

Benjamin Leblanc<sup>1</sup>, Mathieu Bazinet<sup>1</sup>, Nathaniel D'Amours<sup>1</sup>, Alexandre Drouin<sup>1,2</sup>, Pascal Germain<sup>1</sup>

<sup>1</sup>Département d'informatique et de génie logiciel, Université Laval, Québec, Canada

<sup>2</sup>ServiceNow Research, Montréal, Canada

## Contributions

- 1 We generalize sample compression bounds to continuous messages;
- 2 We present an original approach to sample compression by learning the reconstruction function, seeing it as a hypernetwork;
- 3 We propose an algorithm to meta-learn the generation of predictors with sample compression generalization guarantees.

## Sample compression

### The setting

- A data-generating distribution  $\mathcal{D}$  over an instance space  $\mathcal{X} \times \mathcal{Y}$ ;
- A dataset  $S = \{(\mathbf{x}_j, y_j)\}_{j=1}^m \sim \mathcal{D}^m$ ;
- A predictor  $h : \mathcal{X} \rightarrow \mathcal{Y}$  and a learning algorithm  $A(S) \mapsto h$ ;
- The empirical loss  $\hat{\mathcal{L}}_S(h) = \frac{1}{m} \sum_{j=1}^m \ell(h(\mathbf{x}_j), y_j)$ , with  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ ;
- The generalization loss  $\mathcal{L}_{\mathcal{D}}(h) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(h(\mathbf{x}), y)]$ .

### The reconstruction function

In the sample compression framework, a learned predictor  $A(S)$  can be fully defined by a **reconstruction function**  $\mathcal{R}$ , a **compression set**  $S_j$  and a **message**  $\sigma$ , such that  $\mathcal{R}(S_j, \sigma) = A(S)$ .

#### Compression set

$S_j \subseteq S$ , with train indexes  $\mathbf{j} \in \mathcal{P}(m)$ , chosen from the power set of  $\mathbf{m} = \{i\}_{i=1}^m$

#### Message

$\sigma \in \Sigma$ , where  $\Sigma$  is the set of all possible messages.

### First sample compression bound for uncountable sets $\Sigma$

The theorem below bounds the generalization loss  $\mathcal{L}_{\mathcal{D}}(h)$  from the empirical loss  $\hat{\mathcal{L}}_S(h)$  and two complexity terms: the compression set size  $|\mathbf{j}|$  and the KL divergence between a prior  $P_{\Sigma}$  and a posterior  $Q_{\Sigma}$  distributions over messages.

- 1 **Theorem 1** For any distribution  $\mathcal{D}$  over  $\mathcal{X} \times \mathcal{Y}$ , for any set  $J \subseteq \mathcal{P}(m)$ , for any distribution  $P_J$  over  $J$ , for any distribution  $P_{\Sigma}$  over  $\Sigma$ , for any reconstruction function  $\mathcal{R}$ , for any loss  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow [0, 1]$ , for any convex function  $\Delta : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$  and for any  $\delta \in (0, 1]$ , with probability at least  $1 - \delta$  over the draw of  $S \sim \mathcal{D}^m$ , we have that for all  $\mathbf{j} \in J$  and  $Q_{\Sigma}$  over  $\Sigma$ :

$$\Delta \left( \mathbb{E}_{\sigma \sim Q_{\Sigma}} \hat{\mathcal{L}}_{S_j}(\mathcal{R}(S_j, \sigma)), \mathbb{E}_{\sigma \sim Q_{\Sigma}} \mathcal{L}_{\mathcal{D}}(\mathcal{R}(S_j, \sigma)) \right) \leq \frac{1}{m - |\mathbf{j}|} \left[ \text{KL}(Q_{\Sigma} \| P_{\Sigma}) + \ln \left( \frac{\mathcal{J}_{\Delta}(m - |\mathbf{j}|)}{P_J(\mathbf{j}) \cdot \delta} \right) \right],$$

with

$$\mathcal{J}_{\Delta}(m - |\mathbf{j}|) = \mathbb{E}_{\sigma \sim P_{\Sigma}} \mathbb{E}_{T_j \sim \mathcal{D}^{|\mathbf{j}|}} \mathbb{E}_{T_j \sim \mathcal{D}^{m-|\mathbf{j}|}} e^{(m-|\mathbf{j}|) \cdot \Delta(\hat{\mathcal{L}}_{T_j}(\mathcal{R}(T_j, \sigma)), \mathcal{L}_{\mathcal{D}}(\mathcal{R}(T_j, \sigma)))}$$

## Sample compression hypernetworks

- 2 We propose learning the reconstruction function.

Our **reconstruction hypernetwork**  $\gamma = \mathcal{R}_{\theta}(S_j, \sigma)$  takes two inputs:

- A compression set  $S_j$  containing a fixed number  $c$  examples;
- A message  $\sigma$  taking the form of a vector of fixed size  $b$ , either real-valued ( $\sigma \in [-1, 1]^b$ ), or discrete ( $\sigma \in \{-1, 1\}^b$ ).

The output  $\gamma \in \mathbb{R}^{|\mathcal{X}|}$  is the parameters of a **downstream network**:

$$h_{\gamma} : \mathcal{X} \rightarrow \mathcal{Y}.$$

**Objective function:** Minimize the empirical loss of the downstream predictor  $h_{\gamma}$  on the complement set  $S \setminus S_j$ :

$$\min_{\theta} \left\{ \frac{1}{m - |\mathbf{j}|} \sum_{(\mathbf{x}, y) \in S \setminus S_j} \ell(h_{\gamma}(\mathbf{x}), y) \mid \gamma = \mathcal{R}_{\theta}(S_j, \sigma) \right\}.$$

## From sample compression to meta-learning

### The setting

- A meta-distribution  $\mathbf{D}$ , such that  $\mathcal{D}_i \sim \mathbf{D}$ ;
- A meta-dataset  $\mathbf{S} = \{S_i\}_{i=1}^n$ , where  $S_i \sim \mathcal{D}_i^m$ ;
- Each  $S_i$  is split into *support set*  $\hat{S}_i \subset S_i$  and *query set*  $\hat{T}_i = S_i \setminus \hat{S}_i$ .

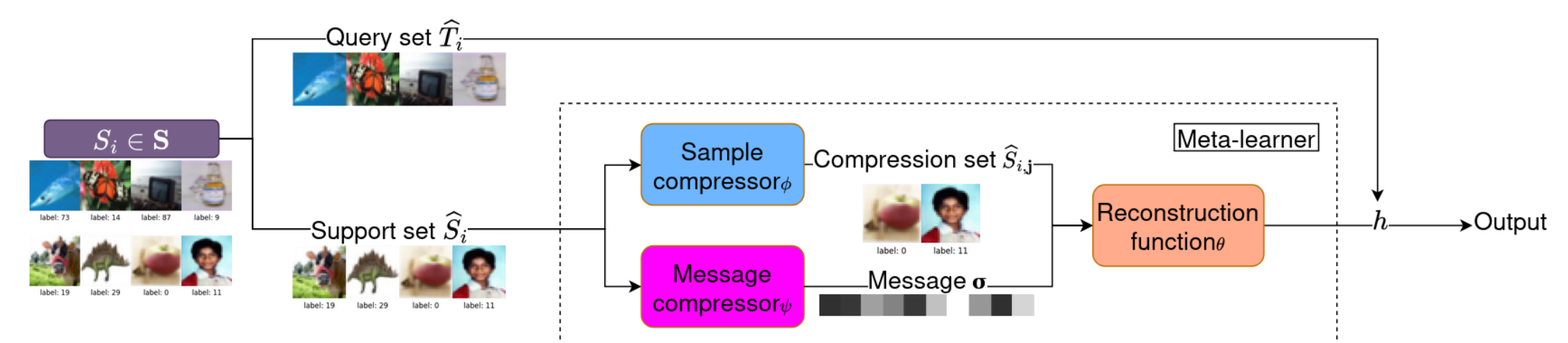
### The proposed meta-learning algorithm

- 3 Let  $\mathcal{C}_{\phi}$  be a *sample compressor* and  $\mathcal{M}_{\psi}$  a *message compressor*. Given a query set  $\hat{S}$ , these respectively yield the compression set and the message used by the reconstruction hypernetwork to obtain the parameters  $\gamma$  of a downstream predictor:

$$\gamma = \mathcal{R}_{\theta}(\mathcal{C}_{\phi}(\hat{S}), \mathcal{M}_{\psi}(\hat{S})).$$

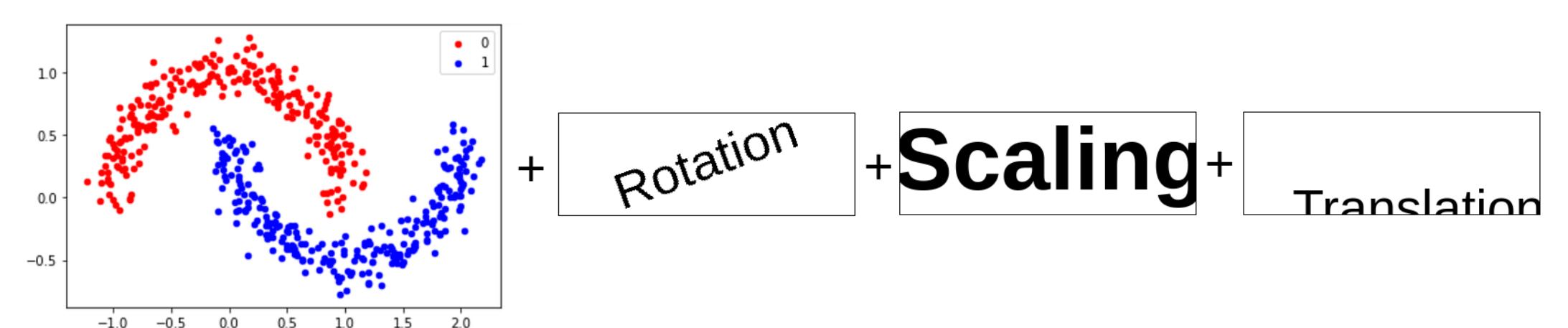
We propose to optimize the following meta-learning objective:

$$\min_{\theta, \phi, \psi} \left\{ \frac{1}{n} \sum_{i=1}^n \frac{1}{m_i - |\hat{S}_i|} \sum_{(\mathbf{x}, y) \in \hat{T}_i} \ell(h_{\gamma_i}(\mathbf{x}), y) \mid \gamma_i = \mathcal{R}_{\theta}(\mathcal{C}_{\phi}(\hat{S}_i), \mathcal{M}_{\psi}(\hat{S}_i)) \right\},$$

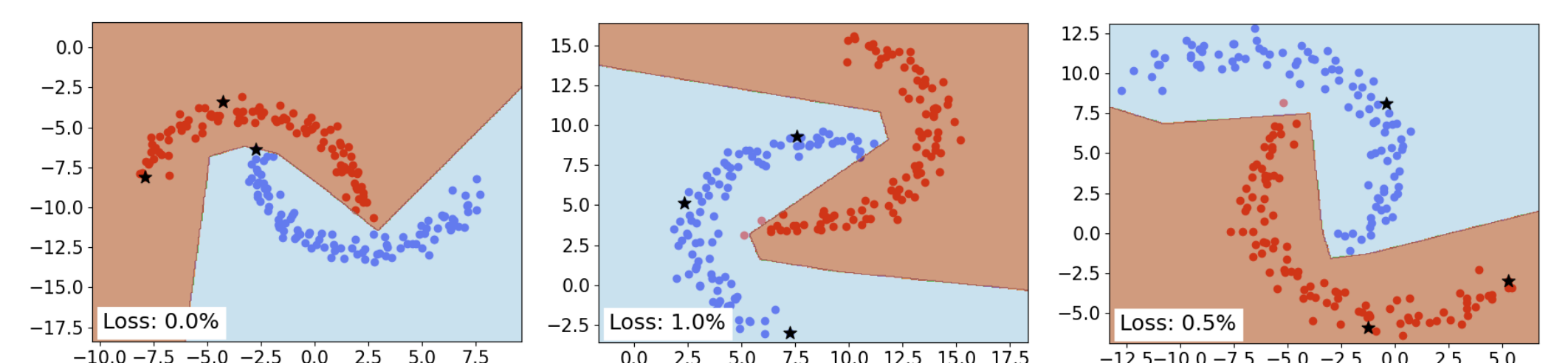


## Numerical experiments

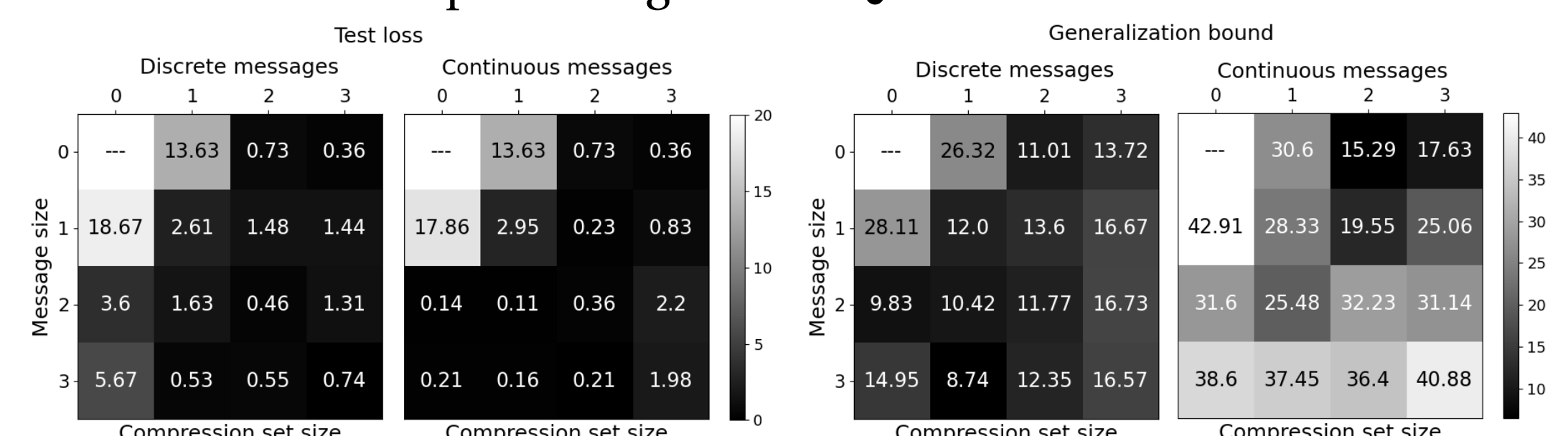
We generated 300 tasks of 200 examples using the following pipeline:



The generated predictor is a single-hidden-layer (of 5 neurons) ReLU MLP. Decision boundaries on three test tasks below (stars \* indicate the compression set examples).



The obtained zero-one loss on 100 test tasks and sample-compressed loss bounds show promising results (stars \*).



- ✓ Message size  $\uparrow \Rightarrow$  test loss  $\downarrow$ , generalization bound  $\searrow \nearrow$
- ✓ Compression set size  $\uparrow \Rightarrow$  test loss  $\downarrow$ , generalization bound  $\searrow \nearrow$