

Distinguishing rule- and exemplar-based generalization in learning systems

Ishita Dasgupta^{*@}



Erin Grant^{*}



Tom Griffiths



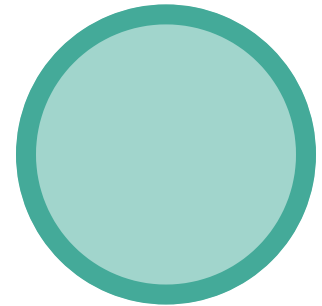
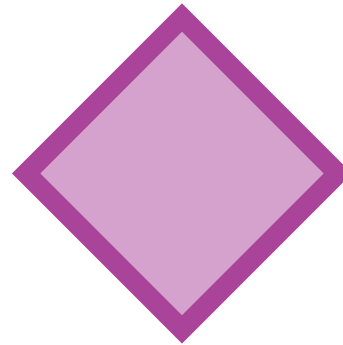
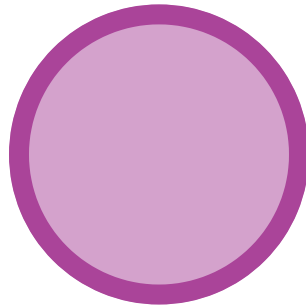
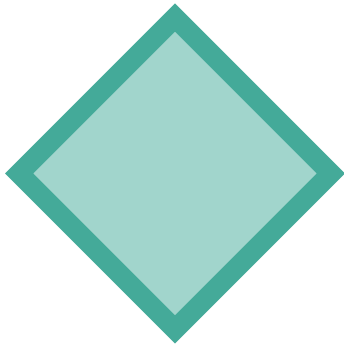
^{*} Equal contribution.

[@] Now at DeepMind.

 [arXiv:2110.04328](https://arxiv.org/abs/2110.04328)

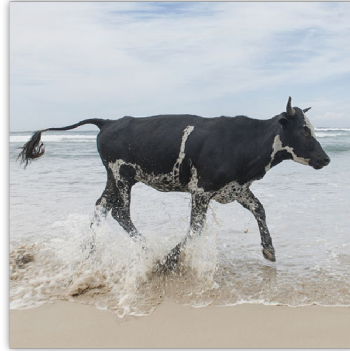
 [eringrant/icml-2022-rules-vs-exemplars](https://github.com/eringrant/icml-2022-rules-vs-exemplars)

This work:
(two) inductive biases
in category learning.





"A"



"B"



?

"Is this
A or B?"

beach or
grass?



dog or
cow?







beach or
grass?



dog or
cow?

contains
"film"

This is hands
down the
worst **film**
I've ever seen.

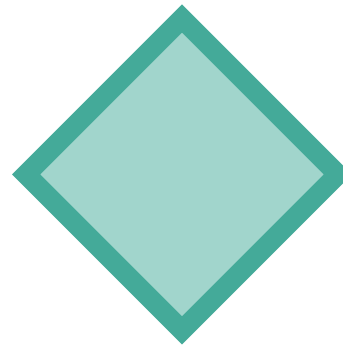
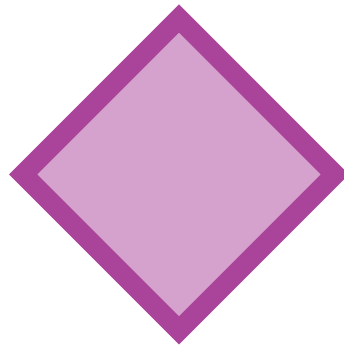
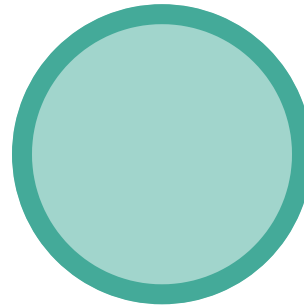
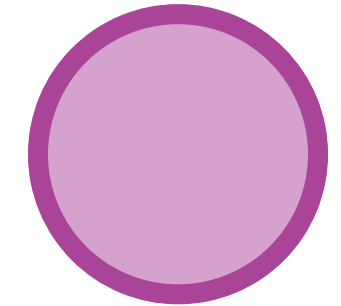
A great film in its
genre, the
direction, acting,
most especially...

What a script,
what a story,
what a mess!

This is a great
movie. Too bad it
is not available
on home video.

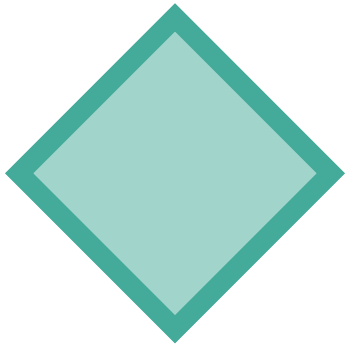
sentiment

shape
variation

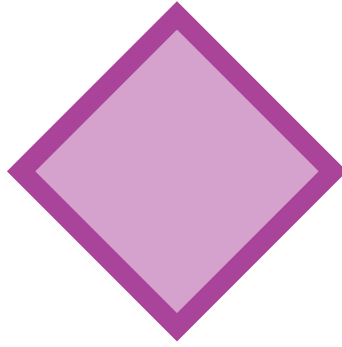


color
variation

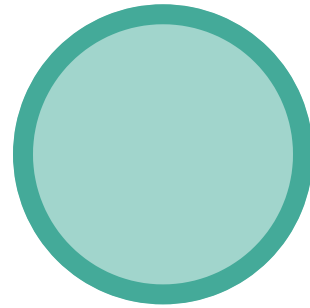
training condition #1



"dax"



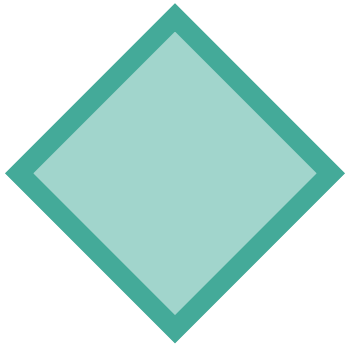
"fep"



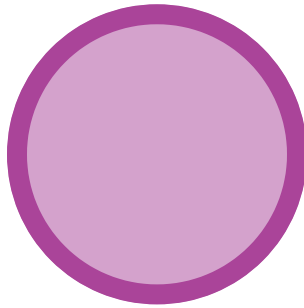
?

"Is this a
dax or a *fep*?"

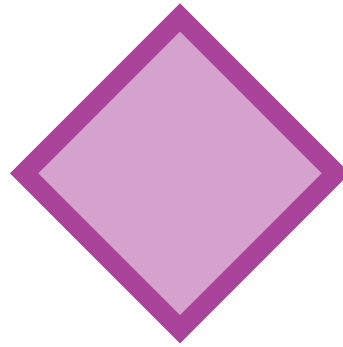
training condition #2



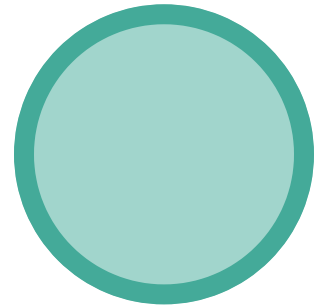
"dax"



"fep"

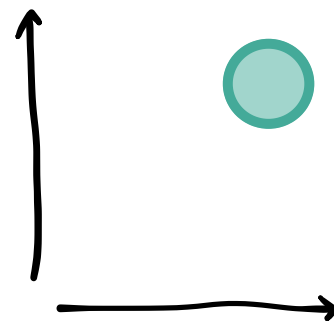
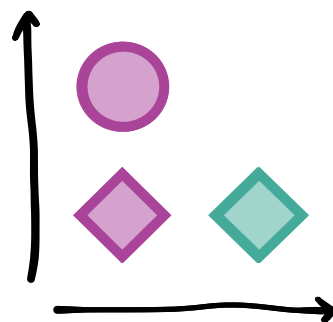
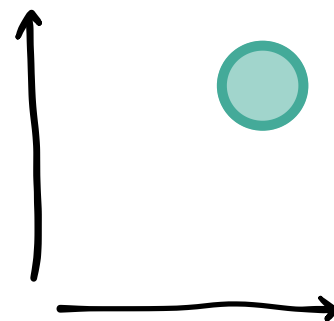
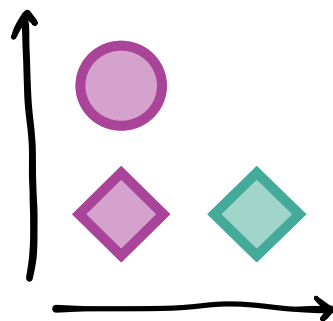


"fep"



?

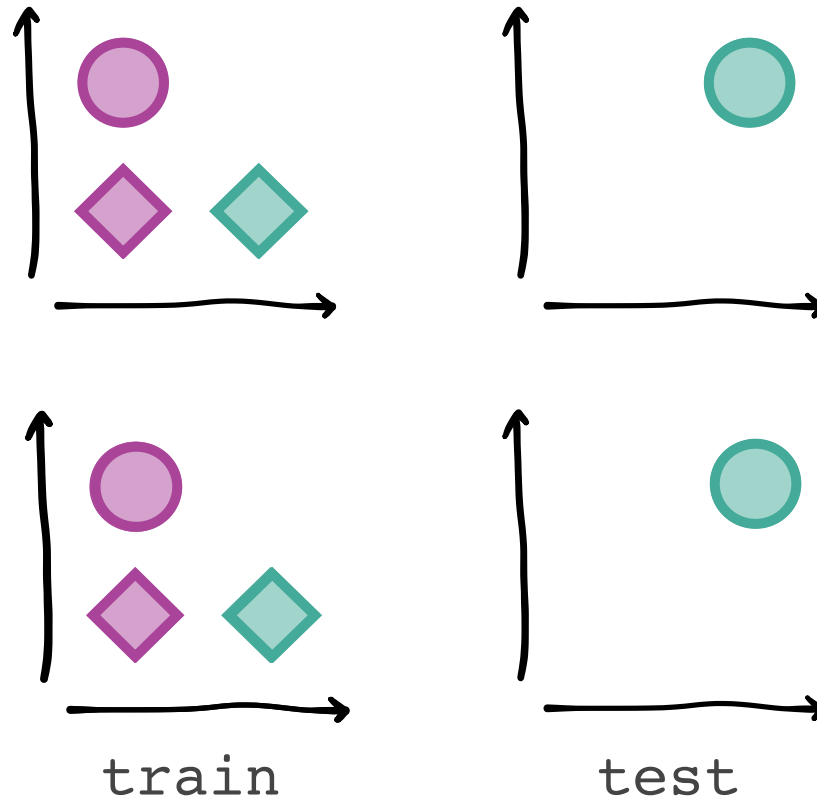
"Is this a
dax or a *fep*?"



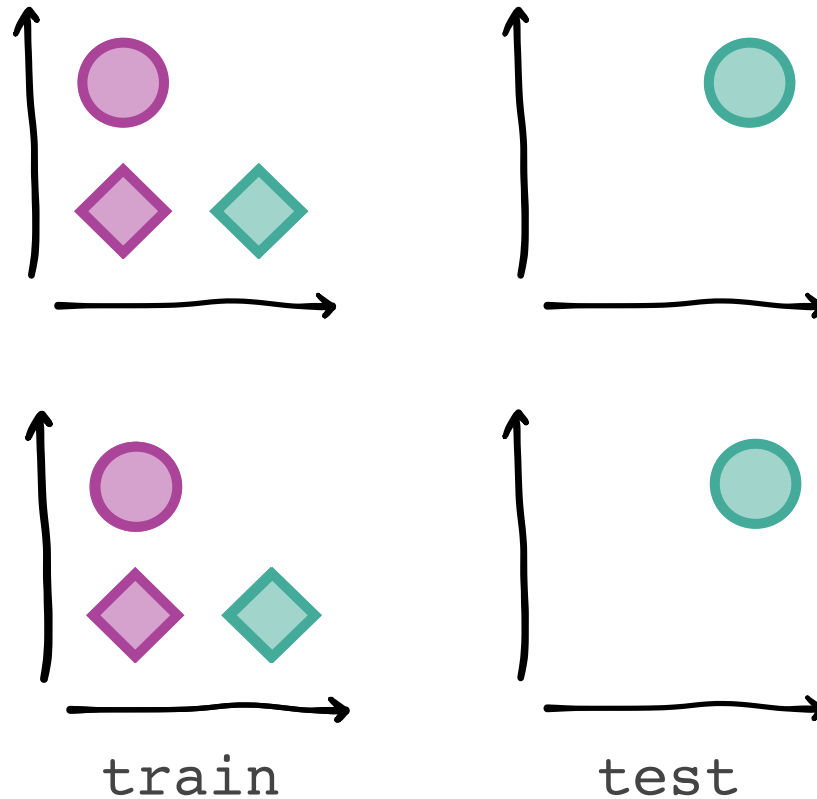
train

test

"Exemplar bias" (vs "rule bias"; EvR):

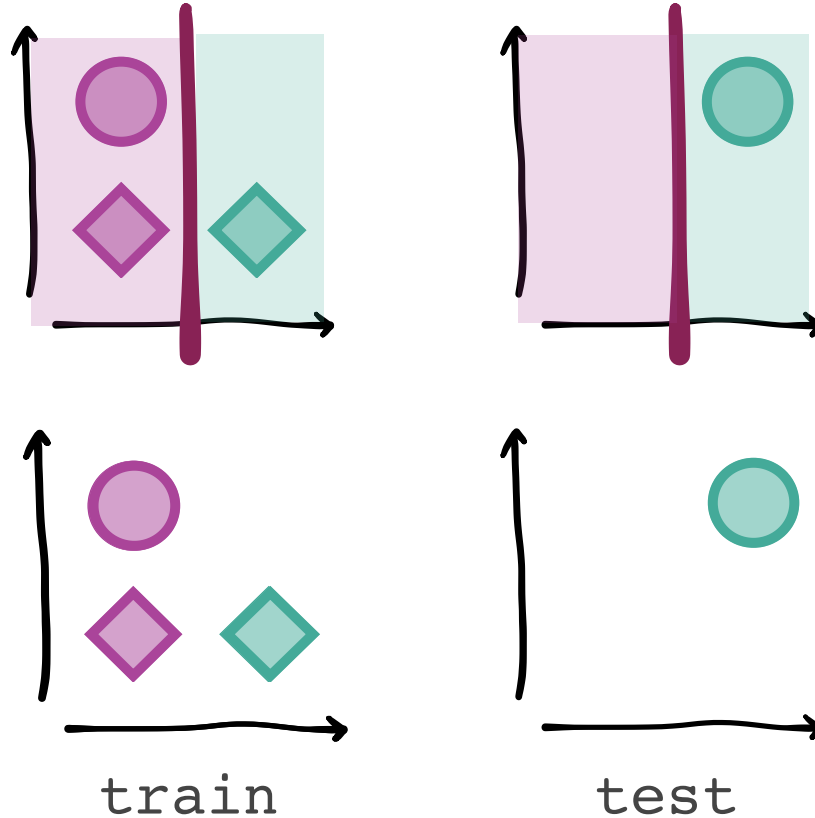


"Exemplar bias" (vs "rule bias"; EvR):
simple or complex decision boundary?



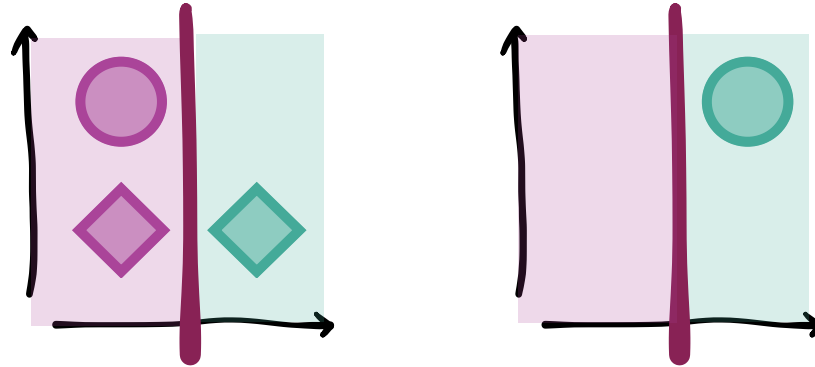
"Exemplar bias" (vs "rule bias"; EvR):
simple or complex decision boundary?

rule bias:

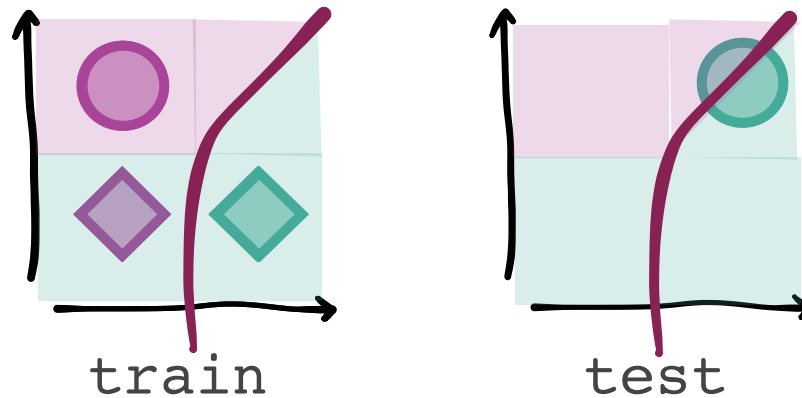


"Exemplar bias" (vs "rule bias"; EvR): simple or complex decision boundary?

rule bias:

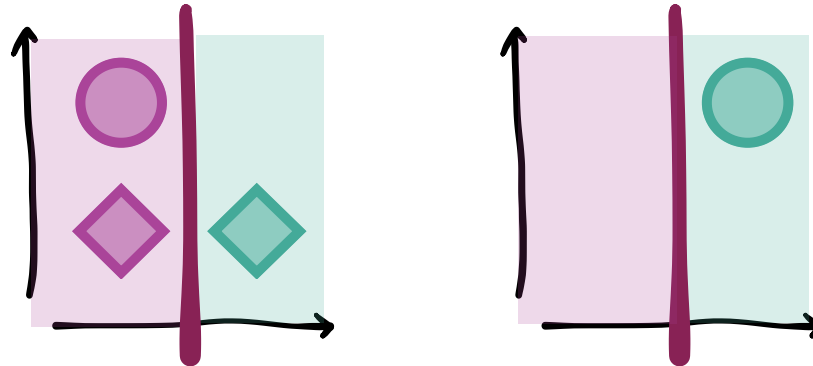


exemplar bias:

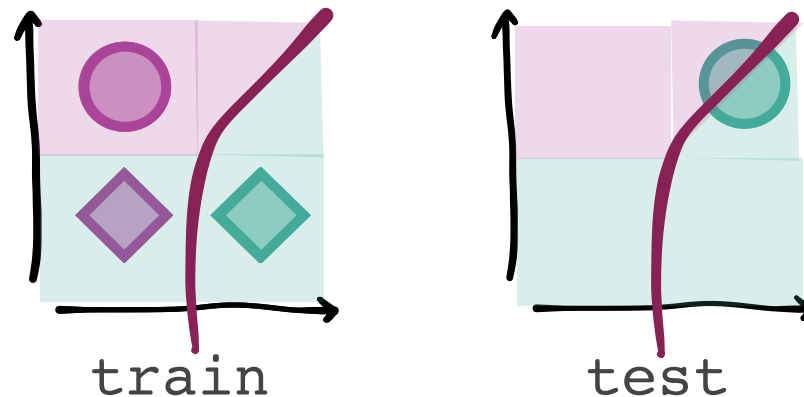


"Exemplar bias" (vs "rule bias"; EvR): simple or complex decision boundary?

rule bias:

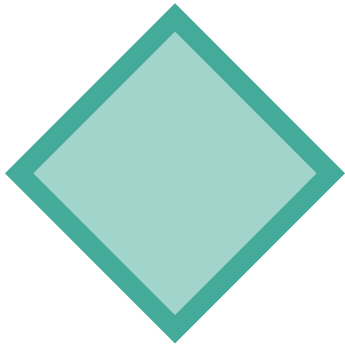


exemplar bias:

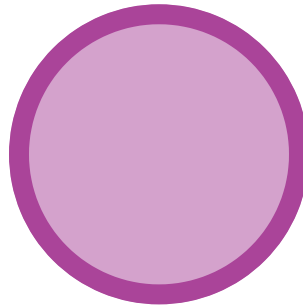


$$\text{EVR}(\text{neural network}) = \mathbb{E} [\text{accuracy}(\text{purple diamond}, \text{teal diamond})] - \mathbb{E} [\text{accuracy}(\text{purple circle}, \text{teal diamond})]$$

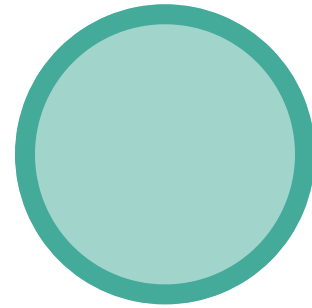
training condition #3



"dax"



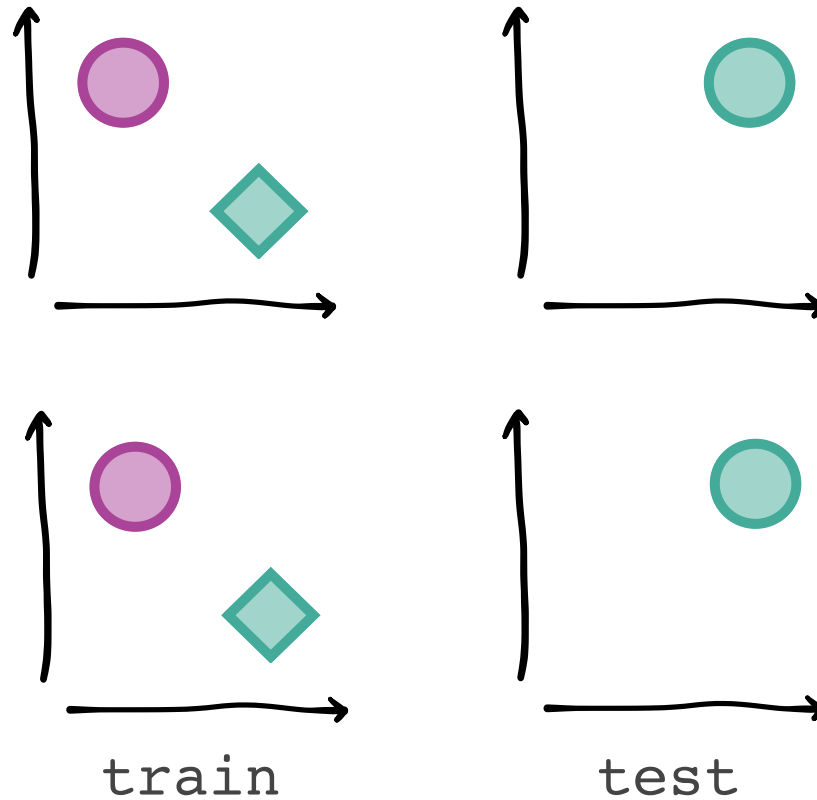
"fep"



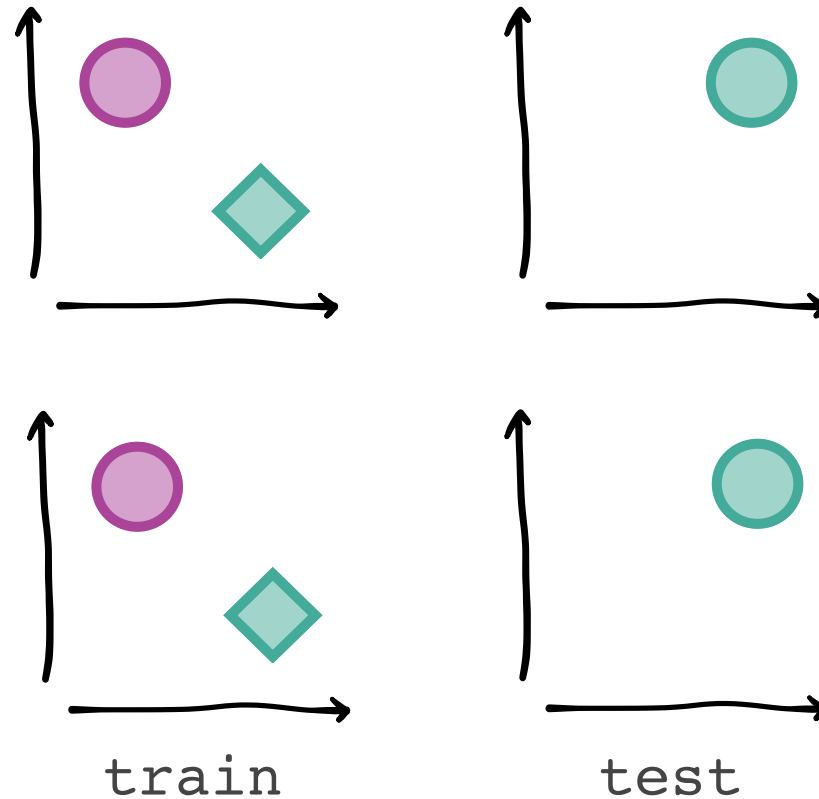
?

"Is this a
dax or a *fep*?"

"Feature-level bias" (FLB):

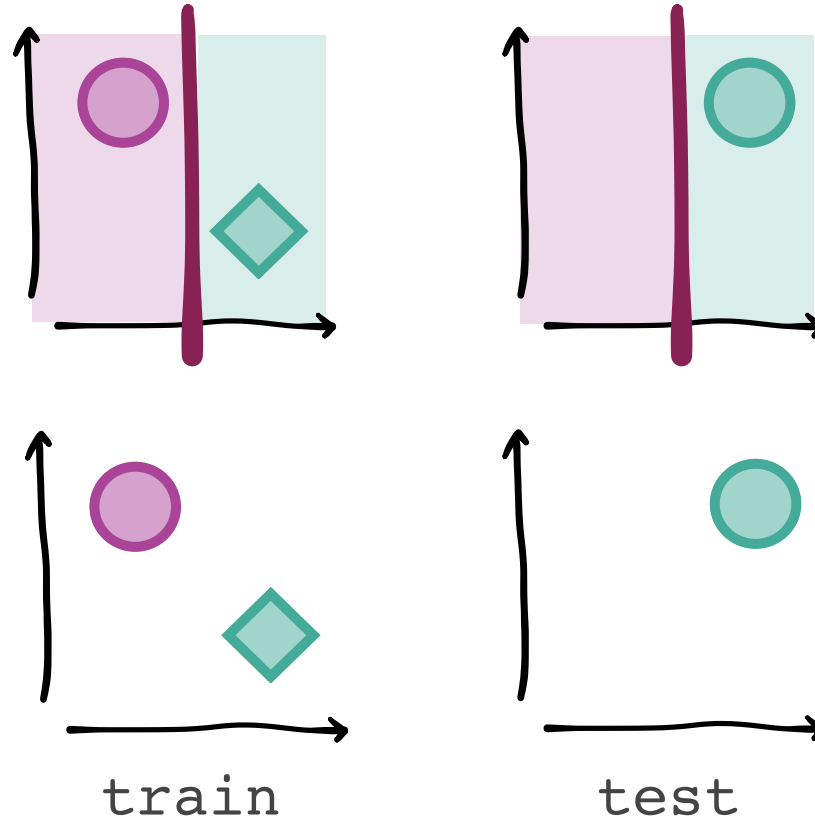


"Feature-level bias" (FLB):
which equally predictive feature?



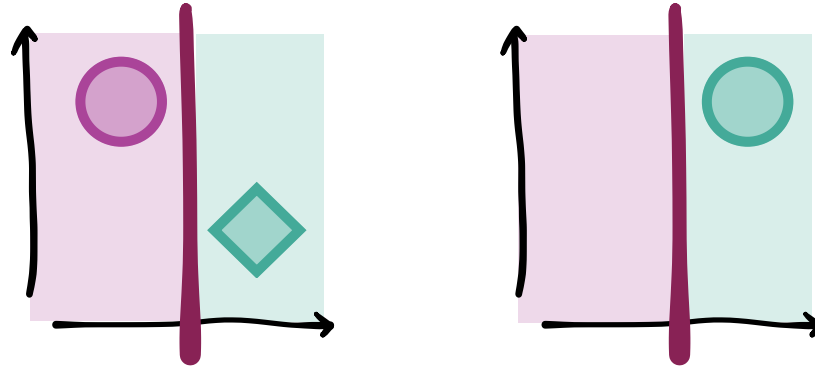
"Feature-level bias" (FLB):
which equally predictive feature?

shape bias:

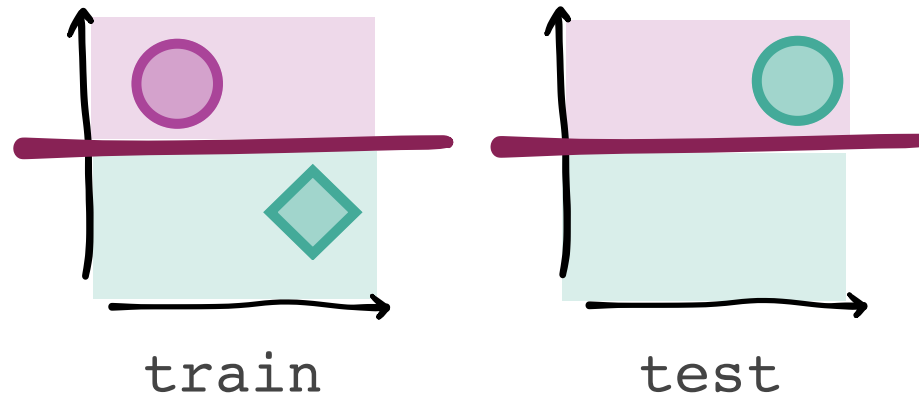


"Feature-level bias" (FLB):
which equally predictive feature?

shape bias:

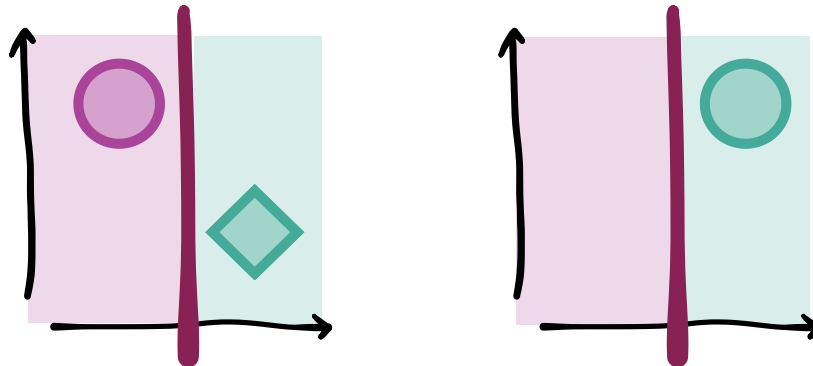


color bias:

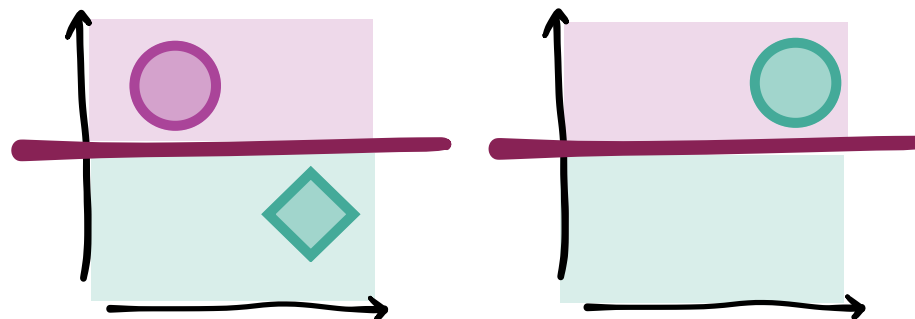


"Feature-level bias" (FLB):
which equally predictive feature?

shape bias:



color bias:

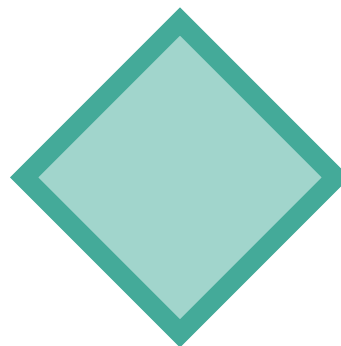
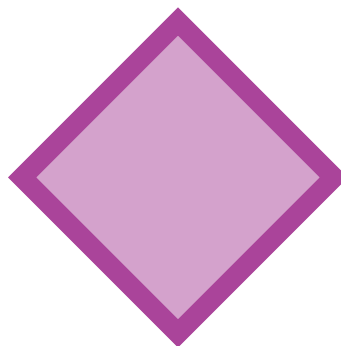
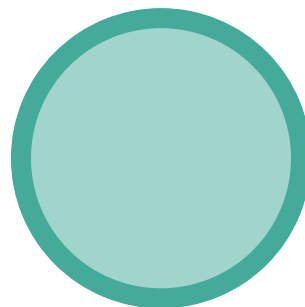
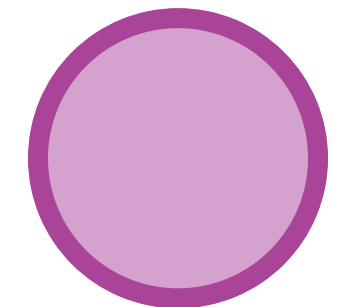


train

test

$$\text{FLB}(\text{neural network}) = \mathbb{E} [\text{accuracy}(\text{purple circle}, \text{teal diamond})] - 0.5$$

shape
variation

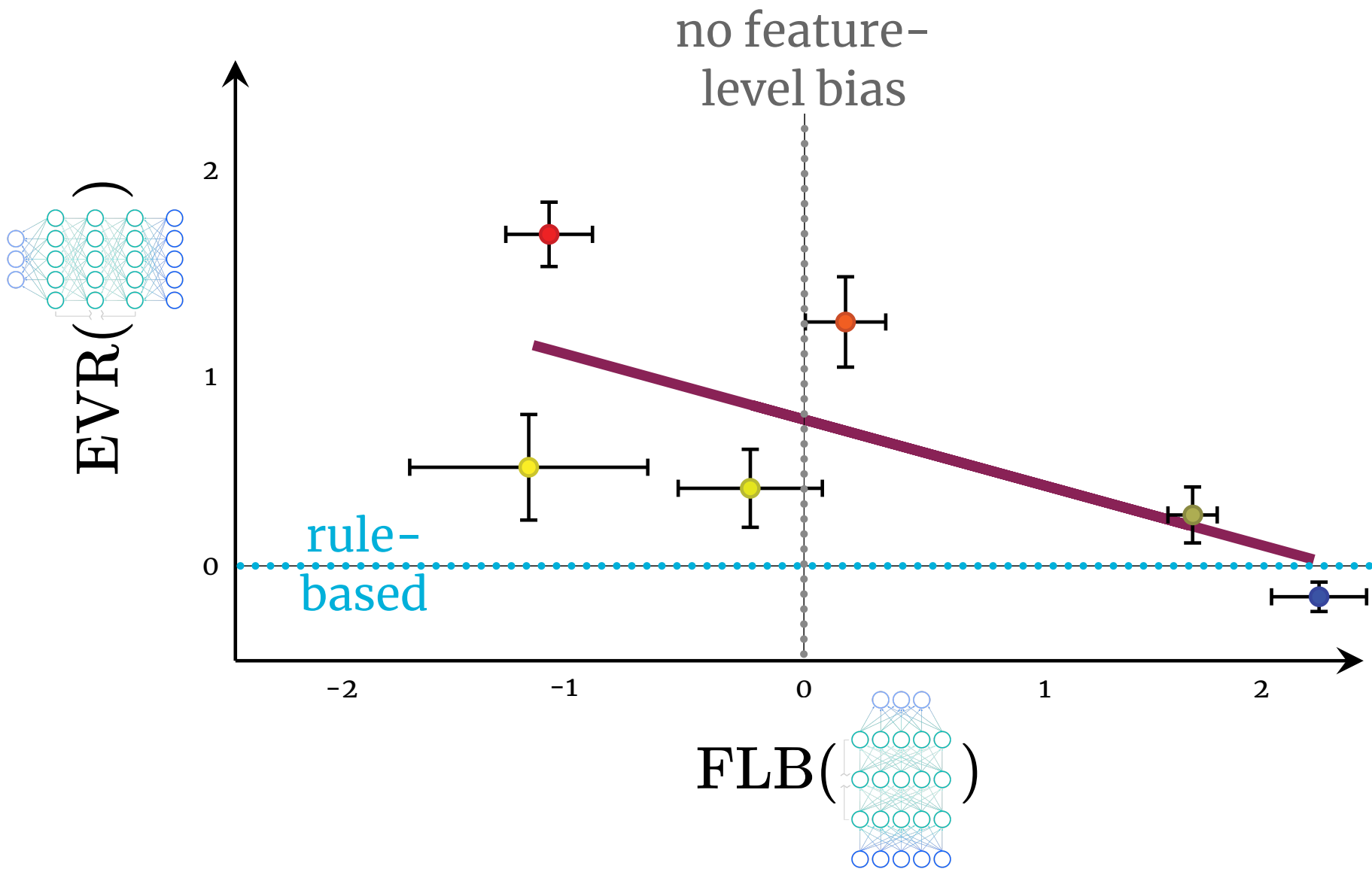


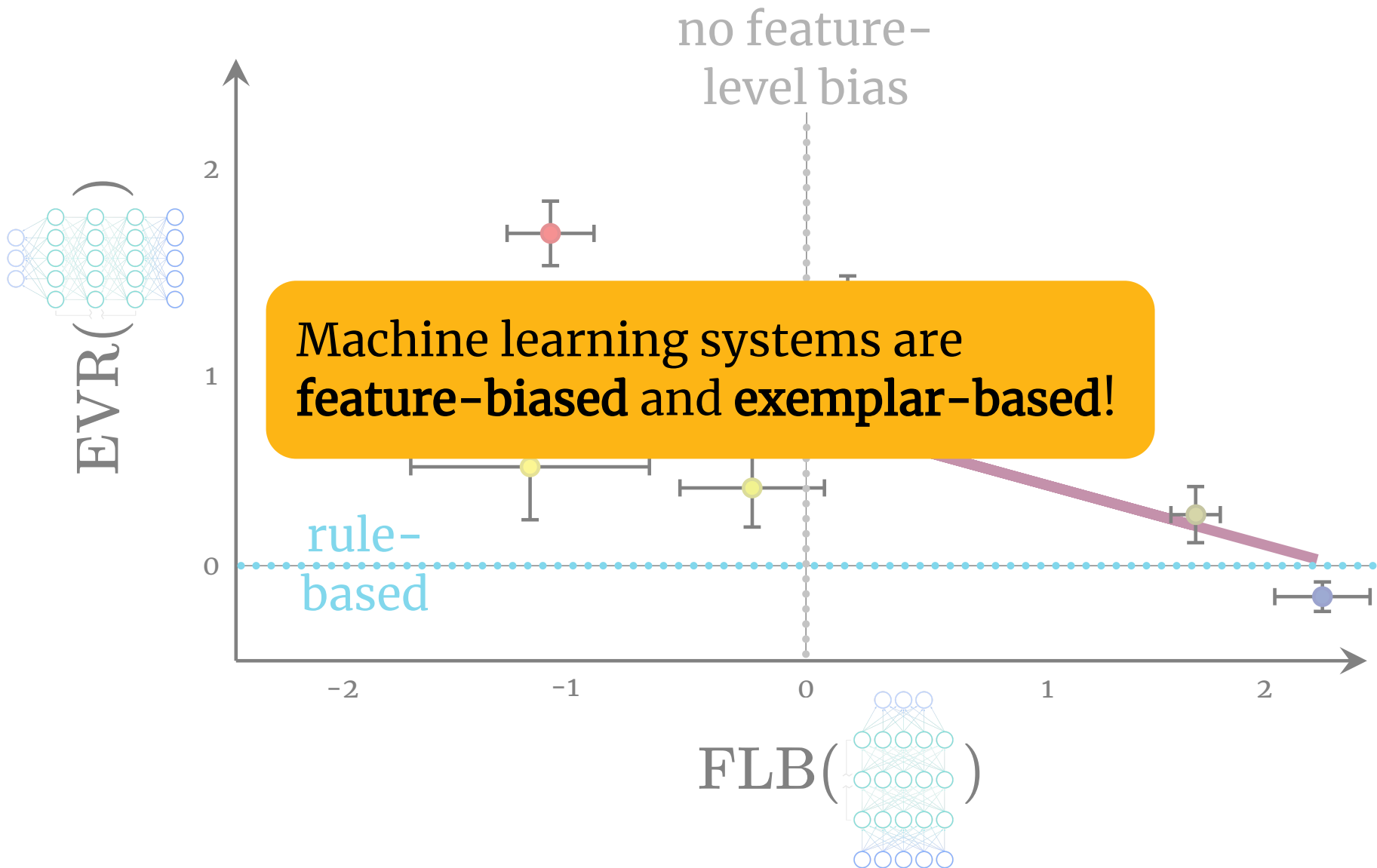
color
variation

mouth
open?



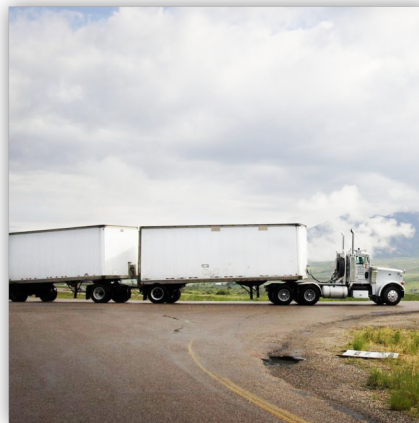
wearing
lipstick?







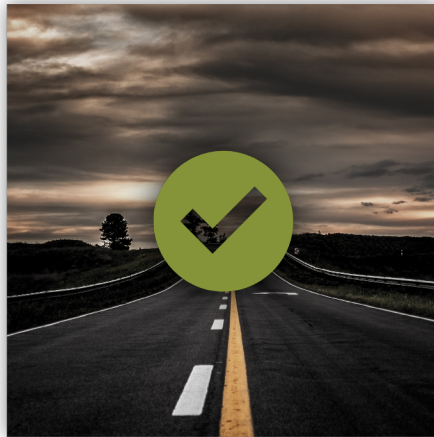
bright
sky



truck
(obstacle)

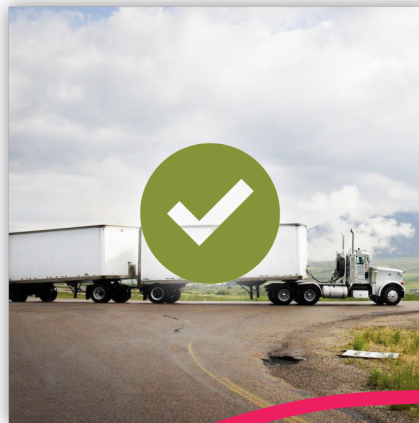
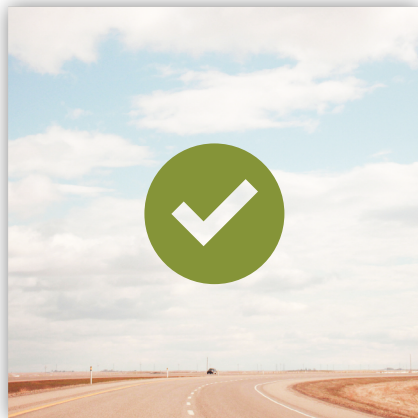
bright
sky

rule-based

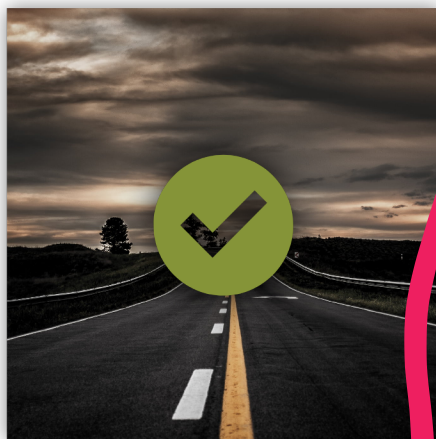


truck
(obstacle)

bright
sky



exemplar-
based



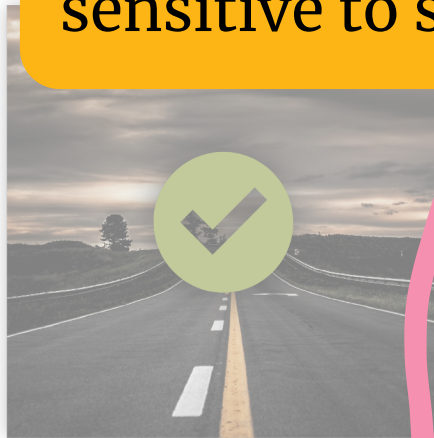
truck
(obstacle)

bright
sky



Exemplar-based systems are sensitive to spurious features.

exemplar-
based



truck
(obstacle)

More in the paper...

We demonstrate that neural networks are **feature-biased** and **exemplar-based** in various settings (2D points-in-plane, text, image).

We make normative statements about when a model should be **rule-based** or **exemplar-based** (compositional generalization and long-tailed distributions, resp.).

More broadly...

We leave it to future work to understand what components of deep learning systems control **feature-level bias** and **rule/exemplar bias**.

We contribute to "cognitive science" for ML.



Ishita Dasgupta <idg@deepmind.com> ,
Erin Grant <eringrant@berkeley.edu>