



Adaptive Bayesian Intelligence (AGI meets ABI)

Mohammad Emtiyaz Khan RIKEN Center for AI Project, Tokyo https://emtiyaz.github.io



Summary of recent research at <u>https://emtiyaz.github.io/papers/symposium_2024.pdf</u> Slides available at <u>https://emtiyaz.github.io/</u>

Al that can learn like us

Quickly adapt & continue to acquire new skills

Human Learning at the age of 6 months.



Converged at the age of 12 months



Transfer skills at the age of 14 months



Teacher-Student Learning?



Current state of Machine Learning



Retraining from Scratch

Even when changes are tiny. It is costly, undemocratic and unsustainable.

Adaptive Intelligence

How do brains adapt quickly? What do they optimize and how?

1. Sternberg. A theory of adaptive intelligence and its relation to general intelligence. *Journal of Intelligence (2019)*

2. Sternberg. Adaptive intelligence. New York: Cambridge University Press (2021)

3. Sternberg. What is intelligence really? the futile search for a holy grail. Learning & Individual Differences (2024)₉

Adaptive Bayesian Intelligence

- Adaptive Intelligence = Bayesian Computation
- Part 1: Bayesian Learning Rule [1]
 - Foundational way to derive learning-algorithms
 - Application to Deep Learning [2]
- Part 2: Posterior Correction [3]
 - Foundational way to derive adaptation-algorithms
 - Application to continual learning [4-5]
 - But also for LLM merging, Federated Learning etc.
- Adaptive Bayesian Intelligence: A roadmap.
- 1. Khan and Rue, The Bayesian Learning Rule, JMLR (2023)
- 2. Shen et al. Variational Learning is Effective for Large Deep Networks, ICML (2024)
- 3. Khan. Knowledge Adaptation as Posterior Correction, arXiv (2025)
- 4. Khan and Swaroop. Knowledge-Adaptation Priors, NeurIPS (2021).
- 5. Pan et al. Continual Deep Learning by Functional Regularisation of Memorable Past, NeurIPS, 2020

"The fact that many different approaches point to the same actual algorithm is a major strength of Bayesianity"

-E. T. Jaynes, discussion of [1]





1. Zellner, Optimal Information Processing and Bayes' Theorem. The American Statistician (1988)

Optimization

Gradient Descent Newton's Method Multimodal Optimization

Deep-Learning

SGD, RMSprop and Adam Sharpness-Aware Minimization Dropout, STE, Label Smoothing Shampoo....

Bayesian Learning Rule [1]

Approximate Inference

Conjugate Bayes Laplace's Method Expectation Maximization Stochastic Variational Inference Variational Message Passing

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

Global-Optimization

Exponential-Weight Aggregation Natural Evolution Strategy Gaussian Homotopy Smoothed Optimization Weight-perturbed Optimization Stochastic Search (annealing) Stochastic Relaxation

Variational Formulation of Bayes' Rule

Bayes' Rule:
$$p_t(\theta) \propto p_0(\theta) \prod_{j=1}^t \text{lik}_j(\theta)$$

Variational Inference to find an approximation $q_t(\theta)$

$$q_{t} = \arg\min_{q \in \mathcal{Q}} \sum_{j=1}^{t} \mathbb{E}_{q} [-\log \operatorname{lik}_{j}] + KL(q \| p_{0})_{\propto e^{-\ell_{0}}}$$
$$= \ell_{j}$$
$$= \arg\min_{q \in \mathcal{Q}} \sum_{j=0}^{t} \mathbb{E}_{q} [\ell_{j}] - \mathcal{H}(q)$$

We will use this variational formulation to discover the inherent Bayesian nature of (non-Bayesian) algorithms.

Exponential Family

Natural
parametersSufficient
StatisticsExpectation
parameters
$$q(\theta) \propto \exp\left[\lambda^{\top}T(\theta)\right]$$
 \downarrow \downarrow $\mathcal{N}(\theta|m, S^{-1}) \propto \exp\left[-\frac{1}{2}(\theta - m)^{\top}S(\theta - m)\right]$
 $\propto \exp\left[(Sm)^{\top}\theta + \operatorname{Tr}\left(-\frac{S}{2}\theta\theta^{\top}\right)\right]$ Gaussian distribution
Natural parameters $q(\theta) := \mathcal{N}(\theta|m, S^{-1})$
 $\lambda := \{Sm, -S/2\}$
Expectation parameters

Wainwright and Jordan, Graphical Models, Exp Fams, and Variational Inference Graphical models 2008
 Malago et al., Towards the Geometry of Estimation of Distribution Algos based on Exp-Fam, FOGA, 2011 14

Bayesian Learning Rule (BLR) [1]



Algorithms (such as SGD/Adam) are special cases of BLR obtained by choosing specific exp-family q_{λ} with natural parameter λ and expectation parameter μ .

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

Deriving Gradient Descent from BLR

Derived by choosing Gaussian with fixed covariance

Gaussian distribution $q(\theta) := \mathcal{N}(m, 1)$ Natural parameters $\lambda := m$ Expectation parameters $\mu := \mathbb{E}_q[\theta] = m$ $\mathcal{H}(q) := \log(2\pi)/2$ Entropy BLR: $\lambda \leftarrow \lambda - \rho \nabla_{\mu} \Big(\mathbb{E}_q[\bar{\ell}] - \mathscr{H}(q) \Big)$ $m \leftarrow m - \rho \ \nabla_m \mathbb{E}_a[\mathscr{C}]$ $m \leftarrow m - \rho \, \mathbb{E}_q[\nabla_\theta \mathscr{C}]$ Bonnet's theorem $m \leftarrow m - \rho \nabla \overline{\ell}(m)$ First-order delta method $\theta \leftarrow \theta - \rho \,\nabla \, \ell(\theta)$

Bayesian learning rule:

Learning Algorithm	Posterior Approx.	Natural-Gradient Approx.				
Optimization Algorithms						
Gradient Descent	Gaussian (fixed cov.)	Delta method				
Newton's method	Gaussian	"				
$Multimodal\ optimization\ {}_{\rm (New)}$	Mixture of Gaussians	"	3.2			
Gradient DescentGaussian (fixed cov.)Delta method1.3Newton's methodGaussian—"—1.3Multimodal optimization (New)Mixture of Gaussians—"—3.2Deep-Learning AlgorithmsStochastic Gradient DescentGaussian (fixed cov.)Delta method, stochastic approx.4.1RMSprop/AdamGaussian (diagonal cov.)Delta method, stochastic approx., square-root scaling, slow-moving scale vectors4.2DropoutMixture of GaussiansDelta method, stochastic approx., square-root scaling, slow-moving scale vectors4.3STEBernoulliDelta method, stochastic approx.4.5Online Gauss-Newton (OGN)Gaussian (diagonal cov.)Gauss-Newton Hessian approx. in Adam & no square-root scaling4.4Nariational OGN (New)—"—Remove delta method from OGN4.4BayesBiNN (New)BernoulliRemove delta method from STE4.5Approx: tract Bayesian Inferetex 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 scal- ing, slow-moving scale vectors				
Dropout	Mixture of Gaussians	Delta method, stochastic approx., responsibility approx.				
STE	Bernoulli	Delta method, stochastic approx.				
Online Gauss-Newton (OGN) $_{(New)}$	Gaussian (diagonal cov.)	Gauss-Newton Hessian approx. in Adam & no square-root scaling				
Variational OGN $_{(New)}$	"	Remove delta method from OGN 4				
$BayesBiNN \ ({\rm New})$	Bernoulli	Remove delta method from STE				
Approximate Bayesian Inference Algorithms						
Conjugate Bayes	Exp-family	Set learning rate $\rho_t = 1$	5.1			
Laplace's method	Gaussian	Delta method				
Expectation-Maximization	Exp-Family + Gaussian	Delta method for the parameters				
Stochastic VI (SVI)	Exp-family (mean-field)	Stochastic approx., local $\rho_t = 1$				
VMP	"	$ \rho_t = 1 $ for all nodes	5.3			
Non-Conjugate VMP	"	"	5.3			
Non-Conjugate VI (New)	Mixture of Exp-family	None	5.4			

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

Improving Adam using BLR

RMSprop/Adam

BLR with diagonal-cov Gaussian [4] (Improved Variational Online Newton)

Differences: sampling in line 1, hessian used in line 2 (not g^2), h > 0 constraint in line 3, no square-root over h in line 4.

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)

Training GPT-2 from Scratch using BLR

Better performance & uncertainty at the same cost [3]



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)

Better Calibration

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



LoRA Finetuning [1] Llama 2 (7 billion)



1. Bai et al. "Variational Low-Rank Adaptation Using IVON", FITML workshop at NeurIPS 2024

Taylor vs Bayes

Why do we recover optimization algorithm from BLR?



Eq. 18 in Khan and Nielsen (2018), Eq. 59 Khan and Rue (2023), Eq. 3 in Section 2 in Khan (2025)

Bayes Generalizes Taylor

BLR with full Cov Gaussian:

2nd-ord

$$\sum_{i} \theta^{\mathsf{T}} \mathbb{E}_{q_{old}} [\nabla \mathscr{E}_{i}]$$

$$+ \frac{1}{2} (\theta - m_{old})^{\mathsf{T}} \mathbb{E}_{q_{old}} [\nabla^{2} \mathscr{E}_{i}] (\theta - m_{old})$$
BLR with exponential-family:
$$\sup_{i \in \mathcal{I}} \sup_{i \in \mathcal{I}} (T(\theta)^{\mathsf{T}} \lambda_{old})$$

$$= \exp\left(-\sum_{i=0}^{t} \frac{T(\theta)^{\mathsf{T}} \nabla_{\mu} \mathbb{E}_{q_{old}}[\mathscr{E}_{i}]}{\operatorname{Site} \widehat{\mathscr{E}}_{i|old}(\theta)}\right)$$

Sites are important for adaptation! $_{\scriptscriptstyle 23}$

 m_{old}

Dual-Representation of the BLR

$$q_{t} \propto \exp(T(\theta)^{\mathsf{T}}\lambda_{t}) = \exp\left(-\sum_{i=0}^{t} T(\theta)^{\mathsf{T}}\nabla_{\mu}\mathbb{E}_{q_{t}}[\mathscr{C}_{i}]\right)$$

$$q_{t} \propto \prod_{i=0}^{t} \exp(-\hat{\mathscr{C}}_{i|t}) \qquad \lambda_{t} = \sum_{i=0}^{t} \nabla_{\mu}\mathbb{E}_{q_{t}}[\mathscr{C}_{i}]$$

Posterior Sites Natural parameters of the second se

Natural Gradients are additive (representation theorem). Largest ones are the most influential.

Khan et al. Fast Dual Variational Inference for Non-Conjugate Latent Gaussian Models. ICML (2013)
 Khan and Nielsen. Fast yet Simple Natural-Gradient Descent for Variational Inference ... ISITA (2018)
 Khan et al. Approximate Inference Turns Deep Networks into Gaussian Processes. NeurIPS (2019)
 Adam et al. Dual Parameterization of Sparse Variational Gaussian Processes. NearIPS (2021)
 Chang et al. Memory-Based Dual Gaussian Processes for Sequential Learning. ICML (2023)
 Moellenhoff et al. Federated ADMM from Bayesian Duality. arXiv (2025)



The site parameters can be used to generalize Influence Estimators.

Binary classification on the Two Moons dataset.

Big markers with red indicate bigger first derivative for IVON [1]

By Rin Intachuen (RIKEN AIP)

Epoch: 0 1. Nickl, Xu, Tailor, Moellenhoff, Khan, The memory-perturbation equation, NeurIPS, 2023





Traffic light (ImageNet)



High Influence

What class is this?

Low Influence



Chihuahua class (ImageNet)



High Influence

Low Influence

Continual Learning

Elastic Weight Consolidation Variational Continual Learning Memory Replay Methods Functional Regularization

Model Merging

Task Arithmetic Fisher/Hessian-Based Merging Ensembles Methods

Posterior Correction [1]

Unlearning and Influence

Student-Teacher Learning

Knowledge Distillation Learning with Privileged information Incremental SVMs

Federated Learning

FedAvg, FedDyn Alternating Direction Method of Multipliers (ADMM) Alternating Minimization Algorithm (AMA) Partition Variational Inference

1. Khan, Knowledge Adaptation as Posterior Correction, arXiv (2025)

Variational Formulation of Online Bayesian Inference

Bayes' Rule:
$$p_{t+1}(\theta) \propto p_0(\theta) \prod_{j=1}^{t+1} e^{-\ell_j(\theta)} \propto p_t(\theta) e^{-\ell_{t+1}(\theta)}$$

Variational formulation:

Batch:
$$q_{t+1} = \arg \min_{q} \sum_{j=1}^{t+1} \mathbb{E}_{q}[\ell_{j}] + KL(q||p_{0})$$

Online [1]: $\hat{q}_{t+1} = \arg \min_{q} \mathbb{E}_{q}[\ell_{t+1}] + KL(q||q_{t})$

How inaccurate is \hat{q}_{t+1} ? Can we correct it to exactly recover q_{t+1} ? This is the goal of posterior correction.

Continual Learning

Model Merging



Posterior Correction [1]

Unlearning and Influence

Federated Learning





Correct the Past due to the Interference Created by the Future



Eq. 4 in Khan (2025)

 $\boldsymbol{\Omega}$

Posterior Correction

We will use the site functions to correct the posterior!

Batch:
$$q_{t+1} = \arg \min_{q} \sum_{j=1}^{t+1} \mathbb{E}_{q}[\ell_{j}] + KL(q || p_{0}) \xrightarrow{q_{t}} \overline{\prod_{i=0}^{t} \exp(-\hat{\ell}_{j|t})}$$

$$= \arg \min_{q} \mathbb{E}_{q}[\ell_{t+1}] + KL(q || q_{t}) + \sum_{j=0}^{t} \mathbb{E}_{q}[\ell_{j} - \hat{\ell}_{j|t}]$$
Correction
Online: $\hat{q}_{t+1} = \arg \min_{q} \mathbb{E}_{q}[\ell_{t+1}] + KL(q || q_{t})$

Very simple proof (3 lines). Exact recovery in general!

Correction as Prediction Mismatch

Linear regression with isotropic Gaussian posterior

$$m_{t+1} = \arg \min_{m} \mathbb{E}_{q} [\frac{1}{2} (y_{t+1} - x_{t+1}^{\top} \theta)^{2}] + KL \left[\mathcal{N}(m, I) \| \mathcal{N}(m_{t}, I) \right]$$

$$+ \sum_{j=1}^{t} \frac{1}{2} (x_{j}^{\top} m_{t} - x_{j}^{\top} m)^{2} + \dots$$

$$+ \sum_{j=1}^{t} \frac{1}{2} (x_{j}^{\top} m_{t} - x_{j}^{\top} m)^{2} + \dots$$
Error due to mean-field is fixed by the correction!
$$\frac{1}{2} (m - m_{t})^{\top} \left(\sum_{j=1}^{t} x_{j} x_{j}^{\top} \right) (m - m_{t})$$

Prediction mismatch is simpler to implement!

Knowledge-Adaptation Prior

Posterior correction with isotropic Gaussian reduces to "prediction or gradient mismatch" (K-priors) [1]

$$\theta_{t+1} = \arg\min_{\theta} \mathscr{C}_{t+1} + \frac{\rho}{2} \|\theta - \theta_t\|^2 + \sum_{j=1}^{l} \mathscr{C}_j \left(\hat{y}_j(\theta_t), \, \hat{y}_j(\theta) \right)$$

Many adaptation methods reduce this mismatch [2-9] and Posterior Correction generalizes it!

- 1. Khan and Swaroop. Knowledge-Adaptation Priors, NeurIPS (2021).
- 2. Kirkpatrick et al. Overcoming catastrophic forgetting in neural networks. PNAS, 2017.
- 3. Benjamin et al. Measuring and regularizing networks in function space. ICLR 2019.
- 4. Hinton et al. Distilling the knowledge in a neural network, arXiv, 2015.
- 5. Buzzega et al. Dark experience for general continual learning: a strong, simple baseline. NeurIPS 2020.
- 6. Cauwenberghs and Poggio. Incremental and decremental SVM learning. NeurIPS, 2001.
- 7. Vapnik and Izmailov. Learning using privileged information: similarity control and JMLR, 2015.
- 8. Lopez-Paz and Ranzato. Gradient episodic memory for continual learning, NIPS'17
- 9. Csató and Opper. Sparse on-line Gaussian processes. Neural computation, 2002.

From Quick to Slow Adaptation

Correction as Information Gain



Quick Adaptation with Compact Memory

Choose memories where interference is more likely. Small correction \implies Small memory \implies Quick adaptation



1. Pan et al. Continual Deep Learning by Functional Regularisation of Memorable Past, NeurIPS, 2020

37

Combine Methods to Reduce Correction

Get 78% accuracy with 7.5% (random) memory



^{1.} Daxberger et al. Improving CL by Accurate Gradient Reconstruction of the Past, TMLR 2023.

Reducing Correction Improves Performance in LLM fine-tuning



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

Summary of Federated Learning, Model Merging, and Memories etc.

Recover
$$q_{jnt}$$
 from q_1 and q_2
 $q_{jnt} = \arg\min_{q} KL(q || q_1 q_2) + \sum_{j=1}^{2} \mathbb{E}_q[\ell_j - \hat{\ell}_{j|j}]$
 $\mathcal{D}_1 \qquad q_{jnt} \qquad \mathcal{D}_2$

By choosing different q, we get different strategies (better q gives better merging) [1,2]. Same is true for federated learning [3,4]. All of them will benefit from compact memories designed to reduce corrections [5].

- 1. Daheim et al. Model merging by uncertainty-based gradient matching, ICLR (2024).
- 2. Monzon et al. How to Weight Multitask Finetuning? Fast Previews via Bayesian Model-Merging, 2024
- 3. Swaroop, Khan, Doshi, Connecting Federated ADMM to Bayes, ICLR 2025
- 4. Moellenhoff et al. Federated ADMM from Bayes Duality, arXiv, 2025
- 5. Nickl, Xu, Tailor, Moellenhoff, Khan, The memory-perturbation equation, NeurIPS (2023)

ADMM as a special case of Bayes (Dual)



Algorithm 1 BayesADMM (Fig. 2b) for Gaussians with diagonal covariance. Additional steps when compared to FederatedADMM are highlighted in red. Implementation details are in App. D.

Hyperparameters: Prior precision $\delta > 0$, step-sizes $\rho > 0$ and $\gamma > 0$. **Initialize:** $\mathbf{v}_k \leftarrow 0$, $\mathbf{u}_k \leftarrow 0$, $\mathbf{\bar{m}} \leftarrow 0$, $\mathbf{\bar{s}} \leftarrow \delta$, $\alpha \leftarrow 1/(1 + \rho K)$.

- 1: while not converged do
- 2: Broadcast $\bar{\mathbf{m}}$ and $\bar{\mathbf{s}}$ to all clients.
- 3: for each client $1, \ldots, K$ in parallel do
- 4: Local training on $\ell_k(\boldsymbol{\theta}) + \boldsymbol{\theta}^\top \mathbf{v}_k \frac{1}{2} \boldsymbol{\theta}^\top (\mathbf{u}_k \boldsymbol{\theta}) + \frac{\rho}{2} \|\boldsymbol{\theta} \bar{\mathbf{m}}\|_{\bar{\mathbf{s}}}^2 > \mathbf{V}_k$
 - ▷ Using IVON [53]

> An additional dual variable.

- 5: $\mathbf{v}_k \leftarrow \mathbf{v}_k + \gamma \left(\mathbf{s}_k \mathbf{m}_k \bar{\mathbf{s}} \bar{\mathbf{m}} \right)$
- 6: $\mathbf{u}_k \leftarrow \mathbf{u}_k + \gamma \left(\mathbf{s}_k \bar{\mathbf{s}} \right)$
- 7: end for
- 8: Gather \mathbf{m}_k , \mathbf{v}_k and \mathbf{s}_k , \mathbf{u}_k from all clients.
- 9: $\bar{\mathbf{m}} \leftarrow (1 \alpha) \operatorname{Mean}(\mathbf{s}_{1:K}\mathbf{m}_{1:K}) + \alpha \operatorname{Sum}(\mathbf{v}_{1:K})$
- 10: $\bar{\mathbf{s}} \leftarrow (1 \alpha) \operatorname{Mean}(\mathbf{s}_{1:K}) + \alpha \left[\delta \mathbf{1} + \operatorname{Sum}(\mathbf{u}_{1:K}) \right]$

```
11: \bar{\mathbf{m}} \leftarrow \bar{\mathbf{m}}/\bar{\mathbf{s}}
```

12: end while

 \triangleright Two additional steps for precision \bar{s}

Method	Test accur 10 rounds	acy († largei 25 rounds	is better) 50 rounds
FedAvg	$72.3{\pm}0.4$	77.7±0.3	80.0±0.2
FedProx	72.2 ± 0.3	77.4 ± 0.1	$80.3 {\pm} 0.1$
FedDyn	$75.3 {\pm} 0.8$	$77.5 {\pm} 0.8$	$78.2{\pm}0.5$
FedLap	72.1 ± 0.2	77.1 ± 0.1	$80.2{\pm}0.1$
FedLap-Cov	$75.0{\pm}0.6$	$79.8 {\pm} 0.4$	$81.8{\pm}0.1$
BayesADMM@m	80.4±0.2	83.1±0.1	83.4±+6%
BayesADMM	80.6 ±0.2	83.5 ±0.1	84.1 ∃
FedAvg	$70.4{\pm}0.9$	74.3±0.5	76.0±0.7
FedProx	$69.9 {\pm} 0.4$	74.7 ± 0.6	$76.9 {\pm} 0.9$
FedDyn	$73.0 {\pm} 0.6$	$74.6 {\pm} 0.4$	$74.6 {\pm} 0.5$
FedLap	$71.3 {\pm} 0.9$	74.3 ± 0.4	$77.6 {\pm} 0.7$
FedLap-Cov	$74.6 {\pm} 0.7$	$78.3 {\pm} 1.0$	80.5 ± 0.6
BayesADMM@m	77.0 ±0.8	81.4 ±0.4	82.1 ∃ <mark>+8%</mark>
BayesADMM	77.0 ±0.8	81.5 ±0.5	82.3 ∃
FedAvg	62.8±3.1	65.4±1.8	66.0±1.5
FedProx	64.3 ±2.0	65.9 ± 1.6	66.3±1.4
FedDyn	63.6±1.1	64.7±0.6	65.4±1.0
FedLap	$60.2 {\pm} 2.4$	$66.4{\pm}1.1$	66.5 ± 1.2
FedLap-Cov	$58.4{\pm}2.4$	$65.4{\pm}1.1$	67.5±1.2
BayesADMM@m	63.8 ±1.4	69.5 ±0.8	70.2∃ <mark>+5%</mark>
BayesADMM	63.8 ±1.4	69.5 ±0.8	70.3 ∃
	Method FedAvg FedProx FedDyn FedLap-Cov BayesADMM@m BayesADMM@m FedProx FedDyn FedLap-Cov BayesADMM@m BayesADMM@m BayesADMM@m FedLap-Cov	Method10 roundsFedAvg 72.3 ± 0.4 FedProx 72.2 ± 0.3 FedDyn 75.3 ± 0.8 FedLap 72.1 ± 0.2 FedLap-Cov 75.0 ± 0.6 BayesADMM@m 80.4 ± 0.2 BayesADMM 80.6 ± 0.2 FedAvg 70.4 ± 0.9 FedProx 69.9 ± 0.4 FedDyn 73.0 ± 0.6 FedLap-Cov 74.6 ± 0.7 BayesADMM@m 77.0 ± 0.8 BayesADMM 77.0 ± 0.8 FedAvg 62.8 ± 3.1 FedAvg 62.8 ± 3.1 FedProx 64.3 ± 2.0 FedDyn 63.6 ± 1.1 FedLap-Cov 58.4 ± 2.4 BayesADMM@m 63.8 ± 1.4 BayesADMM@m 63.8 ± 1.4	Method10 rounds25 roundsFedAvg 72.3 ± 0.4 77.7 ± 0.3 FedProx 72.2 ± 0.3 77.4 ± 0.1 FedDyn 75.3 ± 0.8 77.5 ± 0.8 FedLap 72.1 ± 0.2 77.1 ± 0.1 FedLap-Cov 75.0 ± 0.6 79.8 ± 0.4 BayesADMM@m 80.4 ± 0.2 83.1 ± 0.1 BayesADMM 80.6 ± 0.2 83.5 ± 0.1 FedAvg 70.4 ± 0.9 74.3 ± 0.5 FedProx 69.9 ± 0.4 74.7 ± 0.6 FedDyn 73.0 ± 0.6 74.6 ± 0.4 FedLap-Cov 74.6 ± 0.7 78.3 ± 1.0 BayesADMM@m 77.0 ± 0.8 81.4 ± 0.4 BayesADMM 77.0 ± 0.8 81.5 ± 0.5 FedAvg 62.8 ± 3.1 65.4 ± 1.8 FedProx 64.3 ± 2.0 65.9 ± 1.6 FedDyn 63.6 ± 1.1 64.7 ± 0.6 FedLap-Cov 58.4 ± 2.4 65.4 ± 1.1 BayesADMM@m 63.8 ± 1.4 69.5 ± 0.8 BayesADMM 63.8 ± 1.4 69.5 ± 0.8

Adaptive Bayesian Intelligence

- Adaptive Intelligence = Bayesian Computation
- Part 1: Bayesian Learning Rule [1]
 - Foundational way to derive learning-algorithms
 - Application to Deep Learning [2]
- Part 2: Posterior Correction [3]
 - Foundational way to derive adaptation-algorithms
 - Application to continual learning [4-5]
 - But also for LLM merging, Federated Learning etc.
- Adaptive Bayesian Intelligence: A roadmap.
- 1. Khan and Rue, The Bayesian Learning Rule, JMLR (2023)
- 2. Shen et al. Variational Learning is Effective for Large Deep Networks, ICML (2024)
- 3. Khan. Knowledge Adaptation as Posterior Correction, arXiv (2025)
- 4. Khan and Swaroop. Knowledge-Adaptation Priors, NeurIPS (2021).
- 5. Pan et al. Continual Deep Learning by Functional Regularisation of Memorable Past, NeurIPS, 2020

Questions for the future

- What should the algorithm remember?
- And what new experiences should it seek?
- Memory should be chosen to minimize the corrections that may arise in the future.
- New experiences should be chosen to enable easyenough corrections (not too daunting for the learner)
- Future is unknown but the algorithm has the freedom to explore by "fixing the past & choosing the future"



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

(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 (June 25-27, 2025)

https://bayesduality.github.io/workshop_2025.html



Abeba Birhane

Trinity College Dublin, Ireland







André Martins Instituto Superior

Tecnico, Portugal

Razvan Pascanu



Deepmind, UK

Marcus Rohrbach

TU Darmstadt, Germany



Mark van der Wilk

University of Oxford, UK



David Rügamer

Ludwig-Maximilians-Universität München, Germany Diverse topics: Bayes, Optimization, Information Geometry, Continual Learning, Federated Learning, Active learning, RL, Model understanding, Data Attributions, LLMs, etc.

Adaptive Bayesian Intelligence Team

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



Emtiyaz Khan Team Leader



Thomas Möllenhoff **Research Scientist**



Hugo Monzón Special Postdoctoral Maldonado Postdoctoral Researcher



Christopher Johannes Anders Postdoctoral Researcher



Yohan Jung Postdoctoral Researcher



Part-Time Student

The University of

Tokyo

Bai Cong Part-Time Student Tokvo Institute of Technology



Eiki Shimizu Part-Time Student Institute of Statistical Mathematics



Giulia Lanzillotta Intern ETH Zurich



Researcher

RIKEN BDR

Adrian R. Minut Intern Sapienza, University of Rome



Florian Seligmann Intern Karlsruhe Institute of Technology



Guiomar Pescador Barrios Intern Imperial College London



Henrique Da Silva Gameiro Intern EPFL. Switzerland



Visiting Scientist

University of



Pierre Alguier Visiting Scientist ESSEC Business Winsconsin-Madison School



Geoffrey Wolfer Visiting Scientist Waseda University



Rio Yokota Visiting Scientist Tokvo Institute of Technology



Remote Collaborator University of Amsterdam

And many of our collaborators!

