

# Artificial Intelligence

## The Theory of Learning and Ensemble Learning

1111AI06

MBA, IM, NTPU (M6132) (Fall 2022)

Wed 2, 3, 4 (9:10-12:00) (B8F40)

Min-Yuh Day, Ph.D,  
Associate Professor

Institute of Information Management, National Taipei University

<https://web.ntpu.edu.tw/~myday>



<https://meet.google.com/miy-fbif-max>



# Syllabus

**Week    Date    Subject/Topics**

- |          |                   |   |
|----------|-------------------|---|
| <b>1</b> | <b>2022/09/14</b> | <b>Introduction to Artificial Intelligence</b>  |
| <b>2</b> | <b>2022/09/21</b> | <b>Artificial Intelligence and Intelligent Agents</b>   |
| <b>3</b> | <b>2022/09/28</b> | <b>Problem Solving</b>  |
| <b>4</b> | <b>2022/10/05</b> | <b>Knowledge, Reasoning and Knowledge Representation;<br/>Uncertain Knowledge and Reasoning</b> |
| <b>5</b> | <b>2022/10/12</b> | <b>Case Study on Artificial Intelligence I</b>  |
| <b>6</b> | <b>2022/10/19</b> | <b>Machine Learning: Supervised and Unsupervised Learning</b>                                   |

# Syllabus

**Week    Date    Subject/Topics**

**7   2022/10/26   The Theory of Learning and Ensemble Learning**

**8   2022/11/02   Midterm Project Report**

**9   2022/11/09   Deep Learning and Reinforcement Learning**

**10   2022/11/16   Deep Learning for Natural Language Processing**

**11   2022/11/23   Invited Talk: AI for Information Retrieval**

**12   2022/11/30   Case Study on Artificial Intelligence II**

# Syllabus

**Week    Date    Subject/Topics**

**13    2022/12/07    Computer Vision and Robotics**

**14    2022/12/14    Philosophy and Ethics of AI and the Future of AI**

**15    2022/12/21    Final Project Report I**

**16    2022/12/28    Final Project Report II**

**17    2023/01/04    Self-learning**

**18    2023/01/11    Self-learning**

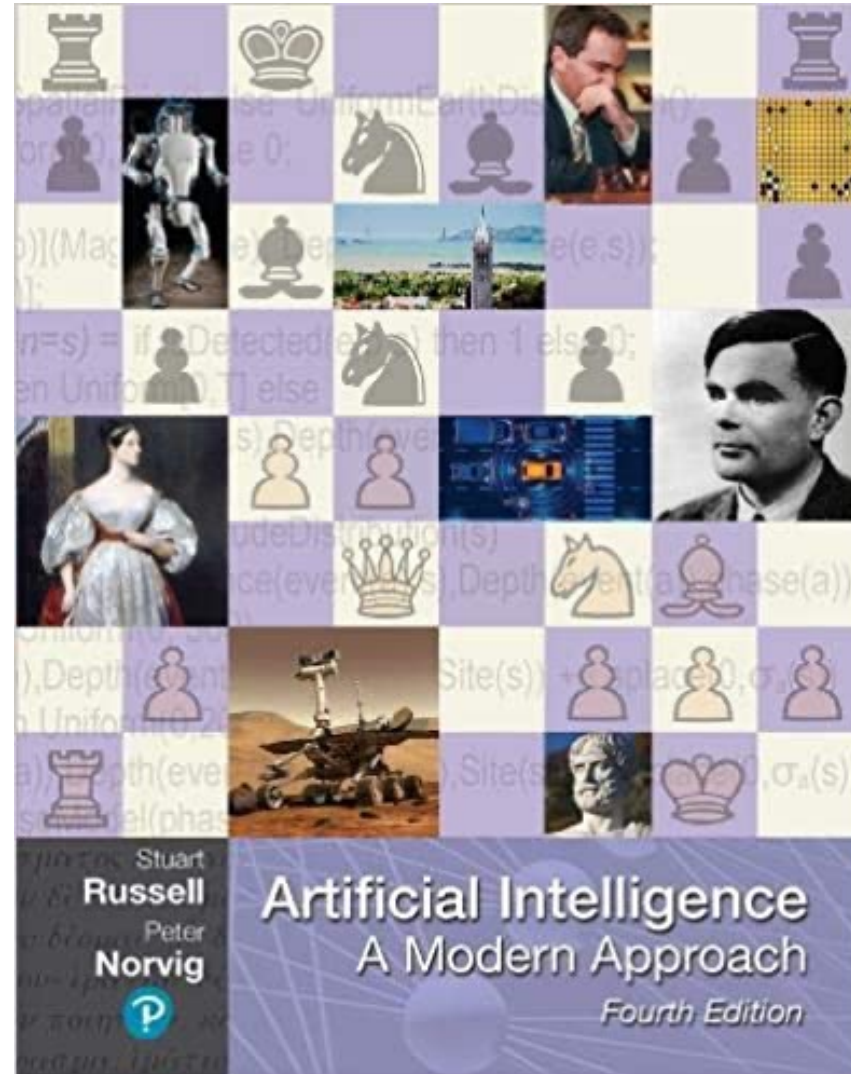


# **The Theory of Learning and Ensemble Learning**

# Outline

- **The Theory of Learning**
  - **Computational Learning Theory**
  - **Probably Approximately Correct (PAC) Learning**
- **Ensemble Learning**
  - **Bagging: Random forests**
  - **Stacking**
  - **Boosting: Gradient boosting**
  - **Online learning**
- **Meta Learning: Learning to Learn**

Stuart Russell and Peter Norvig (2020),  
**Artificial Intelligence: A Modern Approach,**  
4th Edition, Pearson



Source: Stuart Russell and Peter Norvig (2020), Artificial Intelligence: A Modern Approach, 4th Edition, Pearson

<https://www.amazon.com/Artificial-Intelligence-A-Modern-Approach/dp/0134610997/>

# Artificial Intelligence: A Modern Approach

1. Artificial Intelligence
2. Problem Solving
3. Knowledge and Reasoning
4. Uncertain Knowledge and Reasoning
5. Machine Learning
6. Communicating, Perceiving, and Acting
7. Philosophy and Ethics of AI

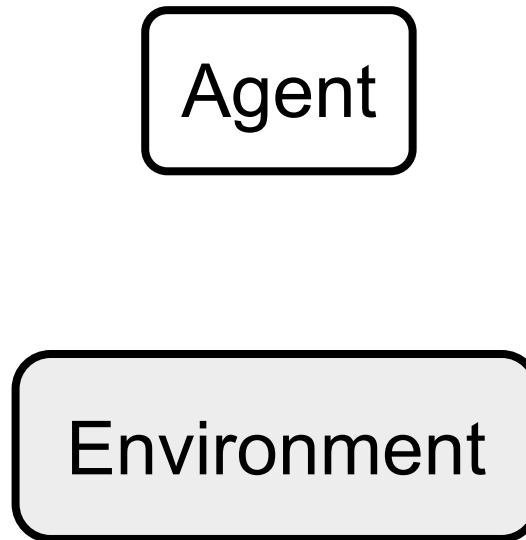
# Artificial Intelligence: Machine Learning

# Artificial Intelligence:

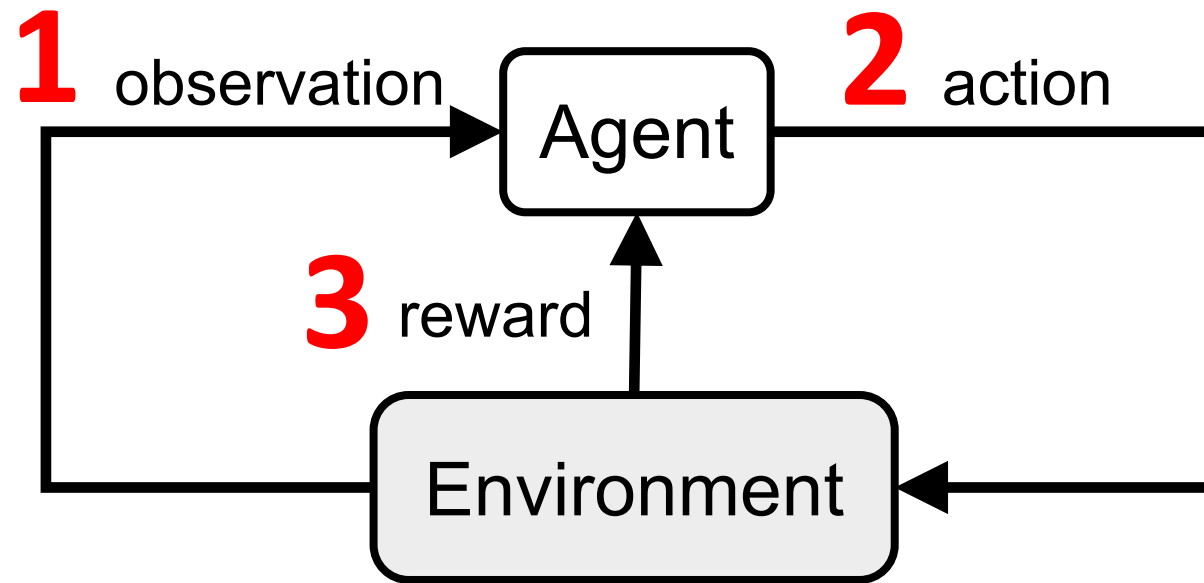
## 5. Machine Learning

- **Learning from Examples**
- **Learning Probabilistic Models**
- **Deep Learning**
- **Reinforcement Learning**

# Reinforcement Learning (DL)

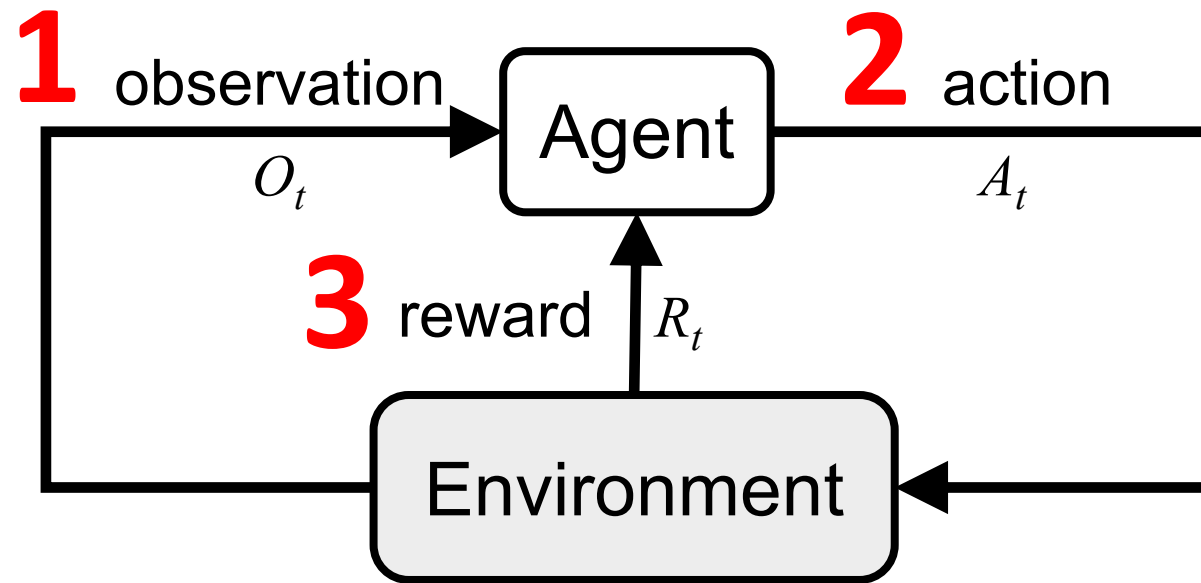


# Reinforcement Learning (DL)

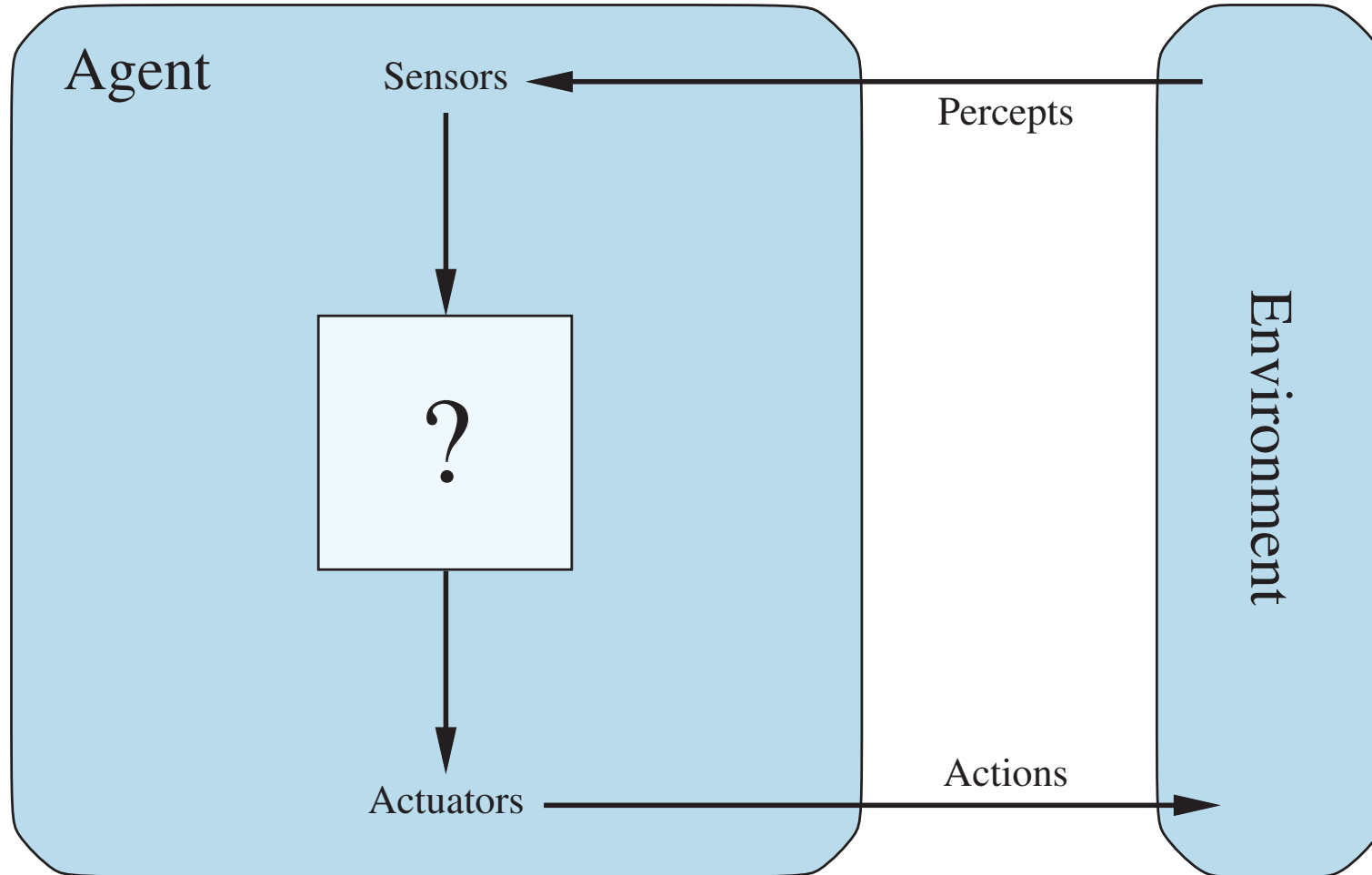




# Reinforcement Learning (DL)



# Agents interact with environments through sensors and actuators

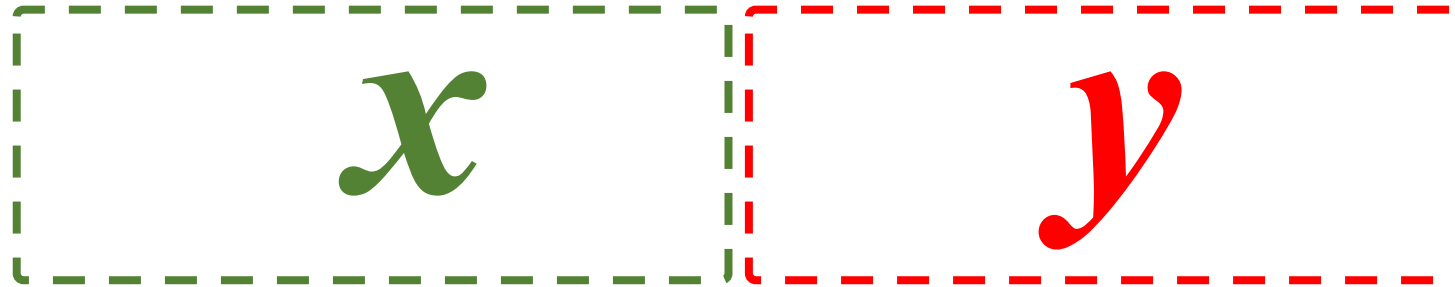


# Machine Learning

## Supervised Learning (Classification)

### Learning from Examples

$$y = f(x)$$

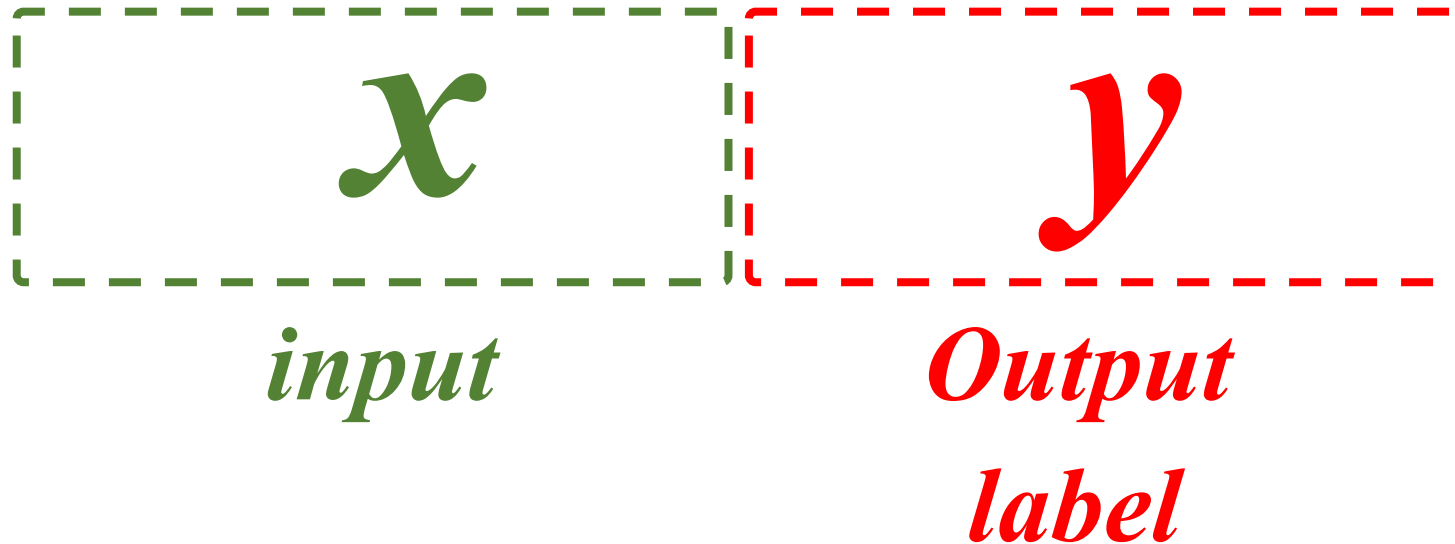


# Machine Learning

## Supervised Learning (Classification)

### Learning from Examples

$$y = f(x)$$



# Machine Learning

## Supervised Learning (Classification)

### Learning from Examples

$$\textcolor{red}{y} = f(\textcolor{green}{x})$$

```
5.1, 3.5, 1.4, 0.2, Iris-setosa
4.9, 3.0, 1.4, 0.2, Iris-setosa
4.7, 3.2, 1.3, 0.2, Iris-setosa
7.0, 3.2, 4.7, 1.4, Iris-versicolor
6.4, 3.2, 4.5, 1.5, Iris-versicolor
6.9, 3.1, 4.9, 1.5, Iris-versicolor
6.3, 3.3, 6.0, 2.5, Iris-virginica
5.8, 2.7, 5.1, 1.9, Iris-virginica
7.1, 3.0, 5.9, 2.1, Iris-virginica
```

# Machine Learning

## Supervised Learning (Classification)

### Learning from Examples

$$\textcolor{red}{y} = f(\textcolor{green}{x})$$

<i>Example</i>	5.1	3.5	1.4	0.2	Iris-setosa
	4.9	3.0	1.4	0.2	Iris-setosa
	4.7	3.2	1.3	0.2	Iris-setosa
	7.0	3.2	4.7	1.4	Iris-versicolor
	6.4	3.2	4.5	1.5	Iris-versicolor
	6.9	3.1	4.9	1.5	Iris-versicolor
	6.3	3.3	6.0	2.5	Iris-virginica
	5.8	2.7	5.1	1.9	Iris-virginica
	7.1	3.0	5.9	2.1	Iris-virginica

# Machine Learning

## Supervised Learning (Classification)

### Learning from Examples

$$y = f(x)$$

*Example*

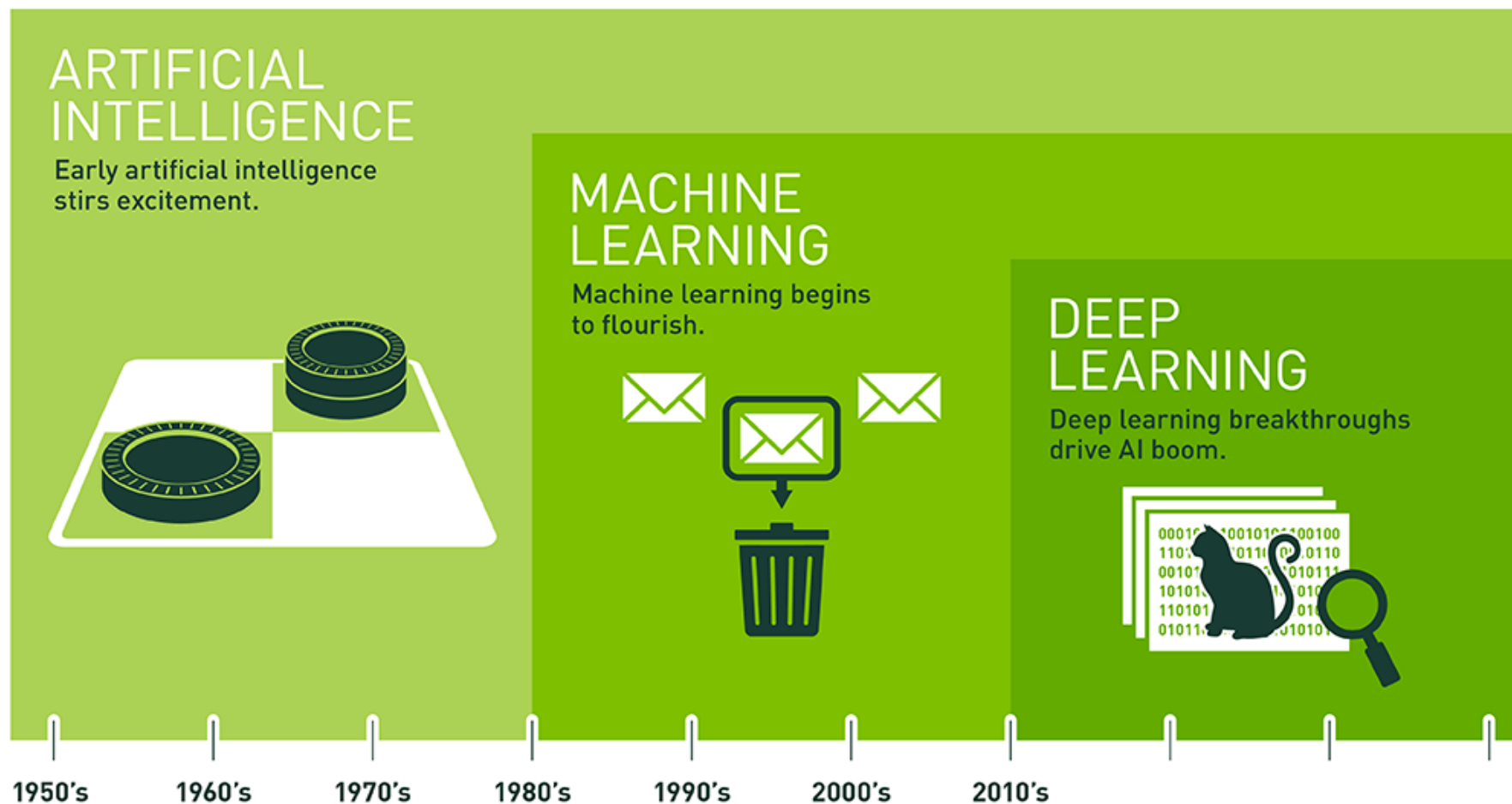
$x$

5.1, 3.5, 1.4, 0.2	Iris-setosa
4.9, 3.0, 1.4, 0.2	Iris-setosa
4.7, 3.2, 1.3, 0.2	Iris-setosa
7.0, 3.2, 4.7, 1.4	Iris-versicolor
6.4, 3.2, 4.5, 1.5	Iris-versicolor
6.9, 3.1, 4.9, 1.5	Iris-versicolor
6.3, 3.3, 6.0, 2.5	Iris-virginica
5.8, 2.7, 5.1, 1.9	Iris-virginica
7.1, 3.0, 5.9, 2.1	Iris-virginica

$y$

# Artificial Intelligence

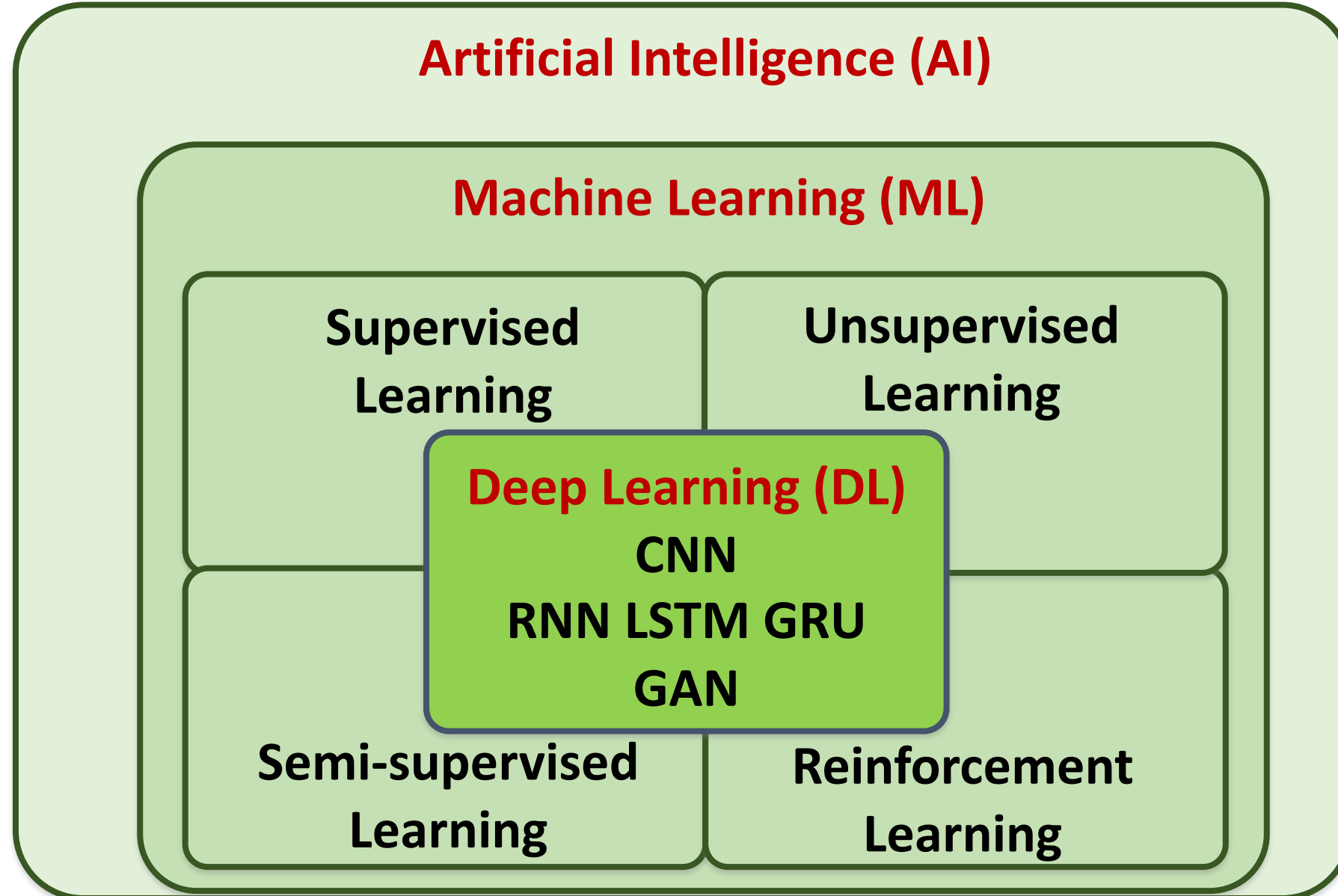
## Machine Learning & Deep Learning



Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.



# AI, ML, DL



# The Theory of Learning

# The Theory of Learning

- **Computational Learning Theory**
- **Probably approximately correct (PAC)**

# The Theory of Learning

- How can we be sure that our **learned hypothesis** will **predict** well for previously unseen inputs?
  - How do we know that the **hypothesis  $h$**  is close to the **target function  $f$**  if we don't know what is?
- How many **examples** do we need to get a good  **$h$** ?
- What **hypothesis space** should we use?
- If the hypothesis space is very complex, can we even find the best  **$h$**  or do we have to settle for a **local maximum**?
- How **complex** should  **$h$**  be?
- How do we avoid **overfitting**?

# Computational Learning Theory

- Intersection of **AI, statistics, and theoretical computer science.**
- Any **hypothesis** that is seriously wrong will almost certainly be “found out” with high probability after a small number of examples.

# Probably approximately correct (PAC)

- Any **hypothesis** that is consistent with a sufficiently large set of training examples is unlikely to be seriously wrong.
- **PAC learning algorithm:**
  - Any learning algorithm that returns hypotheses that are probably approximately correct.

# Linear function

$$y = f(x)$$

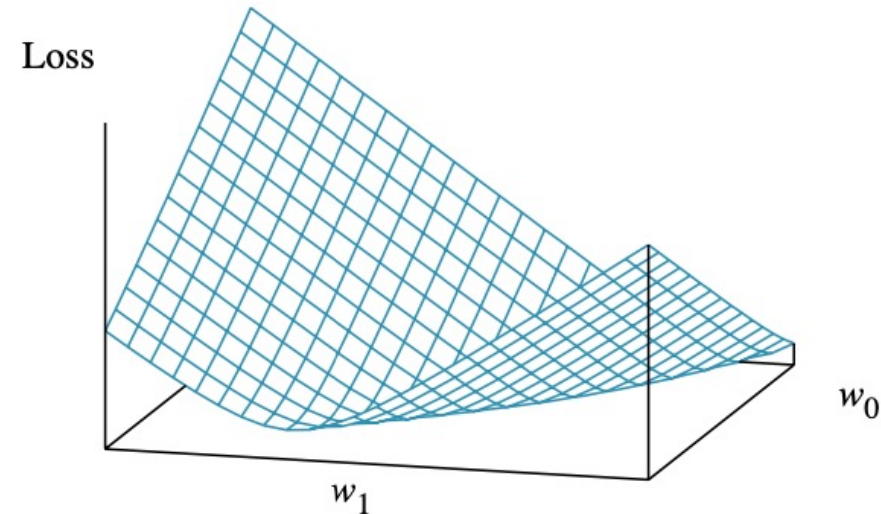
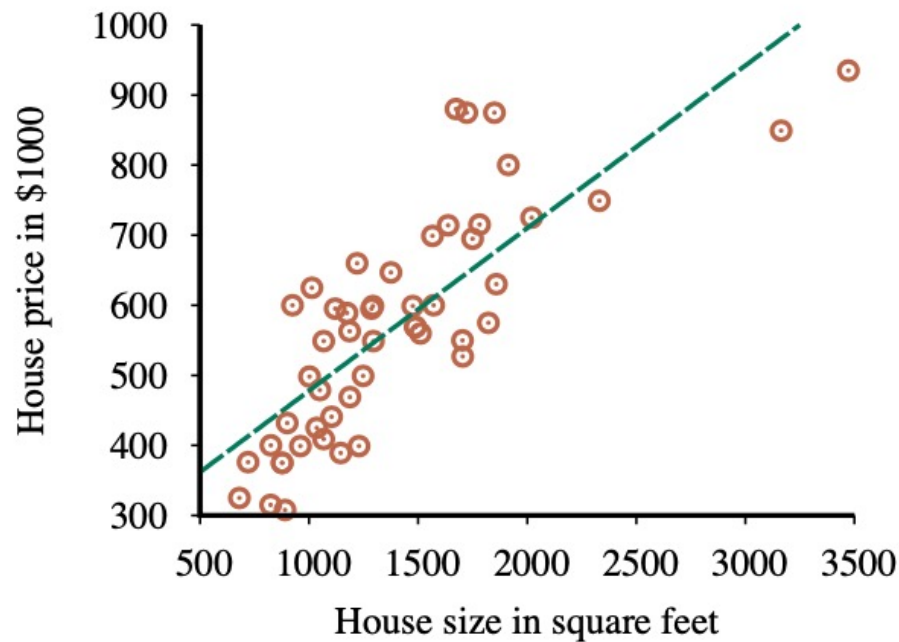
$$y = w_1 x + w_0$$

$$h_w(x) = w_1 x + w_0$$

# Linear Regression Weight Space

$$h_w(x) = w_1 x + w_0$$

$$w^* = \operatorname{argmin}_w \operatorname{Loss}(h_w)$$



$$y = 0.232 x + 246$$

*Loss function for Weights ( $w_1, w_0$ )*



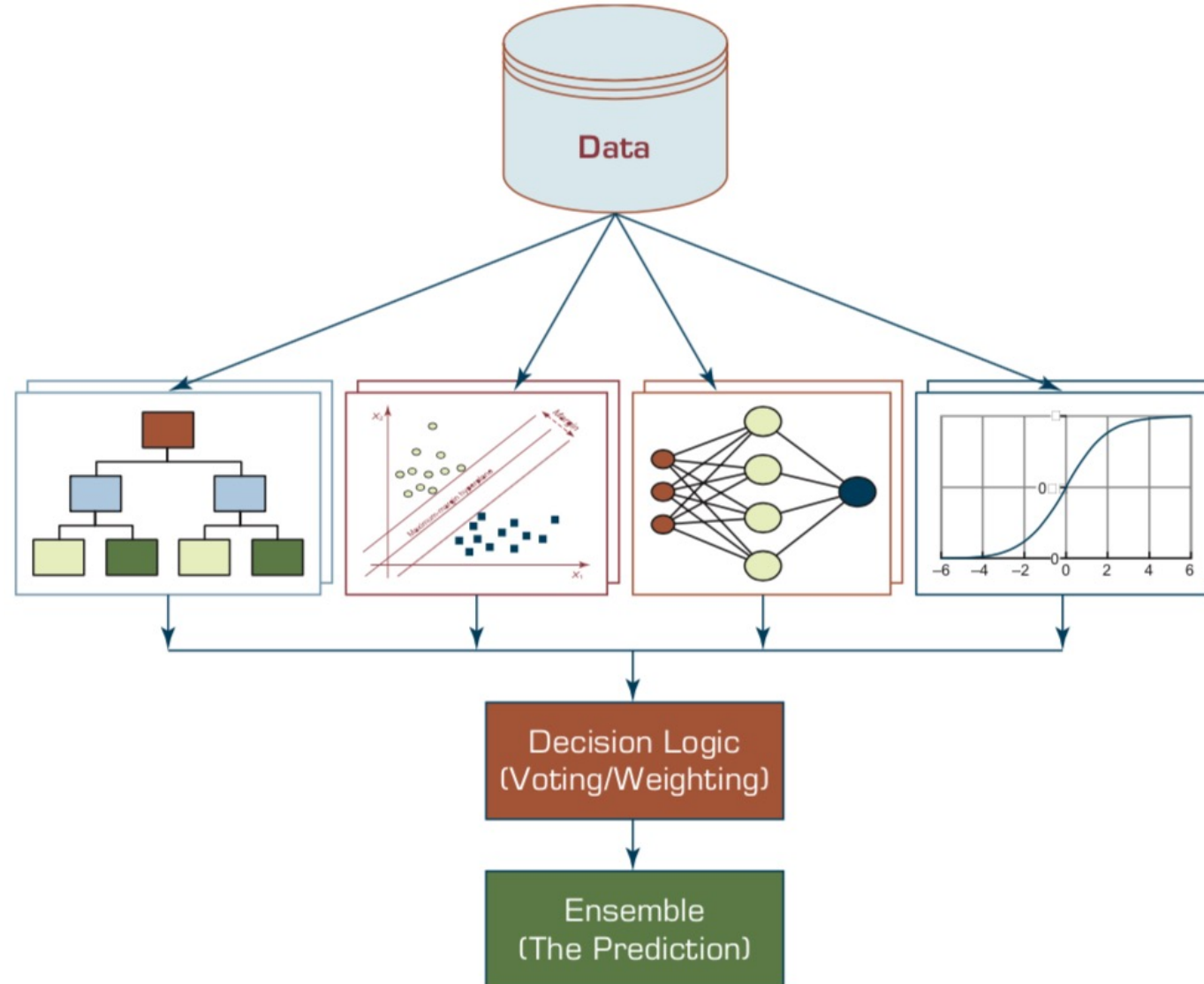
# Ensemble Learning

# Ensemble Learning

- **Select a collection, or ensemble,** of hypotheses,  $h_1, h_2, \dots, h_n$ , **and combine** their predictions **by averaging, voting,** or **by another level of machine learning.**

# Ensemble Models

## Heterogeneous Ensemble



# Ensemble Learning

- **Base model**
  - **individual hypotheses**
    - $h_1, h_2, \dots, h_n$
- **Ensemble model**
  - **hypotheses combination**

# Why Ensemble Learning

- **Reduce bias**
- **Reduce variance**

# Ensemble Learning

- **Bagging**
  - **Random forests**
- **Stacking**
- **Boosting**
  - **Gradient boosting**
- **Online learning**

# Ensemble Learning: Bagging

- **Bagging**
  - **Generate distinct training sets by sampling with replacement from the original training set.**
- **Classification:**
  - **Plurality Vote (Majority Vote)**
- **Regression:**
  - **Average**

# Ensemble Learning: Random forests

- **Random forest** model is a form of **decision tree bagging** in which we take extra steps to make the ensemble of trees more diverse, to reduce variance.
- The key idea is to randomly vary the **attribute** choices (rather than the training examples)



# Ensemble Learning: Random forests

- **Extremely randomized trees (ExtraTrees)**
  - Use randomness in selecting the split point **value**
  - for each selected attribute, we randomly sample several candidate values from a uniform distribution over the attribute's range

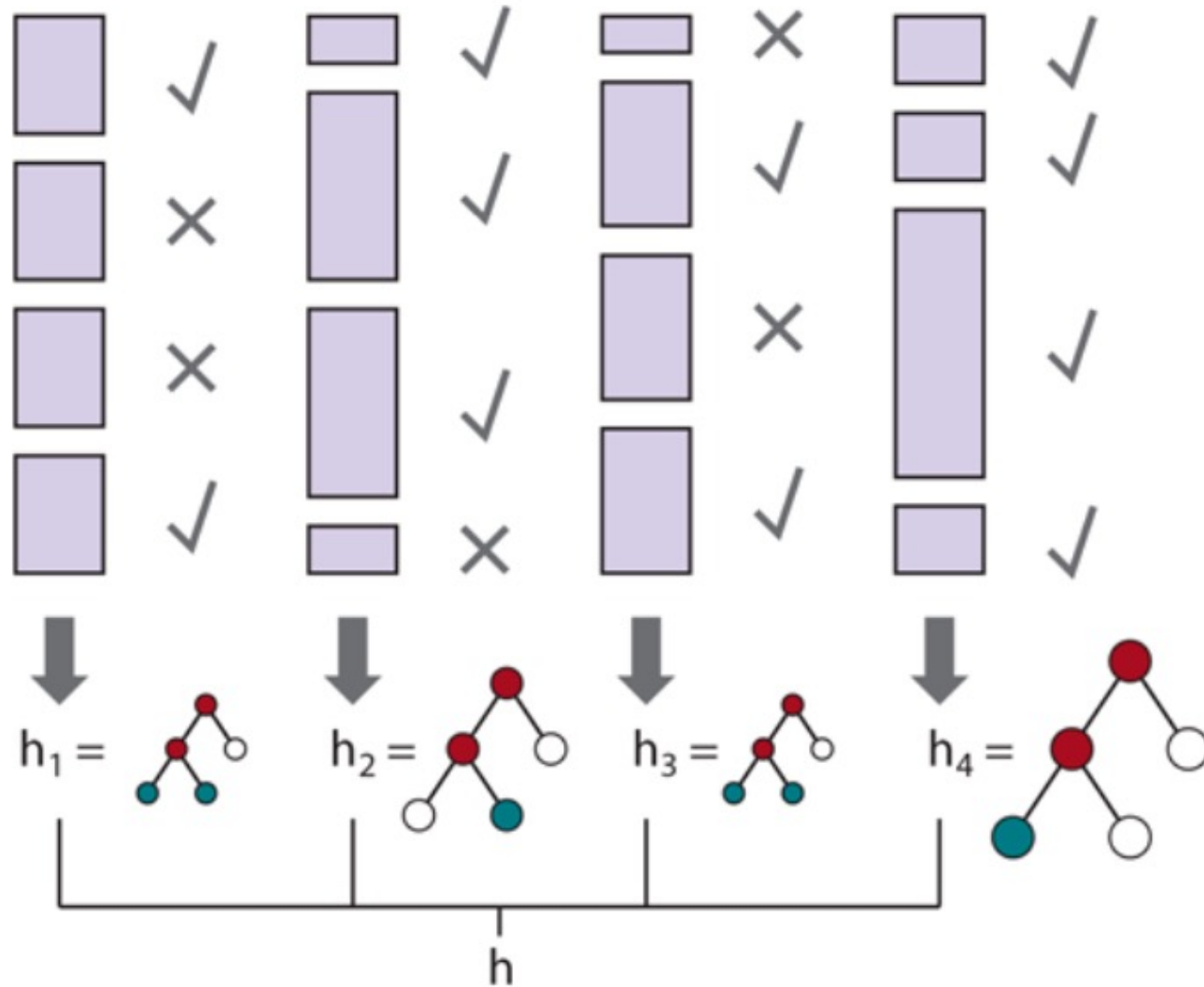
# Ensemble Learning: Stacking

- **Stacking**
  - Stacked generalization combines **multiple base models** from **different model classes** trained on the **same data**.
- **Bagging**
  - Combines multiple base models of the **same model class** trained on **different data**.

# Ensemble Learning: Boosting

- **Boosting**
  - The most popular ensemble method
- **Weighted training set**

# Ensemble Learning: Boosting



# Ensemble Learning:

## Gradient boosting

- **Gradient boosting**
  - Gradient boosting is a form of boosting using gradient descent
- **Gradient boosting machines (GBM)**
- **Gradient boosted regression trees (GBRT)**
- **Popular method for regression and classification of factored tabular data**

# Ensemble Learning:

## Online learning

- **Online learning**
  - **Data are not i.i.d.**  
**(independent and identically distributed)**
  - **An agent receives an input  $x_i$  from nature, predicts the corresponding  $y_i$  and then is told the correct answer.**

# Meta Learning: Learning to Learn

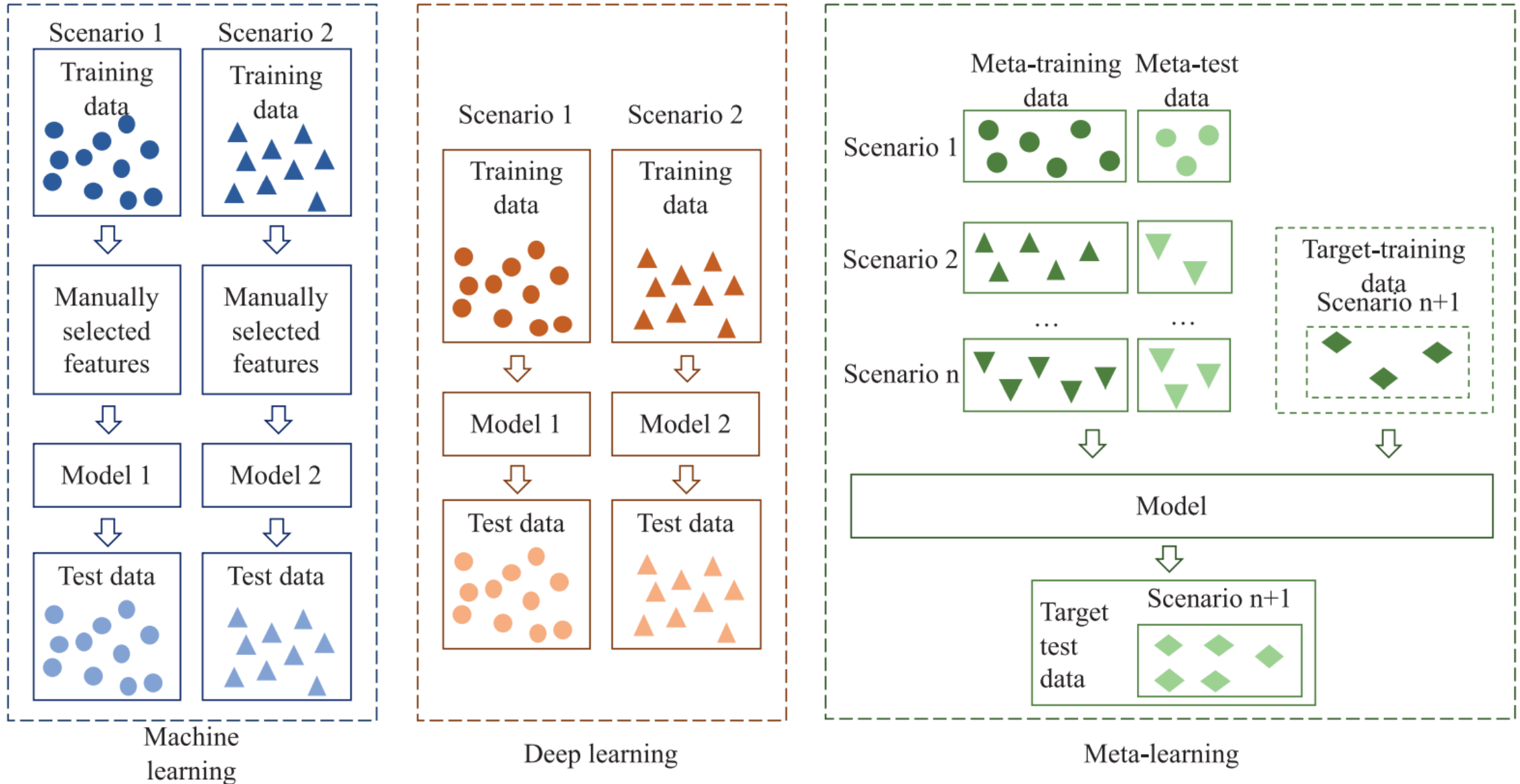
**Deep Learning**  
**Transfer Learning**  
**Few-Shot Learning**  
**Meta Learning**



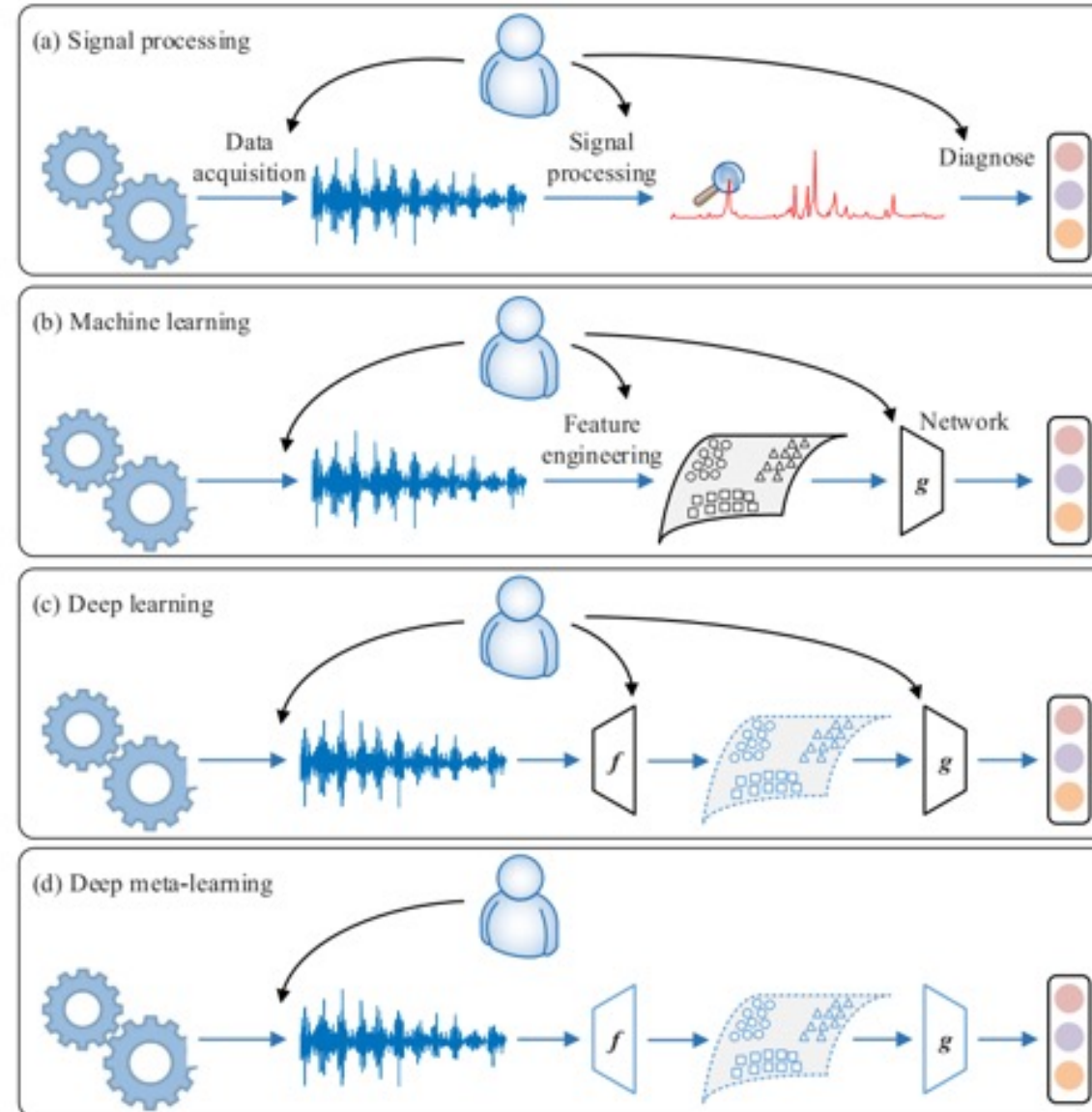
# Deep Learning, Transfer Learning, Few-Shot Learning, Meta Learning

- Deep Learning
  - Transfer Learning
    - Pre-training, Fine-Tuning (FT)
- Meta Learning: Learning to Learn
- Few-Shot Learning (FSL)
- One-Shot Learning (1SL)
- Zero-Shot Learning (0SL)(ZSL)

# Machine Learning, Deep Learning, Meta Learning

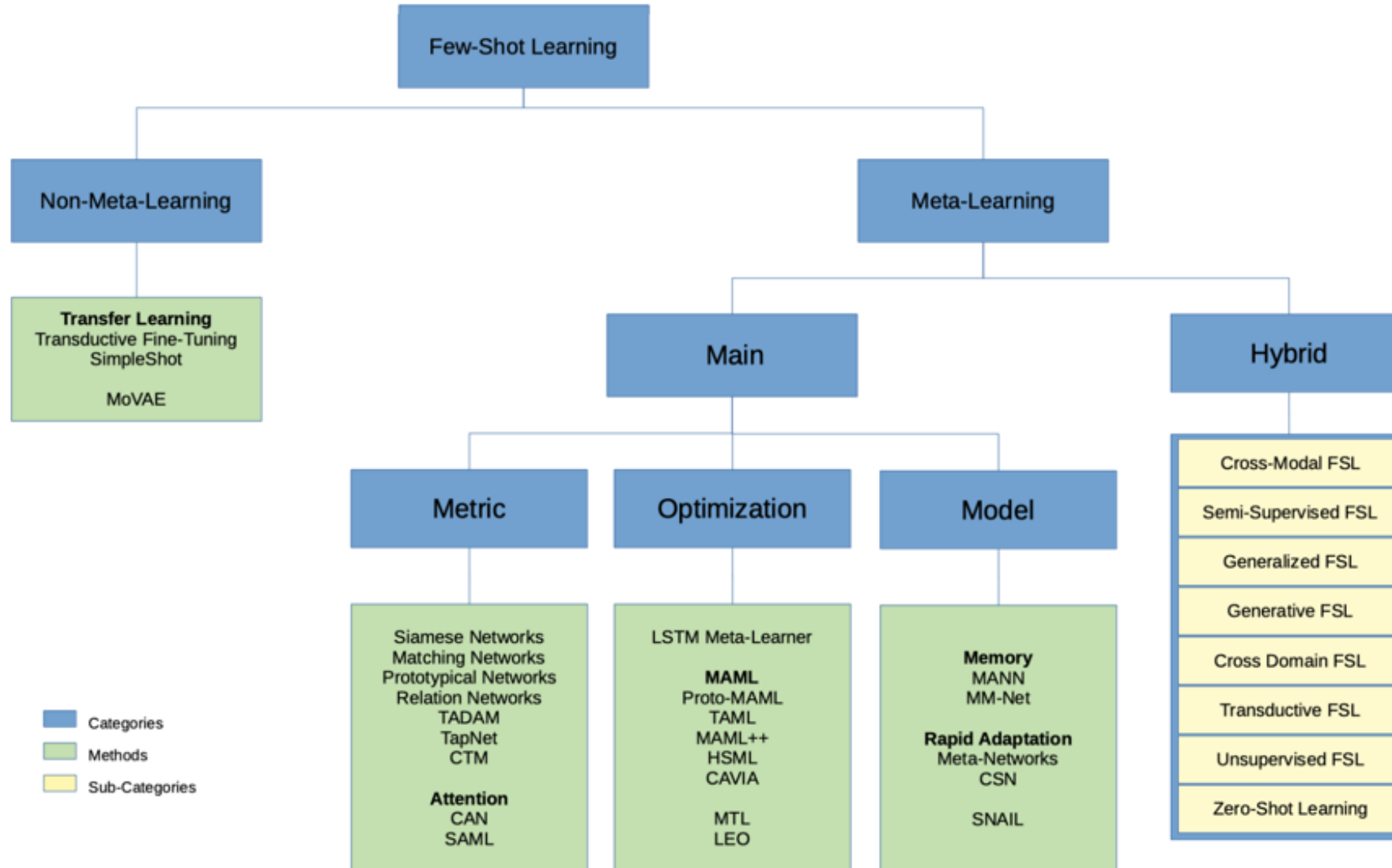


# Machine Learning, Deep Learning, Meta Learning



# Few-Shot Learning (FSL) and Meta Learning

## Machine learning from few training examples



# Meta Learning, Transfer Learning, Ensemble Learning, Continual Learning, Multi-Task Learning

Features	Method					
	Meta-learning	Transfer learning	Ensemble learning	Continual learning	Multi-task learning	Hierarchical Bayesian models
Learning from prior experience	✓	✓	✗	✓	✗	✓
Relationship between source tasks	No limitation	Related	Same	Task streams	Related	Related
Relationship between source tasks and target tasks	No limitation	Related	Same	Related	Related	Related
Considering the requirements of the target task	✓	✗	✗	✗	✗	✗

# Meta-Learning and Few-shot Learning

## Notations and Terms

### Optimization-based Meta-learning

Notation A	Term A
$\mathcal{D}_i^{train}$	Training set for task $\mathcal{T}_i$
$\mathcal{D}_i^{test}$	Test set for task $\mathcal{T}_i$
$\mathcal{D}_{meta-train}$	Meta-training set
$\mathcal{D}_{meta-test}$	Meta-testing set

### Metric-based Meta-learning

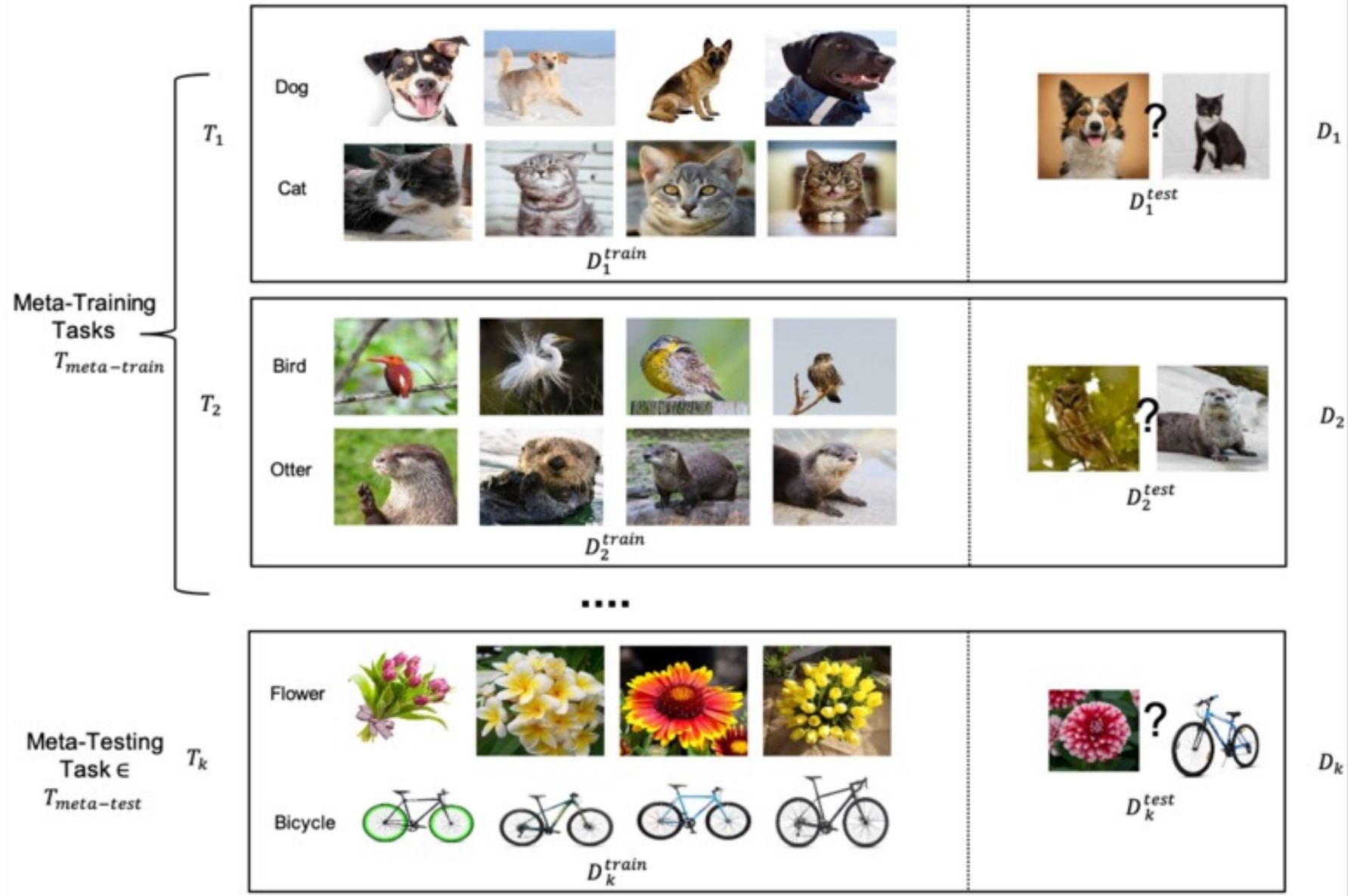
Notation B	Term B
$S_i$	Support Set for task $\mathcal{T}_i$
$Q_i$	Query Set for task $\mathcal{T}_i$
$\mathcal{D}_{train}$	Training Set
$\mathcal{D}_{test}$	Test Set

# Meta-Learning Symbols

Symbol	Meaning
$\mathcal{T}_i$	Task $i$
$\mathcal{L}$	Loss function
$(x_k, y_k)$	Input-Output pair
$f_\theta$	Model (function) with parameters $\theta$
$g_{\theta_1}$	Embedding function
$d_{\theta_2}$ or $d$	Distance function
$g_\phi$	Meta-Learning model with parameters $\phi$
$P_\theta(y x)$	Output probability of $y$ for input $x$ using model parameters $\theta$
$k_\theta(x_1, x_2)$	Kernel function measuring similarity between two vectors $x_1$ and $x_2$
$\sigma$	Softmax function
$\alpha, \beta$	Learning rates
$w$	Weights
$\mathbf{v}_c$	Prototype of class $c$
$\mathcal{C}$	Set of classes present in $\mathcal{S}$
$\mathcal{S}^c$	Subset of $\mathcal{S}$ containing all elements $(x_k, y_k)$ such that $y_k = c$
$\oplus$	Concatenation operator
$B$	Number of batches $(X_b, Y_b)$ sampled in inner-loop for a randomly sampled task $\mathcal{T}_i$
$I$	Number of tasks $\mathcal{T}_i$ sampled in inner-loop
$J$	Number of outer-loop iterations



# Meta-Learning Example Setup

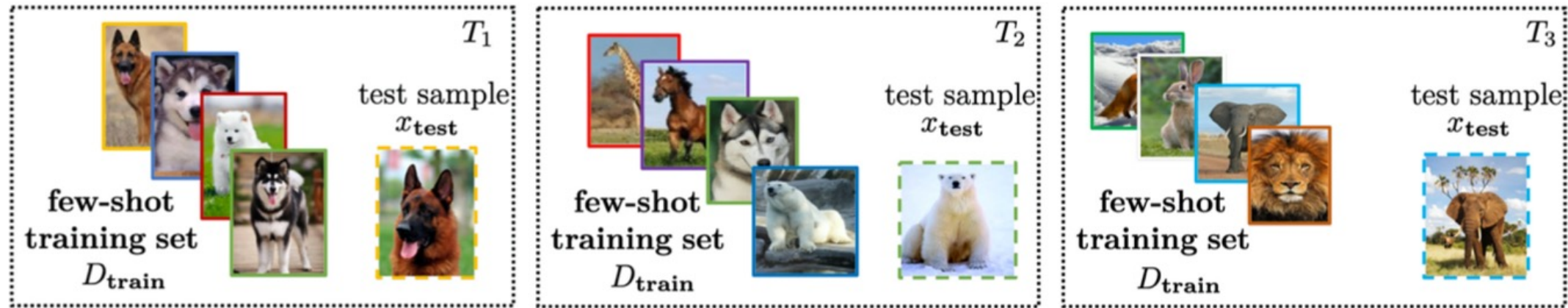




# Few-Shot Learning (FSL)

## Solving the FSL problem by meta-learning

meta-training tasks  $T_s$ 's



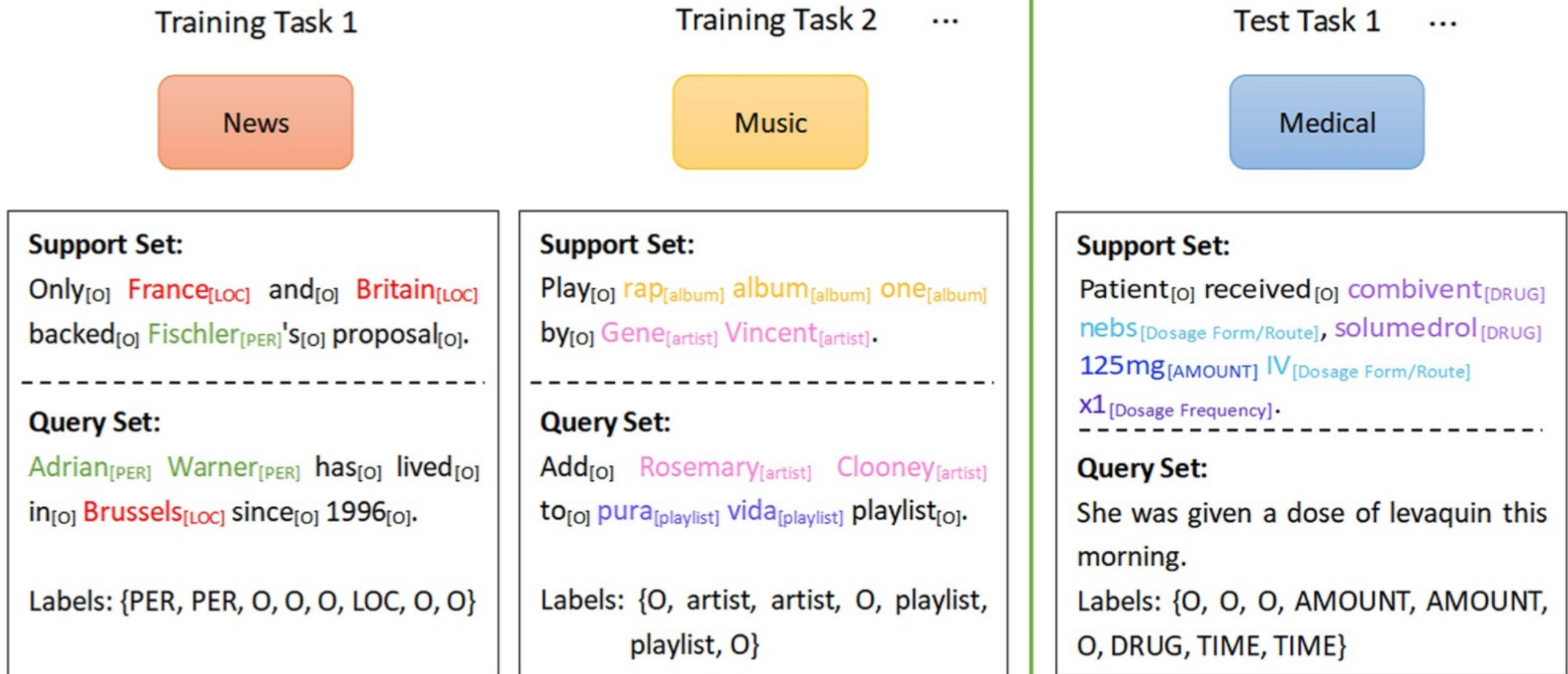
meta-testing tasks  $T_t$ 's



# Few-Shot Learning (FSL)

## Meta-learning

Each task mimics the few-shot scenario, and can be completely non-overlapping.  
Support sets are used to train; query sets are used to evaluate the model



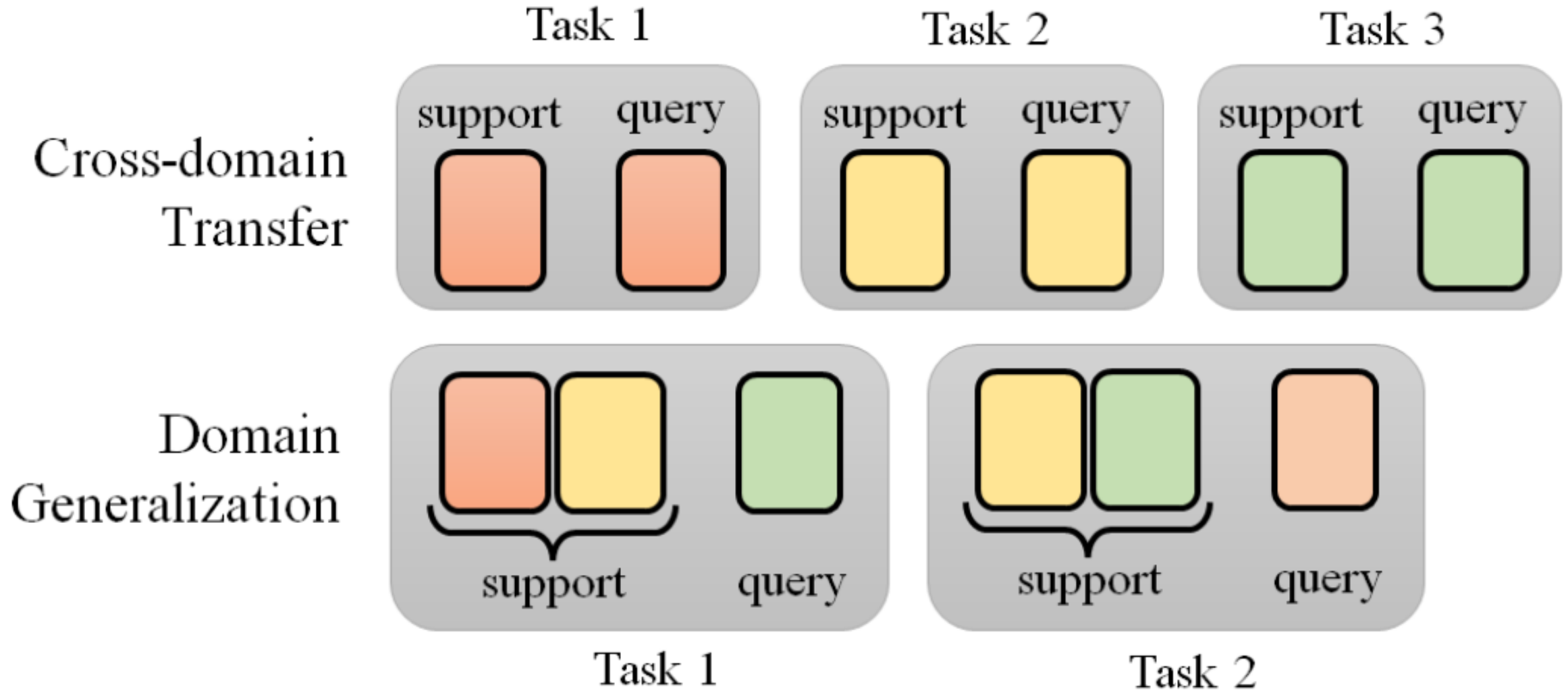
# Meta-Task Learning (MTL)

## Transfer Learning Strategy for Meta-Learning

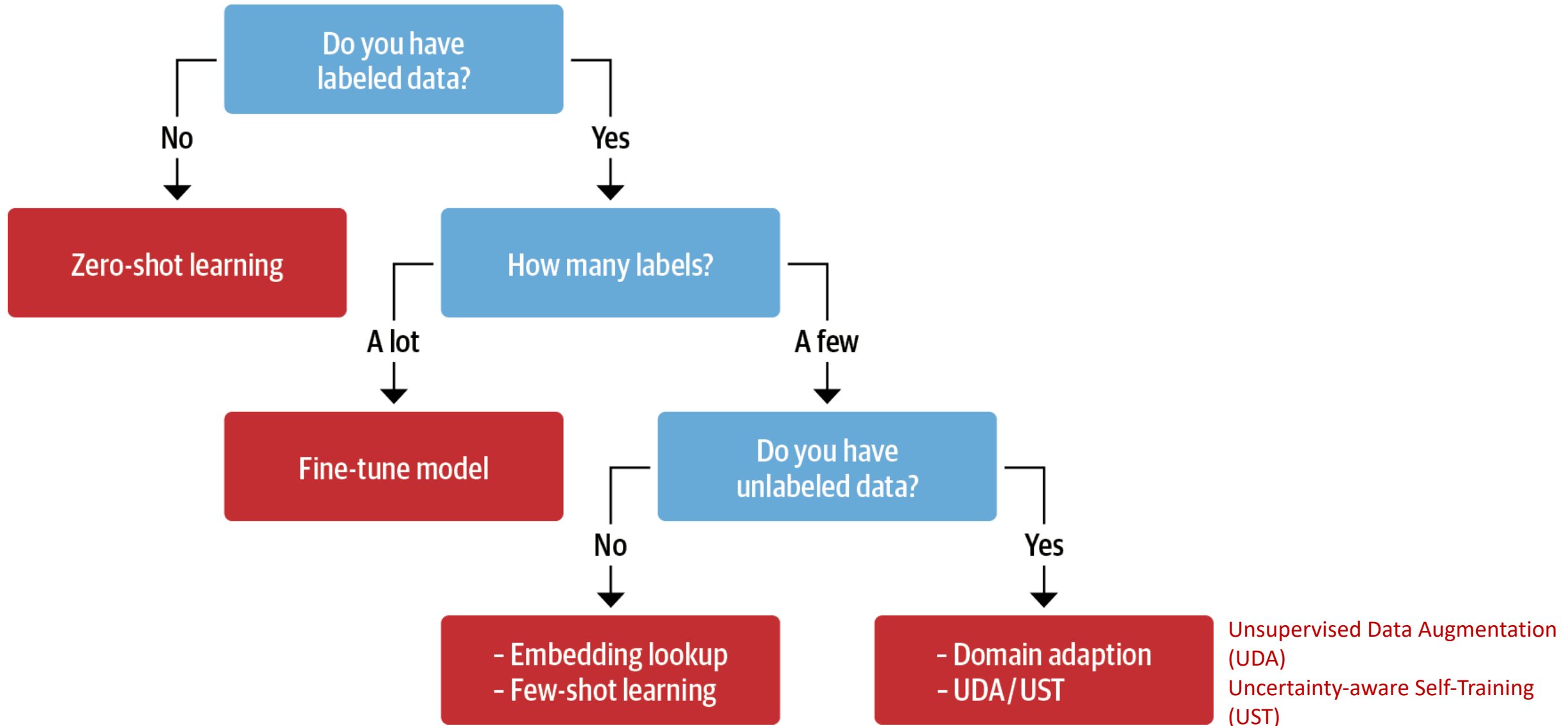
	large-scale training	meta-training	meta-test
Transfer Learning	task <sub>1</sub> model <sub>1</sub>		task <sub>2</sub> model <sub>1</sub> + $FT$
Meta-Learning	task <sub>1</sub> model <sub>1</sub> ...    task <sub>N</sub> model <sub>N</sub>		task <sub>N+1</sub> model <sub>N+1</sub>
Meta-Transfer Learning	task model	task <sub>1</sub> model + $SS_1 + FT_1$ ⋮ task <sub>N</sub> model + $SS_N + FT_N$	task <sub>N+1</sub> model + $SS_N + FT_{N+1}$

# Meta Learning

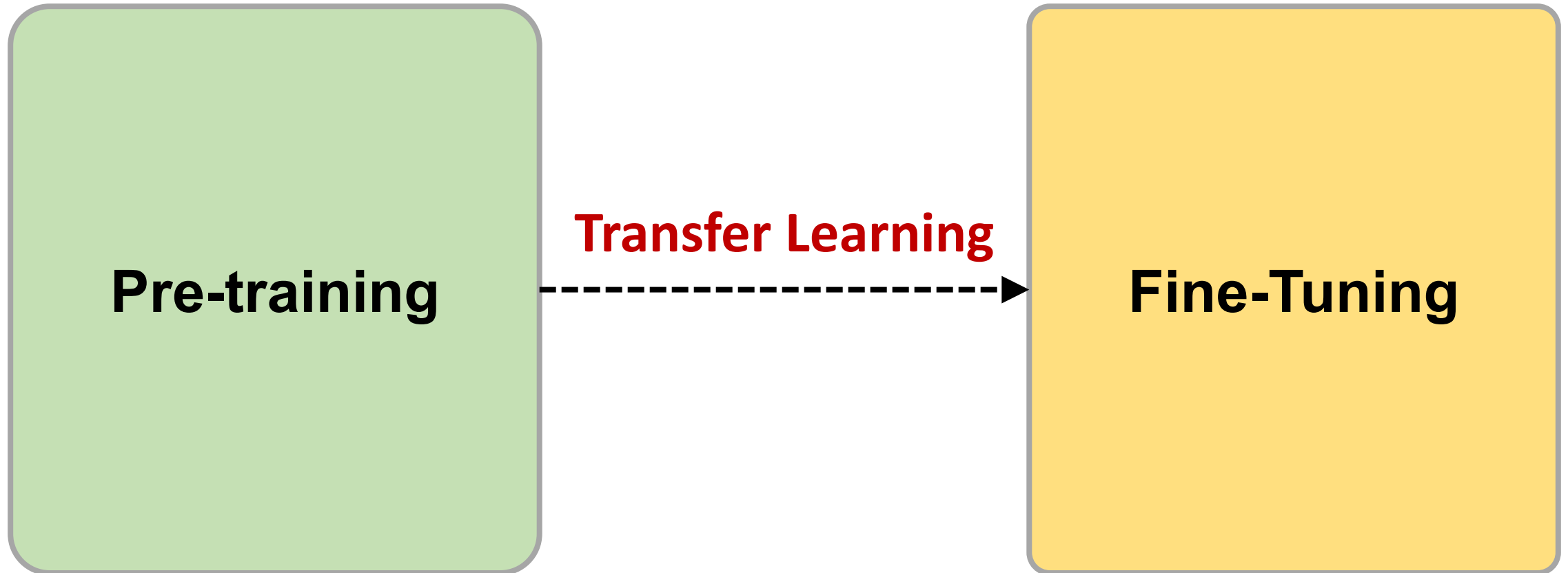
## The task construction of cross-domain transfer and domain generalization



# Transfer Learning, Fine-tuning, Few-shot learning



# Transfer Learning

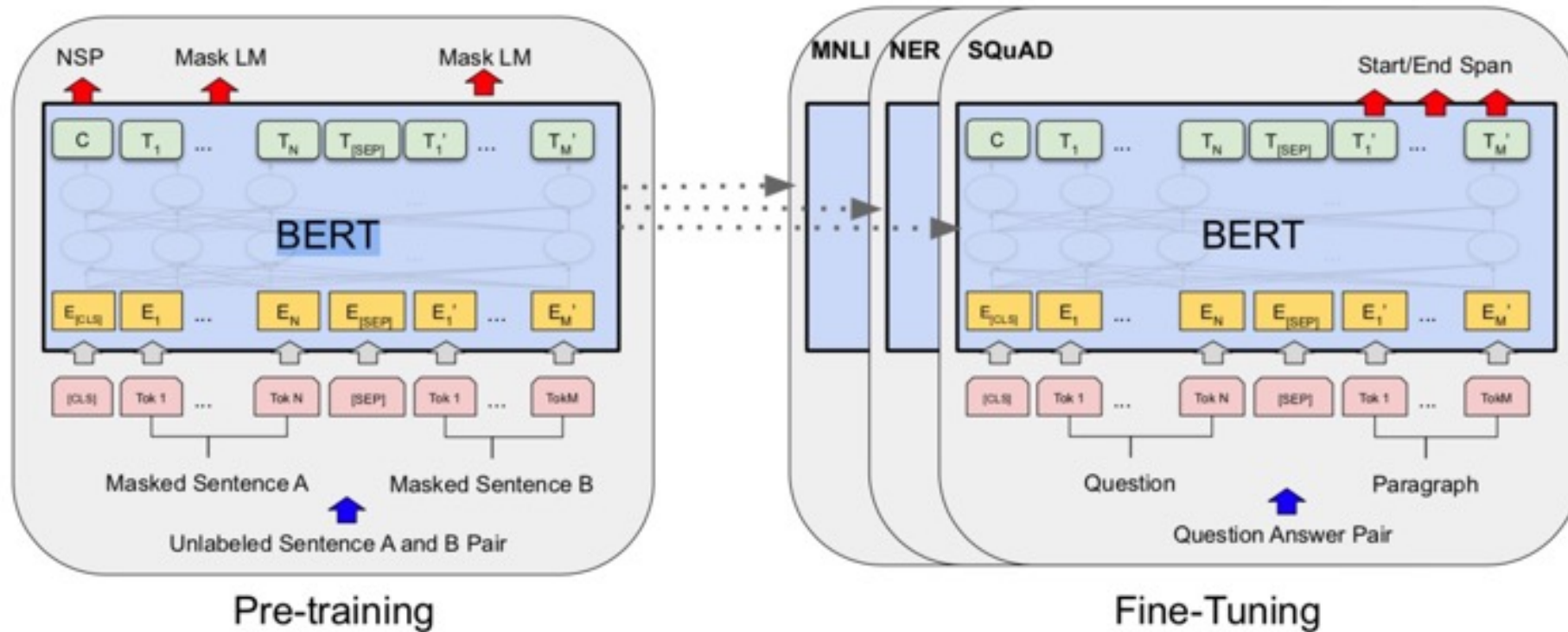




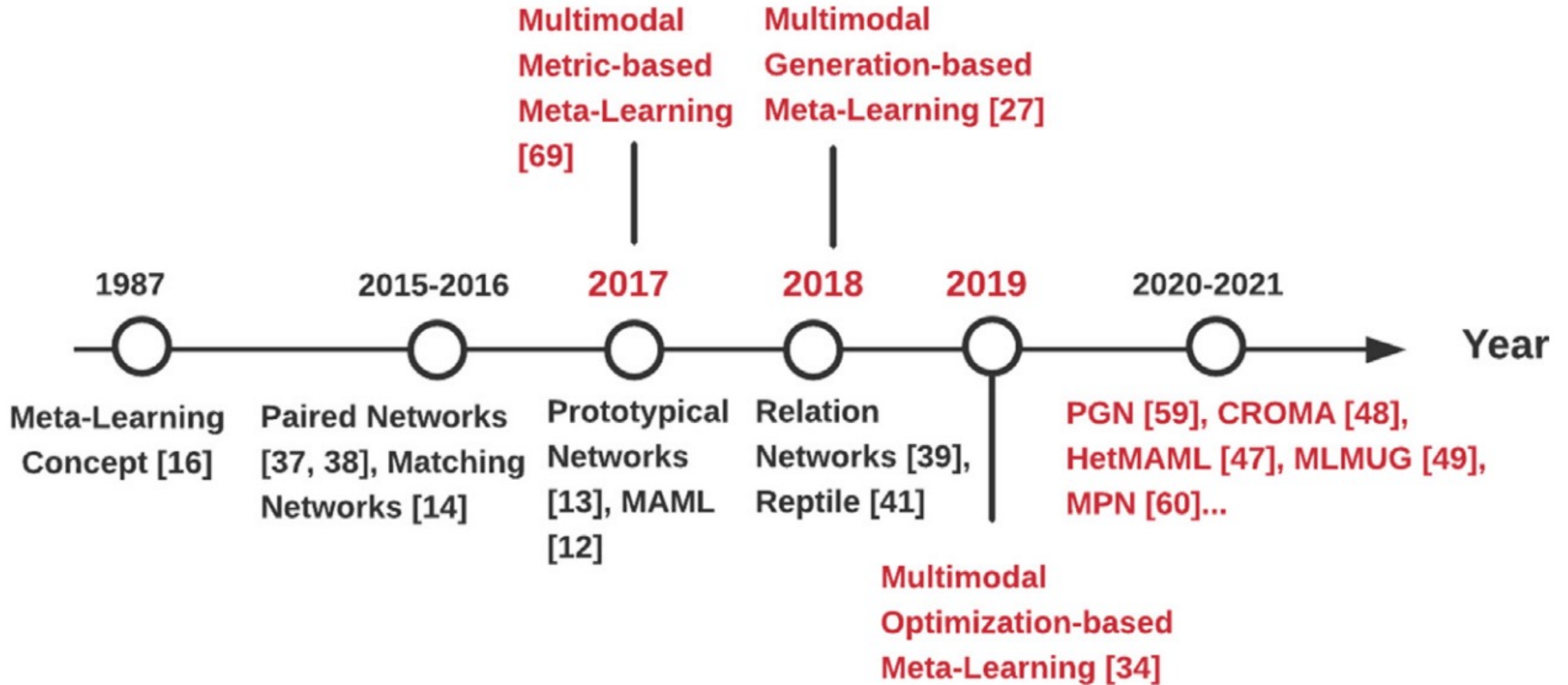
# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

BERT (Bidirectional Encoder Representations from Transformers)

Overall pre-training and fine-tuning procedures for BERT



# Meta Learning





# Meta Learning

Year	Achievement	Ref.
1987	(1) A new framework of “learning how to learn” with self-referential learning was proposed. The neural networks in self-referential learning can regard their weights as inputs and update them continuously. (2) Based on the conventional neural network, two types of wights were used to connect the neurons. Each type of weight presents a different learning speed.	[34,35]
1990	A synaptic learning rule, which is biologically plausible, was proposed to automatically study the learning rules.	[36]
1993	A chain of meta-networks was introduced to improve the learning capacity of a recurrent neural network for a dynamic environment.	[37]
1995	A framework was proposed to optimize the learning rule within a parametric learning rule space.	[38]
1996	An improved self-referential model was proposed. Time ratios were used to measure the effects of learning processes on the later learning processes.	[39]
1998	The term “Learning to learn” was proposed to equally represent the concept of meta-learning.	[40]
2001	Gradient descent methods were firstly used in meta-learning instead of evolutionary methods, which were widely used in previous research.	[41,42]
2003	A biologically plausible meta-reinforcement learning algorithm was proposed to tune the parameters of the meta-learning model dynamically and adaptively.	[43]
2004	A new perspective of meta-learning was proposed: exploring the interaction between the learning mechanism and the specific contexts to which the mechanism applies.	[9]
2008	The zero-data learning problem was addressed.	[44]
2010–2012	The breakthrough of deep neural networks marks the beginning of the era of meta-learning.	[45–47]
2013	The relationship between transfer learning and meta-learning was described.	[48]
2016	A meta-learning algorithm named gradient descent by gradient descent was proposed.	[49]
2017	(1) MAML was proposed. (2) A doctoral thesis systematically introduced the concept of meta-learning and corresponding methods.	[50,51]
2018	Reptile, an improved version of MAML, was proposed.	[52]
2019	The Capsule network provides a new method to improve the learning capacity of meta-learning, especially in computer vision.	[53]
2020	Combining auto-encoder and capsule network to focus on the zero-shot learning problem.	[54]

# Meta-learning Approaches

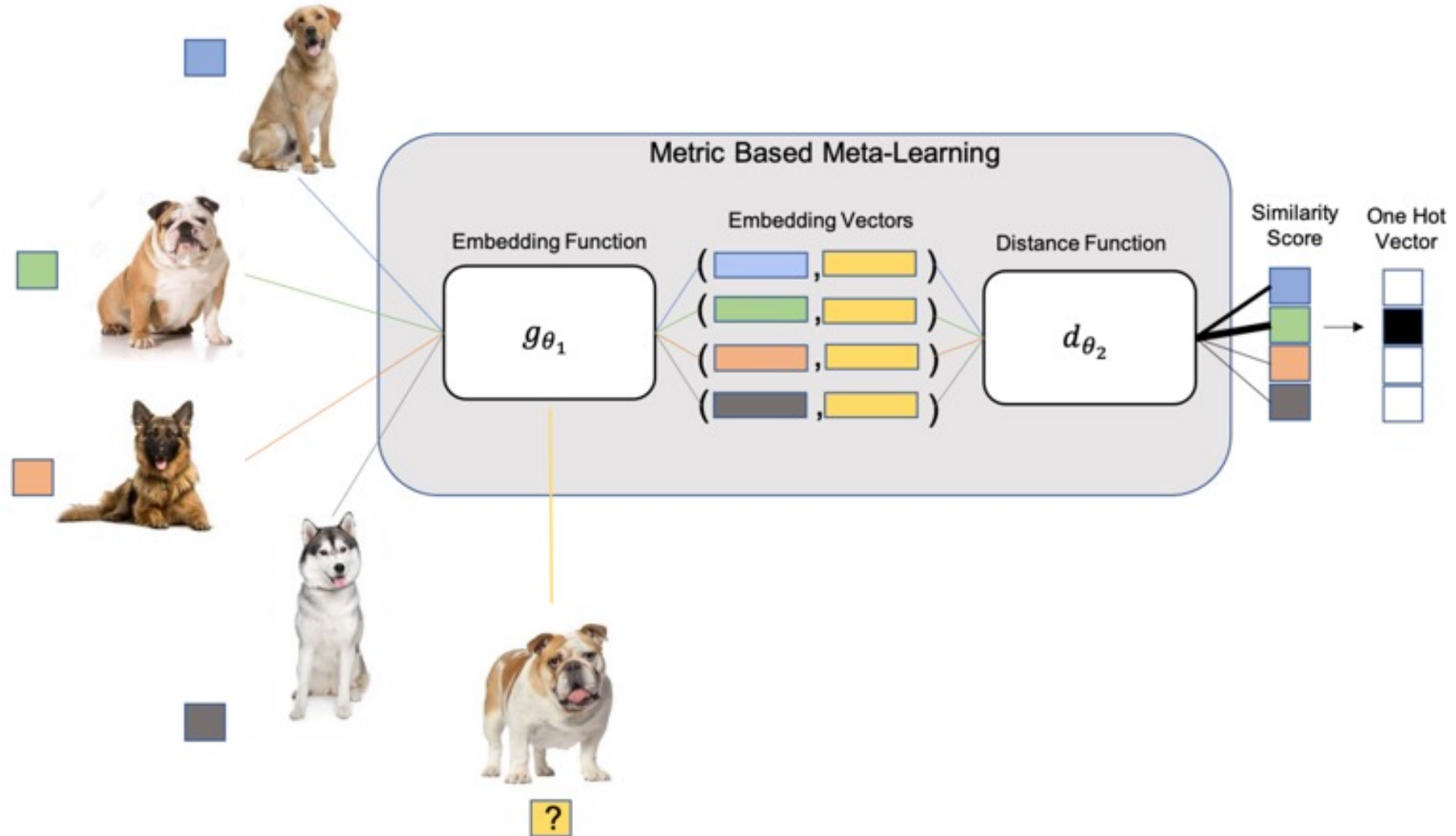
	Metric-based	Optimization-based	Model-based
<b>Key idea</b>	Metric Learning [19]	Gradient Descent	Memory; RNN
<b>How <math>P_\theta(y x)</math> is modeled?</b>	$\sum_{(x_k, y_k) \in S} k_\theta(x, x_k) y_k,$	$P_{\theta'}(y x),$ <p>where <math>\theta' = g_\phi(\theta, S)</math></p>	$f_\theta(x, S).$
<b>Advantages</b>	<p>Faster Inference.</p> <p>Easy to deploy.</p>	<p>Offers flexibility to optimize in dynamic environments.</p> <p><math>S</math> can be discarded post-optimization.</p>	<p>Faster inference with memory models.</p> <p>Eliminates the need for defining a metric or optimizing at test.</p>
<b>Disadvantages</b>	<p>Less adaptive to optimization in dynamic environments.</p> <p>Computational complexity grows linearly with size of <math>S</math> at test.</p>	<p>Optimization at inference is undesirable for real-world deployment.</p> <p>Prone to overfitting.</p>	<p>Less efficient to hold data in memory as <math>S</math> grows.</p> <p>Hard to design.</p>

# Meta Learning: Learning to Learn

Class	Methods	Reference	Summary
Metric-Based	Siamese Neural Networks	[32–36]	We show four metric-based meta-learning algorithms, focusing on feature extractors, similarity metrics, and automatic algorithm selection. However, the metric-based approaches are sensitive to the dataset and increase the computational expenditure when the number of tasks is large.
	Matching Networks	[37–41]	
	Prototype Networks	[42–46]	
	Relation Networks	[47–53]	
Model-Based	Memory-Augmented Neural Networks	[54–56] [57,58]	We display three model-based approaches. MANN combines neural networks with external memory modules, but the model is complex. Meta-Net is computationally intensive and has high memory requirements. SNAIL is relatively simplified, but has to be optimized in terms of automatic parameter tuning and reducing computation.
	Meta Networks	[59–65]	
	Simple Neural Attentive Meta-Learner	[66–71]	
	MAML	[72–80]	
Optimization-Based	META- LSTM	[81–86]	We present three methods of optimization-based meta-learning. MAML is relatively simple to implement, but the capacity of the model is limited. Meta-LSTM has a large capacity, but a complicated training process. Meta-SGD has improved capacity but still has difficulties in generalization ability.
	META- SGD	[87–93]	

# Metric-based Meta-learning

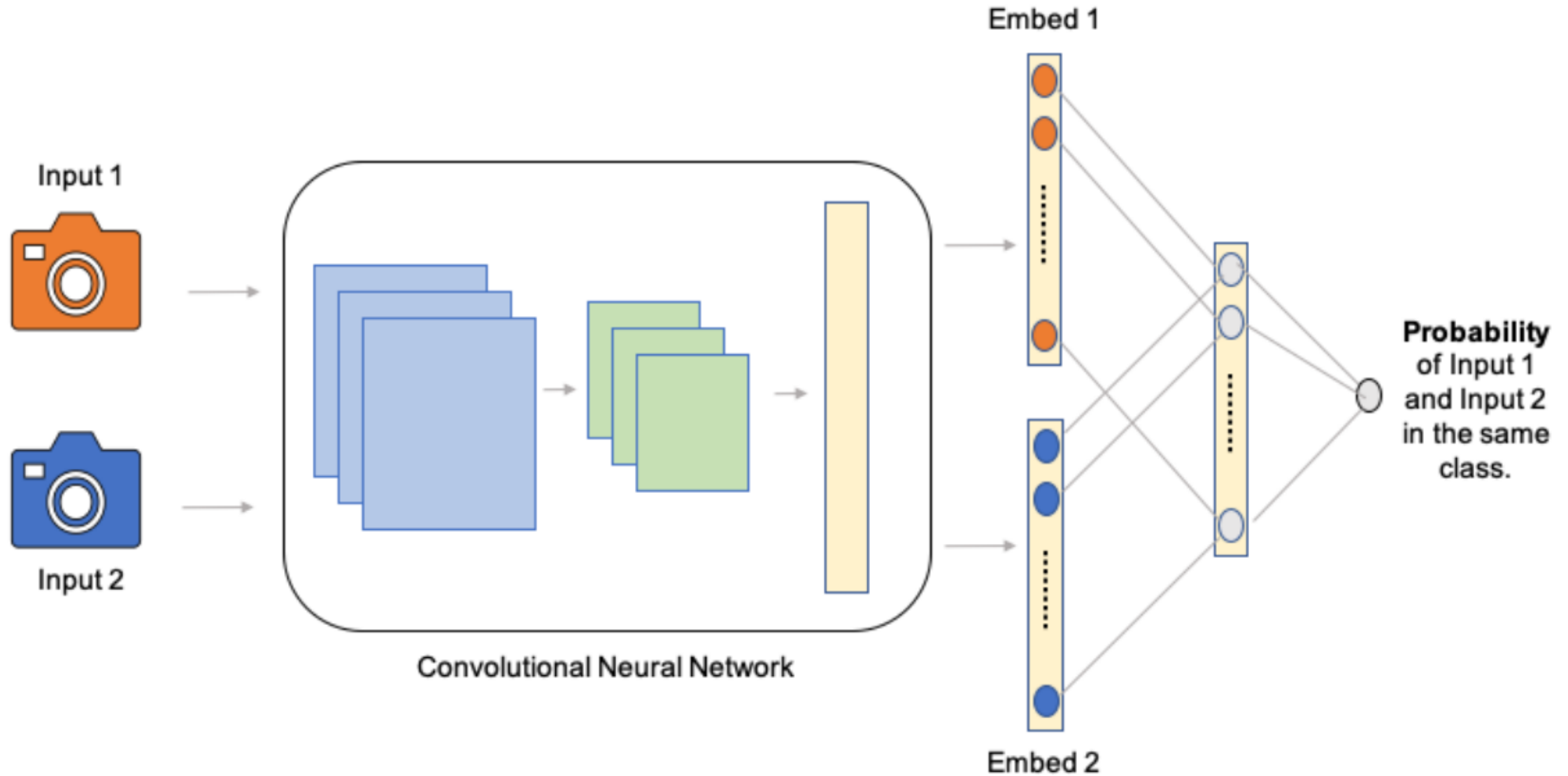
## M-Way K-Shot Task (4-way-1-shot classification task)



# Metric-based Meta-Learning Methods

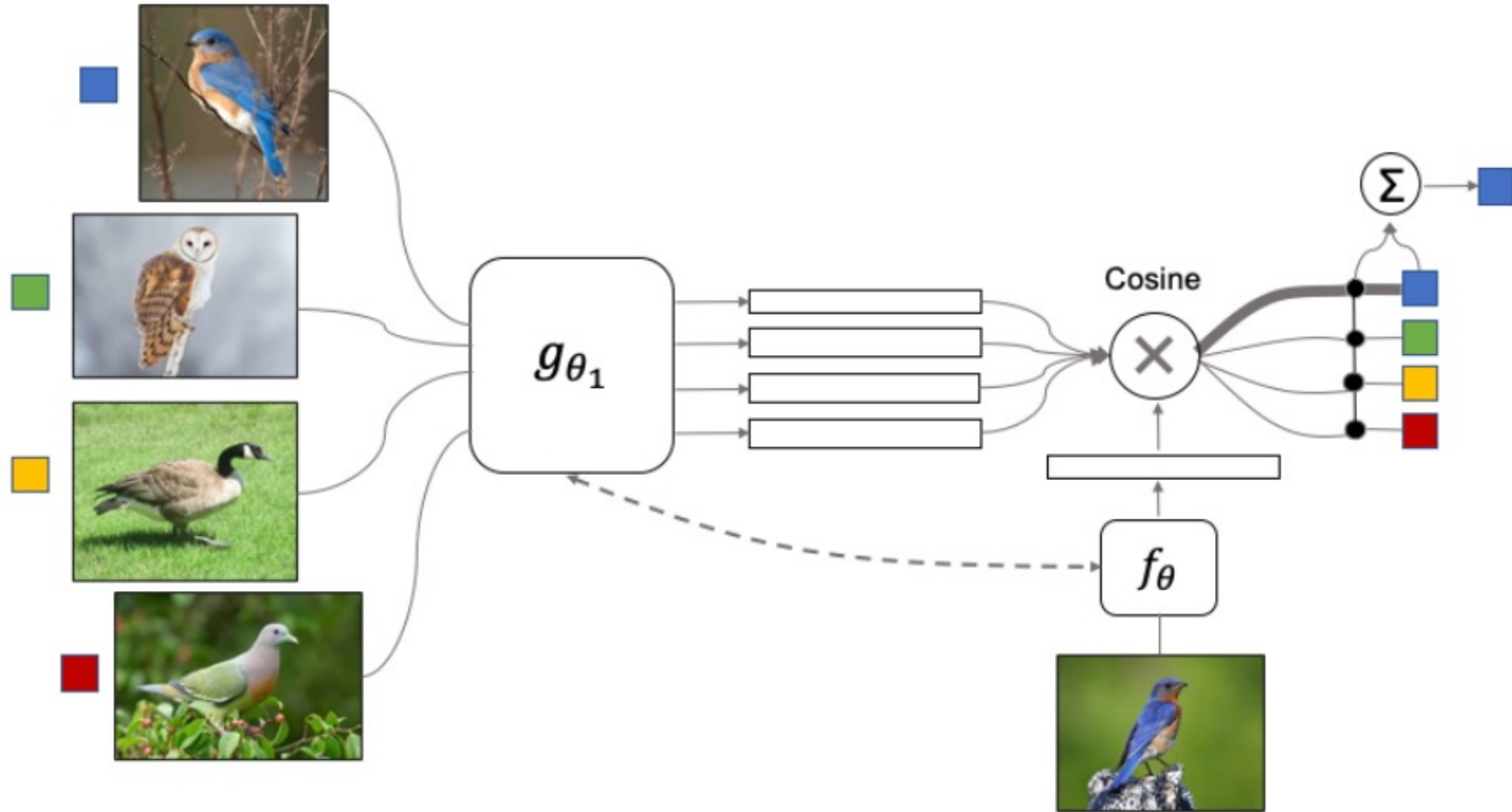
Method	T.I	$g_{\theta_1}$	$d_{\theta_2}$	Prediction	Loss
Siamese Networks [20]	Yes	CNN	L1	$v = w \cdot d(g_{\theta_1}(x_1), g_{\theta_2}(x_2))$ $p = \text{sigmoid}(\sum_j v_j)$	$-(y \log(p) + (1 - y)(\log(1 - p)))$
Matching Networks [13]	Yes	CNN + LSTM w/ attention	Cosine Similarity	$\hat{y} = \frac{\sum_{k=1}^t \sigma(d(f_{\theta}(\hat{x}), g_{\theta_1}(x_k))) y_k}{\sum_{k=1}^t \sigma(d(f_{\theta}(\hat{x}), g_{\theta_1}(x_k)))}$ $P(y = c   \hat{x}) = \hat{y}_c$	$-\log P$
Prototypical Networks [21]	Yes	CNN	Euclidean	$P(y = c   x) = \frac{\exp(-d(g_{\theta_1}(x), \mathbf{v}_c))}{\sum_{c' \in \mathcal{C}} \exp(-d(g_{\theta_1}(x), \mathbf{v}_{c'}))}$	$-\log P$
Relation Networks [22]	Yes	CNN	Learned by CNN	$r_c = d_{\theta_2}(g_{\theta_1}(x) \oplus \mathbf{v}_c)$	$\sum_{c \in \mathcal{C}} (r_c - \mathbf{1}(y == c))^2$
TADAM [16]	No	ResNet-12	Cosine / Euclidean	$P_{\lambda}(y = c   x) = \frac{\exp(-\lambda d(g_{\theta_1}(x, \Gamma), \mathbf{v}_c))}{\sum_{c' \in \mathcal{C}} \exp(-\lambda d(g_{\theta_1}(x, \Gamma), \mathbf{v}_{c'}))}$	$-\log P$
TapNet [23]	No	Resnet-12	Euclidean	$P(y = c   x) = \frac{\exp(-d(\mathbf{M}(g_{\theta_1}(x)), \mathbf{M}(\Phi_c)))}{\sum_{c' \in \mathcal{C}} \exp(-d(\mathbf{M}(g_{\theta_1}(x)), \mathbf{M}(\Phi_{c'})))}$	$-\log P$
CTM [24]	No	Any	Any	-	-

# Convolutional Siamese Network



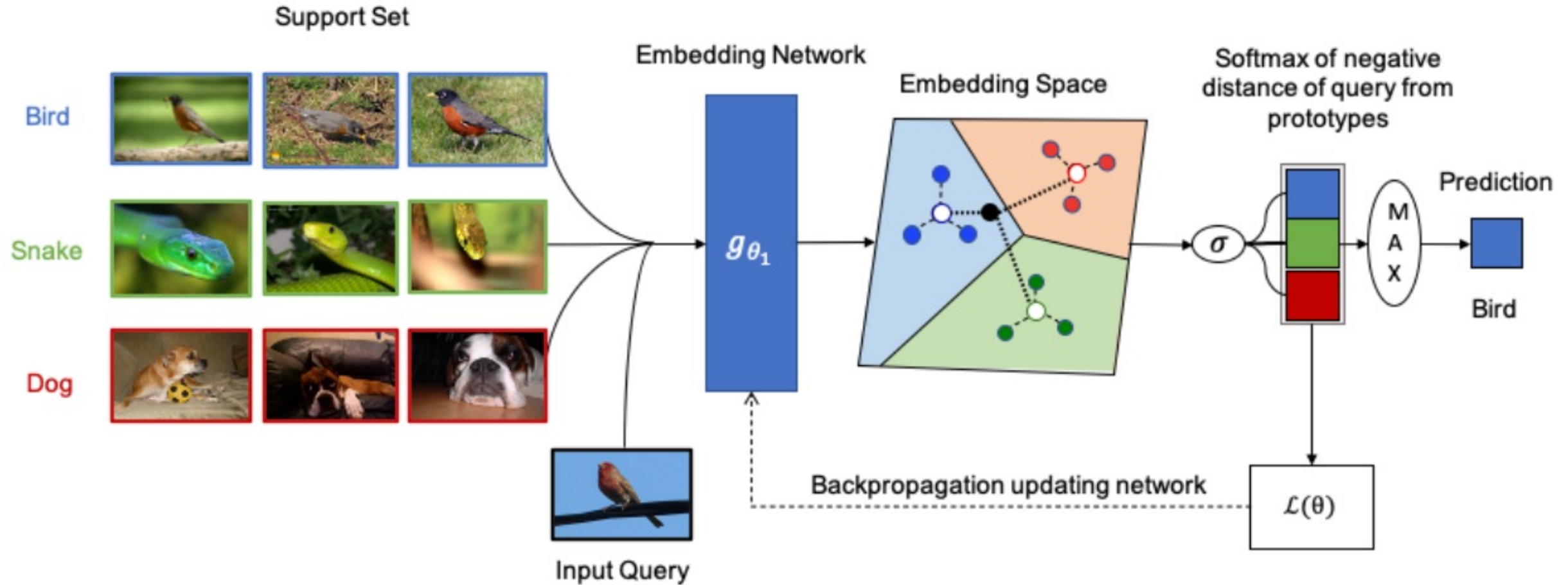


# Meta Learning: Matching Networks



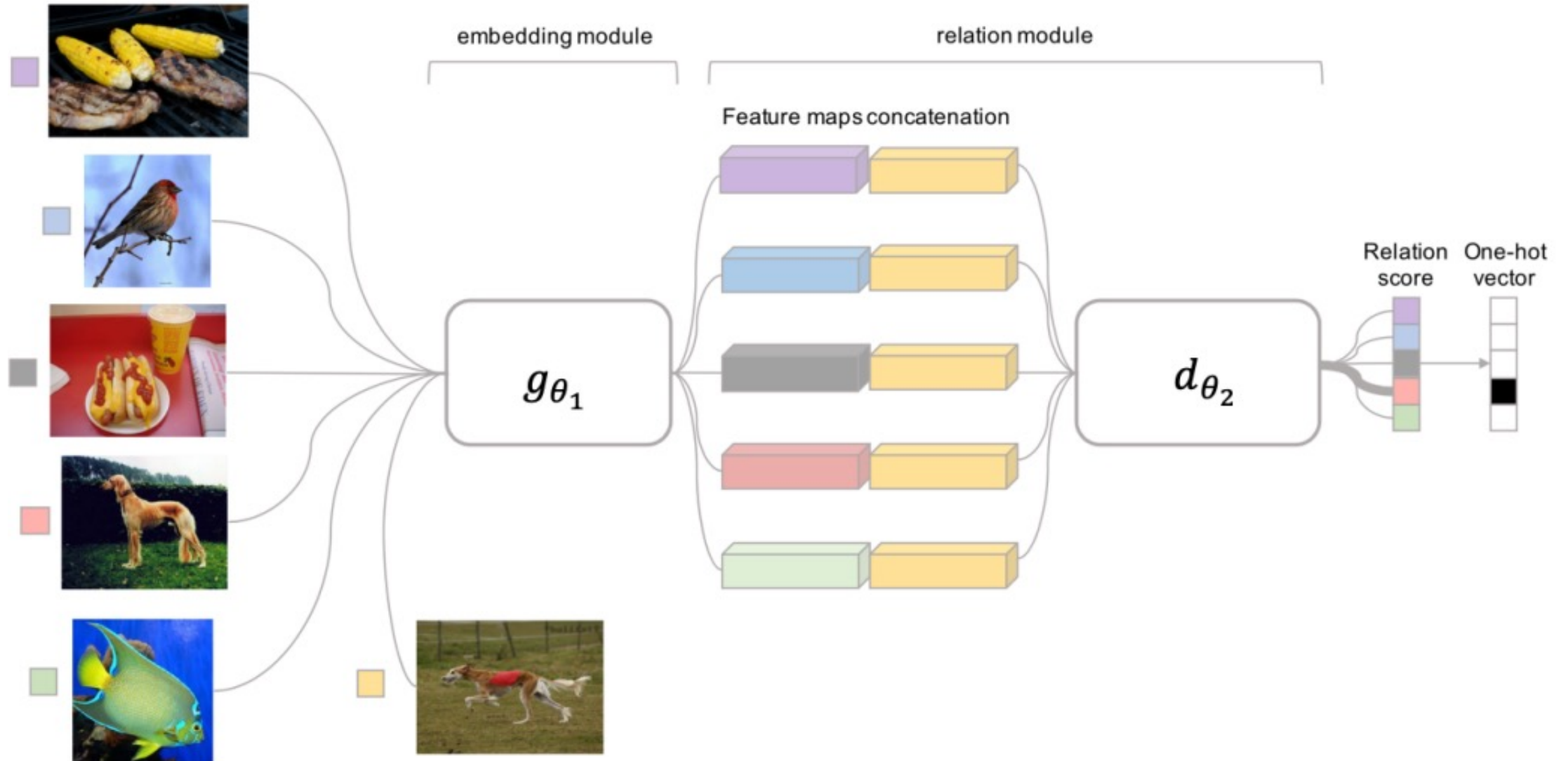
# Few-shot Prototypes

$v_c$  are computed as the mean of embedded support examples for each class

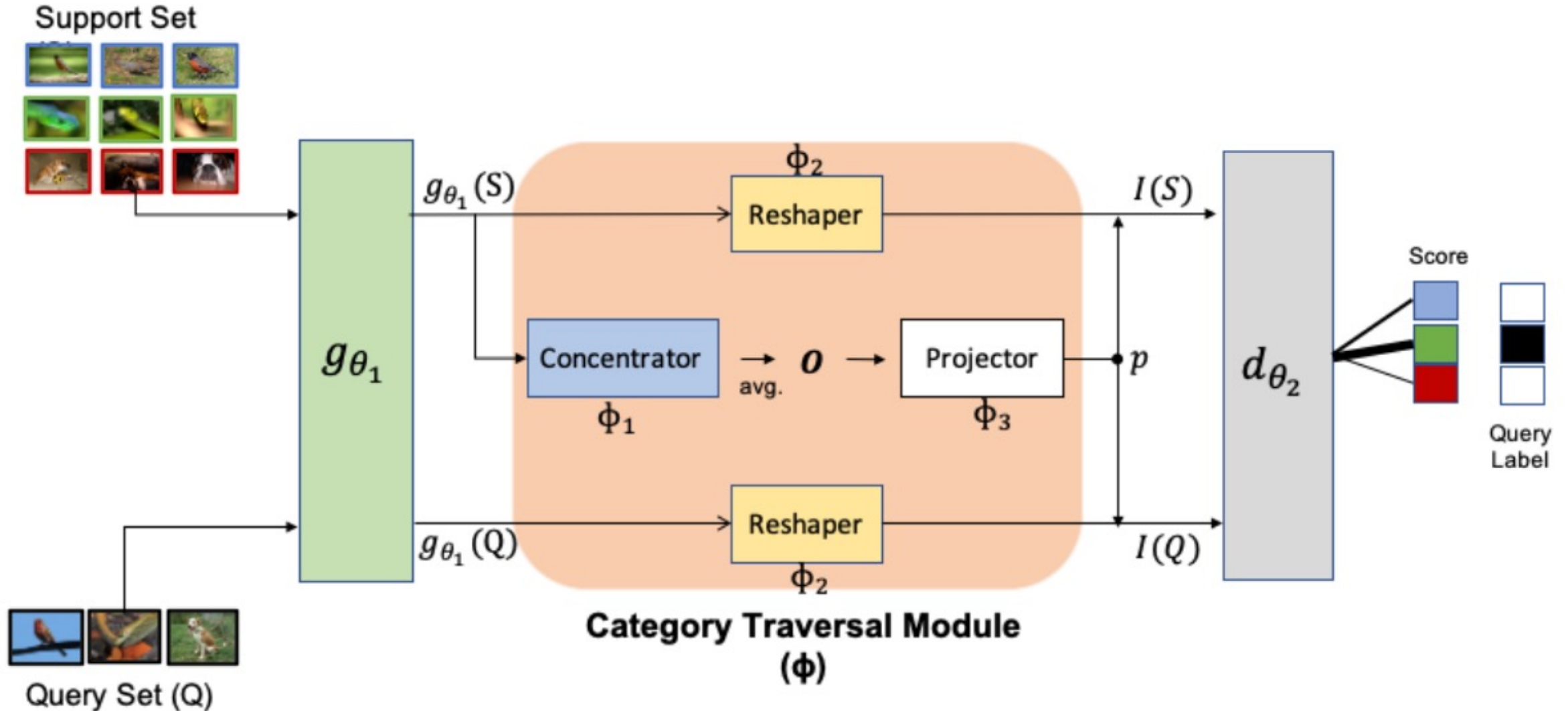




# Meta Learning: Relation Network



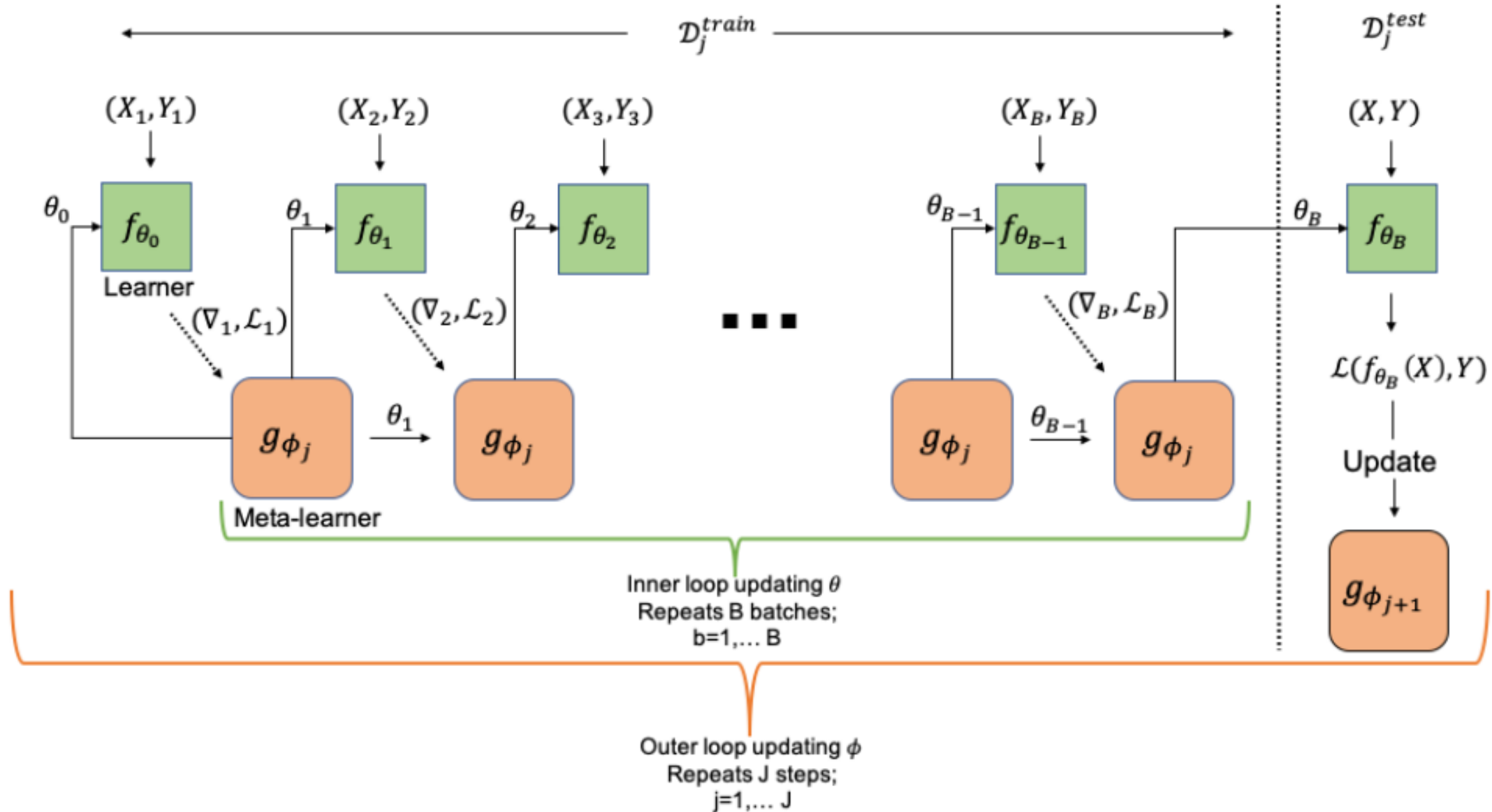
# Meta Learning: Category Traversal Module (CTM)



# Optimization-based Meta-Learning Methods

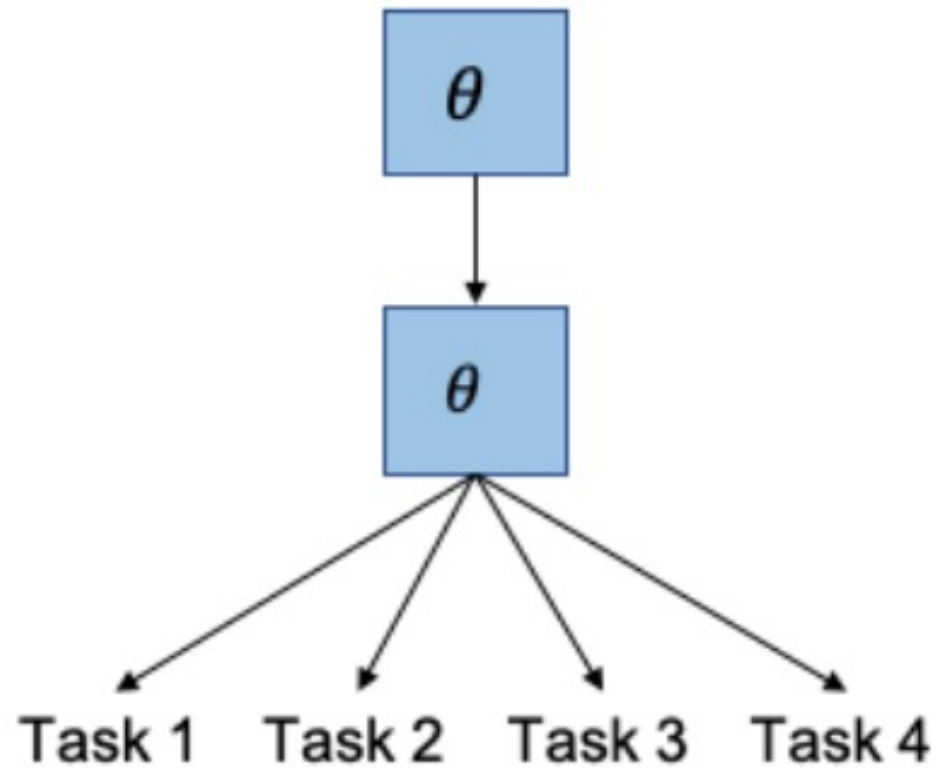
Method	Learner	Meta-Learner
LSTM Meta-Learner [35]	Repeat $\forall b \in [1..B]$ $\mathcal{L}_b \leftarrow \mathcal{L}(f(X_b; \theta_{b-1}), Y_b)$ $\theta_b \leftarrow g((\nabla_{\theta_{b-1}} \mathcal{L}_b, \mathcal{L}_b); \phi_{j-1})$	Repeat $\forall j \in [1..J]$ $\mathcal{L}_j^{test} \leftarrow \mathcal{L}(f(X; \theta_B), Y)$ $\phi_j \leftarrow \phi_{j-1} - \alpha \nabla_{\phi_{j-1}} \mathcal{L}_j^{test}$
MAML [14]	Repeat $\forall i \in [1..I]$ $\mathcal{L}_i^{train} \leftarrow \mathcal{L}(f(\mathcal{D}_i^{train}; \theta_{j-1}))$ $\theta_i^* \leftarrow \theta_{j-1} - \alpha \nabla_{\theta_{j-1}} \mathcal{L}_i^{train}$ $\mathcal{L}_i^{test} \leftarrow \mathcal{L}(f(\mathcal{D}_i^{test}; \theta_i^*))$	Repeat $\forall j \in [1..J]$ $\theta_j \leftarrow \theta_{j-1} - \beta \nabla_{\theta_{j-1}} \sum_{i=1}^I \mathcal{L}_i^{test}$
MTL [37]	$\mathcal{L}_i^{train} \leftarrow \mathcal{L}(f(\mathcal{D}_i^{train}; [\theta_{j-1}, \phi_{j-1}, \Theta]))$ $\theta_i^* \leftarrow \theta_{j-1} - \alpha \nabla_{\theta_{j-1}} \mathcal{L}_i^{train}$ $\mathcal{L}_i^{test} \leftarrow \mathcal{L}(f(\mathcal{D}_i^{test}; \theta_i^*))$	$\theta_j \leftarrow \theta_{j-1} - \beta \nabla_{\theta_{j-1}} \sum_{i=1}^I \mathcal{L}_i^{test}$ $\phi_j \leftarrow \phi_{j-1} - \beta \nabla_{\phi_{j-1}} \sum_{i=1}^I \mathcal{L}_i^{test}$
LEO [38]	$\phi_{j-1} = \{\phi_e, \phi_r, \phi_d, \alpha\}$ $\mathbf{z}_i \leftarrow g(\mathcal{D}_i^{train}; [\phi_e, \phi_r, \Theta])$ $\theta_i \leftarrow g(\mathbf{z}_i; \phi_d)$ $\mathcal{L}_i^{train} \leftarrow \mathcal{L}(f(\mathcal{D}_i^{train}; \theta_i))$ $\mathbf{z}_i^* \leftarrow \mathbf{z}_i - \alpha \nabla_{\mathbf{z}_i} \mathcal{L}_i^{train}$ $\theta_i^* \leftarrow g(\mathbf{z}_i^*; \phi_d)$ $\mathcal{L}_i^{test} \leftarrow \mathcal{L}(f(\mathcal{D}_i^{test}; \theta_i^*))$	$\phi_j \leftarrow \phi_{j-1} - \beta \nabla_{\phi_{j-1}} \sum_{i=1}^I \mathcal{L}_i^{test}$

# Computational Graph for the Forward Pass of the Meta-learner

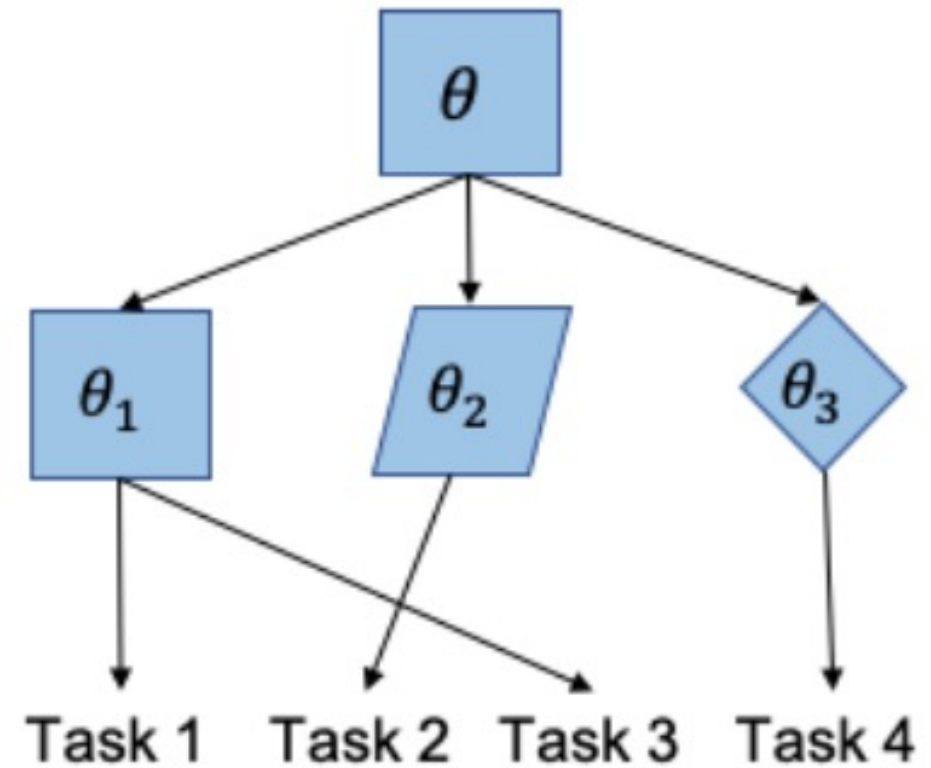


# Model-Agnostic Meta-Learning (MAML)

## Hierarchically Structured Meta-Learning (HSML)

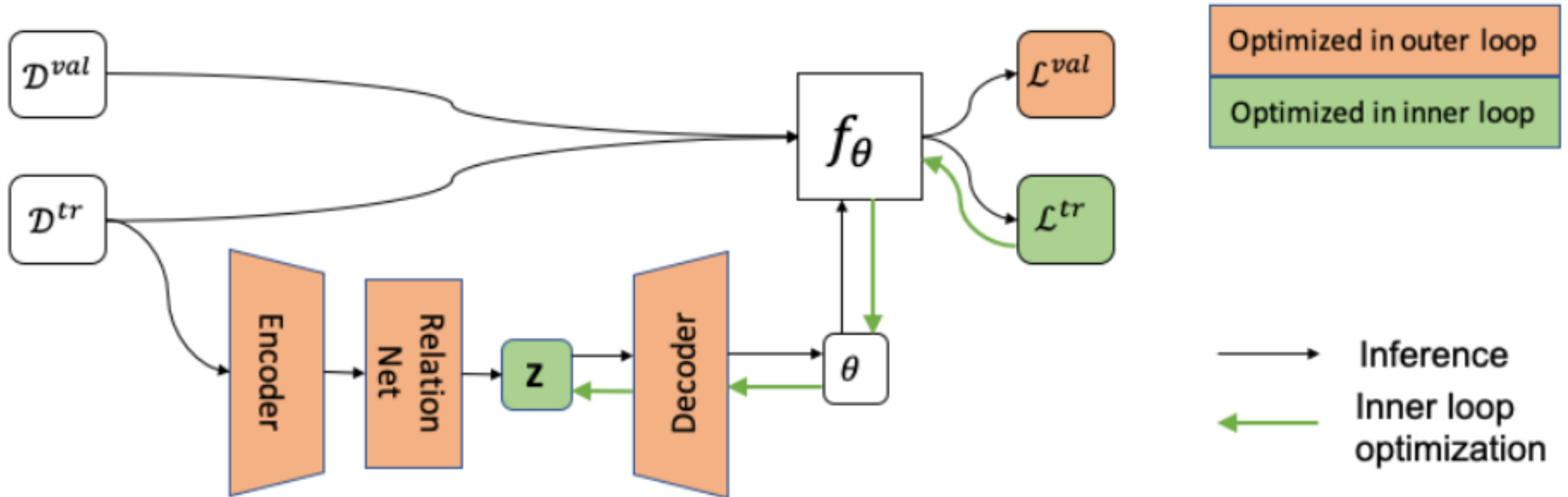


Globally Shared Initialization  
(MAML)

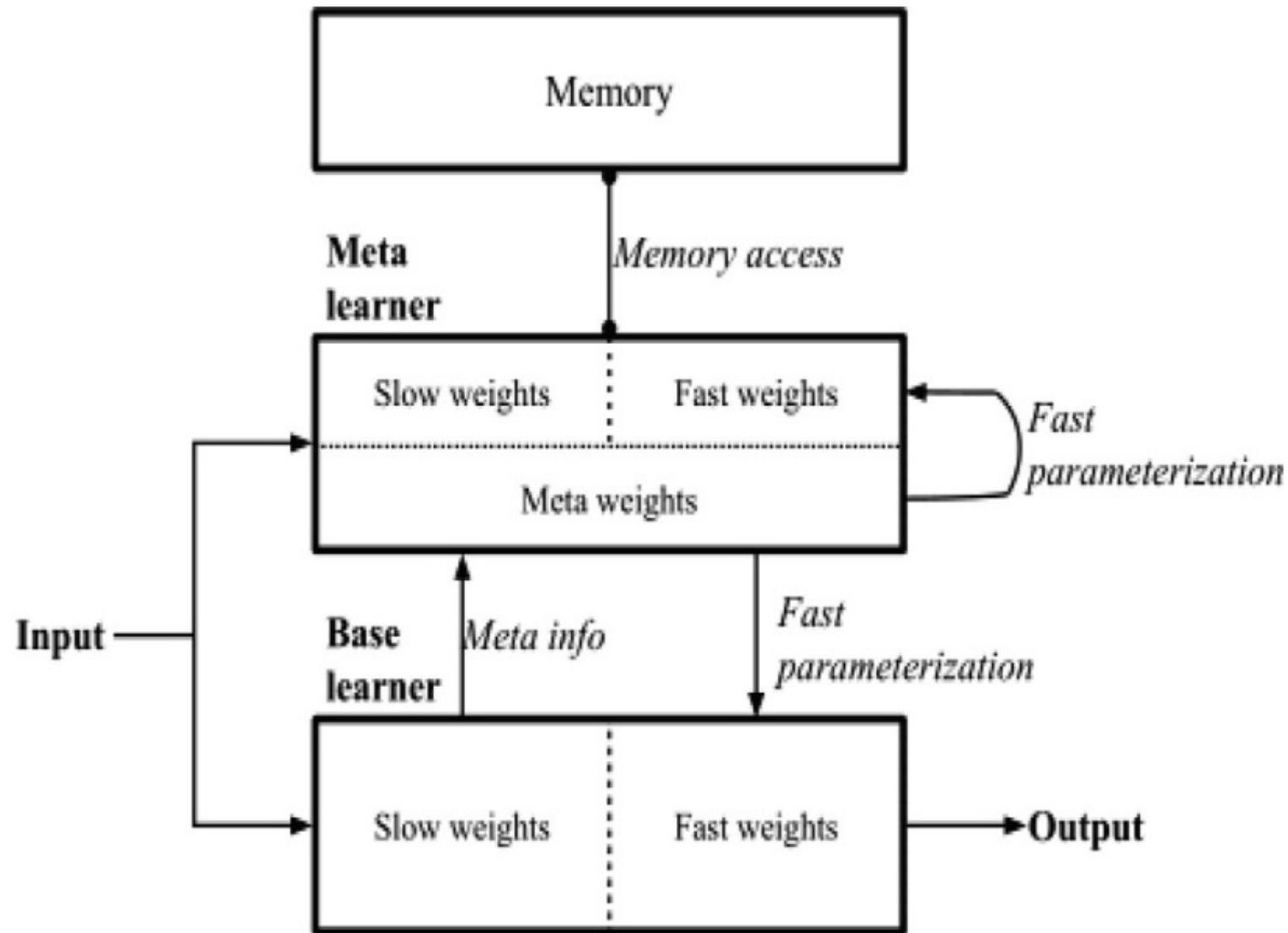


Hierarchically Clustered Initialization  
(HSML)

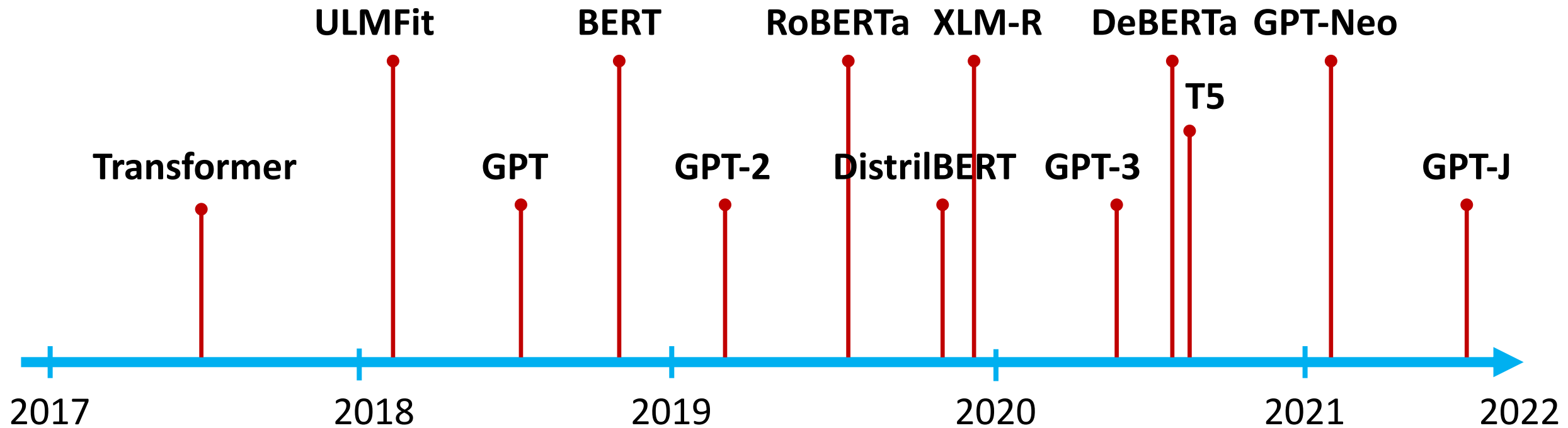
# Latent Embedding Optimization (LEO)



# Overall Architecture of Meta Networks

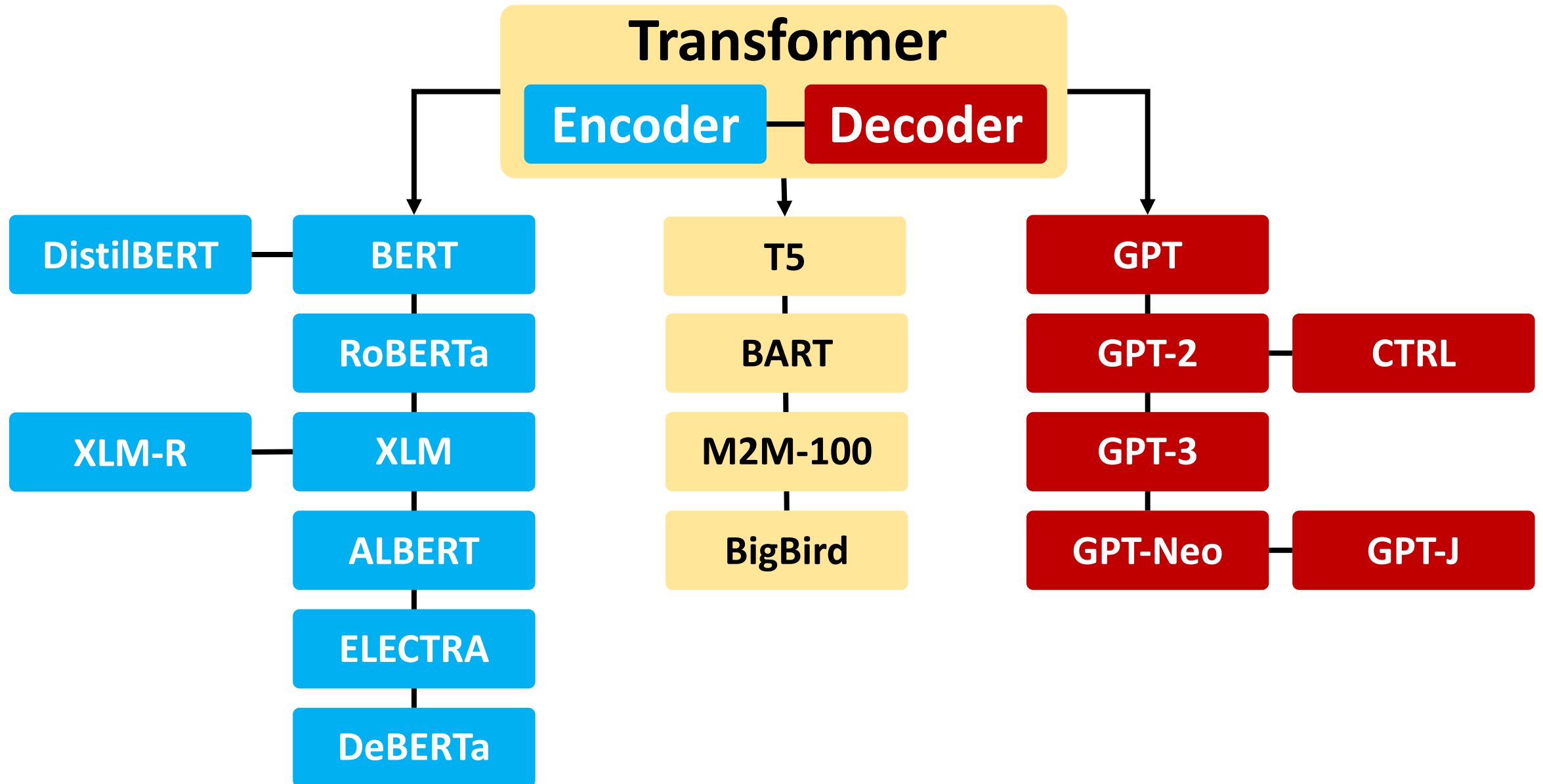


# The Transformers Timeline

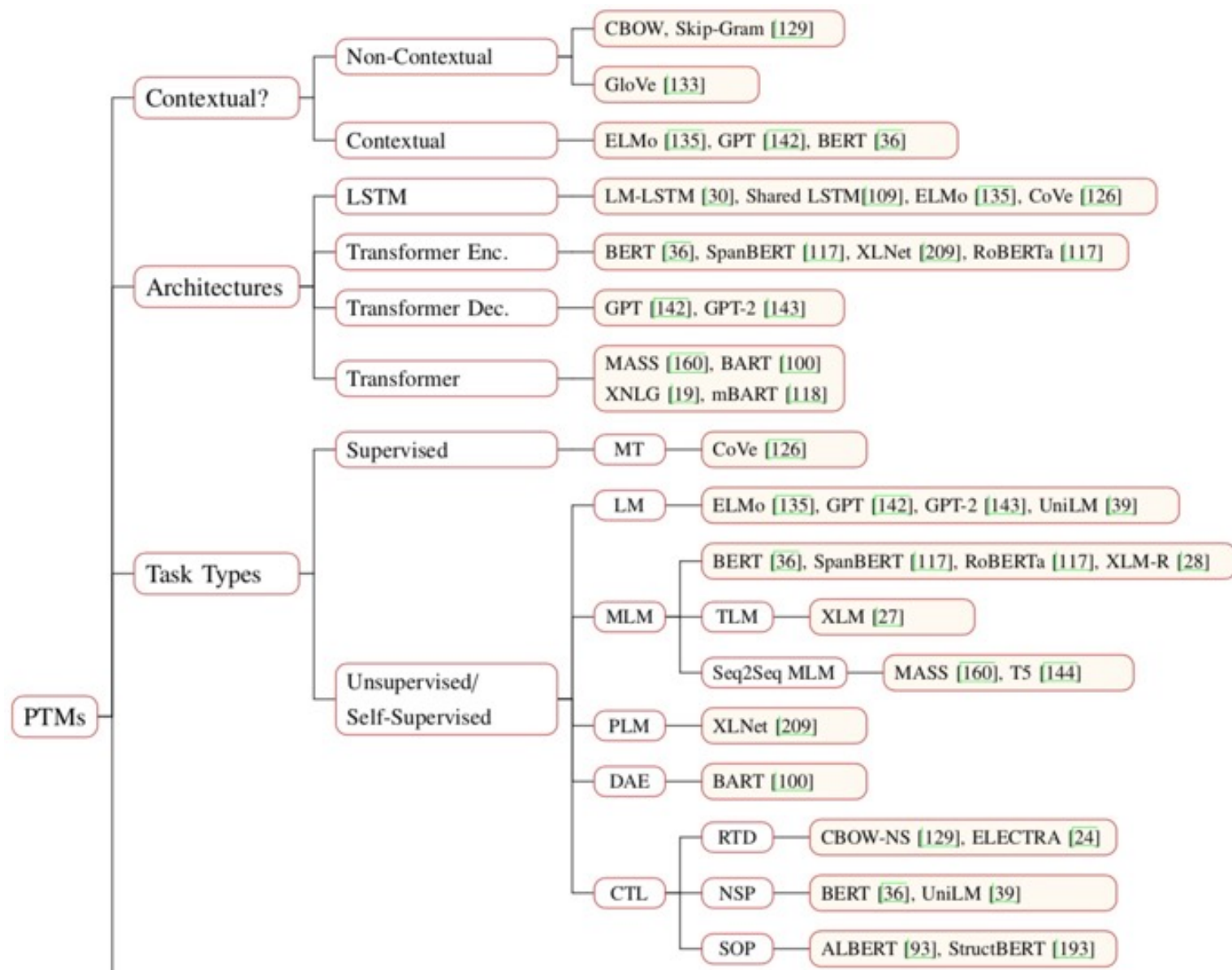




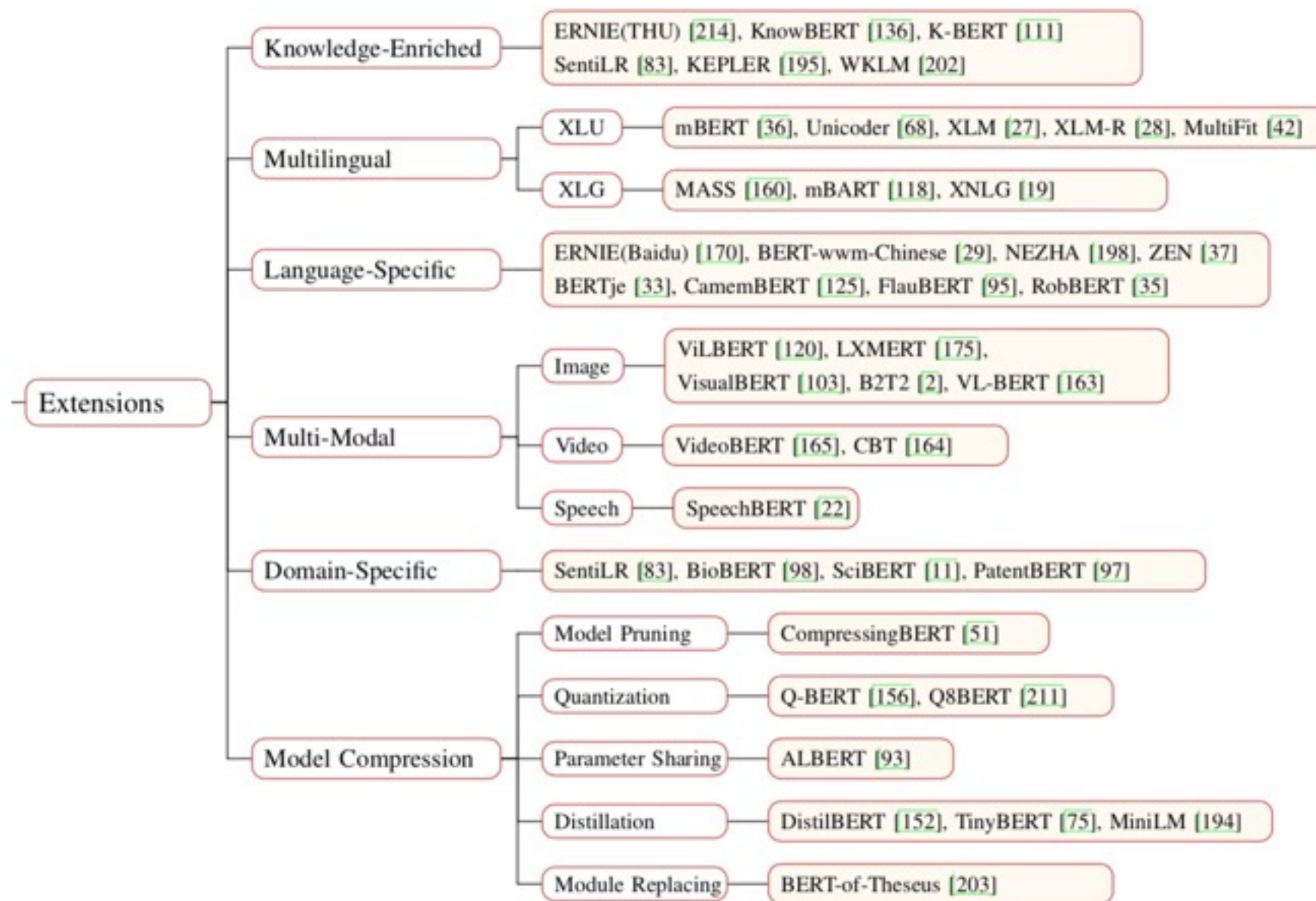
# Transformer Models



# Pre-trained Models (PTM)

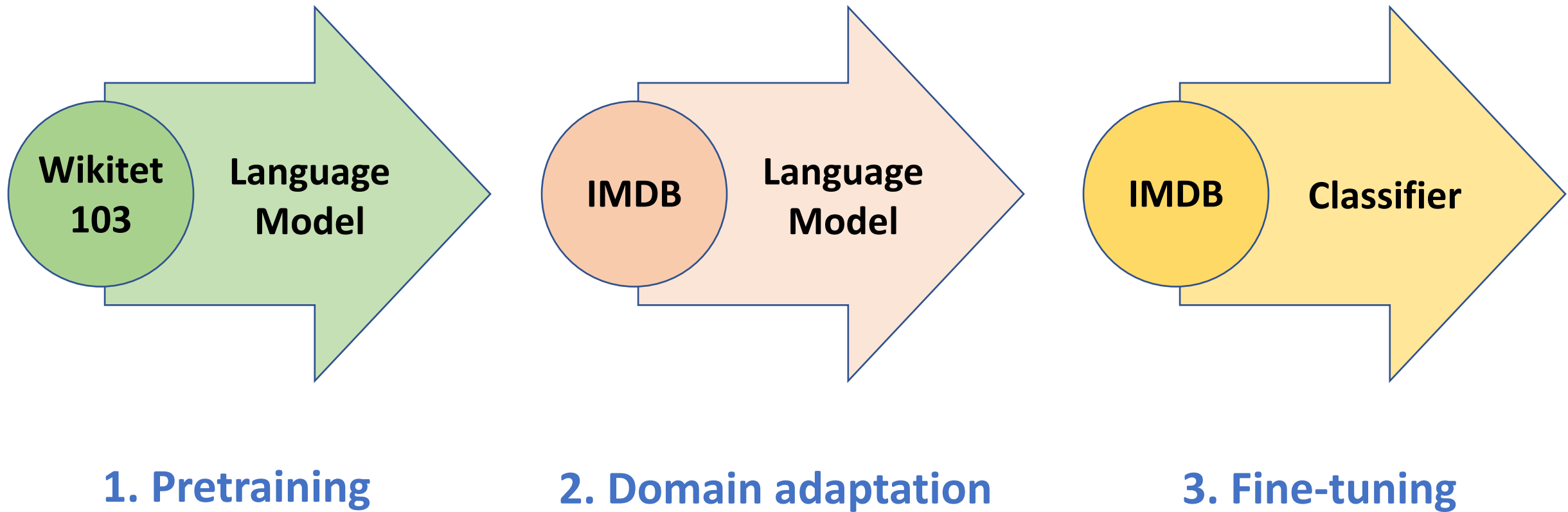


# Pre-trained Models (PTM)



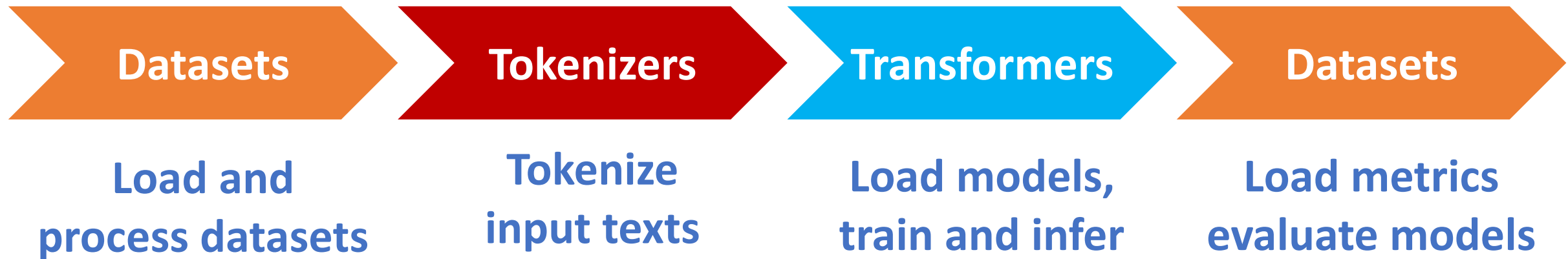
# ULMFiT: 3 Steps

## Transfer Learning in NLP



# A typical pipeline for training transformer models

with the Datasets, Tokenizers, and Transformers libraries



# Few-Shot Learning (FSL)

## Typical Scenarios

- **Acting as a test bed for learning like human**
- **Learning for rare cases**
- **Reducing data gathering effort and computational cost**

# Few-Shot Learning (FSL)

- **Few-Shot Learning (FSL)** is a sub-area in machine learning.
- **Machine Learning Definition**
  - A computer program is said to learn from **experience E** with respect to some classes of **task T** and **performance measure P** if its performance can improve with E on T measured by P.
    - Example: **Image classification task (T)**, a machine learning program can improve its **classification accuracy (P)** through **E** obtained by training on a large number of **labeled images** (e.g., the ImageNet data set).

# Machine Learning

task $T$	experience $E$	performance $P$
image classification [73]	large-scale labeled images for each class	classification accuracy
the ancient game of Go [120]	a database containing around 30 million recorded moves of human experts and self-play records	winning rate



# Few-Shot Learning (FSL)

- **Few-shot Learning (FSL)** is a type of machine learning problems (specified by E, T, and P), where **E contains only a limited number of examples** with supervised information for the target T.
  - Existing FSL problems are mainly supervised learning problems.
  - Few-shot classification learns classifiers given only **a few labeled examples of each class**.
    - image classification
    - sentiment classification from short text
    - object recognition

# Few-Shot Learning (FSL)

- Few-shot classification learns a classifier  $h$ , which predicts label  $y_i$  for each input  $x_i$ .
- Usually, one considers the  *$N$ -way- $K$ -shot* classification, in which  $D_{train}$  contains  $I = KN$  examples from  $N$  classes each with  $K$  examples

# Few-Shot Learning (FSL)

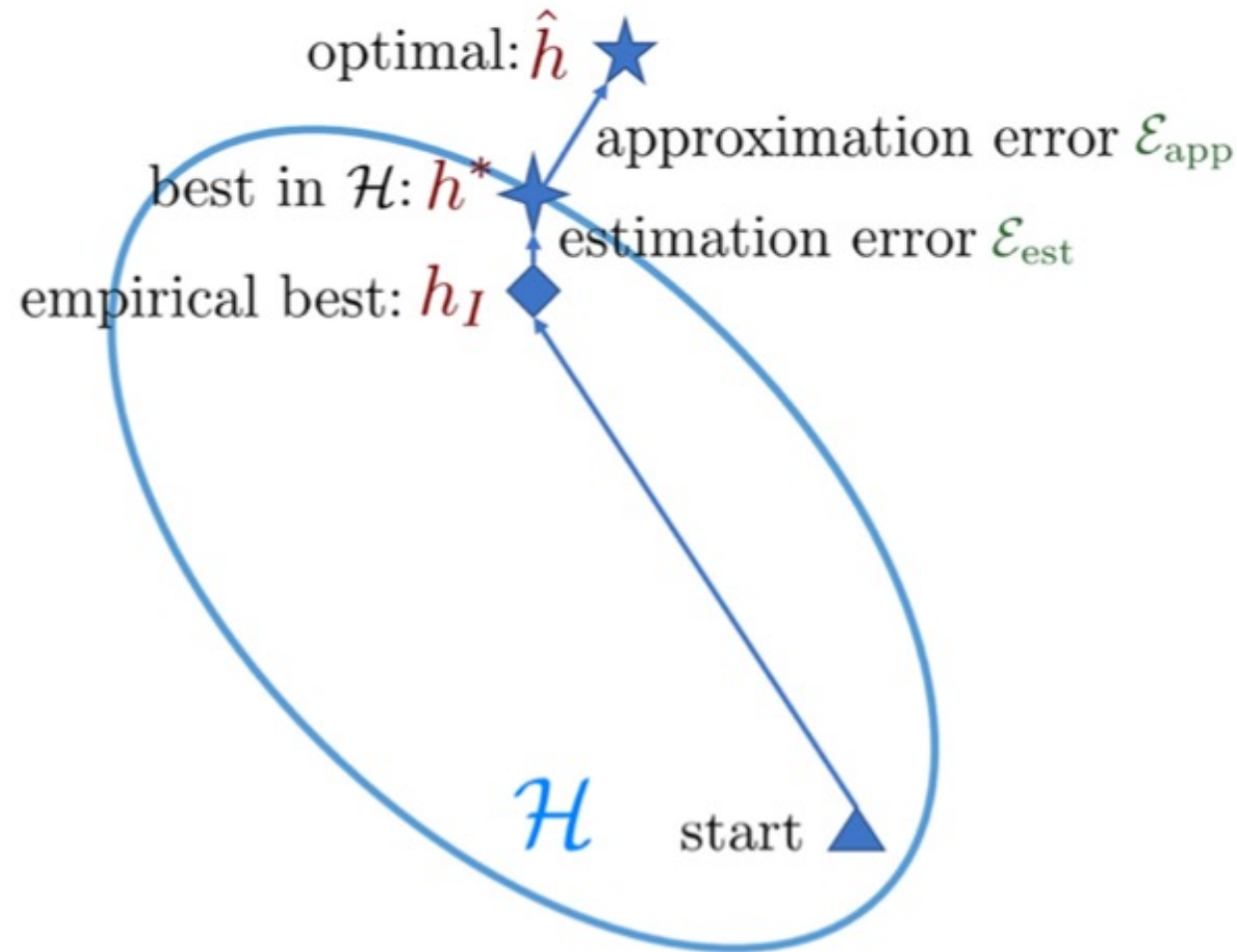
- **Few-Shot Learning (FSL)**
  - $K = 10 \sim 100$  examples
- **One-Shot Learning (1SL)**
  - $K = 1$  example
- **Zero-Shot Learning (0SL)(ZSL)**
  - $K = 0$

# Few-Shot Learning (FSL)

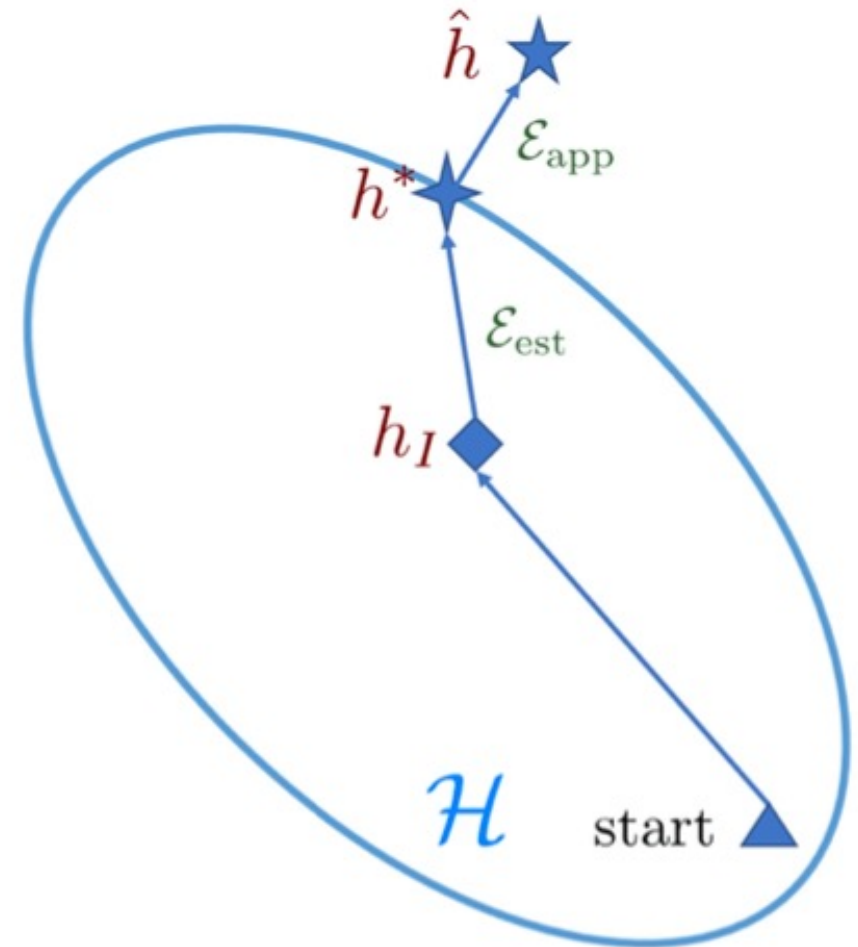
task $T$	experience $E$		performance $P$
	supervised information	prior knowledge	
character generation [76]	a few examples of new character	pre-learned knowledge of parts and relations	pass rate of visual Turing test
drug toxicity discovery [4]	new molecule's limited assay	similar molecules' assays	classification accuracy
image classification [70]	a few labeled images for each class of the target $T$	raw images of other classes, or pre-trained models	classification accuracy

# Few-Shot Learning (FSL)

## Comparison of learning with sufficient and few training samples



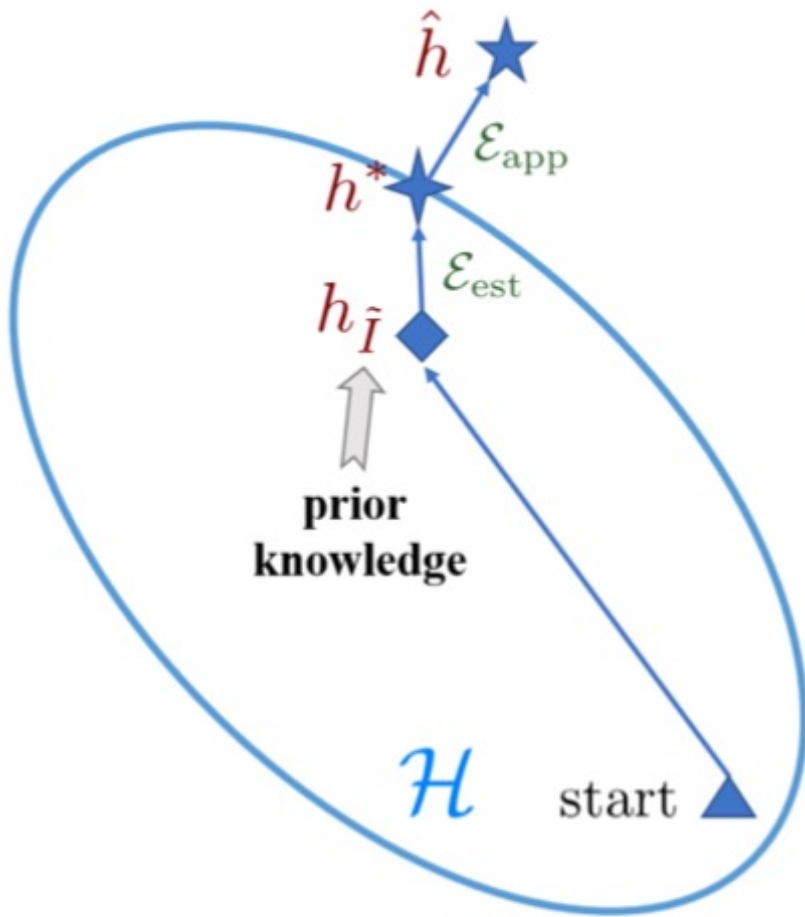
(a) Large  $I$



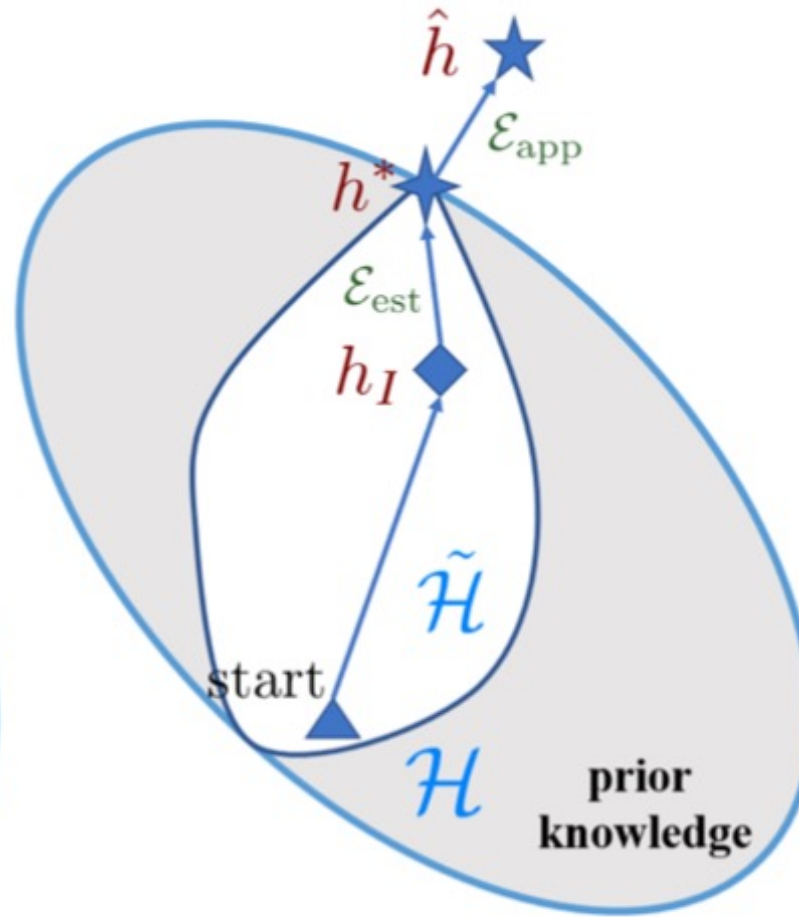
(b) Small  $I$

# Few-Shot Learning (FSL)

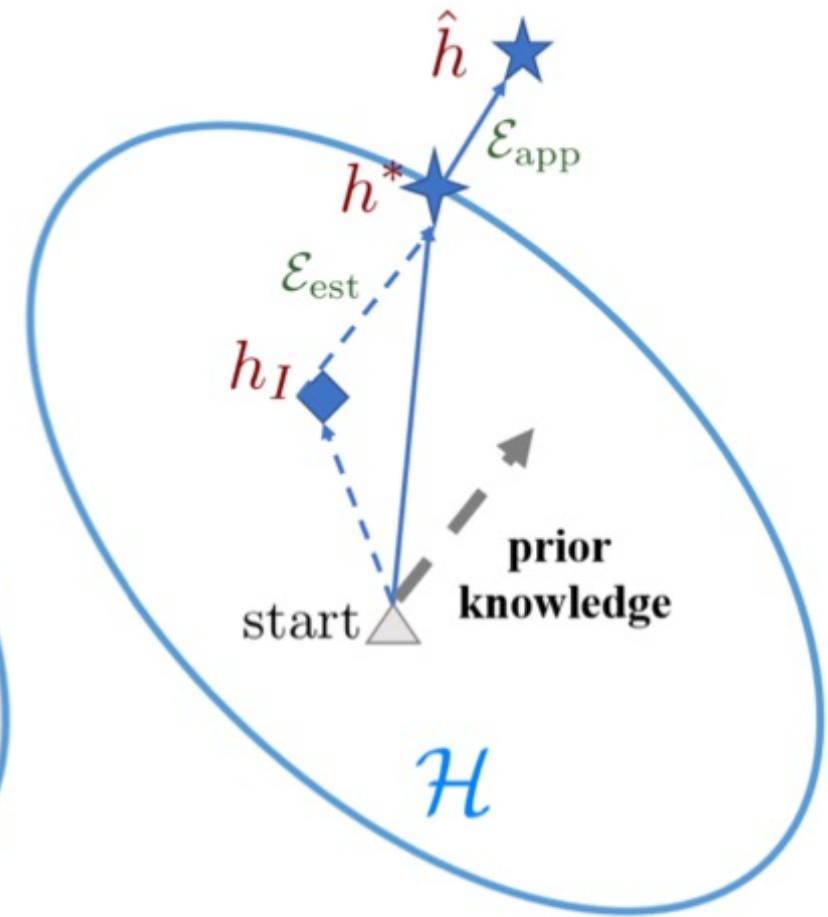
Different perspectives on how FSL methods solve the few-shot problem



(a) Data.



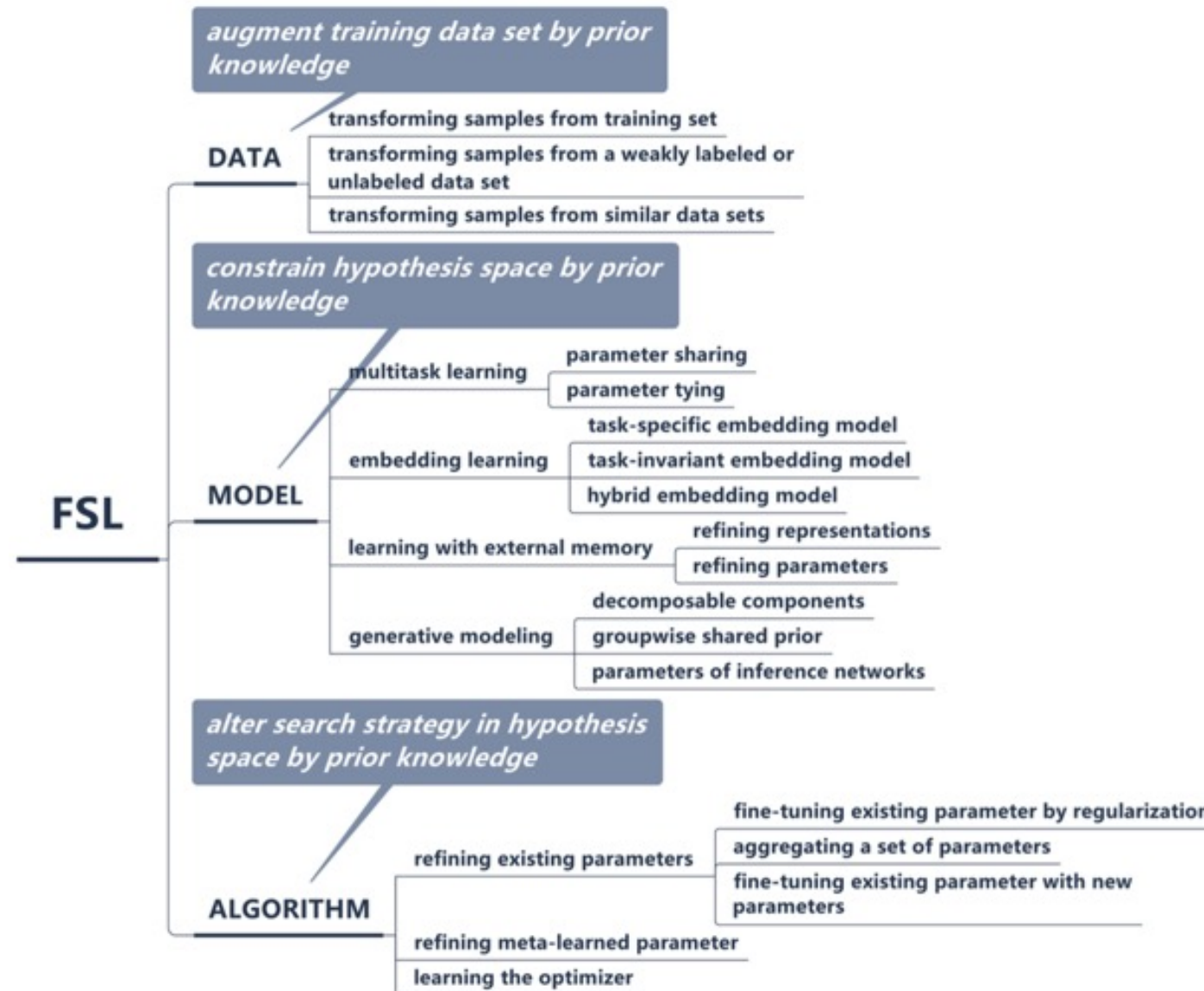
(b) Model.



(c) Algorithm.

# Few-Shot Learning (FSL)

## A taxonomy of FSL methods





# Few-Shot Learning (FSL)

*augment training data set by prior knowledge*

**DATA**

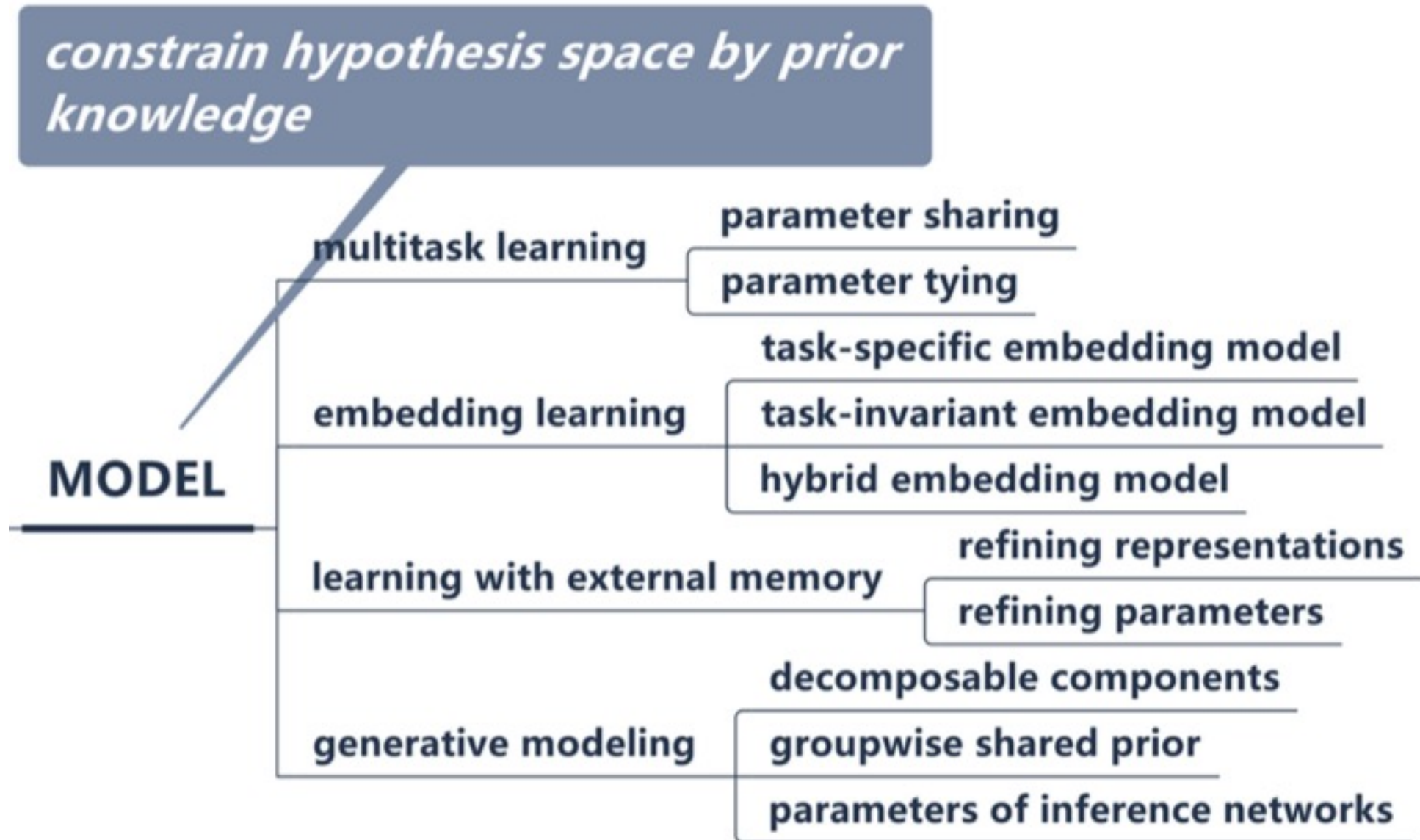
transforming samples from training set

transforming samples from a weakly labeled or unlabeled data set

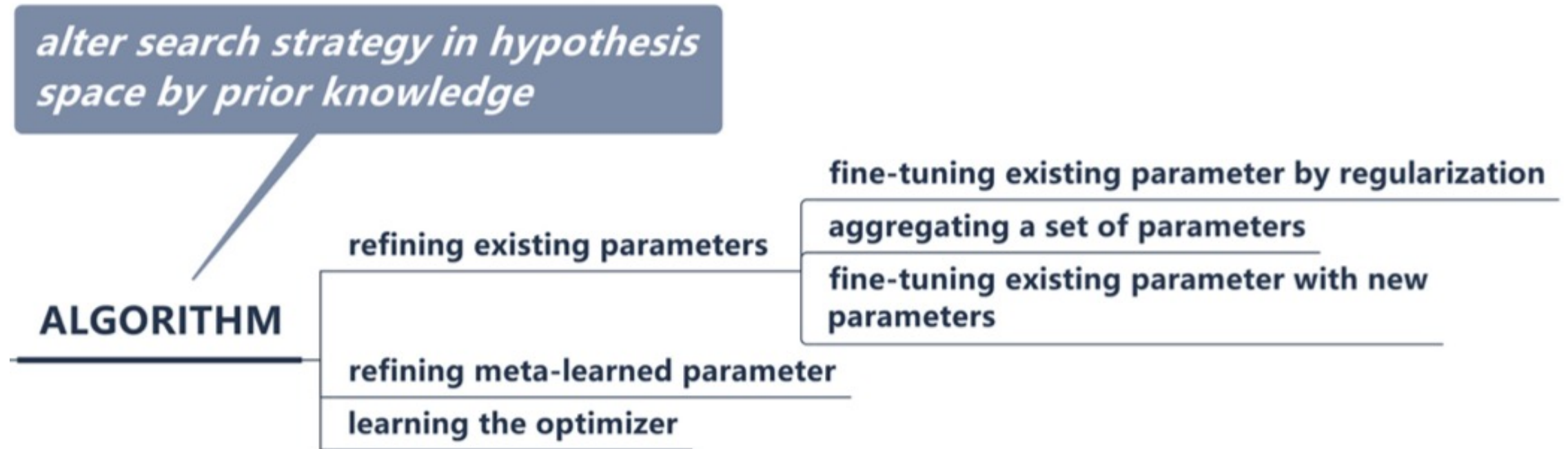
transforming samples from similar data sets



# Few-Shot Learning (FSL)

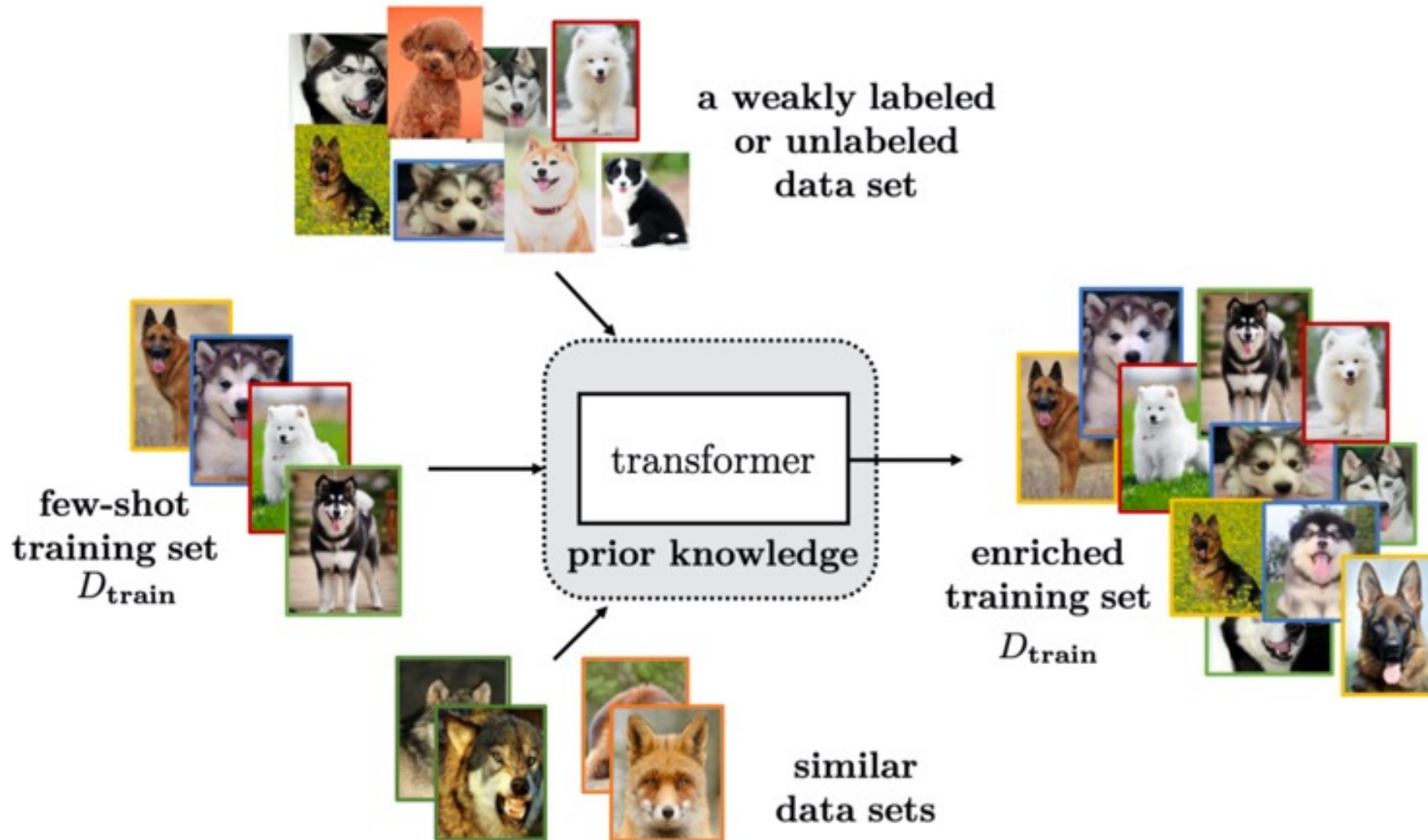


# Few-Shot Learning (FSL)



# Few-Shot Learning (FSL)

## Solving the FSL problem by data augmentation



# Few-Shot Learning (FSL)

## Characteristics for FSL Methods Focusing on the Data Perspective

category	input $(x, y)$	transformer $t$	output $(\tilde{x}, \tilde{y})$
transforming samples from $D_{\text{train}}$	original $(x_i, y_i)$	learned transformation function on $x_i$	$(t(x_i), y_i)$
transforming samples from a weakly labeled or unlabeled data set	weakly labeled or unlabeled $(\bar{x}, -)$	a predictor trained from $D_{\text{train}}$	$(\bar{x}, t(\bar{x}))$
transforming samples from similar data sets	samples $\{(\hat{x}_j, \hat{y}_j)\}$ from similar data sets	an aggregator to combine $\{(\hat{x}_j, \hat{y}_j)\}$	$(t(\{\hat{x}_j\}), t(\{\hat{y}_j\}))$

The transformer  $t(\cdot)$  takes input  $(x, y)$  and returns synthesized sample  $(\tilde{x}, \tilde{y})$  to augment the few-shot  $D_{\text{train}}$ .



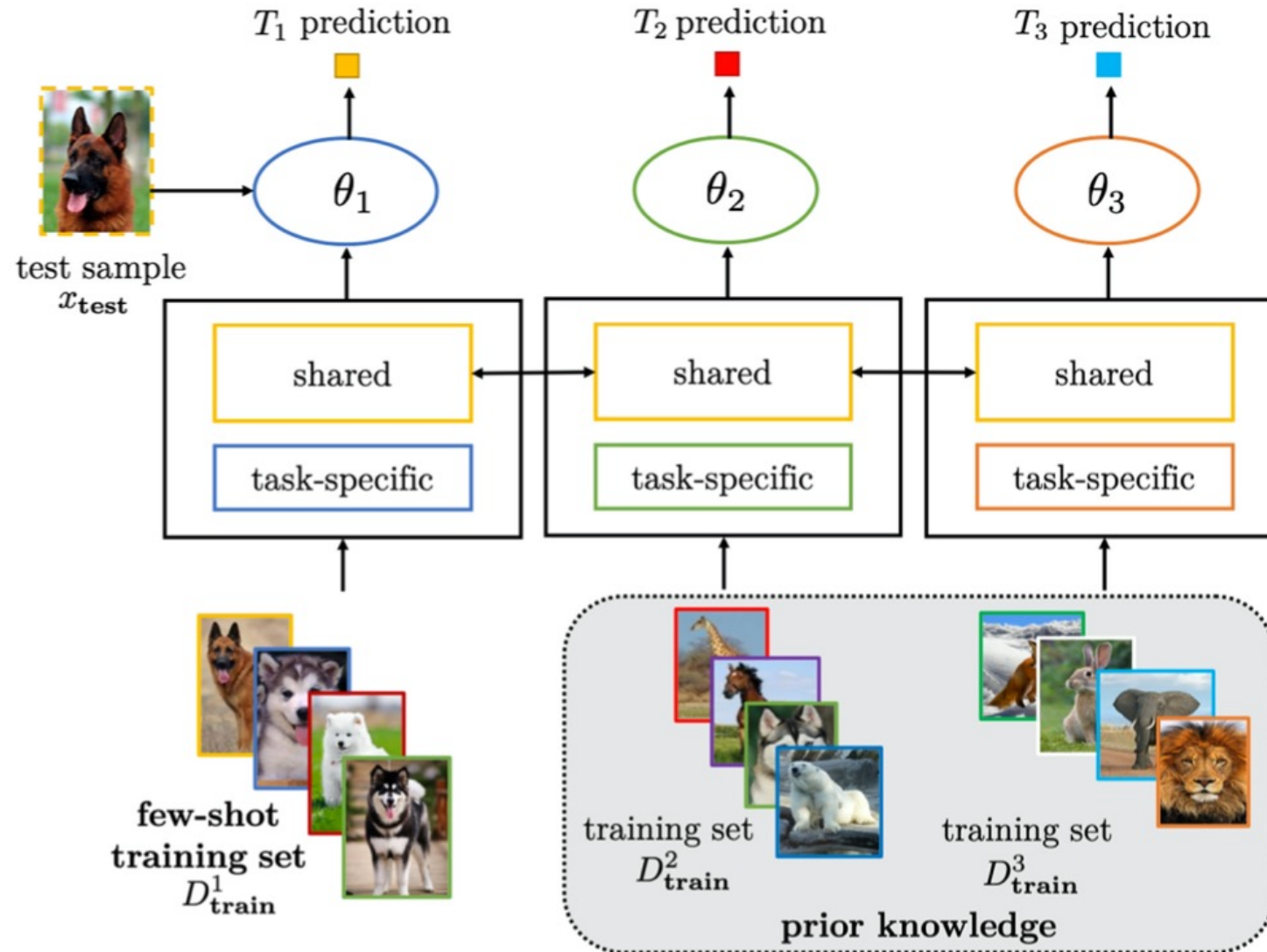
# Few-Shot Learning (FSL)

## Characteristics for FSL Methods Focusing on the Model Perspective

strategy	prior knowledge	how to constrain $\mathcal{H}$
multitask learning	other $T$ 's with their data sets $D$ 's	share/tie parameter
embedding learning	embedding learned from/together with other $T$ 's	project samples to a smaller embedding space in which similar and dissimilar samples can be easily discriminated
learning with external memory	embedding learned from other $T$ 's to interact with memory	refine samples using key-value pairs stored in memory
generative modeling	prior model learned from other $T$ 's	restrict the form of distribution

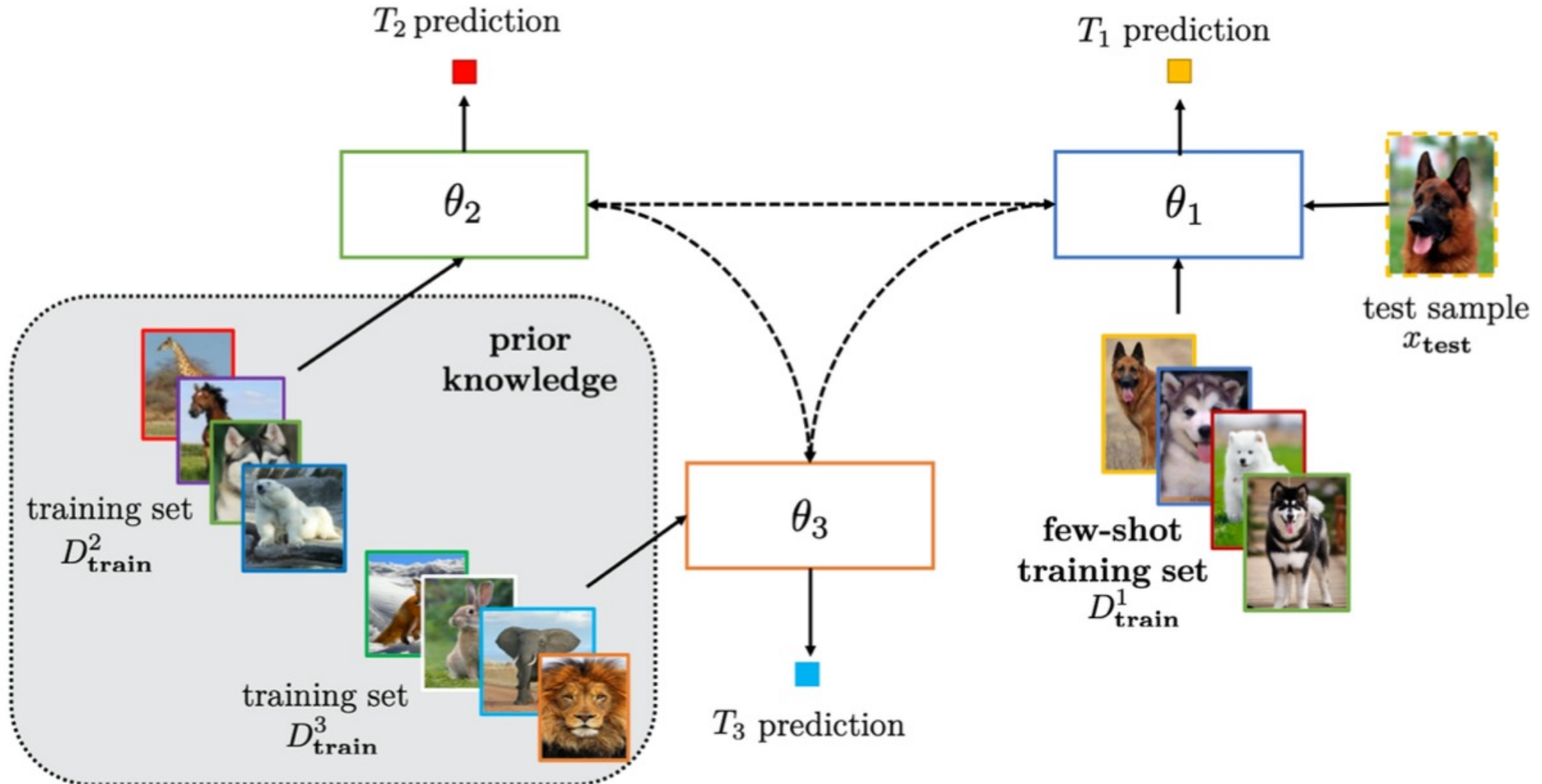
# Few-Shot Learning (FSL)

Solving the FSL problem by multitask learning with parameter sharing



# Few-Shot Learning (FSL)

Solving the FSL problem by multitask learning with parameter tying



# Few-Shot Learning (FSL)

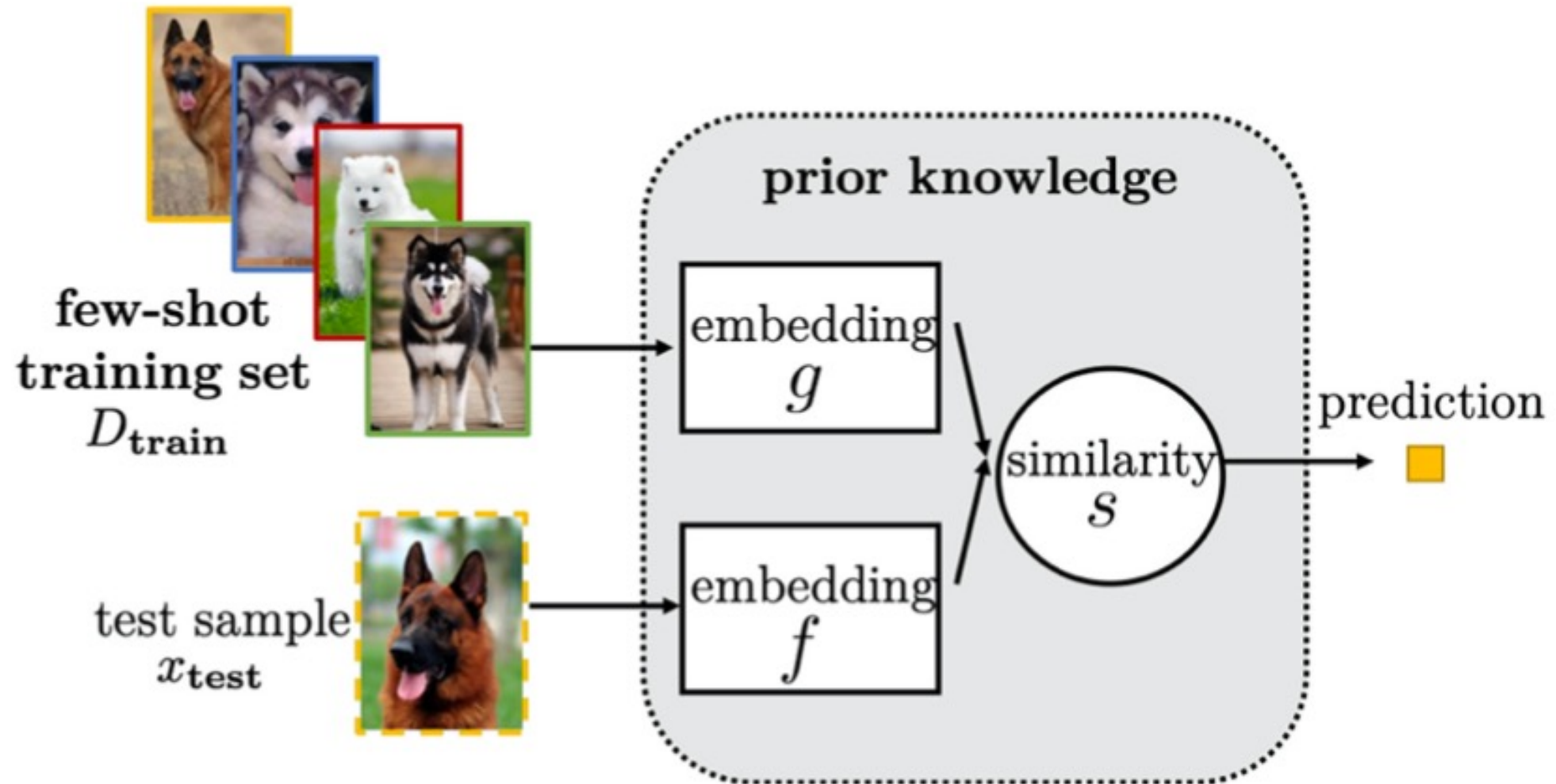
## Characteristics of Embedding Learning Methods

category	method	embedding function $f$ for $x_{\text{test}}$	embedding function $g$ for $D_{\text{train}}$	similarity measure $s$
task-specific	mAP-DLM/SSVM[130]	CNN	the same as $f$	cosine similarity
task-invariant	class relevance pseudo-metric [36]	kernel	the same as $f$	squared $\ell_2$ distance
	convolutional siamese net [70]	CNN	the same as $f$	weighted $\ell_1$ distance
	Micro-Set[127]	logistic projection	the same as $f$	$\ell_2$ distance
	Matching Nets [138]	CNN, LSTM	CNN, biLSTM	cosine similarity
	resLSTM [4]	GNN, LSTM	GNN, LSTM	cosine similarity
	Active MN [8]	CNN	biLSTM	cosine similarity
	SSMN [24]	CNN	another CNN	learned distance
	ProtoNet [121]	CNN	the same as $f$	squared $\ell_2$ distance
	semi-supervised ProtoNet[108]	CNN	the same as $f$	squared $\ell_2$ distance
	PMN [141]	CNN, LSTM	CNN, biLSTM	cosine similarity
	ARC [119]	LSTM, biLSTM	the same as $f$	-
	Relation Net [126]	CNN	the same as $f$	-
	GNN [115]	CNN, GNN	the same as $f$	learned distance
	TPN [84]	CNN	the same as $f$	Gaussian similarity
	SNAIL [91]	CNN	the same as $f$	-
hybrid	Learnet [14]	adaptive CNN	CNN	weighted $\ell_1$ distance
	DCCN [162]	adaptive CNN	CNN	-
	R2-D2 [13]	adaptive CNN	CNN	-
	TADAM [100]	adaptive CNN	the same as $f$	squared $\ell_2$ distance



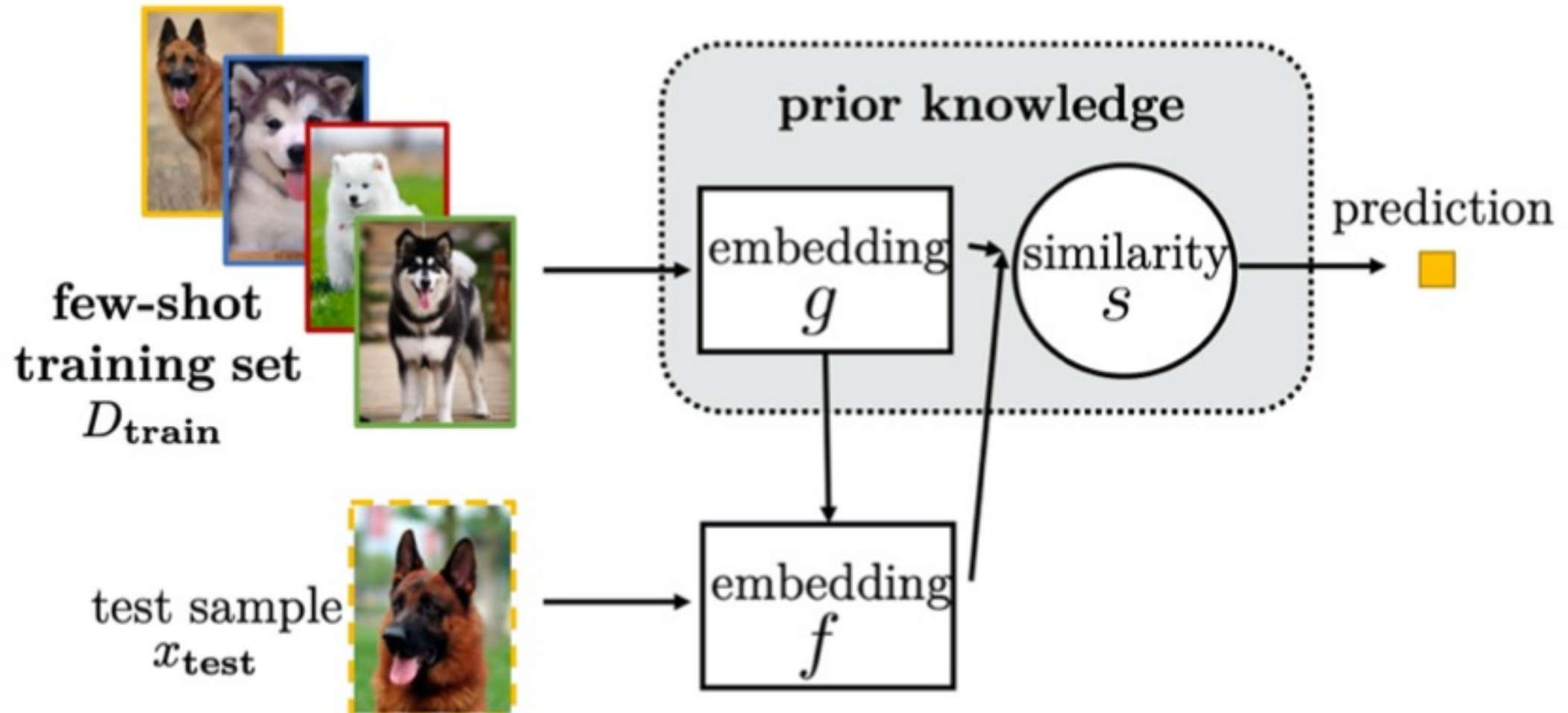
# Few-Shot Learning (FSL)

Solving the FSL problem by task-invariant embedding model



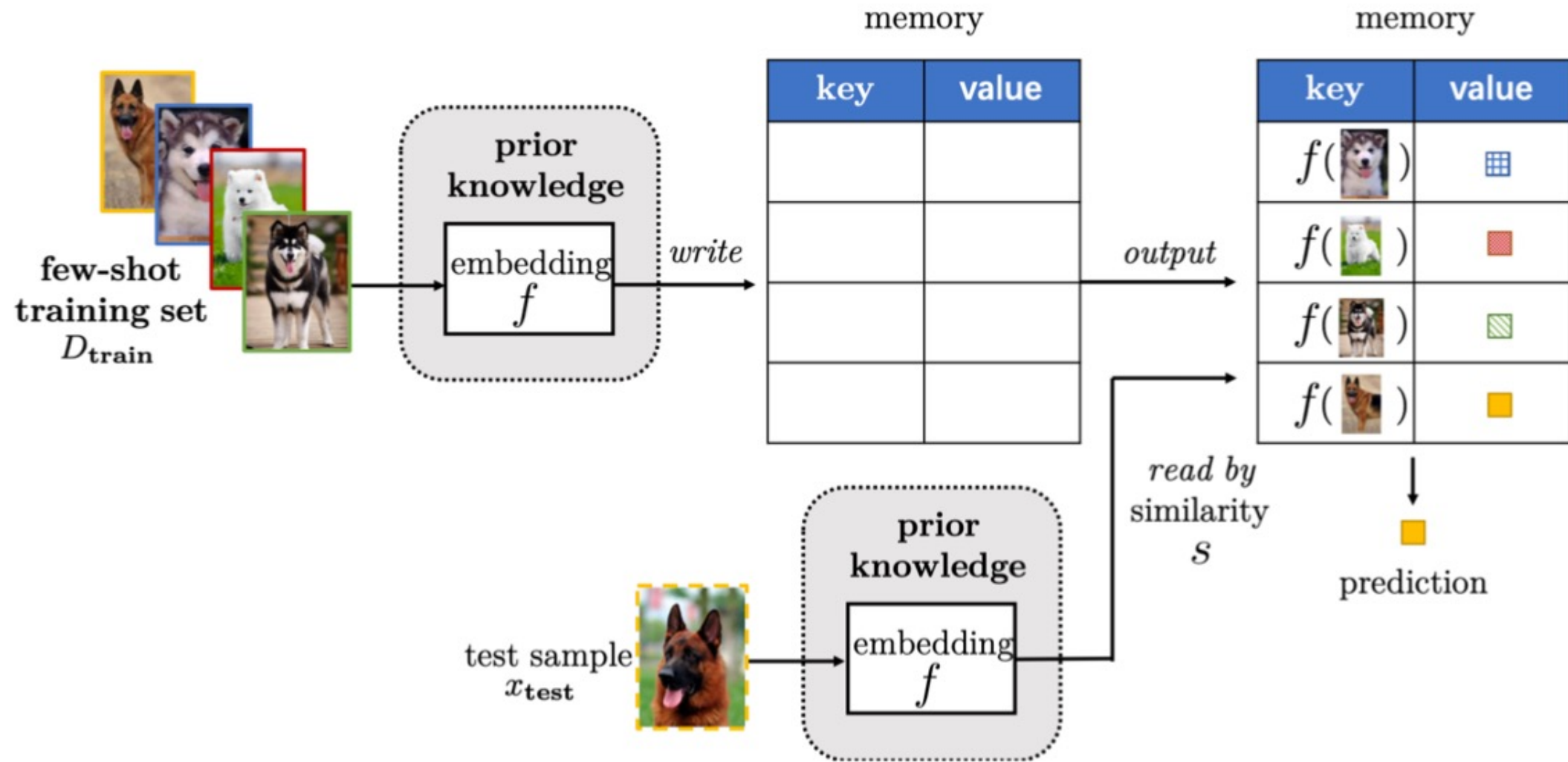
# Few-Shot Learning (FSL)

Solving the FSL problem by hybrid embedding model



# Few-Shot Learning (FSL)

Solving the FSL problem by learning with external memory



# Few-Shot Learning (FSL)

## Characteristics of FSL Methods

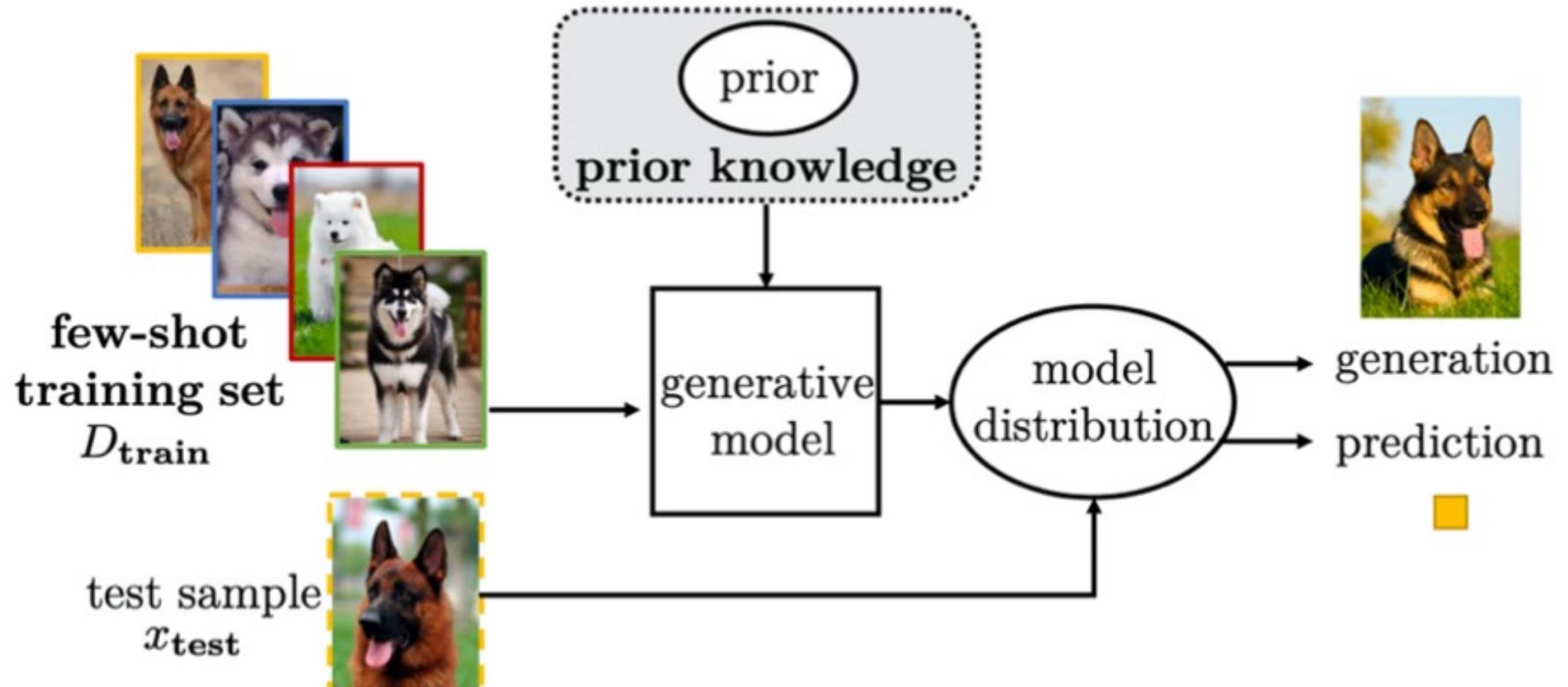
### Based on Learning with External Memory

category	method	memory $M$		similarity $s$
		key $M_{\text{key}}$	value $M_{\text{value}}$	
refining representations	MANN [114]	$f(x_i, y_{i-1})$	$f(x_i, y_{i-1})$	cosine similarity
	APL [104]	$f(x_i)$	$y_i$	squared $\ell_2$ distance
	abstraction memory [149]	$f(x_i)$	word embedding of $y_i$	dot product
	CMN [164]	$f(x_i)$	$y_i, \text{age}$	dot product
	life-long memory [65]	$f(x_i)$	$y_i, \text{age}$	cosine similarity
	Mem2Vec [125]	$f(x_i)$	word embedding of $y_i, \text{age}$	dot product
refining parameters	MetaNet [96]	$f(x_i)$	fast weight	cosine similarity
	CSNs [97]	$f(x_i)$	fast weight	cosine similarity
	MN-Net [22]	$f(x_i)$	$y_i$	dot product

Here,  $f$  is an embedding function usually pre-trained by CNN or LSTM.

# Few-Shot Learning (FSL)

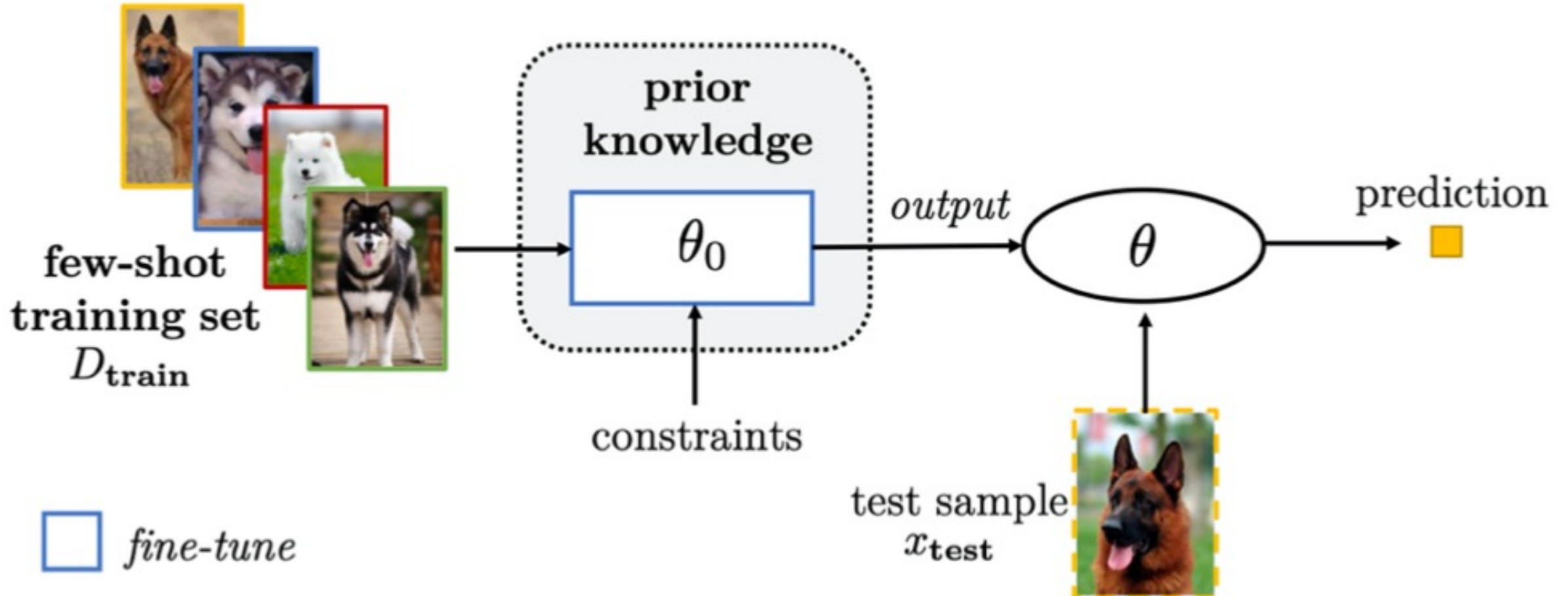
Solving the FSL problem by generative modeling





# Few-Shot Learning (FSL)

Solving the FSL problem by fine-tuning existing parameter  $\theta_0$  by regularization



# Few-Shot Learning (FSL)

## Characteristics for FSL Methods

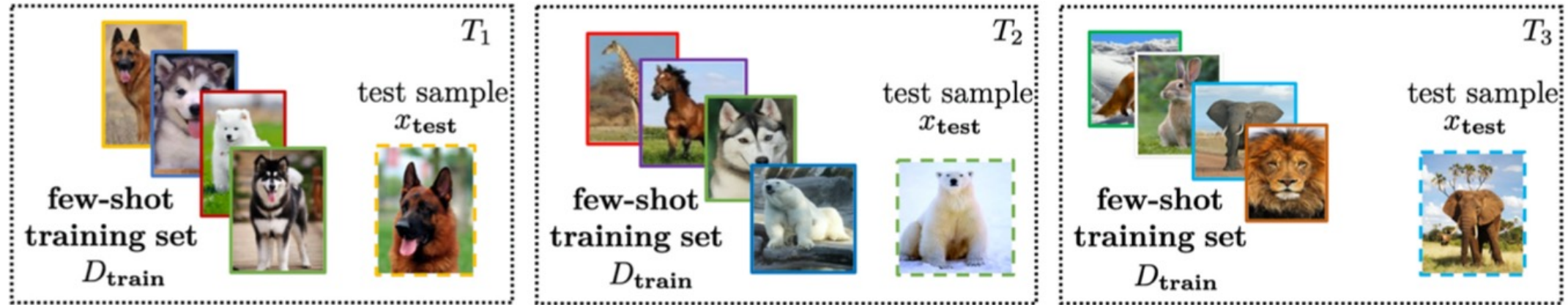
### Focusing on the Algorithm Perspective

strategy	prior knowledge	how to search $\theta$ of the $h^*$ in $\mathcal{H}$
refining existing parameters	learned $\theta_0$	refine $\theta_0$ by $D_{\text{train}}$
refining meta-learned parameters	meta-learner	refine $\theta_0$ by $D_{\text{train}}$
learning the optimizer	meta-learner	use search steps provided by the meta-learner

# Few-Shot Learning (FSL)

## Solving the FSL problem by meta-learning

meta-training tasks  $T_s$ 's



meta-testing tasks  $T_t$ 's

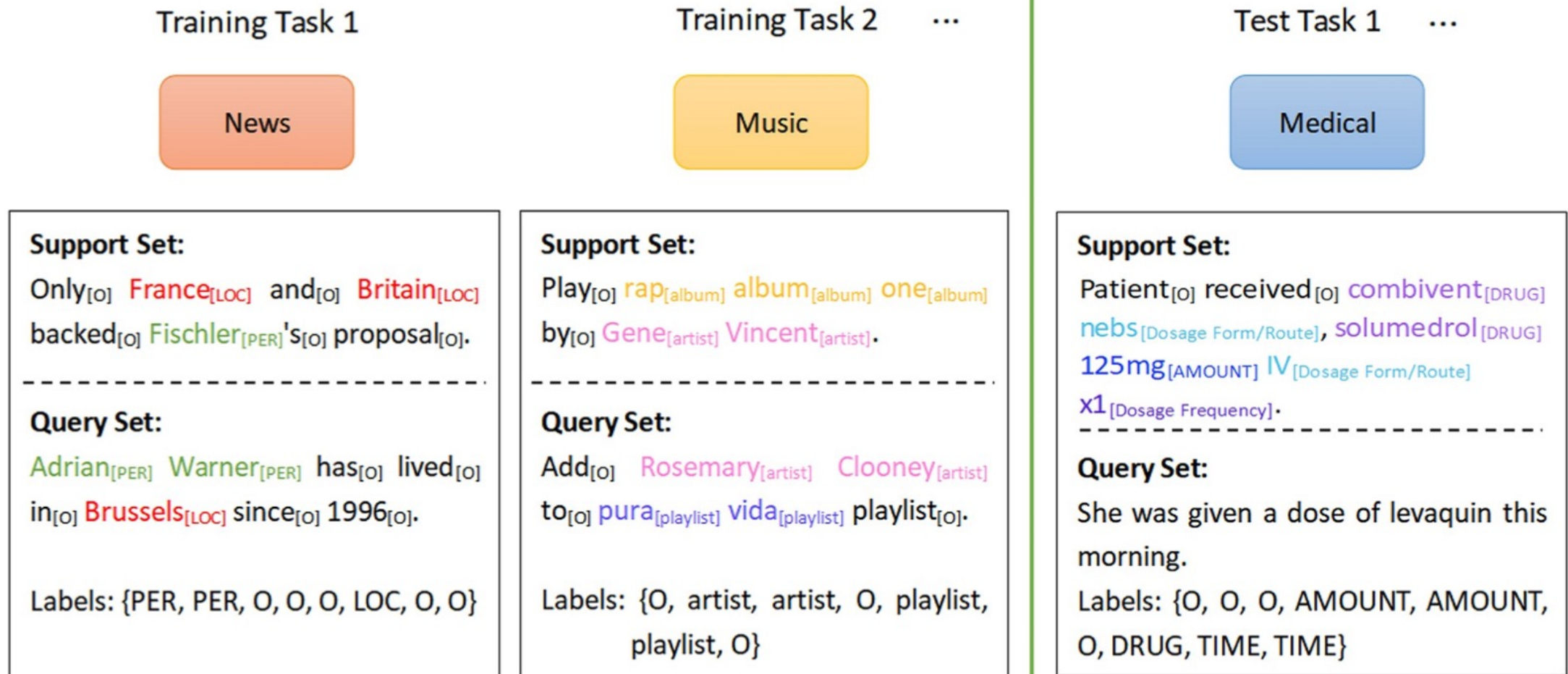




# Few-Shot Learning (FSL)

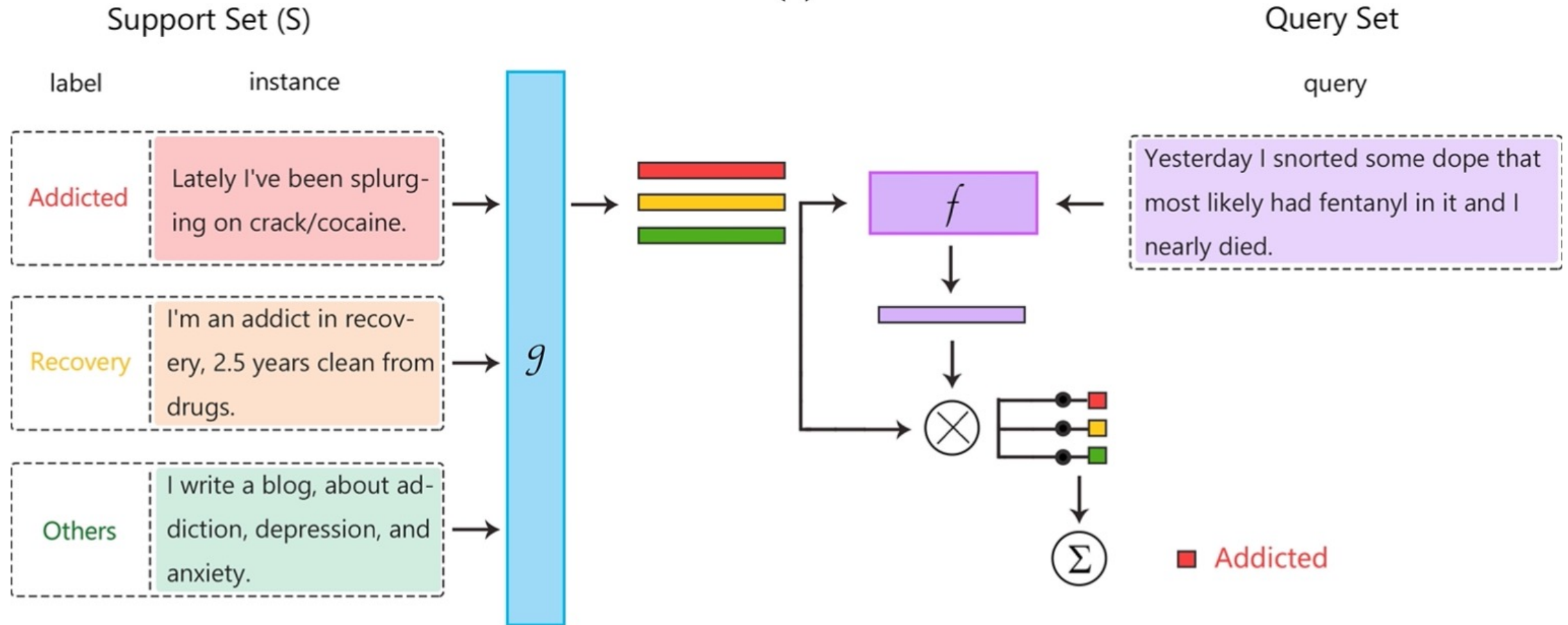
## Meta-learning

Each task mimics the few-shot scenario, and can be completely non-overlapping.  
Support sets are used to train; query sets are used to evaluate the model



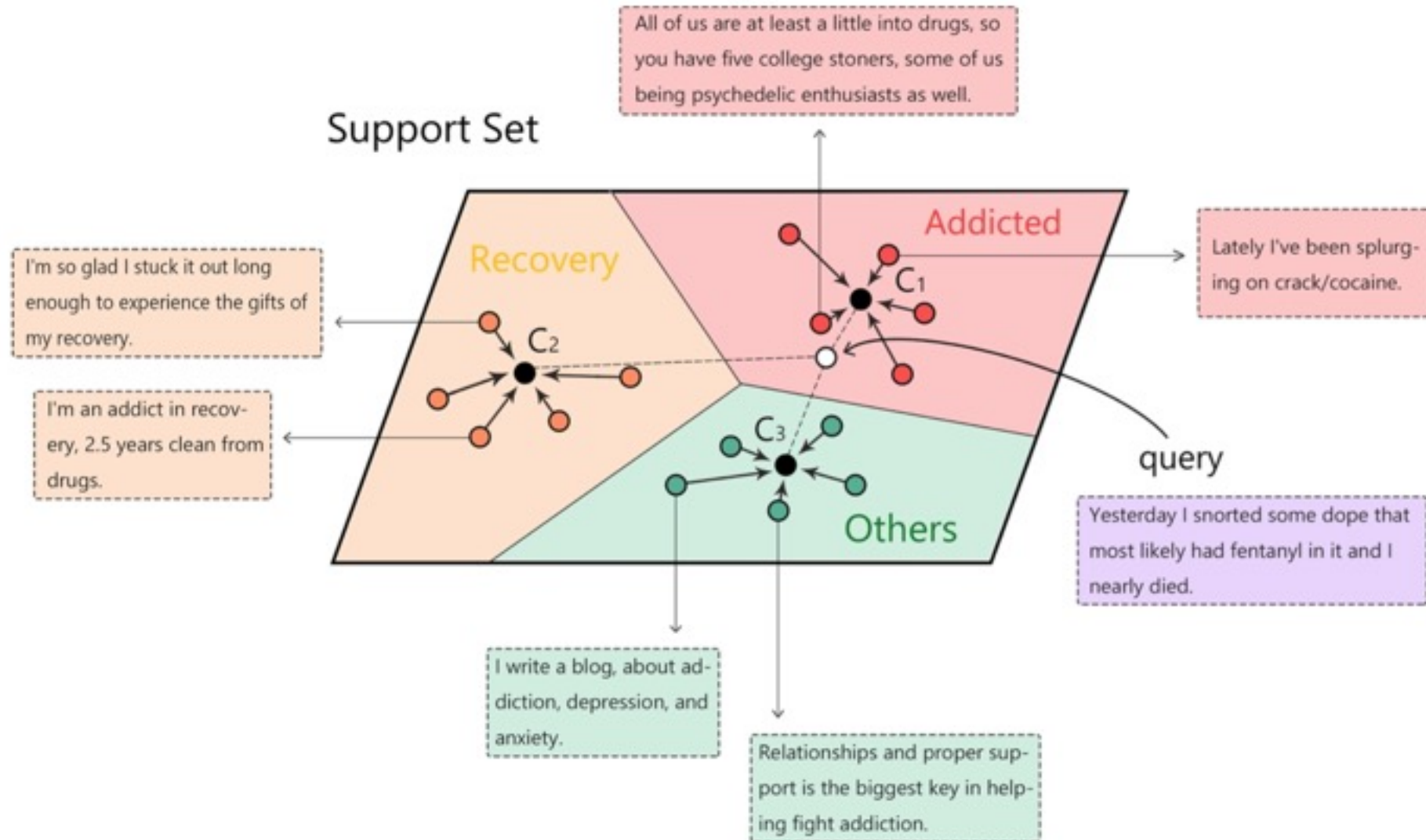
# Few-Shot Learning (FSL)

## Matching networks



# Few-Shot Learning (FSL)

## Prototypical network



# Few-Shot Learning (FSL) for medical text

Study	Year	Data source	Research aim	Size of training set	Number of entities / classes	Entity type of training domain	Entity type of test domain
Alicia Lara-Clares and Ana Garcia-Serrano <sup>44</sup>	2019	MEDDOCAN shared task dataset <sup>45</sup>	NER	500 clinical cases, with no reconstruction	29	Clinical	Clinical
Ferré et al. <sup>46</sup>	2019	BB-norm dataset from the Bacteria Biotope 2019 Task <sup>47</sup>	Entity Normalization	Original dataset with no reconstruction and zero-shot	Not mentioned *	Biological	Biological
Hou et al. <sup>48</sup>	2020	Snips dataset <sup>49</sup>	Slot Tagging (NER)	1-shot and 5-shot	7	Six of Weather, Music, PlayList, Book (including biomedical), Search Screen (including biomedical), Restaurant and Creative Work.	The remaining one
Sharaf et al. <sup>50</sup>	2020	ten different datasets collected from the Open Parallel Corpus (OPUS) <sup>51</sup>	Neural Machine Translation (NMT)	Sizes ranging from 4k to 64k training words (200 to 3200 sentences), but reconstructed	N/A †	Bible, European Central Bank, KDE, Quran, WMT news test sets, Books, European Medicines Agency (EMA), Global Voices, Medical (ufal-Med), TED talks	Bible, European Central Bank, KDE, Quran, WMT news test sets, Books, European Medicines Agency (EMA), Global Voices, Medical (ufal-Med), TED talks
Lu et al. <sup>52</sup>	2020	MIMIC II <sup>22</sup> and MIMIC III <sup>23</sup> , and EU legislation dataset <sup>53</sup>	Multi-label Text Classification	5-shot for MIMIC II and III, 50-shot for EU legislation	MIMIC II: 9 MIMIC III: 15 EU legislation: 5	Medical	Medical

# Few-Shot Learning (FSL) for medical text

Study	Year	Data source	Research aim	Size of training set	Number of entities / classes	Entity type of training domain	Entity type of test domain
Lu et al. <sup>80</sup>	2021	Constructed and shared a novel dataset <sup>††</sup> based on Weibo for the research of few-shot rumor detection, and use PHEME dataset <sup>81</sup>	Rumor Detection (NER)	For the Weibo dataset: 2-way 3-event 5-shot 9-query; for PHEME dataset: 2-way 2-event 5-shot 9-query	Weibo: 14 PHEME: 5	Source posts and comments from Sina Weibo related to COVID-19	Source posts and comments from Sina Weibo related to COVID-19
Ma et al. <sup>82</sup>	2021	CCLE, CERES-correctedCRISPR gene disruption scores, GDSC1000 dataset, PDTC dataset <sup>††</sup> and PDX dataset <sup>††</sup>	Drug-response Predictions	1-shot, 2-shot, 5-shot and 10-shot	N/A <sup>†</sup>	Biomedical	Biomedical
Kormilitzin et al. <sup>83</sup>	2021	MIMIC-III <sup>23</sup> and UK-CRIS datasets <sup>30,31</sup>	NER	25%, 50%, 75% and 100% of the training set, with no reconstruction	7	Electronic health record	Electronic health record
Guo et al. <sup>84</sup>	2021	Abstracts of biomedical literatures (from relation extraction task of BioNLP Shared Task 2011 and 2019 <sup>47</sup> ) and structured biological datasets	NER	100%, 75%, 50%, 25%, 0% of training set, with no reconstruction	Not mentioned *	Biomedical entities	Biomedical entities
Lee et al. <sup>85</sup>	2021	COVID19-Scientific <sup>86</sup> , COVID19-Social <sup>87</sup> (fact-checked by journalists from a website called Politi-fact.com), FEVER <sup>88</sup> (Fact Extraction and Verification, generated by altering sentences extracted from Wikipedia to promote research on fact-checking systems)	Fact-Checking (close to Text Classification)	2-shot, 10-shot and 50-shot	Not mentioned *	Facts about COVID-19	Facts about COVID-19



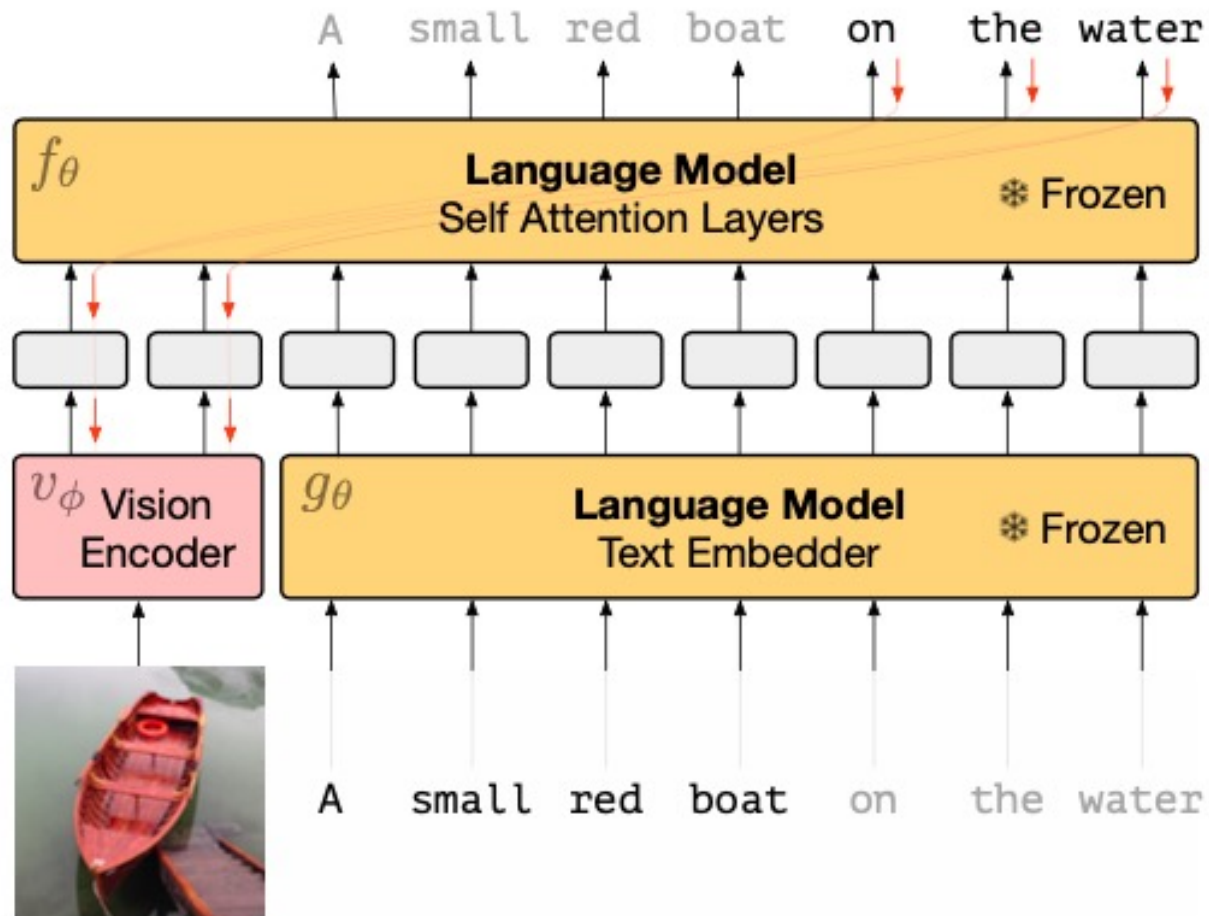
# Multimodal Few-Shot Learning with Frozen Language Models



Curated samples with about five seeds required to get past well-known language model failure modes of either repeating text for the prompt or emitting text that does not pertain to the image.

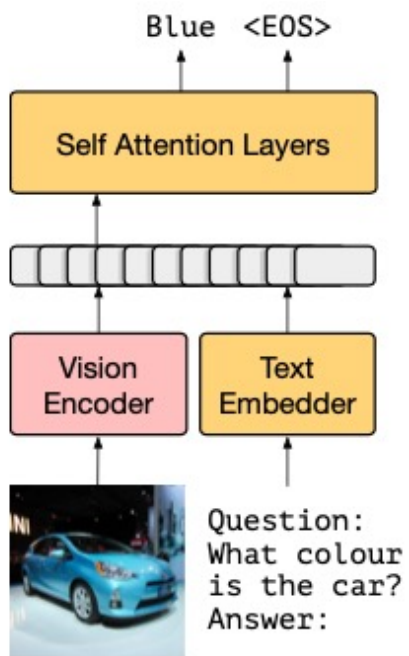
These samples demonstrate the ability to generate open-ended outputs that adapt to both images and text, and to make use of facts that it has learned during language-only pre-training.

# Multimodal Few-Shot Learning with Frozen Language Models

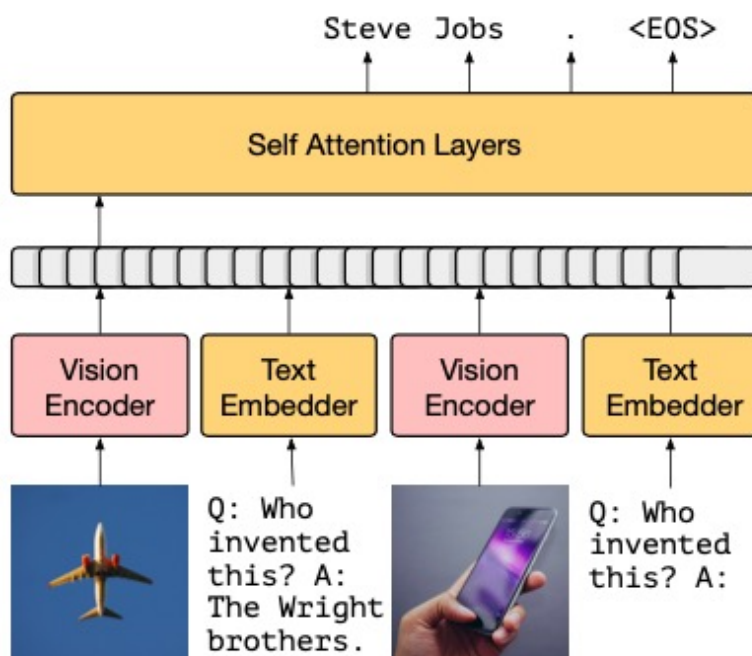


Gradients through a frozen language model's self attention layers are used to train the vision encoder.

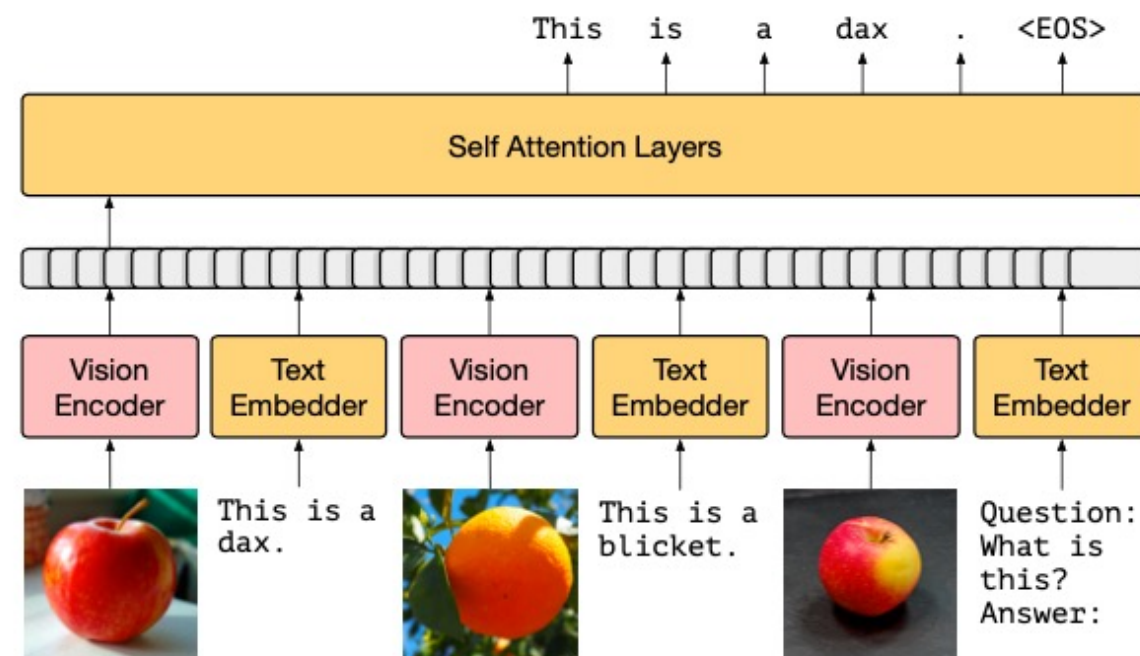
# Multimodal Few-Shot Learning with Frozen Language Models



(a) 0-shot VQA



(b) 1-shot outside-knowledge VQA



(c) Few-shot image classification

Inference-Time interface for *Frozen*. The figure demonstrates how we can support (a) visual question answering, (b) outside-knowledge question answering and (c) few-shot image classification via in-context learning.

Source: Maria Tsimpoukelli, Jacob L. Menick, Serkan Cabi, S. M. Eslami, Oriol Vinyals, and Felix Hill (2021). "Multimodal few-shot learning with frozen language models."

Advances in Neural Information Processing Systems 34 (2021): 200-212.



# Multimodal Few-Shot Learning with Frozen Language Models

(a) minImageNet



(b) Fast VQA



Examples of (a) the Open-Ended minImageNet evaluation (b) the Fast VQA evaluation.

Source: Maria Tsimpoukelli, Jacob L. Menick, Serkan Cabi, S. M. Eslami, Oriol Vinyals, and Felix Hill (2021). "Multimodal few-shot learning with frozen language models."

Advances in Neural Information Processing Systems 34 (2021): 200-212.

# GPT-3: Language Models are Few-Shot Learners

---

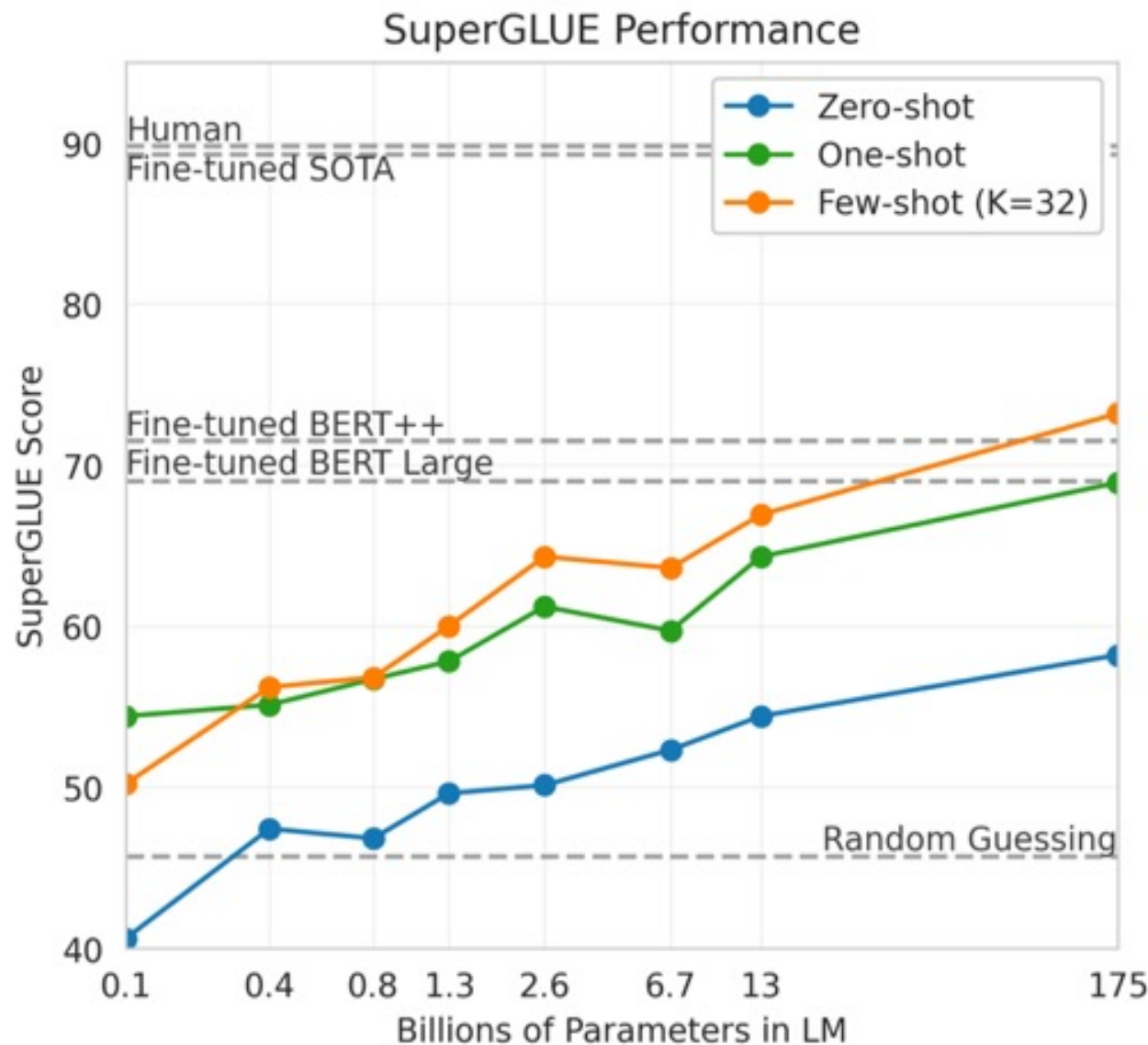
## Language Models are Few-Shot Learners

---

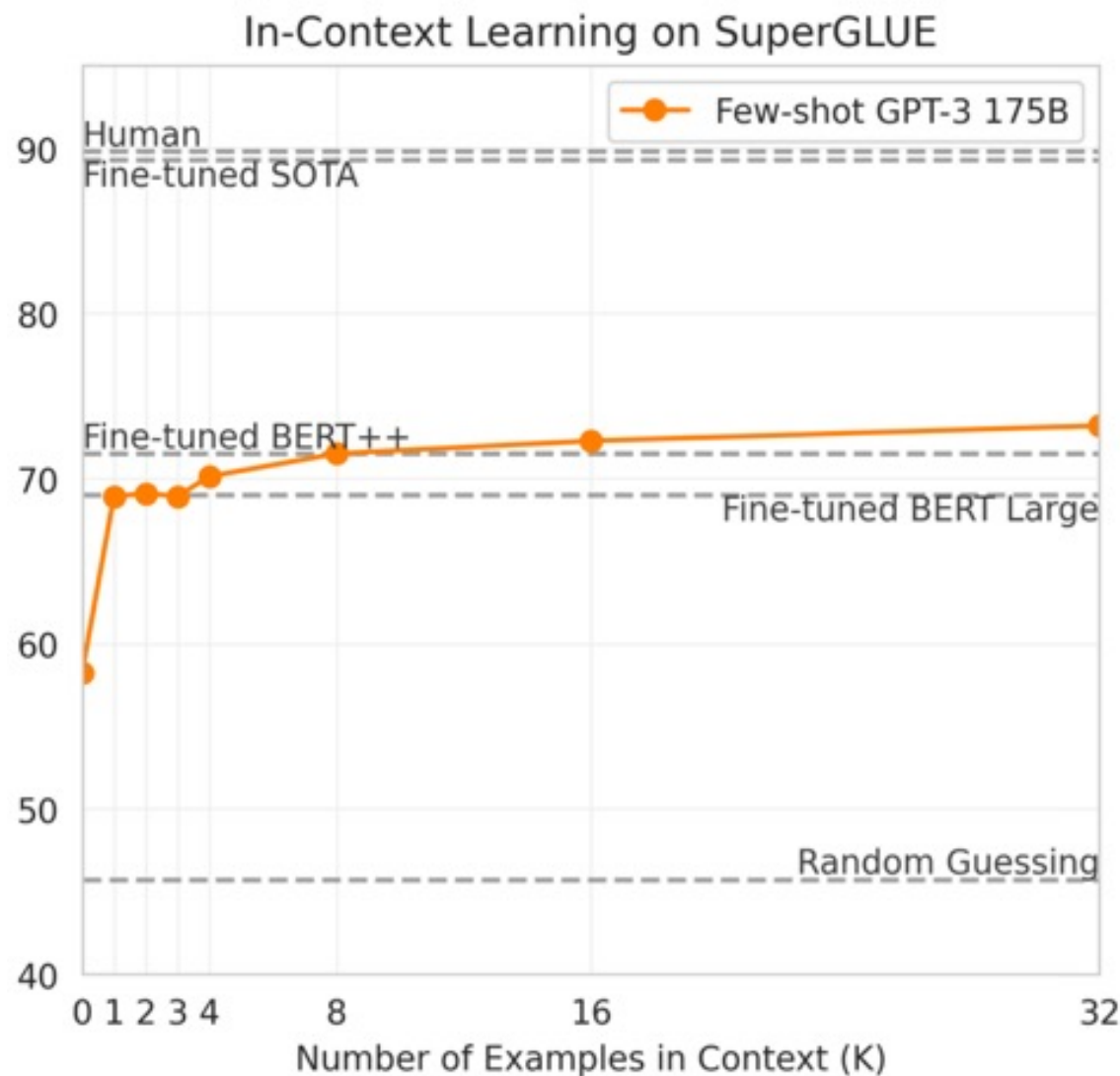
<b>Tom B. Brown*</b>	<b>Benjamin Mann*</b>	<b>Nick Ryder*</b>	<b>Melanie Subbiah*</b>	
<b>Jared Kaplan<sup>†</sup></b>	<b>Prafulla Dhariwal</b>	<b>Arvind Neelakantan</b>	<b>Pranav Shyam</b>	
<b>Girish Sastry</b>	<b>Amanda Askell</b>	<b>Sandhini Agarwal</b>	<b>Ariel Herbert-Voss</b>	
<b>Gretchen Krueger</b>	<b>Tom Henighan</b>	<b>Rewon Child</b>	<b>Aditya Ramesh</b>	
<b>Daniel M. Ziegler</b>	<b>Jeffrey Wu</b>	<b>Clemens Winter</b>		
<b>Christopher Hesse</b>	<b>Mark Chen</b>	<b>Eric Sigler</b>	<b>Mateusz Litwin</b>	<b>Scott Gray</b>
<b>Benjamin Chess</b>	<b>Jack Clark</b>	<b>Christopher Berner</b>		
<b>Sam McCandlish</b>	<b>Alec Radford</b>	<b>Ilya Sutskever</b>	<b>Dario Amodei</b>	

This work was funded by **OpenAI**. All models were trained on **V100** GPU's on part of a high-bandwidth cluster provided by Microsoft.

# GPT-3: Language Models are Few-Shot Learners



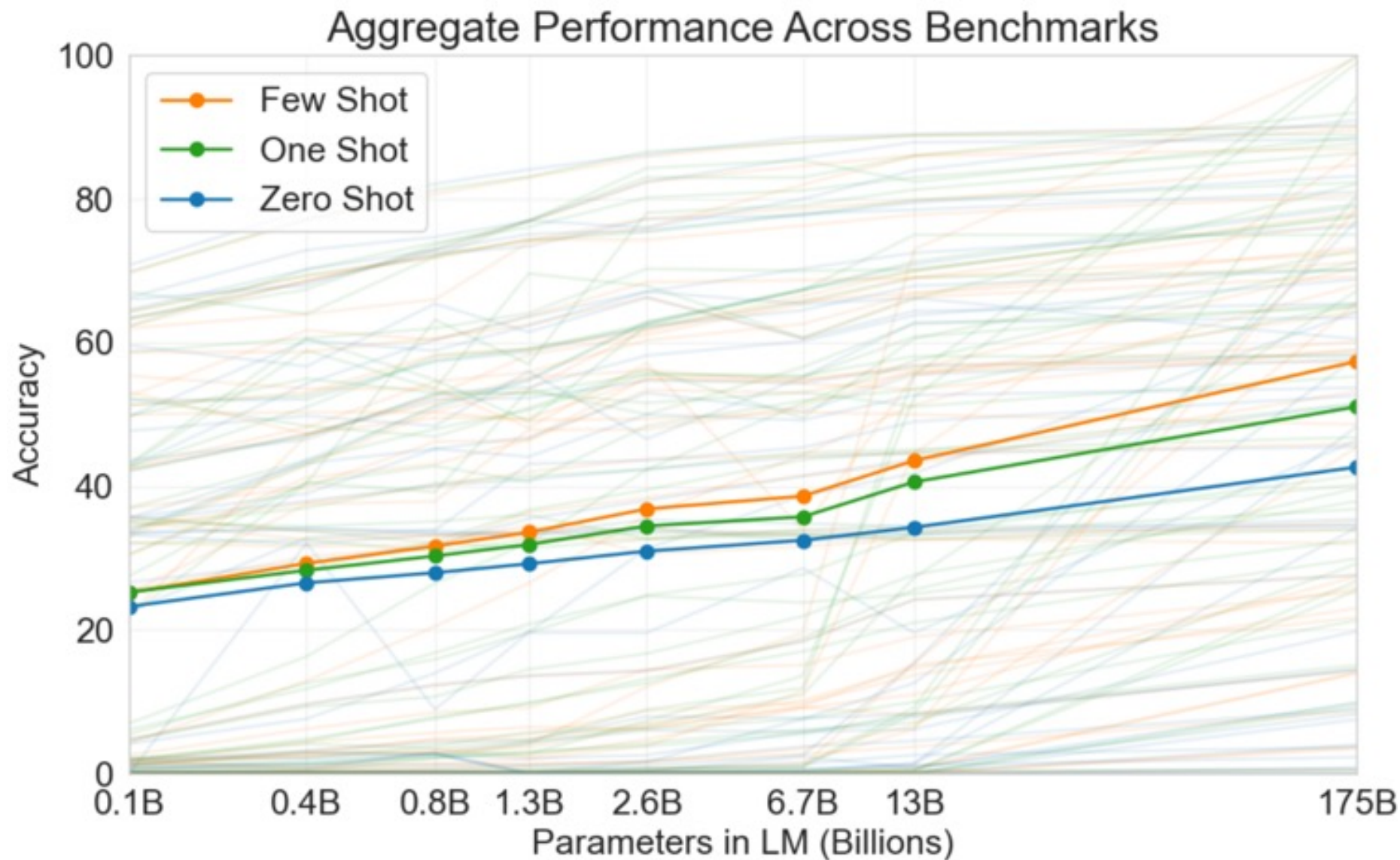
# GPT-3: Language Models are Few-Shot Learners



Performance on SuperGLUE increases with model size. A value of  $K = 32$  means that our model was shown 32 examples per task, for 256 examples total divided across the 8 tasks in SuperGLUE.



# GPT-3: Language Models are Few-Shot Learners



# GPT-3: Language Models are Few-Shot Learners

Performance on cloze and completion tasks.

Setting	LAMBADA (acc)	LAMBADA (ppl)	StoryCloze (acc)	HellaSwag (acc)
SOTA	68.0 <sup>a</sup>	8.63 <sup>b</sup>	<b>91.8<sup>c</sup></b>	<b>85.6<sup>d</sup></b>
GPT-3 Zero-Shot	<b>76.2</b>	<b>3.00</b>	83.2	78.9
GPT-3 One-Shot	<b>72.5</b>	<b>3.35</b>	84.7	78.1
GPT-3 Few-Shot	<b>86.4</b>	<b>1.92</b>	87.7	79.3

GPT-3 significantly improves SOTA on LAMBADA while achieving respectable performance on two difficult completion prediction datasets.

# GPT-3: Language Models are Few-Shot Learners

## Results on three Open-Domain QA tasks

Setting	NaturalQS	WebQS	TriviaQA
RAG (Fine-tuned, Open-Domain) [LPP <sup>+</sup> 20]	<b>44.5</b>	<b>45.5</b>	<b>68.0</b>
T5-11B+SSM (Fine-tuned, Closed-Book) [RRS20]	36.6	44.7	60.5
T5-11B (Fine-tuned, Closed-Book)	34.5	37.4	50.1
GPT-3 Zero-Shot	14.6	14.4	64.3
GPT-3 One-Shot	23.0	25.3	<b>68.0</b>
GPT-3 Few-Shot	29.9	41.5	<b>71.2</b>

GPT-3 is shown in the few-, one-, and zero-shot settings, as compared to prior SOTA results for closed book and open domain settings.

TriviaQA few-shot result is evaluated on the wiki split test server.

# GPT-3: Language Models are Few-Shot Learners

GPT-3 results on a selection of QA / RC tasks.

Setting	ARC (Easy)	ARC (Challenge)	CoQA	DROP
Fine-tuned SOTA	<b>92.0<sup>a</sup></b>	<b>78.5<sup>b</sup></b>	<b>90.7<sup>c</sup></b>	<b>89.1<sup>d</sup></b>
GPT-3 Zero-Shot	68.8	51.4	81.5	23.6
GPT-3 One-Shot	71.2	53.2	84.0	34.3
GPT-3 Few-Shot	70.1	51.5	85.0	36.5

CoQA and DROP are F1 while ARC reports accuracy.

See the appendix for additional experiments. a[KKS+20] b[KKS+20] c[JZC+19] d [JN20]



# GPT-3: Language Models are Few-Shot Learners

Setting	En→Fr	Fr→En	En→De	De→En	En→Ro	Ro→En
SOTA (Supervised)	<b>45.6<sup>a</sup></b>	35.0 <sup>b</sup>	<b>41.2<sup>c</sup></b>	40.2 <sup>d</sup>	<b>38.5<sup>e</sup></b>	<b>39.9<sup>e</sup></b>
XLM [LC19]	33.4	33.3	26.4	34.3	33.3	31.8
MASS [STQ <sup>+</sup> 19]	<u>37.5</u>	34.9	28.3	35.2	<u>35.2</u>	33.1
mBART [LGG <sup>+</sup> 20]	-	-	<u>29.8</u>	34.0	<u>35.0</u>	30.5
GPT-3 Zero-Shot	25.2	21.2	24.6	27.2	14.1	19.9
GPT-3 One-Shot	28.3	33.7	26.2	30.4	20.6	38.6
GPT-3 Few-Shot	32.6	<u>39.2</u>	29.7	<u>40.6</u>	21.0	<u>39.5</u>

**Few-shot GPT-3 outperforms previous unsupervised NMT work by 5 BLEU when translating into English reflecting its strength as an English LM.**

# GPT-3: Language Models are Few-Shot Learners

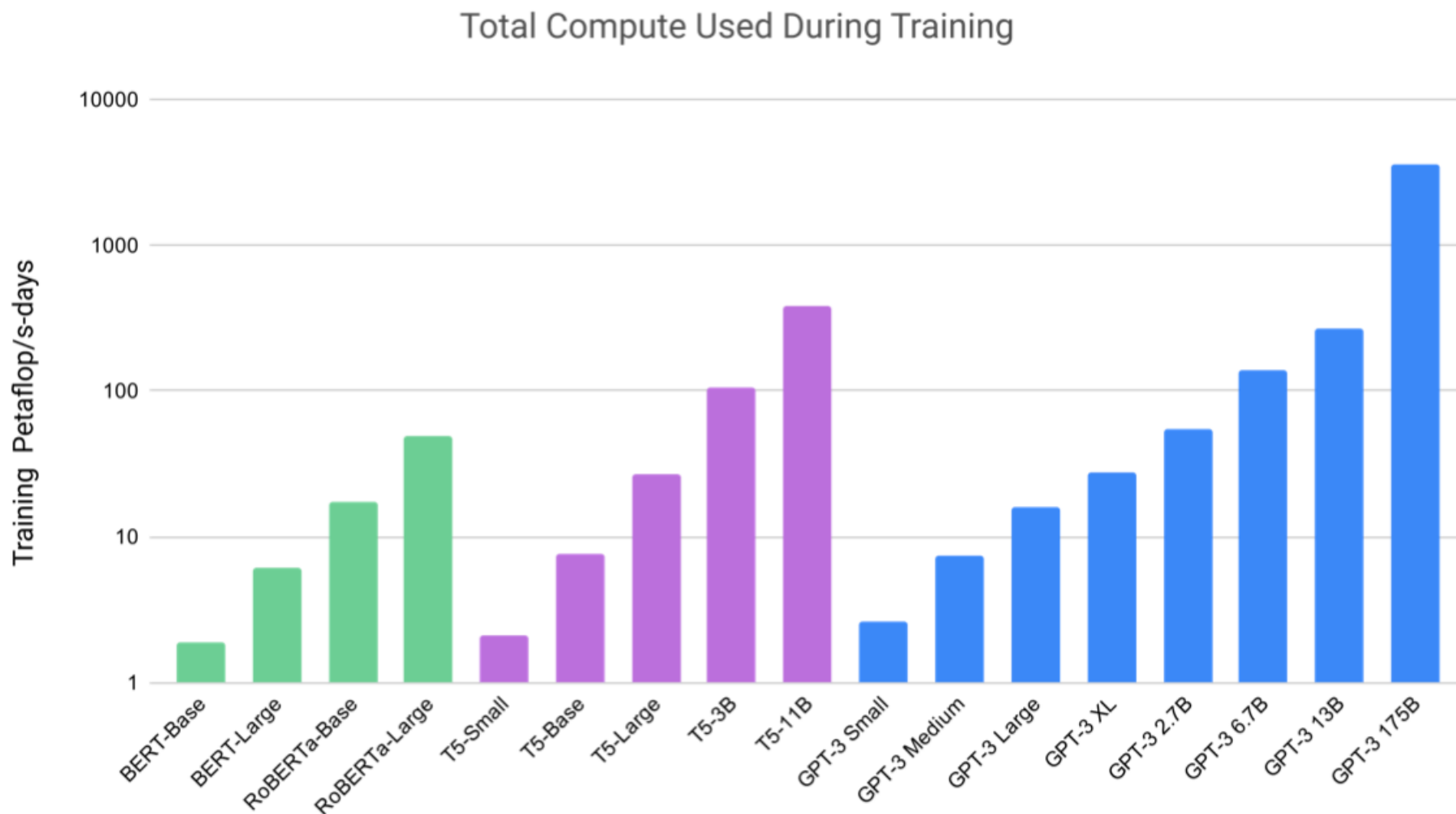
## Performance of GPT-3 on SuperGLUE compared to fine-tuned baselines and SOTA

	SuperGLUE Average	BoolQ Accuracy	CB Accuracy	CB F1	COPA Accuracy	RTE Accuracy
Fine-tuned SOTA	<b>89.0</b>	<b>91.0</b>	<b>96.9</b>	<b>93.9</b>	<b>94.8</b>	<b>92.5</b>
Fine-tuned BERT-Large	69.0	77.4	83.6	75.7	70.6	71.7
GPT-3 Few-Shot	71.8	76.4	75.6	52.0	92.0	69.0

	WiC Accuracy	WSC Accuracy	MultiRC Accuracy	MultiRC F1a	ReCoRD Accuracy	ReCoRD F1
Fine-tuned SOTA	<b>76.1</b>	<b>93.8</b>	<b>62.3</b>	<b>88.2</b>	<b>92.5</b>	<b>93.3</b>
Fine-tuned BERT-Large	69.6	64.6	24.1	70.0	71.3	72.0
GPT-3 Few-Shot	49.4	80.1	30.5	75.4	90.2	91.1

**GPT-3 few-shot is given a total of 32 examples within the context of each task and performs no gradient updates.**

# GPT-3: Language Models are Few-Shot Learners



GPT-3 3B is almost 10x larger than RoBERTa-Large (355M params),  
both models took roughly 50 petaflop/s-days of compute during pre-training

# GPT-3: Language Models are Few-Shot Learners

Human accuracy in identifying  
whether short (~200 word) news articles are model generated

	Mean accuracy	95% Confidence Interval (low, hi)	$t$ compared to control ( $p$ -value)	“I don’t know” assignments
Control (deliberately bad model)	86%	83%–90%	-	3.6 %
GPT-3 Small	76%	72%–80%	3.9 ( $2e-4$ )	4.9%
GPT-3 Medium	61%	58%–65%	10.3 ( $7e-21$ )	6.0%
GPT-3 Large	68%	64%–72%	7.3 ( $3e-11$ )	8.7%
GPT-3 XL	62%	59%–65%	10.7 ( $1e-19$ )	7.5%
GPT-3 2.7B	62%	58%–65%	10.4 ( $5e-19$ )	7.1%
GPT-3 6.7B	60%	56%–63%	11.2 ( $3e-21$ )	6.2%
GPT-3 13B	55%	52%–58%	15.3 ( $1e-32$ )	7.1%
GPT-3 175B	52%	49%–54%	16.9 ( $1e-34$ )	7.8%

This table compares mean accuracy between five different models, and shows the results of a two-sample T-Test for the difference in mean accuracy between each model and the control model (an unconditional GPT-3 Small model with increased output randomness).

# GPT-3: Language Models are Few-Shot Learners

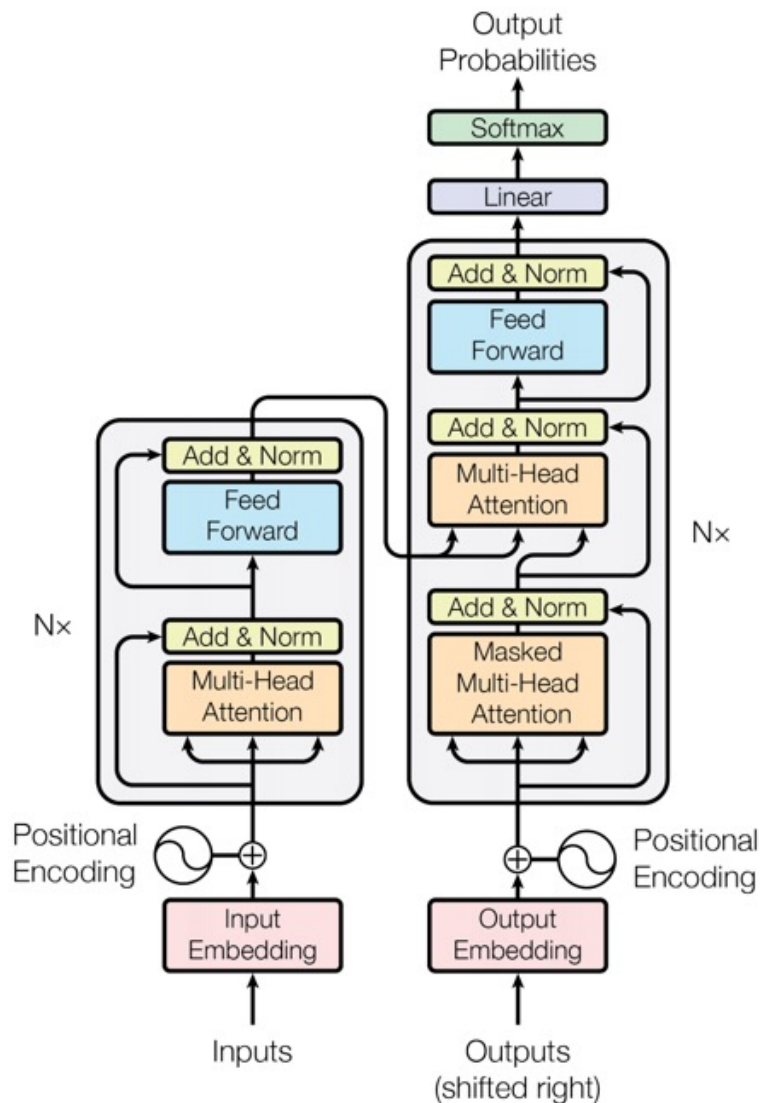
The GPT-3 generated news article that humans had the greatest difficulty distinguishing from a human written article (accuracy: 12%)

Title: United Methodists Agree to Historic Split  
Subtitle: Those who oppose gay marriage will form their own denomination  
Article: After two days of intense debate, the United Methodist Church has agreed to a historic split - one that is expected to end in the creation of a new denomination, one that will be "theologically and socially conservative," according to The Washington Post. The majority of delegates attending the church's annual General Conference in May voted to strengthen a ban on the ordination of LGBTQ clergy and to write new rules that will "discipline" clergy who officiate at same-sex weddings. But those who opposed these measures have a new plan: They say they will form a separate denomination by 2020, calling their church the Christian Methodist denomination. The Post notes that the denomination, which claims 12.5 million members, was in the early 20th century the "largest Protestant denomination in the U.S.," but that it has been shrinking in recent decades. The new split will be the second in the church's history. The first occurred in 1968, when roughly 10 percent of the denomination left to form the Evangelical United Brethren Church. The Post notes that the proposed split "comes at a critical time for the church, which has been losing members for years," which has been "pushed toward the brink of a schism over the role of LGBTQ people in the church." Gay marriage is not the only issue that has divided the church. In 2016, the denomination was split over ordination of transgender clergy, with the North Pacific regional conference voting to ban them from serving as clergy, and the South Pacific regional conference voting to allow them.



# Transformer (Attention is All You Need)

(Vaswani et al., 2017)

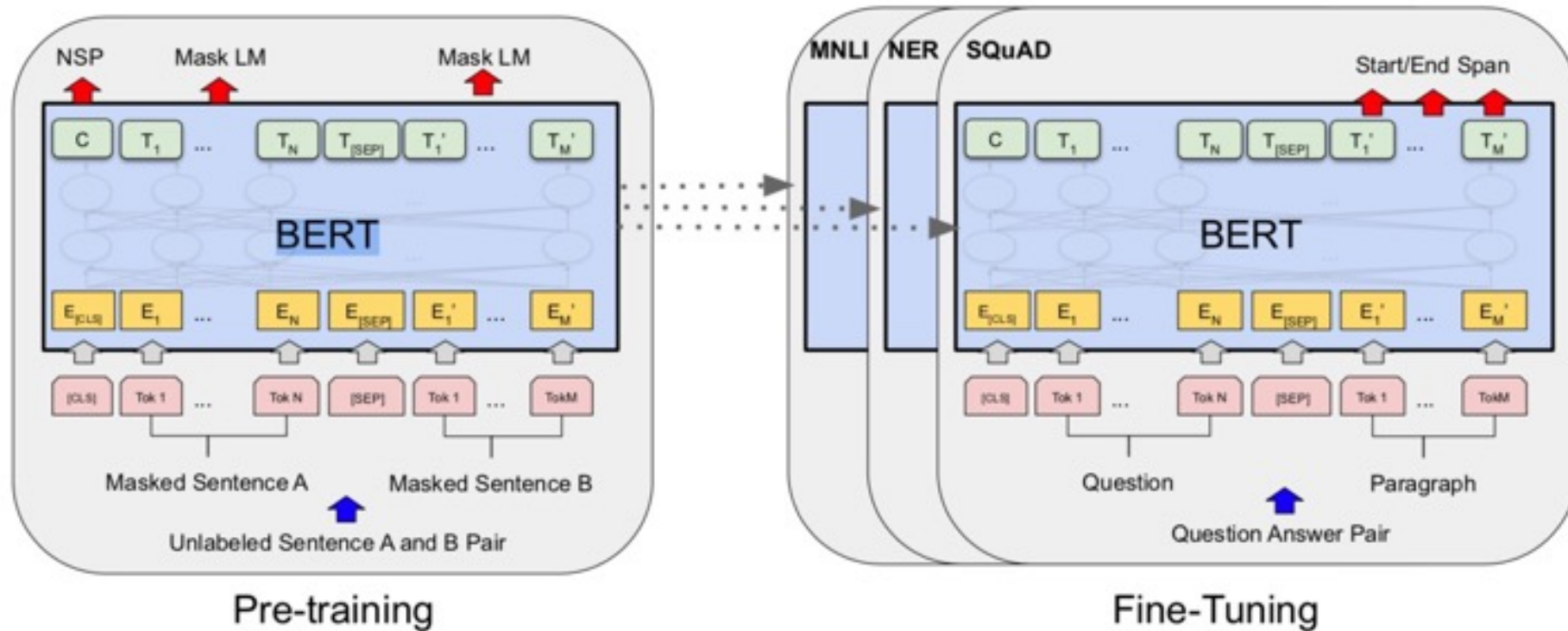


Source: Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "Attention is all you need." In *Advances in neural information processing systems*, pp. 5998-6008. 2017.

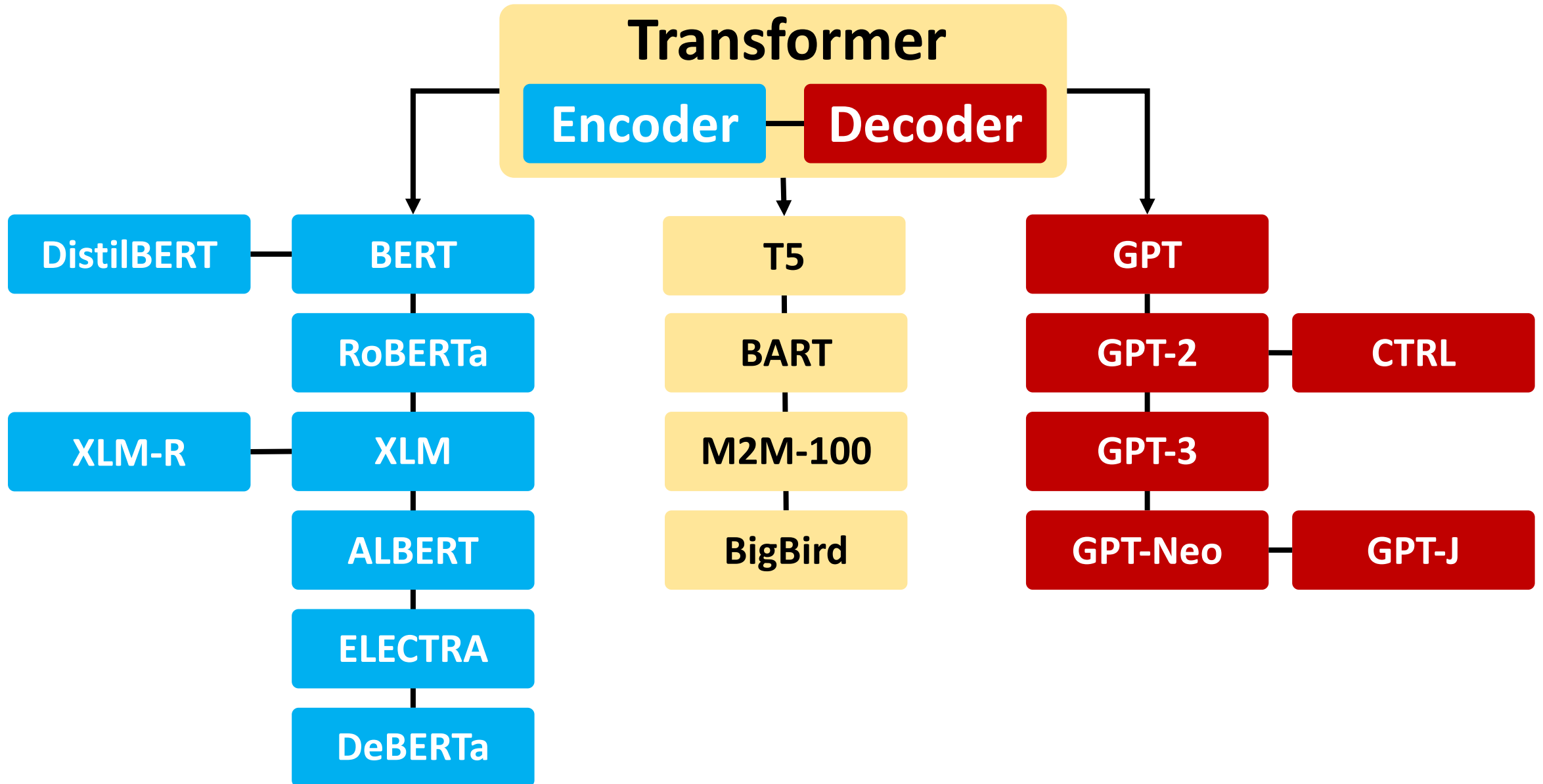
# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

BERT (Bidirectional Encoder Representations from Transformers)

Overall pre-training and fine-tuning procedures for BERT

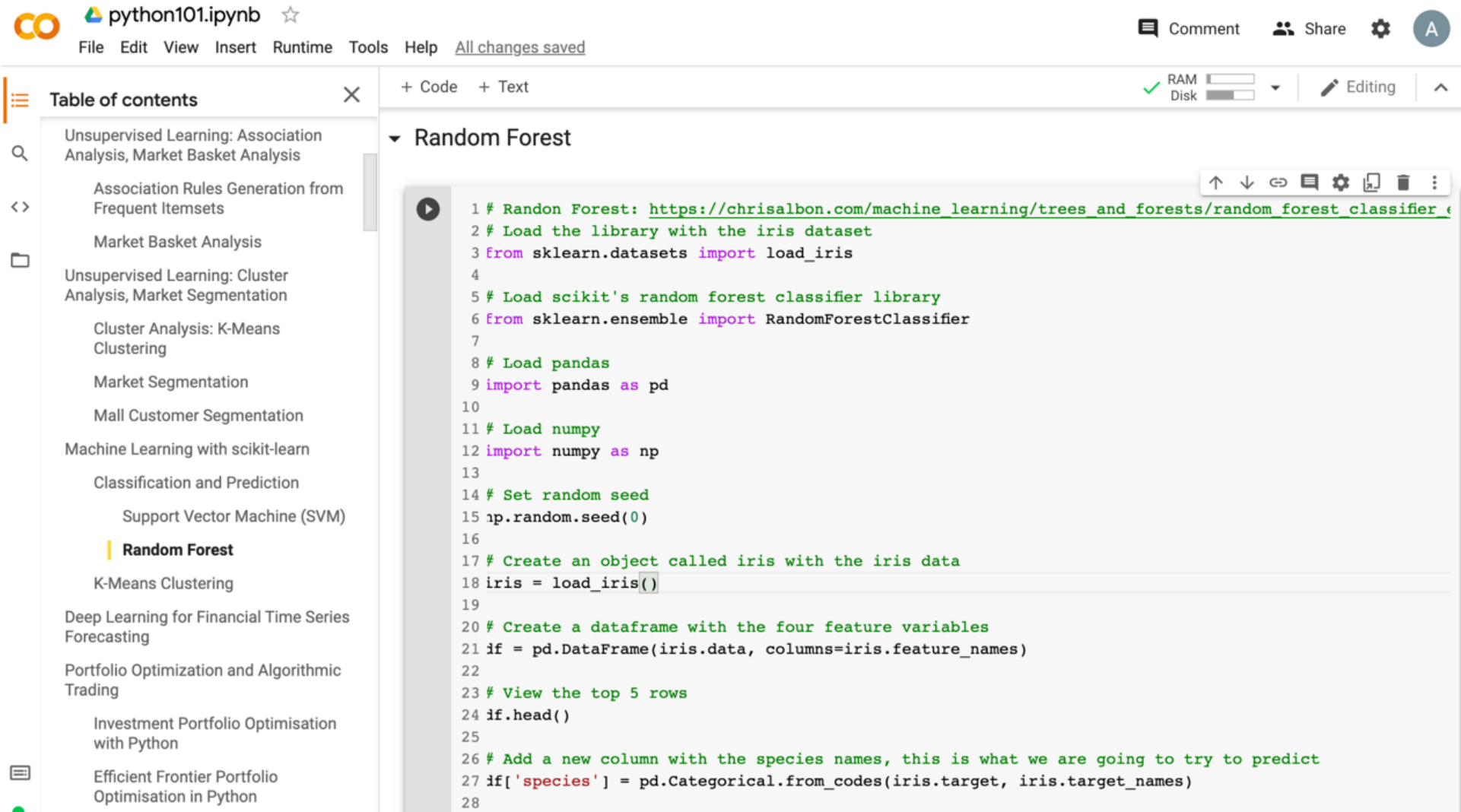


# Transformer Models





# Machine Learning: Ensemble Learning Random Forest



The screenshot displays a Jupyter Notebook interface. On the left, a 'Table of contents' sidebar lists various machine learning topics, with 'Random Forest' highlighted. The main area shows a code cell titled 'Random Forest' containing Python code for loading the Iris dataset and training a Random Forest classifier. The code includes comments in green and Python syntax in black. The interface also shows a top bar with navigation options and a right bar with RAM and disk usage indicators.

python101.ipynb ☆

File Edit View Insert Runtime Tools Help All changes saved

Comment Share ⚙️ A

RAM Disk

Editing

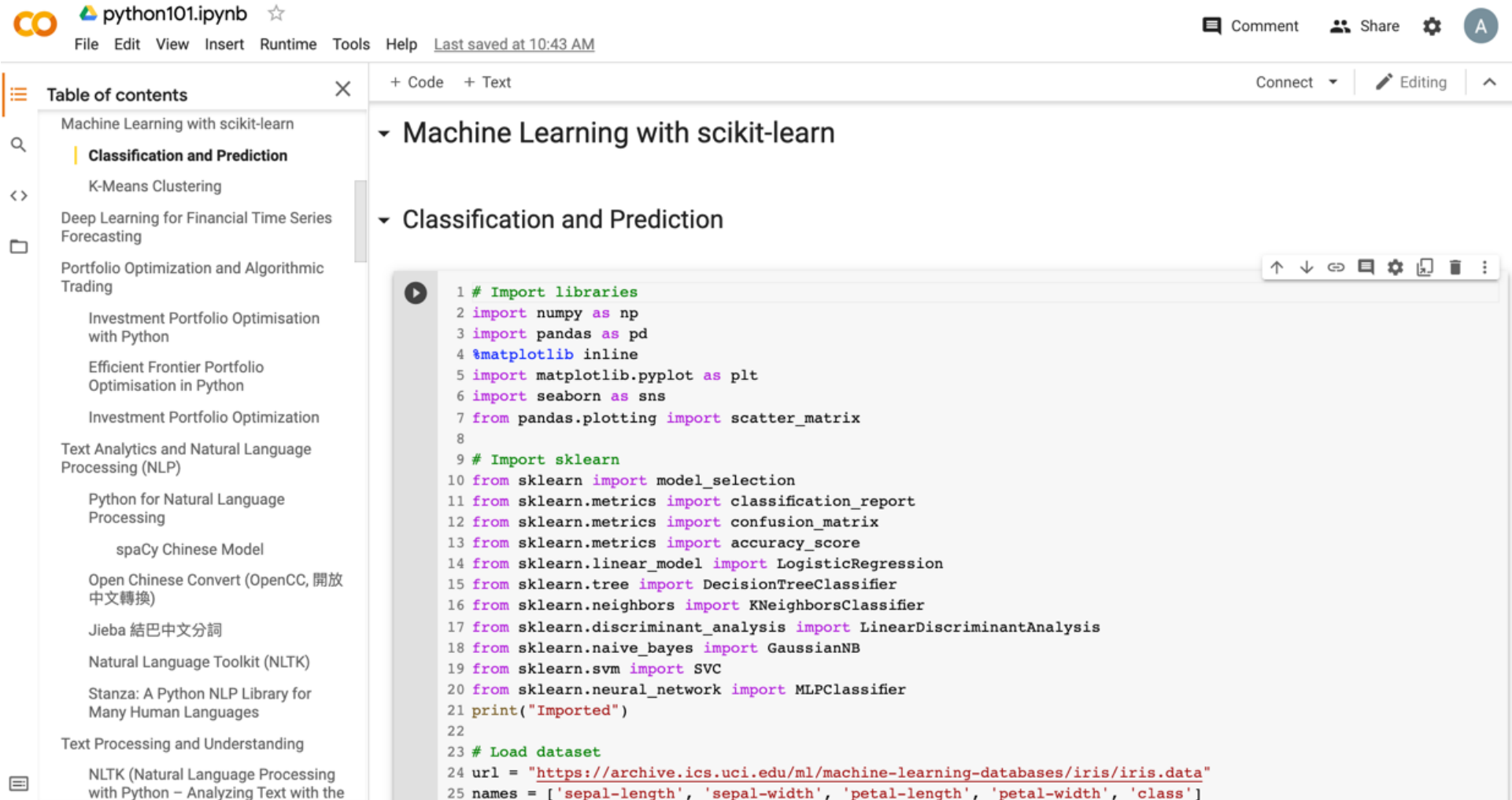
Table of contents

- Unsupervised Learning: Association Analysis, Market Basket Analysis
  - Association Rules Generation from Frequent Itemsets
  - Market Basket Analysis
- Unsupervised Learning: Cluster Analysis, Market Segmentation
  - Cluster Analysis: K-Means Clustering
  - Market Segmentation
  - Mall Customer Segmentation
- Machine Learning with scikit-learn
  - Classification and Prediction
    - Support Vector Machine (SVM)
    - Random Forest**
    - K-Means Clustering
- Deep Learning for Financial Time Series Forecasting
- Portfolio Optimization and Algorithmic Trading
  - Investment Portfolio Optimisation with Python
  - Efficient Frontier Portfolio Optimisation in Python

Random Forest

```
1 # Random Forest: https://chrisalbon.com/machine\_learning/trees\_and\_forests/random\_forest\_classifier/
2 # Load the library with the iris dataset
3 from sklearn.datasets import load_iris
4
5 # Load scikit's random forest classifier library
6 from sklearn.ensemble import RandomForestClassifier
7
8 # Load pandas
9 import pandas as pd
10
11 # Load numpy
12 import numpy as np
13
14 # Set random seed
15 np.random.seed(0)
16
17 # Create an object called iris with the iris data
18 iris = load_iris()
19
20 # Create a dataframe with the four feature variables
21 if = pd.DataFrame(iris.data, columns=iris.feature_names)
22
23 # View the top 5 rows
24 if.head()
25
26 # Add a new column with the species names, this is what we are going to try to predict
27 if['species'] = pd.Categorical.from_codes(iris.target, iris.target_names)
28
```

# Machine Learning: Supervised Learning Classification and Prediction



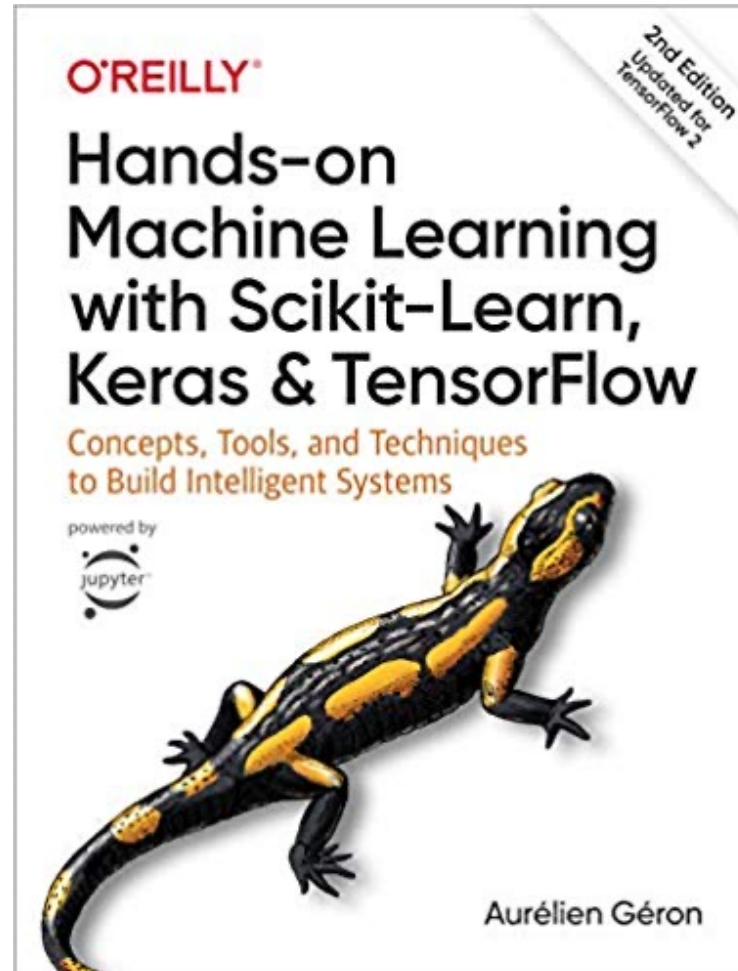
The screenshot displays a Jupyter Notebook interface. At the top, the notebook is titled "python101.ipynb" and shows it was last saved at 10:43 AM. The interface includes a "Table of contents" sidebar on the left, a main content area, and a top navigation bar with options like "File", "Edit", "View", "Insert", "Runtime", "Tools", "Help", "Comment", "Share", and "Settings".

The "Table of contents" sidebar lists various topics, with "Classification and Prediction" highlighted. The main content area shows a table of contents for the current notebook, listing topics such as "Machine Learning with scikit-learn", "K-Means Clustering", "Deep Learning for Financial Time Series Forecasting", "Portfolio Optimization and Algorithmic Trading", "Text Analytics and Natural Language Processing (NLP)", and "Text Processing and Understanding".

The code cell in the main content area contains the following Python code:

```
1 # Import libraries
2 import numpy as np
3 import pandas as pd
4 %matplotlib inline
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 from pandas.plotting import scatter_matrix
8
9 # Import sklearn
10 from sklearn import model_selection
11 from sklearn.metrics import classification_report
12 from sklearn.metrics import confusion_matrix
13 from sklearn.metrics import accuracy_score
14 from sklearn.linear_model import LogisticRegression
15 from sklearn.tree import DecisionTreeClassifier
16 from sklearn.neighbors import KNeighborsClassifier
17 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
18 from sklearn.naive_bayes import GaussianNB
19 from sklearn.svm import SVC
20 from sklearn.neural_network import MLPClassifier
21 print("Imported")
22
23 # Load dataset
24 url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
25 names = ['sepal-length', 'sepal-width', 'petal-length', 'petal-width', 'class']
```

**Aurélien Géron (2019),**  
**Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow:**  
**Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition**  
**O'Reilly Media, 2019**



<https://github.com/ageron/handson-ml2>

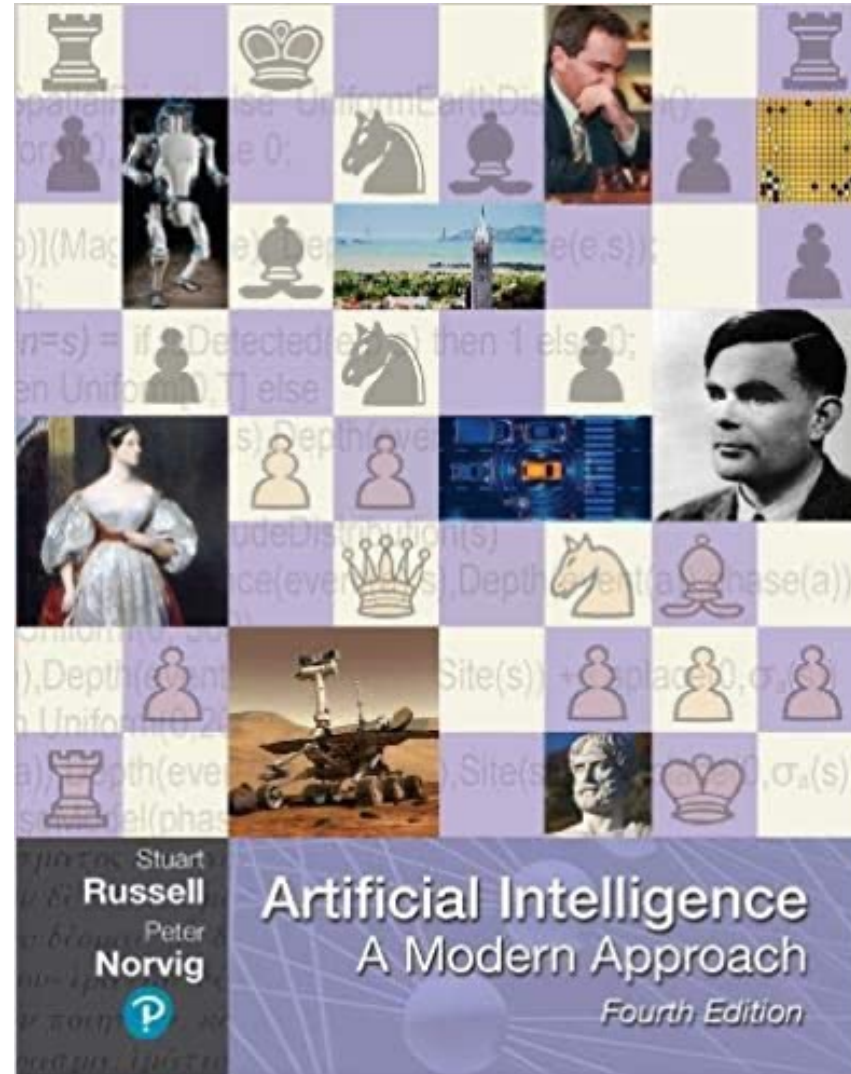
# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow

## Notebooks

- [1.The Machine Learning landscape](#)
- [2.End-to-end Machine Learning project](#)
- [3.Classification](#)
- [4.Training Models](#)
- [5.Support Vector Machines](#)
- [6.Decision Trees](#)
- [7.Ensemble Learning and Random Forests](#)
- [8.Dimensionality Reduction](#)
- [9.Unsupervised Learning Techniques](#)
- [10.Artificial Neural Nets with Keras](#)
- [11.Training Deep Neural Networks](#)
- [12.Custom Models and Training with TensorFlow](#)
- [13.Loading and Preprocessing Data](#)
- [14.Deep Computer Vision Using Convolutional Neural Networks](#)
- [15.Processing Sequences Using RNNs and CNNs](#)
- [16.Natural Language Processing with RNNs and Attention](#)
- [17.Representation Learning Using Autoencoders](#)
- [18.Reinforcement Learning](#)
- [19.Training and Deploying TensorFlow Models at Scale](#)



Stuart Russell and Peter Norvig (2020),  
**Artificial Intelligence: A Modern Approach,**  
4th Edition, Pearson



Source: Stuart Russell and Peter Norvig (2020), Artificial Intelligence: A Modern Approach, 4th Edition, Pearson

<https://www.amazon.com/Artificial-Intelligence-A-Modern-Approach/dp/0134610997/>

# Artificial Intelligence: A Modern Approach (AIMA)

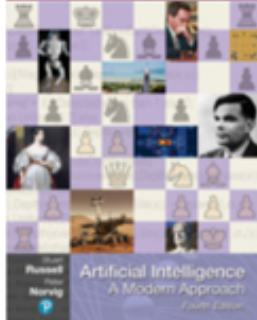
- **Artificial Intelligence: A Modern Approach (AIMA)**
  - <http://aima.cs.berkeley.edu/>
- **AIMA Python**
  - <http://aima.cs.berkeley.edu/python/readme.html>
  - <https://github.com/aimacode/aima-python>
- **Learning**
  - <http://aima.cs.berkeley.edu/python/learning.html>



# Artificial Intelligence: A Modern Approach (AIMA)



△ US Edition
△ Global Edition
Acknowledgements
Code
Courses
Editions
Errata
Exercises
Figures
Instructors Page
Pseudocode
Reviews



## Artificial Intelligence: A Modern Approach, 4th US ed.

by Stuart Russell and Peter Norvig

The authoritative, most-used AI textbook, adopted by over 1500 schools.

Table of Contents for the US Edition (or see the [Global Edition](#))

[Preface \(pdf\)](#); [Contents with subsections](#)

### I Artificial Intelligence

- 1 Introduction ... 1
- 2 Intelligent Agents ... 36

### II Problem-solving

- 3 Solving Problems by Searching ... 63
- 4 Search in Complex Environments ... 110
- 5 Adversarial Search and Games ... 146
- 6 Constraint Satisfaction Problems ... 180

### III Knowledge, reasoning, and planning

- 7 Logical Agents ... 208
- 8 First-Order Logic ... 251
- 9 Inference in First-Order Logic ... 280
- 10 Knowledge Representation ... 314
- 11 Automated Planning ... 344

### IV Uncertain knowledge and reasoning

- 12 Quantifying Uncertainty ... 385
- 13 Probabilistic Reasoning ... 412
- 14 Probabilistic Reasoning over Time ... 461
- 15 Probabilistic Programming ... 500
- 16 Making Simple Decisions ... 528
- 17 Making Complex Decisions ... 562
- 18 Multiagent Decision Making ... 599

### V Machine Learning

- 19 Learning from Examples ... 651
- 20 Learning Probabilistic Models ... 721
- 21 Deep Learning ... 750
- 22 Reinforcement Learning ... 789

### VI Communicating, perceiving, and acting

- 23 Natural Language Processing ... 823
- 24 Deep Learning for Natural Language Processing ... 856
- 25 Computer Vision ... 881
- 26 Robotics ... 925

### VII Conclusions

- 27 Philosophy, Ethics, and Safety of AI ... 981
- 28 The Future of AI ... 1012
- Appendix A: Mathematical Background ... 1023
- Appendix B: Notes on Languages and Algorithms ... 1030
- Bibliography ... 1033 ([pdf](#) and [LaTeX .bib file](#) and [bib data](#))
- Index ... 1069 ([pdf](#))

[Exercises \(website\)](#)

[Figures \(pdf\)](#)

[Code \(website\)](#); [Pseudocode \(pdf\)](#)

Covers: [US](#), [Global](#)

# Papers with Code State-of-the-Art (SOTA)



Search for papers, code and tasks



Browse State-of-the-Art

Follow

Discuss

Trends

About

Log In/Register

## Browse State-of-the-Art

1509 leaderboards • 1327 tasks • 1347 datasets • 17810 papers with code

Follow on Twitter for updates

## Computer Vision



**Semantic Segmentation**

33 leaderboards

667 papers with code



**Image Classification**

52 leaderboards

564 papers with code



**Object Detection**

54 leaderboards

467 papers with code



**Image Generation**

51 leaderboards

231 papers with code



**Pose Estimation**

40 leaderboards

231 papers with code

[▶ See all 707 tasks](#)

## Natural Language Processing



**Machine Translation**



**Language Modelling**



**Question Answering**



**Sentiment Analysis**



**Text Generation**

<https://paperswithcode.com/sota>



# Summary

- **The Theory of Learning**
  - **Computational Learning Theory**
  - **Probably Approximately Correct (PAC) Learning**
- **Ensemble Learning**
  - **Bagging: Random forests**
  - **Stacking**
  - **Boosting: Gradient boosting**
  - **Online learning**
- **Meta Learning: Learning to Learn**

# References

- Stuart Russell and Peter Norvig (2020), Artificial Intelligence: A Modern Approach, 4th Edition, Pearson.
- Aurélien Géron (2019), Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition, O'Reilly Media.
- Steven D'Ascoli (2022), Artificial Intelligence and Deep Learning with Python: Every Line of Code Explained For Readers New to AI and New to Python, Independently published.
- Nithin Buduma, Nikhil Buduma, Joe Papa (2022), Fundamentals of Deep Learning: Designing Next-Generation Machine Intelligence Algorithms, 2nd Edition, O'Reilly Media.
- Min-Yuh Day (2022), Python 101, <https://tinyurl.com/aintpupython101>