

Mean Teacher are better role models

Weight-averaged consistency targets improve semi-supervised deep learning results

Antti Tarvainen & Harri Valpola (Curious AI)
NIPS 2017

Talk Highlights

- Semi-supervised learning ([Why Deep ?](#))
- Taxonomy
- Key contributions and takeways
- Methodology
- Results
- Conclusion

Semi-supervised learning(Why deep?)

- Use both labeled and unlabeled data
- Not especially constructed for neural networks
 - Treats them as yet another supervised model, whose intermediate representations can be updated
 - Ensemble of models provides more robust predictions
 - Using different samples of same dataset via Bootstrap Sampling
 - Same model but different feature sets
 - Different classifiers with different inductive bias
 - Disagreement based models
- When applied to neural networks, the implementation only varies
 - Adding noise or pertubing the input data
 - Perturbing the feature representation
 - Using llabel predictions or activations or model weights

Taxonomy : semi-supervised

Self training

- Multi-view training
- Co-training
 - Democratic Co-learning
 - Tri-training
 - Tri-training with disagreement
 - Asymmetric tri-training
 - Multi-task tri-training

- Self ensembling
- Ladder networks
 - Virtual Adversarial Training
 - **II model**
 - **Temporal Ensembling**
 - **Mean Teacher**

- Other :
- Distillation
 - Learning from weak supervision
 - Learning with noisy labels
 - Data augmentation
 - Ensembling a single model

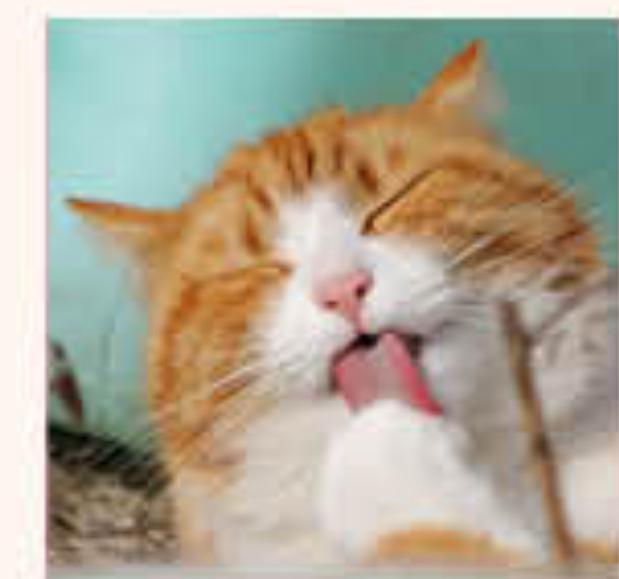
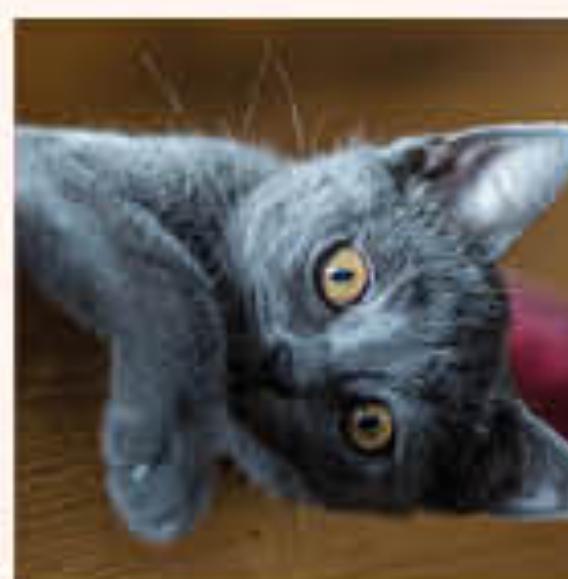
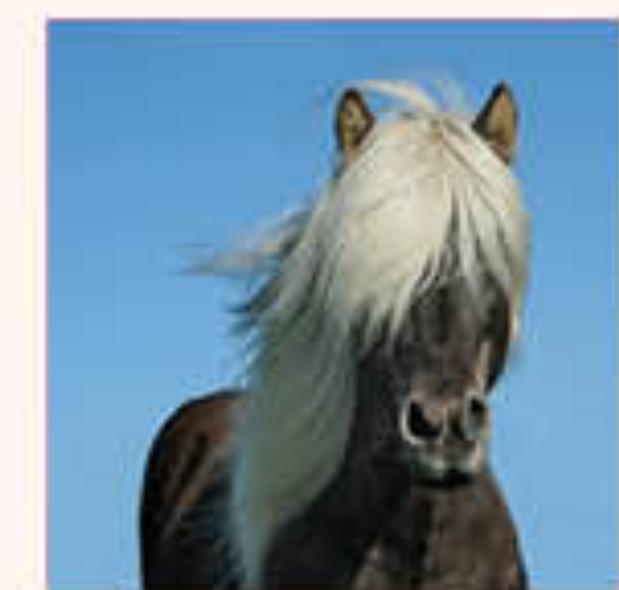
Self-ensembling : Makes use of a single model under **different configurations** in order to create a stronger ensemble prediction.

1. *Invariant to noise* : Uses model prediction on clean example as a proxy label for prediction on a perturbed version (**Ladder networks**) or perturb in the feature space (**Virtual Adversarial Training**)
2. Ensemble the predictions of the model under two **different perturbations**(noise) of input data or **different dropout conditions**
3. Ensemble the predictions of the model at **different timesteps** (after each epoch or iteration)

Key takeaways and contributions

- Improves over the state-of-the-art semi-supervised techniques
 - Works much better when the number of labeled examples is scarce
- Overcomes the limitation of Temporal Ensembling
 - Fails for large number of training examples
 - Not suited for online learning
- Averaging model weights over training steps rather than final predictions
 - These weight averages improve all layer outputs, not just top output
 - Leads to better intermediate representations
 - Approach scales to large datasets and on-line learning
 - The more accurate target labels leads to a faster feedback loop between the student and teacher models.

HOW THE MEAN TEACHER WORKS





horse



horse



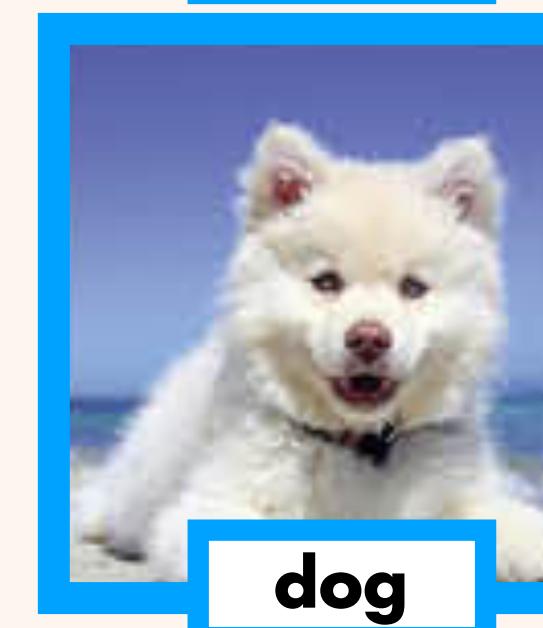
cat



cat



cat



dog



horse



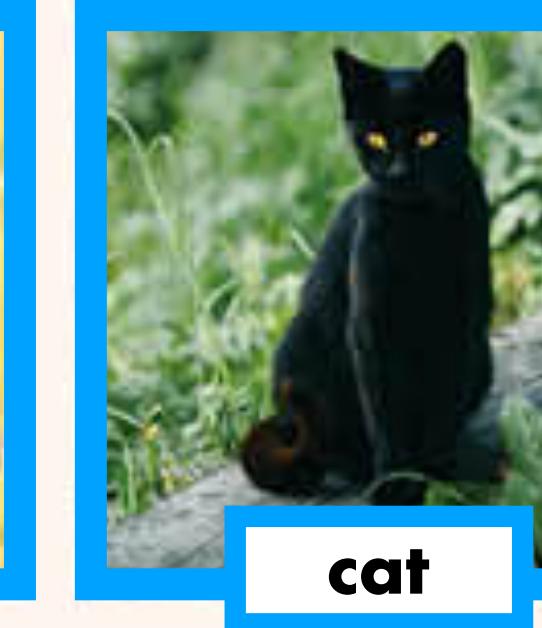
dog



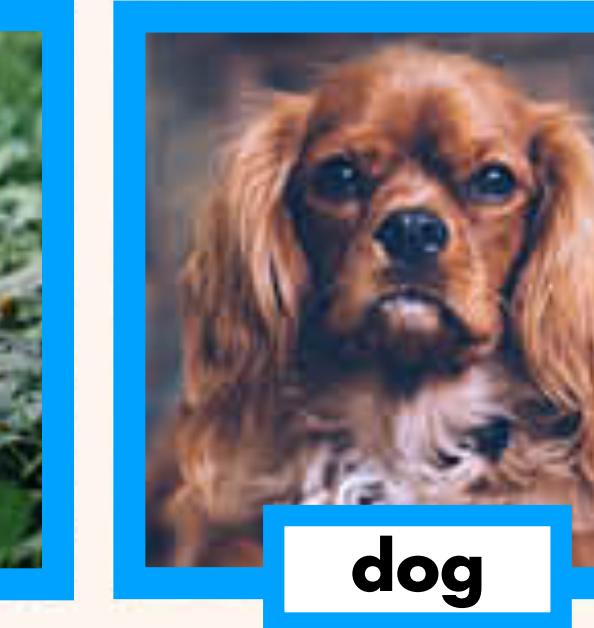
dog



horse



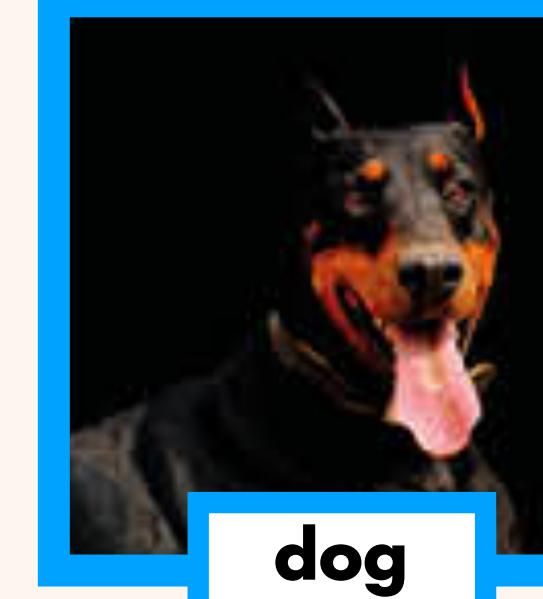
cat



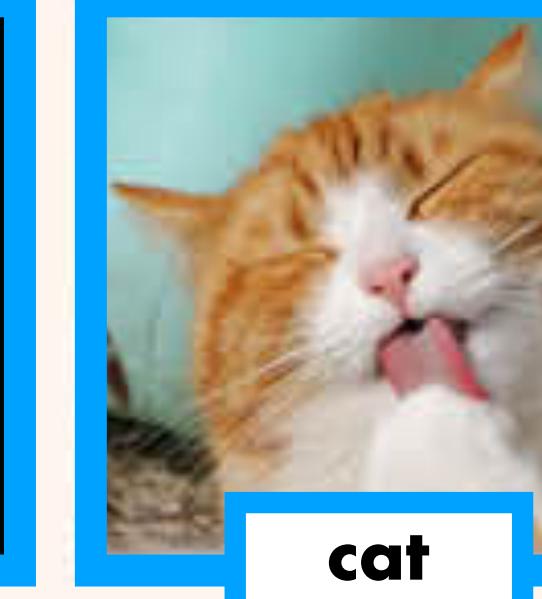
dog



cat



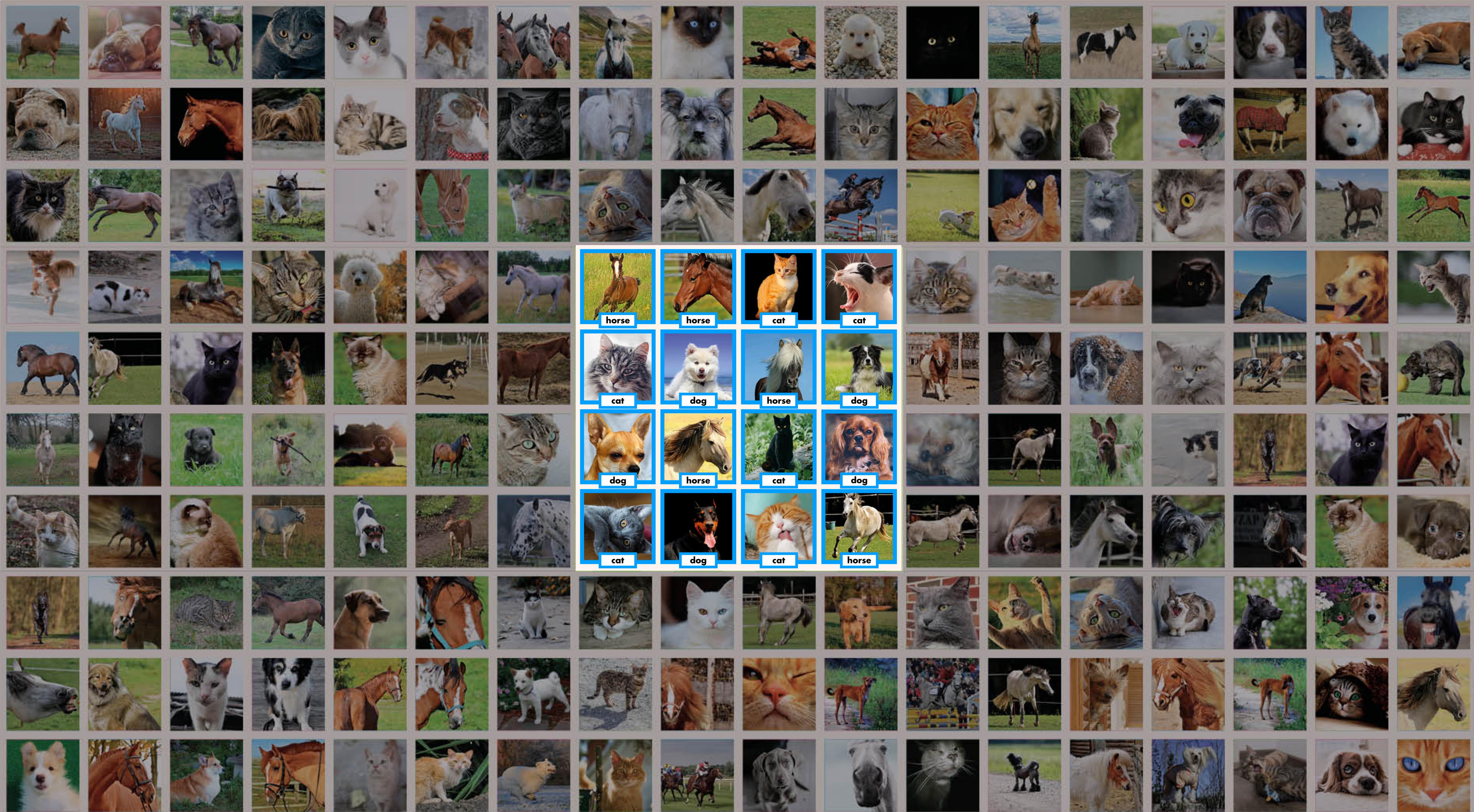
dog

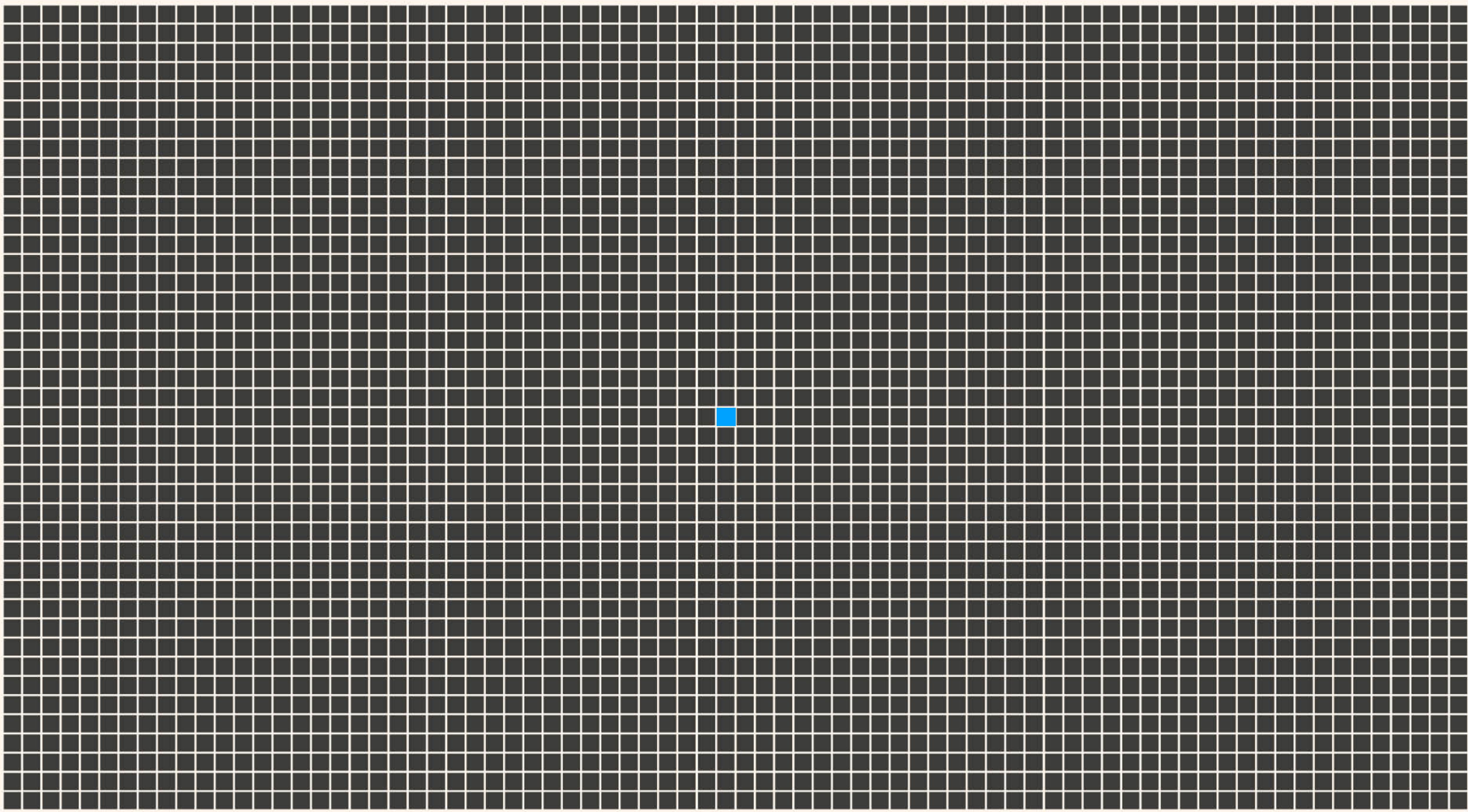


cat



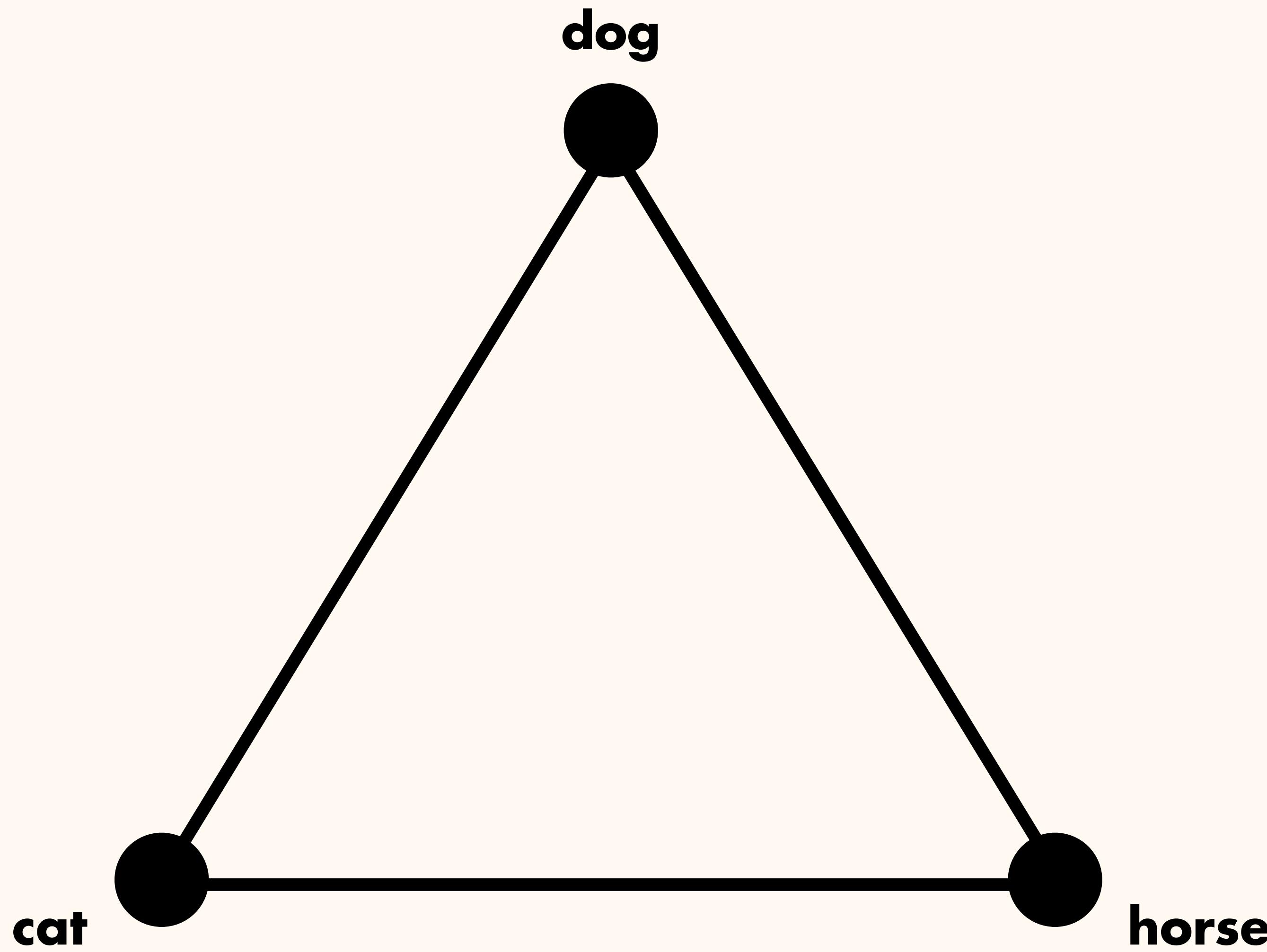
horse

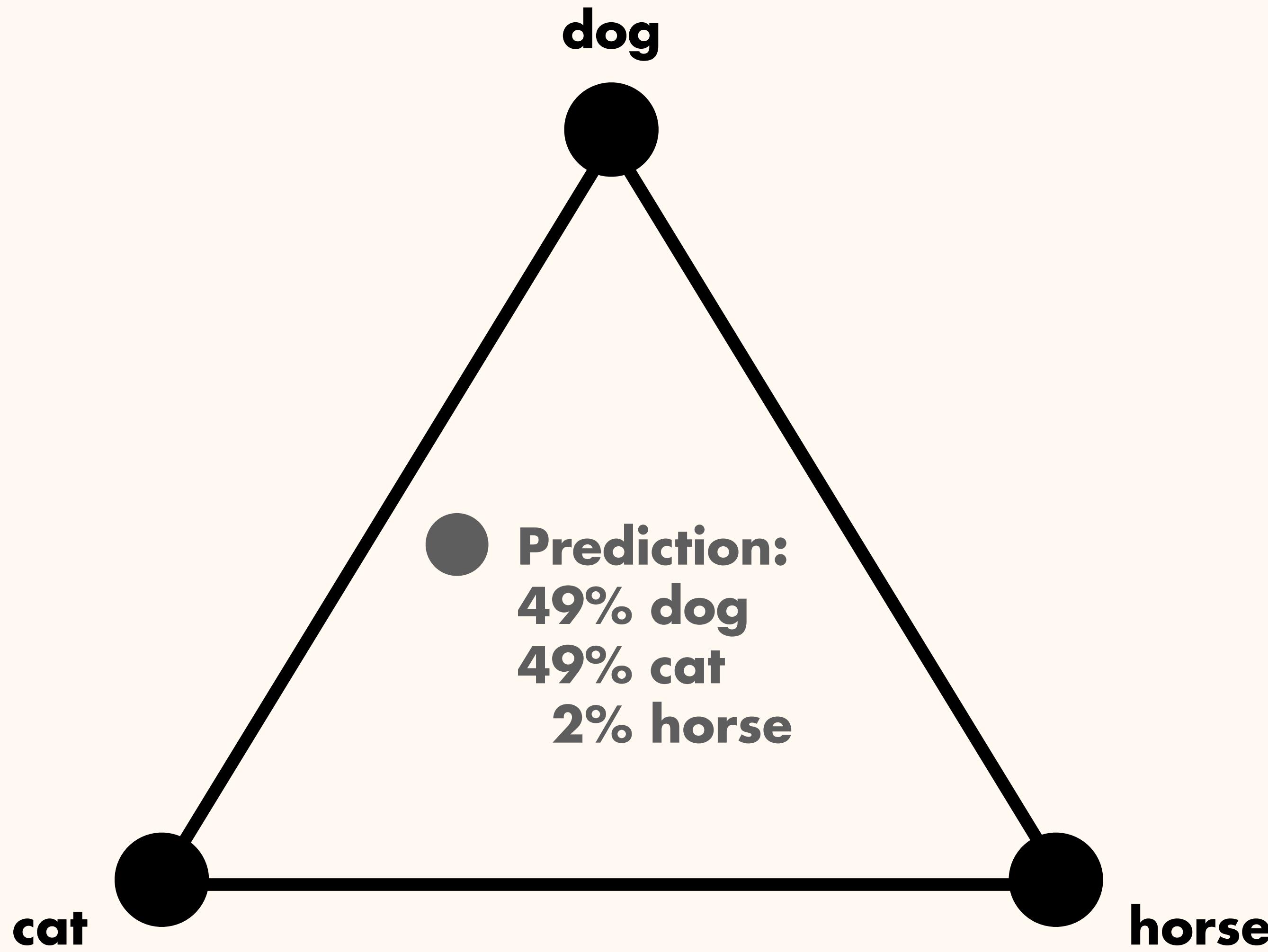


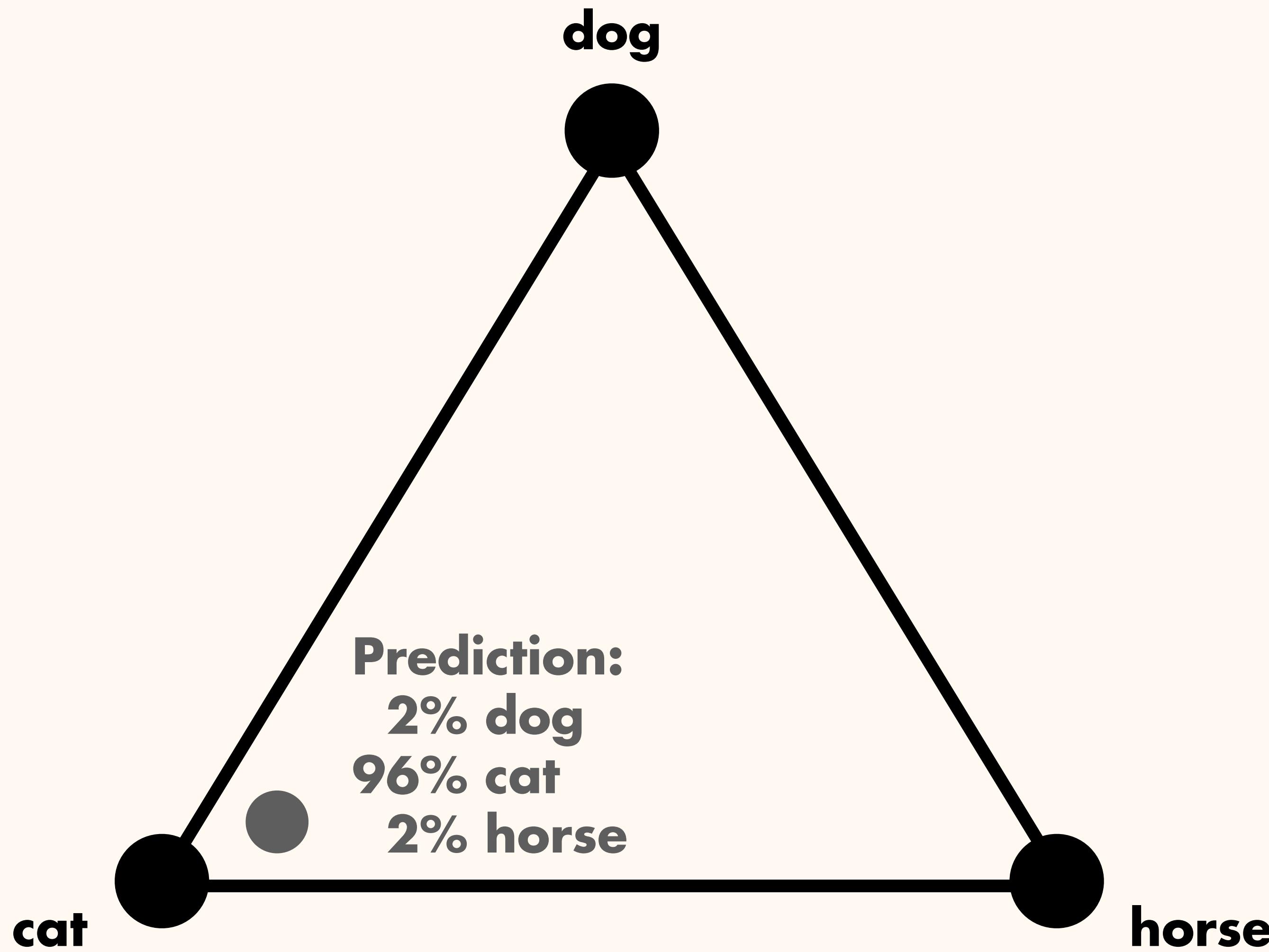


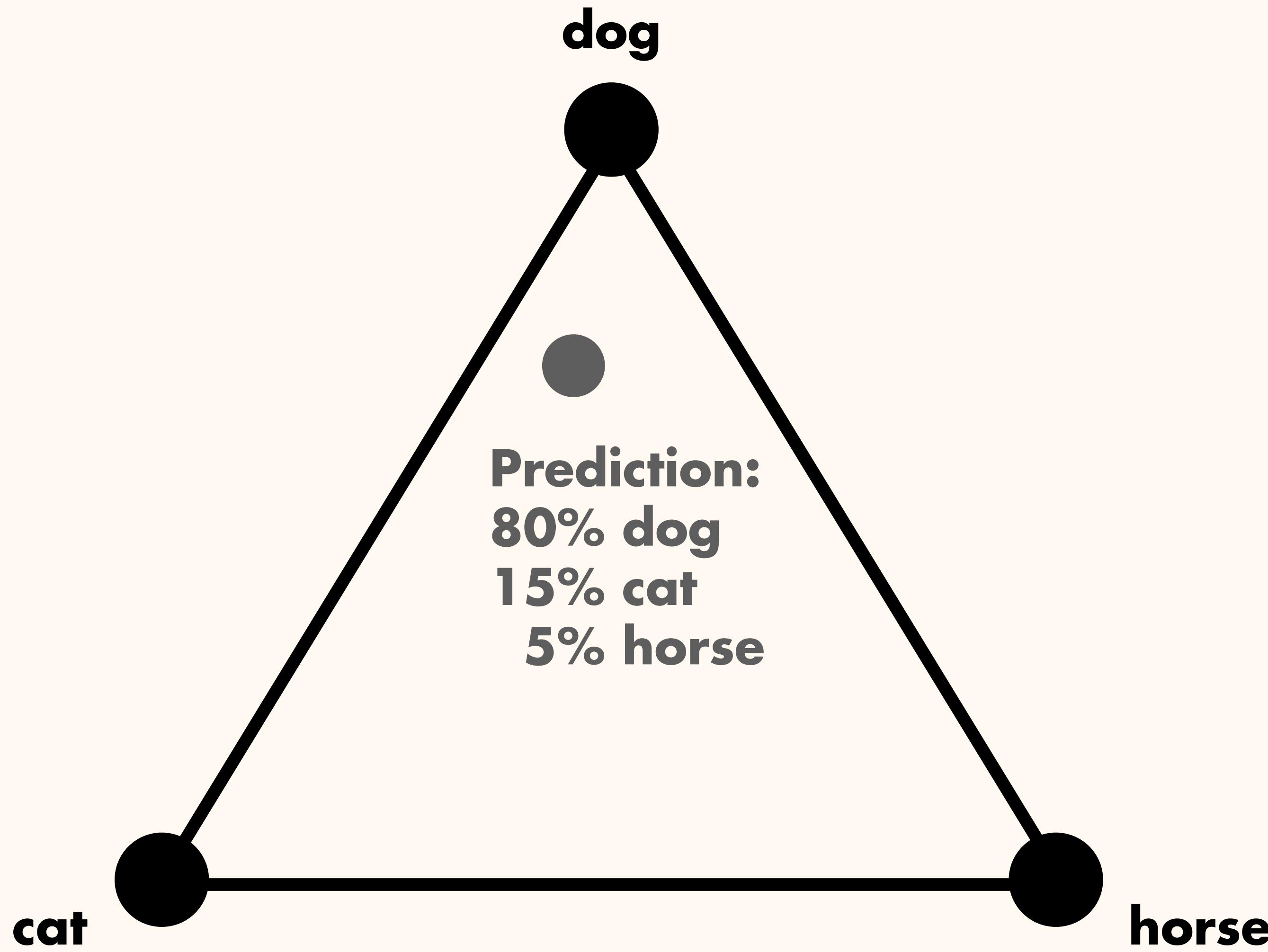


© TWITTER

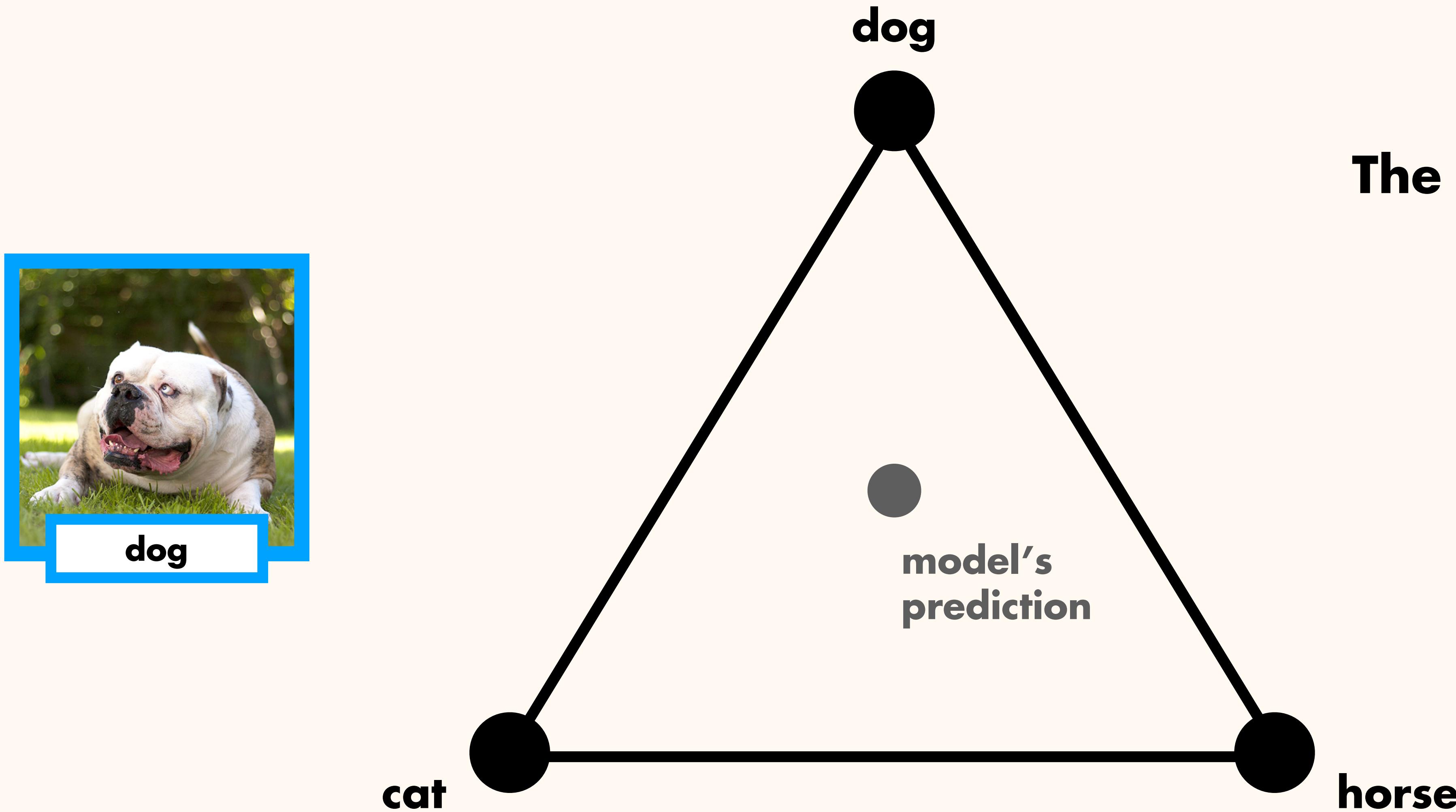






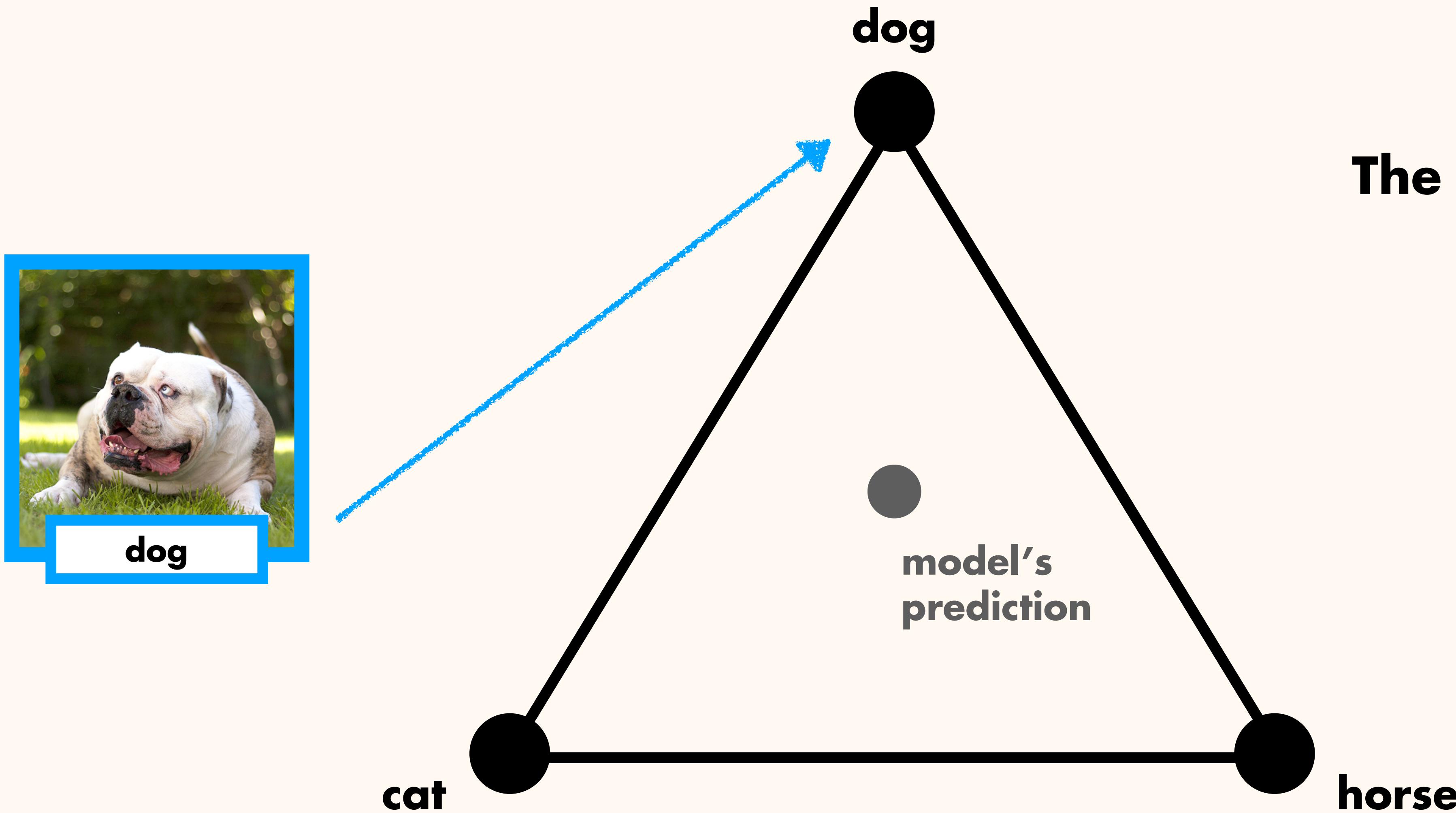


SUPERVISED LEARNING



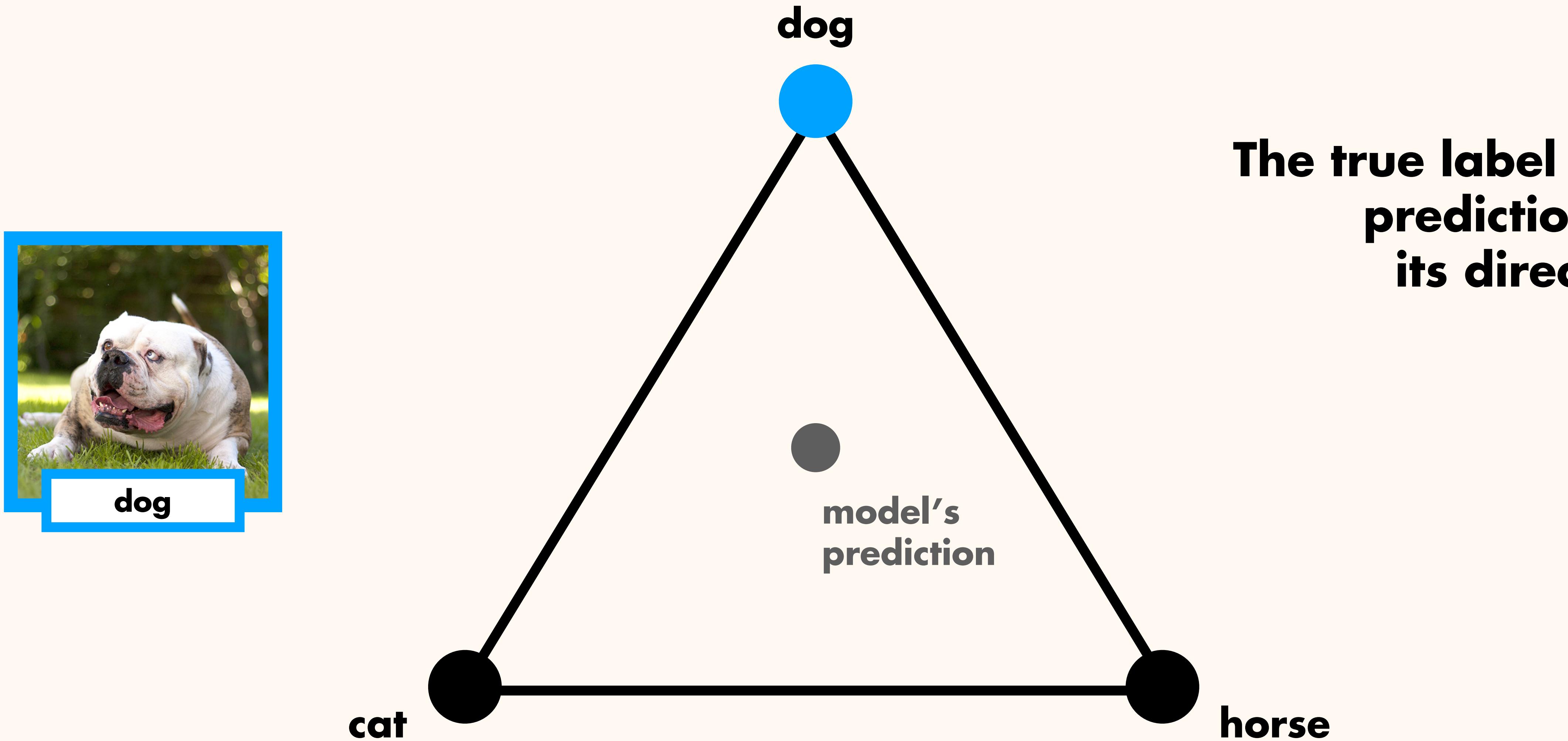
The true label pulls predictions to its direction.

SUPERVISED LEARNING



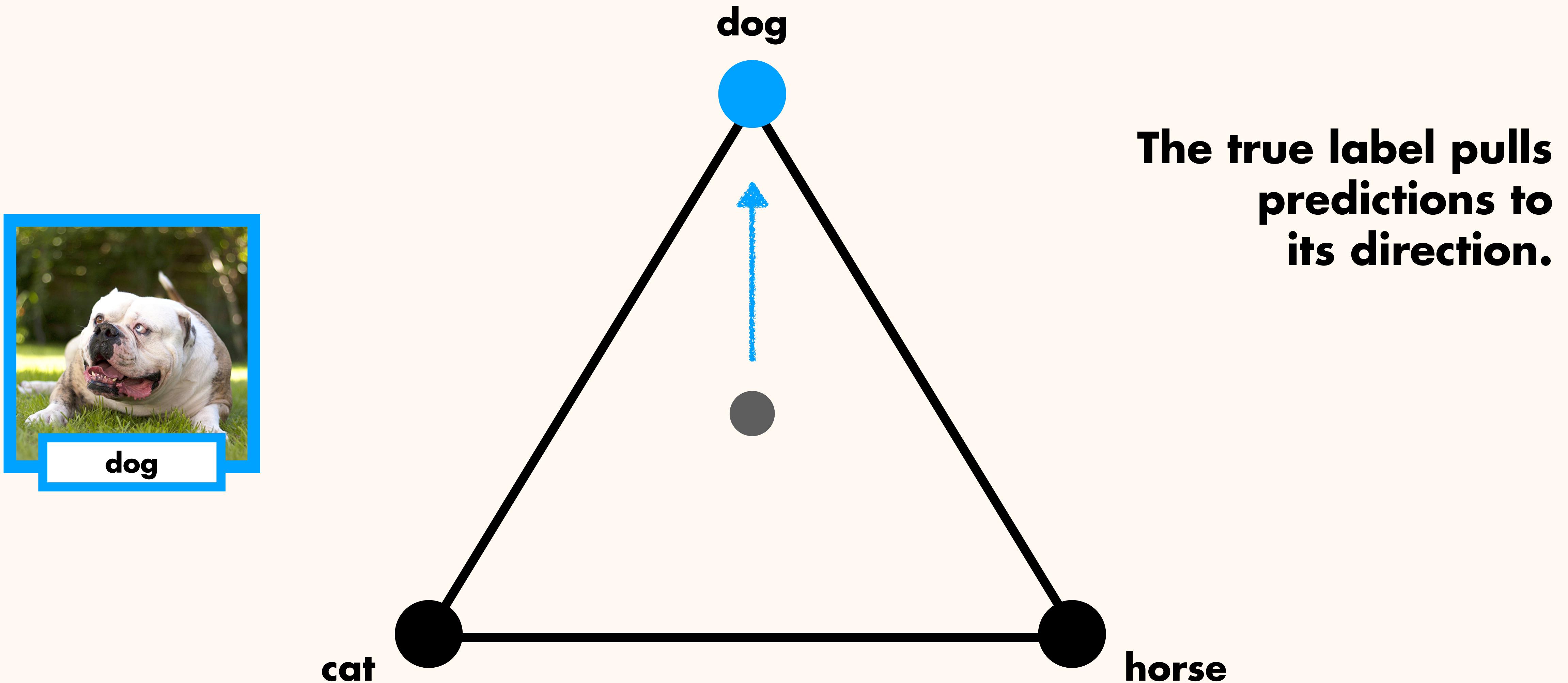
The true label pulls predictions to its direction.

SUPERVISED LEARNING

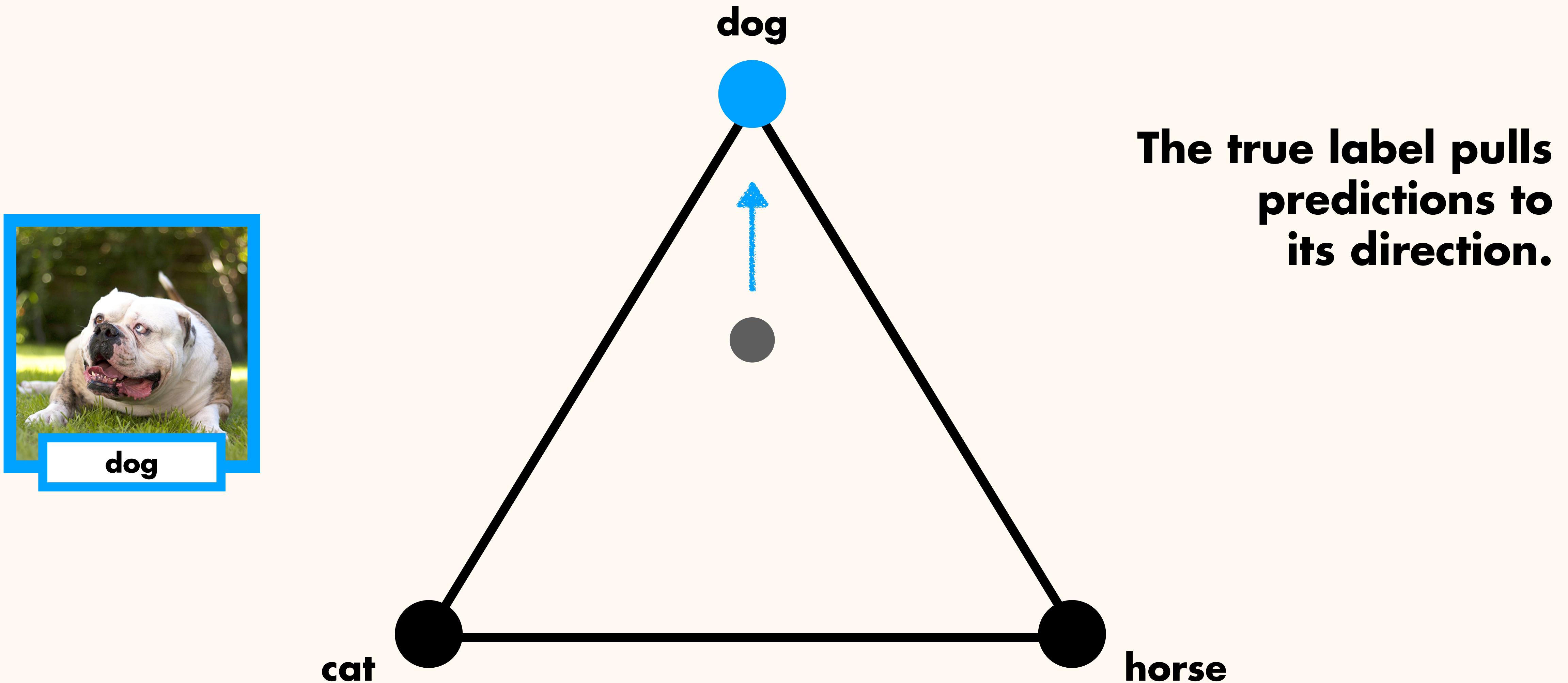


The true label pulls predictions to its direction.

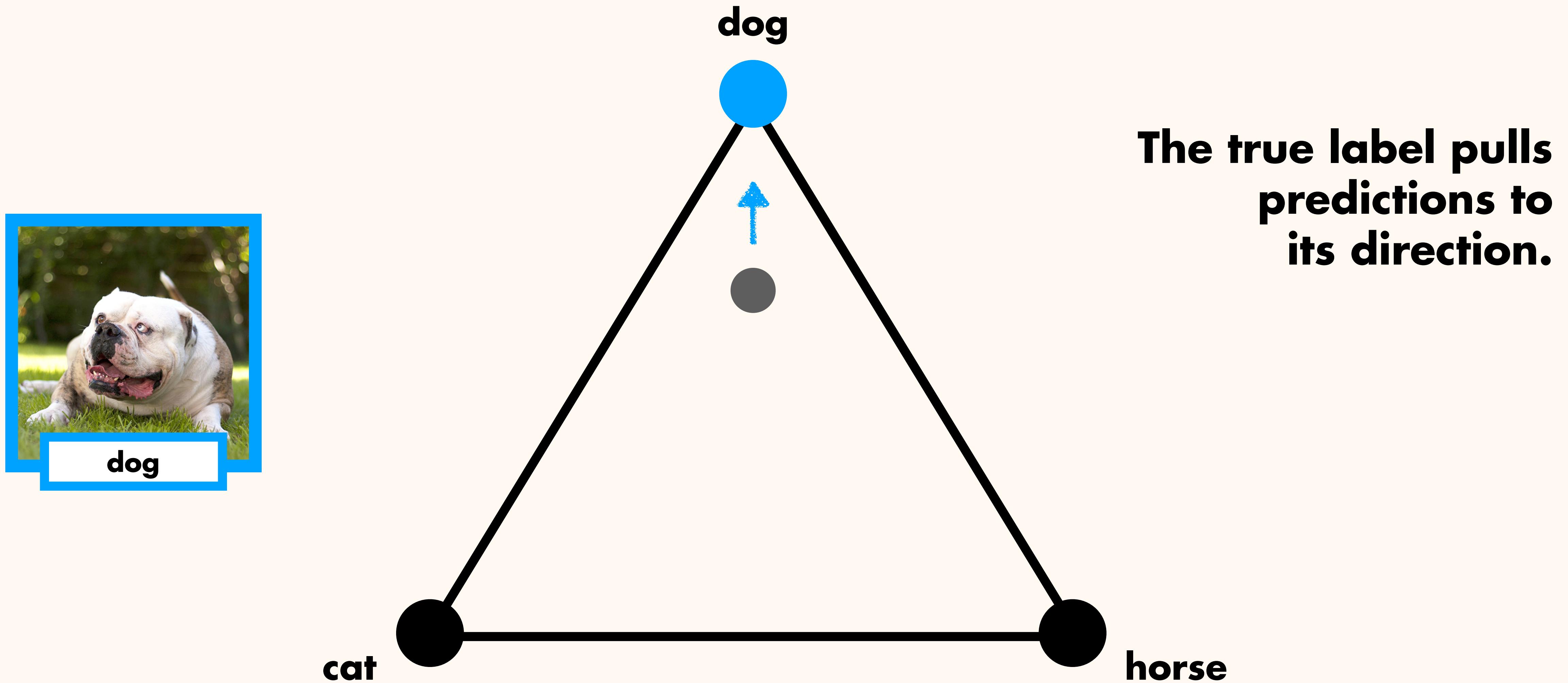
SUPERVISED LEARNING



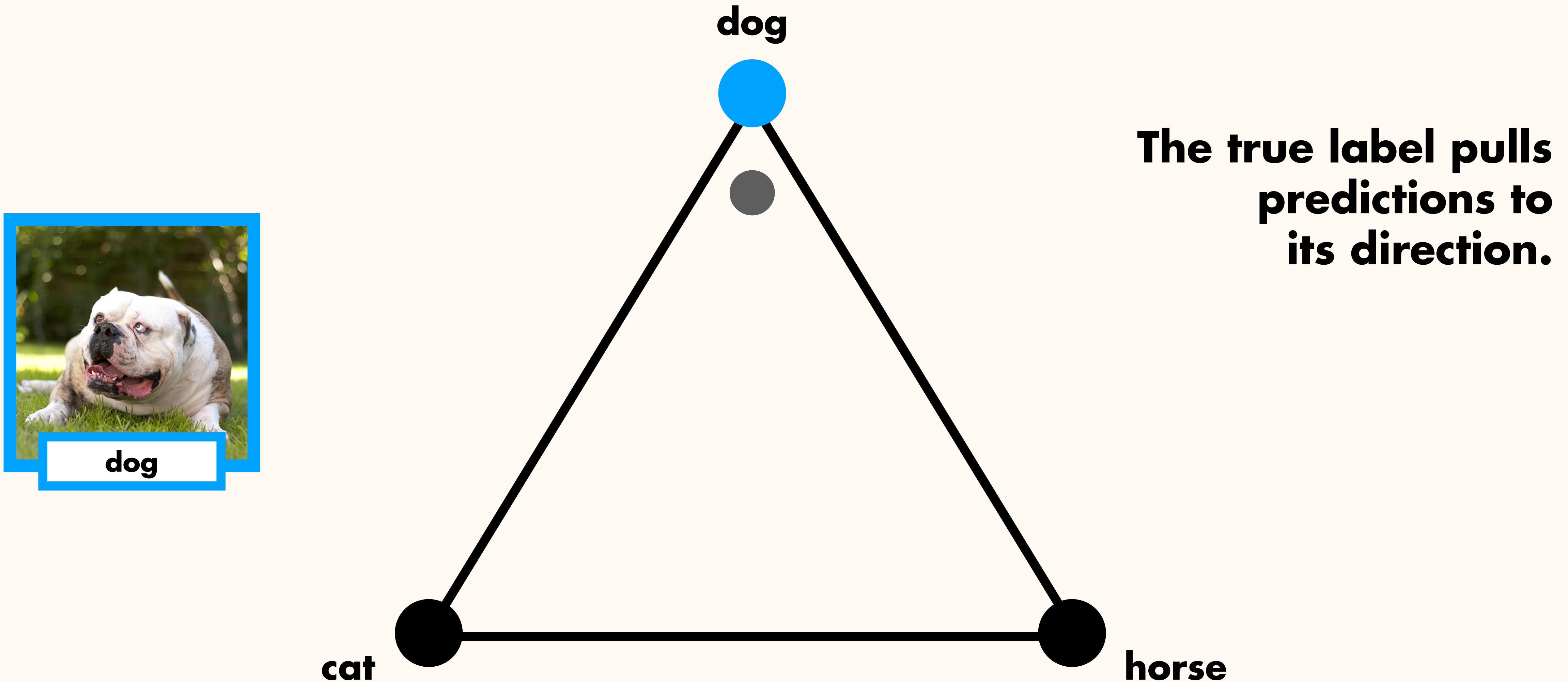
SUPERVISED LEARNING



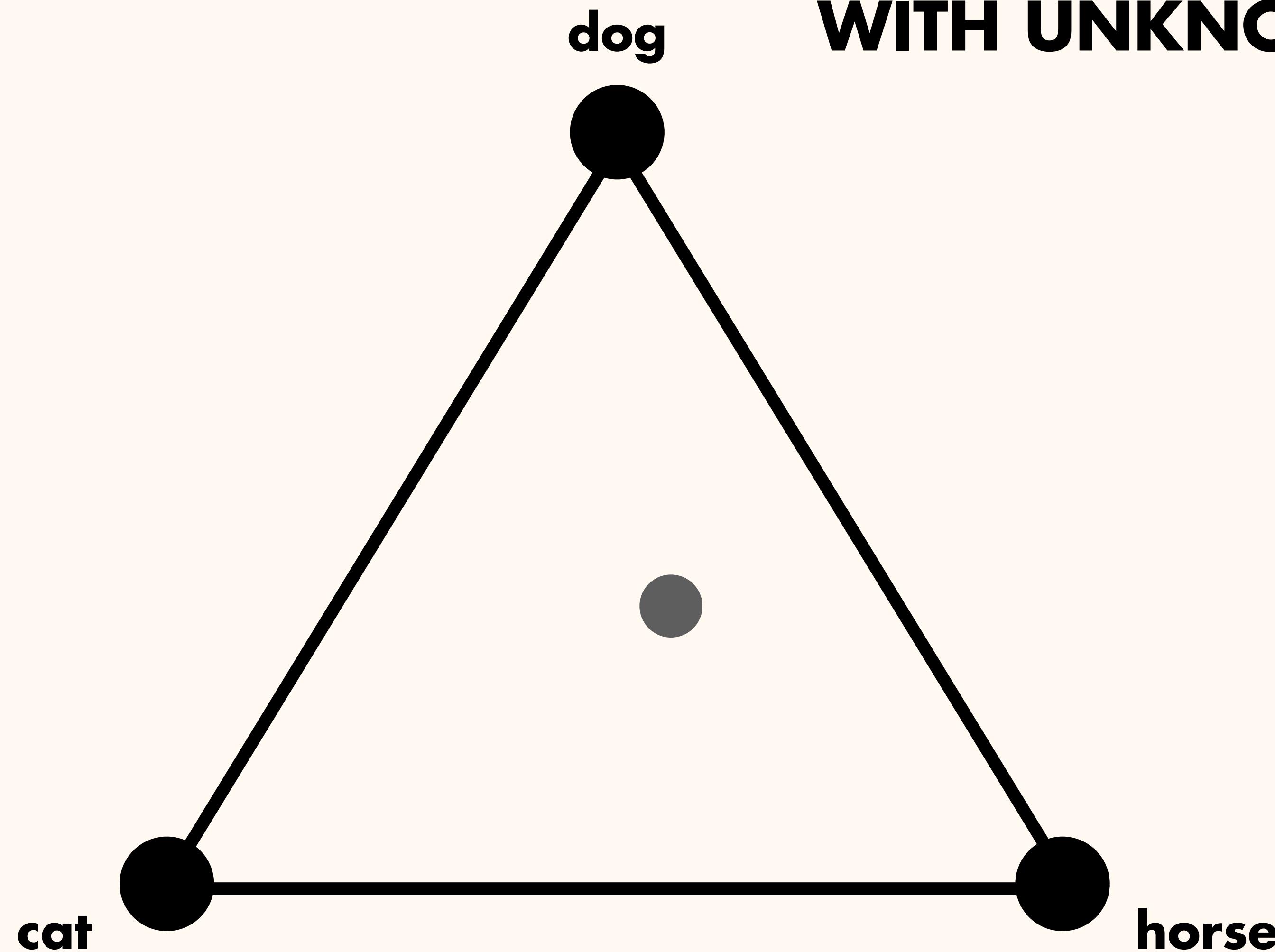
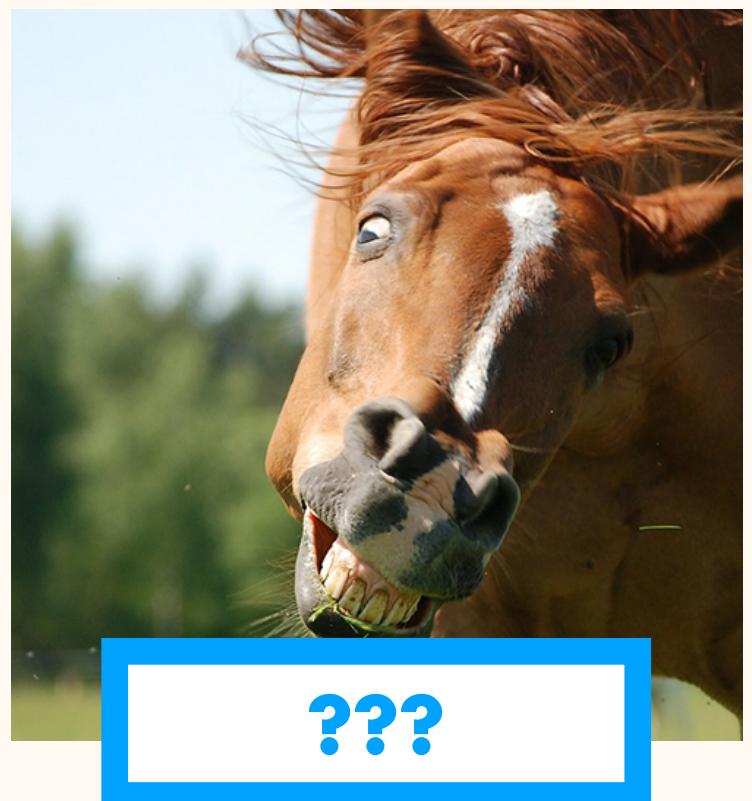
SUPERVISED LEARNING



SUPERVISED LEARNING

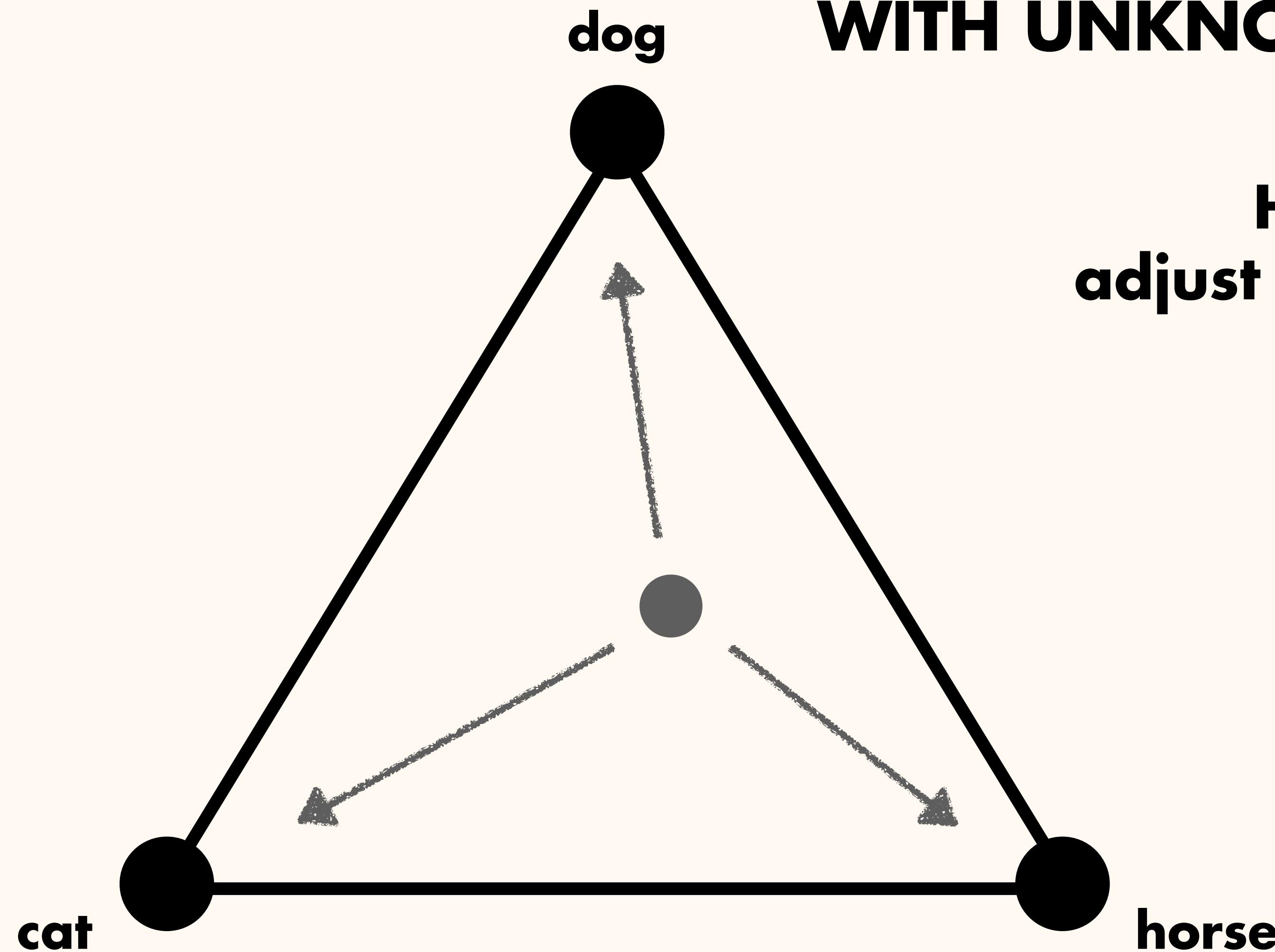
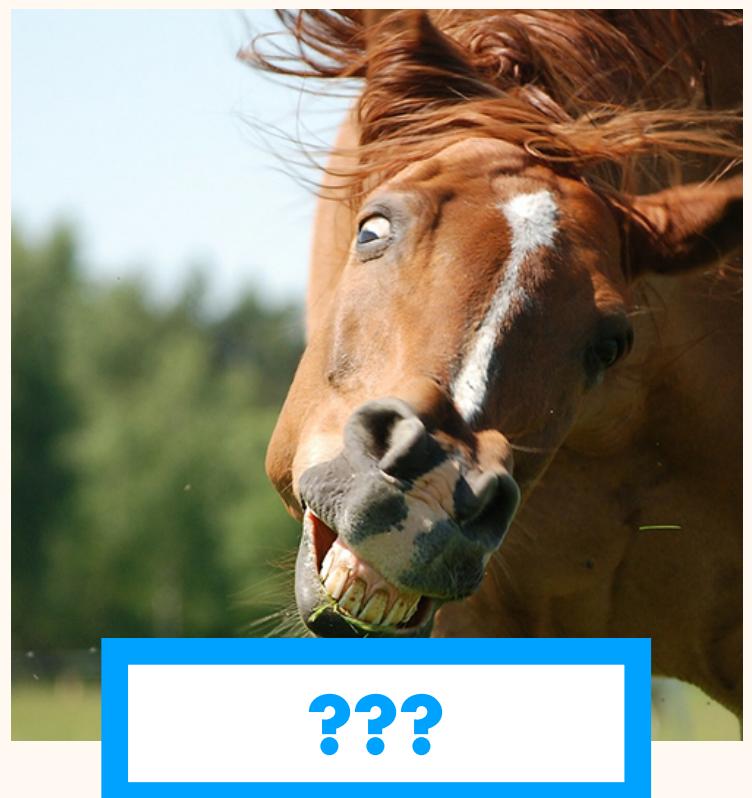


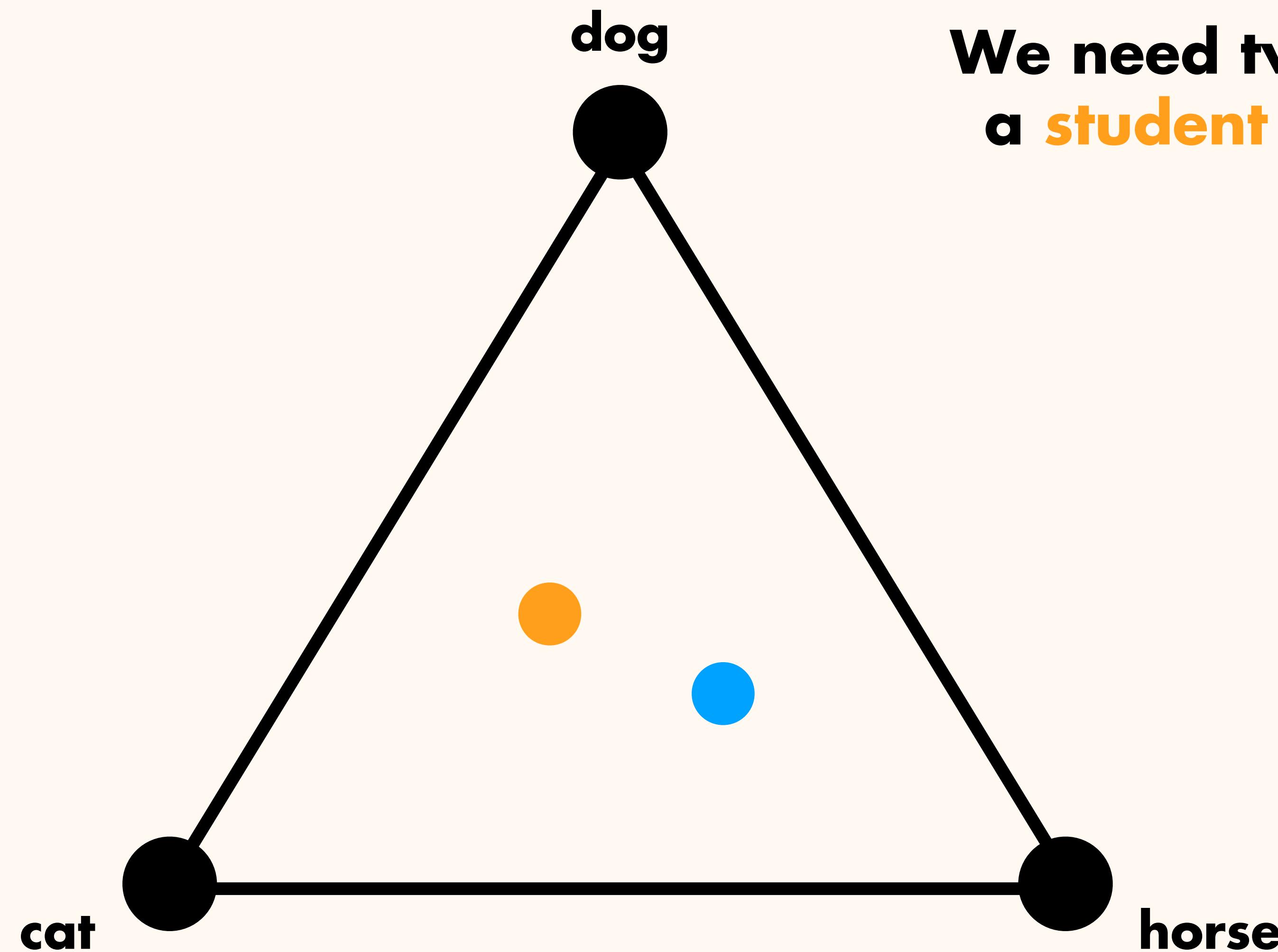
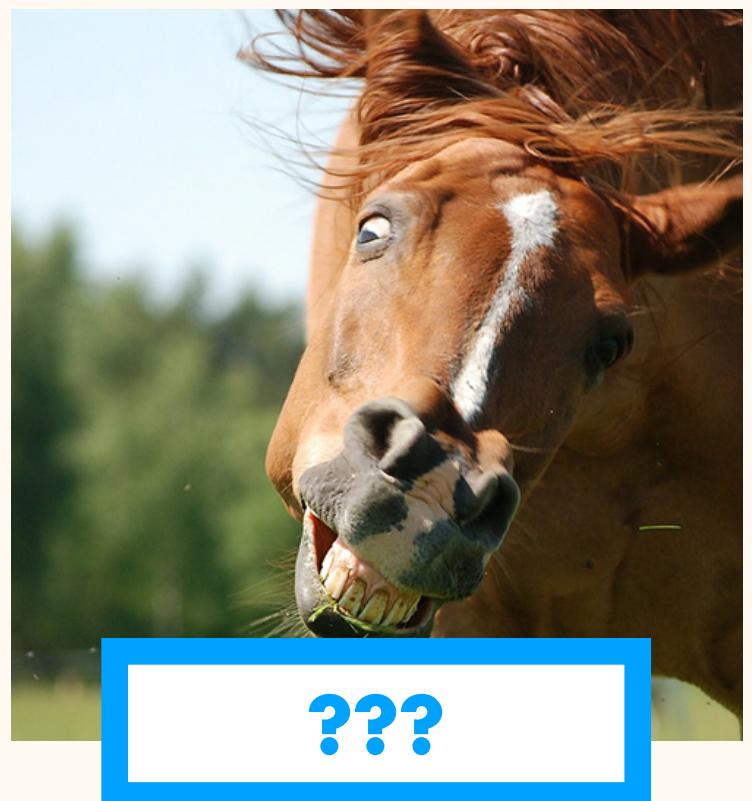
WHAT ABOUT EXAMPLES WITH UNKNOWN LABELS?



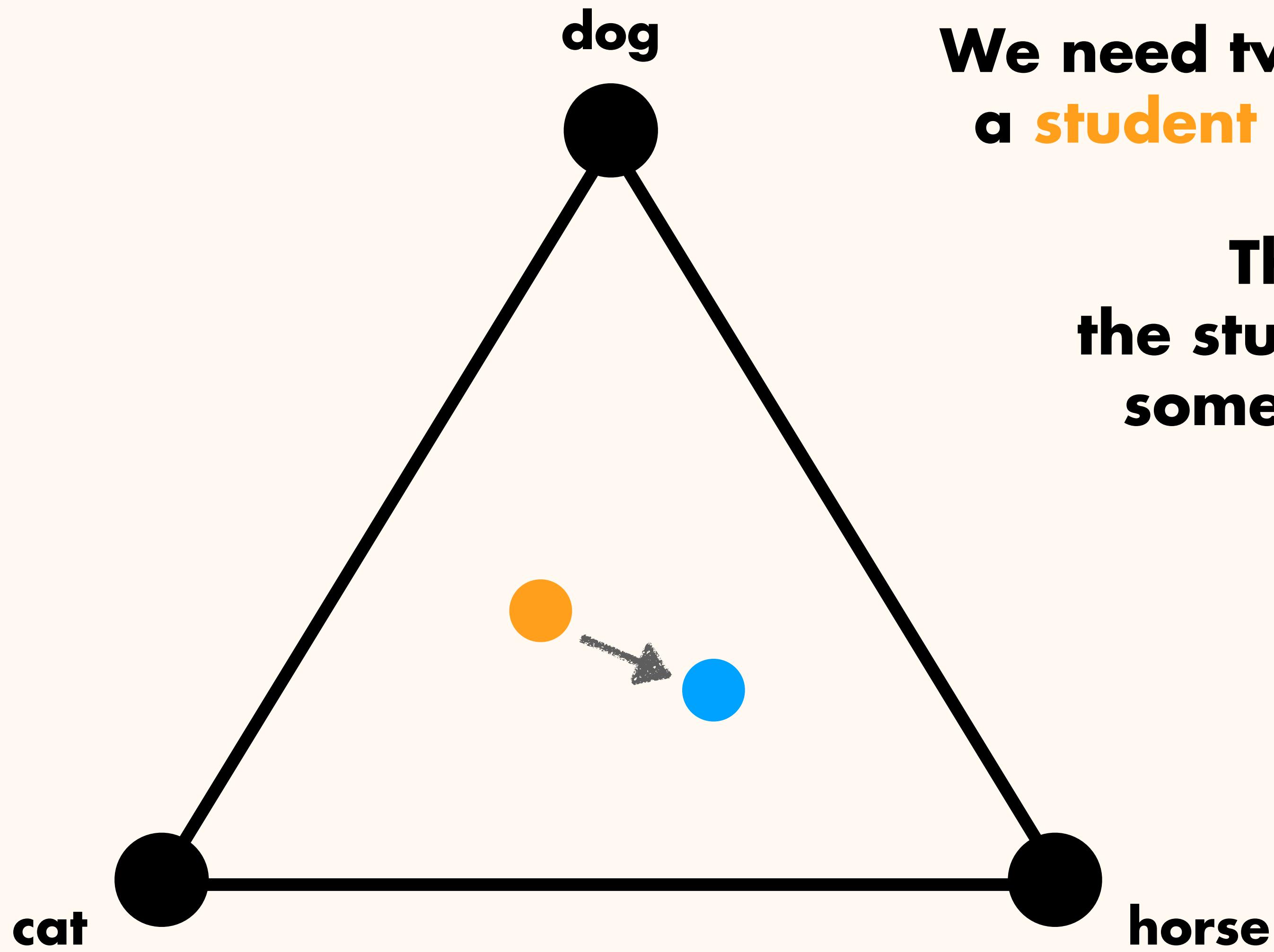
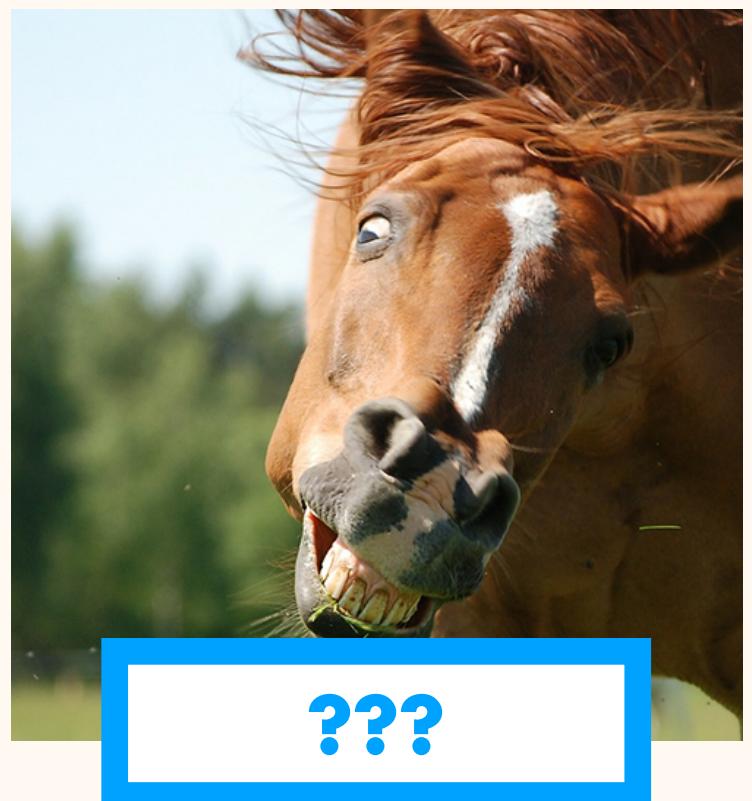
WHAT ABOUT EXAMPLES WITH UNKNOWN LABELS?

How should we
adjust the prediction?



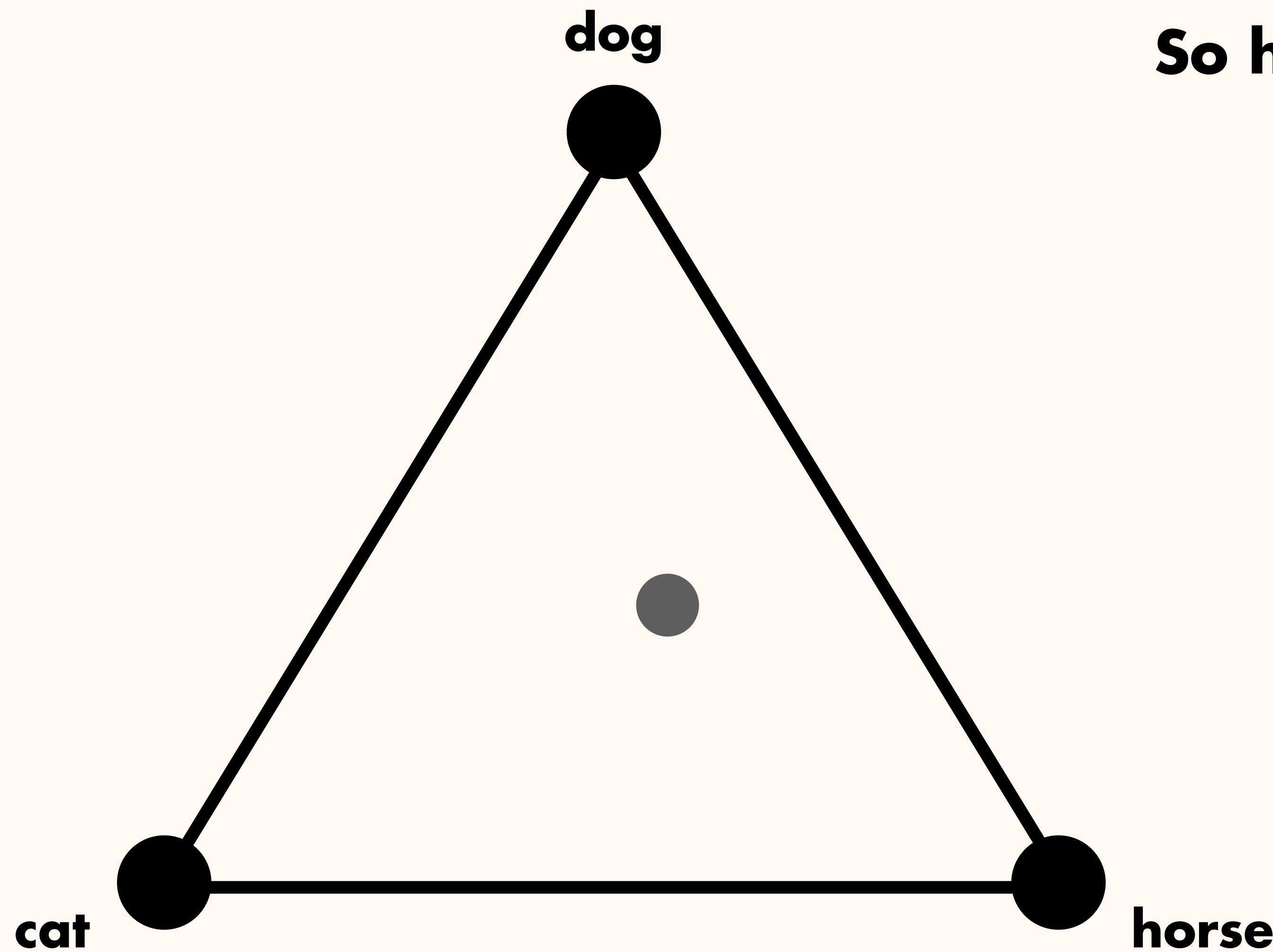
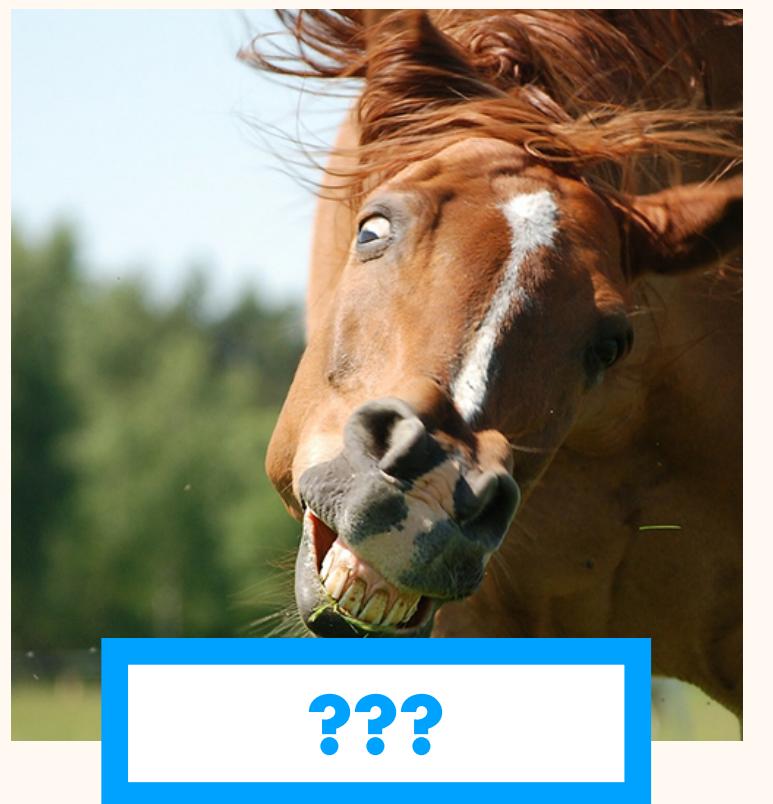


We need two predictions:
a **student** and a **teacher**.



We need two predictions:
a **student** and a **teacher**.

Then, hopefully,
the **student** can learn
something from the
teacher.

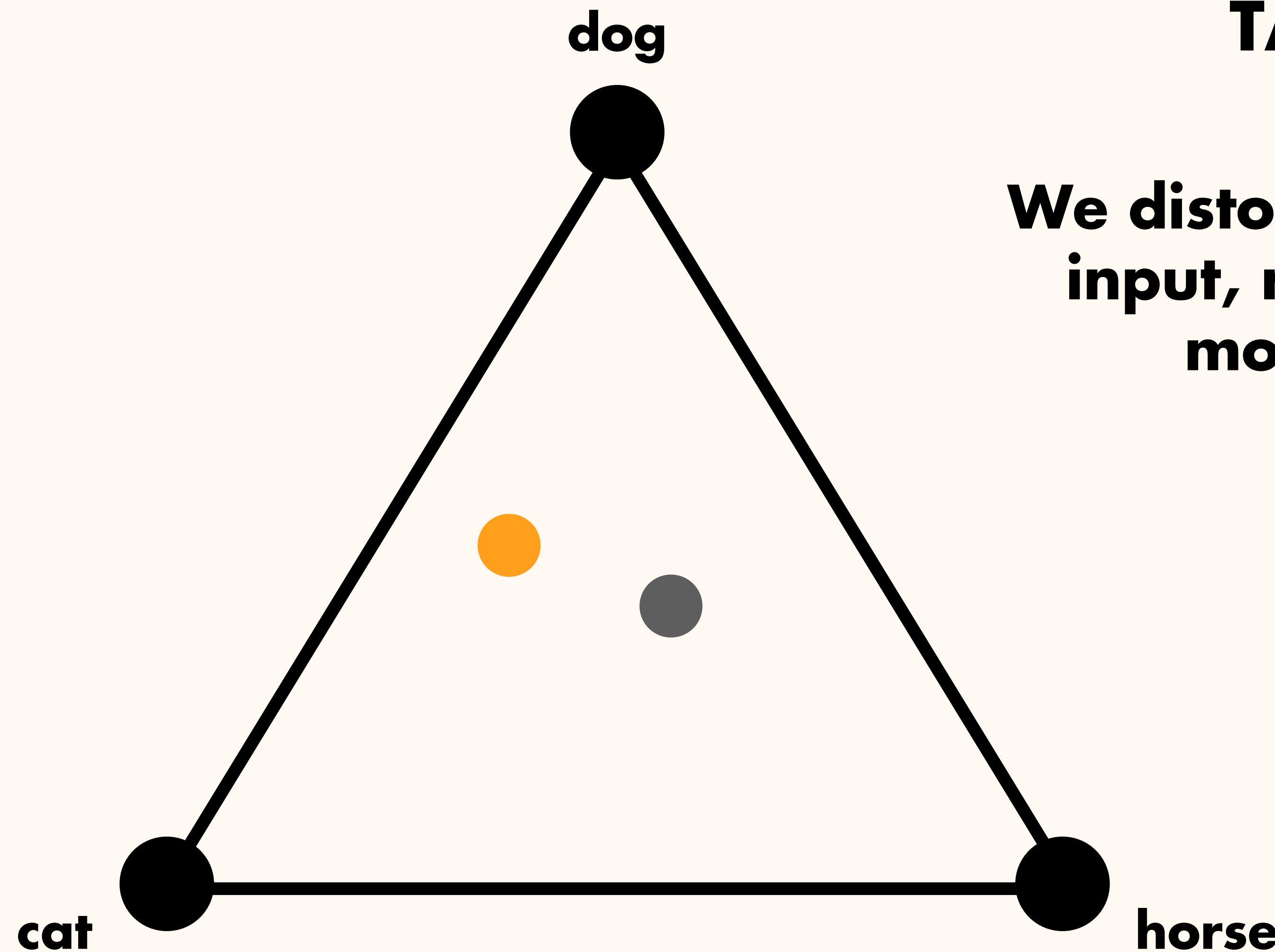


So how to do this?

Two ways.

1. MAKE THE TASK HARDER.

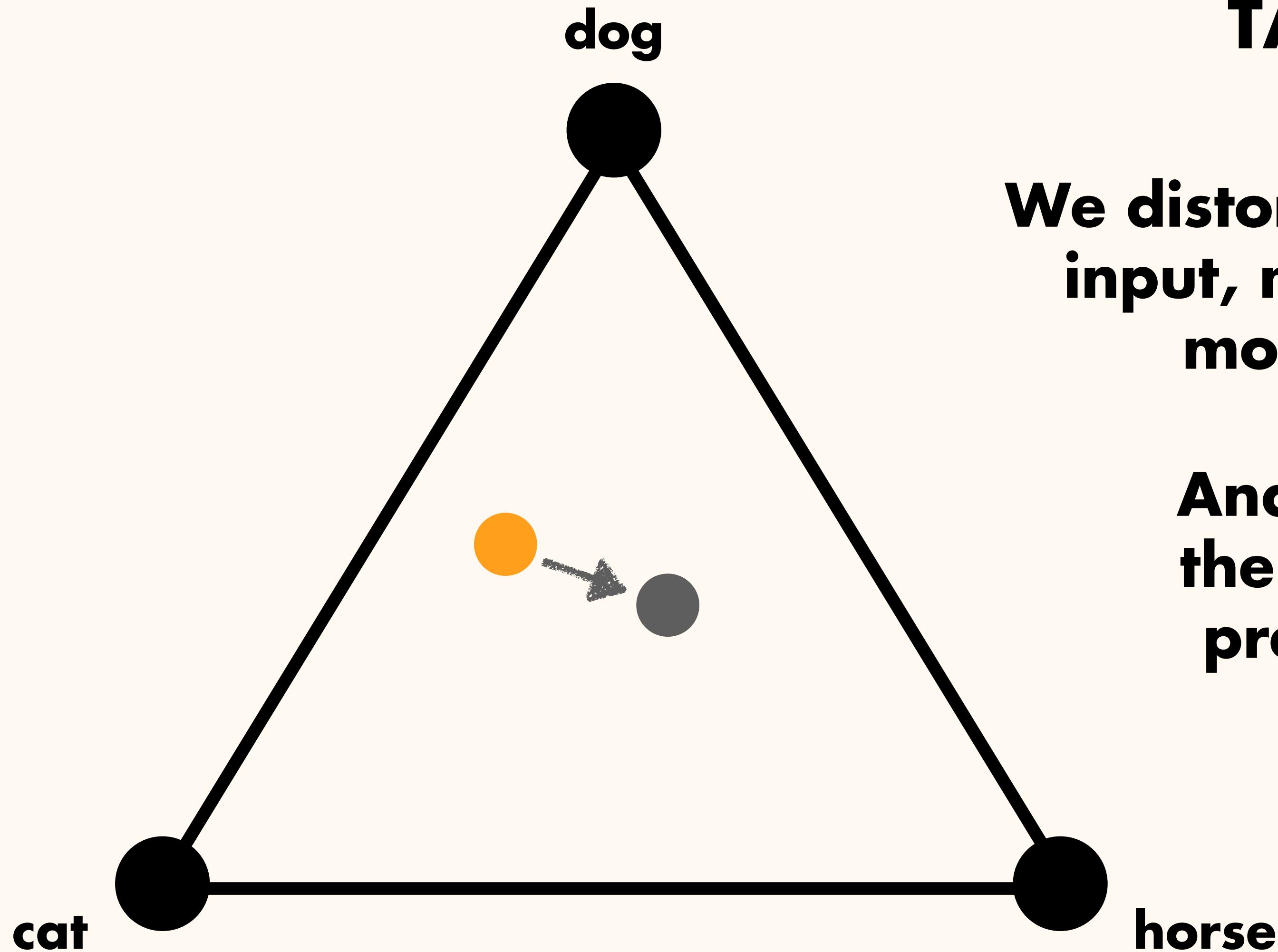
We distort the **student's**
input, making its task
more challenging.



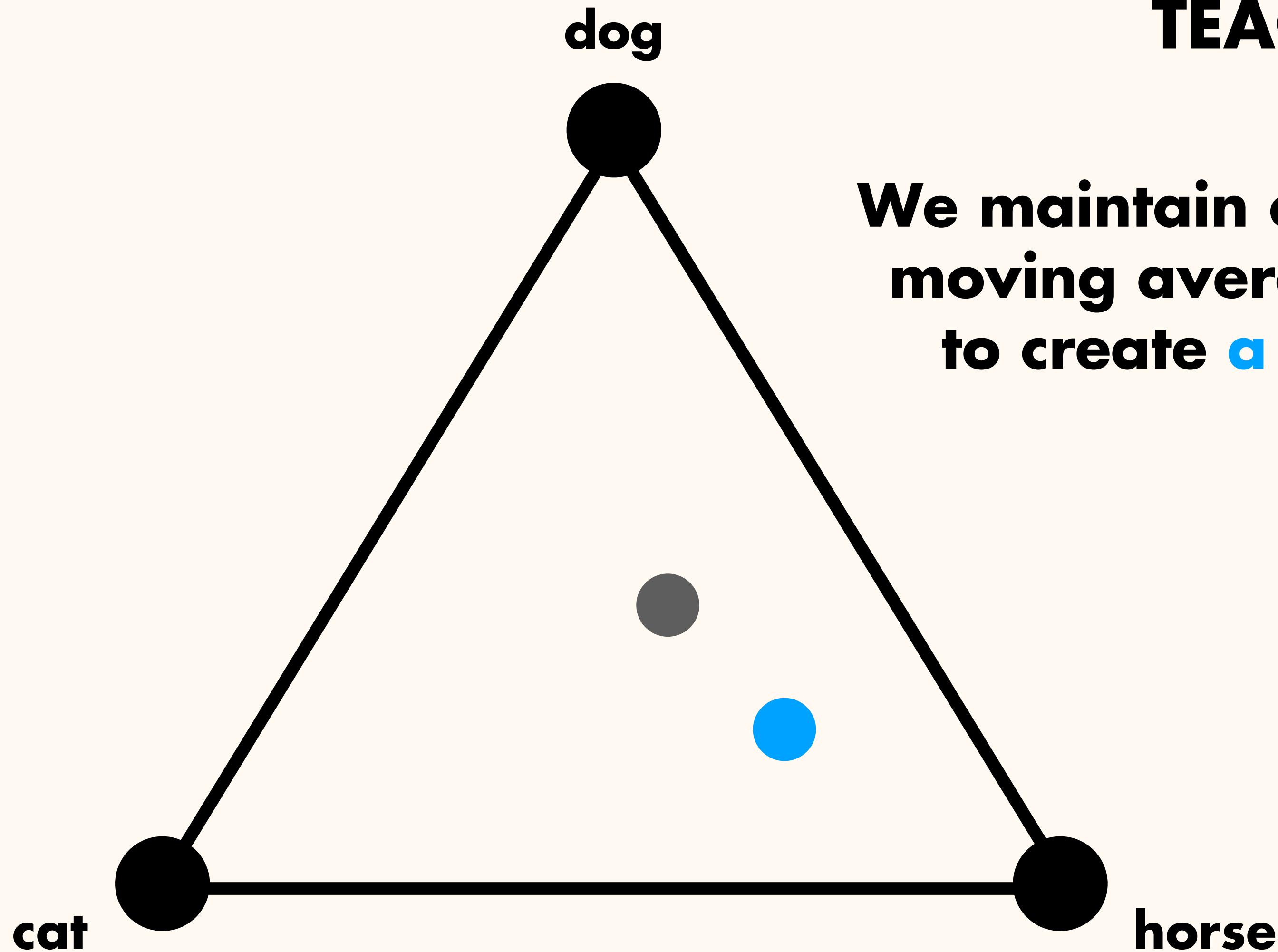
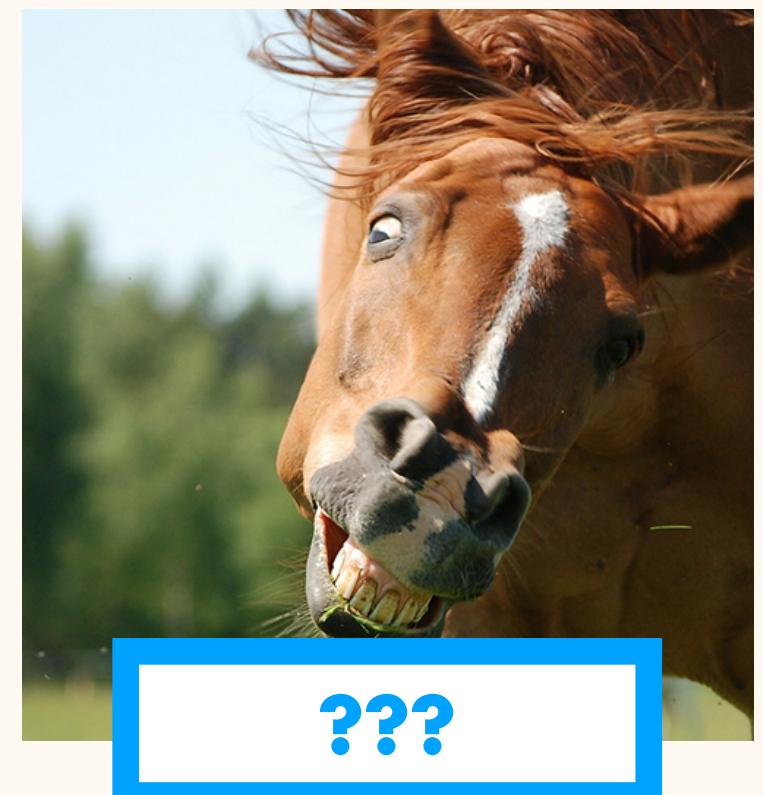
1. MAKE THE TASK HARDER.

We distort the **student's** input, making its task more challenging.

And then we train the harder task to predict the easier tasks' output.

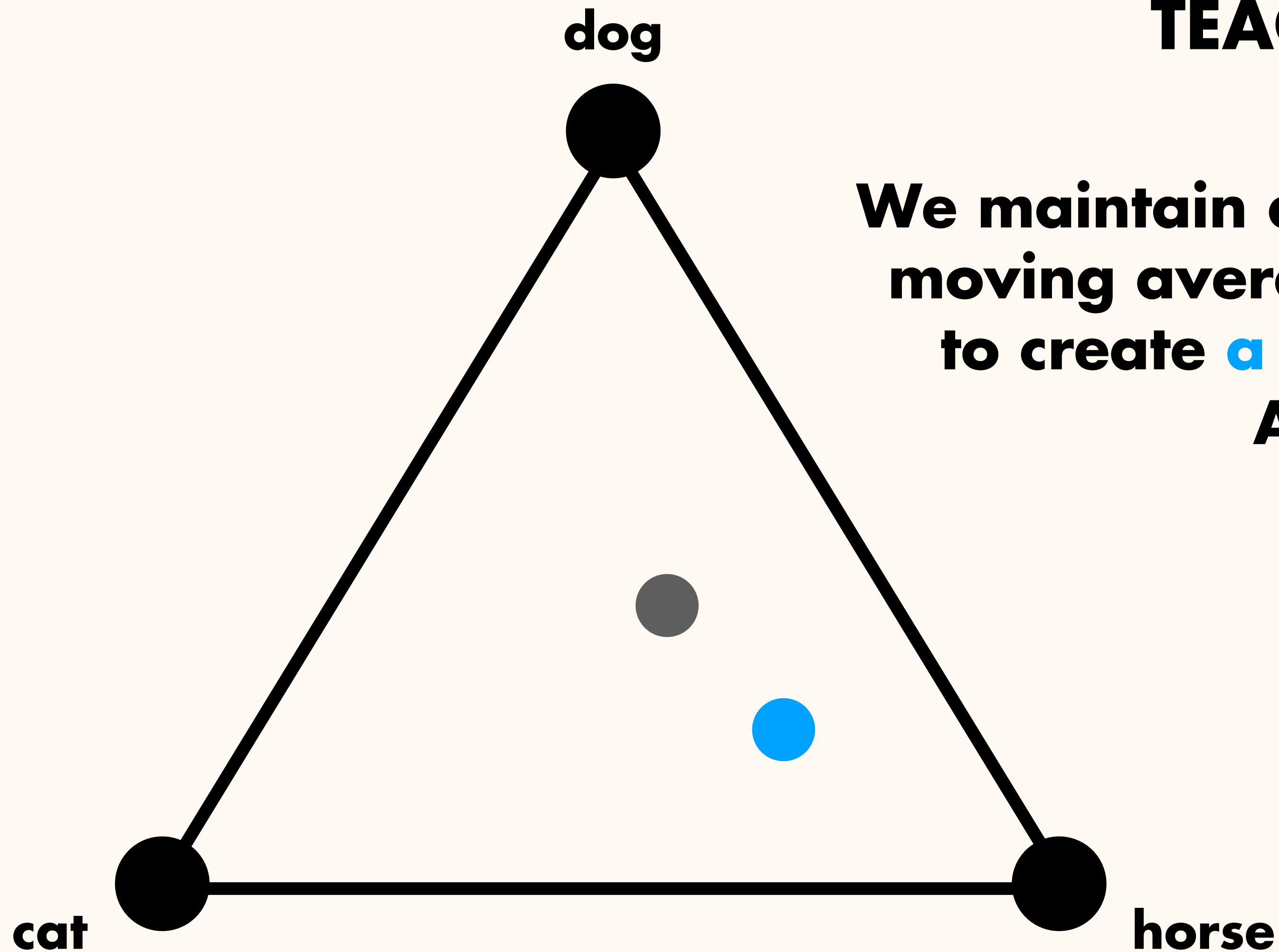
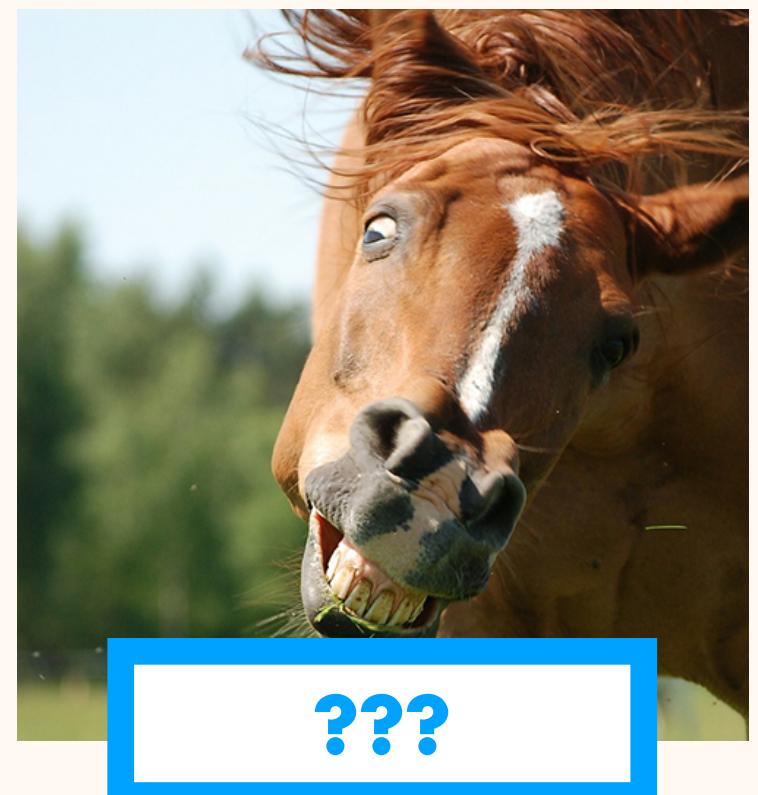


2. MAKE THE TEACHER BETTER.



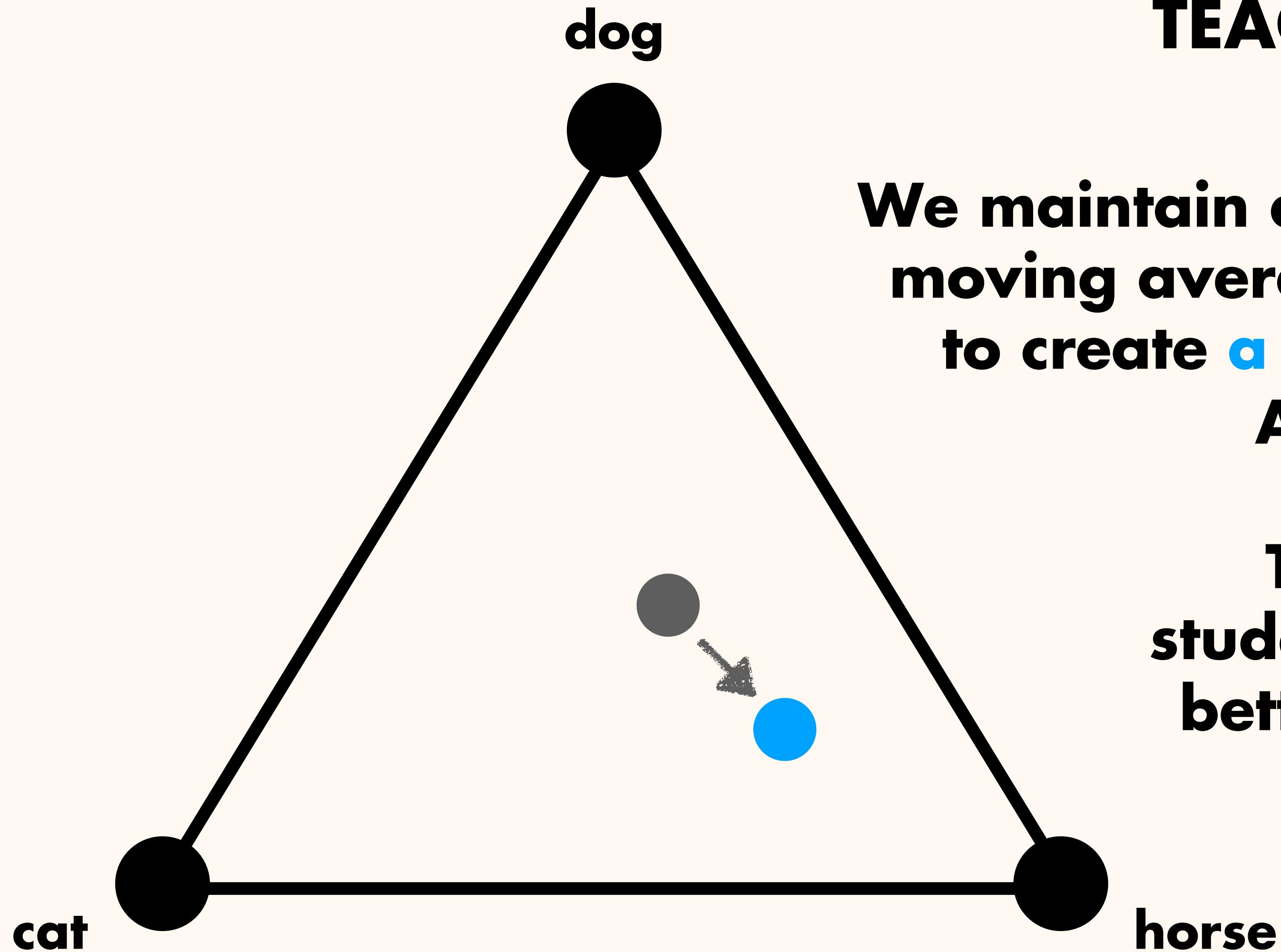
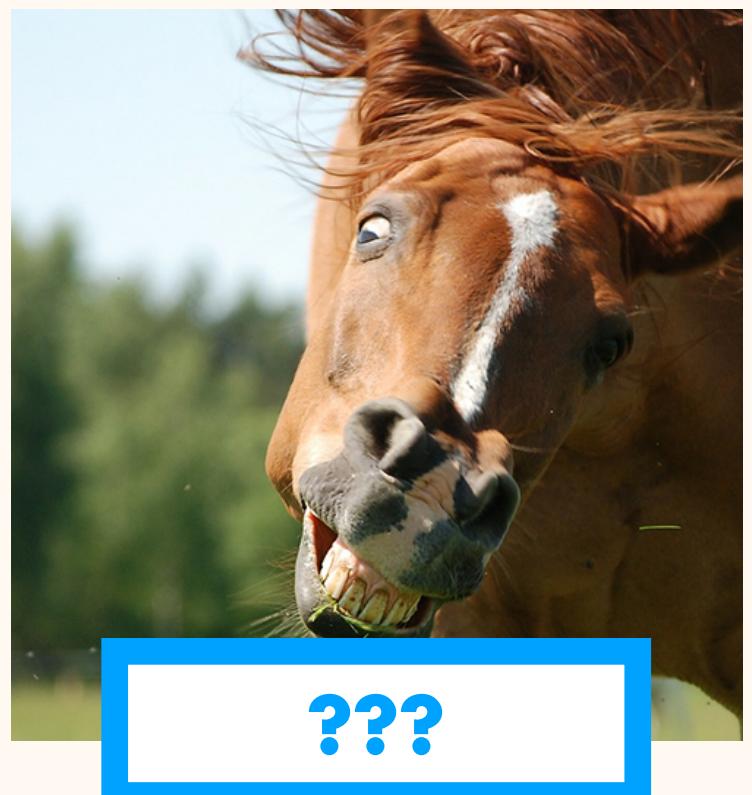
We maintain an exponential moving average of weights to create a better teacher.

2. MAKE THE TEACHER BETTER.



We maintain an exponential moving average of weights to create a better teacher.
A mean teacher.

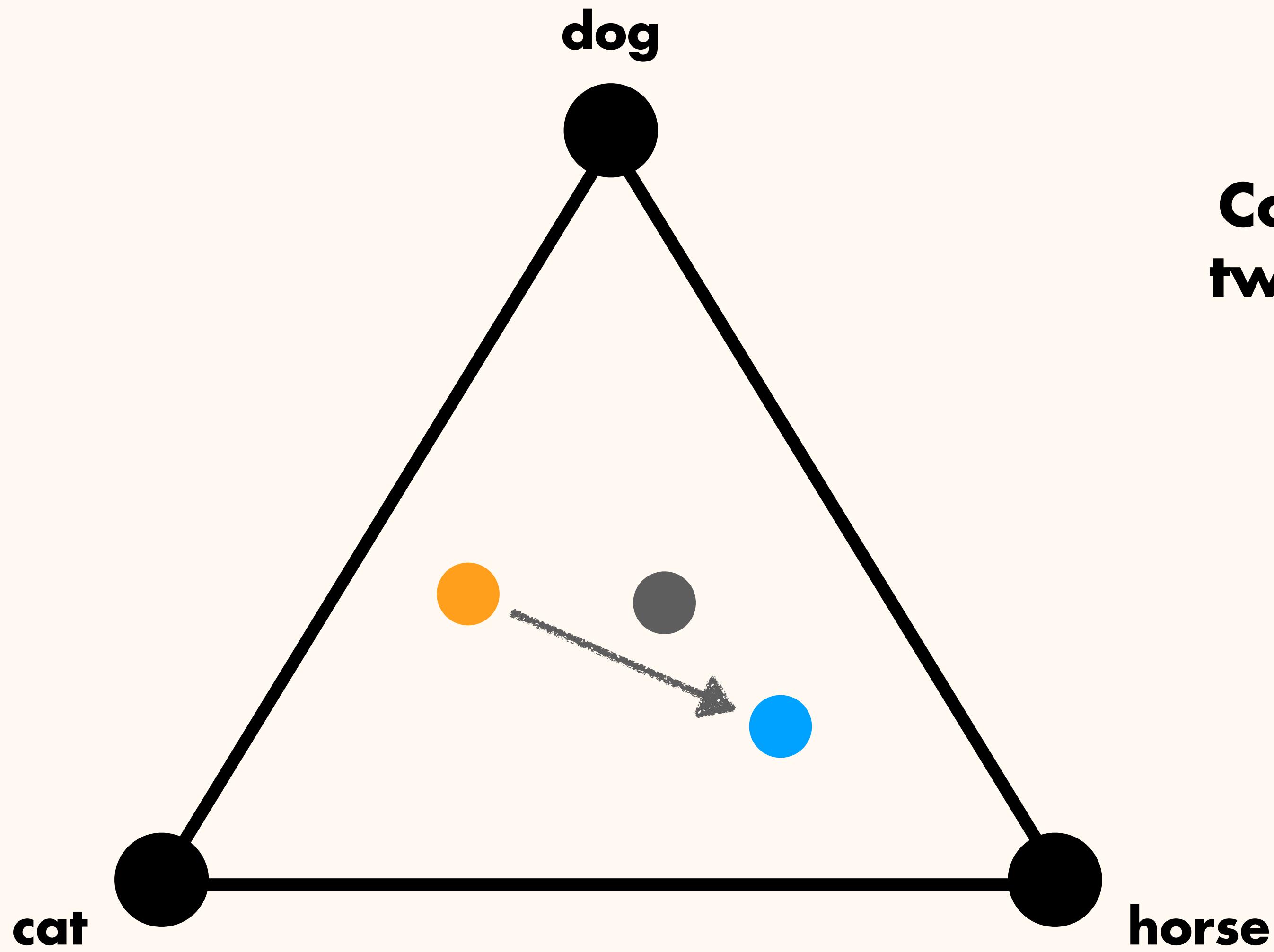
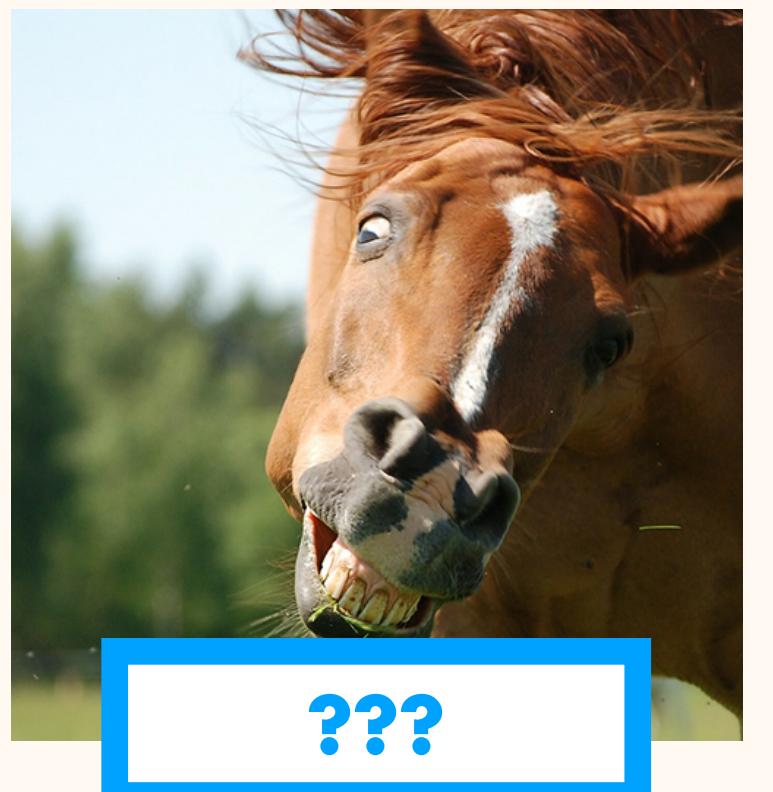
2. MAKE THE TEACHER BETTER.



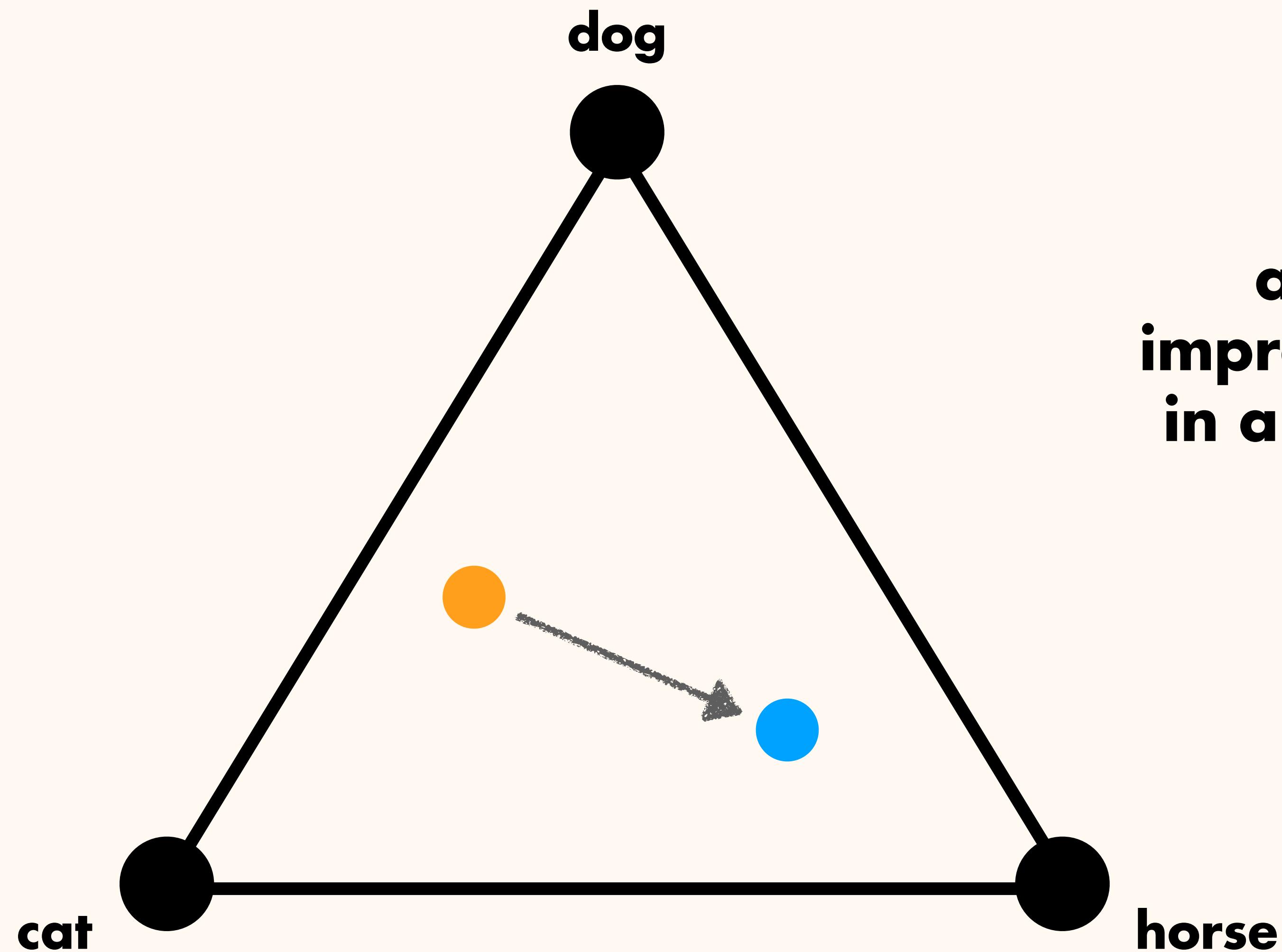
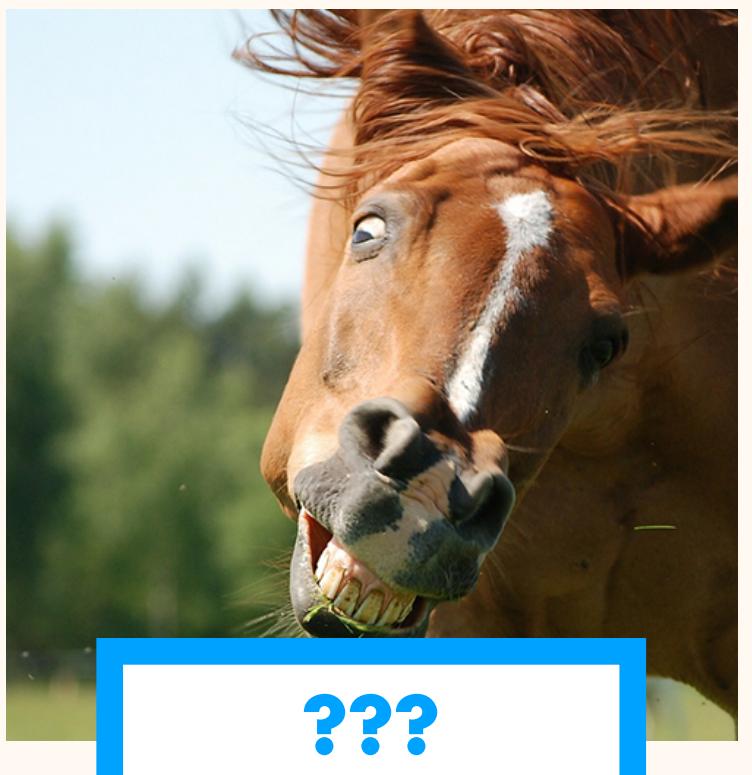
We maintain an exponential moving average of weights to create a better teacher.

A mean teacher.

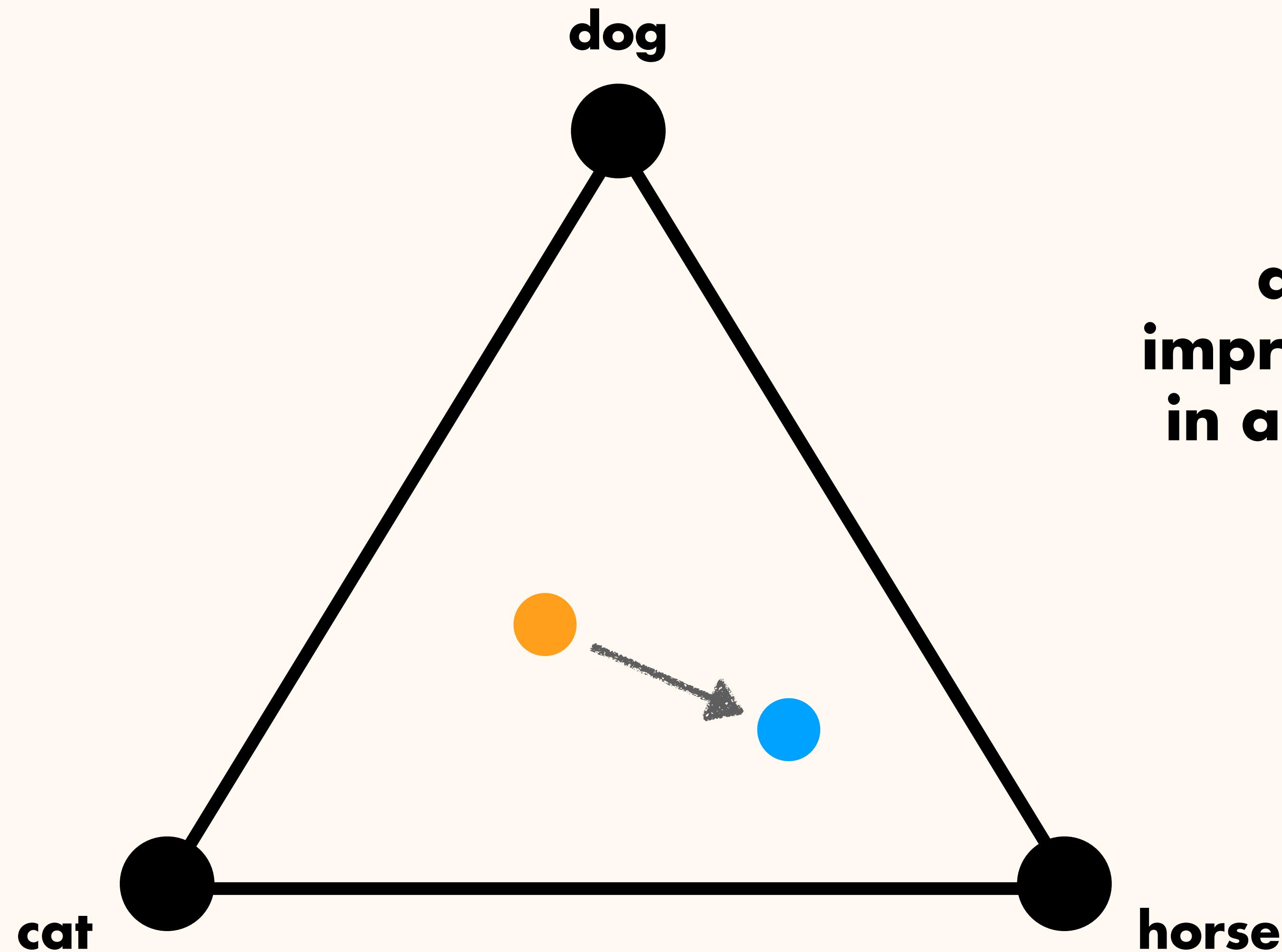
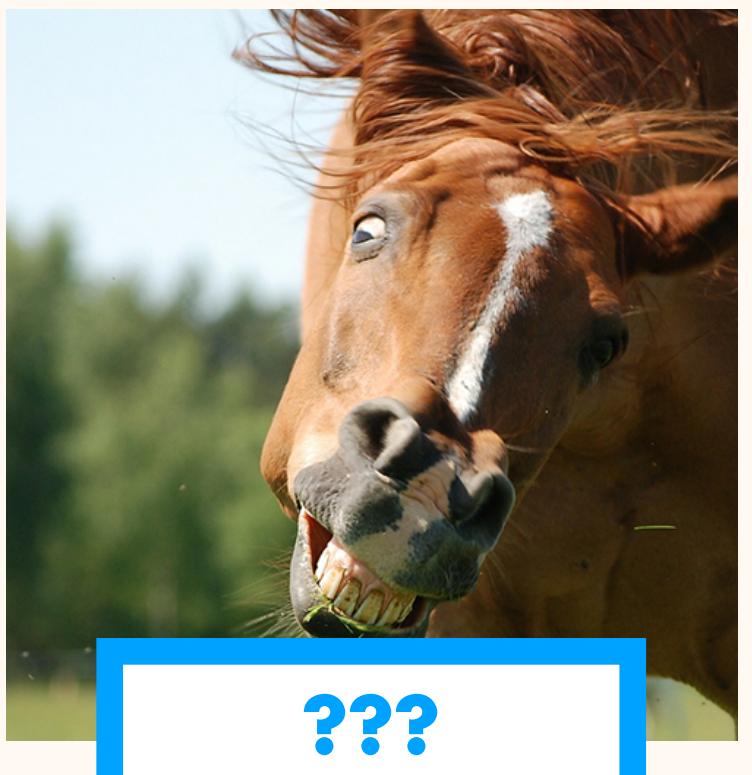
Then we let the student learn these better predictions.



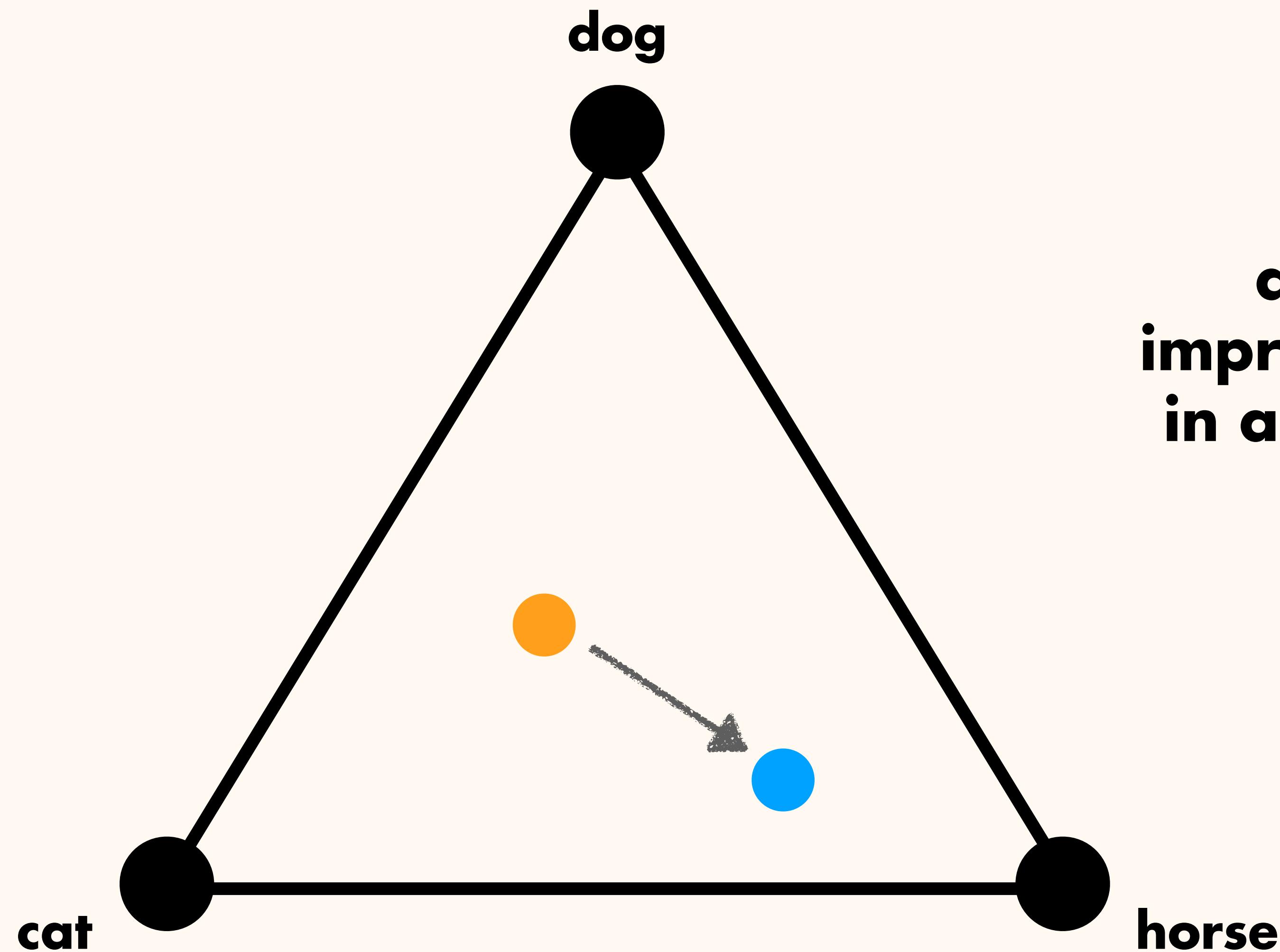
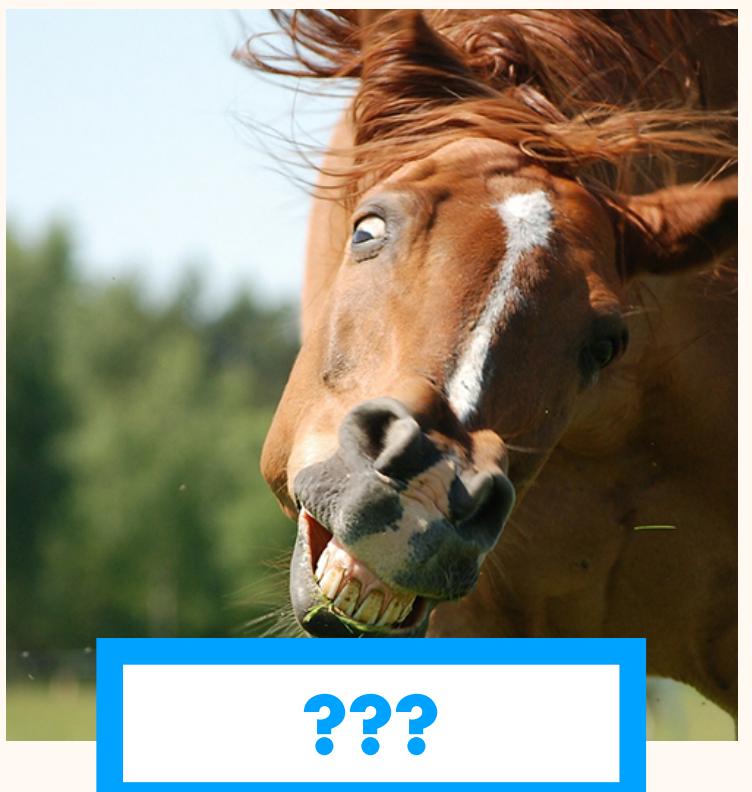
**Combining these
two ways works
even better.**



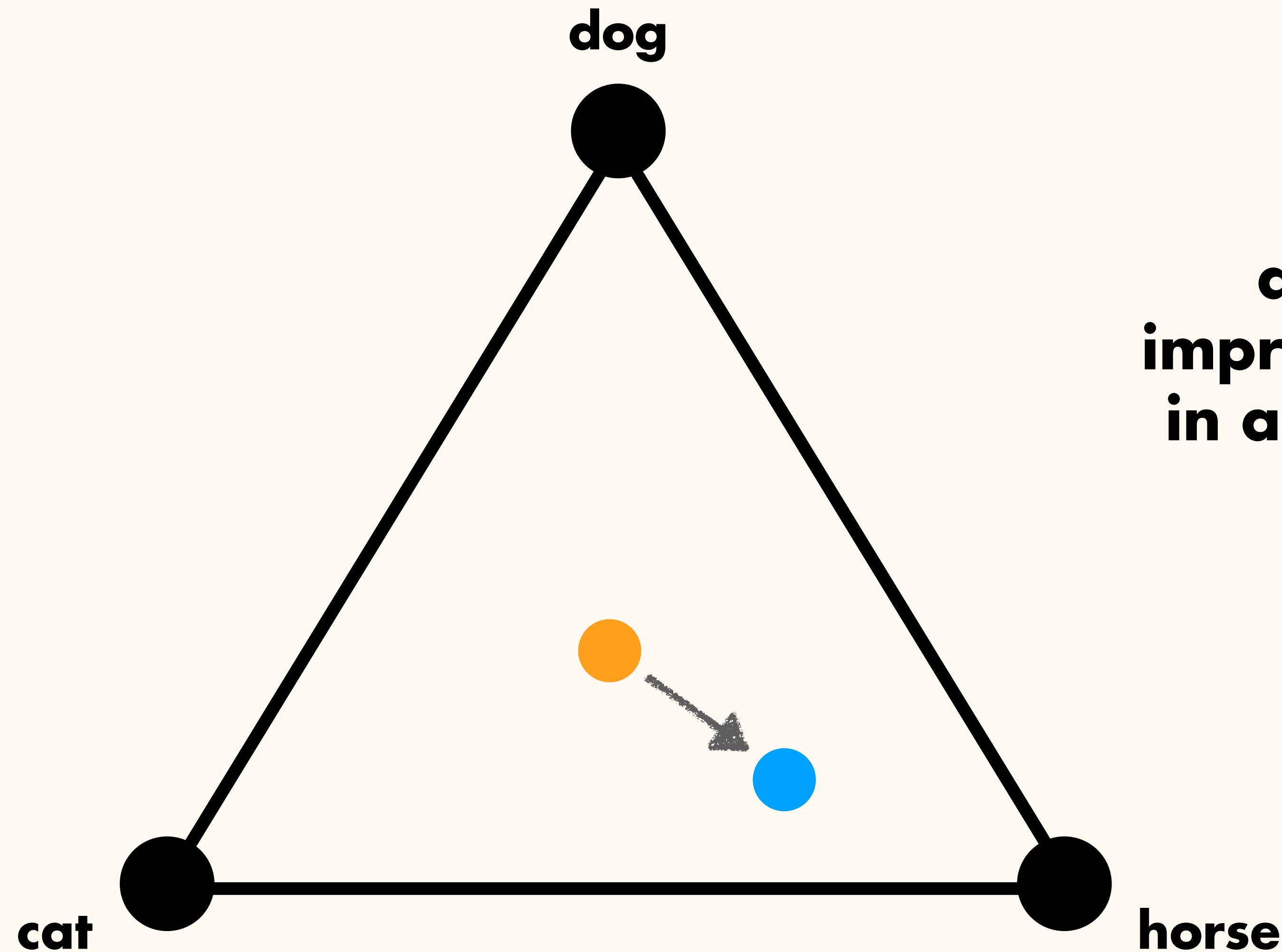
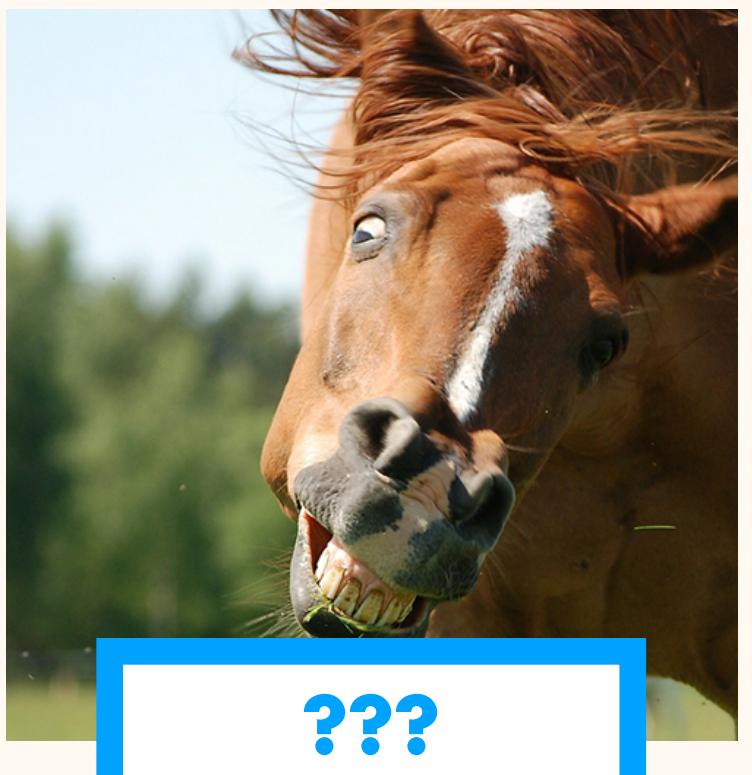
**The student
and the teacher
improve each other
in a *virtuous* cycle.**



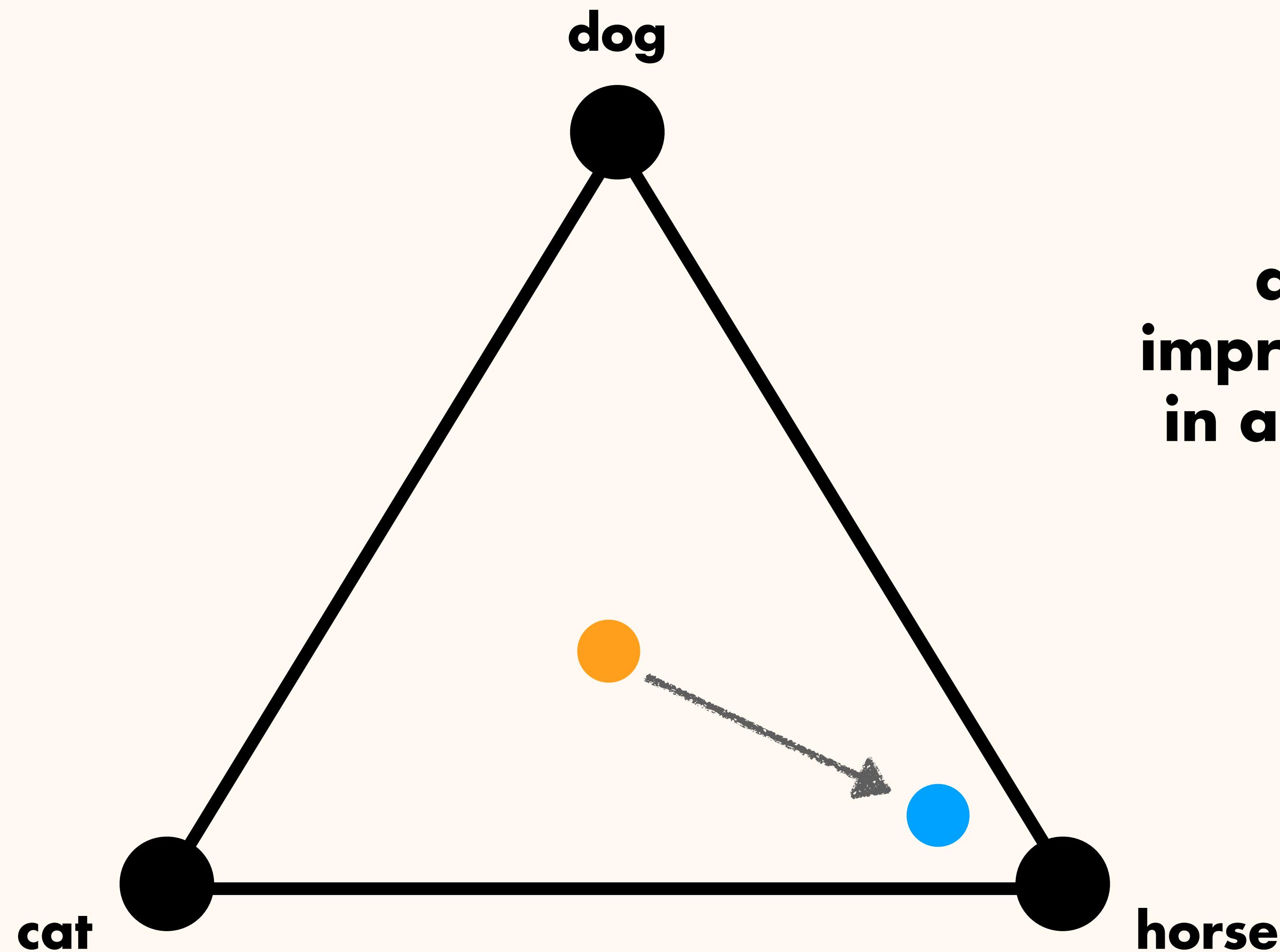
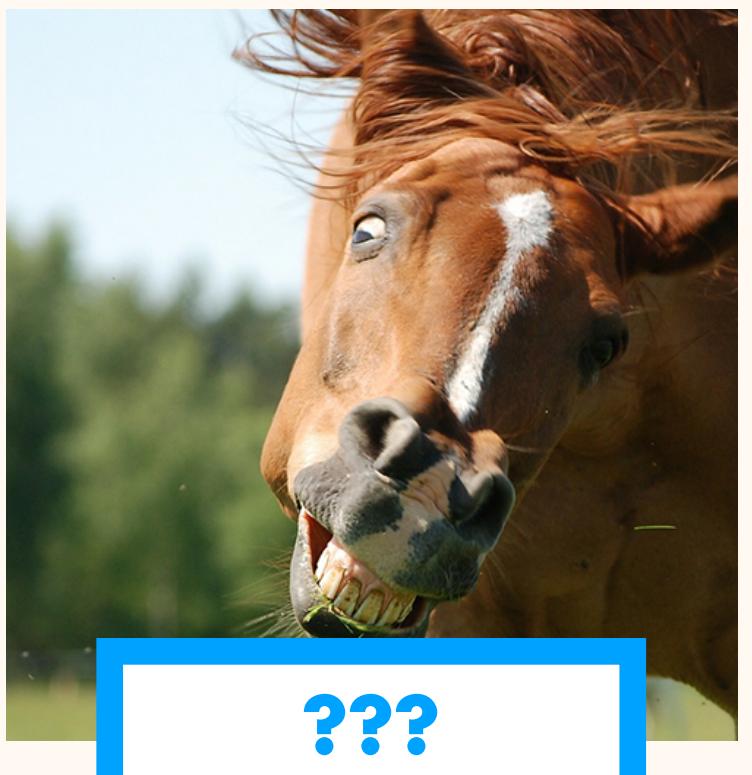
**The student
and the teacher
improve each other
in a **virtuous** cycle.**



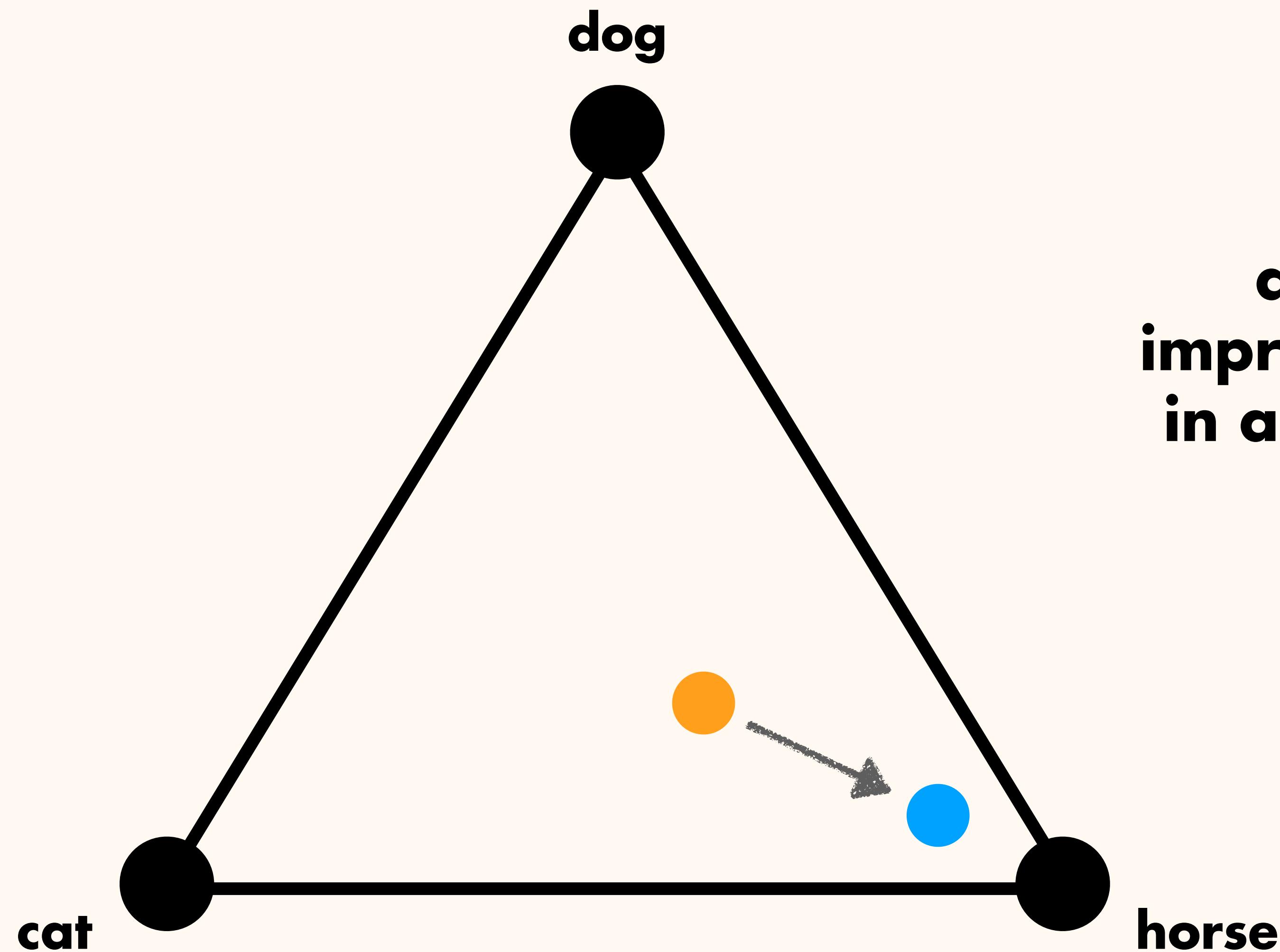
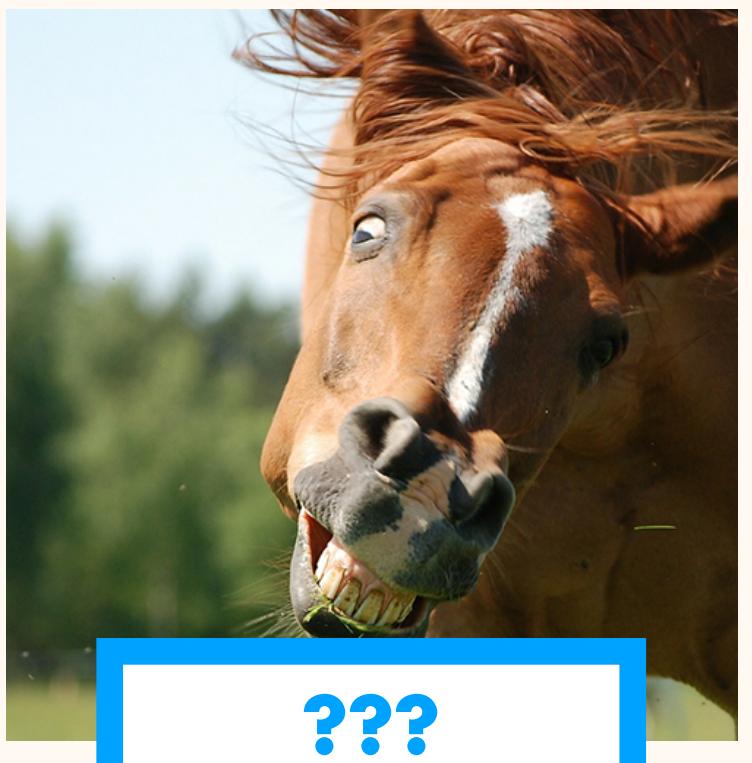
**The student
and the teacher
improve each other
in a virtuous cycle.**



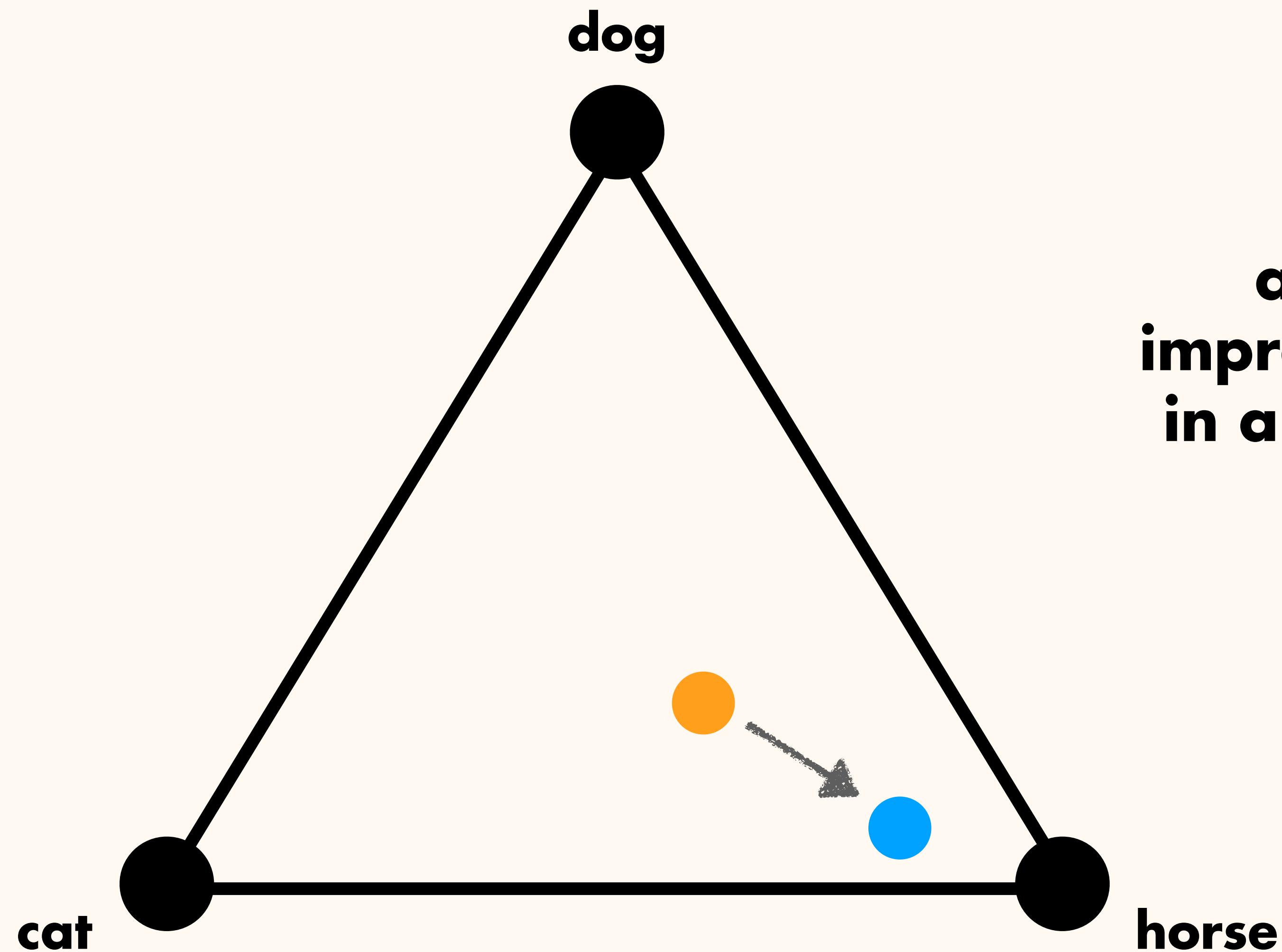
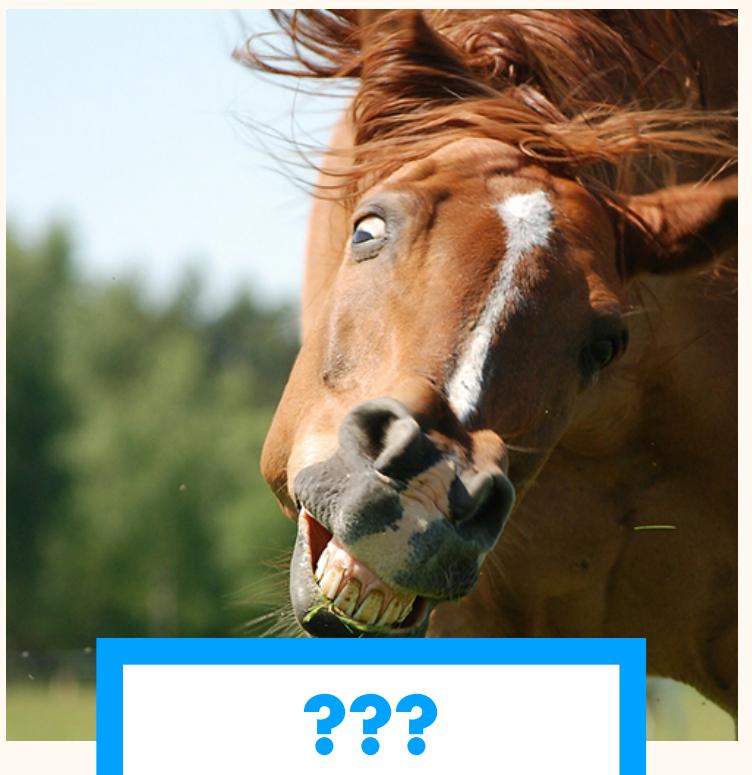
**The student
and the teacher
improve each other
in a virtuous cycle.**



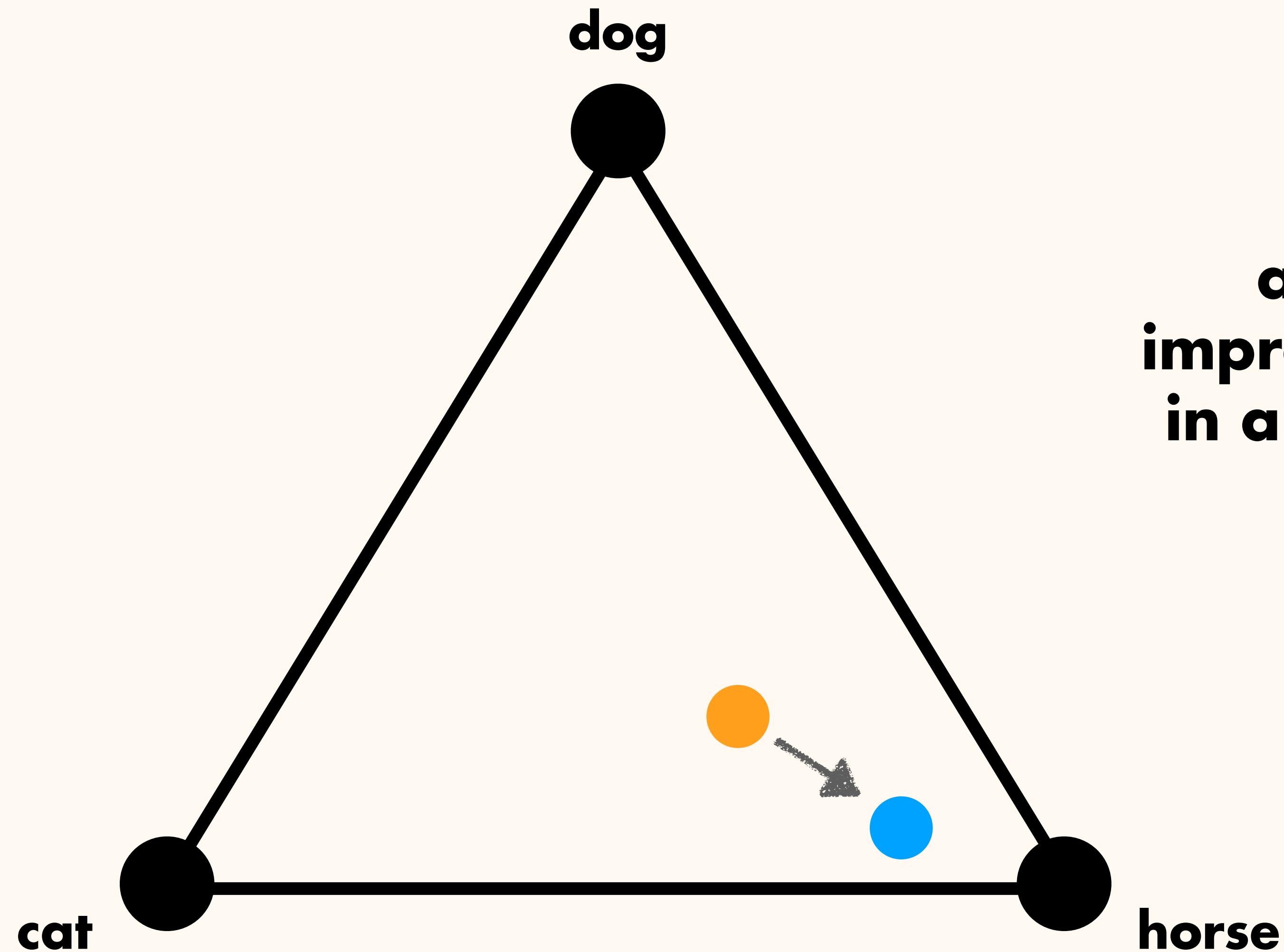
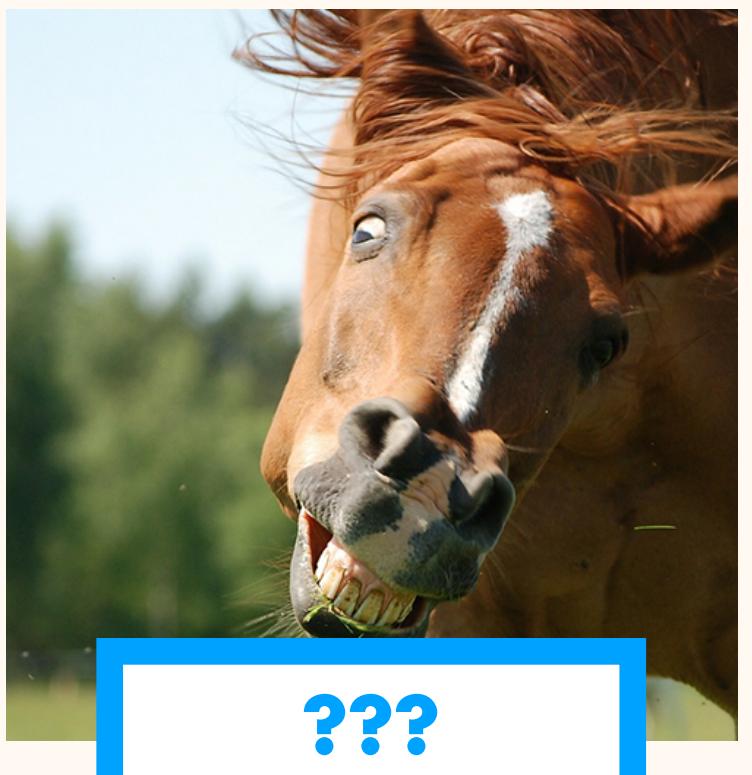
**The student
and the teacher
improve each other
in a virtuous cycle.**



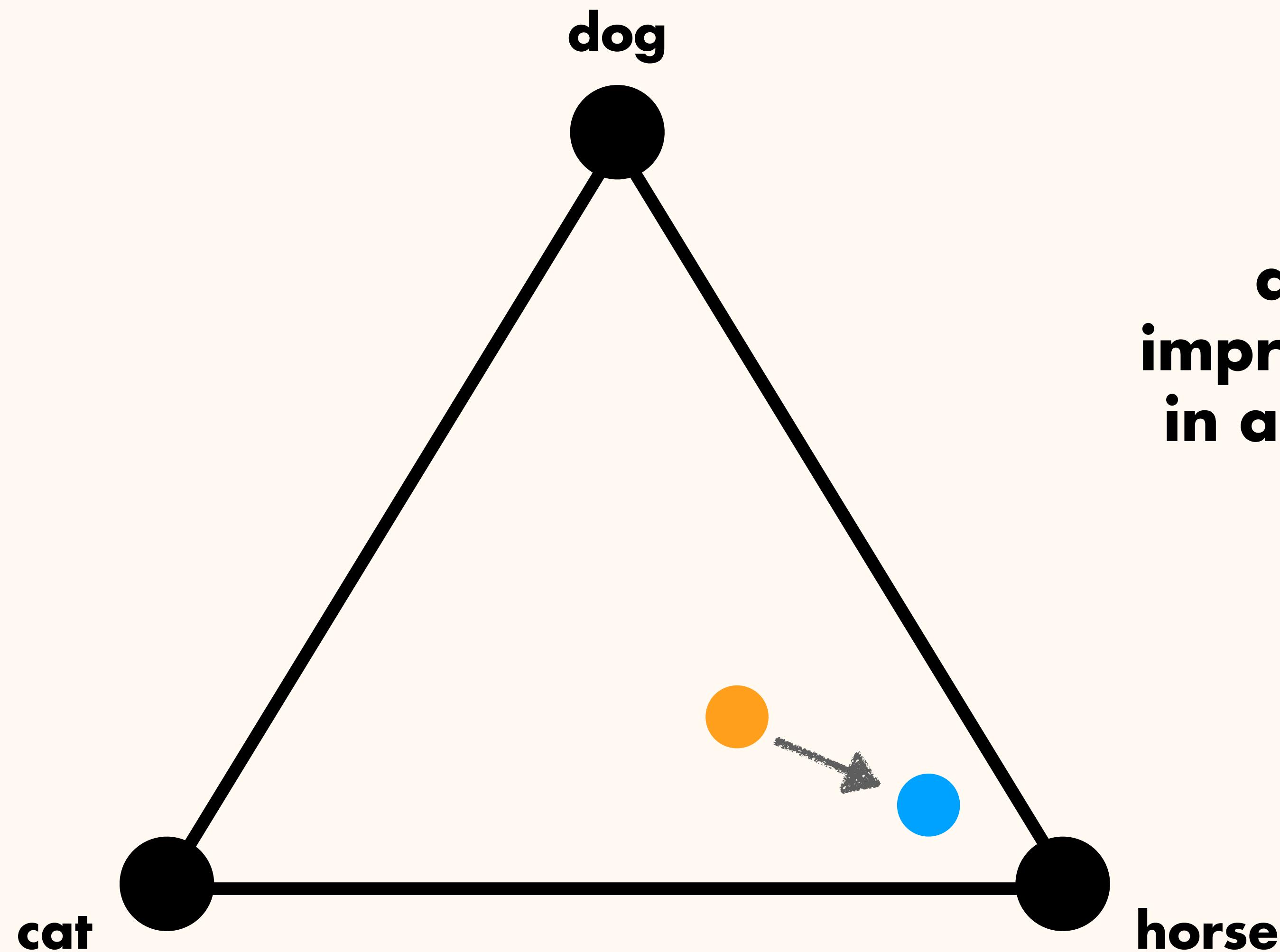
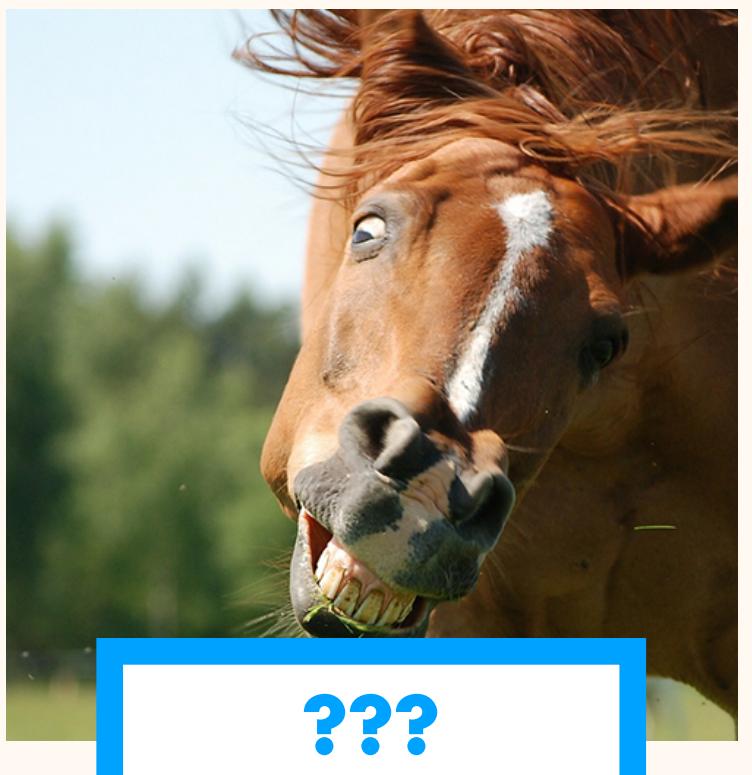
**The student
and the teacher
improve each other
in a virtuous cycle.**



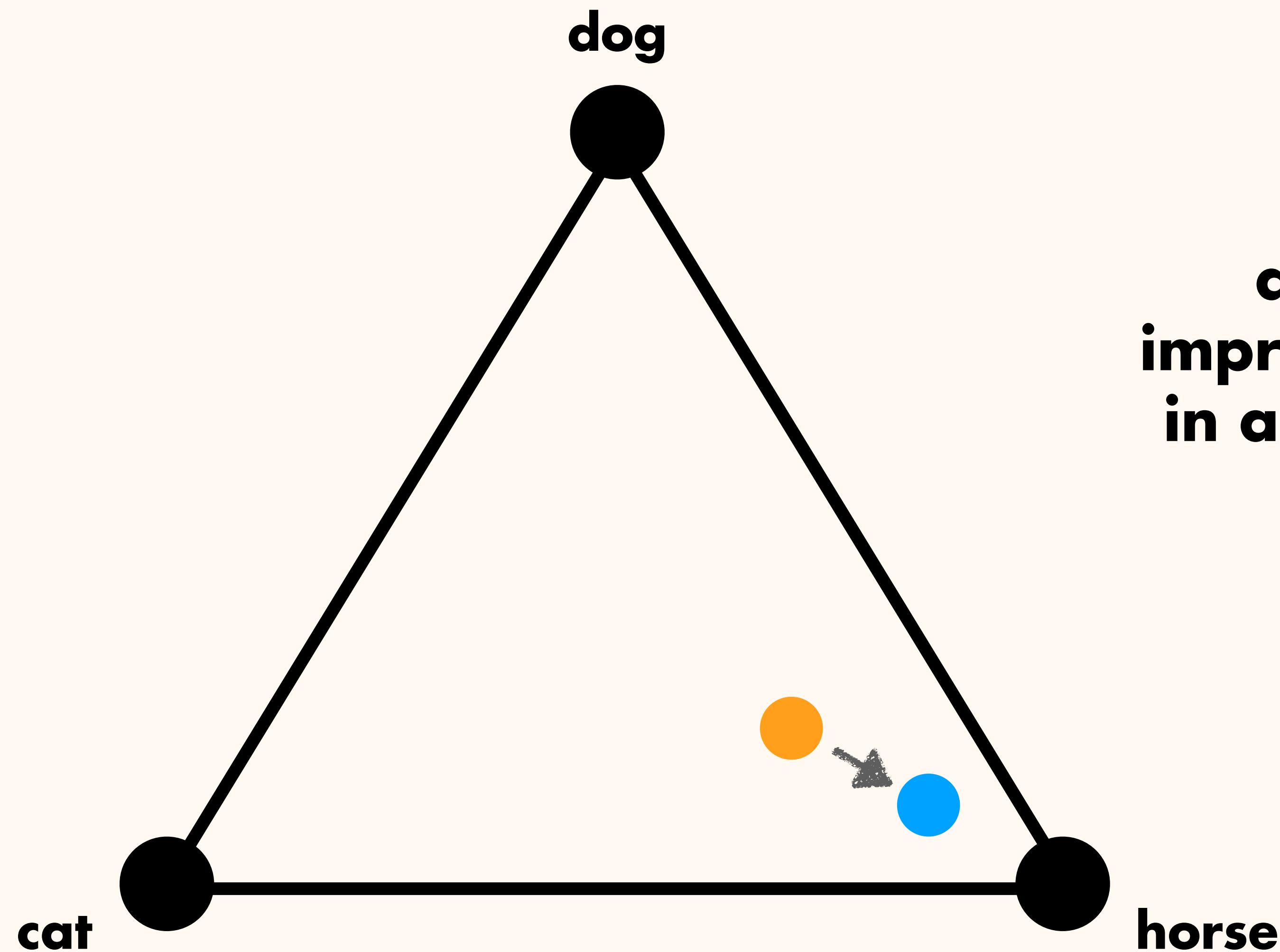
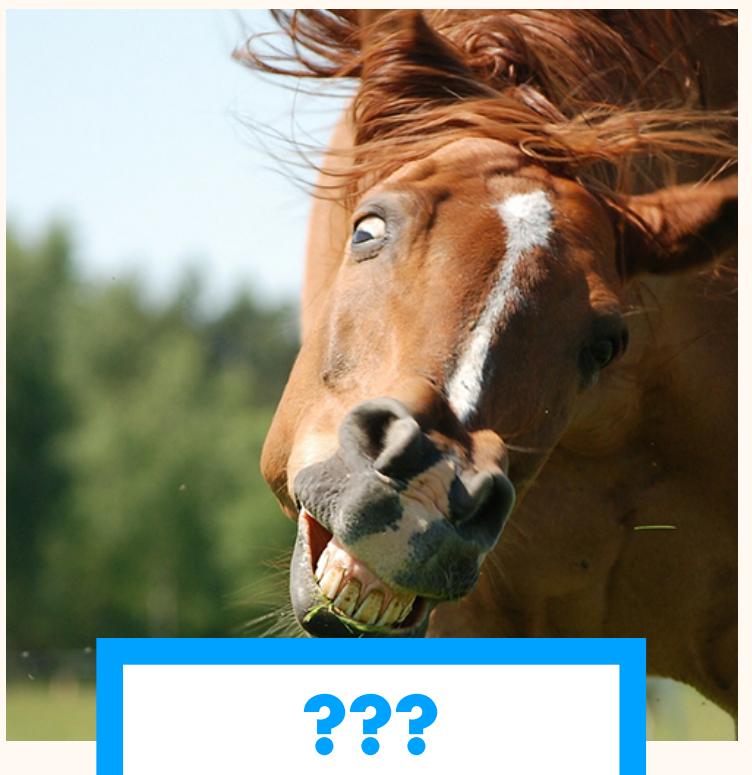
**The student
and the teacher
improve each other
in a virtuous cycle.**



**The student
and the teacher
improve each other
in a virtuous cycle.**



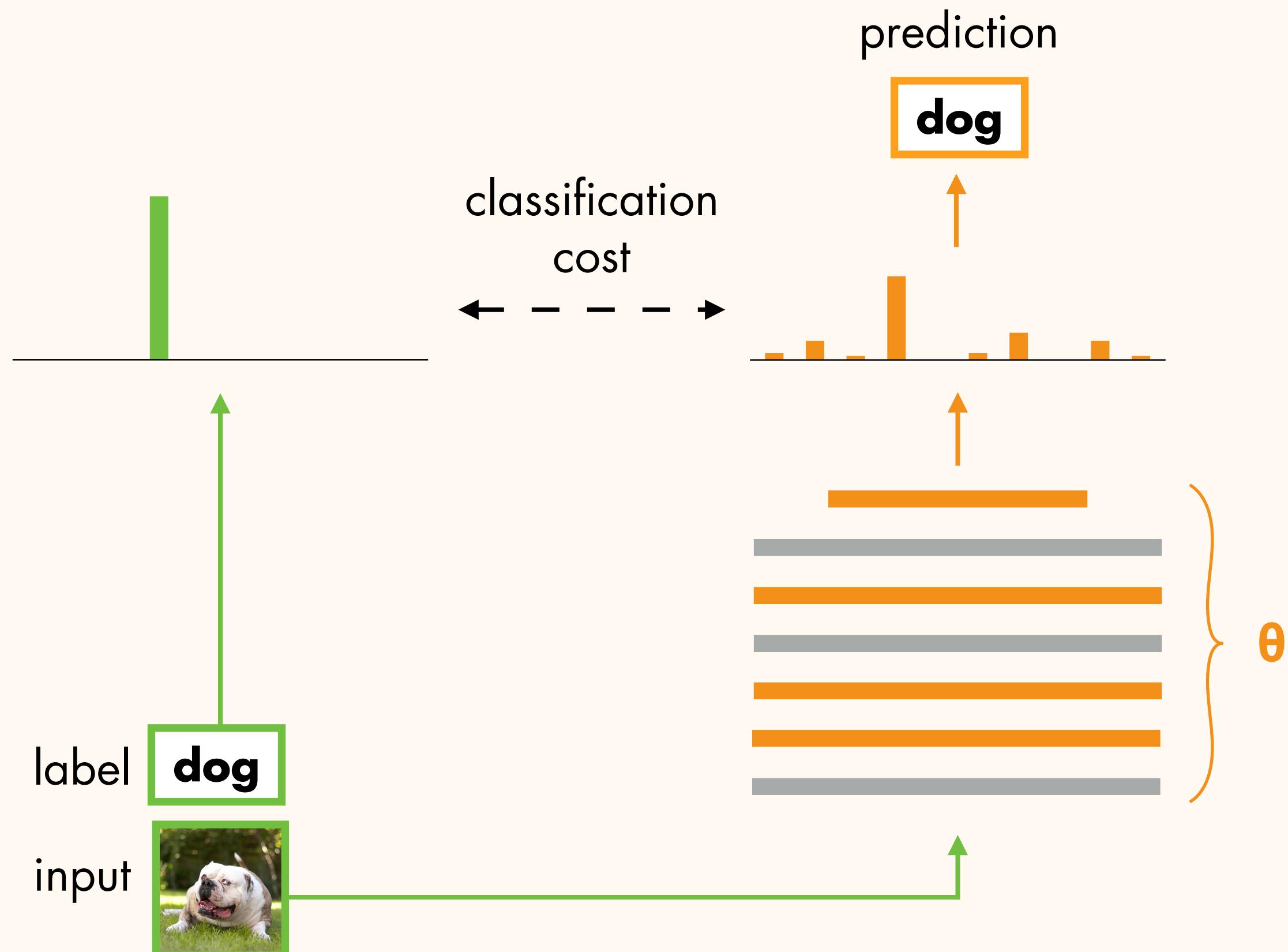
**The student
and the teacher
improve each other
in a virtuous cycle.**



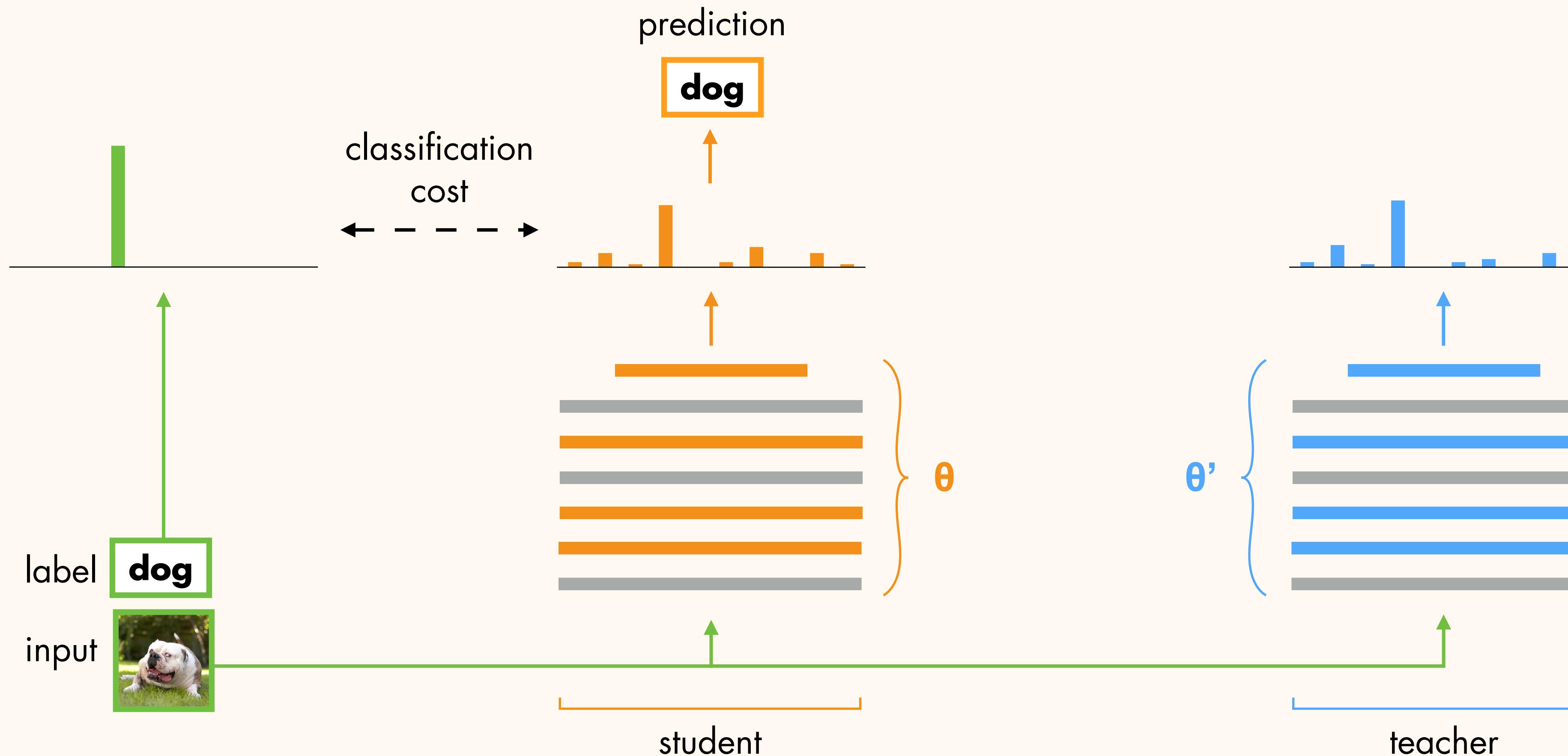
**The student
and the teacher
improve each other
in a virtuous cycle.**

HOW TO START USING THE MEAN TEACHER

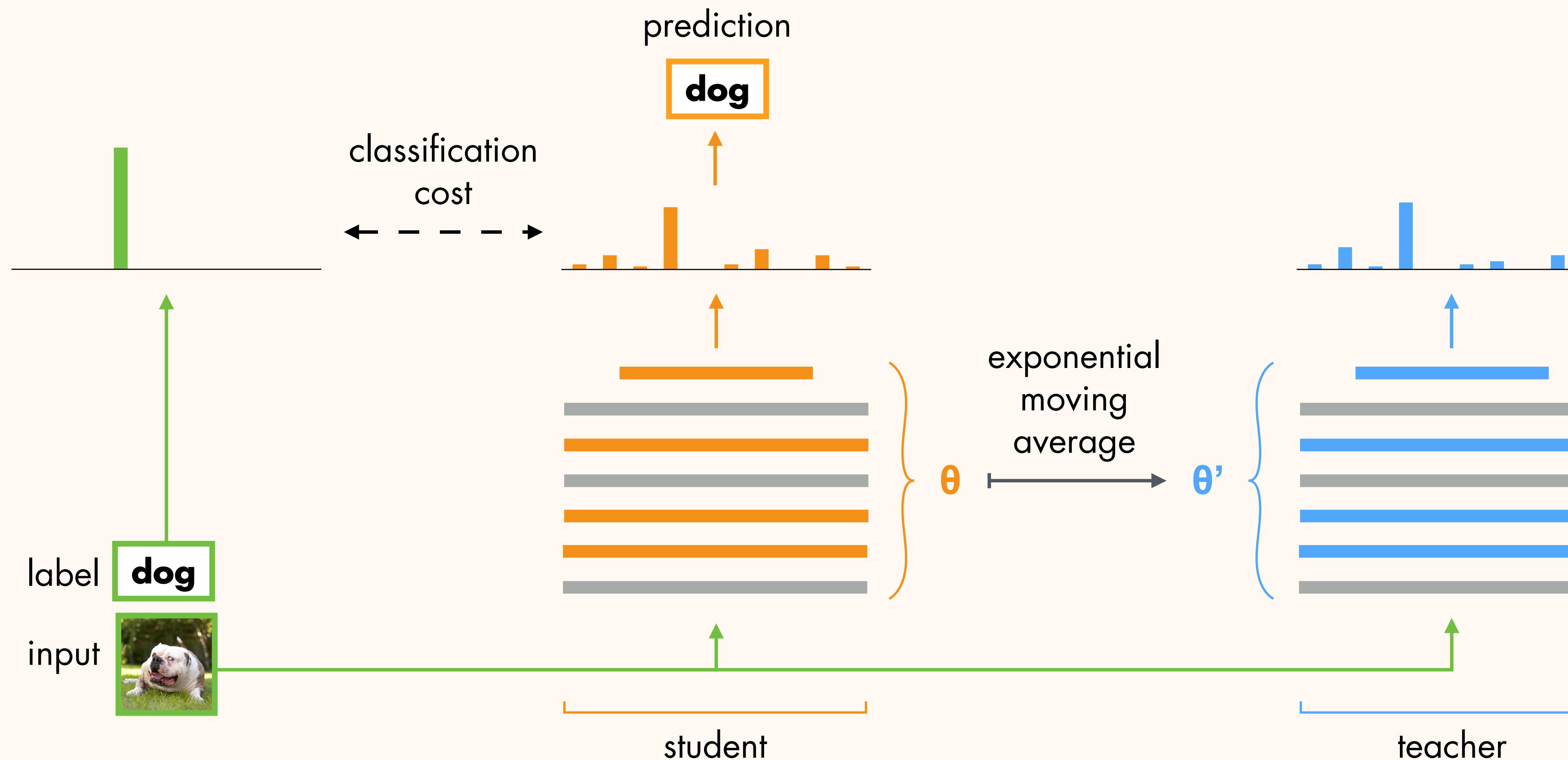
Take a supervised model.



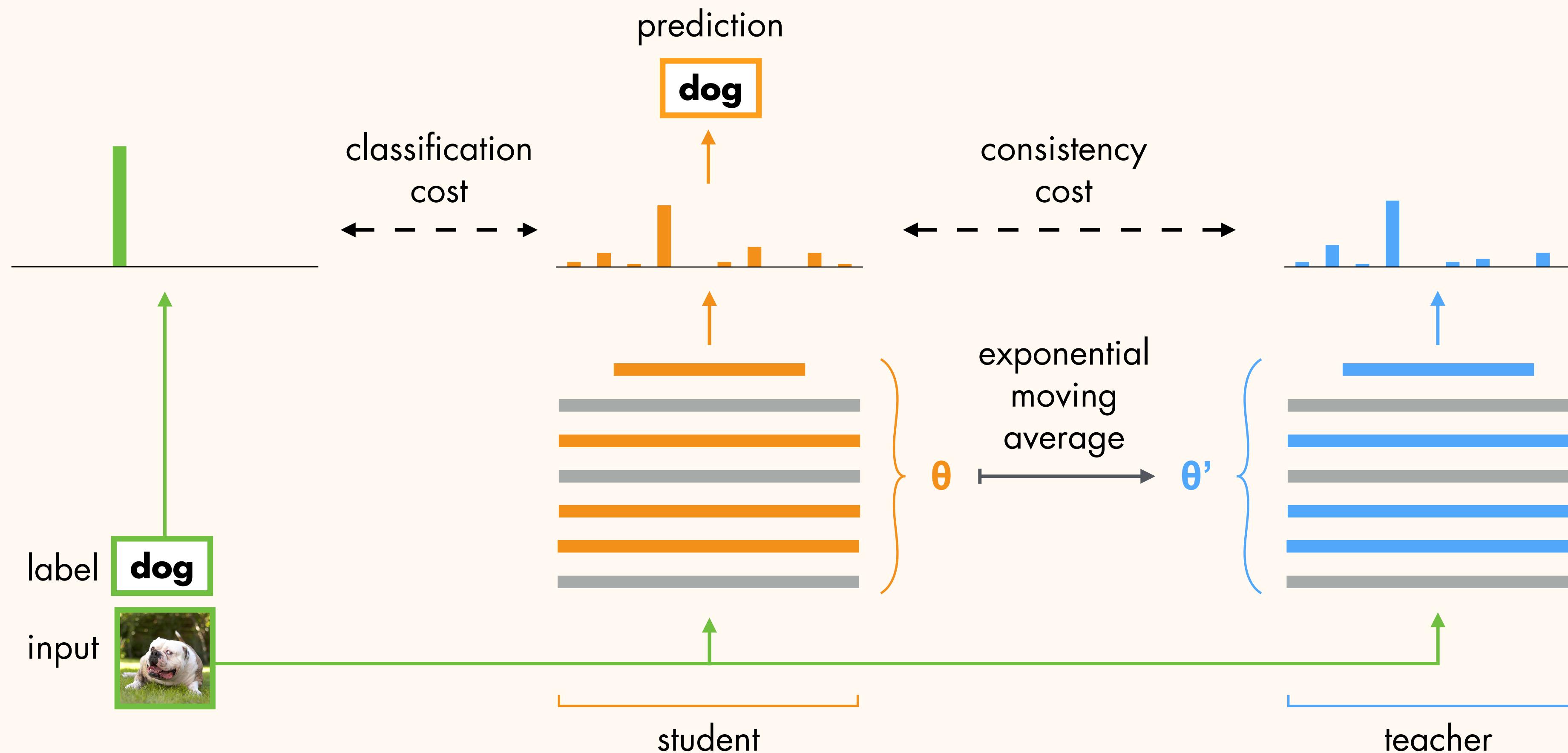
Make a copy of it.



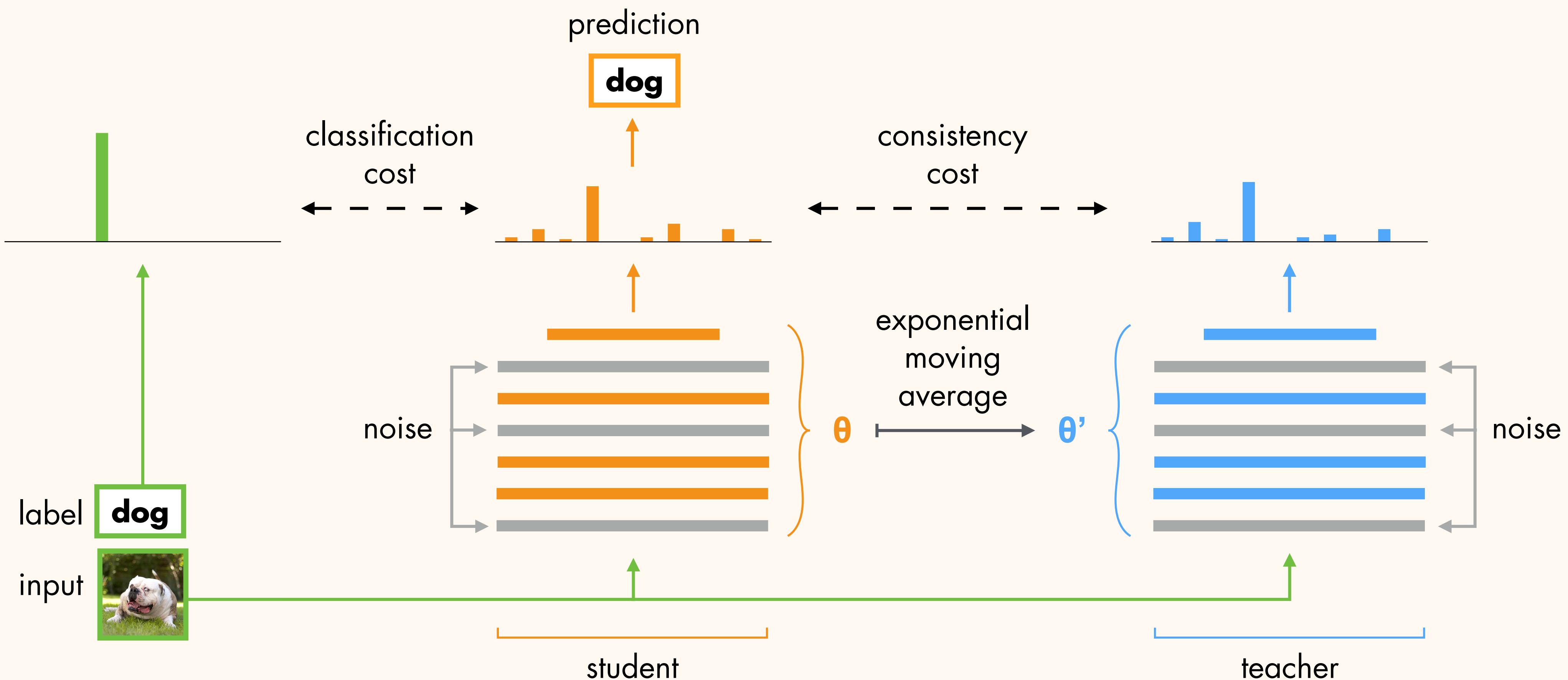
Update teacher weights after each training step.



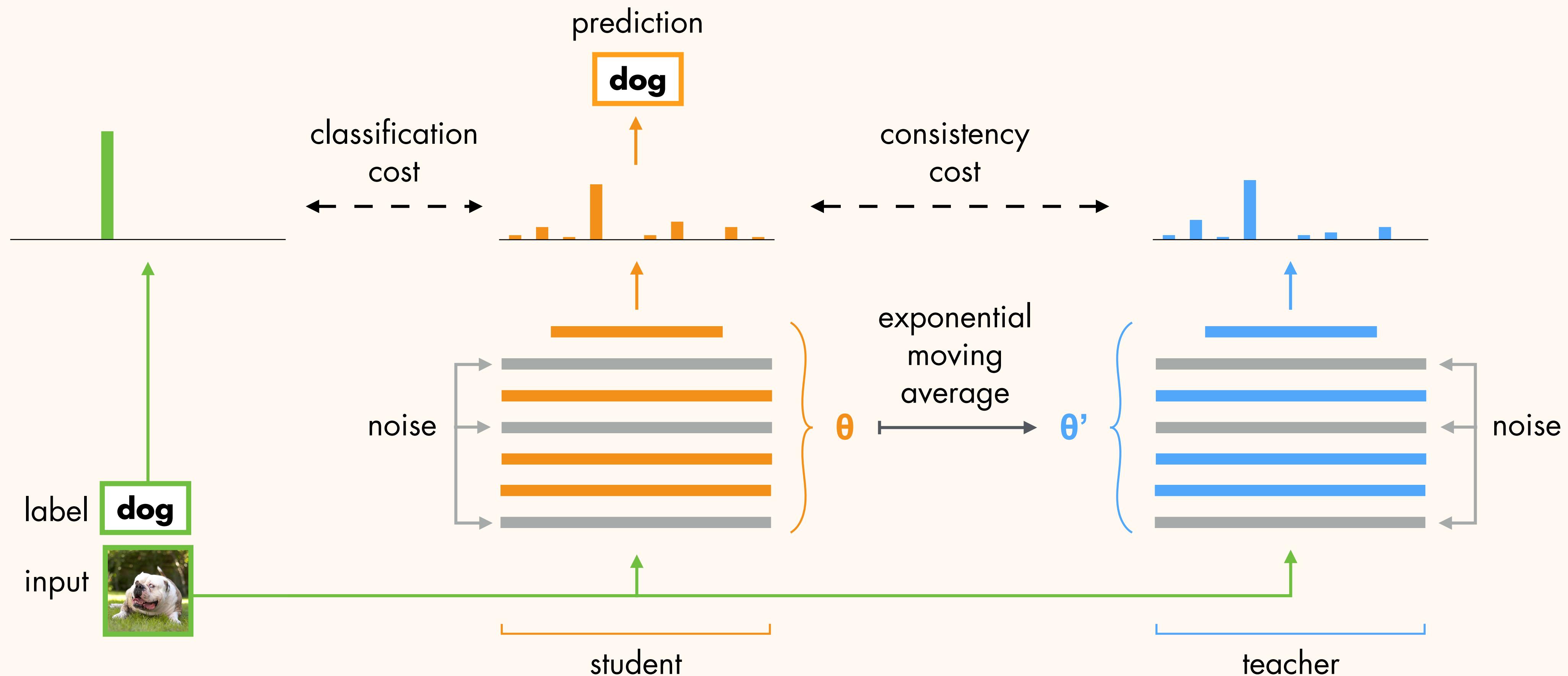
Add a cost between the two predictions.



Maybe add some noise.



Start using it for semi-supervised learning.



Optimization criteria-Mean Teacher

- Classification cost - Cross-entropy loss between student softmax output and one-hot labels
- Consistency cost, $J(\theta)$ - MSE between the student and teacher softmax outputs
 - Approximate $J(\theta)$ by sampling η and η' at each training step with SGD
$$J(\theta) = E_{x, \eta', \eta} \| f(x, \theta', \eta') - f(x, \theta, \eta) \|^2$$
$$\theta_{t'} = \alpha \theta_{t-1} + (1 - \alpha) \theta_t$$
- Total cost : Weighted average of classification cost and consistency cost
 - Weight for classification : Expected number of labeled examples per mini-batch

RESULTS

Experimental setup

- 13-layer ConvNet architecture
 - Applied mean-only batch-normalization and weight normalization on convolutional and softmax layers
 - Used leaky-ReLU with $\alpha = 0.1$
- 12 block(26-layer) ResNet with Shake-shake regularization
- 3 types of Noise
 - random translations and horizontal flips
 - Gaussian noise on the input layer
 - Dropout applied withing the network
- Datasets
 - Street View House Numbers(SVHN) : 73257 training, 26032 test
 - CIFAR-10 : 50000 training, 10000 test

Table 6: The convolutional network architecture we used in the experiments.

Layer	Hyperparameters
Input	32×32 RGB image
Translation	Randomly $\{\Delta x, \Delta y\} \sim [-2, 2]$
Horizontal flip ^a	Randomly $p = 0.5$
Gaussian noise	$\sigma = 0.15$
Convolutional	128 filters, 3×3 , <i>same</i> padding
Convolutional	128 filters, 3×3 , <i>same</i> padding
Convolutional	128 filters, 3×3 , <i>same</i> padding
Pooling	Maxpool 2×2
Dropout	$p = 0.5$
Convolutional	256 filters, 3×3 , <i>same</i> padding
Convolutional	256 filters, 3×3 , <i>same</i> padding
Convolutional	256 filters, 3×3 , <i>same</i> padding
Pooling	Maxpool 2×2
Dropout	$p = 0.5$
Convolutional	512 filters, 3×3 , <i>valid</i> padding
Convolutional	256 filters, 1×1 , <i>same</i> padding
Convolutional	128 filters, 1×1 , <i>same</i> padding
Pooling	Average pool ($6 \times 6 \rightarrow 1 \times 1$ pixels)
Softmax	Fully connected $128 \rightarrow 10$

^a Not applied on SVHN experiments

Baselines

Apply unsupervised loss component to both labelled and unlabeled examples

- Ladder networks - Pi model (Laine & Aila, 2016) : Ensembles predictions of the model under 2 different perturbations of input data and different dropout conditions
- Temporal Ensembling - Ensemble predictions of a model at different timesteps

Baelines	$J(\theta) = E_{x, \eta', \eta} \ f(x, \theta', \eta') - f(x, \theta, \eta) \ ^2$
Π model	$\theta' = \theta$
Temporal Ensembling	weighted average of successive predictions (θ' is not optimized)
Mean Teacher (Our)	Expected moving average of of successive model weights (θ' is not optimized)

Table 5: Error rate percentage on SVHN and CIFAR-10 over 10 runs, including the results without input augmentation. We use exponential moving average weights in the evaluation of all our models. All the comparison methods use a 13-layer ConvNet architecture similar to ours and augmentation similar to ours, except GAN, which does not use augmentation.

SVHN	250 labels	500 labels	1000 labels	all labels ^a
GAN ^b		18.44 \pm 4.8	8.11 \pm 1.3	
II model ^c		6.65 \pm 0.53	4.82 \pm 0.17	2.54 \pm 0.04
Temporal Ensembling ^c		5.12 \pm 0.13	4.42 \pm 0.16	2.74 \pm 0.06
VAT+EntMin ^d			3.86	
<hr/>				
Ours				
Supervised-only ^e	27.77 \pm 3.18	16.88 \pm 1.30	12.32 \pm 0.95	2.75 \pm 0.10
II model	9.69 \pm 0.92	6.83 \pm 0.66	4.95 \pm 0.26	2.50 \pm 0.07
Mean Teacher	4.35 \pm 0.50	4.18 \pm 0.27	3.95 \pm 0.19	2.50 \pm 0.05
<hr/>				
Without augmentation				
Supervised-only ^e	36.26 \pm 3.83	19.68 \pm 1.03	14.15 \pm 0.87	3.04 \pm 0.04
II model	10.36 \pm 0.94	7.01 \pm 0.29	5.73 \pm 0.16	2.75 \pm 0.08
Mean Teacher	5.85 \pm 0.62	5.45 \pm 0.14	5.21 \pm 0.21	2.77 \pm 0.09
<hr/>				
CIFAR-10	1000 labels	2000 labels	4000 labels	all labels ^a
GAN ^b			18.63 \pm 2.32	
II model ^c			12.36 \pm 0.31	5.56 \pm 0.10
Temporal Ensembling ^c			12.16 \pm 0.31	5.60 \pm 0.10
VAT+EntMin ^d			10.55	
<hr/>				
Ours				
Supervised-only ^e	46.43 \pm 1.21	33.94 \pm 0.73	20.66 \pm 0.57	5.82 \pm 0.15
II model	27.36 \pm 1.20	18.02 \pm 0.60	13.20 \pm 0.27	6.06 \pm 0.11
Mean Teacher	21.55 \pm 1.48	15.73 \pm 0.31	12.31 \pm 0.28	5.94 \pm 0.15
Mean Teacher ResNet	10.08 \pm 0.41		6.28 \pm 0.15	
<hr/>				
Without augmentation				
Supervised-only ^e	48.38 \pm 1.07	36.07 \pm 0.90	24.47 \pm 0.50	7.43 \pm 0.06
II model	32.18 \pm 1.33	23.92 \pm 1.07	17.08 \pm 0.32	7.00 \pm 0.20
Mean Teacher	30.62 \pm 1.13	23.14 \pm 0.46	17.74 \pm 0.30	7.21 \pm 0.24

^a 4 runs

^b Salimans et al. [25]

^c Laine & Aila [13]

^d Miyato et al. [16]

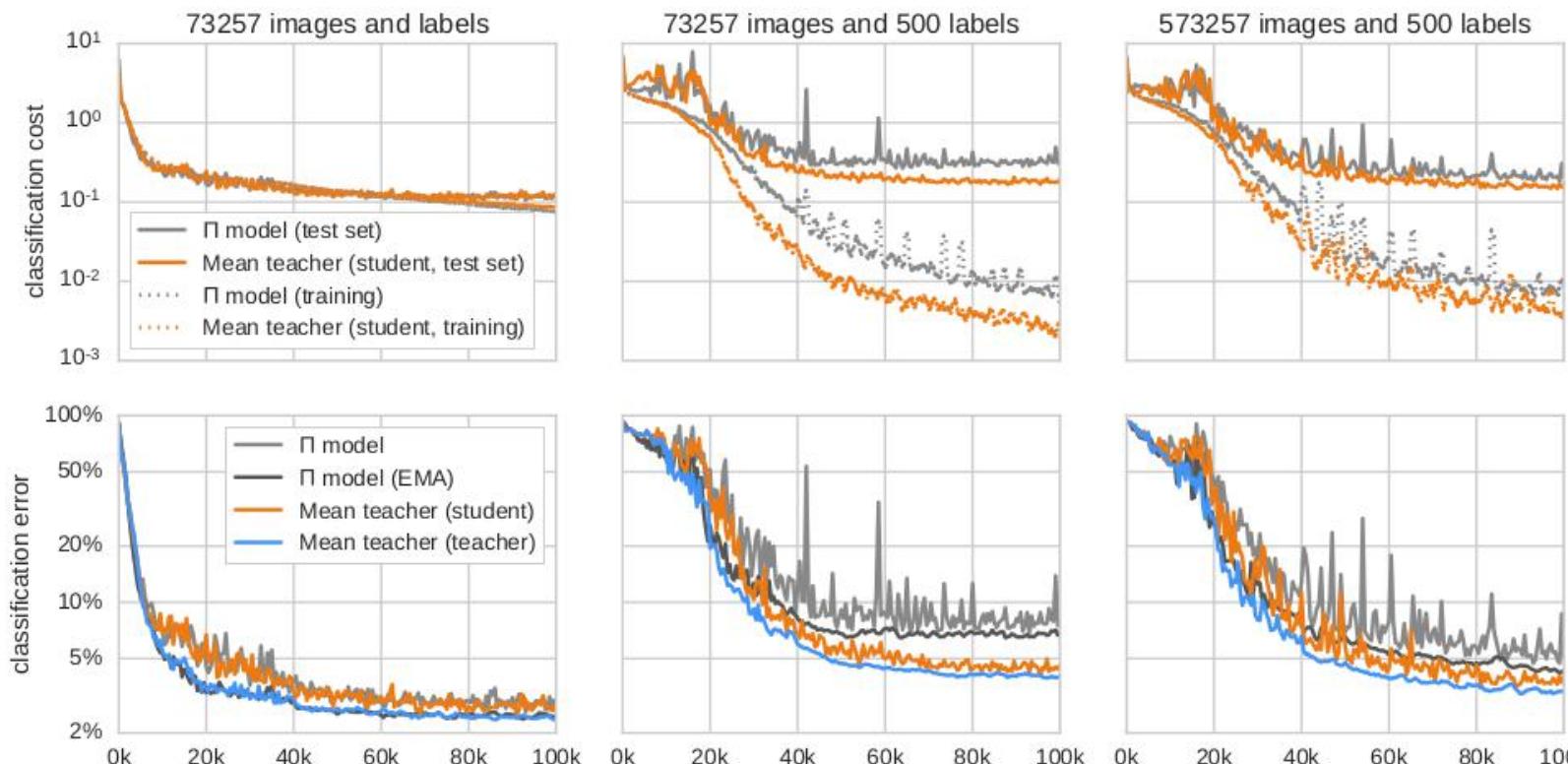
^e Only labeled examples and only classification cost

1. MT performs better when the number of labelled examples is less

2. VAT+EntMin (Miyato et. al., 2017) outperforms MT in certain cases

Table 3: Error percentage over 10 runs on SVHN with extra unlabeled training data.

	500 labels 73257 images	500 labels 173257 images	500 labels 573257 images
Π model (ours)	6.83 ± 0.66	4.49 ± 0.27	3.26 ± 0.14
Mean Teacher	4.18 ± 0.27	3.02 ± 0.16	2.46 ± 0.06



(right column) : Pi-model keeps improving for longer. MT learns faster and eventually converges to a better result.

Fig : Smoothened classification cost (top) and classification error (bottom) of Mean Teacher and our baseline Π model on SVHN over the first 100000 training steps. In the upper row, the training classification costs are measured using only labeled data.

Other experiments

- Instead of using MSE(Laine & Aila, 2016) for calculating consistency cost, we use KL-divergence

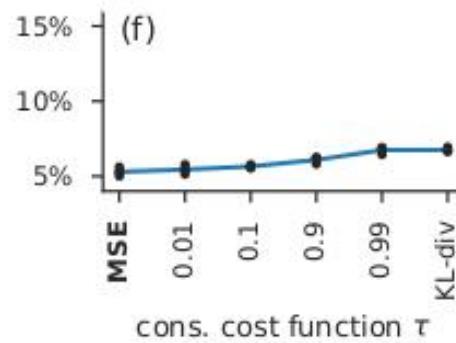


Figure 5: Copy of Figure 4(f) in the main text. Validation error on 250-label SVHN over four runs and their mean, when varying the consistency cost shape hyperparameter τ between mean squared error ($\tau = 0$) and KL-divergence ($\tau = 1$).

MSE performs better than KL-divergence or C_T for different values of T

Related works

- Regularize intermediate representations effectively in DL
 - adversarial training, dropout, DropConnect, Stochastic Depth, Swapout
- Training the model predictions to be consistent to perturbation
 - Denoising Source Separation framework, Ladder Networks(T-model, Pi-model), Virtual Adversarial training
- Teacher model training a student model
 - Model compression and distillation
- Consistency regularization : Like label propagation, but somewhat different !
 - LP requires predefined distance metric
 - Employ a learned distance metric and aid the network to learn a better distance metric

Conclusion

- Mean Teacher(MT) helps when the labels are scarce
 - When using only 500 labels, MT learns faster and continues training for much longer than pi-model
 - In all-labeled case, Mean teacher and pi model perform identically
- MT uses unlabeled training data more efficiently than Pi model
 - With 500k examples extra unlabeled examples, MT learns faster and end eventually converges to a better result
- Works with large datasets and on-line learning
 - Uses moving average of model weights or parameters instead of label predictions

WANT TO KNOW MORE?

PAPER

Mean Teachers Are Better Role Models

Weight-averaged consistency targets improve semi-supervised deep learning results

[arxiv:1703.01780](https://arxiv.org/abs/1703.01780)

POSTER

Today 06:30 – 10:30 PM
@ Pacific Ballroom #13

SOURCE CODE

[http://github.com/CuriousAI/mean-teacher](https://github.com/CuriousAI/mean-teacher)