



# Student-t Variational Autoencoder for Robust Density Estimation

Hiroshi Takahashi<sup>1</sup>, Tomoharu Iwata<sup>2</sup>,  
Yuki Yamanaka<sup>3</sup>, Masanori Yamada<sup>3</sup>,  
Satoshi Yagi<sup>1</sup>

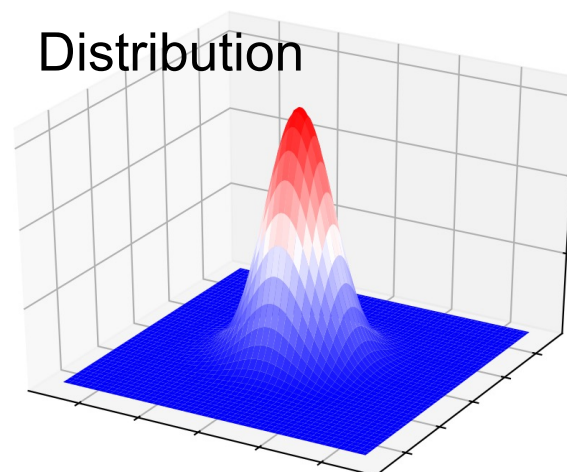
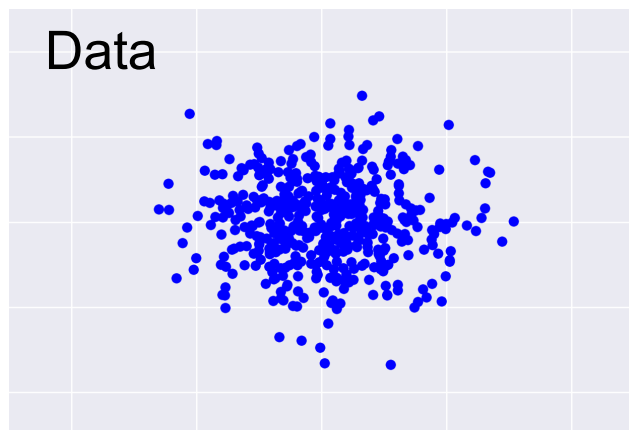
<sup>1</sup>NTT Software Innovation Center

<sup>2</sup>NTT Communication Science Laboratories

<sup>3</sup>NTT Secure Platform Laboratories

**If you use the VAE for continuous data,  
we recommend using the Student-t  
distribution as the decoder!**

- Estimating data distributions is important for AI
  - especially for image, audio, video, and detection tasks



- The VAE is widely used since it can learn the high-dimensional complicated distributions in these tasks
- We focus on estimating distributions of **continuous data** with the VAE

- The VAE estimates the probability of a continuous data point  $\mathbf{x}$  by using latent variable  $\mathbf{z}$ :

$$p_{\theta}(\mathbf{x}) = \int \underbrace{p_{\theta}(\mathbf{x} | \mathbf{z})}_{\text{decoder}} \underbrace{p(\mathbf{z})}_{\text{prior}} d\mathbf{z}$$

- The log marginal likelihood of VAE is bounded below by the evidence lower bound (ELBO):

$$\begin{aligned} \ln p_{\theta}(\mathbf{x}) \\ \geq \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p_{\theta}(\mathbf{x} | \mathbf{z})]}_{\text{encoder}} - D_{KL}(\underbrace{q_{\phi}(\mathbf{z} | \mathbf{x})}_{\text{encoder}} || p(\mathbf{z})) \end{aligned}$$

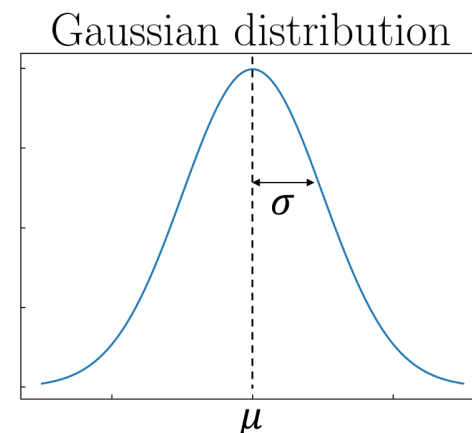
- The VAE is trained to maximize the sum of ELBO

- For continuous data, the encoder, decoder, and prior distributions are usually modeled by a Gaussian:

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I}) \quad : \text{standard Gaussian}$$

$$p_{\theta}(\mathbf{x} | \mathbf{z}) = \mathcal{N}(\mathbf{x} | \underbrace{\mu_{\theta}(\mathbf{z})}_{\text{estimated by neural networks with parameter } \theta}, \underbrace{\sigma_{\theta}^2(\mathbf{z})}_{\text{estimated by neural networks with parameter } \theta})$$

**estimated by neural networks with parameter  $\theta$**



$$q_{\phi}(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \underbrace{\mu_{\phi}(\mathbf{x})}_{\text{estimated by neural networks with parameter } \phi}, \underbrace{\sigma_{\phi}^2(\mathbf{x})}_{\text{estimated by neural networks with parameter } \phi})$$

**estimated by neural networks with parameter  $\phi$**

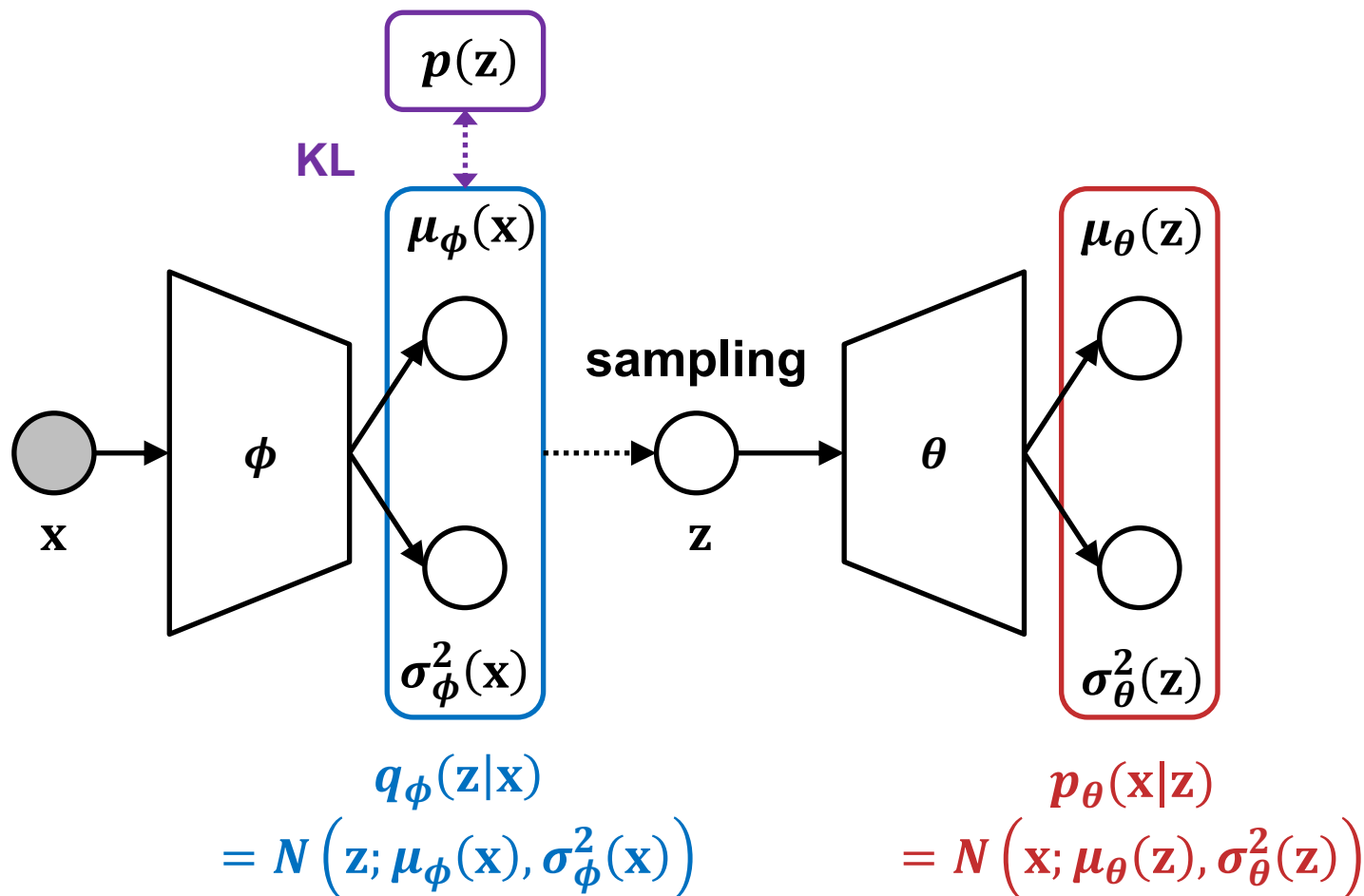
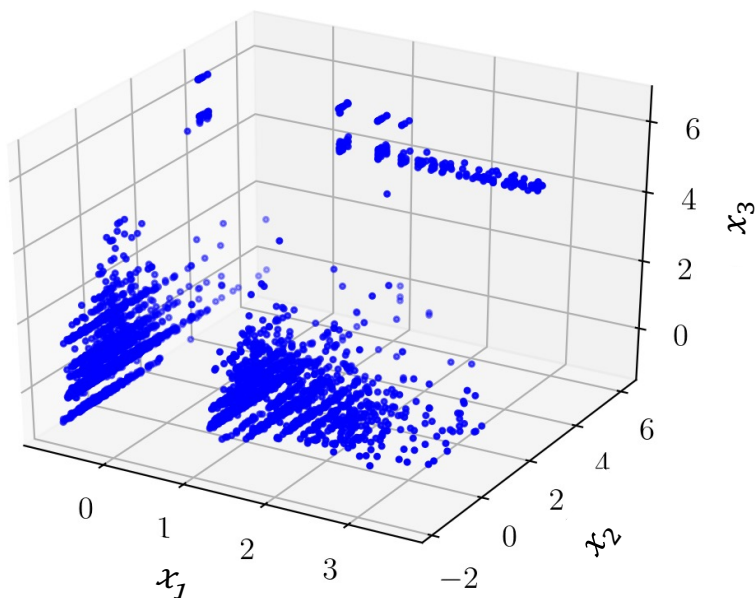


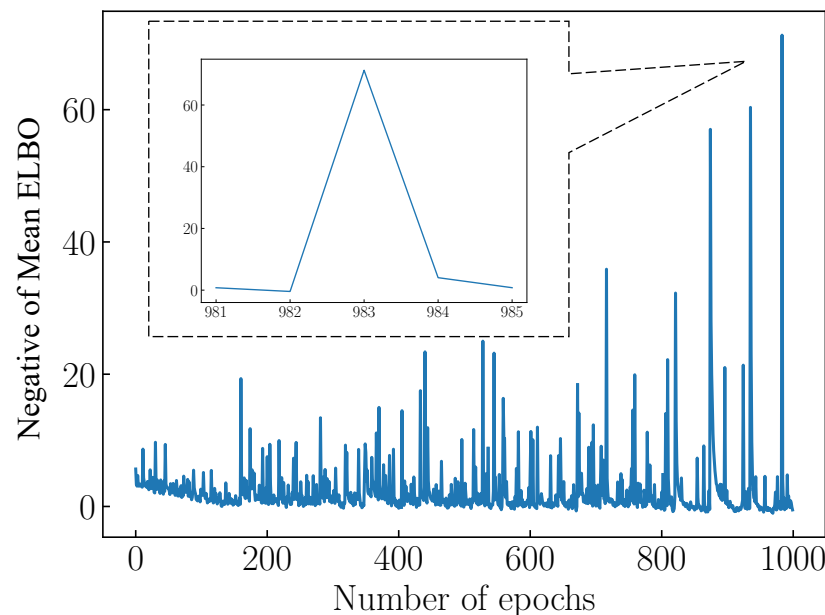
Diagram of VAE for continuous data

- When we use the Gaussian as the decoder, the training of VAE often becomes unstable
  - For example, when we train KDD99 SMTP with VAE, Negative of Mean ELBO sharply jumped up during training



KDD99 SMTP Dataset

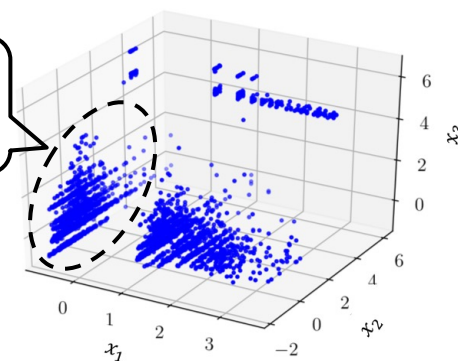
$$\mathbf{x} = (x_1, x_2, x_3)^T$$



Negative of Mean ELBO

- The cause is division by too small variance in ELBO

the variance of these data points is too small along  $x_1$  direction



## First term of ELBO

$$\ln p_{\theta}(\mathbf{x} | \mathbf{z}) = \ln \mathcal{N}(\mathbf{x} | \mu_{\theta}(\mathbf{z}), \sigma_{\theta}^2(\mathbf{z}))$$
$$= \sum_d \left[ \underbrace{-\frac{(\mathbf{x}_d - \mu_{\theta,d}(\mathbf{z}))^2}{2\sigma_{\theta,d}^2(\mathbf{z})}}_{\text{variance}} - \frac{1}{2} \ln 2\pi\sigma_{\theta,d}^2(\mathbf{z}) \right]$$

$d$  : dimension index

When the decoded variance  $\sigma_{\theta}^2(\mathbf{z})$  is almost zero, this term is sensitive to the error between  $\mathbf{x}$  and its decoded mean  $\mu_{\theta}(\mathbf{z})$



- We can avoid this instability problem by preventing the decoded variance  $\sigma_{\theta}^2(\mathbf{z})$  from being too small
- To penalize small variance, we introduce a Gamma distribution as the prior for the decoded variance  $\sigma_{\theta}^2(\mathbf{z})$

$$\text{Gam}(\tau \mid a, b) = \frac{b^a \tau^{a-1} \exp(-b\tau)}{\Gamma(a)}$$

$\tau$  : the inverse of the variance ( $1/\sigma^2$ )

- First, we present the MAP estimation for the VAE
  - To simplify the calculation, we use  $\text{Gam}(\tau|1, b)$  as the prior
  - Then, the objective function of MAP estimation is:

$$\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[ \ln p_{\theta}(\mathbf{x} | \mathbf{z}) - \frac{b}{\sigma_{\theta}^2(\mathbf{z})} \right] - D_{KL}(q_{\phi}(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z}))$$

small variance is penalized with regularization parameter  $b$

- However, there are two drawbacks in MAP estimation
  1. Tuning  $b$  is difficult
  2. The constant  $b$  lacks flexibility in density estimation
    - $b$  should depend on a data point

- We propose a more flexible approach by introducing a Gamma prior that depends on latent variables:

$$\text{Gam}(\tau \mid a(\mathbf{z}), b(\mathbf{z}))$$

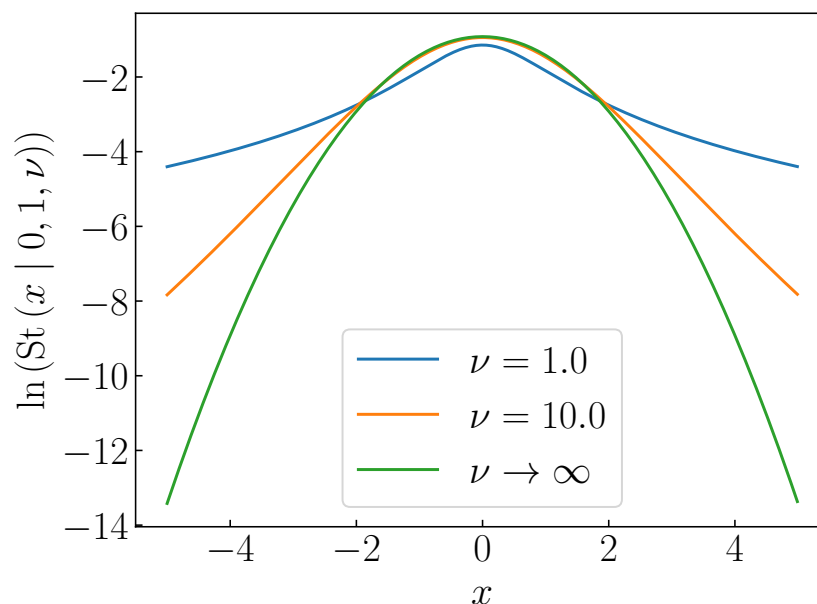
- By analytically integrating out the  $\tau$ , we can obtain a Student-t decoder:

$$\begin{aligned} p_{\theta}(\mathbf{x} \mid \mathbf{z}) &= \int_0^{\infty} \mathcal{N}(\mathbf{x} \mid \mu_{\theta}(\mathbf{z}), \tau^{-1}) \text{Gam}(\tau \mid a(\mathbf{z}), b(\mathbf{z})) d\tau \\ &= \text{St}(\mathbf{x} \mid \mu_{\theta}(\mathbf{z}), \lambda_{\theta}(\mathbf{z}), \nu_{\theta}(\mathbf{z})) \end{aligned}$$

where

$$\lambda_{\theta}(\mathbf{z}) = a(\mathbf{z})/b(\mathbf{z}), \nu_{\theta}(\mathbf{z}) = 2a(\mathbf{z})$$

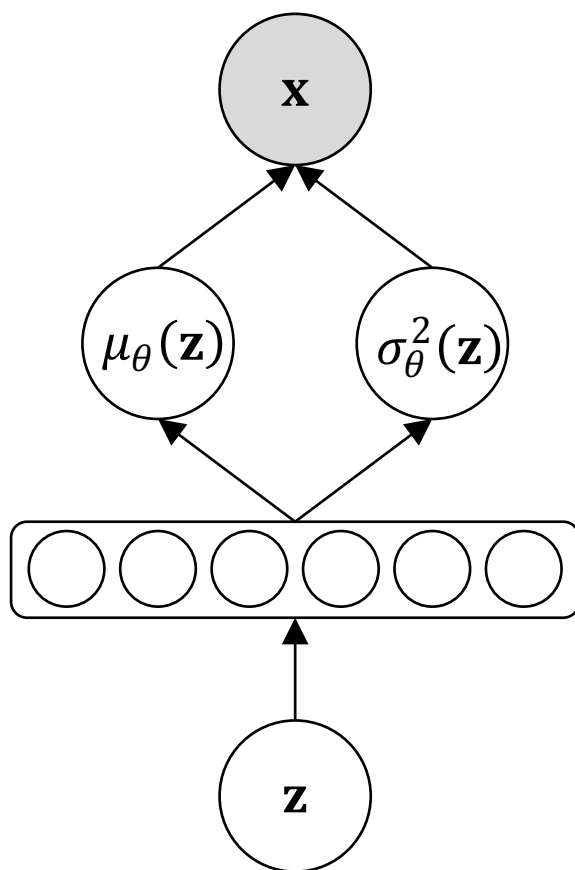
- Since the Student-t distribution is **heavy-tailed** (has large variance), the Student-t decoder is **robust** to the error between the data point and its decoded mean
  - The appropriate robustness is set by  $\nu_{\theta}(\mathbf{z})$ , which makes the training of VAE stable!



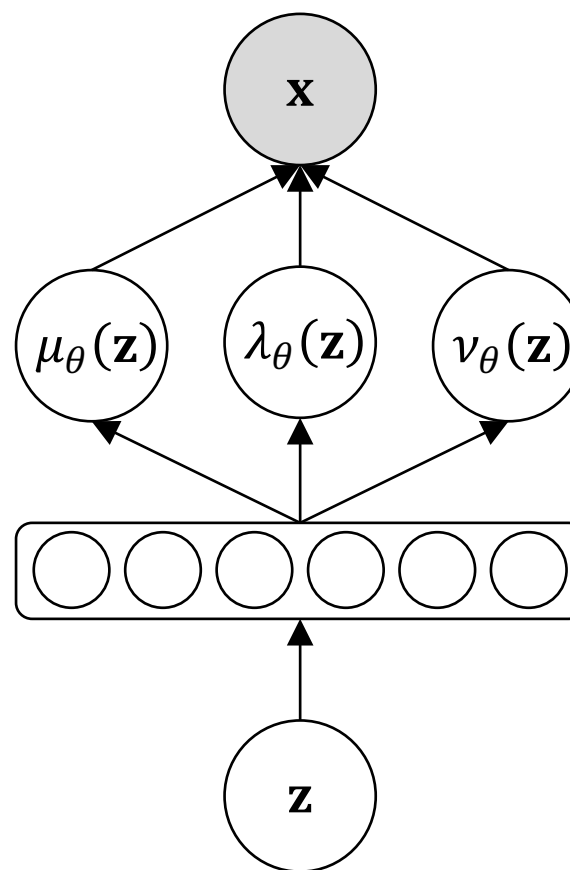
Plot of  $\text{St}(x|0,1,\nu)$  in log scale

# [Proposed method]

## Diagram of Student-t decoder



Gaussian decoder

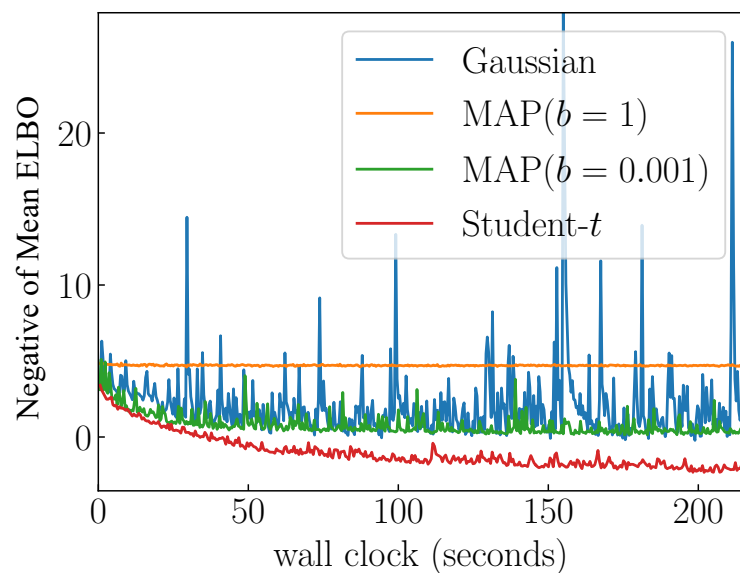


Student-t decoder

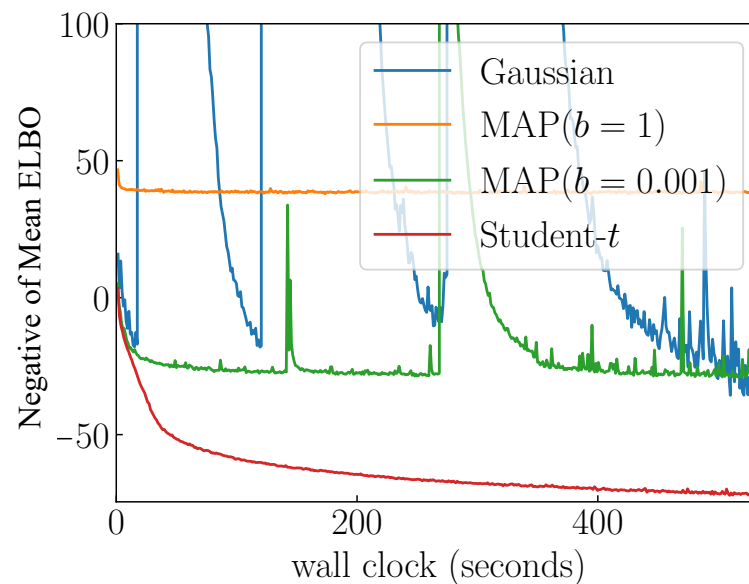
# [Experiments]

## Stability of training

- Our approach reduced the negative ELBO equal to or more stably than other approaches



SMTP



Aloï

Negative of Mean ELBO for each dataset

# [Experiments]

## Test log-likelihoods



- Our approach obtained the equal to or better density estimation performance than that of other approaches.

	Gaussian	MAP( $b = 1$ )	MAP( $b = 0.001$ )	Student- $t$
SMTP	$-1.248 \pm 0.404$	$-4.864 \pm 0.020$	$-1.932 \pm 0.404$	<b><math>0.827 \pm 0.105</math></b>
Aloi	$45.418 \pm 5.457$	$-38.210 \pm 0.156$	$30.406 \pm 0.383$	<b><math>77.022 \pm 0.539</math></b>
Thyroid	$15.519 \pm 4.422$	$-31.266 \pm 0.159$	$18.037 \pm 1.318$	<b><math>69.543 \pm 0.634</math></b>
Cancer	<b><math>-18.668 \pm 3.448</math></b>	$-45.895 \pm 0.843$	<b><math>-19.017 \pm 3.273</math></b>	<b><math>-18.253 \pm 2.629</math></b>
Satellite	<b><math>-1.852 \pm 0.370</math></b>	$-50.895 \pm 0.238$	<b><math>-1.899 \pm 0.372</math></b>	<b><math>-1.811 \pm 0.289</math></b>

### Comparison of test log-likelihoods<sup>1</sup>

<sup>1</sup>We highlighted the best result in bold, and we also highlighted the results in bold which are not statistically different from the best result according to a pair-wise t-test.

## In conclusion

- We proposed the Student-t VAE for robust multivariate density estimation
- We experimentally showed that the stability of the training and the high density estimation performance of the Student-t VAE
- We recommend using the Student-t distribution as the decoder If you use the VAE for continuous data!



# Thank you!



## Thank you for your attention!

If you have any questions,  
email me: [takahashi.hiroshi@lab.ntt.co.jp](mailto:takahashi.hiroshi@lab.ntt.co.jp)

- **Q1: Did you compare this model with GAN?**

- A1: With SMTP dataset, we compared Student-t VAE with Wasserstein GAN, and confirmed that the test log likelihood of the Student-t VAE was better than that of Wasserstein GAN.

- **Q2: What is the limitation of this approach?**

- A2: This approach requires heavier computational cost than Gaussian decoder. (about 1.5 times)

- **Q3: Is this approach useful when the dataset is discrete?**

- A3: If the dataset is binary, we recommend using the Bernoulli distribution as the decoder. Other than that, our approach may be useful.