

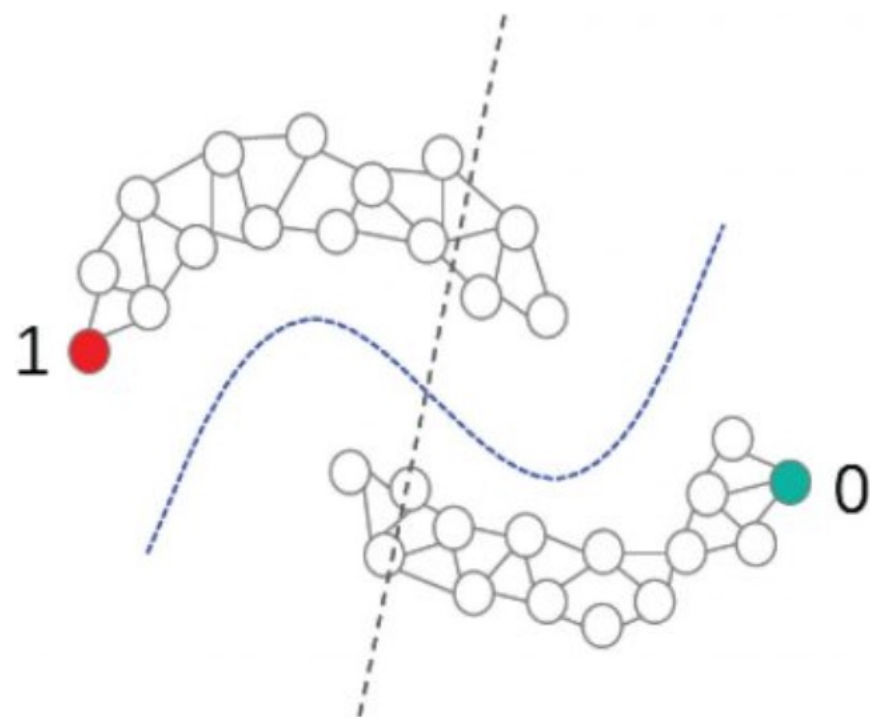
Open Problems in Deep Learning

Neural Networks Design And Application

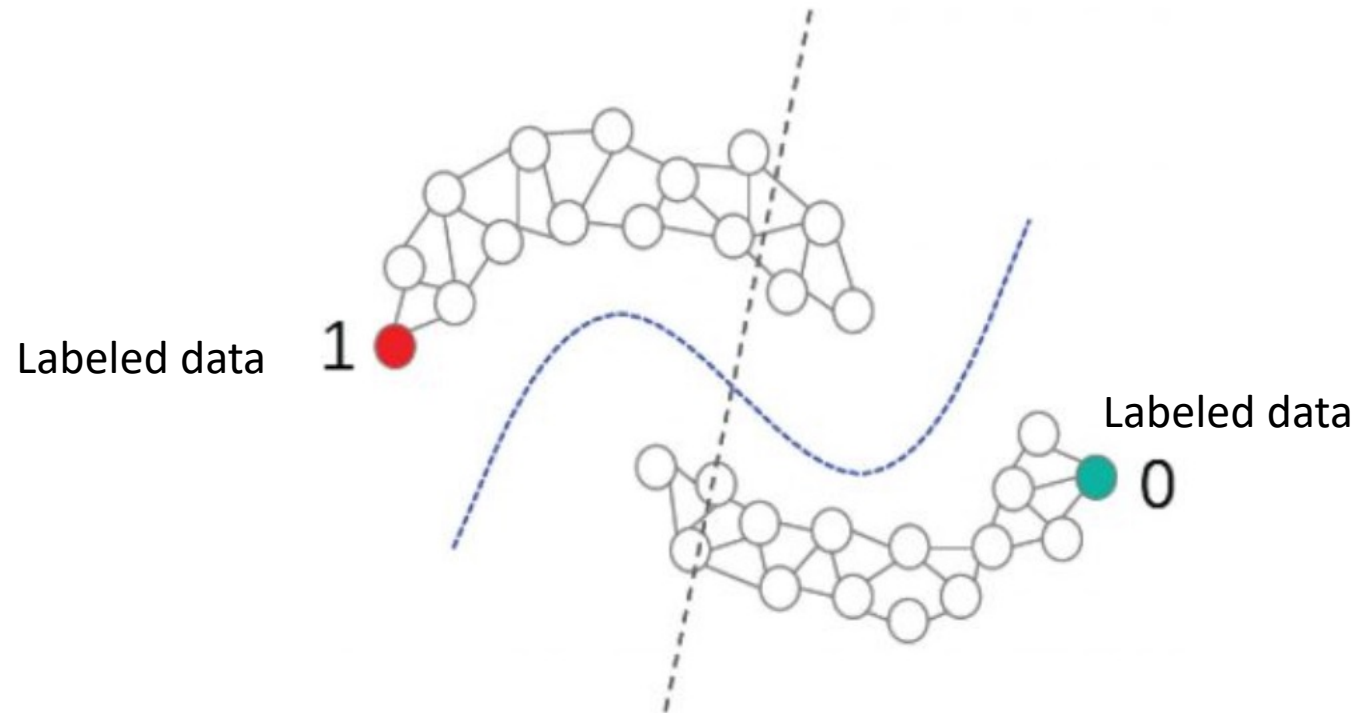
Open problems in deep learning

- Weak supervised learning
- Causal modeling
- Curriculum learning
- Over-parameterized modeling
- Self-supervised learning
- Meta-learning
- Federated learning
-

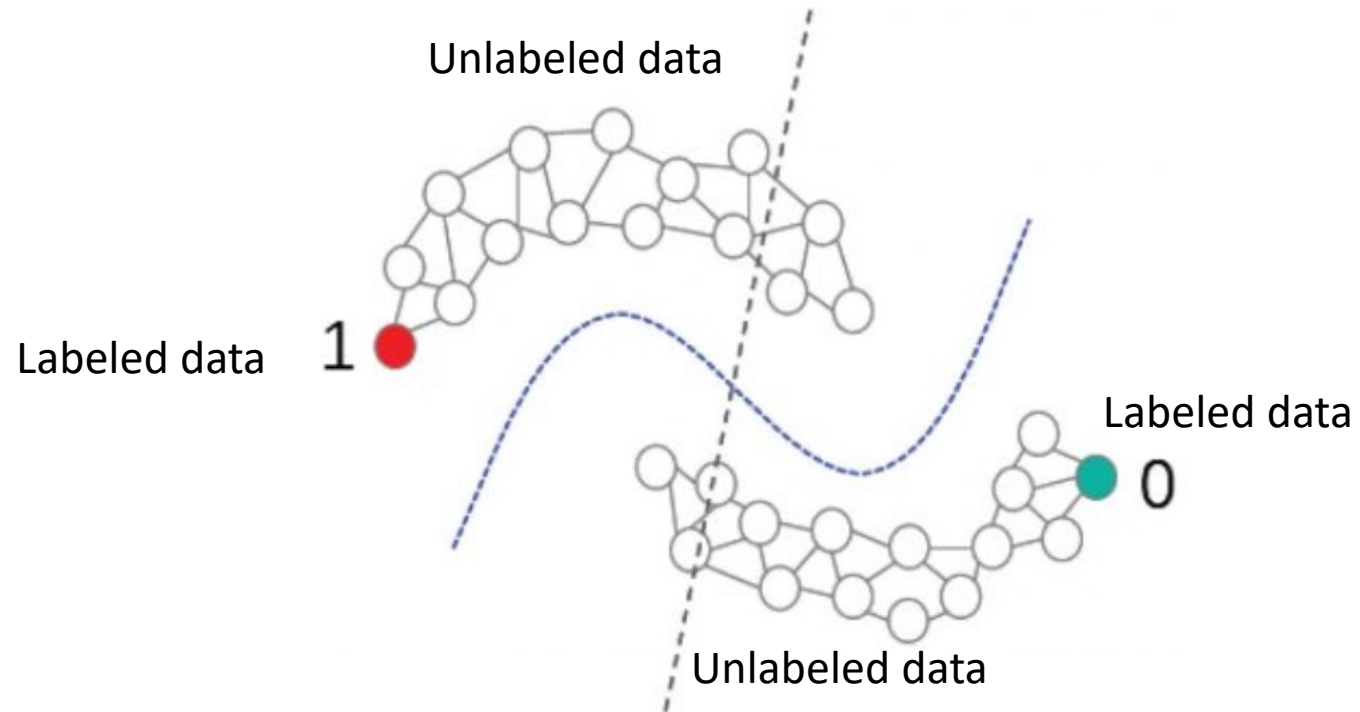
Weak supervision



Weak supervision

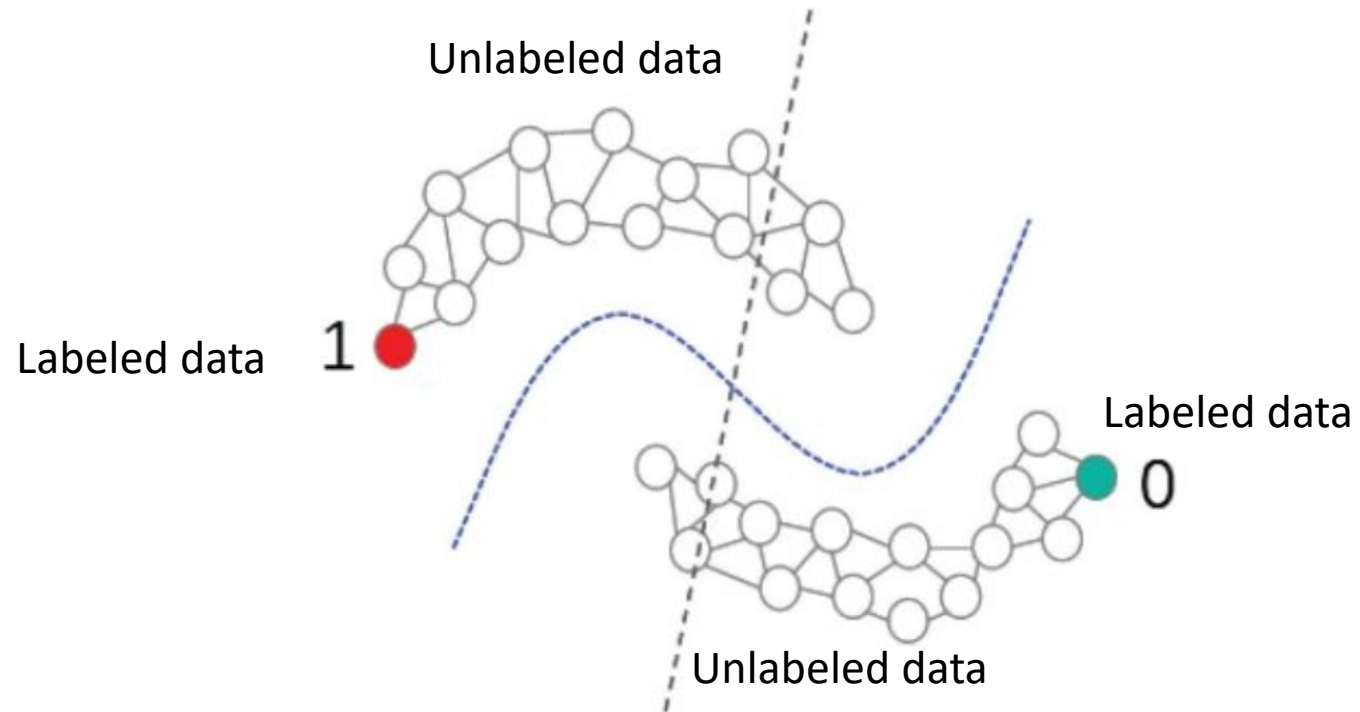


Weak supervision

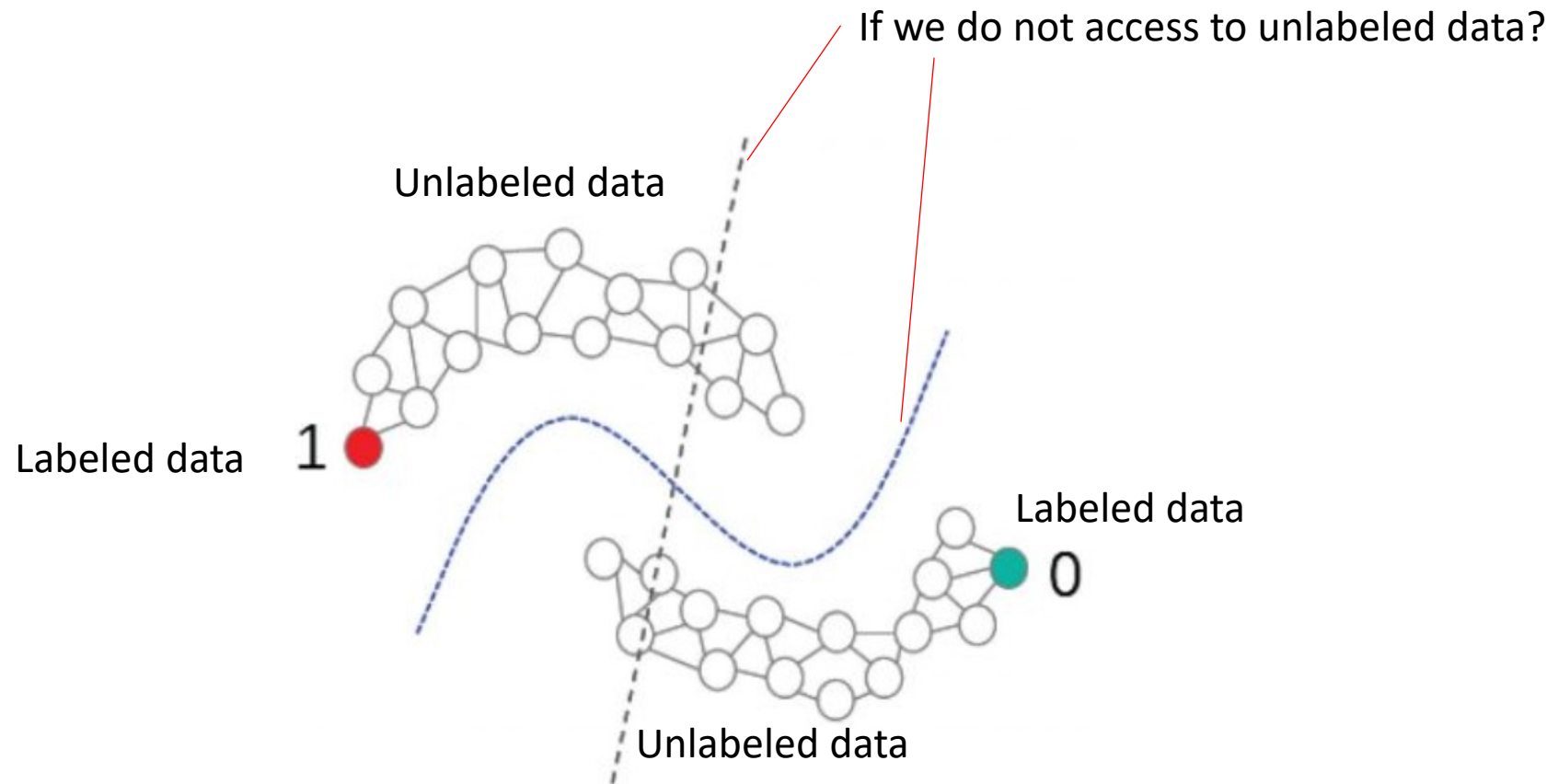


Weak supervision

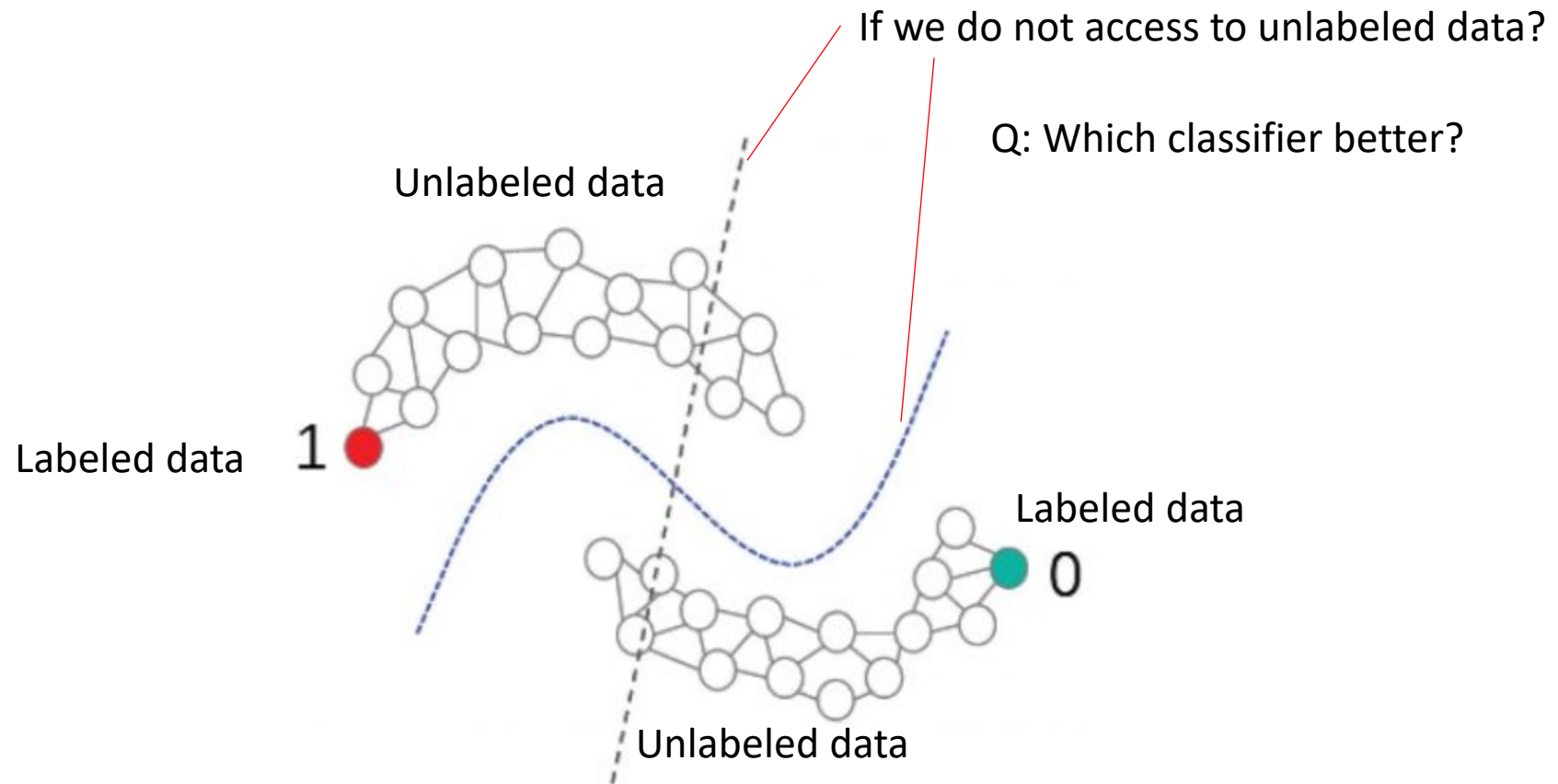
If we do not access to unlabeled data?



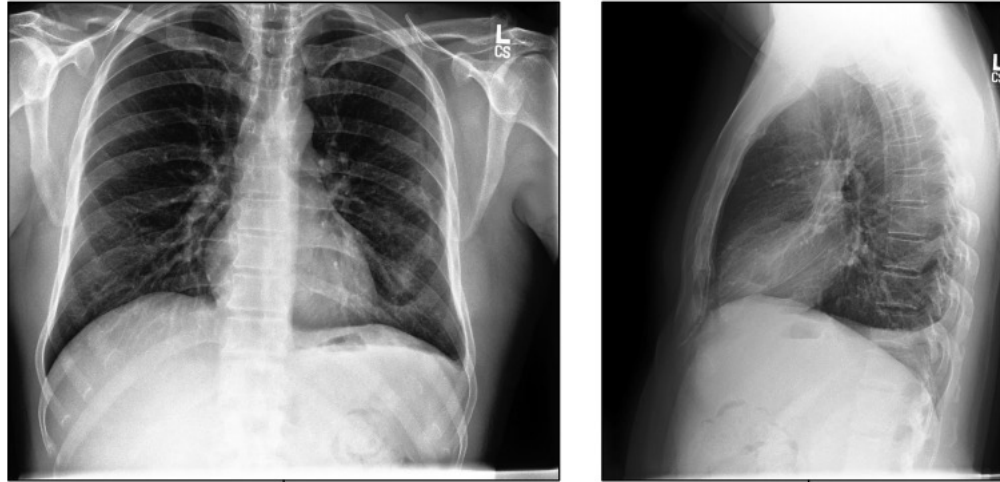
Weak supervision



Weak supervision



Weak supervision



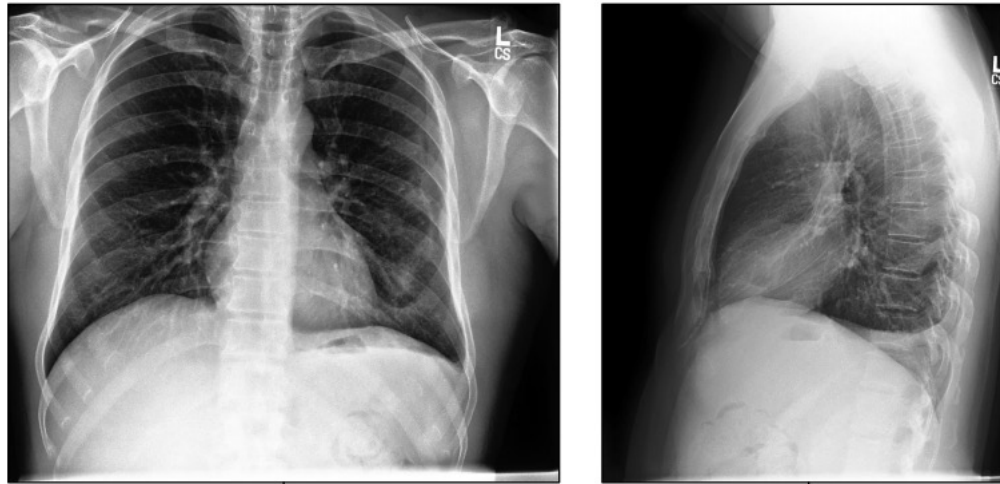
X-ray image analysis

Irvin, Jeremy, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilicus, Chris Chute, Henrik Marklund et al. "Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison." In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, pp. 590-597. 2019.

<https://arxiv.org/pdf/1901.07031.pdf>

Weak supervision

Some classification tasks are very difficult to acquire a lot of **labeled** data



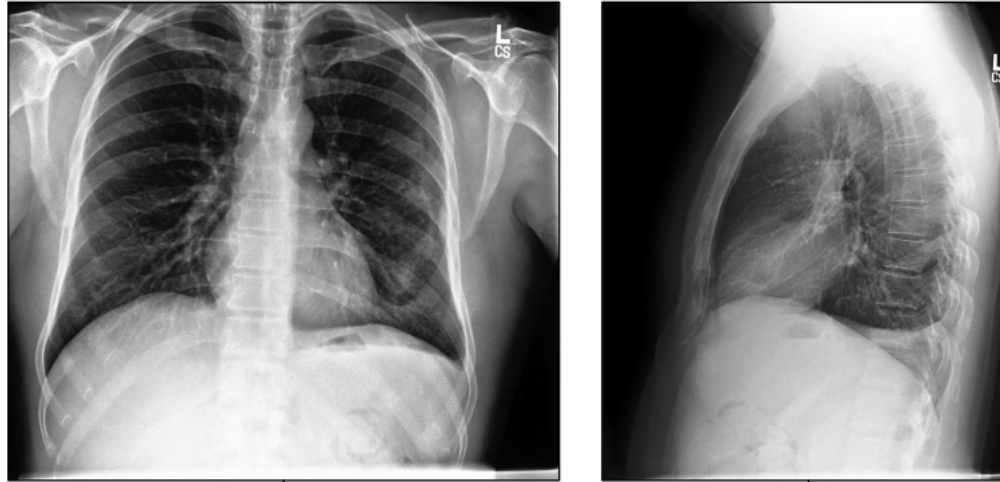
X-ray image analysis

Irvin, Jeremy, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund et al. "Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison." In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, pp. 590-597. 2019.

<https://arxiv.org/pdf/1901.07031.pdf>

Weak supervision

Some classification tasks are very difficult to acquire a lot of **labeled** data, but may provide many **unlabeled** data



X-ray image analysis

Irvin, Jeremy, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund et al. "Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison." In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, pp. 590-597. 2019.

<https://arxiv.org/pdf/1901.07031.pdf>

Weak supervision

Some classification tasks are very difficult to acquire a lot of **labeled** data, but may provide many **unlabeled** data



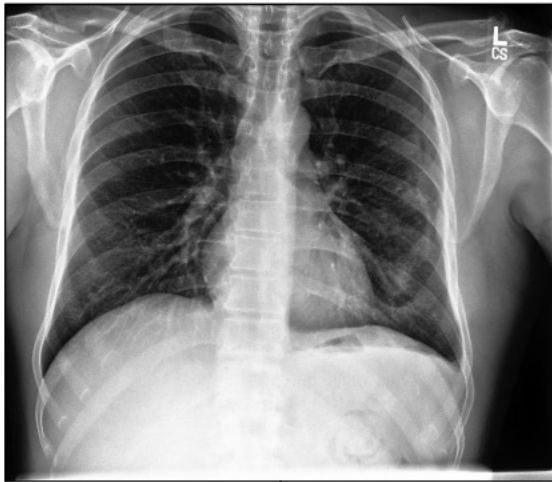
X-ray image analysis

Identify: which x-ray image may imply **life threatening diseases**?

Irvin, Jeremy, Pranav Rajpurkar, Michael Ko, Yifan Yu, Silvana Ciurea-Ilcus, Chris Chute, Henrik Marklund et al. "Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison." In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, pp. 590-597. 2019.

<https://arxiv.org/pdf/1901.07031.pdf>

Weak supervision

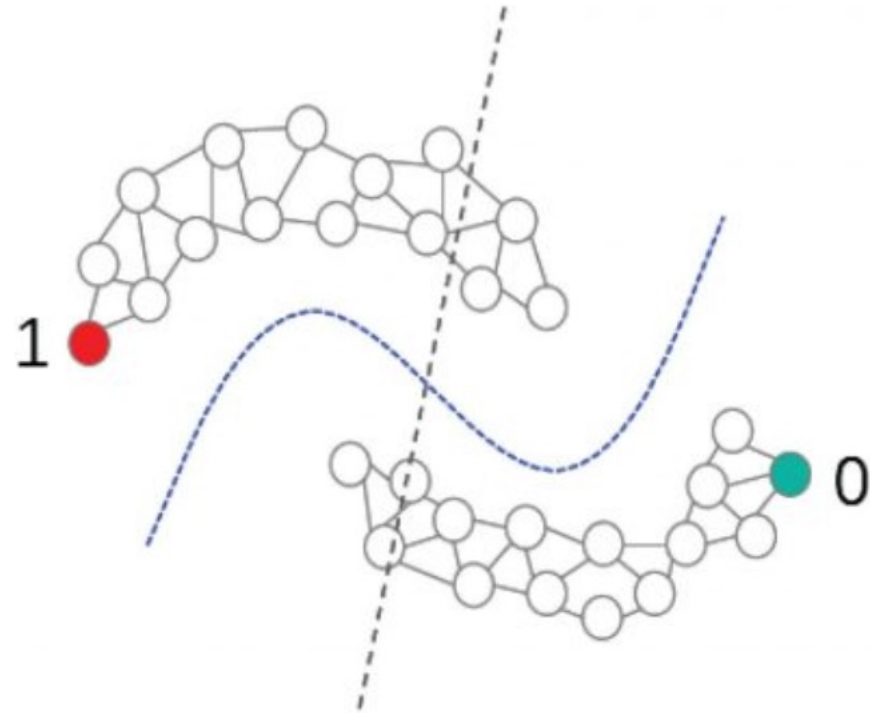


Leaderboard

Will your model perform as well as radiologists in detecting different pathologies in chest X-rays?

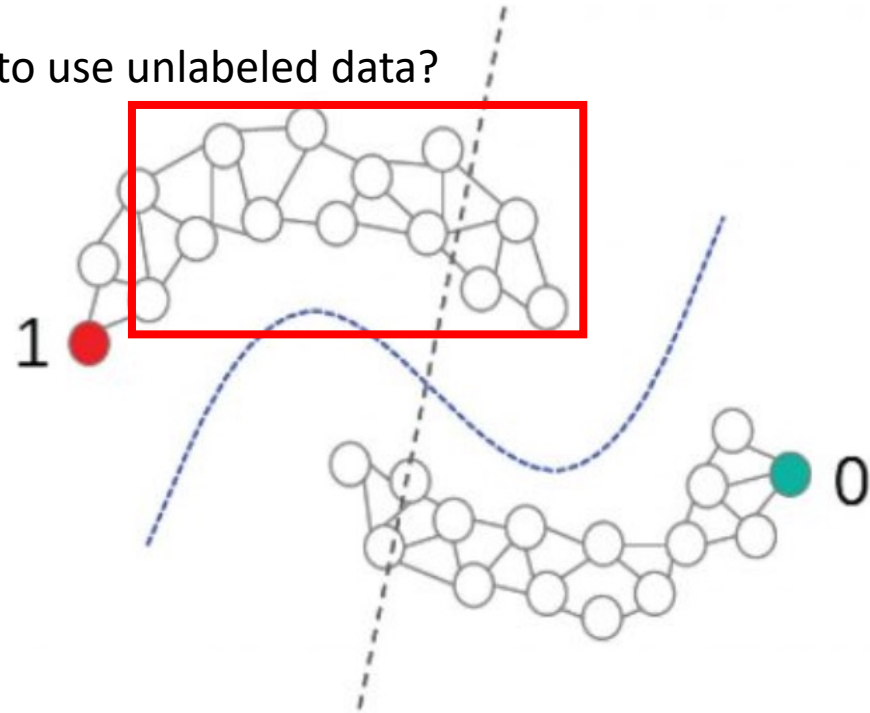
Rank	Date	Model	AUC	Num Rads Below Curve
1	Aug 31, 2020	DeepAUC-v1 <i>ensemble</i> https://arxiv.org/abs/2012.03173	0.930	2.8
2	Sep 01, 2019	Hierarchical-Learning-V1 (ensemble) <i>Vingroup</i> <i>Big Data Institute</i> https://arxiv.org/abs/1911.06475	0.930	2.6
3	Oct 15, 2019	Conditional-Training-LSR <i>ensemble</i>	0.929	2.6
4	Dec 04, 2019	Hierarchical-Learning-V4 (ensemble) <i>Vingroup</i>	0.929	2.6

Weak supervision



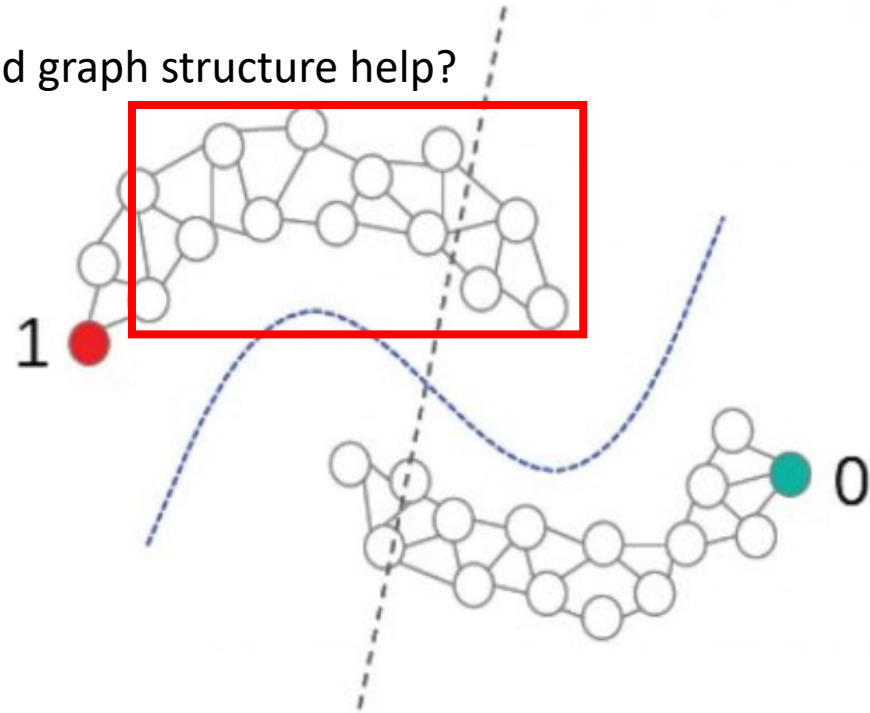
Weak supervision

Q: how to use unlabeled data?



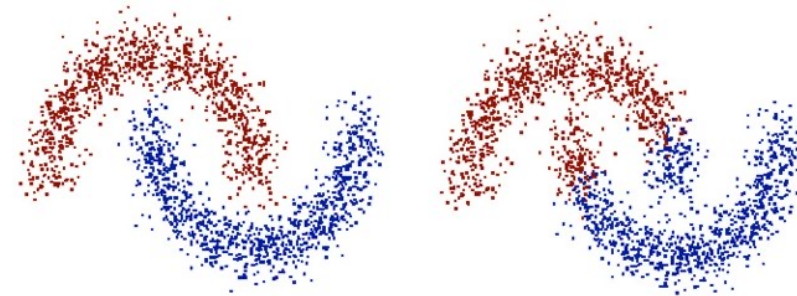
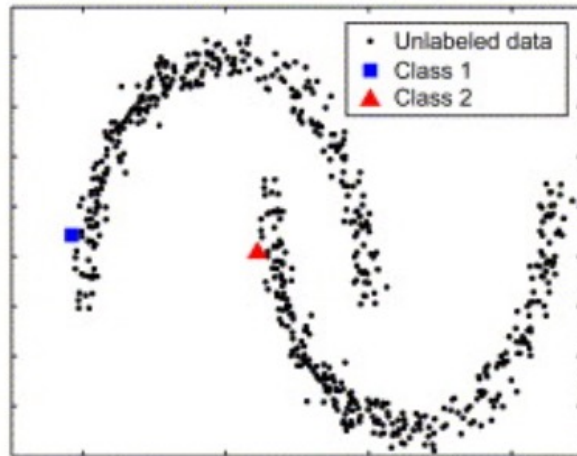
Weak supervision

Q: Would graph structure help?



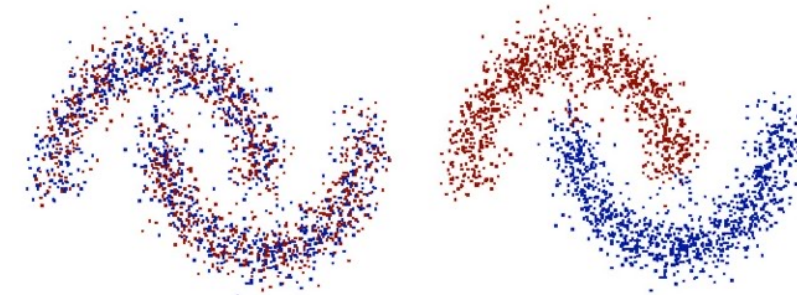
Weak supervision

Two moon dataset



(a) Solution

(b) Shi-Malik [24]

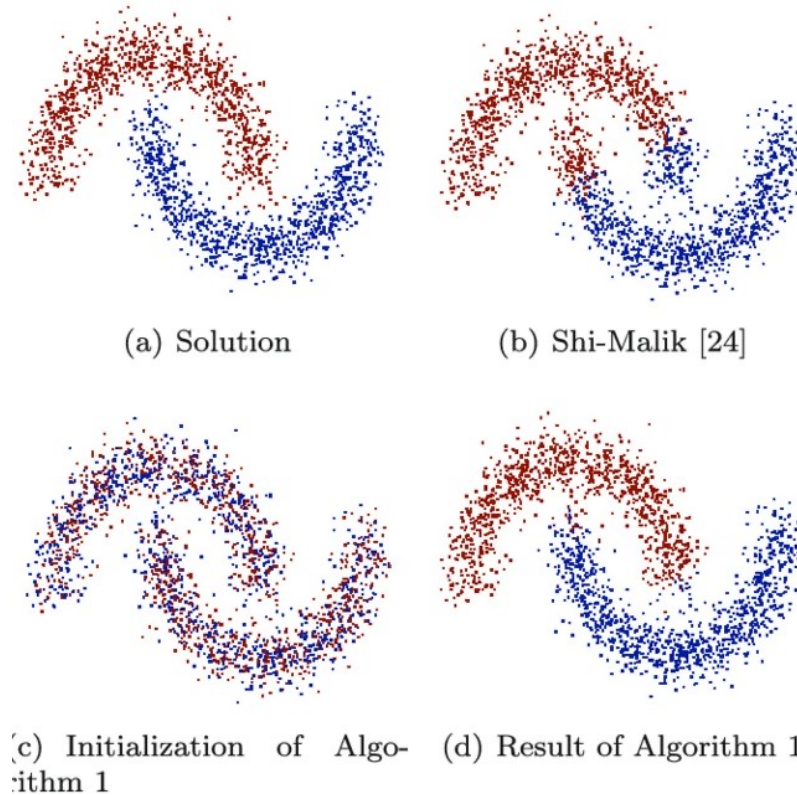
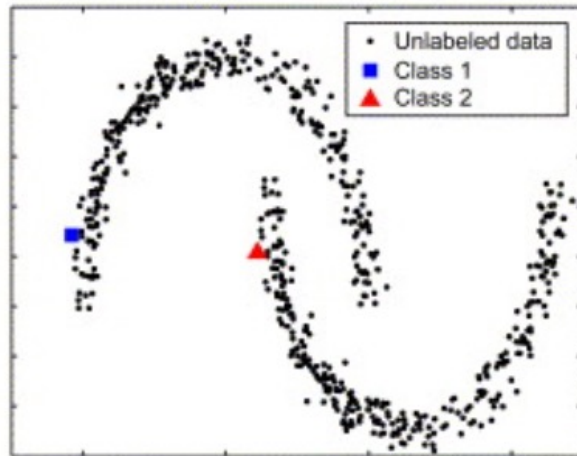


(c) Initialization of Algorithm 1

(d) Result of Algorithm 1

Weak supervision

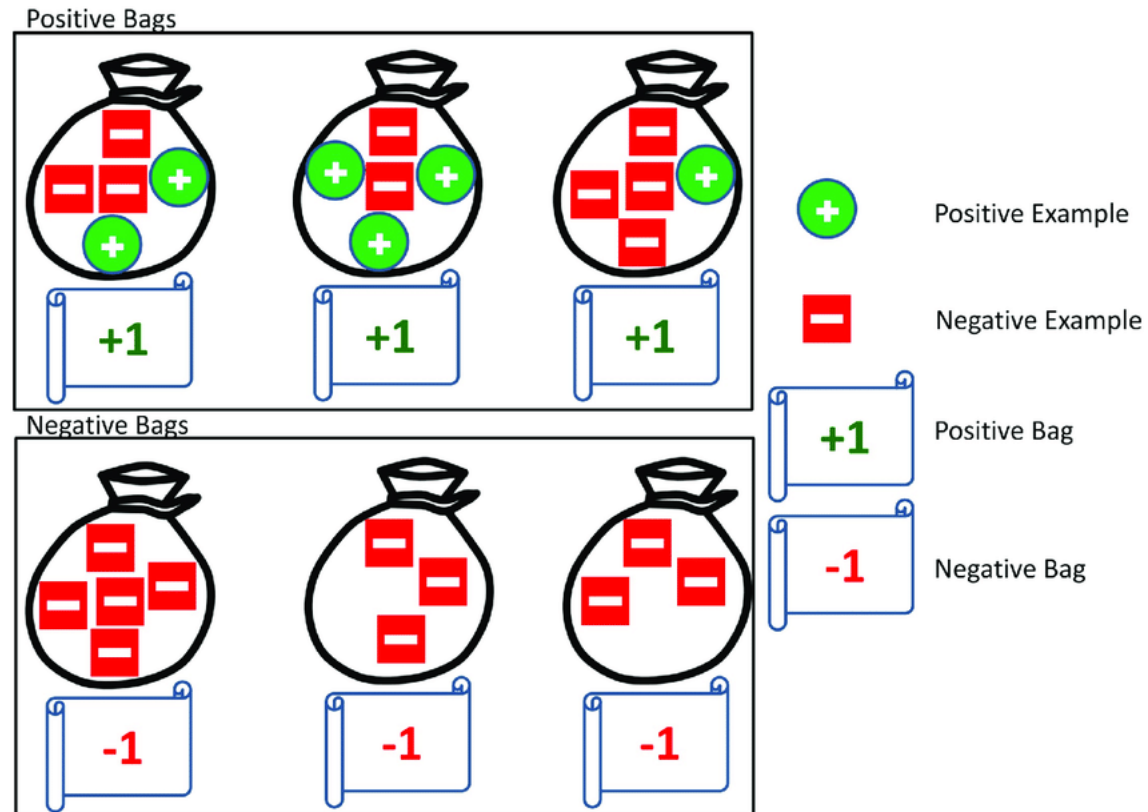
Two moon dataset



Labeled + unlabeled data → semi-supervised learning

Weak supervision

Multi-instance learning (another type of weak supervised learning)

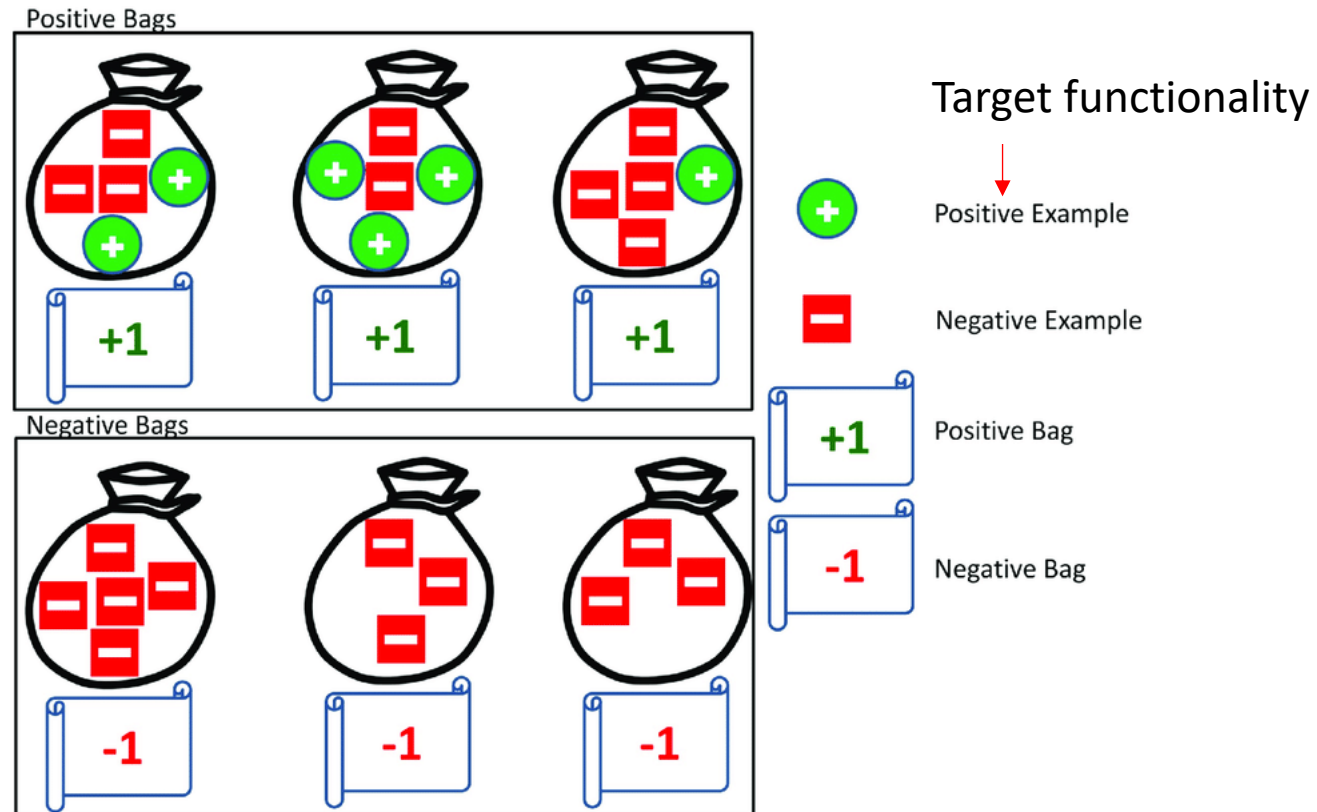


Amino acid composition predicts prion activity. April 2017

[PLoS Computational Biology](https://doi.org/10.1371/journal.pcbi.1005465) 13(4):e1005465. DOI:[10.1371/journal.pcbi.1005465](https://doi.org/10.1371/journal.pcbi.1005465)

Weak supervision

Multi-instance learning (another type of weak supervised learning)

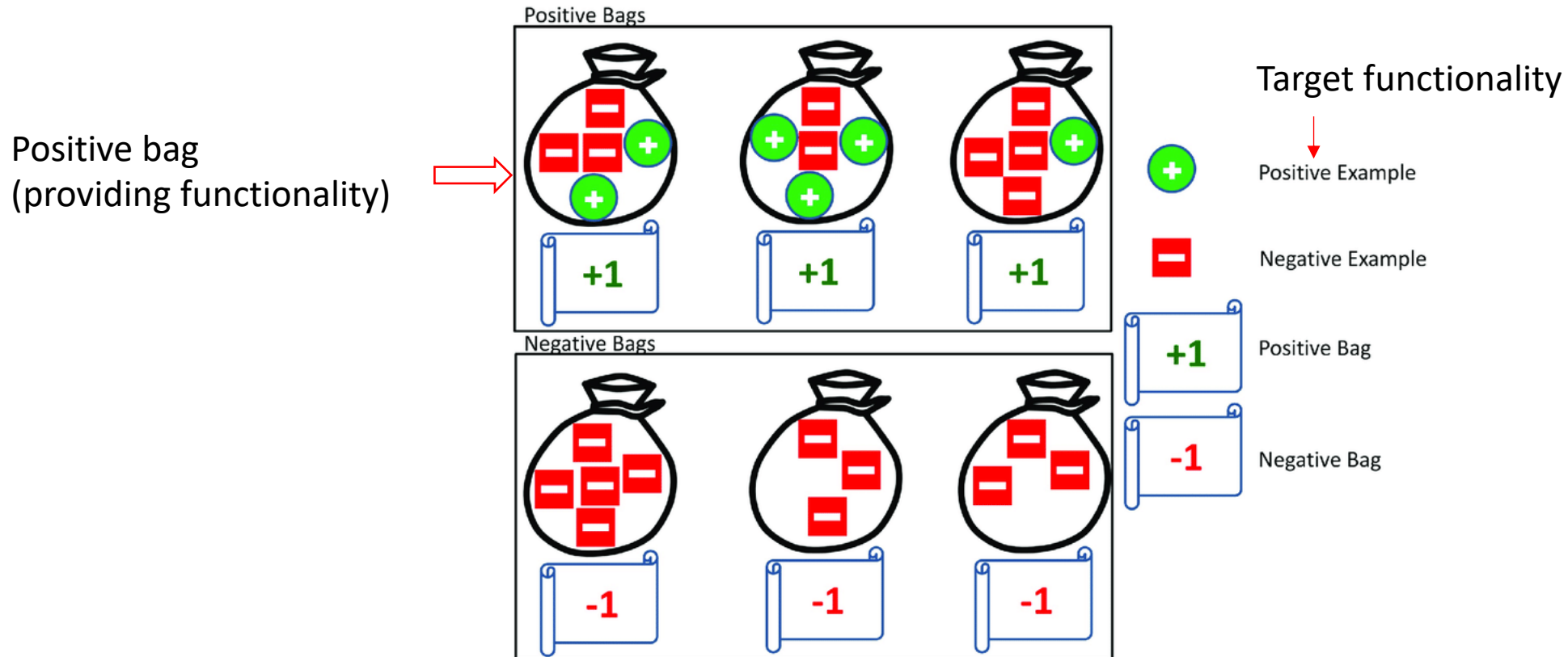


Amino acid composition predicts prion activity. April 2017

[PLoS Computational Biology](https://doi.org/10.1371/journal.pcbi.1005465) 13(4):e1005465. DOI:[10.1371/journal.pcbi.1005465](https://doi.org/10.1371/journal.pcbi.1005465)

Weak supervision

Multi-instance learning (another type of weak supervised learning)

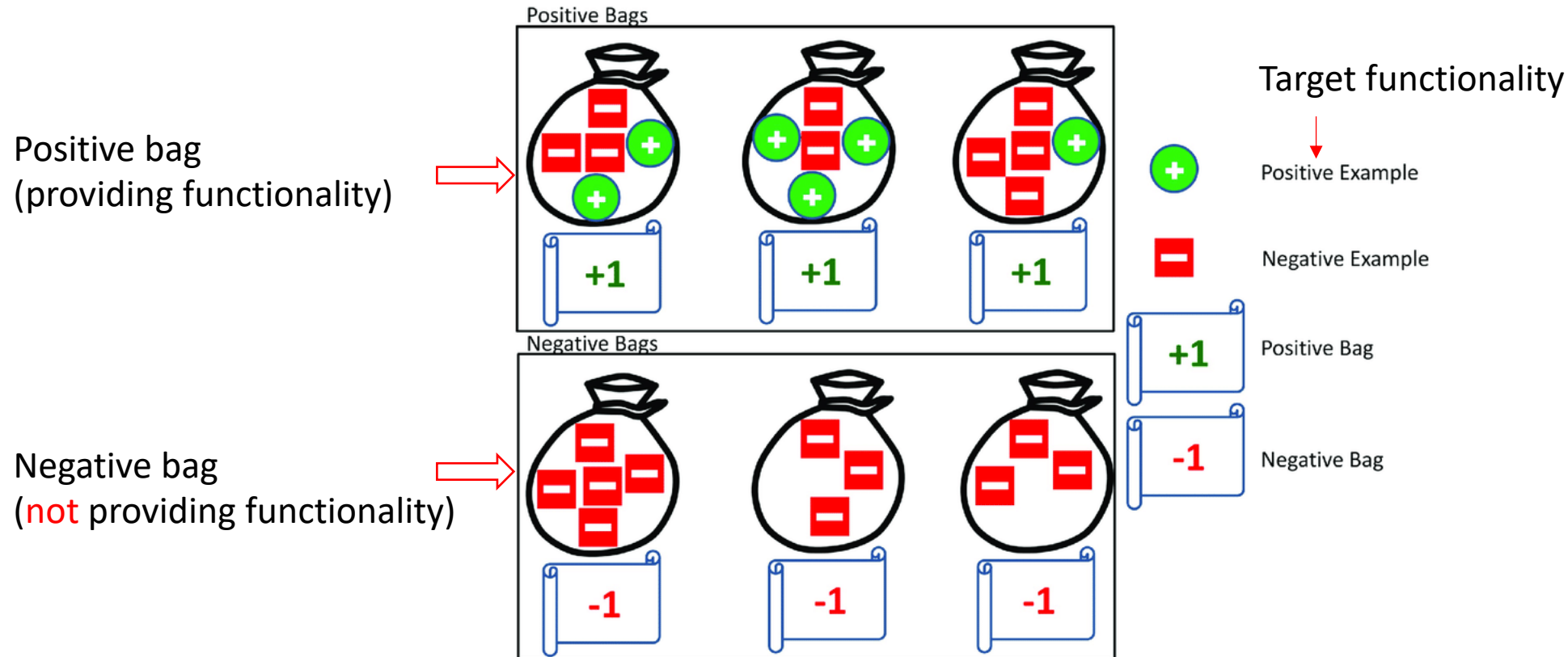


Amino acid composition predicts prion activity. April 2017

[PLoS Computational Biology](https://doi.org/10.1371/journal.pcbi.1005465) 13(4):e1005465. DOI:[10.1371/journal.pcbi.1005465](https://doi.org/10.1371/journal.pcbi.1005465)

Weak supervision

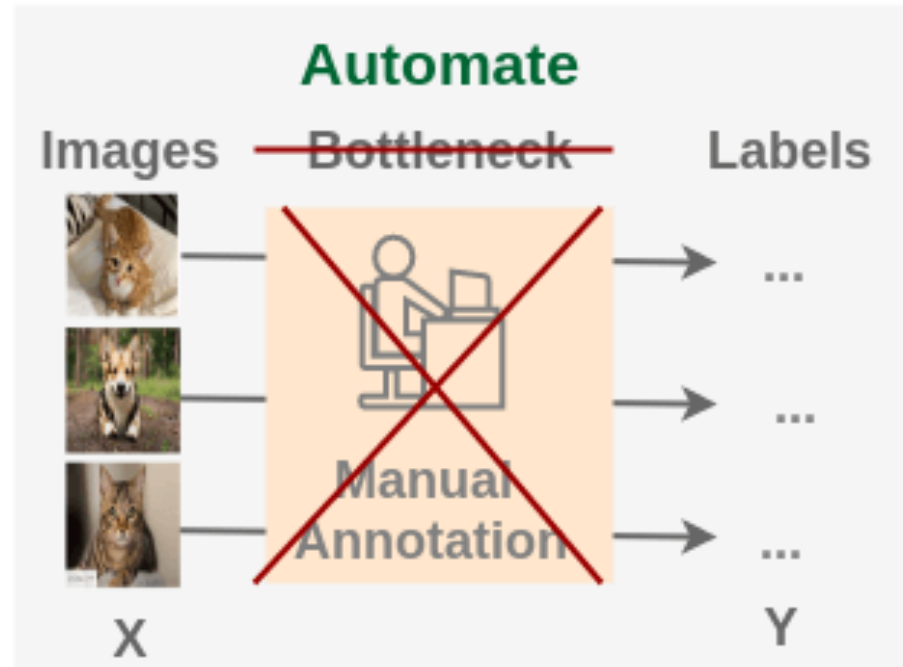
Multi-instance learning (another type of weak supervised learning)



Amino acid composition predicts prion activity. April 2017

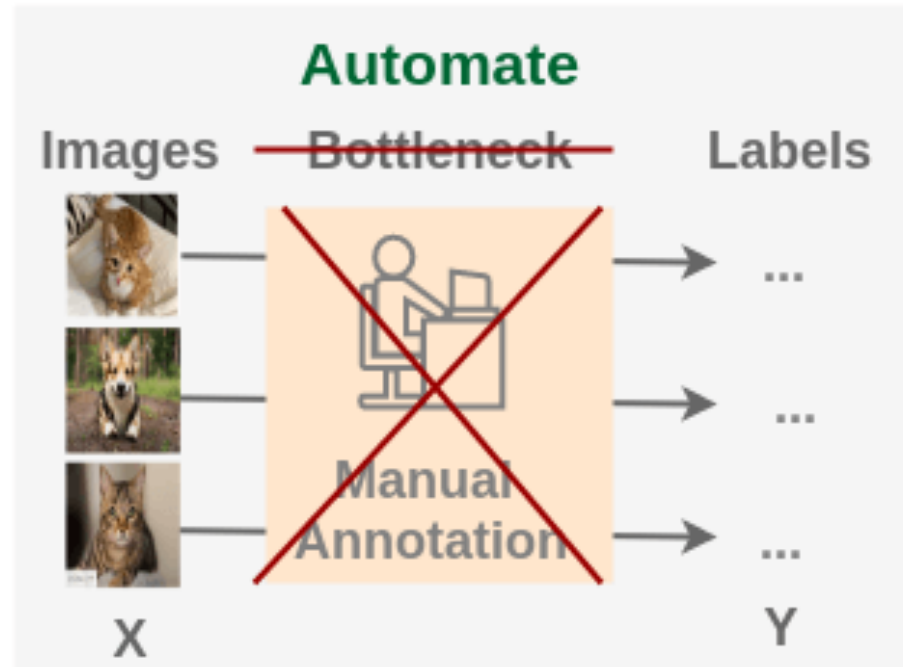
[PLoS Computational Biology](https://doi.org/10.1371/journal.pcbi.1005465) 13(4):e1005465. DOI:[10.1371/journal.pcbi.1005465](https://doi.org/10.1371/journal.pcbi.1005465)

Self-supervision



Without human effort for labeling

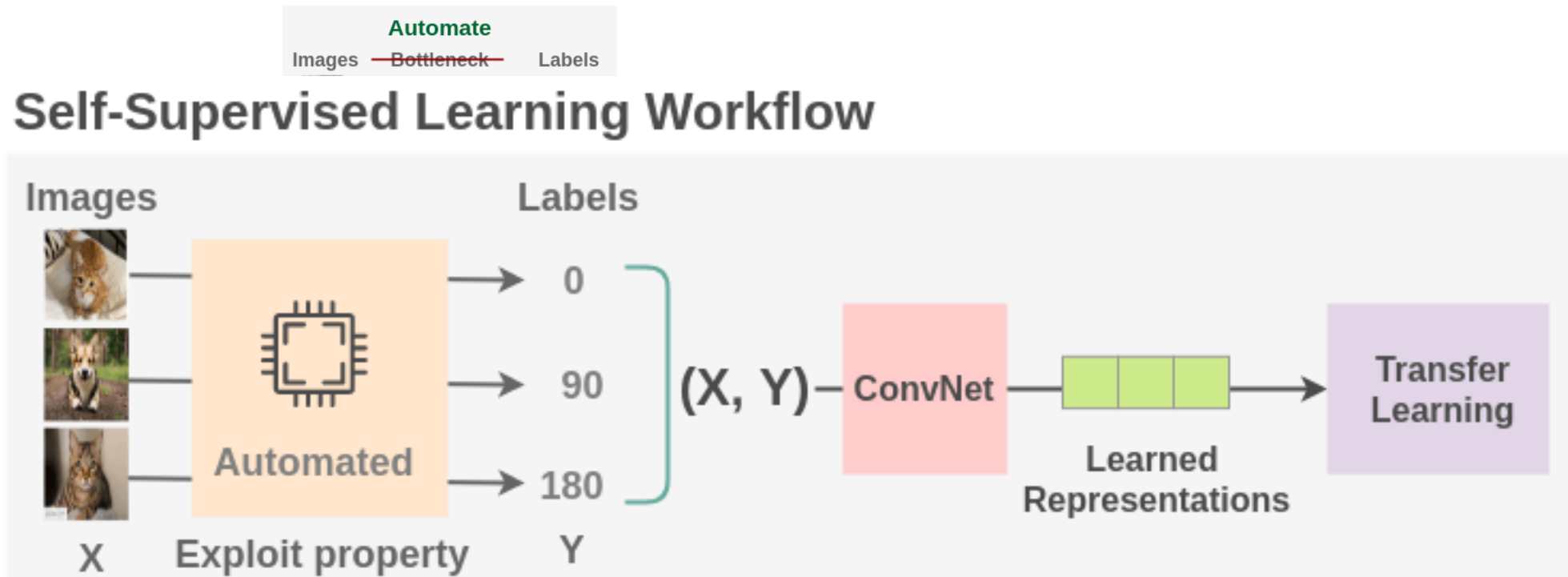
Self-supervision



Without human effort for labeling

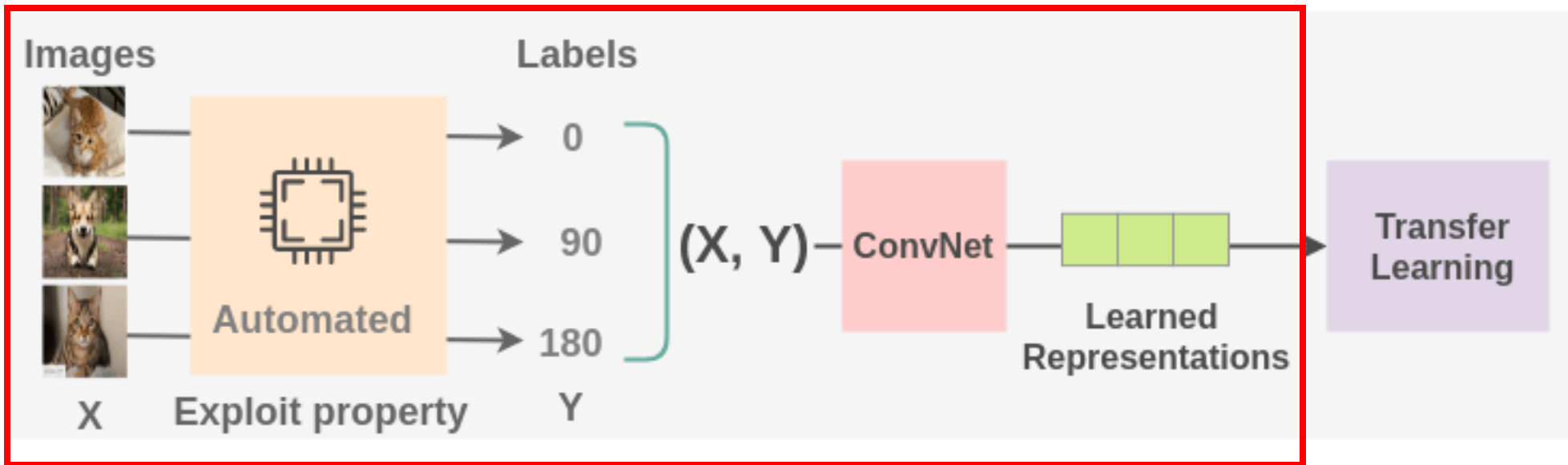
Q: can we automate labeling to get approximation to label information?

Self-supervision



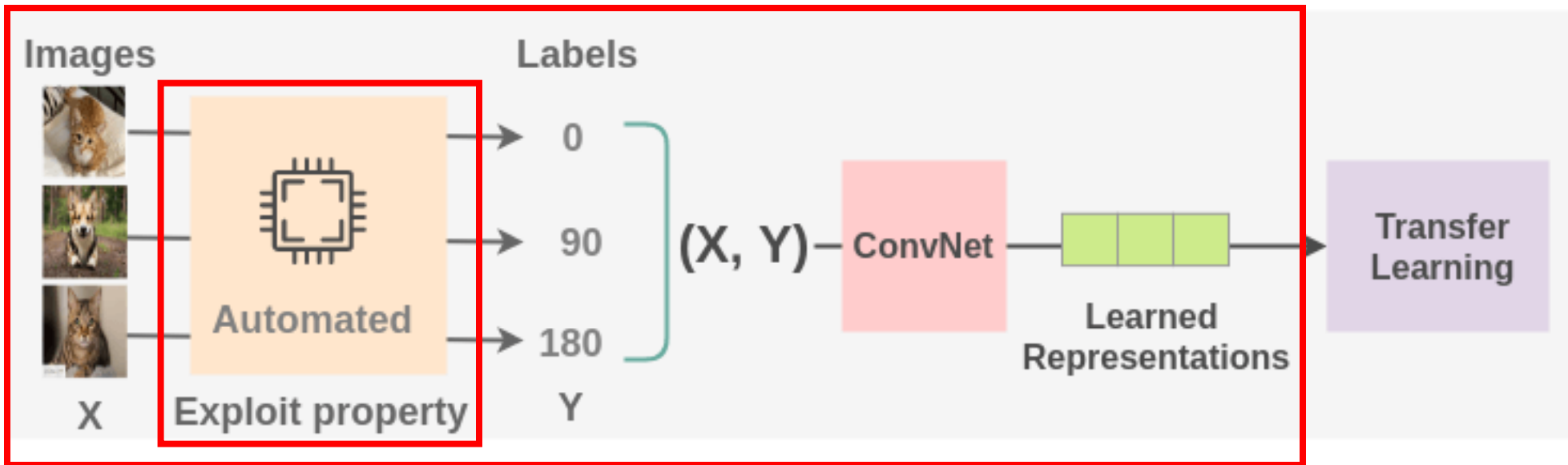
Self-supervision

Self-Supervised Learning Workflow



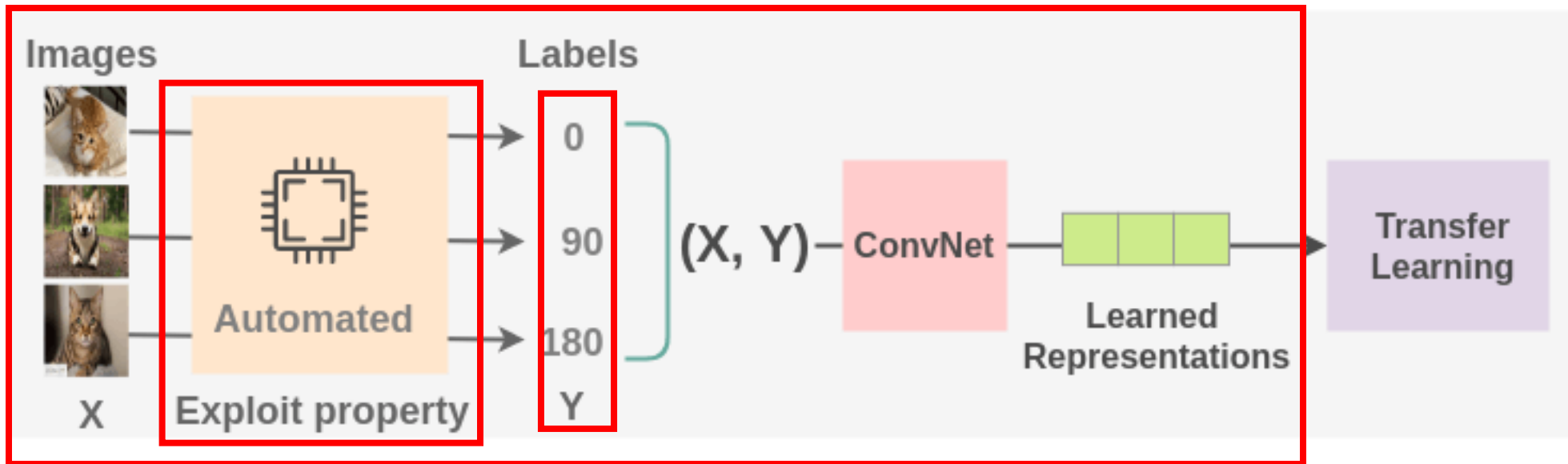
Self-supervision

Self-Supervised Learning Workflow



Self-supervision

Self-Supervised Learning Workflow



Self-supervision

Self-Supervised Learning Workflow



Hypothesis: if the approximated labels behave exactly the same with the underlying labels, then we recover a fully supervised learning

Self-supervision: contrastive learning

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

Self-supervision: contrastive learning

An anchor sample x, x^+, x^- $\mathbb{E} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$

Self-supervision: contrastive learning

An anchor sample x

A positive sample x^+

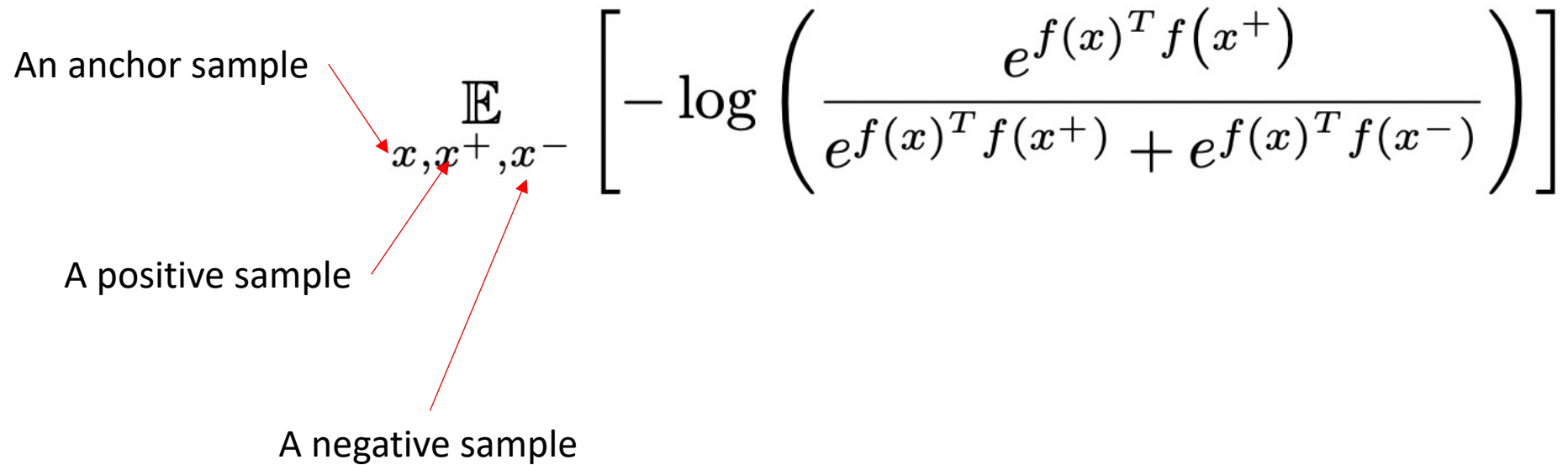
$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

Self-supervision: contrastive learning

An anchor sample x

A positive sample x^+

A negative sample x^-

$$\mathbb{E} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$
The diagram illustrates the expectation formula for contrastive learning. It features three labels on the left: 'An anchor sample' pointing to x , 'A positive sample' pointing to x^+ , and 'A negative sample' pointing to x^- . These labels are part of the expectation \mathbb{E} over the samples x, x^+, x^- . The formula inside the expectation is $-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right)$.

Self-supervision: contrastive learning

An anchor sample x , x^+ , x^-

A positive sample x^+

A negative sample x^-

$$\mathbb{E} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

→ Data augmentation: **small** variation → **not** change semantic → **not** change label

Self-supervision: contrastive learning

An anchor sample x

A positive sample x^+

A negative sample x^-

$$\mathbb{E} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

A negative sample → random sample data and directly treat them as negative samples

Data augmentation: **small** variation → **not** change semantic → **not** change label

Self-supervision: contrastive learning

An anchor sample x

A positive sample x^+

A negative sample x^-

A representation function f

$$\mathbb{E} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

Self-supervision: contrastive learning

An anchor sample x

A positive sample x^+

A negative sample x^-

A representation function f **What is it?**

$$\mathbb{E} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

Self-supervision: contrastive learning

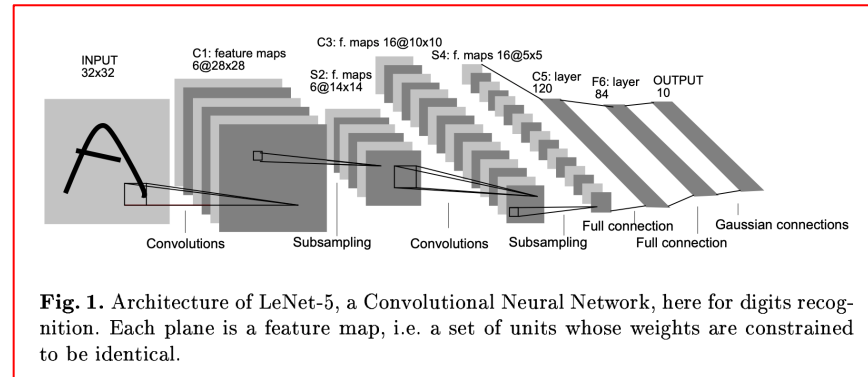
A representation function **What is it?**

An anchor sample x

A positive sample x^+

A negative sample x^-

$$\mathbb{E} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$



Self-supervision: contrastive learning

A representation function **What is it?**

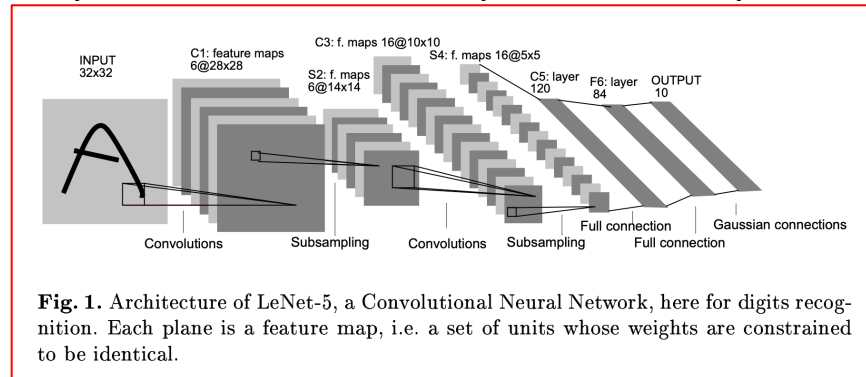
An anchor sample x

A positive sample x^+

A negative sample x^-

$$\mathbb{E} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

Each layer's feature map can be used as representation of an input data sample



Self-supervision: contrastive learning

A representation function **What is it?**

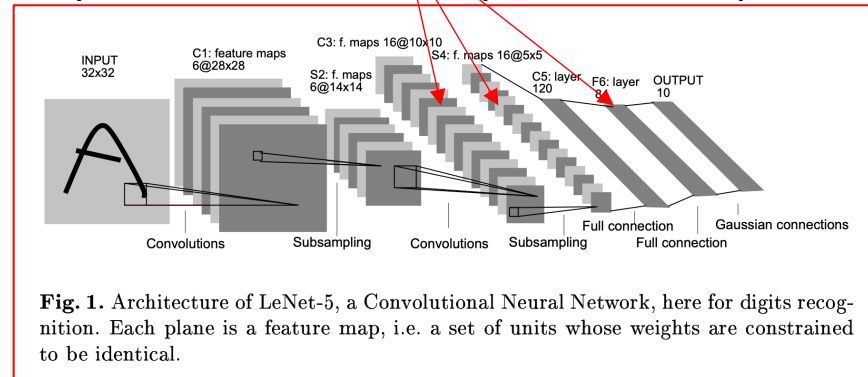
An anchor sample x, x^+, x^-

A positive sample x, x^+, x^-

A negative sample x, x^+, x^-

$$\mathbb{E} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

Each layer's **feature map** can be used as representation of an input data sample



Self-supervision: contrastive learning

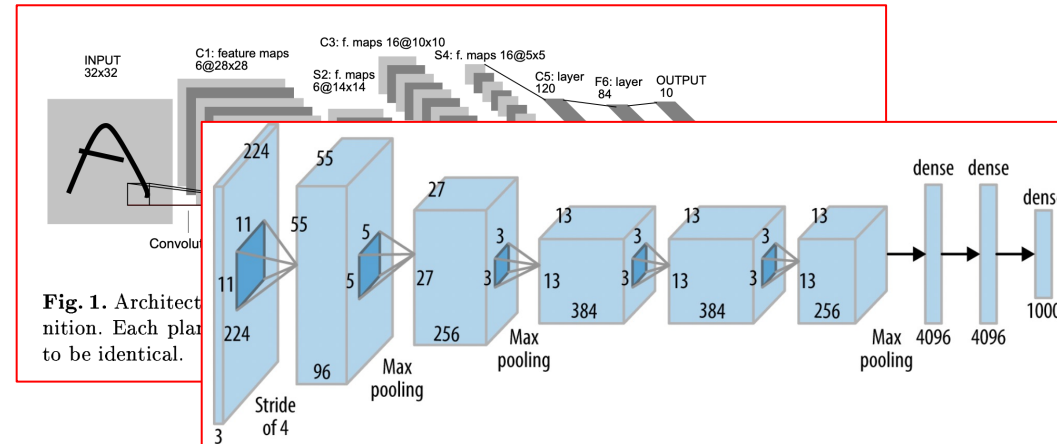
A representation function **What is it?**

An anchor sample x

A positive sample x^+

A negative sample x^-

$$\mathbb{E} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$



Self-supervision: contrastive learning

A representation function **What is it?**

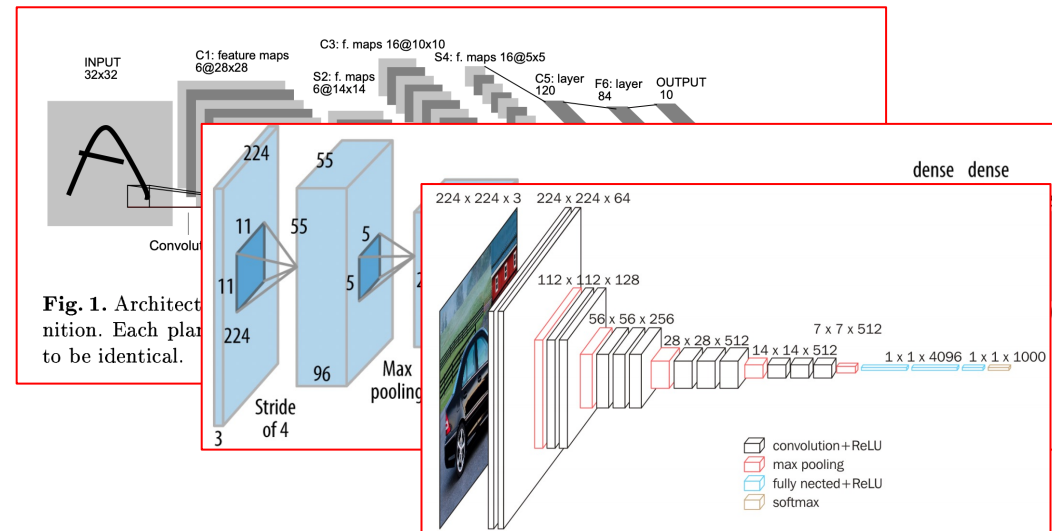
An anchor sample

\mathbb{E}_{x, x^+, x^-}

A positive sample

A negative sample

$$\left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$



Self-supervision: contrastive learning

A representation function **What is it?**

An anchor sample

\mathbb{E}

x, x^+, x^-

A positive sample

A negative sample

$$\left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

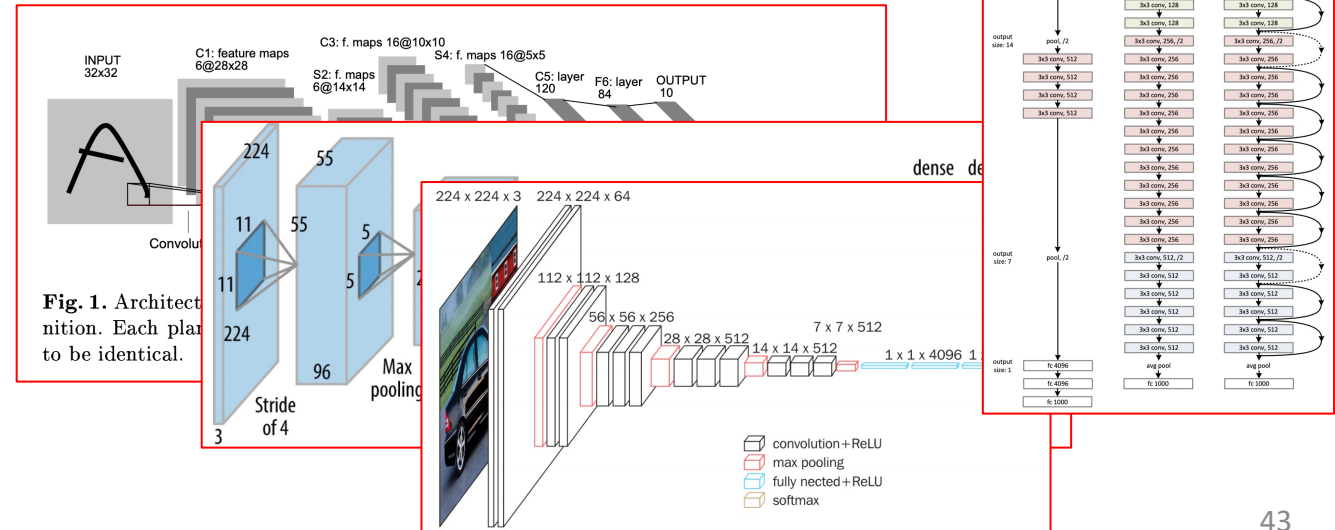


Fig. 1. Architect-
nition. Each plan
to be identical.

Self-supervision: contrastive learning

Go over each data

An anchor sample

A positive sample

A negative sample

A representation function What is it?

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

Self-supervision: contrastive learning

Go over each data

An anchor sample

A positive sample

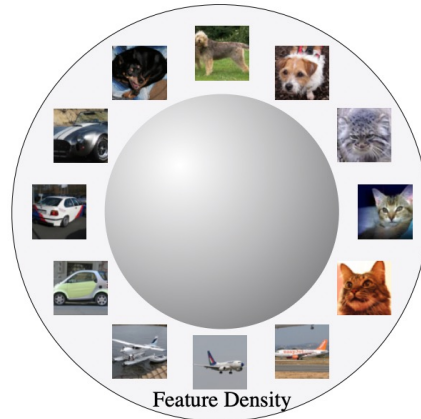
A negative sample

A representation function **What is it?**

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

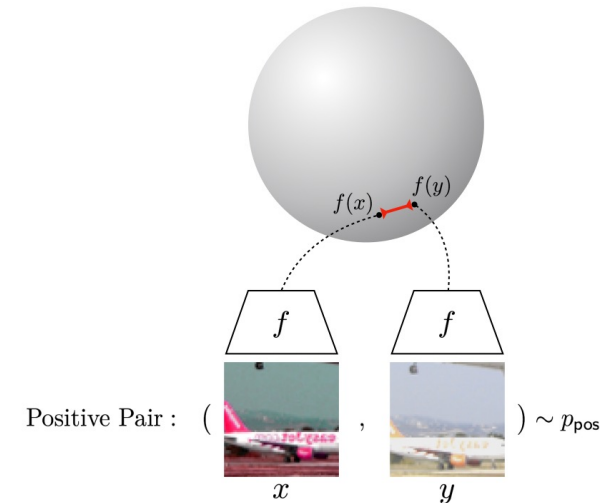
Minimize over **f**?

Self-supervision: contrastive learning



Uniformity: Preserve maximal information.

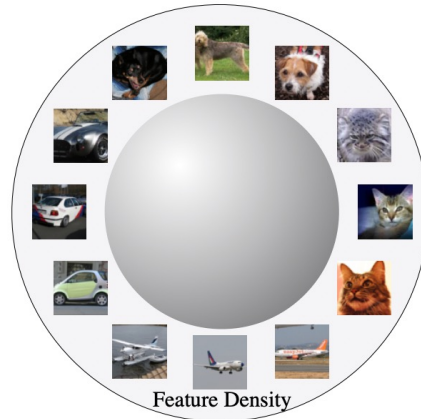
Figure 1: Illustration of alignment and uniformity of feature distributions on the output unit hypersphere. STL-10 (Coates et al., 2011) images are used for demonstration.



Alignment: Similar samples have similar features.
(Figure inspired by Tian et al. (2019).)

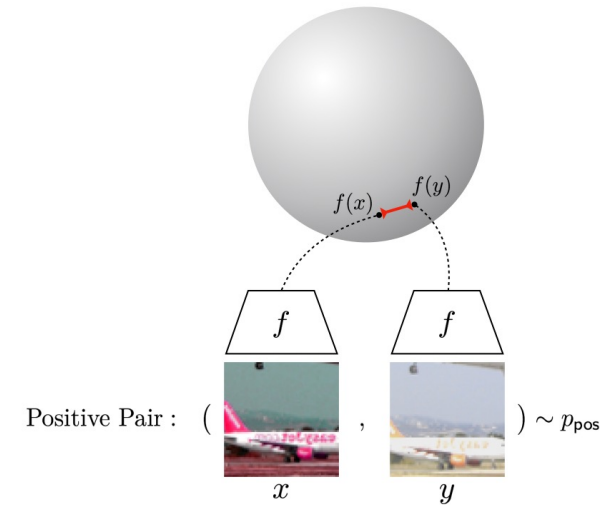
Self-supervision: contrastive learning

Data samples from different classes



Uniformity: Preserve maximal information.

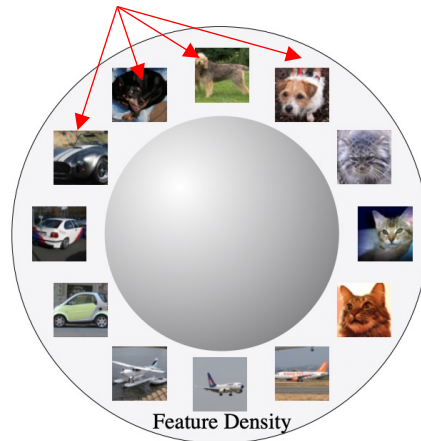
Figure 1: Illustration of alignment and uniformity of feature distributions on the output unit hypersphere. STL-10 (Coates et al., 2011) images are used for demonstration.



Alignment: Similar samples have similar features.
(Figure inspired by Tian et al. (2019).)

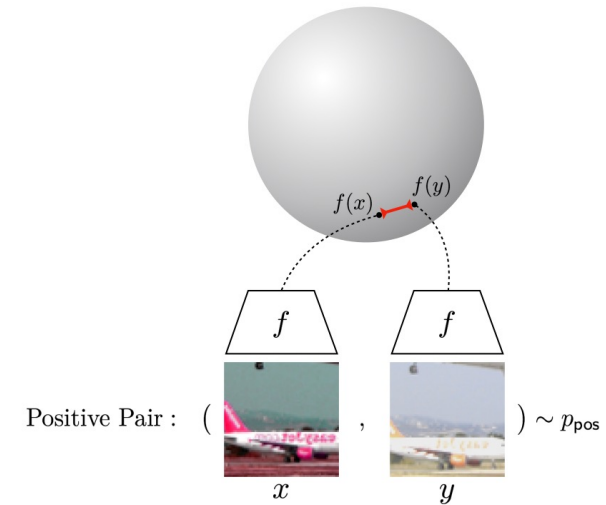
Self-supervision: contrastive learning

Data samples from different classes



Uniformity: Preserve maximal information.

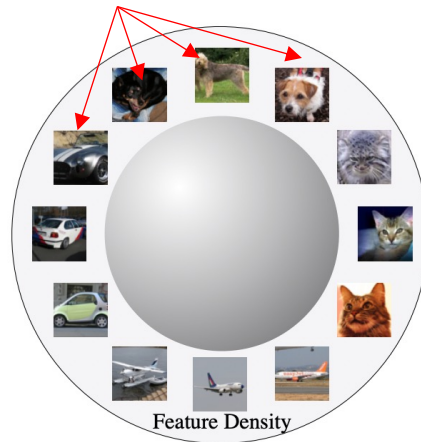
Figure 1: Illustration of alignment and uniformity of feature distributions on the output unit hypersphere. STL-10 (Coates et al., 2011) images are used for demonstration.



Alignment: Similar samples have similar features.
(Figure inspired by Tian et al. (2019).)

Self-supervision: contrastive learning

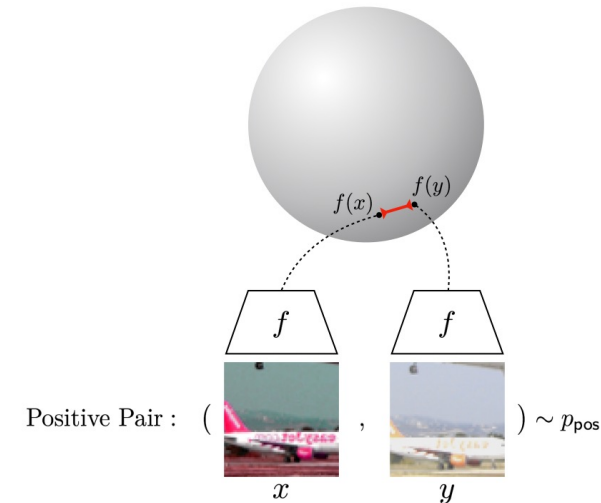
Data samples from different classes



Uniformity: Preserve maximal information.

Figure 1: Illustration of alignment and uniformity of feature distributions on the output unit hypersphere. STL-10 (Coates et al., 2011) images are used for demonstration.

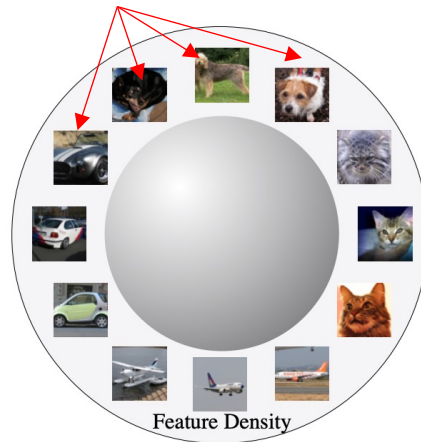
Data samples from the same class:



Alignment: Similar samples have similar features.
(Figure inspired by Tian et al. (2019).)

Self-supervision: contrastive learning

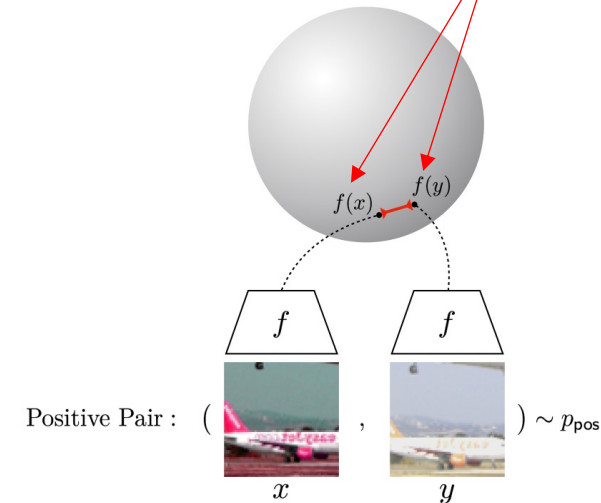
Data samples from different classes



Uniformity: Preserve maximal information.

Figure 1: Illustration of alignment and uniformity of feature distributions on the output unit hypersphere. STL-10 (Coates et al., 2011) images are used for demonstration.

Data samples from the same class:
stay close to each other



Alignment: Similar samples have similar features.
(Figure inspired by Tian et al. (2019).)

Self-supervision: contrastive learning

Random init

Supervised learning

Contrastive learning
(unsupervised, self-supervised)

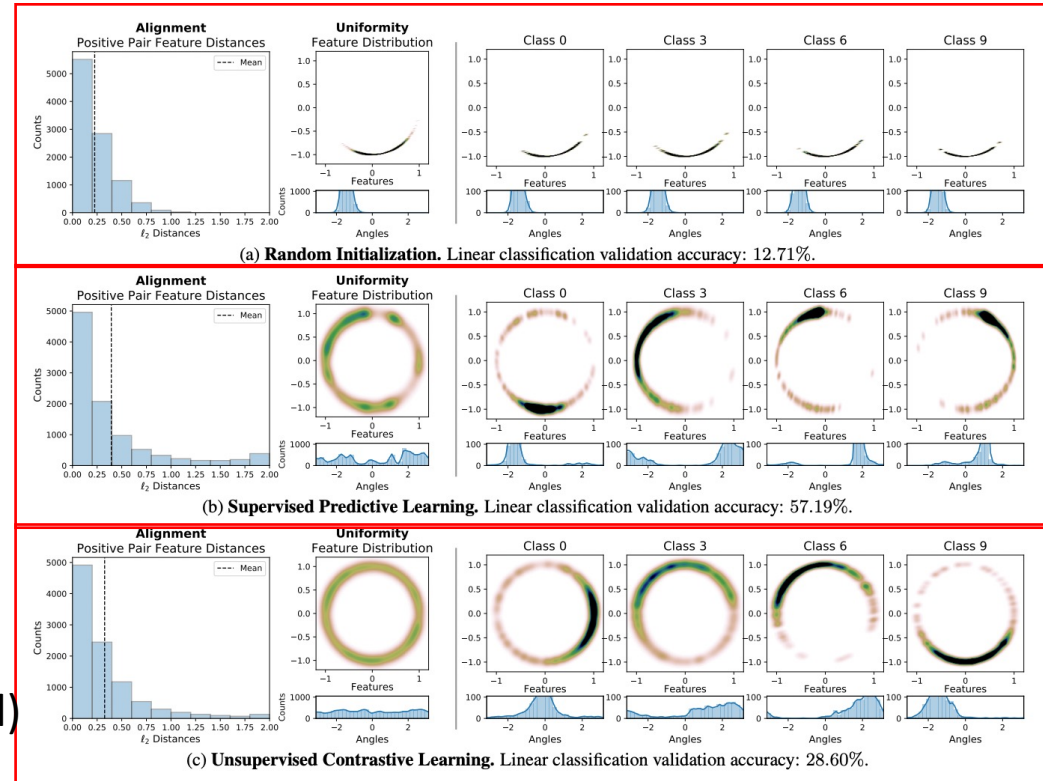


Figure 3: Representations of CIFAR-10 validation set on S^1 . **Alignment analysis:** We show distribution of distance between features of positive pairs (two random augmentations). **Uniformity analysis:** We plot feature distributions with Gaussian kernel density estimation (KDE) in \mathbb{R}^2 and von Mises-Fisher (vMF) KDE on angles (i.e., $\arctan 2(y, x)$ for each point $(x, y) \in S^1$). **Four rightmost plots** visualize feature distributions of selected specific classes. Representation from contrastive learning is both *aligned* (having low positive pair feature distances) and *uniform* (evenly distributed on S^1).

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

Review: A negative sample \rightarrow random sample data and directly treat them as negative samples

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

What if “negative” samples are actually from the same class?

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

What if “negative” samples are actually from the same class?

False negative

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

What if “negative” samples are actually from the same class?

False negative

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

What if “negative” samples are actually from the same class?

False negative

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

What if **“negative” samples** are actually from the same class?

False negative

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

$$\begin{aligned} L_{un}^=(f) &= \mathbb{E}_{\substack{c \sim \nu \\ x, x^+, x^- \sim \mathcal{D}_c^3}} [\ell(f(x)^T (f(x^+) - f(x^-)))] \\ &\geq \mathbb{E}_{c \sim \nu, x \sim \mathcal{D}_c} [\ell(f(x)^T (\mu_c - \mu_c))] = 1 \end{aligned}$$

What if “**negative**” samples are actually from the same class?

False negative

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

$$\begin{aligned} L_{un}^=(f) &= \mathbb{E}_{\substack{c \sim \nu \\ x, x^+, x^- \sim \mathcal{D}_c^3}} [\ell(f(x)^T (f(x^+) - f(x^-)))] \\ &\geq \mathbb{E}_{c \sim \nu, x \sim \mathcal{D}_c} [\ell(f(x)^T (\mu_c - \mu_c))] = 1 \end{aligned}$$

$$\begin{aligned} L_{un}^{\neq}(f) &= \mathbb{E}_{\substack{c^+, c^- \sim \rho^2 \\ x, x^+ \sim \mathcal{D}_{c^+}^2 \\ x^- \sim \mathcal{D}_{c^-}}} [\ell(f(x)^T (f(x^+) - f(x^-))) | c^+ \neq c^-] \end{aligned}$$

What if “**negative**” samples are actually from the same class?

False negative

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

$$\begin{aligned} L_{un}^=(f) &= \mathbb{E}_{\substack{c \sim \nu \\ x, x^+, x^- \sim \mathcal{D}_c^3}} [\ell(f(x)^T (f(x^+) - f(x^-)))] \\ &\geq \mathbb{E}_{c \sim \nu, x \sim \mathcal{D}_c} [\ell(f(x)^T (\mu_c - \mu_c))] = 1 \end{aligned}$$

$$\begin{aligned} L_{un}^{\neq}(f) &= \mathbb{E}_{\substack{c^+, c^- \sim \rho^2 \\ x, x^+ \sim \mathcal{D}_{c^+}^2 \\ x^- \sim \mathcal{D}_{c^-}}} [\ell(f(x)^T (f(x^+) - f(x^-))) | c^+ \neq c^-] \end{aligned}$$

What if “**negative**” samples are actually from the same class?

False negative

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

$$\begin{aligned} L_{un}^=(f) &= \mathbb{E}_{x, x^+, x^- \sim \mathcal{D}_c^3} [\ell(f(x)^T (f(x^+) - f(x^-)))] \\ &\geq \mathbb{E}_{c \sim \nu, x \sim \mathcal{D}_c} [\ell(f(x)^T (\mu_c - \mu_c))] = 1 \end{aligned}$$

$$\begin{aligned} L_{un}^{\neq}(f) &= \mathbb{E}_{\substack{c^+, c^- \sim \rho^2 \\ x, x^+ \sim \mathcal{D}_{c^+}^2 \\ x^- \sim \mathcal{D}_{c^-}}} [\ell(f(x)^T (f(x^+) - f(x^-)))] \end{aligned}$$

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \hat{L}_{un}(f)$$

Theorem 4.5. With probability at least $1 - \delta$, $\forall f \in \mathcal{F}$

$$L_{sup}(\hat{f}) \leq L_{sup}^{\mu}(\hat{f}) \leq L_{un}^{\neq}(f) + \beta s(f) + \eta Gen_M$$

where $\beta = c' \frac{\tau}{1-\tau}$, $\eta = \frac{1}{1-\tau}$ and c' is a constant.

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

$$\begin{aligned} L_{un}^=(f) &= \mathbb{E}_{x, x^+, x^- \sim \mathcal{D}_c} [\ell(f(x)^T (f(x^+) - f(x^-)))] \\ &\geq \mathbb{E}_{c \sim \nu, x \sim \mathcal{D}_c} [\ell(f(x)^T (\mu_c - \mu_c))] = 1 \end{aligned}$$

$$\begin{aligned} L_{un}^{\neq}(f) &= \mathbb{E}_{\substack{c^+, c^- \sim \rho^2 \\ x, x^+ \sim \mathcal{D}_{c^+}^2 \\ x^- \sim \mathcal{D}_{c^-}}} [\ell(f(x)^T (f(x^+) - f(x^-)))] \end{aligned}$$

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \hat{L}_{un}(f)$$

Theorem 4.5. With probability at least $1 - \delta$, $\forall f \in \mathcal{F}$

$$L_{sup}(\hat{f}) \leq L_{sup}^{\mu}(\hat{f}) \leq L_{un}^{\neq}(f) + \beta s(f) + \eta Gen_M$$

where $\beta = c' \frac{\tau}{1-\tau}$, $\eta = \frac{1}{1-\tau}$ and c' is a constant.

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

$$\begin{aligned} L_{un}^=(f) &= \mathbb{E}_{x, x^+, x^- \sim \mathcal{D}_c} [\ell(f(x)^T (f(x^+) - f(x^-)))] \\ &\geq \mathbb{E}_{c \sim \nu, x \sim \mathcal{D}_c} [\ell(f(x)^T (\mu_c - \mu_c))] = 1 \end{aligned}$$

$$\begin{aligned} L_{un}^{\neq}(f) &= \mathbb{E}_{\substack{c^+, c^- \sim \rho^2 \\ x, x^+ \sim \mathcal{D}_{c^+}^2 \\ x^- \sim \mathcal{D}_{c^-}}} [\ell(f(x)^T (f(x^+) - f(x^-)))] \end{aligned}$$

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \hat{L}_{un}(f)$$

Theorem 4.5. With probability at least $1 - \delta$, $\forall f \in \mathcal{F}$

$$L_{sup}(\hat{f}) \leq L_{sup}^{\mu}(\hat{f}) \leq L_{un}^{\neq}(f) + \beta s(f) + \eta Gen_M$$

where $\beta = c' \frac{\tau}{1-\tau}$, $\eta = \frac{1}{1-\tau}$ and c' is a constant.
constant

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

$$\begin{aligned} L_{un}^=(f) &= \mathbb{E}_{x, x^+, x^- \sim \mathcal{D}_c} [\ell(f(x)^T (f(x^+) - f(x^-)))] \\ &\geq \mathbb{E}_{c \sim \nu, x \sim \mathcal{D}_c} [\ell(f(x)^T (\mu_c - \mu_c))] = 1 \end{aligned}$$

$$\begin{aligned} L_{un}^{\neq}(f) &= \mathbb{E}_{\substack{c^+, c^- \sim \rho^2 \\ x, x^+ \sim \mathcal{D}_{c^+} \\ x^- \sim \mathcal{D}_{c^-}}} [\ell(f(x)^T (f(x^+) - f(x^-)))] \end{aligned}$$

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \hat{L}_{un}(f)$$

Theorem 4.5. With probability at least $1 - \delta$, $\forall f \in \mathcal{F}$

$$L_{sup}(\hat{f}) \leq L_{sup}^{\mu}(\hat{f}) \leq L_{un}^{\neq}(f) + \beta s(f) + \eta Gen_M$$

where $\beta = c' \frac{\tau}{1-\tau}$, $\eta = \frac{1}{1-\tau}$ and c' is a constant.

constant

$$o(\sqrt{1/n})$$

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

$$\begin{aligned} L_{un}^=(f) &= \mathbb{E}_{x, x^+, x^- \sim \mathcal{D}_c^3} [\ell(f(x)^T (f(x^+) - f(x^-)))] \\ &\geq \mathbb{E}_{c \sim \nu, x \sim \mathcal{D}_c} [\ell(f(x)^T (\mu_c - \mu_c))] = 1 \end{aligned}$$

$$\begin{aligned} L_{un}^{\neq}(f) &= \mathbb{E}_{\substack{c^+, c^- \sim \rho^2 \\ x, x^+ \sim \mathcal{D}_{c^+}^2 \\ x^- \sim \mathcal{D}_{c^-}}} [\ell(f(x)^T (f(x^+) - f(x^-)))] \end{aligned}$$

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \hat{L}_{un}(f)$$

Theorem 4.5. With probability at least $1 - \delta$, $\forall f \in \mathcal{F}$

$$L_{sup}(\hat{f}) \leq L_{sup}^{\mu}(\hat{f}) \leq L_{un}^{\neq}(f) + \beta s(f) + \eta Gen_M$$

where $\beta = c' \frac{\tau}{1-\tau}$, $\eta = \frac{1}{1-\tau}$ and c' is a constant.

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

$$\begin{aligned} L_{un}^=(f) &= \mathbb{E}_{x, x^+, x^- \sim \mathcal{D}_c} [\ell(f(x)^T (f(x^+) - f(x^-)))] \\ &\geq \mathbb{E}_{c \sim \nu, x \sim \mathcal{D}_c} [\ell(f(x)^T (\mu_c - \mu_c))] = 1 \end{aligned}$$

$$\begin{aligned} L_{un}^{\neq}(f) &= \mathbb{E}_{\substack{c^+, c^- \sim \rho^2 \\ x, x^+ \sim \mathcal{D}_{c^+} \\ x^- \sim \mathcal{D}_{c^-}}} [\ell(f(x)^T (f(x^+) - f(x^-)))] \end{aligned}$$

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \hat{L}_{un}(f)$$

Theorem 4.5. With probability at least $1 - \delta$, $\forall f \in \mathcal{F}$

$$L_{sup}(\hat{f}) \leq L_{sup}^{\mu}(\hat{f}) \leq L_{un}^{\neq}(f) + \beta s(f) + \eta Gen_M$$

where $\beta = c' \frac{\tau}{1-\tau}$, $\eta = \frac{1}{1-\tau}$ and c' is a constant.

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

$$\begin{aligned} L_{un}^=(f) &= \mathbb{E}_{x, x^+, x^- \sim \mathcal{D}_c} [\ell(f(x)^T (f(x^+) - f(x^-)))] \\ &\geq \mathbb{E}_{c \sim \nu, x \sim \mathcal{D}_c} [\ell(f(x)^T (\mu_c - \mu_c))] = 1 \end{aligned}$$

$$\begin{aligned} L_{un}^{\neq}(f) &= \mathbb{E}_{\substack{c^+, c^- \sim \rho^2 \\ x, x^+ \sim \mathcal{D}_{c^+}^2 \\ x^- \sim \mathcal{D}_{c^-}}} [\ell(f(x)^T (f(x^+) - f(x^-)) | c^+ \neq c^-)] \end{aligned}$$

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \hat{L}_{un}(f)$$

Theorem 4.5. With probability at least $1 - \delta$, $\forall f \in \mathcal{F}$

$$L_{sup}(\hat{f}) \leq L_{sup}^{\mu}(\hat{f}) \leq L_{un}^{\neq}(f) + \beta s(f) + \eta Gen_M$$

where $\beta = c' \frac{\tau}{1-\tau}$, $\eta = \frac{1}{1-\tau}$ and c' is a constant.

If **tau** is very small?

Self-supervision: contrastive learning

- The **gap** between supervised learning and contrastive learning?

The prob: we have
true negative

$$\mathbb{E}_{x, x^+, x^-} \left[-\log \left(\frac{e^{f(x)^T f(x^+)}}{e^{f(x)^T f(x^+)} + e^{f(x)^T f(x^-)}} \right) \right]$$

$$L_{un}(f) = \tau L_{un}^=(f) + (1 - \tau) L_{un}^{\neq}(f) \quad (8)$$

$$\begin{aligned} L_{un}^=(f) &= \mathbb{E}_{x, x^+, x^- \sim \mathcal{D}_c} [\ell(f(x)^T (f(x^+) - f(x^-)))] \\ &\geq \mathbb{E}_{c \sim \nu, x \sim \mathcal{D}_c} [\ell(f(x)^T (\mu_c - \mu_c))] = 1 \end{aligned}$$

$$\begin{aligned} L_{un}^{\neq}(f) &= \mathbb{E}_{\substack{c^+, c^- \sim \rho^2 \\ x, x^+ \sim \mathcal{D}_{c^+} \\ x^- \sim \mathcal{D}_{c^-}}} [\ell(f(x)^T (f(x^+) - f(x^-)))] \end{aligned}$$

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \hat{L}_{un}(f)$$

Theorem 4.5. With probability at least $1 - \delta$, $\forall f \in \mathcal{F}$

$$L_{sup}(\hat{f}) \leq L_{sup}^{\mu}(\hat{f}) \leq L_{un}^{\neq}(f) + \beta s(f) + \eta Gen_M$$

where $\beta = c' \frac{\tau}{1-\tau}$, $\eta = \frac{1}{1-\tau}$ and c' is a constant.

If **tau** is very small?

If we have **a large number of classes** and we perform **uniform random sampling**?

Meta-learning

Multi-Task Learning

Solve multiple tasks $\mathcal{T}_1, \dots, \mathcal{T}_T$ at once.

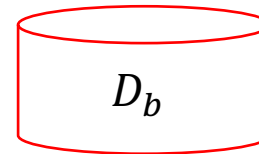
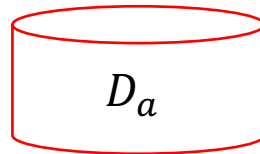
$$\min_{\theta} \sum_{i=1}^T \mathcal{L}_i(\theta, \mathcal{D}_i)$$

Meta-learning

Multi-Task Learning

Solve multiple tasks $\mathcal{T}_1, \dots, \mathcal{T}_T$ at once.

$$\min_{\theta} \sum_{i=1}^T \mathcal{L}_i(\theta, \mathcal{D}_i)$$

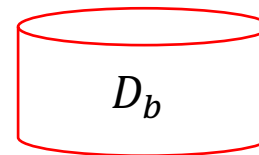
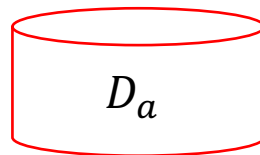


Meta-learning

Multi-Task Learning

Solve multiple tasks $\mathcal{T}_1, \dots, \mathcal{T}_T$ at once.

$$\min_{\theta} \sum_{i=1}^T \mathcal{L}_i(\theta, \mathcal{D}_i)$$



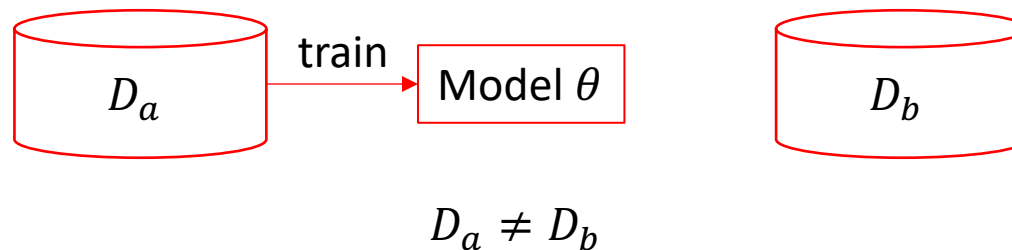
$$D_a \neq D_b$$

Meta-learning

Multi-Task Learning

Solve multiple tasks $\mathcal{T}_1, \dots, \mathcal{T}_T$ at once.

$$\min_{\theta} \sum_{i=1}^T \mathcal{L}_i(\theta, \mathcal{D}_i)$$

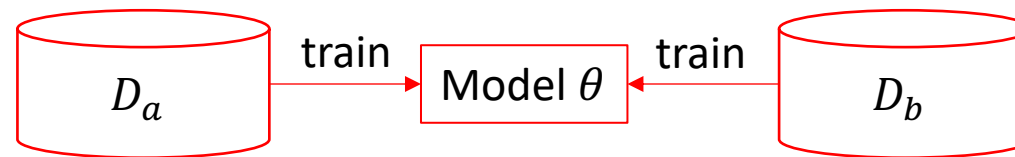


Meta-learning

Multi-Task Learning

Solve multiple tasks $\mathcal{T}_1, \dots, \mathcal{T}_T$ at once.

$$\min_{\theta} \sum_{i=1}^T \mathcal{L}_i(\theta, \mathcal{D}_i)$$



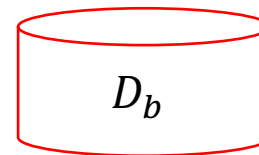
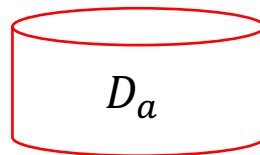
$$D_a \neq D_b$$

Meta-learning

Transfer Learning

Solve target task \mathcal{T}_b after solving source task \mathcal{T}_a
by *transferring* knowledge learned from \mathcal{T}_a

Key assumption: Cannot access data \mathcal{D}_a during transfer.

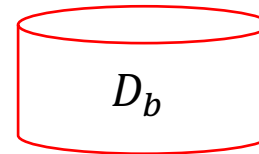
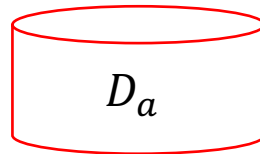


Meta-learning

Transfer Learning

Solve target task \mathcal{T}_b after solving source task \mathcal{T}_a
by *transferring* knowledge learned from \mathcal{T}_a

Key assumption: Cannot access data \mathcal{D}_a during transfer.



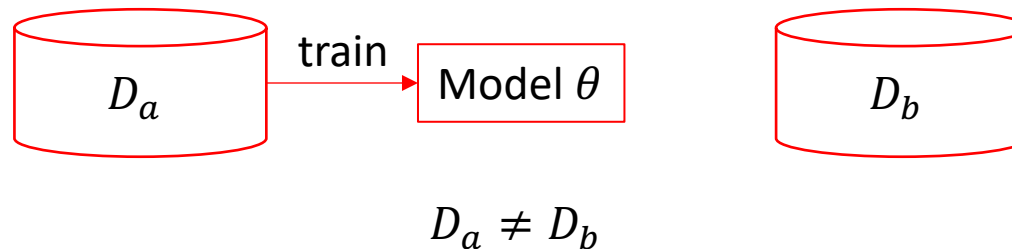
$$D_a \neq D_b$$

Meta-learning

Transfer Learning

Solve target task \mathcal{T}_b after solving source task \mathcal{T}_a
by *transferring* knowledge learned from \mathcal{T}_a

Key assumption: Cannot access data \mathcal{D}_a during transfer.

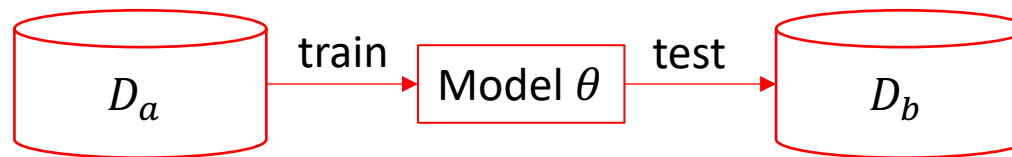


Meta-learning

Transfer Learning

Solve target task \mathcal{T}_b after solving source task \mathcal{T}_a
by *transferring* knowledge learned from \mathcal{T}_a

Key assumption: Cannot access data \mathcal{D}_a during transfer.



$$D_a \neq D_b$$

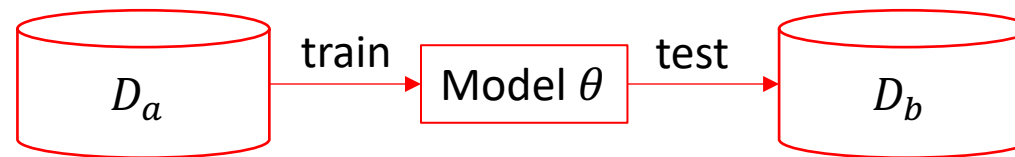
Meta-learning

Q: anything we learned can be regarded as “transfer learning”?

Transfer Learning

Solve target task \mathcal{T}_b after solving source task \mathcal{T}_a
by *transferring* knowledge learned from \mathcal{T}_a

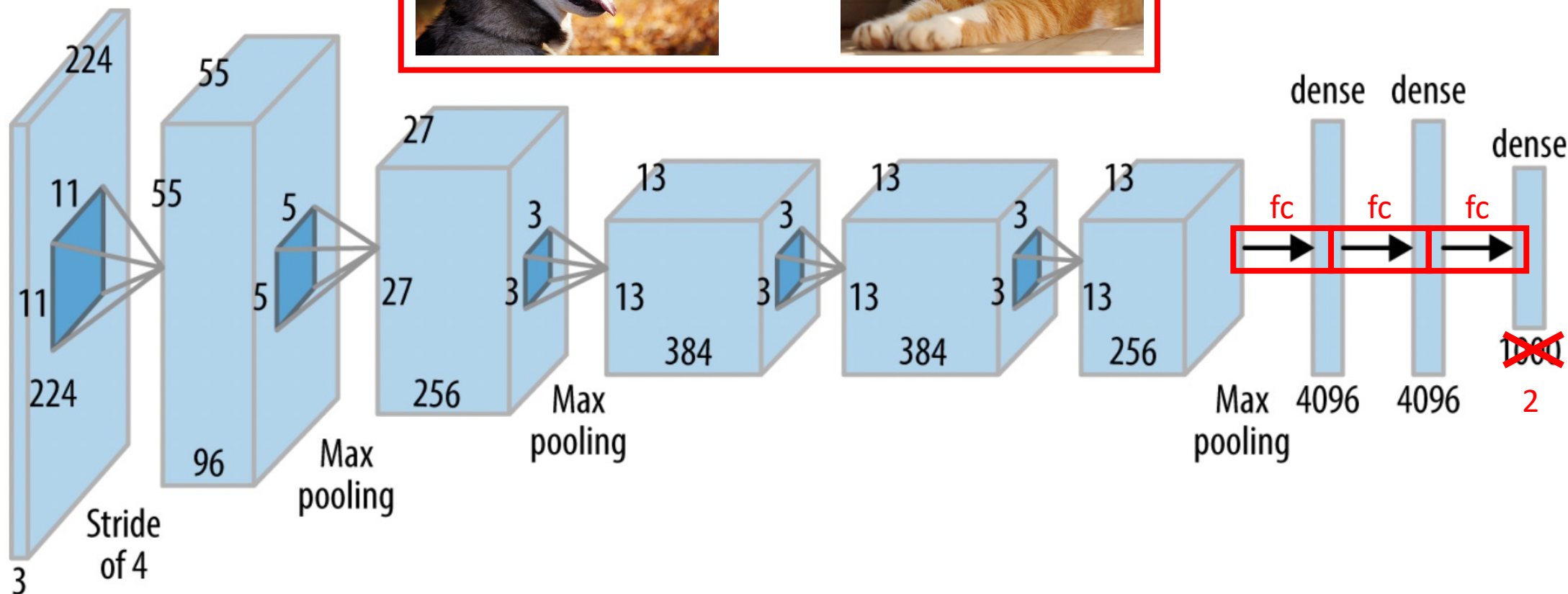
Key assumption: Cannot access data \mathcal{D}_a during transfer.



$$D_a \neq D_b$$

Pre-train?

Two classes



Suppose we have a learned model → weight parameters are determined and fixed

Meta-learning

The Meta-Learning Problem

Given data from $\mathcal{T}_1, \dots, \mathcal{T}_n$, quickly solve new task $\mathcal{T}_{\text{test}}$

Meta-learning

The Meta-Learning Problem

Given data from $\mathcal{T}_1, \dots, \mathcal{T}_n$, quickly solve new task $\mathcal{T}_{\text{test}}$

Meta-learning

The Meta-Learning Problem

Given data from $\mathcal{T}_1, \dots, \mathcal{T}_n$, quickly solve new task $\mathcal{T}_{\text{test}}$

Not directly apply the trained model for prediction

Meta-learning

The Meta-Learning Problem

Given data from $\mathcal{T}_1, \dots, \mathcal{T}_n$, quickly solve new task $\mathcal{T}_{\text{test}}$



Not directly apply the trained model for prediction

Allow some quick update to fit the testing data

Meta-learning

The Meta-Learning Problem

Given data from $\mathcal{T}_1, \dots, \mathcal{T}_n$, quickly solve new task $\mathcal{T}_{\text{test}}$

Not directly apply the trained model for prediction

Allow some quick update to fit the testing data

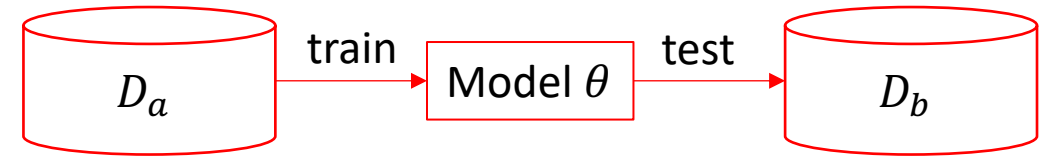
$\mathcal{T}_{\text{test}}$ can be one of many target tasks

Meta-learning

Multi-Task Learning

Solve multiple tasks $\mathcal{T}_1, \dots, \mathcal{T}_T$ at once.

$$\min_{\theta} \sum_{i=1}^T \mathcal{L}_i(\theta, \mathcal{D}_i)$$



Transfer Learning

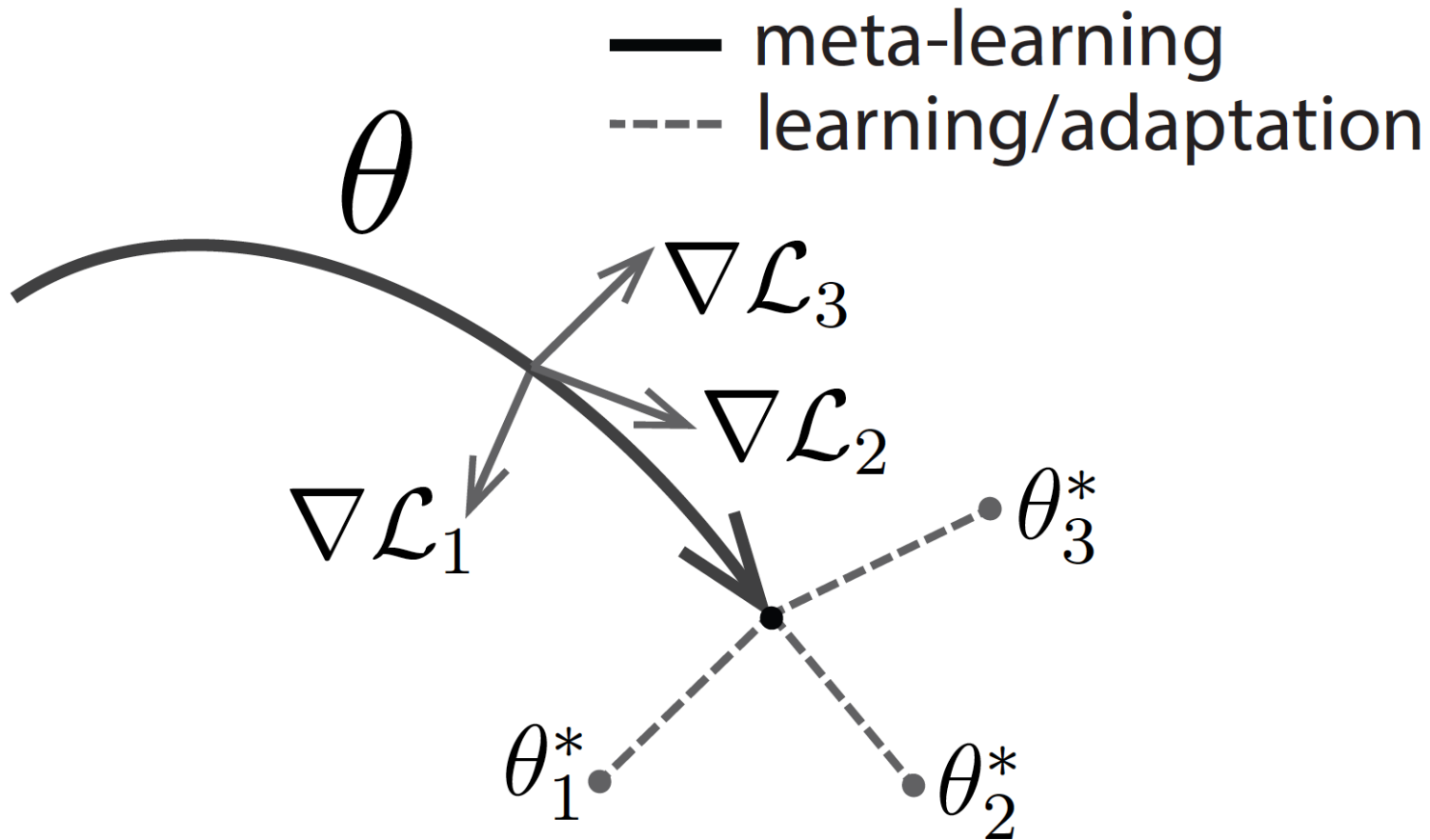
Solve target task \mathcal{T}_b after solving source task \mathcal{T}_a by *transferring* knowledge learned from \mathcal{T}_a

Key assumption: Cannot access data \mathcal{D}_a during transfer.

The Meta-Learning Problem

Given data from $\mathcal{T}_1, \dots, \mathcal{T}_n$, quickly solve new task $\mathcal{T}_{\text{test}}$

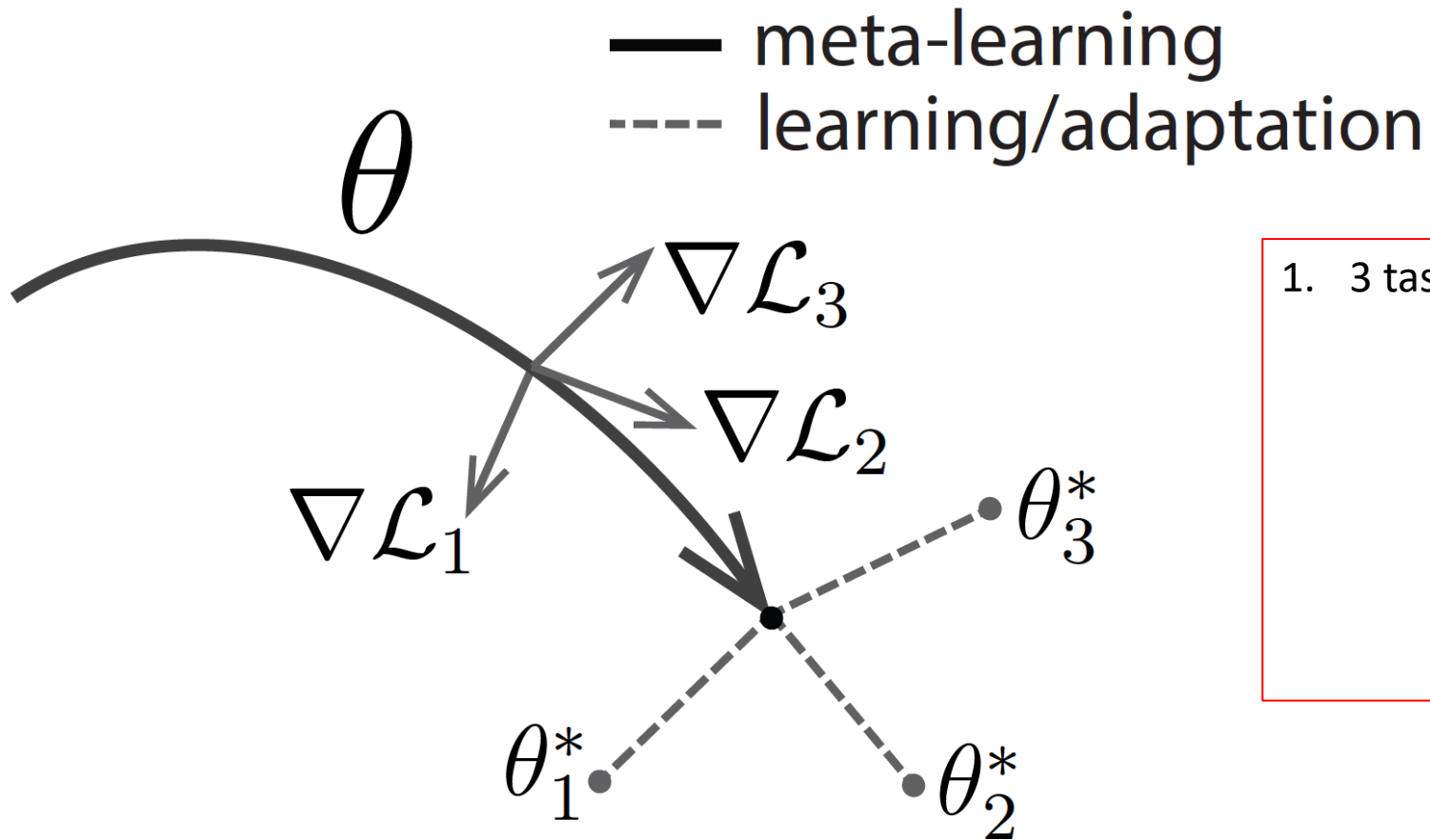
Meta-learning



Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." In *International Conference on Machine Learning*, pp. 1126-1135. PMLR, 2017.

<http://proceedings.mlr.press/v70/finn17a/finn17a.pdf>

Meta-learning

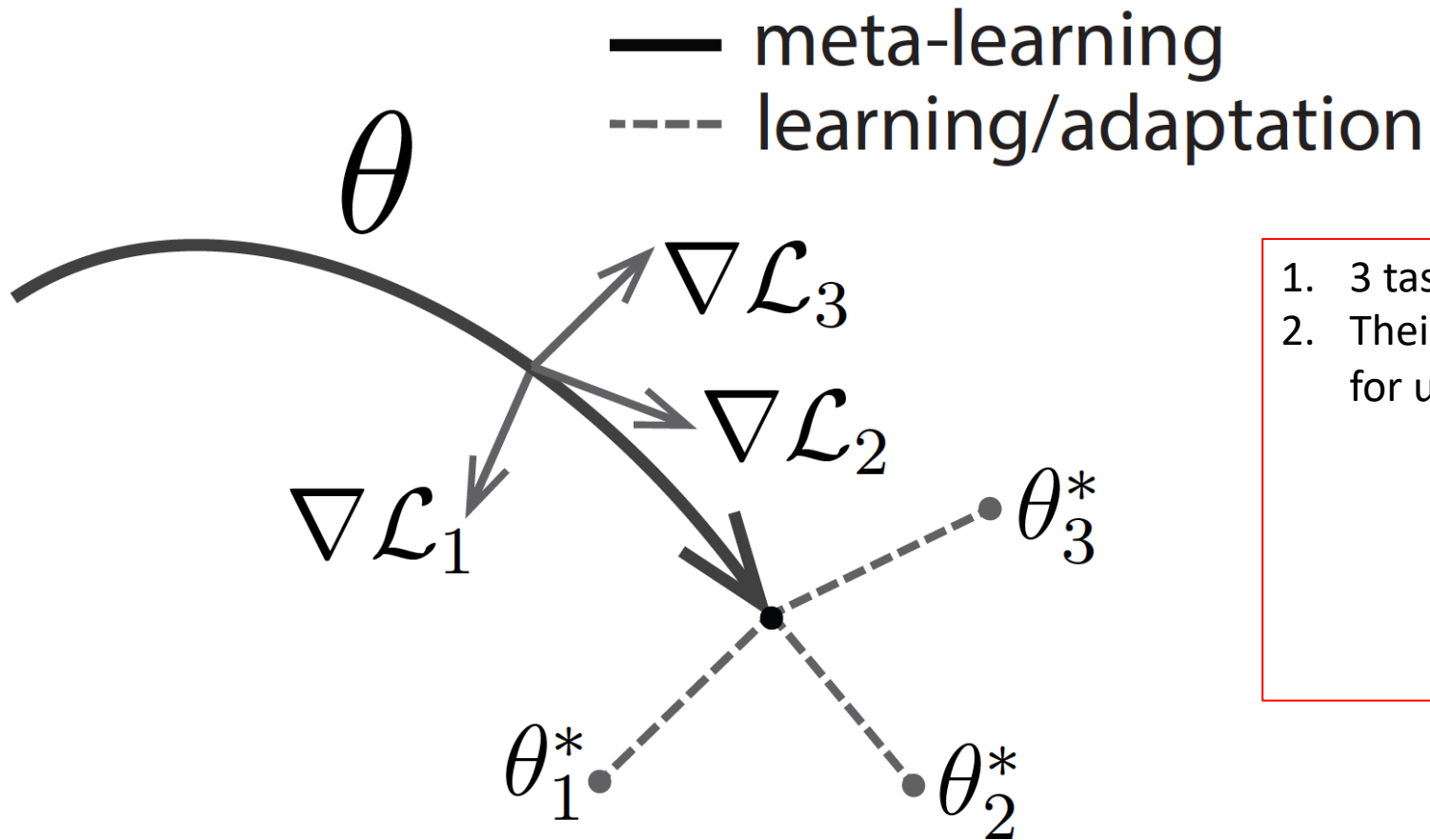


1. 3 tasks

Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." In *International Conference on Machine Learning*, pp. 1126-1135. PMLR, 2017.

<http://proceedings.mlr.press/v70/finn17a/finn17a.pdf>

Meta-learning

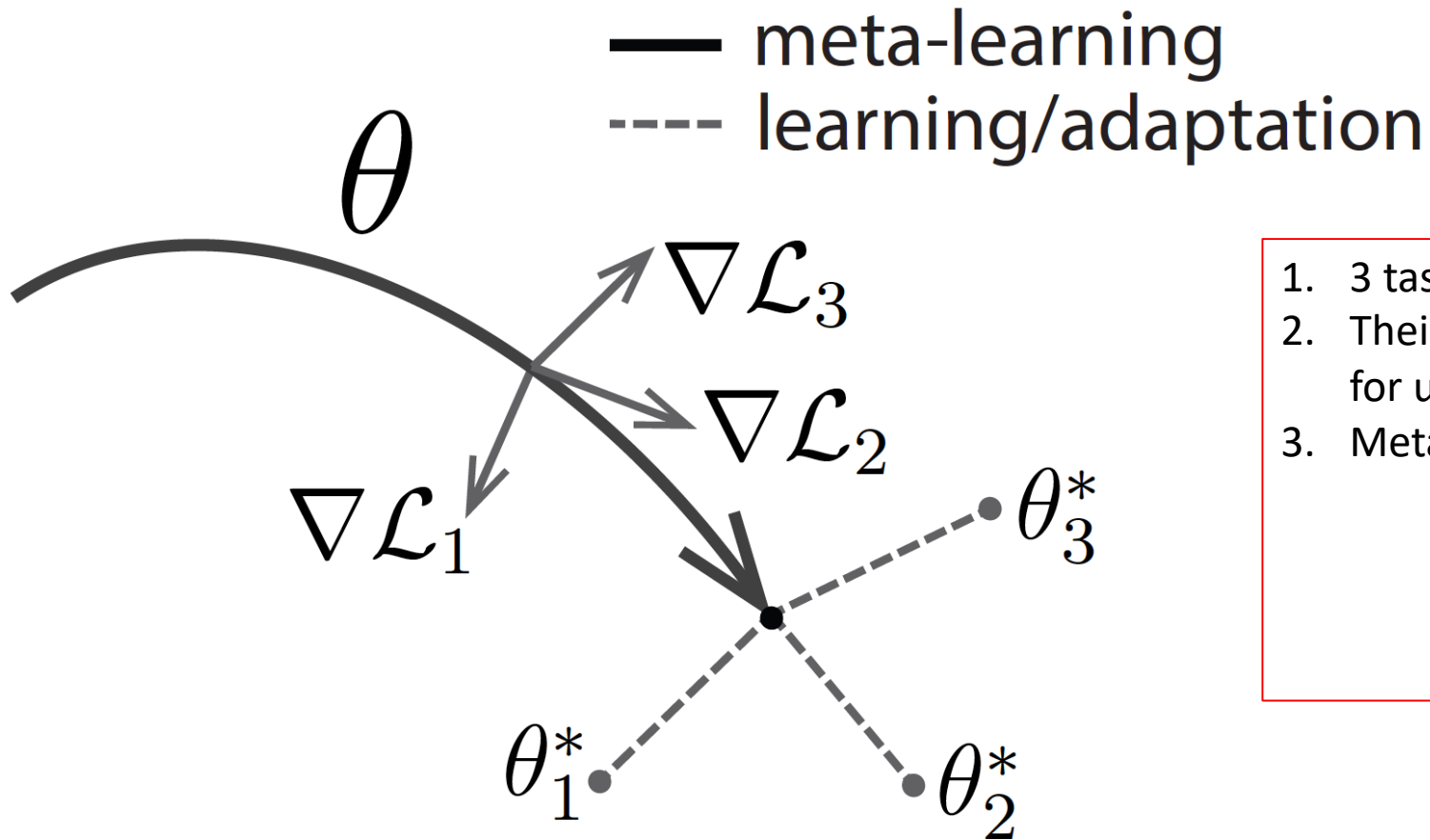


1. 3 tasks
2. Their datasets may have different directions/gradients for updating current model

Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." In *International Conference on Machine Learning*, pp. 1126-1135. PMLR, 2017.

<http://proceedings.mlr.press/v70/finn17a/finn17a.pdf>

Meta-learning

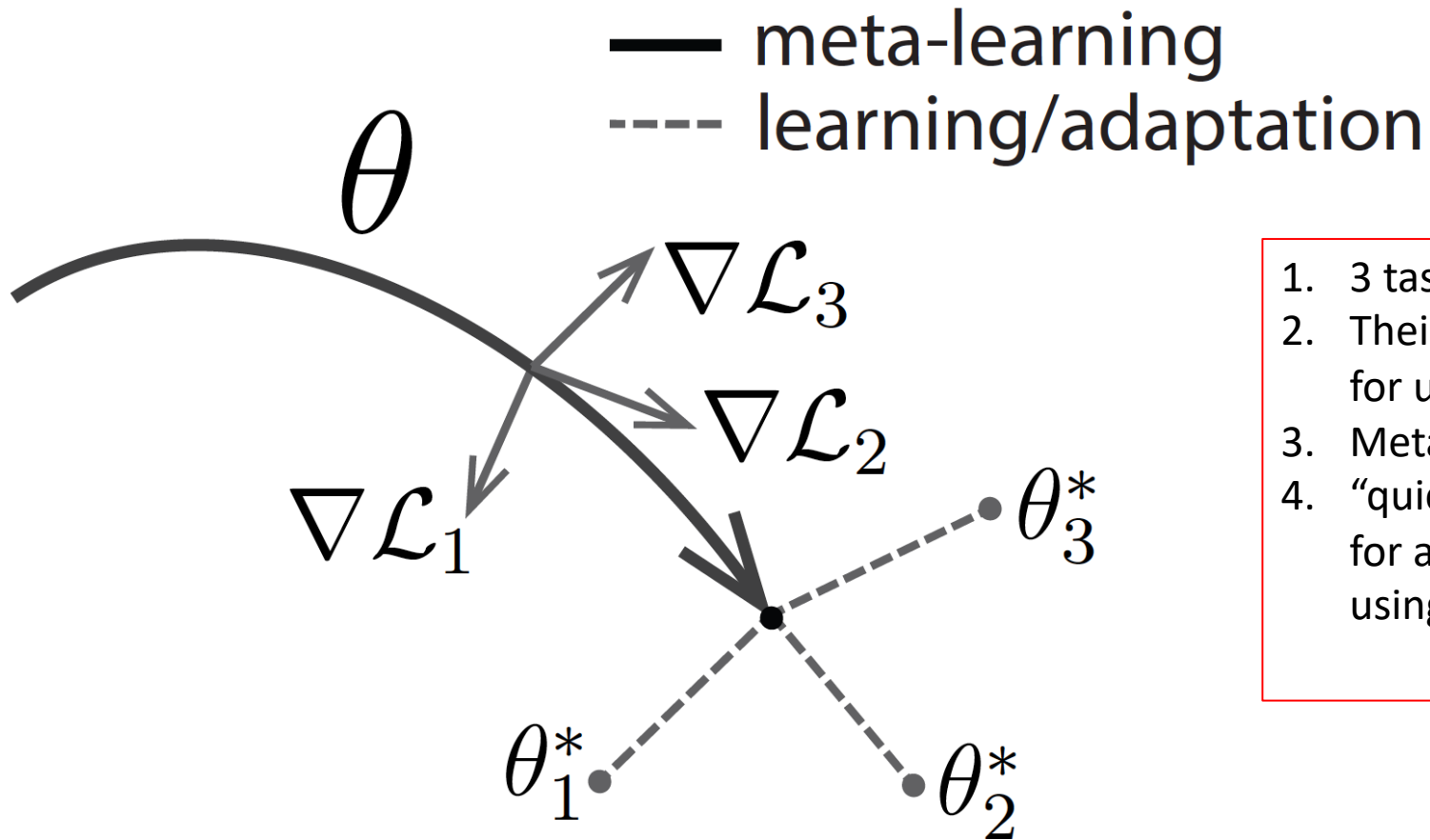


1. 3 tasks
2. Their datasets may have different directions/gradients for updating current model
3. Meta learning determines one update direction

Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." In *International Conference on Machine Learning*, pp. 1126-1135. PMLR, 2017.

<http://proceedings.mlr.press/v70/finn17a/finn17a.pdf>

Meta-learning

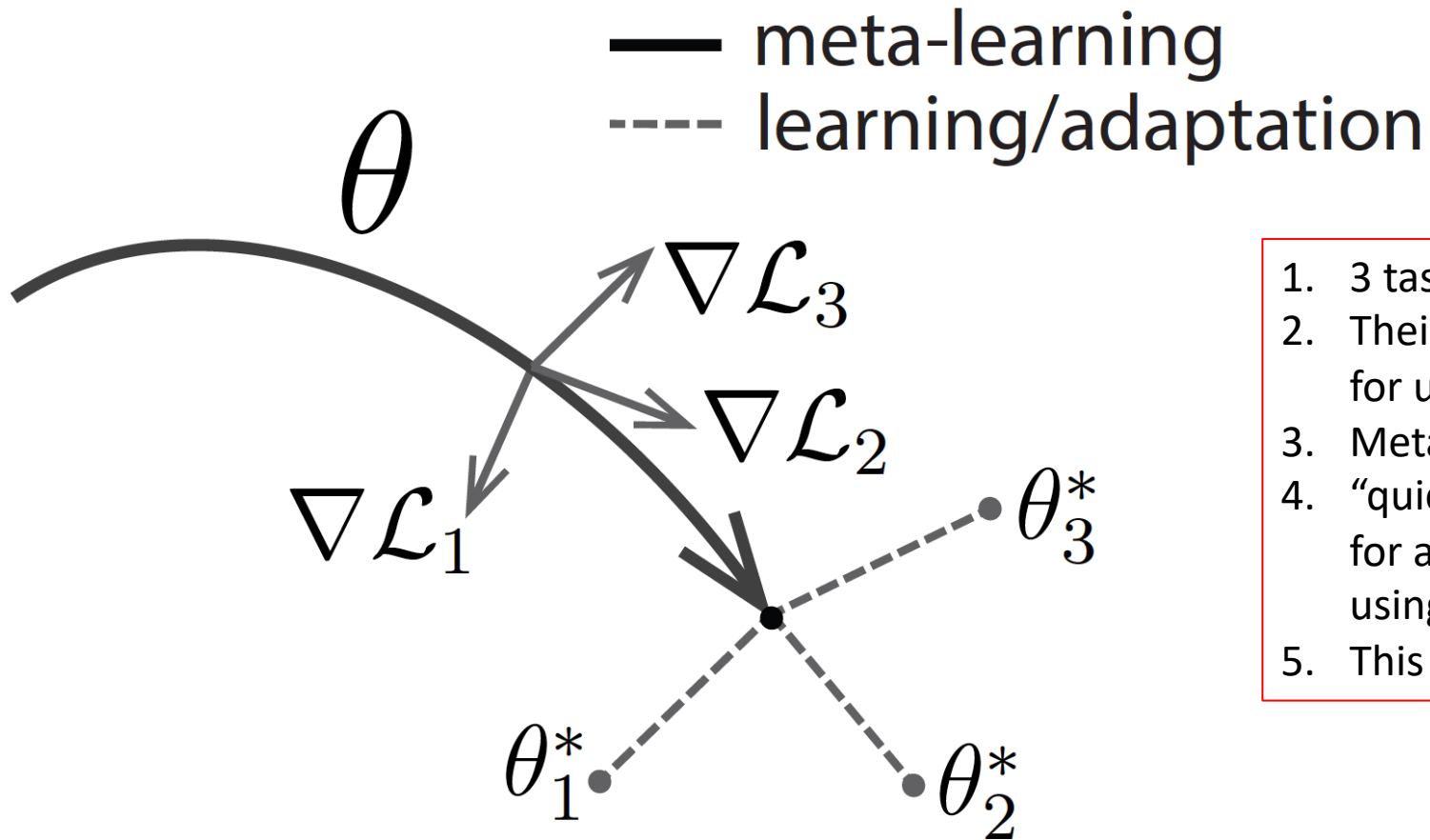


1. 3 tasks
2. Their datasets may have different directions/gradients for updating current model
3. Meta learning determines one update direction
4. "quickly solve" new tasks: after **meta learning update**, for any new task, we can simply update a few steps using the data from the corresponding task

Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." In *International Conference on Machine Learning*, pp. 1126-1135. PMLR, 2017.

<http://proceedings.mlr.press/v70/finn17a/finn17a.pdf>

Meta-learning

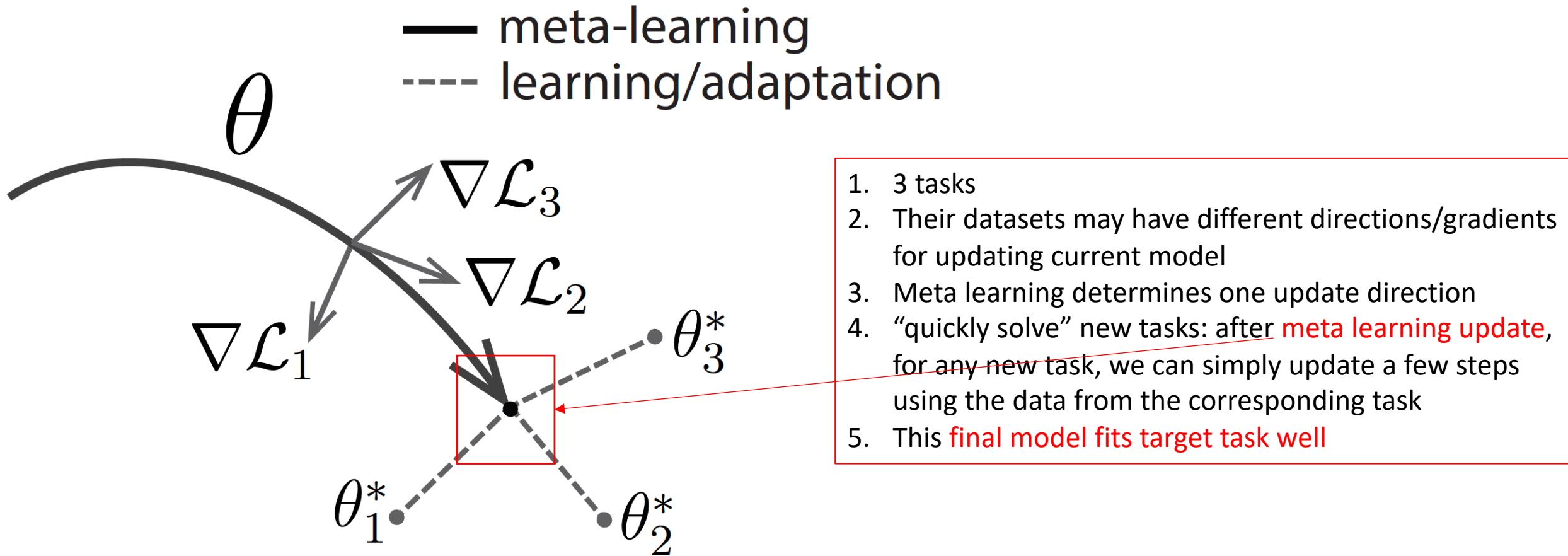


1. 3 tasks
2. Their datasets may have different directions/gradients for updating current model
3. Meta learning determines one update direction
4. "quickly solve" new tasks: after **meta learning update**, for any new task, we can simply update a few steps using the data from the corresponding task
5. This **final model fits target task well**

Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." In *International Conference on Machine Learning*, pp. 1126-1135. PMLR, 2017.

<http://proceedings.mlr.press/v70/finn17a/finn17a.pdf>

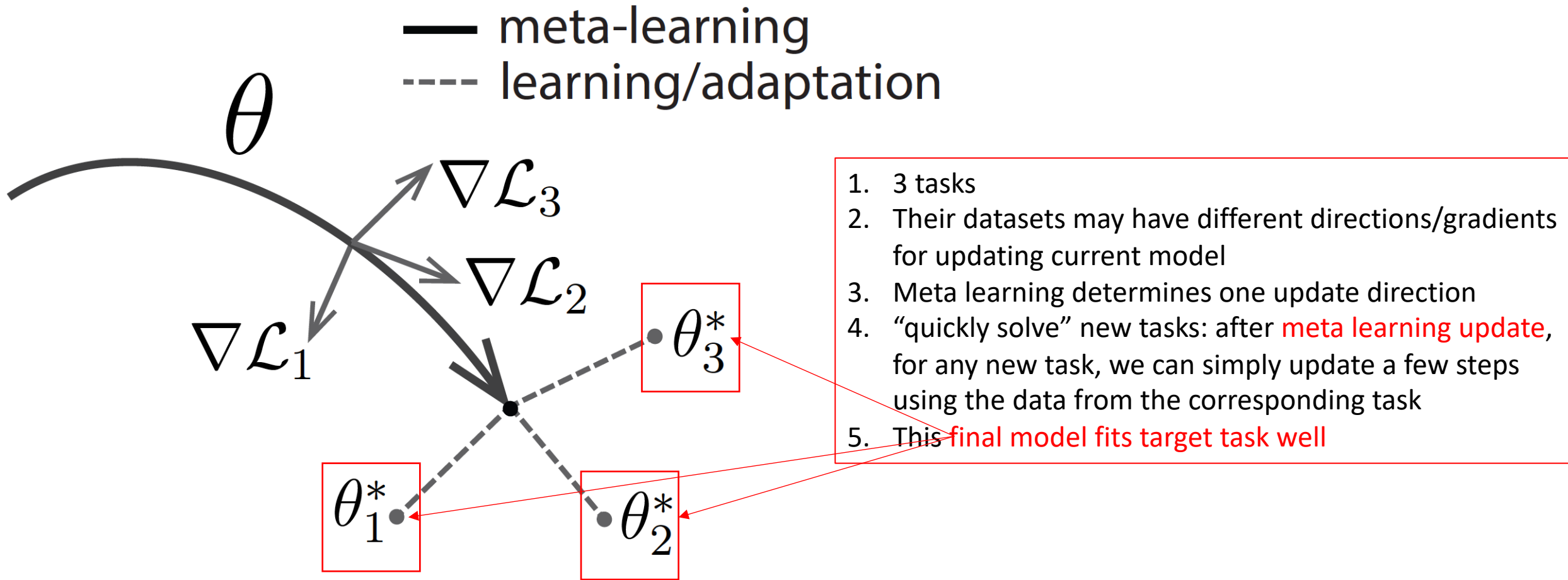
Meta-learning



Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." In *International Conference on Machine Learning*, pp. 1126-1135. PMLR, 2017.

<http://proceedings.mlr.press/v70/finn17a/finn17a.pdf>

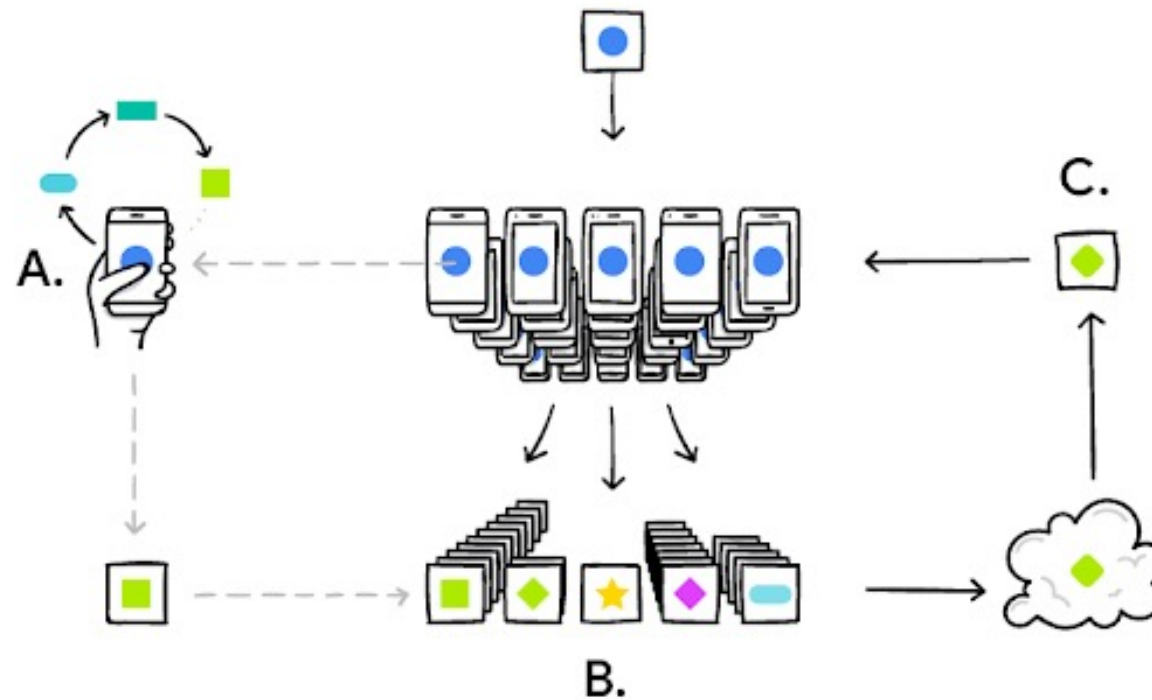
Meta-learning



Finn, Chelsea, Pieter Abbeel, and Sergey Levine. "Model-agnostic meta-learning for fast adaptation of deep networks." In *International Conference on Machine Learning*, pp. 1126-1135. PMLR, 2017.

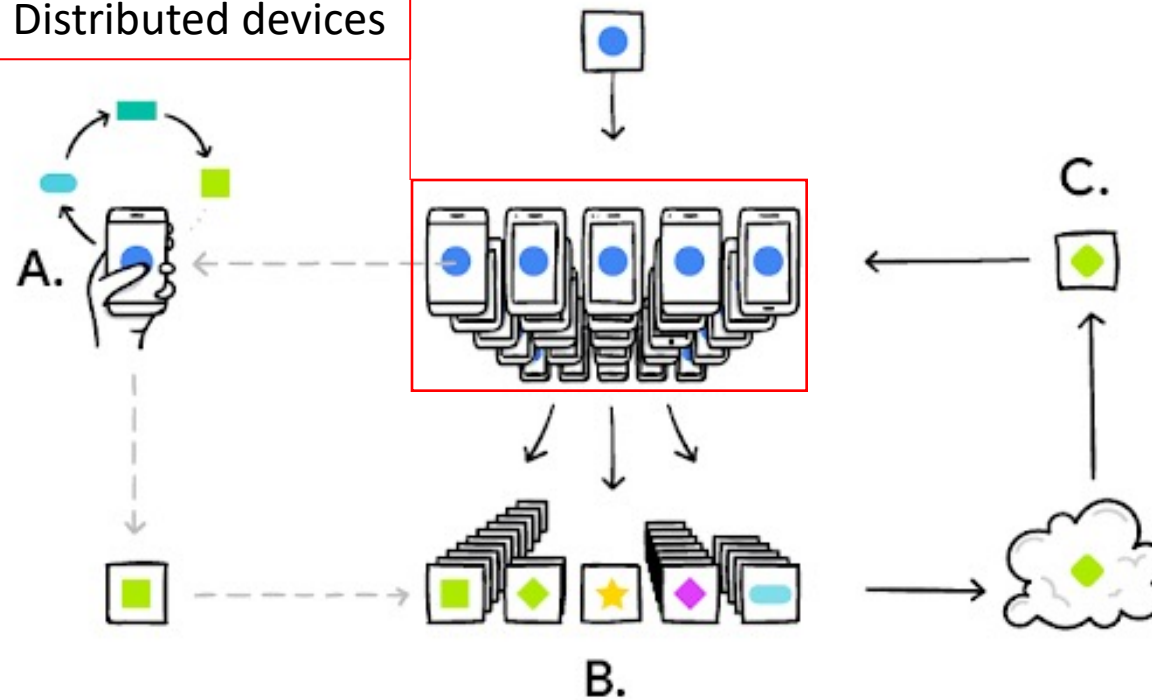
<http://proceedings.mlr.press/v70/finn17a/finn17a.pdf>

Federated learning



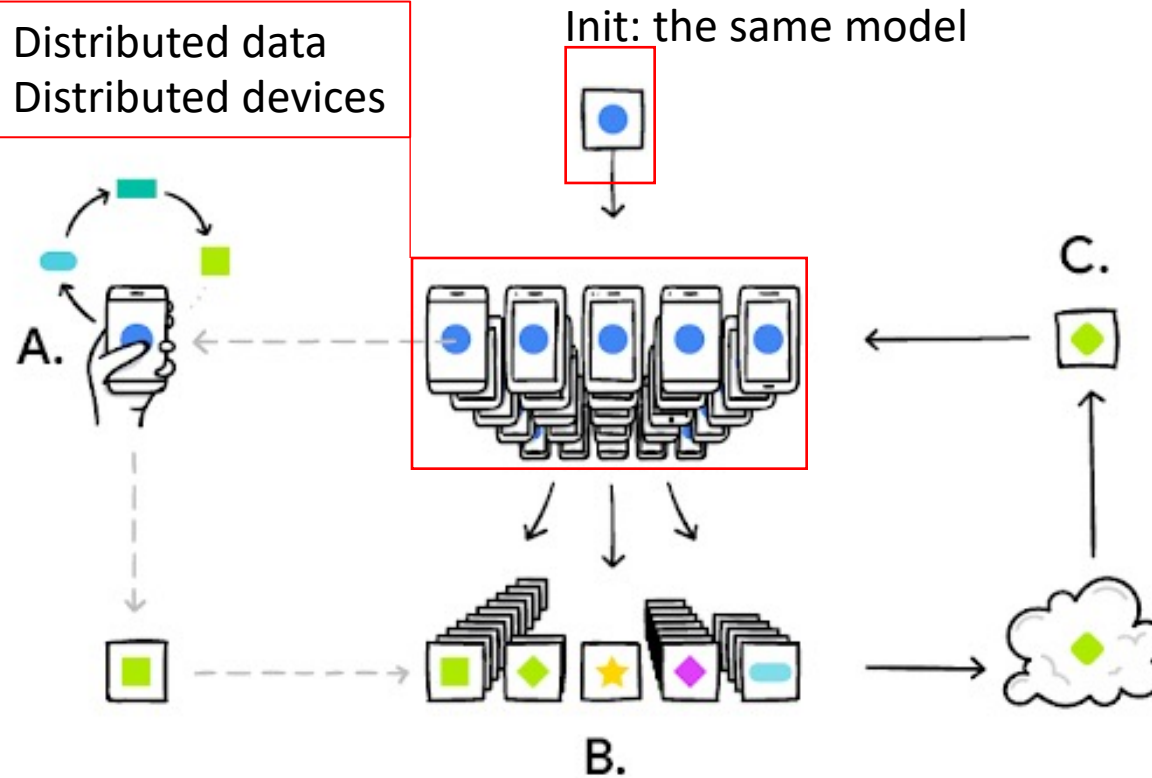
Federated learning

1. Distributed data
2. Distributed devices

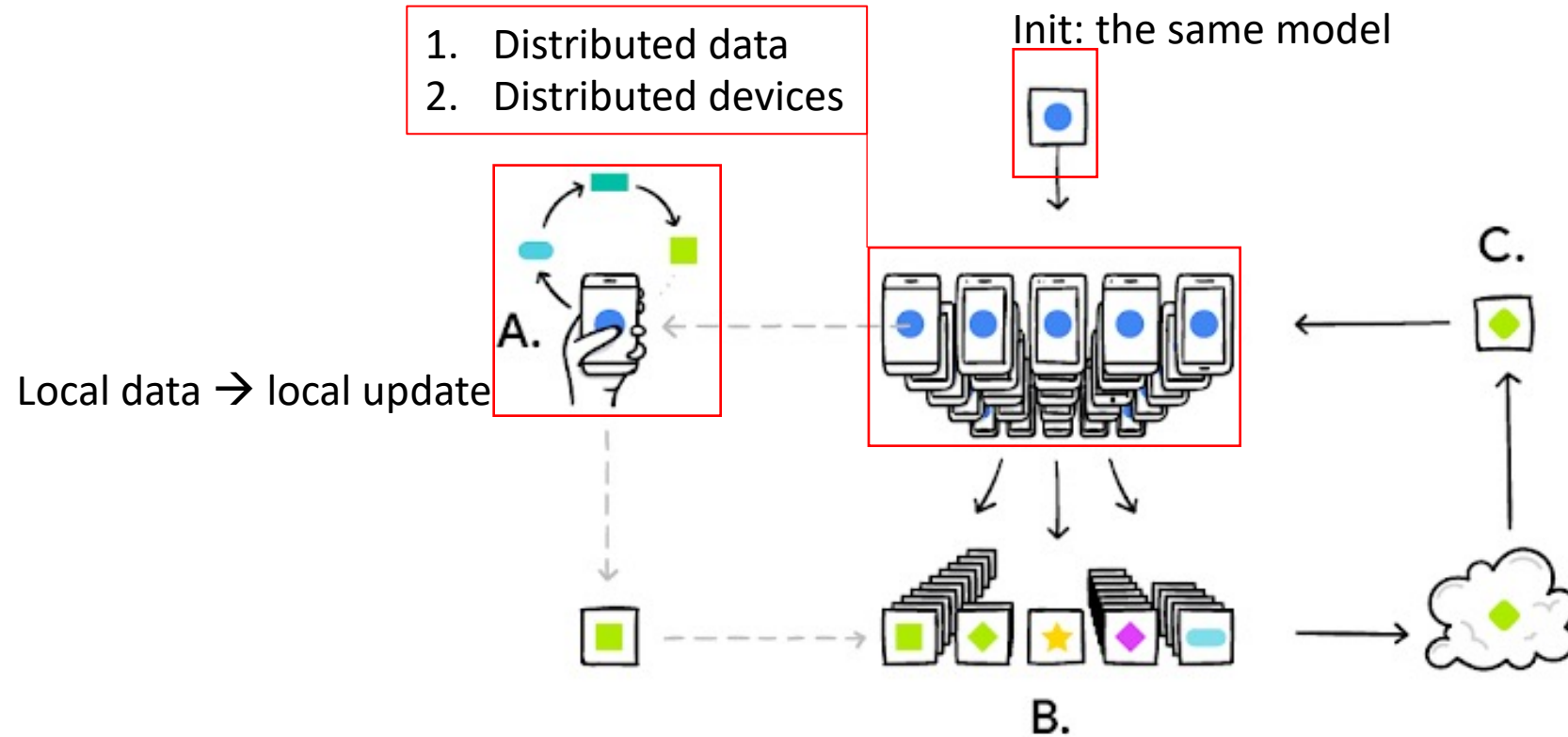


Federated learning

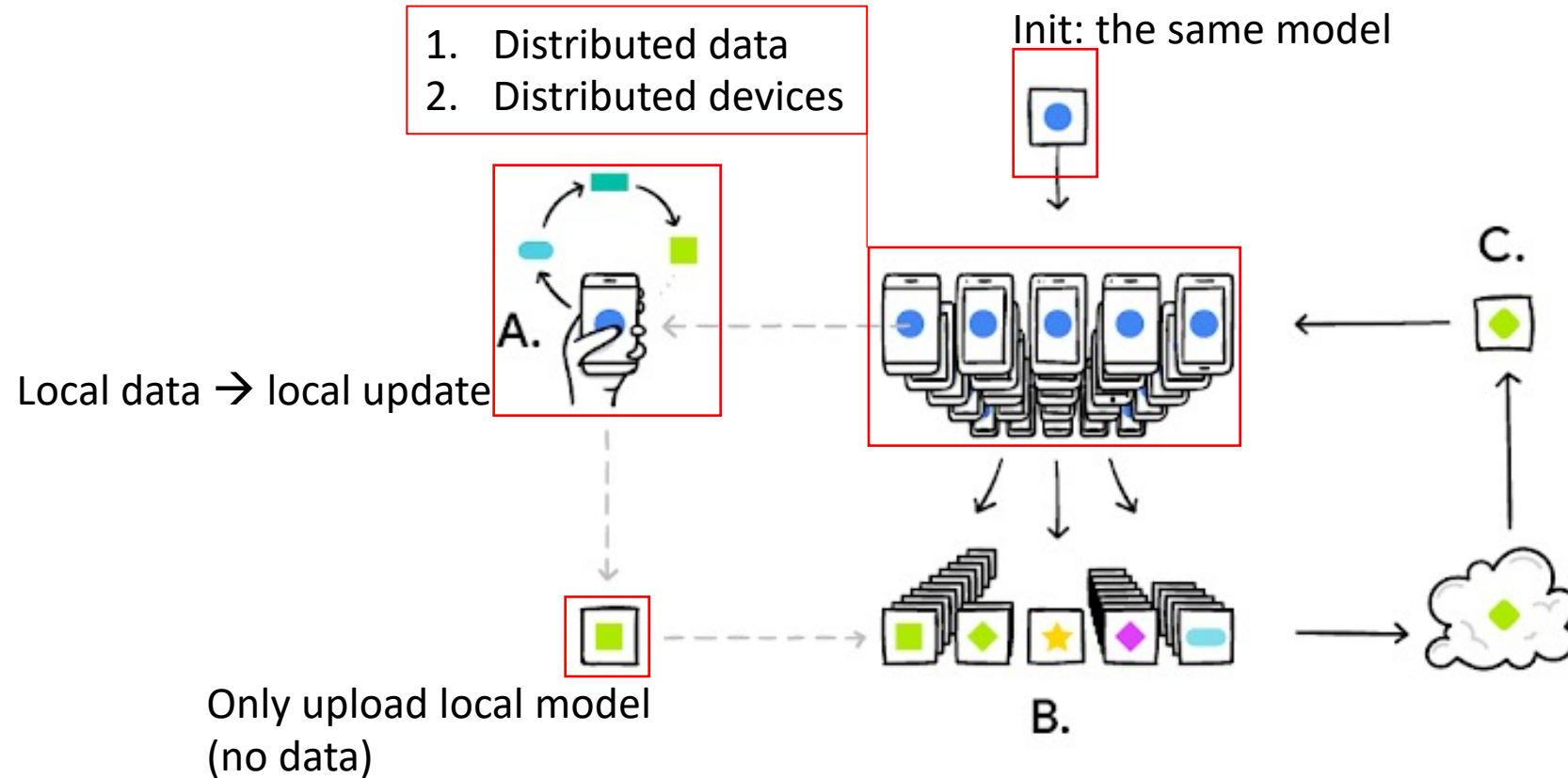
1. Distributed data
2. Distributed devices



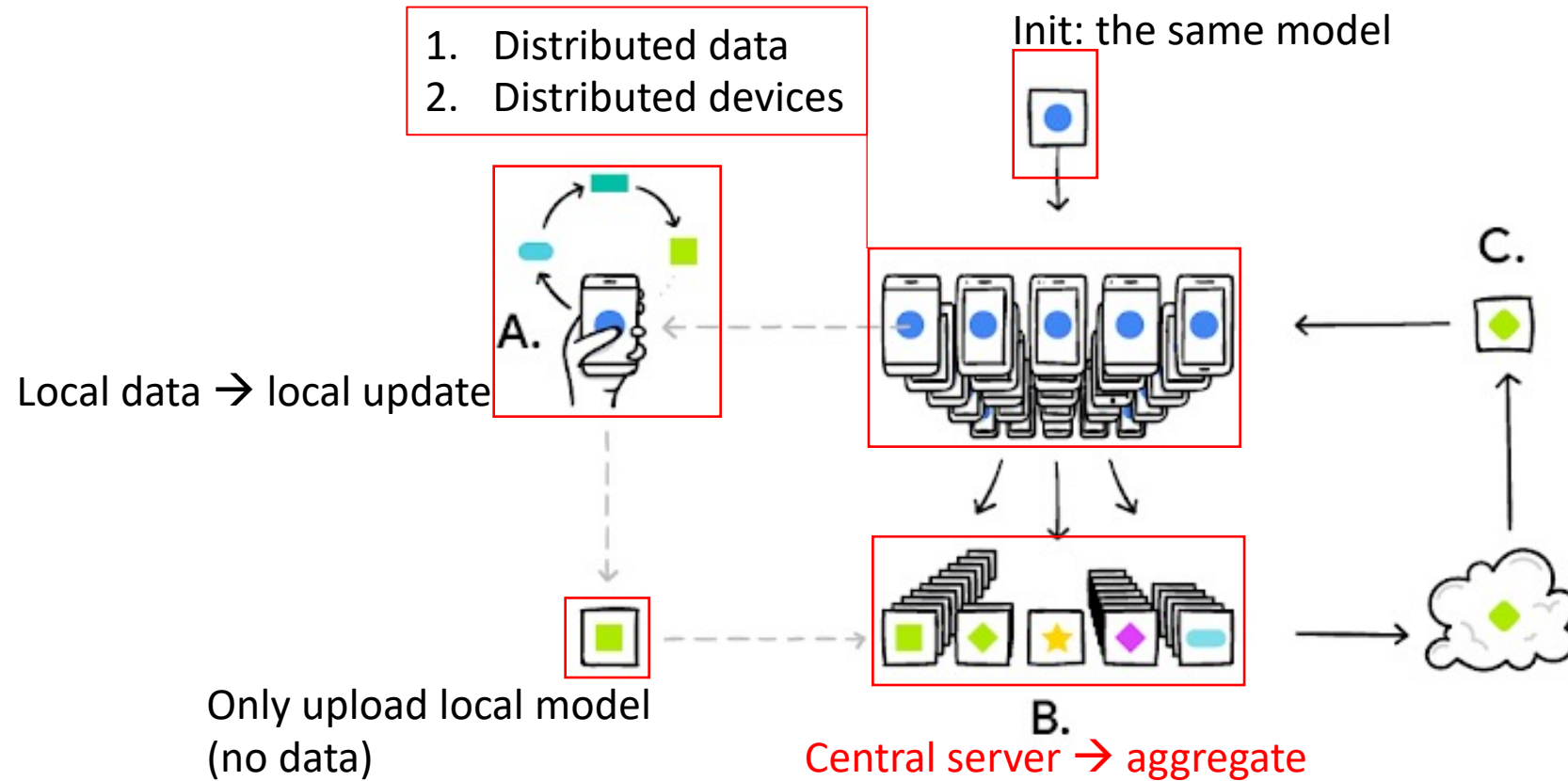
Federated learning



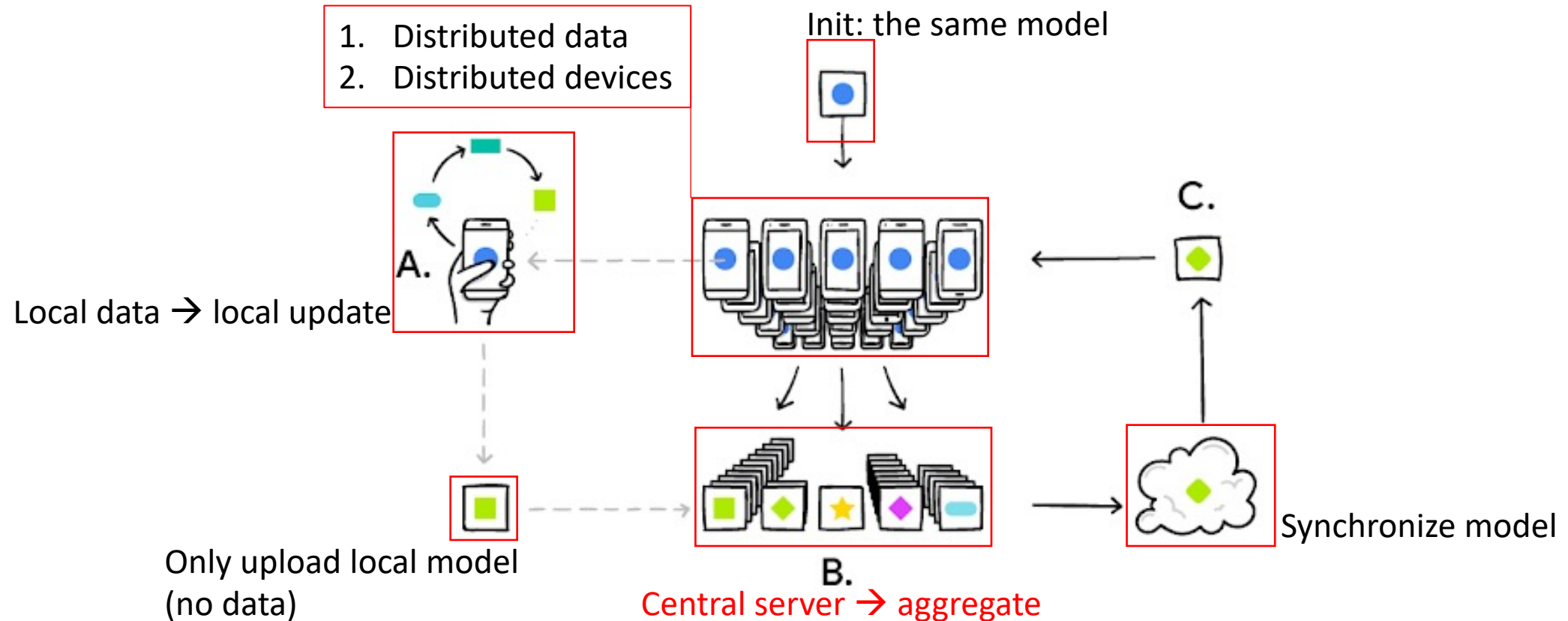
Federated learning



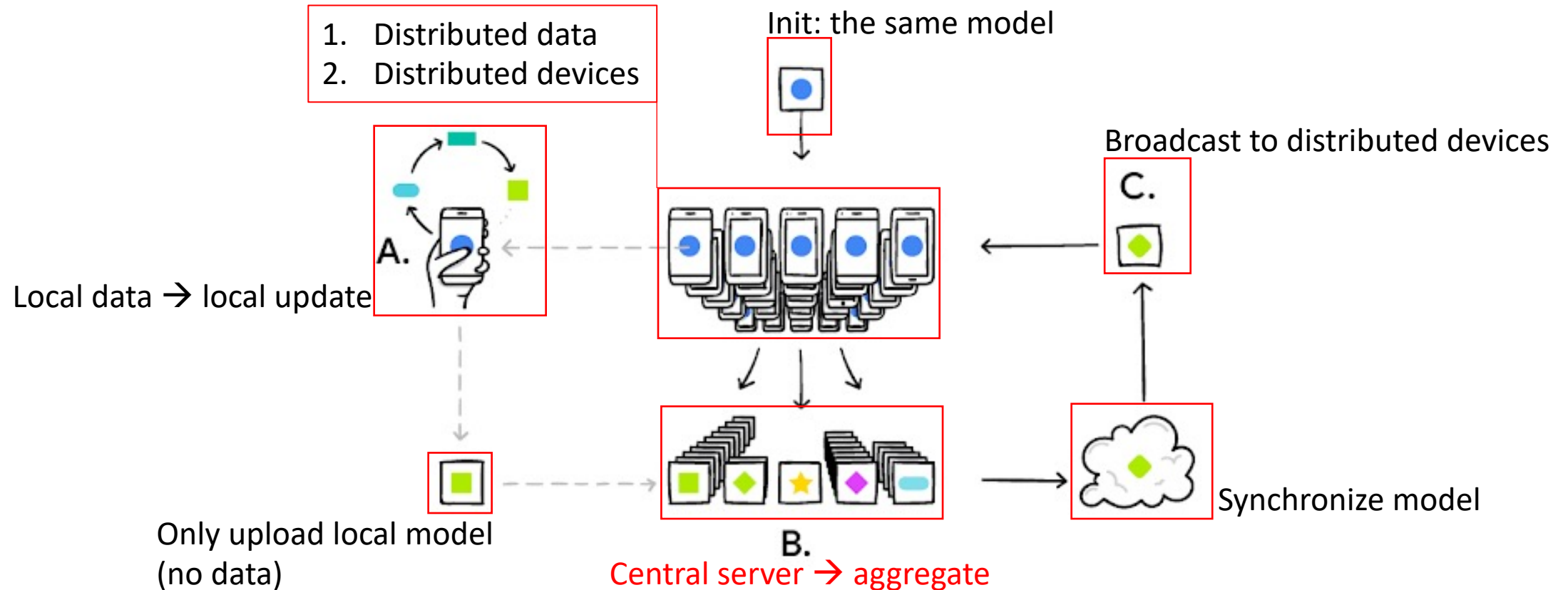
Federated learning



Federated learning



Federated learning



Federated learning via meta learning

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\mathbf{w} - \alpha \nabla \mathcal{L}_i(\mathbf{w})).$$

Federated learning via meta learning

n devices/tasks

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\mathbf{w} - \alpha \nabla \mathcal{L}_i(\mathbf{w})).$$

Federated learning via meta learning

n devices/tasks

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\mathbf{w} - \alpha \nabla \mathcal{L}_i(\mathbf{w})).$$

Federated learning via meta learning

n devices/tasks

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\mathbf{w} - \alpha \nabla \mathcal{L}_i(\mathbf{w})).$$

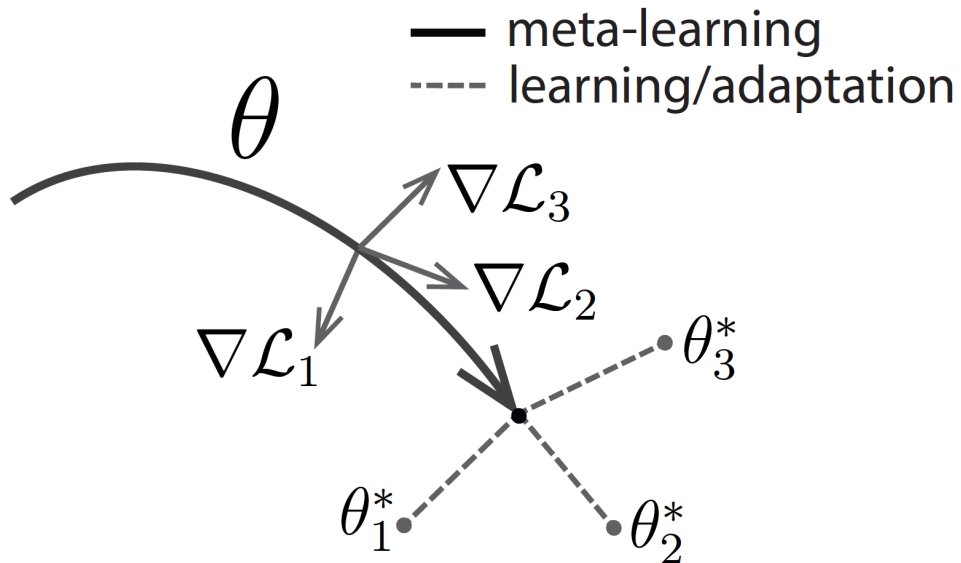
Local model \leftarrow using local data for updating (a single step)

Federated learning via meta learning

n devices/tasks

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\mathbf{w} - \alpha \nabla \mathcal{L}_i(\mathbf{w})).$$

Local model \leftarrow using local data for updating (a single step)

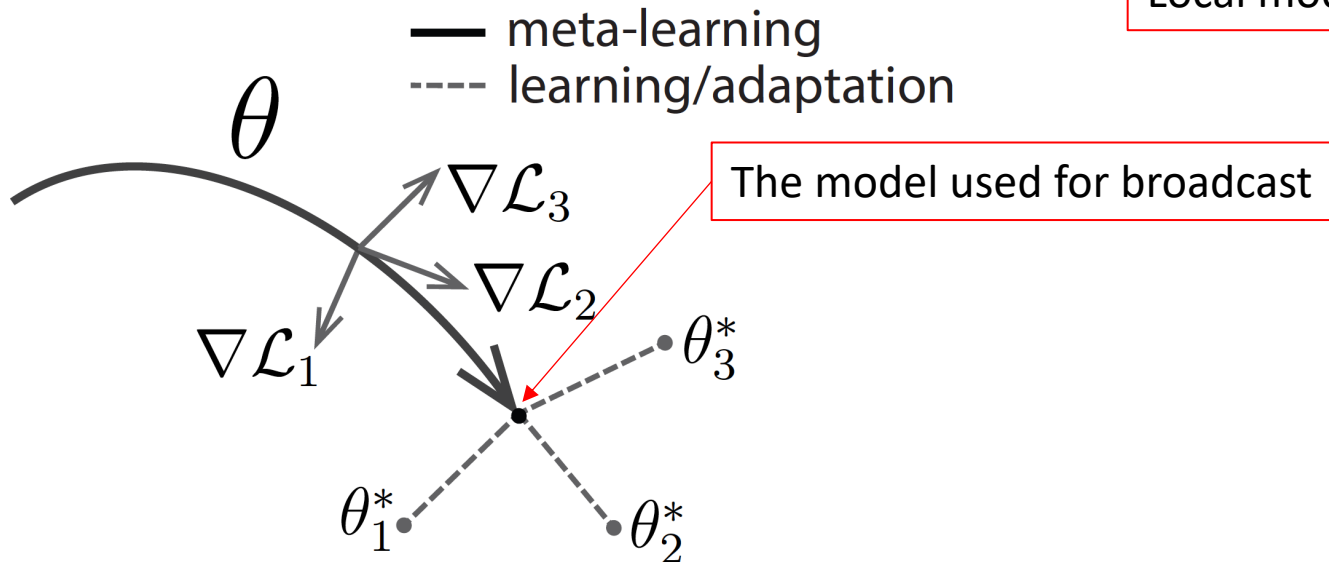


Federated learning via meta learning

n devices/tasks

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\mathbf{w} - \alpha \nabla \mathcal{L}_i(\mathbf{w})).$$

Local model \leftarrow using local data for updating (a single step)

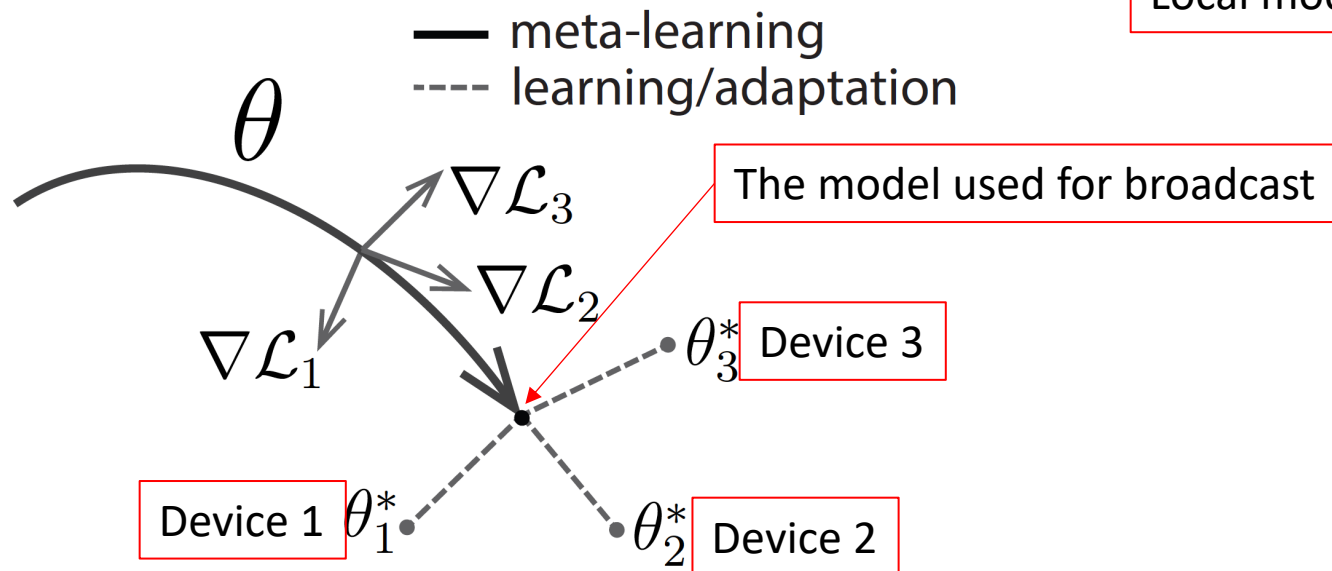


Federated learning via meta learning

n devices/tasks

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\mathbf{w} - \alpha \nabla \mathcal{L}_i(\mathbf{w})).$$

Local model \leftarrow using local data for updating (a single step)

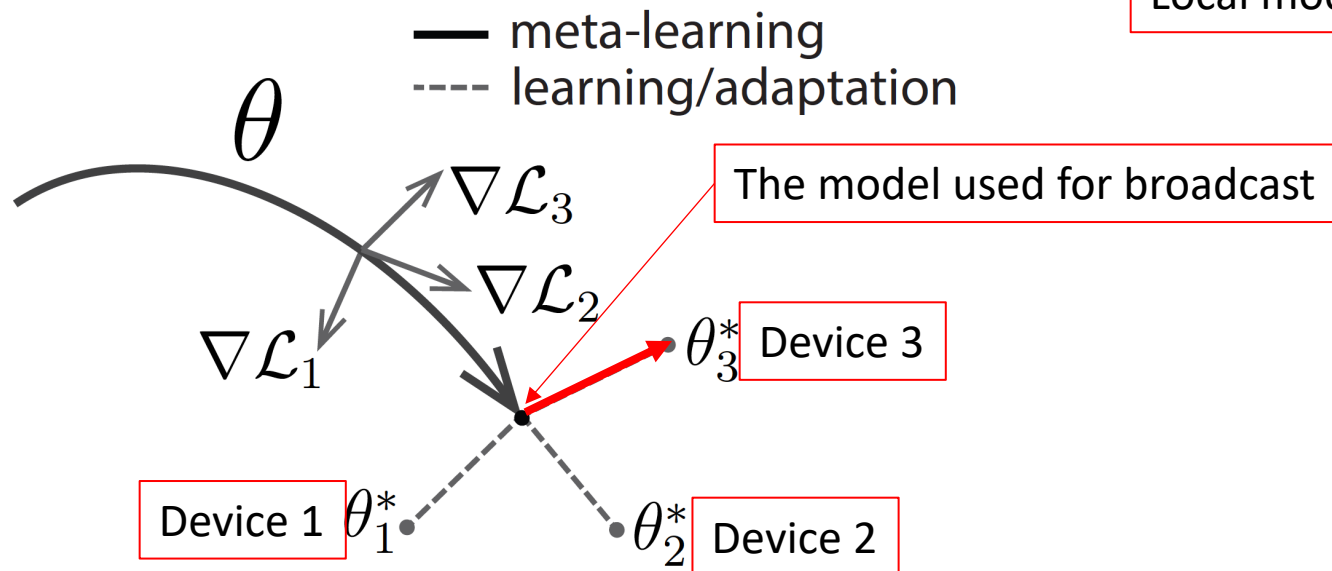


Federated learning via meta learning

n devices/tasks

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\mathbf{w} - \alpha \nabla \mathcal{L}_i(\mathbf{w})).$$

Local model \leftarrow using local data for updating (a single step)

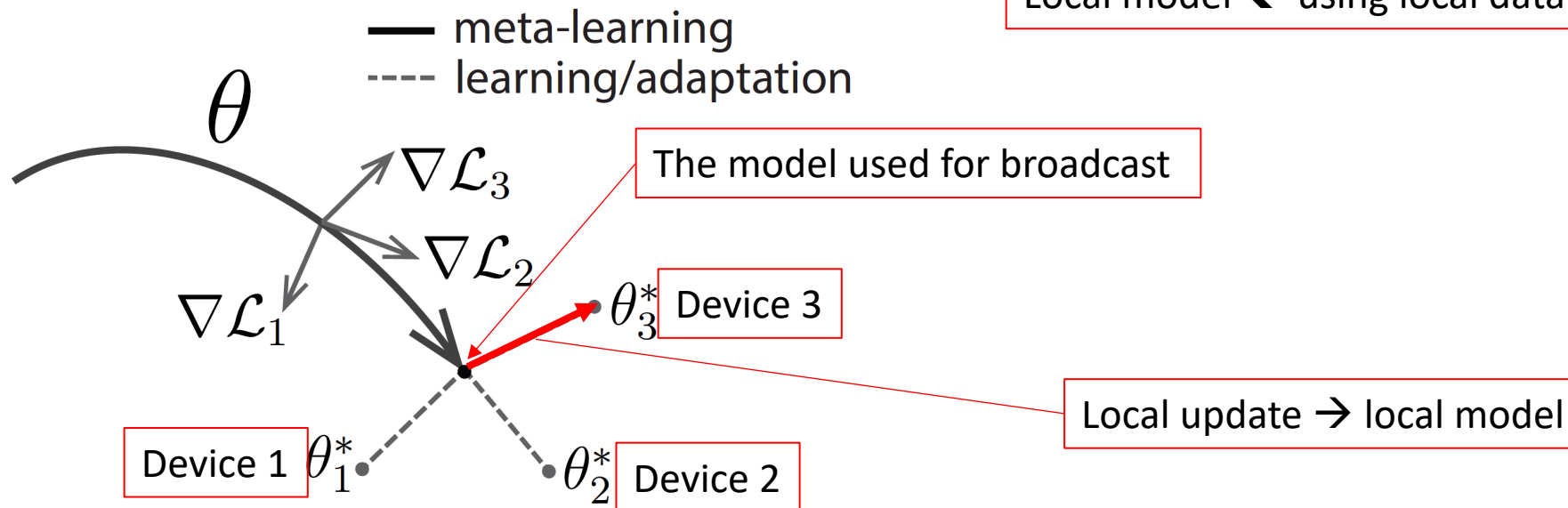


Federated learning via meta learning

n devices/tasks

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\mathbf{w} - \alpha \nabla \mathcal{L}_i(\mathbf{w})).$$

Local model \leftarrow using local data for updating (a single step)

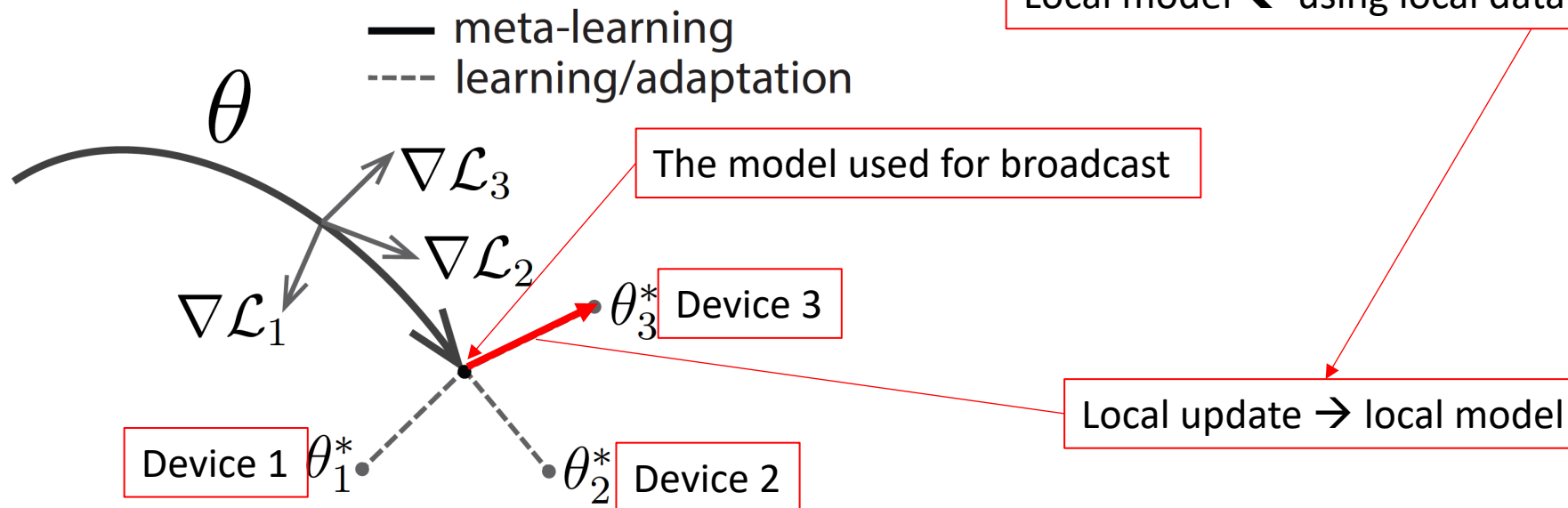


Federated learning via meta learning

n devices/tasks

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\mathbf{w} - \alpha \nabla \mathcal{L}_i(\mathbf{w})).$$

Local model \leftarrow using local data for updating (a single step)



Federated learning via meta learning

Personalized Federated Learning				
Algorithm	Client Sampling	Sample Complexity	Communication Complexity	Avg. #Data points (K) Per Iteration
Per-FedAvg (Fallah et al., 2020b)	\times	$\mathcal{O}(n\epsilon^{-6})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-2})$
Per-FedAvg (This work)	$\checkmark^{(3)}$	$\mathcal{O}(\epsilon^{-7})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$
LocalMOML (This work)	\times	$\mathcal{O}(n\epsilon^{-5})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(1)$
LocalMOML (This work)	\checkmark	$\mathcal{O}(\epsilon^{-6})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(1)$

Federated learning via meta learning

Personalized Federated Learning				
Algorithm	Client Sampling	Sample Complexity	Communication Complexity	Avg. #Data points (K) Per Iteration
Per-FedAvg (Fallah et al., 2020b)	\times	$\mathcal{O}(n\epsilon^{-6})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-2})$
Per-FedAvg (This work)	$\checkmark^{(3)}$	$\mathcal{O}(\epsilon^{-7})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$
LocalMOML (This work)	\times	$\mathcal{O}(n\epsilon^{-5})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(1)$
LocalMOML (This work)	\checkmark	$\mathcal{O}(\epsilon^{-6})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(1)$

Number of data samples (on local devices)

Federated learning via meta learning

Personalized Federated Learning

Algorithm	Client Sampling	Sample Complexity	Communication Complexity	Avg. #Data points (K) Per Iteration
Per-FedAvg (Fallah et al., 2020b)	\times	$\mathcal{O}(n\epsilon^{-6})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-2})$
Per-FedAvg (This work)	$\checkmark^{(3)}$	$\mathcal{O}(\epsilon^{-7})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$
LocalMOML (This work)	\times	$\mathcal{O}(n\epsilon^{-5})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(1)$
LocalMOML (This work)	\checkmark	$\mathcal{O}(\epsilon^{-6})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(1)$

Number of data samples (on local devices)

Number of broadcasts

Federated learning via meta learning

1. drop the dependence on n
2. reduce order of ϵ

Personalized Federated Learning

Algorithm	Client Sampling	Sample Complexity	Communication Complexity	Avg. #Data points (K) Per Iteration
Per-FedAvg (Fallah et al., 2020b)	\times	$\mathcal{O}(n\epsilon^{-6})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-2})$
Per-FedAvg (This work)	$\checkmark^{(3)}$	$\mathcal{O}(\epsilon^{-7})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$
LocalMOML (This work)	\times	$\mathcal{O}(n\epsilon^{-5})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(1)$
LocalMOML (This work)	\checkmark	$\mathcal{O}(\epsilon^{-6})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(1)$

Number of data samples (on local devices)

Number of broadcasts