

Group Meeting - September 11, 2020

Paper review & Research progress

Truong Son Hy *

*Department of Computer Science
The University of Chicago

Ryerson Physical Lab



Sun Tzu

The supreme art of war is to subdue the enemy **without** fighting.



The Great Chinggis Khaan (Genghis Khan)

Violence **never** settles anything.



Buddha

- ① To **understand** everything is to **forgive** everything.
- ② We can never obtain peace in the outer world until we make **peace with ourselves**.



The Art Institute of Chicago



Papers

- ① **Learning Hierarchical Features from Generative Models** (ICML 2017)
<https://arxiv.org/pdf/1702.08396.pdf>
- ② **Improved Variational Inference with Inverse Autoregressive Flow** (NIPS 2016)
<https://arxiv.org/pdf/1606.04934.pdf>
- ③ **NVAE: A Deep Hierarchical Variational Autoencoder**
<https://arxiv.org/pdf/2007.03898.pdf>
- ④ **Hierarchical Implicit Models and Likelihood-Free Variational Inference** (NIPS 2017)
<https://arxiv.org/pdf/1702.08896.pdf>
- ⑤ **Deep Hierarchical Implicit Models**
<https://arxiv.org/pdf/1702.08896v1.pdf>
- ⑥ **Discrete Object Generation with Reversible Inductive Construction**
<https://arxiv.org/pdf/1907.08268.pdf>

Note: Paper 1 criticizes paper 2. But both papers claim valid points, and have some frauds. Paper 3 is tunning of paper 2. Paper 5 is a trivial extension of paper 4.



Paper 1

Learning Hierarchical Features from Generative Models (ICML 2017)

Shengjia Zhao, Jiaming Song, Stefano Ermon (Stanford)

<https://arxiv.org/pdf/1702.08396.pdf>



Motivations:

- A key property of deep feed-forward networks is that they tend to learn increasingly abstract and invariant representations at higher levels in the hierarchy.

Samoyed (1.6); Papillon (5.7); Pomeranian (2.7); Arctic fox (1.0); Eskimo dog (0.6); white wolf (0.4); Siberian husky (0.4)

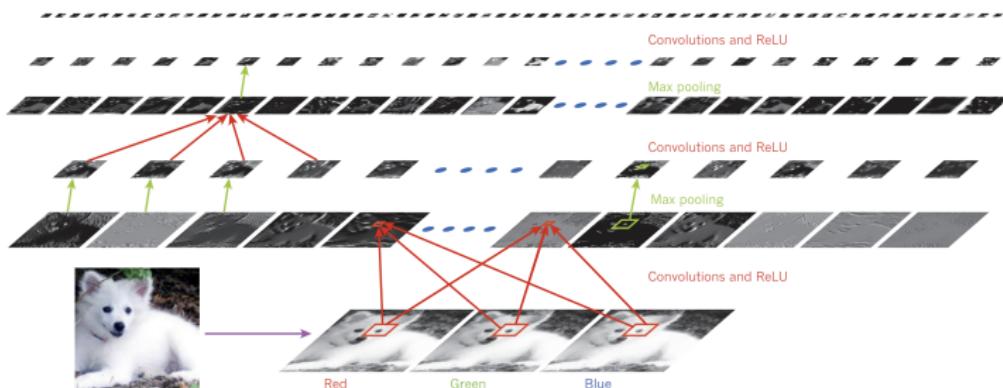


Figure 2 | Inside a convolutional network. The outputs (not the filters) of each layer (horizontally) of a typical convolutional network architecture applied to the image of a Samoyed dog (bottom left; and RGB (red, green, blue) inputs, bottom right). Each rectangular image is a feature map

corresponding to the output for one of the learned features, detected at each of the image positions. Information flows bottom up, with lower-level features acting as oriented edge detectors, and a score is computed for each image class in output. ReLU, rectified linear unit.

Deep Learning (Nature), Yann LeCun, Yoshua Bengio, Geoffrey Hinton



Motivations:

- Generative models with a hierarchical structure, where there are multiple layers of latent variables, have been **less successful** compared to their supervised counter-parts.
- Why:
 - **For classification:** Local features are often sufficient.
 - **For generation (of images):** Preservation of details are needed. The size and location of a sub-part often has to be dependent on the other sub-parts ← For example, an eye should only be generated with the same size as the other eye, at symmetric locations with respect to the center of the face, with appropriate distance between them.



Learning Hierarchical Features from Generative Models (3)

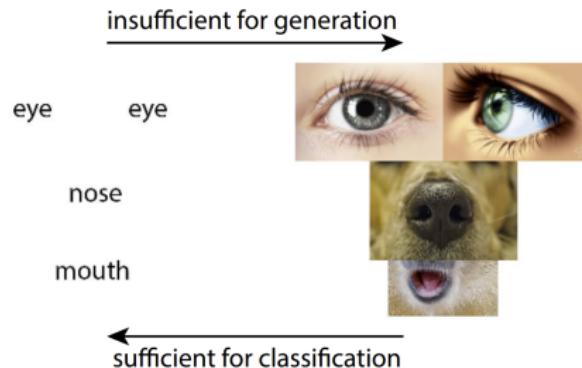


Figure 1. Left: Body parts feature detectors only carry a small amount of information about an underlying image, yet, it is sufficient for a confident classification as a face. Right: if a hierarchical generative model attempts to reconstruct an image based on these high-level features, it could generate inconsistent images, even when each part can be perfectly generated. Even though this “face” is clearly absurd, Google cloud platform classification API can identify with 93% confidence that this is a face.



Proposals

- ① **Stacked hierarchy:** Recursively stacking generative models on top of each other. If these models can be trained to optimality, then the bottom layer alone contains enough information to reconstruct the data distribution already, and the layers above can be ignored.
- ② **Architectural hierarchy:** Single layer latent variable models that prefers to place high-level features on certain part of the latent code, and low-level features in others (**Variational Ladder Autoencoder**).



Consider a family of latent variable models specified by a joint probability distribution $p_\theta(\mathbf{x}, \mathbf{z})$ over a set of observed variables \mathbf{x} and latent variables \mathbf{z} , parameterized by θ . Let $p_\theta(\mathbf{x})$ denote the marginal distribution of \mathbf{x} .

The objective is to maximize the marginal log-likelihood $p(\mathbf{x})$ over a dataset $\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$ drawn from some unknown underlying distribution $p_{\text{data}}(\mathbf{x})$:

$$\log p_\theta(\mathbf{X}) = \sum_{n=1}^N \log p_\theta(\mathbf{x}^{(i)})$$

which is non-convex and often intractable, as it involves marginalization over the latent variables \mathbf{z} .



Variational Autoencoders: For optimizing the intractable marginal likelihood, we optimize the evidence lower bound (ELBO) by introducing an inference model $q_\phi(z|x)$ parameterized by ϕ :

$$\begin{aligned}\log p((x)) &\geq \mathbb{E}_{q(z|x)}[\log p(x, z) - \log q(z|x)] \\ &= \mathbb{E}_{q(z|x)}[\log p(x, z)] - \mathcal{D}(q(z|x)||p(z)) = \mathcal{L}(x; \theta, \phi)\end{aligned}$$



Inference Suboptimality in Variational Autoencoders (ICML 2018)

<https://arxiv.org/pdf/1801.03558.pdf>

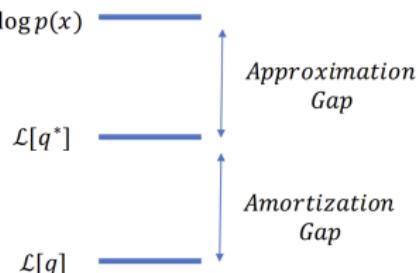


Figure 1. Gaps in Inference

Term	Definition	VAE Formulation
Inference	$\log p(x) - \mathcal{L}[q]$	$\text{KL}(q(z x) p(z x))$
Approximation	$\log p(x) - \mathcal{L}[q^*]$	$\text{KL}(q^*(z x) p(z x))$
Amortization	$\mathcal{L}[q^*] - \mathcal{L}[q]$	$\text{KL}(q(z x) p(z x)) - \text{KL}(q^*(z x) p(z x))$

Table 1. Summary of Gap Terms. The middle column refers to the general case where our variational objective is a lower bound on the marginal log-likelihood. The right most column demonstrates the specific case in VAEs. $q^*(z|x)$ refers to the optimal approximation within a family \mathcal{Q} , i.e. $q^*(z|x) = \arg \min_{q \in \mathcal{Q}} \text{KL}(q(z|x) || p(z|x))$.



Inference Suboptimality in Variational Autoencoders (ICML 2018)

<https://arxiv.org/pdf/1801.03558.pdf>

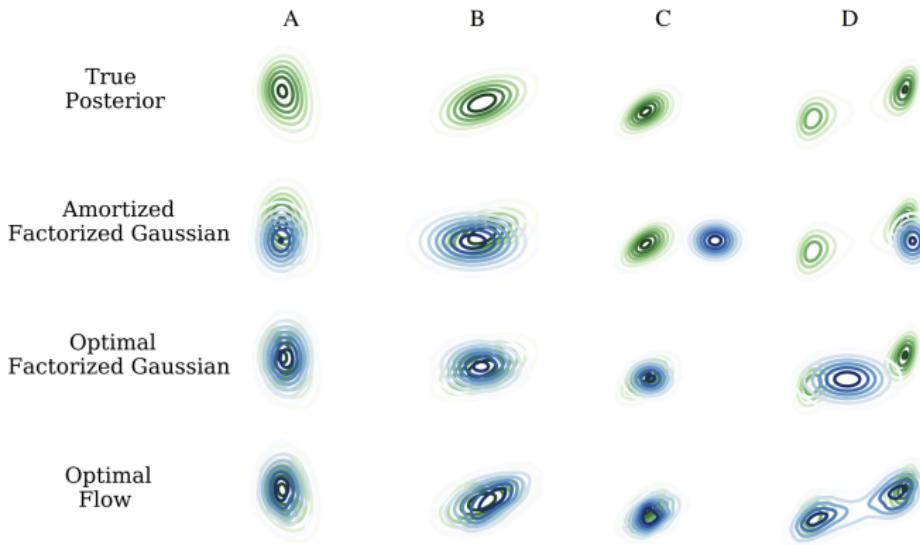


Figure 2. True Posterior and Approximate Distributions of a VAE with 2D latent space. The columns represent four different datapoints. The green distributions are the true posterior distributions, highlighting the mismatch with the blue approximations. Amortized: Variational parameters learned over the entire dataset. Optimal: Variational parameters optimized for each individual datapoint. Flow: Using a flexible approximate distribution.



Hierarchical Variational Autoencoders (HVAE): A series of VAEs stacked on top of each other. Hierarchy of latent variables $\mathbf{z} = \{\mathbf{z}_1, \dots, \mathbf{z}_L\}$, \mathbf{z}_1 represents the lowest layer closest to \mathbf{x} , and \mathbf{z}_L the top layer:

$$p(\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_L) = p(\mathbf{x}|\mathbf{z}_{>0}) \prod_{\ell=1}^{L-1} p(\mathbf{z}_\ell|\mathbf{z}_{>\ell}) p(\mathbf{z}_L)$$

where $\mathbf{z}_{>\ell}$ indicates $(\mathbf{z}_{\ell+1}, \dots, \mathbf{z}_L)$. Assume a Markov independence structure on the hidden variables, leading to the following simpler factorization:

$$p(\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_L) = p(\mathbf{x}|\mathbf{z}_1) \prod_{\ell=1}^{L-1} p(\mathbf{z}_\ell|\mathbf{z}_{\ell+1}) p(\mathbf{z}_L)$$



Hierarchical Variational Autoencoders (HVAE): For the inference distribution $q(\mathbf{z}|\mathbf{x})$, we do not assume any factorized structure. Denote $q(\mathbf{x}, \mathbf{z}) = p_{\text{data}}(\mathbf{x})q(\mathbf{z}|\mathbf{x})$. Both $p(\mathbf{x}|\mathbf{z})$ and $q(\mathbf{z}|\mathbf{x})$ are jointly optimized as maximizing the ELBO objective:

$$\mathcal{L}_{\text{ELBO}} = \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \left[\mathbb{E}_{q(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] - \mathcal{D}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) \right]$$

$$= \sum_{\ell=0}^L \mathbb{E}_{q(\mathbf{z}, \mathbf{x})} [\log p(\mathbf{z}_\ell | \mathbf{z}_{>\ell})] + H(q(\mathbf{z}|\mathbf{x}))$$

where we define $\mathbf{z}_0 \triangleq \mathbf{x}$, and H is the entropy of a distribution.



Representational efficiency

Proposition 1. $\mathcal{L}_{\text{ELBO}}$ is globally maximized as a function of $q(\mathbf{z}|\mathbf{x})$ and $p(\mathbf{x}|\mathbf{z})$ when $\mathcal{L}_{\text{ELBO}} = -H(p_{\text{data}}(\mathbf{x}))$. If $\mathcal{L}_{\text{ELBO}}$ is globally maximized for a Markov HVAE, the following Gibbs sampling chain converges to $p_{\text{data}}(\mathbf{x})$ if it is ergodic:

$$\mathbf{z}_1^{(t)} \sim q(\mathbf{z}_1|\mathbf{x}^{(t)})$$

$$\mathbf{x}^{(t+1)} \sim p(\mathbf{x}|\mathbf{z}_1^{(t)})$$

What does that mean?

One can sample from $p_{\text{data}}(\mathbf{x})$ without using $p(\mathbf{z}_\ell|\mathbf{z}_{>\ell})$, $1 \leq \ell < L$ at all.
The Gibbs chain only used the bottom layer of the model.



We **hope**:

- $p(z|x)$ is a probabilistic feature detector, and $q(z|x)$ as an approximation to it.
- q might learn hierarchical features similarly to a feed-forward network $x \rightarrow z_1 \rightarrow \dots \rightarrow z_L$, where higher layers correspond to higher level features (increasingly abstract).

But:

- If $q(z_{>\ell}|z_\ell)$ maps low-level features to high-level features, then the reverse mapping $q(z_\ell|z_{>\ell})$ maps high-level features into likely low-level sub-features.
- However, the only type of feature hierarchy we can hope to learn is one under which $q(z_\ell|z_{>\ell})$ is also Gaussian. **While the hierarchies we observe for feed-forward models require complex multimodal distributions to be captured.**



Variational Ladder Autoencoders (VLAE):

- Assume that if z_i is more abstract than z_j , then the inference mapping $q(z_i|x)$, and generative mapping when other layers are fixed $p(x|z_i, z_{\neg i} = z_{\neg i}^0)$ requires a more expressive network to capture.
- **No** hierarchical structure over the latent variables. VLAE encourages the latent code z_1, \dots, z_L to learn features with different levels of abstraction by carefully choosing the mapping $p(x|z)$ and $q(z|x)$.



Learning Hierarchical Features from Generative Models (14)

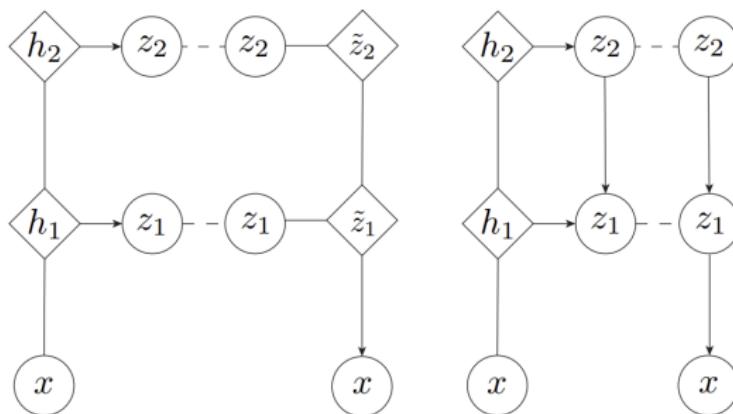


Figure 4. Inference and generative models for VLAE (left) and LVAE (right). Circles indicate stochastic nodes, and squares are deterministically computed nodes. Solid lines with arrows denote conditional probabilities; solid lines without arrows denote deterministic mappings; dash lines indicates regularization to match the prior $p(\mathbf{z})$. Note that in VLAE, we do not attempt to regularize the distance between \mathbf{h} and $\tilde{\mathbf{z}}$.



VLAE's Generative Network: $p(z) = p(z_1, \dots, z_L)$ is a simple prior on all latent variables, chosen as a standard Gaussian $\mathcal{N}(\mathbf{0}, \mathbf{1})$. The conditional distribution $p(x|z_1, \dots, z_L)$ is defined implicitly as:

$$\hat{z}_L = f_L(z_L)$$

$$\hat{z}_\ell = f_\ell(\hat{z}_{\ell+1}, z_\ell) \quad \ell = 1, \dots, L-1$$

$$x \sim r(x; f_0(\hat{z}_1)) \triangleq \mathcal{N}(f_0(\hat{z}_1), \mathbf{1})$$

where f_ℓ is a parameterized neural network:

$$\hat{z}_\ell = f_\ell(\hat{z}_{\ell+1}, z_\ell) = u_\ell([\hat{z}_{\ell+1}; v_\ell(z_\ell)])$$

where u_ℓ and v_ℓ are neural networks, and $[., .]$ denotes the concatenation operation of two vectors.



VLAE's Inference Network: For the inference network, we choose $q(z|x)$ as

$$\mathbf{h}_\ell = \mathbf{g}_\ell(\mathbf{h}_{\ell-1})$$

$$\mathbf{z}_\ell \sim \mathcal{N}(\mu_\ell(\mathbf{h}_\ell), \sigma_\ell(\mathbf{h}_\ell))$$

where $\ell = 1, \dots, L$ and \mathbf{g}_ℓ , μ_ℓ and σ_ℓ are neural networks, and $\mathbf{h}_0 = \mathbf{x}$.

VLAE's Learning:

$$\log p((\mathbf{x})) \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})]$$

$$= \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}, \mathbf{z})] - \mathcal{D}(q(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$



Improved Variational Inference with Inverse Autoregressive Flow (NIPS 2016)

Diederik P. Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya
Sutskever, Max Welling

<https://arxiv.org/pdf/1606.04934.pdf>



Let $p(\mathbf{x}, \mathbf{z})$ be the parametric **generative model**. Perform maximum marginal likelihood learning of its parameters, i.e. to maximize:

$$\log p(\mathbf{X} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}) = \sum_{i=1}^N \log p(\mathbf{x}^{(i)})$$

Introduce $q(z|\mathbf{x})$ a parametric **inference model** and optimize the **variational lower bound** on the marginal log-likelihood of each observation \mathbf{x} :

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(z|\mathbf{x})}[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}|\mathbf{x})] = \mathcal{L}(\mathbf{x}; \boldsymbol{\theta})$$

$\mathcal{L}(\mathbf{x}; \boldsymbol{\theta})$ is clearly a lower bound of $\log p(\mathbf{x})$:

$$\mathcal{L}(\mathbf{x}; \boldsymbol{\theta}) = \log p(\mathbf{x}) - \mathcal{D}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}))$$



Normalizing Flow (NF), introduced by Rezende and Mohamed (2015):

- Start with an initial random variable with a relatively simple distribution with known (and computationally cheap) probability density function: $\mathbf{z}_0 \sim q(\mathbf{z}_0|\mathbf{x})$
- Apply a chain of invertible parameterized transformation \mathbf{f}_t , such that the last iterate \mathbf{z}_T has a more flexible distribution: $\mathbf{z}_t = \mathbf{f}_t(\mathbf{z}_{t-1}, \mathbf{x})$

Probability density function of the last iterate:

$$\log q(\mathbf{z}_T|\mathbf{x}) = \log q(\mathbf{z}_0|\mathbf{x}) - \sum_{t=1}^T \log \det \left| \frac{d\mathbf{z}_t}{d\mathbf{z}_{t-1}} \right|$$

The authors proposed a limited family of **invertible** (why invertible?) transformations with known Jacobian determinant (\mathbf{u} , \mathbf{w} vectors and b non-linearity):

$$\mathbf{f}_t(\mathbf{z}_{t-1}) = \mathbf{z}_{t-1} + \mathbf{u} h(\mathbf{w}^T \mathbf{z}_{t-1} + b)$$



Inverse Autoregressive Flow (IAF):

- ① An initial encoder neural network output μ_0 and σ_0 , in addition to an extra output $\mathbf{h} \leftarrow$ inspired by LSTM.
- ② We draw a random sample $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$ and initialize the chain

$$\mathbf{z}_0 = \mu_0 + \sigma_0 \odot \epsilon$$

- ③ The flow consists of a chain of T of the following transformations:

$$\mathbf{z}_t = \mu_t + \sigma_t \odot \mathbf{z}_{t-1}$$

where at the t -th step of the flow, we use a different autoregressive neural network with inputs \mathbf{z}_{t-1} and \mathbf{h} , and outputs μ_t and σ_t .



Improved Variational Inference with Inverse Autoregressive Flow (4)

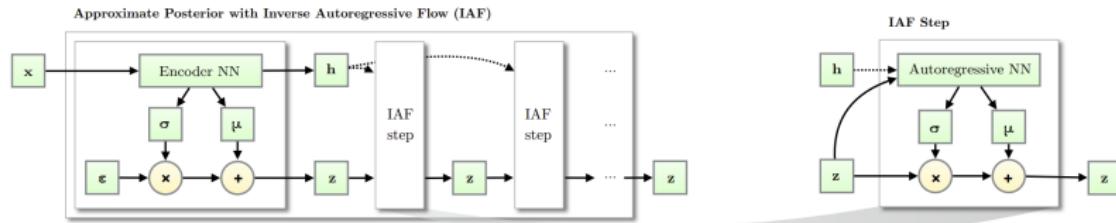


Figure 2: Like other normalizing flows, drawing samples from an approximate posterior with Inverse Autoregressive Flow (IAF) consists of an initial sample z drawn from a simple distribution, such as a Gaussian with diagonal covariance, followed by a chain of nonlinear invertible transformations of z , each with a simple Jacobian determinants.



Algorithm 1: Pseudo-code of an approximate posterior with Inverse Autoregressive Flow (IAF)

Data:

\mathbf{x} : a datapoint, and optionally other conditioning information

θ : neural network parameters

$\text{EncoderNN}(\mathbf{x}; \theta)$: encoder neural network, with additional output \mathbf{h}

$\text{AutoregressiveNN}[*](\mathbf{z}, \mathbf{h}; \theta)$: autoregressive neural networks, with additional input \mathbf{h}

$\text{sum}(\cdot)$: sum over vector elements

$\text{sigmoid}(\cdot)$: element-wise sigmoid function

Result:

\mathbf{z} : a random sample from $q(\mathbf{z}|\mathbf{x})$, the approximate posterior distribution

l : the scalar value of $\log q(\mathbf{z}|\mathbf{x})$, evaluated at sample ' \mathbf{z} '

$[\mu, \sigma, \mathbf{h}] \leftarrow \text{EncoderNN}(\mathbf{x}; \theta)$

$\epsilon \sim \mathcal{N}(0, I)$

$\mathbf{z} \leftarrow \sigma \odot \epsilon + \mu$

$l \leftarrow -\text{sum}(\log \sigma + \frac{1}{2}\epsilon^2 + \frac{1}{2}\log(2\pi))$

for $t \leftarrow 1$ **to** T **do**

$[\mathbf{m}, \mathbf{s}] \leftarrow \text{AutoregressiveNN}[t](\mathbf{z}, \mathbf{h}; \theta)$

$\sigma \leftarrow \text{sigmoid}(\mathbf{s})$

$\mathbf{z} \leftarrow \sigma \odot \mathbf{z} + (1 - \sigma) \odot \mathbf{m}$

$l \leftarrow l - \text{sum}(\log \sigma)$

end



Improved Variational Inference with Inverse Autoregressive Flow (6)

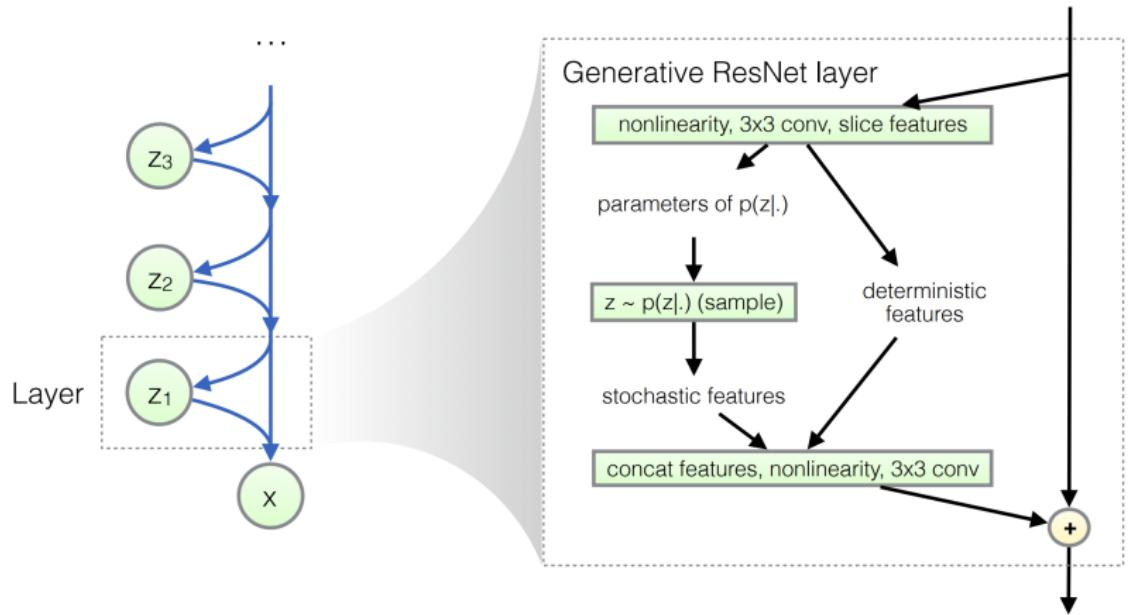


Figure 4: Generative ResNet and detail of layer. This is the generative component of our ResNet VAE.



NVAE: A Deep Hierarchical Variational Autoencoder

Arash Vahdat, Jan Kautz (NVIDIA)

<https://arxiv.org/pdf/2007.03898.pdf>



NVAE: A Deep Hierarchical Variational Autoencoder (1)

I believe there must be someones who look (ever looked) like these generated pictures!



Figure 1: 256×256-pixel samples generated by NVAE, trained on CelebA HQ [28].



NVAE: A Deep Hierarchical Variational Autoencoder (2)

The first row is not based on any scientific evidence, but the Bible (!?)

Year	Population	Births per 1,000	Births Between Benchmarks	Number Ever Born	Percent of Those Ever Born
50,000 B.C.E.	2	-	-	-	-
8000 B.C.E.	5,000,000	80	1,137,789,769	1,137,789,769	0.4
1 C.E.	300,000,000	80	46,025,332,354	47,163,122,125	0.6
1200	450,000,000	60	26,591,343,000	73,754,465,125	0.6
1650	500,000,000	60	12,782,002,453	86,536,467,578	0.6
1750	795,000,000	50	3,171,931,513	89,708,399,091	0.9
1850	1,265,000,000	40	4,046,240,009	93,754,639,100	1.3
1900	1,656,000,000	40	2,900,237,856	96,654,876,956	1.7
1950	2,516,000,000	31-38	3,390,198,215	100,045,075,171	2.5
1995	5,760,000,000	31	5,427,305,000	105,472,380,171	5.5
2011	6,987,000,000	23	2,130,327,622	107,602,707,793	6.5
2019	7,692,000,000	19	1,157,835,998	108,760,543,791	7.1
2030	8,932,000,000	16	2,226,104,131	110,986,647,922	8.0
2050	9,854,000,000	14	2,120,539,877	113,107,187,799	8.7

<https://www.prb.org/howmanypeoplehaveeverlivedonearth/>
(Data from United Nations)



Summary:

- This paper does **not** propose new model, but a carefully designed neural architecture for hierarchical VAEs (named Nouveau VAE or NVAE):
 - ① Depth-wise separable convolutions and batch normalization.
 - ② Residual parameterization of Normal distributions.
 - ③ Spectral regularization to stabilized the training.
- Hierarchical VAEs (Max Welling's work):
 - ① Normalizing flows
 - ② Autoregressive models
 - ③ Deep energy-based models



Insights:

- ① VAEs often respond differently to the over-parameterization in neural networks. Over-parameterization the decoder network may hurt the test log-likelihood, whereas powerful encoders can yield better models because of reducing the amortization gap.

Inference Suboptimality in Variational Autoencoders (ICML 2018)

<https://arxiv.org/pdf/1801.03558.pdf>

- ② VAEs should model long-range correlations in data, requiring the networks to have large receptive fields.



Inspired by **Inverse Autoregressive Flow** or IAF (Paper 2), but not exactly:

- To increase the expressiveness of both the approximate posterior and prior, the latent variables are partitioned into disjoint groups $\mathbf{z} = \{\mathbf{z}_1, \dots, \mathbf{z}_T\}$.
 - The prior is represented by $p(\mathbf{z}) = \prod_t p(\mathbf{z}_t | \mathbf{z}_{<t})$
 - The approximate posterior is represented by $p(\mathbf{z} | \mathbf{x}) = \prod_t q(\mathbf{z}_t | \mathbf{z}_{<t}, \mathbf{x})$
 - Each conditional $p(\mathbf{z}_t | \mathbf{z}_{<t})$ and $q(\mathbf{z}_t | \mathbf{z}_{<t}, \mathbf{x})$ are represented by factorial Normal distributions.

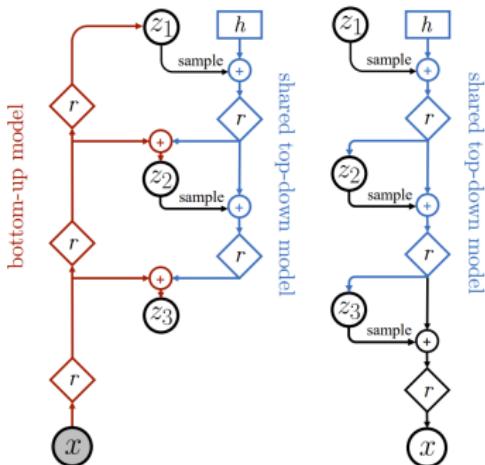
Variational lower bound $\mathcal{L}_{VAE}(\mathbf{x})$ on $\log p(\mathbf{x})$ as:

$$\mathcal{L}_{VAE}(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] - \mathcal{D}(q(\mathbf{z}_1|\mathbf{x})||p(\mathbf{z}_1))$$

$$-\sum_{t=2}^T \mathbb{E}_{q(\mathbf{z}_{<t}|\mathbf{x})}[\mathcal{D}(q(\mathbf{z}_t|\mathbf{z}_{<t}, \mathbf{x}))||p(\mathbf{z}_t|\mathbf{z}_{<t})]$$



NVAE: A Deep Hierarchical Variational Autoencoder (6)

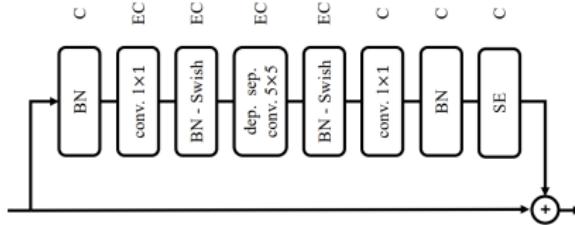


(a) Bidirectional Encoder (b) Generative Model

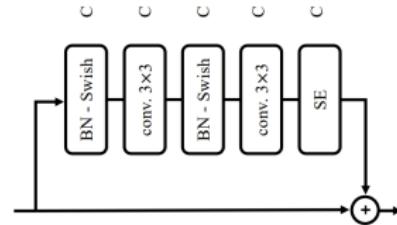
Figure 2: The neural networks implementing an encoder $q(\mathbf{z}|\mathbf{x})$ and generative model $p(\mathbf{x}, \mathbf{z})$ for a 3-group hierarchical VAE. \diamond_r denotes residual neural networks, \oplus denotes feature combination (e.g., concatenation), and $[\mathbf{h}]$ is a trainable parameter.



NVAE: A Deep Hierarchical Variational Autoencoder (7)



(a) Residual Cell for NVAE Generative Model



(b) Residual Cell for NVAE Encoder

Figure 3: The NVAE residual cells for generative and encoder models are shown in (a) and (b). The number of output channels is shown above. The residual cell in (a) expands the number of channels E times before applying the depthwise separable convolution, and then maps it back to C channels. The cell in (b) applies two series of BN-Swish-Conv without changing the number of channels.



Hierarchical Implicit Models and Likelihood-Free Variational Inference (NIPS 2017)

Dustin Tran, Rajesh Ranganath, David M. Blei

<https://arxiv.org/pdf/1702.08896.pdf>



What are implicit probabilistic models?

Implicit probabilistic models are a flexible class of models defined by a simulation process of data.

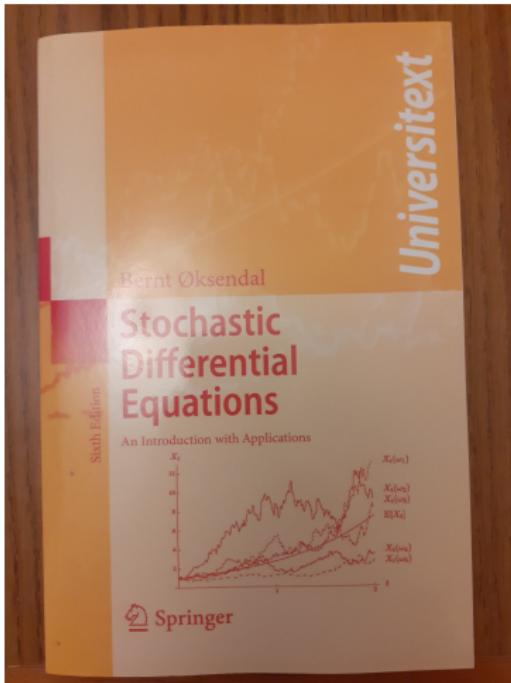
Example: Given initial conditions, physical simulators describe a stochastic process that generates data. In population ecology, the Lotka-Volterra model simulates predator-prey populations ($x_1, x_2 \in \mathbb{R}^+$) over time via a **stochastic differential equation**:

$$\frac{dx_1}{dt} = \beta_1 x_1 - \beta_2 x_1 x_2 + \epsilon_1, \quad \epsilon_1 \sim \mathcal{N}(0, 10)$$

$$\frac{dx_2}{dt} = -\beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon_2, \quad \epsilon_2 \sim \mathcal{N}(0, 10)$$

where Gaussian noises ϵ_1, ϵ_2 are added at each full time step.





Just bought this book



Proposals

- ① **Hierarchical implicit models** (HIMs): combines the idea of implicit densities with hierarchical Bayesian modeling, and defines models via simulators of data.
 - ② **Likelihood-free variational inference** (LFVI): a scalable variational inference algorithm for HIMs.

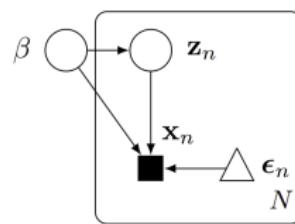
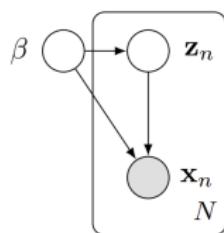


Figure 1: (left) Hierarchical model, with local variables \mathbf{z} and global variables β . **(right) Hierarchical implicit model.** It is a hierarchical model where \mathbf{x} is a deterministic function (denoted with a square) of noise ϵ (denoted with a triangle).



Hierarchical models:

$$p(\mathbf{x}, \mathbf{z}, \boldsymbol{\beta}) = p(\boldsymbol{\beta}) \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\beta}) p(\mathbf{z}_n | \boldsymbol{\beta})$$

where \mathbf{x}_n is an observation, \mathbf{z}_n is the latent variable associated to that observation (local variables), and $\boldsymbol{\beta}$ is the latent variable shared across observations (global variable).

Hierarchical models typically use a tractable likelihood $p(\mathbf{x}_n | \mathbf{z}_n, \boldsymbol{\beta})$. But many likelihoods of interest, such as simulator-based models, and generative adversarial networks do not admit a tractable likelihood. Thus, we develop **hierarchical implicit models** (HIMs) that define a function g that takes random noise $\epsilon_n \sim s(\cdot)$ and outputs \mathbf{x}_n , given \mathbf{z}_n and $\boldsymbol{\beta}$:

$$\mathbf{x}_n = g(\epsilon_n | \mathbf{z}_n, \boldsymbol{\beta}), \quad \epsilon_n \sim s(\cdot)$$



The implicit likelihood of $x_n \in A$ given z_n and β is:

$$\mathcal{P}(x_n \in A | z_n, \beta) = \int_{\{g(\epsilon_n | z_n, \beta) = x_n \in A\}} s(\epsilon_n) d\epsilon_n$$

is typically intractable. Therefore, we need **likelihood-free variational inference (LFVI)**.

The basic idea is to avoid computing the likelihood explicitly, and instead we compute the unbiased gradients with Monte Carlo. **We skip the technical details here.**

My observation

Many papers of David Blei share the same idea of local-global decomposition from **Latent Dirichlet Allocation (LDA)**.

Deep Hierarchical Implicit Models

Dustin Tran, Rajesh Ranganath, David M. Blei
<https://arxiv.org/pdf/1702.08896v1.pdf>



Stacking layers of implicit latent variables produces a **deep implicit model** (DIM). This is basically a trivial extension of the previous one (same authors).

DIMs share the same local-global decomposition as hierarchical implicit models. The key difference is that instead of a single local variable z_n , DIMs define layers of implicit local variables (This is the idea from Hierarchical VAE). Formally, for layers $1, \dots, T$, the local variables for data point x_n are generated via:

$$z_{n,T} = g_T(\epsilon_{n,T} | \beta_T) \quad \epsilon_{n,T} \sim s(.)$$

$$z_{n,t} = g_t(\epsilon_{n,t} | z_{n,t+1}, \beta_t) \quad \epsilon_{n,t} \sim s(.) \quad t < T$$

where g_t is the transformation for layer t , s is random noise, and $\{\beta_1, \dots, \beta_T\}$ are parameters of the T transformations.



Deep Hierarchical Implicit Models (2)

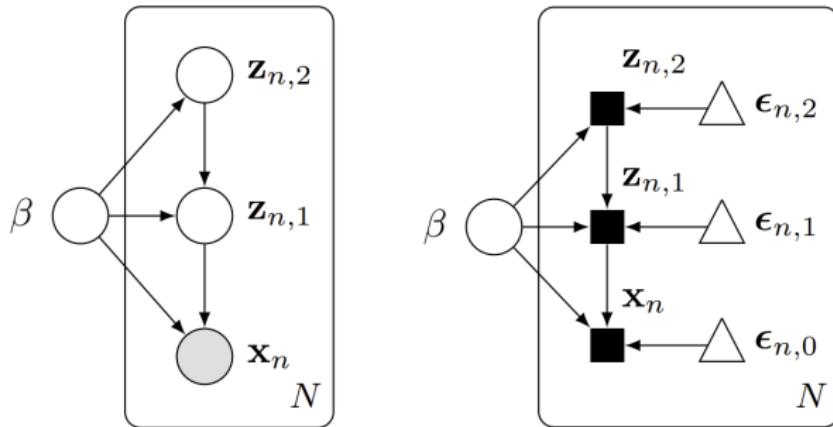


Figure 2. (left) Deep generative model with 2 layers $\{\mathbf{z}_1, \mathbf{z}_2\}$ and global variables β . **(right) Deep implicit model.** It is a deep generative model with implicit likelihood and implicit priors.



Paper 6

Discrete Object Generation with Reversible Inductive Construction

Ari Seff, Wenda Zhou, Farhan Damani, Abigail Doyle, Ryan P. Adams

<https://arxiv.org/pdf/1907.08268.pdf>



Proposal

Generative model (interpreted from denoising autoencoders) for discrete objects (e.g. molecules, source code, graphs, etc.) employing a **Markov chain** where transitions are restricted to a set of **local operations that preserve validity**:

- ① a sequence of corrupted objects that are valid but not from the data distribution.
- ② a learned reconstruction distribution that attempts to fix the corruptions while also preserving validity.

My opinion

What is the difference between this and the **Reinforcement Learning** approach (Markov decision process) then?

The idea is pretty much similar to generating a solvable training example for **Rubik**:

- ① Start from the original state x , we begin to corrupt it iteratively for k steps. Let say operations $s = [s_1, s_2, \dots, s_k]$ done, and result in state \bar{x} .
- ② Each operation is reversible. The solution to solve state \bar{x} back into the original state x is a reverse sequence $[s_k^{-1}, \dots, s_2^{-1}, s_1^{-1}]$.



Let $p(x)$ be an unknown probability mass function. Our goal is to learn a generative model $p_\theta(x)$, approximating $p(x)$. Let $c(\bar{x}|x)$ be a fixed corrupting distribution, and $p_\theta(x|\bar{x})$ be the learned reconstruction distribution:

$$p(x, \bar{x}) = c(\bar{x}|x)p(x)$$

We use supervised learning to train a reconstruction distribution model $p_\theta(x|\bar{x})$ to approximate $p(x|\bar{x})$.

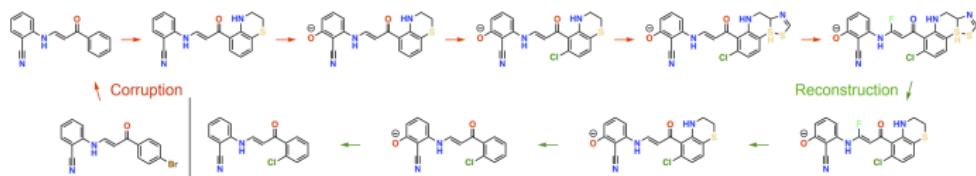


Figure 2: Corruption and subsequent reconstruction of a molecular graph. Our method generates discrete objects by running a Markov chain that alternates between sampling from fixed corruption and learned reconstruction distributions that respect validity constraints.



Discrete Object Generation with Reversible Inductive Construction (4)

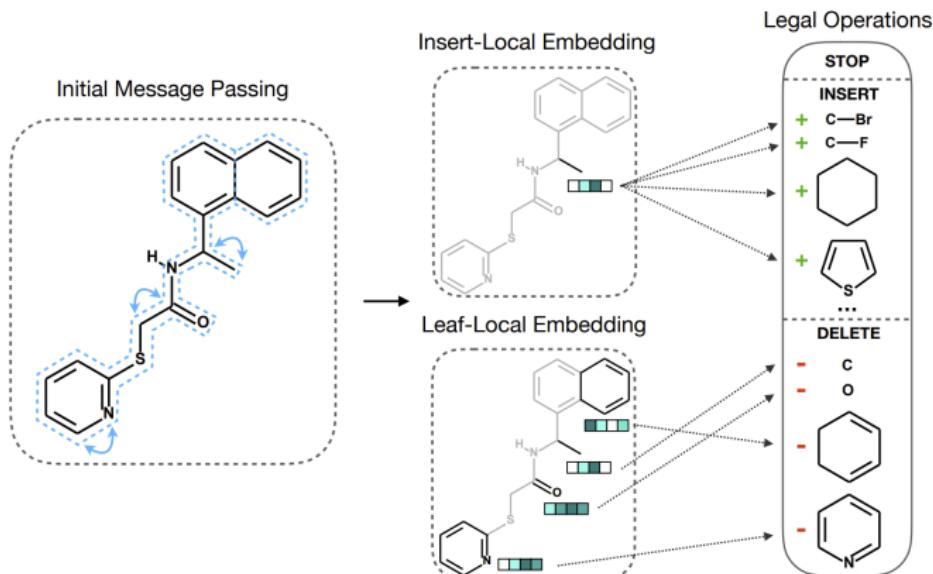


Figure 1: Reconstruction model processing given an input molecule. Location-specific representations computed via message passing are passed through fully-connected layers outputting probabilities for each legal operation.



The probability of arriving at corrupted sample \bar{x} from x is:

$$c(\bar{x}|x) = \sum_s c(\bar{x}, s|x) = \sum_{s \in S(x, \bar{x})} c(s|x)$$

where $S(x, \bar{x})$ denotes the set of all corrupting sequences from x to \bar{x} , and $c(s|x)$ is a fixed conditional distribution over a sequence of corrupting operations $s = [s_1, \dots, s_k]$ in which $s_i \in \text{Ind}(\bar{x}_i)$ – the set of legal inductive moves for a given \bar{x}_i . The joint data-generating distribution:

$$p(x, s, \bar{x}) = c(\bar{x}, s|x)p(x)$$

where $c(\bar{x}, s|x) = 0$ if $s \notin S(x, \bar{x})$.



Optimization problem:

$$\theta^* = \arg \min_{\theta} \mathbb{E}_{p(x, s, \bar{x})} \mathcal{D}_{\text{KL}}(p(s, x | \bar{x}) || p_{\theta}(s, x | \bar{x}))$$

which is to maximize likelihood estimation $p_{\theta}(s, x | \bar{x})$ and likewise $p_{\theta}(x | \bar{x})$. For the expectation over the joint data-generating distribution $p(x, s, \bar{x})$, we can sample a corruption sequence as follows:

$$x \sim p(x)$$

$$\bar{x}, s \sim c(\bar{x}, s | x)$$

The learned reconstruction sequence model can be factorized (represented) as a product of memoryless transitions culminating with a stop token:

$$p_{\theta}(s^{-1} | \bar{x}) = p_{\theta}(\text{stop} | x) p_{\theta}(x | \bar{x}_1) \prod_{i=1}^{k-1} p_{\theta}(\bar{x}_i | \bar{x}_{i+1})$$

where $s^{-1} = [s_k^{-1}, \dots, s_1^{-1}, \text{stop}]$.

