1 - g_y(x_i)) \right. \\ \left. + \sum_{y=1}^{s-1} q_i^y \log g_y(x_i) + (1 - q_i^y) \log(1 - g_y(x_i)) \right]$$

that consists of *classification* and *distillation* terms.

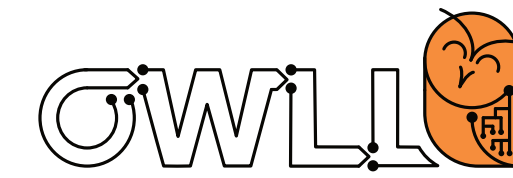
iCaRL & knowledge distillation



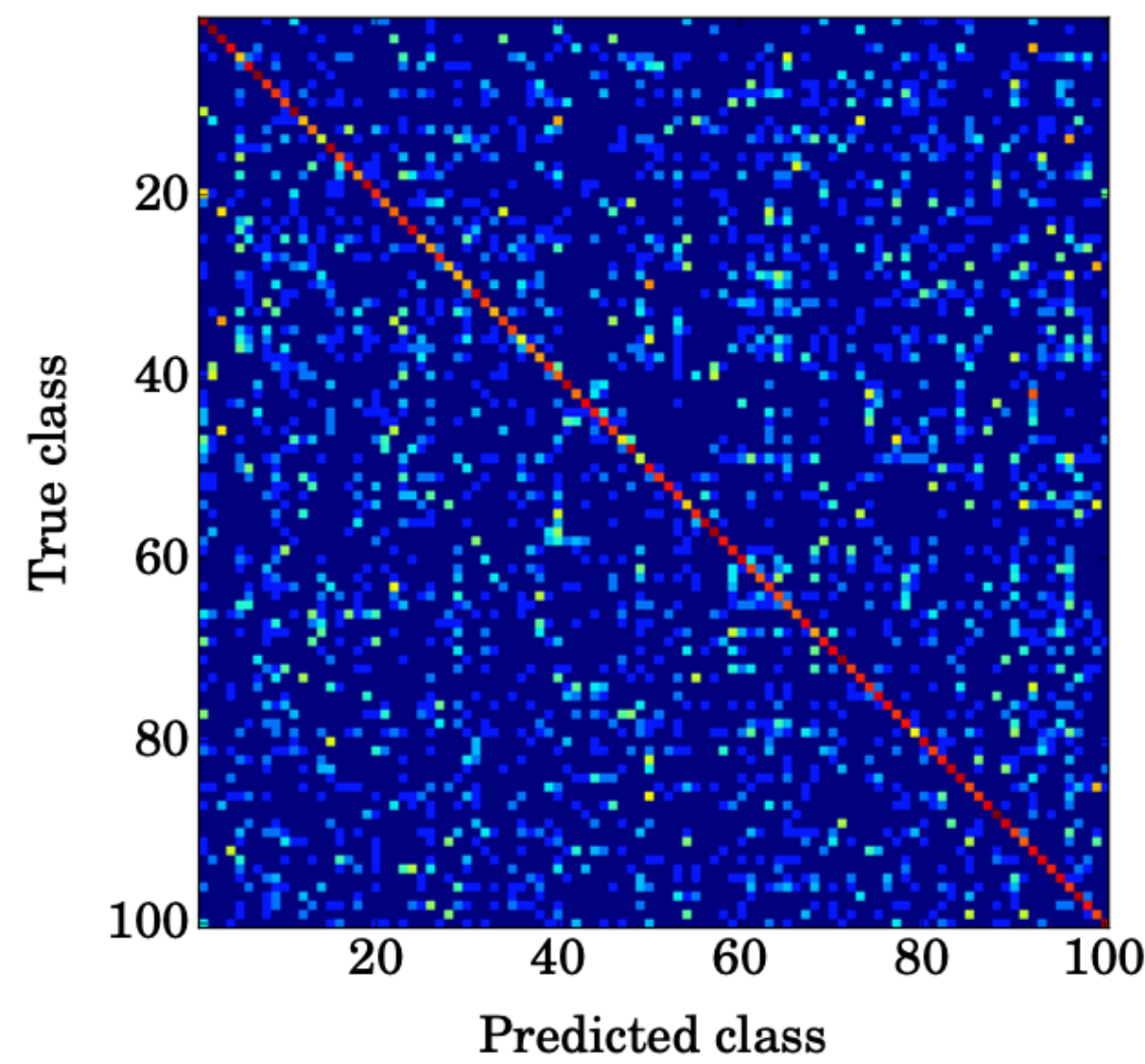
Example where classes of CIFAR100 are learned incrementally: exemplars are crucial



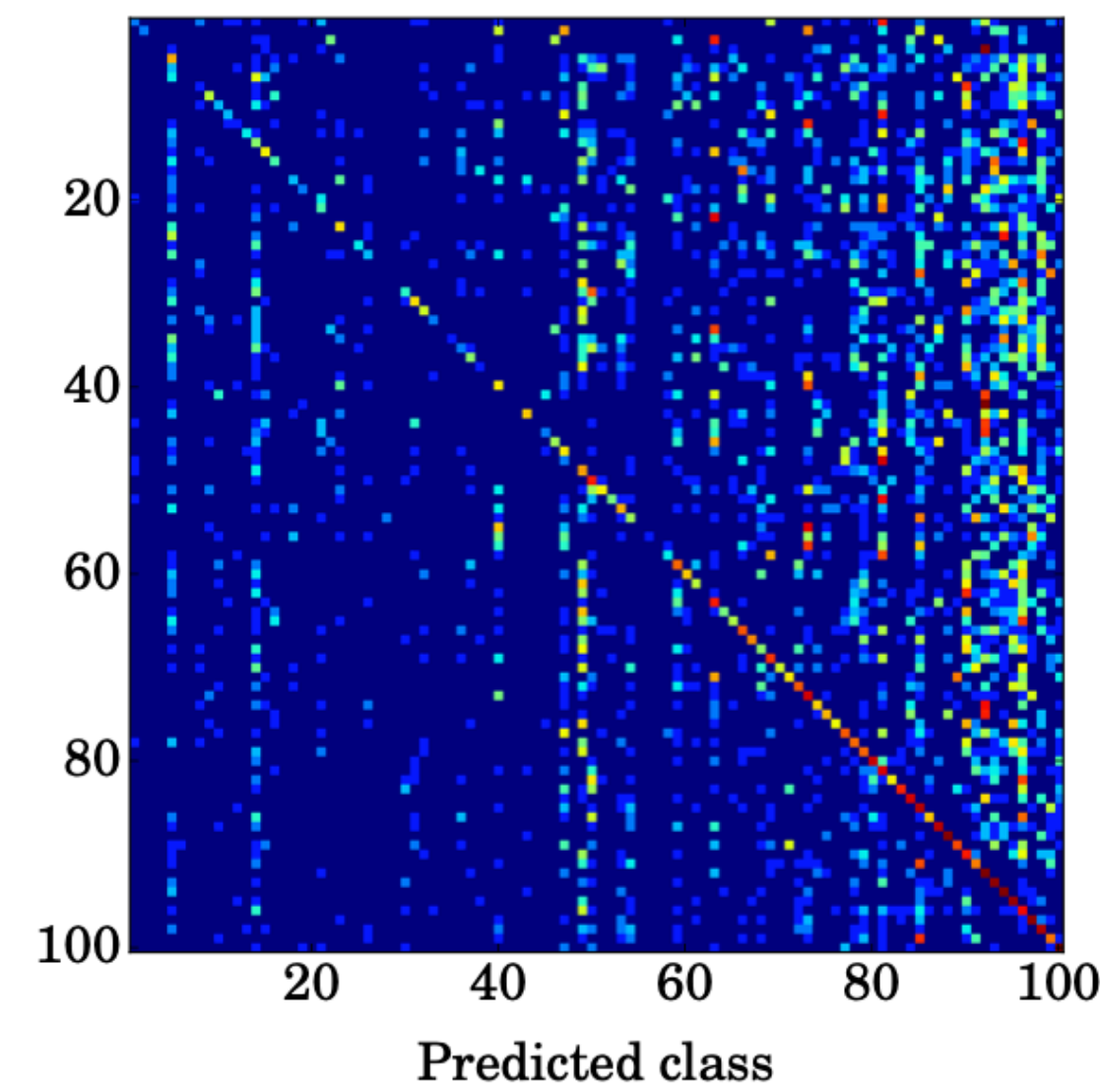
iCaRL & knowledge distillation



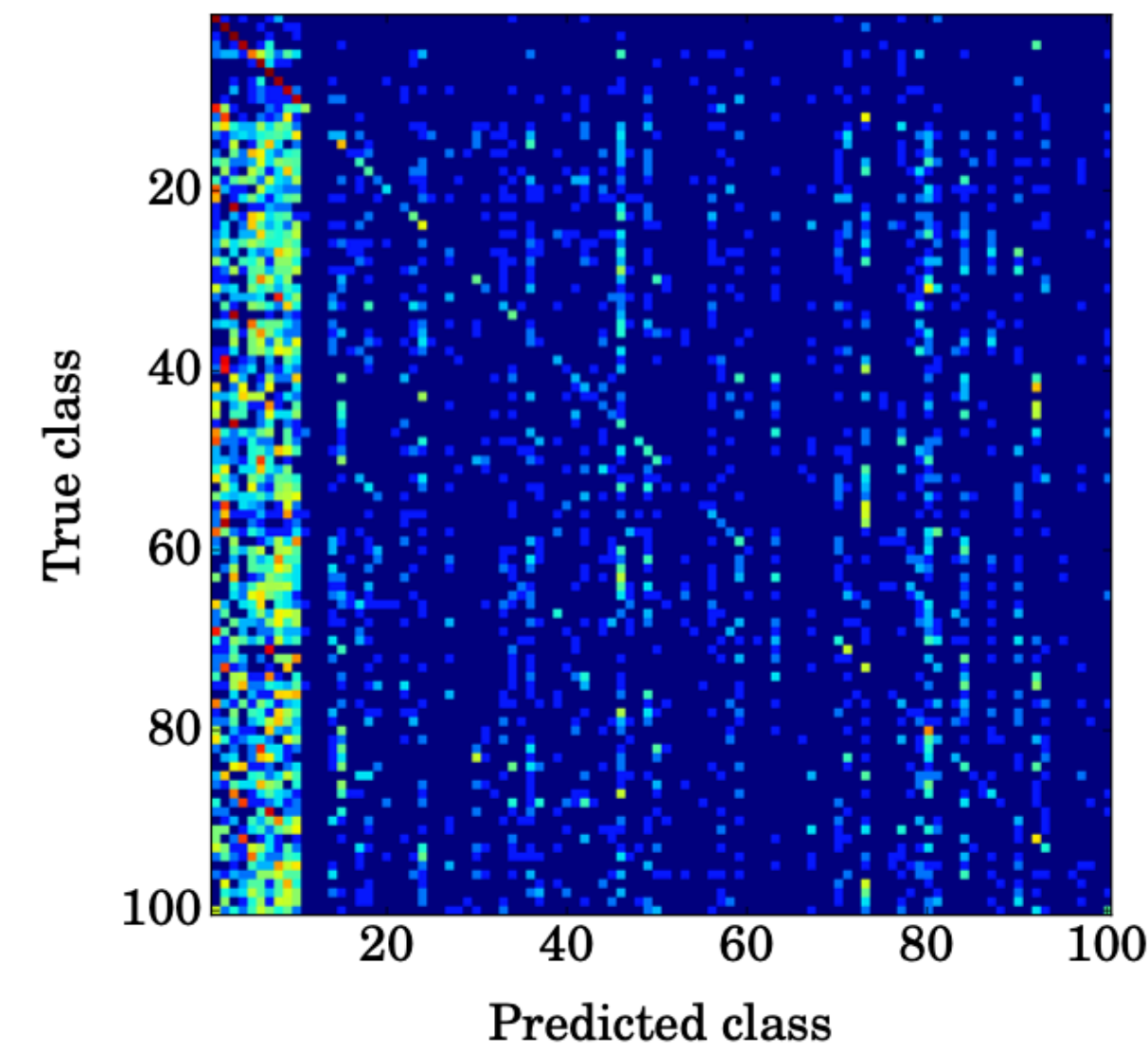
Confusion matrices empirically confirm our intuition



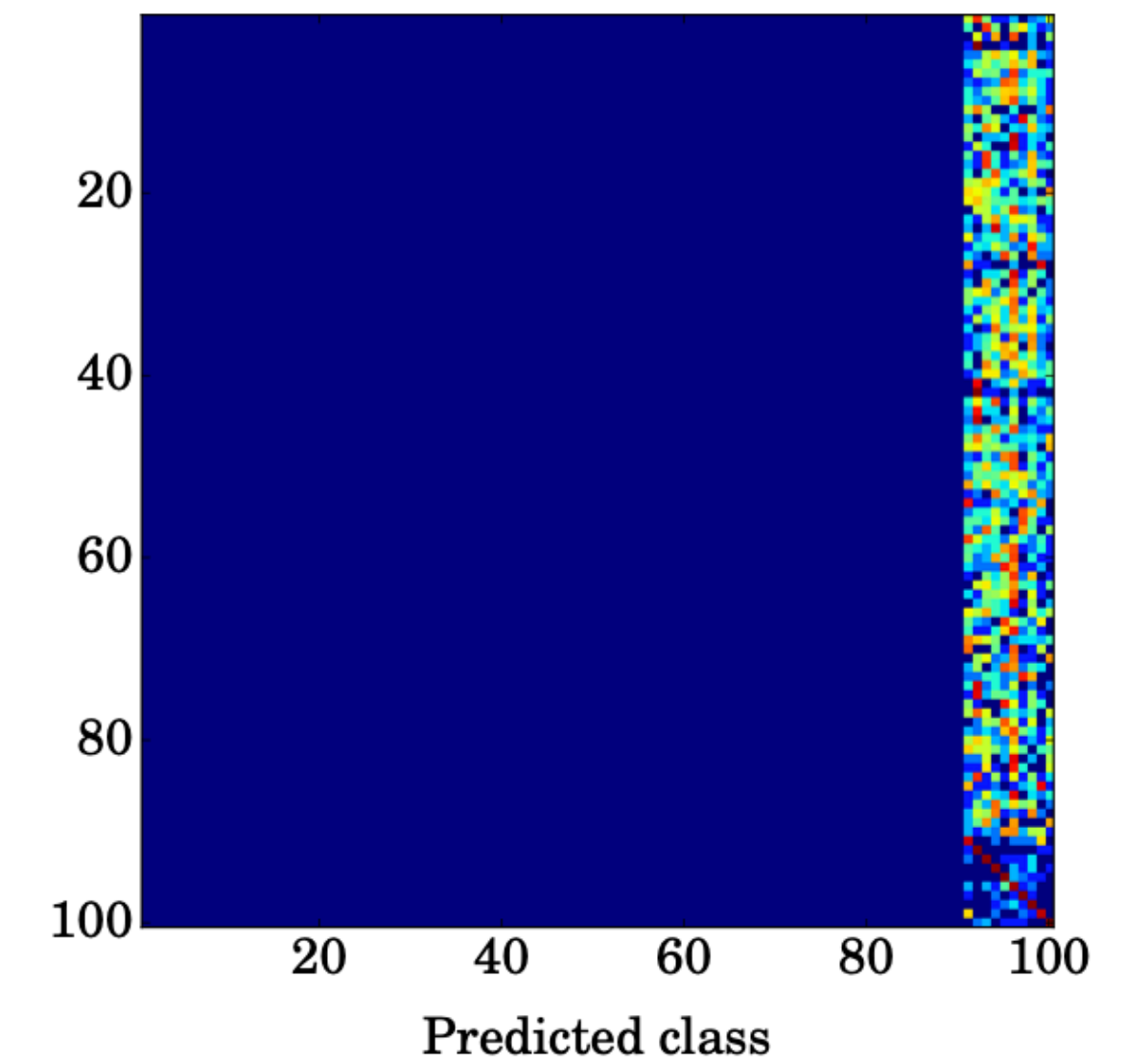
(a) iCaRL



(b) LwFMC

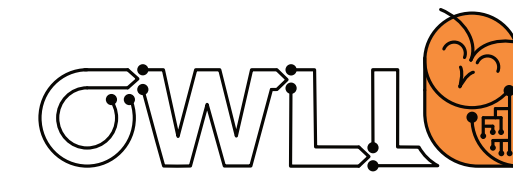


(c) fixed representation

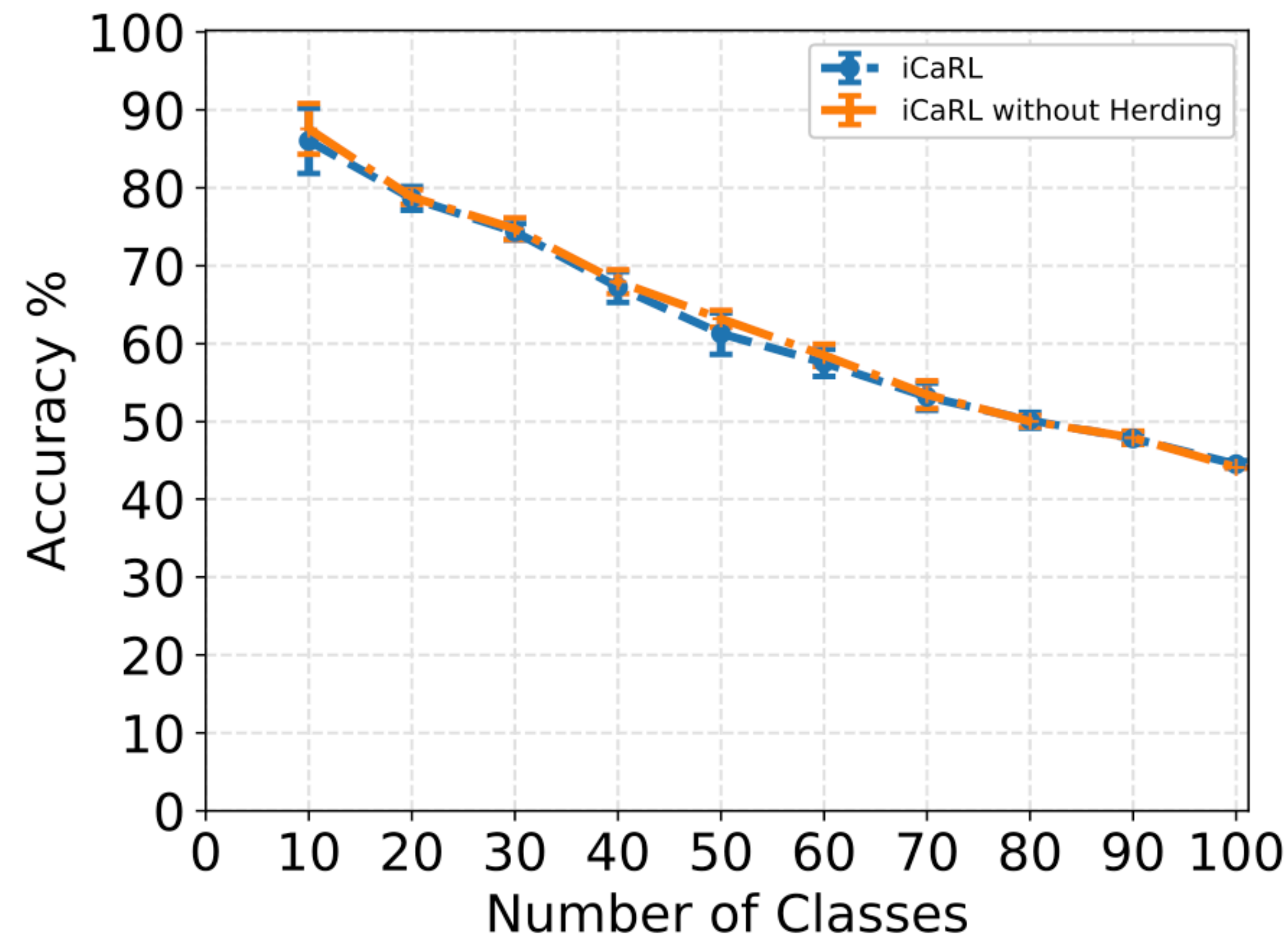


(d) finetuning

Role of extraction algorithm



Does the herding selection algorithm outperform random selection?



Role of memory budget

How expected are our observations?

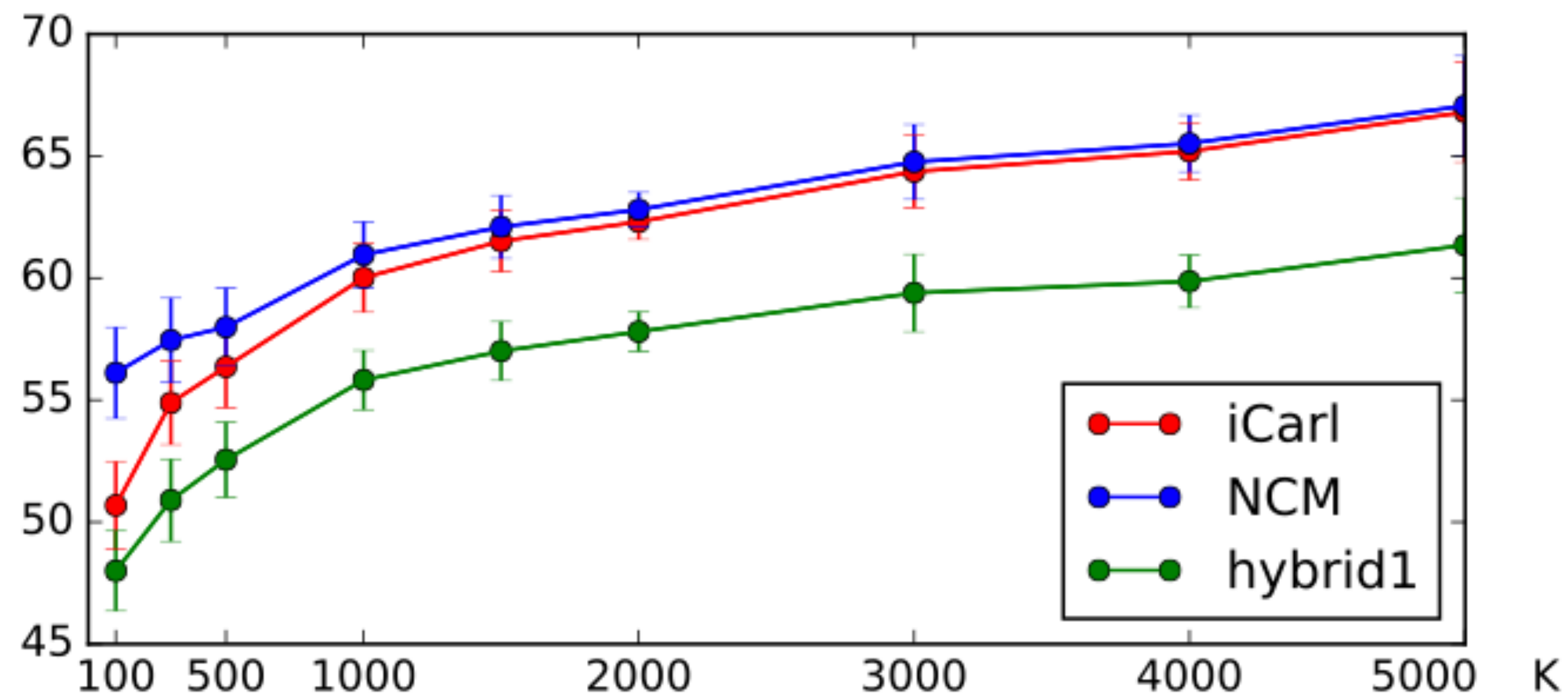
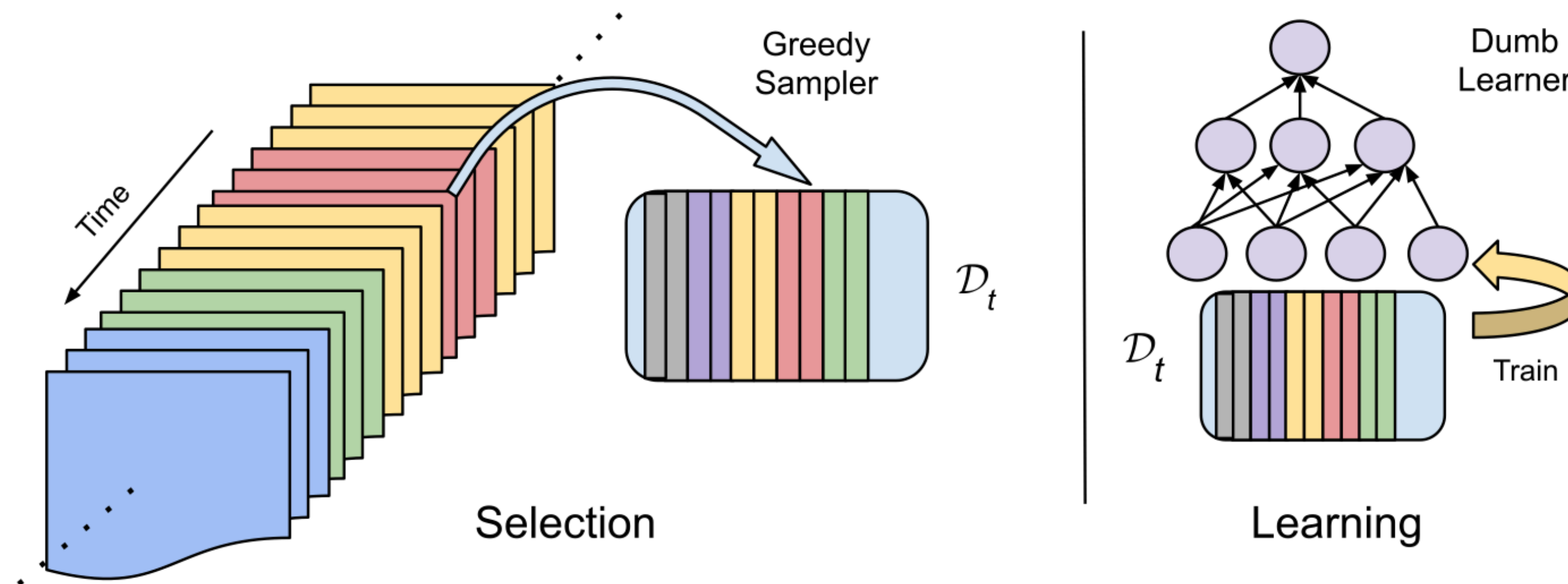


Figure 4: Average incremental accuracy on iCIFAR-100 with 10 classes per batch for different memory budgets K .

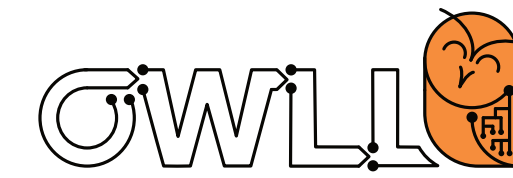
Role of memory

Is it really just the data subset that we retain?

A “dumb learner” comparison suggests that we may get similar performance if we just train on the exemplar subset



Role of memory

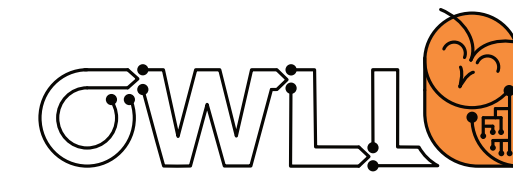


Are memory buffers of real examples the solution to continual learning?

“While it is an effective method in ANNs, rehearsal is unlikely to be a realistic model of biological learning mechanisms, as in this context the actual old information (accurate and complete representation of all items ever learned by the organism) is not available. Pseudorehearsal is significantly more likely to be a mechanism which could actually be employed by organisms as it does not require access to this old information, it just requires a way of approximating it.”

R. French, “Pseudo-recurrent Connectionist Networks:
An Approach to the Sensitivity-Plasticity Dilemma”,
Connection Science 9:4, 1997

(Pseudo-)rehearsal



A rehearsal alternative

“Pseudorehearsal is based on the use in the rehearsal process of artificially constructed populations of “pseudoitems” instead of the “actual” previously learned items. A pseudoitem is constructed by generating a new input vector (setting at random 50% of input elements to 0 and 50% to 1 as usual), and passing it forward through the network in the standard way. Whatever output vector this input generates becomes the associated target output”

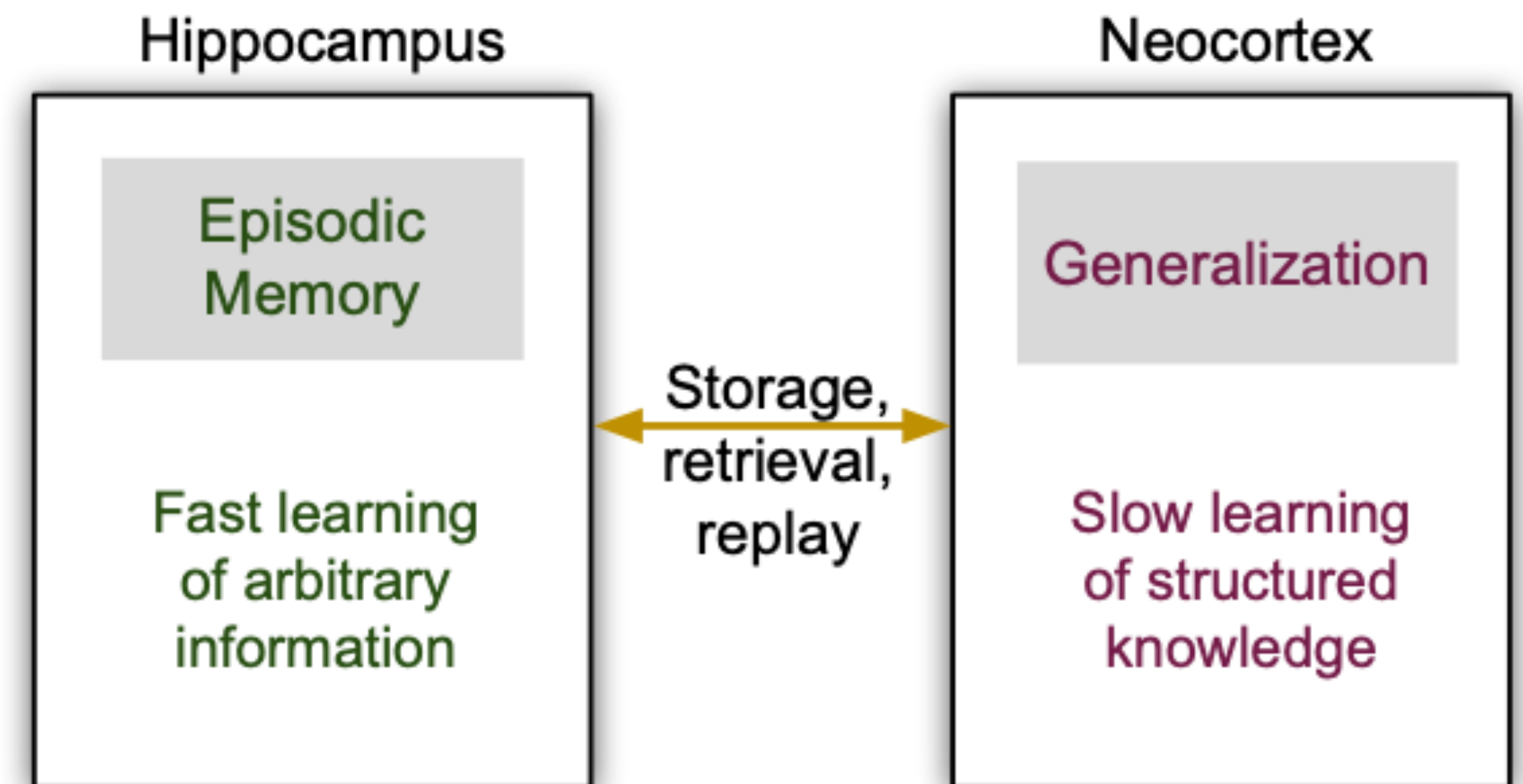
A. Robins, “Catastrophic forgetting, rehearsal and pseudorehearsal”,
Journal of Neural Computing 7: 123-146, 1995

Pseudo-rehearsal

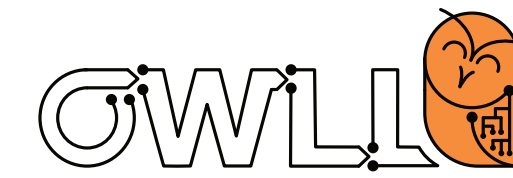
Complementary learning systems theory

(McClelland et al, Psychological Review 102:3, 1995), *simplified picture here*

- Hippocampus: short-term adaptation & rapid learning of novel information
- Neocortex: slow learning, to consolidate & build up overlapping representations
- Hippocampus “plays back” over time to neocortex

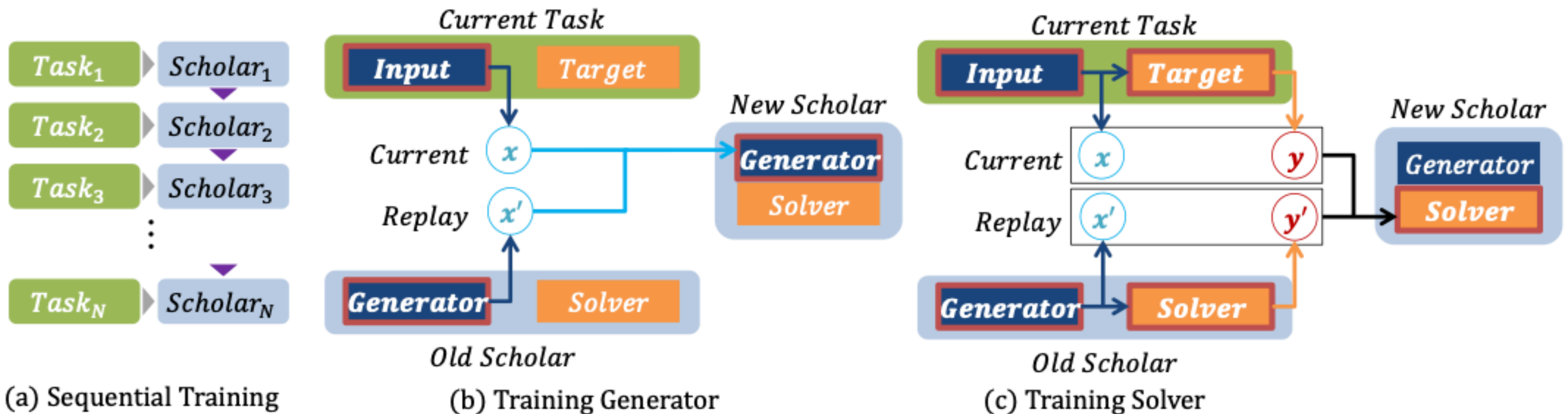


Deep Generative Replay



We could train two machine learning models:

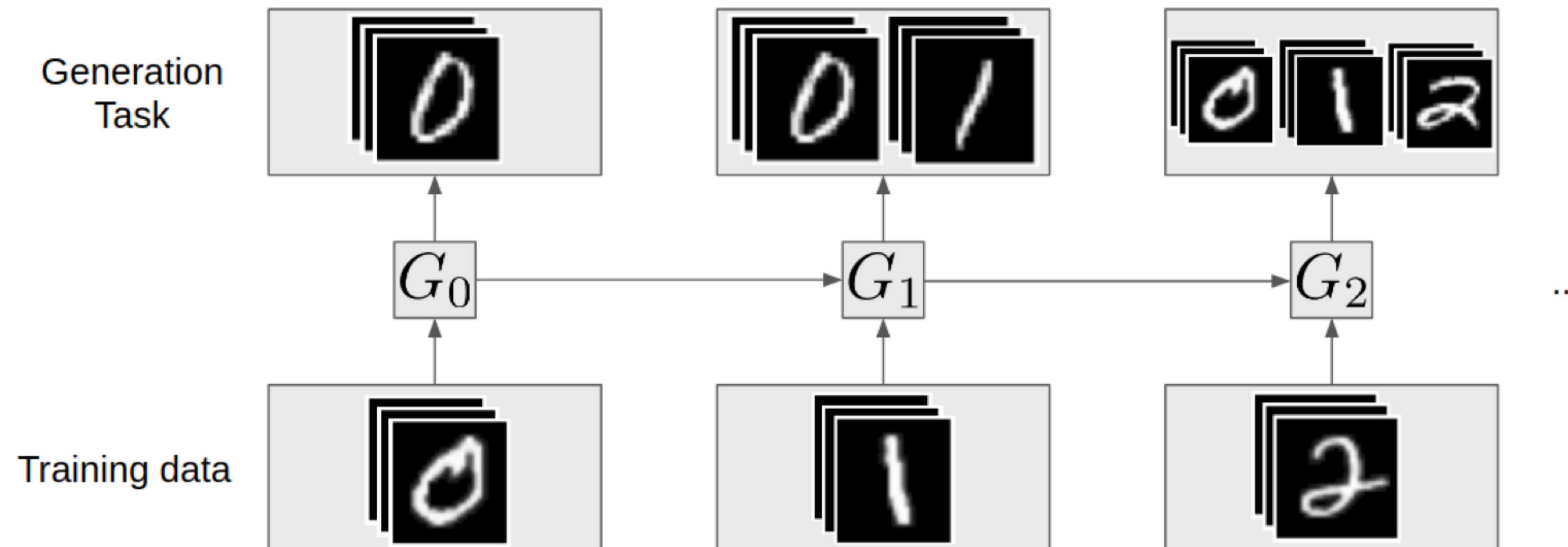
a “generator” (there are many types) + “task solver” -> then alternate training



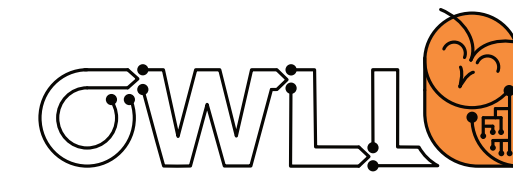
Deep Generative Replay

We could train two machine learning models:

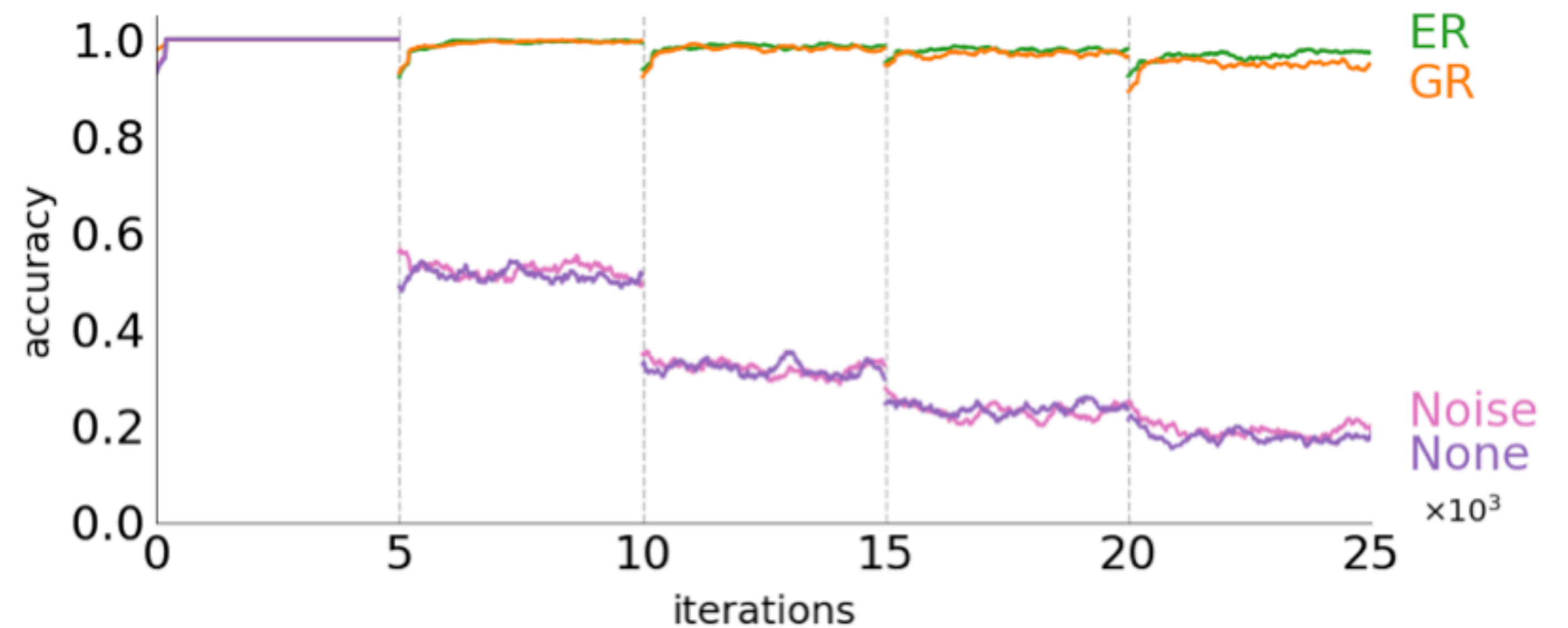
a “generator” (there are many types) + “task solver” -> then alternate training



Exemplar/Generative Rehearsal



- Exemplar Rehearsal (ER) and Generative (Pseudo-)Rehearsal (GR) can work equally well if we have a powerful generator
- In contrast, randomly rehearsed sampled noise patterns will no longer work on complex tasks





A short interlude: generative models

Generative models



- A **discriminative** model typically learns something like $p(y|x)$
- A **generative** model also learns about the data distribution $p(x)$ & the process by which data is created (the generative factors)
- Having a generative model **does not mean we cannot also solve discriminative tasks**

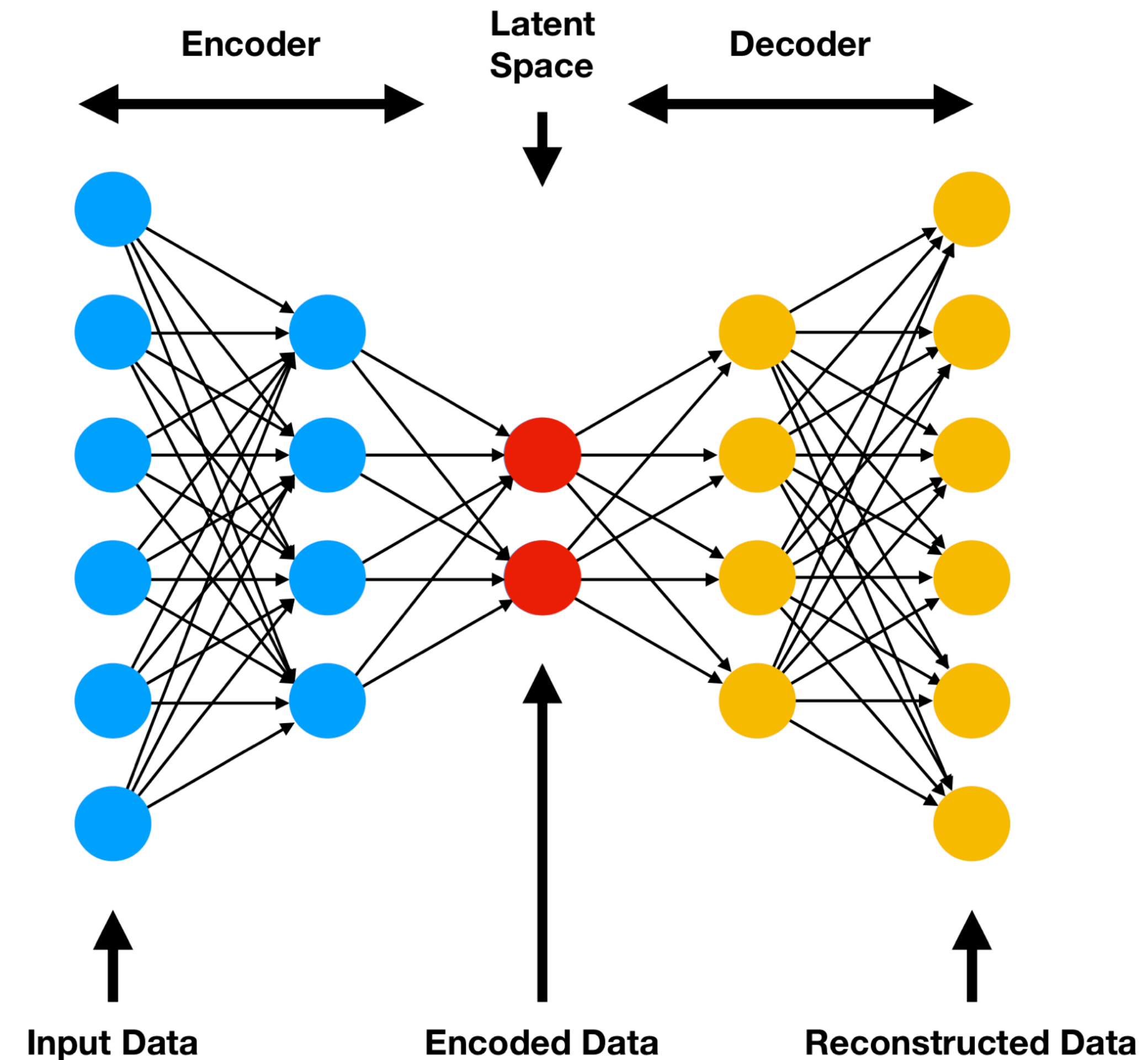
$$p(x,y) = p(y|x)p(x)$$

Autoencoders

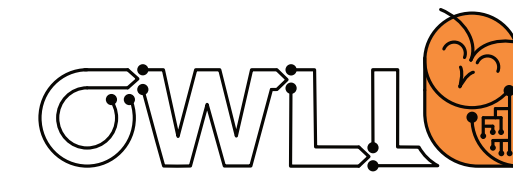
Let's take a look at **autoencoders**

Why? To later see that we don't necessarily need two separate models:

- Learn a representation, an “encoding” of the data
- An **encoder** maps to a “code”: “**latent variables**”
- A **decoder** learns to **reconstruct** the input

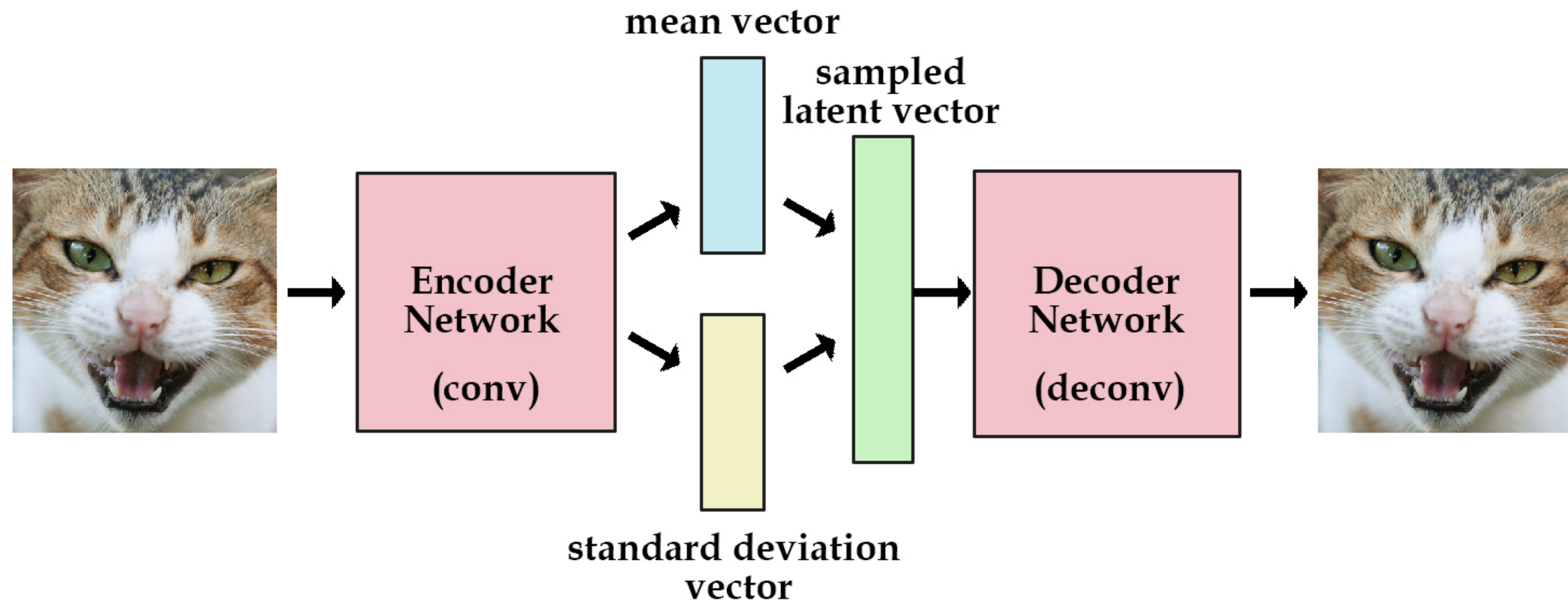


Variational Autoencoders

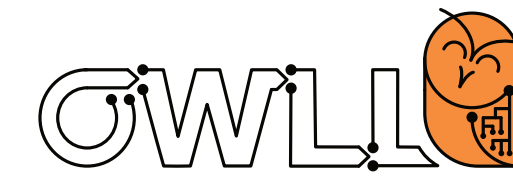


The autoencoders latent embedding/variables may be difficult to grasp if unconstrained

But we could constrain the latent space to follow a specific distribution, e.g. a **Variational Autoencoder**: “Auto-encoding Variational Bayes”, Kingma & Welling, ICLR 2014



Variational Autoencoders



More formally, let us consider:

- A dataset with variable \mathbf{x}
- Data is generated by a random process involving unobserved random variable \mathbf{z}
- \mathbf{z} is generated from some prior distribution $p_{\theta}(\mathbf{z})$
- A value x is generated from some conditional distribution $p_{\theta}(x | \mathbf{z})$

Variational Autoencoders



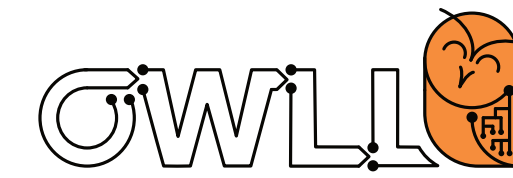
More formally, let us consider:

- A dataset with variable \mathbf{x}
- Data is generated by a random process involving unobserved random variable \mathbf{z}
- \mathbf{z} is generated from some prior distribution $p_{\theta}(\mathbf{z})$
- A value x is generated from some conditional distribution $p_{\theta}(x | \mathbf{z})$

But, the parameters and values of latent variables \mathbf{z} are not known to us.

$$p_{\theta}(x) = \int p_{\theta}(x, z) dz \text{ is intractable}$$

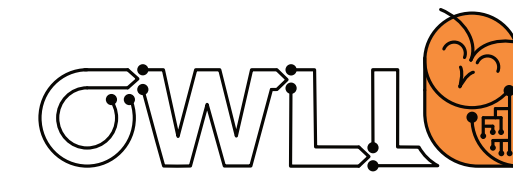
VAE: ELBO derivation



1. The densities of the marginal and joint distribution are related through **Bayes rule**:

$$p_{\theta}(z | x) = \frac{p_{\theta}(x, z)}{p_{\theta}(x)} = \frac{p_{\theta}(x | z)p_{\theta}(z)}{p_{\theta}(x)}$$

VAE: ELBO derivation



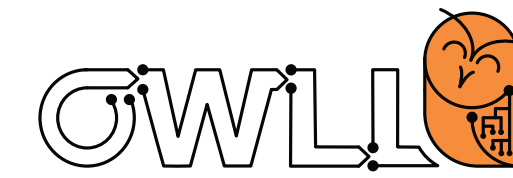
1. The densities of the marginal and joint distribution are related through **Bayes rule**:

$$p_{\theta}(z | x) = \frac{p_{\theta}(x, z)}{p_{\theta}(x)} = \frac{p_{\theta}(x | z)p_{\theta}(z)}{p_{\theta}(x)}$$

2. Make use of the **logarithm** on both sides to write as a sum:

$$\log p_{\theta}(x) = \log p_{\theta}(x | z) + \log p_{\theta}(z) - \log p_{\theta}(z | x)$$

VAE: ELBO derivation



1. The densities of the marginal and joint distribution are related through **Bayes rule**:

$$p_{\theta}(z | x) = \frac{p_{\theta}(x, z)}{p_{\theta}(x)} = \frac{p_{\theta}(x | z)p_{\theta}(z)}{p_{\theta}(x)}$$

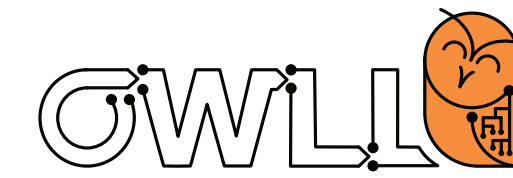
2. Make use of the **logarithm** on both sides to write as a sum:

$$\log p_{\theta}(x) = \log p_{\theta}(x | z) + \log p_{\theta}(z) - \log p_{\theta}(z | x)$$

3. We do not know our real posterior $p_{\theta}(z | x)$. We will make use of **variational inference** and introduce an **approximation** $q_{\phi}(z | x)$:

$$\log p_{\theta}(x) = \int q_{\phi}(z | x) [\log p_{\theta}(x | z) + \log p_{\theta}(z) - \log p_{\theta}(z | x)] dz$$

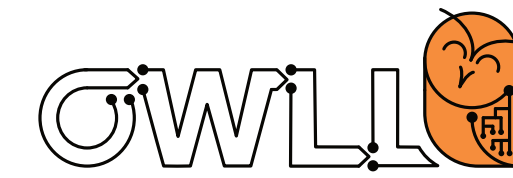
VAE: ELBO derivation



4. We add and subtract $q_\phi(z|x)$:

$$\log p_\theta(x) = \int q_\phi(z|x) \left[\log p_\theta(x|z) + \log p_\theta(z) - \log p_\theta(z|x) + \log q_\phi(z|x) - \log q_\phi(z|x) \right] dz$$

VAE: ELBO derivation



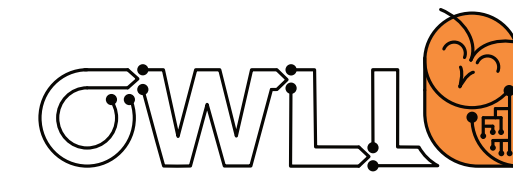
4. We add and subtract $q_\theta(z|x)$:

$$\log p_\theta(x) = \int q_\phi(z|x) \left[\log p_\theta(x|z) + \log p_\theta(z) - \log p_\theta(z|x) + \log q_\phi(z|x) - \log q_\phi(z|x) \right] dz$$

5. We make use of the (reverse) **Kullback Leibler divergence** between distributions:

$$KL(Q||P) = \int_{-\infty}^{\infty} Q(x) \log \frac{Q(x)}{P(x)}$$

VAE: ELBO derivation



4. We add and subtract $q_\theta(z|x)$:

$$\log p_\theta(x) = \int q_\phi(z|x) \left[\log p_\theta(x|z) + \log p_\theta(z) - \log p_\theta(z|x) + \log q_\phi(z|x) - \log q_\phi(z|x) \right] dz$$

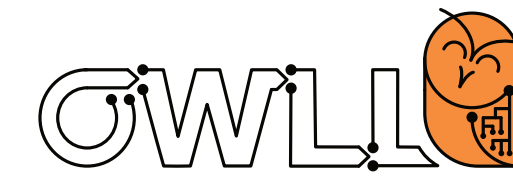
5. We make use of the (reverse) **Kullback Leibler divergence** between distributions:

$$KL(Q||P) = \int_{-\infty}^{\infty} Q(x) \log \frac{Q(x)}{P(x)}$$

and replace the integral with an **expectation over samples**:

$$\log p_\theta(x) = KL [q_\theta(z|x) || p_\theta(z|x)] + \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] - KL [q_\phi(z|x) || p_\theta(z)]$$

VAE: ELBO derivation

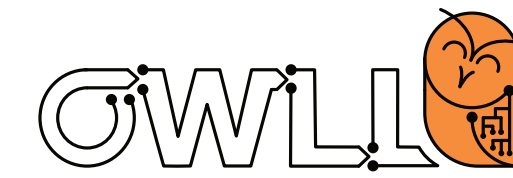


$$\log p_{\theta}(x) = KL [q_{\theta}(z|x) || p_{\theta}(z|x)] + \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - KL [q_{\phi}(z|x) || p_{\theta}(z)]$$

6. We still cannot evaluate the **first term** on the right hand side, but it is **strictly positive**
-> we can optimize the remaining terms

$$\log p_{\theta}(x) = KL [q_{\theta}(z|x) || p_{\theta}(z|x)] + \mathcal{L}(\theta, \phi; x)$$

VAE: ELBO derivation



$$\log p_{\theta}(x) = KL [q_{\theta}(z|x) || p_{\theta}(z|x)] + \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - KL [q_{\phi}(z|x) || p_{\theta}(z)]$$

6. We still cannot evaluate the **first term** on the right hand side, but it is **strictly positive**
-> we can optimize the remaining terms

$$\log p_{\theta}(x) = KL [q_{\theta}(z|x) || p_{\theta}(z|x)] + \mathcal{L}(\theta, \phi; x)$$

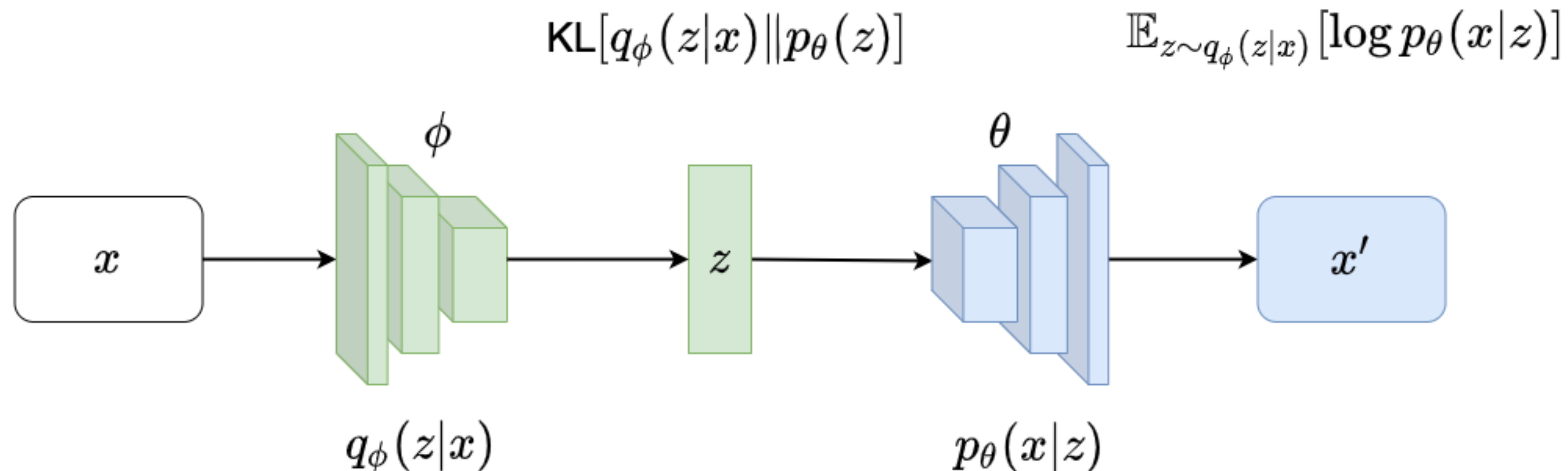
7. We now have a variational **lower-bound** to $p(x)$ (the “evidence lower bound”, or ELBO)

$$\log p_{\theta}(x) \geq \mathcal{L}(\theta, \phi; x) = \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - KL [q_{\phi}(z|x) || p_{\theta}(z)]$$

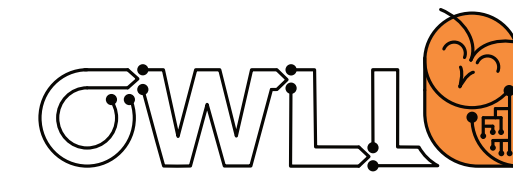
VAE: ELBO derivation

$$\mathcal{L}(\theta, \phi; x) = \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - KL [q_{\phi}(z|x) || p_{\theta}(z)]$$

- The 1. term is the expected **reconstruction error** given by the log-likelihood (with sampling)
- The 2. term is a **KL divergence** encouraging the approximate posterior to be close to a prior

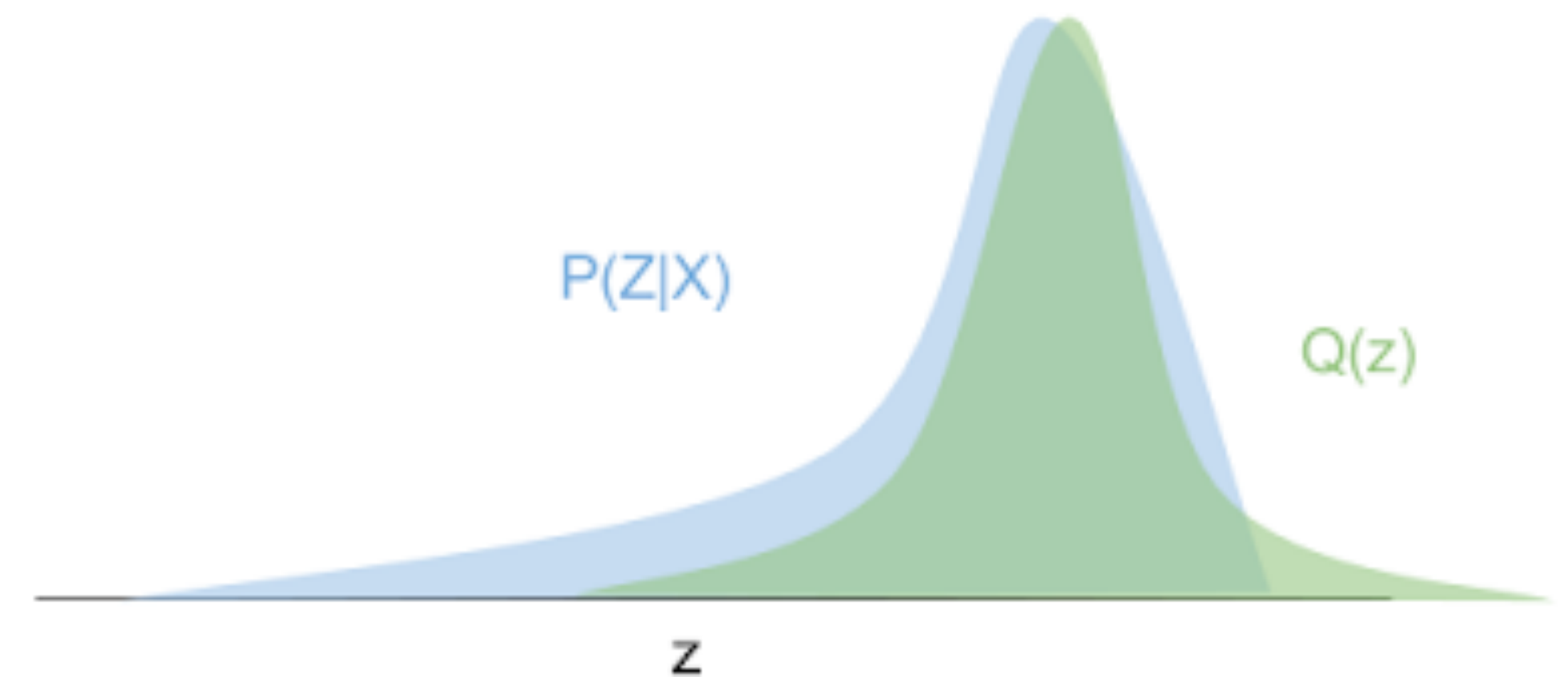


VAE: KL divergence



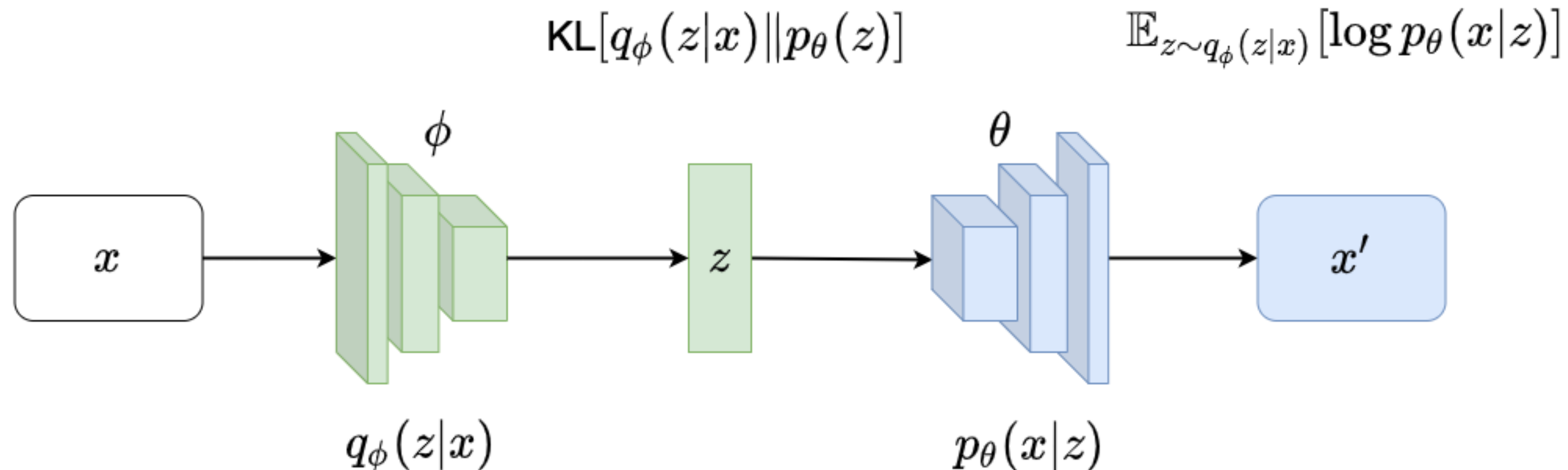
(Reverse) KL divergence measures the **amount of information** (in nats, or units of $1/\log 2$ bits) required to “**distort**” $p(z)$ into $q(z)$

If you want to read more about variational inference and the subtleties of KL divergence: <https://blog.evjang.com/2016/08/variational-bayes.html>



VAE: summary

- The recognition model, a **probabilistic encoder**, $q_\phi(z|x)$ is a variational approximation to the intractable posterior $p_\theta(z|x)$ -> given a datapoint x it produces a distribution over possible values of z from which it could have been generated
- The **probabilistic decoder**, $p_\theta(x|z)$ produces a distribution over possible values of x given z



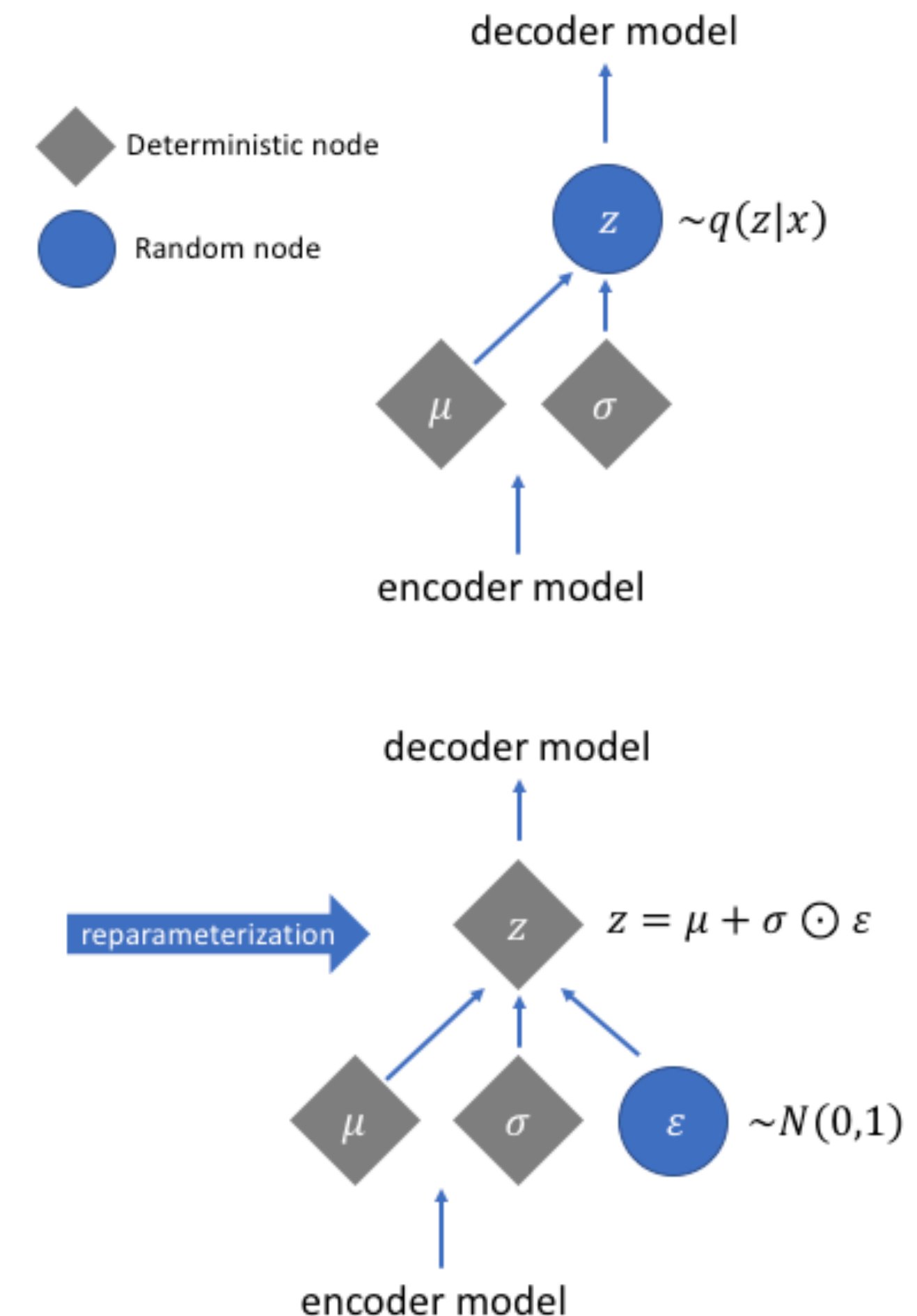
VAE: reparameterization



Example - approx. posterior as multivariate Gaussian:

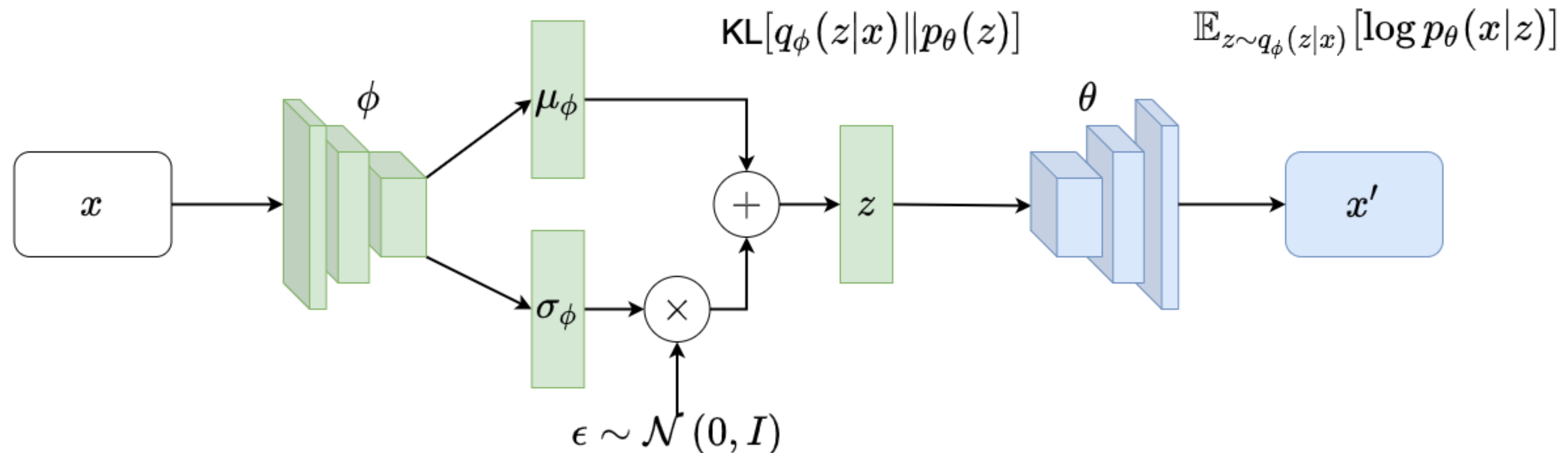
$$\log q_{\phi}(z | x) = \log \mathcal{N}(z; \mu, \sigma I)$$

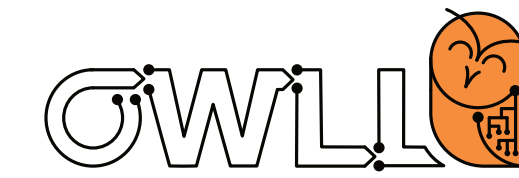
- Use **reparameterization trick** to generate samples from $q_{\phi}(z | x)$ with z as deterministic variable
- Sample posterior $z \sim q_{\phi}(z | x)$ using $z = \mu + \sigma \circ \varepsilon$ with **samples** $\varepsilon \sim \mathcal{N}(0, I)$ (as one choice of prior)



VAE: reparameterization

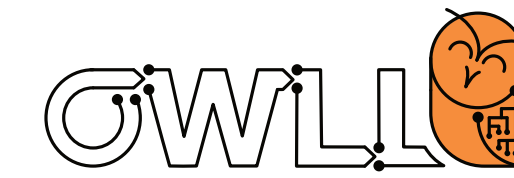
- Use **reparameterization** to generate samples from $q_\phi(z|x)$ with z as deterministic variable
- Sample posterior $z \sim q_\phi(z|x)$ using $z = \mu + \sigma \circ \varepsilon$ with **samples** $\varepsilon \sim \mathcal{N}(0, I)$ (as one choice)





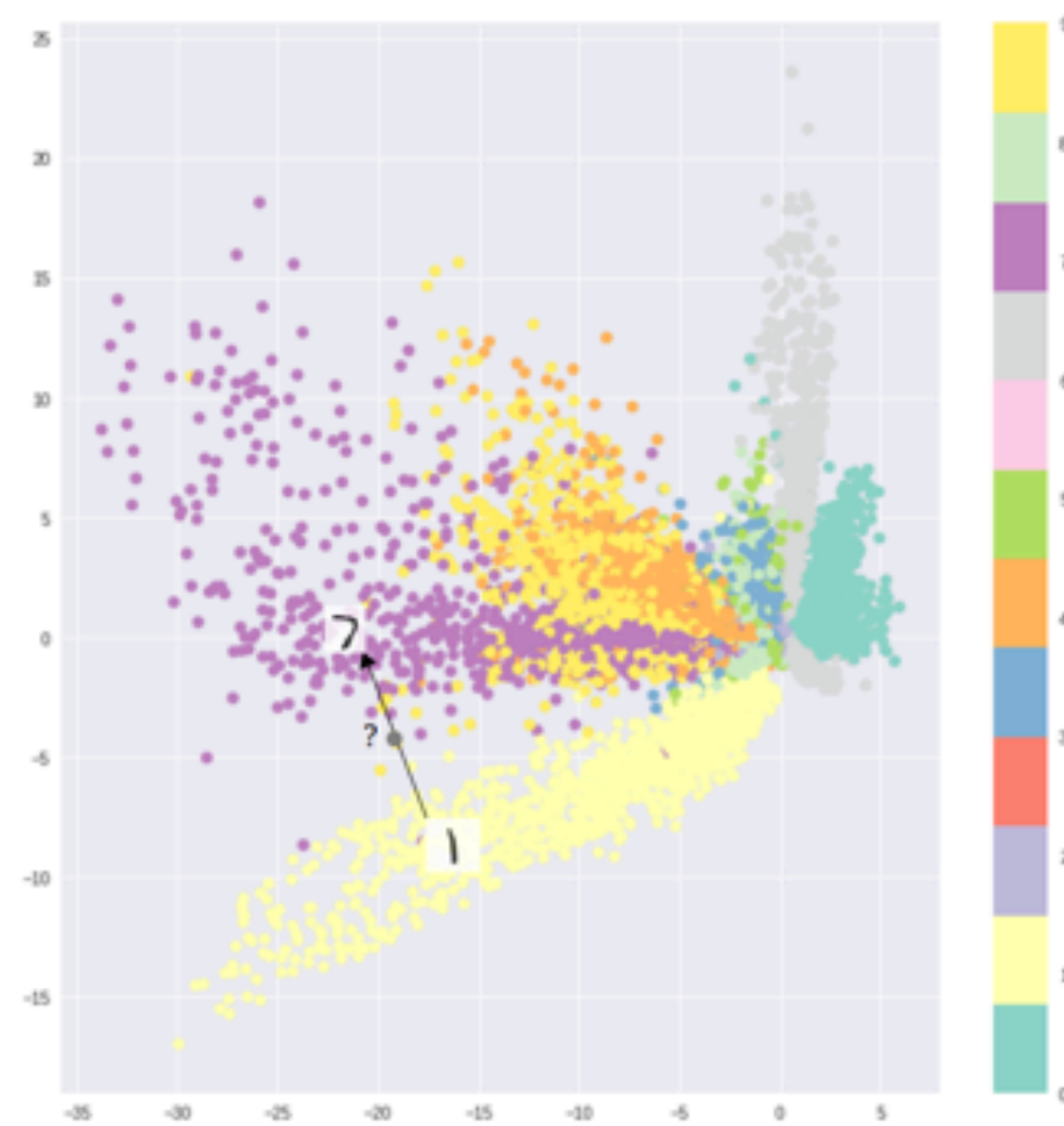
Why was this detour important: back to continual learning

What have we gained?

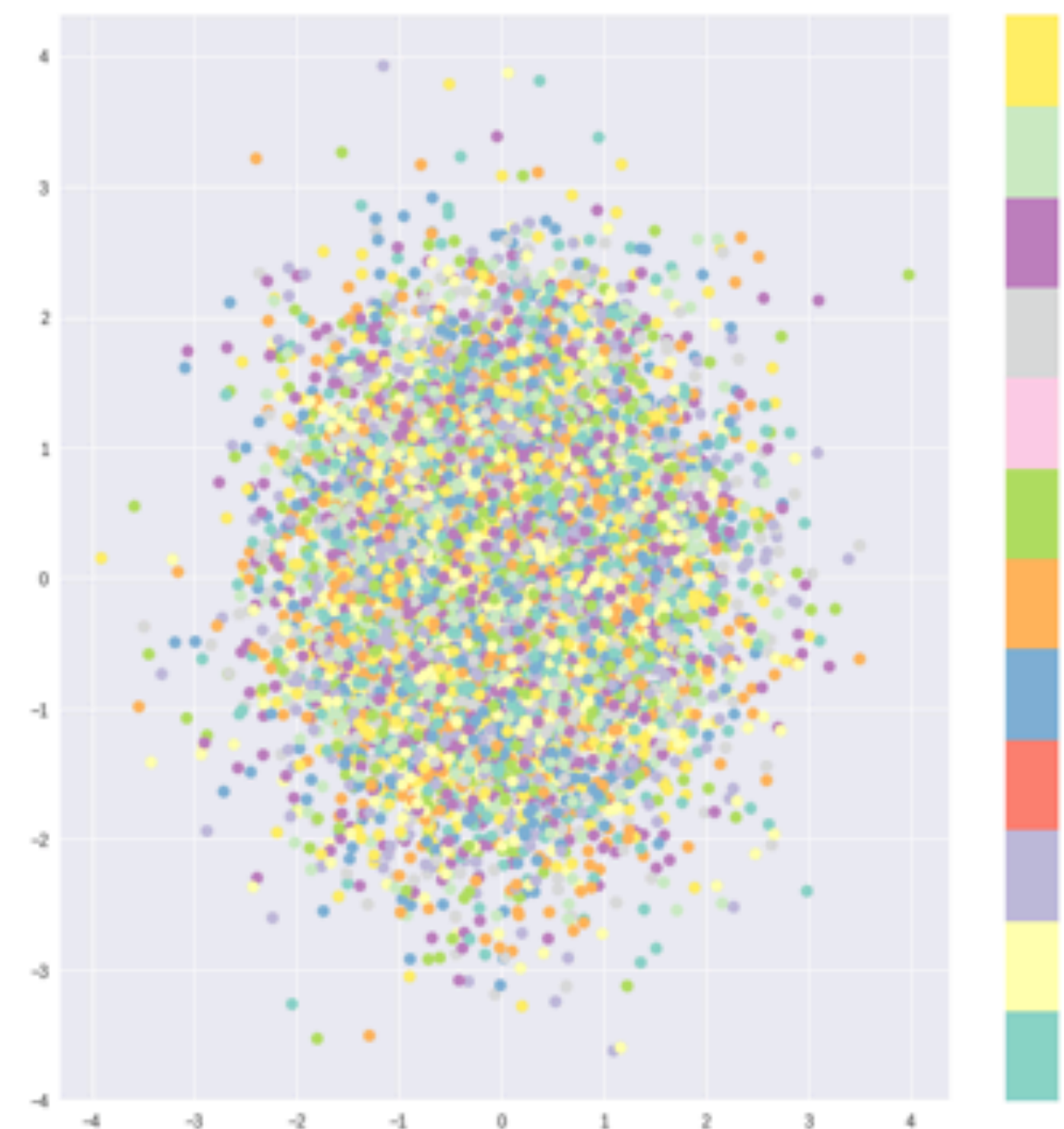


Why did we go through all this math? And why is it important for continual learning?

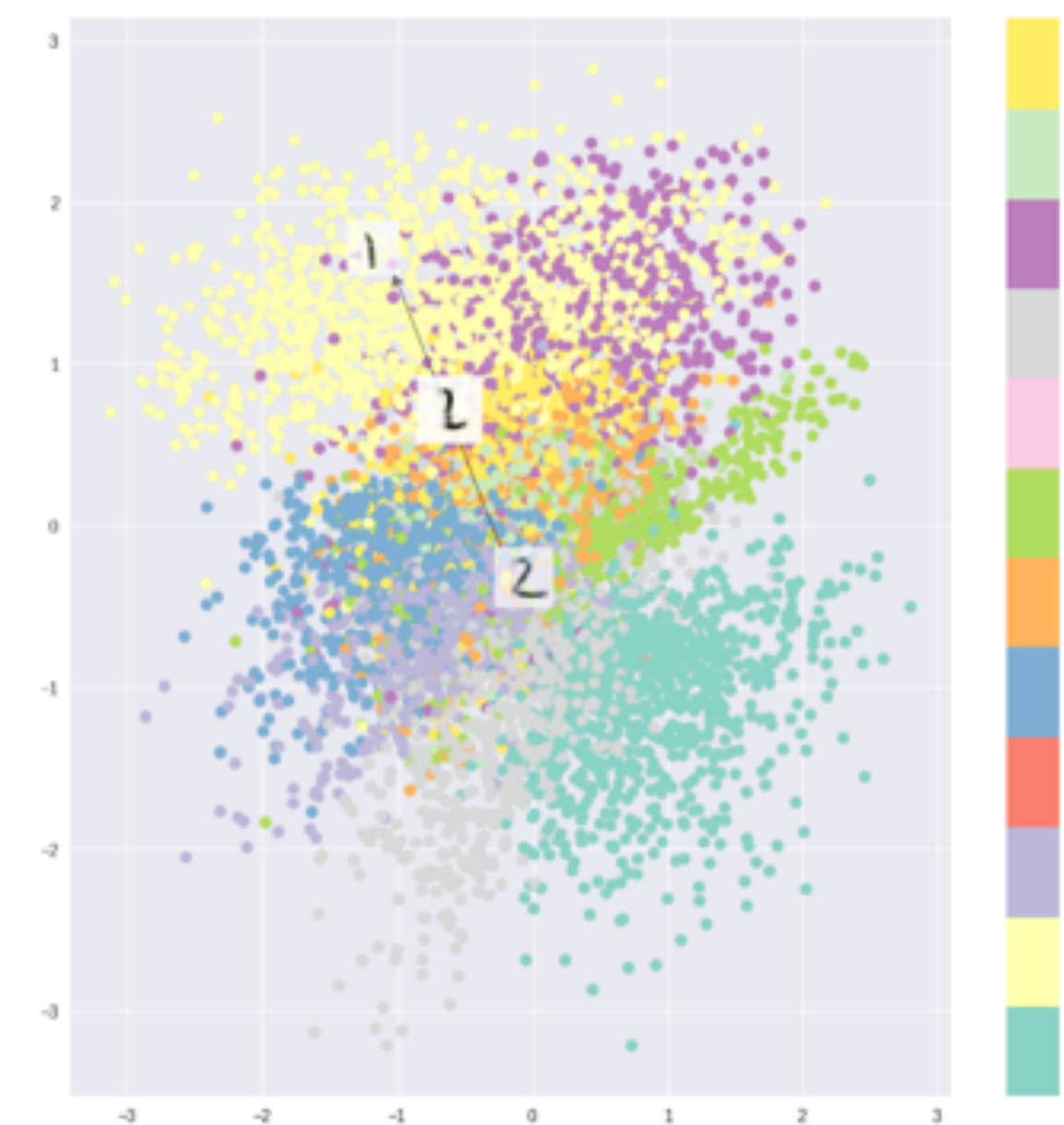
Only reconstruction loss



Only KL divergence



Combination

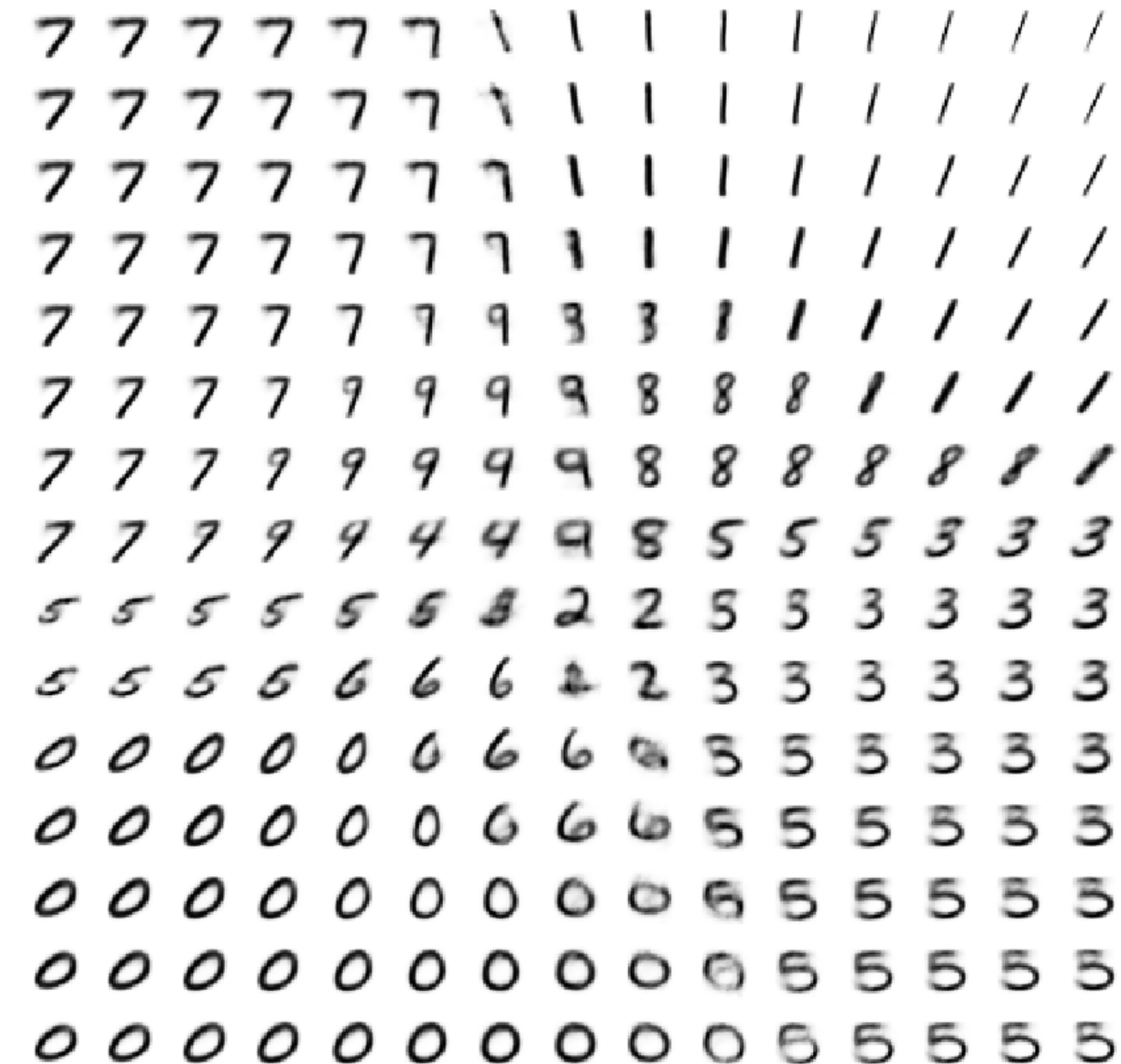


What have we gained?

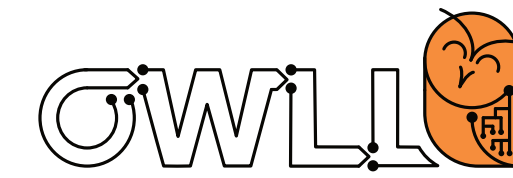


Why did we go through all this math? And why is it important for continual learning?

- We can sample from a trained model: $z \sim p(z)$, here $\mathcal{N}(0, I)$, and generate (decode) x
- We also have the approximate posterior that we could regularize in continual learning



Variational Continual Learning

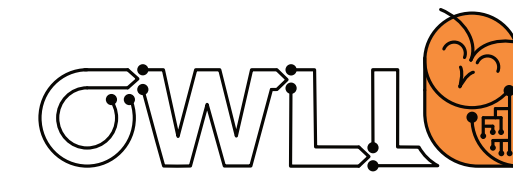


$$\mathcal{L}(\theta, \phi; x) = \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - KL [q_{\phi}(z|x) || p_{\theta}(z)]$$

The “likelihood focused” perspective:
generative/pseudo rehearsal

- Generate old tasks’ data and concatenate it with new task data
- Primarily optimize “the likelihood” (left)

Variational Continual Learning



$$\mathcal{L}(\theta, \phi; x) = \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - KL [q_{\phi}(z|x) || p_{\theta}(z)]$$

The “likelihood focused” perspective:
generative/pseudo rehearsal

- Generate old tasks’ data and concatenate it with new task data
- Primarily optimize “the likelihood” (left)

The “prior focused” perspective:
regularization/distillation

- Only use new task data
- Use the posterior of an old task as the new task’s prior $KL [q_t(z) || q_{t-1}(z)]$

Variational Continual Learning

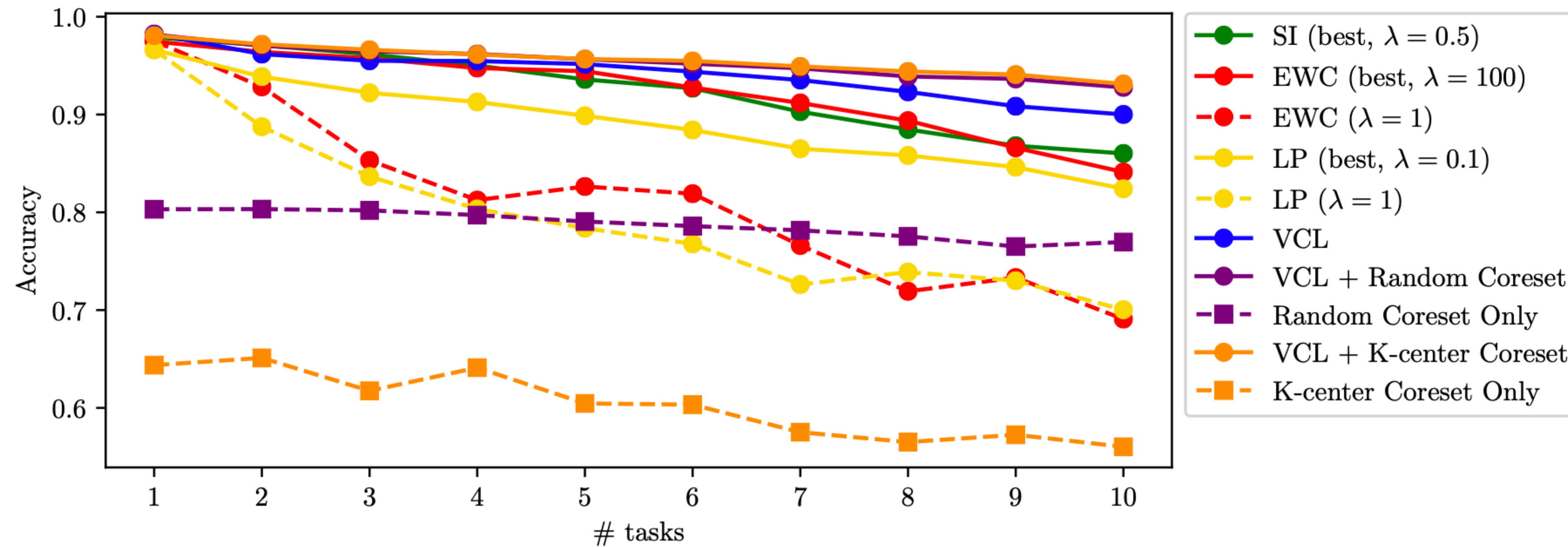
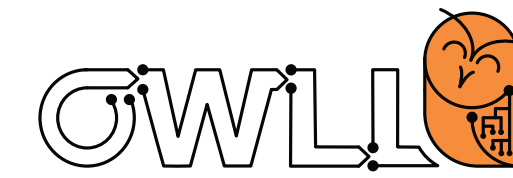
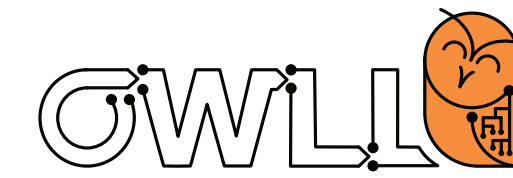


Figure 2: Average test set accuracy on all observed tasks in the Permuted MNIST experiment.

Variational Continual Learning



Optionally also store real data subsets (or a core set)

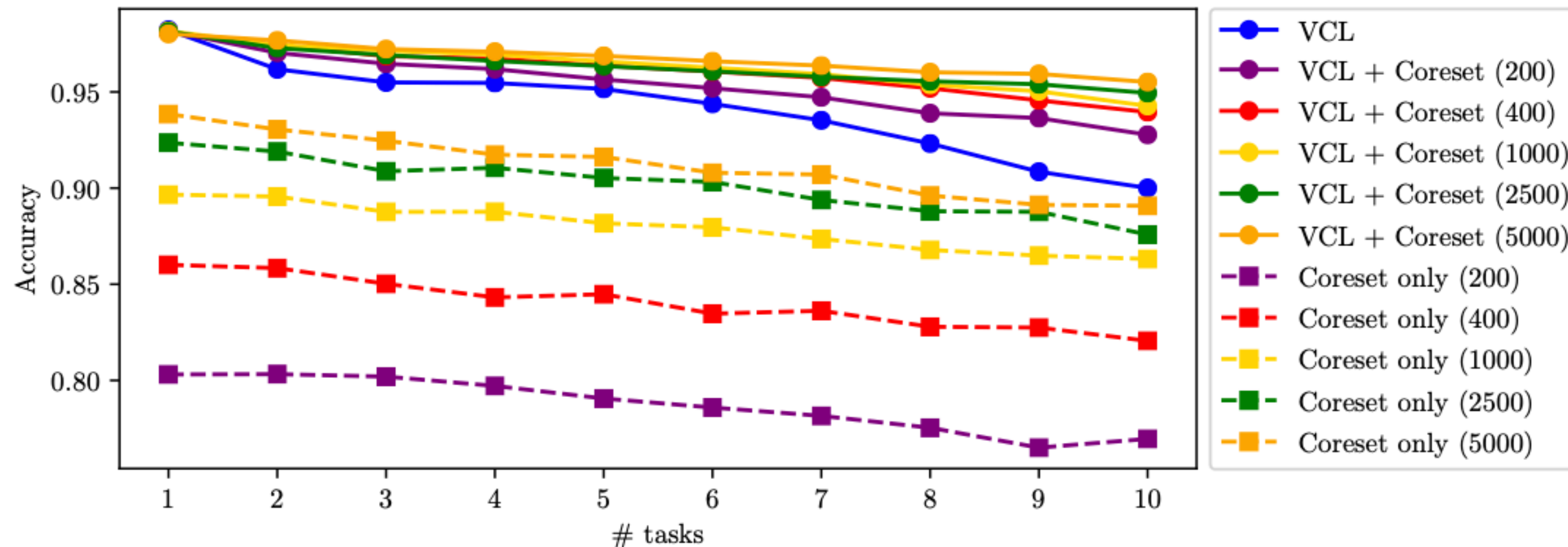
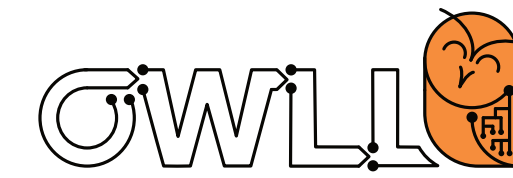


Figure 3: Comparison of the effect of coreset sizes in the Permuted MNIST experiment.

Formally: core sets



What is a core set?

The term core set is often loosely employed in modern literature to be synonymous to exemplars and sub sets of data

*“coresets are small, (weighted) summaries of large data sets such that solutions found on the summary itself are **provably competitive** with solutions found on the full data set”*

$$| \text{cost}(P, Q) - \text{cost}(C, Q) | \leq \varepsilon \cdot \text{cost}(P, Q)$$

Note how this is specific to some data, a set of questions/queries, models + loss/cost functions

Picking a “core set”

Why could we potentially pick better data subsets/exemplars now?

- Example of a 2-D latent space with 4 classes/clusters
- Random or k-means (depending on the amount of k) may not mirror the distribution well

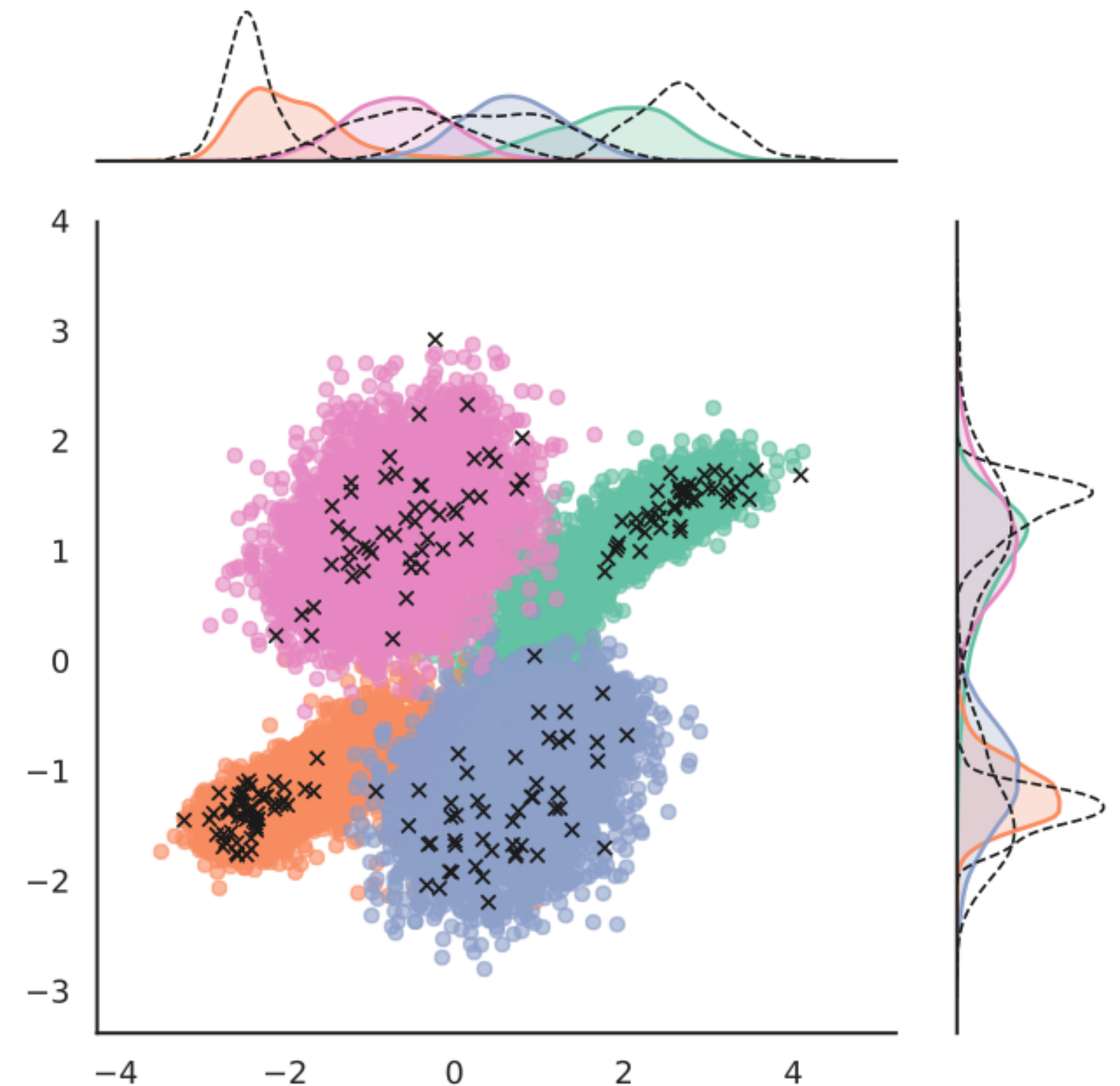


Figure from “A Wholistic View of Deep Neural Networks: Forgotten Lessons and the Bridge to Active and Open World Learning”, Mundt et al 2020

Picking a “core set”

Why could we potentially pick better data subsets/exemplars now?

- If we know our approx. posterior, we could sample and pick instances that lie close the samples
- (It's not actually that easy, for various reasons, but the intuition is that we are somewhat aware of $p(x)$ now)

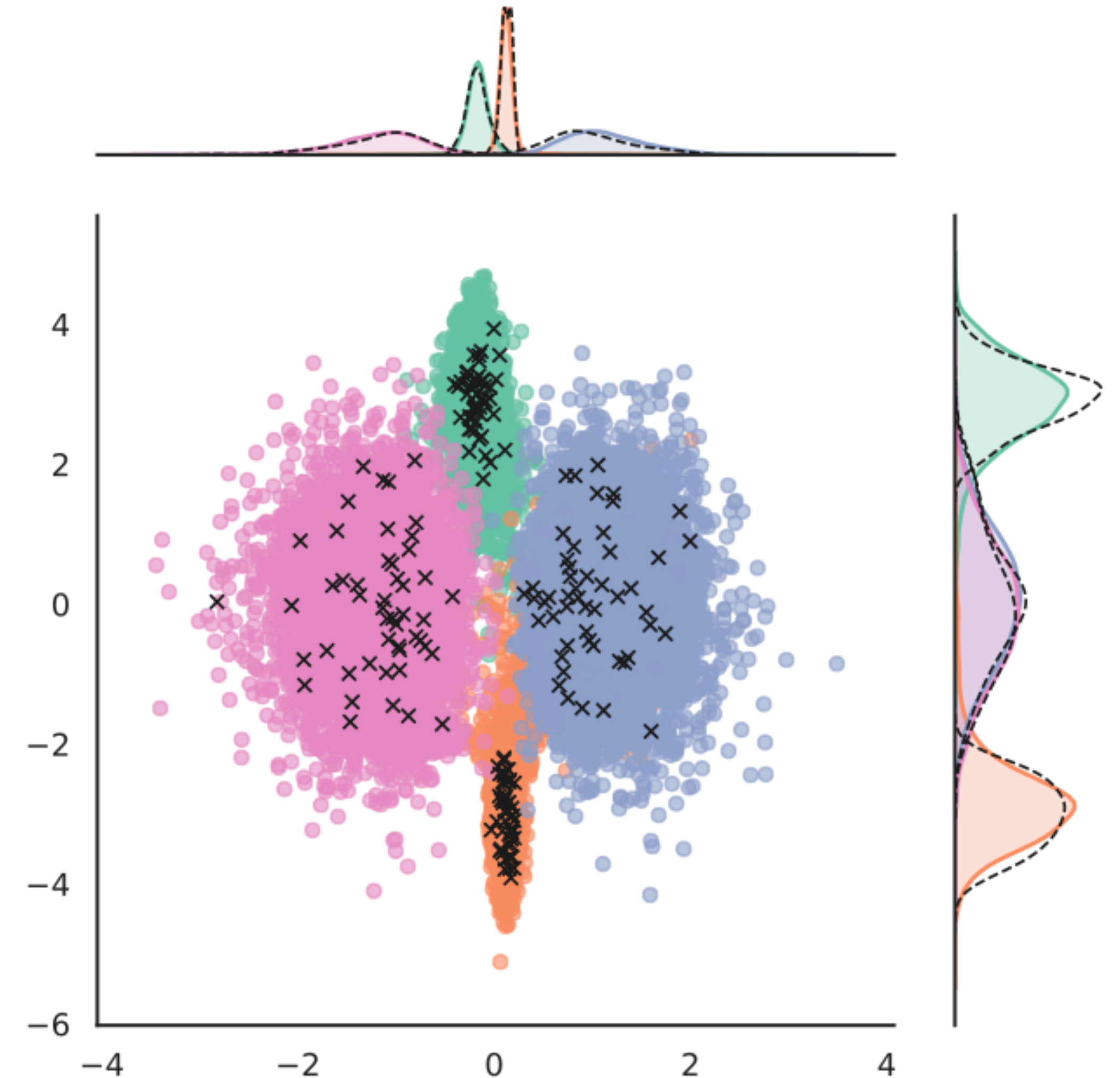
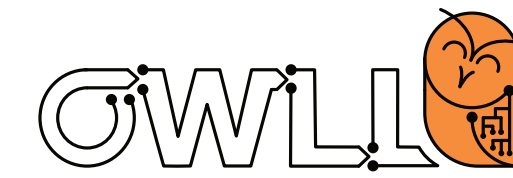


Figure from “A Wholistic View of Deep Neural Networks: Forgotten Lessons and the Bridge to Active and Open World Learning”, Mundt et al 2020

CURL: task specific Gaussians



We could now also use task-specific priors more explicitly

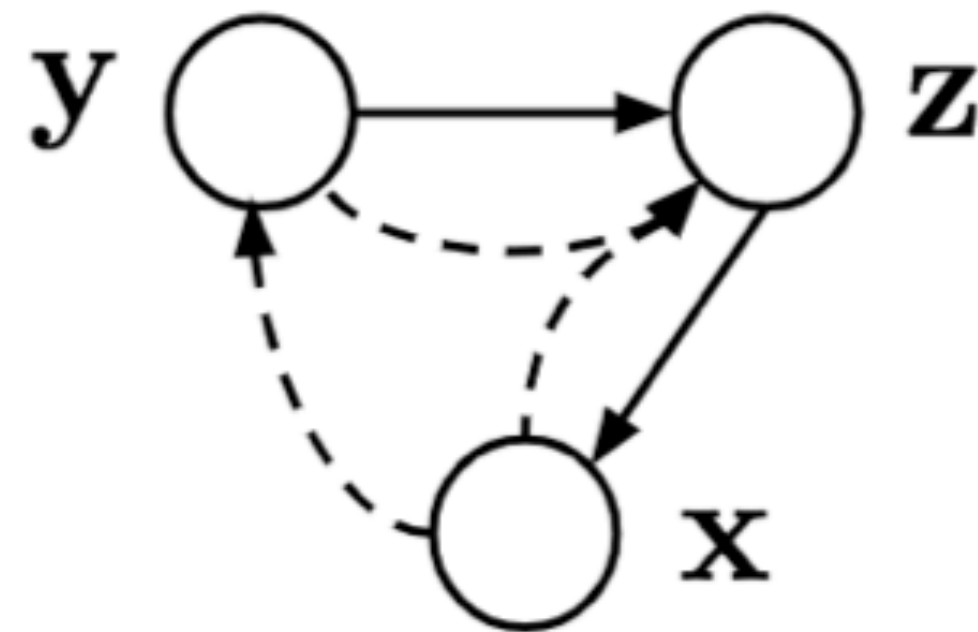


Figure 1: Graphical model for CURL. The categorical task variable y is used to instantiate a latent mixture-of-Gaussians z , which is then decoded to x .

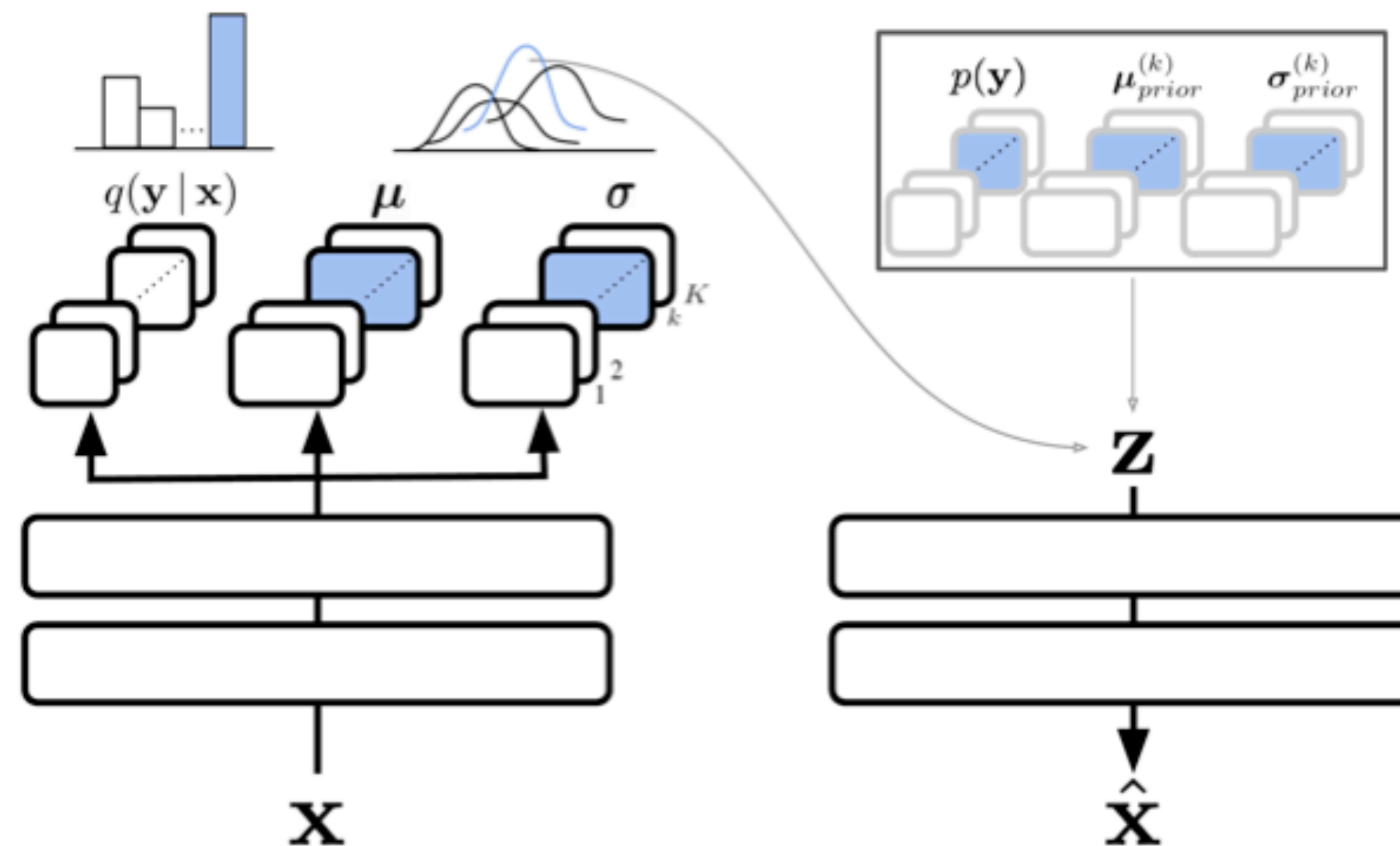
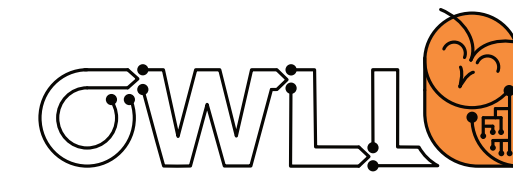


Figure 2: Diagram of the proposed approach, showing the inference procedure and architectural components used.



Have we solved forgetting? Why are we not done?

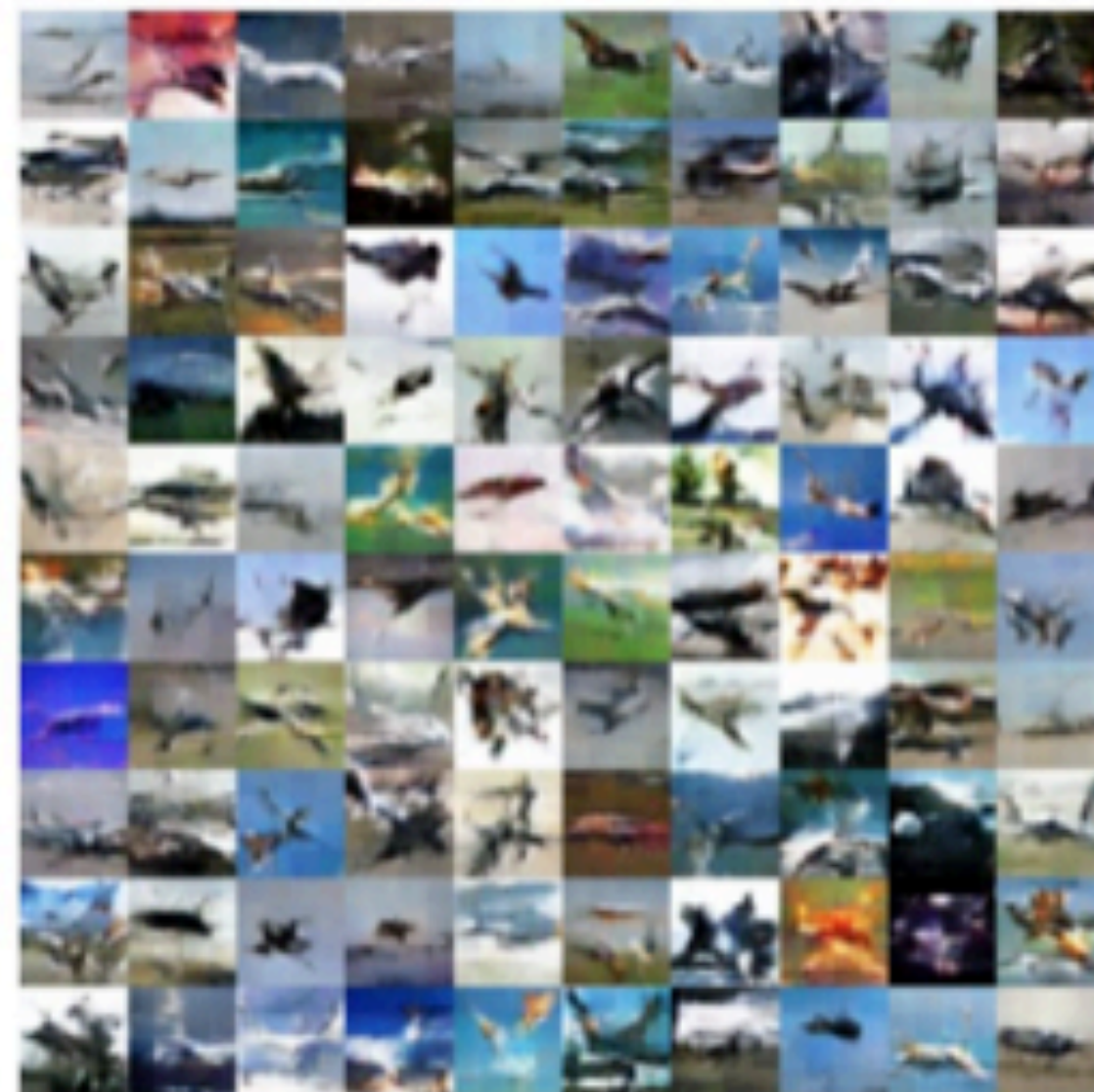
Combining ideas



Why may we need these multiple perspectives?

Generators also suffer from forgetting, errors can “snowball” rapidly

Task 1



Task 2



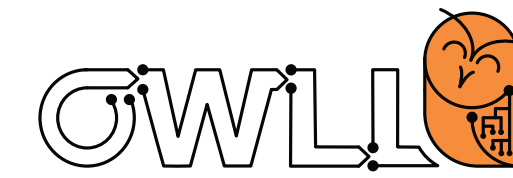
Task 3



Task 4



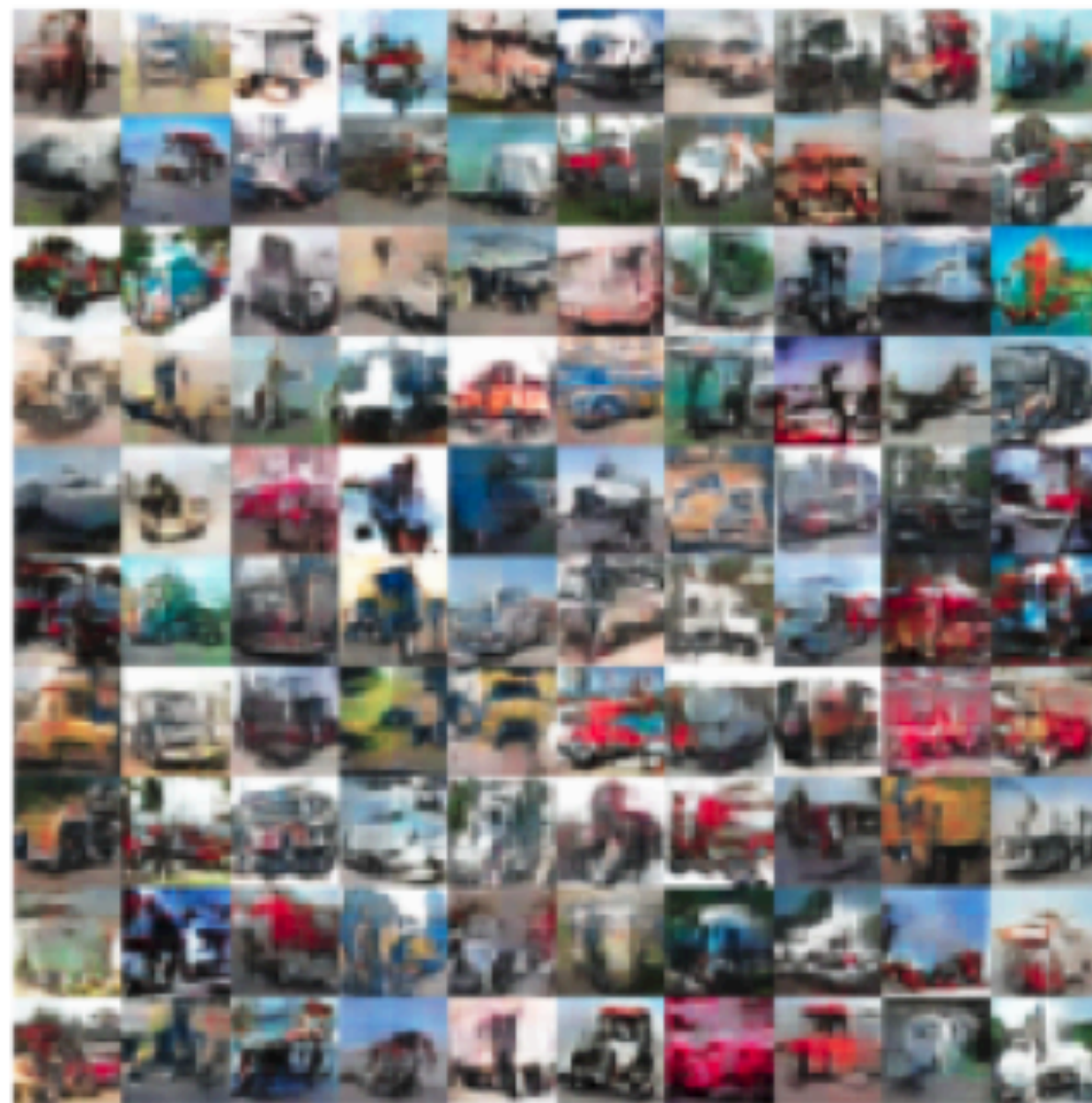
Combining ideas



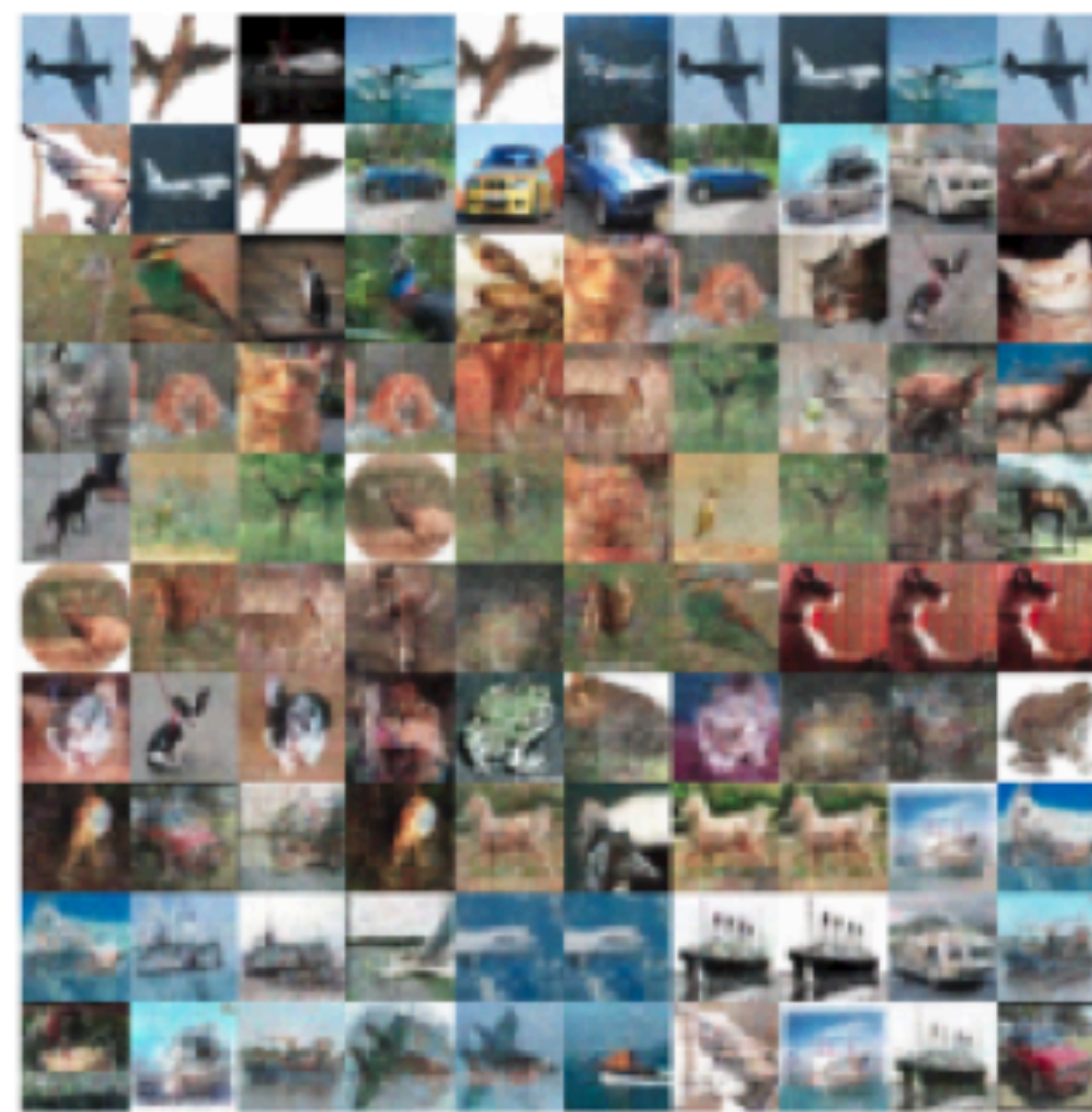
Why may we need these multiple perspectives?

Fine-tuning alone forgets, regularization is a trade-off and data rehearsal can overfit

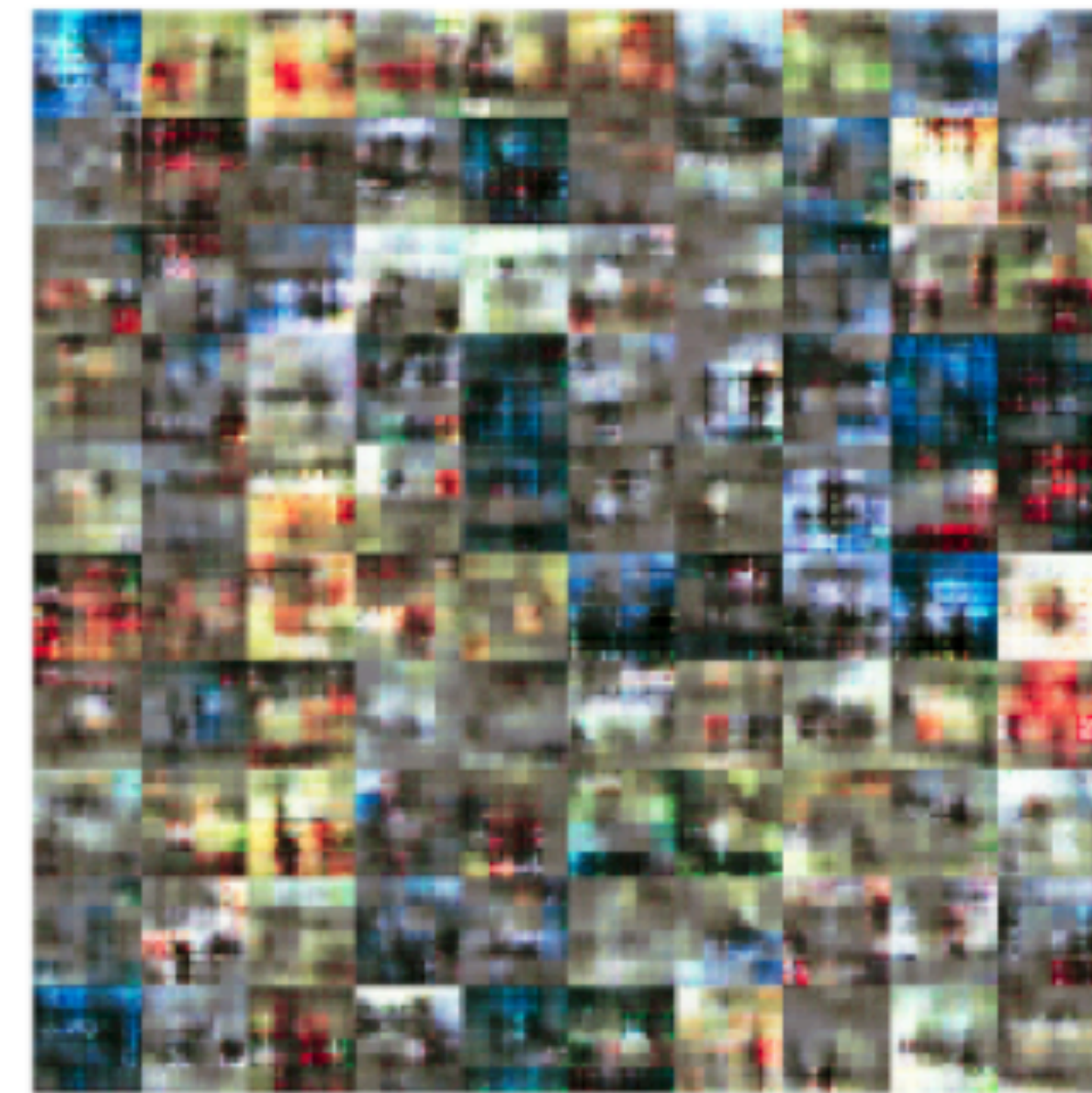
Fine-tuning



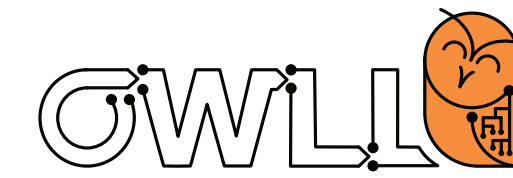
Rehearsal



Generative Replay

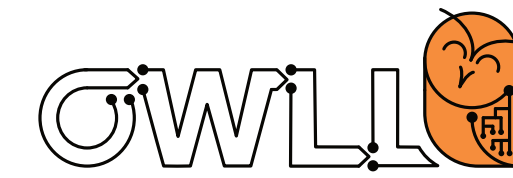


In summary: what do we want?



What could our expectations be, what might we desire?

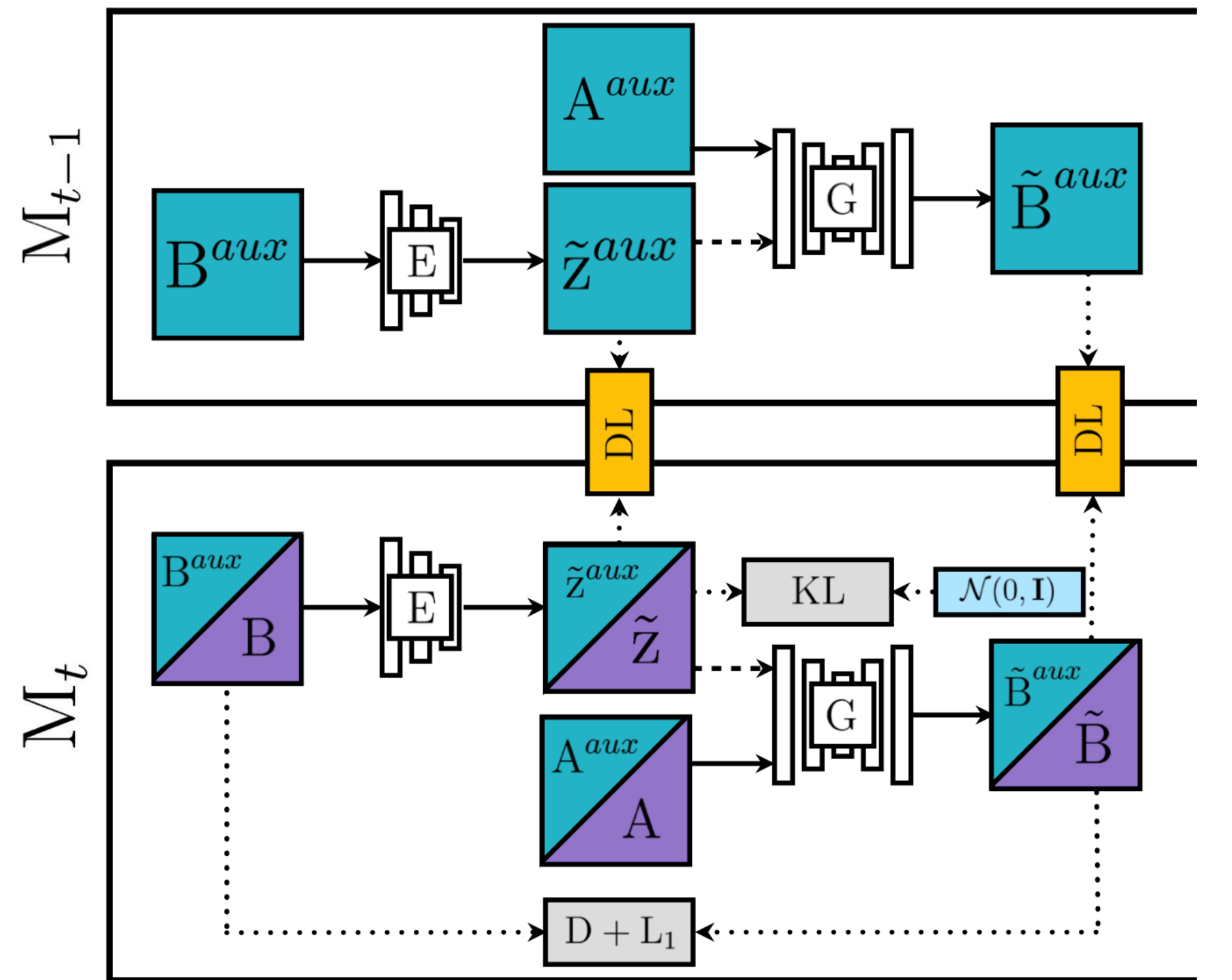
In summary: what do we want?



What could our expectations be, what might we desire?

- Constant memory budget?
- Pragmatically? A selection algorithm that outperforms randomly stored data points?
- A way to shrink the memory buffer to add new tasks, e.g. recursively select exemplars?
- Ideally? Knowledge of the distribution(s) and a subset with approximation guarantees?
- A natural formulation to allow (pseudo-)rehearsal, regularization...?
- many more ...?

Combining ideas



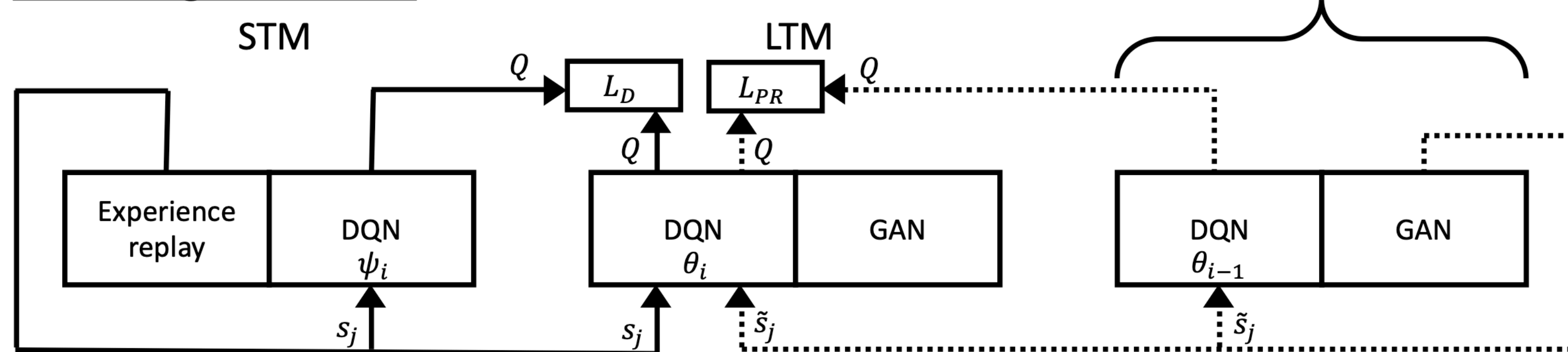
■: auxiliary data ■: data from new task ■: distillation loss
⋯→: input to loss function
□: prior $\mathcal{N}(0, I)$ - - ->: sampling

- We can (naturally) combine ideas
- In most of what we have seen so far however, the “architecture” has remained static.
- The assumption that we can continue accumulate knowledge because we have enough capacity has been prevalent.

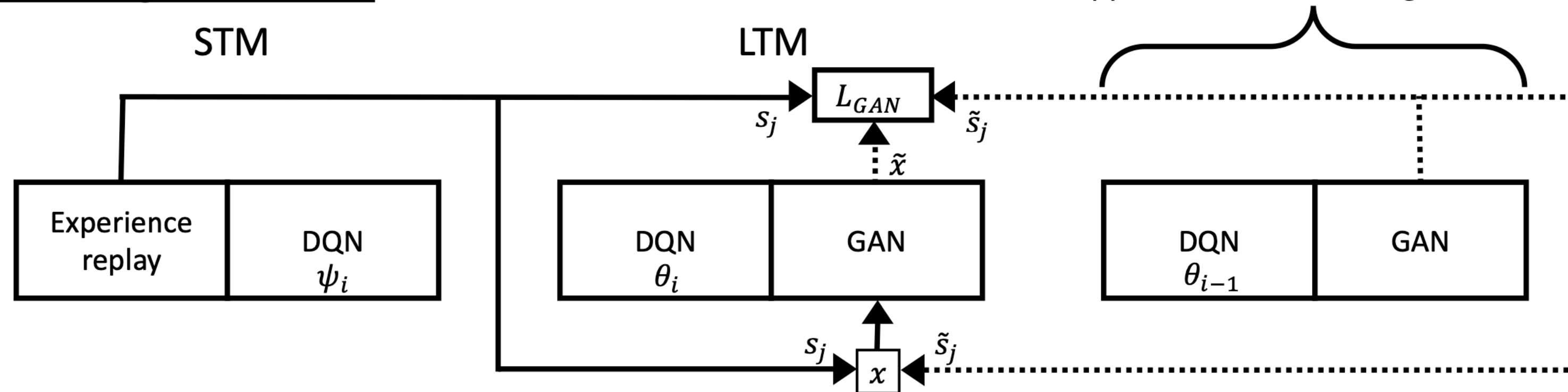
Combining ideas



Train long-term DQN



Train long-term GAN



Disclaimer

- There are A LOT of works that use rehearsal, which are not mentioned here
- They also come in all shapes of supervised/unsupervised/reinforced for different modalities