

# Addressing the lack of reliability in machine learning

## A novel probabilistic viewpoint on self-supervised data representations

Michael Arbel

CRCN and ISFP recruitment campaign at Inria  
Grenoble Rhône-Alpes

23 may 2022



# Academic background

---

2021-present	<b>Starting Research Position at Thoth team</b> Inria Grenoble Rhône-Alpes, France
2016-2021	<b>Ph.D. student at the Gatsby Unit</b> University College London, United Kingdom
2015-2016	<b>Computer vision engineer</b> Prophesee, Paris, France
2014-2015	<b>Student at Ecole Normale Supérieure de Cachan</b> Master M2 (MVA), Cachan, France
2011-2015	<b>Student at Ecole Polytechnique</b> INRIA Grenoble Rhône-Alpes

# Ph.D. contributions

Learning structure from data with generative models

# Learning structure from data with generative models

## Explicit Generative Models

$$p_{\theta}(X) = \frac{1}{Z_{\theta}} e^{-E_{\theta}(x)}$$

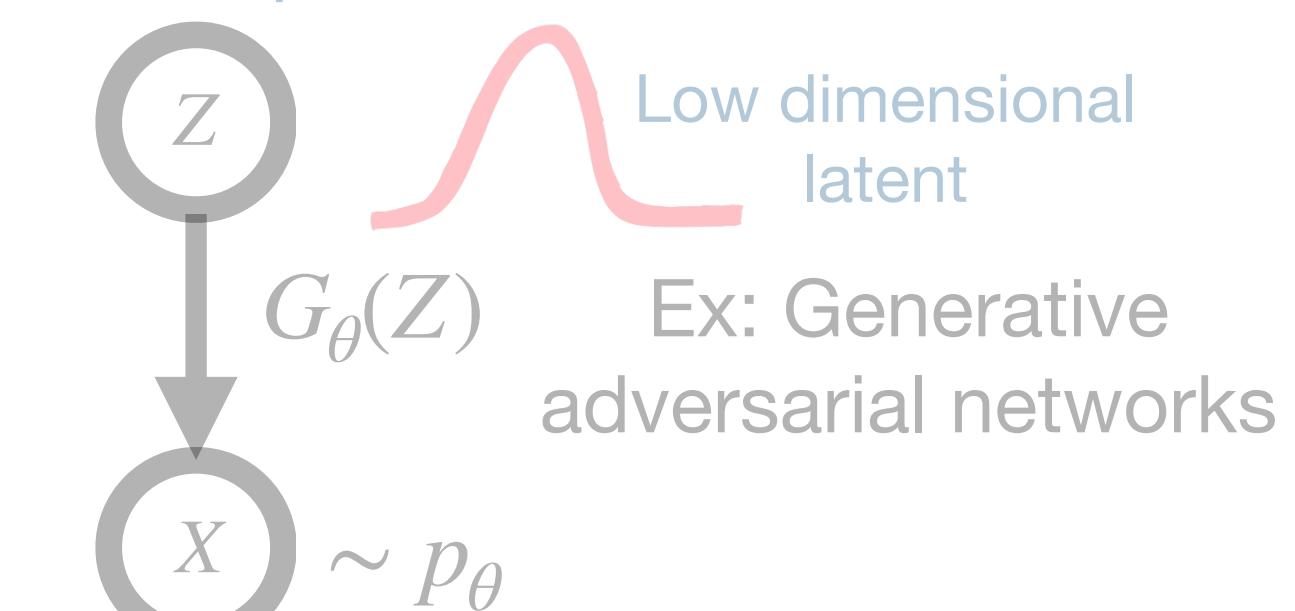
Ex: Energy-based models

- ✓ Capture multimodality
- Hard to sample from



$$\min_{\theta} \mathcal{L}(p_{\theta}) \approx p_{\theta}$$

## Implicit Generative Models



- ✓ Capture low dimensionality
- Hard to optimize

# Learning structure from data with generative models

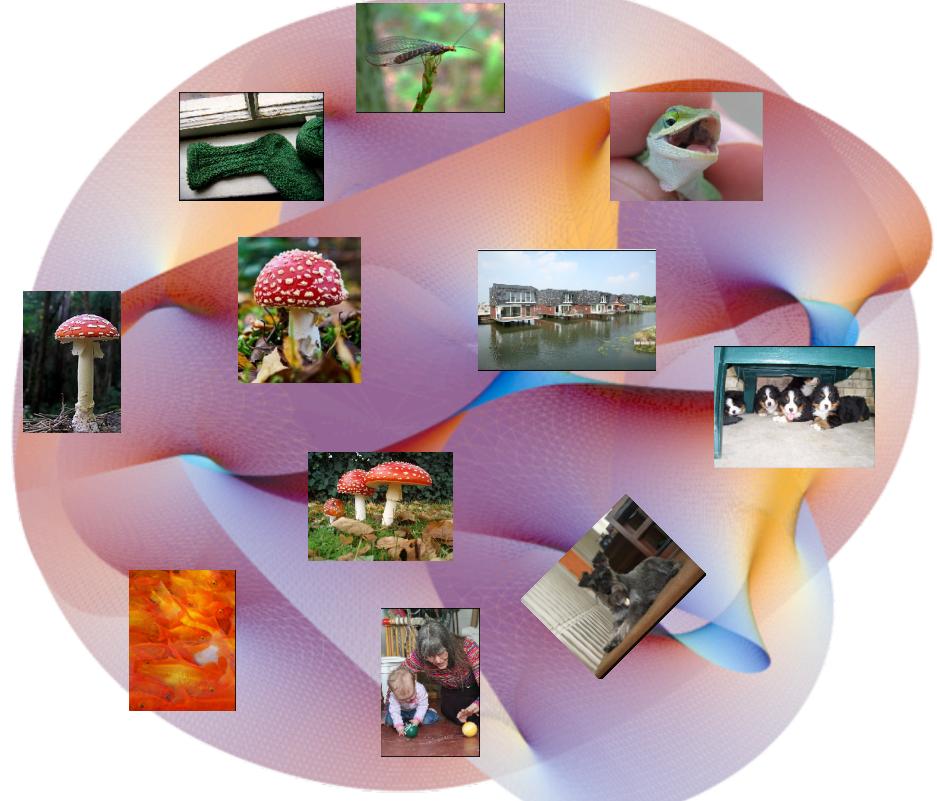
Which model to use?

## Explicit Generative Models

$$p_{\theta}(X) = \frac{1}{Z_{\theta}} e^{-E_{\theta}(x)}$$

Ex: Energy-based models

- ✓ Capture multimodality
- Hard to sample from

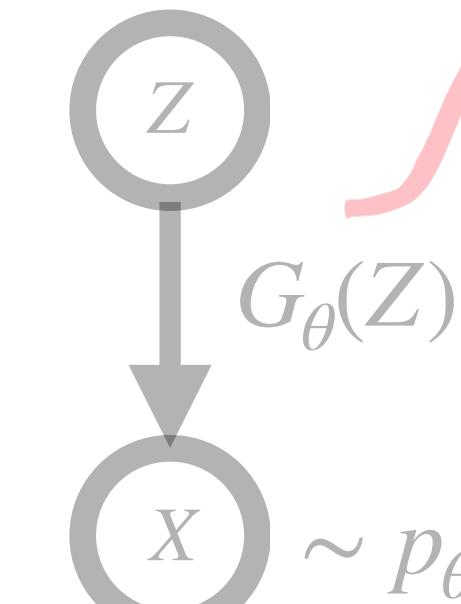


$$\min_{\theta} \mathcal{L}(p_{\theta}) \approx p_{\theta}$$

## Implicit Generative Models

Low dimensional latent

Ex: Generative adversarial networks



- ✓ Capture low dimensionality
- Hard to optimize

# Learning structure from data with generative models

Which model to use?

## Explicit Generative Models

$$p_{\theta}(X) = \frac{1}{Z_{\theta}} e^{-E_{\theta}(x)}$$

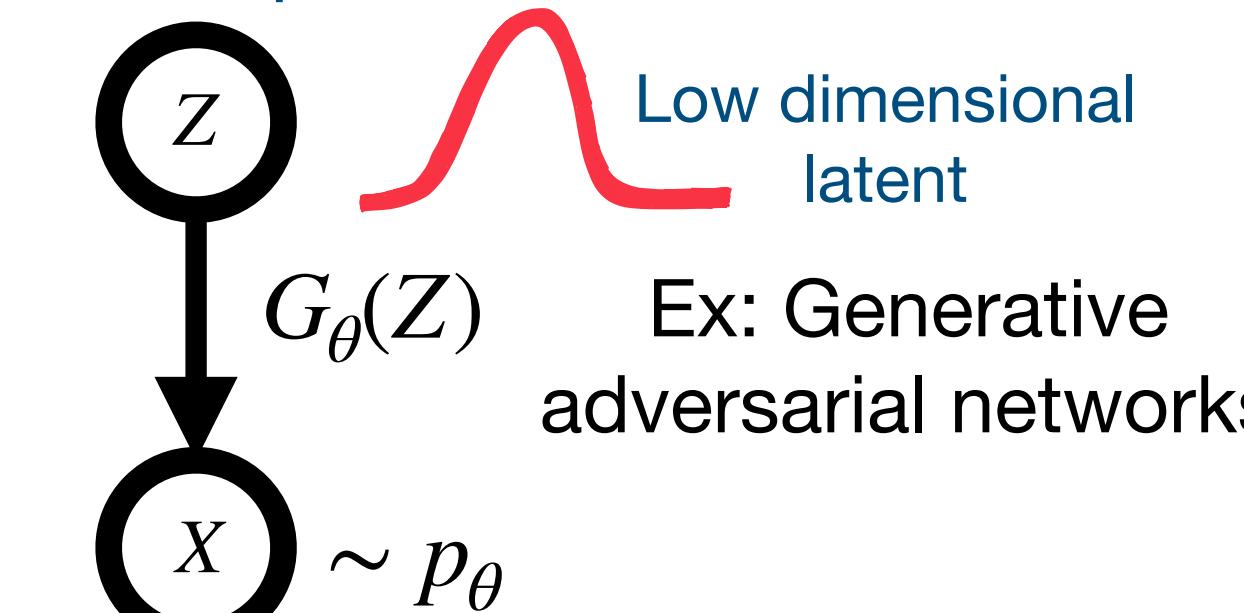
Ex: Energy-based models

- ✓ Capture multimodality
- Hard to sample from



$$\min_{\theta} \mathcal{L}(p_{\theta}) \approx p_{\theta}$$

## Implicit Generative Models



Low dimensional latent

Ex: Generative adversarial networks

- ✓ Capture low dimensionality
- Hard to optimize

# Learning structure from data with generative models

Which model to use?

## Explicit Generative Models

$$p_{\theta}(X) = \frac{1}{Z_{\theta}} e^{-E_{\theta}(x)}$$

Ex: Energy-based models

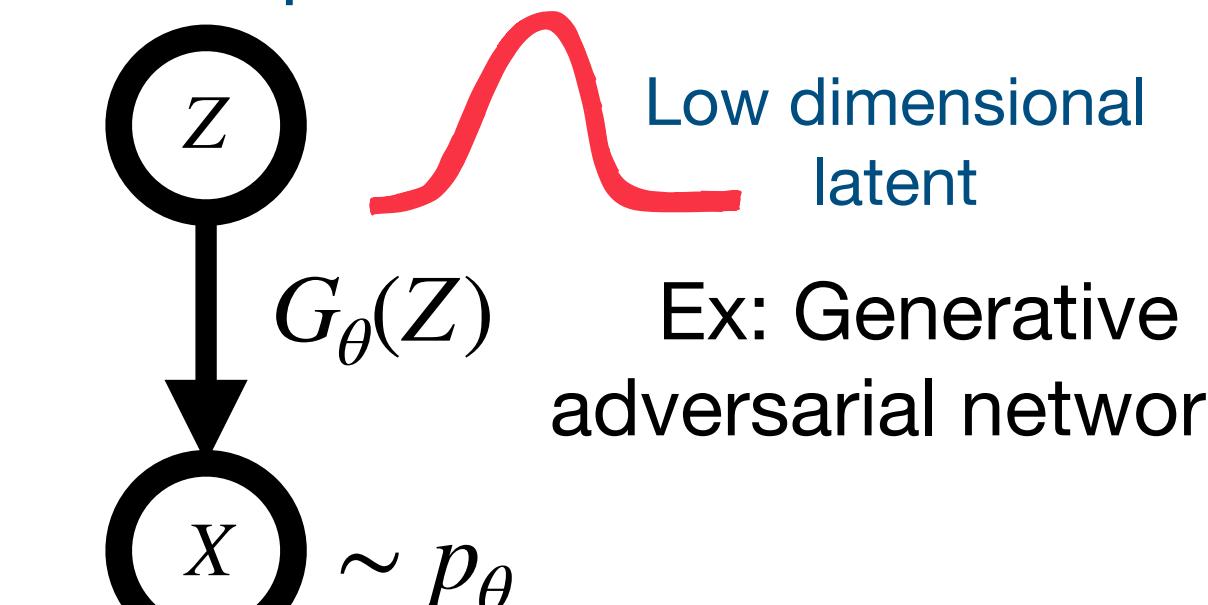
- ✓ Capture multimodality
- ✗ Hard to sample from



$$\min_{\theta} \mathcal{L}(p_{\theta})$$

$$\approx p_{\theta}$$

## Implicit Generative Models



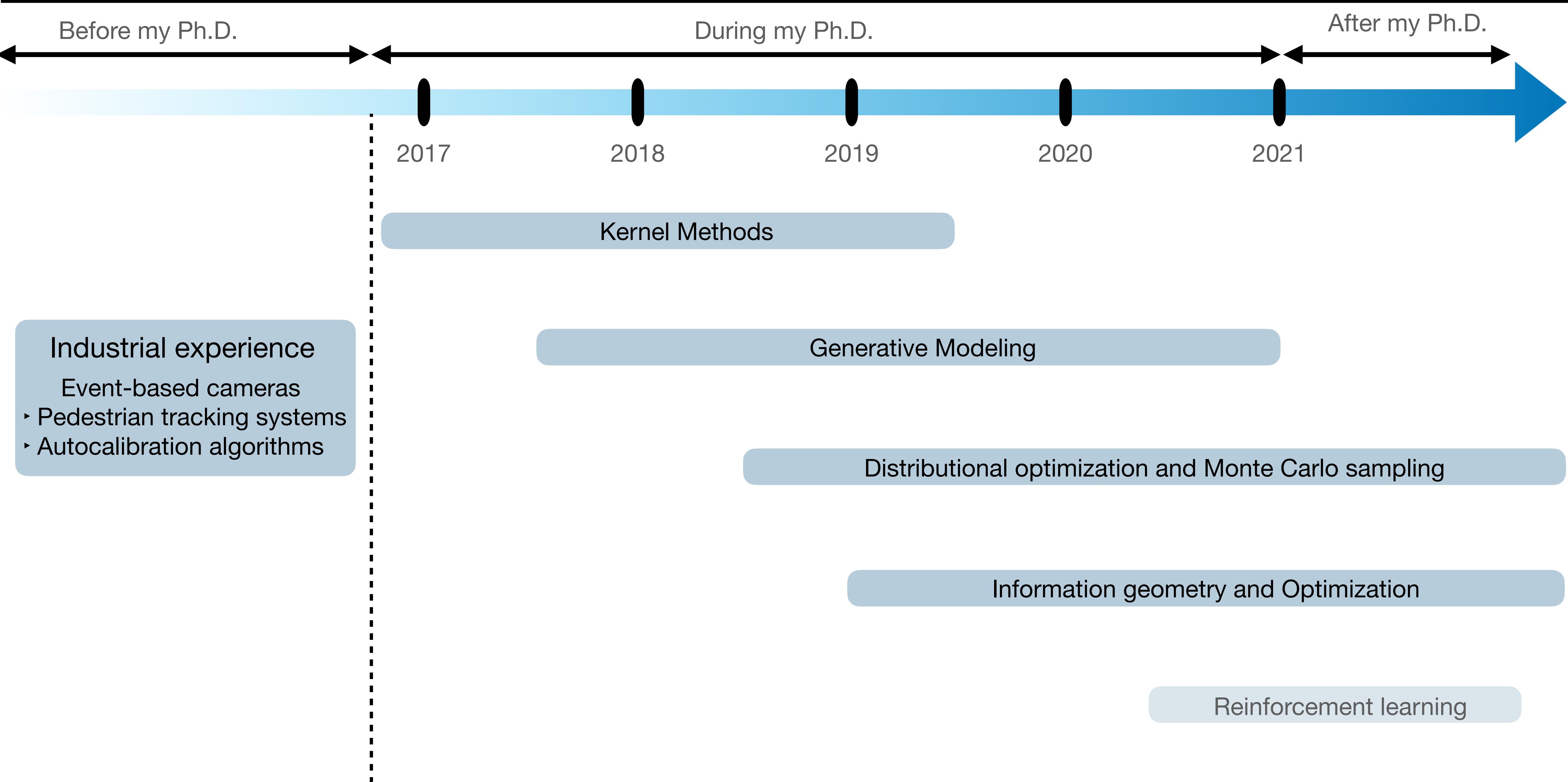
- ✓ Capture low dimensionality
- ✗ Hard to optimize

Contributions in sampling

Contributions in generative modeling

Contributions in optimization

# Major contributions



# Learning structure from data with generative models

## Explicit Generative Models

$$p_{\theta}(X) = \frac{1}{Z_{\theta}} e^{-E_{\theta}(x)}$$

Ex: Energy-based models

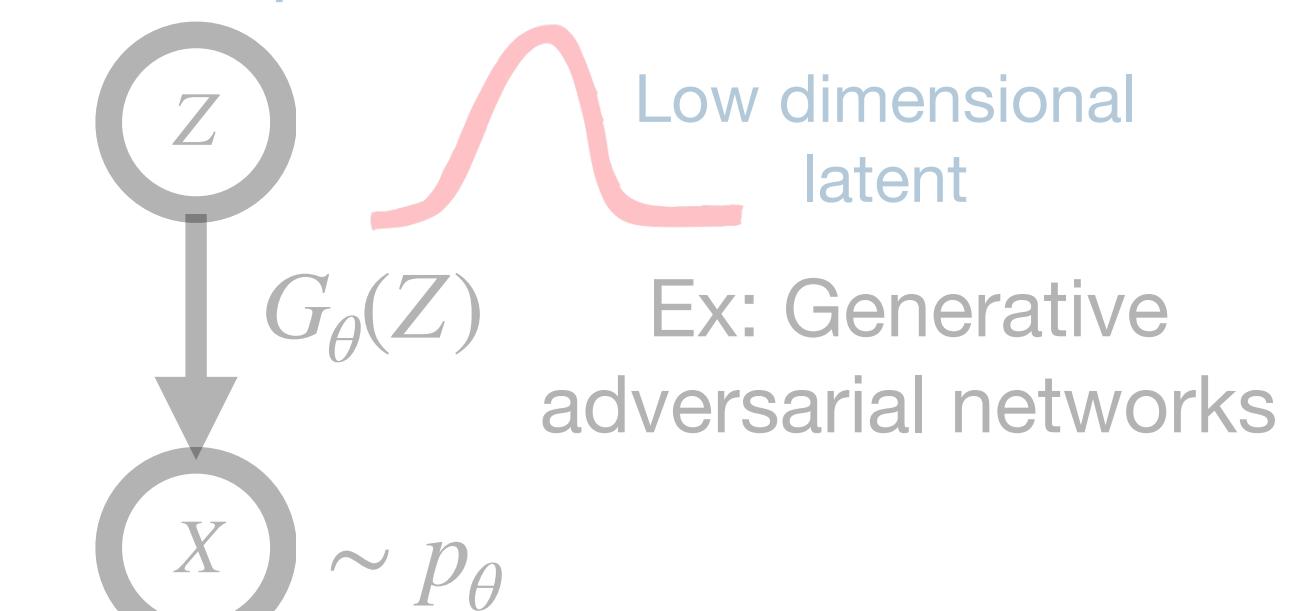
- ✓ Capture multimodality
- Hard to sample from



$$\approx p_{\theta}$$

$$\min_{\theta} \mathcal{L}(p_{\theta})$$

## Implicit Generative Models



- ✓ Capture low dimensionality
- Hard to optimize

Contributions in sampling

Contributions in generative modeling

Contributions in optimization

# Major contributions in generative modelling

Before my Ph.D.

During my Ph.D.

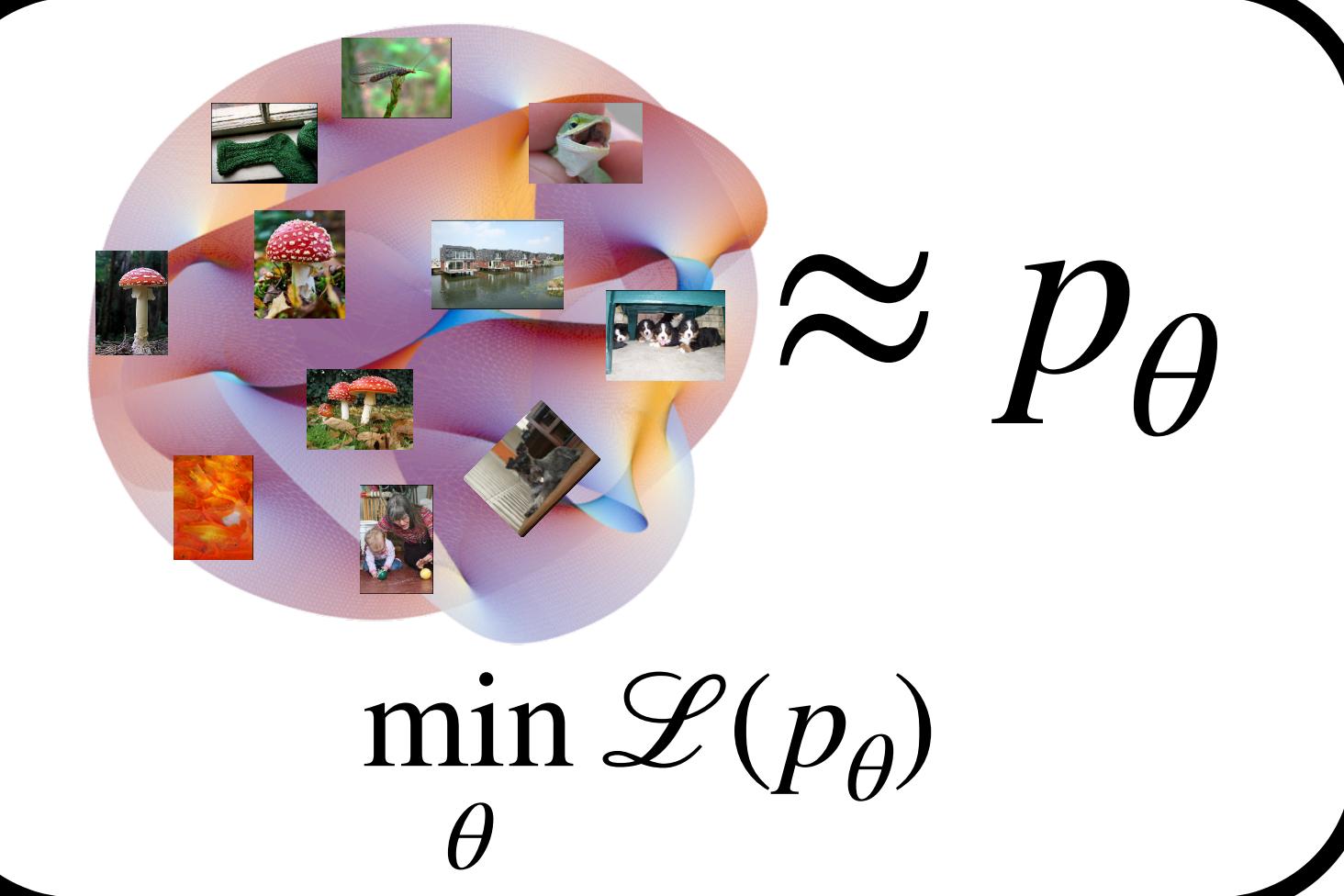
**Unstable** algorithms for learning generative models

Towards **stable and principled** algorithms for learning generative models

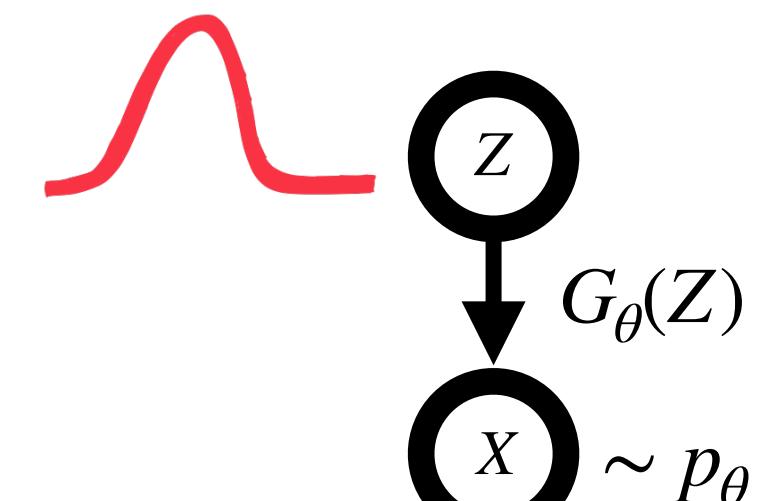
## Explicit Generative Models

$$p_{\theta}(X) = \frac{1}{Z_{\theta}} e^{-E_{\theta}(x)}$$

- ✓ Capture multimodality
- ✗ Hard to sample from



## Implicit Generative Models



- ✓ Capture low dimensionality
- ✗ Hard to optimize

## Generalized Energy-Based models

- ✓ Capture multimodality + low dimensionality
- ✓ A unified theory for learning generative models
- ✓ Practical methods to ensure stable optimization

**Contributions:** A framework unifying both implicit and explicit models

5 publications at top-tier international conferences (NeurIPS, ICLR, AISTATS), more than 600 citations

# Learning structure from data with generative models

## Explicit Generative Models

$$p_{\theta}(X) = \frac{1}{Z_{\theta}} e^{-E_{\theta}(x)}$$

Ex: Energy-based models

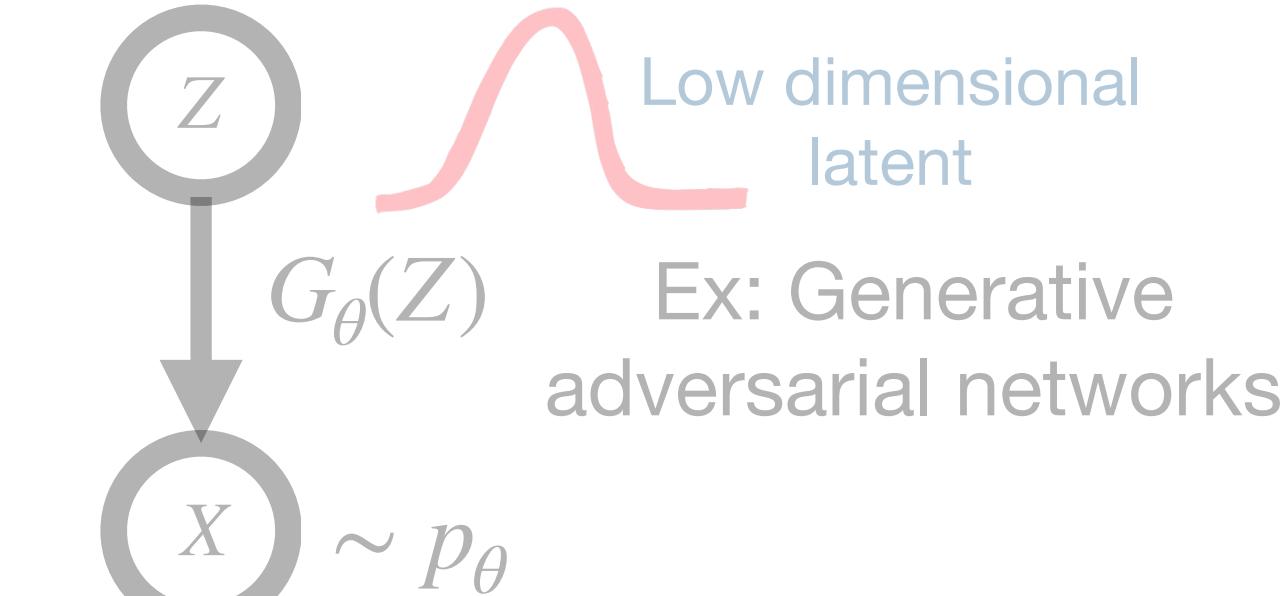
- ✓ Capture multimodality
- Hard to sample from



$$\min_{\theta} \mathcal{L}(p_{\theta})$$

$$\approx p_{\theta}$$

## Implicit Generative Models



- ✓ Capture low dimensionality
- Hard to optimize

Contributions in sampling

Contributions in generative modeling

Contributions in optimization

# Contributions in high dimensional sampling

Before my Ph.D.

**Limited reliability** of  
model-based sampling

During my Ph.D.

**Efficient** model-based samplers with theoretical guarantees

## High dimensional sampling\*

A fundamental problem

$$X \sim \pi$$

## Challenges

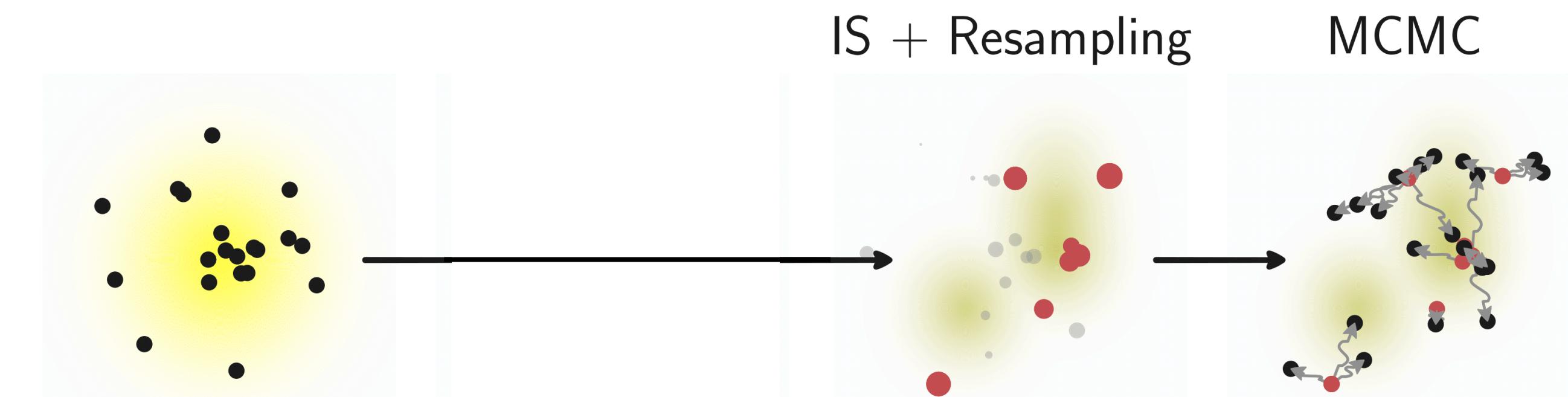
- Complex multi-modal distributions
- Curse of dimensions

## Scientific context

Sequential Monte-Carlo (SMC) methods

✓ Reliable + Strong theoretical guarantees

✗ Agnostic to the target density → High cost in general



**Contributions:** Principled and efficient model-based SMC samplers

1 publication at top-tier international conference (ICML 2021 Long Oral, top 3% submissions)

# Contributions in high dimensional sampling

Before my Ph.D.

**Limited reliability** of  
model-based sampling

During my Ph.D.

**Efficient** model-based samplers with theoretical guarantees

## High dimensional sampling\*

A fundamental problem

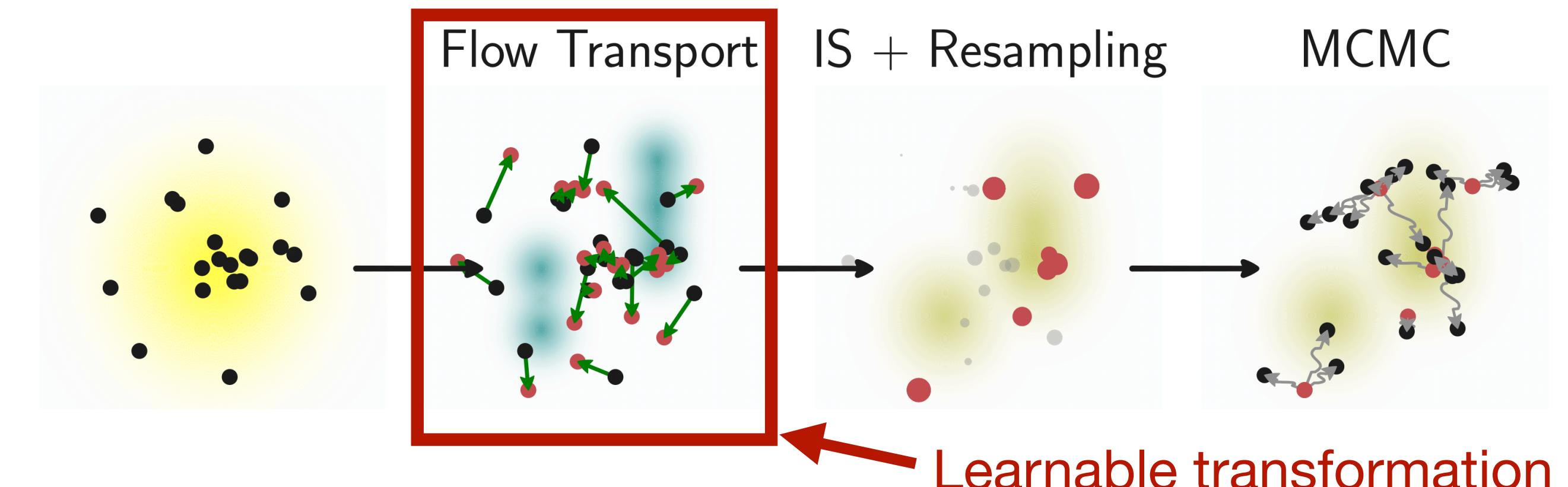
$$X \sim \pi$$

### Key idea

*Exploiting structure in the target density through  
modeling for more efficient sampling*

SMC with a learnable transformation

- ✓ Fast sampling exploiting shared structure
- ✓ Provably improves over than SMC
- ✓ Strong approximation guarantees



**Contributions:** Principled and efficient model-based SMC samplers

1 publication at top-tier international conference (ICML 2021 Long Oral, top 3% submissions)

# Learning structure from data with generative models

## Explicit Generative Models

$$p_{\theta}(X) = \frac{1}{Z_{\theta}} e^{-E_{\theta}(x)}$$

Ex: Energy-based models

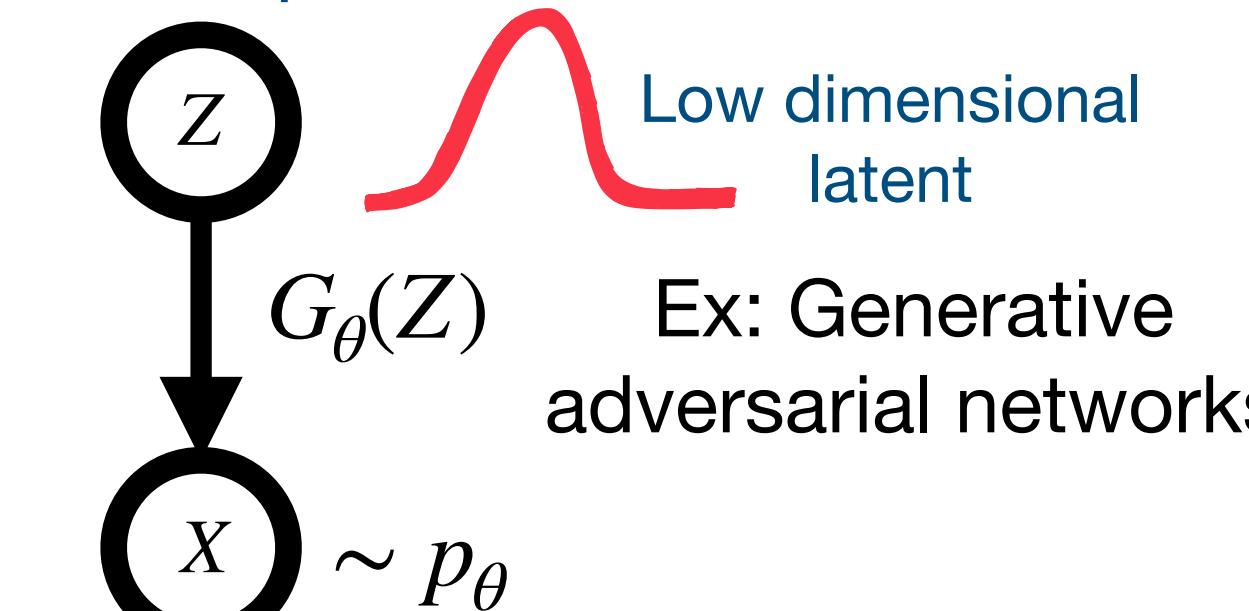
- ✓ Capture multimodality
- Hard to sample from



$$\min_{\theta} \mathcal{L}(p_{\theta})$$

$$\approx p_{\theta}$$

## Implicit Generative Models



- ✓ Capture low dimensionality
- Hard to optimize

Contributions in sampling

Contributions in generative modeling

Contributions in optimization

# Starting Research Position at Inria Thoth

## Bilevel Optimization for Machine Learning

# Contributions in bilevel optimization for machine learning

## Bilevel optimization

$$\begin{aligned} \min_{x \in \mathbb{R}^p} \mathcal{L}(x) &:= f(x, y^*(x)) \\ \text{s.t. } y^*(x) &= \arg \min_{y \in \mathbb{R}^d} g(x, y) \end{aligned}$$

## Applications in machine learning

- Learning hyper-parameters of a learning algorithm
  - Data-augmentation
  - Regularization
  - Architecture search

## Challenges

- Only approximations to  $y^*(x)$  are available
- Stochastic setting: additional errors due to noise

## Scientific context

### Gap between theory and practice

- Theory: - Costly convergent algorithms  
Practice: - Cheaper algorithms
  - Not necessarily convergent

**Contributions:** First optimal convergence results for bilevel problems

1 publication at top-tier international conference (ICLR, 2022)

# Contributions in bilevel optimization for machine learning

## Bilevel optimization

$$\begin{aligned} \min_{x \in \mathbb{R}^p} \mathcal{L}(x) &:= f(x, y^*(x)) \\ \text{s.t. } y^*(x) &= \arg \min_{y \in \mathbb{R}^d} g(x, y) \end{aligned}$$

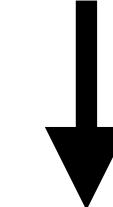
## Applications in machine learning

- Learning hyper-parameters of a learning algorithm
  - Data-augmentation
  - Regularization
  - Architecture search

## Key idea

Introducing tools from dynamical systems:

*Singularly perturbed systems analysis\**



A tight analysis of simple and efficient algorithms

- ✓ The first analysis providing optimal convergence rates
- ✓ Significant reduction of the computational cost

**Contributions:** First optimal convergence results for bilevel problems

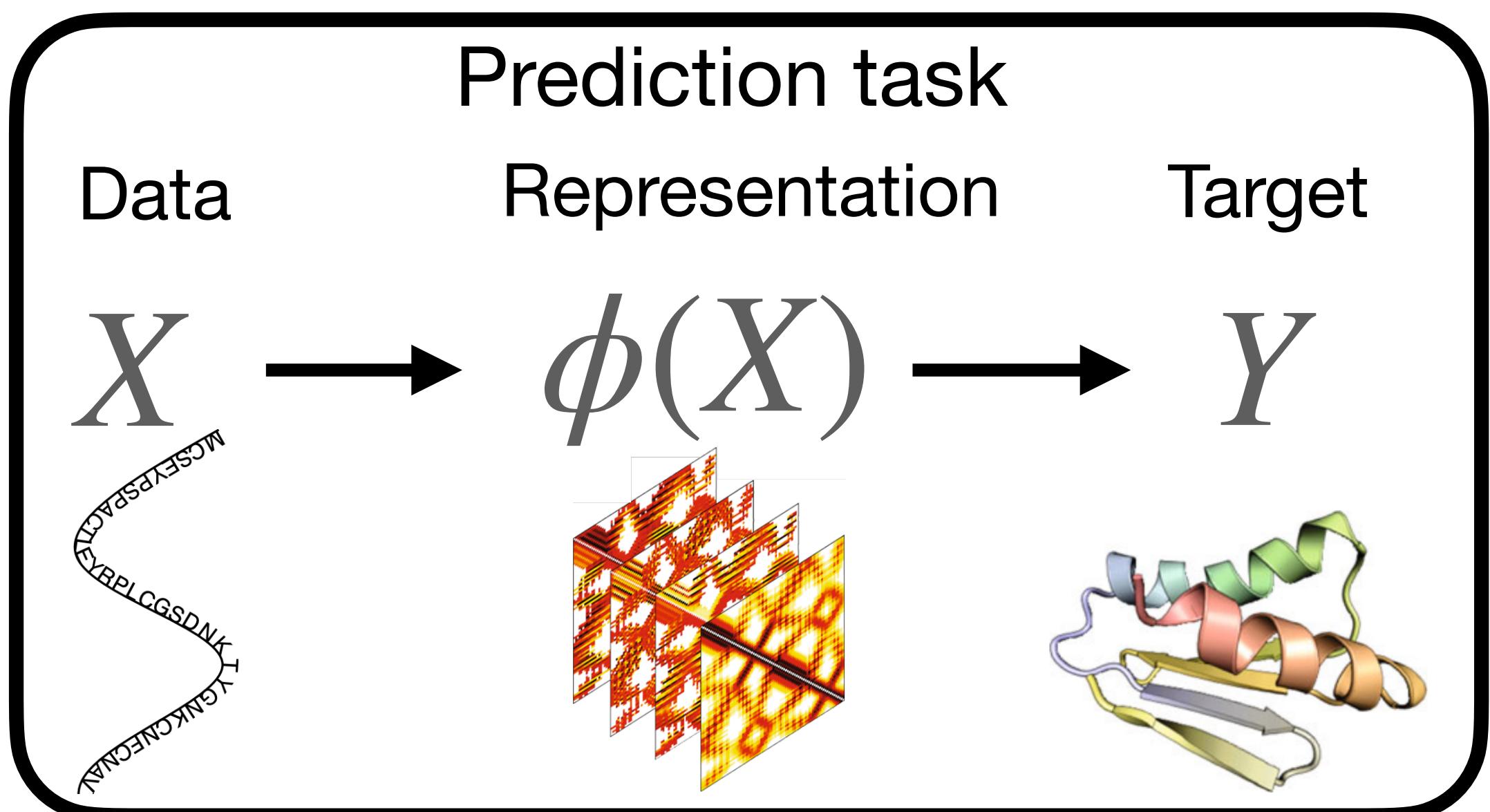
1 publication at top-tier international conference (ICLR, 2022)

\*Stabilité Asymptotique Pour des Problèmes de Perturbations Singulières, Habets, 1974

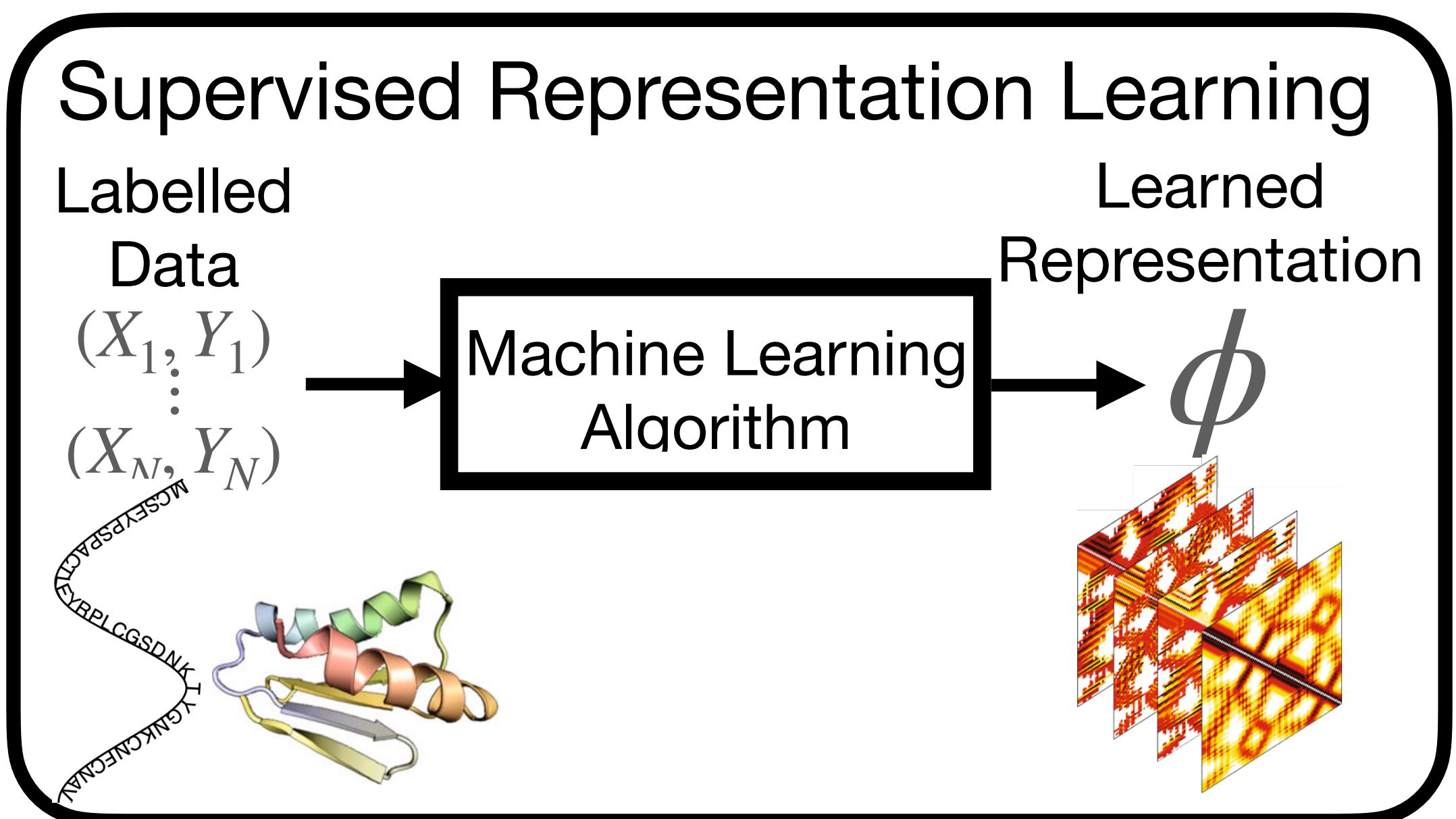
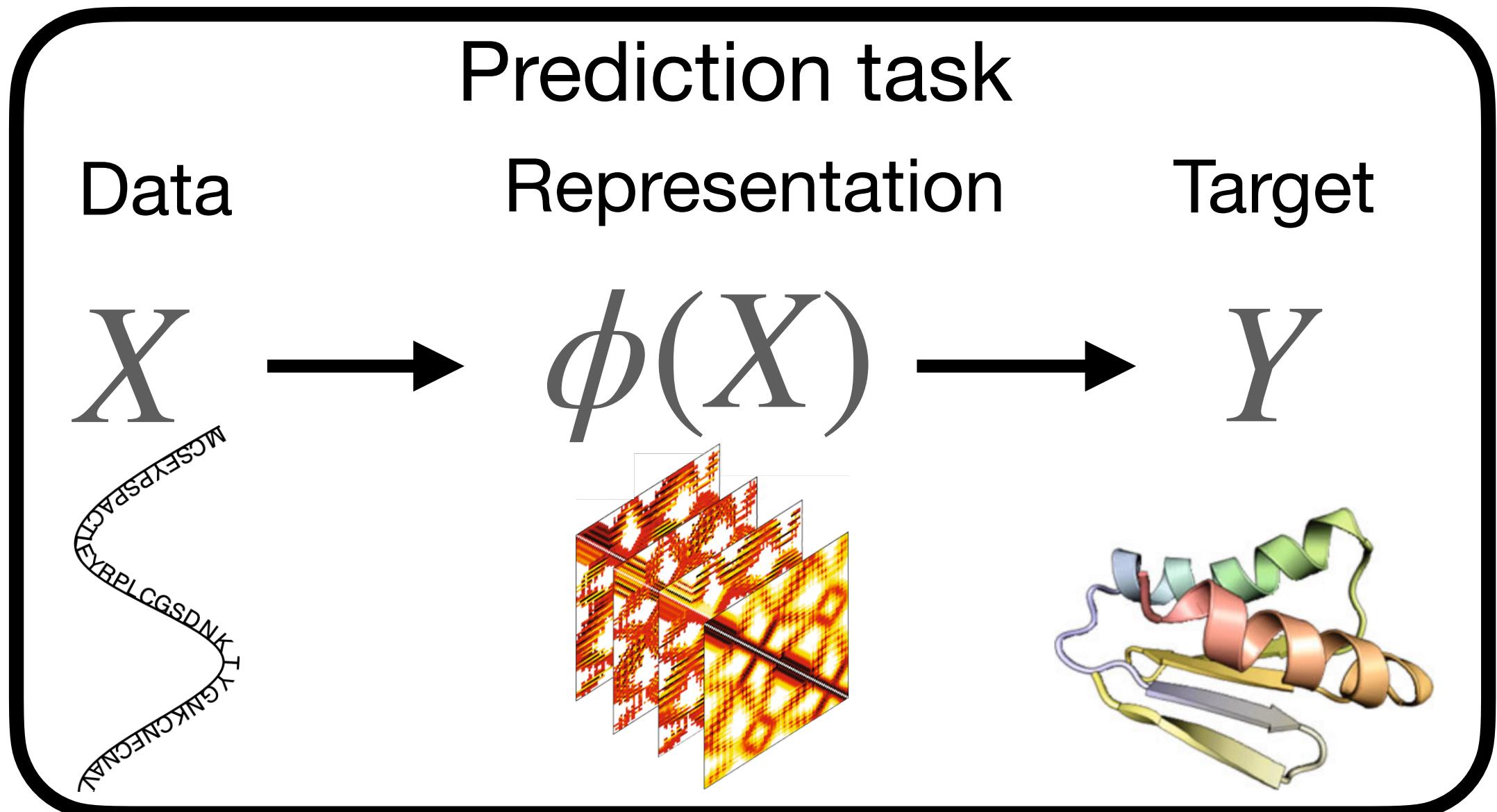
# **Addressing the lack of reliability in Machine learning**

## A novel probabilistic viewpoint on self-supervised data representations

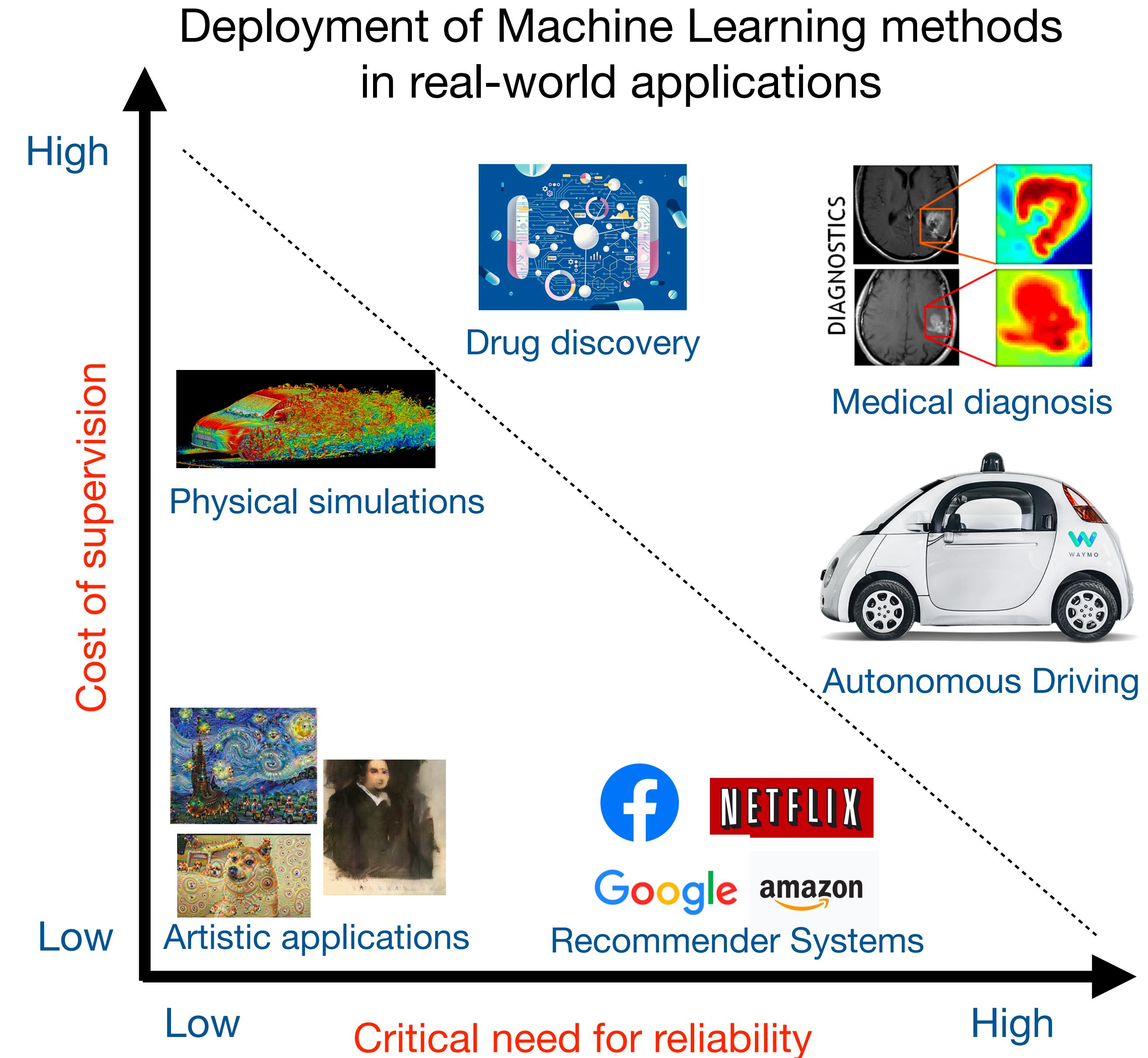
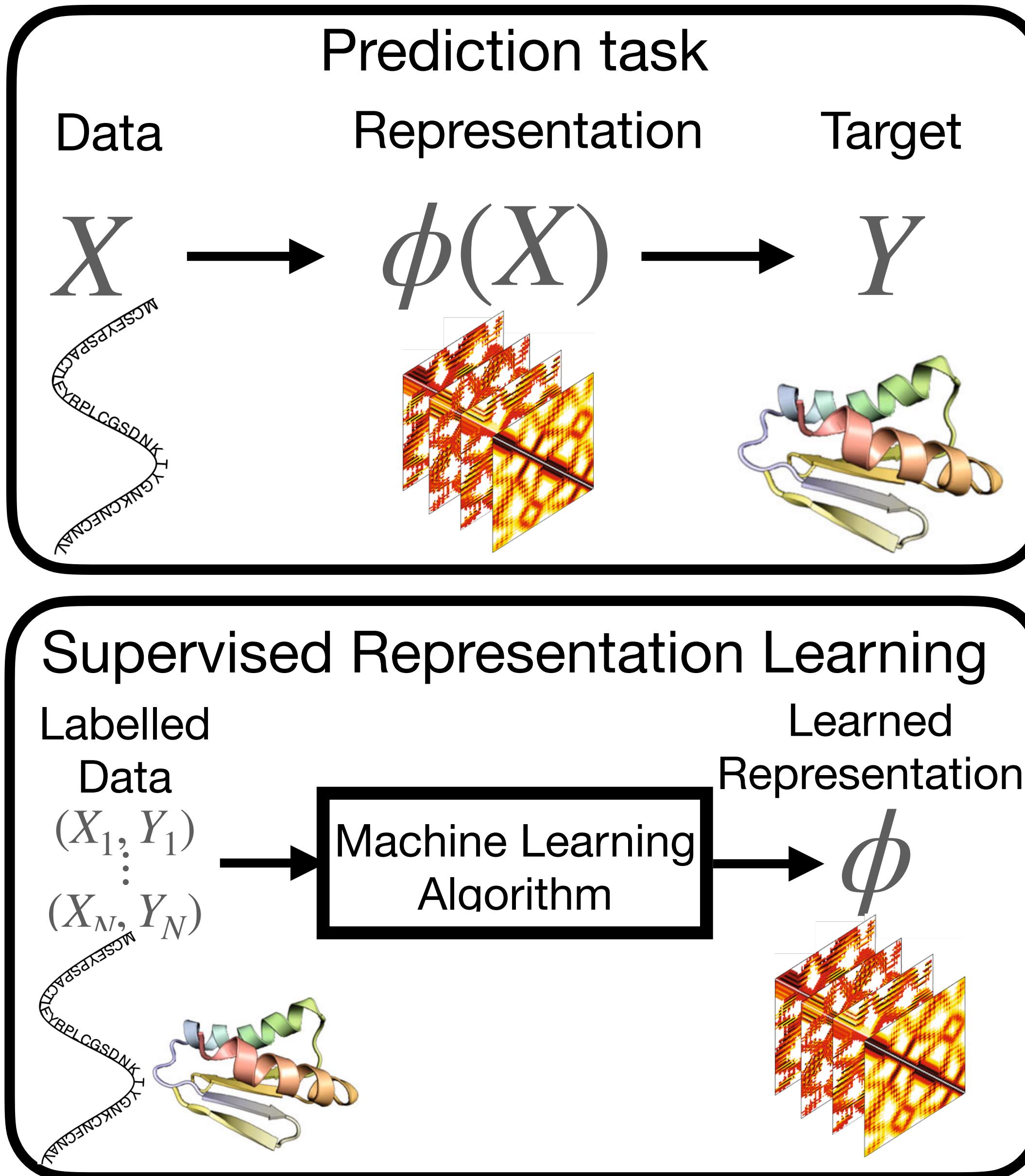
# Key Observation



# Key Observation

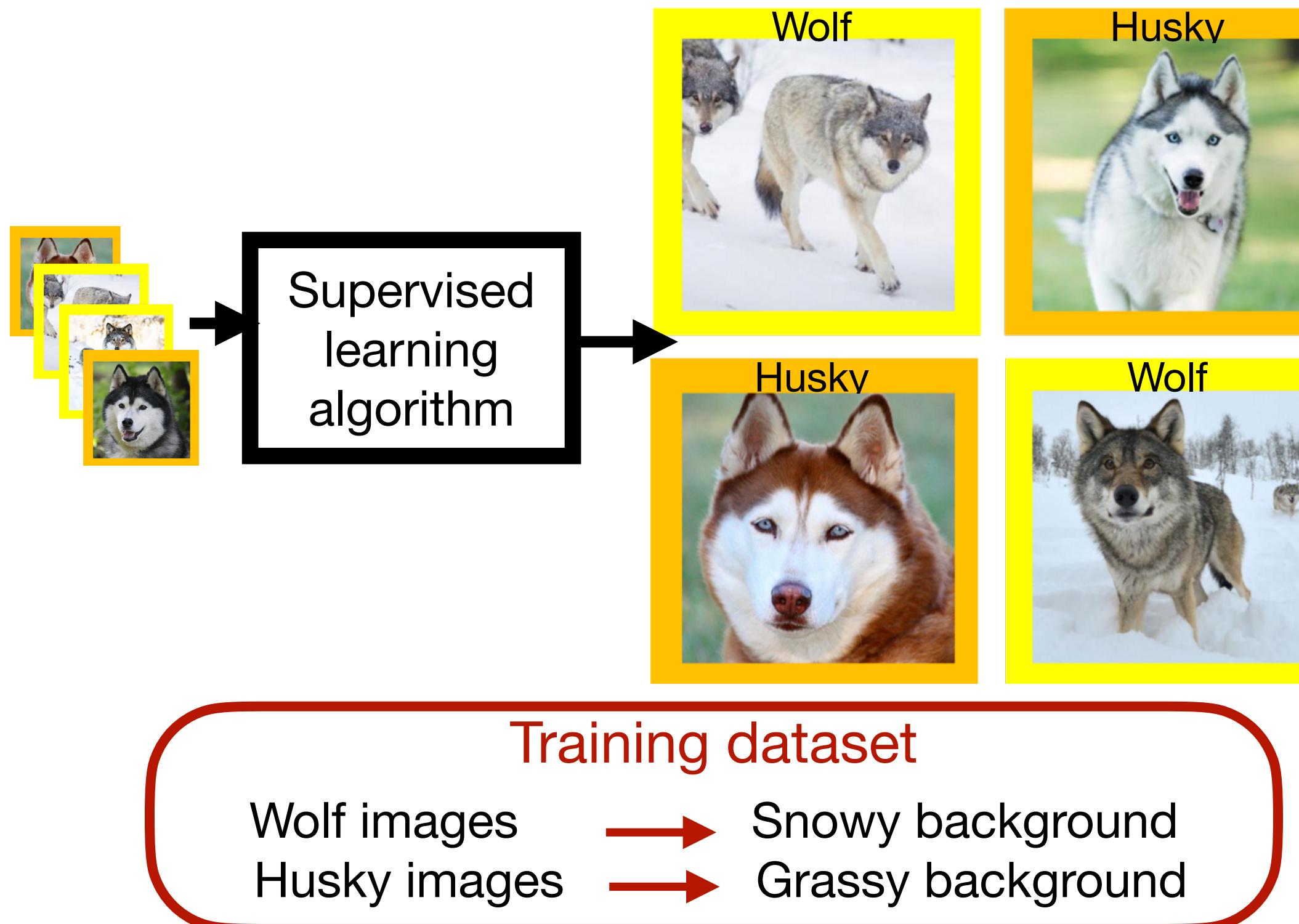


# Key Observation



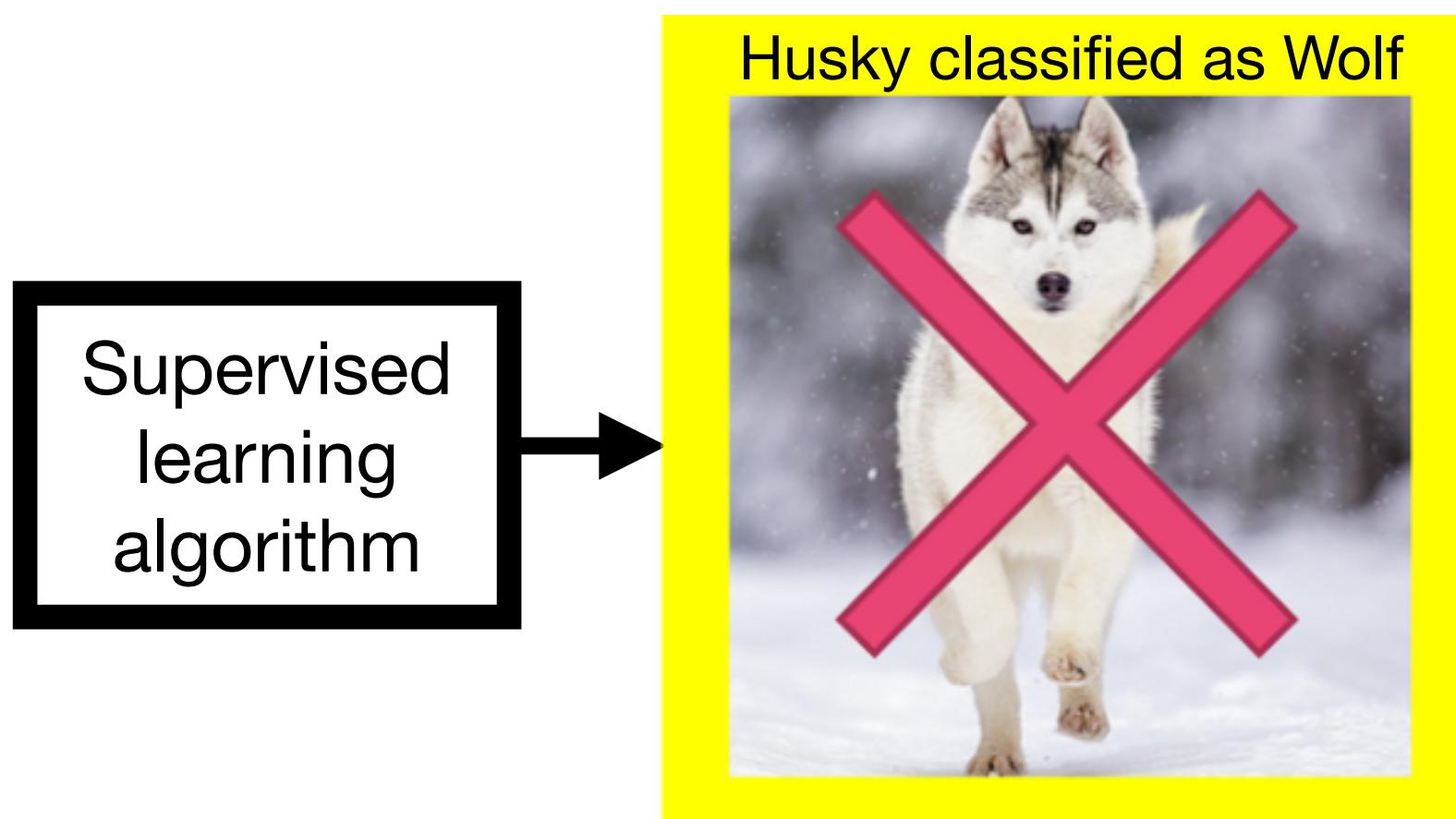
# Some sources of unreliability in Machine Learning

## Lack of robustness to distributional shift



# Some sources of unreliability in Machine Learning

## Lack of robustness to distributional shift



### Training dataset

Wolf images → Snowy background  
Husky images → Grassy background

Systems learn spurious correlations  
Do not generalize to out of distribution shifts

Need to develop robust representations

# Some sources of unreliability in Machine Learning

## Lack of robustness to distributional shift

Supervised learning algorithm



### Training dataset

Wolf images → Snowy background  
Husky images → Grassy background

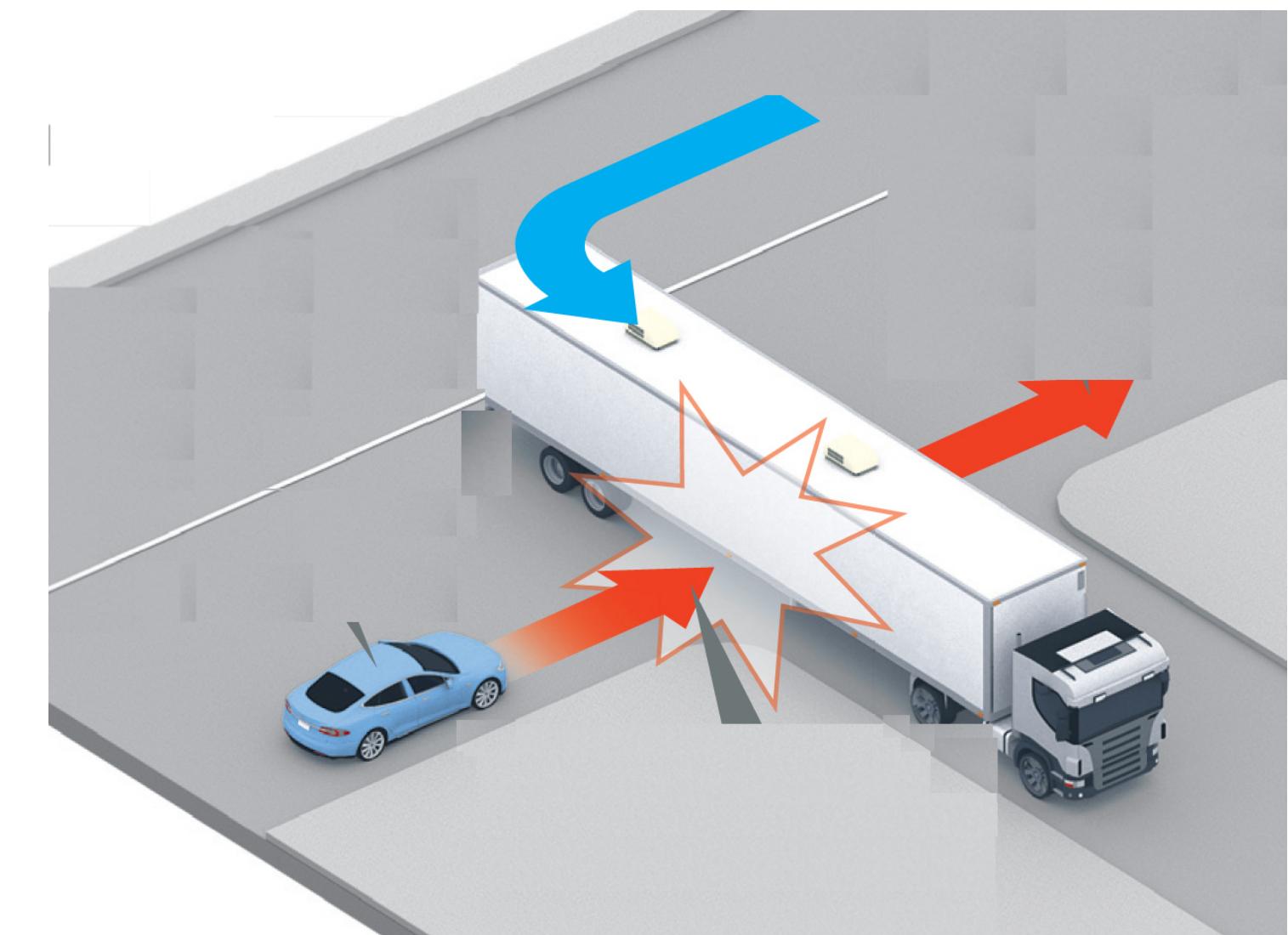
Systems learn spurious correlations  
Do not generalize to out of distribution shifts

Need to develop robust representations

## Lack of uncertainty quantification

Autopilot confuses a truck with a bright sky

↓  
Car continues driving at full speed



Overly-confident systems  
Do not account for alternative outcomes

Need to account for uncertainty in the predictions

# Representation learning



$X$

$\phi(X)$

$Y$

Toxicity

## Self-supervised representation learning

Robustness to spurious correlations\*  
by encoding knowledge relevant for prediction

Prior knowledge

Relevant for the task: Color, Texture, Shape,...

Unlabelled  
data



$X_1$   
 $\vdots$   
 $X_N$

**Prior work**

- DeepCluster
- SwAV
- DINO

Learned  
Representation

$\phi$

Does not account for uncertainty

# Representation learning



$X$

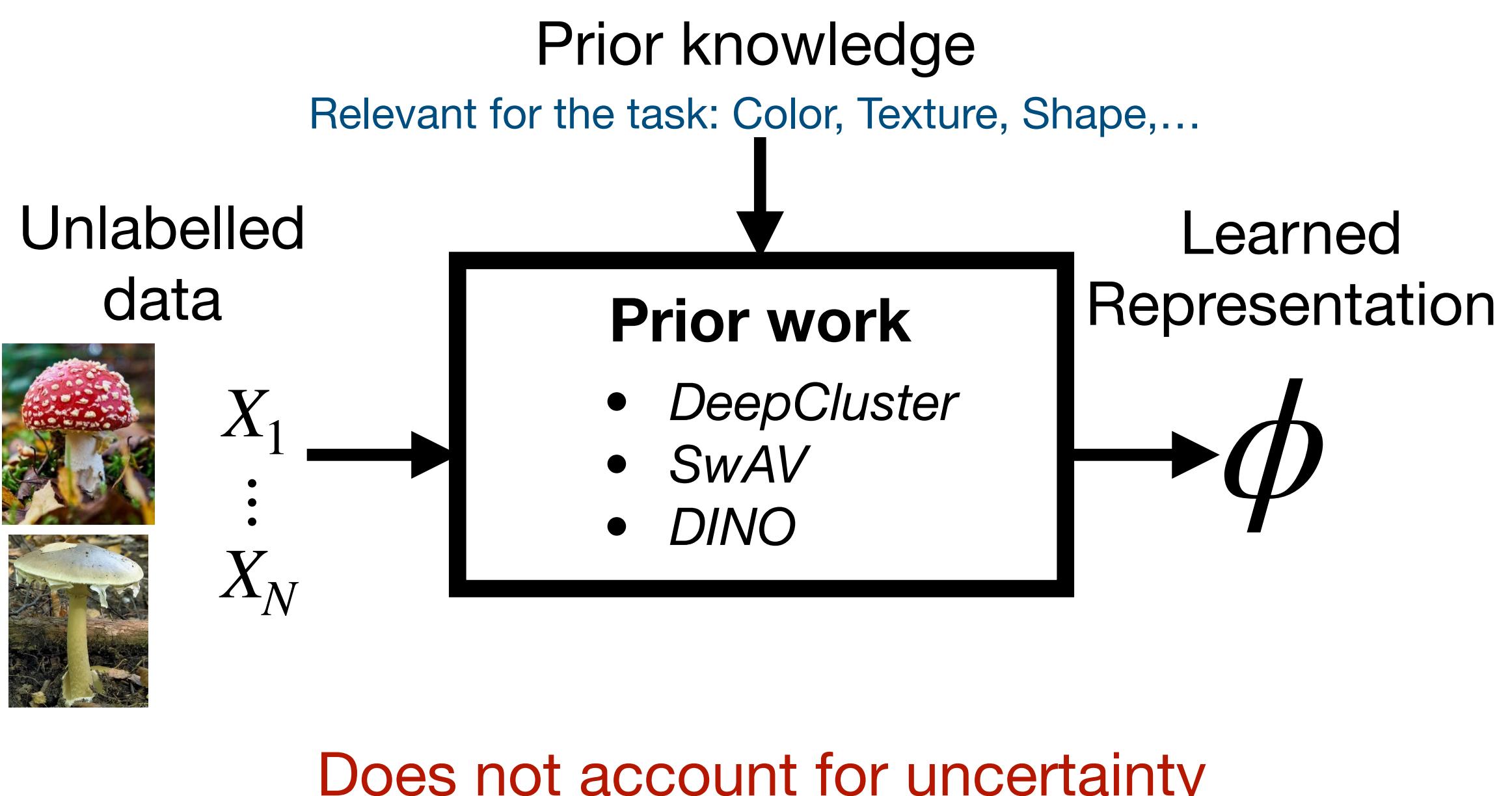
$\phi(X)$

$Y$

Toxicity

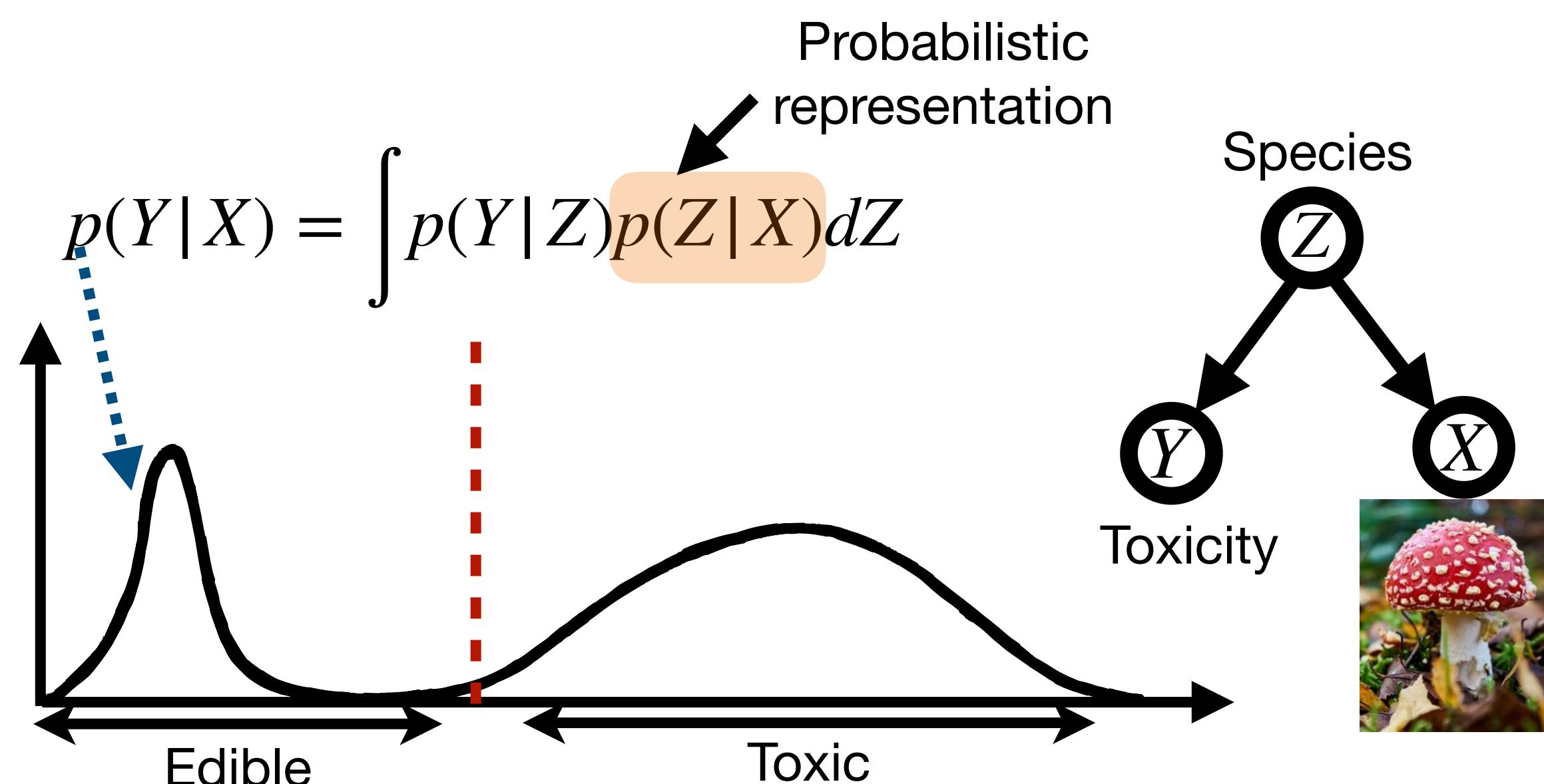
## Self-supervised representation learning

Robustness to spurious correlations\*  
by encoding knowledge relevant for prediction



## Probabilistic unsupervised representation learning

Accounting for uncertainty in the predictions



# Representation learning



$X$

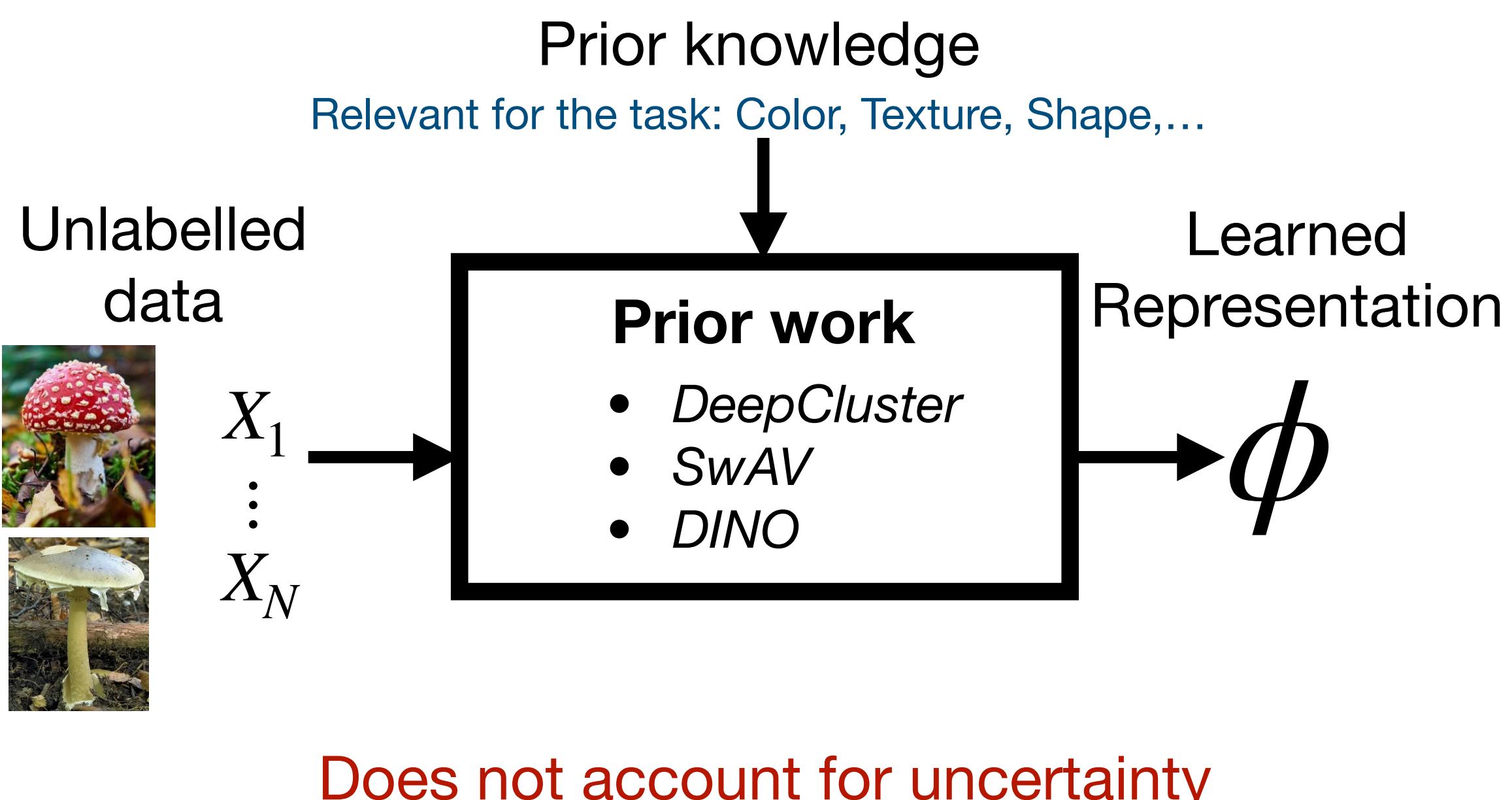
$\phi(X)$

$Y$

Toxicity

## Self-supervised representation learning

Robustness to spurious correlations  
by encoding knowledge relevant for prediction



## Probabilistic unsupervised representation learning

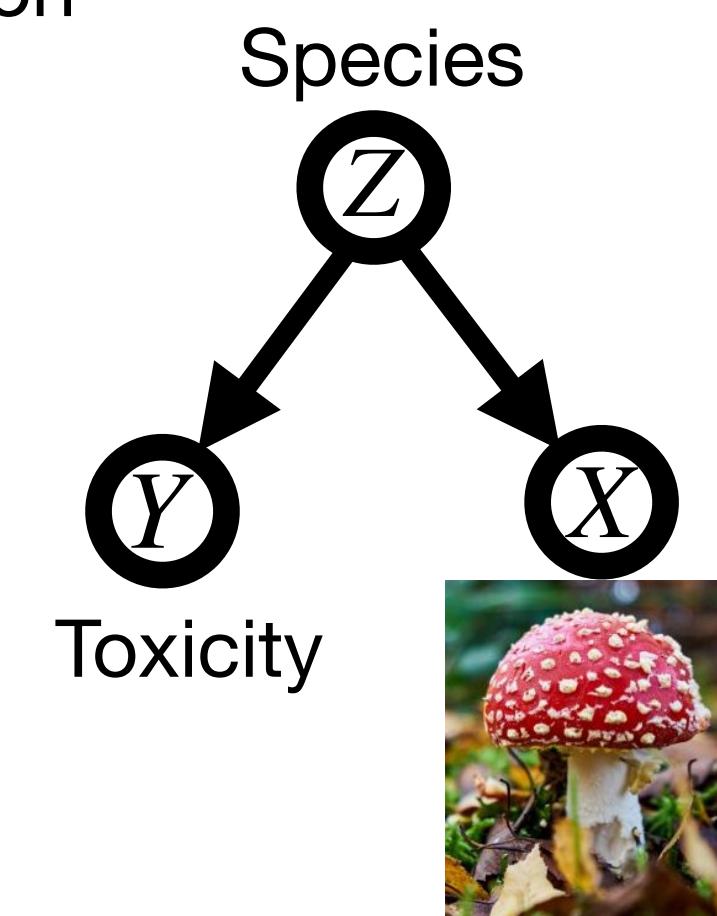
Accounting for uncertainty in the predictions

$$p(Y|X) = \int p(Y|Z)p(Z|X)dZ$$

Probabilistic representation

Prior work

- *Variational Auto-encoders*
- *Deep Bayesian networks*
- *Deep Gaussian processes*
- *Non-linear ICA*



Unsupervised → may not be relevant for prediction\*

## Self-Supervised representation Learning

Encodes knowledge relevant for prediction

Does not account for uncertainty

- Prior work**
- *DeepCluster*
  - *SwAV*
  - *DINO*

## Probabilistic unsupervised representation learning

Accounts for uncertainty

Not always relevant for prediction

### Prior work

- *Variational Auto-encoders*
- *Deep Bayesian networks*
- *Deep Gaussian processes*
- *Non-linear ICA*

# Representation learning

## Self-Supervised representation Learning

Encodes knowledge relevant for prediction

Does not account for uncertainty

### Prior work

- DeepCluster
- SwAV
- DINO

## Probabilistic unsupervised representation learning

Accounts for uncertainty

Not always relevant for prediction

### Prior work

- Variational Auto-encoders
- Deep Bayesian networks
- Deep Gaussian processes
- Non-linear ICA

## My vision

Learning probabilistic representations that reflect the data generating mechanism

## Probabilistic self-supervised representation learning

Axes 1 & 2

### Theory and method

A new probabilistic self-supervised learning framework

✓ Accounts for uncertainty

✓ Robust to distributional shift

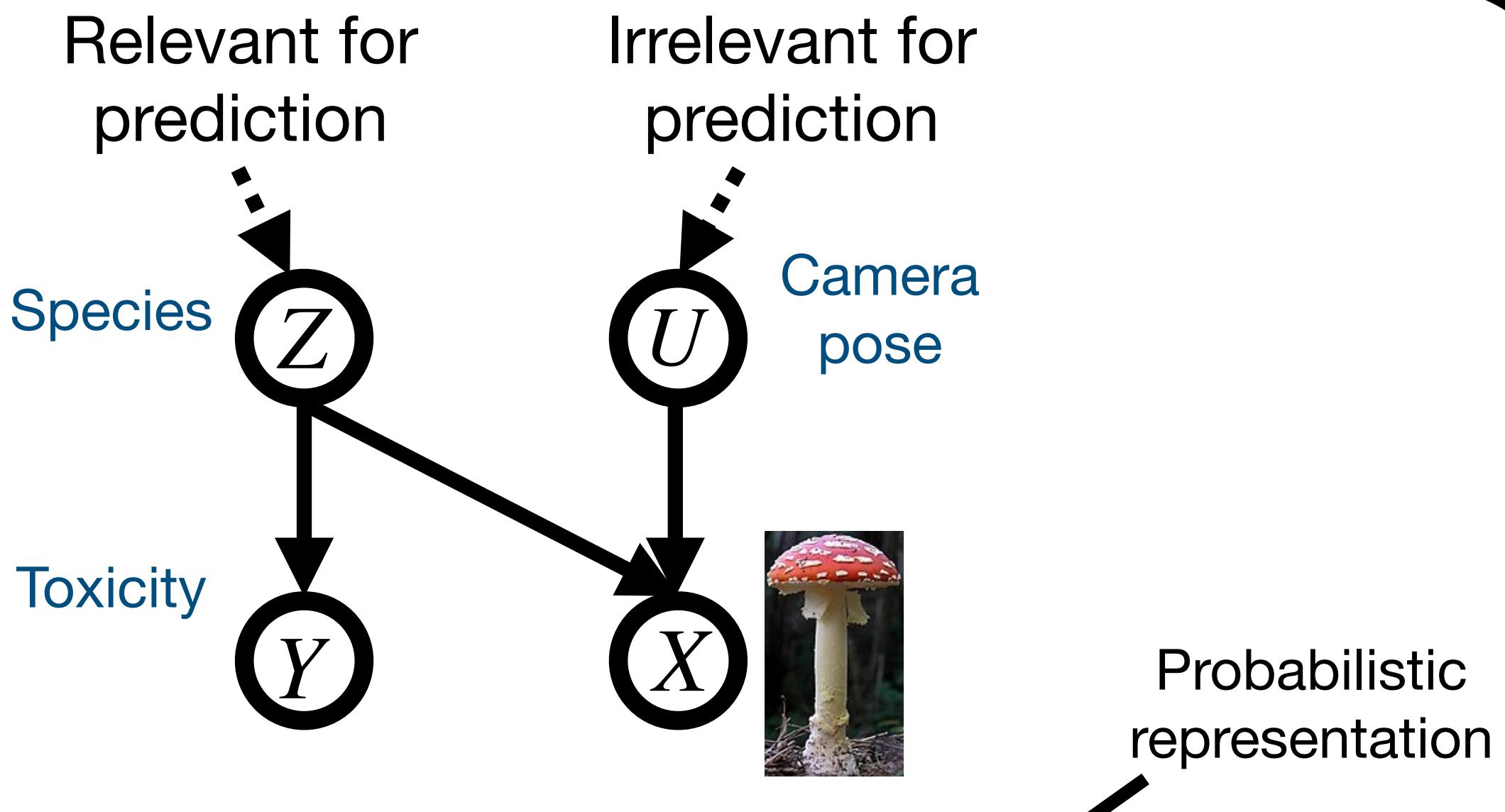
Efficient and reliable machine learning

Axes 3 & 4

### Tools

Efficient and reliable sampling and optimization

# Axis 1: A new probabilistic self-supervised learning framework



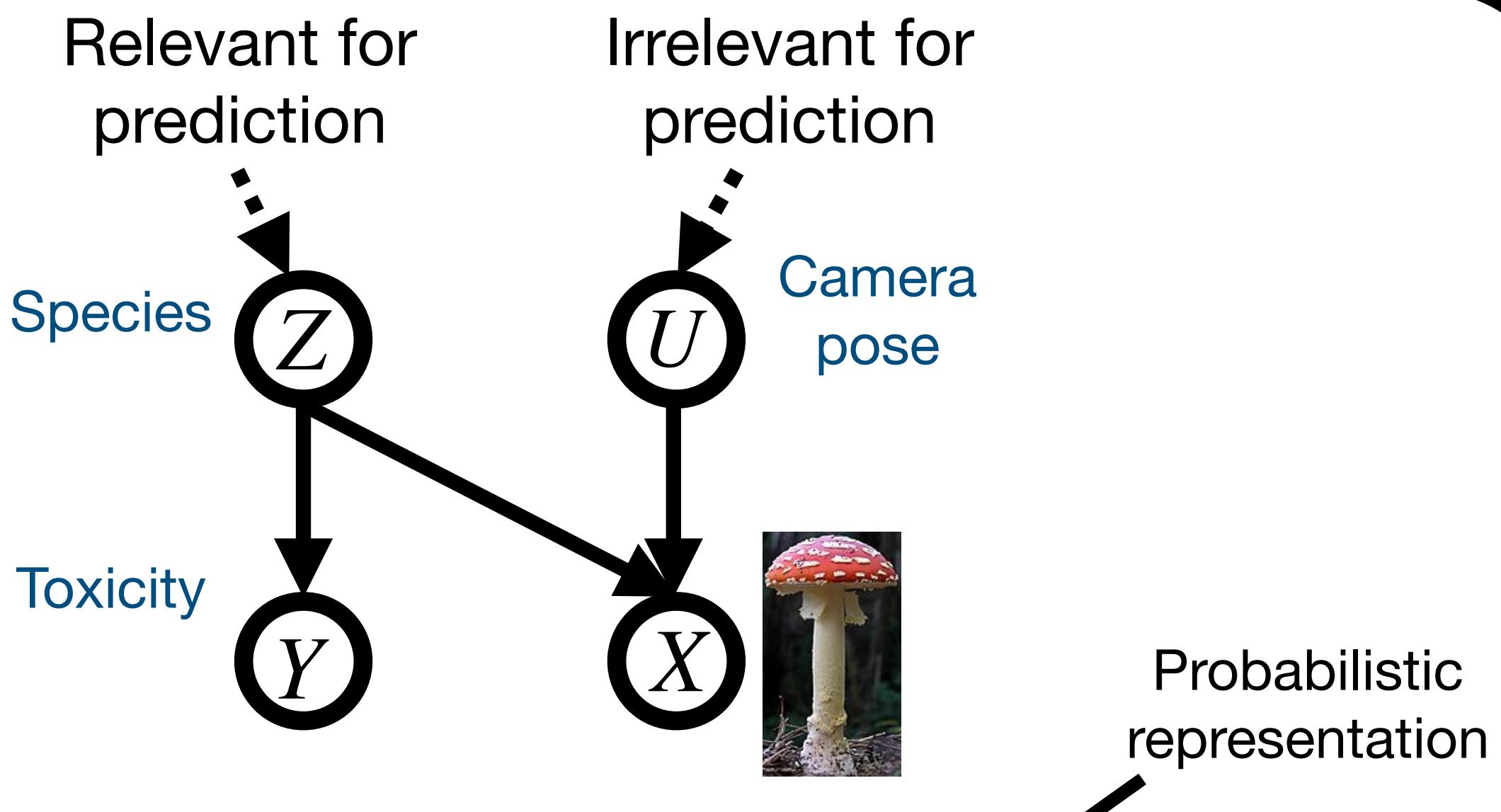
## Challenges

- ❖ Robustness to distributional shift
- ❖ Identifying the relevant latent variable

## Impact

- ❖ Better robustness to distribution shift
- ❖ Modeling uncertainty
- ❖ Provides mechanisms for basic reasoning

# Axis 1: A new probabilistic self-supervised learning framework



## Challenges

- ❖ Robustness to distributional shift
- ❖ Identifying the relevant latent variable

## Impact

- ❖ Better robustness to distribution shift
- ❖ Modeling uncertainty
- ❖ Provides mechanisms for basic reasoning

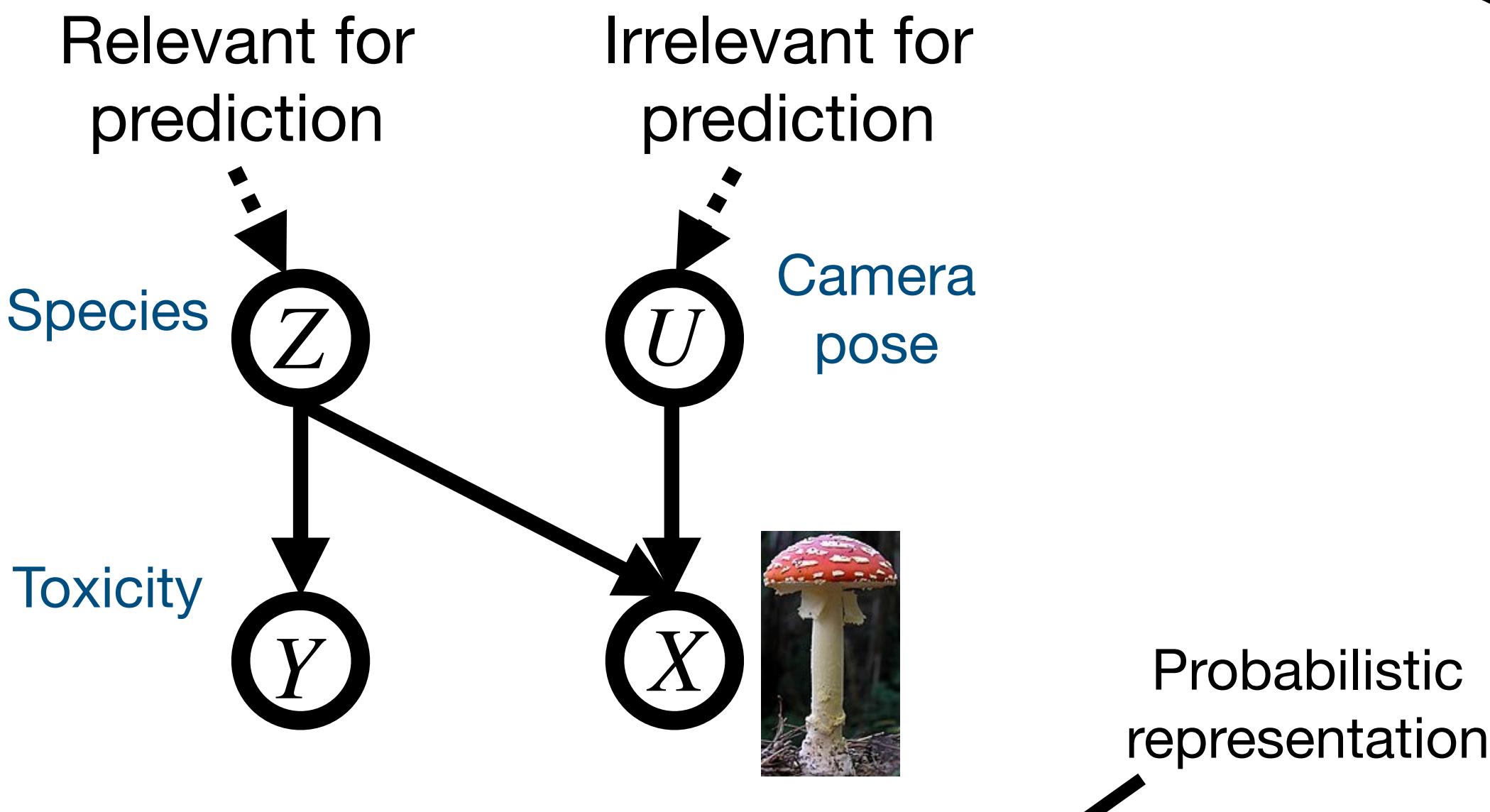
## Self-supervised learning of $p(Z|X)$

### Models robust to distributional shift

Preserving the structure of the generating mechanism\*

$$p(Z|X) \approx p_\theta(X|Z)p_\theta(Z)$$

# Axis 1: A new probabilistic self-supervised learning framework



## Challenges

- ❖ Robustness to distributional shift
- ❖ Identifying the relevant latent variable

## Impact

- ❖ Better robustness to distribution shift
- ❖ Modeling uncertainty
- ❖ Provides mechanisms for basic reasoning

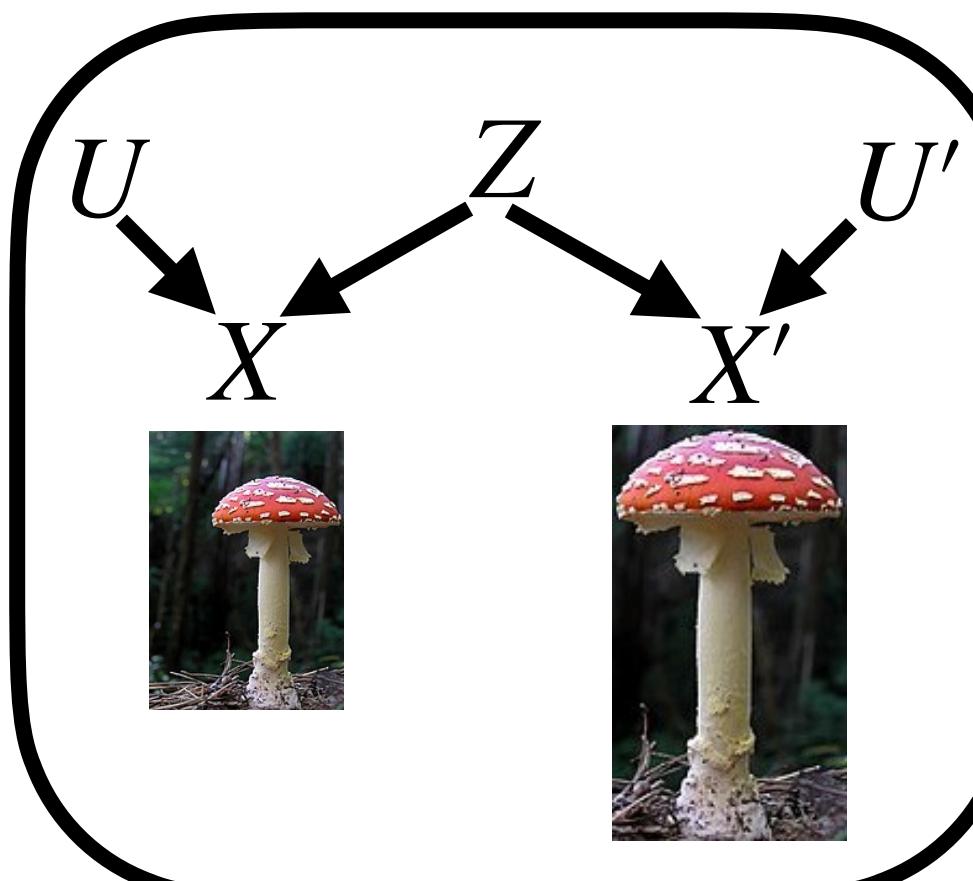
## Self-supervised learning of $p(Z|X)$

### Models robust to distributional shift

Preserving the structure of the generating mechanism\*

$$p(Z|X) \approx p_\theta(X|Z)p_\theta(Z)$$

### Identifying the relevant latent variable

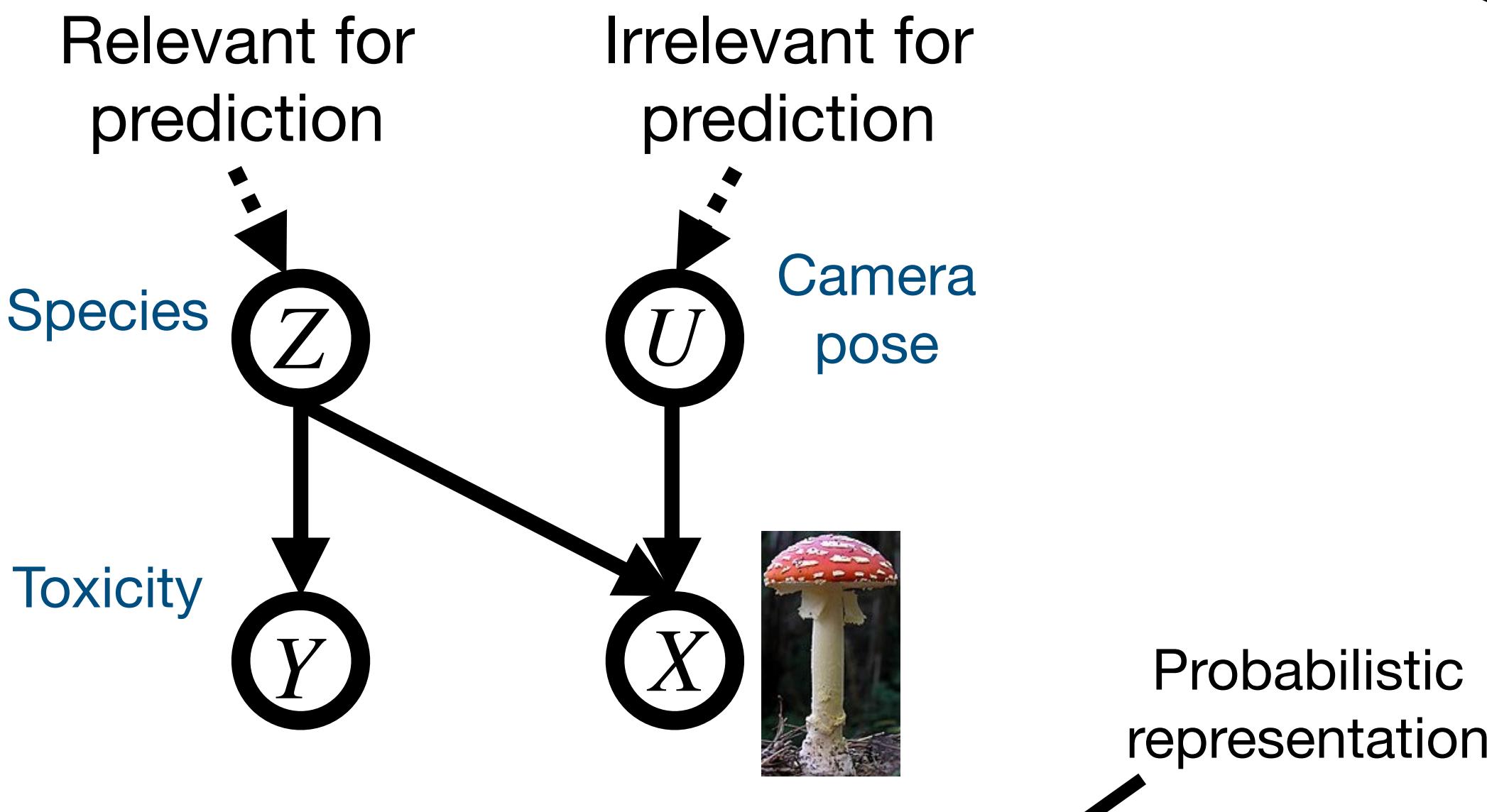


### Examples:

- Data augmentation
- Ordering in time series
- Interventions on causal variables

\* On causal and anticausal learning, Schölkopf et al. , 2012

# Axis 1: A new probabilistic self-supervised learning framework



$$p(Y|X) \propto \int p(Y|Z)p(Z|X)dZ$$

## Challenges

- ❖ Robustness to distributional shift
- ❖ Identifying the relevant latent variable

## Impact

- ❖ Better robustness to distribution shift
- ❖ Modeling uncertainty
- ❖ Provides mechanisms for basic reasoning

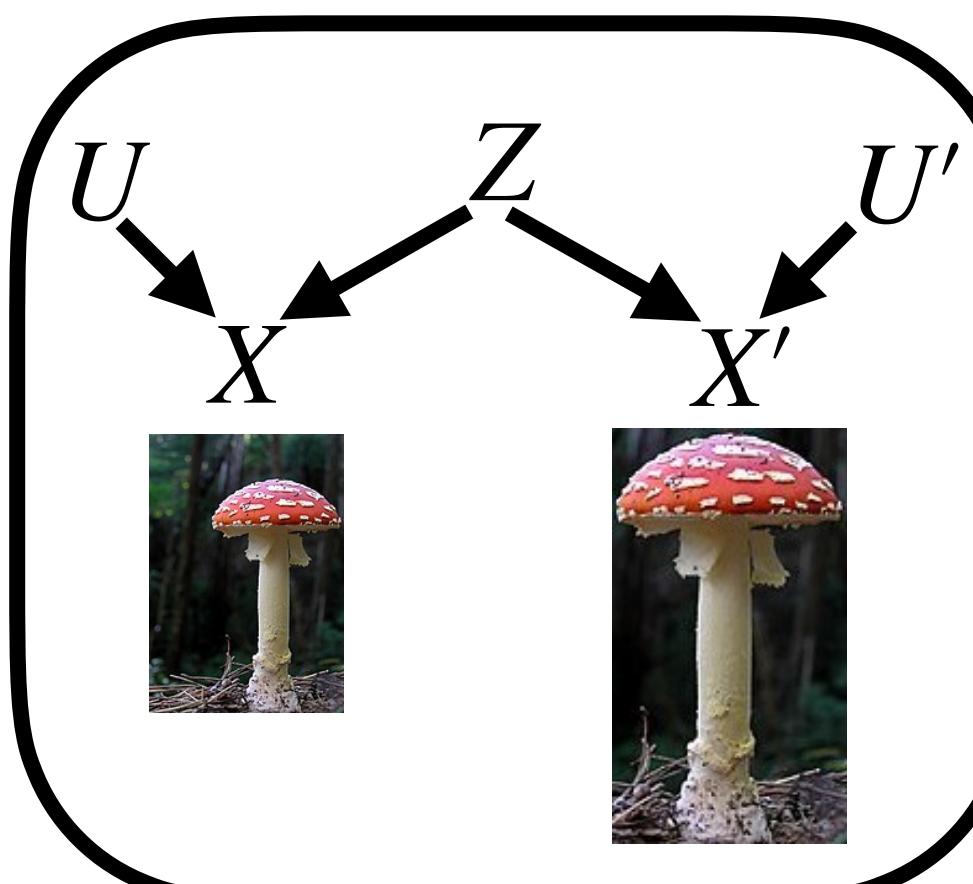
## Self-supervised learning of $p(Z|X)$

### Models robust to distributional shift

Preserving the structure of the generating mechanism\*

$$p(Z|X) \approx p_\theta(X|Z)p_\theta(Z)$$

### Identifying the relevant latent variable



### Examples:

- Data augmentation
- Ordering in time series
- Interventions on causal variables

### Beyond maximum likelihood:

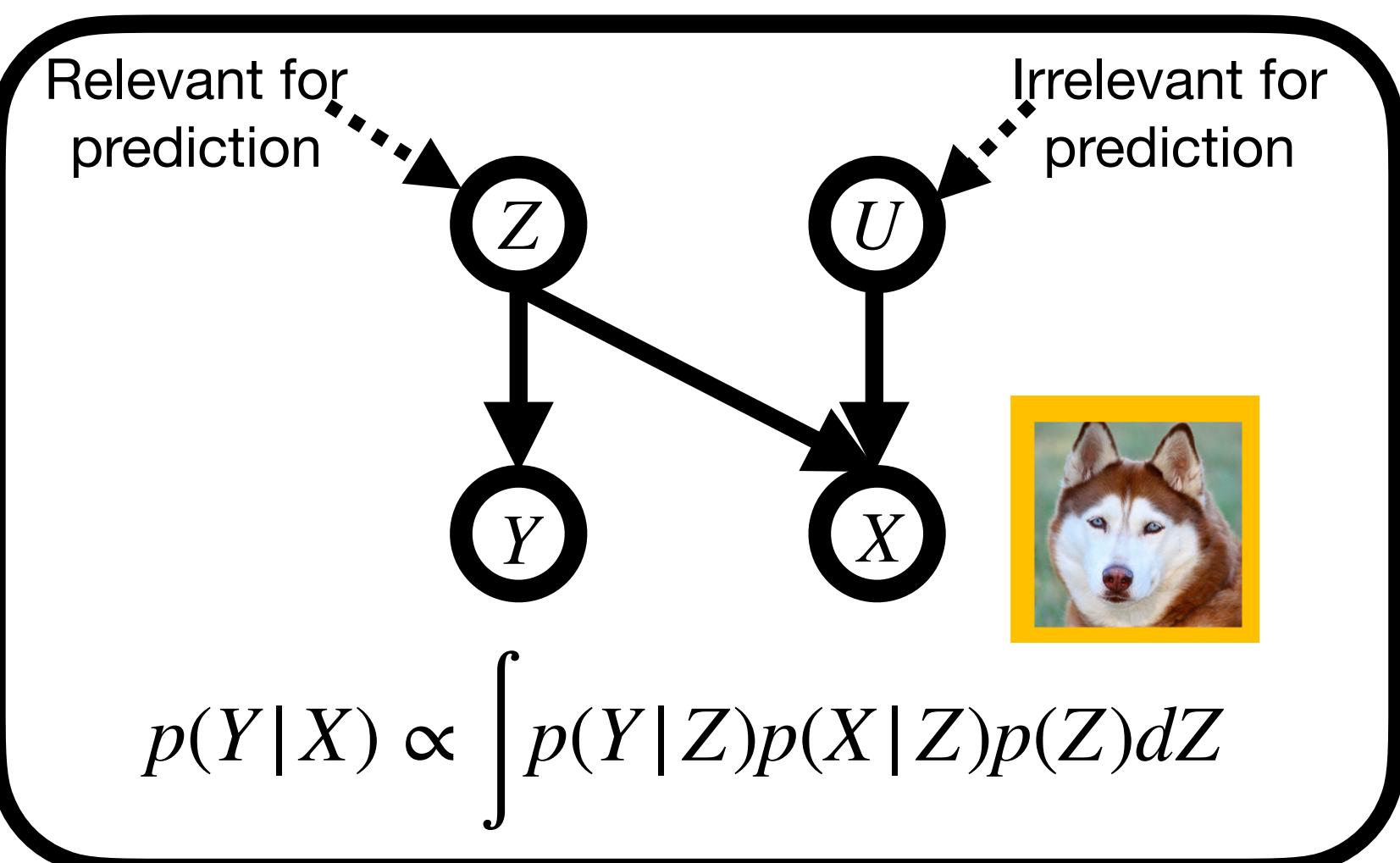
Learning to associate correlated pairs  $(X, X')$

$$p(X, X') \approx \int p_\theta(X|Z)p_\theta(X'|Z)p_\theta(Z)dZ$$

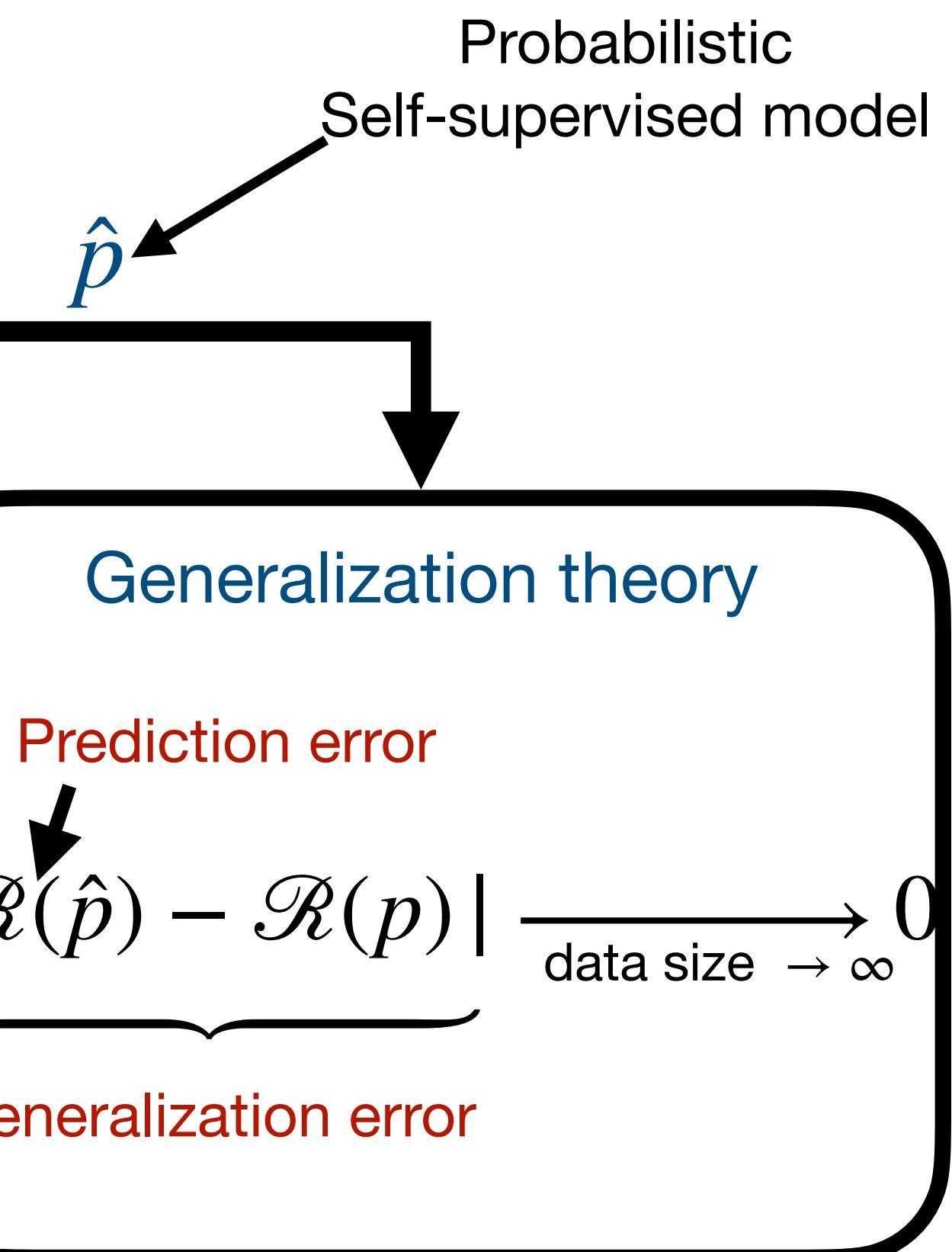
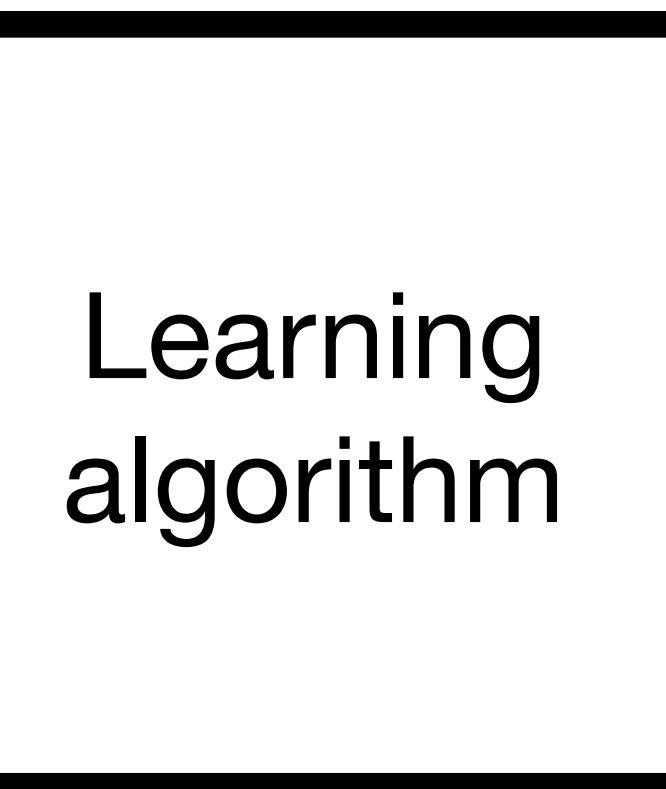
\* On causal and anticausal learning, Schölkopf et al. , 2012

# Axis 2: A learning theory for probabilistic self-supervised learning (PSSL)

## Data generating mechanism

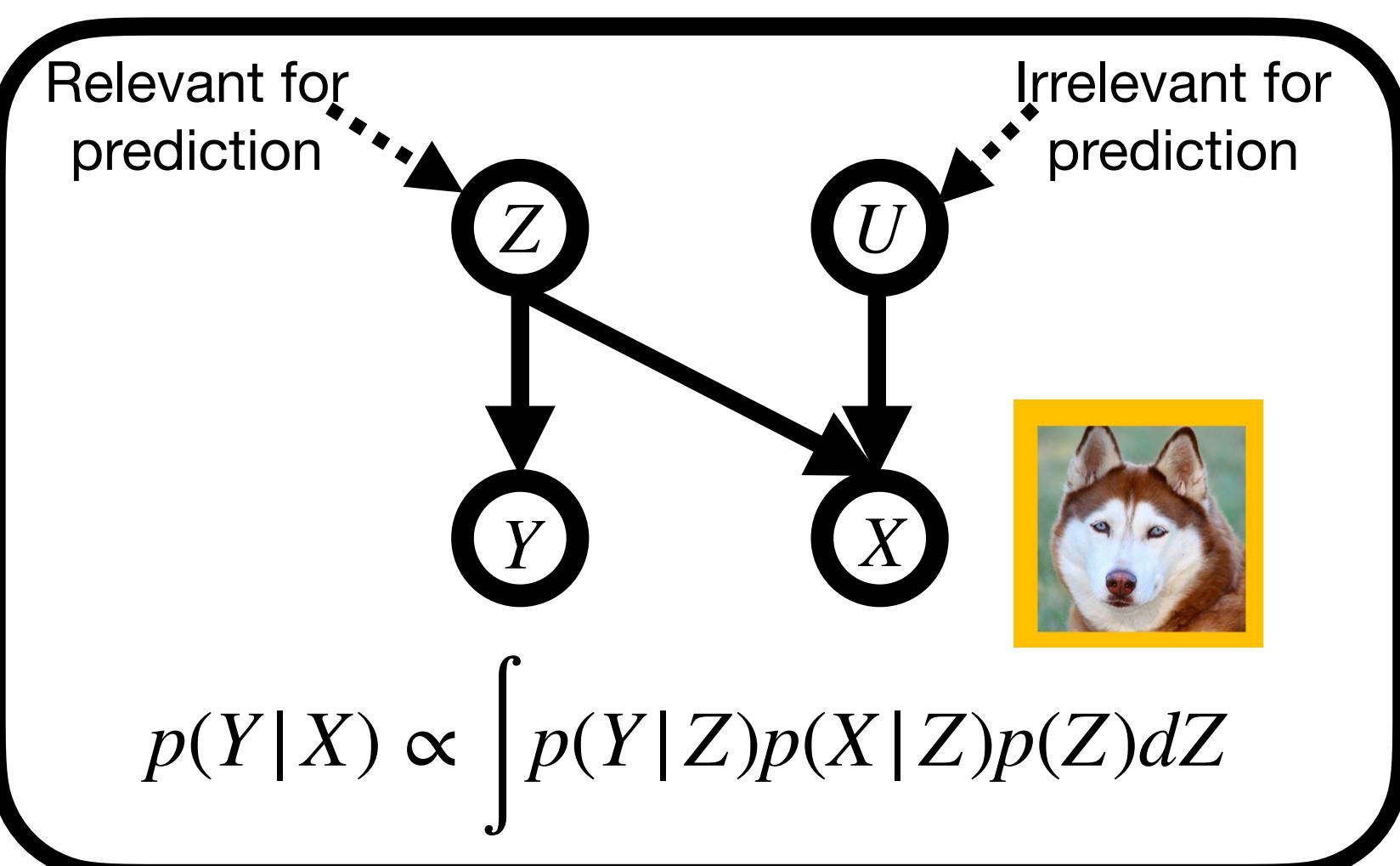


Data

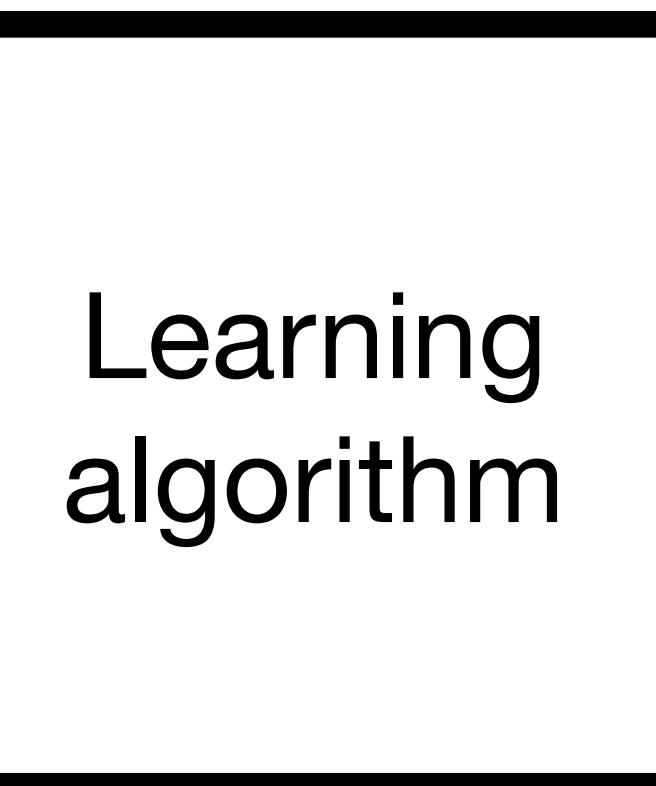


# Axis 2: A learning theory for probabilistic self-supervised learning (PSSL)

## Data generating mechanism



Data



Probabilistic  
Self-supervised model

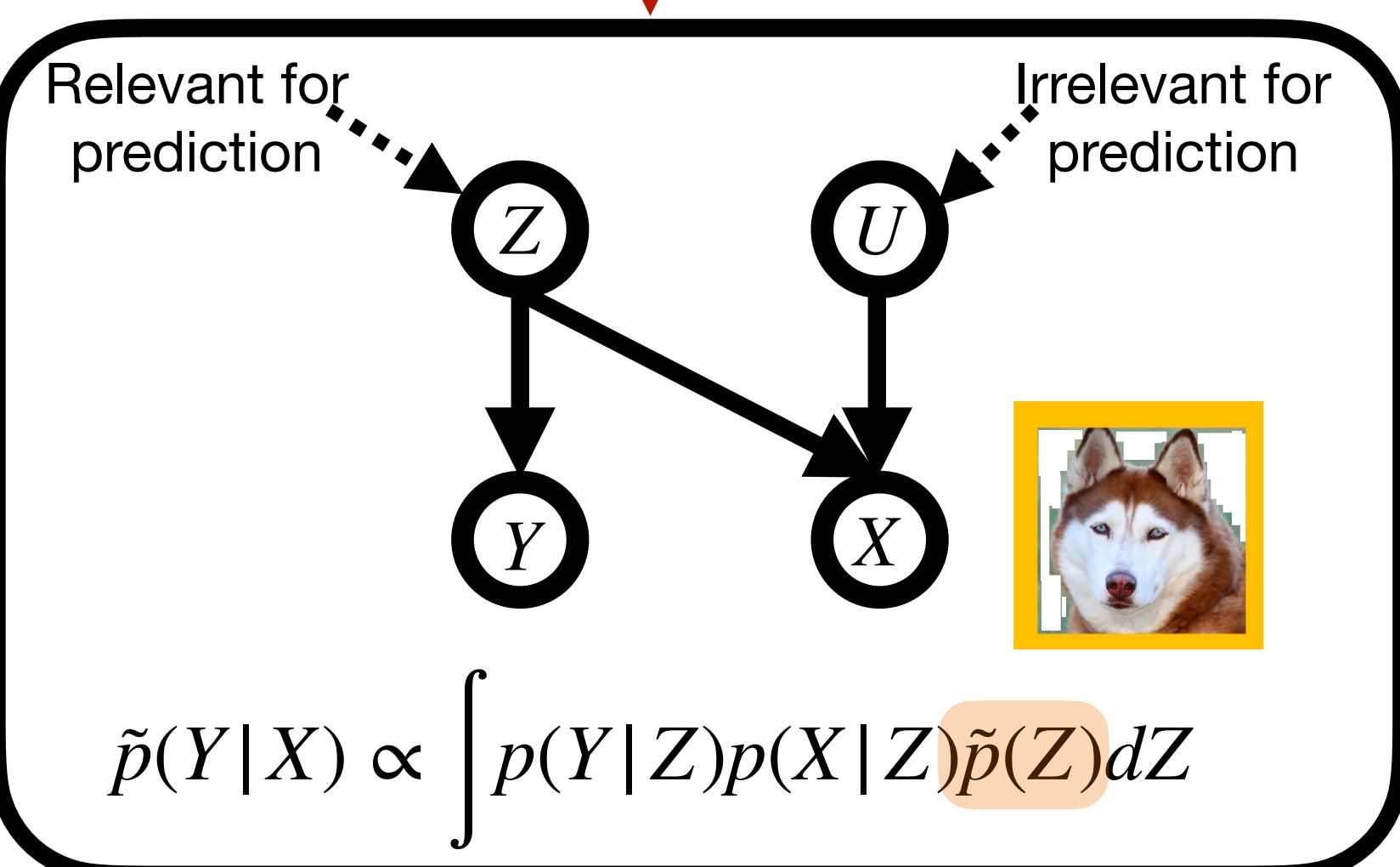
$\hat{p}$

Generalization theory

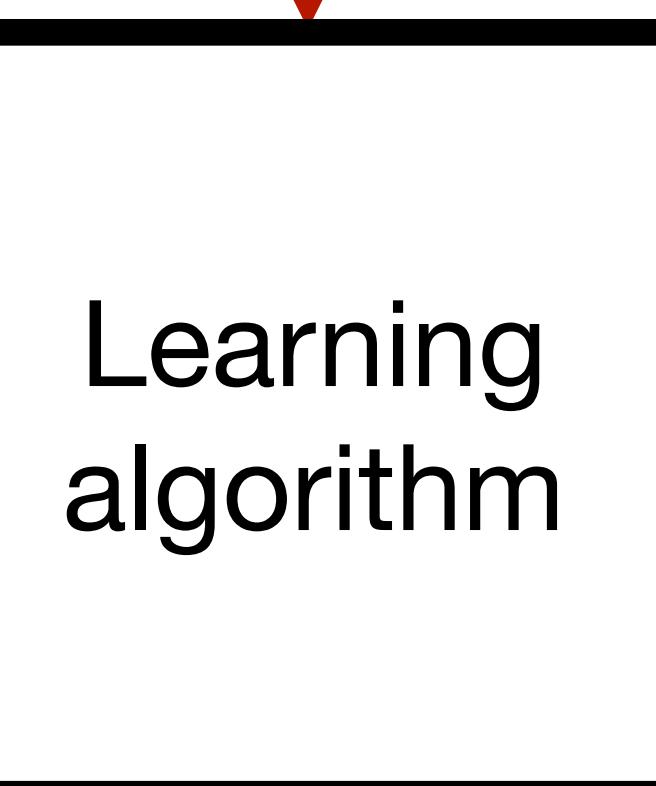
Prediction error

$$\underbrace{|\mathcal{R}(\hat{p}) - \mathcal{R}(p)|}_{\text{Generalization error}} \xrightarrow[\text{data size } \rightarrow \infty]{} 0$$

Distribution shift

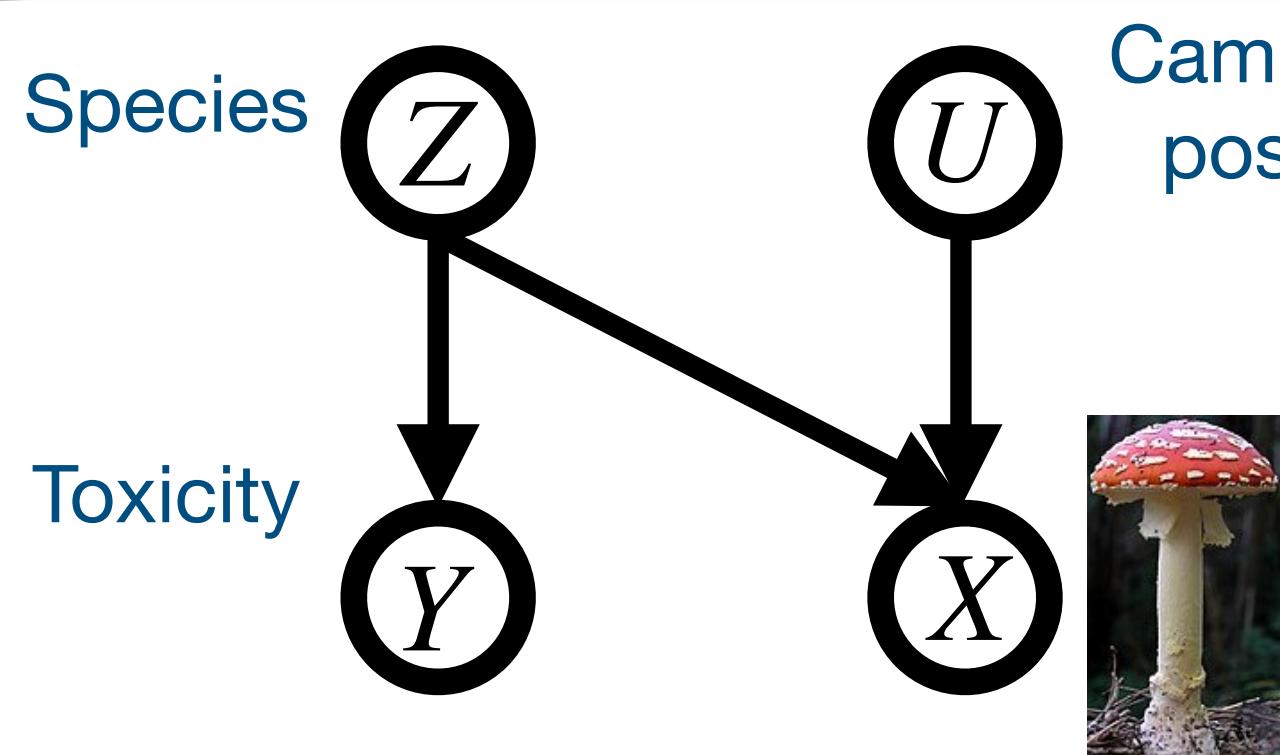


Perturbed  
Data



Adaptation of  $\hat{p}$

# Axis 3: Efficient and reliable sampling from distributions with a shared structure



$$p(Y|X) = \int p(Y|Z)p(Z|X)dZ$$

Sampling from multiple target densities

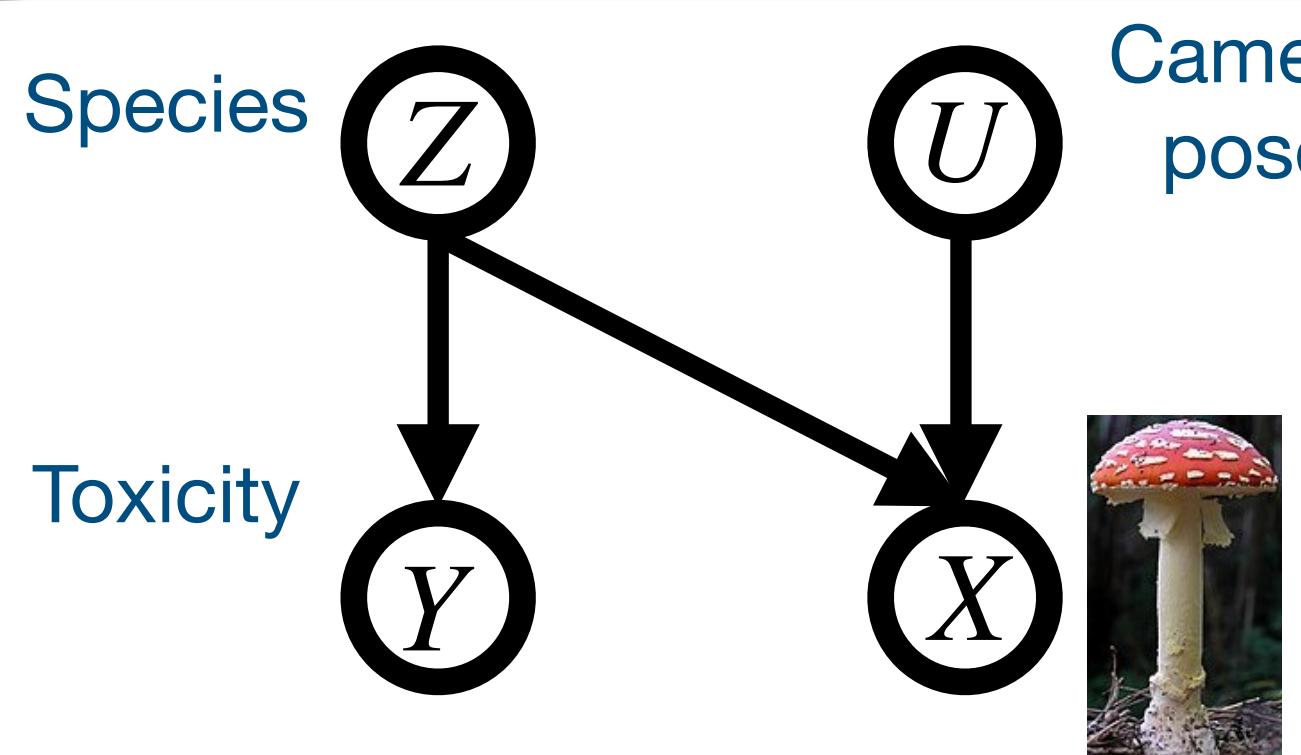
$$Z \sim p(Z|X)$$

Each data point  $X$   
defines its own target

Popular approach: *Amortized Inference*

- ❖ Fast approximate sampling
- ❖ Sample quality can degrade drastically\*
- ❖ **Unreliable:** no approximation guarantees

# Axis 3: Efficient and reliable sampling from distributions with a shared structure



$$p(Y|X) = \int p(Y|Z)p(Z|X)dZ$$

Sampling from multiple target densities

$$Z \sim p(Z|X)$$

Each data point  $X$   
defines its own target

Popular approach: *Amortized Inference*

- ❖ Fast approximate sampling
- ❖ Sample quality can degrade drastically\*
- ❖ **Unreliable:** no approximation guarantees

Model-based samplers  
for fast and reliable inference

Accounting for **shared** structure in  $p(Z|X)$

MCMC methods

**Reliability**

Approximation  
guaranteed

Models of  $p(Z|X)$

**Shared structure**

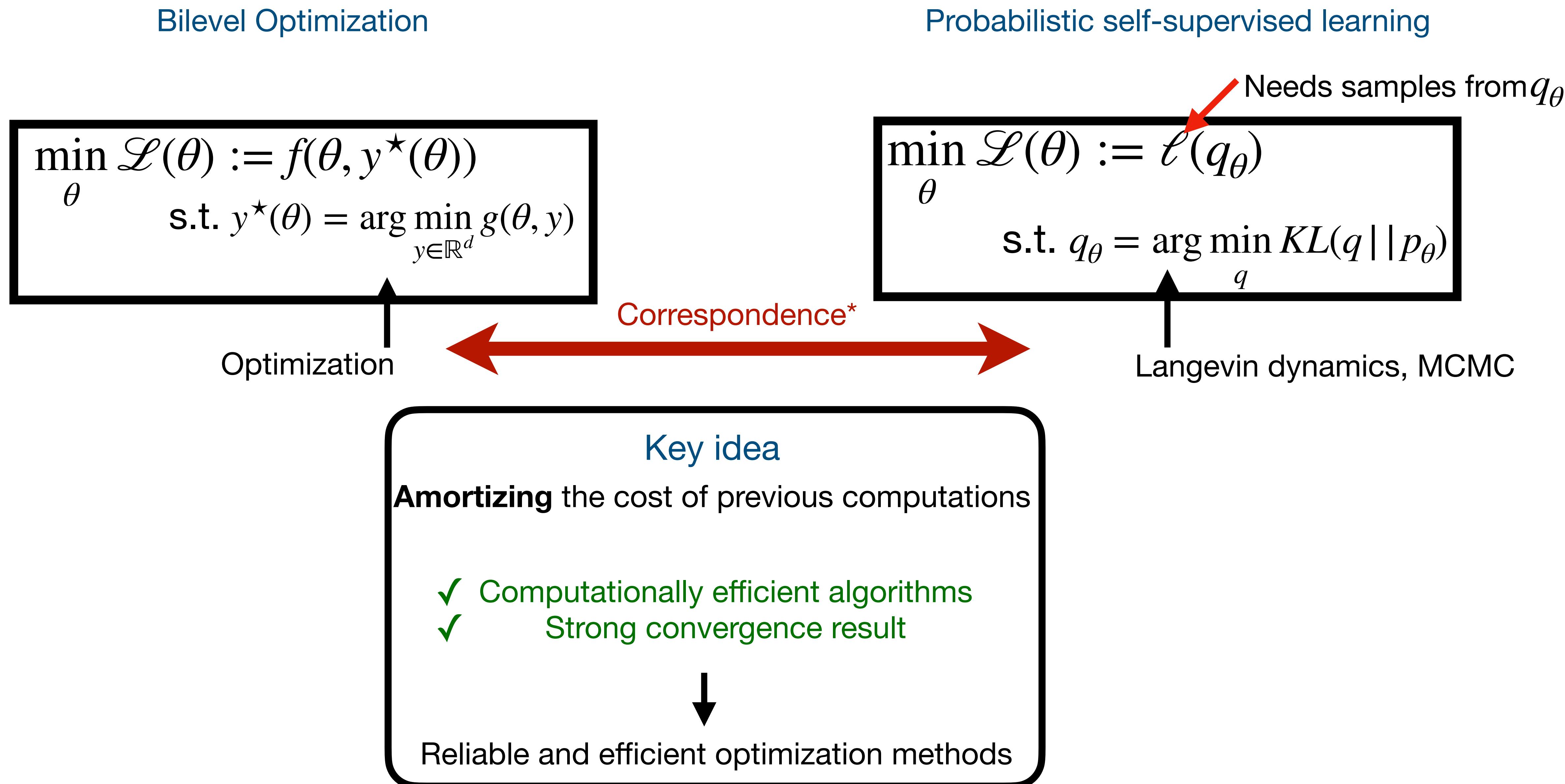
- Smoothness
- Symmetries

Model-based MCMC methods

**Reliability + Efficiency**

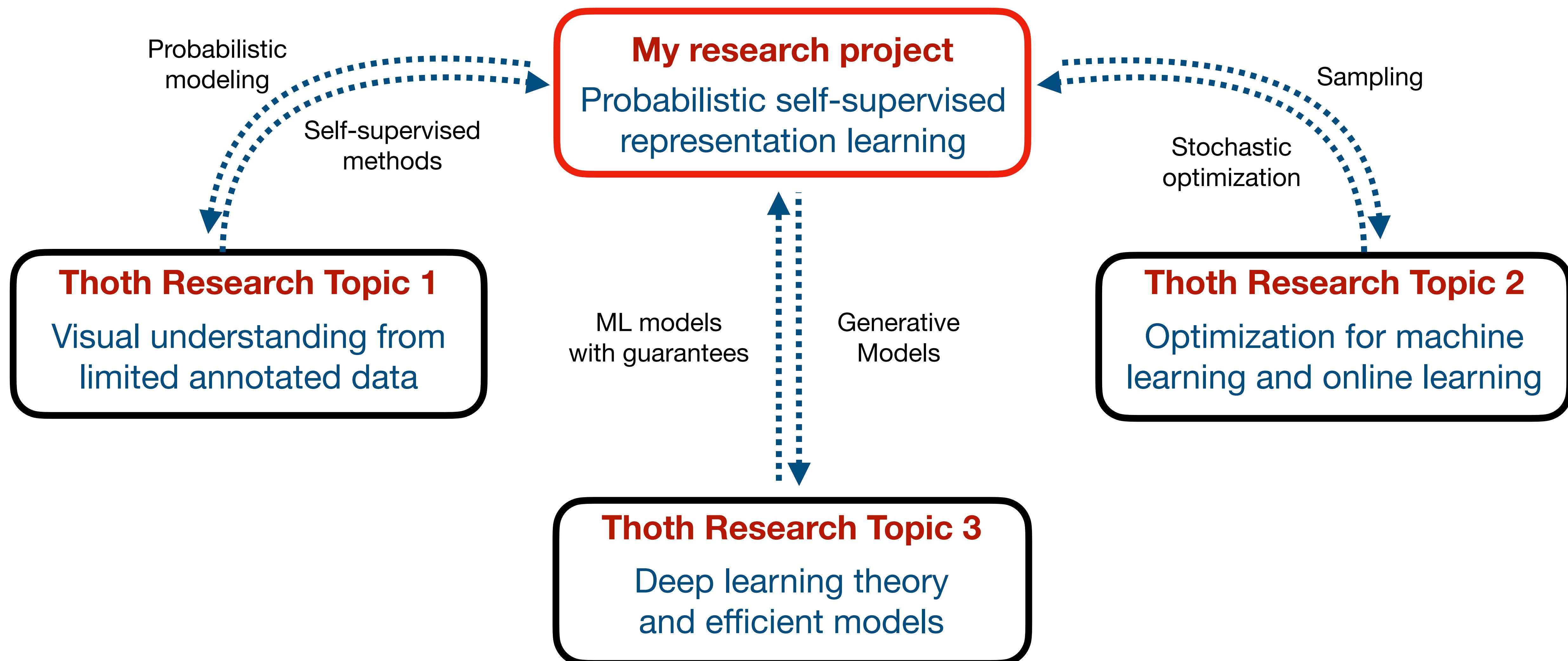
- ❖ Fast sampling exploiting shared structure
- ❖ Asymptotically exact samples
- ❖ **Strong** approximation guarantees

# Axis 4: Non-convex optimization of intractable objectives



\*The variational formulation of the Fokker--Planck equation, Jordan et. al., 1998

# Integration within Thoth team



## Potential collaborations with

### Inria teams:

Statify (Bayesian models)

Polaris (Game theory)

Dantes (Lyon) (Frugal learning)

## Potential collaborations with LJK

DAO: Optimization and  
Learning for Data Sciences

## Potential industrial collaborations in Grenoble

Naver Labs (ongoing)  
Criteo

## Scientific work and collaborations

- ❖ Experience in generative modeling, sampling, and statistics
- ❖ 17 papers at top-tier international conferences in machine learning
- ❖ 17 seminars and invited talks (USA, UK, Germany, France)
- ❖ Academic and industrial collaborations

## Multidisciplinary research project

### Addressing the lack of reliability in Machine learning

A novel probabilistic viewpoint on self-supervised data representations

#### Axes 1 & 2

A new probabilistic self-supervised learning framework

Theory and method

#### Axes 3 & 4

Efficient and reliable sampling and optimization

## ❖ 17 papers at top-tier international conferences in machine learning

- [Moskovitz et al. 2022]: Best Paper Award Finalist (Top 4 papers)
- [Arbel et al. 2021]: Long Oral Presentation (Top 3%)
- [Arbel et al. 2020]: Spotlight Presentation (Top 6%)

## ❖ 17 seminars and invited talks (USA, UK, Germany, France, Switzerland)

- Workshop on Machine Learning Augmented Sampling for the Molecular Sciences, (EPFL), May 2022
- Les journées MAS, (SMAI), August 2021
- Deep Learning Theory Kickoff Meeting, (MPI), March 2019

## ❖ Academic and industrial collaborations

- Optimal transport geometry: W. Li (South Carolina), G. Montufar (MPI)
- Reinforcement learning: A. Pacchiano (Microsoft)
- Computer Vision: T. Birdal (Stanford), U. Simsekli (Inria)
- High dimensional Sampling: A. Matthews (DeepMind), A. Doucet (Oxford)

## ❖ Service to the community

- **Reviewing:** International journals and conferences
- **Organizing seminars:** Ellis seminar series, THOTH seminar series
- **Supervision:** 3 Ph.D. students in co-supervision
- **Teaching:** 50% of a Master course at ENS Paris-Saclay (MVA)

# Major contributions in generative modelling

Before my Ph.D.

During my Ph.D.

**Unstable** algorithms for learning generative models

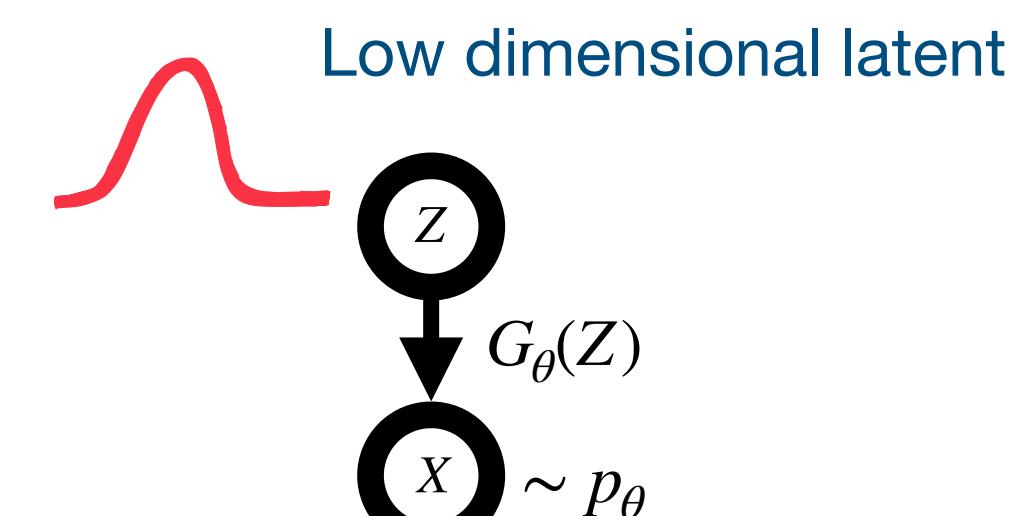
Towards **stable and principled** algorithms for learning generative models

Explicit Generative Models

$$p_\theta(X) = \frac{1}{Z_\theta} e^{-E_\theta(x)}$$

$$\min_{\theta} \mathcal{L}(p_\theta)$$

Implicit Generative Models



Scientific context: **Variational objectives**

- ✓ No need for a density for the model
- ✗ Often leads to unstable algorithms

$$\mathcal{D}(p_{data}, p_\theta) = \sup_{\phi \in \mathcal{F}} \mathbb{E}[\phi(X_{data})] - \mathbb{E}[\phi(X_{model})]$$

**Contributions:** A framework unifying both implicit and explicit models

5 publication at top-tier international conference (NeurIPS, ICLR, AISTATS), more than 600 citations

# Contributions in high dimensional sampling

Before my Ph.D.

Limited reliability of  
model-based sampling

During my Ph.D.

Efficient model-based samplers with theoretical guarantees

## High dimensional sampling\*

A fundamental problem

$$X \sim \pi$$

## Challenges

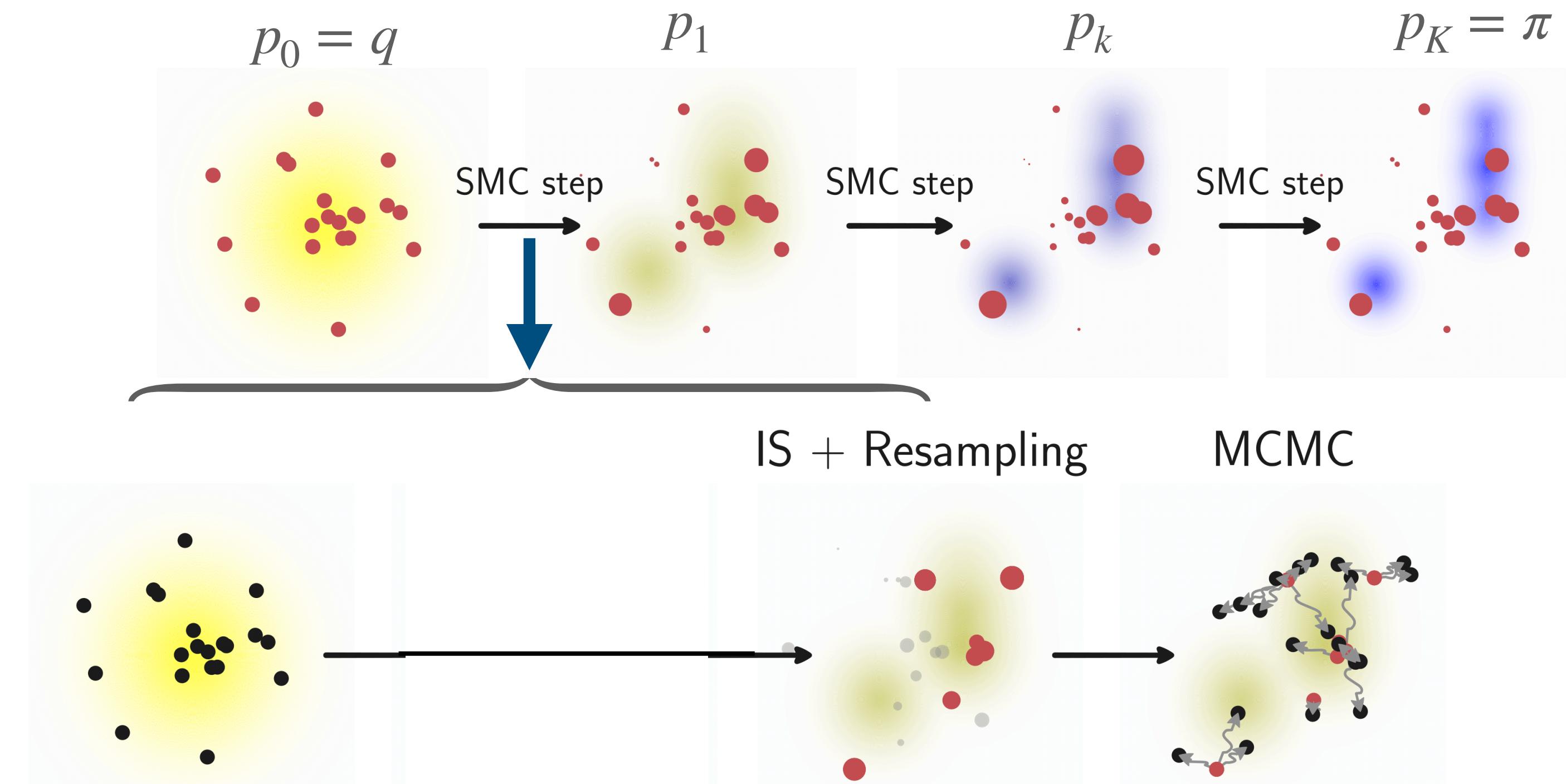
- Complex multi-modal distributions
- Curse of dimensions

## Scientific context

Sequential Monte-Carlo (SMC) methods

✓ Reliable + Strong theoretical guarantees

✗ Agnostic to the target density → High cost in general



**Contributions:** Principled and efficient model-based SMC samplers that capture structure in the target density

1 publication at top-tier international conference (ICML 2021 Long Oral, top 3% submissions)

# Contributions in high dimensional sampling

Before my Ph.D.

Limited reliability of  
model-based sampling

During my Ph.D.

Efficient model-based samplers with theoretical guarantees

## High dimensional sampling\*

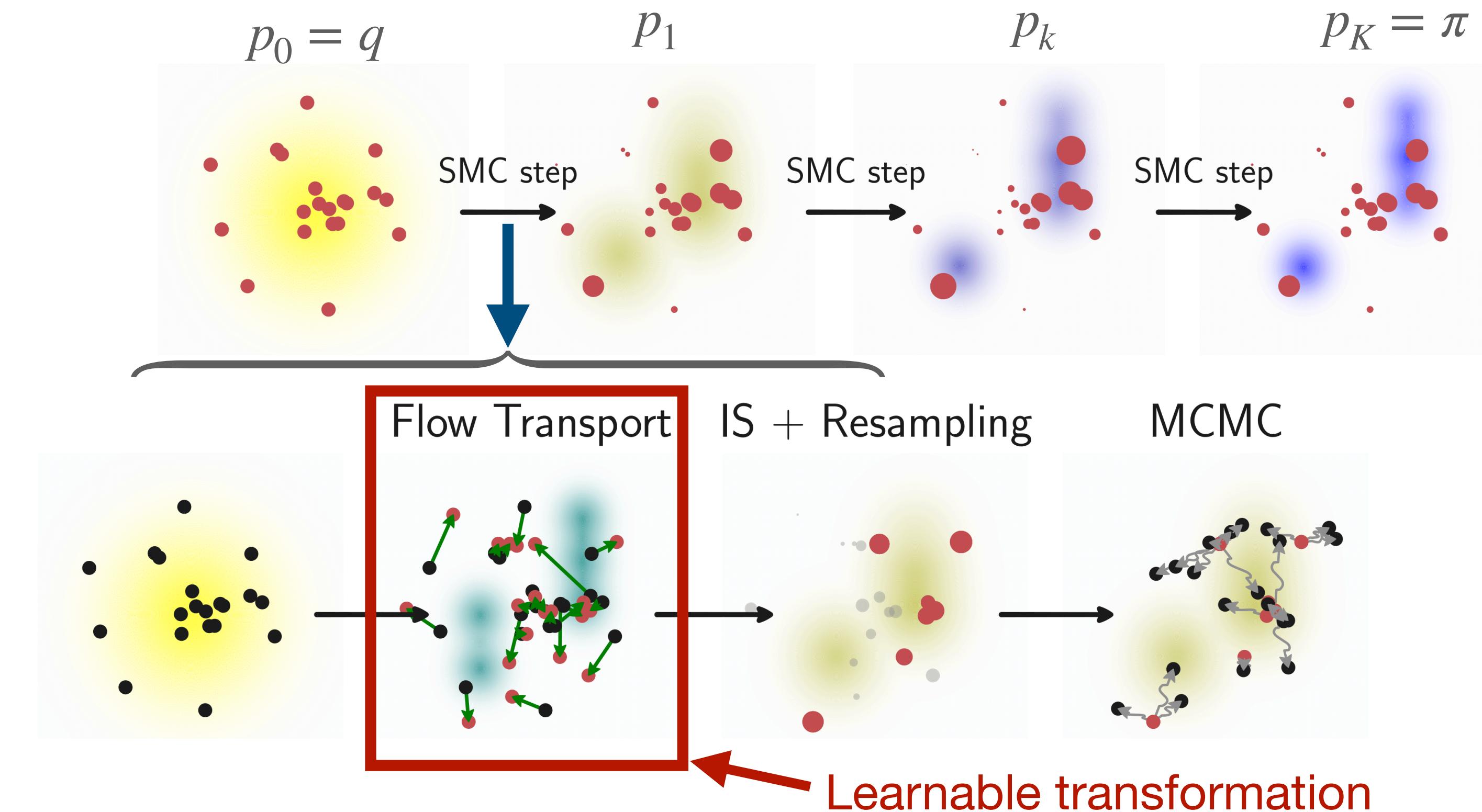
A fundamental problem

$$X \sim \pi$$

## Key idea

*Exploiting structure in the target density through  
modeling for more efficient sampling*

- SMC with a learnable transformation
- ✓ Fast sampling exploiting shared structure
  - ✓ Provably improves over than SMC
  - ✓ Strong approximation guarantees



**Contributions:** Principled and efficient model-based SMC samplers that capture structure in the target density

1 publication at top-tier international conference (ICML 2021 Long Oral, top 3% submissions)

# Contributions in large-scale optimization for generative models

Before my Ph.D.

**Slow convergence** of gradient methods  
for ill-conditioned generative models

During my Ph.D.

**Fast and scalable** optimizers that are **robust** to ill-conditioned models

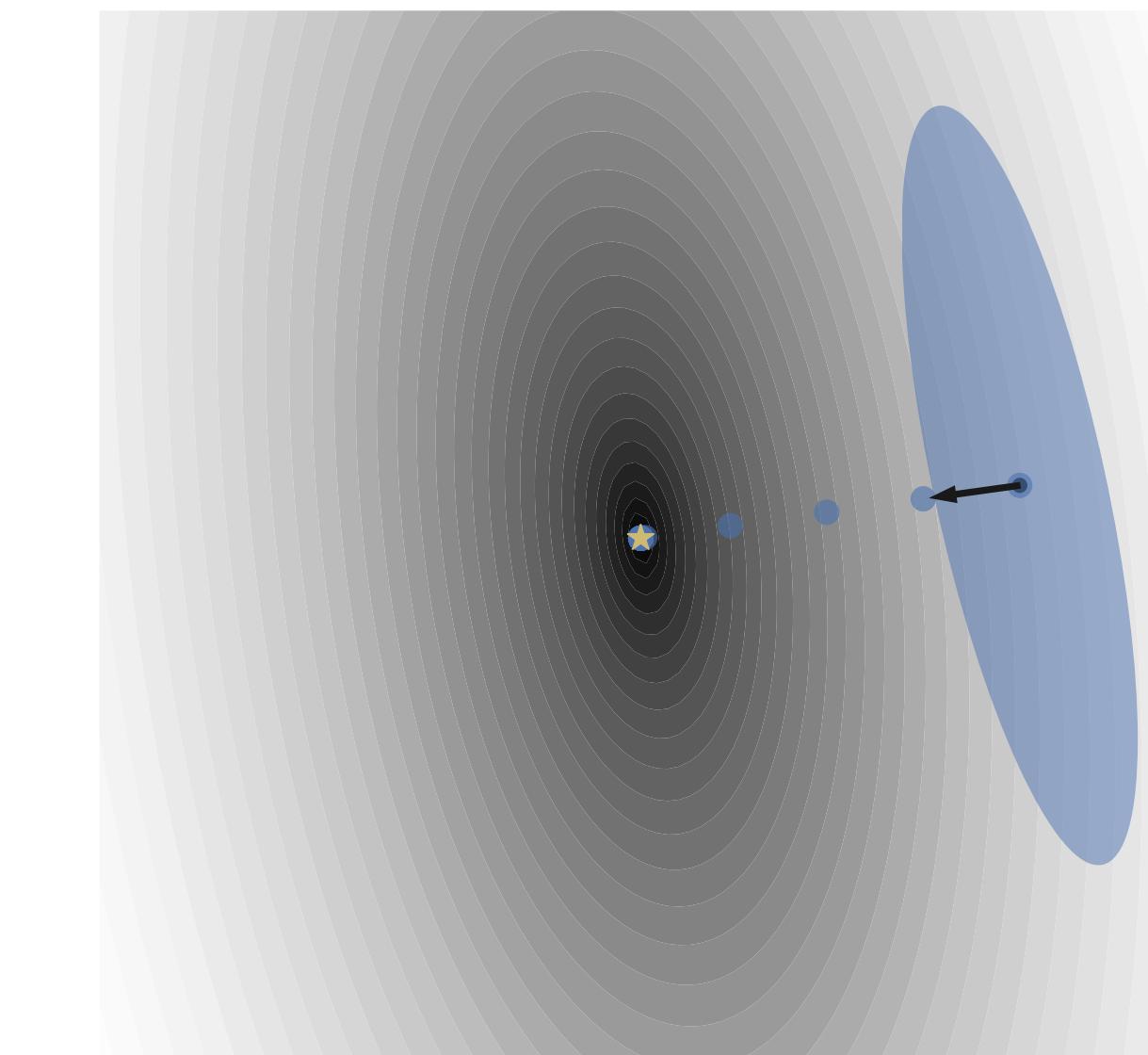
## Optimizing a generative model

$$\min_{\theta} \mathcal{L}(p_{\theta})$$

Dependence on  $\theta$  only through  $p_{\theta}$

Fisher natural gradient descent

$$\theta_{k+1} = \theta_k - \gamma G(\theta_k)^{-1} \nabla \mathcal{L}(p_{\theta_k})$$



## Challenges

Hard to optimize when models are ill-conditioned

## Scientific context

Fisher natural gradient

- ✓ Robust to ill-conditionned models
- ✗ Not defined for Implicit generative models
- ✗ Expensive to compute for large models

**Contributions:** Robust, fast and scalable optimizers based on optimal transport geometry

2 publications at top-tier international conferences (ICLR 2020 Spotlight, top 6% submissions, ICLR 2021)

# Contributions in large-scale optimization for generative models

Before my Ph.D.

During my Ph.D.

**Slow convergence** of gradient methods  
for ill-conditioned generative models

**Fast** and **scalable** optimizers that are **robust** to ill-conditioned models

## Optimizing a generative model

$$\min_{\theta} \mathcal{L}(p_{\theta})$$

Dependence on  $\theta$  only through  $p_{\theta}$

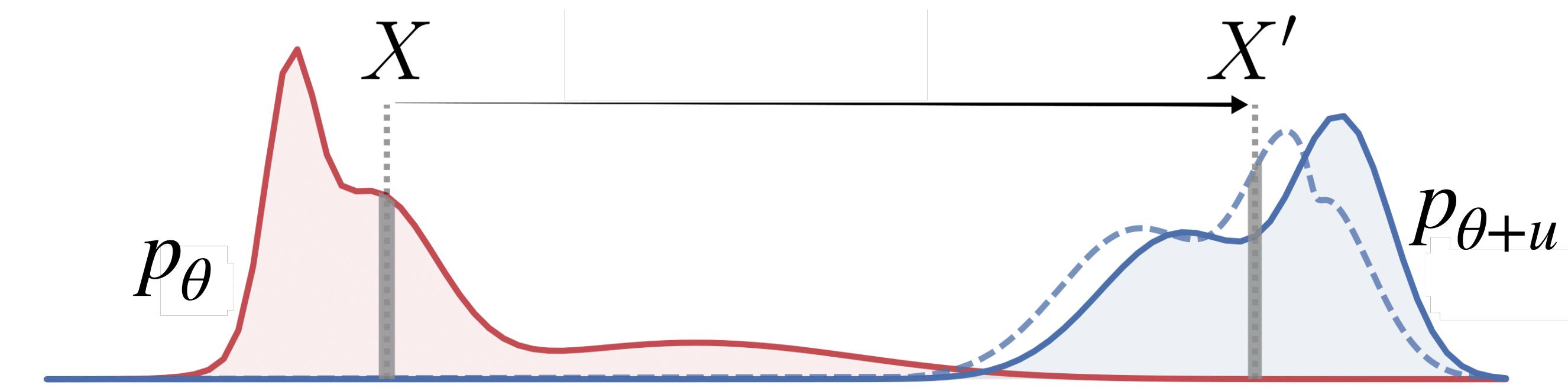
## Key idea

Wasserstein Natural gradient

- ✓ Robust to ill-conditionned models
- ✓ Well-defined implicit generative models
- ✓ Scalable and efficient estimation using kernel methods
- ✓ Statistical guarantees on the estimation

Wasserstein natural gradient descent

$$\theta_{k+1} = \theta_k - \gamma G(\theta_k)^{-1} \nabla \mathcal{L}(p_{\theta_k})$$



$$\frac{1}{2} u^T G(\theta) u \approx W^2(p_{\theta}, p_{\theta+u})$$

**Contributions:** Robust, fast and scalable optimizers based on optimal transport geometry

2 publications at top-tier international conferences (ICLR 2020 Spotlight, top 6% submissions, ICLR 2021)