

Posterior's Sensitivity to Address AI's Uncertainty

Mohammad Emtiyaz Khan

RIKEN Center for AI Project, Tokyo

<http://emtiyaz.github.io>

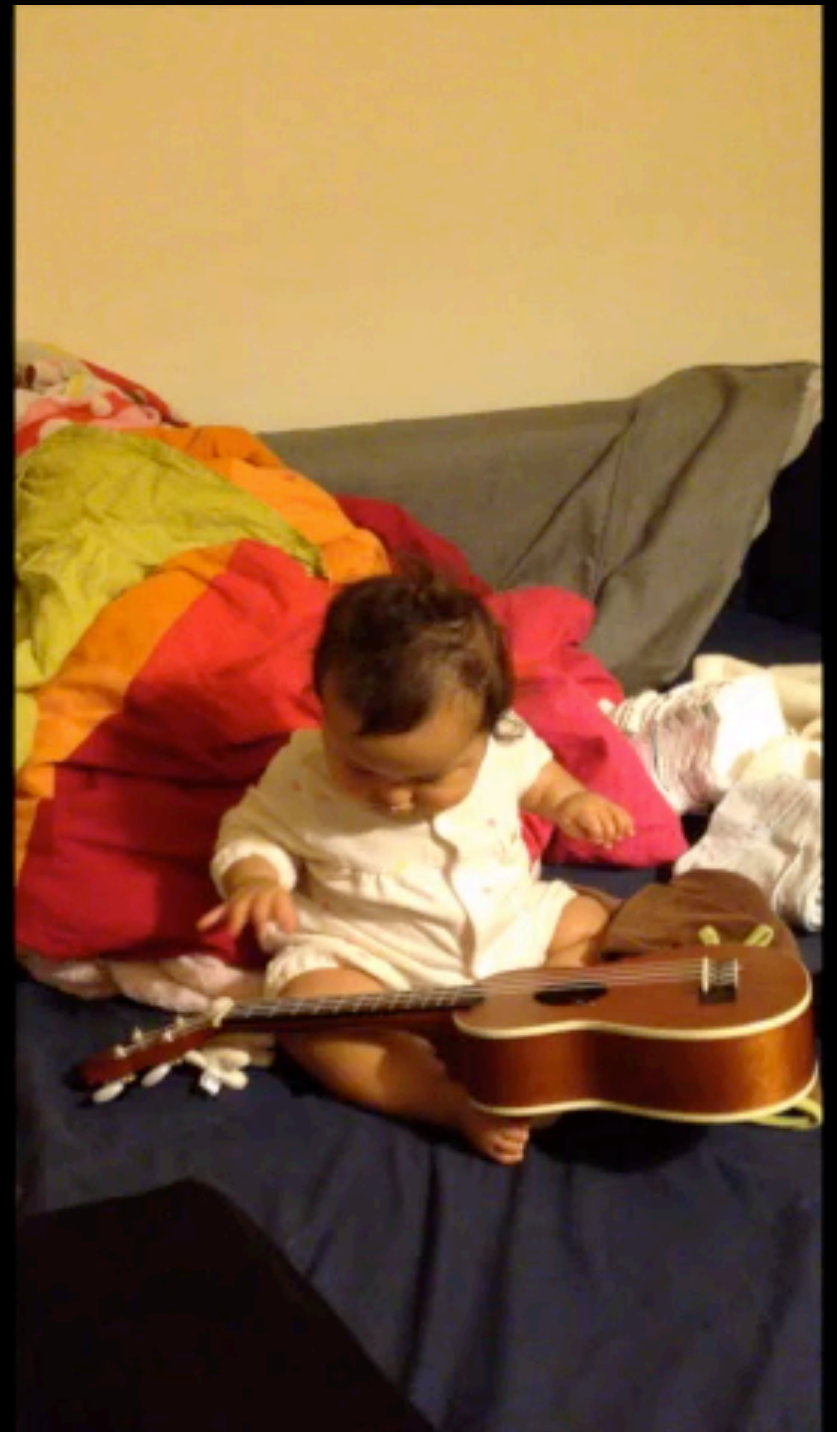


AI that can learn like us

Quickly adapt & continue to acquire new skills.

AI that is low-cost, sustainable, transparent, trustworthy, reliable, composable, modular....

Human Learning at
the age of 6 months.



Converged at the
age of 12 months



Transfer
skills
at the age
of 14
months



Current state of ML



AI that can learn like us

Quickly adapt & continue to acquire new skills.

AI that is low-cost, sustainable, transparent, trustworthy, reliable, composable, modular....

Why haven't we solved it with Bayes?

- In theory, Bayes can solve these problems
 - By using the posterior uncertainty
- But, these are not Bayesian models
- And scale makes it infeasible
- Are there alternatives for Bayes?

Sensitivity and Uncertainty

- Sensitivity of (variational) posteriors to address uncertainty during knowledge transfer
 - Main point: the sensitivity is (essentially) freely available!
- Model sensitivity to data perturbation (addition/removal)
 - Beyond linear regression: conjugate-Bayes [1]
 - Beyond conjugacy [1,2]
 - For large models (VI for GPT-2, ImageNet) [3]
- Model perturbation: LLM model merging [4-5]
 - Federated learning [6] and connections to Bayes-duality

1. Nickl, Xu, Tailor, Moellenhoff, Khan, The memory-perturbation equation, NeurIPS (2023)

2. Khan and Rue, The Bayesian Learning Rule, JMLR (2023).

3. Shen et al. Variational Learning is Effective for Large Deep Networks, ICML (2024)

4. Daheim et al. Model merging by uncertainty-based gradient matching, ICLR (2024).

5. Moldanado et al. How to Weight Multitask Finetuning? Fast Previews via Bayesian Model-Merging, (2024)

6. Swaroop et al. Connecting Federated ADMM to Bayes, ICLR, 2024

How to represent and adapt the knowledge? Perturbation, Sensitivity, and Duality

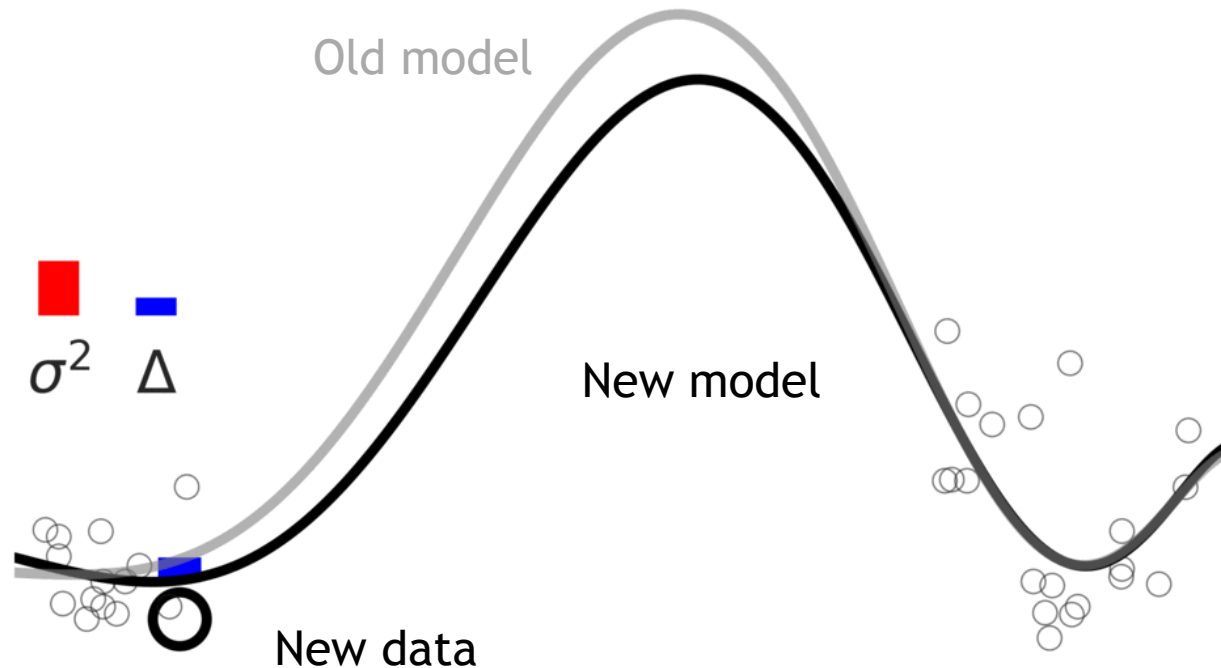


Bayes-Duality

via steampunktendencies.com

<https://tenor.com/view/clockwork-gears-brain-gif-16784329>

Model's Sensitivity to Its Training Data



Model is more sensitive to examples that are “far enough” (in the uncertain territories)

Closed-form Expression for Sensitivity

Linear regression $\ell_i = (y_i - x_i^\top \theta)^2 / 2$

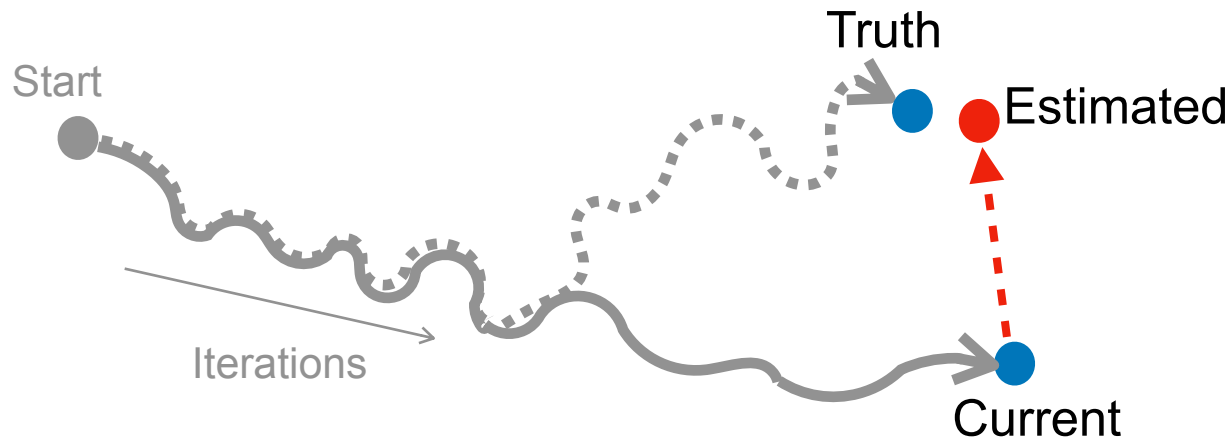
$$\theta_t = \underbrace{H_t^{-1}}_{\text{Hessian}} \sum_{j=1}^t x_j y_j \quad \implies \quad \theta_t - \theta_t^{(i)} = H_t^{-1} x_i (y_i - x_i^\top \theta_t^{(i)})$$
$$x_i^\top (\theta_t - \theta_t^{(i)}) = \underbrace{x_i^\top H_t^{-1} x_i}_{\text{Prediction Variance}} \underbrace{(y_i - x_i^\top \theta_t^{(i)})}_{\text{Prediction Error}}$$
$$= - x_i^\top H_t^{-1} \nabla \ell_i(\theta_t^{(i)})$$

This result is the basis for most works in deep learning [2], but these extensions are too narrow (leave-one-out, at convergence, for data-attribution).

1. Cook. Detection of Influential Observations in Linear Regression. Technometrics. ASA (1977)
2. Koh and Liang. Understanding Black-Box Predictions via Influence Functions. ICML (2017)

A Broader Perspective

- Sensitivity is essential to answer “what-if” questions
- Data Perturbation: What if we add/remove a class? All NY times articles? Continual/active learning
- Model Perturbation: What if we merge separately fine-tuned LLMs? Federated/distributed learning
- Algorithm perturbation, etc. etc.





Peter Nickl



Lu Xu



Dharmesh
Tailor



Thomas
Moellenhoff

Memory-Perturbation

Broadening data-attribution by
using posterior-sensitivity

Exponential Family

Natural
parameters

Sufficient
Statistics

Expectation
parameters

$$q(\theta) \propto \exp \left[\lambda^\top T(\theta) \right]$$

$$\mu := \mathbb{E}_q[T(\theta)]$$

$$\begin{aligned} \mathcal{N}(\theta|m, S^{-1}) &\propto \exp \left[-\frac{1}{2}(\theta - m)^\top S(\theta - m) \right] \\ &\propto \exp \left[(Sm)^\top \theta + \text{Tr} \left(-\frac{S}{2} \theta \theta^\top \right) \right] \end{aligned}$$

Gaussian distribution

$$q(\theta) := \mathcal{N}(\theta|m, S^{-1})$$

Natural parameters

$$\lambda := \{Sm, -S/2\}$$

Expectation parameters

$$\mu := \{\mathbb{E}_q(\theta), \mathbb{E}_q(\theta \theta^\top)\}$$

Conjugate Exp-Fam Models

$$\theta_t - \theta_t^{\setminus i} = H_t^{-1} x_i (y_i - x_i^\top \theta_t^{\setminus i}) = -H_t^{-1} \nabla \ell_i(\theta_t^{\setminus i})$$

We will extend this to posterior's sensitivity

$$q_t \propto \prod_{j=0}^t e^{-\ell_j} \quad q_t^{\setminus i} \propto \prod_{j=0, j \neq i}^t e^{-\ell_j} \quad \frac{q_t}{q_t^{\setminus i}} \propto e^{-\ell_i}$$
$$e^{\lambda_t^\top T(\theta)} \quad e^{(\lambda_t^{\setminus i})^\top T(\theta)} \quad e^{\tilde{\lambda}_i^\top T(\theta)}$$

$$\lambda_t - \lambda_t^{\setminus i} = \tilde{\lambda}_i \quad \text{Lin-reg is a special case [1, Thm. 1]}$$

Linear Regression as a special case

$$q_t = \mathcal{N}(\theta_t, H_t^{-1}) \quad T(\theta) = (\theta, \theta\theta^\top)$$
$$\propto e^{-\frac{1}{2}\theta_t^\top H_t \theta_t + \text{Tr}\left(-\frac{1}{2}H_t \theta\theta^\top\right)} \quad \lambda_t = (H_t \theta_t, -H_t/2)$$

$$q_t^{(i)} = \mathcal{N}(\theta_t^{(i)}, H_t^{(i)-1}) \quad \lambda_t^{(i)} = (H_t^{(i)} \theta_t^{(i)}, -H_t^{(i)}/2)$$
$$e^{-\ell_i} \propto e^{-\frac{1}{2}(y_i - x_i^\top \theta)^2}$$
$$\propto e^{y_i x_i^\top \theta + \text{Tr}\left(-\frac{1}{2}x_i x_i^\top \theta\theta^\top\right)} \quad \tilde{\lambda}_i = (y_i x_i, -x_i x_i^\top / 2)$$

$$\lambda_t - \lambda_t^{(i)} = \tilde{\lambda}_i \quad \implies H_t \theta_t - H_t^{(i)} \theta_t^{(i)} = y_i x_i$$

$$H_t - H_t^{(i)} = x_i x_i^\top$$

$$\theta_t - \theta_t^{(i)} = H_t^{-1} x_i (y_i - x_i^\top \theta_t^{(i)})$$

This addresses all issues!

Group level sensitive (with just addition)

$$\lambda_t - \lambda_t^{\setminus i} = \sum_{i \in \mathcal{M}} \tilde{\lambda}_i$$

Holds at every t during online updating

Can be generalized to neural-network training iterations too (but also to model perturbation and other types of perturbations).

We need “dual” coordinates: $\mu = \mathbb{E}_q[T(\theta)]$

Going beyond conjugacy

$$\theta_t - \theta_t^{\setminus i} = H_t^{-1} x_i (y_i - x_i^\top \theta_t^{\setminus i}) = H_t^{-1} \nabla \ell_i(\theta_t^{\setminus i})$$

$$\lambda_t - \lambda_t^{\setminus i} = \tilde{\lambda}_i = \nabla_{\mu_t} \mathbb{E}_{q_t}[-\ell_i]$$

$$e^{-\ell_i} \propto e^{\tilde{\lambda}_i^\top T(\theta)} \implies -\ell_i = \tilde{\lambda}_i^\top T(\theta) + \text{const.}$$

$$\implies \mathbb{E}_{q_t}[-\ell_i] = \tilde{\lambda}_i^\top \mu_t + \text{const.}$$

$$\implies \nabla_{\mu_t} \mathbb{E}_{q_t}[-\ell_i] = \tilde{\lambda}_i$$

Using this relation we can recover measures used in deep learning (Thm 2-4). Available for free!

Bayesian Learning Rule (BLR) [1]

Many ML algorithms compute the quantity (approx.).
IOW, they are approximately Bayesian!

$$q_t \propto \prod_{j=0}^t e^{-\ell_j} = \arg \min_{q \in \mathcal{Q}} \sum_{j=1}^t \mathbb{E}_q[\ell_j] + KL(q \| p_0)$$

$$\lambda_t = \sum_{j=0}^t \underbrace{\nabla_{\mu_t} \mathbb{E}_{q_t}[-\ell_j]}_{\tilde{\lambda}_{j|t}} \implies \lambda_t = \sum_{j=0}^t \tilde{\lambda}_{j|t}$$

BLR:

$$\lambda_t \leftarrow (1 - \rho)\lambda_t + \rho \sum_{j=0}^t \tilde{\lambda}_{j|t}$$

To estimate sensitivity,
we take a step back

$$\lambda_t^{i} - \lambda_t \approx -\tilde{\lambda}_{i|t}$$

Bayesian learning rule:

Learning Algorithm	Posterior Approx.	Natural-Gradient Approx.	Sec.
Optimization Algorithms			
Gradient Descent	Gaussian (fixed cov.)	Delta method	1.3
Newton's method	Gaussian	—"—	1.3
Multimodal optimization <small>(New)</small>	Mixture of Gaussians	—"—	3.2
Deep-Learning Algorithms			
Stochastic Gradient Descent	Gaussian (fixed cov.)	Delta method, stochastic approx.	4.1
RMSprop/Adam	Gaussian (diagonal cov.)	Delta method, stochastic approx., Hessian approx., square-root scaling, slow-moving scale vectors	4.2
Dropout	Mixture of Gaussians	Delta method, stochastic approx., responsibility approx.	4.3
STE	Bernoulli	Delta method, stochastic approx.	4.5
Online Gauss-Newton <small>(New)</small> (OGN)	Gaussian (diagonal cov.)	Gauss-Newton Hessian approx. in Adam & no square-root scaling	4.4
Variational OGN <small>(New)</small>	—"—	Remove delta method from OGN	4.4
BayesBiNN <small>(New)</small>	Bernoulli	Remove delta method from STE	4.5
Approximate Bayesian Inference Algorithms			
Conjugate Bayes	Exp-family	Set learning rate $\rho_t = 1$	5.1
Laplace's method	Gaussian	Delta method	4.4
Expectation-Maximization	Exp-Family + Gaussian	Delta method for the parameters	5.2
Stochastic VI (SVI)	Exp-family (mean-field)	Stochastic approx., local $\rho_t = 1$	5.3
VMP	—"—	$\rho_t = 1$ for all nodes	5.3
Non-Conjugate VMP	—"—	—"—	5.3
Non-Conjugate VI <small>(New)</small>	Mixture of Exp-family	None	5.4

Improved Variational Online Newton

$$\lambda_t \leftarrow (1 - \rho)\lambda_t + \rho \sum_{j=0}^t \tilde{\lambda}_{j|t}$$

RMSprop/Adam

BLR [1] variant called IVON [5]
(Improved Variational Online Newton)

```
1  $\hat{g} \leftarrow \hat{\nabla} \ell(\theta)$ 
2  $\hat{h} \leftarrow \hat{g}^2$ 
3  $h \leftarrow (1 - \rho)h + \rho \hat{h}$ 
4  $\theta \leftarrow \theta - \alpha(\hat{g} + \delta m) / (\sqrt{h} + \delta)$ 
```

```
1  $\hat{g} \leftarrow \hat{\nabla} \ell(\theta)$  where  $\theta \sim \mathcal{N}(m, \sigma^2)$ 
2  $\hat{h} \leftarrow \hat{g} \cdot (\theta - m) / \sigma^2$ 
3  $h \leftarrow (1 - \rho)h + \rho \hat{h} + \rho^2 (h - \hat{h})^2 / (2(h + \delta))$ 
4  $m \leftarrow m - \alpha(\hat{g} + \delta m) / (h + \delta)$ 
5  $\sigma^2 \leftarrow 1 / (N(h + \delta))$ 
```

Only tune initial value of h (a scalar)

Check out the blog: <https://team-approx-bayes.github.io/blog/ivon/>

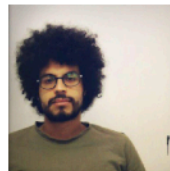
1. Khan, et al. "Fast and scalable Bayesian deep learning by weight-perturbation in Adam." *ICML* (2018).
2. Osawa et al. "Practical Deep Learning with Bayesian Principles." *NeurIPS* (2019).
3. Lin et al. "Handling the positive-definite constraints in the BLR." *ICML* (2020).
4. Shen et al. "Variational Learning is Effective for Large Deep Networks." *ICML* (2024)

IVON got 1st prize in NeurIPS 2021 Approximate Inference Challenge

Watch **Thomas Moellenhoff's** talk at
<https://www.youtube.com/watch?v=LQInIN5EU7E>.

Mixture-of-Gaussian Posteriors with an Improved Bayesian Learning Rule

Thomas Möllenhoff¹, Yuesong Shen², Gian Maria Marconi¹
Peter Nickl¹, Mohammad Emtiyaz Khan¹



1 Approximate Bayesian Inference Team
RIKEN Center for AI Project, Tokyo, Japan

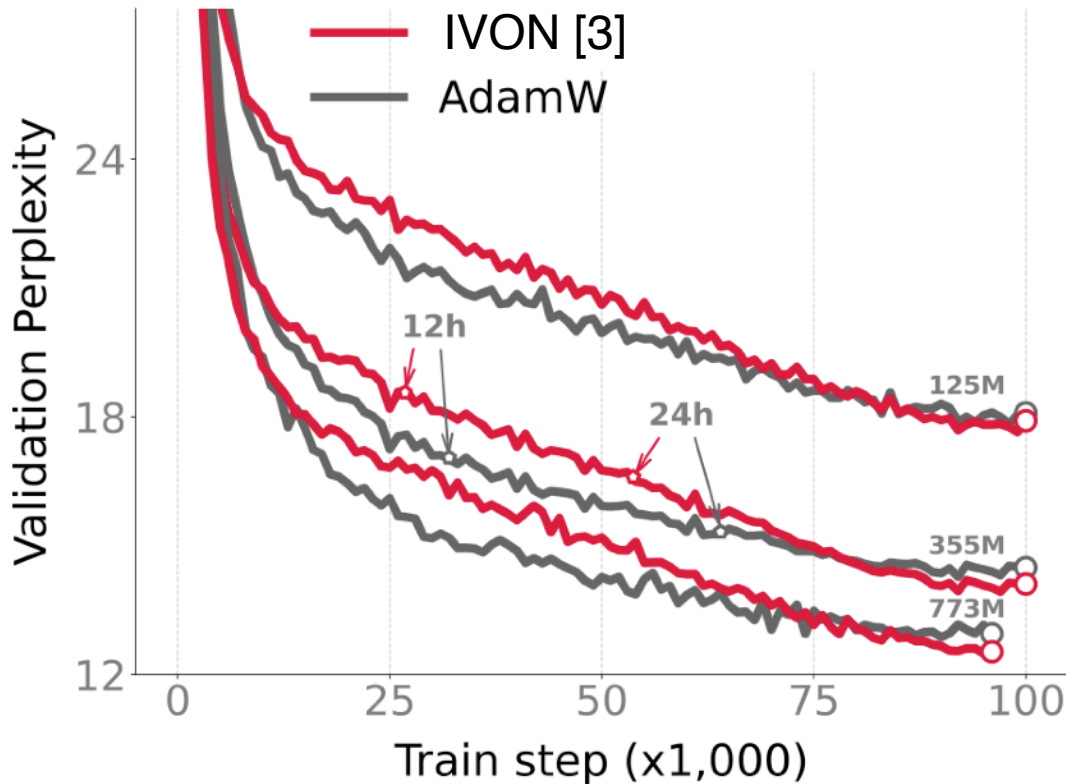
2 Computer Vision Group
Technical University of Munich, Germany

Dec 14th, 2021 — NeurIPS Workshop on Bayesian Deep Learning

1. Khan, et al. "Fast and scalable Bayesian deep learning by weight-perturbation in Adam." *ICML* (2018).
2. Osawa et al. "Practical Deep Learning with Bayesian Principles." *NeurIPS* (2019).
3. Lin et al. "Handling the positive-definite constraints in the BLR." *ICML* (2020).

GPT-2 with IVON

Better performance & uncertainty at the same cost



Trained on OpenWebText data (49.2B tokens).

On 773M, we get a gain of 0.5 in perplexity.

On 355M, we get a gain of 0.4 in perplexity.

1. Khan, et al. "Fast and scalable Bayesian deep learning by weight-perturbation in Adam." *ICML* (2018).
2. Osawa et al. "Practical Deep Learning with Bayesian Principles." *NeurIPS* (2019).
3. Shen et al. Variational Learning is Effective for Large Deep Networks, *ICML* (2024)

Drop-in replacement of Adam

<https://github.com/team-approx-bayes/ivon>

```
import torch
+import ivon

train_loader = torch.utils.data.DataLoader(train_dataset)
test_loader = torch.utils.data.DataLoader(test_dataset)
model = MLP()

-optimizer = torch.optim.Adam(model.parameters())
+optimizer = ivon.IVON(model.parameters())

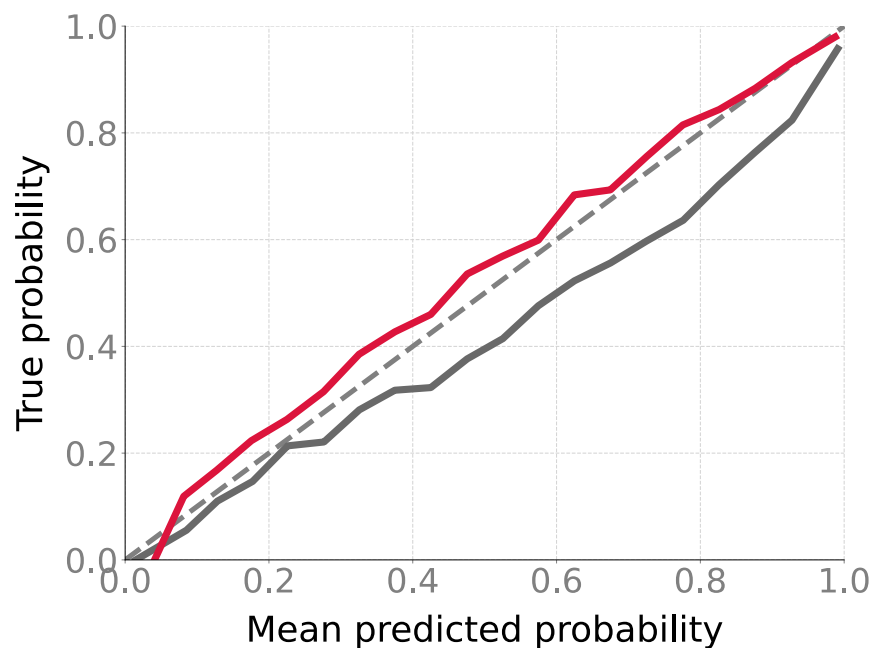
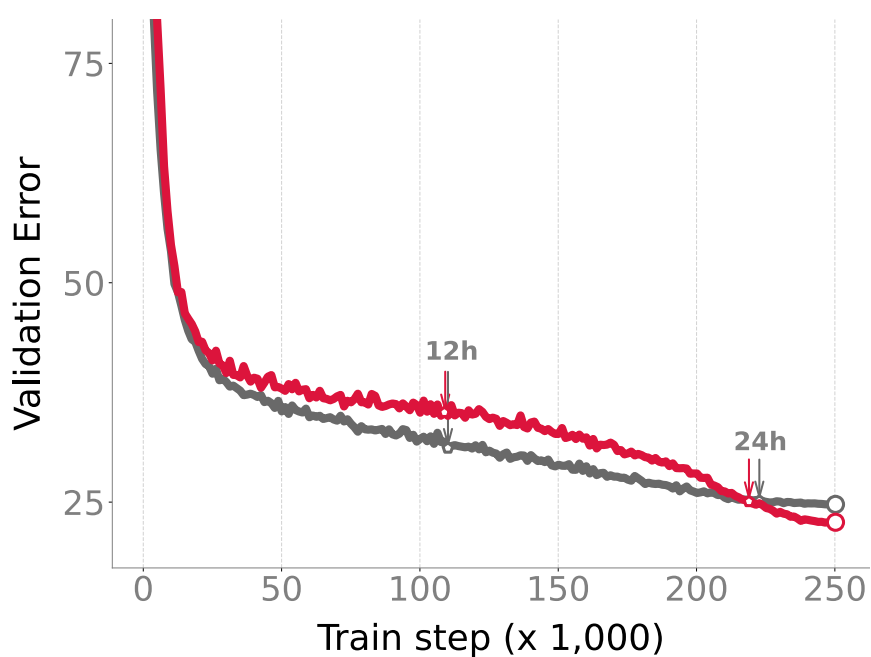
for X, y in train_loader:
+   for _ in range(train_samples):
+       with optimizer.sampled_params(train=True):
           optimizer.zero_grad()
           logit = model(X)
           loss = torch.nn.CrossEntropyLoss(logit, y)
           loss.backward()

optimizer.step()
```

Don't use BBB
Use IVON!

Better Calibration

2% better accuracy over AdamW and 1% over SGD. Better calibration (ECE of 0.022 vs 0.066)



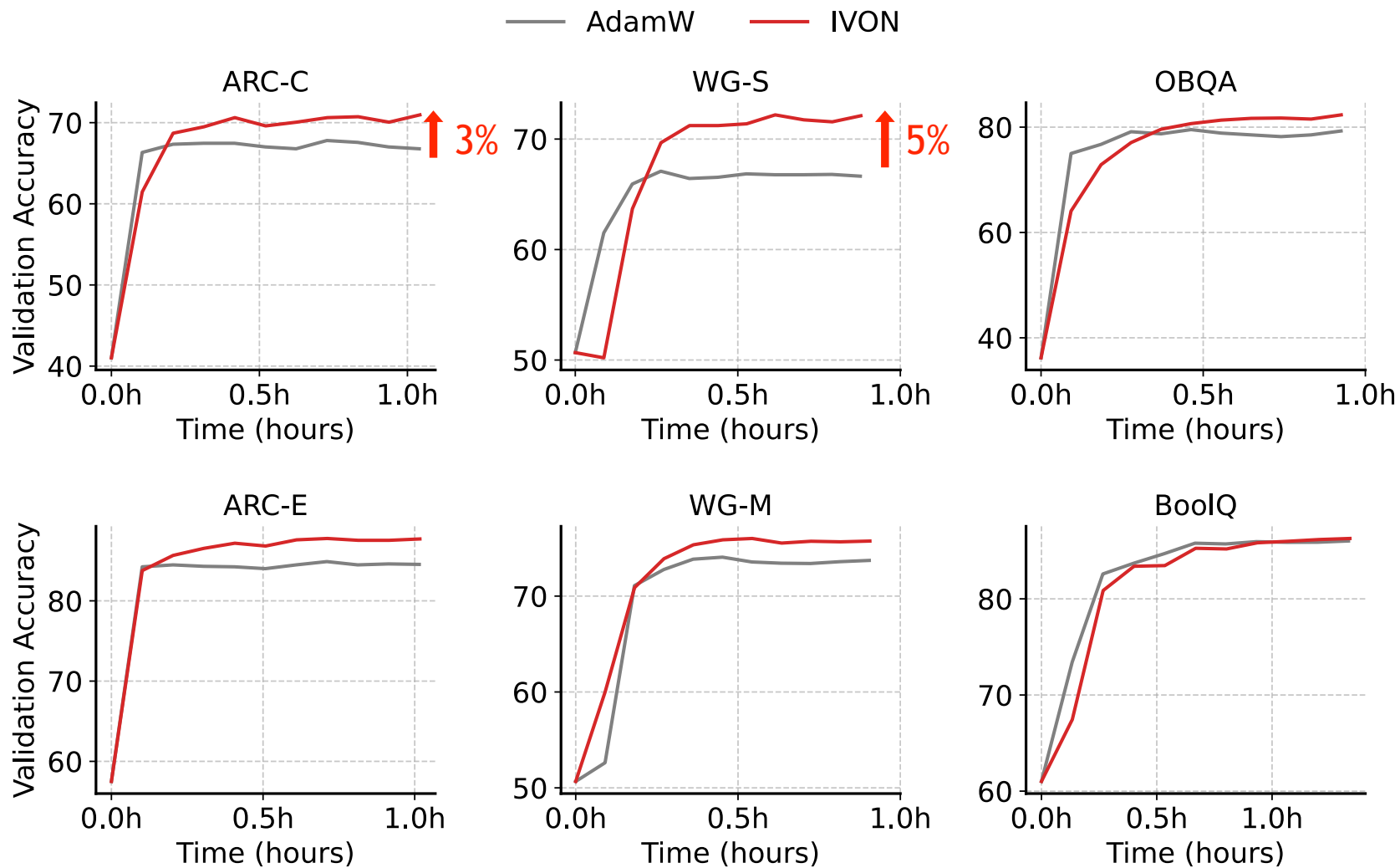
No Severe Overfitting

....like AdamW while improving accuracy over SGD consistently & better uncertainty

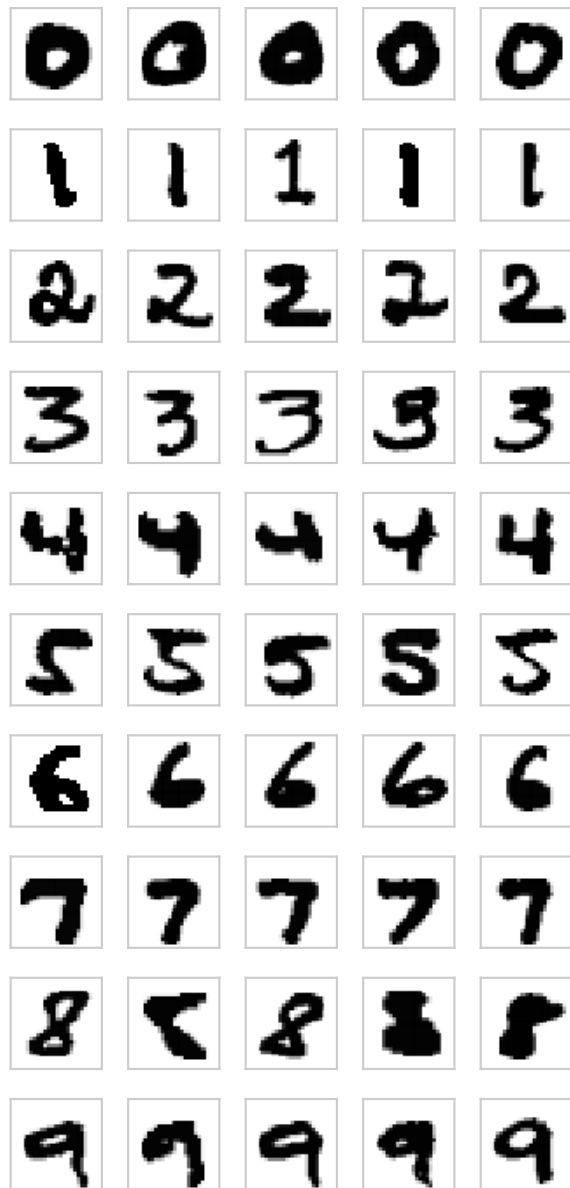
Dataset & Model	Epochs	Method	Top-1 Acc. \uparrow	Top-5 Acc. \uparrow	NLL \downarrow	ECE \downarrow	Brier \downarrow
ImageNet-1k ResNet-50 (25.6M params)	100	AdamW	74.56 \pm 0.24	92.05 \pm 0.17	1.018 \pm 0.012	0.043 \pm 0.001	0.352 \pm 0.003
		SGD	76.18 \pm 0.09	92.94 \pm 0.05	0.928 \pm 0.003	0.019 \pm 0.001	0.330 \pm 0.001
		IVON@mean	76.14 \pm 0.11	92.83 \pm 0.04	0.934 \pm 0.002	0.025 \pm 0.001	0.330 \pm 0.001
		IVON	76.24 \pm 0.09	92.90 \pm 0.04	0.925 \pm 0.002	0.015 \pm 0.001	0.330 \pm 0.001
	200	AdamW	+2% 75.16 \pm 0.14	92.37 \pm 0.03	1.018 \pm 0.003	0.066 \pm 0.002	0.349 \pm 0.002
		SGD	+1% 76.63 \pm 0.45	93.21 \pm 0.25	0.917 \pm 0.026	0.038 \pm 0.009	0.326 \pm 0.006
		IVON@mean	77.30 \pm 0.08	93.58 \pm 0.05	0.884 \pm 0.002	0.035 \pm 0.002	0.316 \pm 0.001
		IVON	77.46 \pm 0.07	93.68 \pm 0.04	0.869 \pm 0.002	0.022 \pm 0.002	0.315 \pm 0.001
TinyImageNet ResNet-18 (11M params, wide)	200	AdamW	+15% 47.33 \pm 0.90	71.54 \pm 0.95	6.823 \pm 0.235	0.421 \pm 0.008	0.913 \pm 0.018
		SGD	+1% 61.39 \pm 0.18	82.30 \pm 0.22	1.811 \pm 0.010	0.138 \pm 0.002	0.536 \pm 0.002
		IVON@mean	62.41 \pm 0.15	83.77 \pm 0.18	1.776 \pm 0.018	0.150 \pm 0.005	0.532 \pm 0.002
		IVON	62.68 \pm 0.16	84.12 \pm 0.24	1.528 \pm 0.010	0.019 \pm 0.004	0.491 \pm 0.001
TinyImageNet PreResNet-110 (4M params, deep)	200	AdamW	+10% 50.65 \pm 0.0*	74.94 \pm 0.0*	4.487 \pm 0.0*	0.357 \pm 0.0*	0.812 \pm 0.0*
		AdaHessian	55.03 \pm 0.53	78.49 \pm 0.34	2.971 \pm 0.064	0.272 \pm 0.005	0.690 \pm 0.008
		SGD	+2% 59.39 \pm 0.50	81.34 \pm 0.30	2.040 \pm 0.040	0.176 \pm 0.006	0.577 \pm 0.007
		IVON @mean	60.85 \pm 0.39	83.89 \pm 0.14	1.584 \pm 0.009	0.053 \pm 0.002	0.514 \pm 0.003
		IVON	61.25 \pm 0.48	84.13 \pm 0.17	1.550 \pm 0.009	0.049 \pm 0.002	0.511 \pm 0.003
CIFAR-100 ResNet-18 (11M params, wide)	200	AdamW	+11% 64.12 \pm 0.43	86.85 \pm 0.51	3.357 \pm 0.071	0.278 \pm 0.005	0.615 \pm 0.008
		SGD	+7% 74.46 \pm 0.17	92.66 \pm 0.06	1.083 \pm 0.007	0.113 \pm 0.001	0.376 \pm 0.001
		IVON@mean	74.51 \pm 0.24	92.74 \pm 0.19	1.284 \pm 0.013	0.152 \pm 0.003	0.399 \pm 0.002
		IVON	75.14 \pm 0.34	93.30 \pm 0.19	0.912 \pm 0.009	0.021 \pm 0.003	0.344 \pm 0.003

LoRA Finetuning

Llama 2 (7 billion)



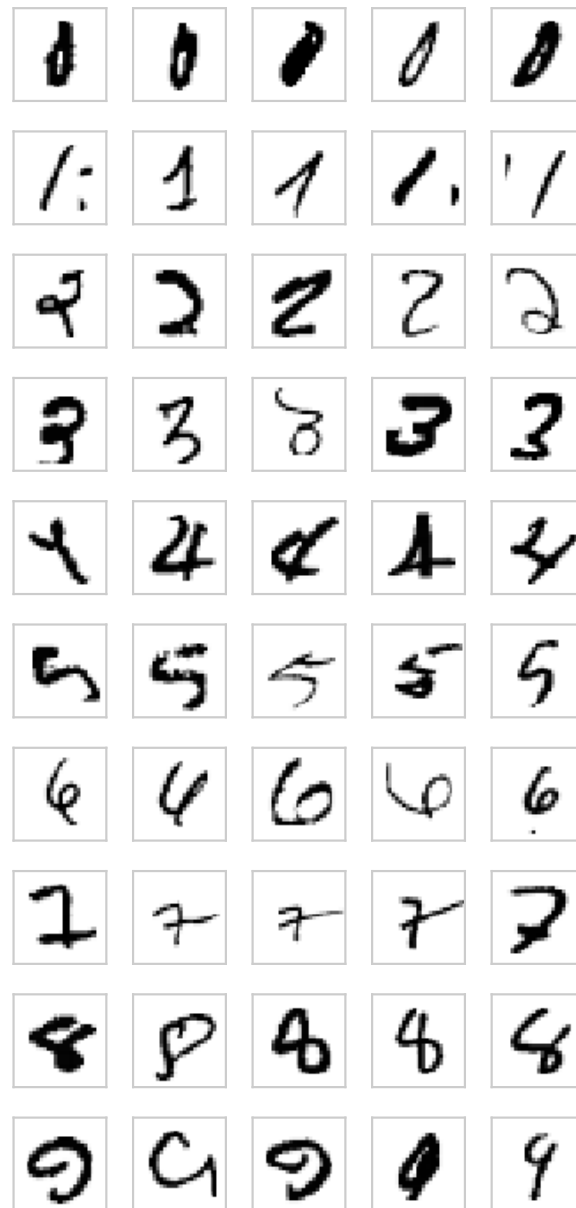
Low Sensitivity



To estimate sensitivity,
just take a step back

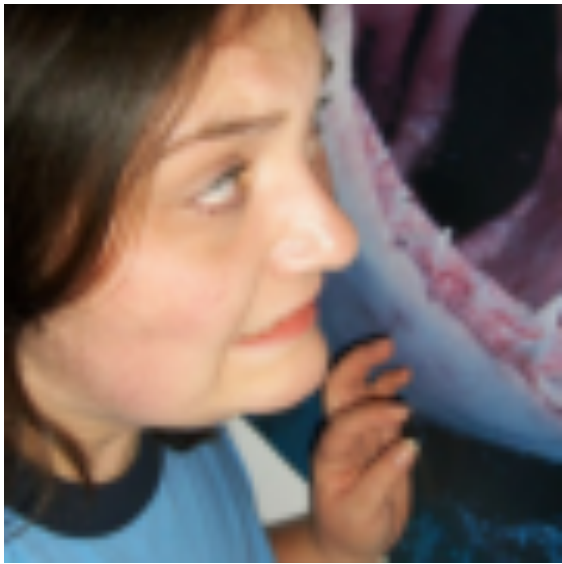
$$\lambda_t^{i|t} - \lambda_t \approx -\tilde{\lambda}_{i|t}$$

High Sensitivity



1. Khan et al. Approximate Inference Turns Deep Networks into Gaussian Process, NeurIPS, 2019
2. Pan et al. Continual Deep Learning by Functional Regularisation of Memorable Past, NeurIPS, 2020

Guess the ImageNet class [1]



High Sensitivity



Traffic light (ImageNet)

What class is this?



Low Sensitivity

High Sensitivity



Chihuahua class (ImageNet)

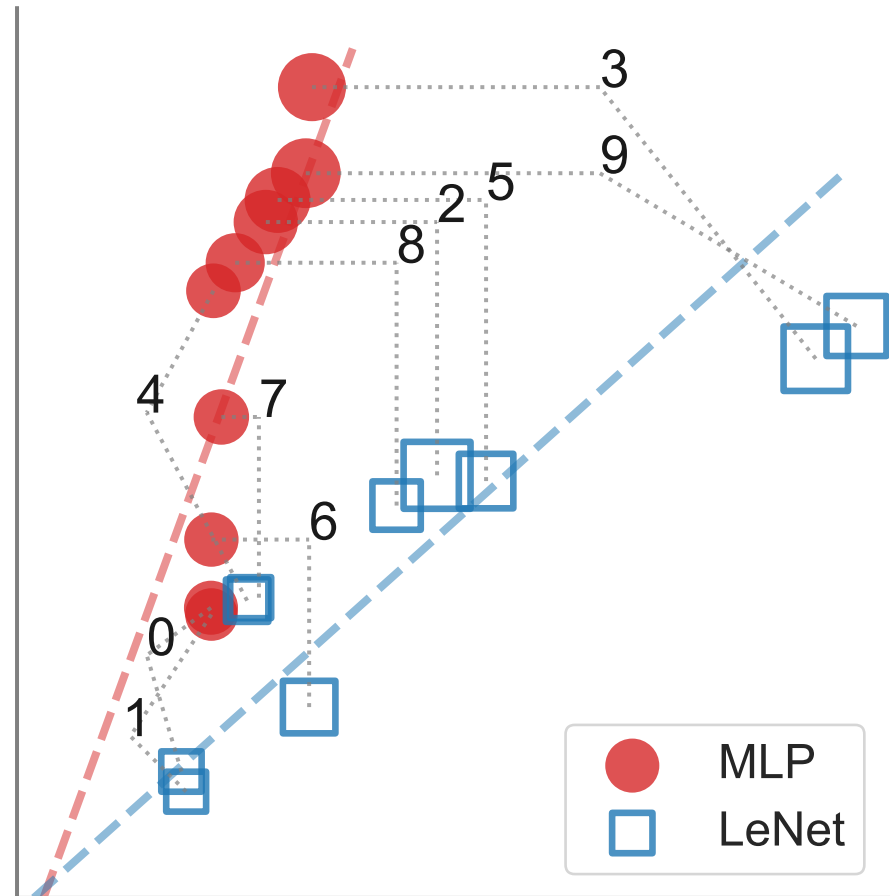


Low Sensitivity

Answering “What-If” Questions

What if we removed a class from MNIST?

Estimates on training data (no retraining)



Test Performance (NLL) by
brute-force retraining

Model Merging

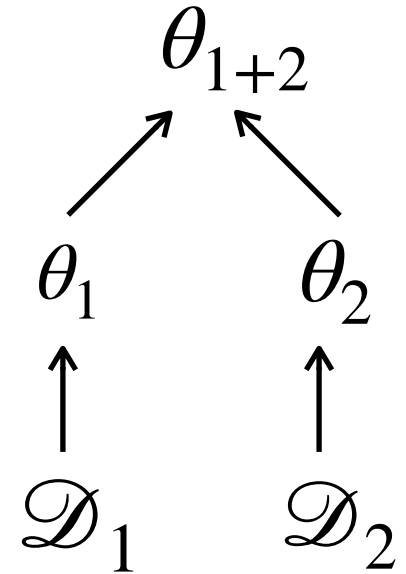
Given θ_1 fine-tuned on \mathcal{D}_1 and θ_2 fine-tuned on \mathcal{D}_2 , merge them (to estimate θ_{1+2}).

Simplest strategy: $\alpha_1\theta_1 + \alpha_2\theta_2$ [1].

A generalization is to use $\alpha_1\lambda_1 + \alpha_2\lambda_2$ [3], eg, use Hessian which is necessarily better [2]

$$H_{1+2}\theta_{1+2} \approx \alpha_1 H_1 \theta_1 + \alpha_2 H_2 \theta_2$$

$$\implies \theta_{1+2} - \theta_1 \approx H_{1+2}^{-1} \nabla \ell_1(\theta_1) \text{ (Thm 1, [2])}$$

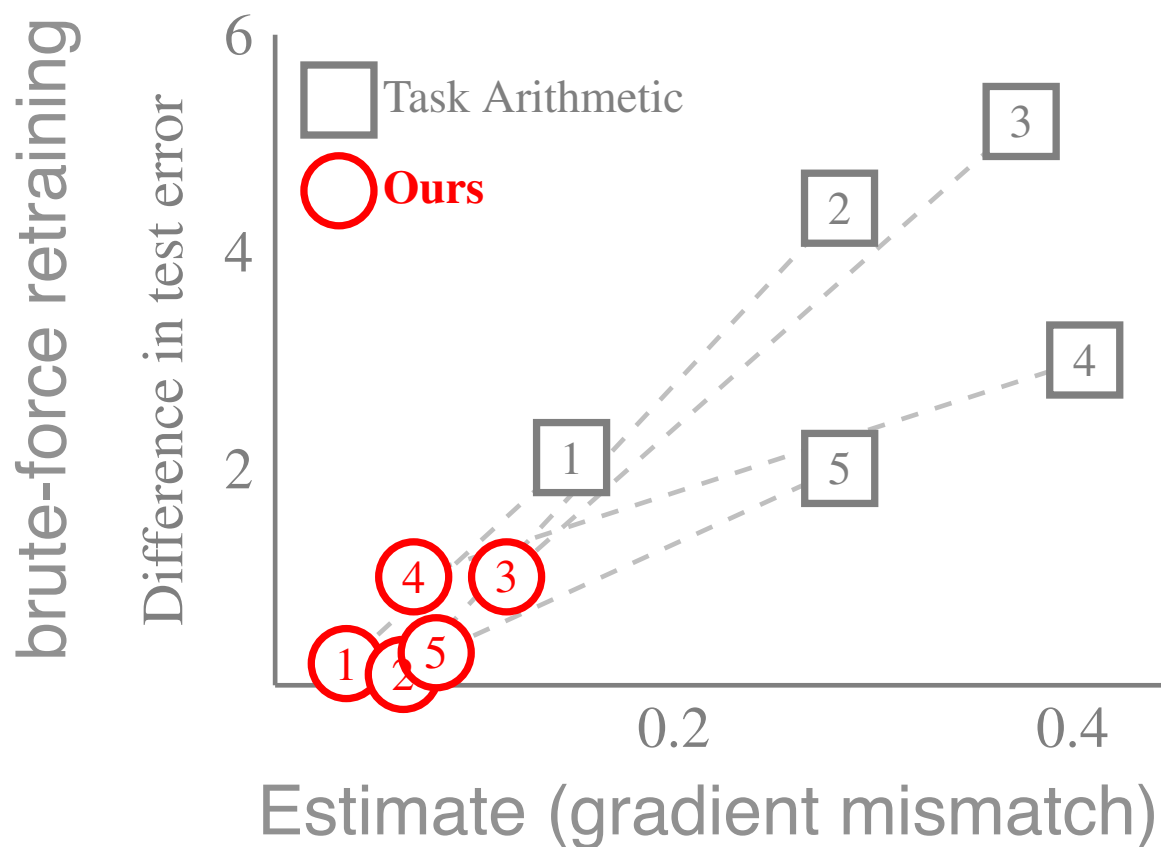


1. Wortsman et al. Robust fine-tuning of zero-shot models, CVPR 2022

2. Daheim et al. Model merging by uncertainty-based gradient matching, ICLR (2024).

3. Maldonado et al. Fast Previews via Bayesian Model-Merging (under review, 2024)

“What-if” we merged models



RoBERTa
on IMDB

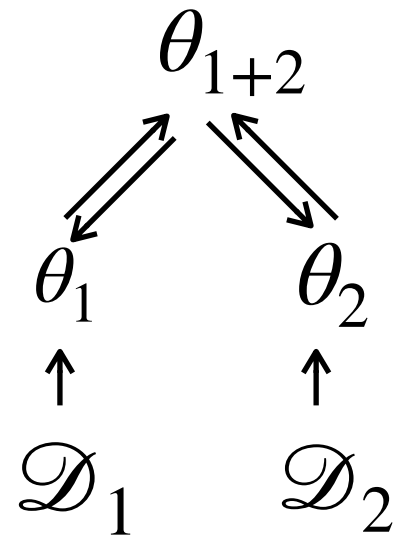
Bayesian Duality

The variables $\tilde{\lambda}_i$ are dual variables (Lagrange multipliers). In fact, variational posteriors have an equivalent dual representation in terms of $\tilde{\lambda}_i$ [1-4]

Federated Learning

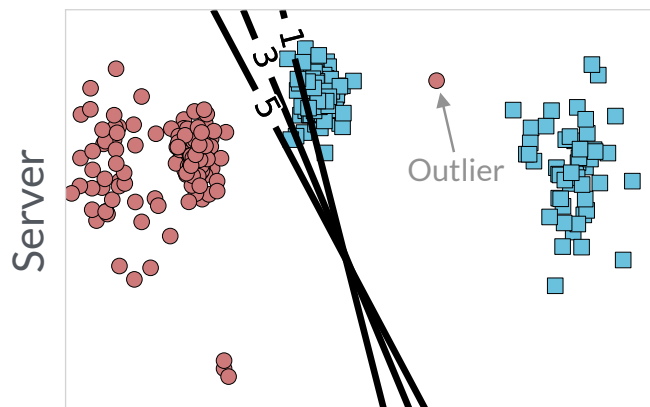
Eg, dual variables in federated ADMM automatically emerges through $\tilde{\lambda}_i$ in variational Bayes [4]

$$\lambda_{1+2} \leftarrow \tilde{\lambda}_1 + \tilde{\lambda}_2$$

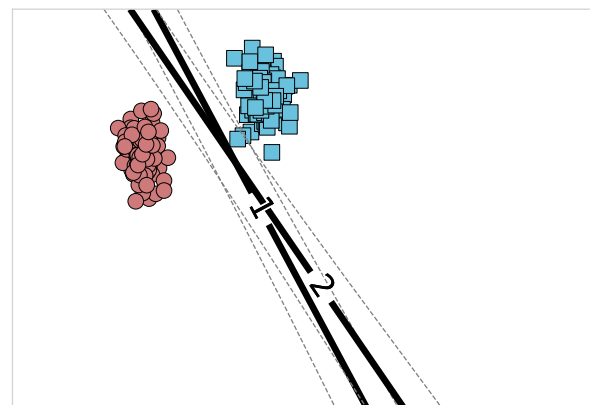
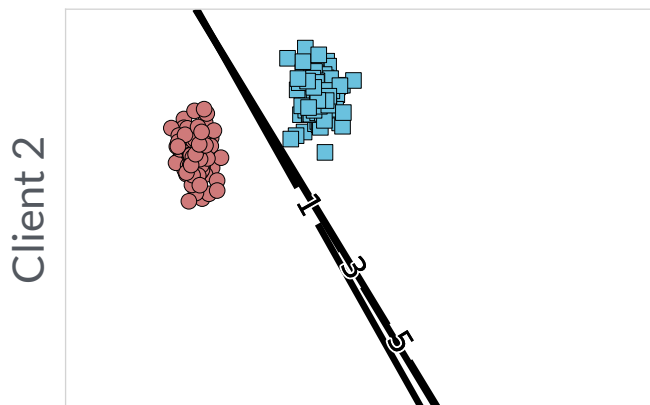
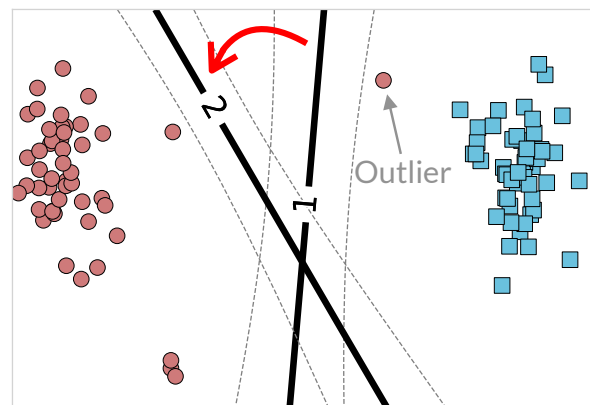
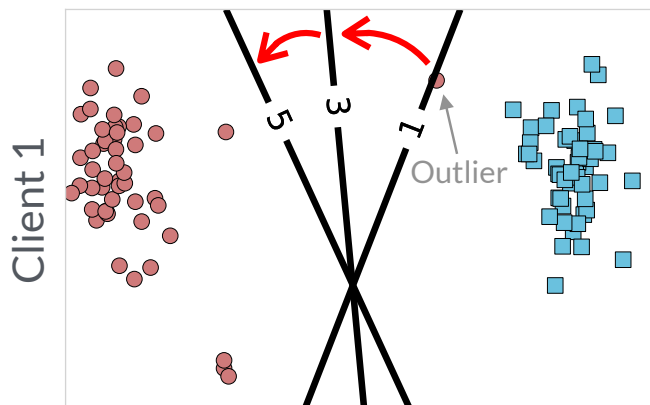
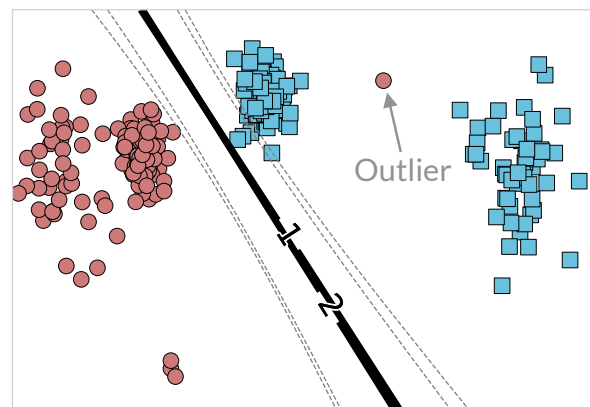


1. Khan et al. Fast Dual Variational Inference for Non-Conjugate Latent Gaussian Models, ICML, 2013
2. Khan et al. Approximate Inference Turns Deep Networks into Gaussian Processes, NeurIPS, 2019
3. Adam et al. Dual Parameterization of Sparse Variational Gaussian Processes, NeurIPS, 2021
4. Swaroop et al. Connecting Federated ADMM to Bayes, ICLR, 2024

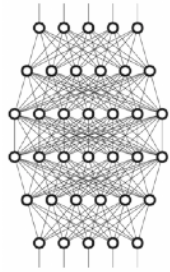
FedADMM



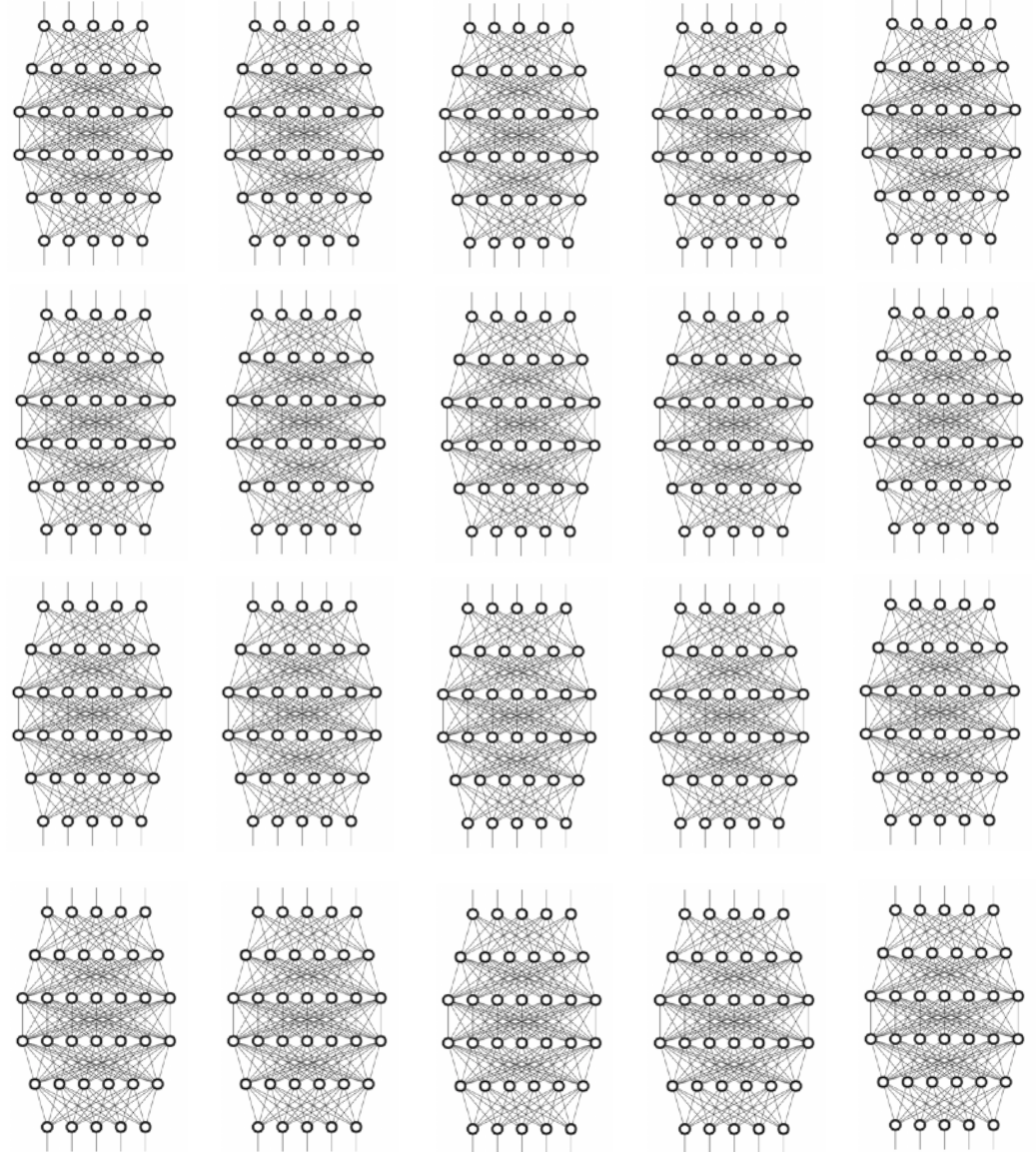
BayesADMM



Standard



Bayes



$$\log \text{Partition} = \sum_{\text{all } s} \text{Leave-S-Out-CV}$$

Sensitivity and Uncertainty

- Sensitivity of (variational) posteriors to address uncertainty during knowledge transfer
 - Main point: it is essentially available for free!
- Model sensitivity to training-data perturbation
 - Beyond linear regression: conjugate-Bayes [1]
 - Beyond conjugacy [2]
 - For large models (GPT-2, ImageNet) [3]
- Model perturbation: LLM model merging [4-5]
 - Federated learning [6] and connections to duality

1. Nickl, Xu, Tailor, Moellenhoff, Khan, The memory-perturbation equation, NeurIPS (2023)

2. Khan and Rue, The Bayesian Learning Rule, JMLR (2023).

3. Shen et al. Variational Learning is Effective for Large Deep Networks, ICML (2024)

4. Daheim et al. Model merging by uncertainty-based gradient matching, ICLR (2024).

5. Moldanado et al. How to Weight Multitask Finetuning? Fast Previews via Bayesian Model-Merging, (2024)

6. Swaroop et al. Connecting Federated ADMM to Bayes, ICLR, 2024

The Bayes-Duality Project

Toward AI that learns adaptively, robustly, and continuously, like humans



Emtiyaz Khan

Research director
(Japan side)

Approx-Bayes team at
RIKEN-AIP and OIST



Julyan Arbel

Research director
(France side)

Statify-team, Inria
Grenoble Rhône-Alpes



Kenichi Bannai

Co-PI (Japan side)

Math-Science Team at
RIKEN-AIP and Keio
University



Rio Yokota

Co-PI
(Japan side)

Tokyo Institute of
Technology

Received total funding of JPY 220M + EUR 500K through the CREST-ANR grant! Thanks to JST for their generous funding!

Bayes-Duality Workshop

https://bayesduality.github.io/workshop_2024.html



Adam White

University of Alberta,
Canada



Alexander Immer

ETH, Switzerland



Arindam Banerjee

University of Illinois
Urbana-Champaign,
US



Daiki Chijiwa

NTT Corporation,
Japan



Ehsan Amid

Google DeepMind,
US



Eugene Ndiaye

Apple, France



Frank Nielsen

Sony Computer
Science Laboratories,
Japan



Jonghyun Choi

Seoul National
University, South
Korea



Juho Lee

KAIST, South Korea



Haavard Rue

KAUST, Saudi Arabia



Hossein Mobahi

Google Research, US



Martin Mundt

TU Darmstadt,
Germany



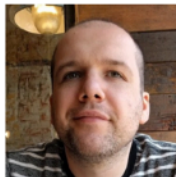
Matt Jones

University of
Colorado, US



Nico Daheim

TU Darmstadt,
Germany



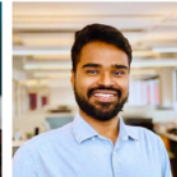
Razvan Pascanu

Google DeepMind,
US



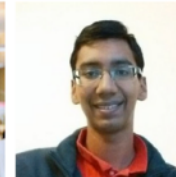
Rupam Mahmood

University of Alberta,
Canada



Sarath Chandar

École Polytechnique
de Montréal, Canada



Siddharth Swaroop

Harvard University,
US



Tom Rainforth

University of Oxford,
UK



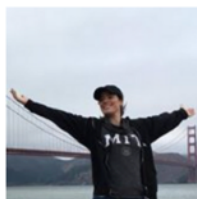
Vincent Fortuin

Helmholtz AI,
Germany



Yingzhen Li

Imperial College
London, UK



Zelda Mariet

Bioptimus, US

Every June in Tokyo (June 25-27, 2025)
Attendees are from a diverse research
interests: Bayes, Duality, Continual/
Federated/Active learning,
RL, Experiment Design etc.

Team Approx-Bayes

<https://team-approx-bayes.github.io/>



Emtiyaz Khan
Team Leader



Thomas Möllenhoff
Research Scientist



Keigo Nishida
Special Postdoctoral
Researcher
RIKEN BDR



**Hugo Monzón
Maldonado**
Postdoctoral
Researcher



**Christopher Johannes
Anders**
Postdoctoral
Researcher



Yohan Jung
Postdoctoral
Researcher



Sin-Han Yang
Technical Staff



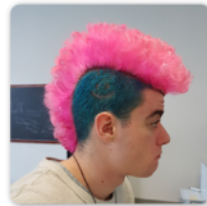
Anita Yang
Part-Time Student
The University of
Tokyo



Bai Cong
Part-Time Student
Tokyo Institute of
Technology



Eiki Shimizu
Part-Time Student
Institute of Statistical
Mathematics



Marco Miani
Intern
Technical University of
Denmark



Rin Intachuen
Intern
Mahidol University



Alexander Timans
Intern
University of
Amsterdam



Masaki Adachi
Intern
University of Oxford



Adrian R. Minut
Intern
Sapienza, University of
Rome



Joseph Austerweil
Visiting Scientist
University of
Wisconsin-Madison



Pierre Alquier
Visiting Scientist
ESSEC Business
School



Geoffrey Wolfer
Visiting Scientist
Waseda University



Rio Yokota
Visiting Scientist
Tokyo Institute of
Technology



Dharmesh Tailor
Remote Collaborator
University of
Amsterdam

Visit us! Let's collaborate!
Also see open (post-doc)
positions on the webpage