

Deep Learning

MPHY0041 Machine Learning in Medical Imaging

Yipeng Hu
yipeng.hu@ucl.ac.uk

Research

Selected Research Topic

Multiple instance learning

Adversarial learning and generative models

Domain adaptation

Meta-learning

Unsupervised learning*

Reinforcement learning*

Graph networks*

Research | Multiple Instance Learning

Weakly-supervised learning and Weak labels

- Inaccurate supervision, e.g. noisy labels
- Incomplete supervision, e.g. partial labels
- Inexact supervision, e.g. multiple instance learning

Bags, instances and labels in MIL

Serge's key-chain



Serge **cannot** enter
the *Secret Room*

Sanjoy's key-chain



Sanjoy **can** enter
the *Secret Room*

Lawrence's key-chain



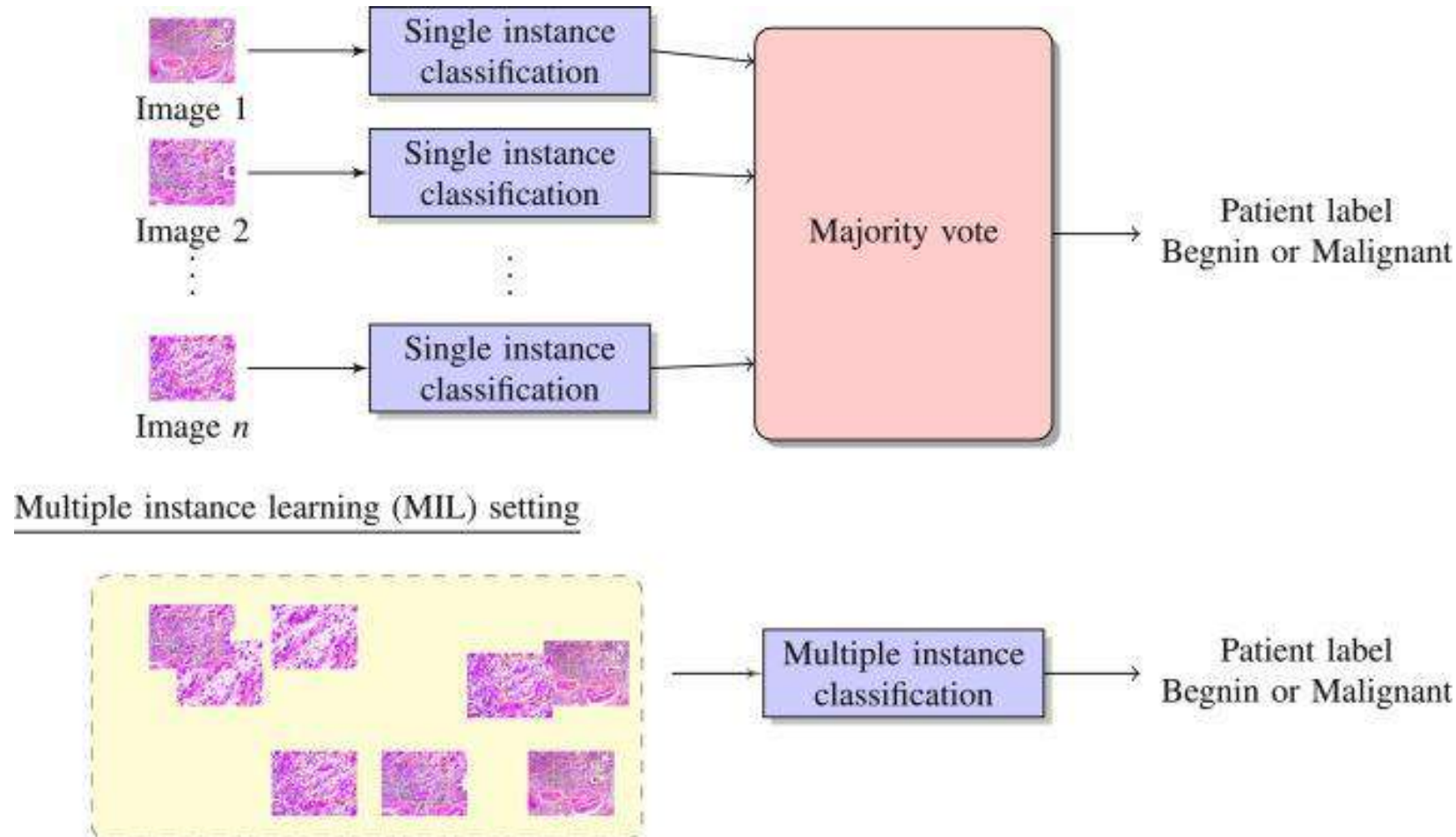
Lawrence **can** enter
the *Secret Room*

Linkage function between bag-level (known) and instance-level (to infer) relations

$$y_i = \begin{cases} 1 & \text{if } \exists j \text{ s.t. } y_{ij} = 1 \\ 0 & \text{otherwise} \end{cases} \qquad y_i = \max_j \{y_{ij}\}$$

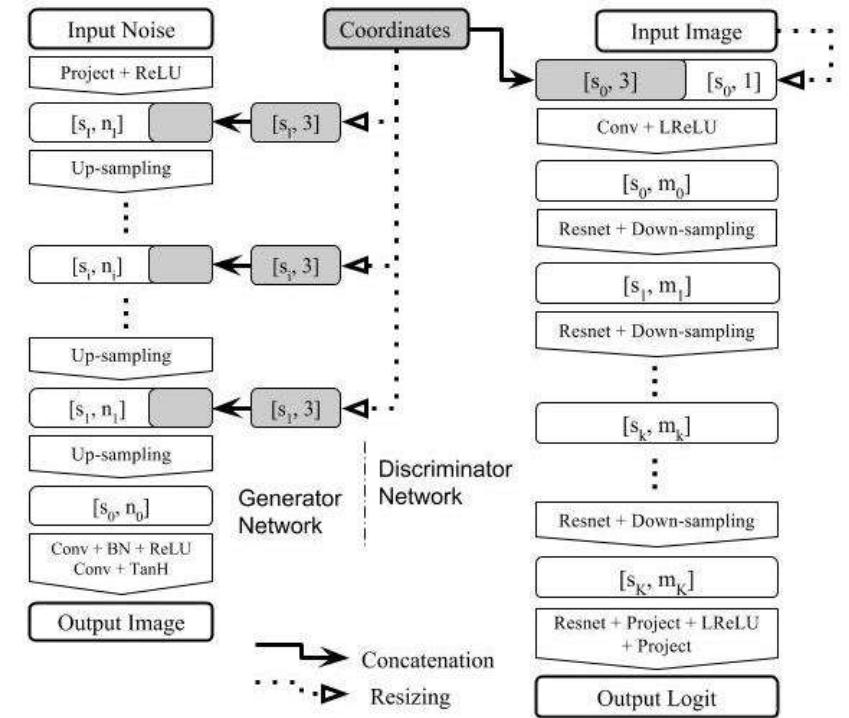
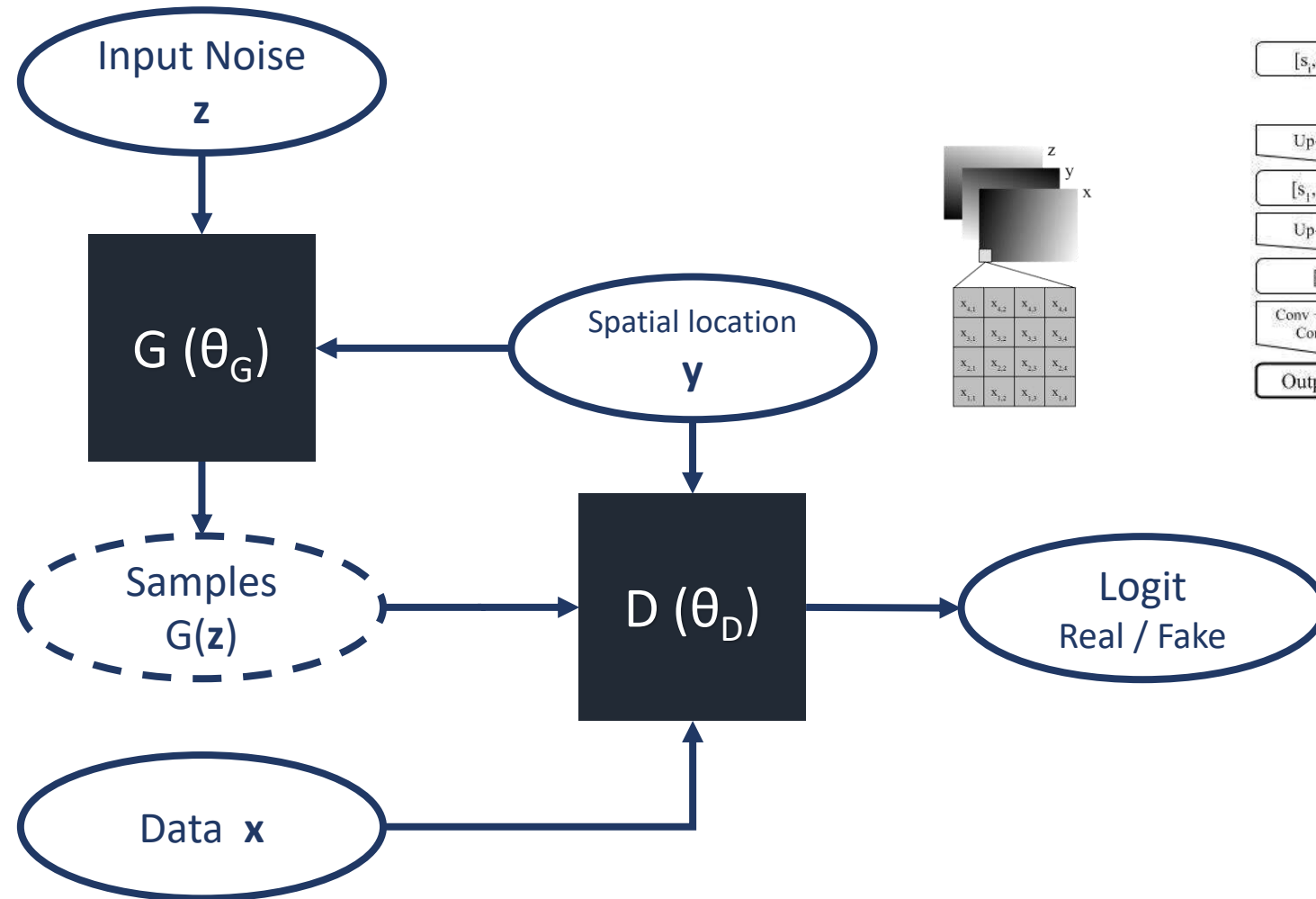
$x_i \in \mathcal{X}$	i^{th} instance (supervised)
$X_i = \{x_{i1}, x_{i2} \dots, x_{im}\}$	i^{th} bag with m instances x_{ij} (MIL)
$y_i \in \mathcal{Y}$	label of i^{th} instance (supervised) or i^{th} bag (MIL)
y_{ij}	true label of j^{th} instance in i^{th} bag
y'	$(2y - 1)$
n	instances (supervised) or number of bags (MIL)
m	number of instances per bag
d	dimensionality
$h(x)$	instance classifier
$H(X)$	bag classifier

Segmentation on whole-slide imaging in digital pathology

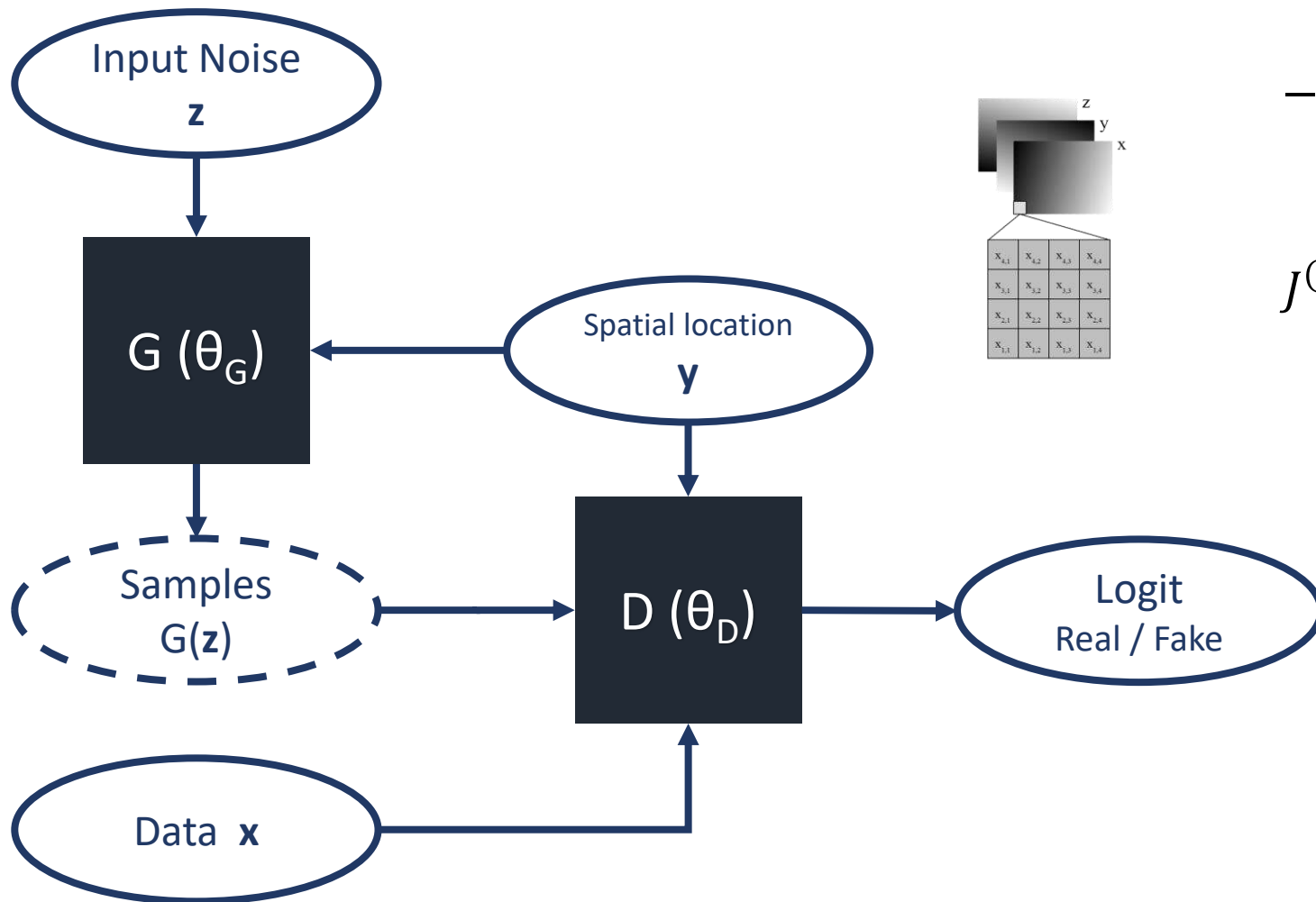


Research | Adversarial Learning and Generative Models

Generative adversarial networks (GANs)



Generative adversarial networks (GANs)



$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim P_{data}} \log D(\mathbf{x}, \mathbf{y}) - \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim N, \mathbf{y} \sim P_{loc}} \log(1 - D(G(\mathbf{z}, \mathbf{y}), \mathbf{y}))$$

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z} \sim N, \mathbf{y} \sim P_{loc}} \log D(G(\mathbf{z}, \mathbf{y}), \mathbf{y})$$

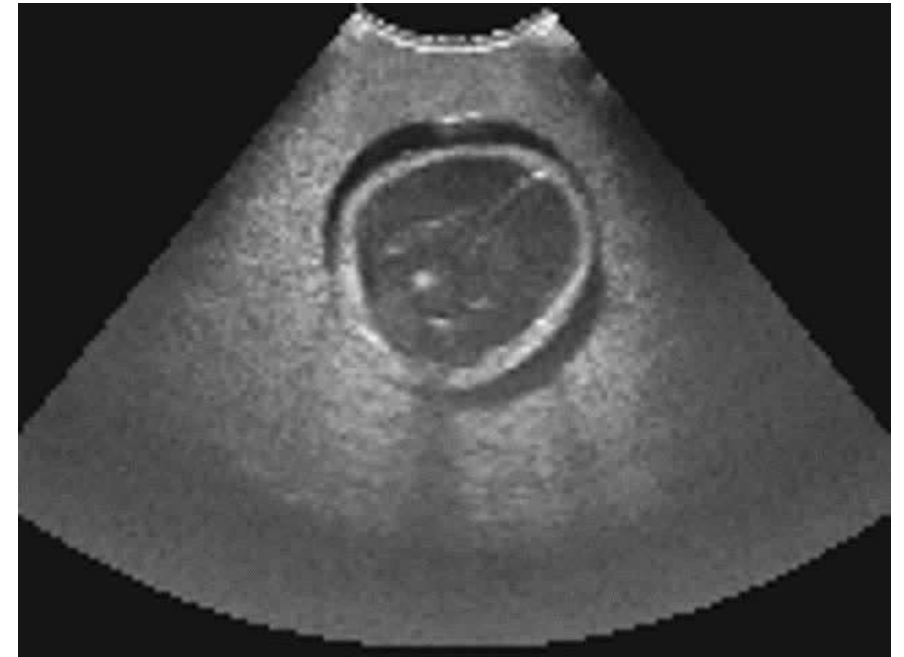
Generative adversarial networks (GANs)

Image synthesis

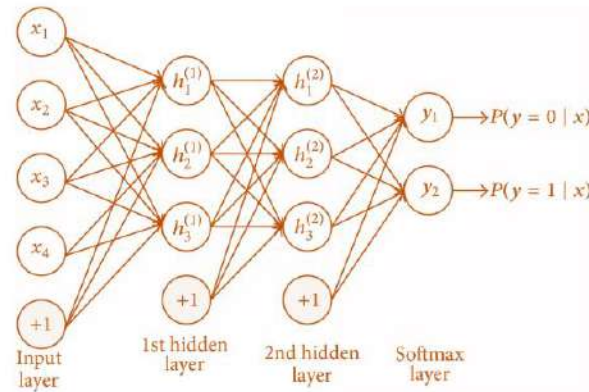
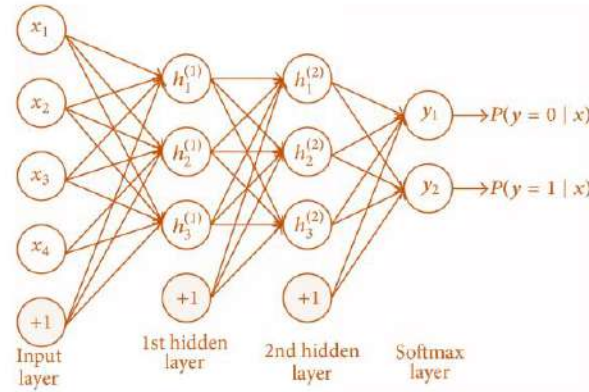
Regularisation by minimising divergence

Image-to-image translation

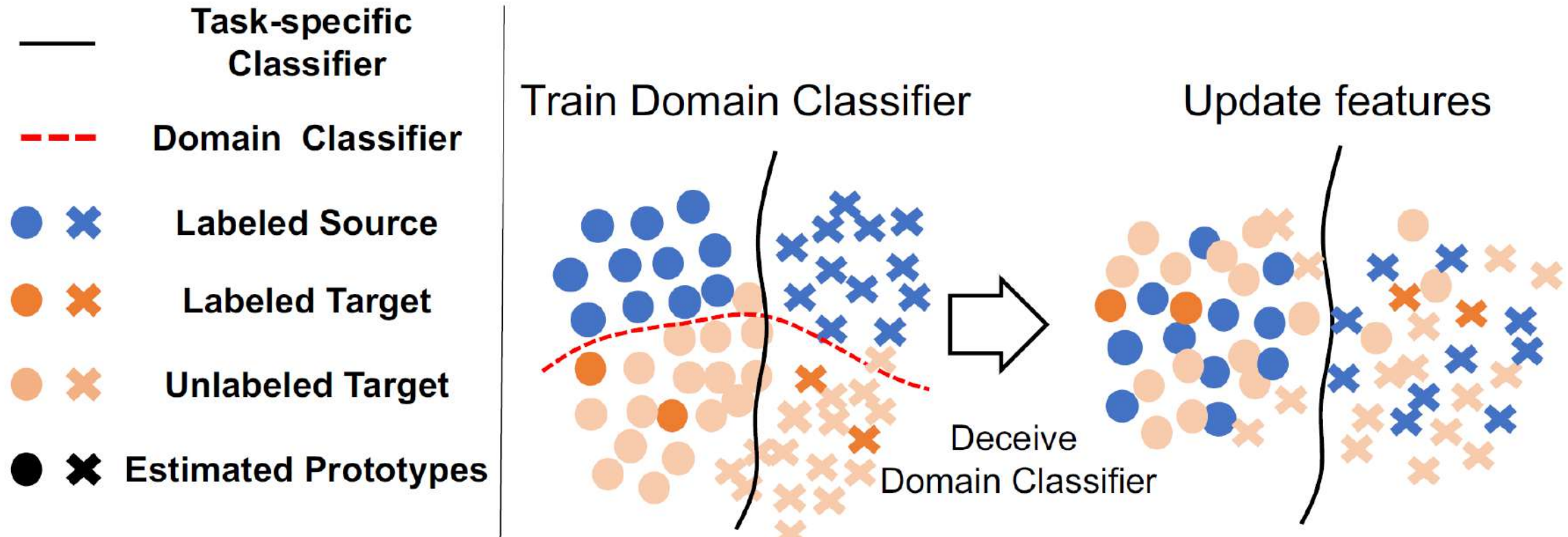
Domain adaptation



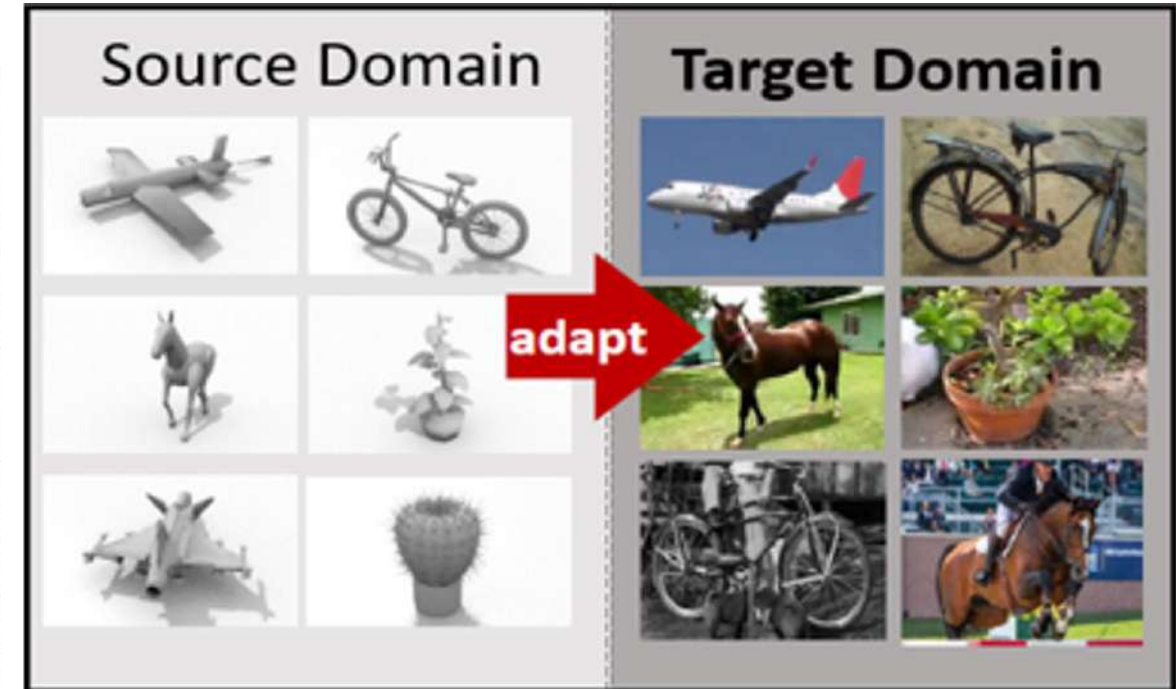
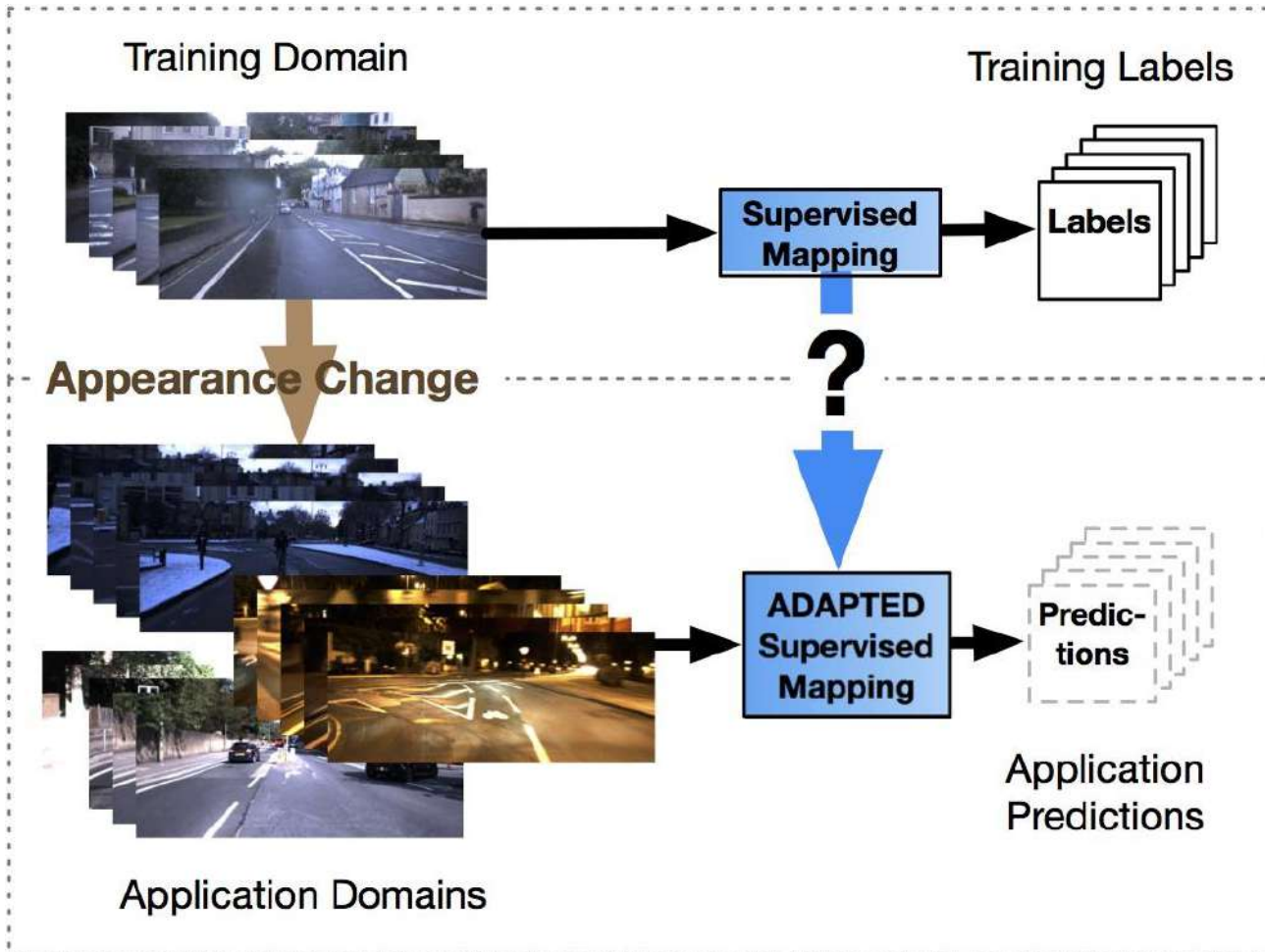
Research | Domain Adaptation



- Unsupervised domain adaptation
- Semi-supervised domain adaptation



- Availability between domains

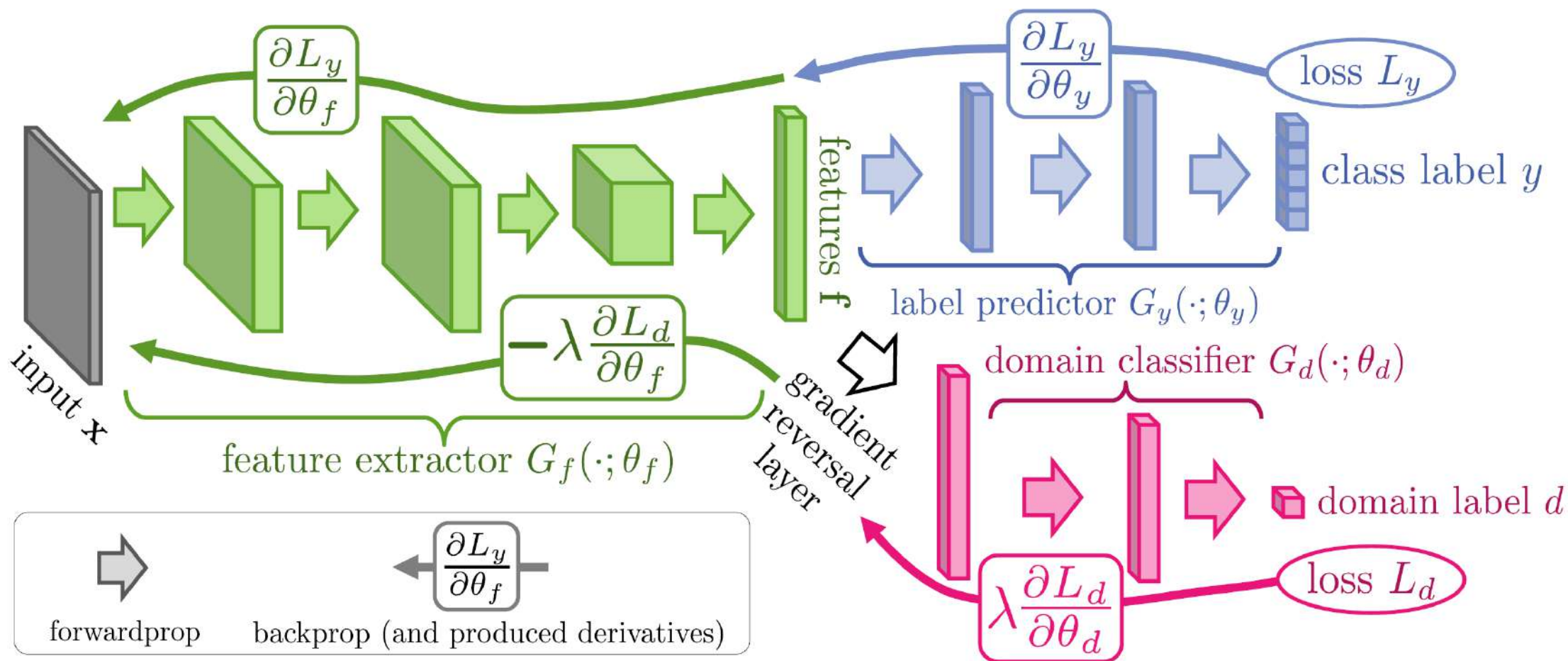


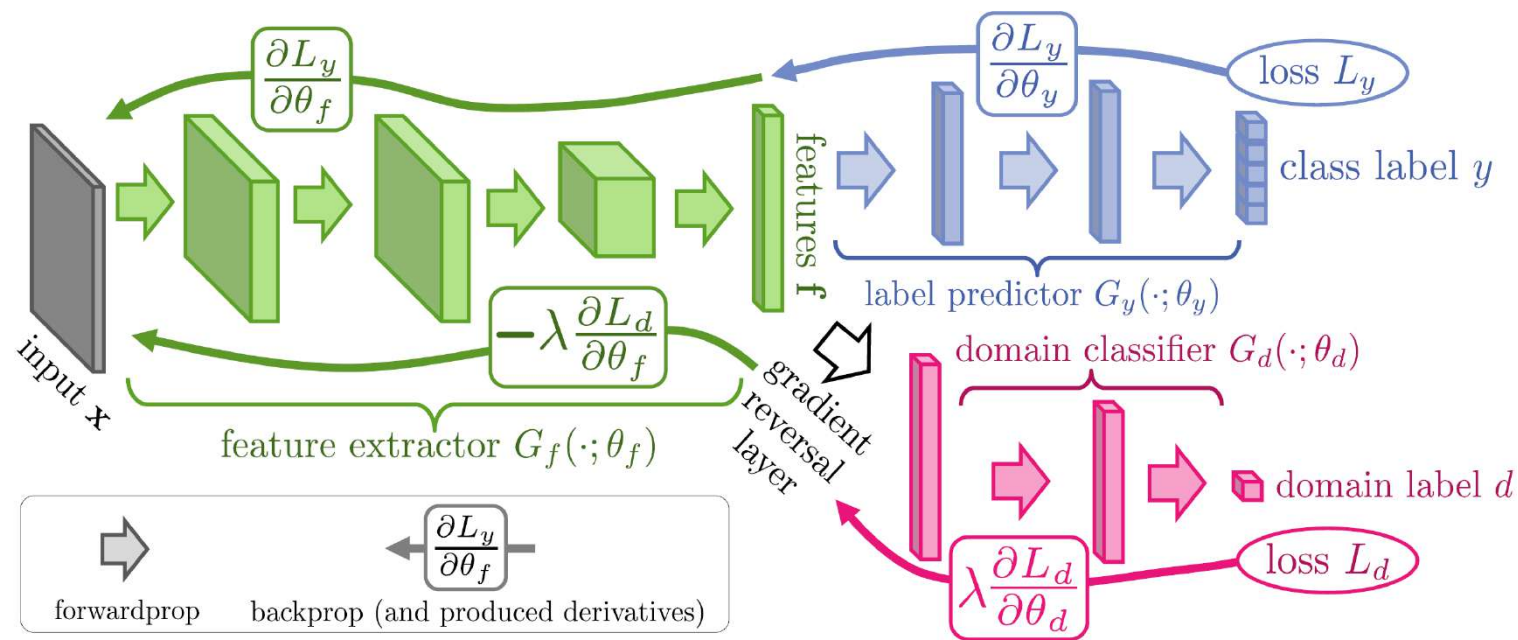
- Labelled data are expensive or infeasible
- Distribution “shift” between training and test data sets
e.g. synthetic or semi-synthetic (image) data – different data distribution
- Different types of training data
e.g. movie vs. books, MR vs. CT
- Y. Gannin et al. 2016: Unsupervised domain adaptation
 - Find a classifier:
 - a) Discriminateness (the ability to classify)
 - b) Domain-invariance (the ability to use ONLY domain-independent features to classify)
 - By:
 - 1) Label classifier (to predict the correct labels)
 - 2) Domain classifier (to find the features that can discriminate between source and target domain and NOT use them)

- Find a classifier:
 - a) Discriminativeness (the ability to classify)
 - b) Domain-invariance (the ability to use ONLY domain-independent features to classify)
- By:
 - 1) Label classifier (to predict the correct labels)
 - 2) Domain classifier (to find the features that can discriminate between source and target domain and NOT use them)
- That is:
 - 1) Minimise the loss of the label classifier
 - 2) Maximise the loss of the domain classifier
 - 3) Minimise the loss of domain classifier to train the domain classifier



Adversarial learning





$$(\hat{\theta}_f, \hat{\theta}_y) = \underset{\theta_f, \theta_y}{\operatorname{argmin}} E(\theta_f, \theta_y, \hat{\theta}_d),$$

$$\hat{\theta}_d = \underset{\theta_d}{\operatorname{argmax}} E(\hat{\theta}_f, \hat{\theta}_y, \theta_d).$$

$$\theta_f \leftarrow \theta_f - \mu \left(\frac{\partial \mathcal{L}_y^i}{\partial \theta_f} - \lambda \frac{\partial \mathcal{L}_d^i}{\partial \theta_f} \right),$$

$$\theta_y \leftarrow \theta_y - \mu \frac{\partial \mathcal{L}_y^i}{\partial \theta_y},$$

$$\theta_d \leftarrow \theta_d - \mu \lambda \frac{\partial \mathcal{L}_d^i}{\partial \theta_d},$$

$$\mathcal{L}_y^i(\theta_f, \theta_y) = \mathcal{L}_y(G_y(G_f(\mathbf{x}_i; \theta_f); \theta_y), y_i),$$

$$\mathcal{L}_d^i(\theta_f, \theta_d) = \mathcal{L}_d(G_d(G_f(\mathbf{x}_i; \theta_f); \theta_d), d_i).$$

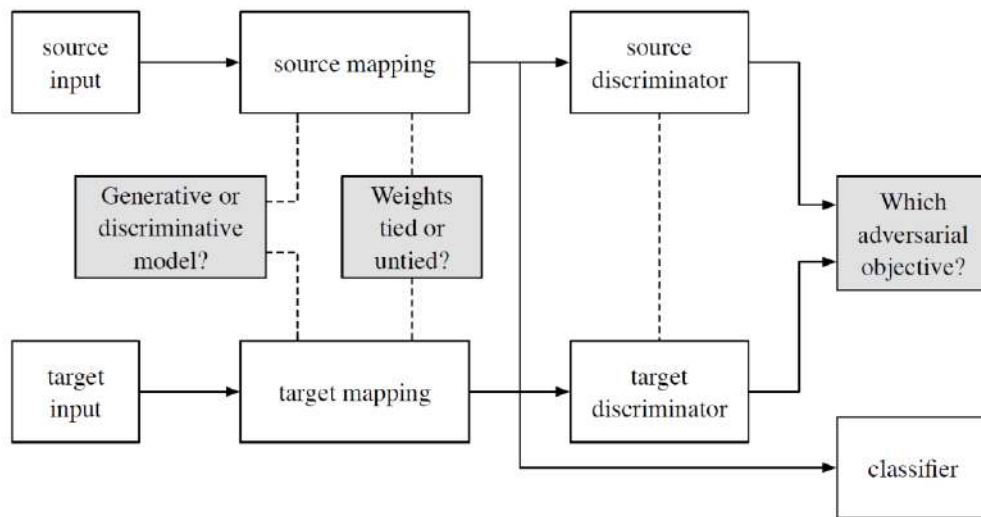
$$E(\theta_f, \theta_y, \theta_d) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_y^i(\theta_f, \theta_y) - \lambda \left(\frac{1}{n} \sum_{i=1}^n \mathcal{L}_d^i(\theta_f, \theta_d) + \frac{1}{n'} \sum_{i=n+1}^N \mathcal{L}_d^i(\theta_f, \theta_d) \right),$$



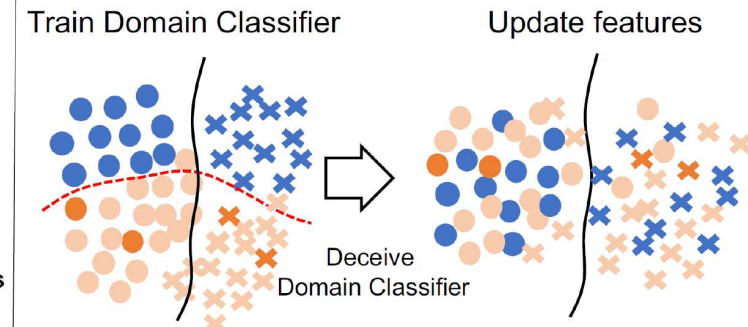
Distance between source and target distribution, e.g. H-divergence ©

$$\hat{d}_{\mathcal{H}}(S(G_f), T(G_f)) = 2 \left(1 - \min_{\eta \in \mathcal{H}} \left[\frac{1}{n} \sum_{i=1}^n I[\eta(G_f(\mathbf{x}_i))=0] + \frac{1}{n'} \sum_{i=n+1}^N I[\eta(G_f(\mathbf{x}_i))=1] \right] \right).$$

$$\mathcal{L}_d^i(\mathbf{W}, \mathbf{b}, \mathbf{u}, z) = \mathcal{L}_d(G_d(G_f(\mathbf{x}_i; \mathbf{W}, \mathbf{b}); \mathbf{u}, z), d_i)$$



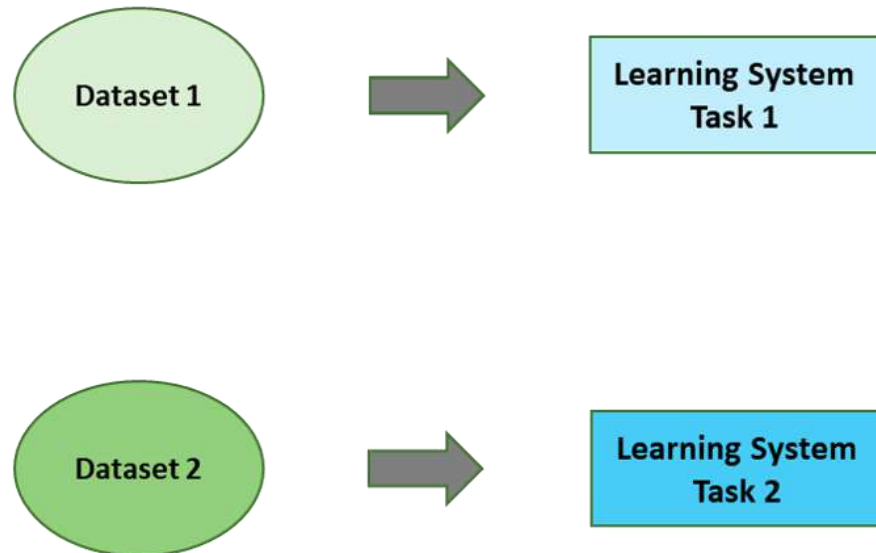
- Task-specific Classifier
- - - Domain Classifier
- × Labeled Source
- × Labeled Target
- × Unlabeled Target
- × Estimated Prototypes



Research | Meta-Learning

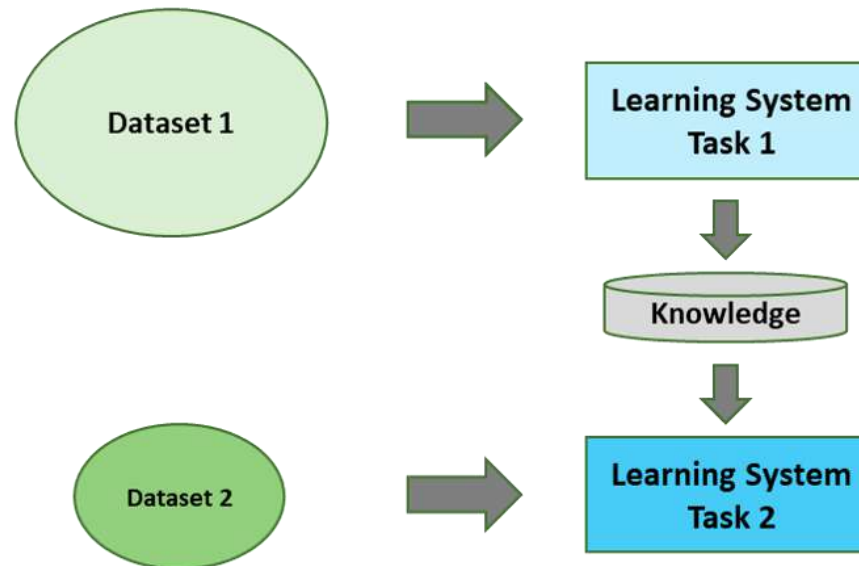
Transfer learning

Traditional ML

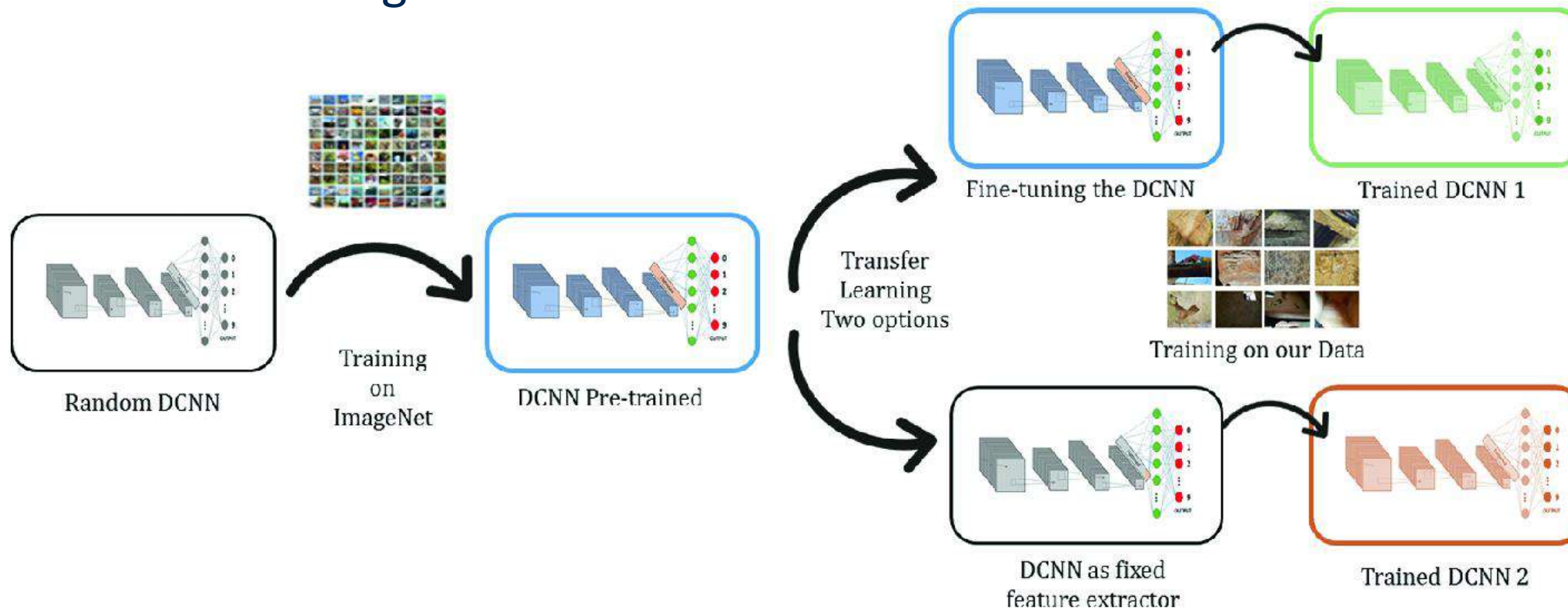


vs

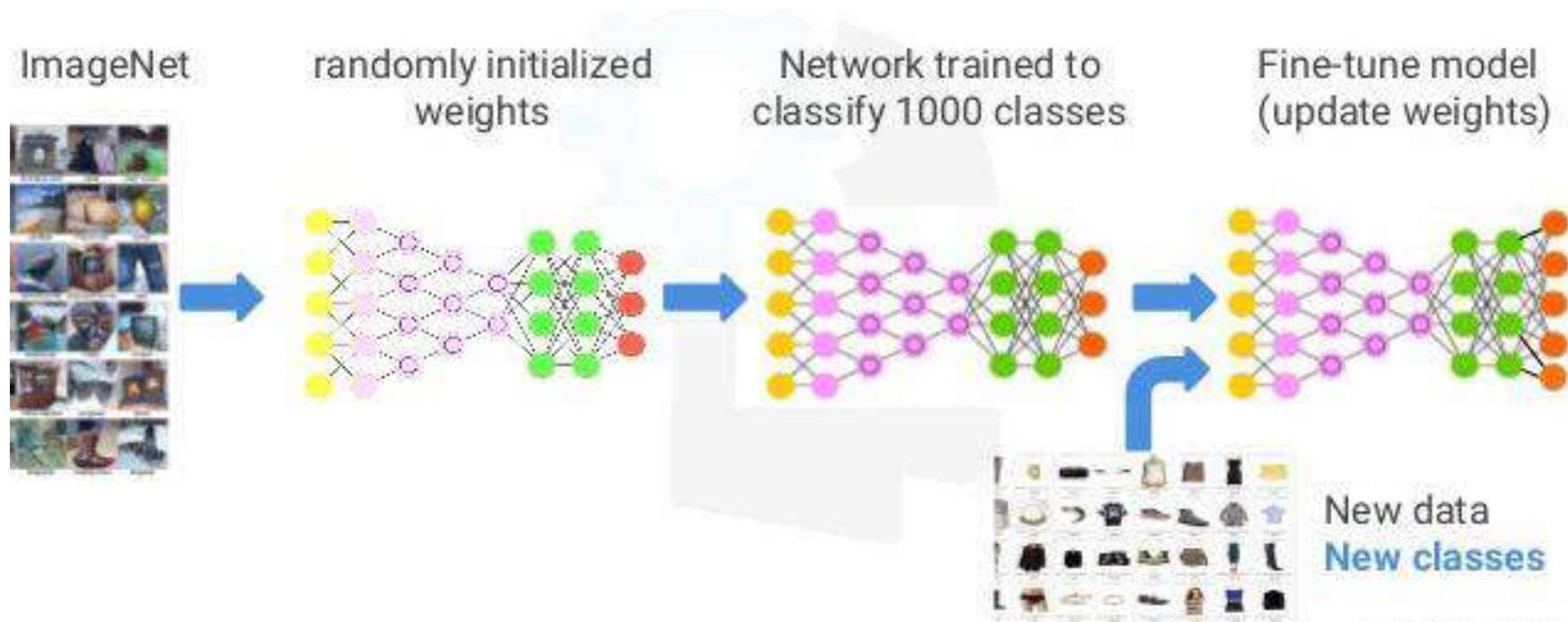
Transfer Learning



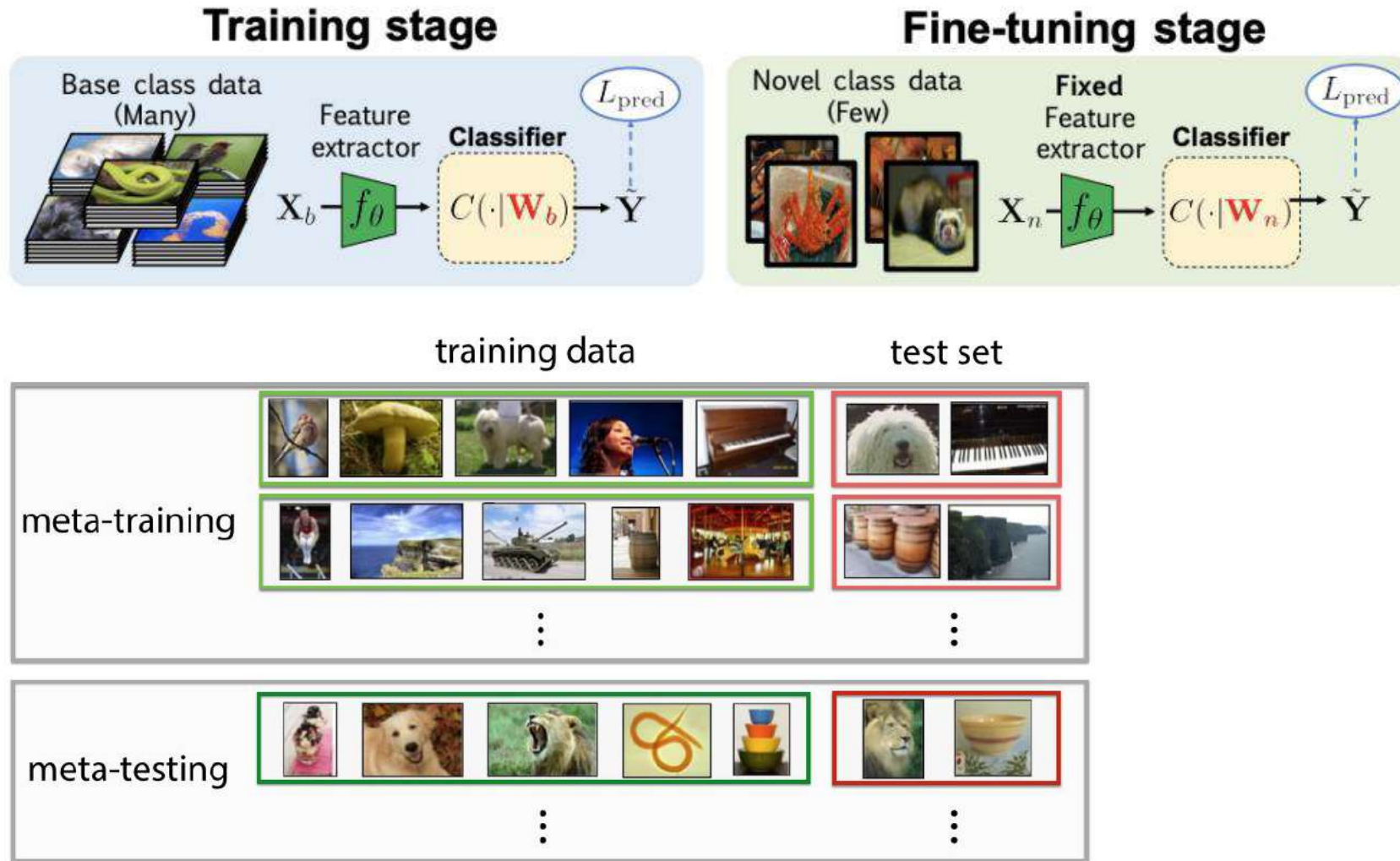
Transfer learning



Transfer learning for few-shot learning



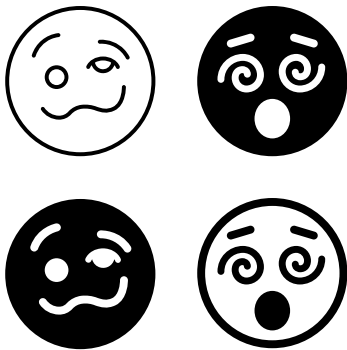
Few-shot learning and meta-learning



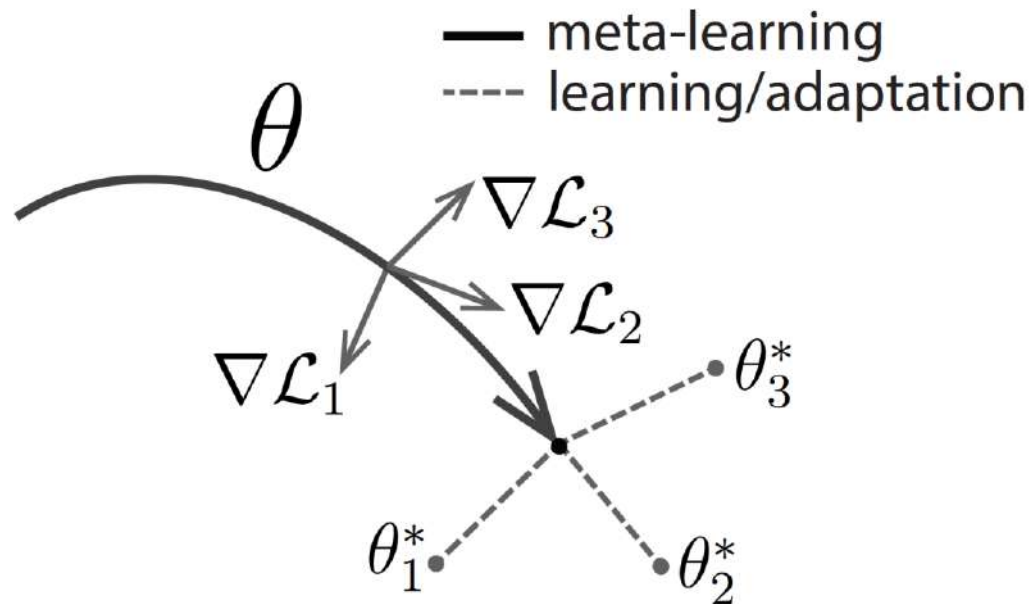
Terminology and taxonomy

Learning Settings		Source and Target Domains	Source and Target Tasks
Traditional Machine Learning		the same	the same
Transfer Learning	<i>Inductive Transfer Learning /</i>	the same	different but related
	<i>Unsupervised Transfer Learning</i>	different but related	different but related
	<i>Transductive Transfer Learning</i>	different but related	the same

Transfer Learning Settings	Related Areas	Source Domain Labels	Target Domain Labels	Tasks
<i>Inductive Transfer Learning</i>	Multi-task Learning	Available	Available	Regression, Classification
	Self-taught Learning	Unavailable	Available	Regression, Classification
<i>Transductive Transfer Learning</i>	Domain Adaptation, Sample Selection Bias, Co-variate Shift	Available	Unavailable	Regression, Classification
<i>Unsupervised Transfer Learning</i>		Unavailable	Unavailable	Clustering, Dimensionality Reduction



Model-agnostic meta-learning (MAML)



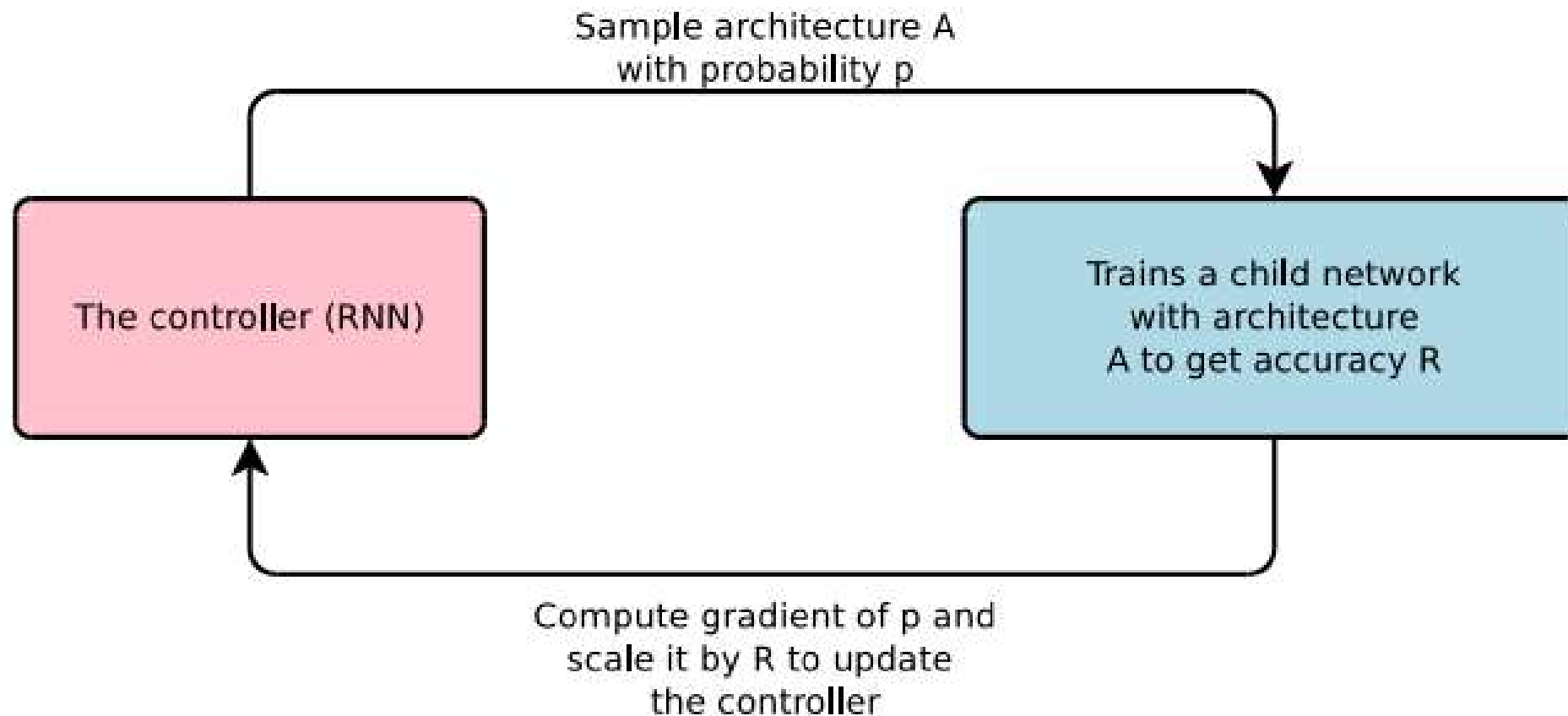
Algorithm 1 Model-Agnostic Meta-Learning

Require: $p(\mathcal{T})$: distribution over tasks

Require: α, β : step size hyperparameters

- 1: randomly initialize θ
 - 2: **while** not done **do**
 - 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
 - 4: **for all** \mathcal{T}_i **do**
 - 5: Evaluate $\nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$ with respect to K examples
 - 6: Compute adapted parameters with gradient descent: $\theta'_i = \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
 - 7: **end for**
 - 8: Update $\theta \leftarrow \theta - \beta \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$
 - 9: **end while**
-

Searching for hyper-parameters (architecture, loss, data augmentation...)



“Transferability” from computer vision to medical imaging

- Applications:
 - Encoder-decoder U-Net - segmentation
 - Optical flow – registration
 - GANs – synthesis
 - ...
- Data size
- Data availability
- Data variability, equipment, protocols, demography etc.
- Ground-truth, inc. label uncertainty, fidelity etc.
- Test significance, e.g. 1% effect size, 20 subjects, ~10% variance
- Clinical relevance