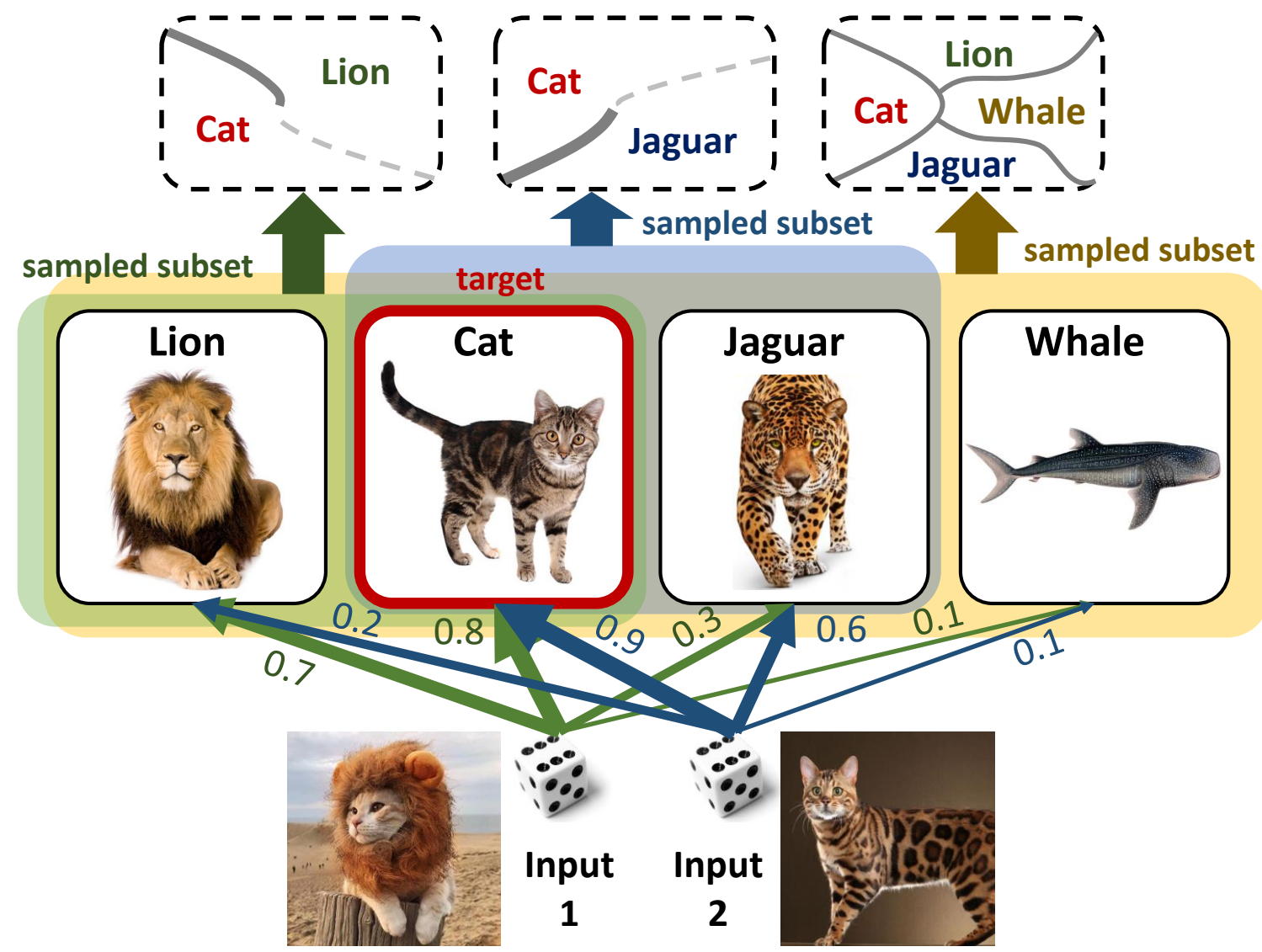


Hae Beom Lee^{1 2} Juho Lee^{3 2}, Saehoon Kim², Eunho Yang^{1 2}, and Sung Ju Hwang^{1 2}
KAIST¹, Altrics², University of Oxford³

Motivation

What if we apply **dropout** to **softmax** function?



- Ensemble learning with exponentially many different **sub-classification** problems.
- Stochastic classifier can consider **confusing classes** more often than others in **input-adaptive** manner.

The goal of this work is to figure out if such way of **noise injection** can improve the **generalization performance** of the classifier.

Approach

Softmax function:

$$p(y_t = 1 | \mathbf{x}) = \frac{\exp(o_t(\mathbf{x}; \theta))}{\sum_k \exp(o_k(\mathbf{x}; \theta))}$$

Define DropMax:

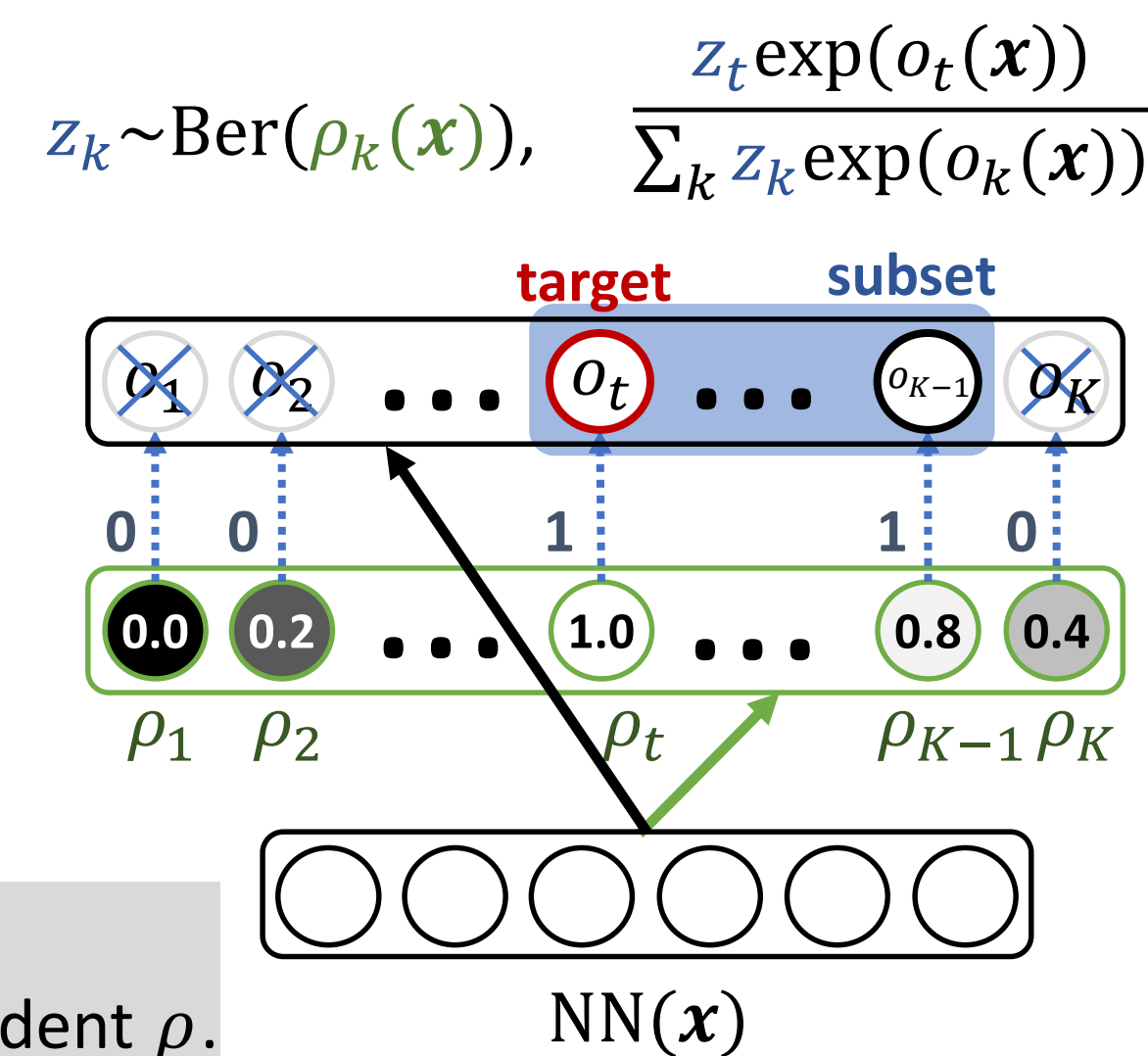
$$z_k | \mathbf{x} \sim \text{Ber}(z_k; \rho_k(\mathbf{x}; \theta))$$

$$p(y_t = 1 | \mathbf{x}, \mathbf{z}; \theta) = \frac{(z_t + \varepsilon) \exp(o_t(\mathbf{x}; \theta))}{\sum_k (z_k + \varepsilon) \exp(o_k(\mathbf{x}; \theta))}$$

- Class k is completely excluded when $z_k = 0$.
- Target class t should not be dropped \rightarrow input dependent ρ .
- The method can be applied to any neural network.

Related to Adaptive dropout [1], but our focus is on the output layer and we train in more principled way using reparameterization trick with concrete relaxation:

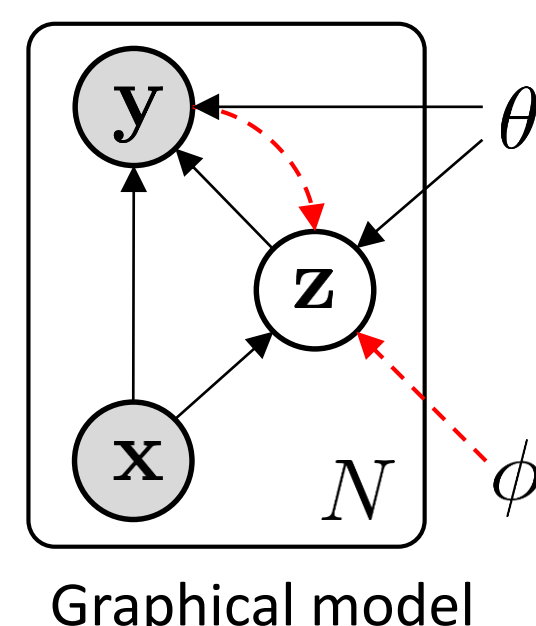
$$z_k = \text{sgm} \left\{ \frac{1}{\tau} \left(\log \frac{\rho_k(\mathbf{x}; \theta)}{1 - \rho_k(\mathbf{x}; \theta)} + \log \frac{u}{1 - u} \right) \right\}, \quad u \sim \mathcal{U}(0, 1)$$



Standard Variational Inference

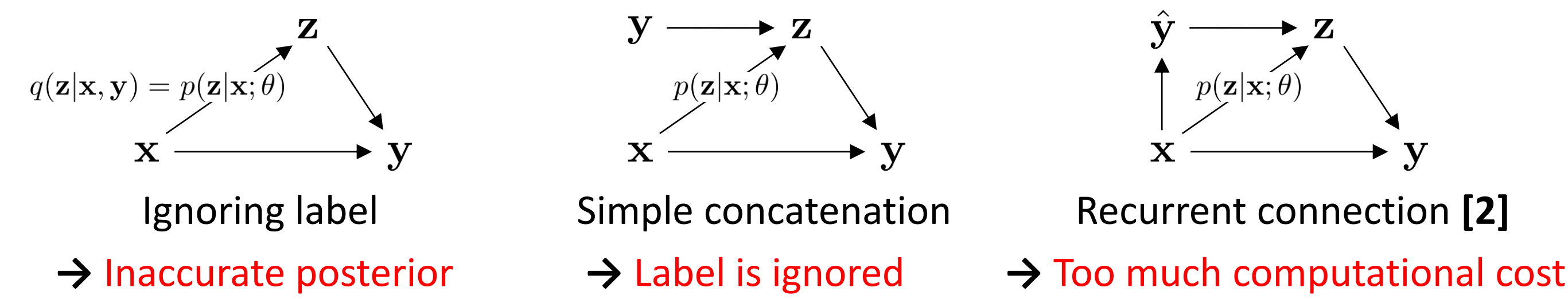
We maximize the evidence lower bound (ELBO).

$$\begin{aligned} & \log p(\mathbf{Y} | \mathbf{X}; \theta) \\ &= \log \int p(\mathbf{Y}, \mathbf{Z} | \mathbf{X}; \theta) d\mathbf{Z} \\ &\geq \int q(\mathbf{Z} | \mathbf{X}, \mathbf{Y}; \phi) \log \frac{p(\mathbf{Y} | \mathbf{X}, \mathbf{Z}; \theta) p(\mathbf{Z} | \mathbf{X}; \theta)}{q(\mathbf{Z} | \mathbf{X}, \mathbf{Y}; \phi)} d\mathbf{Z} \\ &= \sum_{i=1}^N \left\{ \mathbb{E}_{q(\mathbf{z}_i | \mathbf{x}_i, \mathbf{y}_i; \phi)} \left[\log p(\mathbf{y}_i | \mathbf{z}_i, \mathbf{x}_i; \theta) \right] - \text{KL} \left[q(\mathbf{z}_i | \mathbf{x}_i, \mathbf{y}_i; \phi) \| p(\mathbf{z}_i | \mathbf{x}_i; \theta) \right] \right\} \end{aligned}$$



Approximate Posterior

In modeling the approximate posterior $q(\mathbf{z} | \mathbf{x}, \mathbf{y}; \phi)$, how to utilize the label \mathbf{y} is not a straightforward matter.



Structural form of true posterior

- The relationship between \mathbf{z} and \mathbf{y} is relatively simple in DropMax.
- We encode this property into the approximate posterior.

True posterior is decomposed as : $p(\mathbf{z} | \mathbf{x}, \mathbf{y}) = \underbrace{p(\mathbf{z} | \mathbf{x})}_{\text{without label}} \times \underbrace{p(\mathbf{y} | \mathbf{z}, \mathbf{x}) / p(\mathbf{y} | \mathbf{x})}_{\text{with label}}$

$$\mathbf{g}(\mathbf{x}; \phi) = \text{sgm}(\bar{\mathbf{W}}_{\theta}^{\top} \mathbf{h} + \bar{\mathbf{b}}_{\theta} + \mathbf{r}(\mathbf{x}; \phi)), \quad \mathbf{r}(\mathbf{x}; \phi) = \mathbf{W}_{\phi}^{\top} \mathbf{h} + \mathbf{b}_{\phi}$$

Encoding label information

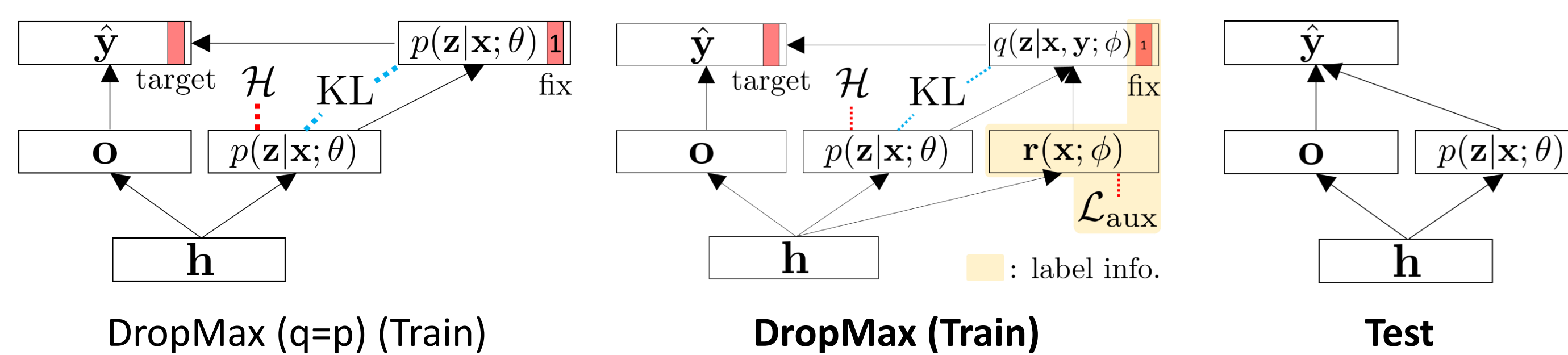
- [Obs 1] \mathbf{z}_t is positively and $\mathbf{z}_{k \neq t}$ is negatively correlated with \mathbf{y} .
- [Obs 2] With mean field approx., $p(z_t = 1 | \mathbf{x}, \mathbf{y}) = 1$, excluding $z_1 = \dots = z_K = 0$.

From [Obs 1], we simply regress $\text{sgm}(\mathbf{r})$ to \mathbf{y} :

$$\mathcal{L}_{\text{aux}}(\phi) = - \sum_i \sum_k \left\{ y_{i,k} \log \text{sgm}(r_k(\mathbf{x}_i; \phi)) + (1 - y_{i,k}) \log(1 - \text{sgm}(r_k(\mathbf{x}_i; \phi))) \right\}$$

From [Obs 2], we have

$$q(\mathbf{z} | \mathbf{x}, \mathbf{y}; \phi) = \text{Ber}(z_t; 1) \prod_{k \neq t} \text{Ber}(z_k; g_k(\mathbf{x}; \phi))$$



Regularized Variational Inference

$p(\mathbf{z} | \mathbf{x}; \theta)$ collapses into $q(\mathbf{z} | \mathbf{x}, \mathbf{y}; \phi)$ too easily, as $p(\mathbf{z} | \mathbf{x}; \theta)$ is parametric with input \mathbf{x} . We add entropy regularizer to $p(\mathbf{z} | \mathbf{x}; \theta)$.

$$\mathcal{H}(p(\mathbf{z} | \mathbf{x}; \theta)) = \sum_k \rho_k \log \rho_k + (1 - \rho_k) \log(1 - \rho_k)$$

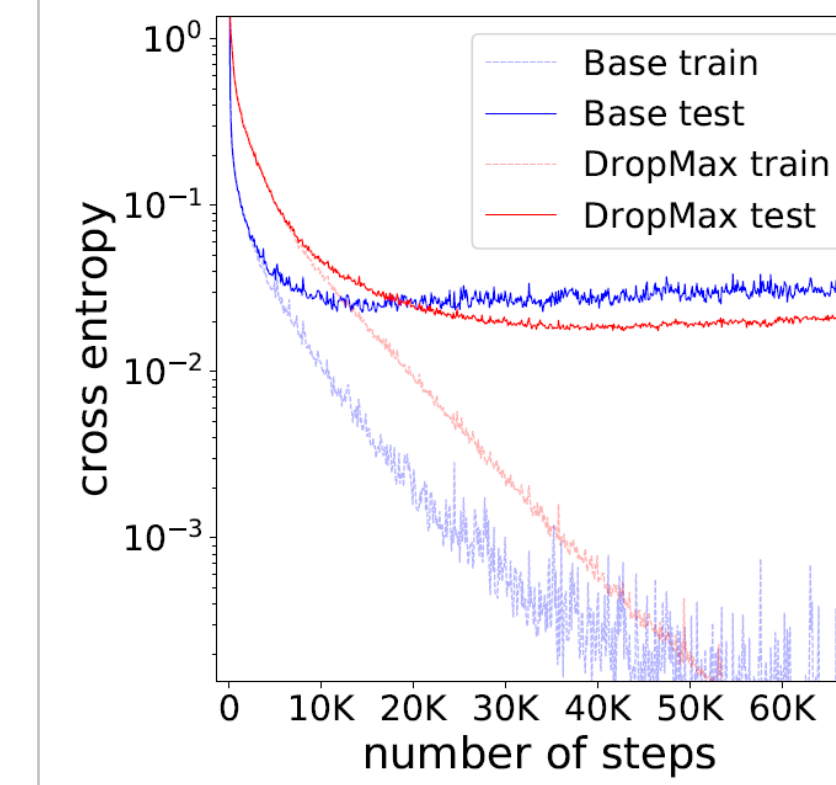
The KL divergence and the final objective is

$$\begin{aligned} \text{KL}[q(\mathbf{z} | \mathbf{x}, \mathbf{y}; \phi) \| p(\mathbf{z} | \mathbf{x}; \theta)] &= \log \frac{1}{\rho_t} + \sum_{k \neq t} g_k \log \frac{g_k}{\rho_k} + (1 - g_k) \log \frac{1 - g_k}{1 - \rho_k} \\ \mathcal{L}(\theta, \phi) &= \sum_{i=1}^N \left\{ \frac{1}{S} \sum_{s=1}^S - \log p(\mathbf{y}_i | \mathbf{x}_i, \mathbf{z}_i^{(s)}; \theta) + \text{KL}[q(\mathbf{z} | \mathbf{x}, \mathbf{y}; \phi) \| p(\mathbf{z} | \mathbf{x}; \theta)] - \mathcal{H} \right\} + \mathcal{L}_{\text{aux}} \end{aligned}$$

Experiments

1. Generalization Performance

Learning curve (MNIST)



Classification error (%)

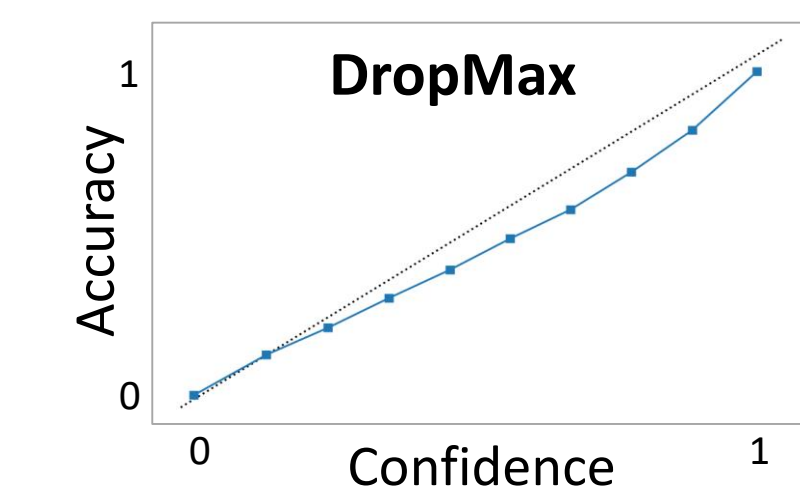
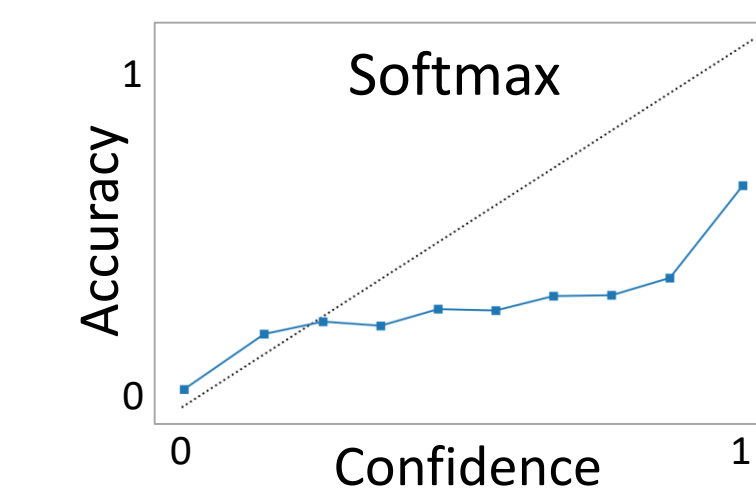
Method	MNIST-1K	CIFAR-100	AWA	CUB
Softmax	7.09	30.60	30.29	48.84
Sparsemax [3]	6.57	31.41	36.06	64.41
Sampled Softmax [4]	7.36	30.87	29.81	49.90
Class Dropout	7.19	30.78	31.11	48.87
Deterministic Attention	6.91	30.60	30.98	49.97
DropMax (q=p)	7.52	29.98	29.27	42.08
DropMax	5.32	29.87	26.91	41.07

- Class Dropout : Randomly drops out non-target classes with a predefined probability.
- Deterministic (sigmoid) Attentions : are multiplied to the softmax exponentiations.
- DropMax (q=p) : A variant of DropMax where we let $q(\mathbf{z} | \mathbf{x}, \mathbf{y}) = p(\mathbf{z} | \mathbf{x}; \theta)$.

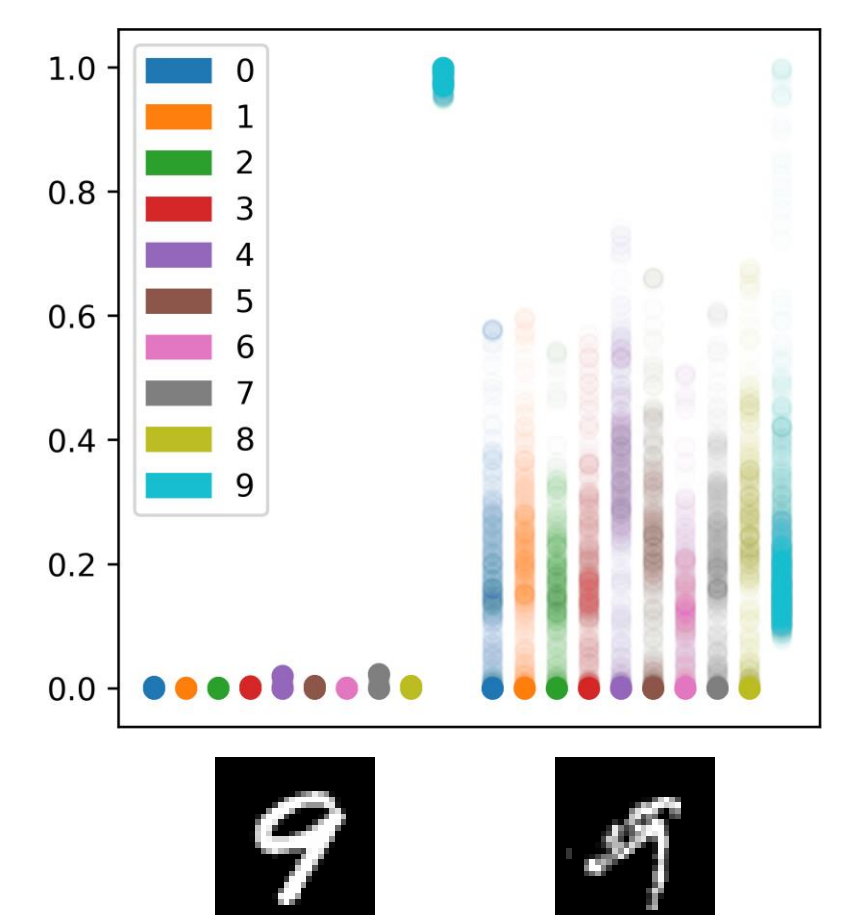
2. Reliability / Calibration

Classification error / Expected calibration error (ECE)

	MNIST-1K Random Background	MNIST-1K Background & Rotation
Error (%)	12.89	52.06
ECE (%)	1.52	7.69
DropMax	10.86	48.81
ECE (%)	0.683	2.38

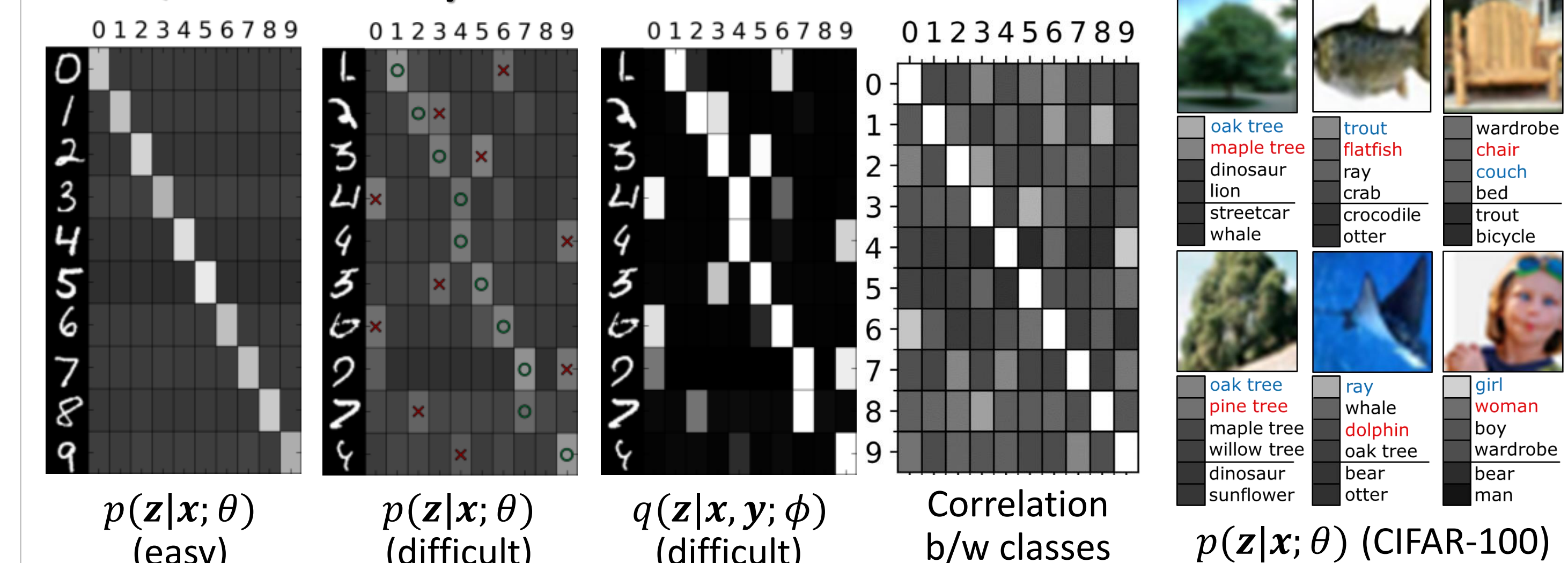


Predictive distribution



- DropMax improves calibration of output probabilities.

3. Qualitative Analysis



- When easy, $p(\mathbf{z} | \mathbf{x}; \theta)$ pre-classifies with high confidence.
- When difficult, $p(\mathbf{z} | \mathbf{x}; \theta)$ selects multiple relevant classes for each instance differently.

References

- J. Ba and B. Frey. Adaptive dropout for training deep neural networks. In *NIPS*, 2013.
- K. Sohn, H. Lee, and X. Yan. Learning structured output representation using deep conditional generative models. In *NIPS*, 2015.
- A. F. T. Martins and R. Fernandez Astudillo. From Softmax to Sparsemax: A Sparse Model of Attention and Multi-Label Classification. In *ICML*, 2016.
- S. Jean, K. Cho, R. Memisevic, and Y. Bengio. On Using Very Large Target Vocabulary for Neural Machine Translation. In *ACL*, 2015.