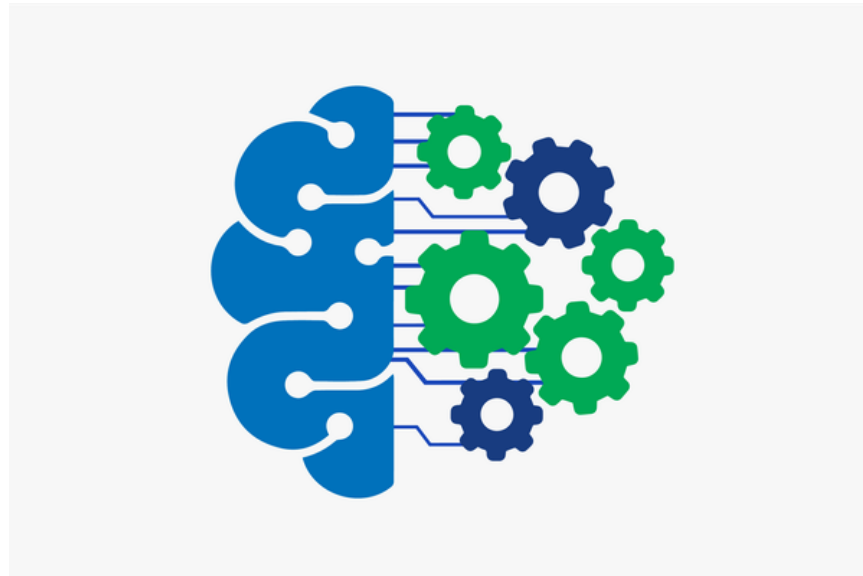


Learning Concepts

Narada Warakagoda



Introduction

- What is learning?
 - Adjust model parameters according to data
- Why is it challenging?
 - We do not have enough data
 - Learning takes time
 - Learning demands computational resources
- What is the remedy?
 - Learn efficiently with as little data as possible

Learning buzzwords

Deep Learning

Supervised learning

Unsupervised learning

Reinforcement learning

Semi-supervised learning

Transfer learning

Curriculum learning

Self-supervised learning

Self taught learning

Machine Learning

Hebbian learning

Active learning

Feature learning

Metric learning

Multi-task learning

Multi-modal learning

Few-shot learning

One-shot learning

Zero-shot learning

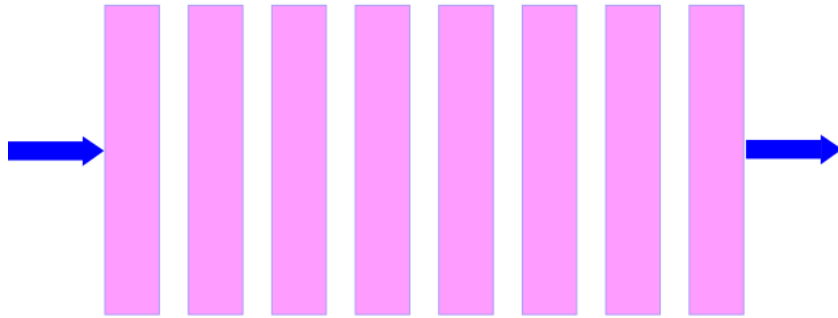
Meta learning

Transfer Learning

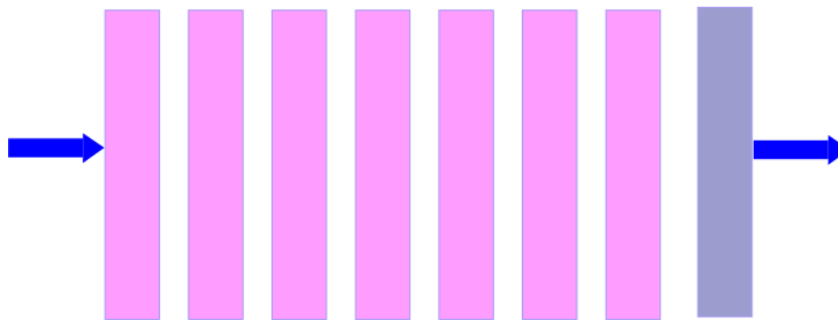
- General definition:
 - Using the knowledge gathered in one learning task in another learning task
- In deep learning:
 - Re-train a previously trained network with new data
- Advantage:
 - Can train a large network with relatively small (new) data set.

Possible Approaches

Original pretrained

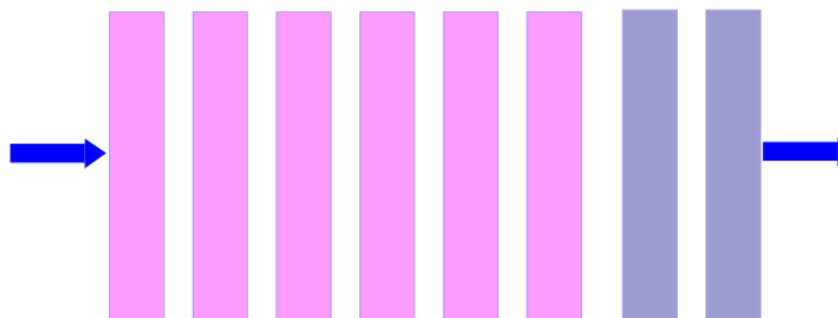


Transfer learning last layer



- New data set is very small
- New data and original data from the same domain
- Different class structure

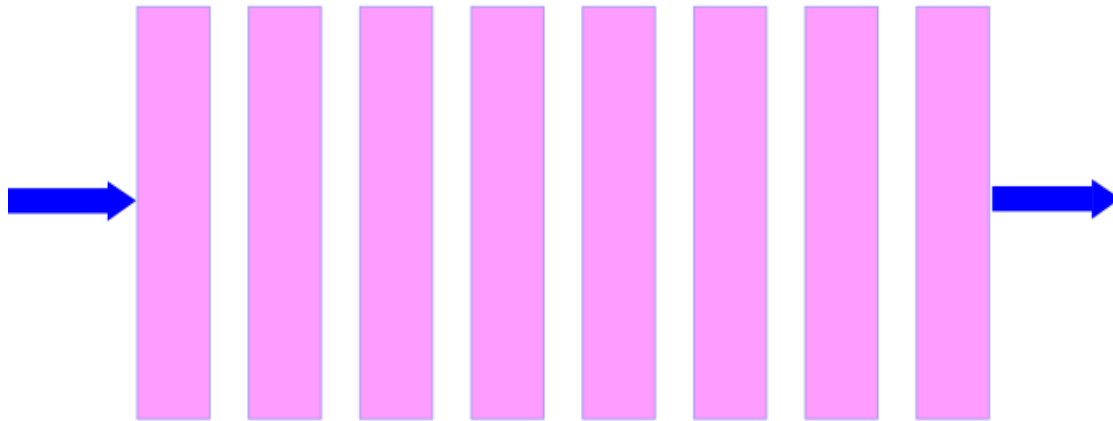
Transfer learning last two layers



- New data set is less small
- Same domain, different class structure

Possible Approaches

Original pretrained



Transfer learning all layers

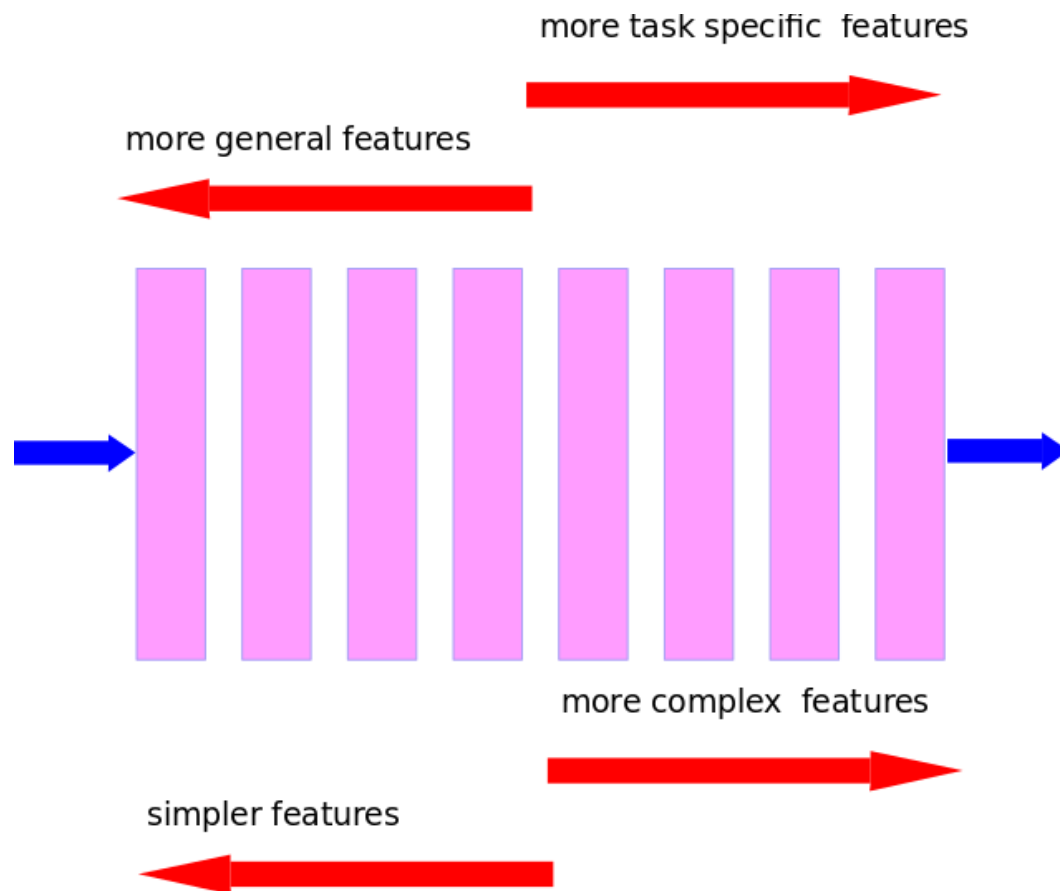


New data is from a different domain than original data.

Eg: Optical images and sonar images.

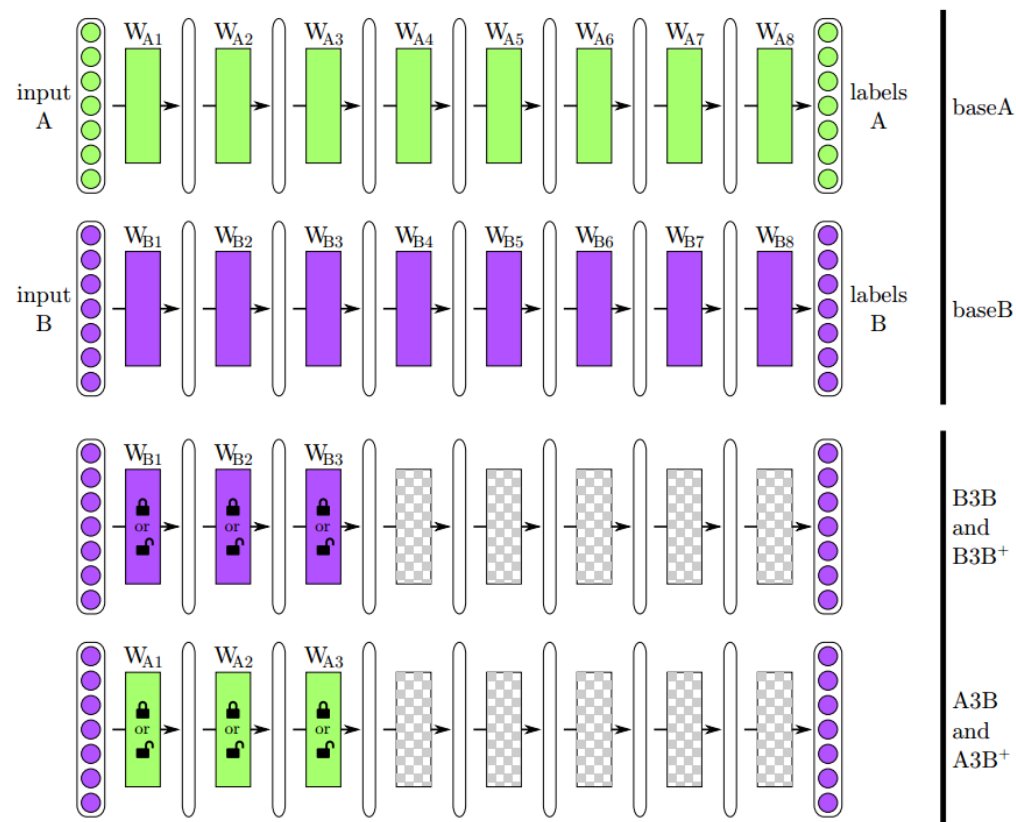
Why is Transfer Learning Possible?

- The network learns general feature representations
- More general, more transferable?

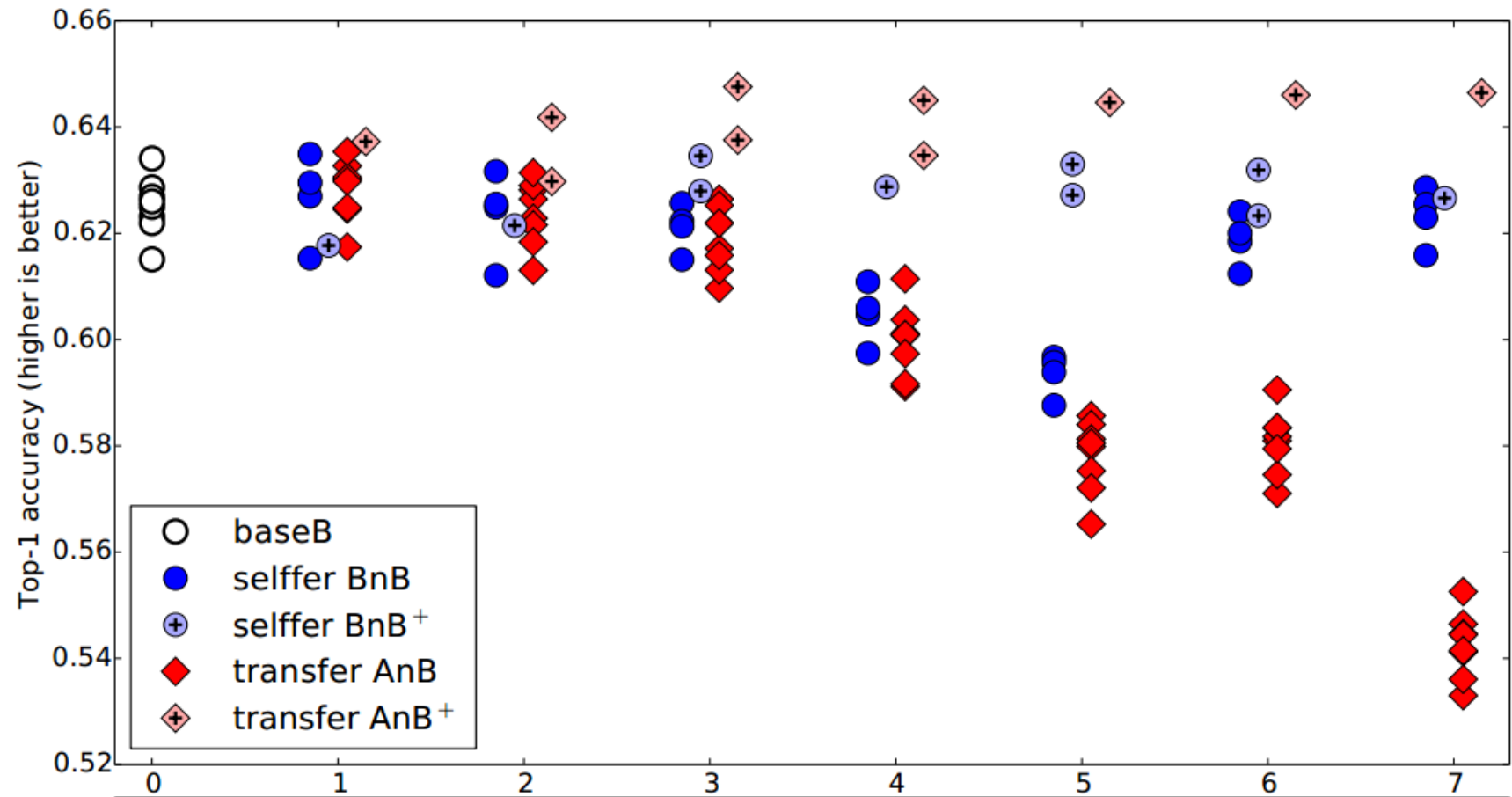


How Transferable are Features?

- Experiment reported in *How transferable are features in deep neural networks?*, Yosinski et.al 2014
 - Two pre-trained networks baseA and baseB
 - Same architecture trained with data sets A and B
- AnB = First n layers from baseA, rest are retrained with data set B
 - AnB+ = Same approach but the whole network is re-trained
- BnB = First n layers from baseB, rest are retrained with data set B
 - BbB+ = Same approach but the whole network is retrained

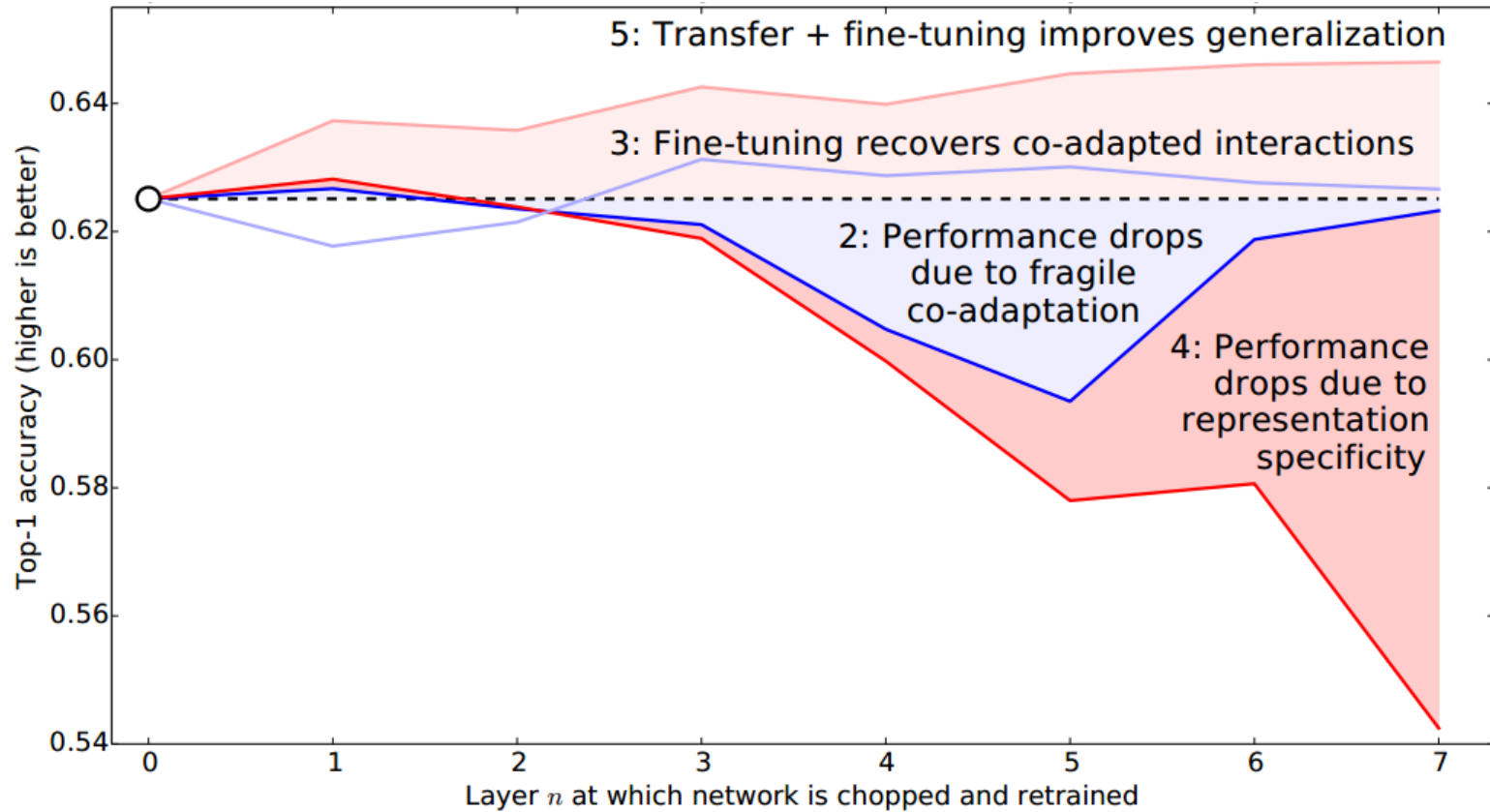


How Transferable are Features?



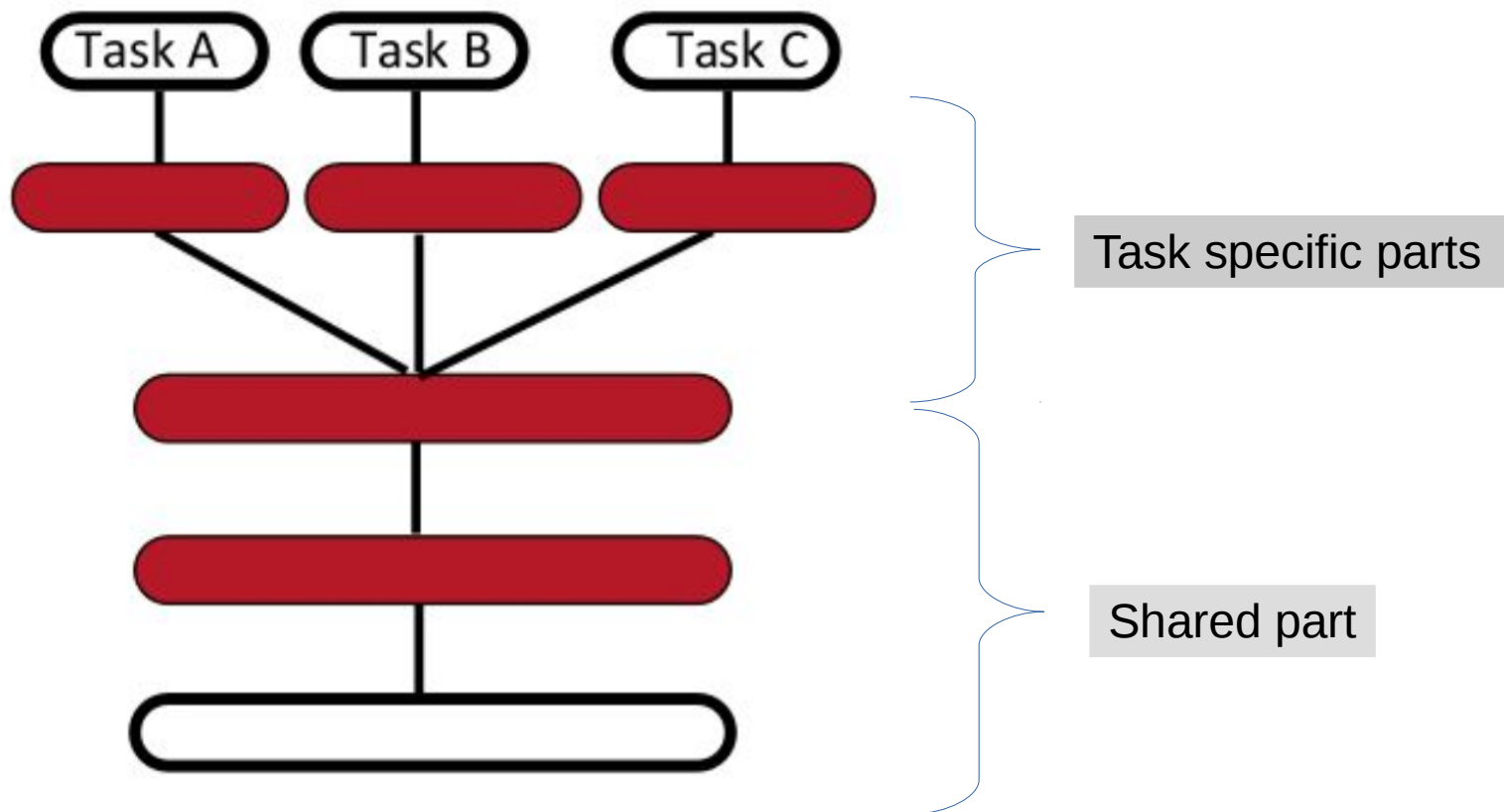
How Transferable are Features?

• C
-
-

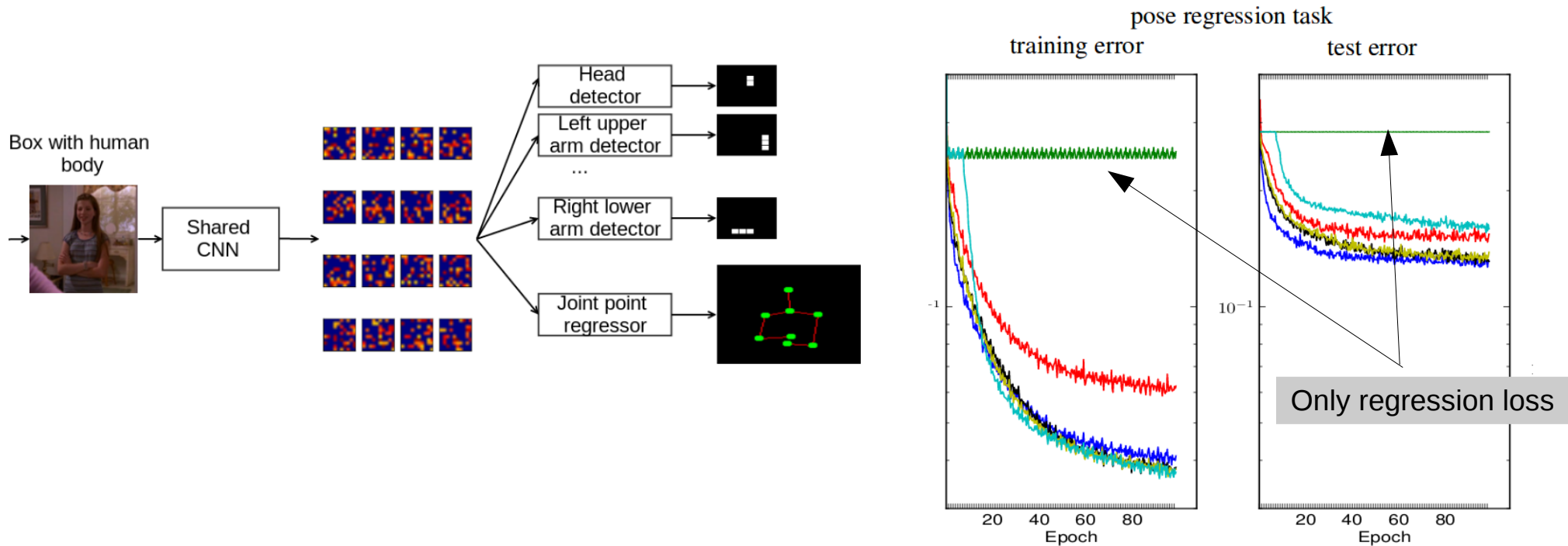


Multi-task Learning

- Sharing a single network for several different tasks
- Enhanced regularization effect



Multi-task Learning Example

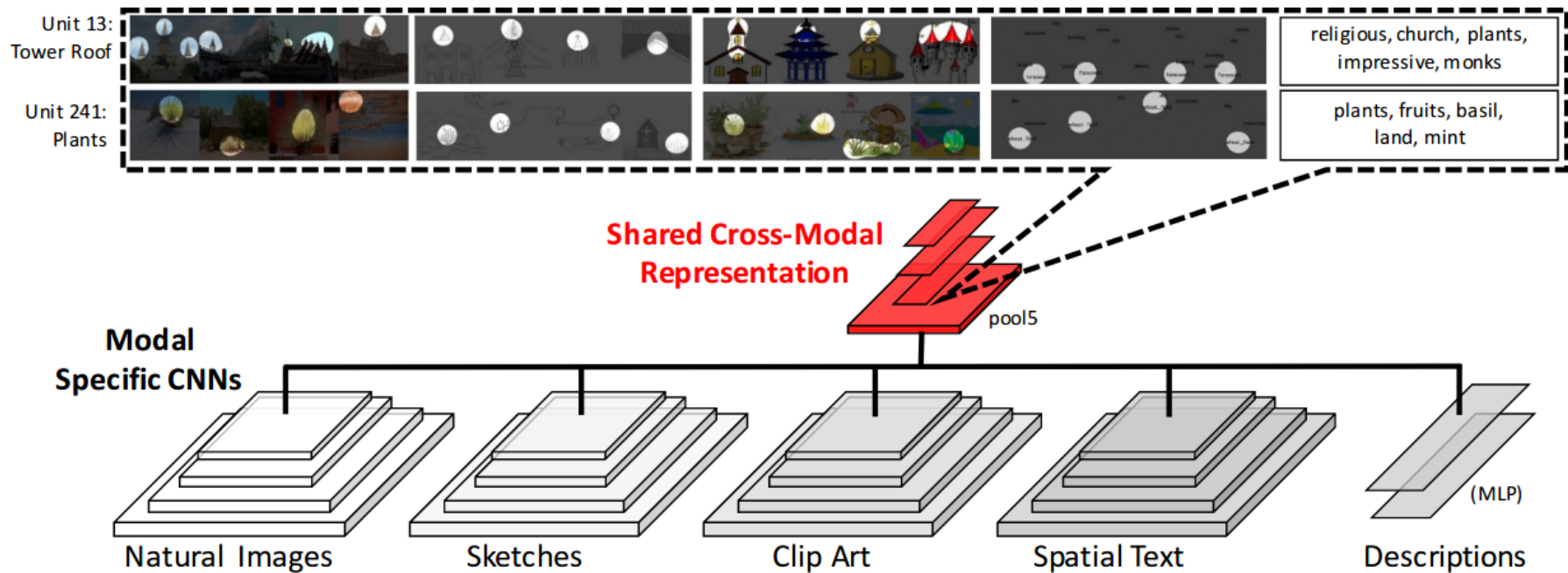


- Loss function $L_{\text{tot}} = \lambda_r L_{\text{regression}} + \lambda_d L_{\text{detection}}$
- Detection task helps regression task

Multi-modal (Cross-modal) Learning

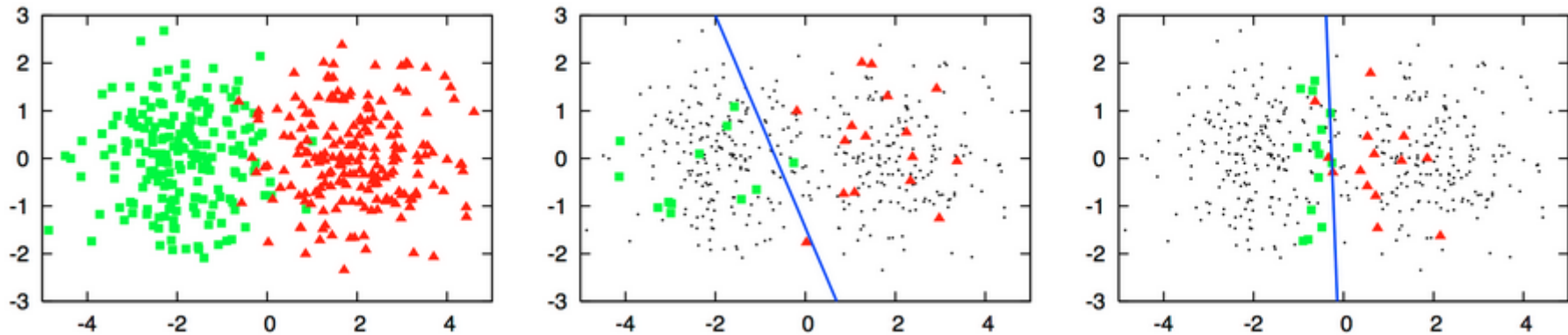
- Several input types sharing a single task

3

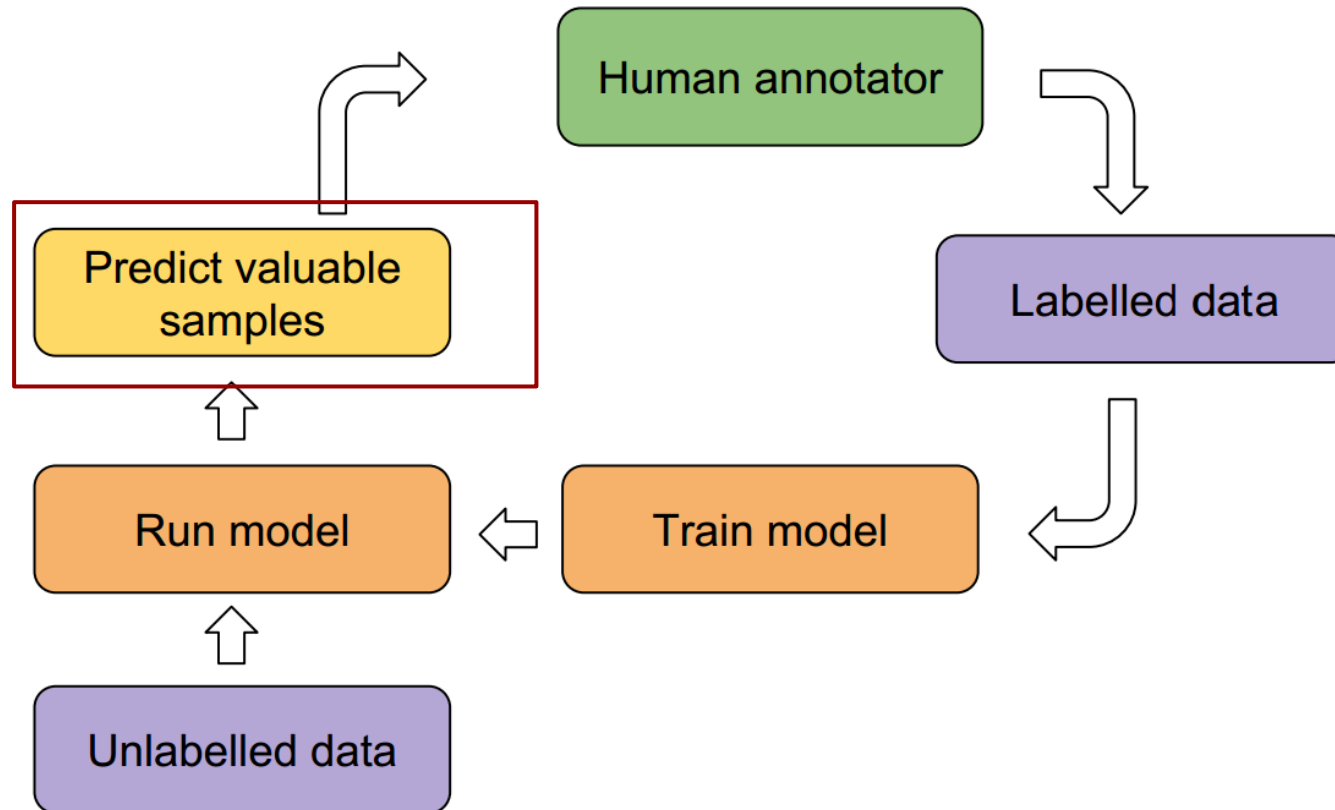


Active Learning

- Actively choose the samples in training



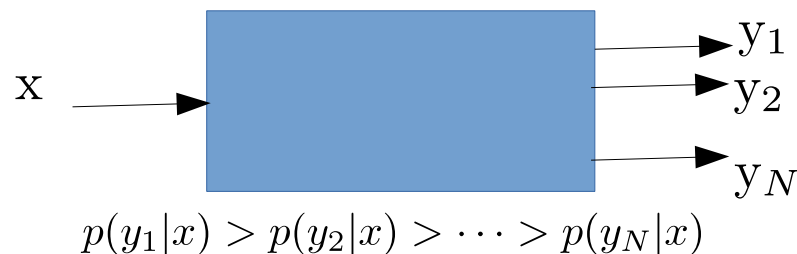
Active Learning Typical Procedure



- The most important element of active learning is to predict which samples are the most “valuable” ones

Predicting the “Valuable” Samples

- Learner can learn more from samples that are difficult to classify now.
- Measures to identify difficult-to-classify samples
 - Confidence $x_{\text{lowest confidence}} = \arg \min_x p(y_1|x)$
 - Classification margin $x_{\text{lowest margin}} = \arg \min_x |p(y_1|x) - p(y_2|x)|$
 - Entropy $x_{\text{highest entropy}} = \arg \max_x - \sum_i p(y_i|x) \log p(y_i|x)$
 - Variance
 - Bayesian techniques (eg: Monte-Carlo dropout, ensembles)



Self-supervised Learning

- Related to Transfer learning
- How to get a pre-trained model?
 - Alt1: Use an existing big labeled database (eg: ImageNet) and supervised learning
 - Alt2: Use own unlabeled data and self-supervised learning
- Supervised vs Self-supervised learning

- Supervised:

$$\text{loss}(D) = \min_{\theta} \frac{1}{N} \sum_{i=1}^N \text{loss}(X_i, Y_i).$$

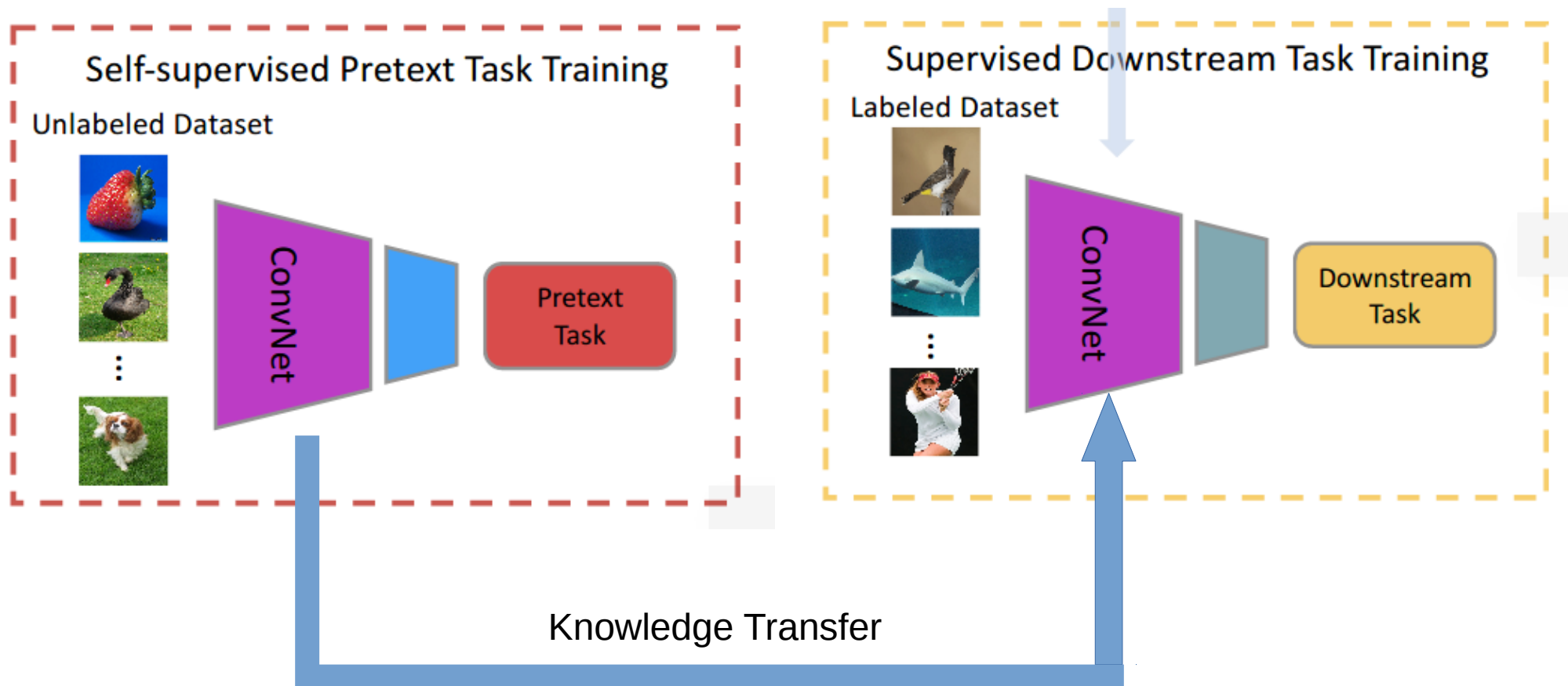
Manually created labels

- Self-Supervised

$$\text{loss}(D) = \min_{\theta} \frac{1}{N} \sum_{i=1}^N \text{loss}(X_i, P_i).$$

Automatically created labels

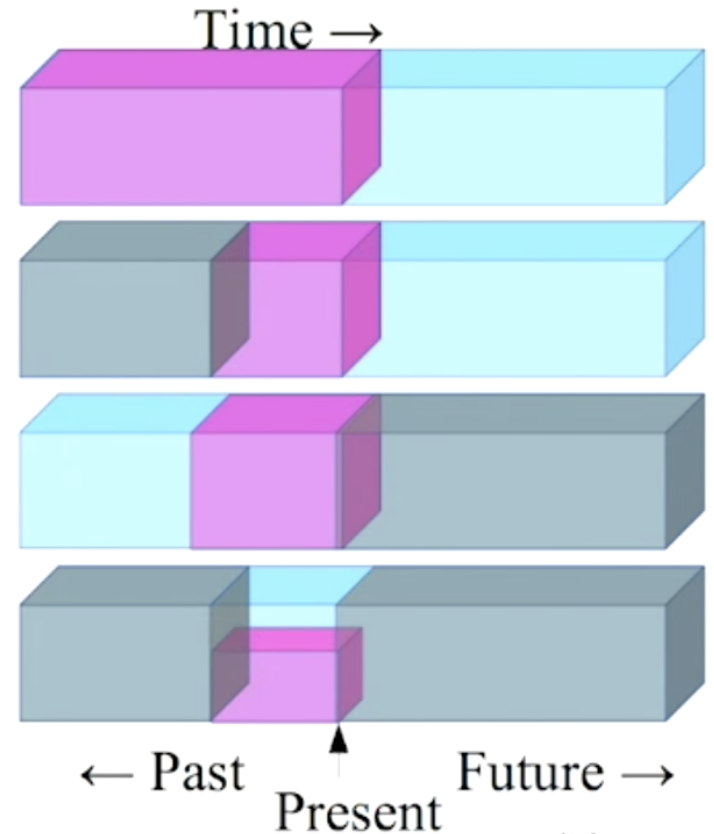
Pretext and Downstream Tasks



Pretext task is defined so that automatic generation of labels is possible

Pretext Task Approaches

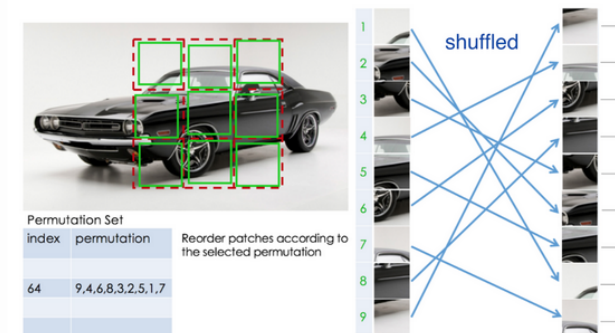
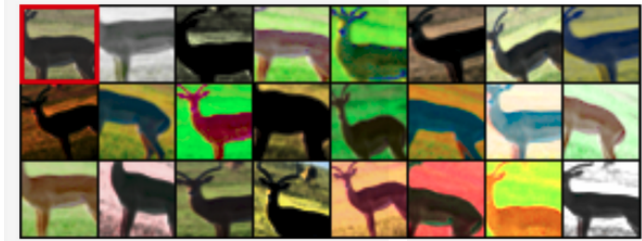
- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the occluded from the visible
- ▶ **Pretend there is a part of the input you don't know and predict that.**



Slide: LeCun

Examples

- Eg 1 (ULMFiT):
 - Predict the next word given the previous words (pretext)
 - Predict sentiments for given sentences (down-stream)
- Eg 2: (BERT)
 - Predict a randomly masked word in a sentence (Pretext)
 - Question answering, language inference etc. (down-stream)
- Eg 3: (Exemplar-CNN)
 - Predict the class of distorted images (Pretext)
 - Image classification (down-stream)
- Eg 4: (Jigsaw Puzzle)
 - Predict the correct order of image patches extracted using a 3x3 grid. (Pretext)
 - Image retrieval (Down-stream)
- Eg 5: (Colorization)
 - Predict the Color given the grayscale image (Pretext)
 - Image classification, detection and segmentation (down-stream)
- And many more



<https://lilianweng.github.io/lil-log/2019/11/10/self-supervised-learning.html>

<https://amitnss.com/2020/05/self-supervised-learning-nlp/>

Discriminative Unsupervised Feature Learning with Exemplar Convolutional Neural Networks, A. Dosovitskiy et.al , 2015

Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles, Mehdi Noroozi, Paolo Favaro, 2017

Metric Learning

- Learn a distance measure (metric) which reflects the semantic similarity
- In practice, we learn a transform $f(\cdot)$ such that the Euclidean distance reflects the similarity.

$$\|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|^2 = [f(\mathbf{x}_i) - f(\mathbf{x}_j)]^T [f(\mathbf{x}_i) - f(\mathbf{x}_j)]$$

- “Similarity” is expressed through side information
 - Eg1: x_i and x_j are similar if they can be linked (i.e. same class) and vice versa
 - Eg2: x_a is more similar to x_p than x_n if a and p can be linked while a and n cannot be linked

- A loss function can be defined for extracting similarity. Examples:

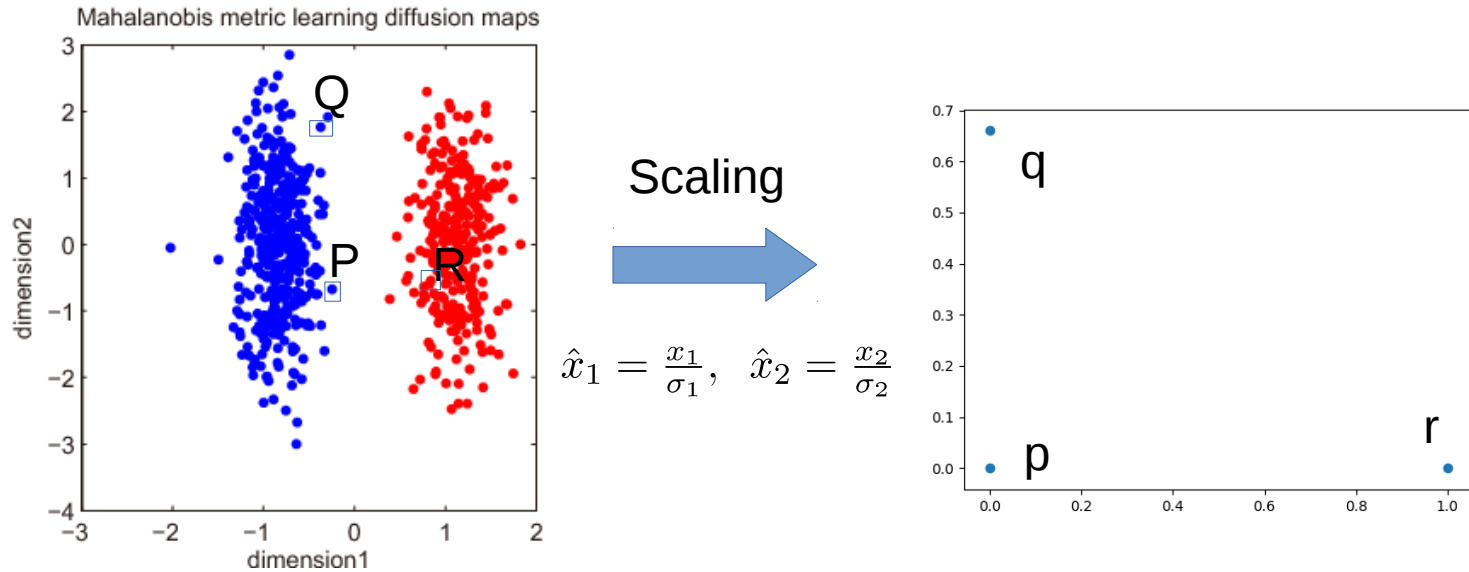
- Contrastive loss:

$$L_c = \begin{cases} \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|^2, & \text{if } \mathbf{x}_i, \mathbf{x}_j \in P \\ -\|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|^2, & \text{if } \mathbf{x}_i, \mathbf{x}_j \in N \end{cases}$$

- Triplet loss

$$L_t = \sum_{\mathbf{x}_i \in A, \mathbf{x}_j \in P, \mathbf{x}_k \in N} \|f(\mathbf{x}_i) - f(\mathbf{x}_j)\|^2 - \|f(\mathbf{x}_i) - f(\mathbf{x}_k)\|^2$$

Mahalanobis Metric (example)



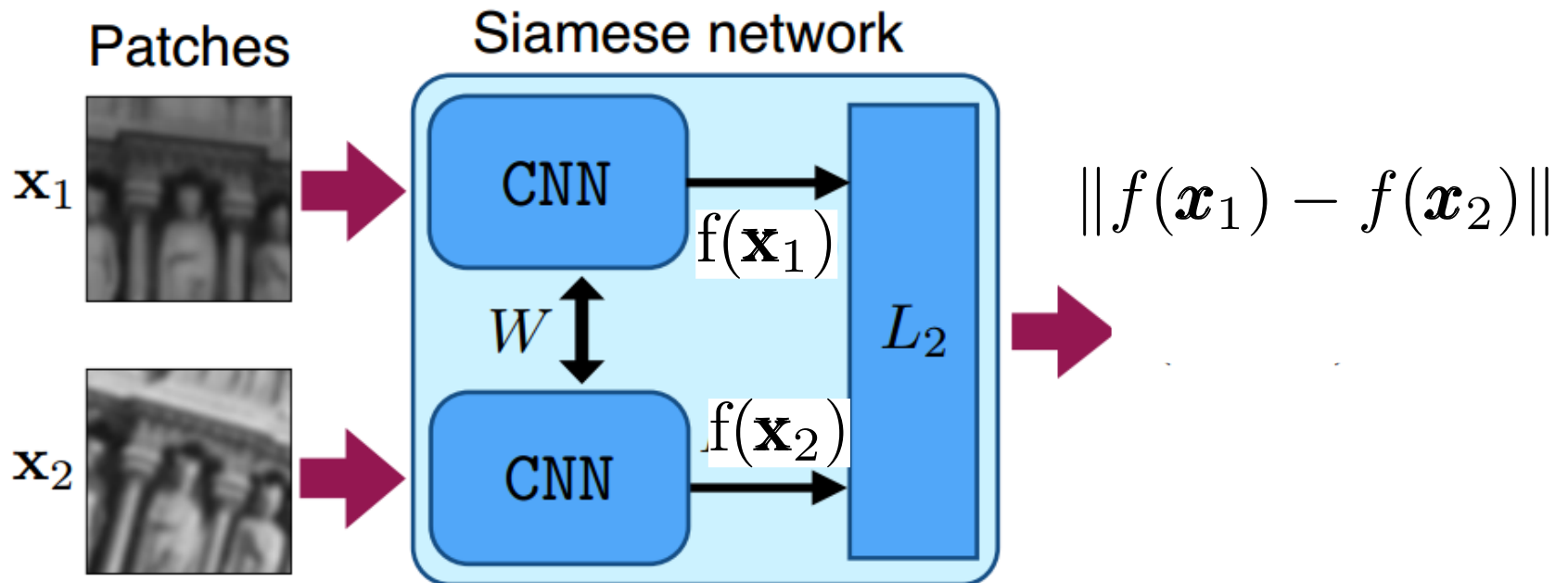
- Consider P, Q and R
 - Q should be more similar to P than R (P and Q belongs to the same class)
 - But Euclidean distance $d(P,Q) > d(P,R)$,
 - Correct this by scaling variance along each dimension

➔ Mahalanobis distance/metric

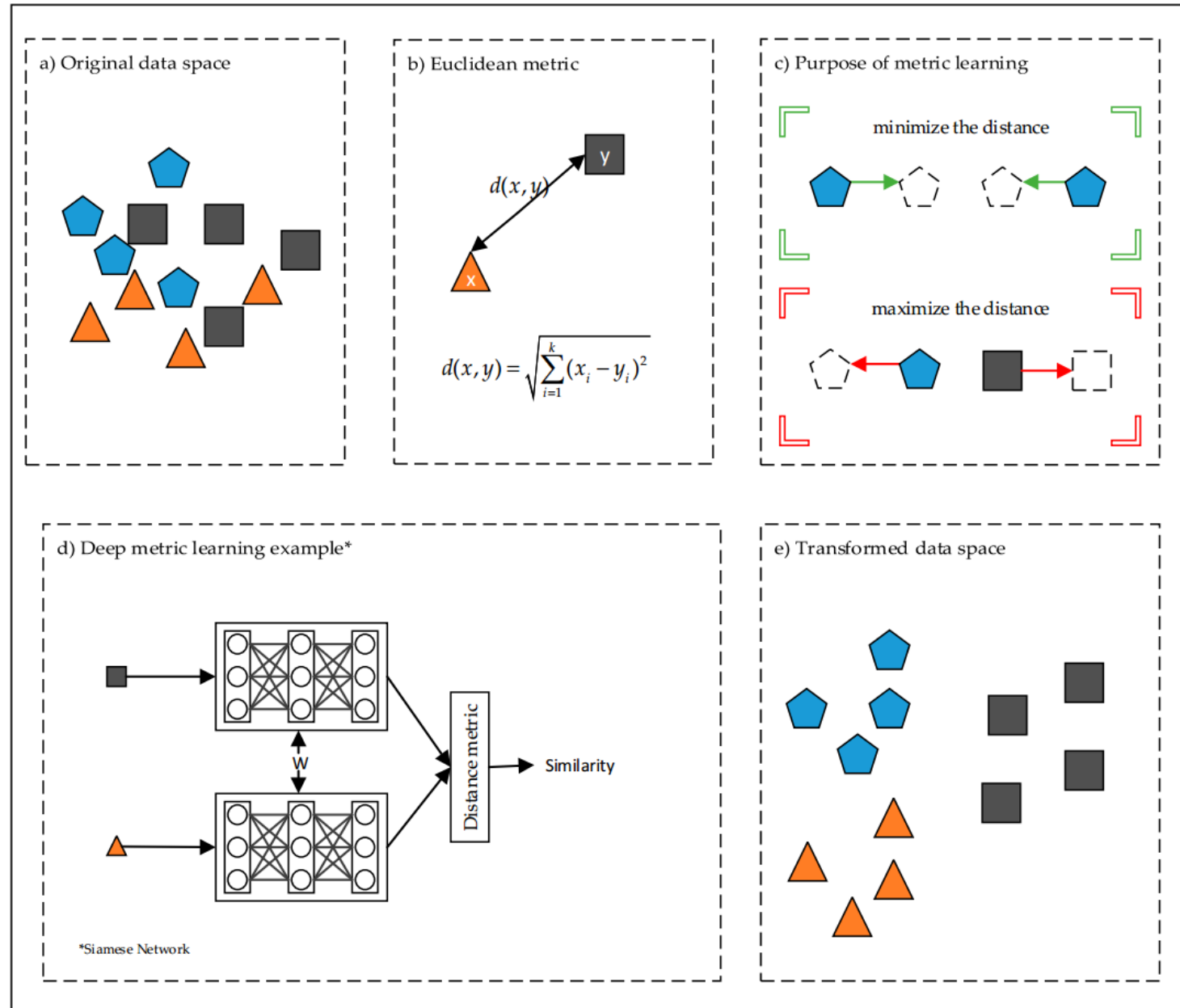
$$\begin{aligned}
 D_m(P, Q) &= \left\| \left[\frac{x_P(1)}{\sigma_1}, \frac{x_P(2)}{\sigma_2} \right] - \left[\frac{x_Q(1)}{\sigma_1}, \frac{x_Q(2)}{\sigma_2} \right] \right\| \\
 &= \left\| [x_P(1), x_P(2)] \begin{bmatrix} \frac{1}{\sigma_1} & 0 \\ 0 & \frac{1}{\sigma_2} \end{bmatrix} - [x_Q(1), x_Q(2)] \begin{bmatrix} \frac{1}{\sigma_1} & 0 \\ 0 & \frac{1}{\sigma_2} \end{bmatrix} \right\| \\
 &= \left\| \mathbf{x}_P^T \Sigma^{-\frac{1}{2}} - \mathbf{x}_Q^T \Sigma^{-\frac{1}{2}} \right\|
 \end{aligned}$$

Metric Learning in Deep Learning

- Replace the linear transform with a neural net and train using a suitable loss



Metric Learning Summary



Meta Learning

- Also known as “Learning to Learn”
- It is about learning what a system learns in different tasks
- Then the meta learner can predict the system parameters for a new task (without going through full training).

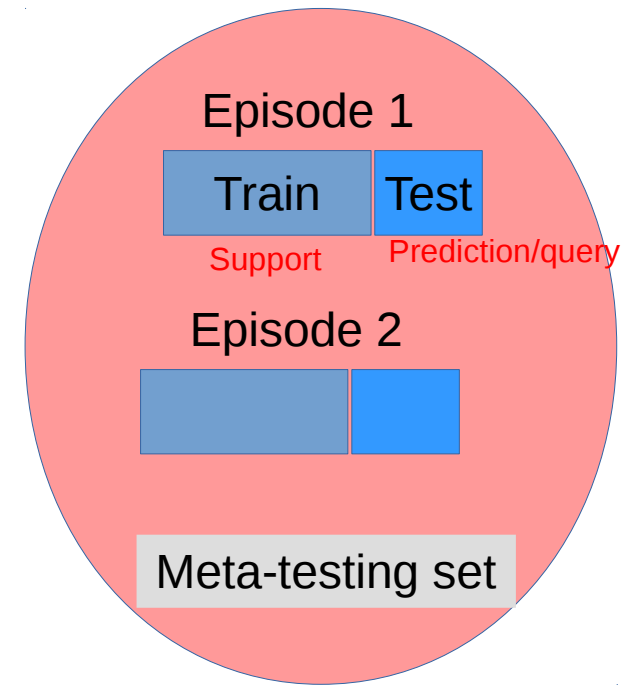
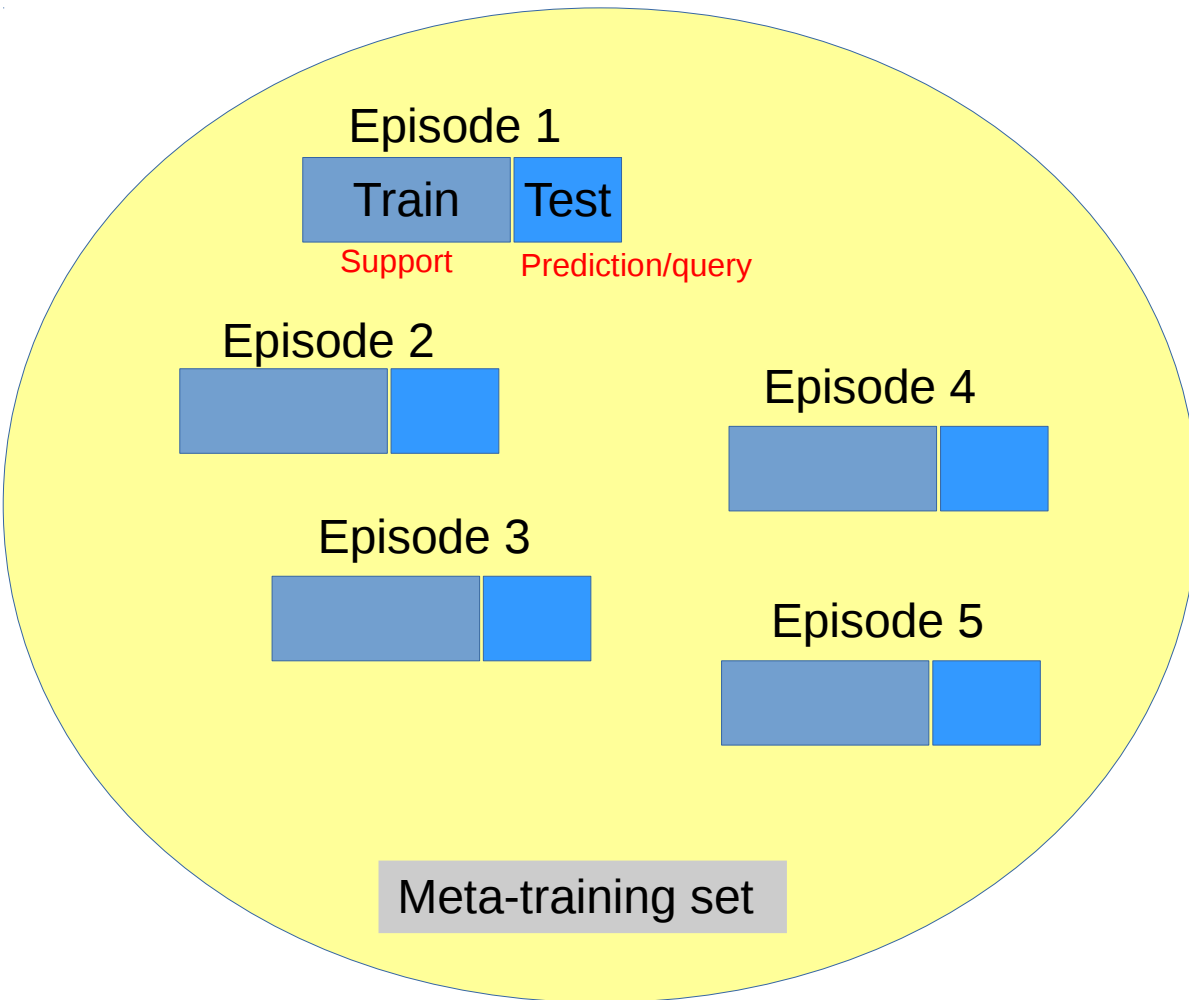


Learning vs Meta-learning

- **Learning algorithm (Learner)**
 - **Input:** Training set $D_{\text{train}} = \{(\mathbf{x}_i, \mathbf{y}_i)\}$
 - **Output:** Parameters θ of the model M
 - **Objective:** Performance on the test set $D_{\text{test}} = \{(\hat{\mathbf{x}}_i, \hat{\mathbf{y}}_i)\}$

- **Meta-Learning algorithm (Meta-learner)**
 - **Input:** Meta-training set $D_{\text{meta-train}} = \{(D_{\text{train}}^{(n)}, D_{\text{test}}^{(n)})\}$
 - **Output:** Parameters Θ of the learner
 - **Objective:** Performance on the meta-test set $D_{\text{meta-test}} = \{(\hat{D}_{\text{train}}^{(m)}, \hat{D}_{\text{test}}^{(m)})\}$

Terminology



Optimization Objectives

- Formulation 1

$$\theta^* = \arg \max_{\theta} \frac{1}{N} \sum_{n=1}^N \sum_{\mathbf{x}, \mathbf{y} \in D_{\text{test}}^{(n)}} p_{\theta}(\mathbf{y} | \mathbf{x}, D_{\text{train}}^{(n)})$$

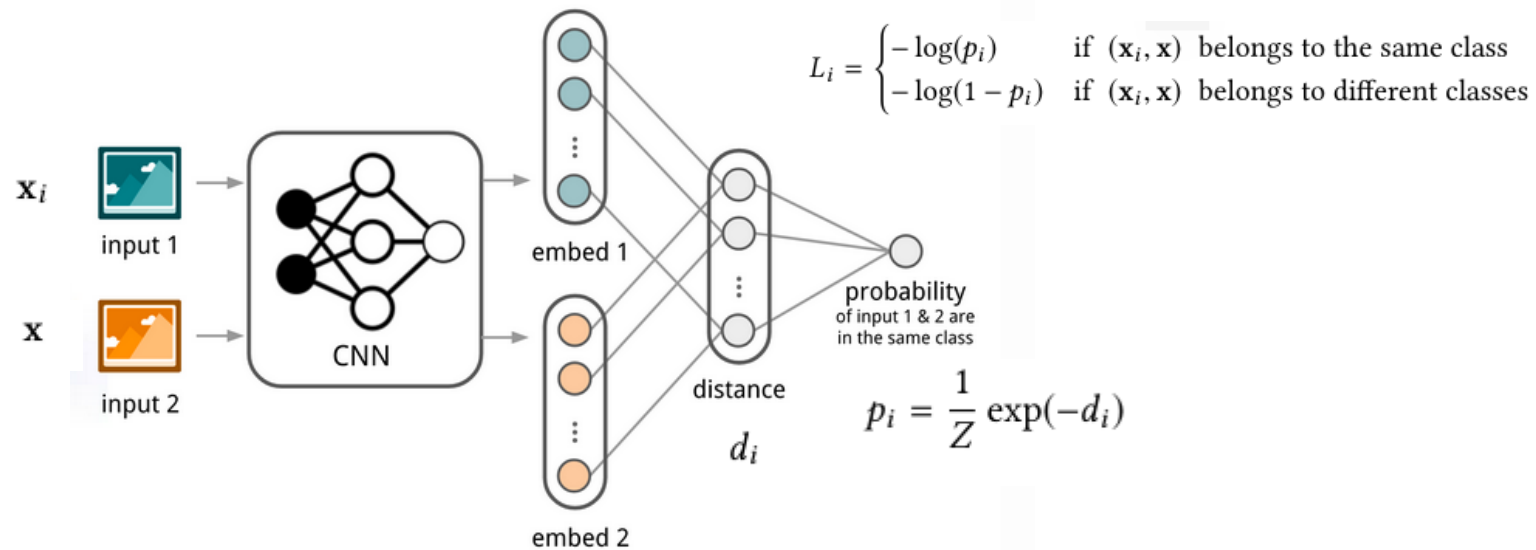
- Formulation 2

$$[\theta^*, \phi^*] = \arg \max_{\theta, \phi} \frac{1}{N} \sum_{n=1}^N \sum_{\mathbf{x}, \mathbf{y} \in D_{\text{test}}^{(n)}} p_{g_{\phi}(\theta, D_{\text{train}}^{(n)})}(\mathbf{y} | \mathbf{x})$$

Possible Approaches

- Model dependent approaches
 - Metric based
 - Siamese networks
 - Matching network
 - Relation network
 - Prototypical network
 - RNN based
 - Memory augmented neural network (MANN)
 - LSTM meta learner
- Model independent approaches
 - Optimization based
 - Model Agnostic Meta-learning (MAML)
 - Reptile

Siamese networks



- Training

- Meta-training set

- $D_{train} = \{(x_1, c_1), (x_2, c_2), (x_3, c_3), \dots, (x_N, c_N)\}$
 - $D_{test} = \{x, c\}$

- Make image pairs $\{(x_1, x), (x_2, x), \dots, (x_N, x)\}$

- Trained to optimize $L = \sum_i L_i$

- Testing

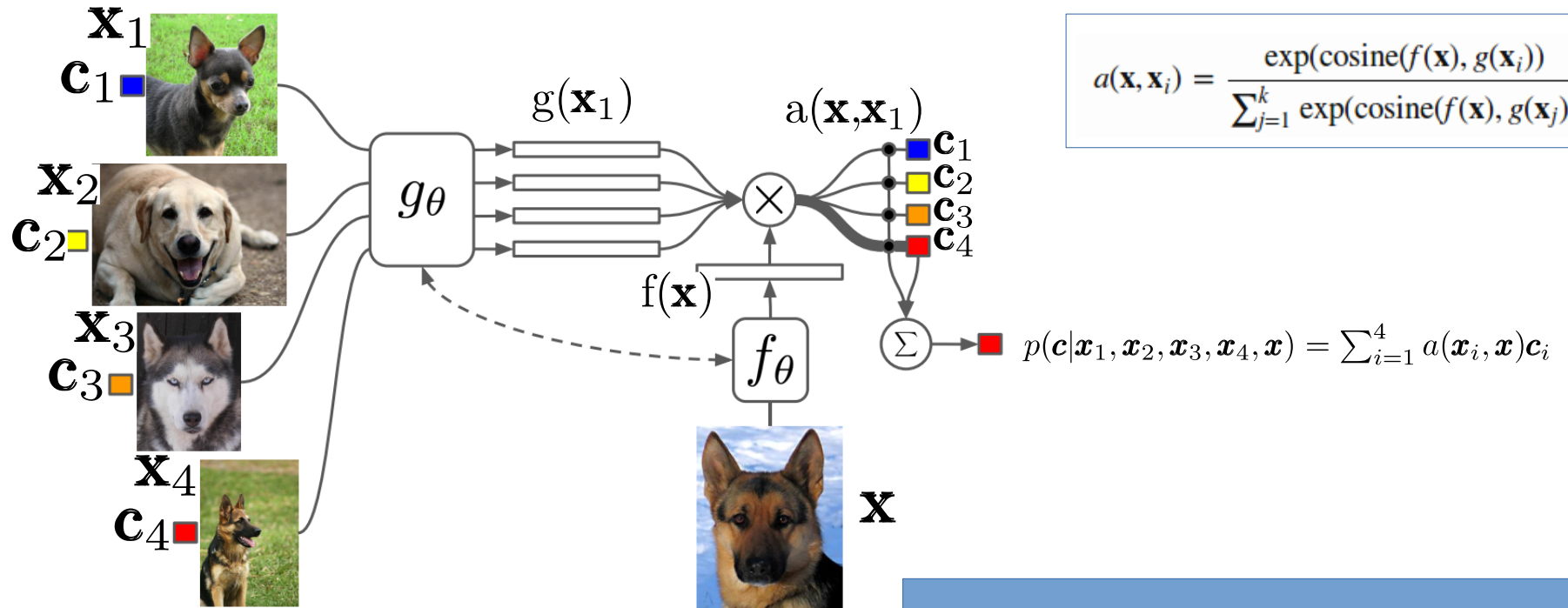
- Meta-testing set

- $\hat{D}_{train} = \{(\hat{x}_1, \hat{c}_1), (\hat{x}_2, \hat{c}_2), (\hat{x}_3, \hat{c}_3), \dots, (\hat{x}_N, \hat{c}_N)\}$
 - and a test image \hat{x}

- Make a set of image pairs $\{(\hat{x}_1, \hat{x}), (\hat{x}_2, \hat{x}), \dots, (\hat{x}_N, \hat{x})\}$

- Feed each image pair and register the output probability. Class of support set image which outputs the highest probability is the class of the test image. $c = \arg \max_i p(\hat{x}, \hat{x}_i)$

Matching network



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

Note: Labels are in one-hot encoding format

- Training

- Meta-training set

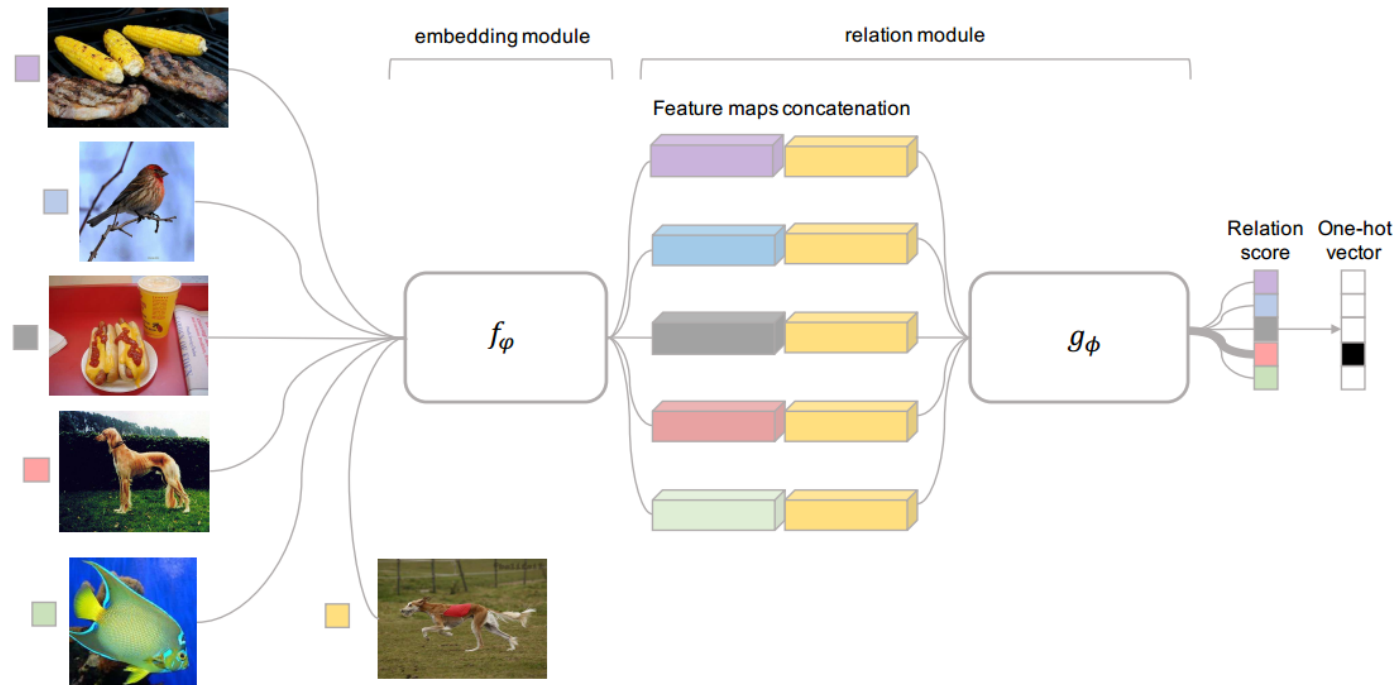
- $D_{\text{train}} = \{(\mathbf{x}_1, \mathbf{c}_1), (\mathbf{x}_2, \mathbf{c}_2), (\mathbf{x}_3, \mathbf{c}_3), (\mathbf{x}_4, \mathbf{c}_4)\}$
 - $D_{\text{test}} = \{(\mathbf{x}, \mathbf{c})\}$
 - Train to maximize $p(\mathbf{c} | \mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}) = \sum_{i=1}^4 a(\mathbf{x}_i, \mathbf{x}) \mathbf{c}_i$

- Testing

- Meta-testing set

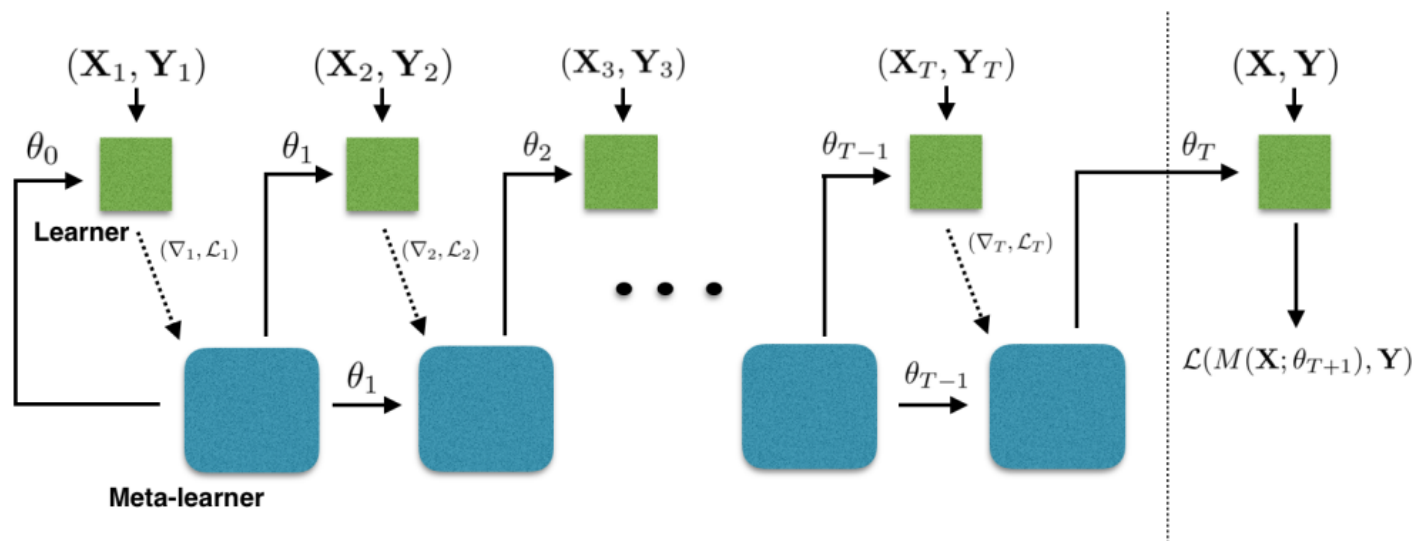
- $\hat{D}_{\text{train}} = \{(\hat{\mathbf{x}}_1, \hat{\mathbf{c}}_1), (\hat{\mathbf{x}}_2, \hat{\mathbf{c}}_2), (\hat{\mathbf{x}}_3, \hat{\mathbf{c}}_3), (\hat{\mathbf{x}}_4, \hat{\mathbf{c}}_4)\}$
 - $\hat{D}_{\text{test}} = \{(\hat{\mathbf{x}}, \hat{\mathbf{c}})\}$
 - Find $\arg \max_k c_k \in p(\hat{\mathbf{c}} | \hat{\mathbf{x}}, \hat{\mathbf{x}}_k, k = 1, 2, 3, 4)$

Relation network



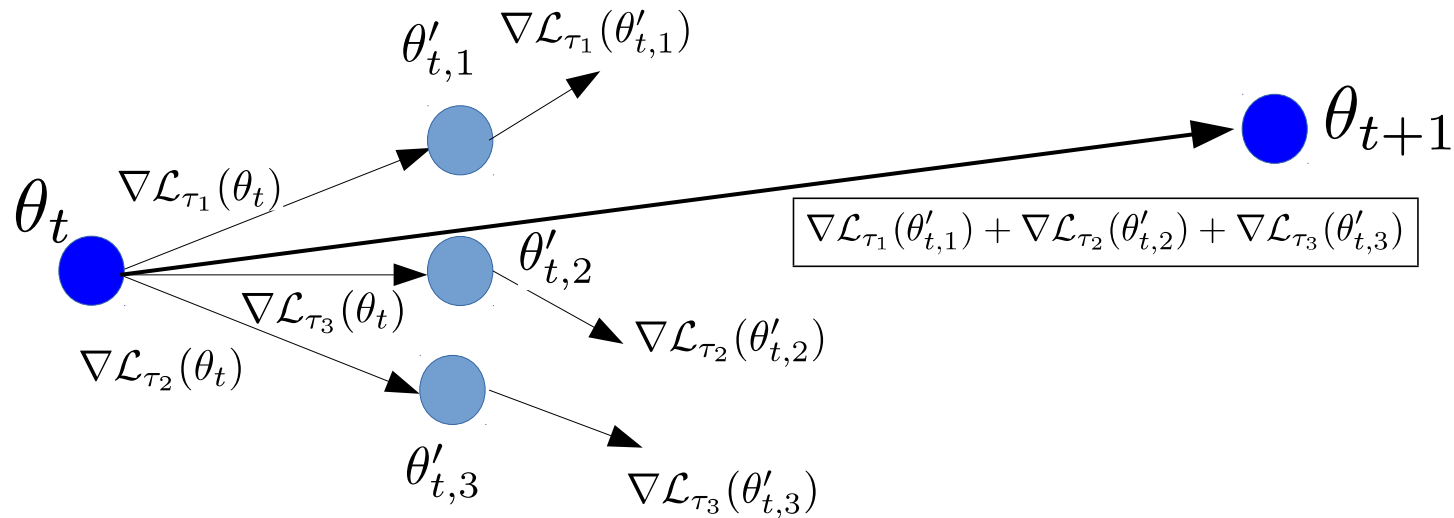
- Idea is similar to matching network
- Similarity is learned with a different network architecture

LSTM Meta-learner



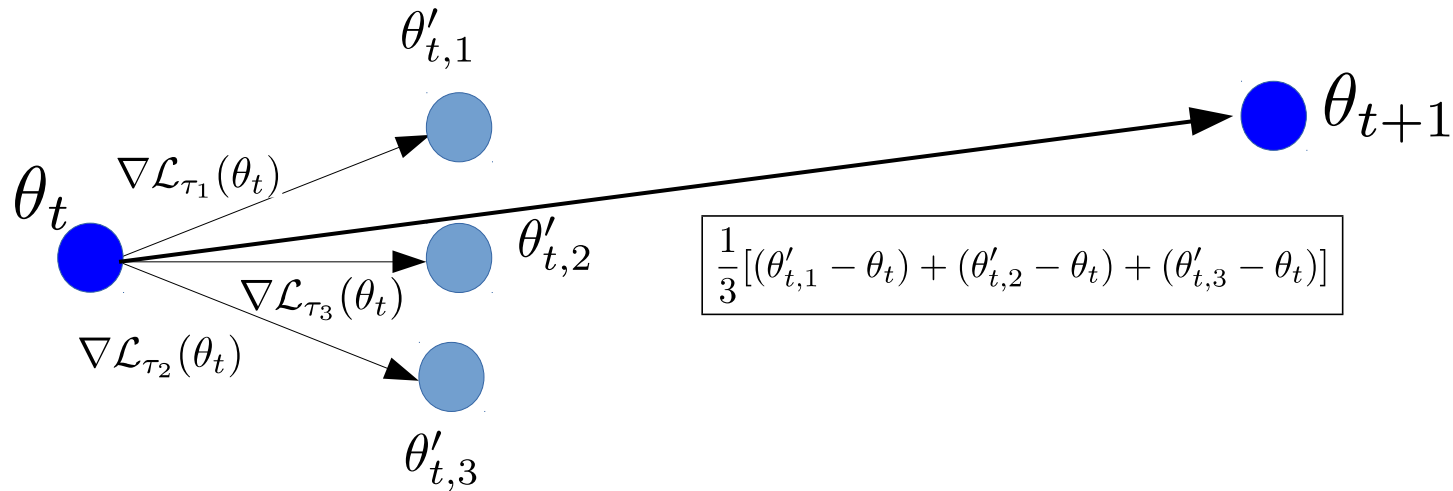
- Learner is a neural network with parameters θ
- Meta-learner is an LSTM with parameters ϕ
- Data pairs from D_{train} are fed sequentially
- At each iteration, meta-learner delivers better parameter set θ_{t+1} given previous parameters θ_t , learner loss \mathcal{L}_t and its gradient $\nabla \mathcal{L}_t$.
- Finally feed data from D_{test} and find a loss
- Loss is back-propagated through the both learner and meta-learner and both θ and ϕ can be updated.

Model Agnostic Meta Learning (MAML)



- Optimization approach related to gradient descent
- Can be applied to any model
- Main idea:
 - Divide the meta-training set D_{train} into a set of tasks $\tau_i, i = 1, 2, \dots, N$
 - Update the model parameters, using the average of gradients over the tasks
 - The updated model parameter set is kept close to optimum model parameter set for each individual task.

Reptile



- Very similar to MAML
- Averaging strategy is different
- Idea is to minimize the Euclidean distance between θ_i and $\theta'_{t,i}$ $i = 1, 2, \dots, N$

$$\nabla \sum_{i=1}^N (\theta_t - \theta'_{t,i})^2 = 2 \sum_{i=1}^N (\theta_t - \theta'_{t,i})$$

One/Few Shot Learning

- Given a support set of one/few training examples of new classes, learn to classify a test input
 - K-class (k-way) one shot learning
 - one example per class and k classes
 - K-class (k-way) n-shot learning
 - n examples per class and k classes
- A straight-forward application of meta learning.



Support set (Training)



query set (Training)



Support set (Testing)



Test input

Zero Shot Learning

- Example: Given a semantic description, learn to classify an image



Okapi is " zebra-striped four legged animal with a brown torso and a deer-like face". Which of these images is of Okapi?

Zero-shot Learning Problem

- Given

- Training set $D_{tr} = \{(\mathbf{x}_i, c_i^{tr}), i = 1, 2, \dots, n\}$

- Training example is \mathbf{x}_i

- Corresponding class label is c_i^{tr}

- Class label is drawn from the **seen** training set classes $C_{tr} = \{c_i^{tr} | i = 1, 2, \dots, n_C^{tr}\}$

- **Unseen** test set classes $C_{te} = \{c_i^{te} | i = 1, 2, \dots, n_C^{te}\}$

- Semantic/auxiliary representation vectors for each class in both training and test set $\{\mathbf{v}_i^{tr} | i = 1, 2, \dots, n_C^{tr}\} \cup \{\mathbf{v}_j^{te} | j = 1, 2, \dots, n_C^{te}\}$

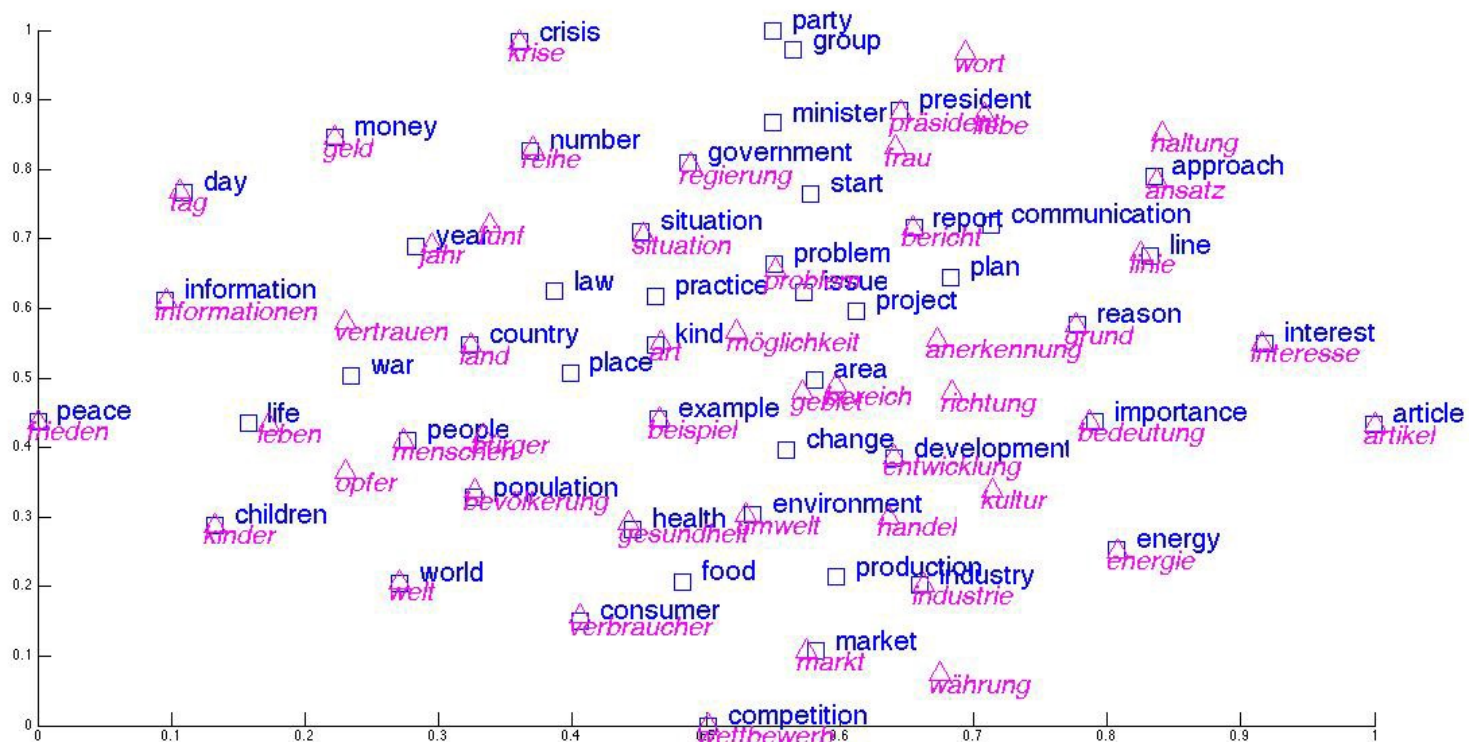
- Training and testing classes are disjoint $C_{tr} \cap C_{te} = \emptyset$

- Find

- A function which maps test input examples into the corresponding class $f : \mathbf{x}_i^{te} \rightarrow c_i^{te}$

Semantic Representation

- Attribute description:
 - Example
 - Attributes: small, cute, furry, horns
 - Dog=[1,1,1,0], Bull=[0,1,1,1]
- Word co-occurrence count
- Word embedding vectors
 - Trained vector representations for each class name



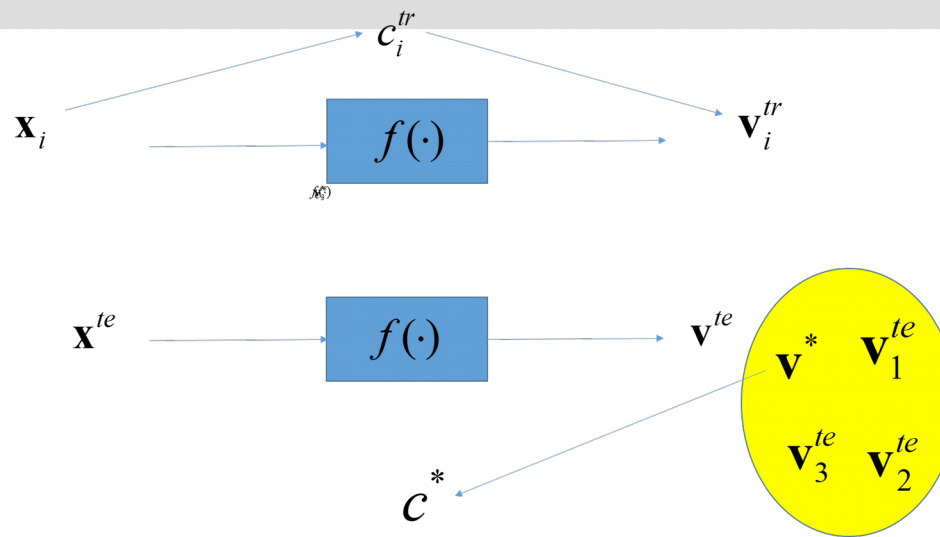
Zero-shot Learning Approaches-1

- Training:

- For each training example (\mathbf{x}_i, c_i^{tr}) , get the semantic representation vector (or category vector) \mathbf{v}_i^{tr} corresponding to the class c_i^{tr} .
- Train a network $f(\cdot)$ to map \mathbf{x}_i to \mathbf{v}_i^{tr} .

- Testing:

- Get the test input \mathbf{x}^{te} and map it to the semantic representation vector $\mathbf{v}^{te} = f(\mathbf{x}^{te})$
- Find the nearest neighbour to \mathbf{v}^{te} and call it \mathbf{v}^*
- Return the corresponding class c^* to \mathbf{v}^*



Zero-shot Learning Approaches-2

- Training:
 - For each training example (\mathbf{x}_i, c_i^{tr}) , get the semantic representation vector (or category vector) \mathbf{v}_i^{tr} corresponding to the class c_i^{tr} .
 - Train a network $f(\cdot, \cdot)$ to map $(\mathbf{x}_i, \mathbf{v}_i^{tr})$ to a similarity value E .
- Testing:
 - Get the test input \mathbf{x}^{te} and couple it with semantic representation vector \mathbf{v}_i^{te} corresponding to each of the possible classes c_i^{te} . This will create $\{(\mathbf{x}^{te}, \mathbf{v}_i^{te}) \mid i = 1, 2, \dots, n_C^{te}\}$
 - For each pair $(\mathbf{x}^{te}, \mathbf{v}_i^{te})$ find $E_i = f(\mathbf{x}^{te}, \mathbf{v}_i^{te})$.
 - Return the corresponding class to maximum E_i

