

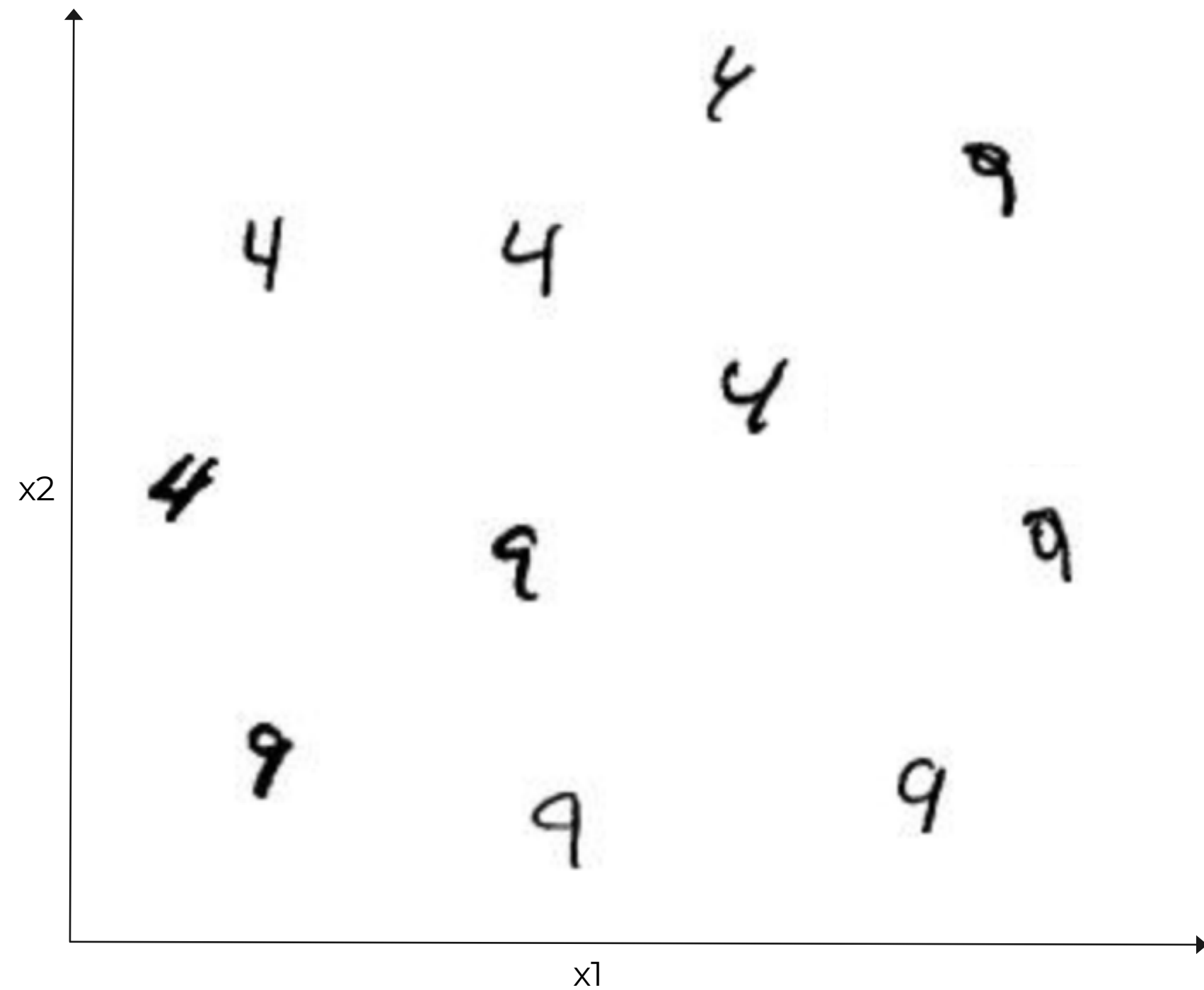
# Stress Testing Models to Understand their Decisions

---

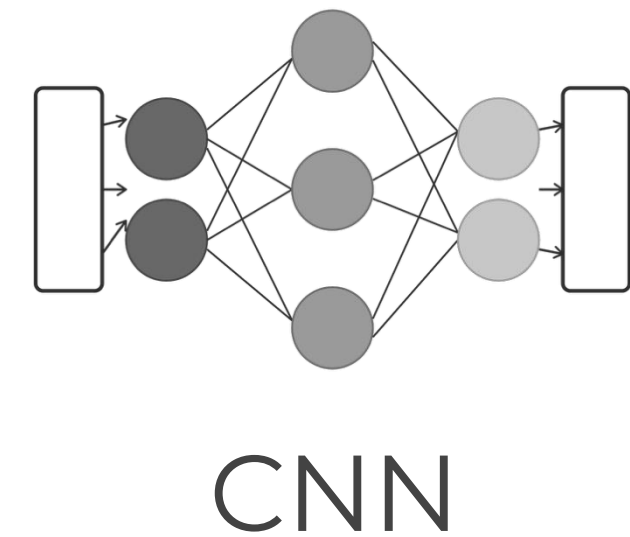
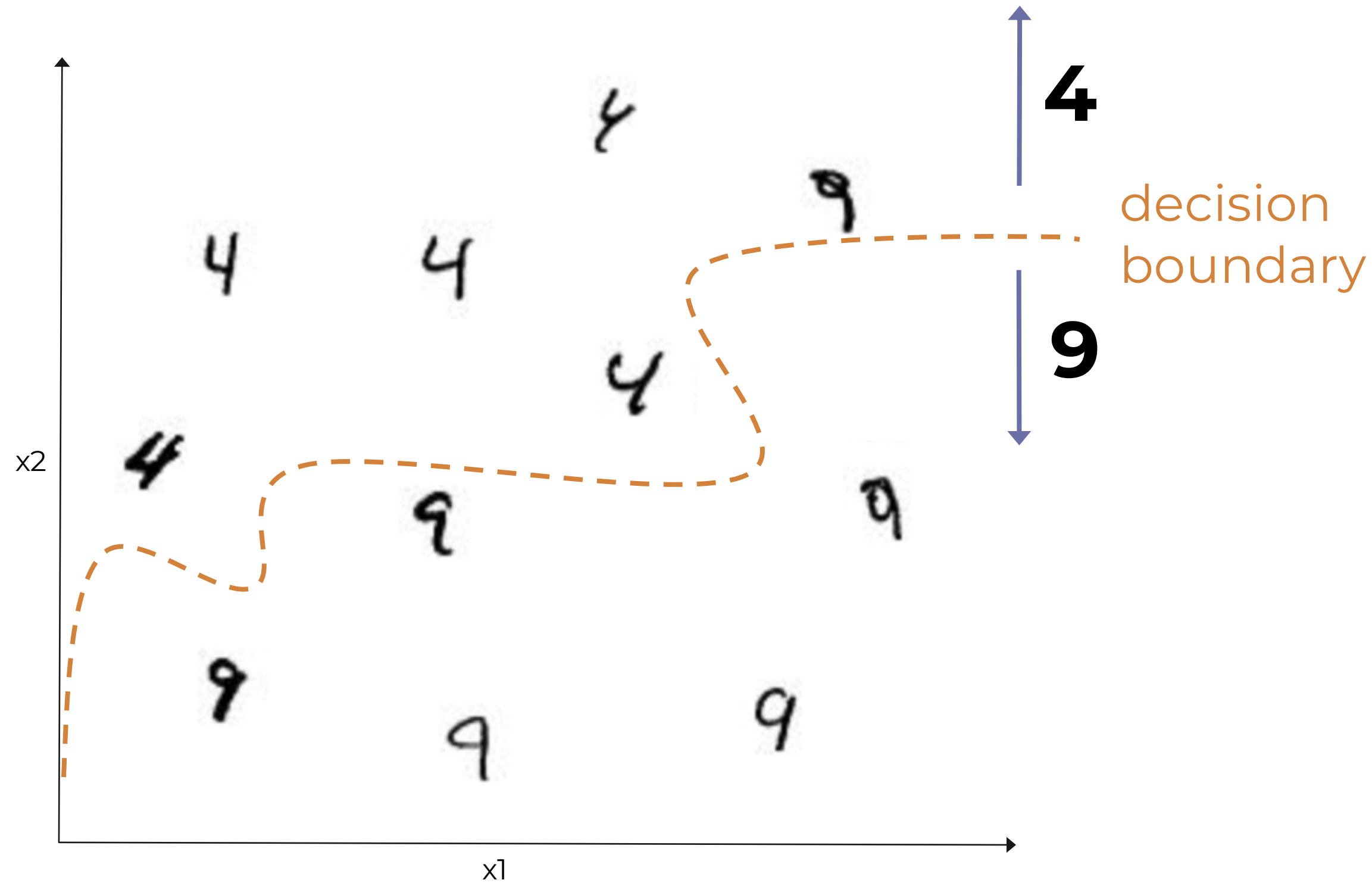
**Inês Gomes**

Data Makers Fest, 5 May 2026

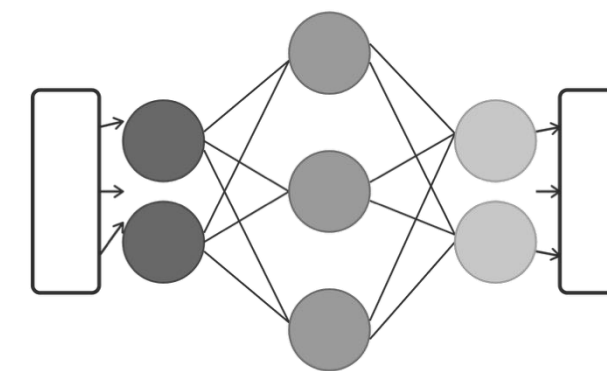
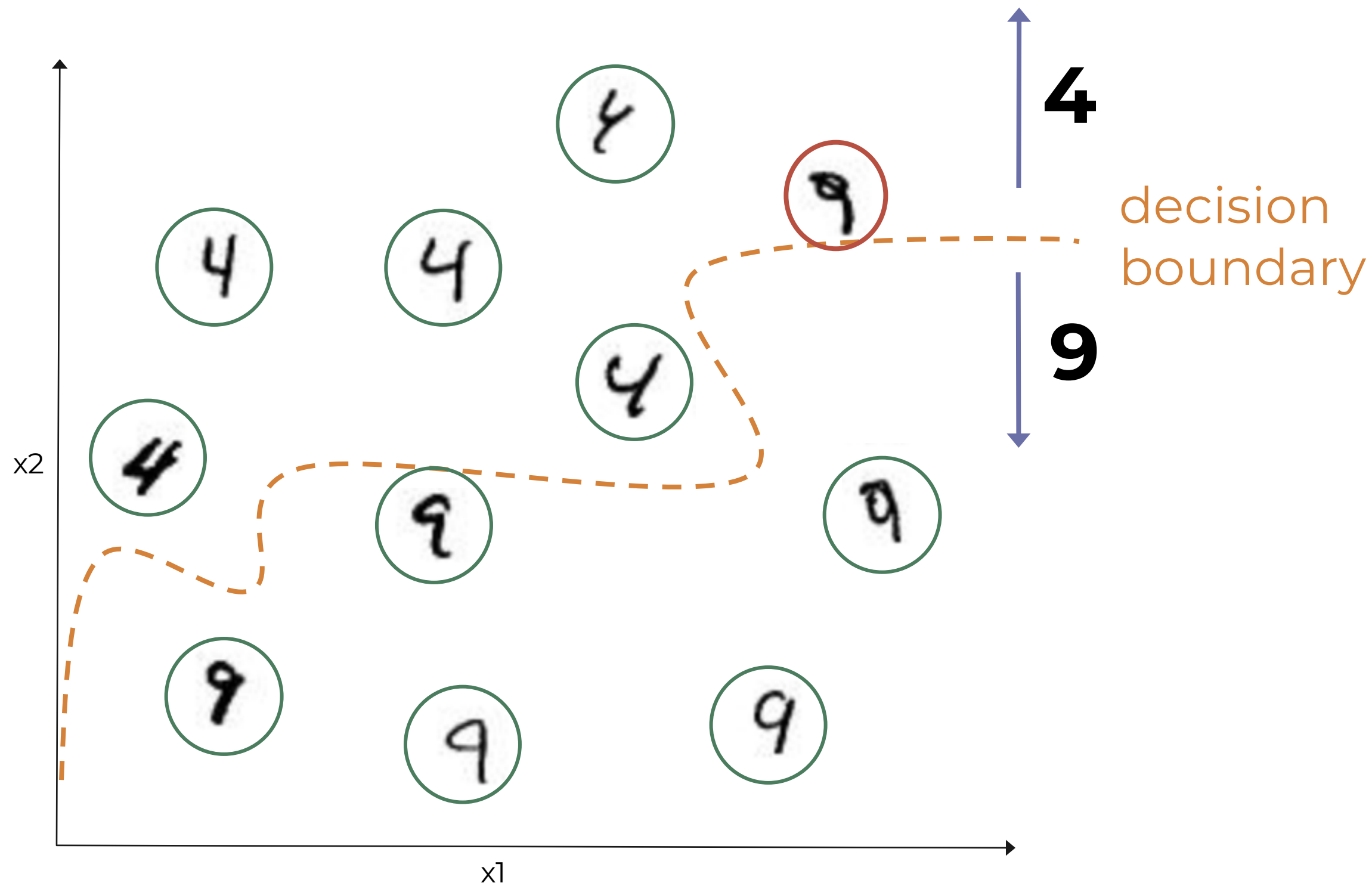
# Example: distinguish handwritten digits



# Example: distinguish handwritten digits



# Example: distinguish handwritten digits



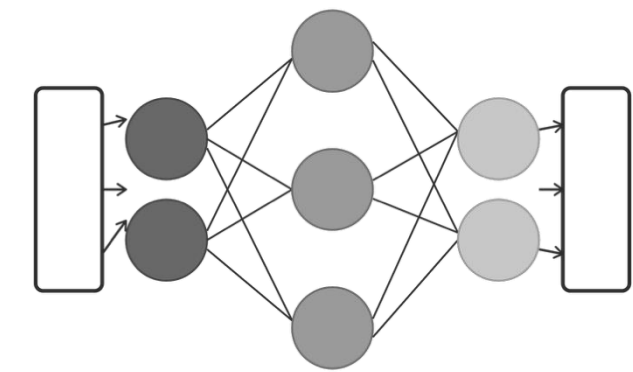
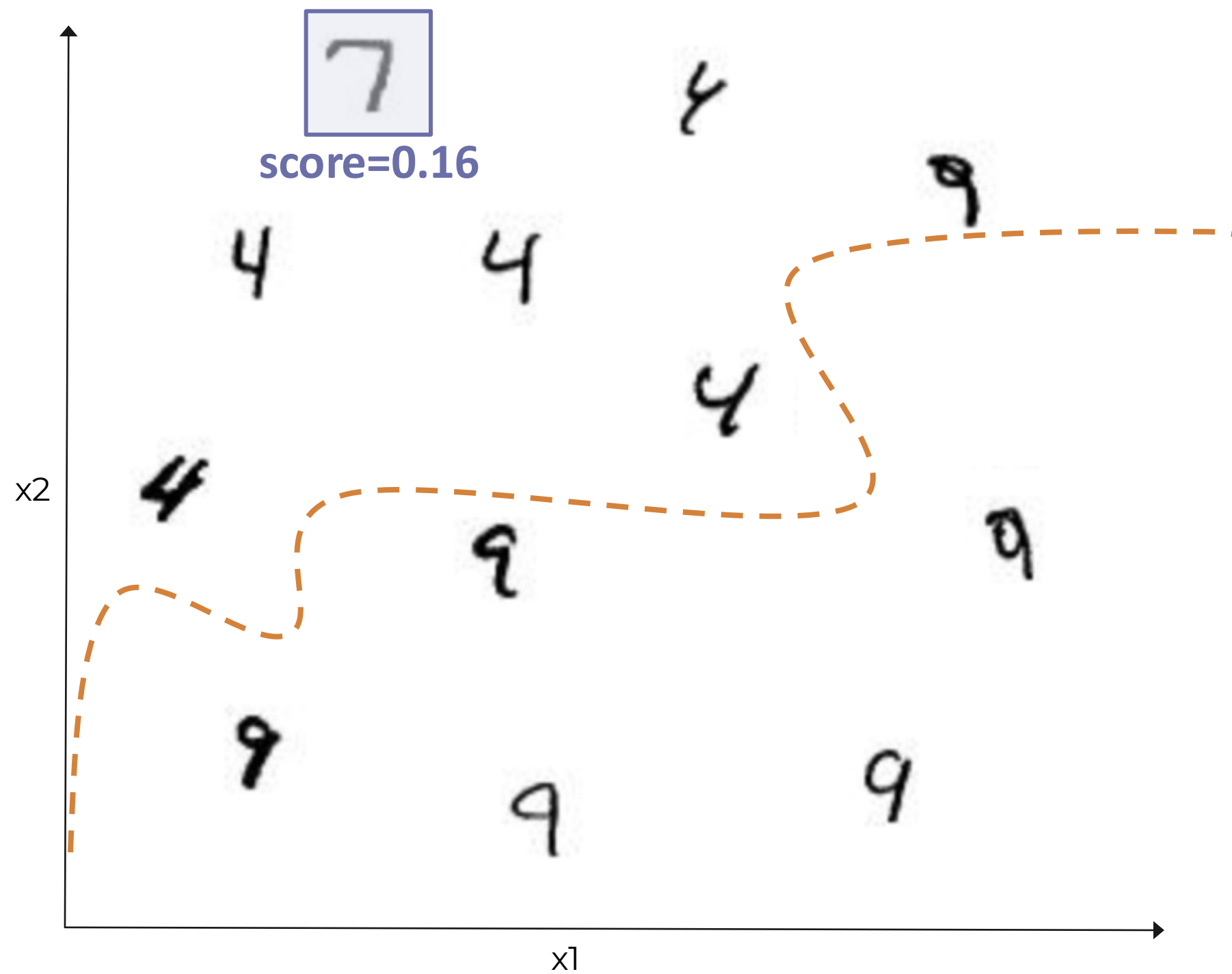
CNN

accuracy = 98%

BUT...

does 98% accuracy tell the  
whole story?

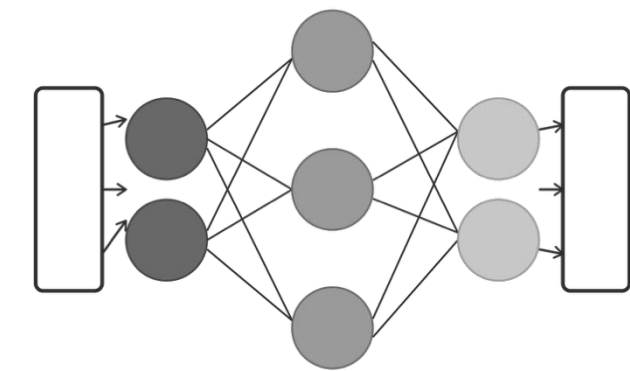
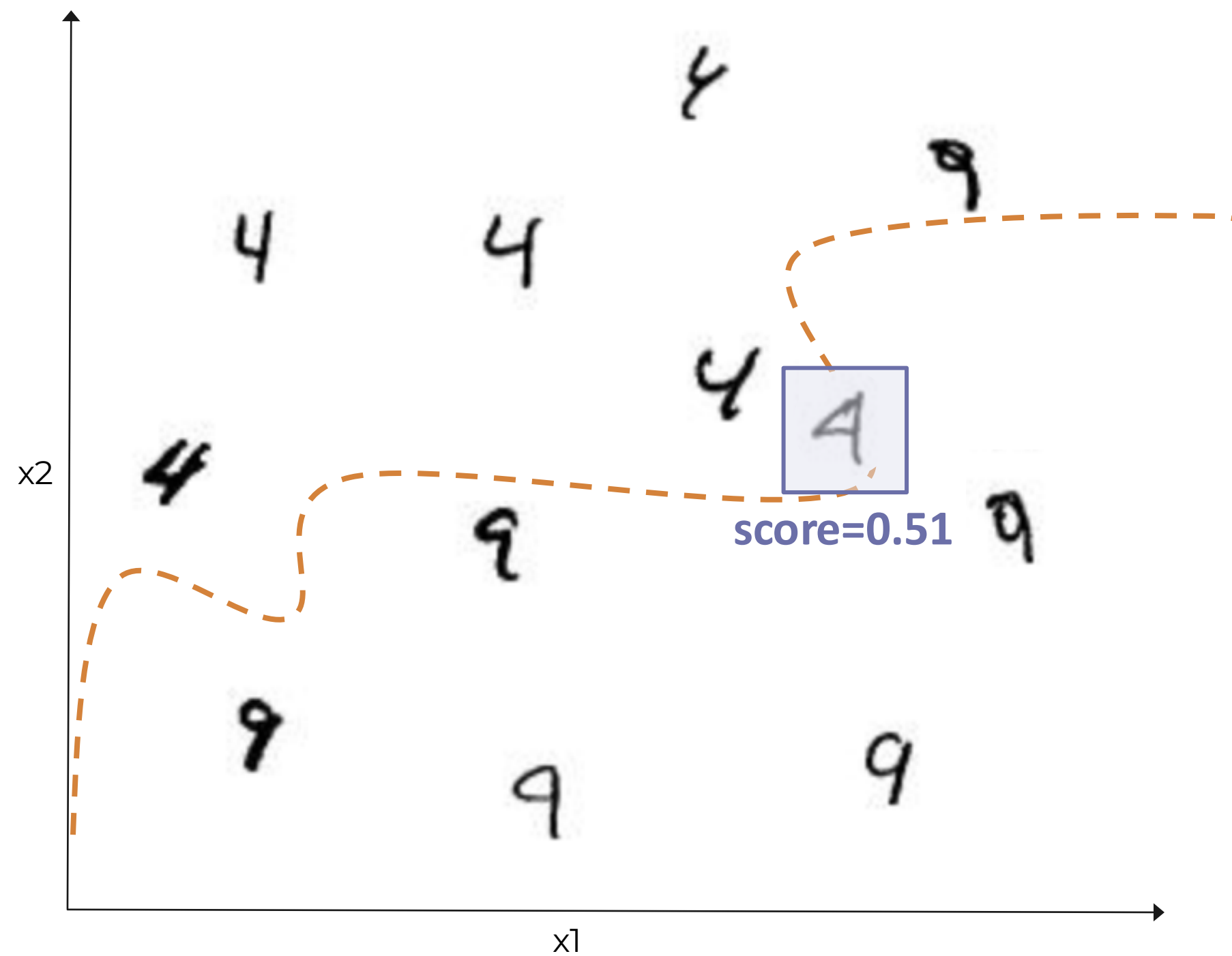
# Out-of-distribution: visually unlike any training sample



CNN

accuracy = 98%

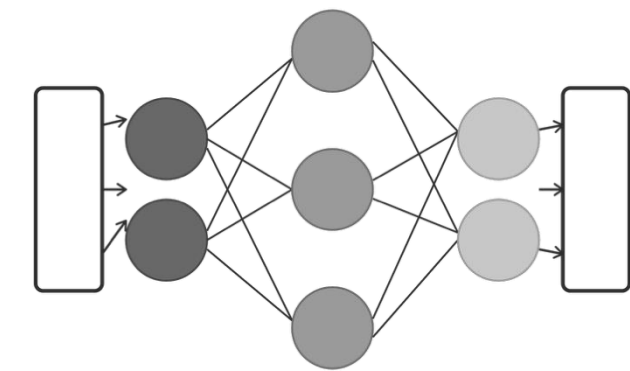
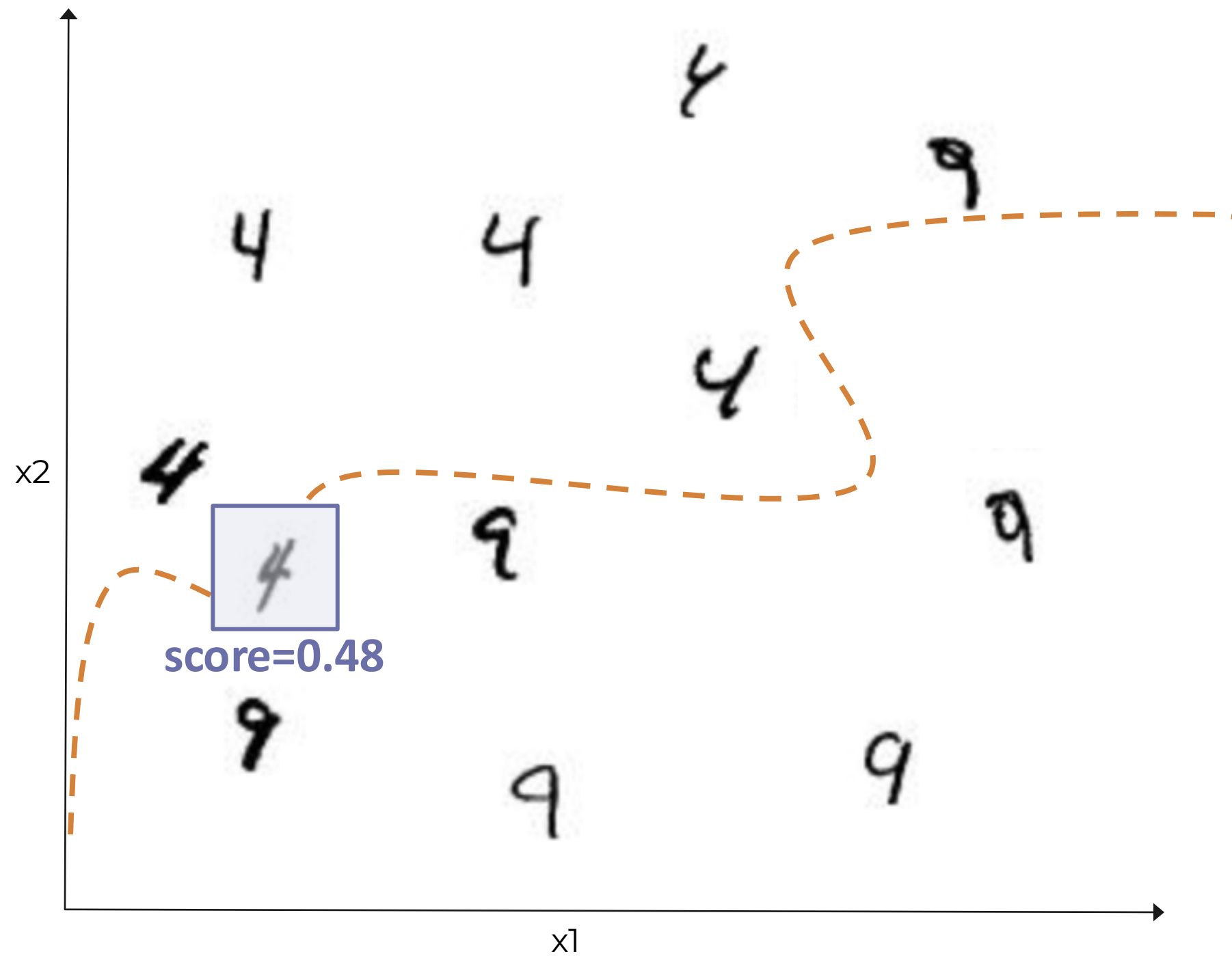
# Visual ambiguity: image sits between two or more classes



CNN

accuracy = 98%

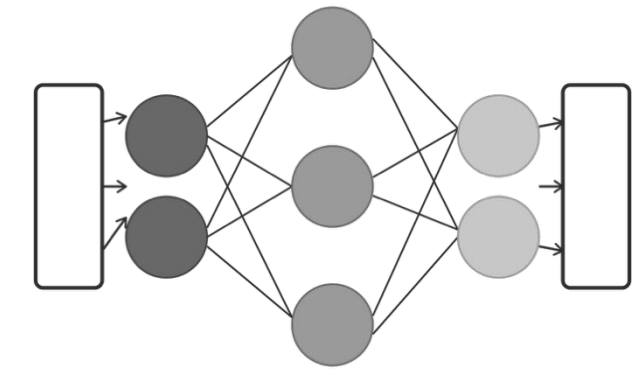
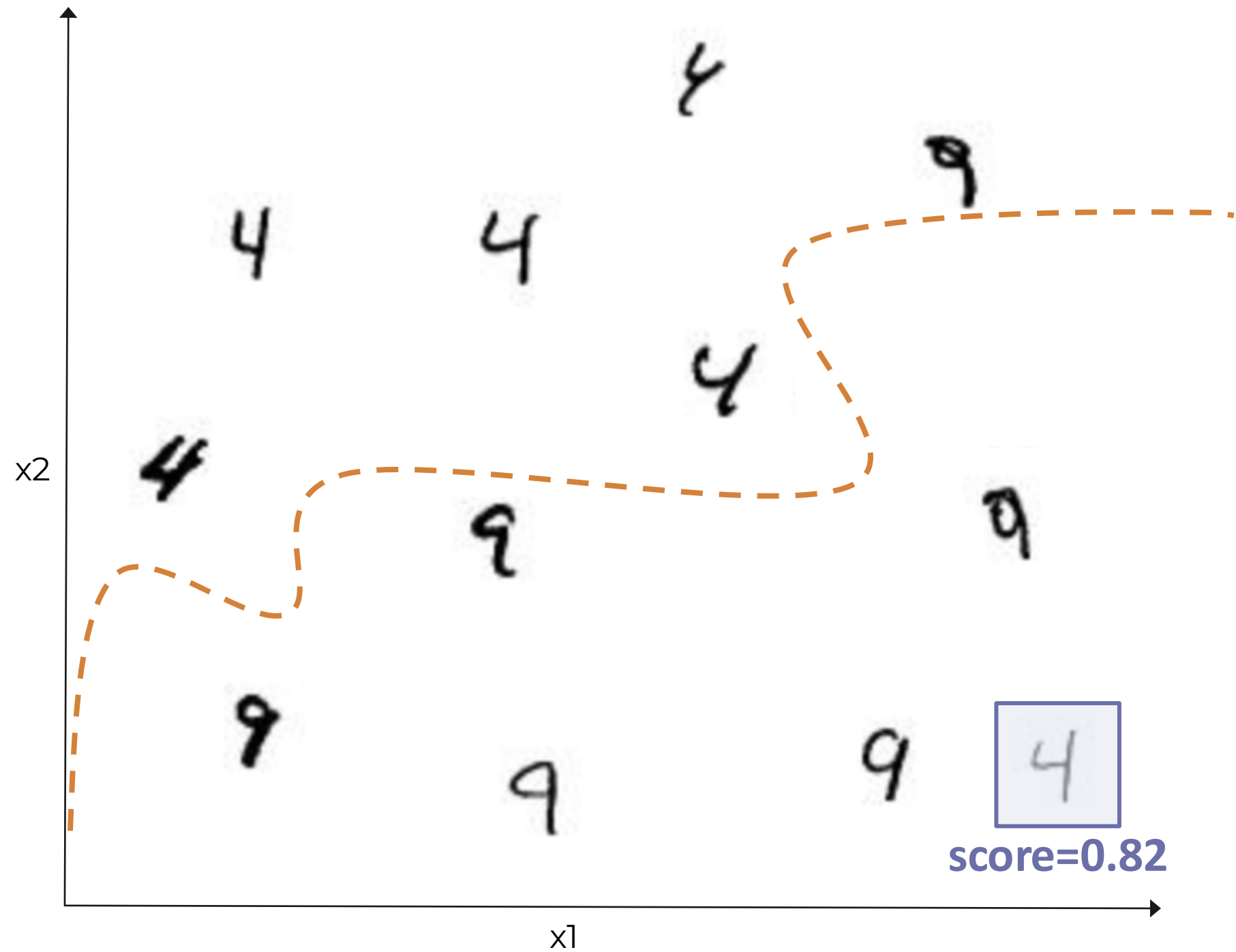
# Underrepresented subgroup: rare style that leads to not confident predictions



CNN

accuracy = 98%

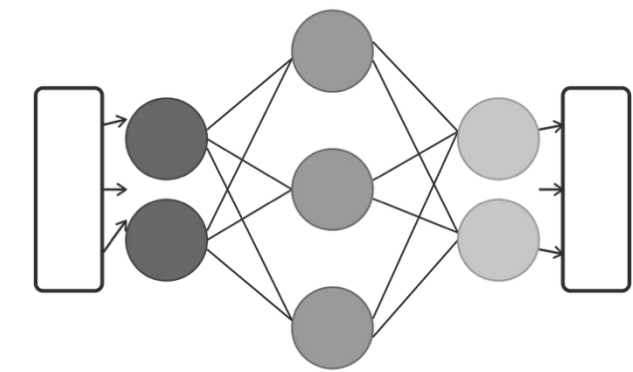
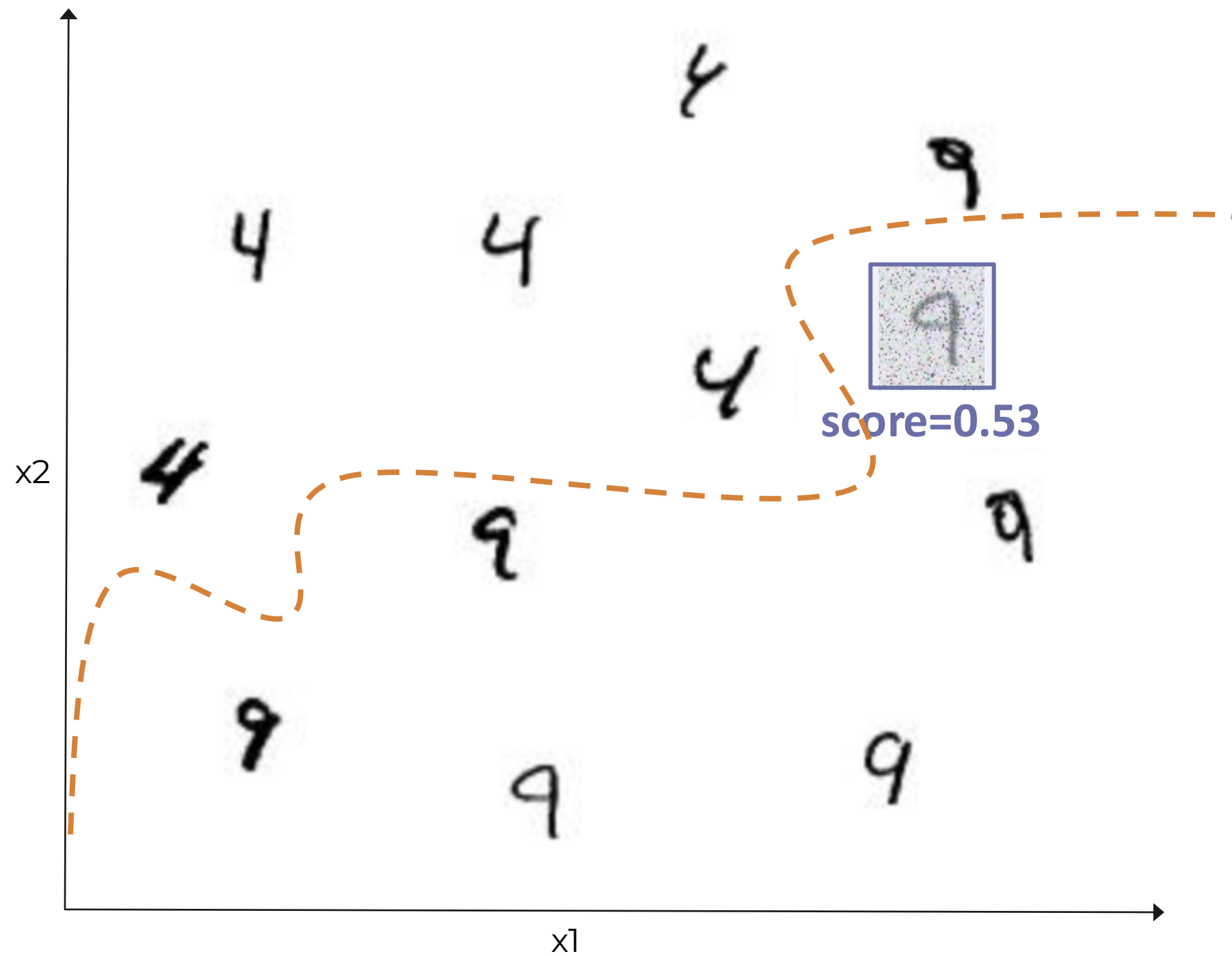
# Adversarial perturbation: imperceptible noise that flips the prediction



CNN

accuracy = 98%

# Corruption & noise: blur or artifacts that degrade input quality



CNN

accuracy = 98%

high performance  
on test set

≠

trustworthy  
behavior

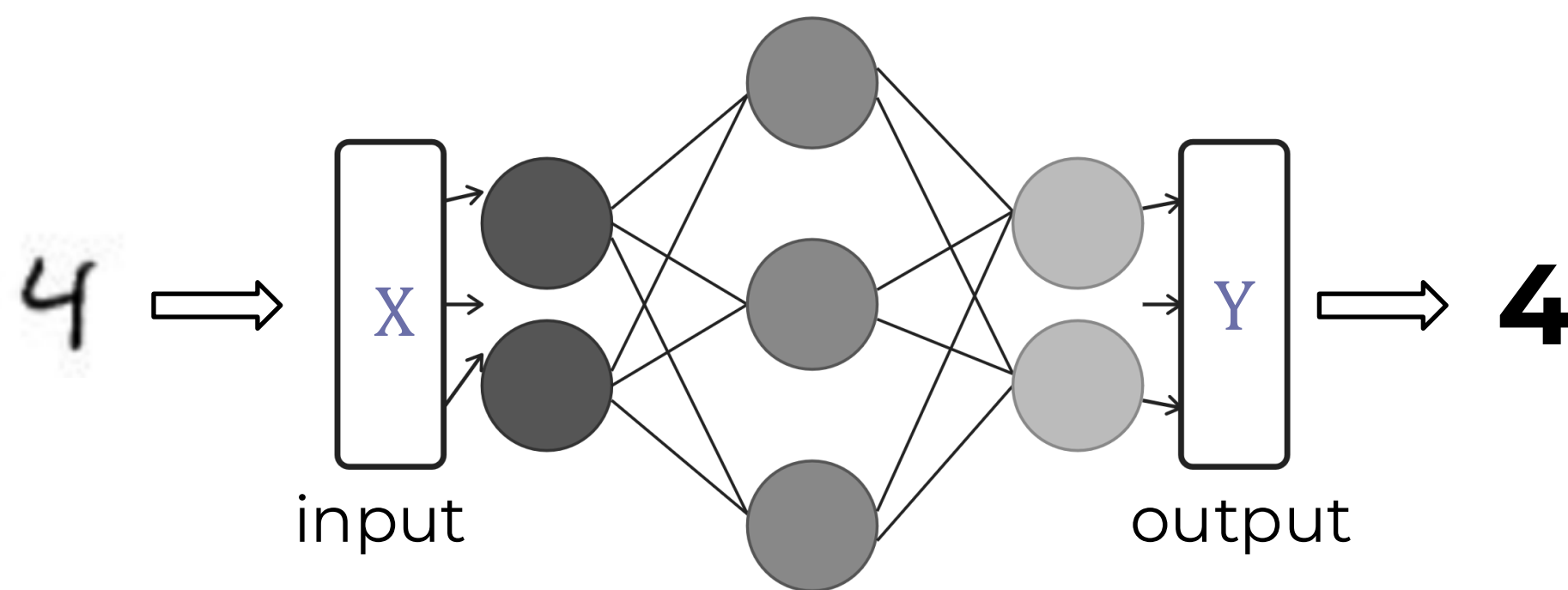
... because limitations may go undetected!

# Why do models fail?

---

How we formalize it...

A model (f) is a mapping between input and output



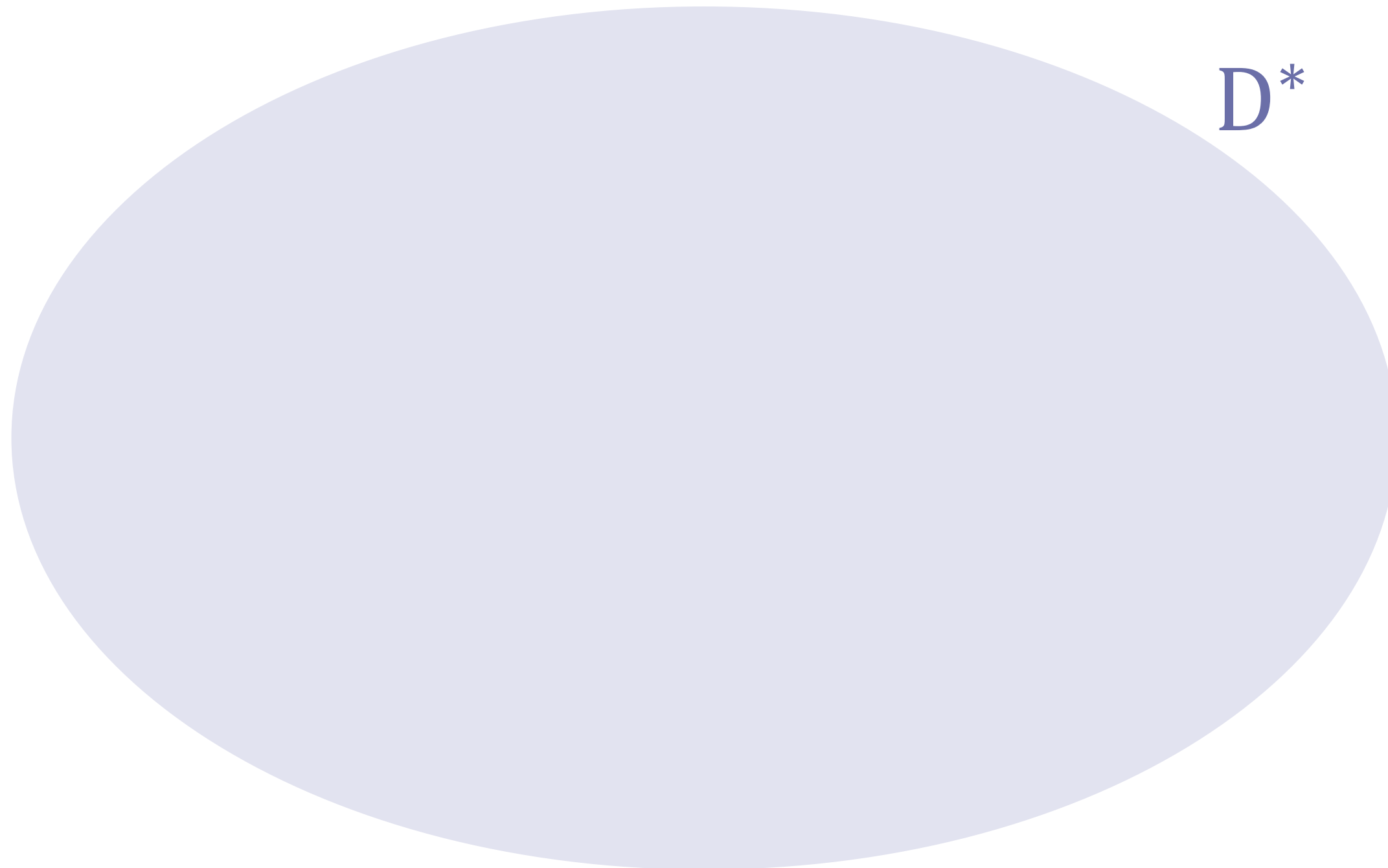
$$f: X \rightarrow Y$$

$$(x, y) \sim D_{\text{train}}$$

---

## In a perfect world...

An ideal model  $f^*$  would map perfectly the inputs and outputs from the perfect distribution  $D^*$



---

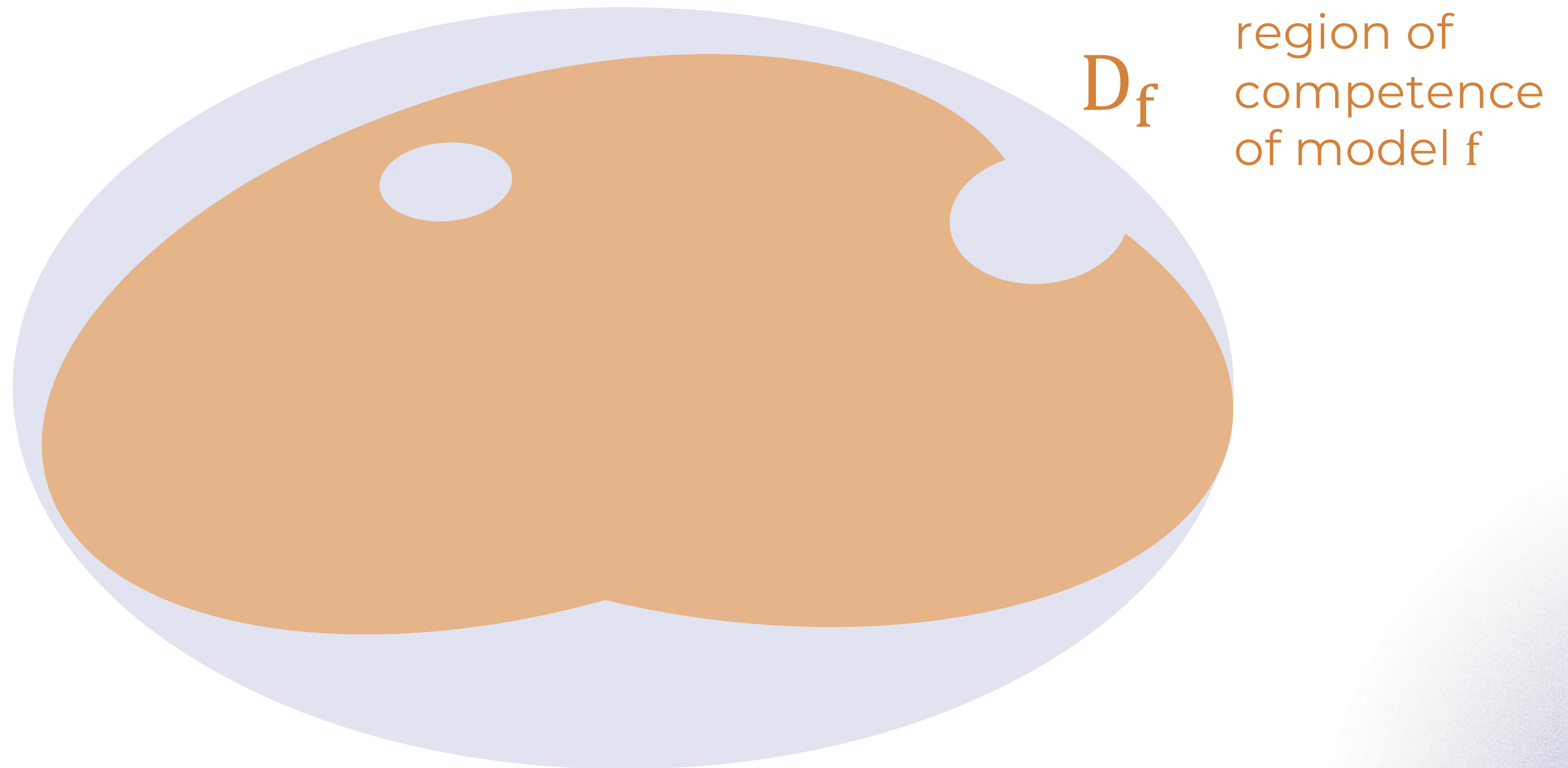
# Real world: training/test distributions

Finite samples drawn from the perfect distribution  $D^*$

$$D_{\text{train}}, D_{\text{test}} \sim D^*$$

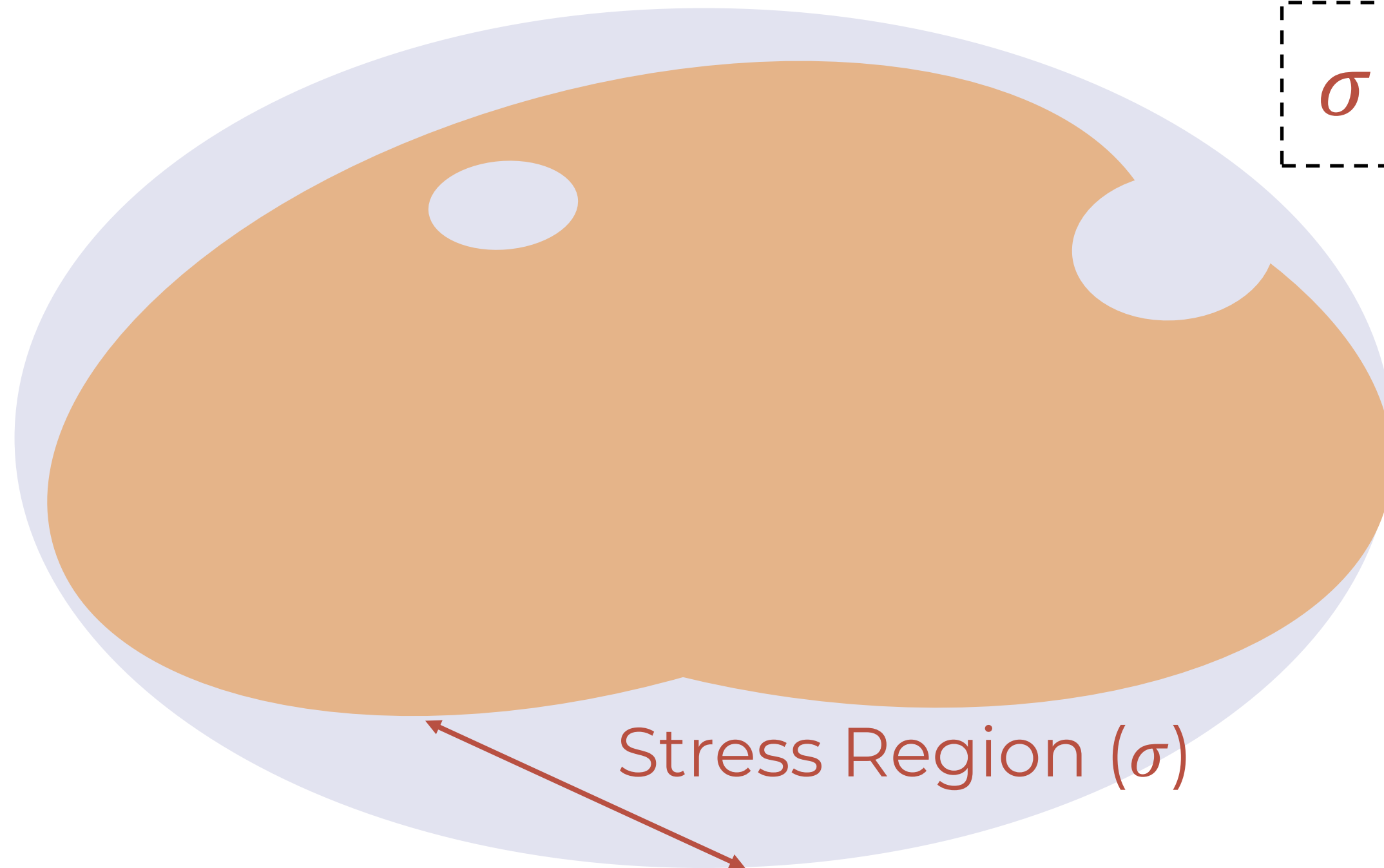
---

What the algorithm learned **correctly** from training data



→

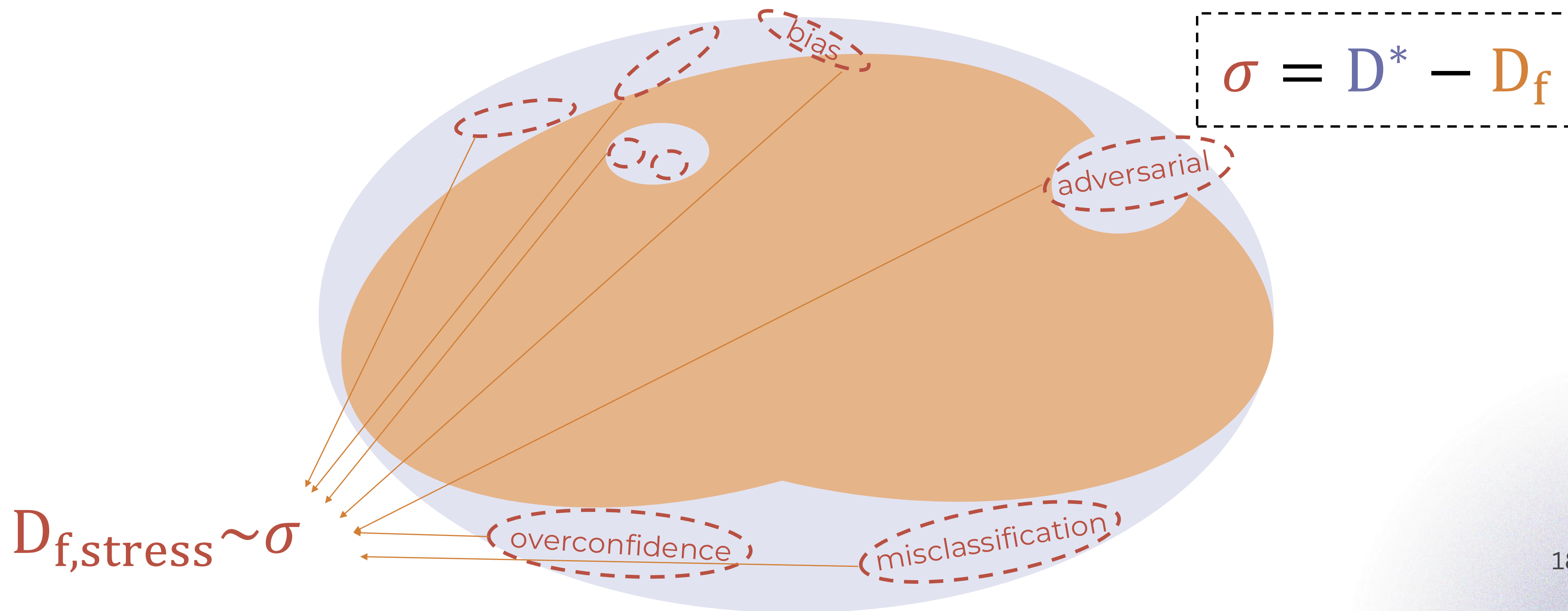
Stress region: where limitations emerge



$$\sigma = D^* - D_f$$

# Stress region: where limitations emerge

Set of distributions under which the model does not behave as expected



How can we expose model limitations  
in the stress region?

# Stress Testing: a framework

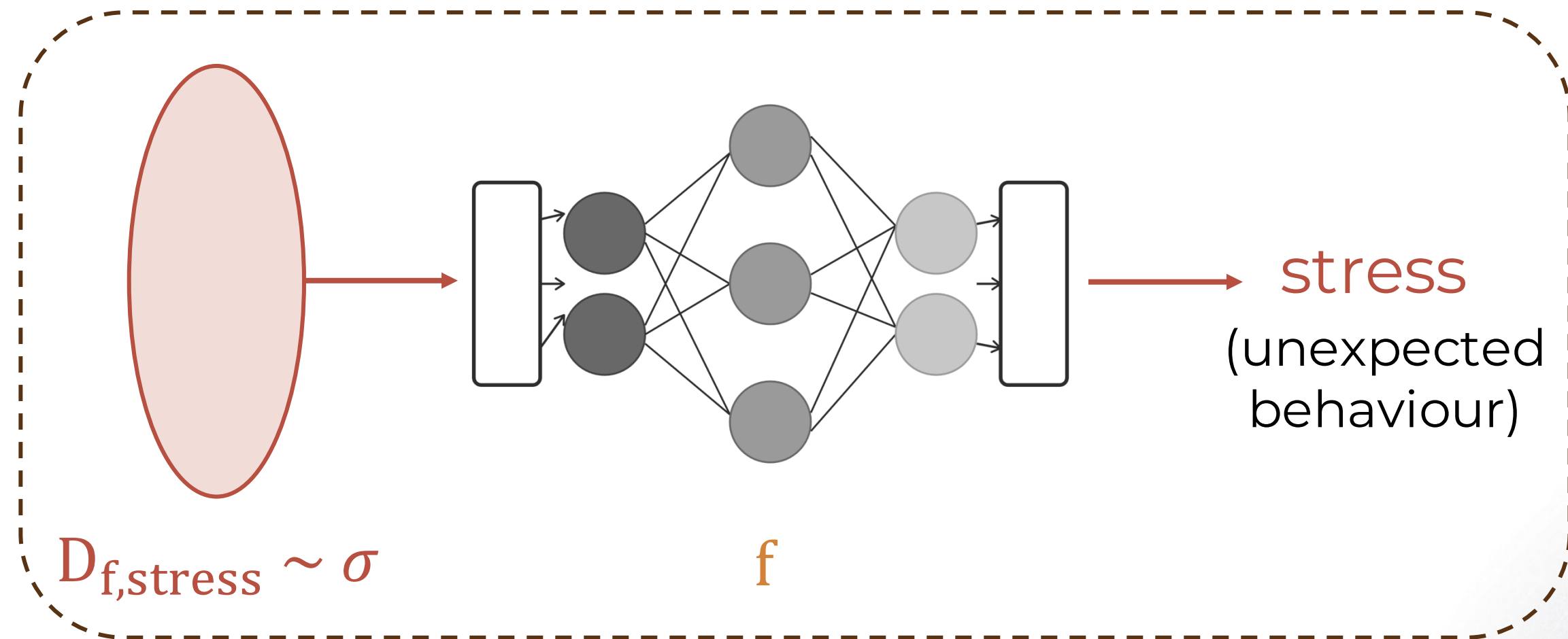


---

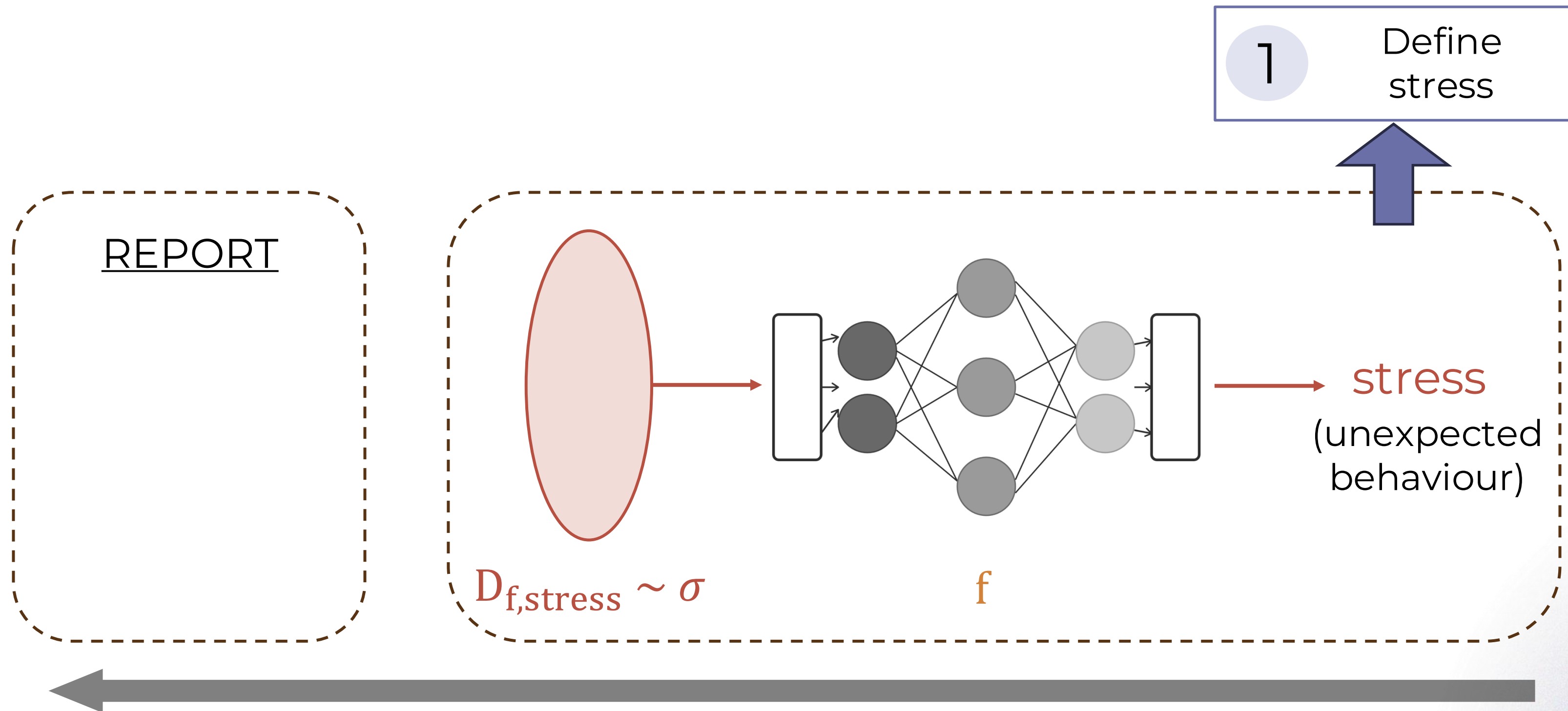
# Stress Testing: map the limits of the model competence

Data-driven way to report model limitations

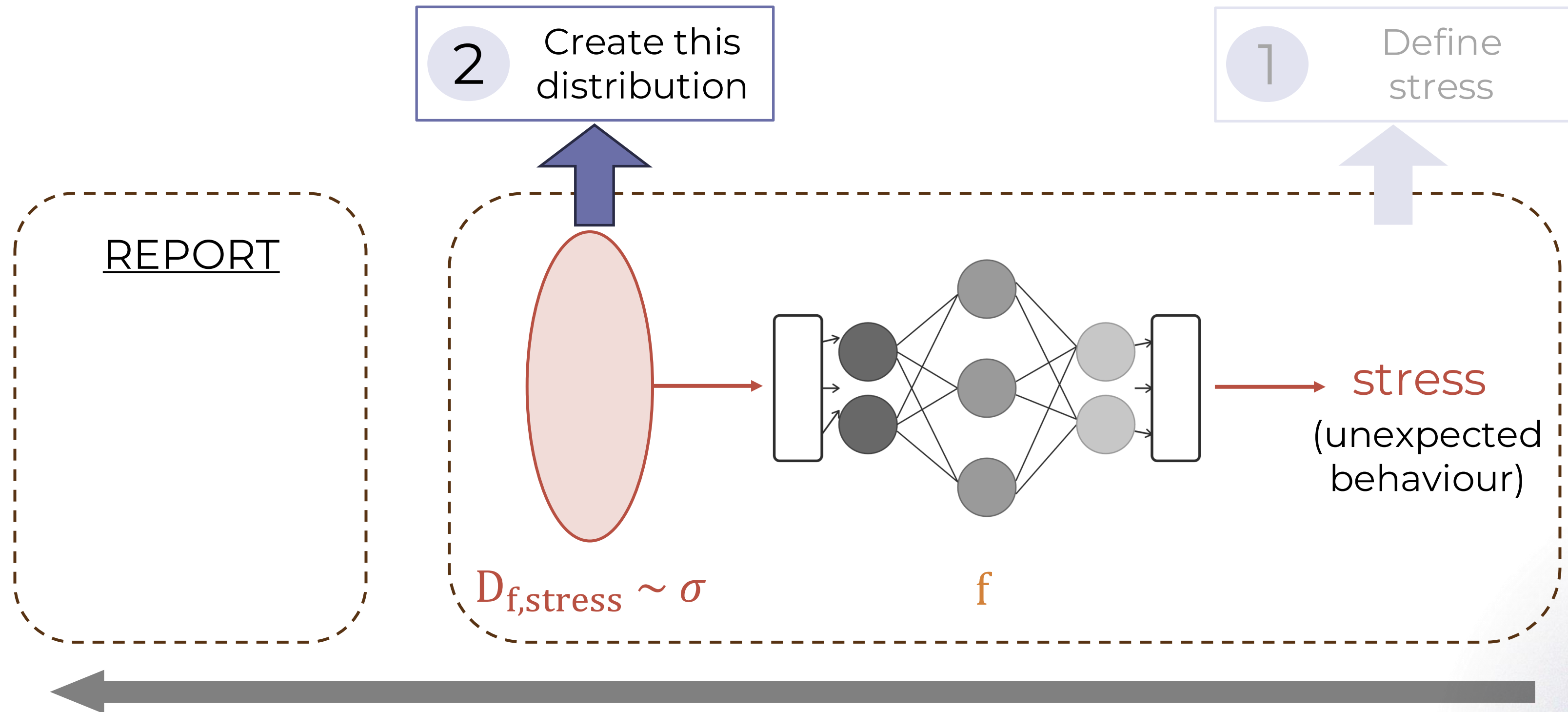
REPORT



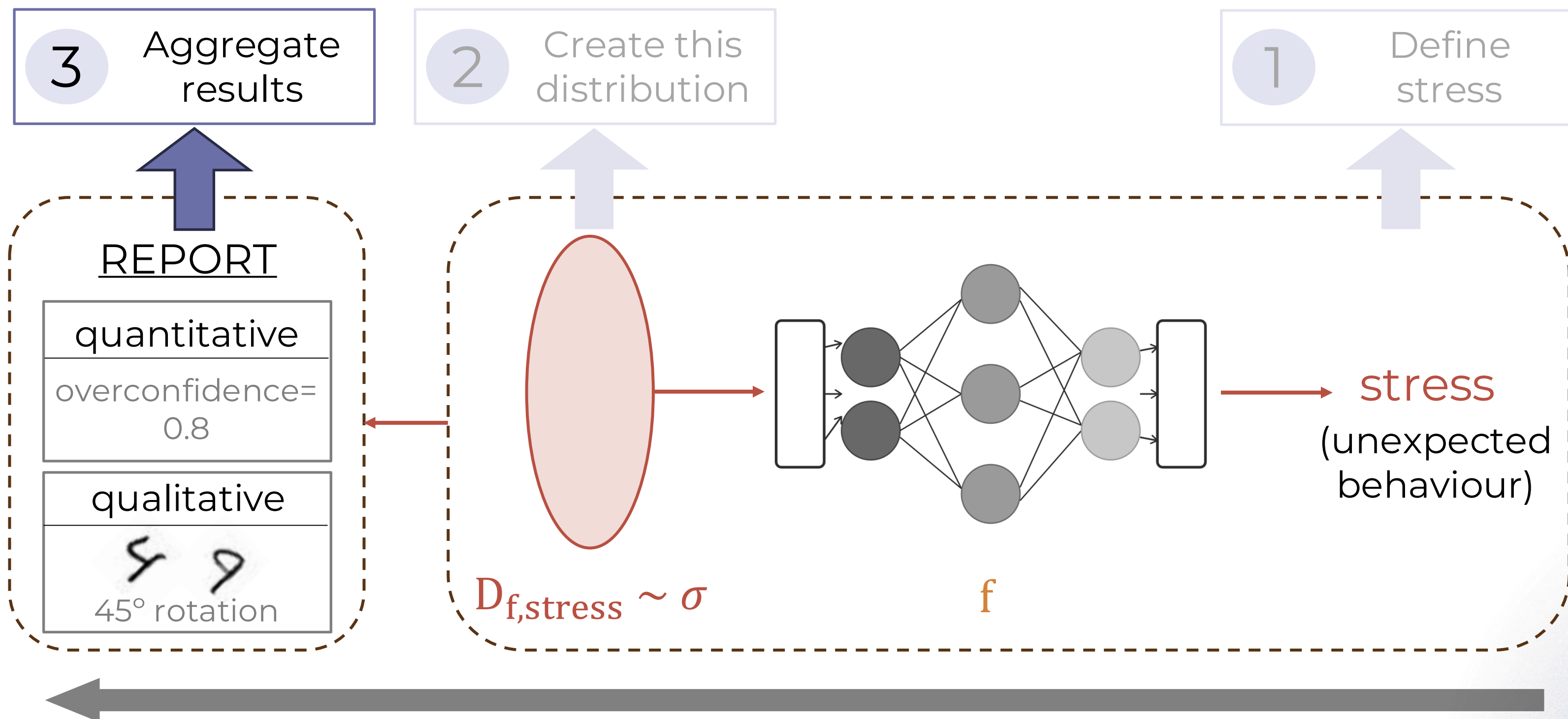
# Stress Testing



# Stress Testing



# Stress Testing: understand model decisions



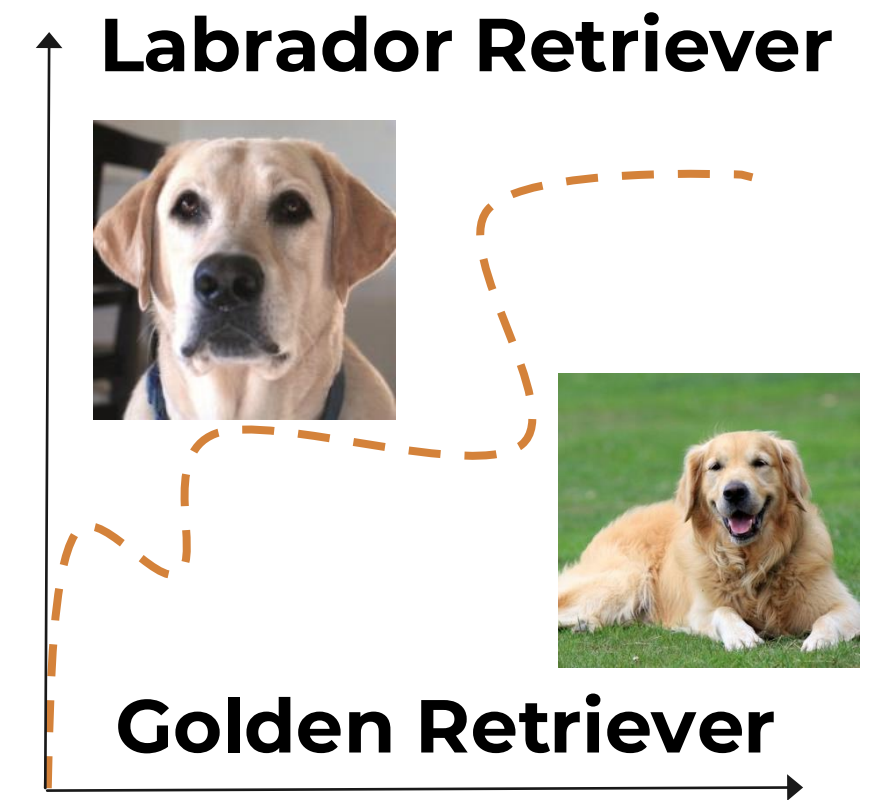
Application:

Stress Testing Decision

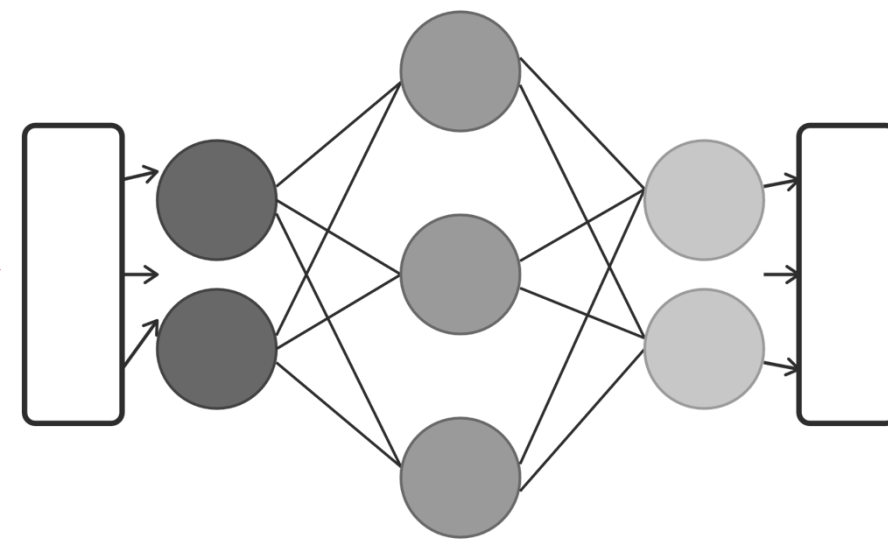
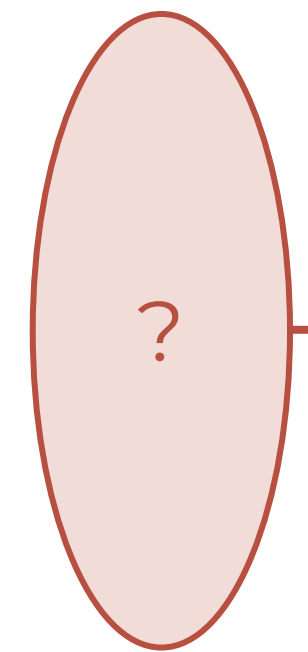
Boundaries of Image Classifiers

---

# What happens at decision boundaries?



REPORT



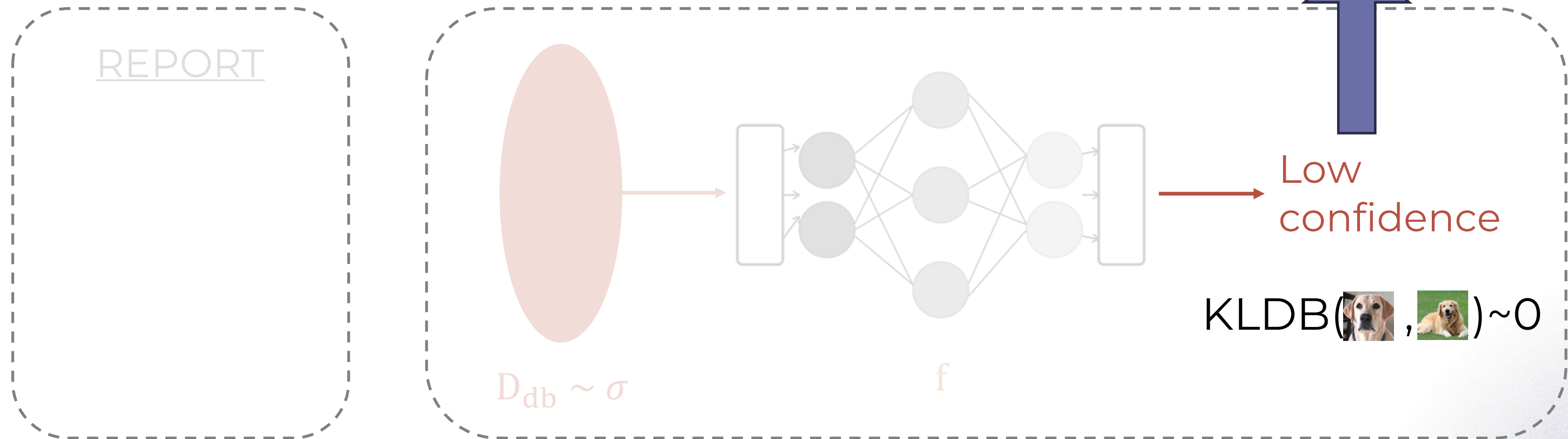
$f \rightarrow \text{ViT-B/16}$

Low  
confidence

$D_{\text{train}}$ : ImageNet     $D^*$ : visual world

# Stress Testing Decision Boundaries

KLDB: measures if an instance is in a decision boundary of N classes



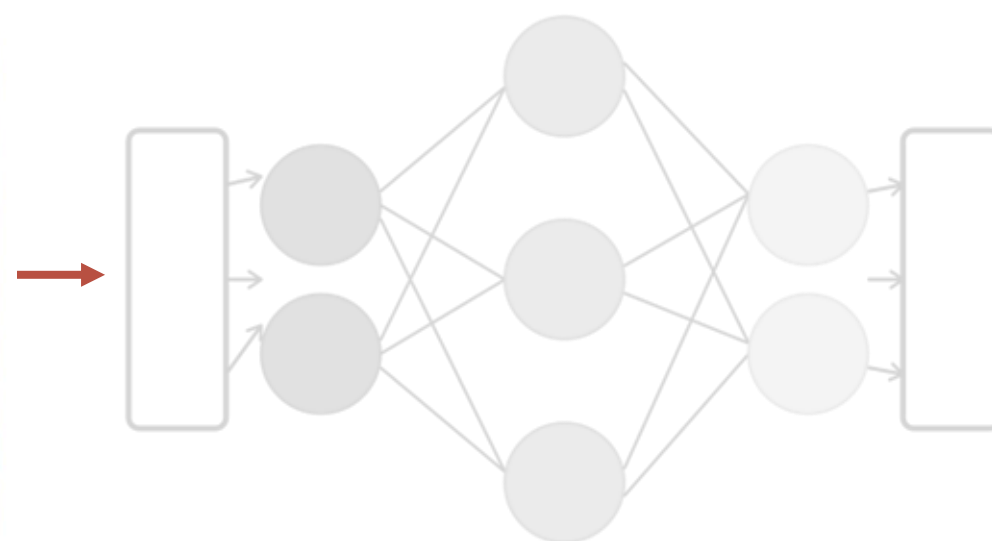
# Stress Testing Decision Boundaries

Synthetic Data Generation  
Stable Diffusion + classifier  
guidance

REPORT



$$D_{db} \sim \sigma$$

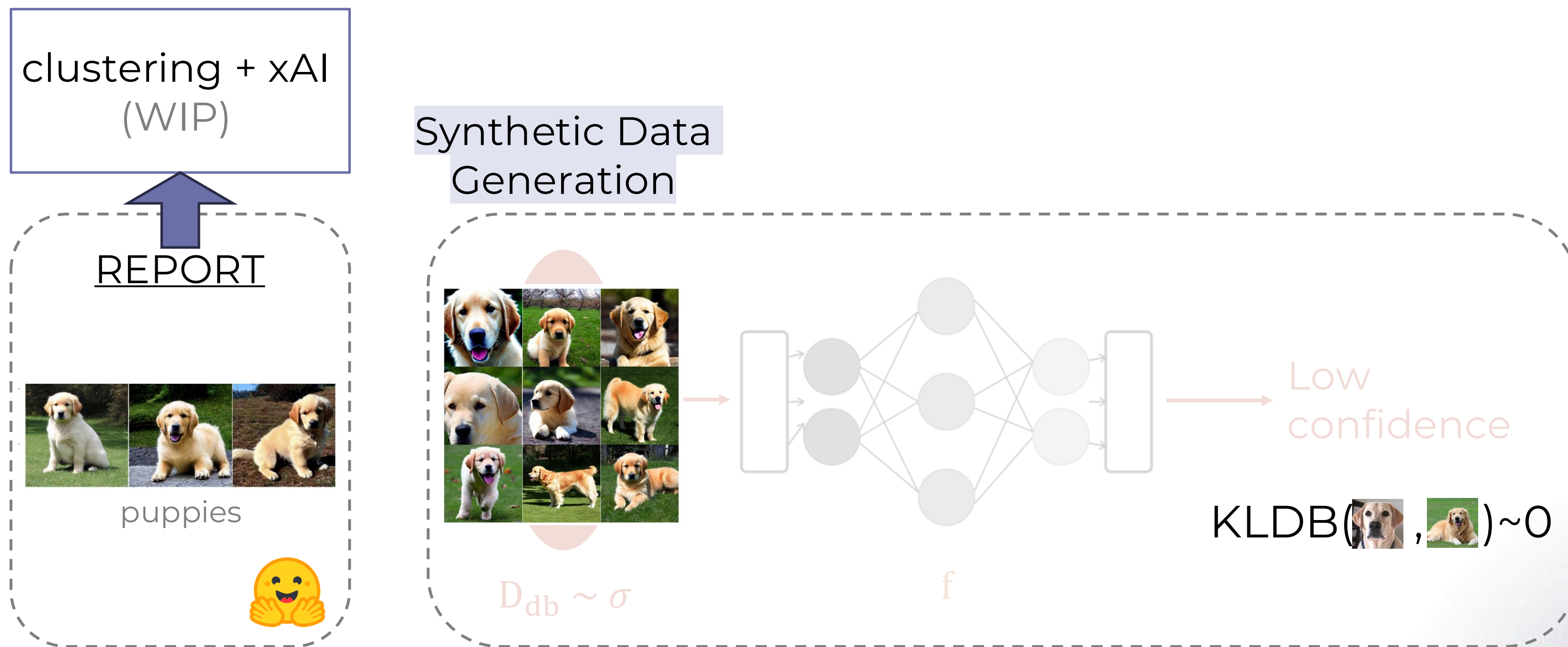


f

Low  
confidence

$$\text{KLDB}(\text{img}_1, \text{img}_2) \sim 0$$

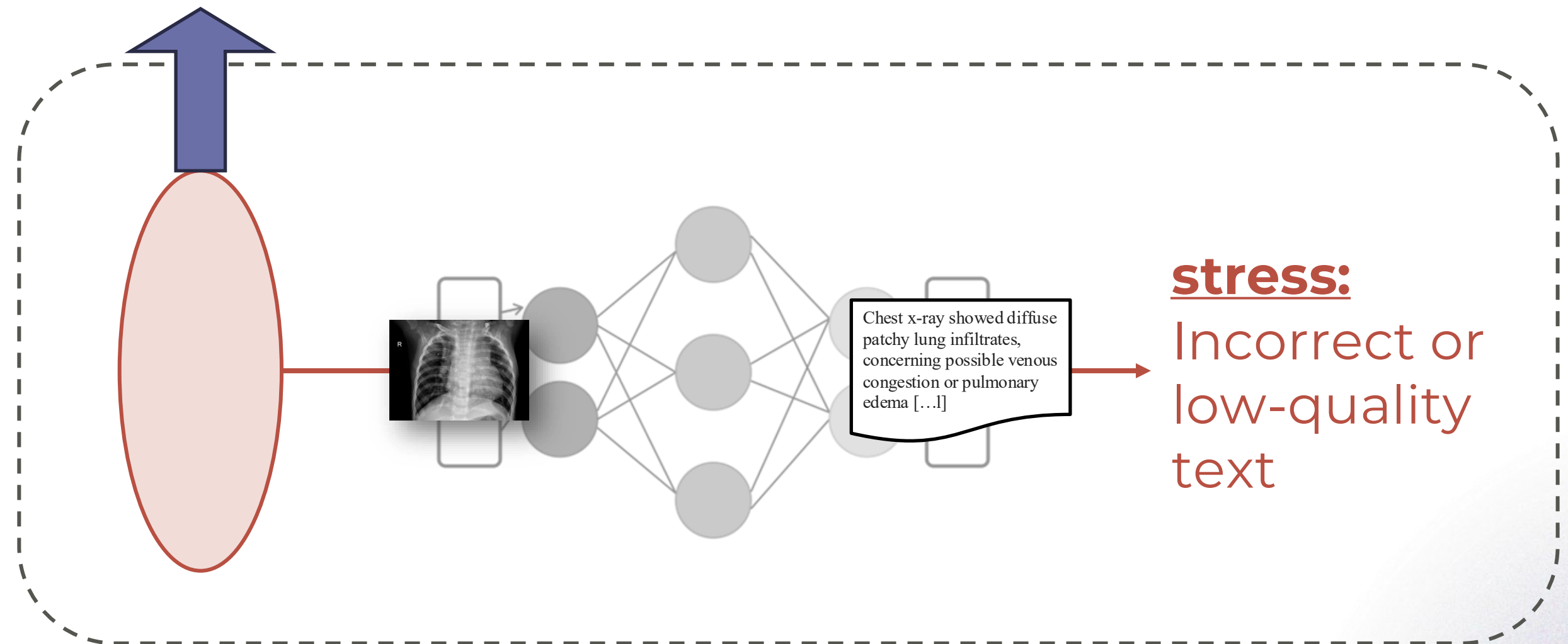
# Stress Testing: understand model decisions



# Other applications: medical report generation

synthetic data  
generation + image  
perturbations

REPORT  
entropy=0.4  
“patients with  
pacemakers cause a  
drop in the  
identification of lung  
masses of 30%”



# Other applications: time series forecasting

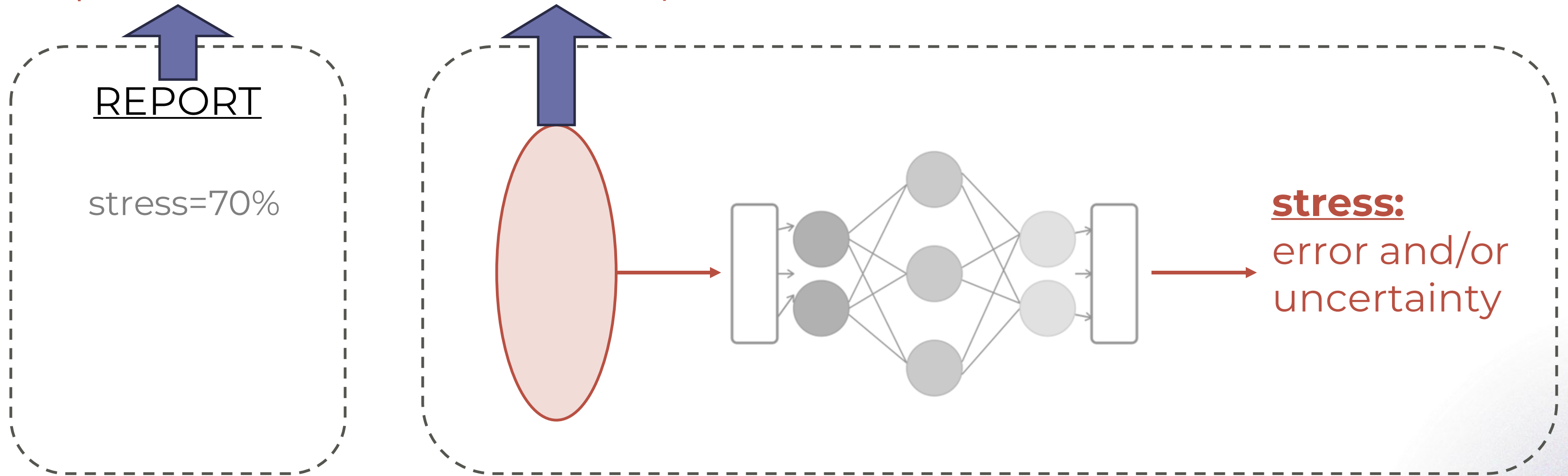
metaclassifier  
for stress  
prediction

REPORT

stress=70%

data augmentation  
(meta features,  
time series)

**stress:**  
error and/or  
uncertainty



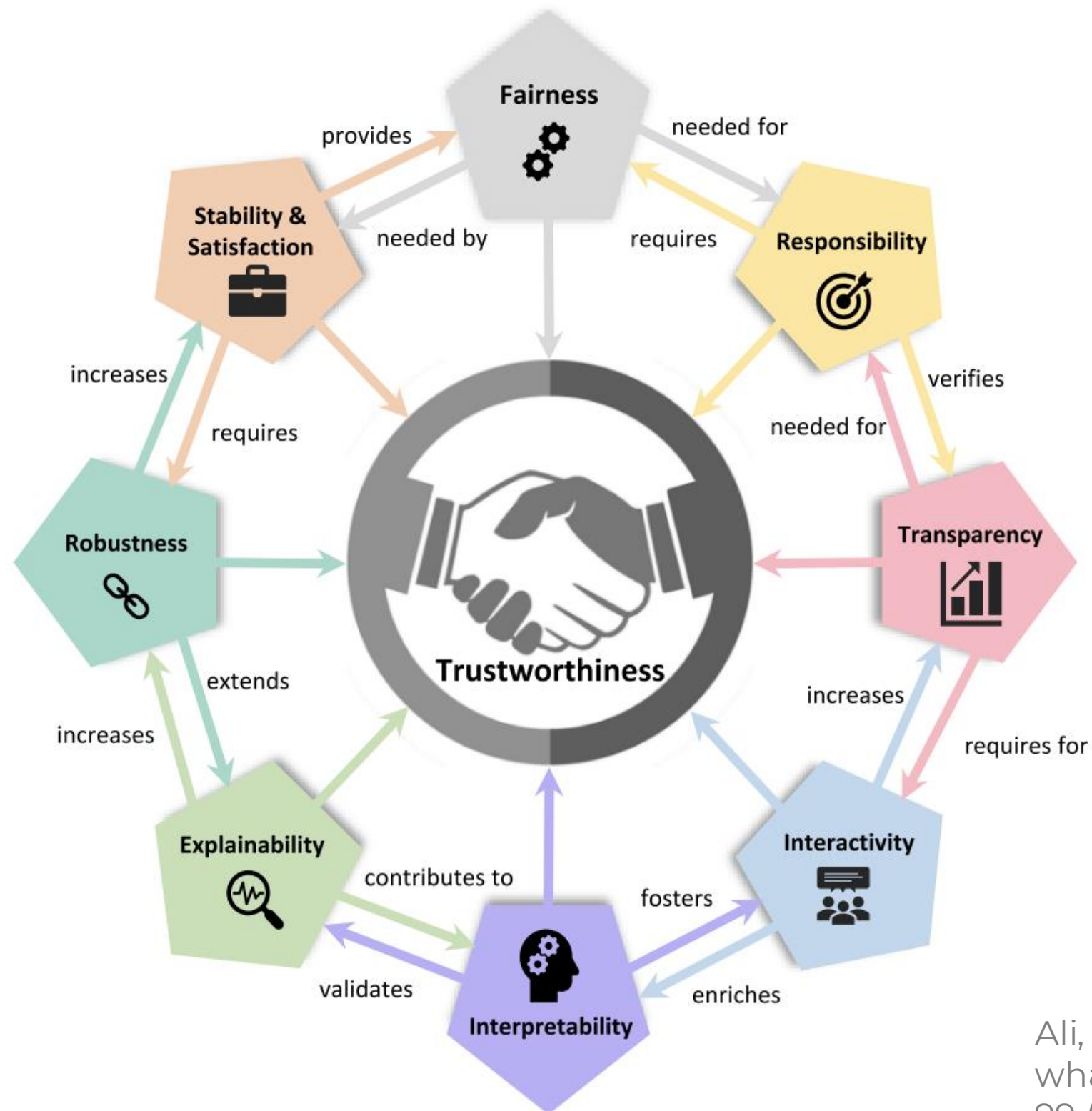
# The bigger picture: positioning stress testing



---

# Stress Testing & Trustworthiness

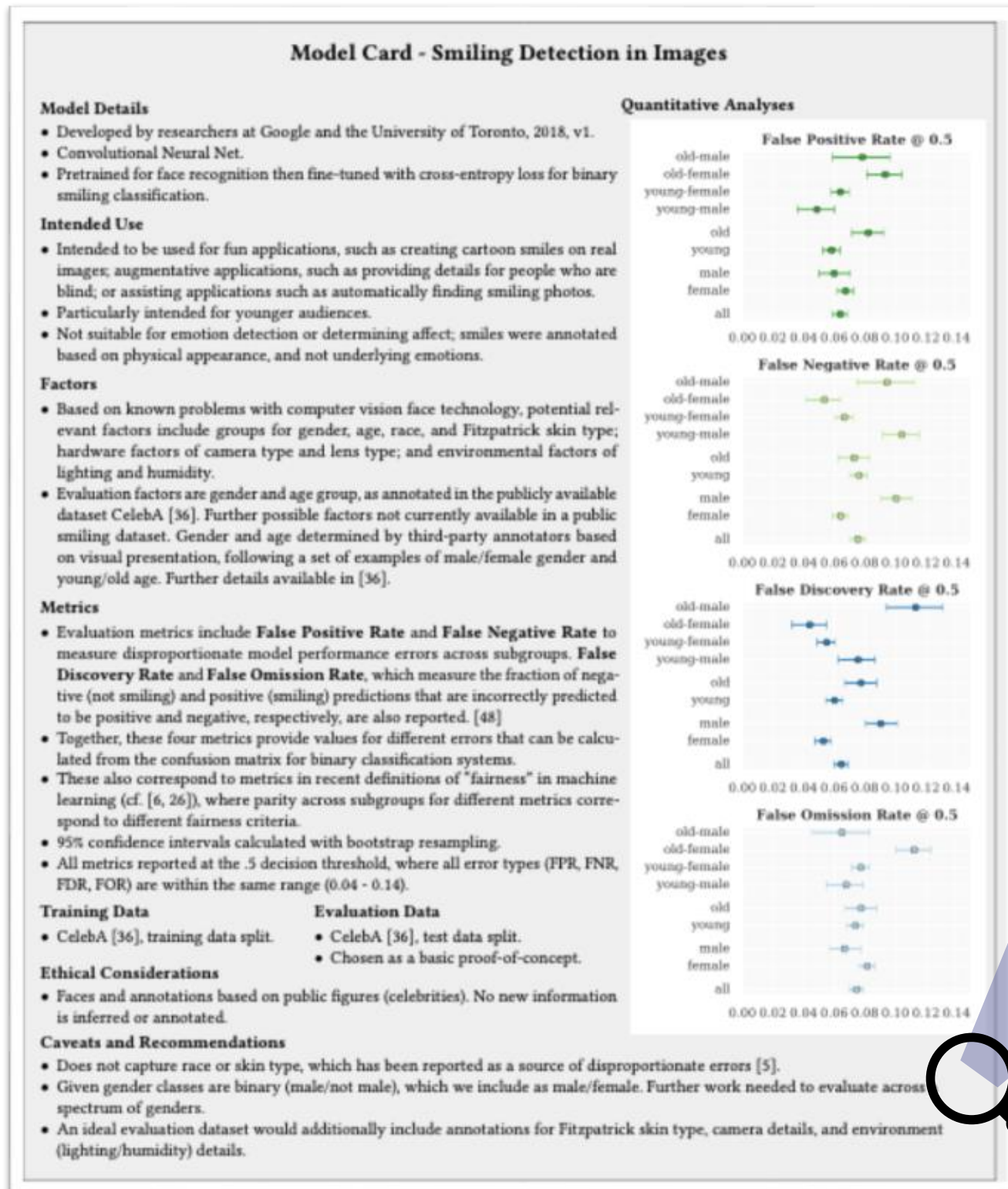
stress testing complements standard evaluation



Ali, Sajid, et al. "Explainable Artificial Intelligence (XAI): What we know and what is left to attain Trustworthy Artificial Intelligence." *Information fusion* 99 (2023)

# Stress Testing & Responsible AI

stress testing generates the evidence to fill model cards



**Caveats and Recommendations**

- Does not capture race or skin type, which has been reported as a source of disproportionate errors [5].
- Given gender classes are binary (male/not male), which we include as male/female. Further work needed to evaluate across a spectrum of genders.
- An ideal evaluation dataset would additionally include annotations for Fitzpatrick skin type, camera details, and environment (lighting/humidity) details.

Mitchell, Margaret, et al. "Model cards for model reporting." *Proceedings of the conference on fairness, accountability, and transparency*. 2019

---

# Stress Testing: framework limitations

- Stress definition shapes the framework
  - poor specification may lead to wrong finding or missed failures
- Finding/creating the stress distribution is non-trivial
  - options such as synthetic data generation may introduce new biases
  - the probe may be imperfect
- General concept with costly implementation

---

## Key takeaways

- High performance on a test set does not guarantee trustworthy behaviour
  - limitations may go undetected...
- Stress testing probes the stress region to map the limits of a model's competence
  - supporting responsible AI
- General framework: applicable to any model, any domain, any definition of stress



Inês Gomes  
(image)



Flávia Carvalhido  
(text-image)



Ricardo Inácio  
(time series)



Ons Zammel  
(text/text-image)



Moisés Santos  
(time series)

André Restivo

Carlos Soares

Cátia Teixeira

Henrique L. Cardoso

Jan van Rijn

Luís F. Teixeira

Marília Barandas

Thomas Bäck

Gomes, I., Teixeira, L. F., Van Rijn, J. N., Soares, C., Restivo, A., Cunha, L., & Santos, M. (2024). Finding patterns in ambiguity: interpretable stress testing in the decision boundary. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops* (pp. 8316-8321).

Brito, A., Santos, M., Folgado, D., & Soares, C. (2025). tsMIST: Model Sensitivity Analysis with Time Series Morphing. In *International Conference on Discovery Science* (pp. 270-285).

Carvalhido, F., Cardoso, H. L., Cerqueira, V., & Soares, C. (2025). XAI in Medical Image Report Generation: Unlocking Stress Testing as a Responsible AI Practice for Multimodal Models. *Proceedings of the Second Multimodal, Affective and Interactive eXplainable AI Workshop (MAI-XAI 2025) co-located with the 28th European Conference on Artificial*.

Inácio, R., Cerqueira, V., Barandas, M., & Soares, C. (2025). Mast: interpretable stress testing via meta-learning for forecasting model robustness evaluation. *Machine Learning*, 114(11), 251.

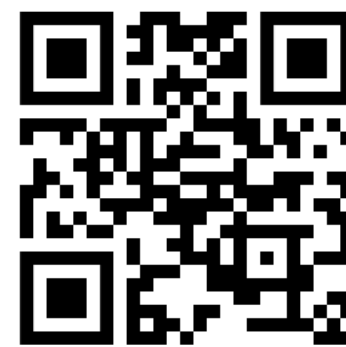
Teixeira, C., Gomes, I., Soares, C., & van Rijn, J. N. (2025). Hubris Benchmarking with AmbiGANs: Assessing Model Overconfidence with Synthetic Ambiguous Data. In *International Conference on Discovery Science* (pp. 476-491).

Gomes, I., Santos, M., Restivo, A., Soares, C., Teixeira, L. F., Van Rijn, J. N., Bäck, T. (2026). Stress Testing the Decision Boundaries of Image Classifiers via Latent Diffusion. Submitted to *Machine Learning*.

# Thank you!

## Do you have any questions?

ines.gomes@fe.up.pt



contacts



feedback

Logos of partner institutions and funding bodies:

- Universiteit Leiden The Netherlands
- liacs Leiden Institute of Advanced Computer Science
- FEUP FACULDADE DE ENGENHARIA UNIVERSIDADE DO PORTO
- LIACC
- Center for Responsible AI
- PRR Plano de Recuperação e Resiliência
- REPÚBLICA PORTUGUESA
- Financiado pela União Europeia NextGenerationEU
- UK Research and Innovation
- AISYN4MED
- Project funded by:
  - Schweizerische Eidgenossenschaft Confédération suisse Confederazione Svizzera Confederaziun svizra Swiss Confederation
  - Federal Department of Economic Affairs, Education and Research EAER State Secretariat for Education, Research and Innovation SERI

Funded by the European Union under the **Horizon Europe Framework Programme** Grant Agreement N°: 101095387. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or European Commission. Neither the European Union nor the European Health and Digital Executive Agency can be held responsible for them.