

# Uncertainty Estimation and Generalization Bounds for Modern Deep Learning

Advances in Function-Space Variational Inference, Linearized Laplace Approximation, Deep Ensembles, and Chernoff-Based Generalization Bounds

---

Luis Antonio Ortega Andrés

Supervisor: Daniel Hernández Lobato

March 7, 2026

Department of Computer Science and Telecommunications  
Escuela Politécnica Superior  
Universidad Autónoma de Madrid



# Table of contents

1. Part I: Uncertainty estimation via function-space methods
  - DVIP: Deep Variational Implicit Processes
  - Post-hoc Bayesian deep learning: VaLLA and FMGP
2. Part II: Generalization theory via diversity and large deviations
  - Diversity and generalization in ensembles
  - PAC–Chernoff bounds and the rate function
  - Implicit bias of Stochastic Gradient Descent
3. Conclusions and Future Work

## Deep learning is accurate, but:

- **Uncertainty** is often miscalibrated
  - overconfidence on OOD inputs
  - unreliable risk estimates
- **Generalization** at scale is not fully understood
  - over-parameterization / interpolation
  - why SGD solutions generalize?

## Thesis perspective

Unify three tools:

- Bayesian inference
- Function-space modeling
- Large-deviation theory

## Methodological goal

- Scalable Bayesian uncertainty for deep models.
- Function-space inference as the main tool, rather than parameter-space.

## Theoretical goal

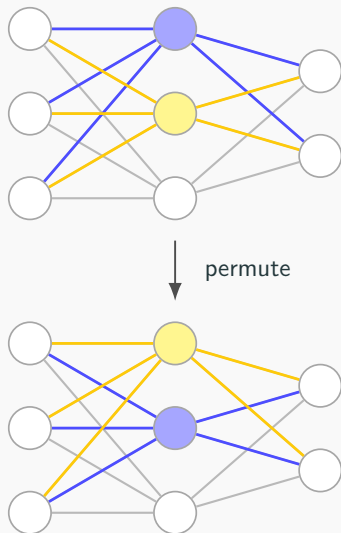
- Distribution-dependent view of generalization.
- Informative beyond capacity-only bounds (including interpolation).

# **Part I: Uncertainty estimation via function-space methods**

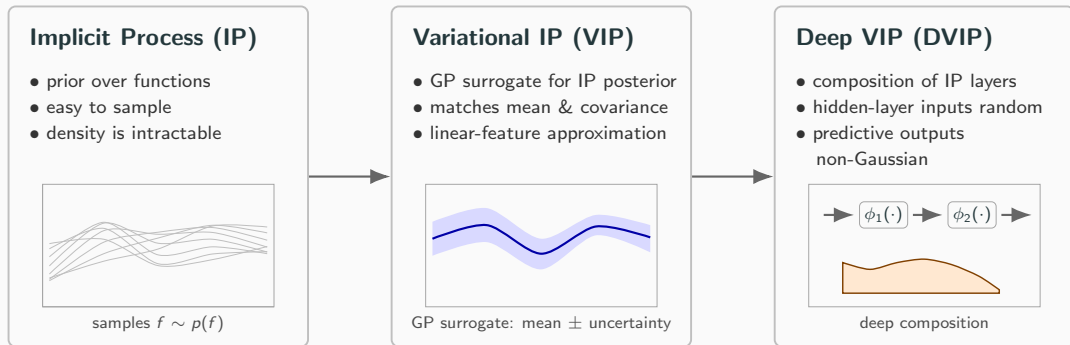
---

# Why function-space inference?

- Parameter-space posteriors are difficult in deep nets:
  - symmetries / non-identifiability
  - strong parameter dependencies
- Function-space inference is often better behaved:
  - directly targets predictions
  - avoids many parameterization pathologies
- Reference models in function space
  - Gaussian processes (GPs): canonical Bayesian prior on functions.
  - Implicit processes: sample-based priors that generalize GPs.



# Implicit Processes $\rightarrow$ VIP $\rightarrow$ DVIP



# DVIP definition

- **Implicit Process (IP).** An implicit process is a distribution over functions  $f : \mathcal{X} \rightarrow \mathbb{R}$  specified by a *sampling procedure* rather than a tractable density:

$$f \sim \text{IP}(g, P) \iff \exists z \sim P \text{ s.t. } f(\cdot) = g(\cdot, z).$$

(Sampling is easy via  $z \sim P$ ; the function-density is generally intractable)

- **Deep composition of IP layers.** For datapoint  $\mathbf{x}_n$  and layer  $l$ :

$$\mathbf{f}_{h,n}^l = f_h^l(\mathbf{f}_{\cdot,n}^{l-1}), \quad \mathbf{f}_{\cdot,n}^0 = \mathbf{x}_n.$$

where each hidden unit is an independent IP draw

$$f_h^l(\cdot) \sim \text{IP}(g_h^l, P_h^l), \quad h = 1, \dots, H_l.$$

Randomness propagates through layers because the inputs  $\mathbf{f}_{\cdot,n}^{l-1}$  are themselves random.

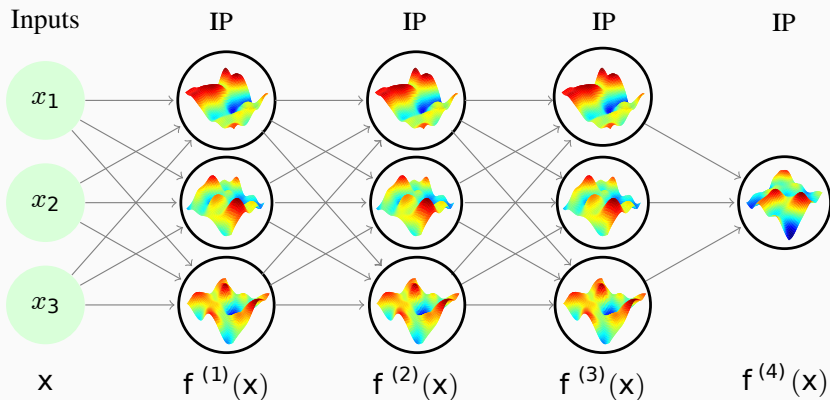
- **VIP/GP-style linear surrogate per unit (practical inference).** Approximate each unit-function by a linear model in learned features:

$$\mathbf{f}_{h,n}^l = \phi_h^l(\mathbf{f}_{\cdot,n}^{l-1})^\top \mathbf{a}_h^l + m_h^l(\mathbf{f}_{\cdot,n}^{l-1}),$$

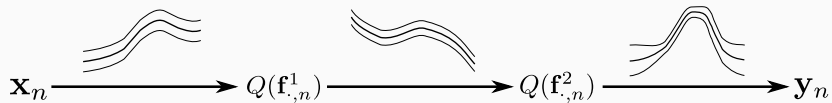
where  $\phi_h^l(\cdot)$  and  $m_h^l(\cdot)$  are constructed from the sampled IP, and  $\mathbf{a}$  are Gaussian weights.

# DVIP architecture

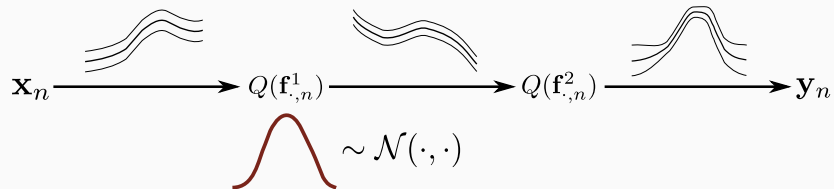
Fully connected DVIP: input propagated through  $L$  IP layers (each layer has  $H_i$  independent IP units).



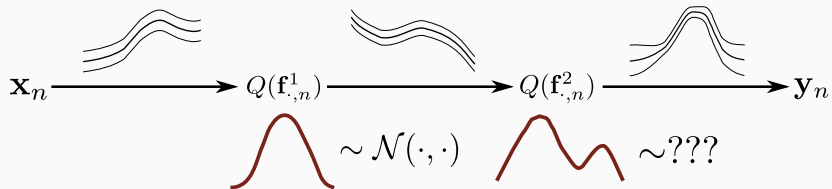
# Stochasticity Propagation



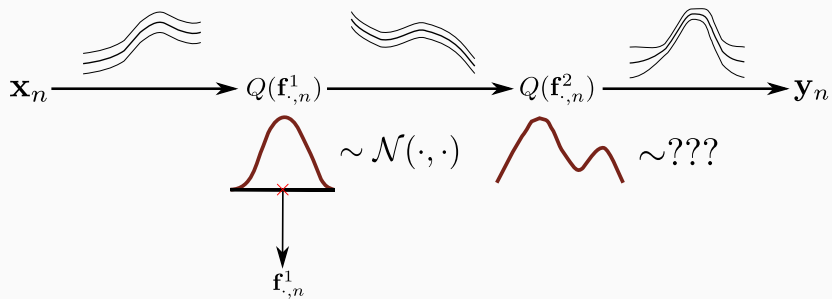
# Stochasticity Propagation



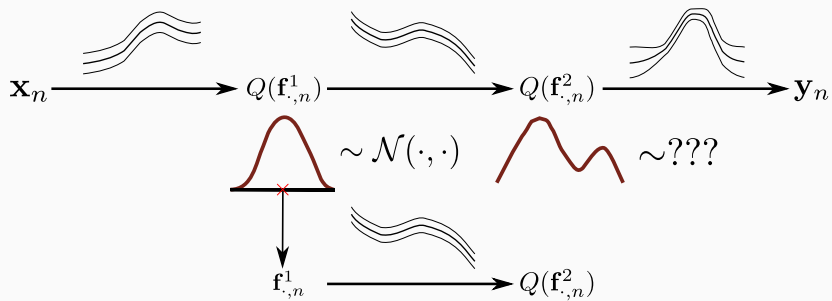
# Stochasticity Propagation



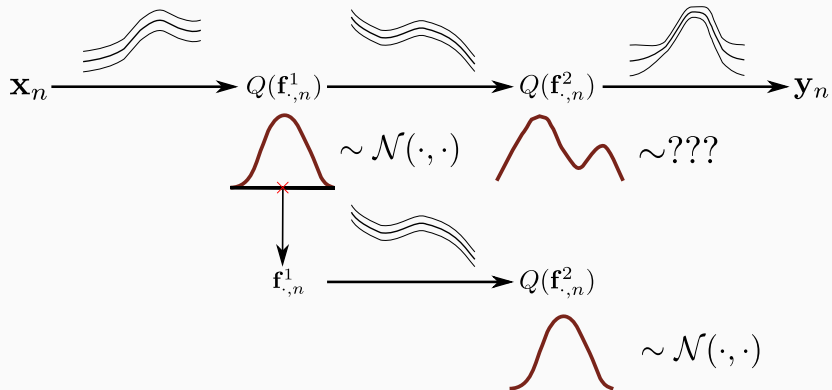
# Stochasticity Propagation



# Stochasticity Propagation

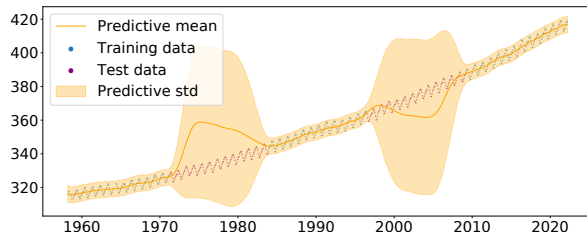
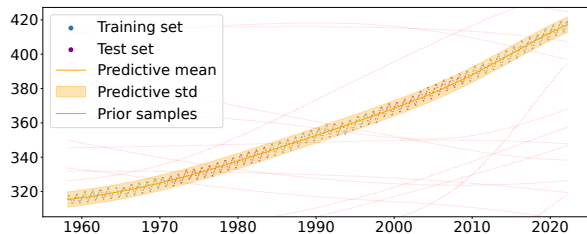


# Stochasticity Propagation



- **Inference objective (function space)**
  - optimize a variational approximation to the posterior over functions/layer outputs
  - gradients estimated with Monte Carlo samples
- **Scalable optimization**
  - mini-batches over data points
  - stochastic estimates of the ELBO
- **Why predictions are non-Gaussian**
  - each layer receives *random* inputs (previous layer outputs)
  - composition of random mappings breaks Gaussianity
- **Prediction in practice**
  - sample functions per layer
  - forward-propagate samples through the network
  - use a Gaussian mixture as predictive distribution

## CO2 Prediction Dataset - DVIP Top and DGP Bottom



- **Objective**
  - start from a strong pretrained deterministic DNN
  - add calibrated predictive uncertainty
- **Baseline: Linearized Laplace Approximation (LLA)**
  - strong performance in practice
  - scalability bottleneck: Jacobian / Hessian-related computations
- **This thesis: scalable function-space alternatives**
  - **VaLLA**: sparse variational GP surrogate of LLA (function space)
  - **FMGP**: fixed mean = pretrained predictor; learn covariance via VI

- **Network / MAP:**  $f(\mathbf{x}; \theta)$  is the network output at input  $\mathbf{x}$ ,  $\theta_{\text{MAP}}$  are MAP parameters.
- **Local Gaussian posterior (Laplace):**

$$P(\theta|D) \approx \mathcal{N}(\theta_{\text{MAP}}, H^{-1}),$$

where  $H = \nabla_{\theta}^2 [-\log P(\theta|D)]|_{\theta_{\text{MAP}}}$  (often replaced by GGN Approximation).

- **Jacobian:**

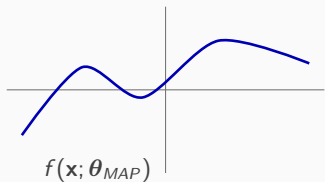
$$J_{\theta_{\text{MAP}}}(\mathbf{x}) := \nabla_{\theta} f(\mathbf{x}; \theta)|_{\theta_{\text{MAP}}}.$$

- **Linearization induces function-space GP:**

$$f(\cdot) \sim \text{GP}(m(x), k(x, x')), \quad m(\mathbf{x}) := f(\mathbf{x}; \theta_{\text{MAP}}), \quad k(\mathbf{x}, \mathbf{x}') := J_{\theta_{\text{MAP}}}(\mathbf{x})H^{-1}J_{\theta_{\text{MAP}}}(\mathbf{x}')^{\top}.$$

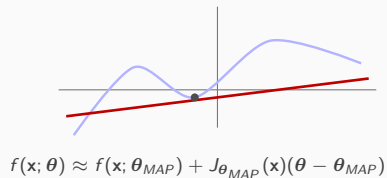
- **Dimensions (typical):**  $J_{\theta_{\text{MAP}}}(\mathbf{x}) \in \mathbb{R}^{d_{\text{out}} \times d_{\theta}}$ ,  $H \in \mathbb{R}^{d_{\theta} \times d_{\theta}}$ , so  $k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}^{d_{\text{out}} \times d_{\text{out}}}$ .

## 1) Pre-trained model



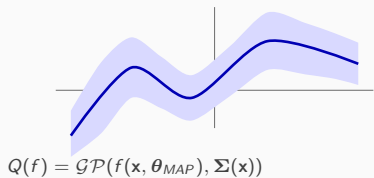
linearize

## 2) Linearization at $\theta_{MAP}$



decoupled sparse Gaussian process

## 4) Prediction (VaLLA)



use  $Q$

## 3) Variational inference

maximize (ELBO):

$$\mathcal{L} = \mathbb{E}_{Q(f)}[\log P(\mathcal{D}|f)] - \text{KL}(Q(f) \| P(f))$$

# Function-space inference with Decoupled Sparse Gaussian Processes

- Notation

- Decoupled inducing sets:  $\mathbf{Z}_\alpha = \{\mathbf{z}_{\alpha,j}\}_{j=1}^{m_\alpha}$  (mean),  $\mathbf{Z}_\beta = \{\mathbf{z}_{\beta,j}\}_{j=1}^{m_\beta}$  (covariance)
- GP prior:  $P(f) = \mathcal{GP}(0, k)$

- Decoupled sparse variational posterior in function space

$$Q(f) = \mathcal{GP}(m_Q(\cdot), k_Q(\cdot, \cdot)),$$

with the standard decoupled parametrization

$$m_Q(\mathbf{x}) = K_{\mathbf{x}, \mathbf{z}_\alpha} \mathbf{a}, \quad k_Q(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}, \mathbf{x}') + k_{\mathbf{x}, \mathbf{z}_\beta} K_{\mathbf{z}_\beta, \mathbf{z}_\beta}^{-1} (\mathbf{S} - K_{\mathbf{z}_\beta, \mathbf{z}_\beta}) K_{\mathbf{z}_\beta, \mathbf{z}_\beta}^{-1} k_{\mathbf{z}_\beta, \mathbf{x}'},$$

where  $\mathbf{a} \in \mathbb{R}^{m_\alpha \times d_{\text{out}}}$  controls the mean, and  $\mathbf{S} \succeq 0$  controls the posterior covariance.

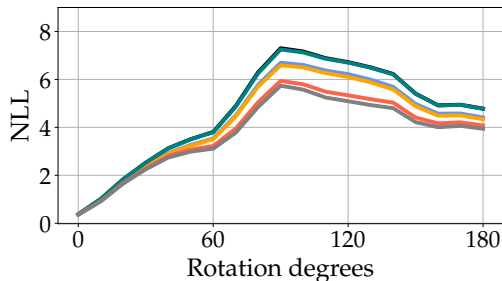
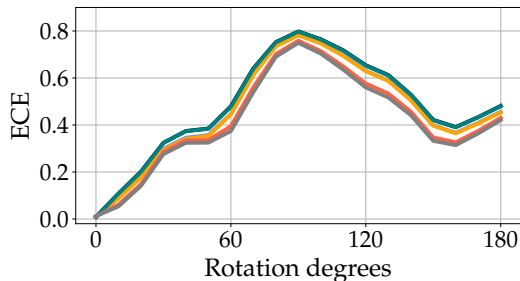
$$\mathbf{Z}_\alpha \text{ and } \mathbf{a} \text{ are fixed} \rightarrow m_Q(\mathbf{x}) \approx f(\mathbf{x}, \boldsymbol{\theta}_{\text{MAP}}).$$

- Function-space variational objective (ELBO)

$$\log P(\mathcal{D}) \geq \mathcal{L}(Q) := \mathbb{E}_{Q(f)} \left[ \sum_{i=1}^n \log p(y_i | f(\mathbf{x}_i)) \right] - \text{KL}(Q(f) \| P(f)).$$

# VaLLA Results on Fashion MNIST

Model	ACC	NLL	ECE	BRIER	OOD-AUC
MAP	86.6	0.373	<b>0.009</b>	0.193	0.874
LLA* KFAC	86.6	0.373	<b>0.008</b>	0.193	0.880
LLA*	86.6	0.373	<b>0.008</b>	0.193	0.882
ELLA	86.6	0.373	<b>0.008</b>	0.193	0.874
VaLLA 100	<b>87.4</b>	<b>0.335</b>	0.011	<b>0.182</b>	<b>0.923</b>
VaLLA 200	<b>87.6</b>	<b>0.332</b>	0.013	<b>0.181</b>	<b>0.933</b>



- **Function prior as a Gaussian measure**

- Work on a separable Hilbert space  $\mathcal{H}$  of functions.
- Gaussian measure on (a suitable completion of)  $\mathcal{H}$ :

$$P = \mathcal{N}(g, \mathcal{C}),$$

where  $g \in \mathcal{H}$  is a fixed mean and  $\mathcal{C} : \mathcal{H} \rightarrow \mathcal{H}$  is a positive, self-adjoint covariance operator.

- **Infer only uncertainty (operator view)**

- keep the mean fixed at  $g$  and update only the covariance operator
- sparse/low-rank uncertainty updates via subspaces of  $\mathcal{H}$

- **Practical benefit**

- Bayesian uncertainty on top of a frozen network predictor (mean  $f(\cdot, \theta_{MAP})$ )
- avoids Jacobians / curvature of the DNN (inference happens in function space)

- **Posterior family as Gaussian measures**
  - Approximate posterior:  $Q = \mathcal{N}(m^*, \mathcal{C}^*)$ .
- **Mean constraint (Hilbert norm)**

$$\|m^* - f(\cdot, \boldsymbol{\theta}_{MAP})\|_{\mathcal{H}} \leq \varepsilon \quad (\text{FMGP enforces } m^* \approx f(\cdot, \boldsymbol{\theta}_{MAP})).$$

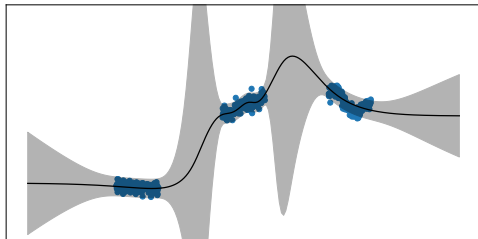
- **Sparse covariance update via a subspace**
  - Choose a finite-dimensional subspace  $U = \text{span}\{\phi_1, \dots, \phi_{m_\beta}\} \subset \mathcal{H}$
  - Variational parameter:  $\mathbf{A} \succeq 0$  acting on  $\mathbb{R}^{m_\beta}$

$$K_{\mathbf{Z}_\beta, \mathbf{A}}^*(\mathbf{x}, \mathbf{x}') = K(\mathbf{x}, \mathbf{x}') - K(\mathbf{x}, \mathbf{Z}_\beta)(\mathbf{A} + K(\mathbf{Z}_\beta, \mathbf{Z}_\beta))^{-1}K(\mathbf{Z}_\beta, \mathbf{x}').$$

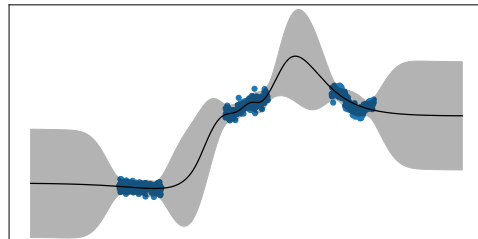
- **Learning:** optimize (VI) over  $(\mathbf{A}, \mathbf{Z}_\beta)$  and kernel hyperparameters.

## Predictive distribution of LLA and FMGP

LLA



FMGP



## **Part II: Generalization theory via diversity and large deviations**

---

## Ensemble diversity: core decomposition

For a distribution  $\rho$  over predictors  $\Theta$ :

$$L(\rho) \leq \alpha \left( \mathbb{E}_\rho[L(\boldsymbol{\theta})] - \mathbb{D}(\rho) \right),$$

where  $\alpha = 1$  (squared or cross-entropy) and  $\alpha = 4$  (0/1), and for squared loss the inequality is exact.

Diversity is captured by variance across ensemble members.

Examples (under data distribution  $\nu$ ):

$$\mathbb{D}_{\text{sq}}(\rho) := \mathbb{E}_{\nu} [\text{Var}_{\rho}(h_R(x; \boldsymbol{\theta}))],$$

$$\mathbb{D}_{\text{ce}}(\rho) := \mathbb{E}_{\nu} \left[ \text{Var}_{\rho} \left( \frac{P(y|\mathbf{x}, \boldsymbol{\theta})}{\sqrt{2} \max_{\boldsymbol{\theta}} P(y|\mathbf{x}, \boldsymbol{\theta})} \right) \right],$$

$$\mathbb{D}_{0/1}(\rho) := \mathbb{E}_{\nu} [\text{Var}_{\rho}(\mathbf{1}(h_W(x; \boldsymbol{\theta}) \neq y))],$$

and, more generally,

$$\mathbb{D}(\rho) = \mathbb{E}_{\nu} [\text{Var}_{\rho}(f(y, \mathbf{x}; \boldsymbol{\theta}))].$$

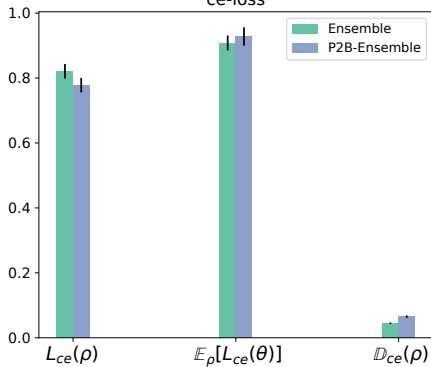
Key message: beneficial diversity reduces ensemble loss by reducing correlated errors.

**Theorem (PAC–Bayes bound).** For any prior  $\pi$  over  $\Theta$  independent of  $D$ , any  $\delta \in (0, 1)$  and any  $\lambda > 0$ , with probability at least  $1 - \delta$  over draws  $D \sim \nu^n$ , for all posteriors  $\rho$  over  $\Theta$  (simultaneously),

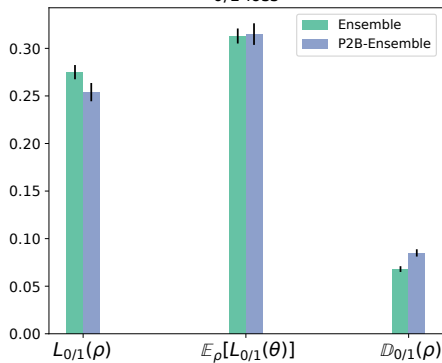
$$L(\rho) \leq \alpha \left( \mathbb{E}_\rho[\hat{L}(\theta, D)] - \hat{D}(\rho, D) + \frac{2 \text{KL}(\rho \| \pi) + \varepsilon}{\lambda n} \right).$$

**Interpretation:** to generalize we want (i) low empirical risk, (ii) *high diversity*  $\hat{D}$ , and (iii) a posterior  $\rho$  *close to the prior*  $\pi$  (small KL).

CIFAR-10 LeNet5  
ce-loss



CIFAR-10 LeNet5  
0/1-loss



# Rate function and inverse rate function

Define cumulant-generating function:

$$J_{\theta}(\lambda) = \log \mathbb{E}_{\nu} [\exp (\lambda(L(\theta) - \ell(y, \mathbf{x}, \theta)))] .$$

Rate function (Legendre transform):

$$\mathcal{I}_{\theta}(a) = \sup_{\lambda > 0} (\lambda a - J_{\theta}(\lambda)) .$$

Inverse rate:

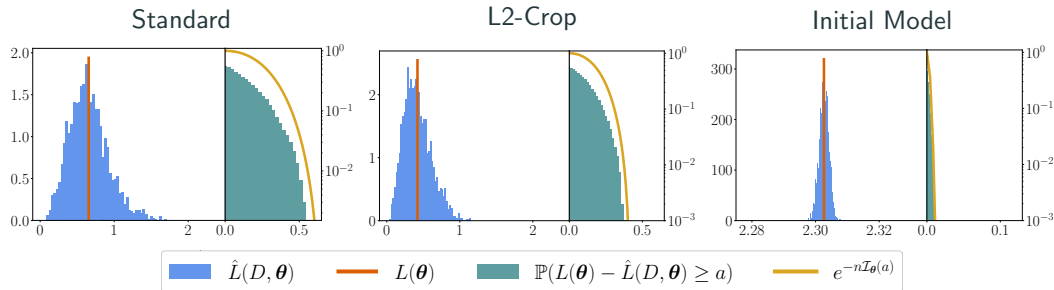
$$\mathcal{I}_{\theta}^{-1}(s) = \inf_{\lambda > 0} \frac{J_{\theta}(\lambda) + s}{\lambda} .$$

# Chernoff bound

For fixed  $\theta$  and  $a > 0$ :

$$\mathbb{P}_{D \sim \nu^n} \left( L(\theta) - \hat{L}(D, \theta) \geq a \right) \leq \exp \left( - n \mathcal{I}_\theta(a) \right).$$

Key message: if  $\mathcal{I}_\theta(\cdot)$  is large, large deviations are exponentially unlikely.

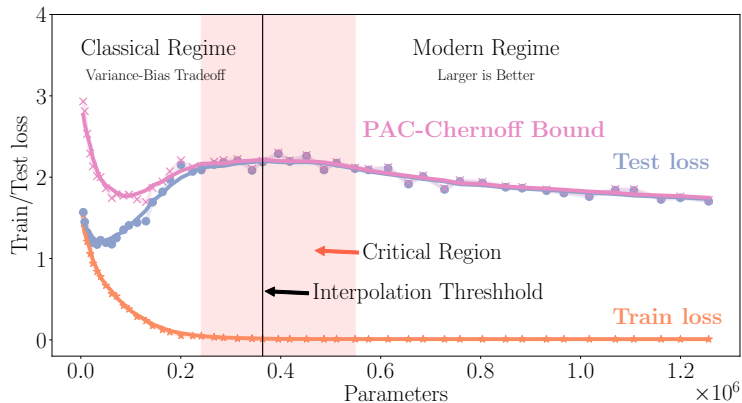


Distribution of empirical loss, tail probabilities and Chernoff bound for three Inception models on CIFAR-10 ( $n = 50$ ).

# PAC-Chernoff Bound

With h.p.  $1 - \delta$  over  $D \sim \nu^n$ , for all  $\theta \in \Theta$ , simultaneously,

$$L(\theta) \leq \hat{L}(D, \theta) + \mathcal{I}_{\theta}^{-1}\left(\frac{1}{n} \log \frac{k^p}{\delta}\right).$$



# Optimal regularization at interpolation

Define (approximate) interpolators with  $\hat{L}(D, \theta) \leq \varepsilon$ . Let:

$$\theta_\varepsilon^* = \arg \min_{\hat{L}(D, \theta) \leq \varepsilon} L(\theta), \quad \theta_\varepsilon^\times = \arg \min_{\hat{L}(D, \theta) \leq \varepsilon} \left( \hat{L}(D, \theta) + \mathcal{I}_\theta^{-1} \left( \frac{1}{n} \log \frac{k^p}{\delta} \right) \right).$$

With high probability,

$$|L(\theta_\varepsilon^*) - L(\theta_\varepsilon^\times)| \leq \varepsilon.$$

Interpretation: minimizing the inverse rate function is an optimal regularization principle for interpolating regimes.

# Understanding Existing Regularizers

Many common regularization techniques are **approximations** to the **optimal regularizer**:

- **Distance from initialization and  $\ell_2$ -norm:**

$$\mathcal{I}_{\theta}^{-1}\left(\frac{1}{n} \ln \frac{k^p}{\delta}\right) \leq \sqrt{2Ma} \|\theta\|_2,$$

- **Input-gradient norm:**

$$\mathcal{I}_{\theta}^{-1}\left(\frac{1}{n} \ln \frac{k^p}{\delta}\right) \leq \sqrt{\frac{1}{n} \ln \frac{k^p}{\delta}} \sqrt{M \mathbb{E}_{\nu} \left[ \|\nabla_{\mathbf{x}} \ell(\mathbf{y}, \mathbf{x}, \theta)\|_2^2 \right]}.$$

The distribution-dependent PAC-Chernoff Bound can be used to obtain **bounds over the number of parameters** of **interpolators**:

For any  $\epsilon \in (0, L^*)$  and any  $\delta \in (0, 1)$ , with high probability  $1 - \delta$  over  $D \sim \nu^n$ , for all  $\theta \in \Theta$ , simultaneously,

$$\text{if } \hat{L}(D, \theta) \leq \epsilon \text{ then } p \geq \frac{n\mathcal{I}_{\theta}(L^* - \epsilon) + \ln \delta}{\ln k}.$$

where  $L^* = \arg \min_{\theta} L(\theta)$ .

## SGD implicit bias via large deviations

Defining the **abnormality** rate of a model as

$$\begin{aligned}\alpha(D, \theta) &:= \mathcal{I}_\theta(L(\theta) - \hat{L}(D, \theta)), \\ &\asymp -\frac{1}{n} \ln \mathbb{P}_{S \sim \nu^n} \left( L(\theta) - \hat{L}(S, \theta) \geq L(\theta) - \hat{L}(D, \theta) \right)\end{aligned}$$

raises a decomposition on the empirical loss given by

$$\hat{L}(D, \theta) = L(\theta) + \mathcal{I}_\theta^{-1}(\alpha(D, \theta)).$$

The cumulative distribution of  $\alpha(D, \theta)$  satisfies

$$\begin{aligned}\forall s > 0 \quad \mathbb{P}_{D \sim \nu^n}(\alpha(D, \theta) \geq s) &\leq e^{-n|s|}, \\ \forall s < 0 \quad \mathbb{P}_{D \sim \nu^n}(\alpha(D, \theta) \leq s) &\leq e^{-n|s|},\end{aligned}$$

In conclusion, the distribution of empirical loss for *large*  $n$  values can be expressed as:

$$\hat{L}(D, \theta) \approx L(\theta) - \mathcal{I}_\theta^{-1}(s), \quad s \sim \text{Laplace}(0, n).$$

# Implicit Bias of Gradient Descent (GD)

$$\hat{L}(D, \theta) = L(\theta) - \mathcal{I}_\theta^{-1}(\alpha(D, \theta))$$

- **Two implicit incentives beyond lowering  $L(\theta)$ :**
  - Increase **abnormality**  $\alpha(D, \theta)$  (move toward rare/left-tail realizations).
  - Increase  $\mathcal{I}_\theta^{-1}(\cdot)$  for a given  $\alpha$ , i.e. implicitly favor **less concentrated** empirical-loss distributions (smaller  $\mathcal{I}_\theta$ ).
- **Consequence:** GD can drive  $\hat{L}(D, \theta)$  down by exploiting unusually favorable deviations, often yielding **large generalization error**.

# Implicit Bias of Stochastic Gradient Descent (SGD)

- SGD objective (mini-batches  $B$ , size  $m$ ):

$$\min_{\theta} \hat{L}(B, \theta) \quad \text{and} \quad \mathbb{E}_{B \sim D} [\hat{L}(B, \theta)] = \hat{L}(D, \theta).$$

# Implicit Bias of Stochastic Gradient Descent (SGD)

- **SGD objective (mini-batches  $B$ , size  $m$ ):**

$$\min_{\theta} \hat{L}(B, \theta) \quad \text{and} \quad \mathbb{E}_{B \sim D} [\hat{L}(B, \theta)] = \hat{L}(D, \theta).$$

- **Batch-level decomposition:**

$$\hat{L}(B, \theta) = L(\theta) - \mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B) \Rightarrow \hat{L}(D, \theta) = L(\theta) - \mathbb{E}_{B \sim D} [\mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B)],$$

so dynamics depend on the batch-induced distribution  $Q(\alpha_{\theta}^B | D, \theta)$ .

# Implicit Bias of Stochastic Gradient Descent (SGD)

- **SGD objective (mini-batches  $B$ , size  $m$ ):**

$$\min_{\theta} \hat{L}(B, \theta) \quad \text{and} \quad \mathbb{E}_{B \sim D} [\hat{L}(B, \theta)] = \hat{L}(D, \theta).$$

- **Batch-level decomposition:**

$$\hat{L}(B, \theta) = L(\theta) - \mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B) \Rightarrow \hat{L}(D, \theta) = L(\theta) - \mathbb{E}_{B \sim D} [\mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B)],$$

so dynamics depend on the batch-induced distribution  $Q(\alpha_{\theta}^B | D, \theta)$ .

- **Why SGD differs from GD:** a Taylor expansion around  $\mathbb{E}[\alpha_{\theta}^B]$  introduces an explicit **variance correction**:

# Implicit Bias of Stochastic Gradient Descent (SGD)

- **SGD objective (mini-batches  $B$ , size  $m$ ):**

$$\min_{\theta} \hat{L}(B, \theta) \quad \text{and} \quad \mathbb{E}_{B \sim D} [\hat{L}(B, \theta)] = \hat{L}(D, \theta).$$

- **Batch-level decomposition:**

$$\hat{L}(B, \theta) = L(\theta) - \mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B) \Rightarrow \hat{L}(D, \theta) = L(\theta) - \mathbb{E}_{B \sim D} [\mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B)],$$

so dynamics depend on the batch-induced distribution  $Q(\alpha_{\theta}^B | D, \theta)$ .

- **Why SGD differs from GD:** a Taylor expansion around  $\mathbb{E}[\alpha_{\theta}^B]$  introduces an explicit **variance correction**:
  - larger  $\mathbb{E}[\alpha_{\theta}^B]$  tends to **increase** the deviation effect,

# Implicit Bias of Stochastic Gradient Descent (SGD)

- **SGD objective (mini-batches  $B$ , size  $m$ ):**

$$\min_{\theta} \hat{L}(B, \theta) \quad \text{and} \quad \mathbb{E}_{B \sim D} [\hat{L}(B, \theta)] = \hat{L}(D, \theta).$$

- **Batch-level decomposition:**

$$\hat{L}(B, \theta) = L(\theta) - \mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B) \Rightarrow \hat{L}(D, \theta) = L(\theta) - \mathbb{E}_{B \sim D} [\mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B)],$$

so dynamics depend on the batch-induced distribution  $Q(\alpha_{\theta}^B | D, \theta)$ .

- **Why SGD differs from GD:** a Taylor expansion around  $\mathbb{E}[\alpha_{\theta}^B]$  introduces an explicit **variance correction**:
  - larger  $\mathbb{E}[\alpha_{\theta}^B]$  tends to **increase** the deviation effect,
  - larger  $\text{Var}(\alpha_{\theta}^B)$  tends to **decrease** it.

# Implicit Bias of Stochastic Gradient Descent (SGD)

- **SGD objective (mini-batches  $B$ , size  $m$ ):**

$$\min_{\theta} \hat{L}(B, \theta) \quad \text{and} \quad \mathbb{E}_{B \sim D} [\hat{L}(B, \theta)] = \hat{L}(D, \theta).$$

- **Batch-level decomposition:**

$$\hat{L}(B, \theta) = L(\theta) - \mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B) \Rightarrow \hat{L}(D, \theta) = L(\theta) - \mathbb{E}_{B \sim D} [\mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B)],$$

so dynamics depend on the batch-induced distribution  $Q(\alpha_{\theta}^B | D, \theta)$ .

- **Why SGD differs from GD:** a Taylor expansion around  $\mathbb{E}[\alpha_{\theta}^B]$  introduces an explicit **variance correction**:
  - larger  $\mathbb{E}[\alpha_{\theta}^B]$  tends to **increase** the deviation effect,
  - larger  $\text{Var}(\alpha_{\theta}^B)$  tends to **decrease** it.
- **Batch size:** when  $m = n$ , variance  $\rightarrow 0$  and SGD reduces to GD; smaller  $m$  increases  $\text{Var}(\alpha_{\theta}^B)$ .

# Implicit Bias of Stochastic Gradient Descent (SGD)

- **SGD objective (mini-batches  $B$ , size  $m$ ):**

$$\min_{\theta} \hat{L}(B, \theta) \quad \text{and} \quad \mathbb{E}_{B \sim D} [\hat{L}(B, \theta)] = \hat{L}(D, \theta).$$

- **Batch-level decomposition:**

$$\hat{L}(B, \theta) = L(\theta) - \mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B) \Rightarrow \hat{L}(D, \theta) = L(\theta) - \mathbb{E}_{B \sim D} [\mathcal{I}_{\theta}^{-1}(\alpha_{\theta}^B)],$$

so dynamics depend on the batch-induced distribution  $Q(\alpha_{\theta}^B | D, \theta)$ .

- **Why SGD differs from GD:** a Taylor expansion around  $\mathbb{E}[\alpha_{\theta}^B]$  introduces an explicit **variance correction**:
  - larger  $\mathbb{E}[\alpha_{\theta}^B]$  tends to **increase** the deviation effect,
  - larger  $\text{Var}(\alpha_{\theta}^B)$  tends to **decrease** it.
- **Batch size:** when  $m = n$ , variance  $\rightarrow 0$  and SGD reduces to GD; smaller  $m$  increases  $\text{Var}(\alpha_{\theta}^B)$ .
- **Resulting bias:** smaller batches weaken the push toward highly abnormal deviations and toward broad (high-variance) loss distributions, typically improving **generalization**.

## Conclusions and Future Work

---

## Conclusion: A Unified Probabilistic View of Deep Learning

- **Thesis goal:** make uncertainty estimation scalable *and* make generalization analysis informative in modern (over-parameterized / interpolating) regimes.

## Conclusion: A Unified Probabilistic View of Deep Learning

- **Thesis goal:** make uncertainty estimation scalable *and* make generalization analysis informative in modern (over-parameterized / interpolating) regimes.
- **Function-space Bayesian uncertainty (Part I):**

# Conclusion: A Unified Probabilistic View of Deep Learning

- **Thesis goal:** make uncertainty estimation scalable *and* make generalization analysis informative in modern (over-parameterized / interpolating) regimes.
- **Function-space Bayesian uncertainty (Part I):**
  - **DVIP:** deep implicit-process priors + variational inference in function space  $\Rightarrow$  expressive *non-Gaussian* predictive uncertainty, competitive accuracy, and improved efficiency vs. deep GPs.

# Conclusion: A Unified Probabilistic View of Deep Learning

- **Thesis goal:** make uncertainty estimation scalable *and* make generalization analysis informative in modern (over-parameterized / interpolating) regimes.
- **Function-space Bayesian uncertainty (Part I):**
  - **DVIP:** deep implicit-process priors + variational inference in function space  $\Rightarrow$  expressive *non-Gaussian* predictive uncertainty, competitive accuracy, and improved efficiency vs. deep GPs.
  - **Post-hoc Bayesianization: VaLLA** (function-space surrogate of linearized Laplace) and **FMGP** (freeze mean to pretrained network, learn covariance)  $\Rightarrow$  calibrated uncertainty without performance loss.

# Conclusion: A Unified Probabilistic View of Deep Learning

- **Thesis goal:** make uncertainty estimation scalable *and* make generalization analysis informative in modern (over-parameterized / interpolating) regimes.
- **Function-space Bayesian uncertainty (Part I):**
  - **DVIP:** deep implicit-process priors + variational inference in function space  $\Rightarrow$  expressive *non-Gaussian* predictive uncertainty, competitive accuracy, and improved efficiency vs. deep GPs.
  - **Post-hoc Bayesianization: VaLLA** (function-space surrogate of linearized Laplace) and **FMGP** (freeze mean to pretrained network, learn covariance)  $\Rightarrow$  calibrated uncertainty without performance loss.
- **Distribution-dependent generalization theory (Part II):**

# Conclusion: A Unified Probabilistic View of Deep Learning

- **Thesis goal:** make uncertainty estimation scalable *and* make generalization analysis informative in modern (over-parameterized / interpolating) regimes.
- **Function-space Bayesian uncertainty (Part I):**
  - **DVIP:** deep implicit-process priors + variational inference in function space  $\Rightarrow$  expressive *non-Gaussian* predictive uncertainty, competitive accuracy, and improved efficiency vs. deep GPs.
  - **Post-hoc Bayesianization: VaLLA** (function-space surrogate of linearized Laplace) and **FMGP** (freeze mean to pretrained network, learn covariance)  $\Rightarrow$  calibrated uncertainty without performance loss.
- **Distribution-dependent generalization theory (Part II):**
  - **Diversity:** ensemble risk decomposes into average member error + a diversity (de-correlation) gain.

# Conclusion: A Unified Probabilistic View of Deep Learning

- **Thesis goal:** make uncertainty estimation scalable *and* make generalization analysis informative in modern (over-parameterized / interpolating) regimes.
- **Function-space Bayesian uncertainty (Part I):**
  - **DVIP:** deep implicit-process priors + variational inference in function space  $\Rightarrow$  expressive *non-Gaussian* predictive uncertainty, competitive accuracy, and improved efficiency vs. deep GPs.
  - **Post-hoc Bayesianization: VaLLA** (function-space surrogate of linearized Laplace) and **FMGP** (freeze mean to pretrained network, learn covariance)  $\Rightarrow$  calibrated uncertainty without performance loss.
- **Distribution-dependent generalization theory (Part II):**
  - **Diversity:** ensemble risk decomposes into average member error + a diversity (de-correlation) gain.
  - **PAC-Chernoff / rate-function view:** bounds remain informative at interpolation.

# Conclusion: A Unified Probabilistic View of Deep Learning

- **Thesis goal:** make uncertainty estimation scalable *and* make generalization analysis informative in modern (over-parameterized / interpolating) regimes.
- **Function-space Bayesian uncertainty (Part I):**
  - **DVIP:** deep implicit-process priors + variational inference in function space  $\Rightarrow$  expressive *non-Gaussian* predictive uncertainty, competitive accuracy, and improved efficiency vs. deep GPs.
  - **Post-hoc Bayesianization: VaLLA** (function-space surrogate of linearized Laplace) and **FMGP** (freeze mean to pretrained network, learn covariance)  $\Rightarrow$  calibrated uncertainty without performance loss.
- **Distribution-dependent generalization theory (Part II):**
  - **Diversity:** ensemble risk decomposes into average member error + a diversity (de-correlation) gain.
  - **PAC-Chernoff / rate-function view:** bounds remain informative at interpolation.
  - **SGD:** mini-batch noise induces an implicit regularization bias toward more stable (better-concentrated) solutions.

## Future Work

- DVIP extensions

- **DVIP extensions**
  - richer function-space posteriors beyond Gaussian/GP surrogates (e.g., flows or diffusion-based variational families);

- **DVIP extensions**

- richer function-space posteriors beyond Gaussian/GP surrogates (e.g., flows or diffusion-based variational families);
- stronger base approximations (e.g., sparse implicit processes) and *structured priors* (attention/recurrent for sequences, convolutional for spatial data);

- **DVIP extensions**

- richer function-space posteriors beyond Gaussian/GP surrogates (e.g., flows or diffusion-based variational families);
- stronger base approximations (e.g., sparse implicit processes) and *structured priors* (attention/recurrent for sequences, convolutional for spatial data);
- integration into autoencoders / generative architectures to unify representation learning and uncertainty.

- **DVIP extensions**
  - richer function-space posteriors beyond Gaussian/GP surrogates (e.g., flows or diffusion-based variational families);
  - stronger base approximations (e.g., sparse implicit processes) and *structured priors* (attention/recurrent for sequences, convolutional for spatial data);
  - integration into autoencoders / generative architectures to unify representation learning and uncertainty.
- **Post-hoc uncertainty (VaLLA / FMGP)**

- **DVIP extensions**
  - richer function-space posteriors beyond Gaussian/GP surrogates (e.g., flows or diffusion-based variational families);
  - stronger base approximations (e.g., sparse implicit processes) and *structured priors* (attention/recurrent for sequences, convolutional for spatial data);
  - integration into autoencoders / generative architectures to unify representation learning and uncertainty.
- **Post-hoc uncertainty (VaLLA / FMGP)**
  - VaLLA: relax the Gaussian linearized-posterior assumption for better calibration in highly non-linear regions;

- **DVIP extensions**

- richer function-space posteriors beyond Gaussian/GP surrogates (e.g., flows or diffusion-based variational families);
- stronger base approximations (e.g., sparse implicit processes) and *structured priors* (attention/recurrent for sequences, convolutional for spatial data);
- integration into autoencoders / generative architectures to unify representation learning and uncertainty.

- **Post-hoc uncertainty (VaLLA / FMGP)**

- VaLLA: relax the Gaussian linearized-posterior assumption for better calibration in highly non-linear regions;
- FMGP: design task-structured kernels (images, spatio-temporal) and pursue efficient approximations to Jacobian-induced kernels.

- **DVIP extensions**
  - richer function-space posteriors beyond Gaussian/GP surrogates (e.g., flows or diffusion-based variational families);
  - stronger base approximations (e.g., sparse implicit processes) and *structured priors* (attention/recurrent for sequences, convolutional for spatial data);
  - integration into autoencoders / generative architectures to unify representation learning and uncertainty.
- **Post-hoc uncertainty (VaLLA / FMGP)**
  - VaLLA: relax the Gaussian linearized-posterior assumption for better calibration in highly non-linear regions;
  - FMGP: design task-structured kernels (images, spatio-temporal) and pursue efficient approximations to Jacobian-induced kernels.
- **Theory directions**

- **DVIP extensions**
  - richer function-space posteriors beyond Gaussian/GP surrogates (e.g., flows or diffusion-based variational families);
  - stronger base approximations (e.g., sparse implicit processes) and *structured priors* (attention/recurrent for sequences, convolutional for spatial data);
  - integration into autoencoders / generative architectures to unify representation learning and uncertainty.
- **Post-hoc uncertainty (VaLLA / FMGP)**
  - VaLLA: relax the Gaussian linearized-posterior assumption for better calibration in highly non-linear regions;
  - FMGP: design task-structured kernels (images, spatio-temporal) and pursue efficient approximations to Jacobian-induced kernels.
- **Theory directions**
  - relax i.i.d. assumptions in large-deviation analyses;

- **DVIP extensions**

- richer function-space posteriors beyond Gaussian/GP surrogates (e.g., flows or diffusion-based variational families);
- stronger base approximations (e.g., sparse implicit processes) and *structured priors* (attention/recurrent for sequences, convolutional for spatial data);
- integration into autoencoders / generative architectures to unify representation learning and uncertainty.

- **Post-hoc uncertainty (VaLLA / FMGP)**

- VaLLA: relax the Gaussian linearized-posterior assumption for better calibration in highly non-linear regions;
- FMGP: design task-structured kernels (images, spatio-temporal) and pursue efficient approximations to Jacobian-induced kernels.

- **Theory directions**

- relax i.i.d. assumptions in large-deviation analyses;
- connect PAC–Chernoff/rate-function analysis with information-theoretic bounds and algorithmic stability;

- **DVIP extensions**

- richer function-space posteriors beyond Gaussian/GP surrogates (e.g., flows or diffusion-based variational families);
- stronger base approximations (e.g., sparse implicit processes) and *structured priors* (attention/recurrent for sequences, convolutional for spatial data);
- integration into autoencoders / generative architectures to unify representation learning and uncertainty.

- **Post-hoc uncertainty (VaLLA / FMGP)**

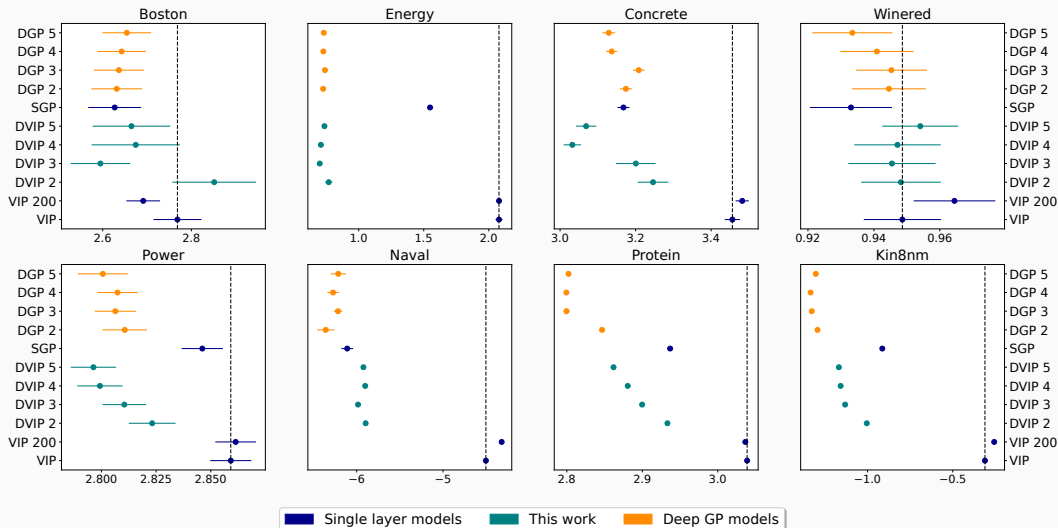
- VaLLA: relax the Gaussian linearized-posterior assumption for better calibration in highly non-linear regions;
- FMGP: design task-structured kernels (images, spatio-temporal) and pursue efficient approximations to Jacobian-induced kernels.

- **Theory directions**

- relax i.i.d. assumptions in large-deviation analyses;
- connect PAC–Chernoff/rate-function analysis with information-theoretic bounds and algorithmic stability;
- apply the rate-function formalism to transfer/continual learning and large language models.

**Questions?**

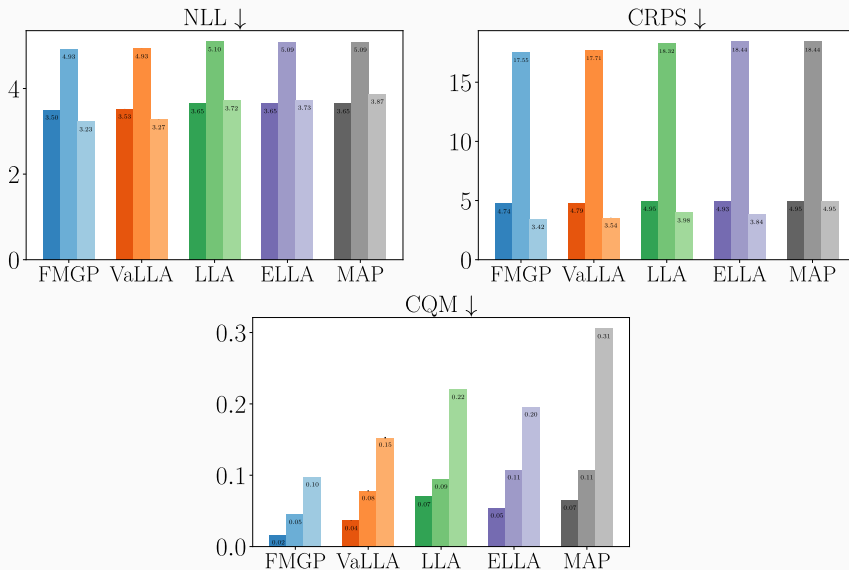
# DVIP Regression Results



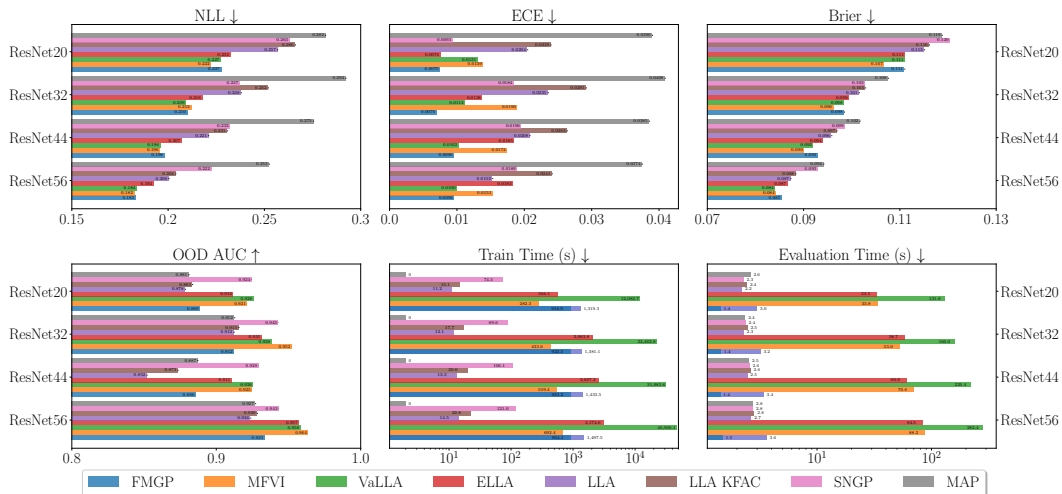
# VaLLA Results

Method	ResNet-20			ResNet-32			ResNet-44			ResNet-56			Rank
	ACC	NLL	ECE	ACC	NLL	ECE	ACC	NLL	ECE	ACC	NLL	ECE	
MAP	<b>92.6</b>	0.282	0.039	<b>93.5</b>	0.292	0.041	<b>94.0</b>	0.275	0.039	<b>94.4</b>	0.252	0.037	—
MF-VI	<b>92.7</b>	<b>0.231</b>	0.016	<b>93.5</b>	0.222	0.020	<b>93.9</b>	0.206	0.018	<b>94.4</b>	0.188	0.016	—
SNGP	92.4	0.266	0.024	93.2	0.256	0.025	93.8	0.242	0.028	93.8	0.229	0.022	—
GP - Subset	<b>92.6</b>	0.555	0.299	<b>93.4</b>	0.462	0.247	93.6	0.424	0.225	<b>94.4</b>	0.403	0.221	—
LLA Diag	92.2	0.728	0.404	92.7	0.755	0.430	92.8	0.778	0.445	<b>92.9</b>	0.843	0.480	—
LLA KFAC	92.0	0.852	0.467	91.8	1.027	0.547	91.4	1.091	0.566	89.8	1.174	0.579	—
LLA*	<b>92.6</b>	0.269	0.034	<b>93.5</b>	0.259	0.033	<b>94.0</b>	0.237	0.028	<b>94.4</b>	0.213	0.022	—
LLA* KFAC	<b>92.6</b>	0.271	0.035	<b>93.5</b>	0.260	0.033	<b>94.0</b>	0.232	0.028	<b>94.4</b>	0.202	0.024	—
ELLA	92.5	0.233	0.009	<b>93.5</b>	<b>0.215</b>	<b>0.008</b>	<b>93.9</b>	0.204	<b>0.007</b>	<b>94.4</b>	0.187	<b>0.007</b>	2.37
Sampled LLA	92.5	<b>0.231</b>	<b>0.006</b>	<b>93.5</b>	0.217	<b>0.008</b>	<b>94.0</b>	<b>0.200</b>	<b>0.007</b>	<b>94.4</b>	<b>0.185</b>	0.015	<b>2.00</b>
VaLLA	<b>92.6</b>	<b>0.228</b>	<b>0.007</b>	<b>93.5</b>	<b>0.211</b>	<b>0.007</b>	<b>94.0</b>	<b>0.198</b>	<b>0.008</b>	<b>94.4</b>	<b>0.183</b>	<b>0.009</b>	<b>1.37</b>

# FMGP Results



# FMGP Results



Model	Method	NLL	ECE	Train Time	Test Time
ResNet18	MAP	<b>1.247±0.000</b>	0.026±0.000	<b>0.0±0.0</b>	<b>5.058±0.029</b> ×10 <sup>2</sup>
	ELLA	1.248±0.000	<b>0.025±0.000</b>	<b>7.890±0.275</b> ×10 <sup>3</sup>	8.060±0.010×10 <sup>2</sup>
	FMGP	1.248±0.001	<b>0.015±0.001</b>	1.835±0.099×10 <sup>4</sup>	<b>7.324±0.001</b> ×10 <sup>2</sup>
	MFVI	<b>1.242±0.001</b>	0.040±0.000	7.602±0.032×10 <sup>4</sup>	3.773±0.308×10 <sup>4</sup>
ResNet34	MAP	<b>1.081±0.000</b>	0.035±0.000	<b>0.0±0.0</b>	<b>5.088±0.004</b> ×10 <sup>2</sup>
	ELLA	1.082±0.000	<b>0.034±0.000</b>	<b>1.201±0.373</b> ×10 <sup>4</sup>	1.087±0.018×10 <sup>3</sup>
	FMGP	<b>1.077±0.000</b>	<b>0.016±0.000</b>	1.942±0.103×10 <sup>4</sup>	<b>8.563±0.011</b> ×10 <sup>2</sup>
ResNet50	MAP	<b>0.962±0.000</b>	0.037±0.000	<b>0.0±0.0</b>	<b>4.954±0.010</b> ×10 <sup>2</sup>
	ELLA	<b>0.962±0.000</b>	<b>0.036±0.000</b>	2.997±1.215×10 <sup>4</sup>	1.954±0.018×10 <sup>3</sup>
	FMGP	<b>0.958±0.001</b>	<b>0.018±0.001</b>	<b>2.543±0.046</b> ×10 <sup>4</sup>	<b>1.100±0.010</b> ×10 <sup>3</sup>
ResNet101	MAP	<b>0.912±0.000</b>	0.049±0.000	<b>0.0±0.0</b>	<b>5.059±0.001</b> ×10 <sup>2</sup>
	ELLA	0.913±0.000	<b>0.048±0.000</b>	4.464±1.649×10 <sup>4</sup>	2.808±0.001×10 <sup>3</sup>
	FMGP	<b>0.900±0.000</b>	<b>0.030±0.001</b>	<b>2.654±0.064</b> ×10 <sup>4</sup>	<b>1.134±0.001</b> ×10 <sup>3</sup>
ResNet152	MAP	<b>0.876±0.000</b>	0.050±0.000	<b>0.0±0.0</b>	<b>6.324±0.004</b> ×10 <sup>2</sup>
	ELLA	0.877±0.000	<b>0.048±0.000</b>	6.820±0.526×10 <sup>4</sup>	3.877±0.007×10 <sup>3</sup>
	FMGP	<b>0.865±0.001</b>	<b>0.024±0.001</b>	<b>2.973±0.069</b> ×10 <sup>4</sup>	<b>1.267±0.002</b> ×10 <sup>3</sup>