# PAC-Bayesian Hypernetworks

Pascal Germain

www.pascalgermain.info

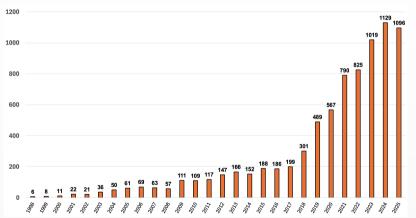
Département d'informatique et de génie logiciel

Université Laval

October 28, 2025

# PAC-Bayesian Theory

Pioneered by Shawe-Taylor and Williamson (1997), McAllester (1999), and Catoni (2003).



Number of search results per year for "PAC-Bayes(ian)" keywords on Google Scholar.

### This talk: The Mechanization of PAC-Bayes

Computer scientists [...] must create abstractions of real-world problems that can be understood by computer users and, at the same time, that can be represented and manipulated inside a computer.

— Foundations of Computer Science (Aho and Ullman 1992).

Chapter 1. Computer Science: The Mechanization of Abstraction

- Focus on PAC-Bayes for machine learning algorithm design;
- Propose to craft neural network architectures inspired from learning theory principles;
- A (perfectible) step towards guarantees for nowadays large/multimodal/foundation "models".

### Plan

- PAC-Bayesian Learning
- Meta-Learning Framework and PAC-Bayes Hypernetworks
- 3 Sample Compress Theory and Sample Compress Hypernetworks
- 4 PAC-Bayes Sample Compress Hypernetworks
- 5 Perspectives

#### **Definitions**

A learning example  $z := (x, y) \in \mathcal{X} \times \mathcal{Y} = \mathcal{Z}$  is a description-label pair.

### Data generating distribution

Each example is an **observation from distribution** D on  $\mathcal{Z}$ .

#### Learning sample

$$S := \{ z_1, z_2, \ldots, z_n \} \sim D^n$$

### Predictors (or hypothesis)

$$h: \mathcal{X} \to \mathcal{Y}, \quad h \in \mathcal{H}$$

### Learning algorithm

$$A(S) \longrightarrow h$$

#### Loss function

$$\ell: \mathcal{H} \times \mathcal{Z} \to \mathbb{R}$$

#### Empirical loss

$$\widehat{\mathcal{L}}_{\mathcal{S}}(h) = \frac{1}{n} \sum_{i=1}^{n} \ell(h, z_i)$$

#### Generalization loss

$$\mathcal{L}_D(h) = \mathop{\mathbf{E}}_{z \sim D} \ell(h, z)$$

## A Classical PAC-Bayesian Theorem

#### PAC-Bayesian theorem

(adapted from McAllester 1999; McAllester 2003)

For any distribution P on  $\mathcal{H}$ , for any  $\delta \in (0,1]$ , we have,

$$\Pr_{S \sim D^n} \left( \forall Q \text{ on } \mathcal{H} : \underbrace{\mathsf{E} \mathcal{L}_D(h)}_{\substack{h \sim Q \\ \text{loss}}} \leq \underbrace{\underbrace{\mathsf{E} \widehat{\mathcal{L}}_S(h)}_{\substack{h \sim Q \\ \text{loss}}} + \underbrace{\sqrt{\frac{1}{2n} \left[ \text{KL}(Q \| P) + \ln \frac{2\sqrt{n}}{\delta} \right]}}_{\text{complexity term}} \right) \geq 1 - \delta,$$

where 
$$\mathrm{KL}(Q||P) = \mathop{\mathbf{E}}_{f \sim Q} \ln \frac{Q(f)}{P(f)}$$
 is the Kullback-Leibler divergence.

### Valid for all posterior Q on ${\cal H}$

Inspiration for conceiving new learning algorithms.

# Tighter bounds for the [0,1]-loss (Classical PAC-Bayes theorems)

$$\operatorname{kl}\left(\underbrace{\frac{\mathbf{E}\widehat{\mathcal{L}}_{S}(h)}_{h\sim Q}, \underbrace{\frac{\mathbf{E}\mathcal{L}_{D}(h)}{h\sim Q}}_{h\sim Q}(h)\right) \leq \frac{1}{n}\left[\operatorname{KL}(Q\|P) + \ln\frac{2\sqrt{n}}{\delta}\right]_{\operatorname{kl}(q,\,p) := \,q\ln\frac{q}{p} + (1-q)\ln\frac{1-q}{1-p}}$$

From an algorithm design perspective, the "kl bound" promotes the minimization of

$$\mathrm{kl}^{-1}\left(\underset{h\sim Q}{\mathsf{E}\widehat{\mathcal{L}}_{\mathcal{S}}(h)}, \tfrac{1}{n}\left[\mathrm{KL}(Q\|P) + \ln\frac{2\sqrt{n}}{\delta}\right]\right) \coloneqq \underset{0\leq p\leq 1}{\operatorname{argsup}}\left\{\mathrm{kl}\left(\underset{h\sim Q}{\mathsf{E}\widehat{\mathcal{L}}_{\mathcal{S}}(h)}, p\right) \leq \frac{1}{n}\left[\mathrm{KL}(Q\|P) + \ln\frac{2\sqrt{n}}{\delta}\right]\right\}$$

#### The function $kl^{-1}$ is differentiable (see Reeb et al. 2018)

pyTorch implementation (Viallard et al. 2021):

https://github.com/paulviallard/ECML21-PB-CBound/blob/master/core/kl\_inv.py

#### Lemma (see Letarte et al. 2019)

$$\mathrm{kl}^{-1}\left(\underset{h\sim Q}{\overset{\mathbf{E}\widehat{\mathcal{L}}_{S}}(h)}, \tfrac{1}{n}\left[\mathrm{KL}(Q\|P) + \ln\frac{2\sqrt{n}}{\delta}\right]\right) = \inf_{c>0}\left\{\frac{1}{1-e^{-c}}\left(c \cdot \underset{h\sim Q}{\overset{\mathbf{E}\widehat{\mathcal{L}}_{S}}{(h)}} + \tfrac{1}{n}\left[\mathrm{KL}(Q\|P) + \ln\frac{2\sqrt{n}}{\delta}\right]\right)\right\}$$

Pascal Germain (Université Laval)

# Distribution over parameters

### Given a model / predictor $h_{\theta}$ , where $\theta$ are parameters.

Consider P and Q as distributions over the set of parameters  $\Theta$ .

$$\forall Q \text{ on } \Theta: \quad \mathrm{kl}\Big(\underset{\theta \sim Q}{\overset{\textstyle \mathbf{E}\widehat{\mathcal{L}}}{\mathcal{S}}} \mathsf{S}\big(h_{\theta}\big), \underset{\theta \sim Q}{\overset{\textstyle \mathbf{E}}{\mathcal{L}}} \mathcal{L}_{D}\big(h_{\theta}\big)\Big) \, \leq \, \tfrac{1}{n} \, \Big[\mathrm{KL}\big(Q\|P\big) + \ln \tfrac{2\sqrt{n}}{\delta}\Big].$$

### Typical approach for (stochastics) neural networks

(Dziugaite and Roy 2017; Neyshabur, Bhojanapalli, and Srebro 2018; Nozawa, Germain, and Guedj 2020; Pérez-Ortiz et al. 2021, among many others.)

•  $P = \mathcal{N}(\mathbf{W}_p, \sigma_p \mathbf{I})$ 

where  $\mathbf{W}_{p}$  are the random/pre-learned weights initialization.

•  $Q = \mathcal{N}(\mathbf{W}, \sigma \mathbf{I}),$ 

where  $\boldsymbol{W}$  are the learned/fine-tuned neural network weights.

Then, 
$$KL(Q||P) = \frac{1}{2} ||W - W_p||^2$$
.

Pérez-Ortiz, Rivasplata, Shawe-Taylor and Szepesvári

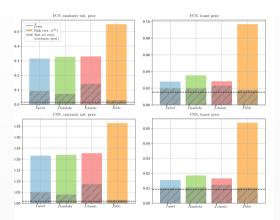


Figure 3: Tightness of the risk certificates for MNIST across different architectures, priors and training objectives. The bottom shaded areas correspond to the test set 0-1 error of the stochastic classifier. The coloured areas on top indicate the tightness of the risk certificate (smaller is better). The horizontal dashed line corresponds to the test set 0-1 error of f<sub>em.</sub> i.e. the deterministic classifier learnt

- Build on the pioneer work of Dziugaite and Roy (2017).
- Tight guarantees!

risk 
$$\leq 1.55\%$$
 on MNIST (CNN) with probability  $\geq 95\%$ .

Easy to train.

Source code (pyTorch):

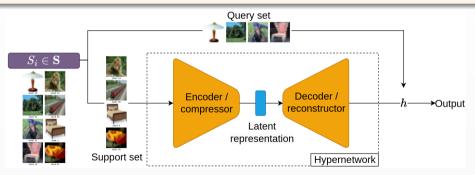
https://github.com/mperezortiz/PBB

### Plan

- PAC-Bayesian Learning
- Meta-Learning Framework and PAC-Bayes Hypernetworks
- 3 Sample Compress Theory and Sample Compress Hypernetworks
- 4 PAC-Bayes Sample Compress Hypernetworks
- 5 Perspectives

# Overview of our meta-learning framework

Benjamin Leblanc et al. (2025). "Generalization Bounds via Meta-Learned Model Representations: PAC-Bayes and Sample Compression Hypernetworks". In: ICML. PMLR



### **Definitions**

### Meta-Learning dataset

$$\mathcal{S} \coloneqq \{S_1, S_2, \dots, S_m\}$$
 such that  $S_i \sim \left(D_i\right)^{n_i}$ .

### Meta-Learning algorithm

$$\mathcal{A}(\mathcal{S}) \longrightarrow (\phi, \psi)$$

#### Encoder

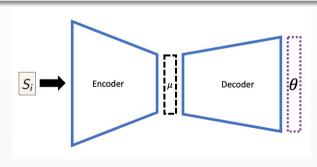
$$\mathcal{E}_{\phi}(S) \longrightarrow \boldsymbol{\mu}$$

#### Decoder

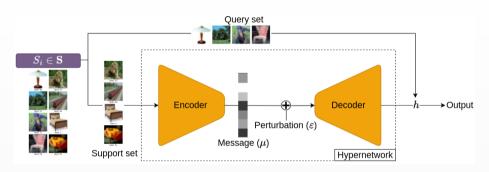
$$\mathfrak{D}_{\psi}(\boldsymbol{\mu}) \longrightarrow \theta$$

#### (Downstream) Predictor

$$h_{\theta}(\mathbf{x}) \longrightarrow \mathbf{y}$$



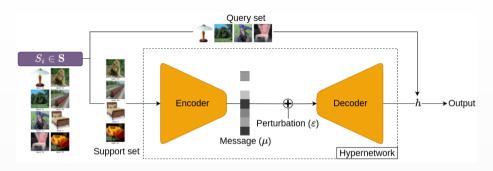
# PAC-Bayes Hypernetworks



**Training objective.** Each dataset  $S_i \in \mathcal{S}$  is split into a *support set*  $\hat{S}_i$  and a *query set*  $\hat{T}_i = S_i \setminus \hat{S}_i$ .

$$\min_{\psi,\phi} \left\{ \frac{1}{m} \sum_{i=1}^{m} \mathbf{E} \, \widehat{\mathcal{L}}_{\widehat{T}_i}(h_{\theta_i}) \, \middle| \, \theta_i = \mathcal{D}_{\psi}(\boldsymbol{\mu}_i + \boldsymbol{\epsilon}); \, \boldsymbol{\mu}_i = \mathcal{E}_{\phi}(\hat{S}_i) \right\}, \quad \text{ with } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

## PAC-Bayes Hypernetworks

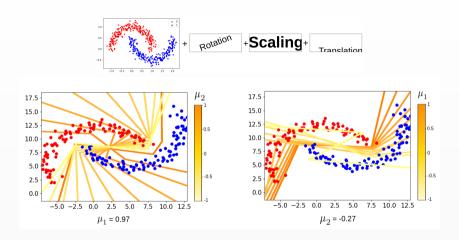


**Bound computation.** Given a new task sample  $S' \sim (\mathcal{D}')^{n'}$ , let  $\mu = \mathcal{E}_{\phi}(S')$  and  $\theta' = \mathcal{D}_{\psi}(\mu + \epsilon)$ .

$$\mathbf{E}\,\mathcal{L}_{D'}(h_{\theta'}) \leq \operatorname*{argsup}_{0 \leq \rho \leq 1} \left\{ \mathrm{kl}\Big(\,\mathbf{E}\,\widehat{\mathcal{L}}_{\mathcal{S}'}(h_{\theta'}), \rho\Big) \leq \frac{\frac{1}{2} \|\boldsymbol{\mu}\|^2 + \ln \frac{2\sqrt{n'}}{\delta}}{n'} \right\},\,$$

using a prior  $P_0 = \mathcal{N}(\mathbf{0}, \mathbf{I})$  and a posterior  $Q = \mathcal{N}(\mu, \mathbf{I})$  over the latent representation space.

## Toy Experiments



$$\theta = \mathfrak{D}_{\psi}([\mu_1, \mu_2]^T)$$

## MNIST pixel swap

Multiclass classification

Meta-train set: 10 tasks of 60 000 examples

Meta-test set: 20 tasks of 2000 examples









Image: Benjamin Leblanc & Claude

Algorithm	100 Pixels swap		200 Pixels swap		300 Pixels swap	
	Bound	Test error	Bound	Test error	Bound	Test error
Pentina and Lampert (2014)	0.190	0.019	0.240	0.026	0.334	0.038
Amit and Meir (2018)	0.138	0.016	0.161	0.020	0.329	0.040
Guan and Lu (2022)	0.093	0.015	0.128	0.019	0.210	0.024
Zakerinia, Behjati, and Lampert (2024)	0.053	0.019	0.108	0.026	0.149	0.035
Our PAC-Bayes Hypernetwork	0.068	0.027	0.112	0.076	0.219	0.186
Opaque encoder		0.037		0.087		0.159

### MINIST and CIFAR100 Binary Pairs

Binary classification

 Meta-train set: 56 tasks of 2000 examples for MNIST; 100 tasks of 1200 examples for CIFAR

 Meta-test set: 34 tasks of 2000 examples for MNIST; 50 tasks of 200 examples for CIFAR

Task #1: 3 VS 6 Task #2: 8 VS ()













Algorithm	MN	IIST	CIFAR100		
Algorithm	Bound	Test error	Bound	Test error	
Pentina and Lampert (2014)	$0.767 \pm 0.001$	$0.369 \pm 0.223$	$0.801 \pm 0.001$	$0.490 \pm 0.070$	
Amit and Meir (2018)	$1372 \pm 23.36$	$0.351\pm0.212$	$950.9 \pm 343.1$	$0.284\pm0.120$	
Guan and Lu (2022)	$0.754 \pm 0.003$	$0.366\pm0.221$	$0.802 \pm 0.001$	$0.489 \pm 0.073$	
Zakerinia, Behjati, and Lampert (2024)	$0.684 \pm 0.021$	$0.351 \pm 0.212$	$0.953 \pm 0.315$	$0.281 \pm 0.125$	
Our PAC-Bayes Hypernetwork	$0.597 \pm 0.107$	$0.150 \pm 0.114$	$0.974 \pm 0.022$	$0.295\pm0.103$	
Opaque encoder		$0.497 \pm 0.134$		$0.506 \pm 0.101$	

### Plan

- PAC-Bayesian Learning
- 2 Meta-Learning Framework and PAC-Bayes Hypernetworks
- 3 Sample Compress Theory and Sample Compress Hypernetworks
- 4 PAC-Bayes Sample Compress Hypernetworks
- 5 Perspectives

# History

Inception: Littlestone and Warmuth (1986): "Relating Data Compression and Learnability"

### The Set Covering Machine

Marchand and Shawe-Taylor (2002): "The Set Covering Machine"

Marchand and Sokolova (2005): "Learning with Decision Lists of Data-Dependent Features"

Laviolette, Marchand, and Shah (2005): "Margin-Sparsity Trade-Off for the Set Covering Machine"

Hussain et al. (2007): "Revised Loss Bounds for the Set Covering Machine and Sample-Compression Loss Bounds for Imbalanced Data"

Drouin et al. (2019): "Interpretable genotype-to-phenotype classifiers with performance guarantees"

#### Pick-To-Learn

Campi and Garatti (2023): "Compression, Generalization and Learning"

Paccagnan, Campi, and Garatti (2023): "The Pick-to-Learn Algorithm: Empowering Compression for Tight Generalization Bounds and Improved Post-training Performance"

Marks and Paccagnan (2025): "Pick-to-Learn and Self-Certified Gaussian Process Approximations"

Bazinet, Zantedeschi, and Germain (2025): "Sample Compression Unleashed: New Generalization Bounds for Real Valued Losses"

#### **Definitions**

A **compressed predictor**  $h_{\mathbf{i}}^{\mu}$  is a data-dependant predictor encoded by two quantities

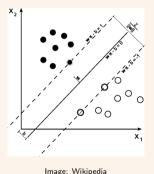
- A **compression set**  $S_i$  is a subset of the training set S:
  - $\mathbf{i} \in \mathcal{I}_n := \mathcal{P}(\{1, 2, \dots, n\})$
- A **message**  $\mu \in \mathcal{M}_{\mathbf{i}}$  contains additional information to describe the predictor  $h_{\mathbf{i}}^{\mu}$ .
  - The message  $\mu$  is chosen among a (discrete) set  $\mathcal{M}_i$  of predefined messages given  $S_i$ .
  - For simplicity, we sometime use a message set  $\mathcal{M}$  that does not rely on  $S_i$ .

Given  $S_i \in \mathcal{Z}^{|i|}$  et  $\mu \in \mathcal{M}_i$ , a reconstruction function  $\mathcal{R}$  outputs a predictor :

$$h_{\mathbf{i}}^{\mu} = \mathcal{R}(S_{\mathbf{i}}, \mu).$$

# A Classical Sample Compressed Classifiers

### SVM: Support Vector Machine (hard margin)



The SVM's learning algorithm acts as its own reconstruction function

$$SVM(S) = h_{\mathbf{i}}^{\mu} = SVM(S_{\mathbf{i}})$$

with  $S_i = \{\text{support vectors}\}\$  and  $\mu = \emptyset$ 

# Generalization Bounds for sample compressed binary classifiers

Given  $h^\mu_{\mathbf{i}}: \mathcal{X} \to \{-1, +1\}$ , and  $\mathcal{L}^{\text{Ol}}_D(h^\mu_{\mathbf{i}}) = \mathop{\mathbf{E}}_{(x,y) \sim D} I(h^\mu_{\mathbf{i}}(x) \neq y)$  the zero-one loss.

### Theorem (Marchand and Sokolova 2005; Laviolette, Marchand, and Shah 2005)

Let  $\mathcal R$  be a reconstruction function,  $P_{\mathcal M_i}$  a distribution over messages, and  $\delta \in (0,1]$ . With high probability  $(\geq 1-\delta)$  over  $S \sim D^n$ , we have

 $\forall \mathbf{i} \in \mathcal{I}_n, \mu \in \mathcal{M}_{\mathbf{i}}$ :

$$\mathcal{L}_{D}^{01}(\textit{h}_{\textbf{i}}^{\mu}) \leq 1 - \exp \left( \frac{-1}{\textit{n} - |\textbf{i}| - \textit{k}_{\textit{S}_{\textbf{i}^{c}}}} \left[ \ln \binom{\textit{n} - |\textbf{i}|}{\textit{k}_{\textit{S}_{\textbf{i}^{c}}}} \right) + \ln \binom{\textit{n}}{|\textbf{i}|} + \ln \left( \frac{1}{\textit{P}_{\mathcal{M}_{\textbf{i}}}(\mu) \cdot \xi(|\textbf{i}|) \cdot \delta} \right) \right] \right)$$

where  $k_{S_{\mathbf{i}^c}} := |\mathbf{i}^c| \widehat{\mathcal{L}}_{S_{\mathbf{i}^c}}^{01}(h_{\mathbf{i}}^{\mu})$  is the error count on  $S_{\mathbf{i}^c} := S \setminus S_{\mathbf{i}}$  and  $\xi(a) := \frac{6}{\pi^2}(a+1)^{-2}$ .

### New bound for real-valued losses

### Assumption

The loss  $\ell: \mathcal{H} \times \mathcal{Z} \to \mathbb{R}$  is  $\sigma$ -sub-Gaussian, *i.e.*, for all  $\mathbf{i} \in \mathcal{I}_n, \mu \in \mathcal{M}_{\mathbf{i}}$ :

$$\sum_{z\sim D} e^{\lambda(\ell(h_{\mathbf{i}}^{\mu},z)-\mathcal{L}_{D}(h_{\mathbf{i}}^{\mu}))} \leq e^{rac{\lambda^{2}\sigma^{2}}{2}}, \quad orall \lambda \in \mathbb{R}.$$

#### Theorem

Let  $\mathcal{R}$  be a reconstruction function,  $P_{\mathcal{M}_i}$  a distribution over messages, and  $\delta \in (0,1]$ . With high probability  $(\geq 1-\delta)$  over  $S \sim D^n$ , we have

 $\forall \mathbf{i} \in \mathcal{I}_n, \mu \in \mathcal{M}_{\mathbf{i}}$ :

$$\mathcal{L}_D(\textit{h}_{\textbf{i}}^{\mu}) \leq \frac{\widehat{\mathcal{L}}_{\textit{S}_{\textbf{i}c}}(\textit{h}_{\textbf{i}}^{\mu})}{\sqrt{n-|\textbf{i}|}} \left[ \frac{\sigma^2}{2} + \ln \binom{n}{|\textbf{i}|} + \ln \left( \frac{1}{P_{\mathcal{M}_{\textbf{i}}}(\mu) \cdot \xi(|\textbf{i}|) \cdot \delta} \right) \right].$$

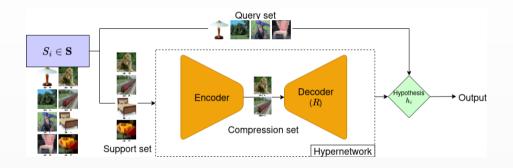
**Proof idea.** Each sample compressed classifier  $h_{\mathbf{i}}^{\mu}$  is defined independently of  $S_{\mathbf{i}^c} \sim D^{n-|\mathbf{i}|}$ .

For all  $\mathbf{i} \in \mathcal{I}_n, \mu \in \mathcal{M}_{\mathbf{i}}$  and  $\delta_{\mathbf{i},\mu} \in (0,1]$ , with probability at least  $1 - \delta_{\mathbf{i},\mu}$ :

$$\begin{split} \mathcal{L}_D(h_{\mathbf{i}}^\mu) - \widehat{\mathcal{L}}_{S_{\mathbf{i}^c}}(h_{\mathbf{i}}^\mu) & \leq & \frac{1}{t} \left[ \ln \left( \sum_{S_{\mathbf{i}^c} \sim D^{n-|\mathbf{i}|}} e^{t(\mathcal{L}_D(h_{\mathbf{i}}^\mu) - \widehat{\mathcal{L}}_{S_{\mathbf{i}^c}}(h_{\mathbf{i}}^\mu))} \right) + \ln \frac{1}{\delta_{\mathbf{i},\mu}} \right] \ \, \langle \ \, \mathsf{Chernoff} \ \, (t > 0) \ \, \rangle \\ & = & \frac{1}{t} \left[ \ln \left( \prod_{i=1}^{n-|\mathbf{i}|} \mathop{\mathbf{E}}_{z \sim D} e^{\frac{t}{n-|\mathbf{i}|} (\ell(h_{\mathbf{i}}^\mu, z) - \widehat{\mathcal{L}}_{S_{\mathbf{i}^c}}(h_{\mathbf{i}}^\mu))} \right) + \ln \frac{1}{\delta_{\mathbf{i},\mu}} \right] \\ & \leq & \frac{1}{t} \left[ \ln \left( \prod_{i=1}^{n-|\mathbf{i}|} e^{\frac{t^2 \sigma^2}{2(n-|\mathbf{i}|)^2}} \right) + \ln \frac{1}{\delta_{\mathbf{i},\mu}} \right] \ \, \langle \ \, \mathsf{sub\text{-}Gaussian \ \, loss, \ \, with \ \, \lambda \coloneqq \frac{t}{n-|\mathbf{i}|} \ \, \rangle \\ & = & \frac{1}{t} \left[ \frac{t^2 \sigma^2}{2(n-|\mathbf{i}|)} + \ln \frac{1}{\delta_{\mathbf{i},\mu}} \right]. \end{split}$$

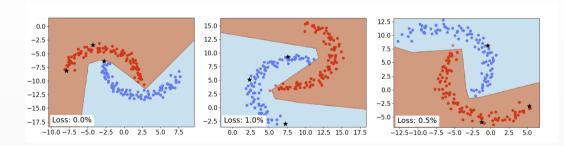
The final result is obtained by choosing  $t := \sqrt{n - |\mathbf{i}|}$  and by an <u>union bound</u> over the (discrete) set of all possible sample compress classifiers, with  $\delta_{\mathbf{i},\mu} = \frac{\xi(|\mathbf{i}|)}{\binom{n}{i}} \cdot P_{\mathcal{M}_{\mathbf{i}}}(\mu) \cdot \delta$ , since  $\sum_{i=1}^{\infty} \xi(i) = 1$ .

# Sample Compress Hypernetworks



# Toy Experiments





### Plan

- PAC-Bayesian Learning
- 2 Meta-Learning Framework and PAC-Bayes Hypernetworks
- Sample Compress Theory and Sample Compress Hypernetworks
- 4 PAC-Bayes Sample Compress Hypernetworks
- 5 Perspectives

# New bound for continuous messages

#### PAC-Bayes to the rescue!

We now consider a posterior distribution  $Q_{\mathcal{M}}$  over the (possibly continuous) message set  $\mathcal{M}$ .

#### Theorem

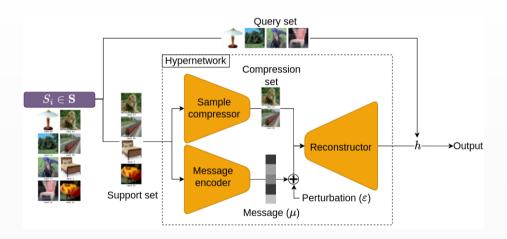
Let  $\mathcal{R}$  be a reconstruction function,  $P_{\mathcal{M}}$  a distribution over messages, and  $\delta \in (0,1]$ . With high probability  $(\geq 1 - \delta)$  over  $S \sim D^n$ , we have

 $\forall \mathbf{i} \in \mathcal{I}_n, Q_{\mathcal{M}} \text{ over } \mathcal{M} :$ 

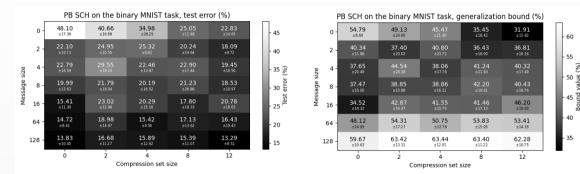
$$\underset{\mu \sim Q_{\mathcal{M}}}{\text{\textbf{E}}} \mathcal{L}_{D}(h_{\mathbf{i}}^{\mu}) \leq \underset{\mu \sim Q_{\mathcal{M}}}{\text{\textbf{E}}} \hat{\mathcal{L}}_{S_{\mathbf{i}^{c}}}(h_{\mathbf{i}}^{\mu}) + \frac{1}{\sqrt{n - |\mathbf{i}|}} \left[ \text{KL}(Q_{\mathcal{M}} \| P_{\mathcal{M}}) + \frac{\sigma^{2}}{2} + \ln \binom{n}{|\mathbf{i}|} + \ln \left( \frac{1}{\xi(|\mathbf{i}|) \cdot \delta} \right) \right].$$

**Proof idea.** For a fixed  $\mathbf{i} \in \mathcal{I}_n$ , get a typical PAC-Bayes bound with prior/posterior distribution over messages  $\mathcal{M}$ . Then, use a <u>union bound</u> to get a bound uniformly valid over the compression sets  $\mathcal{I}_n$ .

# PAC-Bayes Sample Compress Hypernetworks



# **MNIST Binary Pairs**



### Plan

- PAC-Bayesian Learning
- Meta-Learning Framework and PAC-Bayes Hypernetworks
- 3 Sample Compress Theory and Sample Compress Hypernetworks
- 4 PAC-Bayes Sample Compress Hypernetworks
- 5 Perspectives

### Perspectives

- Dynamically select the message/compression set size during learning;
- Fine-tune large language models.

### Collaborators

Benjamin Leblanc



Nathaniel D'Amours



Mathieu Bazinet



Alexandre Drouin



- Aho, Alfred V. and Jeffrey D. Ullman (1992). Foundations of Computer Science, C Edition. Computer Science Press / W. H. Freeman.
- Amit, Ron and Ron Meir (2018). "Meta-Learning by Adjusting Priors Based on Extended PAC-Bayes Theory". In: ICML.
  - Bazinet, Mathieu, Valentina Zantedeschi, and Pascal Germain (2025). "Sample Compression Unleashed: New Generalization Bounds for Real Valued Losses". In: *AISTATS*. Vol. 258. Proceedings of Machine Learning Research. PMLR, pp. 3286–3294.
- Campi, Marco C. and Simone Garatti (2023). "Compression, Generalization and Learning". In: *JMLR* abs/2301.12767.
- Catoni, Olivier (2003). "A PAC-Bayesian approach to adaptive classification". In: preprint 840.
- Drouin, Alexandre et al. (2019). "Interpretable genotype-to-phenotype classifiers with performance guarantees". In: *Scientific reports* 9.1, pp. 1–13.
- Dziugaite, Gintare Karolina and Daniel M. Roy (2017). "Computing Nonvacuous Generalization Bounds for Deep (Stochastic) Neural Networks with Many More Parameters than Training Data". In: *UAI*. AUAI Press.
- Guan, Jiechao and Zhiwu Lu (2022). "Fast-Rate PAC-Bayesian Generalization Bounds for Meta-Learning". In: ICML.

- Hussain, Zakria et al. (2007). "Revised Loss Bounds for the Set Covering Machine and Sample-Compression Loss Bounds for Imbalanced Data". In: *J. Mach. Learn. Res.* 8, pp. 2533–2549.
- Laviolette, François, Mario Marchand, and Mohak Shah (2005). "Margin-Sparsity Trade-Off for the Set Covering Machine". In: *ECML*. Vol. 3720. Lecture Notes in Computer Science. Springer, pp. 206–217.
- Leblanc, Benjamin et al. (2025). "Generalization Bounds via Meta-Learned Model Representations: PAC-Bayes and Sample Compression Hypernetworks". In: ICML. PMLR.
- Letarte, Gaël et al. (2019). "Dichotomize and Generalize: PAC-Bayesian Binary Activated Deep Neural Networks". In: *NeurIPS*, pp. 6869–6879.
- Littlestone, Nick and Manfred K. Warmuth (1986). "Relating Data Compression and Learnability". In: *Technical Report*.
- Marchand, Mario and John Shawe-Taylor (2002). "The Set Covering Machine". In: JMLR 3.
- Marchand, Mario and Marina Sokolova (2005). "Learning with Decision Lists of Data-Dependent Features". In: *JMLR* 6.

- Marks, Daniel and Dario Paccagnan (2025). "Pick-to-Learn and Self-Certified Gaussian Process Approximations". In: *AISTATS*. Vol. 258. Proceedings of Machine Learning Research. PMLR, pp. 2656–2664.
- McAllester, David (1999). "Some PAC-Bayesian Theorems". In: Machine Learning 37.3.
- - Neyshabur, Behnam, Srinadh Bhojanapalli, and Nathan Srebro (2018). "A PAC-Bayesian Approach to Spectrally-Normalized Margin Bounds for Neural Networks". In: *ICLR* (*Poster*). OpenReview.net.
- Nozawa, Kento, Pascal Germain, and Benjamin Guedj (2020). "PAC-Bayesian Contrastive Unsupervised Representation Learning". In: *UAI*. Vol. 124. Proceedings of Machine Learning Research. AUAI Press, pp. 21–30.
- Paccagnan, Dario, Marco C. Campi, and Simone Garatti (2023). "The Pick-to-Learn Algorithm: Empowering Compression for Tight Generalization Bounds and Improved Post-training Performance". In: NeurIPS.
- Pentina, Anastasia and Christoph H. Lampert (2014). "A PAC-Bayesian bound for Lifelong Learning". In: *ICML*.

- Pérez-Ortiz, María et al. (2021). "Tighter Risk Certificates for Neural Networks". In: *J. Mach. Learn. Res.* 22, 227:1–227:40.
- Reeb, David et al. (2018). "Learning Gaussian Processes by Minimizing PAC-Bayesian Generalization Bounds". In: *NeurIPS*, pp. 3341–3351.
- Shawe-Taylor, John and Robert C. Williamson (1997). "A PAC Analysis of a Bayesian Estimator". In: COLT.
- Viallard, Paul et al. (2021). "Self-bounding Majority Vote Learning Algorithms by the Direct Minimization of a Tight PAC-Bayesian C-Bound". In: ECML/PKDD (2). Vol. 12976. Lecture Notes in Computer Science. Springer, pp. 167–183.
- Zakerinia, Hossein, Amin Behjati, and Christoph H. Lampert (2024). "More flexible PAC-Bayesian meta-learning by learning learning algorithms". In: ICML.