

# PAC-Bayesian Contrastive Unsupervised Representation Learning



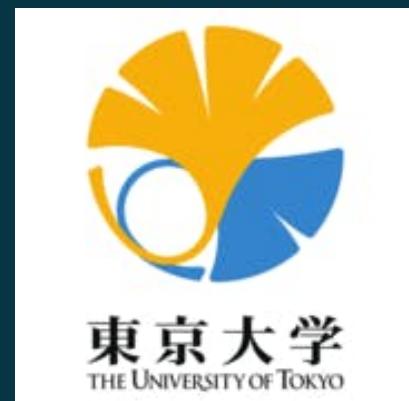
Kento Nozawa



Pascal Germain



Benjamin Guedj



Paper: [http://auai.org/uai2020/proceedings/24\\_main\\_paper.pdf](http://auai.org/uai2020/proceedings/24_main_paper.pdf)  
Code: <https://github.com/nzw0301/pb-contrastive>

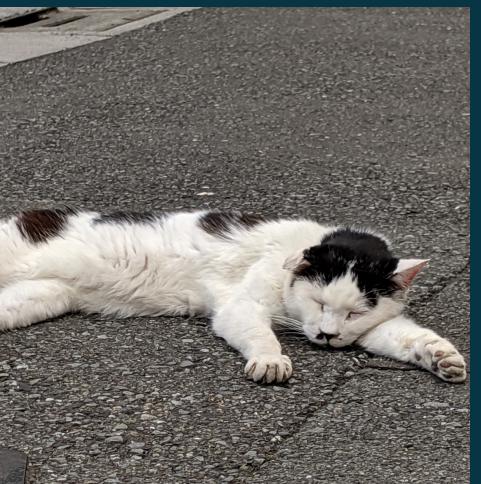
# Contrastive unsupervised representation learning (CURL)

Goal: learn a good feature extractor  $\mathbf{f}$ , e.g. DNNs.

$\mathbf{x}$



similar  $\mathbf{x}^+$

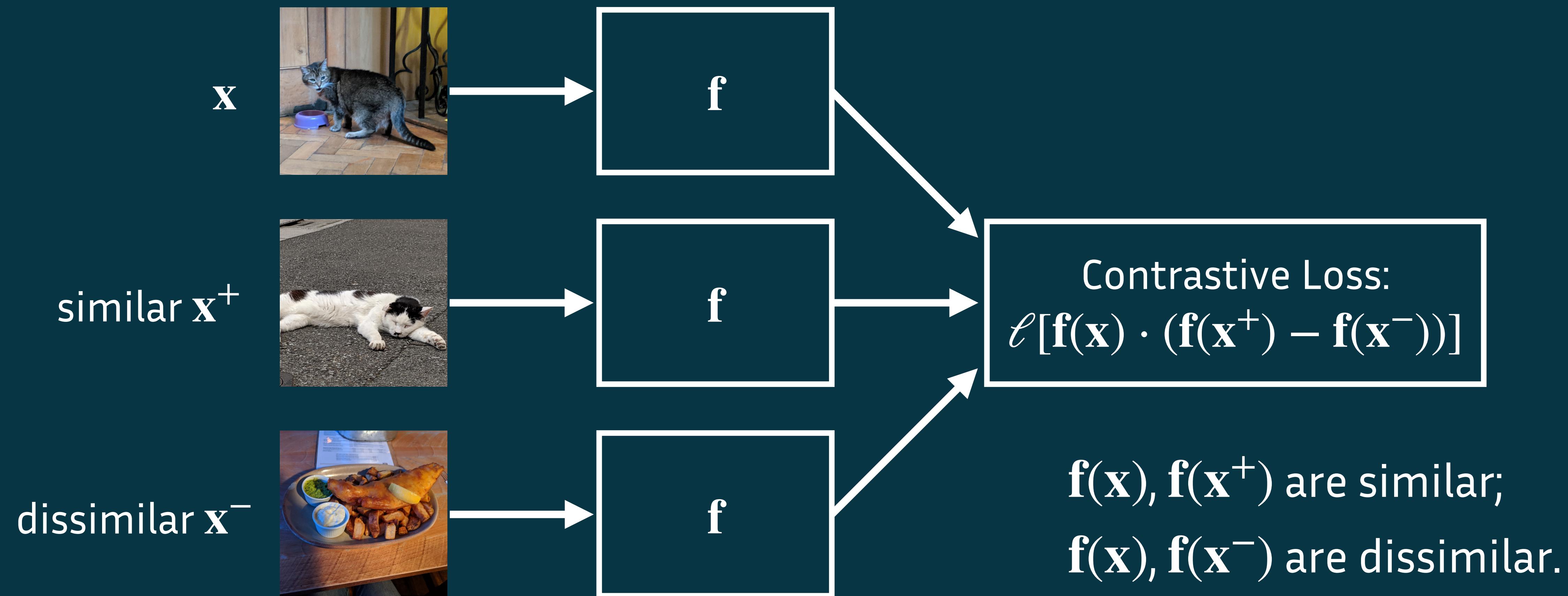


dissimilar  $\mathbf{x}^-$

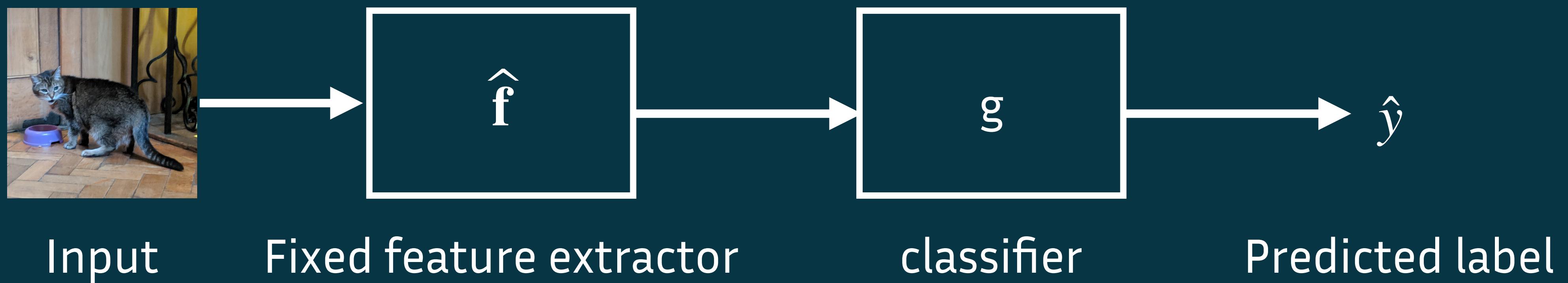


# Contrastive unsupervised representation learning (CURL)

Goal: learn a good feature extractor  $\mathbf{f}$ , e.g. DNNs.



# Learnt representation works for supervised tasks



Why does CURL perform well?

# The first theoretical guarantees for CURL (Arora et al. 2019)

$$\text{Informal bound: } L_{\text{sup}}(\hat{\mathbf{f}}) \leq \underbrace{\alpha L_{\text{un}}(\mathbf{f}) + \mathcal{O}(\mathcal{R}(\mathcal{F}), \delta)}_{\text{Complexity}} \quad \forall \mathbf{f} \in \mathcal{F}$$

- $\alpha$ : Constant
- $\mathcal{R}$  : Rademacher complexity of function class  $\mathcal{F}$ .
- $\delta$ : Confidence of PAC learning

Finding a good representation  $\hat{\mathbf{f}}$  guarantees to generalise well.

# Our contributions

- We show PAC-Bayes bounds for CURL and derive new algorithms by minimising the bounds.
  - We replace the Rademacher complexity term from Arora et al. (2019) with a Kullback–Leibler divergence term, which is easier to compute in general.
  - The PAC-Bayes bound directly suggests a (theory driven) learning algorithm.
- We also show a PAC-Bayes bound for non-iid contrastive data.
  - The iid assumption seems unrealistic in many settings and is unlikely to hold with contrastive datasets.

# General PAC-Bayes

- $Q$ : Posterior. Probability distribution over a function class  $\mathcal{F}$ . It can depend on training data.
- $P$ : Prior. Probability distribution over a function class  $\mathcal{F}$ . It cannot depend on training data.
- $R(Q) = \mathbb{E}_{(\mathbf{x},y) \sim \mathcal{D}} \mathbb{E}_{f \sim Q} \ell(y, f(\mathbf{x}))$ : Expected risk of  $Q$  on test data
- $\widehat{R}(Q) = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{f \sim Q} \ell(y_i, f(\mathbf{x}_i))$ : Expected risk of  $Q$  on train data

Informal bound:

$$R(Q) \leq \underbrace{\alpha \widehat{R}(Q) + \mathcal{O}(\text{KL}(Q\|P), \delta)}_{\text{Complexity}} \quad \forall Q \text{ over } \mathcal{F}, \text{ w.h.p. } 1 - \delta$$

# The first PAC-Bayesian generalisation bound for CURL

Informal bound:

$$L_{\text{Sup}}(Q) \leq \alpha \widehat{L}_{\text{un}}(Q) + \underbrace{O\left(\text{KL}(Q||P), \delta\right)}_{\text{Complexity}} \quad \forall Q \text{ over } \mathcal{F}, \text{ w.h.p. } 1 - \delta$$

- The complexity term is easier to compute than Rademacher one.
- Since all terms in the right-hand side are explicit or easy to approximate, we can minimise the bound directly.

# Learning algorithms & Experiments

- Minimising  $\hat{L}_{\text{un}}(Q) + \mathcal{O}(\text{KL}(Q||P), \delta)$  w.r.t.  $Q$ .
  - $P$  and  $Q$  are multivariate Gaussians with diagonal covariance.
  - We optimise  $Q$ 's mean and covariance by using SGD (Dziugaite and Roy. 2017).
    - Approximate  $\hat{L}_{\text{un}}$  by sampling weights of neural networks from  $Q$ .
- Evaluation procedures:
  - Learning:  $Q$  on contrastive unsupervised data.
  - Evaluation: test 0-1 risk on supervised data by using centroid classifier.

# Experimental results: Supervised performance & bound

	supervised	Arora et al. (2019)				PAC-Bayes based methods					
		$\mu$		$\mu-5$		$\mu$		$\mu-5$		$\mu$	
		$\mu$	$\mu-5$	$\mu$	$\mu-5$	$\mu$	$\mu-5$	$\mu$	$\mu-5$	$\mu$	$\mu-5$
CIFAR-100											
AVG-2 risk ↓	0.086	0.125	0.106	0.144	<b>0.100</b>	<b>0.128</b>	0.246	0.292			
TOP-5 risk ↓	0.422	0.540	0.471	0.574	<b>0.460</b>	<b>0.548</b>	0.766	0.806			
Contrastive test risk $R_{\text{un}}(Q) \downarrow$	–	–	–	–		<b>0.197</b>	0.327				
PAC-Bayes upper bound ↓	–	–	–	–		0.718	0.437				
AUSLAN											
AVG-2 risk ↓	0.198	0.249	<b>0.144</b>	<b>0.167</b>	0.147	0.171	0.174	0.209			
TOP-5 risk ↓	0.643	0.759	<b>0.433</b>	0.518	0.435	<b>0.509</b>	0.494	0.616			
Contrastive test risk $R_{\text{un}}(Q) \downarrow$	–	–	–	–		<b>0.185</b>	0.220				
PAC-Bayes upper bound ↓	–	–	–	–		0.417	0.361				

- **AVG-2 risk**: averaged 0-1 risk over all combination of two classes in supervised data.
- **PAC-Bayes bound**: computed on the stochastic neural networks.

# Conclusion

- We provide the first PAC-Bayes generalisation bounds for CURL.
  - This allows to derive new algorithms by directly optimising the bound.
- More results in the paper:
  - General PAC-Bayes bound for multiple dissimilar samples.
  - Bounds and learning algorithm for the non-iid case.

Paper: [http://auai.org/uai2020/proceedings/24\\_main\\_paper.pdf](http://auai.org/uai2020/proceedings/24_main_paper.pdf)

Code: <https://github.com/nzw0301/pb-contrastive>