

Meta-Learning to Compositionally Generalize

Henry Conklin^{1*}, Bailin Wang^{1*}, Kenny Smith¹,
Ivan Titov^{1,2}

¹University of Edinburgh ²University of Amsterdam *equal contribution



**the meaning of a sentence is constructed from
the meaning of its parts and the way in which
they are combined**

(Cann, 1993)

Compositionality

Generalization

**Enables robust generalization outside of
prior experience**

Compositionality

Generalization

**Enables robust generalization outside of
prior experience**

The deer ran across the road last night.

Compositionality

Generalization

**Enables robust generalization outside of
prior experience**

The deer ran across the road last night.

The chicken walked into town.

Compositionality

Generalization

**Enables robust generalization outside of
prior experience**

The deer ran across the road last night.

The deer walked into town.

The chicken walked into town.

Compositionality

Generalization

**Enables robust generalization outside of
prior experience**

The deer ran across the road last night.

The **deer** walked into town.

The chicken walked into town.

I don't like pears, I find them sinister.

Compositionality

Generalization

**Enables robust generalization outside of
prior experience**

The deer ran across the road last night.

The deer walked into town.

The chicken walked into town.

I don't like pears, I find them sinister.

apples

Compositionality

Generalization

**Enables robust generalization outside of
prior experience**

The deer ran across the road last night.

The `deer` walked into town.

The chicken walked into town.

I don't like pears, I find them sinister.

I don't like `apples`, I find them sinister.

apples

Compositionality

Current Limitations

State of the art neural models struggle to generalize outside of their training distribution

Compositionality

Current Limitations

State of the art neural models struggle to generalize outside of their training distribution

- They struggle to use new words in a compositional context (Lake and Baroni, 2018)

Compositionality

Current Limitations

State of the art neural models struggle to generalize outside of their training distribution

- They struggle to use new words in a compositional context (Lake and Baroni, 2018)
- Difficulty interpreting known words in new contexts (Keysers et al., 2020)

Compositionality

Current Limitations

State of the art neural models struggle to generalize outside of their training distribution

- They struggle to use new words in a compositional context (Lake and Baroni, 2018)
- Difficulty interpreting known words in new contexts (Keysers et al., 2020)
- Issues generalizing known words to new syntactic structures (Kim and Linzen, 2020)

Compositionality

Current Limitations

State of the art neural models struggle to generalize outside of their training distribution

Compositionality

Current Limitations

State of the art neural models struggle to generalize outside of their training distribution

- Models may prefer memorization over generalization (Liška et al., 2018)

Compositionality

Current Limitations

State of the art neural models struggle to generalize outside of their training distribution

- Models may prefer memorization over generalization (Liška et al., 2018)
- Models may memorize sections of their input (Hupkes et al., 2019)

Compositionality

Current Limitations

State of the art neural models struggle to generalize outside of their training distribution

- Models may prefer memorization over generalization (Liška et al., 2018)
- Models may memorize sections of their input (Hupkes et al., 2019)
- Limited memory may be key to why humans arrive at robust solutions (Griffiths, 2020)

Compositionality

Current Limitations

State of the art neural models struggle to generalize outside of their training distribution

- Models may prefer memorization over generalization (Liška et al., 2018)
- Models may memorize sections of their input (Hupkes et al., 2019)
- Limited memory may be key to why humans arrive at robust solutions (Griffiths, 2020)

Can we inhibit these models' ability to memorize?

Explicitly testing compositionality

COGS (Kim and Linzen, 2020)

Explicitly testing compositionality

COGS (Kim and Linzen, 2020)

The sailor dusted a boy . → * sailor (x _ 1) ; dust . agent (x _ 2 , x _ 1) AND
dust . theme (x _ 2 , x _ 4) AND boy (x _ 4)

Explicitly testing compositionality

COGS (Kim and Linzen, 2020)

The sailor dusted a boy . \longrightarrow * sailor (x _ 1) ; dust . agent (x _ 2 , x _ 1) AND
dust . theme (x _ 2 , x _ 4) AND boy (x _ 4)

Novel Combination of Familiar Primitives and Grammatical Roles

Subject \longrightarrow Object

A **hedgehog** ate the cake.

The baby liked the **hedgehog**.

Deeper Recursion

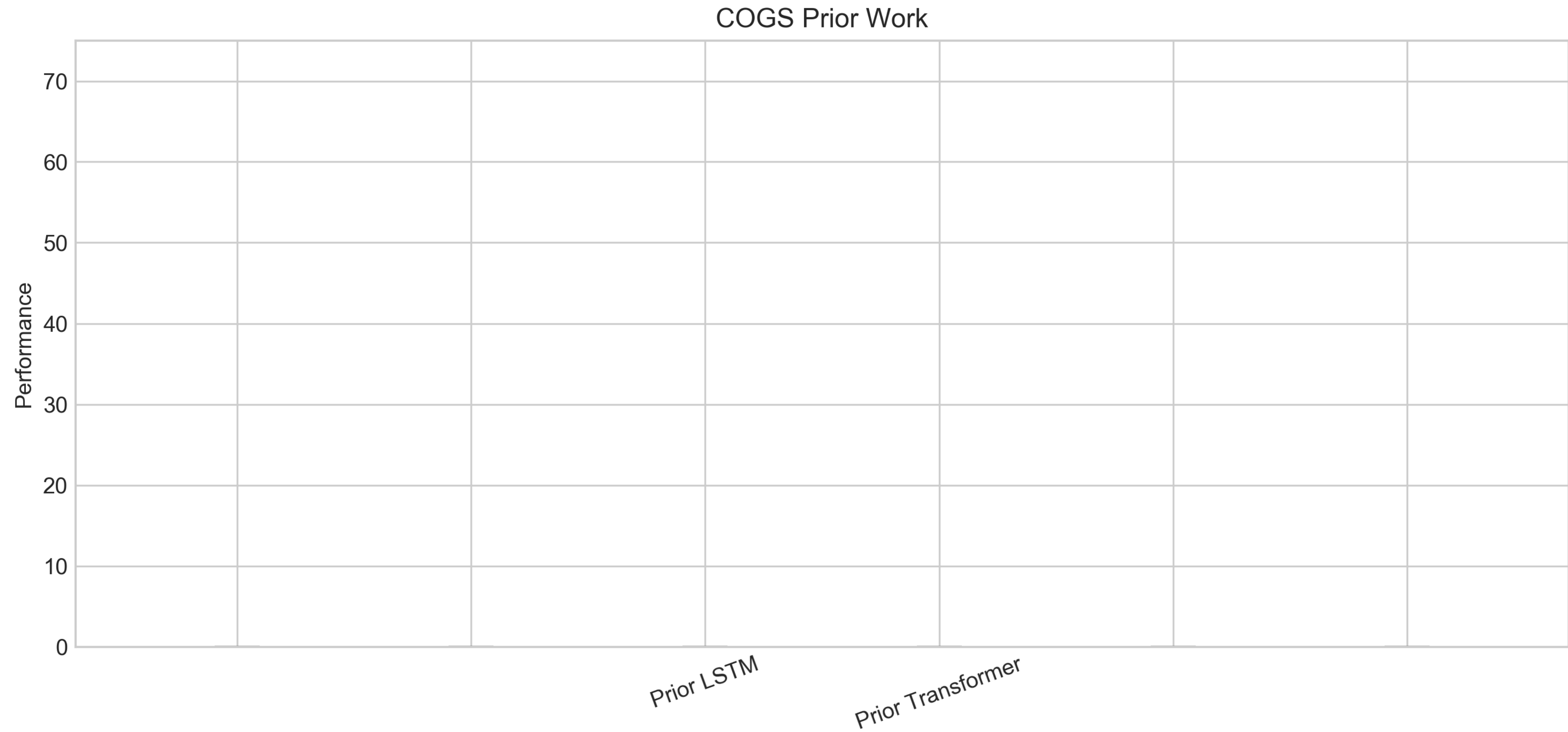
Sentential complements

Emma said **that** Noah knew **that** the cat danced.

Emma said **that** Noah knew **that** Lucas saw **that** the cat danced.

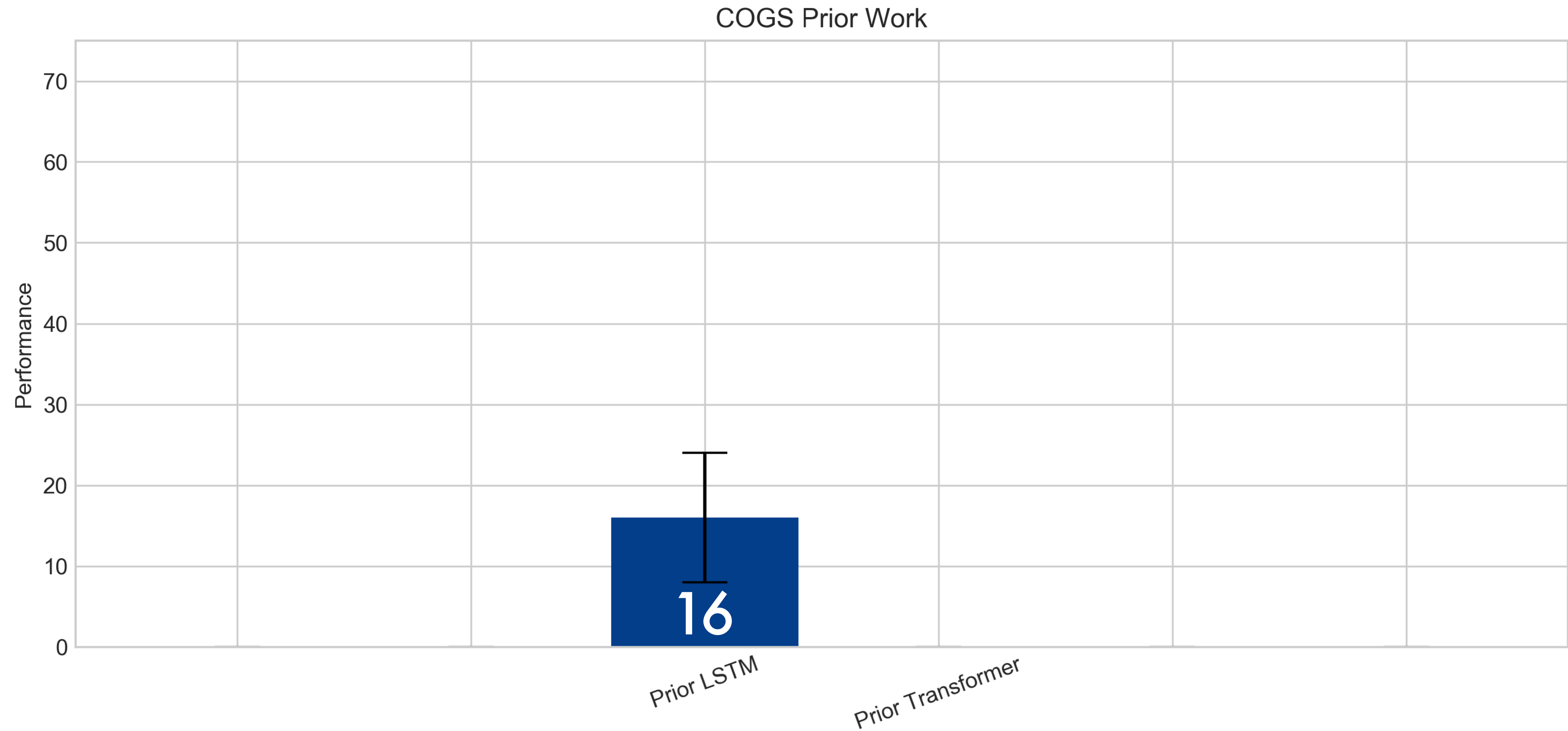
Explicitly testing compositionality

COGS (Kim and Linzen, 2020)



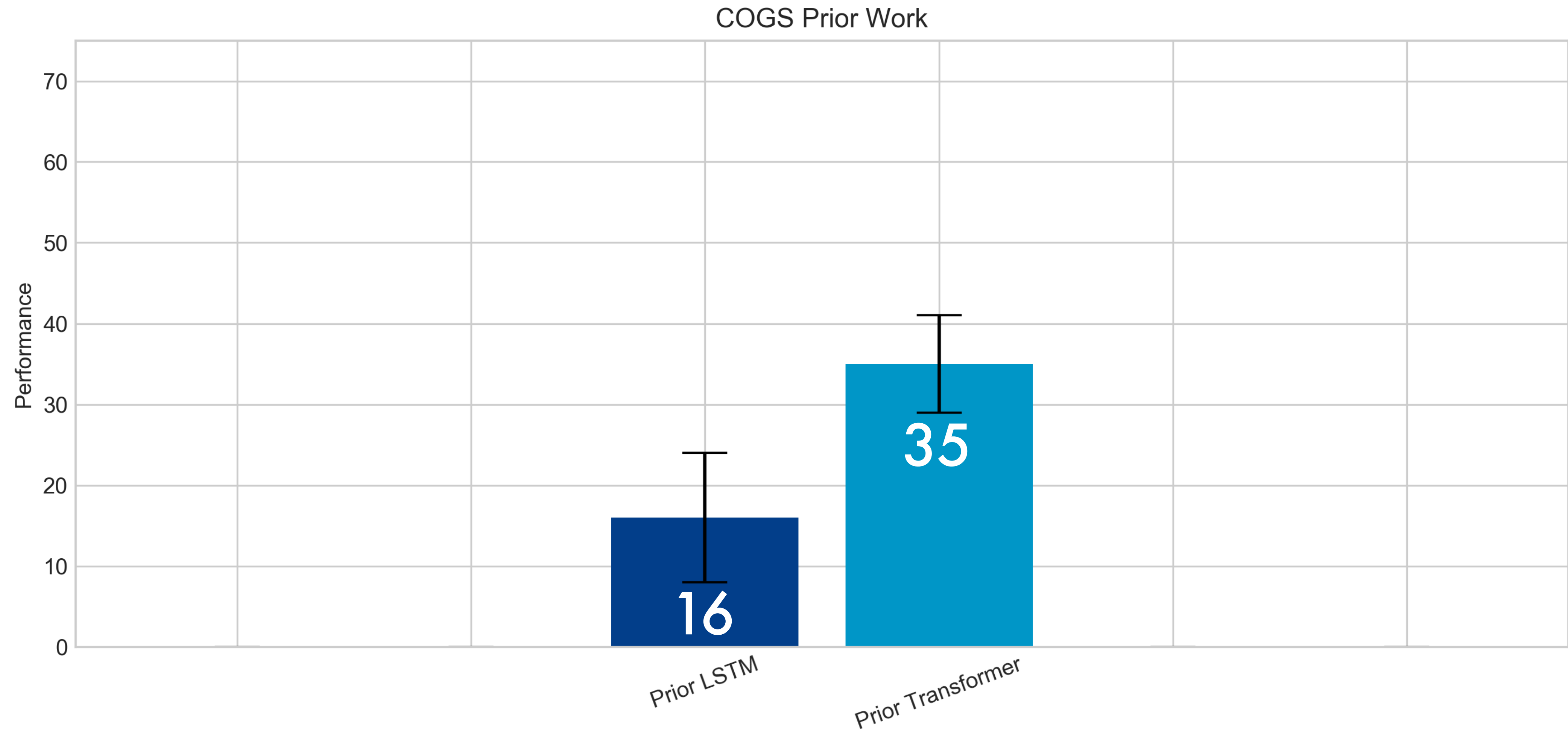
Explicitly testing compositionality

COGS (Kim and Linzen, 2020)



Explicitly testing compositionality

COGS (Kim and Linzen, 2020)



**is it surprising that these
tasks are hard?**

is it surprising that this is hard?

Underspecification & Supervised learning

is it surprising that this is hard?

Underspecification & Supervised learning

- Data under-specifies for the generalizations that produced it (Goodman, 1955)

is it surprising that this is hard?

Underspecification & Supervised learning

- Data under-specifies for the generalizations that produced it (Goodman, 1955)
- Models are trained on these tasks using supervised learning.

is it surprising that this is hard?

Underspecification & Supervised learning

- Data under-specifies for the generalizations that produced it (Goodman, 1955)
- Models are trained on these tasks using supervised learning.
- Independent and Identically Distributed Assumption

is it surprising that this is hard?

Underspecification & Supervised learning

- Data under-specifies for the generalizations that produced it (Goodman, 1955)
- Models are trained on these tasks using supervised learning.
- Independent and Identically Distributed Assumption
 - explicitly not met here.

is it surprising that this is hard?

Underspecification & Supervised learning

- Data under-specifies for the generalizations that produced it (Goodman, 1955)
- Models are trained on these tasks using supervised learning.
- Independent and Identically Distributed Assumption
 - explicitly not met here.
- **Underspecified data + Supervised learning may fail to consistently extract robust strategies**

Improving Generalization

DG-MAML (Li et al., 2018)

Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that

Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that



dataset

Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that

no-embedding

dataset

lots of embedding

Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that

short

dataset

long

Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that

short

dataset

long

I like the cat.

Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that

short

dataset

long

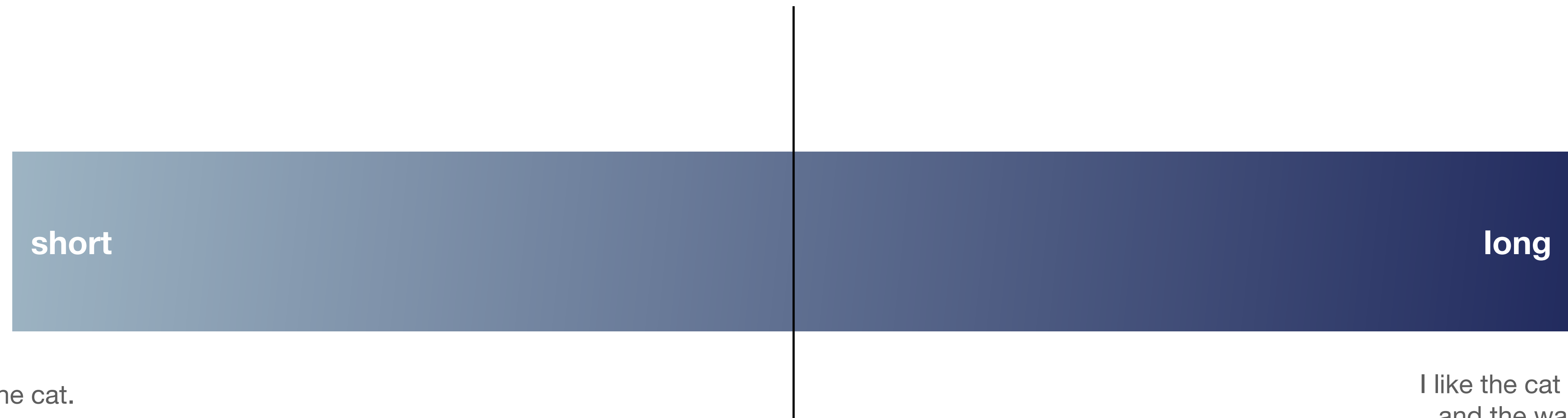
I like the cat.

I like the cat and the dog
and the way that they
seem to be friends.

Improving Generalization

DG-MAML (Li et al., 2018)

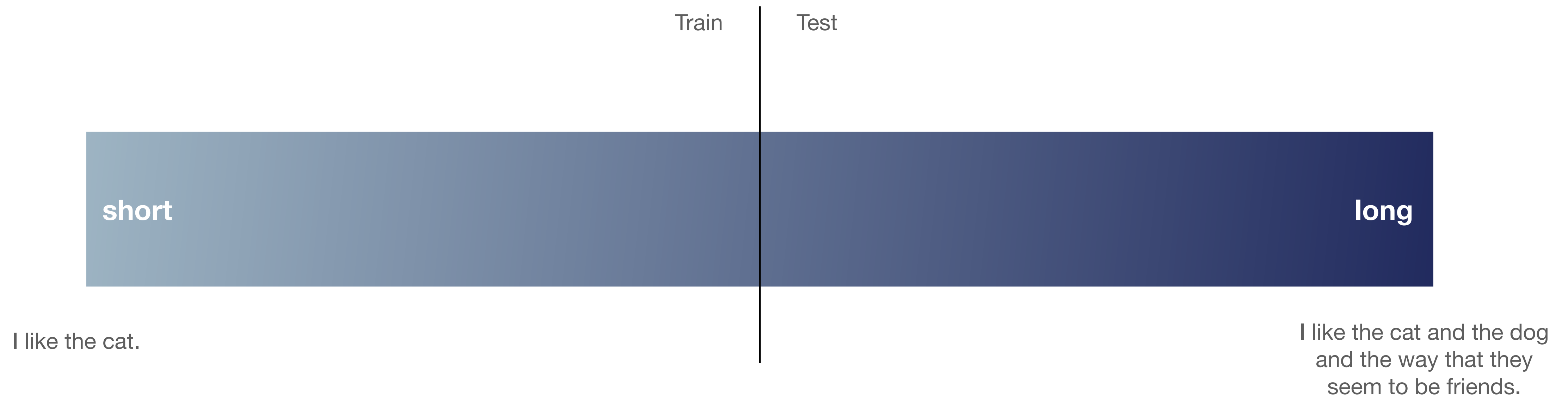
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

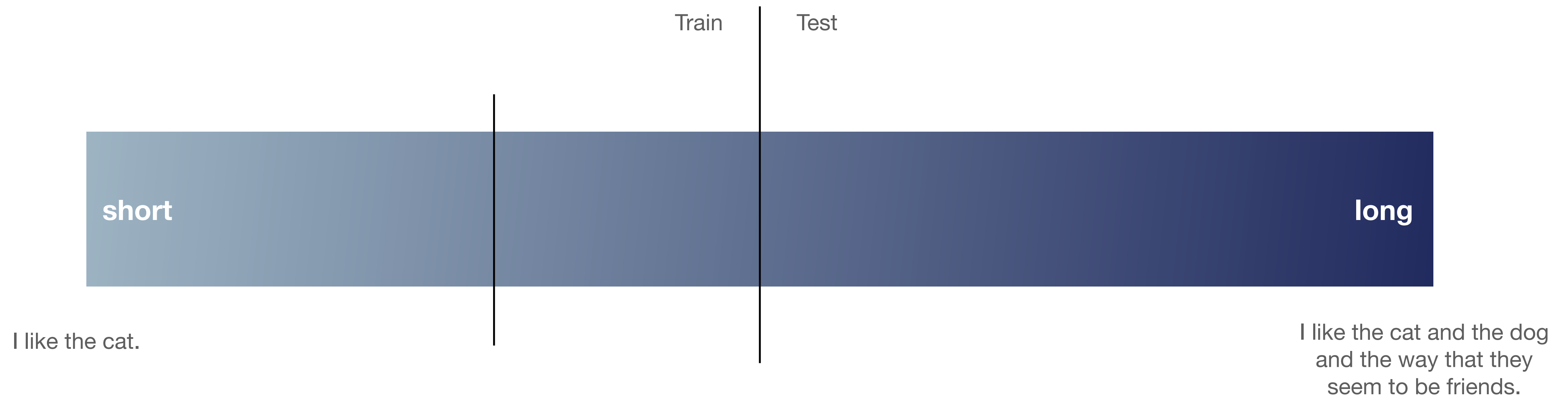
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

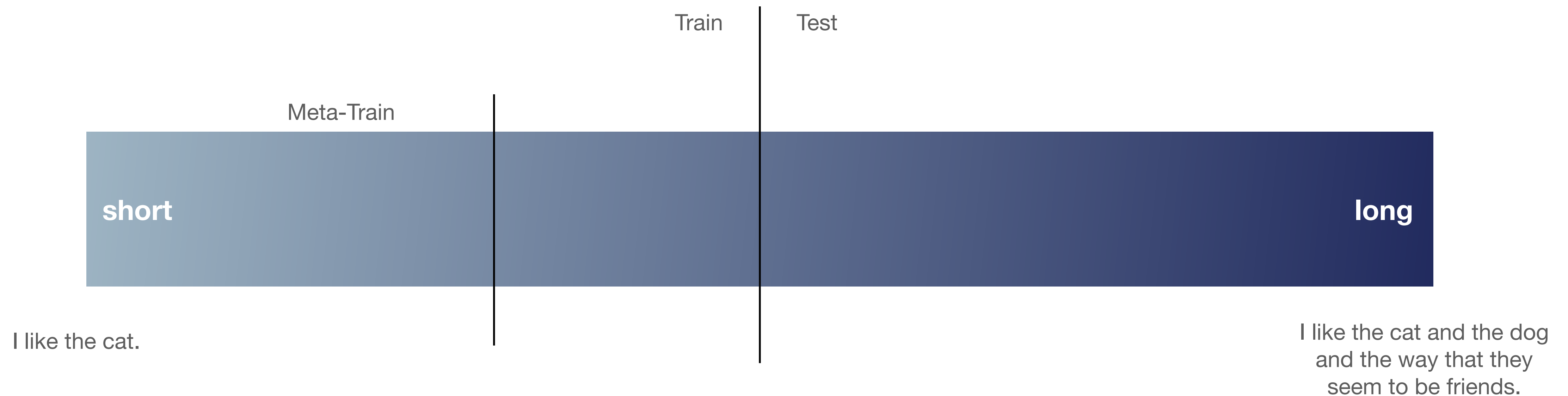
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

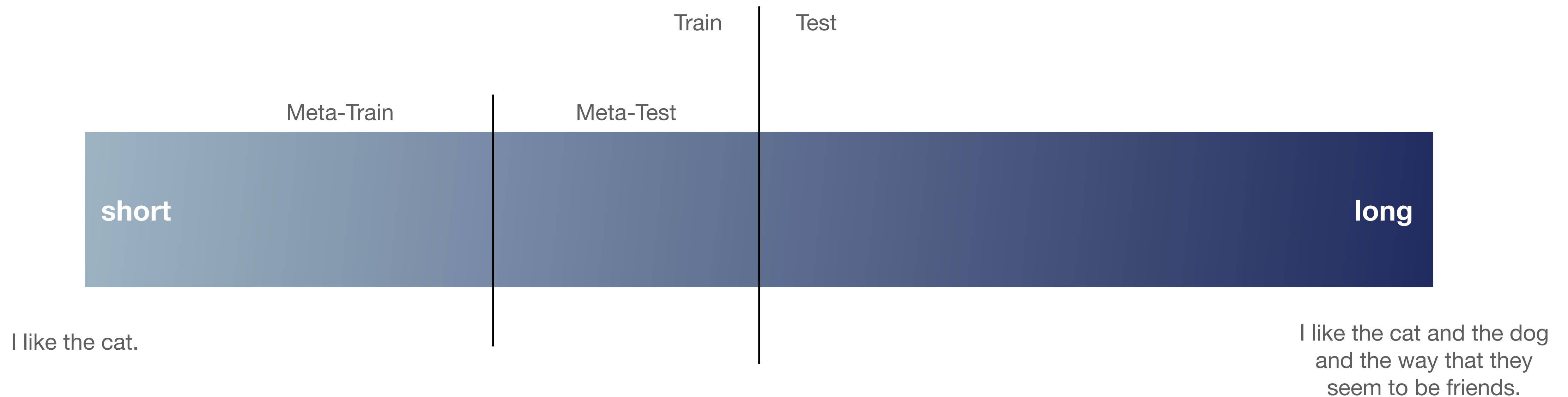
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that

Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that

\mathcal{T}

Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that

$\mathcal{T} <$

Improving Generalization

DG-MAML (Li et al., 2018)

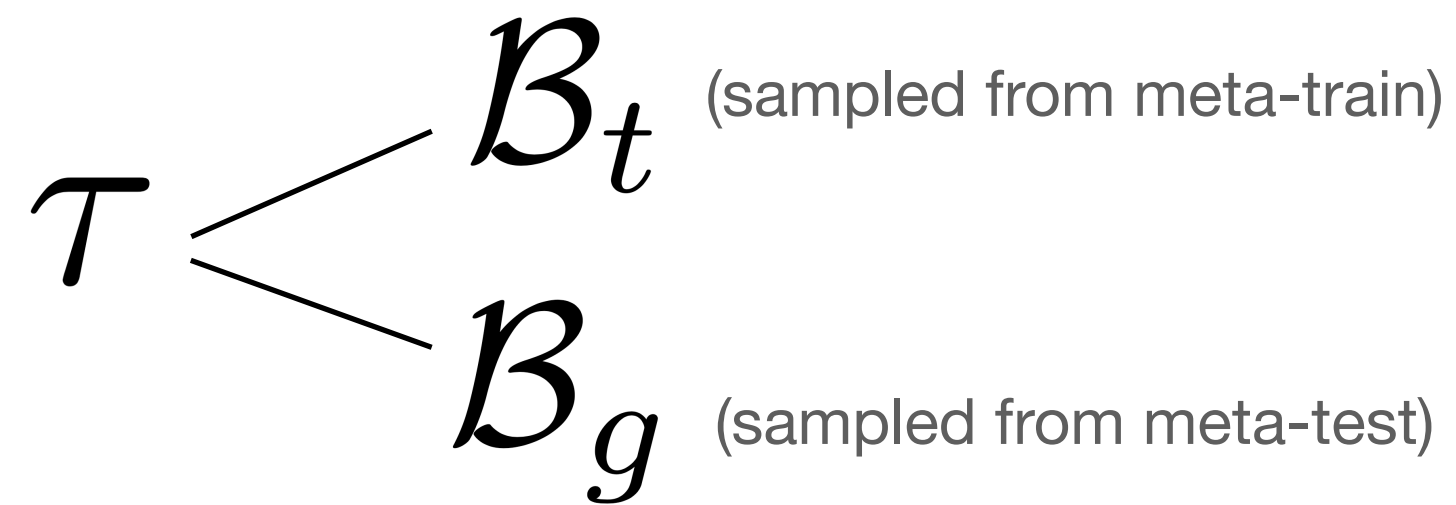
- if we know we're going to be tested on something different, let's train for that

$$\mathcal{T} \leftarrow \mathcal{B}_t \text{ (sampled from meta-train)}$$

Improving Generalization

DG-MAML (Li et al., 2018)

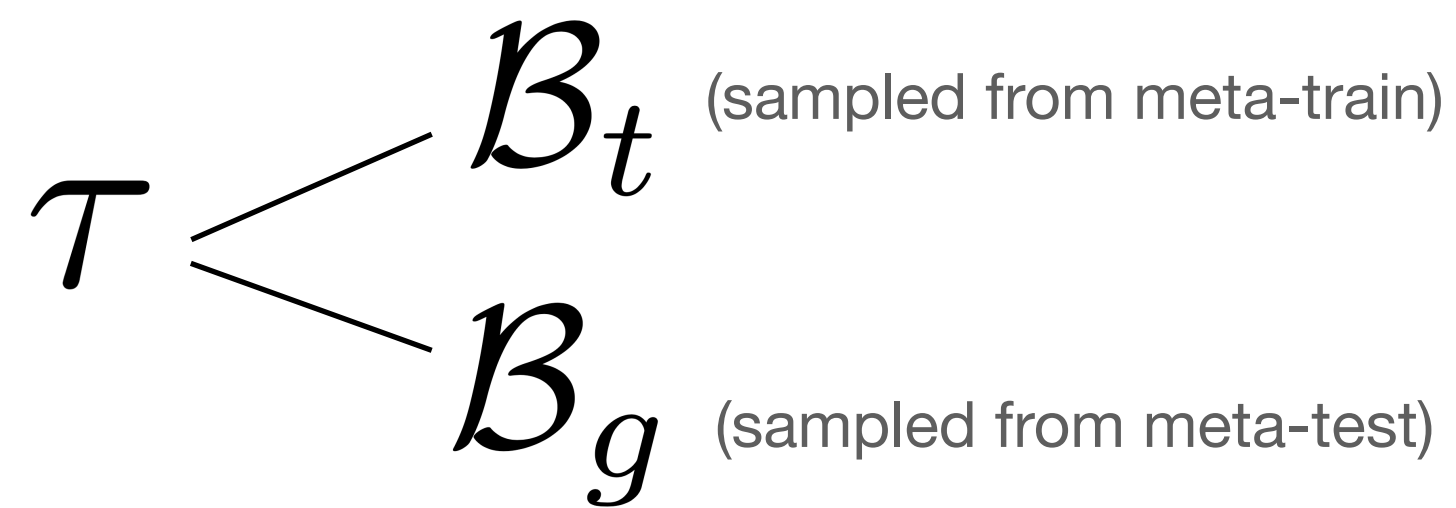
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that

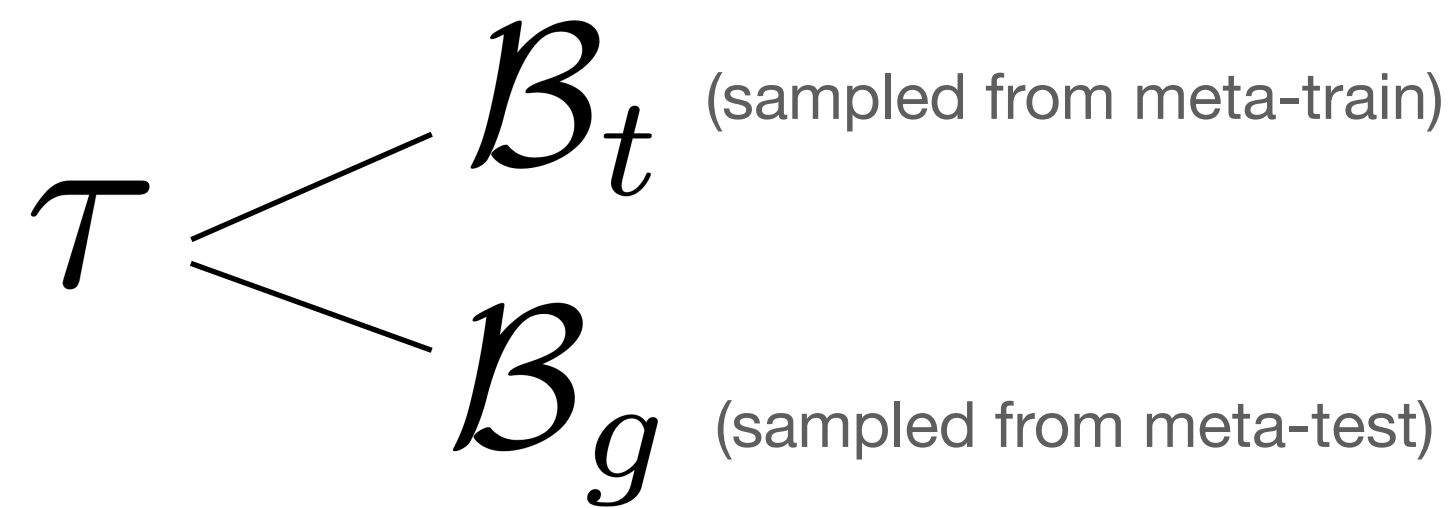


$$\mathcal{L}_{\mathcal{B}_t}(\theta)$$

Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that

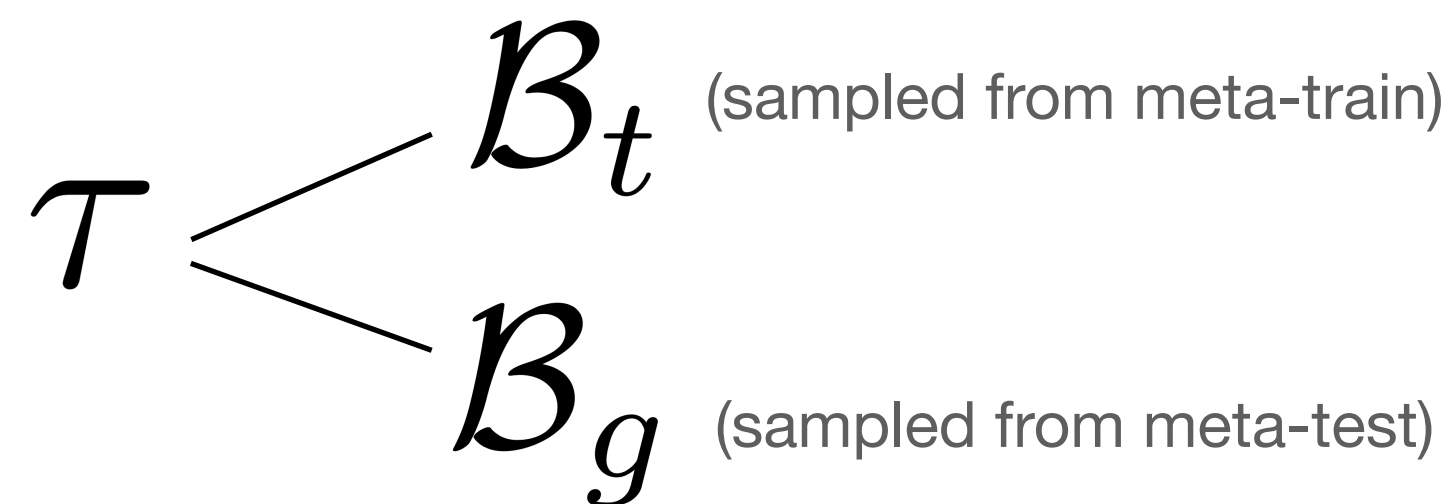


$$\mathcal{L}_{\mathcal{B}_t}(\theta)$$
$$\theta' \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{B}_t}(\theta)$$

Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that

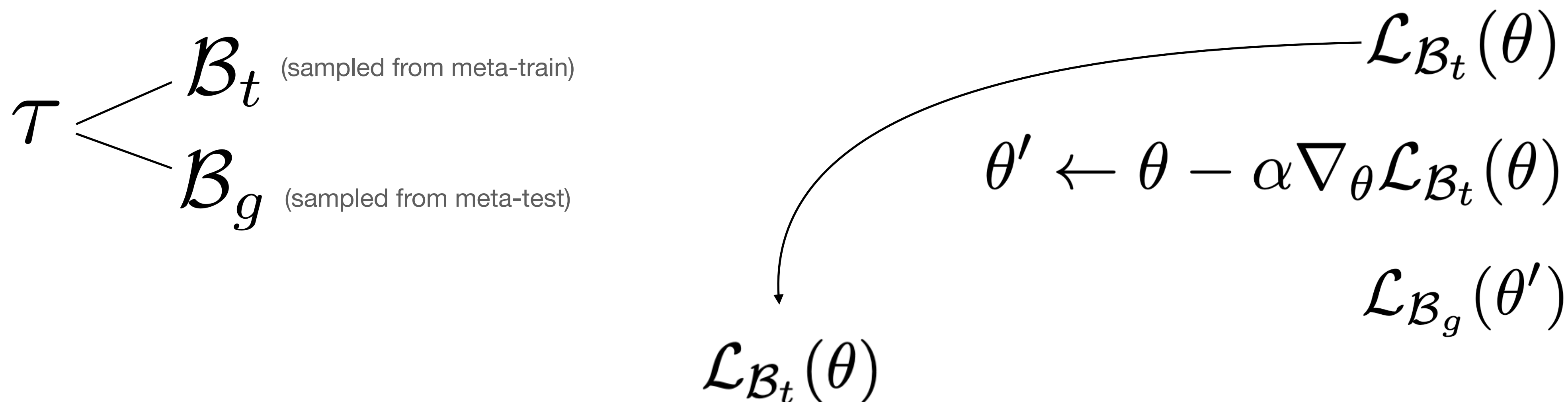


$$\begin{aligned} & \mathcal{L}_{\mathcal{B}_t}(\theta) \\ \theta' & \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}_{\mathcal{B}_t}(\theta) \\ & \mathcal{L}_{\mathcal{B}_g}(\theta') \end{aligned}$$

Improving Generalization

DG-MAML (Li et al., 2018)

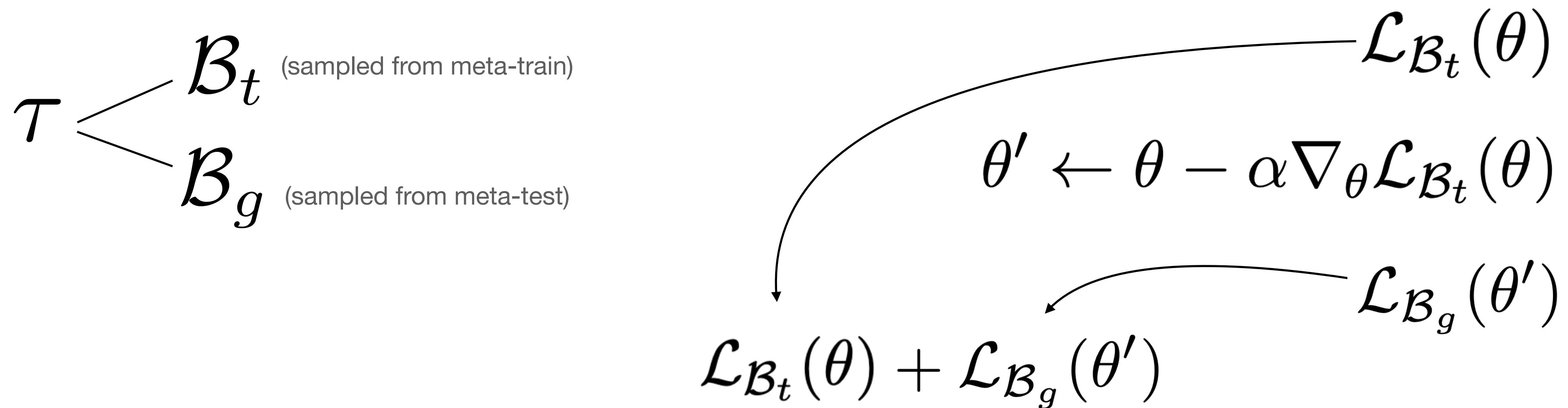
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

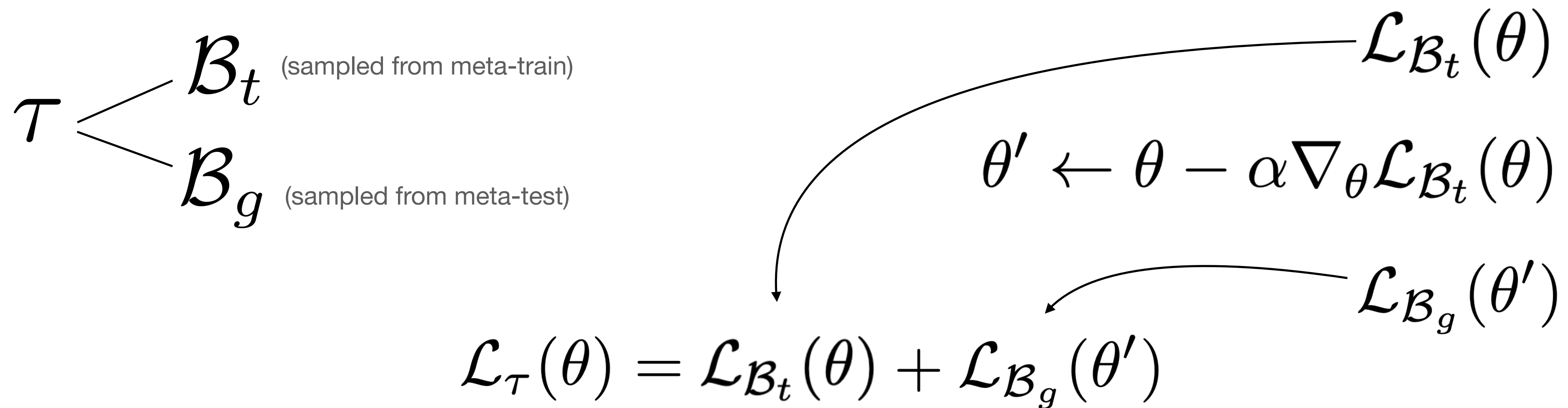
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

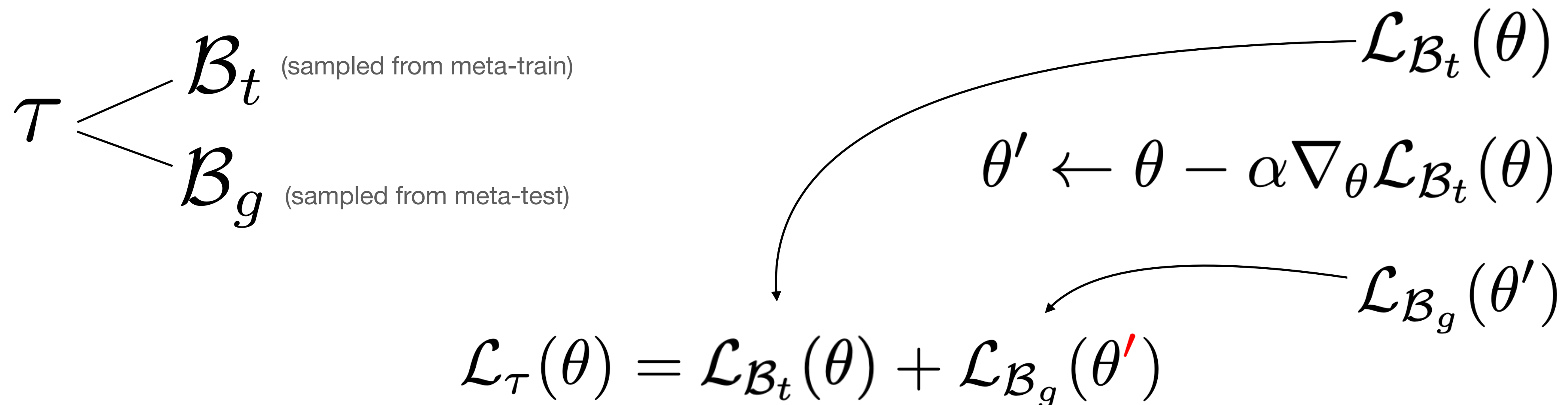
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

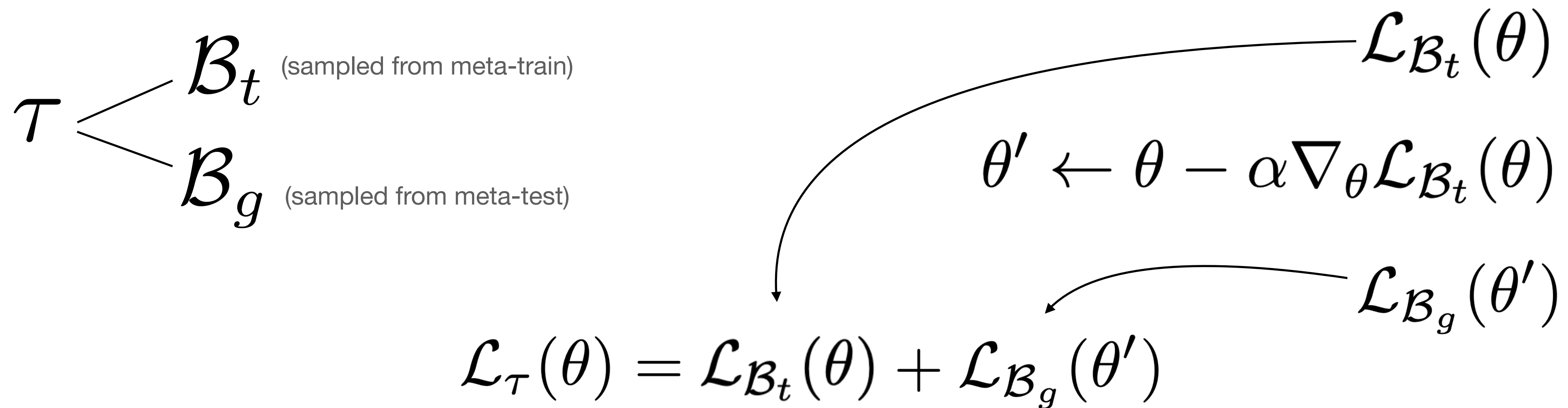
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

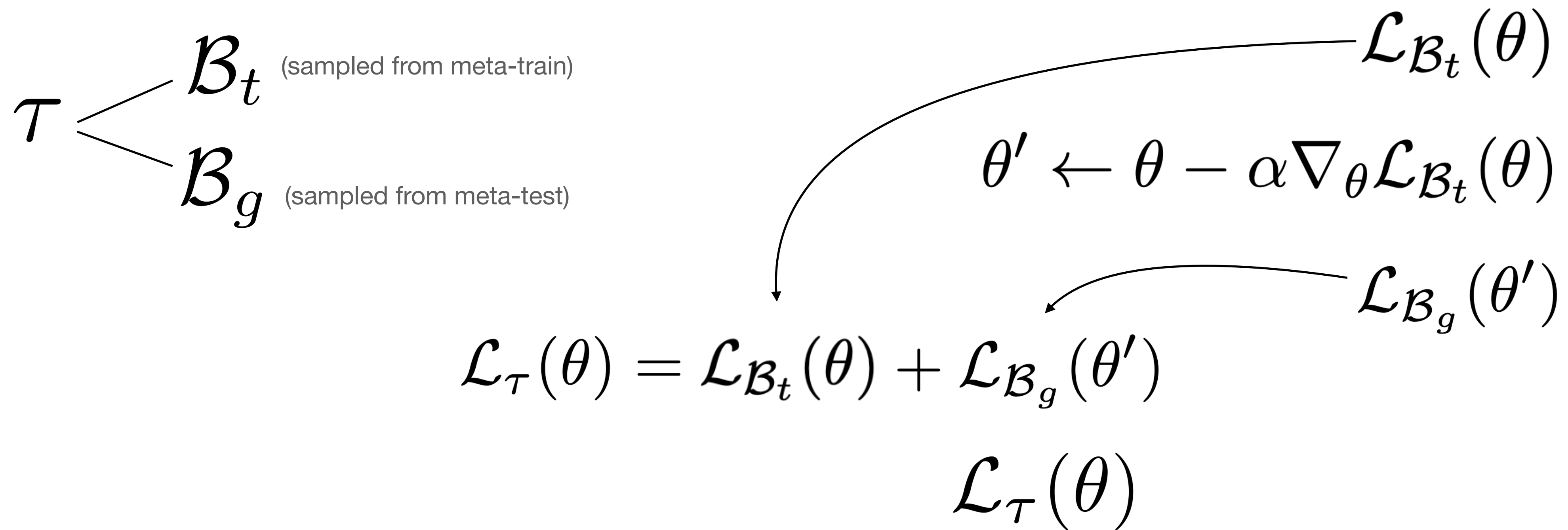
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

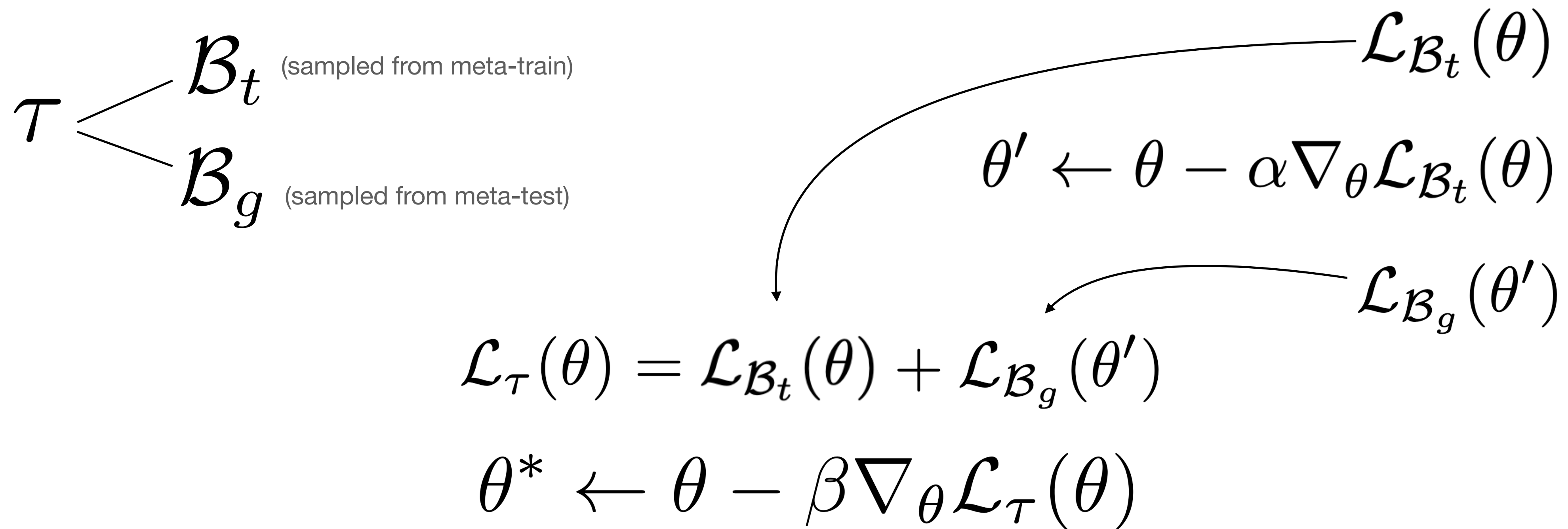
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

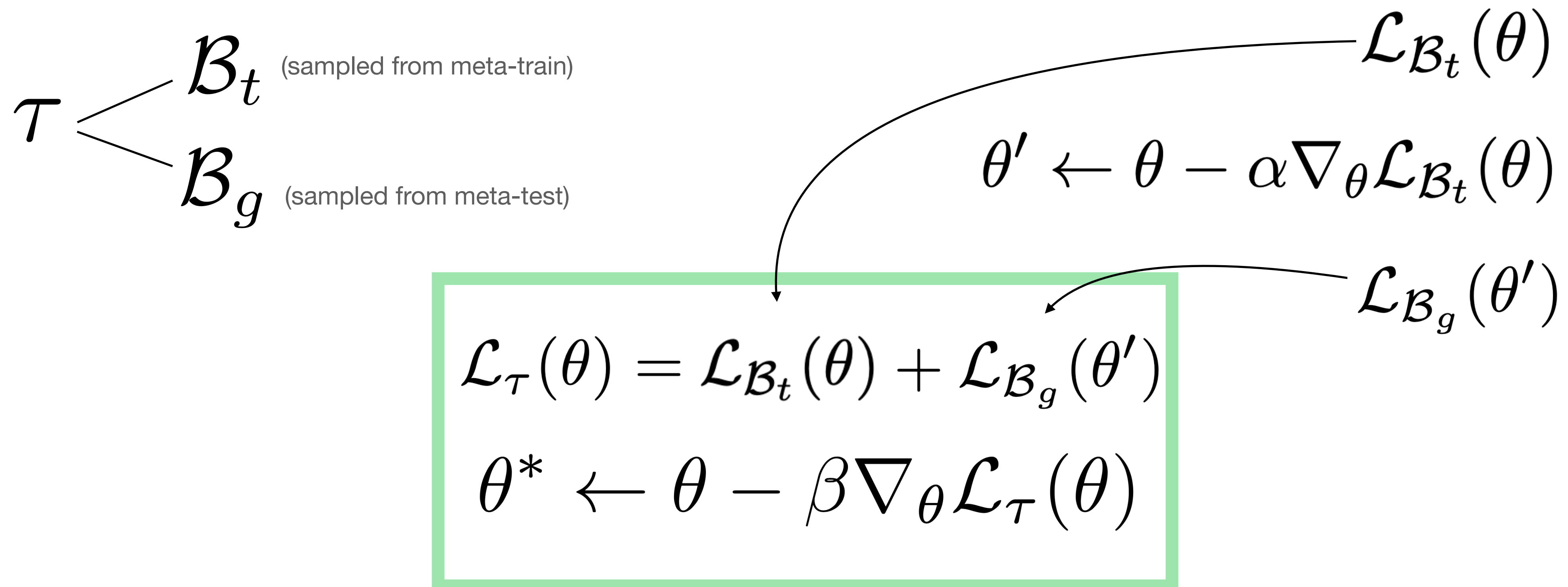
- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML (Li et al., 2018)

- if we know we're going to be tested on something different, let's train for that



Improving Generalization

DG-MAML

||||| Possible Strategies

Improving Generalization

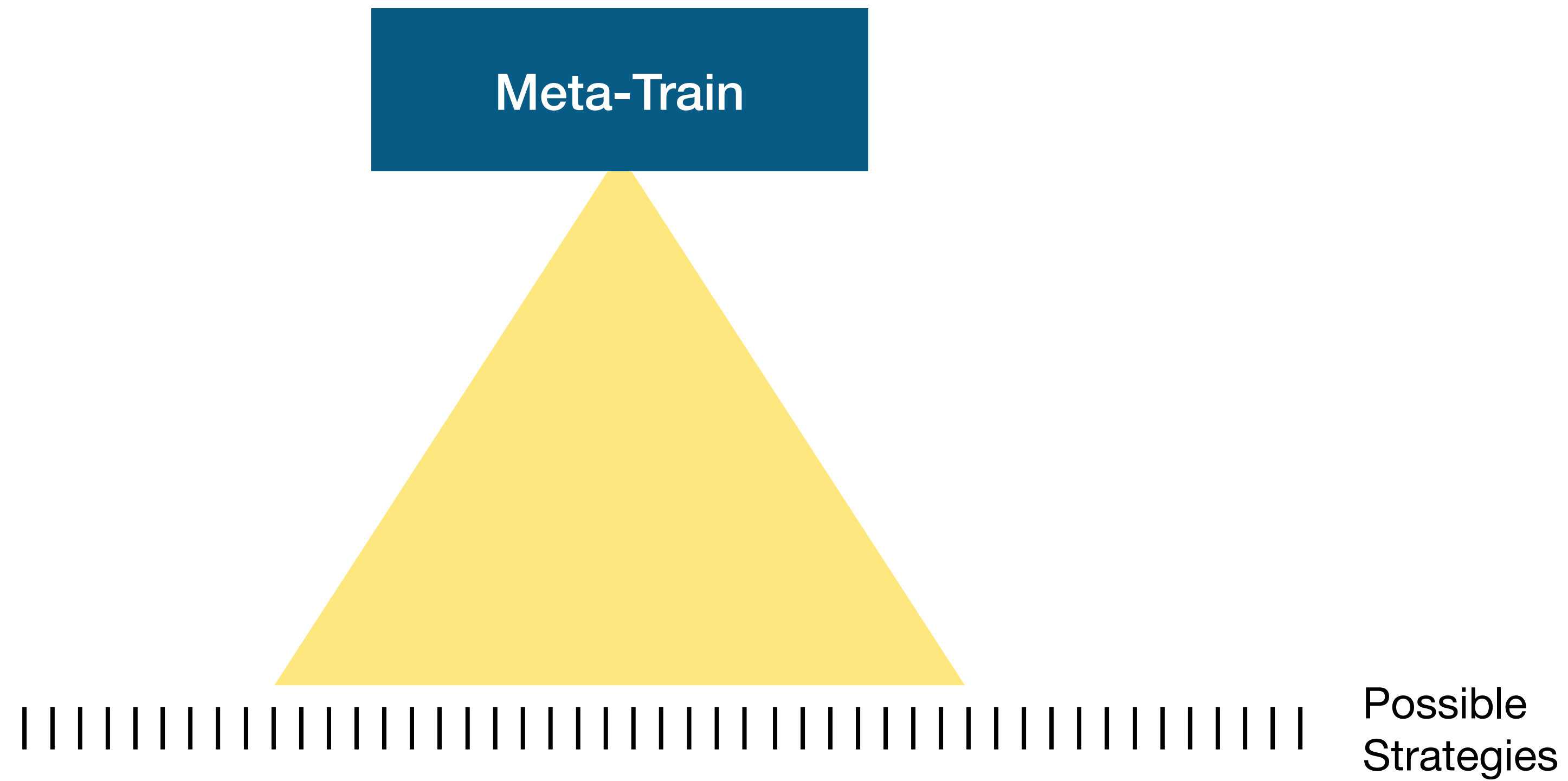
DG-MAML

Meta-Train

||||| Possible Strategies

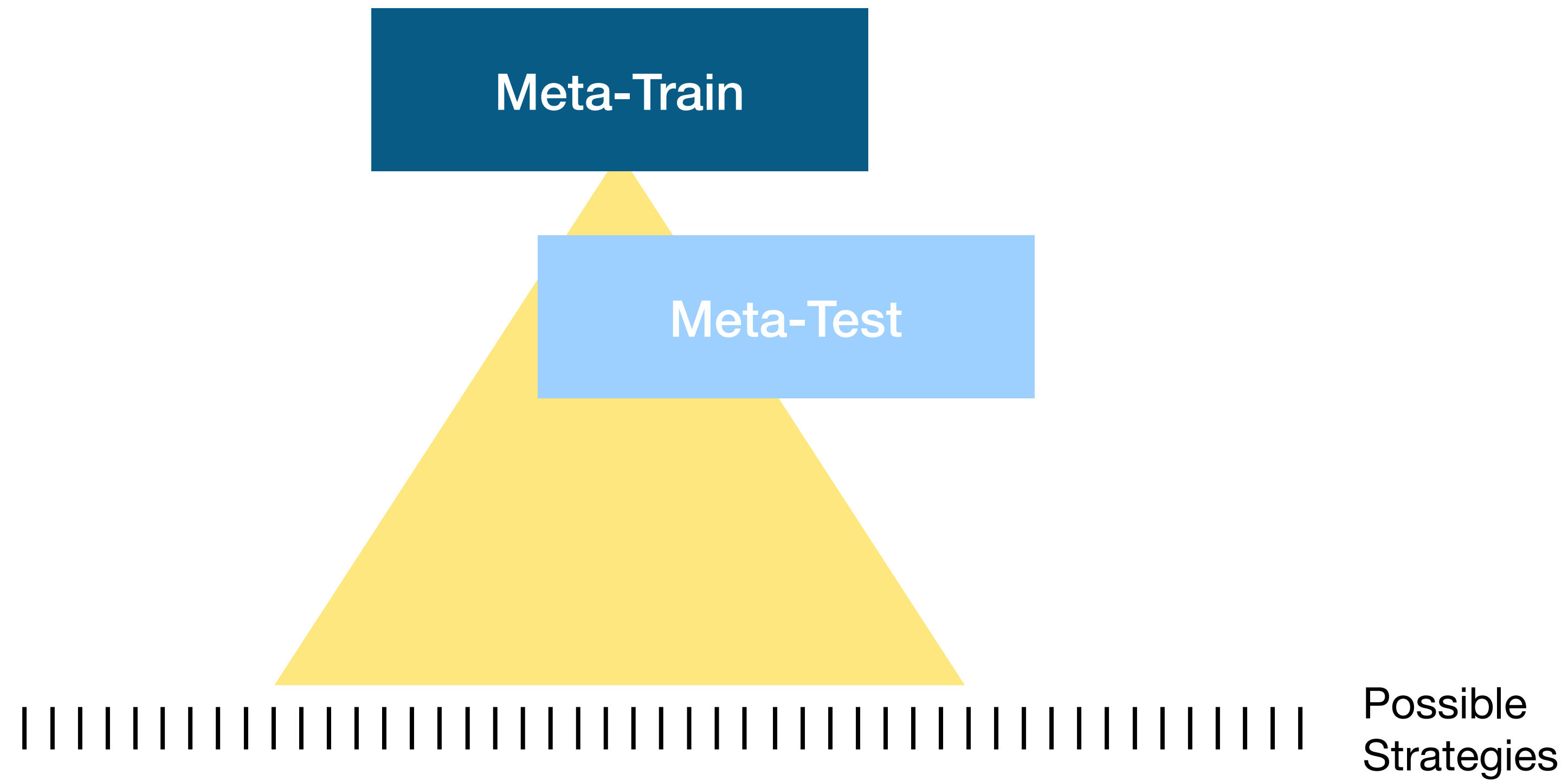
Improving Generalization

DG-MAML



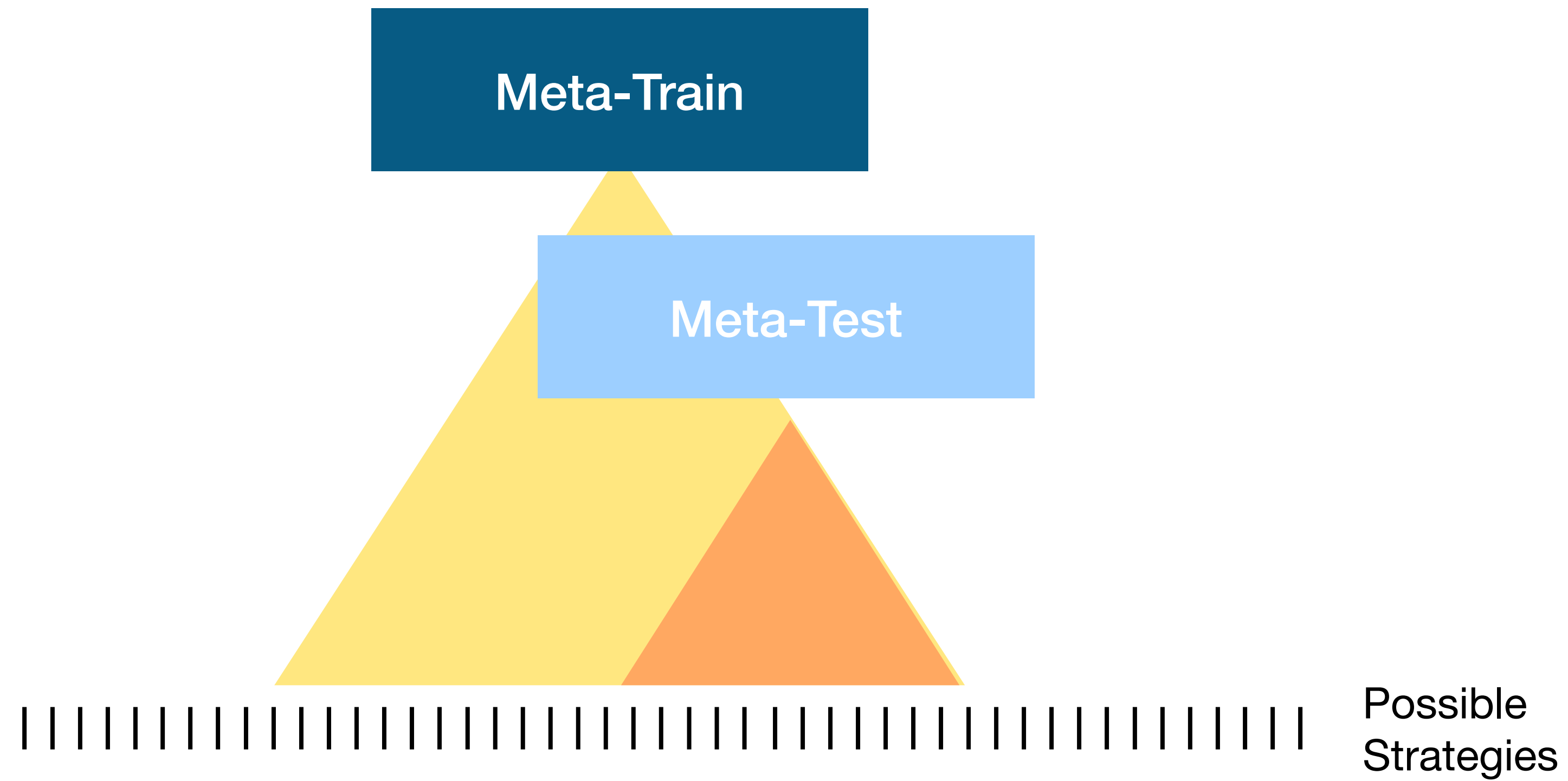
Improving Generalization

DG-MAML



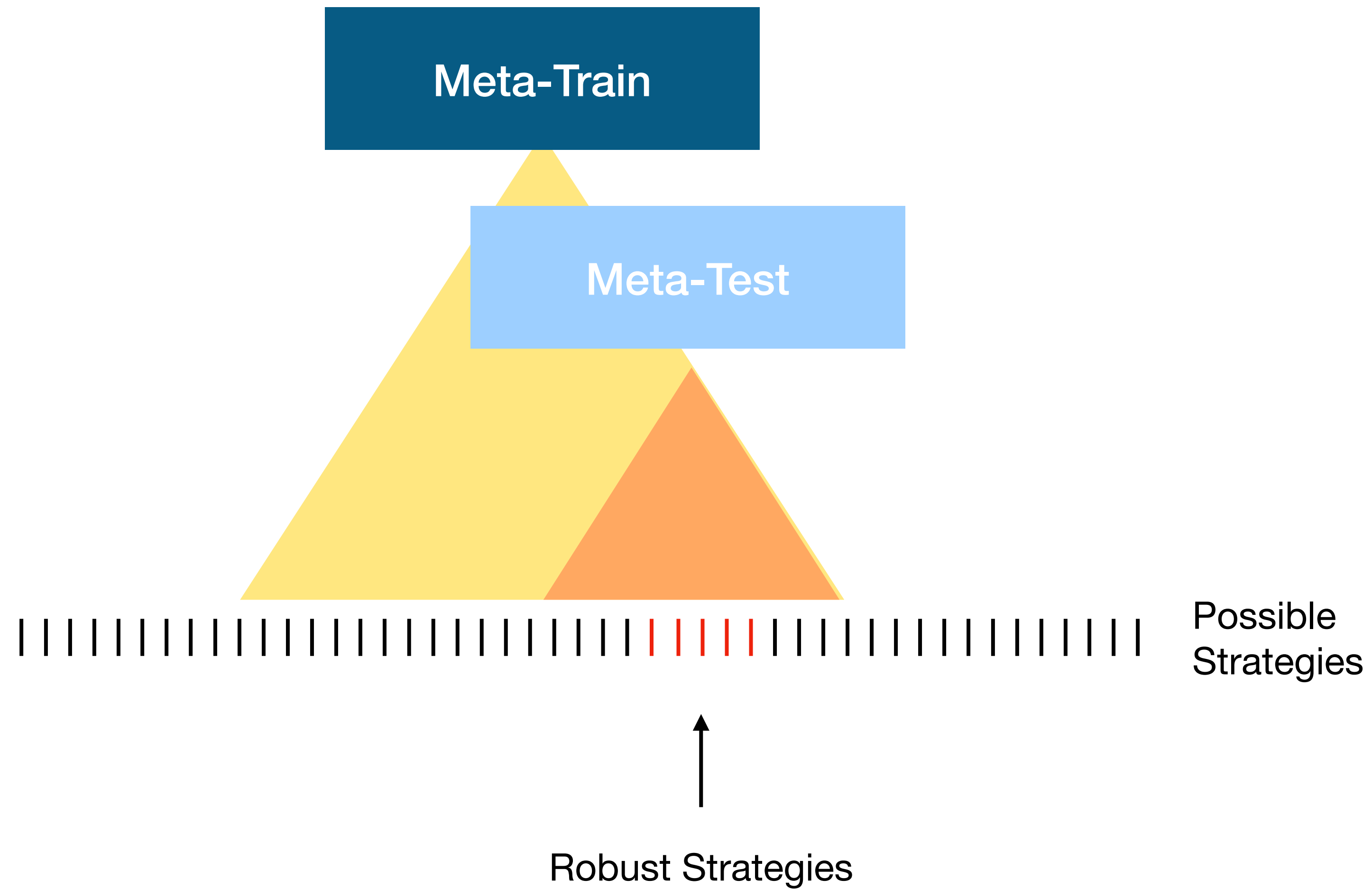
Improving Generalization

DG-MAML



Improving Generalization

DG-MAML



Improving Generalization

Prior Knowledge

Improving Generalization

Prior Knowledge

- But this requires prior knowledge

Improving Generalization

Prior Knowledge

- But this requires prior knowledge
 - not of the test distribution, but of the family of distributions from which the test distribution will be drawn

**how can we help to resolve this
underspecification more domain-generally?**

Domain-General Bias

DG-MAML

Domain-General Bias

DG-MAML

- DG-MAML presents a way to introduce a bias during training

Domain-General Bias

DG-MAML

- DG-MAML presents a way to introduce a bias during training

$$\mathcal{L}_{\tau}(\theta) = \mathcal{L}_{\mathcal{B}_t}(\theta) + \mathcal{L}_{\mathcal{B}_g}(\theta')$$

Domain-General Bias

DG-MAML

- DG-MAML presents a way to introduce a bias during training

$$\mathcal{L}_{\tau}(\theta) = \mathcal{L}_{\mathcal{B}_t}(\theta) + \mathcal{L}_{\mathcal{B}_g}(\theta')$$

Domain-General Bias

DG-MAML

- DG-MAML presents a way to introduce a bias during training

$$\mathcal{L}_\tau(\theta) = \mathcal{L}_{\mathcal{B}_t}(\theta) + \mathcal{L}_{\mathcal{B}_g}(\theta')$$

- Whatever we put in the meta-test batch constrains our update step on meta-train

Domain-General Bias

DG-MAML

- DG-MAML presents a way to introduce a bias during training

$$\mathcal{L}_\tau(\theta) = \mathcal{L}_{\mathcal{B}_t}(\theta) + \mathcal{L}_{\mathcal{B}_g}(\theta')$$

- Whatever we put in the meta-test batch constrains our update step on meta-train

Let's use this to introduce a general bias rather than a task specific one

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

Meta-Test

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

Meta-Test

The girl changed a sandwich by the bed.

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Meta-Test

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Uniform Sampling | Uni-MAML

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Uniform Sampling | Uni-MAML

The penguin ate a donut.

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Uniform Sampling | Uni-MAML

The penguin ate a donut.

Amelia gave Emma a strawberry.

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Uniform Sampling | Uni-MAML

The penguin ate a donut.

Amelia gave Emma a strawberry.

A cat disintegrated a girl.

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Uniform Sampling | Uni-MAML

The penguin ate a donut.

Amelia gave Emma a strawberry.

A cat disintegrated a girl.

A visitor was posted a rose by a turtle.

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Lev Distance | Lev-MAML

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Lev Distance | Lev-MAML

The girl rolled a drink beside the table.

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Lev Distance | Lev-MAML

The girl rolled a drink beside the table.
The girl liked a dealer beside the table .

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Lev Distance | Lev-MAML

The girl rolled a drink beside the table.
The girl liked a dealer beside the table .

The sailor dusted a girl.

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Lev Distance | Lev-MAML

The girl rolled a drink beside the table.
The girl liked a dealer beside the table .

The sailor dusted a girl.
The girl dusted a boy.

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Convolutional String Kernel | Str-MAML

The girl rolled a drink beside the table.
The girl liked a dealer beside the table .

The sailor dusted a girl.
The girl dusted a boy.

Domain-General Bias

Memorization

- let's introduce a bias that impairs the model's ability to memorize whole sentences

Meta-Train

The girl changed a sandwich by the bed.

The sailor dusted a boy.

Partial Tree Kernel | Tree-MAML

Mateo dusted a boy .
dust

A sandwich changed.
A block was changed by the girl.

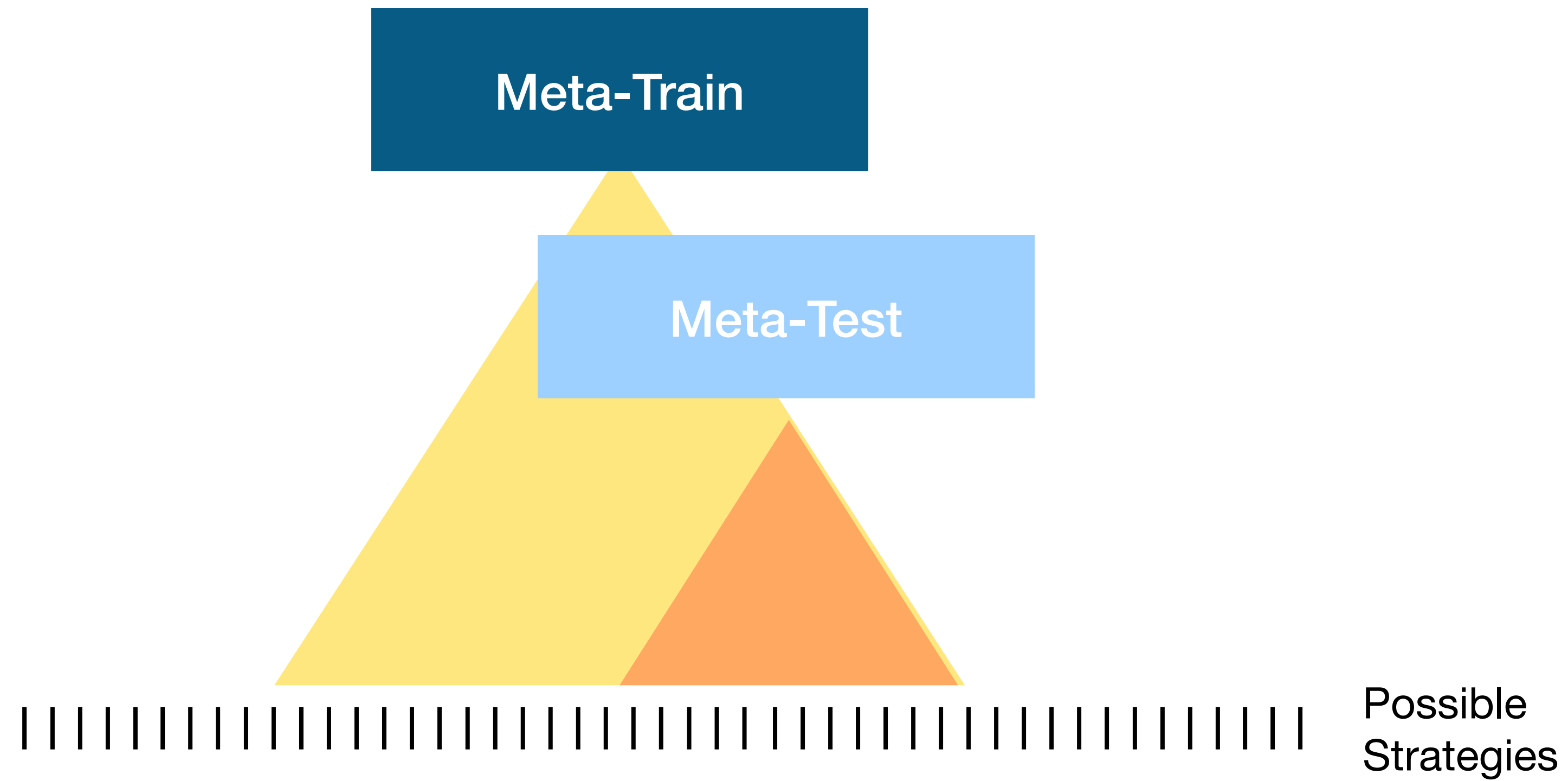
**Does inhibiting models' memory
improve generalization?**

Experiments

COGS (& SCAN)

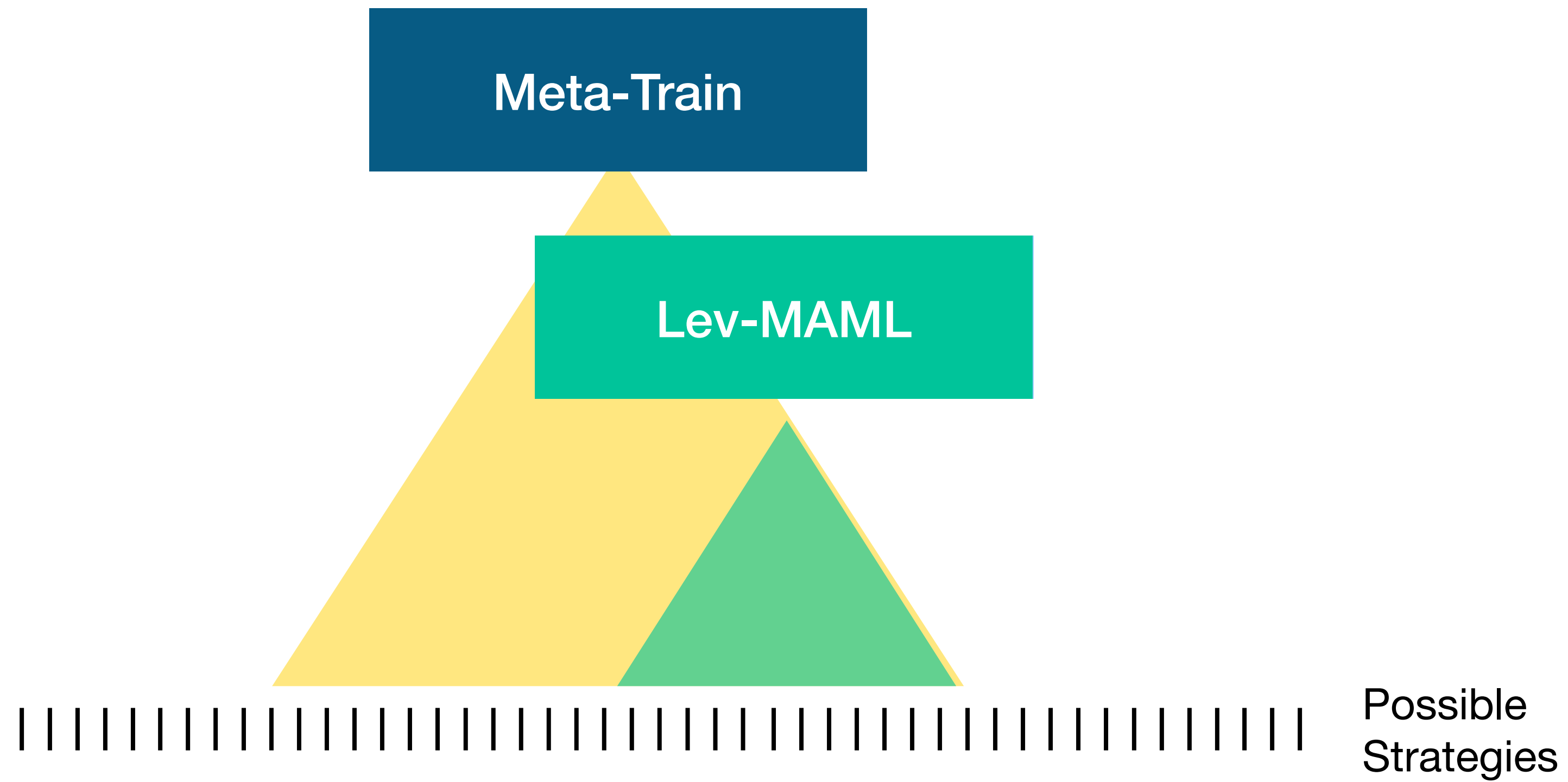
Experiments

COGS (& SCAN)



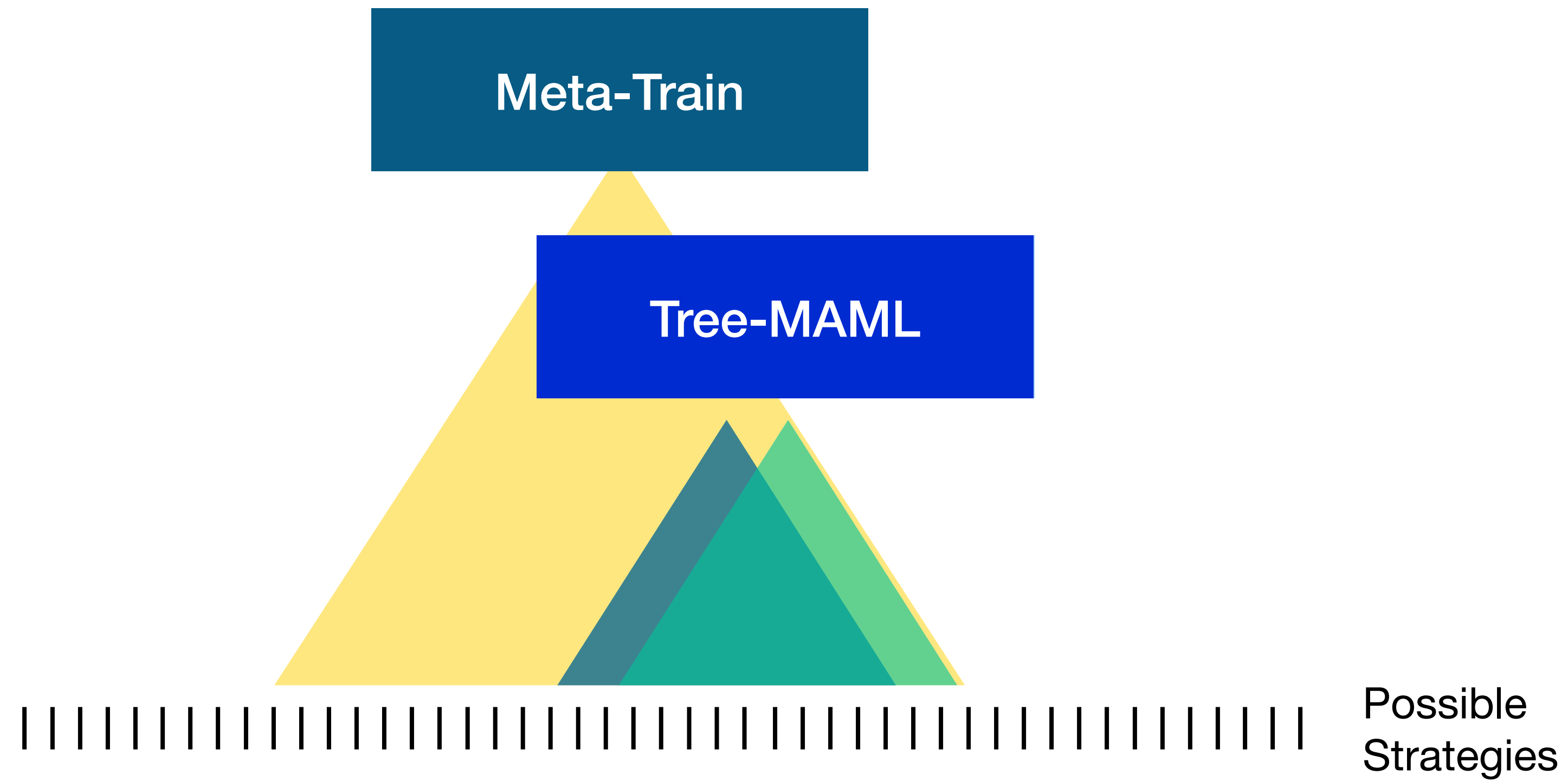
Experiments

COGS (& SCAN)



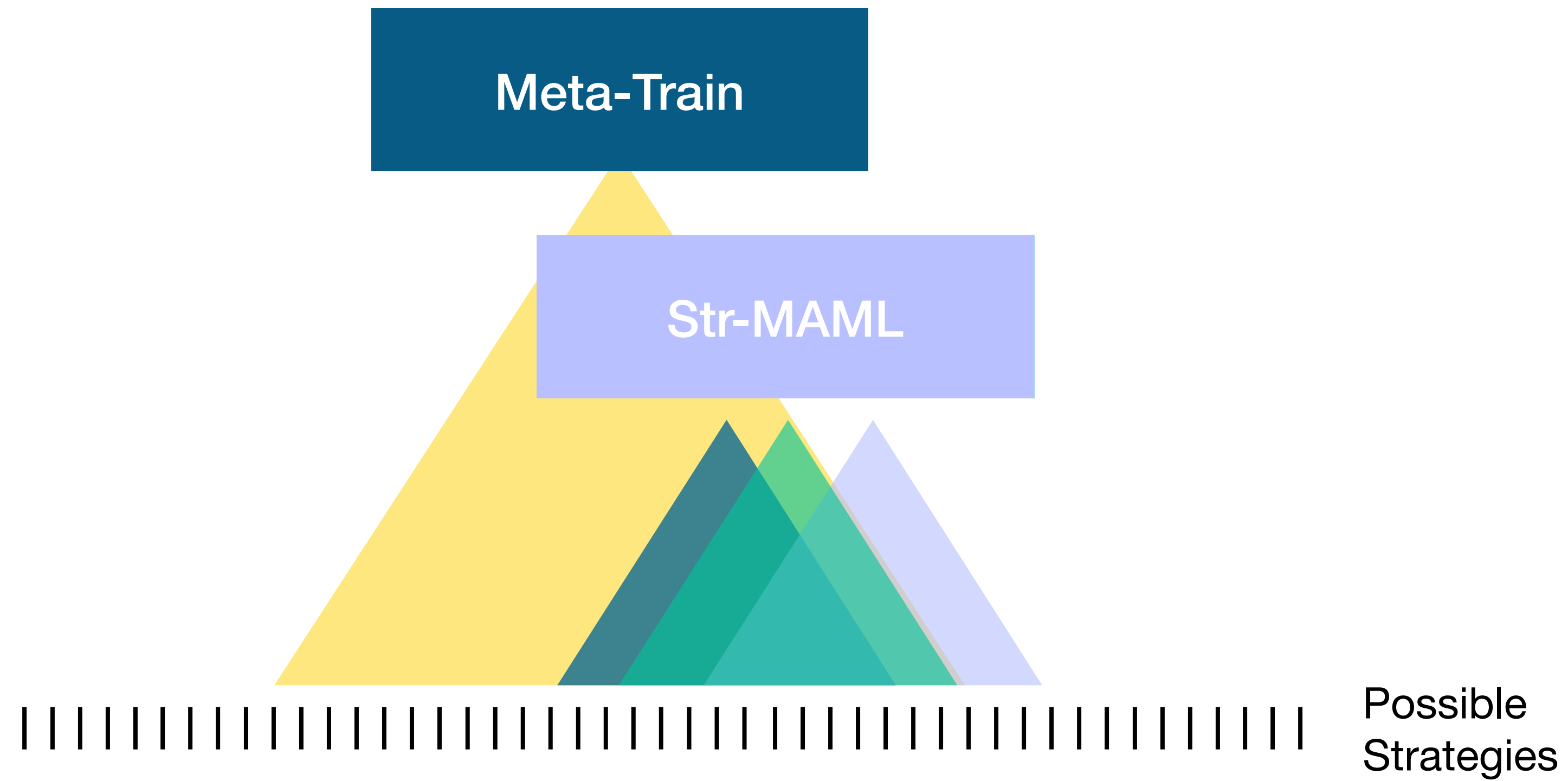
Experiments

COGS (& SCAN)



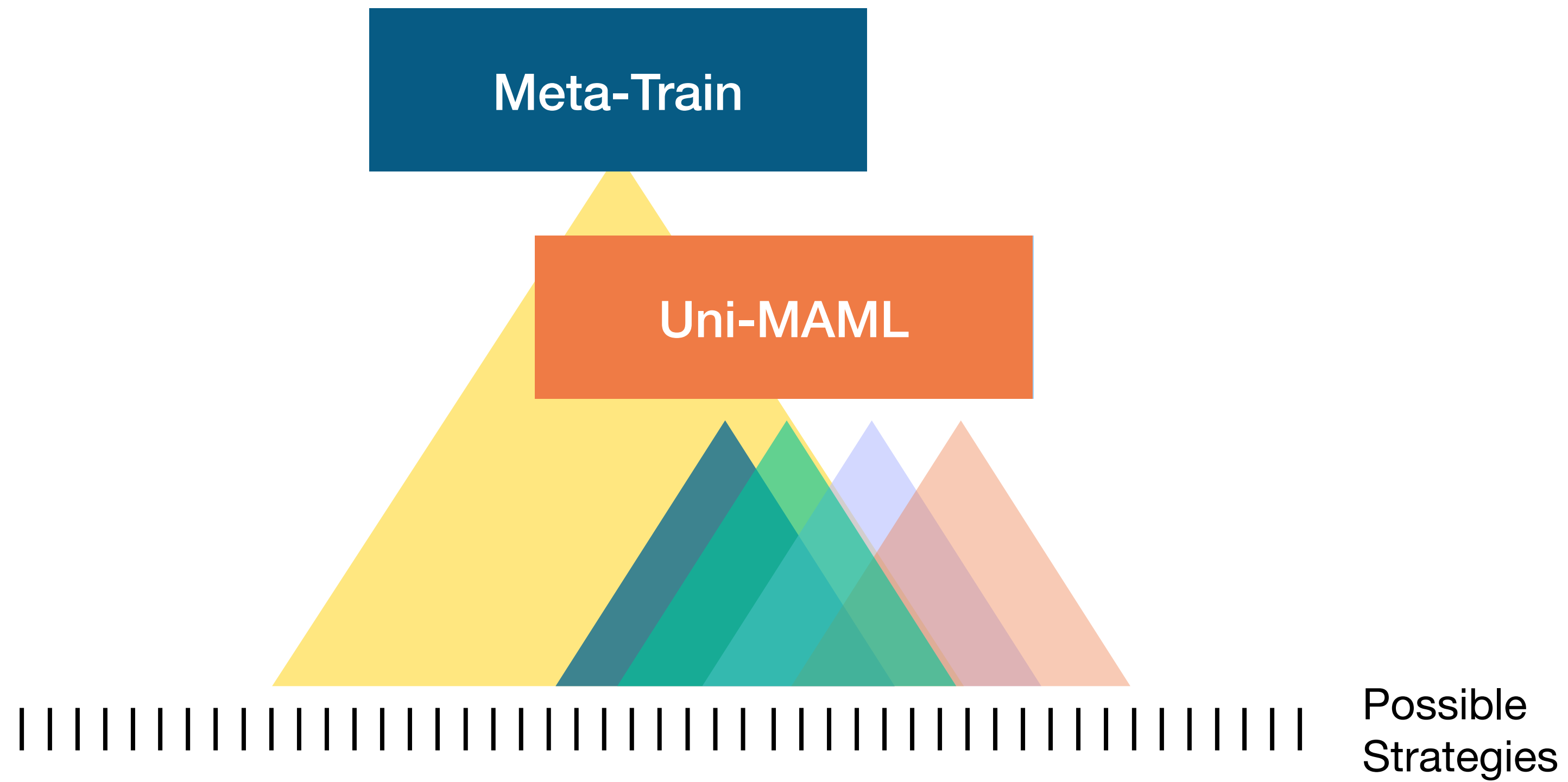
Experiments

COGS (& SCAN)



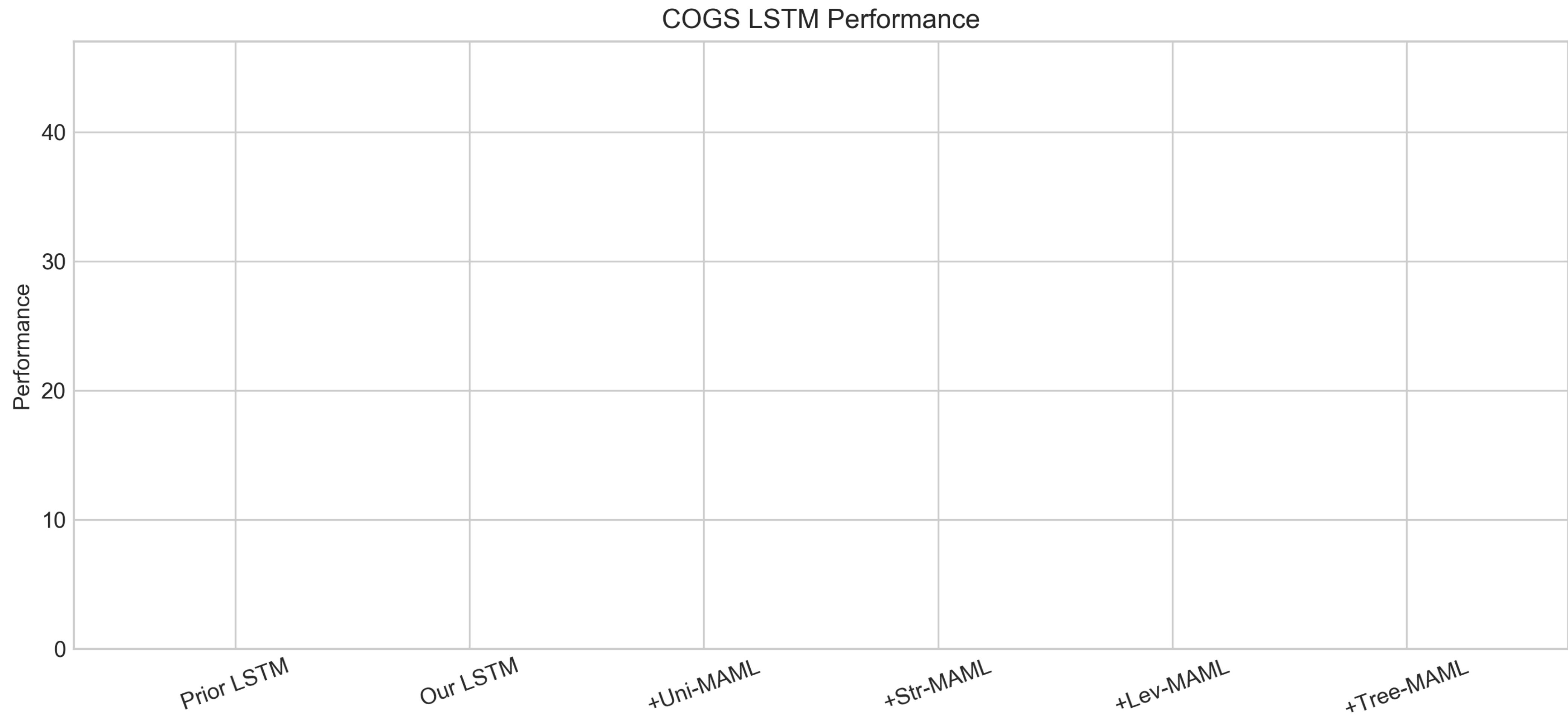
Experiments

COGS (& SCAN)



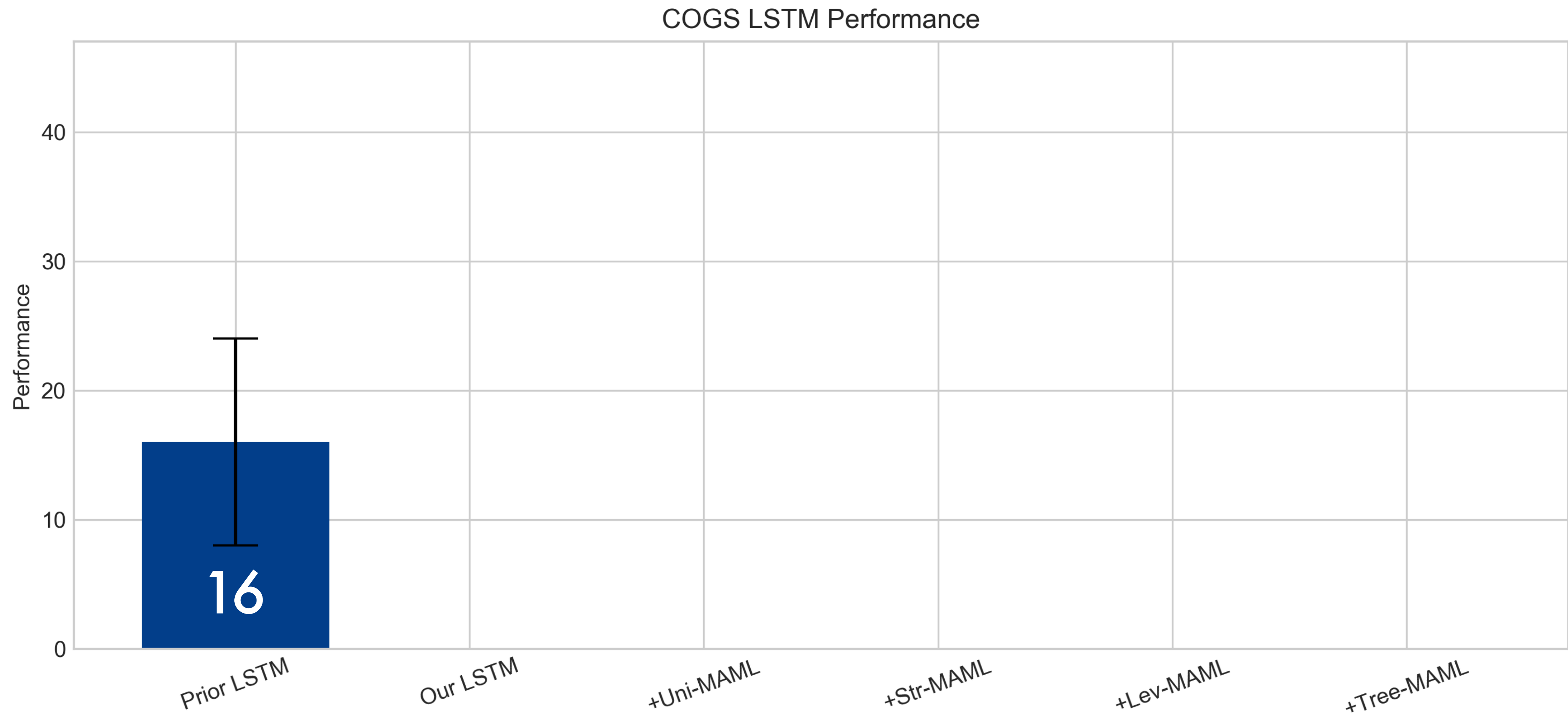
Experiments

COGS Dataset



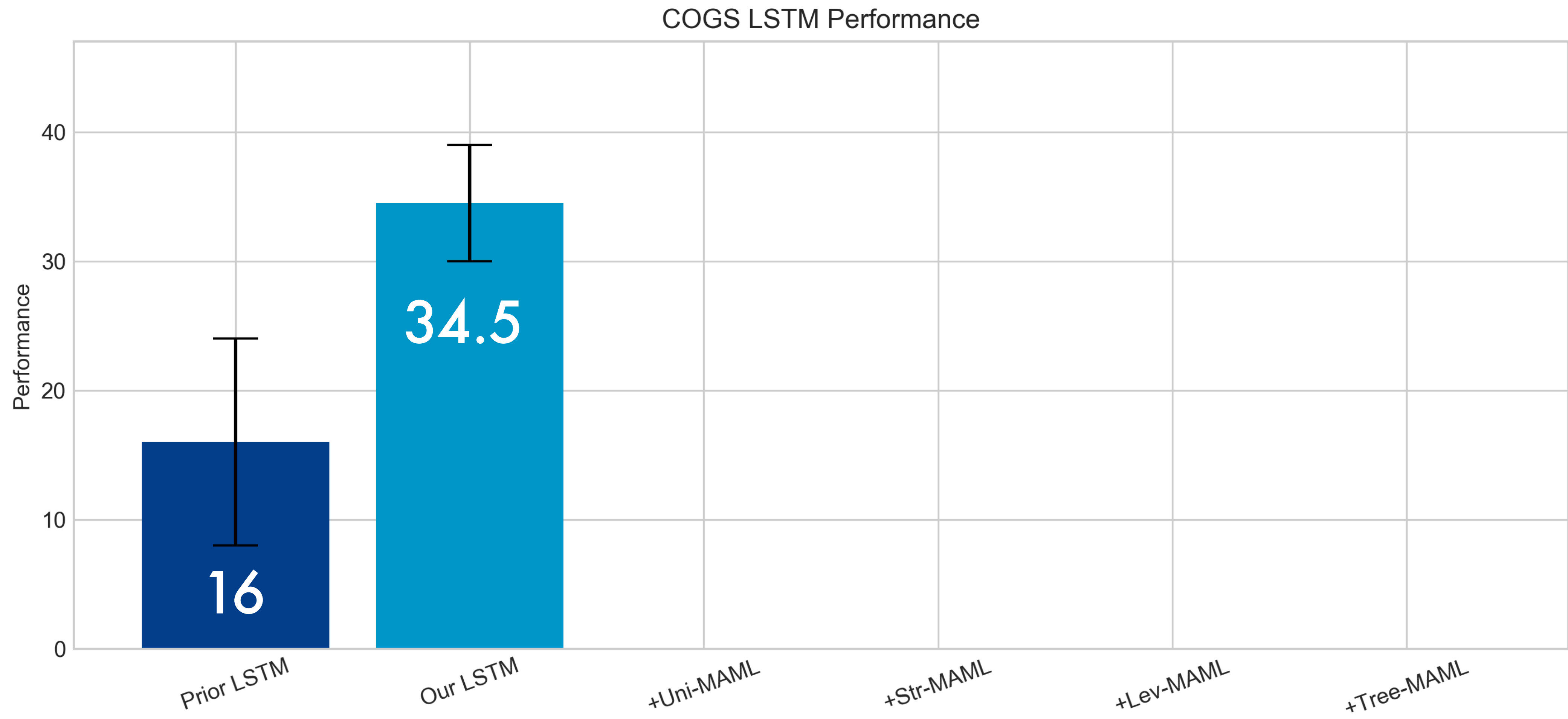
Experiments

COGS Dataset



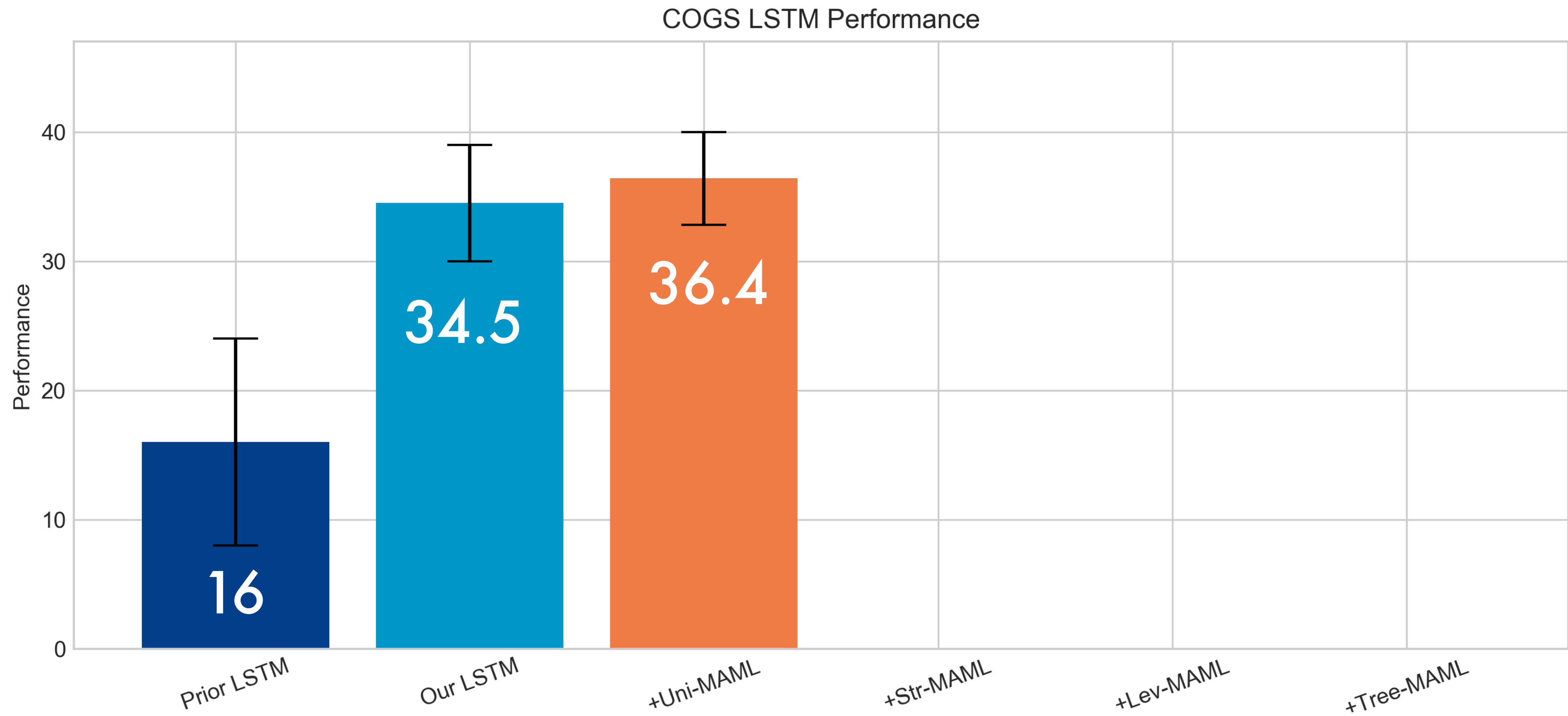
Experiments

COGS Dataset



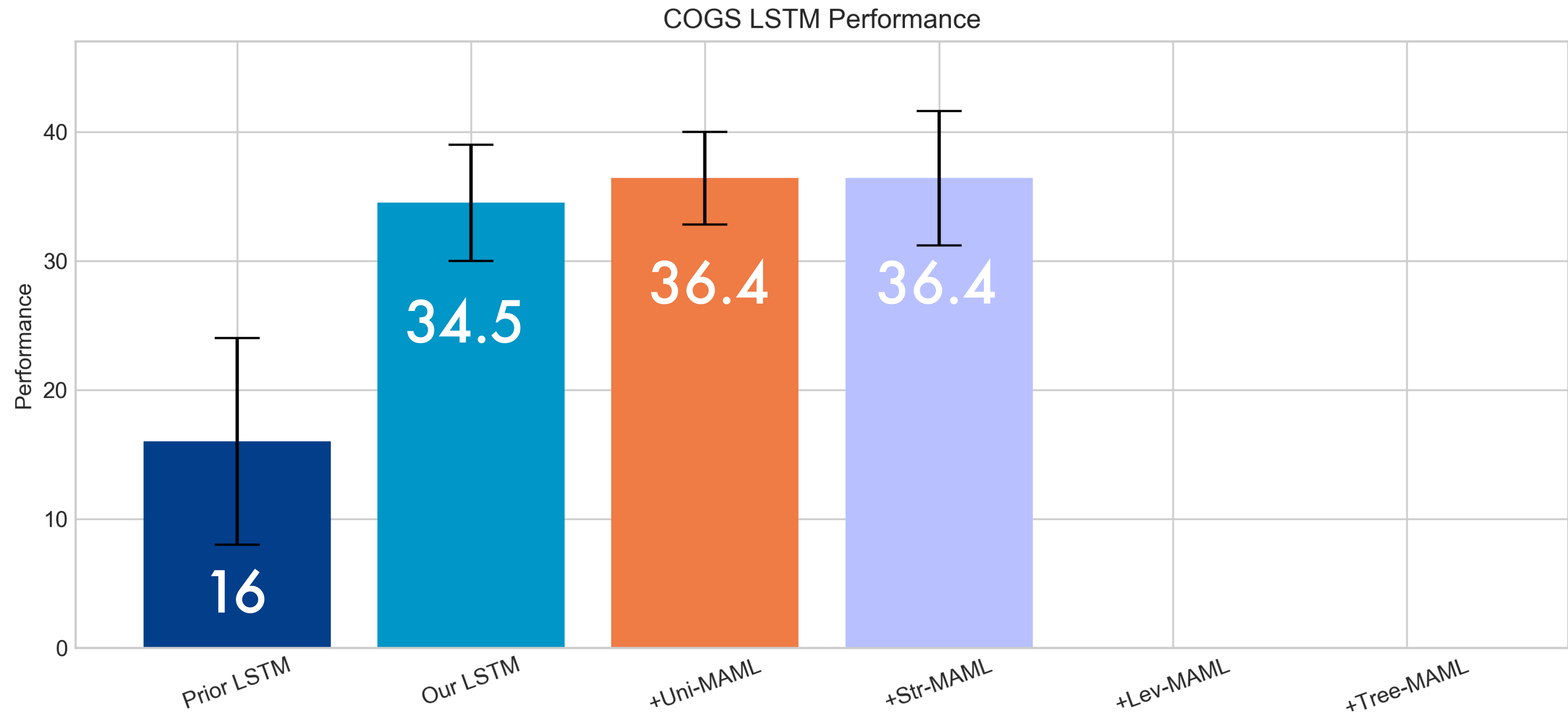
Experiments

COGS Dataset



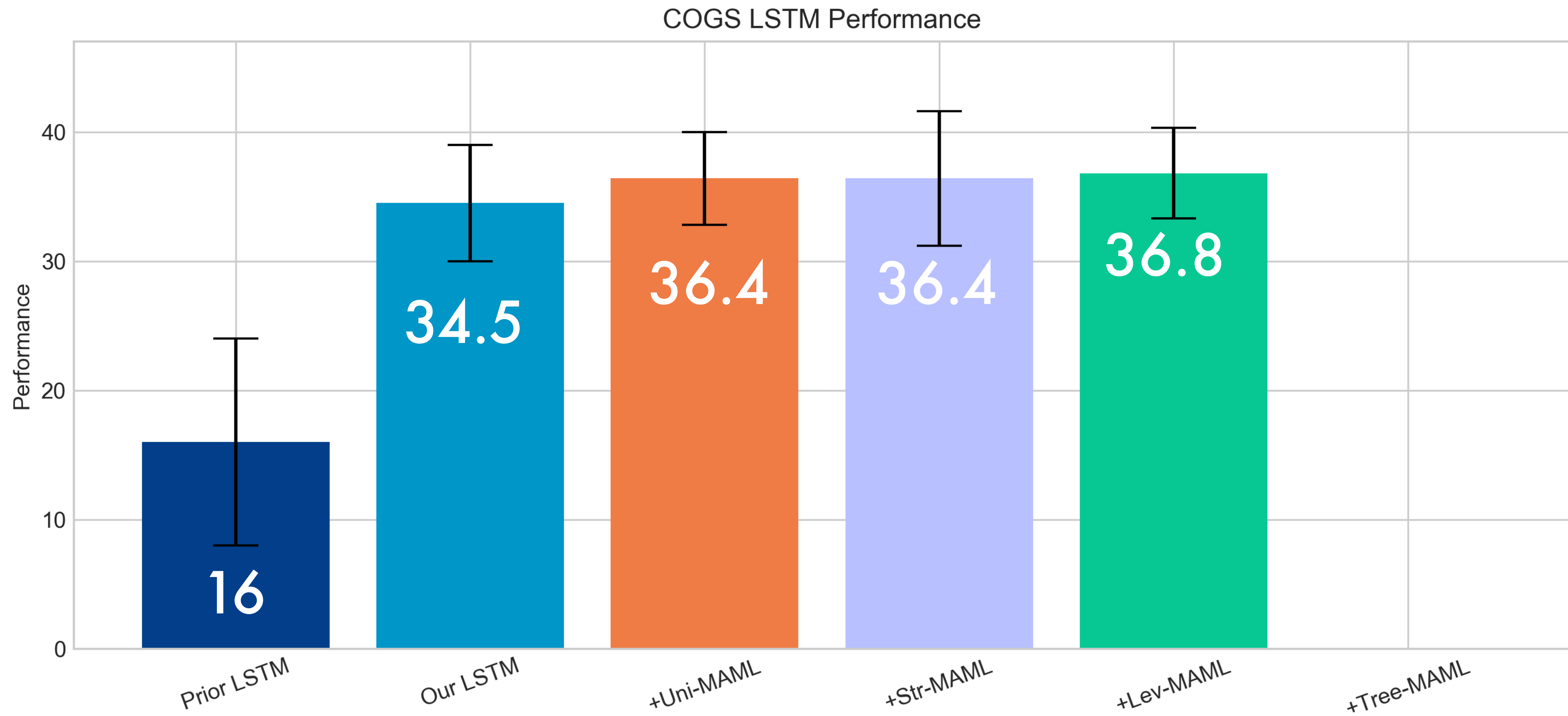
Experiments

COGS Dataset



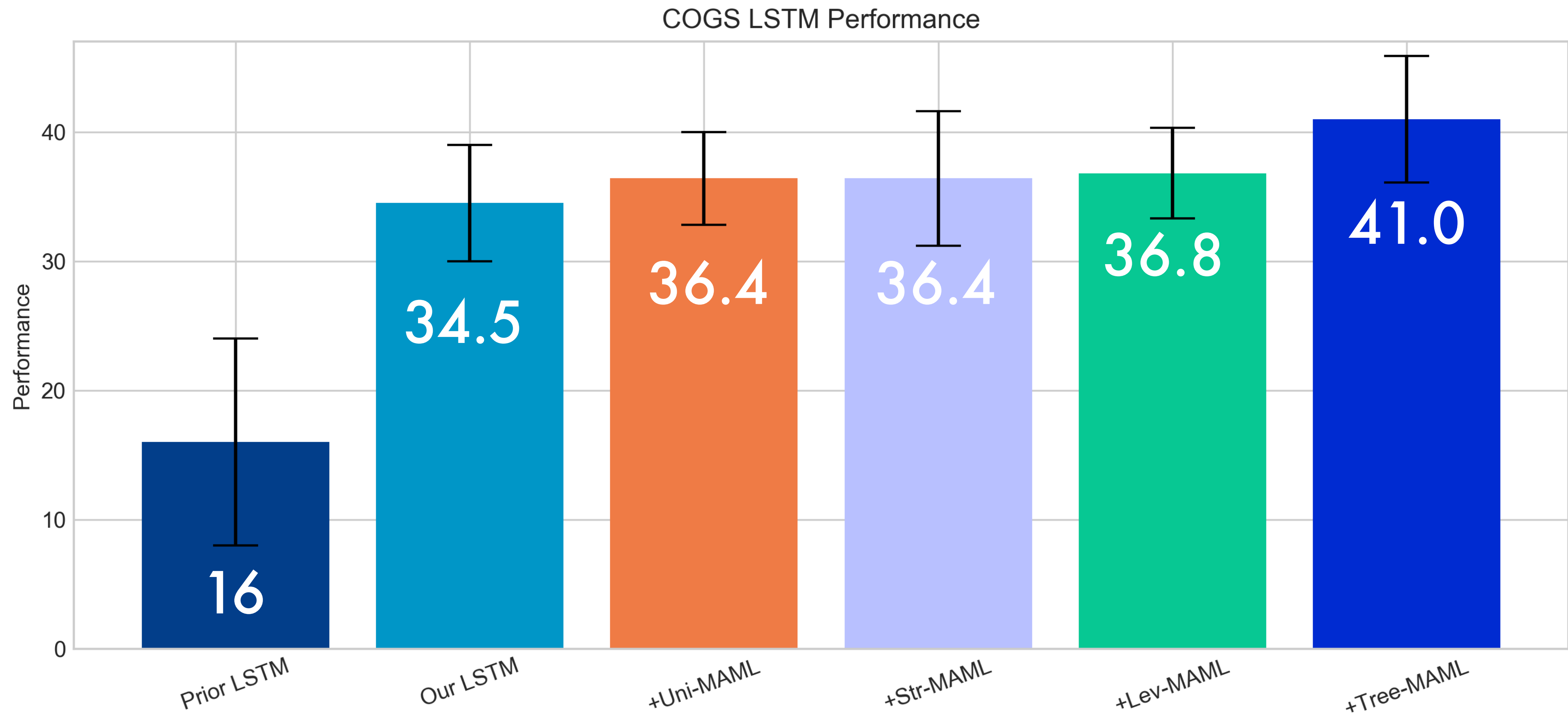
Experiments

COGS Dataset



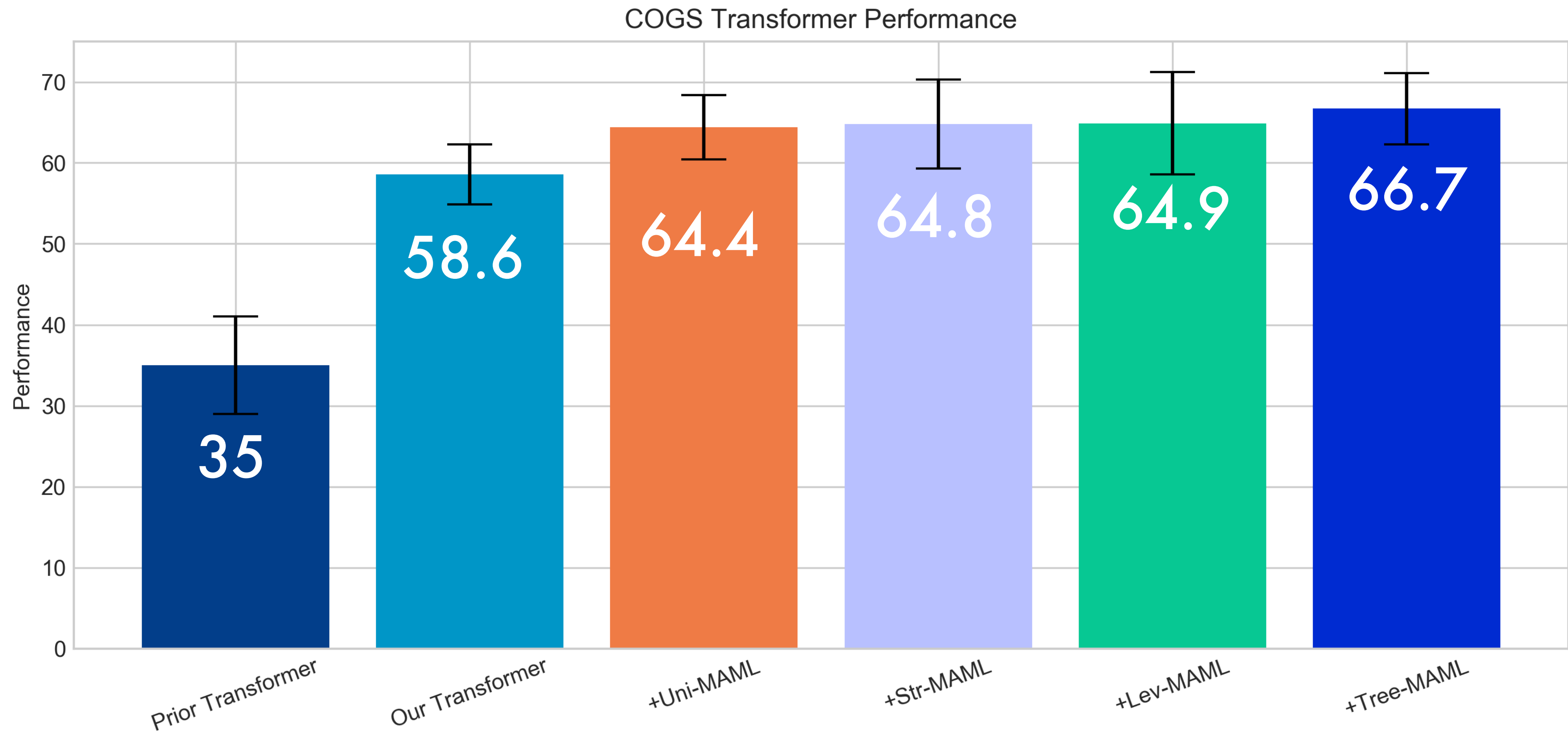
Experiments

COGS Dataset



Experiments

COGS Dataset



Conclusions

Conclusions

- Biasing the model against memorization **does improve generalization performance** on COGS and SCAN

Conclusions

- Biasing the model against memorization **does improve generalization performance** on COGS and SCAN
- Unlike more task specific methods this approach is **model** and in many cases **data-set agnostic**

Conclusions

- Biasing the model against memorization **does improve generalization performance** on COGS and SCAN
- Unlike more task specific methods this approach is **model** and in many cases **data-set agnostic**
- The design of the Meta-Test task allows for the **design of the bias applied** during training

thank
you

Henry Conklin*

@henryconklin

Bailin Wang*

@bailin_28

Kenny Smith

@kennysmithed

Ivan Titov

@iatitov