

# Bayesian Deep Learning: Insights, Methods, and Applications

Zhijie Deng

Tsinghua University



# Deep learning success

Google search results for "deep learning". The search bar shows "deep learning" and the results page displays a featured snippet from Investopedia, a "People also ask" section with four questions, and a "Deep learning" knowledge panel with images and a definition. Below the knowledge panel are "Deep learning books" and "People also search for".

**Deep learning** is a subset of machine learning in artificial intelligence (AI) that has networks capable of learning unsupervised from data that is unstructured or unlabeled. Also known as **deep neural learning** or **deep neural network**. Apr 30, 2019

[Deep Learning Definition - Investopedia](https://www.investopedia.com/terms/d/deep-learning)  
<https://www.investopedia.com/terms/d/deep-learning>

**People also ask**

- Why is it called deep learning?
- What is deep learning examples?
- What is deep learning vs Machine Learning?
- What is deep learning and how it works?

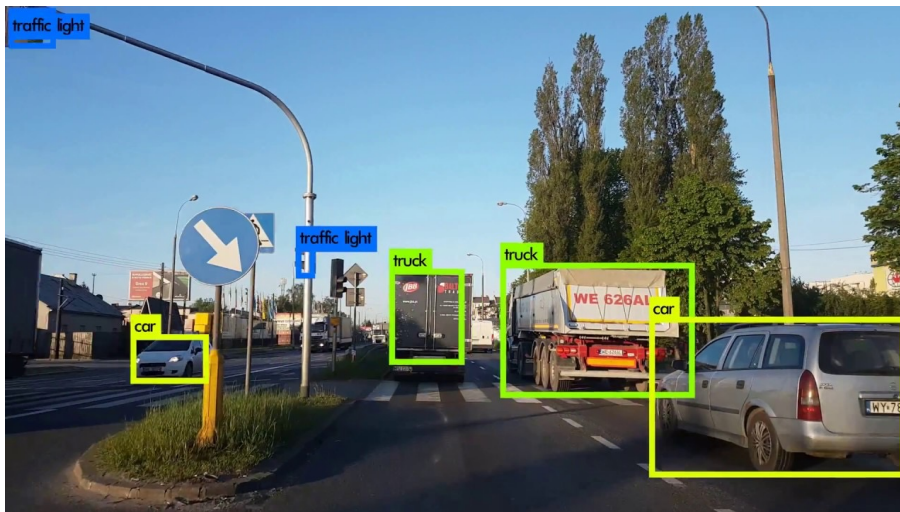
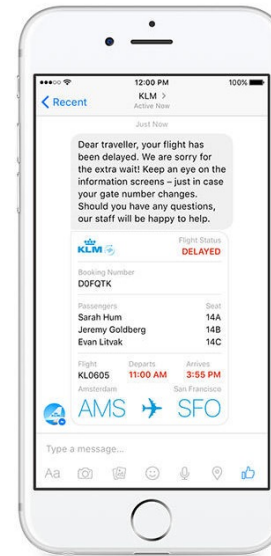
**Deep learning**

Deep learning is part of a broader family of machine learning methods based on artificial neural networks. Learning can be supervised, semi-supervised or unsupervised. [Wikipedia](#)

**Deep learning books**

- Deep Learning with Pyth...
- Deep Learning
- Deep Learning: A Practiti...
- Python Machine Learning
- Hands-On Machine Learning...

**People also search for**



Google Translate interface showing the translation of "city" from English to French, resulting in "ville". The interface includes a text input field, a dropdown menu for the target language (French), and a list of definitions for "ville".

**Google Translate**

ENGLISH - DETECTED | ENGLISH | SPANISH | FRENCH | FRENCH | ENGLISH | SPANISH

city | ville

Definitions of city

- Noun**
- ① a large town. "But we do not accept this fate with the torpor of other city dwellers."
- ② a place or situation characterized by a specified attribute. "panic city"
- ③ the financial and commercial district of London, England. "Reaction in the City was on the cool side, as it also tended to be in Europe."

Translations of city

**Noun**

to ville city, town, place, burgh

# Problems remain

- **Lack uncertainty**

## 每小时预报



- **Data driven** (label-eager, poor robustness, etc.)

- CV systems trained on ImageNet (**1M+** images)
- ASR (speech) systems trained on **11,000+ hrs** of annotated data
- OntoNotes (English) NER dataset contains **625,000** annotated words



# The promise of probabilistic (Bayesian) modeling

Uncertainty

Prior belief

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



Thomas Bayes (1702 – 1761)



U. Cambridge  
Fellow of the Royal Society (FRS)

## REVIEW

doi:10.1038/nature14541

### Probabilistic machine learning and artificial intelligence

Zoubin Ghahramani<sup>1</sup>

How can a machine learn from experience? Probabilistic modelling provides a framework for understanding what learning is, and has therefore emerged as one of the principal theoretical and practical approaches for designing machines that learn from data acquired through experience. The probabilistic framework, which describes how to represent and manipulate uncertainty about models and predictions, has a central role in scientific data analysis, machine learning, robotics, cognitive science and artificial intelligence. This Review provides an introduction to this framework, and discusses some of the state-of-the-art advances in the field, namely, probabilistic programming, Bayesian optimization, data compression and automatic model discovery.

# Uncertainty: let models knowing their limits

All models are wrong, but some models that know when they are wrong, are useful.

 **Charlie Langton**   
@charlielangton

A Tesla is trapped under a semi truck on Fort Street near Waterman in Southwest Detroit. Two passengers were taken to the hospital and they are in critical condition. The driver of the semi is not hurt. [@WWJ950](#)  
[@FOX2News](#)



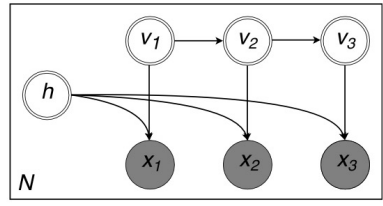
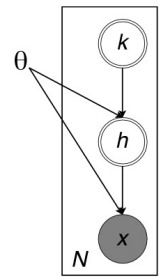
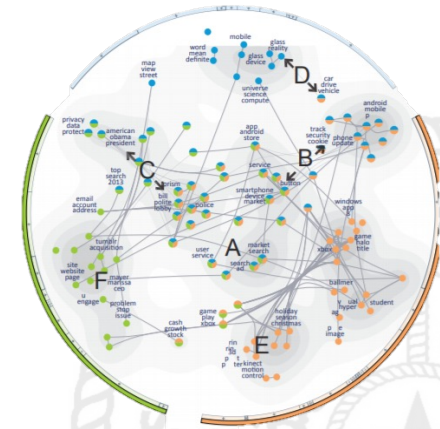
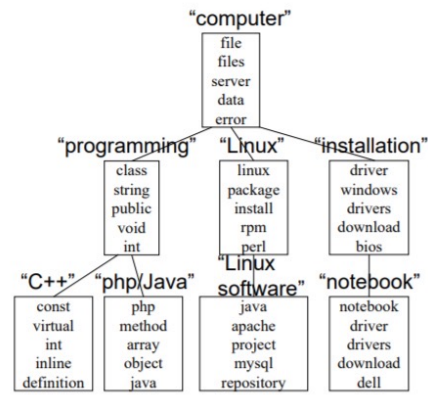
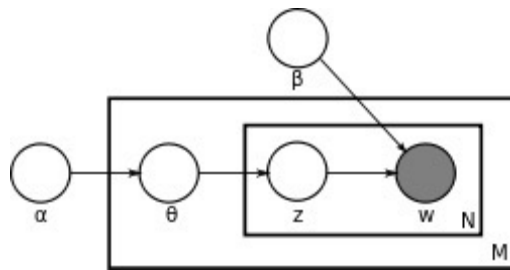
下午8:14 · 2021年3月11日 来自 Detroit, MI · Twitter for iPhone



costly mistakes

# Prior knowledge: deconstruct the black-box of the model

Tell the machine what we know and let them focus on what we do not know.



# The success of probabilistic (Bayesian) modeling

$$\begin{array}{ccc} \text{Uncertainty} & & \text{Prior belief} \\ \downarrow & & \downarrow \\ p(\theta|D) = \frac{p(D|\theta)\pi(\theta)}{p(D)} \end{array}$$

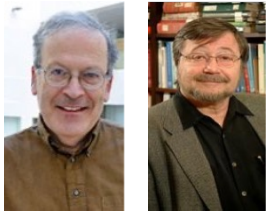


Thomas Bayes (1702 – 1761)

- Prediction:  $p(x|D, H) = \int p(x|\theta, D, H)p(\theta|D)d\theta$
- Model selection  $p(D|H_1) \geq ? \text{ (or } \leq ?) p(D|H_2)$   $p(D|H) = \int p(D|\theta)p(\theta|H)d\theta$
- Regularization  $\max_{\mathbf{w}} \log p(y|X, \mathbf{w}) + \log p(\mathbf{w})$
- Modeling with latent var.  $p(z|x) = \frac{p(x, z)}{p(x)} \propto p(z)p(x|z; \theta)$

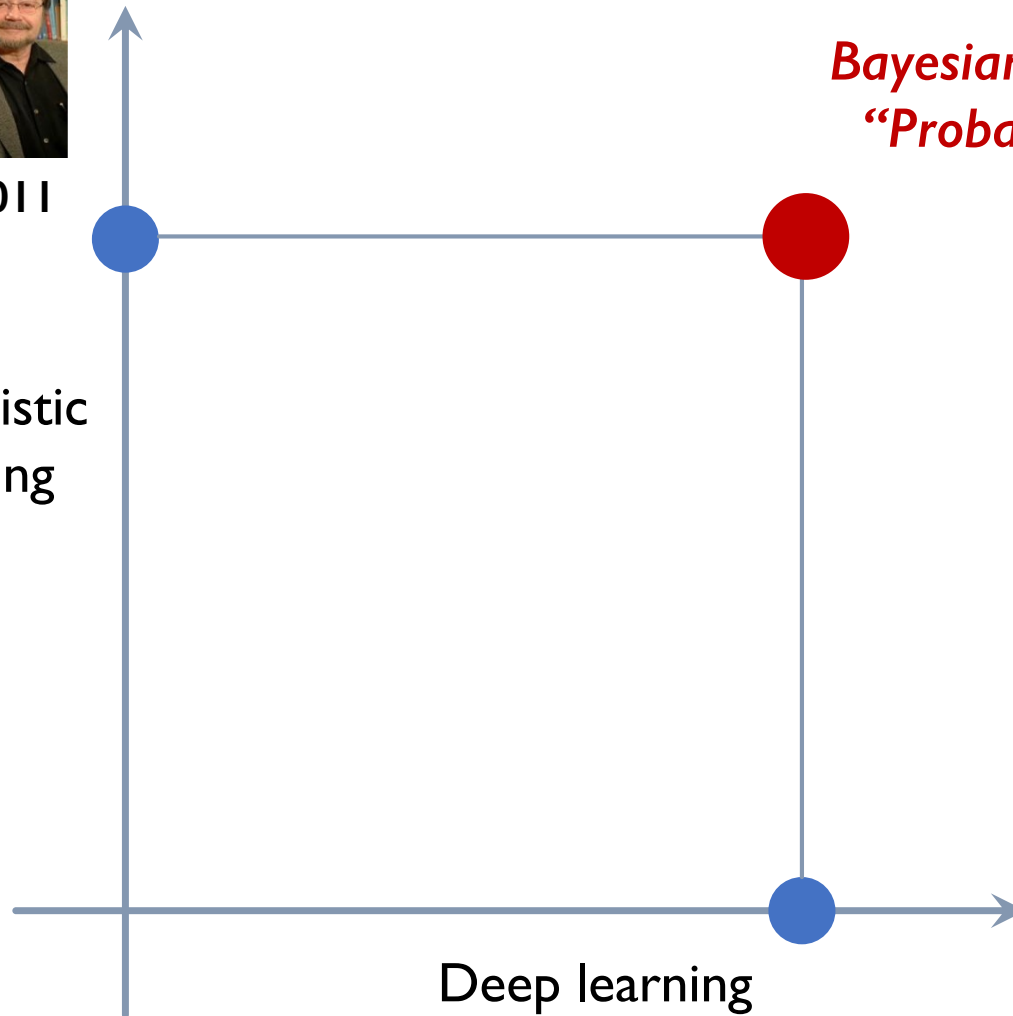
# Research focus: Bayesian deep learning (BDL)

Probabilistic modeling *meets* deep learning



2010 2011

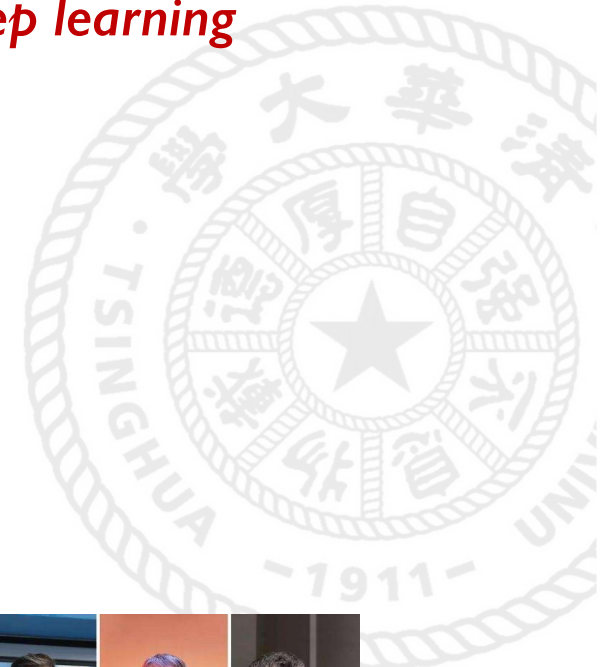
Probabilistic modeling



*Bayesian deep learning (BDL):  
“Probabilistic” revolution of  
deep learning*

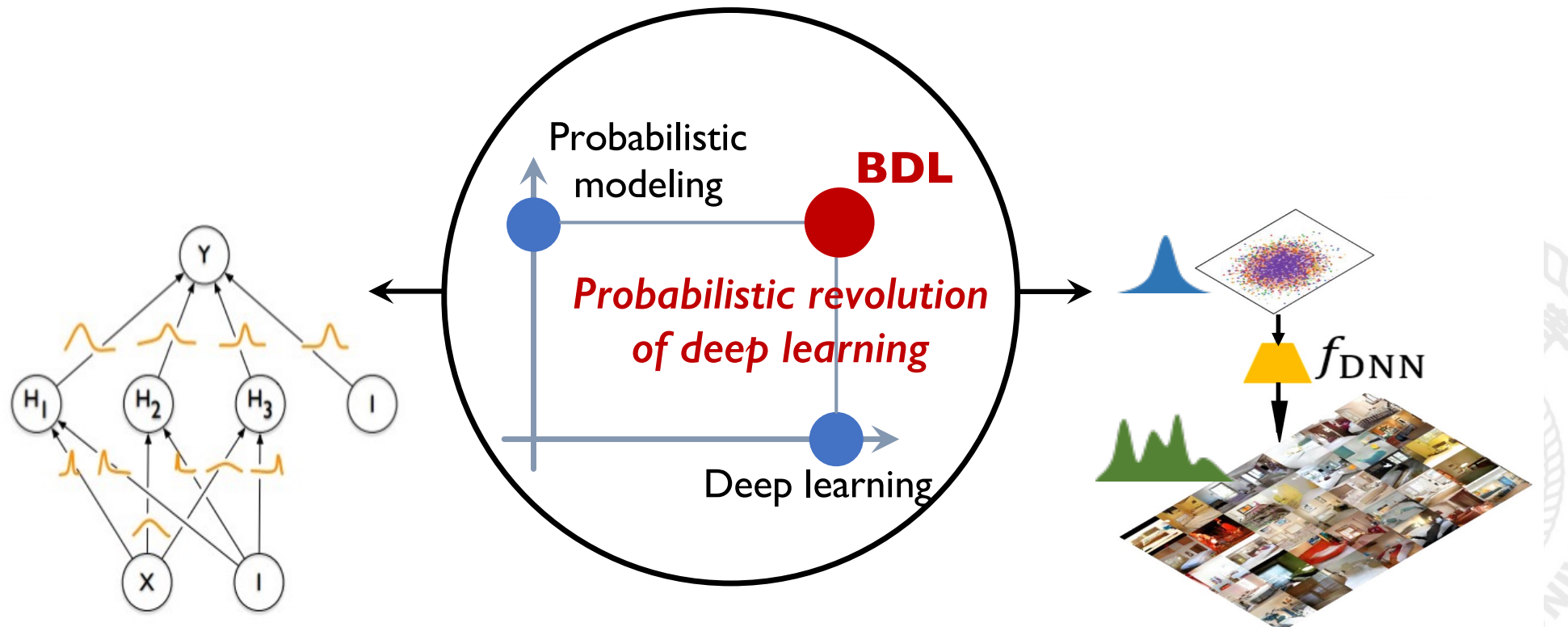


2018





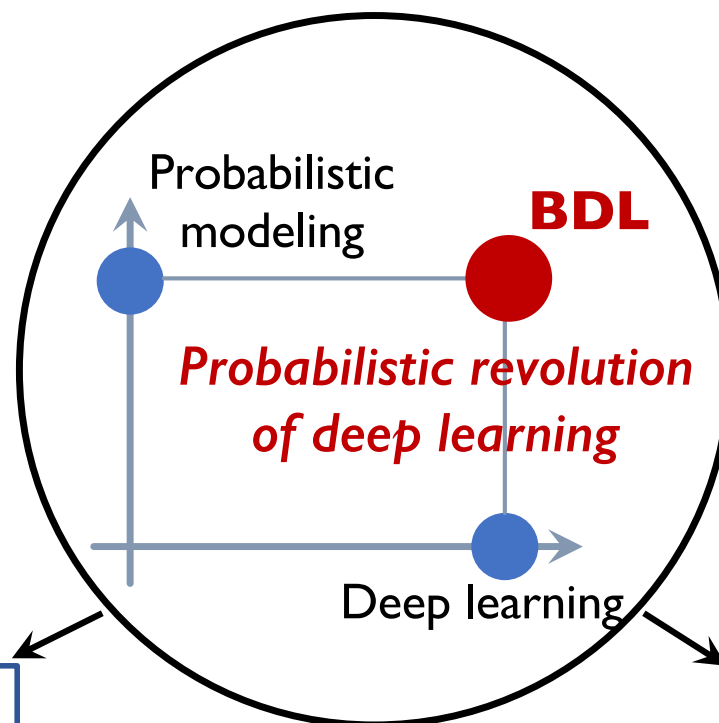
# Research focus: Bayesian deep learning (BDL)



*Probabilistic modeling of DNNs :*  
Bayesian neural networks (BNNs)  
**Benefits:** uncertainty estimation

*DNNs enrich probabilistic models:* deep generative models (DGMs)  
**Benefits:** incorporating prior knowledge

# Research focus: Bayesian deep learning (BDL)



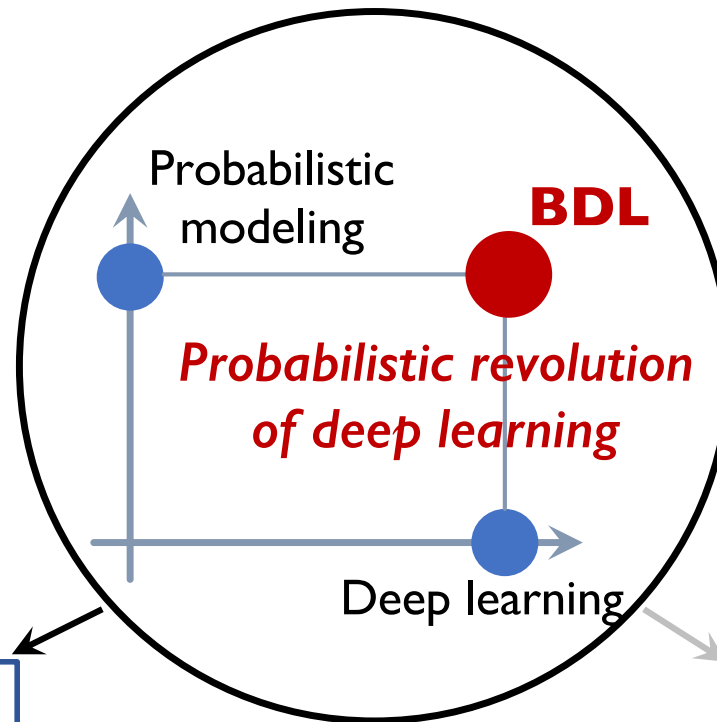
## Theme I: Uncertainty quantification via BNNs

(NeurIPS 2020; CVPR 2021; ICLR 2021 workshop; NeurIPS sub.)

## Theme II: Knowledge incorporation via DGMs

(NeurIPS 2017; ICCV 2019, NeurIPS 2020; ICML 2019 workshop)

# Research focus: Bayesian deep learning (BDL)



## Theme I: Uncertainty quantification via BNNs

(NeurIPS 2020; CVPR 2021; ICLR 2021 workshop; NeurIPS sub.)

## Theme II: Knowledge incorporation via DGMs

(NeurIPS 2017; ICCV 2019, NeurIPS 2020; ICML 2019 workshop)

# Uncertainty is ubiquitous



Road conditions



Traffics



Pedestrian behaviors



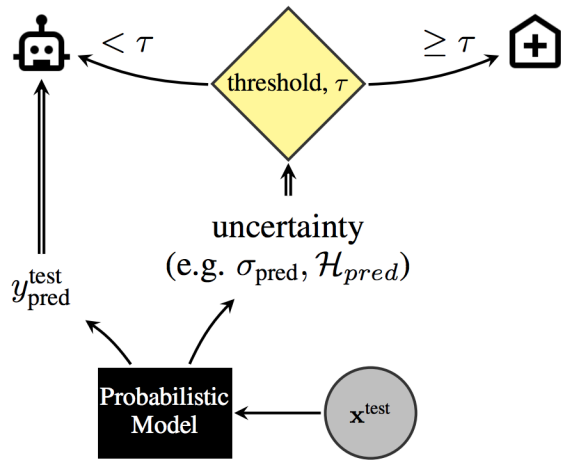
Even **malicious**



AlexNet: lionfish, confidence 81.3%  
VGG-16: lionfish, confidence 93.3%  
ResNet-18: lionfish, confidence 95.6%

# Uncertainty is the key to bringing DL to the masses

## Automated diagnosis



## Self-driving cars

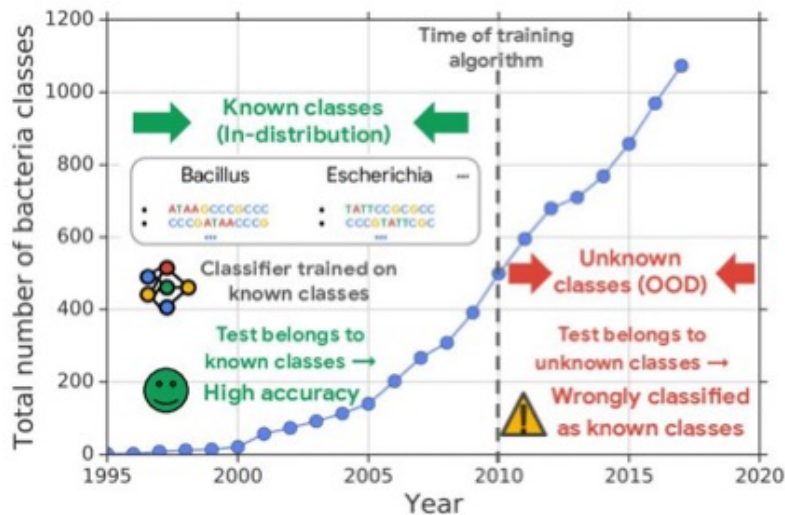


## Conversational dialog systems

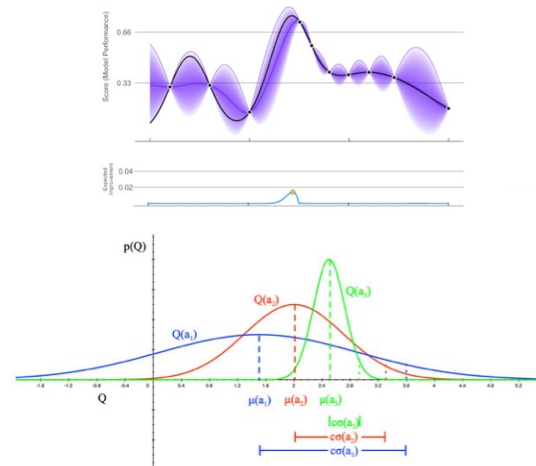


Figure 1: Example exchanges between a user (blue, right side) and a task-driven dialog system for personal finance (grey, left side). The system correctly identifies the user's query in ①, but in ② the user's query is mis-identified as in-scope, and the system gives an unrelated response. In ③ the user's query is correctly identified as out-of-scope and the system gives a fallback response.

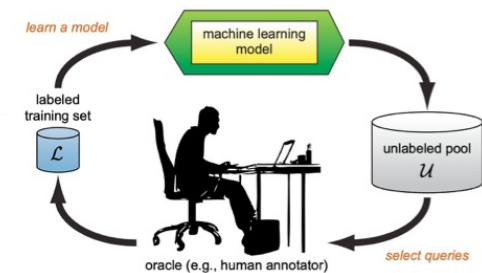
## Open set recognition



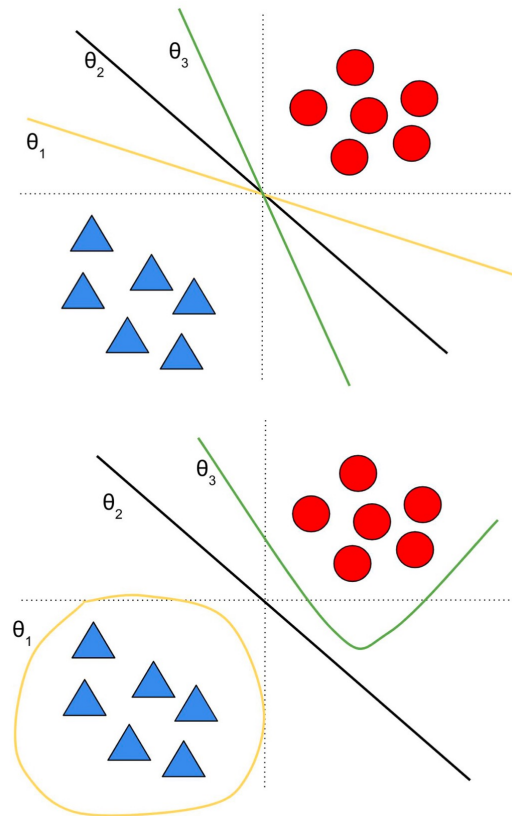
## Bayesian optimization and reinforcement learning



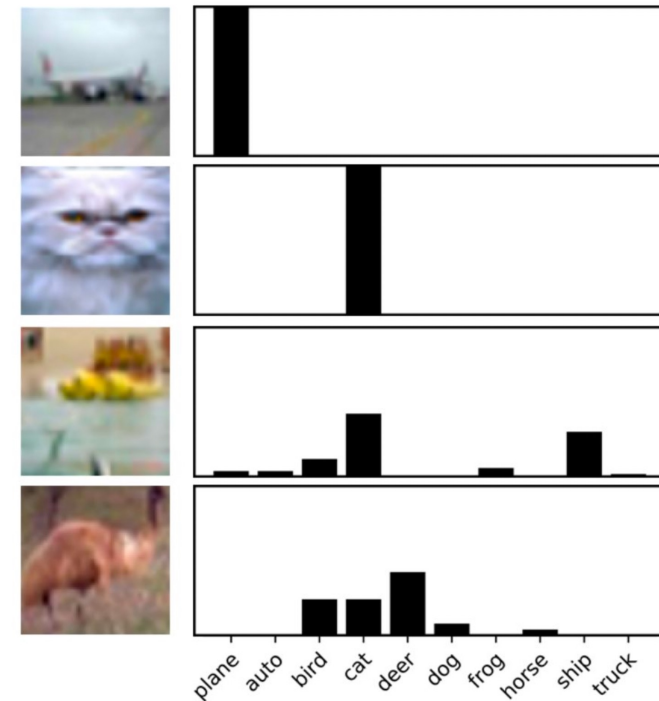
## Active Learning



# Types of Bayesian uncertainty



- Model (**epistemic**) uncertainty
- Various interpretations for the data
- Reducible
- Models can be from same hypotheses class or not

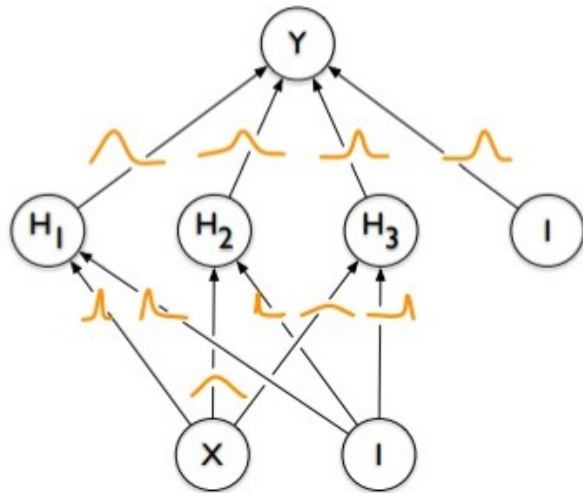


- Data (**aleatoric**) uncertainty
- Stem from labeling noise, measurement noise, or missing data
- Irreducible\*

# Bayesian uncertainty in DL

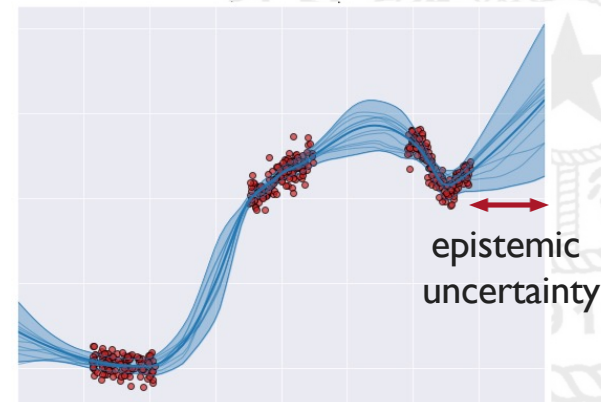
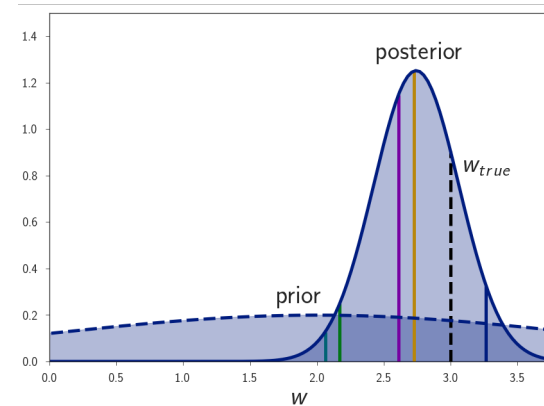
## Bayesian neural networks (BNNs)

Bayesian treatment of DNNs (weights) captures uncertainty



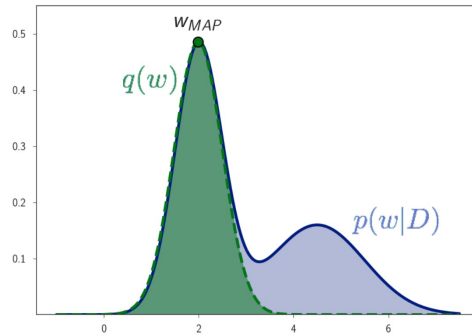
$$p(w|D) = \frac{p(D|w)p(w)}{p(D)}$$

$$p(y^*|x^*, D) = \int_w p(y^*|x^*, w)p(w|D)dw$$

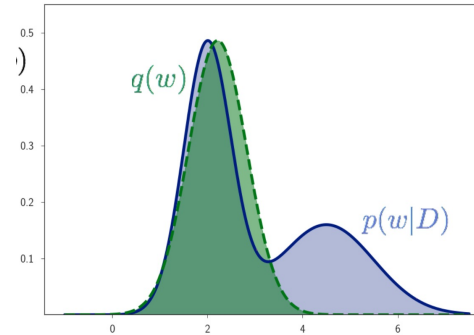


# Core of BNNs: posterior inference

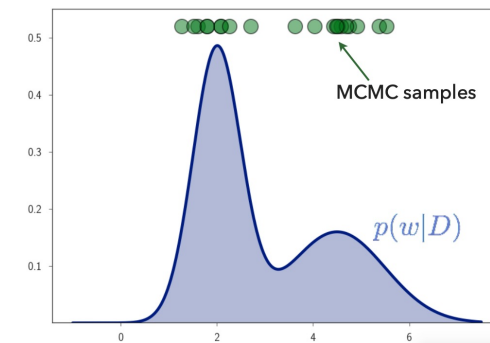
## Methods and challenges



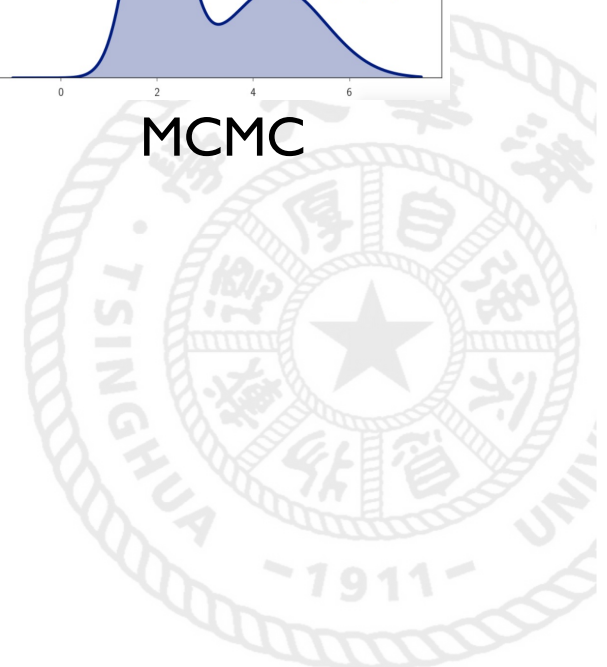
Laplace approx.



Variational inference



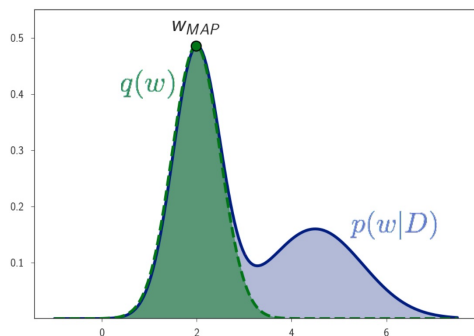
MCMC



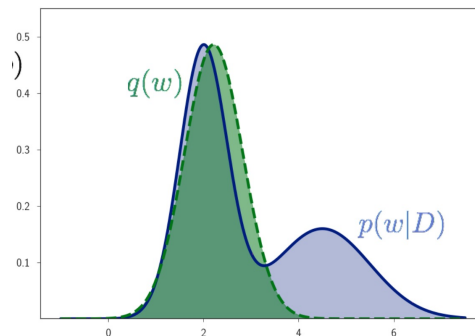


# Core of BNNs: posterior inference

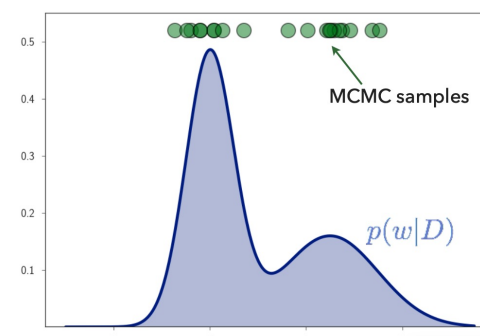
## Methods and challenges



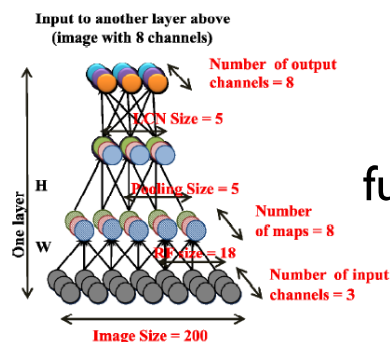
Laplace approx.



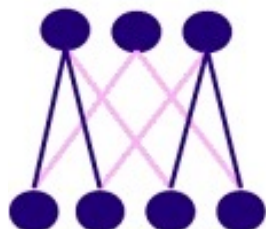
Variational inference



MCMC



**High-dim. weight space** poses fundamental obstacles for from-scratch inference



**Over-parameterization** nature of DNNs leads to collapsed weight uncertainty

Existing BNNs

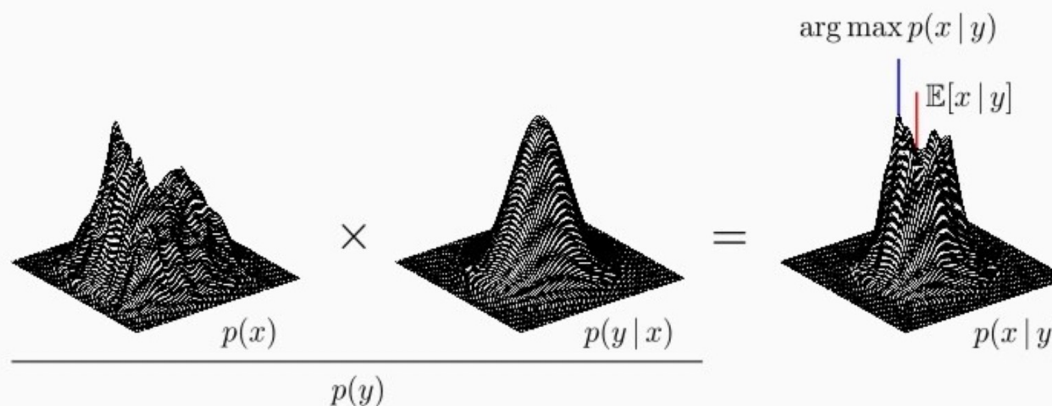
- **Poor performance**
- **Less calibrated uncertainty estimates**
- **Poor scalability**

# The Bayesian viewpoint of deterministic training

## Maximum a Posteriori (MAP)

- Take a mode of the posterior, or **Maximum a Posteriori**

$$x^* \in \operatorname{argmax}_x p(x | y)$$

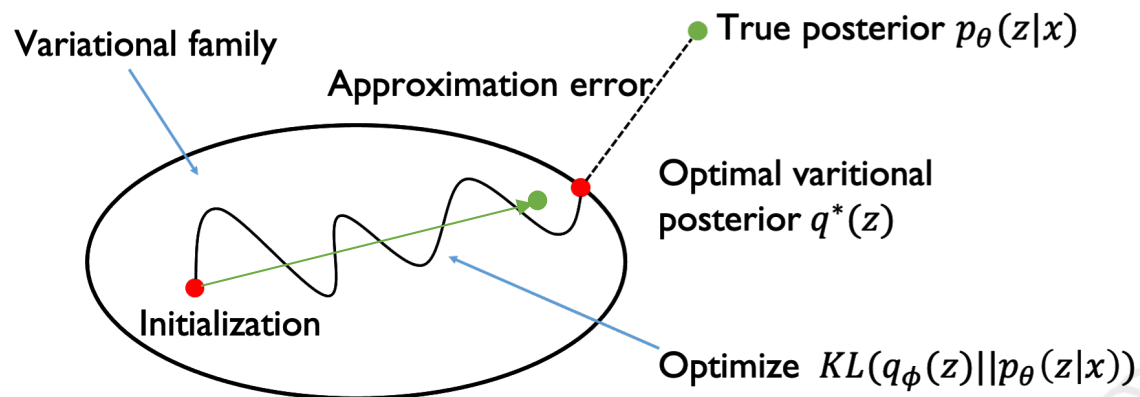


$$\text{MAP: } \max_{\mathbf{w}} \frac{1}{n} \sum_i [\log p(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w})] + \frac{1}{n} \log p(\mathbf{w})$$

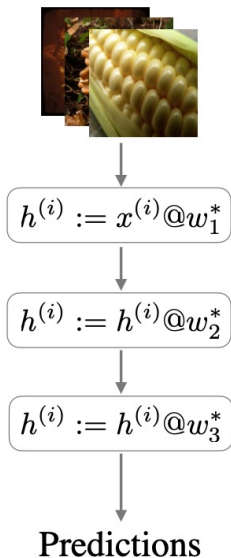
$$\text{Variational inference (VI): } \underbrace{\max_{\theta} \mathbb{E}_{q(\mathbf{w}|\theta)} \left[ \frac{1}{n} \sum_i \log p(y^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}) \right]}_{\mathcal{L}_{ell}} - \underbrace{\frac{1}{n} D_{\text{KL}}(q(\mathbf{w}|\theta) \| p(\mathbf{w}))}_{\mathcal{L}_c}$$

# BayesAdapter: variational inference by *Bayesian fine-tuning*

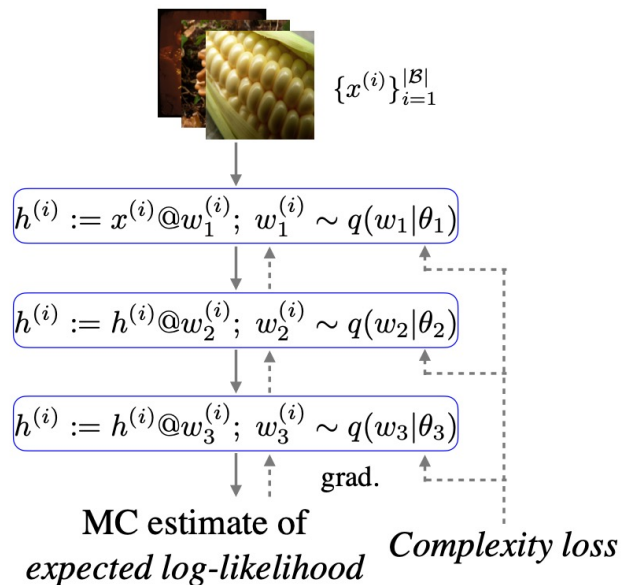
[Deng et al., NeurIPS sub.; Deng et al., CVRP 2021]



Pre-trained DNN

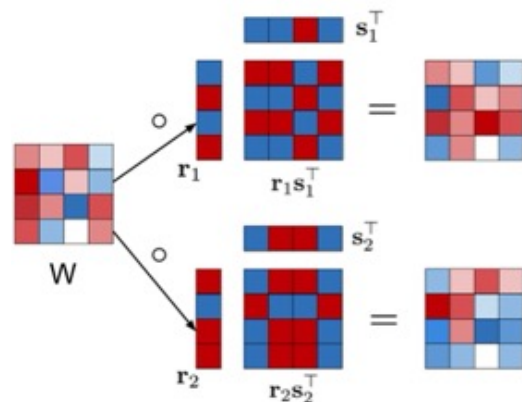
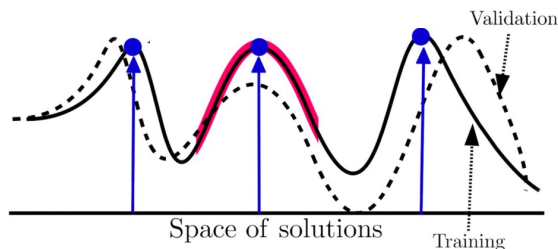


Variational BNN



# BayesAdapter: variational inference by *Bayesian fine-tuning*

[Deng et al., NeurIPS sub.; Deng et al., CVRP 2021]



```
def BayesAdapter_conv(x, theta, stride, padding, groups):
    b = x.shape[0]
    # sample a batch of parameters w: [b, o, i, k, k]
    w = mc_sample(theta, num_mc_samples=b)
    # reshape w to have shape [b*o, i, k, k]
    w = w.flatten(start_dim=0, end_dim=1)
    # reshape x to have shape [1, b*i, h, w]
    x = x.flatten(start_dim=0, end_dim=1).unsqueeze(0)
    # perform b convs in parallel
    y = conv2d(x, w, stride, padding, groups*b)
    # reshape the result to standard format
    return y.view(b, -1, y.shape[2], y.shape[3])
```

- To capture the multi-mode DNN posterior mixture of delta (Gaussian) variational
- To maintain parameter efficiency parameter sharing
- To reduce the variance of stochastic gradients Exemplar reparameterization

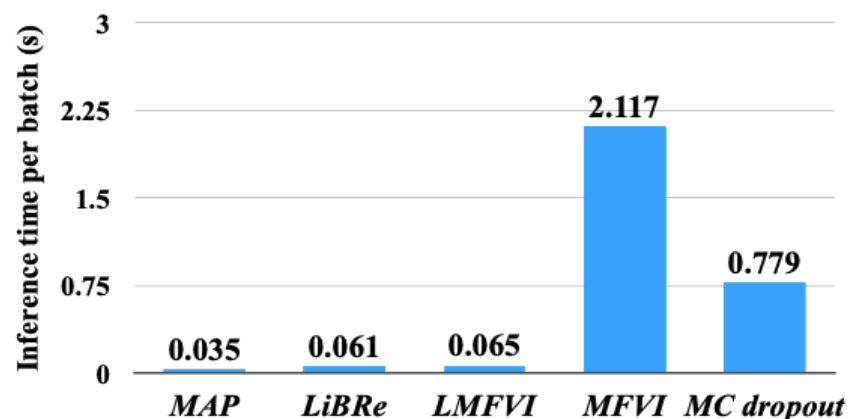
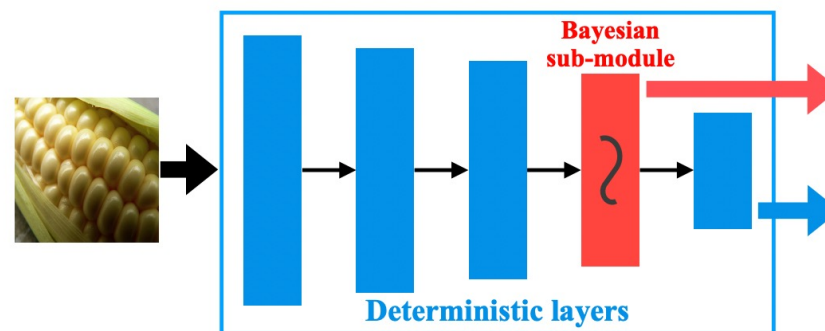


# BayesAdapter: variational inference by *Bayesian fine-tuning*

[Deng et al., NeurIPS sub.; Deng et al., CVRP 2021]

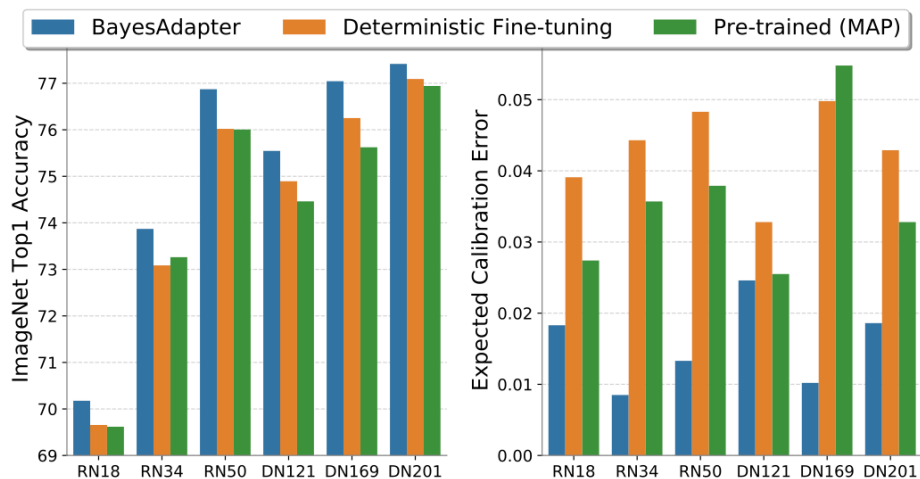
Make Bayesian modeling *lightweight*

- A practical and theoretical sound BDL approach
- Need minimal added training cost
- Promising Bayesian model average speed

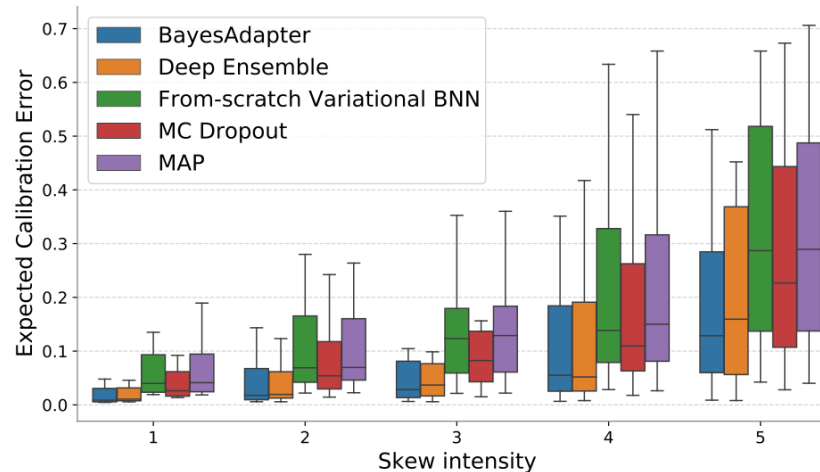


(a) Inference speed comparison

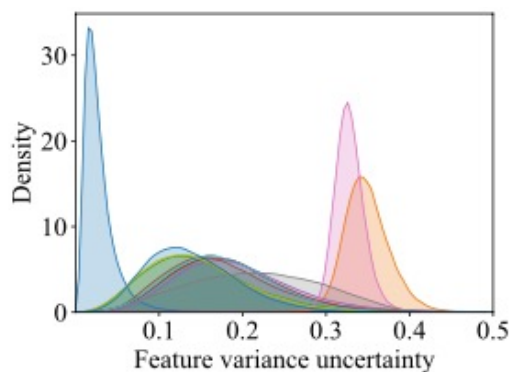
# BayesAdapter: results



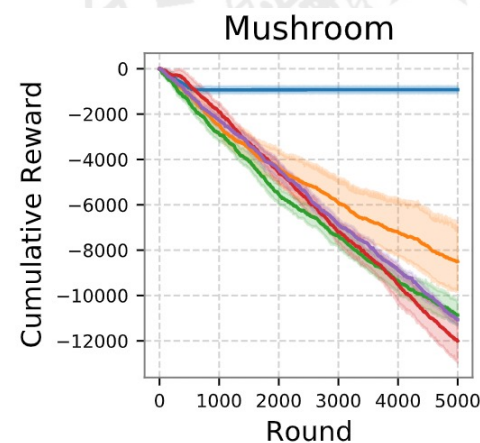
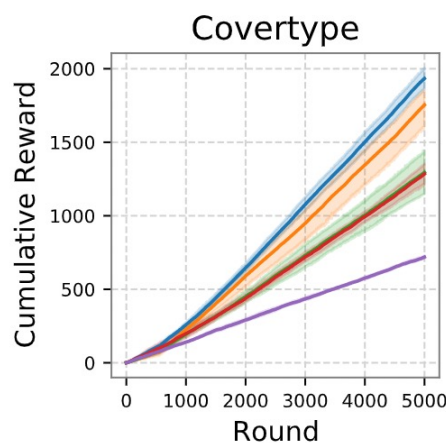
Bayesian model ensemble: one of the **first** variational BNNs that beat DNNs on *ImageNet*



OOD robustness (**resistance to over-confidence**)



Uncertainty based **detection of adversarial examples**



**Exploration** in reinforcement learning (contextual bandit)

# BayesAdapter: a Python library

thudzj / ScalableBDL

<> Code Pull requests Actions Projects Wiki Security Insights Settings

master 3 branches 0 tags Go to file Add file Code

thudzj Update readme	243b8b8 yesterday	82 commits
docs	Update bib.txt	9 months ago
reproduction	Update finetune_imagenet.py	9 months ago
scalablebdl	Update readme	yesterday
.gitignore	Release 0.0	9 months ago
README.md	Update readme	yesterday
demo.py	Update readme	yesterday
license.txt	v0.0.1	9 months ago
requirements.txt	U	9 months ago
setup.py	Update setup.py	9 months ago

☰ README.md

**A plug-and-play implementation for *Bayesian fine-tuning* to practically learn Bayesian Neural Networks**



# Uncertainty over the DNN structure?

Is there a more scalable alternative to the weight uncertainty?

## On weights

- Hard to specifying **sensible priors**
- Using flexible variational posterior for high-dim weights is **expensive**
- Over parameterization nature of DNNs may lead to **degenerated weight posterior**

VS

## On structure

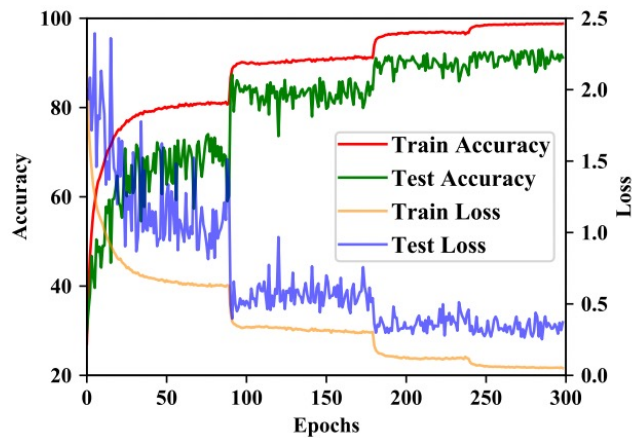
- Impose prior beliefs more **explicitly**
- As shown by NAS, the network structure can be defined in a **compact** manner
- Learning network structure can **boost performance**



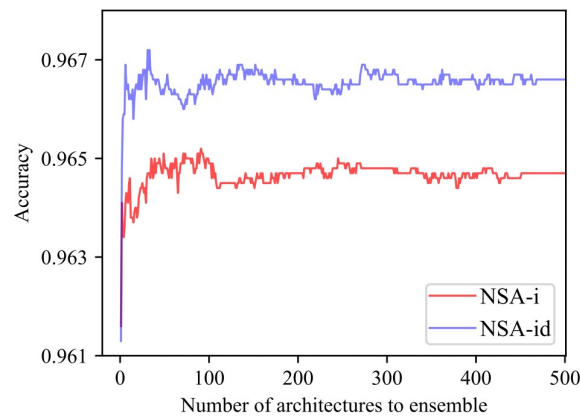
# Structure uncertainty: a new BNN paradigm

Deng et al., NeurIPS 2020; Deng et al., ICLR 2021 NAS workshop

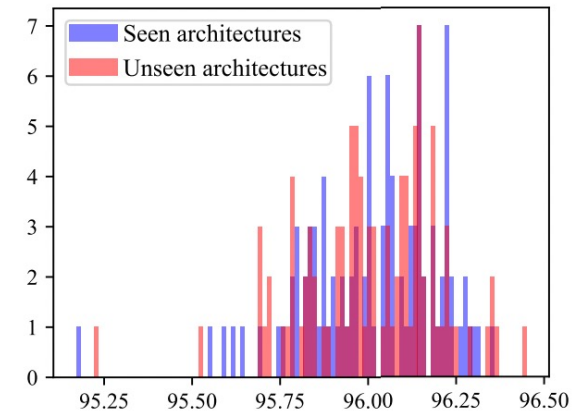
- We pre-specify the structural uncertainty and perform a **first** investigation/understanding on DNNs with such structure uncertainty



training/test disparity



function mode collapse



generalization

# Structure uncertainty: a new BNN paradigm

Deng et al., NeurIPS 2020; Deng et al., ICLR 2021 NAS workshop

Structure uncertainty meets the advance in NAS

- Assume priors and define variationals:

$$p(\alpha, w) = p(\alpha)p(w) \quad q(\alpha, w) = q(\alpha|\theta)\delta(w - w_0)$$

$$p(\alpha) = \prod_{i < j} p(\alpha^{(i,j)}) \quad q(\alpha|\theta) = \prod_{i < j} q(\alpha^{(i,j)}|\theta^{(i,j)})$$

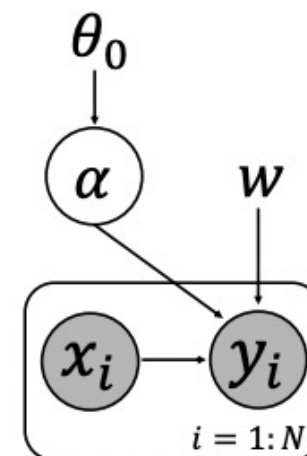
- A unified training objective (**ELBO**):

$$\min_{q \in \mathcal{Q}} D_{\text{KL}}(q(\alpha, w) \| p(\alpha, w | \mathcal{D})) = -\mathbb{E}_{q(\alpha, w)}[\log p(\mathcal{D} | \alpha, w)] + D_{\text{KL}}(q(\alpha, w) \| p(\alpha, w)) + \text{constant}.$$

- Continuous relaxation and reparameterization

$$\alpha^{(i,j)} = g(\theta^{(i,j)}, \beta^{(i,j)}, \epsilon^{(i,j)}) = \text{softmax}((\theta^{(i,j)} + \beta^{(i,j)} \epsilon^{(i,j)}) / \tau).$$

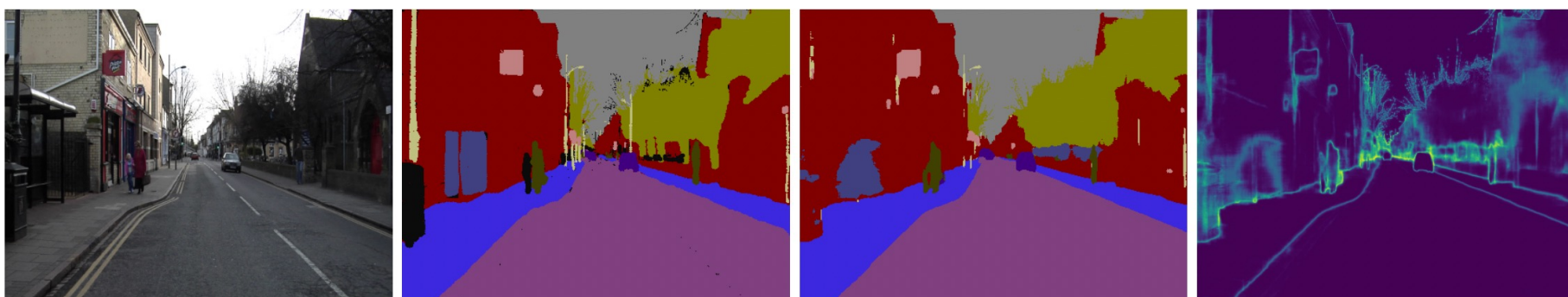
- “Cold posterior”: sharpened concrete distribution



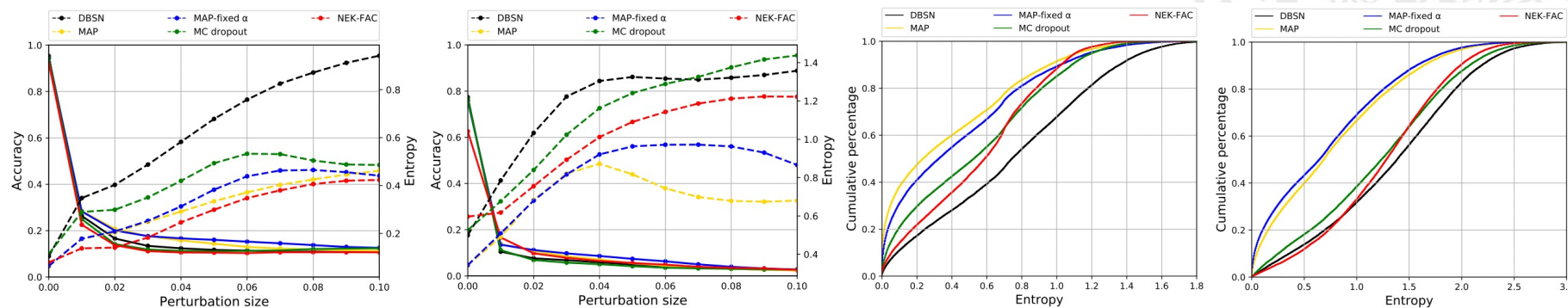
# Structure uncertainty: results

Method	DBSN	MAP	MAP-fixed $\alpha$	MC dropout	BBB	FBN	NEK-FAC
<b>CIFAR-10</b>	<b>0.0109</b>	0.0339	0.0327	0.0150	0.0745	0.0966	0.0434
<b>CIFAR-100</b>	<b>0.0599</b>	0.1240	0.1259	0.0617	0.0700	0.1091	0.1665

Less over-confidence than BNNs with weight uncertainty



Meaningful uncertainty estimates in semantic segmentation problems

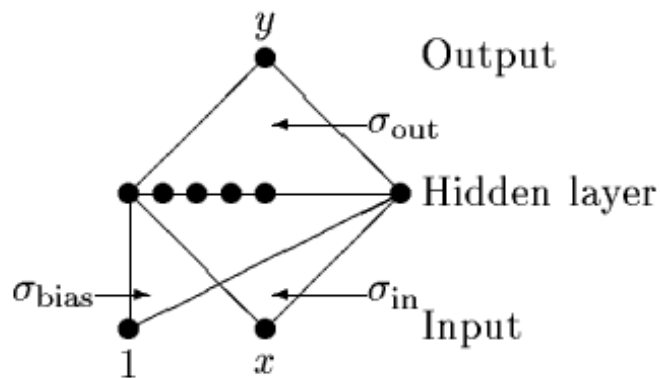


Better OOD robustness

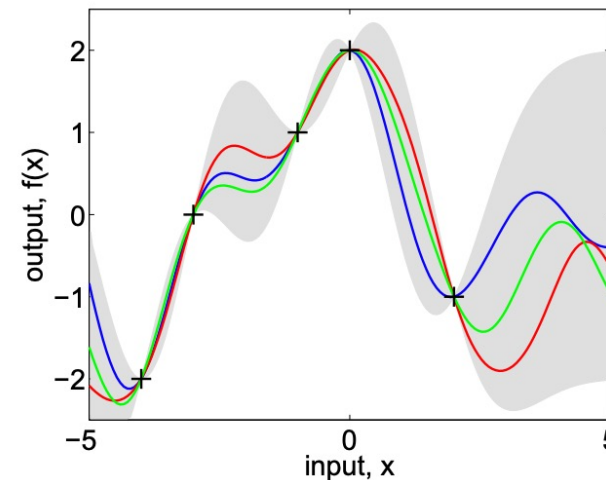
# BNNs are Gaussian processes (GPs) in the width limit

## An exciting perspective

Neal, 1995; Lee et al., 2017



$width \rightarrow \infty$



$$\mathbf{f}_* | X, \mathbf{y}, X_* \sim \mathcal{N}(\bar{\mathbf{f}}_*, \text{cov}(\mathbf{f}_*)), \text{ where}$$

$$\bar{\mathbf{f}}_* \triangleq \mathbb{E}[\mathbf{f}_* | X, \mathbf{y}, X_*] = K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} \mathbf{y},$$

$$\text{cov}(\mathbf{f}_*) = K(X_*, X_*) - K(X_*, X)[K(X, X) + \sigma_n^2 I]^{-1} K(X, X_*)$$

Allows us to *understand* neural networks (e.g. generalization properties) without practically training them

# Can deep ensemble be understood in this spirit?

## Deep ensemble:

- One of the **most performant** prediction & uncertainty modeling approaches
- **Lack** a proper **Bayesian** justification

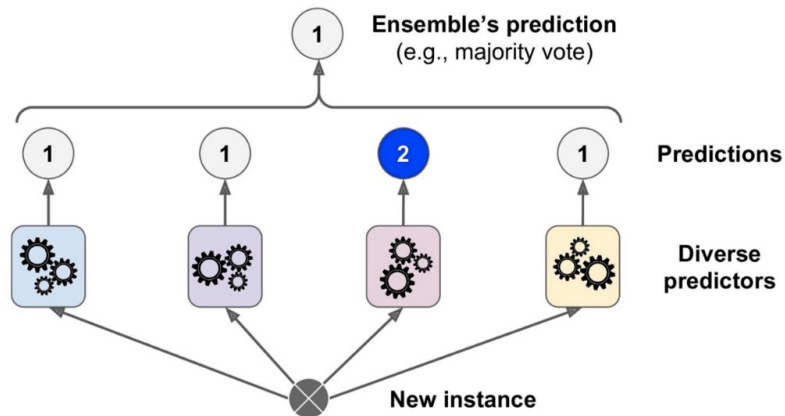


Figure 7-2. Hard voting classifier predictions



# Deep ensemble defines a GP posterior

Deng et al., NeurIPS sub.

**The form:**  $q(f|\mathbf{w}_1, \dots, \mathbf{w}_M) = \mathcal{GP}(f|m_q(\mathbf{x}), k_q(\mathbf{x}, \mathbf{x}')),$

$$m_q(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M g(\mathbf{x}, \mathbf{w}_i),$$

$$k_q(\mathbf{x}, \mathbf{x}') = \frac{1}{M} \sum_{i=1}^M (g(\mathbf{x}, \mathbf{w}_i) - m_q(\mathbf{x})) (g(\mathbf{x}', \mathbf{w}_i) - m_q(\mathbf{x}'))^T + \lambda \mathbf{I}_C.$$



# Deep ensemble defines a GP posterior

Deng et al., NeurIPS sub.

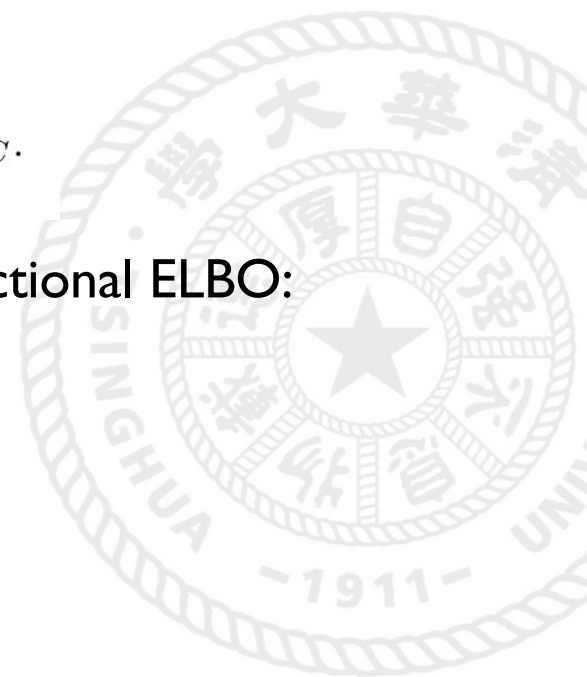
**The form:**  $q(f|\mathbf{w}_1, \dots, \mathbf{w}_M) = \mathcal{GP}(f|m_q(\mathbf{x}), k_q(\mathbf{x}, \mathbf{x}')),$

$$m_q(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M g(\mathbf{x}, \mathbf{w}_i),$$

$$k_q(\mathbf{x}, \mathbf{x}') = \frac{1}{M} \sum_{i=1}^M (g(\mathbf{x}, \mathbf{w}_i) - m_q(\mathbf{x})) (g(\mathbf{x}', \mathbf{w}_i) - m_q(\mathbf{x}'))^T + \lambda \mathbf{I}_C.$$

Bayesian inference in function space: **theorem** on the functional ELBO:

$$\begin{aligned} \mathcal{L}' &= \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}} \mathbb{E}_{q(f)}[\log p(\mathbf{y}_i | f(\mathbf{x}_i))] - D_{\text{KL}}[q(\mathbf{f}^{\tilde{\mathbf{X}}}) \| p(\mathbf{f}^{\tilde{\mathbf{X}}})] \\ &= \log p(\mathcal{D}) - D_{\text{KL}}[q(\mathbf{f}^{\tilde{\mathbf{X}}}) \| p(\mathbf{f}^{\tilde{\mathbf{X}}} | \mathcal{D})] \leq \log p(\mathcal{D}), \end{aligned}$$



# Deep ensemble defines a GP posterior

Deng et al., NeurIPS sub.

The form:  $q(f|\mathbf{w}_1, \dots, \mathbf{w}_M) = \mathcal{GP}(f|m_q(\mathbf{x}), k_q(\mathbf{x}, \mathbf{x}')),$

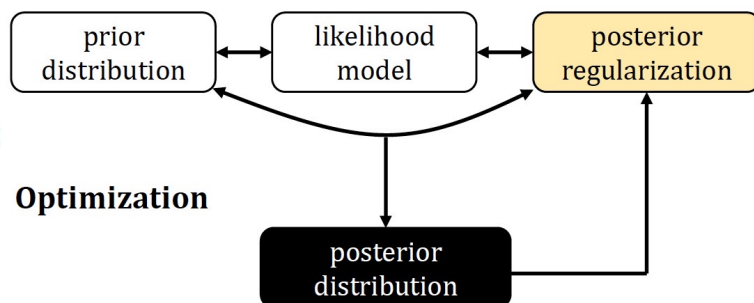
$$m_q(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M g(\mathbf{x}, \mathbf{w}_i),$$

$$k_q(\mathbf{x}, \mathbf{x}') = \frac{1}{M} \sum_{i=1}^M (g(\mathbf{x}, \mathbf{w}_i) - m_q(\mathbf{x})) (g(\mathbf{x}', \mathbf{w}_i) - m_q(\mathbf{x}'))^T + \lambda \mathbf{I}_C.$$

Bayesian inference in function space: **theorem** on the functional ELBO:

$$\begin{aligned} \mathcal{L}' &= \sum_{(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}} \mathbb{E}_{q(f)}[\log p(\mathbf{y}_i|f(\mathbf{x}_i))] - D_{\text{KL}}[q(\mathbf{f}^{\tilde{\mathbf{X}}})||p(\mathbf{f}^{\tilde{\mathbf{X}}})] \\ &= \log p(\mathcal{D}) - D_{\text{KL}}[q(\mathbf{f}^{\tilde{\mathbf{X}}})||p(\mathbf{f}^{\tilde{\mathbf{X}}}|\mathcal{D})] \leq \log p(\mathcal{D}), \end{aligned}$$

RegBayes:

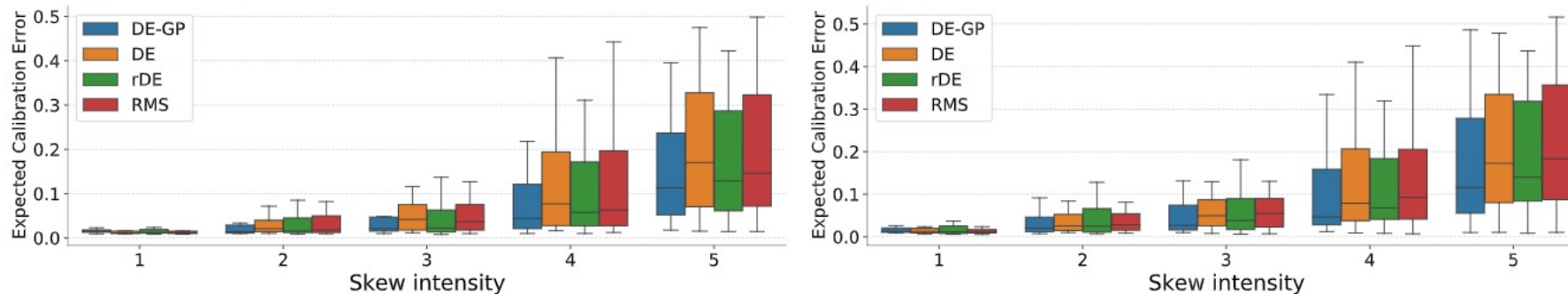
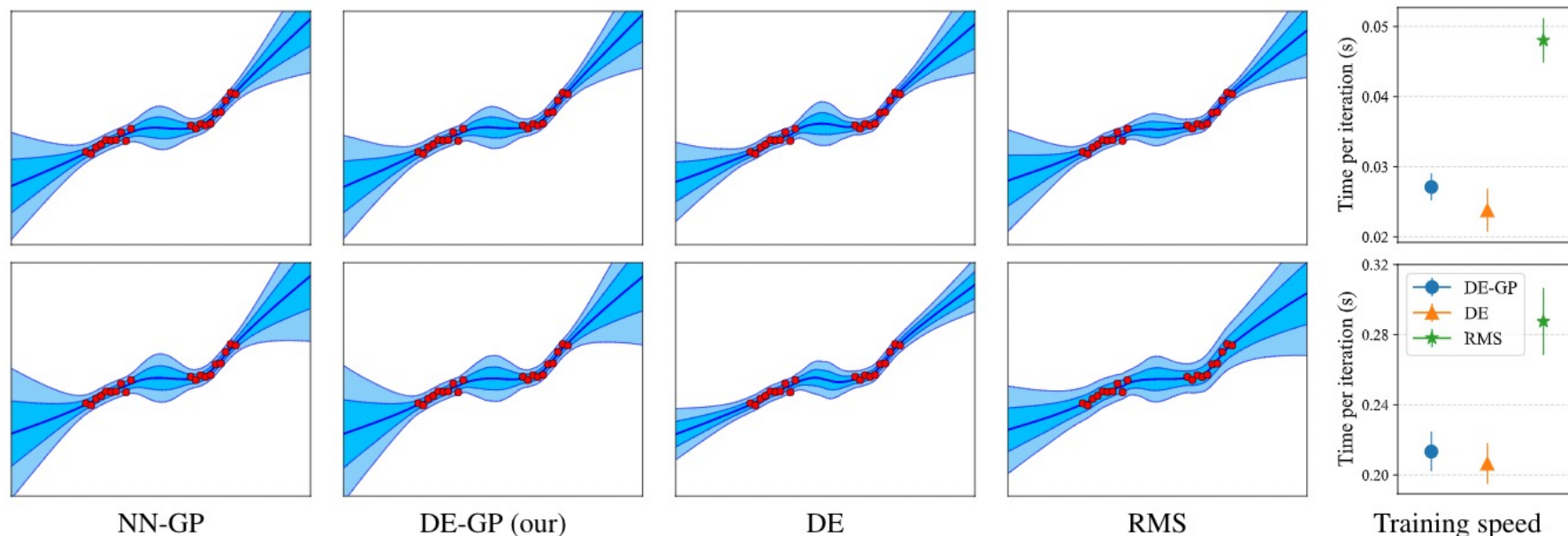


One can encode *any differentiable constraints* on the functional posterior



# Deep ensemble defines a GP posterior: results

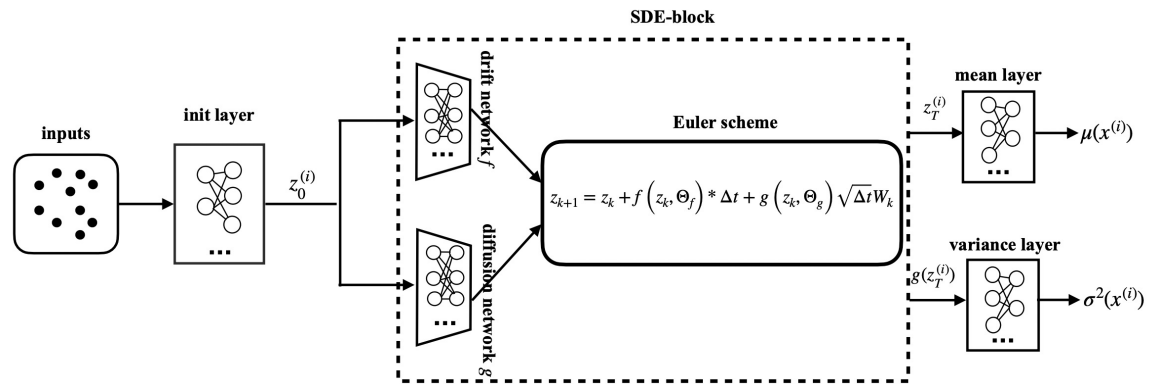
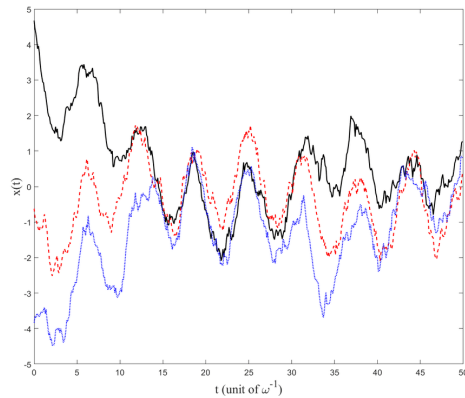
Deng et al., NeurIPS sub.



More calibrated/reliable uncertainty estimates than standard deep ensemble

# Uncertainty quantification methods beyond BNNs

## NeurIPS sub.

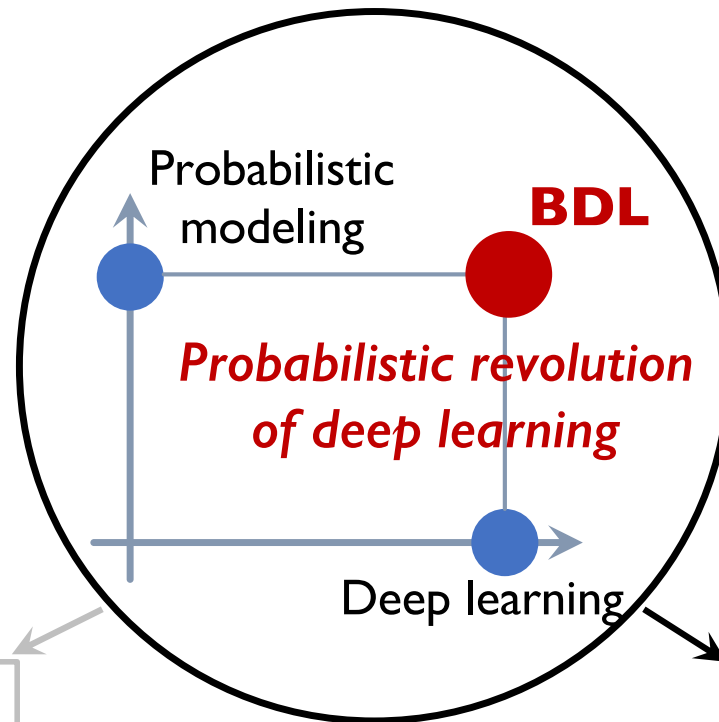


Stochastic differential equations (SDEs)

SDE based heteroscedastic neural networks

Dataset	Metric	MCD	DGP	BNN	Deep-ens	HNN	Proposed
Metro-traffic	RMSE	697.021	651.341	786.694	533.426	559.354	<b>483.639±2.657</b>
	$R^2 \uparrow$	0.877	0.892	0.843	0.928	0.920	<b>0.939±0.011</b>
	CWCE	52.152	10.552	21.486	9.078	9.305	<b>2.894±0.085</b>
	EPIW	<del>167.859</del>	1168.044	610.662	814.143	883.475	<b>539.254±19.334</b>
	R-CWCE	6.428	1.136	3.373	0.655	0.747	<b>0.177±0.014</b>
Pickups	RMSE	625.812	523.041	720.013	428.032	421.752	<b>340.331±5.072</b>
	$R^2 \uparrow$	0.878	0.914	0.838	0.943	0.945	<b>0.964±0.012</b>
	CWCE	34.441	22.799	42.570	4.878	6.043	<b>2.925±0.758</b>
	EPIW	<del>313.432</del>	1872.481	247.229	684.381	688.989	<b>438.324±19.222</b>
	R-CWCE	4.205	1.951	6.904	0.280	0.335	<b>0.173±0.012</b>

# Research focus: Bayesian deep learning (BDL)



Theme I: Uncertainty quantification via BNNs

(NeurIPS 2020; CVPR 2021; ICLR 2021 workshop; NeurIPS sub.)

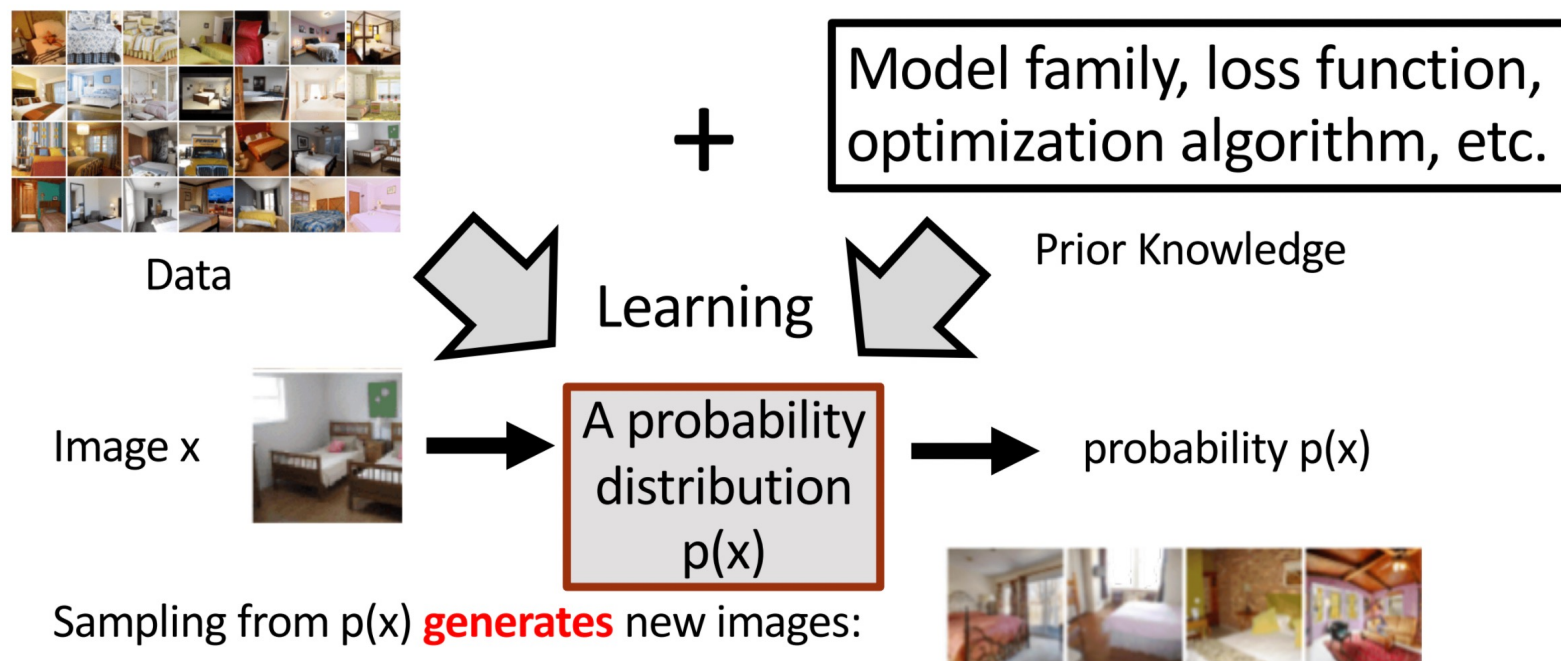
Theme II: Knowledge incorporation via DGMs

(NeurIPS 2017; ICCV 2019, NeurIPS 2020; ICML 2019 workshop)

# Deep generative models (DGMs)

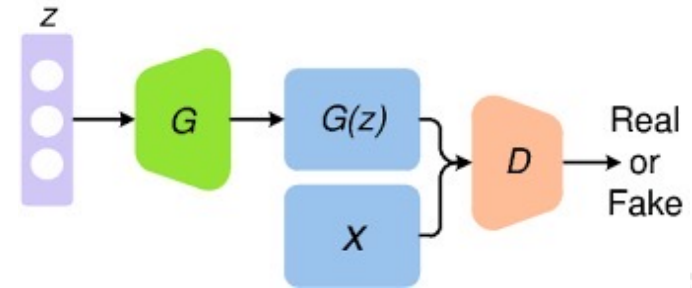
## DNNs enrich probabilistic modeling

Richard Feynman: “What I cannot create, I do not understand”



# Conditional Generative Adversarial Nets (GANs)

Generative models with implicit density



- GANs -- a two-player minimax game:

$$\min_G \max_D \mathcal{L}(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log(D(\mathbf{x}))] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [1 - \log(D(G(\mathbf{z})))]$$

(Minimizing Jensen-Shannon divergence)

- cGANs – label-aware GANs:  $\mathbf{z}$  - noise,  $\mathbf{y}$  - label

$$\min_G \max_{D_{xy}} \mathcal{L}_{xy} = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\mathbf{x}, \mathbf{y})} [\log(D_{xy}(\mathbf{x}, \mathbf{y}))] + \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}), \mathbf{z} \sim p(\mathbf{z})} [\log(1 - D_{xy}(G(\mathbf{y}, \mathbf{z}), \mathbf{y}))]$$

Modeling the **joint** between data and label

# Conditional generative modeling with few labels is non-trivial

## SSL meets cGANs

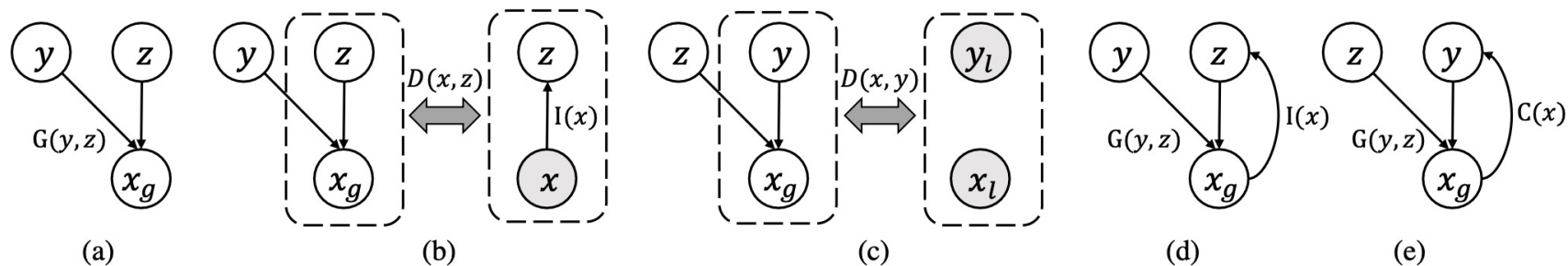
- The conditional generators in existing works exhibit inadequate **controllability** – the generator’s ability to conditionally generate samples that have structures strictly agreeing with the condition



- Reason: noise  $z$  encodes some semantic info., confounding  $G$
- Solution: **disentangle** the semantics of our interest and other variations

# Structured GANs: cGANs with a structured hidden space

Deng et al., NeurIPS 2017



## A simple prior knowledge

$$\left. \begin{aligned} & \min_{C,G} \mathbb{E}_{y \sim p(y)} \|p_c(y|G(y, z_1)), p_c(y|G(y, z_2))\|, \forall z_1, z_2 \sim p(z) \\ & \min_{I,G} \mathbb{E}_{z \sim p(z)} \|p_i(z|G(y_1, z)), p_i(z|G(y_2, z))\|, \forall y_1, y_2 \sim p(y) \end{aligned} \right\} \text{Implemented by optimizing reconstruction error in hidden space}$$

## Adversarial games for aligning joint distribution

$$\left. \begin{aligned} & \min_G \max_{D_{xy}} \mathcal{L}_{xy} = \mathbb{E}_{(x,y) \sim p(x,y)} [\log(D_{xy}(x, y))] + \mathbb{E}_{y \sim p(y), z \sim p(z)} [\log(1 - D_{xy}(G(y, z), y))] \\ & \min_{I,G} \max_{D_{xz}} \mathcal{L}_{xz} = \mathbb{E}_{x \sim p(x)} [\log(D_{xz}(x, I(x)))] + \mathbb{E}_{z \sim p(z), y \sim p(y)} [\log(1 - D_{xz}(G(y, z), z))] \end{aligned} \right\} \text{Adversarial training}$$

## Main theorem: unbiased equilibrium

**Theorem 3.3** *Minimizing  $\mathcal{R}_z$  w.r.t.  $I$  will keep the equilibrium of the adversarial game  $\mathcal{L}_{xz}$ . Similarly, minimizing  $\mathcal{R}_y$  w.r.t.  $C$  will keep the equilibrium of the adversarial game  $\mathcal{L}_{xy}$  unchanged.*

# Structured GANs: results

Deng et al., NeurIPS 2017

Method	MNIST			SVHN	CIFAR-10
	$n = 20$	$n = 50$	$n = 100$	$n = 1000$	$n = 4000$
Ladder [22]	-	-	<b>0.89(±0.50)</b>	-	20.40(±0.47)
VAE [12]	-	-	3.33(±0.14)	36.02(±0.10)	-
CatGAN [28]	-	-	1.39(±0.28)	-	19.58(±0.58)
ALI [5]	-	-	-	7.3	18.3
ImprovedGAN [27]	16.77(±4.52)	2.21(±1.36)	0.93 (±0.07)	8.11(±1.3)	18.63(±2.32)
TripleGAN [15]	5.40(±6.53)	1.59(±0.69)	0.92(±0.58)	5.83(±0.20)	18.82(±0.32)
SGAN	<b>4.0(±4.14)</b>	<b>1.29(±0.47)</b>	<b>0.89(±0.11)</b>	<b>5.73(±0.12)</b>	<b>17.26(±0.69)</b>

Impressive  
SSL  
classification  
accuracy

Table 2: Comparisons of semi-supervised classification errors (%) on MNIST, SVHN and CIFAR-10 test sets.



Fixed  
label  
in each  
row

Fixed style in each column



(a) MNIST

(b) SVHN



(c) CIFAR-10

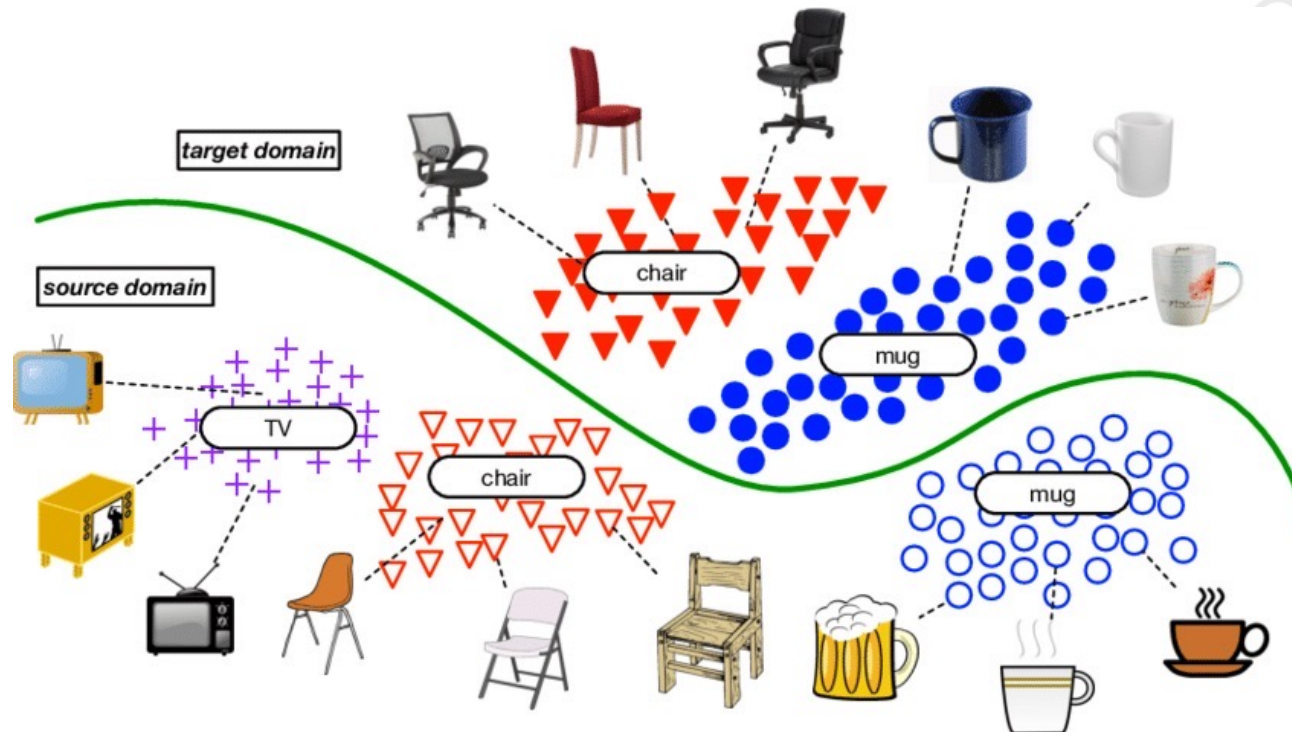
Style  
transfer



# A more extreme discriminative learning scenario: UDA

*Unsupervised domain adaptation (UDA):*

- the concerned domain (target domain) is **unlabeled**. We have only access to labeled data from a **related domain** (source domain)



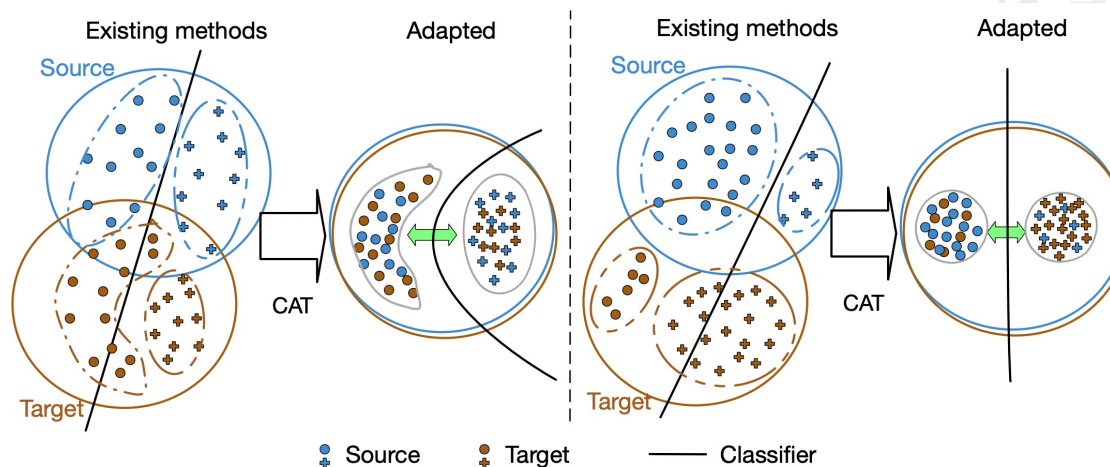
# Marginal distribution alignment is not inadequate for UDA

The generalization bound for UDA

$$\begin{aligned} \epsilon_t(h) &\leq \epsilon_s(h) + \frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(s, t) \\ &\quad + \min_{\hat{h} \in \mathcal{H}} (\epsilon_s(\hat{h}, l_s) + \epsilon_t(\hat{h}, l_t)) \\ &\leq \epsilon_s(h) + \boxed{\frac{1}{2}d_{\mathcal{H}\Delta\mathcal{H}}(s, t)} + \boxed{\epsilon_t(l_s, l_t)} \\ &\quad + \min_{\hat{h} \in \mathcal{H}} (\epsilon_s(\hat{h}, l_s) + \epsilon_t(\hat{h}, l_t)). \end{aligned}$$

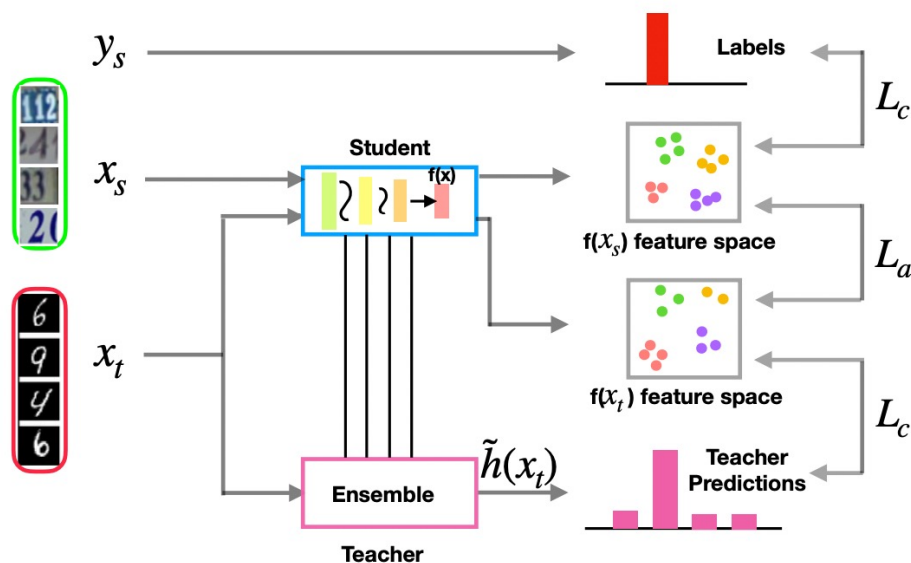
Mismatch between marginal dist.

Mismatch between labeling functions (class-conditional dist.)



# Cluster alignment with a teacher for UDA

Deng et al., ICCV 2019



$$\min_{\theta} \mathcal{L}_c(\mathcal{X}_s, \mathcal{X}_t) = \mathcal{L}_c(\mathcal{X}_s) + \mathcal{L}_c(\mathcal{X}_t),$$

$$\mathcal{L}_c(\mathcal{X}) = \frac{1}{|\mathcal{X}|^2} \sum_{i=1}^{|\mathcal{X}|} \sum_{j=1}^{|\mathcal{X}|} [\delta_{ij} d(h(x^i), h(x^j)) + (1 - \delta_{ij}) \max(0, m - d(h(x^i), h(x^j)))] ,$$

$$\min_{\theta} \mathcal{L}_a(\mathcal{X}_s, \mathcal{Y}_s, \mathcal{X}_t) = \frac{1}{K} \sum_{k=1}^K [d(\lambda_{s,k}, \lambda_{t,k})]$$

$$\min_{\theta} \max_{\phi} \mathcal{L}_{cd}(\mathcal{X}_s, \mathcal{X}_t) = \frac{1}{N} \sum_{i=1}^N [\log c(h(x_s^i; \theta); \phi)] +$$

$$\frac{1}{\tilde{M}} \sum_{i=1}^{\tilde{M}} [\log(1 - c(h(x_t^i; \theta); \phi)) \gamma_i]$$

Distribution alignment with **class-conditional structure awareness**:

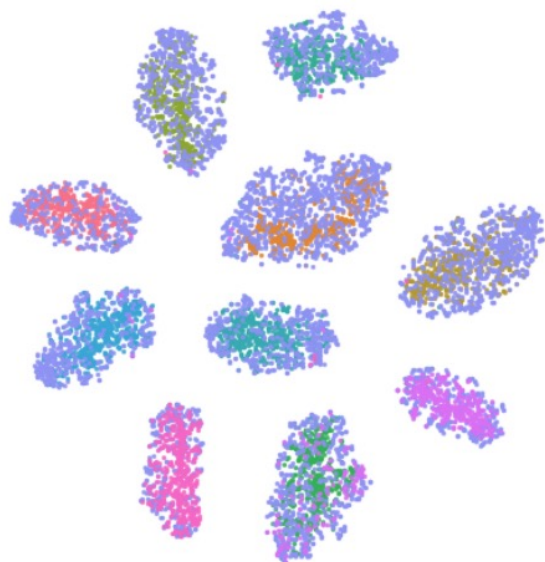
- implement the **cluster assumption** of discriminative learning

# Cluster alignment with a teacher for UDA: results

Deng et al., ICCV 2019

Method	SVHN to MNIST	MNIST to USPS	USPS to MNIST
<b>RevGrad</b> [7]	27.4 ± 6.3	26.7 ± 2.0	17.9 ± 1.4
<b>MSTN</b> [49]	25.8 ± 3.6	30.3 ± 1.0	29.4 ± 0.5
<b>CAT</b>	<b>100.0 ± 0.05</b>	<b>100.0 ± 0.0</b>	<b>99.9 ± 0.2</b>

*Epecially effective for class imbalanced tasks*



**Separated** clusters in the feature space

Method	SVHN to MNIST	MNIST to USPS	USPS to MNIST
<b>Source Only</b>	60.1 ± 1.1	75.2 ± 1.6	57.1 ± 1.7
<b>DDC</b> [45]	68.1 ± 0.3	79.1 ± 0.5	66.5 ± 3.3
<b>CoGAN</b> [20]	-	91.2 ± 0.8	89.1 ± 0.8
<b>DRCN</b> [8]	82.0 ± 0.1	91.8 ± 0.09	73.7 ± 0.04
<b>ADDA</b> [44]	76.0 ± 1.8	89.4 ± 0.2	90.1 ± 0.8
<b>LEL</b> [26]	81.0 ± 0.3	-	-
<b>AssocDA</b> [11]	97.6	-	-
<b>MSTN</b> [49]	91.7 ± 1.5	92.9 ± 1.1	-
<b>CAT</b>	98.1 ± 1.3	90.6 ± 2.3	80.9 ± 3.1
<b>RevGrad</b> [7]	73.9	77.1 ± 1.8	73.0 ± 2.0
<b>RevGrad+CAT</b>	98.0 ± 0.8	93.7 ± 1.1	95.7 ± 1.3
<b>rRevGrad+CAT</b>	<b>98.8 ± 0.02</b>	94.0 ± 0.7	<b>96.0 ± 0.9</b>
<b>MCD</b> [37]	96.2 ± 0.4	94.2 ± 0.7	94.1 ± 0.3
<b>MCD+CAT</b>	97.1 ± 0.2	<b>96.3 ± 0.5</b>	95.2 ± 0.4
<b>VADA</b> [41]	94.5	-	-
<b>VADA+CAT</b>	95.2	-	-

*SOTA UDA performance*

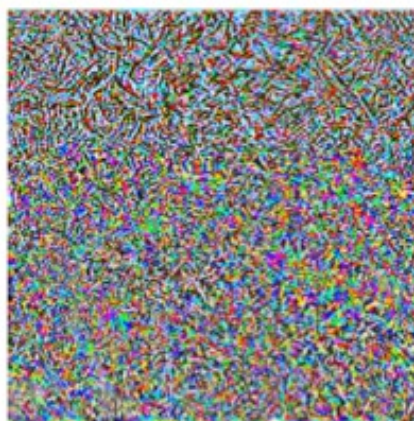
# DNNs are vulnerable to adversarial examples

Clean images

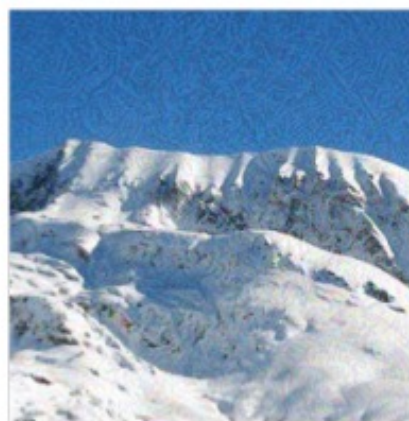


Alps: 94.39%

Adversarial noise



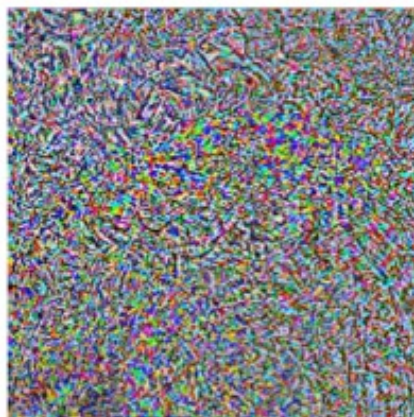
Adversarial examples



Dog: 99.99%



Puffer: 97.99%



Crab: 100.00%



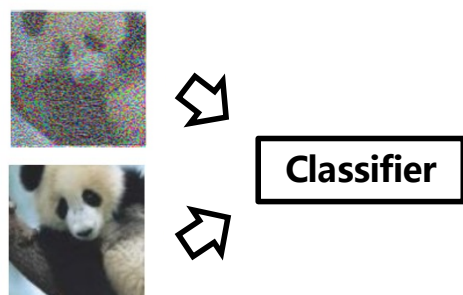
Dong et al. 2018

*What is the underlying distribution of adversarial examples?*

# Modeling adversarial distribution may be helpful

## In the sense of improving adversarially robustness

Adversarial training (AT):



Outer minimization: train a robust classifier

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{\delta_i \in \mathcal{S}} L(f_{\theta}(x_i + \delta_i), y_i)$$

Inner maximization: generate an adversarial example

**Generalization issue** of AT under point-estimate attacker

Model	$\mathcal{A}_{\text{nat}}$	FGSM	PGD-20	PGD-100	MIM	C&W	FeaAttack	$\mathcal{A}_{\text{rob}}$
Standard	<b>94.81%</b>	12.05%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
AT <sub>FGSM</sub>	93.80%	<b>79.86%</b>	0.12%	0.04%	0.06%	0.13%	0.01%	0.01%
AT <sub>PGD</sub> <sup>†</sup>	87.25%	56.04%	45.88%	45.33%	47.15%	46.67%	46.01%	44.89%
AT <sub>PGD</sub>	86.91%	58.30%	50.03%	49.40%	51.40%	50.23%	50.46%	48.26%
ALP	86.81%	56.83%	48.97%	48.60%	50.13%	49.10%	48.51%	47.90%
FeaScatter	89.98%	<b>77.40%</b>	<b>70.85%</b>	<b>68.81%</b>	<b>72.74%</b>	<b>58.46%</b>	37.45%	<b>37.40%</b>

# Adversarial distributional training

Deng et al., NeurIPS 2020

Outer minimization: train a robust classifier

A probabilistic modeling of heterogeneous adversarial examples

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{p(\delta_i) \in \mathcal{P}} \mathbb{E}_{p(\delta_i)} [L(f_{\theta}(x_i + \delta_i), y_i)] + \lambda H(p(\delta_i))$$

Inner maximization: learn an adversarial distribution

- Under mild assumption, we theoretically prove **iterative optimization** method can still be used for solving the minimax problem

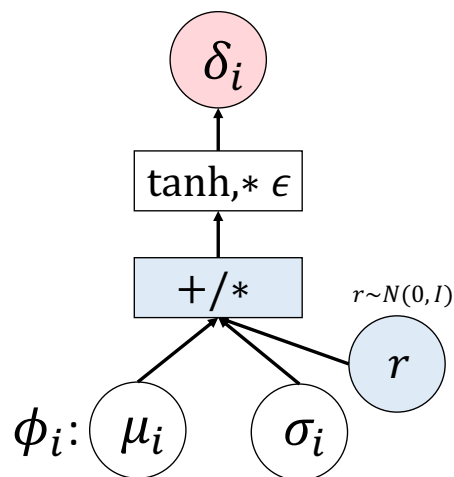
**Theorem 1.** Suppose Assumptions 1 and 2 hold. We define  $\rho(\theta) = \max_{p(\delta_i) \in \mathcal{P}} \mathcal{J}(p(\delta_i), \theta)$ , and  $\mathcal{P}^*(\theta) = \{p(\delta_i) \in \mathcal{P} : \mathcal{J}(p(\delta_i), \theta) = \rho(\theta)\}$ . Then  $\rho(\theta)$  is directionally differentiable, and its directional derivative along the direction  $\mathbf{v}$  satisfies

$$\rho'(\theta; \mathbf{v}) = \sup_{p(\delta_i) \in \mathcal{P}^*(\theta)} \mathbf{v}^{\top} \nabla_{\theta} \mathcal{J}(p(\delta_i), \theta). \quad (6)$$

Particularly, when  $\mathcal{P}^*(\theta) = \{p^*(\delta_i)\}$  only contains one maximizer,  $\rho(\theta)$  is differentiable at  $\theta$  and

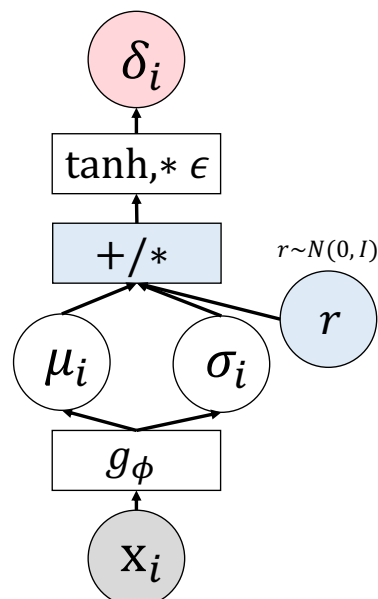
$$\nabla_{\theta} \rho(\theta) = \nabla_{\theta} \mathcal{J}(p^*(\delta_i), \theta). \quad (7)$$

# Use **DGMs** to instantiate the adversarial distributions



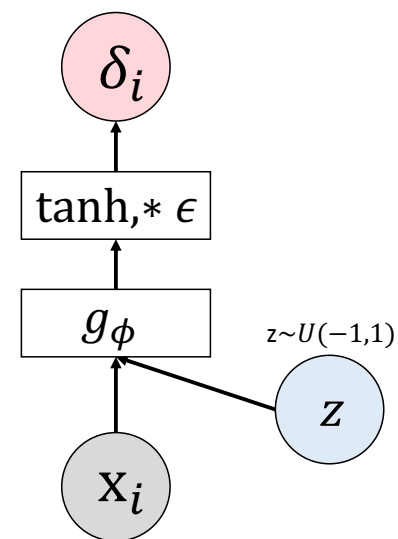
(a)

DGM with  
**explicit** density



(b)

DGM with  
**explicit** density  
(**amortized** version)



(c)

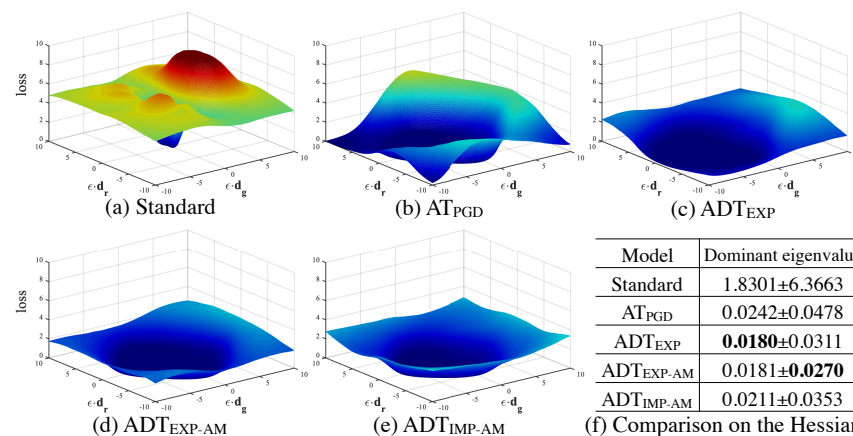
DGM with  
**implicit** density  
(**amortized** version)



# Adversarial distributional training: results



The distribution captures *more diverse modes* of adv. examples



ADT leads to *flatter* loss surfaces

Model	$\mathcal{A}_{\text{nat}}$	FGSM	PGD-20	PGD-100	MIM	C&W	FeaAttack	$\mathcal{A}_{\text{rob}}$
Standard	<b>94.81%</b>	12.05%	0.00%	0.00%	0.00%	0.00%	0.00%	0.00%
AT <sub>FGSM</sub>	93.80%	<b>79.86%</b>	0.12%	0.04%	0.06%	0.13%	0.01%	0.01%
AT <sub>PGD</sub> <sup>†</sup>	87.25%	56.04%	45.88%	45.33%	47.15%	46.67%	46.01%	44.89%
AT <sub>PGD</sub>	86.91%	58.30%	50.03%	49.40%	51.40%	50.23%	50.46%	48.26%
ALP	86.81%	56.83%	48.97%	48.60%	50.13%	49.10%	48.51%	47.90%
FeaScatter	89.98%	<b>77.40%</b>	<b>70.85%</b>	<b>68.81%</b>	<b>72.74%</b>	<b>58.46%</b>	37.45%	<b>37.40%</b>
ADT <sub>EXP</sub>	86.89%	60.41%	52.18%	<b>51.69%</b>	<b>53.27%</b>	52.49%	<b>52.38%</b>	<b>50.56%</b>
ADT <sub>EXP-AM</sub>	87.82%	62.42%	51.95%	51.26%	52.99%	51.75%	52.04%	<b>50.04%</b>
ADT <sub>IMP-AM</sub>	88.00%	<b>64.89%</b>	<b>52.28%</b>	51.23%	52.64%	<b>52.65%</b>	51.89%	<b>49.81%</b>

**Superior** adversarial **robustness** over baselines with clear margins

Thanks!

