

# Uncertainty in Deep Learning: A Probabilistic Robotics Perspective

**Jongseok Lee**

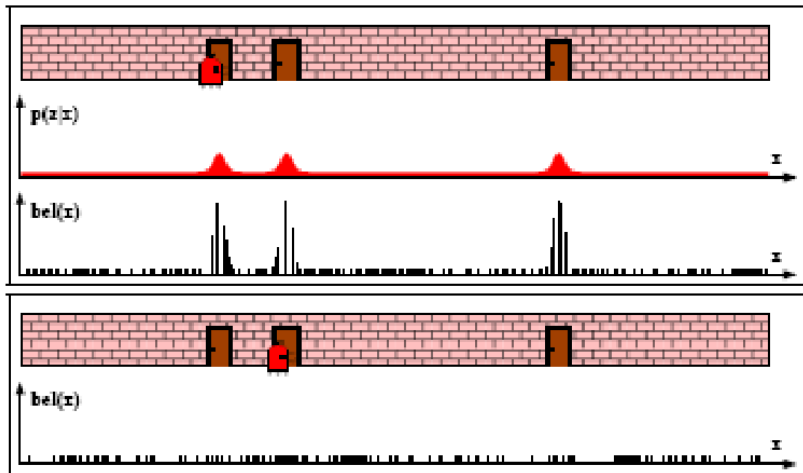
**PhD Defense**



# Is this our Sisyphean climb?

Since 1990s

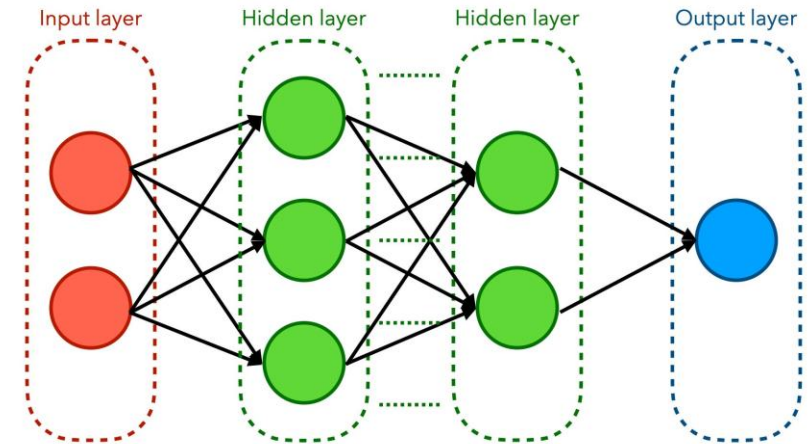
## Probabilistic Robotics



- Do not optimize but average.
- Pay tribute to uncertainty.

Since 2010s

## AI-based Robotics



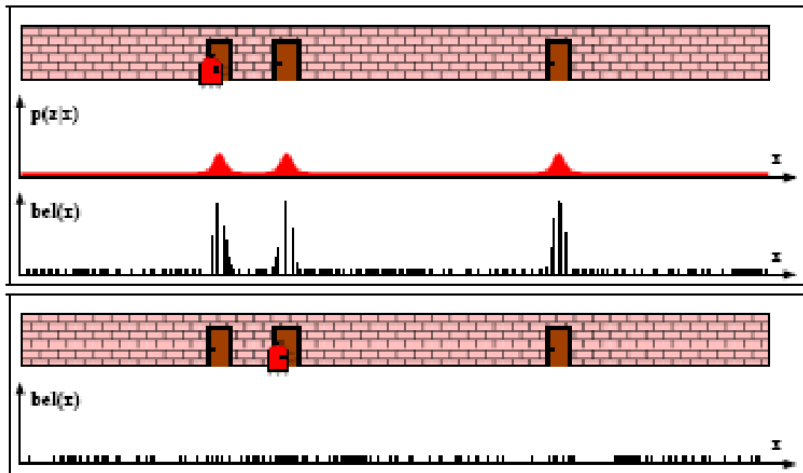
- Optimize as much as possible.
- Worships performance, not doubt.



# No Sisyphus – Hercule’s bridge

Since 1990s

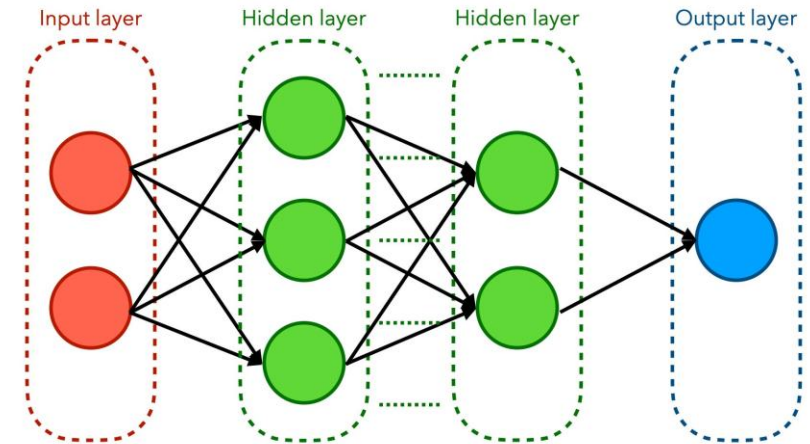
## Probabilistic Robotics



- Do not optimize but average.
- Pay tribute to uncertainty.

Since 2010s

## AI-based Robotics



- Optimize as much as possible.
- Worships performance, not doubt.

**AI-based robotics that reason about uncertainty – no maximum likelihood**

# Why should we care?

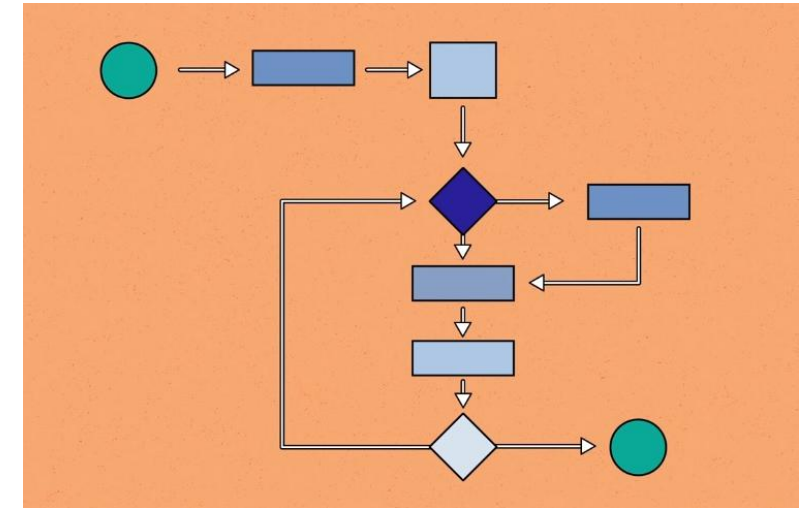
## Safety & Risk



## Trustworthy AI



## Better Algorithms



# Preamble: Bayesian statistics

**Recap 1:** Probability is defined as degree of one's belief about an event.

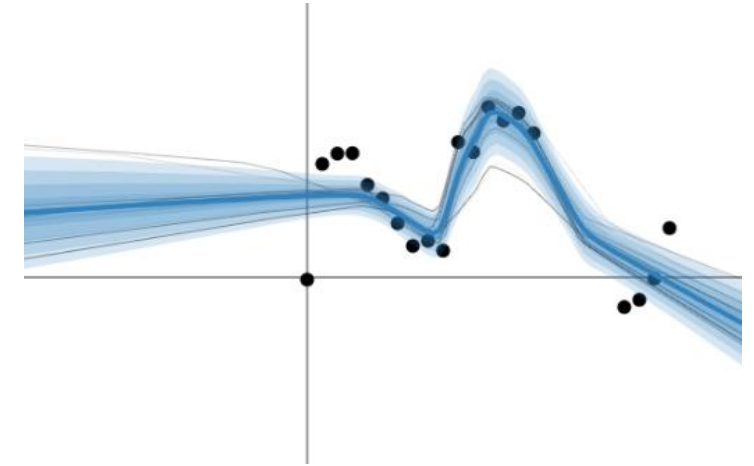
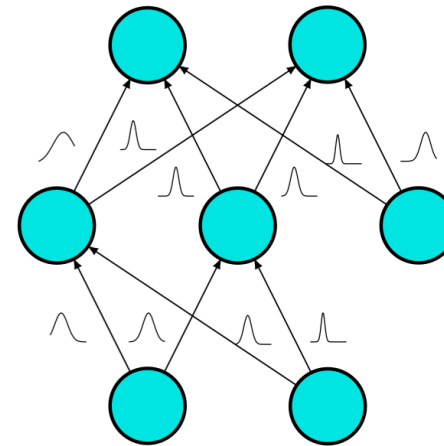
$$\begin{array}{ccc} \text{posterior} & & \text{prior} \\ \downarrow & & \downarrow \\ p(\boldsymbol{\theta} | \boldsymbol{D}) = & \frac{p(\boldsymbol{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\boldsymbol{D})} & \end{array}$$

**Recap 2:** Bayes Theorem updates prior belief to posterior given data.

# Preamble: Bayesian Deep Learning

Intractable!

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$



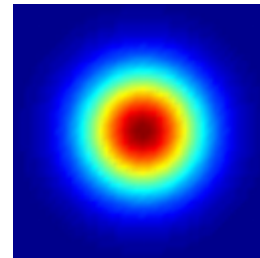
Intractable!

$$p(y^*|x^*, D) = \int p(y^*|x^*, \theta)p(\theta|D)d\theta$$

- + Inclusion of domain knowledge via prior.
- + Quantifies uncertainty about the model.
- + Average different hypothesis about events.

# Preamble: Being Bayesian is hard

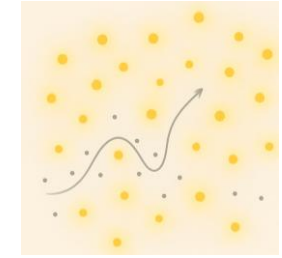
- Tisby et al. 1989
- Mackay 1992
- Hinton and Van Camp 1993
- Neal 1995
- AI winter**
- Graves 2011
- Hernandez-Lobato and Adams 2015
- **MC-dropout**  
(Gal and Ghahramani 2016)
- Deep Ensemble  
(Lakshminarayanan et al. 2017)
- **When I started**



Uninformative prior



Bernoulli posterior



Sampling

- Prior specification is non-trivial.
- Posterior is crudely approximated.
- Sampling is computationally expensive.

# Key research objective

## Challenges:

- Prior specification is non-trivial.
- Posterior is crudely approximated.
- Sampling is more expensive.

**Bayesian methods that are well applicable in AI-based robotics.**

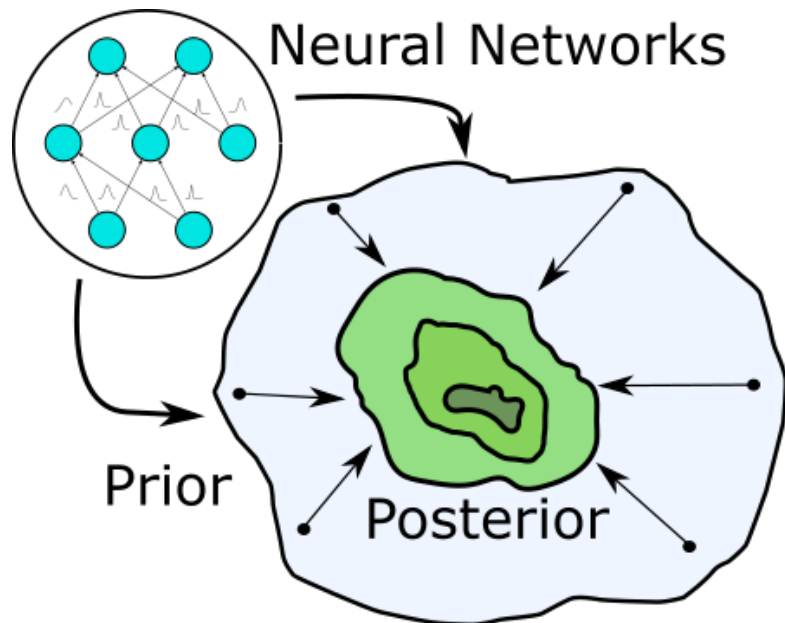
## Promises:

- + Inclusion of domain knowledge via prior.
- + Quantifies uncertainty about the model.
- + Average different hypothesis about events.



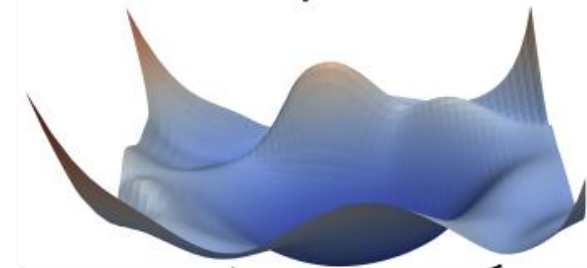
# Contributions

## On Priors



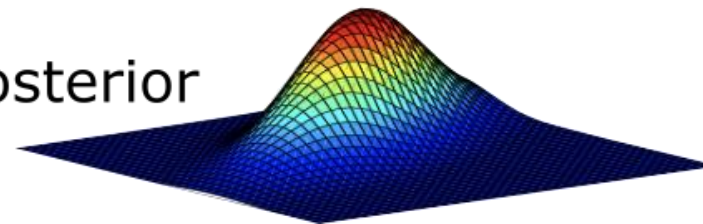
## On Posteriors

Loss landscape



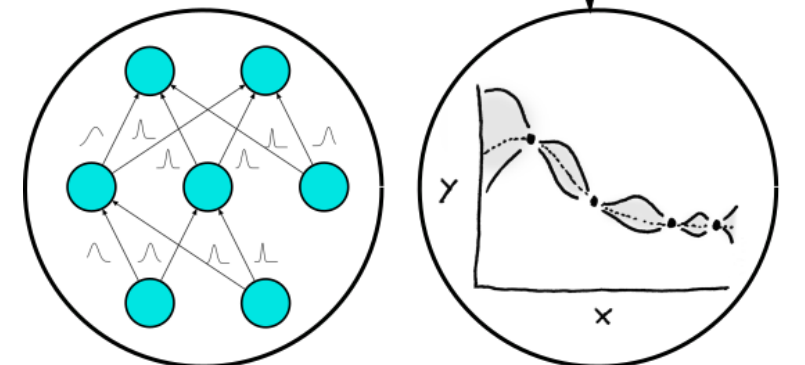
Curvature

Posterior



## On Predictions

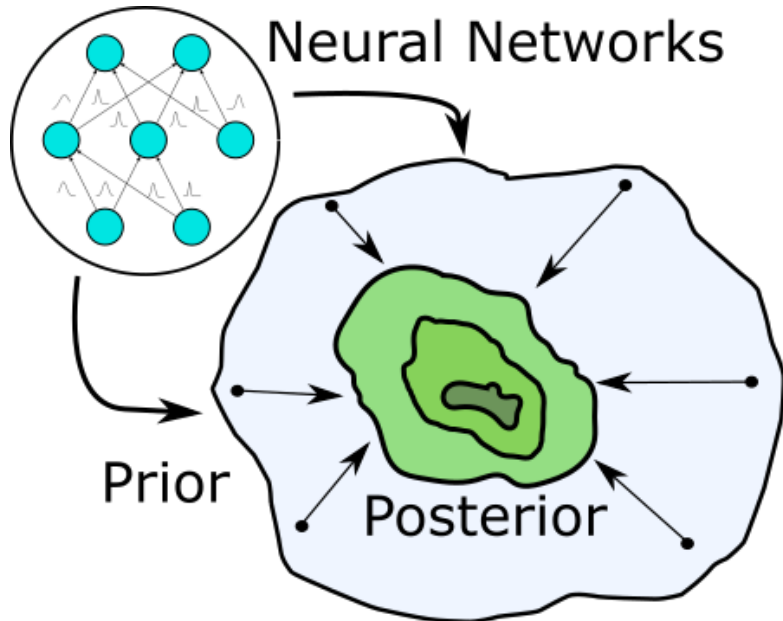
Neural Networks



Gaussian Processes

# Overview

## On Priors

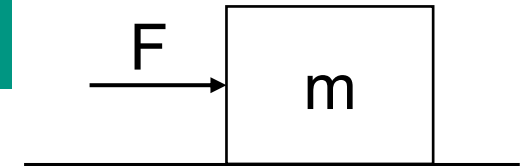


How to specify more meaningful priors for uncertainty and generalization?

# Challenge: Issue of interpretability

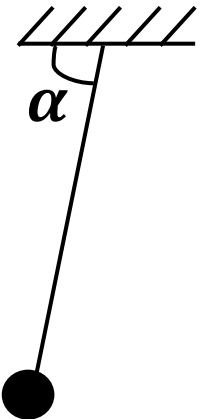
**$F=ma$  where a constant  $m$  is a random variable. Prior  $p(m)$ ?**

- Gamma distribution. Mass  $m$  is a quantity of positive values.



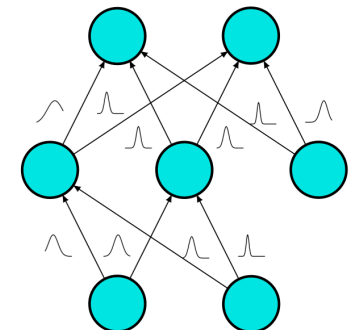
**A pendulum where an angle  $\alpha$  is a random variable. Prior  $p(\alpha)$ ?**

- Von Mises distribution to reflect that angle  $\alpha$  ranges from zero to  $2\pi$ .



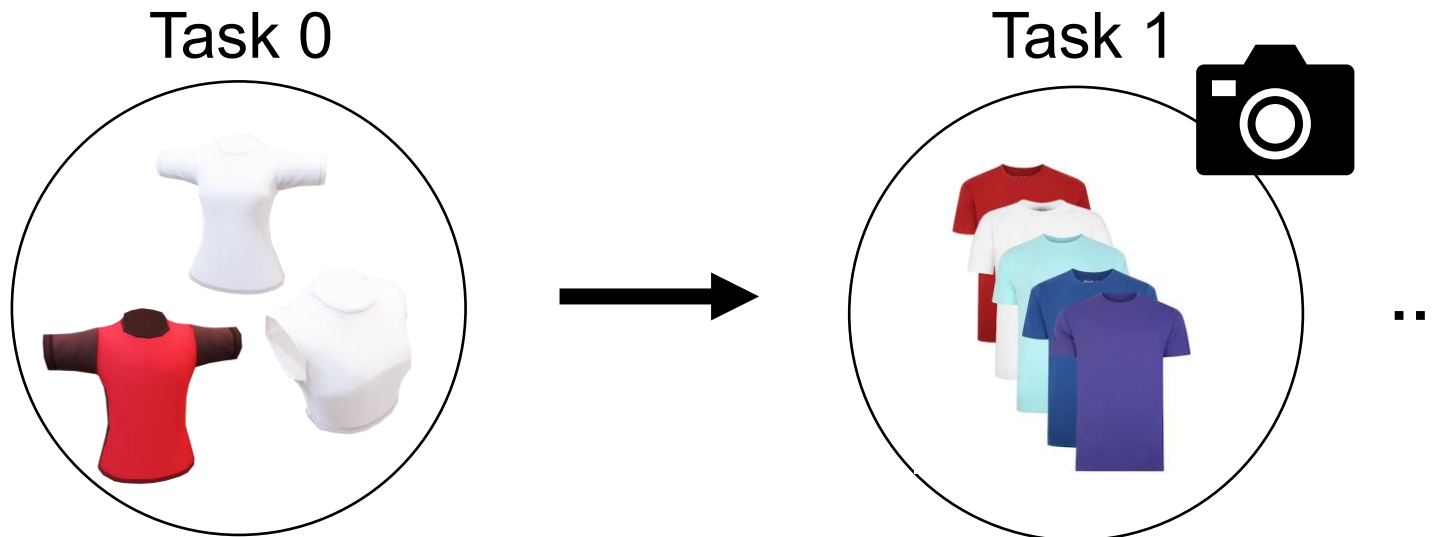
**A network with the weights  $\theta$  as random variable. Prior  $p(\theta)$ ?**

- Zero mean isotropic Gaussian:  $p(\theta) = \mathcal{N}(\mathbf{0}, \gamma I)$ .
- Despite reported signs of prior misspecification (Fortuin, 2022).

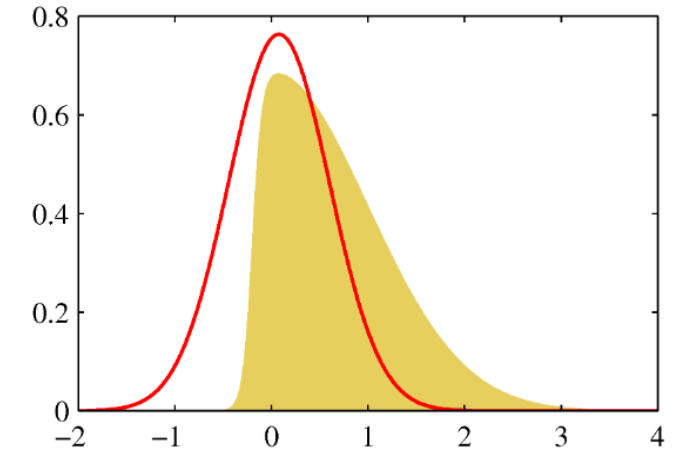


# Idea: Sequential Bayes for neural networks

## Learning priors from data:



$$\begin{aligned} 0. & p(\theta^{(0)}) = \mathcal{N}(\theta^{(0)} | \mathbf{0}, \tau I) \\ 1. & p(\theta^{(0)} | D^{(0)}) \\ 2. & p(\theta^{(1)}) = p(\theta^{(0)} | D^{(0)}) \\ 3. & p(\theta^{(1)} | D^{(1)}) \end{aligned}$$

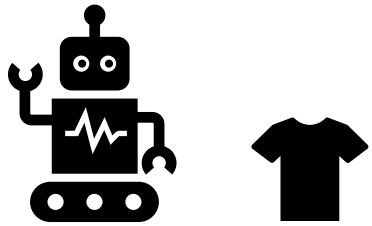


$$p(\theta | D) \approx \mathcal{N}(\theta | \hat{\theta}, H^{-1})$$

$$\hat{\theta} \in \operatorname{argmax} p(\theta | D)$$

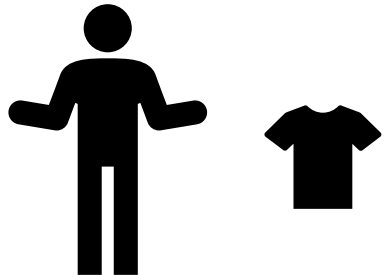
$$H = H_{\text{likelihood}} + H_{\text{prior}}$$

# Evaluation: Robot asks help and learn



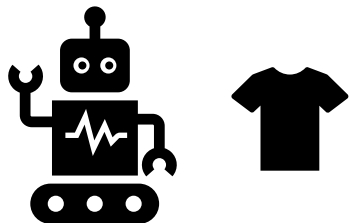
I don't know what this object is.  
Can you help?

- **Task:** Uncertainty estimation.  
→ Evaluates prior specification.



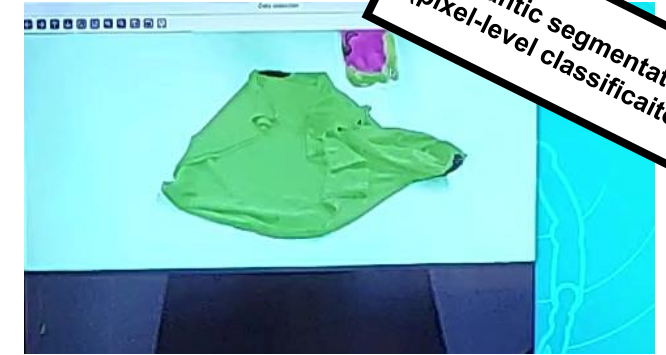
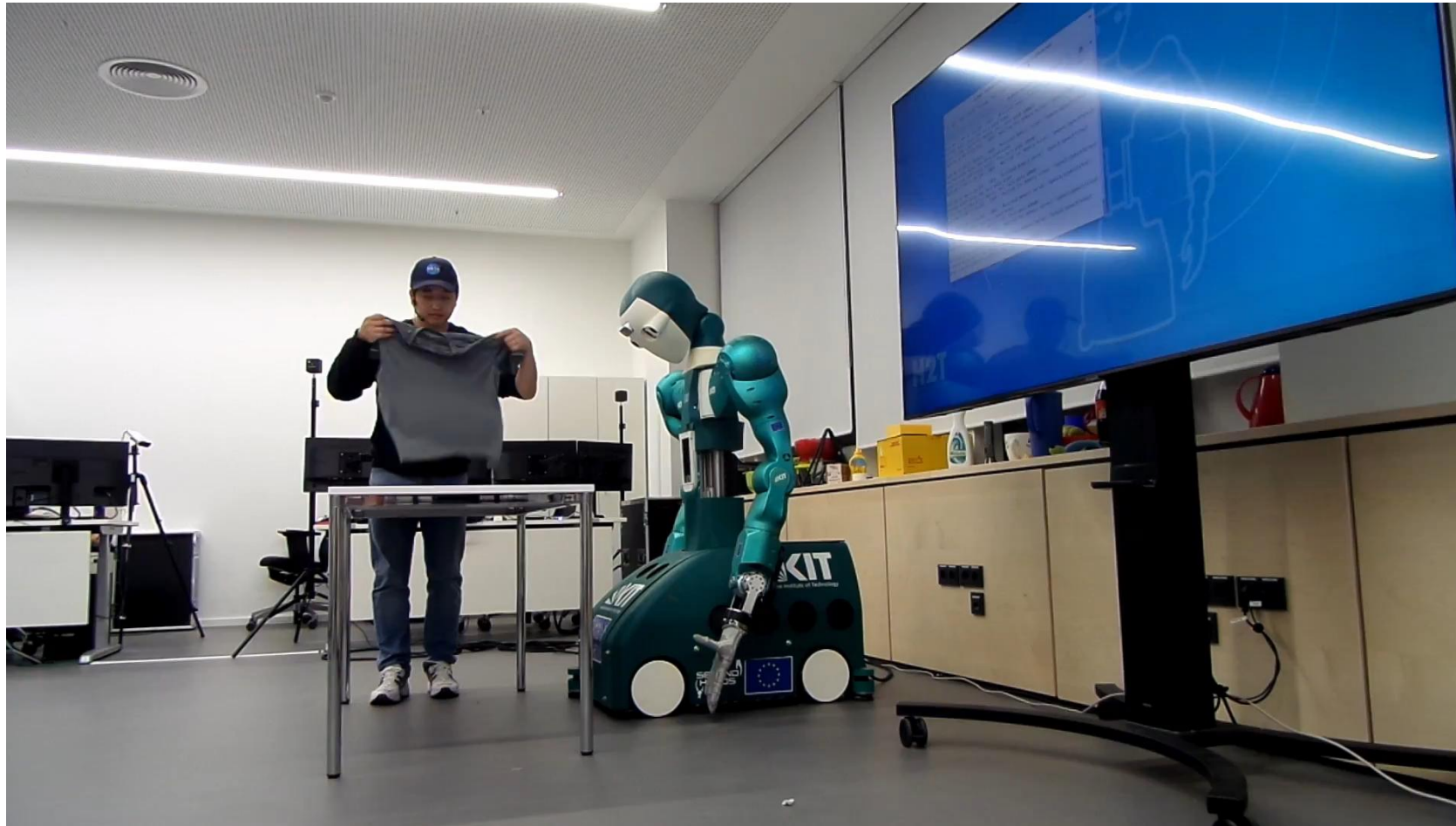
This is a T-shirt.  
Let me demonstrate to you.

- **Task:** Learning from small data.  
→ Evaluates generalization.



Let me train.  
Now I know about this object.

# Demonstration on a humanoid robot

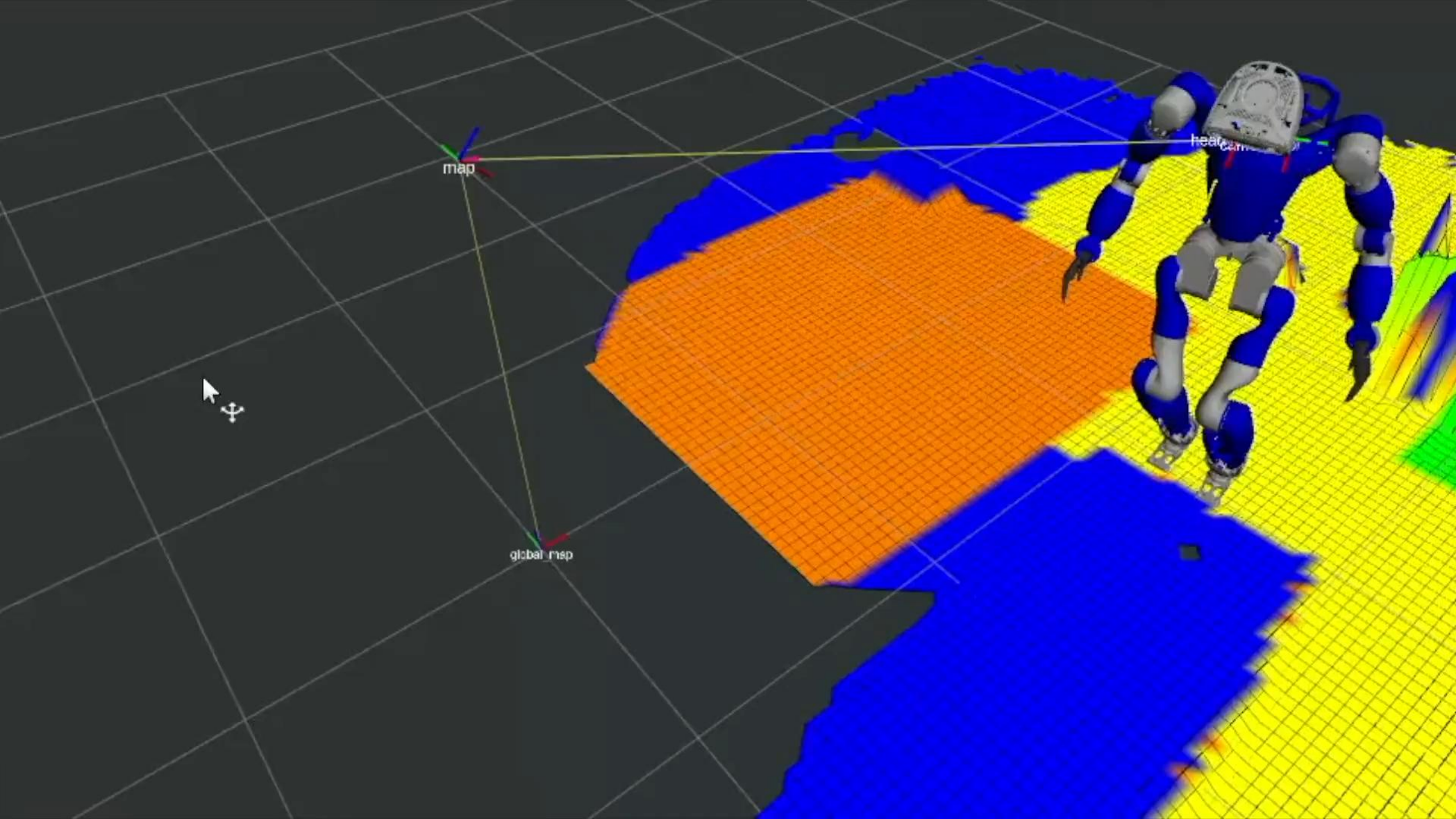


	Success rate	ECE
CLEVERv1	$0.887 \pm 0.049$	$0.060 \pm 0.017$
CLEVERv2	$0.826 \pm 0.053$	$0.092 \pm 0.036$
Vanilla	$0.801 \pm 0.054$	$0.177 \pm 0.021$

	Precision	Nr. Queries
CLEVERv1	$0.933 \pm 0.020$	$1.539 \pm 0.544$
CLEVERv2	$0.900 \pm 0.034$	$2.440 \pm 0.866$
Vanilla	$0.817 \pm 0.039$	$4.320 \pm 0.992$

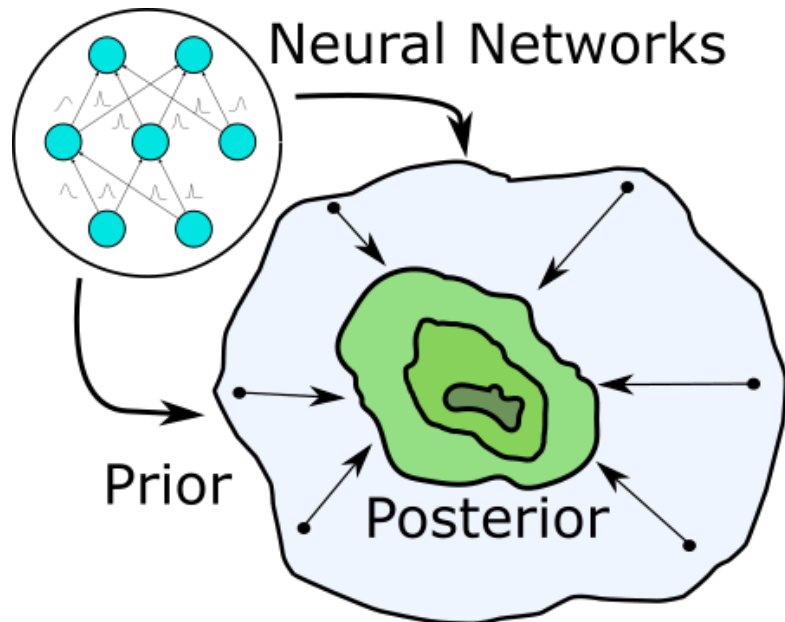
Prior helps uncertainty & generalization.





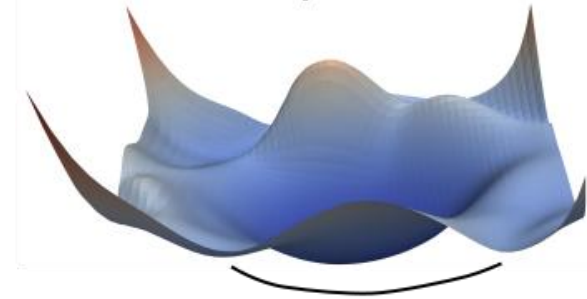
# Overview

## On Priors



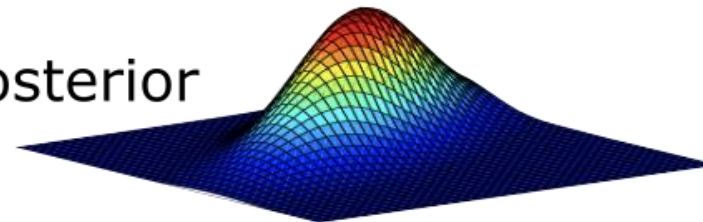
## On Posteriors

Loss landscape



Curvature

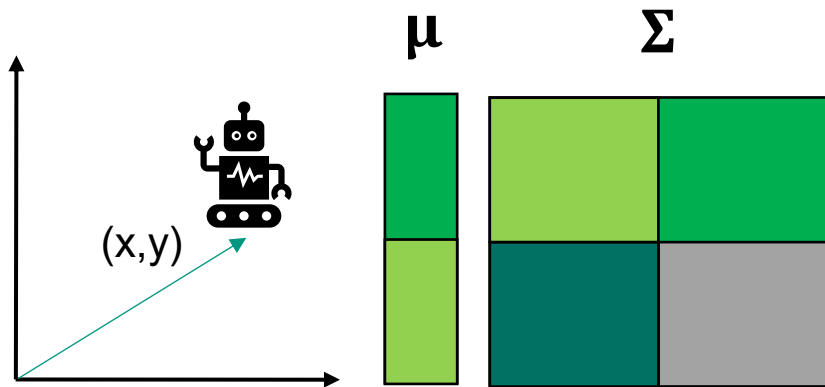
Posterior



How to infer complex posteriors while scaling to large dataset and models?

# Challenge: Curse of Dimensionality

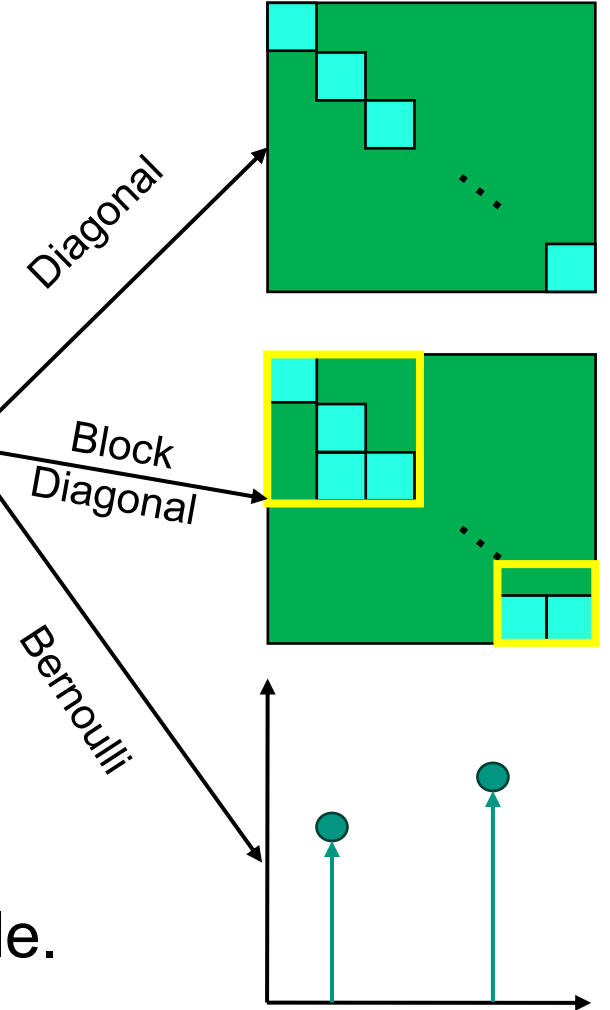
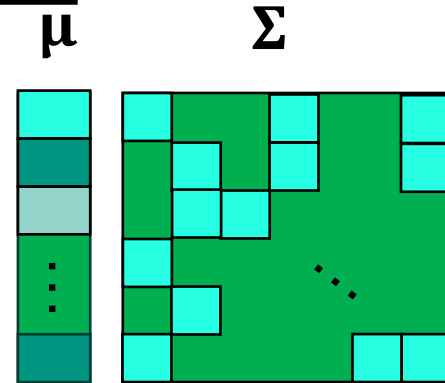
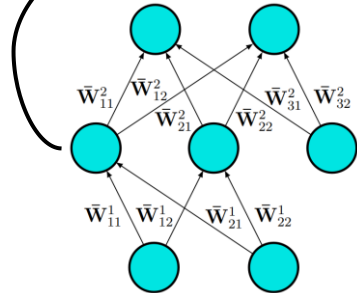
## Simple 2D problem



$$\theta \propto e^{-\frac{1}{2}(\theta - \mu)^T \Sigma^{-1} (\theta - \mu)}$$

## Deep Learning

Millions of parameters!



Computational costs:

■ Storage:  $O(n^2)$

■ Inversion:  $O(n^3)$

Note:  $n$  is dim of random variable.

# Idea: Information Form

- Estimate uncertainty in information form:

$$p(\boldsymbol{\theta}|\mathcal{D}) \approx \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \mathcal{N}^{-1}(\boldsymbol{\eta}, \boldsymbol{\Lambda}) \text{ with } \boldsymbol{\eta} = \boldsymbol{\mu}^T \boldsymbol{\Sigma} \text{ and } \boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}.$$

- Why information form?

$$p(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \mathcal{N} \left( \begin{bmatrix} \boldsymbol{\mu}_\alpha \\ \boldsymbol{\mu}_\beta \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{\alpha\alpha} & \boldsymbol{\Sigma}_{\alpha\beta} \\ \boldsymbol{\Sigma}_{\beta\alpha} & \boldsymbol{\Sigma}_{\beta\beta} \end{bmatrix} \right) = \mathcal{N}^{-1} \left( \begin{bmatrix} \boldsymbol{\eta}_\alpha \\ \boldsymbol{\eta}_\beta \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Lambda}_{\alpha\alpha} & \boldsymbol{\Lambda}_{\alpha\beta} \\ \boldsymbol{\Lambda}_{\beta\alpha} & \boldsymbol{\Lambda}_{\beta\beta} \end{bmatrix} \right)$$

- Goal: compute conditioning or posterior calculation  $p(\boldsymbol{\alpha}|\boldsymbol{\beta}) = \frac{p(\boldsymbol{\alpha}, \boldsymbol{\beta})}{p(\boldsymbol{\beta})}$ :

Covariance form:  $\boldsymbol{\Sigma}' = \boldsymbol{\Sigma}_{\alpha\alpha} - \boldsymbol{\Sigma}_{\alpha\beta} \boldsymbol{\Sigma}_{\beta\beta}^{-1} \boldsymbol{\Sigma}_{\beta\alpha}$ .

Information form:  $\boldsymbol{\Lambda}' = \boldsymbol{\Lambda}_{\alpha\alpha}$ .

# Idea: Information Form

$$O(n^3) \rightarrow O(L^3) \text{ for } L \ll n$$

- Estimate uncertainty in information form:

$$p(\theta|D) \approx \mathcal{N}(\mu, \Sigma) = \mathcal{N}^{-1}(\eta, \Lambda) \text{ with } \eta = \mu^T \Sigma \text{ and } \Lambda = \Sigma^{-1}.$$

- Why information form?

$$p(\alpha, \beta) = \mathcal{N} \left( \begin{bmatrix} \mu_\alpha \\ \mu_\beta \end{bmatrix}, \begin{bmatrix} \Sigma_{\alpha\alpha} & \Sigma_{\alpha\beta} \\ \Sigma_{\beta\alpha} & \Sigma_{\beta\beta} \end{bmatrix} \right) = \mathcal{N}^{-1} \left( \begin{bmatrix} \eta_\alpha \\ \eta_\beta \end{bmatrix}, \begin{bmatrix} \Lambda_{\alpha\alpha} & \Lambda_{\alpha\beta} \\ \Lambda_{\beta\alpha} & \Lambda_{\beta\beta} \end{bmatrix} \right)$$

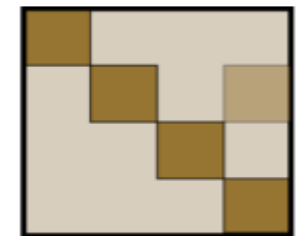
- Goal: compute marginalization or prediction calculation  $p(\alpha) = \int p(\alpha, \beta) d\beta$ :

Covariance form:  $\Sigma = \Sigma_{\alpha\alpha}$

Information form:  $\Lambda = \Lambda_{\alpha\alpha} - \Lambda_{\alpha\beta} \Lambda_{\beta\beta}^{-1} \Lambda_{\beta\alpha}$



Vs

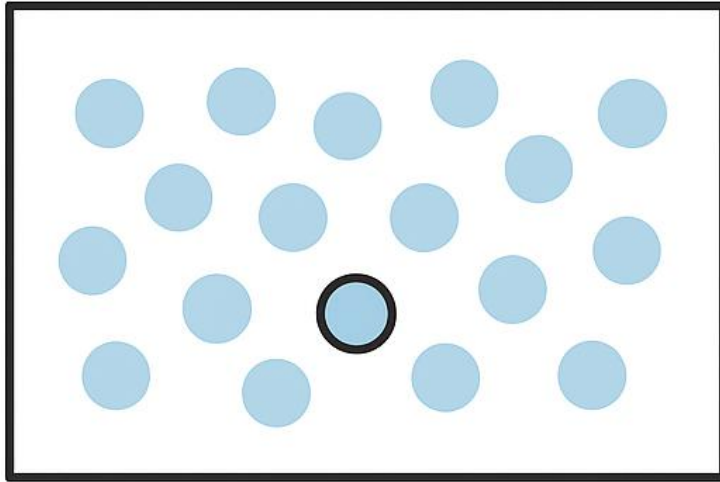


Covariance matrix

Information matrix

(Sagun et al. 2018)

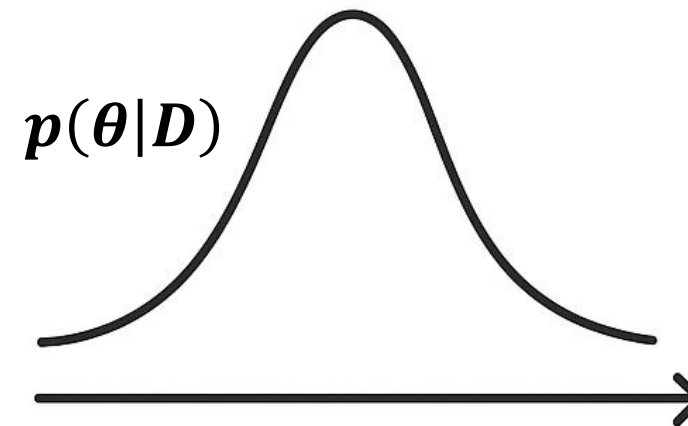
# Evaluation: Pool-based Active Learning



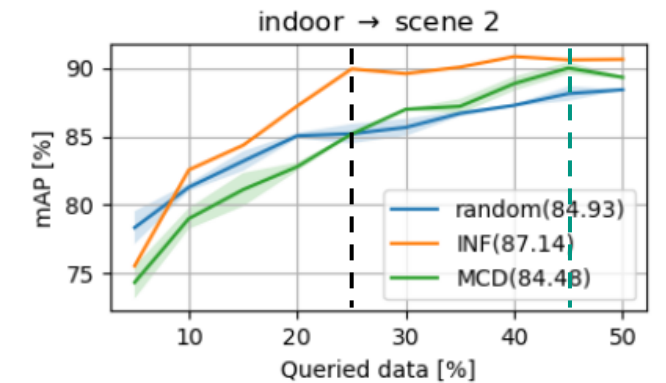
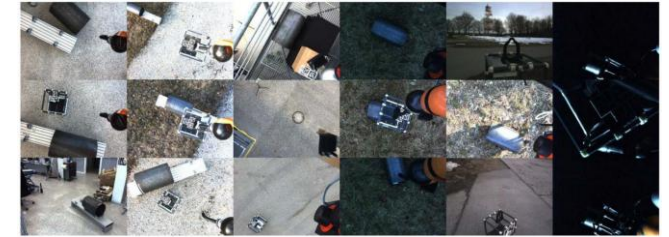
Pool of unlabelled data

**Goal:** select the most informative data in order to reduce labelling costs.

Pick data that my model doesn't know!

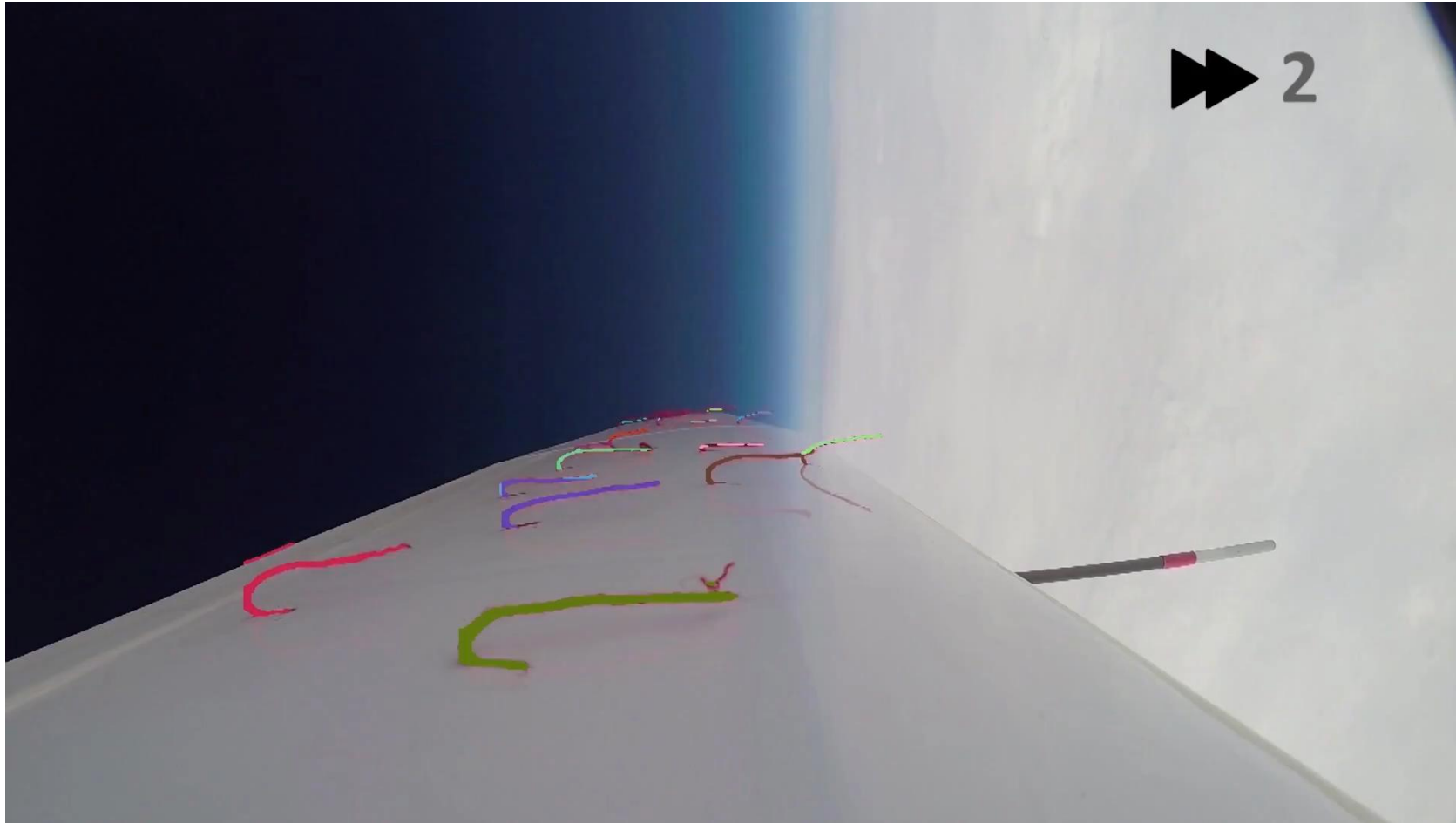


# Outside lab: Reducing labelling costs



Expressive posterior →  
pool-based active learning  
performance.

# AI4Science: Reducing labelling costs





A man in a black jacket and blue jeans stands at a white workstation, operating a computer terminal with multiple monitors and a keyboard. The workstation is positioned on the left side of the exhibition area.

A white robotic arm with orange joints is mounted on a white and yellow workstation. The arm is positioned over a table with various components and tools. A man in a dark suit is standing next to the workstation, speaking into a microphone.

A man in a dark suit is speaking into a microphone, addressing the crowd. He is standing in the center of the exhibition area, facing the audience.

A large crowd of people, including men and women, is gathered around the workstation. Many are holding up their smartphones to take photos or videos of the robotic arm. The crowd is diverse in age and attire, with many wearing lanyards.



IBG

karmetal

MOTION

ALW



KURA Innoventory Award

UNIVERSITY WATER

OMNO

OMNO

OMNO

OMNO

OMNO

OMNO

OMNO

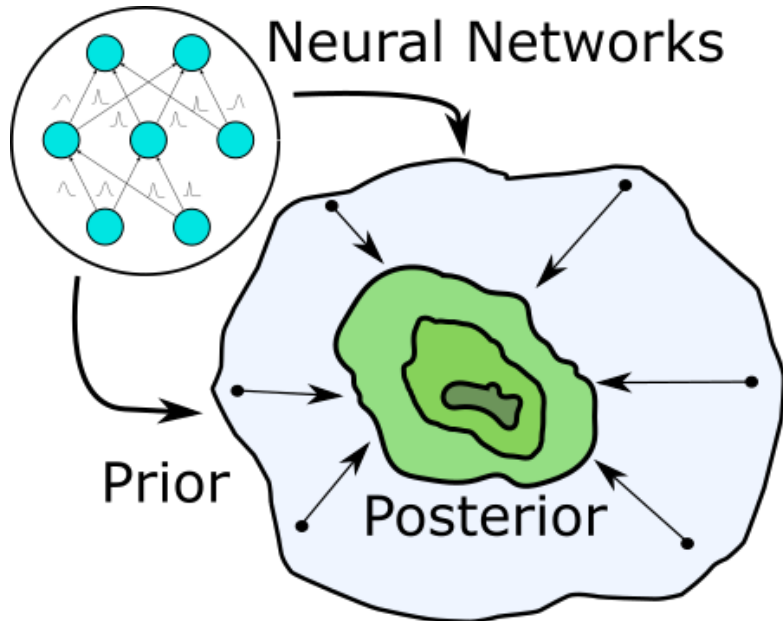
OMNO

OMNO

OMNO

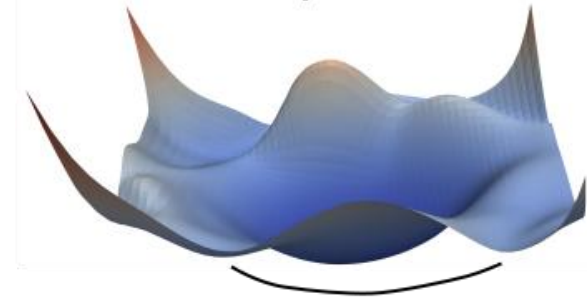
# Overview

## On Priors



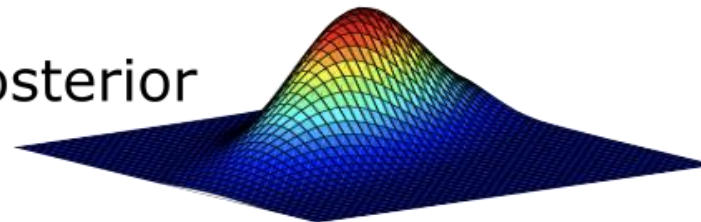
## On Posteriors

Loss landscape



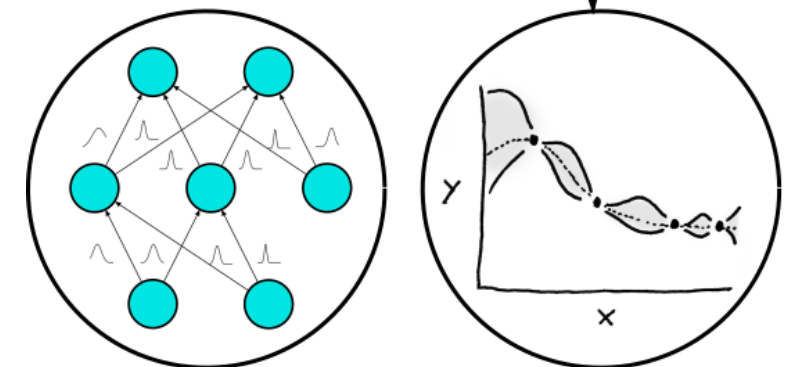
Curvature

Posterior



## On Predictions

Neural Networks

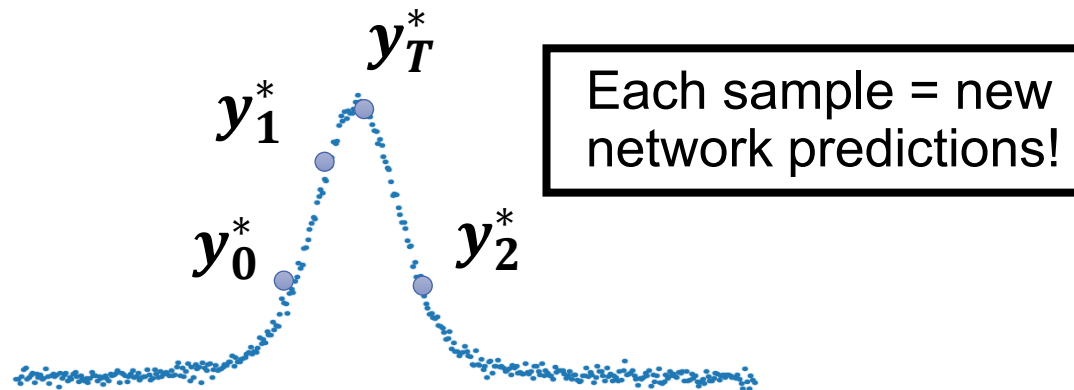


Gaussian Processes

How to efficiently predict through marginalization by avoiding any sampling?

# Challenge: Sampling is expensive

## Sampling methods



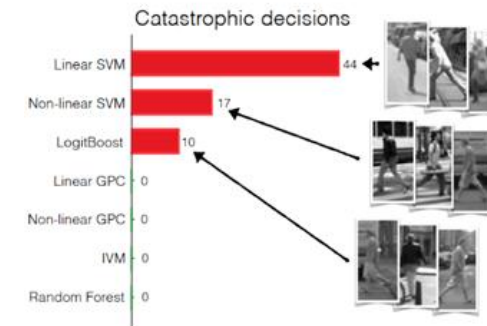
- Problem: Marginalization is intractable!

$$p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{D}) = \int p(\mathbf{y}^* | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{D}) d\boldsymbol{\theta}$$

- Solution: Monte-carlo integration.

$$\approx \frac{1}{T} \sum_{t=1}^T f_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}(\mathbf{x}^*) \quad \text{for } \boldsymbol{\theta}_t \sim p(\mathbf{w} | \mathbf{D})$$

## Sampling-free methods



(Grimmett et al, 2016 IJRR)

- Not every algorithm needs sampling.
- Gaussian Processes (GP) marginalized analytically and still golden standards.

**Idea:** Use GPs for Deep Learning?

# Idea: Theory of Neural Tangent Kernel

## Early works (1990s)

*Radford Neal, "Priors for Infinite Networks", 1995.*

- Pioneered the GP & neural network connection.
- Increasing width, single hidden layer and independent priors on neural network weights.
- Neural networks converge to GPs with Neural Tangent Kernel (NTK):  $\mathbf{K}(\mathbf{x}, \mathbf{x}) = \mathbf{J}_f^T(\mathbf{x})\mathbf{J}_f(\mathbf{x})$ .

## Recent breakthroughs

*Modern extensions of Neal 1995.*

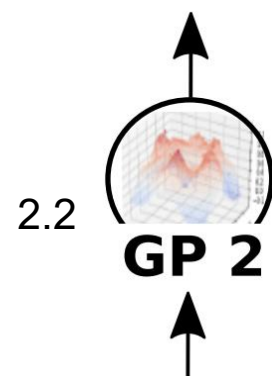
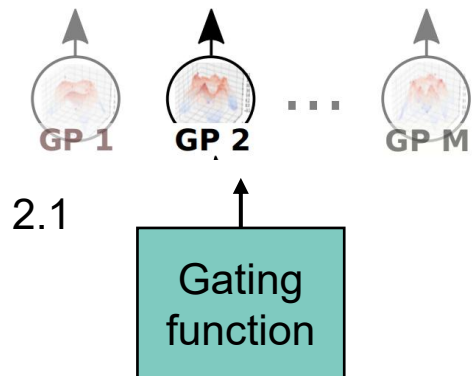
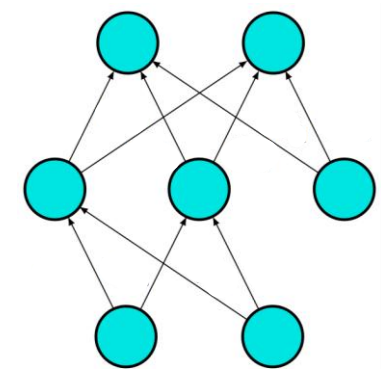
- Multiple hidden layers.  
→ J. Lee et al (2018).  
→ Matthews et al (2018).
- Convolution layers.  
→ Alonso et al (2019).
- Bayesian inference.  
→ Khan et al (2019).
- Finite width.  
→ Novak et al (2022).

**These works greatly advance the learning theory of deep learning!**

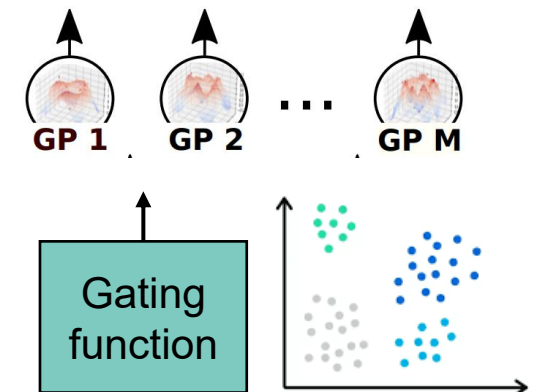
# Idea: Theory of Neural Tangent Kernel

$$p(\mathbf{y}^* | \mathbf{x}^*, D) = \mathcal{N}(\mathbf{f}_{\theta=\hat{\theta}}, \Sigma)$$

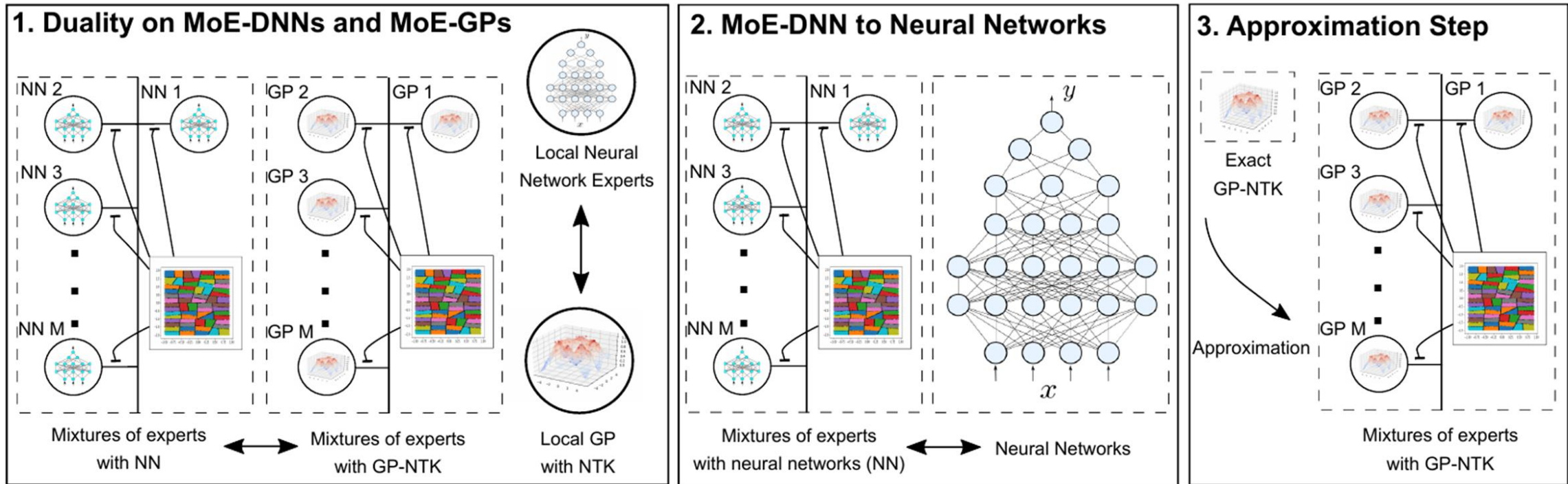
1. Compute neural network predictions:  $\mathbf{y}^* = \mathbf{f}_{\theta=\hat{\theta}}(\mathbf{x}^*)$ .
2. Compute covariance  $\Sigma$  with mixtures of GP experts:
  - 2.1. Gating function assigns  $m$ -th GP expert.
  - 2.2. Using  $m$ -th GP expert:  $\Sigma = K_{**} - K_*^T [K + \sigma_n^2 I]^{-1} K_*$ .



- + Sampling-free.
- + Winning combo.
- + Aids scaling.



# Idea: Theory of Neural Tangent Kernel



**Contribution:** theoretic foundation behind the proposed method with the NTK.

# Evaluation: Perceptive Shared Autonomy



**High uncertainty**

High uncertainty

Level 0: teleoperation

This block illustrates high uncertainty. It features a green header 'High uncertainty' and a diagram of a neural network with a red target icon overlaid. A photograph shows a person in a white shirt and blue jeans teleoperating a robotic arm. The text 'High uncertainty' and 'Level 0: teleoperation' is positioned below the diagram and photo respectively.

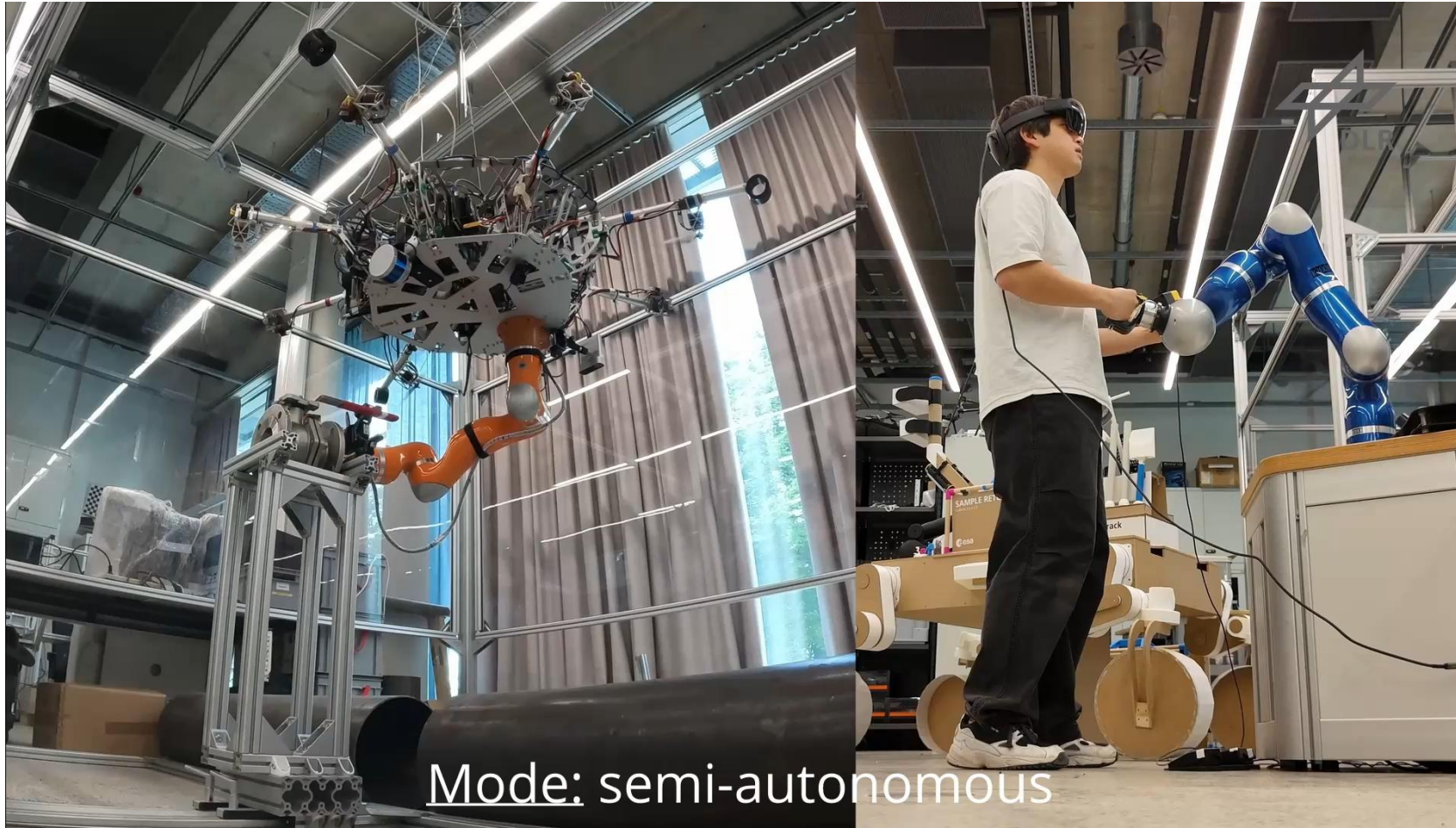
**Low uncertainty**

Low uncertainty

Level 3: assistance

This block illustrates low uncertainty. It features a green header 'Low uncertainty' and a diagram of a neural network with a green target icon overlaid. A photograph shows a person in a white shirt and blue jeans assisting with a robotic arm. The text 'Low uncertainty' and 'Level 3: assistance' is positioned below the diagram and photo respectively.

# Robust Aerial Manipulation



Set 1	Vanila-teleop	SPIRITv2
Success rate [%] ↑	86.66±0.000	<b>100.0±0.000</b>
Time [s] ↓	160.46±144.1	<b>61.86±46.03</b>
Forces [N] ↓	11.66±1.325	11.72±1.414
Torques [Nm] ↓	6.674±0.748	7.442±0.810
NASA TLX [-] ↓	11.58±3.715	<b>7.422±2.613</b>
SUS score [-] ↑	51.00±25.73	<b>71.33±16.77</b>

Shared autonomy improves the performance.

Set 2	Vanila-VF	SPIRIT
Success rate [%] ↑	40.00±0.000	<b>100.0±0.000</b>
Time [s] ↓	335.0±202.2	<b>78.33±72.77</b>
Forces [N] ↓	13.70±1.018	<b>10.86±1.586</b>
Torques [Nm] ↓	7.202±0.449	<b>6.679±0.903</b>
NASA TLX [-] ↓	14.18±3.416	<b>8.288±4.137</b>
SUS score [-] ↑	35.83±20.92	<b>65.50±24.10</b>

Uncertainty-awareness improves system reliability.

iiQKA – Robots for the People

XITASO

Consulting, Development, Implementation

Future-ready production

XITASO

Unleash the potential of the Digital Twin

automox Robotics

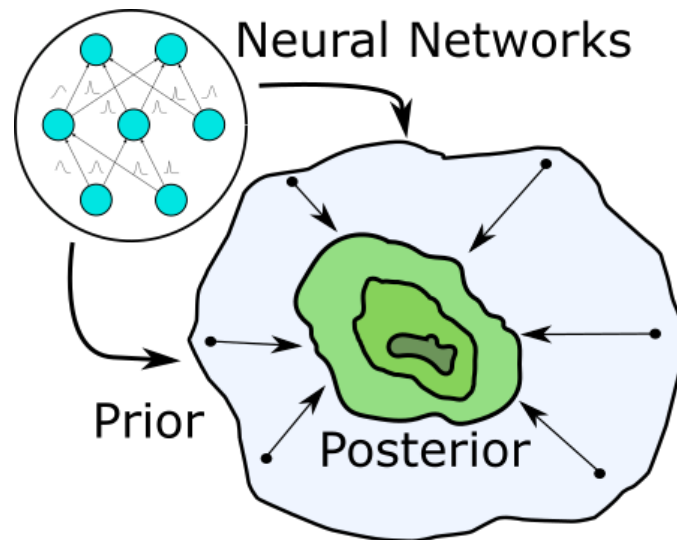
ROBOT MECHANICS FOR EVERY CONTROLLER

KUKA

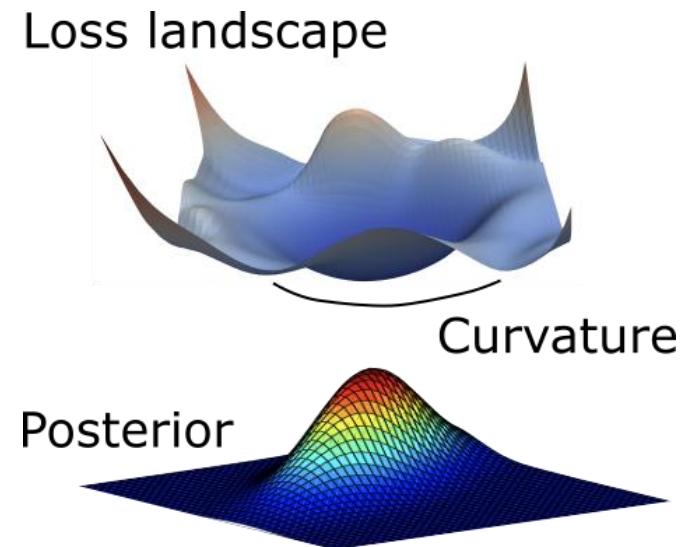
# Conclusion

AI-based robotics that reason about uncertainty – no maximum likelihood.

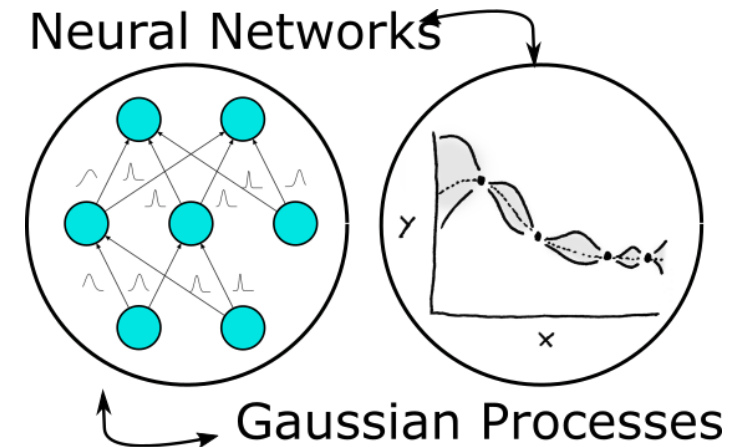
## On Priors



## On Posteriors



## On Predictions



Novel probabilistic representations suitable for applications in robotics.

# Achievements

- Tisby et al. 1989
- Mackay 1992
- Hinton and Van Camp 1993
- Neal 1995

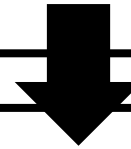
## AI winter

- Graves 2011
- Hernandez-Lobato & Adams 2015
- Gal and Ghahramani 2016
- Lakshminarayanan et al. 2017

## When I started!

### Promises:

- ⊕ Inclusion of domain knowledge via prior.
- ⊕ Quantifies uncertainty about the model.
- ⊕ Average different hypothesis about events.



### Future work:

- **Prior:** Relevant data and task to learn the prior from.
  - Foundational priors that generalize more broadly.
- **Posterior:** Missing different modes of the posterior.
  - Combination with efficient variant of deep ensemble.
- **Predictions:** Limitations in scalability of GPs.
  - Hybrid methods for sparsification.
  - E.g., partitioning + inducing points + random features.

# Dissemination

## Research



ICRA 2020, **ICML 2020**,  
CoRL 2021, RA-L 2023\*\*,  
**ICML 2023\*\***, T-FR 2024,  
RA-L2025a, RA-L2025b\*

**Thesis: 8 lead-author**

10 > projects

9 lead author / 25 total \*

## Grants and Awards



Major industrial award\*

## Media Coverage



ProRobots

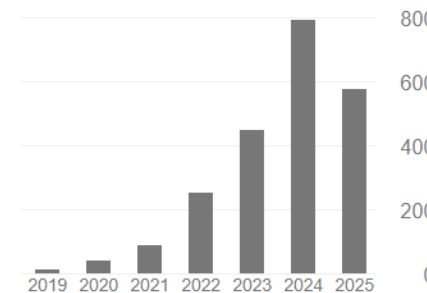


IEEE  
Video Friday



MIT review

## Academic Recognition



2400 > citations



Invited dissemination

\*Out of those 1 lead / 4 co-author are under review / 2 of 9 co-lead.

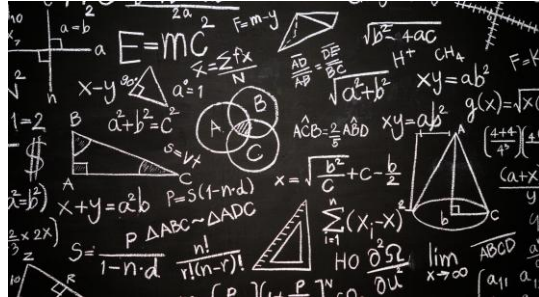
\*Finalist of KUKA Innovation Award 2023, 2024.

# Dissemination

## Theoretical contributions

ICRA2020, ICML2020,  
CoRL2021, RA-L2023\*\*,  
ICML2023\*\*, T-FR2024,  
RA-L2025, RSS2026\*

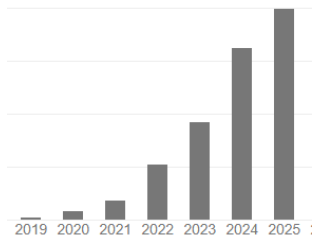
**Thesis: 8 lead-author**



9 lead author / 25 total \*

10 key lemmas

## Community Recognition



2700>citations



International dissemination



Media coverages

## Technology transfer potential



Major industrial awards\*\*\*



PI of market research project

## Multidisciplinary



5 perception tasks

7 different robots

6 outside-lab settings

\*Out of those 1 lead / 2 co-author are under review.

\*\*co-lead.

\*\*\*Finalist of KUKA Innovation Award 2023, 2024.

# Mentors and collaborators





No Sisyphus.  
No maximum  
likelihood!

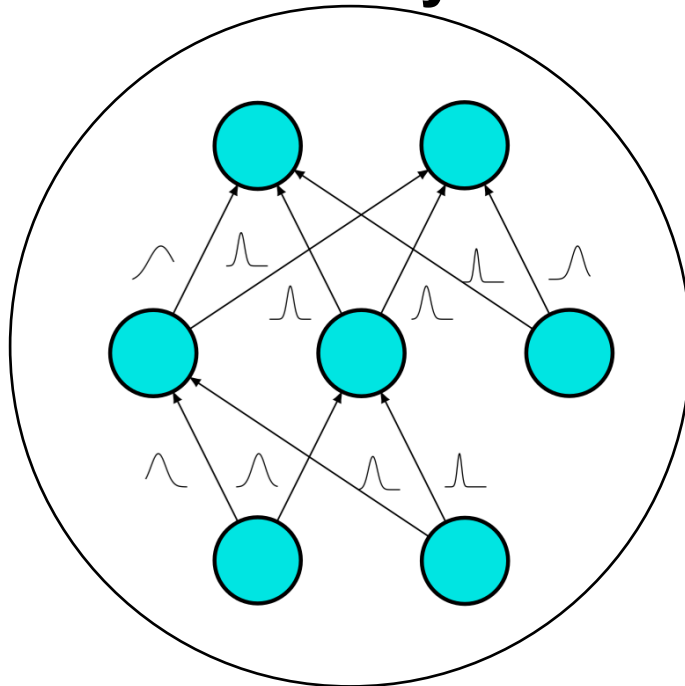
Thank you for listening!

# A: Overview

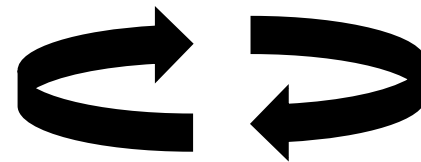
# Major novelty

## Uncertainty in Deep Learning: A Probabilistic Robotics Perspective

Theory



Practice



# Priors, Posteriors and Predictions

## Priors

## Posteriors

## Predictions

Research gap

Existing works choose uninformative prior despite the signs of prior misspecification.

Existing works simplifies the covariance matrix for scalability in lieu of expressivity.

Existing works require combining multiple predictions from the model's distribution.

Contributions

A novel method for learning scalable and structured posteriors of neural networks as information priors with generalization theory.

A novel information theoretic formulation that exploits the sparsity of inverse space of posteriors and its inference procedures.

A novel neural tangent kernel theory for sampling-free uncertainty estimation with mixtures of Gaussian Process experts.

Ramification

Domain knowledge can be incorporated for generalization and uncertainty estimation

More expressive but memory efficient approximation of posterior can be obtained

Real-time uncertainty estimation becomes possible with Gaussian Process formulation

Use cases

Excels in small data regime where the prior term dominates over the likelihood.

Excels in large data regime, large architectures but requires small memory consumption.

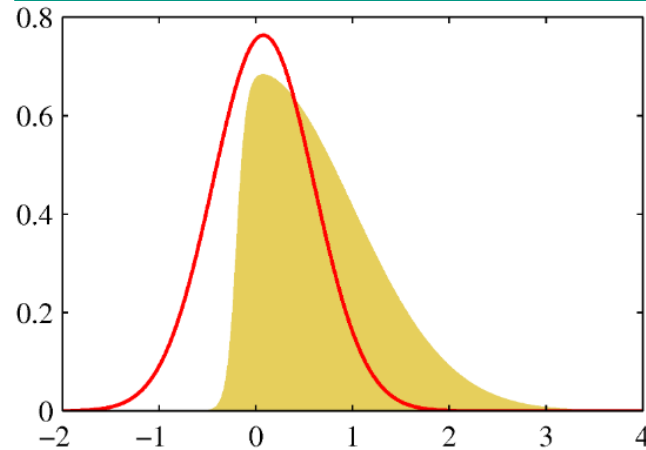
Excels in data regime upto 2 million. Permits large input but smaller output space.

# Priors, Posteriors and Predictions

	Priors	Posteriors	Predictions
Accuracy	No design guarantees on preservation of accuracy (a-priori method).	Preserves accuracy of deep learning by design (post-hoc method).	Preserves accuracy of deep learning by design (post-hoc method).
Memory efficiency	Deployable on embedded GPUs. Memory grows with tasks	Deployable on embedded GPUs. Most memory efficient.	Deployable on embedded GPUs. Memory grows with data.
Run-time efficiency	Requires sampling for predictions and efficient only when parallelized.	Requires sampling for predictions and efficient only when parallelized.	Does not require sampling and thus efficient without parallelization.
Scalability	ImageNet-1K data-set within 24 hours on NVIDIA 1080 GPU.	ImageNet-1K data-set within 24 hours on NVIDIA 1080 GPU.	Two million data points within 15 hours on NVIDIA 1080 GPU.
Inference scheme	Informative prior Laplace Approximation Sampling	Uninformative prior Laplace Approximation Sampling	Uninformative prior Laplace Approximation Gaussian Processes

# Bayesian Deep Learning

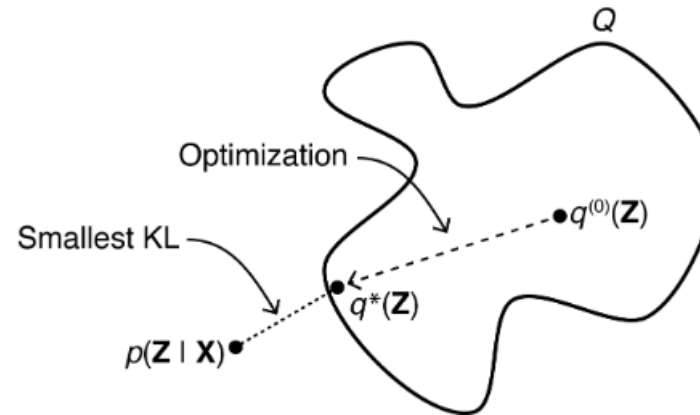
## Laplace Approximation



Hessian approach

- + Computationally scalable.
- + **Ease of implementation/training-free.**
- + Competitive performance.
- Unimodal posterior.

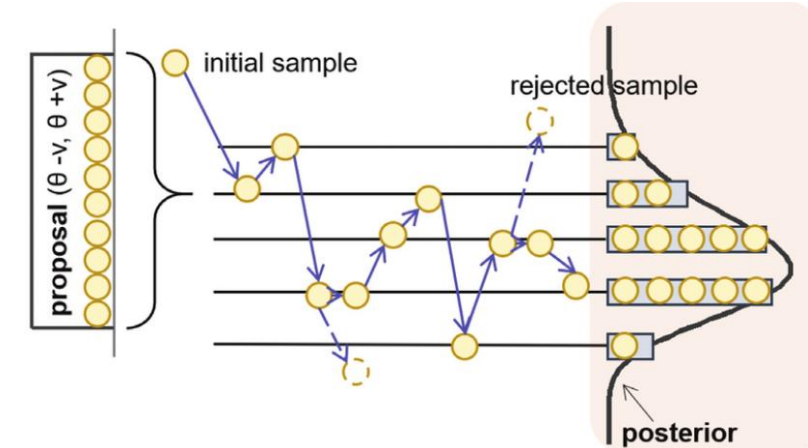
## Variational Inference



Optimization approach

- + More scalable than sampling.
- + Theoretic guarantees.
- **Optimization challenges.**

## MCMC Sampling



Sampling approach

- + Best uncertainty estimates.
- + Theoretic guarantees.
- **Do not scale to large data and architectures.**



More and more agreement that LA is practical alternative.  
 Post-hoc methods – rely on standard neural network training.

# Why not Expectation Propagation?

## Recap

Given posterior:

$$p(\boldsymbol{\theta} \mid \mathcal{D}) = \frac{1}{p(\mathcal{D})} \prod_{i=1}^M f_i(\boldsymbol{\theta})$$

We want to obtain:

$$q(\boldsymbol{\theta}) = \frac{1}{Z} \prod_{i=1}^M \tilde{f}_i(\boldsymbol{\theta})$$

By minimizing:

$$\text{KL} \left( \frac{1}{p(\mathcal{D})} \prod_{i=1}^M f_i(\boldsymbol{\theta}) \parallel \frac{1}{Z} \prod_{i=1}^M \tilde{f}_i(\boldsymbol{\theta}) \right)$$

## Application

We can formulate BNN posterior as:

$$p(W \mid D) \propto p(W) \prod_n p(y_n \mid x_n, W)$$

EP algorithm (roughly):

1. Initialize all factors
2. Pick a factor  $i$  to update
3. Compute tilted distribution
4. Approximate the factor by matching the moments with tilted distribution.
5. Get new  $q$ .
6. Repeat until convergence.

## Why not?

1. Computing moment is not tractable.
2. Non-Gaussian likelihood of BNN with Gaussian likelihood.
3. Many iterations.
4. **Convergence issues.**

Assumed density filtering (PBP) is better but:

1. Accuracy is limited.
  1. Update once per factor.
  2. Mean-field approximation.
2. No factor revisiting – too fast shrinkage.
3. Ignores correlation between data points.

Easier to train neural networks with stochastic gradient descents.

# Data and Model Uncertainty

## Aleatoric Uncertainty

- Uncertainty that comes from inherent noise in the data itself. **Irreducible (always there)**.
- Source: intrinsic randomness, ambiguity in labelling and label noises.

## Epistemic Uncertainty

- Uncertainty that comes from lack of knowledge in the model. Reduced with more data.
- Source: limited or sparse training data, underrepresented regions in the input.

Model (epistemic) uncertainty is often more important than data uncertainty because it reveals when the model encounters **unknown or unexpected inputs**, allowing safer decisions and targeted learning.

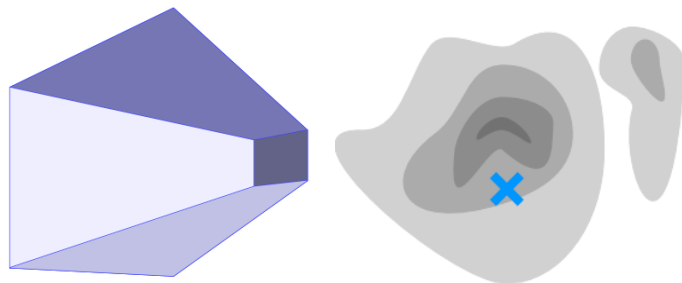
# Uncertainty in Deep Learning

## Conformal Predictions

Prediction sets guaranteed to contain the true outcome with a chosen probability. Conformal prediction instead produces a **set of plausible predictions** whose size reflects uncertainty. If the model is **uncertain**, the set is **larger** (e.g., {"cat", "fox", "dog"}).

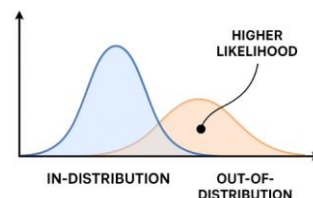
- + Easy to use / model agnostic.
- + Guaranteed coverage.
- Need for calibration data.
- Doesn't provide uncertainty.

## OOD Detectors

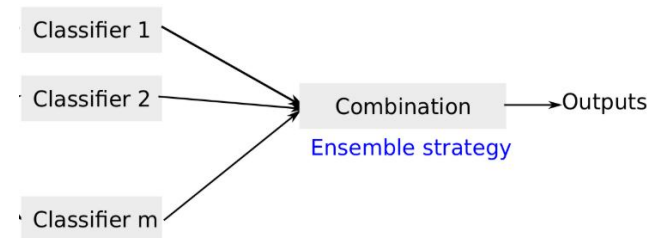


Latent space + generative model

- + Easy to use / model agnostic.
- + Fast at run-time.
- Doesn't provide uncertainty.
- Need good latent space.
- Likelihood paradox.



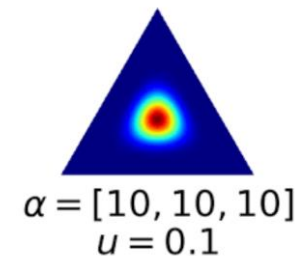
## Deep Ensembles



- + Easy to use / model agnostic.
- + Superiority to MC-dropout.
- Computational costs.

**Note:** deep ensemble can be combined with Bayesian models. \* state-of-the-art.

## Evidential Learning



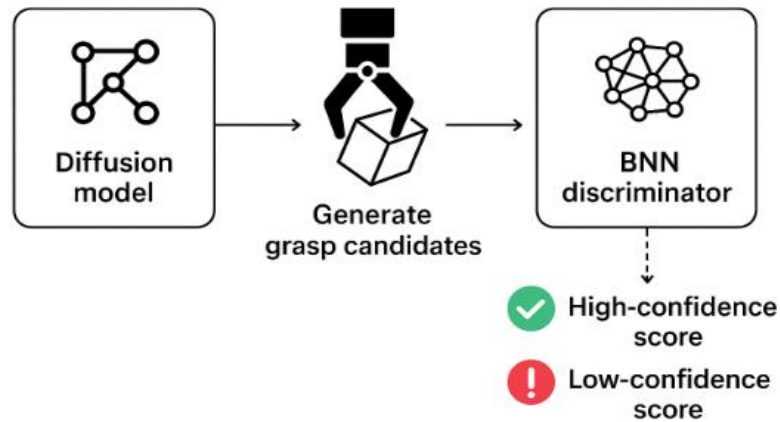
Network outputs parameters of dirichlet or normal-inverse-gamma distribution.

- + Fast at run-time.
- + Superiority to MC-dropout.
- Need to optimize / loss change.
- Need for out-of-distribution data.
- Indirect model uncertainty.

Use cases  
 LLM – CP.  
 YOLO – OOD detection.  
 Semantic segmentation – Evidential  
 Offline – DE + BNN

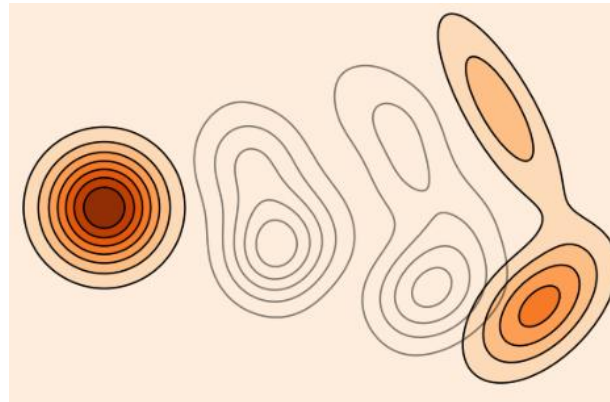
# Uncertainty & GenAI

## GenAI + BNN



- Diffusion model generate grasp candidates.
- A **discriminator** evaluates the quality of candidates.
- Likelihood  $\neq$  success probability (more to encode).

## GenAI in Bayesian inference



- Diffusion model to generate samples from posterior.
- Likelihood-Free Bayesian inference to simulate data.
- **Flow models to capture multi-modal posterior.**

## Visual foundational model

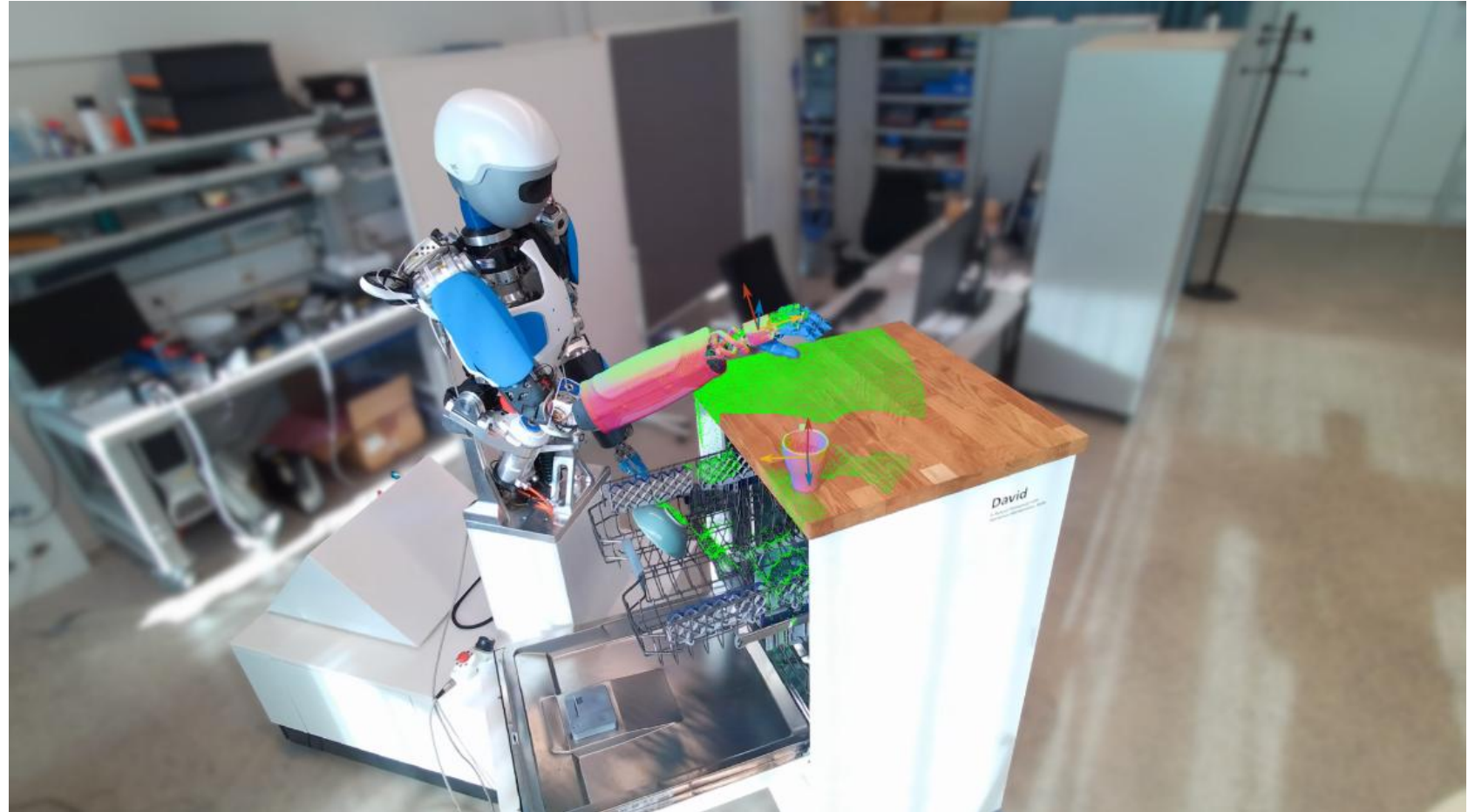


- Visual foundational models are discriminative models, trained like to generate.
- Masked Autoencoder (MAE) and Segment-anything is a discriminative model.

# Rigorous Perception

## Rigorous keywords:

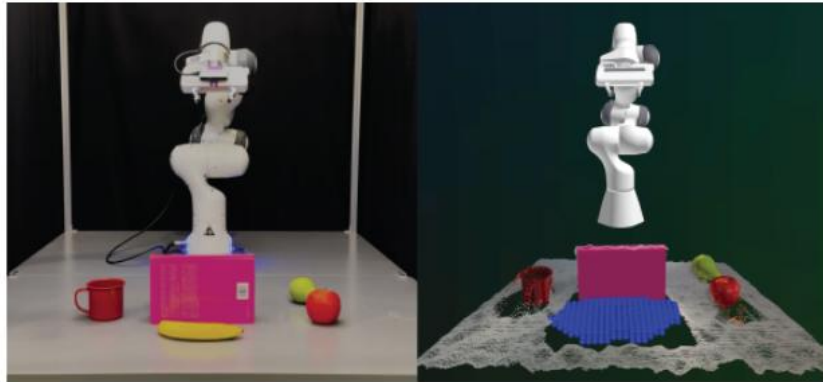
- Robust
- Interactive
- Interpretable
- Generalizable
- Omnimodal
- Resource-efficient
- Open
- Uncertainty-aware
- Semantic



# B: Extensions

# Task and Motion Planning 1

## ■ Partially Observable Task and Motion Planning with Uncertainty and Risk Awareness (RSS 2024)



- If high uncertainty → gather information first.
- If low uncertainty → act directly.

```

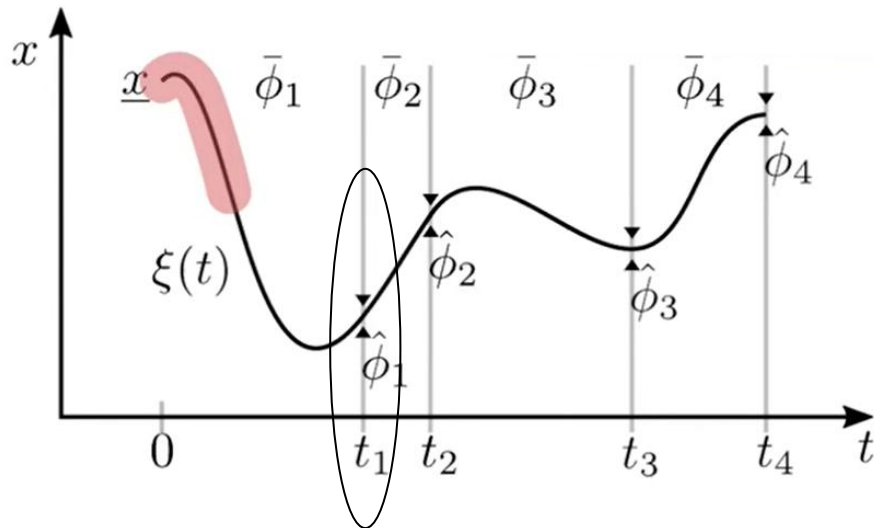
look_behind(box1):
  preconditions: can_see(box1)
  effects: reveal(object_location)
  uncertain_effects: {object_found, object_not_found}
  
```

Used Bayes3D for robot perception – 6D pose.



# Task and Motion Planning 2

- Deep Visual Constraints: Neural Implicit Models for Manipulation Planning from Visual Input (RAL2021)



- Task and motion planning with logic geometric programming.

$$\begin{aligned} \min_{x:[0,t_K] \rightarrow x, t_1:t_K} \quad & t_K + \alpha \int_0^{t_K} c(\bar{x}(t)) dt \\ \text{s.t.} \quad & \bar{x}(0) = \bar{x}_0, \dot{\bar{x}}(t_K) = 0, \\ & \forall_k : 0 < t_k < t_{k+1}, \\ & \forall_k : \hat{\phi}_k(x(t_k)) \leq 0, \\ & \forall_{t \in [t_{k-1}, t_k]} : \bar{\phi}_k(\bar{x}(t)) \leq 0. \end{aligned}$$

Optimal control with constraints.

- $h = \phi(q; V)$ , where  $q$  is robot pose (or joint state) relative to object, and  $V$  the image views. Zero when feasible.
- This mapping is learned – Unet + MLP type architecture.

**Rough idea:** estimate uncertainty for constraints of neural networks.

- Chance constrained formulation:  $\mu(x, u) + k_\alpha \sigma(x, u) \leq 0$
- Probabilistic cost as soft constraints:  $J = \sum_{t=0}^{H-1} \ell(x_t, u_t) + \lambda [\mu(x_t, u_t) + k_\alpha \sigma(x_t, u_t)]_+^2$

The probability of violating the constraint is bounded.

Optimization penalizes also the risky actions.

# C: Important Details

# Uncertainty Evaluation

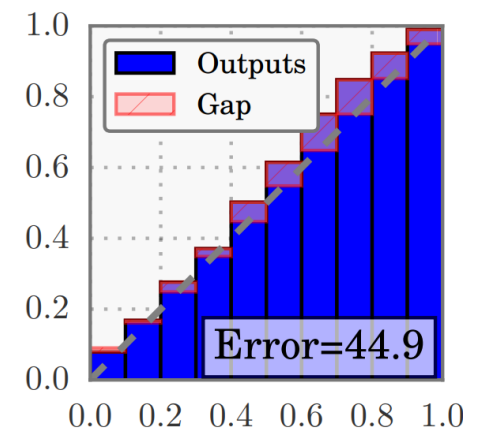
<b>Tasks</b>	Classification <sup>[1-3,5,7,10-13,16-19,21-25]</sup> & Regression <sup>[1-11,13,14,16,21,23]</sup> (incl. calibration, open set recognition, distribution shifts)
<b>Downstream</b>	Reinforcement learning <sup>[1,2,8,22]</sup> , Active learning <sup>[4,7,13,18,22,23]</sup> , Adversarial attacks <sup>[3,5,23]</sup> , Bayesian optimization <sup>[8]</sup> , Language modeling <sup>[11,17]</sup> , Continual and transfer learning <sup>[12,22,24,25]</sup> , Few-shot learning <sup>[18,25]</sup> , Generative model <sup>[19]</sup> , & Pruning <sup>[15]</sup>
<b>Metrics</b>	Test LL or NLL <sup>[1,4,6-13,16-18,20-26]</sup> , ECE or Reliability diagram <sup>[7,11-13,17,18,20,22-26]</sup> , Empirical CDF for entropy <sup>[3,5,7,12,18,21,23]</sup> , AUROC <sup>[12,17-19,20,23-25]</sup> , and Brier score <sup>[16,17,18,23]</sup> & Latency <sup>[20,21,23]</sup>
<b>Dataset</b>	Toy example <sup>[1-10,14,16,19,21]</sup> , PLEX <sup>[18,25]</sup> , CO2 <sup>[1]</sup> , Sarcos & KUKA <sup>[21]</sup> , fashionMNIST <sup>[19,24,26]</sup> , UCI <sup>[1,2,4,6-11,13,16,23,26]</sup> , ImageNet <sup>[7,11,12,16-18,23]</sup> , 20newsgroup <sup>[17]</sup> , Criteo Display Advertising Challenge <sup>[17]</sup> , RETINA <sup>[18]</sup> , CLINC OOS intent detection <sup>[20,23]</sup> , CIFAR10 <sup>[3,7,12,13,16-20,22-24]</sup> , SHVN <sup>[7,16,19,20,24]</sup> , CIFAR100 <sup>[5,10,11,12,18,20,23,24,26]</sup> , Wilds <sup>[22]</sup> , Diabetic Retinopathy Detection <sup>[23]</sup> , Wikipedia Toxicity <sup>[23]</sup> , MNIST <sup>[1-3,5,7,16,17,19,22,25,26]</sup> , notMNIST <sup>[3,5,7,16]</sup> , YCB <sup>[25]</sup> & GPT2 <sup>[26]</sup> , TinyStories <sup>[26]</sup>

## Proper score rule

A scoring rule SSS is **proper** if the expected score is **maximized** when the predicted probabilities match the true probabilities.

$$\mathbb{E}_{Y \sim p}[S(p, Y)] \geq \mathbb{E}_{Y \sim p}[S(q, Y)] \quad \forall q$$

## Calibration



# Laplace Approximation

Neural Networks:  $f_{\theta} : \mathbb{R}^d \mapsto \mathbb{R}^k$ ,  $f_{\theta}(\mathbf{x}) = \varphi(\mathbf{W}_L \varphi(\mathbf{W}_{L-1}(\dots \varphi(\mathbf{W}_1 \mathbf{x}))))$

$\theta = [\mathbf{W}_1, \dots, \mathbf{W}_L]$

Prior:  $p(\theta|\alpha) = \mathcal{N}(\theta|\mathbf{0}, \alpha^{-1}\mathbf{I})$

Likelihood:  $p(\mathcal{D}|\theta, \beta) = \prod_{N=1}^{n=1} \mathcal{N}(\mathbf{y}_i|f_{\theta}(\mathbf{x}_i)), \beta^{-1}$

Posterior:  $p(\theta|\mathcal{D}, \alpha, \beta) \propto p(\theta|\alpha)p(\mathcal{D}|\theta, \beta)$

Laplace Approximation – Posterior:

$$\ln p(\theta|\mathcal{D}) \approx \ln p(\hat{\theta}|\mathcal{D}) - \frac{1}{2}(\theta - \hat{\theta})^T \mathbf{H}(\theta - \hat{\theta})$$

Taylor series around the mode

$$\begin{aligned} & \frac{\partial^2}{\partial \theta^2} \ln p(\theta|\mathcal{D}) \Big|_{\theta=\hat{\theta}} \\ &= \frac{\partial^2}{\partial \theta^2} \ln p(\mathcal{D}|f_{\theta}) \Big|_{\theta=\hat{\theta}} + \frac{\partial^2}{\partial \theta^2} \ln p(\theta) \Big|_{\theta=\hat{\theta}} \\ &=: \mathbf{H}_{\text{likelihood}} + \mathbf{H}_{\text{prior}}. \end{aligned}$$

Rewriting Hessian term

$$p(\theta|\mathcal{D}) \approx \mathcal{N}(\theta|\hat{\theta}, \mathbf{H}^{-1})$$

Exponential on both sides and reverse engineering the density

# Laplace Approximation

Same setting as before. Then, consider predictive distribution:

$$p(\mathbf{y}^*, |\mathbf{x}^*, \mathcal{D}) = \int p(\mathbf{y}^* | \mathbf{x}^*, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta}$$

$$\mathbb{E}_{\boldsymbol{\theta} \sim p(\boldsymbol{\theta} | \mathcal{D})} [p(\mathbf{y}^* | \mathbf{x}^*, \boldsymbol{\theta})] \approx \frac{1}{T} \sum_{t=1}^T p(\mathbf{y}^* | \mathbf{x}^*, \boldsymbol{\theta}^{(t)}) \quad \text{for } \boldsymbol{\theta}^{(t)} \sim p(\boldsymbol{\theta} | \mathcal{D})$$

Taylor series approximation around the output.

$$f_{\boldsymbol{\theta}}(\mathbf{x}) \approx f_{\hat{\boldsymbol{\theta}}}(\mathbf{x}) + \mathbf{J}_f^T(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})$$

Rules of Gaussian – closure under marginalization.

$$p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}, \beta) \approx \mathcal{N}(\mathbf{y} | f_{\hat{\boldsymbol{\theta}}}(\mathbf{x}) + \mathbf{J}_f^T(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}), \beta^{-1})$$

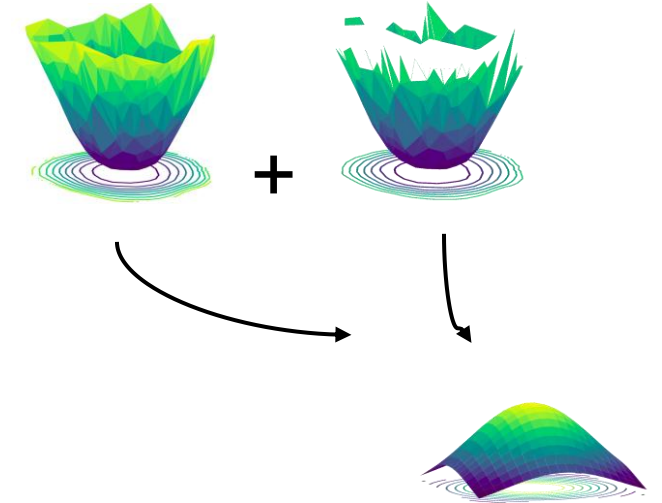
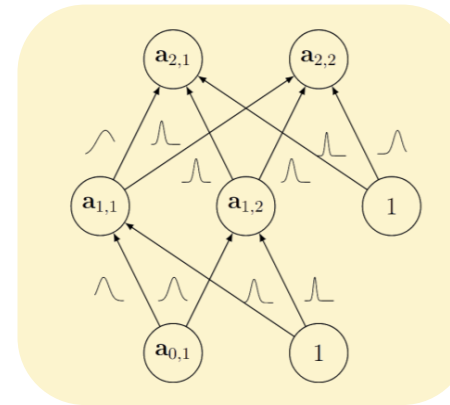
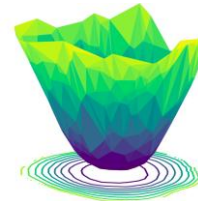
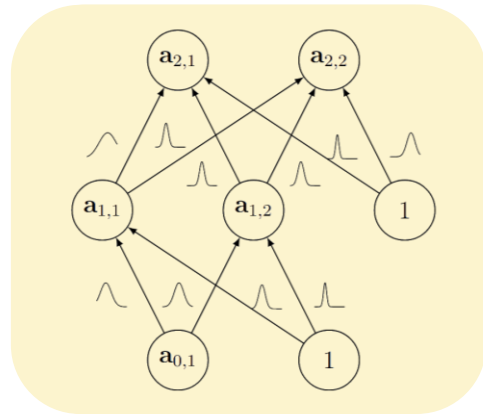
$$p(\mathbf{y} | \mathbf{x}, \mathcal{D}, \alpha, \beta) = \mathcal{N}(\mathbf{y} | f_{\hat{\boldsymbol{\theta}}}(\mathbf{x}), \beta^{-1} + \mathbf{J}_f^T \mathbf{H}^{-1} \mathbf{J}_f)$$

# Assumptions

#	Assumption	Explanation
1	<b>Posterior is unimodal (one main peak)</b>	The approximation is centered at a single mode (the MAP). If there are multiple modes, the Gaussian can't capture the shape well.
2	<b>Log-posterior is approximately quadratic near the mode</b>	The method assumes the second-order Taylor expansion is accurate — i.e. higher-order terms are negligible.
3	<b>Posterior is well-peaked (not flat or heavy-tailed)</b>	The curvature (Hessian) around the mode should dominate; otherwise, the Gaussian spread misrepresents uncertainty.
4	<b>Negative Hessian of the log-posterior is positive definite</b>	Ensures the local curvature corresponds to a maximum (not a saddle point), and defines a valid covariance matrix.
5	<b>Sufficient sample size / asymptotic normality</b>	When data are large, the posterior tends to become Gaussian around the true parameter (Bernstein–von Mises theorem). Laplace works best in that regime.
6	<b>Continuity and differentiability</b>	The log posterior must be twice continuously differentiable with respect to parameters (needed for the Taylor expansion).

# Prior

# Posterior as prior



At task 0 (broad data-set and large architecture)

- Maximum a posteriori estimation

$$\hat{\boldsymbol{\theta}}^{(0)} \in \operatorname{argmax} p(\boldsymbol{\theta}^{(0)} | \mathbf{D}^{(0)}).$$

- Laplace Approximation

$$p(\boldsymbol{\theta}^{(0)} | \mathbf{D}^{(0)}) \approx \mathcal{N}(\boldsymbol{\theta}^{(0)} | \hat{\boldsymbol{\theta}}^{(0)}, (\mathbf{F}^{(0)})^{-1}).$$

- Kronecker-factorized information matrix

$$\mathbf{F}^{(0)} = \mathbf{F}_{\text{likelihood}}^{(0)} + \mathbf{F}_{\text{prior}}^{(0)} = \mathbf{L}^{(0)} \otimes \mathbf{R}^{(0)} + \gamma \mathbf{I}.$$

At task 1 (our robotic data of interest):

- Prior learned on task 0

$$\pi(\boldsymbol{\theta}^{(1)}) = \mathcal{N}(\hat{\boldsymbol{\theta}}^{(0)}, \mathbf{F}^{(0)-1}).$$

- Posterior update on task 1

$$p(\boldsymbol{\theta}^{(1)} | \mathbf{D}) = \mathcal{N}(\hat{\boldsymbol{\theta}}^{(1)}, (\mathbf{F}^{(1)} + \mathbf{F}^{(0)})^{-1}).$$

- Kronecker-factorized information matrix

$$\mathbf{F}^{(1)} = \mathbf{L}^{(1)} \otimes \mathbf{R}^{(1)} + \mathbf{L}^{(0)} \otimes \mathbf{R}^{(0)} + \gamma \mathbf{I}.$$

# Sums-of-Kronecker product computations

## Problem: sums-of-kronecker product as kronecker product

- Sums-of-kronecker product is not a kronecker product:

$$F = L^{(1)} \otimes R^{(1)} + L^{(0)} \otimes R^{(0)} + \gamma I \neq \hat{L}^{(1)} \otimes \hat{R}^{(1)}$$

$$F = L \otimes R = \begin{pmatrix} l_{11}R & \cdots & l_{1n}R \\ \vdots & \ddots & \vdots \\ l_{m1}R & \cdots & l_{mn}R \end{pmatrix}$$

- But people still need Kronecker factorization for efficiency:

Storage:  $F = L \otimes R, \quad F \in \mathbb{R}^{pm \times qn}, L \in \mathbb{R}^{m \times n}, R \in \mathbb{R}^{p \times q}$

Inversion:  $(L \otimes R)^{-1} = L^{-1} \otimes R^{-1}$

- Hence, prior works make assumptions that:  $L \otimes R + \sqrt{\tau}I_{m \times n} \otimes \sqrt{\tau}I_{p \times q} = (L + \sqrt{\tau}I_{m \times n}) \otimes (R + \sqrt{\tau}I_{p \times q}) = \hat{L} \otimes \hat{R}$

## Idea: power method converges to rank-1 optimal solution

- Sums-of-kronecker product as an optimization problem:

$$\hat{L}, \hat{R} \in \operatorname{argmin}_{L,R} \left| \sum_{k=1}^K L^k \otimes R^k - L \otimes R \right|_F$$

- Reorder the elements leads to rank 1 approximation problem:

$$\left| \sum_{k=1}^K L^k \otimes R^k - L \otimes R \right|_F = \left| \sum_{k=1}^K \operatorname{vec}(L^k) \operatorname{vec}(R^k)^T - \operatorname{vec}(L) \operatorname{vec}(R)^T \right|_F$$

- Power method converges to rank-1 optimal solution.

$$R^{(n)} \leftarrow \frac{\sum_{k=1}^K \langle L^k, L^{(n-1)} \rangle_F R^k}{\left| \sum_{k=1}^K \langle L^k, L^{(n-1)} \rangle_F R^k \right|_F}$$

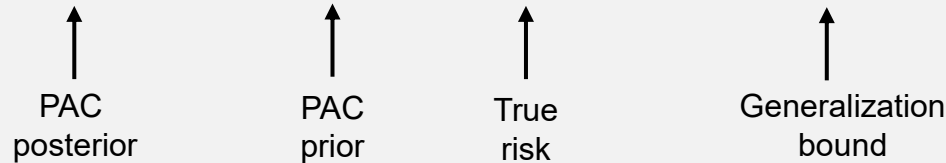
↑ Proposed power iterations

$$L^{(n)} \leftarrow \frac{\sum_{k=1}^K \langle R^k, R^{(n-1)} \rangle_F L^k}{\left| \sum_{k=1}^K \langle R^k, R^{(n-1)} \rangle_F L^k \right|_F}$$

# Priors with generalization guarantees

PAC-Bayes Theorem:

$$P_{D \sim p^N} \left( \forall \rho \in \mathcal{M}, \rho \ll \pi : \mathbb{E}_{g \sim \rho} [L_P^l(g)] \leq \delta(\rho, \pi, D, \varepsilon) \right) \geq 1 - \varepsilon$$



For example, McAllester bound states:  $\delta(\rho, \pi, D, \varepsilon) = \mathbb{E}_{g \sim \rho} [\hat{L}_D^l(g)] + \sqrt{\frac{\text{KL}(\rho || \pi) + \ln \frac{2\sqrt{N}}{\varepsilon}}{2N}}$

Formulation (Humt et al., 2020)

$$\begin{aligned}
 & p(\boldsymbol{\theta}^{(1)} | \mathbf{D}^{(1)}) \\
 &= \mathcal{N} \left( \hat{\boldsymbol{\theta}}^{(1)}, (\mathbf{F}^{(1)} + \mathbf{F}^{(0)})^{-1} \right) \\
 &\rightarrow \mathcal{N} \left( \hat{\boldsymbol{\theta}}^{(1)}, \boldsymbol{\tau} \left( \boldsymbol{\beta} \mathbf{F}^{(1)} + \boldsymbol{\alpha} (\mathbf{F}^{(0)}) \right)^{-1} \right)
 \end{aligned}$$

Find  $\boldsymbol{\tau}$ ,  $\boldsymbol{\alpha}$ ,  $\boldsymbol{\beta}$  per layer?

Idea: PAC-Bayes for probabilistic inference (P. Germain et al, 2016)

$$\pi \rightarrow \pi(\boldsymbol{\theta}^{(1)}) \quad \rho \rightarrow p(\boldsymbol{\theta}^{(1)} | \mathbf{D}^{(1)})$$

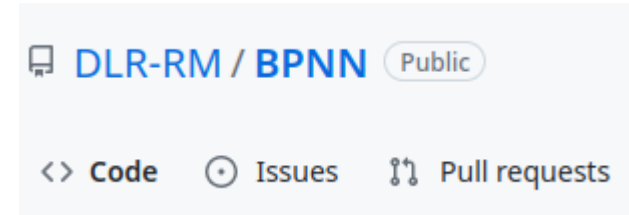
Propose an approximate upper bound for empirical risk

$$\begin{aligned}
 \mathbb{E}_{g \sim \rho} [\hat{L}_D^l(g)] &\leq \mathbb{E}_{\boldsymbol{\theta} \sim \rho} \left[ \frac{1}{N} \sum_{i=1}^N -\frac{\ln p(\mathbf{y}_i | \mathbf{x}_i, f_{\boldsymbol{\theta}})}{\ln 2} \right] \\
 &\approx \frac{-\ln p(\mathbf{D} | f_{\hat{\boldsymbol{\theta}}}) + \frac{1}{2} \sum_{l \in \{L\}} \boldsymbol{\tau} \text{tr} \left( \mathbf{F}^{(1)} + (\boldsymbol{\beta} \mathbf{F}^{(1)} + \boldsymbol{\alpha} \mathbf{F}^{(0)})^{-1} \right)}{N \ln 2}
 \end{aligned}$$

# Algorithmic overview

## Approximate Bayesian inference

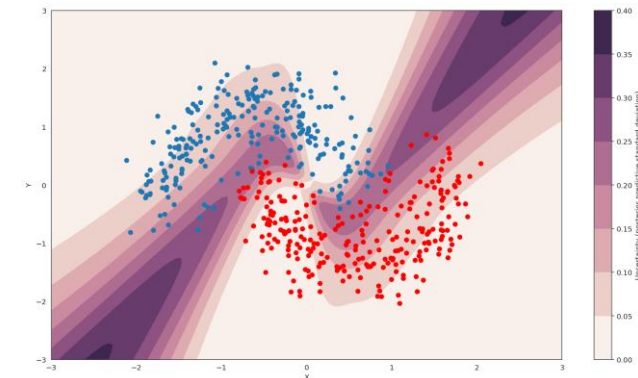
1. Download your pre-trained posterior as prior.
2. Fine tune your network using that prior.
3. Apply Kronecker-factorized Laplace Approximation.
  - 3.1 Hyperparameter tuning with differentiable PAC-Bayes.
  - 3.2 Sums-of-Kronecker product computations.



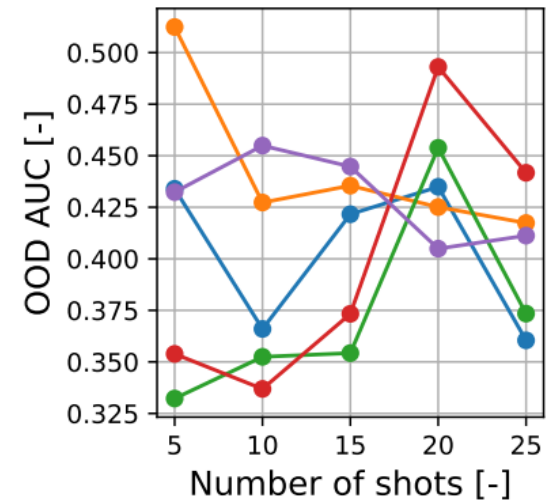
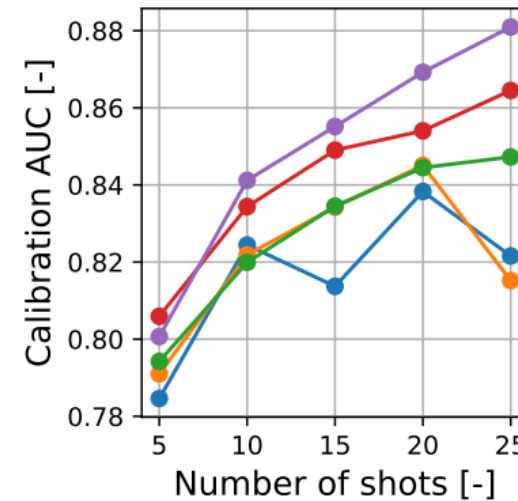
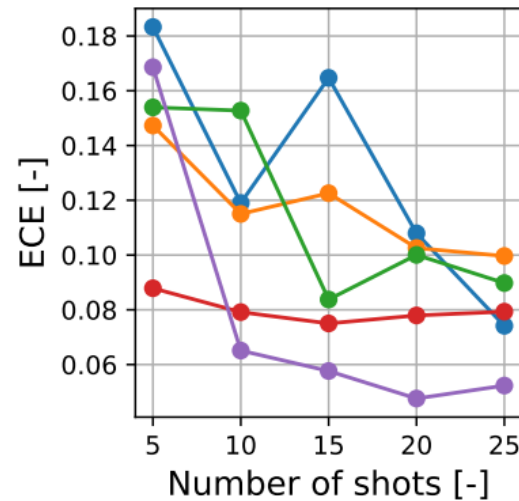
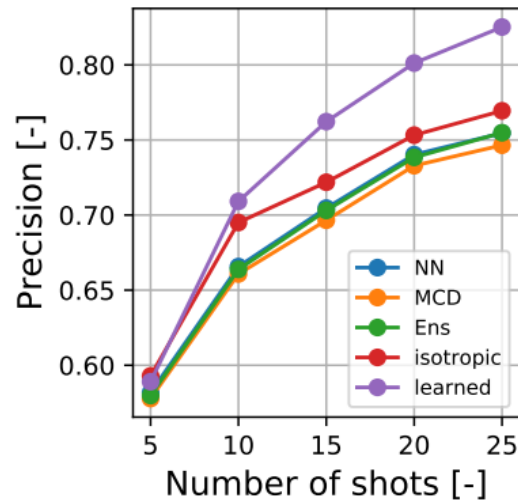
$$p(\boldsymbol{\theta}^{(1)} | \mathbf{D}^{(1)}) = \mathcal{N}(\hat{\boldsymbol{\theta}}^{(1)}, (\hat{\mathbf{L}}^{(1)} \otimes \hat{\mathbf{R}}^{(1)})^{-1}).$$

## Predictive distributions

$$p(\mathbf{y}^* | \mathbf{D}^{(1)}, \mathbf{x}^*) = \int p(\mathbf{y}^* | \mathbf{x}^*, \boldsymbol{\theta}^{(1)}) p(\boldsymbol{\theta}^{(1)} | \mathbf{D}^{(1)}) d\boldsymbol{\theta}^{(1)}$$



# Few-shot classification



## Experiment setting

- Standardized benchmark from uncertainty baselines (Z. Nado et al, 2021).
- Average results over 8 data-sets.
- Few-shot image classification.

## Key observations

- A working example learning prior from posteriors of Resnet-50 on ImageNet.
- Better generalization and uncertainty.
- Competitive on OOD detection.

# Posterior

# Inference in information form

- Approximate inference using Laplace Approximation:

$$\mathbf{p}(\boldsymbol{\theta}|\mathbf{D}) \sim \mathcal{N}(\boldsymbol{\theta}_{\text{map}}, \mathbf{H}^{-1}) \quad \text{or} \quad \sim \mathcal{N}^{-1}(\boldsymbol{\eta}_{\text{map}}, \mathbf{H})^*$$

- Employ EFB [George et al 2018] to estimate the Hessian  $\mathbf{H}$ :

$$\mathbf{H} \approx \mathbf{I}_{\text{efb}} = (\mathbf{U}_A \otimes \mathbf{U}_G) \boldsymbol{\Lambda} (\mathbf{U}_A \otimes \mathbf{U}_G)^T$$

EFB Fisher Information matrix

$$\mathbf{U}_A \text{ is an eigenvector of } \mathbf{A} = \mathbb{E}[\mathbf{a}\mathbf{a}^T]$$

Forward pass

$$\mathbf{U}_G \text{ is an eigenvector of } \mathbf{G} = \mathbb{E}[\mathbf{g}\mathbf{g}^T]$$

Backward pass

$$\Lambda_{ii} = \mathbb{E}\left[\left((\mathbf{U}_A \otimes \mathbf{U}_G)^T \boldsymbol{\delta}\boldsymbol{\theta}\right)_i^2\right]$$

Eigenvalue updates

# Inference in information form

- Diagonal elements of true Information matrix is known and easy to compute!

$$\mathbf{I} = \mathbb{E}[\boldsymbol{\delta}\boldsymbol{\theta}\boldsymbol{\delta}\boldsymbol{\theta}^T] \text{ by definition, and } \mathbf{I}_{ii} = \mathbb{E}[\boldsymbol{\delta}\boldsymbol{\theta}_i^2] \quad \forall i \quad \text{Not true for the covariance matrix}$$

- Resulting Kronecker-factored Eigen-decomposition plus diagonal structured information matrix:

$$\mathbf{I}_{\text{inf}} = (\mathbf{U}_A \otimes \mathbf{U}_G)\boldsymbol{\Lambda}(\mathbf{U}_A \otimes \mathbf{U}_G)^T + \mathbf{D}^* \quad \text{Exact on the diagonals}$$

- This step brings a theoretical guarantee on improvements:

**Lemma 1: theoretical guarantees regardless of the chosen data-set and architecture**

Let  $\mathbf{I}$  be the real information matrix, and let  $\mathbf{I}_{\text{inf}}$  and  $\mathbf{I}_{\text{efb}}$  be the INF and EFB estimates of it respectively.

Then, it is guaranteed to have  $\|\mathbf{I} - \mathbf{I}_{\text{efb}}\|_F \geq \|\mathbf{I} - \mathbf{I}_{\text{inf}}\|_F$

# Low rank sampling computations

- Computing the predictive uncertainty requires samples from the posterior

$$\mathbf{p}(\mathbf{y}^* | \mathbf{x}^*, \mathbf{D}) = \int \mathbf{p}(\mathbf{y}^* | \mathbf{x}^*, \boldsymbol{\theta}) \mathbf{p}(\boldsymbol{\theta} | \mathbf{D}) d\boldsymbol{\theta}$$

- A naive approach is not sufficient if there are many parameters (e.g. millions)

$$\approx \frac{1}{T} \sum_{t=1}^T \mathbf{y}^*(\mathbf{x}^*, \boldsymbol{\theta}_t^s) \quad \text{for } \boldsymbol{\theta}_t^s \sim \mathbf{p}(\boldsymbol{\theta} | \mathbf{D})$$

1. Evaluate the matrix:

$$\mathbf{I}_{\text{inf}} = (\mathbf{U}_A \otimes \mathbf{U}_G) \boldsymbol{\Lambda} (\mathbf{U}_A \otimes \mathbf{U}_G)^T + \mathbf{D}$$

$O(N^2)$ : infeasible

2. Perform Cholesky decomposition:

$$\mathbf{I}_{\text{inf}}^{-1} = \mathbf{F}_c \mathbf{F}_c^T$$

$O(N^3)$ : infeasible

3. Draw samples from the distribution:

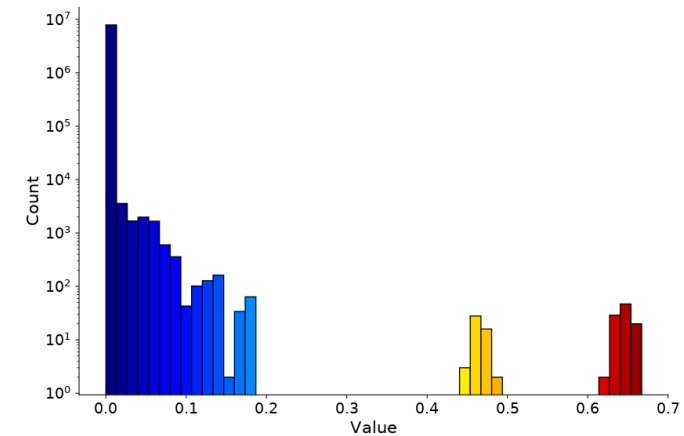
$$\boldsymbol{\theta}_t^s = \boldsymbol{\theta}_{\text{MAP}} + \mathbf{F}_c \mathbf{X}^1 \quad \text{with } \mathbf{X}^1 \text{ the samples of a standard Gaussian}$$

# Low rank sampling computations

- **Step 1:** low rank approximation while saving Kronecker products in eigenvectors:

$$\left( \begin{matrix} \text{Green} \\ \otimes \\ \text{Purple} \end{matrix} \right) \overbrace{\begin{matrix} N \times N \\ \text{Red} \end{matrix}} \left( \begin{matrix} \text{Green} \\ \otimes \\ \text{Purple} \end{matrix} \right)^T \approx \left( \begin{matrix} \text{Green} \\ \otimes \\ \text{Purple} \end{matrix} \right) \overbrace{\begin{matrix} L \times L \\ \text{Red} \end{matrix}} \left( \begin{matrix} \text{Green} \\ \otimes \\ \text{Purple} \end{matrix} \right)^T$$

$$(\mathbf{U}_A \otimes \mathbf{U}_G) \Lambda (\mathbf{U}_A \otimes \mathbf{U}_G)^T \approx (\mathbf{U}_a \otimes \mathbf{U}_g) \Lambda_{1:L} (\mathbf{U}_a \otimes \mathbf{U}_g)^T *$$



- **Step 2:** samples with much lower cost and insignificant errors!

$$\theta_t^s = \theta_{\text{MAP}} + \mathbf{F}_c \mathbf{X}^l$$

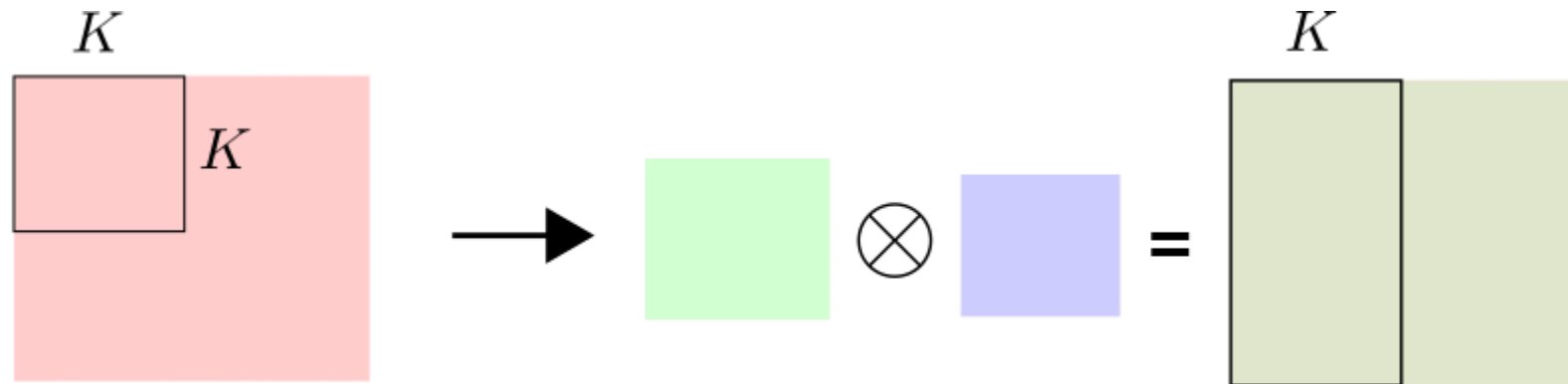
$$\mathbf{F}_c = \mathbf{D}^{-\frac{1}{2}} \left( \mathbf{I}_{nm} - \mathbf{D}^{-\frac{1}{2}} (\mathbf{U}_a \otimes \mathbf{U}_g) \Lambda_{1:L}^{\frac{1}{2}} (\mathbf{C}^{-1} + \mathbf{V}_s^T \mathbf{V}_s)^{-1} \Lambda_{1:L}^{\frac{1}{2}} (\mathbf{U}_a \otimes \mathbf{U}_g)^T \mathbf{D}^{-\frac{1}{2}} \right)$$

# Sparsification algorithm

- How to perform low rank approximation on the Kronecker-factored eigendecomposition?

$$(\mathbf{U}_A \otimes \mathbf{U}_G) \Lambda (\mathbf{U}_A \otimes \mathbf{U}_G)^T \approx (\mathbf{U}_a \otimes \mathbf{U}_g) \Lambda_{1:L} (\mathbf{U}_a \otimes \mathbf{U}_g)^T$$

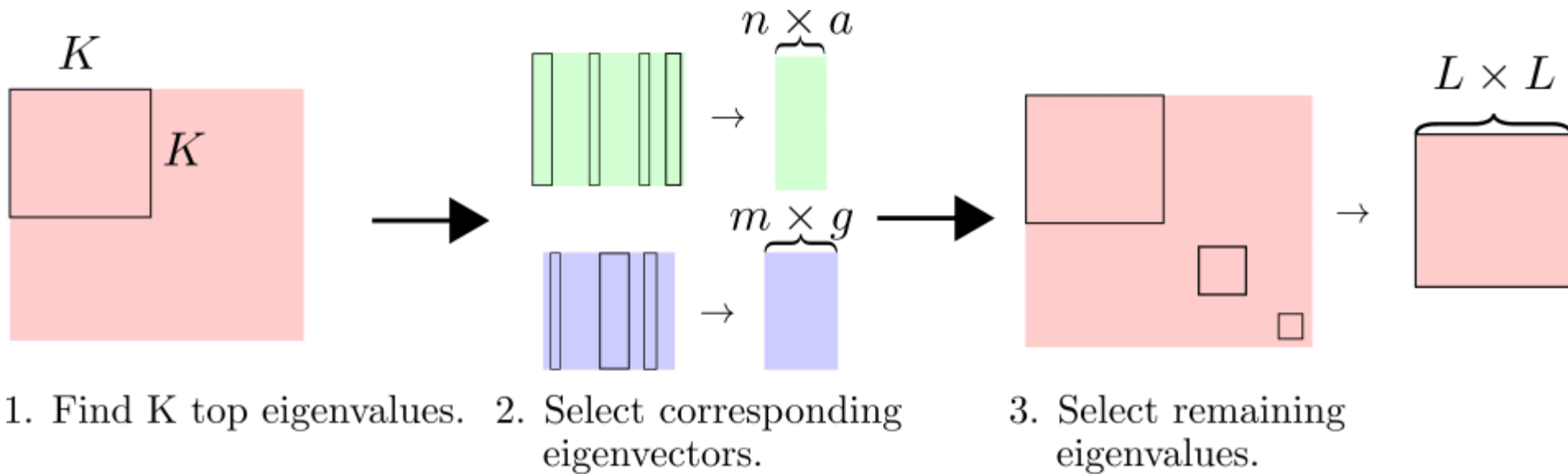
- Conventional low rank approximation such as singular value decomposition:



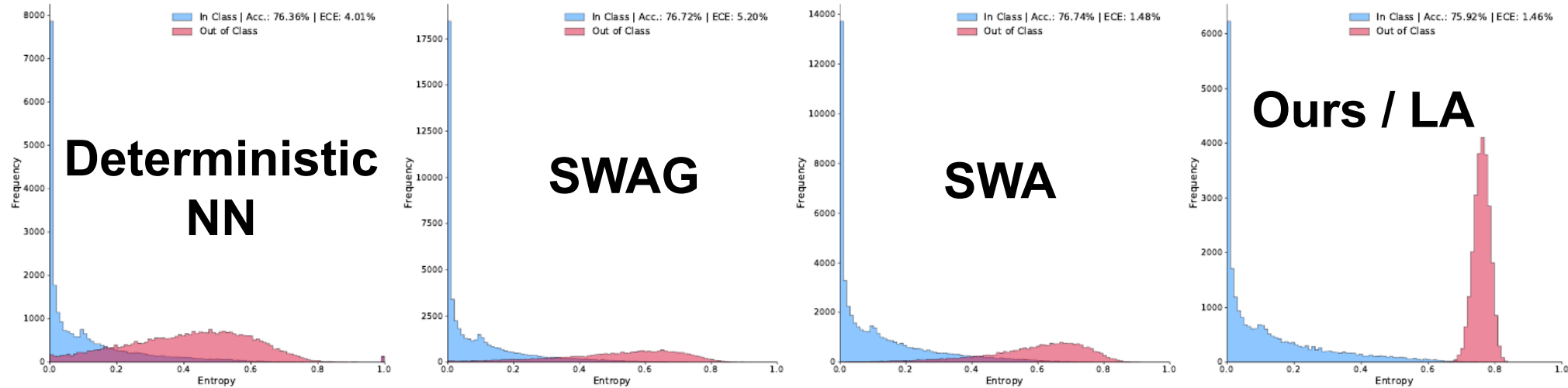
1. Find  $K$  top eigenvalues.
2. Select corresponding eigenvectors.

1. Select the top  $L$  eigenvalues and then:  $\Lambda \approx \Lambda_{1:L}$
2. Using the indices of  $L$  eigenvalues,  $V = (\mathbf{U}_A \otimes \mathbf{U}_G)$  and  $V \approx V_{1:L}$

# Sparsification algorithm



# Results: ImageNet



**Table 3. Network Space Complexity Comparison:** The total number of information matrix parameters and its size in MB are reported for ResNet and DenseNet variants. Lower the better. Here, we also check if the methods take into account the weight correlations (corr).

Model	Diag			KFAC			EFB			INF		
	#Parameters	Size	Corr	#Parameters	Size	Corr	#Parameters	Size	Corr	#Parameters	Size	Corr
<b>ResNet18</b>	11,679,912	44.6	X	95,013,546	362.4	✓	106,693,458	407.0	✓	12,317,373	<b>47.0</b>	✓
<b>ResNet50</b>	25,503,912	97.3	X	153,851,562	586.9	✓	179,355,474	684.2	✓	27,614,896	<b>105.3</b>	✓
<b>ResNet152</b>	60,041,384	229.0	X	389,519,018	1485.9	✓	449,560,402	1714.9	✓	65,558,402	<b>250.1</b>	✓
<b>DenseNet121</b>	7,895,208	30.1	X	103,094,954	393.3	✓	110,990,162	423.4	✓	9,711,081	<b>37.0</b>	✓
<b>DenseNet161</b>	28,461,064	108.6	X	379,105,514	1446.2	✓	407,566,578	1554.7	✓	32,329,191	<b>123.3</b>	✓

# Predictions

# Idea: Theory of Neural Tangent Kernel

## Preliminaries

- Mixtures of experts (MoE):

- M learners:  $f_m(x)$
- A gating function:  $g_m(x)$

- A strict division of data:  $y = \sum_{m=1}^M g_m(x) f_m(x)$

Either 0 or 1



## Bayesian duality

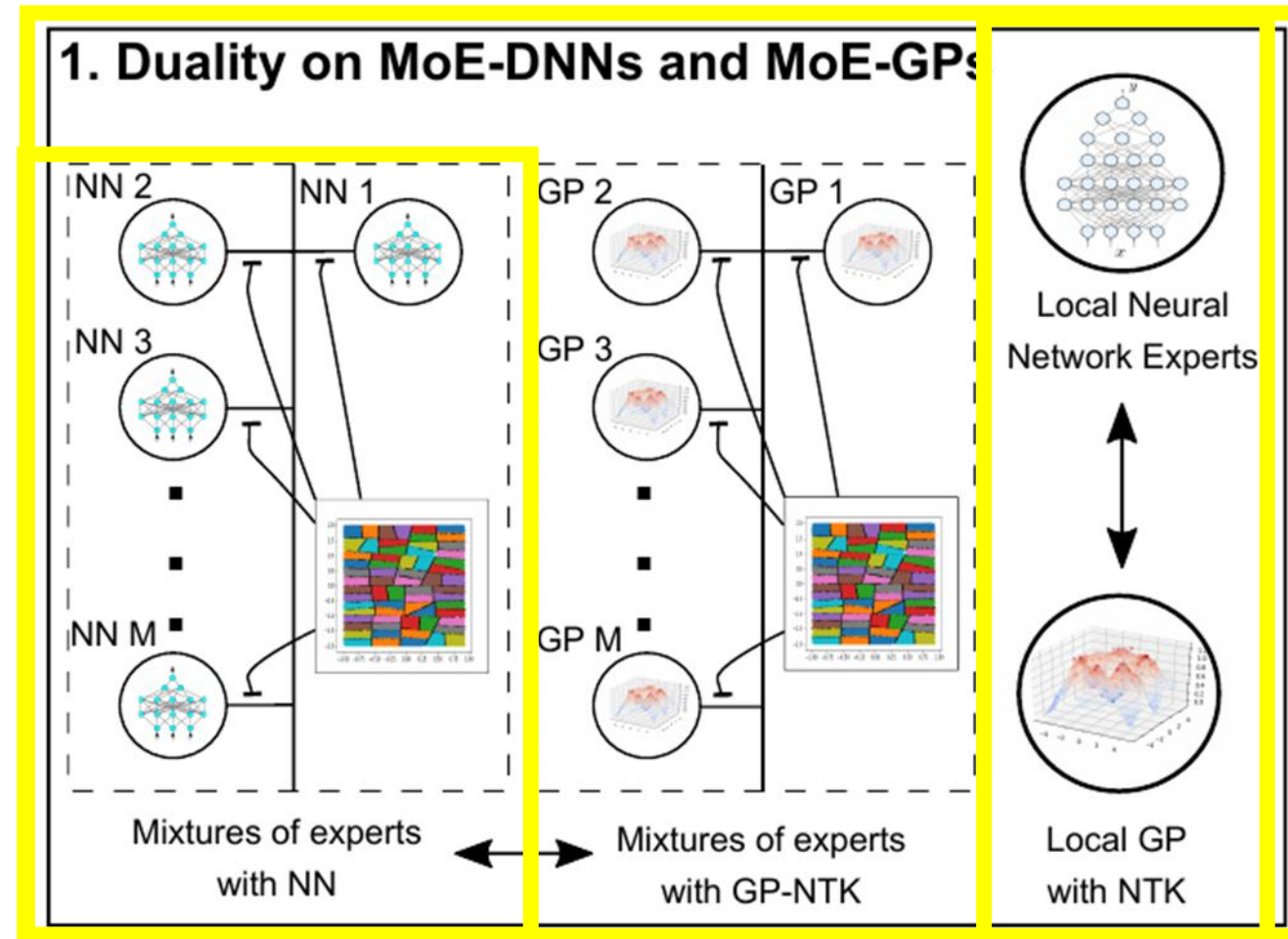
- Local NN  $\leftrightarrow$  Local GPs with NTK:

$$p(\theta_m; \mathbf{D}_m) \leftrightarrow p(f_m; \tilde{\mathbf{D}}_m)$$

- Direct application of Khan et al (2019).

- MoE NN  $\leftrightarrow$  MoE GP with NTK:

$$\prod p(\theta_m; \mathbf{D}_m) \leftrightarrow \prod p(f_m; \tilde{\mathbf{D}}_m)$$



Note: NN is any neural networks with valid Jacobians!

# Idea: Theory of Neural Tangent Kernel

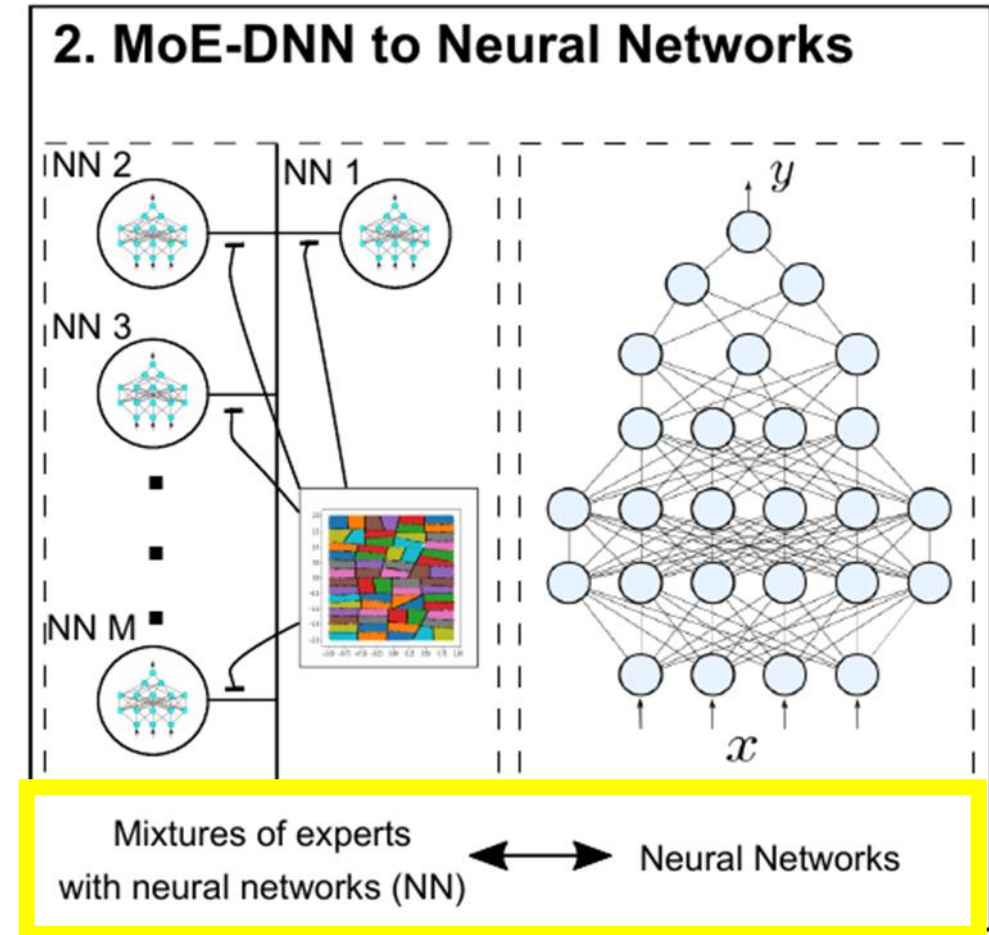
## Problem

- MoE NN  $\leftrightarrow$  MoE GP (not the goal)
- A single NN  $\leftrightarrow$  MoE GP (original goal)

## Insight

- Imagine an already trained NN:  $f_{\theta}(x)$
- Assume all NN experts are the same:  
 $f_{\theta}(x) = f_{\theta_1}(x) = f_{\theta_2}(x) \dots = f_{\theta_M}(x)$
- Due to hard partitioning, input-output relationship are the same:

$$y = \sum_{m=1}^M g_m(x) f_{\theta}(x) = f_{\theta}(x)$$



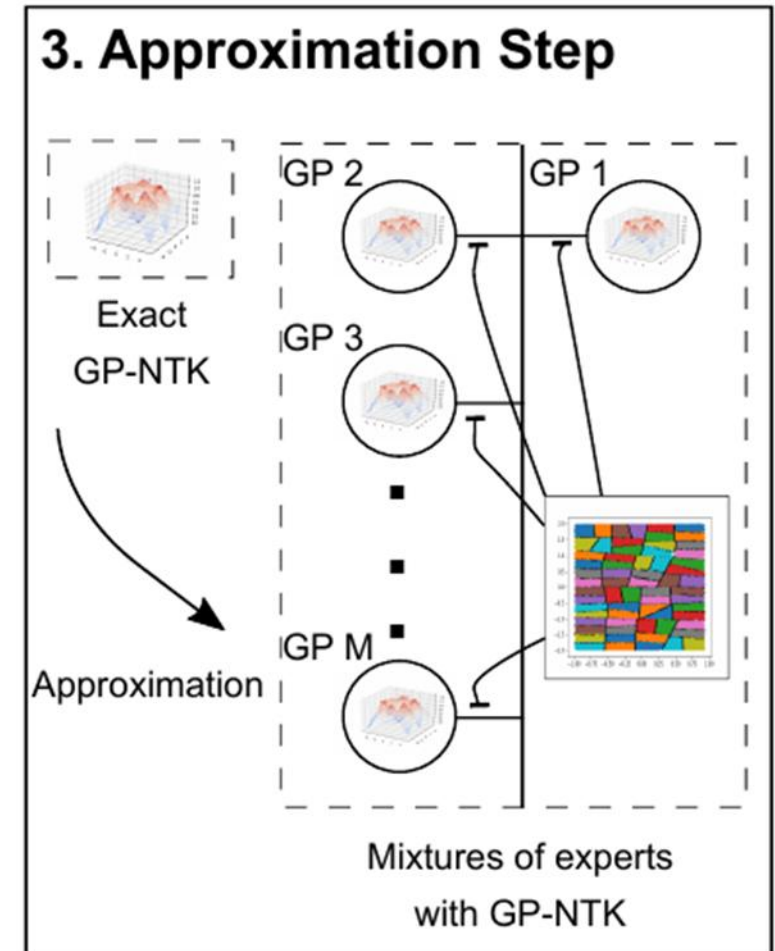
# Idea: Theory of Neural Tangent Kernel

## Result

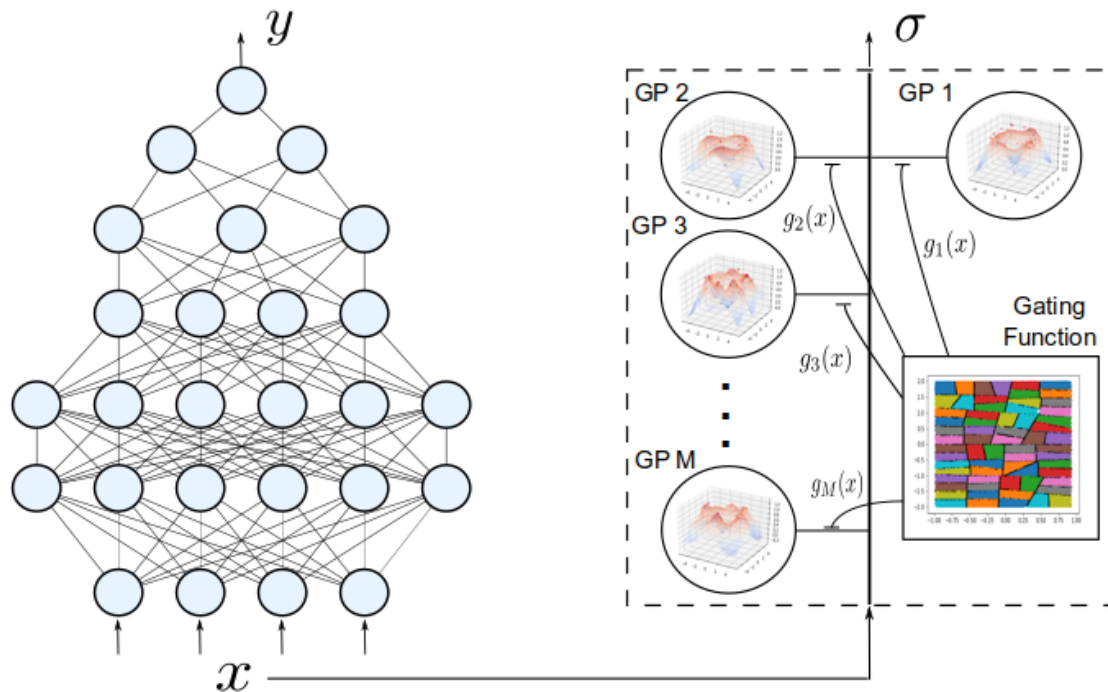
- A single NN imagination  $\rightarrow$  GPs with  $K = K_1 = K_2 \dots = K_M$ .
- A full GP  $\approx$  MoE GPs (block diagonal approximation)

## Implications

- Neural networks for mean predictions:  $y = f_w(\mathbf{x})$
- MoE GP for covariance estimation:  $\Sigma = K_{**} - K_*^T [K + \sigma_n^2 I]^{-1} K_*$ 
  - The gating function has picked a GP expert.
  - The NTK is a neural network based kernel:  $\text{GP}(\mathbf{0}, \frac{1}{\delta_m} \mathbf{J}_{f_m}^T(\mathbf{x}) \mathbf{J}_{f_m}(\mathbf{x}))$
  - Note:  $y \neq \tilde{y}$  but  $\text{COV}(y) \approx \text{COV}(\tilde{y})$ .
- Closed-form solution for uncertainty:  $p(y^* | x^*, D) = \mathcal{N}(f_{\theta=\hat{\theta}(x^*)}, \Sigma)$



# Idea: Theory of Neural Tangent Kernel



- Due to the mixtures of experts idea (Jacobs et al 1991), the complexity of full GP reduces (Tresp 2001):

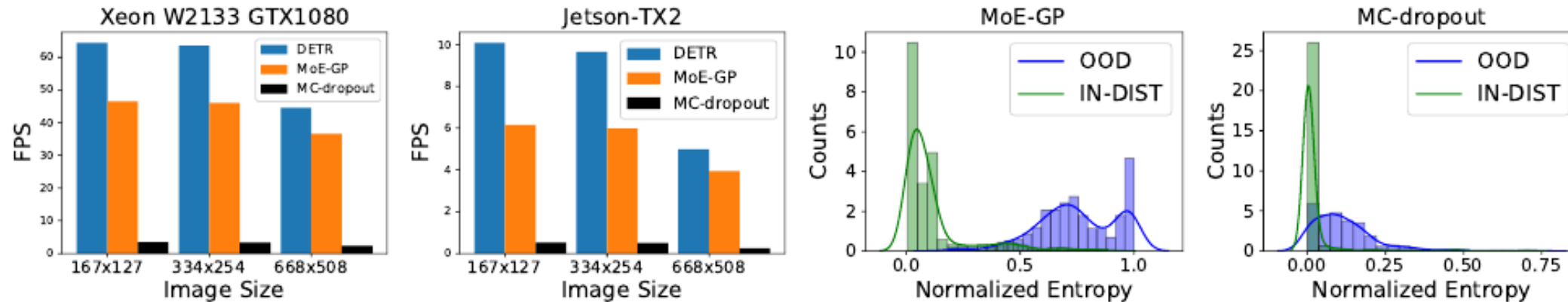
$$\begin{bmatrix} \tilde{y}_m \\ \tilde{f}_m \end{bmatrix} \sim N \left( \mathbf{0}, \begin{bmatrix} \mathbf{K}_m + \sigma_{0,m} \mathbf{I} & \mathbf{k}_{m,*} \\ \mathbf{k}_{m,*}^T & k_{m,**} \end{bmatrix} \right)$$

- Analytical method for uncertainty Estimates (“sampling-free”):

$$\Sigma_m = k_{m,**} - \mathbf{k}_{m,*}^T (\mathbf{K}_m + \sigma_0 \mathbf{I})^{-1} \mathbf{k}_{m,*} + \sigma_0$$

$$p(c|\mathbf{z}_m) = \left( \frac{\mathbf{z}_m}{\sqrt{\mathbf{1} + \lambda_{m,0} \Sigma_m(\mathbf{x}^*)}} \right) \quad (\text{Lu et al 2020})$$

# Results: Planetary scenario

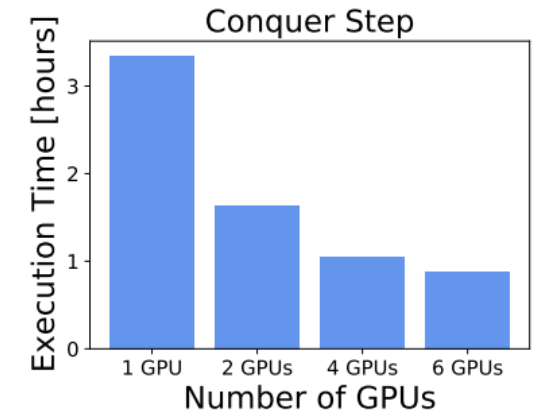
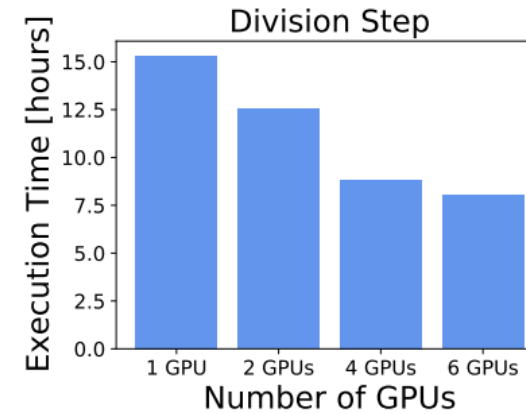


- Run-time comparison on a GPU-desktop and an embedded GPU. Higher FPS the faster.
- Entropy histogram. More separable, better calibrated the uncertainty estimates.

**Main take-away/use-cases:** when sparse GPs can scale, real-time uncertainty estimates from a GP formulation of neural networks can be obtained, improving over the state-of-the-art methods.

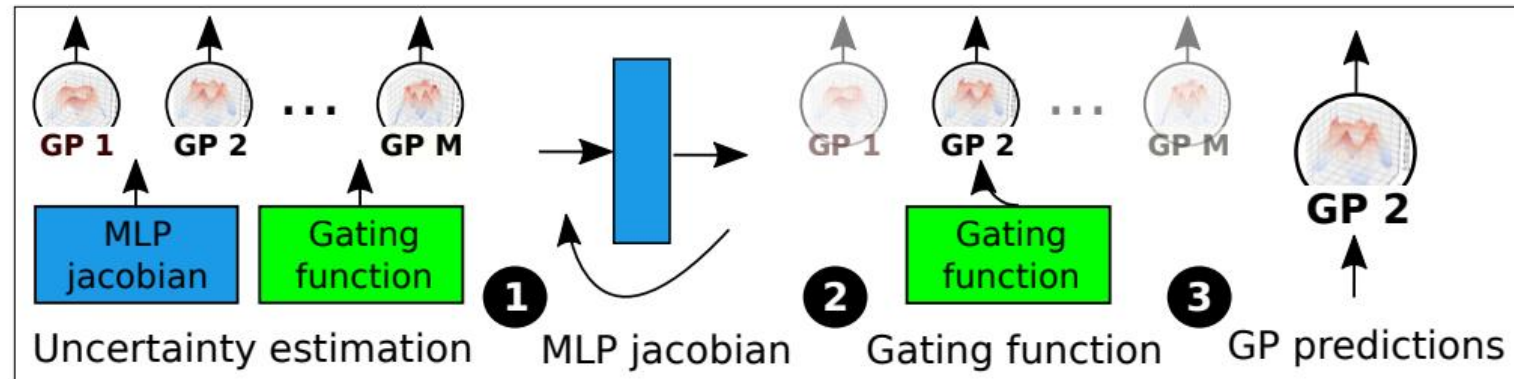
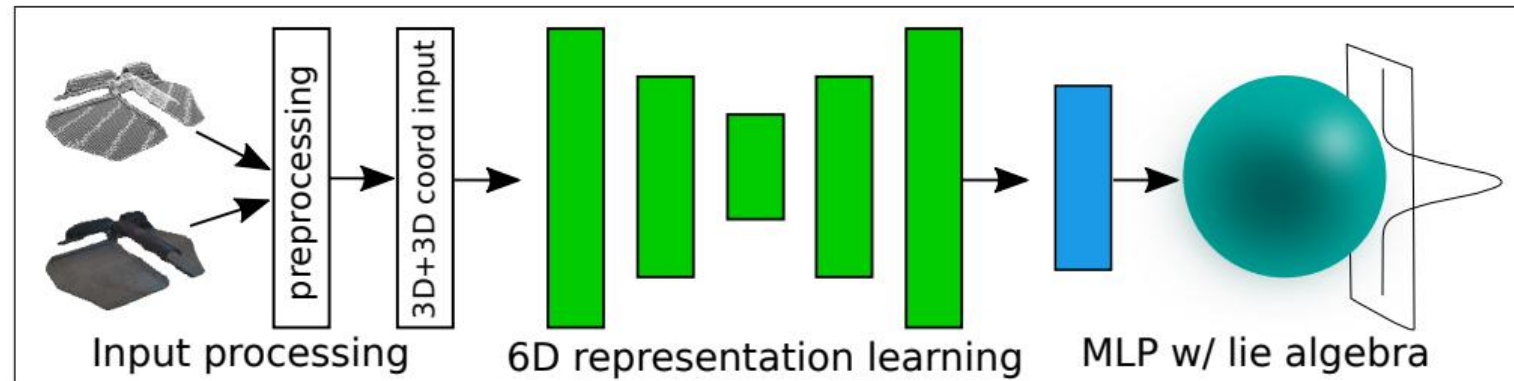
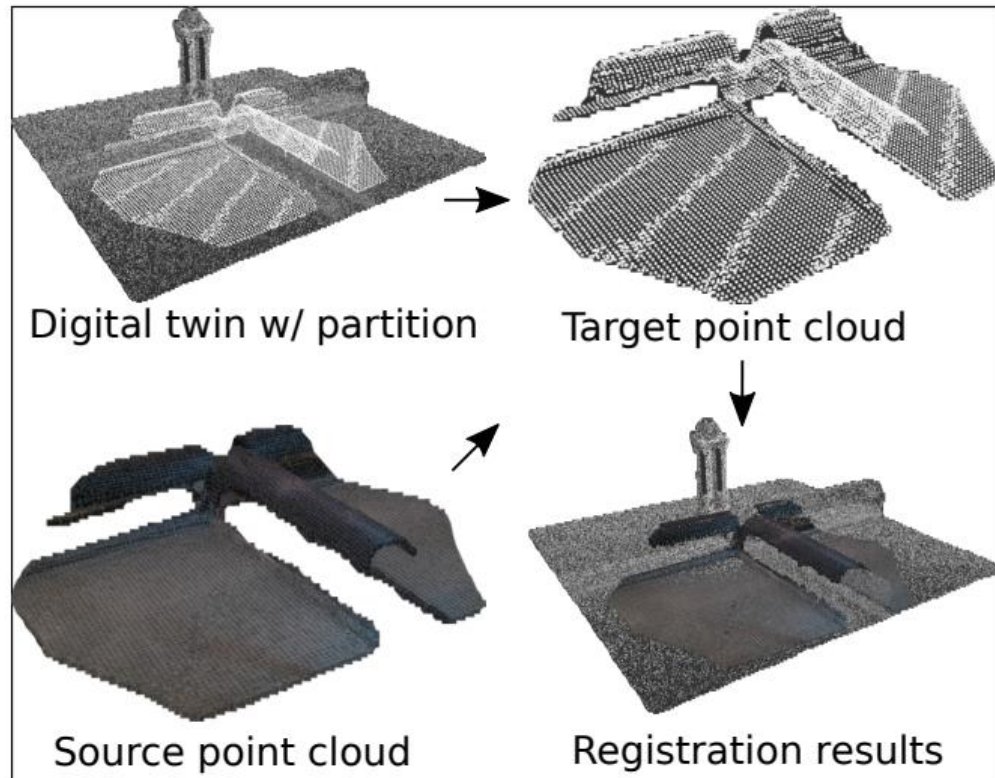
# Results: Scalability

- Time to get a GP model for 2 million data points.
  - Evaluated on inverse dynamics tasks.
- E.g., 180360 train points on kuka1 and kuka2



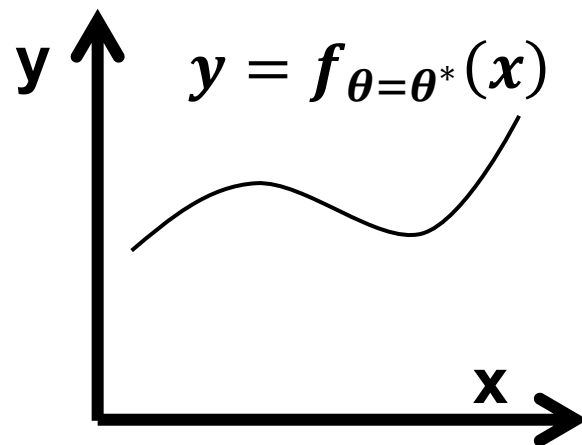
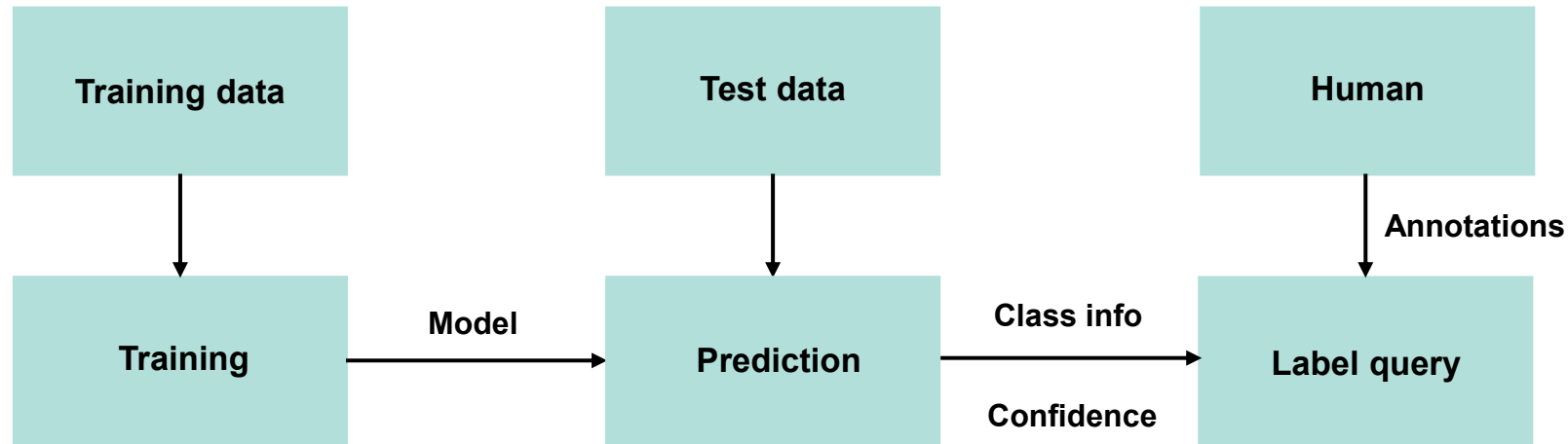
Train	Test	Ensemble	MC-dropout	rMC-dropout	Laplace	rLaplace	Ours
sarcos	sarcos	$1.261 \pm 0.105$	$1.398 \pm 0.012$	$1.398 \pm 0.012$	$1.392 \pm 0.014$	$1.392 \pm 0.014$	<b><math>1.035 \pm 0.039</math></b>
sarcos	kuka1	$17.86 \pm 2.995$	<b><math>16.21 \pm 0.866</math></b>	$20.58 \pm 1.144$	$25.88 \pm 1.354$	$24.55 \pm 2.746$	$21.53 \pm 2.268$
sarcos	kuka2	$17.50 \pm 2.895$	<b><math>15.63 \pm 0.858</math></b>	$20.17 \pm 0.988$	$25.59 \pm 1.079$	$24.42 \pm 2.257$	$20.92 \pm 2.236$
sarcos	kukasim	<b><math>23.06 \pm 2.649</math></b>	$51.09 \pm 11.14$	$61.40 \pm 12.921$	$77.13 \pm 15.40$	$74.03 \pm 13.23$	$68.95 \pm 4.003$
kuka1	kuka1	$2.013 \pm 0.020$	$1.611 \pm 0.007$	$1.620 \pm 0.006$	$1.676 \pm 0.013$	$1.719 \pm 0.071$	<b><math>1.347 \pm 0.013</math></b>
kuka1	kuka2	$2.085 \pm 0.005$	$1.349 \pm 0.019$	$1.330 \pm 0.006$	<b><math>1.310 \pm 0.005</math></b>	<b><math>1.310 \pm 0.008</math></b>	$1.315 \pm 0.003$
kuka1	kukasim	$60.44 \pm 1.108$	$42.14 \pm 2.869$	$48.76 \pm 0.566$	$60.37 \pm 0.799$	$59.74 \pm 1.571$	<b><math>1.348 \pm 0.012</math></b>
kuka1	sarcos	$8.128 \pm 0.169$	$22.24 \pm 4.288$	$42.36 \pm 2.398$	$122.68 \pm 13.32$	$1837 \pm 2431.3$	<b><math>1.356 \pm 0.014</math></b>
kuka2	kuka2	$2.042 \pm 0.009$	$1.443 \pm 0.005$	$1.423 \pm 0.006$	$1.429 \pm 0.006$	$1.423 \pm 0.006$	<b><math>1.354 \pm 0.013</math></b>
kuka2	kukasim	$63.07 \pm 0.741$	$38.94 \pm 0.561$	$48.93 \pm 0.779$	$60.36 \pm 1.005$	$59.05 \pm 1.158$	<b><math>1.333 \pm 0.012</math></b>
kuka2	sarcos	$8.509 \pm 0.461$	$18.24 \pm 1.599$	$53.24 \pm 6.287$	$141.7 \pm 17.14$	$211.3 \pm 180.2$	<b><math>1.355 \pm 0.014</math></b>
kuka2	kuka1	$2.106 \pm 0.010$	$1.461 \pm 0.004$	$1.395 \pm 0.004$	$1.373 \pm 0.005$	$1.374 \pm 0.004$	<b><math>1.331 \pm 0.013</math></b>

# SPIRIT: proposed method

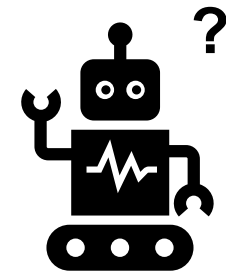


# Active Learning

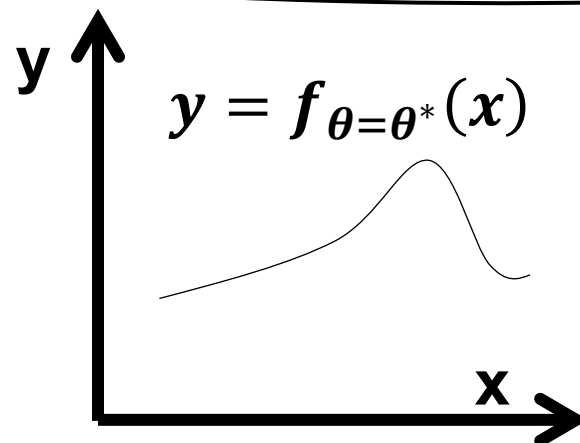
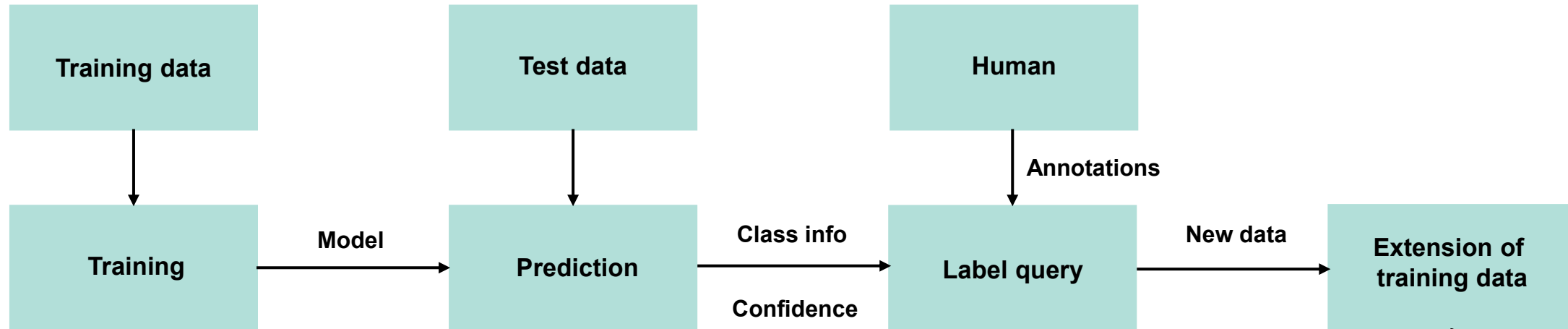
# Stream-based active learning



This is a t-shirt!



# Stream-based active learning



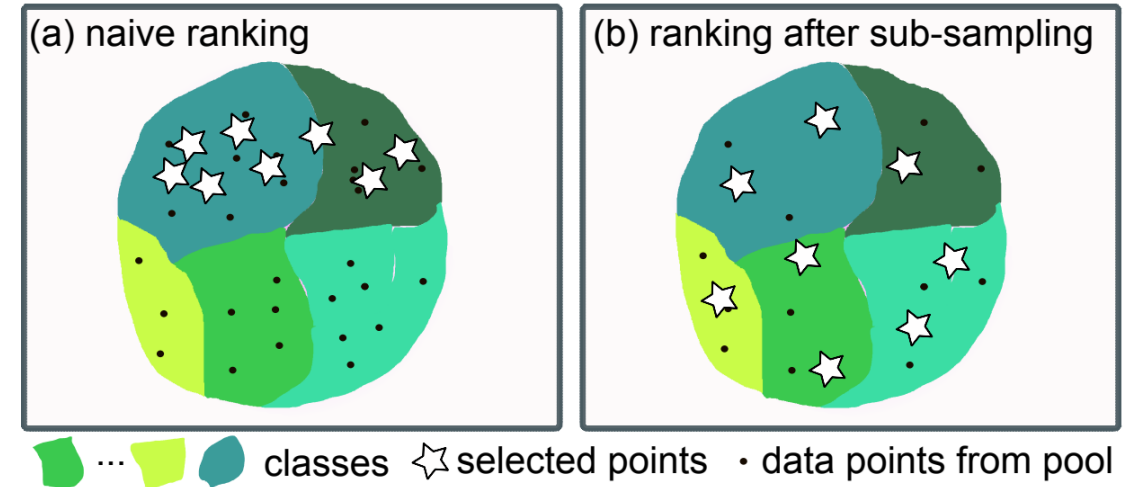
- Learning from limited data.
  - Need good generalization
- Ask for help in a right way.
  - Need good uncertainty.

# Active learning criteria

$$I[y; \theta \mid x, D] = H[y \mid x, D] - \mathbb{E}_{p(\theta|D)}[H[y \mid x, \theta]]$$

Mutual information:

This measures how much learning the label  $y$  at  $x$  would reduce uncertainty about the model parameters.



$H[y \mid x, D]$  Predictive entropy — how uncertain the model's predictions are overall.

$\mathbb{E}_{p(\theta|D)}[H[y \mid x, \theta]]$  Expected entropy if we knew which parameter setting  $\theta$  is true (average over model posterior).

Measures **how much uncertainty comes from the parameters** (epistemic).

Diversity was also important for DL.

# Open-set evaluation with users

