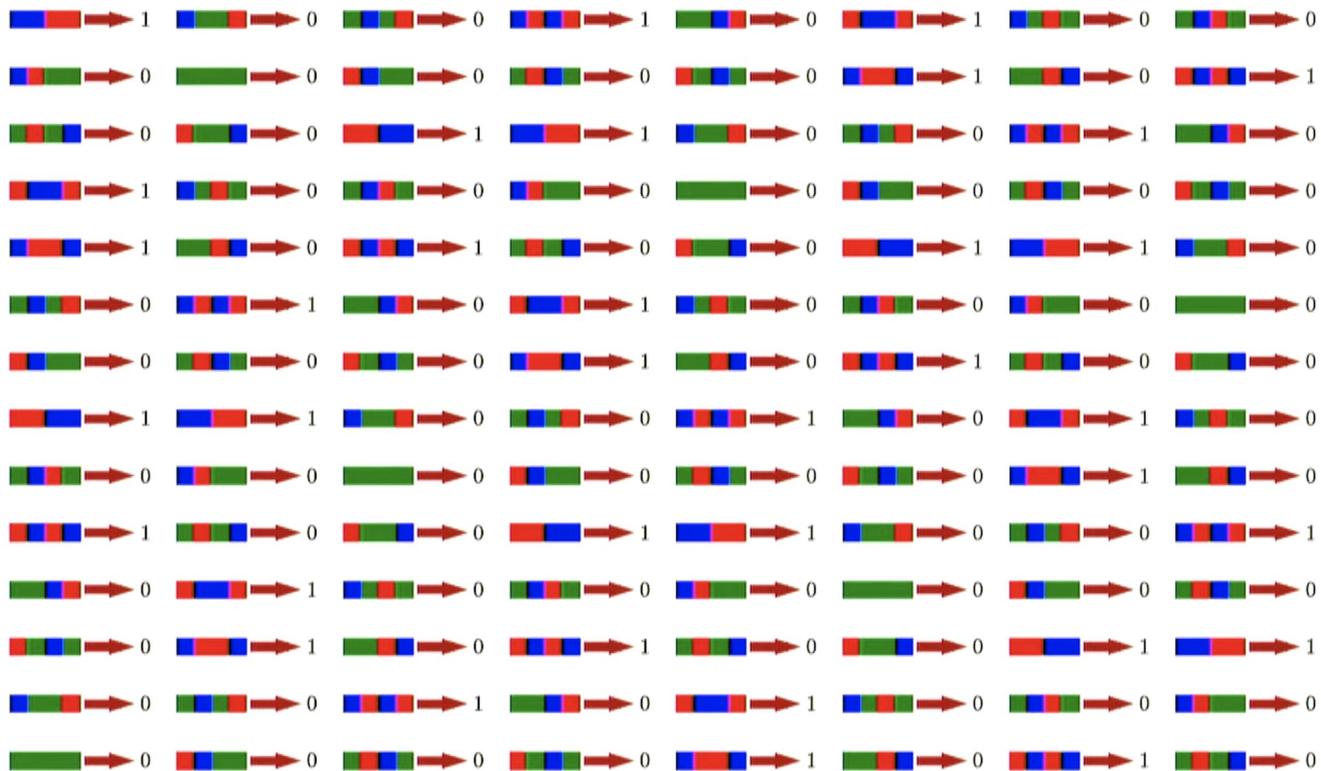


Shaping Visual Representations with Language for Few-Shot Classification



Jesse Mu Percy Liang Noah Goodman

ACL 2020



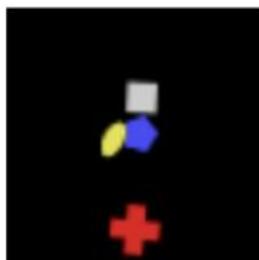
(courtesy Percy Liang)

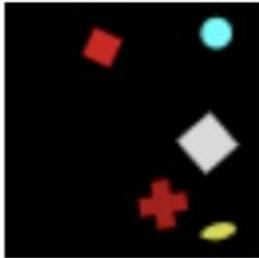
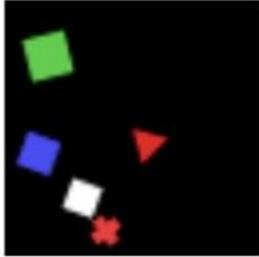
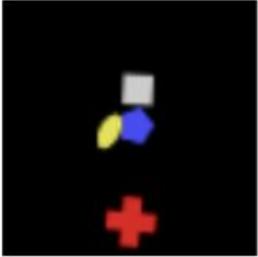
at least 2 red squares

How can language **guide representation learning**, especially when data is scarce?

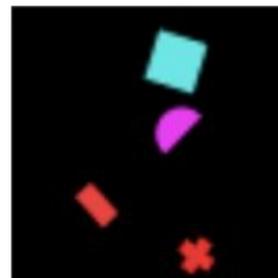
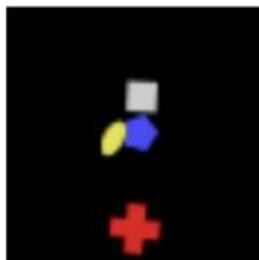
How can language **guide representation learning**, especially when data is scarce?

We study the (underexplored!) setting where language is **available** at train time, but **unavailable** for new tasks at test time

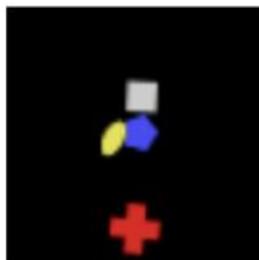




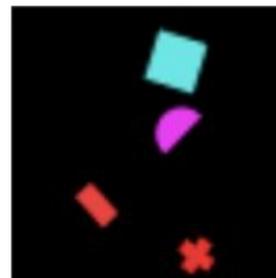
?



true



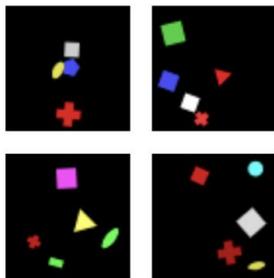
a red cross is below
a square



true

Meta-Train

Task 1

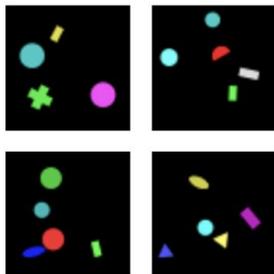


true



a red cross is below
a square

Task 2



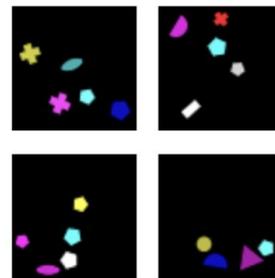
false



a cyan circle is to
the left of a
rectangle

Language
descriptions

Meta-Test

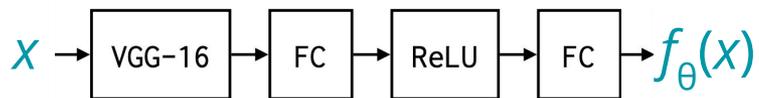
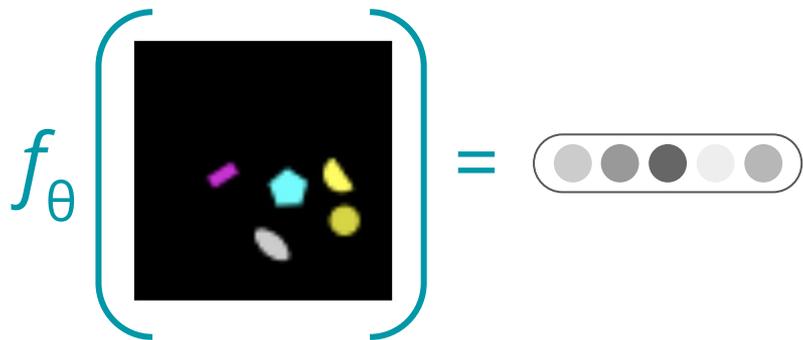


true



Prototype networks (Snell et al., 2017)

Prototype networks (Snell et al., 2017)

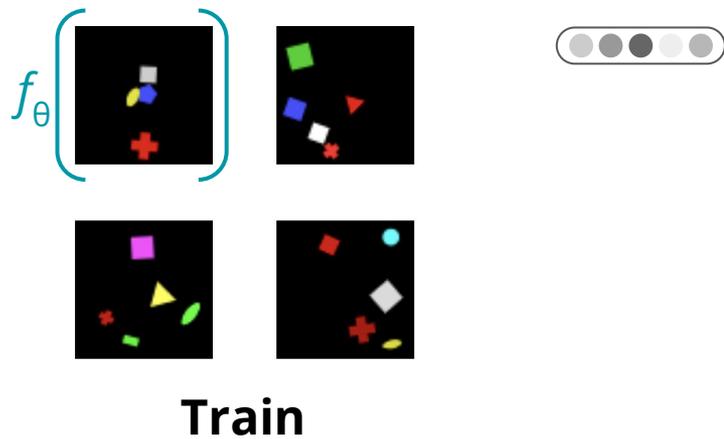


Prototype networks (Snell et al., 2017)



Train

Prototype networks (Snell et al., 2017)

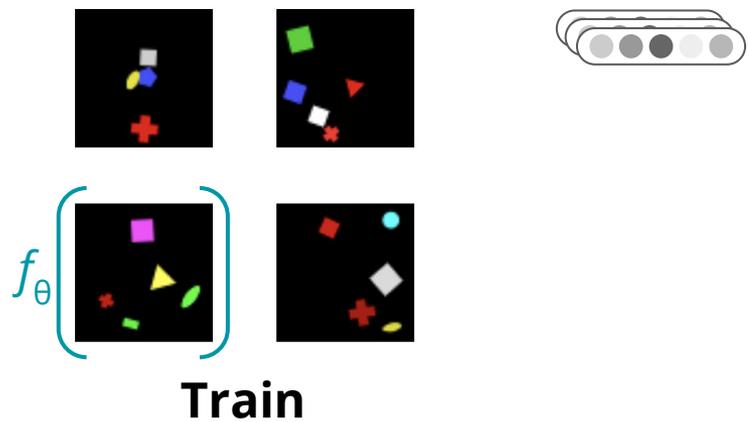


Prototype networks (Snell et al., 2017)

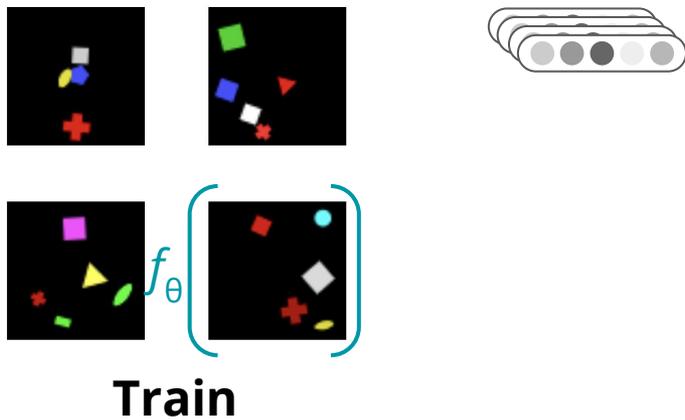


Train

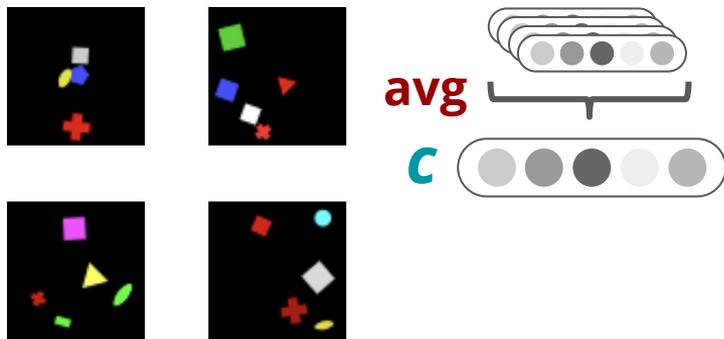
Prototype networks (Snell et al., 2017)



Prototype networks (Snell et al., 2017)

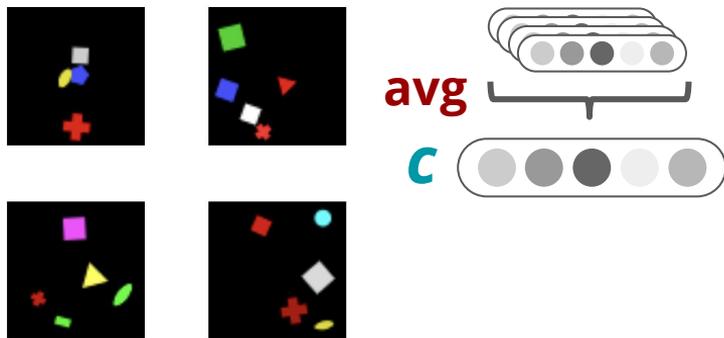


Prototype networks (Snell et al., 2017)



Train

Prototype networks (Snell et al., 2017)

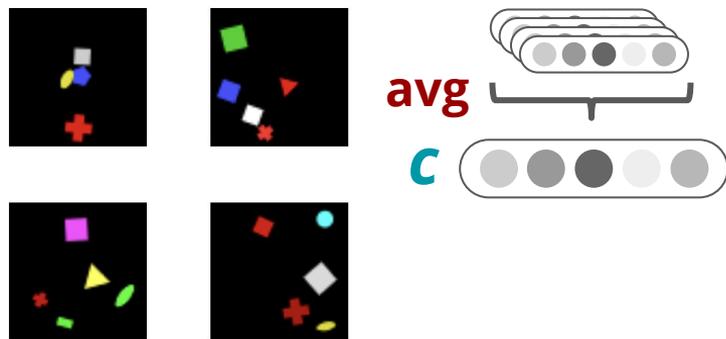


Train



Test

Prototype networks (Snell et al., 2017)

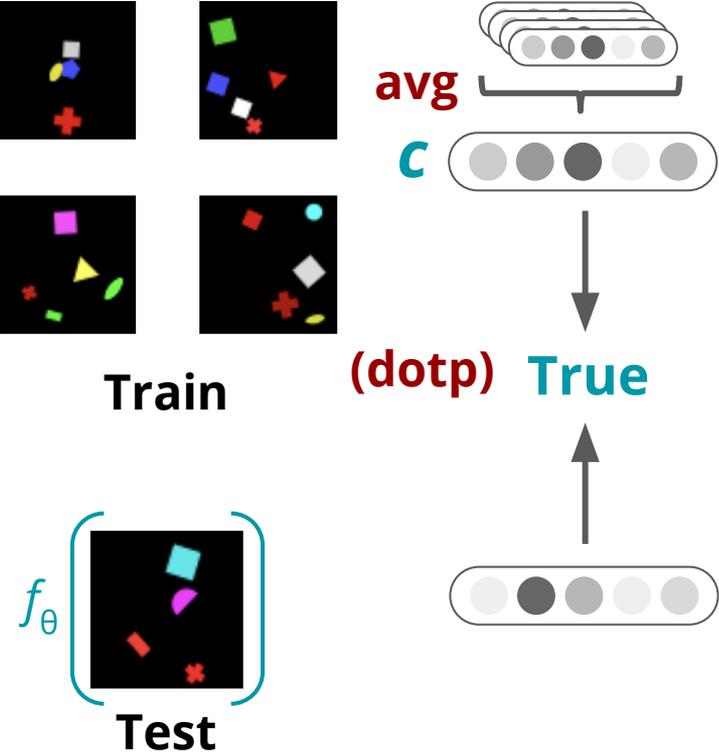


Train

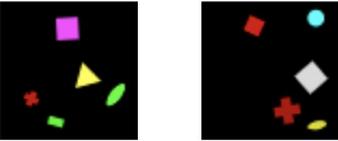
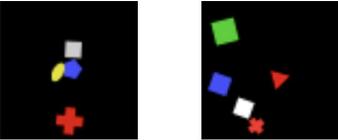


Test

Prototype networks (Snell et al., 2017)



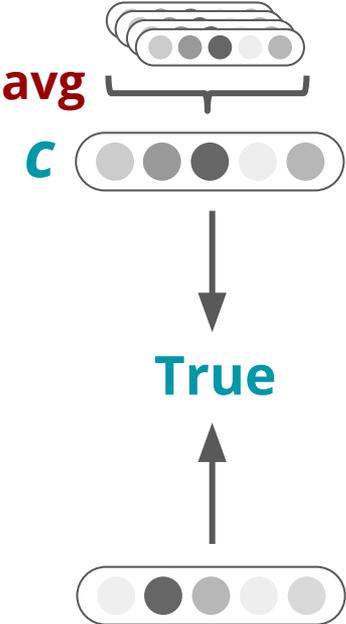
Prototype networks (Snell et al., 2017)



Train



Test



Minimize

$$\arg \min_{\theta} \mathcal{L}_{CLS}(\theta)$$

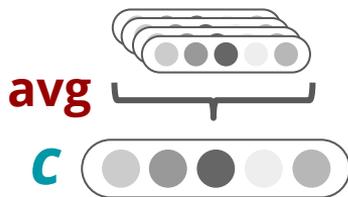
Language-shaped learning (LSL): **Train**



Train



Test



True



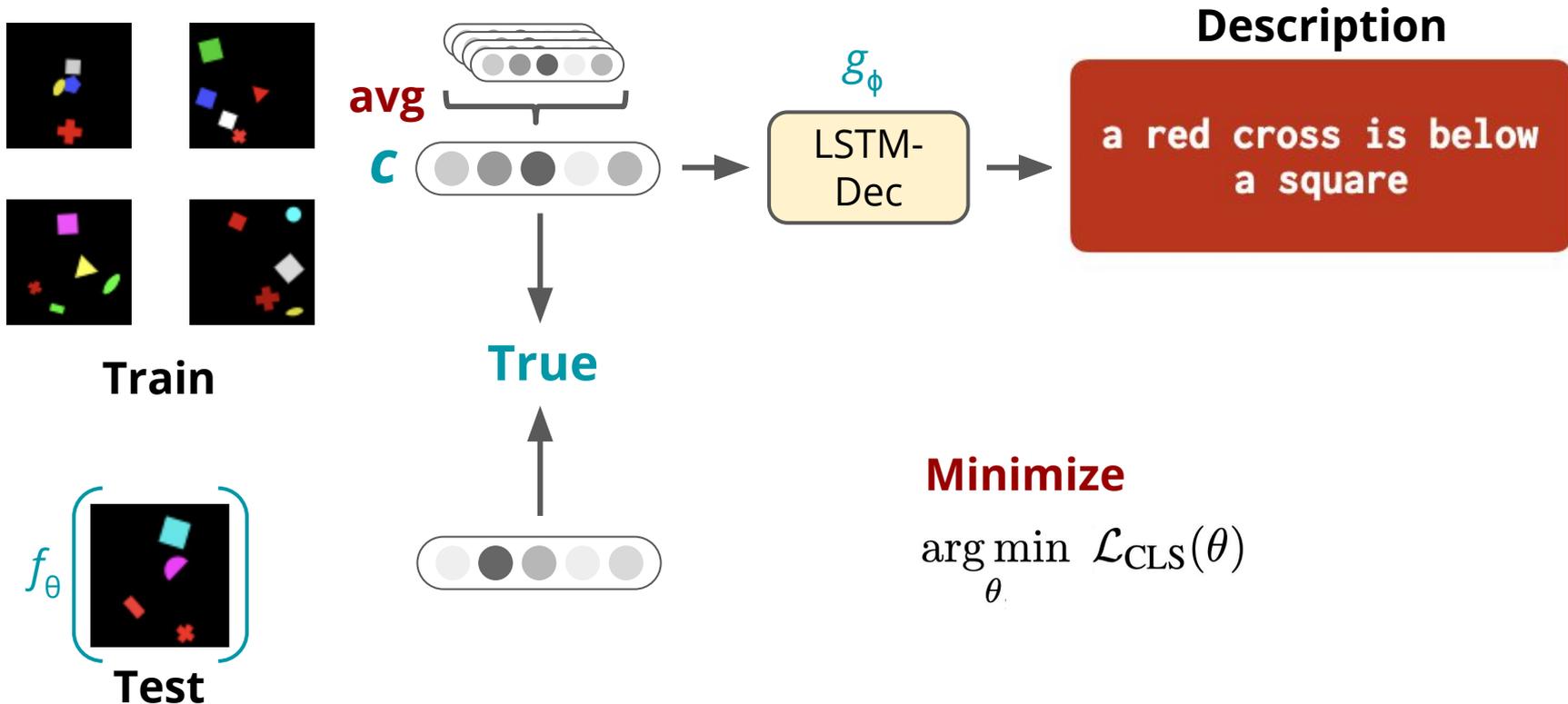
Description

a red cross is below
a square

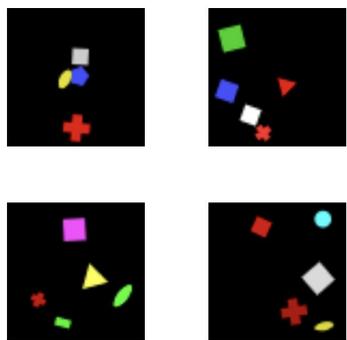
Minimize

$$\arg \min_{\theta} \mathcal{L}_{\text{CLS}}(\theta)$$

Language-shaped learning (LSL): **Train**



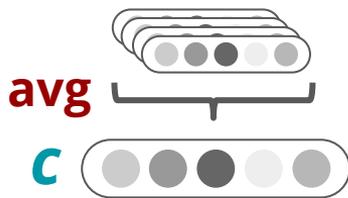
Language-shaped learning (LSL): Train



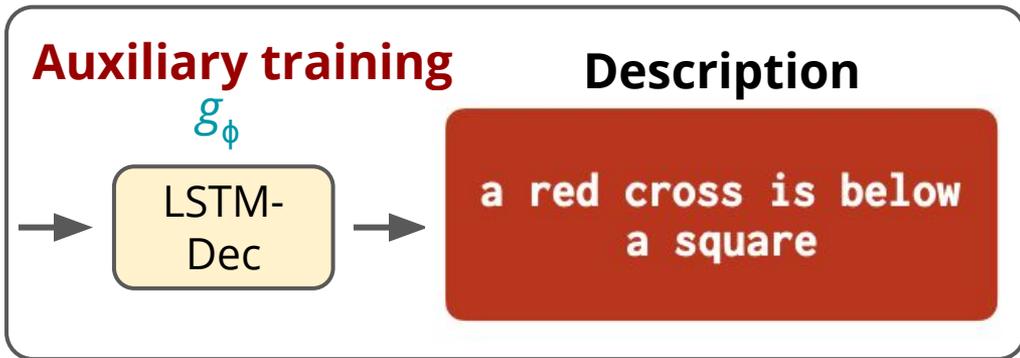
Train



Test



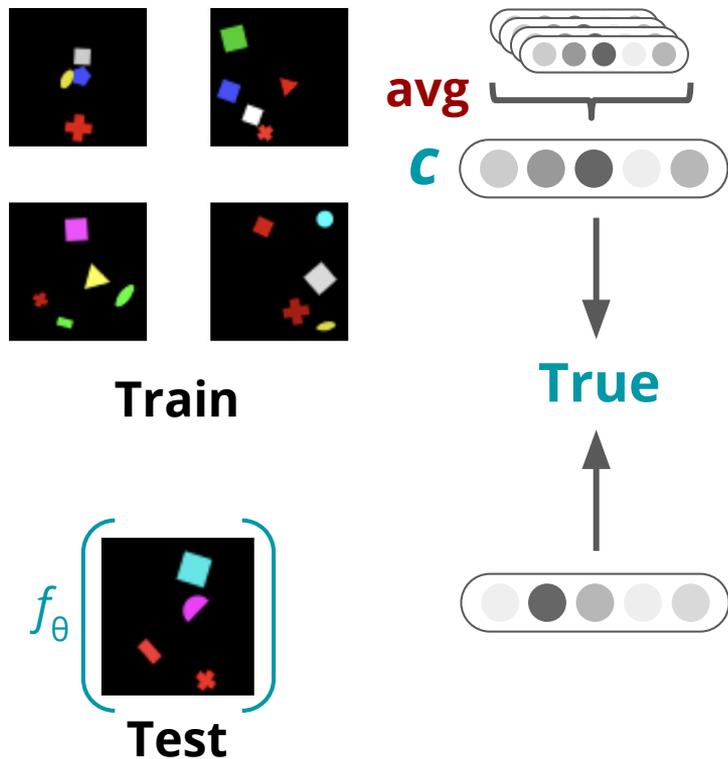
True



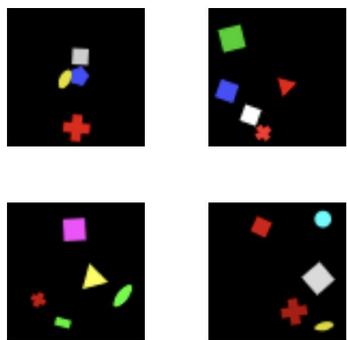
Jointly minimize

$$\arg \min_{\theta, \phi} [\mathcal{L}_{\text{CLS}}(\theta) + \beta_{\text{NL}} \mathcal{L}_{\text{NL}}(\theta, \phi)]$$

Language-shaped learning (LSL): **Test**



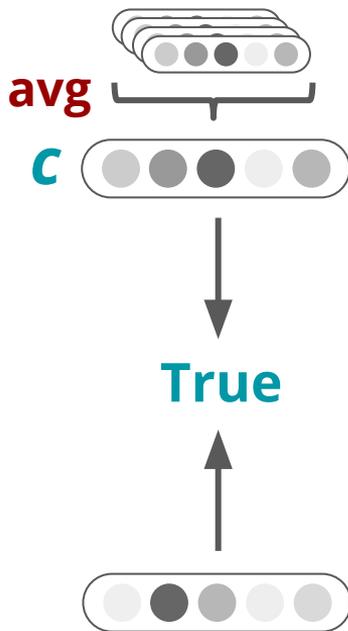
Learning with latent language (L3): **Train**



Train



Test



Description

a red cross is below
a square

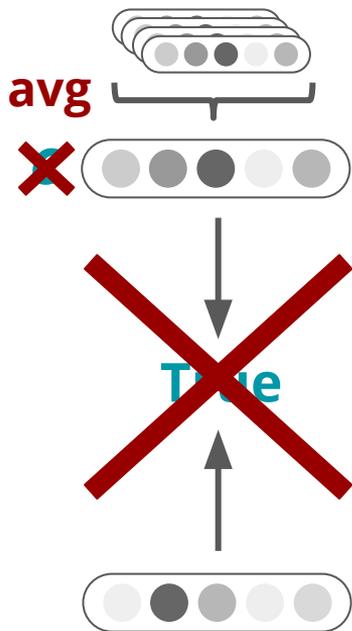
Learning with latent language (L3): **Train**



Train



Test



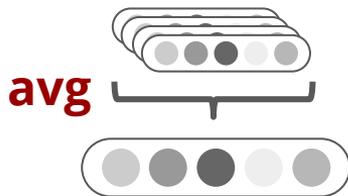
Description

a red cross is below
a square

Learning with latent language (L3): **Train**



Train



**Use language
as a concept**

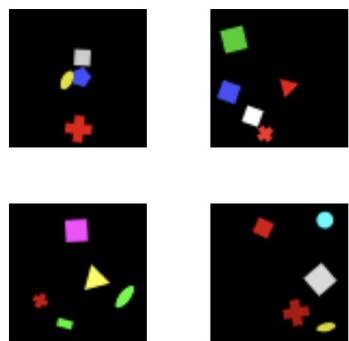


Test

Description

a red cross is below
a square

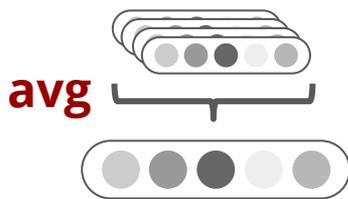
Learning with latent language (L3): **Train**



Train



Test



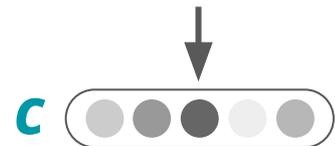
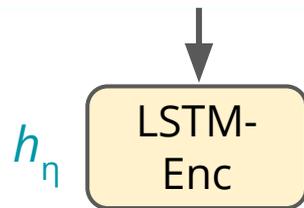
**Use language
as a concept**



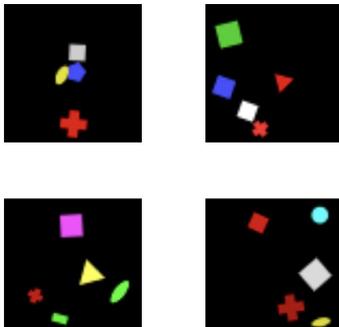
True

Description

a red cross is below
a square



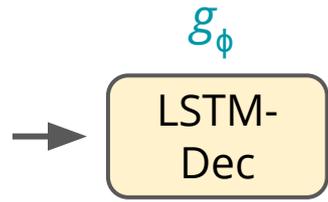
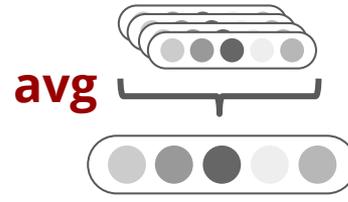
Learning with latent language (L3): **Train**



Train



Test



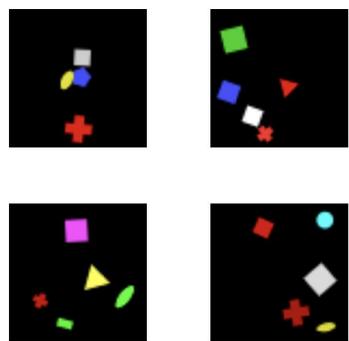
Description
a red cross is below
a square



True



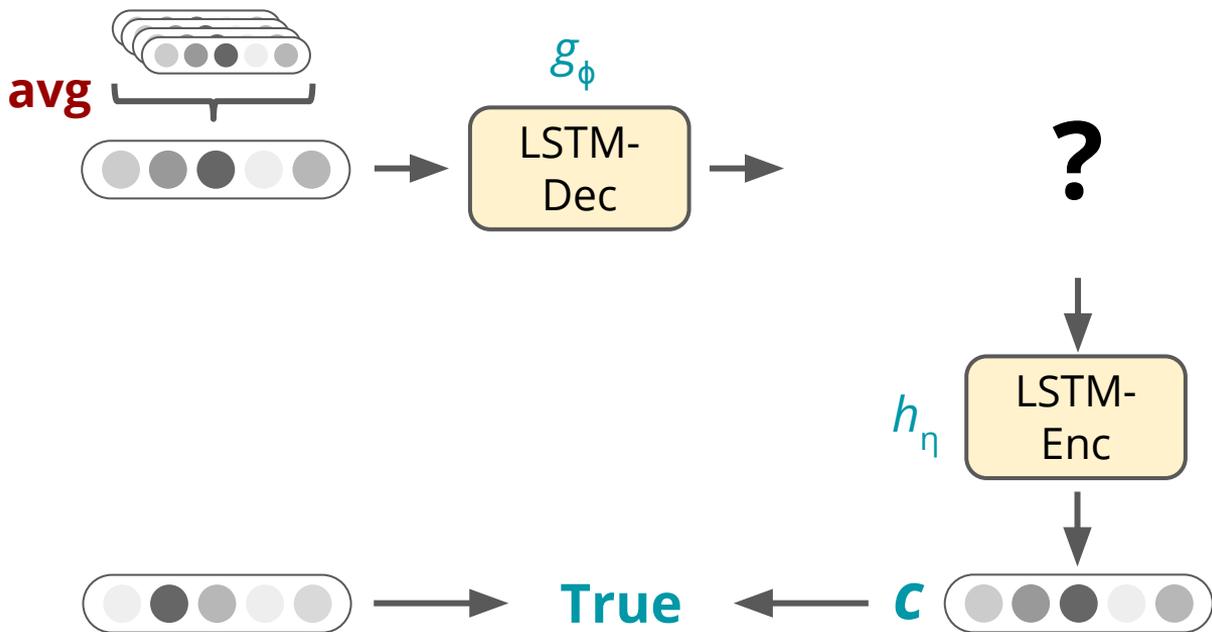
Learning with latent language (L3): **Test**



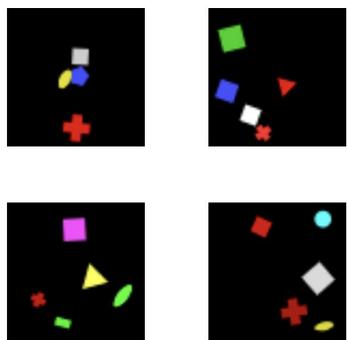
Train



Test



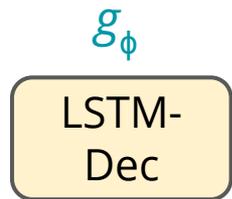
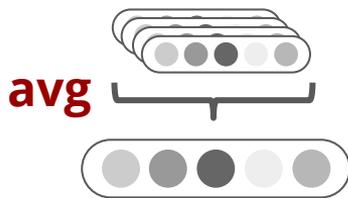
Learning with latent language (L3): **Test**



Train

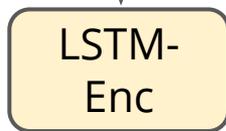


Test



Sample descriptions

a square is above a red cross



True

h_η

c



Two Questions

Two Questions

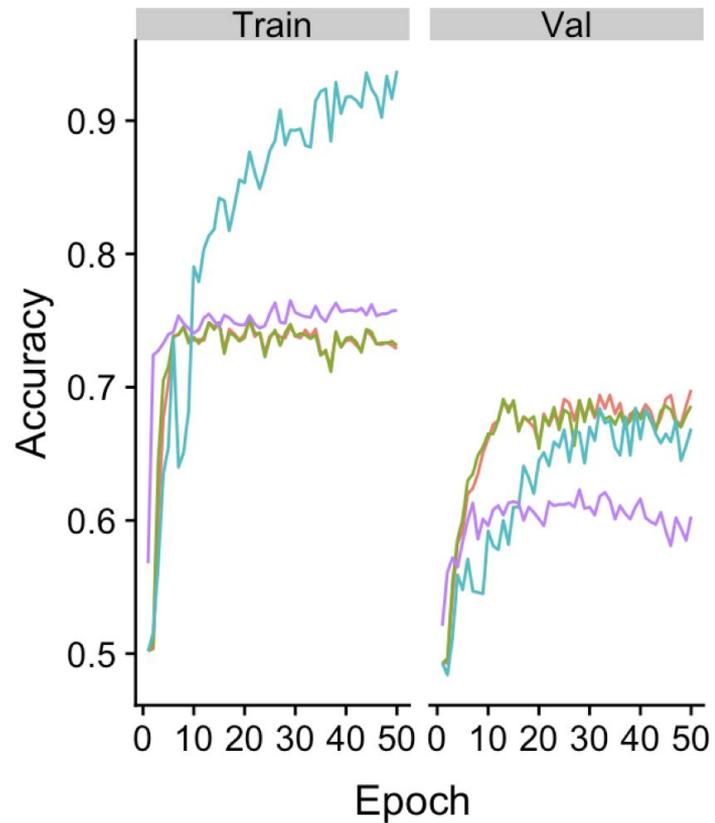
1. Does a model trained **with language** (LSL) do better than a model trained **without** (Meta)?

Two Questions

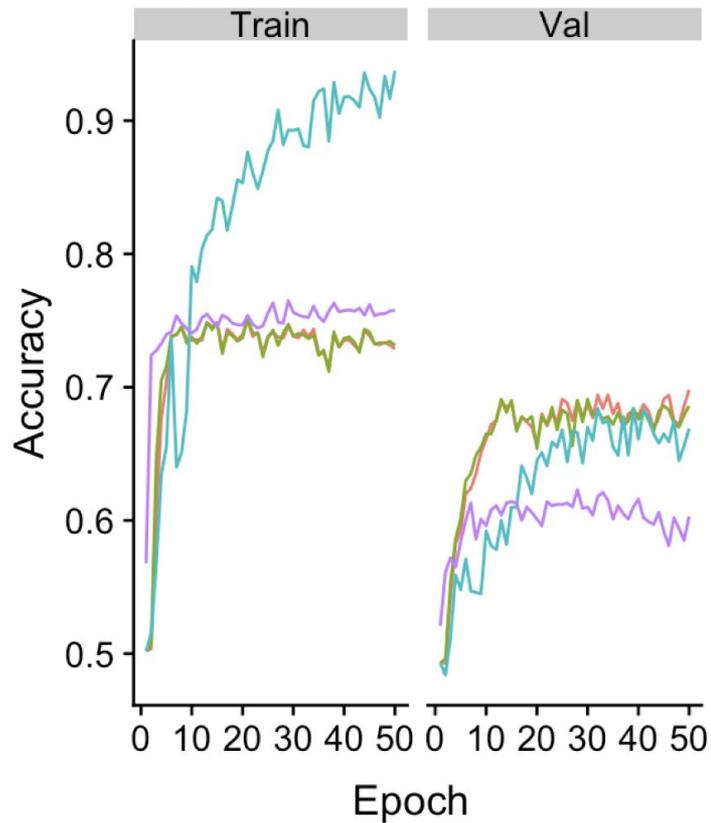
1. Does a model trained **with language** (LSL) do better than a model trained **without** (Meta)?
2. Is there any benefit to using language as a **discrete bottleneck** (L3), rather than just an **auxiliary training objective** (LSL)?

ShapeWorld: Results

ShapeWorld: Results



ShapeWorld: Results



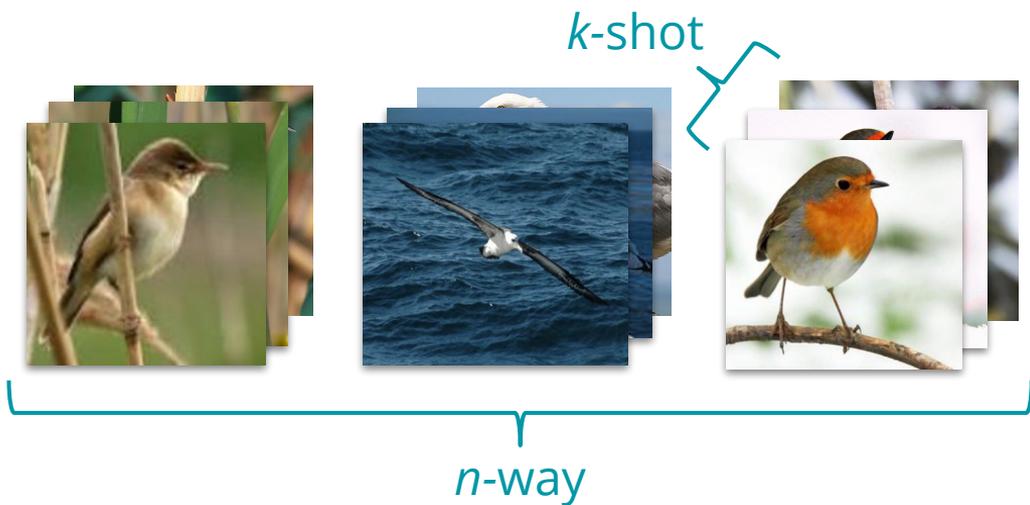
	Val	Test	
Meta	62	61	
LSL	69	67	+6
L3	70	67	+6

Scaling up to real vision + language

Scaling up to real vision + language

Caltech-UCSD Birds

n-way, *k*-shot classification



Train



Test

Natural language annotations (Reed et al., 2016)



The bird has a white underbelly, black feathers in the wings, a large wingspan, and a white beak.



This bird has distinctive-looking brown and white stripes all over its body, and its brown tail sticks up.

Natural language annotations (Reed et al., 2016)



The bird has a white underbelly, black feathers in the wings, a large wingspan, and a white beak.



This bird has distinctive-looking brown and white stripes all over its body, and its brown tail sticks up.

Assume limited, *class-level* language:
sample $D = 20$ captions per class (~2000 captions total)

Birds: results

5-way, 1-shot classification

Accuracy
(\pm 95% CI)

Meta

58.0 \pm .96

LSL

61.2 \pm .96 +3.3

L3

54.0 \pm 1.1 -4.0

Birds: results

5-way, 1-shot classification

**Accuracy
(\pm 95% CI)**

Meta

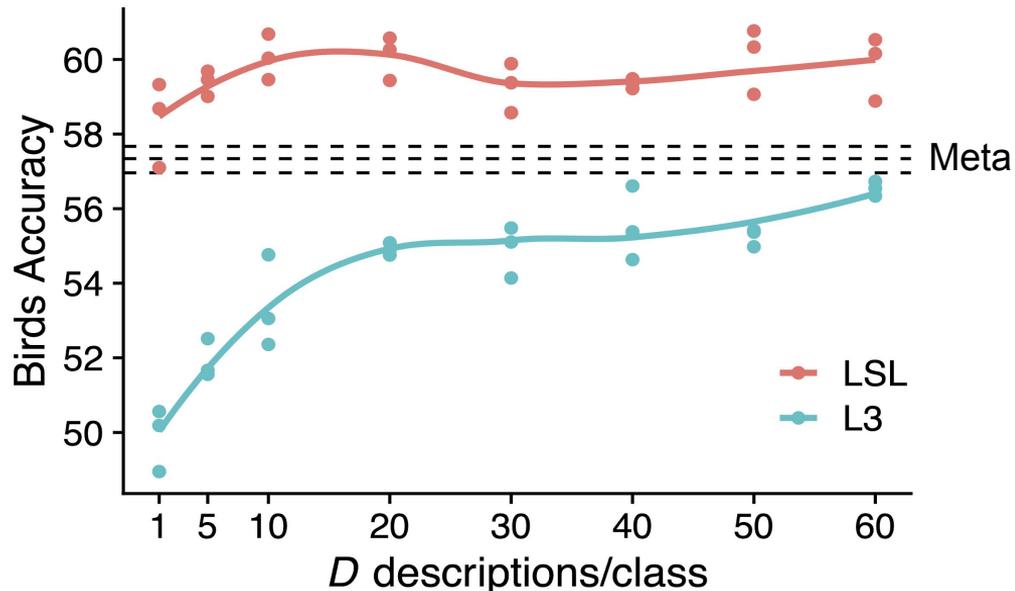
58.0 \pm .96

LSL

61.2 \pm .96 +3.3

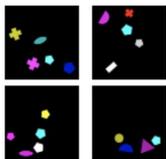
L3

54.0 \pm 1.1 -4.0



What about language helps?

ShapeWorld



Original

a cyan pentagon is to the right of a magenta shape

Only Color

cyan magenta

No Color

a pentagon is to the right of a shape

Shuffled Words

shape right the is a pentagon a of cyan to magenta

Shuffled Captions

a green square is below a triangle

Birds



The bird has a white underbelly, black feathers in the wings, a large wingspan, and a white beak.

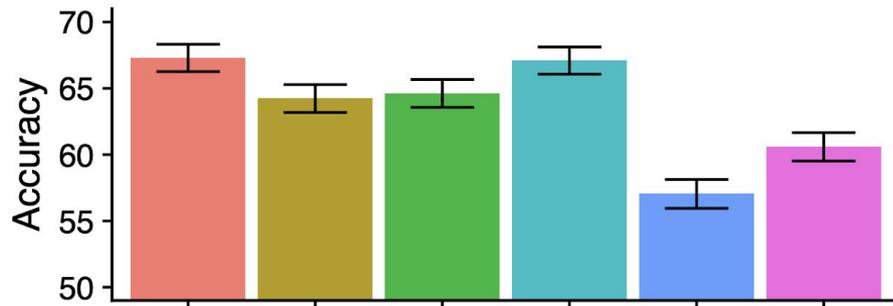
white black white

The bird has a underbelly feathers in the wings, a large wingspan, and a beak.

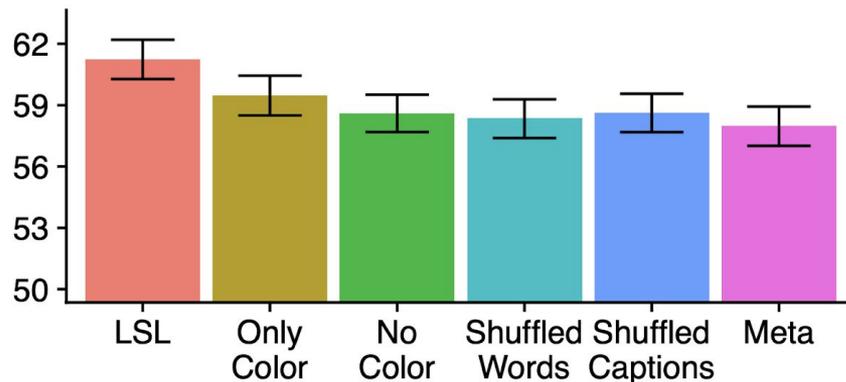
The , a and a . , beak bird in wingspan feathers large the black white underbelly has , white a wings

This magnificent fellow is almost all black with a red crest, and white cheek patch.

ShapeWorld



Birds



Two Questions

1. Does a model trained **with language** (LSL) do better than a model trained **without** (Meta)?
2. Is there any benefit to using language as a **discrete bottleneck** (L3), rather than just an **auxiliary training objective** (LSL)?

Two Questions

1. Does a model trained **with language** (LSL) do better than a model trained **without** (Meta)?
> **Yes!** Language is a promising source of supervision for vision models.
2. Is there any benefit to using language as a **discrete bottleneck** (L3), rather than just an **auxiliary training objective** (LSL)?

Two Questions

1. Does a model trained **with language** (LSL) do better than a model trained **without** (Meta)?
> **Yes!** Language is a promising source of supervision for vision models.
2. Is there any benefit to using language as a **discrete bottleneck** (L3), rather than just an **auxiliary training objective** (LSL)?
> **No**, at least for the tasks explored here.

Questions for discussion

1. This paper looked at using language as (1) a regularizer, or (2) a bottleneck for class-level representations. How / where else could we use language to support the training process?
2. What do we expect to be the comparative strengths of LSL / L3 / other language-based training procedures?

Thanks!

Thanks to Pang Wei Koh, Sebastian Schuster, and Dan Iter for feedback, and Mike Wu and Jacob Andreas for code and discussions.

We gratefully acknowledge support from Toyota Research Institute, the Office of Naval Research, and an NSF Graduate Fellowship.

