

# Understanding prompt engineering does not require rethinking generalization

Victor Akinwande<sup>1</sup>, Yiding Jiang<sup>1</sup>, Dylan Sam<sup>1</sup>, Zico Kolter<sup>1,2</sup>  
ICLR 2024.

<sup>1</sup>Carnegie Mellon University

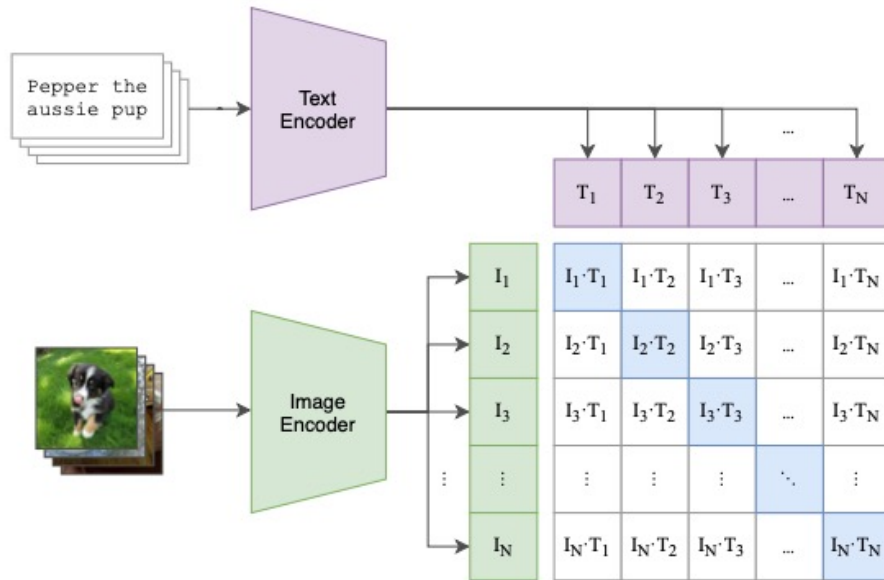
<sup>2</sup>Bosch Center for AI

## Key Takeaway

- The hypothesis space of visual-language models is relatively small – **prompt engineering** or **greedy search** of the set of tokens will not overfit.
- Uniform convergence or **PAC-Bayes bounds** on the discrete space of natural language tokens are remarkably tight and useful for model selection.

# Background: Learning models from natural language supervision

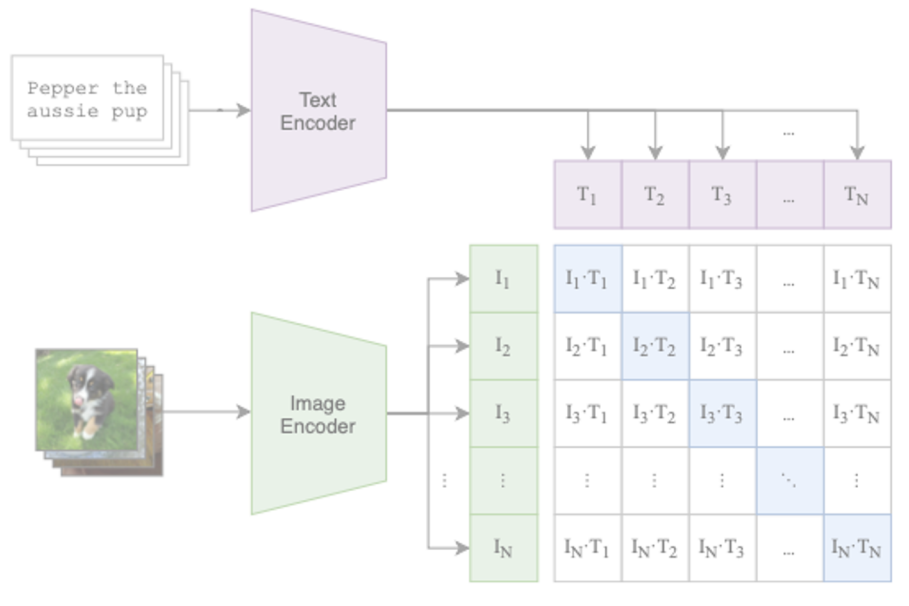
- Obtaining labels is expensive – can we use the vast amount of natural language on the web to build vision models?



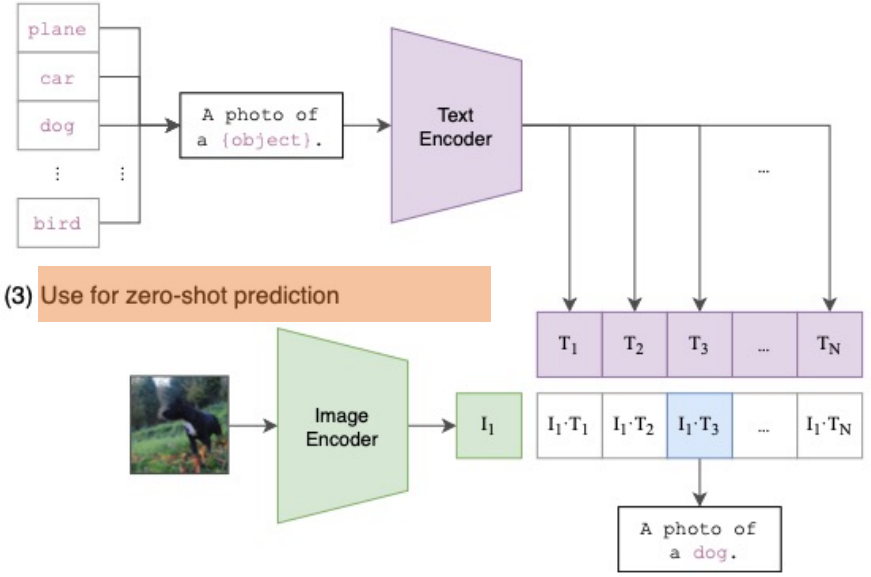
# Background: Learning models from natural language supervision

- Obtaining labels is expensive – can we use the vast amount of natural language on the web to build vision models?

**CLIP**



(2) Create dataset classifier from label text



(3) Use for zero-shot prediction

75% zero-shot on ImageNet

# Background: PAC-Bayes bounds

Posterior over  $\mathcal{H}$  denoted by  $Q$  (defines a randomized classification rule)

Given new instance  $x$ , randomly pick  $h \in \mathcal{H}$  according to  $Q$  and predict  $h(x)$

Generalization loss  $L_D(Q) \stackrel{\text{def}}{=} \mathbb{E}_{h \sim Q} [L_D(h)]$  and training loss  $L_S(Q) \stackrel{\text{def}}{=} \mathbb{E}_{h \sim Q} [L_S(h)]$

Theorem ingredients:  $D$ : arbitrary distribution over domain

0 – 1 loss function

$P$ : Prior distribution over  $\mathcal{H}$

$$L_D(Q) \leq L_S(Q) + \sqrt{\frac{D(Q \parallel P) + \ln m/\delta}{2(m-1)}}$$

Where

$$D(Q \parallel P) \stackrel{\text{def}}{=} \mathbb{E}_{h \sim Q} [\ln(Q(h)/P(h))]$$

is the KL divergence

# Background: PAC-Bayes bounds

Given a prior  $p$ , return a posterior  $Q$  that minimizes the function

$$L_S(Q) + \sqrt{\frac{D(Q \parallel P) + \ln m/\delta}{2(m-1)}}$$

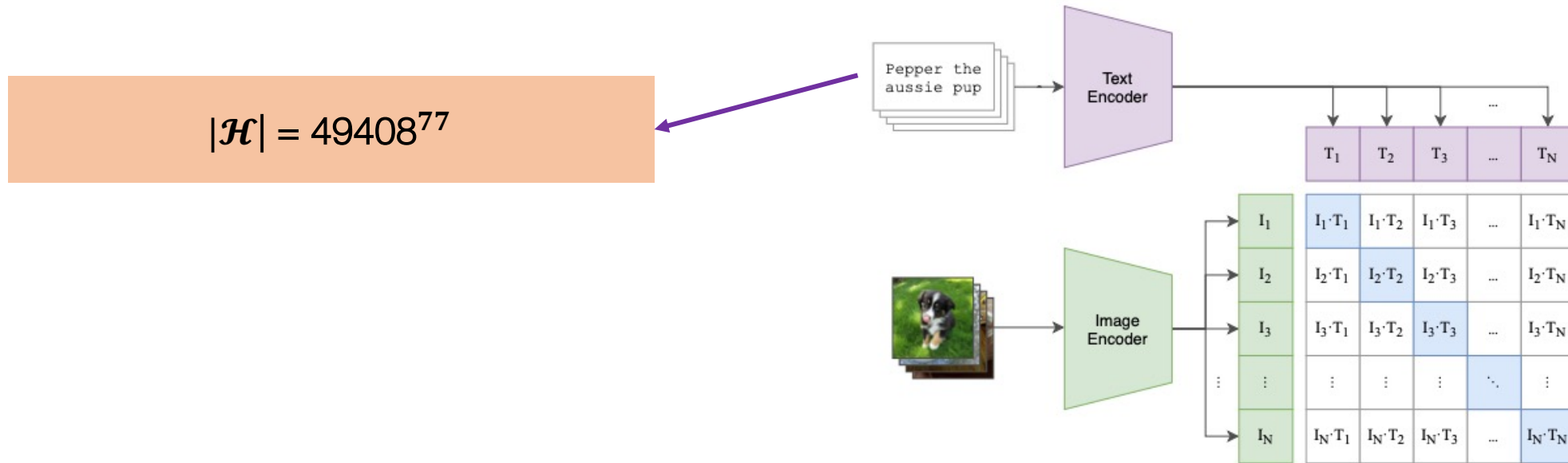
- See (Alquier, 2021 for survey of recent extensions)

Application in (Dziugaite, Roy 2017)

- The posterior is constrained to be Gaussian
- Upper bounds the empirical risk by a convex, Lipschitz upper bound that only has to be minimized w.r.t to the parameters
- Data dependent prior  $\mathcal{N}(\mathbf{w}_0, \sigma^2 \mathbf{I})$  where  $\sigma^2$  is chosen to minimize the bound
- Empirical bounds between 0.16 and 0.22 on MNIST

# The hypothesis class of Prompts

- The vocab size of CLIP is 49408 tokens
- The context length is 77 max



## The hypothesis class of Prompts

- Suppose  $\mathcal{H}$  is a finite hypothesis class, and we set the prior to be uniform over  $\mathcal{H}$  and posterior  $Q(h_S) = 1$  for some  $h_S$  and  $Q(h) = 0$  for other  $h \in \mathcal{H}$

$$|\mathcal{H}| = 49408^{77}$$

$$L_D(h_S) \leq L_S(h) + \sqrt{\frac{\ln(|\mathcal{H}|) + \ln m/\delta}{2(m-1)}}$$



# The hypothesis class of Prompts

- Suppose  $\mathcal{H}$  is a finite hypothesis class, and we set the prior to be uniform over  $\mathcal{H}$  and posterior  $Q(h_S) = 1$  for some  $h_S$  and  $Q(h) = 0$  for other  $h \in \mathcal{H}$

$$|\mathcal{H}| = 49408^{77}$$

$$L_D(h_S) \leq L_S(h) + \sqrt{\frac{\ln(|\mathcal{H}|) + \ln m/\delta}{2(m-1)}}$$

$$\sqrt{\frac{770 * \ln(49408) + \ln(50000/.01)}{2(50000 - 1)}} = 0.29$$

\*Non-vacuous bound on CIFAR-10

# An improved bound with PAC-Bayes

- We know how to model natural language well
- Therefore if

$$D(Q \parallel P) = \sum_{h \in \mathcal{H}} Q(h) \ln \frac{Q(h)}{P(h)} = \ln \frac{1}{P(\hat{h})} = - \sum_{i=1}^K \sum_{j=1}^L \ln p_{\text{LM}}(\hat{h}_j^i \mid \hat{h}_{\leq j}^i)$$

$$L_{\mathcal{D}}(Q) \leq L_S(Q) + \sqrt{\frac{- \sum_{i=1}^K \sum_{j=1}^L \ln p_{\text{LM}}(\hat{h}_j^i \mid \hat{h}_{\leq j}^i) + \ln m / \delta}{2(m-1)}}$$

# An improved bound with PAC-Bayes

- We know how to model natural language well
- Therefore if

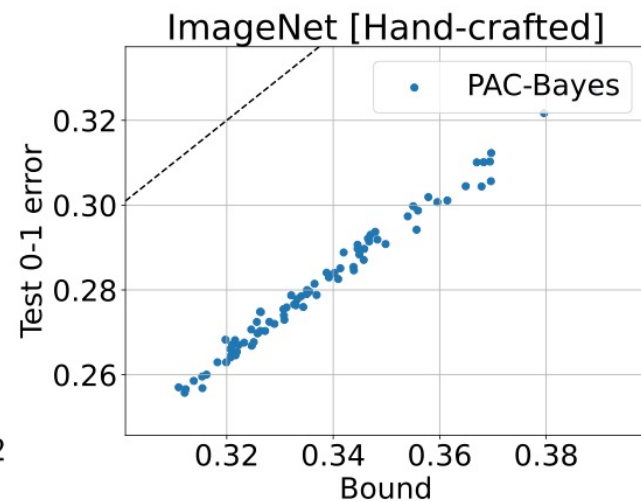
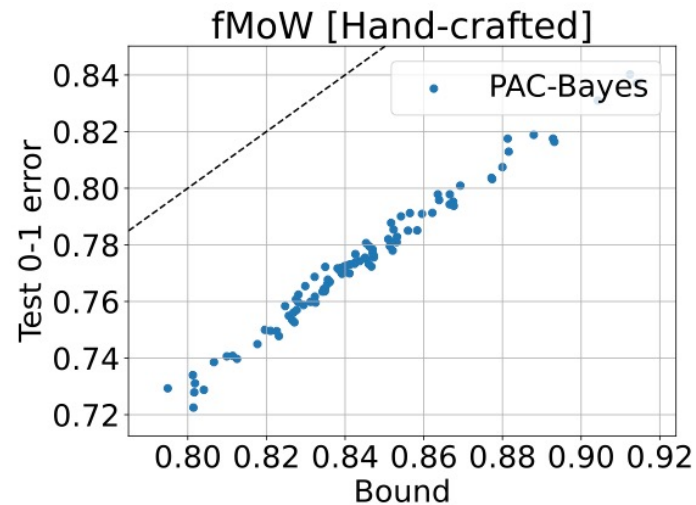
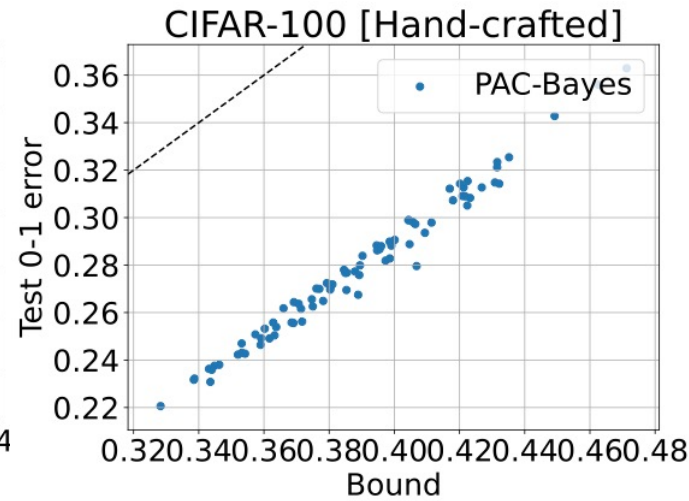
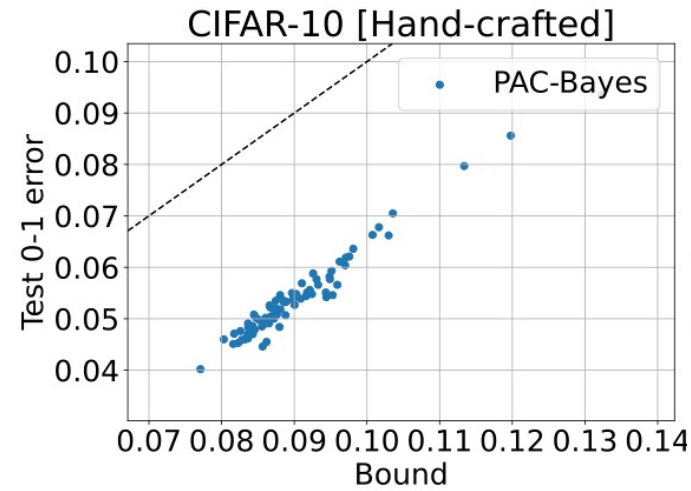
$$D(Q \parallel P) = \sum_{h \in \mathcal{H}} Q(h) \ln \frac{Q(h)}{P(h)} = \ln \frac{1}{P(\hat{h})} = - \sum_{i=1}^K \sum_{j=1}^L \ln p_{\text{LM}}(\hat{h}_j^i \mid \hat{h}_{\leq j}^i)$$

$$L_{\mathcal{D}}(Q) \leq L_S(Q) + \sqrt{\frac{- \sum_{i=1}^K \sum_{j=1}^L \ln p_{\text{LM}}(\hat{h}_j^i \mid \hat{h}_{\leq j}^i) + \ln m / \delta}{2(m-1)}}$$

For a fixed length  $k$ ,

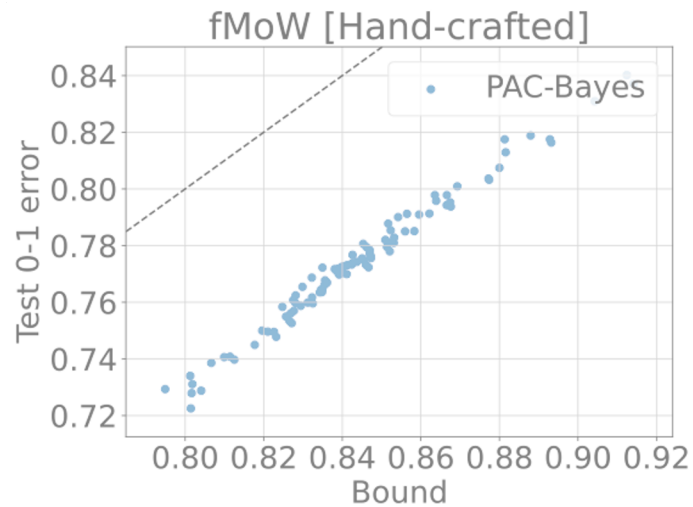
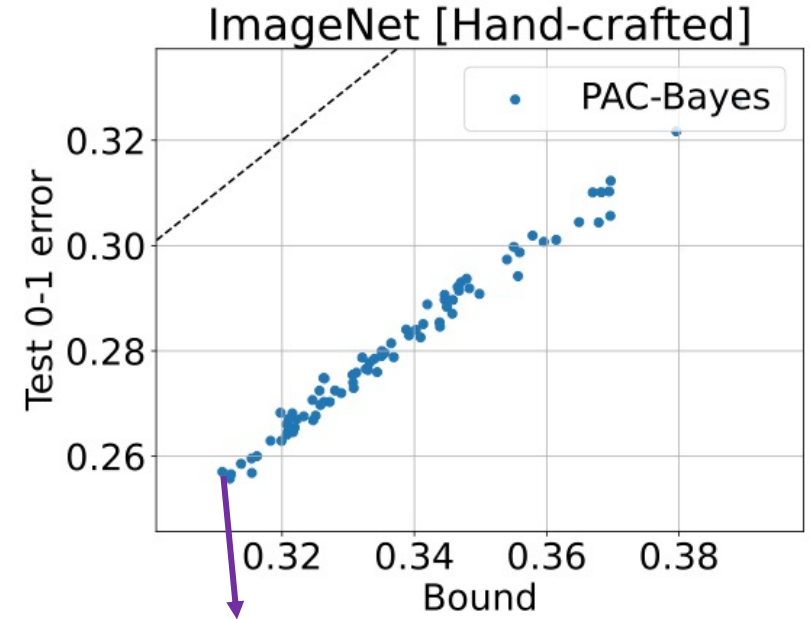
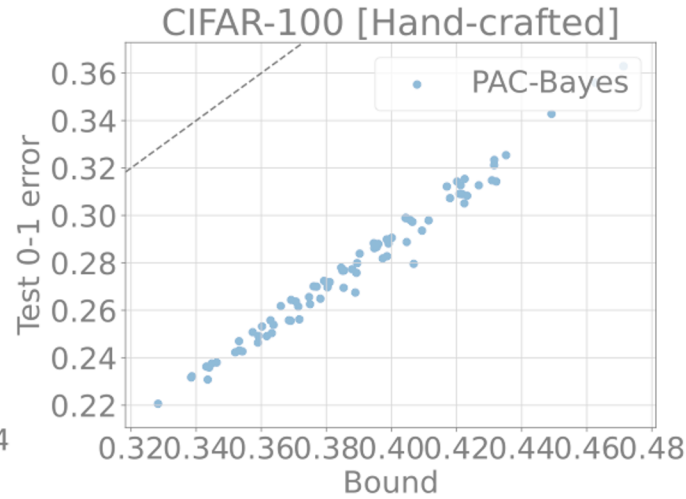
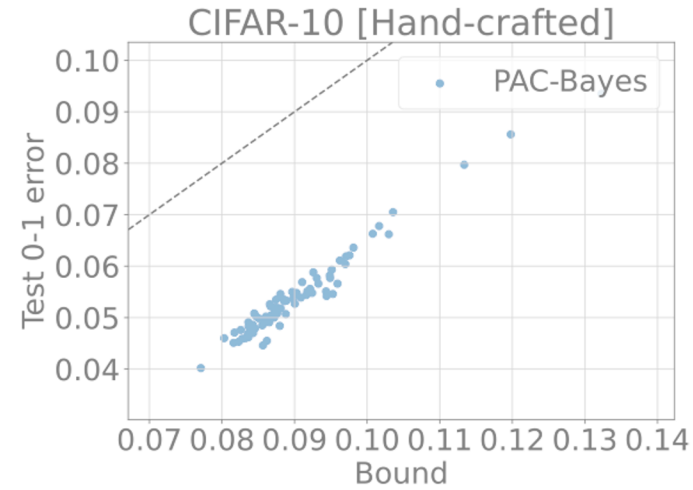
$$\mathbb{E}_{h \sim Q} [L_{\mathcal{D}}(h)] = L_{\mathcal{D}}(h)$$

# Risk bounds for hand-crafted prompts



dotted line:  $y = x$  | Llama -7b prior

# Risk bounds for hand-crafted prompts



A photo of a nice {class\_name}

# Greedy search yields a good posterior

---

**Algorithm 1** Sequential Prompt Search (Greedy)

---

```
1:  $\theta \leftarrow [\text{initial\_prompt}] \times K$ 
2: for  $l = 0$  to  $L - 1$  do
3:    $\text{class\_order} \leftarrow$  randomly sampled order of class indices
4:   for  $k$  in  $\text{class\_order}$  do
5:      $\text{criteria} \leftarrow -\infty$ 
6:     for  $v$  in  $\widehat{\mathcal{V}}(\theta)$  do
7:        $\text{score} \leftarrow \mathcal{J}(v, \theta_{\leq l}^k, \theta^{-k})$ 
8:       if  $\text{score} > \text{criteria}$  then
9:          $\text{criteria} \leftarrow \text{score}$ 
10:         $\theta_{l+1}^k \leftarrow v$ 
11:       end if
12:     end for
13:   end for
14: end for
15: Return  $\theta$ 
```

▷ This step is vectorized in practice.  
▷ Evaluate the score of  $v$ .  
▷ Keep the prompt with best performance.  
▷ Update the current best score.  
▷ Update  $\theta^k$  with the better token.

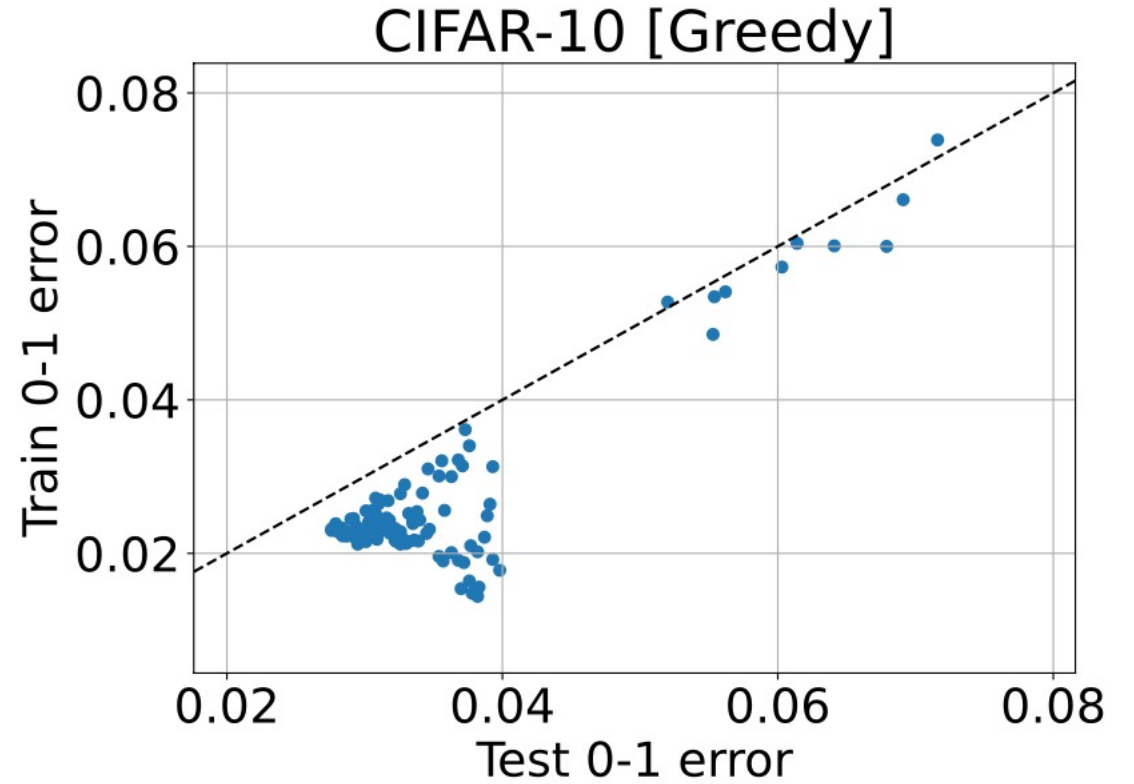
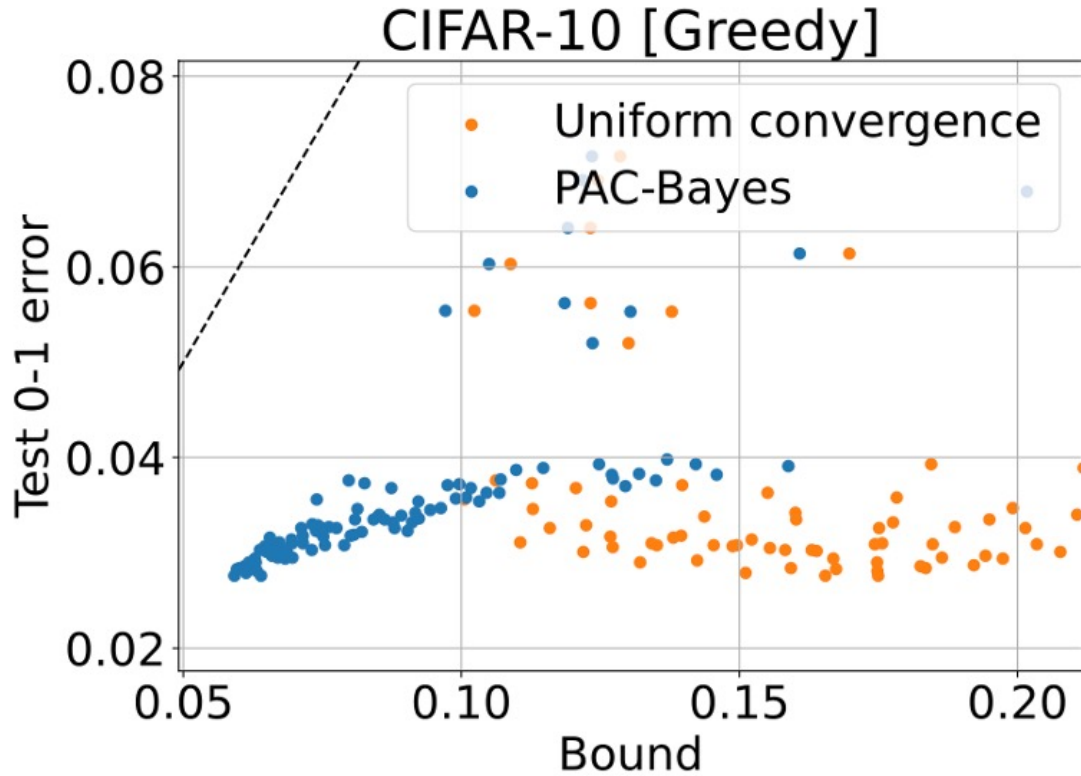
---

# Comparison with state-of-the-art generalization bounds

<b>Dataset</b>	Zhou et al. [2019]	Dziugaite et al. [2021]	Lotfi et al. [2022]	<b>Ours</b>
CIFAR-10	—	0.230*	0.582 / 0.166*	<b>0.059</b>
CIFAR-100	—	—	0.946 / 0.444*	<b>0.251</b>
ImageNet	0.965	—	0.930 / 0.409*	<b>0.312</b>

Data dependent prior involves using a validation set to estimate the bound

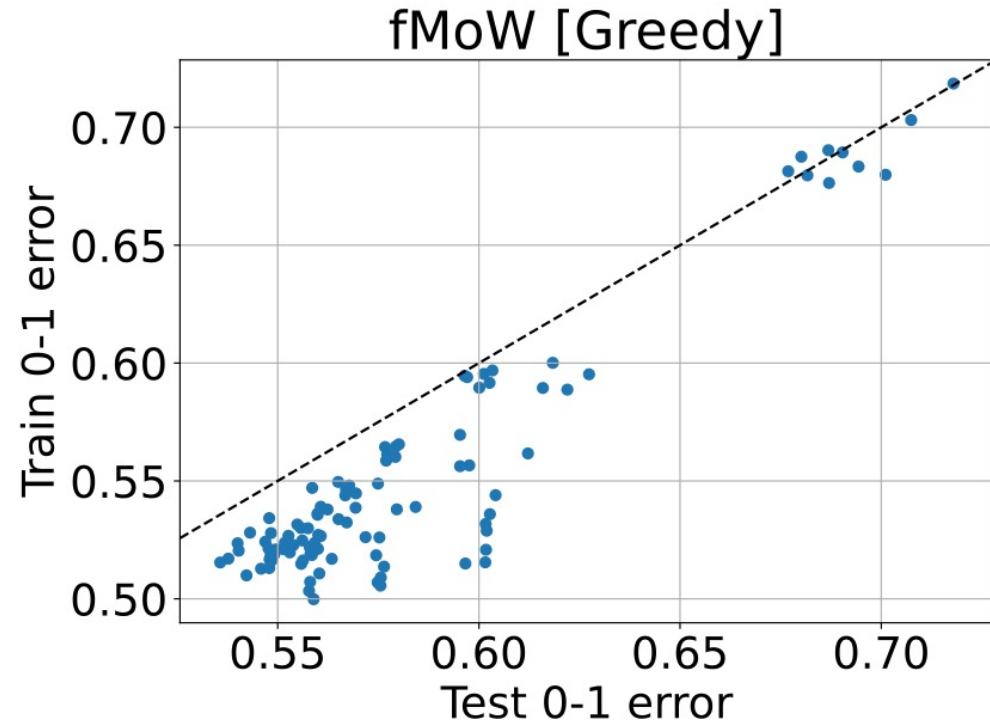
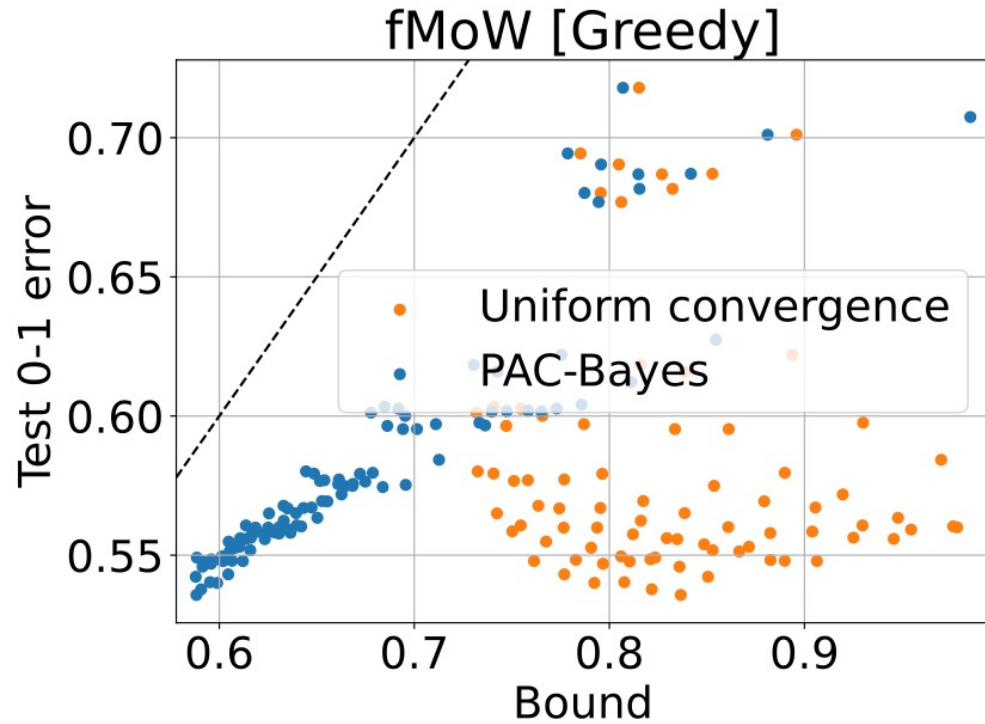
# PAC-Bayes bounds are useful for model selection



dotted line:  $y = x$

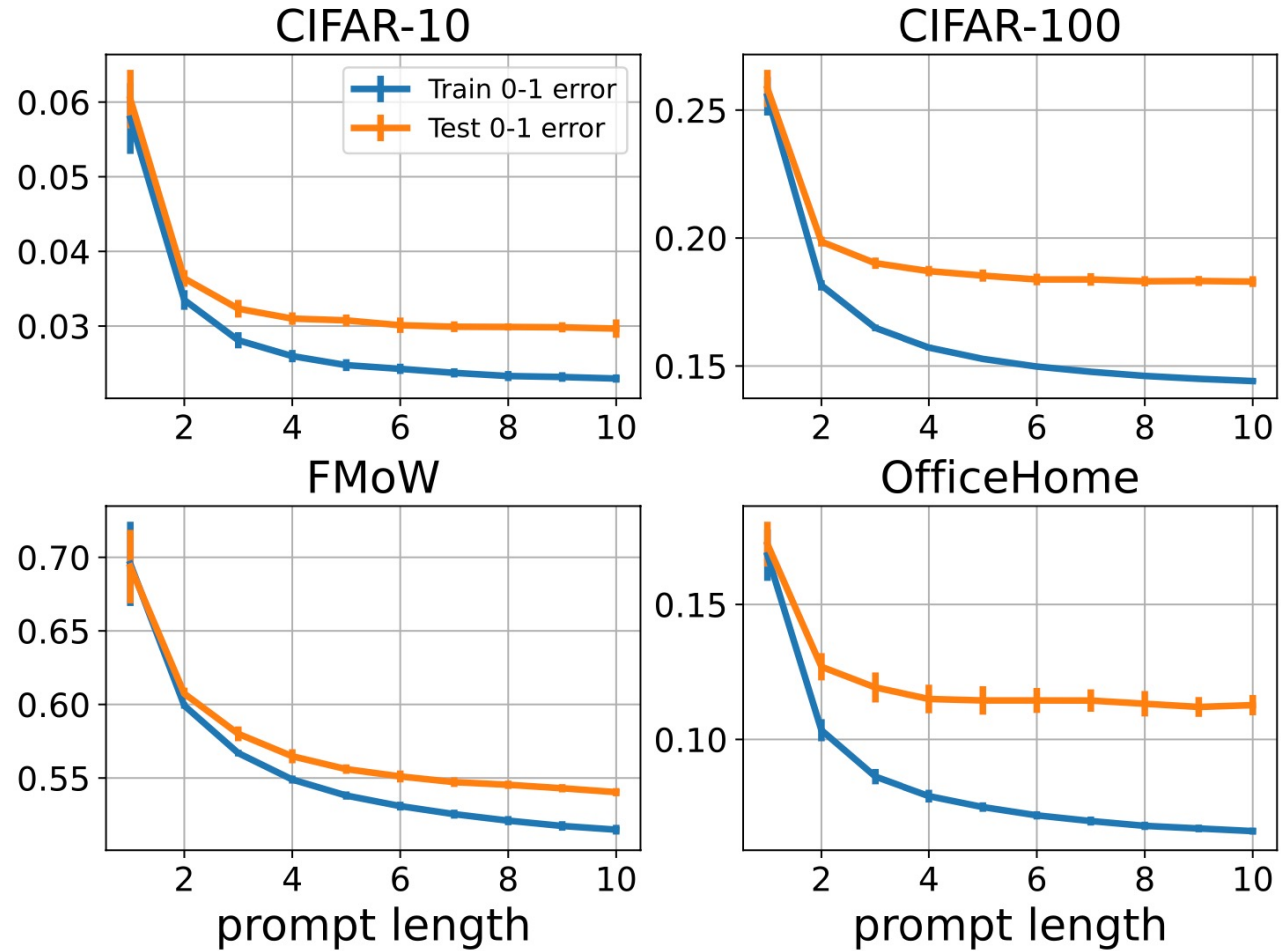


# PAC-Bayes bounds are useful for model selection

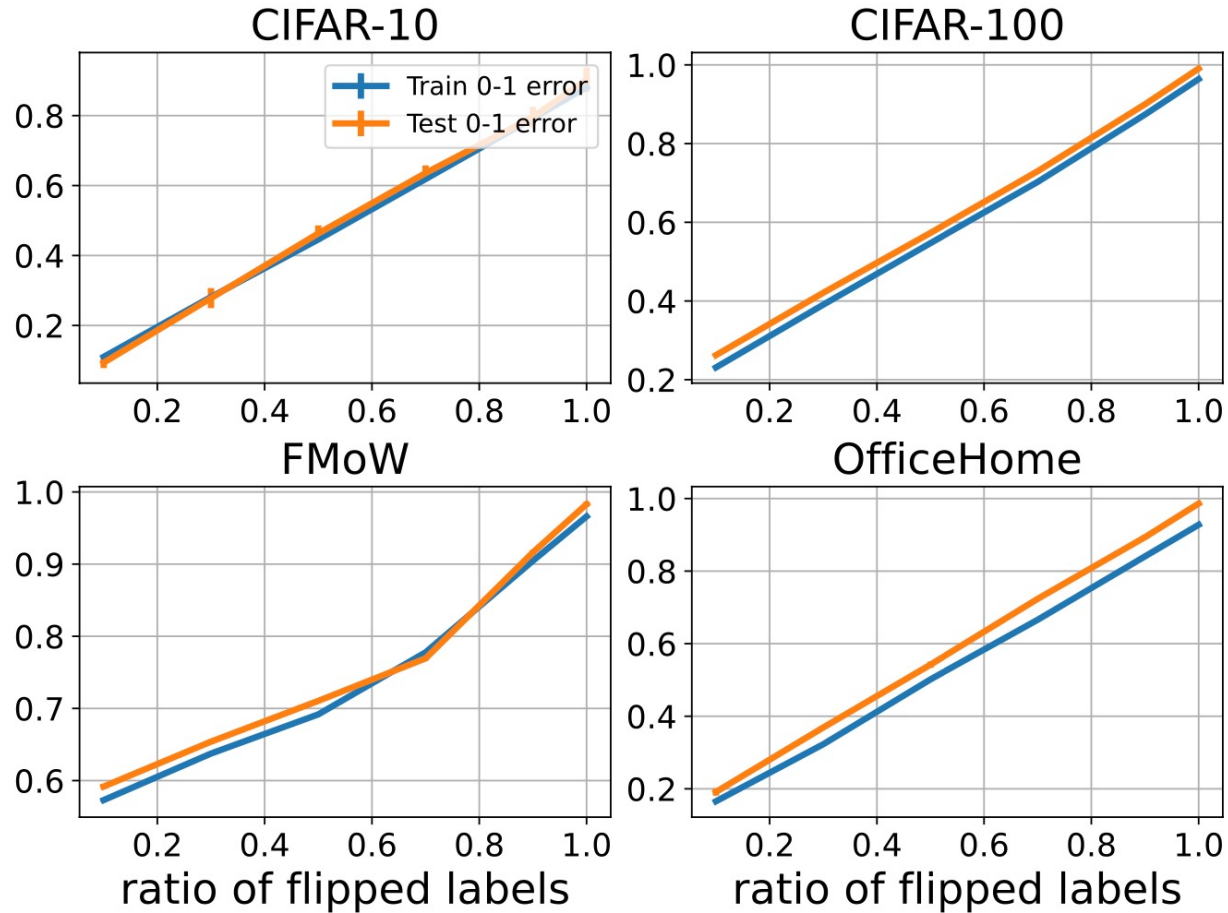


62 class dataset of satellite images

# Greedy rarely overfits



# Greedy does not fit random labels



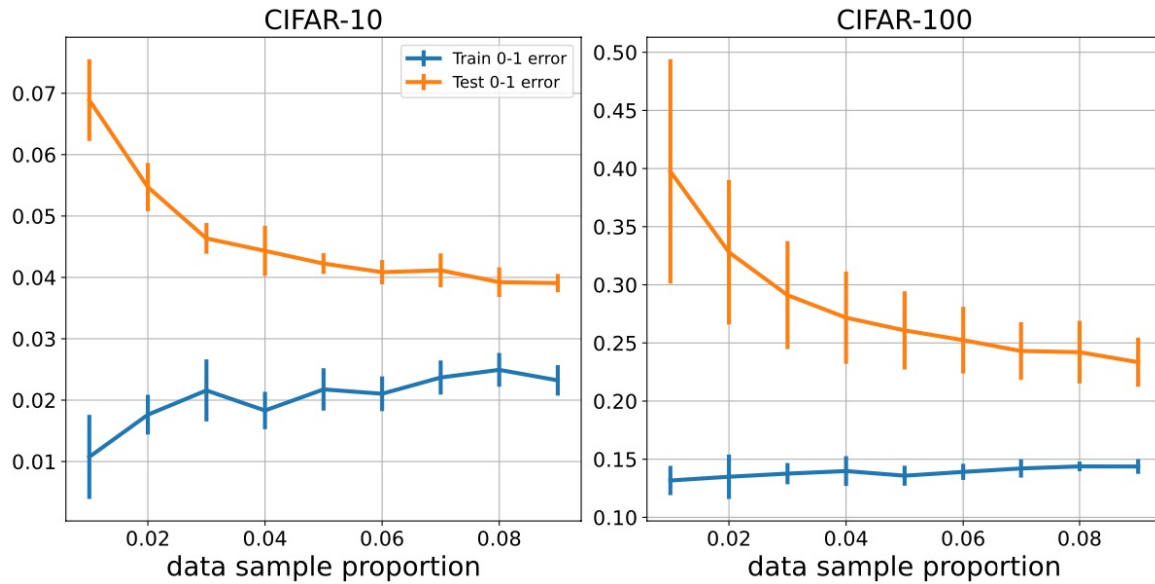
Recall:

Conventional deep neural networks trained with SGD can fit both random labels and random data

Zhang et. al. 2017

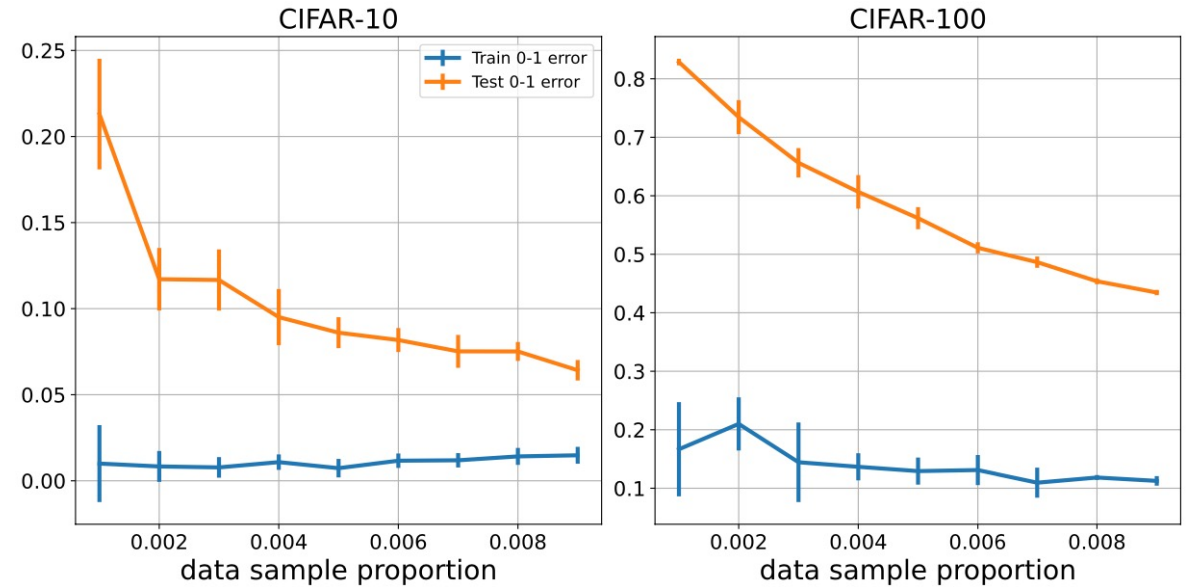
# Greedy can be remarkably data efficient

Learning w' 1 - 10% of the data



<2% increase in error

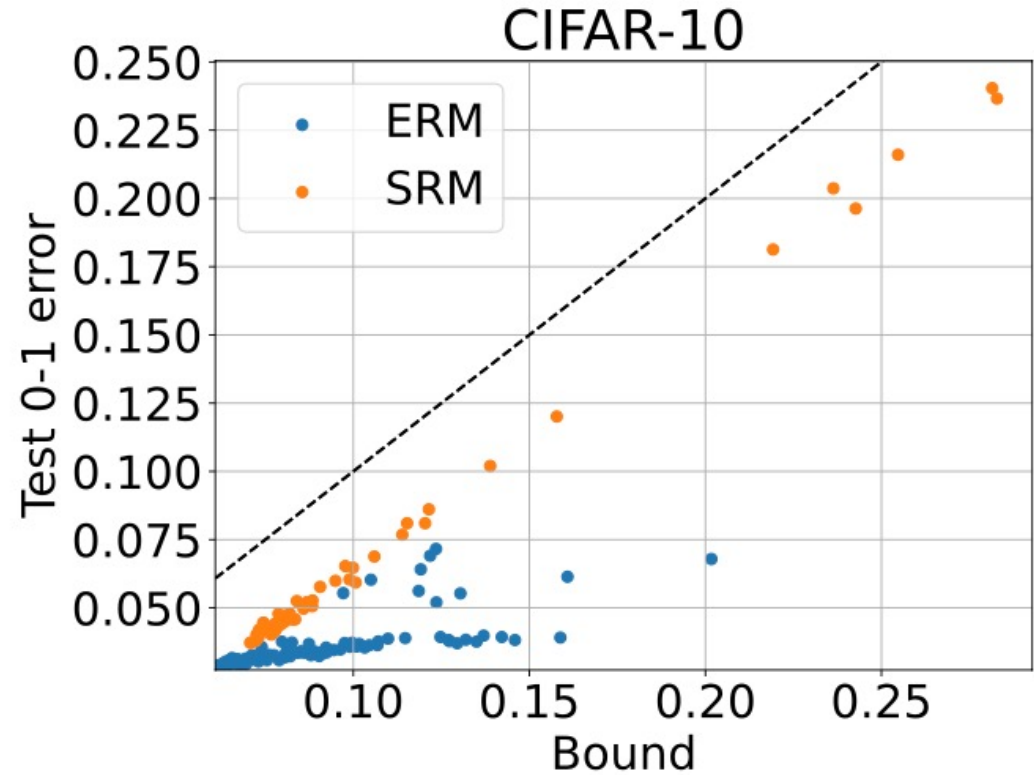
Learning w' 0.1 - 1% of the data



# Structural risk minimization

Given a prior, return a posterior  $Q$  that minimizes

$$L_S(Q) + \sqrt{\frac{-\sum_{i=1}^K \sum_{j=1}^L \ln p_{\text{LM}}(\hat{h}_j^i | \hat{h}_{\leq j}^i) + \ln m/\delta}{2(m-1)}}$$



Using the Llama vocabulary

# Learned prompts may not interpretable

[airplane, automobile, bird, cat, deer, dog, frog, horse, ship, truck]

aviarist nonsonant confirmment establishmentism hemiteratics

nonmotoring known vaticinal allot nth

ornithophile slimsy renishly redivive muchness

wheencat compearant stintedness osiery thisness

stagnature unchawed lophobranh primariness primariness

dogrib babooism pneumococcic kaoliang bogusness

froggerly phthalid auhuhu rippit hideousness

horsemonger fooyoung inordinary spreadingness forthbring

seafaring rumness tralatition babeship knocking

truckling phthartolatrae semantology waywarden decess

96.74% accuracy on CIFAR10

vocab: pip install english-words

## Summary & Key takeaway

1. Given pretrained models, manual prompt engineering (even when “overfitting” to a test set) often exhibits surprisingly strong generalization behaviour.
2. Uniform convergence or PAC-Bayes bounds on the discrete space of natural language tokens are remarkably tight and useful for model selection.