

Noisy Channel Language Model Prompting for Few-shot Text Classification

Sewon Min, Mike Lewis, Hannaneh Hajishirzi, Luke Zettlemoyer
University of Washington, Meta AI, Allen Institute for AI



LM Prompting

Using a frozen language model for a downstream task

LM Prompting

Using a frozen language model for a downstream task

“Why are boolean values capitalized in Python?”

LM Prompting

Using a frozen language model for a downstream task

“Why are boolean values capitalized in Python?”

<options>

Society & Culture
Computer & Internet
Business & Finance
Entertainment & Music
Politics & Government

⋮

(Brown et al. 2020)

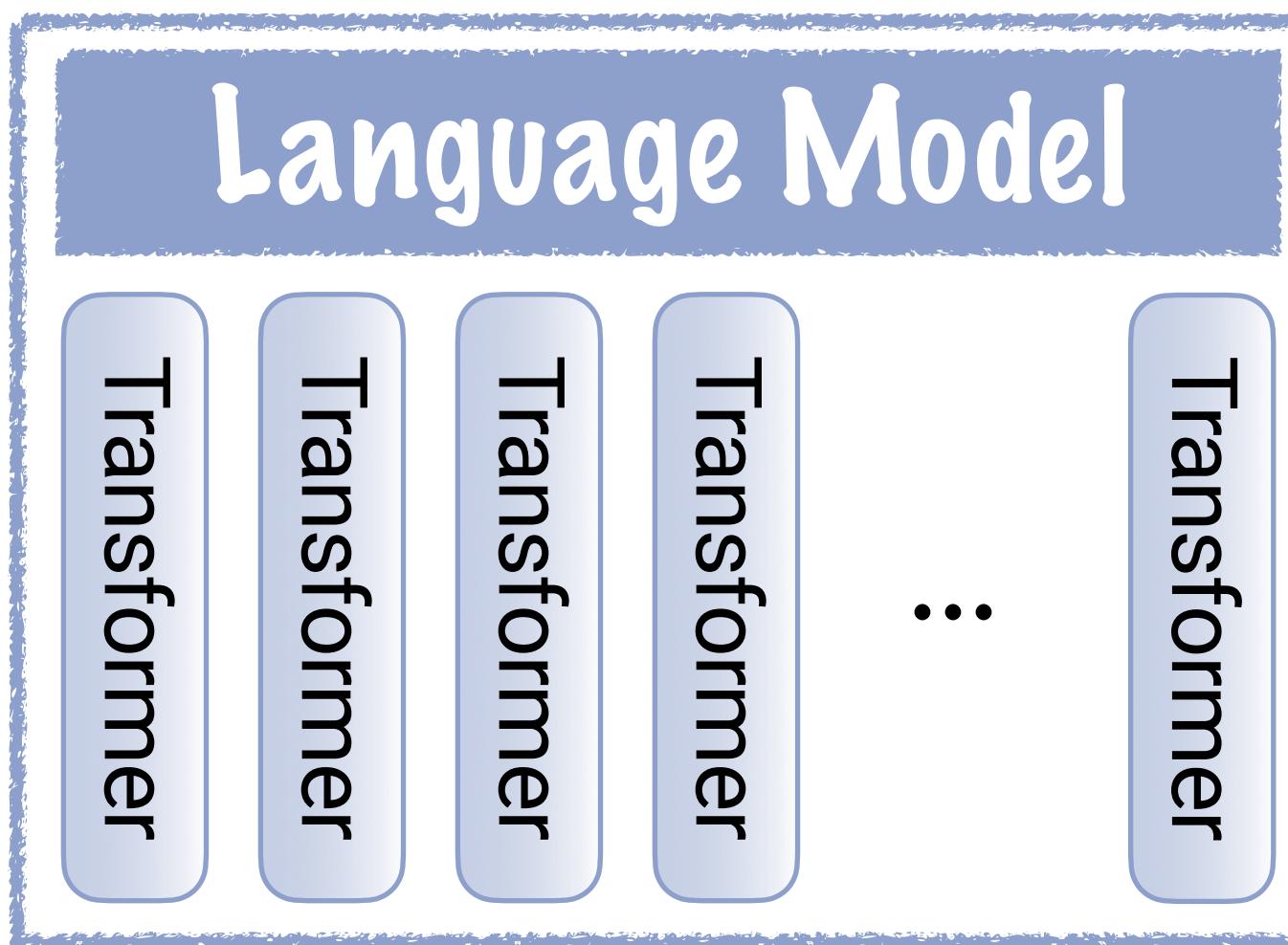
LM Prompting

Using a frozen language model for a downstream task

Why are boolean values capitalized in Python?

<options>

- Society & Culture
 - Computer & Internet
 - Business & Finance
 - Entertainment & Music
 - Politics & Government
- ⋮



(Brown et al. 2020)

LM Prompting

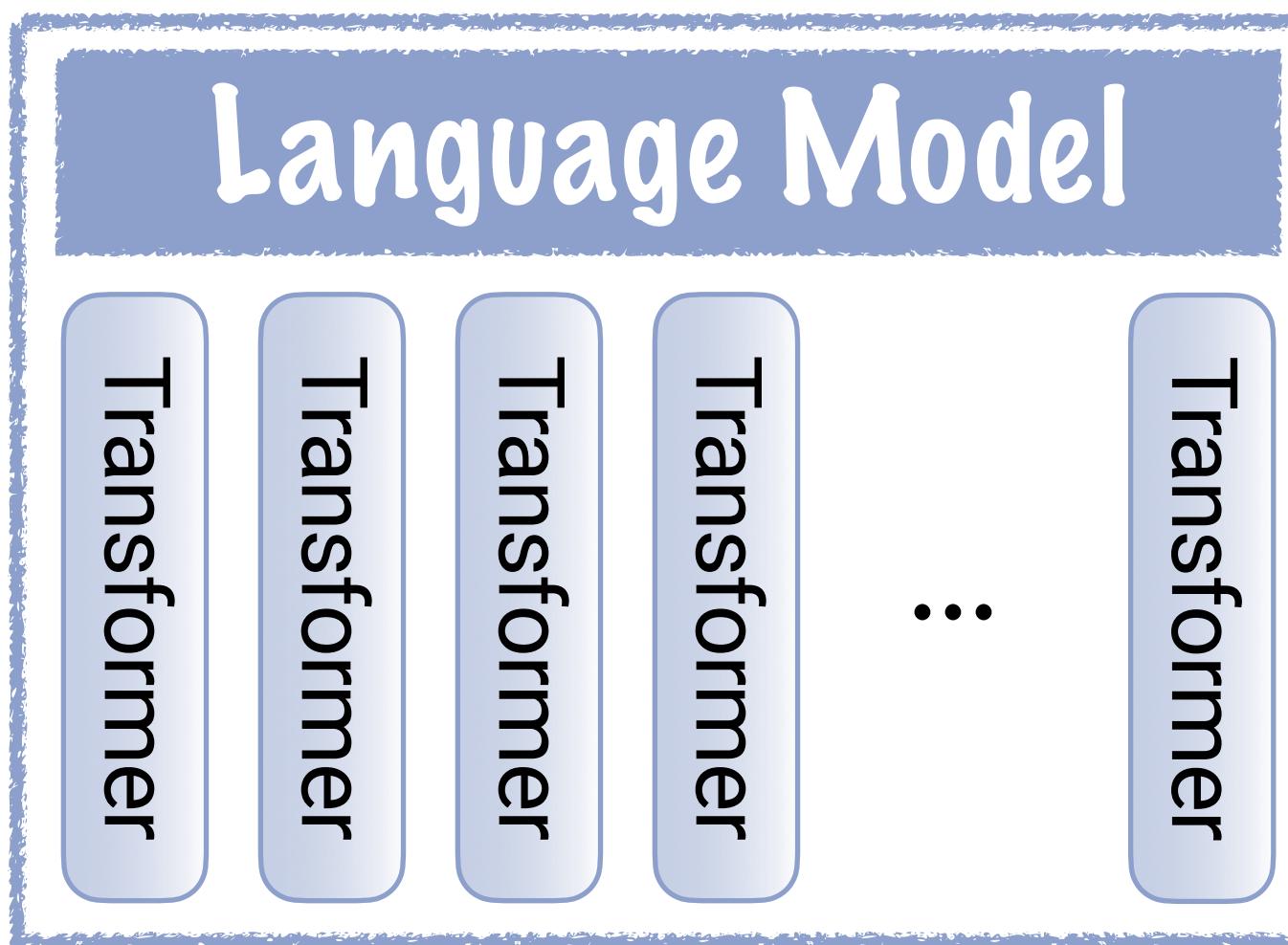
Using a frozen language model for a downstream task

Why are boolean values capitalized in Python?

<options>

- Society & Culture
- Computer & Internet
- Business & Finance
- Entertainment & Music
- Politics & Government

⋮



It is about Computer & Internet

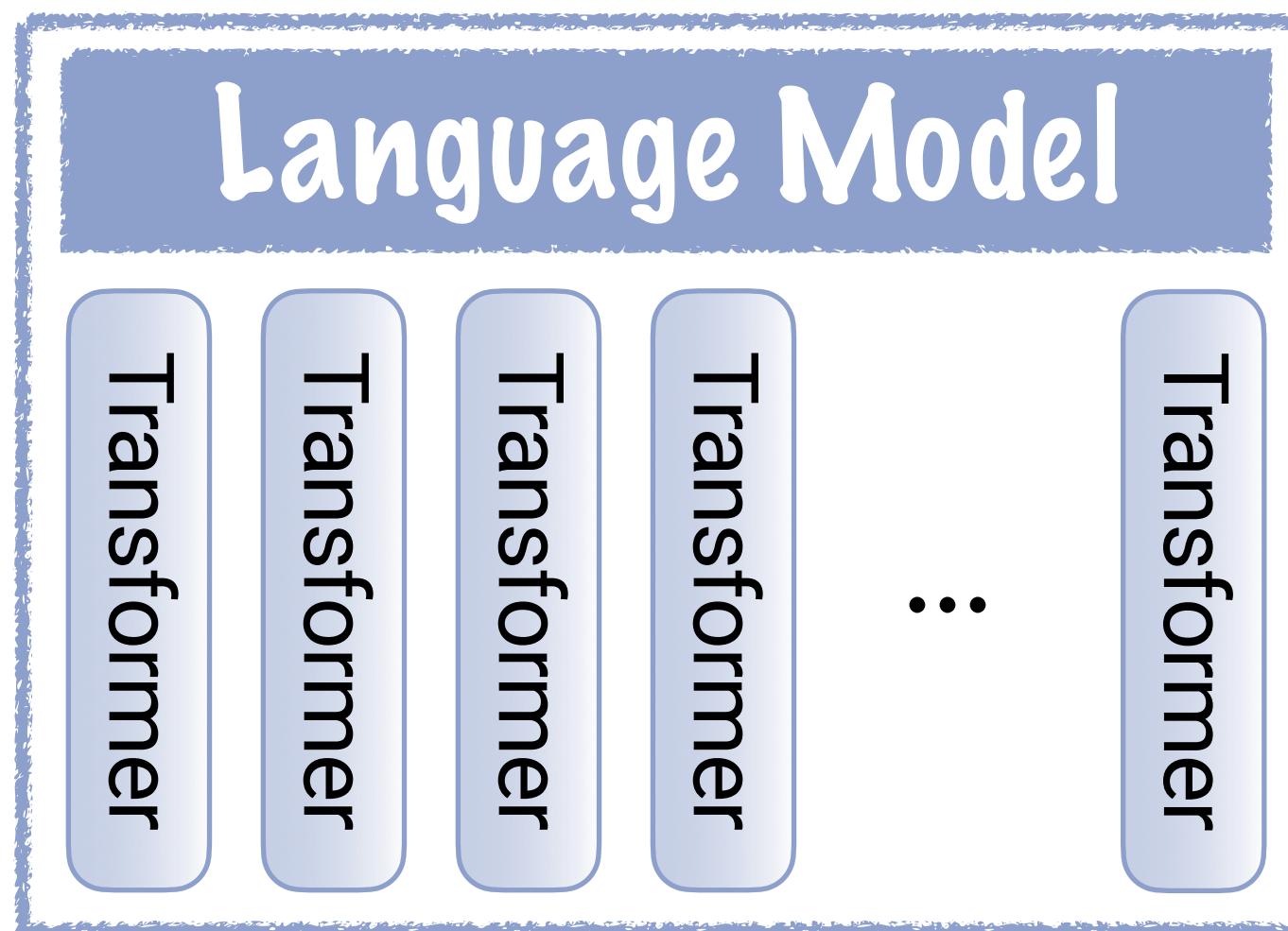
LM Prompting

Using a frozen language model for a downstream task

Why are boolean values capitalized in Python?

<options>

- Society & Culture
- Computer & Internet
- Business & Finance
- Entertainment & Music
- Politics & Government
- ⋮



It is about Computer & Internet

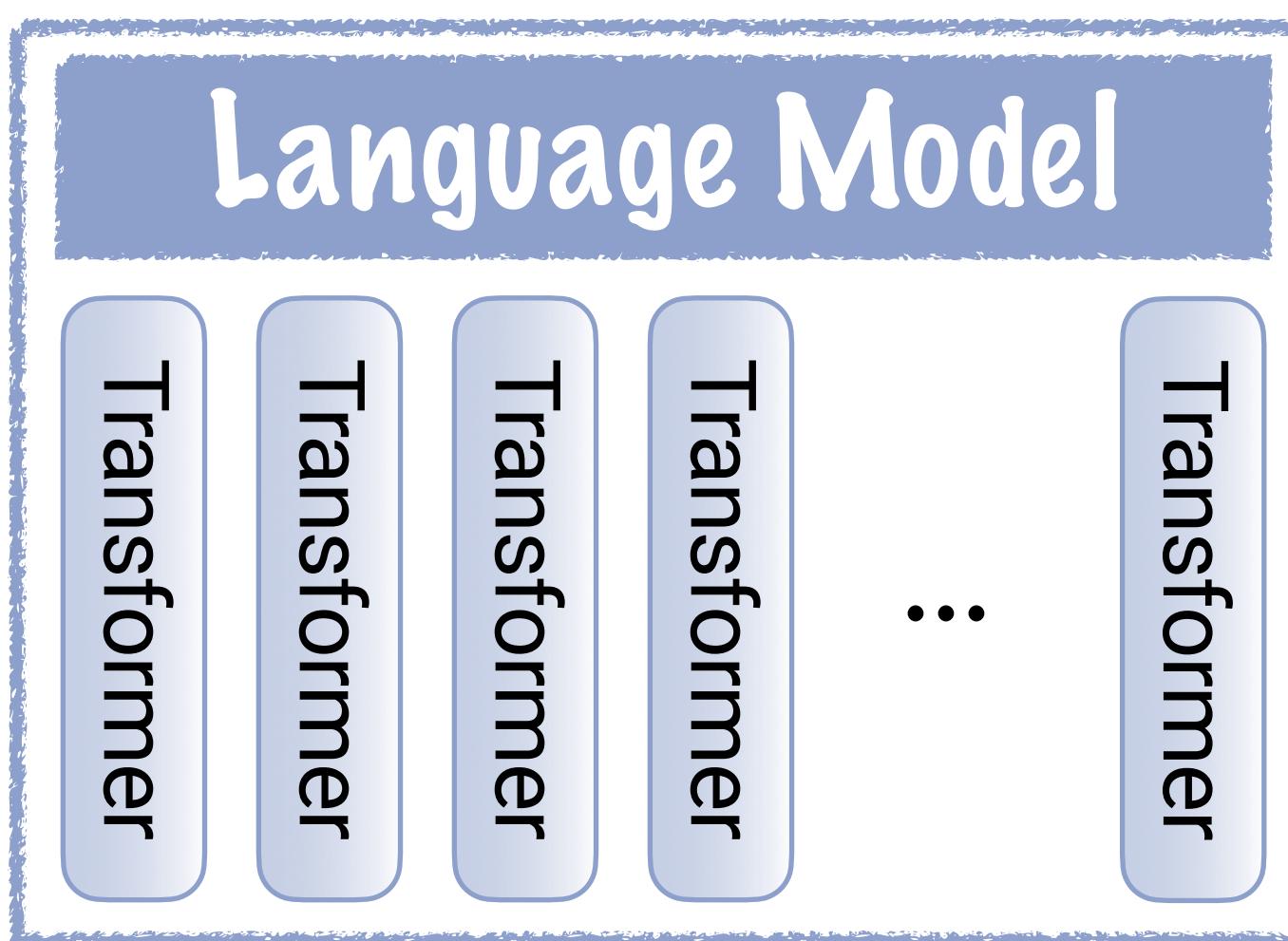
LM Prompting

Using a frozen language model for a downstream task

Why are boolean values capitalized in Python?

<options>

- Society & Culture
- Computer & Internet
- Business & Finance
- Entertainment & Music
- Politics & Government
- ⋮



It is about Computer & Internet

∨

It is about Society & Culture



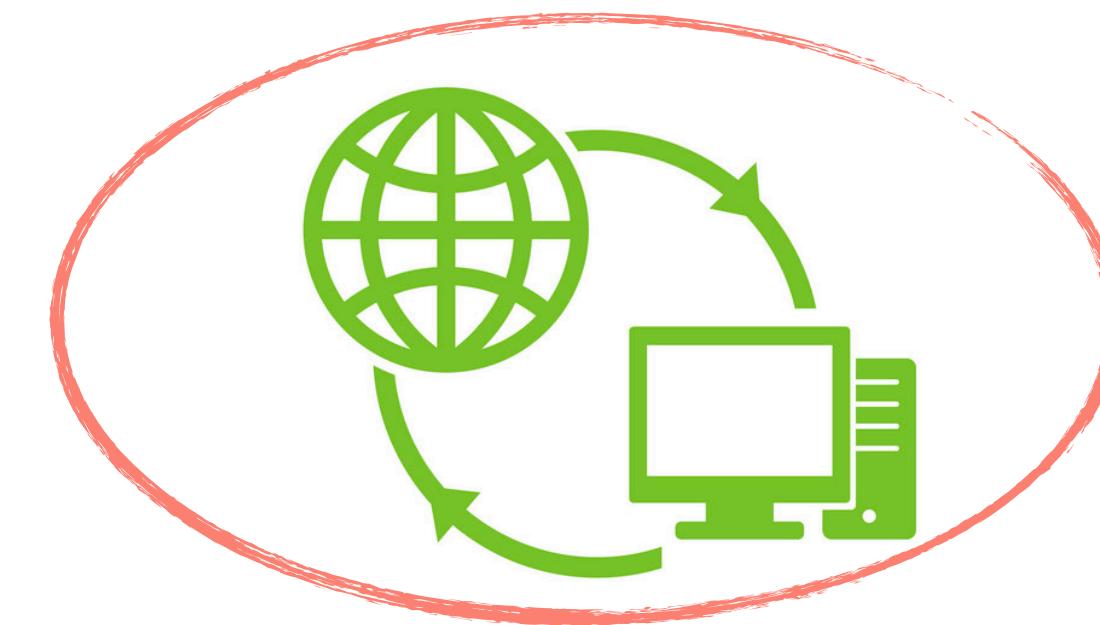
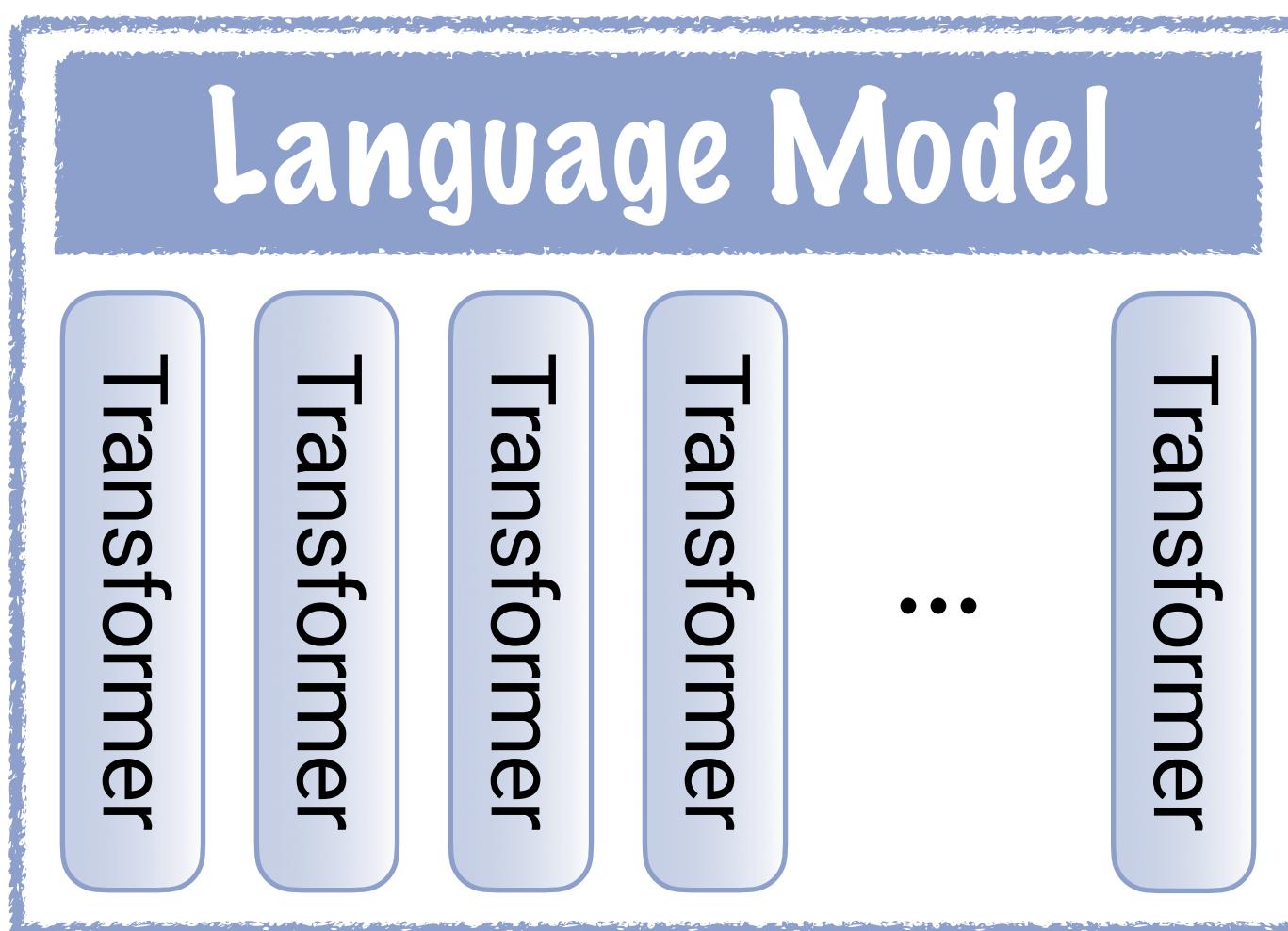
LM Prompting

Using a frozen language model for a downstream task

Why are boolean values capitalized in Python?

<options>

- Society & Culture
- Computer & Internet ✓
- Business & Finance
- Entertainment & Music
- Politics & Government
- ⋮



It is about Computer & Internet

∨

It is about Society & Culture



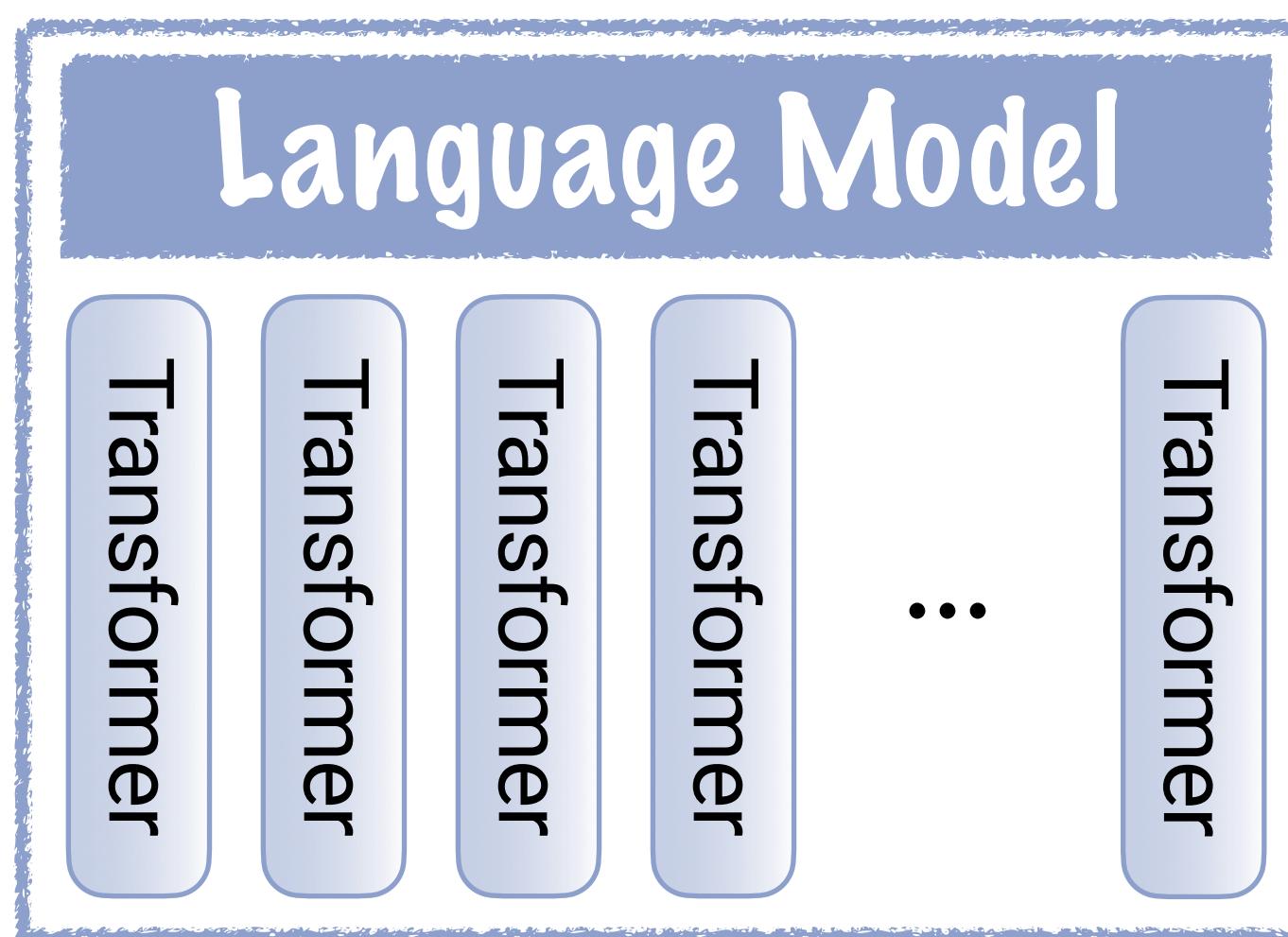
LM Prompting

Using a frozen language model for a downstream task

Why are boolean values capitalized in Python?

<options>

- Society & Culture
- Computer & Internet 
- Business & Finance
- Entertainment & Music
- Politics & Government
- ⋮



It is about Computer & Internet

∨

It is about Society & Culture



*More variations
(in-context learning, prompt tuning, etc..)*

LM Prompting



(Brown et al. 2020)



LM Prompting



(Brown et al. 2020)



Promising results with a 175B model

And by Zhao et al. 2021, Holtzman et al. 2021, Shin et al. 2020, Jiang et al. 2020, Gao et al. 2021, Li & Liang 2021, Lester et al. 2021, etc...

LM Prompting



(Brown et al. 2020)



Promising results with a 175B model

And by Zhao et al. 2021, Holtzman et al. 2021, Shin et al. 2020, Jiang et al. 2020, Gao et al. 2021, Li & Liang 2021, Lester et al. 2021, etc...

No or very limited number of trainable params

LM Prompting



(Brown et al. 2020)



Promising results with a 175B model

And by Zhao et al. 2021, Holtzman et al. 2021, Shin et al. 2020, Jiang et al. 2020, Gao et al. 2021, Li & Liang 2021, Lester et al. 2021, etc...

No or very limited number of trainable params

LM Prompting



(Brown et al. 2020)



Promising results with a 175B model

And by Zhao et al. 2021, Holtzman et al. 2021, Shin et al. 2020, Jiang et al. 2020, Gao et al. 2021, Li & Liang 2021, Lester et al. 2021, etc...

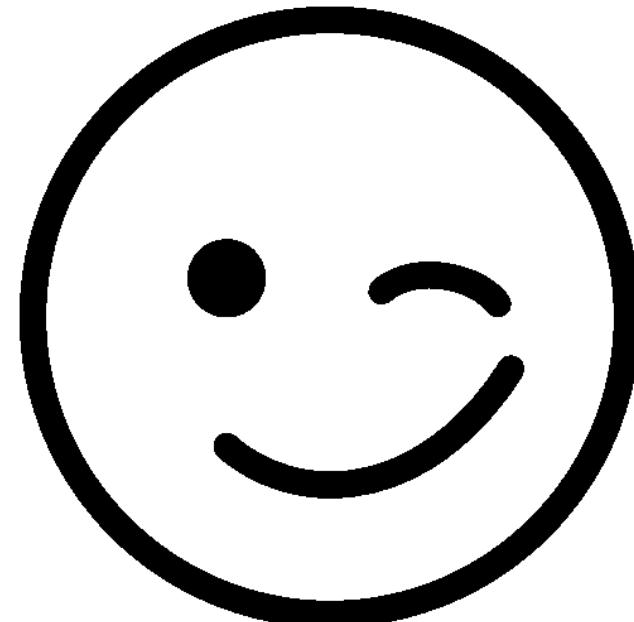
No or very limited number of trainable params

Instability – high variance

LM Prompting



(Brown et al. 2020)



Promising results with a 175B model

And by Zhao et al. 2021, Holtzman et al. 2021, Shin et al. 2020, Jiang et al. 2020, Gao et al. 2021, Li & Liang 2021, Lester et al. 2021, etc...

No or very limited number of trainable params

Instability – high variance

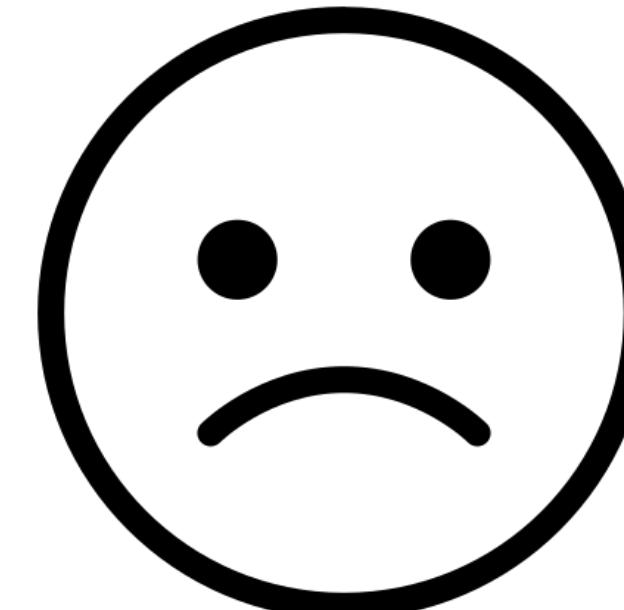
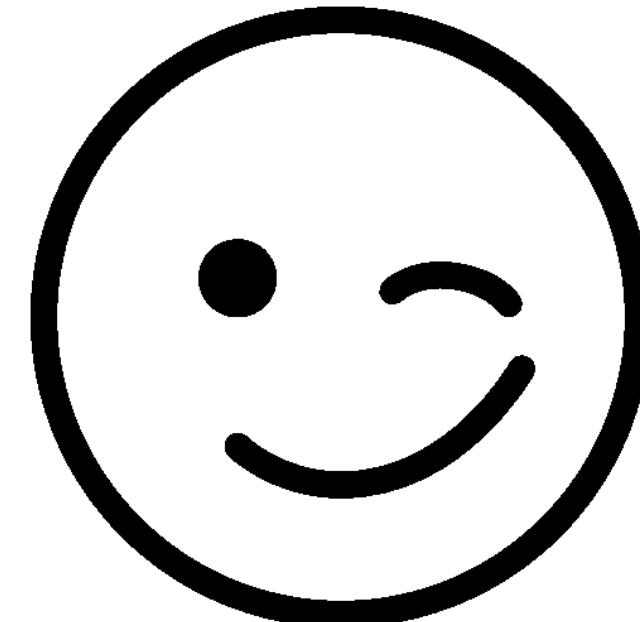
Low worst-case accuracy

Zhao et al. 2021, Perez et al. 2021, Lu et al. 2021

LM Prompting



(Brown et al. 2020)



Promising results with a 175B model

And by Zhao et al. 2021, Holtzman et al. 2021, Shin et al. 2020, Jiang et al. 2020, Gao et al. 2021, Li & Liang 2021, Lester et al. 2021, etc...

No or very limited number of trainable params

Instability – high variance

Low worst-case accuracy

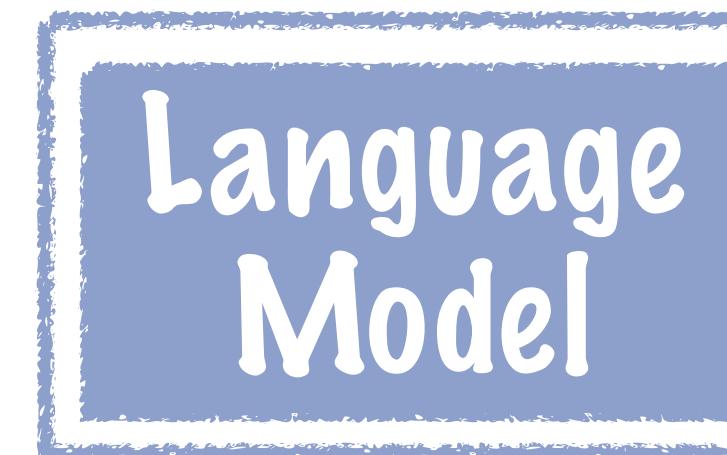
Zhao et al. 2021, Perez et al. 2021, Lu et al. 2021

How can we do better?

Noisy Channel Approach

Why are boolean values capitalized in Python?

Input

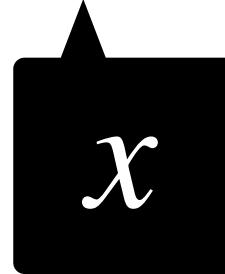


Output

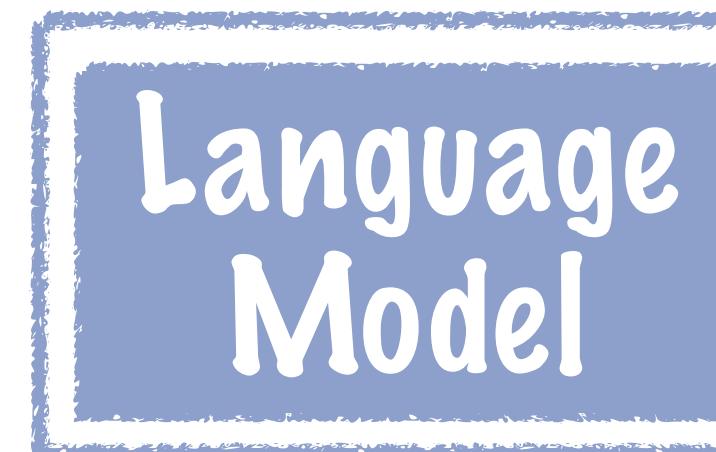
It is about Computer & Internet.

Noisy Channel Approach

Why are boolean values capitalized in Python?



Input



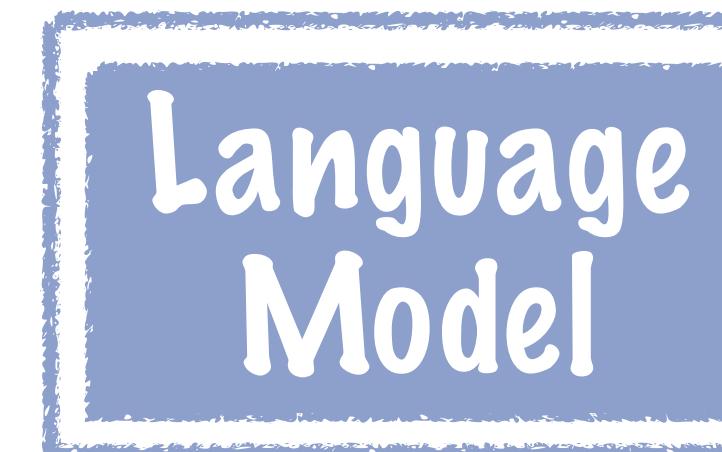
Output

It is about Computer & Internet.

Noisy Channel Approach

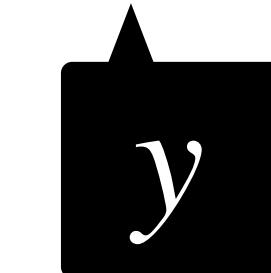
Why are boolean values capitalized in Python?

Input

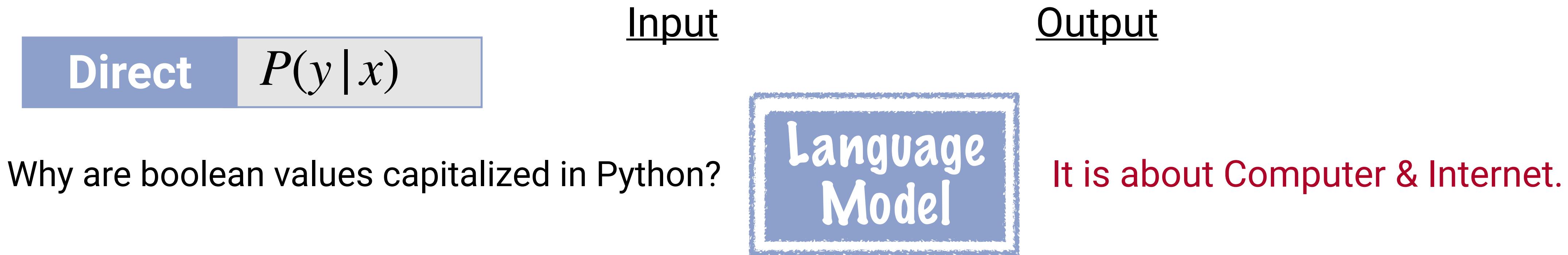


Output

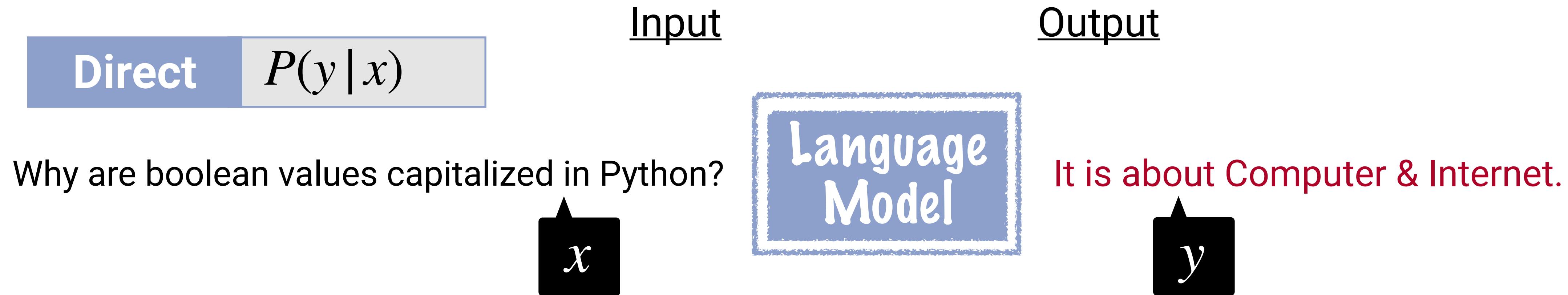
It is about Computer & Internet.



Noisy Channel Approach



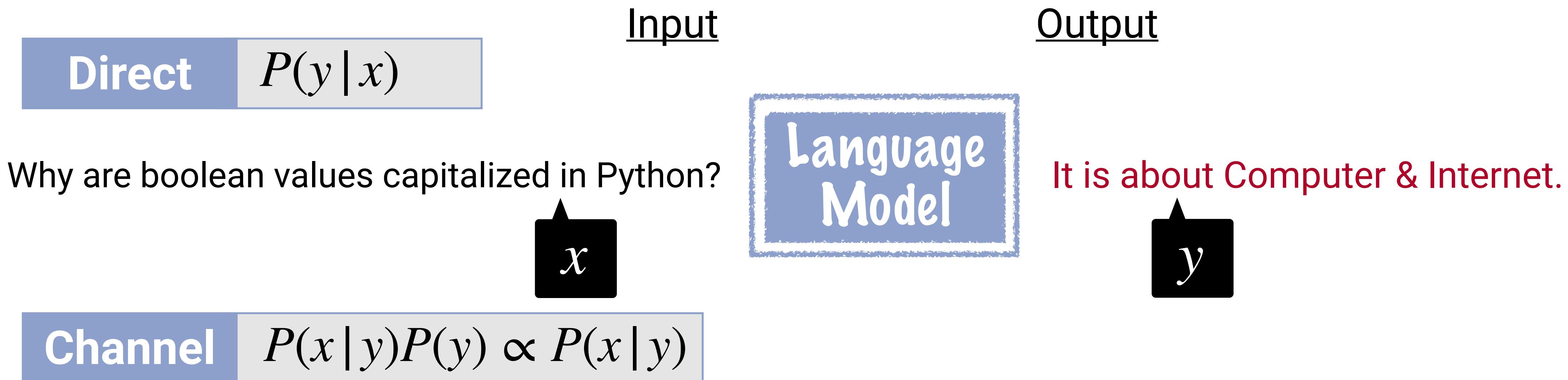
Noisy Channel Approach



Noisy channel model:

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} \propto P(x|y)P(y)$$

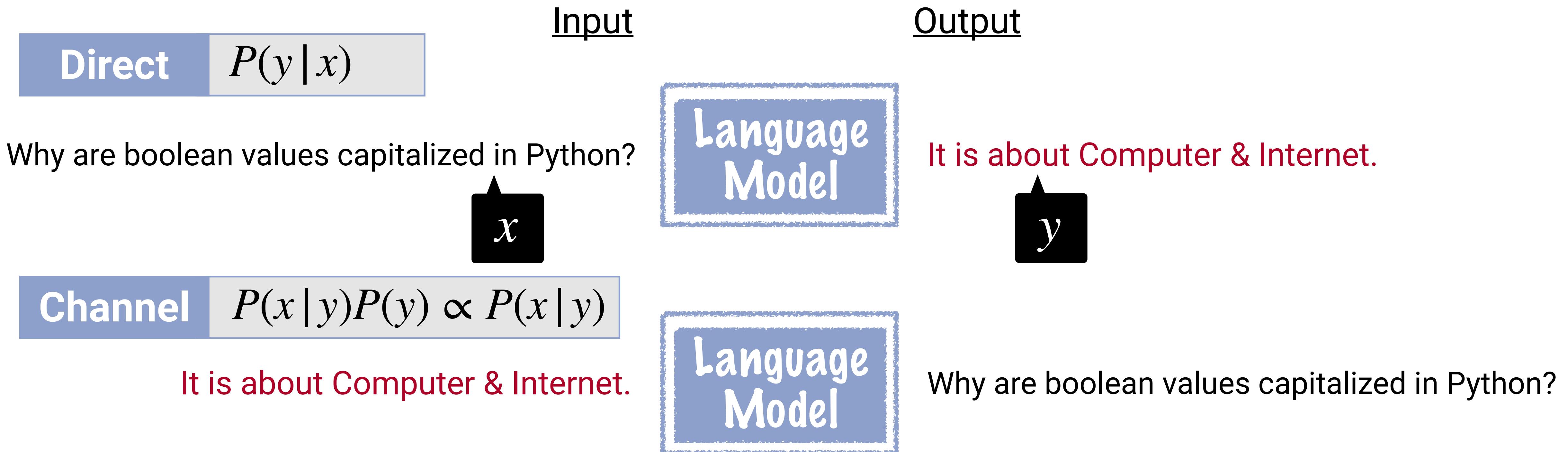
Noisy Channel Approach



Noisy channel model:

$$P(y | x) = \frac{P(x | y)P(y)}{P(x)} \propto P(x | y)P(y)$$

Noisy Channel Approach



Noisy channel model:

$$P(y | x) = \frac{P(x | y)P(y)}{P(x)} \propto P(x | y)P(y)$$

Method

x : Test input

y : Test output

$(x_1, y_1), \dots, (x_k, y_k)$: Training data

Direct

Channel

Method

x : Test input

y : Test output

$(x_1, y_1), \dots, (x_k, y_k)$: Training data

Direct

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(y \mid x)$$

Channel

Zero-shot

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(x \mid y)$$

Method

x : Test input

y : Test output

$(x_1, y_1), \dots, (x_k, y_k)$: Training data

Direct

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(y \mid x)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(y \mid x_1, y_1, \dots, x_k, y_k, x)$$

Channel

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(x \mid y)$$

Zero-shot

In-context
learning

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(x \mid y_1, x_1, \dots, y_k, x_k, y)$$

Method

x : Test input

y : Test output

$(x_1, y_1), \dots, (x_k, y_k)$: Training data

Direct

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(y \mid x)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(y \mid x_1, y_1, \dots, x_k, y_k, x)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} \prod_{i=1}^k P(y \mid x_i, y_i, x)$$

Channel

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(x \mid y)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(x \mid y_1, x_1, \dots, y_k, x_k, y)$$

Zero-shot

In-context
learning

Ens. In-context
learning

$$\operatorname{argmax}_{y \in \mathcal{Y}} \prod_{i=1}^k P(x \mid y_i, x_i, y)$$

Method

x : Test input

y : Test output

$(x_1, y_1), \dots, (x_k, y_k)$: Training data

Direct

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(y | x)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(y | x_1, y_1, \dots, x_k, y_k, x)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} \prod_{i=1}^k P(y | x_i, y_i, x)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(y | \mathbf{V}_{\text{prompt}}, x)$$

Zero-shot

In-context learning

Ens. In-context learning

Prompt tuning (Lester et al. 2021)

($\mathbf{V}_{\text{prompt}}$ is trainable)

Channel

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(x | y)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(x | y_1, x_1, \dots, y_k, x_k, y)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} \prod_{i=1}^k P(x | y_i, x_i, y)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(x | \mathbf{V}_{\text{prompt}}, y)$$

Method

x : Test input

y : Test output

$(x_1, y_1), \dots, (x_k, y_k)$: Training data

Direct

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(y | x)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(y | x_1, y_1, \dots, x_k, y_k, x)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} \prod_{i=1}^k P(y | x_i, y_i, x)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(y | \mathbf{V}_{\text{prompt}}, x)$$

Zero-shot

In-context learning

Ens. In-context learning

Prompt tuning (Lester et al. 2021)

($\mathbf{V}_{\text{prompt}}$ is trainable)

Channel

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(x | y)$$

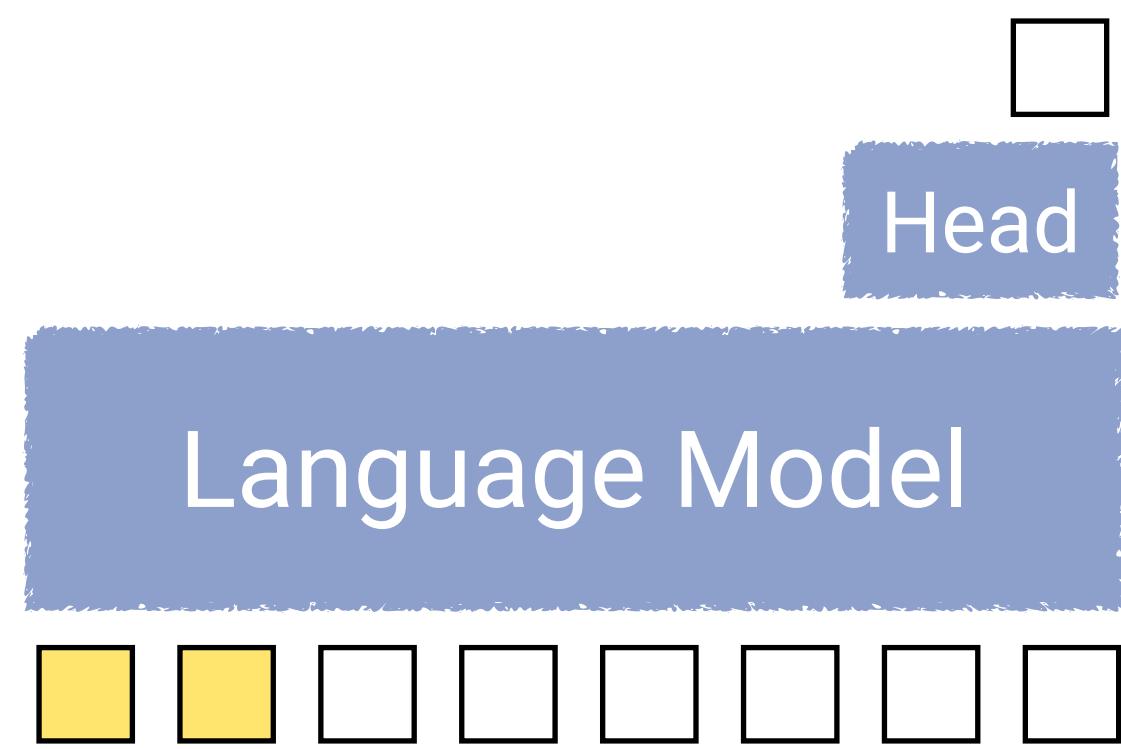
$$\operatorname{argmax}_{y \in \mathcal{Y}} P(x | y_1, x_1, \dots, y_k, x_k, y)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} \prod_{i=1}^k P(x | y_i, x_i, y)$$

$$\operatorname{argmax}_{y \in \mathcal{Y}} P(x | \mathbf{V}_{\text{prompt}}, y)$$

Methods – Prompt tuning setup

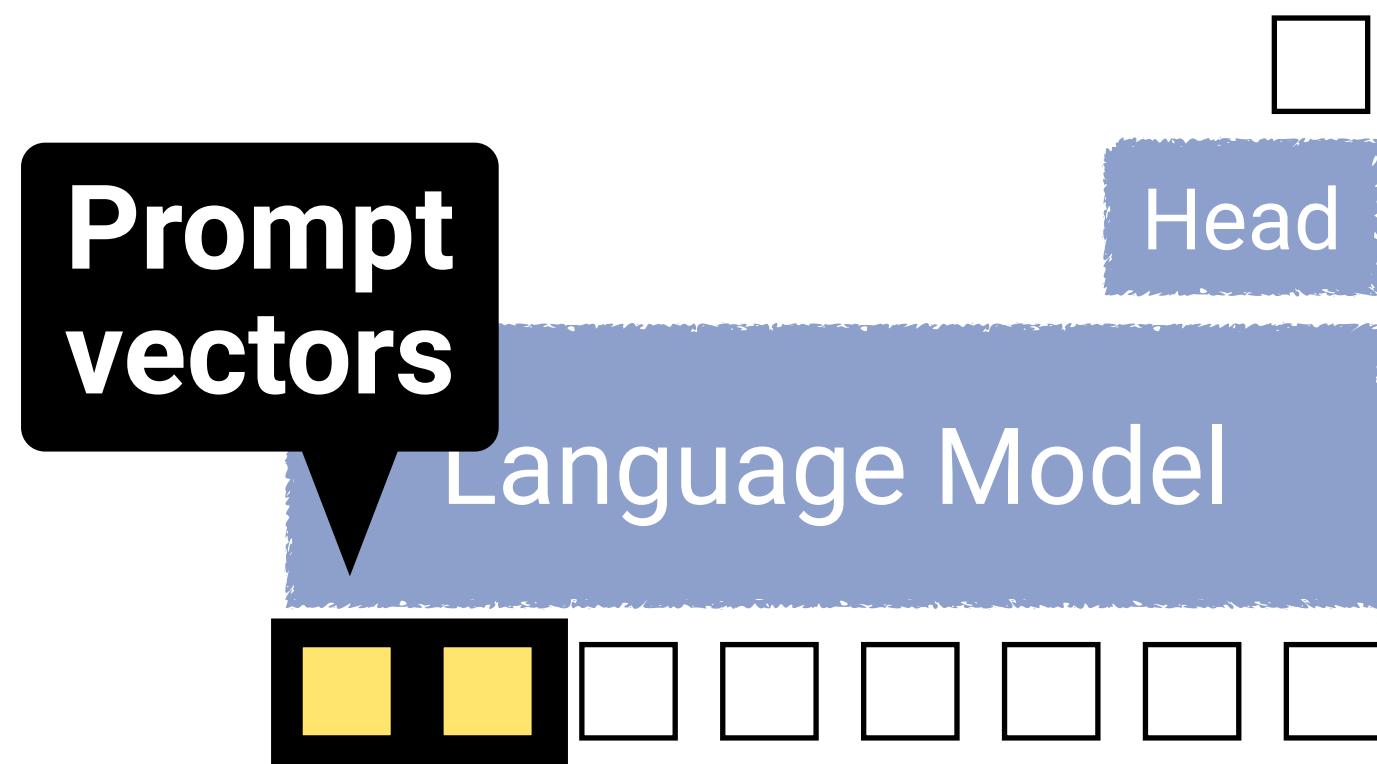
Methods – Prompt tuning setup



Prompt tuning (Lester et al. 2021)

(Yellow box = trainable)

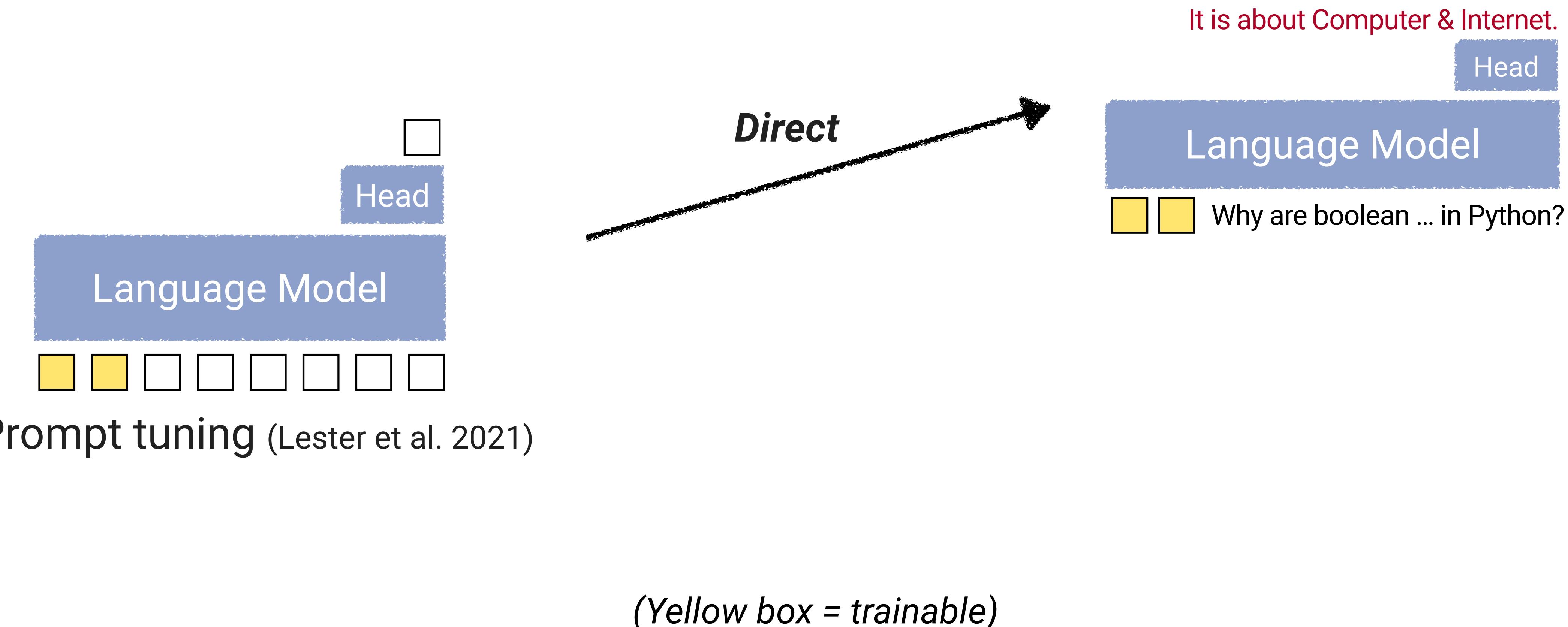
Methods – Prompt tuning setup



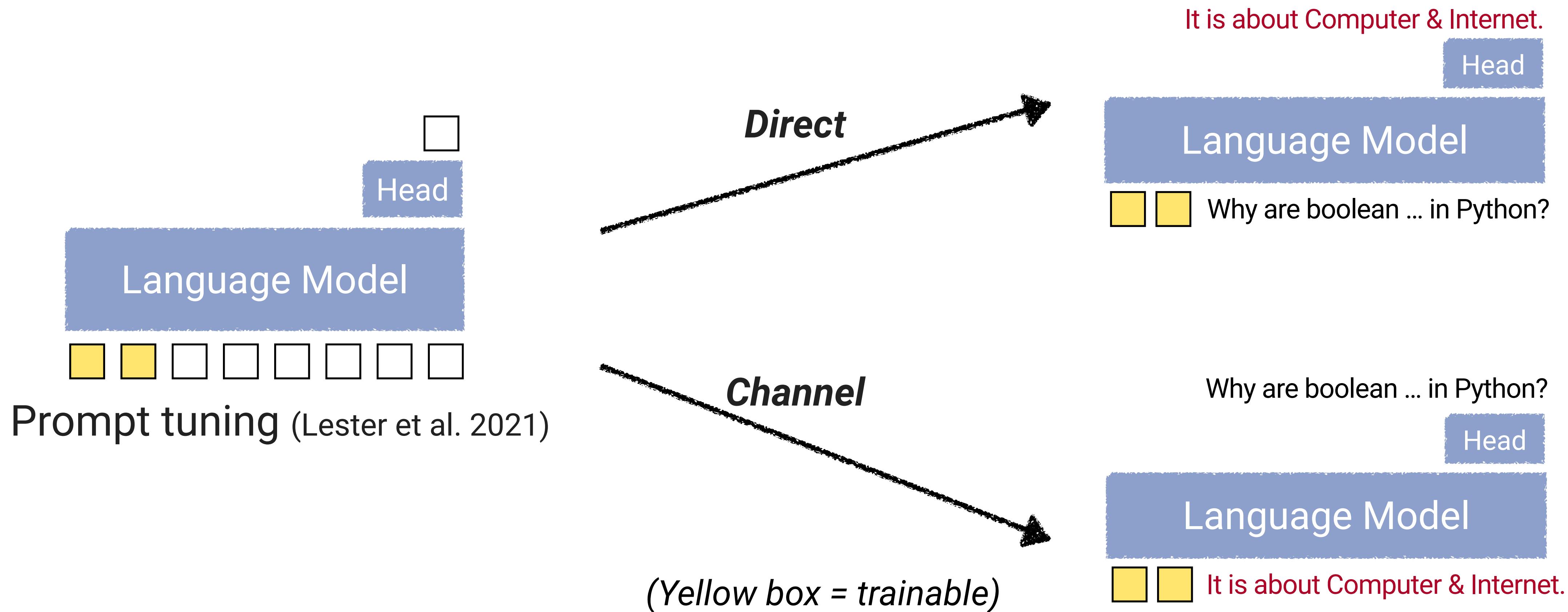
Prompt tuning (Lester et al. 2021)

(Yellow box = trainable)

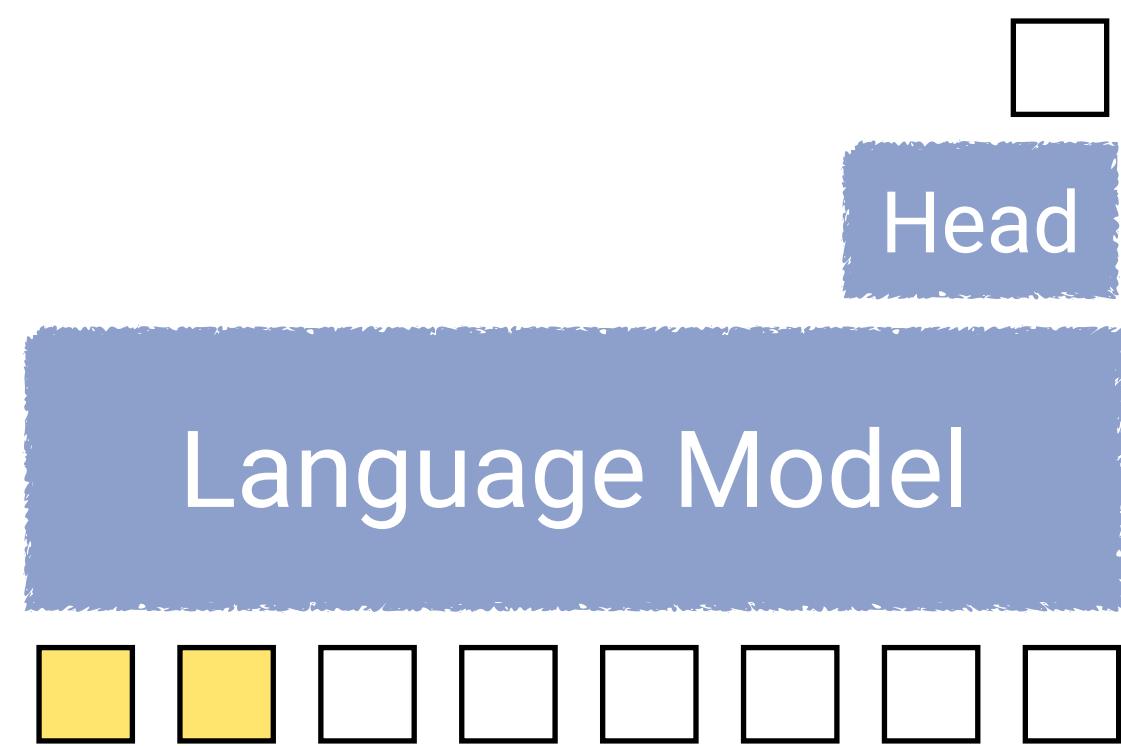
Methods – Prompt tuning setup



Methods – Prompt tuning setup



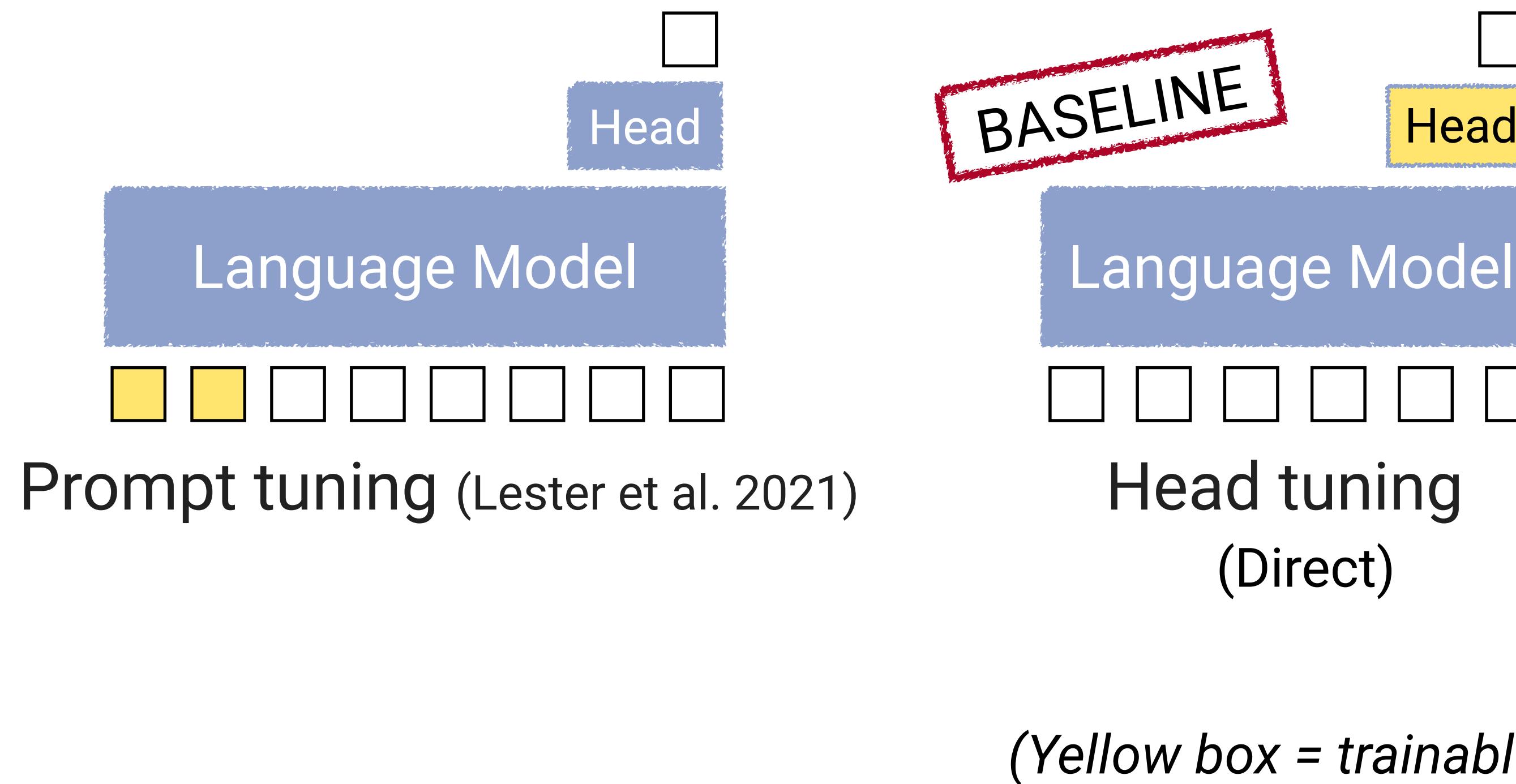
Methods – Prompt tuning setup



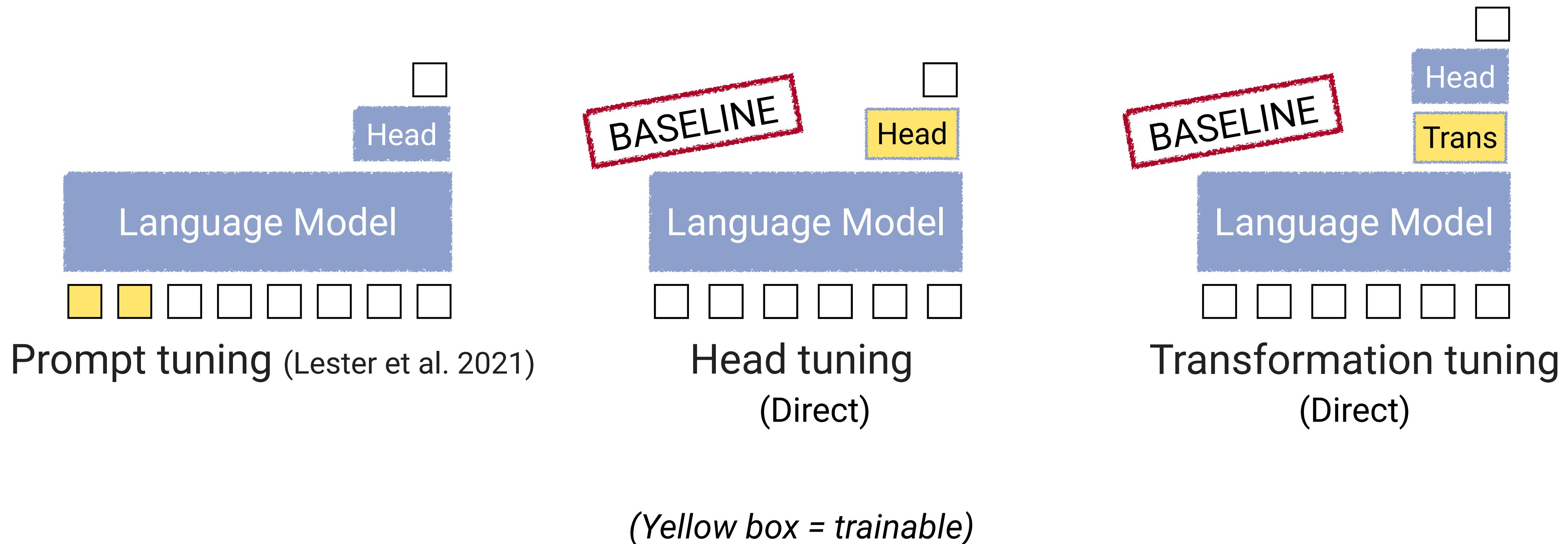
Prompt tuning (Lester et al. 2021)

(Yellow box = trainable)

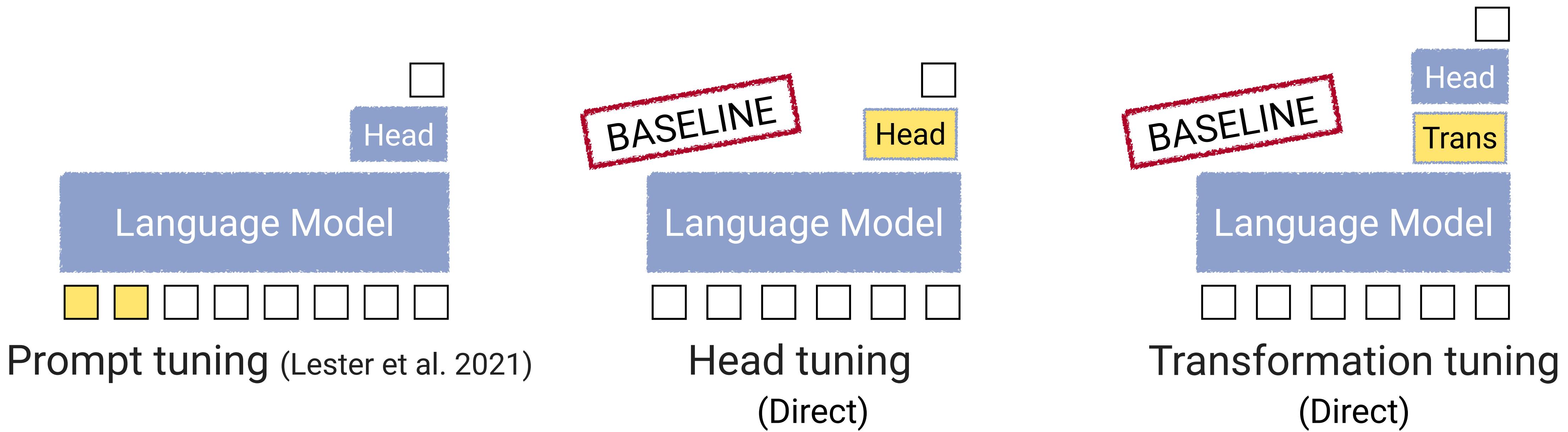
Methods – Prompt tuning setup



Methods – Prompt tuning setup



Methods – Prompt tuning setup



(Yellow box = *trainable*)

tune ~0.002% of the total parameters

Why would it work better?

Why would it work better?

1. Few-shot learning

Noisy channel models are shown to work better in a **few-shot** setup – theoretically (Ng and Jordan, 2002) and empirically (Ding and Gimpel, 2019)

Why would it work better?

1. Few-shot learning

Noisy channel models are shown to work better in a **few-shot** setup – theoretically (Ng and Jordan, 2002) and empirically (Ding and Gimpel, 2019)

Predicting the entire input *amplifies* training signals (Lewis and Fan, 2018)

Why would it work better?

1. Few-shot learning

Noisy channel models are shown to work better in a **few-shot** setup – theoretically (Ng and Jordan, 2002) and empirically (Ding and Gimpel, 2019)

Direct

The company anticipated its operating profit to improve.

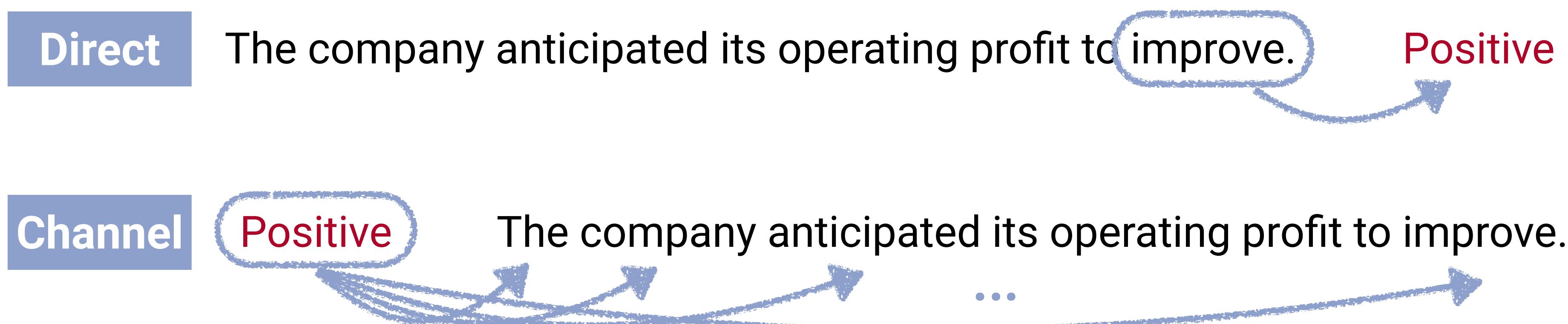
Positive

Predicting the entire input *amplifies* training signals (Lewis and Fan, 2018)

Why would it work better?

1. Few-shot learning

Noisy channel models are shown to work better in a **few-shot** setup – theoretically (Ng and Jordan, 2002) and empirically (Ding and Gimpel, 2019)



Predicting the entire input *amplifies* training signals (Lewis and Fan, 2018)

Why would it work better?

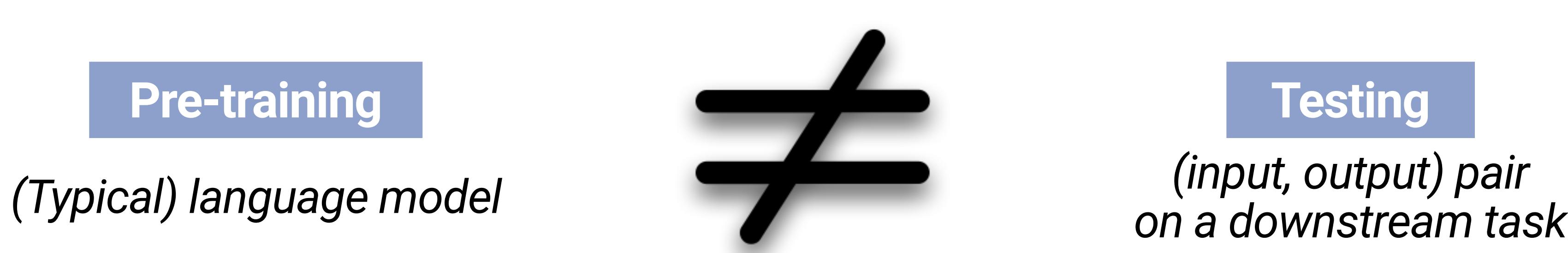
2. Distribution shifts

Noisy channel models are shown to work better under **distribution shifts** (Yogtama et al. 2017, Lewis and Fan, 2018)

Why would it work better?

2. Distribution shifts

Noisy channel models are shown to work better under **distribution shifts** (Yogtama et al. 2017, Lewis and Fan, 2018)



- 1) LM prompting always needs to deal with **distribution shifts**

Why would it work better?

2. Distribution shifts

Noisy channel models are shown to work better under **distribution shifts** (Yogtama et al. 2017, Lewis and Fan, 2018)

- 1) LM prompting always needs to deal with **distribution shifts**
- 2) Channel model may be better at predicting **unseen** labels

Why would it work better?

2. Distribution shifts

Noisy channel models are shown to work better under **distribution shifts** (Yogtama et al. 2017, Lewis and Fan, 2018)

$$\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^N \quad \{(y_i)\}_{i=1}^N = \{\text{positive, negative}\}$$

- 1) LM prompting always needs to deal with **distribution shifts**
- 2) Channel model may be better at predicting **unseen** labels

Why would it work better?

2. Distribution shifts

Noisy channel models are shown to work better under **distribution shifts** (Yogtama et al. 2017, Lewis and Fan, 2018)

$$\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^N \quad \{(y_i)\}_{i=1}^N = \{\text{positive, negative}\}$$

Panostaja did not disclose the purchase price.  neutral

- 1) LM prompting always needs to deal with **distribution shifts**
- 2) Channel model may be better at predicting **unseen** labels

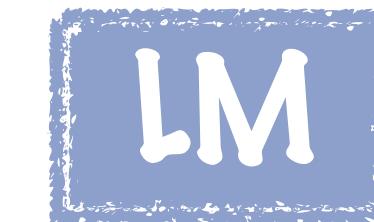
Why would it work better?

2. Distribution shifts

Noisy channel models are shown to work better under **distribution shifts** (Yogtama et al. 2017, Lewis and Fan, 2018)

$$\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^N \quad \{(y_i)\}_{i=1}^N = \{\text{positive, negative}\}$$

Panostaja did not disclose the purchase price.



neutral

← *likely very low*

- 1) LM prompting always needs to deal with **distribution shifts**
- 2) Channel model may be better at predicting **unseen** labels

Why would it work better?

2. Distribution shifts

Noisy channel models are shown to work better under **distribution shifts** (Yogtama et al. 2017, Lewis and Fan, 2018)

$$\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^N \quad \{(y_i)\}_{i=1}^N = \{\text{positive, negative}\}$$

Panostaja did not disclose the purchase price.  neutral \leftarrow *likely very low*

neutral  Panostaja did not disclose the purchase price. $\leftarrow ?$

- 1) LM prompting always needs to deal with **distribution shifts**
- 2) Channel model may be better at predicting **unseen** labels

Why would it work better?

2. Distribution shifts

Noisy channel models are shown to work better under **distribution shifts** (Yogtama et al. 2017, Lewis and Fan, 2018)

$$\mathcal{D}_{\text{train}} = \{(x_i, y_i)\}_{i=1}^N \quad \{(y_i)\}_{i=1}^N = \{\text{positive, negative}\}$$

Panostaja did not disclose the purchase price.  neutral \leftarrow *likely very low*

neutral  Panostaja did not disclose the purchase price. $\leftarrow ?$

1) LM prompting always needs to deal with **distribution shifts**

2) Channel model may be better at predicting **unseen** labels
(*experiments showing this later!*)

Experimental Setup

Dataset	Task	# labels
SST-2	Sentiment analysis (movie)	2
SST-5	Sentiment analysis (movie)	5
MR	Sentiment analysis (movie)	2
CR	Sentiment analysis (electronics)	2
Amazon	Sentiment analysis (Amazon)	5
Yelp	Sentiment analysis (Yelp)	5
TREC	Question classification (answer type)	6
AGNews	News classification (topic)	4
Yahoo	Question classification (topic)	10
DBPedia	Ontology classification	14
Subj	Subjectivity classification	2

Experimental Setup

Dataset	Task	# labels
SST-2	Sentiment analysis (movie)	2
SST-5	Sentiment analysis (movie)	5
MR	Sentiment analysis (movie)	2
CR	Sentiment analysis (electronics)	2
Amazon	Sentiment analysis (Amazon)	5
Yelp	Sentiment analysis (Yelp)	5
TREC	Question classification (answer type)	6
AGNews	News classification (topic)	4
Yahoo	Question classification (topic)	10
DBPedia	Ontology classification	14
Subj	Subjectivity classification	2



GPT-2
LARGE

(# params = 770M)

Experimental Setup

Dataset	Task	# labels
SST-2	Sentiment analysis (movie)	2
SST-5	Sentiment analysis (movie)	5
MR	Sentiment analysis (movie)	2
CR	Sentiment analysis (electronics)	2
Amazon	Sentiment analysis (Amazon)	5
Yelp	Sentiment analysis (Yelp)	5
TREC	Question classification (answer type)	6
AGNews	News classification (topic)	4
Yahoo	Question classification (topic)	10
DBPedia	Ontology classification	14
Subj	Subjectivity classification	2



GPT-2
LARGE

(# params = 770M)

Training data: K=16

Experimental Setup

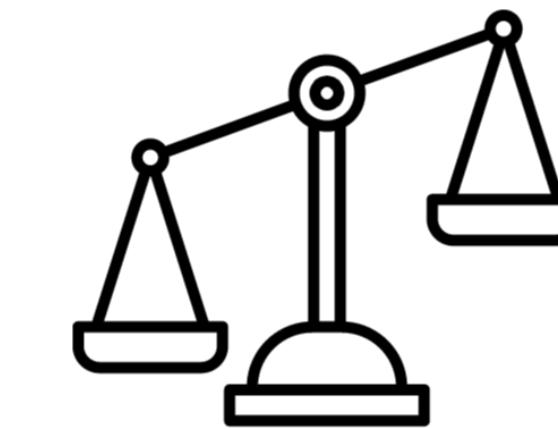
Dataset	Task	# labels
SST-2	Sentiment analysis (movie)	2
SST-5	Sentiment analysis (movie)	5
MR	Sentiment analysis (movie)	2
CR	Sentiment analysis (electronics)	2
Amazon	Sentiment analysis (Amazon)	5
Yelp	Sentiment analysis (Yelp)	5
TREC	Question classification (answer type)	6
AGNews	News classification (topic)	4
Yahoo	Question classification (topic)	10
DBPedia	Ontology classification	14
Subj	Subjectivity classification	2



GPT-2
LARGE

(# params = 770M)

Training data: K=16



No Guarantee

Experimental Setup

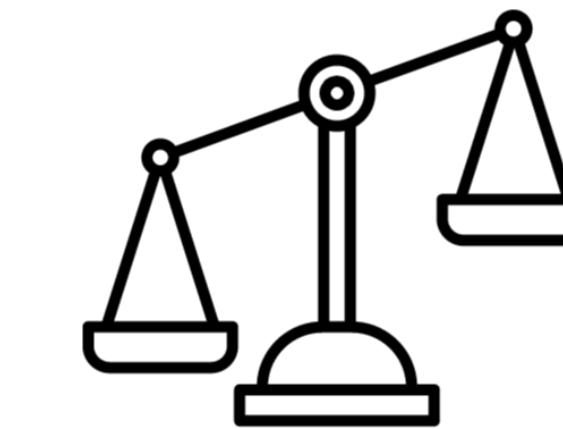
Dataset	Task	# labels
SST-2	Sentiment analysis (movie)	2
SST-5	Sentiment analysis (movie)	5
MR	Sentiment analysis (movie)	2
CR	Sentiment analysis (electronics)	2
Amazon	Sentiment analysis (Amazon)	5
Yelp	Sentiment analysis (Yelp)	5
TREC	Question classification (answer type)	6
AGNews	News classification (topic)	4
Yahoo	Question classification (topic)	10
DBPedia	Ontology classification	14
Subj	Subjectivity classification	2



GPT-2
LARGE

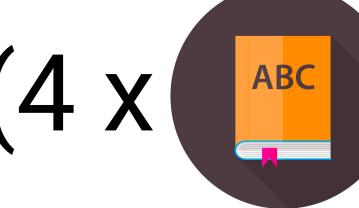
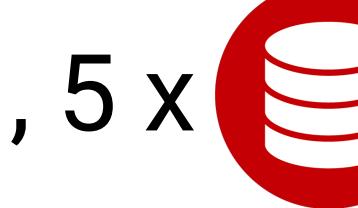
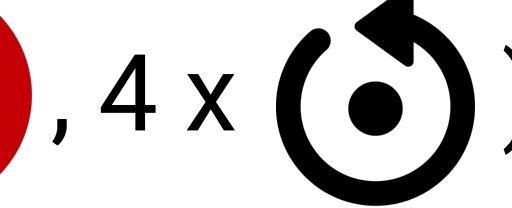
(# params = 770M)

Training data: K=16

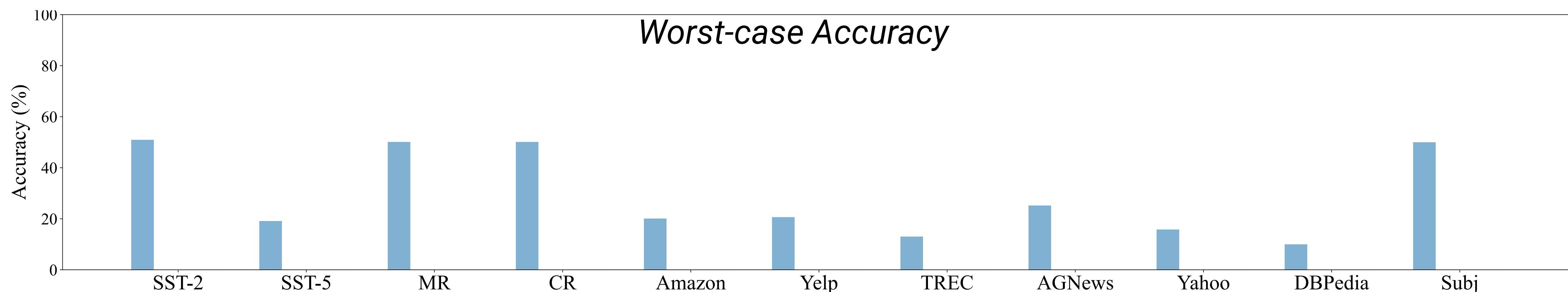
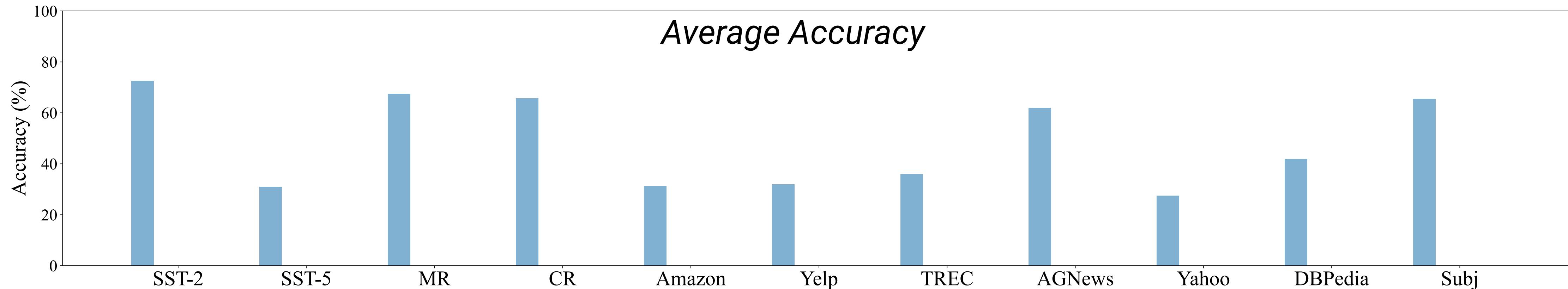
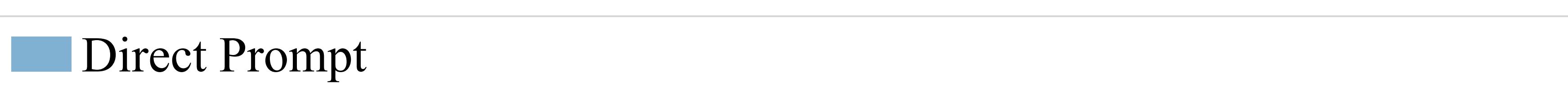


No Guarantee

80 different runs

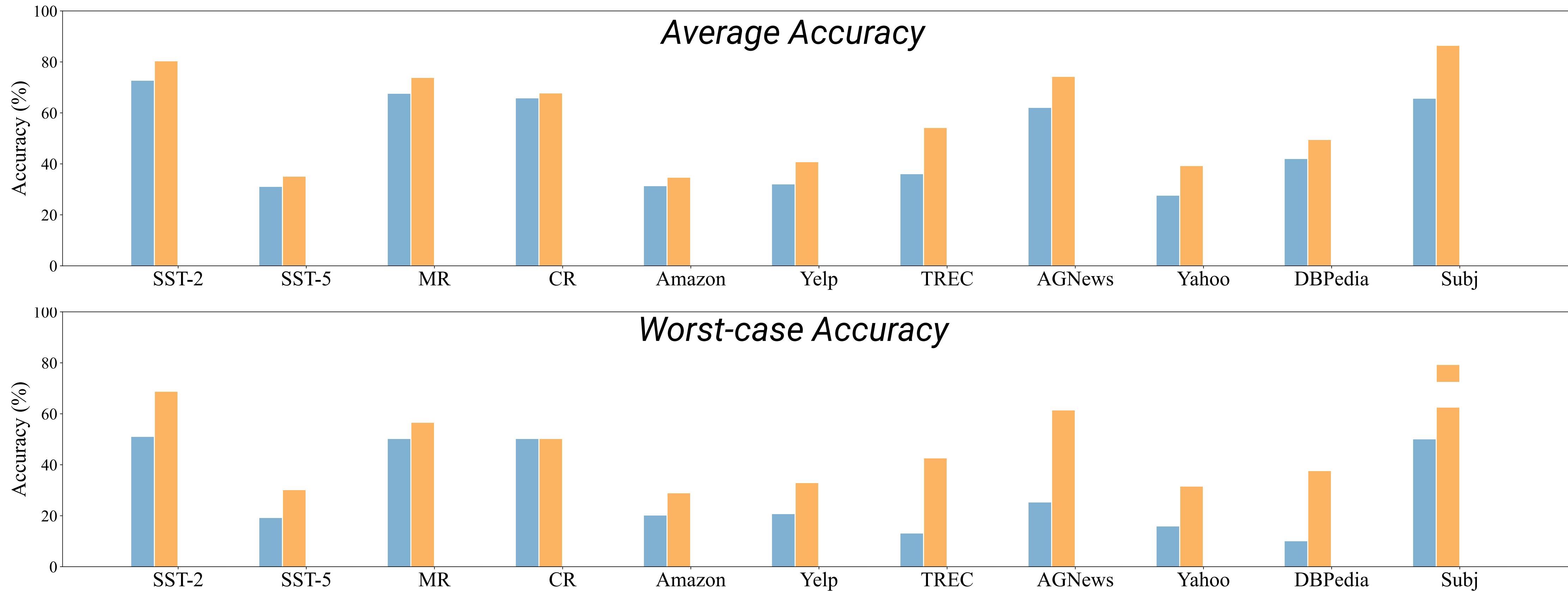
(4 x , 5 x , 4 x )

Results – Prompt tuning setup

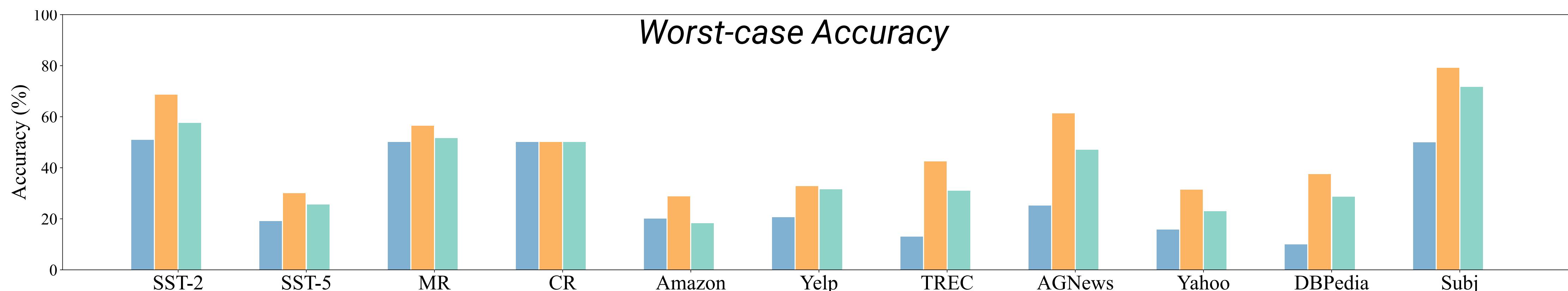
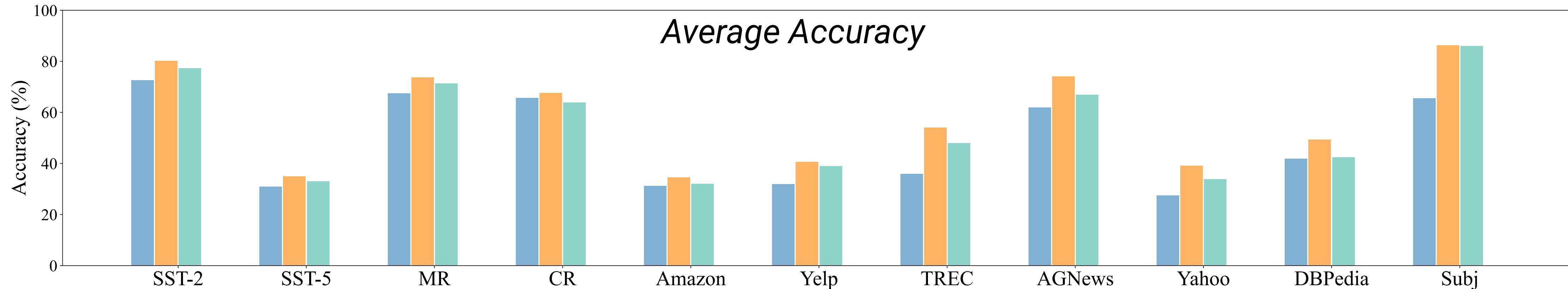
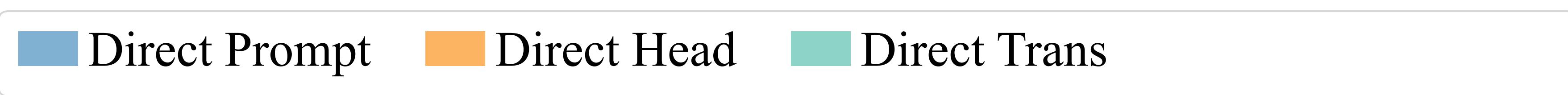


Results – Prompt tuning setup

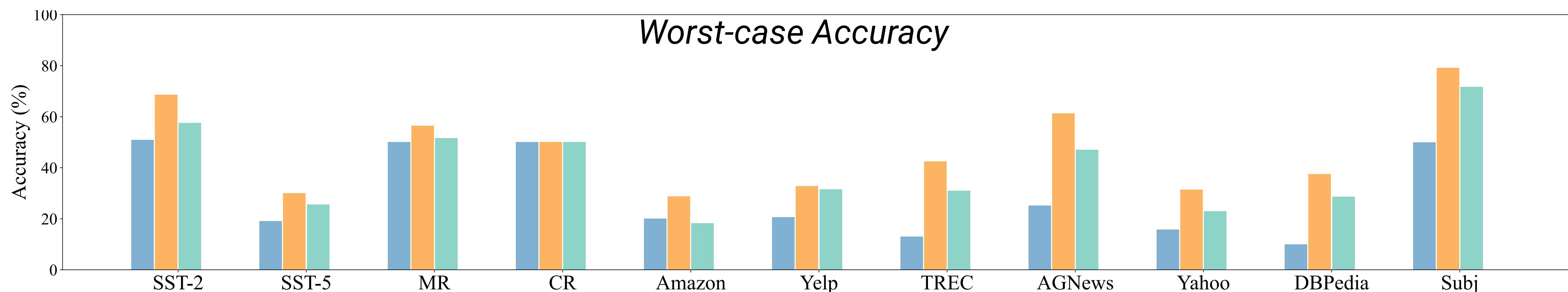
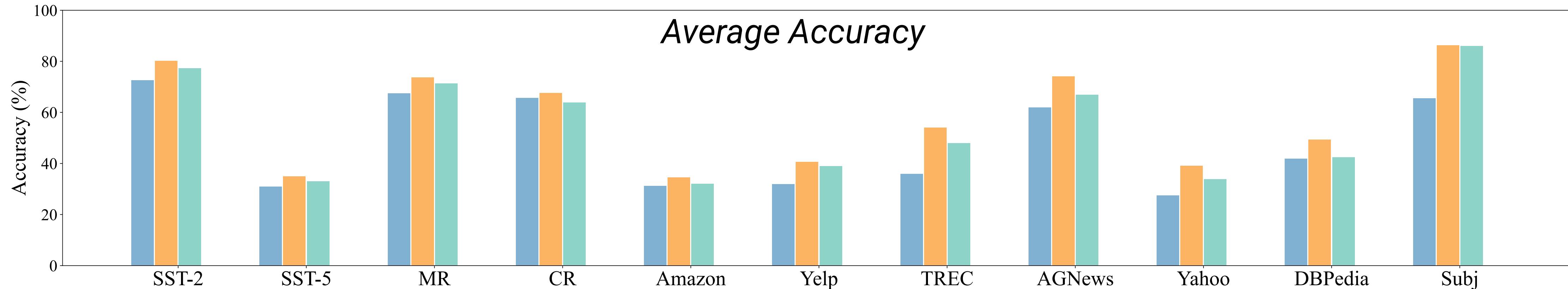
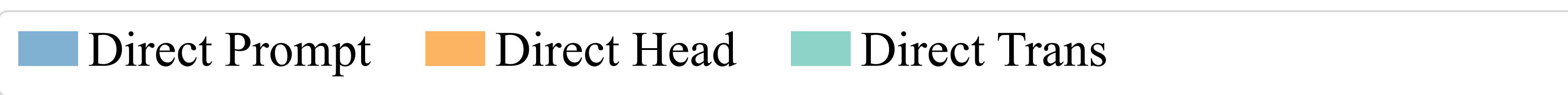
■ Direct Prompt ■ Direct Head



Results – Prompt tuning setup

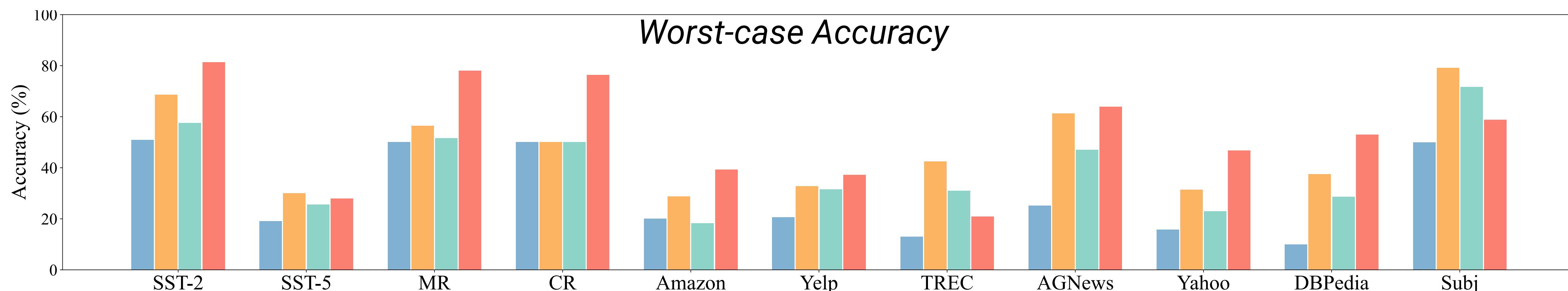
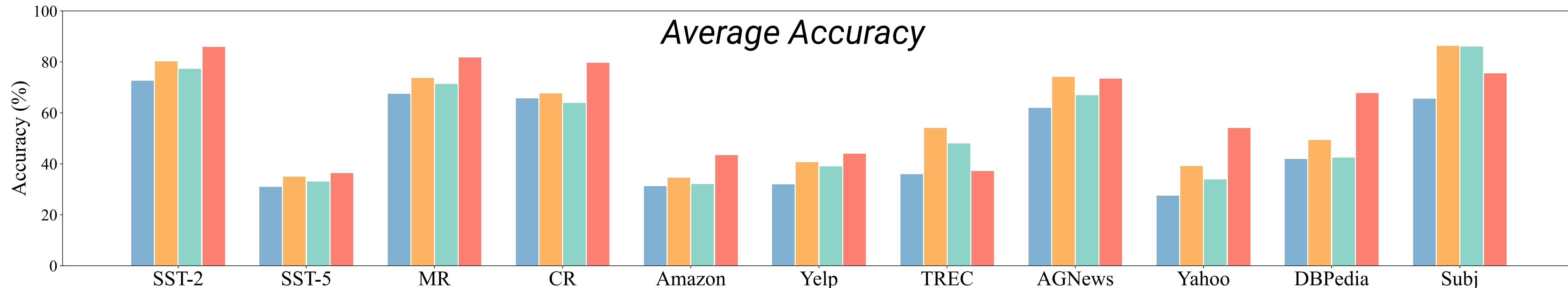
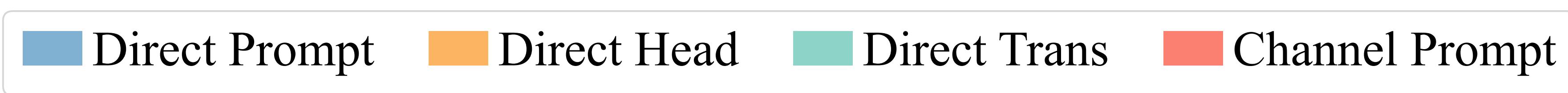


Results – Prompt tuning setup



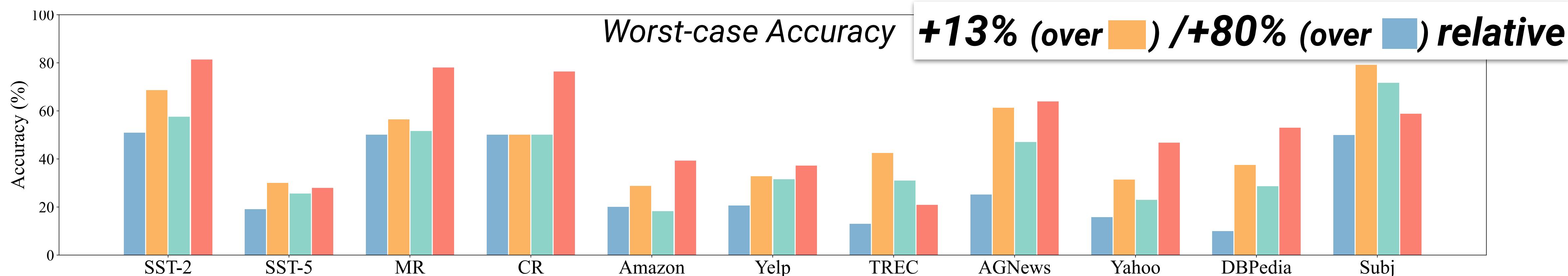
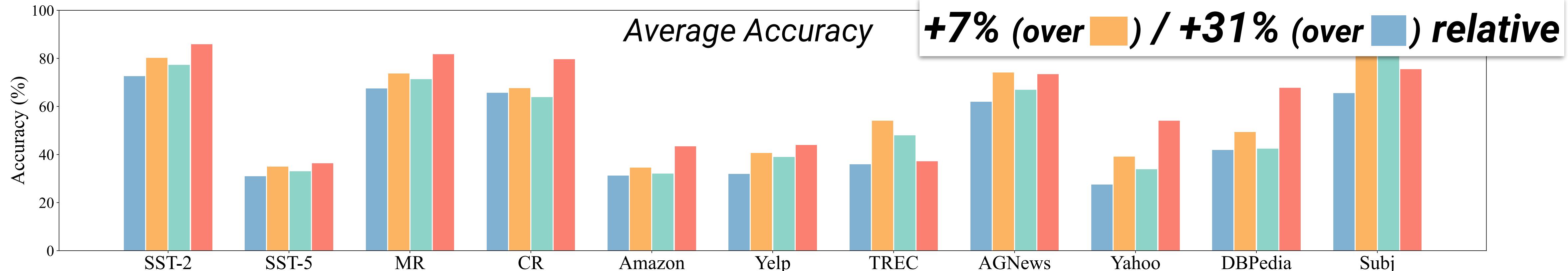
Direct Head Tuning is a powerful baseline

Results – Prompt tuning setup



Channel Prompt Tuning outperforms Direct models

Results – Prompt tuning setup



Channel Prompt Tuning outperforms Direct models
(more significantly in worst-case accuracy)

Ablations

Impact of imbalance in training data

Ablations

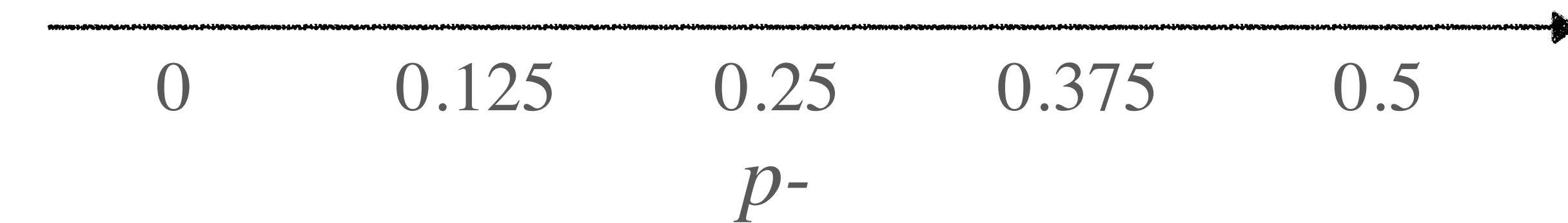
Impact of imbalance in training data

Binary classification datasets

Ablations

Impact of imbalance in training data

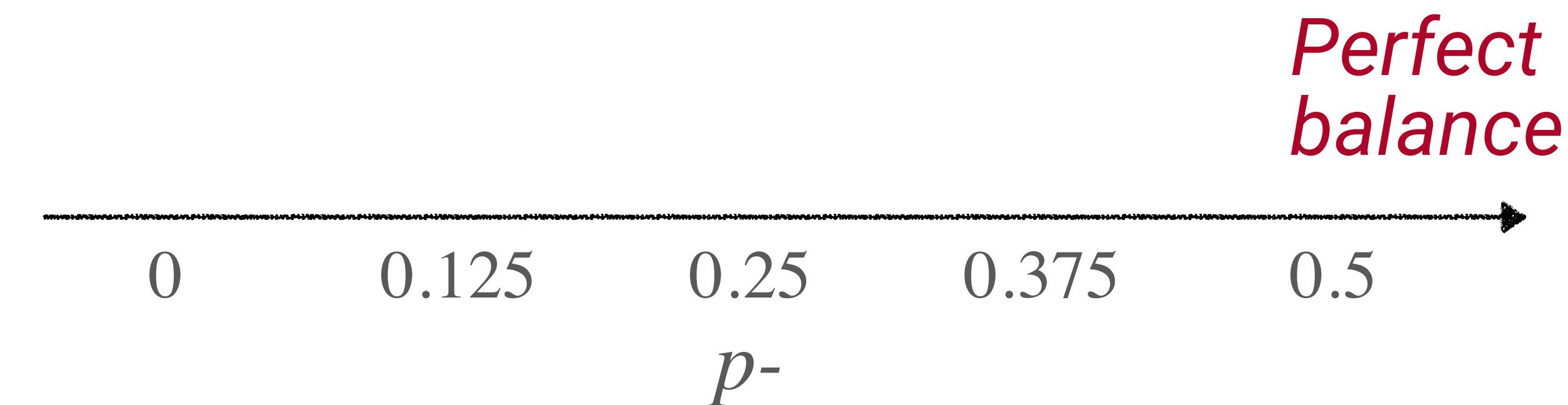
Binary classification datasets



Ablations

Impact of imbalance in training data

Binary classification datasets



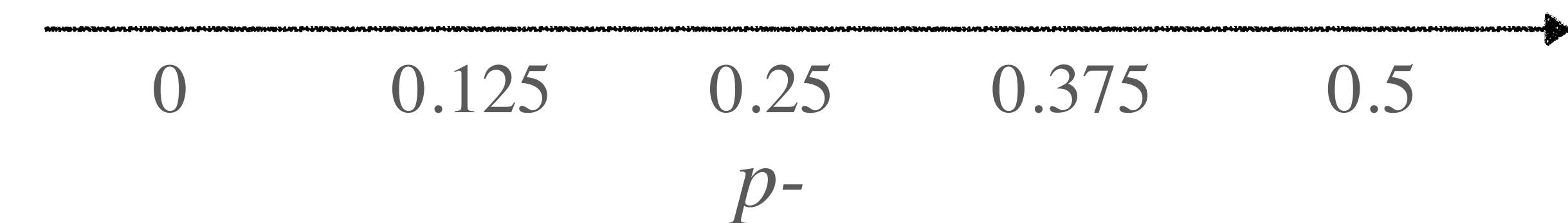
Ablations

Impact of imbalance in training data

Binary classification datasets

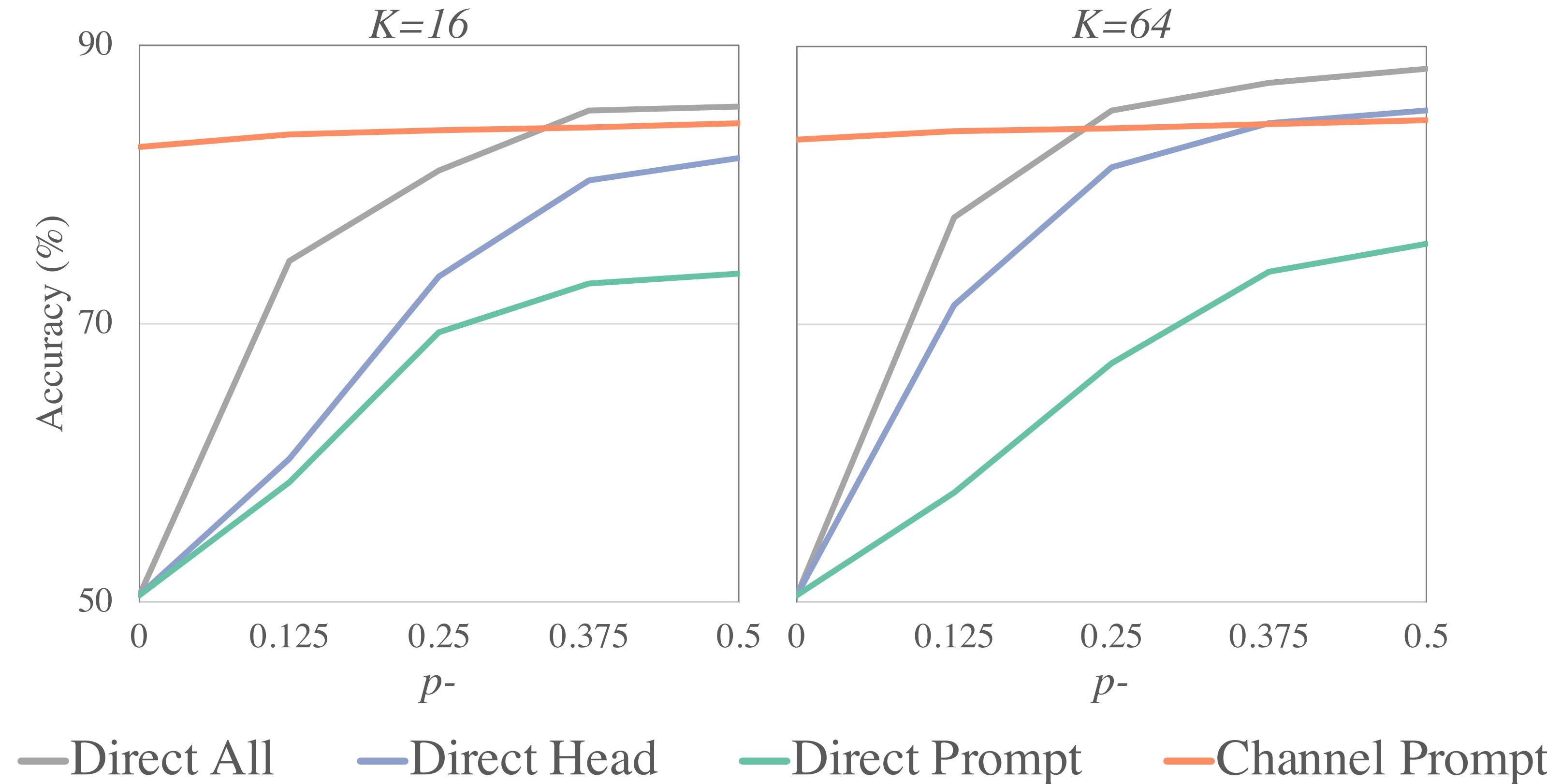
*A label is
unseen*

*Perfect
balance*



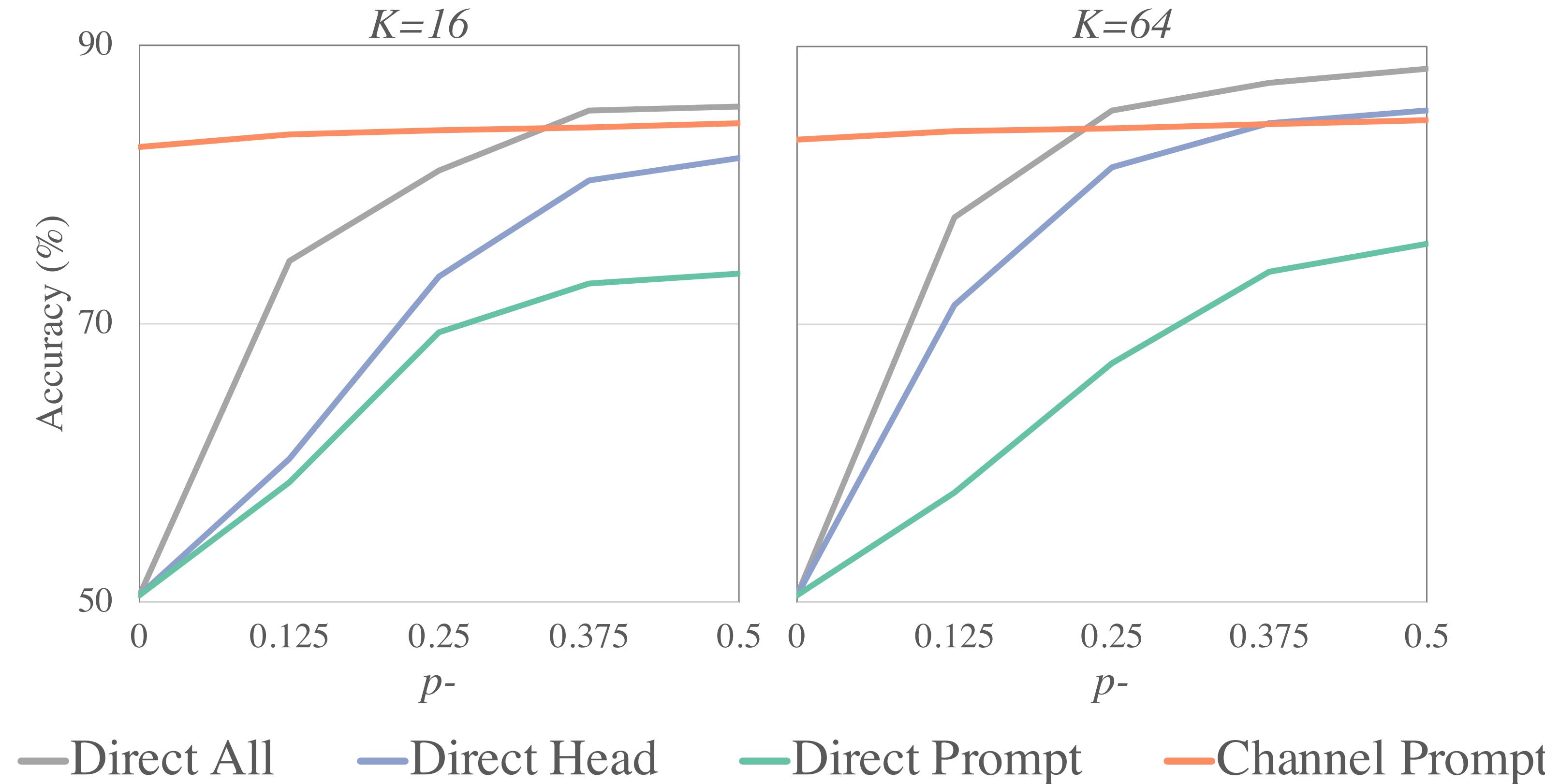
Ablations

Impact of imbalance in training data



Ablations

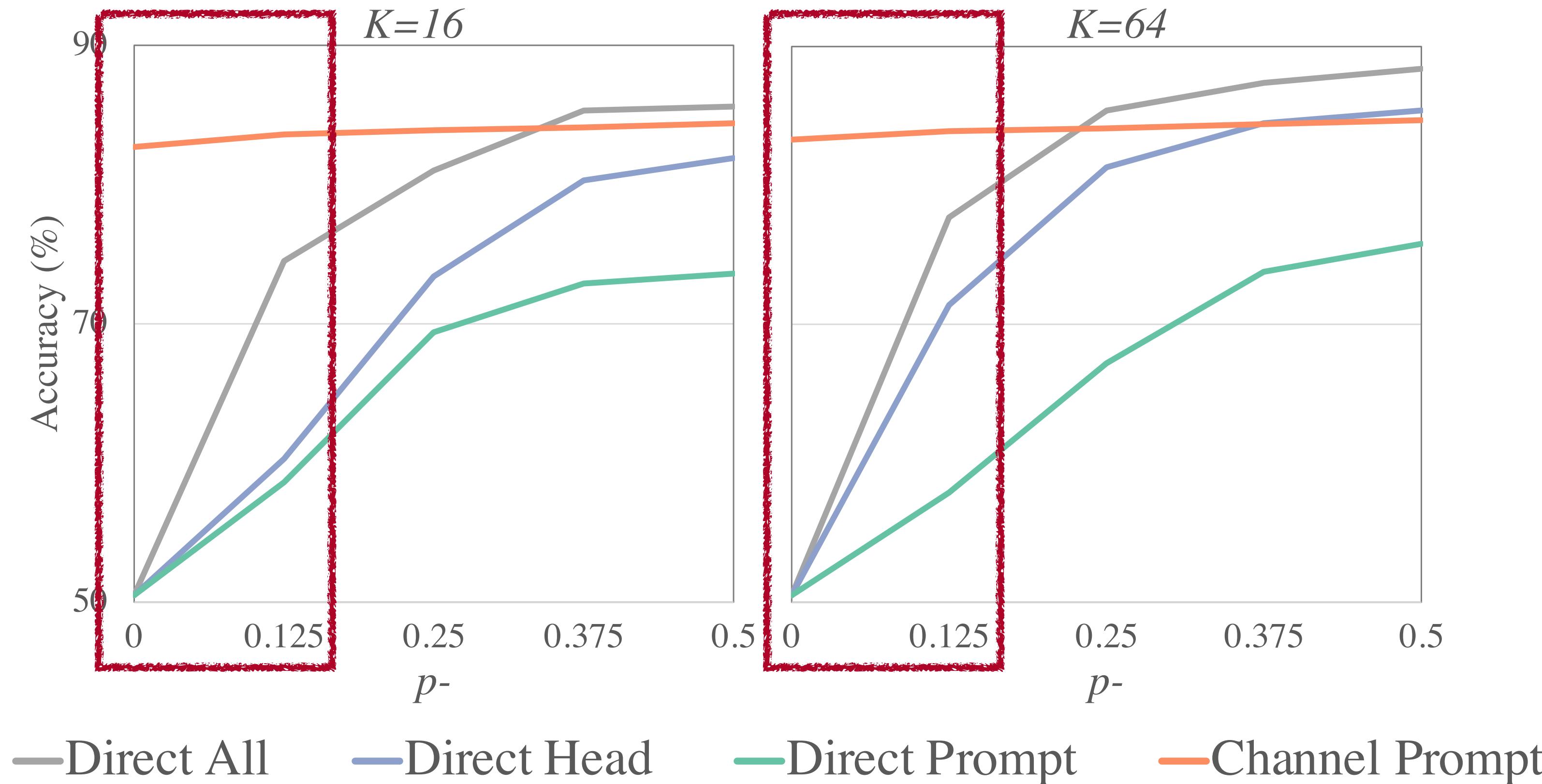
Impact of imbalance in training data



Channel models are much less sensitive to imbalance in training data

Ablations

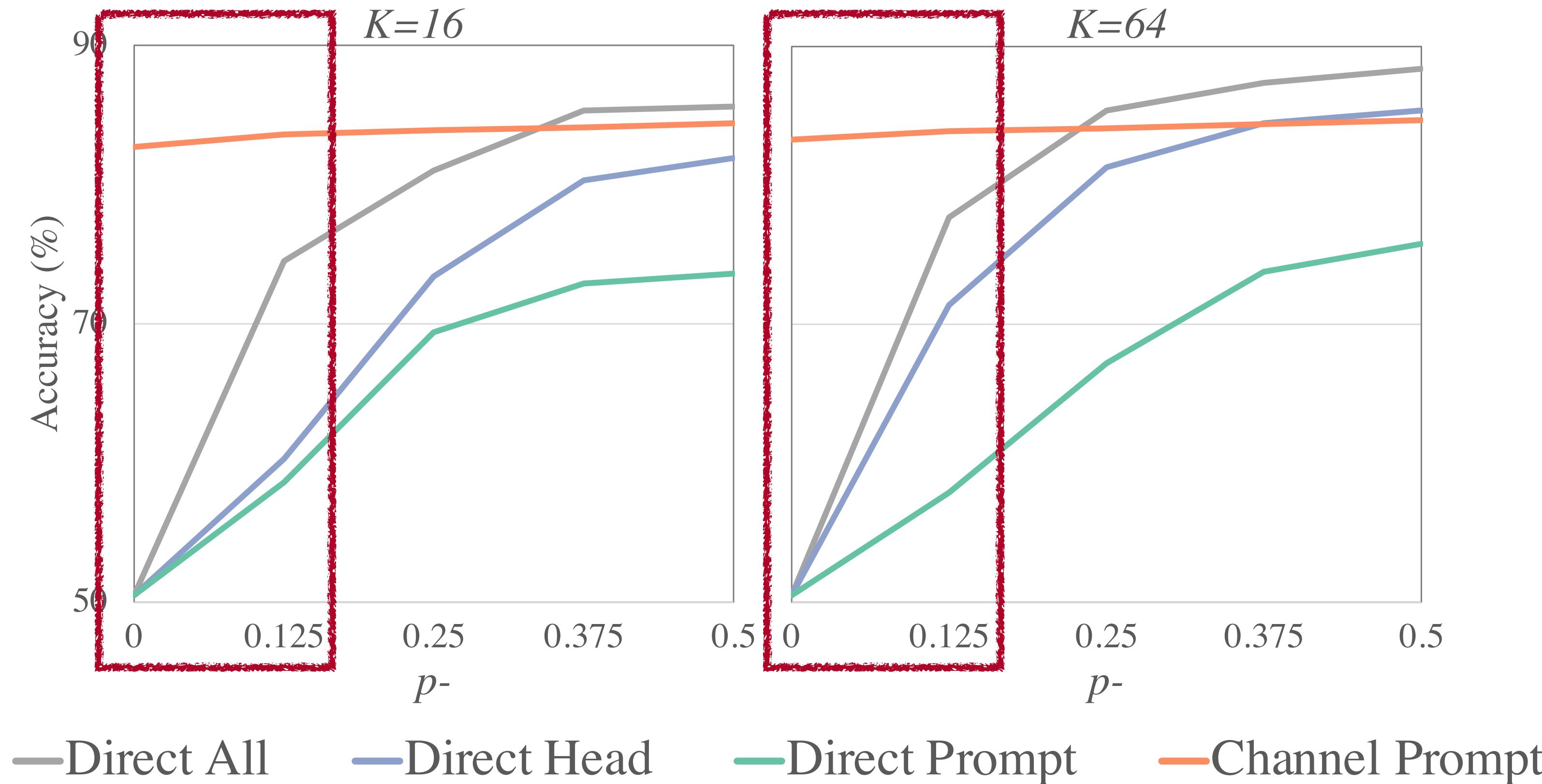
Impact of imbalance in training data



Channel models are much less sensitive to imbalance in training data
Channel models do well even with unseen labels

Ablations

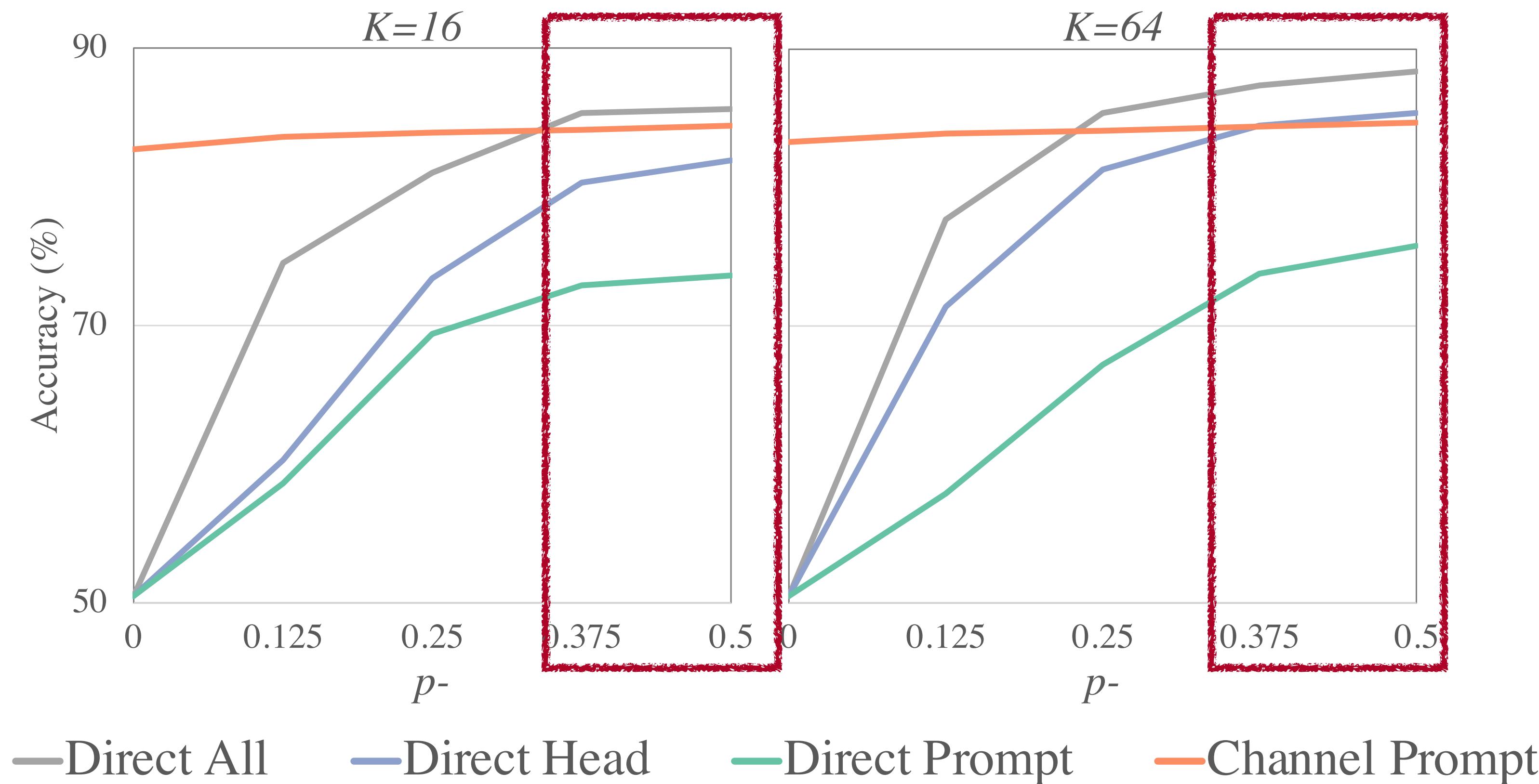
Impact of imbalance in training data



Channel models are much less sensitive to imbalance in training data
Channel models do well even with unseen labels
(*Likely because labels are conditioning variables*)

Ablations

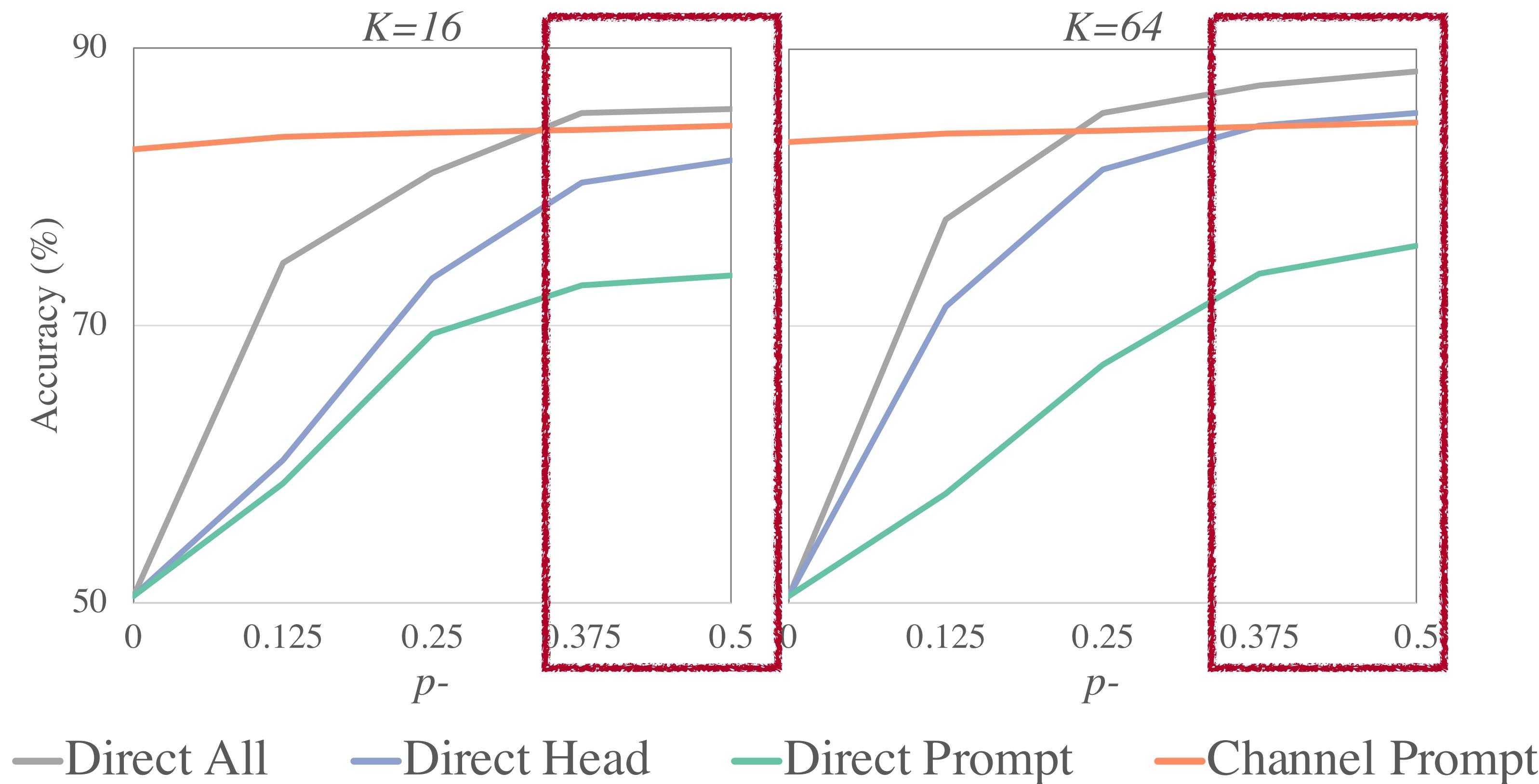
Impact of imbalance in training data



Direct models can be better with perfect balanced data

Ablations

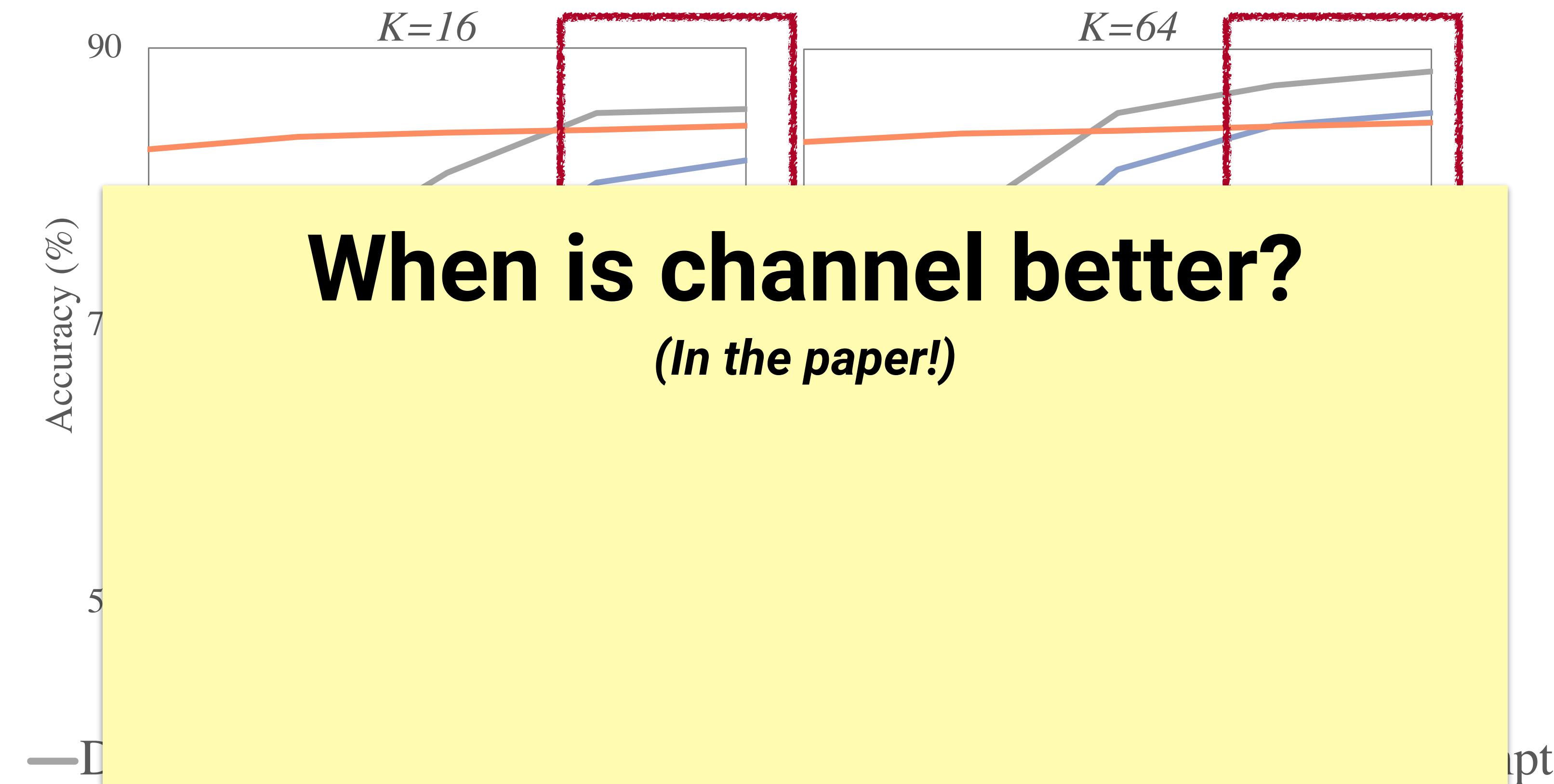
Impact of imbalance in training data



Direct models can be better with perfect balanced data

Ablations

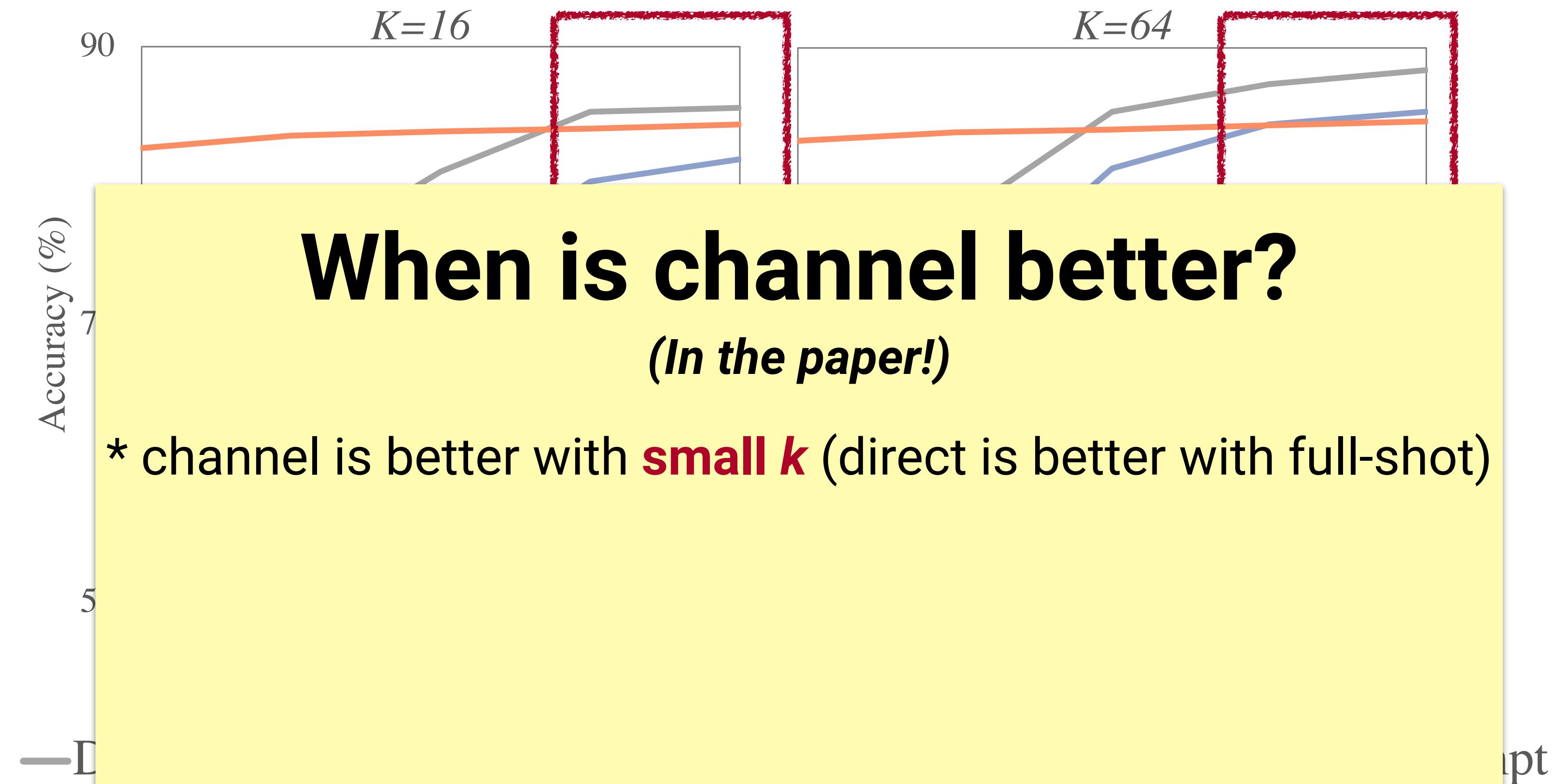
Impact of imbalance in training data



Direct models can be better with perfect balanced data

Ablations

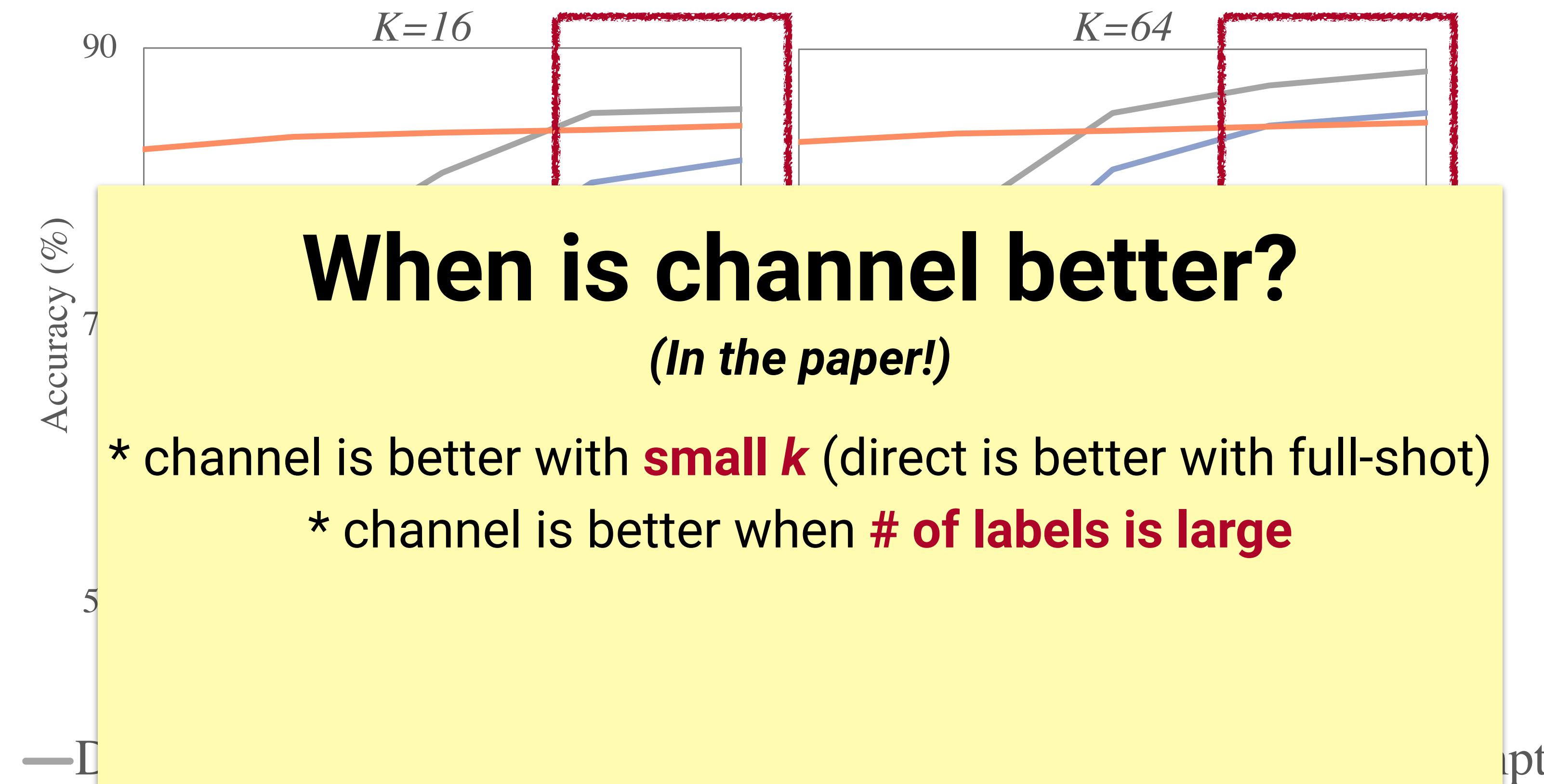
Impact of imbalance in training data



Direct models can be better with perfect balanced data

Ablations

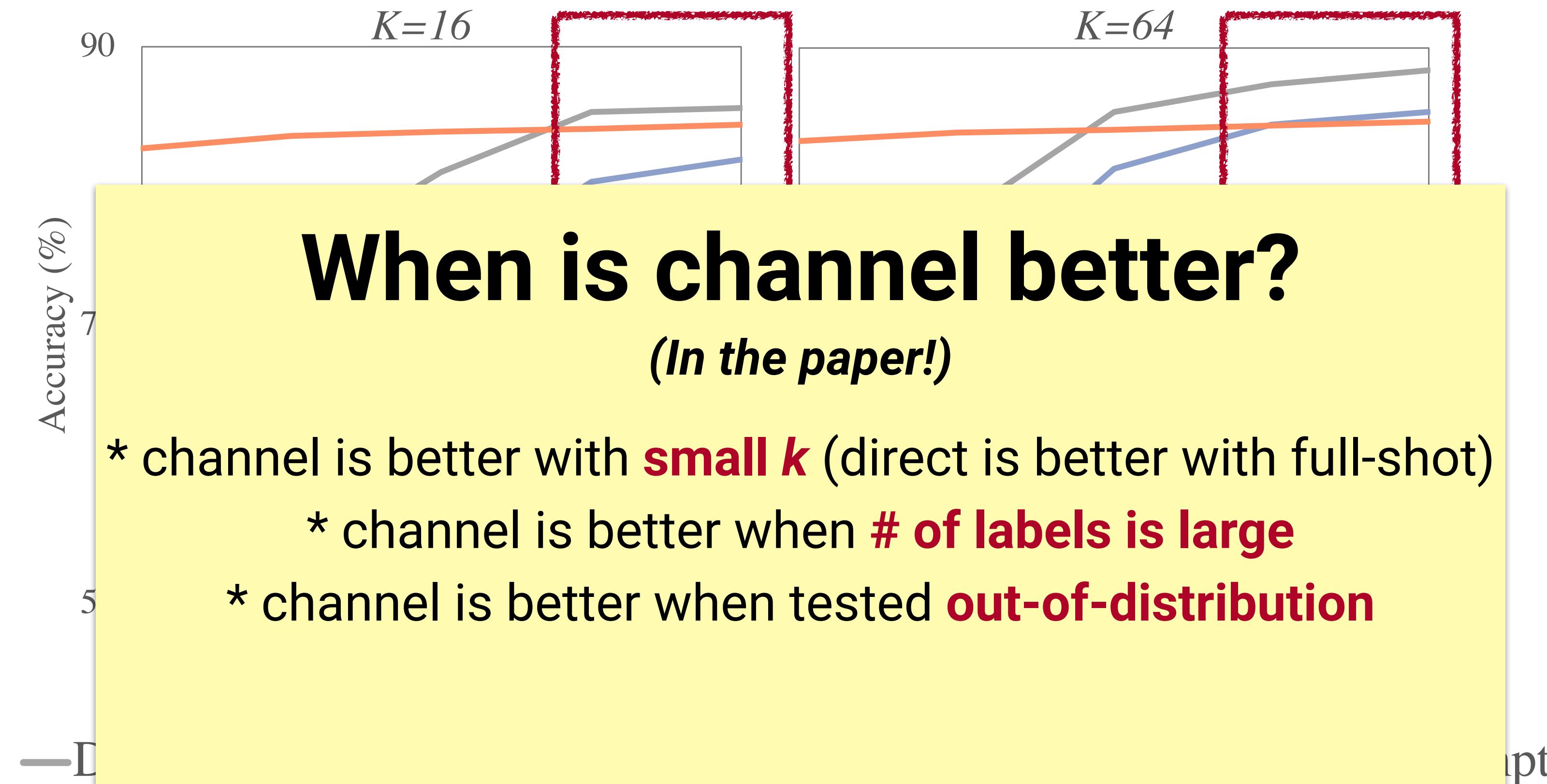
Impact of imbalance in training data



Direct models can be better with perfect balanced data

Ablations

Impact of imbalance in training data



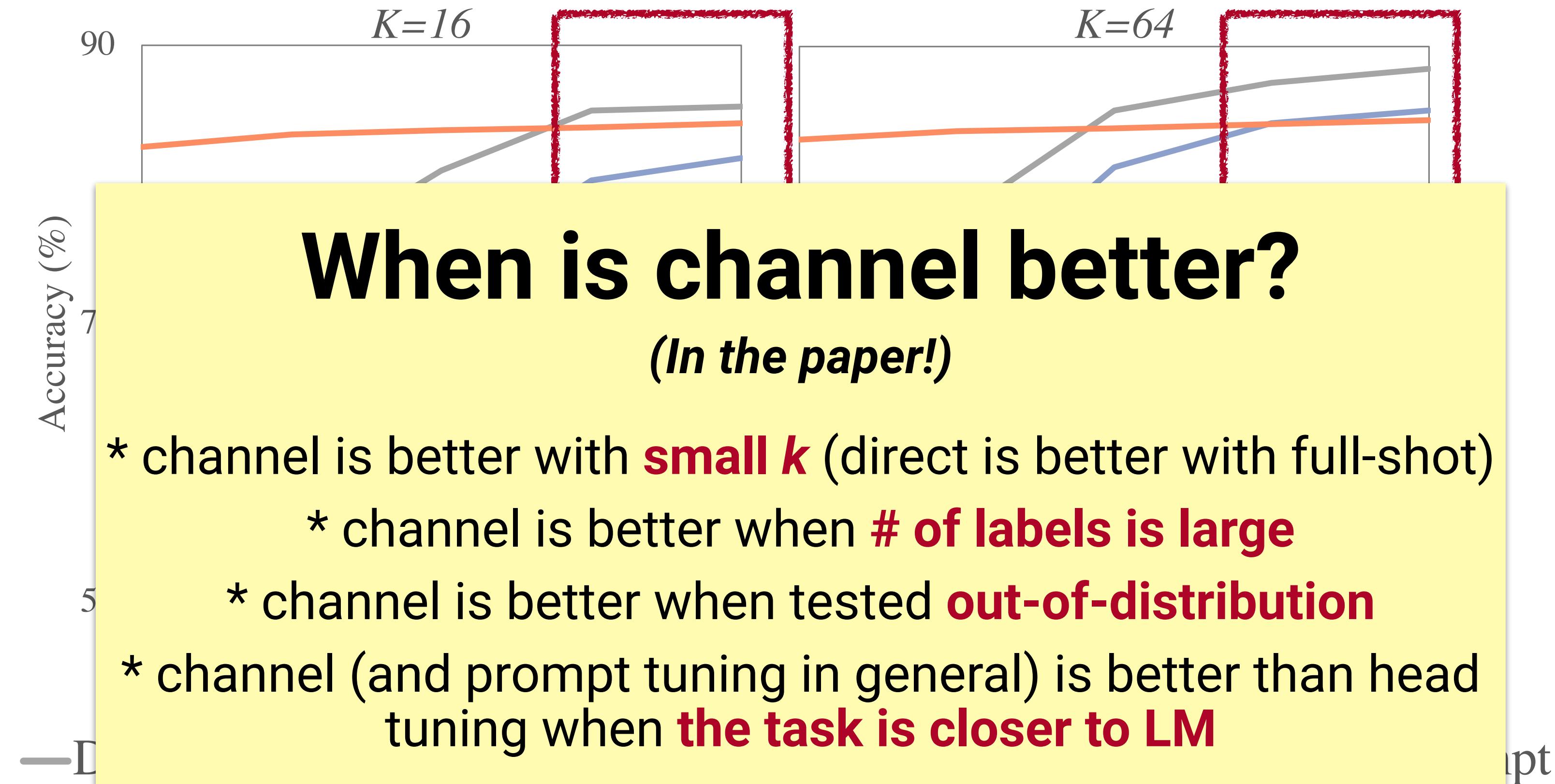
When is channel better?
(In the paper!)

- * channel is better with **small k** (direct is better with full-shot)
 - * channel is better when **# of labels is large**
 - * channel is better when tested **out-of-distribution**

Direct models can be better with perfect balanced data

Ablations

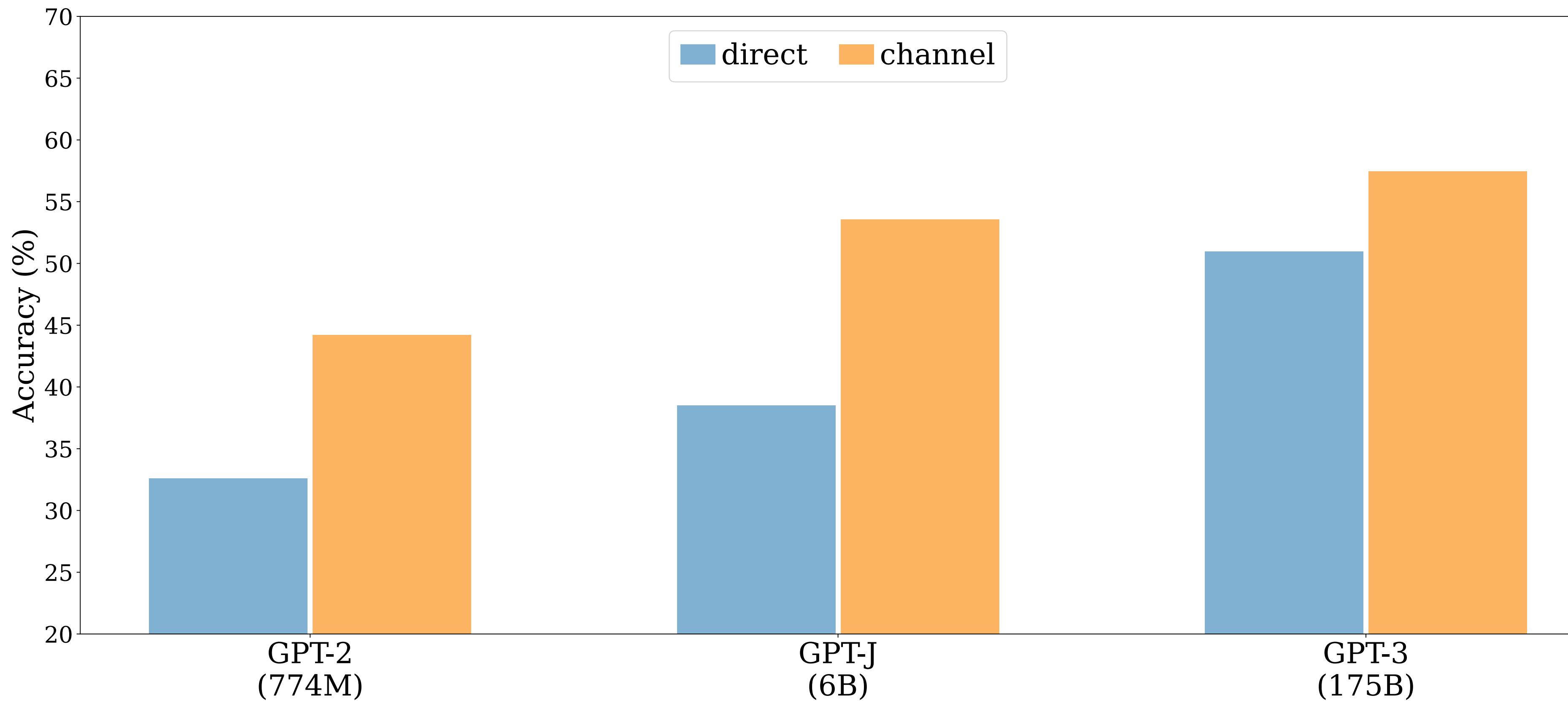
Impact of imbalance in training data



Direct models can be better with perfect balanced data

Extension to larger models

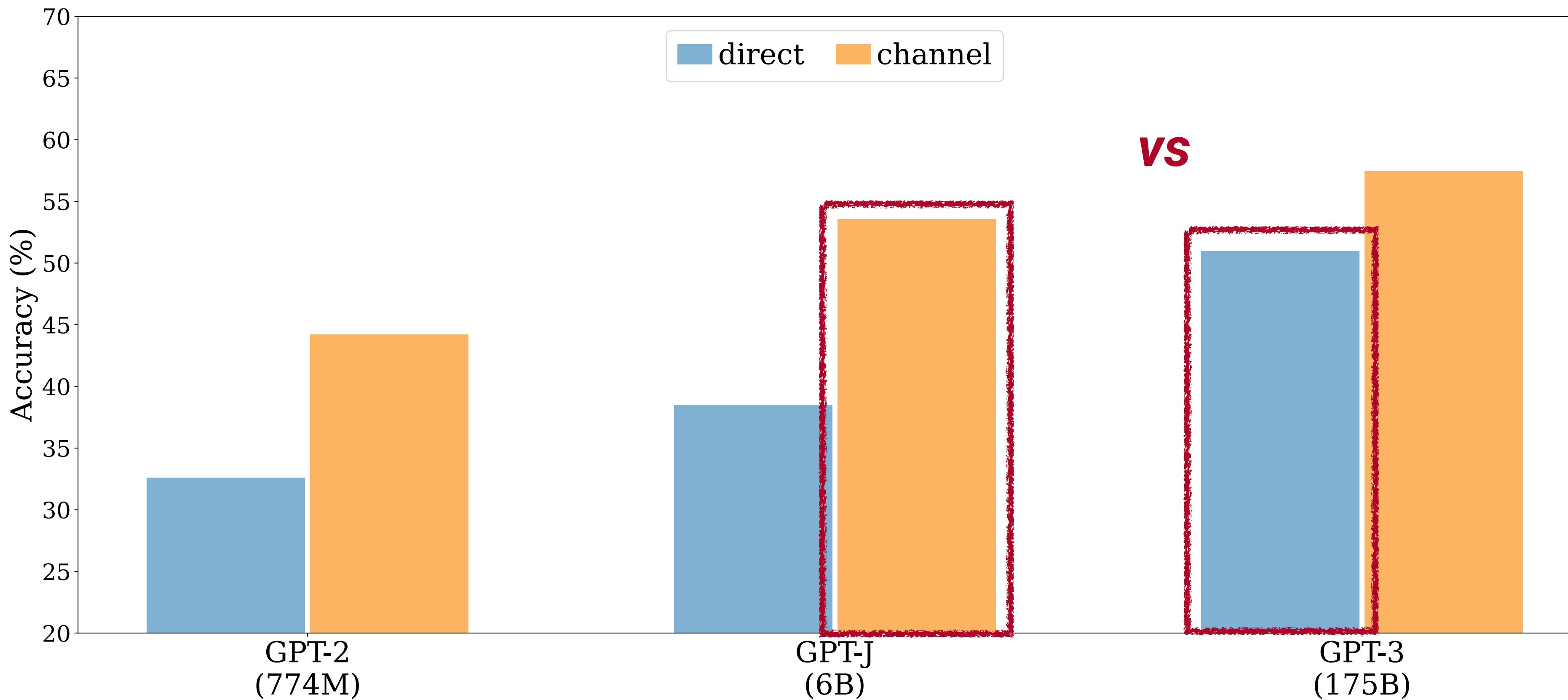
*Channel outperforms direct across **774M to 175B** models*



(on three classification datasets)

Extension to larger models

*Channel outperforms direct across **774M to 175B** models*
*Channel method allows outperforming **x30** bigger models*



(on three classification datasets)

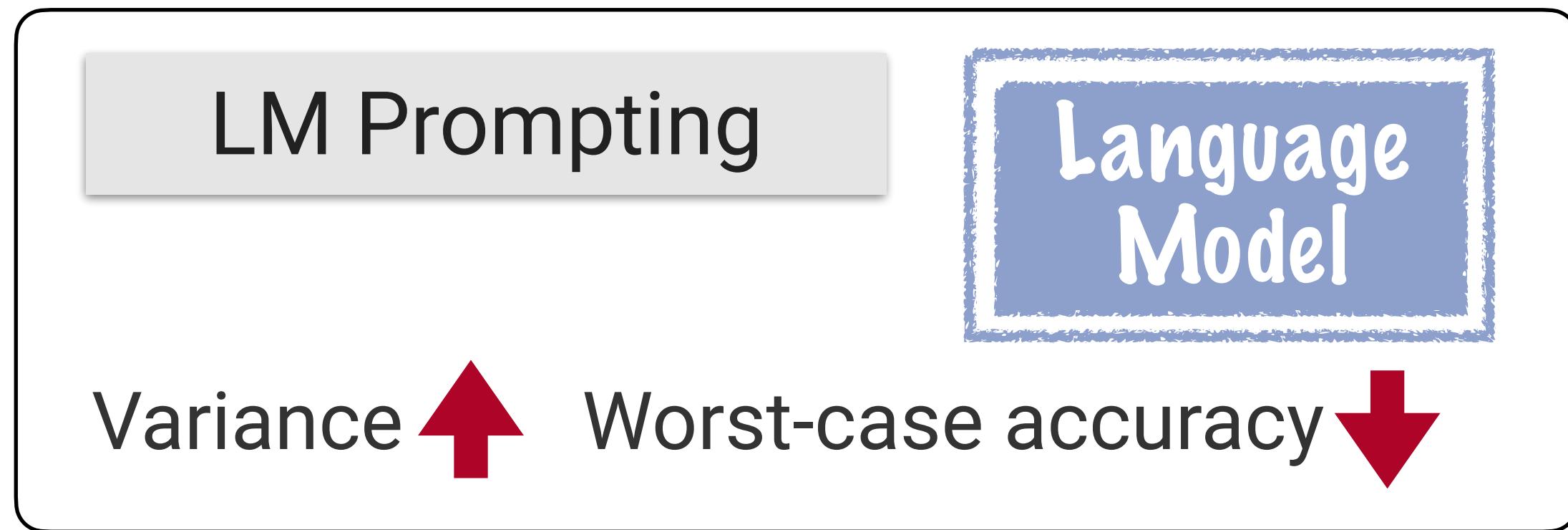
Takeaways

Takeaways

LM Prompting

Language
Model

Takeaways



Takeaways

LM Prompting

Language
Model

Variance 

Worst-case accuracy 

Noisy Channel Approach

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} \propto P(x|y)P(y)$$

Takeaways

LM Prompting

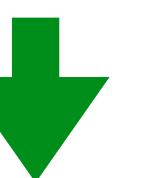
Language
Model

Variance 

Worst-case accuracy 

Noisy Channel Approach

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} \propto P(x|y)P(y)$$

Variance  Worst-case accuracy 

on 10 classification datasets

Takeaways

LM Prompting

Language
Model

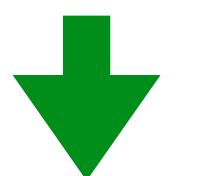
Variance 

Worst-case accuracy 

When is it better?

Noisy Channel Approach

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} \propto P(x|y)P(y)$$

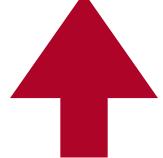
Variance  Worst-case accuracy 

on 10 classification datasets

Takeaways

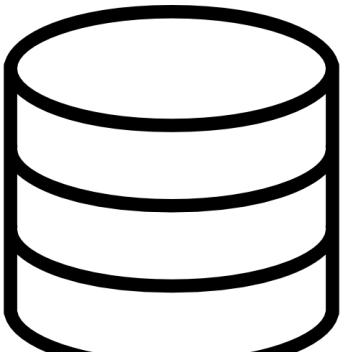
LM Prompting

Language
Model

Variance 

Worst-case accuracy 

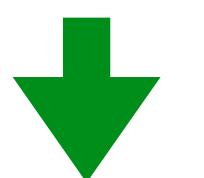
When is it better?



Small k

Noisy Channel Approach

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} \propto P(x|y)P(y)$$

Variance  Worst-case accuracy 

on 10 classification datasets

Takeaways

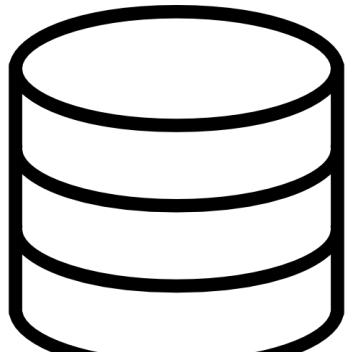
LM Prompting

Language
Model

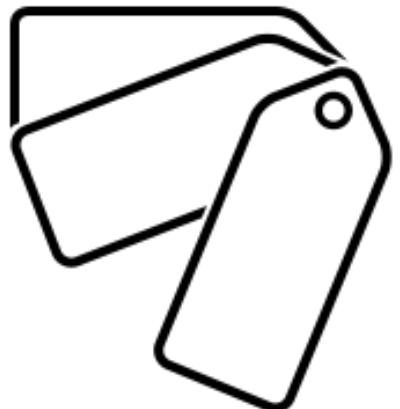
Variance 

Worst-case accuracy 

When is it better?



Small k



*Large
of labels*

Noisy Channel Approach

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} \propto P(x|y)P(y)$$

Variance  Worst-case accuracy 

on 10 classification datasets

Takeaways

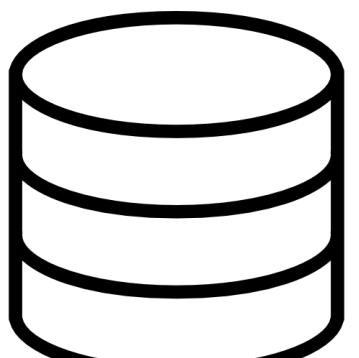
LM Prompting

Language
Model

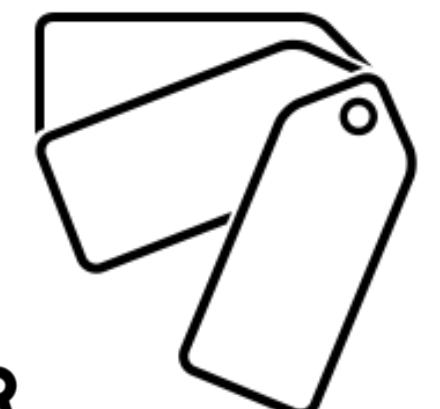
Variance 

Worst-case accuracy 

When is it better?



Small k



*Large
of labels*



Imbalanced data

Noisy Channel Approach

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} \propto P(x|y)P(y)$$

Variance  Worst-case accuracy 
on 10 classification datasets

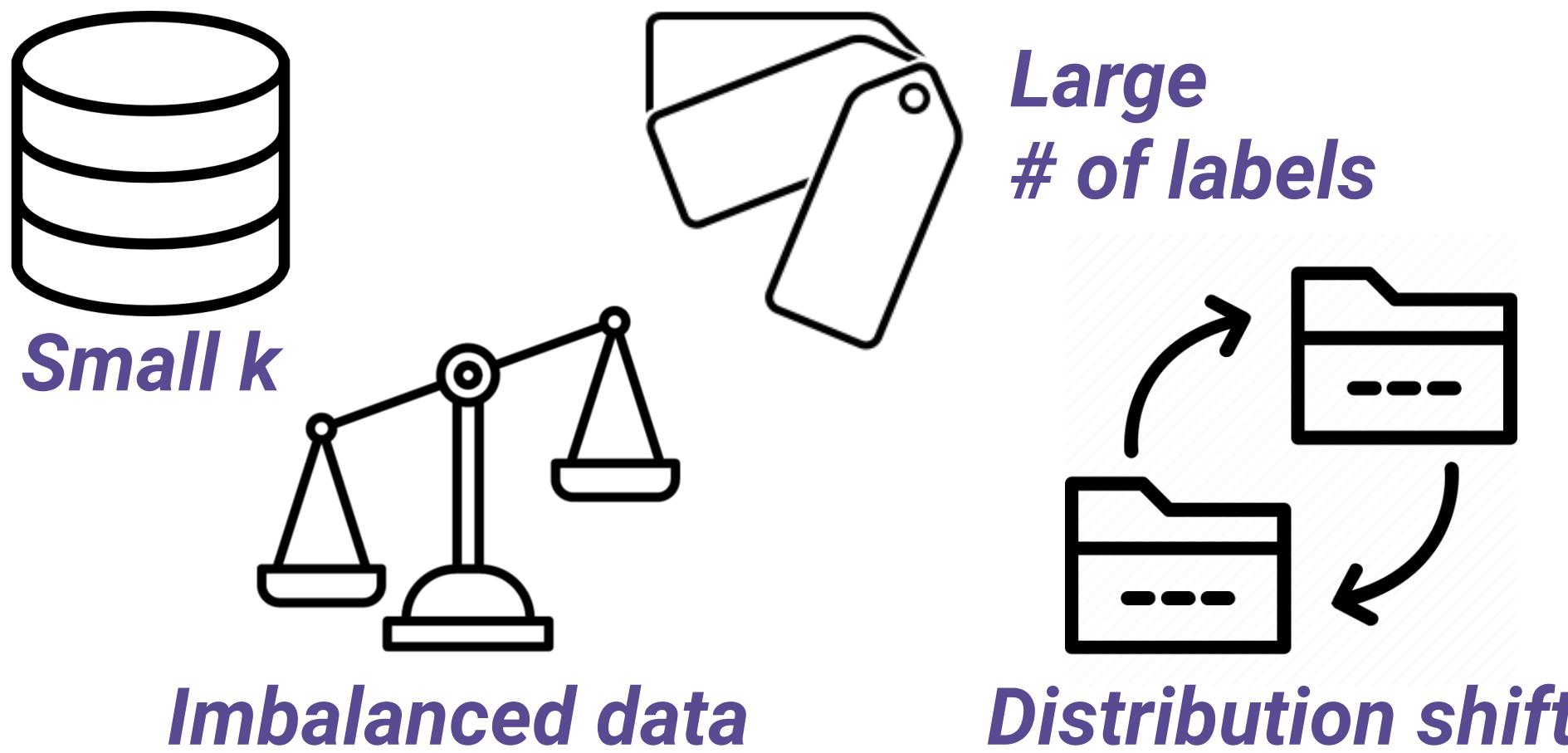
Takeaways

LM Prompting

Language
Model

Variance  Worst-case accuracy 

When is it better?



Noisy Channel Approach

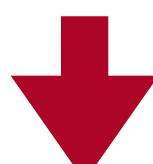
$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} \propto P(x|y)P(y)$$

Variance  Worst-case accuracy 
on 10 classification datasets

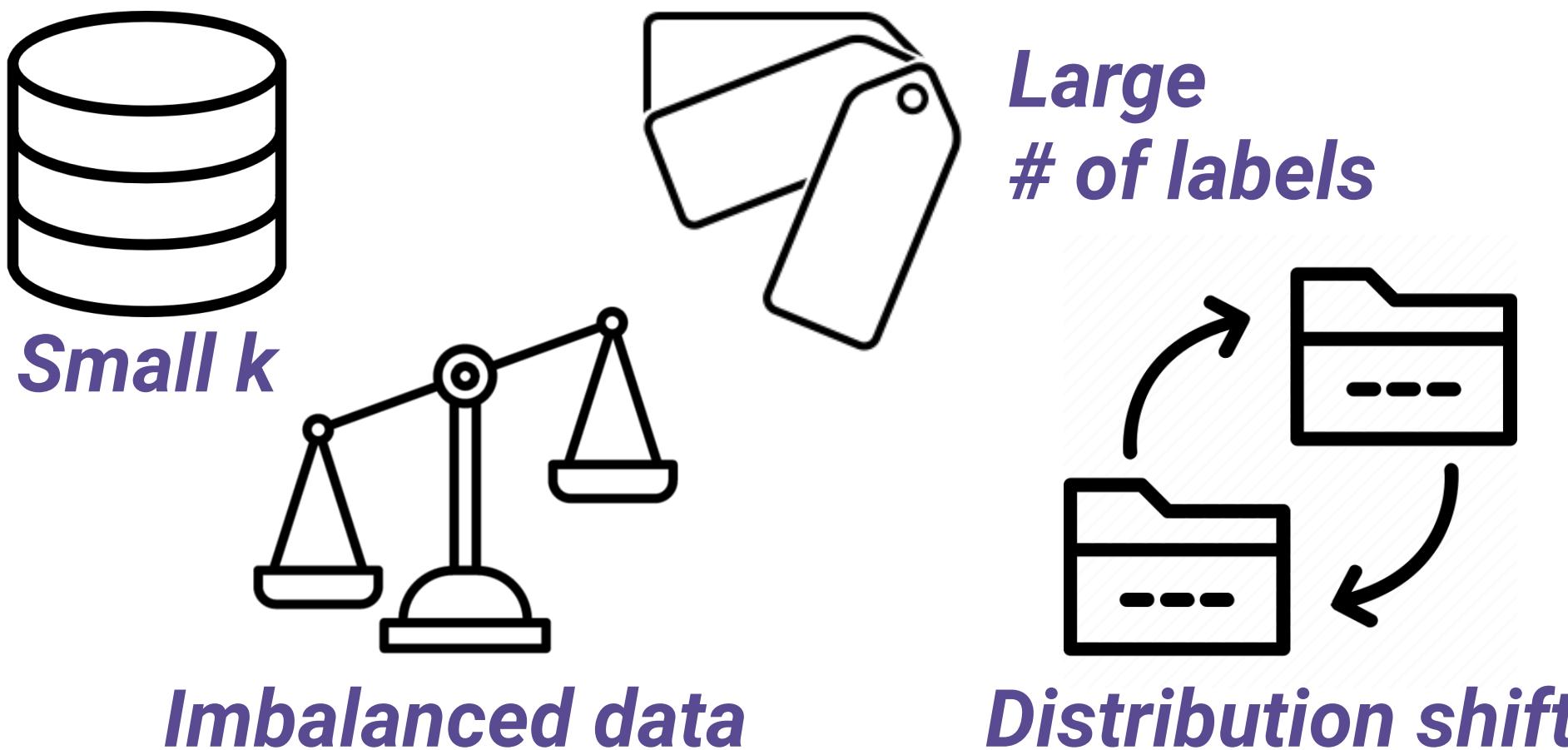
Takeaways

LM Prompting

Language Model

Variance  Worst-case accuracy 

When is it better?



Noisy Channel Approach

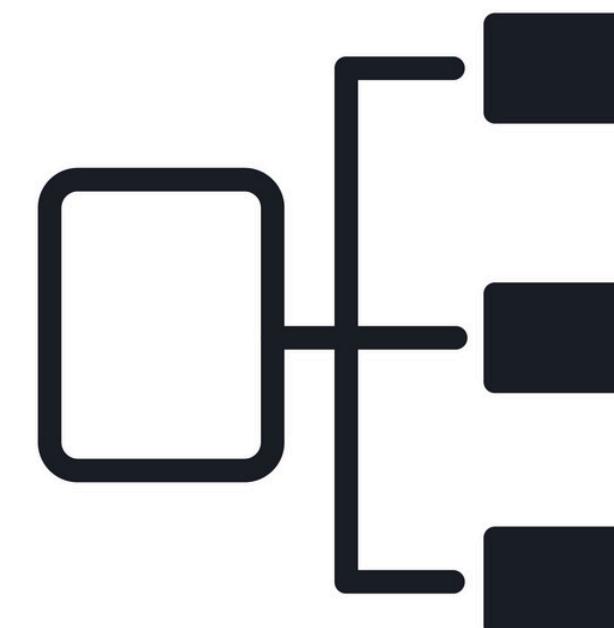
$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} \propto P(x|y)P(y)$$

Variance  Worst-case accuracy 
on 10 classification datasets

Limitations

Limited to classification

Extension needed!

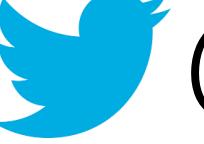


Thank you for listening

Paper: <https://arxiv.org/abs/2108.04106>

Code: <https://github.com/shmsw25/Channel-LM-Prompting/>

Demo: <http://qa.cs.washington.edu:2021/>

Contact:  sewon@cs.washington.edu /  [@sewon_min](https://twitter.com/sewon_min)