

Artificial Intelligence

The Theory of Learning and Ensemble Learning

1131AI05

MBA, IM, NTPU (M5276) (Fall 2024)

Tue 2, 3, 4 (9:10-12:00) (B3F17)

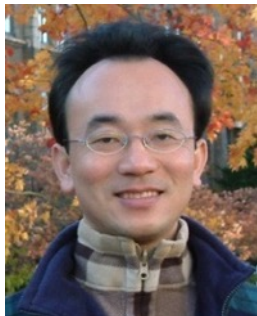
Min-Yuh Day, Ph.D,
Professor

Institute of Information Management, National Taipei University

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



<https://meet.google.com/paj-zhhi-mya>



Syllabus

Week Date Subject/Topics

1 2024/09/10 Introduction to Artificial Intelligence

2 2024/09/17 Mid-Autumn Festival (Day off)

3 2024/09/24 Artificial Intelligence and Intelligent Agents; Problem Solving

**4 2024/10/01 Knowledge, Reasoning and Knowledge Representation;
Uncertain Knowledge and Reasoning**

5 2024/10/08 Case Study on Artificial Intelligence I

6 2024/10/15 Machine Learning: Supervised and Unsupervised Learning

Syllabus

Week Date Subject/Topics

7 2024/10/22 The Theory of Learning and Ensemble Learning

8 2024/10/29 Midterm Project Report

9 2024/11/05 Self-Learning

10 2024/11/12 Deep Learning, Reinforcement Learning

11 2024/11/19 Case Study on Artificial Intelligence II

12 2024/11/26 Deep Learning for Natural Language Processing

Syllabus

Week Date Subject/Topics

13 2024/12/03 Computer Vision and Robotics

**14 2024/12/10 Generative AI,
Philosophy and Ethics of AI and the Future of AI**

15 2024/12/17 Final Project Report I

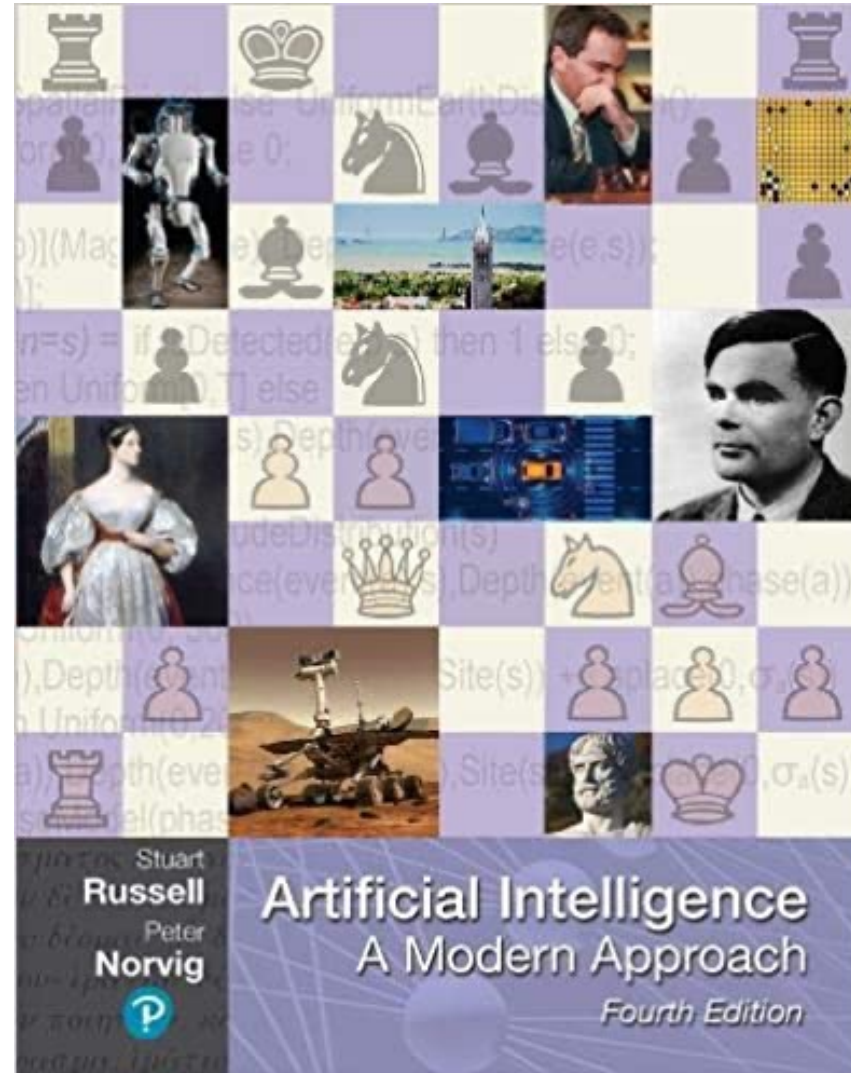
16 2024/12/24 Final Project Report II

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 (RF)**
 - **Boosting: Gradient Boosting, XGBoost, LightGBM, CatBoost**
 - **Stacking**
 - **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

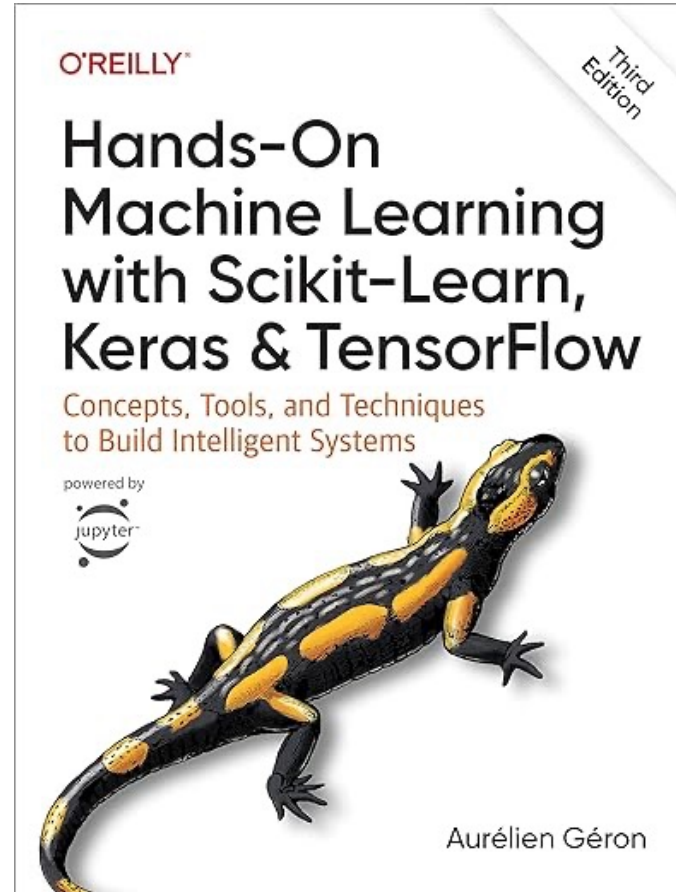
Artificial Intelligence:

5. Machine Learning

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

Aurélien Géron (2022),

Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 3rd Edition, O'Reilly Media



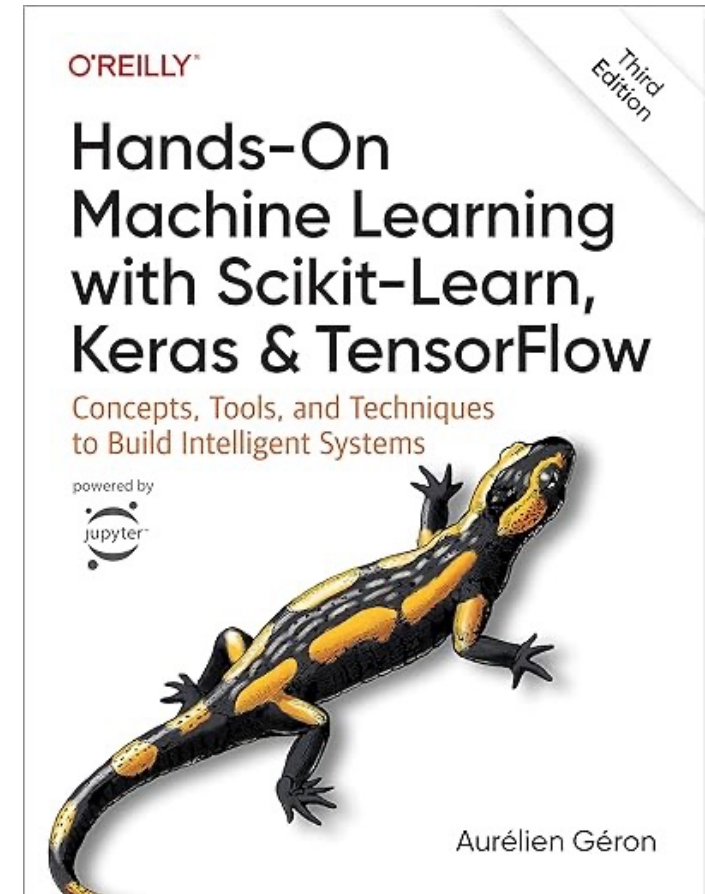
<https://github.com/ageron/handson-ml3>

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. Autoencoders, GANs, and Diffusion Models](#)
- [18. Reinforcement Learning](#)
- [19. Training and Deploying TensorFlow Models at Scale](#)

<https://github.com/ageron/handson-ml3>

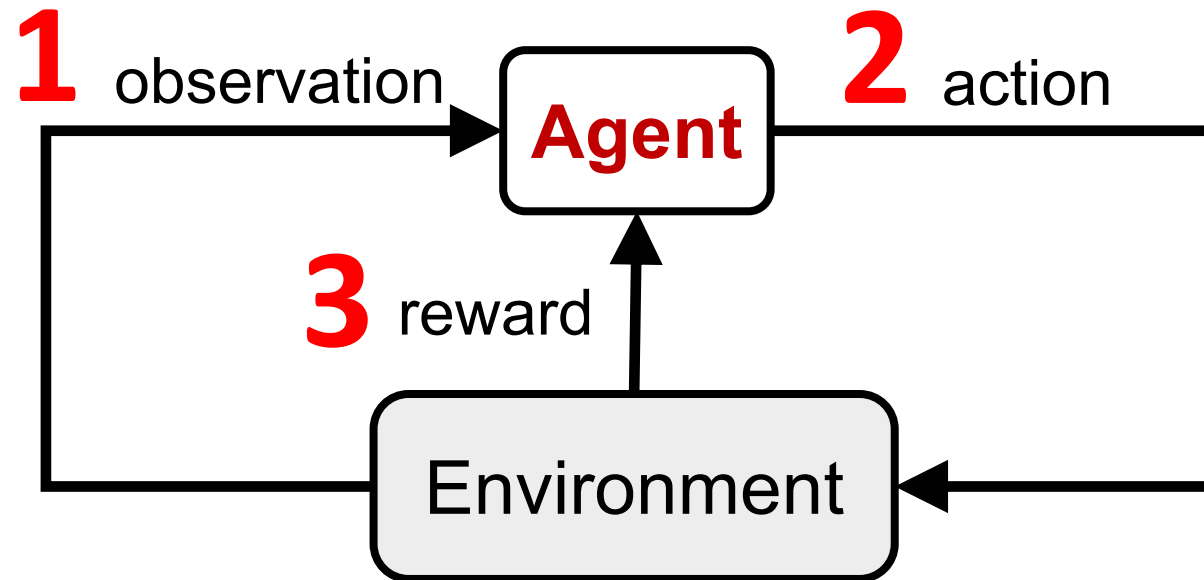


Reinforcement Learning (DL)

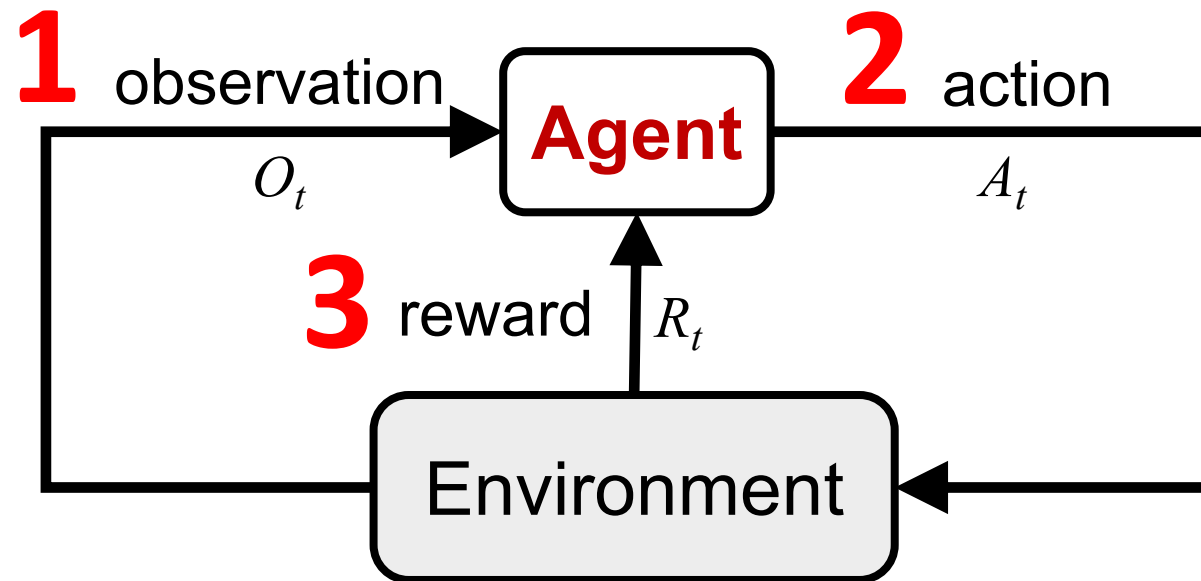
Agent

Environment

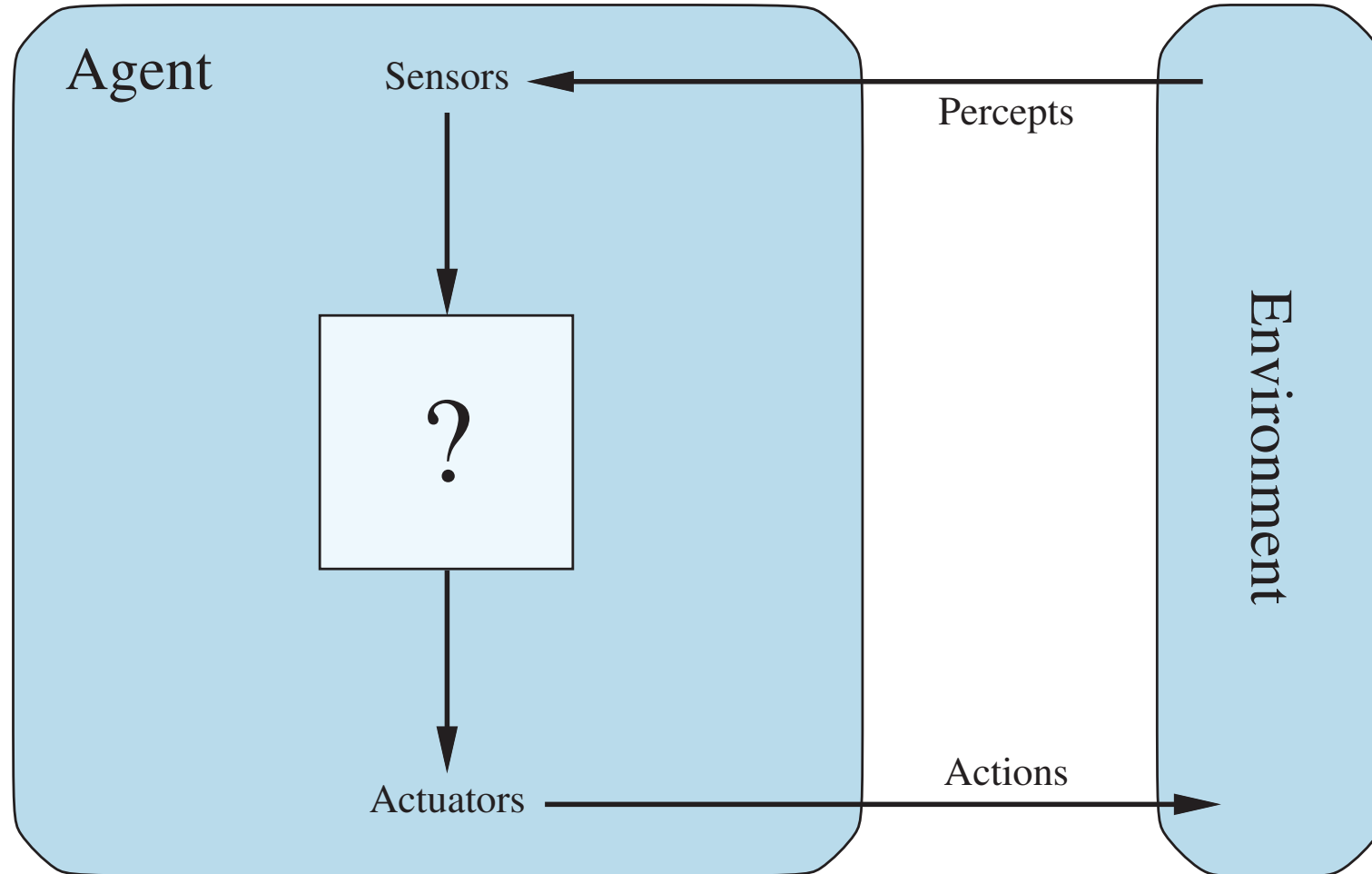
Reinforcement Learning (DL)



Reinforcement Learning (DL)



Agents interact with environments through sensors and actuators

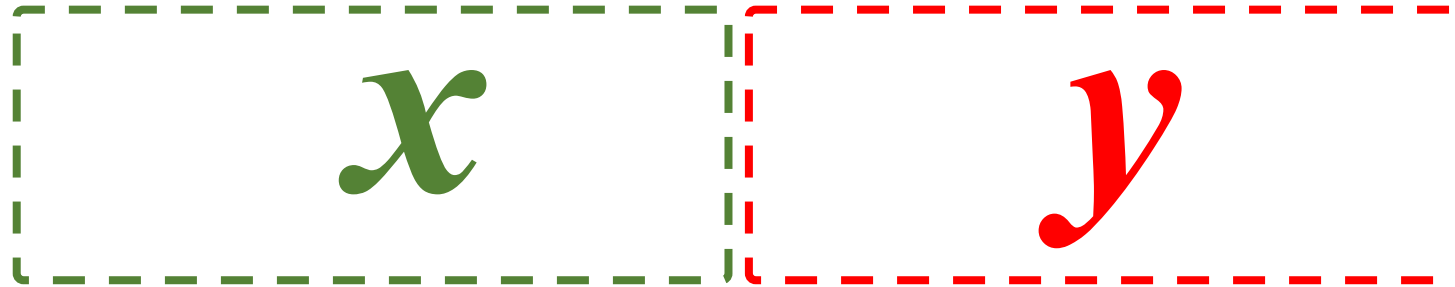


Machine Learning

Supervised Learning (Classification)

Learning from Examples

$$y = f(x)$$

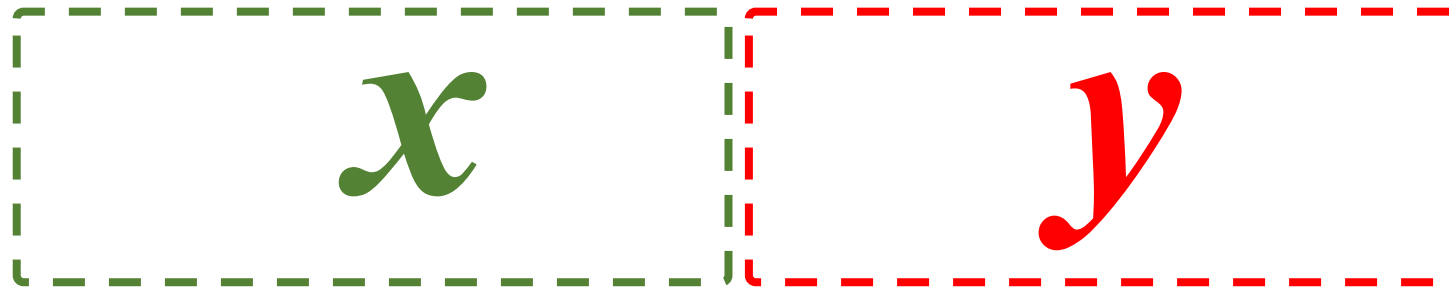


Machine Learning

Supervised Learning (Classification)

Learning from Examples

$$y = f(x)$$



input

Output
label

Machine Learning

Supervised Learning (Classification)

Learning from Examples

$$y = f(x)$$

```
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

$$y = f(x)$$

Example

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

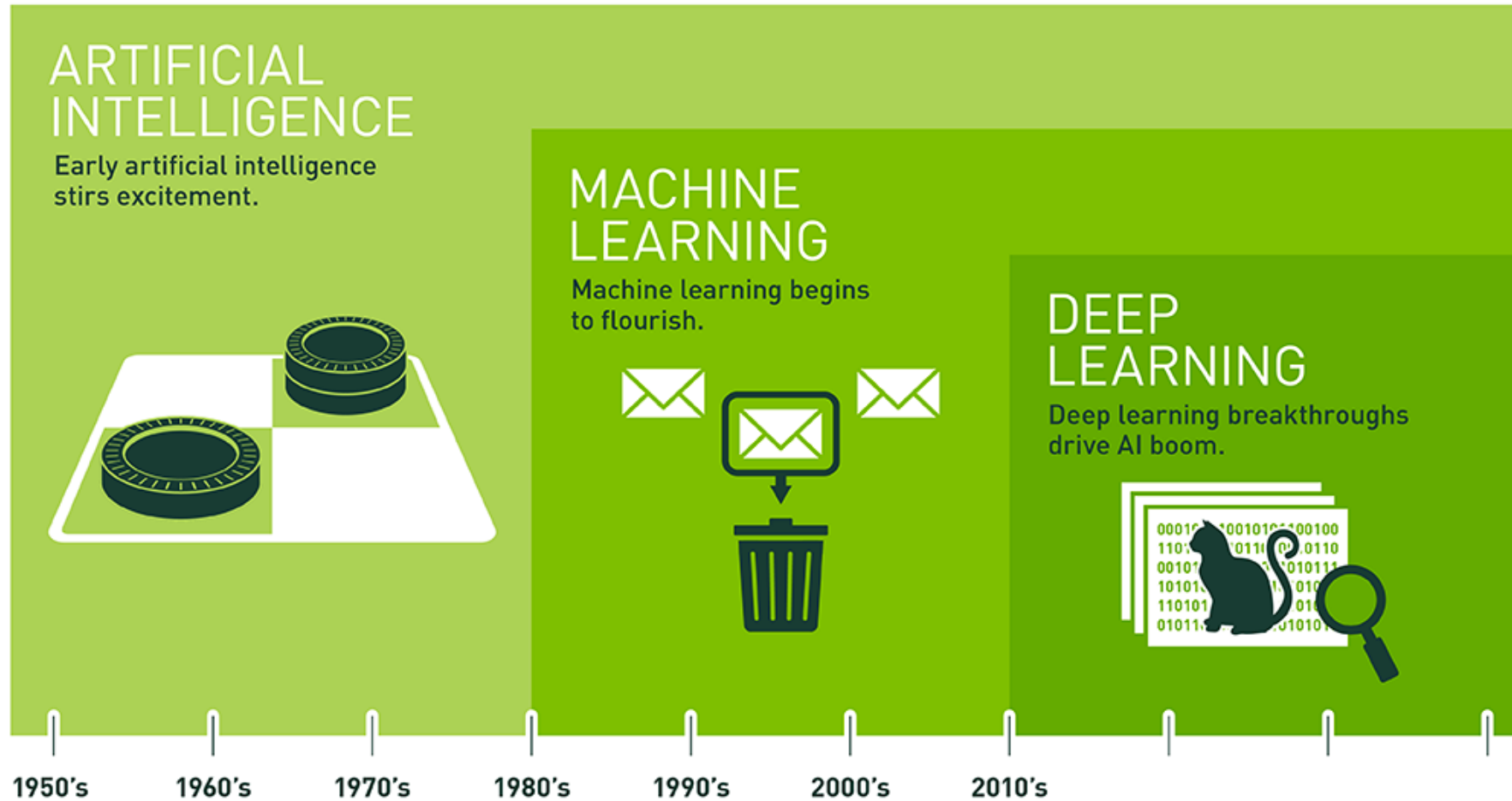
x

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

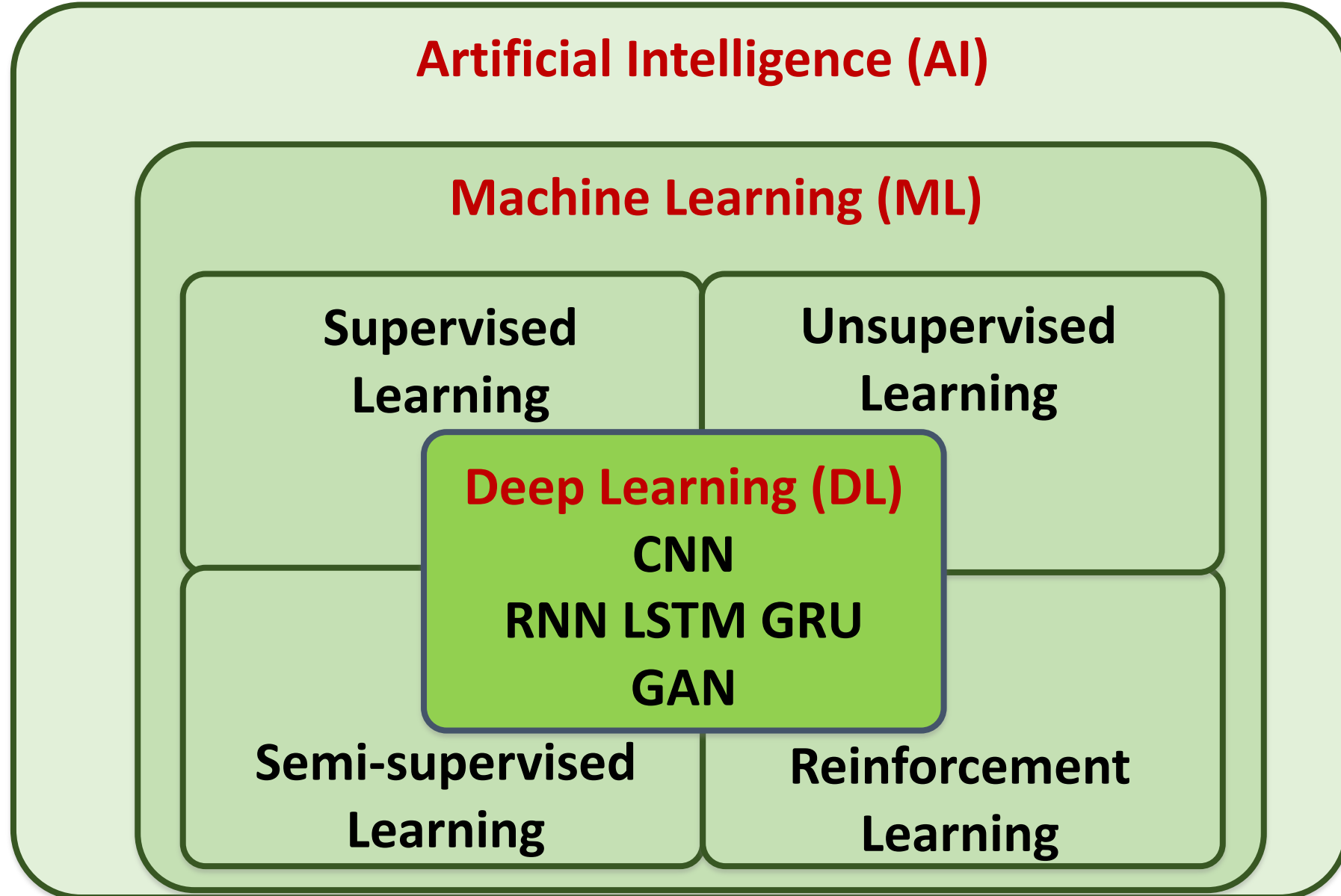
Artificial Intelligence

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) Learning**
 - **Vapnik-Chervonenkis (VC) Dimension**
 - **Bias-Variance Trade-off**
- **Overfitting and Underfitting**
 - **Avoid overfitting:
Regularization, Cross-Validation**

What is Learning in Machine Learning?

- **Learning** involves **finding patterns** from **data** to make **predictions**.
- Performance is measured through **generalization** to unseen data.
- Three paradigms of learning:
 - **Supervised Learning** (e.g., classification, regression)
 - **Unsupervised Learning** (e.g., clustering, dimensionality reduction)
 - **Reinforcement Learning** (e.g., learning through **rewards**)

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) Learning

- 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.

Probably Approximately Correct (PAC) Learning

- **PAC Learning** provides a way to understand **generalization**
- **Key Components:**
 - Error threshold (ϵ)
 - Confidence level ($1 - \delta$)
- **Example: Linear models are PAC-learnable with enough data**

Linear function

$$y = f(x)$$

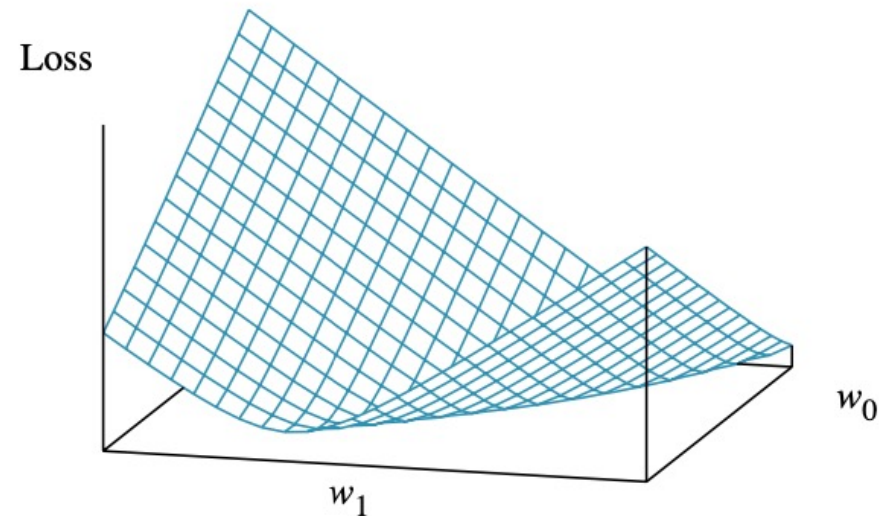
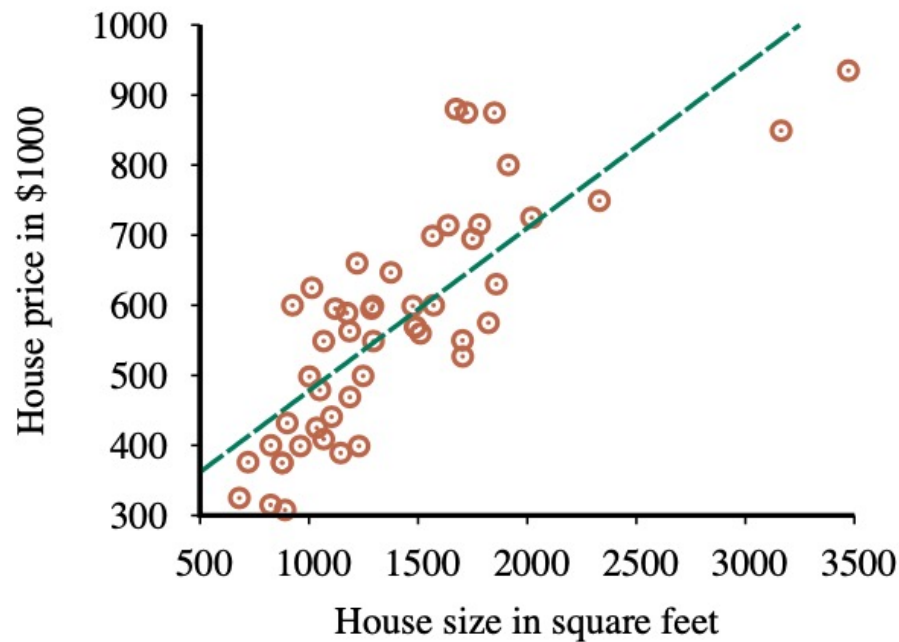
$$y = w_1 x + w_0$$

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

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)

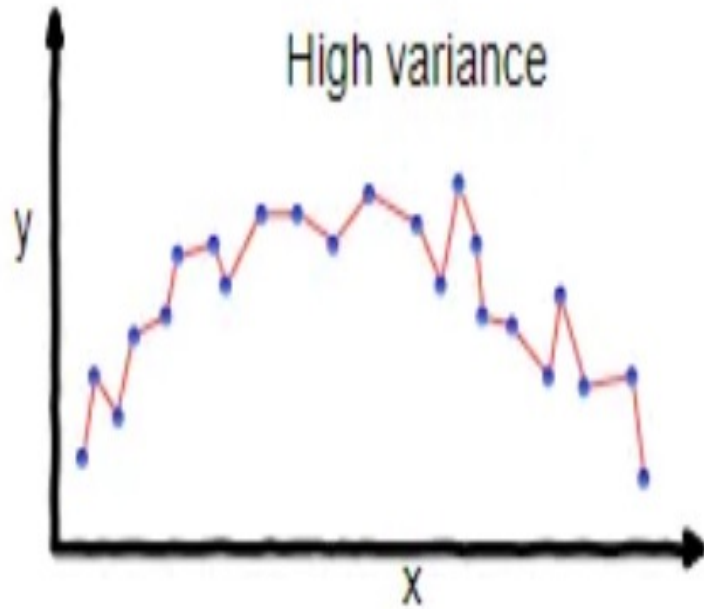
Vapnik-Chervonenkis (VC) Dimension

- **VC Dimension** measures a model's capacity to **fit various patterns**
- **Higher VC Dimension**
Higher model complexity
- **Balancing model complexity** helps improve **generalization**

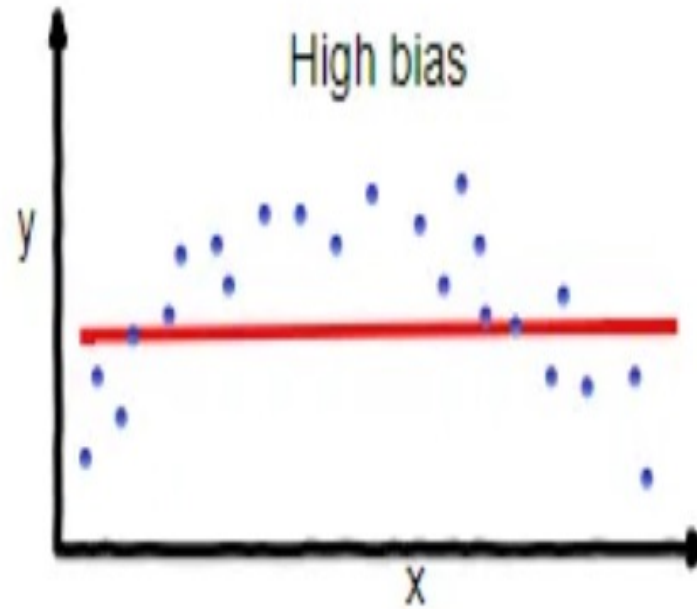
Bias-Variance Trade-off

- **Bias**: Error due to simple models (e.g., linear models)
- **Variance**: Error from models too sensitive to noise (e.g., deep trees)
- **Bias-Variance Trade-off**:
 - Low bias-high variance: Risk of **overfitting**
 - High bias-low variance: Risk of **underfitting**
- Selecting the right model **balances bias and variance**

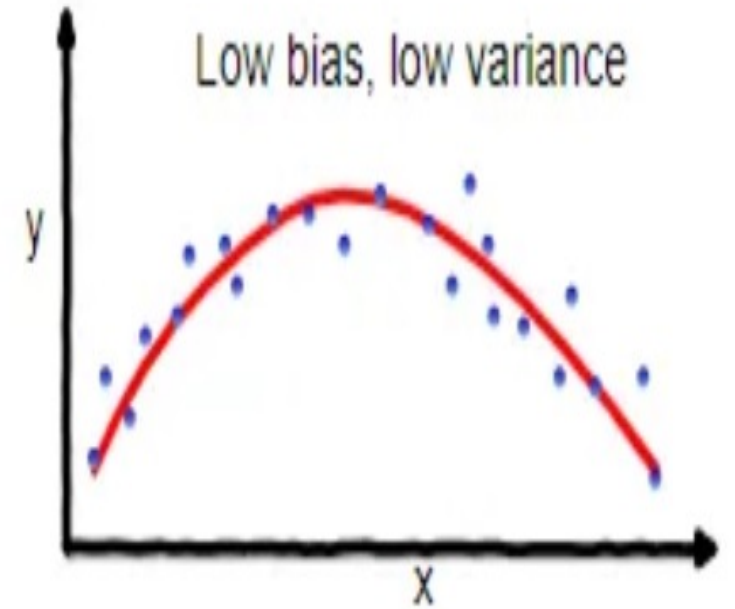
Overfitting vs. Underfitting



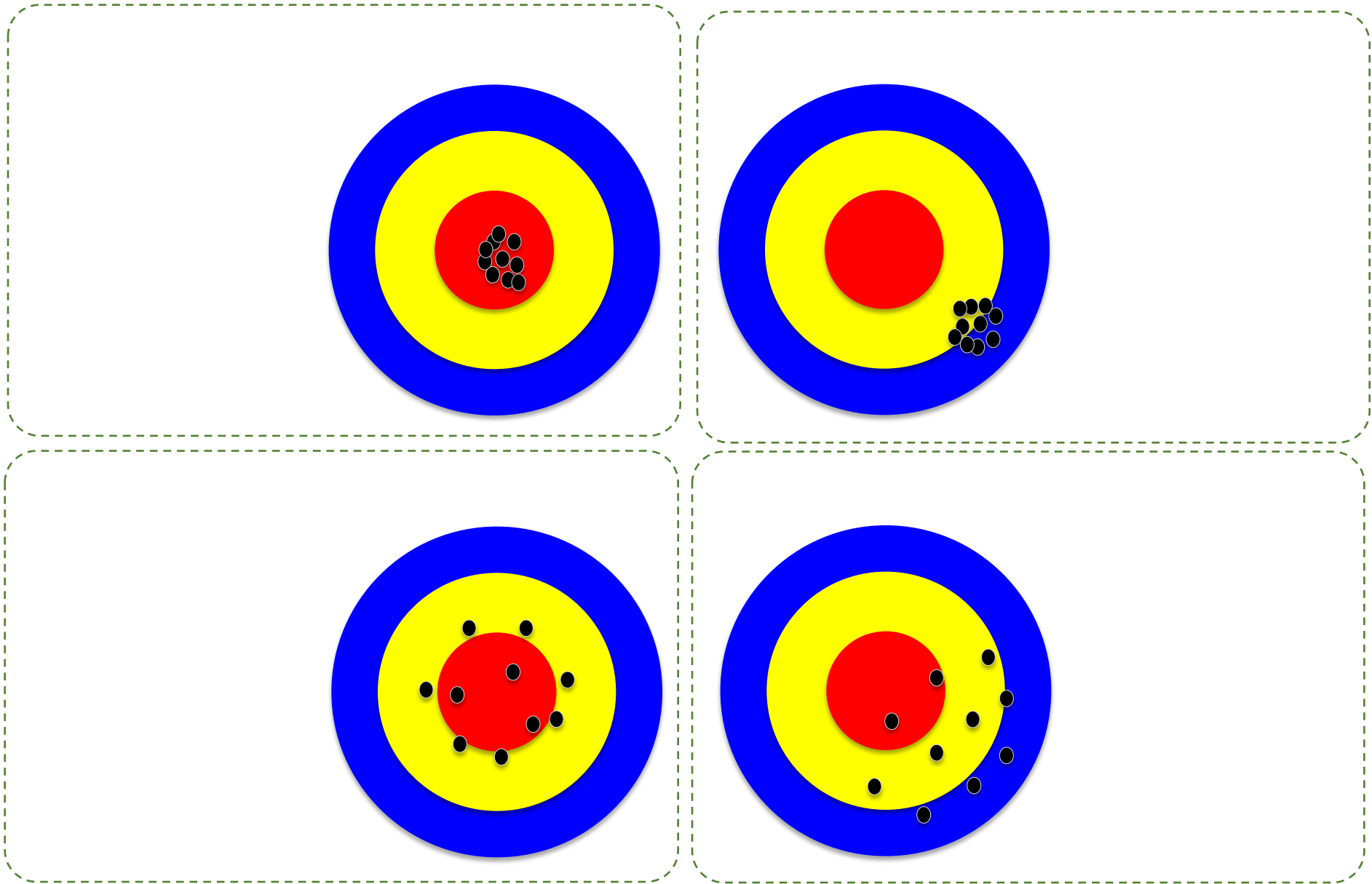
Overfitting



Underfitting



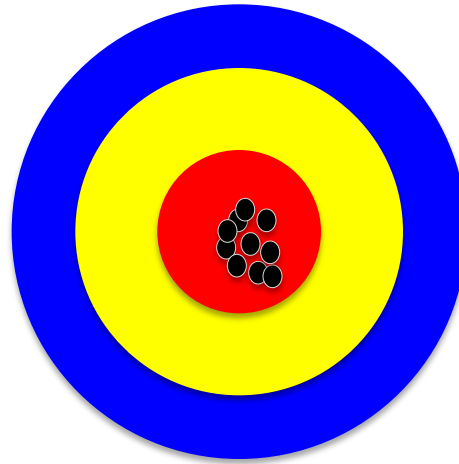
Good Balance



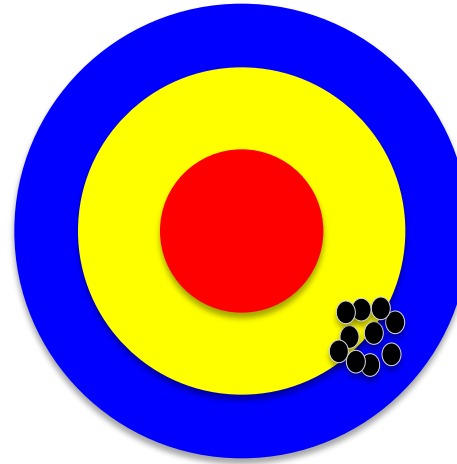
Bias vs. Variance

Low
Variance

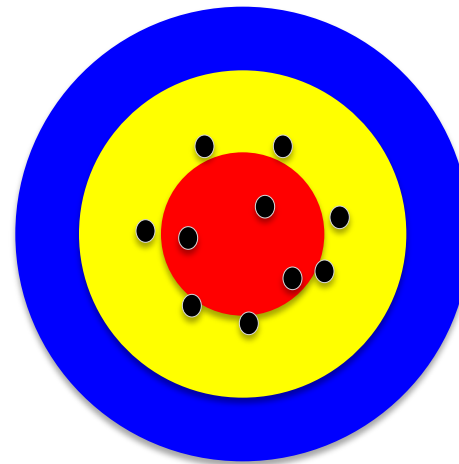
Low Bias
Low Variance



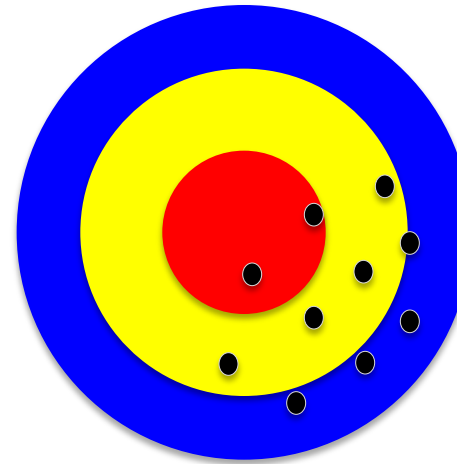
High Bias
Low Variance



Low Bias
High Variance



High Bias
High Variance

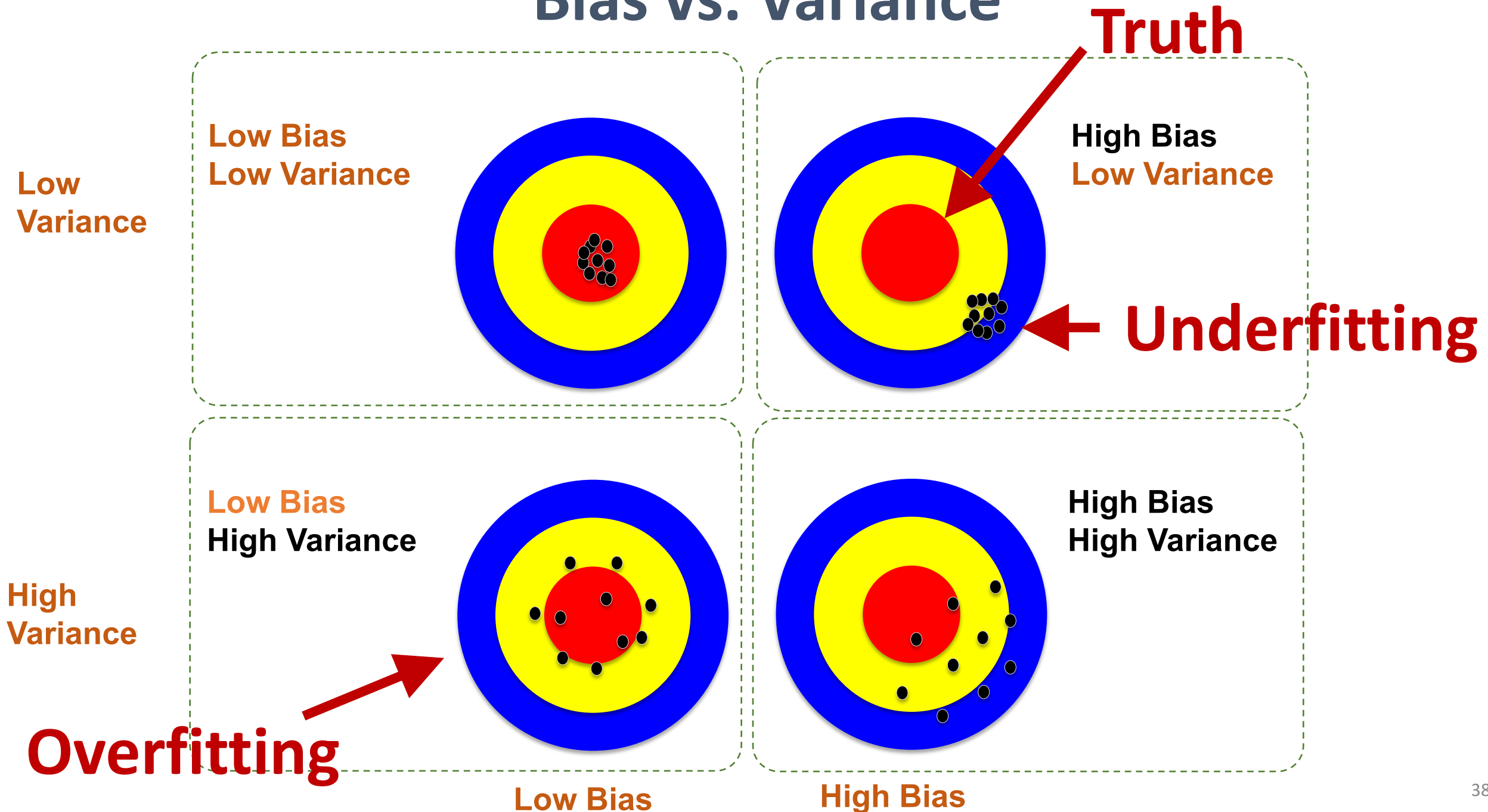


High
Variance

Low Bias

High Bias

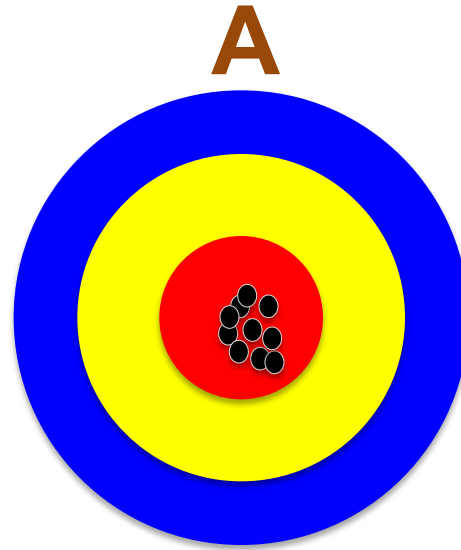
Bias vs. Variance



Bias vs. Variance

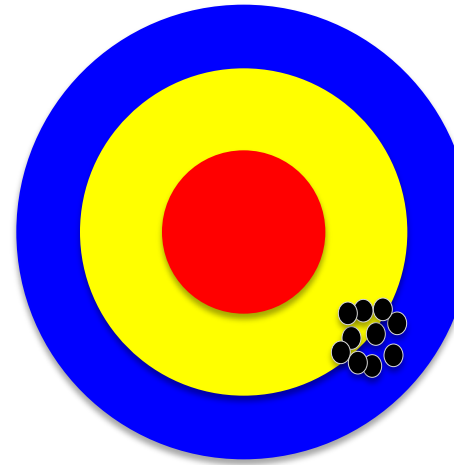
Low
Variance

Low Bias
Low Variance
High Accuracy
High Precision



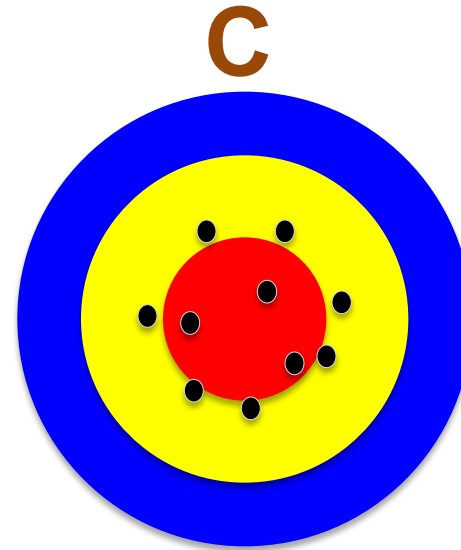
B

High Bias
Low Variance
Low Accuracy
High Precision



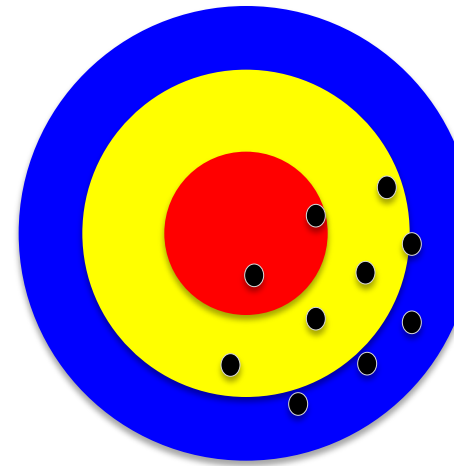
High
Variance

Low Bias
High Variance
High Accuracy
Low Precision



D

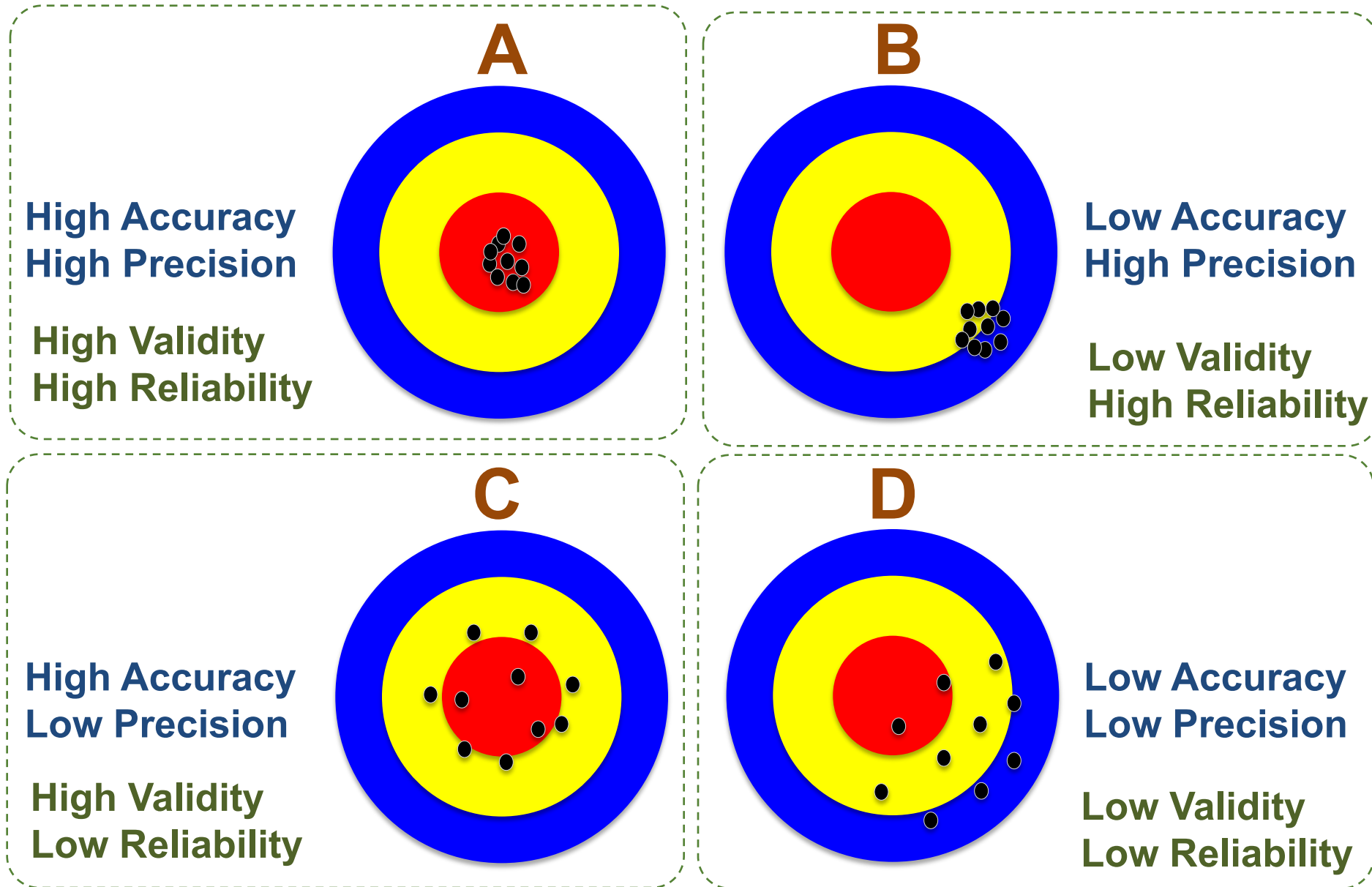
High Bias
High Variance
Low Accuracy
Low Precision



Low Bias

High Bias

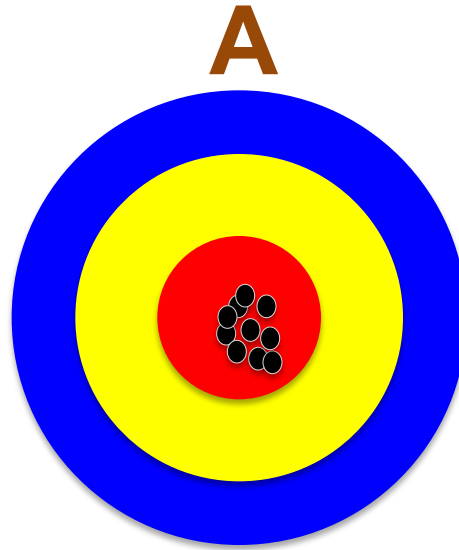
Accuracy vs. Precision



Accuracy (Validity) vs. Precision (Reliability)

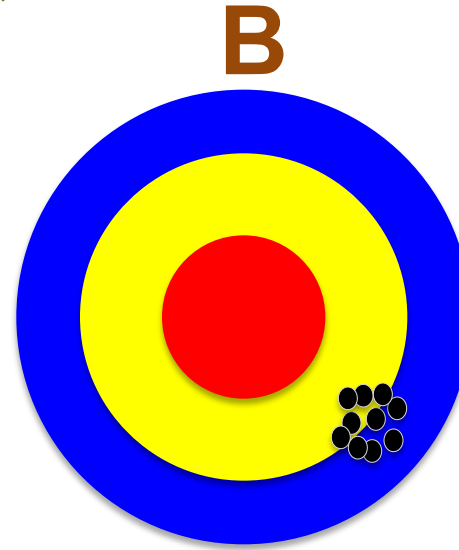
High Accuracy
High Precision

High Validity
High Reliability



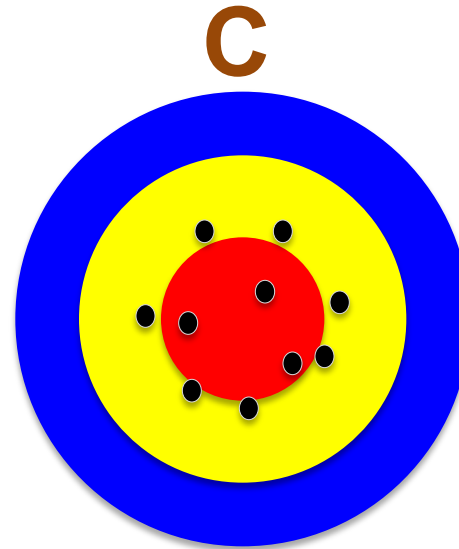
Low Accuracy
High Precision

Low Validity
High Reliability



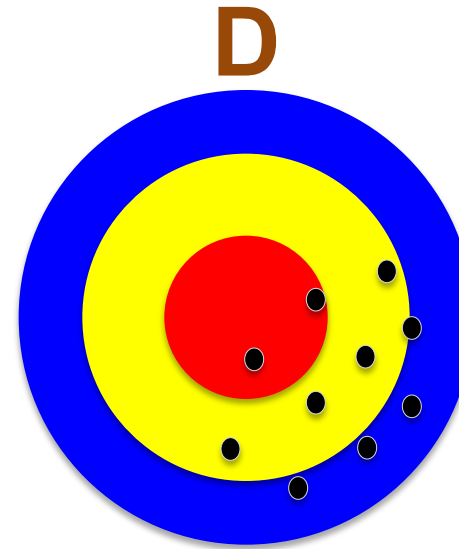
High Accuracy
Low Precision

High Validity
Low Reliability



Low Accuracy
Low Precision

Low Validity
Low Reliability



Learning Theory

- **PAC Learning:**
 - Focuses on **generalization**
- **VC Dimension:**
 - Measures model capacity and impacts **generalization**
- **Bias-Variance Trade-off:**
 - Helps in **model selection and tuning**

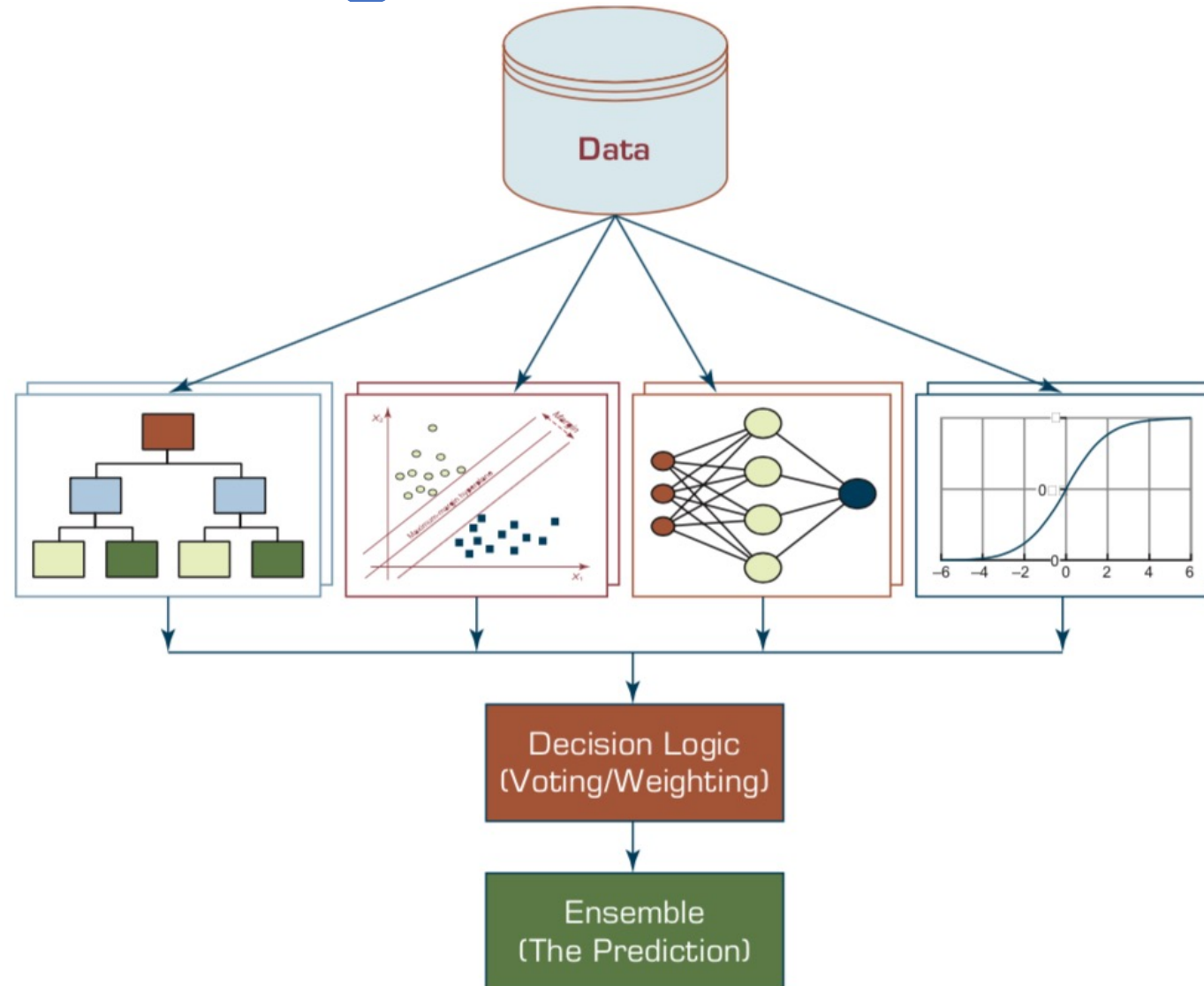
Ensemble Learning

Ensemble Learning

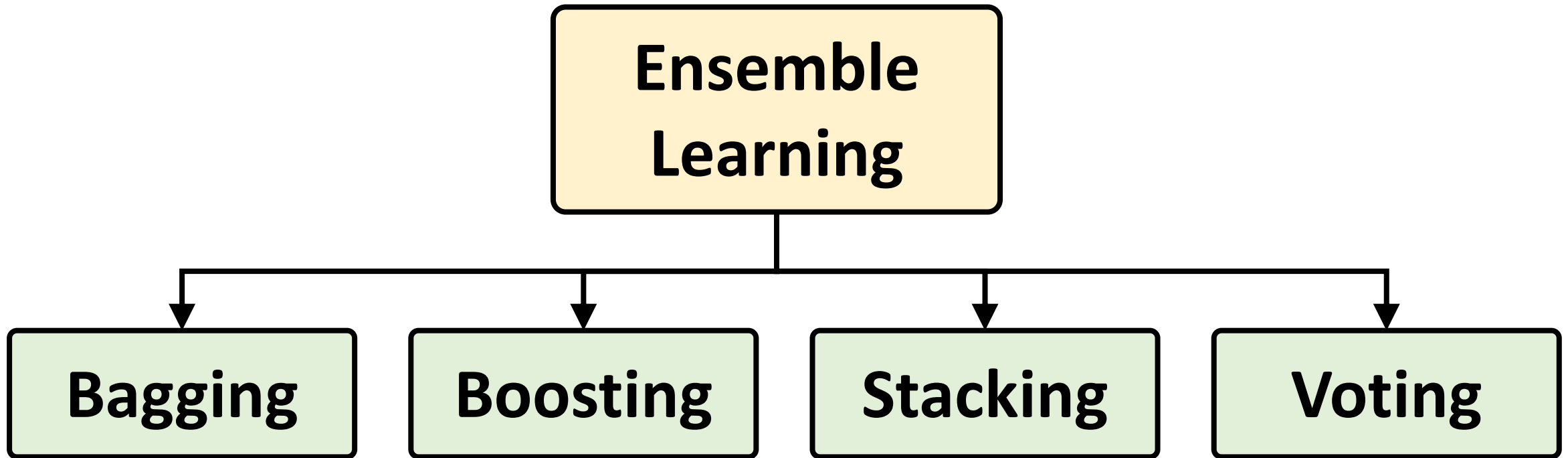
- **Select a collection, or ensemble,** of hypotheses, h_1, h_2, \dots, h_n , and **combine** their predictions by **averaging, voting,** or by another level of machine learning.

Ensemble Models

Heterogeneous Ensemble



Ensemble Learning: Bagging, Boosting, Stacking, Voting



Ensemble Learning

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

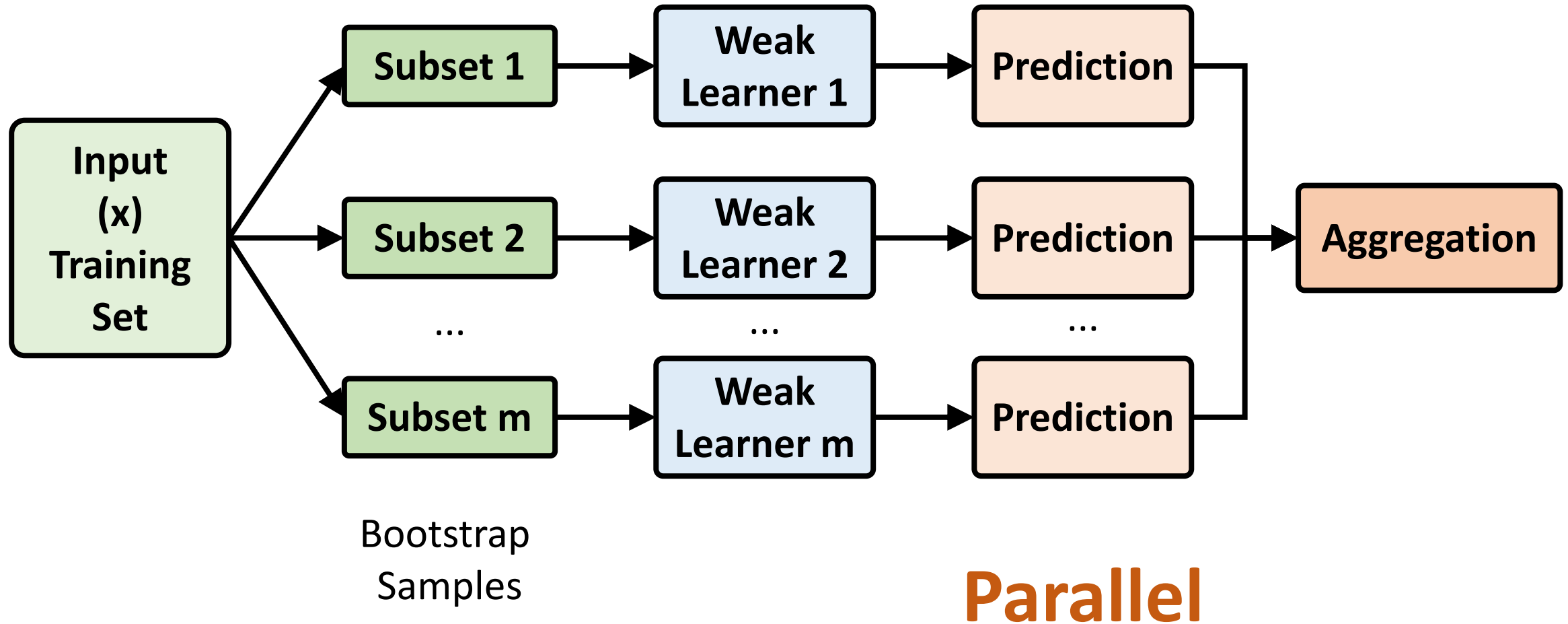
Why Ensemble Learning

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

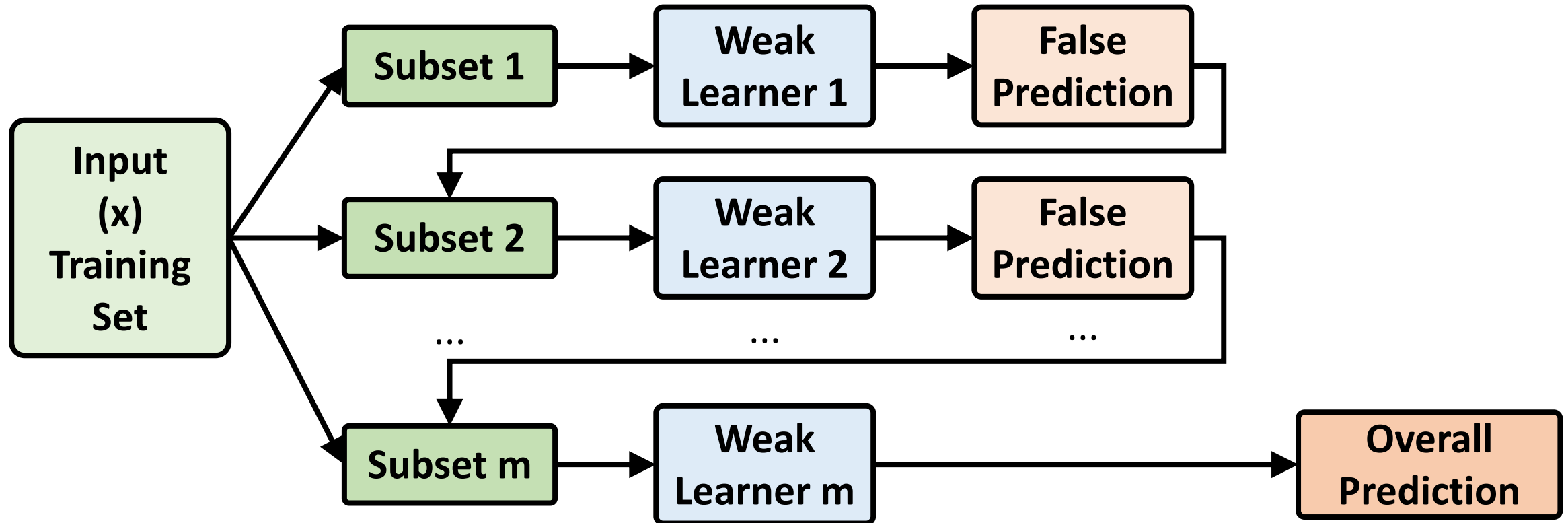
Ensemble Learning

- **Bagging**
 - **Random Forests (RF)**
- **Boosting**
 - **Gradient Boosting, XGBoost, LightGBM, CatBoost**
- **Stacking**
- **Online learning**

Ensemble Learning: Bagging (Bootstrap Aggregation)

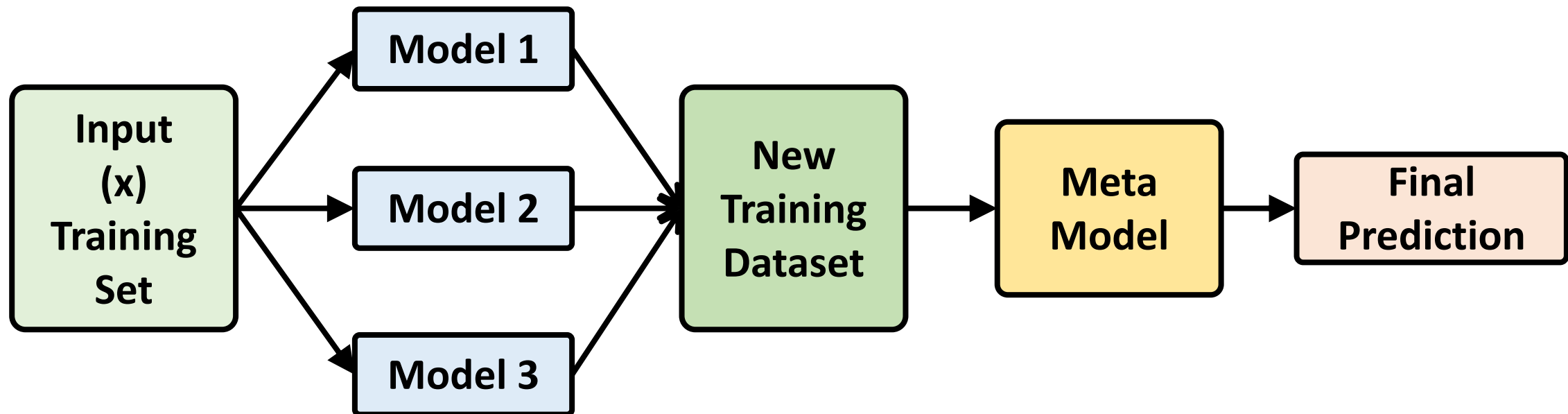


Ensemble Learning: Boosting

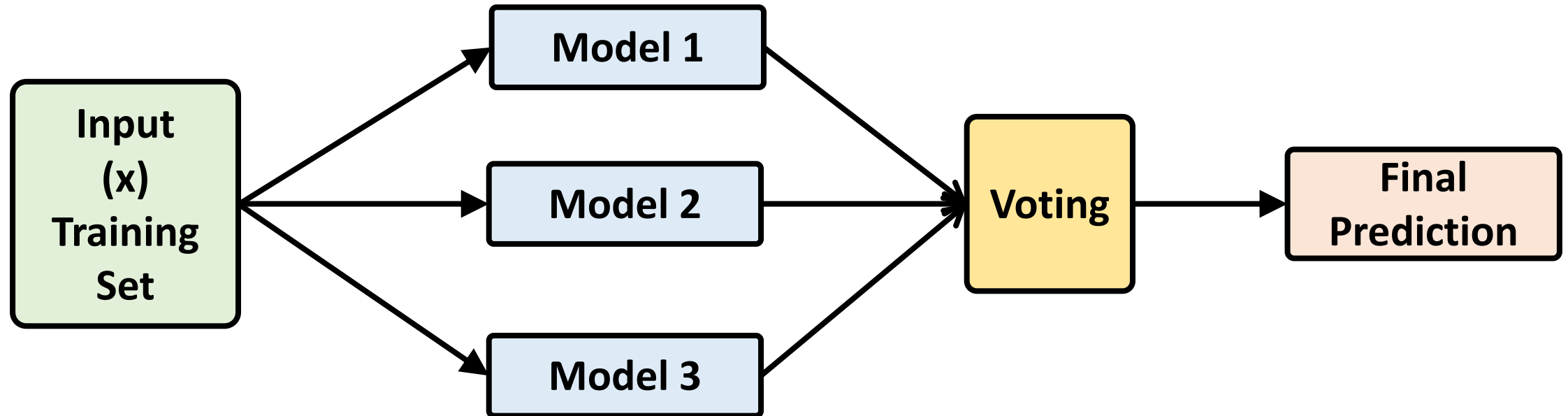


Sequential

Ensemble Learning: Stacking



Ensemble Learning: Voting



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)

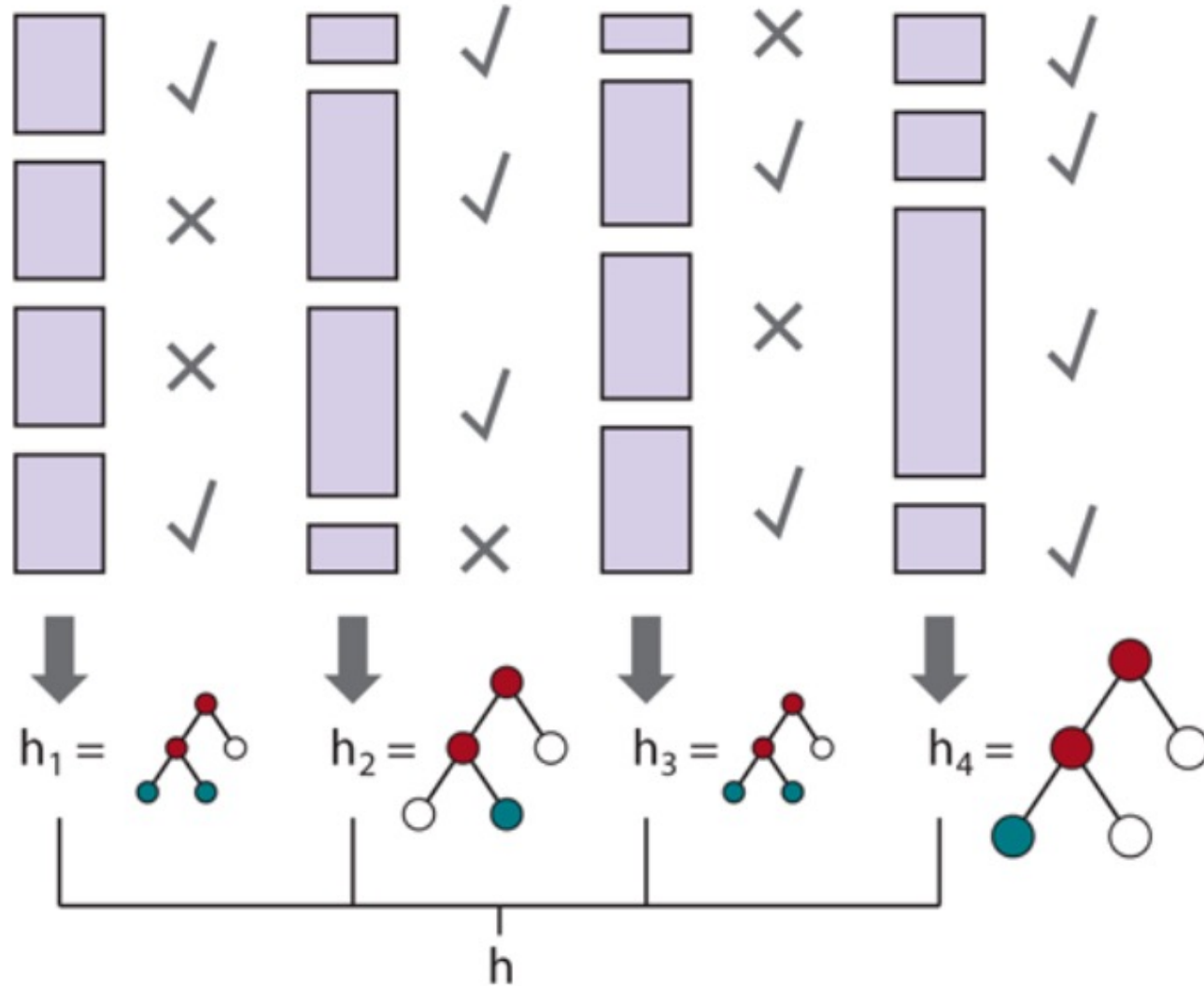
Ensemble Learning: Random forests

- **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: 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: Stacking

- **Staking**

- **Stacked generalization combines multiple base models from different model classes trained on the same data.**

- **Bagging**

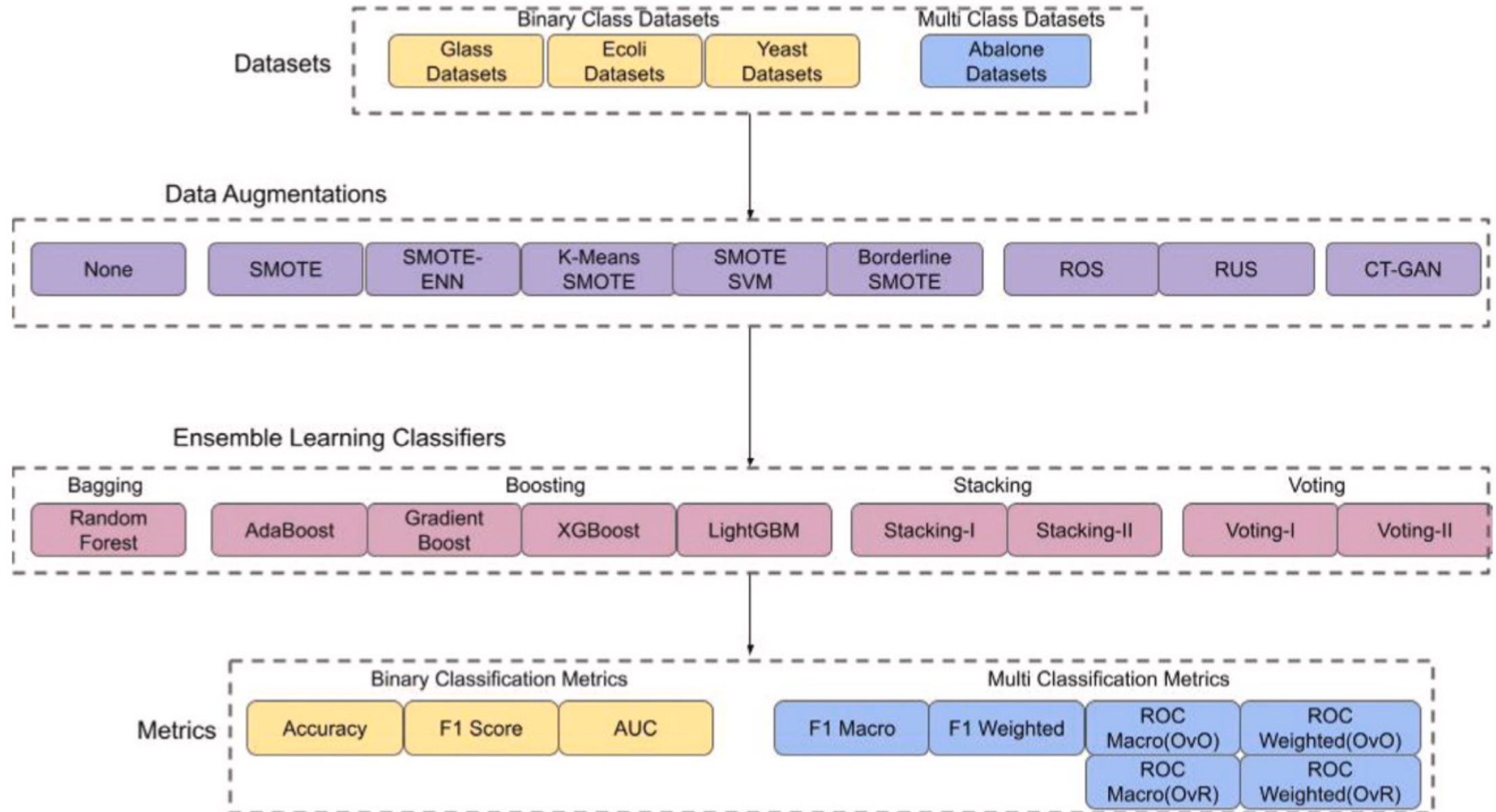
- **Combines multiple base models of the same model class trained on different data.**

Ensemble Learning:

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.**

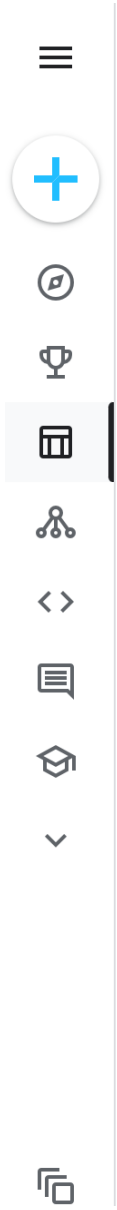
Ensemble Learning and Data Augmentations (DA)



ML Evaluation of Imbalanced Dataset: Ensemble Learning and Data Augmentations (DA)

Data augmentation	Ensemble model	Accuracy (Best, Std)	F1 (Best, Std)	AUC (Best, Std)
Yeast-v6				
Borderline-SMOTE	LightGBM	98.490(98.653, 0.150)	99.230(99.314, 0.077)	76.668(76.751, 0.077)
No augmentation	Stacking-I	98.185(98.653, 0.320)	99.076(99.315, 0.165)	69.757(76.579, 3.658)
SVM-SMOTE	AdaBoost	98.072(98.653, 0.466)	99.015(99.314, 0.240)	75.577(80.253, 2.014)
ROS	Stacking-I	97.997(98.822, 0.415)	98.977(99.401, 0.214)	74.884(76.837, 2.607)
Borderline-SMOTE	Voting-Soft	97.949(98.653, 0.501)	98.951(99.314, 0.259)	75.854(80.253, 1.923)
ROS	XGBoost	97.866(98.485, 0.303)	98.908(99.227, 0.157)	76.348(76.665, 0.155)
SMOTE	Stacking-II	97.539(98.485, 0.708)	98.737(99.227, 0.370)	77.385(83.841, 2.273)
SMOTE	Voting-Hard	97.386(98.485, 0.707)	98.657(99.227, 0.370)	77.607(83.841, 2.492)
SMOTE-ENN	Random Forests	92.917(97.306, 1.932)	96.273(98.618, 1.048)	73.380(78.962, 1.694)
RUS	Stacking-II	87.345(94.444, 3.405)	93.087(97.093, 2.014)	81.085(88.581, 3.112)

Kaggle Datasets for Machine Learning



Datasets

+ New Dataset

Your Work

Filters

- All datasets X
- Computer Science
- Education
- Classification
- Computer Vision
- NLP
- Data Visualization
- Pre-Trained Model

89 Datasets

Hotness ▾



S&P 500 ESG Risk Ratings
Pritish Dugar · Updated 4 months ago
Usability 10.0 · 1 File (CSV) · 252 kB

▲ 50

🥉 Bronze ...



Public Company ESG Ratings Dataset
Alistair King · Updated 8 months ago
Usability 10.0 · 1 File (CSV) · 43 kB

▲ 67

🥇 Gold ...




US Funds dataset from Yahoo Finance
Stefano Leone · Updated 3 years ago
Usability 10.0 · 7 Files (CSV) · 371 MB

▲ 246

🥈 Silver ...

Kaggle Datasets for Machine Learning



Search 










Datasets

[+ New Dataset](#) [Your Work](#)

[Filters](#)

- All datasets X
- Computer Science
- Education
- Classification
- Computer Vision
- NLP
- Data Visualization
- Pre-Trained Model

 **1,125 Datasets** [Most Votes](#)  

	Credit Card Fraud Detection Machine Learning Group - ULB · Updated 7 years ago Usability 8.5 · 1 File (CSV) · 69 MB	 11606  Gold ...
	Supermarket sales Aung Pyae · Updated 5 years ago Usability 8.8 · 1 File (CSV) · 37 kB	 2580  Gold ...
	Credit Card customers Sakshi Goyal · Updated 4 years ago Usability 10.0 · 1 File (CSV) · 388 kB	 2256  Gold ...

Kaggle Datasets: Credit Card Fraud Detection



Search



MACHINE LEARNING GROUP - ULB AND 1 COLLABORATOR · UPDATED 7 YEARS AGO

▲ 11606

New Notebook

Download



Credit Card Fraud Detection

Anonymized credit card transactions labeled as fraudulent or genuine



Data Card Code (4966) Discussion (107) Suggestions (0)

About Dataset

Context

It is important that credit card companies are able to recognize fraudulent credit card transactions so that customers are not charged for items that they did not purchase.

Content

The dataset contains transactions made by credit cards in September 2013 by European cardholders.

Usability ⓘ

8.53

License

Database: Open Database, Cont...

Expected update frequency

Not specified

Tags

Kaggle Code: Credit Card Fraud Detection

Credit Card Fraud Detection

▲ 11606

New Notebook

Download



Data Card Code (4966) Discussion (107) Suggestions (0)

All Your Work Shared With You Bookmarks

Most Votes ▼



Credit Fraud || Dealing with Imbalanced Datasets

Updated 5y ago

[657 comments](#) · Credit Card Fraud Detection

▲ 5453

Gold ...



Outlier!!! The Silent Killer

Updated 3y ago

[138 comments](#) · Titanic - Machine Learning from Disaster +16

▲ 1070

Gold ...



In depth skewed data classif. (93% recall acc now)

Updated 8y ago

[125 comments](#) · Credit Card Fraud Detection

▲ 796

Gold ...



Credit Card Fraud Detection Predictive Models

Updated 4y ago

[41 comments](#) · Credit Card Fraud Detection

▲ 489

Gold ...

Kaggle Code: Credit Card Fraud Detection



JANIO MARTINEZ BACHMANN · 5Y AGO · 922,344 VIEWS

▲ 5453

Edit My Copy 12207



Credit Fraud | | Dealing with Imbalanced Datasets

Python · Credit Card Fraud Detection

Notebook Input Output Logs Comments (657)

Run
861.1s

🕒 Version 70 of 70

- Finance
- Banking
- Data Visualization
- Classification
- Dimensionality Reduction

Credit Fraud Detector



Note: There are still aspects of this kernel that will be subjected to changes. I've noticed a recent increase of interest towards this kernel so I will focus more on the steps I took and why I took them to make it clear why I took those steps.

Before we Begin:

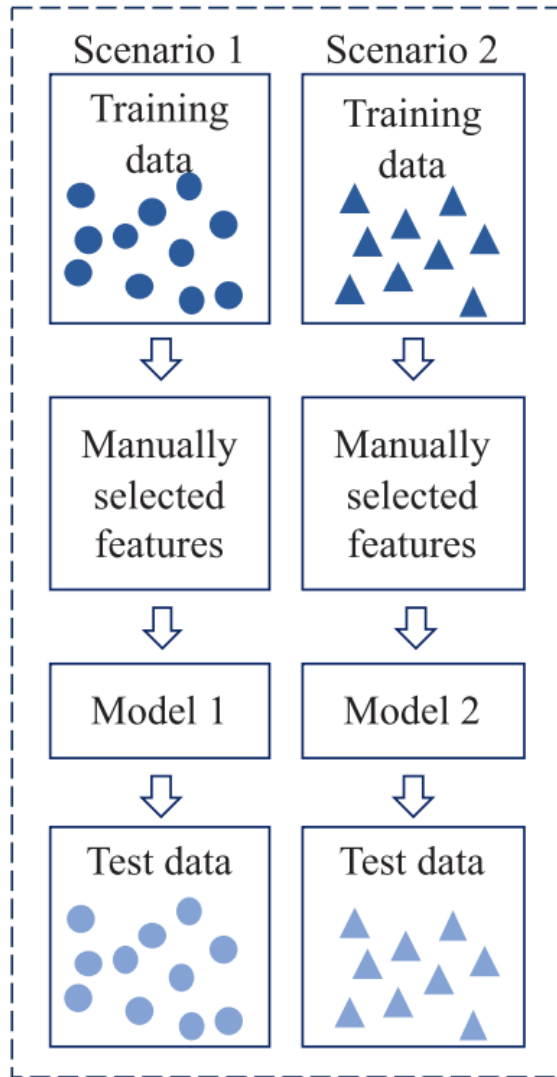
Meta Learning: Learning to Learn

Deep Learning
Transfer Learning
Few-Shot Learning
Meta Learning

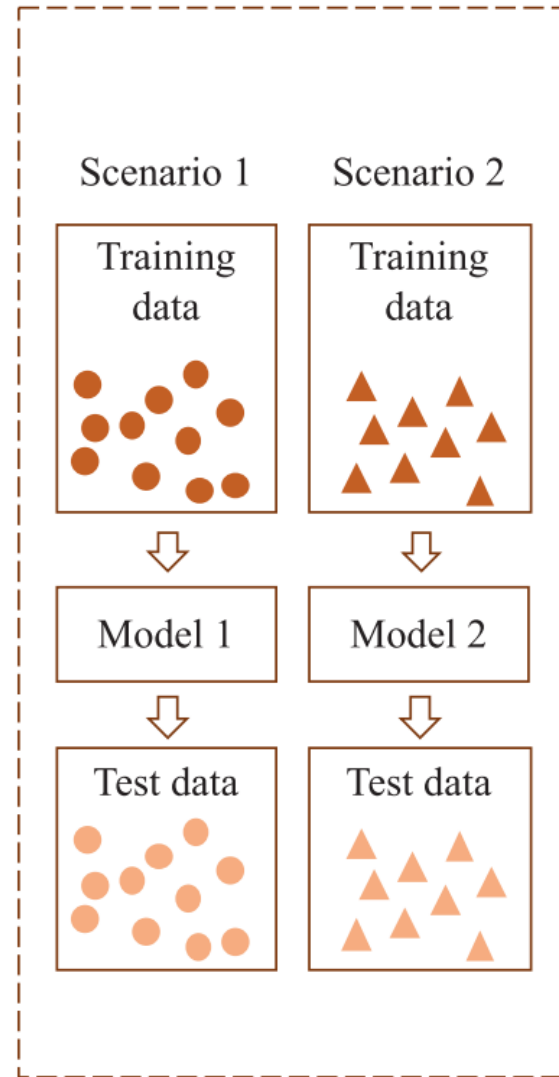
Deep Learning, Transfer Learning, Few-Shot Learning, Meta Learning

- **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 (OSL)(ZSL)**

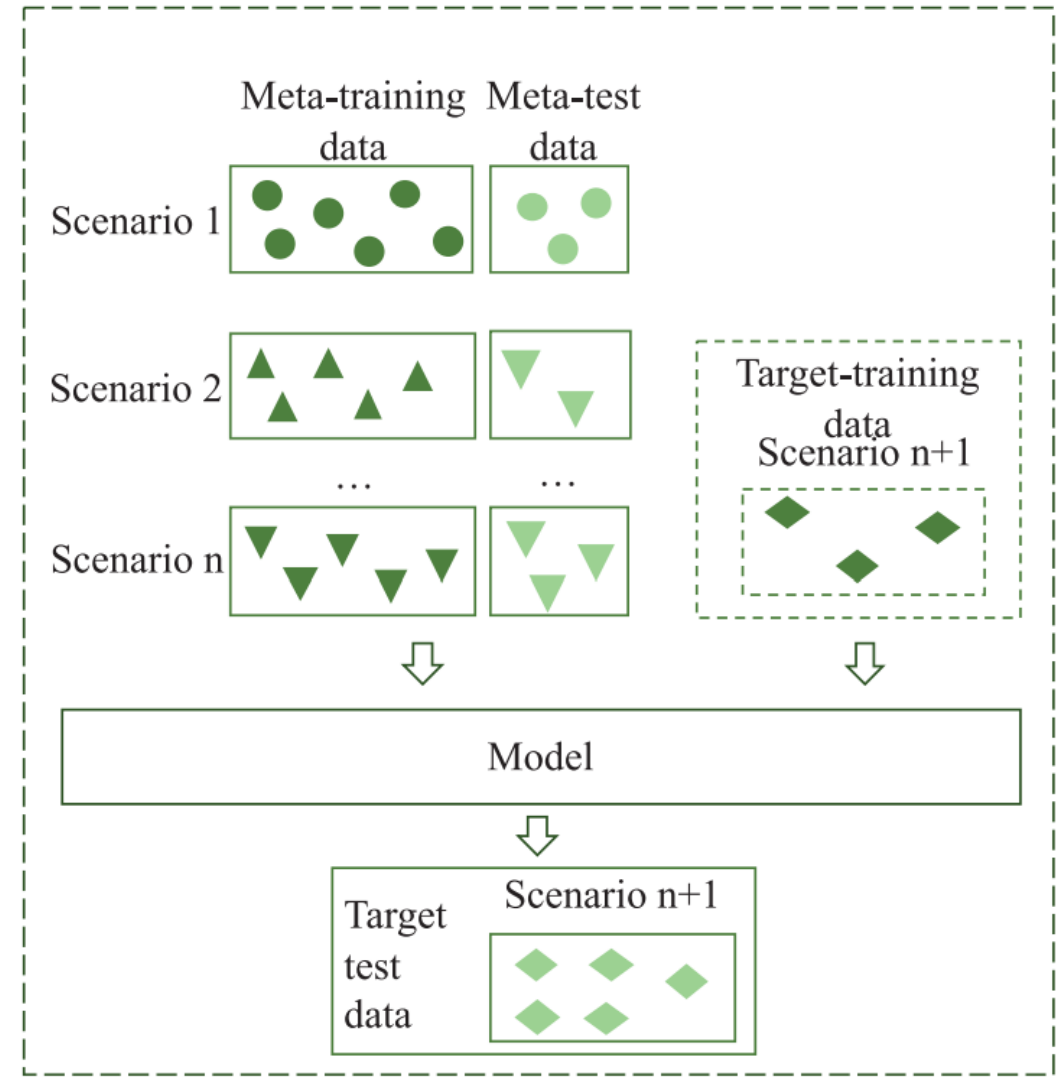
Machine Learning, Deep Learning, Meta Learning



Machine learning

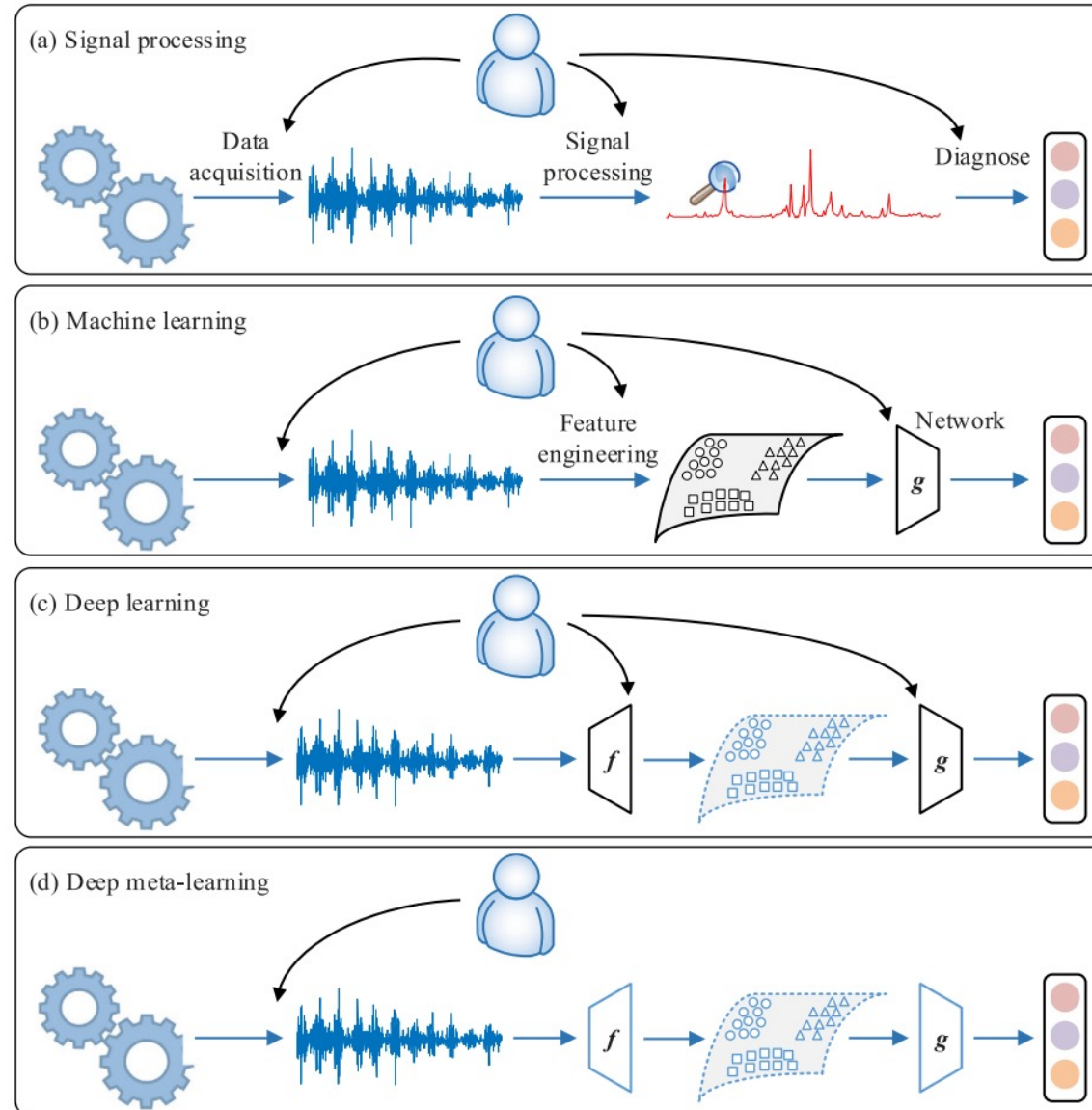


Deep learning



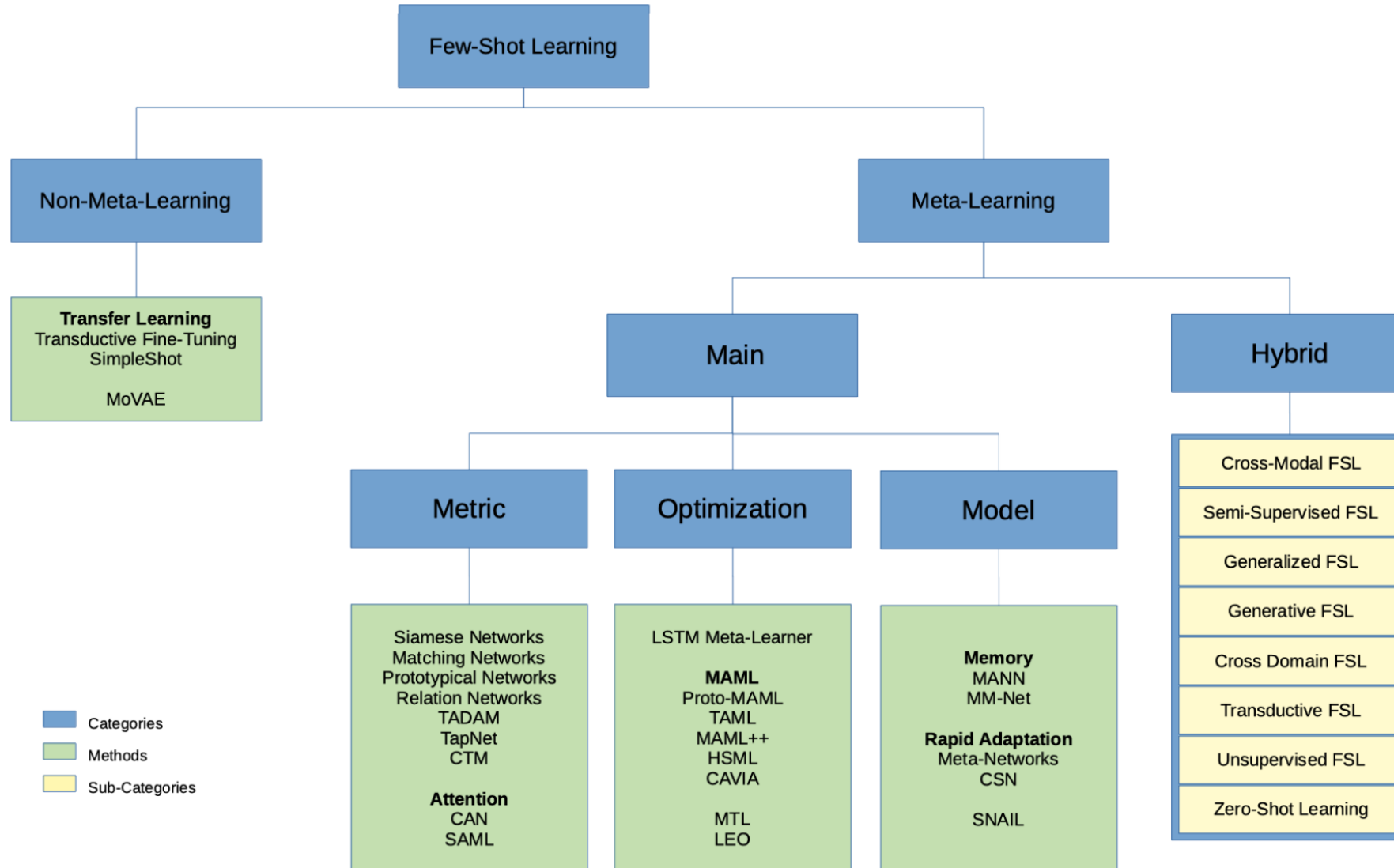
Meta-learning

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

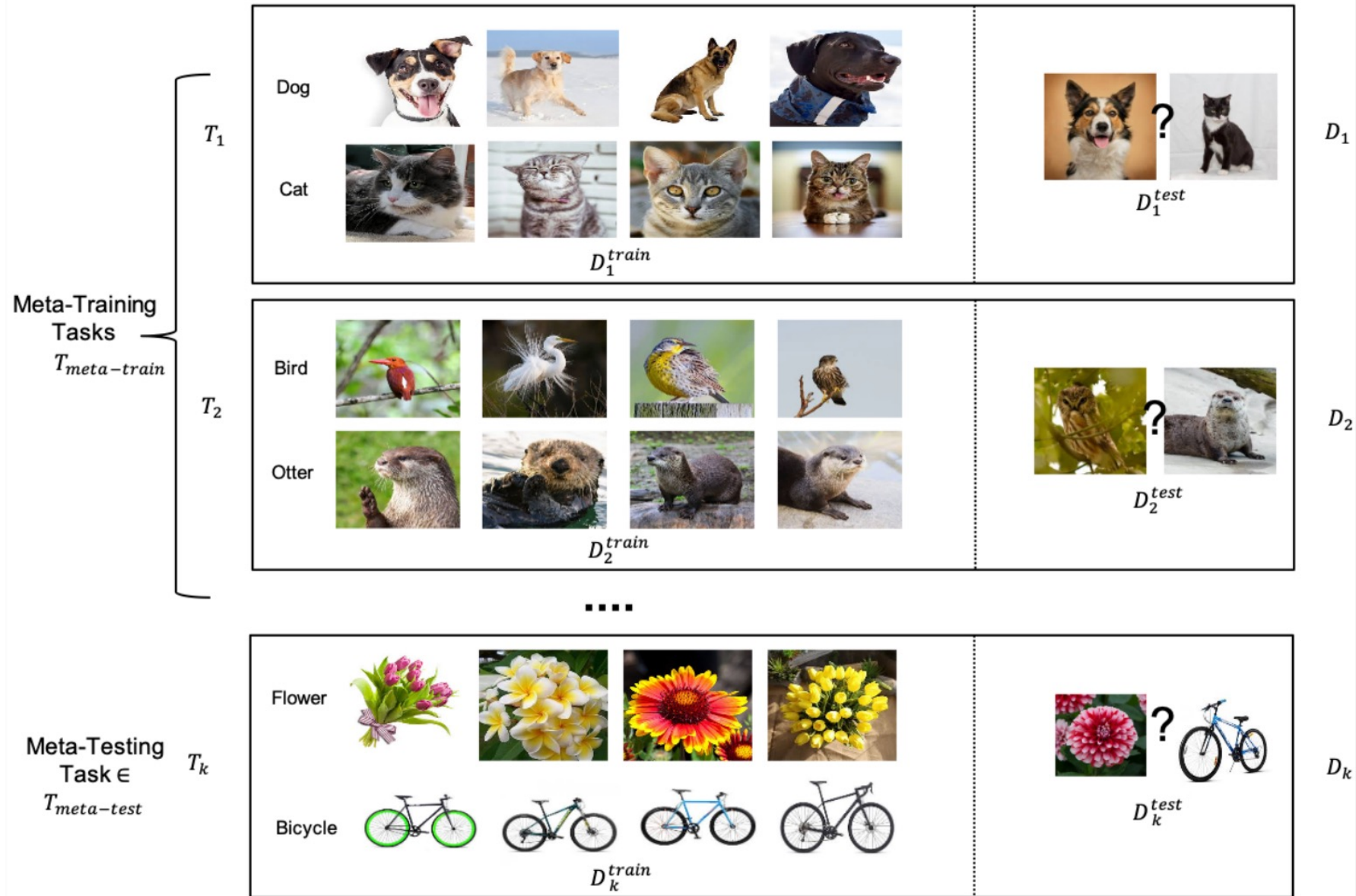
Metric-based Meta-learning

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_θ	Model (function) with parameters θ
g_{θ_1}	Embedding function
d_{θ_2} or d	Distance function
g_ϕ	Meta-Learning model with parameters ϕ
$P_\theta(y x)$	Output probability of y for input x using model parameters θ
$k_\theta(x_1, x_2)$	Kernel function measuring similarity between two vectors x_1 and x_2
σ	Softmax function
α, β	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

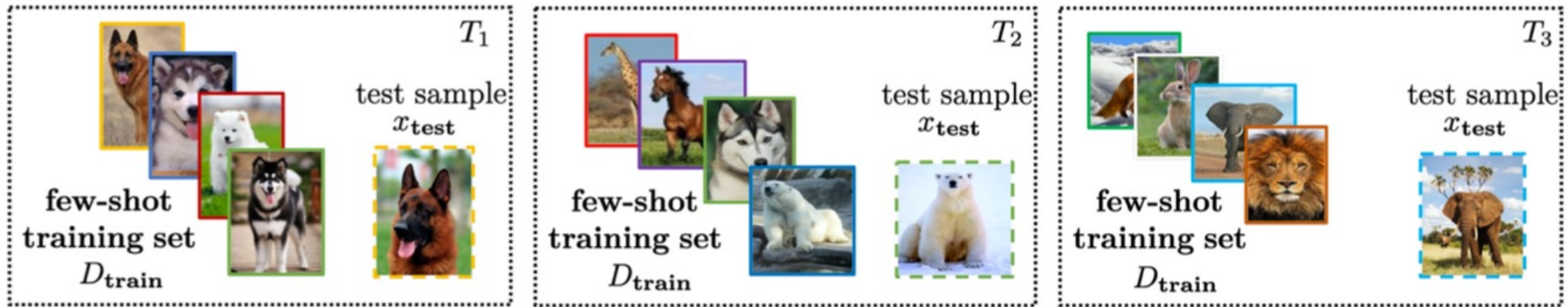
Meta-Learning Example Setup



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)

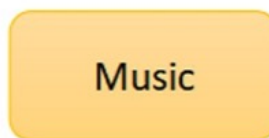
Meta-learning

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

Training Task 1



Training Task 2 ...



Test Task 1 ...



Support Set:

Only_[O] France_[LOC] and_[O] Britain_[LOC]
backed_[O] Fischler_[PER]'s_[O] proposal_[O].

Query Set:

Adrian_[PER] Warner_[PER] has_[O] lived_[O]
in_[O] Brussels_[LOC] since_[O] 1996_[O].

Labels: {PER, PER, O, O, O, LOC, O, O}

Support Set:

Play_[O] rap_[album] album_[album] one_[album]
by_[O] Gene_[artist] Vincent_[artist].

Query Set:

Add_[O] Rosemary_[artist] Clooney_[artist]
to_[O] pura_[playlist] vida_[playlist] playlist_[O].

Labels: {O, artist, artist, O, playlist,
playlist, O}

Support Set:

Patient_[O] received_[O] combivent_[DRUG]
nebs_[Dosage Form/Route], solumedrol_[DRUG]
125mg_[AMOUNT] IV_[Dosage Form/Route]
X1_[Dosage Frequency].

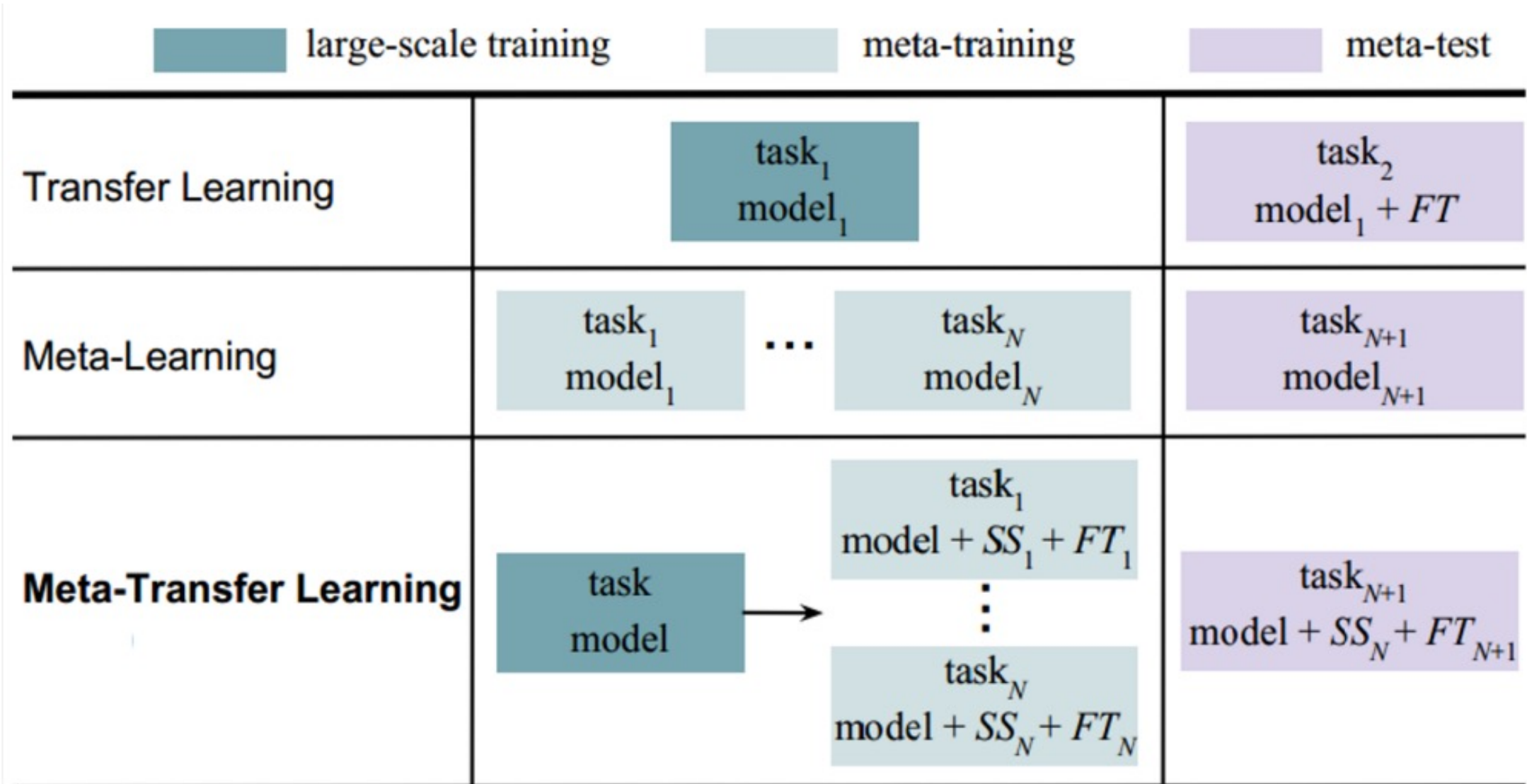
Query Set:

She was given a dose of levaquin this
morning.

Labels: {O, O, O, AMOUNT, AMOUNT,
O, DRUG, TIME, TIME}

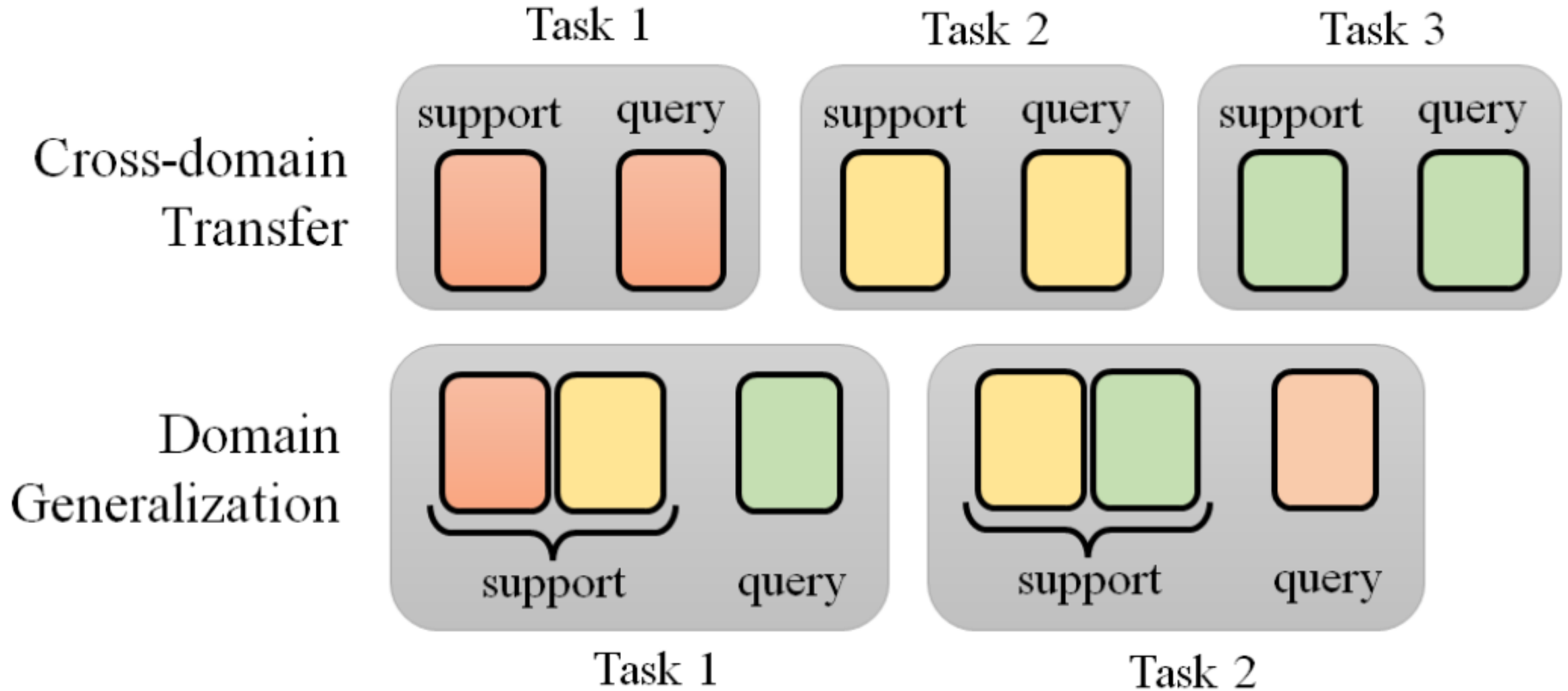
Meta-Task Learning (MTL)

Transfer Learning Strategy for Meta-Learning

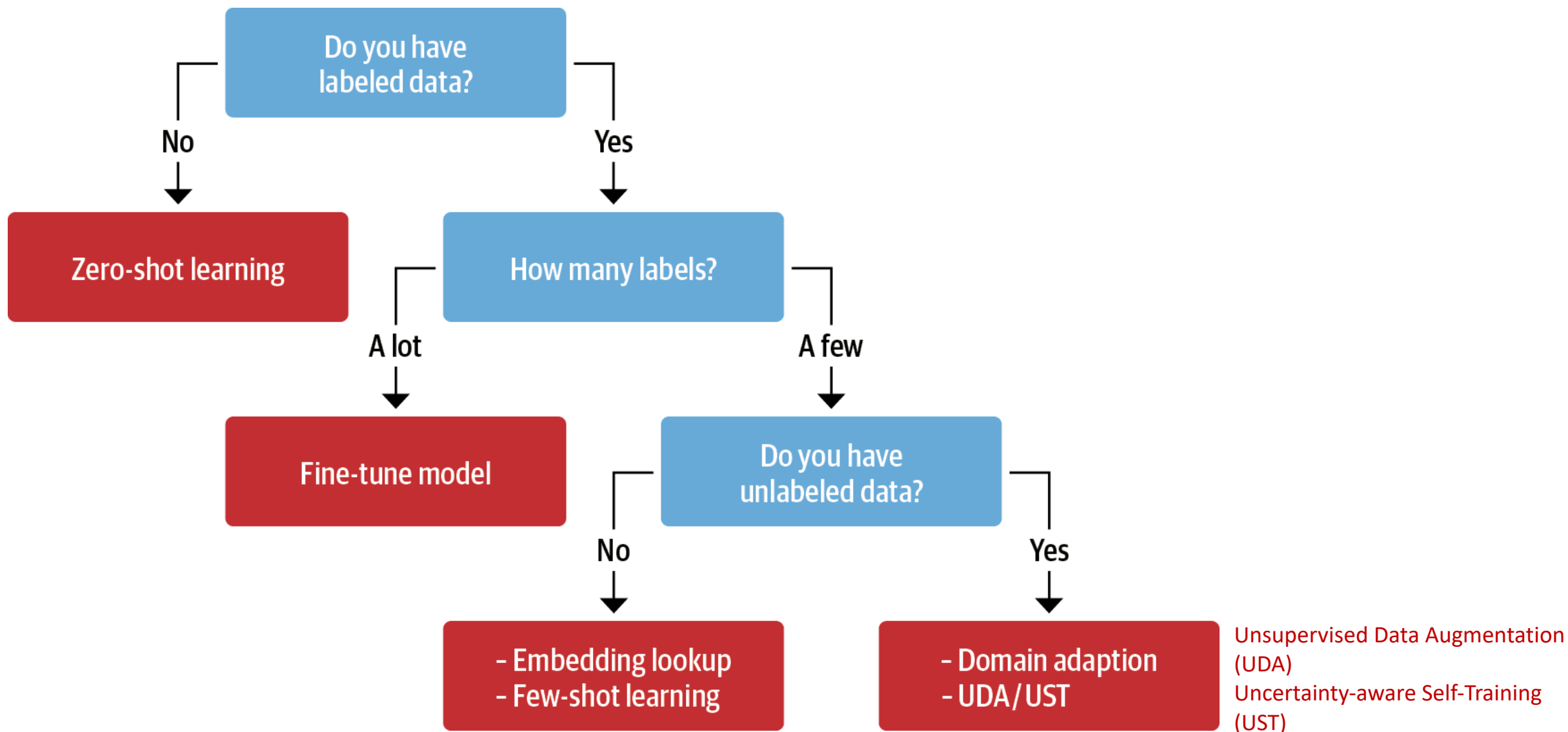


Meta Learning

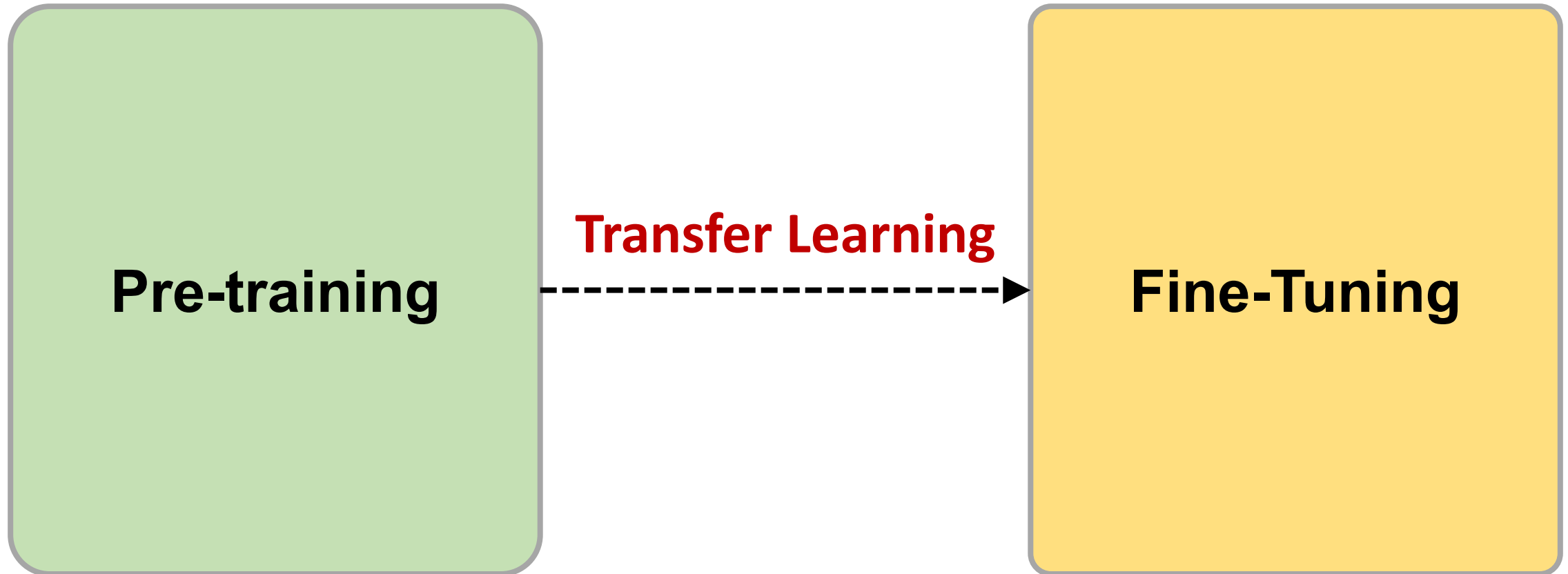
The task construction of cross-domain transfer and domain generalization



Transfer Learning, Fine-tuning, Few-shot learning



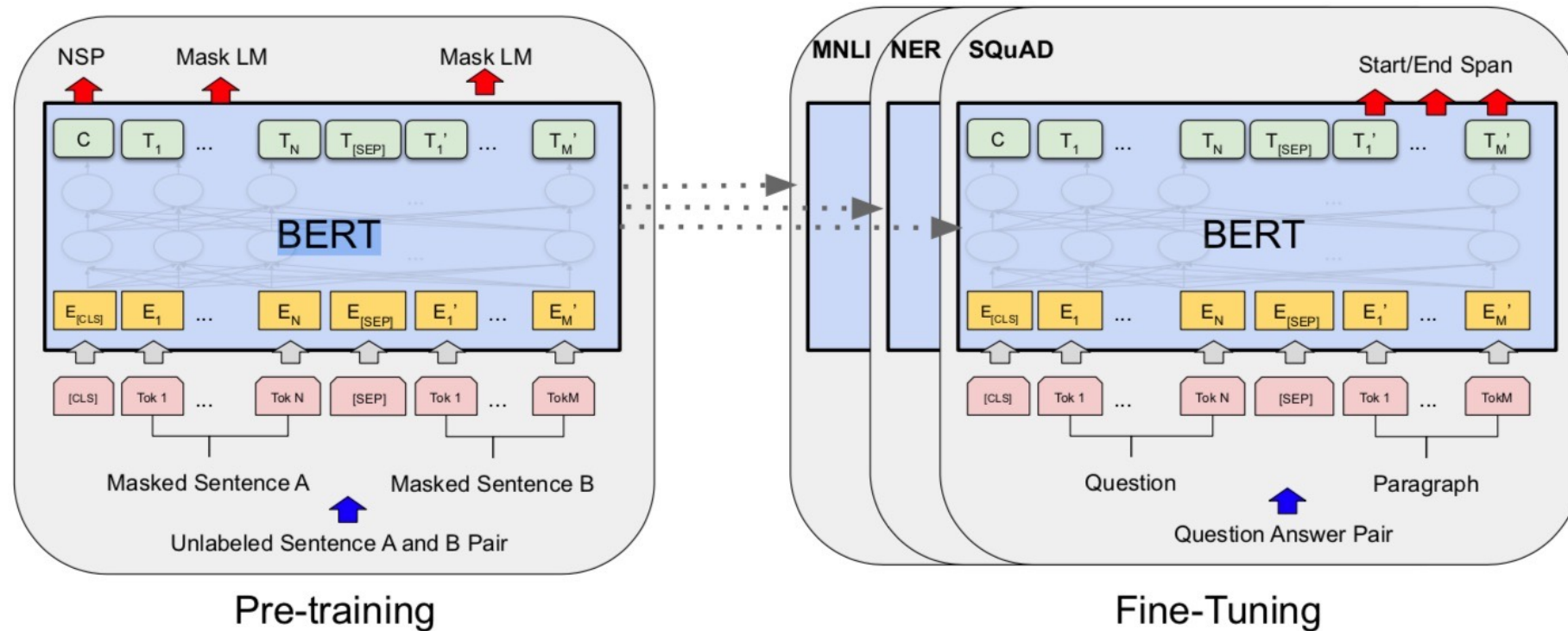
Transfer Learning



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

BERT (Bidirectional Encoder Representations from Transformers)

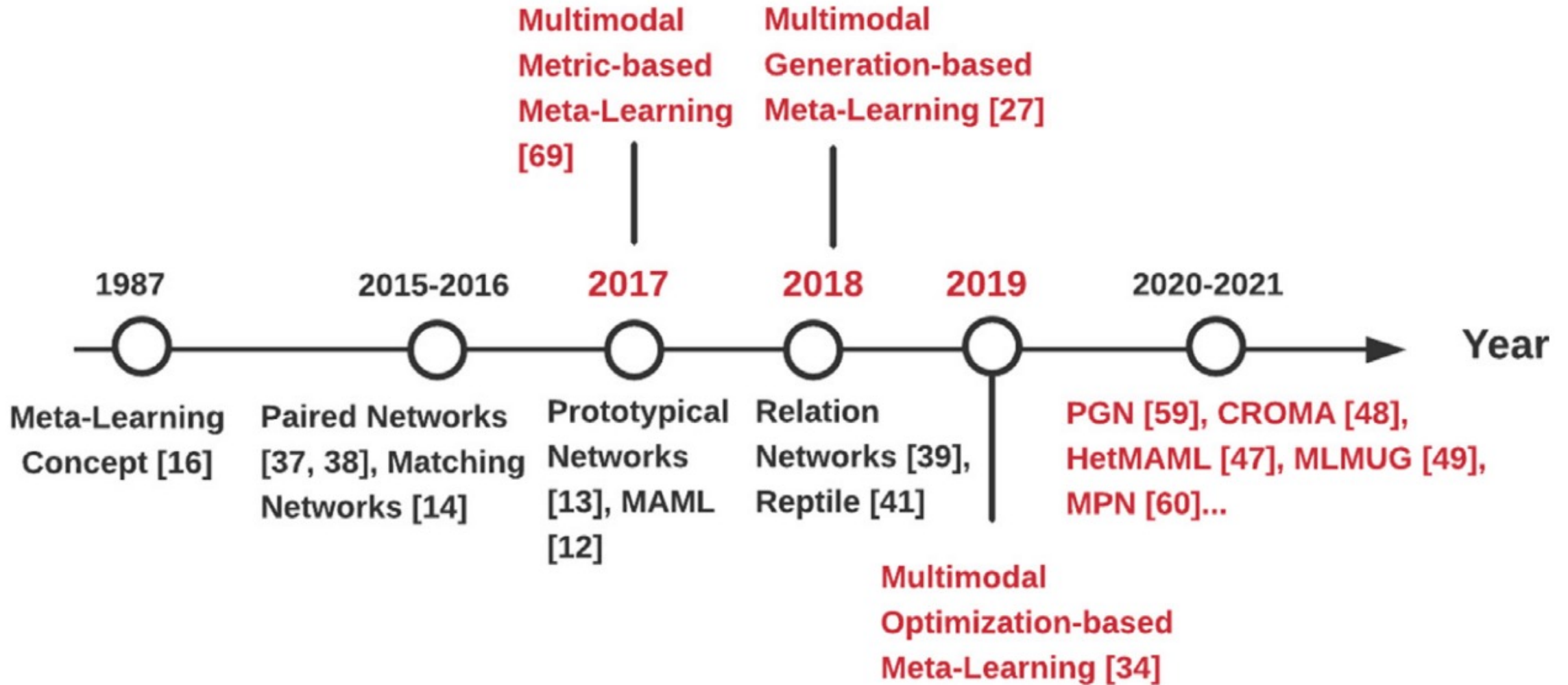
Overall pre-training and fine-tuning procedures for BERT



Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).

"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

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

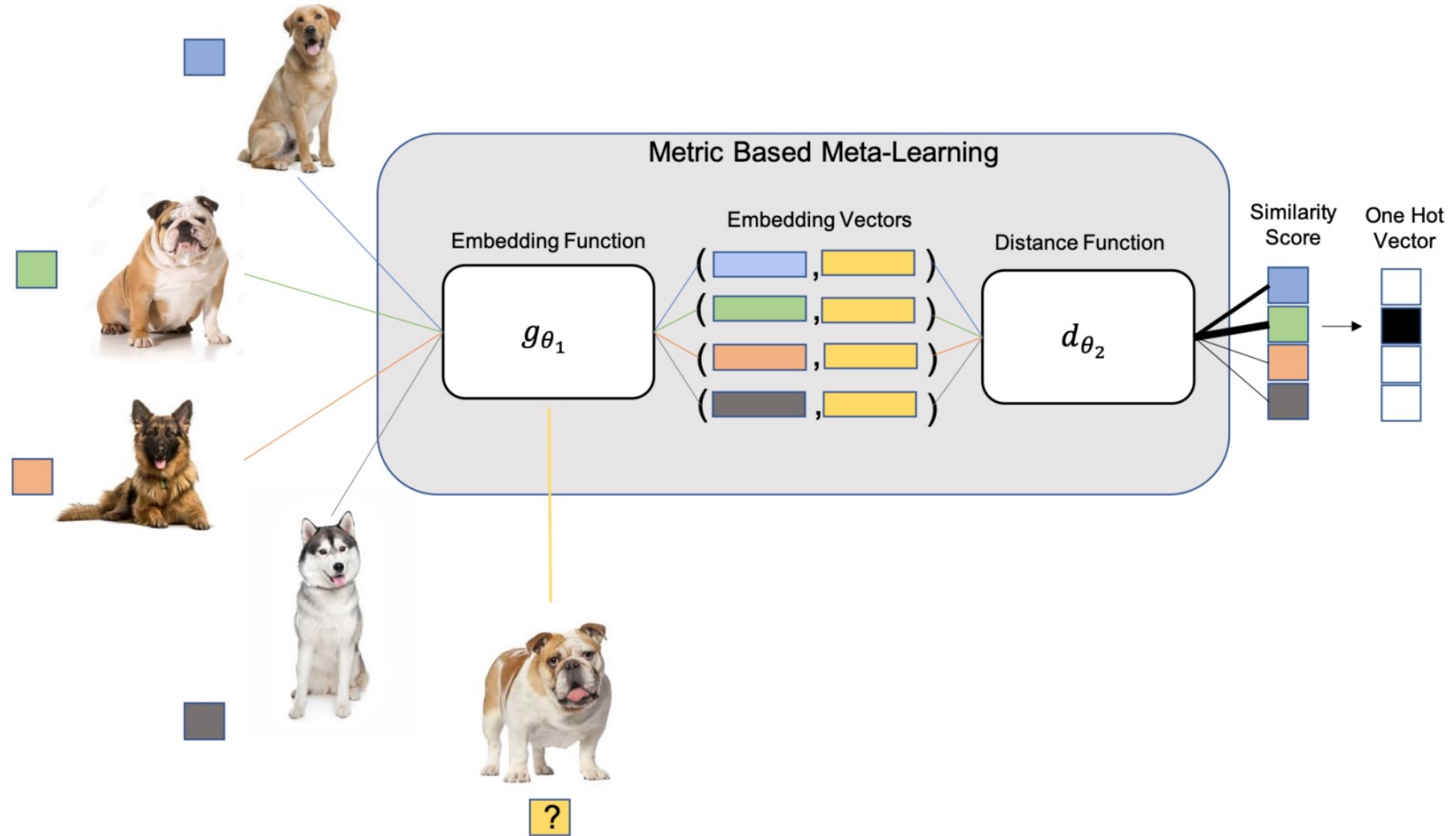
	Metric-based	Optimization-based	Model-based
Key idea	Metric Learning [19]	Gradient Descent	Memory; RNN
How $P_\theta(y x)$ is modeled?	$\sum_{(x_k, y_k) \in S} k_\theta(x, x_k) y_k,$	$P_{\theta'}(y x),$ <p>where $\theta' = g_\phi(\theta, S)$</p>	$f_\theta(x, S).$
Advantages	<p>Faster Inference.</p> <p>Easy to deploy.</p>	<p>Offers flexibility to optimize in dynamic environments.</p> <p>S 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>
Disadvantages	<p>Less adaptive to optimization in dynamic environments.</p> <p>Computational complexity grows linearly with size of S 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 S 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]	
Optimization-Based	MAML	[72–80]	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- LSTM	[81–86]	
	META- SGD	[87–93]	

Metric-based Meta-learning

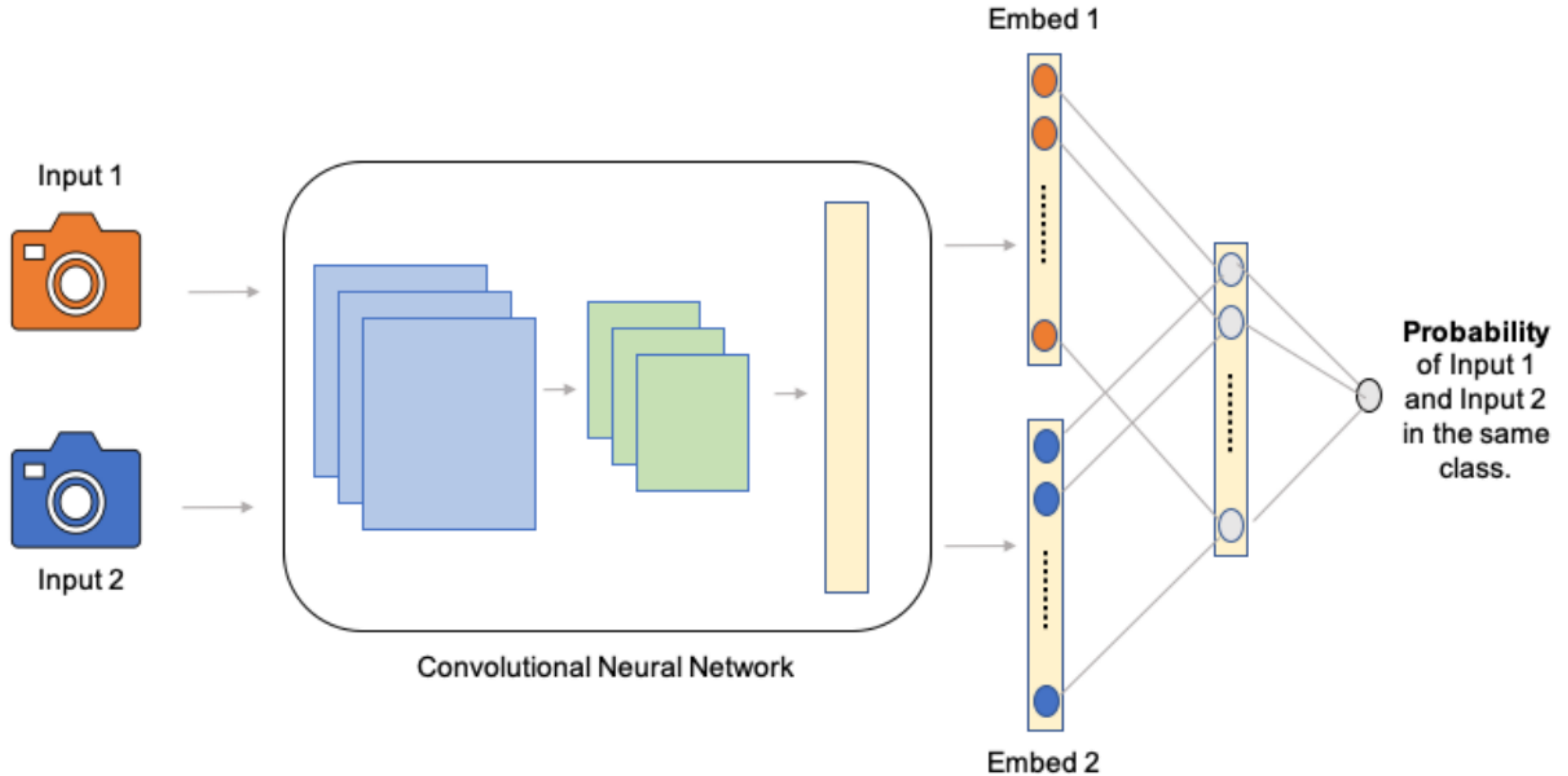
M-Way K-Shot Task (4-way-1-shot classification task)



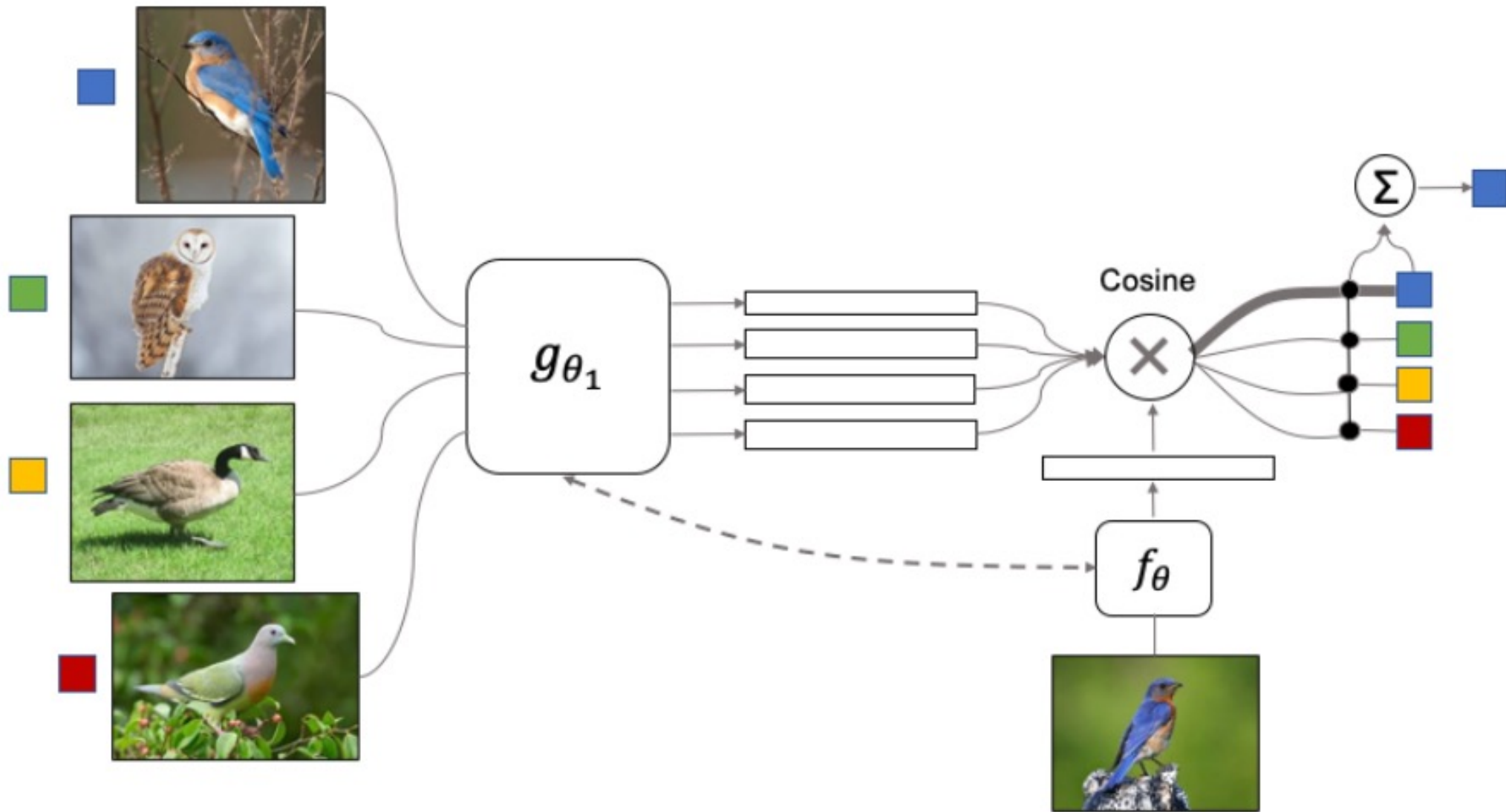
Metric-based Meta-Learning Methods

Method	T.I	g_{θ_1}	d_{θ_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	$\hat{y} = \sum_{k=1}^t \sigma(d(f_{\theta}(\hat{x}), g_{\theta_1}(x_k))) y_k$ $P(y = c \hat{x}) = \hat{y}_c$	$-\log P$
Prototypical Networks [21]	Yes	CNN	Euclidean	$P(y = c x) = \sigma(-d(g_{\theta_1}(x), \mathbf{v}_c))$	$-\log P$
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	$P_{\lambda}(y = c x) = \sigma(-\lambda d(g_{\theta_1}(x, \Gamma), \mathbf{v}_c))$	$-\log P$
TapNet [23]	No	Resnet-12	Euclidean	$P(y = c x) = \sigma(-d(\mathbf{M}(g_{\theta_1}(x)), \mathbf{M}(\Phi_c)))$	$-\log P$
CTM [24]	No	Any	Any	-	-

Convolutional Siamese Network

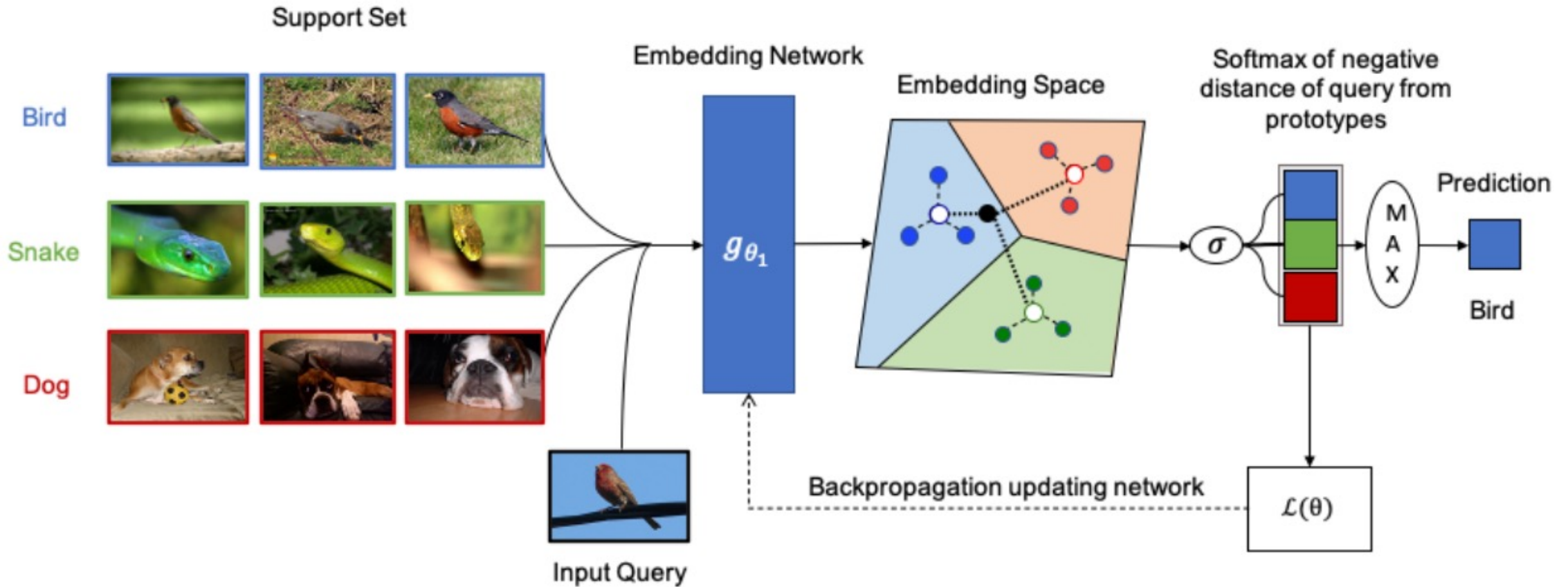


Meta Learning: Matching Networks

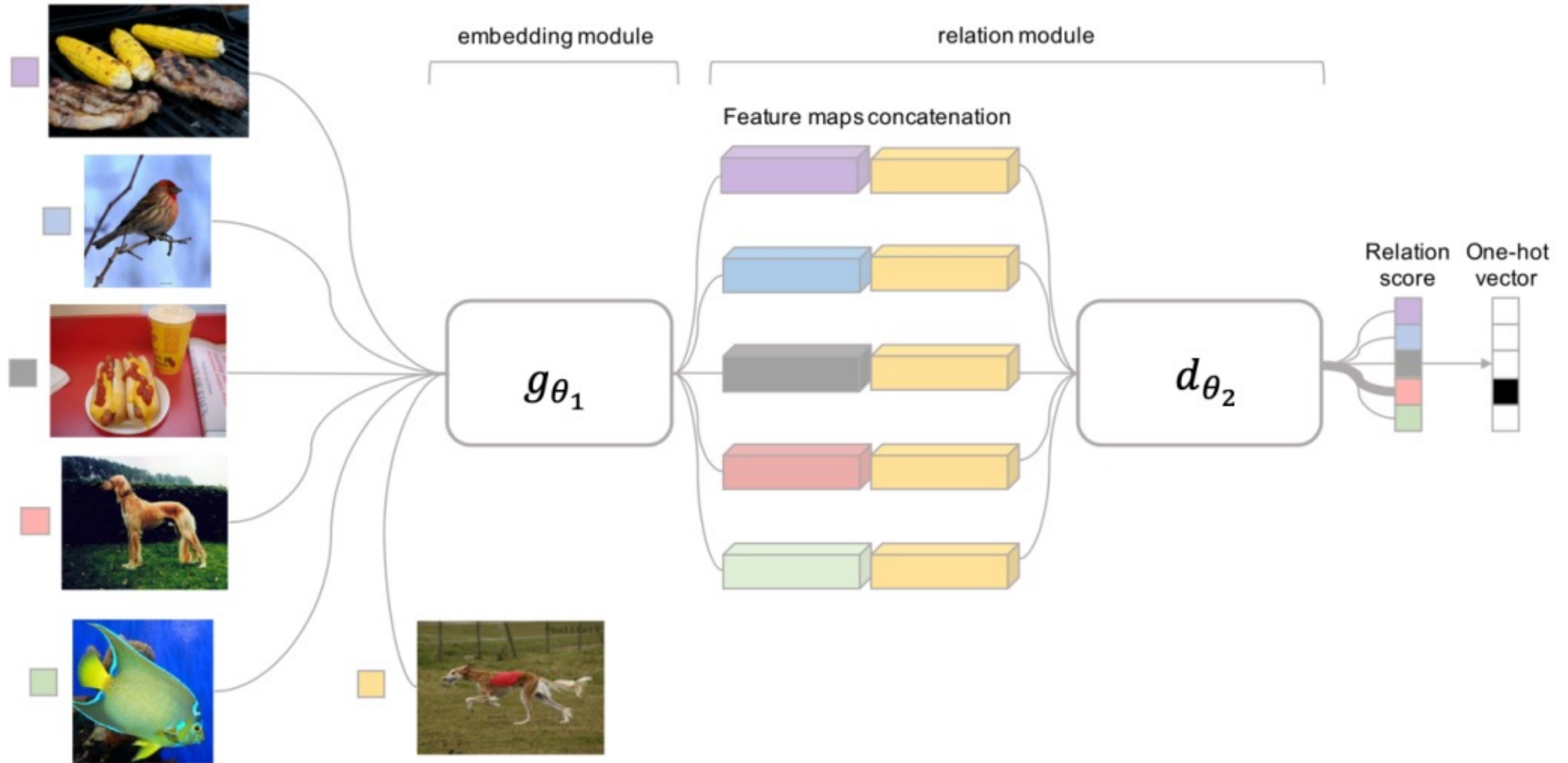


Few-shot Prototypes

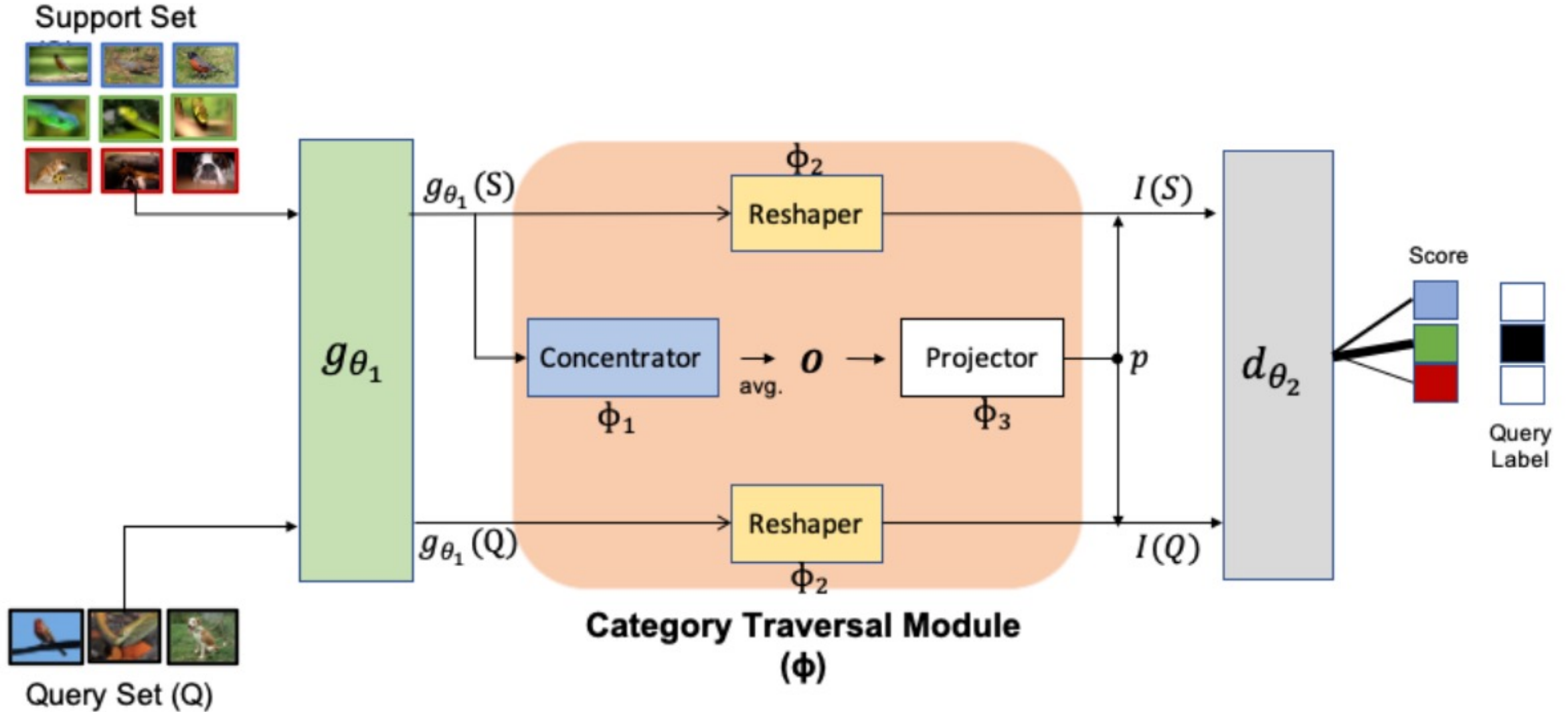
v_c are computed as the mean of embedded support examples for each class



Meta Learning: Relation Network



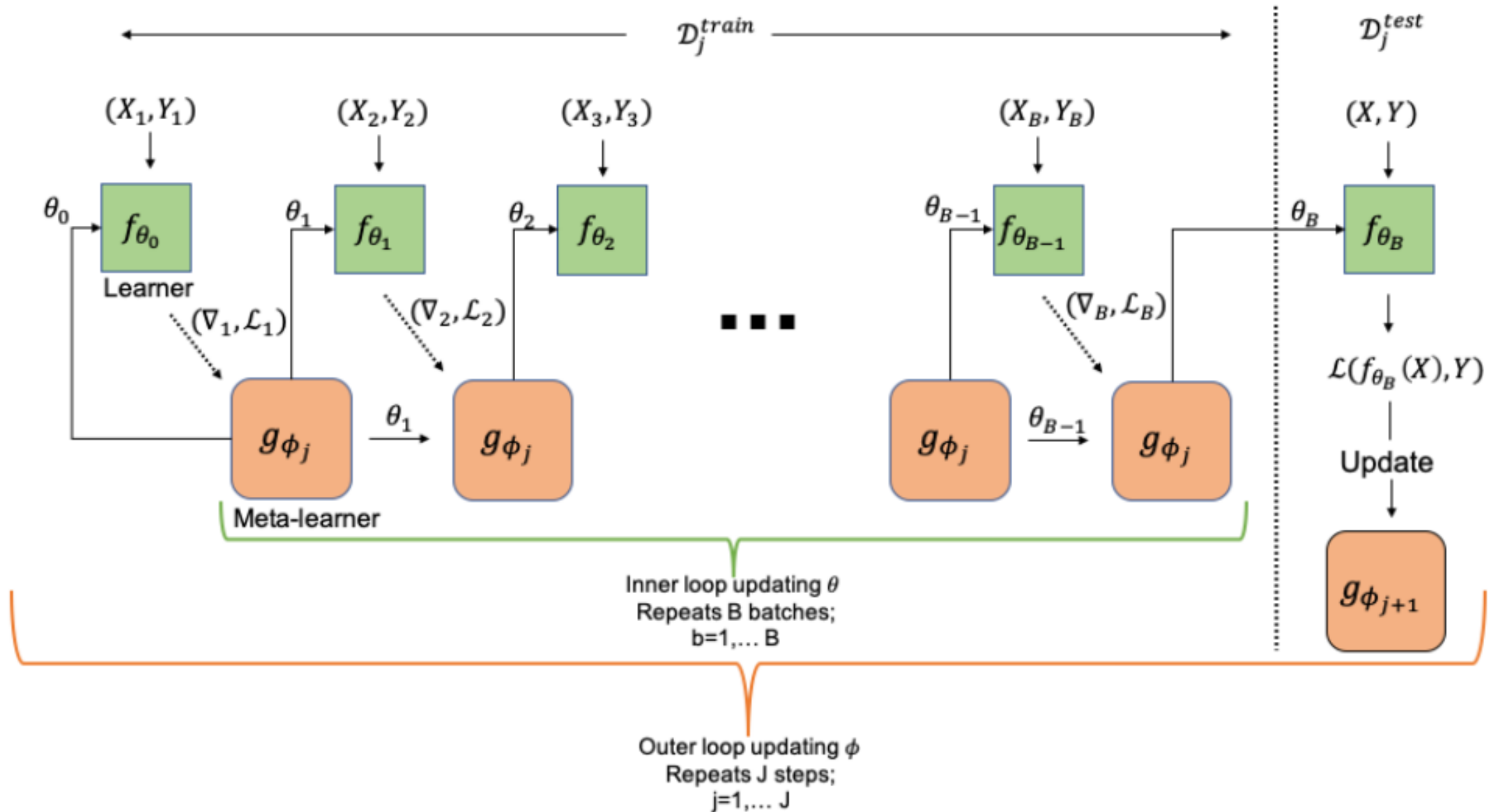
Meta Learning: Category Traversal Module (CTM)



Optimization-based Meta-Learning Methods

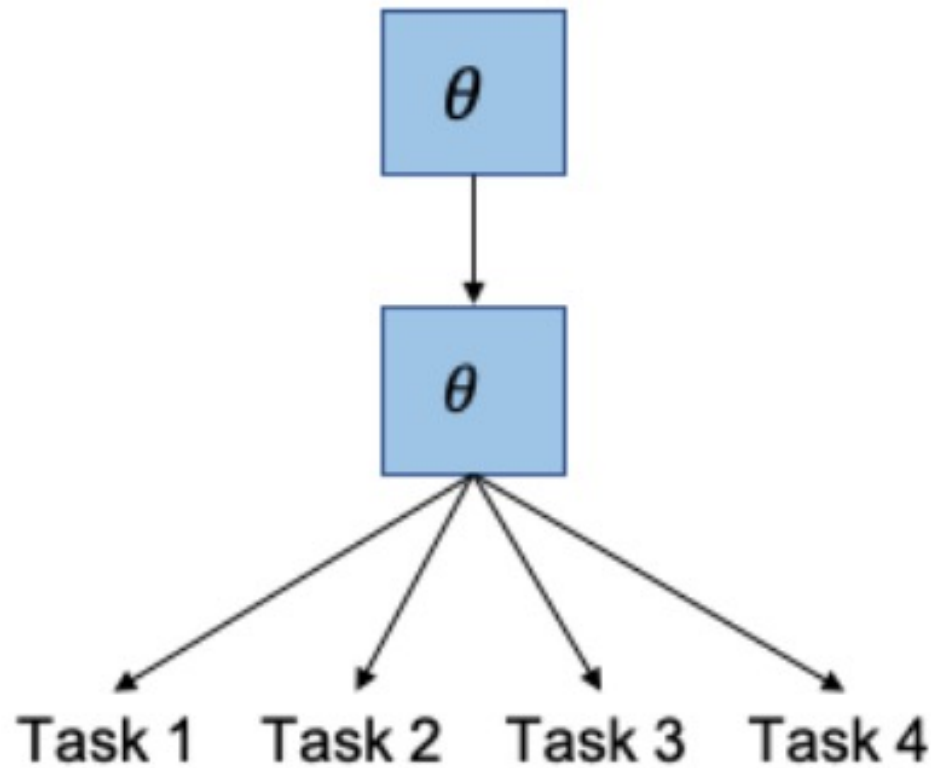
Method	Learner	Meta-Learner
LSTM Meta-Learner [35]	Repeat $\forall b \in [1..B]$ $\mathcal{L}_b \leftarrow \mathcal{L}(f(X_b; \theta_{b-1}), Y_b)$ $\theta_b \leftarrow g((\nabla_{\theta_{b-1}} \mathcal{L}_b, \mathcal{L}_b); \phi_{j-1})$	Repeat $\forall j \in [1..J]$ $\mathcal{L}_j^{test} \leftarrow \mathcal{L}(f(X; \theta_B), Y)$ $\phi_j \leftarrow \phi_{j-1} - \alpha \nabla_{\phi_{j-1}} \mathcal{L}_j^{test}$
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

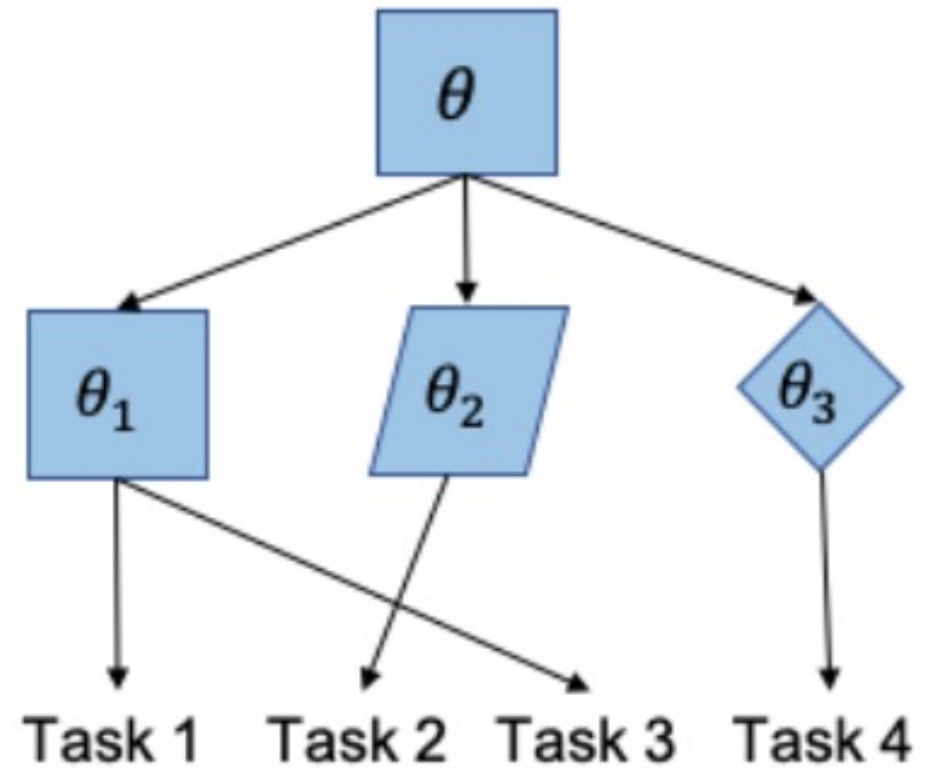


Model-Agnostic Meta-Learning (MAML)

Hierarchically Structured Meta-Learning (HSML)

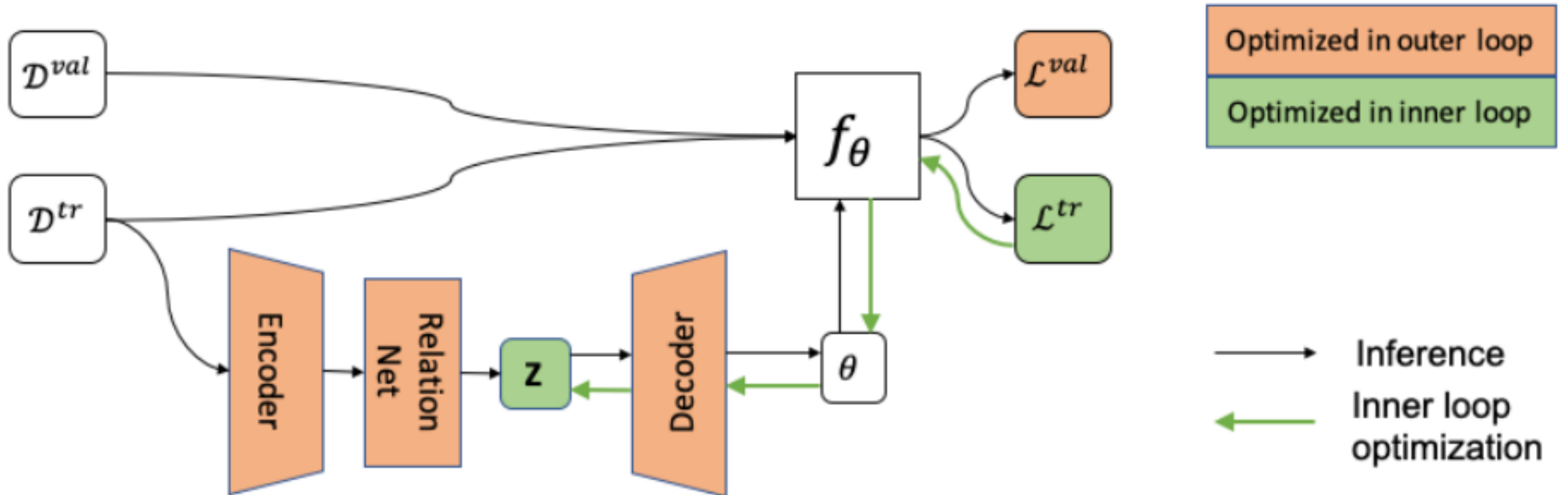


Globally Shared Initialization
(MAML)

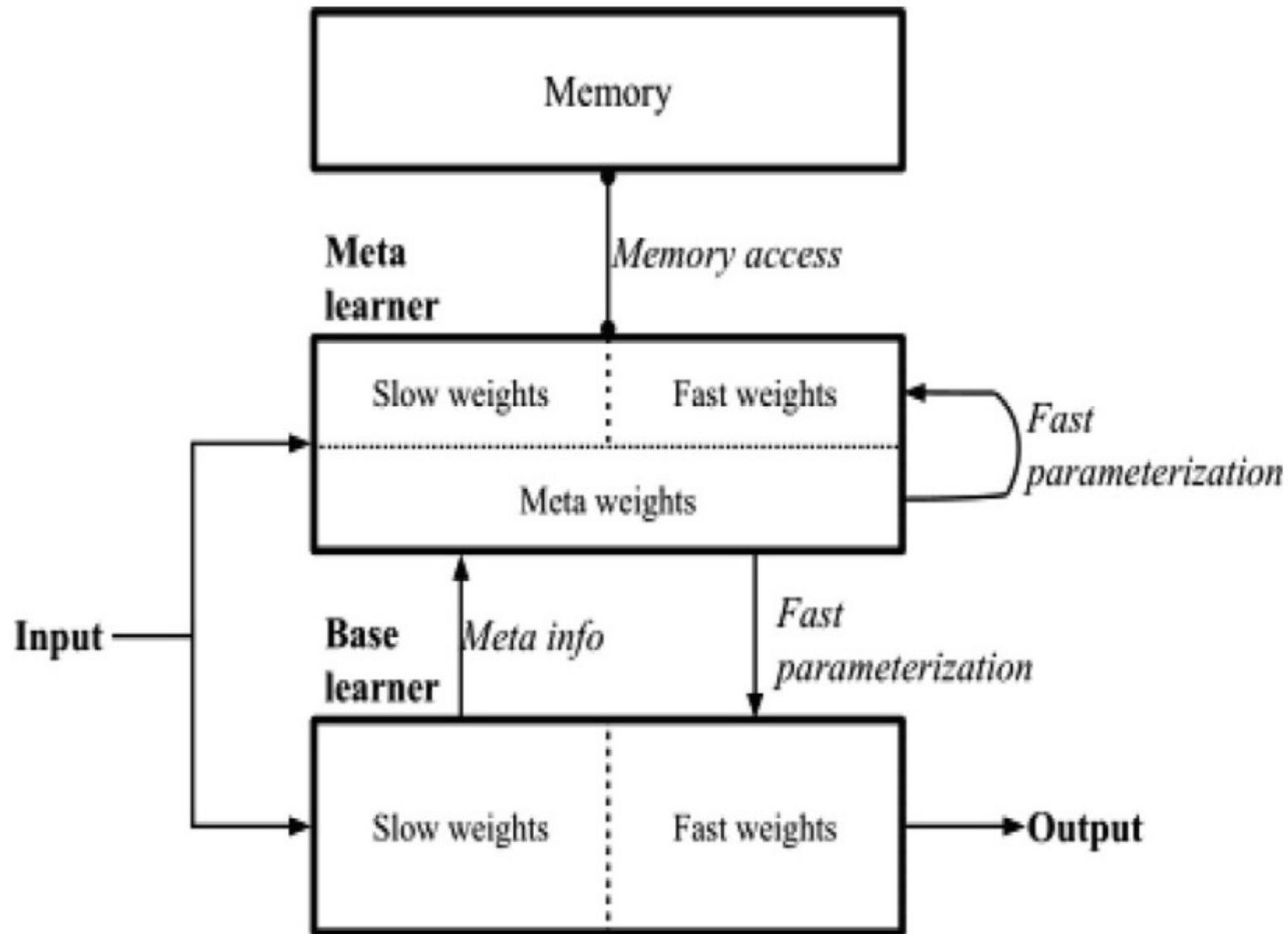


Hierarchically Clustered Initialization
(HSML)

Latent Embedding Optimization (LEO)

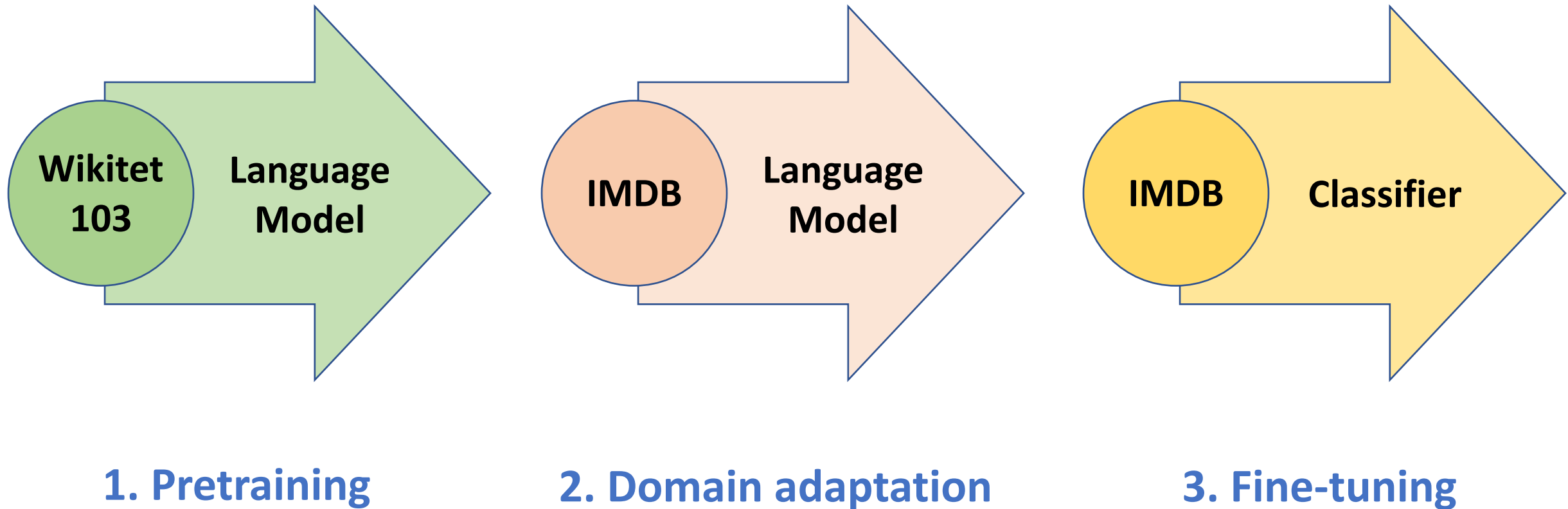


Overall Architecture of Meta Networks



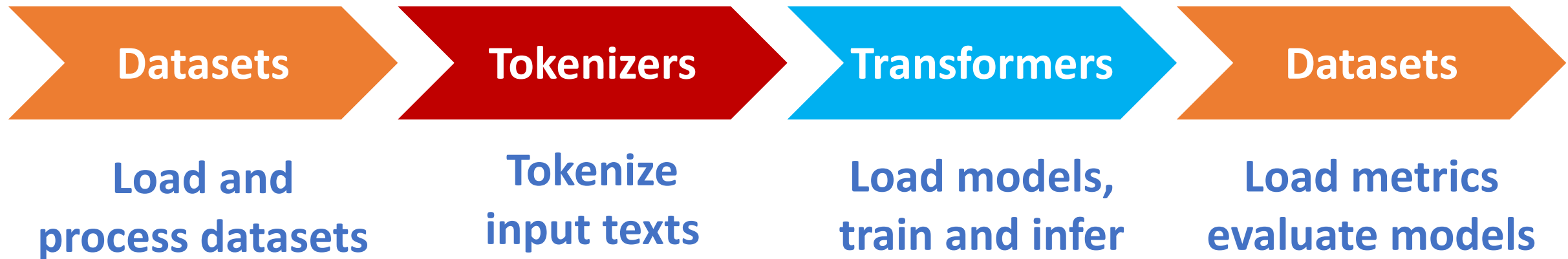
ULMFiT: 3 Steps

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)

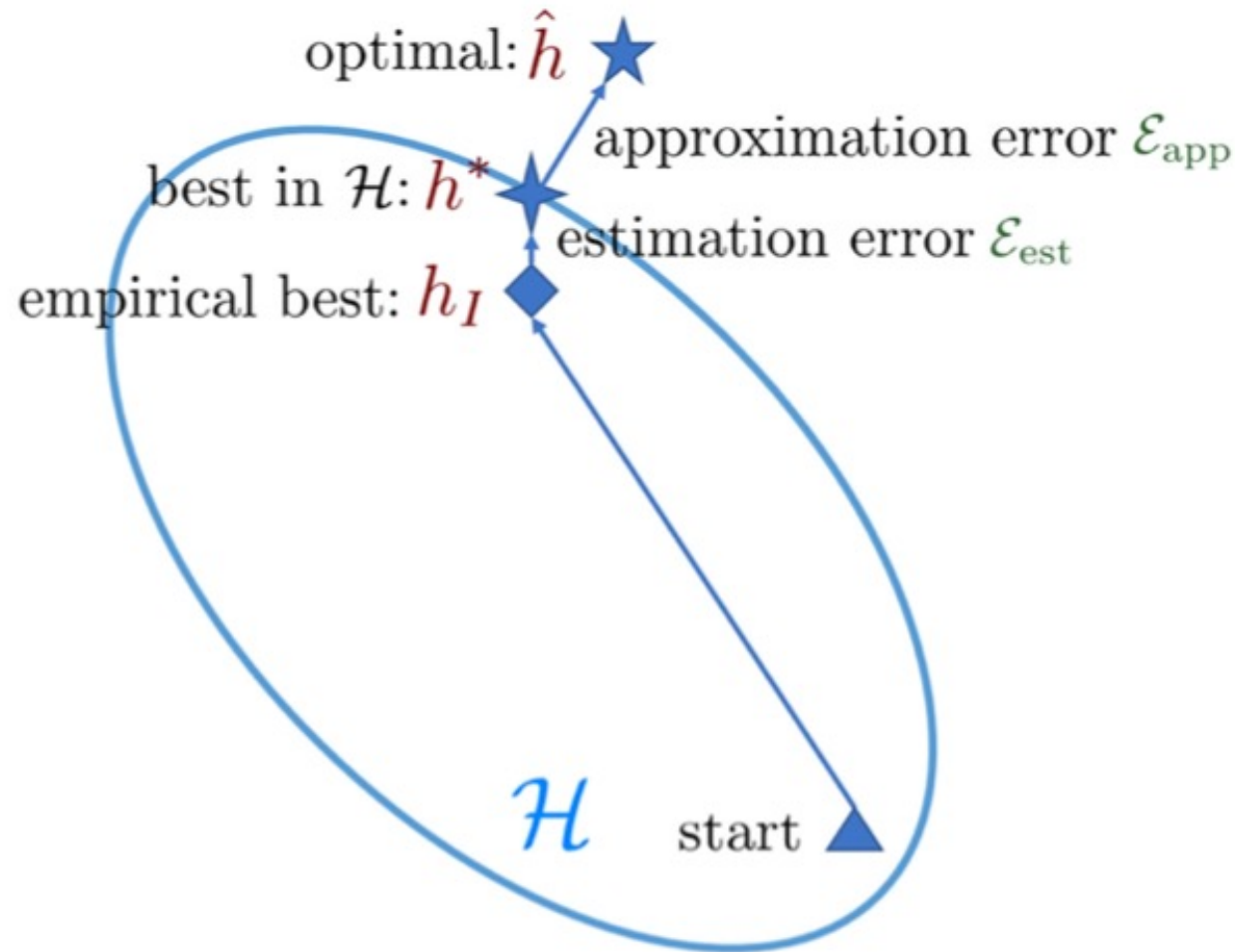
- **Few-Shot Learning (FSL)**
 - $K = 10 \sim 100$ examples
- **One-Shot Learning (1SL)**
 - $K = 1$ example
- **Zero-Shot Learning (OSL)(ZSL)**
 - $K = 0$

Few-Shot Learning (FSL)

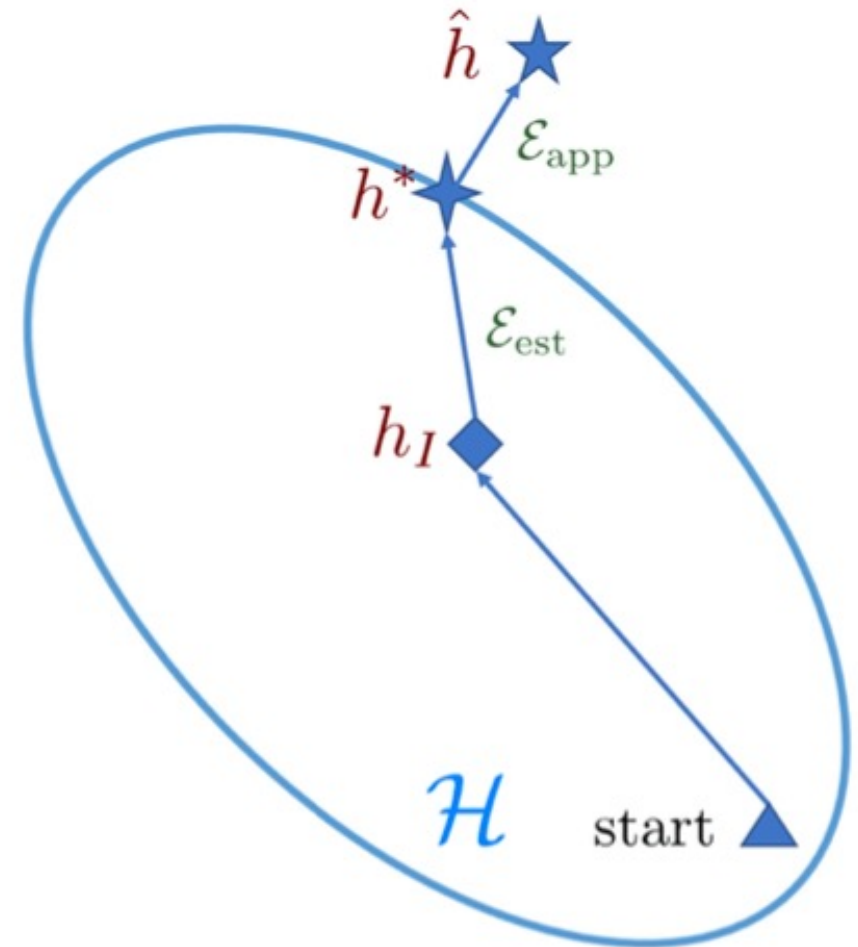
task T	experience E		performance P
	supervised information	prior knowledge	
character generation [76]	a few examples of new character	pre-learned knowledge of parts and relations	pass rate of visual Turing test
drug toxicity discovery [4]	new molecule's limited assay	similar molecules' assays	classification 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



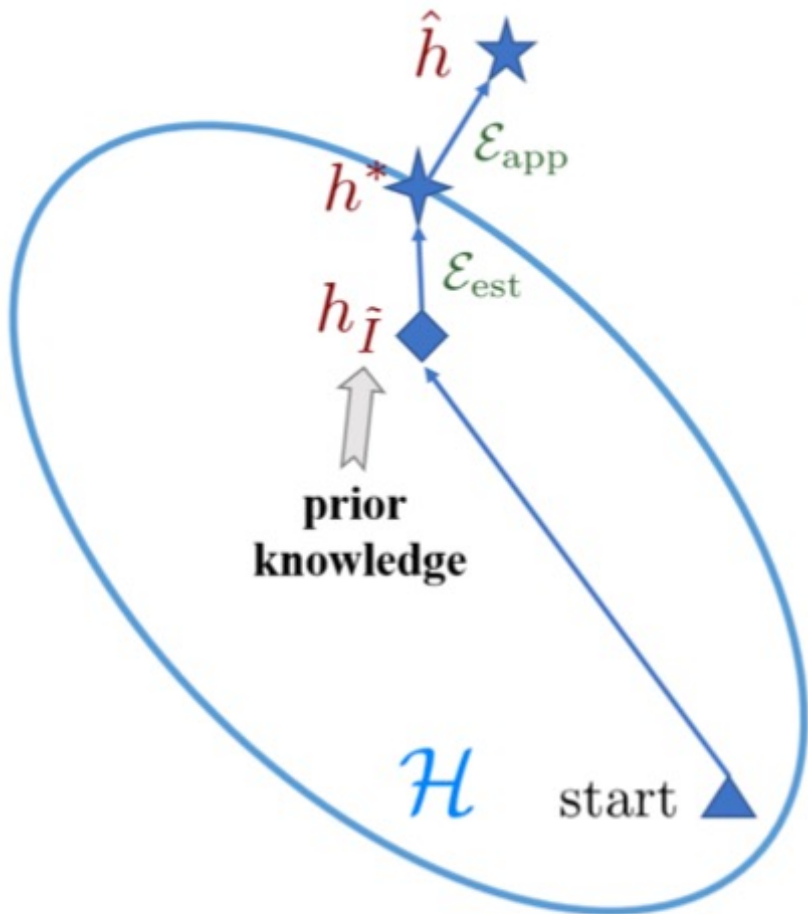
(a) Large I



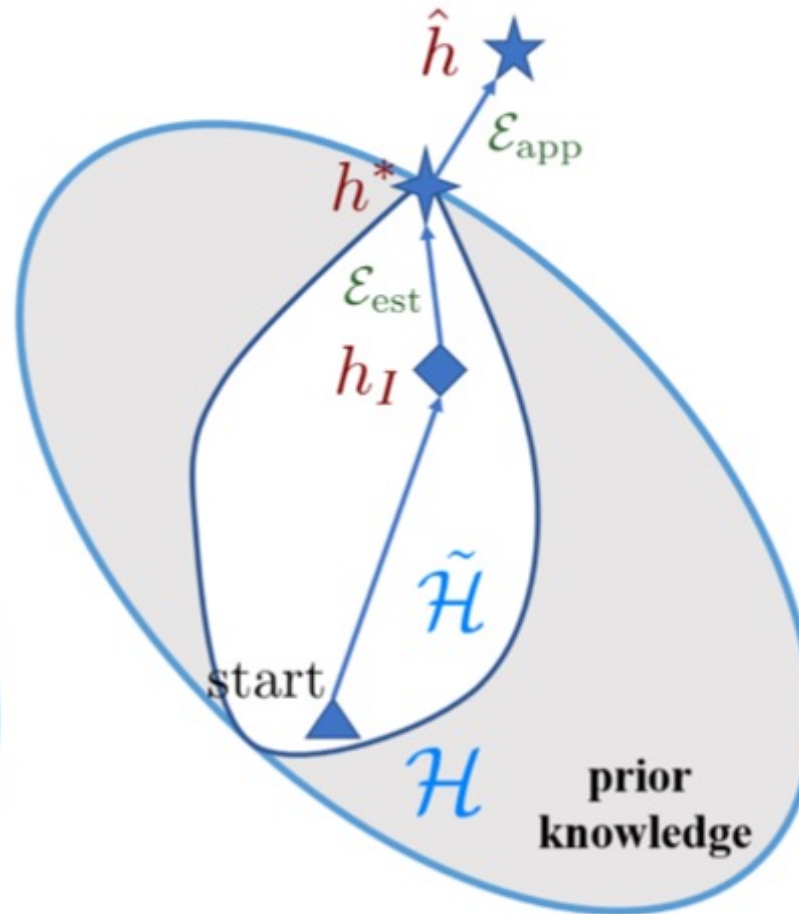
(b) Small I

Few-Shot Learning (FSL)

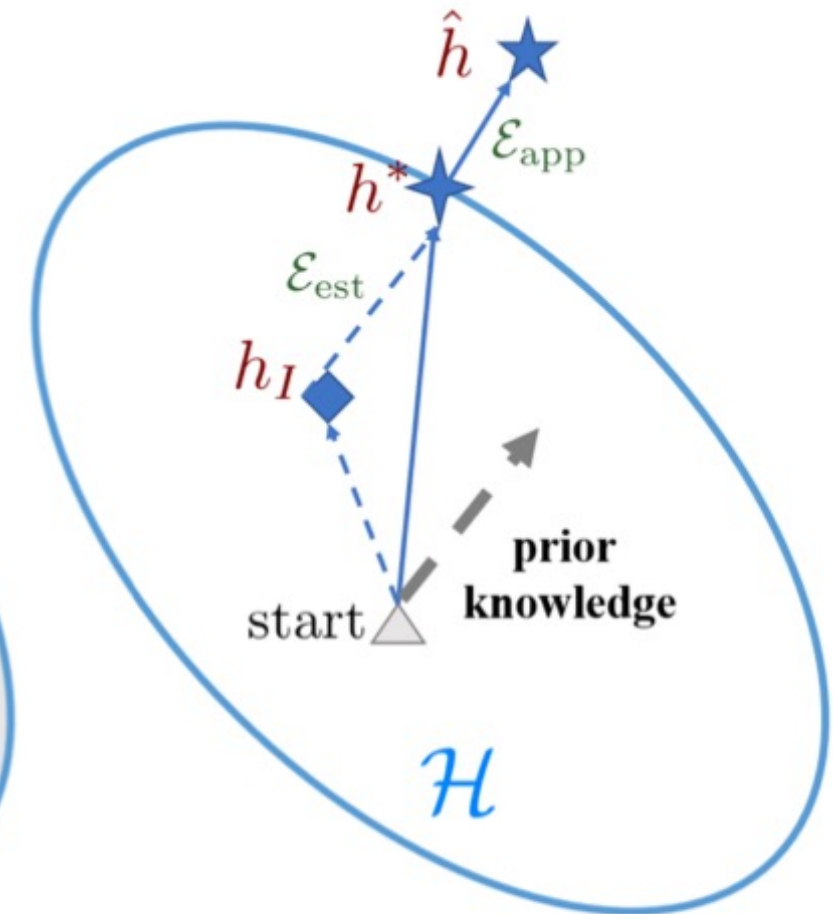
Different perspectives on how FSL methods solve the few-shot problem



(a) Data.



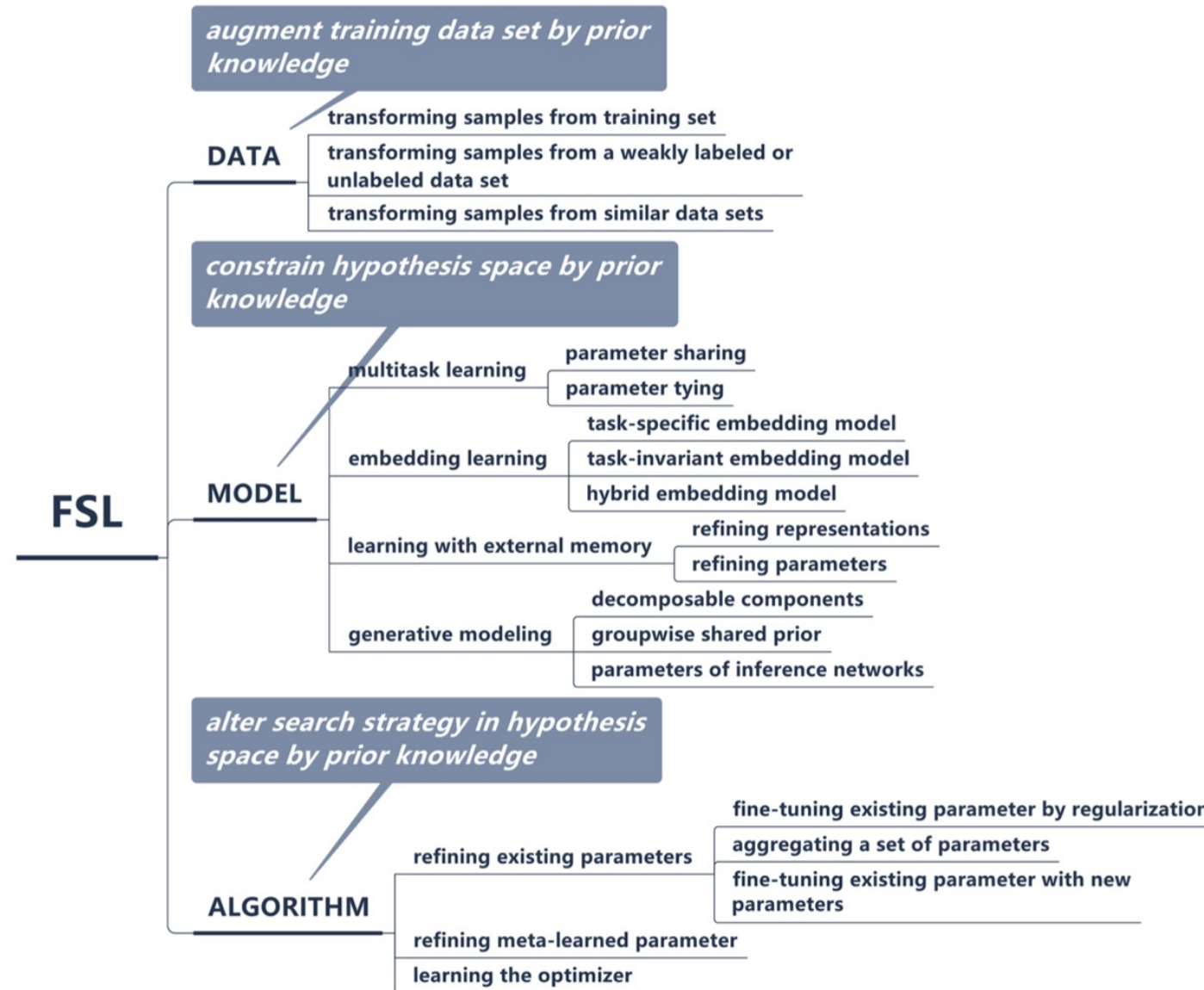
(b) Model.



(c) Algorithm.

Few-Shot Learning (FSL)

A taxonomy of FSL methods



Few-Shot Learning (FSL)

augment training data set by prior knowledge

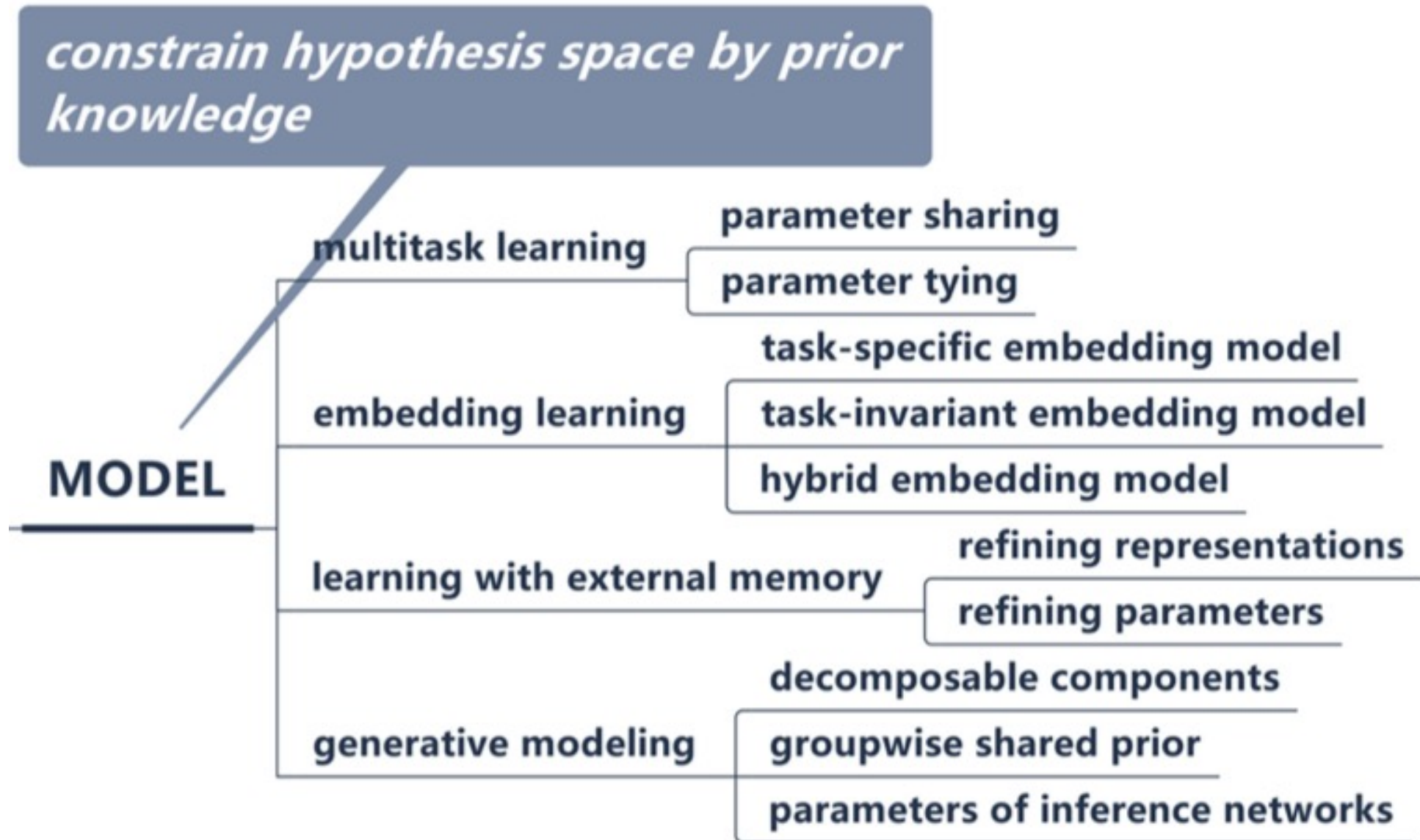
DATA

transforming samples from training set

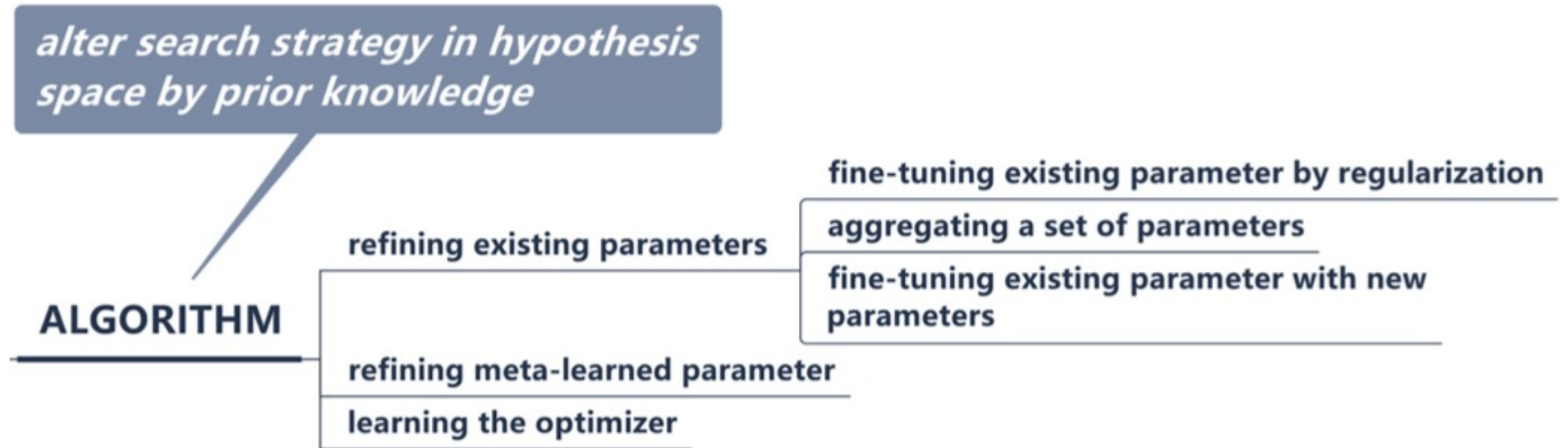
transforming samples from a weakly labeled or unlabeled data set

transforming samples from similar data sets

Few-Shot Learning (FSL)

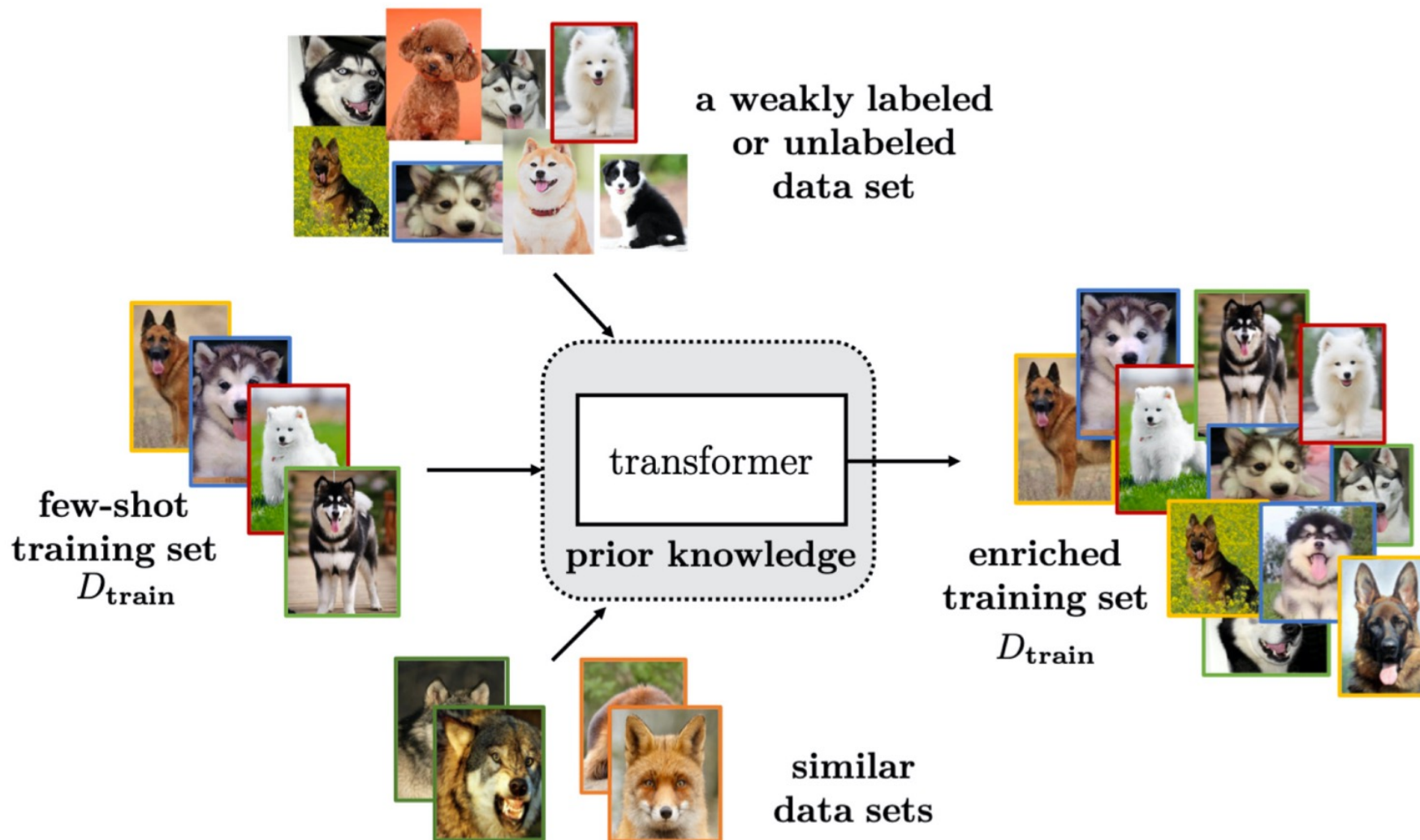


Few-Shot Learning (FSL)



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_{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_{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_{train} .

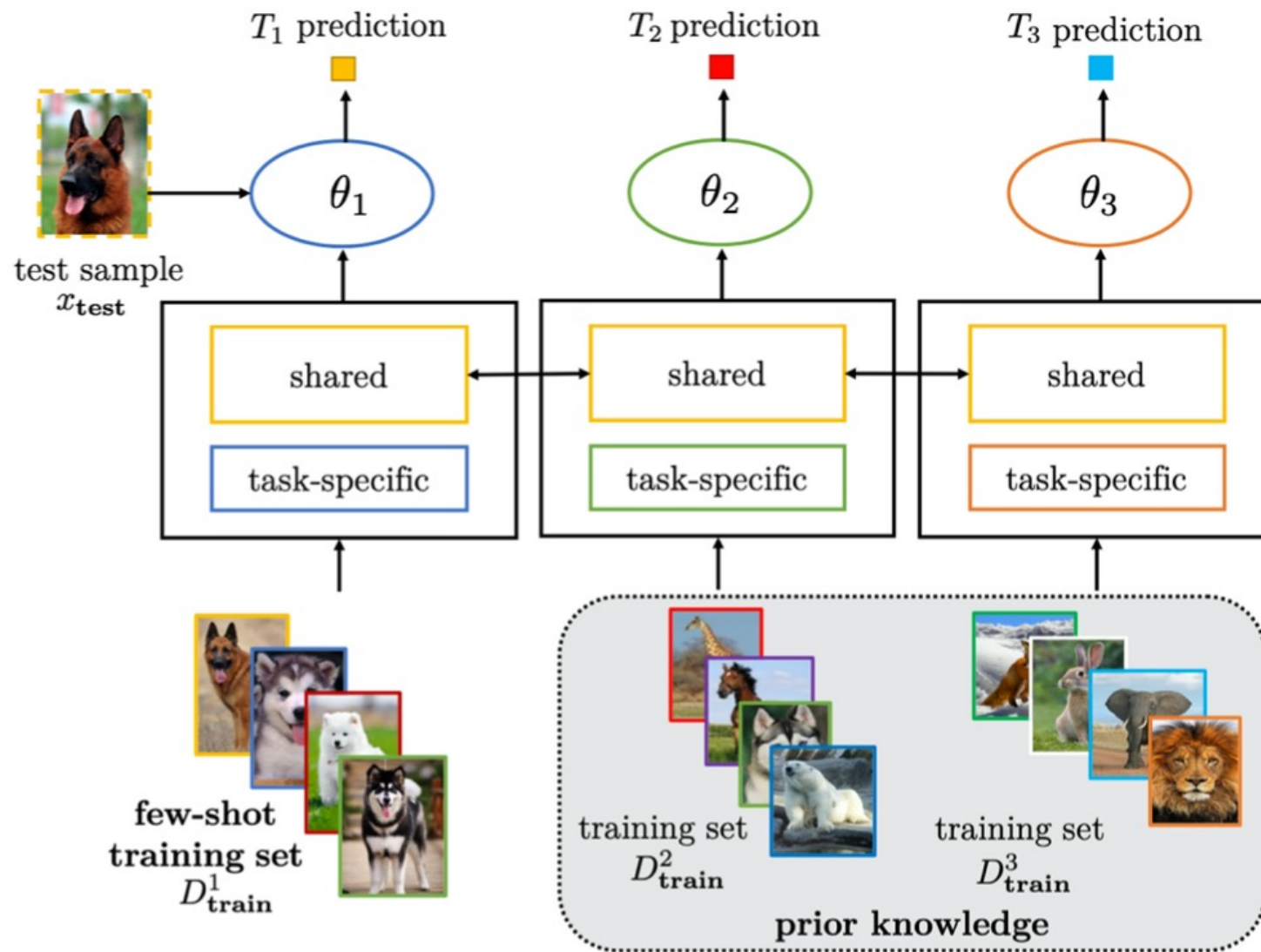
Few-Shot Learning (FSL)

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

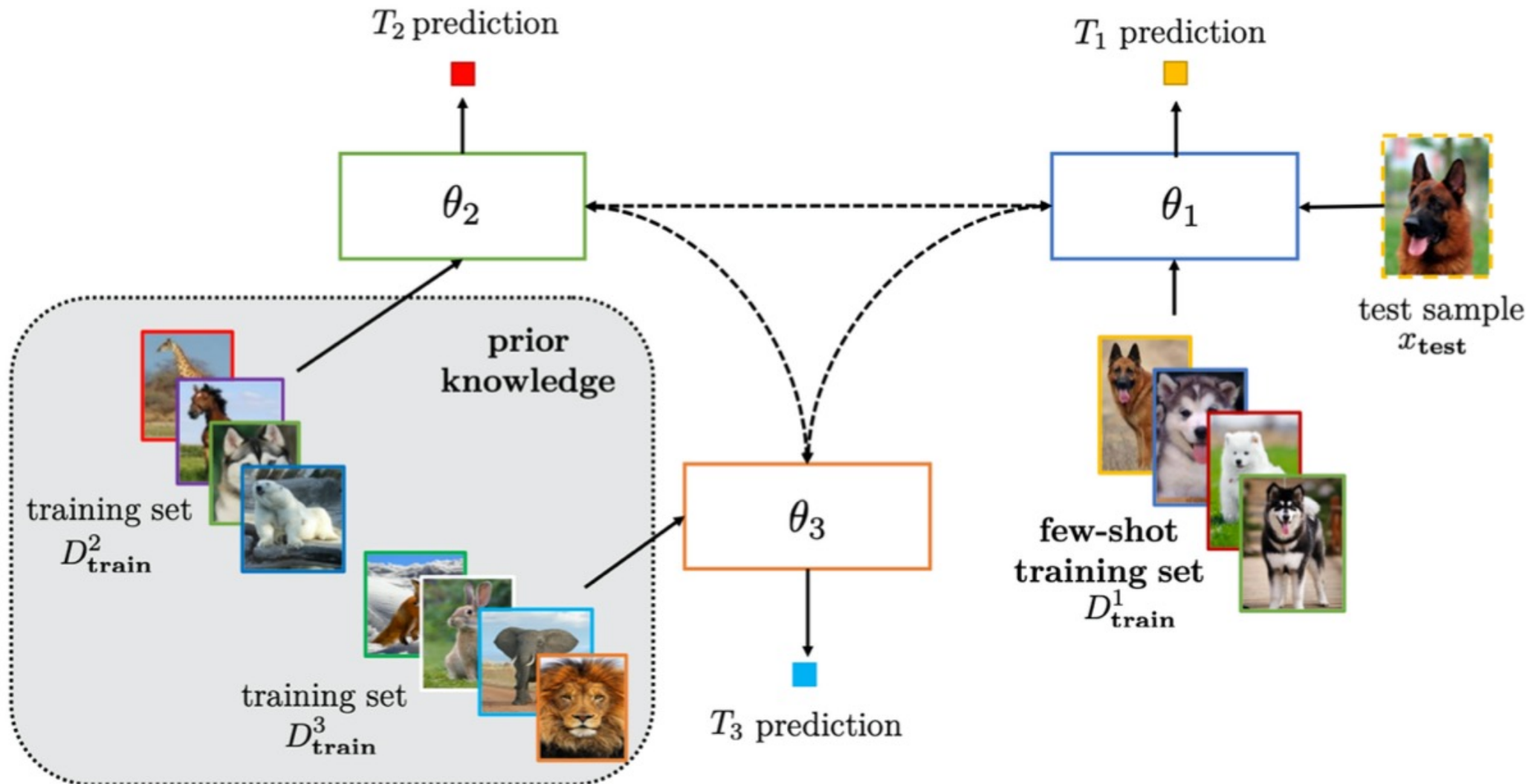
Few-Shot Learning (FSL)

Solving the FSL problem by multitask learning with parameter sharing



Few-Shot Learning (FSL)

Solving the FSL problem by multitask learning with parameter tying



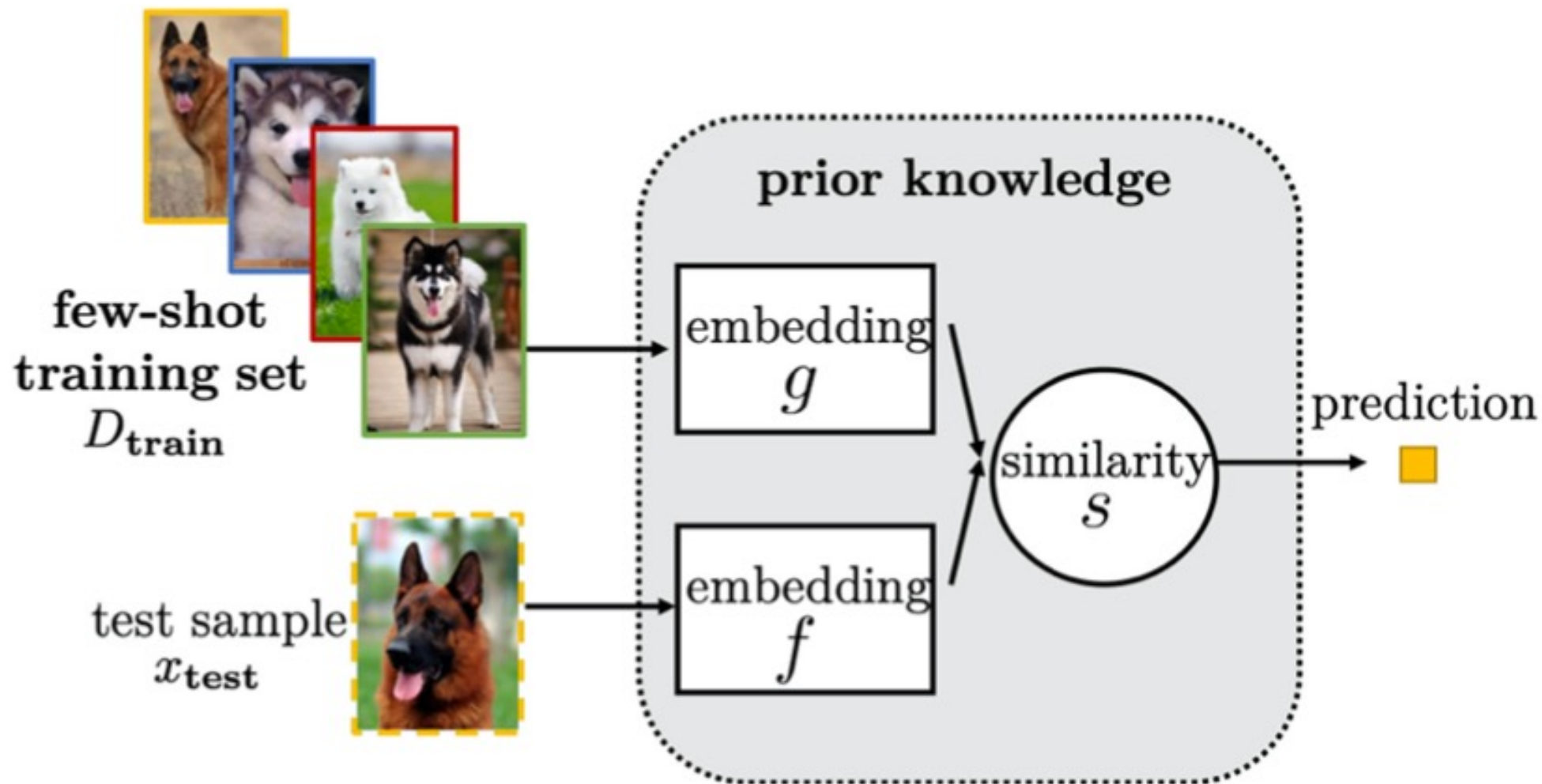
Few-Shot Learning (FSL)

Characteristics of Embedding Learning Methods

category	method	embedding function f for x_{test}	embedding function g for D_{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 ℓ_2 distance
	convolutional siamese net [70]	CNN	the same as f	weighted ℓ_1 distance
	Micro-Set[127]	logistic projection	the same as f	ℓ_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 ℓ_2 distance
	semi-supervised ProtoNet[108]	CNN	the same as f	squared ℓ_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 ℓ_1 distance
	DCCN [162]	adaptive CNN	CNN	-
	R2-D2 [13]	adaptive CNN	CNN	-
	TADAM [100]	adaptive CNN	the same as f	squared ℓ_2 distance

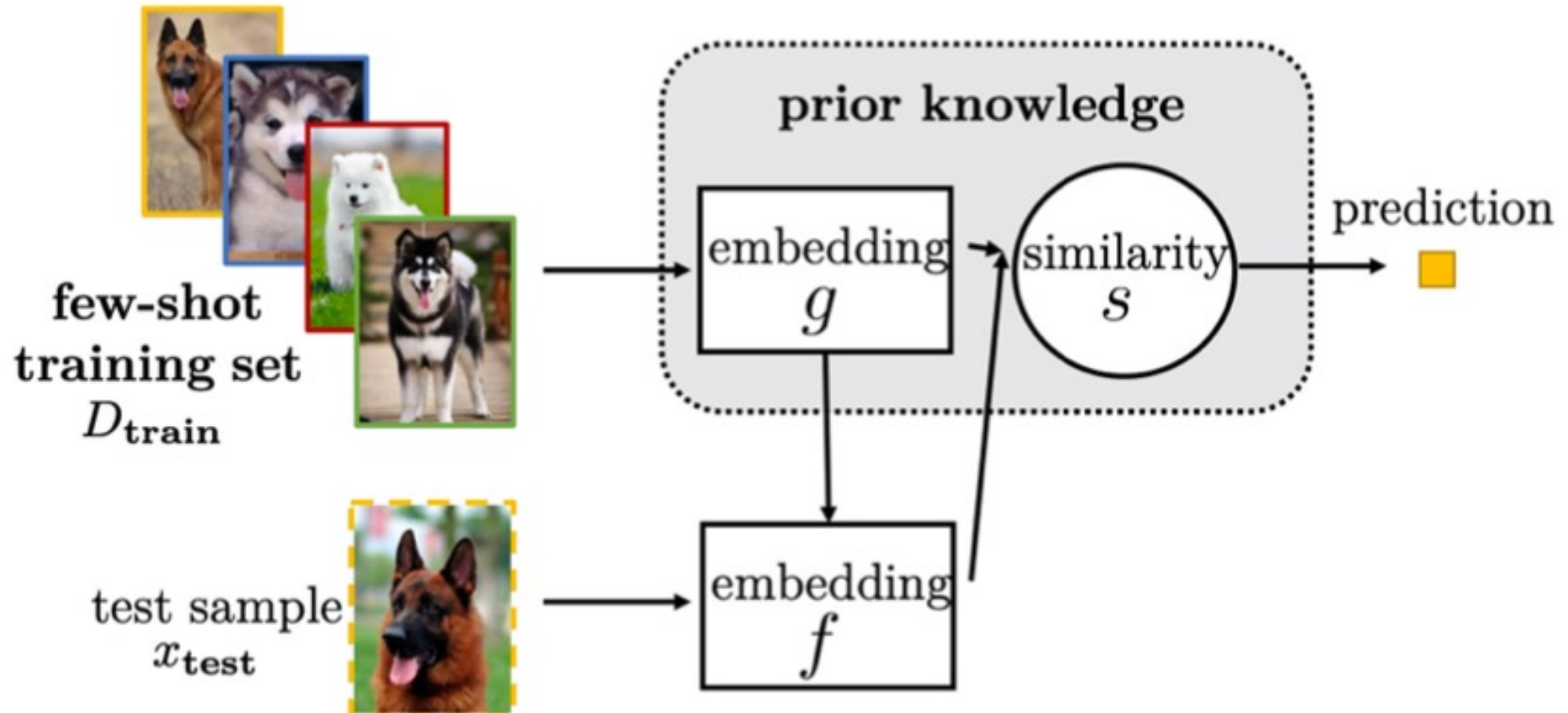
Few-Shot Learning (FSL)

Solving the FSL problem by task-invariant embedding model



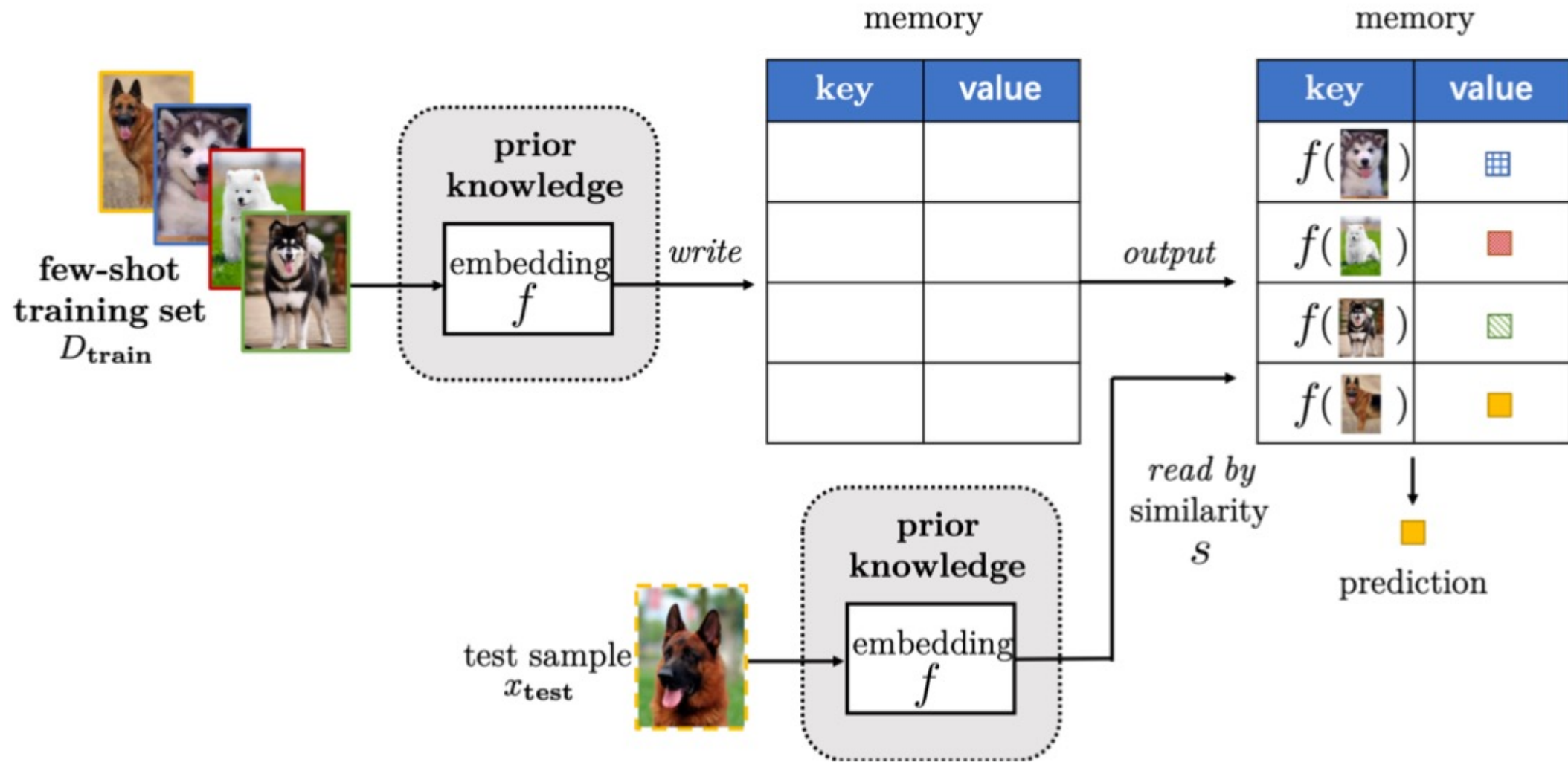
Few-Shot Learning (FSL)

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

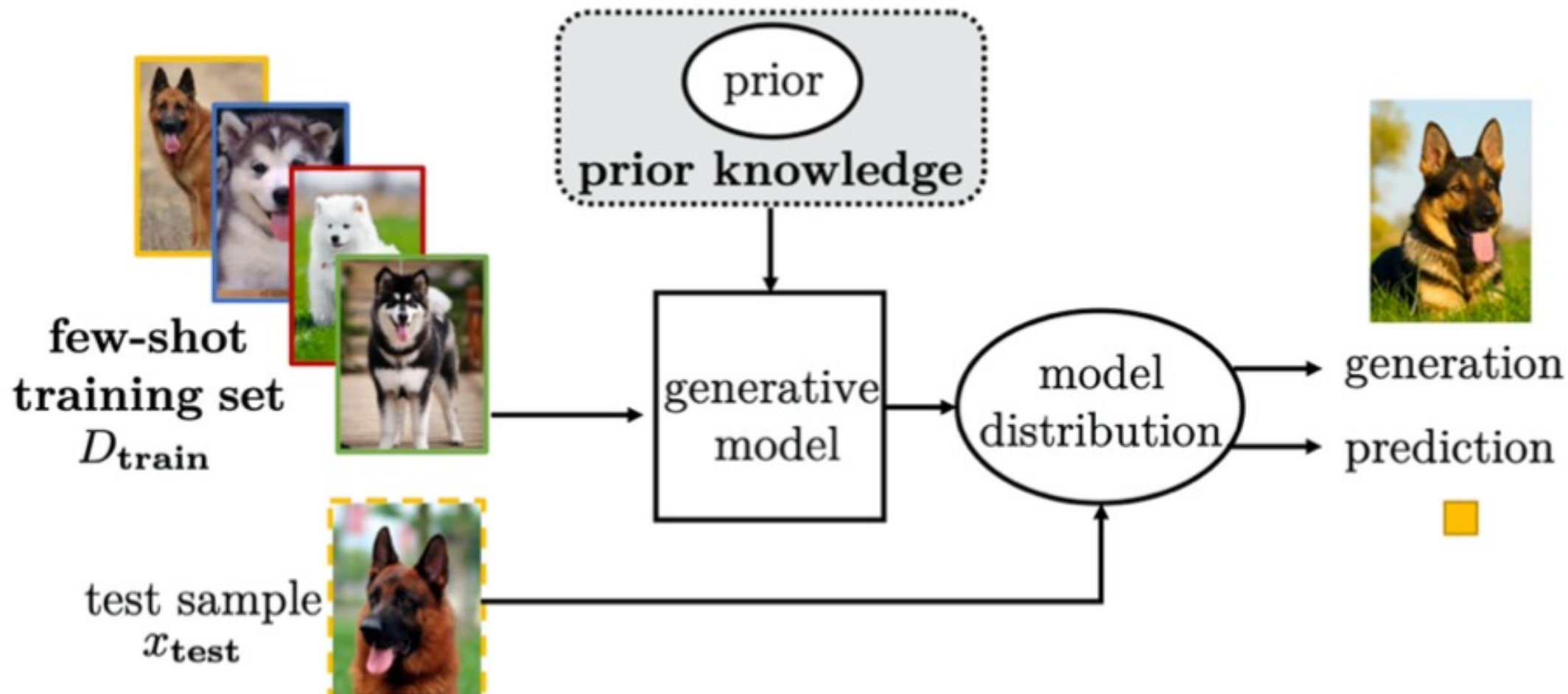
Based on Learning with External Memory

category	method	memory M		similarity s
		key M_{key}	value M_{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 ℓ_2 distance
	abstraction memory [149]	$f(x_i)$	word embedding of y_i	dot product
	CMN [164]	$f(x_i)$	y_i, age	dot product
	life-long memory [65]	$f(x_i)$	y_i, age	cosine similarity
	Mem2Vec [125]	$f(x_i)$	word embedding of y_i, 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.

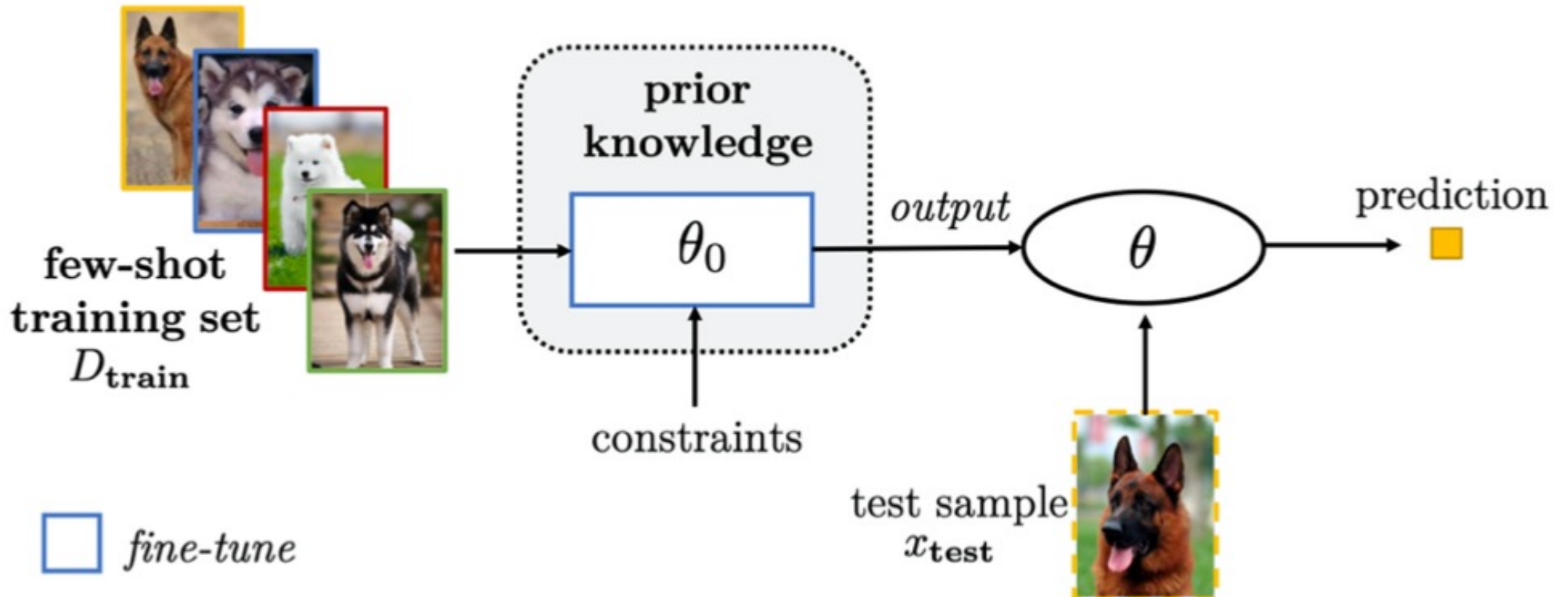
Few-Shot Learning (FSL)

Solving the FSL problem by generative modeling



Few-Shot Learning (FSL)

Solving the FSL problem by fine-tuning existing parameter θ_0 by regularization



Few-Shot Learning (FSL)

Characteristics for FSL Methods

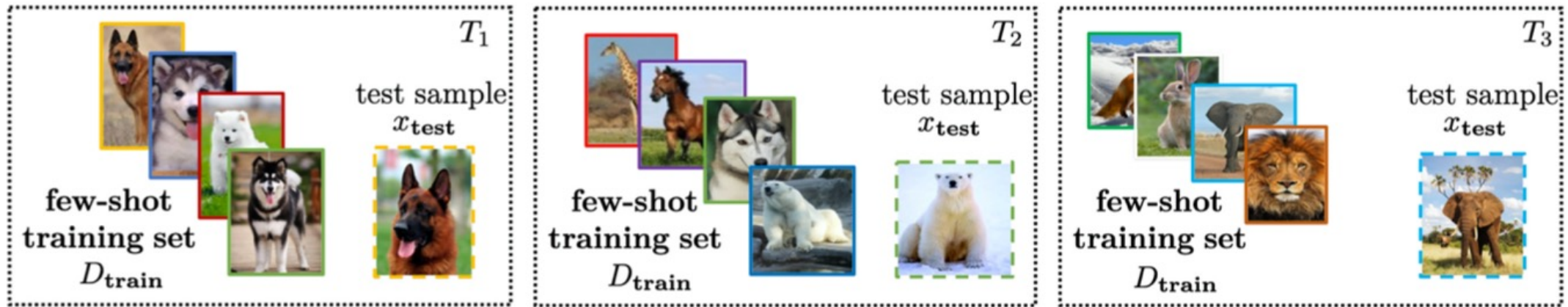
Focusing on the Algorithm Perspective

strategy	prior knowledge	how to search θ of the h^* in \mathcal{H}
refining existing parameters	learned θ_0	refine θ_0 by D_{train}
refining meta-learned parameters	meta-learner	refine θ_0 by D_{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)

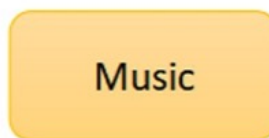
Meta-learning

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

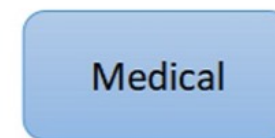
Training Task 1



Training Task 2 ...



Test Task 1 ...



Support Set:

Only_[O] France_[LOC] and_[O] Britain_[LOC]
backed_[O] Fischler_[PER]'s_[O] proposal_[O].

Query Set:

Adrian_[PER] Warner_[PER] has_[O] lived_[O]
in_[O] Brussels_[LOC] since_[O] 1996_[O].

Labels: {PER, PER, O, O, O, LOC, O, O}

Support Set:

Play_[O] rap_[album] album_[album] one_[album]
by_[O] Gene_[artist] Vincent_[artist].

Query Set:

Add_[O] Rosemary_[artist] Clooney_[artist]
to_[O] pura_[playlist] vida_[playlist] playlist_[O].

Labels: {O, artist, artist, O, playlist,
playlist, O}

Support Set:

Patient_[O] received_[O] combivent_[DRUG]
nebs_[Dosage Form/Route], solumedrol_[DRUG]
125mg_[AMOUNT] IV_[Dosage Form/Route]
X1_[Dosage Frequency].

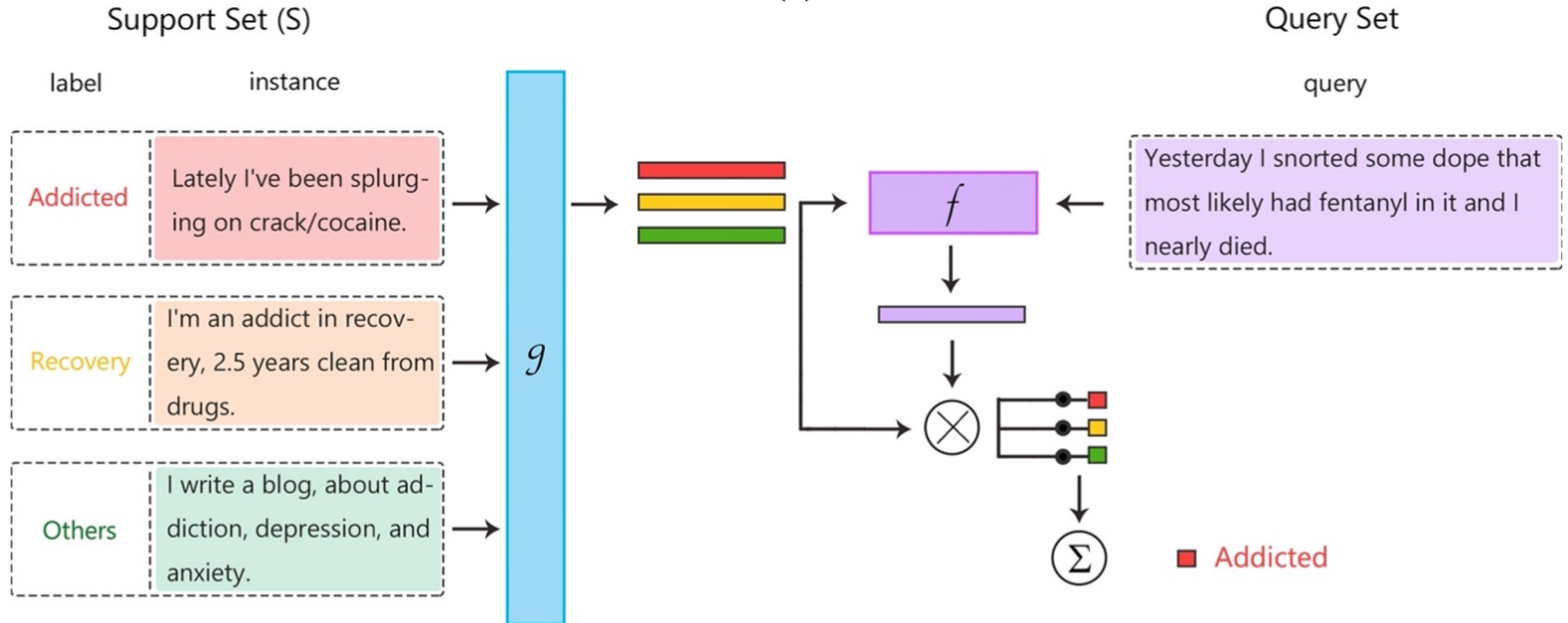
Query Set:

She was given a dose of levaquin this
morning.

Labels: {O, O, O, AMOUNT, AMOUNT,
O, DRUG, TIME, TIME}

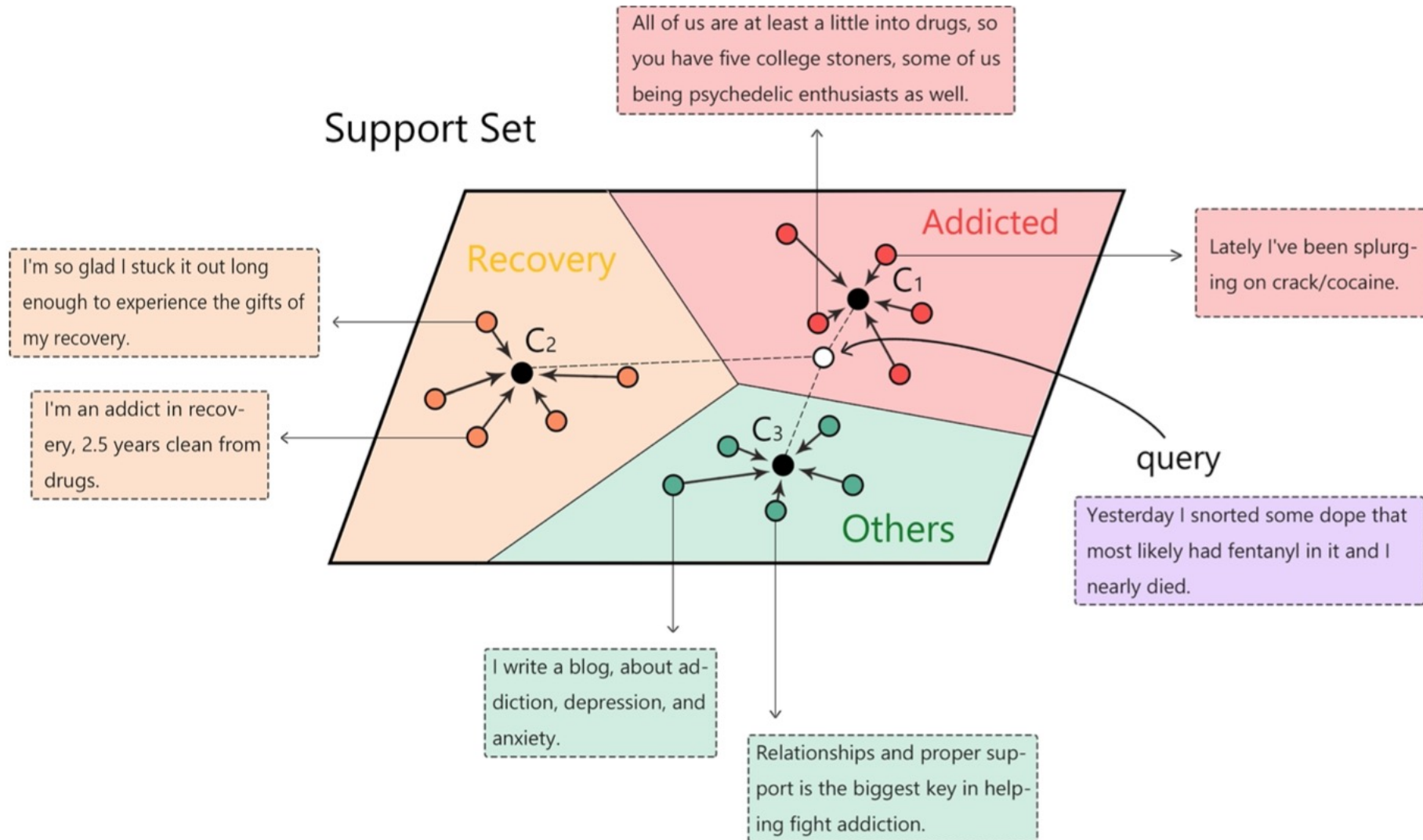
Few-Shot Learning (FSL)

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 ⁴⁴	2019	MEDDOCAN shared task dataset ⁴⁵	NER	500 clinical cases, with no reconstruction	29	Clinical	Clinical
Ferré et al. ⁴⁶	2019	BB-norm dataset from the Bacteria Biotope 2019 Task ⁴⁷	Entity Normalization	Original dataset with no reconstruction and zero-shot	Not mentioned *	Biological	Biological
Hou et al. ⁴⁸	2020	Snips dataset ⁴⁹	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. ⁵⁰	2020	ten different datasets collected from the Open Parallel Corpus (OPUS) ⁵¹	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. ⁵²	2020	MIMIC II ²² and MIMIC III ²³ , and EU legislation dataset ⁵³	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. ⁸⁰	2021	Constructed and shared a novel dataset ^{††} based on Weibo for the research of few-shot rumor detection, and use PHEME dataset ⁸¹	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. ⁸²	2021	CCLE, CERES-correctedCRISPR gene disruption scores, GDSC1000 dataset, PDTC dataset and PDX dataset ^{††}	Drug-response Predictions	1-shot, 2-shot, 5-shot and 10-shot	N/A [†]	Biomedical	Biomedical
Kormilitzin et al. ⁸³	2021	MIMIC-III ²³ and UK-CRIS datasets ^{30,31}	NER	25%, 50%, 75% and 100% of the training set, with no reconstruction	7	Electronic health record	Electronic health record
Guo et al. ⁸⁴	2021	Abstracts of biomedical literatures (from relation extraction task of BioNLP Shared Task 2011 and 2019 ⁴⁷) 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. ⁸⁵	2021	COVID19-Scientific ⁸⁶ , COVID19-Social ⁸⁷ (fact-checked by journalists from a website called Politi-fact.com), FEVER ⁸⁸ (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

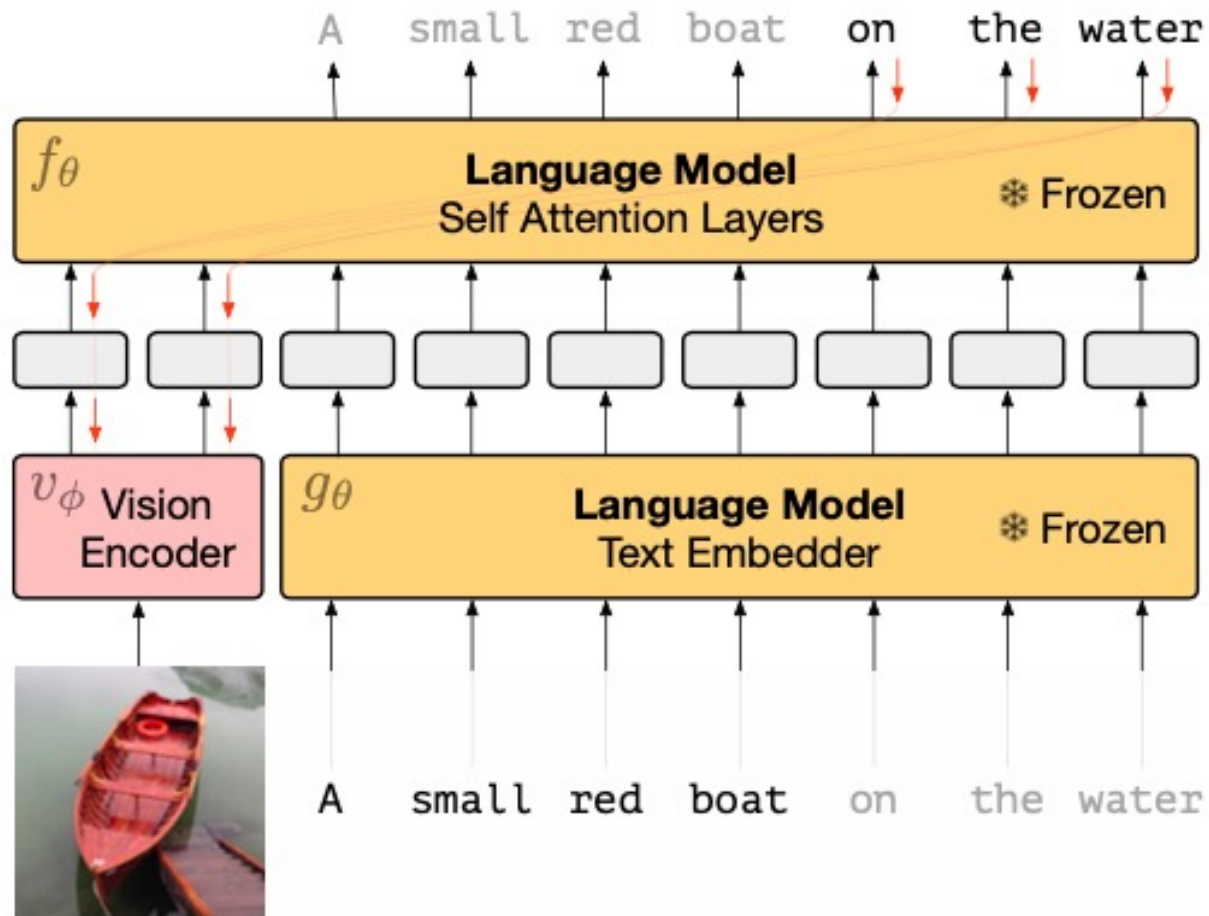
Multimodal Few-Shot Learning with Frozen Language Models



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.

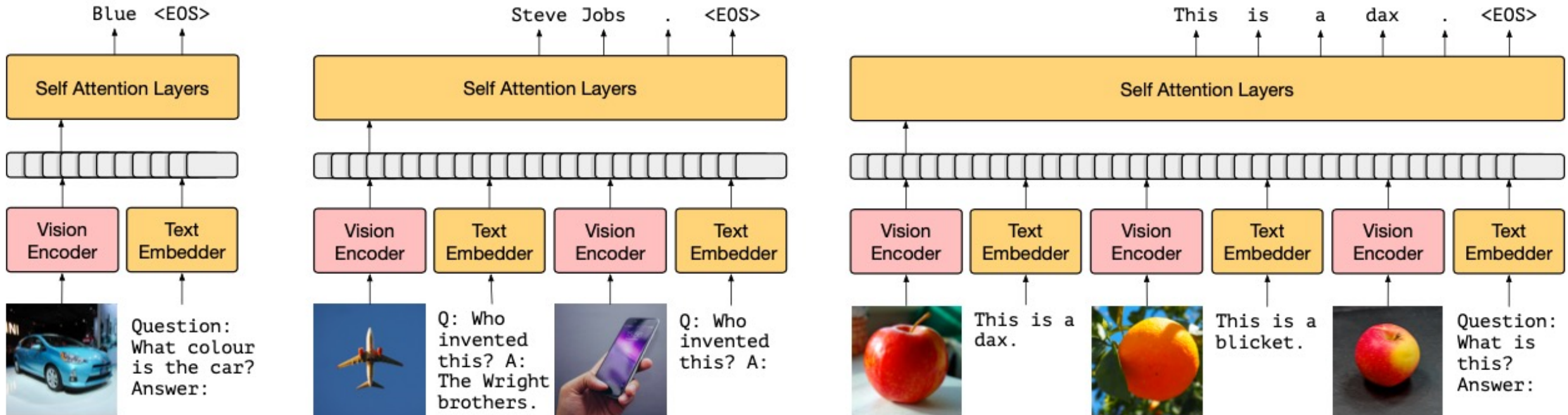
These samples demonstrate the ability to generate open-ended outputs that adapt to both images and text, and to make use of facts that it has learned during language-only pre-training.

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 Tsimgoukelli, 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) miniImageNet

0-repeats
0-shots
2-way
0-repeats
2-inner-shots

Task Induction

Answer with dax
or blicket.

**Support
from ImageNet**

inner-shot 1



This is a
blicket.

inner-shot 1



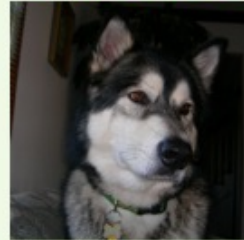
This is a dax.

inner-shot 2



This is a
blicket.

inner-shot 2



This is a dax.

**Question
from ImageNet**



Q: What is this?
A: This is a

Model Completion

blicket.

(b) Fast VQA

0-repeats
0-shots
2-way
0-repeats
2-inner-shots

**Support
from ImageNet**

inner-shot 1



This is a
blicket.

inner-shot 1



This is a dax.

inner-shot 2



This is a
blicket.

inner-shot 2



This is a dax.

**Question
from VisualGenome**



Q: What is the
dax made of? A:

blicket (vase)

dax (table)

Model Completion

wood

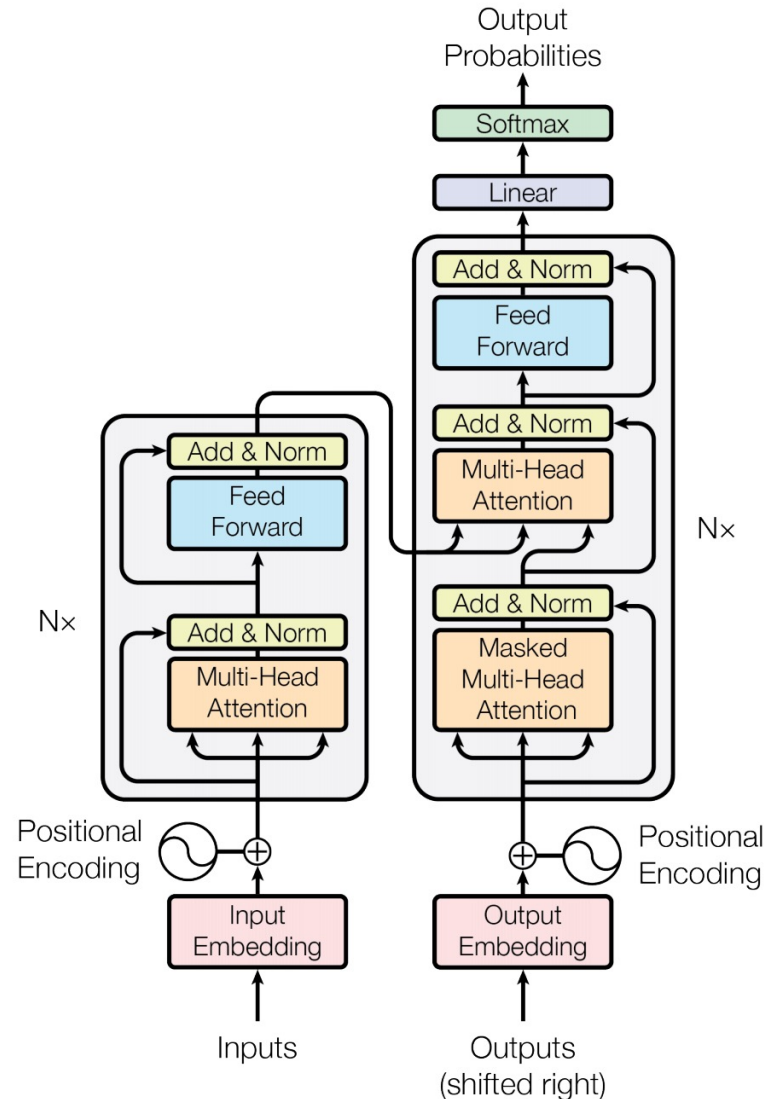
Examples of (a) the Open-Ended miniImageNet evaluation (b) the Fast VQA evaluation.

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.

Transformer (Attention is All You Need)

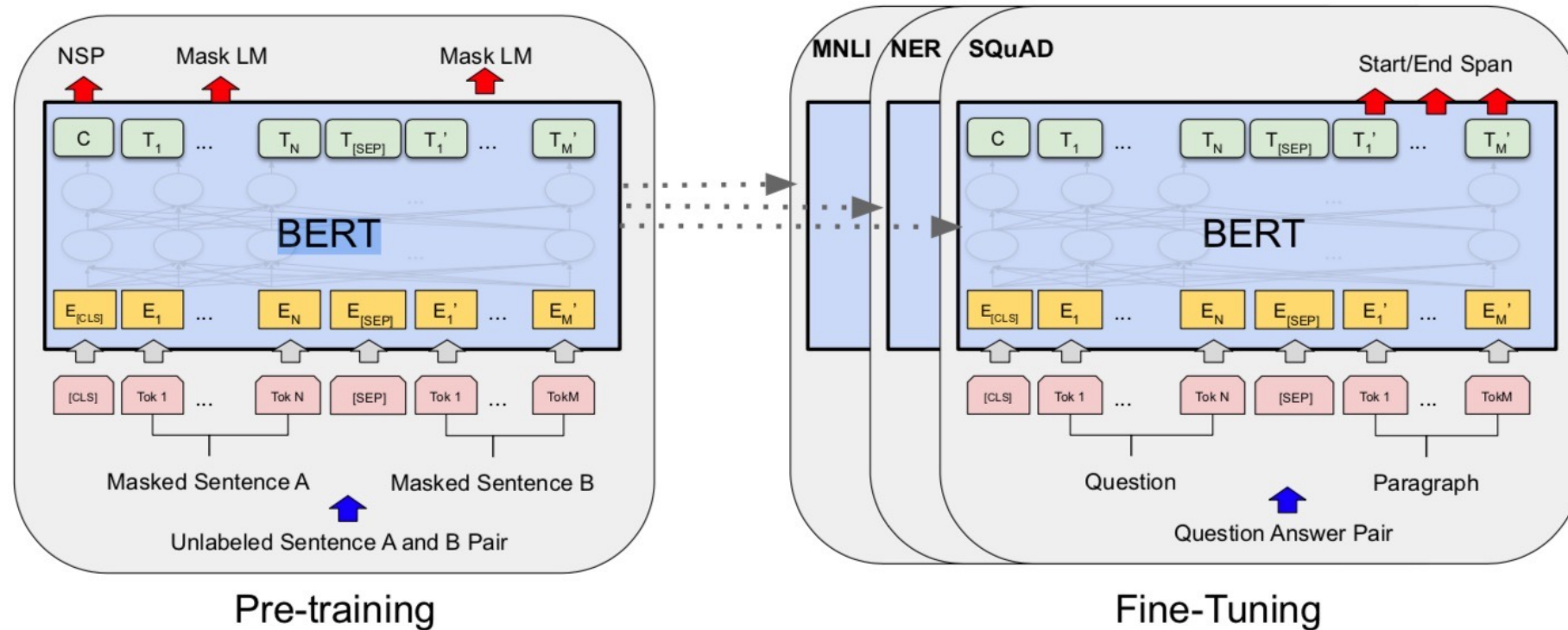
(Vaswani et al., 2017)



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



Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).

"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

Machine Learning: Ensemble Learning

Random Forest

The screenshot shows a Jupyter Notebook interface. On the left is a 'Table of contents' sidebar with a search icon and a list of topics. The 'Random Forest' topic is highlighted with a vertical bar. The main area shows a code cell with the following Python code:

```
1 # Randon 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

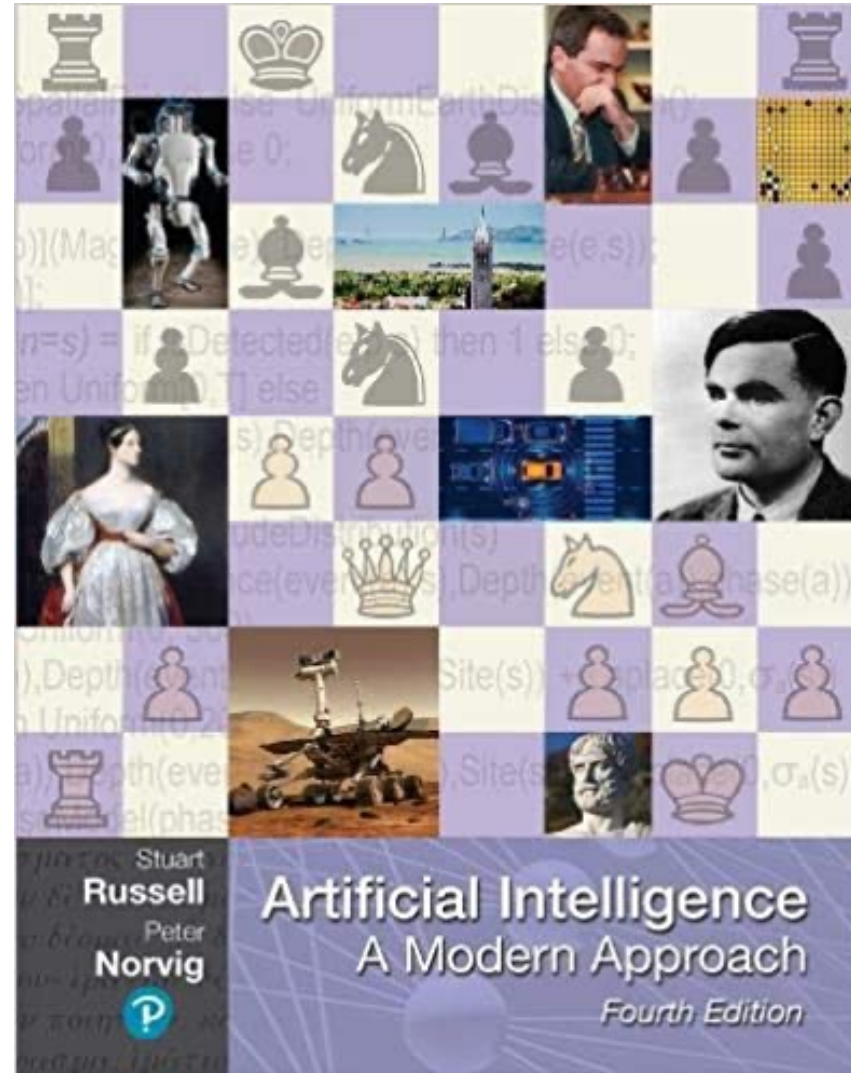
The screenshot shows a Jupyter Notebook interface. At the top, the notebook is titled 'python101.ipynb' and has a star icon. The top bar includes 'File', 'Edit', 'View', 'Insert', 'Runtime', 'Tools', and 'Help' menus, along with 'Last saved at 10:43 AM'. On the right, there are icons for 'Comment', 'Share', 'Settings', and a user profile 'A'.

On the left, a 'Table of contents' sidebar is visible, listing various topics under 'Machine Learning with scikit-learn', with 'Classification and Prediction' highlighted. Other topics include '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', 'Investment Portfolio Optimization', 'Text Analytics and Natural Language Processing (NLP)', 'Python for Natural Language Processing', 'spaCy Chinese Model', 'Open Chinese Convert (OpenCC, 開放中文轉換)', 'Jieba 結巴中文分詞', 'Natural Language Toolkit (NLTK)', 'Stanza: A Python NLP Library for Many Human Languages', and 'Text Processing and Understanding'.

The main area shows a code cell with 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']
```

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**

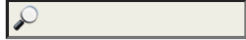
- <http://aima.cs.berkeley.edu/python/readme.html>

- <https://github.com/aimacode/aima-python>

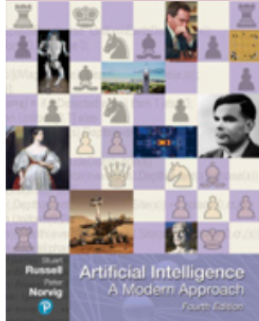
- **Learning**

- <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 (RF)**
 - **Boosting: Gradient Boosting, XGBoost, LightGBM, CatBoost**
 - **Stacking**
 - **Online learning**
- **Meta Learning: Learning to Learn**

References

- Stuart Russell and Peter Norvig (2020), Artificial Intelligence: A Modern Approach, 4th Edition, Pearson.
- Numa Dhamani and Maggie Engler (2024), Introduction to Generative AI, Manning
- Denis Rothman (2024), Transformers for Natural Language Processing and Computer Vision - Third Edition: Explore Generative AI and Large Language Models with Hugging Face, ChatGPT, GPT-4V, and DALL-E 3, 3rd ed. Edition, Packt Publishing
- Ben Auffarth (2023), Generative AI with LangChain: Build large language model (LLM) apps with Python, ChatGPT and other LLMs, Packt Publishing.
- Aurélien Géron (2022), Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 3rd 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.
- Aske Plaat, Annie Wong, Suzan Verberne, Joost Broekens, Niki van Stein, and Thomas Back. (2024) "Reasoning with Large Language Models, a Survey." arXiv preprint arXiv:2407.11511.
- Madaan, Aman, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon et al. (2024) "Self-refine: Iterative refinement with self-feedback." Advances in Neural Information Processing Systems 36.