

# A Generalization Theory for Zero-Shot Prediction

Paper



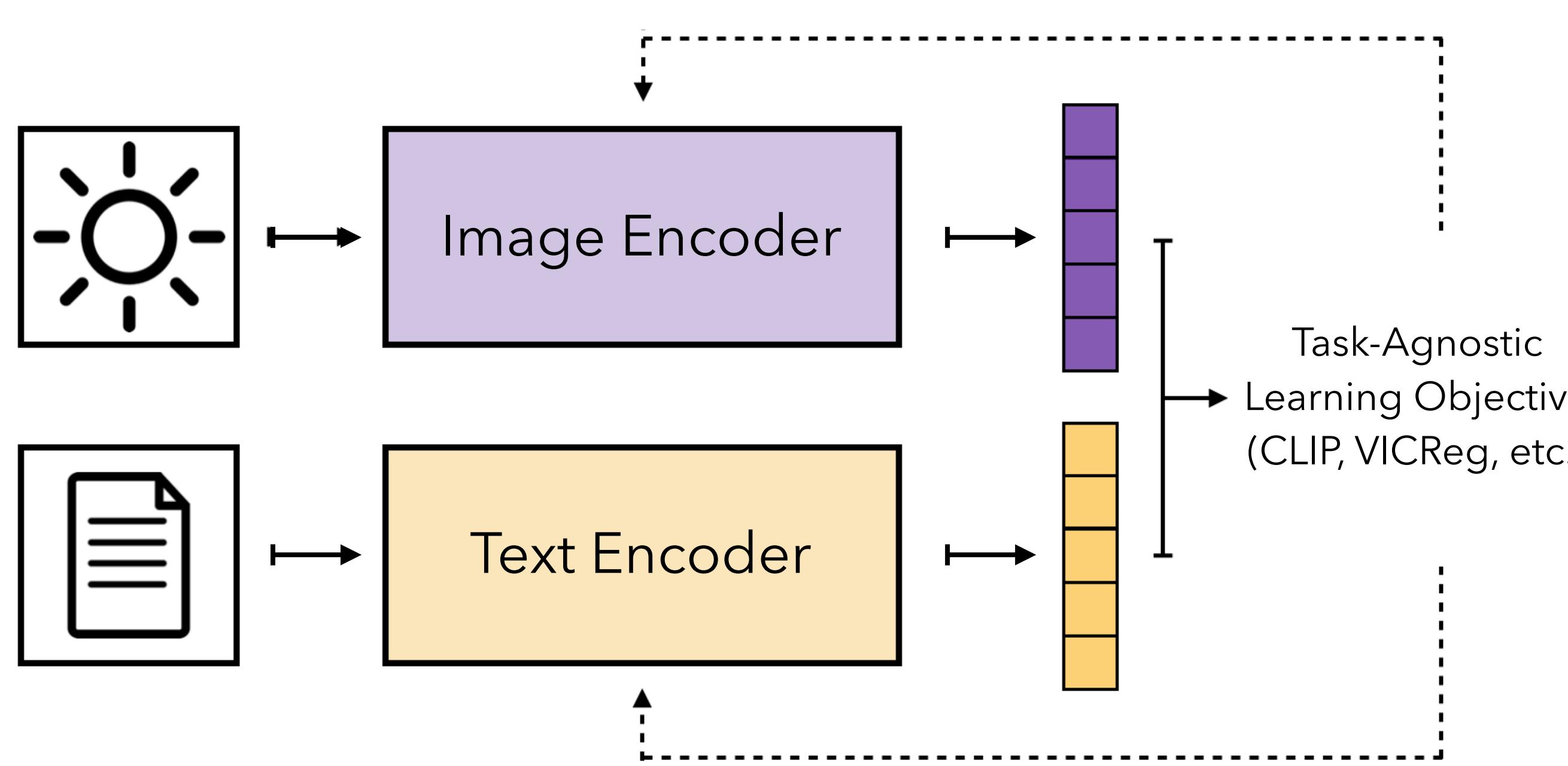
Ronak Mehta and Zaid Harchaoui

## Zero-Shot Prediction (ZSP)

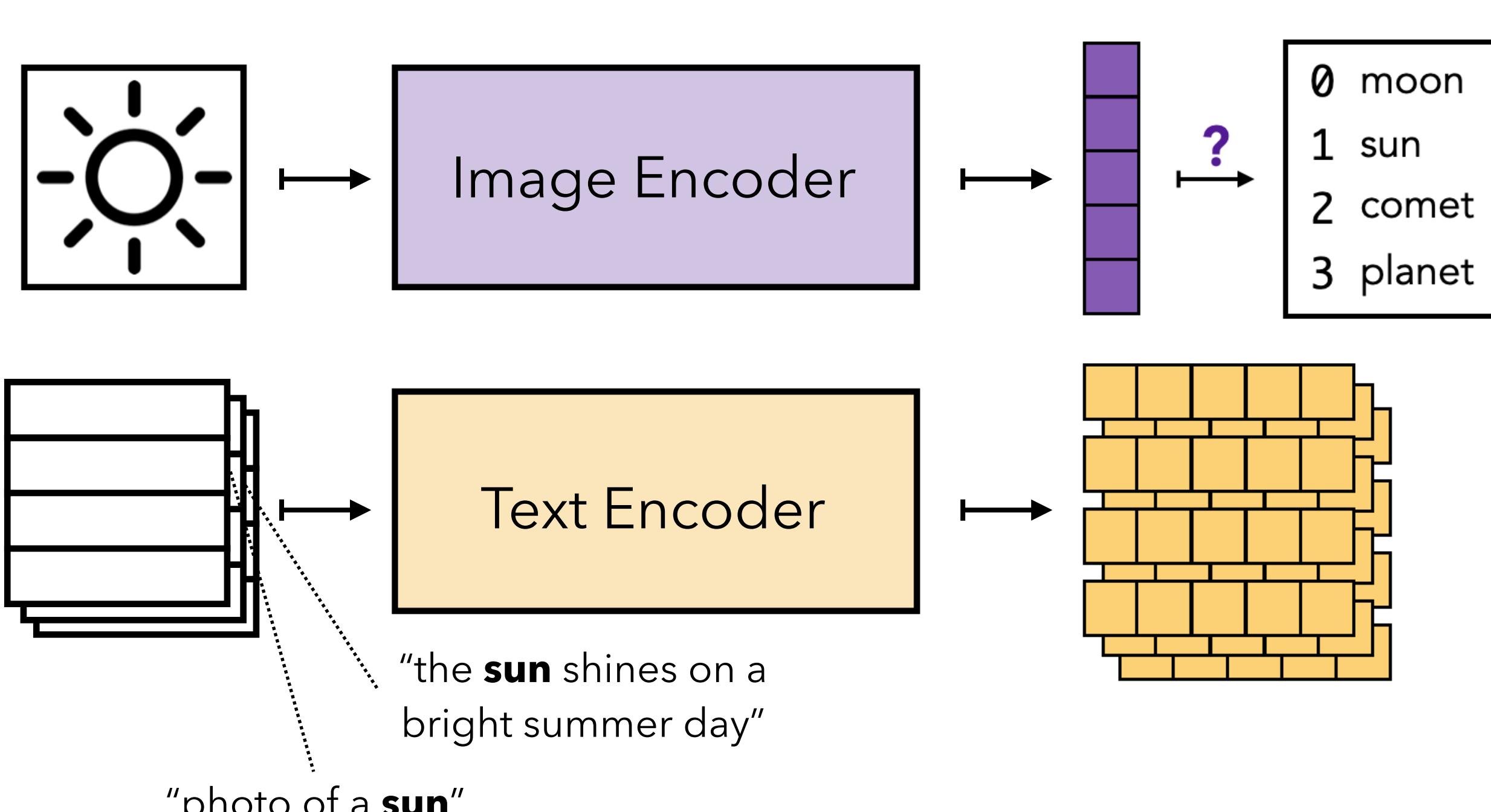
**Motivation:** Zero-shot prediction is a modern method that reuses foundation models to build classifiers for tasks without seeing *any* directly labeled training data.

Need for theoretical understanding has arisen.

### Contrastive Pre-Training



### Evaluation



**Research Question:** How does the downstream performance of ZSP depend on the pre-training distributions, downstream task distribution, and prompting strategy?

## Theoretical Framework

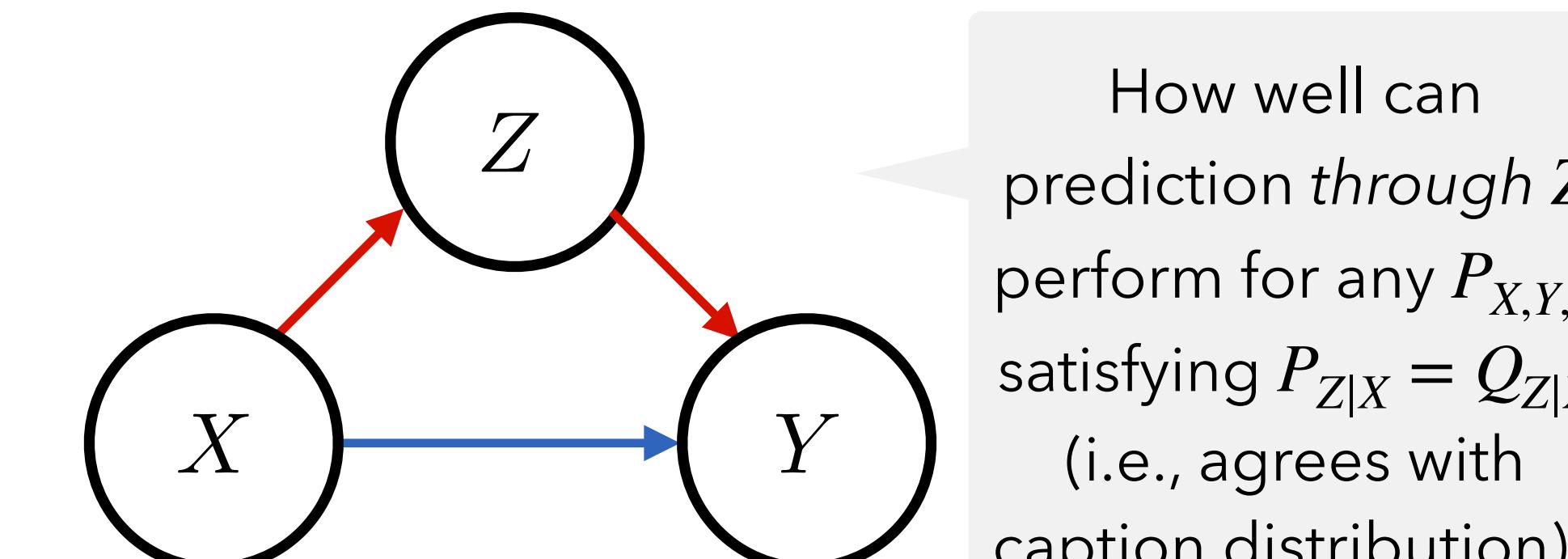
Fundamental limits of ZSP rely on the compatibility of three distributions.

$$P_{X,Y} \quad \text{Evaluation}$$

$$Q_{X,Z} \quad \text{Pre-Training}$$

$$\rho_{Y,Z} \quad \text{Prompting}$$

$X$  = image  
 $Y$  = label  
 $Z$  = caption



How well can prediction through  $Z$  perform for any  $P_{X,Y,Z}$  satisfying  $P_{Z|X} = Q_{Z|X}$  (i.e., agrees with caption distribution)?

Direct Predictor  $f_*(x) = \mathbb{E}_{P_{X,Y}} [Y|X = x]$

Indirect Predictor (Population Version of ZSP)  $\bar{f}(x) = \mathbb{E}_{Q_{X,Z}} [\mathbb{E}_{\rho_{Y,Z}} [Y|Z] | X = x]$

## Main Results

Error decomposition for ZSP procedures.

$$\mathbb{E}_{X \sim P_X} [(f_*(X) - \hat{f}(X))^2] \leq 2\mathbb{E}_{X \sim P_X} [(f_*(X) - \bar{f}(X))^2] + 2\mathbb{E}_{X \sim P_X} [(\bar{f}(X) - \hat{f}(X))^2]$$

information-theoretic error                                    learning error

### Theorem.

$$\mathbb{E}_{X \sim P_X} [(f_*(X) - \bar{f}(X))^2] \lesssim I(X, Y|Z) + \text{err}(P_{Y,Z}, \rho_{Y,Z})$$

**Interpretation:** Residual dependence between image/label not explained by text.

**Interpretation:** Bias of prompt distribution.

### Theorem.

$$\mathbb{E}_{X \sim P_X} [(\bar{f}(X) - \hat{f}(X))^2] \lesssim C_N(Q_{X,Z}) + C_M(\rho_{Y,Z})$$

**Interpretation:** Complexity of learning foundation model (e.g., CLIP) from  $N$  pre-training examples.

**Interpretation:** Complexity of approximating prompt distribution with  $M$  prompts.

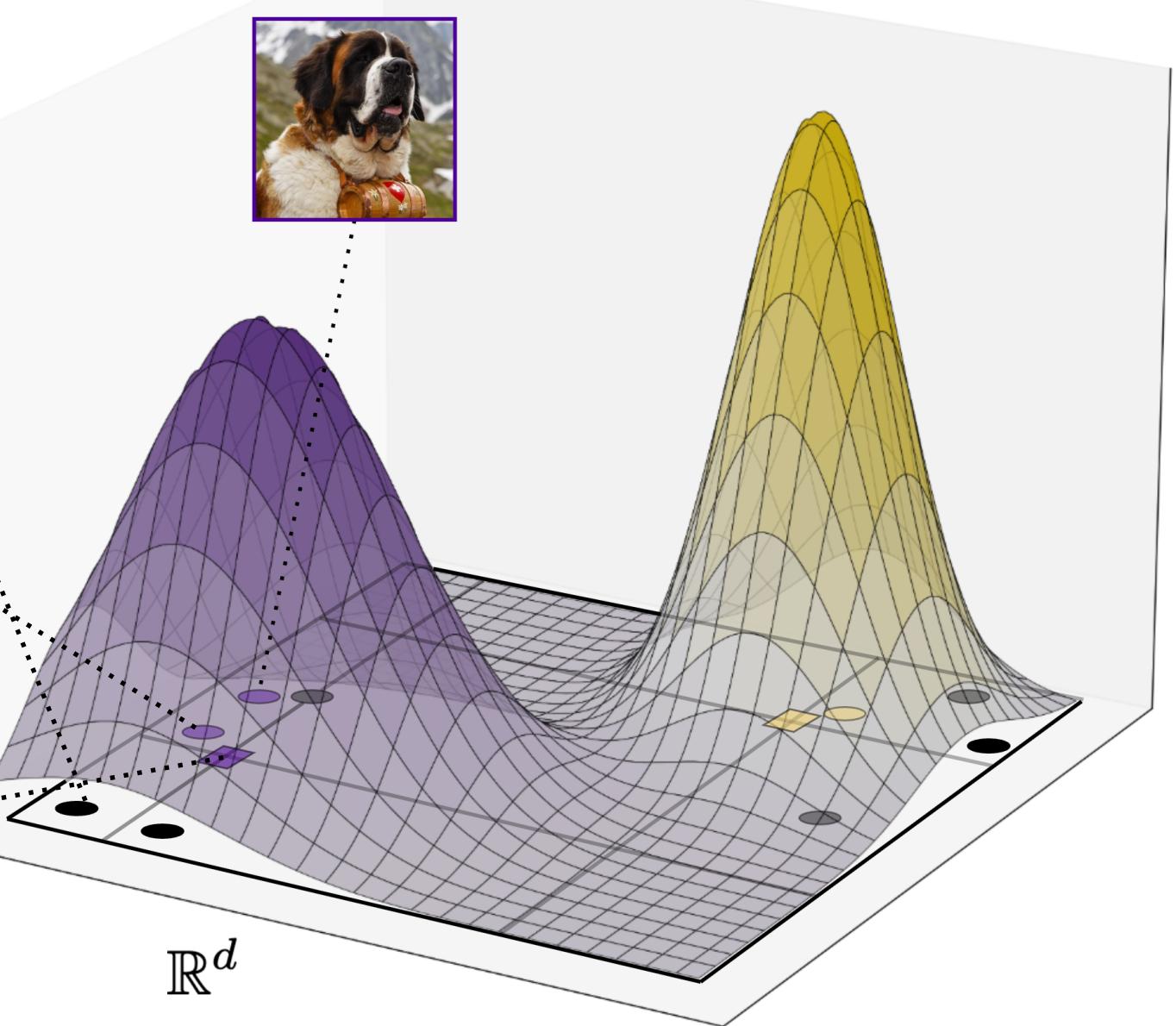
## Prompting Strategies

Distribution of Text Embeddings

Template-Based  
"photo of a dog"

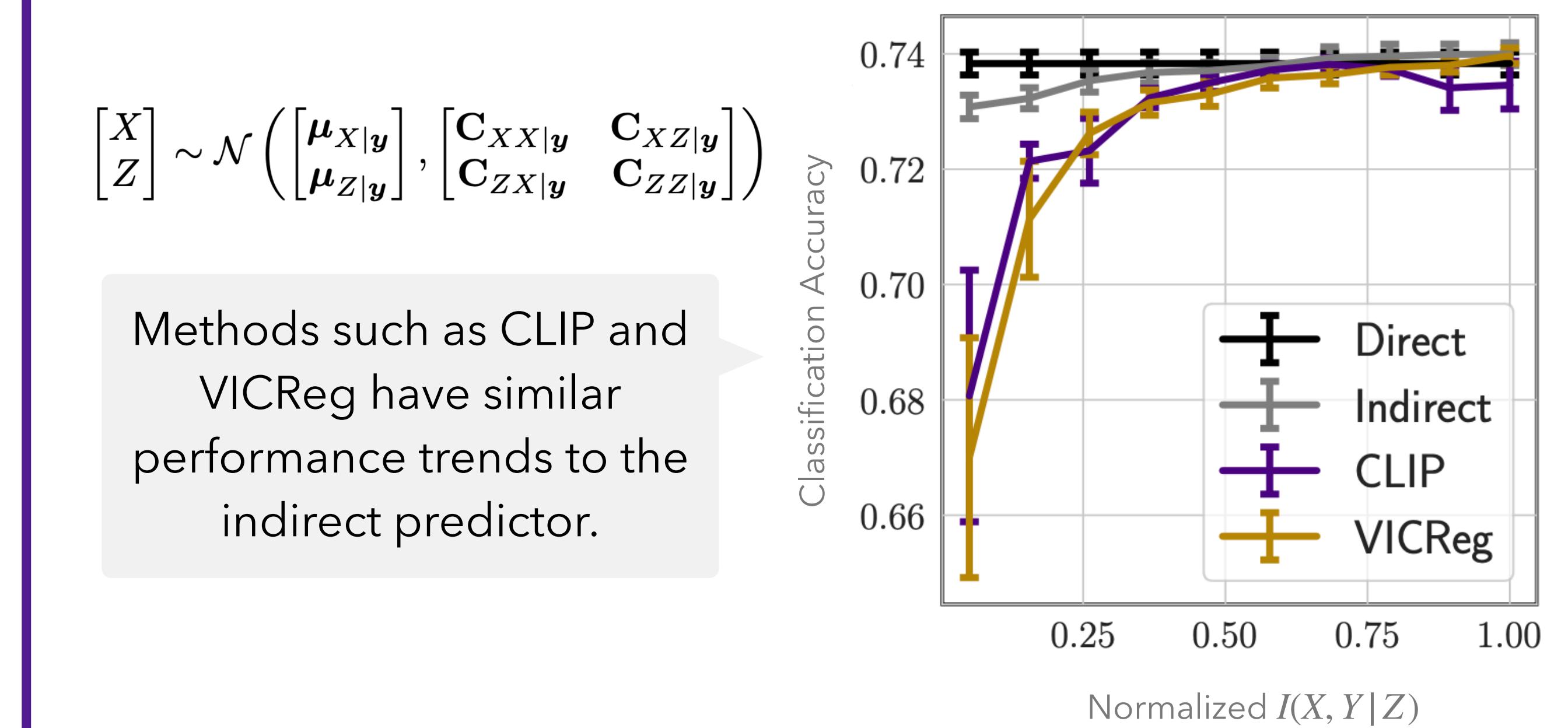
Class-Conditional  
"st. bernard rescue near me"

Unbiased  
sample directly from  $P_{Y,Z}$

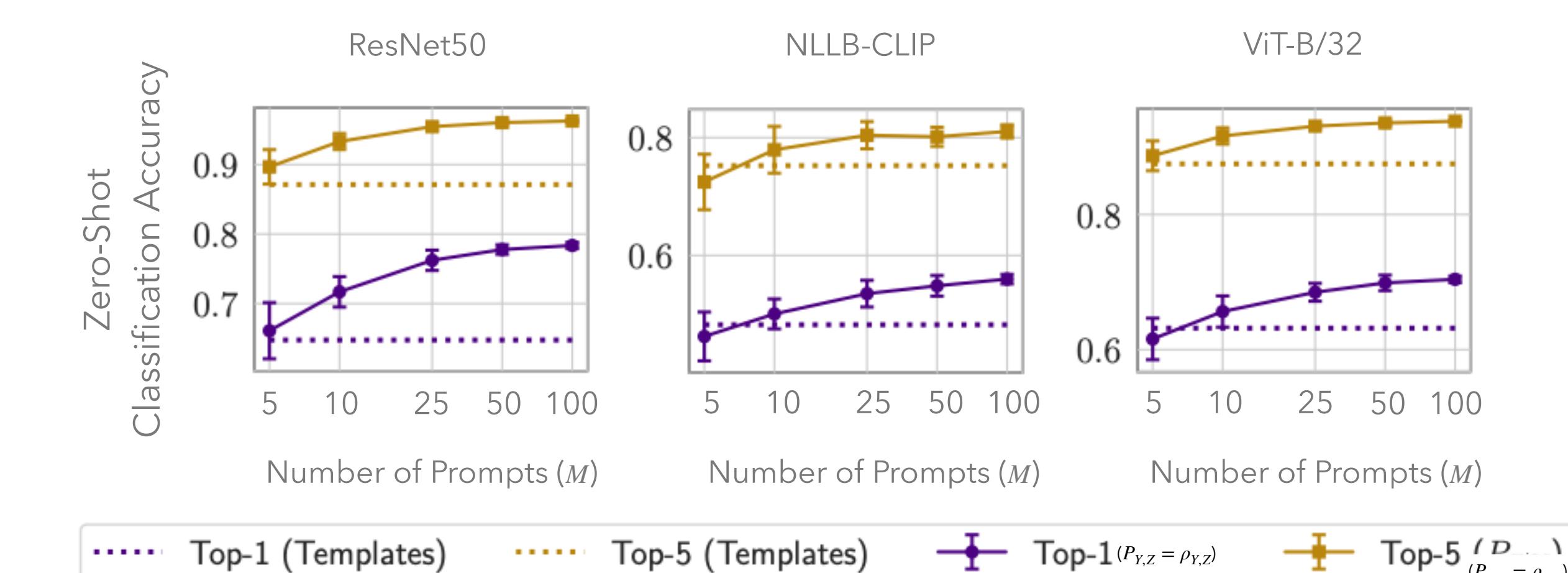


## Experiments

**Synthetic Data:** Controllable Residual Dependence and Prompt Bias



**Real Data:** Language-Image Pre-Training and Zero-Shot Image Classification



When  $P_{Y,Z}$  is observable, unbiased prompting strategy outperforms template-based.