

## \* Introduction.

VAE는 prior, tighter bound, posterior estimate 층층을 갖는다.

AEC는 classification이나 대별을 고려하지 않기, 기본적으로 VAE의 arch은 요구 사항이 조금 다르다.

VAE는 input과 latent 간의 상호 정보 maximize.

↳ input의 특징에는 판별하는 classification의 차를

Decoder generator가 대응되며, overfit 되면 안되지만, encoder는 험난하다.

encoder-based의 marginal LL처럼 encoder overfit이 민감하게 된다.

2개의 receptive field가 필요. ELBO 때문에 불안정 (일반 VAE는 2개의 MN이상)

CE2L의 NVAE 개선. → non-autoregressive likelihood. 풍 SOTA

↳ depthwise conv 사용.

↳ BNE 풍도, hierarchical 이론적 속도, unstable.

↳ KL minimize를 위한 residual parameterization은 spectral Norm이 stable한 풍도.

즉, 새로운 계층 VAE 개선. (NVAE)

而后은 posterior approximel residual parameterization.

SN으로 stable.

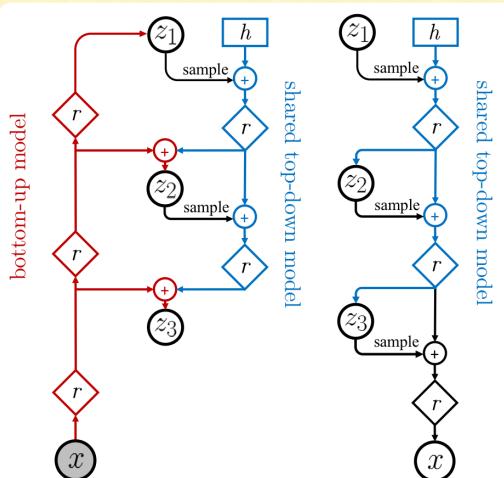
제일의 부정을 줄이는 solution.

기본 VAE의 MN 차수도 성능↑.

IAF-VAE의 한계 있지만, 성능↑.

BIVIA보다 전연성 있지만 성능↑.

## \* Related Work.



(a) Bidirectional Encoder (b) Generative Model

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

## \* Method.

### NVAE challenge.

1) VAE의 불안정한 네트워크.

2) Scale up, stability를 유지하는 hierarchical.

### \* Residual cell for VAE

VAE의 불안정한 long-term dependency.

NVAE 같은 Unconditional하게 hierarchical로 global長い time delay.

### \* Residual Cell for Generative Model.

kernel size를 늘리는 것은 receptive  $\leftrightarrow$  param trade off.

Depthwise로 쭉 풀면, channel 별 개별이기 때문에 효율적.

∴ MobileNet V2 와 같이,  $1 \times 1$  channel conv는 짧음.

- BN.

이 BN의 경우 학습 단계를 매우 빠르다.

단지 몇몇 Terning 단계를 포함한, 훨씬 더 적은 계산 단계로 대체됨.

∴ BN의 momentum을 조절하고, 훨씬 빠른 학습 단계로 대체되는 scaling parameterization의 reg.

- Swish.

- Squeeze and Excitation

- Final cell.

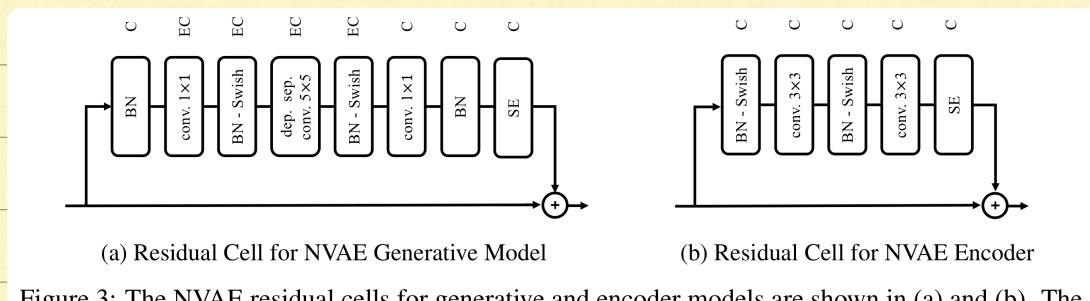


Figure 3: The NVAE residual cells for generative and encoder models are shown in (a) and (b). The number of output channels is shown above. The residual cell in (a) expands the number of channels  $E$  times before applying the depthwise separable convolution, and then maps it back to  $C$  channels.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

\* Residual cell for Encoder.

Depthwise bottom-up Encoder는 이런 식

∴ 같은 채널

인 경우 같은 BN-conv-activate 단계를 짧음.

\* Reducing Memory require.

Depthwise는 메모리 짧아 짧음

channel level batch

1) APEXel mixed precision.

2) 그림 3번의 각 feature map의 뒤로 backward를 하여 가능함.

다만 모델이 초기화 되어 GPU switch를 허용하지, backward의 다른 batched feature map은 가능.

$L_i$  = gradient checking points.

\* Taming the Unbounded KL Term

Hierarchical prior  $p(z_{<i}|x)$ 는 고려된다.

Encoder dec가 훈련 중 나온 먼 distribution을 생성하면 KL이 무한대가 가능해 update 불가능, unstable하다.

- Residual Normal Distribution.

$$p(z_i^i|z_{<i}) = N(\mu_i(z_{<i}), \sigma_i(z_{<i}))$$

$$q(z_i^i|z_{<i}, x) = N(\mu_i(z_{<i}) + \Delta\mu_i(z_{<i}, x), \sigma_i(z_{<i}) \cdot \Delta\sigma_i(z_{<i}, x))$$

→ prior와 posterior approximation의 scale, location.

→ prior와 posterior의 업데이트.

$$\text{KL}(q(z^i|x)||p(z^i)) = \frac{1}{2} \left( \frac{\Delta\mu_i^2}{\sigma_i^2} + \Delta\sigma_i^2 - \log \Delta\sigma_i^2 - 1 \right),$$

↳ Decoder의 의존  $\sigma_i$ 가 below bound 일 때, encoder의 의존한 CL.

↳ Weight averaging prior, posterior의 같은 맥락.

- Spectral Regularization (SR)

KL은 bound 하지 못해, encoder의 훈련은 input의 차원으로 규제하기 바람지 않아야 한다.

↳ Lipschitz regularization.

∴ SR 적용

$$L_{SR} = \lambda \sum_i s^{(i)}, s^{(i)} = \text{largest singular value}.$$

- More expressive Approximate Posteriors with Normalizing flow.

NVAE는  $p(z)$ 와  $q(z|x)$ 의 그룹간 autoregressive 분포와 각 그룹의 독립적 분포로 modeling

celeba 빙글 sampling 위한 representation.

∴ 추가적인 normalization flow 추가.

encoder 비단 적용도(?)가 때문에 explicit inversion이 필요없고, IAF 사용 가능. Sampling 시간이 증가하지 않음.

## \* Experiments.

### \* Main quantitative results.

Table 1: Comparison against the state-of-the-art likelihood-based generative models. The performance is measured in bits/dimension (bpd) for all the datasets but MNIST in which negative log-likelihood in nats is reported (lower is better in all cases). NVAE outperforms previous non-autoregressive models on most datasets and reduces the gap with autoregressive models.

Method	MNIST 28×28	CIFAR-10 32×32	ImageNet 32×32	CelebA 64×64	CelebA HQ 256×256	FFHQ 256×256
NVAE w/o flow	<b>78.01</b>	2.93	-	2.04	-	0.71
NVAE w/ flow	78.19	<b>2.91</b>	<b>3.92</b>	<b>2.03</b>	<b>0.70</b>	<b>0.69</b>
<b>VAE Models with an Unconditional Decoder</b>						
BIVA [36]	78.41	3.08	3.96	2.48	-	-
IAF-VAE [4]	79.10	3.11	-	-	-	-
DVAE++ [20]	78.49	3.38	-	-	-	-
Conv Draw [42]	-	3.58	4.40	-	-	-
<b>Flow Models without any Autoregressive Components in the Generative Model</b>						
VFlow [59]	-	2.98	-	-	-	-
ANF [60]	-	3.05	3.92	-	0.72	-
Flow++ [61]	-	3.08	<b>3.86</b>	-	-	-
Residual flow [50]	-	3.28	4.01	-	0.99	-
GLOW [62]	-	3.35	4.09	-	1.03	-
Real NVP [63]	-	3.49	4.28	3.02	-	-
<b>VAE and Flow Models with Autoregressive Components in the Generative Model</b>						
δ-VAE [25]	-	2.83	3.77	-	-	-
PixelVAE++ [35]	78.00	2.90	-	-	-	-
VampPrior [64]	78.45	-	-	-	-	-
MAE [65]	77.98	2.95	-	-	-	-
Lossy VAE [66]	78.53	2.95	-	-	-	-
MaCow [67]	-	3.16	-	-	0.67	-
<b>Autoregressive Models</b>						
SPN [68]	-	-	3.85	-	0.61	-
PixelSNAIL [34]	-	2.85	3.80	-	-	-
Image Transformer [69]	-	2.90	3.77	-	-	-
PixelCNN++ [70]	-	2.92	-	-	-	-
PixelRNN [41]	-	3.00	3.86	-	-	-
Gated PixelCNN [71]	-	3.03	3.83	-	-	-

= Qualitative Results.

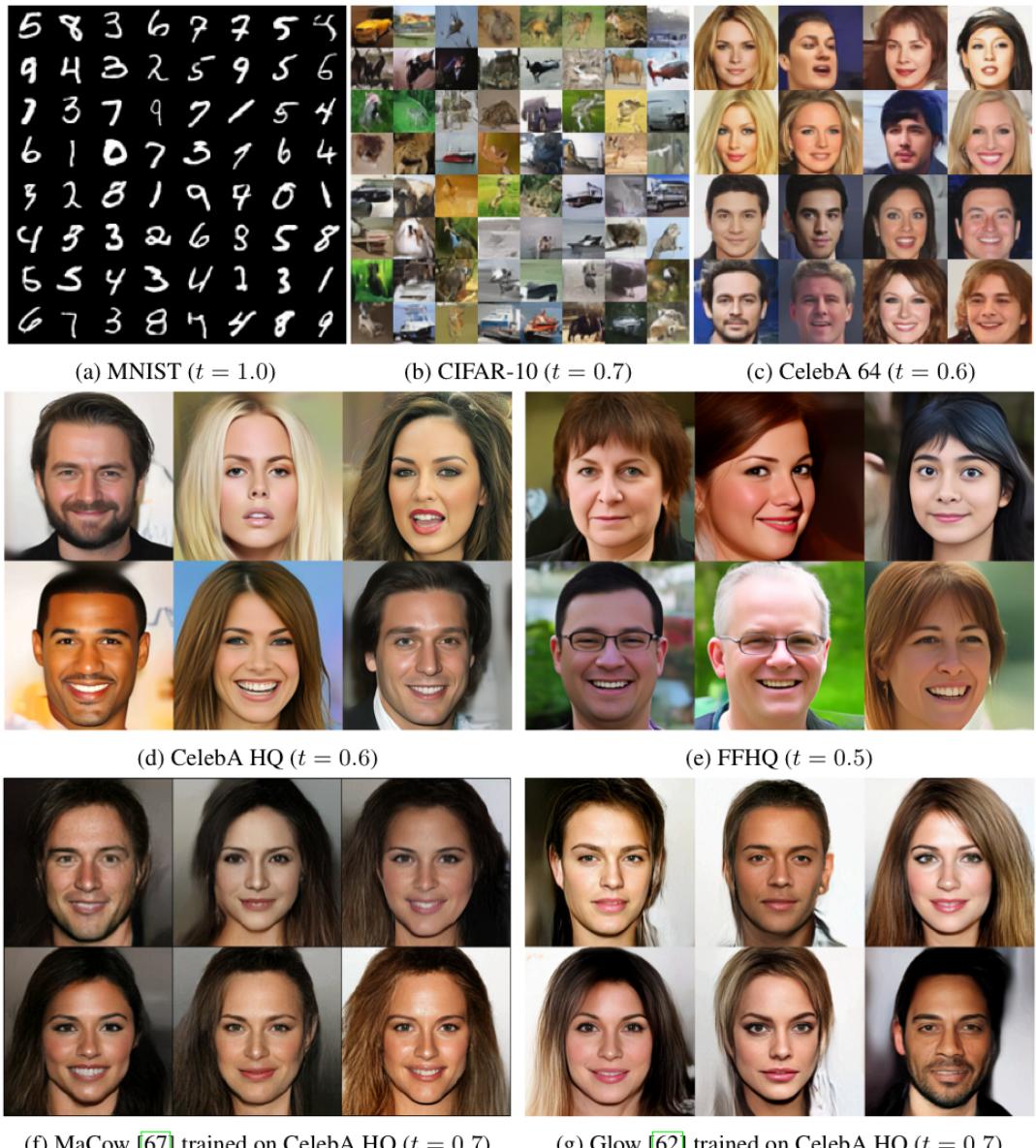


Figure 4: (a)-(e) Sampled images from NVAE with the temperature in prior ( $t$ ). (f)-(g) A few images generated by MaCow [67] and Glow [62] are shown for comparison (images are from the original publications). NVAE generates diverse high quality samples even with a small temperature, and it exhibits remarkably better hair details and diversity (best seen when zoomed in).

\* Ablation studies.

- Norm and Activation flow.

Table 2: Normalization & activation

Functions	$L = 10$	$L = 20$	$L = 40$
WN + ELU	3.36	3.27	3.31
BN + ELU	3.36	3.26	3.22
BN + Swish	<b>3.34</b>	<b>3.23</b>	<b>3.16</b>

- Residual Cell.

Table 3: Residual cells in NVAE

Bottom-up model	Top-down model	Test (bpd)	Train time (h)	Mem. (GB)
Regular	Regular	3.11	43.3	6.3
Separable	Regular	3.12	49.0	10.6
Regular	Separable	<b>3.07</b>	48.0	10.7
Separable	Separable	<b>3.07</b>	50.4	14.9

- Residual Normal Distribution.

Residual은 NLL posterior collapse의 나쁜 영향을 가진다.

- The effect of SR and SE

Table 5: SR & SE

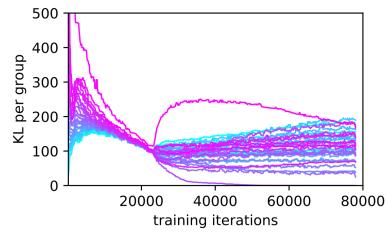
Model	Test NLL
NVAE	<b>3.16</b>
NVAE w/o SR	3.18
NVAE w/o SE	3.22

- Sampling Speed.

빠름.

### - Reconstruction

느낌.



(a) Reconstruction results (best seen when zoomed in).

(b) Average KL per group.

Figure 5: (a) Input input on the left and reconstructed images on the right for CelebA HQ. (b) KL per group on CIFAR-10.

### \* Conclusion.

#### NVAE 장점

↳ Generator에는 depthwise 를 사용하고, Encoder에는 일반 residual 사용.

↳ Encoder normal distribution의 residual parameterization 과, stable 를 위한 SN 적용

↳ 또한 memory 줄임.

↳ 고화질의 큰 image,

↳ 신축적인 architecture 지원.