

Unsupervised Anomaly Detection via Variational Auto-Encoder for Seasonal KPIs in Web Applications

Haowen Xu¹ Wenxiao Chen¹ Nengwen Zhao¹ Zeyan Li¹
Jiahao Bu¹ Zhihan Li¹ Ying Liu¹ Youjian Zhao¹ Dan Pei¹
Yang Feng² Jie Chen² Zhaogang Wang² Honglin Qiao²

¹Tsinghua University

²Alibaba Group

April 26, 2018

Outline

- 1 Background
- 2 Architecture
- 3 Evaluation
- 4 Analysis
- 5 Conclusion

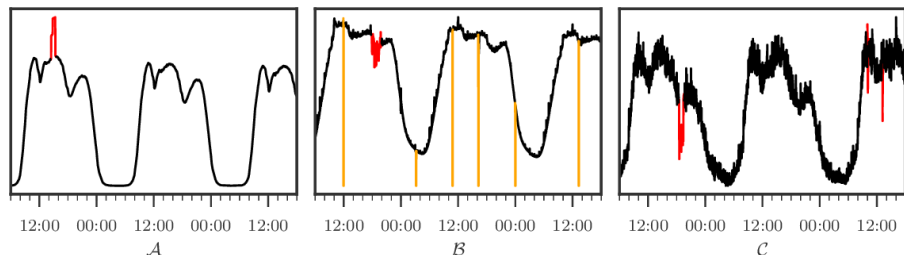
Outline

- 1 Background
- 2 Architecture
- 3 Evaluation
- 4 Analysis
- 5 Conclusion

Problem Scenario: Anomaly Detection for Seasonal KPIs

KPIs are time sequences, yet one of the most fundamental system monitoring indicators. A failure usually causes more or less anomalies on at least one KPI. Thus anomaly detection for KPIs are very useful in **Artificial Intelligence for IT Operations (AIOps)**.

For web applications, the user activities are usually seasonal, so are the KPIs, including high level KPIs like the trading volumes, and low level KPIs like the CPU consumptions. We thus focus on **anomaly detection for seasonal KPIs in this work**.



Problem Formulation: Detection of “Abnormal” Patterns

Since KPIs are time sequences, and since in most real cases human operators are willing to see a detection output every time a new observation arrives, the anomaly detection for KPIs can be formulated as:

Anomaly Detection for KPIs

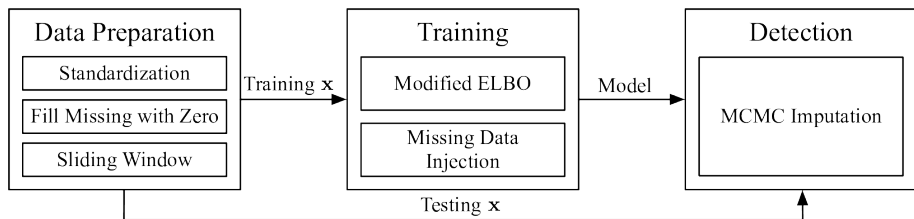
For each time t , given the on-time KPI observation x_t and historical observations $x_{t-W+1}, \dots, x_{t-1}$, determine whether an “abnormal” pattern has occurred (denoted by $y_t = 1$).

Detection algorithms are often designed to compute a **real-valued score** $s(y_t = 1)$ (“**anomaly score**” hereafter), e.g., $p(y_t = 1 | x_{t-W+1}, \dots, x_t)$, leaving the final decision of triggering alerts to the operators.

Outline

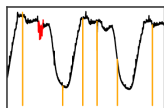
- 1 Background
- 2 Architecture**
- 3 Evaluation
- 4 Analysis
- 5 Conclusion

Overall Architecture of *Donut*



Data Preparation

- Fill Missing with Zero:



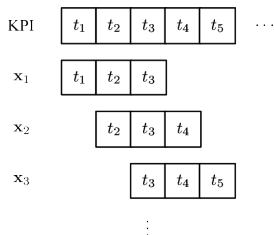
“Missing” are special anomalies, *always* known beforehand. We fill missing points with zeros (orange points in the left figure), and let our model to handle them afterwards.

- Standardization: $\hat{x}_t = (x_t - \mu_x) / \sigma_x$.

x_t are the original KPI values, μ_x and σ_x are the mean and std of x_t .

We shall use \hat{x}_t to denote \hat{x}_t and neglect the original values x_t *hereafter*.

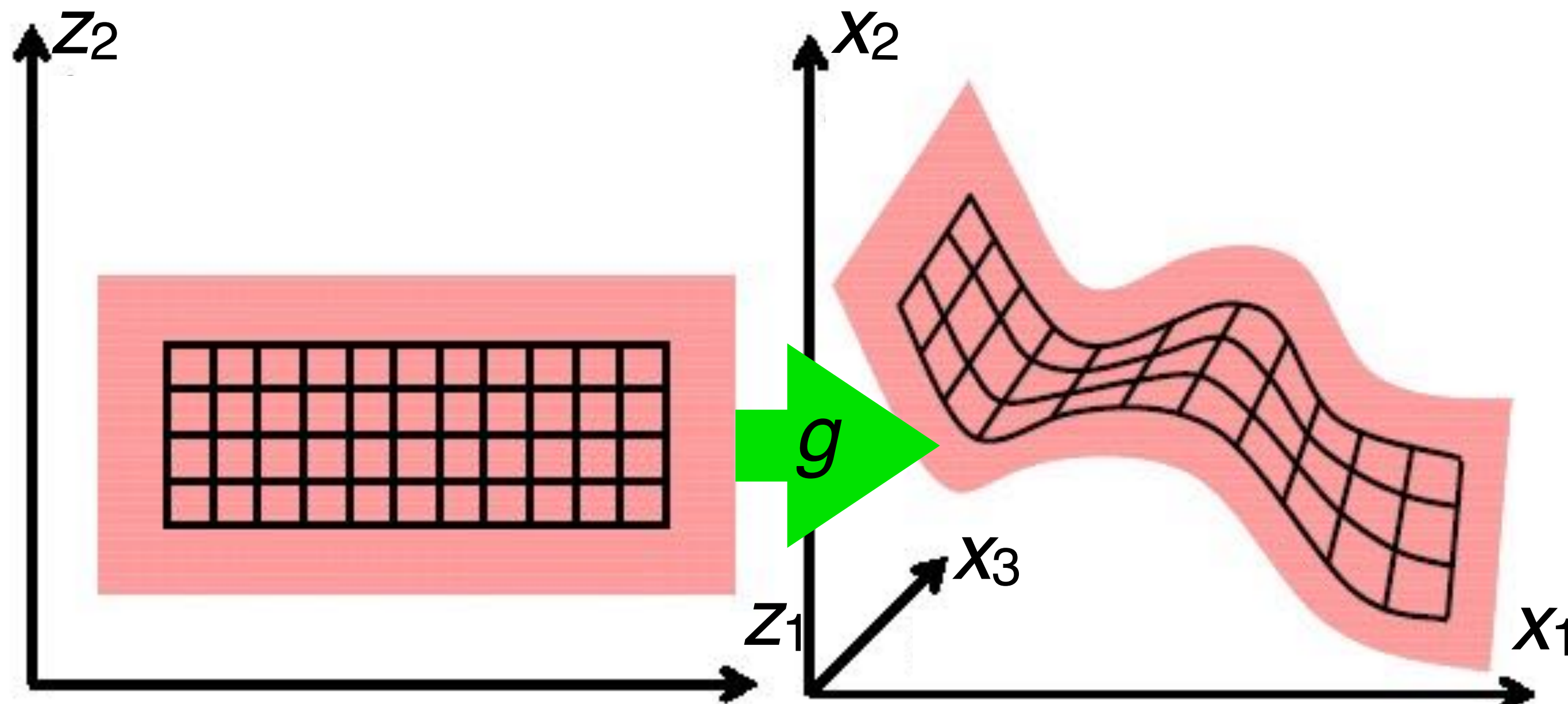
- Sliding Window:



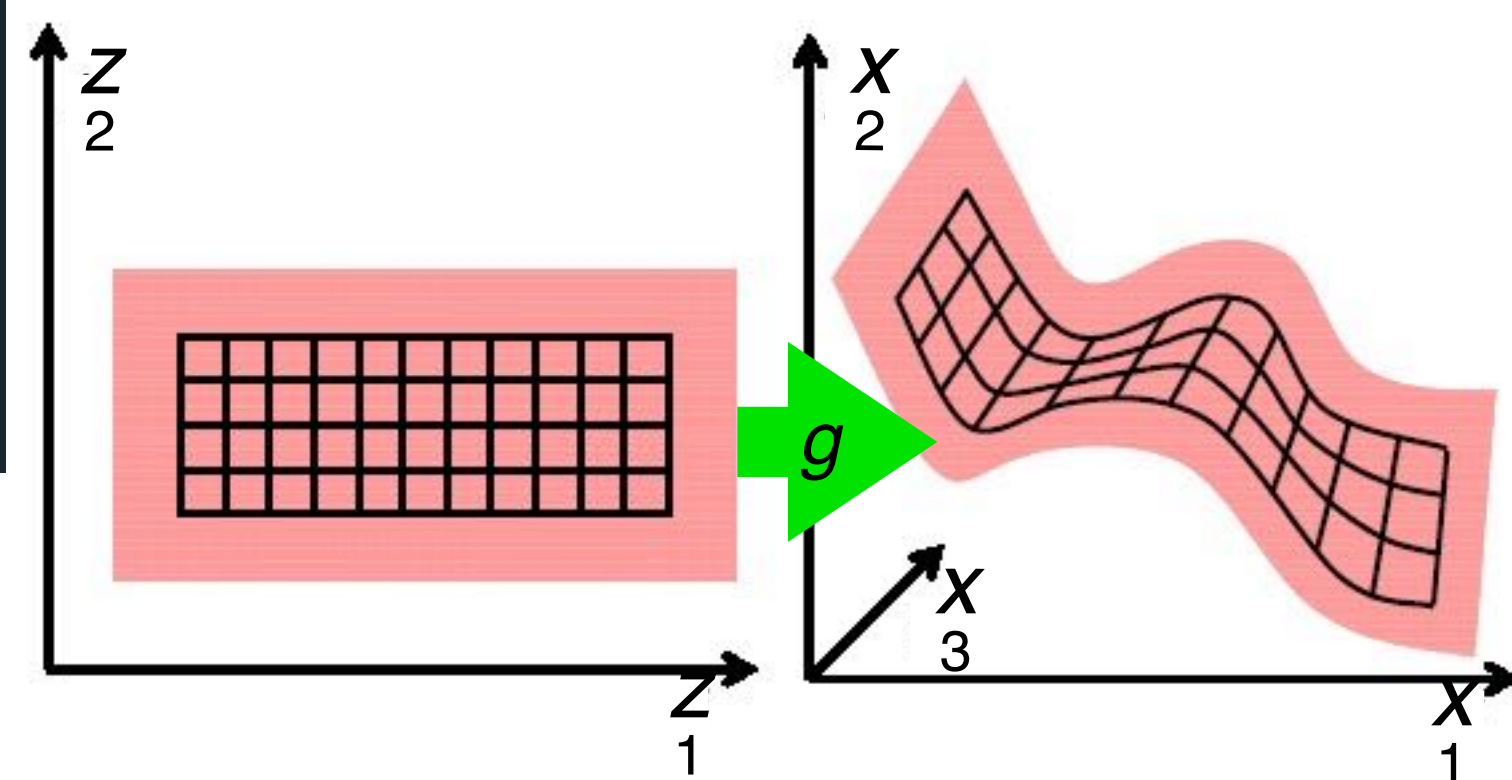
We split the KPIs into fixed-length **sliding windows** \mathbf{x}_t , which are assumed to be *i.i.d.*, and are used as **the input \mathbf{x} of VAE** at every time t . For simplicity, we shall omit the subscript t , using \mathbf{x} to denote **the window of “current time”**, and x_1, \dots, x_W to denote **each point in \mathbf{x}** afterwards.

Latent Variable Models

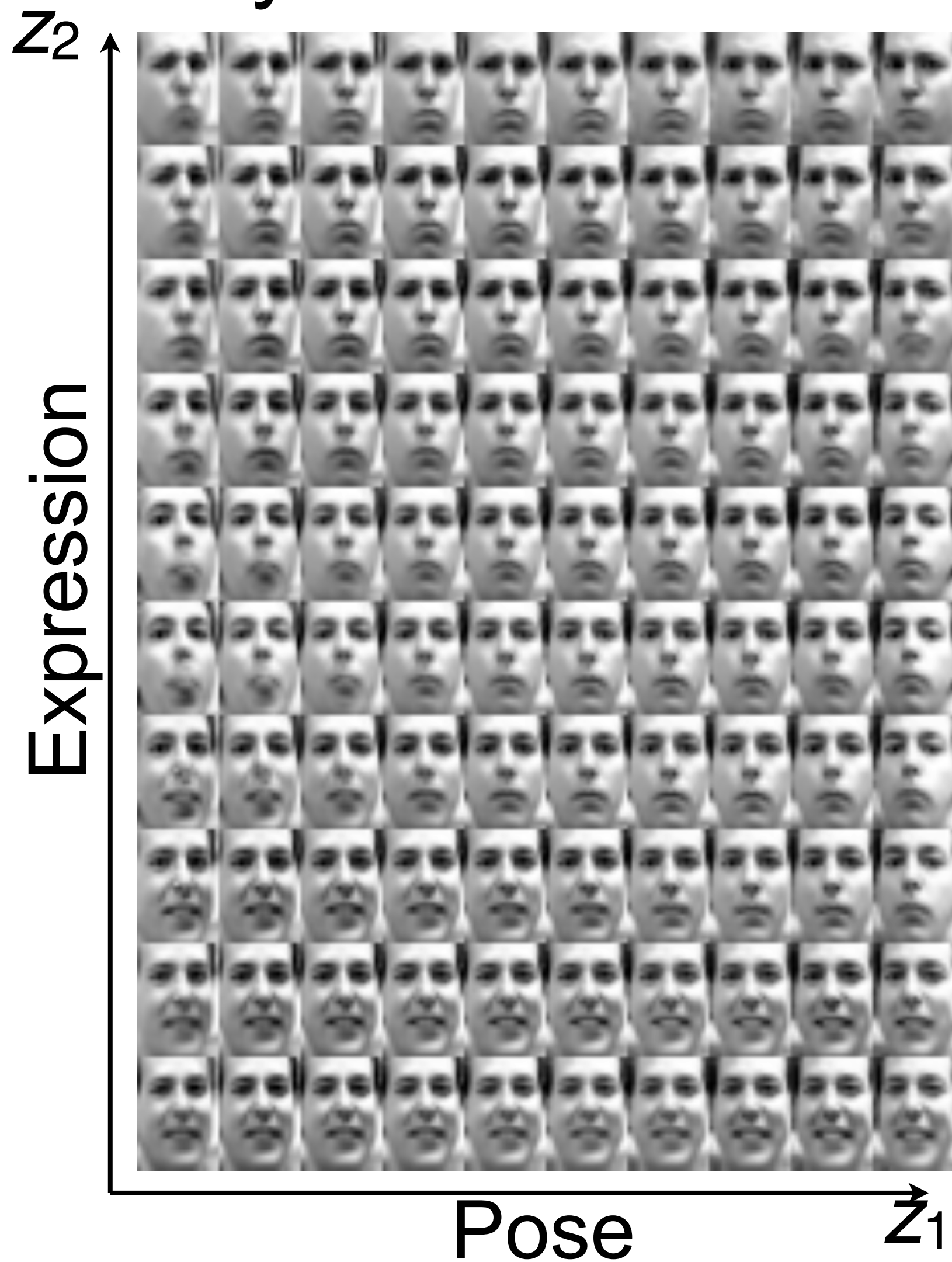
- The Variational Autoencoder model:
 - Kingma and Welling, *Auto-Encoding Variational Bayes*, *International Conference on Learning Representations (ICLR) 2014*.
 - Rezende, Mohamed and Wierstra, *Stochastic back-propagation and variational inference in deep latent Gaussian models*. ICML 2014.



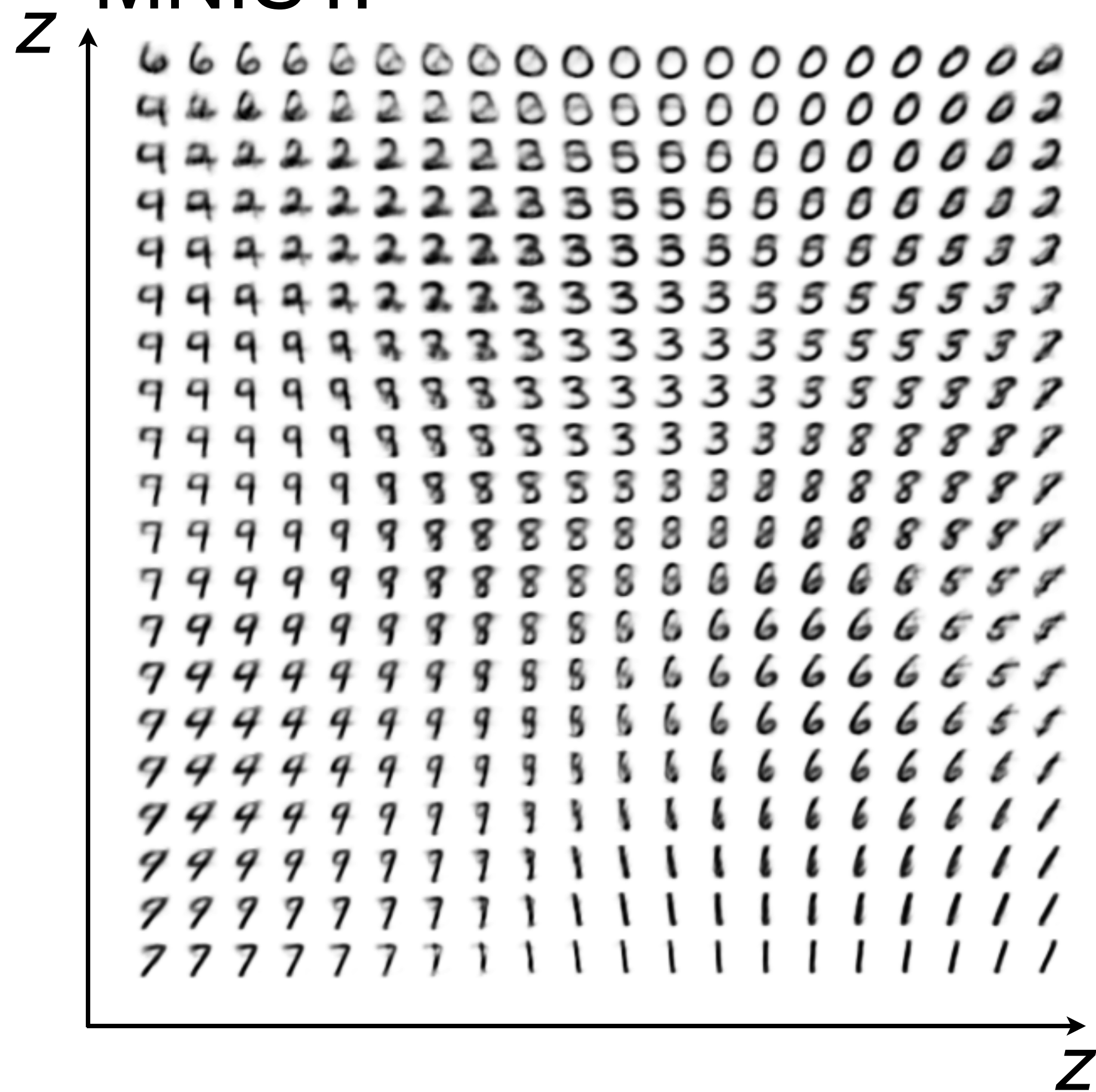
Latent Variable Models



Frey Faces:



MNIST:



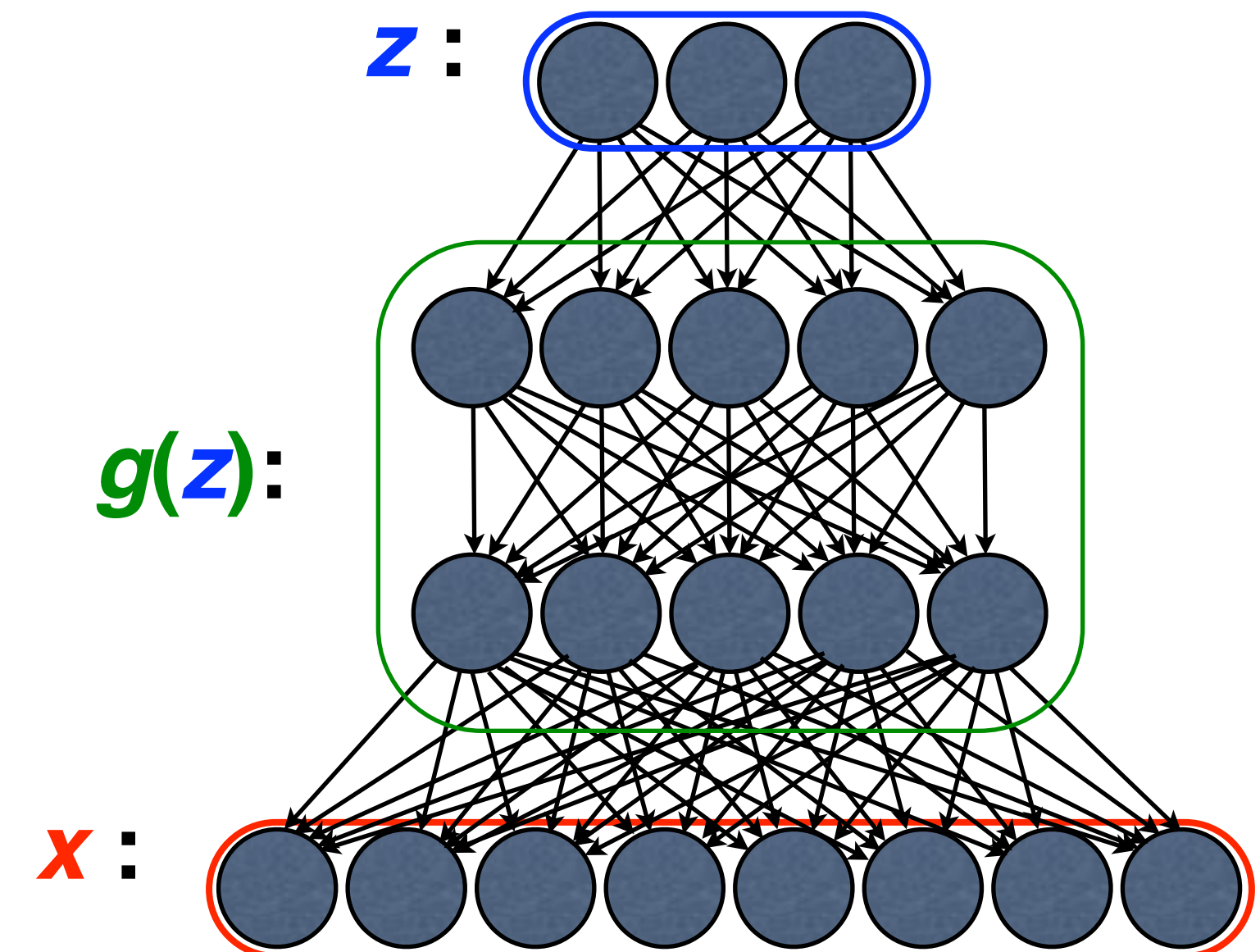
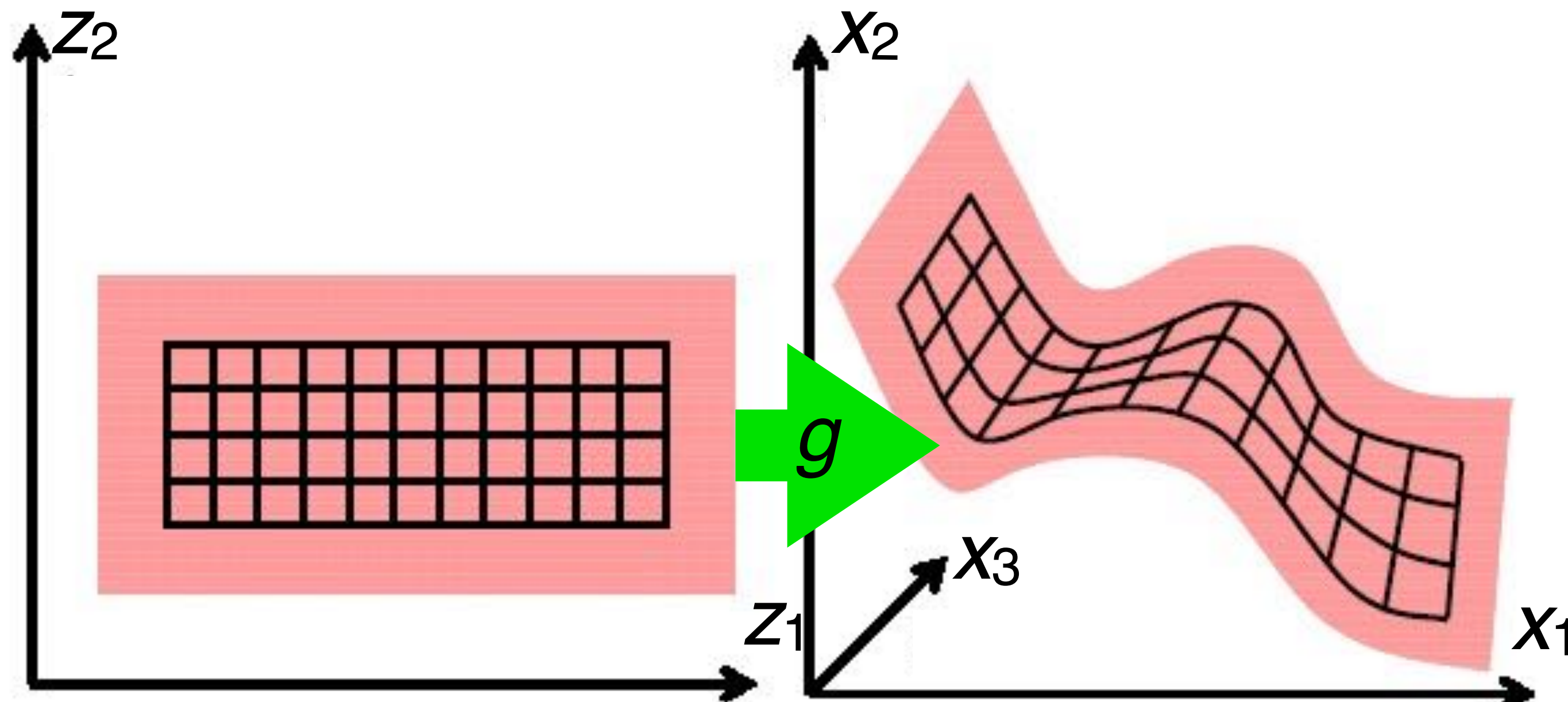
Latent Variable Models

- **latent variable model**: learn a mapping from some latent variable z to a complicated distribution on x .

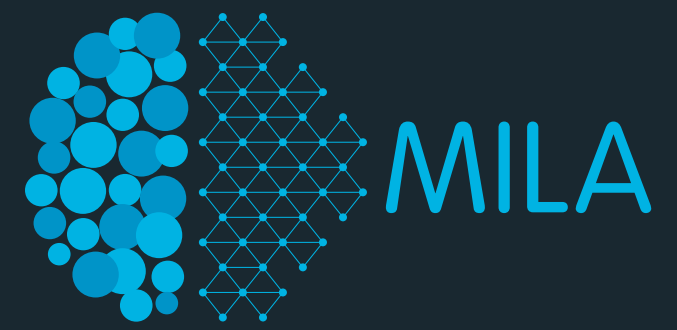
$$p(x) = \int p(x, z) dz \quad \text{where } p(x, z) = p(x | z)p(z)$$

prior $p(z) = \text{something simple}$ $p(x | z) = g(z)$

- Can we learn to decouple the true **explanatory factors** underlying the data distribution?
E.g. separate identity and expression in face images



Variational Auto-Encoder (VAE)



- Where does z come from? — The classic DAG problem.
- The VAE approach: introduce an inference machine $q_\phi(z | x)$ that **learns** to approximate the posterior $p_\theta(z | x)$.
- Define a **variational lower bound** on the data likelihood: $p_\theta(x) \geq \mathcal{L}(\theta, \phi, x)$

$$\begin{aligned}\mathcal{L}(\theta, \phi, x) &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z | x)] \\ &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z) + \log p_\theta(z) - \log q_\phi(z | x)] \\ &= -D_{\text{KL}}(q_\phi(z | x) || p_\theta(z)) + \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)]\end{aligned}$$

regularization term **reconstruction term**

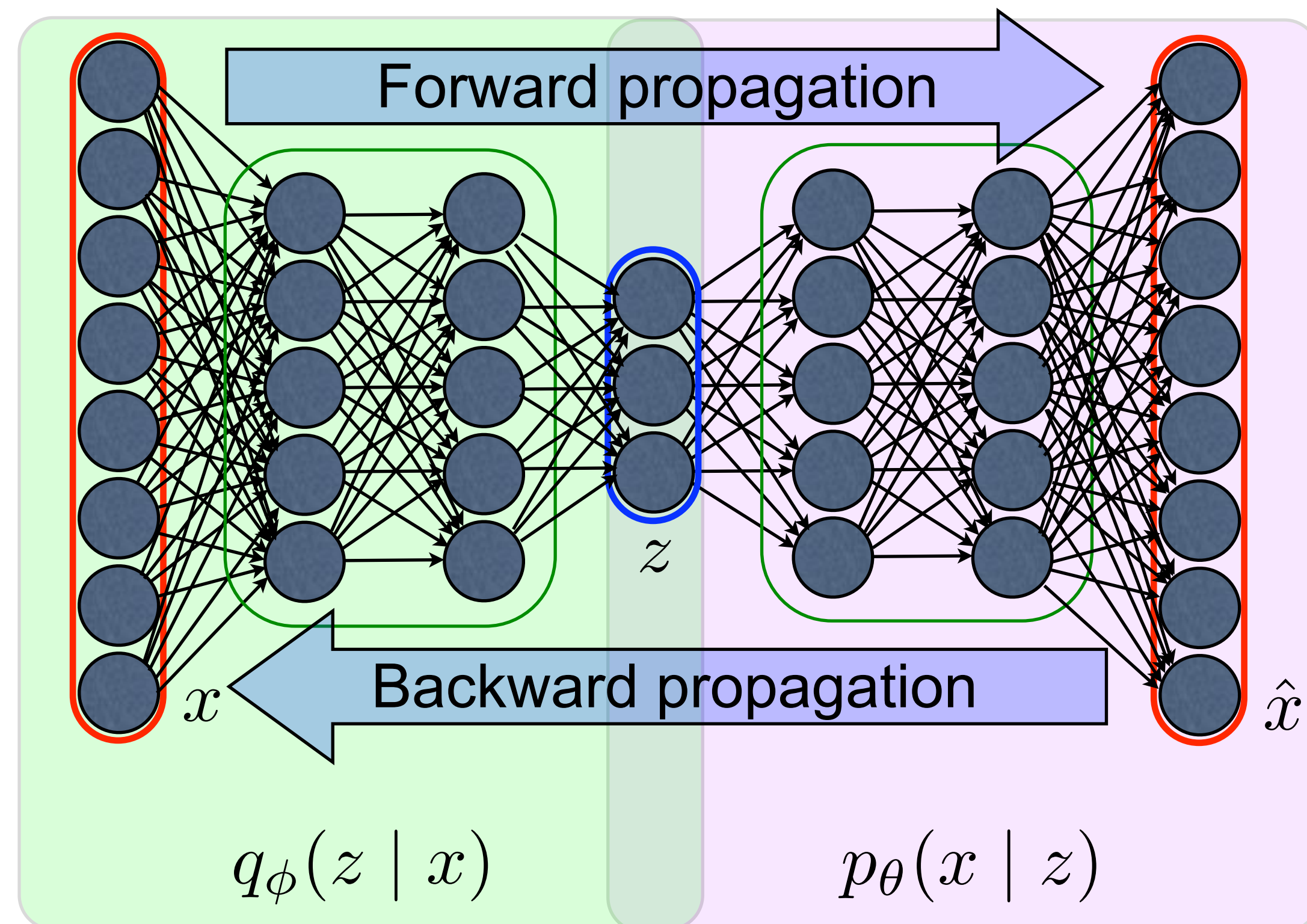
- What is $q_\phi(z | x)$?

Training with backpropagation!

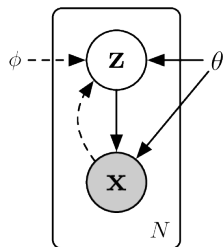


- Due to a **reparametrization** trick, we can simultaneously train both the generative model $p_{\theta}(x | z)$ and the inference model $q_{\phi}(z | x)$ by optimizing the variational bound using gradient **backpropagation**.

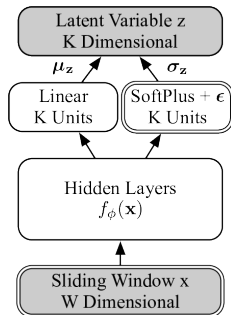
Objective function: $\mathcal{L}(\theta, \phi, x) = -D_{\text{KL}}(q_{\phi}(z | x) || p_{\theta}(z)) + \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x | z)]$



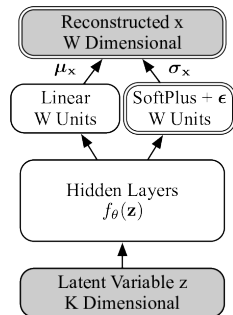
Network Structure



(a) VAE General Structure



(b) $q_\phi(\mathbf{z}|\mathbf{x})$ of Donut



(c) $p_\theta(\mathbf{x}|\mathbf{z})$ of Donut

- Variational net: $q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z^2 \mathbf{I})$.
- Generative net: $p_\theta(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$, $p_\theta(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\sigma}_x^2 \mathbf{I})$.
- SoftPlus Trick: $\boldsymbol{\sigma}_z = \text{SoftPlus}[\mathbf{W}_{\sigma_z}^\top f_\phi(\mathbf{x}) + \mathbf{b}_{\sigma_z}] + \epsilon$, $\text{SoftPlus}[a] = \log[\exp(a) + 1]$. Similar for $\boldsymbol{\sigma}_x$. (otherwise, unbounded)

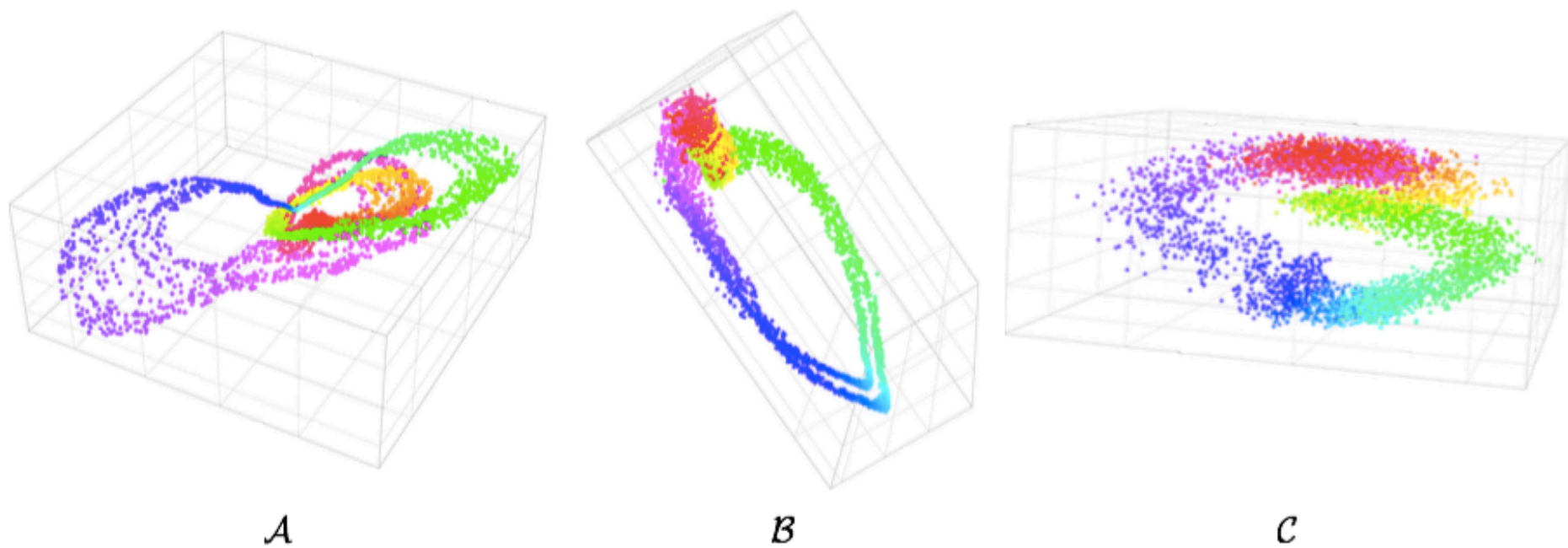
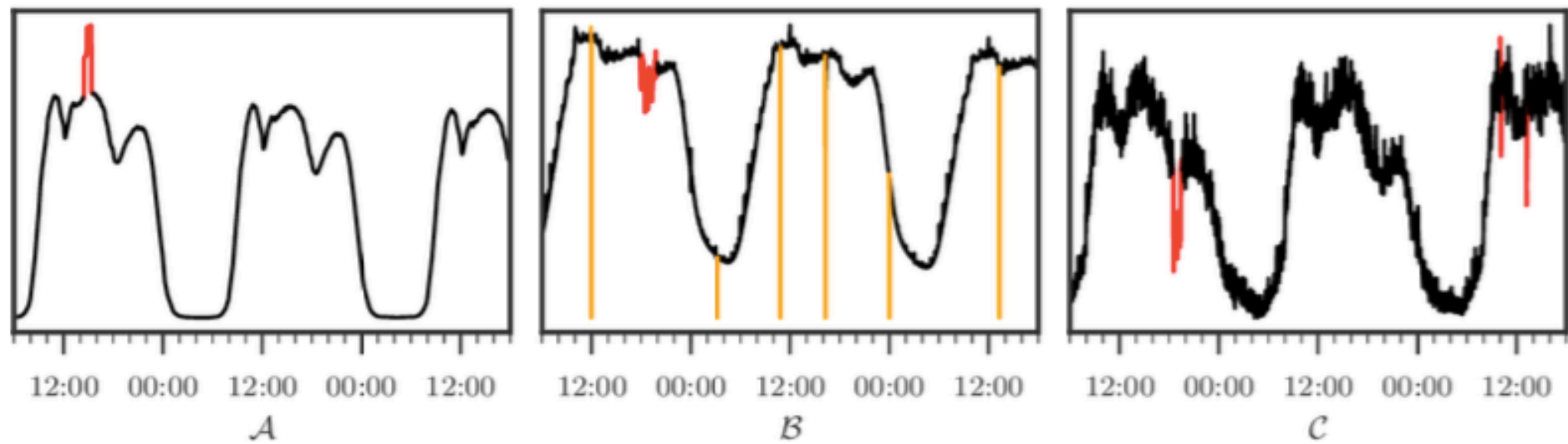


Figure 12: 3-d latent space of all three datasets.

Dealing with Missing and Anomaly

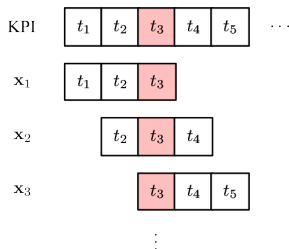


Figure: The anomaly at t_3 shall affect t_4 and t_5 , potentially causing trouble in training and detection.

We use three techniques to handle such “historical anomalies”.

- 1 **M-ELBO:** We modify the ELBO (objective function) of VAE into **M-ELBO** $\tilde{\mathcal{L}}(\mathbf{x})$:

$$\tilde{\mathcal{L}}(\mathbf{x}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\sum_{w=1}^W \alpha_w \log p_{\theta}(x_w|\mathbf{z}) + \beta \log p_{\theta}(\mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}) \right]$$

$\alpha_w=0$ when anomalous or missing, and $\alpha_w = 1$ otherwise; β is the ratio of normal data

- 2 **Missing Data Injection:** We randomly set 1% points to be missing at **every epoch** in training, to compensate for having no anomaly labels in the unsupervised scenario.
- 3 **MCMC Imputation (Rezende et al., 2014):** In detection, we adopt this technique on known missing points.

The Time Gradient

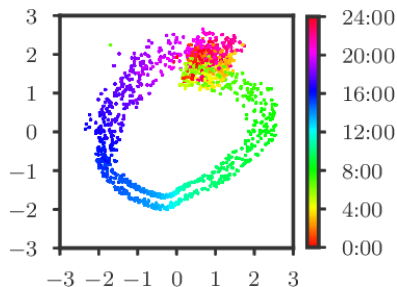


Figure: The \mathbf{z} layout of dataset \mathcal{B} . Figure is plotted by sampling \mathbf{z} from $q_\phi(\mathbf{z}|\mathbf{x})$, corresponding to normal \mathbf{x} randomly chosen from the testing set. K is chosen as 2, so the x - and y -axis are the two dimensions of \mathbf{z} samples. The color of a \mathbf{z} sample denotes its time of the day.

- Time gradient: $q_\phi(\mathbf{z}|\mathbf{x})$ are organized in smooth transition: \mathbf{x} at contiguous time are mapped to nearby $q_\phi(\mathbf{z}|\mathbf{x})$.
- Contiguous \mathbf{x} are highly similar in the KPIs of our interest, since they are smooth in general.
- Transition of $q_\phi(\mathbf{z}|\mathbf{x})$ in the shape of \mathbf{x} , rather than time, is the cause of time gradient, since *Donut* consumes no time information.
- *Donut* encodes the “shape” or “normal patterns” of \mathbf{x} by \mathbf{z} , as shown by the time gradient.
- The time gradient can benefit generalization.

The Causes of Time Gradient

$$\begin{aligned}\mathcal{L}(\mathbf{x}) &= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z}) + \log p_\theta(\mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E} [\log p_\theta(\mathbf{x}|\mathbf{z})] + \mathbb{E} [\log p_\theta(\mathbf{z})] + H[\mathbf{z}|\mathbf{x}]\end{aligned}$$

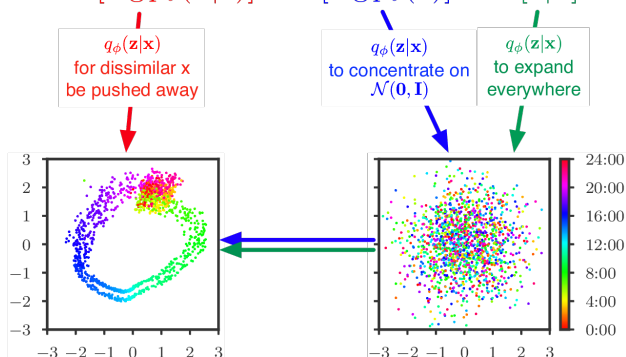


Figure: Causes of the time gradient. Surprisingly, we find no term in ELBO directly pulling $q_\phi(\mathbf{z}|\mathbf{x})$ for similar \mathbf{x} together. The time gradient is likely to be caused mainly by **expansion** ($H(\mathbf{z}|\mathbf{x})$), **squeezing** ($\mathbb{E}[\log p_\theta(\mathbf{z})]$), **pushing** ($\mathbb{E}[\log p_\theta(\mathbf{x}|\mathbf{z})]$), and the **training dynamics** (random initialization and SGVB).

The KDE Interpretation

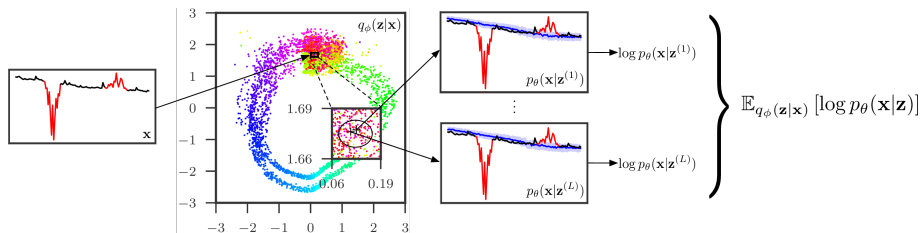


Figure: Illustration of the KDE interpretation. **For a given \mathbf{x} potentially with anomalies, Donut tries to recognize what normal pattern it follows, encoded as $q_\phi(\mathbf{z}|\mathbf{x})$.** The black ellipse in the middle figure denotes the $3\text{-}\sigma_{\mathbf{z}}$ region of $q_\phi(\mathbf{z}|\mathbf{x})$. **L samples of \mathbf{z} are then taken from $q_\phi(\mathbf{z}|\mathbf{x})$, denoted as the crosses in the middle figure. Each \mathbf{z} is associated with a density estimator kernel $\log p_\theta(\mathbf{x}|\mathbf{z})$.** The blue curves in the right two figures are $\mu_{\mathbf{x}}$ of each kernel, while the surrounding stripes are $\sigma_{\mathbf{x}}$. Finally, the **values of $\log p_\theta(\mathbf{x}|\mathbf{z})$ are computed from each kernel, and further averaged together as the reconstruction probability.**

The Anomaly Score

An and Cho (2015) has already adopted VAE in anomaly detection tasks of other domain¹. They use the **reconstruction probability** (1) of truly *i.i.d.* samples \mathbf{x} (e.g., image pixel vectors) as the anomaly score:

$$s(y = 1) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] \approx \frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}|\mathbf{z}^{(l)}), \mathbf{z}^{(l)} \sim q_\phi(\mathbf{z}|\mathbf{x}) \quad (1)$$

Since the KPIs are time sequences, and the operators are willing to see on-time detection outputs each time a new point arrives, we compute the **element-wise reconstruction probability** (2) for the last point x_W in \mathbf{x} , as the anomaly score for the time being:

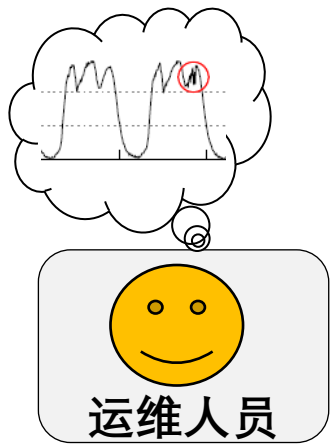
$$s(y_W = 1) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(x_W|\mathbf{z})] \approx \frac{1}{L} \sum_{l=1}^L \log p_\theta(x_W|\mathbf{z}^{(l)}), \mathbf{z}^{(l)} \sim q_\phi(\mathbf{z}|\mathbf{x}) \quad (2)$$

¹An and Cho (2015) uses vanilla VAE, without developing techniques like ours to improve performance. We shall compare *Donut* against their vanilla VAE in evaluation.

Outline

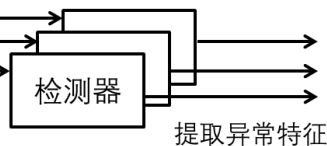
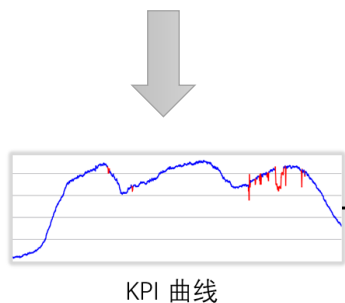
- 1 Background
- 2 Architecture
- 3 Evaluation**
- 4 Analysis
- 5 Conclusion

关键技术: 自动挑选算法和参数组合

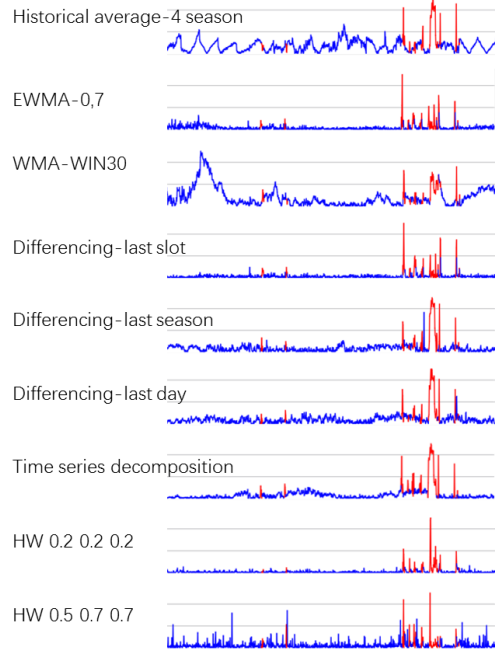


解决问题：大量差异化的指标，大量复杂算法及参数，导致无从选择

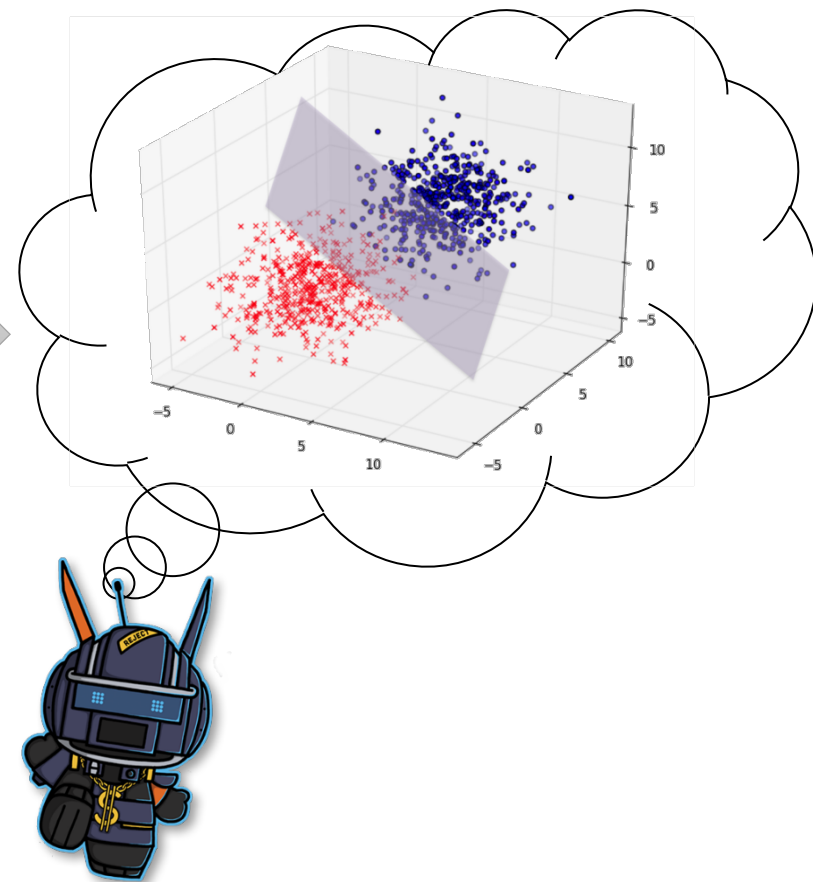
核心思想：多检测器作为特征+集成学习，避免人工挑选算法和参数



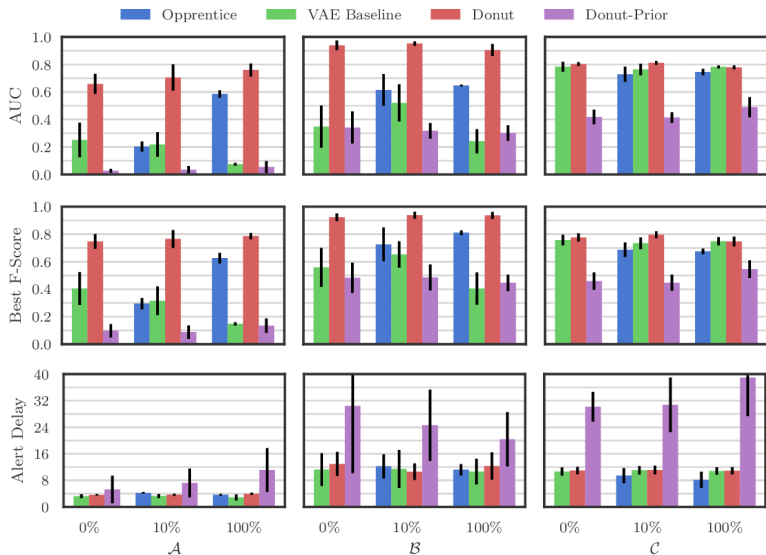
133种算法和参数组合



异常特征空间分类



Overall Performance



Effects of *Donut* Techniques

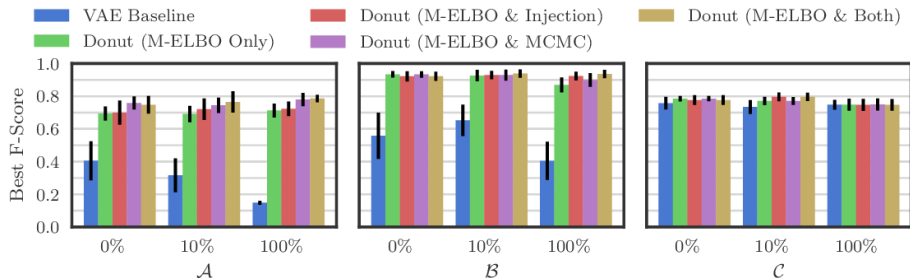


Figure: Best F-score of (1) VAE Baseline, (2) *Donut* with M-ELBO, (3) M-ELBO + missing data injection, (4) M-ELBO + MCMC, and (5) M-ELBO + both MCMC and injection.

The **M-ELBO** alone contributes most of the improvement over VAE Baseline, while the **missing data injection** and the **MCMC imputation** can further benefit the performance.

Impact of Z Dimension Number K

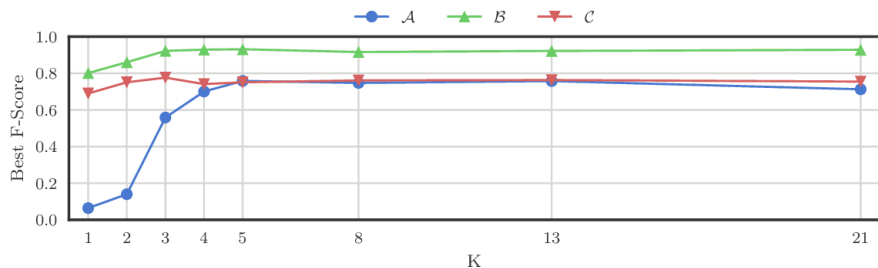


Figure: The best F-score of unsupervised *Donut* with different K on testing set.

- The essential of **Dimension reduction**: W (the dimension of \mathbf{x}) is 120, while the best K (the dimension of \mathbf{z}) is no larger than 10.
- It should be quite easy to empirically choose K .
 - 1 The best performance could be achieved with fairly small K .
 - 2 The performance does not drop too heavily for K up to 21.
- Smoother KPIs seem to demand larger K .

Outline

- 1 Background
- 2 Architecture
- 3 Evaluation
- 4 Analysis
- 5 Conclusion**

Conclusion

Our unsupervised anomaly detection algorithm *Donut* for seasonal KPIs, based on VAE, greatly outperforms state-of-art supervised and vanilla VAE anomaly detection algorithms. The best F-scores range from 0.75 to 0.90 for the studied KPIs. The key factors of *Donut* to be successful are:

- **Dimension Reduction**: forces *Donut* to focus on the overall shape of normal patterns, and gain the ability of resisting abnormal points.
- **M-ELBO, Missing Data Injection and MCMC Imputation**: further improves *Donut*'s ability to resist abnormal points.

Furthermore, we made the **KDE Interpretation**, which provides a new perspective of VAE-based KPI anomaly detection. All of the above factors can be verified by such interpretation. The KDE Interpretation potentially has more theoretical value in the further development of deep generative models for KPI anomaly detection.

Donut source code published at: <https://github.com/korepwx/donut>.
Full slide: <https://github.com/korepwx/donut/tree/slide>.

Thank you

Q & A

- An, J. and Cho, S. (2015). Variational autoencoder based anomaly detection using reconstruction probability. Technical report, SNU Data Mining Center.
- Rezende, D. J., Mohamed, S., and Wierstra, D. (2014). Stochastic Backpropagation and Approximate Inference in Deep Generative Models. In *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32, ICML'14*, pages II-1278-II-1286, Beijing, China. JMLR.org.