

# Machine Learning Beyond Static Datasets

ESSAI 2023



**Dr. Martin Mundt,**

Research Group Leader, TU Darmstadt & hessian.AI

Board Member of Directors, ContinualAI



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

Course: <http://owll-lab.com/teaching/essai-23>



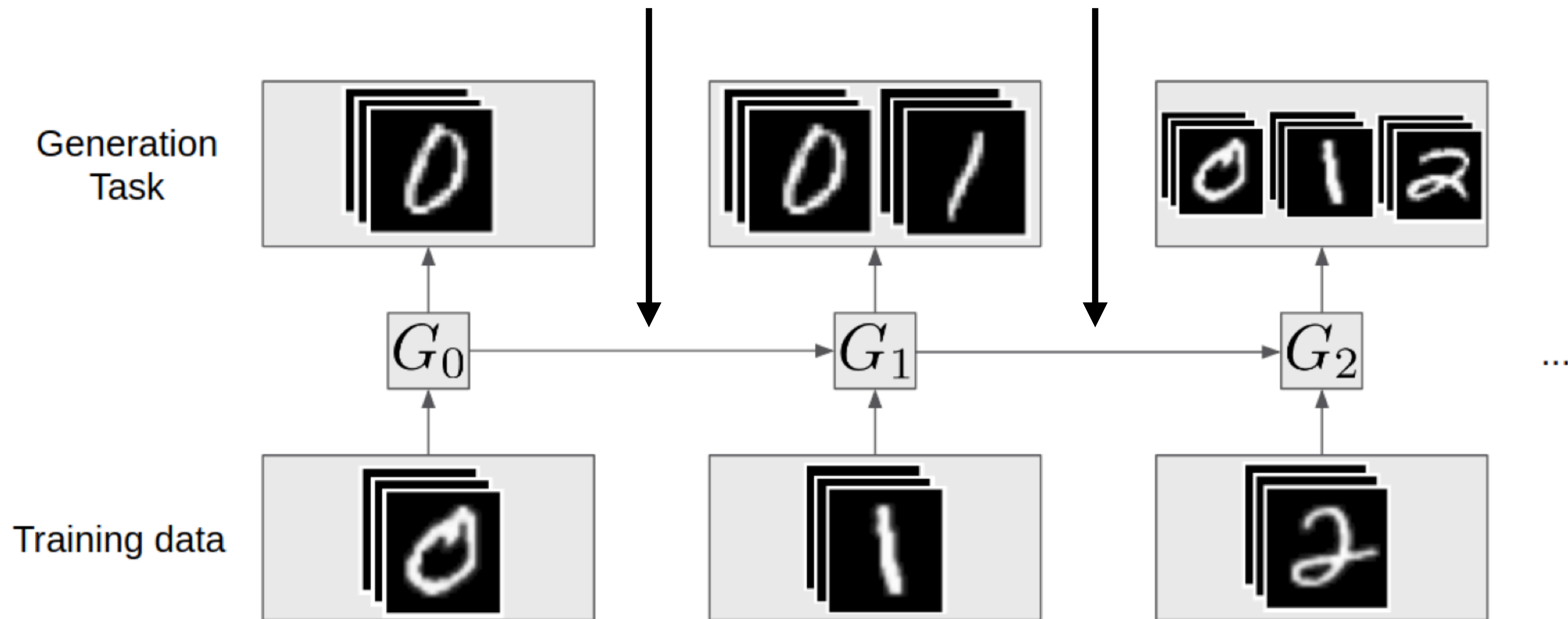
Day 4: The Future  
Data Selection & Learning Curricula



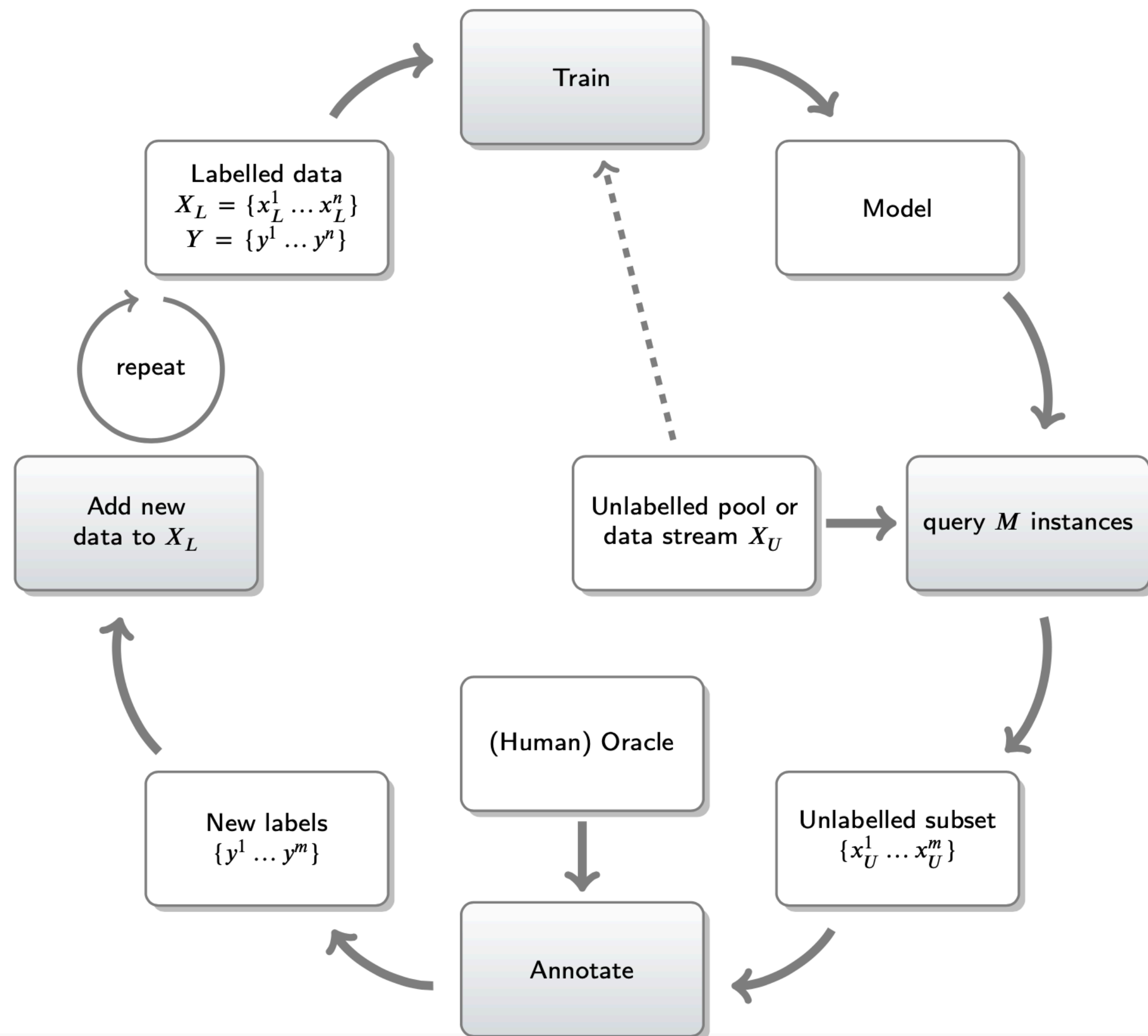
# The future: picking what comes next



Who decides what comes next? A stream? A human? The model? How?



# Active learning



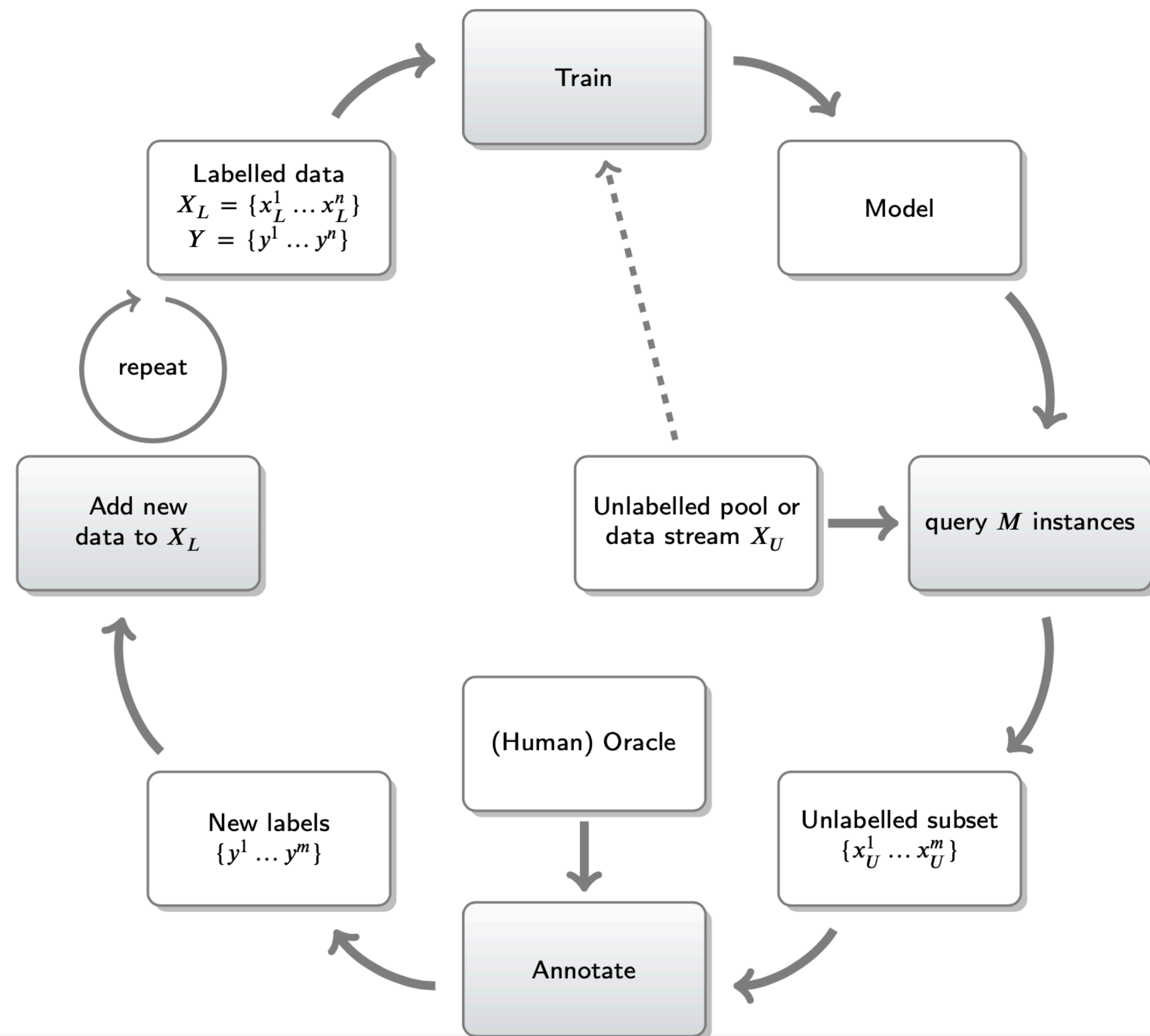
In essence: deciding what new data points are most informative

Also called “**query learning**” with the underlying mechanism called “**acquisition function**”



When querying new data, what are some assumptions & considerations on set-up we can make?

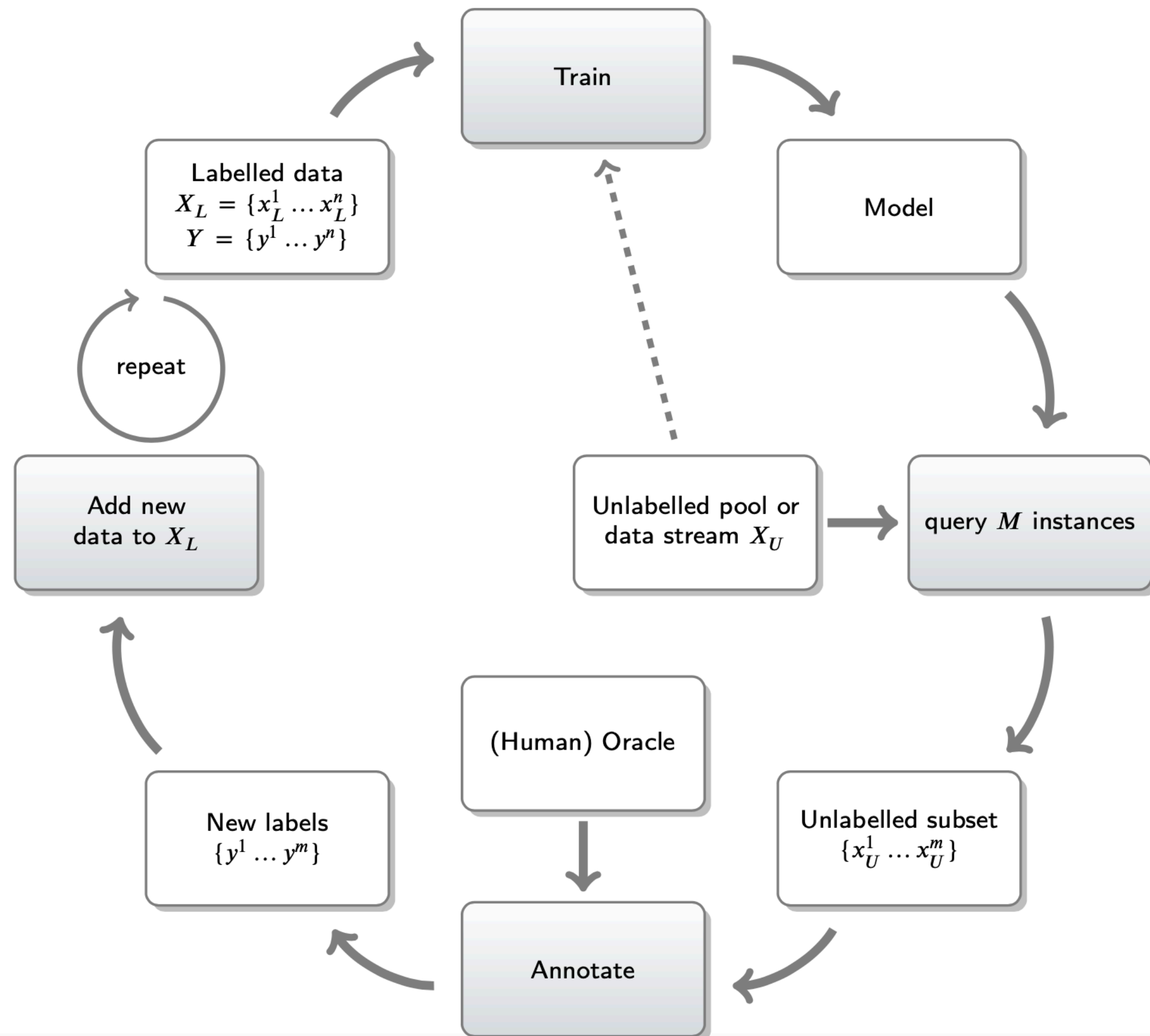
# Active learning



## Many assumptions (non-exhaustive)

- Data is cheap vs. labeling is not?
- Pool of data upfront vs. stream?
- 1 data point vs. batches vs. tasks?
- Accumulate data after selection?
- Re-train vs. continued training?
- Oracle: infallible vs. noisy?

# Majority of “traditional” active learning



## Many assumptions (non-exhaustive)

- Data is cheap vs. labeling is not?
- Pool of data upfront vs. stream?
- 1 data point vs. batches vs. tasks?
- Accumulate data after selection?
- Re-train vs. continued training?
- Oracle: infallible vs. noisy?



What techniques to query data can you think of?

Again, a small tangent:  
discriminative or generative models?



**Discriminative** models could allow for natural ways to assess “novelty” of a new example -> *But caution*: overconfidence phenomena (tomorrow)

**Generative** models learn about the data distribution

-> *But caution*: our parameters only reflect the distribution seen so far!  
(do we make use of a pool that is always available?)

We will see that the choice also depends on the set-up assumption!



# Active learning perspectives



## Version space reduction

reduce the set/space of possible hypotheses  $h : \mathcal{X} \rightarrow \mathcal{Y}$  by removing the ones that are inconsistent with the data

## Uncertainty & heuristics

use the predictions, or maybe even better, uncertainty in the predictions



# Version Space

## Version space (Mitchel 1978)

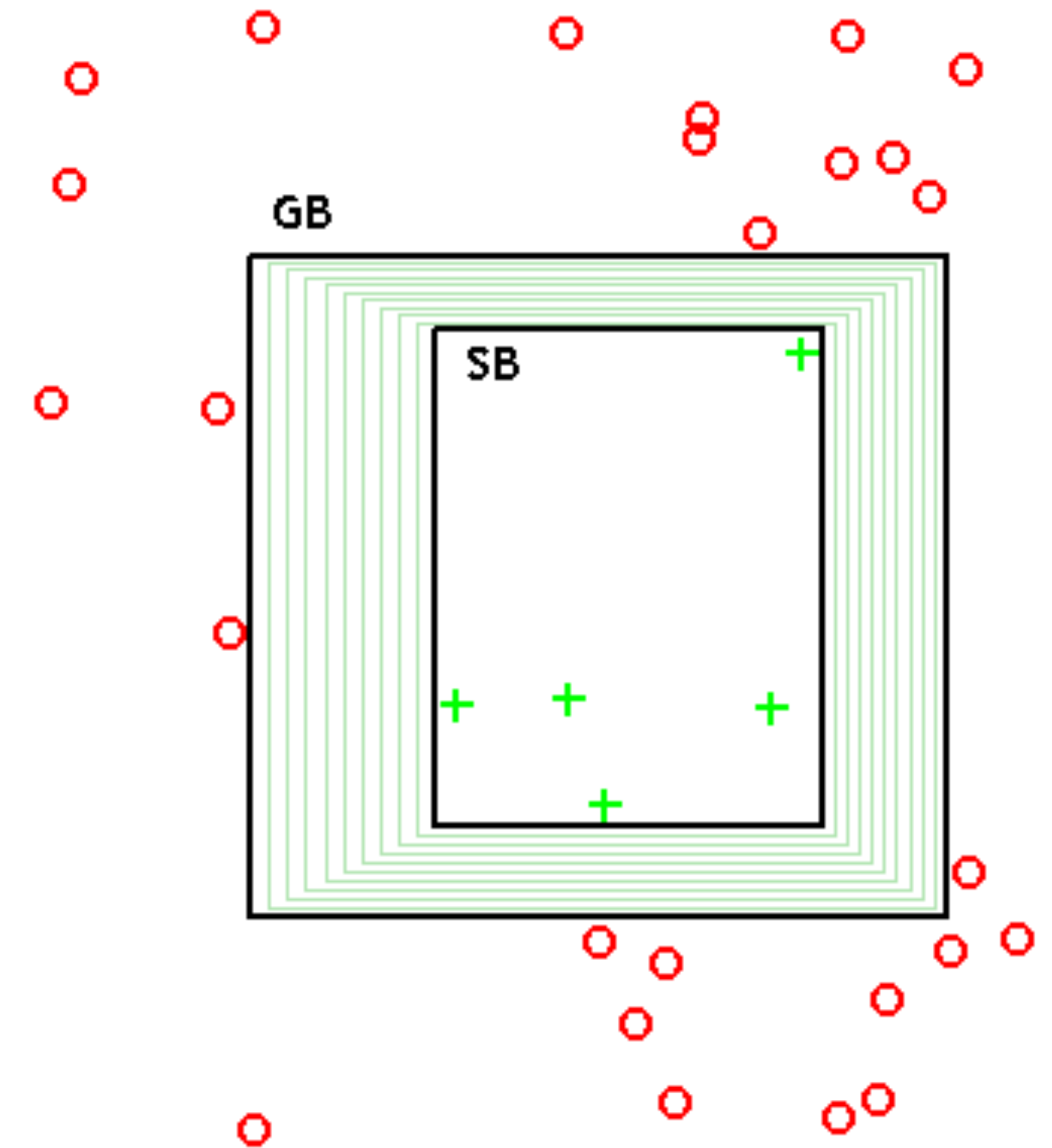


- Assume that there exist hypotheses consistent with the labeled data points  $h : \mathcal{X} \rightarrow \mathcal{Y}$   
**version space:**  $VS(D) = \{h \in H \mid \text{cons}(h, D)\}$

# Version space (Mitchel 1978)



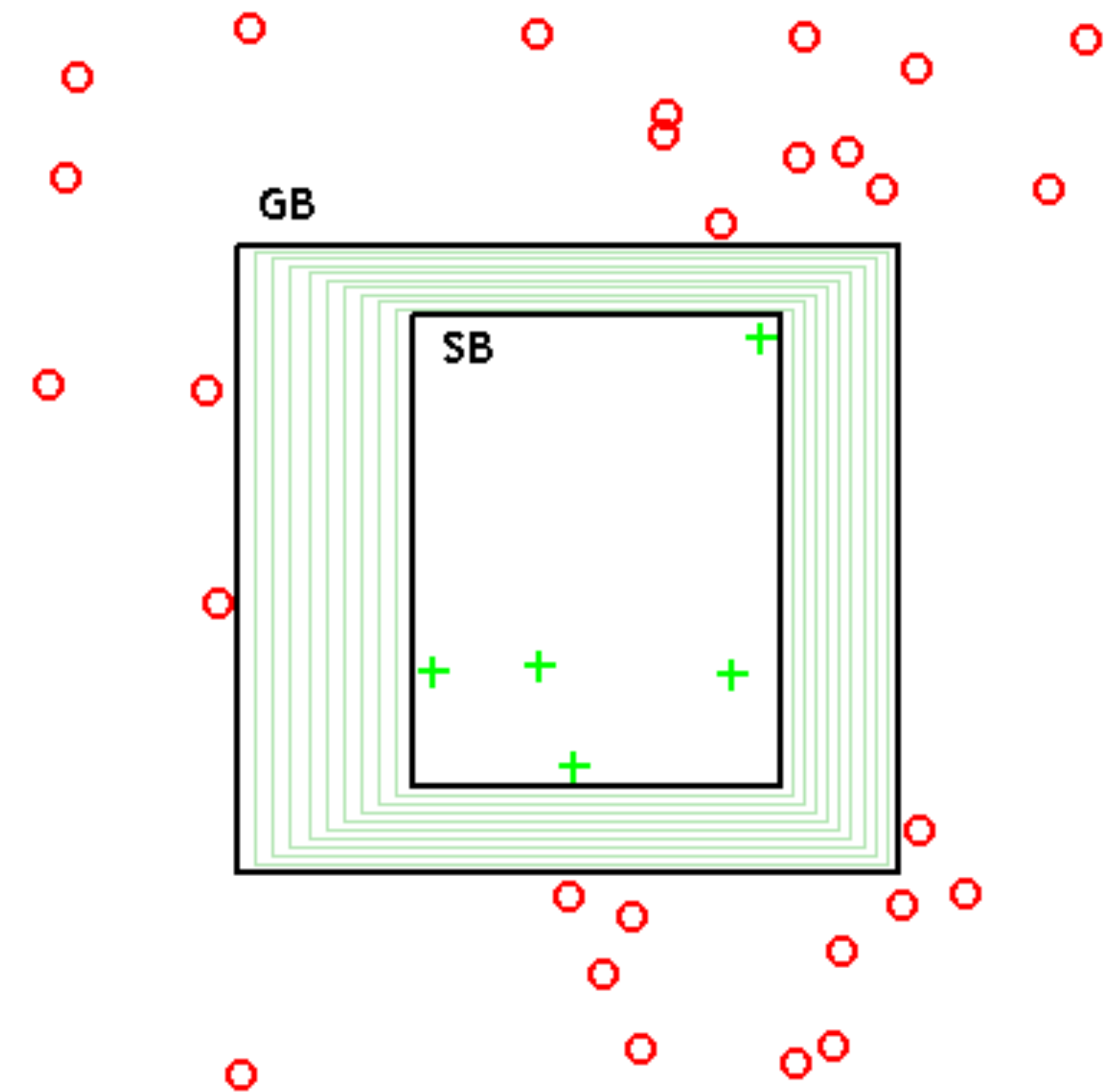
- Assume that there exist hypotheses consistent with the labeled data points  $h : \mathcal{X} \rightarrow \mathcal{Y}$   
**version space:**  $VS(D) = \{h \in H \mid \text{cons}(h, D)\}$



## Version space (Mitchel 1978)



- Assume that there exist hypotheses consistent with the labeled data points  $h : \mathcal{X} \rightarrow \mathcal{Y}$   
**version space:**  $VS(D) = \{h \in H \mid \text{cons}(h, D)\}$
- **Specific hypotheses:** cover positive examples & as little remaining feature space
- **General hypotheses:** cover positive examples & as much of the remaining feature space
- **Version space:** represented as green rectangles



# Version space reduction



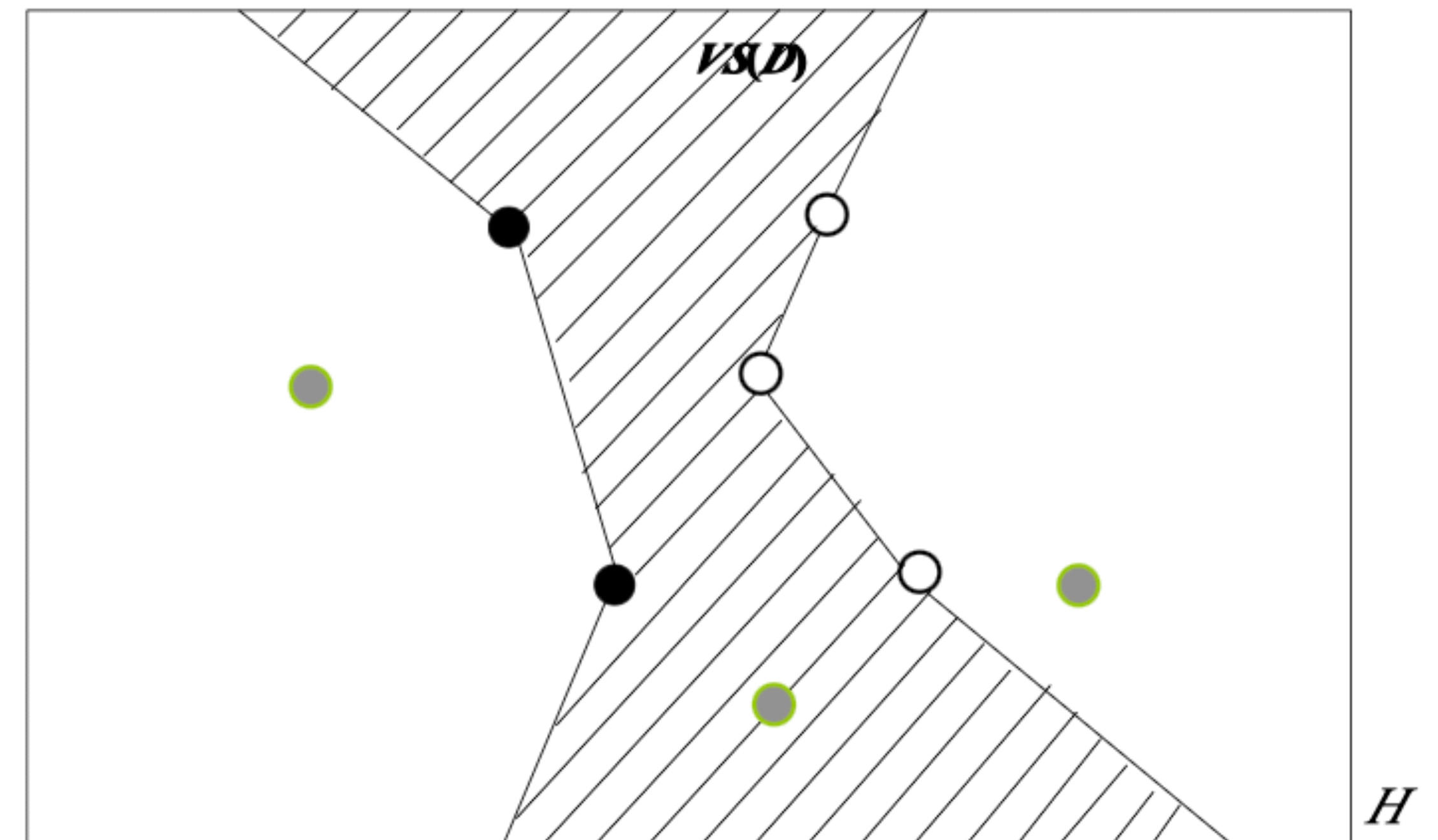
“Generalization as Search”, Mitchell 1982

We could query such that the **version**

**space**:  $VS(D) = \{h \in H \mid \text{cons}(h, D)\}$  ,i.e.

the set of consistent hypotheses,

quickly **gets reduced**

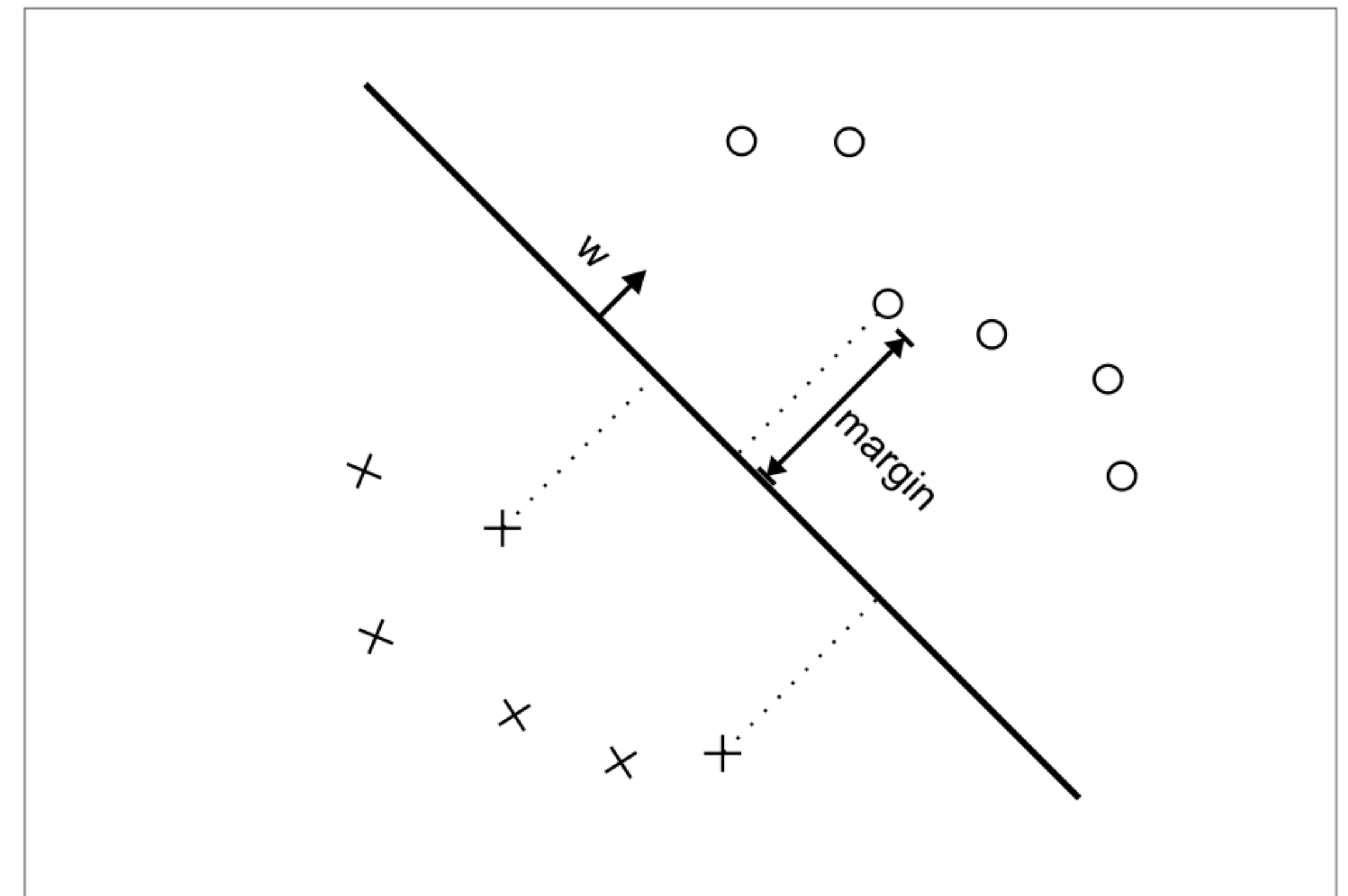


# Active learning with support vector machines (SVM)

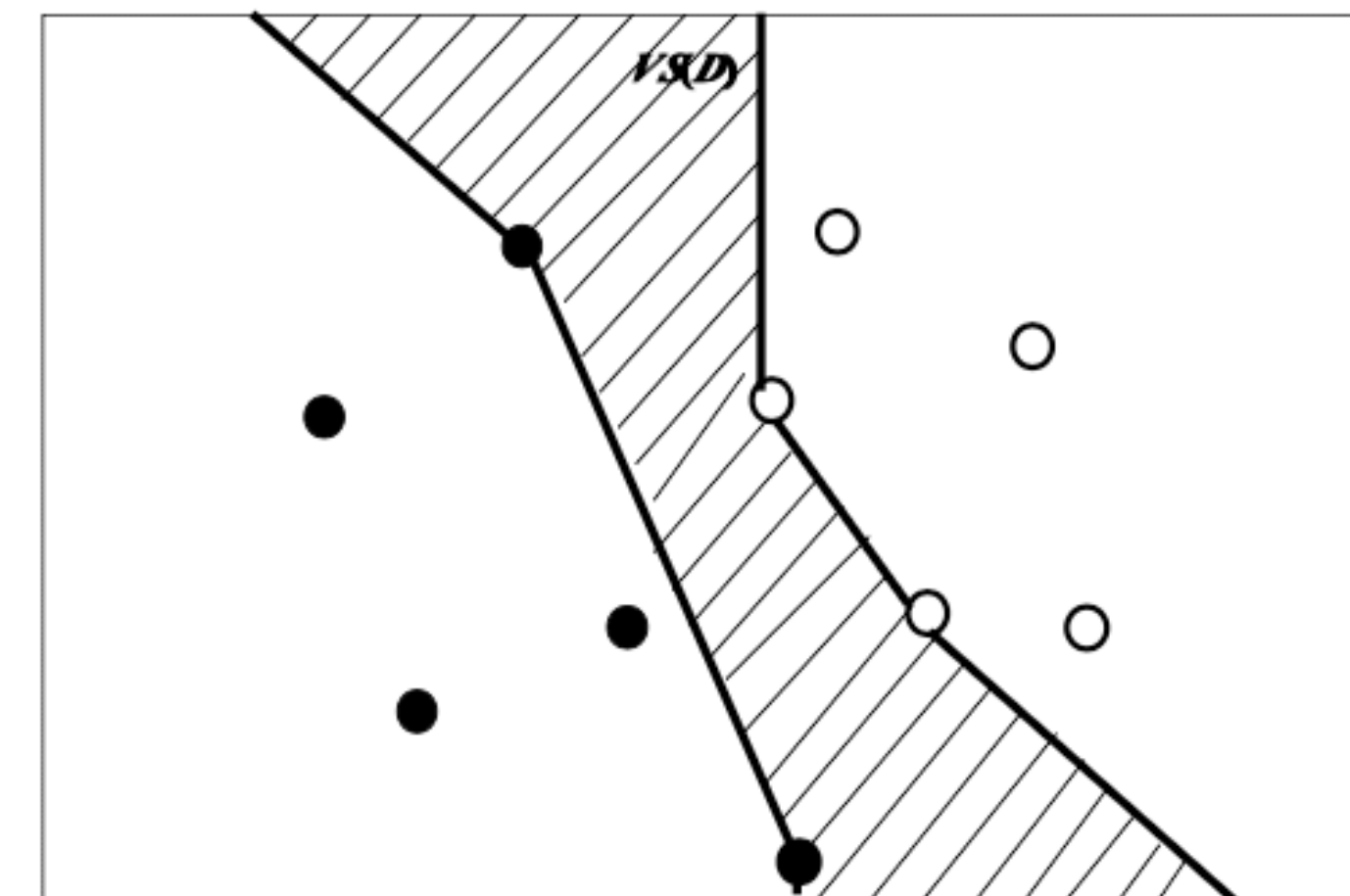
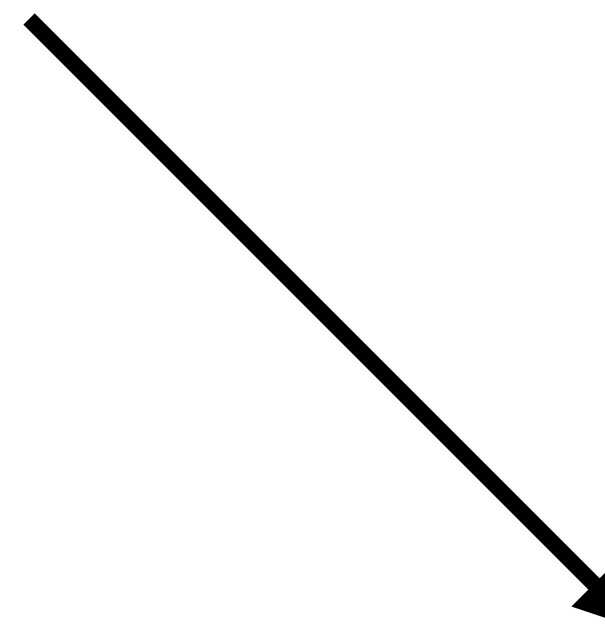
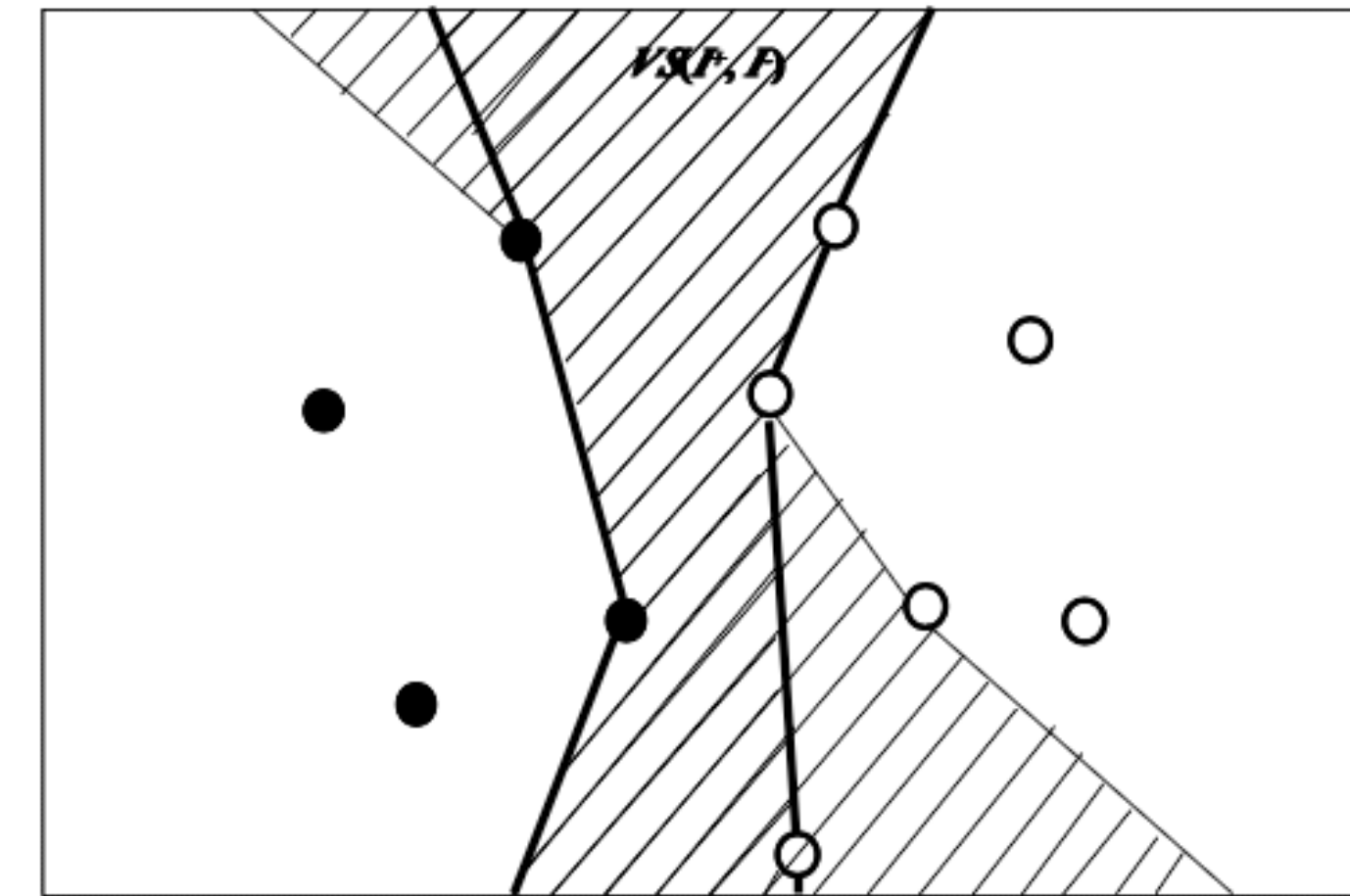
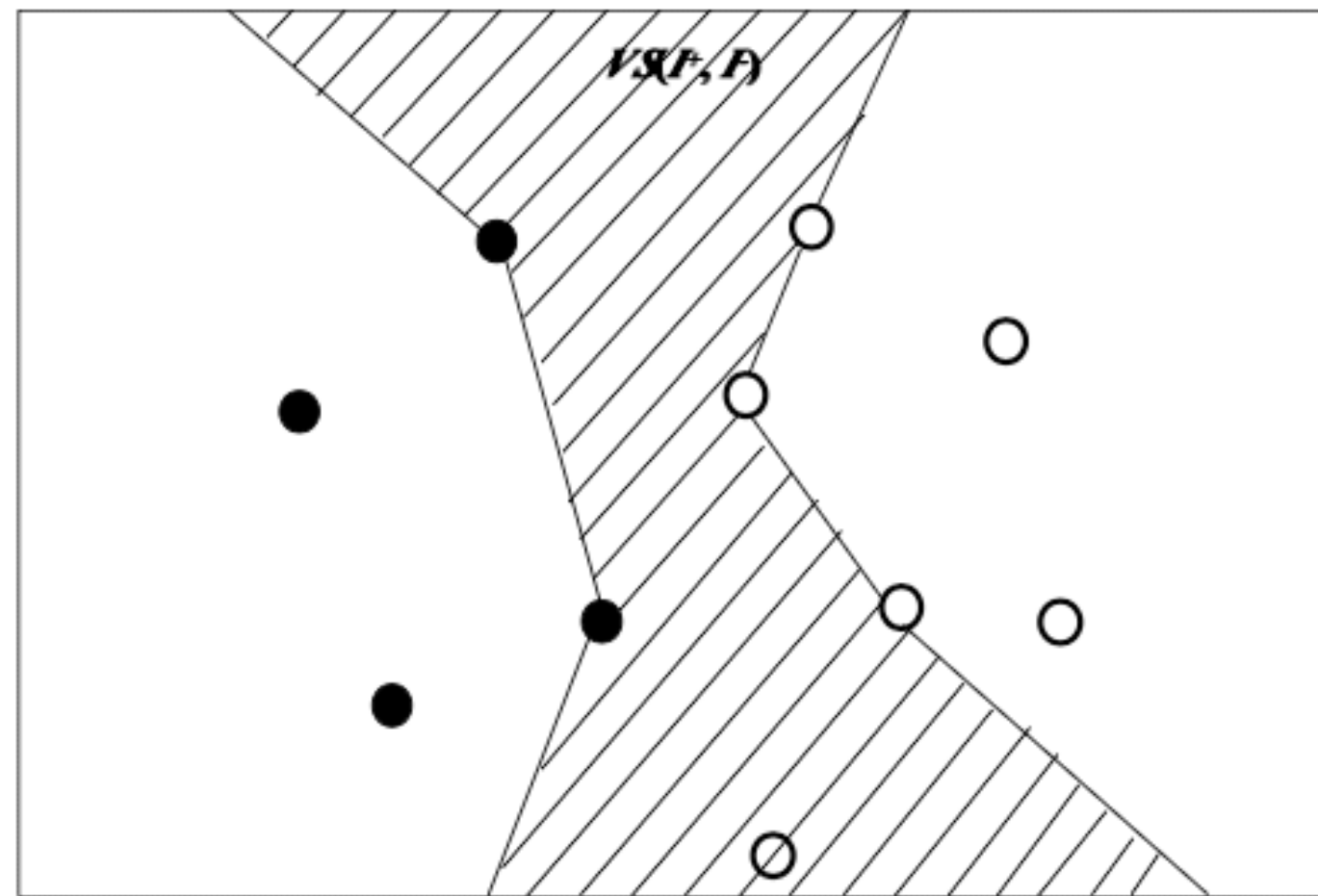


There are some models in which we can do this. Why?

- Hyperplane chosen to **maximize margin to closest instances**: the support vectors



# Active learning with SVM version space



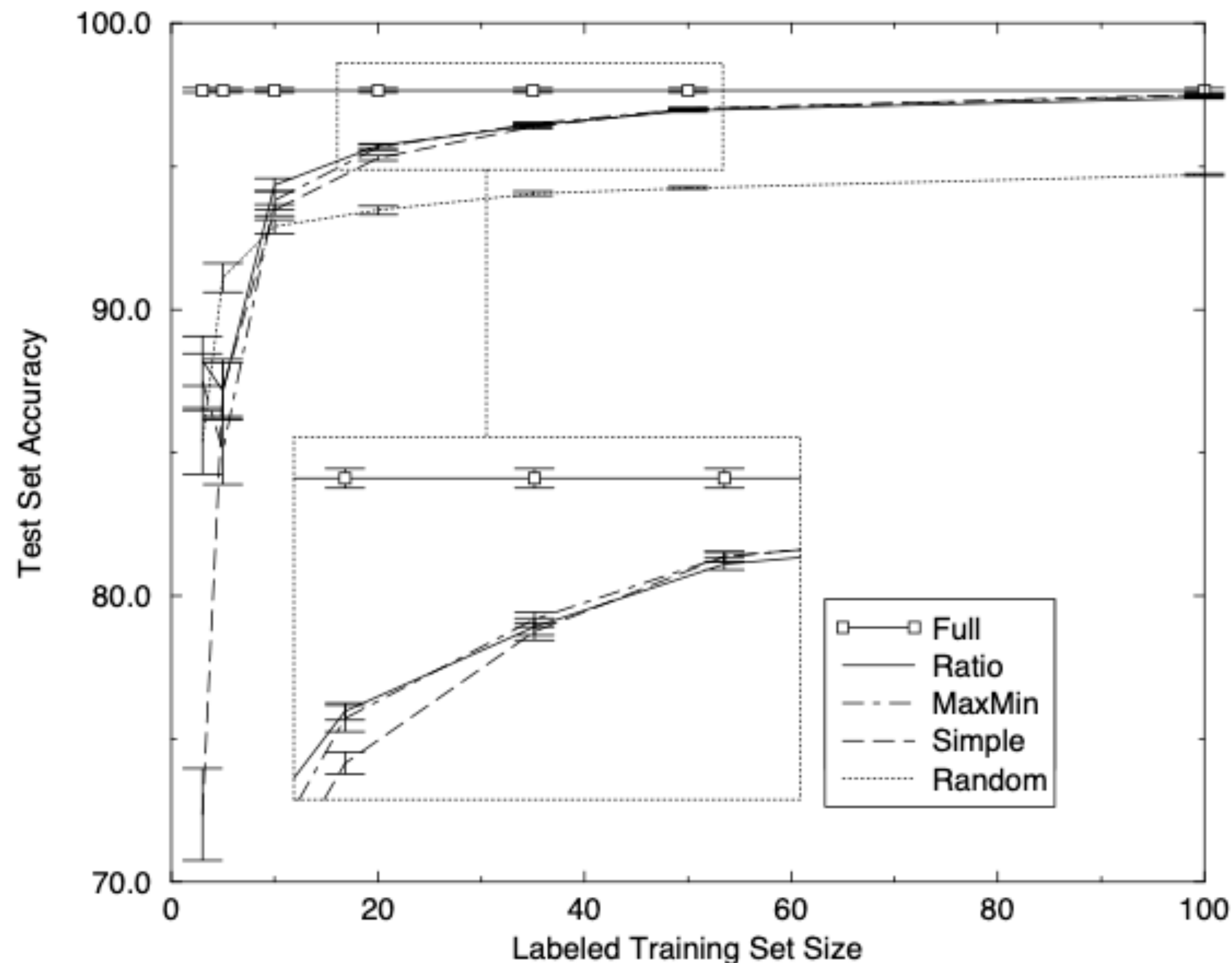
Version space is **set of hyperplanes**  
(or could be redefined through vectors  $W$ )



# Active learning with SVM version space



- Rapidly reduce version space
- Intuitively: choose queries that halve the version space
- Various approximations: is version space symmetric? Estimates of the size? etc.





Reducing the set of consistent hypotheses does not regard the evaluation metric

## An alternative to version space (the ML way)



We could also take a look at the machine learning **loss** and include points that would:

- **most reduce the expected error**
- **most change the current model**

## An alternative to version space (the ML way)



We could also take a look at the machine learning **loss** and include points that would:

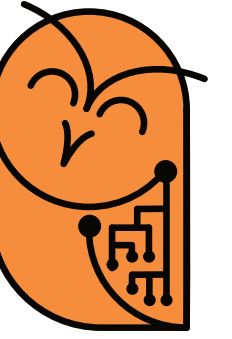
- **most reduce the expected error**
- **most change the current model**

*“First-order Markov active learning aims to select a query  $x^*$ , such that when the query is given label  $y^*$  and added to the training set, the learner trained on the resulting set  $D_+(x^*, y^*)$  has lower error than any other  $x$ ”*

*Roy & McCallum, “Toward Optimal Active Learning through Monte Carlo Estimation of Error Reduction”, ICML 2001)*

*(See also Cohn et al, “Active learning with statistical models”, JAIR 4, 1996)*

## The simplest (?) approach



Version spaces & expected error reduction are hard (& heavy to compute).

**Simple heuristics** are thus popular

## The simplest (?) approach



Version spaces & expected error reduction are hard (& heavy to compute).

**Simple heuristics** are thus popular, but have lots of caveats (tomorrow)

1. Create an initial classifier
2. While teacher is willing to label examples
  - (a) Apply the current classifier to each unlabeled example
  - (b) Find the  $b$  examples for which the classifier is least certain of class membership
  - (c) Have the teacher label the subsample of  $b$  examples
  - (d) Train a new classifier on all labeled examples

## Query by committee



We could maximize information gain between multiple models: **ensembles**

# Query by committee



We could maximize information gain between multiple models: **ensembles**

## Query by a committee of two

Repeat the following until  $n$  queries have been accepted

1. Draw an unlabeled input  $x \in X$  at random from  $\mathcal{D}$ .
2. Select two hypotheses  $h_1, h_2$  from the posterior distribution. In other words, pick two hypotheses that are consistent with the labeled examples seen so far.
3. If  $h_1(x) \neq h_2(x)$  then query the teacher for the label of  $x$ , and add it to the training set.



# Query by committee



We could maximize information gain between multiple models: **ensembles**

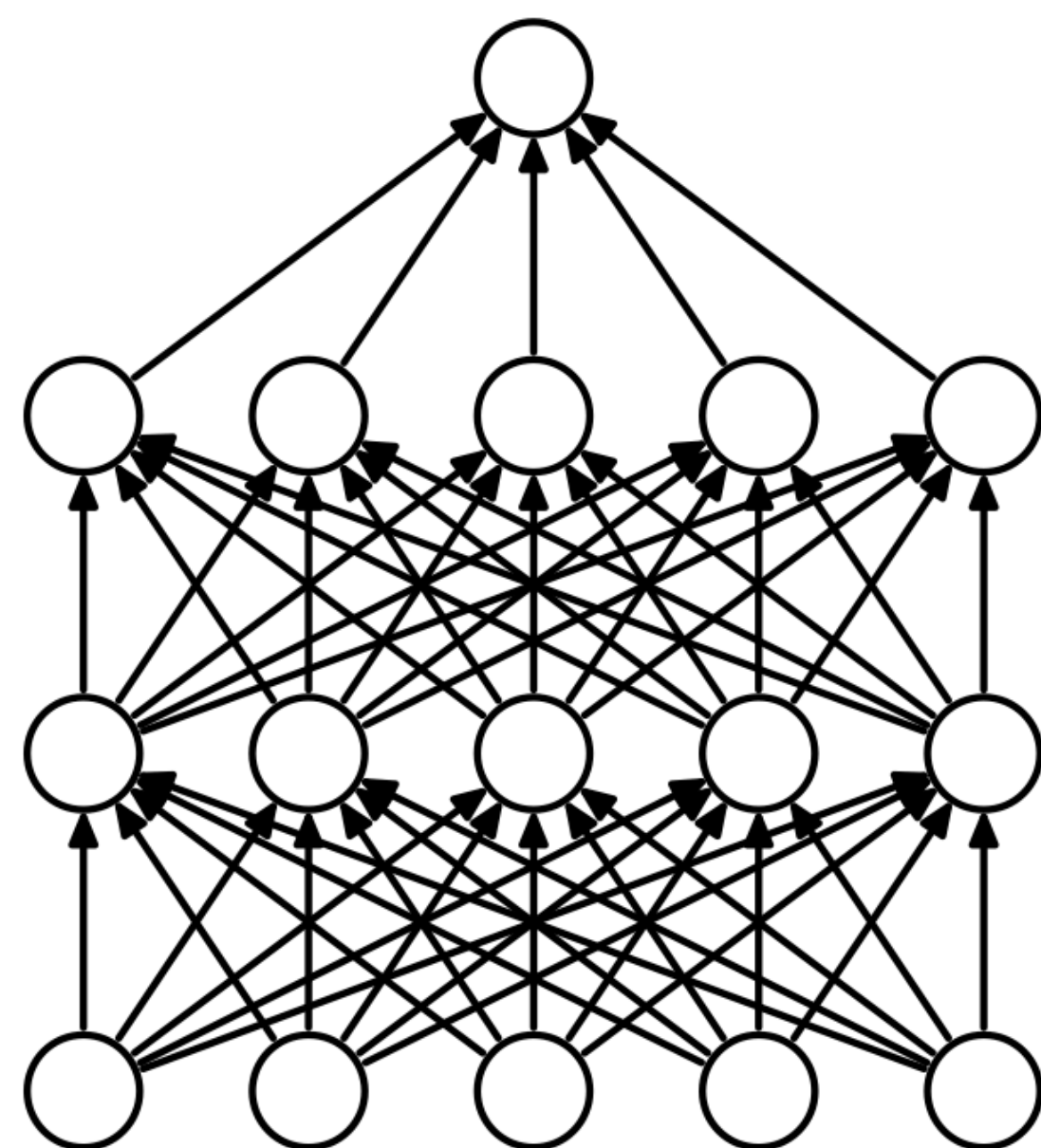
## Query by a committee of two

Repeat the following until  $n$  queries have been accepted

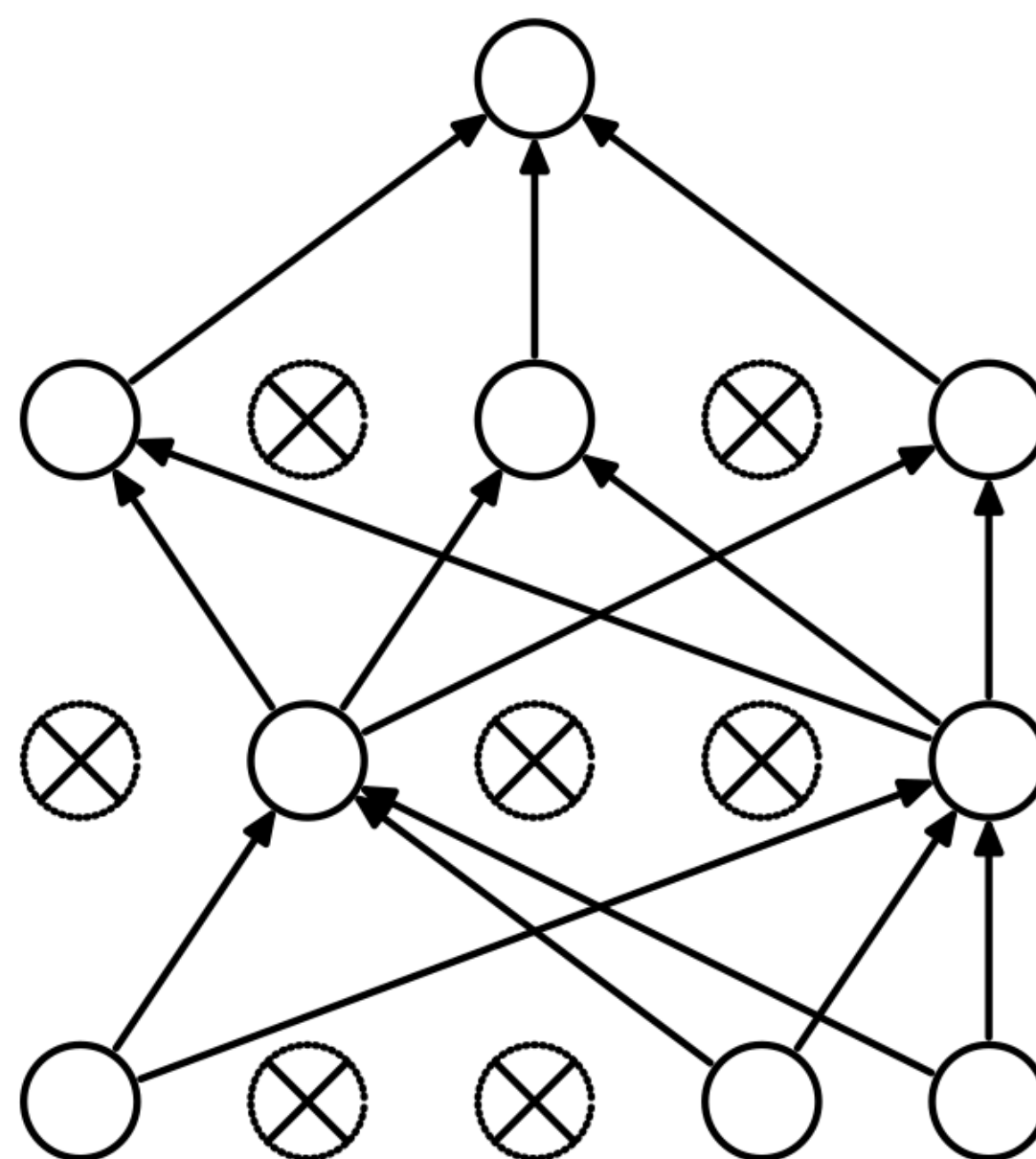
1. Draw an unlabeled input  $x \in X$  at random from  $\mathcal{D}$ .
2. Select two hypotheses  $h_1, h_2$  from the posterior distribution. In other words, pick two hypotheses that are consistent with the labeled examples seen so far.
3. If  $h_1(x) \neq h_2(x)$  then query the teacher for the label of  $x$ , and add it to the training set.

Could also be interpreted as reducing the version space across models

# Monte Carlo Dropout (Gal et al, ICML 2016)



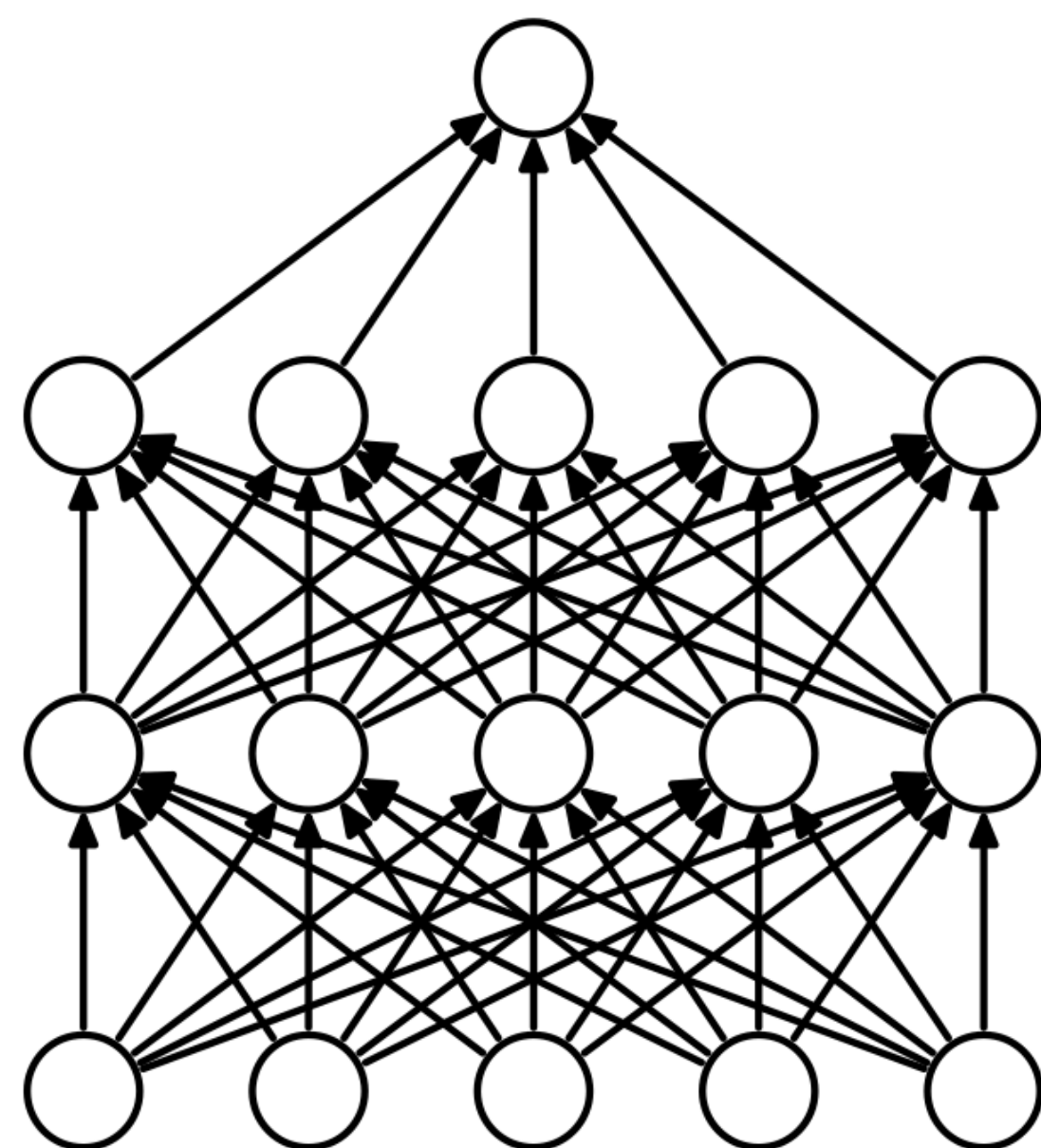
(a) Standard Neural Net



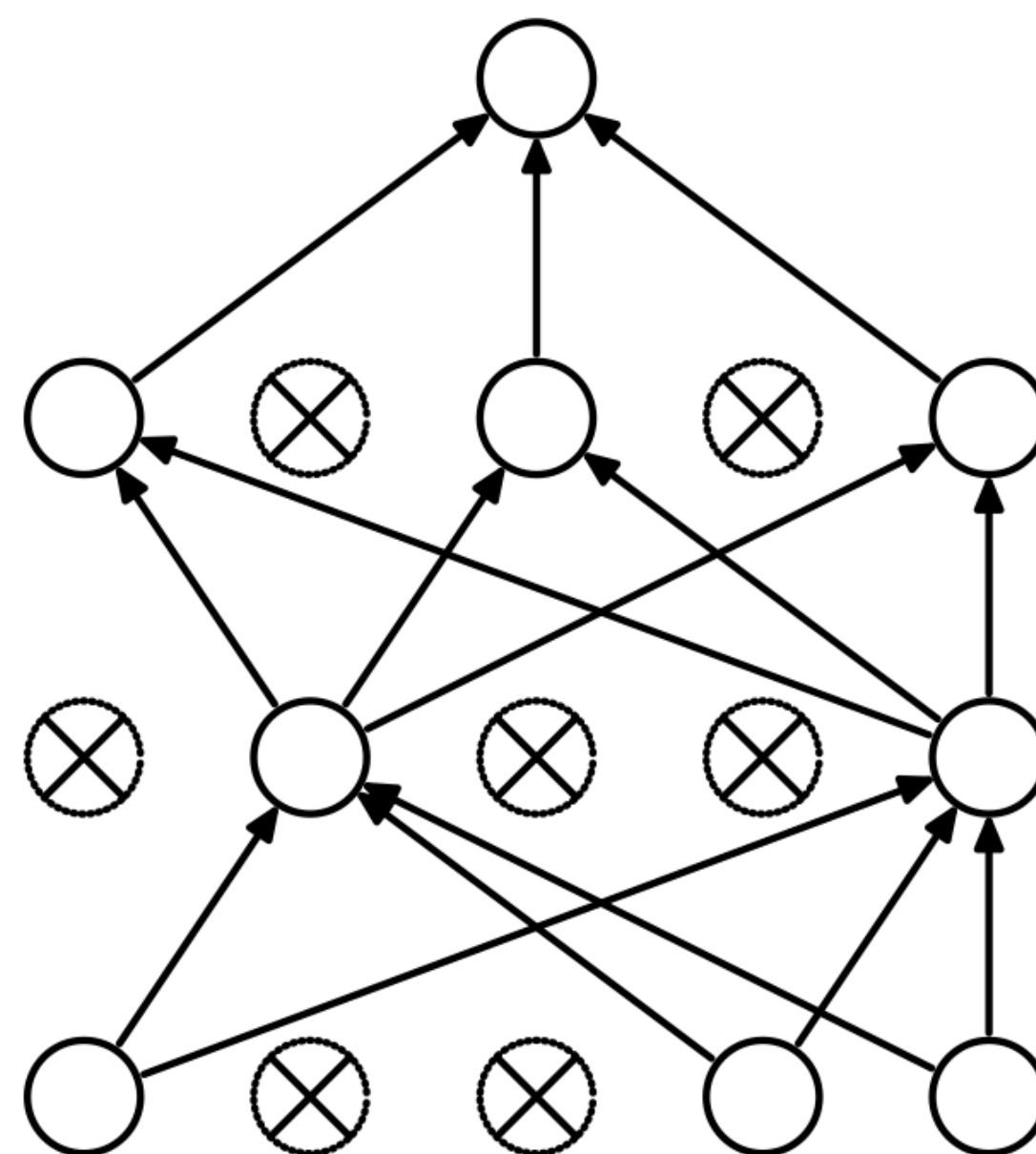
(b) After applying dropout.

- Make use of dropout: randomly turning off units in a model

# Monte Carlo Dropout (Gal et al, ICML 2016)



(a) Standard Neural Net



(b) After applying dropout.

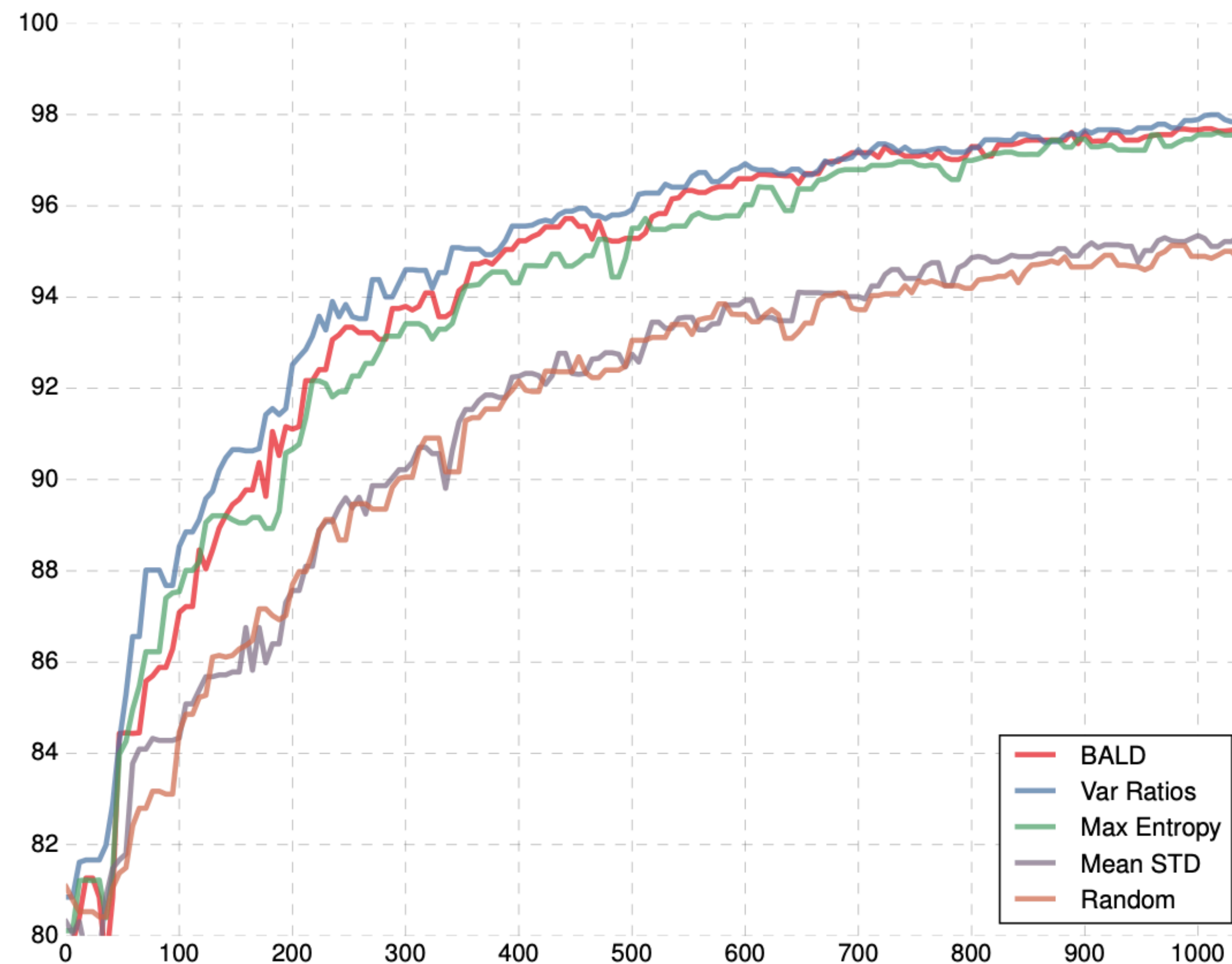
- Make use of dropout: randomly turning off units in a model
- Bayesian interpretation: Bernoulli distribution on the parameters
- Stochastic forward passes to get variation in predictions

# Monte Carlo Dropout (Gal et al, ICML 2016)



MCD could be useful as an approximation to using multiple model based ensembles

The acquisition function can still be entropy, standard deviation in output confidence etc.



# Random is hard to beat.

## Why aren't these approaches a lot better?

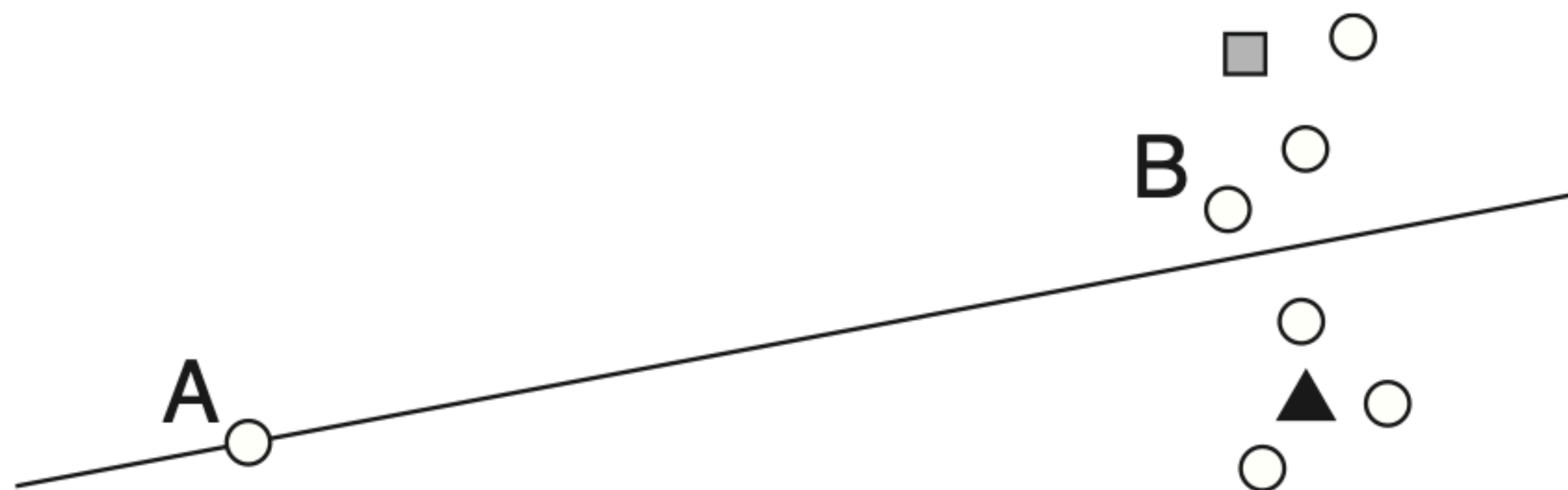


Figure 2: An illustration of when uncertainty sampling can be a poor strategy for classification. Shaded polygons represent labeled instances ( $\mathcal{L}$ ), and circles represent unlabeled instances ( $\mathcal{U}$ ). Since  $A$  is on the decision boundary, it will be queried as the most uncertain. However, querying  $B$  is likely to result in more information about the data as a whole.



There are more challenges with data that is “far away” (tomorrow). Let us first complete the picture

# Active learning perspectives



## Version space reduction

reduce the set/space of possible hypotheses  $h : \mathcal{X} \rightarrow \mathcal{Y}$  by removing the ones that are inconsistent with the data

## Uncertainty & heuristics

use the predictions, or maybe even better, uncertainty in the predictions for the queries

## Core sets & representation learning

Maximize distribution coverage instead of reducing the possible set of hypotheses

## Representations & core sets



**What if we allow to use & even train on the unlabelled pool?**

Assumption: a **“teacher” information source is allowed**, e.g. generative model

We wouldn't necessarily get a lot of advantage of generative models in active learning, unless **we also train on the (unlabelled pool)**

We could then also **make use of core sets**, as discussed for memory

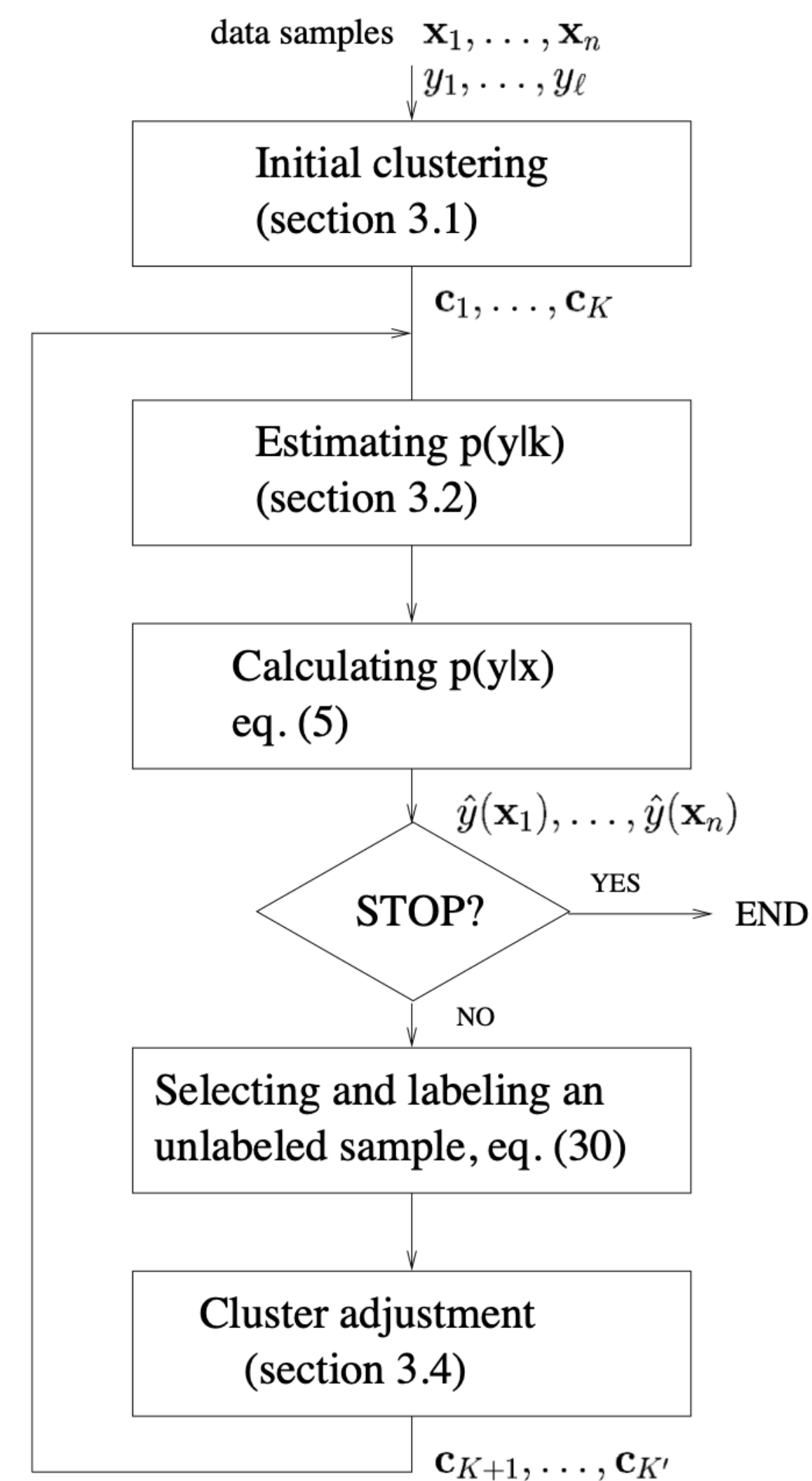


# Representations & core sets

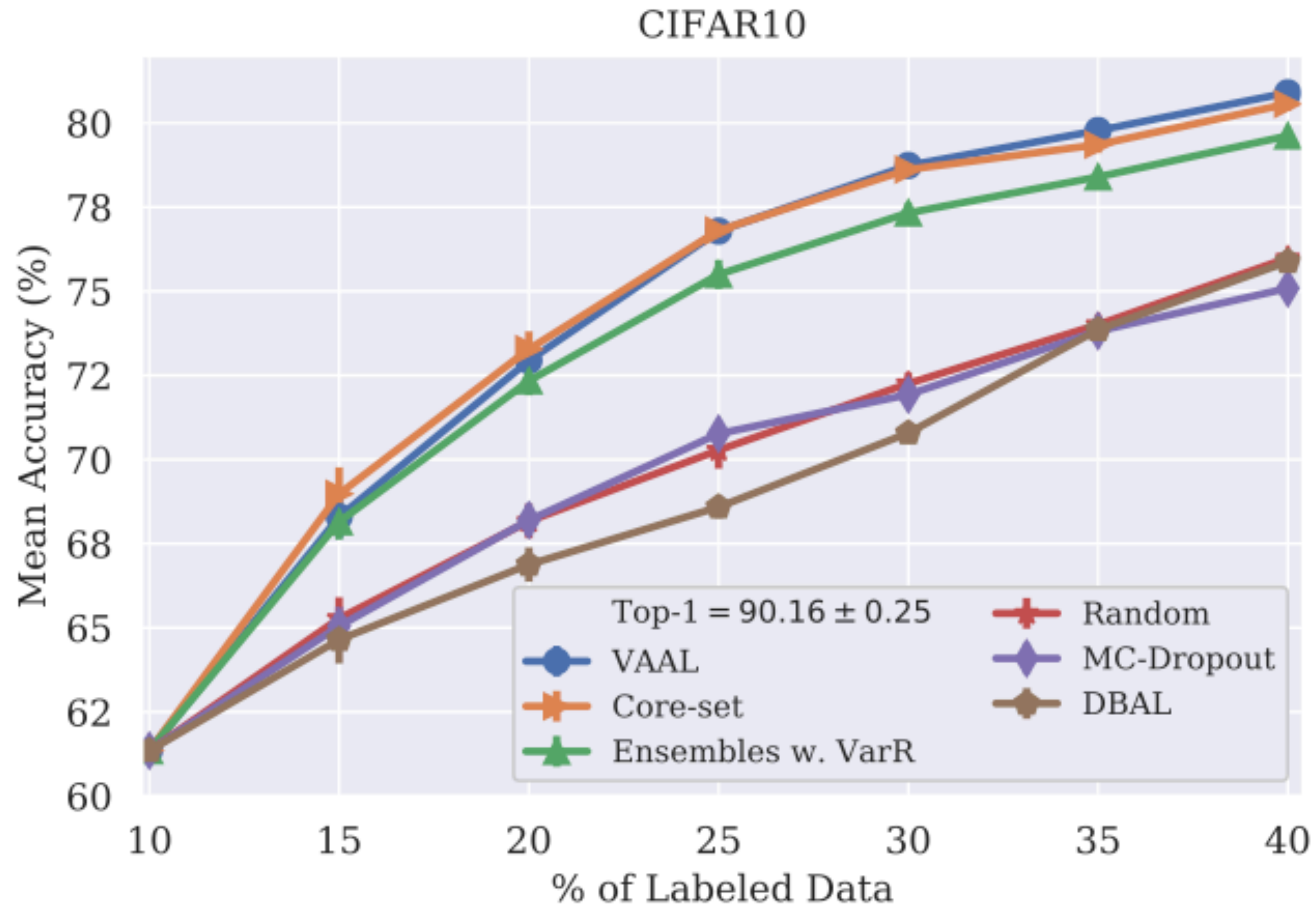


We could now try to:

- Pre-cluster our unlabelled data pool
- Compute core sets of the unlabelled data pool
- Learn a generative model & representations on the unlabelled data pool



# Representations & core sets





Intermediate summary: assumptions & trade-offs

## Intermediate summary: active learning perspectives



### **Version space reduction (*Hypotheses*)**

reduce the set/space of possible hypotheses  $h : \mathcal{X} \rightarrow \mathcal{Y}$  by removing the ones that are inconsistent with the data

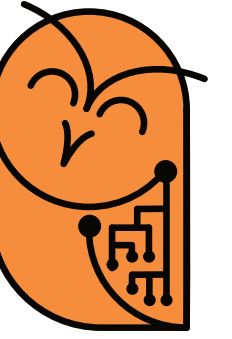
### **Uncertainty & heuristics (*Novelty*)**

use the predictions, or maybe even better, uncertainty in the predictions for the queries

### **Core sets & representation learning - accessing the entire pool (*Diversity*)**

maximize distribution coverage instead of reducing the possible set of hypotheses

# Intermediate summary: active learning perspectives



## Techniques

- Version space reduction

## & (some of) their assumptions

- Set of hypotheses is clear

# Intermediate summary: active learning perspectives



## Techniques

- Version space reduction
- Minimum confidence
- Maximum entropy

## & (some of) their assumptions

- Set of hypotheses is clear
- No overconfidence phenomenon and out-of-distribution/task data

# Intermediate summary: active learning perspectives



## Techniques

- Version space reduction
- Minimum confidence
- Maximum entropy
- Model “uncertainty” (output variability)
- Ensembles/query by committee

## & (some of) their assumptions

- Set of hypotheses is clear
- No overconfidence phenomenon and out-of-distribution/task data
- Accurate uncertainty everywhere
- Training of multiple models

# Intermediate summary: active learning perspectives



## Techniques

- Version space reduction
- Minimum confidence
- Maximum entropy
- Model “uncertainty” (output variability)
- Ensembles/query by committee
- Representation learning on the pool
- Core sets

## & (some of) their assumptions

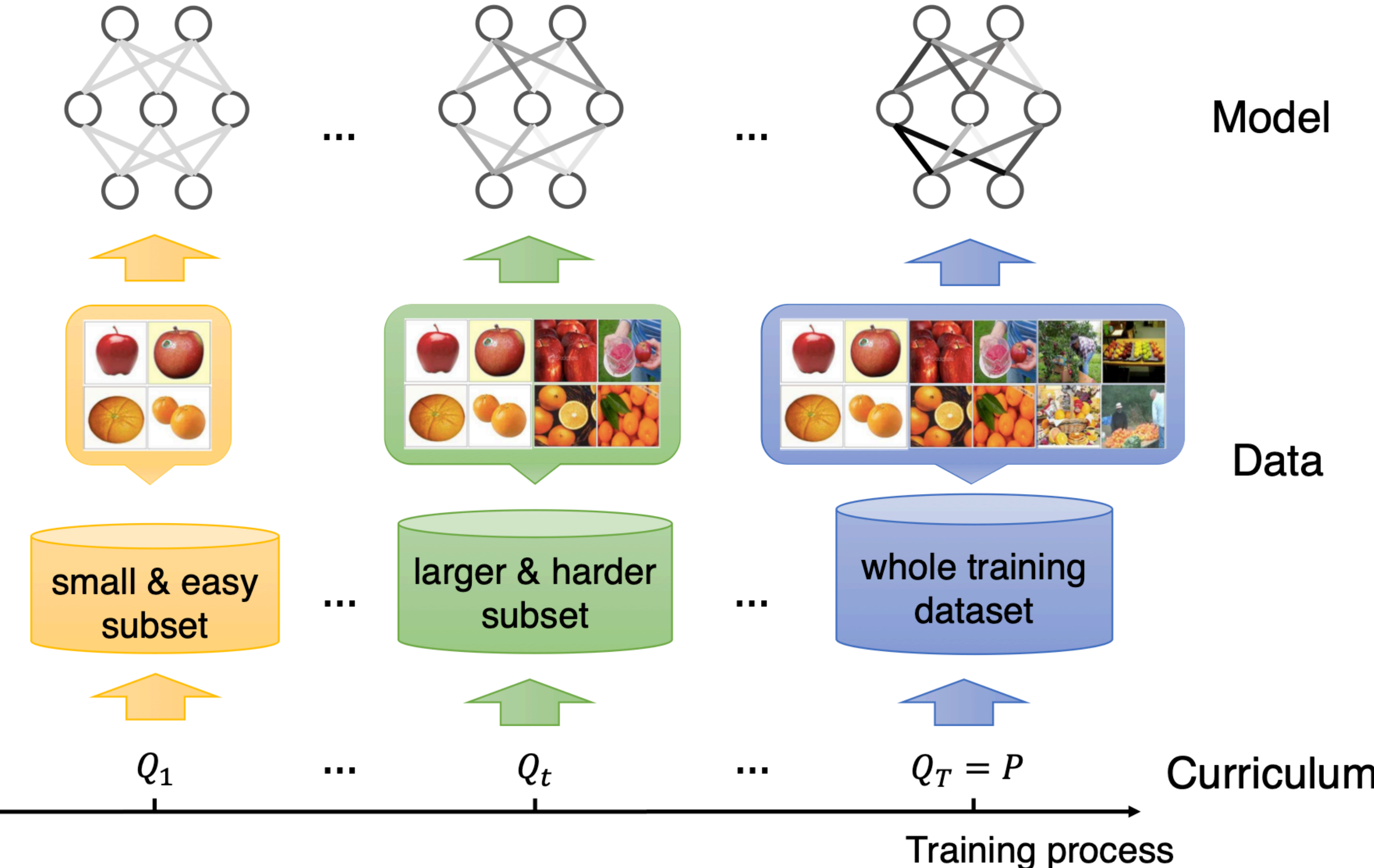
- Set of hypotheses is clear
- No overconfidence phenomenon and out-of-distribution/task data
- Accurate uncertainty everywhere
- Training of multiple models
- Upfront training of entire pool (no data stream)  
(access + computational expense)





There is another aspect to consider:  
the informativeness of data + difficulty to learn

# Curriculum Learning



what is “**easy**” & what is a “**harder**” subset/dataset?

And what is the difference to **informativeness**?

Let's start with an intuitive example: Ranking language model trained with vs without curriculum on Wikipedia



“Error” is log of the rank of the next word (within 20k-word vocabulary).

# Let's start with an intuitive example: Ranking language model trained with vs without curriculum on Wikipedia



“Error” is log of the rank of the next word (within 20k-word vocabulary).

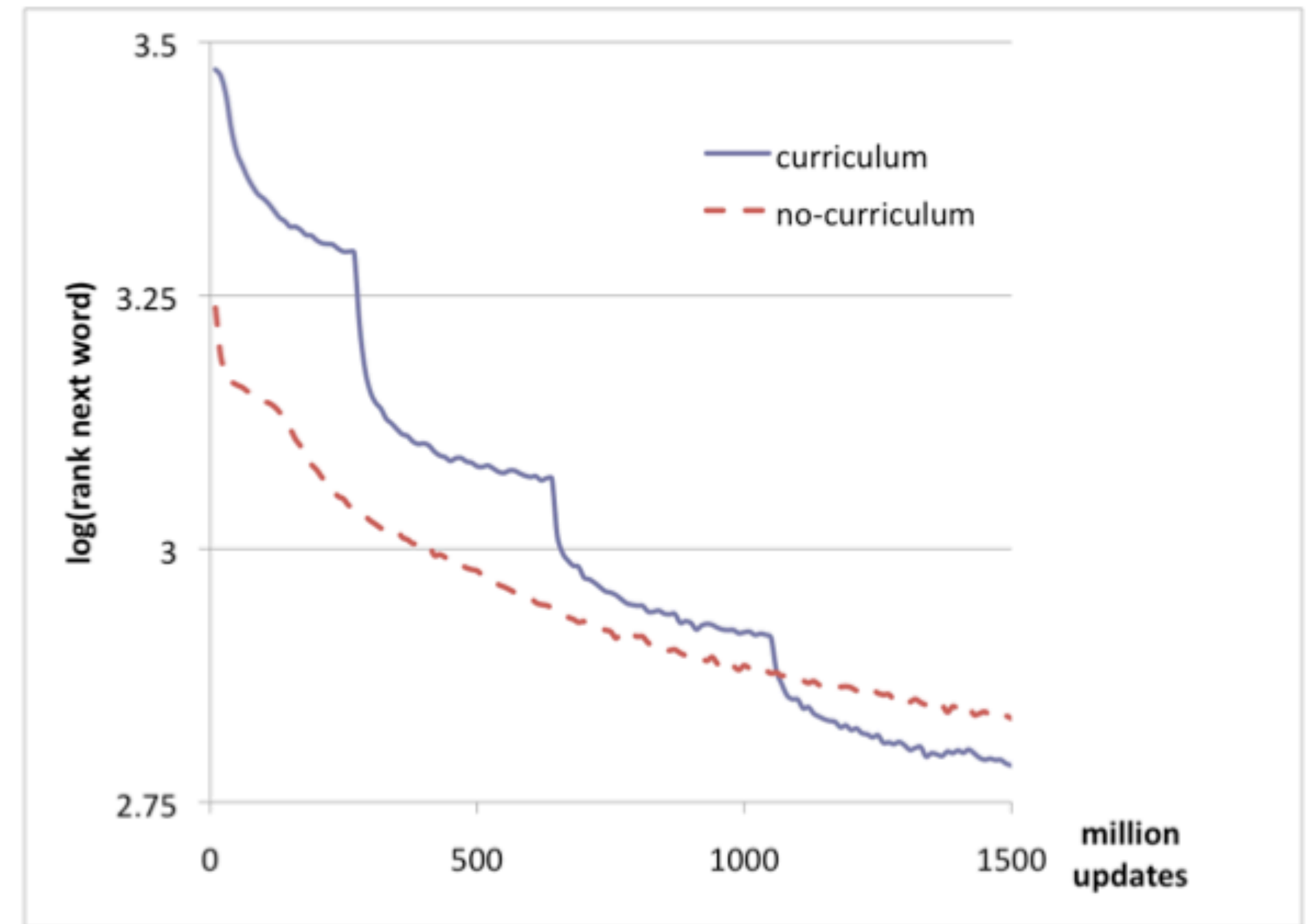
1. The curriculum-trained model skips examples with words outside of 5k most frequent words
2. Then skips examples outside 10k most frequent words and so on

# Let's start with an intuitive example: Ranking language model trained with vs without curriculum on Wikipedia



“Error” is log of the rank of the next word (within 20k-word vocabulary).

1. The curriculum-trained model skips examples with words outside of 5k most frequent words
2. Then skips examples outside 10k most frequent words and so on



Bengio et al, “Curriculum Learning”, ICML 2009

# Curriculum learning: the two key challenges



## **Scoring function (or difficulty measurer):**

Any function that provides us with an estimate of the difficulty of the instances in our dataset(s).

# Curriculum learning: the two key challenges



## Scoring function (or difficulty measurer):

Any function that provides us with an estimate of the difficulty of the instances in our dataset(s).

## Pacing function (or training scheduler):

(sometimes also called competence, as we'll see)

The function that tells us how to interleave samples into the training process over time.

# Curriculum Learning



**Definition 1: Original Curriculum Learning [6].** A curriculum is a sequence of training criteria over  $T$  training steps:  $\mathcal{C} = \langle Q_1, \dots, Q_t, \dots, Q_T \rangle$ . Each criterion  $Q_t$  is a reweighting of the target training distribution  $P(z)$ :

$$Q_t(z) \propto W_t(z)P(z) \quad \forall \text{example } z \in \text{training set } D, \quad (1)$$

such that the following three conditions are satisfied:

- 1) The entropy of distributions gradually increases, i.e.,  $H(Q_t) < H(Q_{t+1})$ .
- 2) The weight for any example increases, i.e.,  $W_t(z) \leq W_{t+1}(z) \quad \forall z \in D$ .
- 3)  $Q_T(z) = P(z)$ .



# Curriculum Learning



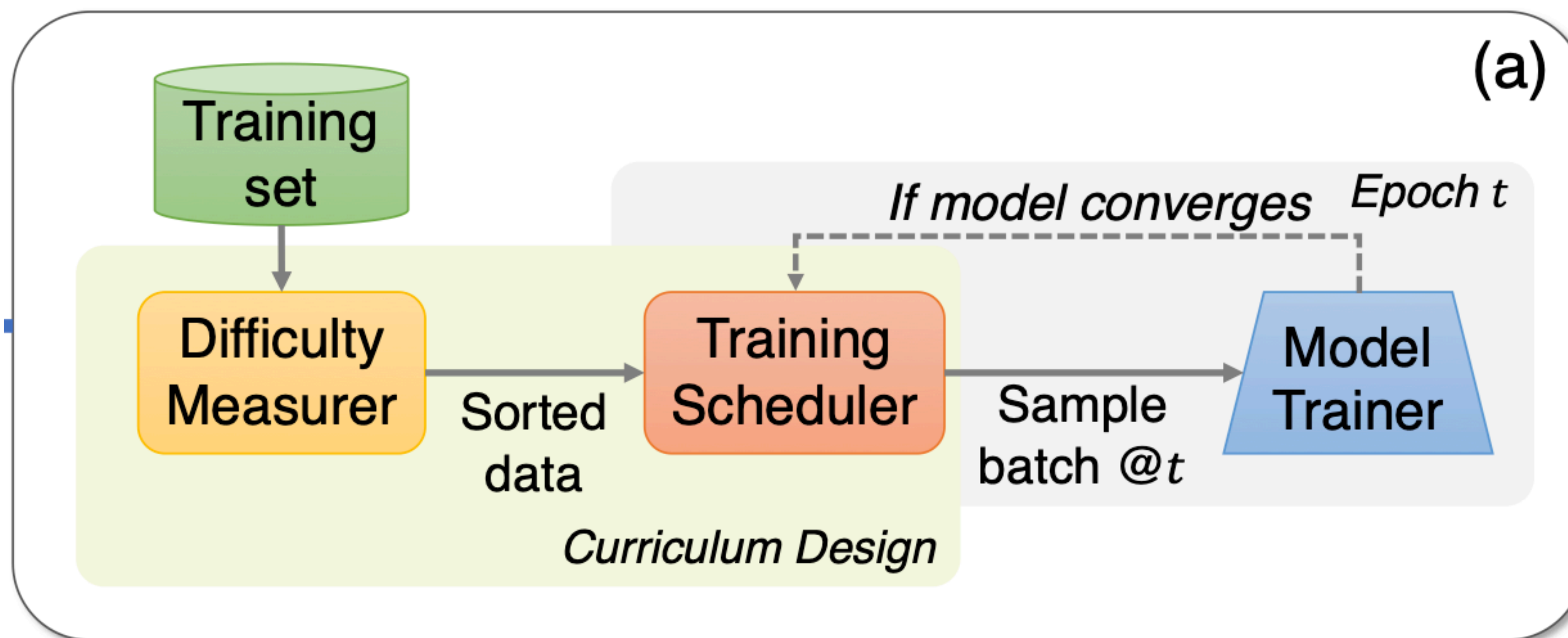
Curriculum learning: the more intuitive definition  
(with a little bit of a tautology)

**Definition 3: Generalized Curriculum Learning.** Discarding the definition of  $Q_t$  (Eq. 1) and its three conditions in Definition 1, a curriculum is a sequence of training criteria over  $T$  training steps. Each criterion  $Q_t$  includes the design for all the elements in training a machine learning model, e.g., data/tasks, model capacity, learning objective, etc. Curriculum learning is the strategy that trains a model with such a curriculum.

# Curriculum learning



Let's start by considering a **pre-defined curriculum**, inspired by learning from “textbook style” content





Can you think of ways to define “difficulty”?

# How to define difficulty: it is task & model specific



TABLE 2

Common types of predefined Difficulty Measurer. The “+” in  $\propto$ Easy means the higher the measured value, the easier the data example, and the “-” has the opposite meaning.

Difficulty Measurer*	Angle	Data Type	$\propto$ Easy
Sentence length [86], [107]	Complexity	Text	-
Number of objects [122]	Complexity	Images	-
# conj. [50], #phrases [113]	Complexity	Text	-
Parse tree depth [113]	Complexity	Text	-
Nesting of operations [131]	Complexity	Programs	-
Shape variability [6]	Diversity	Images	-
Word rarity [50], [86]	Diversity	Text	-
POS entropy [113]	Diversity	Text	-
Mahalanobis distance [14]	Diversity	Tabular	-
Cluster density [11], [31]	Noise	Images	+
Data source [10]	Noise	Images	/
SNR / SND [7], [89]	Noise	Audio	-
Grammaticality [66]	Domain	Text	+
Prototypicality [113]	Domain	Text	+
Medical based [44]	Domain	X-ray film	/
Retrieval based [18], [82]	Domain	Retrieval	/
Intensity [30] / Severity [111]	Intensity	Images	+
Image difficulty score [106], [114]	Annotation	Images	-
Norm of word vector [68]	Multiple	Text	-

# How to define difficulty: it is task & model specific



We have already seen that specific tasks allow for specific definitions of difficulty. Example: natural language translation (sentence length)

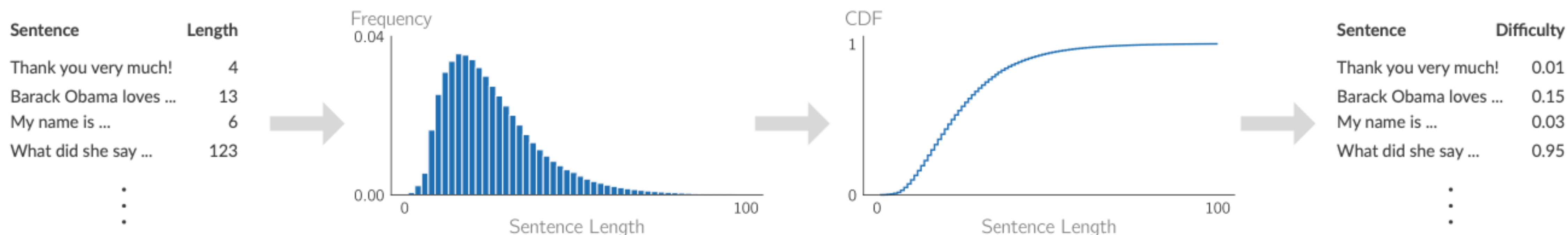


Figure 2: Example visualization of the preprocessing sequence used in the proposed algorithm. The histogram shown is that of sentence lengths from the WMT-16  $E_n \rightarrow D_e$  dataset used in our experiments. Here sentence lengths represent an example difficulty scoring function,  $d$ . “CDF” stands for the empirical “cumulative density function” obtained from the histogram on the left plot.

# How to define difficulty: it is task & model specific



Another example: image segmentation (entropy/clutter)



Figure 1. Images with difficulty scores predicted by our system in increasing order of their difficulty.

# How to define difficulty: it is task & model specific



There are various dimensions to difficulty, not just (basic) data statistics.  
Especially if we think about factors that relate to what humans find difficult

## Compositional factors:

Size



"A sail boat on the ocean."

Location



"Two men standing on beach."

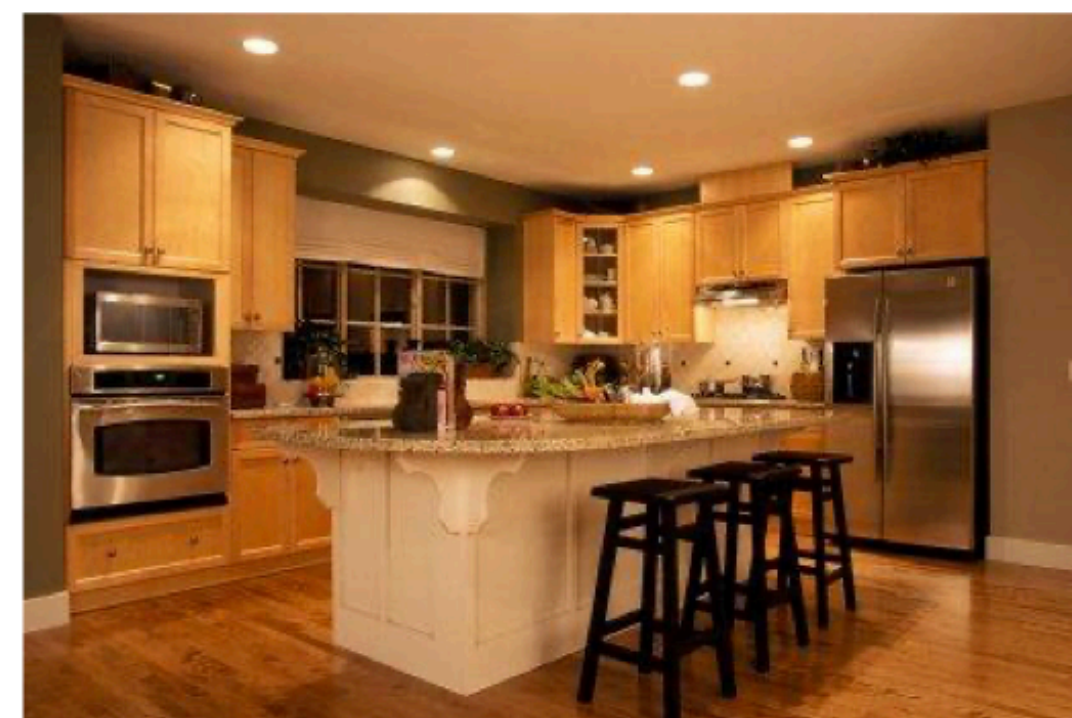
## Semantic factors:

Object Type



"Girl in the street"

Scene Type & Depiction Strength



"kitchen in house"

## Context factors:

Unusual object-scene Pair



"A tree in water and a boy with a beard"

## What is difficult for ML models?



But what is difficult for ML models & is this related to human perception?

Example: human response time

**Collecting response times.** We collected ground-truth difficulty annotations by human evaluators using the following protocol: (i) we ask each annotator a question of the type “Is there an *{object class}* in the next image?”, where *{object class}* is one of the 20 classes included in the PASCAL VOC 2012; (ii) we show the image to the annotator; (iii) we record the time spent by the annotator to answer the question by “Yes” or “No”. Finally, we use this response time to estimate the visual search difficulty.



## What is difficult for ML models?



Average human ranks about 80% image pairs in the same order as given by the mean response time of all annotators -> compared to Pascal "difficulty"

**Collecting response times.** We collected ground-truth difficulty annotations by human evaluators using the following protocol: (i) we ask each annotator a question of the type “Is there an *{object class}* in the next image?”, where *{object class}* is one of the 20 classes included in the PASCAL VOC 2012; (ii) we show the image to the annotator; (iii) we record the time spent by the annotator to answer the question by “Yes” or “No”. Finally, we use this response time to estimate the visual search difficulty.

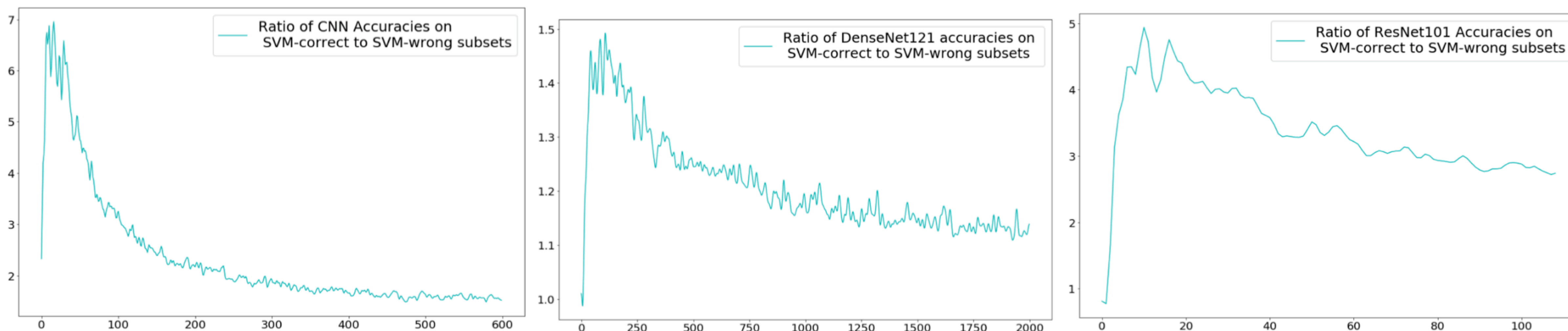
	Image property	Kendall $\tau$
(i)	number of objects	0.32
(ii)	mean area covered by objects	-0.28
(iii)	non-centeredness	0.29
(iv)	number of different classes	0.33
(v)	number of truncated objects	0.22
(vi)	number of occluded objects	0.26
(vii)	number of difficult objects	0.20

# What is difficult for ML models?



Example: shallow embeddable examples seem to be learned first

A deep network in comparison to a SVM (random forest also in the paper)



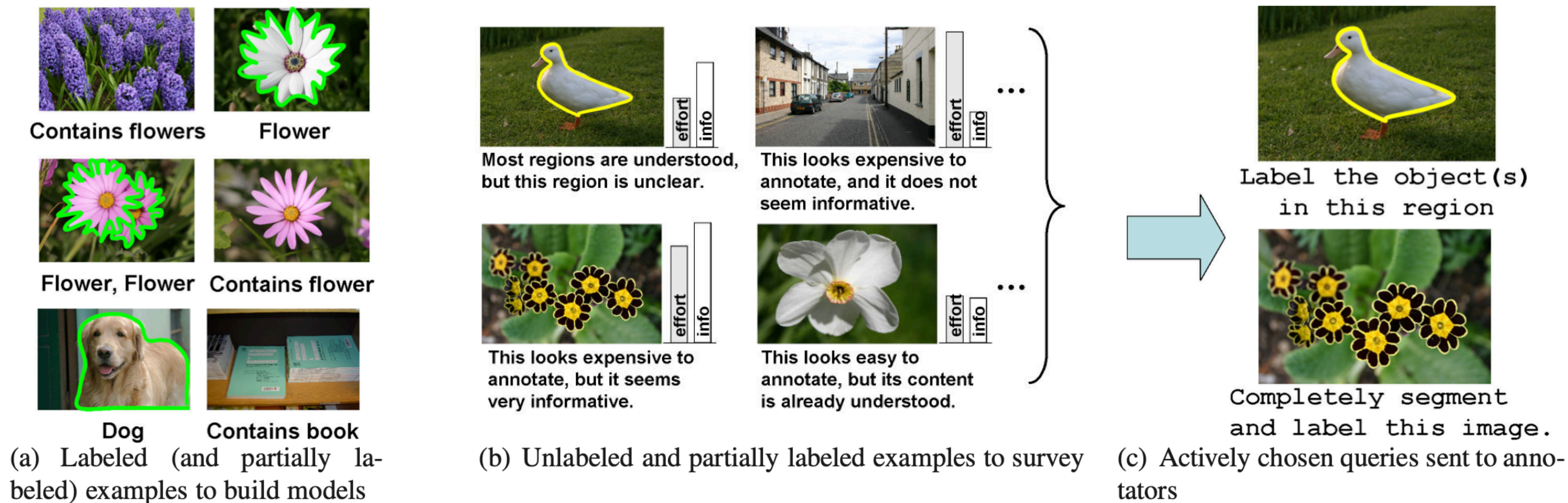
Ratio of Accuracies  $R^i$  plotted against  $i$  with  $\mathcal{M}$  being a Support Vector Machine.

# Difficulty beyond (curriculum) learning



Assessing difficulty is interesting beyond curriculum learning

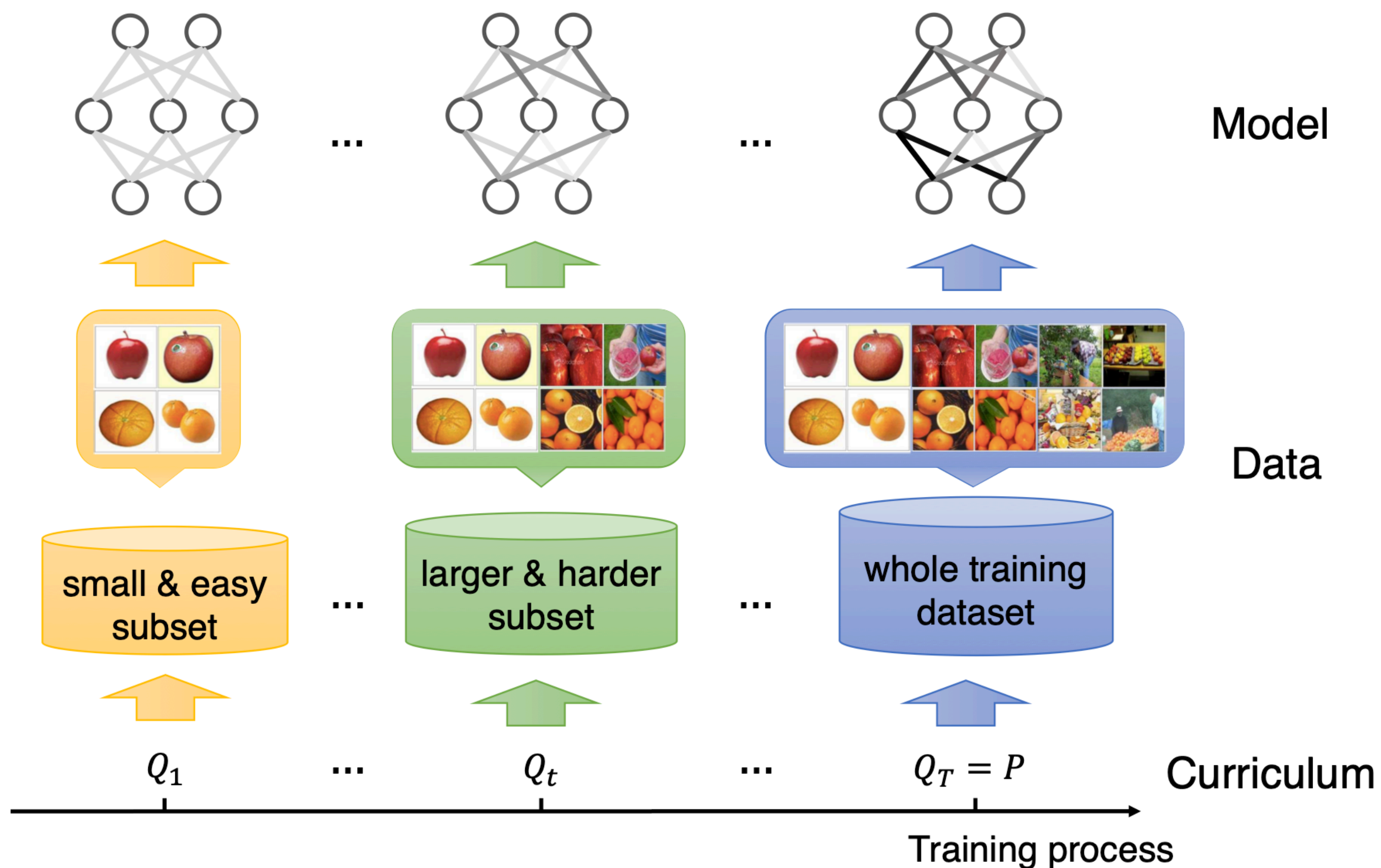
Example: estimating the difficulty with respect to annotation cost





## Pacing: how to schedule the training

# Scheduling training



If we want to define the curriculum up-front, according to prior knowledge, then:

when do we introduce more difficult examples?

# Pacing functions



Various options & heuristics are conceivable

---

## Algorithm 1 One-Pass Curriculum

---

```
1: procedure OP-CURRICULUM( $M, \mathcal{D}, \mathcal{C}$ )
2:    $\mathcal{D}' = \text{sort}(\mathcal{D}, \mathcal{C})$ 
3:    $\{\mathcal{D}^1, \mathcal{D}^2, \dots, \mathcal{D}^k\} = \mathcal{D}'$  where  $\mathcal{C}(d_a) < \mathcal{C}(d_b)$   $d_a \in \mathcal{D}^i, d_b \in \mathcal{D}^j, \forall i < j$ 
4:   for  $s = 1 \dots k$  do
5:     while not converged for  $p$  epochs do
6:       train( $M, \mathcal{D}^s$ )
7:     end while
8:   end for
9: end procedure
```

---

Algorithm from Cirik et al, "Visualizing and understanding curriculum learning for long short-term memory networks", arXiv, 2016

Based on the procedure described in Bengio et al, "Curriculum Learning", ICML 2009

# Pacing functions



Various options & heuristics are conceivable

---

## Algorithm 1 One-Pass Curriculum

---

```
1: procedure OP-CURRICULUM( $M, \mathcal{D}, \mathcal{C}$ )
2:    $\mathcal{D}' = \text{sort}(\mathcal{D}, \mathcal{C})$ 
3:    $\{\mathcal{D}^1, \mathcal{D}^2, \dots, \mathcal{D}^k\} = \mathcal{D}'$  where  $\mathcal{C}(d_a) < \mathcal{C}(d_b)$   $d_a \in \mathcal{D}^i, d_b \in \mathcal{D}^j, \forall i < j$ 
4:   for  $s = 1 \dots k$  do
5:     while not converged for  $p$  epochs do
6:       train( $M, \mathcal{D}^s$ )
7:     end while
8:   end for
9: end procedure
```

---

---

## Algorithm 2 Baby Steps Curriculum

---

```
1: procedure BS-CURRICULUM( $M, \mathcal{D}, \mathcal{C}$ )
2:    $\mathcal{D}' = \text{sort}(\mathcal{D}, \mathcal{C})$ 
3:    $\{\mathcal{D}^1, \mathcal{D}^2, \dots, \mathcal{D}^k\} = \mathcal{D}'$  where  $\mathcal{C}(d_a) < \mathcal{C}(d_b)$   $d_a \in \mathcal{D}^i, d_b \in \mathcal{D}^j, \forall i < j$ 
4:    $\mathcal{D}^{\text{train}} = \emptyset$ 
5:   for  $s = 1 \dots k$  do
6:      $\mathcal{D}^{\text{train}} = \mathcal{D}^{\text{train}} \cup \mathcal{D}^s$ 
7:     while not converged for  $p$  epochs do
8:       train( $M, \mathcal{D}^{\text{train}}$ )
9:     end while
10:  end for
11: end procedure
```

---

Algorithm from Cirik et al, "Visualizing and understanding curriculum learning for long short-term memory networks", arXiv, 2016

Based on the procedure described in Bengio et al, "Curriculum Learning", ICML 2009

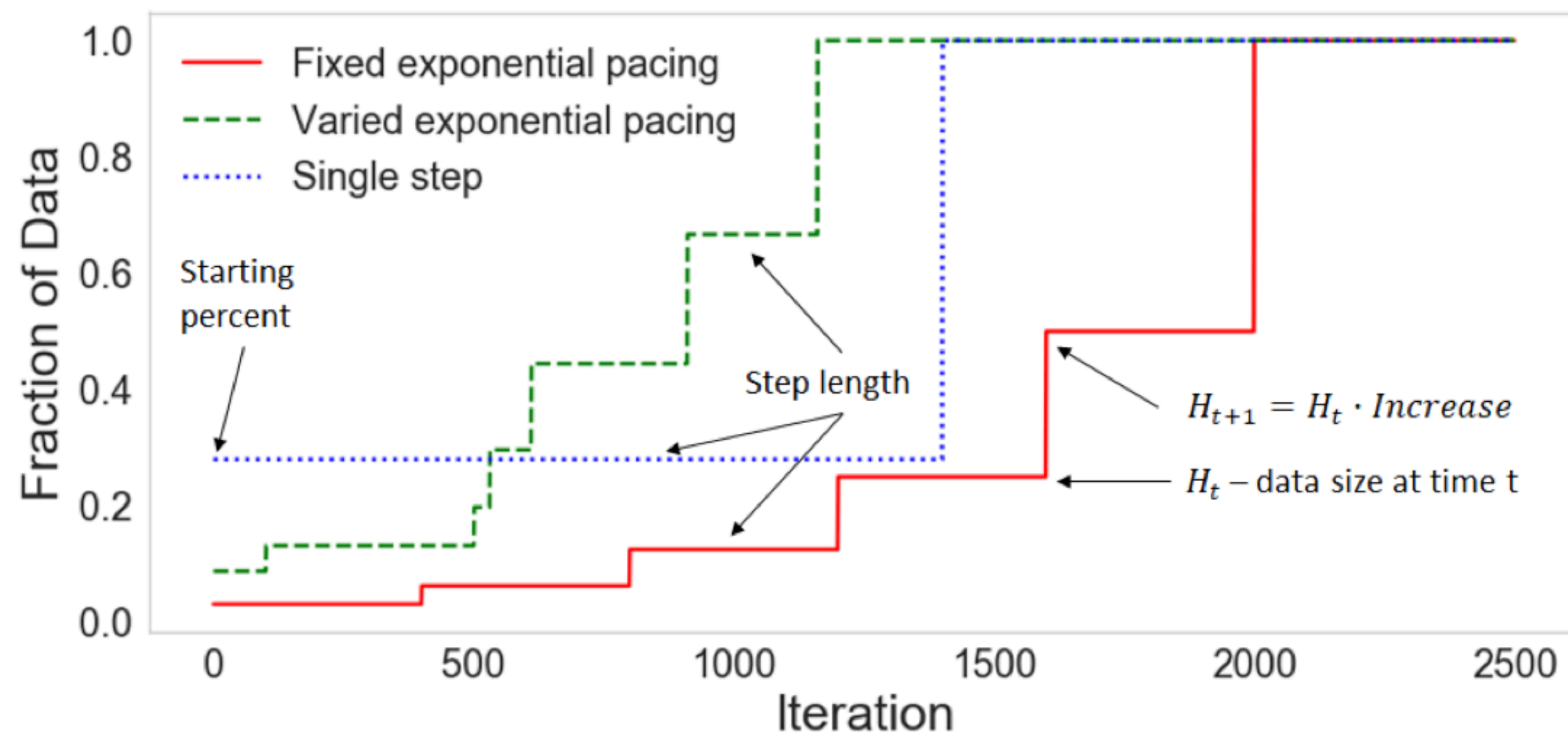
Algorithm from Cirik et al, "Visualizing and understanding curriculum learning for long short-term memory networks", arXiv, 2016

Based on the procedure described in Spitzkovsky et al, "From baby steps to leapfrogs: how less is more in unsupervised dependency parsing", NAACL-HLT, 2010

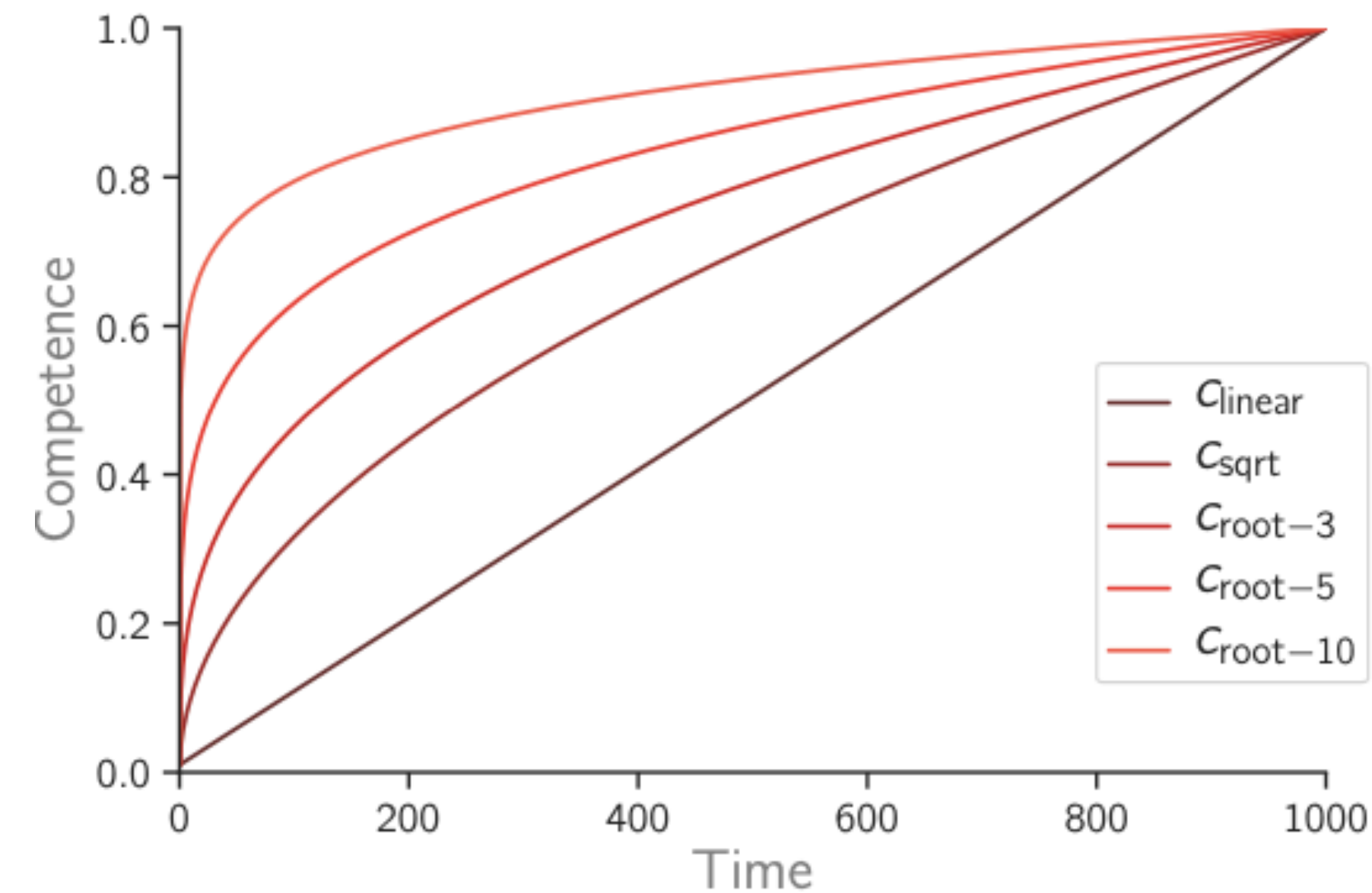
# Pacing functions



Various options & heuristics are conceivable



Hacohen & Weinshall, "On the power of curriculum learning in deep networks", ICML 2019



Platanios et al, "Competence based curriculum learning for neural machine translation", NAACL-HLT 2019

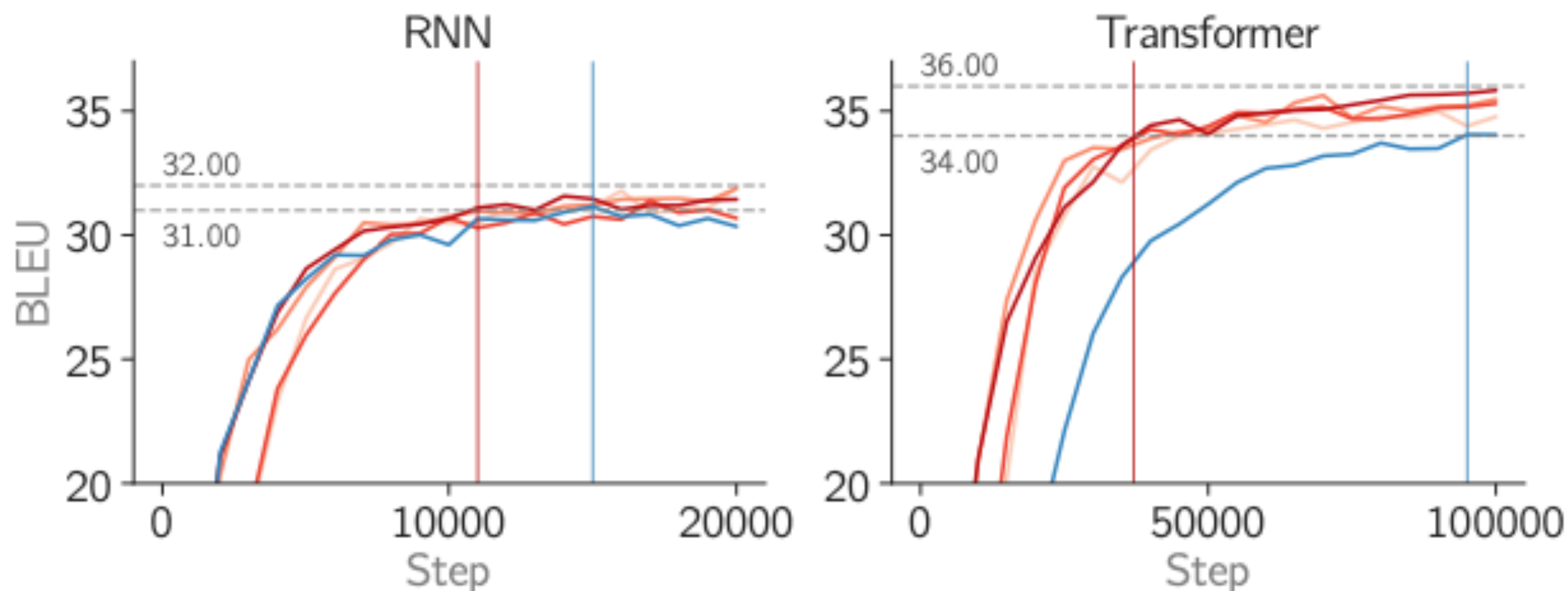


# Pacing functions



It's not straightforward to choose, especially due to model/task dependency

## IWSLT16 : Fr → En



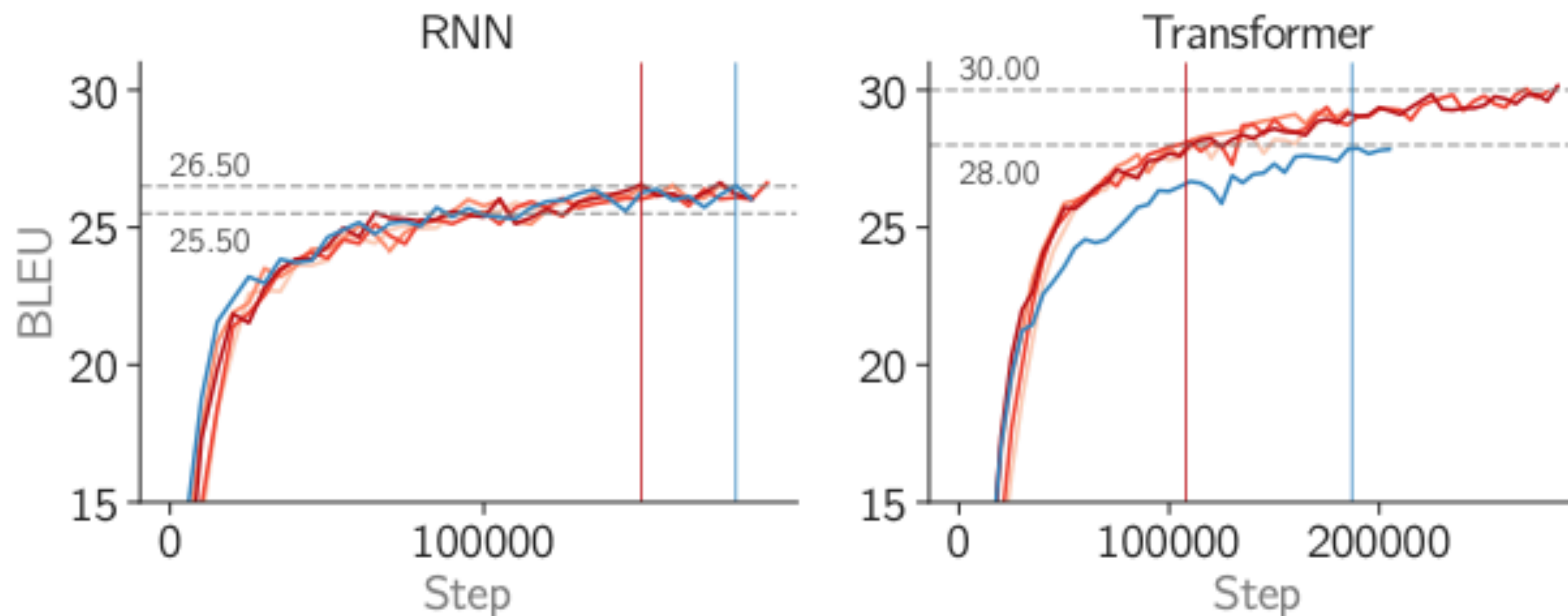
— Plain — SL Linear — SL Sqrt — SR Linear — SR Sqrt

# Pacing functions



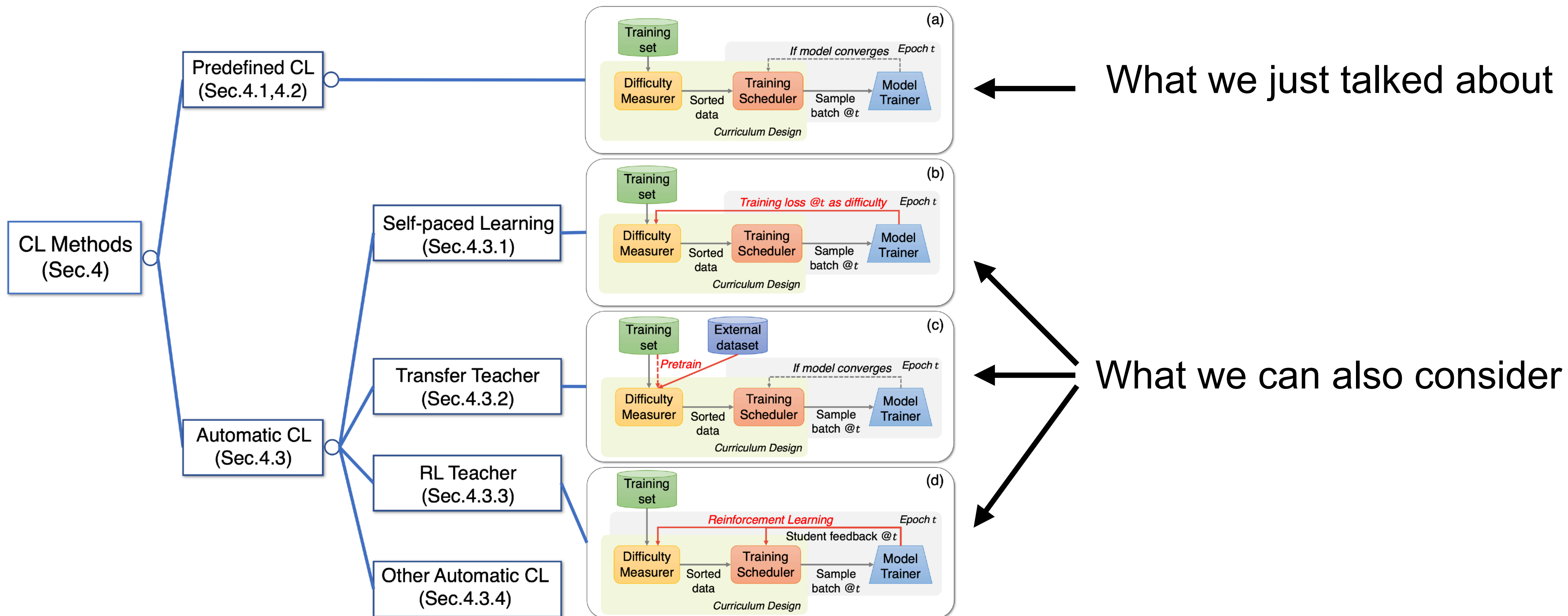
It's not straightforward to choose, especially due to model/task dependency

## WMT16 : En → De



— Plain — SL Linear — SL Sqrt — SR Linear — SR Sqrt

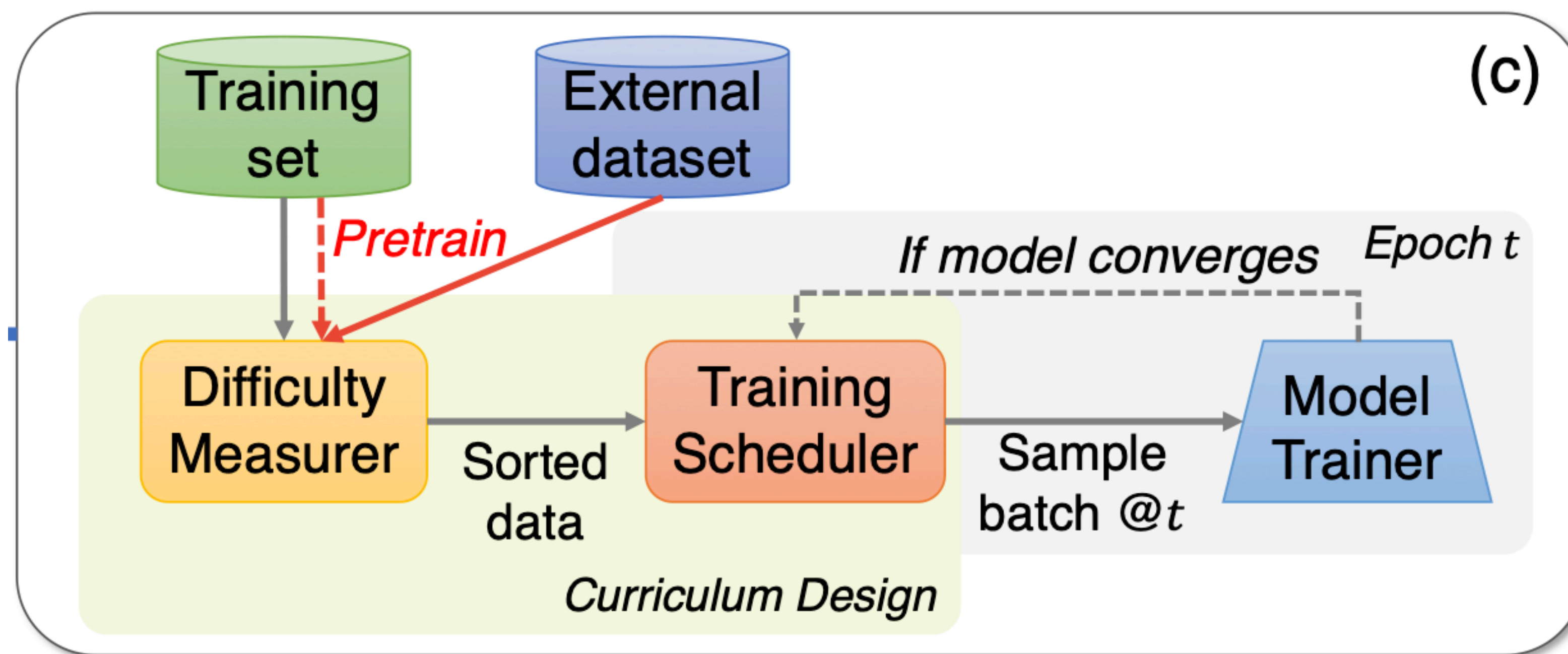
# Beyond pre-defined curricula



# Transfer-teacher curricula



Instead of defining the curriculum, we could use a **pre-trained teacher** model (based on a different related dataset) **based difficulty** measure



# Transfer-teacher curricula



Instead of defining the curriculum, we could use a **pre-trained teacher** model (based on a different related dataset) **based difficulty** measure

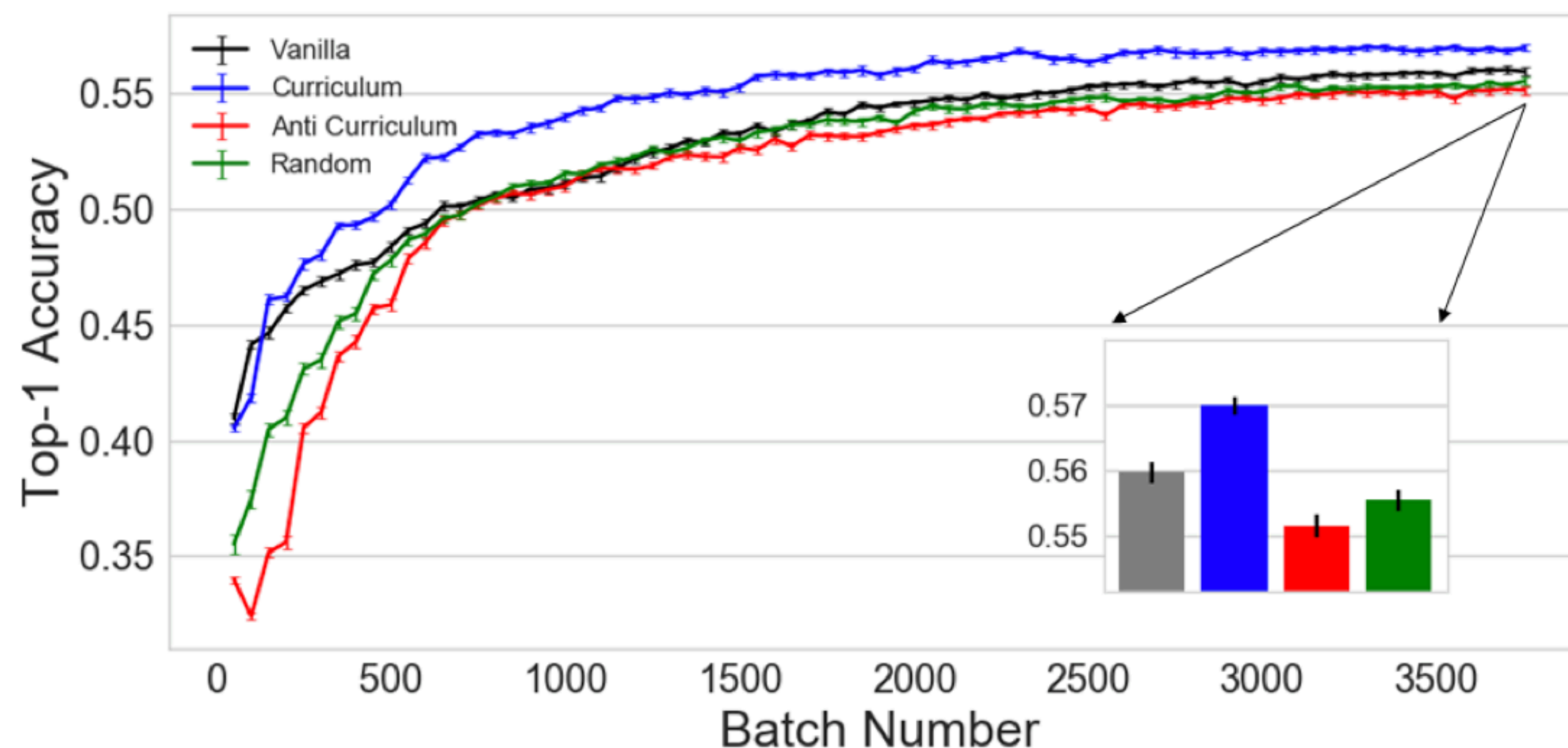
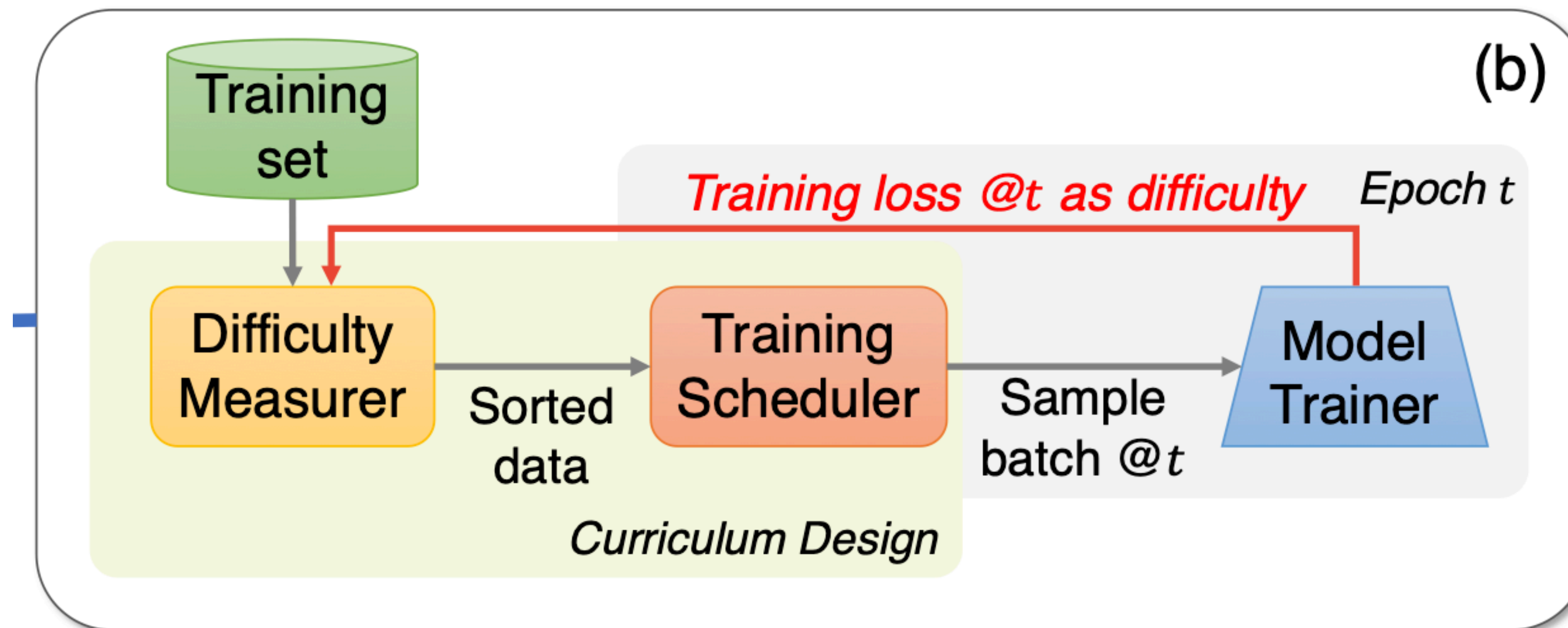


Figure 2. Results in **case 1**, with Inception-based *transfer scoring* function and *fixed exponential pacing* function.

# From pre-defined to self-paced



Using a teacher is a form of **pre-defined curriculum**, what if we want to have an **adaptive measure of difficulty**, based on our current model?  
Moving away from a pre-defined curriculum towards model “**competence**”

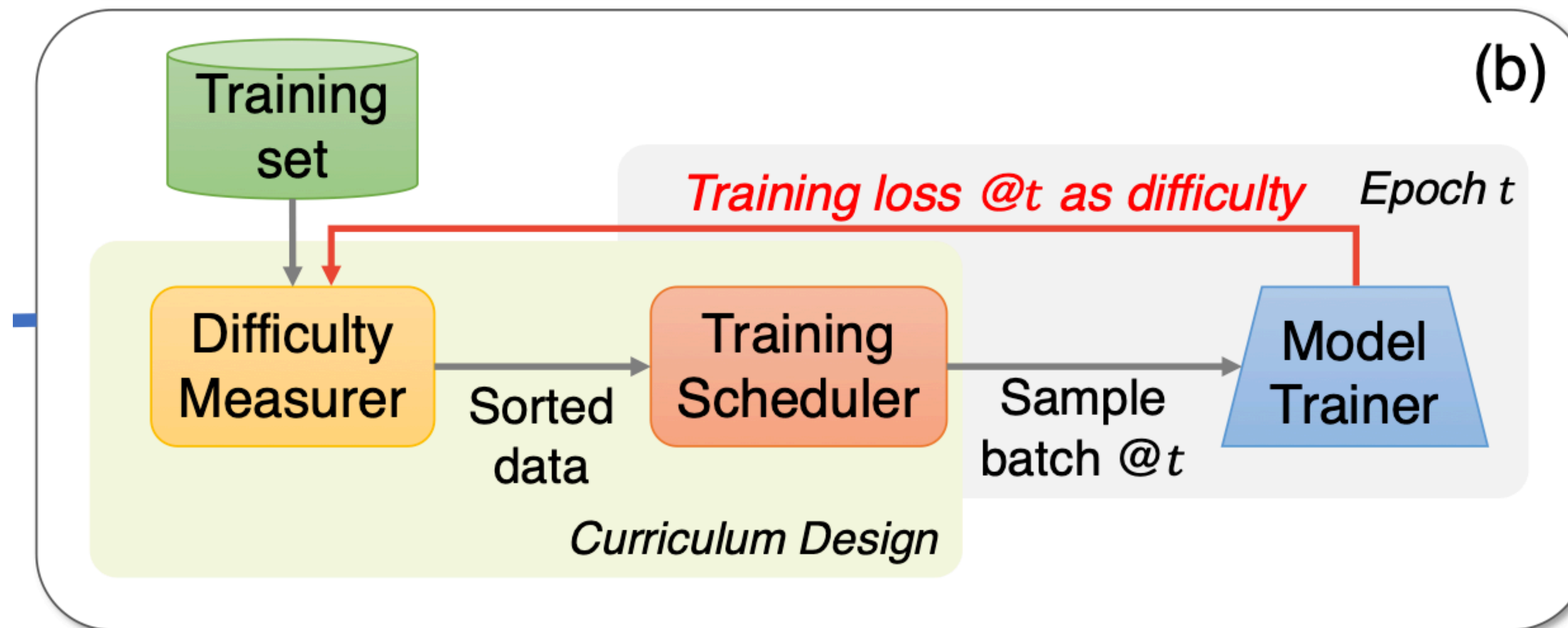


# From pre-defined to self-paced



Often this is called **self-paced learning**

Now rely on a model's **current hypothesis** at each point in time to **assign difficulty** to the training data, rather than ranking according to the target



# Self-paced & self-taught



## Self-paced learning:

Measure the difficulty of an instance according to current loss/predictions etc. (related to the ideas in *active learning*)

## Self-taught learning:

Train a model fully, measure each instance according to final model, assign difficulty score and start over with curriculum -> repeat (related to the ideas in *boosting*)



# Self-paced & self-taught

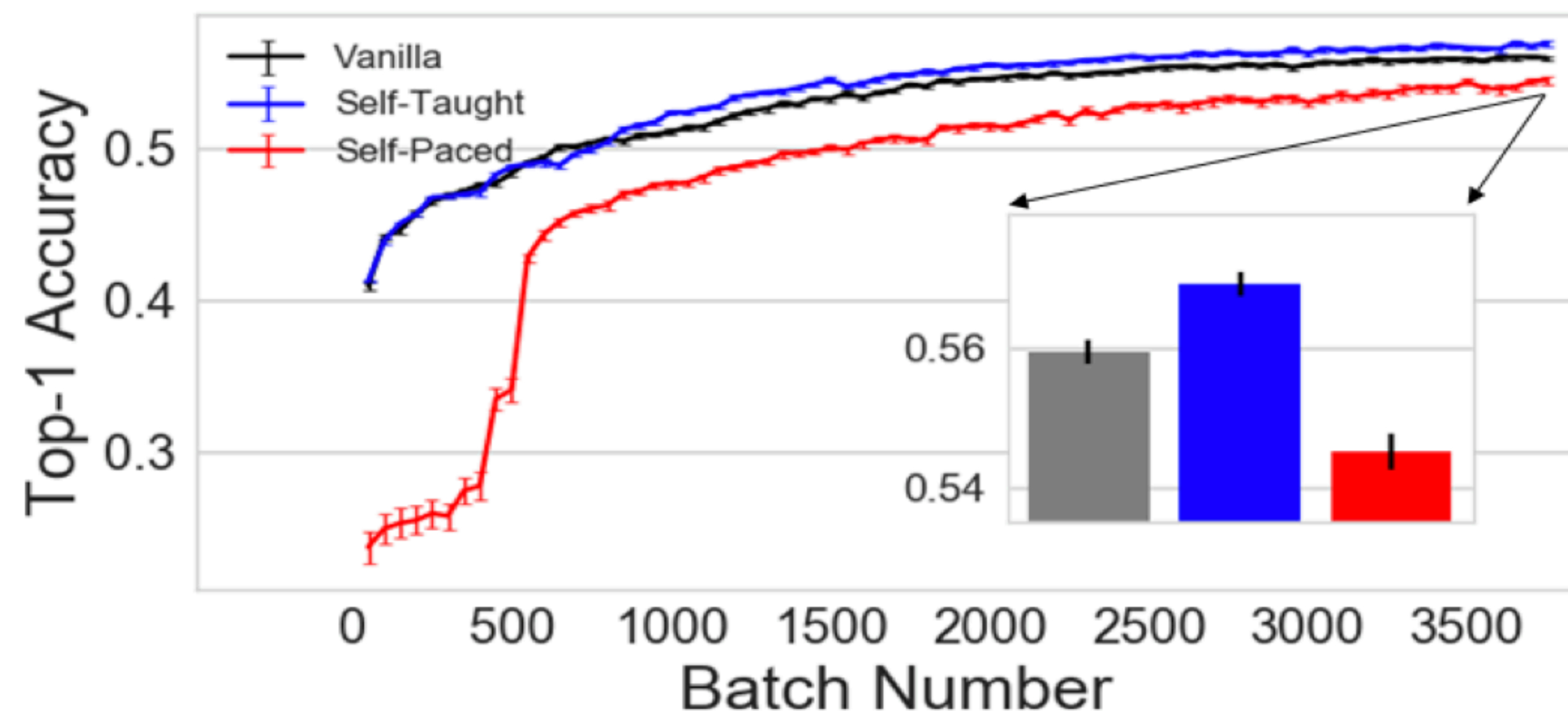


## Self-paced learning:

Measure the difficulty of an instance according to current loss/predictions etc. (related to the ideas in *active learning*)

## Self-taught learning:

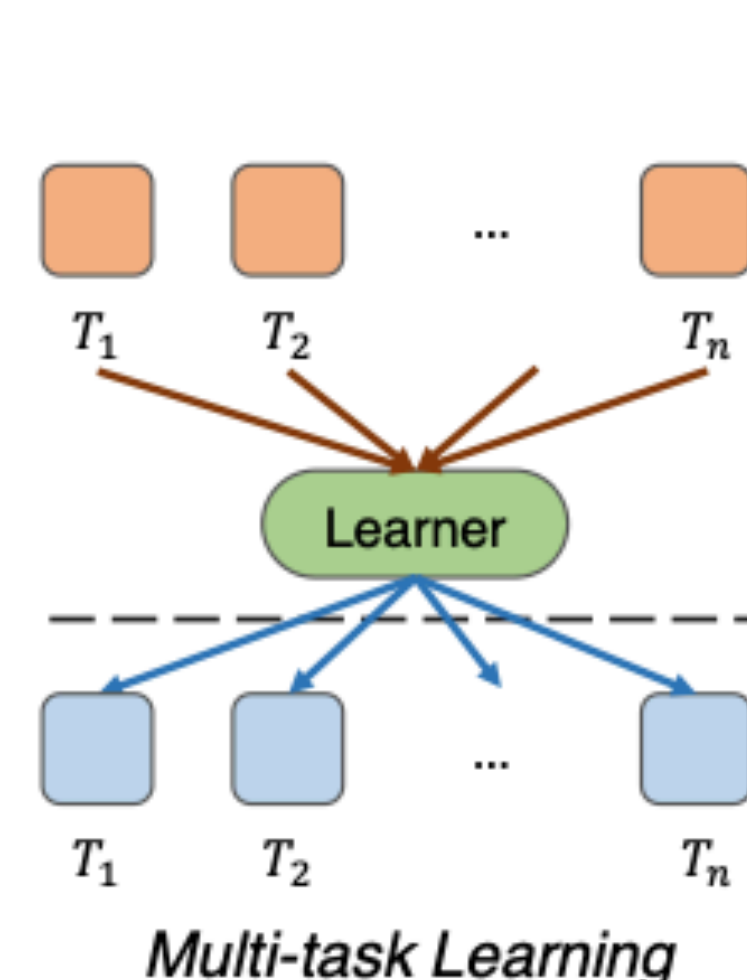
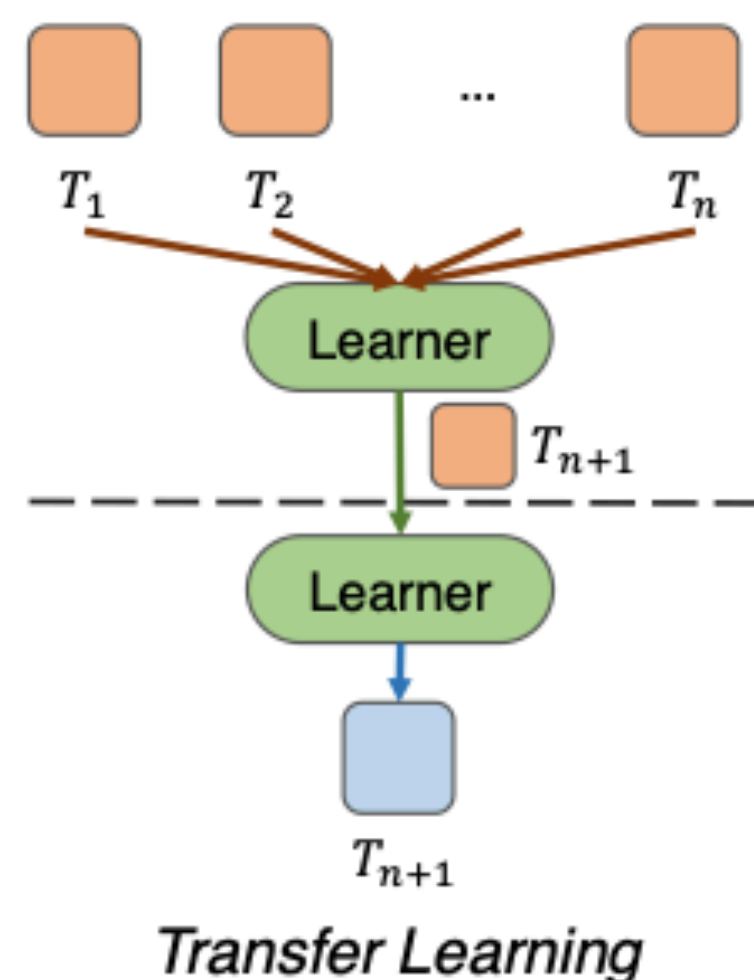
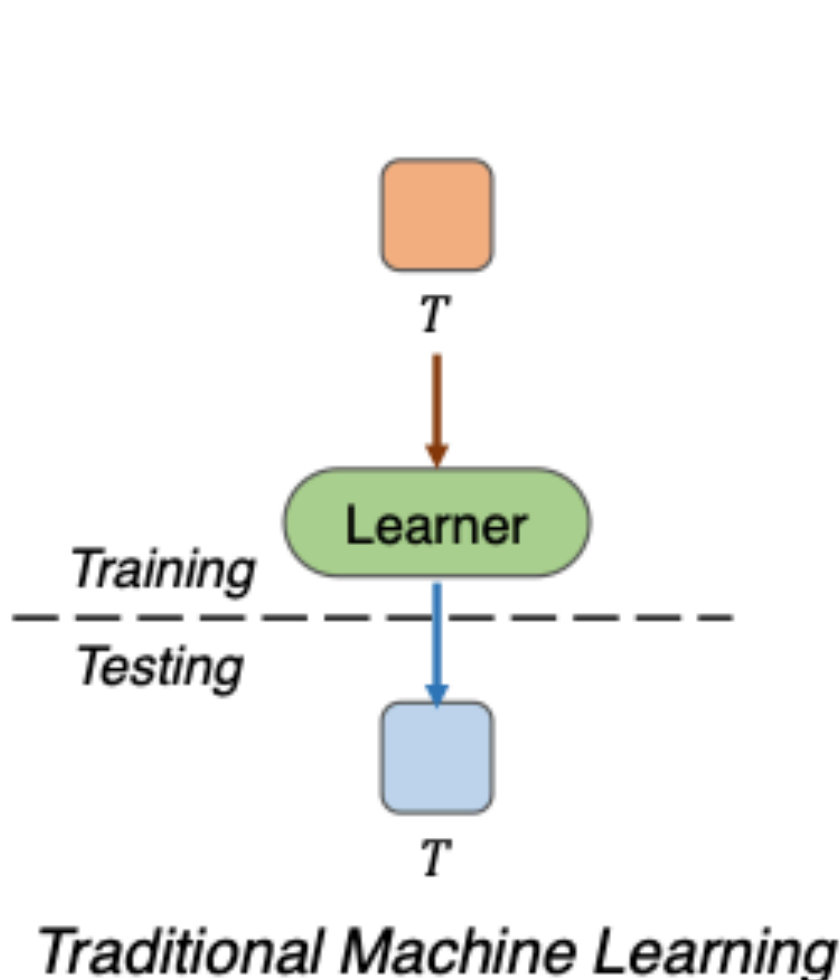
Train a model fully, measure each instance according to final model, assign difficulty score and start over with curriculum -> repeat (related to the ideas in *boosting*)



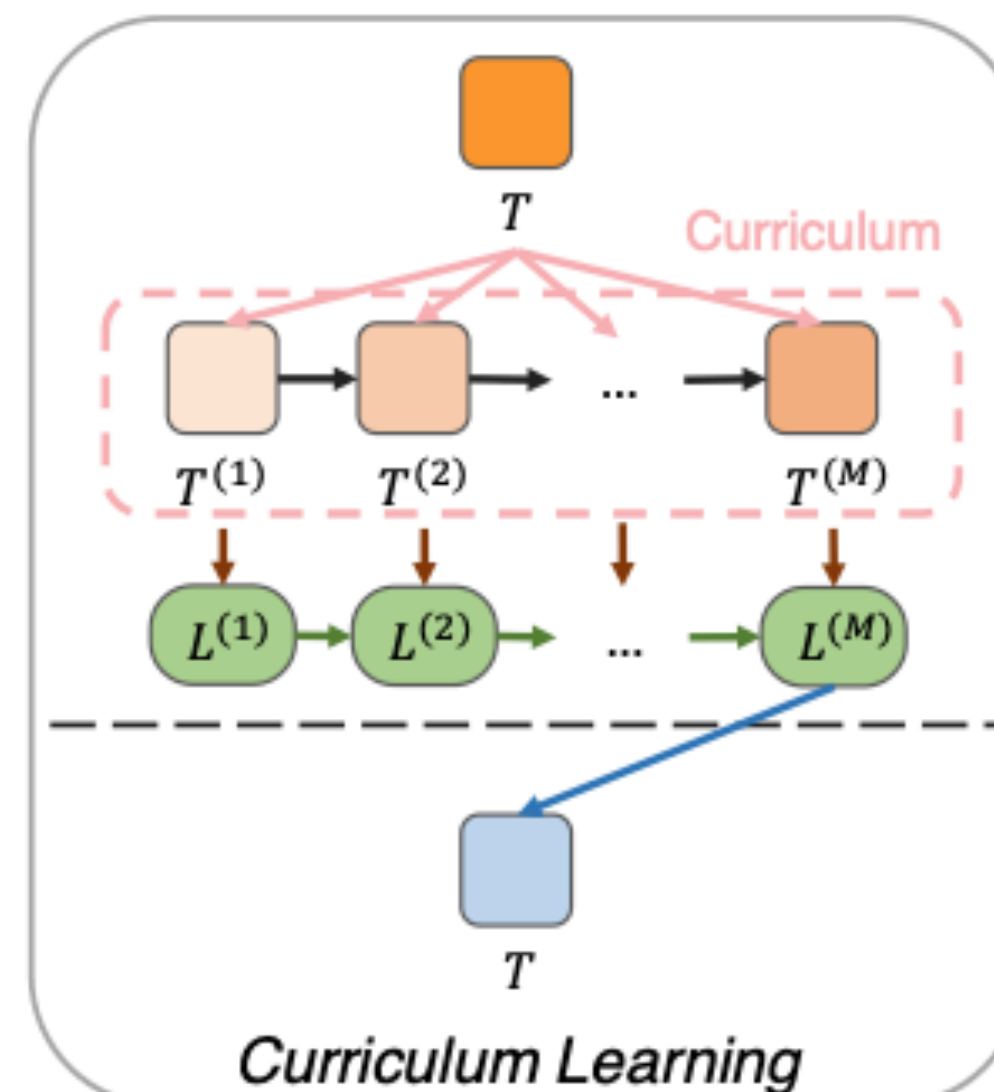
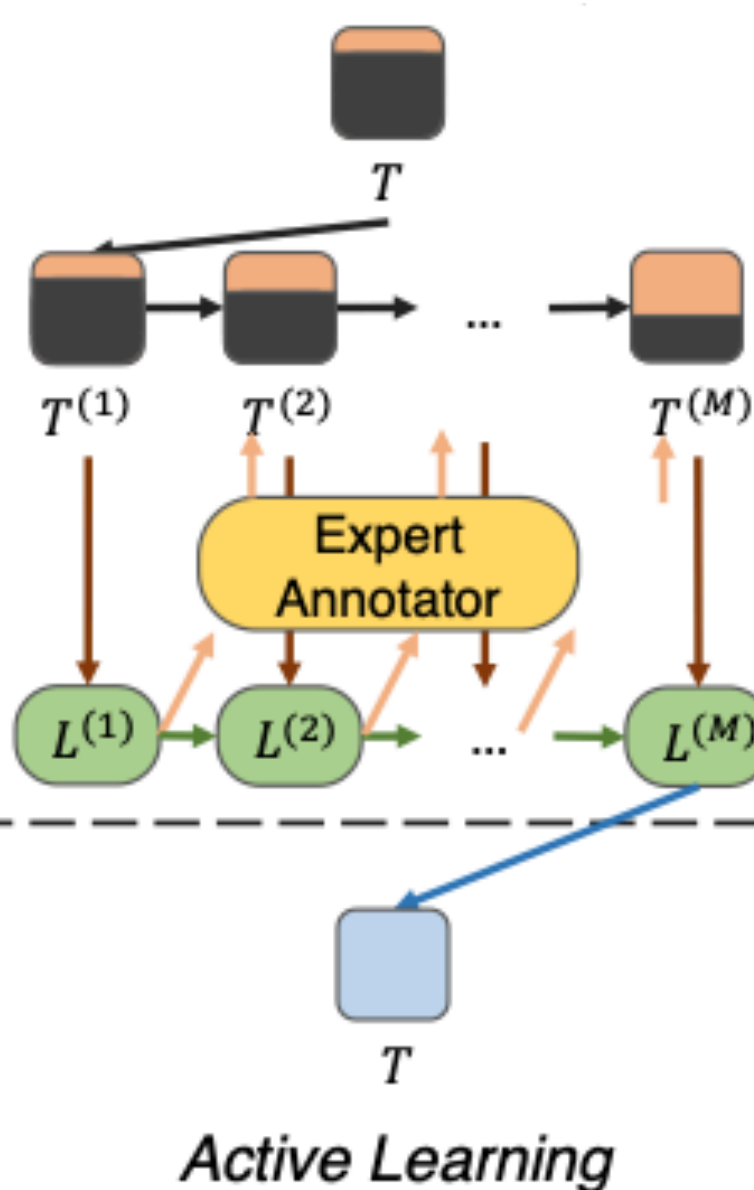
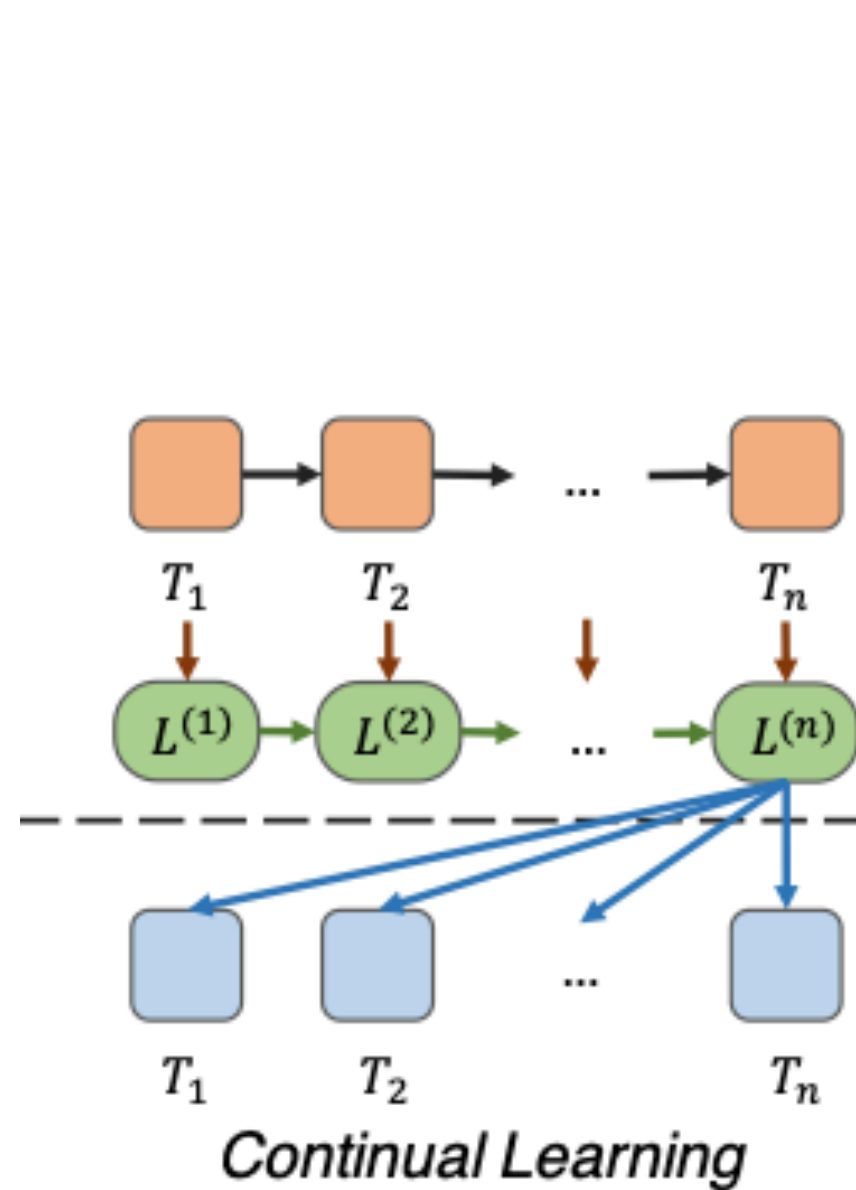


Again: why is it so hard to beat “random”?  
"wrong" things to measure & constrained evaluation

# It's about set-up & evaluation (our topic tomorrow)

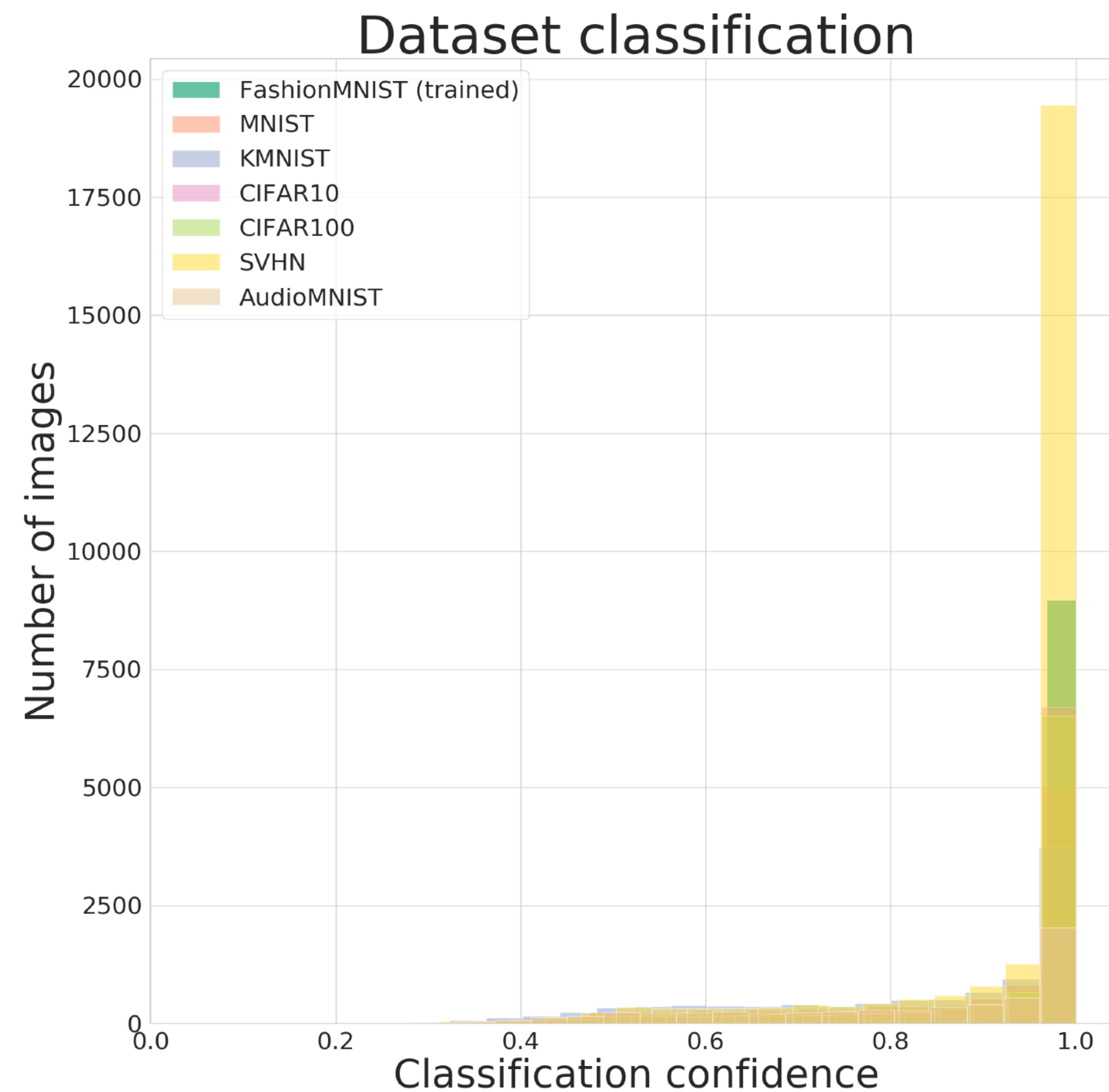
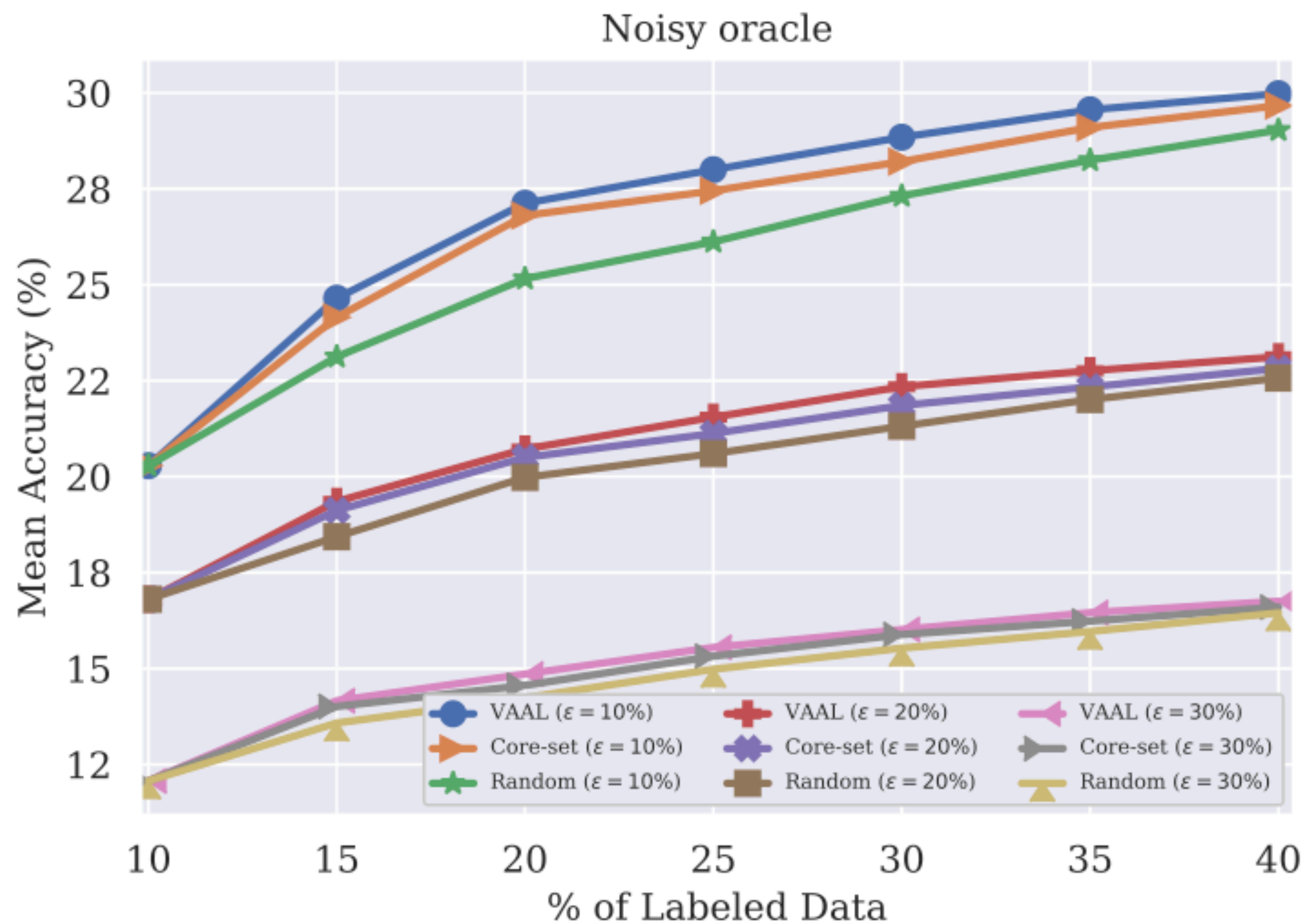


- Training
- Testing
- Model update / Finetune
- Annotation path in AL
- Sequence (seq.) of tasks
- Training / Testing data
- Unlabeled training data
- L<sup>(i)</sup> Learner at step i in seq.
- L<sub>j</sub> Specific learner for task j



Wang et al, "A Survey on Curriculum Learning", TPAMI 2021

# We have consistently assumed A LOT! Tomorrow's essence: opening "Pandora's box" of evaluation



Sinha et al, "Variational Adversarial Active Learning", ICCV 2019

Mundt et al "Open Set Recognition Through Deep Neural Network Uncertainty, Does Out-of-Distribution Detection Require Generative Classifiers?", ICCV Statistical Deep Learning Workshop 2019 (Based on a long-known problem, Matan1990)