

# MetalCL: Learning to Learn In Context

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



# In-Context Learning

# In-Context Learning

## *k shot data*

Input: Circulation revenue has increased by 5% in Finland. Output : Positive

Input: Panostaja did not disclose the purchase price. Output: Neutral

Input: Paying off the national debt will be extremely painful. Output: Negative

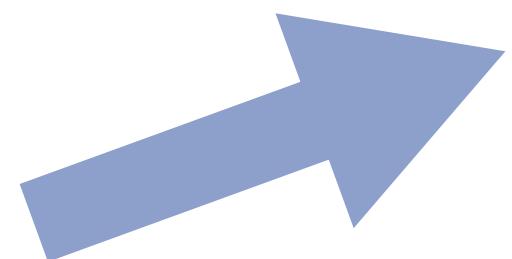
# In-Context Learning

*k* **shot data**

Input: Circulation revenue has increased by 5% in Finland. Output : Positive

Input: Panostaja did not disclose the purchase price. Output: Neutral

Input: Paying off the national debt will be extremely painful. Output: Negative



**Concatenate**

$x =$  Circulation revenue has increased by 5% in Finland. // Positive

Panostaja did not disclose the purchase price. // Neutral

Paying off the national debt will be extremely painful. // Negative

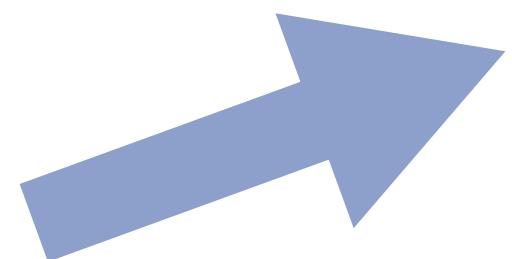
# In-Context Learning

*k* **shot data**

Input: Circulation revenue has increased by 5% in Finland. Output : Positive

Input: Panostaja did not disclose the purchase price. Output: Neutral

Input: Paying off the national debt will be extremely painful. Output: Negative



**Concatenate**

$x =$  Circulation revenue has increased by 5% in Finland. // Positive

Panostaja did not disclose the purchase price. // Neutral

Paying off the national debt will be extremely painful. // Negative

The company anticipated its operating profit to improve. //

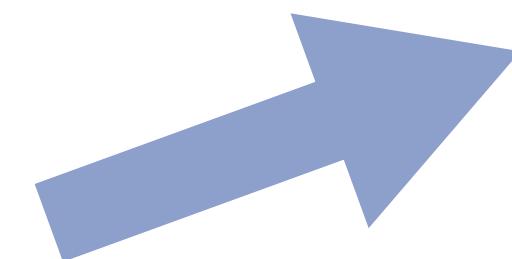
# In-Context Learning

*k* **shot data**

Input: Circulation revenue has increased by 5% in Finland. Output : Positive

Input: Panostaja did not disclose the purchase price. Output: Neutral

Input: Paying off the national debt will be extremely painful. Output: Negative



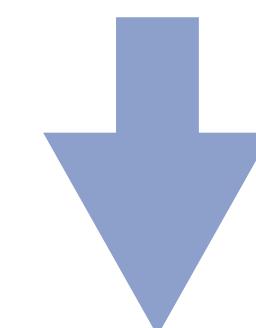
**Concatenate**

$x =$  Circulation revenue has increased by 5% in Finland. // Positive

Panostaja did not disclose the purchase price. // Neutral

Paying off the national debt will be extremely painful. // Negative

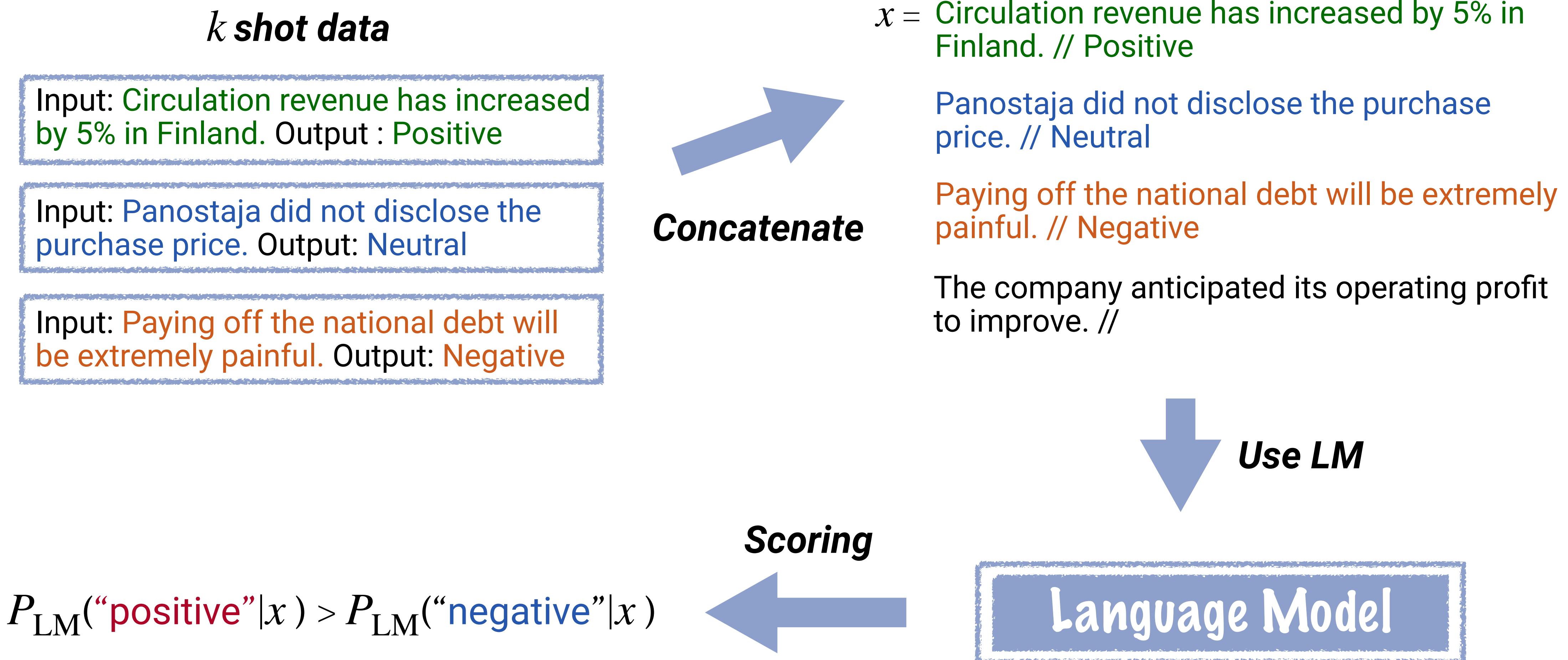
The company anticipated its operating profit to improve. //



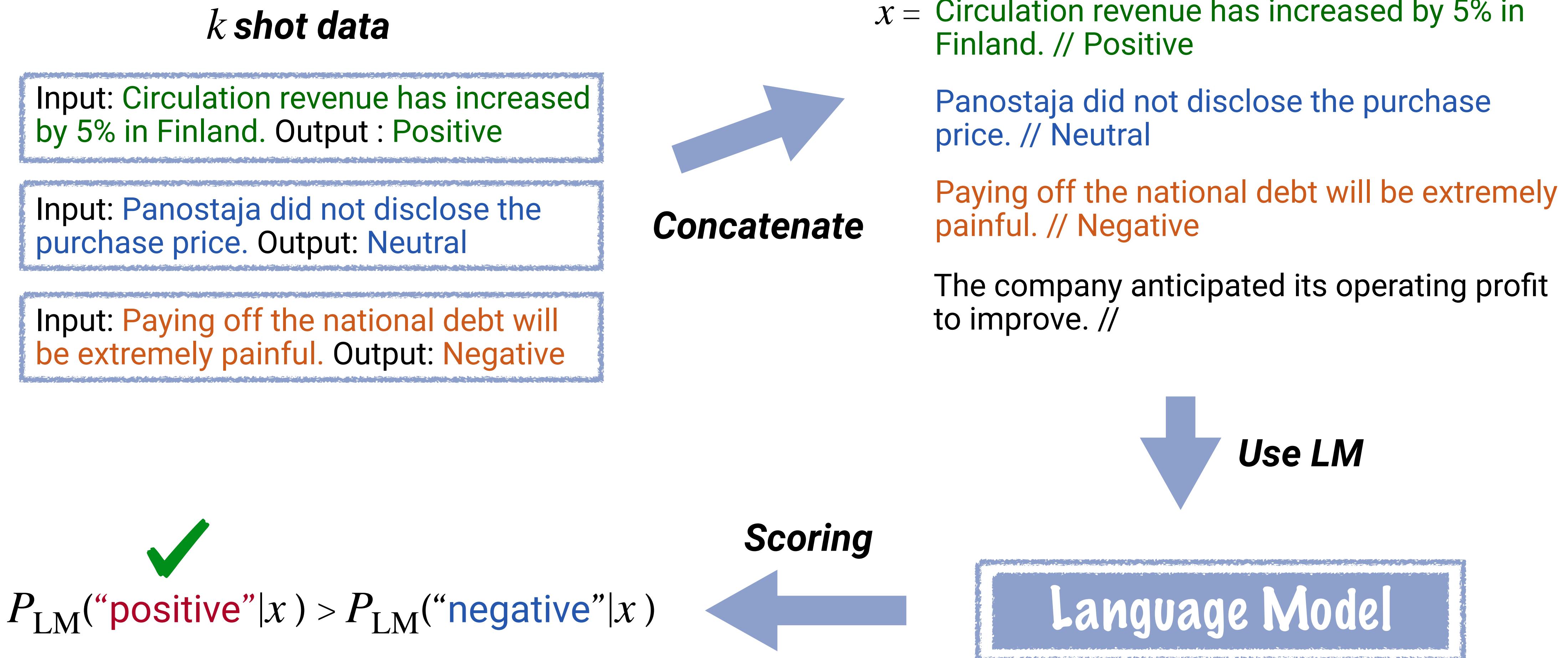
**Use LM**

Language Model

# In-Context Learning



# In-Context Learning

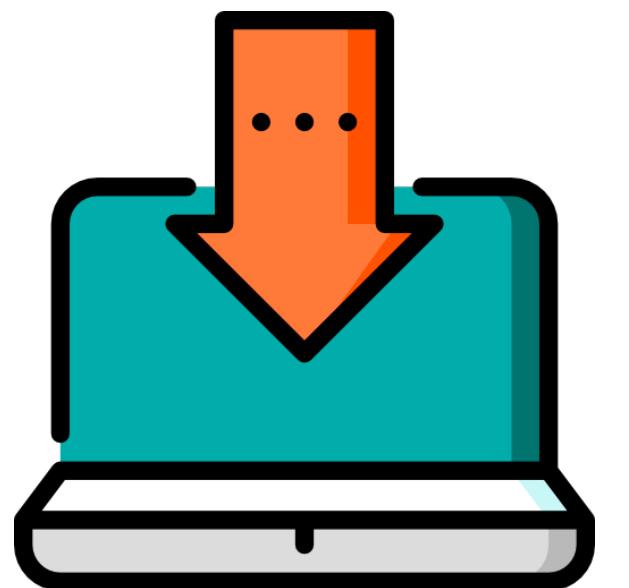




**In-context learning**  
(Brown et al. 2020)



**In-context learning**  
(Brown et al. 2020)



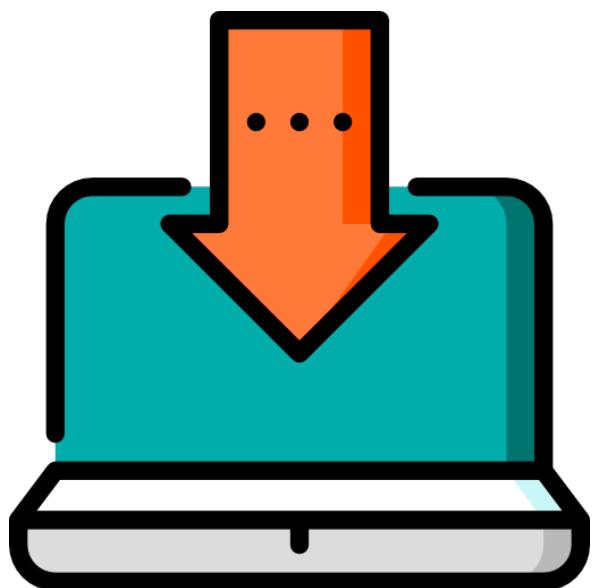
**Better tuning of prompt**  
(Shin et al. 2020, Jiang et al. 2020, Gao et al. 2021)



**In-context learning**  
(Brown et al. 2020)



**Better scoring**  
(Zhao et al. 2021, Holtzman et al. 2021)



**Better tuning of prompt**  
(Shin et al. 2020, Jiang et al. 2020, Gao et al. 2021)



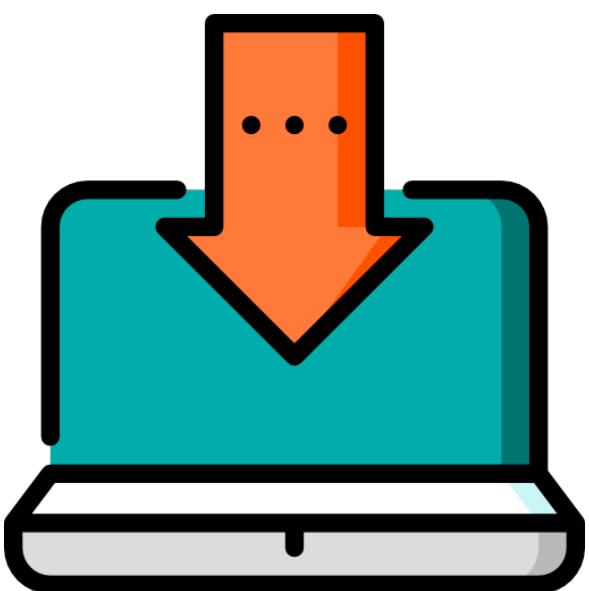
## In-context learning

(Brown et al. 2020)



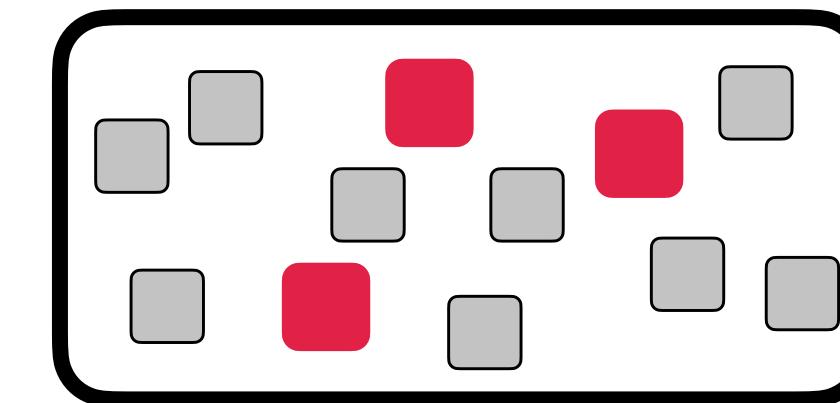
## Better scoring

(Zhao et al. 2021, Holtzman et al. 2021)



## Better tuning of prompt

(Shin et al. 2020, Jiang et al. 2020, Gao et al. 2021)



## Better choice of in-context examples

(Liu et al. 2021, Lu et al. 2021, Rubin et al. 2021)

*Next presentation!*

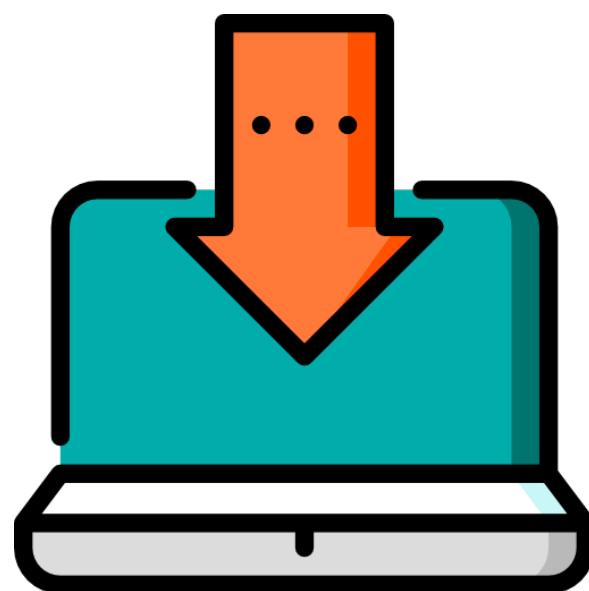
*We always used the plain language model*



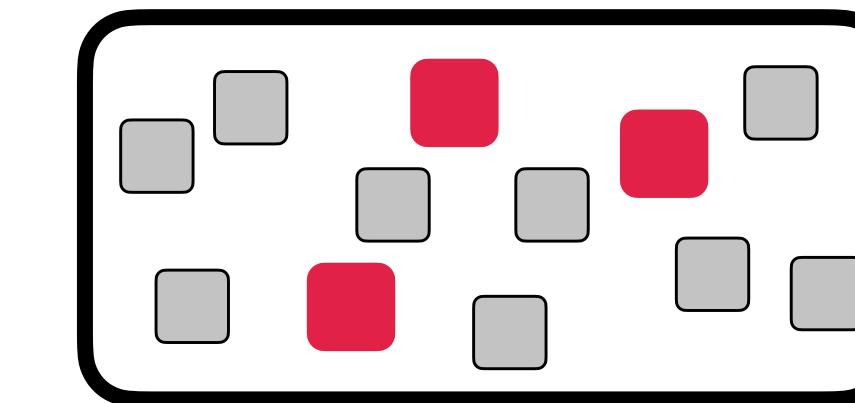
**In-context learning**  
(Brown et al. 2020)



**Better scoring**  
(Zhao et al. 2021, Holtzman et al. 2021)



**Better tuning of prompt**  
(Shin et al. 2020, Jiang et al. 2020, Gao et al. 2021)



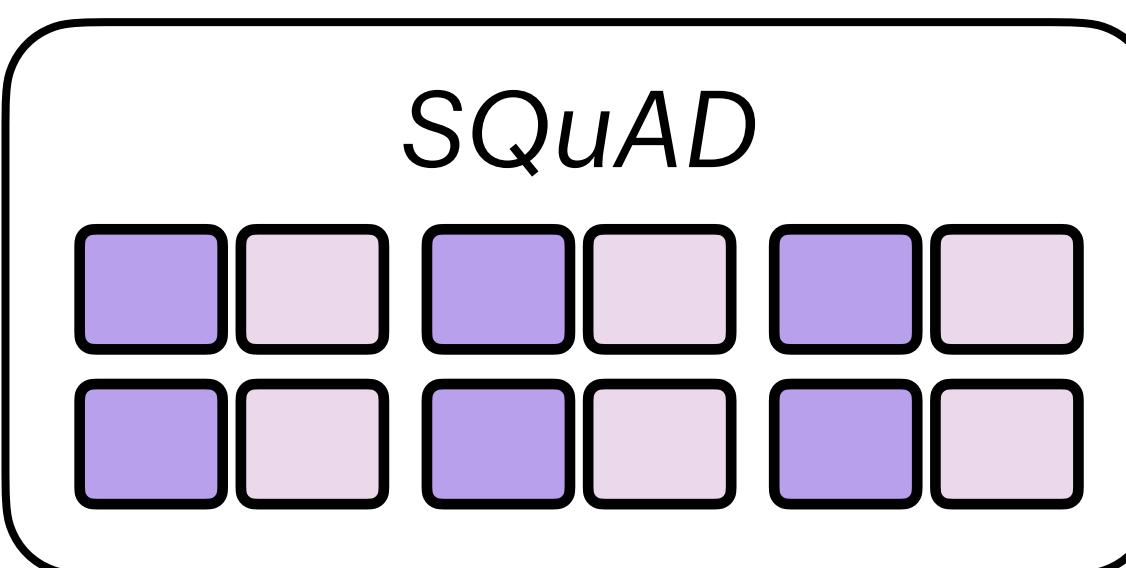
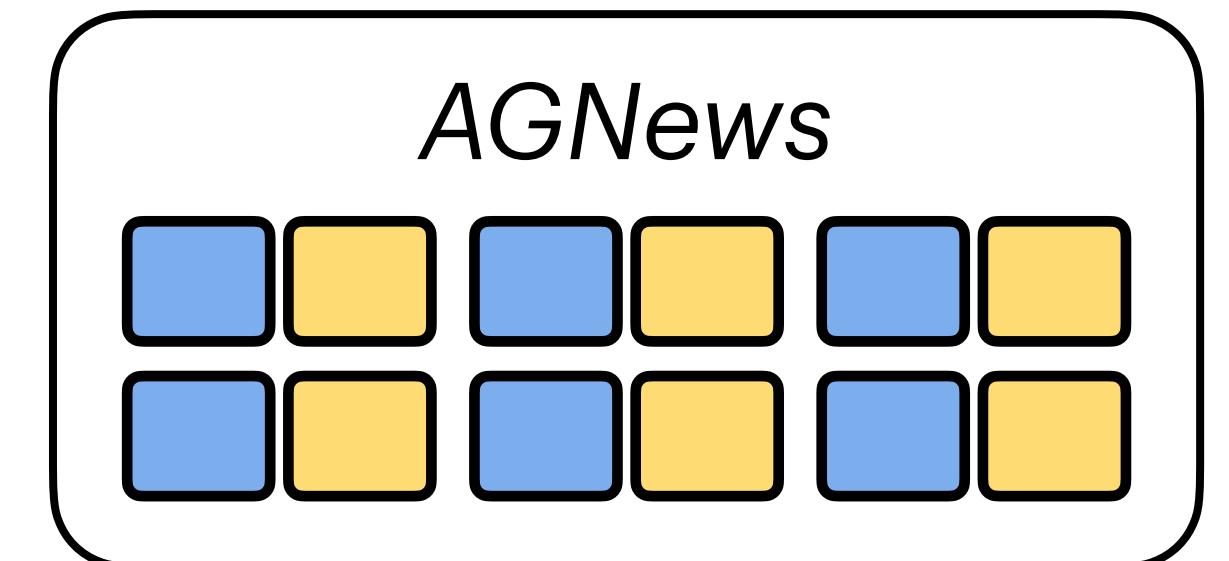
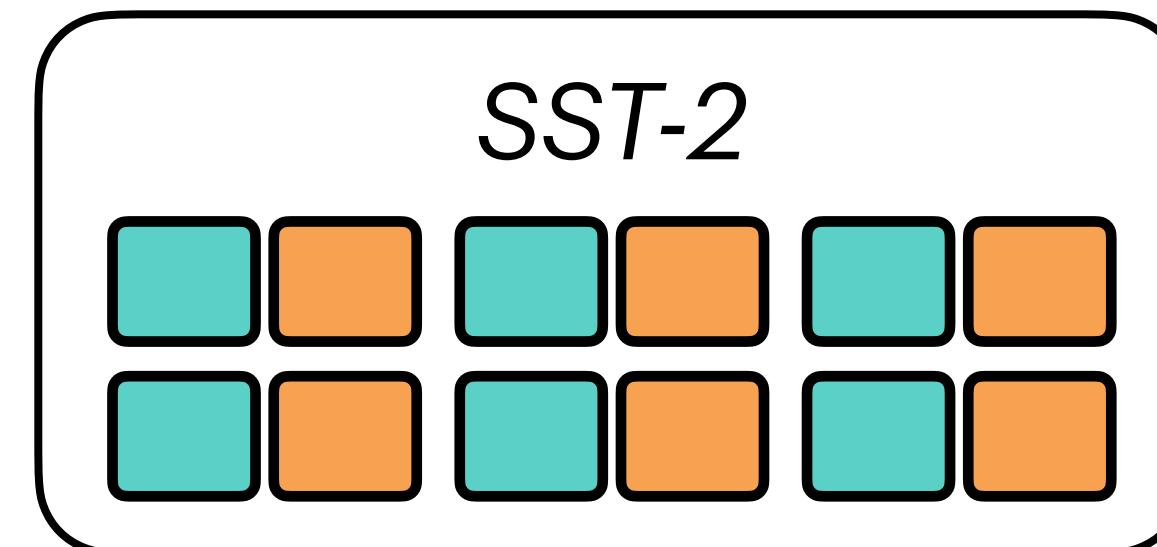
**Better choice of in-context examples**  
(Liu et al. 2021, Lu et al. 2021, Rubin et al. 2021)

*Next presentation!*

*We always used the plain language model*  
**“Can we *meta-train* the model to be a *better in-context learner*?”**

*We always used the plain language model*

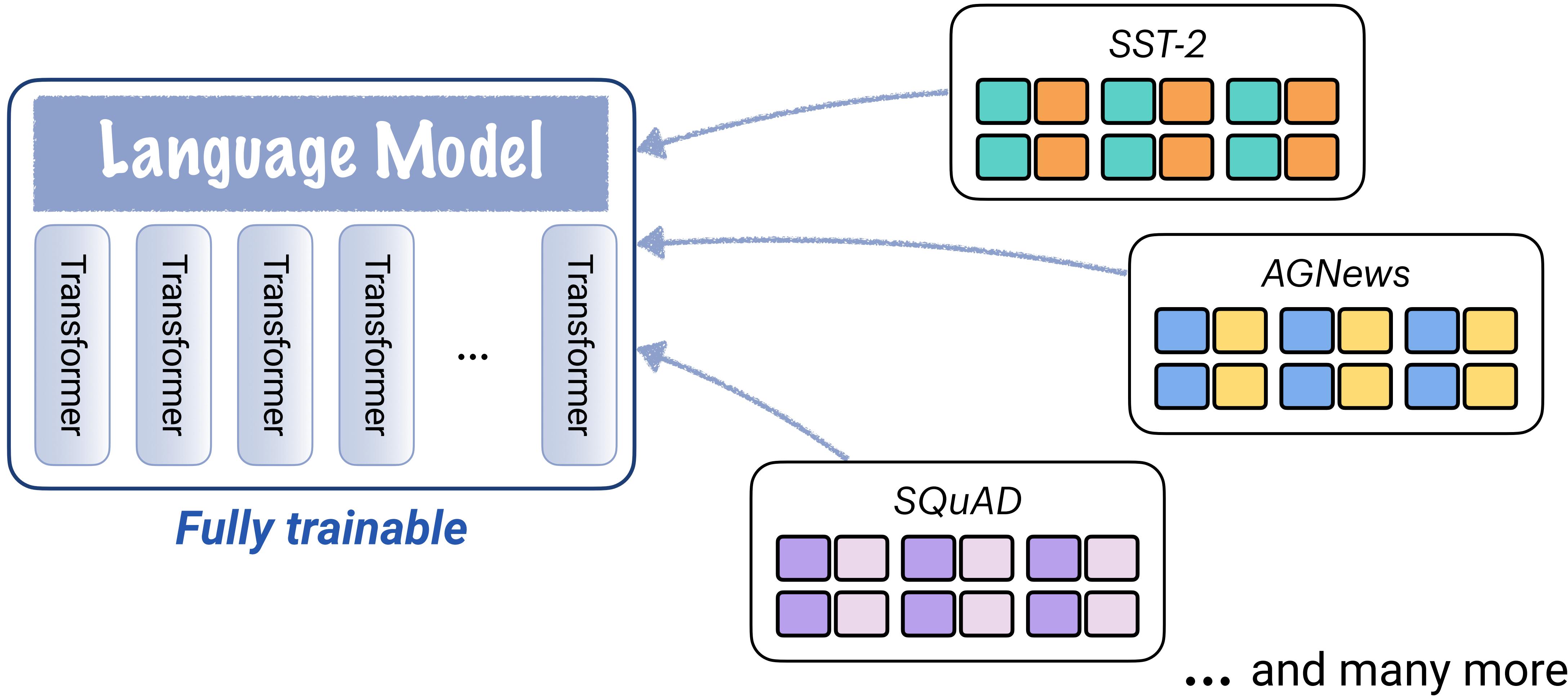
**“Can we *meta-train* the model to be a *better in-context learner*?”**



... and many more

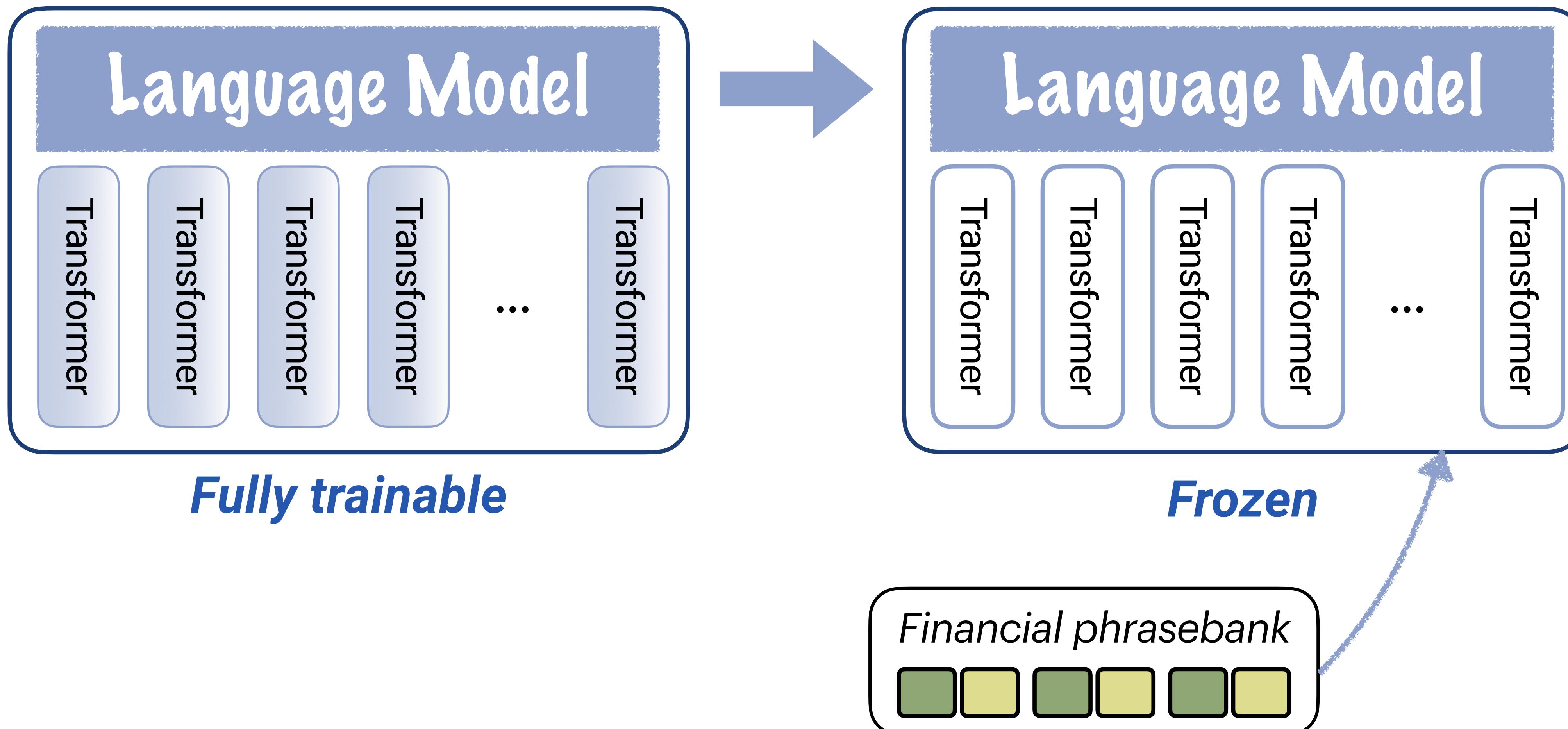
*We always used the plain language model*

**“Can we *meta-train* the model to be a better *in-context learner*?”**



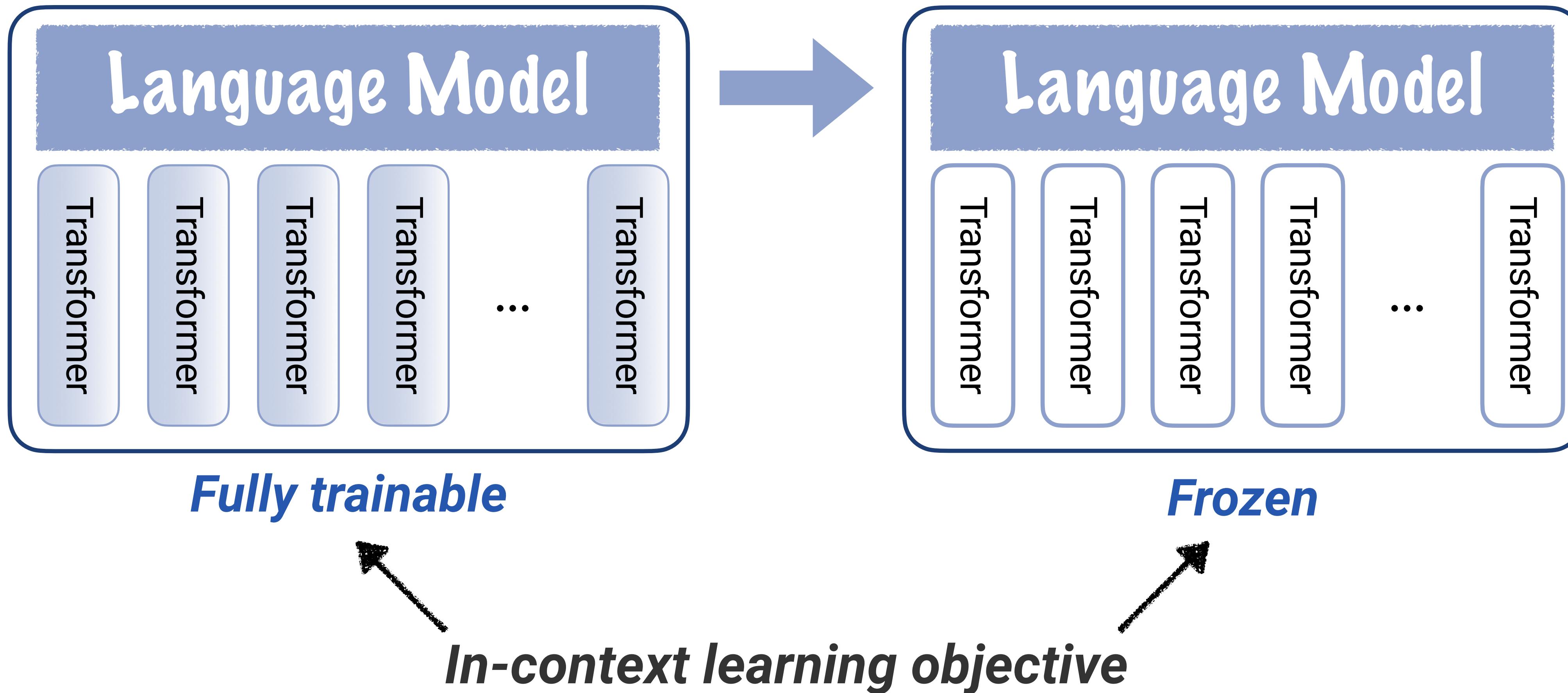
*“We always used the plain language model”*

# “Can we *meta-train* the model to be a *better in-context learner*? ”



*“We always used the plain language model”*

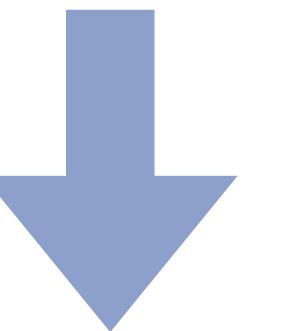
*“Can we **meta-train** the model to be a **better in-context learner**?”*



# MetalCL

## Meta-training

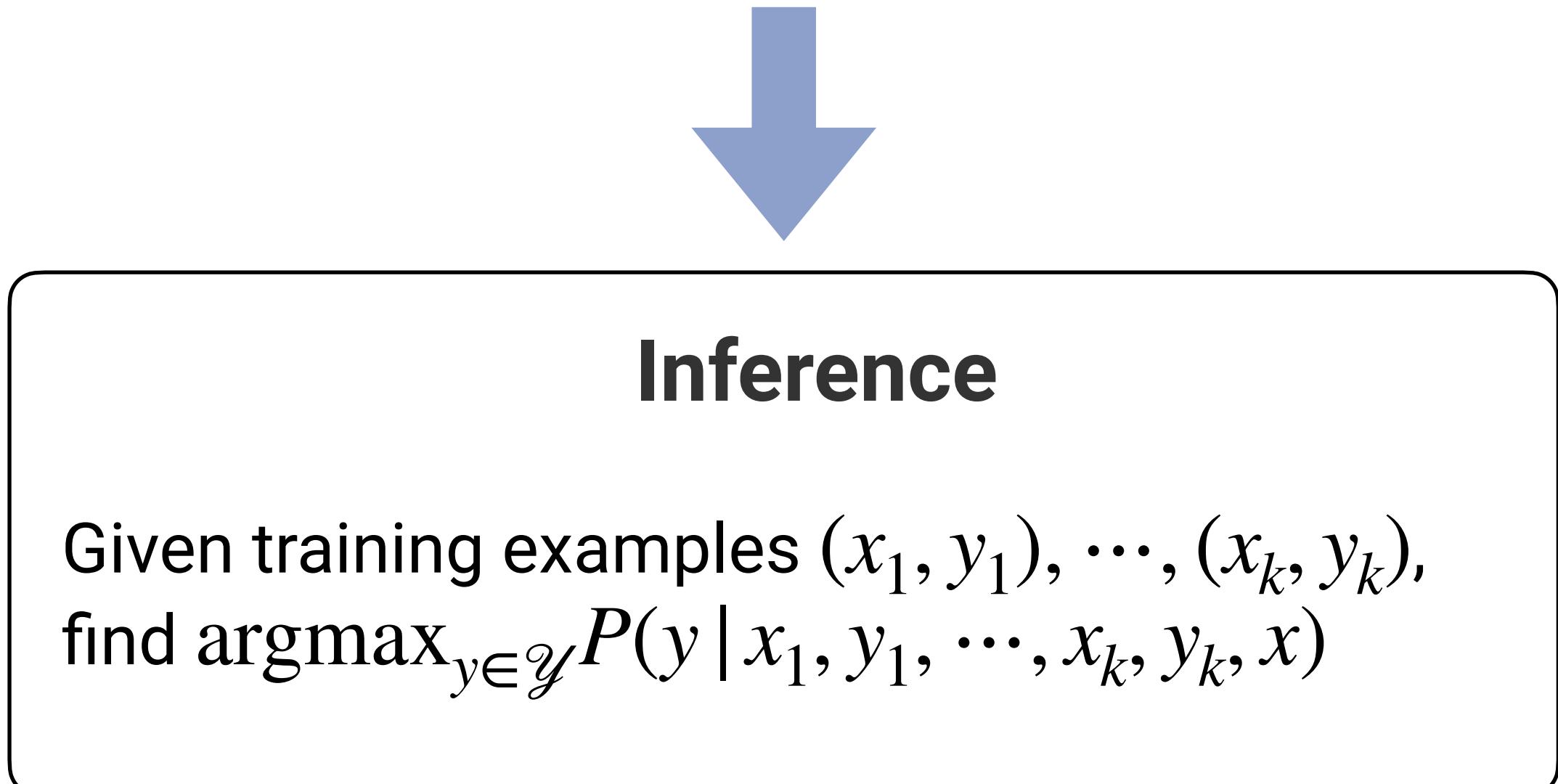
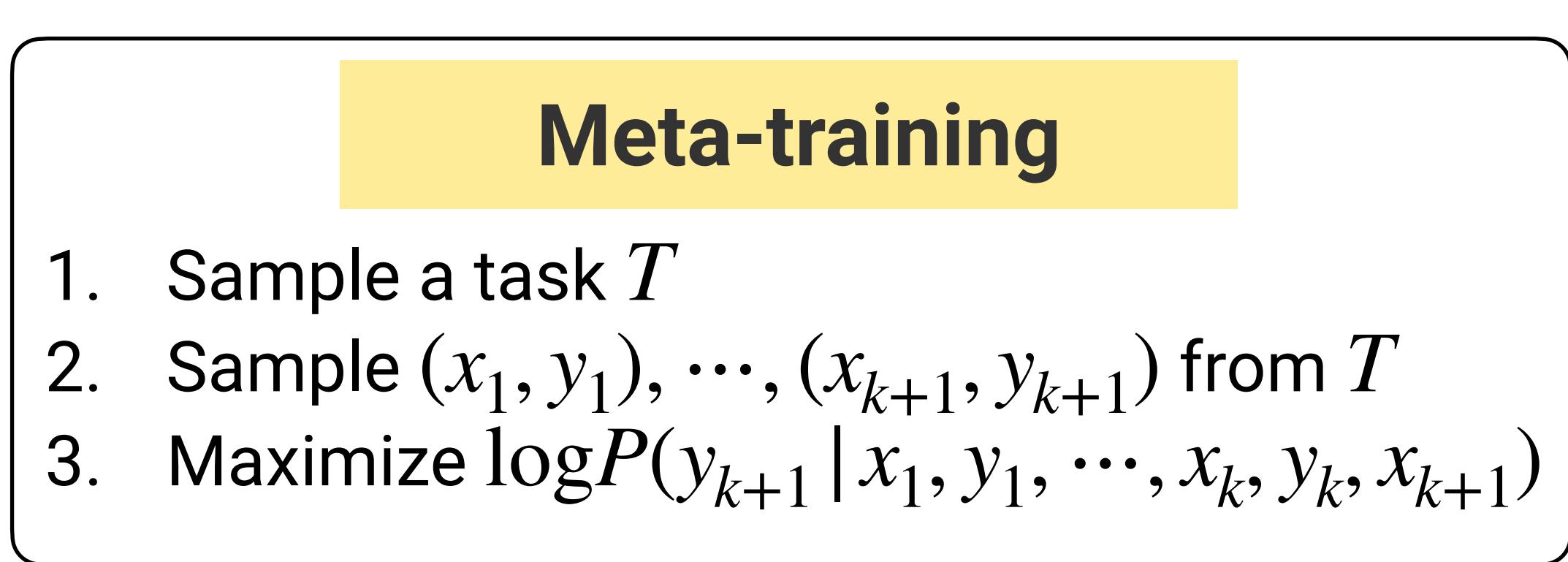
1. Sample a task  $T$
2. Sample  $(x_1, y_1), \dots, (x_{k+1}, y_{k+1})$  from  $T$
3. Maximize  $\log P(y_{k+1} | x_1, y_1, \dots, x_k, y_k, x_{k+1})$



## Inference

Given training examples  $(x_1, y_1), \dots, (x_k, y_k)$ ,  
find  $\text{argmax}_{y \in \mathcal{Y}} P(y | x_1, y_1, \dots, x_k, y_k, x)$

# MetalCL



A collection of meta-training tasks

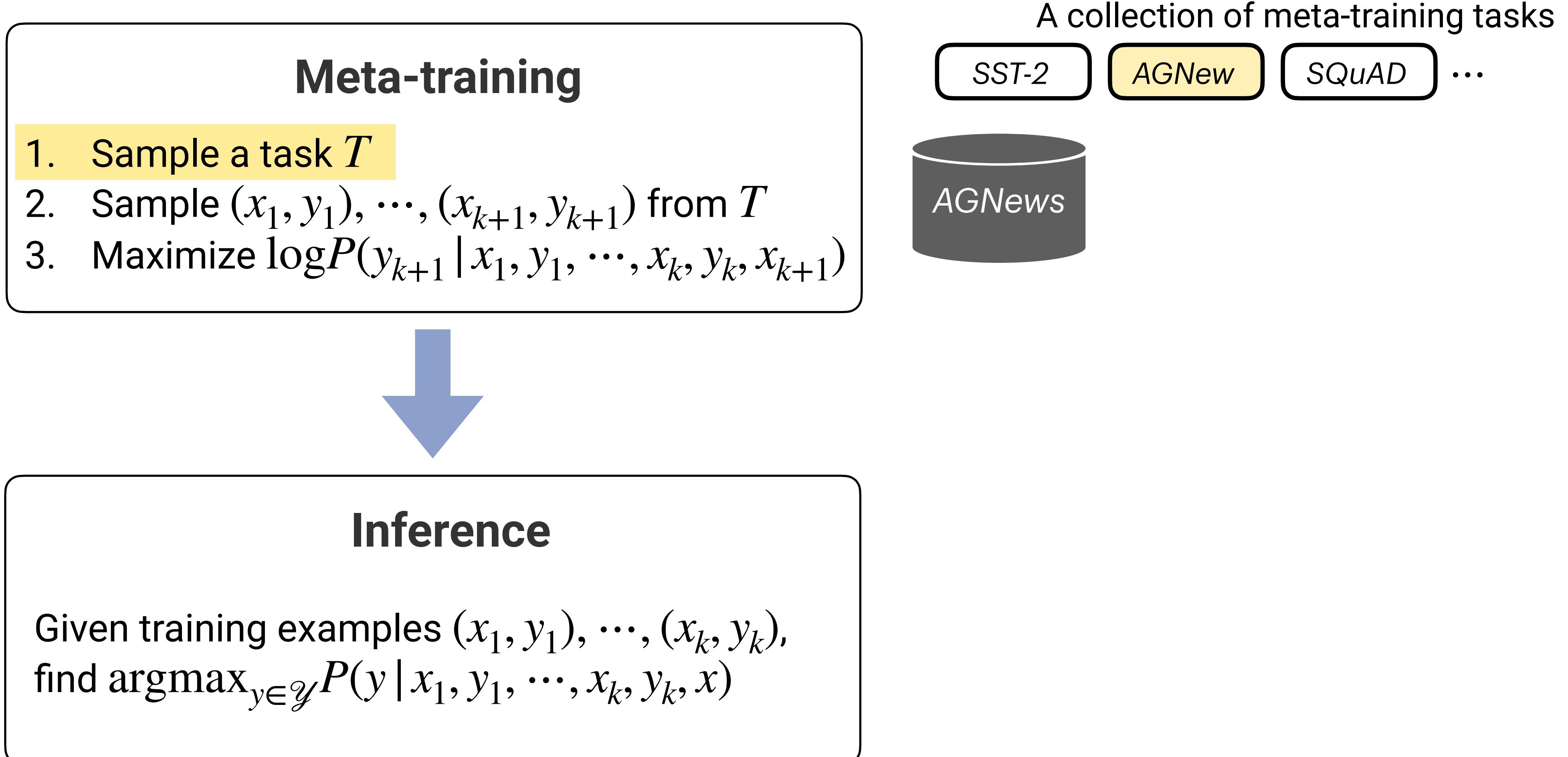
SST-2

AGNews

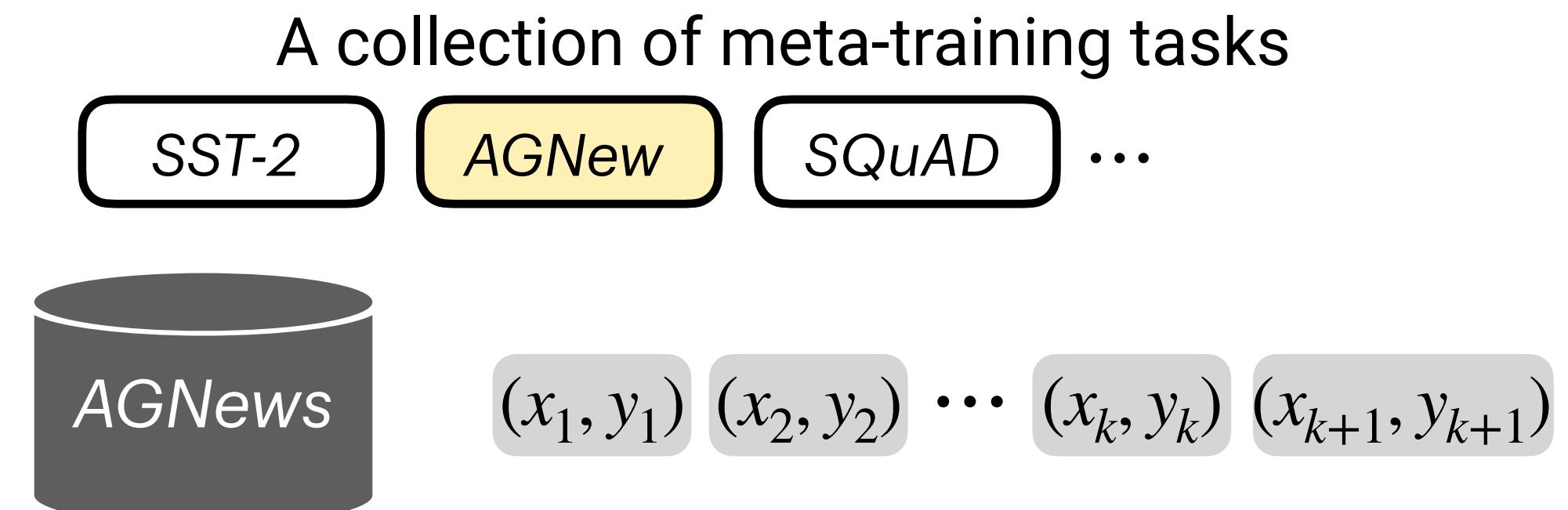
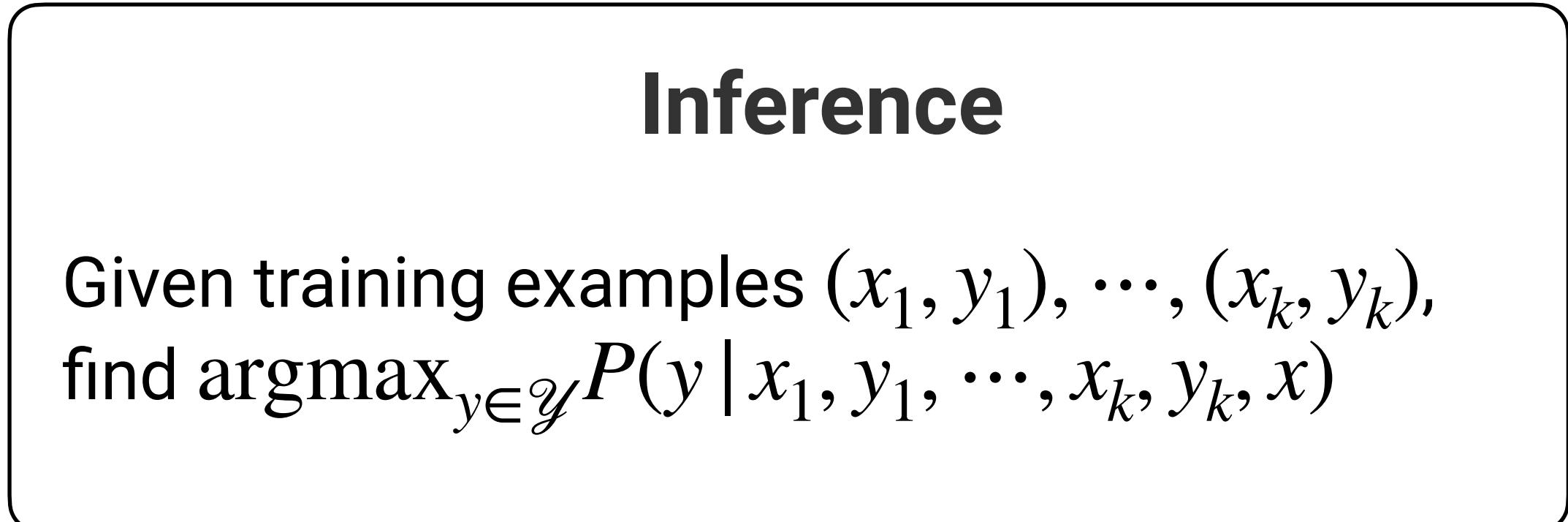
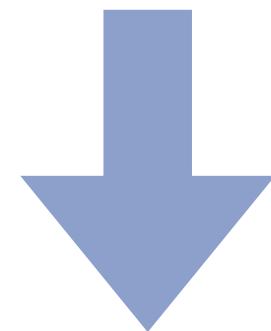
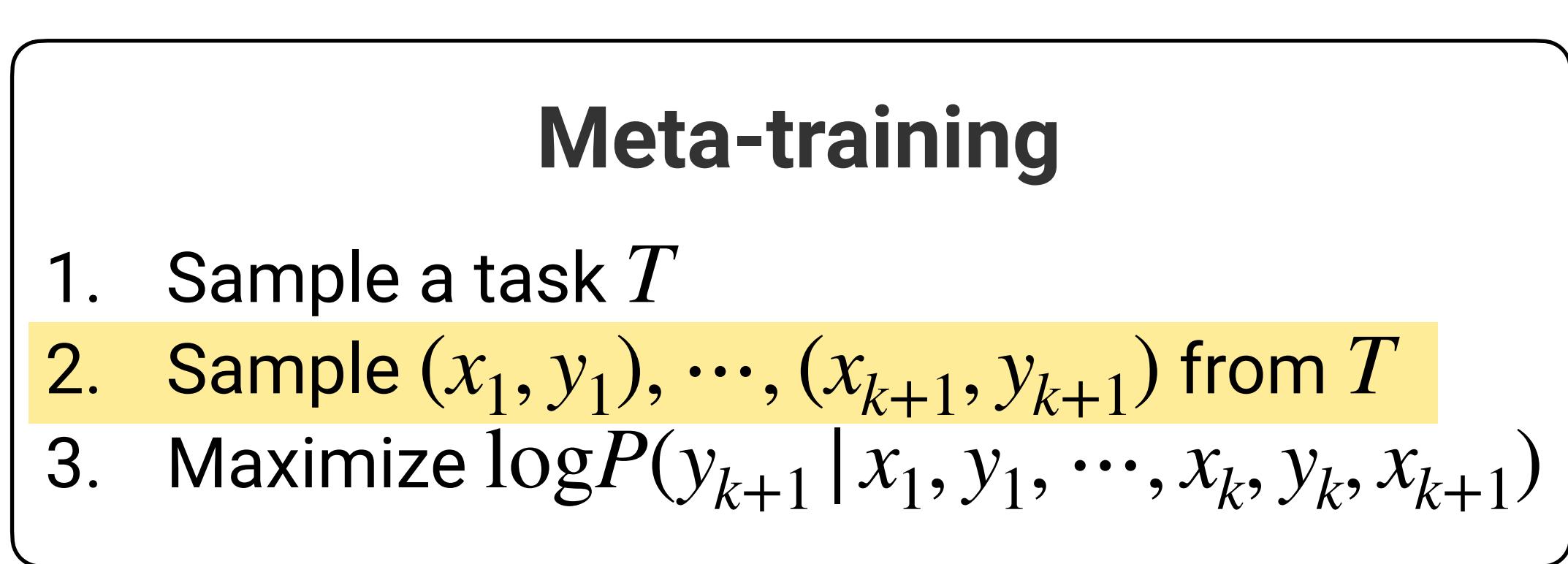
SQuAD

...

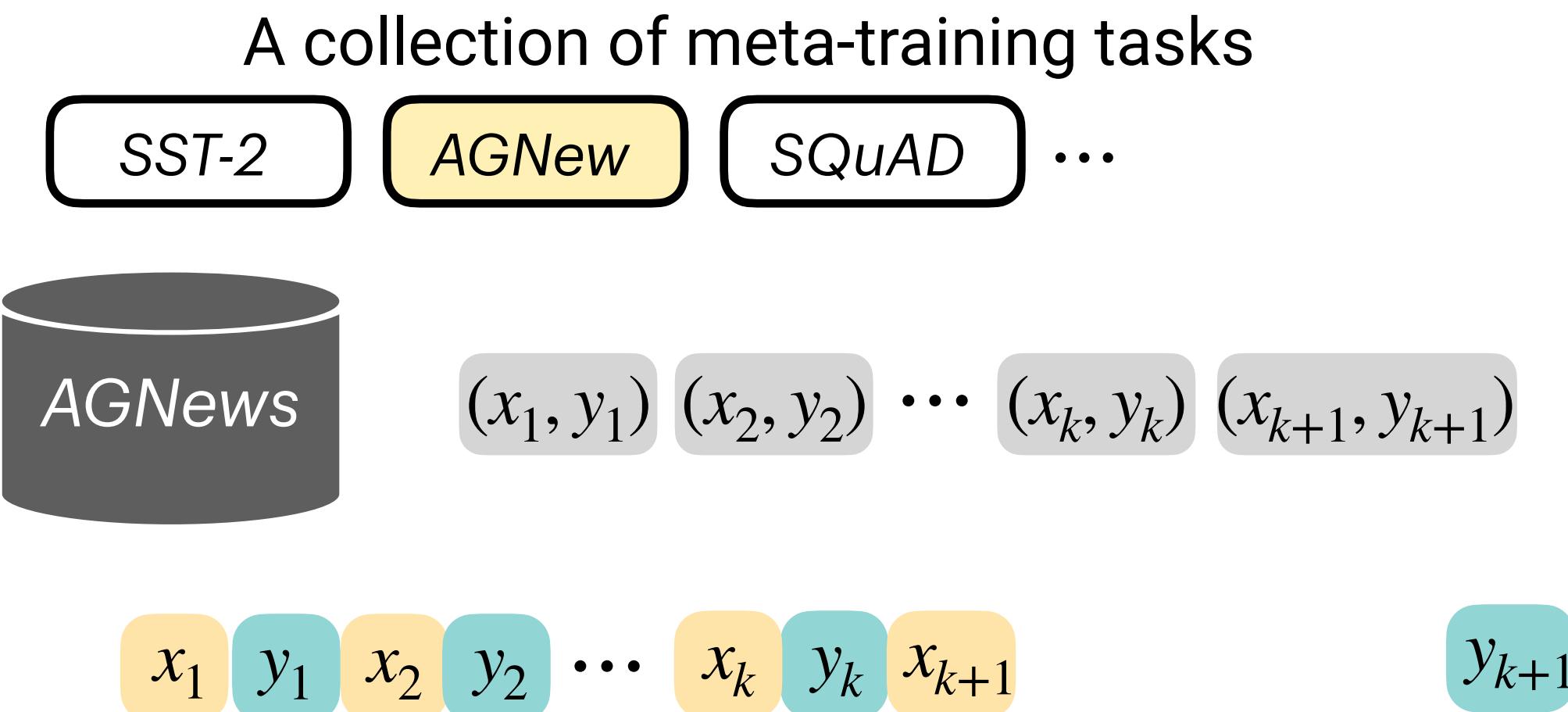
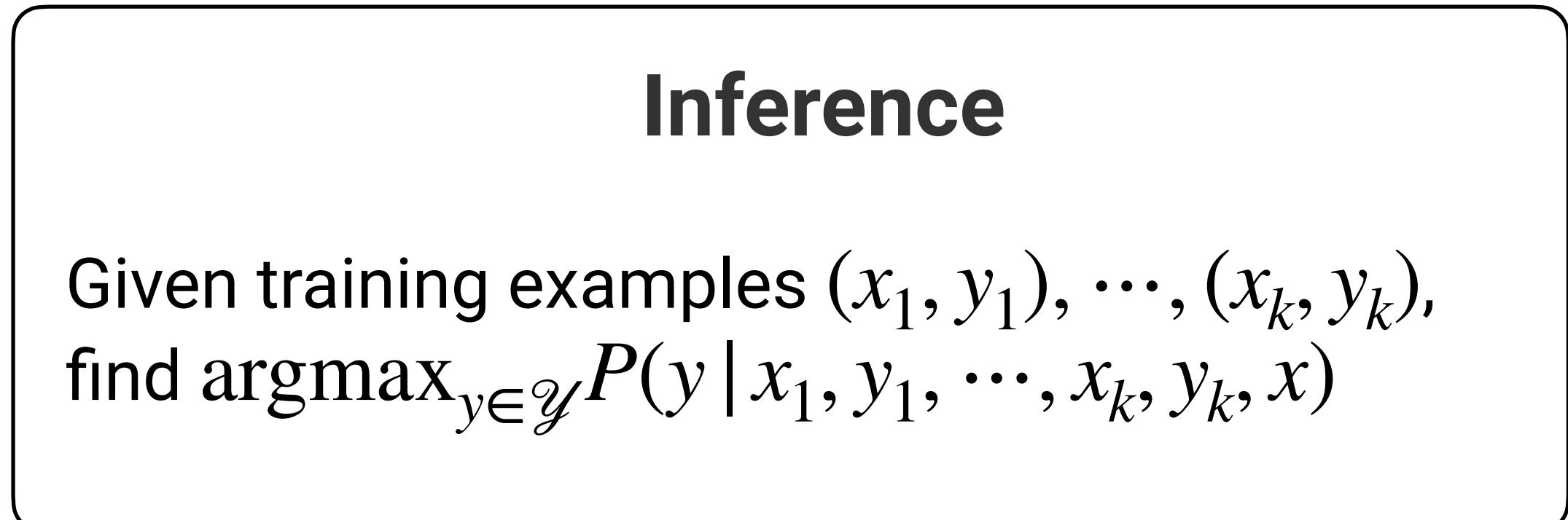
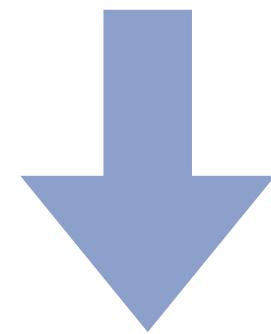
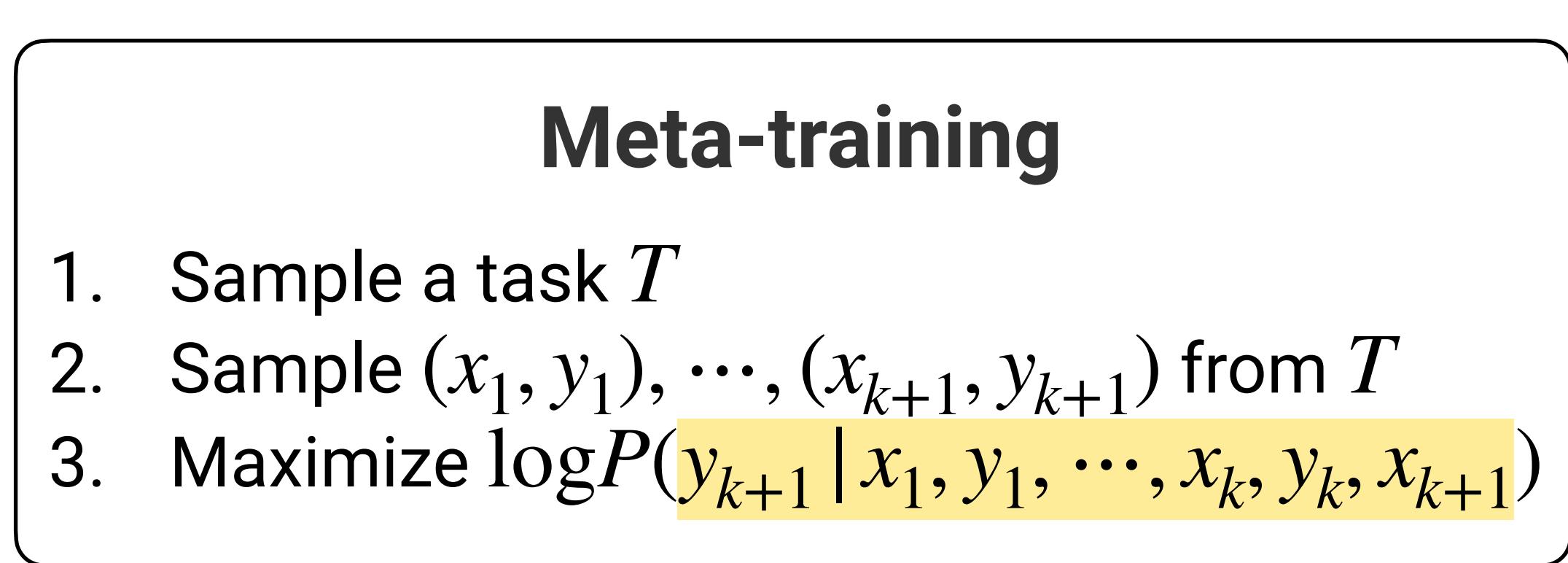
# MetalCL



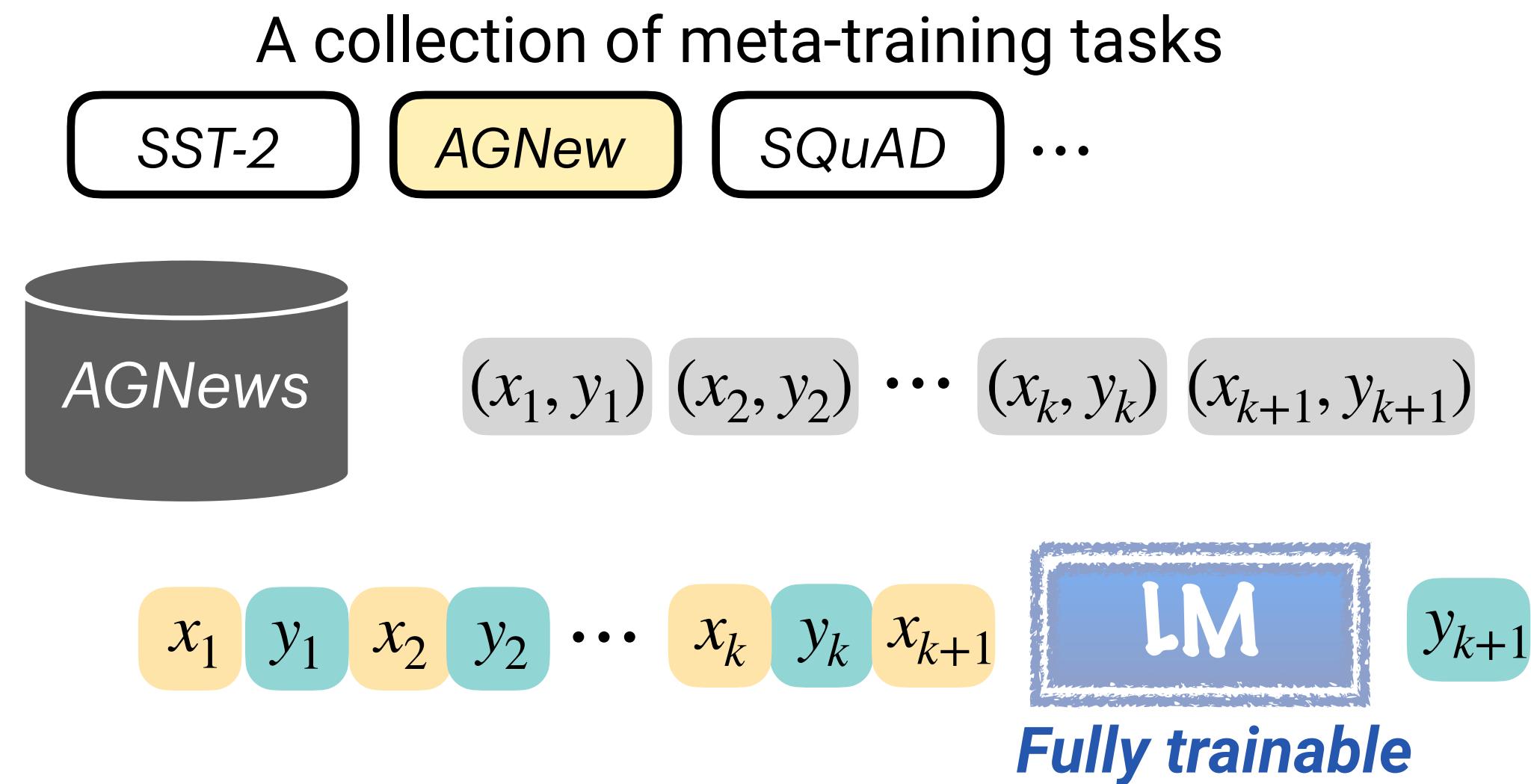
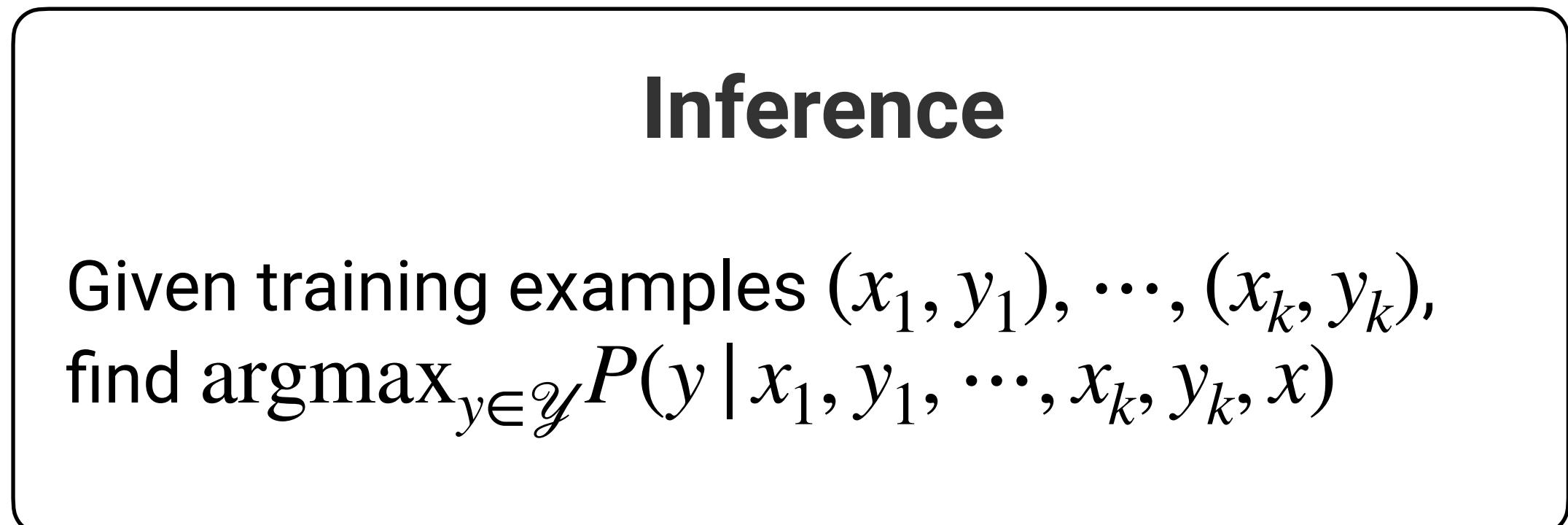
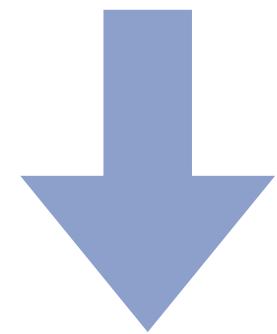
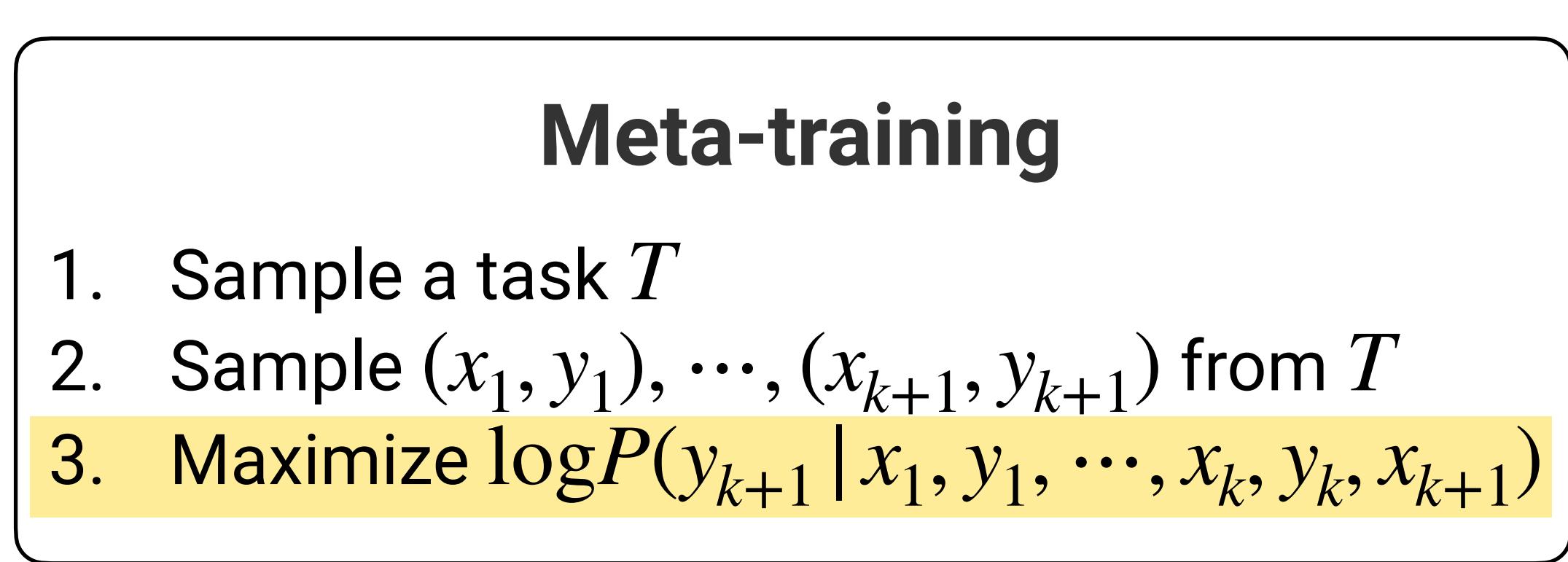
# MetalCL



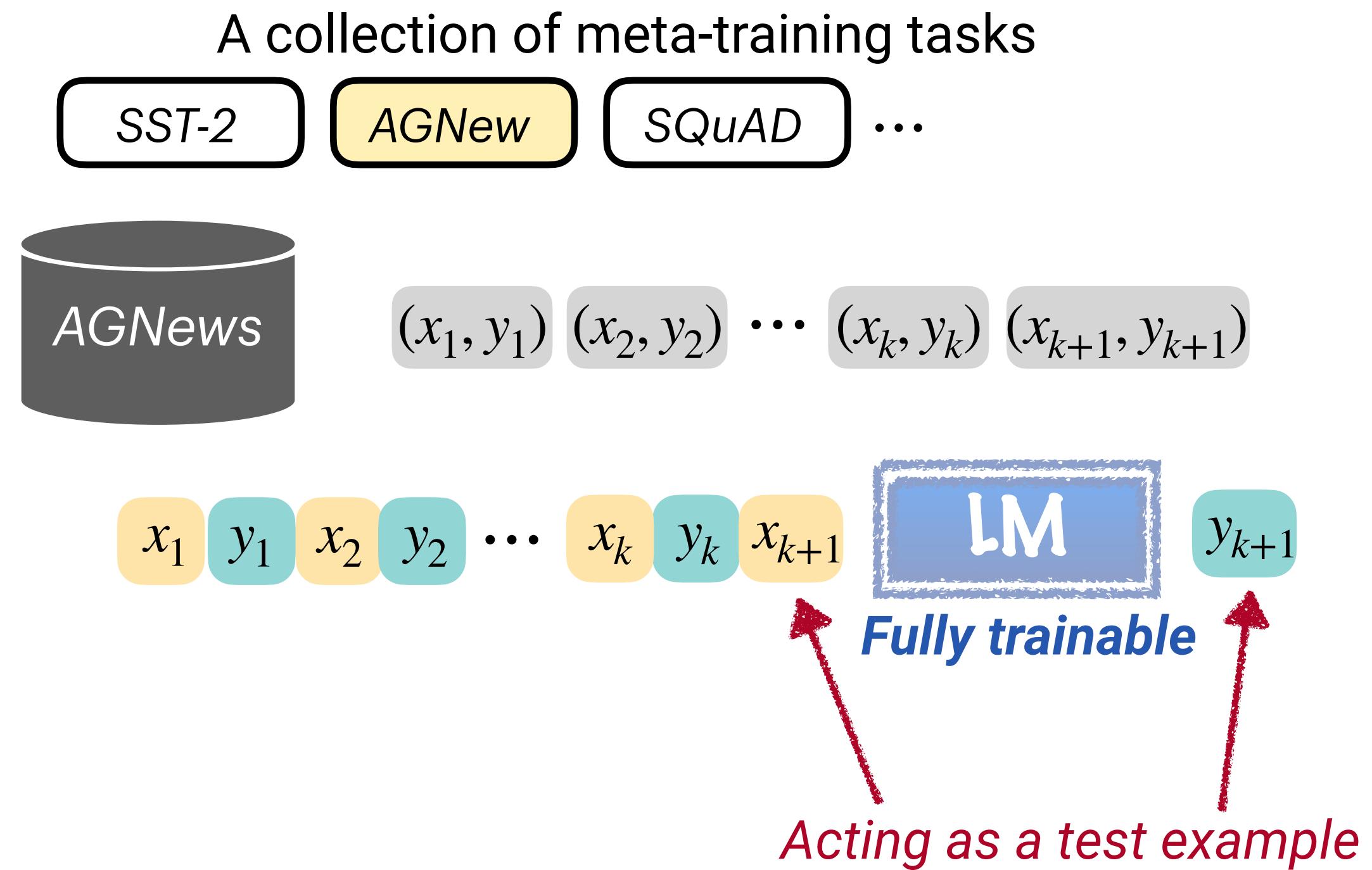
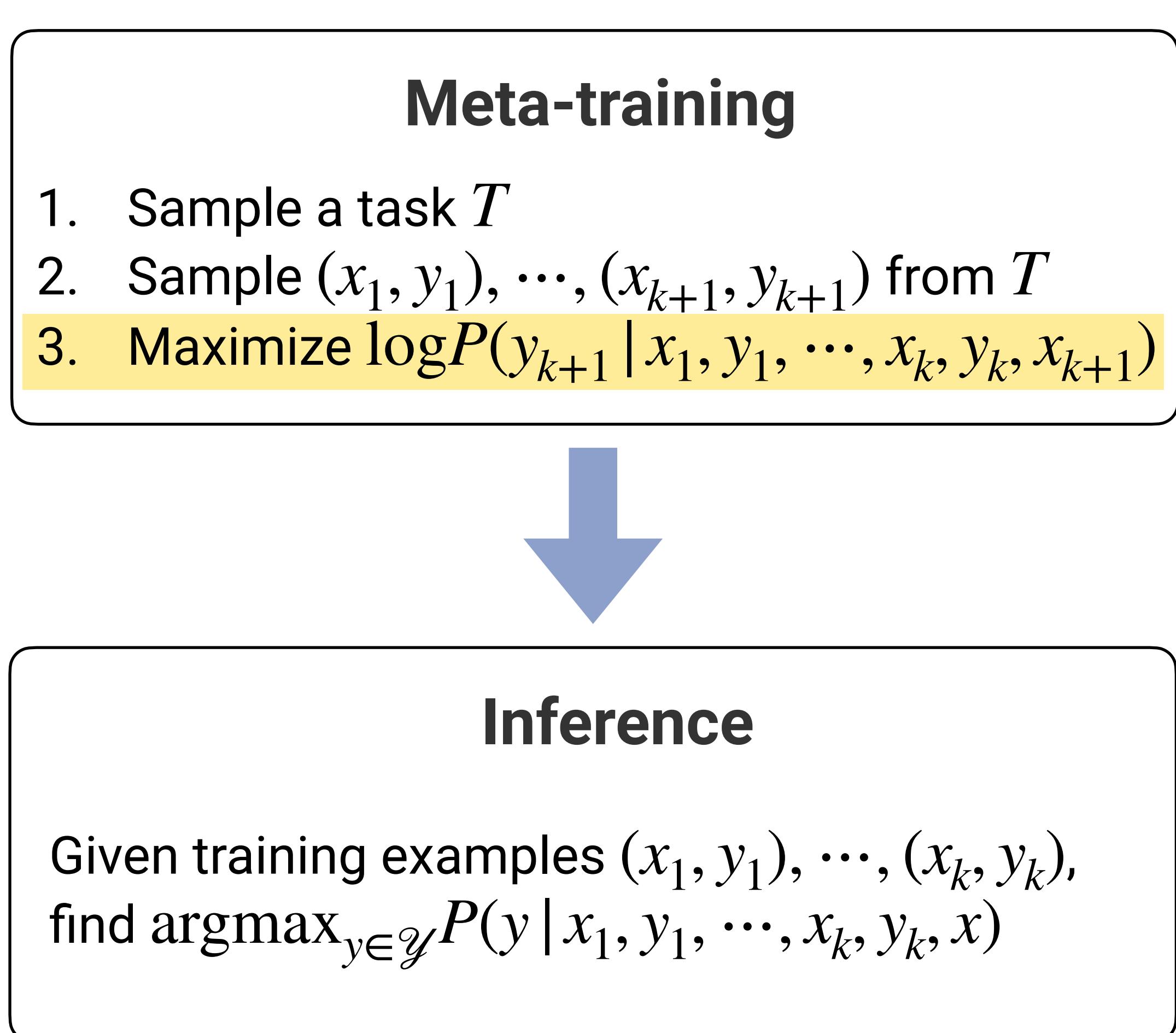
# MetalCL



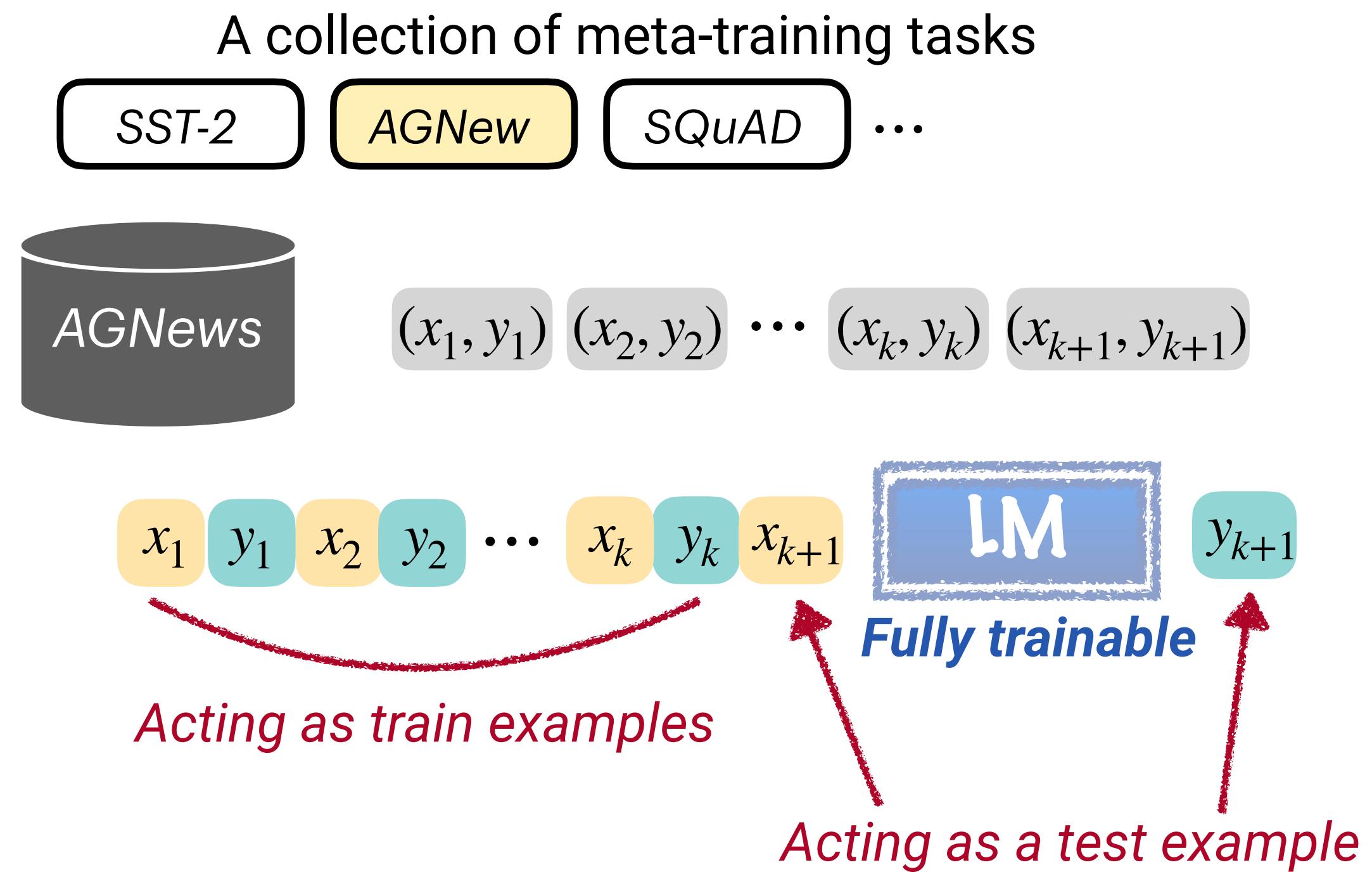
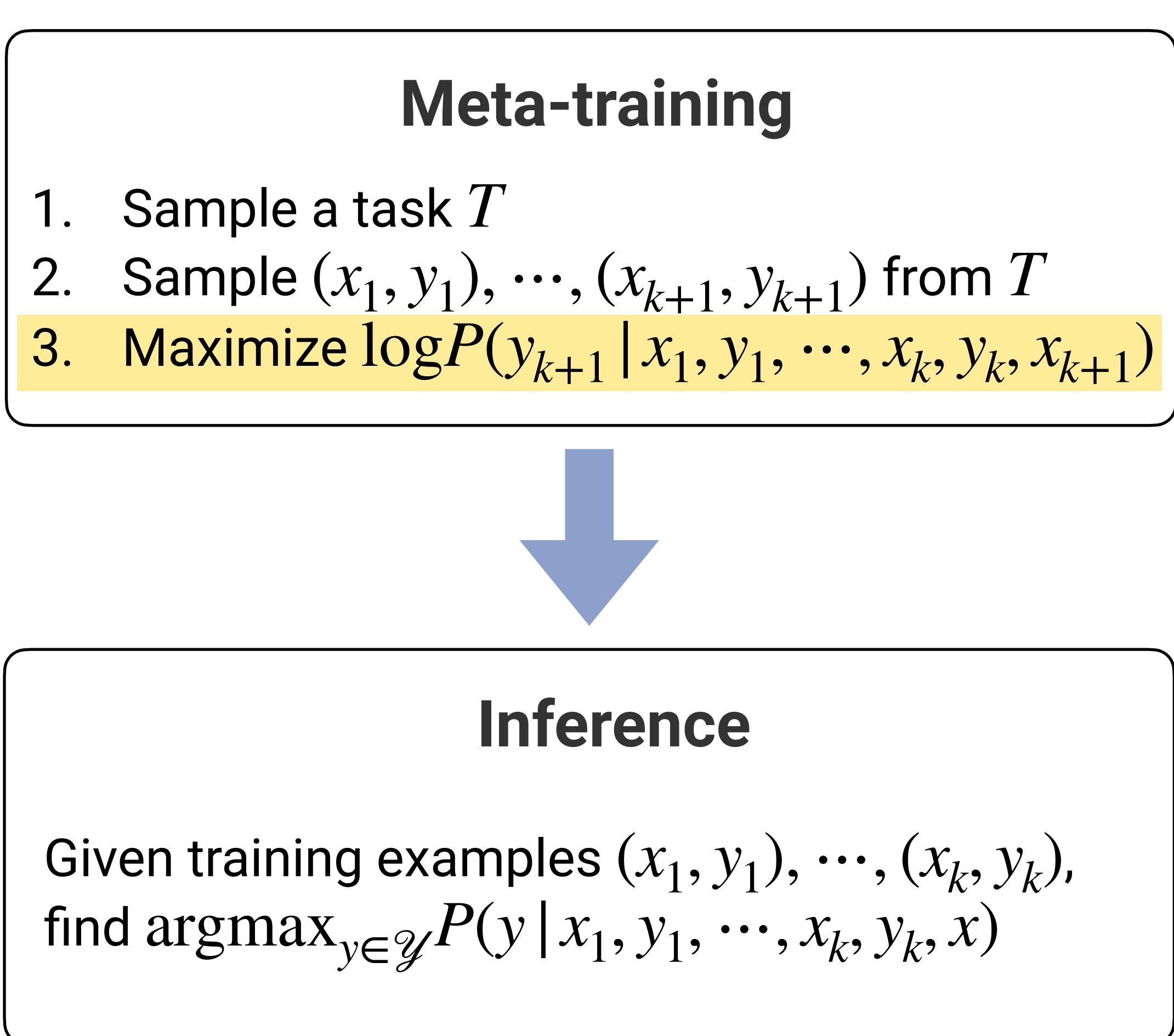
# MetalCL



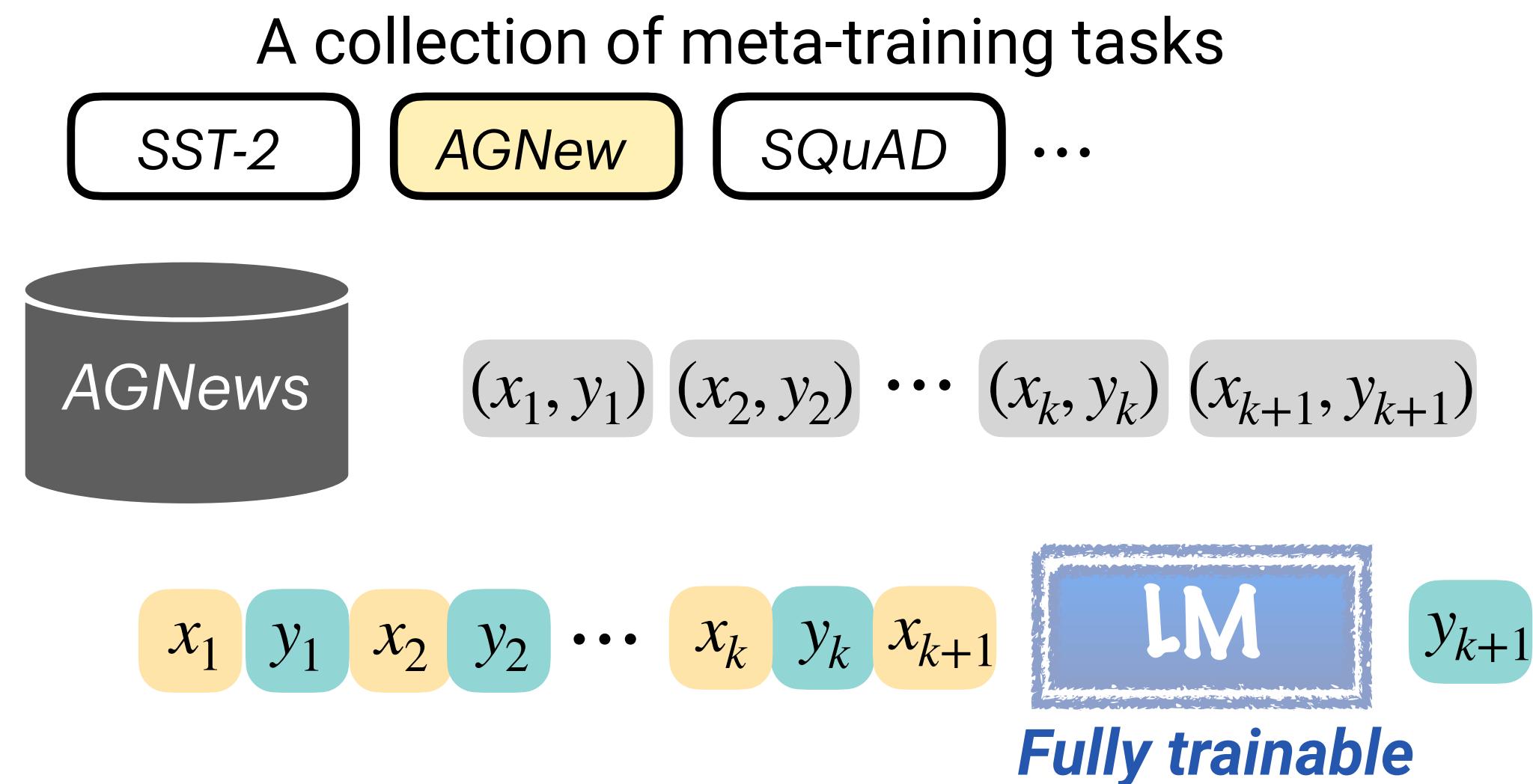
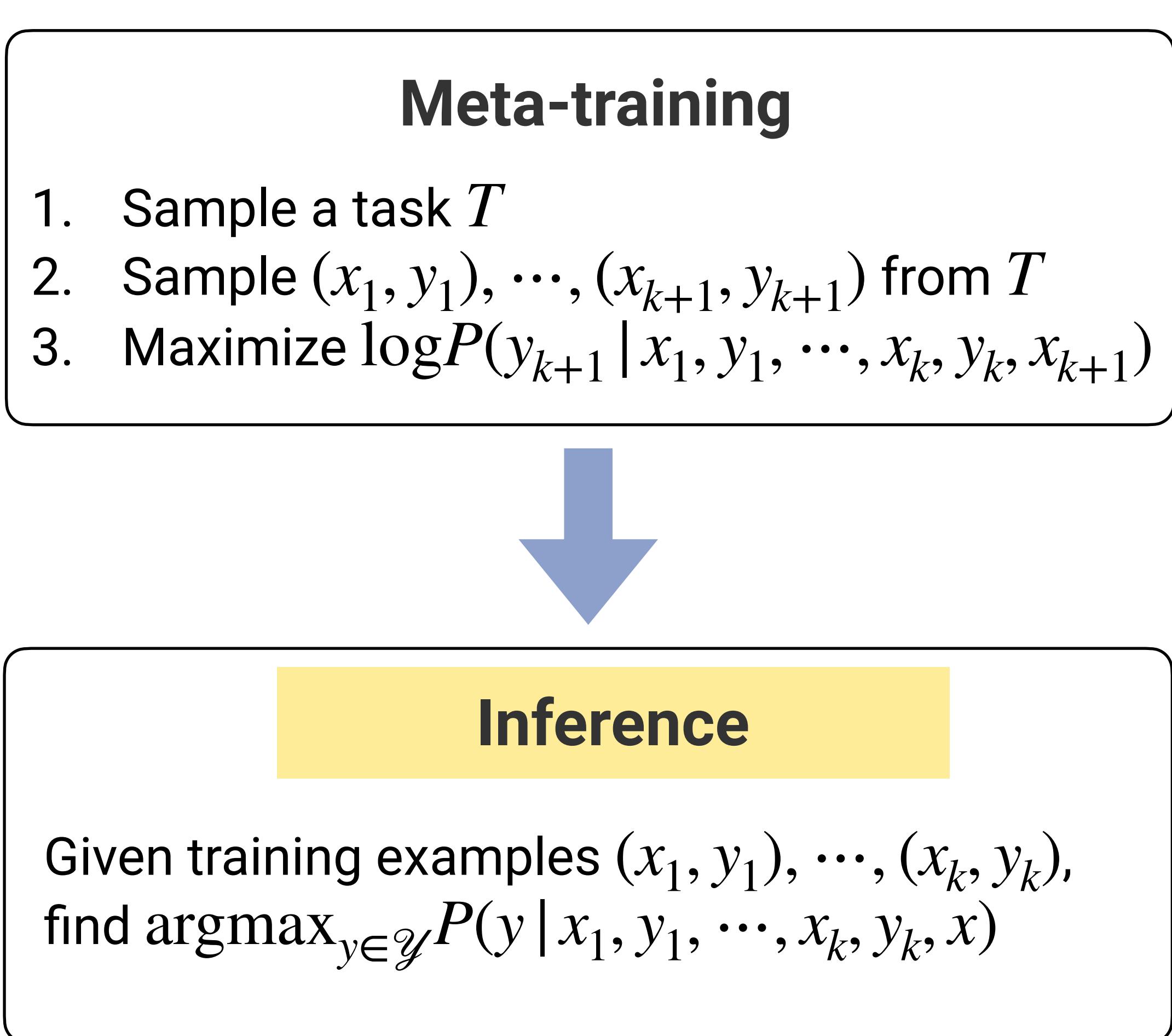
# MetalCL



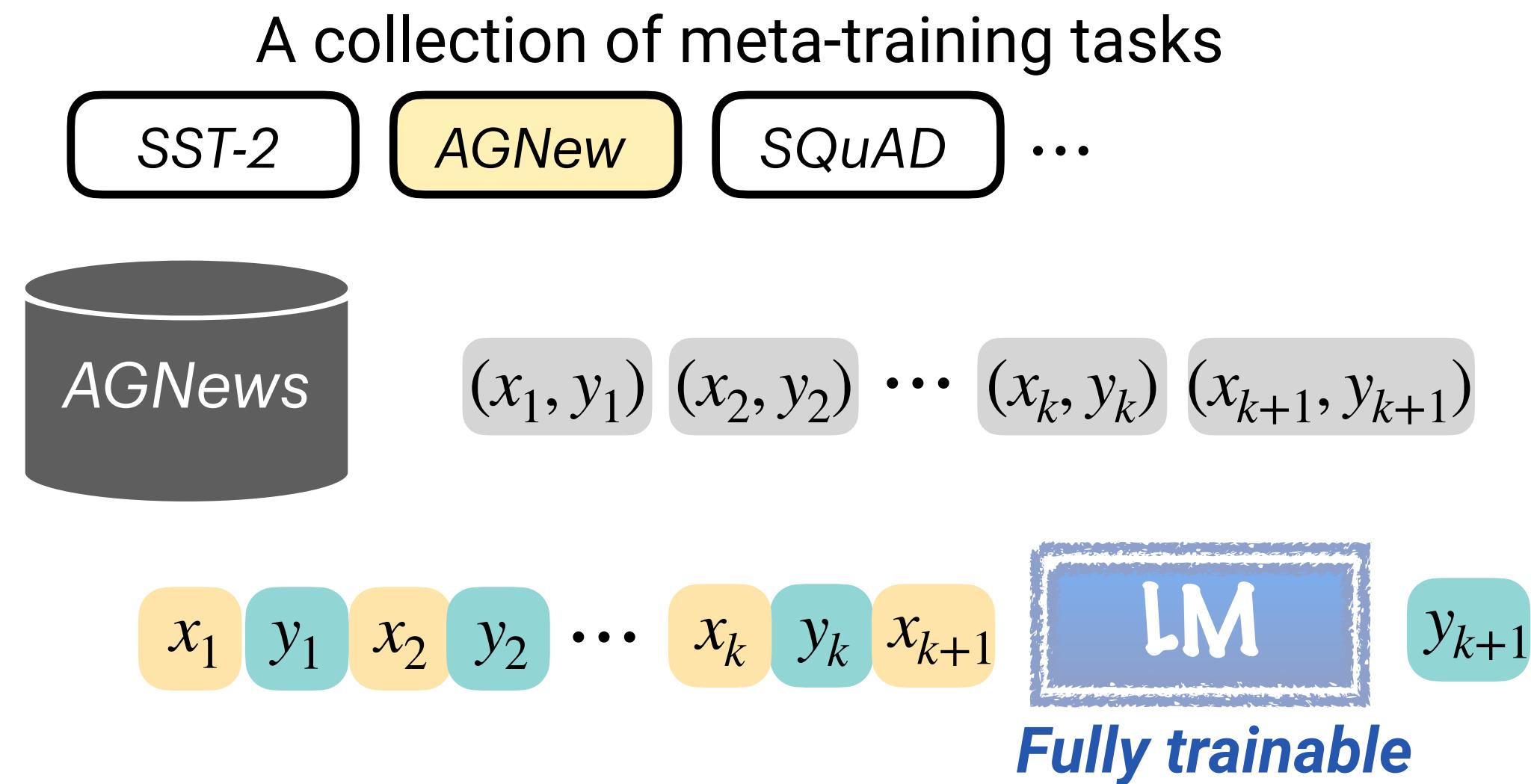
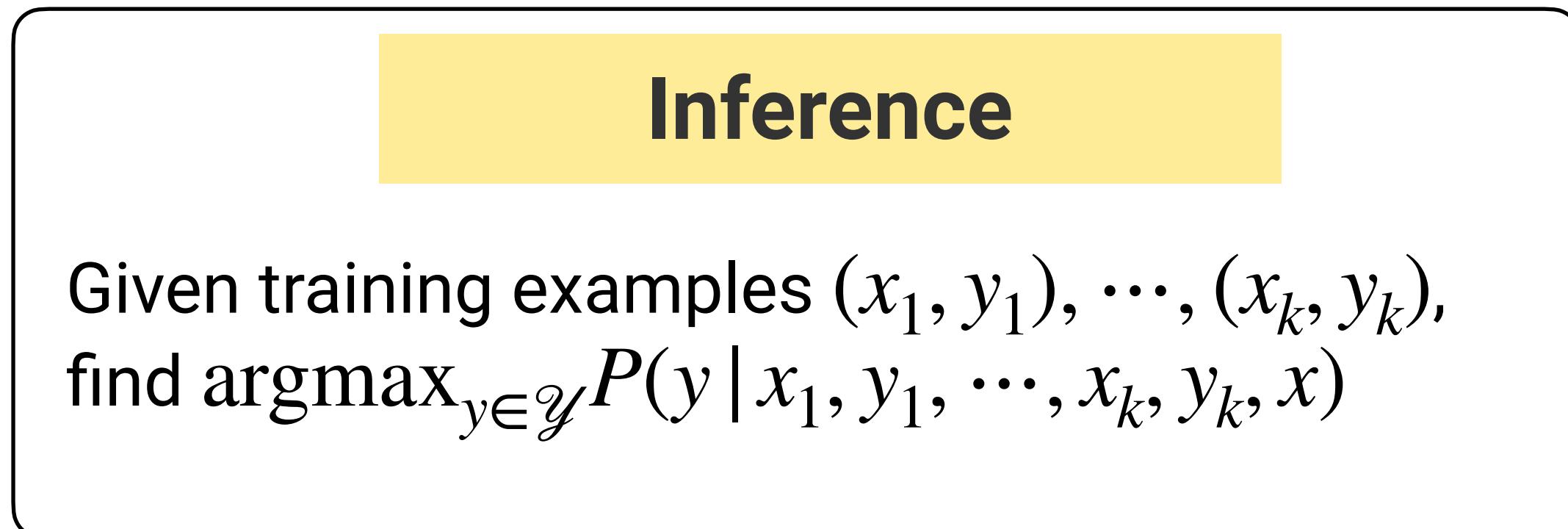
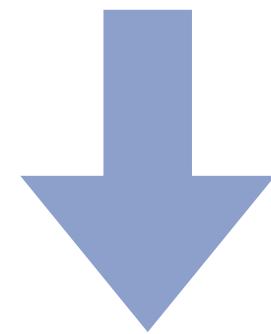
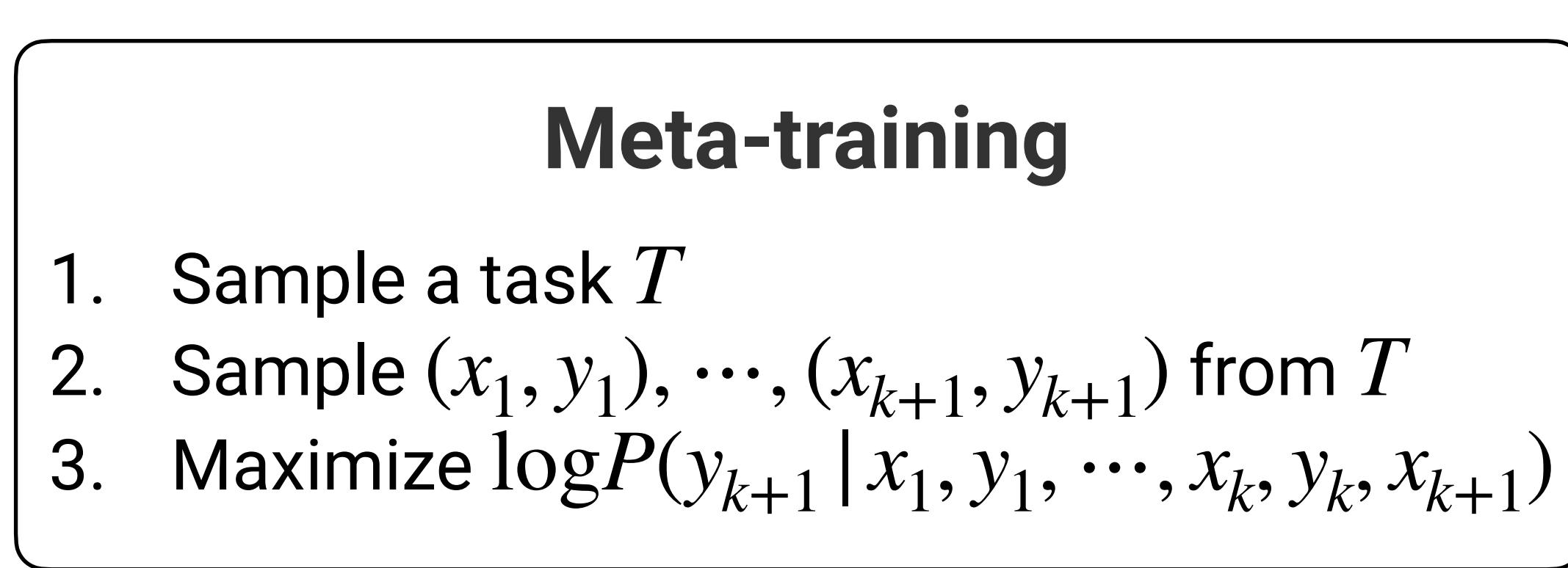
# MetalCL



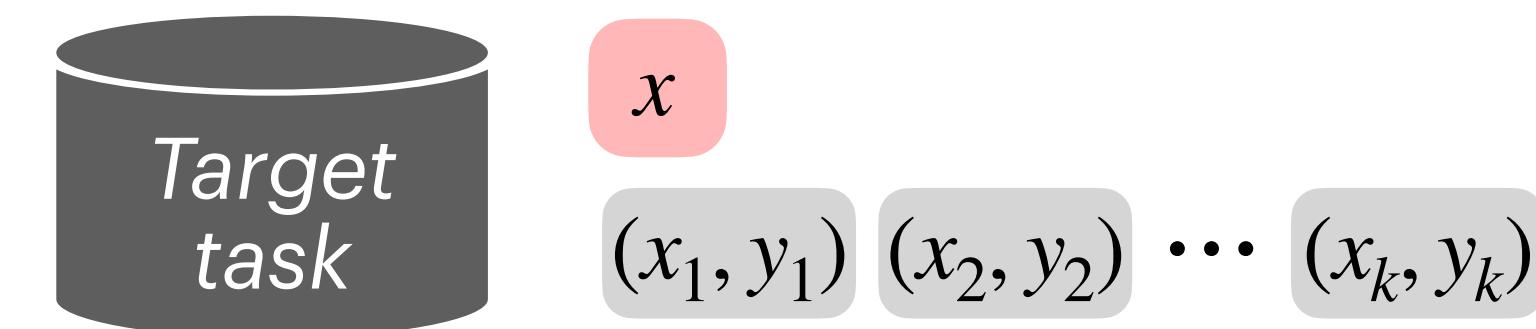
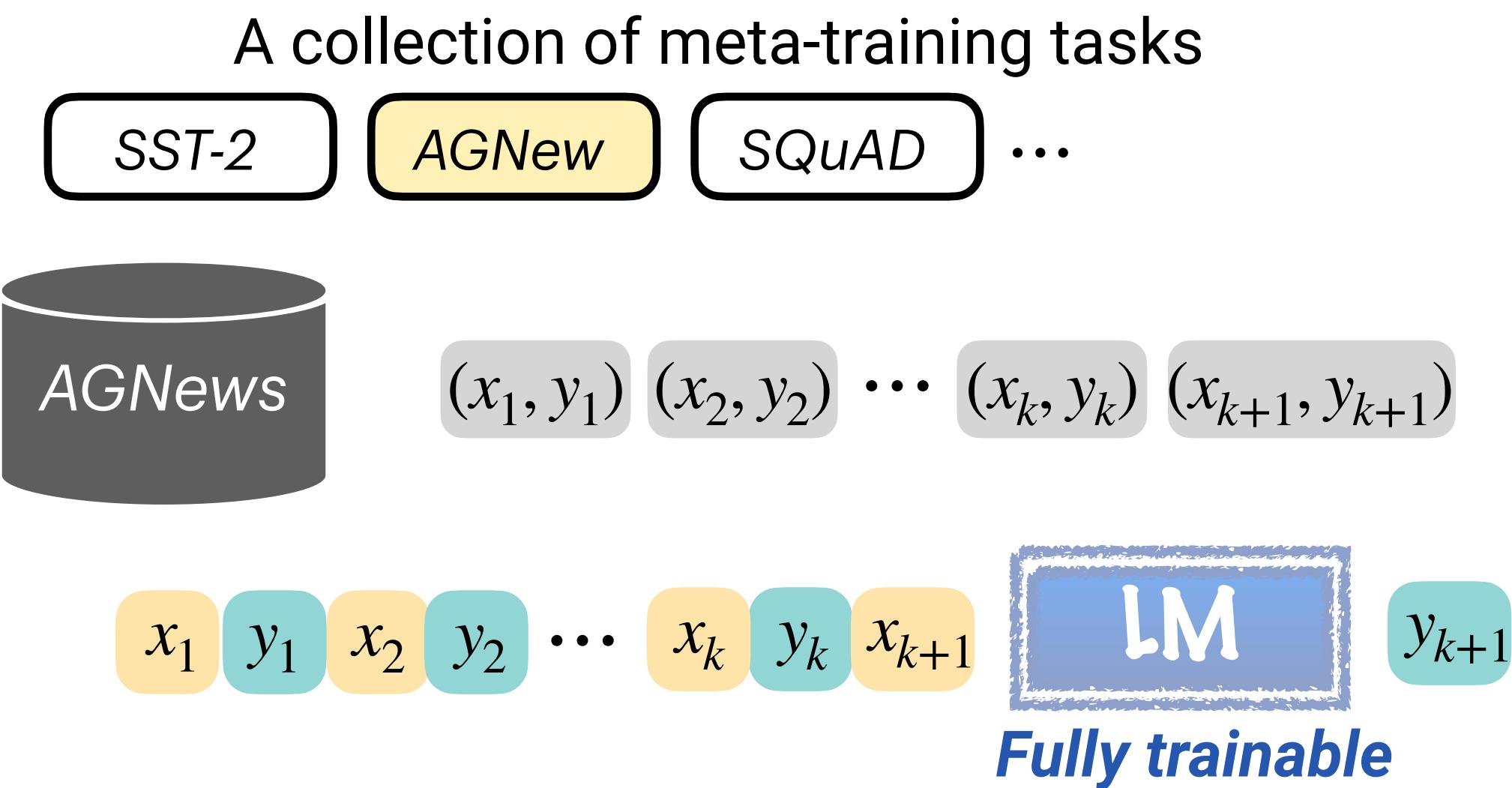
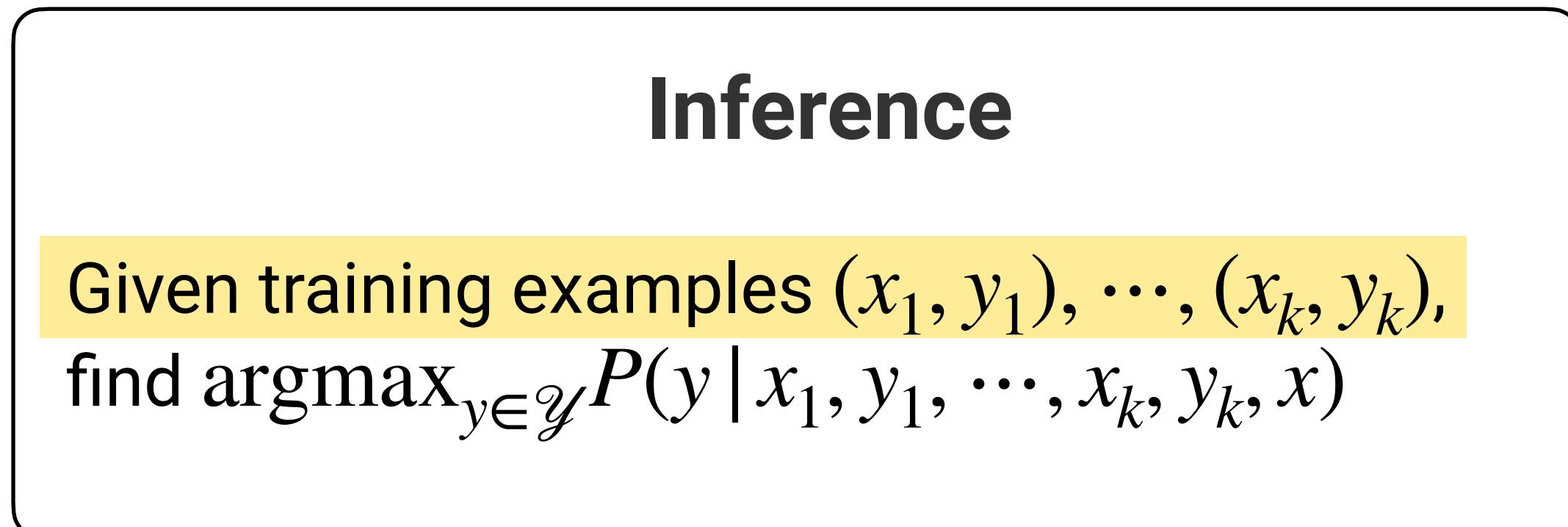
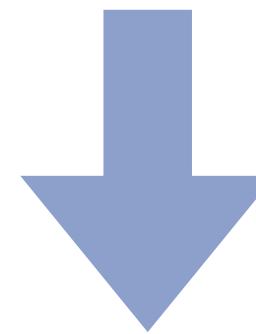
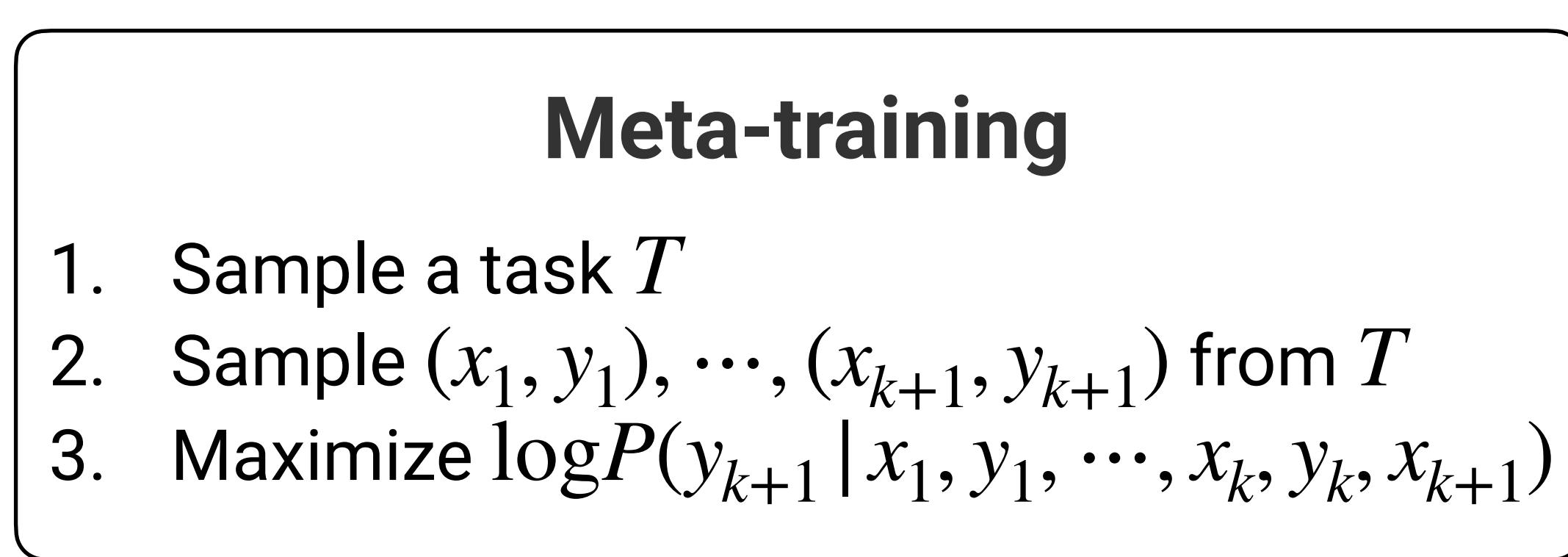
# MetalCL



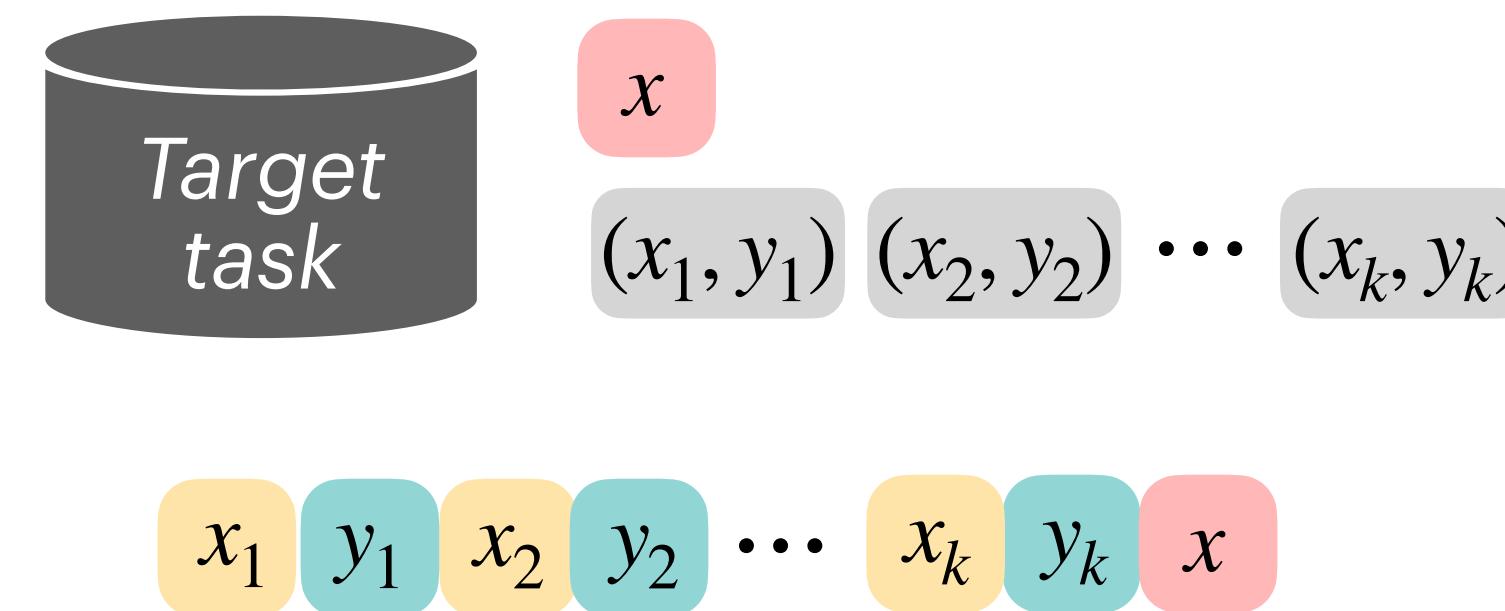
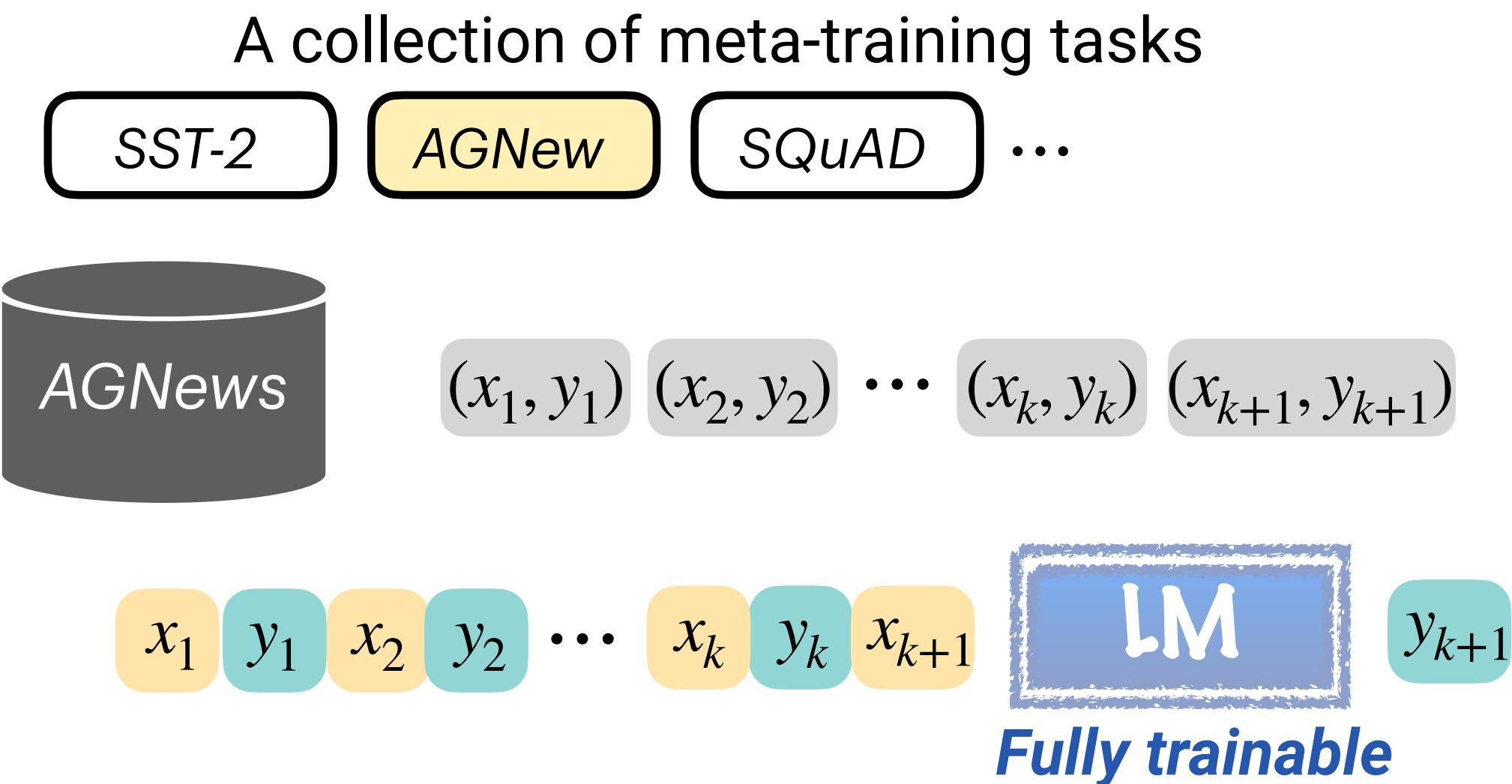
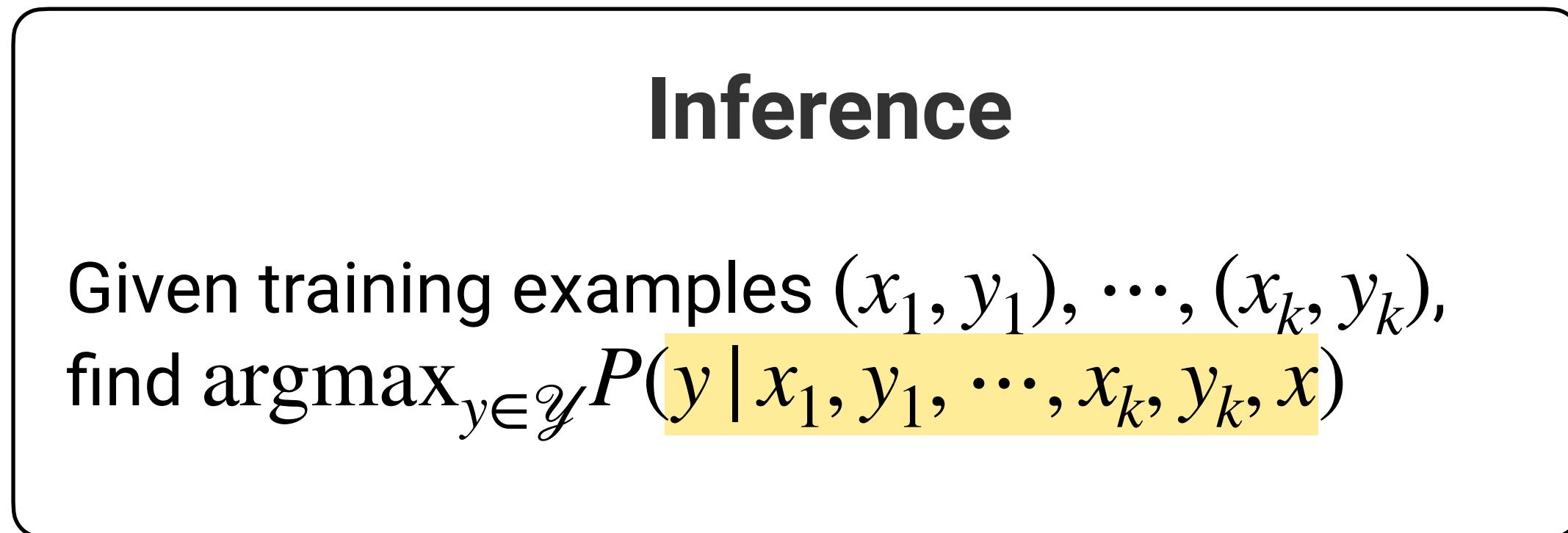
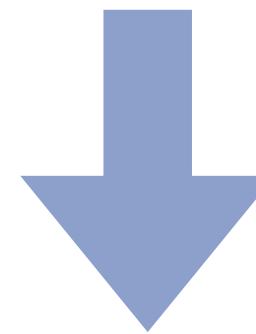
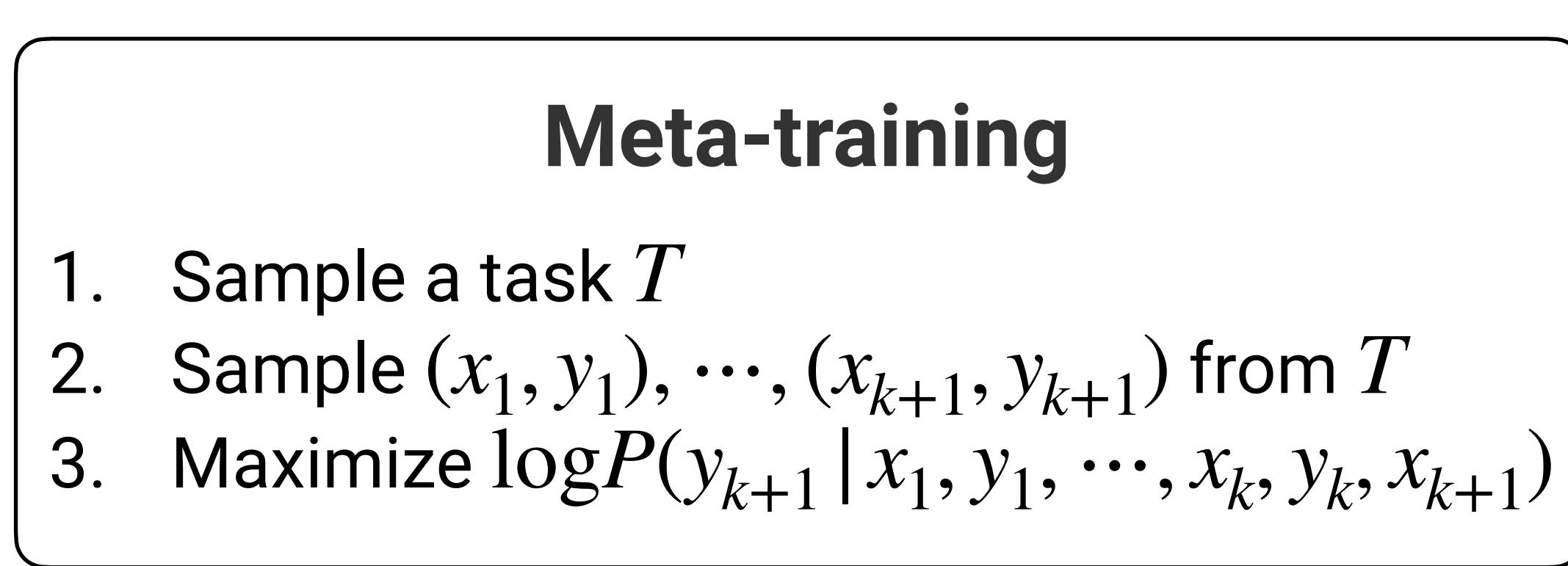
# MetalCL



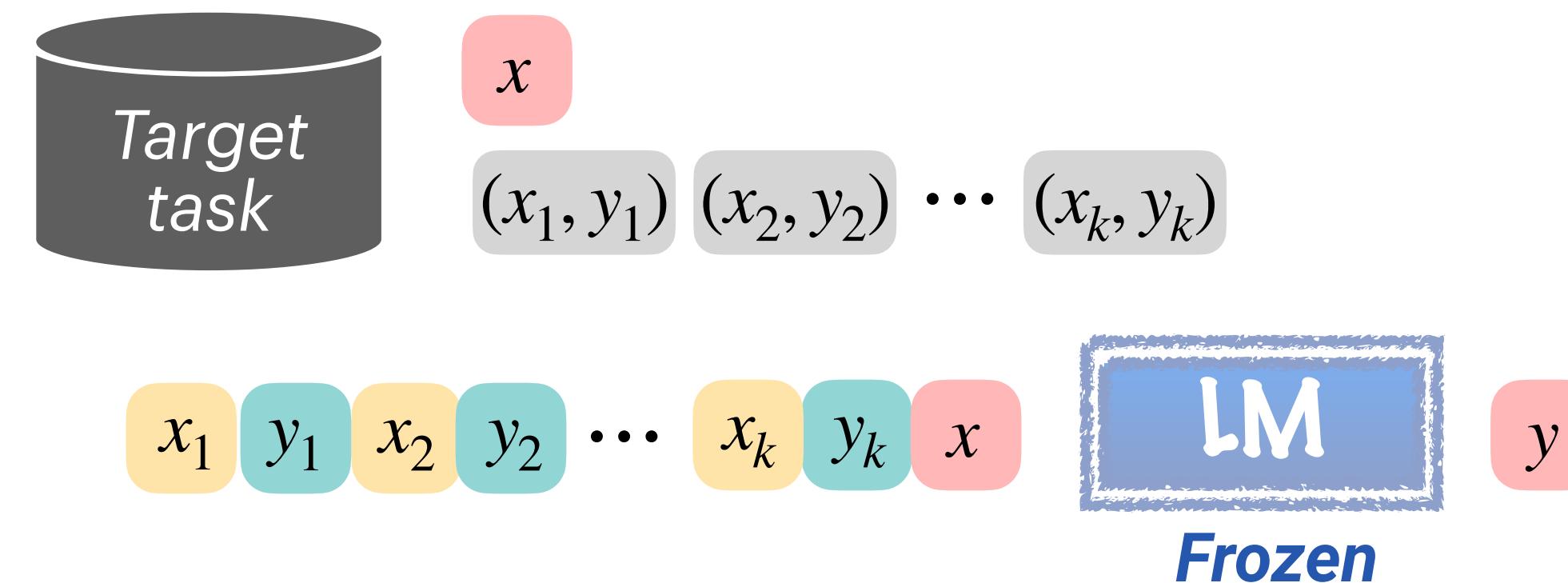
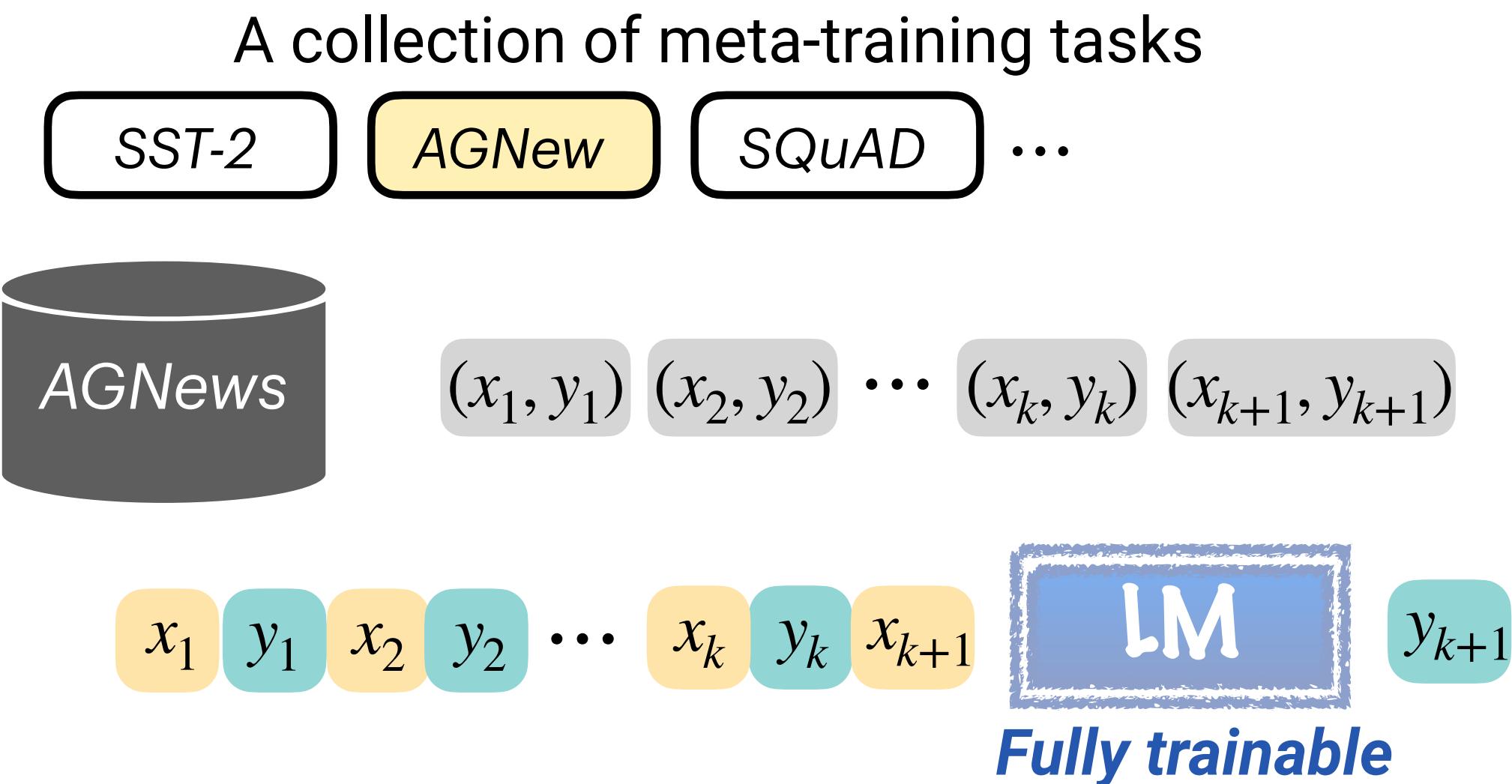
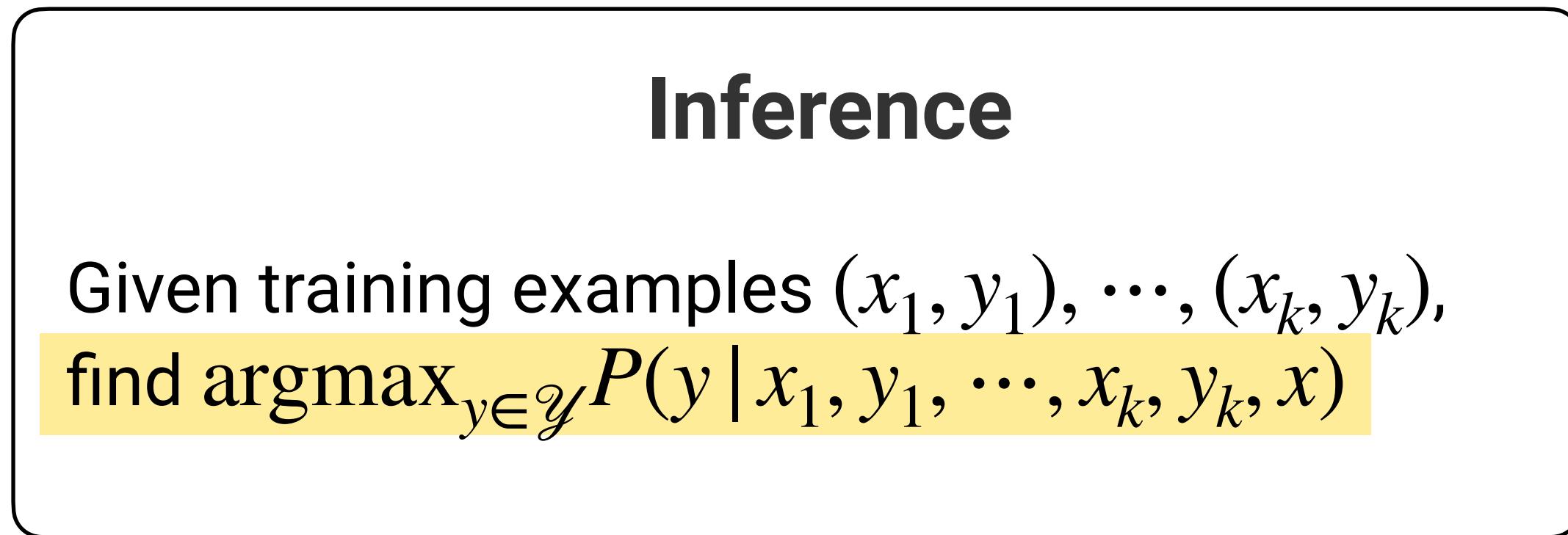
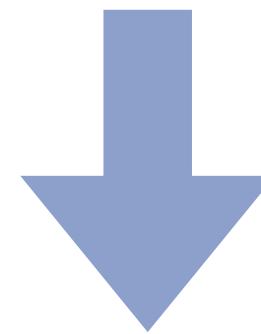
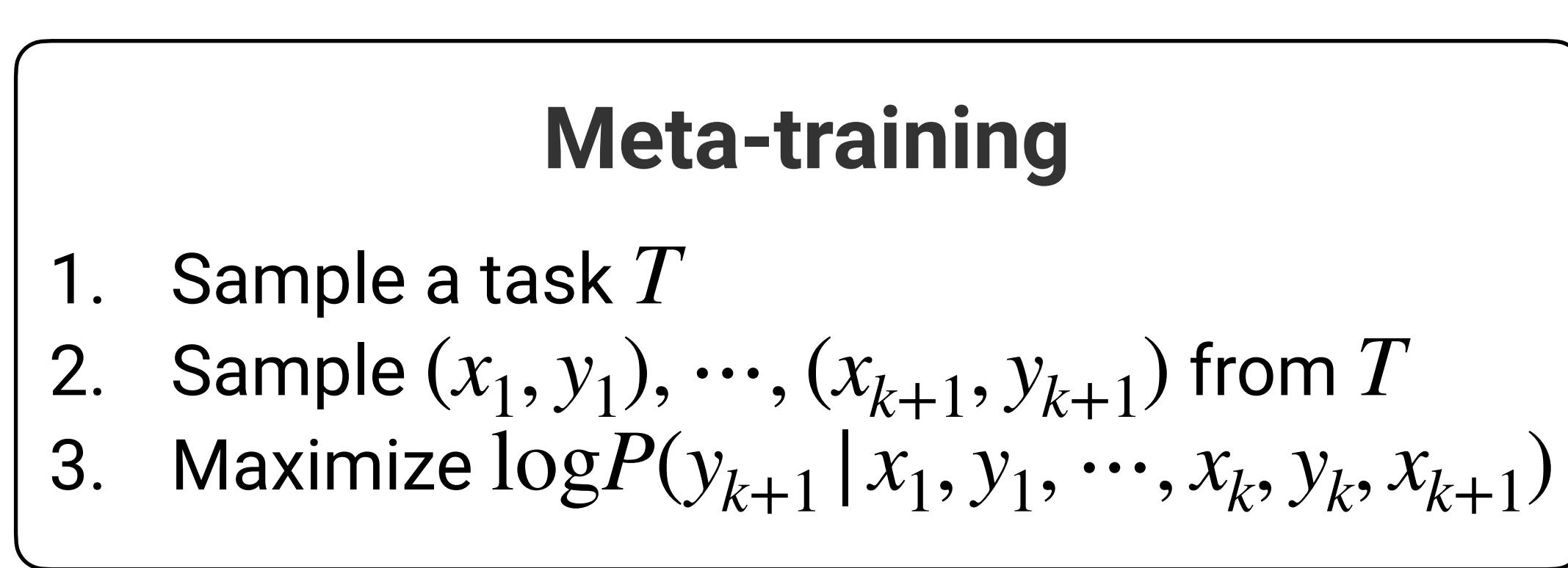
# MetalCL



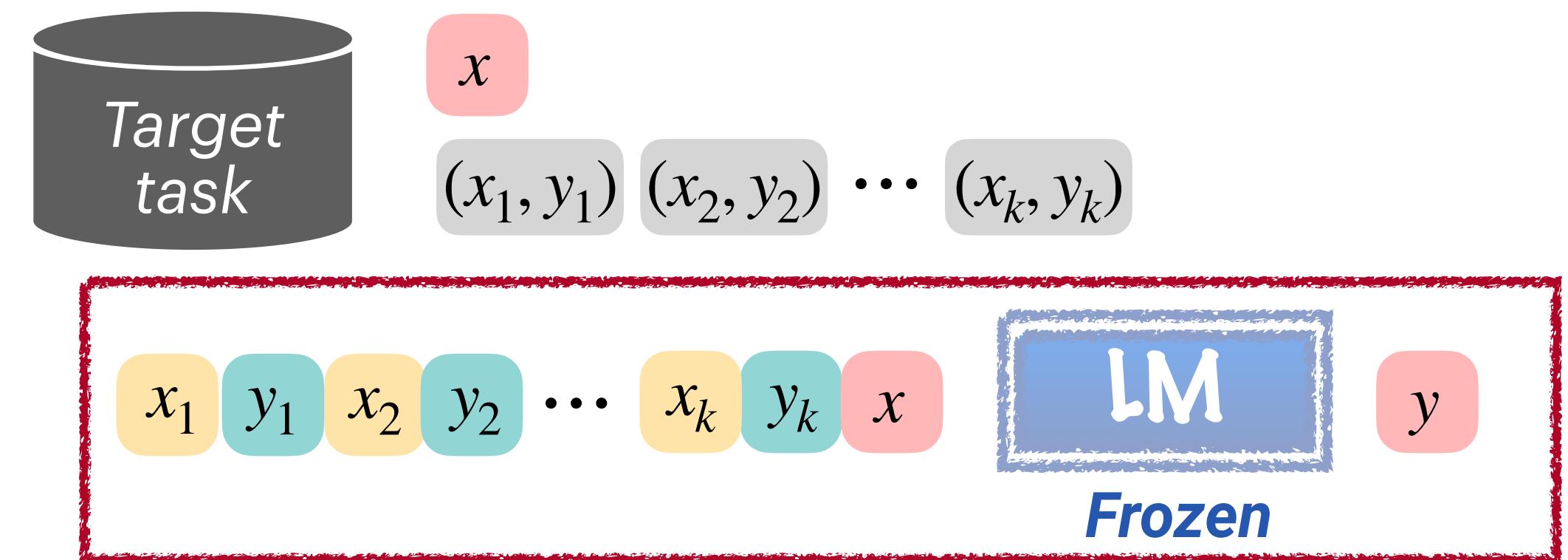
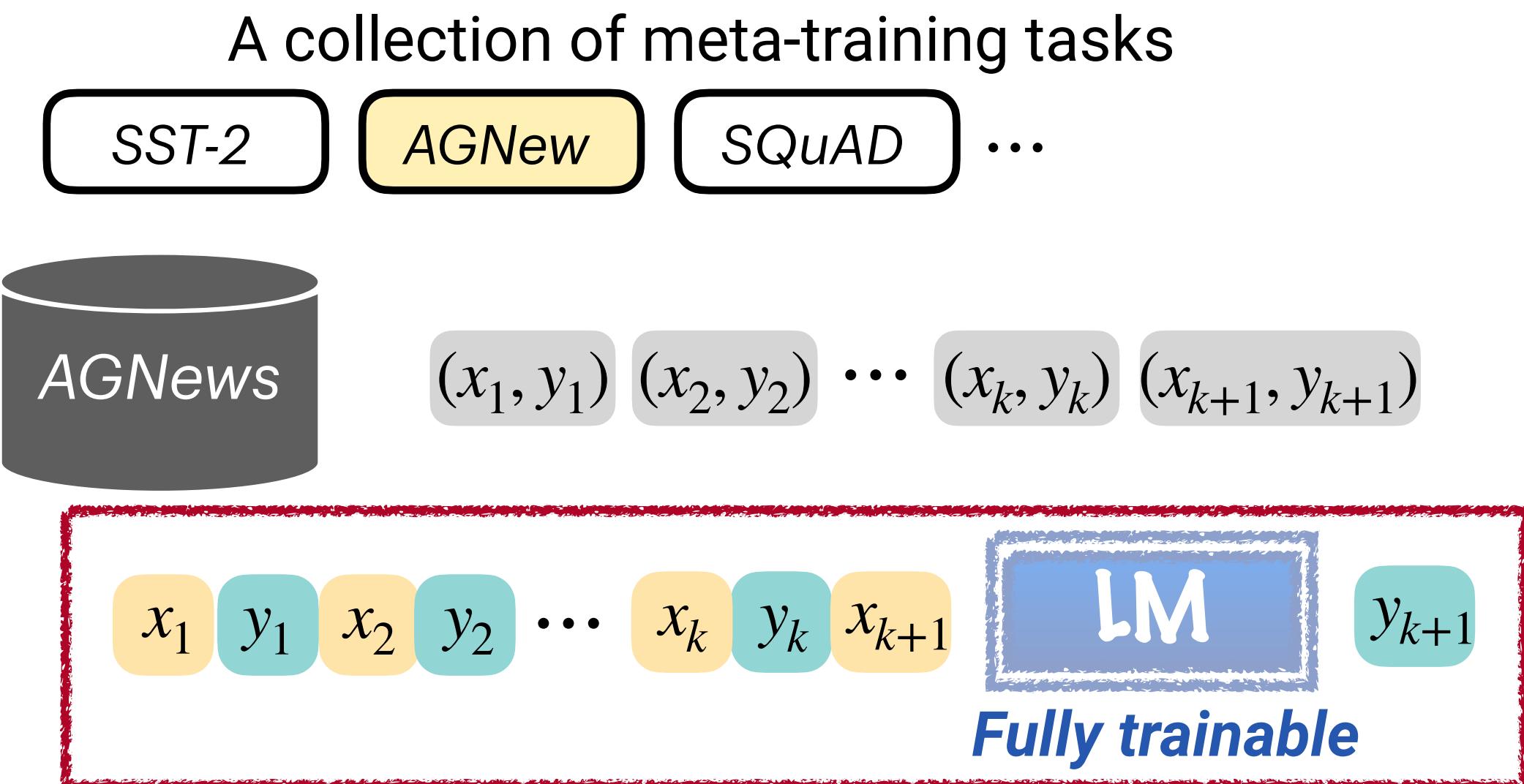
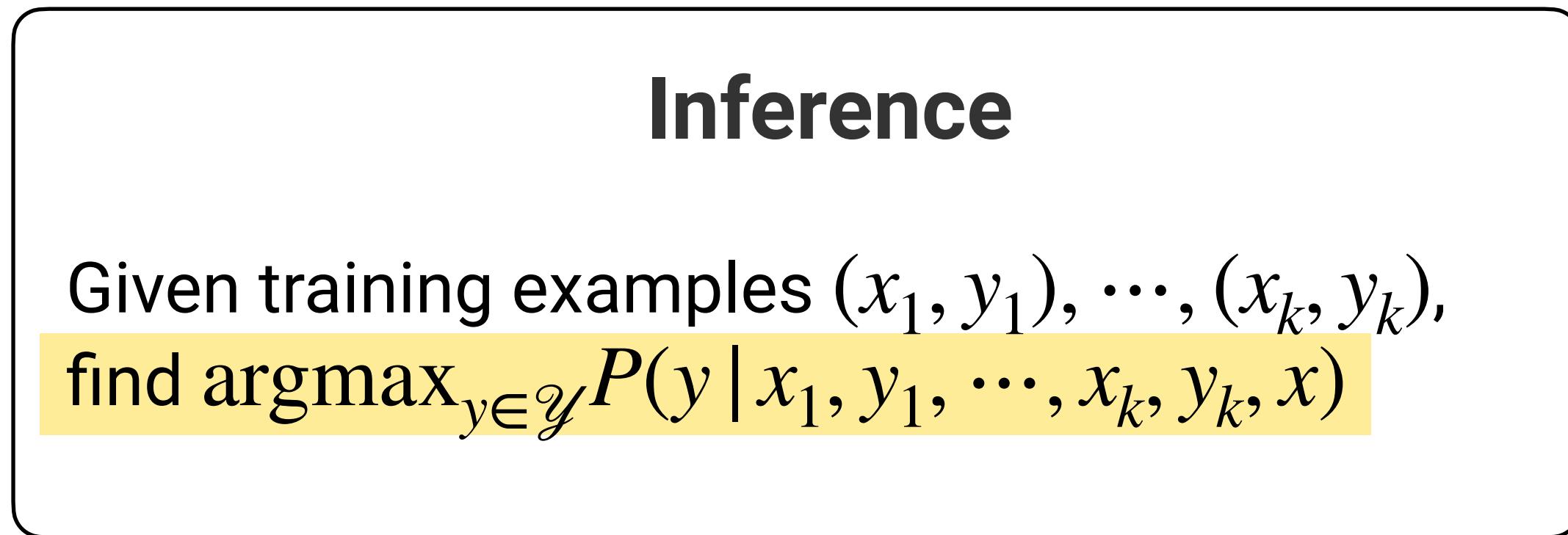
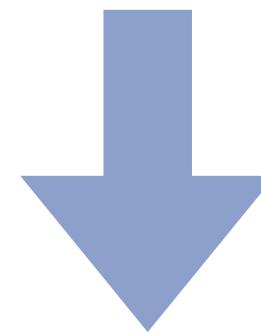
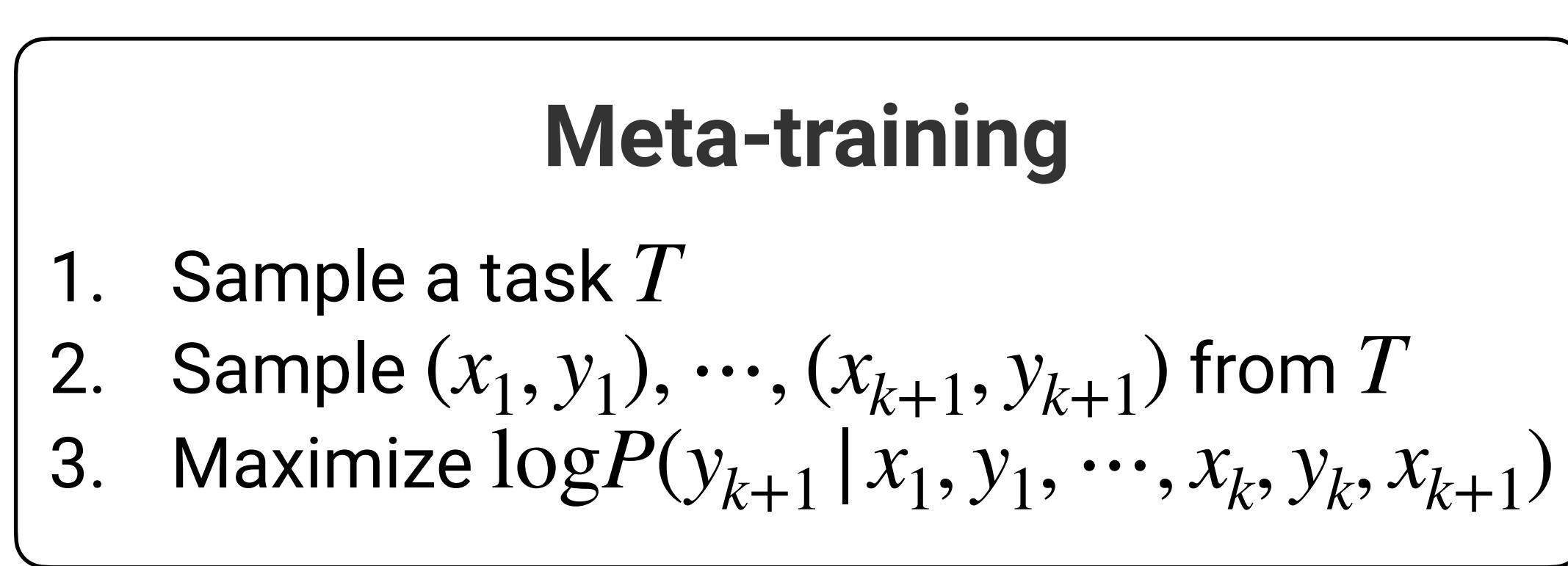
# MetalCL



# MetalCL



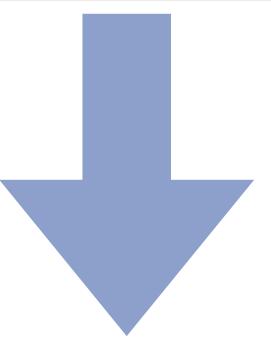
# MetalCL



# MetalCL

## Meta-training

1. Sample a task  $T$
2. Sample  $(x_1, y_1), \dots, (x_{k+1}, y_{k+1})$  from  $T$
3. Maximize  $\log P(y_{k+1} | x_1, y_1, \dots, x_k, y_k, x_{k+1})$



## Inference

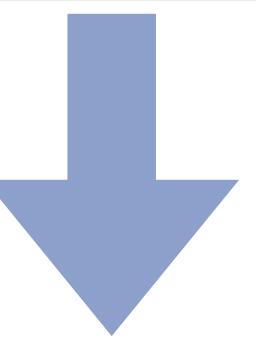
Given training examples  $(x_1, y_1), \dots, (x_k, y_k)$ ,  
find  $\operatorname{argmax}_{y \in \mathcal{Y}} P(y | x_1, y_1, \dots, x_k, y_k, x)$

***Direct MetalCL***

# MetalCL

## Meta-training

1. Sample a task  $T$
2. Sample  $(x_1, y_1), \dots, (x_{k+1}, y_{k+1})$  from  $T$
3. Maximize  $\log P(y_{k+1} | x_1, y_1, \dots, x_k, y_k, x_{k+1})$



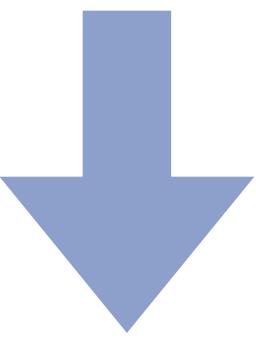
## Inference

Given training examples  $(x_1, y_1), \dots, (x_k, y_k)$ ,  
find  $\text{argmax}_{y \in \mathcal{Y}} P(y | x_1, y_1, \dots, x_k, y_k, x)$

**Direct MetalCL**

## Meta-training

1. Sample a task  $T$
2. Sample  $(x_1, y_1), \dots, (x_{k+1}, y_{k+1})$  from  $T$
3. Maximize  $\log P(x_{k+1} | y_1, x_1, \dots, y_k, x_k, y_{k+1})$



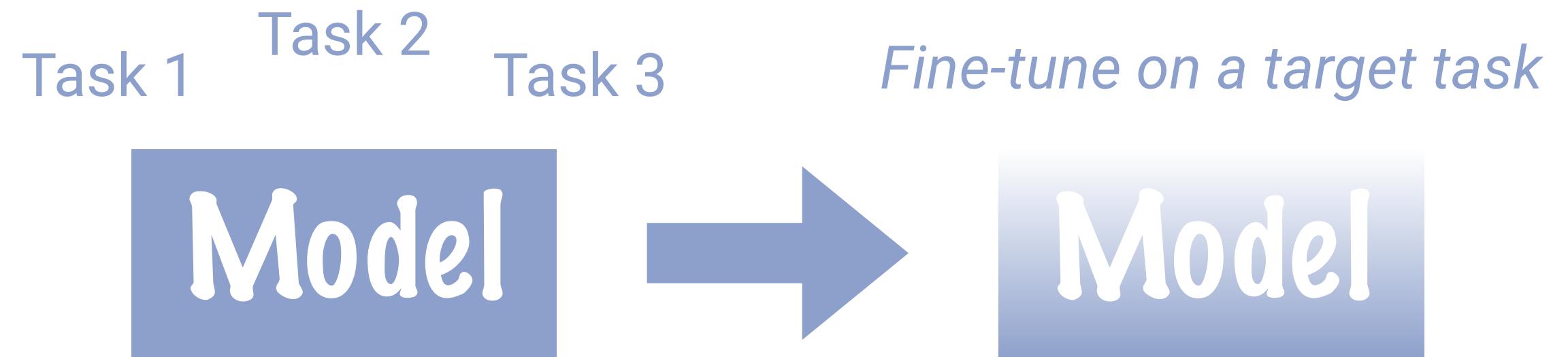
## Inference

Given training examples  $(x_1, y_1), \dots, (x_k, y_k)$ ,  
find  $\text{argmax}_{y \in \mathcal{Y}} P(x | y_1, x_1, \dots, y_k, x_k, y)$

**Channel MetalCL**

# Related Work

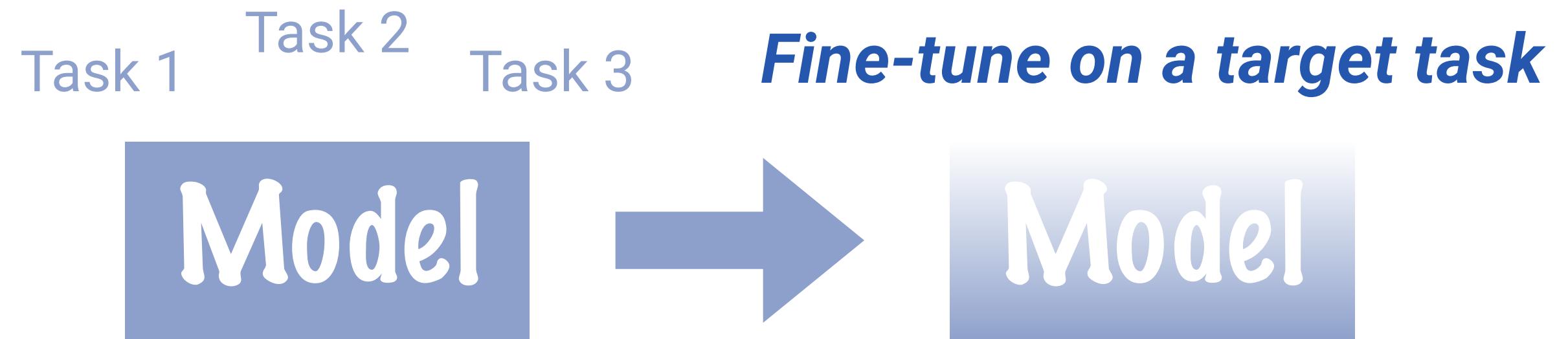
# Related Work



## Multi-task learning/Meta-learning

(Villa and Drissi 2002, Evgeniou and Pontil 2004, Ruder 2017,  
Wang et al. 2020, Aghajanyan et al. 2021)

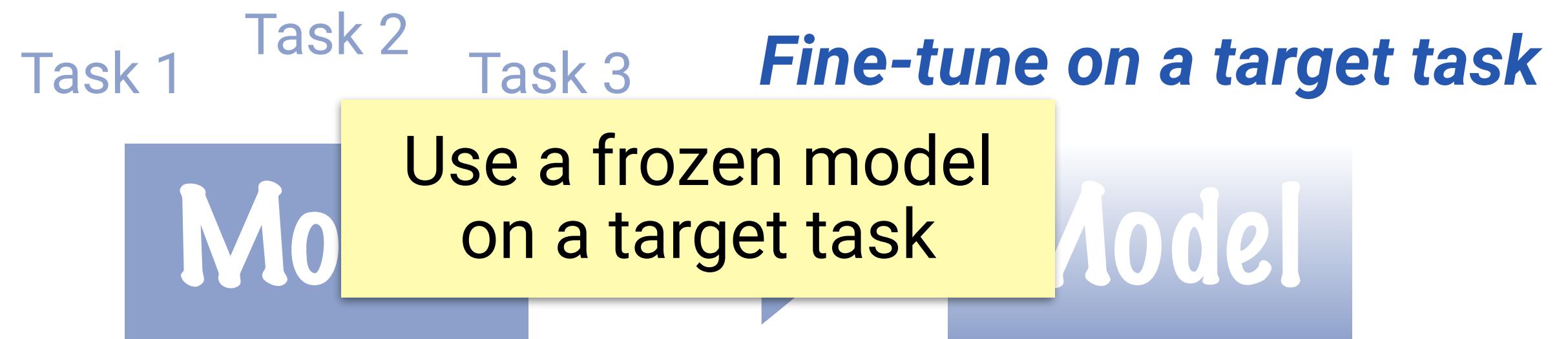
# Related Work



## Multi-task learning/Meta-learning

(Villa and Drissi 2002, Evgeniou and Pontil 2004, Ruder 2017,  
Wang et al. 2020, Aghajanyan et al. 2021)

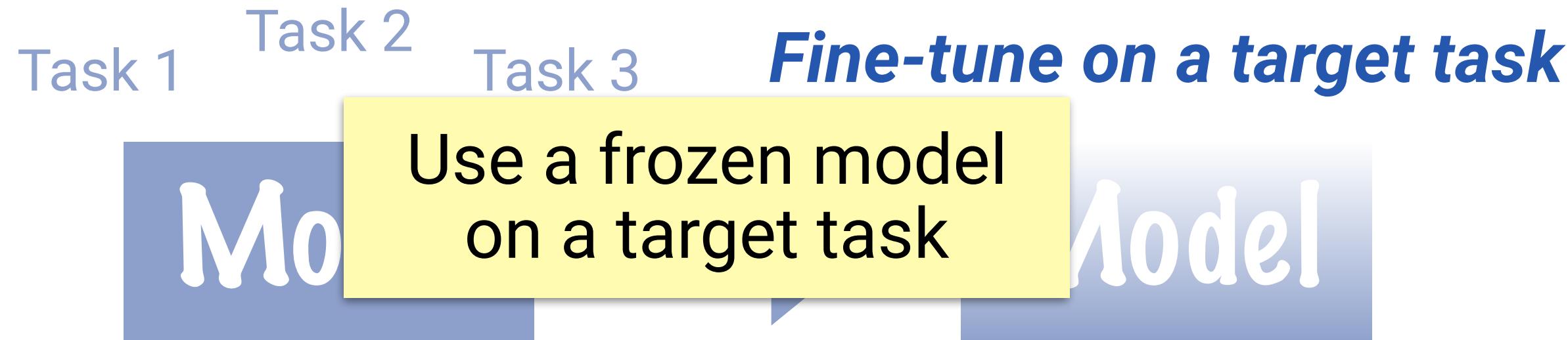
# Related Work



## Multi-task learning/Meta-learning

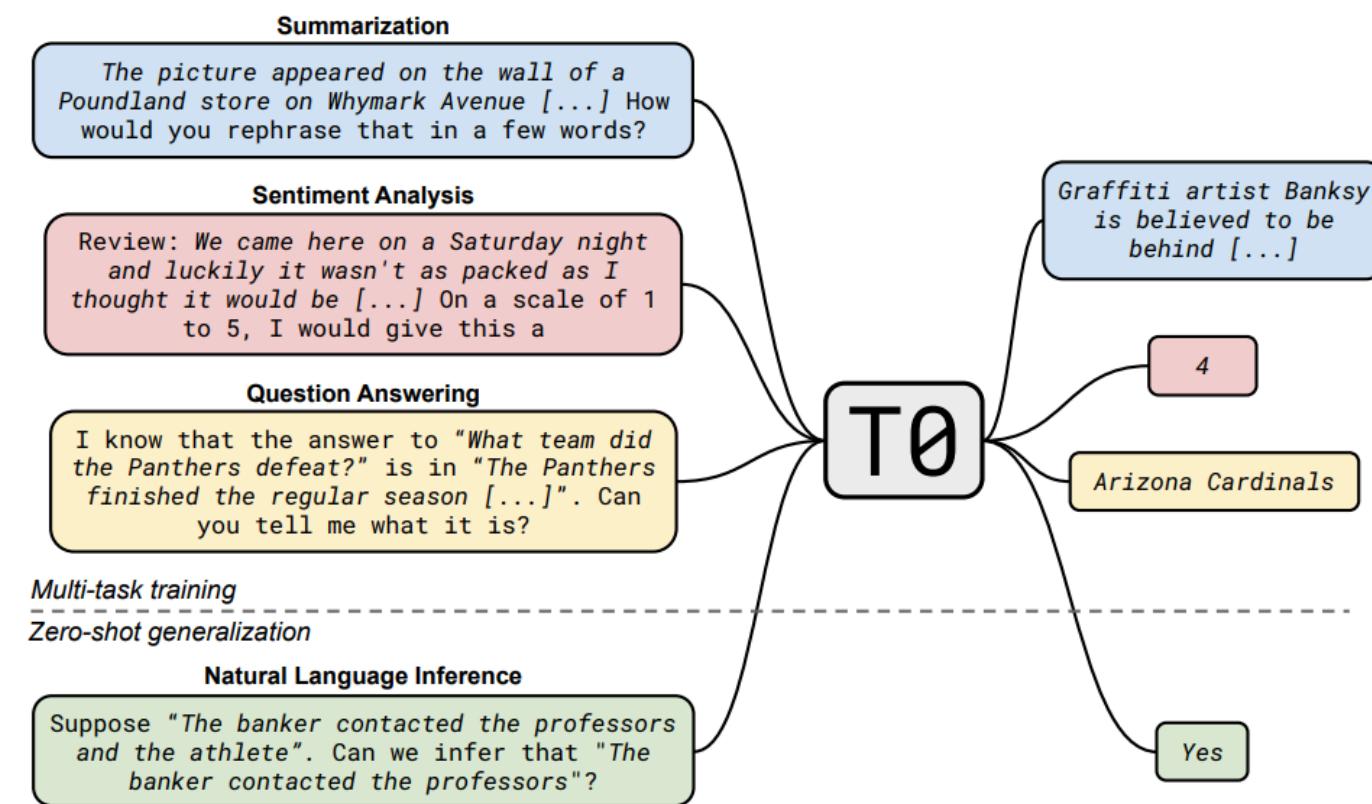
(Villa and Drissi 2002, Evgeniou and Pontil 2004, Ruder 2017,  
Wang et al. 2020, Aghajanyan et al. 2021)

# Related Work



## Multi-task learning/Meta-learning

(Villa and Drissi 2002, Evgeniou and Pontil 2004, Ruder 2017, Wang et al. 2020, Aghajanyan et al. 2021)

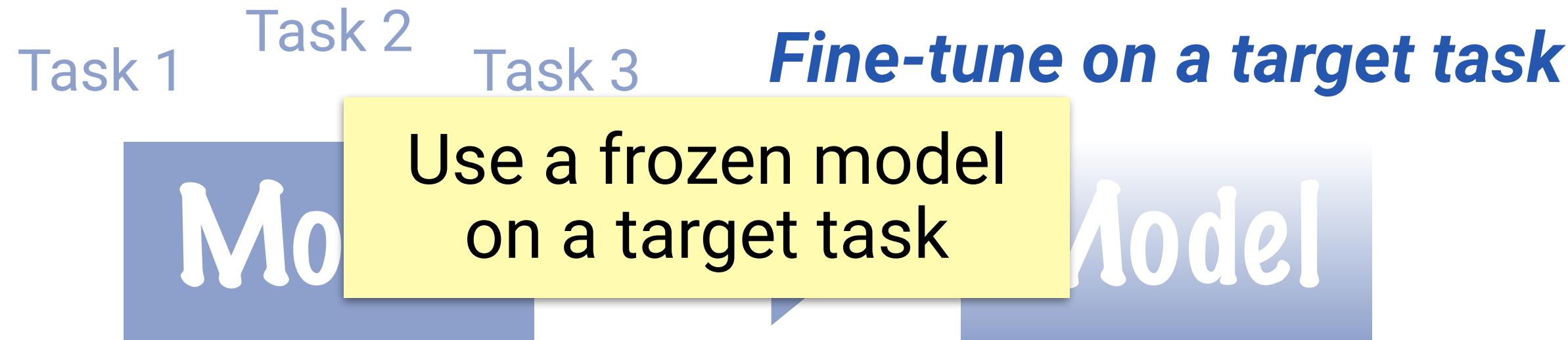


## Instruction-reading models

(Mishra et al. 2021, Wei et al. 2021, Sanh et al. 2021, Zhong et al. 2021)

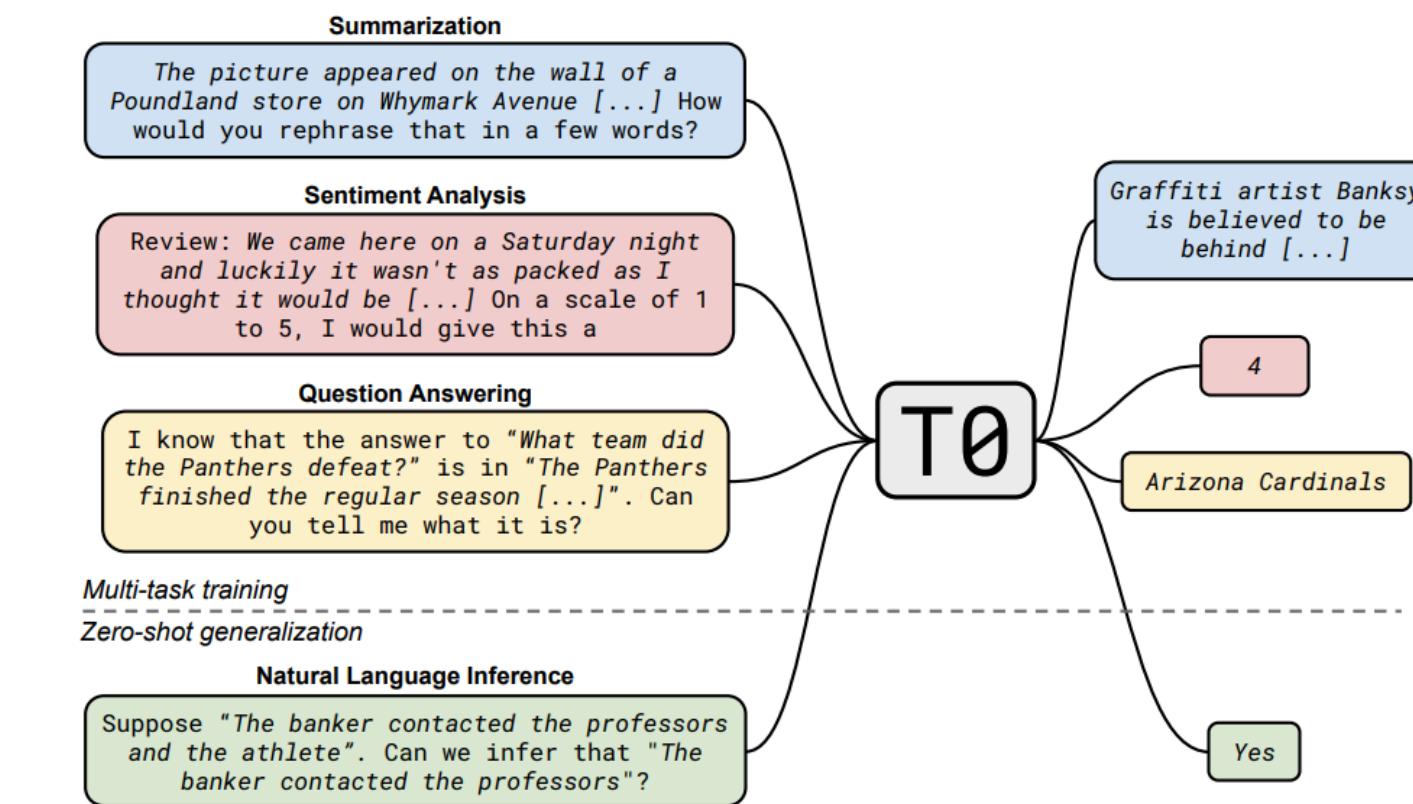
Figure from Sanh et al. 2021

# Related Work



## Multi-task learning/Meta-learning

(Villa and Drissi 2002, Evgeniou and Pontil 2004, Ruder 2017, Wang et al. 2020, Aghajanyan et al. 2021)



## Instruction-reading models

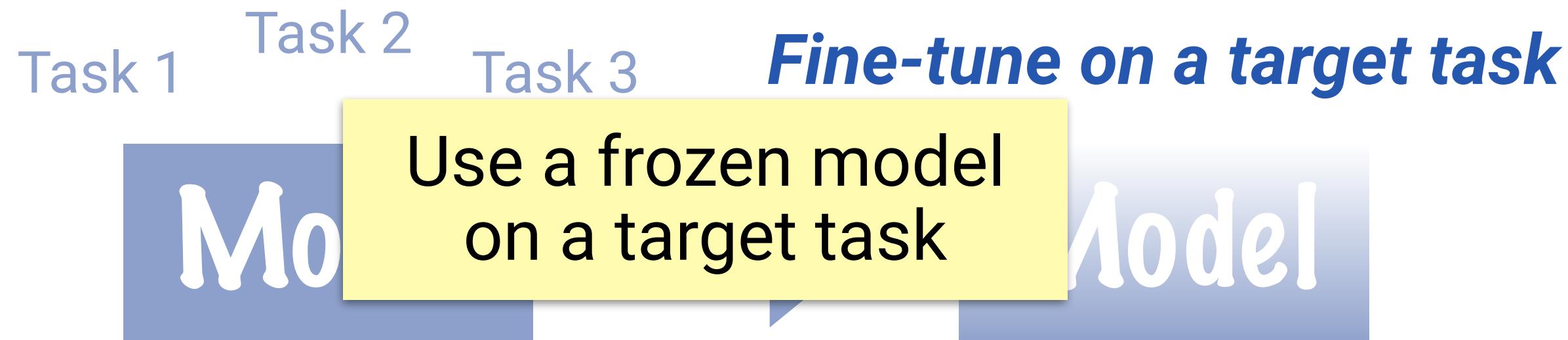
(Mishra et al. 2021, Wei et al. 2021, Sanh et al. 2021, Zhong et al. 2021)

Figure from Sanh et al. 2021



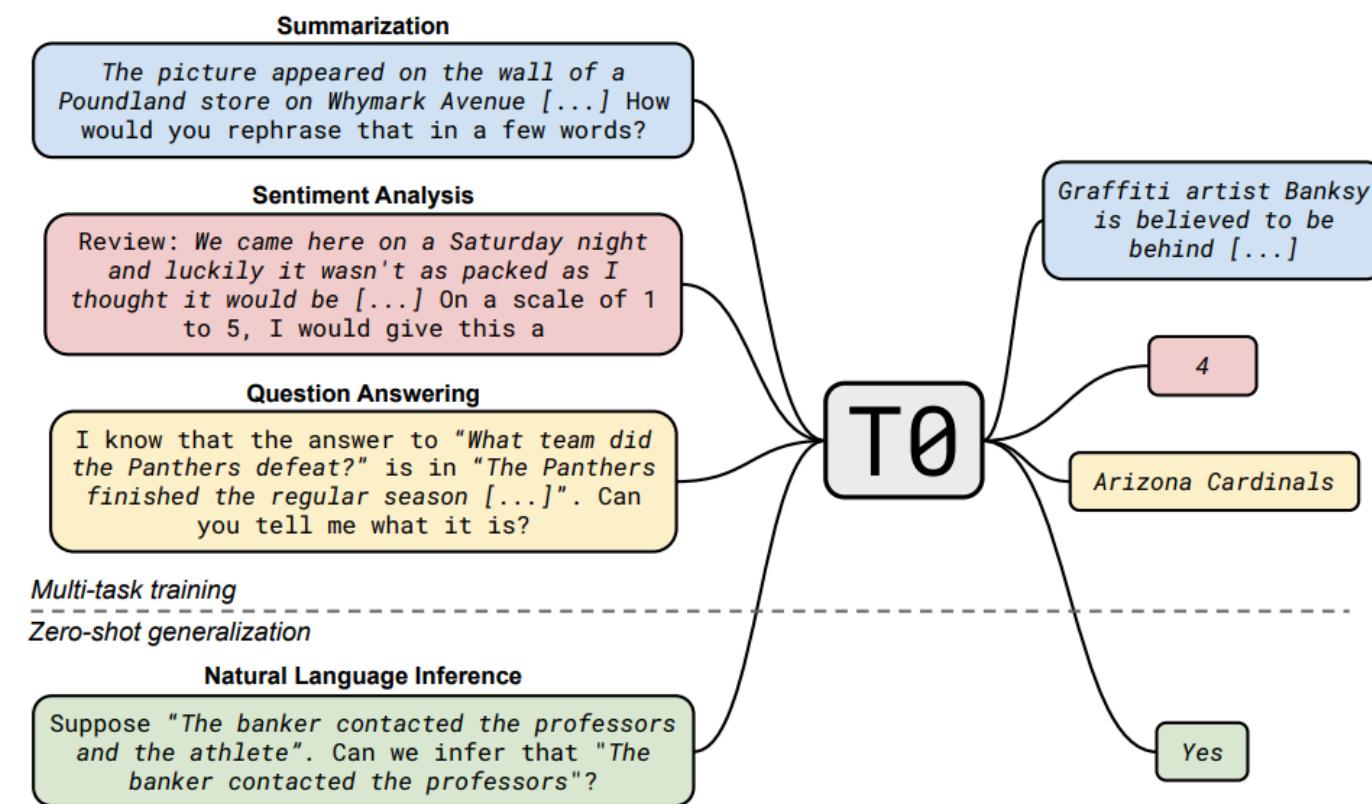
Define the task as the input

# Related Work



## Multi-task learning/Meta-learning

(Villa and Drissi 2002, Evgeniou and Pontil 2004, Ruder 2017, Wang et al. 2020, Aghajanyan et al. 2021)



## Instruction-reading models

(Mishra et al. 2021, Wei et al. 2021, Sanh et al. 2021, Zhong et al. 2021)  
Figure from Sanh et al. 2021

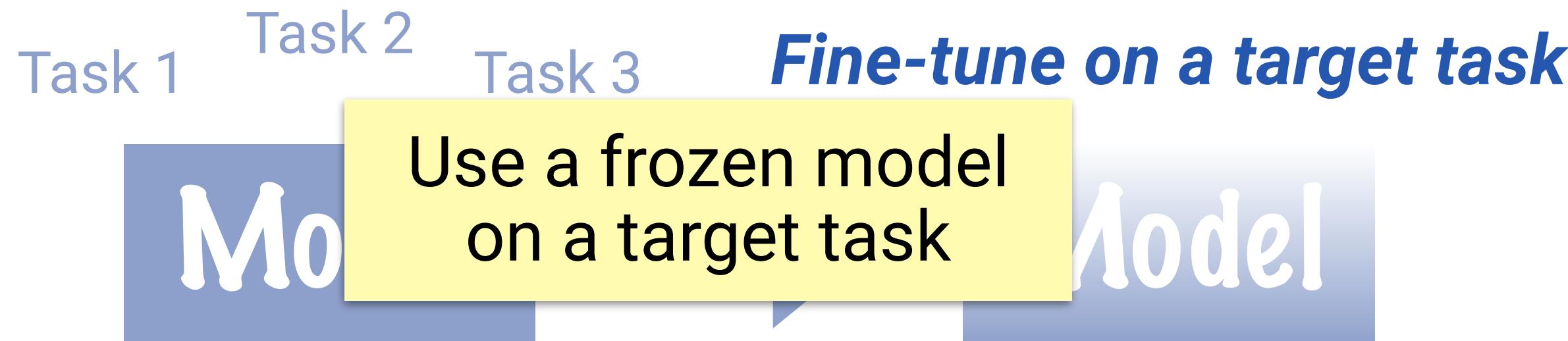


Define the task as the input



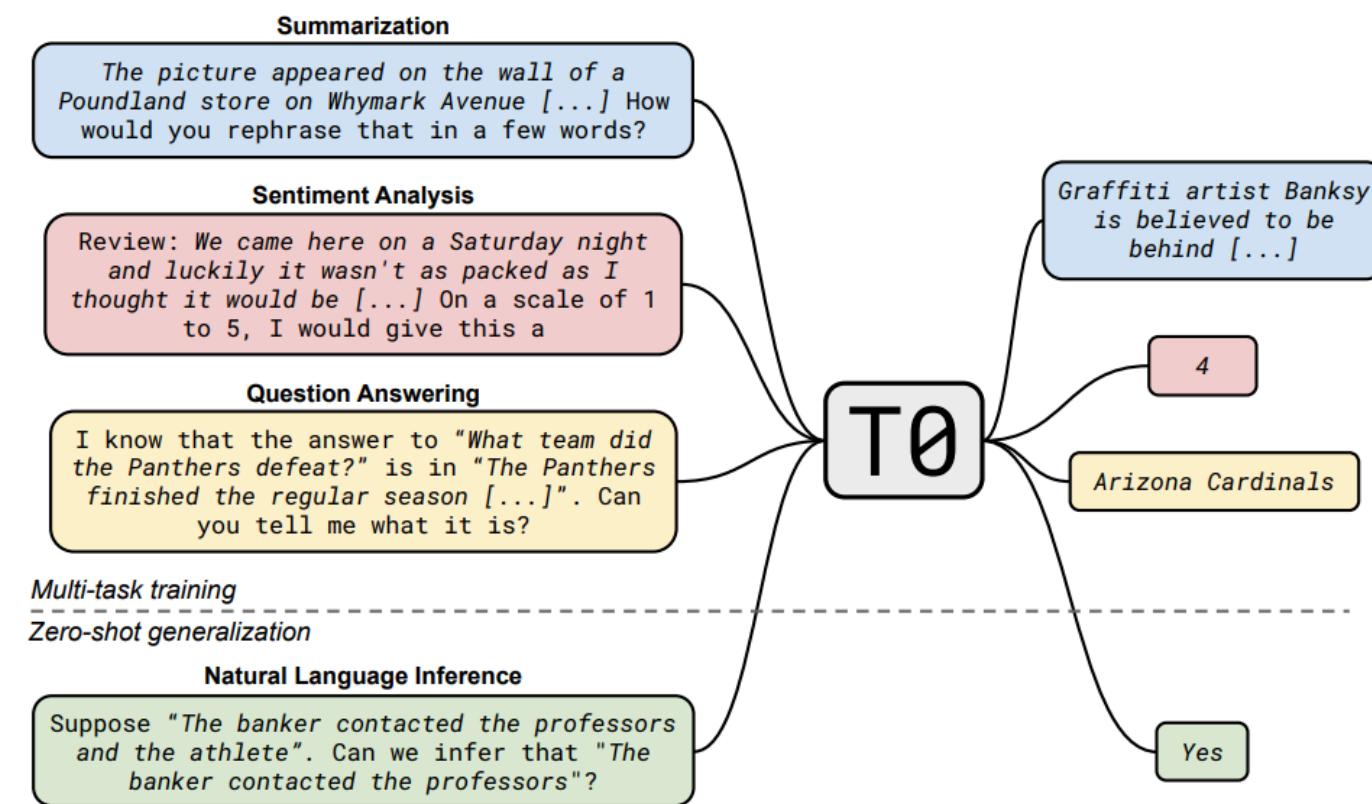
Natural language description  
vs. k-shot data

# Related Work



## Multi-task learning/Meta-learning

(Villa and Drissi 2002, Evgeniou and Pontil 2004, Ruder 2017, Wang et al. 2020, Aghajanyan et al. 2021)



## Instruction-reading models

(Mishra et al. 2021, Wei et al. 2021, Sanh et al. 2021, Zhong et al. 2021)

Figure from Sanh et al. 2021



Define the task as the input

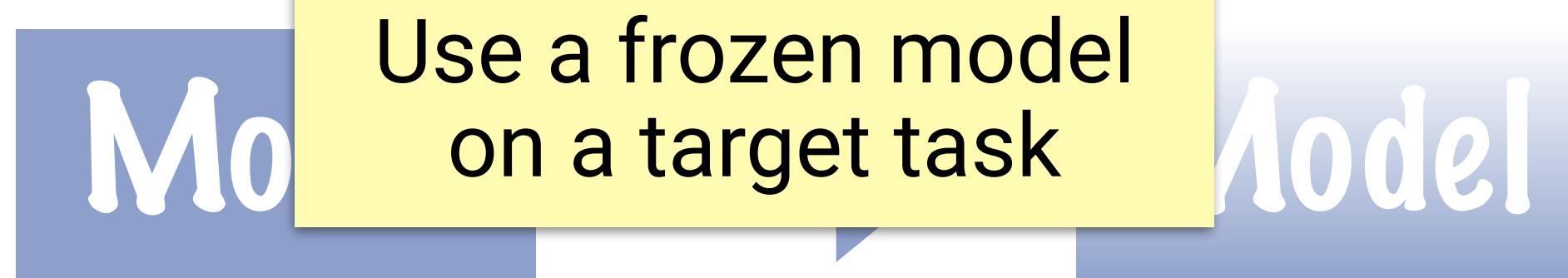


Natural language description  
vs.  $k$ -shot data

Empirically, MetaICL can be  
better/complementary

# Related Work

Task 1 Task 2 Task 3 **Fine-tune on a target task**



## Multi-task learning/Meta-learning

(Villa and Drissi 2002, Evgeniou and Pontil 2004, Ruder 2017, Wang et al. 2020, Aghajanyan et al. 2021)

Instruction: “*Is the comment positive?*”

$x_1$ : “*Good movie!*”  $y_1$ : “*yes*”

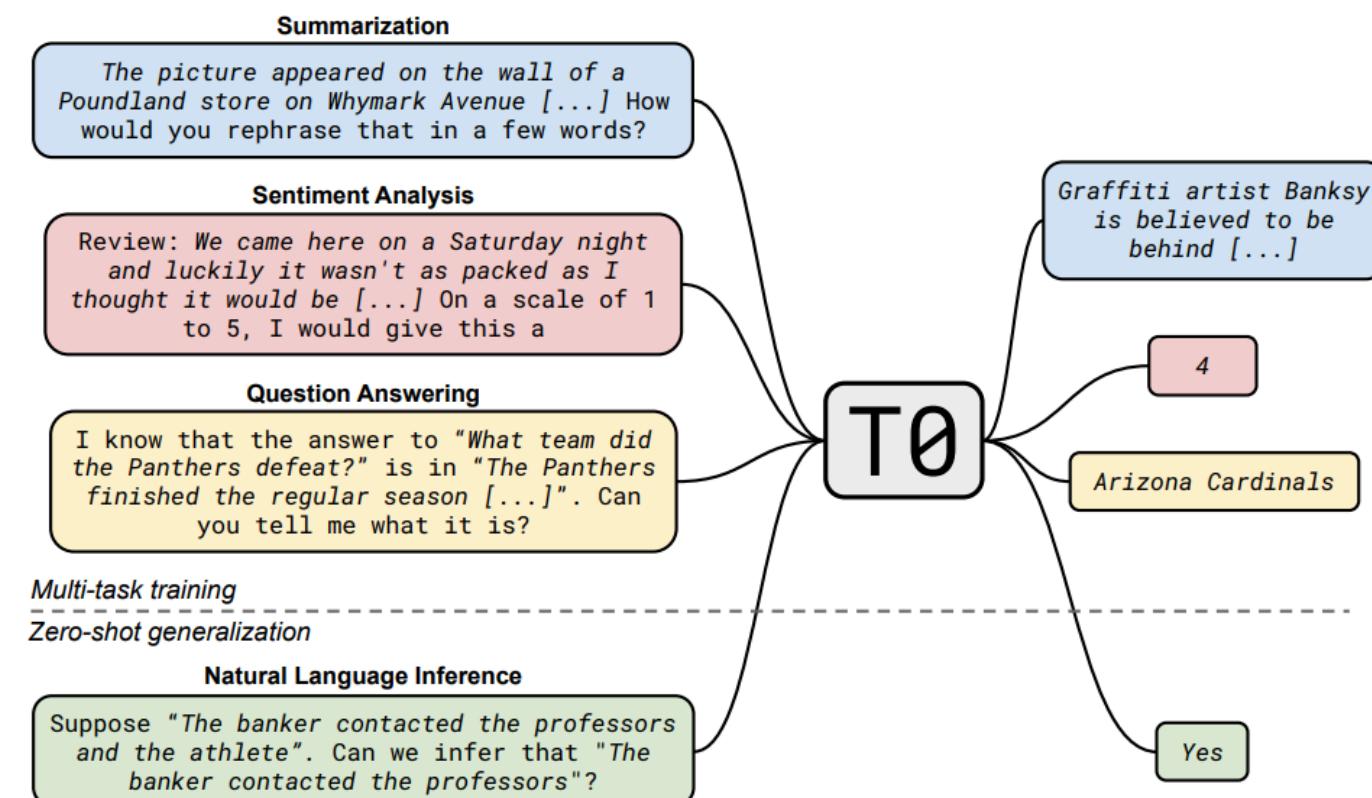
$x_2$ : “*Bad movie!*”  $y_2$ : “*no*”



## Concurrent work

Chen et al. 2021

“Meta-learning via Language Model In-context Tuning”



## Instruction-reading models

(Mishra et al. 2021, Wei et al. 2021, Sanh et al. 2021, Zhong et al. 2021)

Figure from Sanh et al. 2021



Define the task as the input



Natural language description  
vs. k-shot data

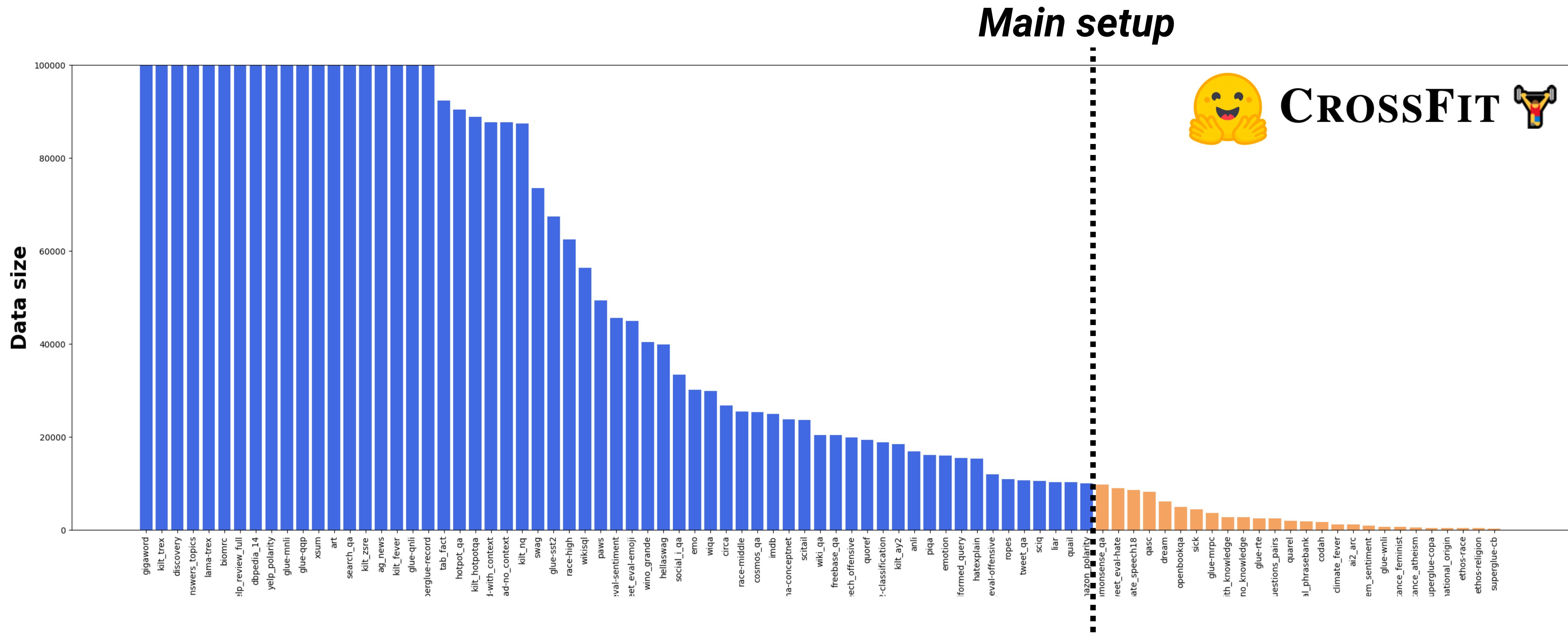
Empirically, MetaICL can be  
better/complementary

# Experimental Setup

In total, experiment with **142** NLP datasets and **7** different metatrain-target splits

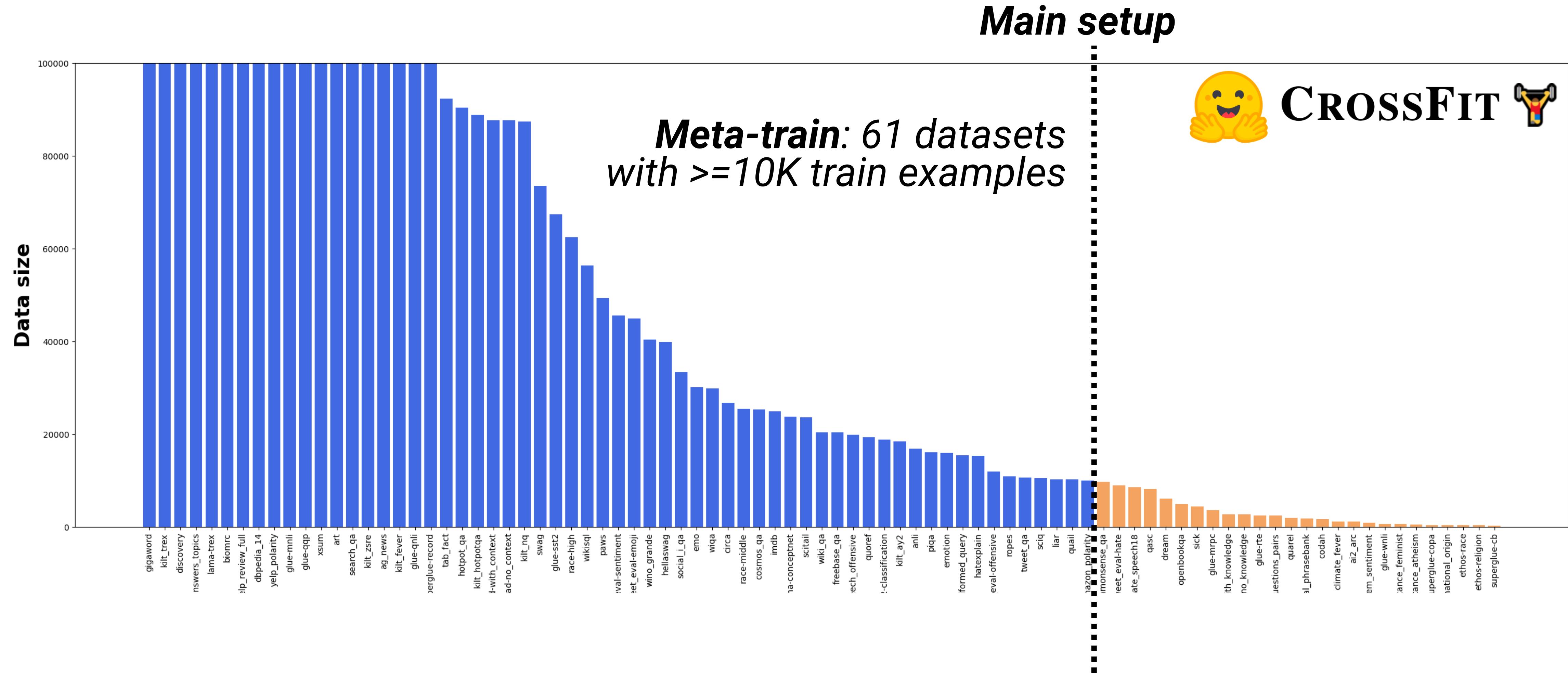
# Experimental Setup

In total, experiment with 142 NLP datasets and 7 different metatrain-target splits



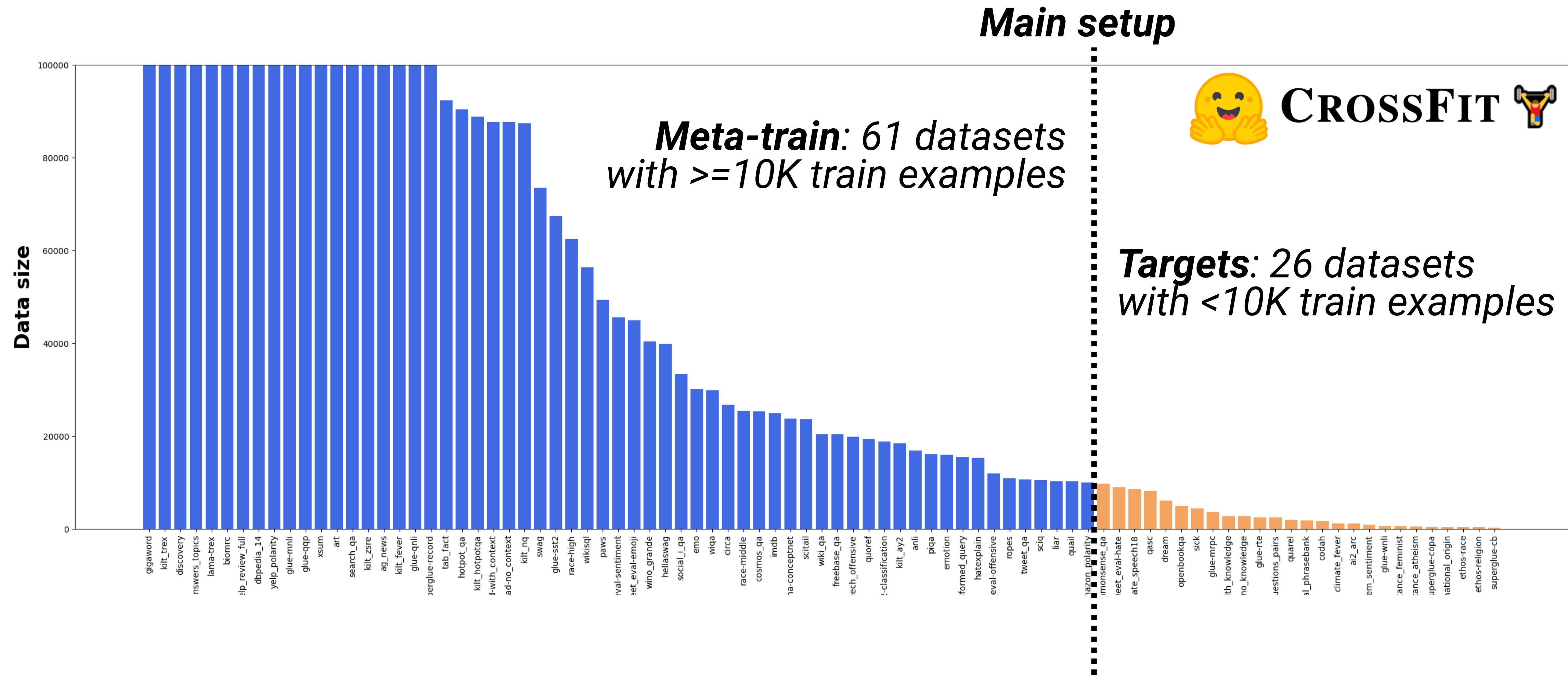
# Experimental Setup

In total, experiment with 142 NLP datasets and 7 different metatrain-target splits



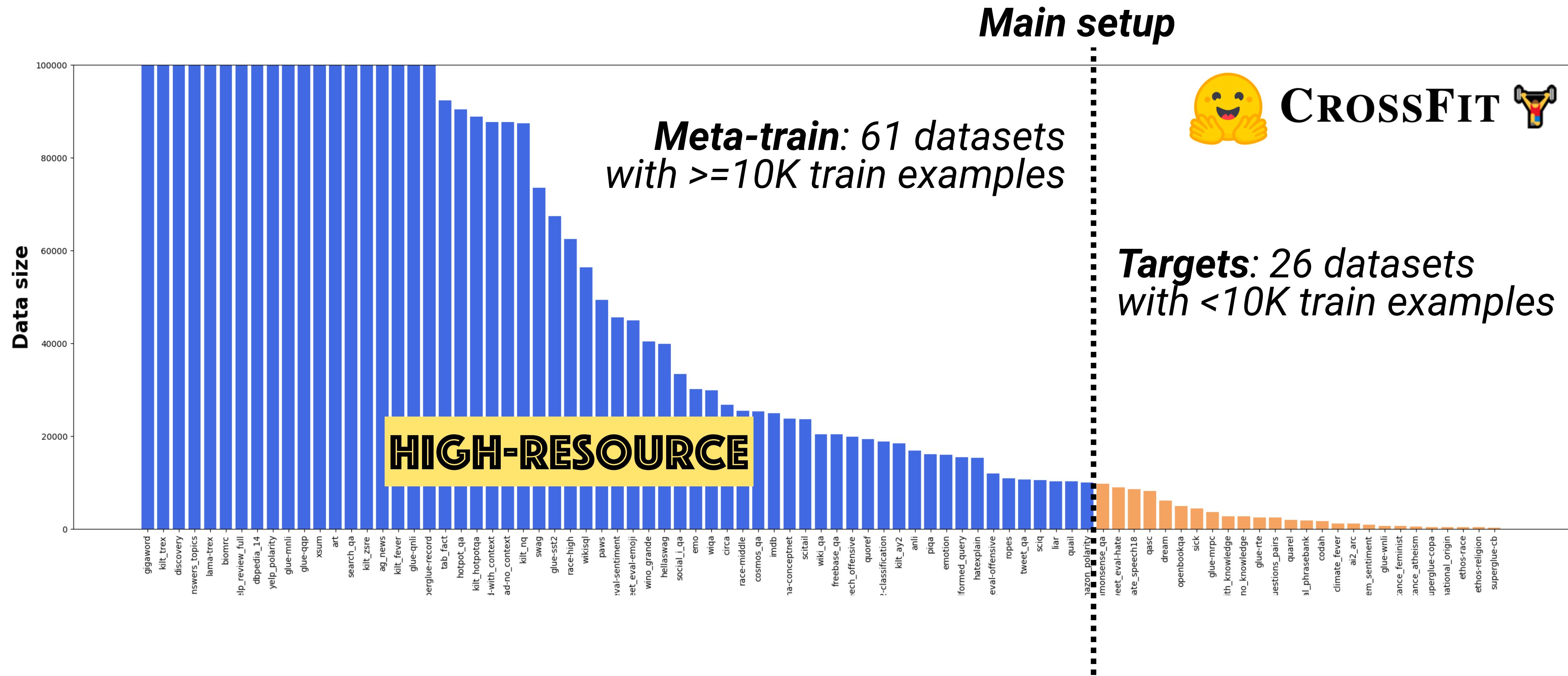
# Experimental Setup

In total, experiment with 142 NLP datasets and 7 different metatrain-target splits



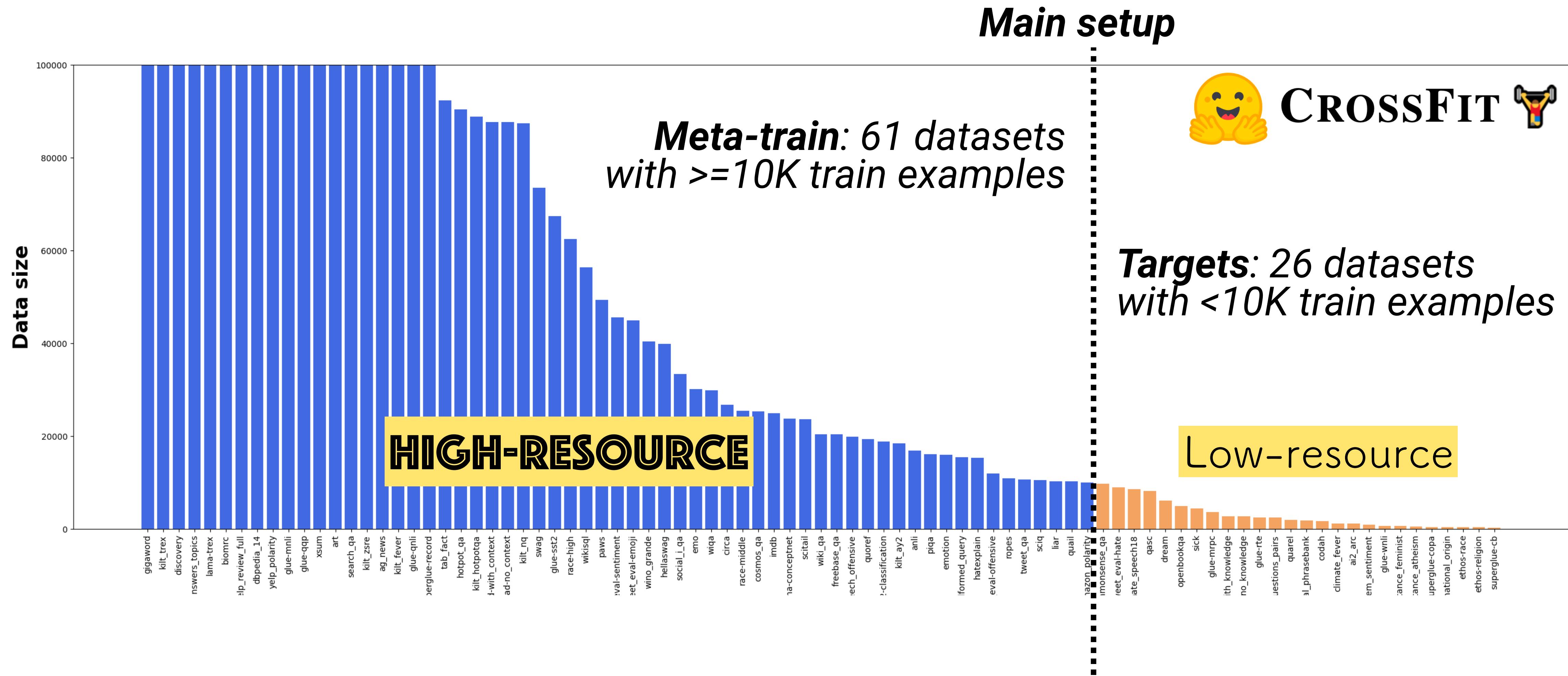
# Experimental Setup

In total, experiment with 142 NLP datasets and 7 different metatrain-target splits

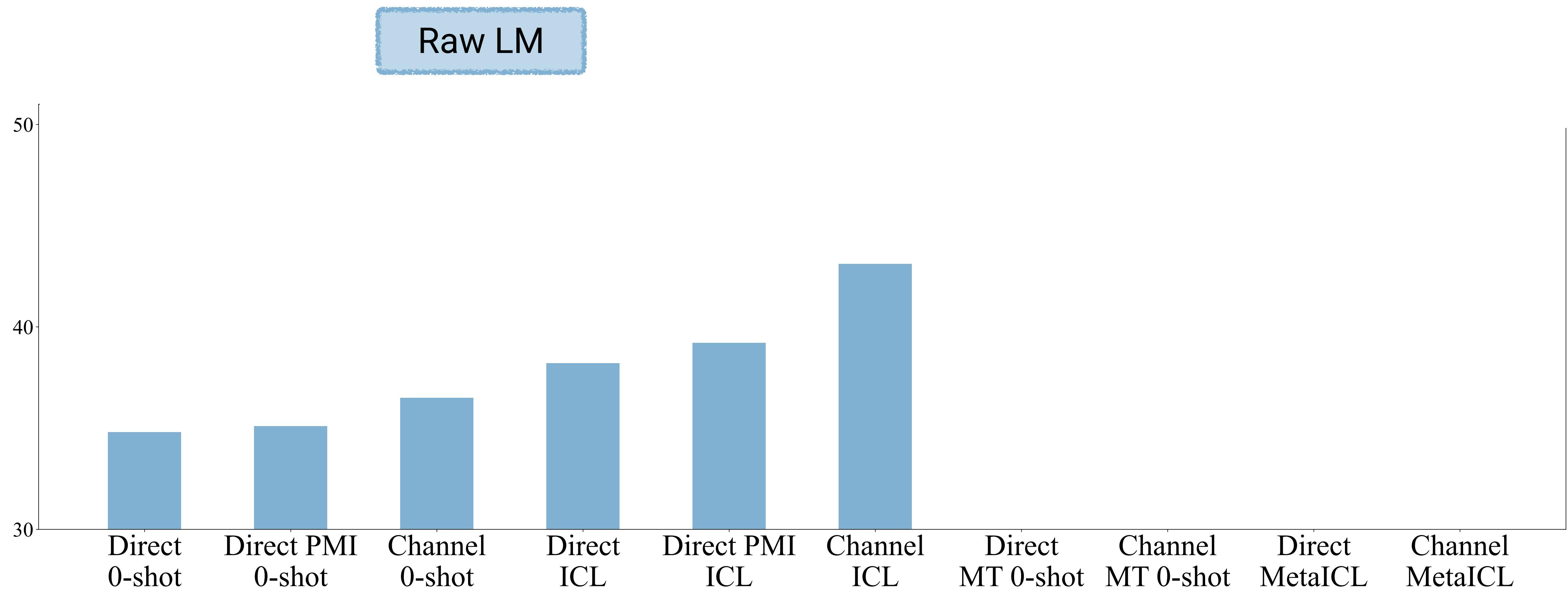


# Experimental Setup

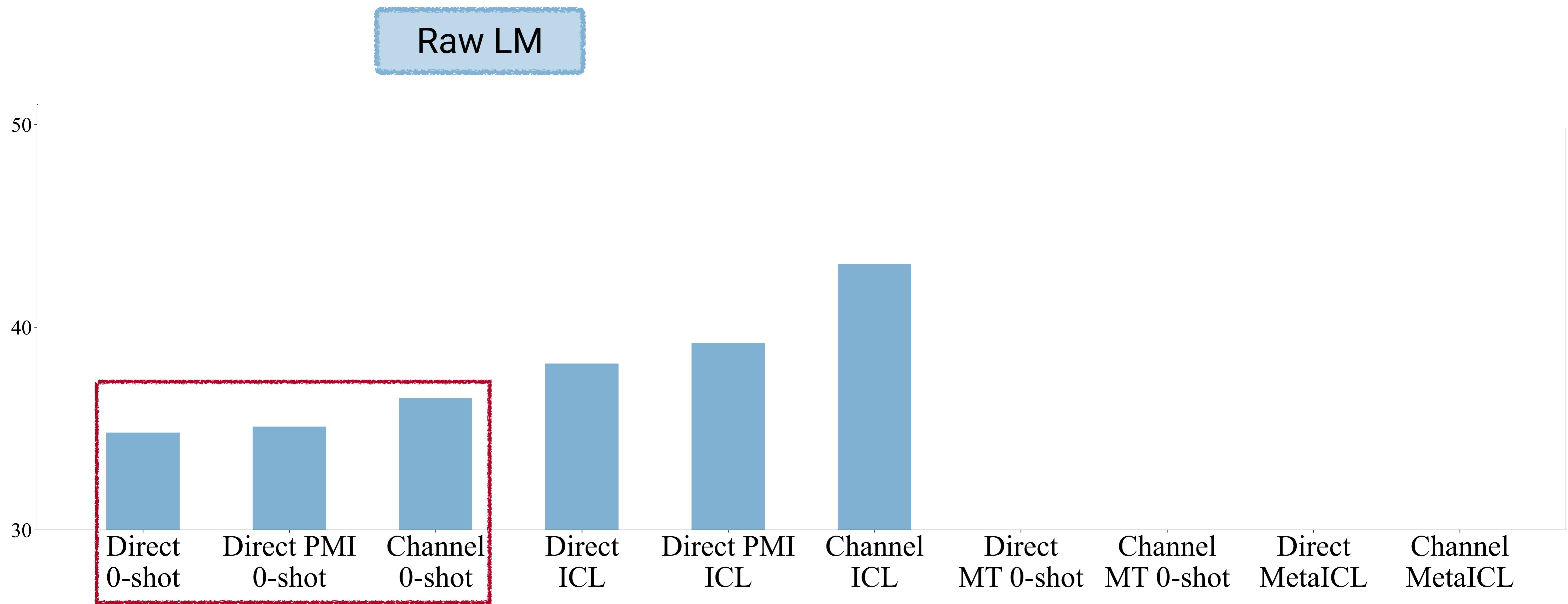
In total, experiment with 142 NLP datasets and 7 different metatrain-target splits



# Results



# Results

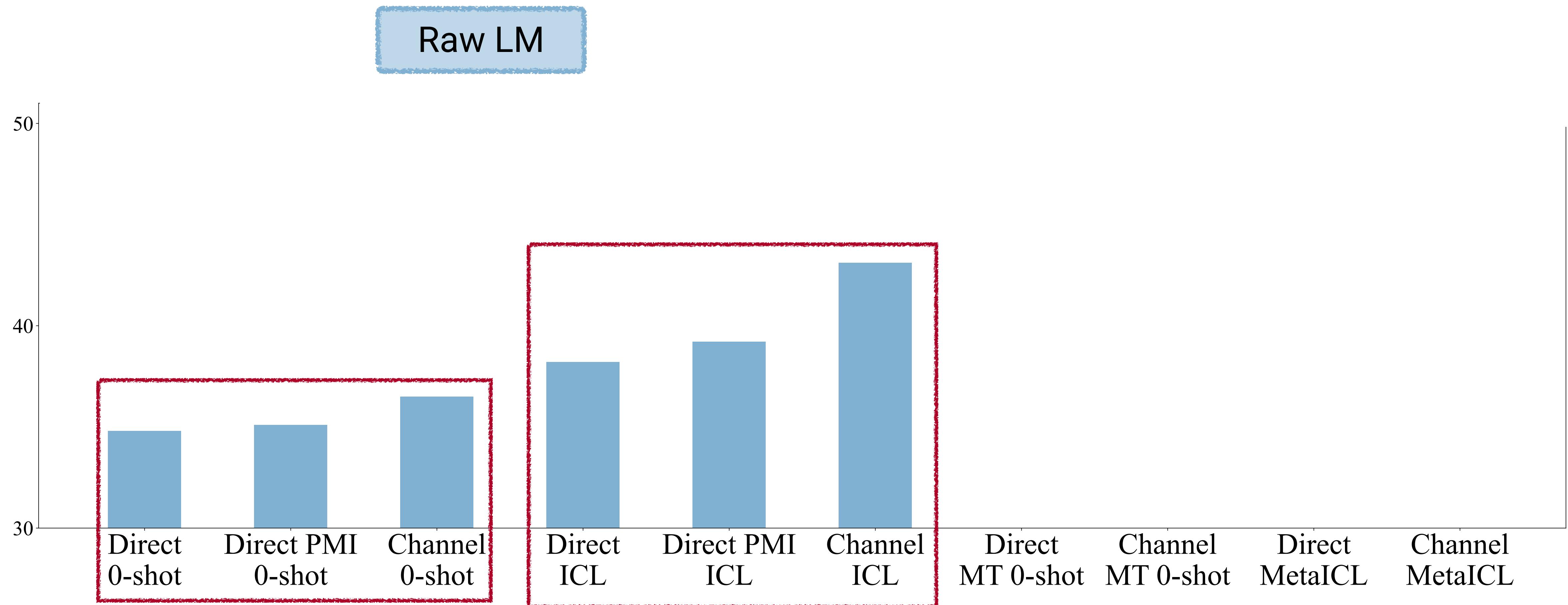


\* Direct: from Brown et al. 2020

\* Direct PMI: from Holtzman et al 2021, Zhao et al. 2021

\* Channel: from Min et al. 2022

# Results

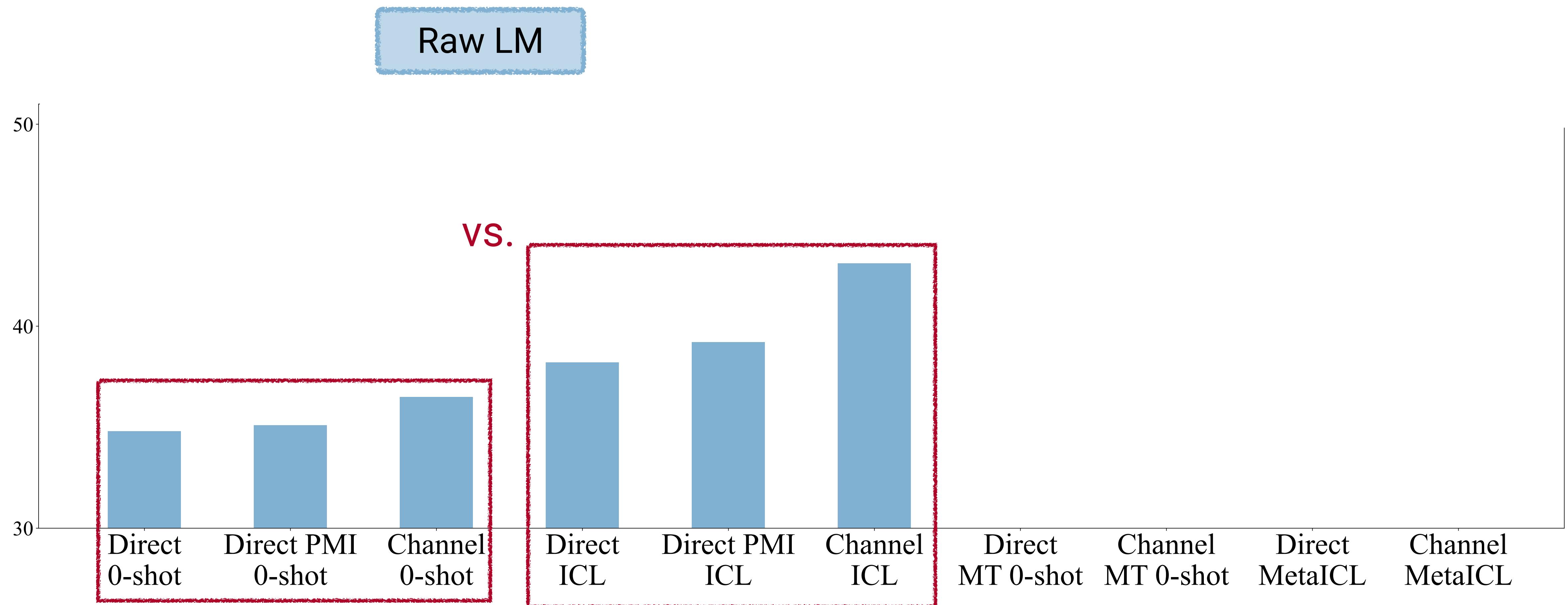


\* Direct: from Brown et al. 2020

\* Direct PMI: from Holtzman et al 2021, Zhao et al. 2021

\* Channel: from Min et al. 2022

# Results

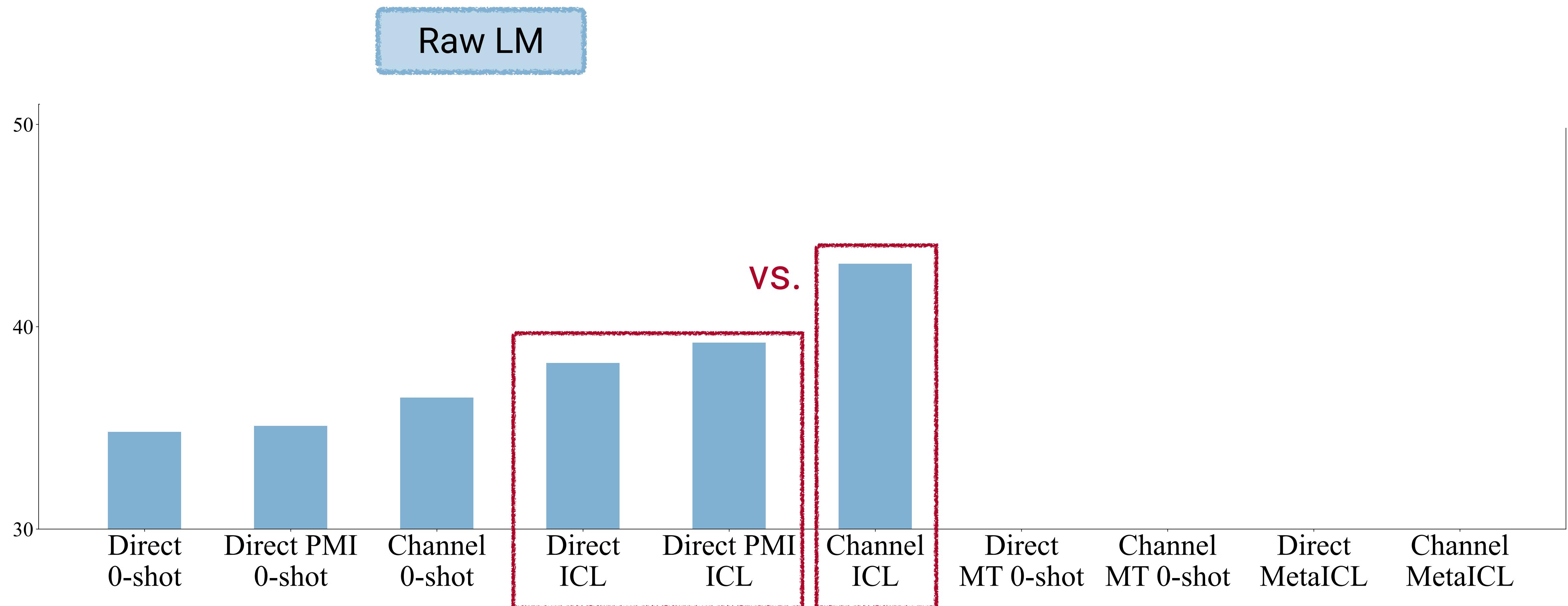


\* Direct: from Brown et al. 2020

\* Direct PMI: from Holtzman et al 2021, Zhao et al. 2021

\* Channel: from Min et al. 2022

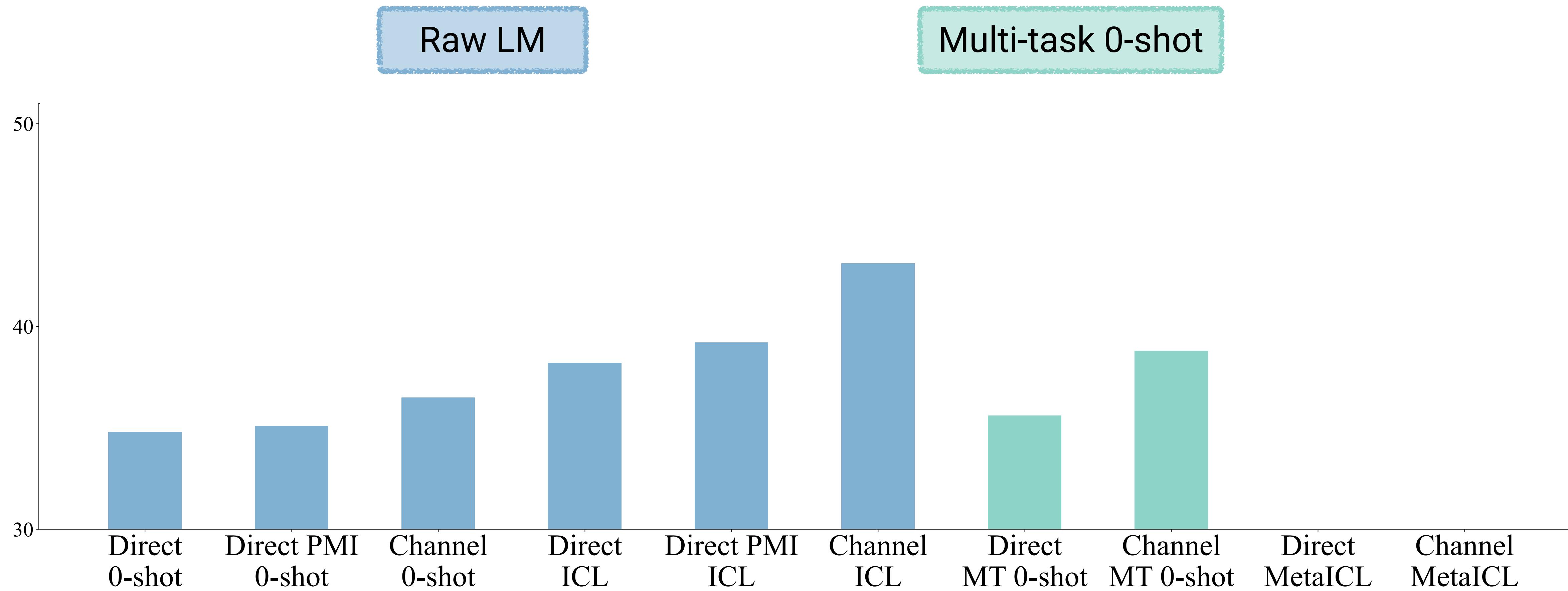
# Results



Channel ICL achieves the best result when no meta-training

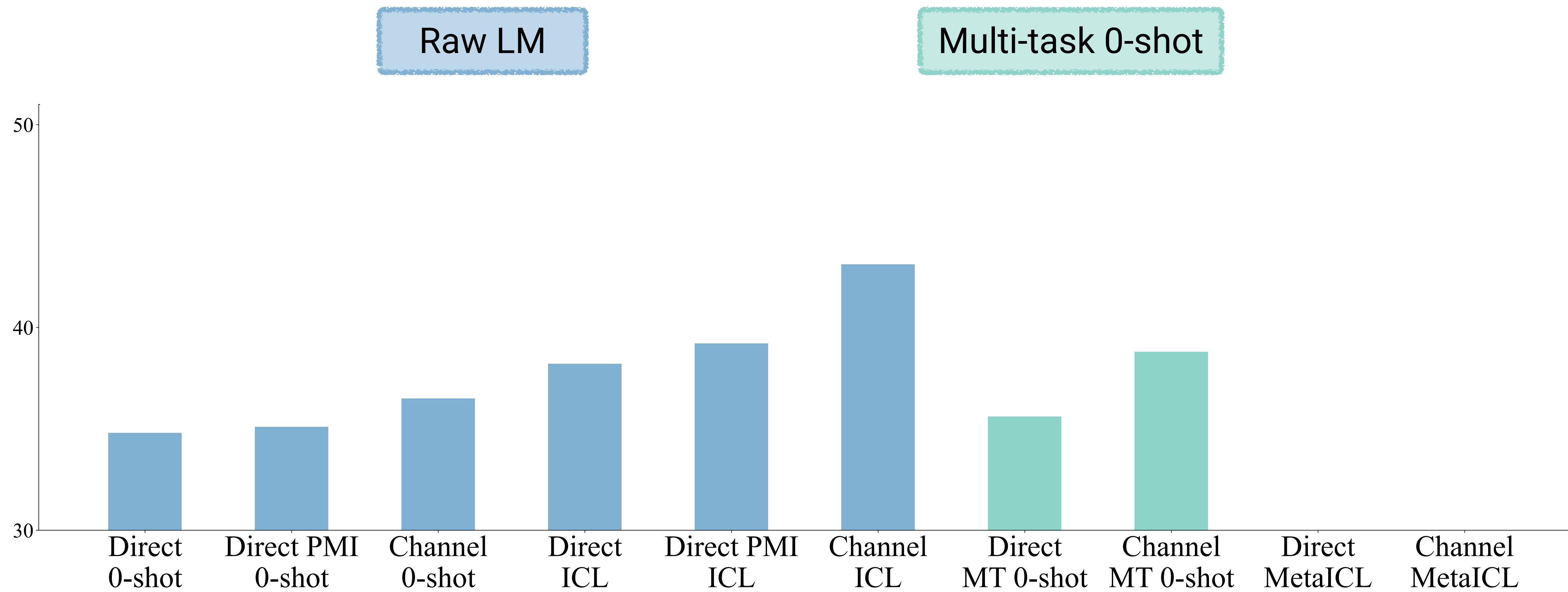
# Results

Multi-task training  
but **w/o** in-context learning



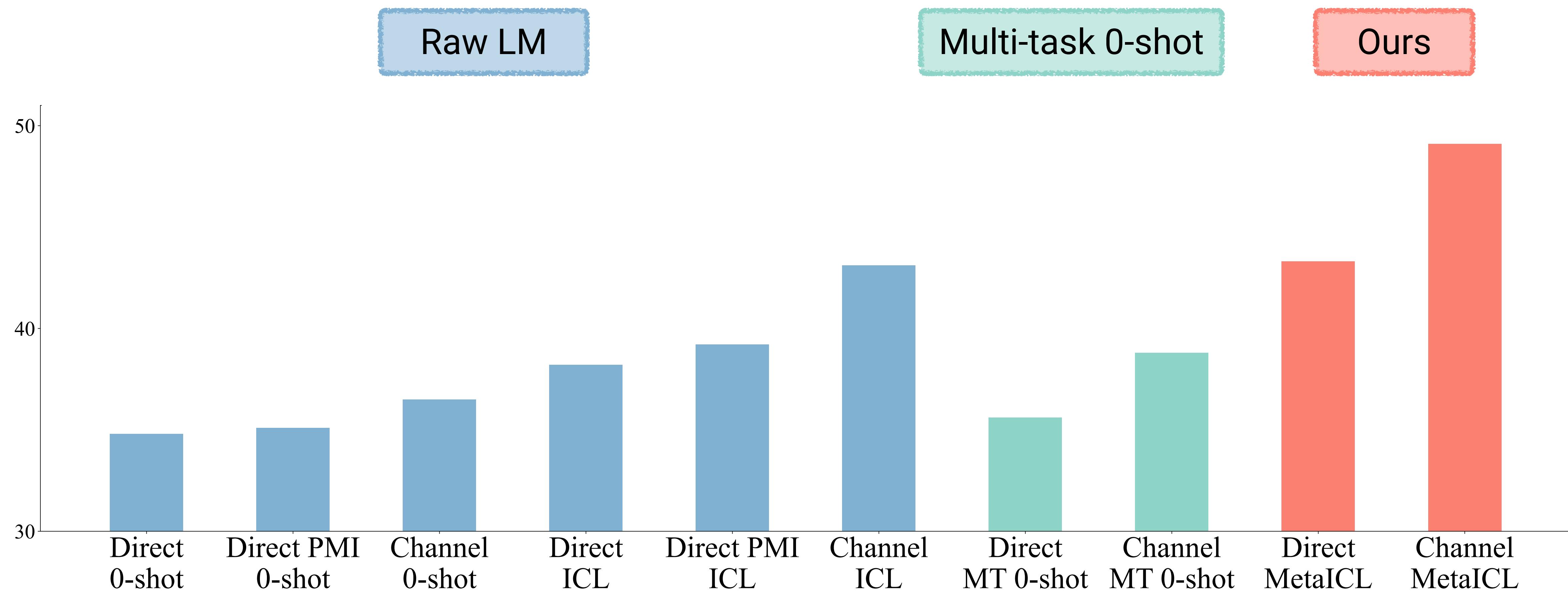
# Results

Multi-task training  
but *w/o* in-context learning



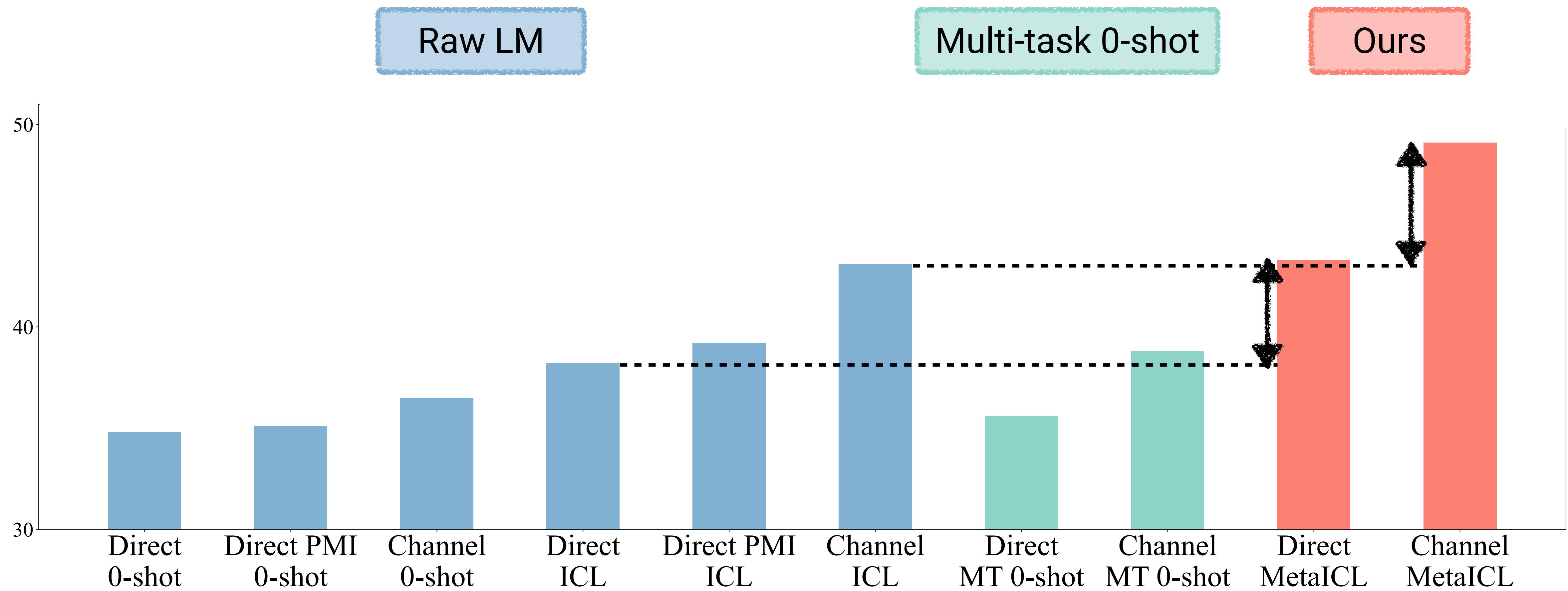
Multi-task 0-shot is not as competitive

# Results

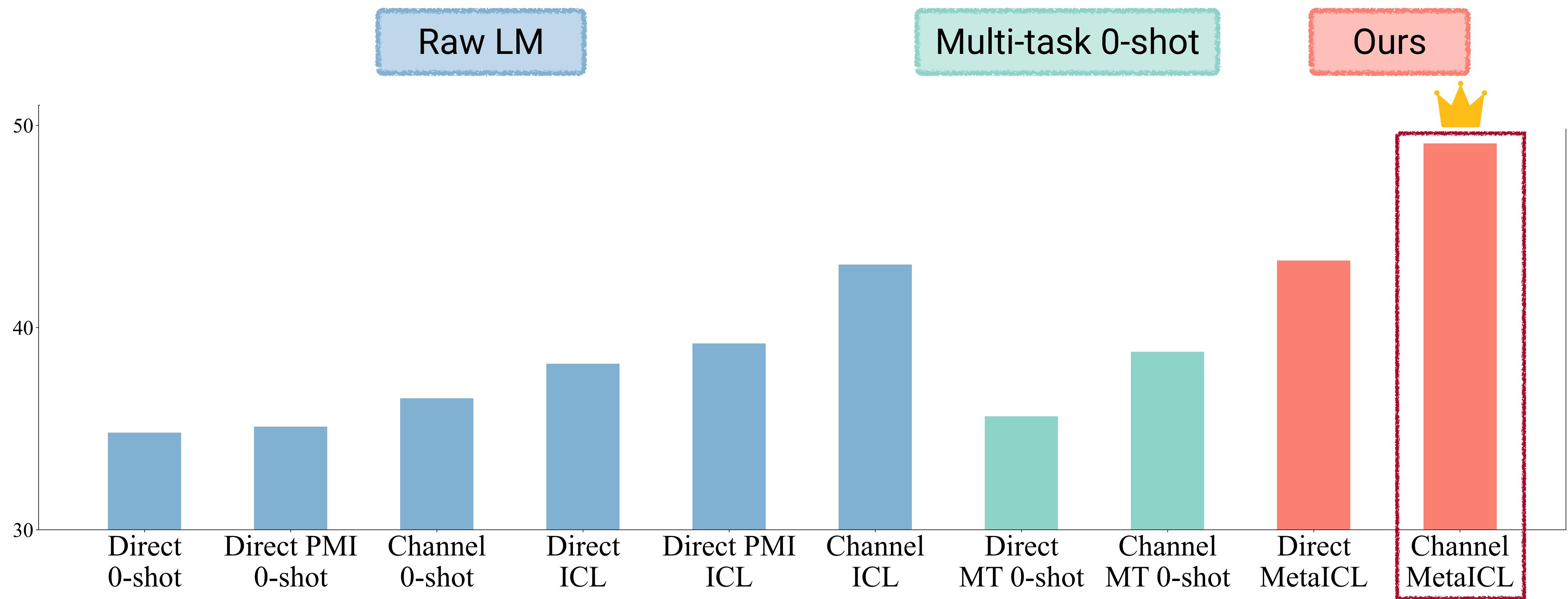


# Results

Base LM: GPT-2 Large



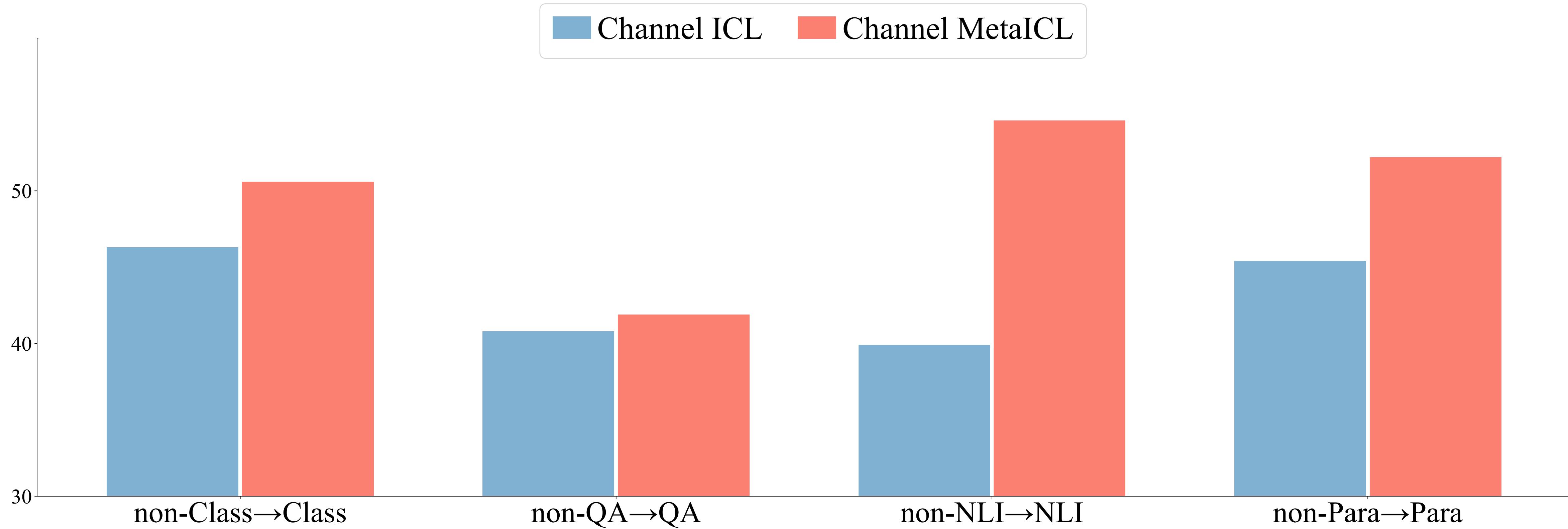
# Results



Channel MetaICL achieves the best result

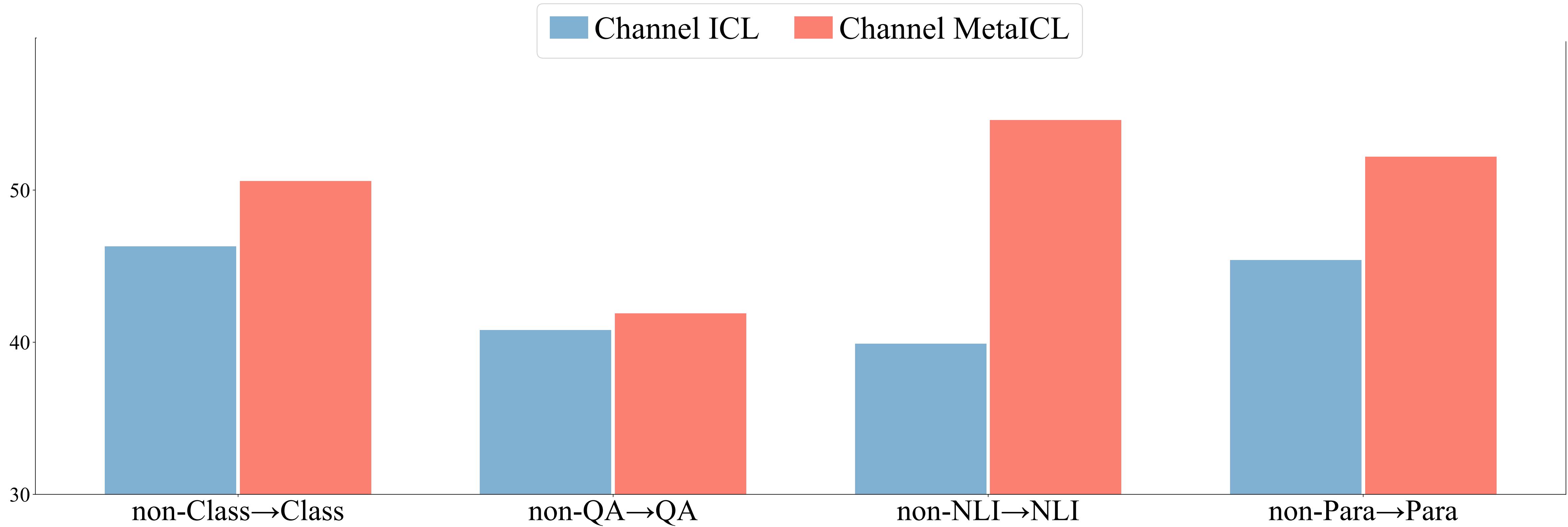
# Results

## Other meta-training/targets splits



# Results

## Other meta-training/targets splits



Large gains – even with no similar tasks at meta-training

# Results

## vs. Instruction reading models

# Results

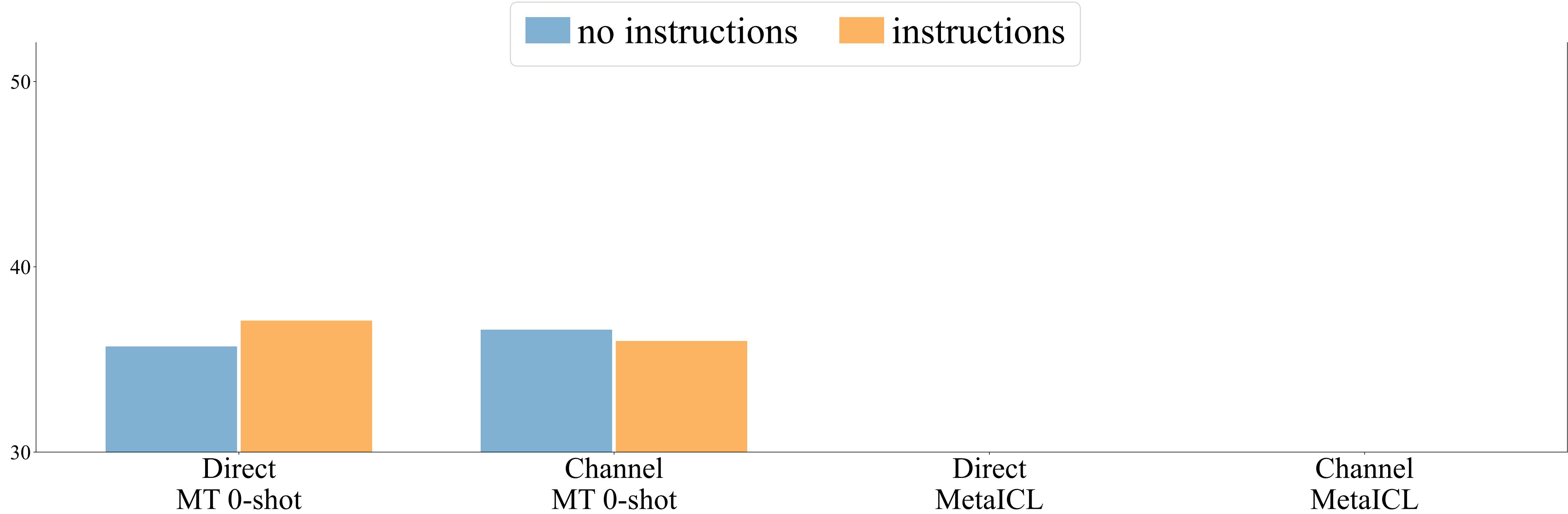
## vs. Instruction reading models

Used the data from T0 (San et al 2021) – 8.3 instructions per task

# Results

## vs. Instruction reading models

Used the data from T0 (San et al 2021) – 8.3 instructions per task

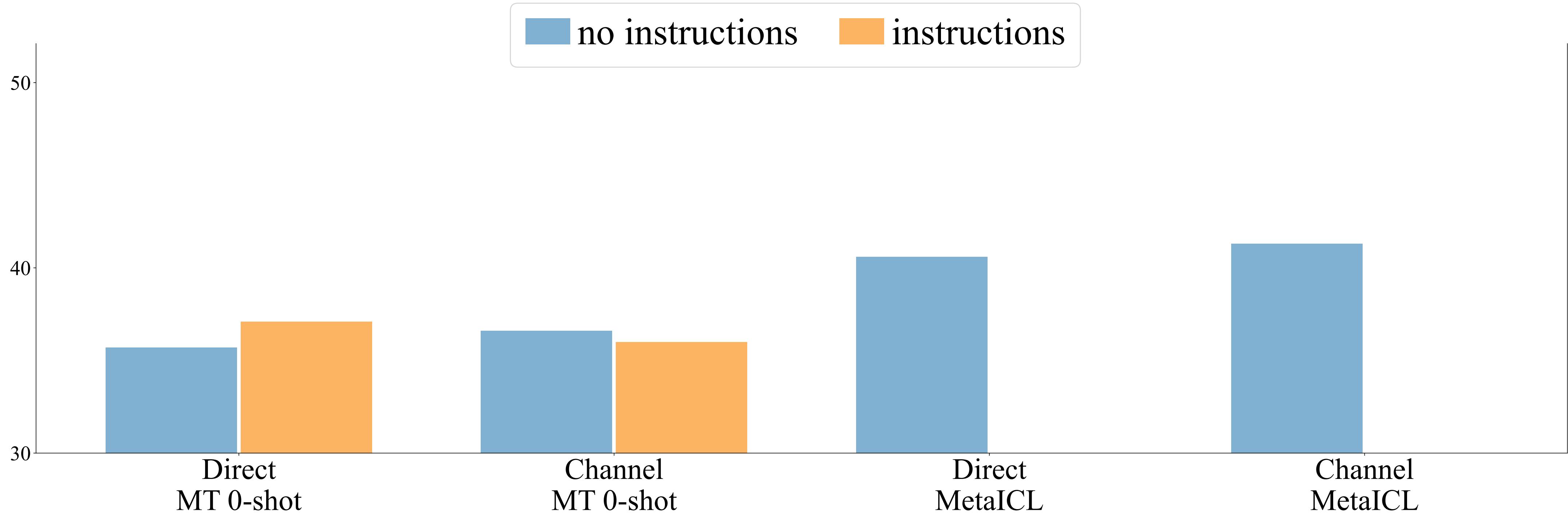


Multi-task model w/o in-context learning benefits from instructions

# Results

## vs. Instruction reading models

Used the data from T0 (San et al 2021) – 8.3 instructions per task

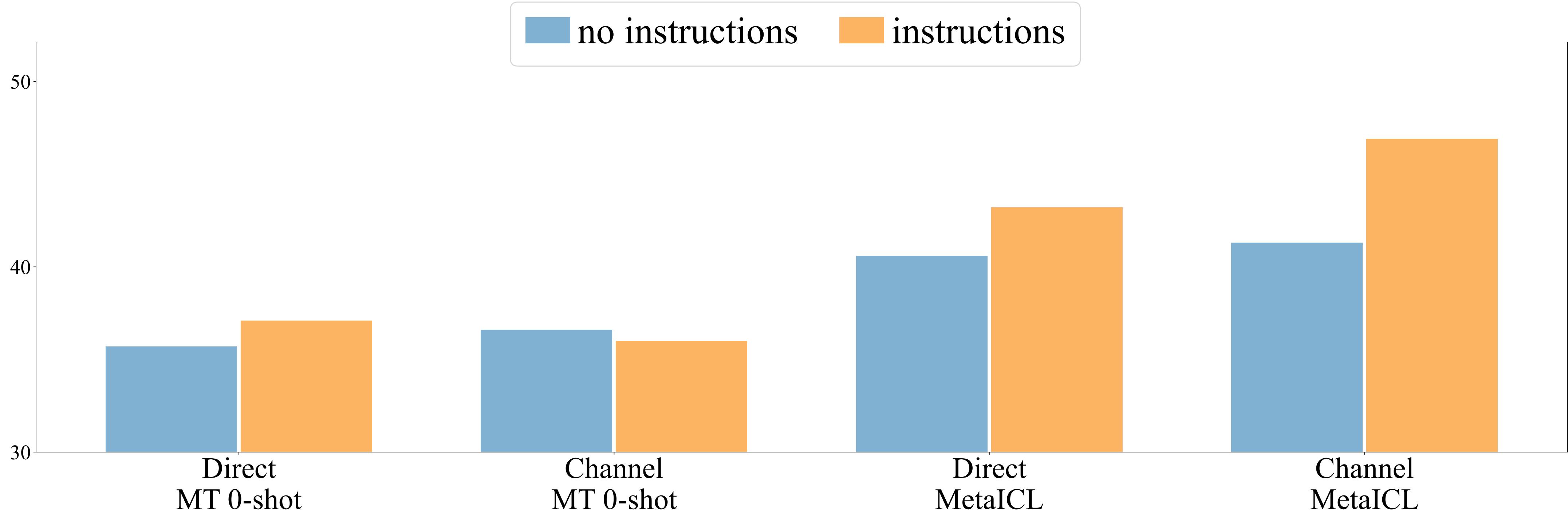


Training to in-context learn is better than training to read instructions

# Results

## vs. Instruction reading models

Used the data from T0 (San et al 2021) – 8.3 instructions per task



MetaICL + instructions reading achieves the best performance

# Why does MetalCL work?

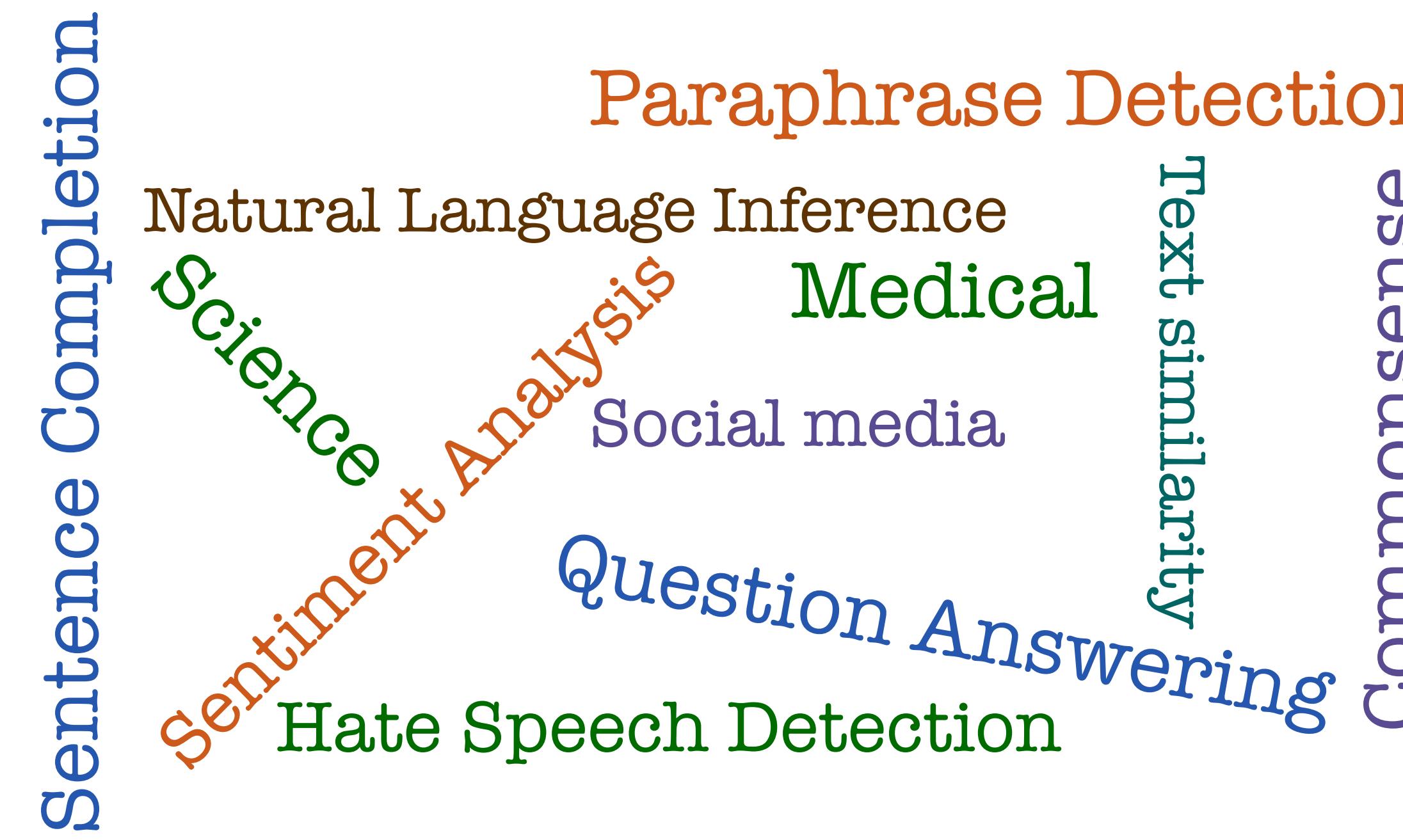
# Why does MetalCL work?

1) The model benefits from *related* tasks during meta-training



# Why does MetalCL work?

1) The model benefits from *related* tasks during meta-training



Already verified by a large body of work in multi-task learning

(Villa and Drissi 2002, Evgeniou and Pontil 2004, Ruder 2017, Aghajanyan et al. 2021)

# Why does MetalCL work?

1) The model benefits from *related* tasks during meta-training



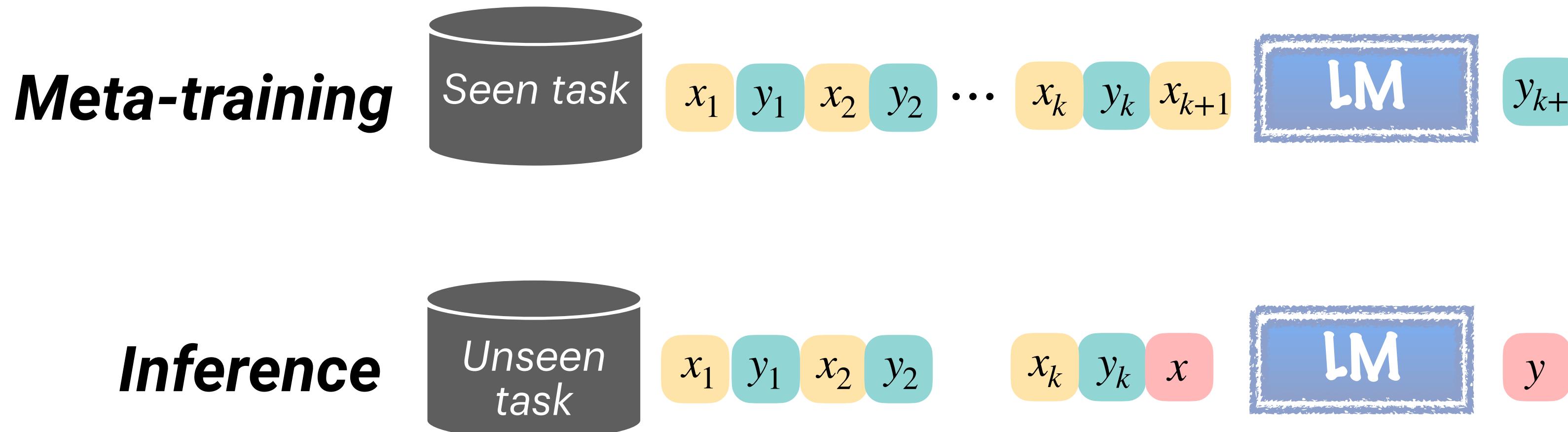
Already verified by a large body of work in multi-task learning

(Villa and Drissi 2002, Evgeniou and Pontil 2004, Ruder 2017, Aghajanyan et al. 2021)

However, Multi-task 0-shot is not enough

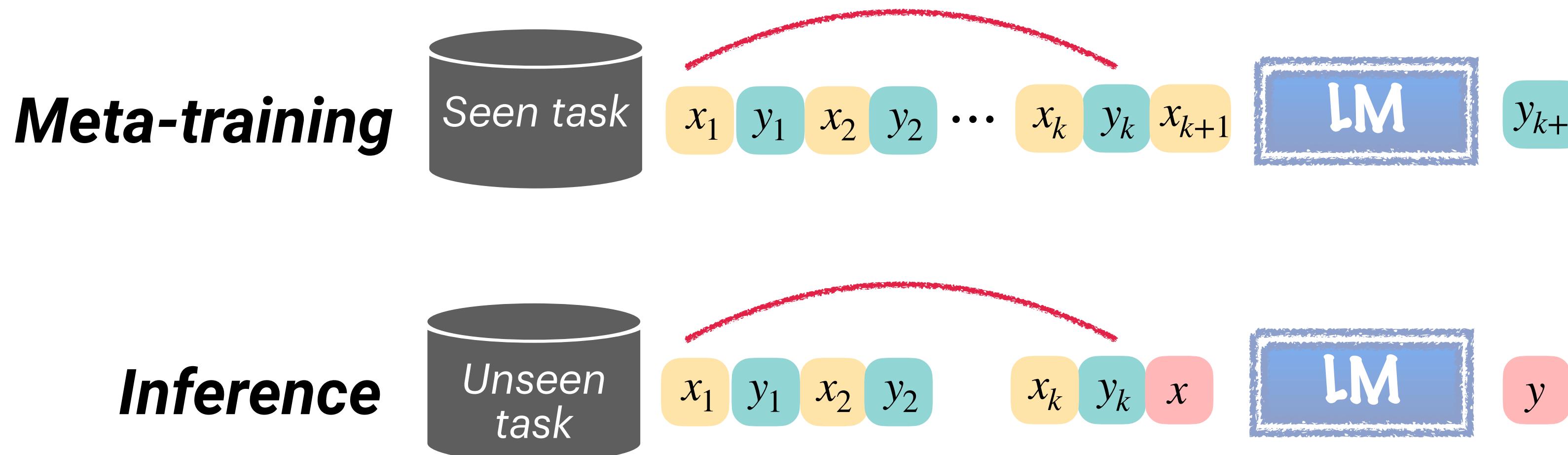
# Why does MetalCL work?

2) The model learns the ability to *in-context learn a given task*



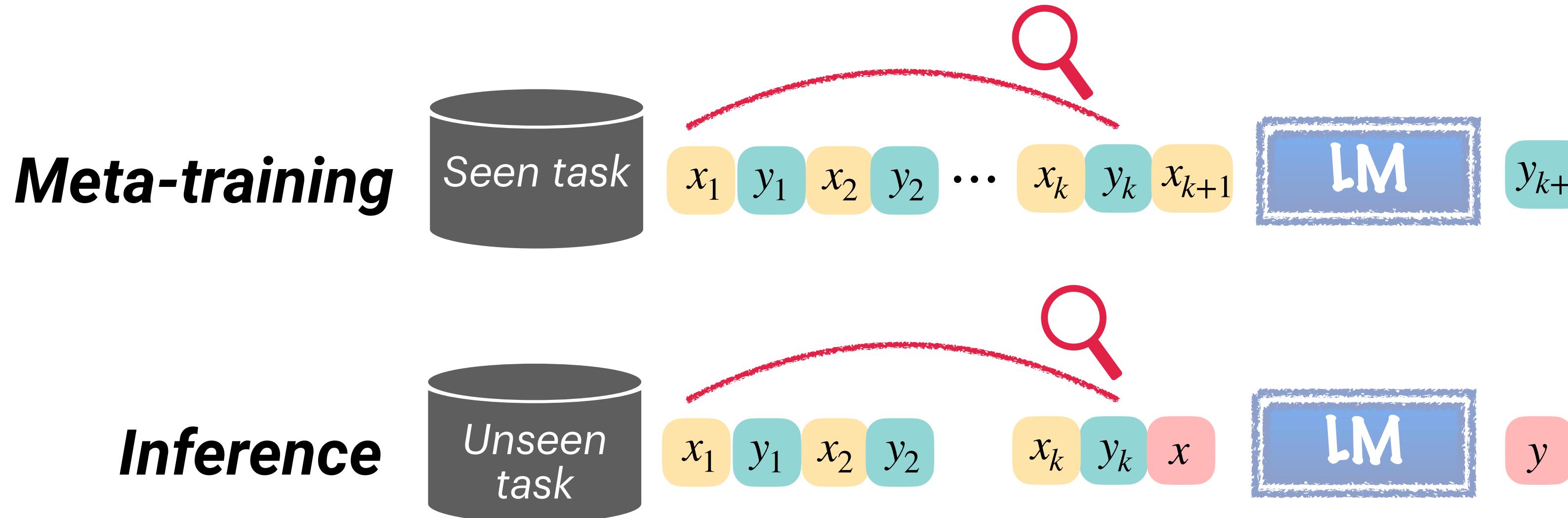
# Why does MetalCL work?

2) The model learns the ability to *in-context learn a given task*



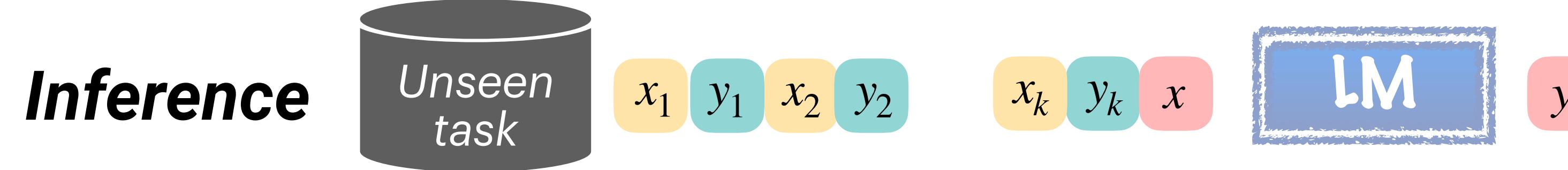
# Why does MetalCL work?

2) The model learns the ability to *in-context learn a given task*



# Why does MetalCL work?

2) The model learns the ability to *in-context learn a given task*

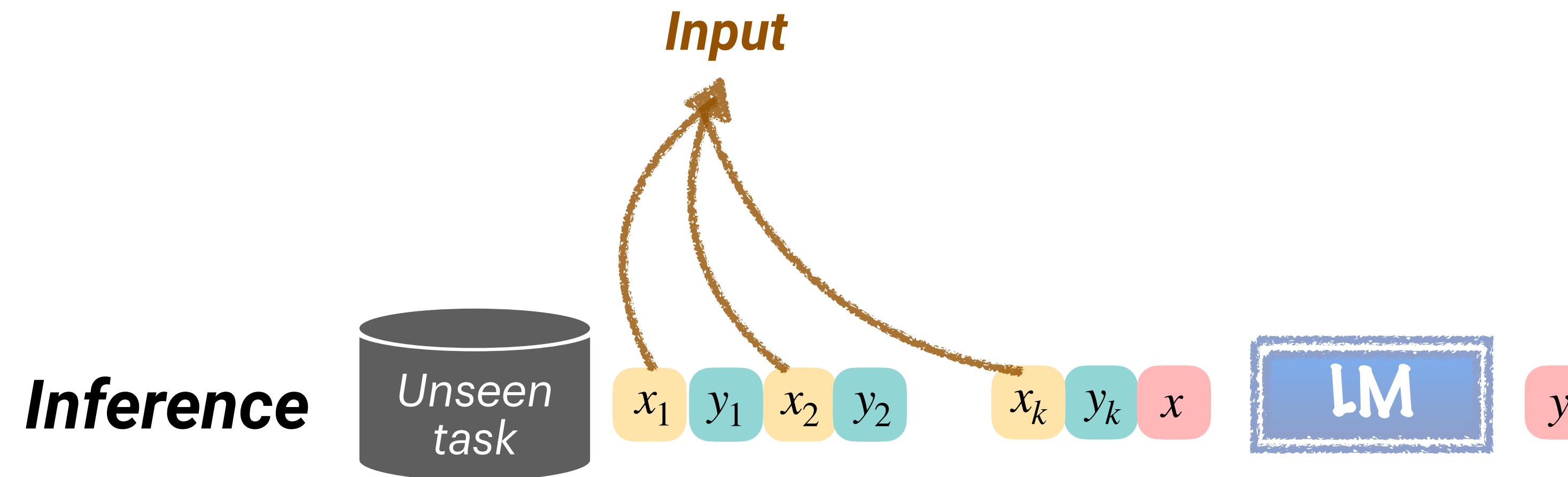


“*Meta-training encourages exploiting **input distribution & output space**,*

– Min et al. 2022. “Rethinking the Role of Demonstrations: What Makes In-Context Learning Work?”

# Why does MetalCL work?

2) The model learns the ability to *in-context learn a given task*

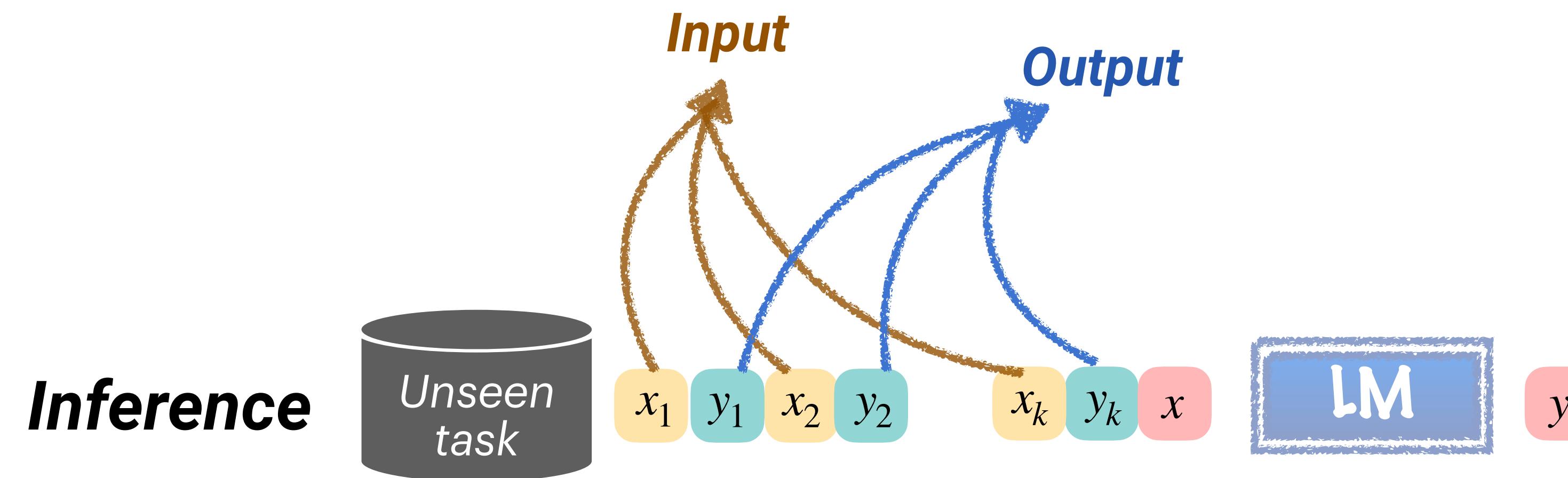


“Meta-training encourages exploiting **input distribution & output space**,

– Min et al. 2022. “Rethinking the Role of Demonstrations: What Makes In-Context Learning Work?”

# Why does MetalCL work?

2) The model learns the ability to *in-context learn a given task*

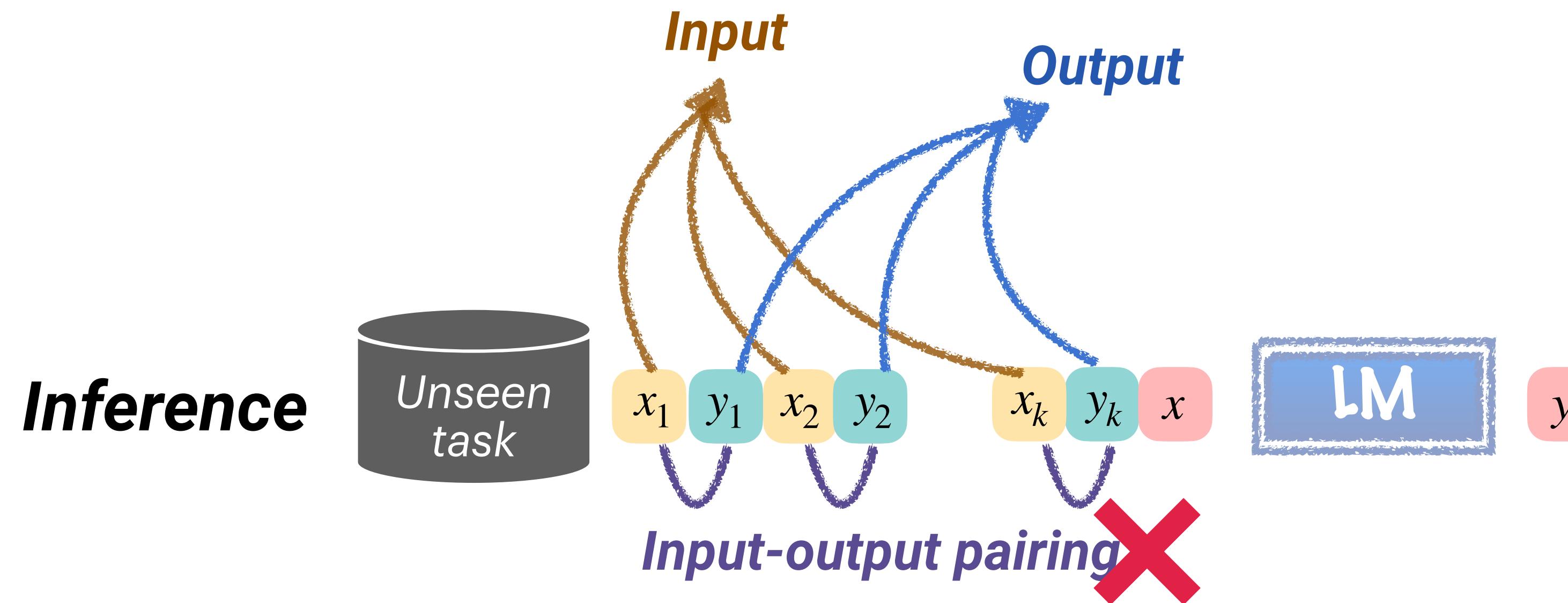


“Meta-training encourages exploiting **input distribution** & **output space**,

– Min et al. 2022. “Rethinking the Role of Demonstrations: What Makes In-Context Learning Work?”

# Why does MetalCL work?

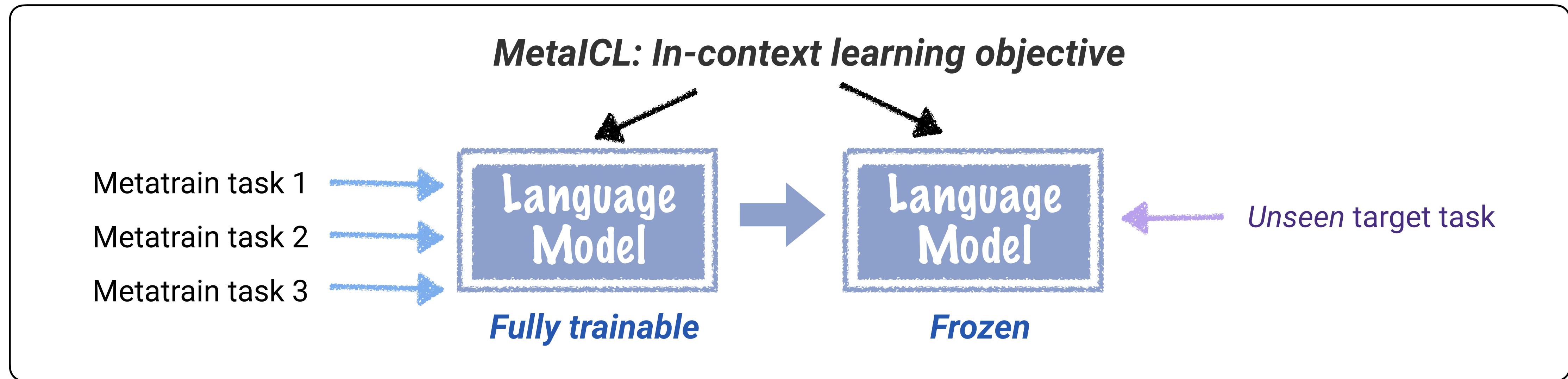
2) The model learns the ability to *in-context learn a given task*



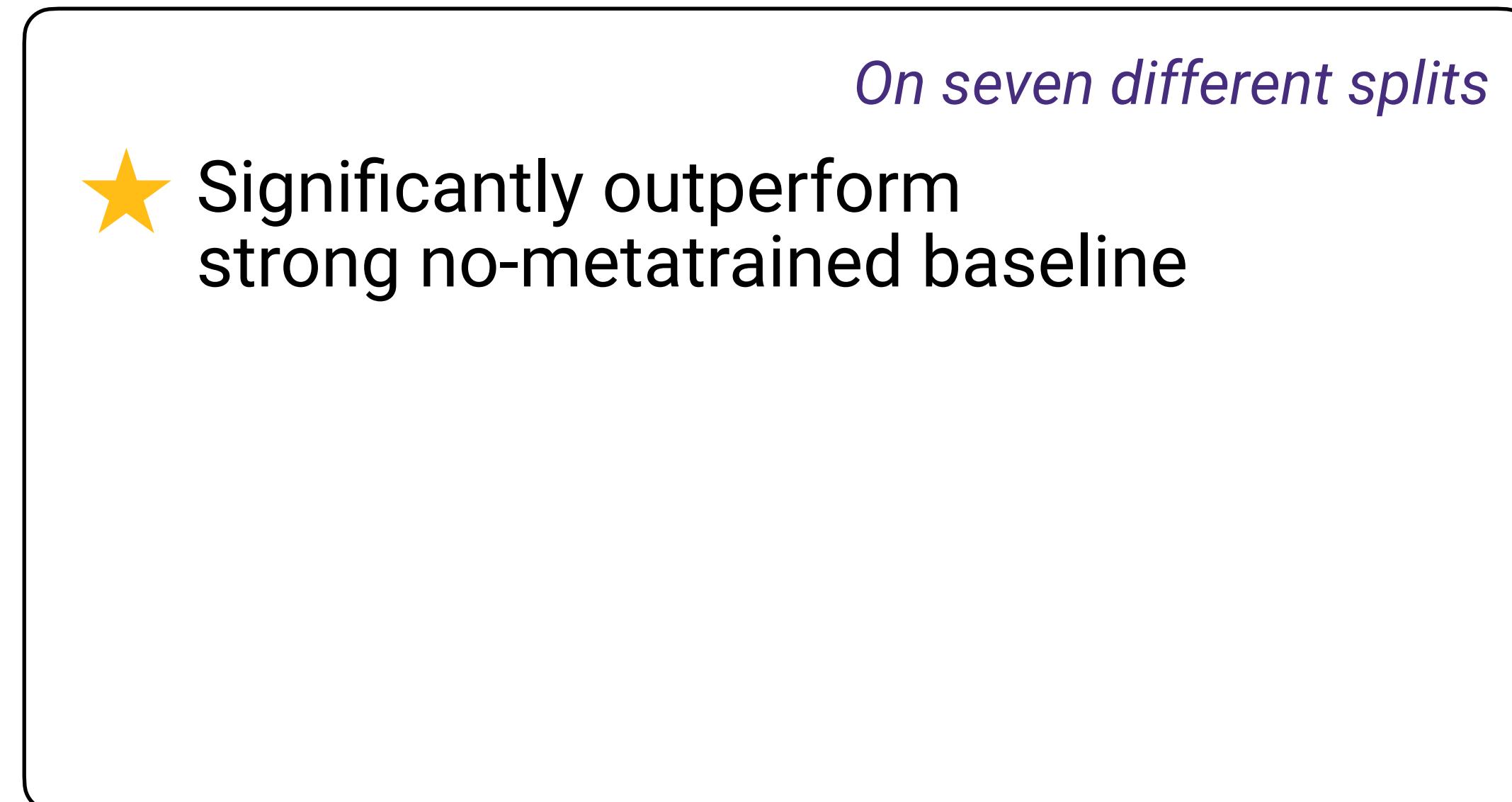
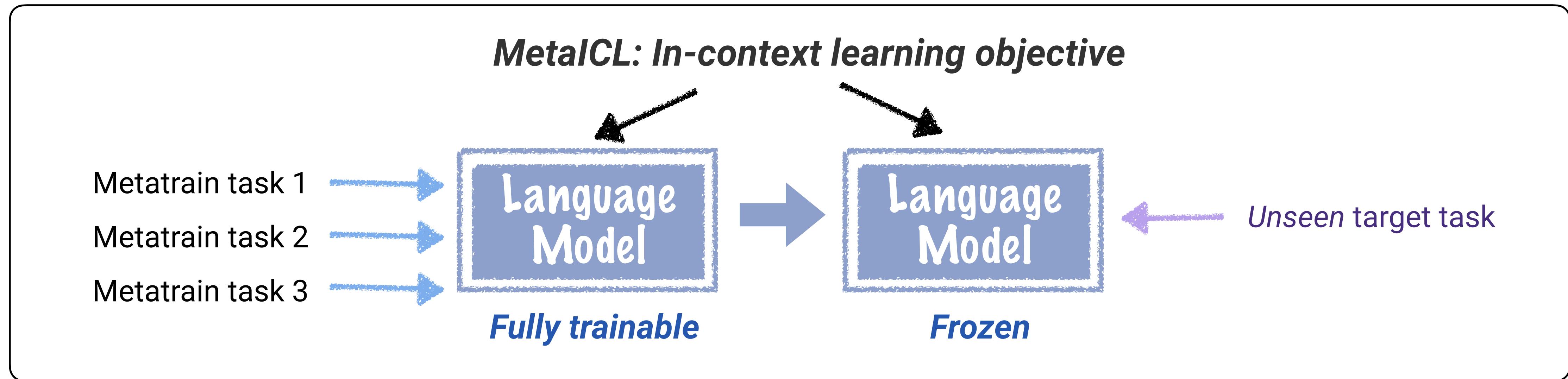
“Meta-training encourages exploiting **input distribution & output space**, and ignoring **input-output pairing**”

– Min et al. 2022. “Rethinking the Role of Demonstrations: What Makes In-Context Learning Work?”

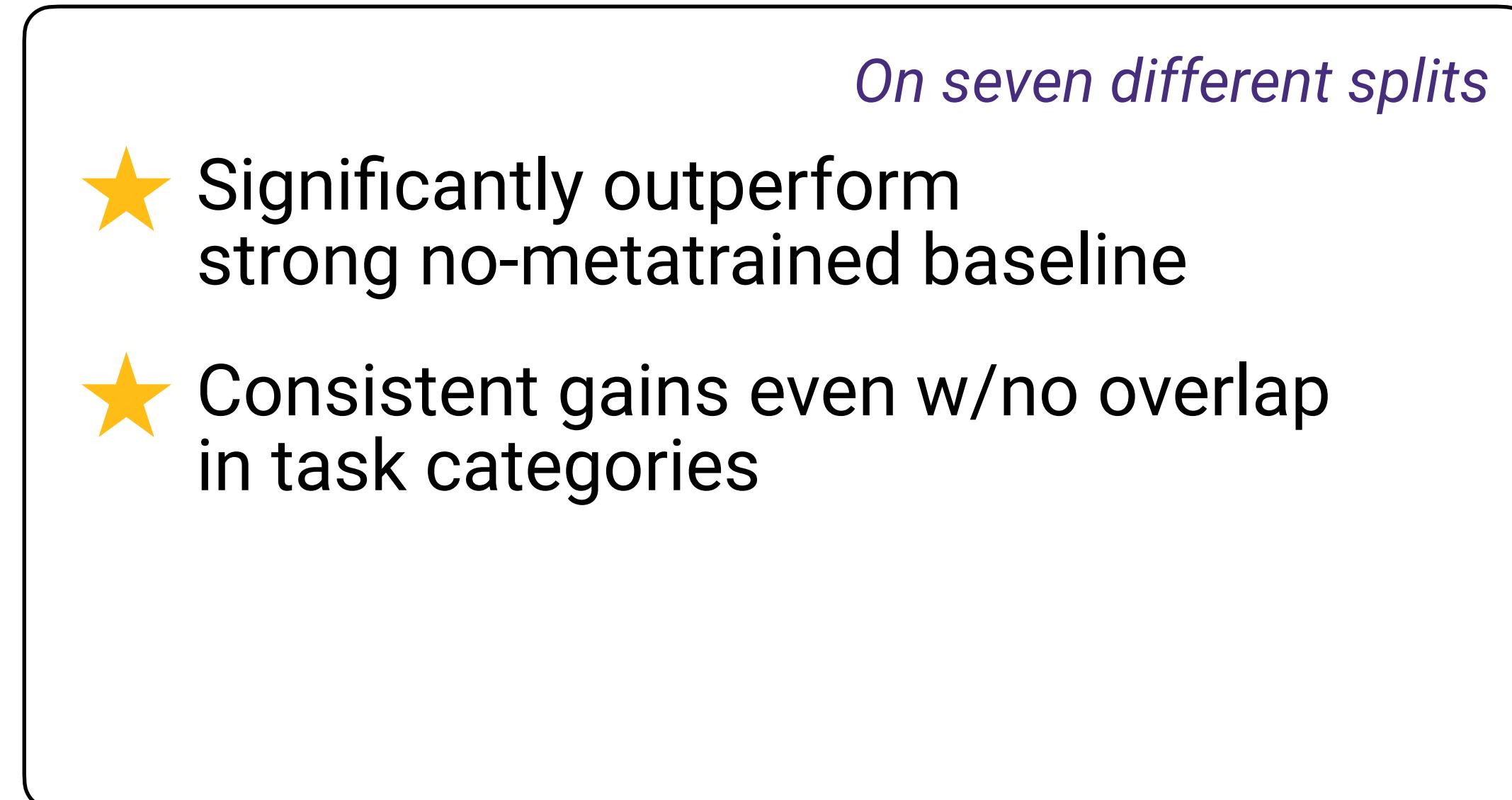
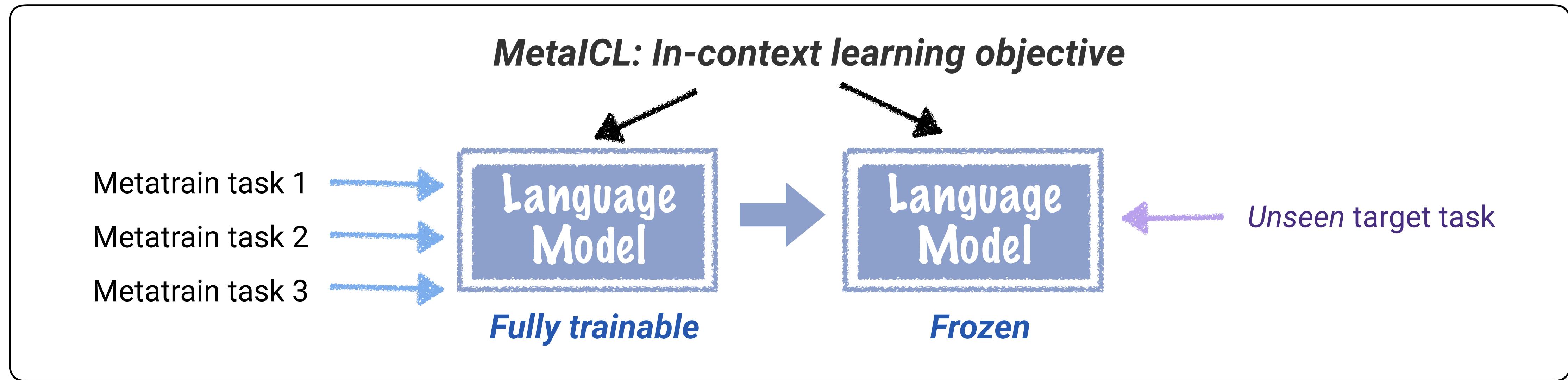
# Takeaways



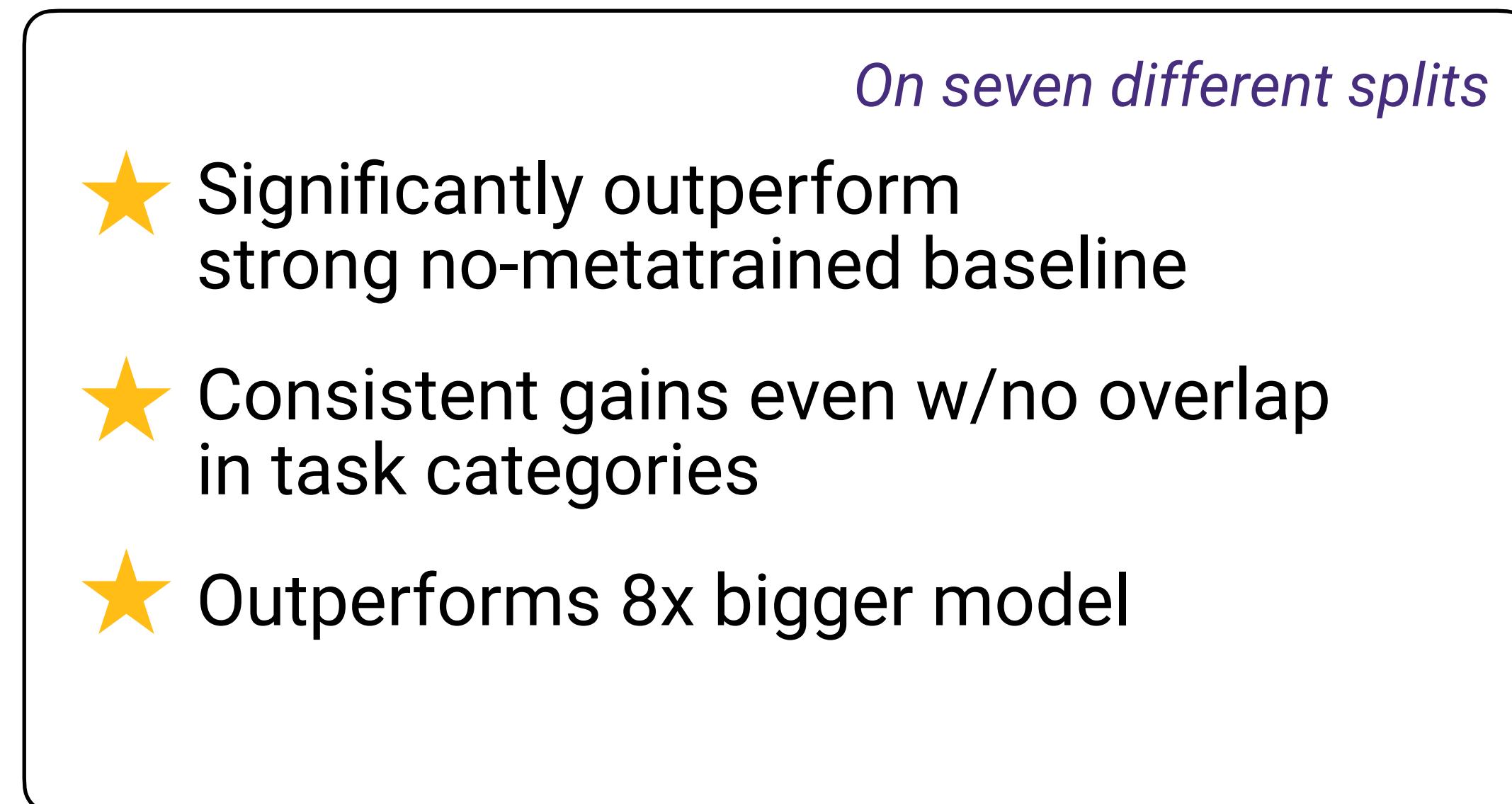
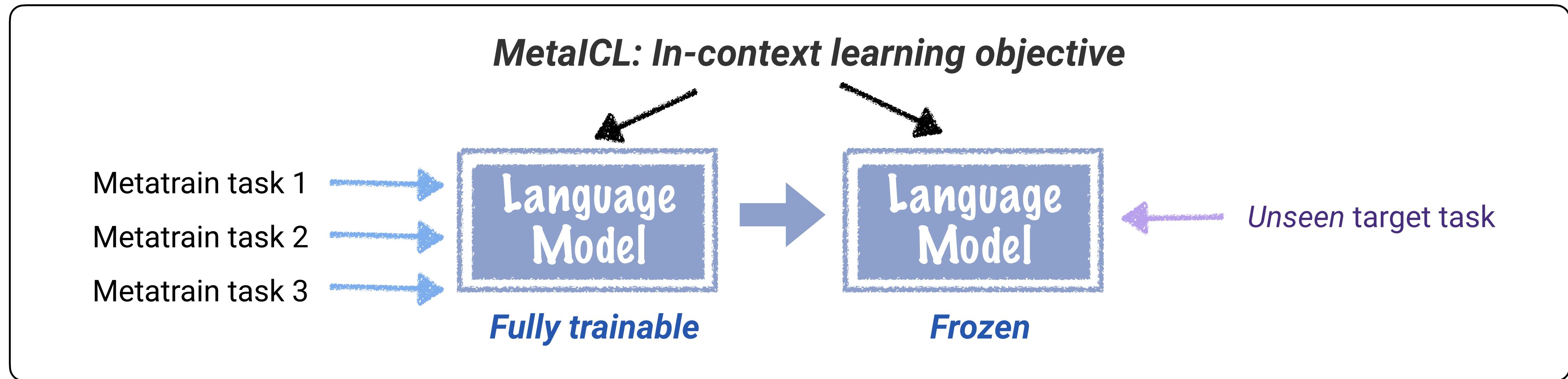
# Takeaways



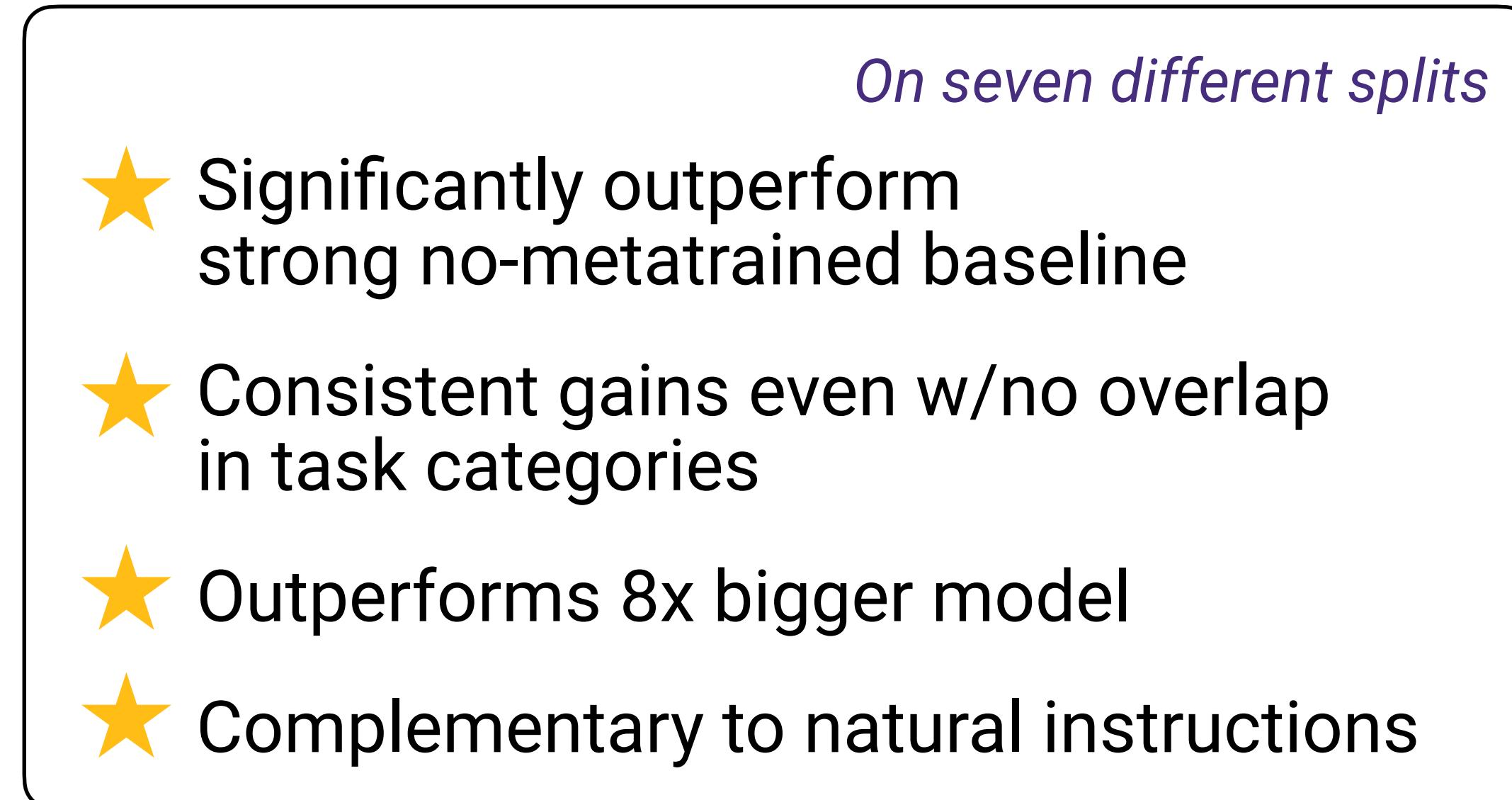
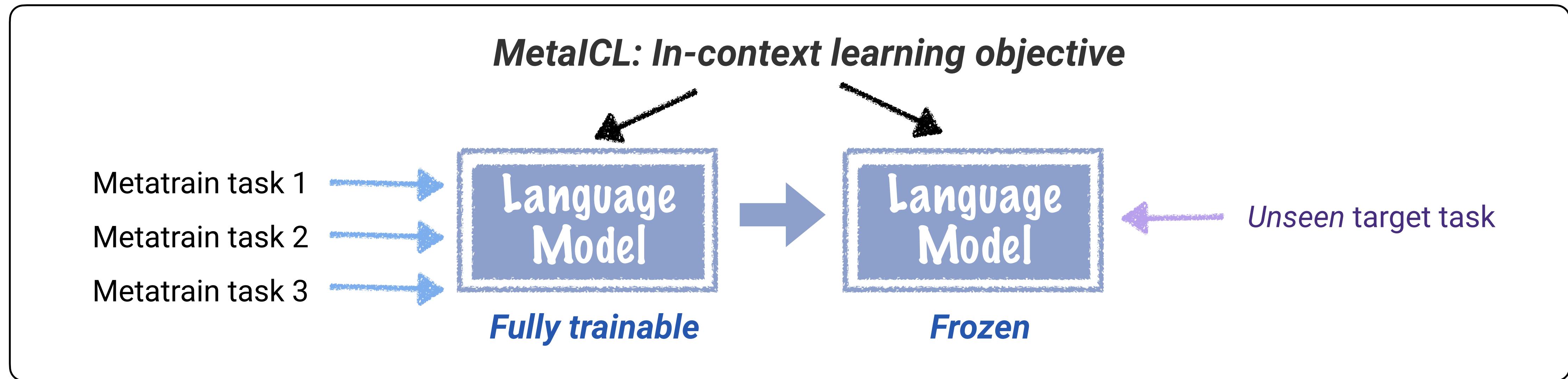
# Takeaways



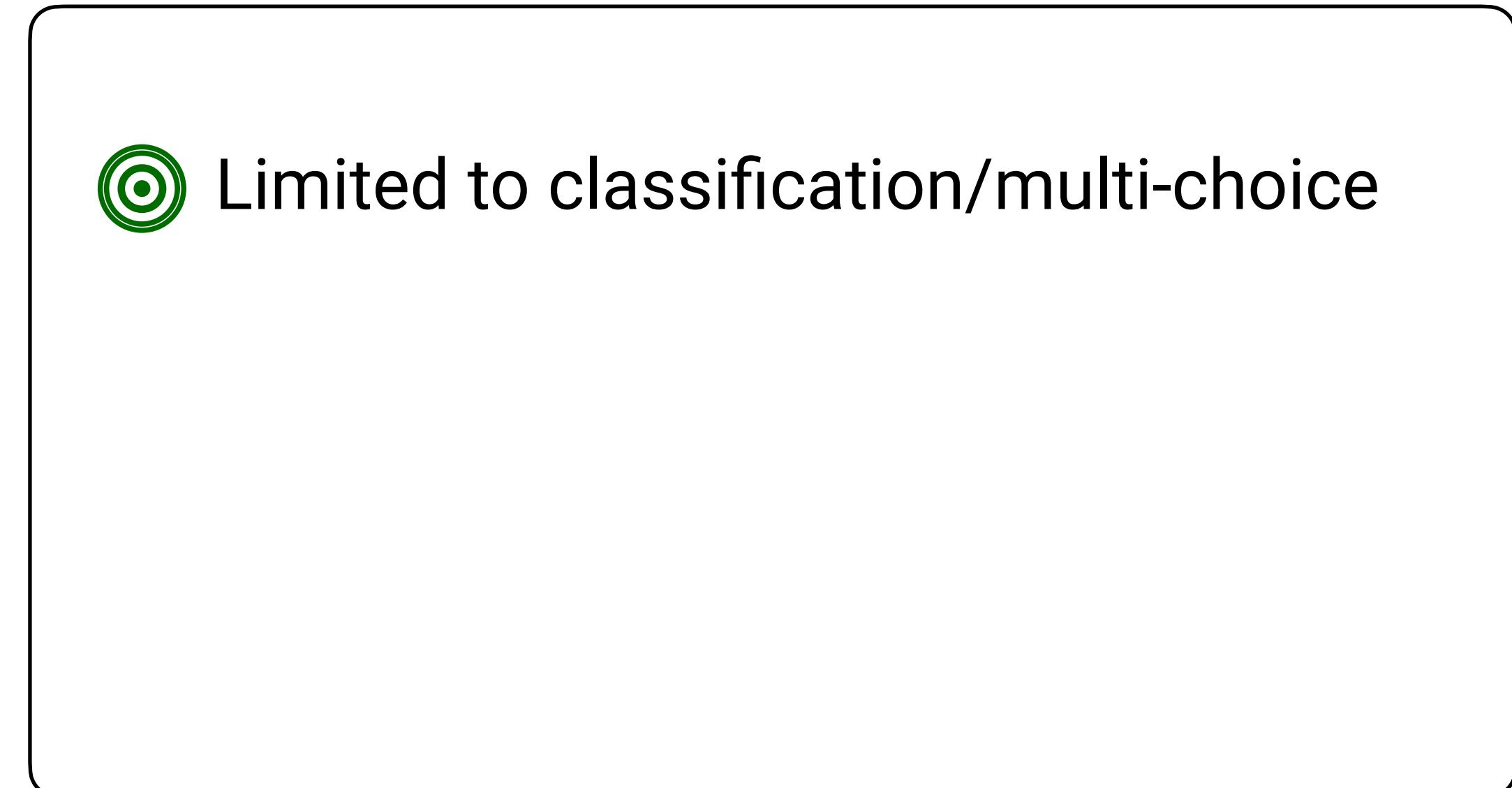
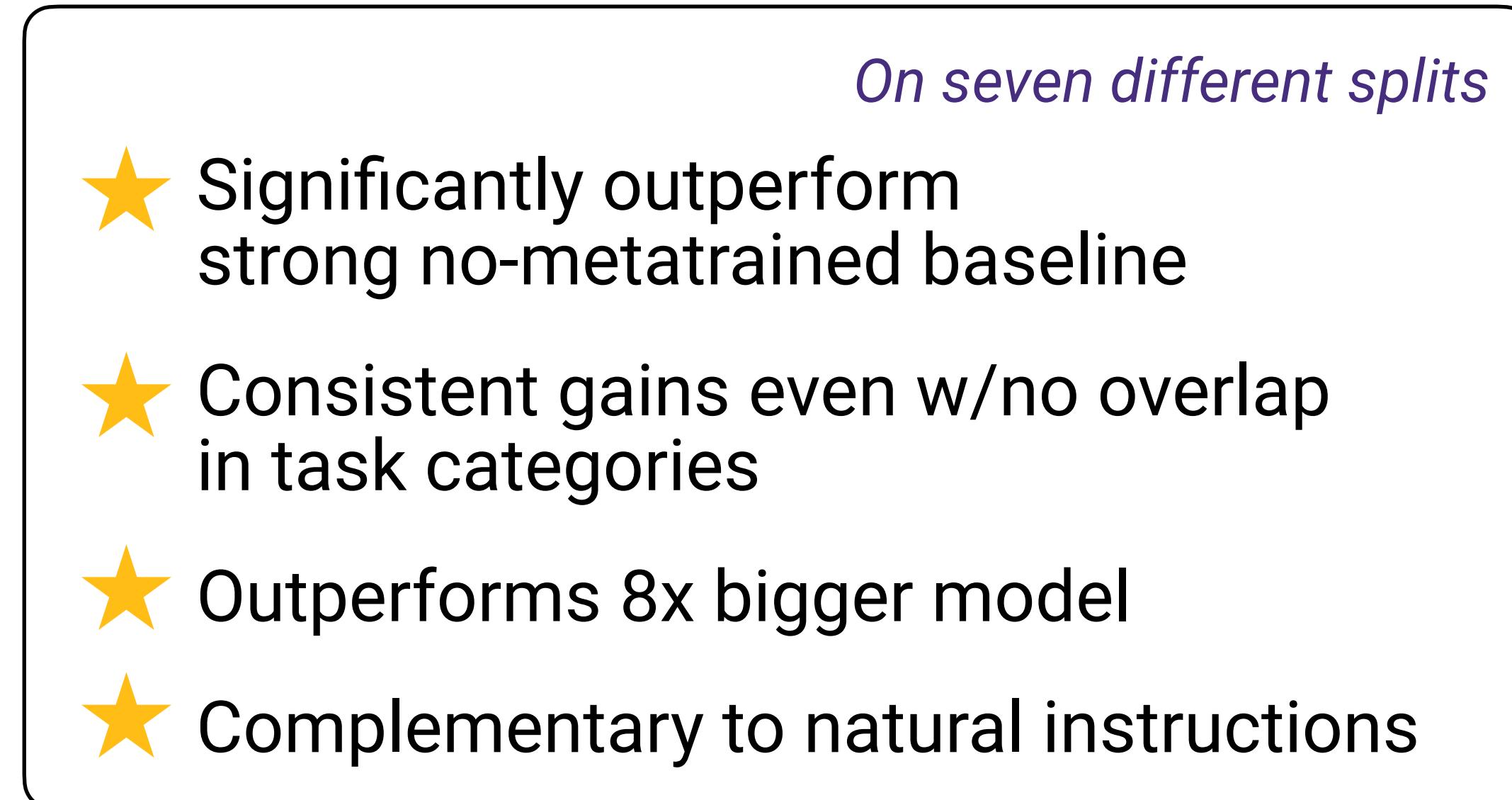
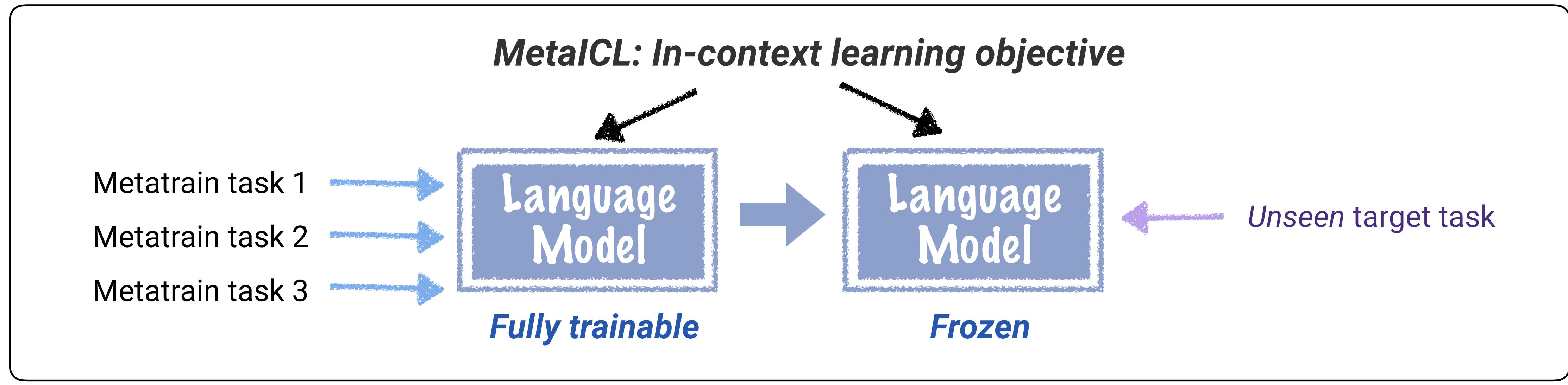
# Takeaways



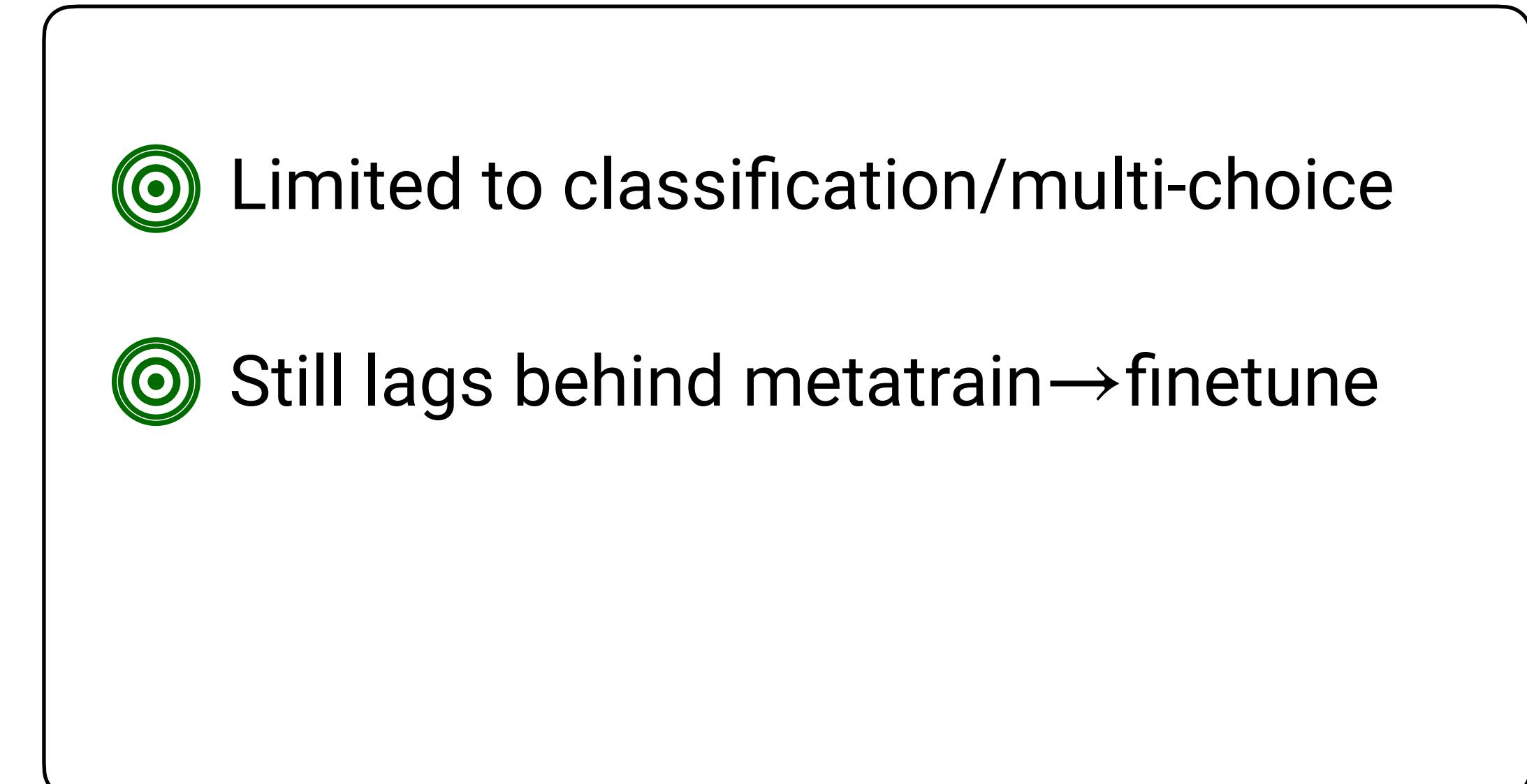
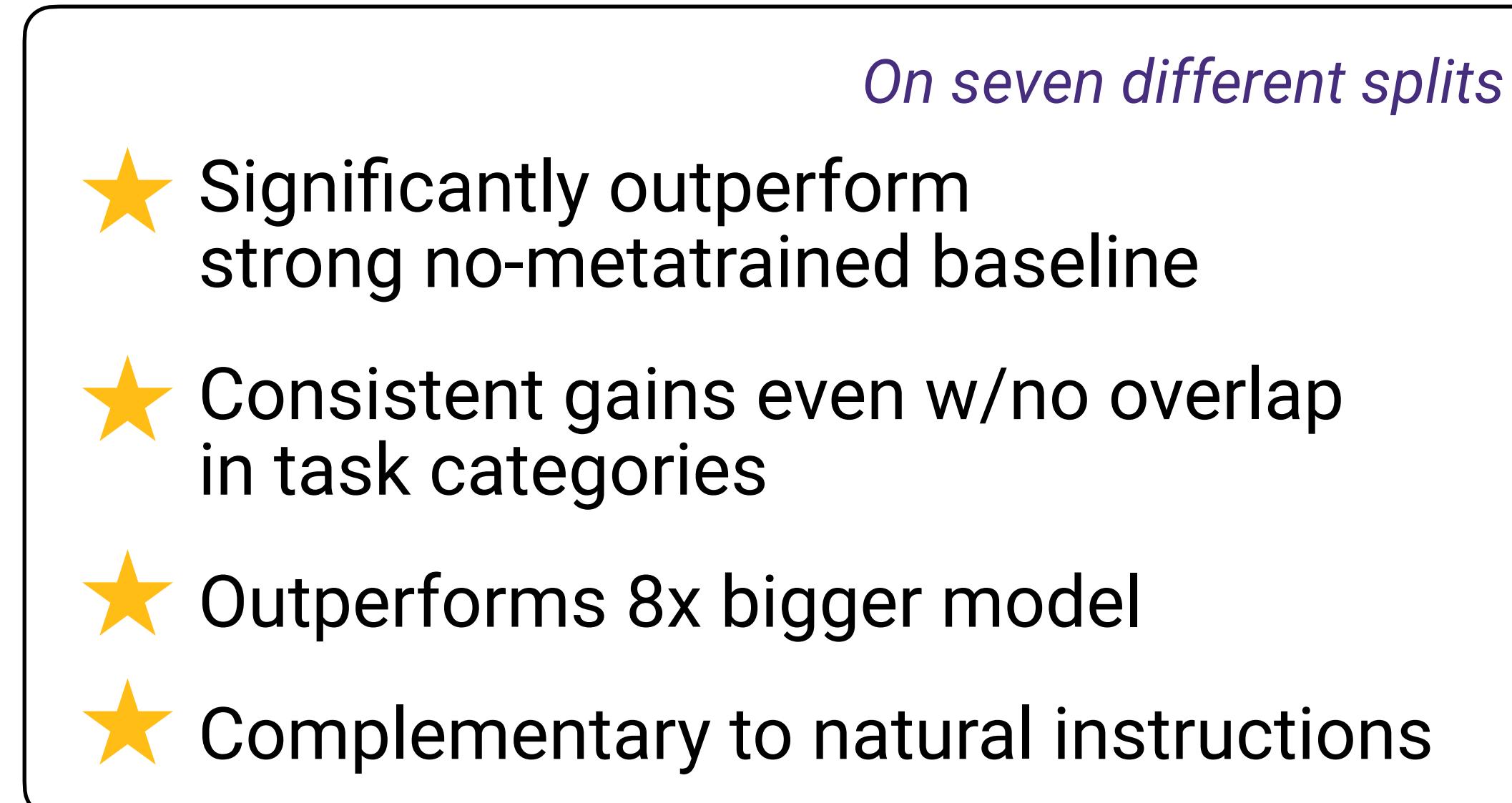
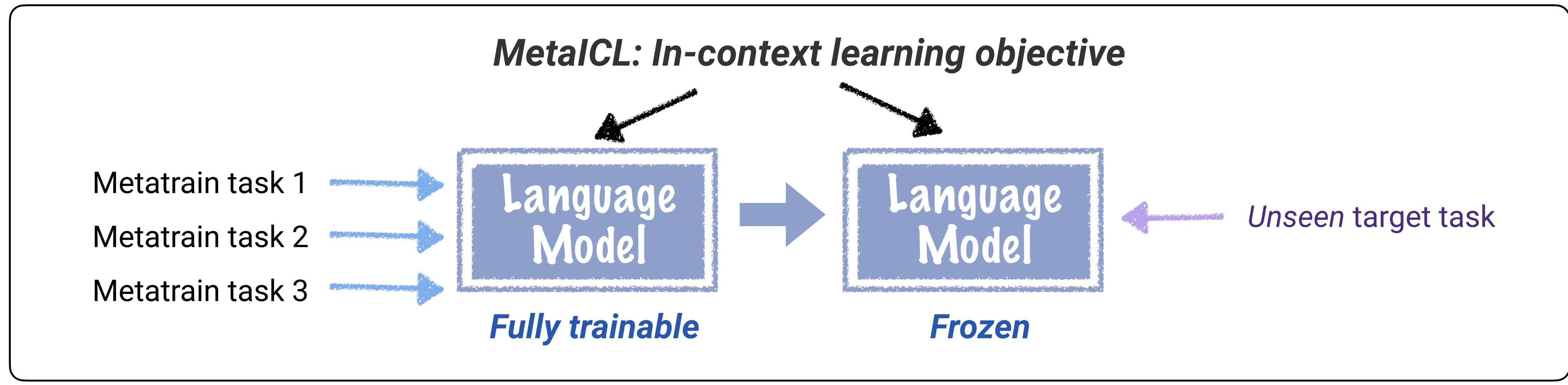
# Takeaways



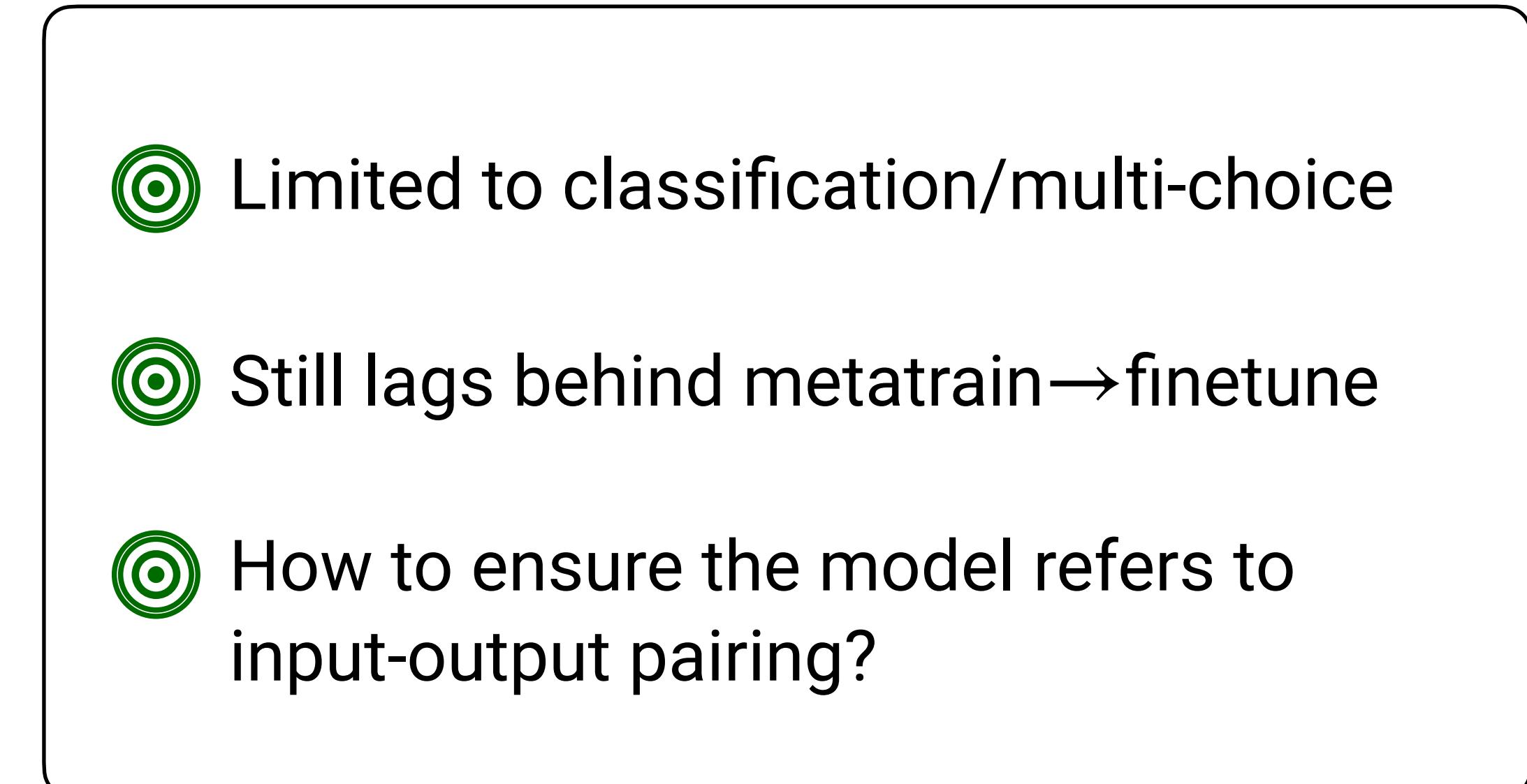
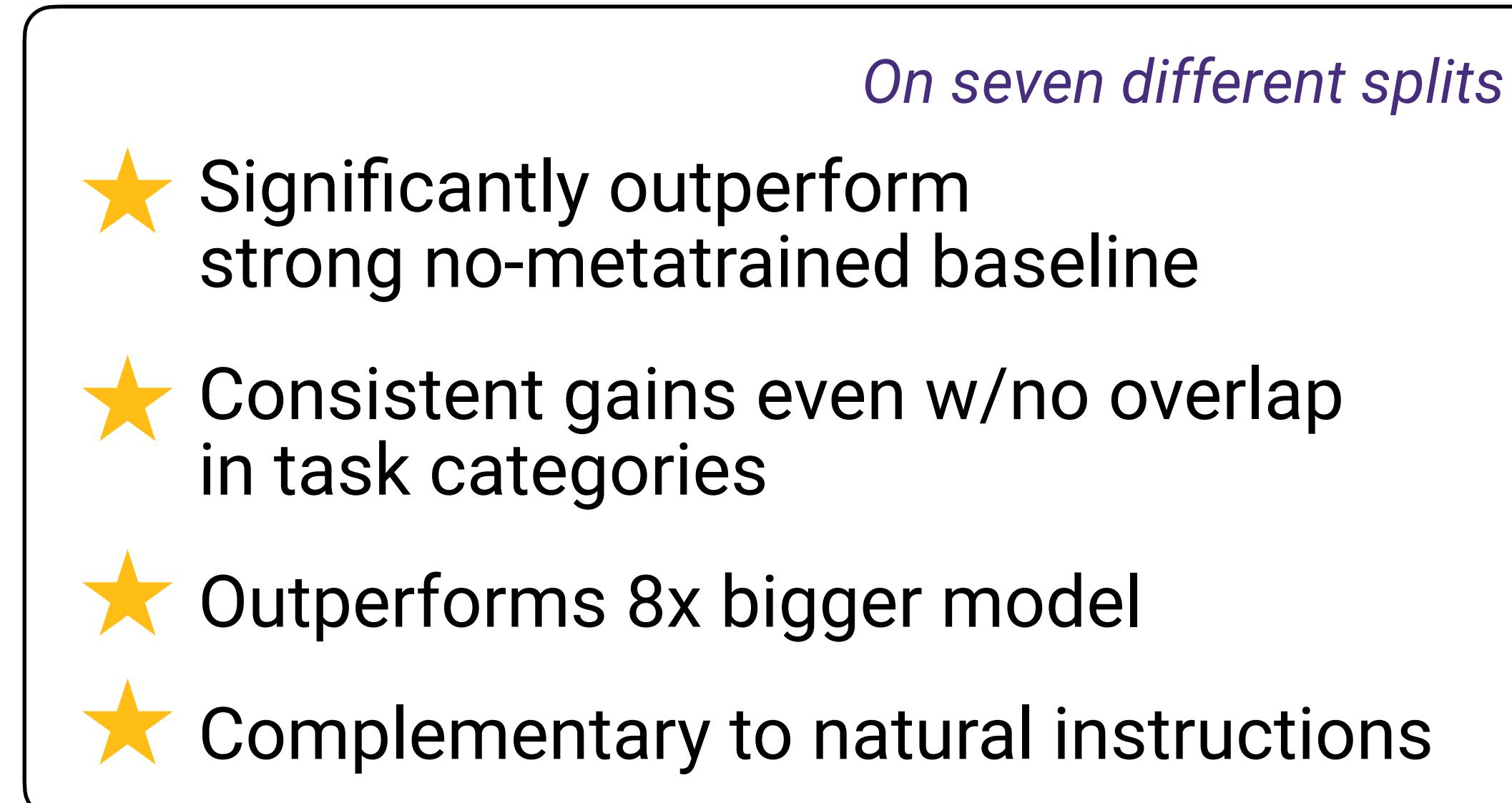
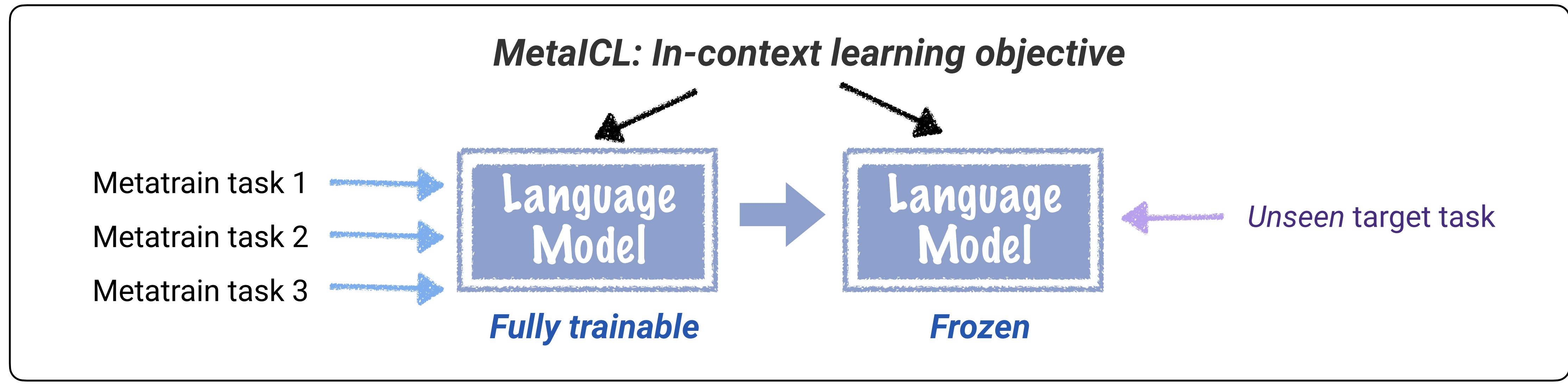
# Takeaways



# Takeaways



# Takeaways



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

## Few-shot Learning through Channel MetaICL

Sewon Min, Mike Lewis, Luke Zettlemoyer, Hannaneh Hajishirzi. 2021. "[MetaICL: Learning to Learn In Context](#)"

Sewon Min, Mike Lewis, Hannaneh Hajishirzi, Luke Zettlemoyer. 2021. "[Noisy Channel Language Model Prompting for Few-Shot Text Classification](#)"

Contact/bug report: [Sewon Min](#) ([✉](#), [🐦](#))

[\[Show Me Details!\]](#)

### Instruction:

Give the model 16 training examples (by writing yourself or choosing from an example task) and the test input, and see what model returns!

*You do not need to use any templates or instructions to tell what the task is.*

Channel LM  Channel MetaICL

### Recommended Tasks

=====Your Own Task=====

Get Random Examples!

### Training examples

Input: Our model outperforms a range of competitive baselines.

Output: positive

Input: Our model achieves the new state-of-the-art.

### Test Input

Write your own input

Run

Please enter the test input

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

## Few-shot Learning through Channel MetaICL

Sewon Min, Mike Lewis, Luke Zettlemoyer, Hannaneh Hajishirzi. 2021. "[MetaICL: Learning to Learn In Context](#)"

Sewon Min, Mike Lewis, Hannaneh Hajishirzi, Luke Zettlemoyer. 2021. "[Noisy Channel Language Model Prompting for Few-Shot Text Classification](#)"

Contact/bug report: [Sewon Min](#) ([✉](#), [🐦](#))

[\[Show Me Details!\]](#)

### Instruction:

Give the model 16 training examples (by writing yourself or choosing from an example task) and the test input, and see what model returns!

*You do not need to use any templates or instructions to tell what the task is.*

Channel LM  Channel MetaICL

### Recommended Tasks

=====Your Own Task=====

Get Random Examples!

### Training examples

Input: Our model outperforms a range of competitive baselines.

Output: positive

Input: Our model achieves the new state-of-the-art.

### Test Input

Write your own input

Run

Please enter the test input

# Thank you for listening

**Paper:** [arxiv.org/abs/2110.15943](https://arxiv.org/abs/2110.15943)

**Code:** [github.com/facebookresearch/MetaCL](https://github.com/facebookresearch/MetaCL)

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

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