

Few-shot learning

State of the Art

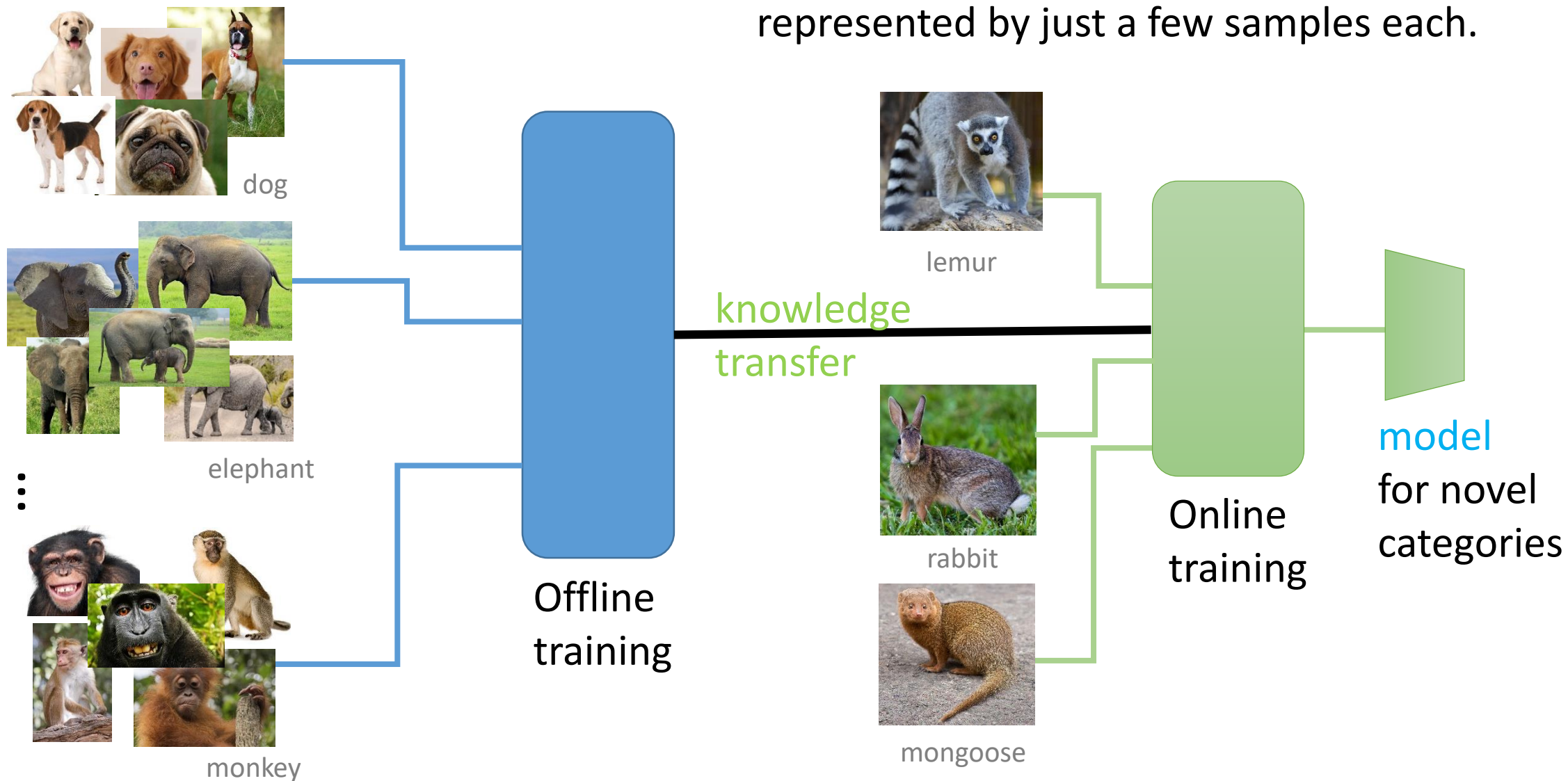
Joseph Shtok

IBM Research AI

Problem statement



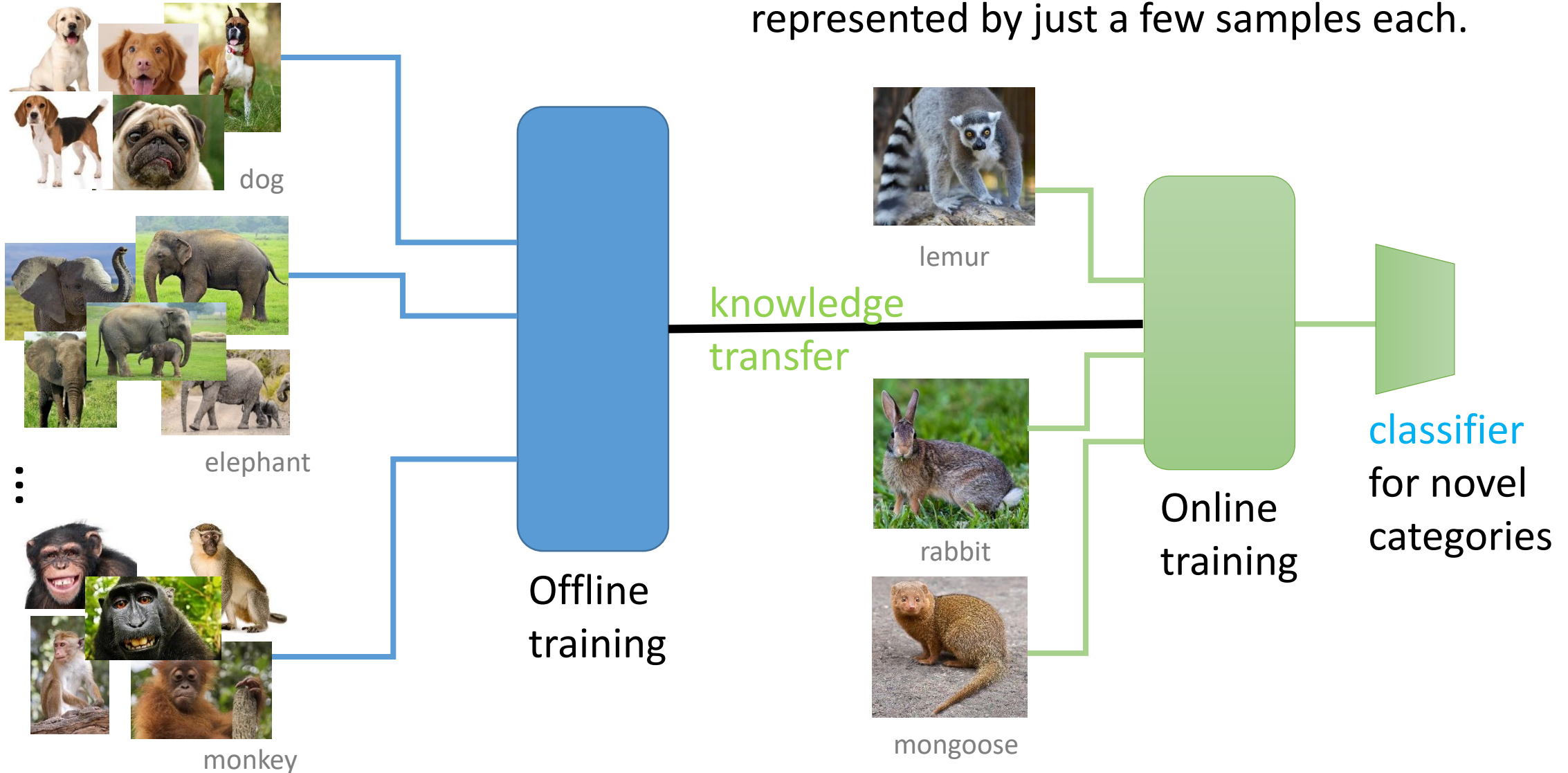
Using a large annotated offline dataset, perform **given task** for novel categories, represented by just a few samples each.



Problem statement



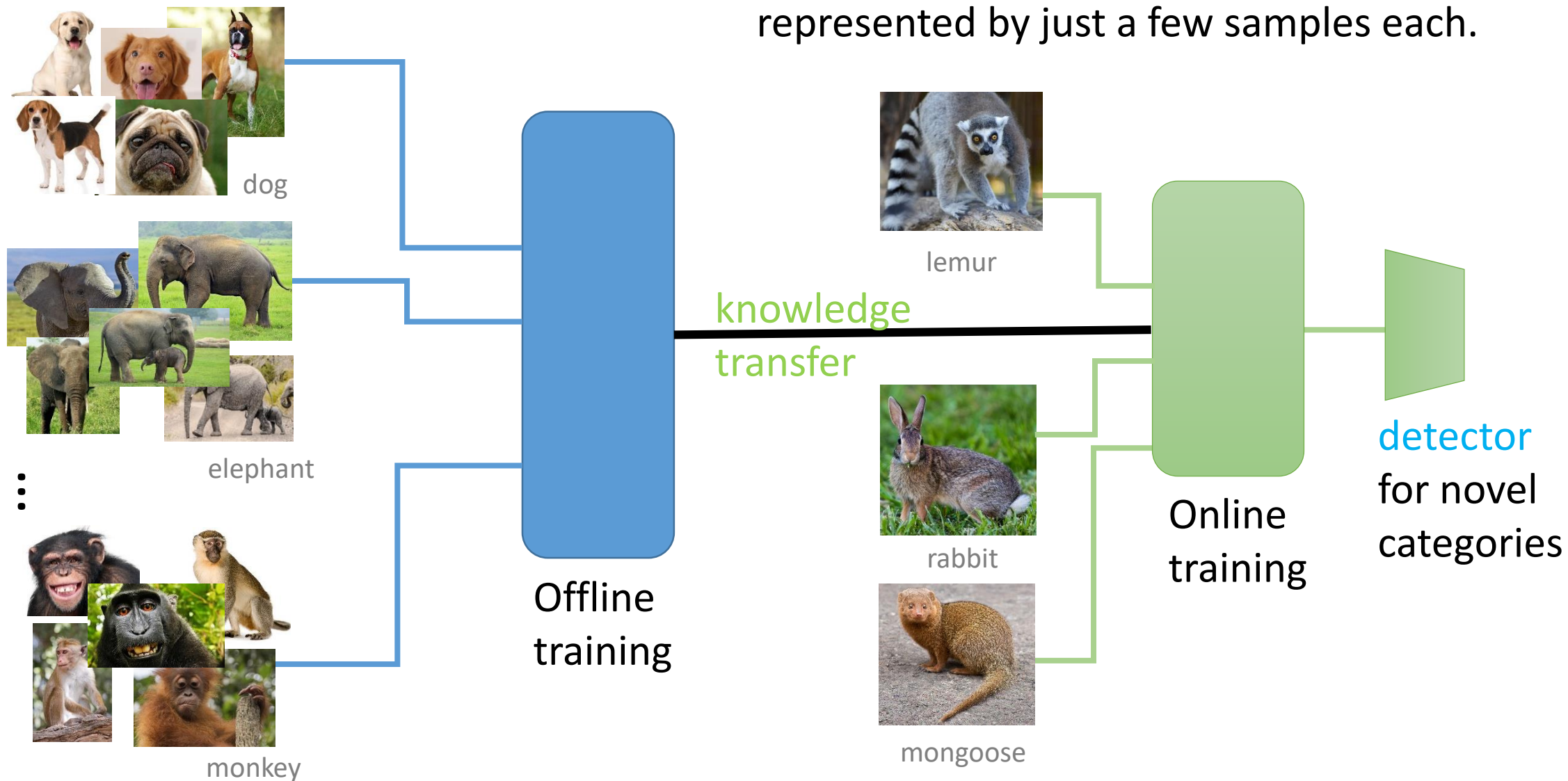
Using a large annotated offline dataset, perform **classification** for novel categories, represented by just a few samples each.



Problem statement



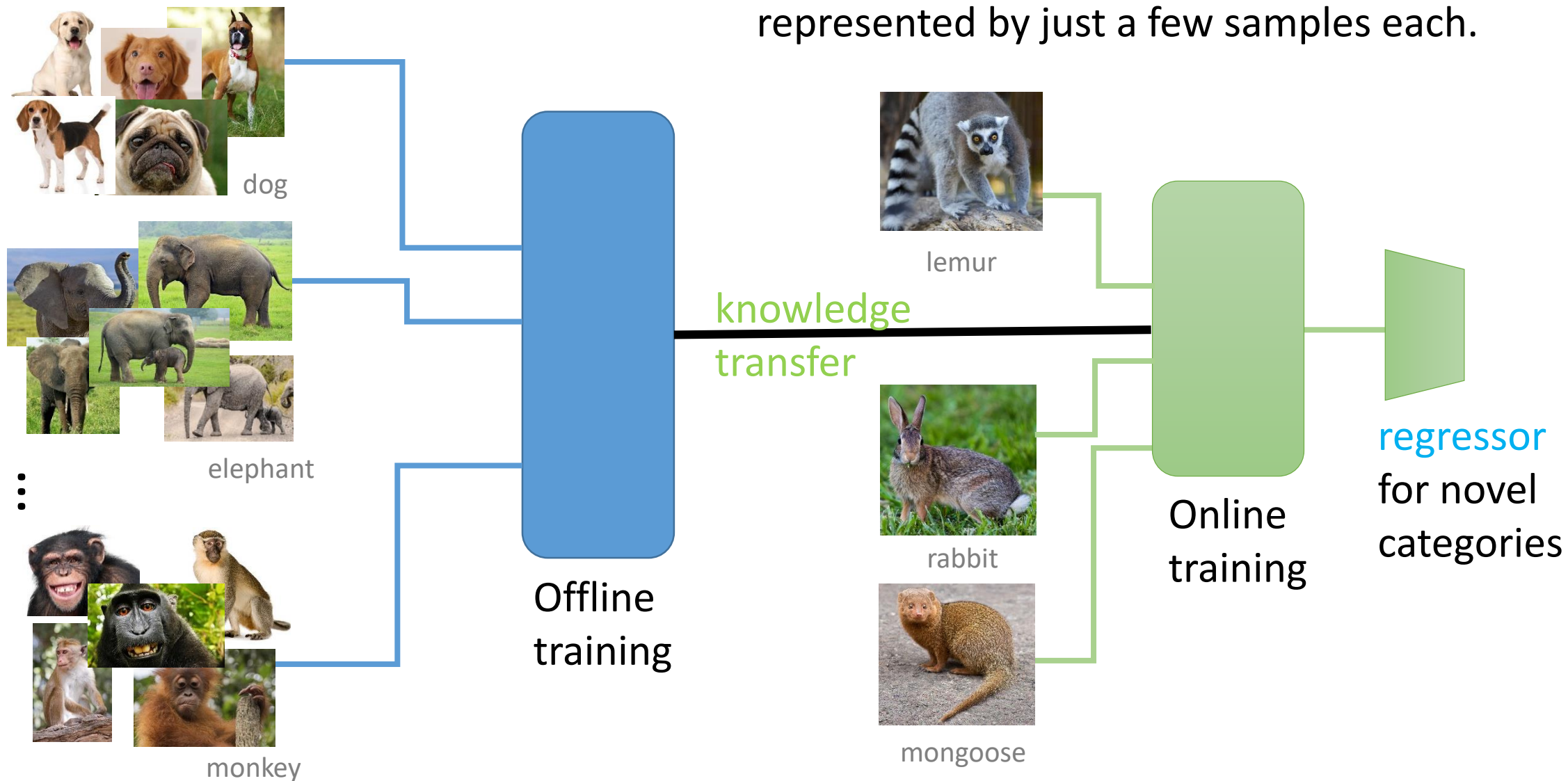
Using a large annotated offline dataset, perform **detection** for novel categories, represented by just a few samples each.



Problem statement



Using a large annotated offline dataset, perform **regression** for novel categories, represented by just a few samples each.



Why work on few-shot learning?



2. It brings the DL closer to exciting world business technologies.

- Companies hesitate to spend much money on a solution that they can't reuse. **Meta-learning methods** and **Data synthesizers** help address this.
- Relevant objects are constantly changing with new ones. DL has to be agile. **Neural Turing Machines** and **GANs** help address this.

Graph neural networks

Semantic metric spaces

Networks generating networks

Few-shot learning

Meta-learning

Learn a learning strategy to adjust well to a new few-shot learning task

Metric learning

Learn a `semantic` embedding space using a distance loss function

Learn to perform classification, detection, regression



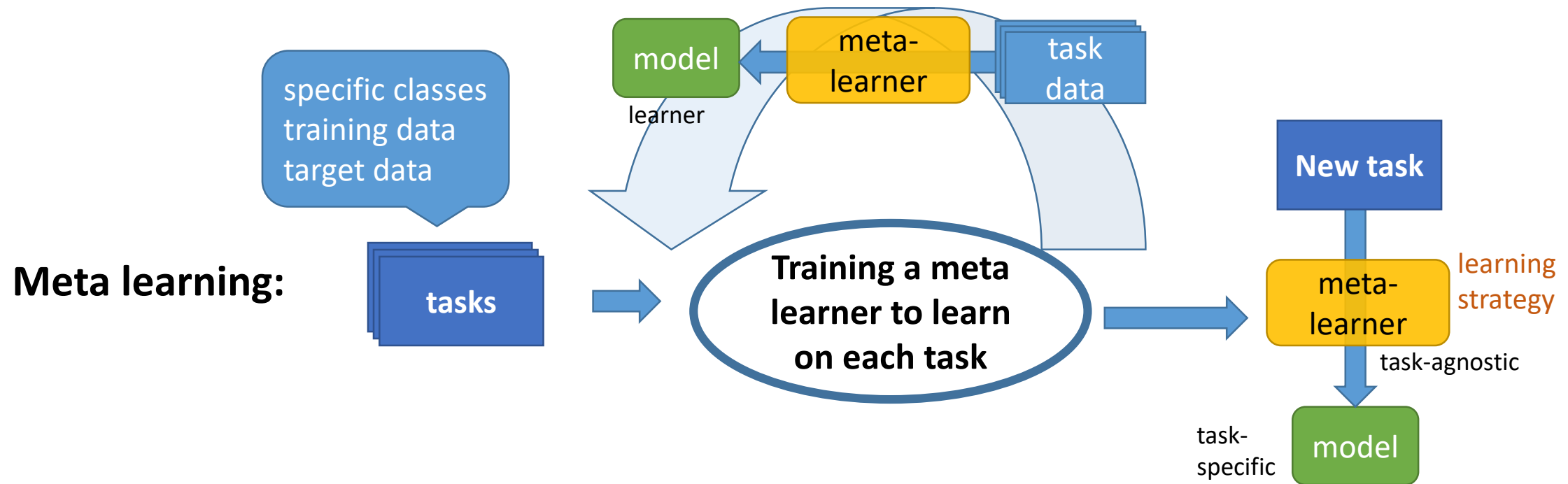
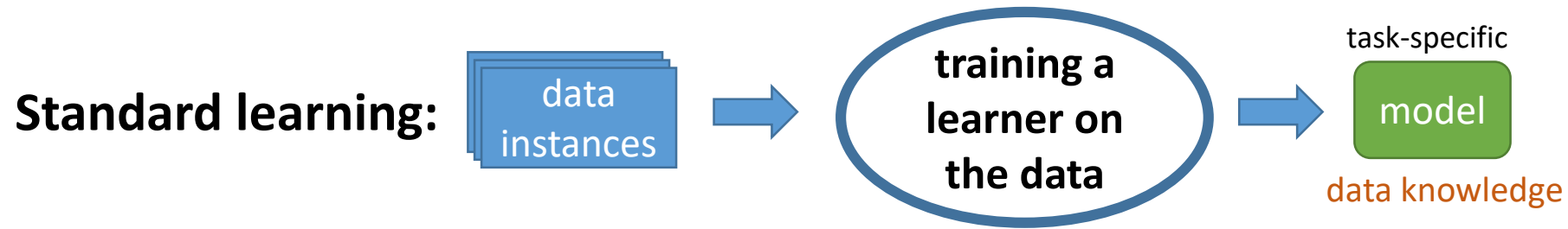
Each category is represented by just a few examples



Data augmentation

Synthesize more data from the novel classes to facilitate the regular learning

Meta-Learning



Recurrent meta-learners

Matching Networks in Vinyals et.al., NIPS 2016

Distance-based classification: based on similarity between the query and support samples in the embedding space (*adaptive metric*):

$$\hat{y} = \sum_i a(\hat{x}, x_i) y_i, \quad a(\hat{x}, x_i) = \text{similarity}(f(\hat{x}, S), g(x_i, S))$$

to be elaborated later

f, g - LSTM embeddings of x dependent on the support set S

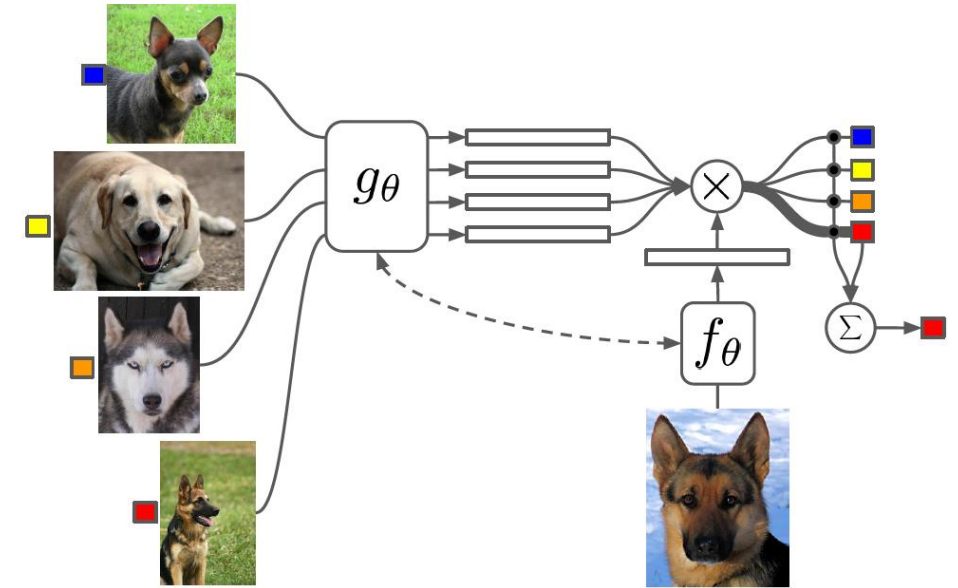


Figure 1: Matching Networks architecture
reprinted from Vinyals et.al., 2016

- Embedding space is class-agnostic
- LSTM attention mechanism

Memory-augmented neural networks
ICML 2016

- Neural Turing Machine = differentiable neural computer
- Learn to predict the distribution $p(y_t | x_t, S_{1:t-1}; \theta)$
- Explicitly store the support samples in the external memory

Method	minImageNet classification accuracy 1/5 shot
Matching networks	43.56 / 55.31

Concept of episodes: test conditions in the training.

N new categories
M training examples per category
one query example in $\{1..N\}$ categories.
Typically, $N=5, M=1, 5$.

Optimizers



Optimize the learner to perform well **after fine-tuning** on the task data done by a single (or few) step(s) of Gradient Descent.

MAML (Model-Agnostic Meta-Learning) [Finn et.al., ICML 2017](#)

Standard objective (task-specific, for task T):

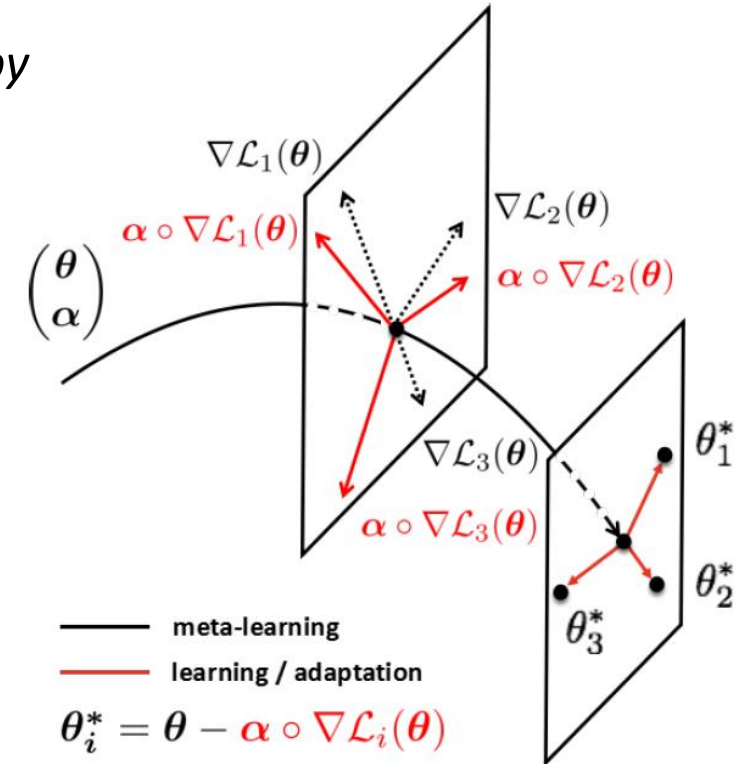
$$\min_{\theta} \mathcal{L}_T(\theta), \text{ learned via update } \theta' = \theta - \alpha \cdot \nabla_{\theta} \mathcal{L}_T(\theta)$$

Meta-objective (across tasks):

$$\min_{\theta} \sum_{T \sim p(\mathcal{T})} \mathcal{L}_T(\theta'), \text{ learn } (\theta')$$

Method	minImageNet classification accuracy 1/5 shot
Matching networks	43.56 / 55.31
MAML	48.70 / 63.11
Meta-SGD	54.24 / 70.86

Meta-SGD [Li et.al., 2017](#)



reprinted from Li et.al., 2017

More interestingly, the training process can continue forever, thus enabling life-long learning, and at training time, the meta-optimizer optimizes the loss $\mathcal{L}_p(\theta')$ with respect to the parameters θ for weight initialization, α = update direction and scale, across the tasks.

Optimizers

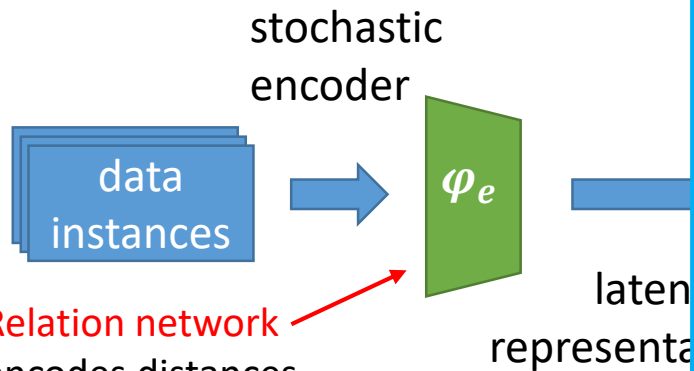
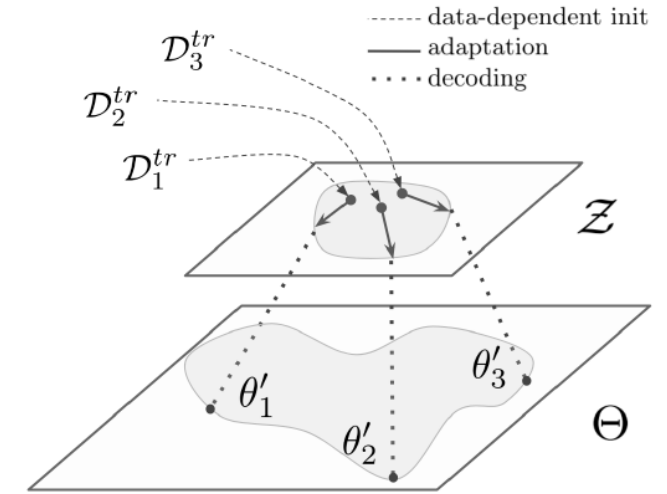


LEO Rusu et.al., NeurIPS 2018

Latent Embedding Optimization: take the optimization problem from high-dim. space of weights θ to a low-dim. space, for robust Meta-Learning.

Learn a **generative distribution** of model parameters θ , by learning a stochastic **latent space** with an information bottleneck.

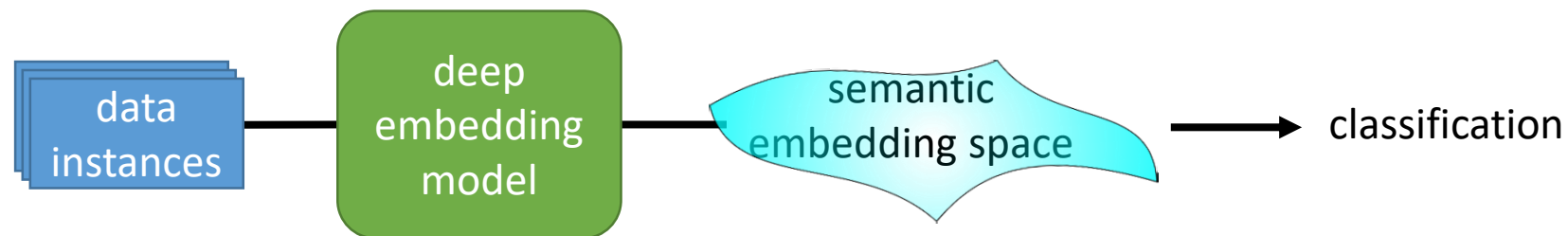
$$\varphi^* = \arg \min_{\varphi} \sum_{T \sim p(\mathcal{Z})} \mathcal{L}_T(\theta' = g_{\varphi}(z')) \quad z' = z - \alpha \cdot \nabla_z \mathcal{L}_T(\theta_z)$$



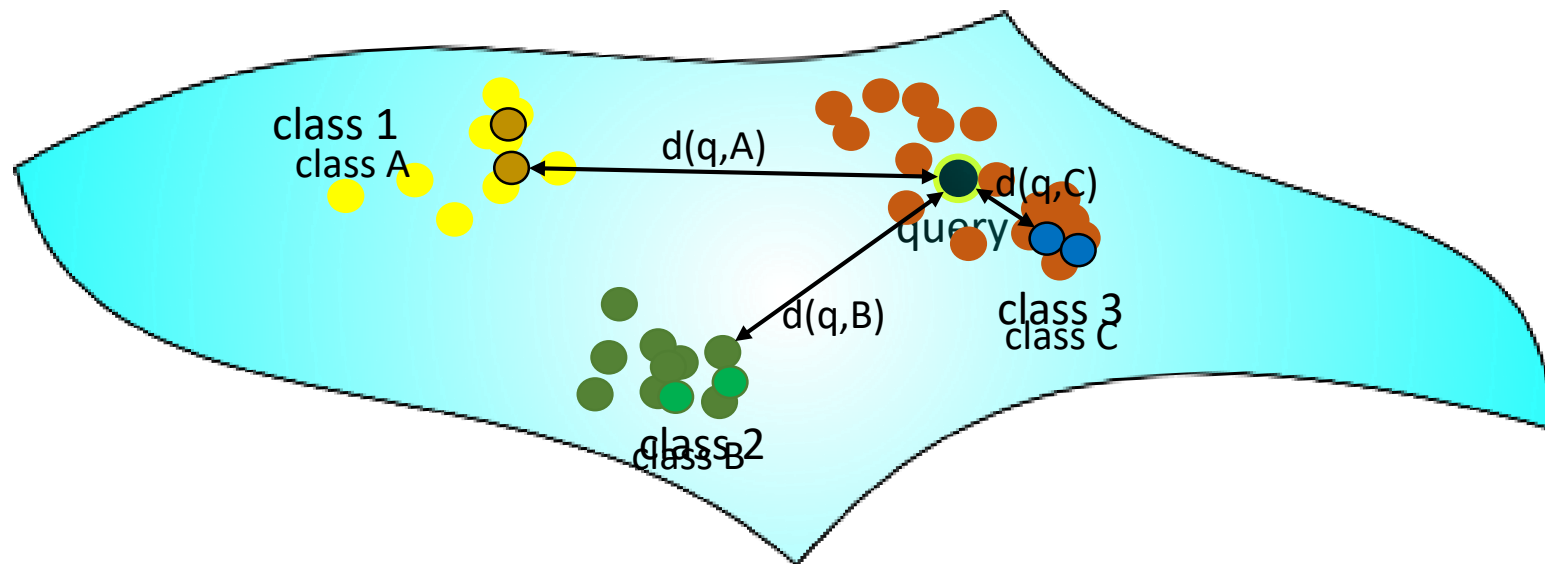
Relation network encodes distances between elements of support set for a new task

Method	minImageNet classification accuracy 1/5 shot
Matching networks	43.56 / 55.31
MAML	48.70 / 63.11
Meta-SGD	54.24 / 70.86
LEO	61.76 / 77.59

Metric Learning



Offline training
New task data



Training: achieve good distributions for offline categories

Inference: Nearest Neighbour in the embedding space

Metric Learning

Relation networks, Sung et.al., CVPR 2018

Use the Siamese Networks principle :

- Concatenate embeddings of query and support samples
- Relation module is trained to produces score 1 for correct class and 0 for others
- Extends to zero-shot learning by replacing support embeddings with semantic features.

Method	minImageNet classification accuracy 1/5 shot
Matching networks	43.56 / 55.31
MAML	48.70 / 63.11
Relation networks	50.44 / 65.32
Meta-SGD	54.24 / 70.86
LEO	61.76 / 77.59

One-hot vector

Sung et.al., Learning to compare relation network for few-shot learning, CVPR 2018

Metric Learning



Matching Networks, Vinyals et.al., NIPS 2016

Objective: maximize the cross-entropy for the non-parametric softmax

classifier $\sum_{(x,y)} \log P_{\theta}(y|x, S)$, with

$$P_{\theta}(y|x, S) = \text{soft}$$

Method	minImageNet classification accuracy 1/5 shot
Matching networks	43.56 / 55.31
MAML	48.70 / 63.11
Relation networks	50.44 / 65.32
Prototypical Networks	49.42 / 68.20
Meta-SGD	54.24 / 70.86
LEO	61.76 / 77.59

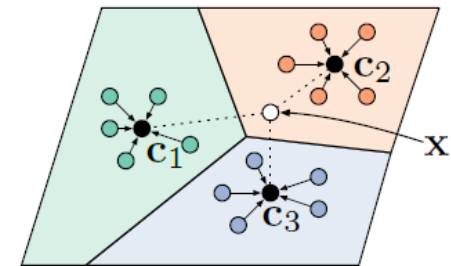
Prototypical Networks, Snell et

Each category is represented by

Objective: maximize the cross-entropy

probability expression:

$$P_{\theta}(y|x, C) = \text{soft}$$



Each category is represented by a single prototype c_i .

Metric Learning



Large Margin Meta-Learning

Regularize the cross-entropy objective:

$$\mathcal{L} =$$

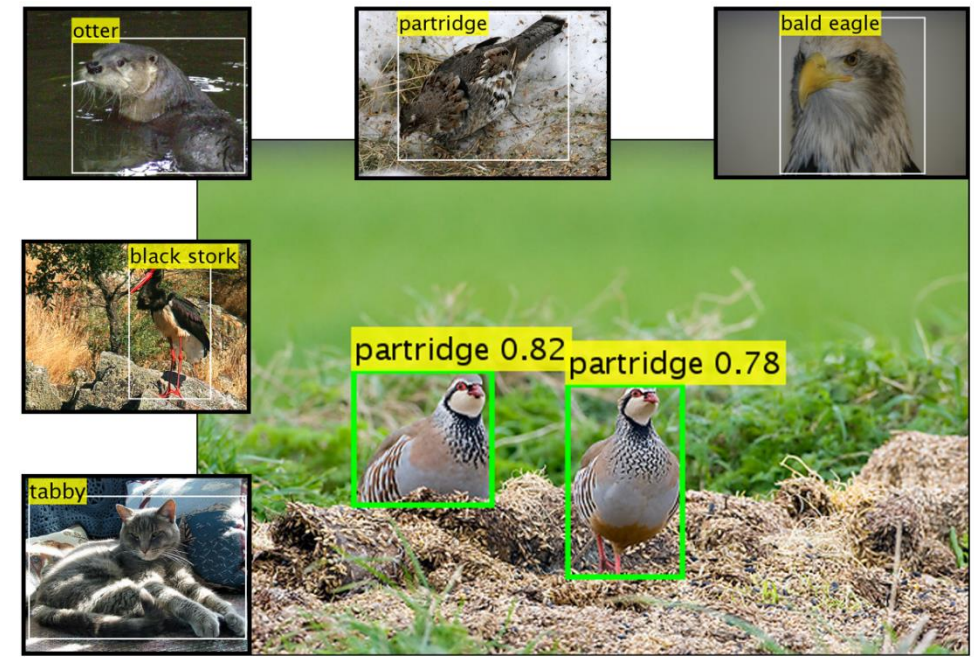
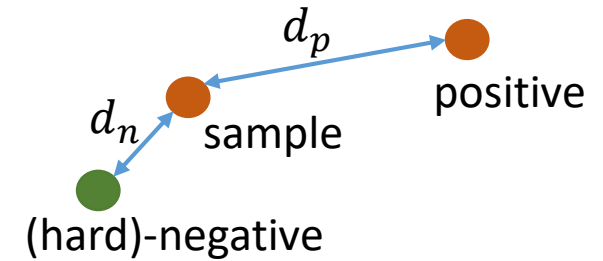
Method	minImageNet classification accuracy 1/5 shot
Matching networks	43.56 / 55.31
MAML	48.70 / 63.11
Relation networks	50.44 / 65.32
Prototypical Networks	49.42 / 68.20
Large-margin	51.08 / 67.57
Meta-SGD	54.24 / 70.86
LEO	61.76 / 77.59

RepMet: Few-shot detection

- Equip a standard
- Introduce class re
- Learn an embedding space using the objective

$$\mathcal{L} = \mathcal{L}_{CE} + \left| \min_j d(E, R_{ij}) - \min_{j, k \neq i} d(E, R_{kj}) + \alpha \right|_+$$

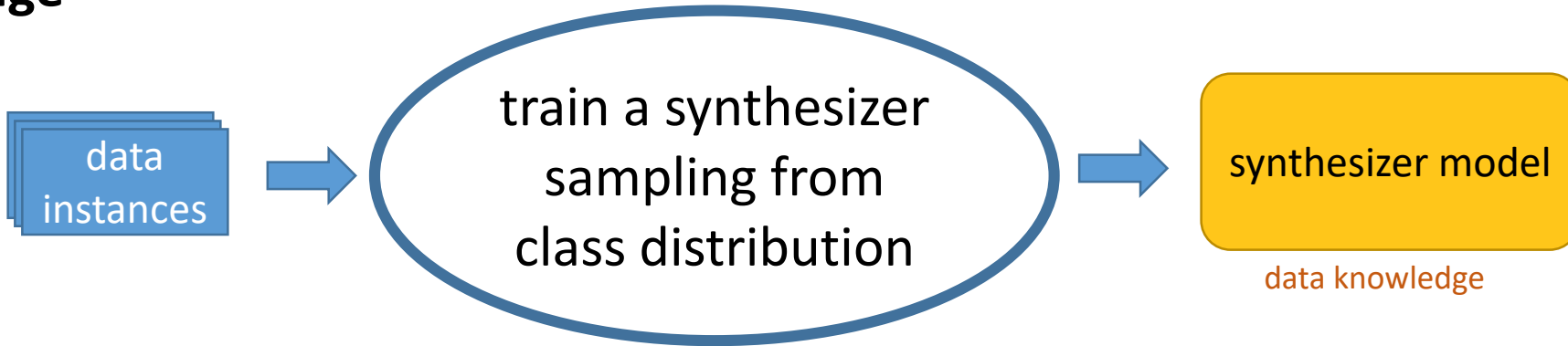
Distances for triplet loss



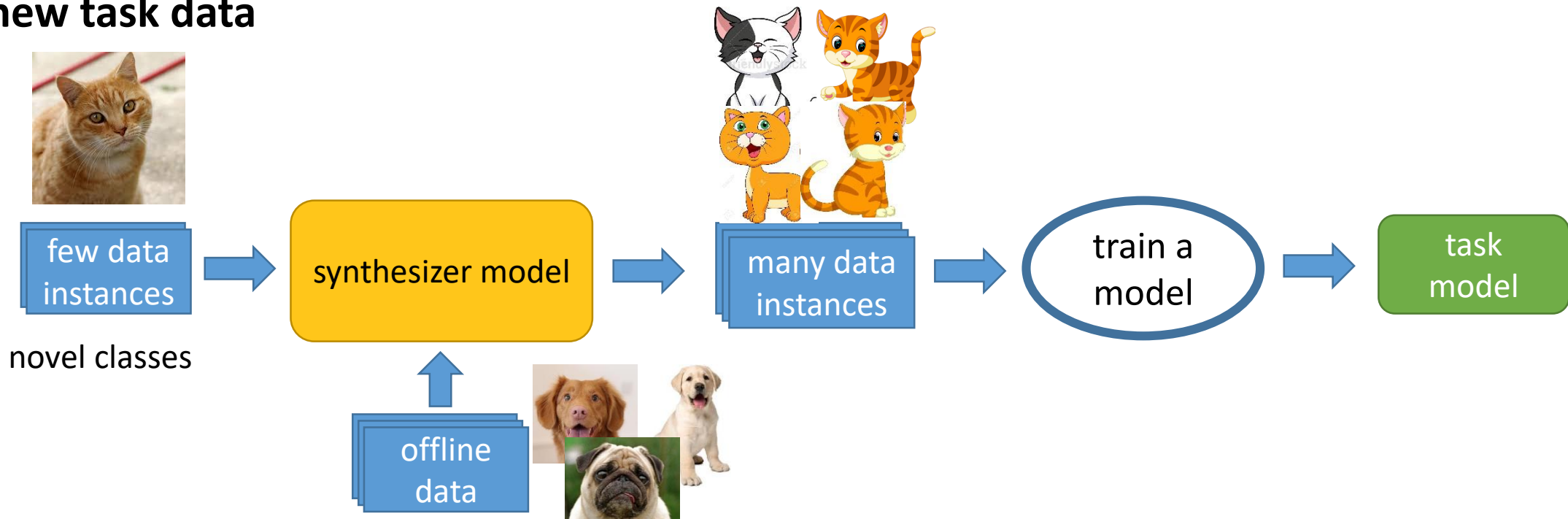
Sample synthesis



Offline stage



On new task data



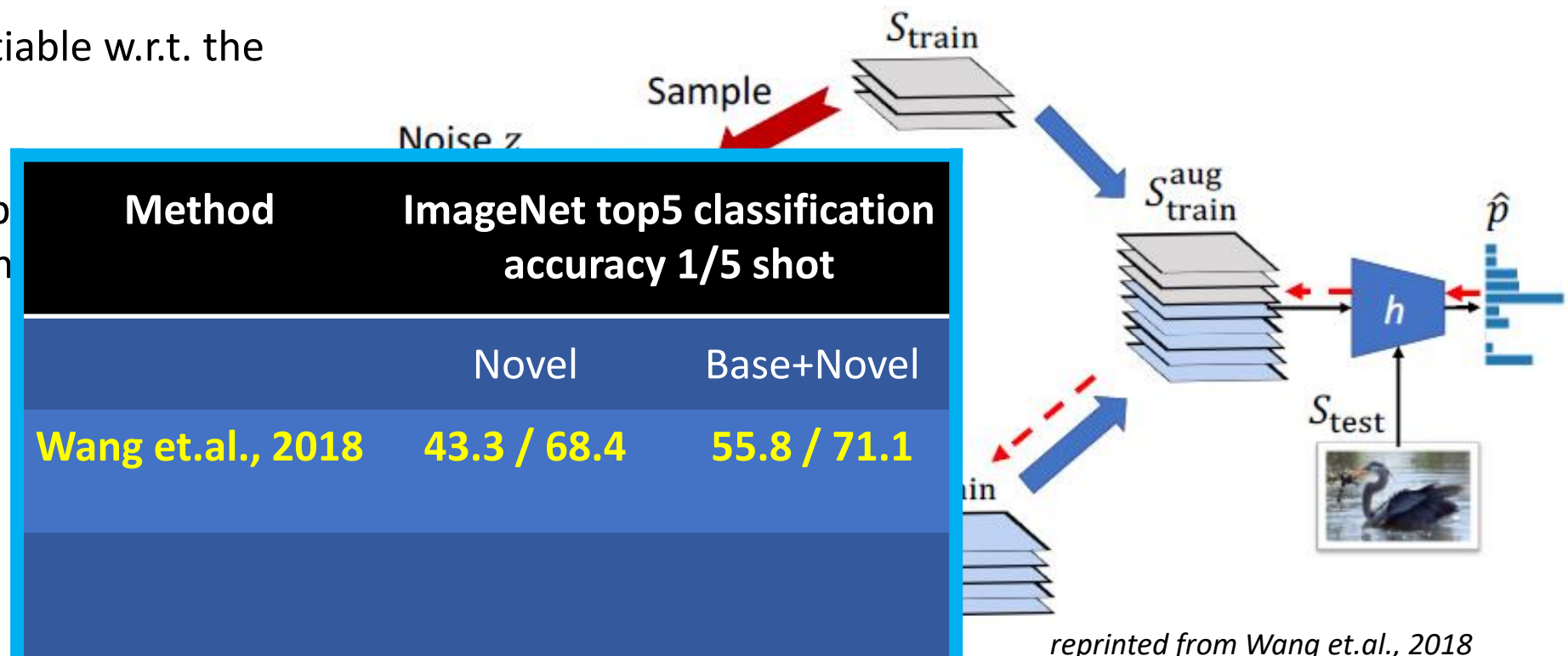
Synthesizer optimized for classification



Low-Shot Learning from Imaginary Data

Wang et.al., 2018

- The synthesizer is a part of classifier pipeline, trained end-to-end
- The classifier h is differentiable w.r.t. the training data
- In each episode, the backprop gradient updates the synthesizer



Method	ImageNet top5 classification accuracy 1/5 shot	
	Novel	Base+Novel
Wang et.al., 2018	43.3 / 68.4	55.8 / 71.1

reprinted from Wang et.al., 2018

More augmentation approaches



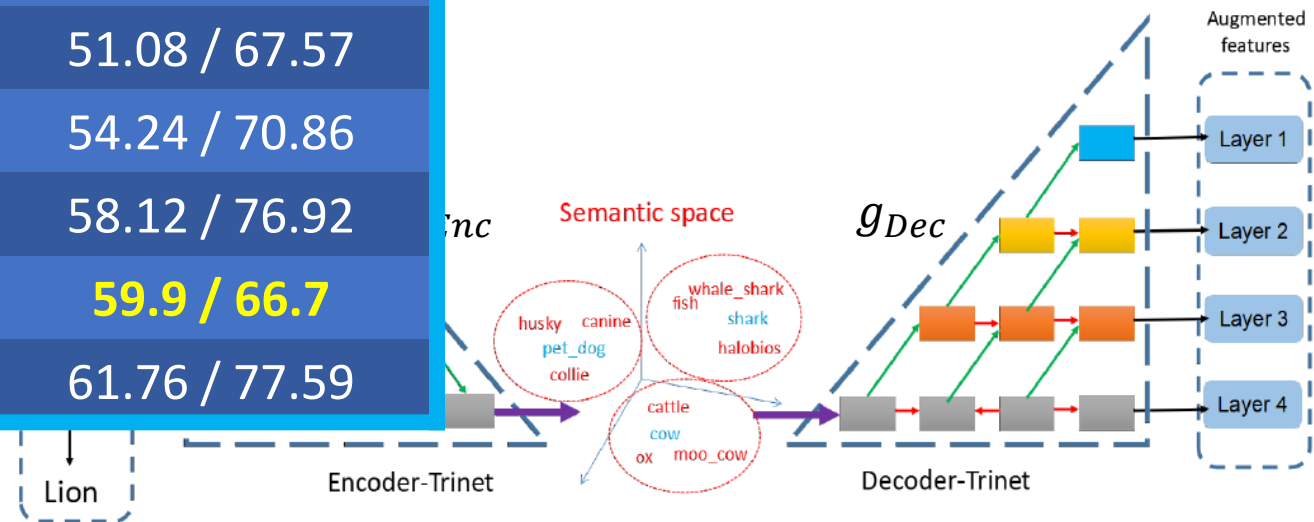
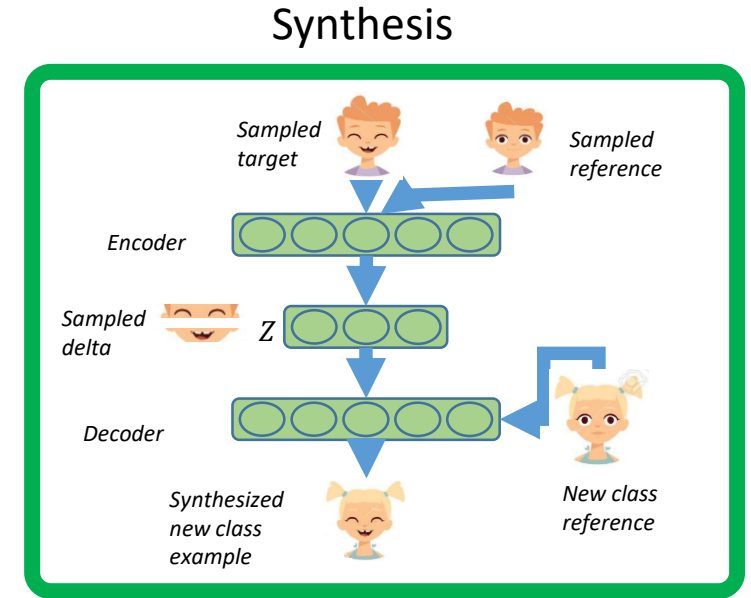
Δ -encoder Schwartz et.al., NeurIPS 2019

- Use a variant of autoencoder between two class samples
- Transfer class distributions

Method	minImageNet classification accuracy 1/5 shot
Matching networks	43.56 / 55.31
MAML	48.70 / 63.11
Relation networks	50.44 / 65.32
Prototypical Networks	49.42 / 68.20
Large-margin	51.08 / 67.57
Meta-SGD	54.24 / 70.86
Semantic Feat. Aug.	58.12 / 76.92
Δ-encoder	59.9 / 66.7
LEO	61.76 / 77.59

Semantic Feature Augmentation
Learning, Chen et.al., 2018

- Synthesize samples by adding semantic space to autoencoder's bottleneck
- Make it into a semantic space by adding semantic embeddings or visual attributes to the objective's fidelity term.



Augmentation with GANs



Covariance-Preserving Adversarial Augmentation Networks

Gao et.al., NeurIPS 2018

- Act in in feature space. Generate samples for novel categories from offline categories, selected by proximity of samples.

- Discriminate by classical

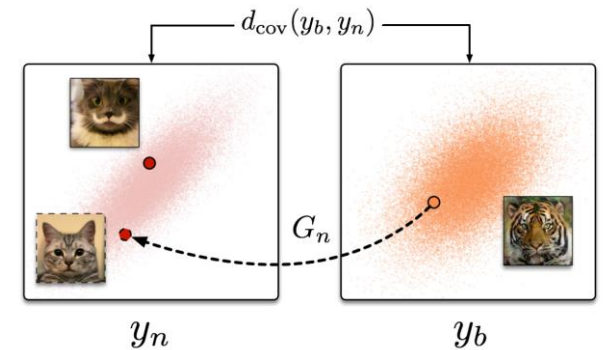
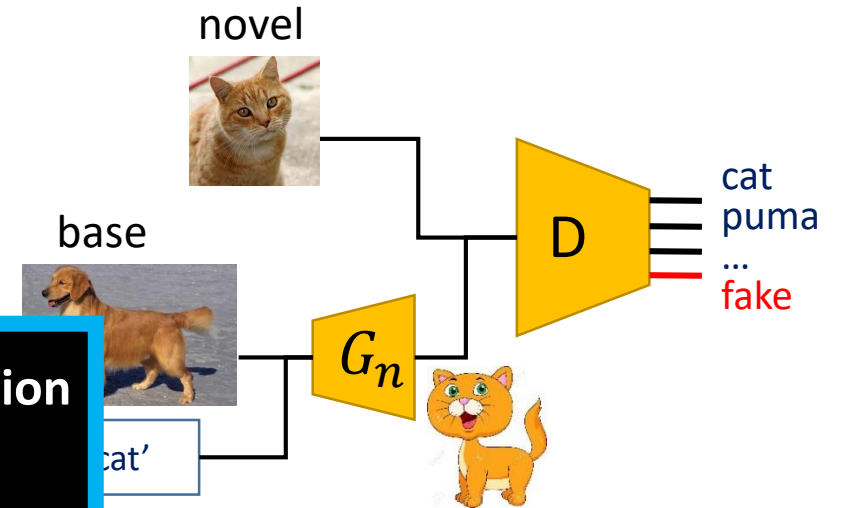
Method	ImageNet top5 classification accuracy 1/5 shot	
	Novel	Base+Novel
Wang et.al., 2018	43.3 / 68.4	55.8 / 71.1
Gao et.al., 2018	48.4 / 70.2	58.5 / 73.5

- Cycle-consistency const

- Preserve the category d

Objective: penalize the difference between the two covariance matrices via *Ky Fan m-norm*, i.e., the sum of singular values of m-truncated SVD:

$$d_{\text{cov}}(y_b, y_n) = \left\| \left[\Sigma_{\mathbf{x}}(\mathbb{P}_{y_b}) - \Sigma_G(\mathbb{P}_{y_n}) \right]_m \right\|_*$$



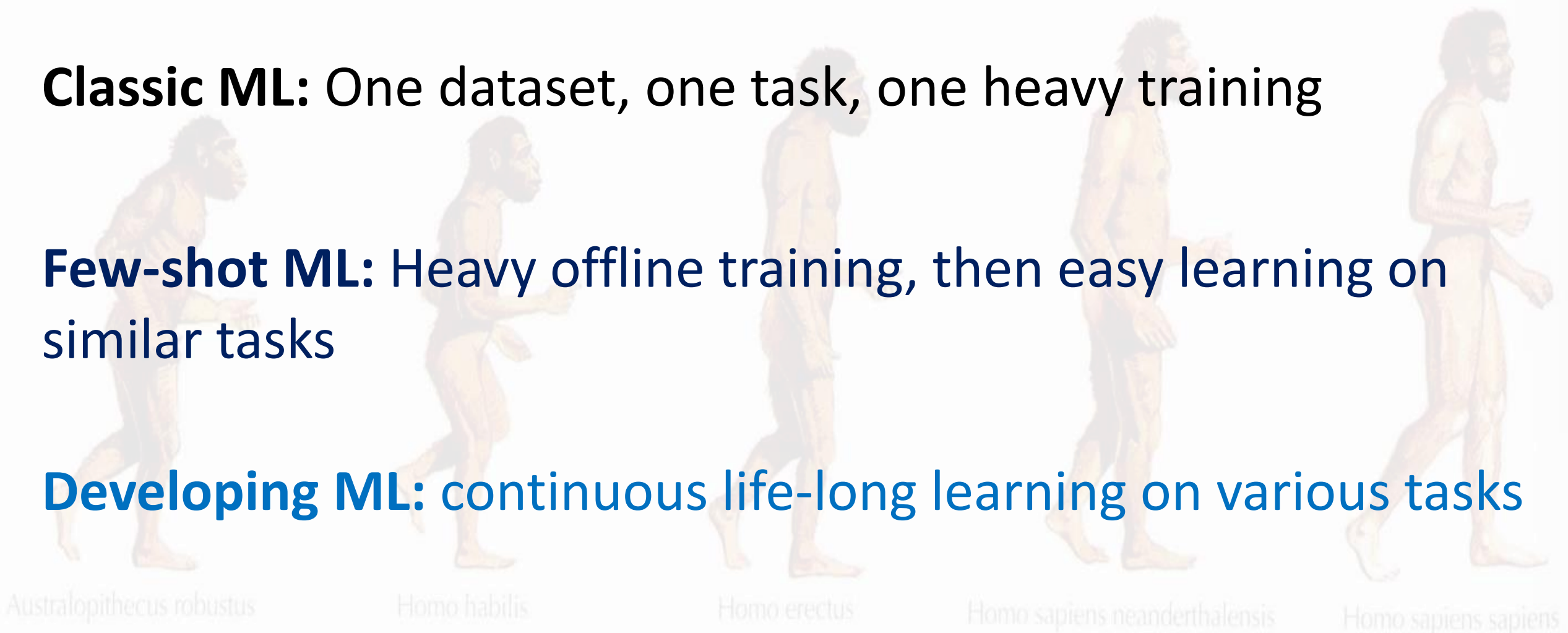
My personal view on the evolution of Machine Learning



Classic ML: One dataset, one task, one heavy training

Few-shot ML: Heavy offline training, then easy learning on similar tasks

Developing ML: continuous life-long learning on various tasks



Australopithecus robustus

Homo habilis

Homo erectus

Homo sapiens neanderthalensis

Homo sapiens sapiens

THANK YOU

The presentation is available at http://www.research.ibm.com/haifa/dept/imt/ist_dm.shtml