# Continually Learning Deep Machines that Understand What They Don't Know

Dr. Martin Mundt

TUDa & hessian.AI - Junior Research Group Leader
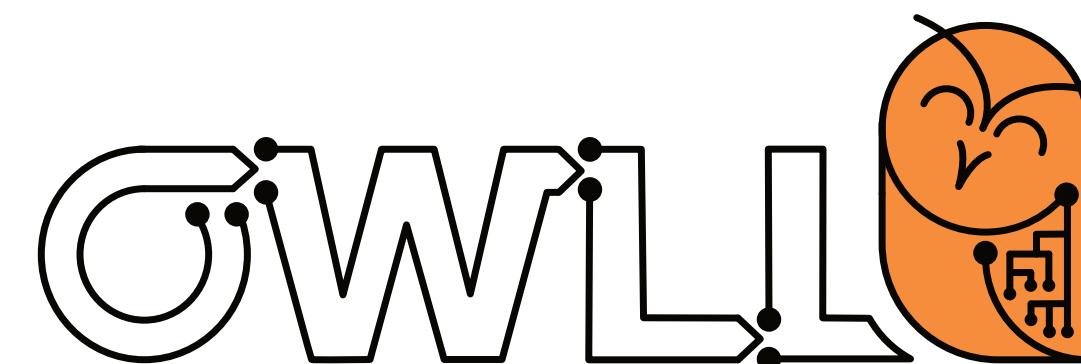
ContinualAI Board Member

http://owll-lab.com

Prof. Dr. Kristian Kersting

TECHNISCHE UNIVERSITÄT DARMSTADT

hessian.AI

OWLL

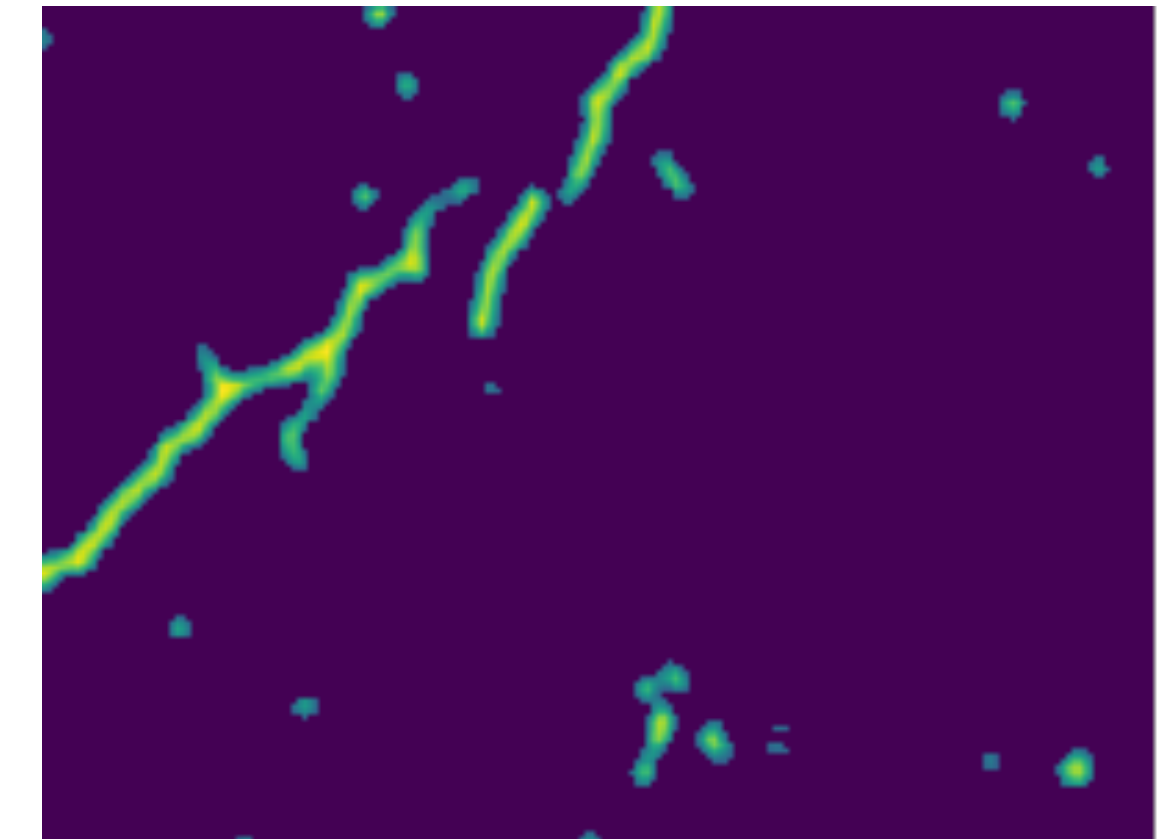ContinualAI

We could talk about AI applications...

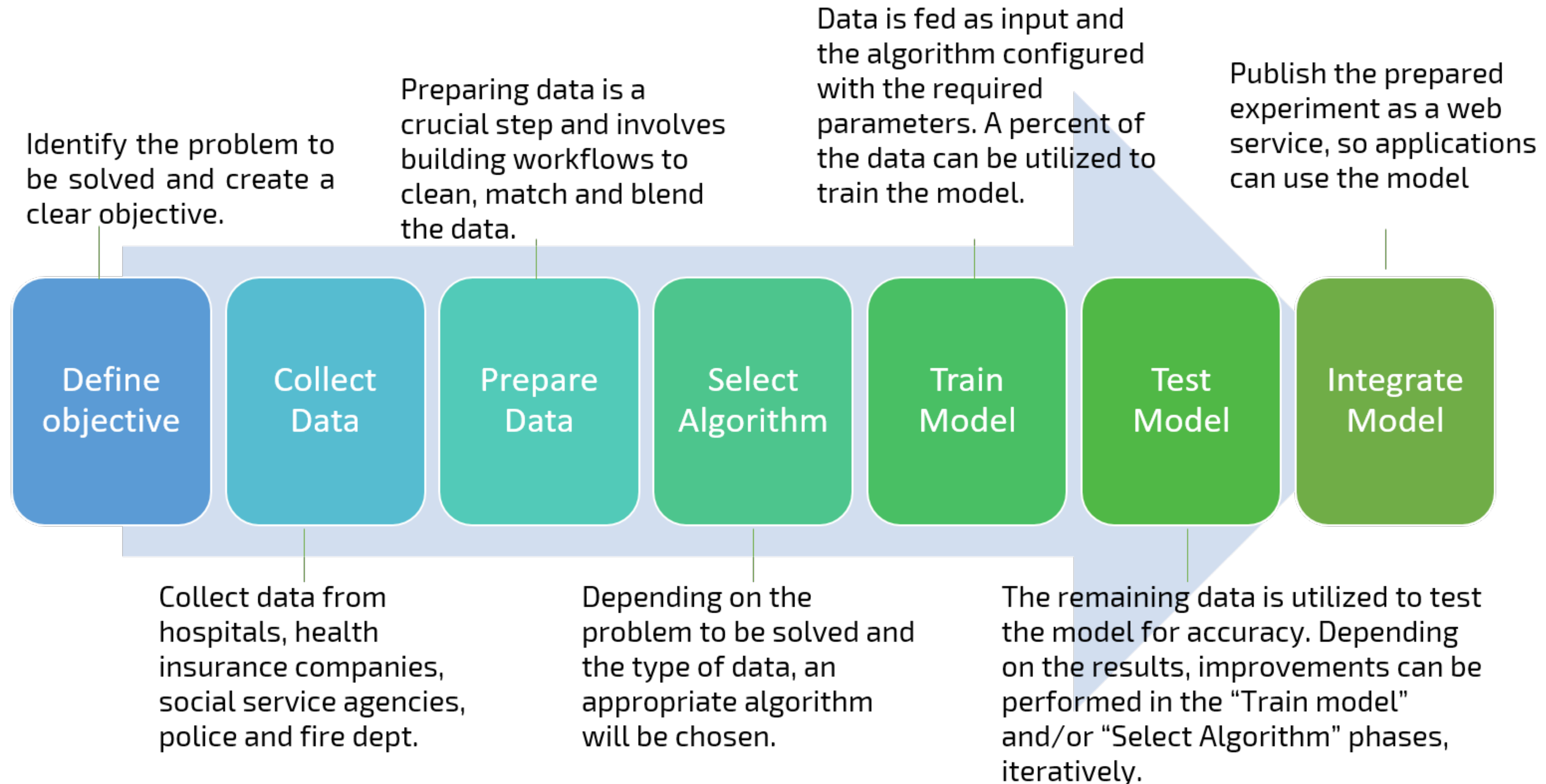*Fly drone*  *Scan bridge*  *Inspect surface*  *Measure defects*

Many factors: low amounts of data, experts are rare, annotation is cumbersome, predictions need to be robust (safety critical), tons of variation when system is deployed
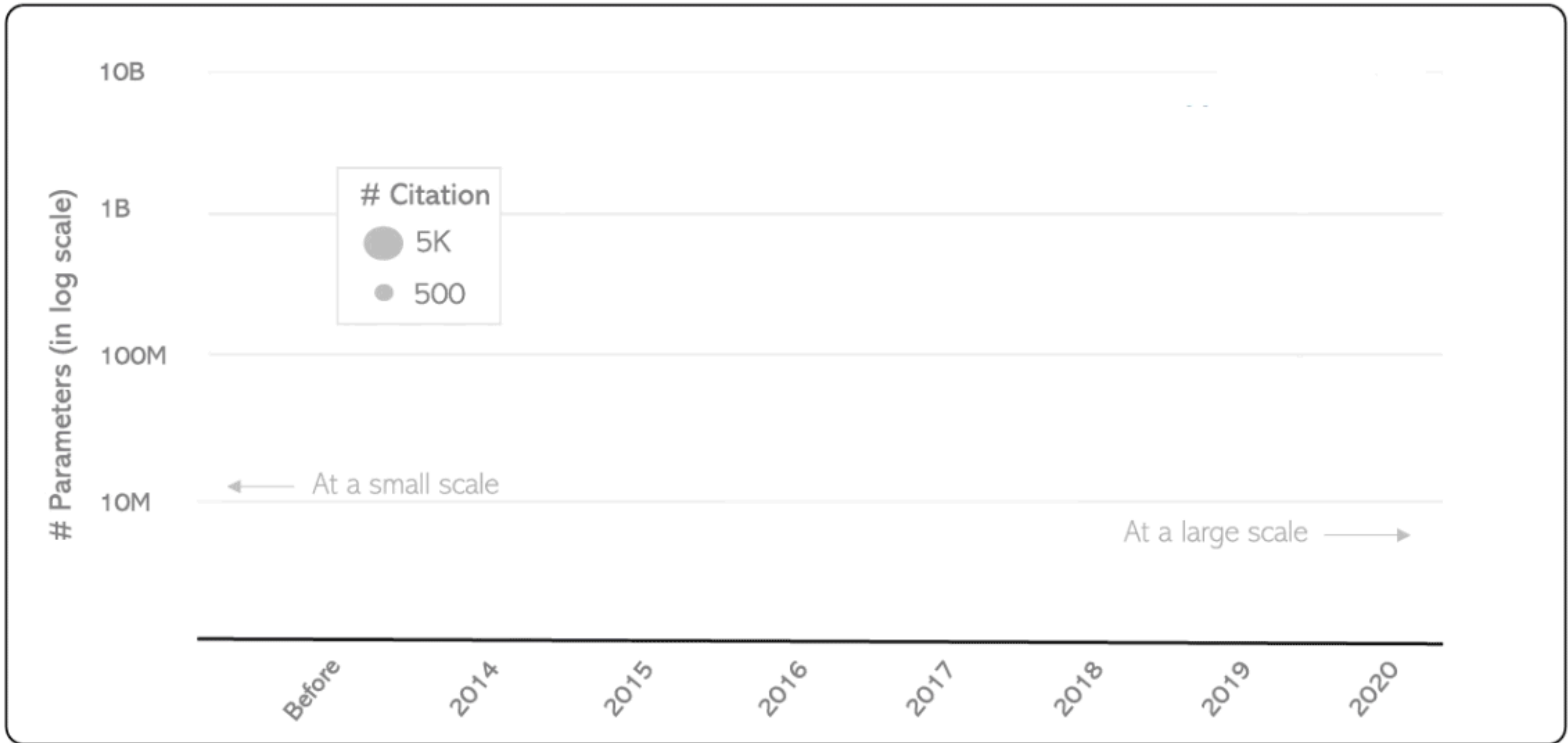
*Mundt, CVPR 2019*

# The standard machine learning workflow

Identify the problem to be solved and create a clear objective.

Preparing data is a crucial step and involves building workflows to clean, match and blend the data.

Data is fed as input and the algorithm configured with the required parameters. A percent of the data can be utilized to train the model.

Publish the prepared experiment as a web service, so applications can use the model

**Define objective** → **Collect Data** → **Prepare Data** → **Select Algorithm** → **Train Model** → **Test Model** → **Integrate Model**

Collect data from hospitals, health insurance companies, social service agencies, police and fire dept.

Depending on the problem to be solved and the type of data, an appropriate algorithm will be chosen.

The remaining data is utilized to test the model for accuracy. Depending on the results, improvements can be performed in the "Train model" and/or "Select Algorithm" phases, iteratively.

Figure from https://www.congrelate.com/get-workflow-machine-learning-images/

# Is a static machine learning workflow + <u>scale</u> all we need?



Research Director at Deepmind says all we need now is scaling

**Nando de Freitas** @Nando... · 4 t.
Someone's opinion article. My opinion: It's all about scale now! The Game is Over! It's about making these models bigger, safer, compute efficient, faster at sampling, smarter memory, more modalities, INNOVATIVE DATA, on/offline, ... 1/N

thenextweb.com
DeepMind's new Gato AI makes me fear humans will never achieve AGI

💬 10   🔁 22   ♡ 78   ⬄

Li & Gao, "A deep generative model trifecta: three advances that work towards harnessing large-scale power, Microsoft Research Blog, 2020:
https://www.microsoft.com/en-us/research/blog/a-deep-generative-model-trifecta-three-advances-that-work-towards-harnessing-large-scale-power/

We have "foundation" models now, but **humans learn & reason**



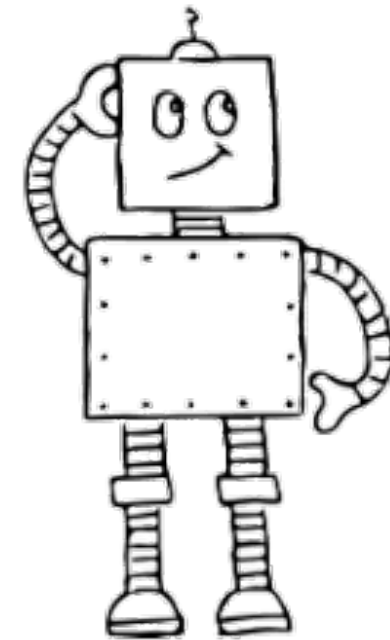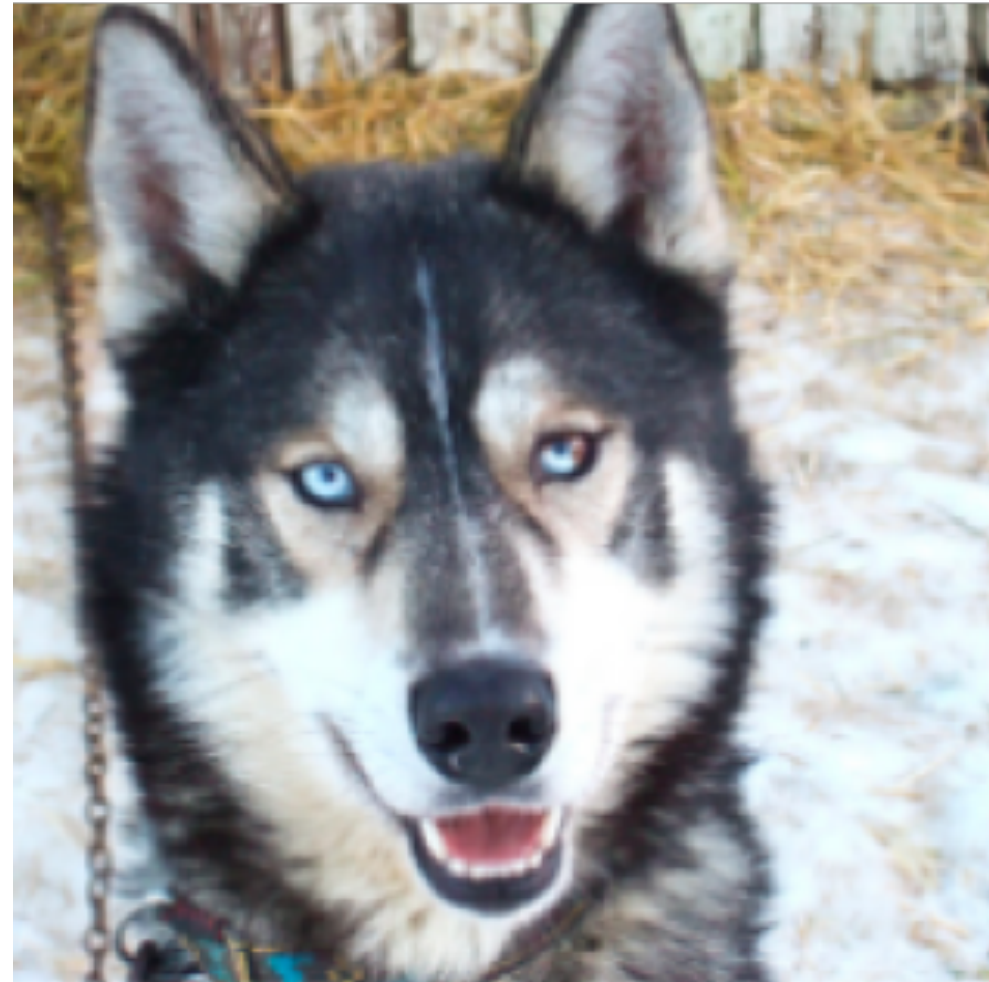Importantly, humans **revise** their knowledge & **continue** adapting

Why should we care: can we **trust** deep neural networks?
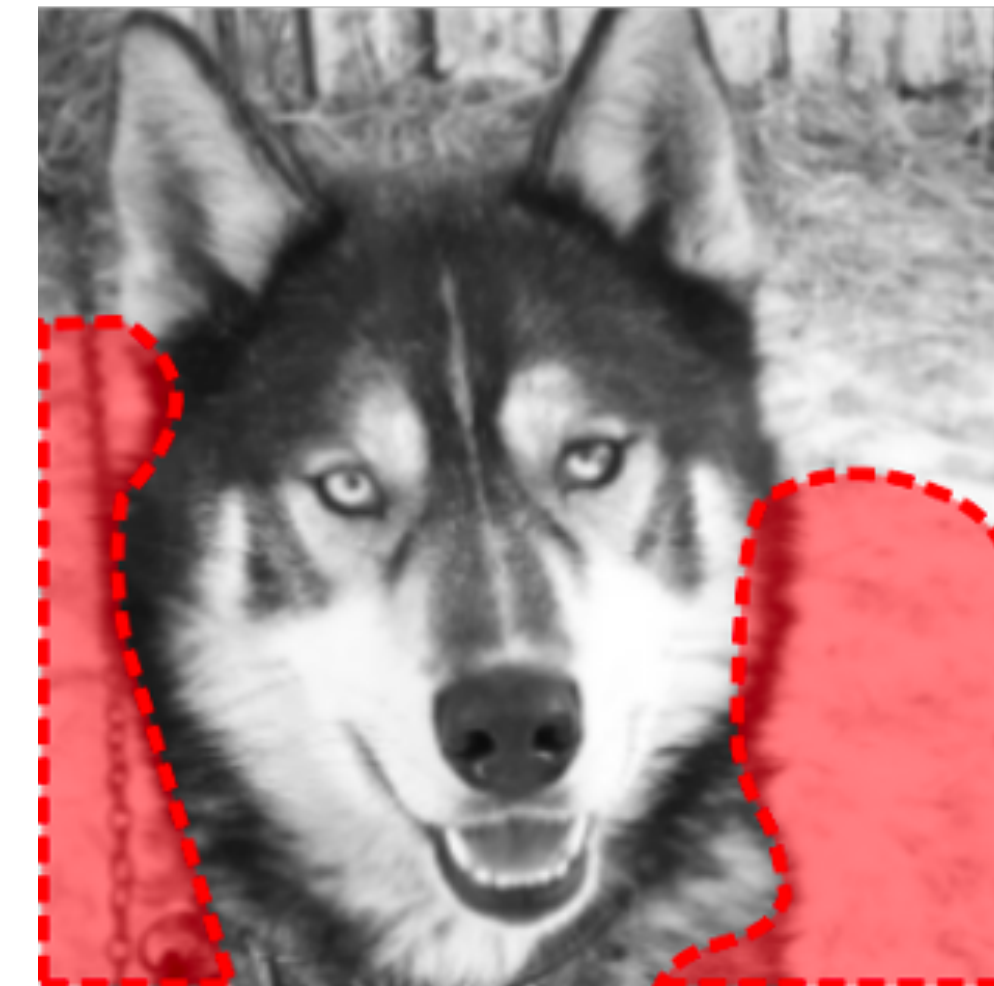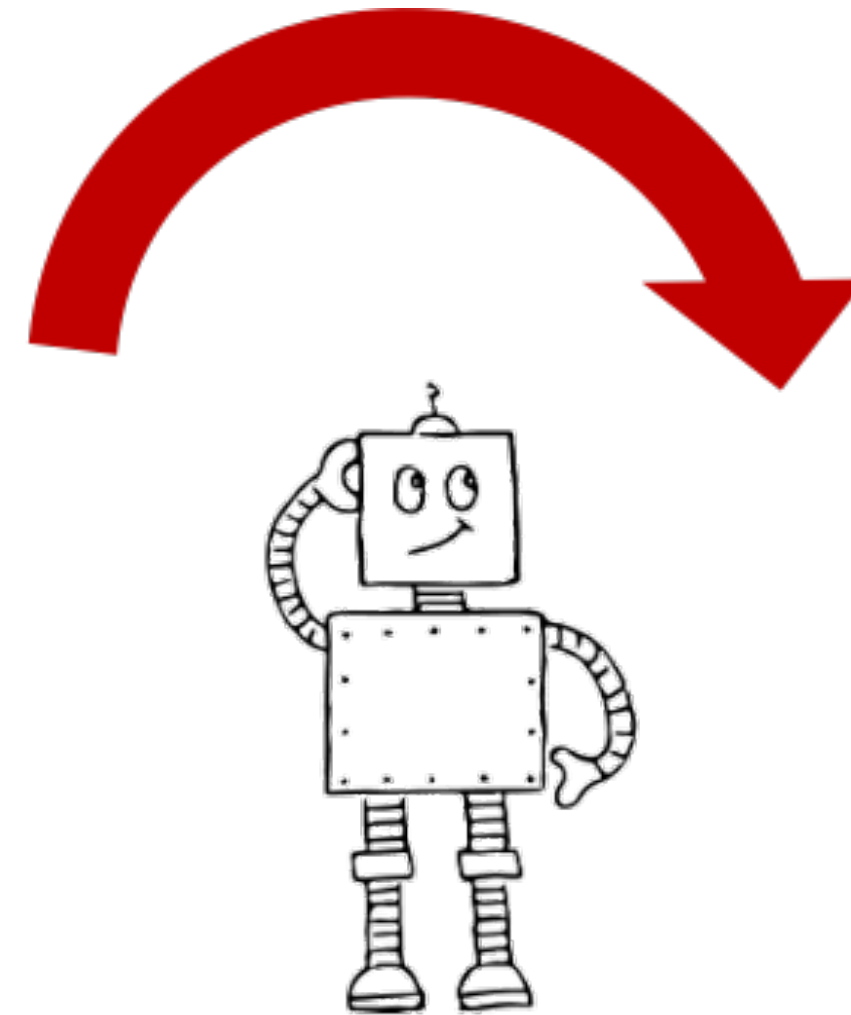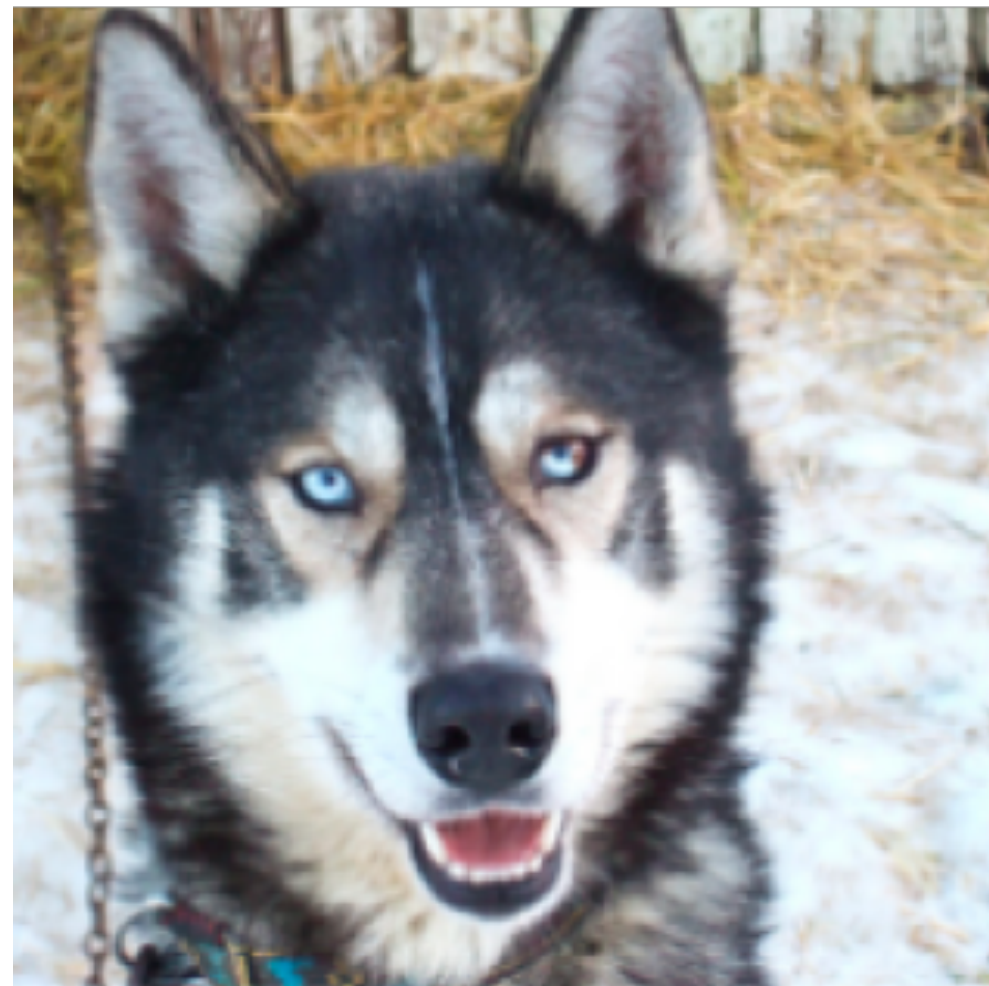
XAI suggests ways to **detect**, but not **fix** the issue!

Consider an example image classification task about distinguishing between **husky dogs** and **wolves**
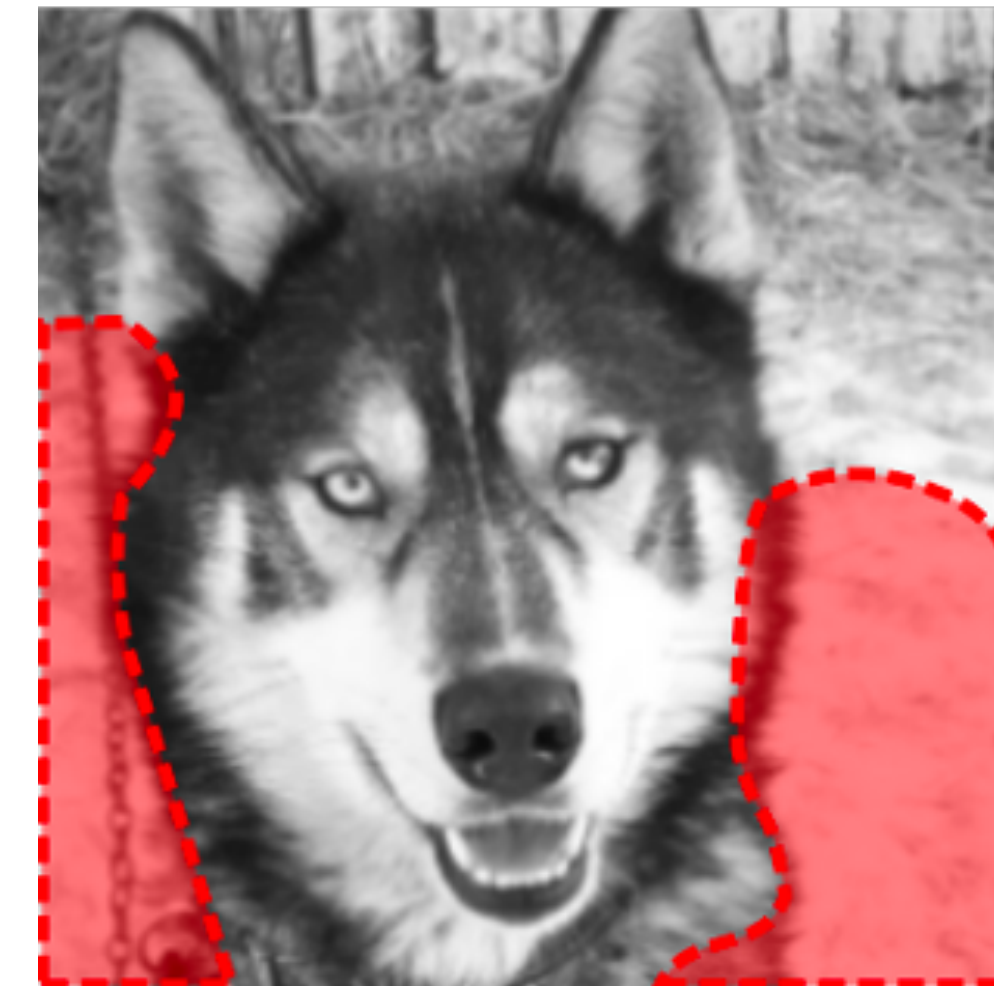
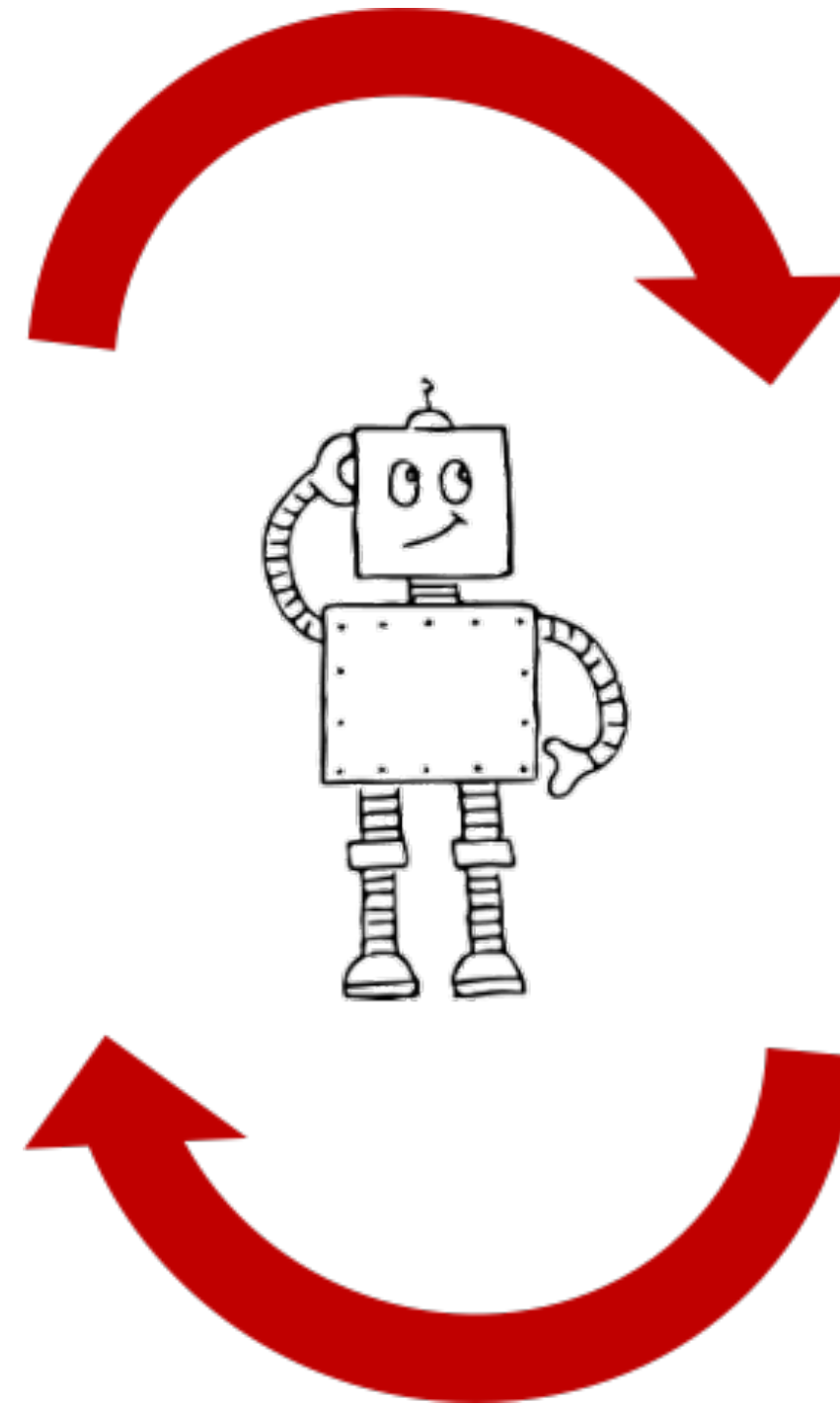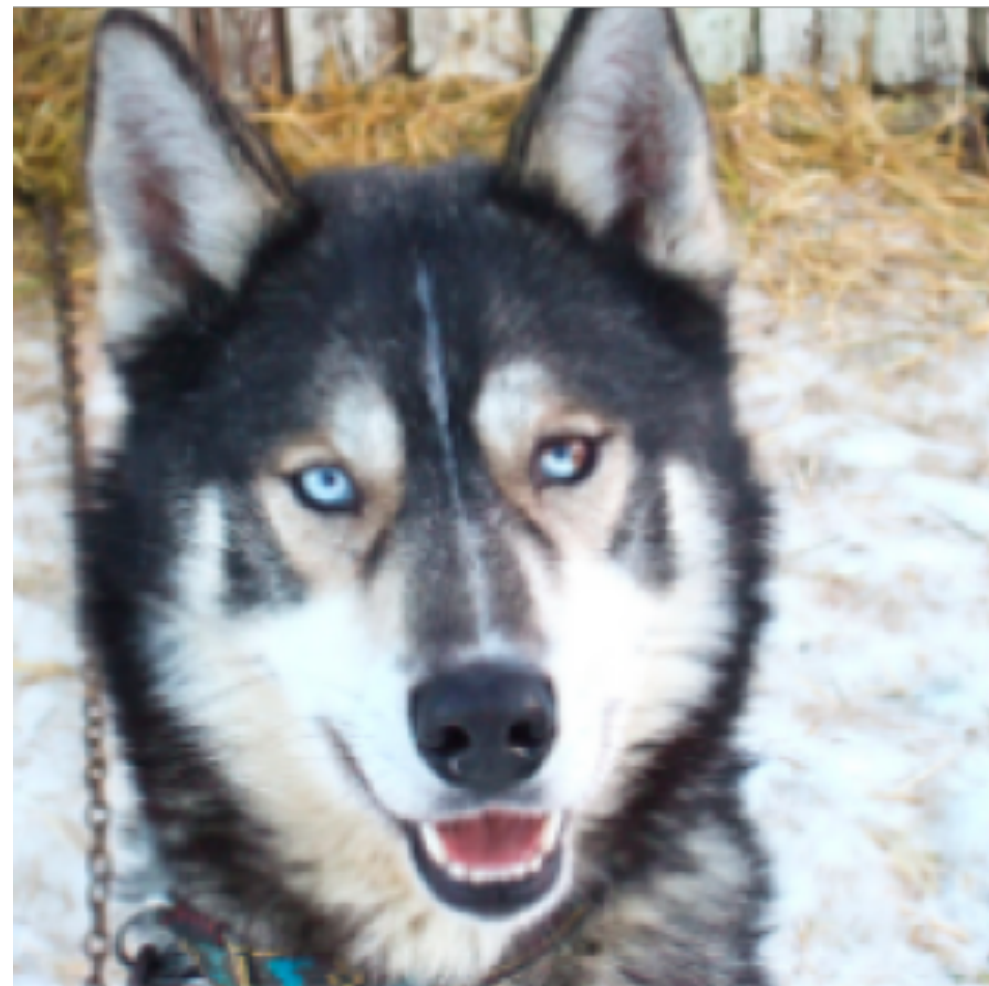**Example: Husky or Wolf? … and why?**

Consider an example image classification task about distinguishing between **husky dogs** and **wolves**

Local explanations allow to spot cases where the model is **right** for the **wrong reasons**
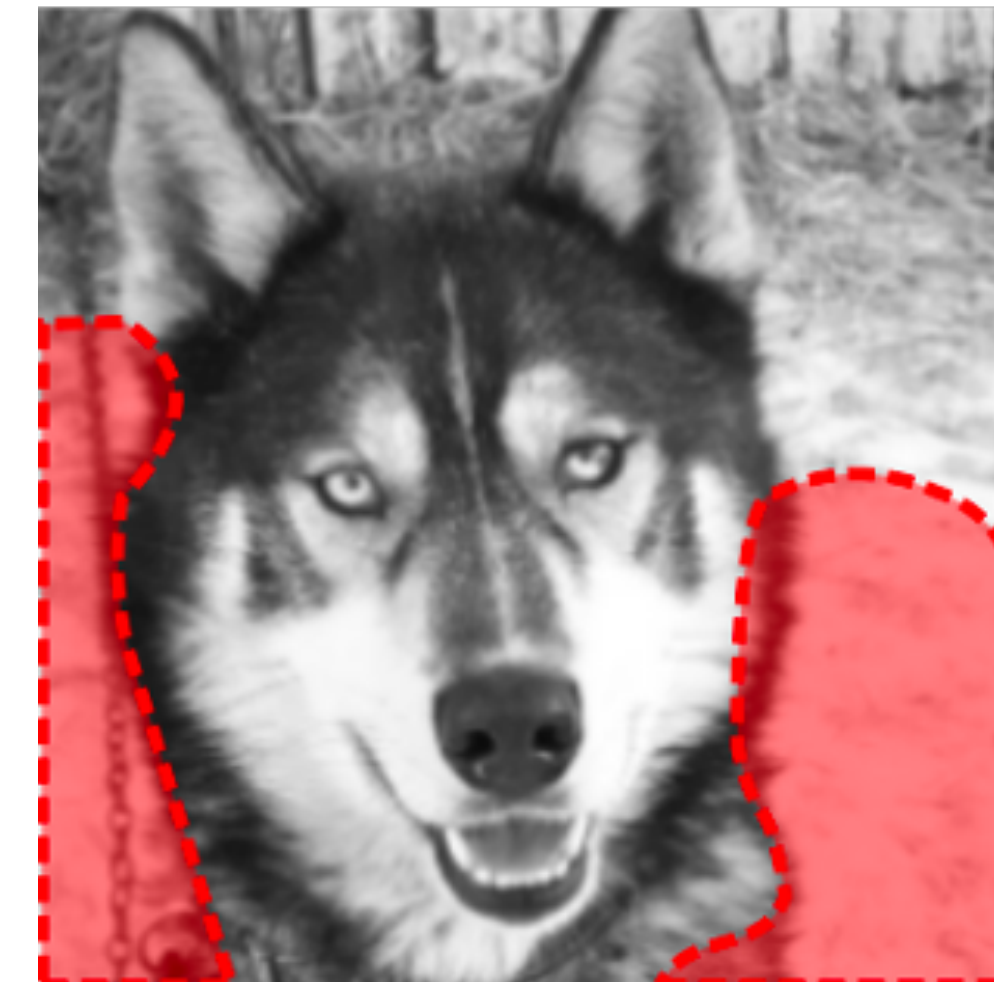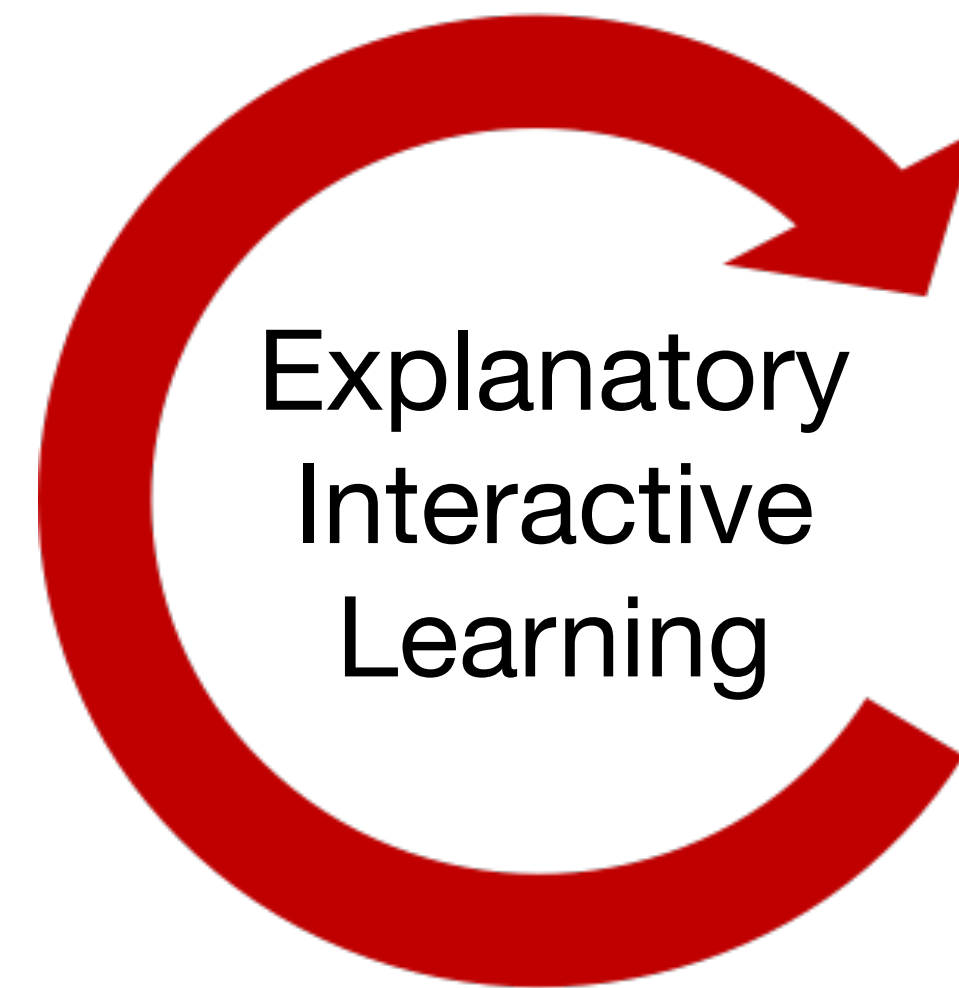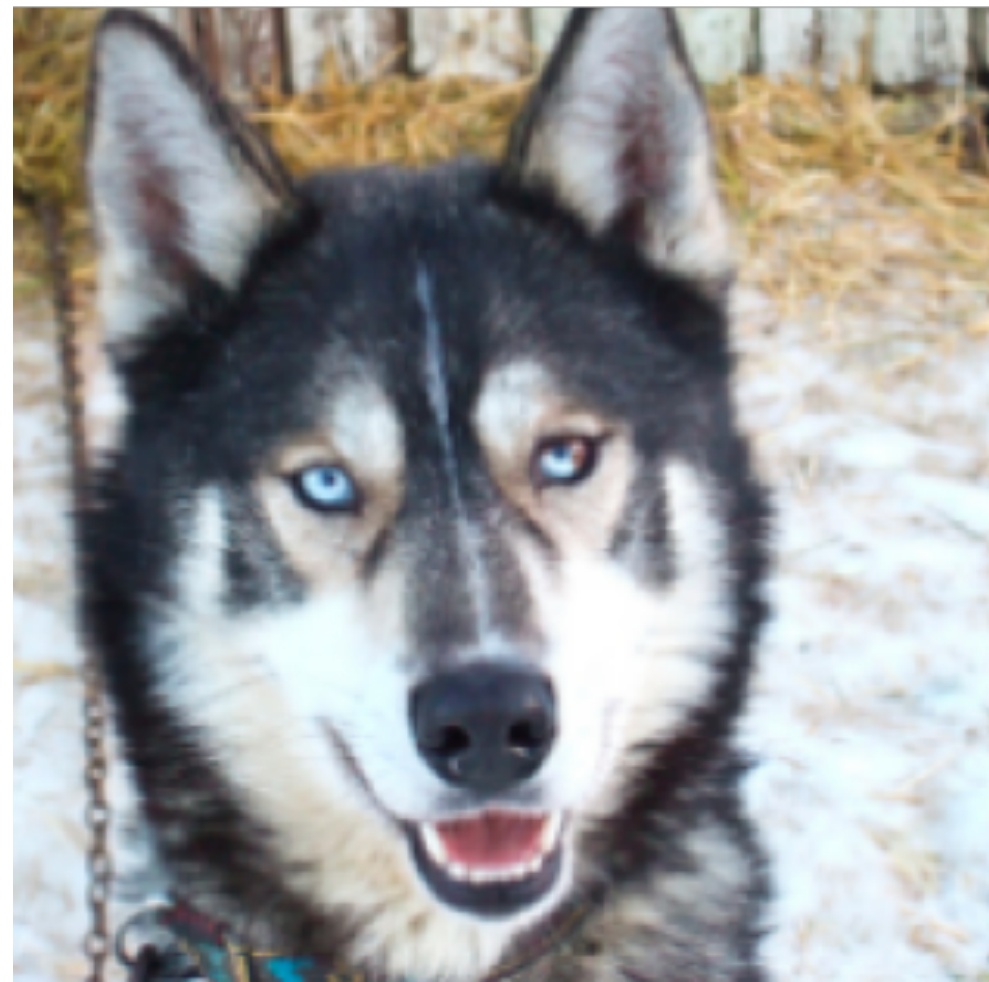
Example: Husky or Wolf? … and **why**? … and **feedback**!

It is a husky, but **not because** of the **highlighted pixels**!

Example: Husky or Wolf? … and why? … and feedback!

Explanatory Interactive Learning

A.    **Explain predictions** to users (*competence*, *understandability*)
B.    Allow user to **correct explanations** (*directability*)

Of course, the pattern is not exclusive to images

Example: plant phenotyping



(a) Scientific Task

Predict if a plant suffers from biotic stress using hyperspectral images.

Plant tissue & expected reasons

Hyperspectral cube

Spectral signatures

Reflectance

nm

Leaf
Diseased Spot

(b) Machine Learning: Often no interaction. Expert just provides the hyperspectral data and the labels.

*Schramowski et al. Nature Machine Intelligence 2020*

Of course, the pattern is not exclusive to images

Example: plant phenotypic



(a) Scientific Task

Plant tissue & expected reasons

Hyperspectral cube

Spectral signatures

Predict if a plant suffers from biotic stress using hyperspectral images.

(b) Machine Learning: Often no interaction. Expert just provides the hyperspectral data and the labels.

(c) "Clever Hans"-like behavior: Explainable AI (XAI) methods may reveal "Clever Hans" behavior of the learnt machine, in particular when using deep neural networks.
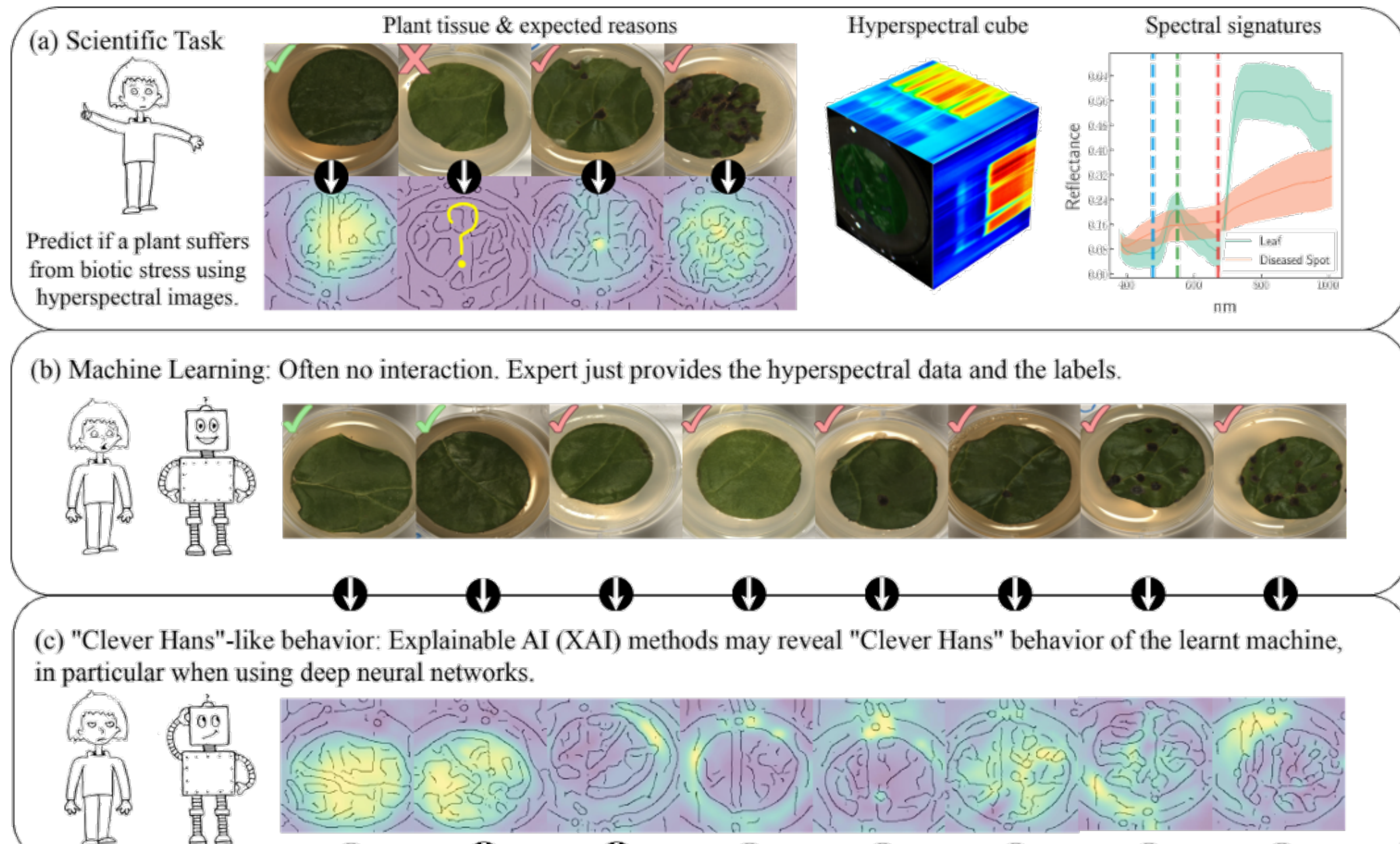
*Schramowski et al. Nature Machine Intelligence 2020*

nature machine intell

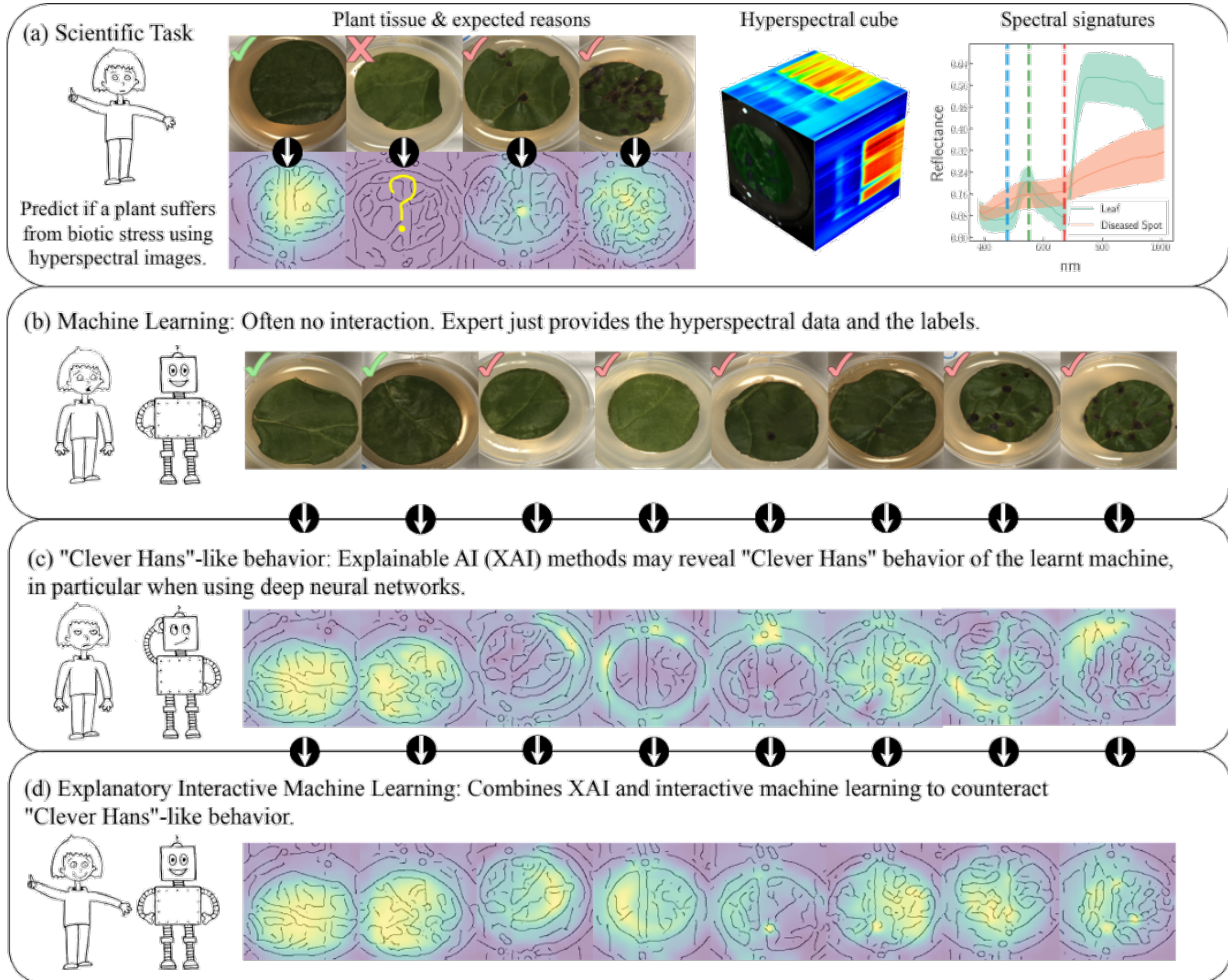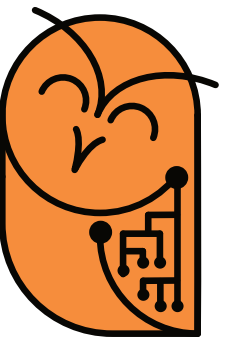Of course, the pattern is not exclusive to images

Example: plant phenotypic

*Schramowski et al. Nature Machine Intelligence 2020*



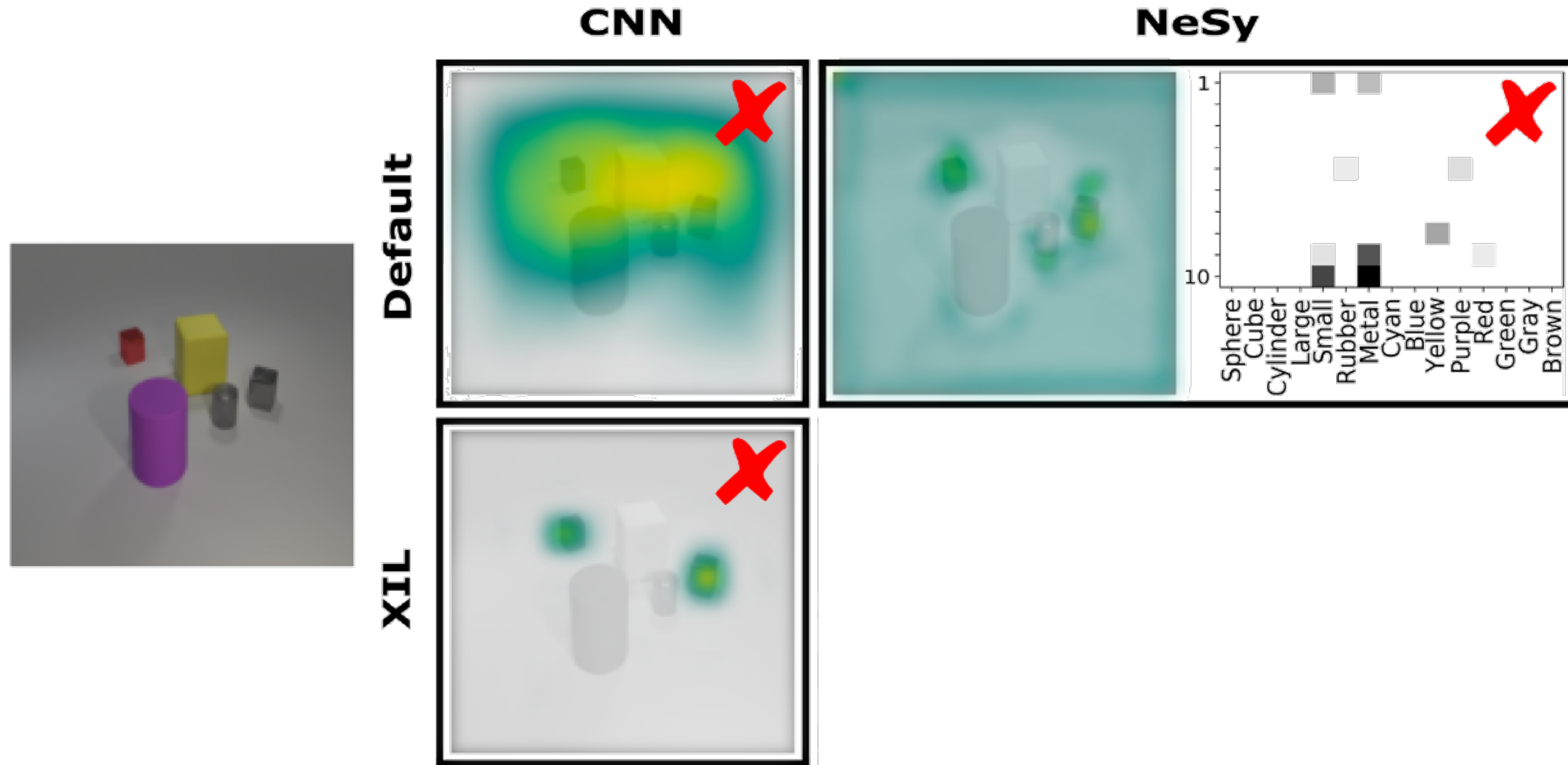(a) Scientific Task

Predict if a plant suffers from biotic stress using hyperspectral images.

Plant tissue & expected reasons

Hyperspectral cube

Spectral signatures

(b) Machine Learning: Often no interaction. Expert just provides the hyperspectral data and the labels.

(c) "Clever Hans"-like behavior: Explainable AI (XAI) methods may reveal "Clever Hans" behavior of the learnt machine, in particular when using deep neural networks.

(d) Explanatory Interactive Machine Learning: Combines XAI and interactive machine learning to counteract "Clever Hans"-like behavior.

# Unfortunately, visual explanations alone are not all we need either
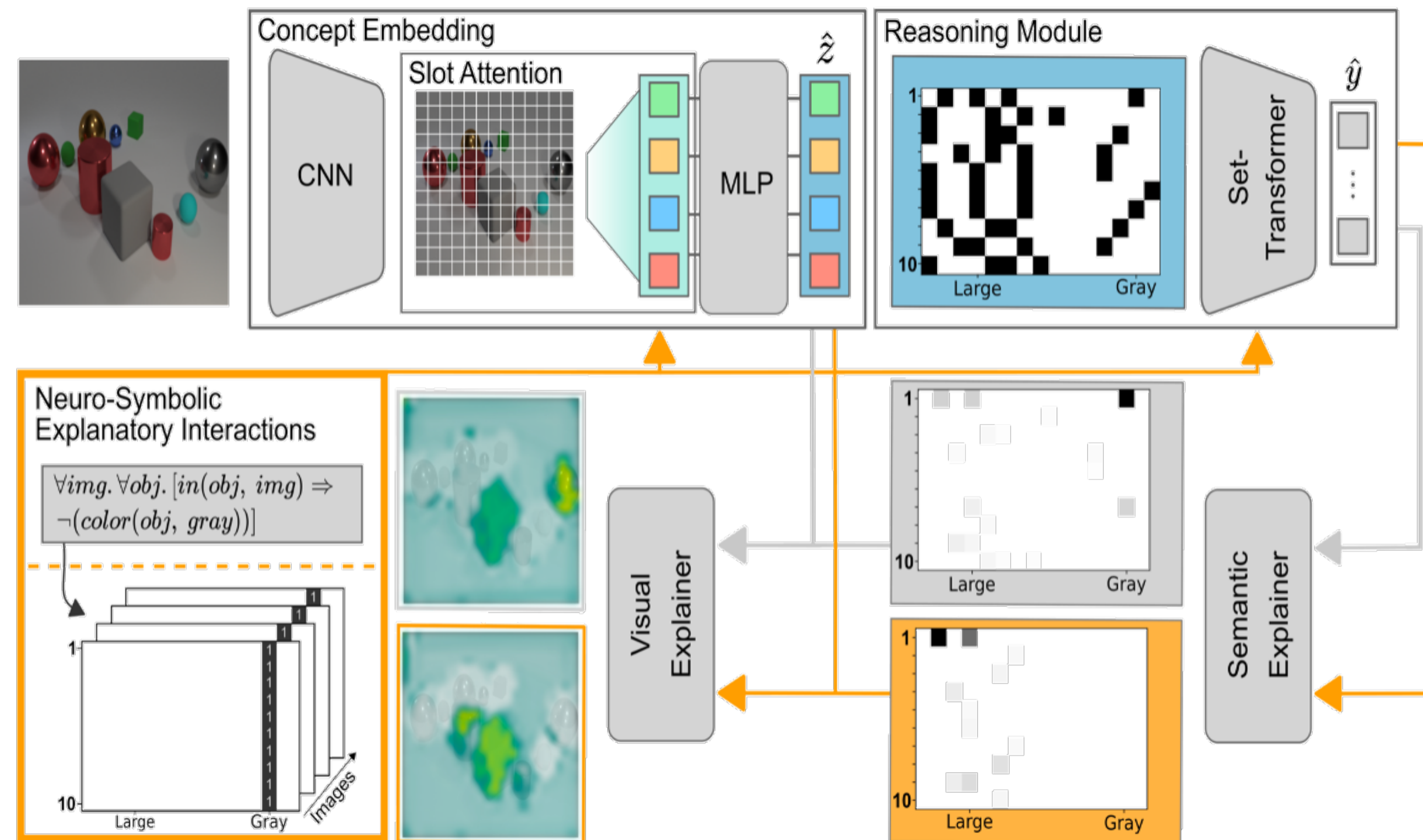


**CNN**

Default

XIL

**Underlying concept:** the image contains a **large cube** & a **large cylinder**

Unfortunately, visual explanations alone are not all we need either
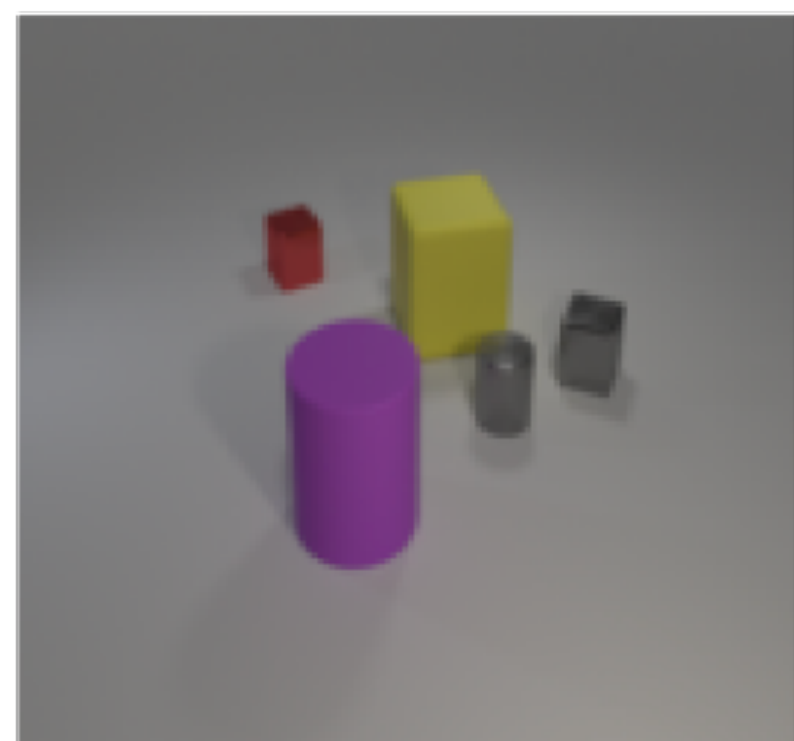
# Right for the right neuro-symbolic reasons!



**Combine human & machine intelligence via an explanatory loss term**

$$\lambda \sum_{i=1}^{N} r(A_i^v, \hat{e}_i^h) + (1 - \lambda) \sum_{i=1}^{N} r(A_i^s, \hat{e}_i^g)$$

*Ross et al IJCAI 2017*
*Teso, Kersting AIES 2019*
*Selvaraju et al ICCV 2019*
*Schramowski et al Nature MI 2020*

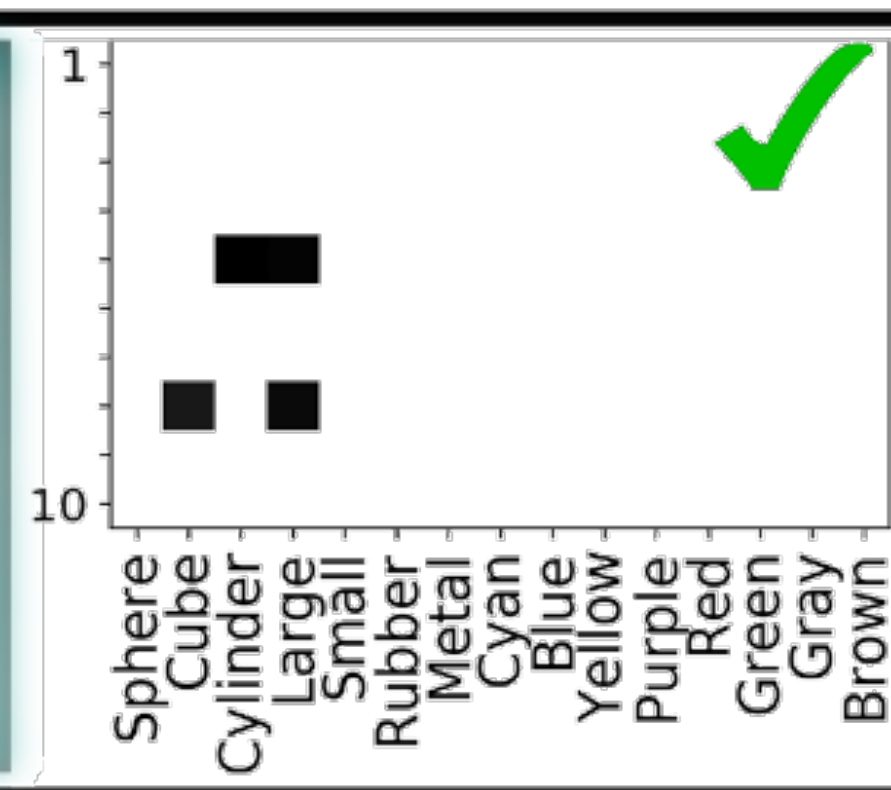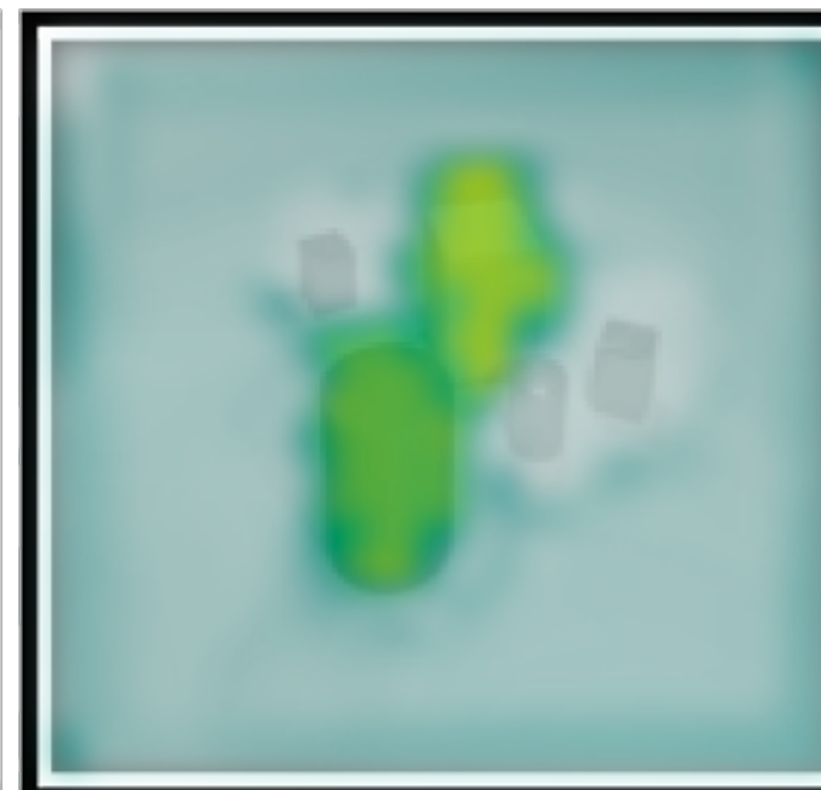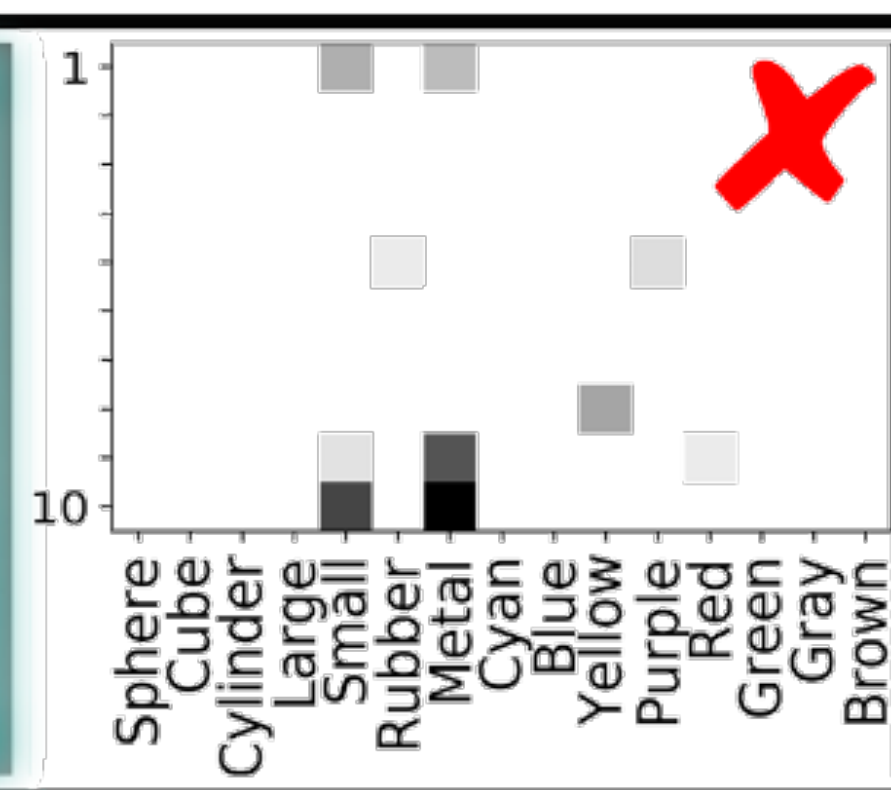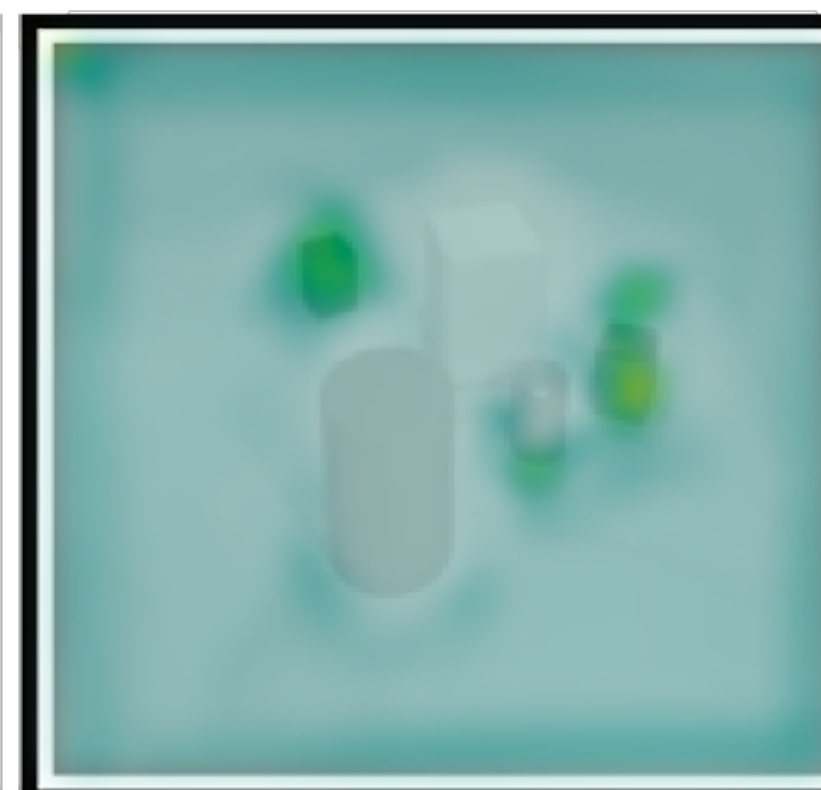Symbolic feedback:
"Never base your decision on gray cubes!"

# What if there is a new concept

# Or even a new task

*Welinder et al, Caltech-UCSD Birds 200*   *Kudithipudi et al, Nature MI 2022*

# Another challenge: deep models **don't know when they don't know**

## Dataset classification



**(Deep neural) models are overconfident on datasets they were never trained on**

*Matan 1990*
*Mundt et al. ICCV 2019*
*Mundt et al. Journal of Imaging 2022*

…and even worse, neural nets fail to learn sequentially. Humans don't!



Training order

Blocked training

Concept One

Concept Two

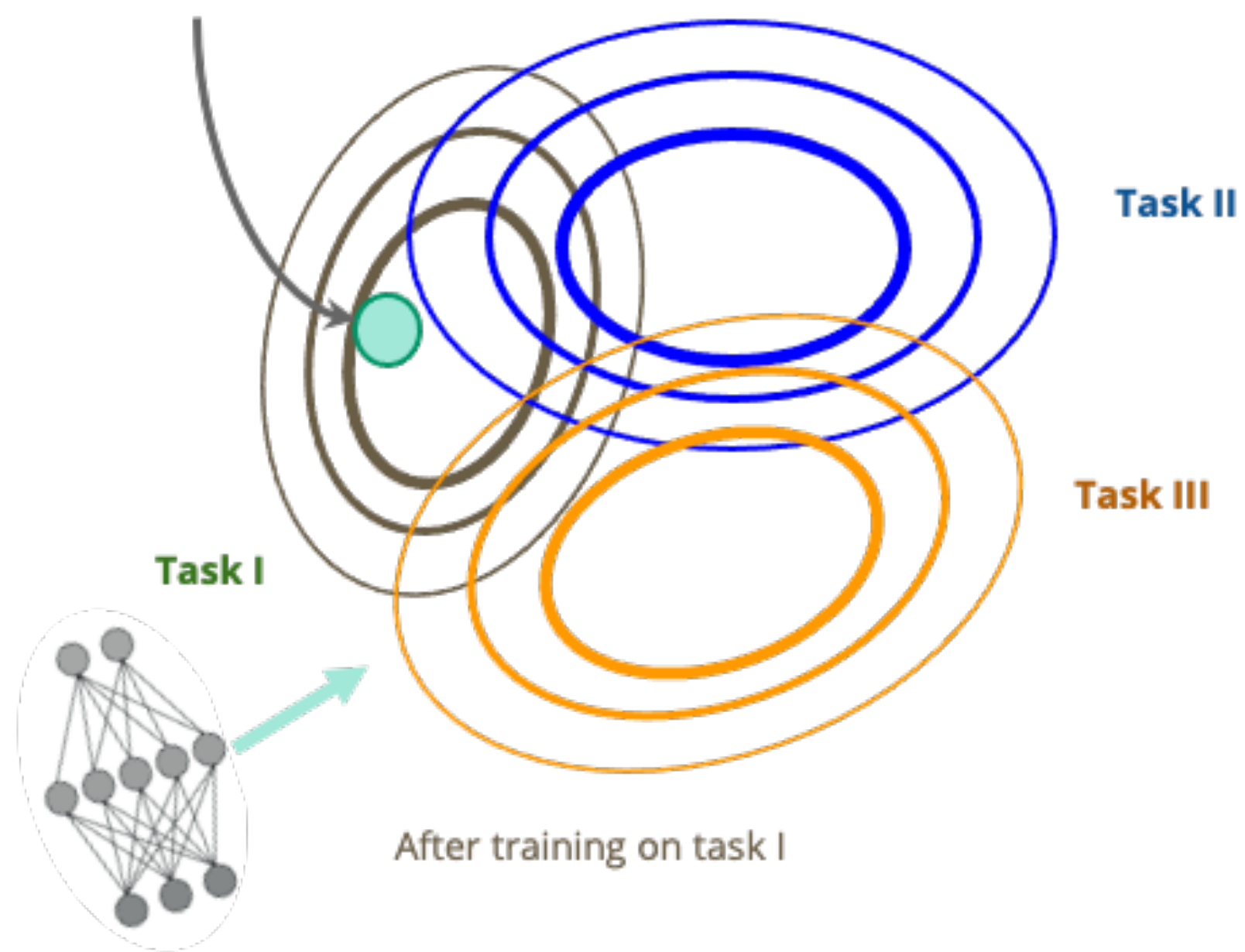Training Accuracy 100 0

Training Accuracy 100 0

Training Epoch

Catastrophic Interference (McCloskey & Cohen 89)

*Adapted from Flesch et al 2022*

# Why might neural networks be so forgetful? Is it surprising?



Task I

Task II

Task III

After training on task I

# Why might neural networks be so forgetful? Is it surprising?



Task I

Task II

Task III

After training on task I

Task I

Task II

Task III

After training on task II

# Why might neural networks be so forgetful? Is it surprising?



Task II

Task III

Task I

After training on task I

Task II

Task III

Task I

After training on task II

Task II

Task III

Task I

# We can mitigate this problem with generative models



Continual Machine Learning That Can Identify What It Doesn't Yet Know

$\epsilon \sim \mathcal{N}(0, I)$

Encoder — $\theta$ — $\sigma_\theta$, $\mu_\theta$

EVT open set — $q_\theta(z|x)$

Decoder — $\phi_d$ — $p_\phi(x|z)$

$x$ — $\times$ — $+$ — $z$ — $x'$

outlier free generation

$\tau \; \kappa \; \lambda$

$\Omega_\rho$

reject or query

$p_\xi(y|z)$

inverse sampling $\Omega_\rho^{-1}(z)$ to pick exemplars

Weibull Distribution $\rho_c$

$\xi$ — $y'$

Linear Classifier

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}, \boldsymbol{\xi}) = \mathbb{E}_{q_{\boldsymbol{\theta}}(\boldsymbol{z}|\boldsymbol{x})}\left[\log p_{\boldsymbol{\phi}}(\boldsymbol{x}|\boldsymbol{z}) + \log p_{\boldsymbol{\xi}}(\boldsymbol{y}|\boldsymbol{z})\right] - \beta KL(q_{\boldsymbol{\theta}}(\boldsymbol{z}|\boldsymbol{x}) \,||\, p(\boldsymbol{z}))]$$

-> we learn how to encode data into generative factors & in turn how to decode (generate) these into data

# We can mitigate this problem with generative models





*Continual Machine Learning That Can Identify What It Doesn't Yet Know*

**We can then measure similarity to factors we have already observed & replay knowledge from already seen ones**

We can then measure similarity to factors we have already observed & replay knowledge from already seen ones

Accept prediction & generate known data

Avoid ambiguous & reject unknown data

Assigned outlier percentage

1.31  1.96  4.84  10.22  80.53  100
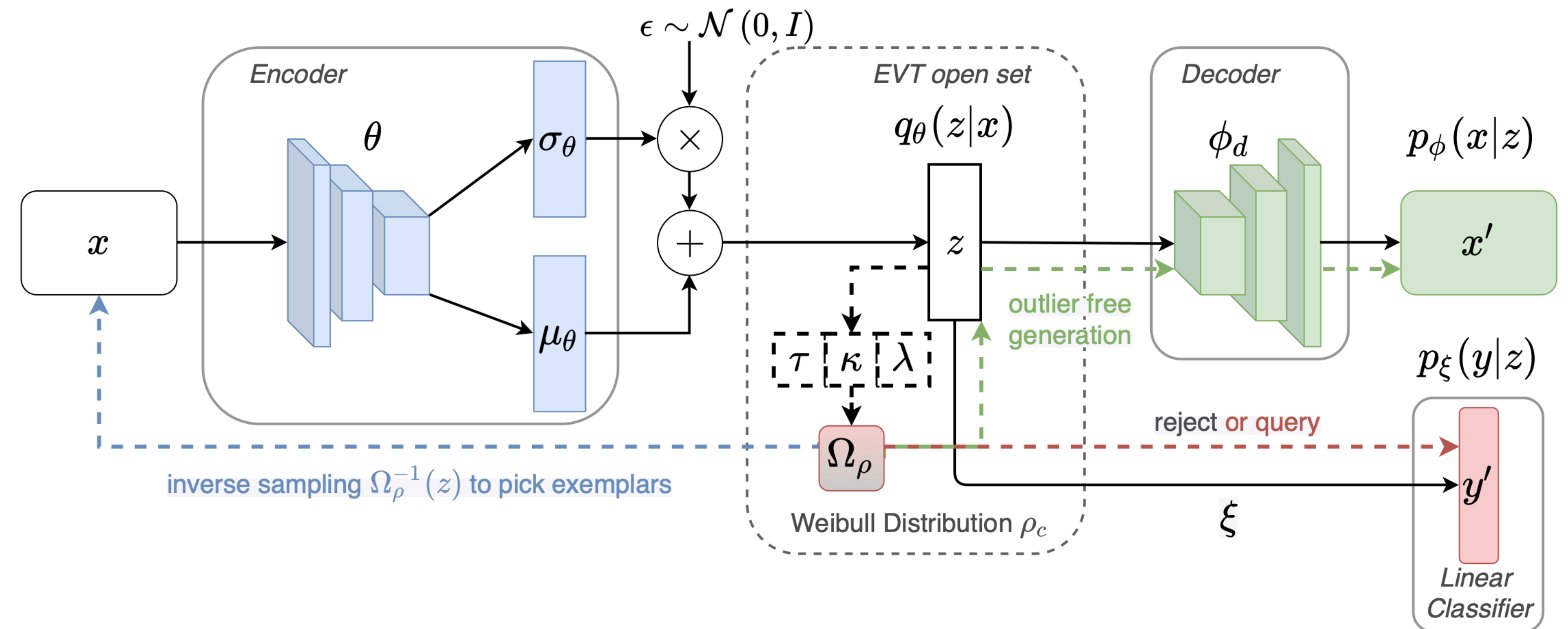
0.73  1.15  4.62  11.70  76.19  99.9

*Mundt et al. Journal of Imaging 2022*
*Hong & Mundt Neural Networks 2022*

# So ultimately, the prevalent common ML pipeline is unrealistic

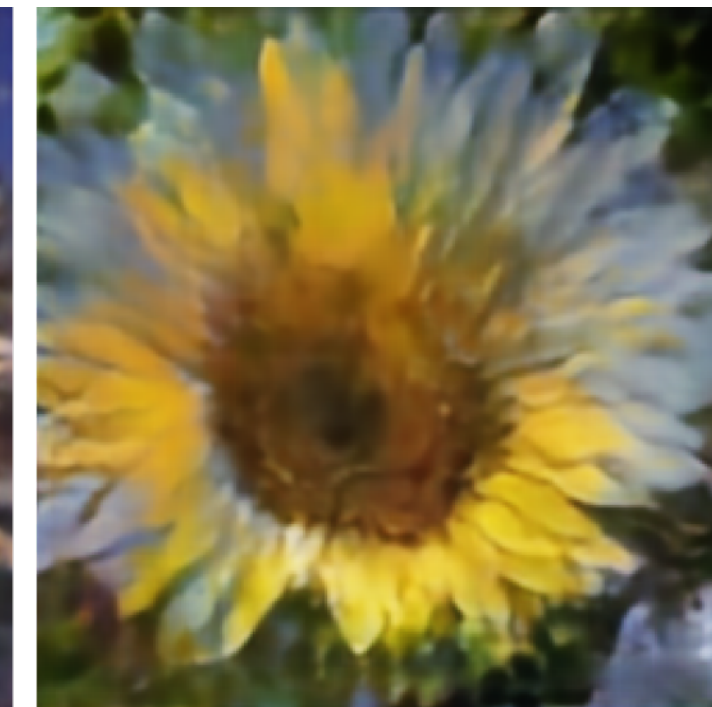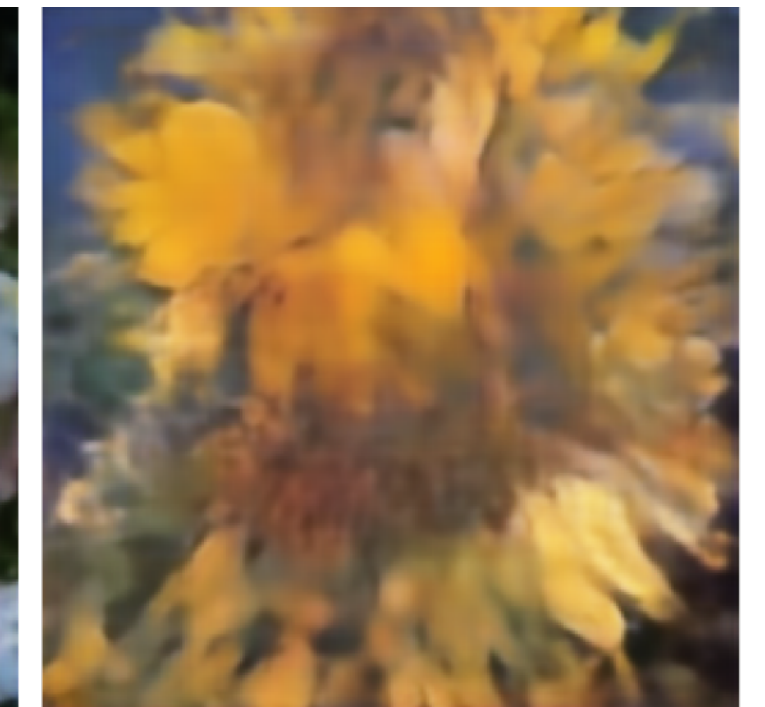Identify the problem to be solved and create a clear objective.

Preparing data is a crucial step and involves building workflows to clean, match and blend the data.

Data is fed as input and the algorithm configured with the required parameters. A percent of the data can be utilized to train the model.

Publish the prepared experiment as a web service, so applications can use the model

| Define objective | Collect Data | Prepare Data | Select Algorithm | Train Model | Test Model | Integrate Model |

Collect data from hospitals, health insurance companies, social service agencies, police and fire dept.

Depending on the problem to be solved and the type of data, an appropriate algorithm will be chosen.

The remaining data is utilized to test the model for accuracy. Depending on the results, improvements can be performed in the "Train model" and/or "Select Algorithm" phases, iteratively.

Figure from https://www.congrelate.com/get-workflow-machine-learning-images/

in reality it may look much more like this with sensor drifts, novel concepts in data, focus on wrong reasons, or even new tasks

Label Data

- World with Knowns (K) & Unknowns Unknowns (UU)
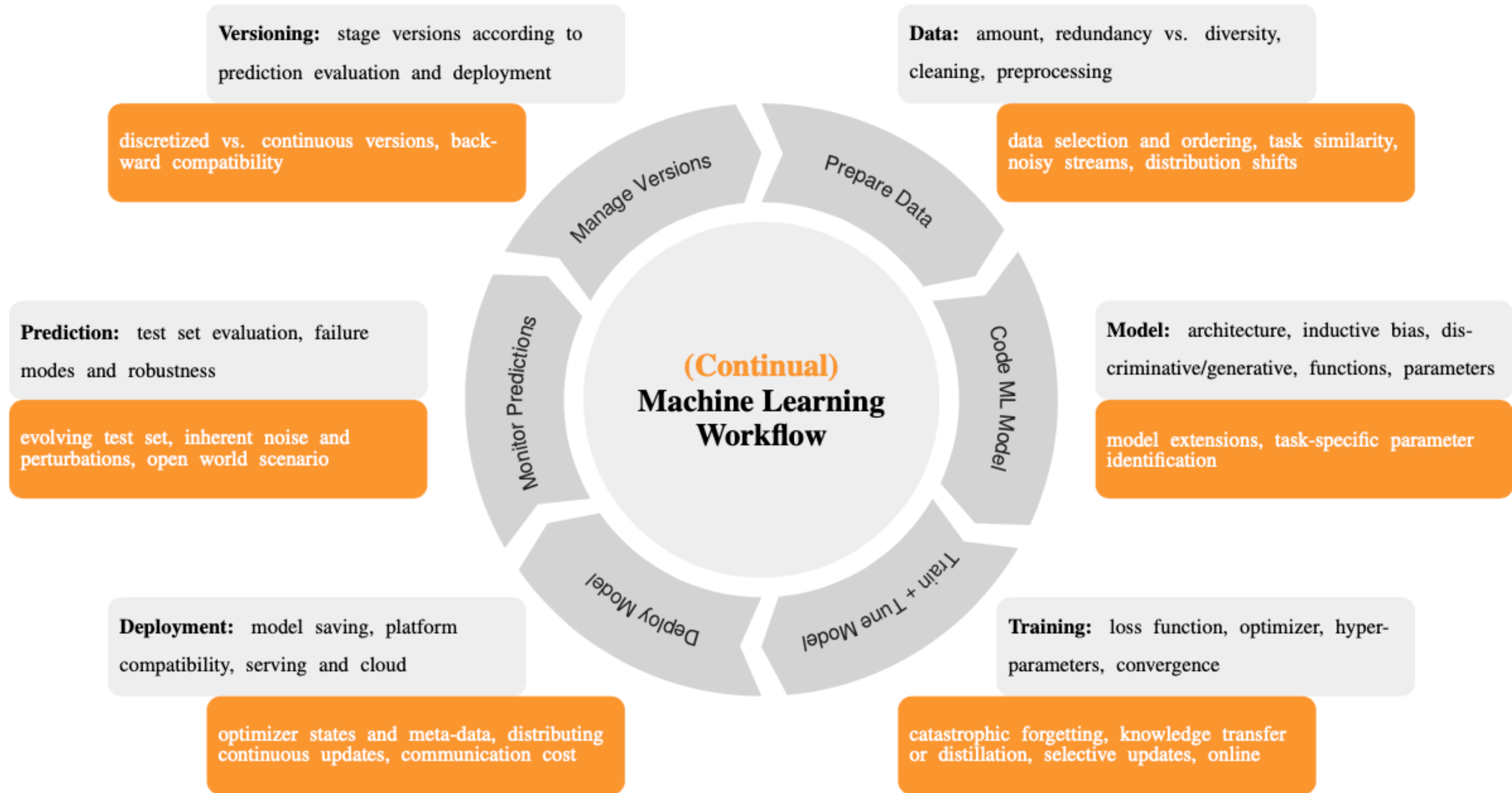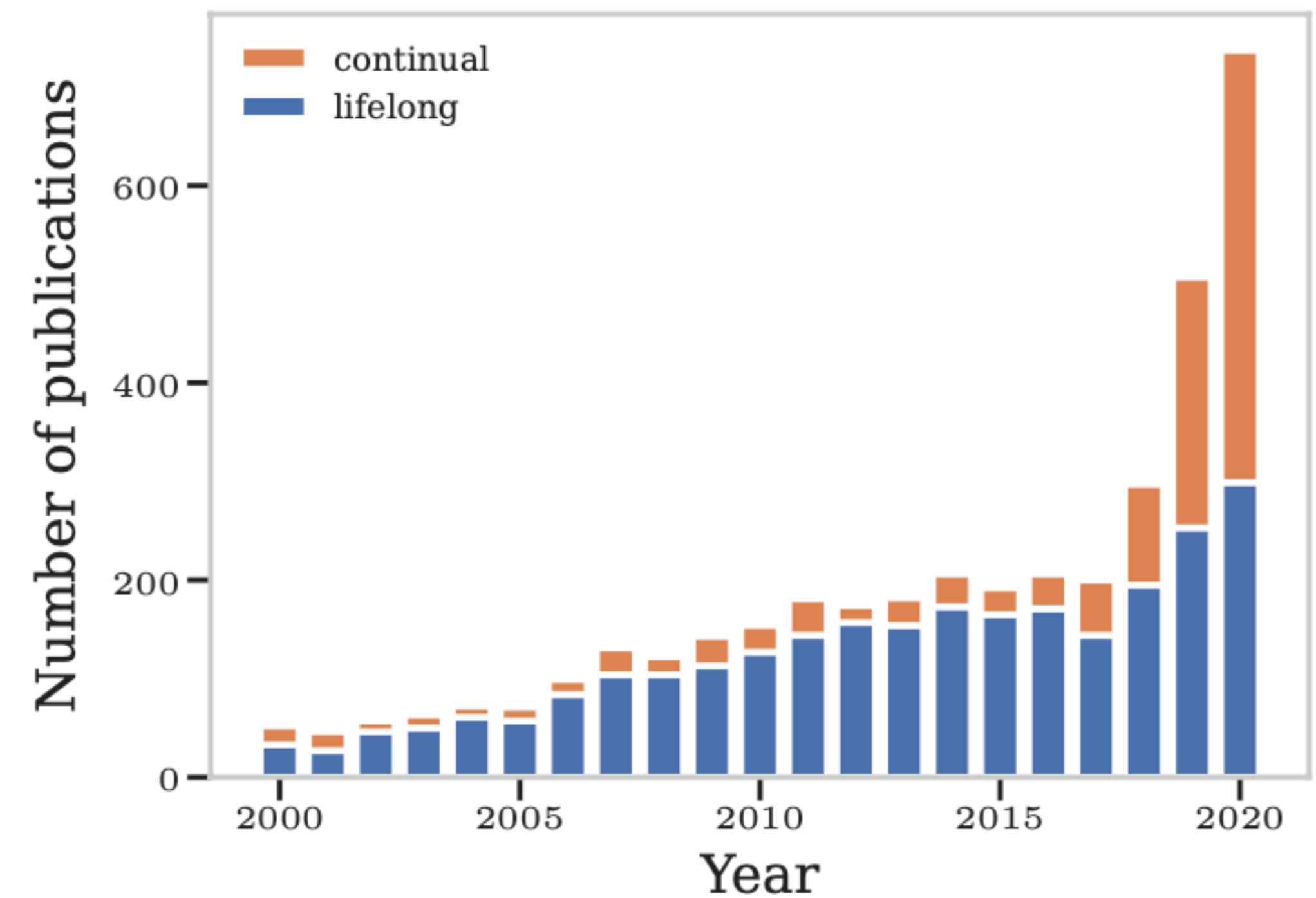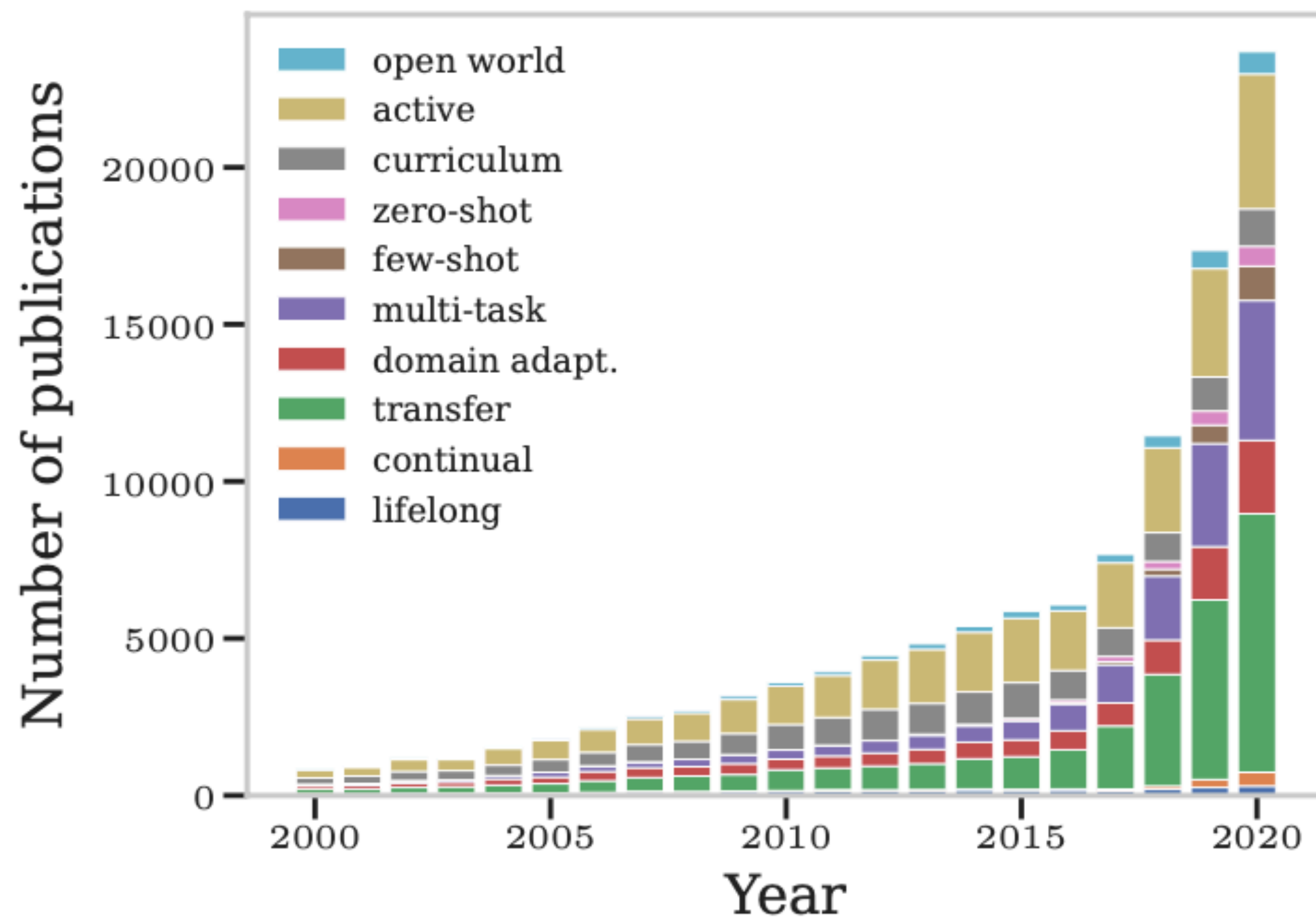
Recognize as Known

Detect as Unknown

- NU: Novel Unknowns

- LU: Labeled Unknowns

Incremental Learning

- K: Known

Scale

And perhaps even further
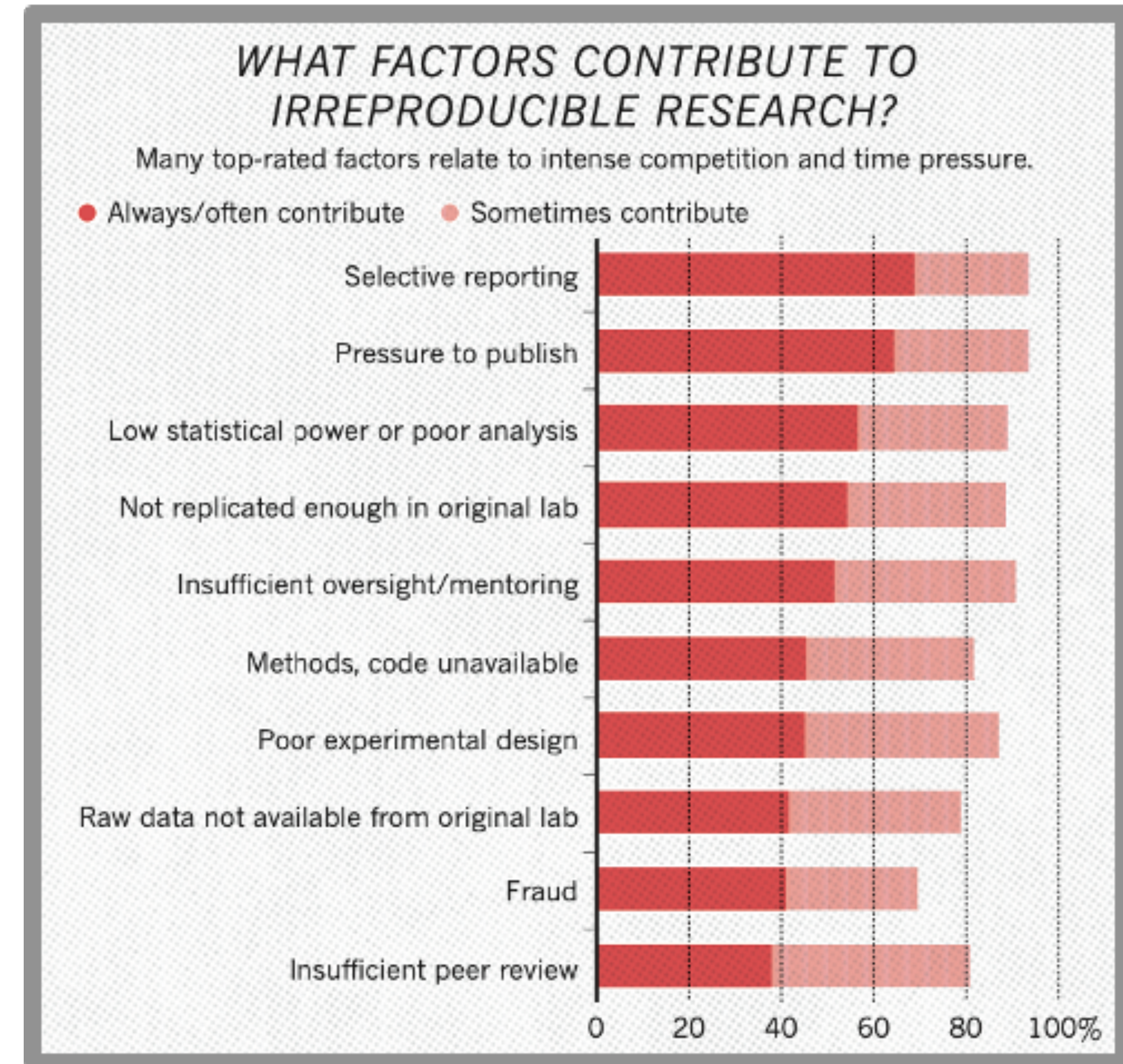
Bendale & Boult, CVPR 2015
Mundt et al, Neural Networks 2023

# And in reality, likely much much more complex than six simple steps



**Versioning:** stage versions according to prediction evaluation and deployment

discretized vs. continuous versions, backward compatibility

**Data:** amount, redundancy vs. diversity, cleaning, preprocessing

data selection and ordering, task similarity, noisy streams, distribution shifts

**Prediction:** test set evaluation, failure modes and robustness

evolving test set, inherent noise and perturbations, open world scenario

**Model:** architecture, inductive bias, discriminative/generative, functions, parameters

model extensions, task-specific parameter identification

**Deployment:** model saving, platform compatibility, serving and cloud

optimizer states and meta-data, distributing continuous updates, communication cost

**Training:** loss function, optimizer, hyperparameters, convergence

catastrophic forgetting, knowledge transfer or distillation, selective updates, online

Manage Versions

Prepare Data

Monitor Predictions

Code ML Model

Deploy Model

Train + Tune Model

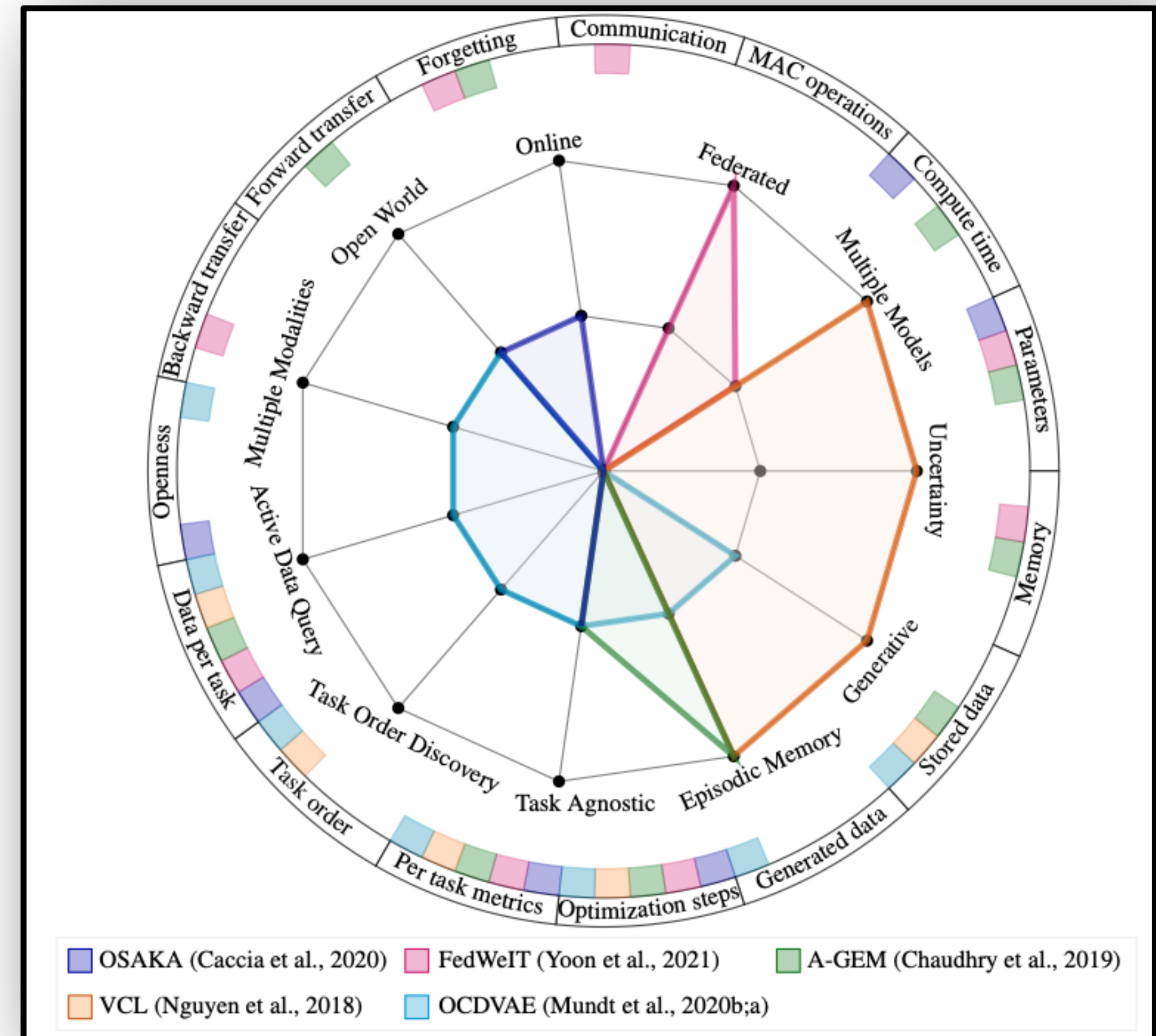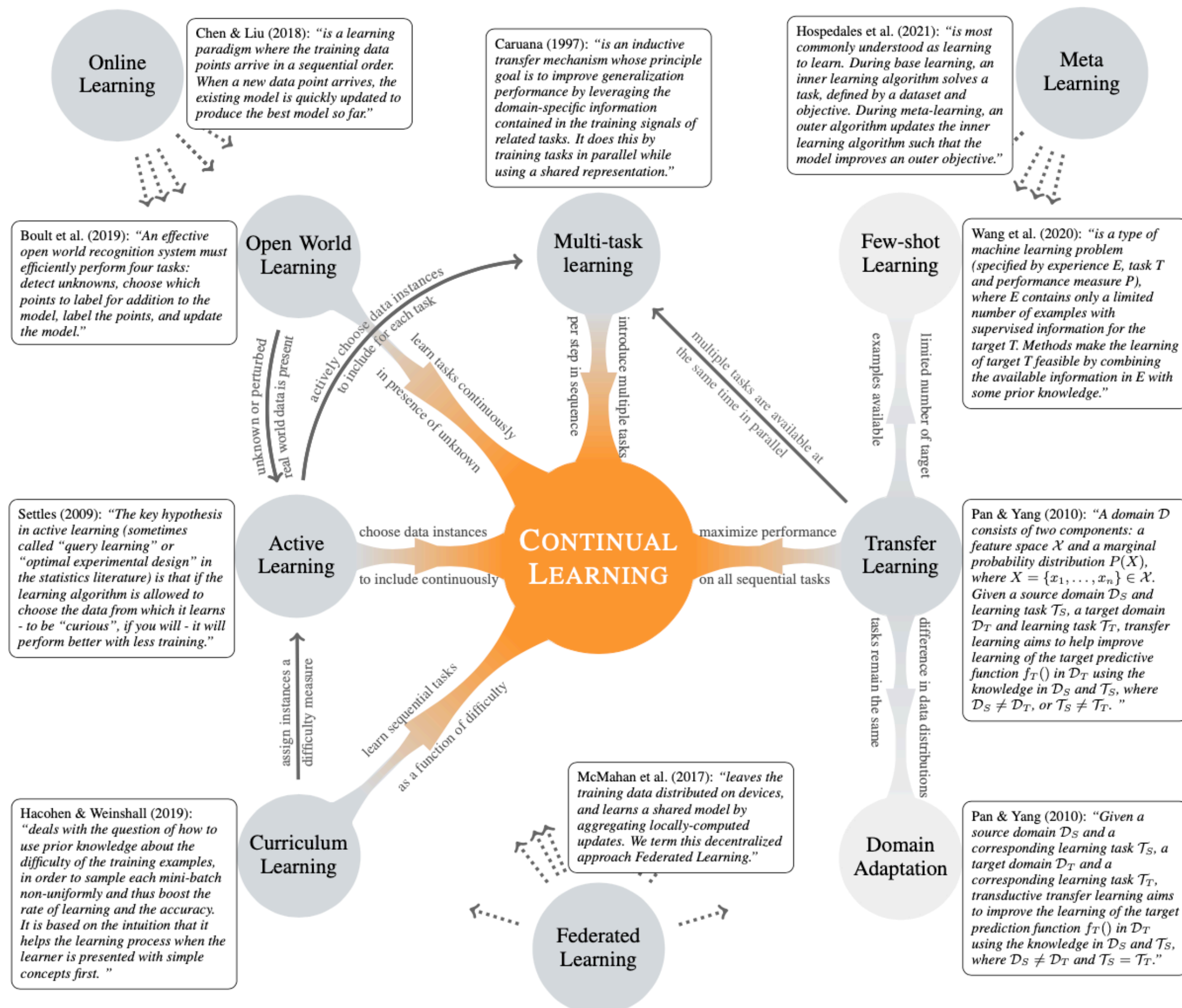**(Continual) Machine Learning Workflow**

Mundt et al, ICLR 2022

**Not yet convinced? What if our <u>application doesn't require *all* of these factors?</u> Pragmatically: Why should we still care?**

Mundt et al, ICLR 2022

Baker, Nature 2016

# Systems are complex! Assumptions & evaluation setups often collapse aspects into scalar measured quantities. Let's acknowledge this fact & make it transparent



Mundt et al, ICLR 2022

# Summary & take-aways

1. Standard deep neural networks are not **<u>right for the right reasons</u>**

2. Standard deep neural networks **<u>don't know what they don't know</u>**

3. Standard deep neural networks are **<u>bad at learning sequentially/continually</u>**

But very powerful -> **generative + symbolic + human**

With this we can enable **<u>explanatory interactive</u>** and **<u>continual learning</u>** & in the process **<u>make the machine learning workflow reflect our real-world desiderata</u>** more accurately than current static benchmarking
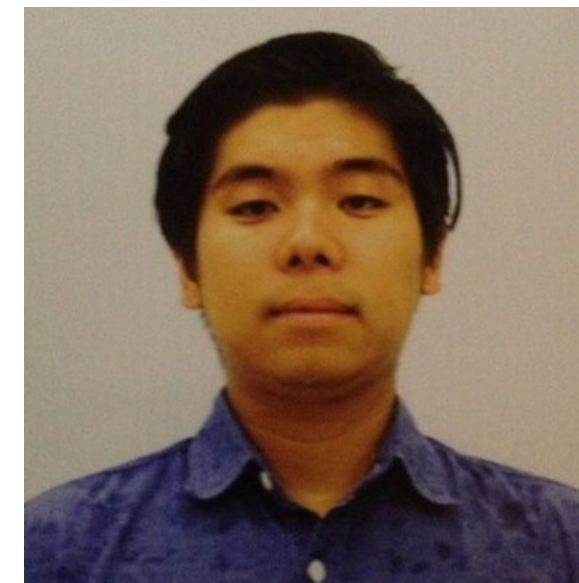
# Thanks to ... and many more!

Visvanathan Ramesh
Uni Frankfurt, hessian.AI

Kristian Kersting
TU Darmstadt, hessian.AI
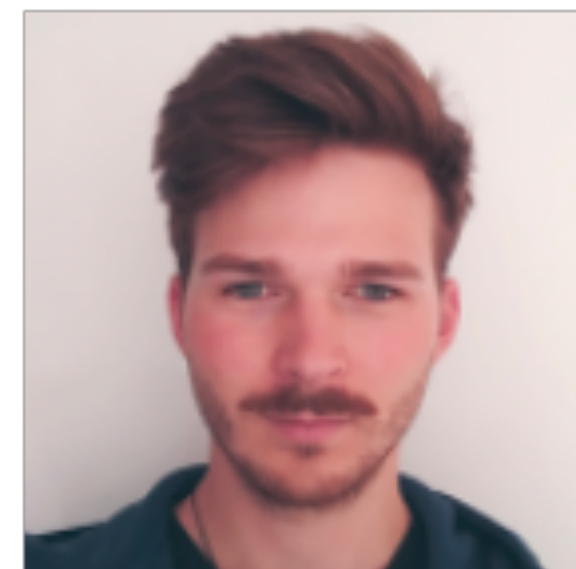
Iuliia Pliushch
Uni Frankfurt

Yongwon Hong,
Yonsei University

Vincenzo Lomonaco
Uni Pisa, ContinualAI

Patrick Schramowski
TU Darmstadt, DFKI

Wolfgang Stammer
TU Darmstadt

Xiaoting Shao
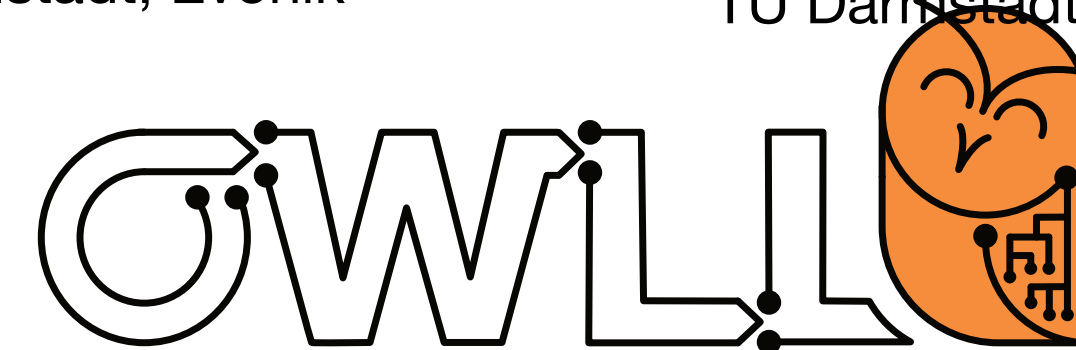TU Darmstadt, Evonik

Steven Braun
TU Darmstadt

Tyler Hayes
NAVER, ContinualAI

TECHNISCHE
UNIVERSITÄT
DARMSTADT

hessian.AI

OWLI

ContinualAI