



# Variational Adversarial Active Learning

---

Samarth Sinha \*  
University of Toronto

samarth.sinha@mail.utoronto.ca

Sayna Ebrahimi \*  
UC Berkeley

sayna@eecs.berkeley.edu

Trevor Darrell  
UC Berkeley

trevor@eecs.berkeley.edu

---

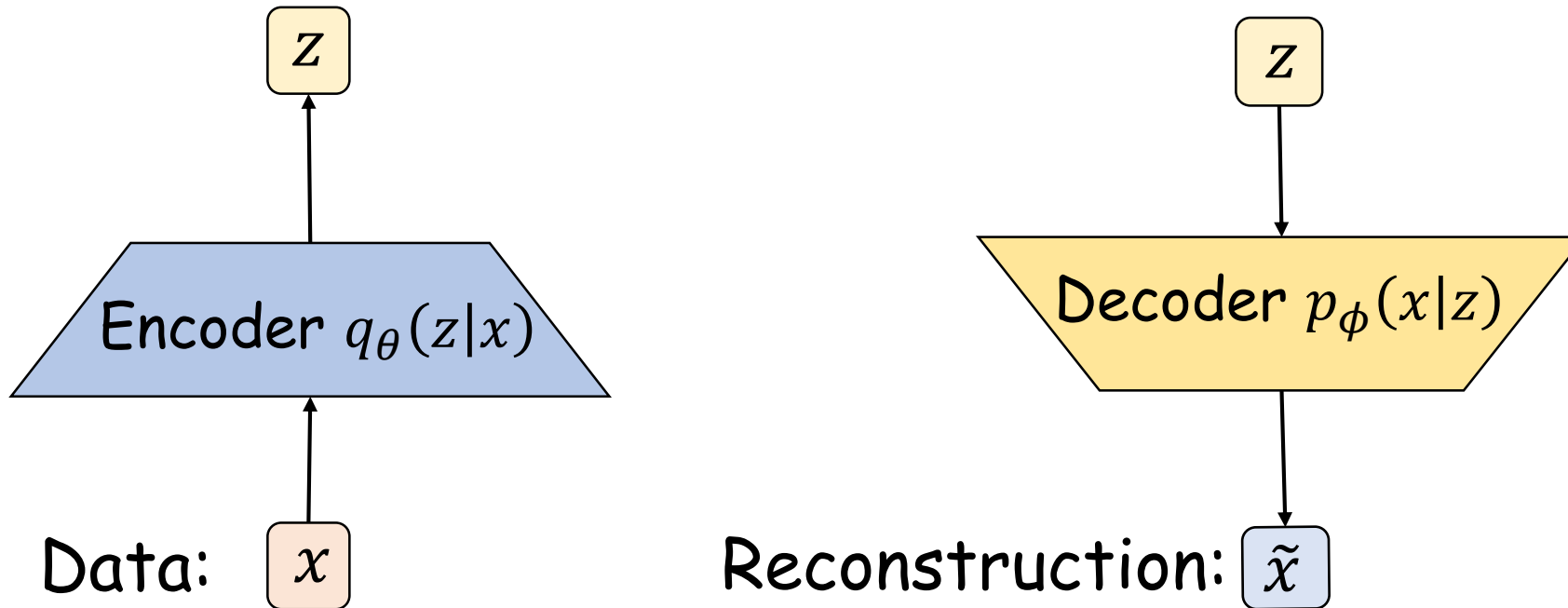
# Outline

---

- Variational Auto-Encoder (VAE)
  - A Neural Net Perspective
  - A Probability Model Perspective
- Variational Adversarial Active Learning (VAAL)

# VAE - A Neural Net Perspective

## Encoder and Decoder



How much information is lost?

Reconstruction log-likelihood  $\log p_{\phi}(x|z)$

# VAE - A Neural Net Perspective

Loss function: the negative log-likelihood with a regularizer

reconstruction loss

regularizer

$$l_i(\theta, \phi) = \left[ -\mathbb{E}_{z \sim q_\theta(z|x_i)} [\log p_\phi(x_i | z)] \right] + \left[ \text{KL}(q_\theta(z | x_i) || p(z)) \right]$$

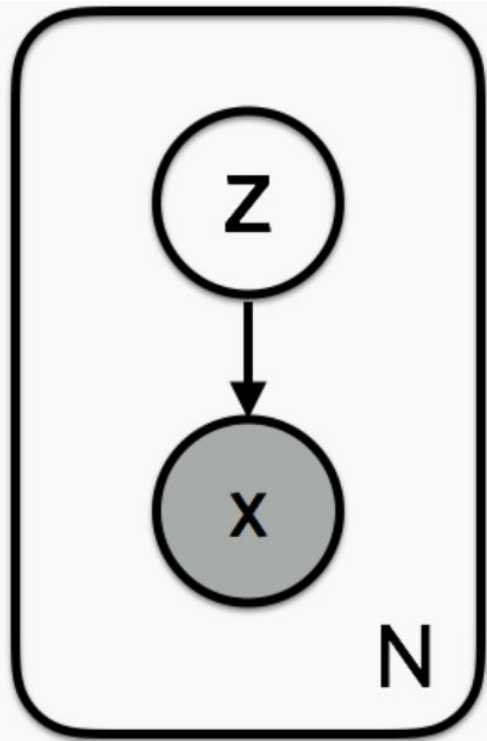
Encourages the decoder parameterizes a likelihood distribution that place much probability mass on the true data

Measures how much information is lost when using  $q$  to represent  $p$

The total loss:  $\sum_{i=1}^N l_i$

# VAE - A Probability Model Perspective

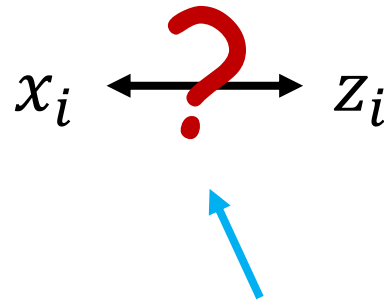
## The generative process



PRML. Chapter 8

For each datapoint  $i$ :

- Draw latent variable  $z_i \sim p(z)$
- Draw datapoint  $x_i \sim p(x|z)$



Inference:  $p(z|x)$



Infer good values of the latent variables  
given observed data

# VAE - A Probability Model Perspective

How to calculate  $p(z|x)$  ?

Bayes tells us: 
$$p(z | x) = \frac{p(x | z)p(z)}{p(x)}$$

Evidence 
$$p(x) = \int p(x | z)p(z)dz$$

requires exp time to compute

## Variational Inference

Use a family of distribution  $q_\lambda(z|x)$  to approximate the true posterior  $p(z|x)$

$$q_\lambda^*(z | x) = \arg \min_\lambda \text{KL}(q_\lambda(z | x) || p(z | x))$$

where

$$\text{KL}(q_\lambda(z | x) || p(z | x)) = \mathbf{E}_q[\log q_\lambda(z | x)] - \mathbf{E}_q[\log p(x, z)] + \log p(x)$$



# VAE - A Probability Model Perspective

The connection to neural net language

$q_\lambda(z | x_i)$   $\xrightarrow{\text{parameterize}}$  Encoder with parameters  $\theta$

$p(x | z)$   $\xrightarrow{\text{parameterize}}$  Decoder with parameters  $\phi$

$$ELBO_i(\theta, \phi) = -l_i(\theta, \phi)$$

$$ELBO_i(\theta, \phi) = \underbrace{\mathbb{E}_{q_\theta(z | x_i)}[\log p_\phi(x_i | z)]}_{\text{reconstruction loss}} - \underbrace{\mathbb{KL}(q_\theta(z | x_i) || p(z))}_{\text{regularizer}}$$

The probability model makes clear why these terms exist:  
to minimize the KL divergence between the approximate posterior  $q_\lambda(z|x)$  and the model posterior  $p(z|x)$ .



# Variational Adversarial Active Learning

---

Samarth Sinha \*  
University of Toronto

samarth.sinha@mail.utoronto.ca

Sayna Ebrahimi \*  
UC Berkeley

sayna@eecs.berkeley.edu

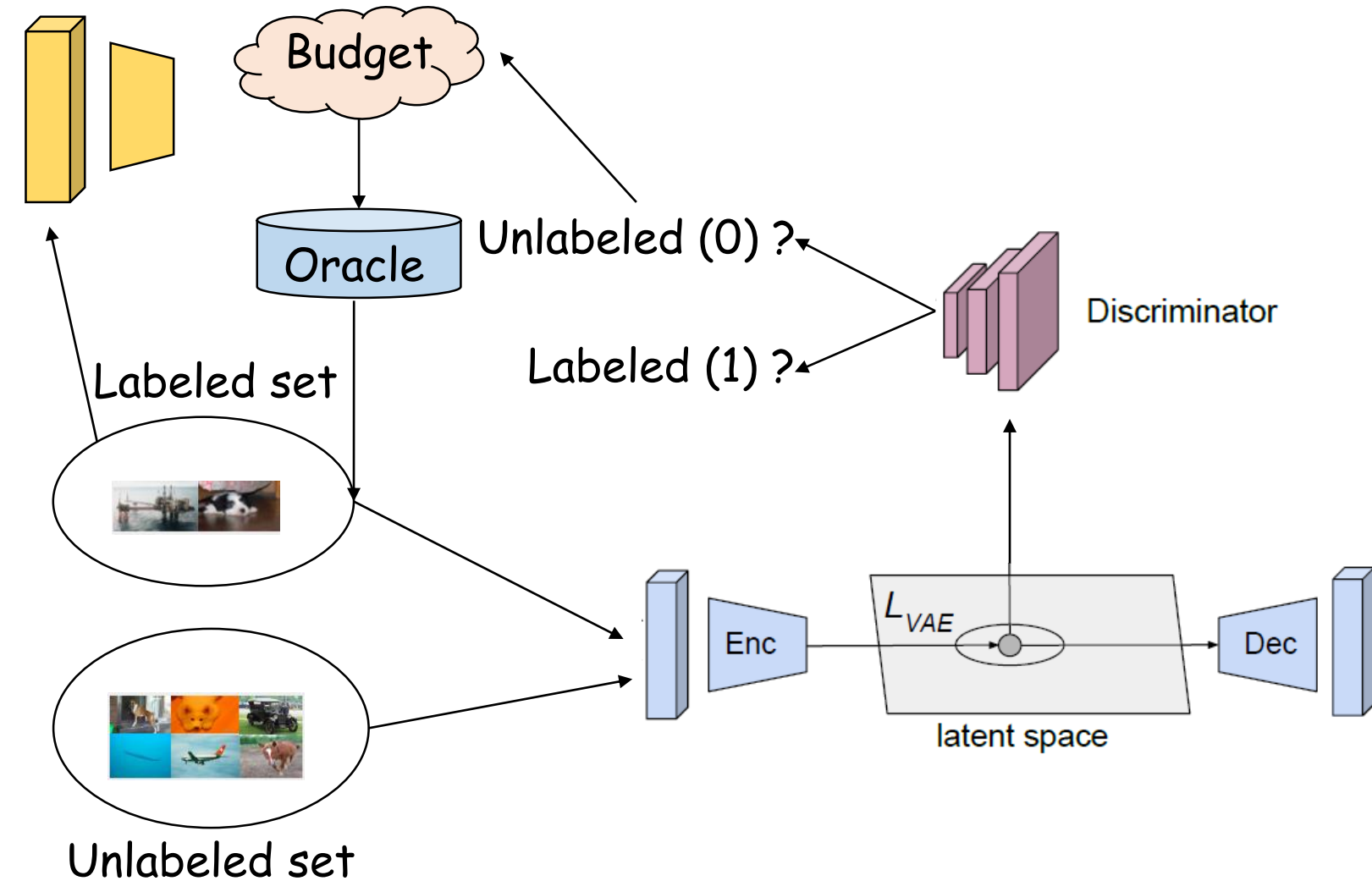
Trevor Darrell  
UC Berkeley

trevor@eecs.berkeley.edu

---

# Idea

Classifier



## Adversarial learning

VAE fools the D to classify all inputs as labeled

**VAE and D are learned together**

D attempts to estimate the probability that the data comes from the unlabeled data

# Variational Adversarial Active Learning (VAAL)

The objective function of the  $\beta$ -VAE

$$\mathcal{L}_{\text{VAE}}^{\text{trd}} = \mathbb{E}[\log p_{\theta}(x_L|z_L)] - \beta \text{D}_{\text{KL}}(q_{\phi}(z_L|x_L)||p(z)) \\ + \mathbb{E}[\log p_{\theta}(x_U|z_U)] - \beta \text{D}_{\text{KL}}(q_{\phi}(z_U|x_U)||p(z))$$

The objective function for the adversarial role of the VAE

$$\mathcal{L}_{\text{VAE}}^{\text{adv}} = \mathcal{L}_{\text{BCE}}(q_{\phi}(z_L|x_L), \mathbb{1}) + \mathcal{L}_{\text{BCE}}(q_{\phi}(z_U|x_U), \mathbb{1})$$

The full objective function for the VAE:  $\mathcal{L}_{\text{VAE}} = \lambda_1 \mathcal{L}_{\text{VAE}}^{\text{trd}} + \lambda_2 \mathcal{L}_{\text{VAE}}^{\text{adv}}$

The objective to train the discriminator

$$\mathcal{L}_{\text{D}} = \mathcal{L}_{\text{BCE}}(q_{\phi}(z_L|x_L), \mathbb{1}) + \mathcal{L}_{\text{BCE}}(q_{\phi}(z_U|x_U), \mathbb{0})$$

# Variational Adversarial Active Learning (VAAL)

---

## Algorithm 1 Variational Adversarial Active Learning

---

**Input:** Labeled pool  $(X_L, Y_L)$ , Unlabeled pool  $(X_U)$ , Initialized models for  $\theta_T$ ,  $\theta_{VAE}$ , and  $\theta_D$

**Input:** Hyperparameters: epochs,  $\lambda_1$ ,  $\lambda_2$ ,  $\alpha_1$ ,  $\alpha_2$ ,  $\alpha_3$

- 1: **for**  $e = 1$  to epochs **do**
- 2: sample  $(x_L, y_L) \sim (X_L, Y_L)$  mini-batch sampling
- 3: sample  $x_U \sim X_U$
- 4: Compute  $\mathcal{L}_{VAE}^{trd}$  by using Eq. 1
- 5: Compute  $\mathcal{L}_{VAE}^{adv}$  by using Eq. 2
- 6:  $\mathcal{L}_{VAE} \leftarrow \lambda_1 \mathcal{L}_{VAE}^{trd} + \lambda_2 \mathcal{L}_{VAE}^{adv}$
- 7: Update VAE by descending stochastic gradients:
- 8:  $\theta'_{VAE} \leftarrow \theta_{VAE} - \alpha_1 \nabla \mathcal{L}_{VAE}$
- 9:  $\mathcal{L}_D \leftarrow \mathcal{L}_{BCE}(q_\phi(z_L|x_L), \mathbb{1}) + \mathcal{L}_{BCE}(q_\phi(z_U|x_U), \mathbb{0})$
- 10: Update  $D$  by ascending its stochastic gradient:
- 11:  $\theta'_D \leftarrow \theta_D - \alpha_2 \nabla \mathcal{L}_D$
- 12: Train and update  $T$ : update classifier
- 13:  $\theta'_T \leftarrow \theta_T - \alpha_3 \nabla \mathcal{L}_T$
- 14: **end for**
- 15: **return** Trained  $\theta_T, \theta_{VAE}, \theta_D$

---

---

## Algorithm 2 Sampling Strategy in VAAL

---

**Input:**  $b, X_L, X_U$

**Output:**  $X_L, X_U$

- 1: Select samples  $(X_s)$  with  $\min_b \{\theta_D(z_U)\}$
  - 2:  $y_o \leftarrow \text{ORACLE}(X_s)$
  - 3:  $X_L \leftarrow X_L \cup (X_s, y_o)$
  - 4:  $X_U \leftarrow X_U - X_s$
  - 5: **return**  $X_L, X_U$
- 

VAE training



Dis training

# Experiments

---

## Datasets

- Image Classification
  - CIFAR10 and CIFAR100
  - Caltech-256
  - ImageNet
- Semantic Segmentation
  - BDD100K
  - Cityscapes

# Experiments

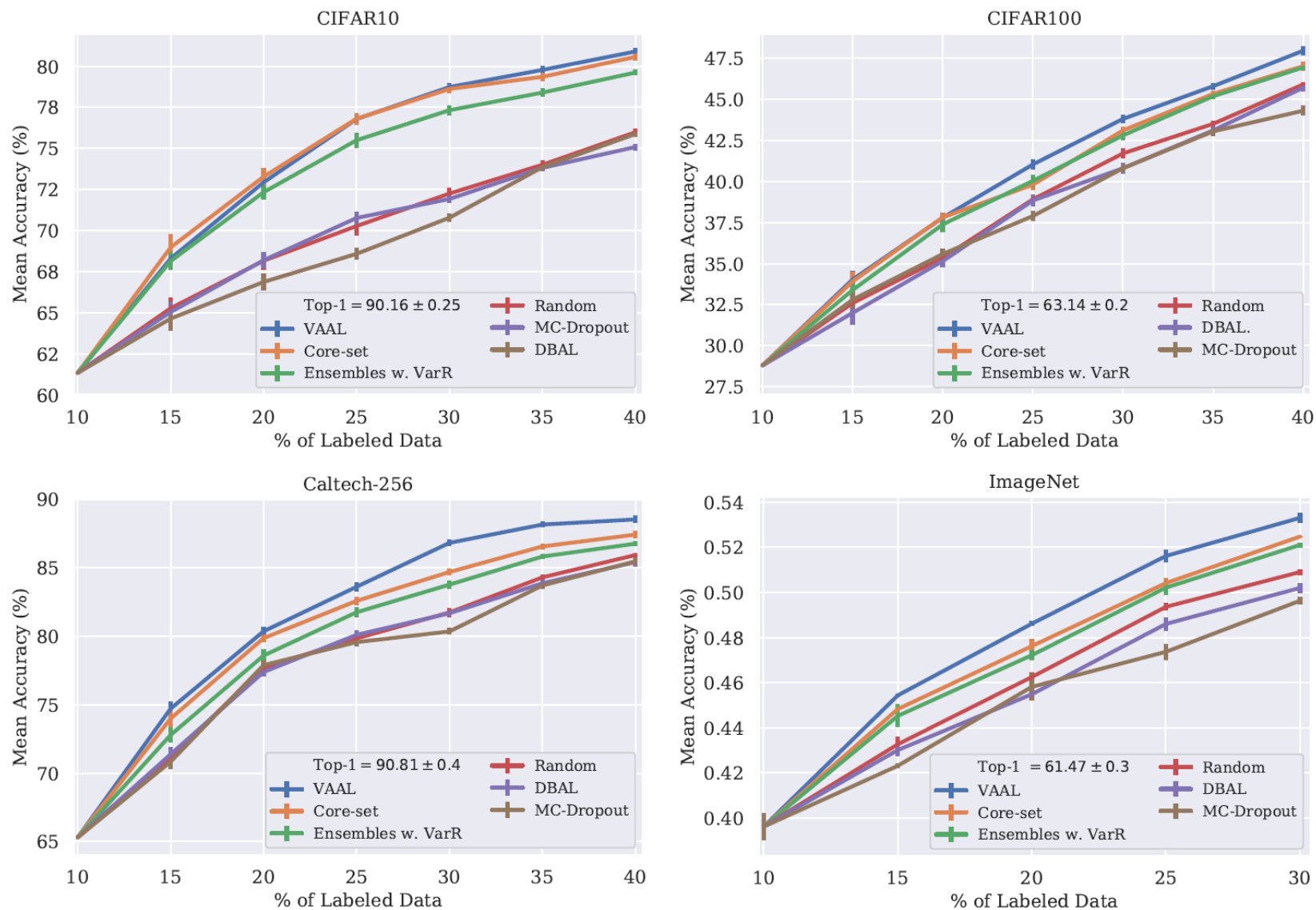
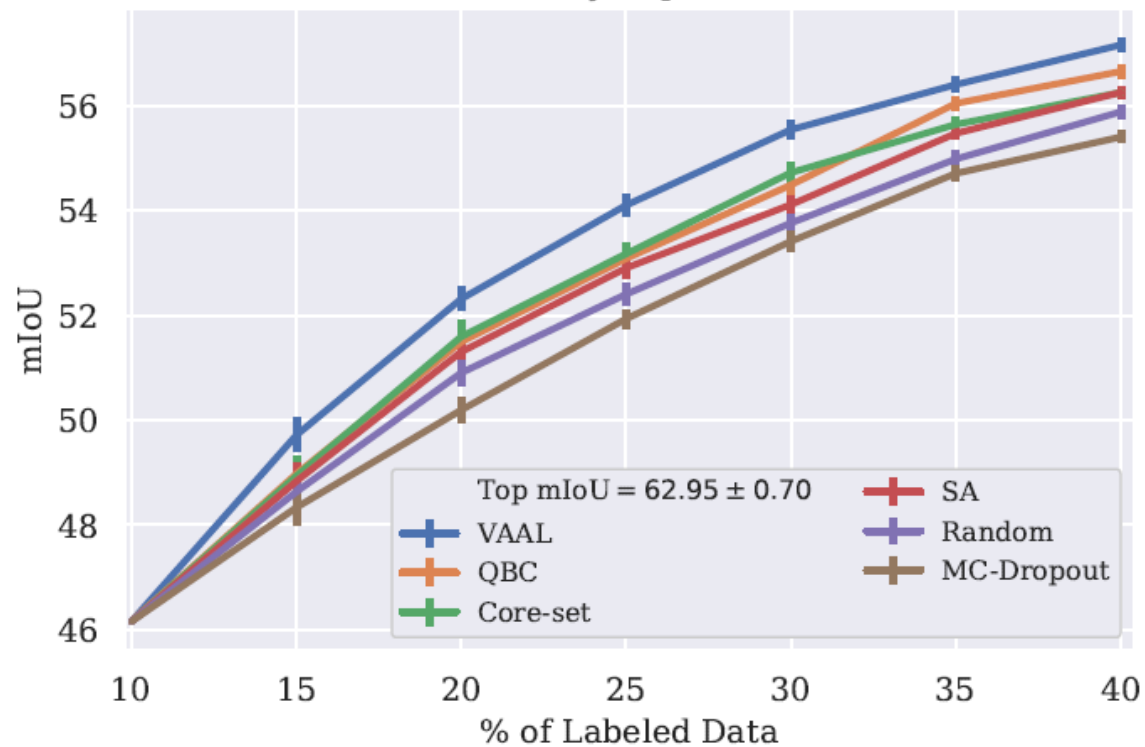


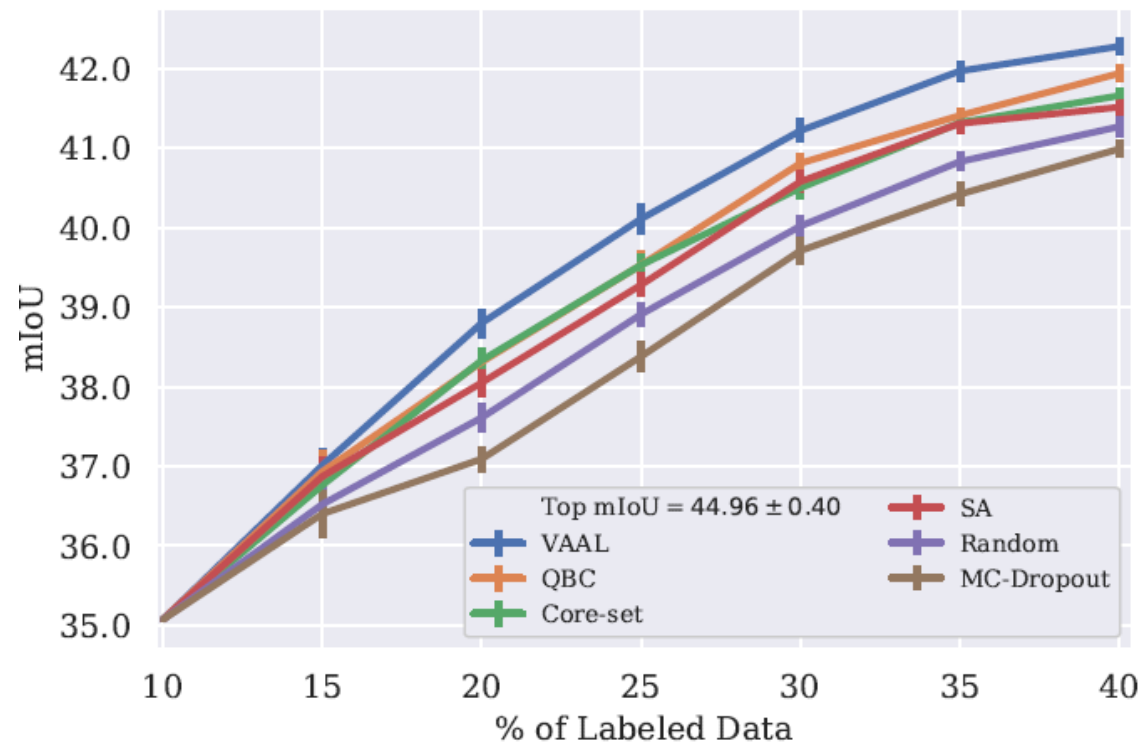
Figure 2. VAAL performance on classification tasks using CIFAR10, CIFAR100, Caltech-256, and ImageNet compared to Core-set [40], Ensembles w. VarR [1], MC-Dropout [13], DBAL [14], and Random Sampling.

# Experiments

Cityscapes

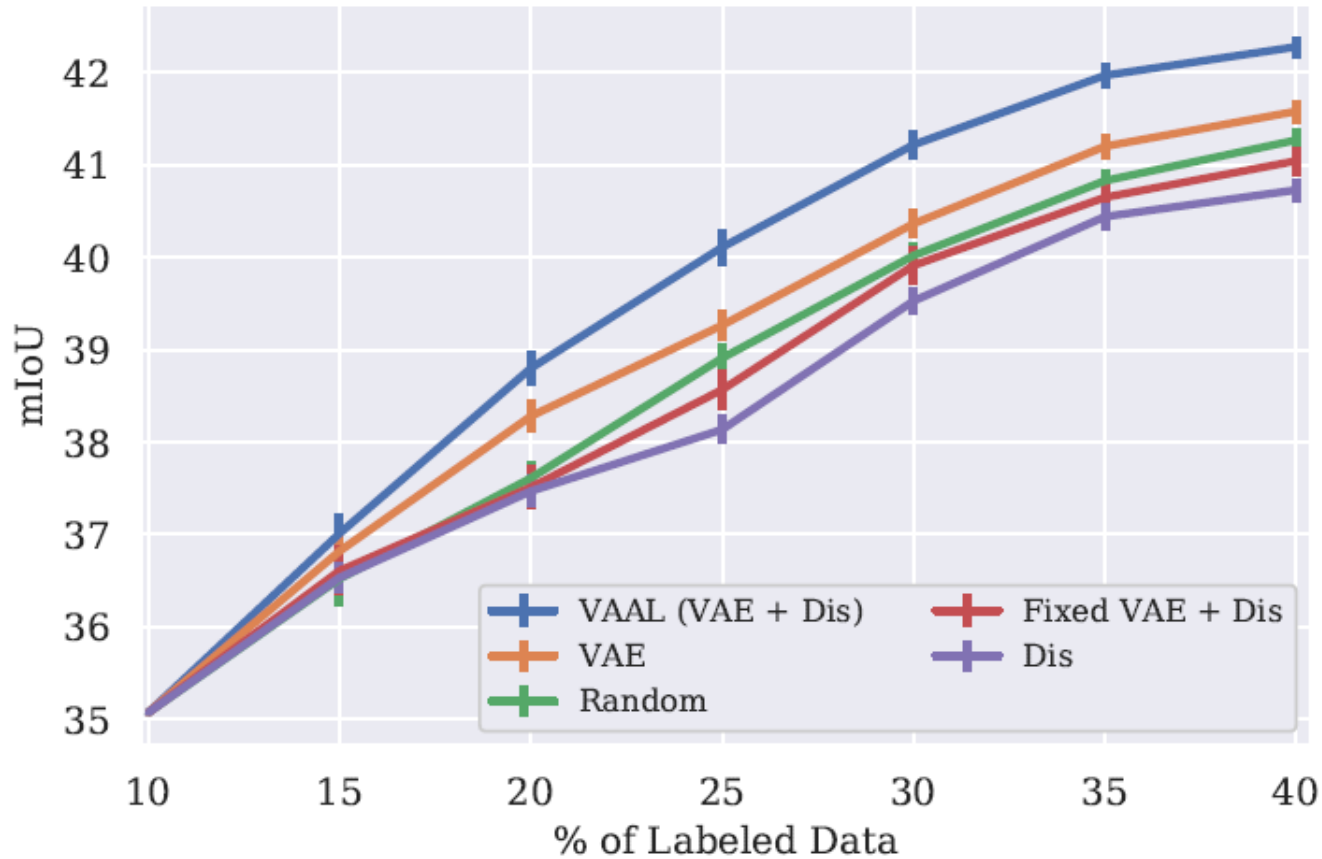


BDD100K



# Experiments

Ablation on BDD100K



Dis: eliminate VAE

Fixed VAE + Dis: frozen VAE with D

VAE: eliminate D

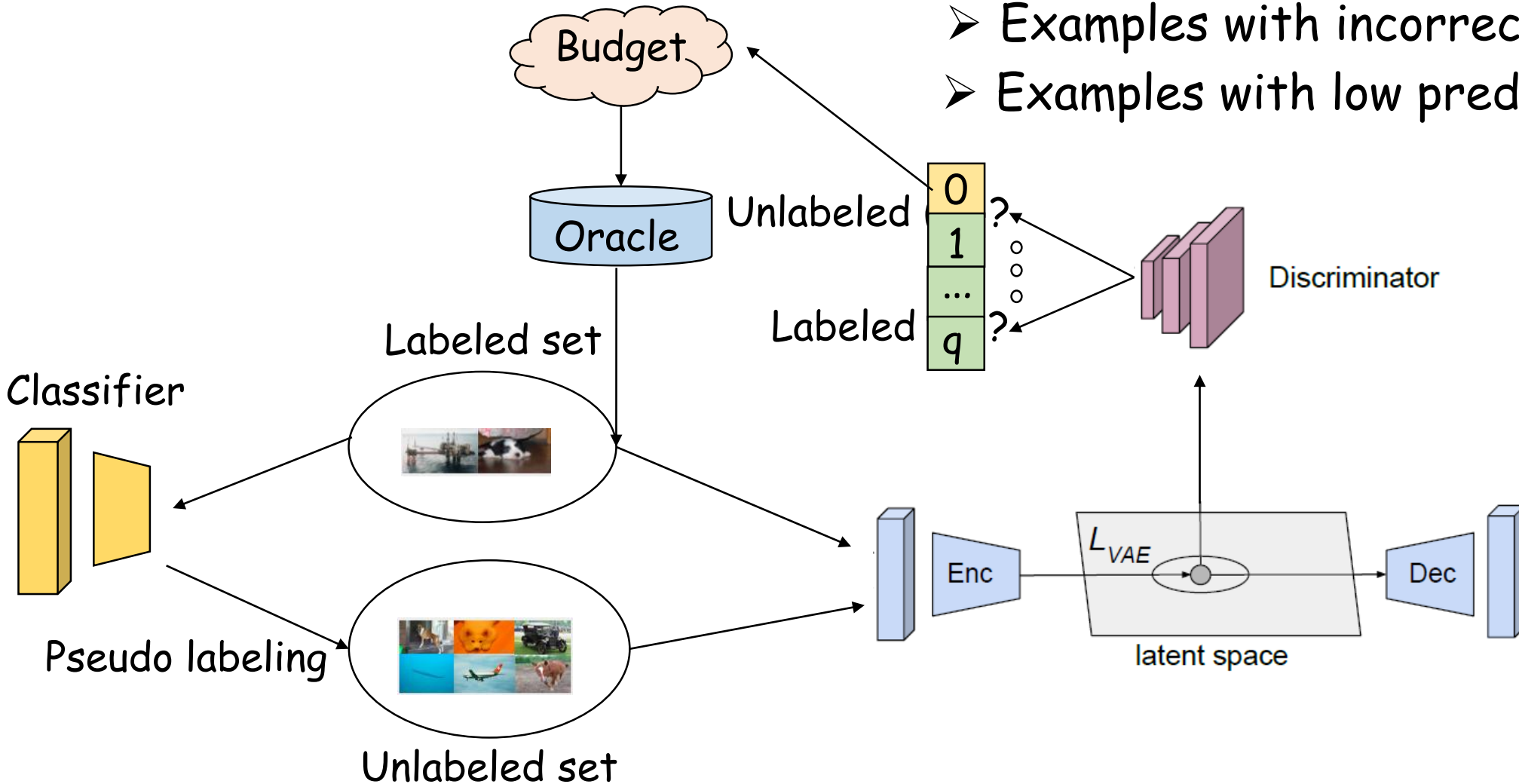
Use 2-Wasserstein distance from the cluster-centroid of the labeled dataset as a heuristic to explicitly measure uncertainty.

# Improvement

## Uncertainty

Two kinds of examples will be identified:

- Examples with incorrect pseudo labels
- Examples with low prediction confidence



Thanks

---