

# Not all Textual Instances are Alike: Efficient NLP by Better Understanding of our Data

**Roy Schwartz**

Hebrew University of Jerusalem  
SustainNLP 2021

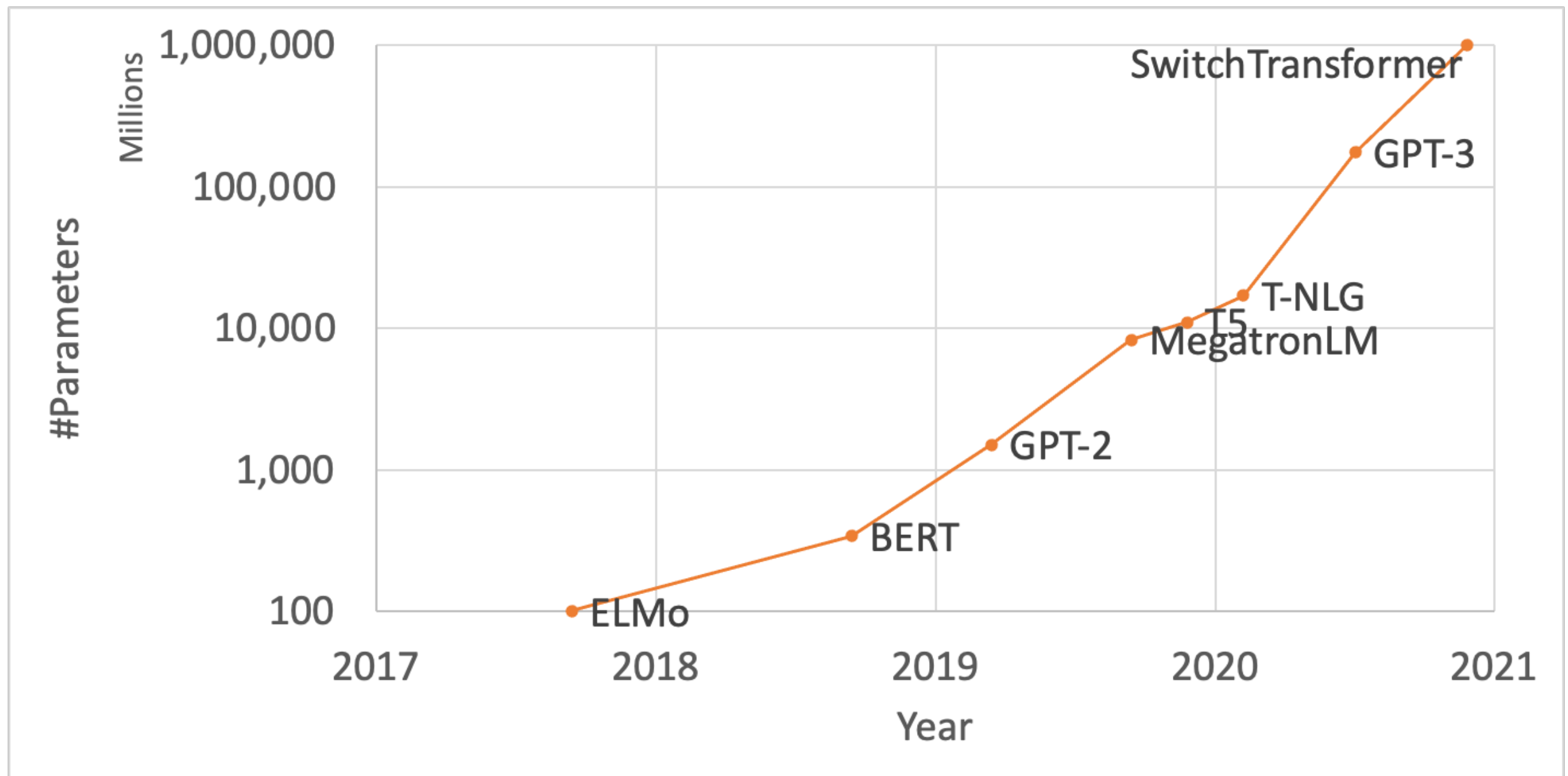


THE HEBREW  
UNIVERSITY  
OF JERUSALEM



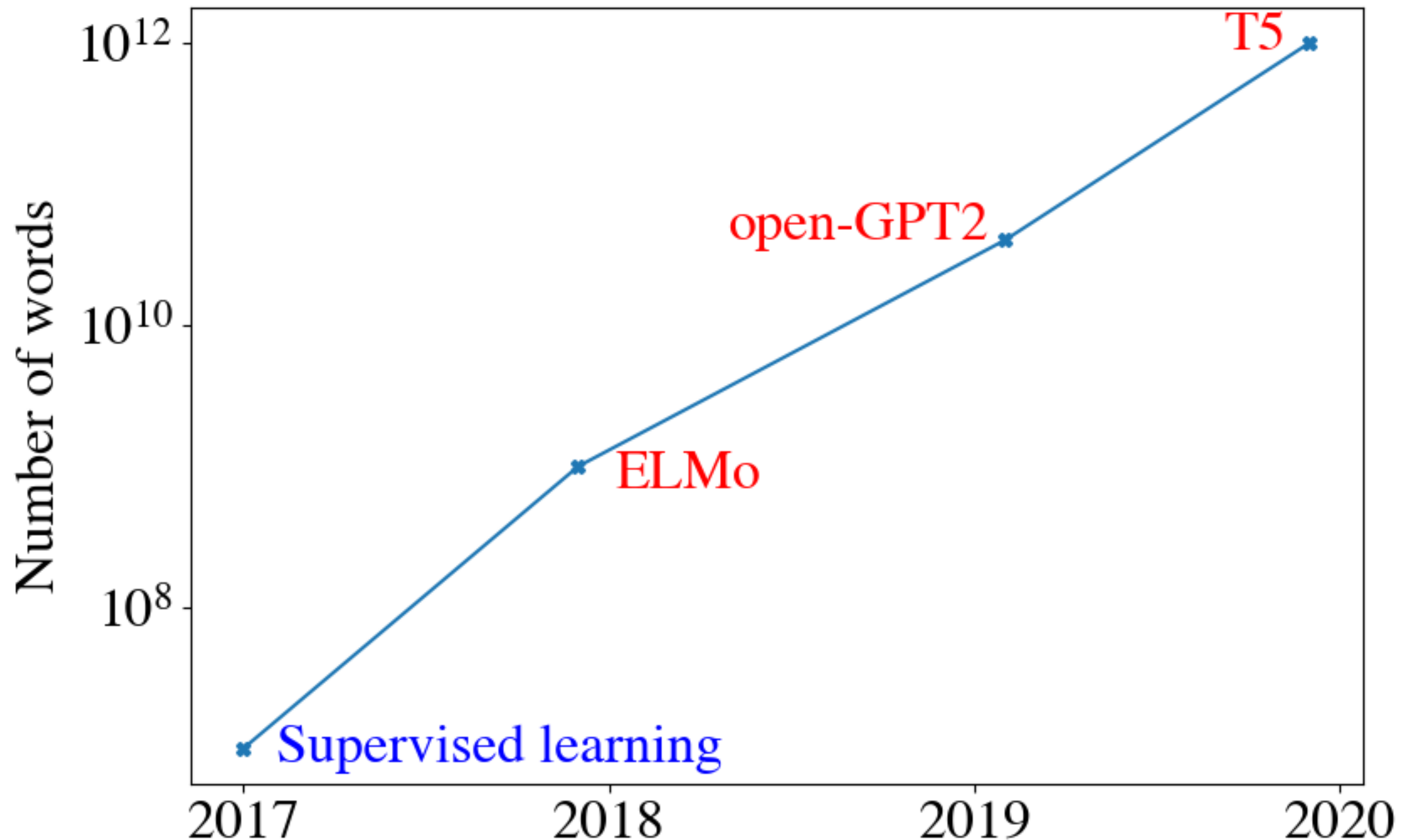
# Premise: **Big** Models

## 10,000X in 3 Years



# Large Datasets

100,000X in 3 Years



# Efficiency

## Current Approaches



# Efficiency

## Current Approaches



# Efficiency

## Current Approaches

- **Model distillation**

- Hinton et al. (2015); MobileBERT (Sun et al., 2019); DistilBERT (Sanh et al., 2019)



- **Pruning / Structural Pruning**

- Han et al. (2016); SNIP (Lee et al., 2019); LTH (Frankle & Corbin, 2019); MorphNet (Gordon et al., 2018); Michel et al. (2019); LayerDrop (Fan et al., 2020); Dodge, **Schwartz** et al. (2019)

- **Quantization**

- Gong et al. (2014); Q8BERT (Zafrir et al., 2019); Q-BERT (Shen et al., 2019)

# Data in NLP

**Basic Assumption: Instances are IID**





# Not all Instances are Alike

1. *The movie was awesome.*
2. *I could definitely see why this movie received such great critiques, but at the same time I can't help but wonder whether the plot was written by a 12 year-old or by an award-winning writer.*



# Not all Instances are Alike

1. *The movie was awesome.*
2. *I could definitely see why this movie received such great critiques, but at the same time I can't help but wonder whether the plot was written by a 12 year-old or by an award-winning writer.*

**What is the capital of Italy?**

**Which country won the largest number of swimming medals in the 2016 summer olympics?**

**Would a glass of water that falls from 10 feet down to a trampoline break?**

# Outline

## Not all Instances are Alike

- Efficient inference
  - Schwartz et al., ACL 2020
- Efficient training
  - Swayamdipta et al., EMNLP 2020
- Better masked language modeling for vision and language
  - Bitton et al., Findings of EMNLP 2021

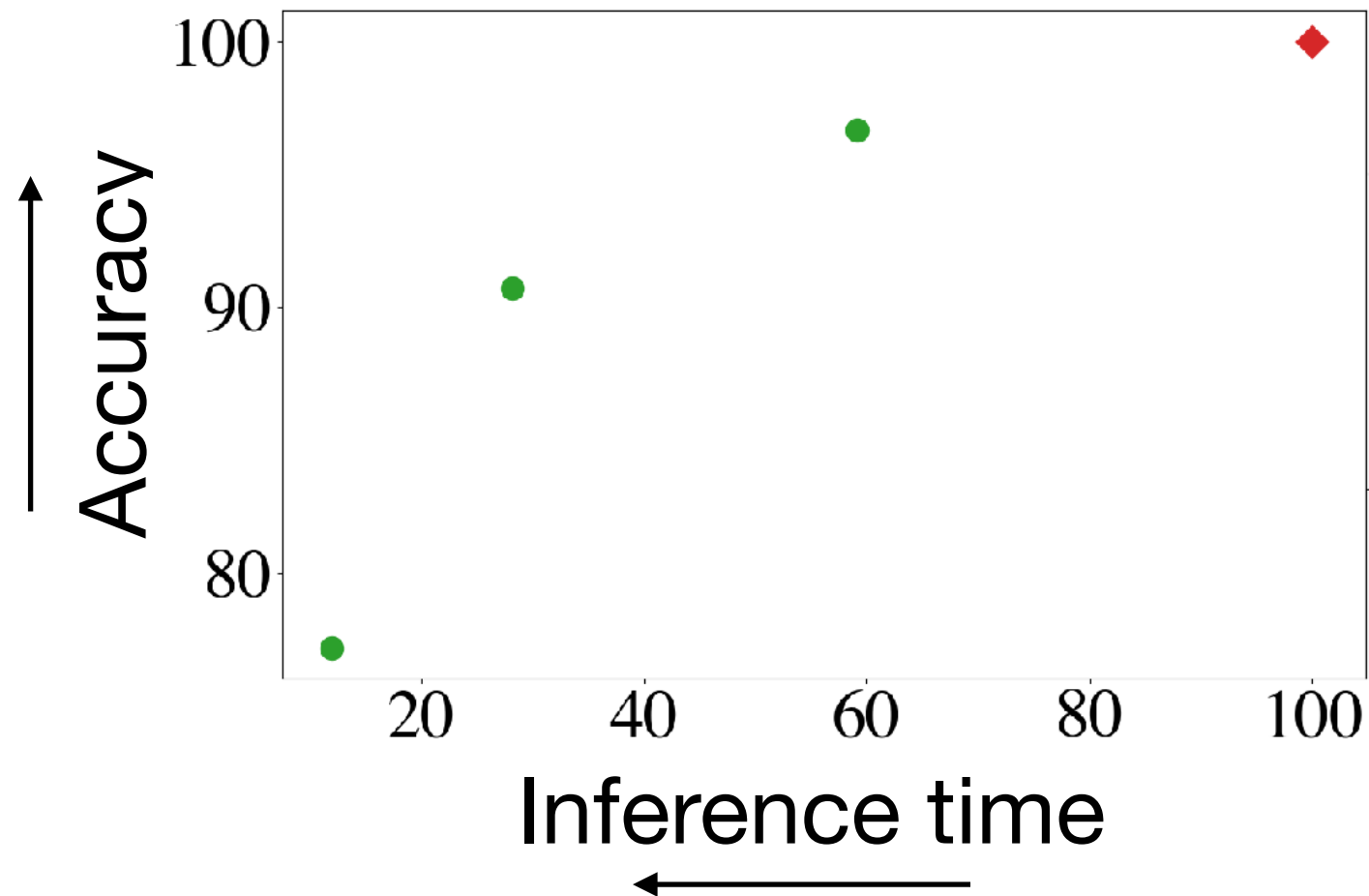
# Case Study 1: Efficient Inference

Schwartz et al., ACL 2020

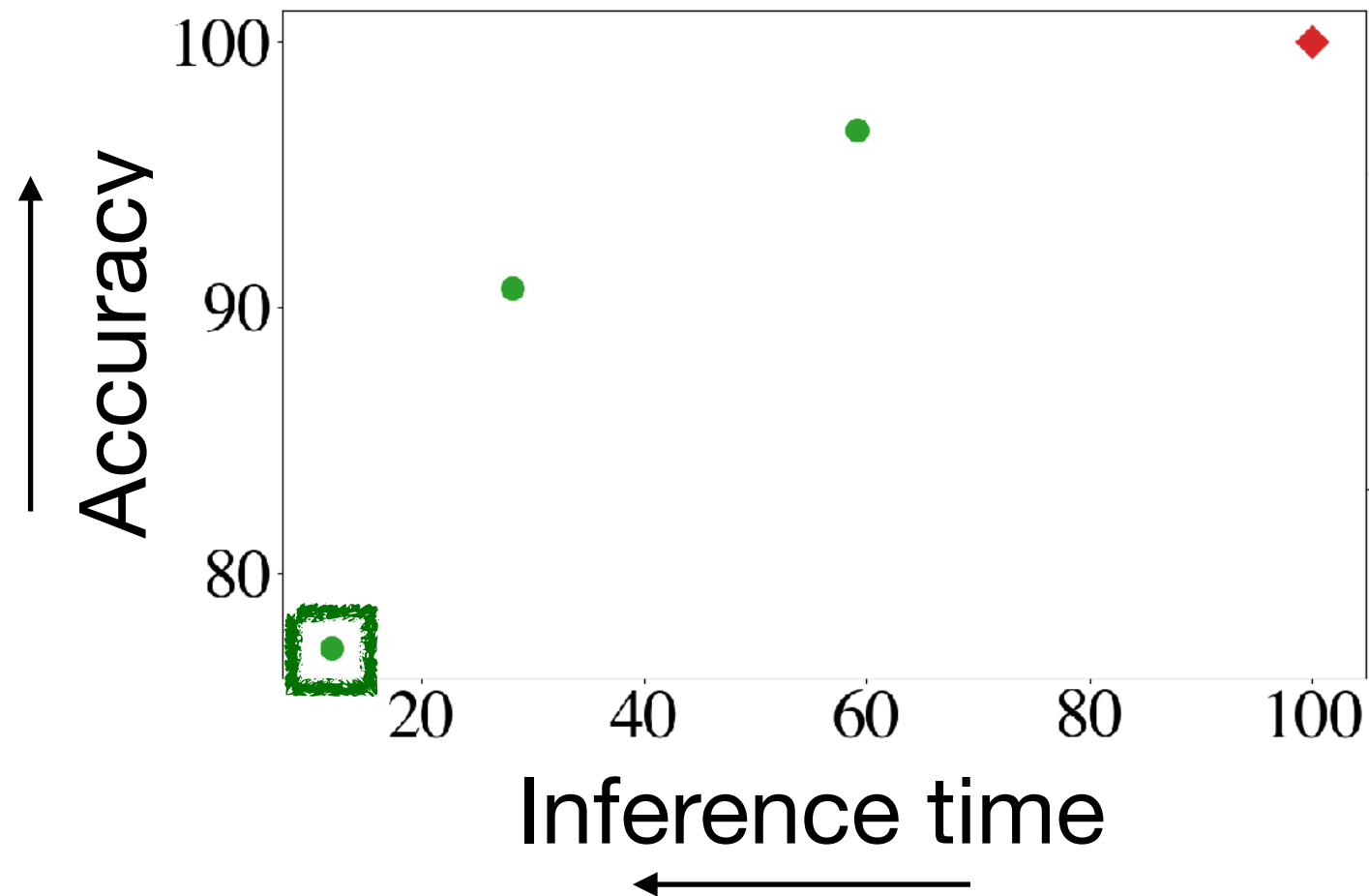
*Some instances require **less  
processing** than others*



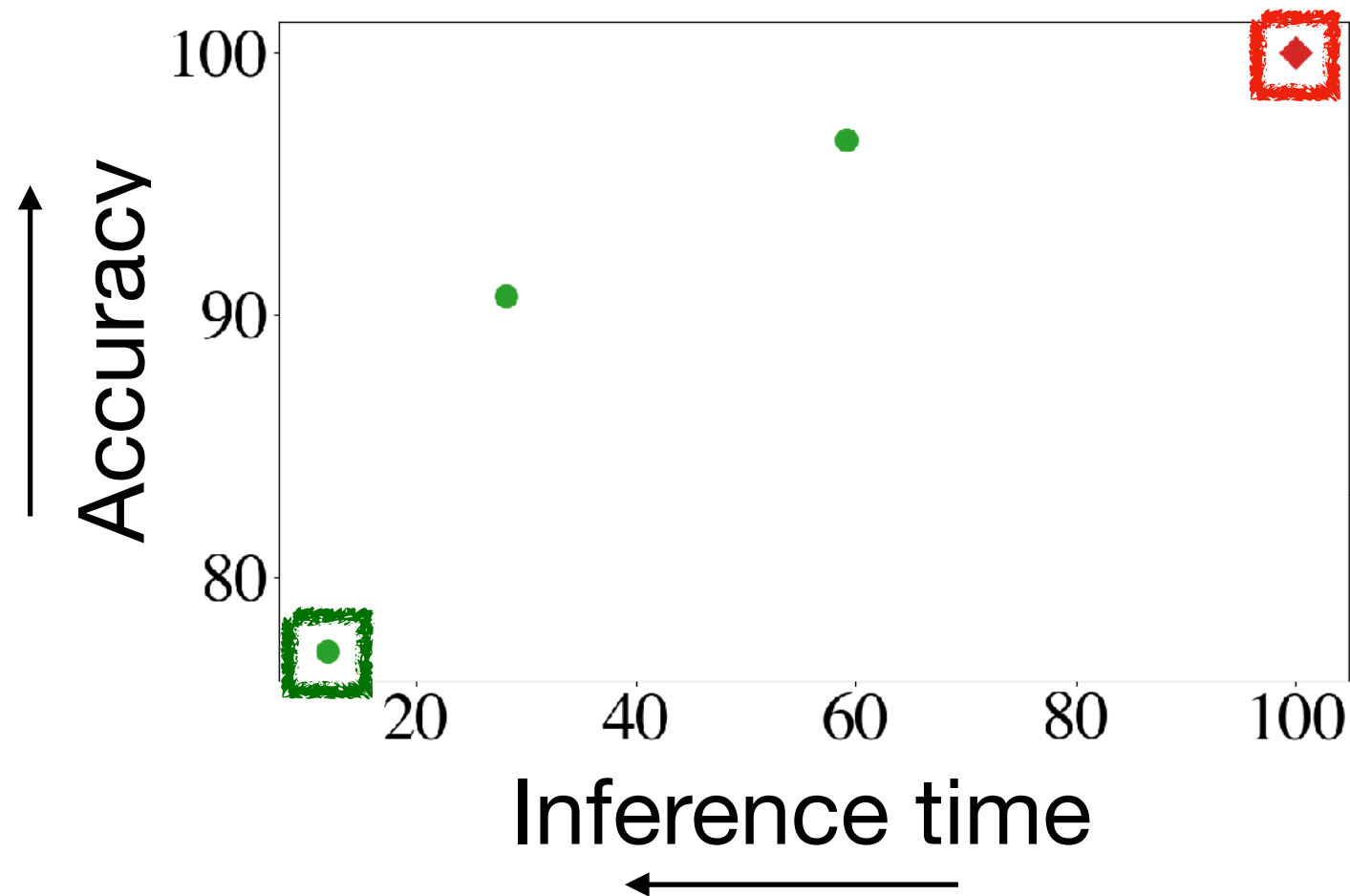
# High-Level Idea



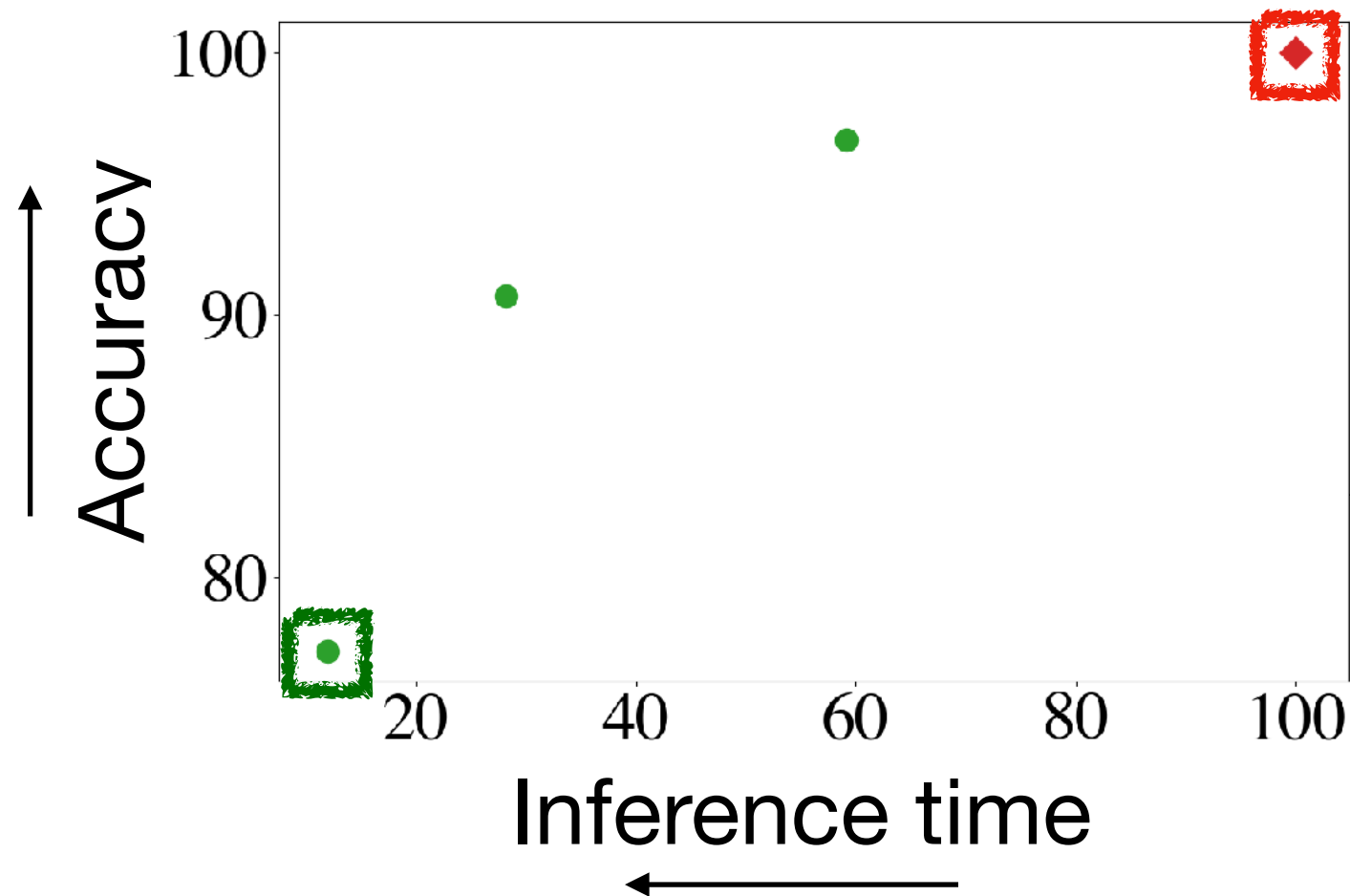
# High-Level Idea



# High-Level Idea



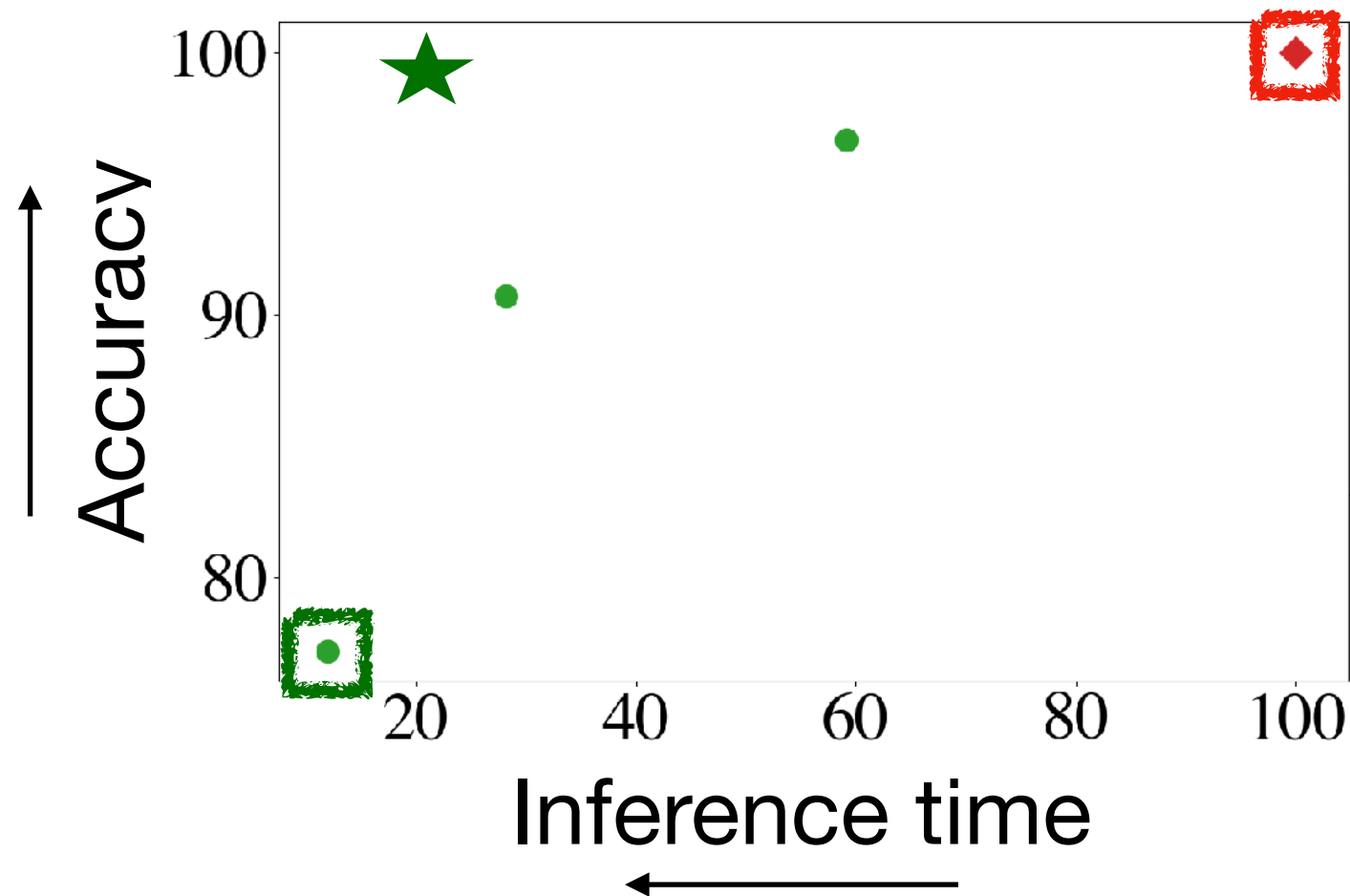
# High-Level Idea



*Run an **efficient** model on “**easy**” instances,  
and an **expensive** model on “**hard**” instances*

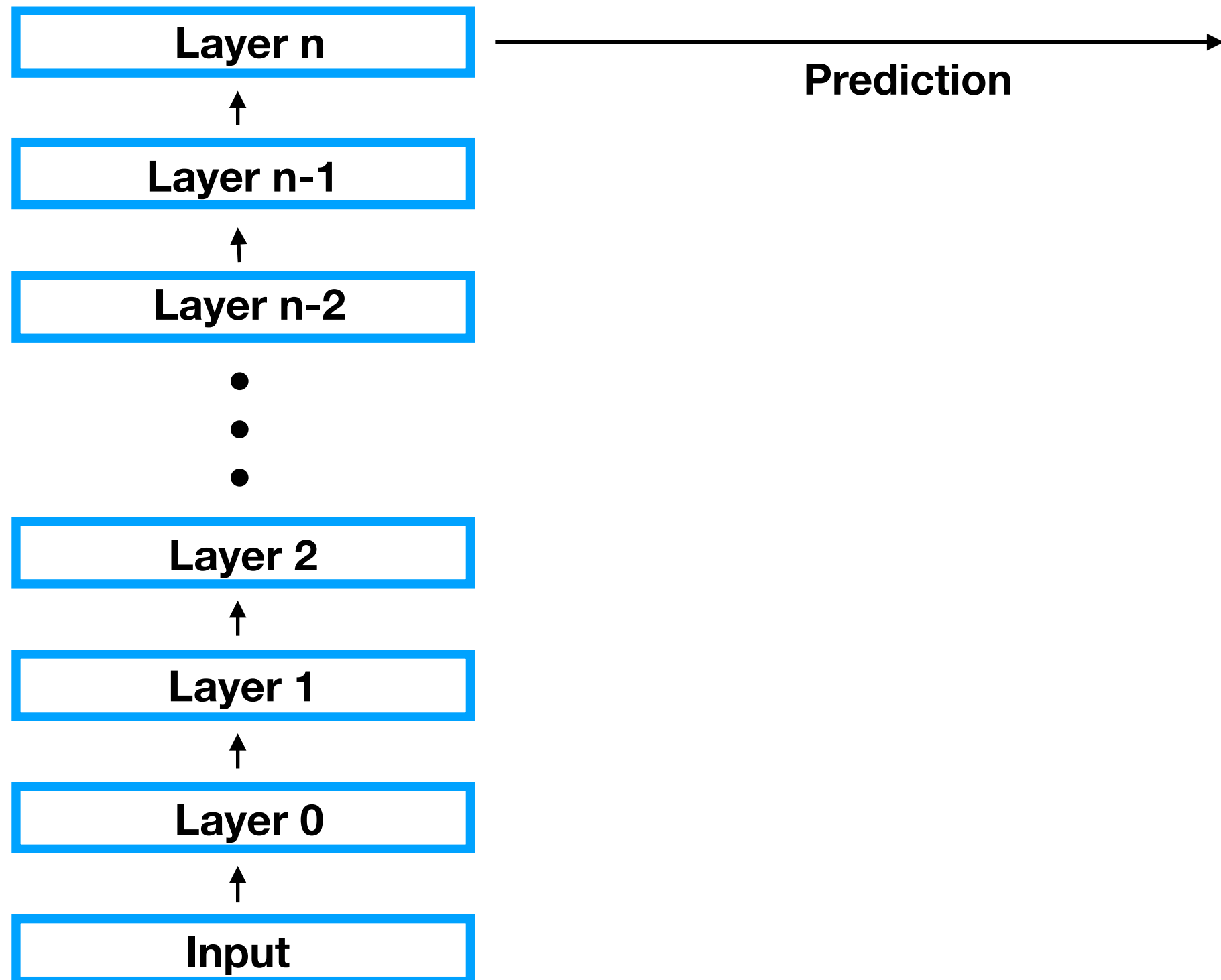


# High-Level Idea

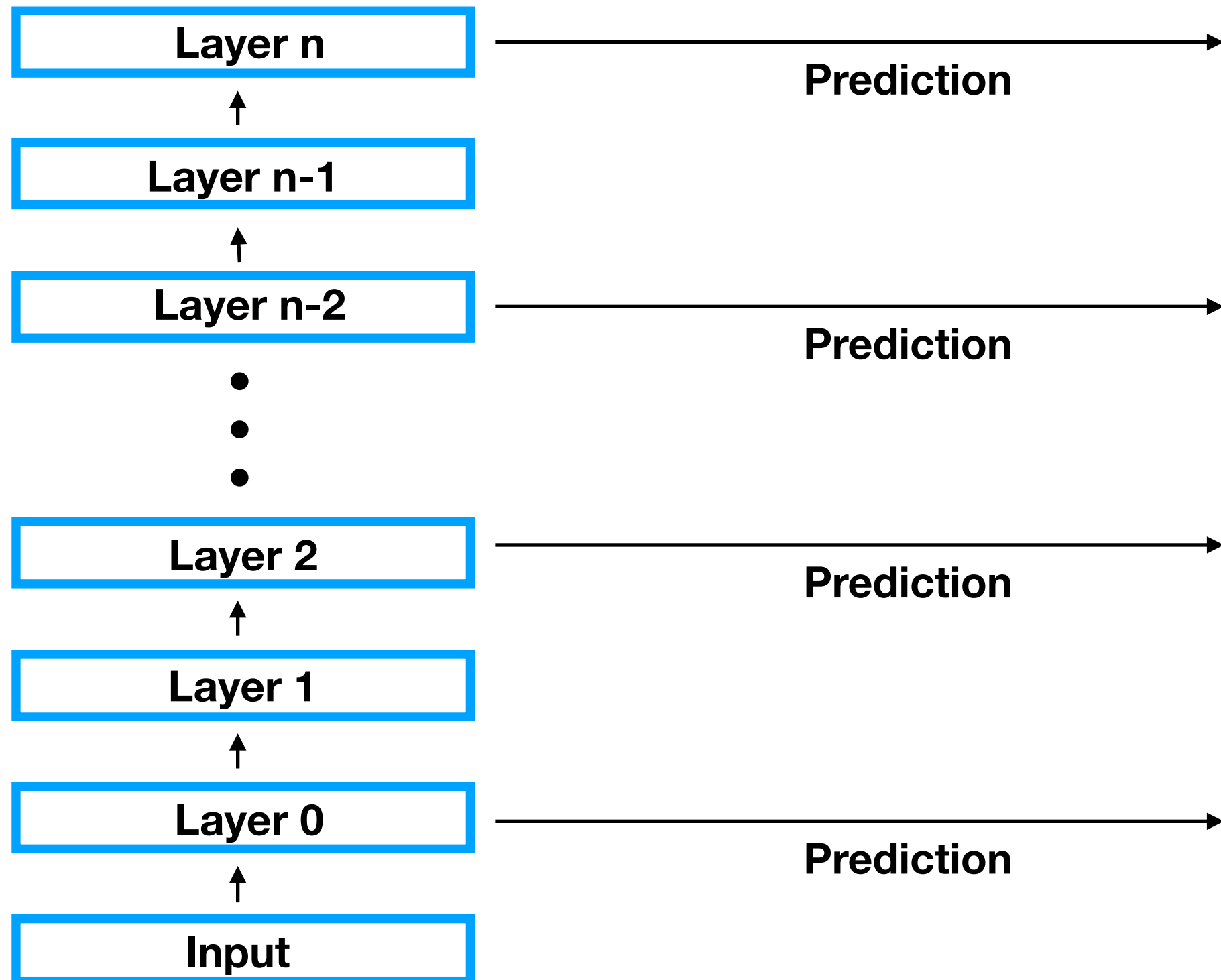


*Run an **efficient** model on “**easy**” instances,  
and an **expensive** model on “**hard**” instances*

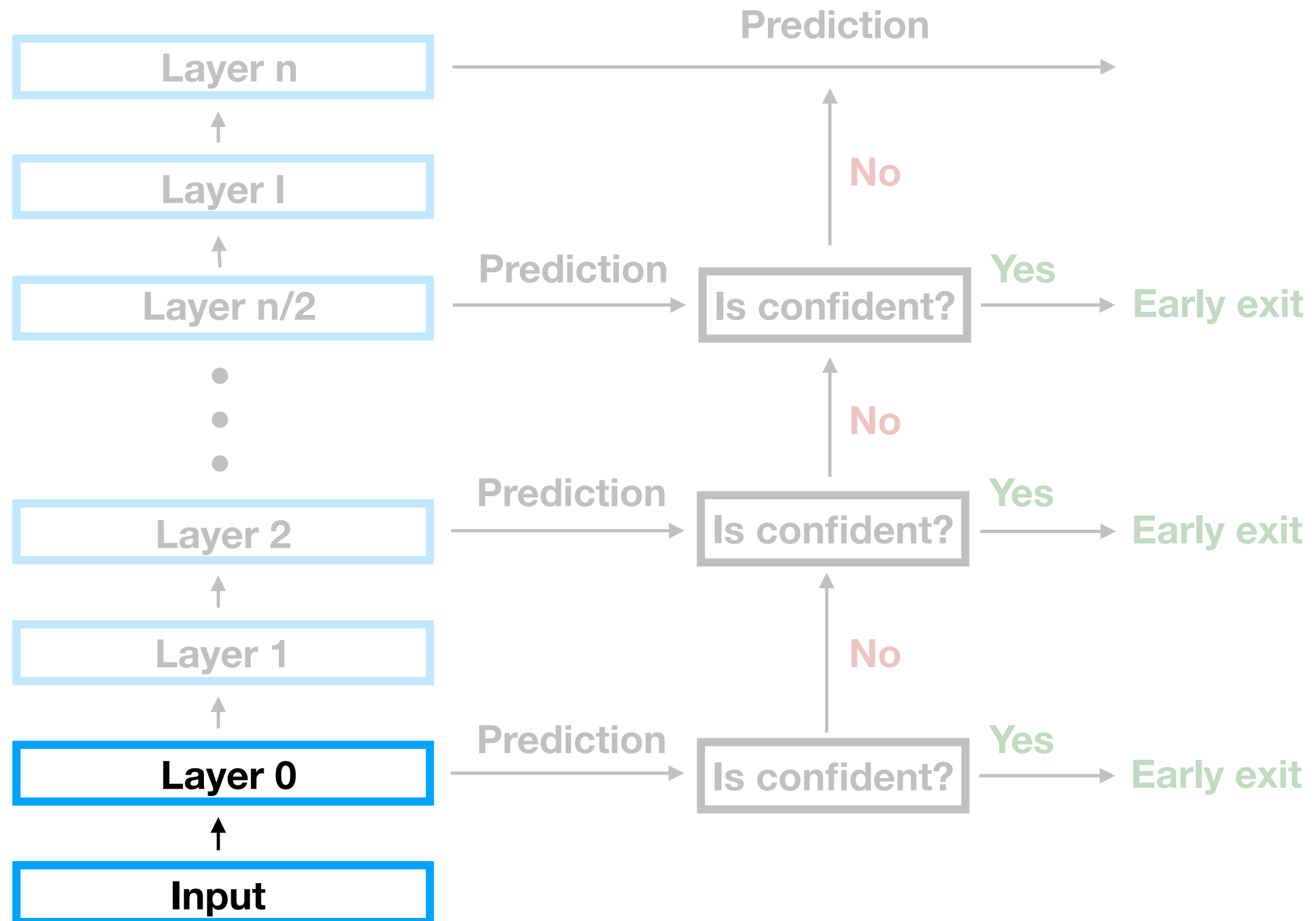
# Our Approach: Training Time



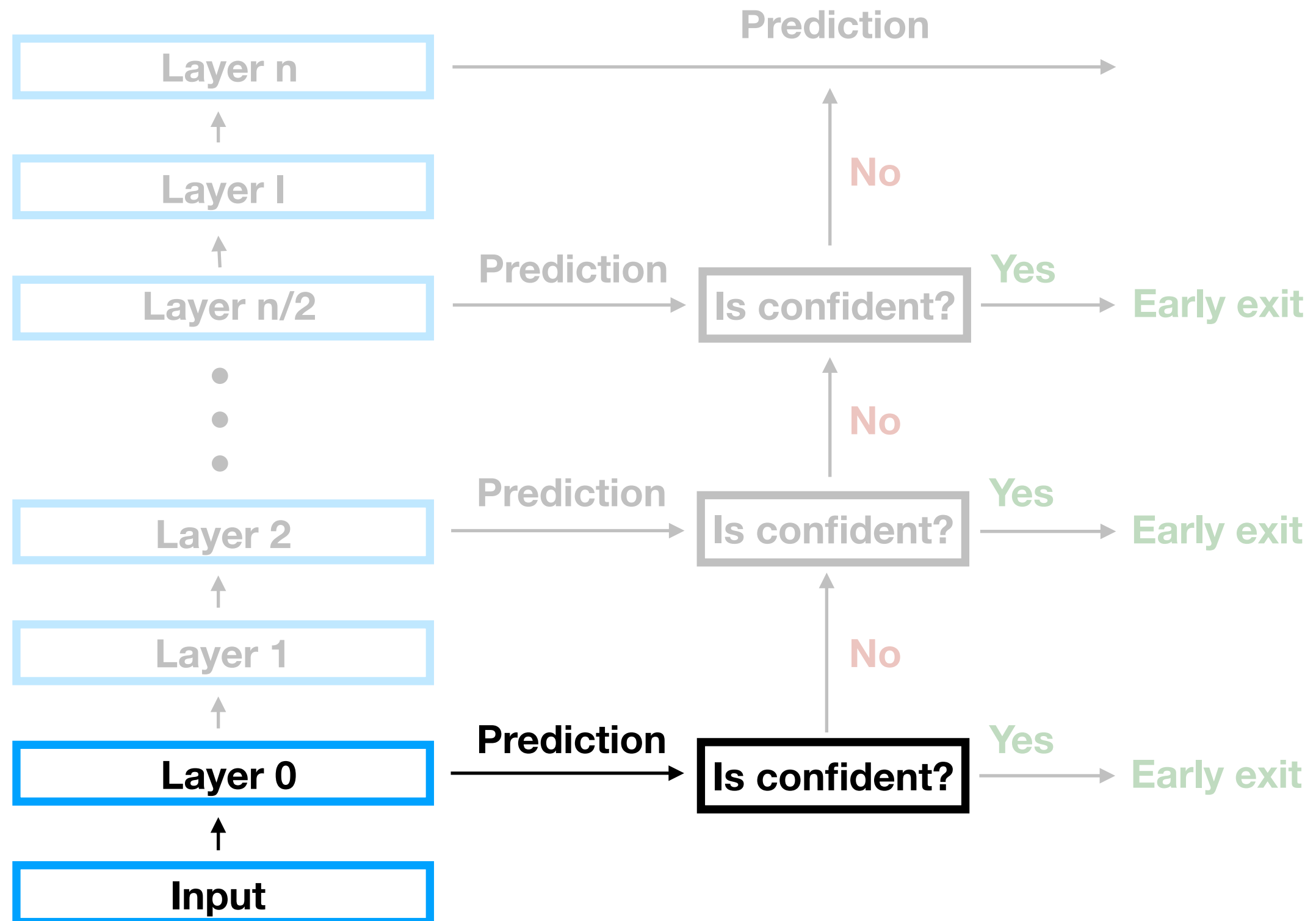
# Our Approach: Training Time



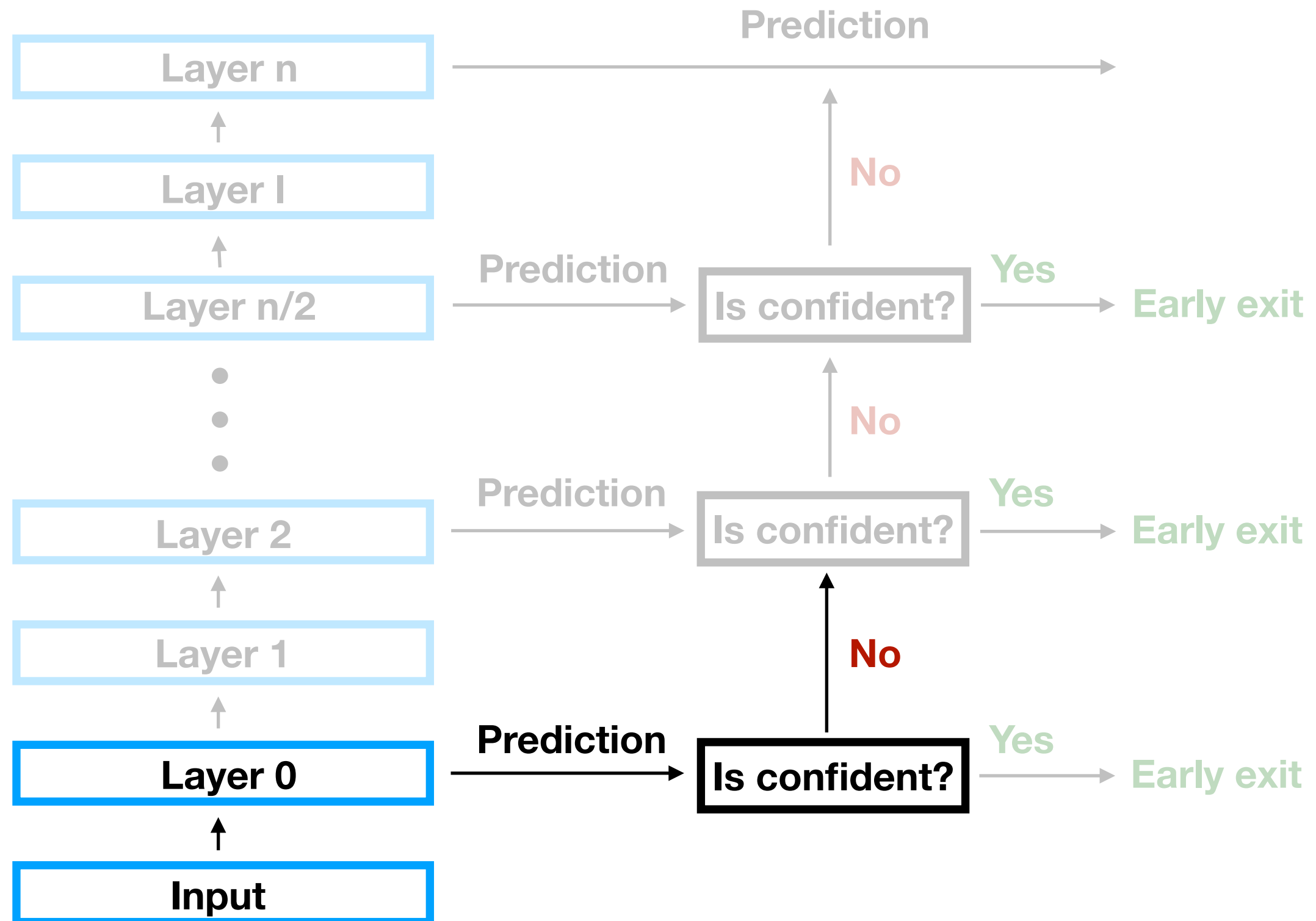
# Our Approach: Test Time



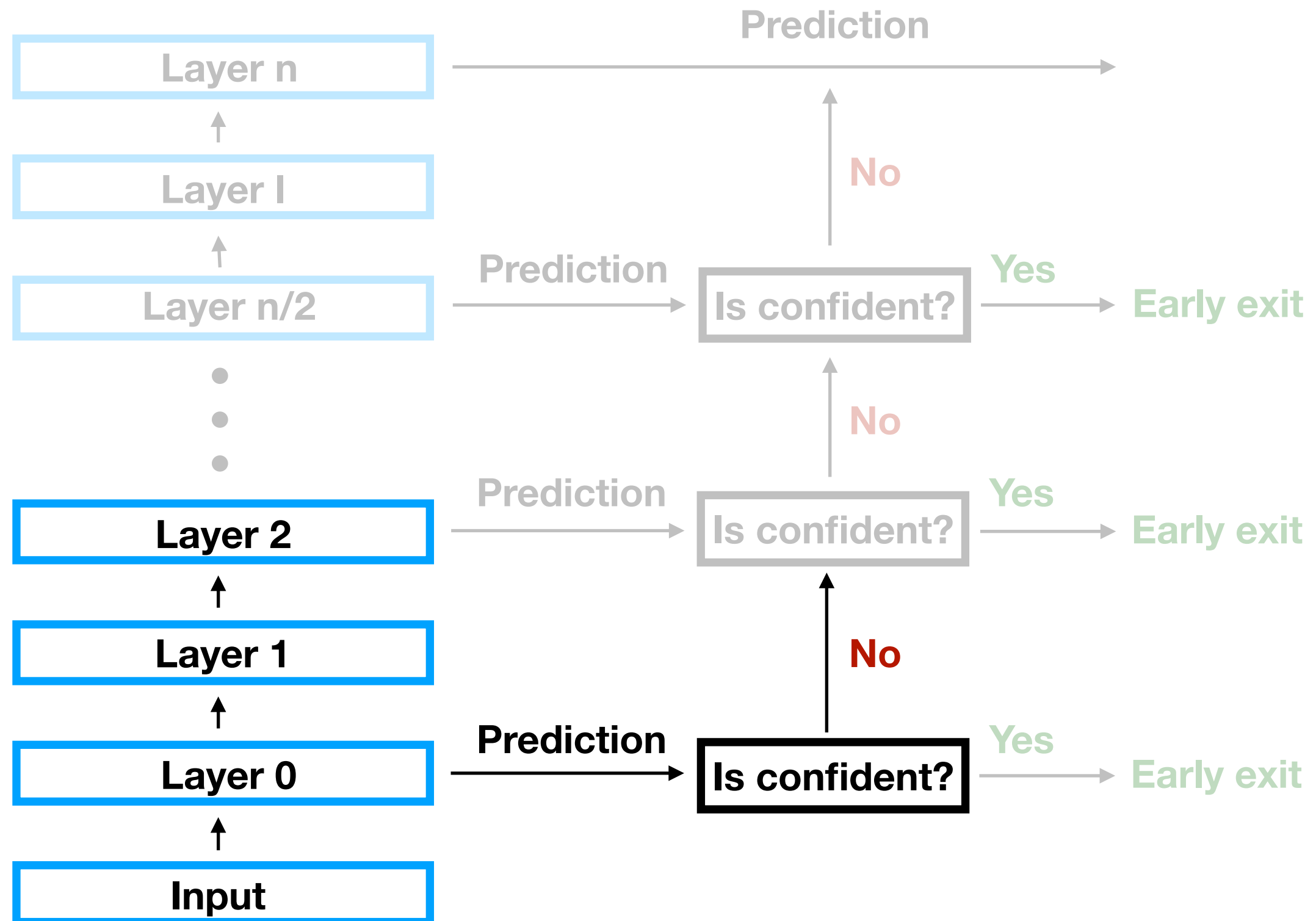
# Our Approach: Test Time



# Our Approach: Test Time

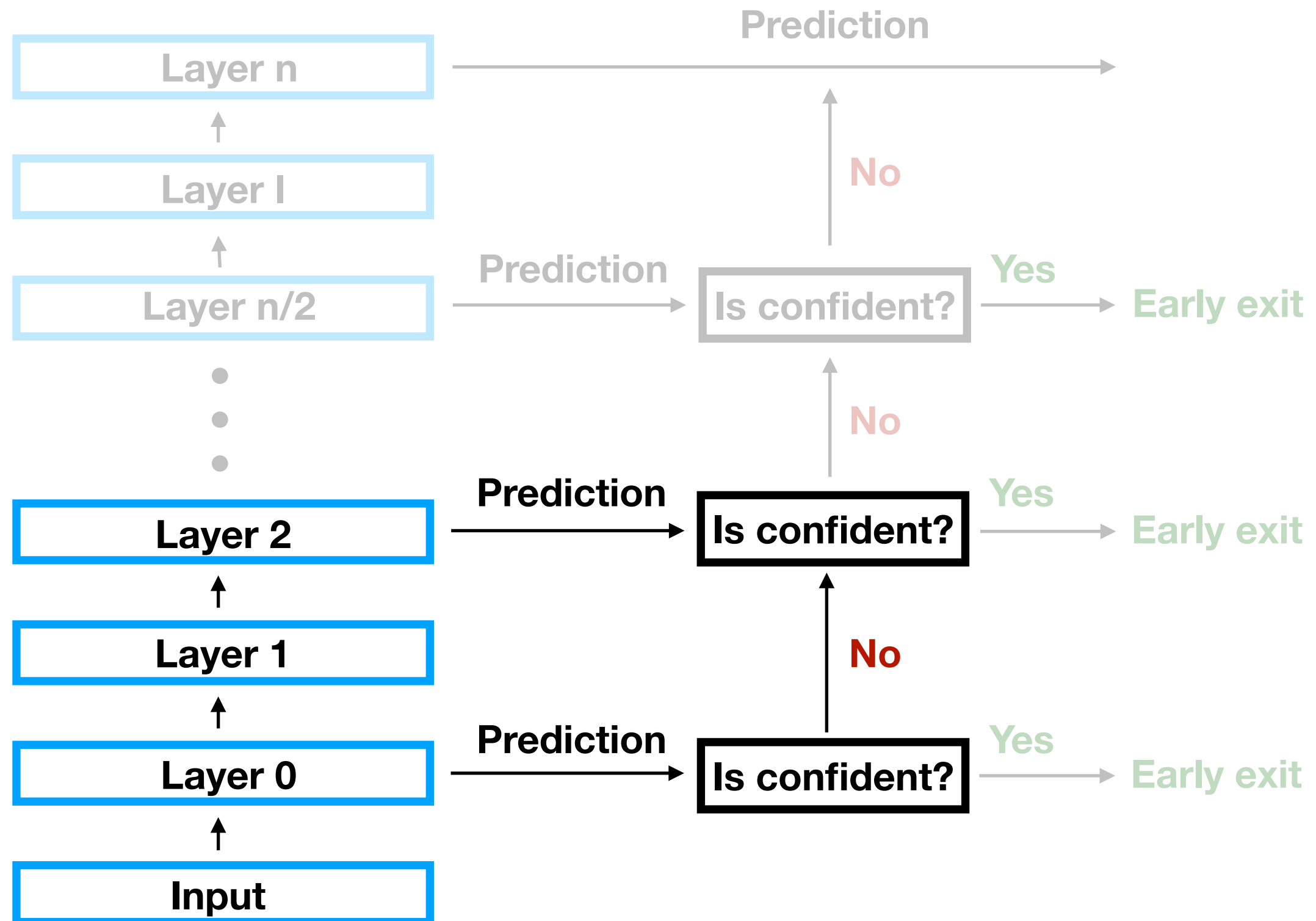


# Our Approach: Test Time

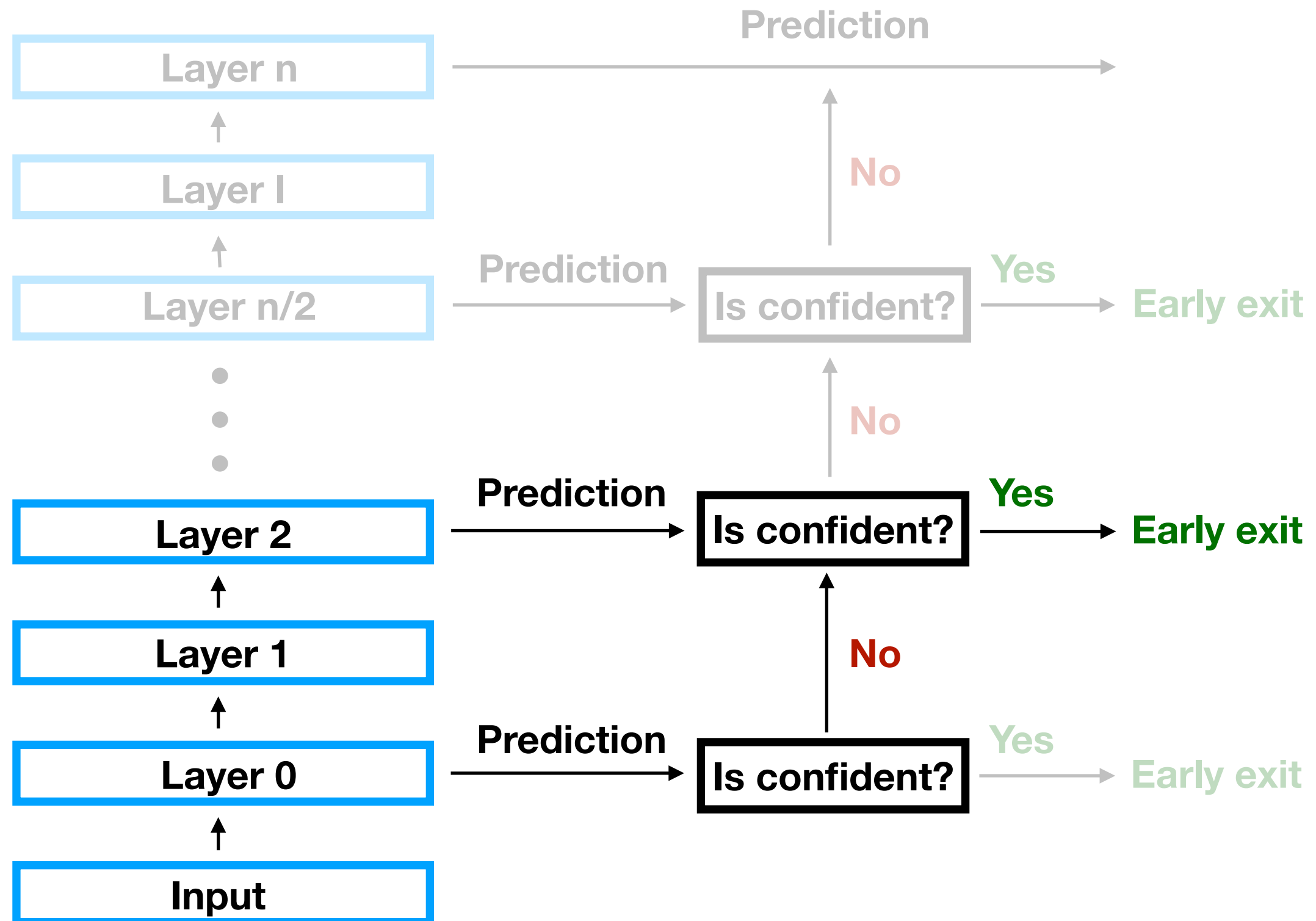




# Our Approach: Test Time



# Our Approach: Test Time



# Calibrated Confidence Scores

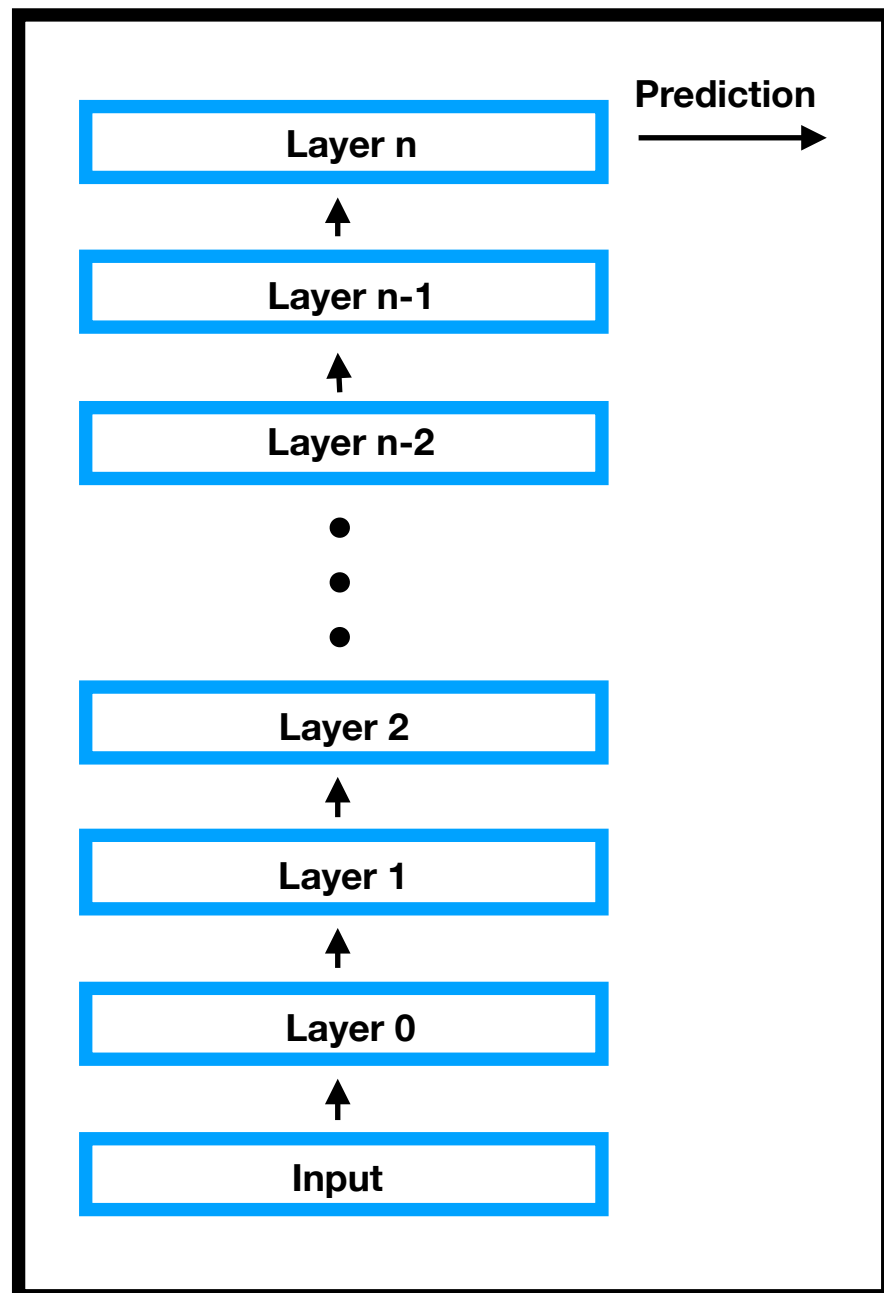
- Interpret the calibrated softmax label scores as model confidence
  - We use temperature calibration (Guo et al., 2017)

# Calibrated Confidence Scores

- Interpret the calibrated softmax label scores as model confidence
  - We use temperature calibration (Guo et al., 2017)
- Speed/accuracy tradeoff controlled by a **single early-exit confidence threshold**

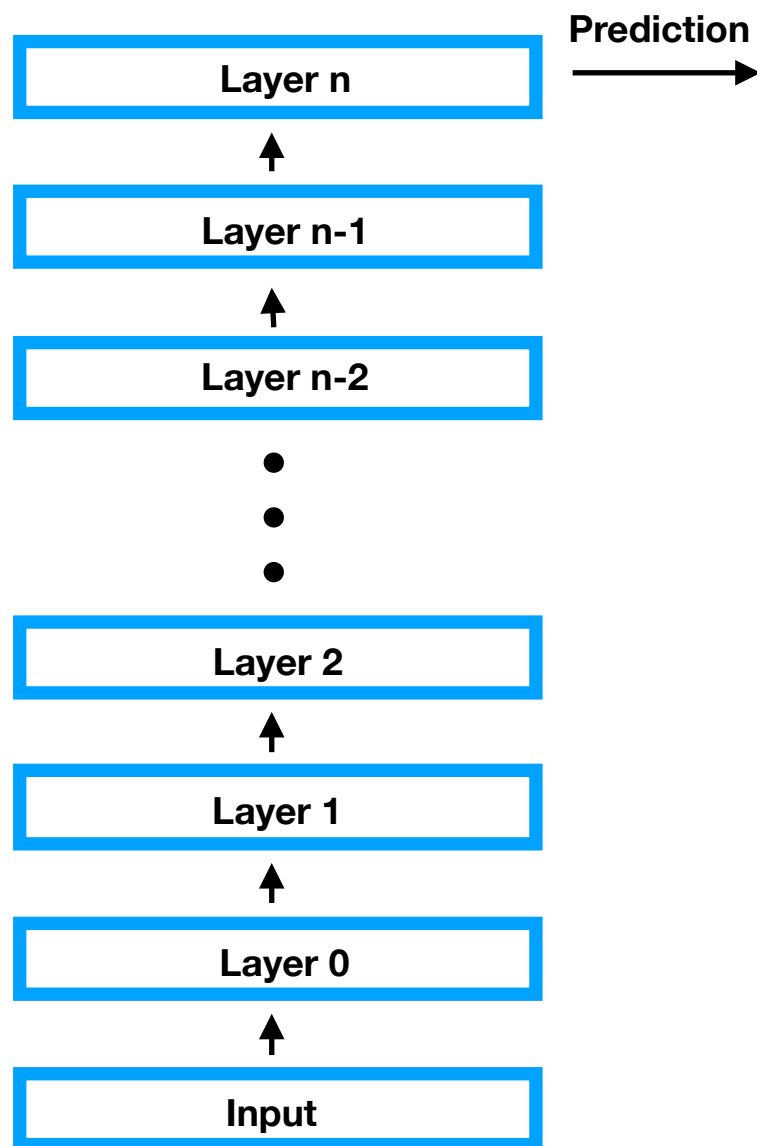
# Baselines

## Standard baseline

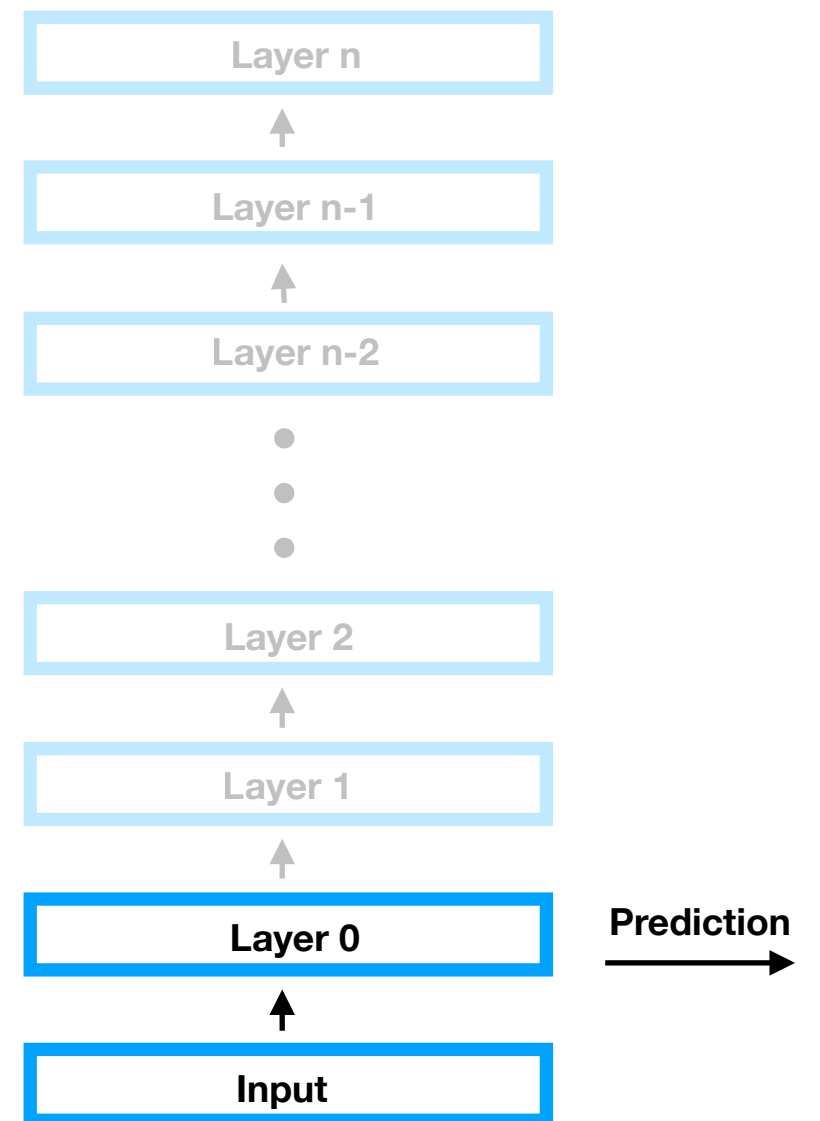
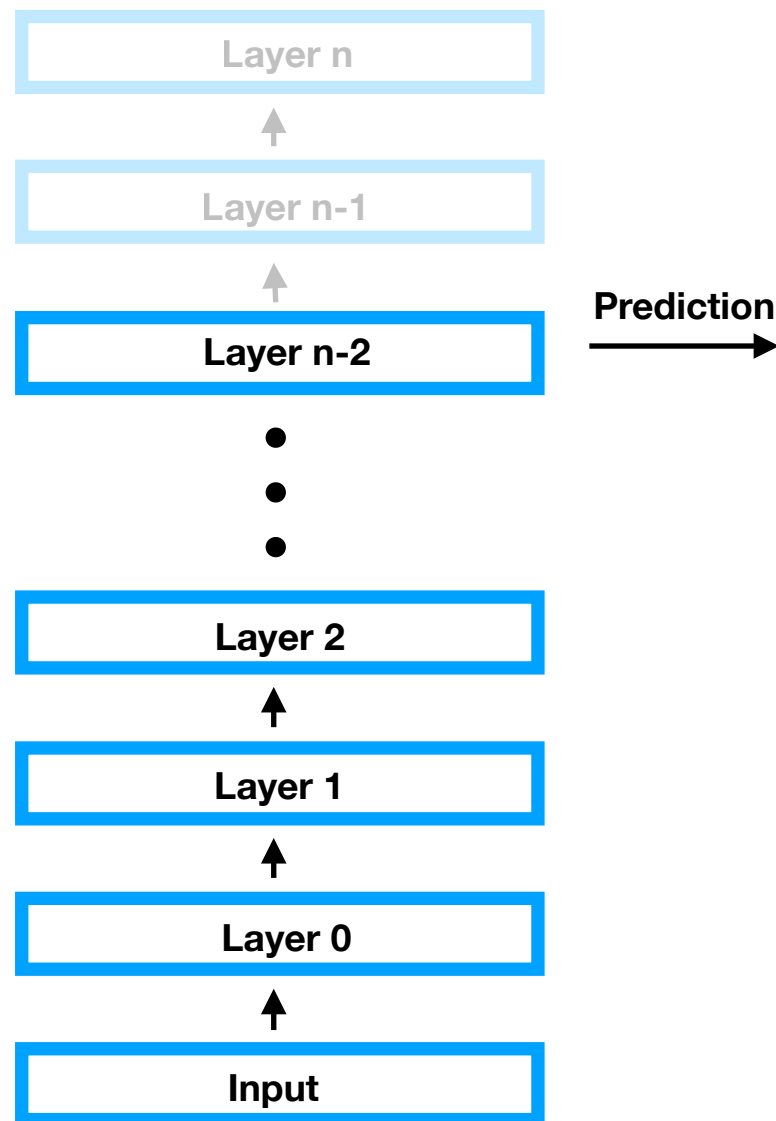


# Baselines

Standard baseline

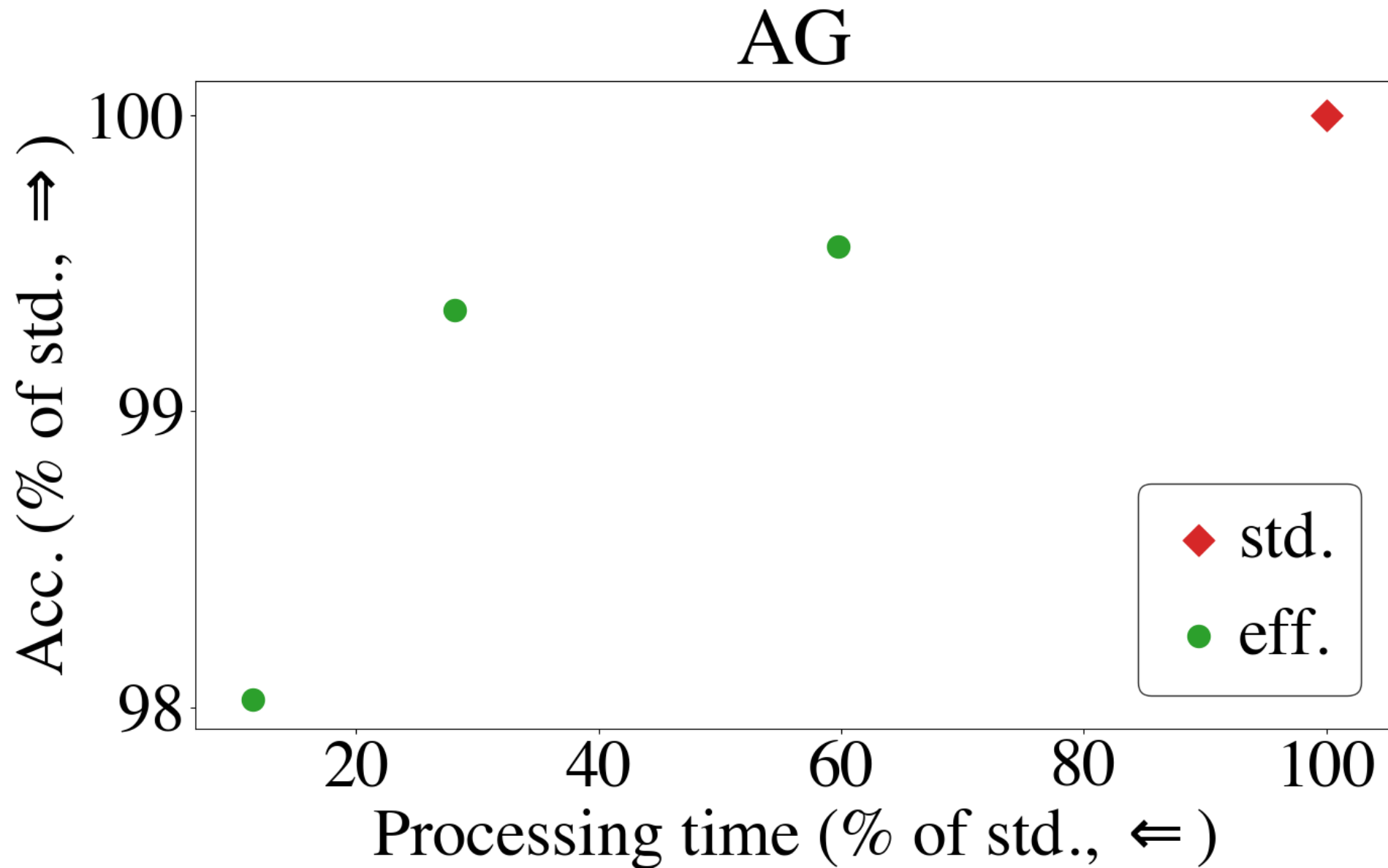


Efficient baselines



# Experimental Results:

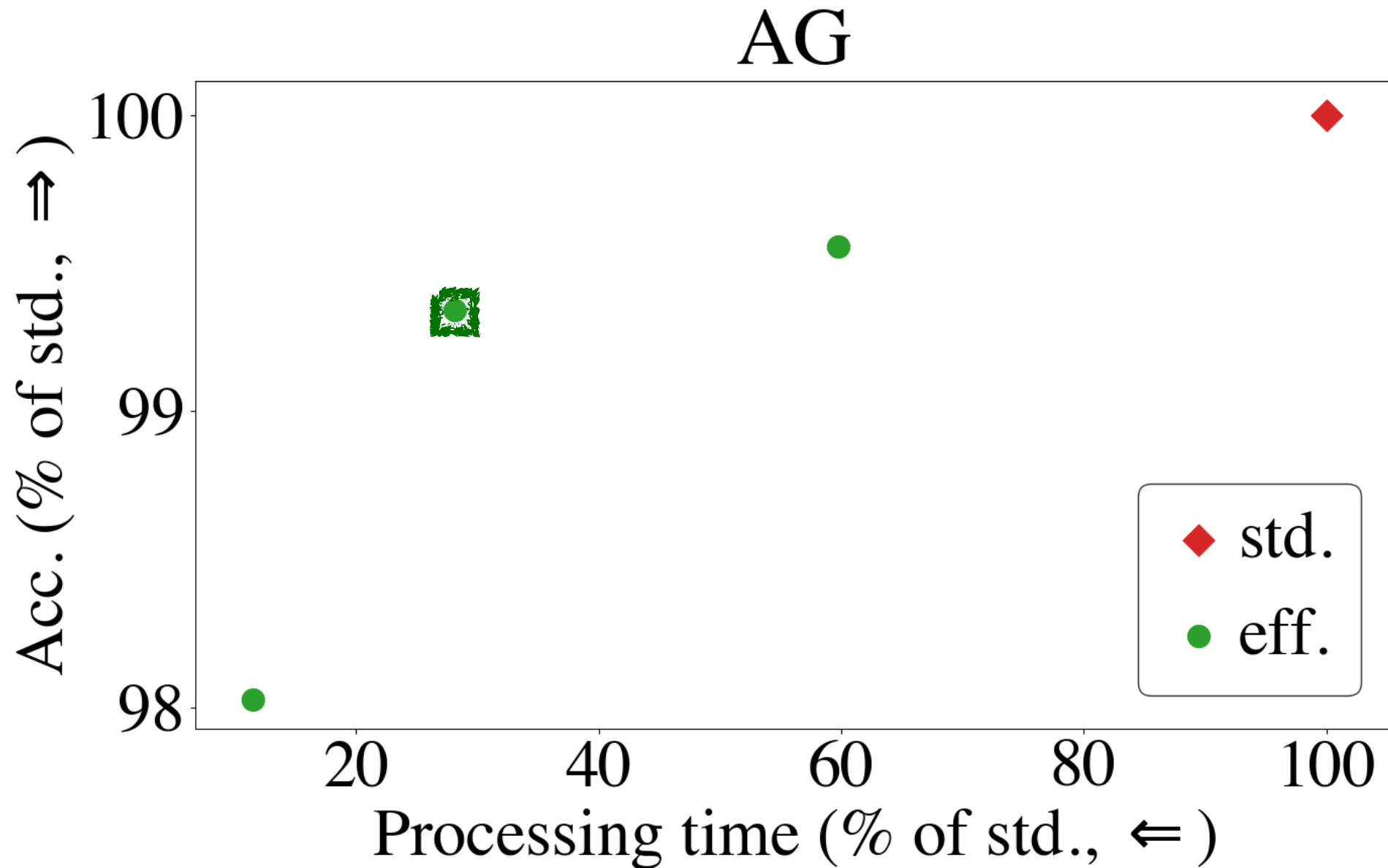
## Strong Baselines!





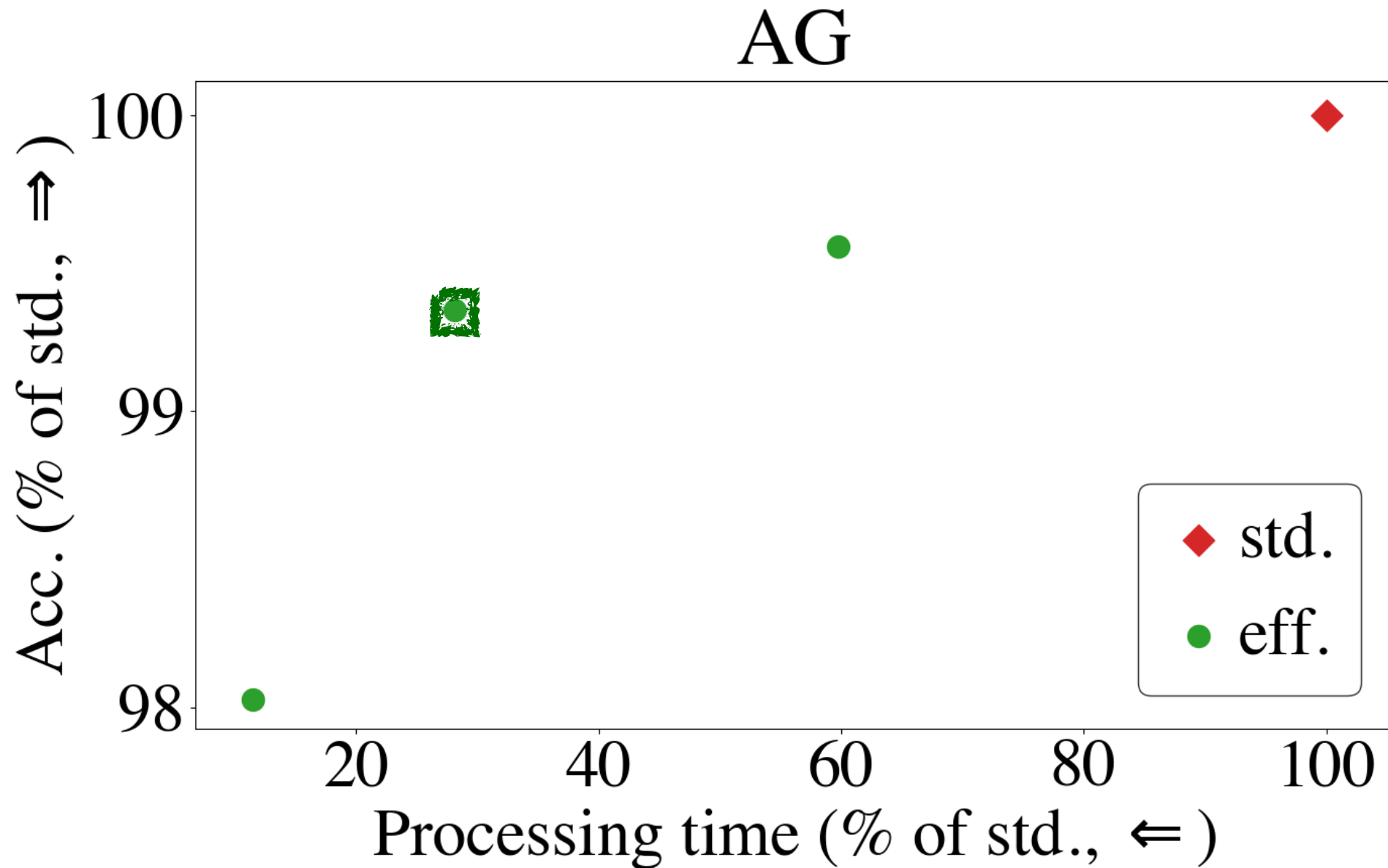
# Experimental Results:

## Strong Baselines!



# Experimental Results:

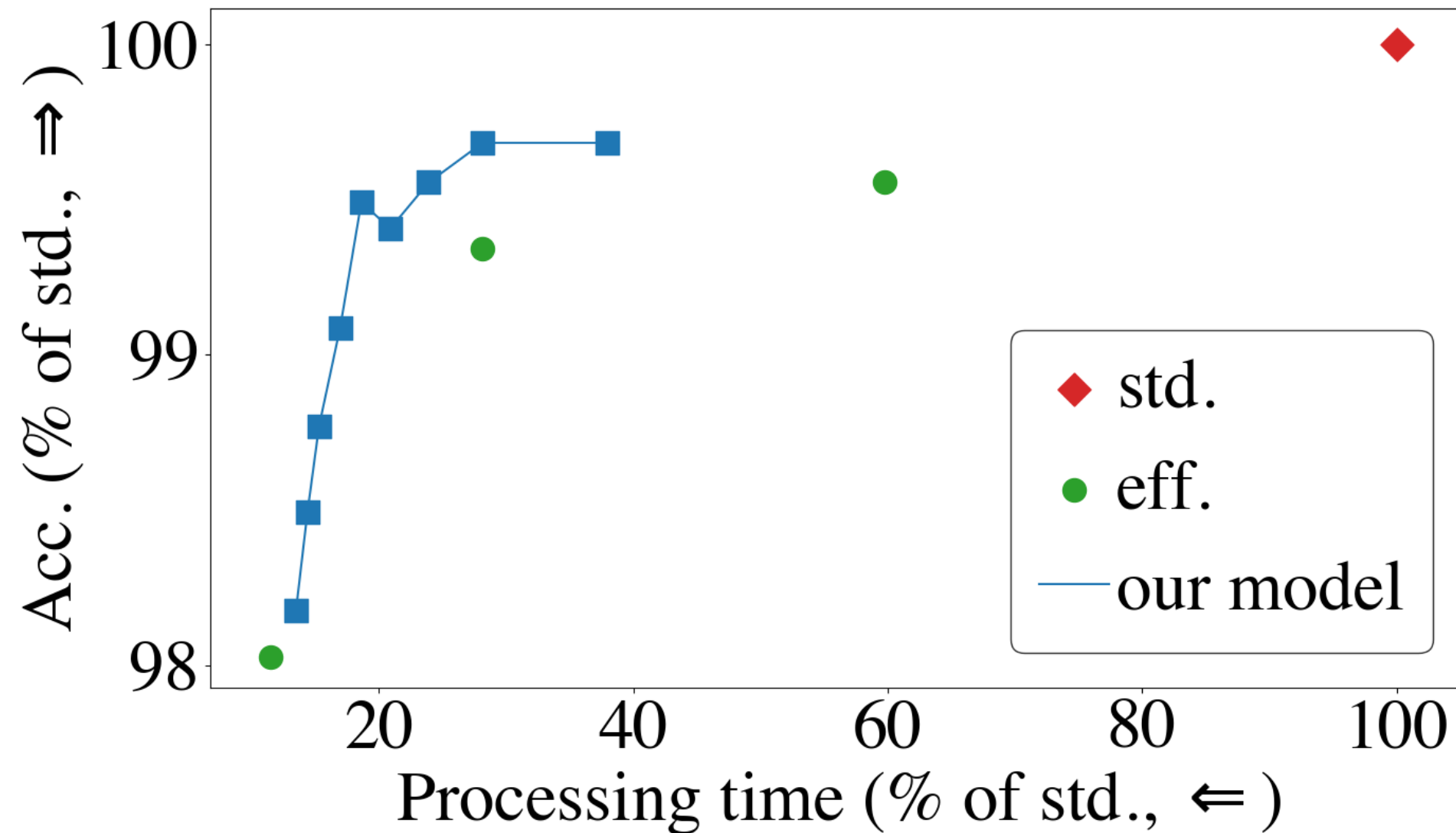
## Strong Baselines!



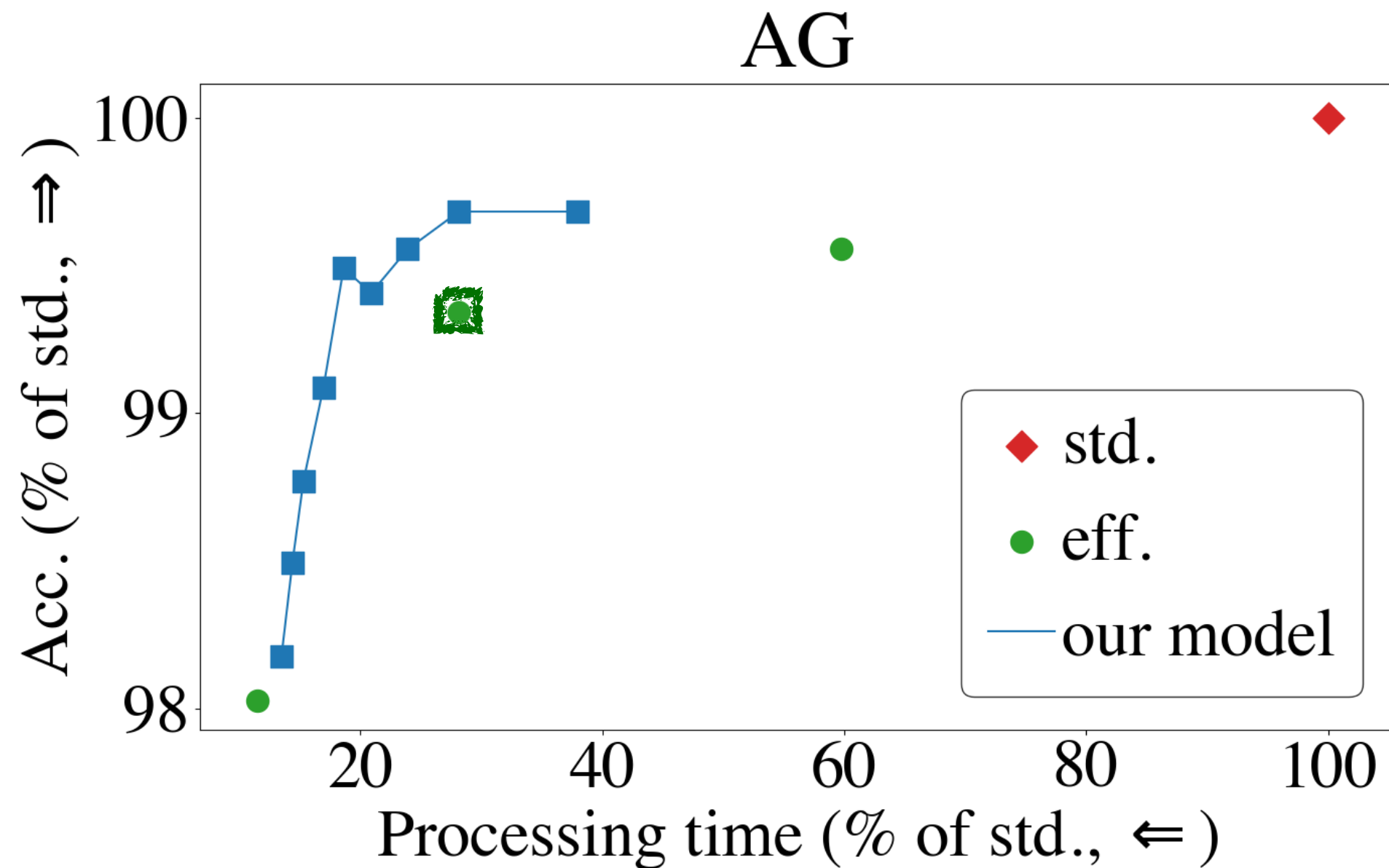
**3** times faster, within 1% of full model

# Better Speed/Accuracy Tradeoff

AG

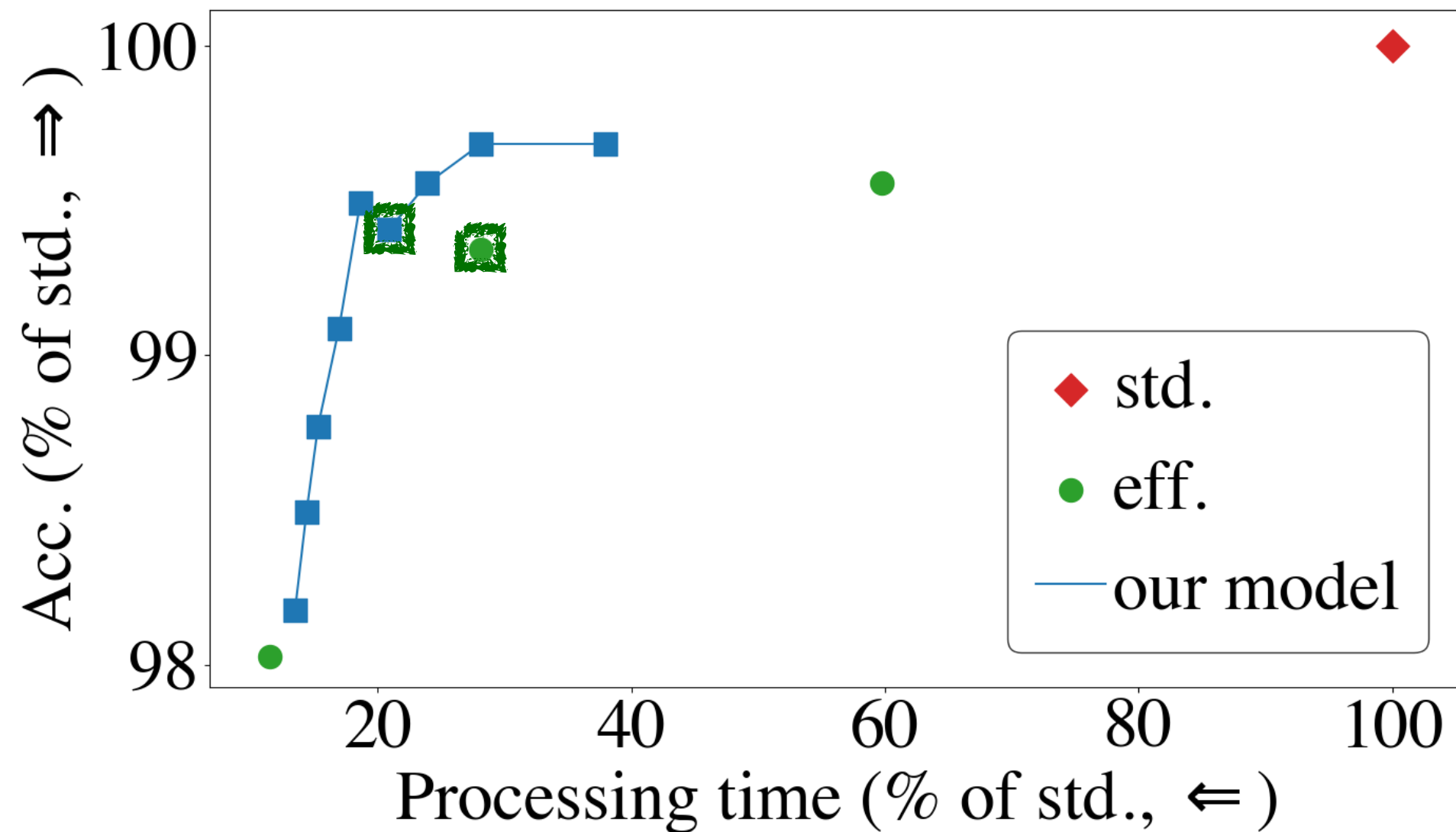


# Better Speed/Accuracy Tradeoff



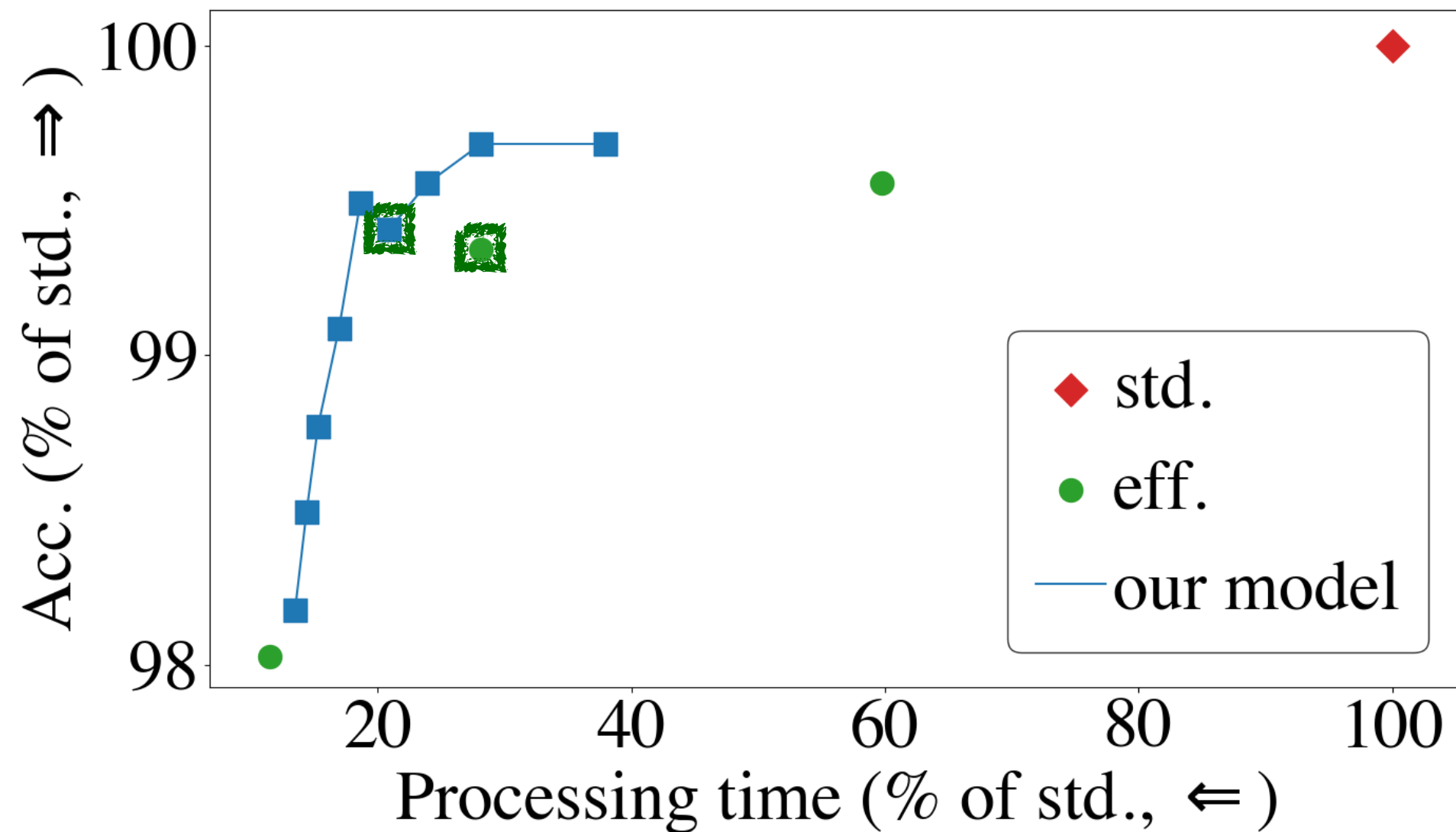
# Better Speed/Accuracy Tradeoff

AG



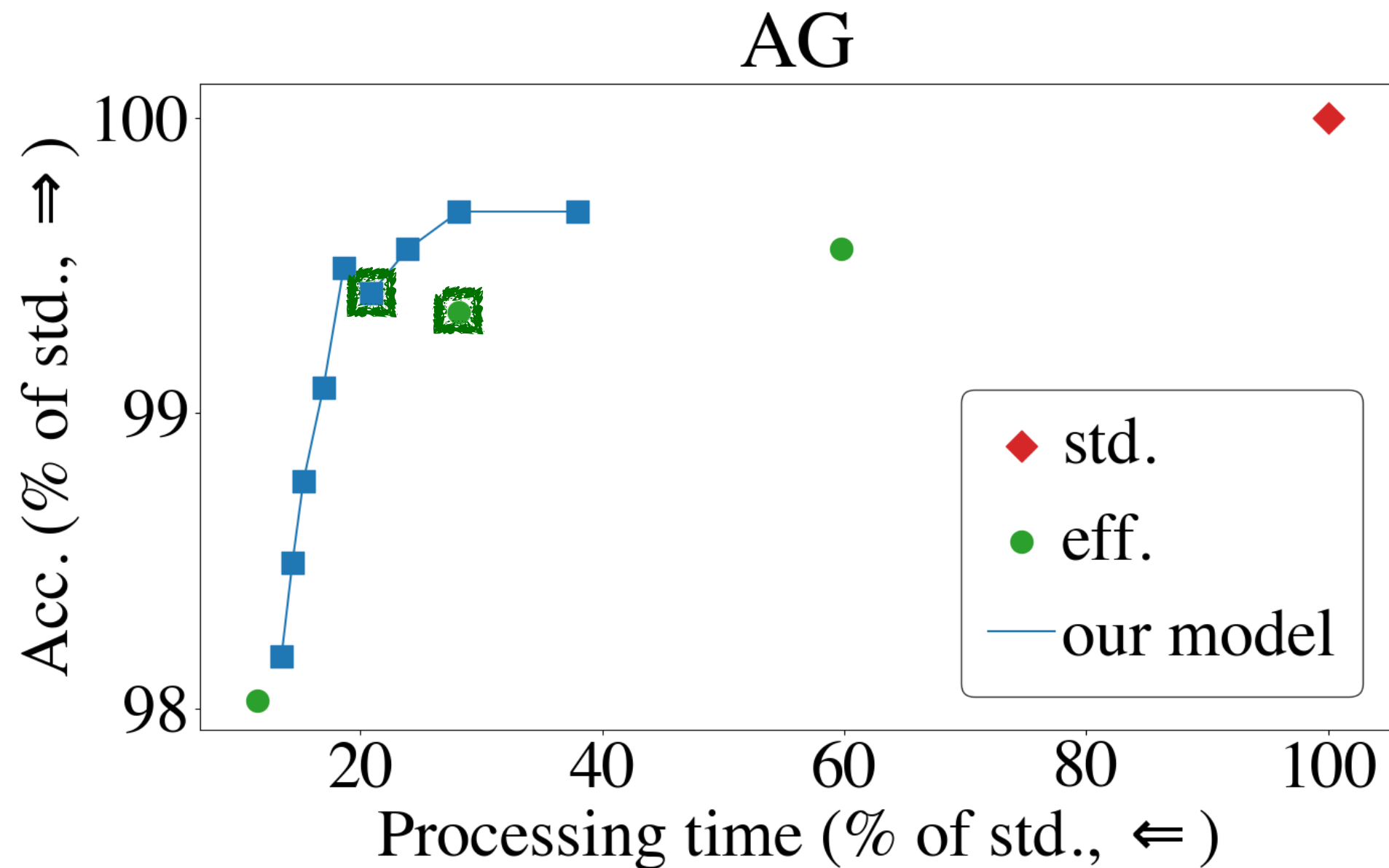
# Better Speed/Accuracy Tradeoff

AG



**5** times faster, within 1% of full model

# Better Speed/Accuracy Tradeoff



5 times faster, within 1% of full model

Similar results on  
SST, IMDB



# More about our Approach

- No effective growth in parameters
  - $< 0.005\%$  additional parameters

# More about our Approach

- No effective growth in parameters
  - $< 0.005\%$  additional parameters
- Training is **not** slower

# More about our Approach

- No effective growth in parameters
  - $< 0.005\%$  additional parameters
- Training is **not** slower
- A **single** trained model provides multiple options along the speed/accuracy tradeoff
  - A single parameter: confidence threshold

# More about our Approach

- No effective growth in parameters
  - $< 0.005\%$  additional parameters
- Training is **not** slower
- A **single** trained model provides multiple options along the speed/accuracy tradeoff
  - A single parameter: confidence threshold
- Caveat: requires batch size=1 during inference

# Case Study 2: Efficient Training

Swayamdipta, Schwartz et al., EMNLP 2020

*Some instances are **more valuable** for training than others*



# High-Level Idea

- *Divide the instances in a dataset into different groups*
- *Identify the groups that are **most valuable** for learning*
- *Train on those groups only, leading to **substantially faster training***

# Training Dynamics

- Assume a model trained for  $K$  epochs

# Training Dynamics

- Assume a model trained for  $K$  epochs
- At each epoch, the model makes predictions on each training sample
  - This leads to a vector of size  $K$  for each training instance

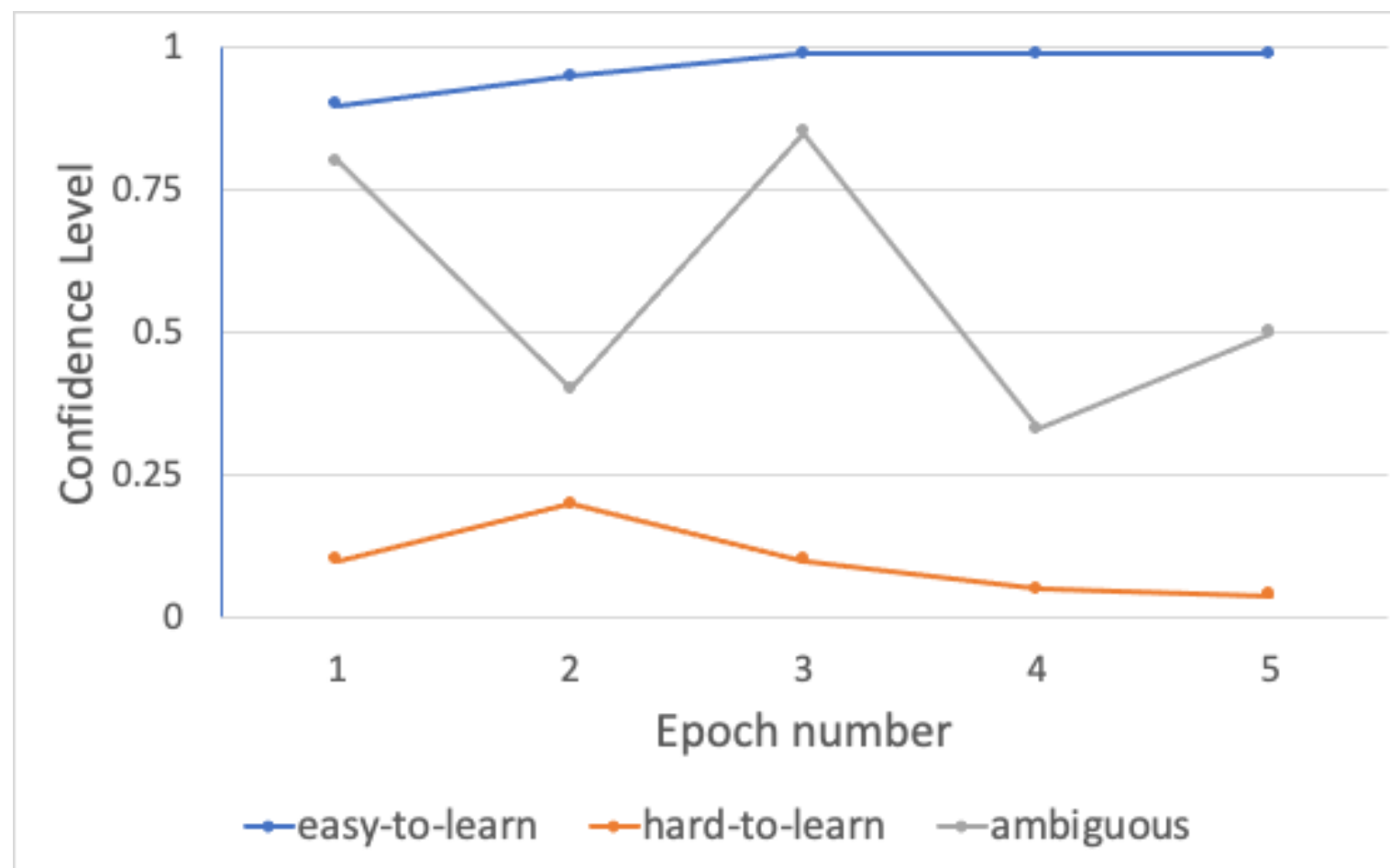


# Training Dynamics

- Assume a model trained for  $K$  epochs
- At each epoch, the model makes predictions on each training sample
  - This leads to a vector of size  $K$  for each training instance
- We compute two measures on each vector:
  - Mean
  - Variability

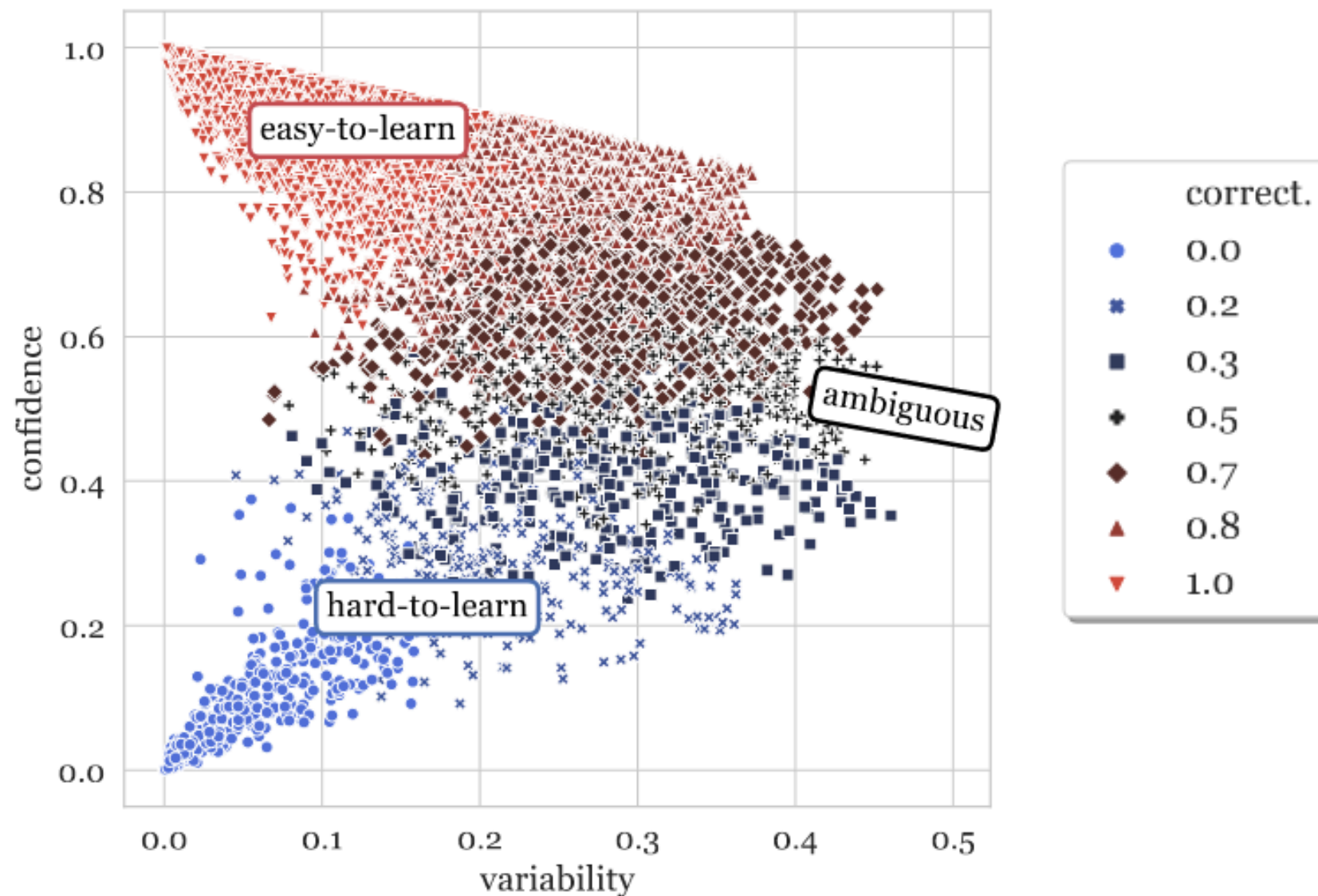
# Training Dynamics

## Toy Example



# Dataset Map Example

## SNLI, RoBERTa-Large



# Sample-Efficient Training

- *easy-to-learn* instances provide little value to training
- Can we use training dynamics to select the *most valuable* instances?

# Experiments

## WinoGrande, RoBERTa-Large

	WINO G. Val. (ID)
100% train	79.7 <sub>0.2</sub>
random	73.3 <sub>1.3</sub>
forgetting	75.5 <sub>1.3</sub>
AL-uncertainty	75.7 <sub>0.8</sub>
AL-greedyK	74.2 <sub>0.4</sub>
AFLite	76.8 <sub>0.8</sub>
<i>ambiguous</i>	<b>78.7</b> <sub>0.4</sub>

# Experiments

## WinoGrande, RoBERTa-Large

	WINOGr. Val. (ID)
100% train	79.7 <sub>0.2</sub>
random	73.3 <sub>1.3</sub>
forgetting	75.5 <sub>1.3</sub>
AL-uncertainty	75.7 <sub>0.8</sub>
AL-greedyK	74.2 <sub>0.4</sub>
AFLite	76.8 <sub>0.8</sub>
<i>ambiguous</i>	<b>78.7</b> <sub>0.4</sub>

# Experiments

WinoGrande, RoBERTa-Large

	WINO G. Val. (ID)
100% train	79.7 <sub>0.2</sub>
random	73.3 <sub>1.3</sub>
forgetting	75.5 <sub>1.3</sub>
AL-uncertainty	75.7 <sub>0.8</sub>
AL-greedyK	74.2 <sub>0.4</sub>
AFLite	76.8 <sub>0.8</sub>
<i>ambiguous</i>	<b>78.7</b> <sub>0.4</sub>

33%

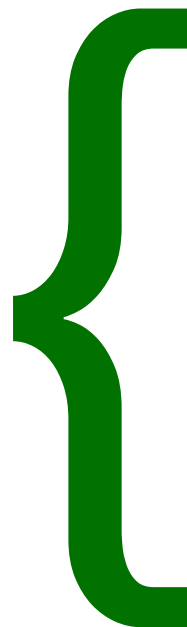


# Experiments

## WinoGrande, RoBERTa-Large

	WINOGr. Val. (ID)
100% train	79.7 <sub>0.2</sub>
random	73.3 <sub>1.3</sub>
forgetting	75.5 <sub>1.3</sub>
AL-uncertainty	75.7 <sub>0.8</sub>
AL-greedyK	74.2 <sub>0.4</sub>
AFLite	76.8 <sub>0.8</sub>
<i>ambiguous</i>	<b>78.7<sub>0.4</sub></b>

33%





# Experiments

WinoGrande, RoBERTa-Large

	WINOGr. Val. (ID)
100% train	79.7 <sub>0.2</sub>
random	73.3 <sub>1.3</sub>
forgetting	75.5 <sub>1.3</sub>
AL-uncertainty	75.7 <sub>0.8</sub>
AL-greedyK	74.2 <sub>0.4</sub>
AFLite	76.8 <sub>0.8</sub>
<i>ambiguous</i>	<b>78.7<sub>0.4</sub></b>

33%

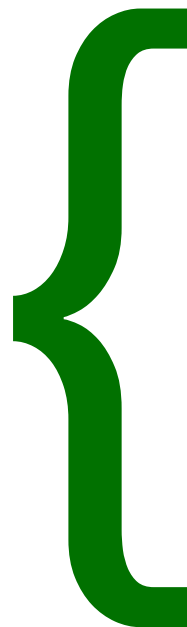


# Experiments

## WinoGrande, RoBERTa-Large

	WINOGr. Val. (ID)
100% train	79.7 <sub>0.2</sub>
random	73.3 <sub>1.3</sub>
forgetting	75.5 <sub>1.3</sub>
AL-uncertainty	75.7 <sub>0.8</sub>
AL-greedyK	74.2 <sub>0.4</sub>
AFLite	76.8 <sub>0.8</sub>
<i>ambiguous</i>	<b>78.7<sub>0.4</sub></b>

33%



# Experiments

## WinoGrande, RoBERTa-Large

	WINOGr. Val. (ID)	WSC (OOD)
100% train	79.7 <sub>0.2</sub>	86.0 <sub>0.1</sub>
random	73.3 <sub>1.3</sub>	85.6 <sub>0.4</sub>
forgetting	75.5 <sub>1.3</sub>	84.8 <sub>0.7</sub>
AL-uncertainty	75.7 <sub>0.8</sub>	85.7 <sub>0.8</sub>
AL-greedyK	74.2 <sub>0.4</sub>	86.5 <sub>0.5</sub>
AFLite	76.8 <sub>0.8</sub>	86.6 <sub>0.6</sub>
<i>ambiguous</i>	<b>78.7<sub>0.4</sub></b>	<b>87.6<sub>0.6</sub></b>

33%

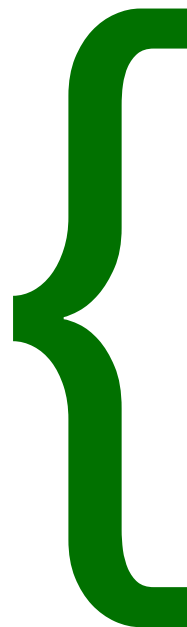


# Experiments

## WinoGrande, RoBERTa-Large

	WINOGr. Val. (ID)	WSC (OOD)
100% train	79.7 <sub>0.2</sub>	86.0 <sub>0.1</sub>
random	73.3 <sub>1.3</sub>	85.6 <sub>0.4</sub>
forgetting	75.5 <sub>1.3</sub>	84.8 <sub>0.7</sub>
AL-uncertainty	75.7 <sub>0.8</sub>	85.7 <sub>0.8</sub>
AL-greedyK	74.2 <sub>0.4</sub>	86.5 <sub>0.5</sub>
AFLite	76.8 <sub>0.8</sub>	86.6 <sub>0.6</sub>
<i>ambiguous</i>	78.7 <sub>0.4</sub>	87.6 <sub>0.6</sub>

33%



# Experiments

## WinoGrande, RoBERTa-Large

	WINOGr. Val. (ID)	WSC (OOD)
100% train	79.7 <sub>0.2</sub>	86.0 <sub>0.1</sub>
random	73.3 <sub>1.3</sub>	85.6 <sub>0.4</sub>
forgetting	75.5 <sub>1.3</sub>	84.8 <sub>0.7</sub>
AL-uncertainty	75.7 <sub>0.8</sub>	85.7 <sub>0.8</sub>
AL-greedyK	74.2 <sub>0.4</sub>	86.5 <sub>0.5</sub>
AFLite	76.8 <sub>0.8</sub>	86.6 <sub>0.6</sub>
<i>ambiguous</i>	78.7 <sub>0.4</sub>	87.6 <sub>0.6</sub>

33% {

Similar results on  
SNLI, MNLI, QNLI

# Recap

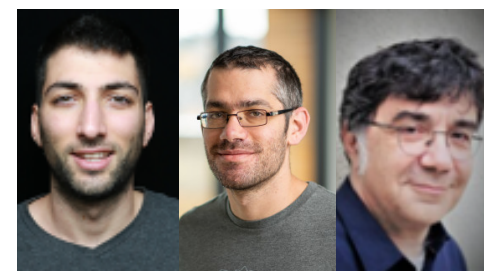
- Some instances contribute more to learning
- We select the ones with the *highest variance* in confidence level across training
- 3x reduction in training time
  - Minimal reduction in ID performance
  - **Improvement** on OOD performance
- Limitations
  - Model-dependent
  - Requires training on the full dataset first

# Case Study 3:

## Efficient Pre-training for Vision and Language

Bitton, Stanovsky, Elhadad & Schwartz, Findings of EMNLP 2021

*Some **words** are **more**  
**valuable** for pre-training than  
others*



# MLM in Vision and Language

- Virtually all vision pre-training works (Shin et al., 2021) follow BERT and randomly mask 15% of the tokens



# MLM in Vision and Language

- Virtually all vision pre-training works (Shin et al., 2021) follow BERT and randomly mask 15% of the tokens
- Of the masked tokens, roughly one half are *stop-words* or *punctuation*

# MLM in Vision and Language

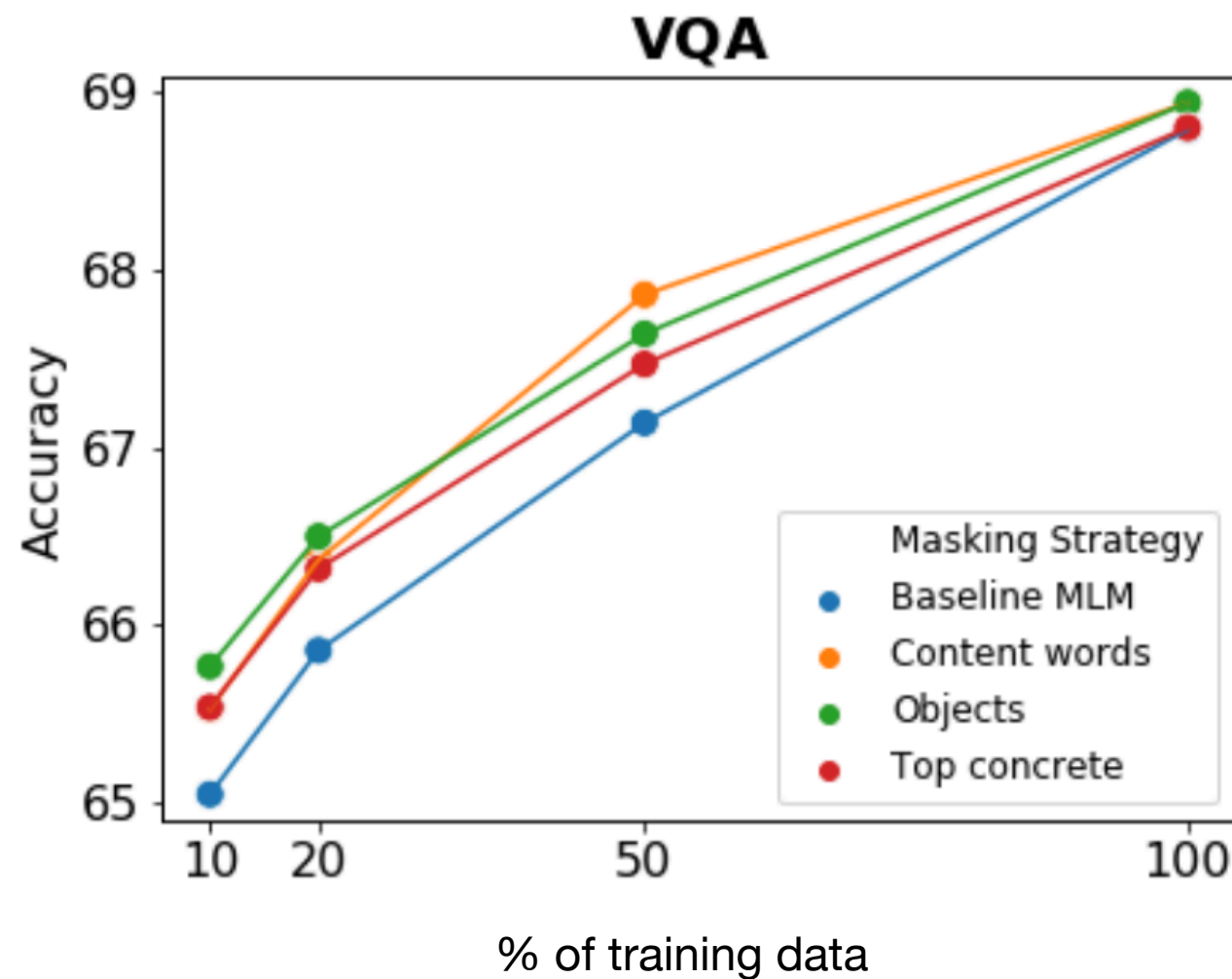
- Virtually all vision pre-training works (Shin et al., 2021) follow BERT and randomly mask 15% of the tokens
- Of the masked tokens, roughly one half are *stop-words* or *punctuation*
- We propose better masking strategies for V&L MLM

# MLM in Vision and Language

- Virtually all vision pre-training works (Shin et al., 2021) follow BERT and randomly mask 15% of the tokens
- Of the masked tokens, roughly one half are *stop-words* or *punctuation*
- We propose better masking strategies for V&L MLM
- *See Yonatan's talk for more details!*

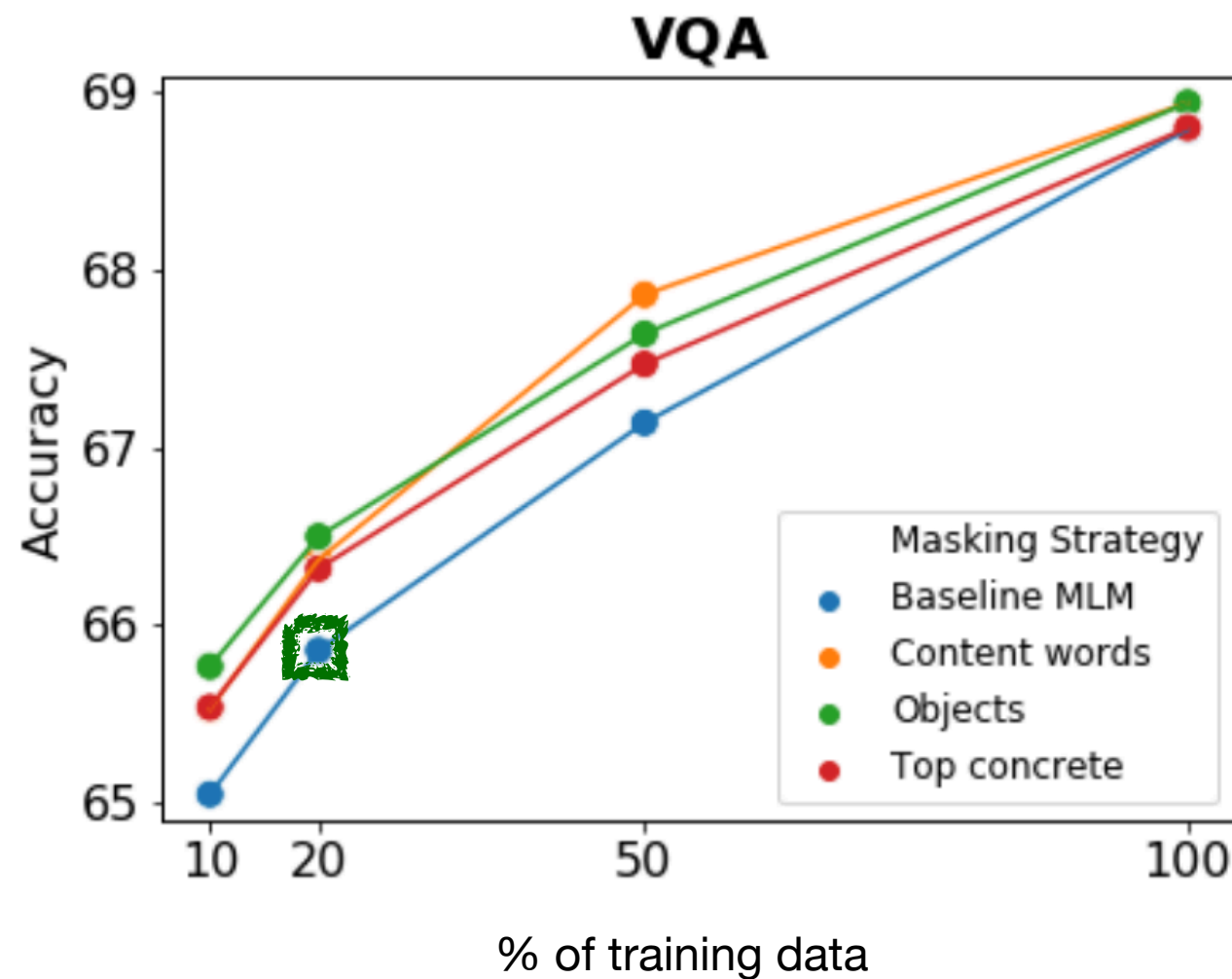
# Improved Downstream Performance

## Especially on Low-Resource Settings



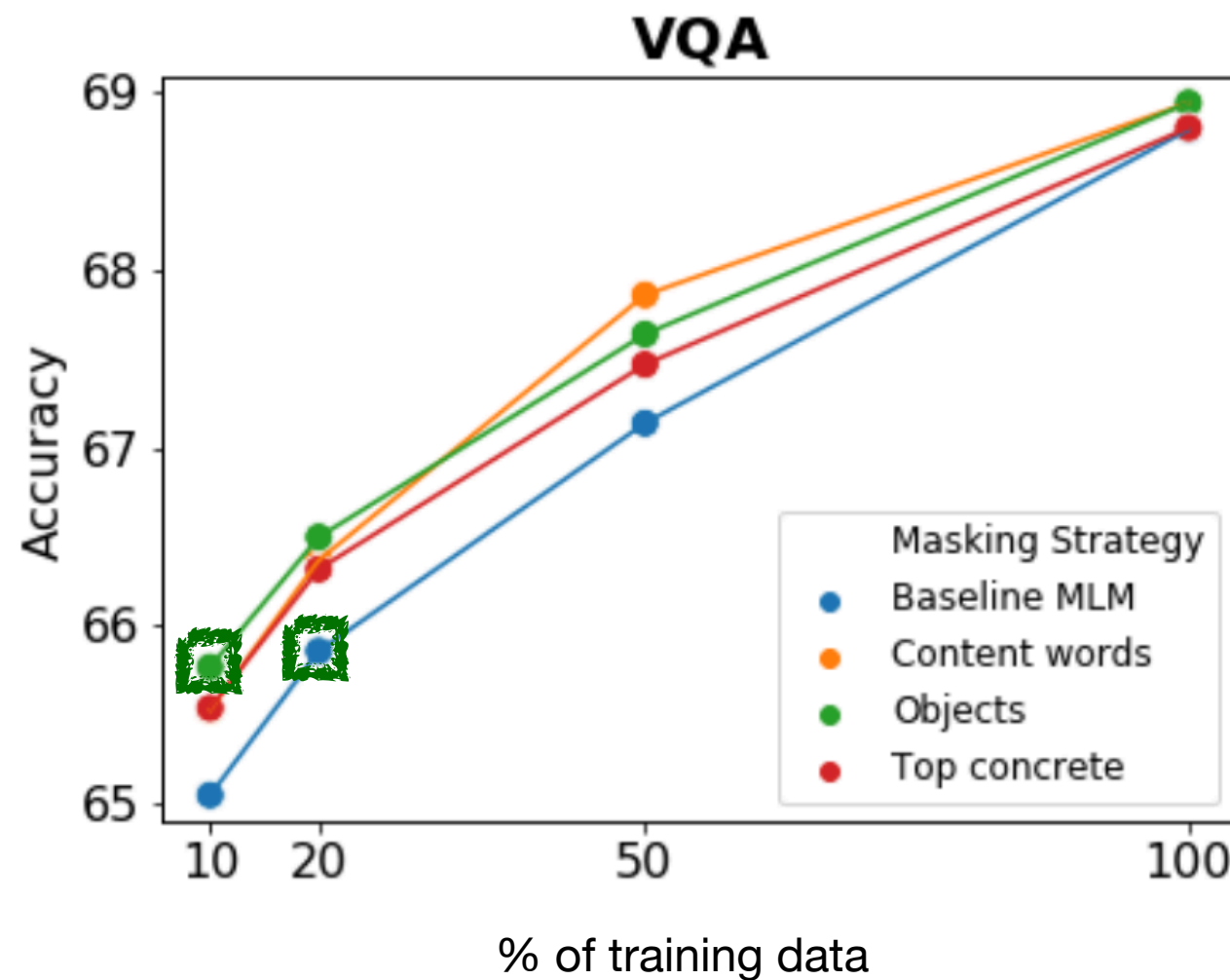
# Improved Downstream Performance

## Especially on Low-Resource Settings



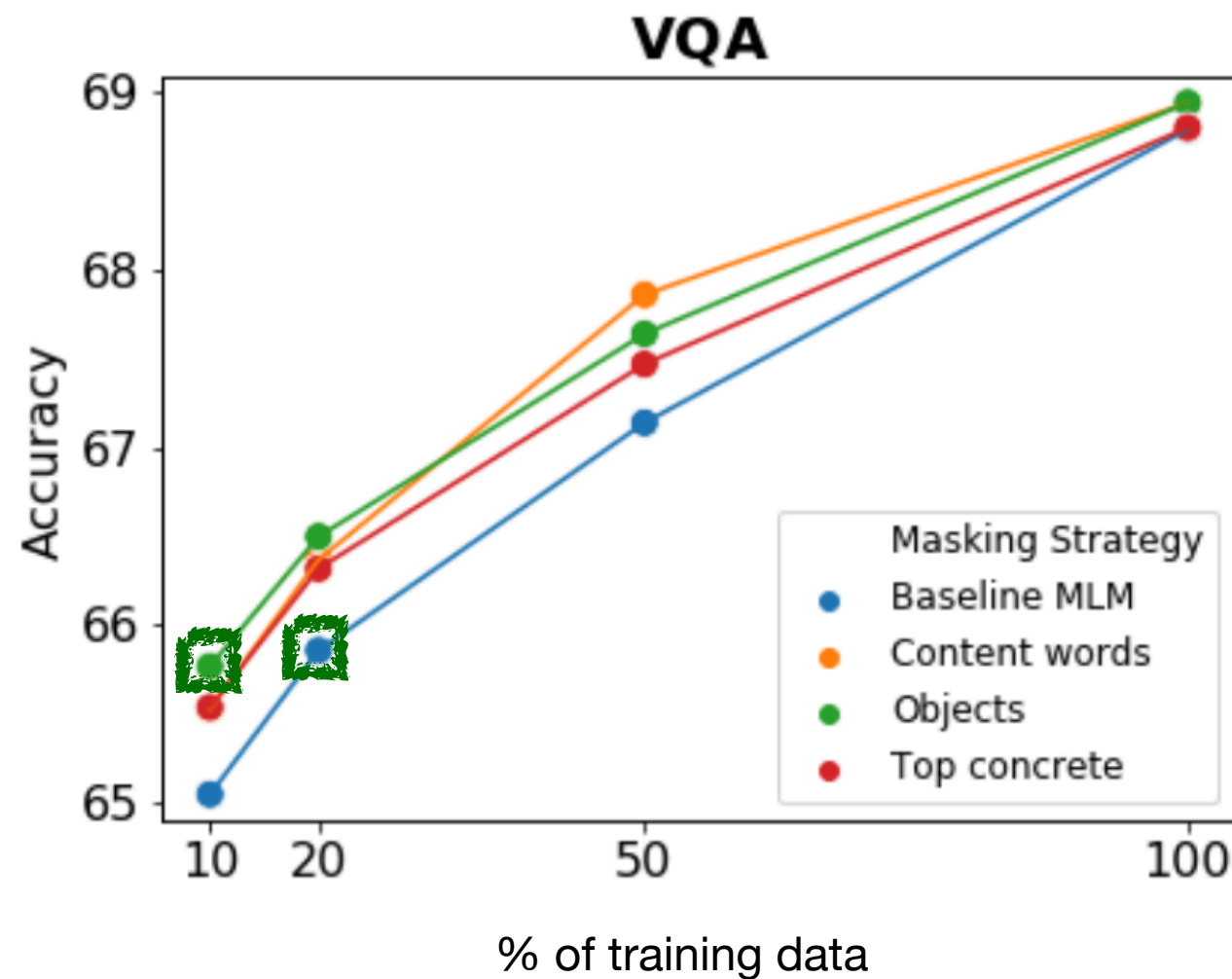
# Improved Downstream Performance

## Especially on Low-Resource Settings



# Improved Downstream Performance

## Especially on Low-Resource Settings

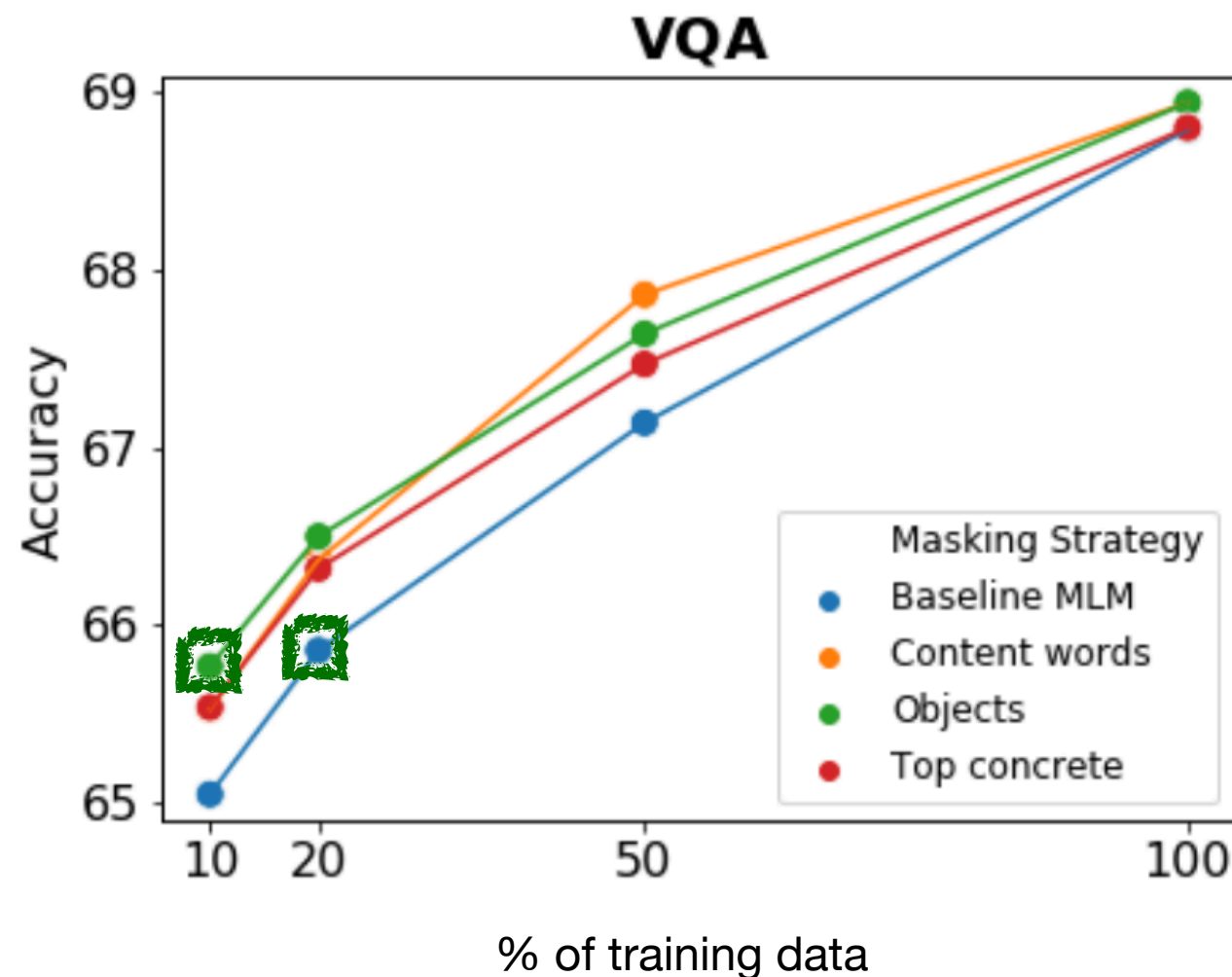


**Similar accuracy, twice as fast**

# Improved Downstream Performance

## Especially on Low-Resource Settings

Similar results on  
GQA, NLVR2



Similar accuracy, twice as fast



# Not all Instances are Alike

## Recap

- Efficient **inference** by selecting *the right tool for the job*
- Efficient **fine-tuning** by selecting the most *ambiguous* examples
- Efficient multi-modal **pre-training** by better *masking strategies*

# Amazing Collaborators!



# Not all instances are alike

## Recap

- Efficient **inference** by selecting *the right tool for the job*
- Efficient **fine-tuning** by selecting the most *ambiguous* examples
- Efficient multi-modal **pre-training** by better *masking strategies*

