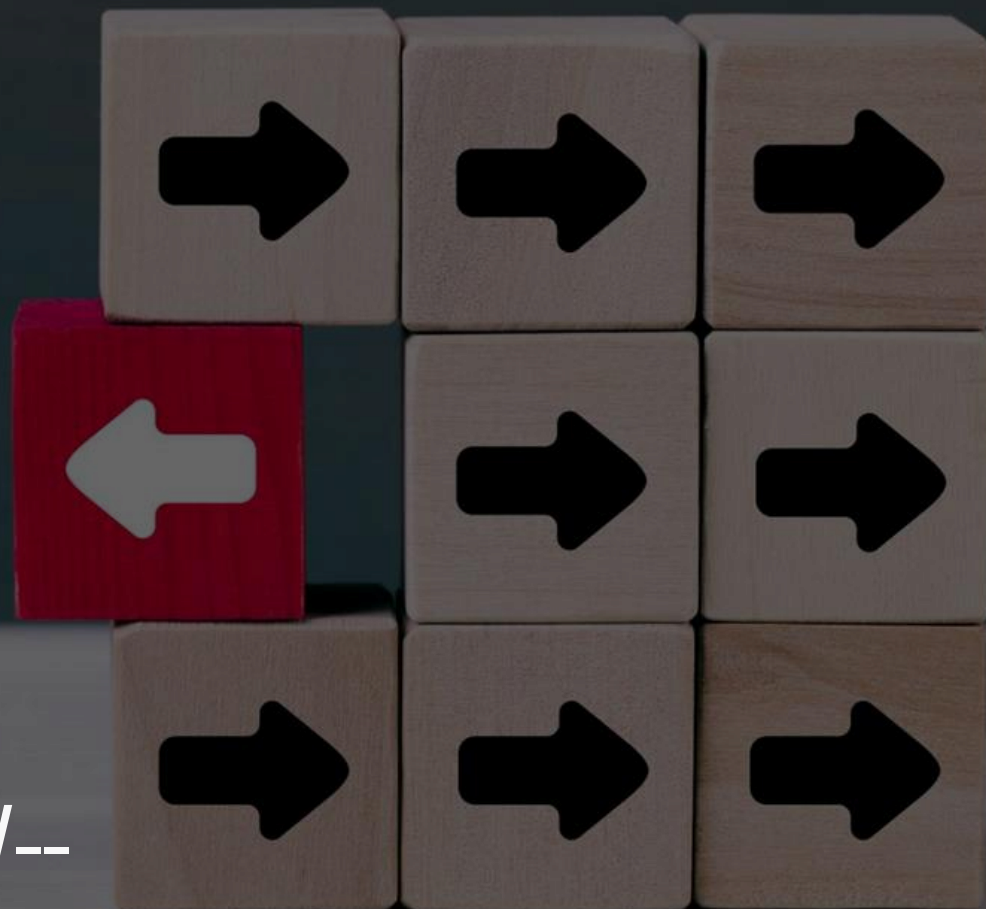


# WARSZTATY BADAWCZE: „METODY PRZENOSZENIA WIEDZY W SIECIACH NEURONOWYCH”

-- PRZEGLĄD ALGORYTMÓW--

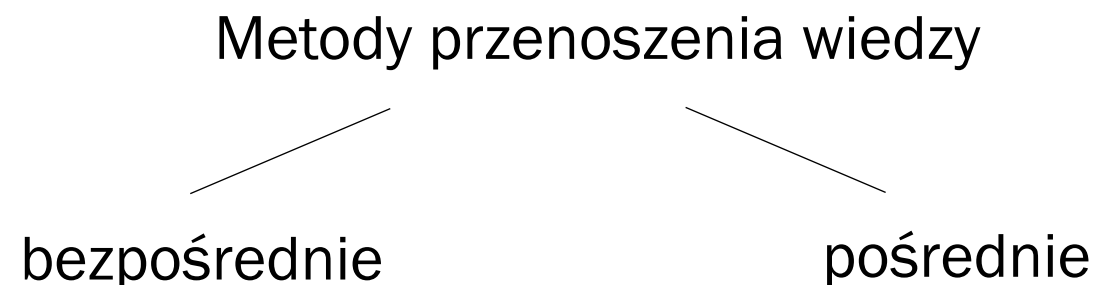


Paulina Tomaszewska

# Metody przenoszenia wiedzy vs. transfer learning

Metody przenoszenia wiedzy\*:

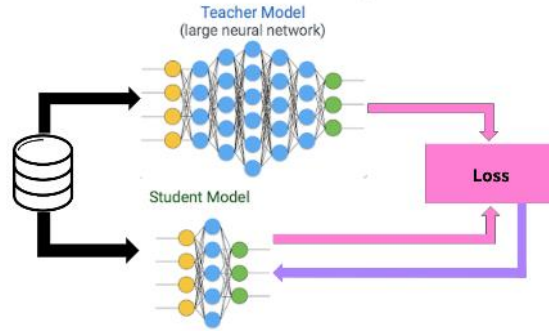
- transfer learning\*\*
- *multi-task learning*
- *continual learning*
- *knowledge distillation*
- *meta-learning*
- *domain adaptation*
- *etc.*



\*definicja własna, \*\* w j. polskim określenie „transfer learning” jest używane w wąskim kontekście

## Knowledge distillation

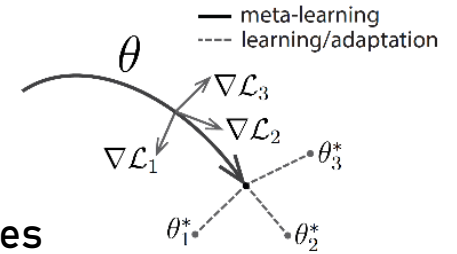
Goal: distil knowledge from big neural network to small one to achieve similar performance



## Meta Learning (learning to learn)

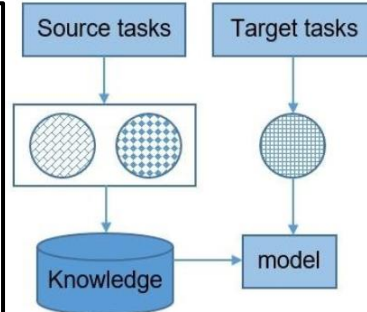
Example: MAML

Goal: learn weights that can be easily (by only few gradient updates) adjusted to new scenarios using only few samples



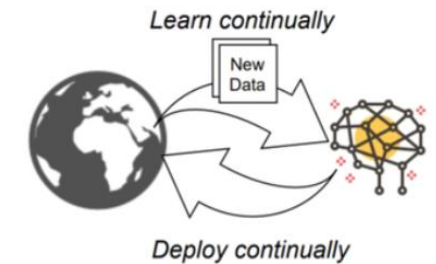
## Transfer Learning

Goal: train one model on a big dataset so that the model can be adjusted to the new dataset after less computationally demanding fine-tuning (during pretraining downstream data is not used)



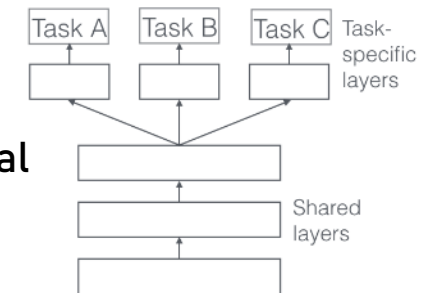
## Continual Learning

Goal: train model in an incremental way using non-stationary data sequence  
**Challenge**: catastrophic forgetting due to concept drifts



## Multi-task Learning

Goal: share the common knowledge between different tasks that leads to more universal representations



# KNOWLEDGE TRANSFER – ARTIFICIAL NEURAL NETWORKS

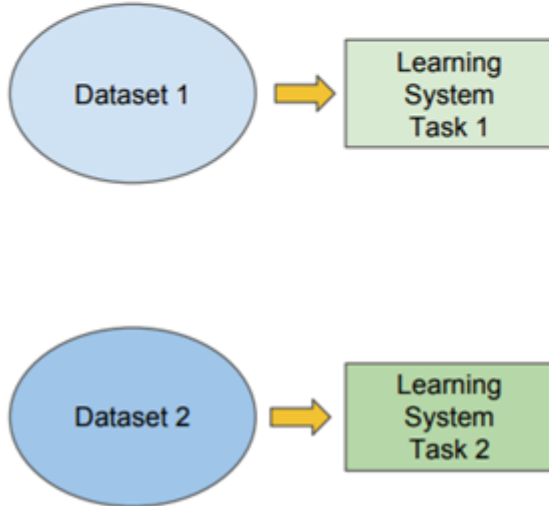


# TRANSFER LEARNING

# Transfer learning

## Traditional ML

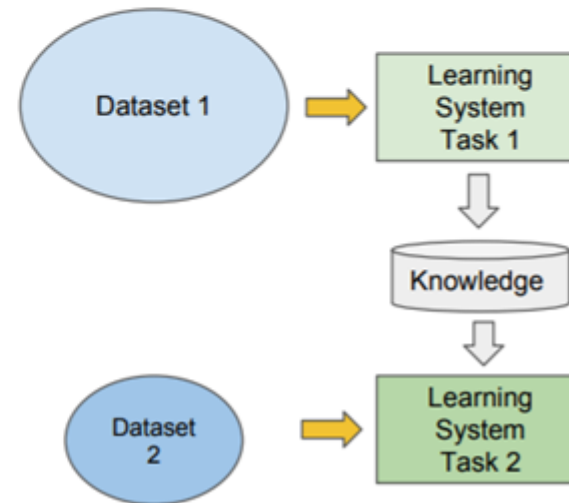
- Isolated, single task learning:
  - Knowledge is not retained or accumulated. Learning is performed w.o. considering past learned knowledge in other tasks



vs

## Transfer Learning

- Learning of a new task relies on the previous learned tasks:
  - Learning process can be faster, more accurate and/or need less training data



# Freeze or fine-tune?

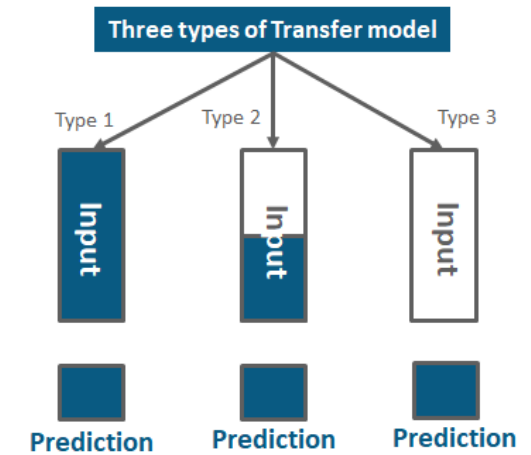
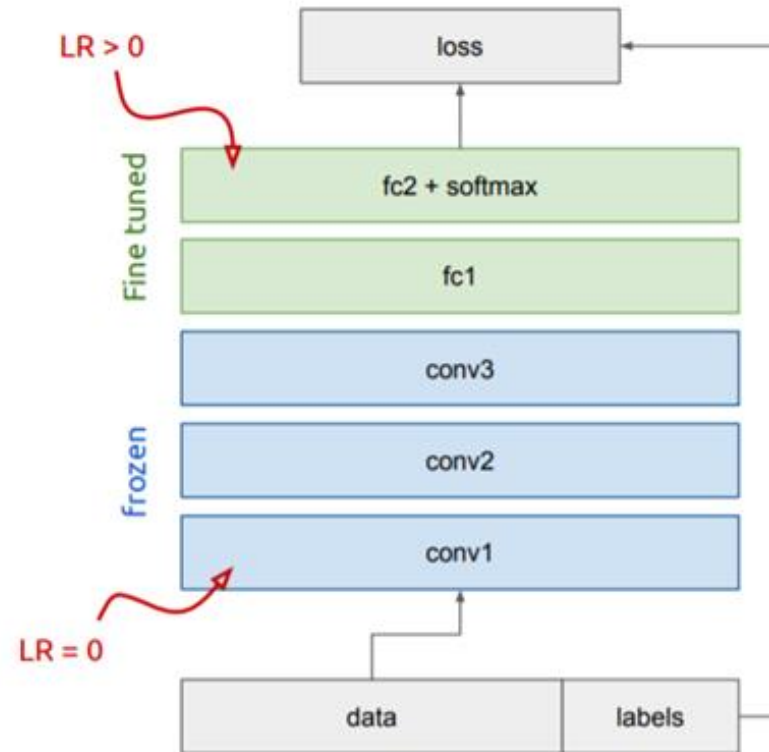
Bottom  $n$  layers can be frozen or fine tuned.

- **Frozen:** not updated during backprop
- **Fine-tuned:** updated during backprop

Which to do depends on target task:

- **Freeze:** target task labels are scarce, and we want to avoid overfitting
- **Fine-tune:** target task labels are more plentiful

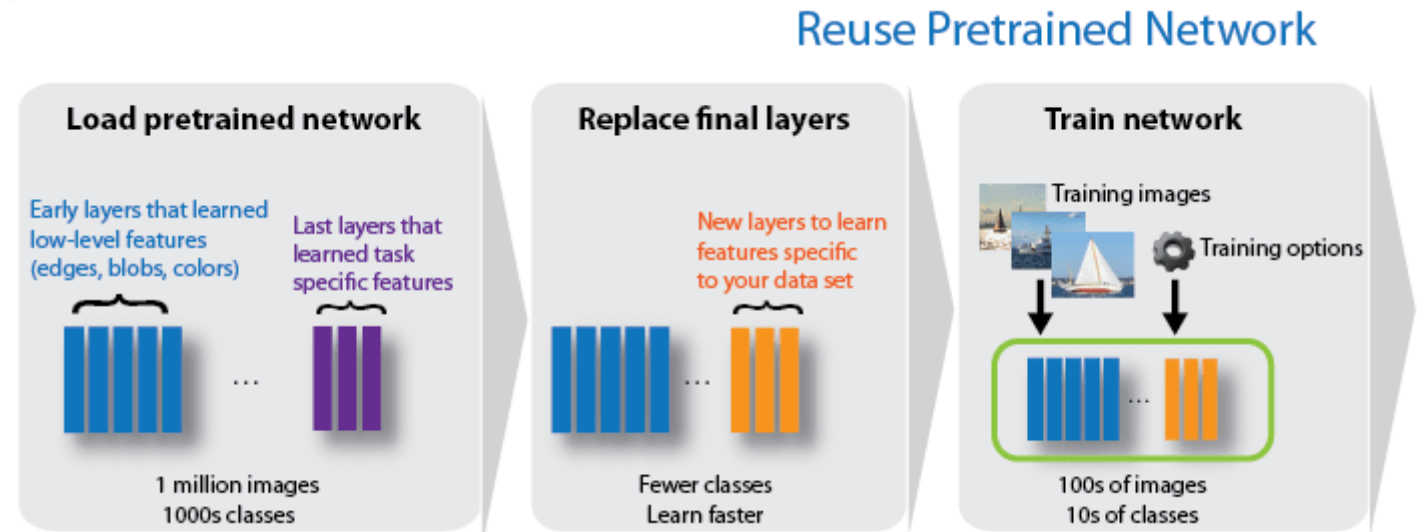
In general, we can set learning rates to be different for each layer to find a tradeoff between freezing and fine tuning



# Transfer learning in practice

Pretrained networks for Computer Vision:

- *ResNet*
- *VGG*
- *DenseNet*





# MULTI-TASK LEARNING (MTL)



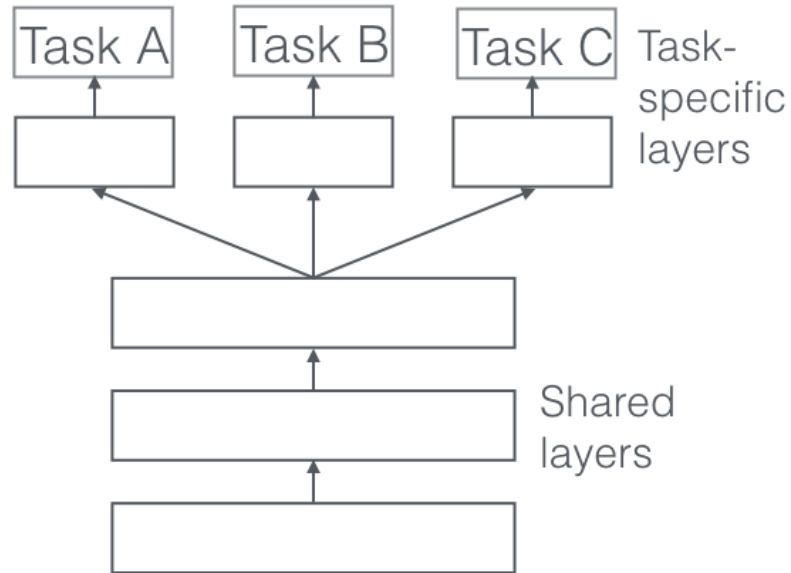


# Motivation

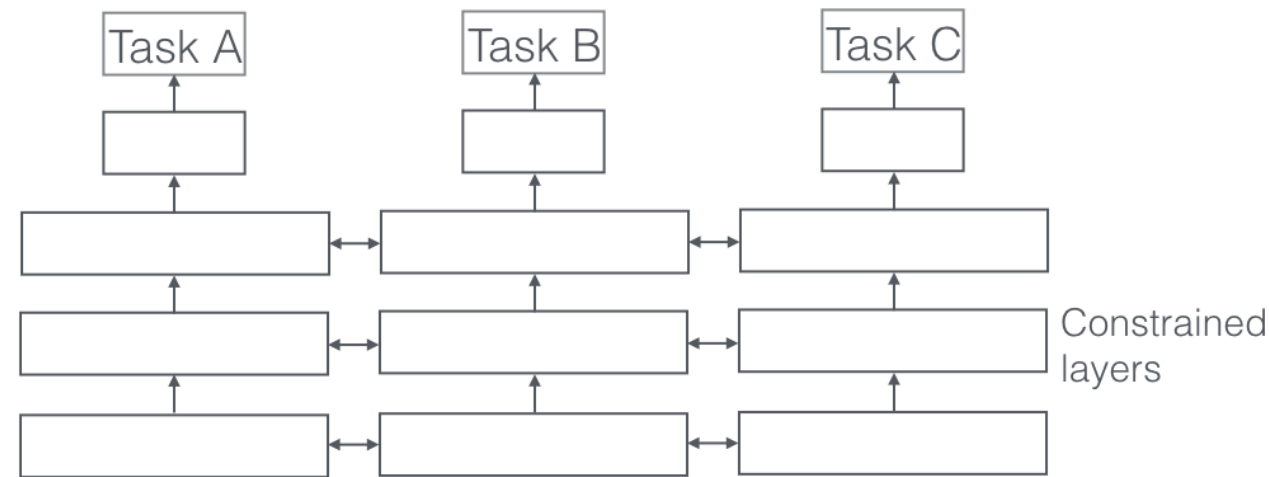
- compact models inspired by human ability of multi-tasking
- better representations (more universal)
- increased performance
- decreased inference time (autonomous cars)

Main problem: negative transfer

# MTL architectures – *shared trunk*

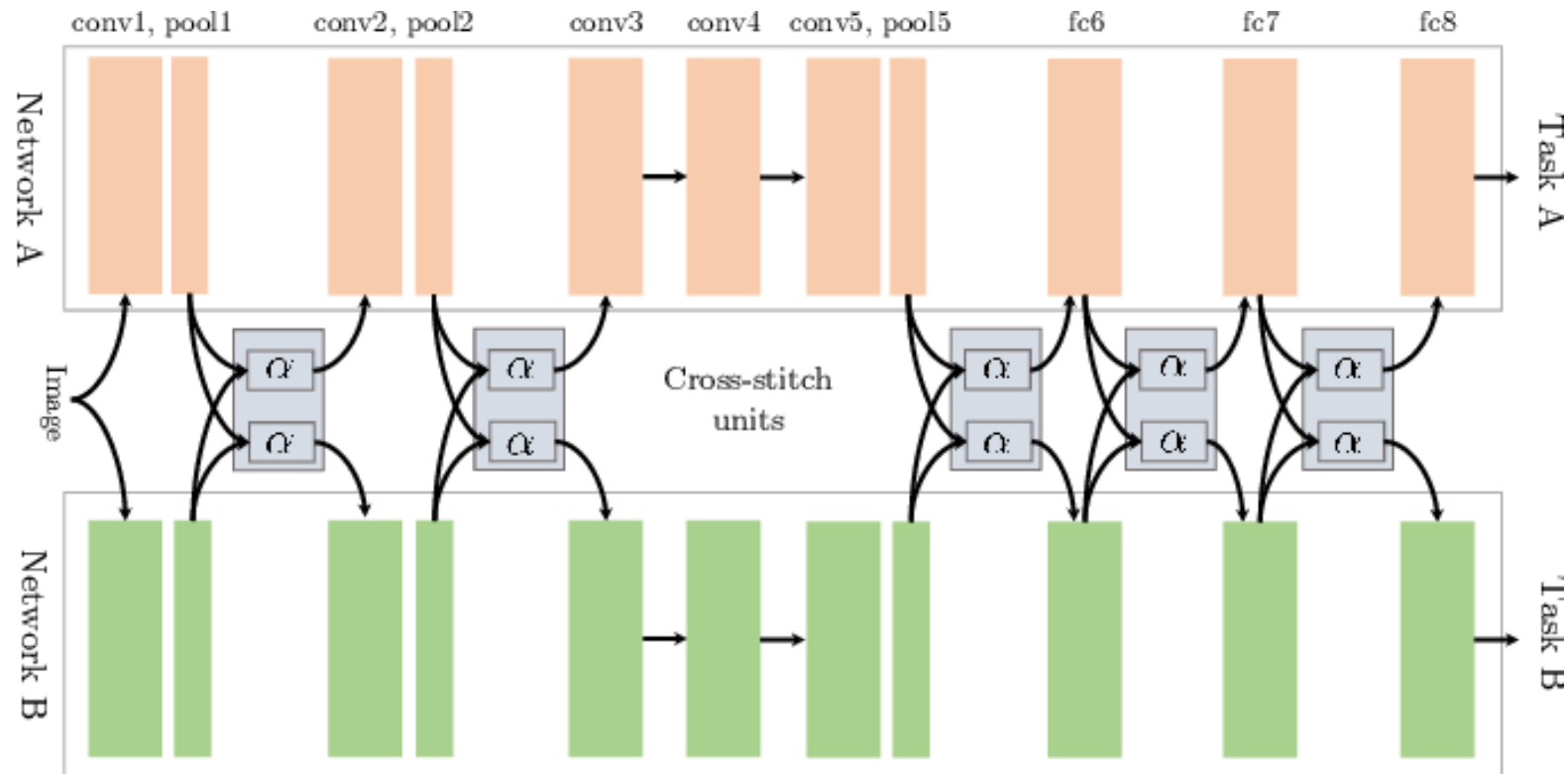


Hard parameter sharing



Soft parameter sharing

# MTL architectures – *cross talk*



Misra, Ishan et al. "Cross-Stitch Networks for Multi-task Learning." *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016): 3994-4003.

# Other research topics in MTL

- Optimization (e.g. GradNorm)
- Dataset (which tasks should be learnt together)

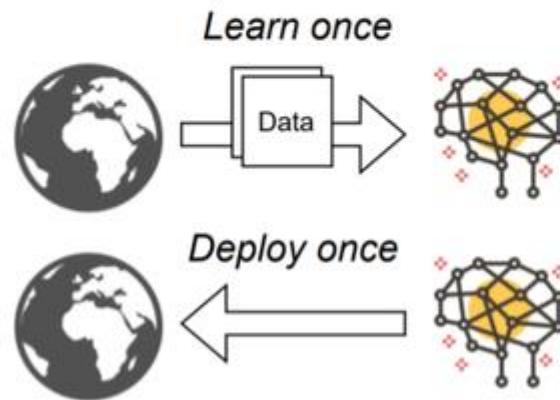


# CONTINUAL LEARNING

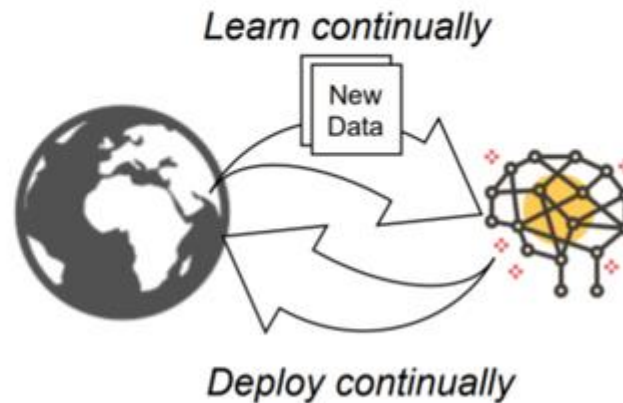
# Applications

- Used in case of nonstationary (=changing statistics over time) streaming data
- Similar terms: *Incremental learning*, *Life-long learning*

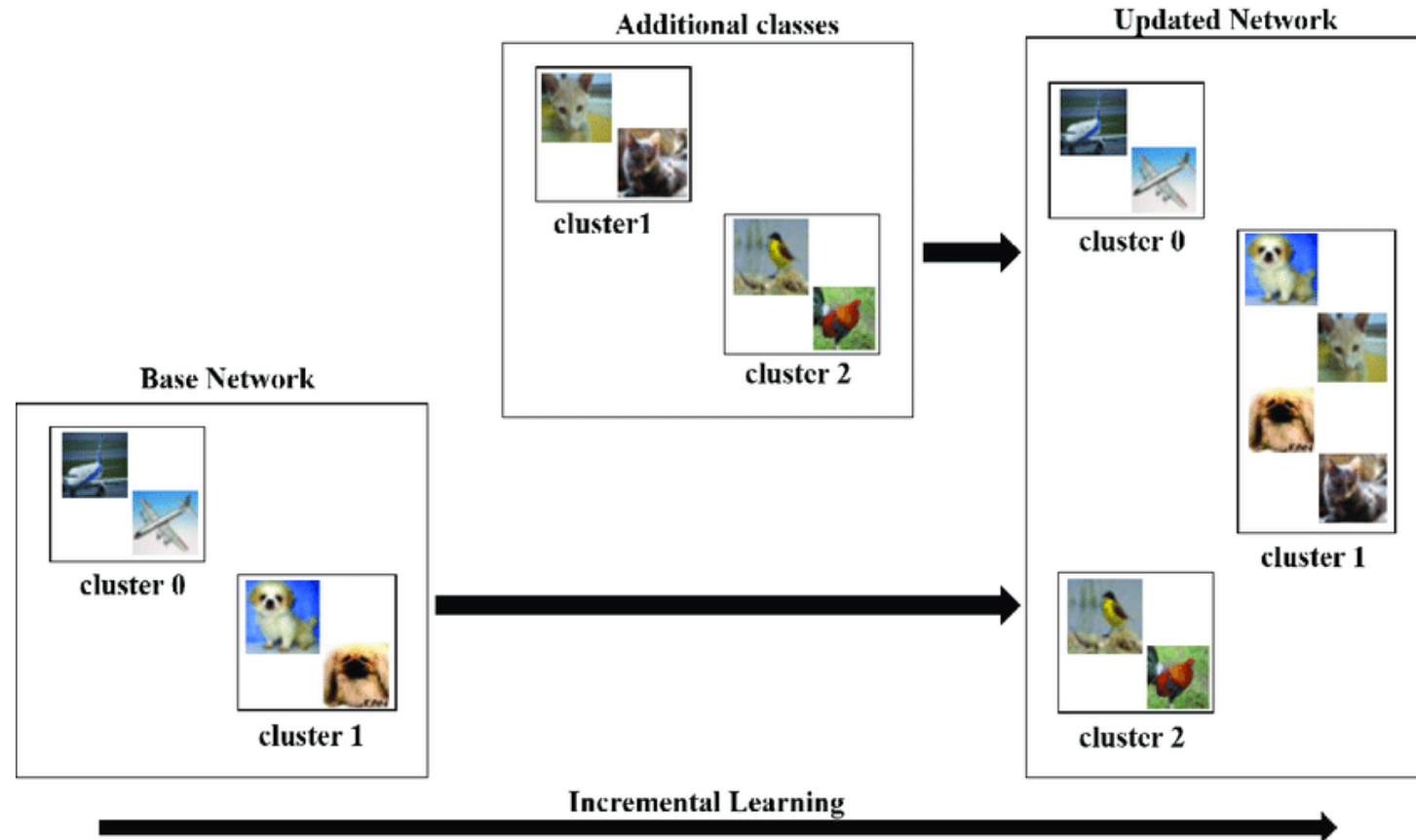
## Static ML



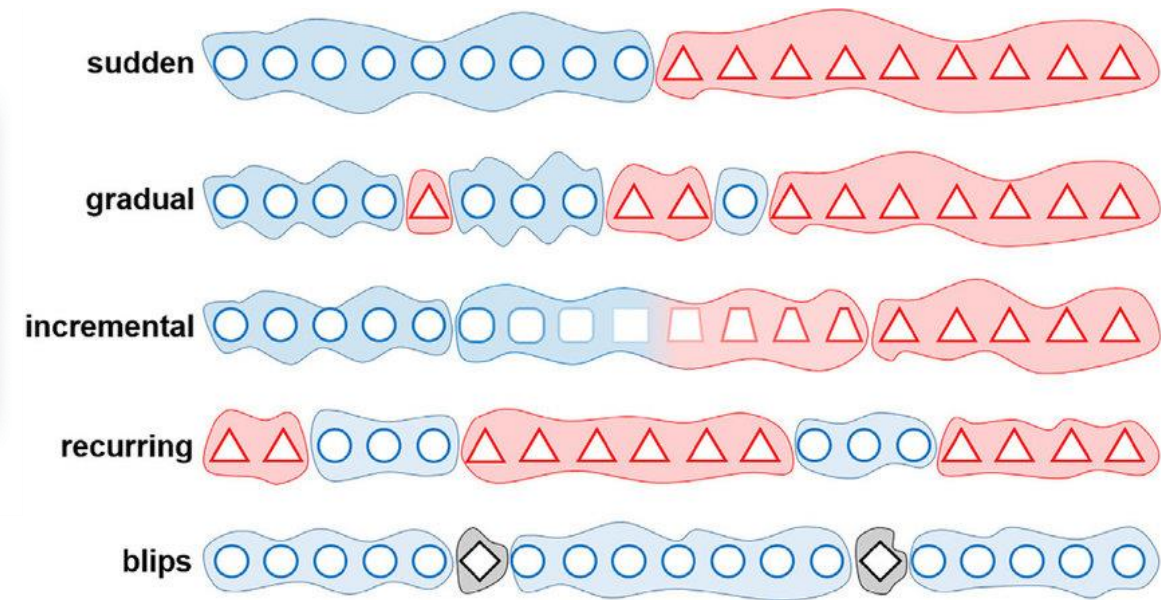
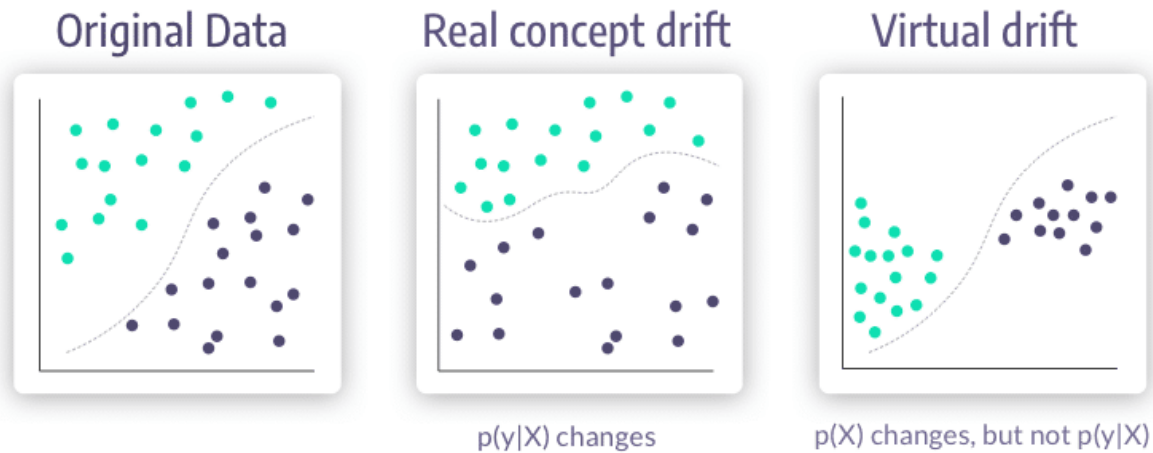
## Adaptive ML



# Challenges – additional classes



# Challenges - data drift

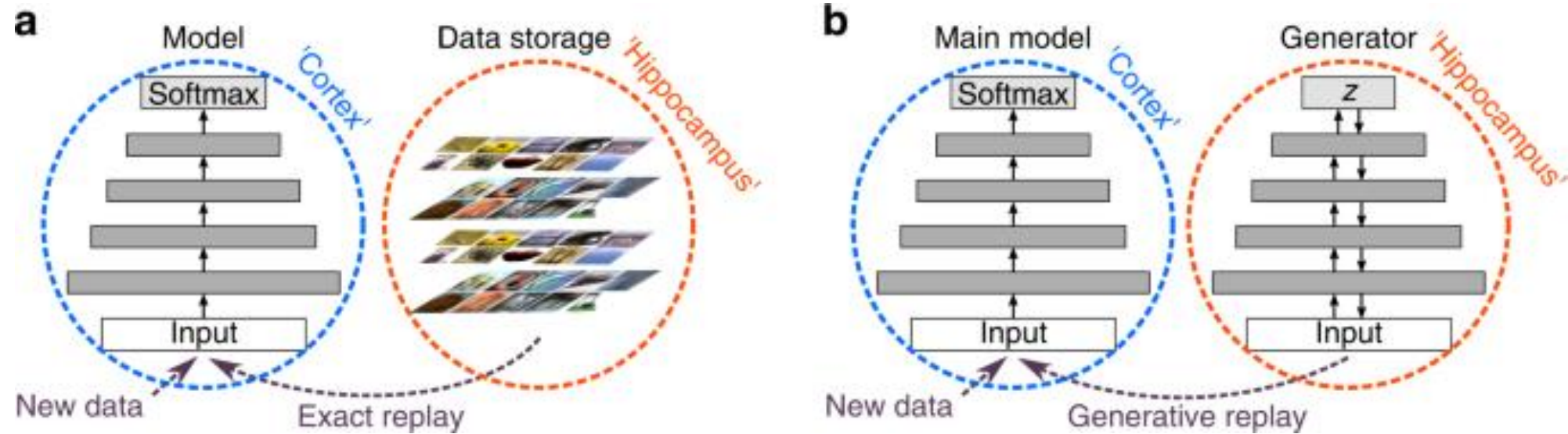


<https://www.iguazio.com/blog/concept-drift-deep-dive-how-to-build-a-drift-aware-ml-system/>

Krawczyk, Bartosz & Cano, Alberto. (2018). Online Ensemble Learning with Abstaining Classifiers for Drifting and Noisy Data Streams. Applied Soft Computing. 68. 677-692. 10.1016/j.asoc.2017.12.008.



# Possible remedy



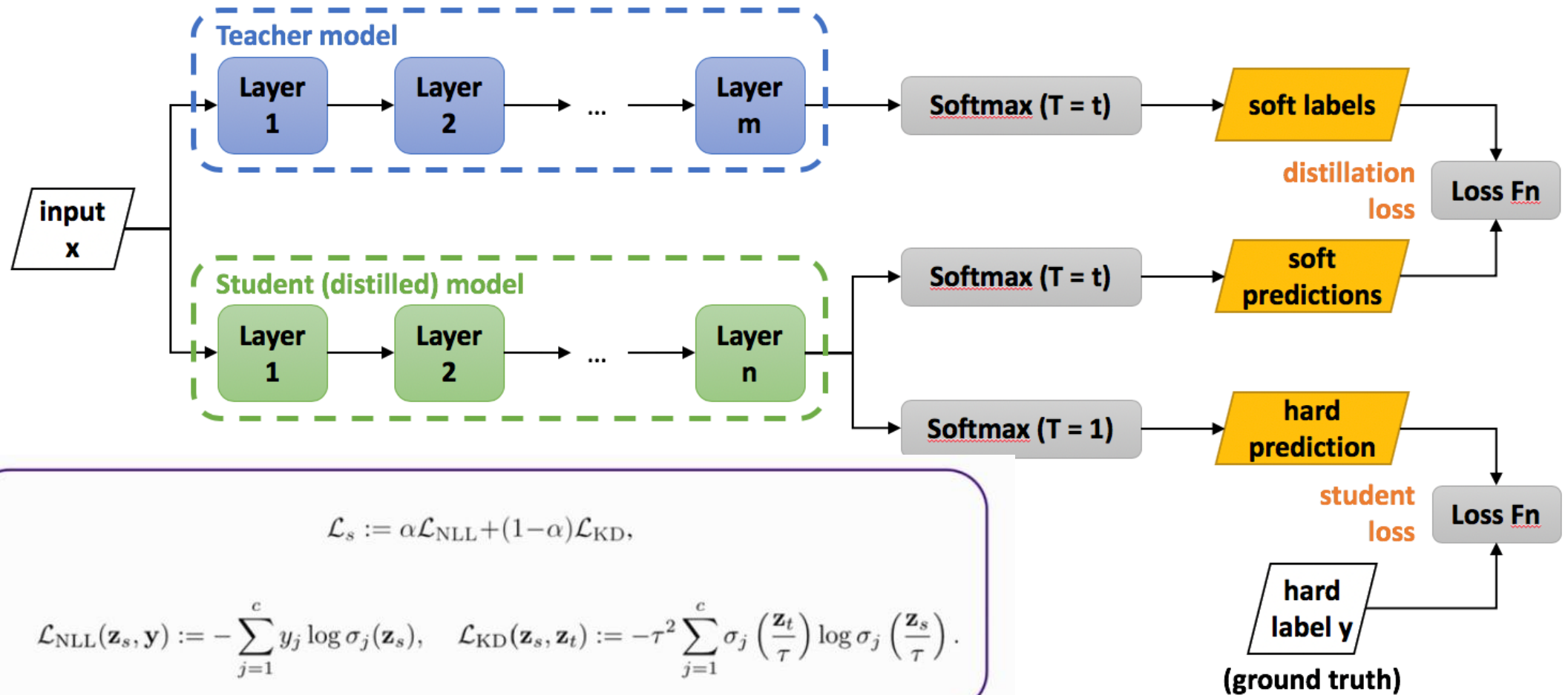


# KNOWLEDGE DISTILLATION



# Architecture

Compress big models so that they can be used at the edge devices (IoT)





META-LEARNING

# Motivation

- Other term: Learn to learn
- Goals:
  - *learning new concepts and skills fast with a few training examples*
  - *well adapting or generalizing to new tasks that have never been encountered during training time*
- The adaptation process (a mini learning session), happens during test but with a limited exposure to the new task configurations.
- Examples:
  - *A classifier trained on non-cat images can tell whether a given image contains a cat after seeing a handful of cat pictures*
  - *A game bot is able to quickly master a new game (meta reinforcement learning)*
  - *A mini robot completes the desired task on an uphill surface during test even though it was only trained in a flat surface environment (meta reinforcement learning)*

# Problem statement

A good meta-learning model should be trained over a variety of learning tasks and optimized for the best performance on a distribution of tasks, including potentially unseen tasks. Each task is associated with a dataset  $\mathcal{D}$ , containing both feature vectors and true labels. The optimal model parameters are:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{\mathcal{D} \sim p(\mathcal{D})} [\mathcal{L}_{\theta}(\mathcal{D})]$$

It looks very similar to a normal learning task, but *one dataset* is considered as *one data sample*.

# Few-shot classification

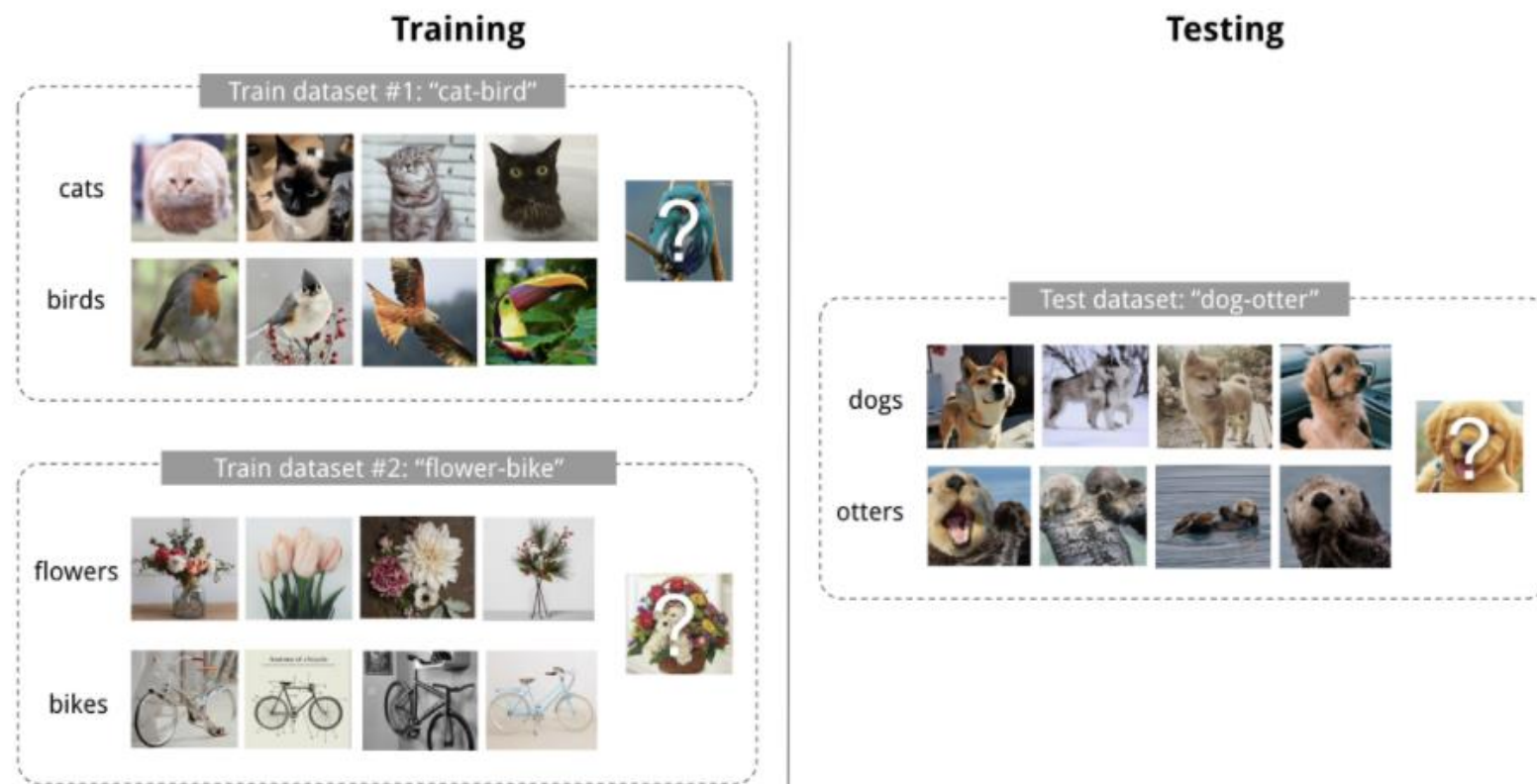


Fig. 1. An example of 4-shot 2-class image classification. (Image thumbnails are from Pinterest)

## Learner and Meta-Learner

Another popular view of meta-learning decomposes the model update into two stages:

- A classifier  $f_\theta$  is the “learner” model, trained for operating a given task;
- In the meantime, a optimizer  $g_\phi$  learns how to update the learner model’s parameters via the support set  $S$ ,  $\theta' = g_\phi(\theta, S)$ .

Then in final optimization step, we need to update both  $\theta$  and  $\phi$  to maximize:

$$\mathbb{E}_{L \subset \mathcal{L}} [\mathbb{E}_{S^L \subset \mathcal{D}, B^L \subset \mathcal{D}} [\sum_{(\mathbf{x}, y) \in B^L} P_{g_\phi(\theta, S^L)}(y|\mathbf{x})]]$$



# Common approaches in meta-learning



Model-based

Metric-based

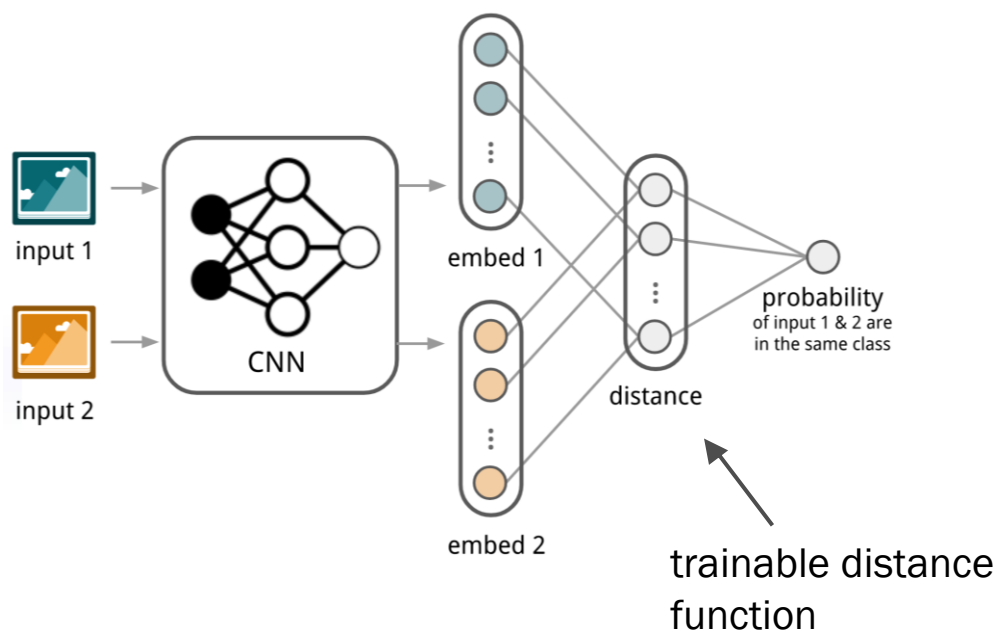
Optimization-based

A thick black L-shaped frame is positioned on the left and bottom edges of the slide, framing the central text.

# METRIC-BASED META-LEARNING

Analogy to k-nearest neighbours

# Siamese networks



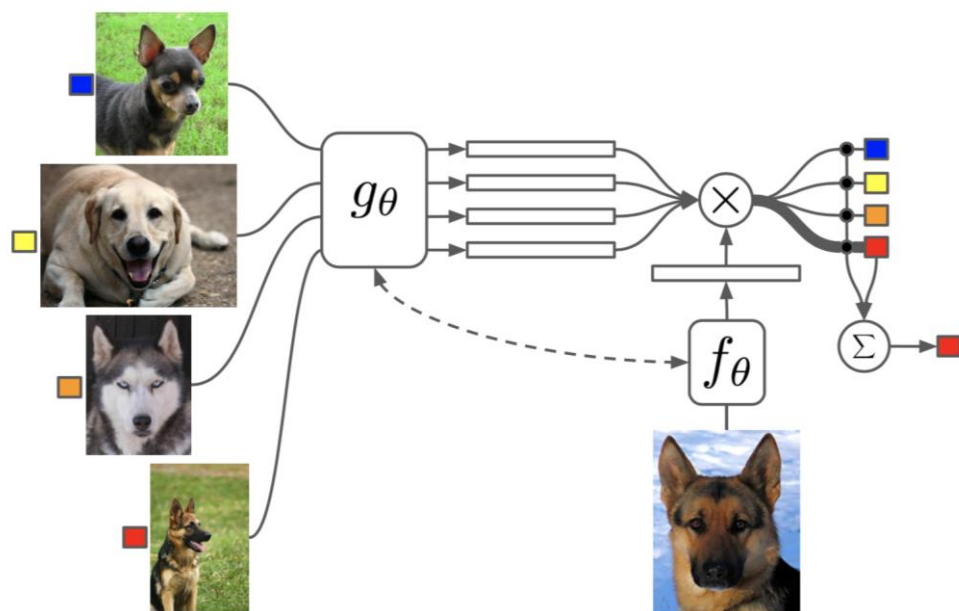
Applications:

- recognizing a single person at a train station or airport.
- verifying whether the photo in a pass is the same as the person claiming he or she is the same person.

1. First, convolutional siamese network learns to encode two images into feature vectors via a embedding function  $f_{\theta}$  which contains a couple of convolutional layers.
2. The L1-distance between two embeddings is  $|f_{\theta}(\mathbf{x}_i) - f_{\theta}(\mathbf{x}_j)|$ .
3. The distance is converted to a probability  $p$  by a linear feedforward layer and sigmoid. It is the probability of whether two images are drawn from the same class.
4. Intuitively the loss is cross entropy because the label is binary.

$$p(\mathbf{x}_i, \mathbf{x}_j) = \sigma(\mathbf{W}|f_{\theta}(\mathbf{x}_i) - f_{\theta}(\mathbf{x}_j)|)$$
$$\mathcal{L}(B) = \sum_{(\mathbf{x}_i, \mathbf{x}_j, y_i, y_j) \in B} \mathbf{1}_{y_i=y_j} \log p(\mathbf{x}_i, \mathbf{x}_j) + (1 - \mathbf{1}_{y_i=y_j}) \log(1 - p(\mathbf{x}_i, \mathbf{x}_j))$$

# Matching networks

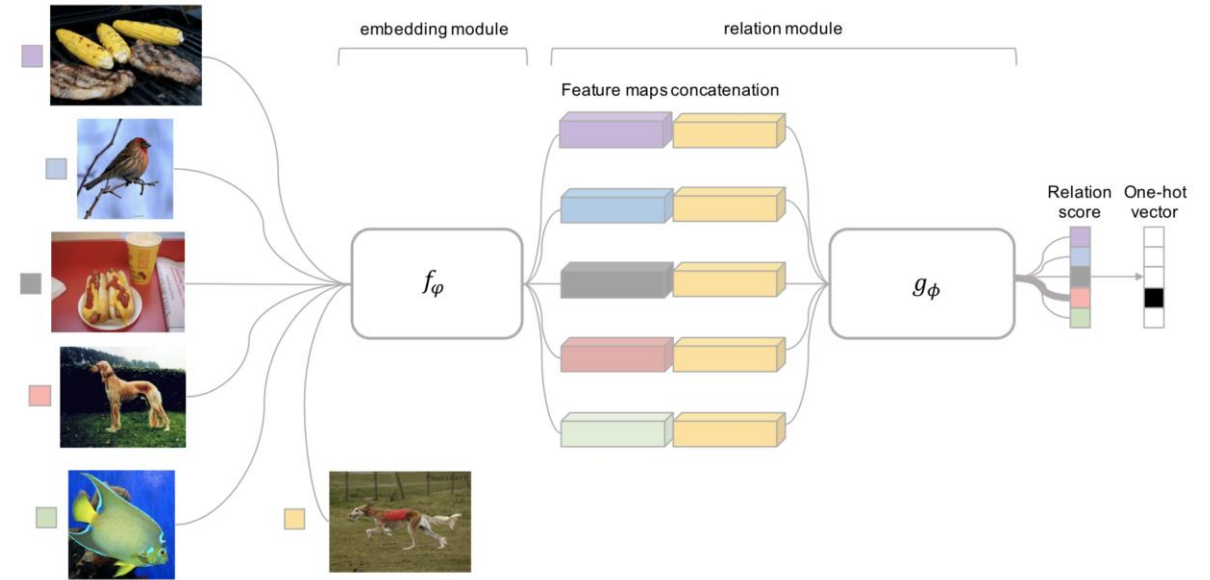


$$c_S(\mathbf{x}) = P(y|\mathbf{x}, S) = \sum_{i=1}^k a(\mathbf{x}, \mathbf{x}_i) y_i, \text{ where } S = \{(\mathbf{x}_i, y_i)\}_{i=1}^k$$

The attention kernel depends on two embedding functions,  $f$  and  $g$ , for encoding the test sample and the support set samples respectively. The attention weight between two data points is the cosine similarity,  $\text{cosine}(\cdot)$ , between their embedding vectors, normalized by softmax:

$$a(\mathbf{x}, \mathbf{x}_i) = \frac{\exp(\text{cosine}(f(\mathbf{x}), g(\mathbf{x}_i)))}{\sum_{j=1}^k \exp(\text{cosine}(f(\mathbf{x}), g(\mathbf{x}_j)))}$$

# Relation Network

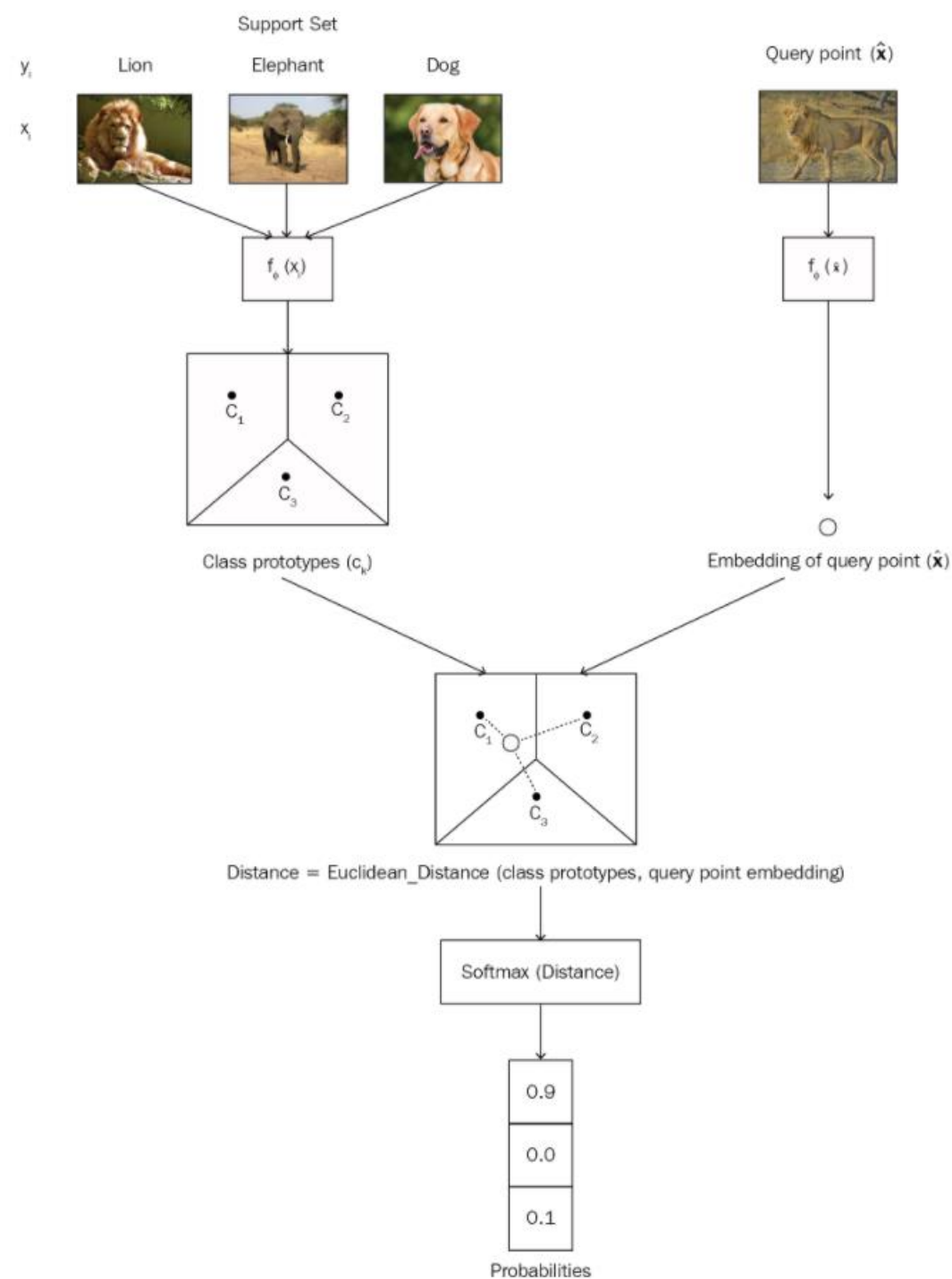
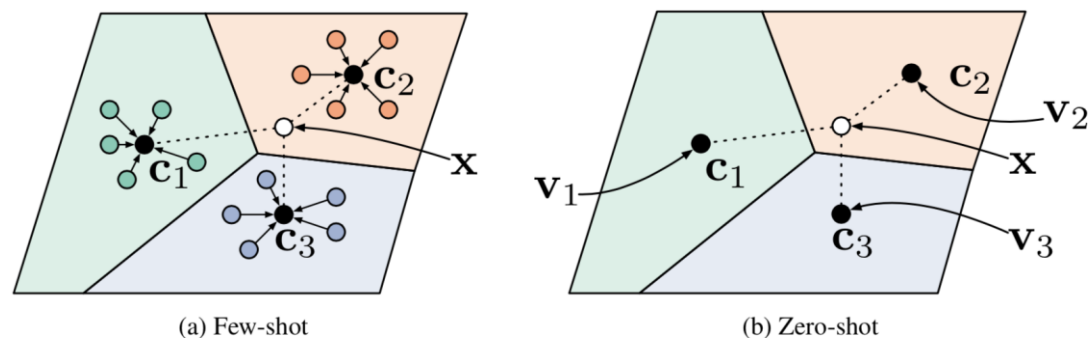


**Relation Network (RN)** (Sung et al., 2018) is similar to siamese network but with a few differences:

1. The relationship is not captured by a simple L1 distance in the feature space, but predicted by a CNN classifier  $g_\phi$ . The relation score between a pair of inputs,  $\mathbf{x}_i$  and  $\mathbf{x}_j$ , is  $r_{ij} = g_\phi([\mathbf{x}_i, \mathbf{x}_j])$  where  $[\cdot, \cdot]$  is concatenation.
2. The objective function is MSE loss instead of cross-entropy, because conceptually RN focuses more on predicting relation scores which is more like regression, rather than binary classification,

$$\mathcal{L}(B) = \sum_{(\mathbf{x}_i, \mathbf{x}_j, y_i, y_j) \in B} (r_{ij} - \mathbf{1}_{y_i=y_j})^2.$$

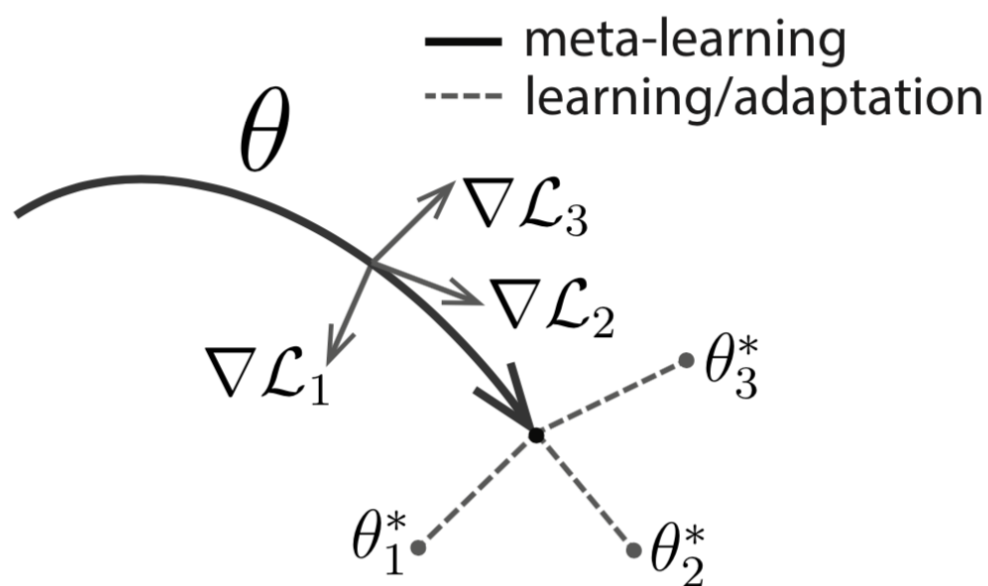
# Prototypical networks



A thick black L-shaped frame is positioned on the left and bottom edges of the slide, framing the central text.

# OPTIMIZATION-BASED META-LEARNING

# Model-Agnostic Meta-Learning (MAML)



---

## 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** **Note: the meta-update is using different set of data.**
  - 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**
-

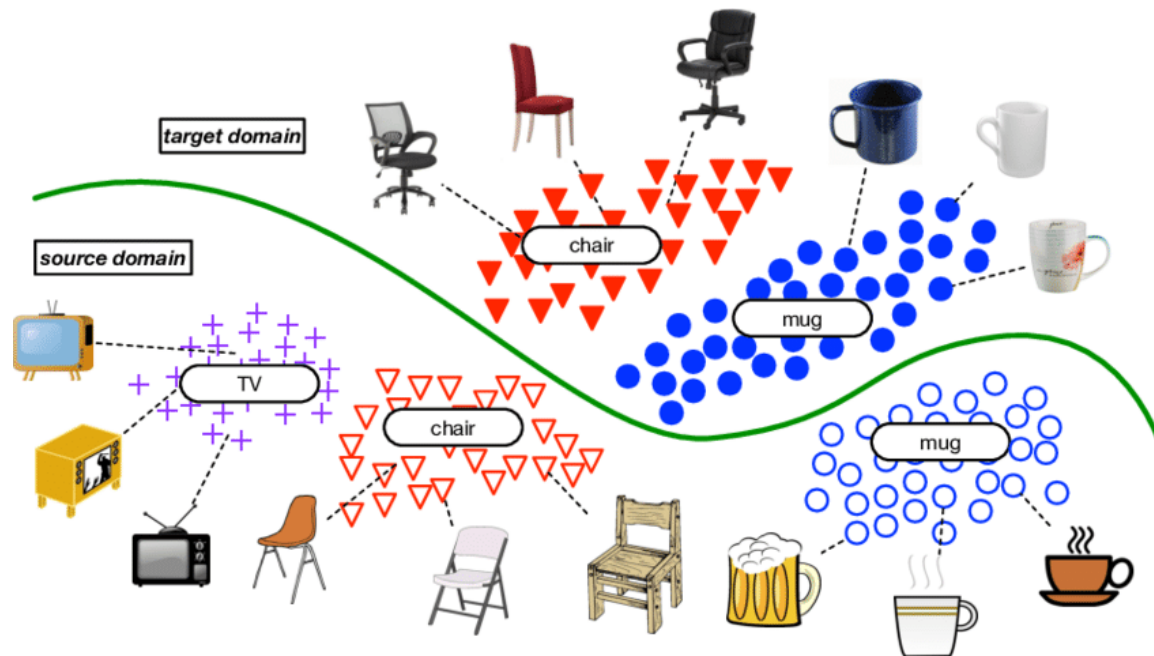




# DOMAIN ADAPTATION



# Motivation



- learning **fair representations** (insensitive to some given attributes), while retaining enough information to solve the task

More information at: [https://ameroyer.github.io/reading-notes/representation%20learning/2019/04/26/gradient\\_reversal\\_against\\_discrimination.html](https://ameroyer.github.io/reading-notes/representation%20learning/2019/04/26/gradient_reversal_against_discrimination.html)

The proposed model builds on the **Domain Adversarial Network (DANN)** [1], originally introduced for unsupervised domain adaptation. Given some labeled data  $(x, y) \sim \mathcal{X} \times \mathcal{Y}$ , and some unlabeled data  $\tilde{x} \sim \tilde{\mathcal{X}}$ , the goal is to learn a network that solves both classification tasks  $\mathcal{X} \rightarrow \mathcal{Y}$  and  $\tilde{\mathcal{X}} \rightarrow \mathcal{Y}$  while learning a shared representation between  $\mathcal{X}$  and  $\tilde{\mathcal{X}}$ .

# Omówienie kodu – rozwiązań z kanonu

- Siemese networks: <https://github.com/sudharsan13296/Hands-On-Meta-Learning-With-Python/blob/master/02.%20Face%20and%20Audio%20Recognition%20using%20Siamese%20Networks/2.4%20Face%20Recognition%20Using%20Siamese%20Network.ipynb>
- Prototypical networks: <https://github.com/PacktPublishing/Hands-On-Meta-Learning-with-Python/tree/master/Chapter03>
- Relation networks, Matching networks: <https://github.com/PacktPublishing/Hands-On-Meta-Learning-with-Python/tree/master/Chapter04>
- MAML: <https://github.com/sudharsan13296/Hands-On-Meta-Learning-With-Python/tree/master/06.%20MAML%20and%20it's%20Variants>