# Modeling Multimodal Aleatoric Uncertainty in Segmentation with Mixture of Stochastic Experts

Zhitong Gao, Yucong Chen, Chuyu Zhang, Xuming He

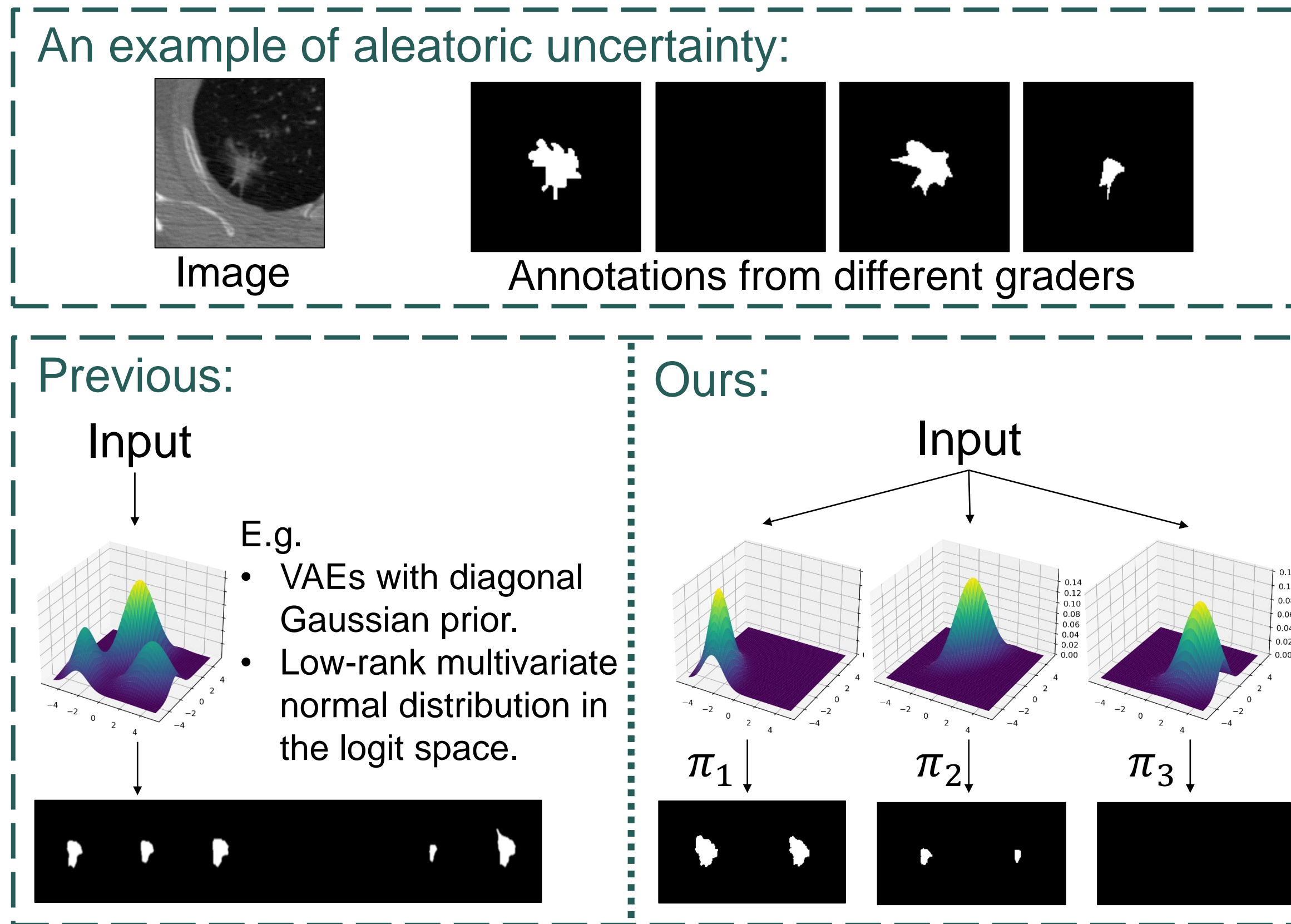ShanghaiTech University

Paper

Code

## INTRODUCTION

- **Problem:** Images are often ambiguous, leading to multiple plausible ground truth segmentation results.

- **Goal:** We aim to capture this data-inherent uncertainty (aka Aleatoric uncertainty) by learning the latent segmentation distribution.

- **Motivation:** The segmentation distribution is typically multi-modal. However, most previous methods have restricted capacity in capturing multi-modality, and rely on inefficient sampling to represent the predictive distribution.

- **Main idea:** We propose to explicitly model the multimodal characteristics of the distribution and provide a more efficient representation of the uncertainty.
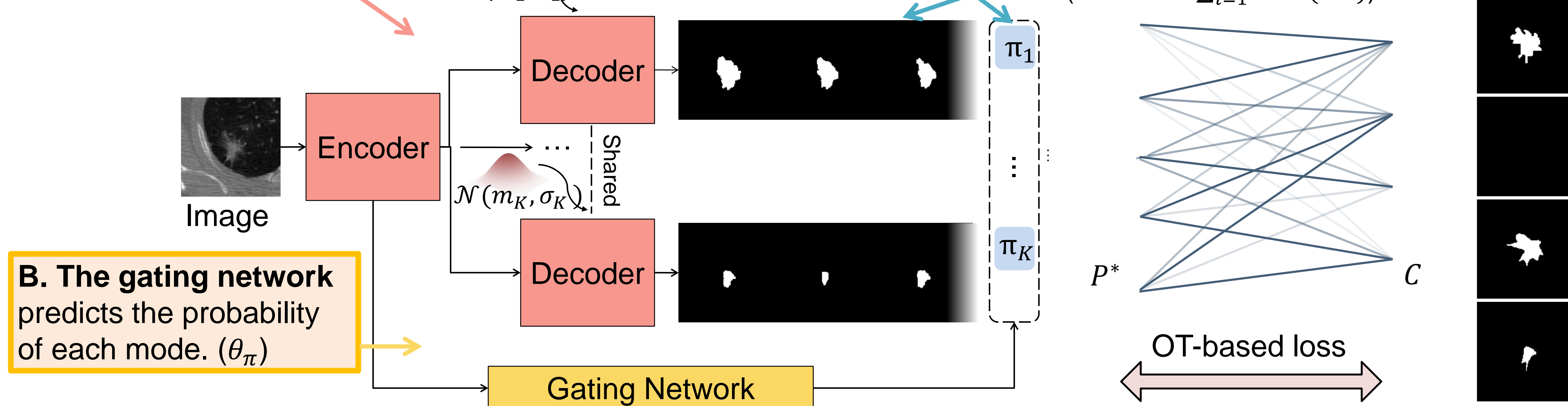
An example of aleatoric uncertainty:

Image   Annotations from different graders

Previous:

Input

E.g.
- VAEs with diagonal Gaussian prior.
- Low-rank multivariate normal distribution in the logit space.

Ours:

Input

$\pi_1$   $\pi_2$   $\pi_3$

## METHOD

- **Overview:** We introduce a conditional probabilistic model on segmentation $\mu_\theta(y|x)$, and design a distributional loss $\ell$ to represent and learn the aleatoric uncertainty.

$$\min_\theta \sum_n \ell(\mu_\theta(y|x_n), \nu_n)$$

where $\nu_n := \sum_{i=1}^M v_n^{(i)} \delta(y_n^{(i)})$ is the empirical GT distribution ($M \geq 1$).

### 1. Mixture of Stochastic Experts (MoSE) Framework

**A. Each expert** encodes a distinct mode of aleatoric uncertainty.($\theta_s$)

**C. Compact uncertainty representation**
$$\hat{\mu}_\theta(y|x) = \sum_{k=1}^K \frac{\pi_k(x)}{S} \sum_{i=1}^S \delta\left(s_k^{(i)}\right)$$
(rewrite as $\sum_{i=1}^N u^{(i)}\delta(s^{(i)})$)

$\mathcal{N}(m_1, \sigma_1)$

Decoder

$\pi_1$

$\mathcal{N}(m_K, \sigma_K)$   Shared

Image   Encoder

Decoder

$\pi_K$

$P^*$   $C$

GT

OT-based loss

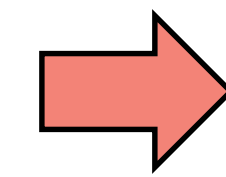**B. The gating network** predicts the probability of each mode. ($\theta_\pi$)

Gating Network

### 2. Optimal-Transport-Based (OT) Loss

- Formulate the model learning as an OT problem.

$$\min_{\theta_s,\theta_\pi} \sum_n \langle P^*, C_n(\theta_s)\rangle$$
$$P^* = \underset{P\in R_+}{\operatorname{argmin}}\langle P, C_n\rangle \ s.t. P1_M = u_n(\theta_\pi); \ P^T 1_N = v_n$$

- Efficient bi-level optimization with constraint relaxation.

$$\min_{\theta_s,\theta_\pi} \sum_n \langle P^*, C_n(\theta_s)\rangle + \beta\, KL(P^*1_M || u_n(\theta_\pi))$$
$$P^* = \underset{P\in R_+}{\operatorname{argmin}}\langle P, C_n\rangle \ s.t. P1_M \leq \gamma; \ P^T 1_N = v_n$$
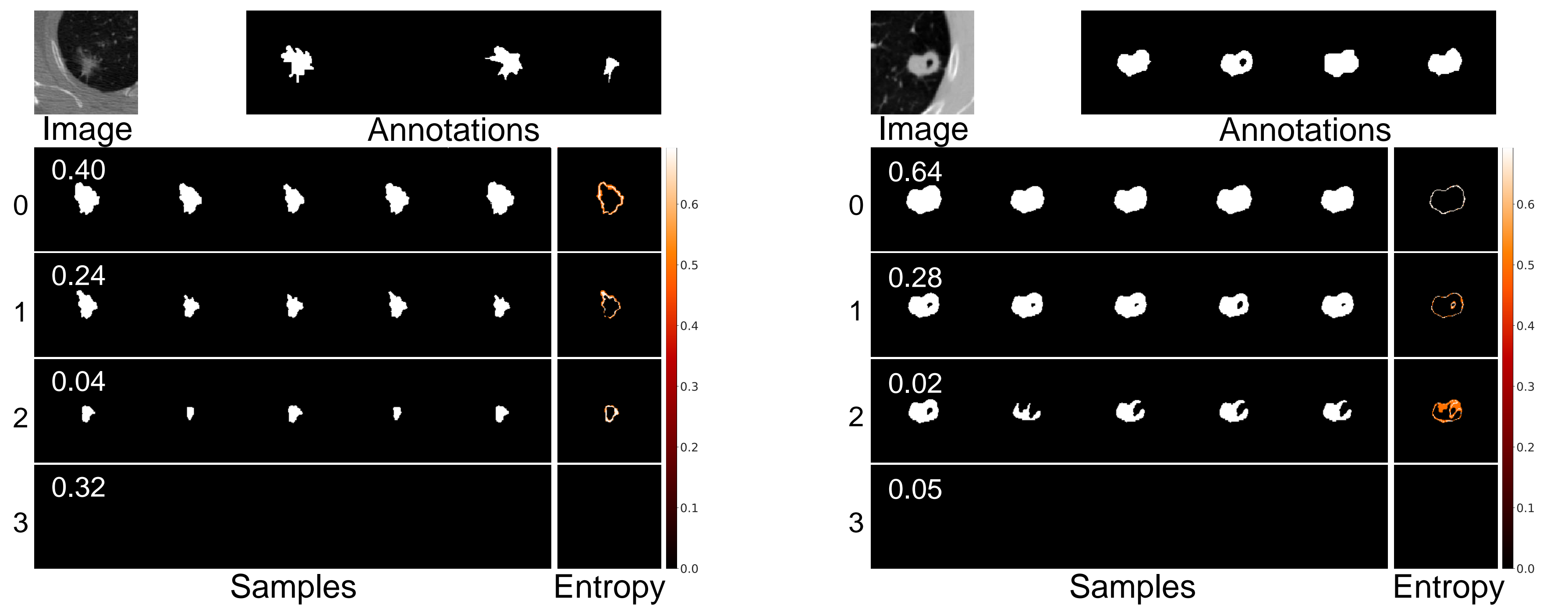
Annealing to 1.

## RESULTS

### 1. Results on the LIDC dataset

- Compare with previous SOTA models, (.) denotes number of sampled outputs.

| Method | # label | GED ↓ (16) | GED ↓ (100) | M-IoU ↑ (16) | ECE ↓ (%) (16) | # param. |
|---|---|---|---|---|---|---|
| Kohl et al. (2018) | | $0.320 \pm 0.030$ | $0.239 \pm$ N/A [†] | $0.500 \pm 0.030$ | - | 76.15M |
| Kohl et al. (2019) | | $0.270 \pm 0.010$ | - | $0.530 \pm 0.010$ | - | 87.51M |
| Baumgartner et al. (2019) | | - | $0.224 \pm$ N/A | - | - | 74.82M |
| Monteiro et al. (2020) | All | - | $0.225 \pm 0.002$ | - | - | 41.28M |
| Kassapis et al. (2021) | | $0.264 \pm 0.002$ | $0.243 \pm 0.004$ | $0.592 \pm 0.005$ | $0.214$ * | 175.36M |
| **Ours** | | $\mathbf{0.218 \pm 0.003}$ | $\mathbf{0.189 \pm 0.002}$ | $\mathbf{0.624 \pm 0.004}$ | $\mathbf{0.064 \pm 0.015}$ | 41.60M |
| **Ours - compact** | | $\mathbf{0.195 \pm 0.005}$ | $\mathbf{0.186 \pm 0.002}$ | $\mathbf{0.635 \pm 0.003}$ | $\mathbf{0.054 \pm 0.015}$ | 41.60M |
| Kohl et al. (2018) | | - | $0.445 \pm$ N/A [†] | - | - | 76.15M |
| Baumgartner et al. (2019) | | - | $0.323 \pm$ N/A | - | - | 74.82M |
| Monteiro et al. (2020) | One | - | $0.365 \pm 0.005$ | - | - | 41.28M |
| **Ours** | | $\mathbf{0.252 \pm 0.004}$ | $\mathbf{0.223 \pm 0.005}$ | $\mathbf{0.596 \pm 0.003}$ | $\mathbf{0.105 \pm 0.009}$ | 41.60M |
| **Ours - compact** | | $\mathbf{0.228 \pm 0.004}$ | $\mathbf{0.220 \pm 0.005}$ | $\mathbf{0.605 \pm 0.003}$ | $\mathbf{0.090 \pm 0.011}$ | 41.60M |

Image   Annotations

Image   Annotations

Samples   Entropy

Samples   Entropy

### 2. Ablation study

- Evaluate the impact of each component on the full-labeled LIDC dataset. (16 samples)

| Expert type | Expert weights | loss | GED ↓ | M-IoU ↑ | ECE ↓ (%) |
|---|---|---|---|---|---|
| stochastic | learnable / uniform | IoU loss | $0.533 \pm 0.001$ | $0.533 \pm 0.001$ | $0.277 \pm 0.017$ |
| stochastic | uniform | OT-based | $0.282 \pm 0.002$ | $0.545 \pm 0.007$ | $0.215 \pm 0.006$ |
| deterministic | learnable | OT-based | $0.246 \pm 0.006$ | $0.591 \pm 0.001$ | $0.142 \pm 0.003$ |
| stochastic | learnable | OT-based | $\mathbf{0.218 \pm 0.003}$ | $\mathbf{0.624 \pm 0.004}$ | $\mathbf{0.064 \pm 0.015}$ |

### 3. Results on the synthetic multimodal Cityscapes dataset

- Constructed by randomly flipping five classes with certain probabilities (GT distribution known).
- Quantitatively, our model achieves the SOTA or comparable performance on three metrics. (Please refer to our paper for more detailed information.)

GTs

Preds