

long-range sequence w/ transformer-L SOT/LM.

CNN의 inductive bias는 model을 흔히 찾지만, long-range interactive의 비현실적인 compute cost가 ECL.

∴ CNN과 Transformer의 장점 결합하고 싶음

↳ CNN으로 context $\frac{2}{3}$, transformer $\frac{1}{3}$ 구성을 드로일 뿐은.

Conditional句도 동작을 잘합니다.

* Introduction.

Transformer는 NLP ain SOTA고, 다른 분야에도 적용되고 있다.

(CNN과 다른 이유인 inductive bias)가 없기 때문에 입력의 complex한 관계를 학습할 수 있다.

그러나, 이때문에 모든 pair간 계산이 필요해 quadratic cost가 된다.

∴ 고체수송의 예는 1인당 100kg을 많이 이동다.

Transformer가 Conv의 구조를 학습하는 경향이 있다.

-> 10 epoch마다 image의 구조를 scratch로 다시 배운다.

- 'oldboy transformer'이 무엇을 알고 있는지 image의 inductive bias를 활용해 encoding 가능할까?

기저는 low-level image structure에 CMU 등은 흐르게 이지하는 structural assumption에 대해서는 아울러 같다.

즉 CNN은 locality inductive bias를 갖는다. spatial invariance를 갖는다. (weight sharing)

CNN은 image를 읽기 위한 `read_hdf5`을 만들었고, 그나마 네트워크를 만들 때 `transformer`을 활용된다.

这些 dice 的 local part 的 low-level distributions 与 capture 之数的分布 D₂ 完全相同。

72(2) conditional Σ 가정합니다.

* Related Work

- The Transformer Family

Transformer는 attention을 input 사이의 상호작용에 사용된다.

NLP에서 token은 input을 텍스트, image는 audio는 비슷한 signal로 취급된다.

Transformer는 attention을 fc로 구현된다.

Attention은 Q, k, V로 구성된다.

$$Q \in \mathbb{R}^{N \times d_k}, k \in \mathbb{R}^{N \times d_k}, V \in \mathbb{R}^{N \times d_v}$$

$$\text{Attn}(Q, k, V) = \text{softmax}\left(\frac{Qk^T}{\sqrt{d_k}}\right)V$$

Attention은 모든 input seq 사이 계산량이 quadratic cost가 된다.

↳ long-range interaction은 주제와 속도의 대신 trade-off가 있다.

이런 문제를 해결하기 위해 attention receptive field를 제한하는데, pixel independent assumption을 넣는다.

- Convolutional approach.

CNN은 image size를 줄여 compute cost가 linear이지만, kernel은 quadratic인 kernel size를 갖는다.

Receptive field는 computation 효율성이지만, 계산량이 매우 많다.

low-res 이미지 transformer보다 더 강한 성능을 보임.

∴ 두 층을 섞기 병.

- Two-stage approach.

먼저 image를 encoding하고, encoding된 대로 probability model을 만들고자.

VAE를 사용하여, data의 representation을 학습하고, 그걸 VAE의 블록을 학습하는 것의 이점을 넣어논문이 있다.

2-stage의 conditional normalizing flow와 unconditional normalizing flow가 비슷한 성능을 보인다.

GAN의 훈련흐름을 살피기 위해, AE의 representation을

각각의 G에 대한 image를 decoding 시킨 low-resolution wavelet coefficients를 GAN에 학습한다.

VQVAE에서 image의 discrete한 representation을 학습하기 위해 convol 층이 활용했다.

제작적인 손으로 VQVAE를 발전시킨 model도 있다.

이제는 long-range의 다른 물체가 있는 scene의 영상이 더운다.

32x32 image를 encoding할 때 cost가 부담스러워졌기 때문에, VQVAE는 96x96를 encoding 한다.

가장 중요한 작은 pixel의 다른 spatially invariant 차이로 작은 receptive field를 갖는 shallow VQVAE를 찾다.

이제 encoding의 compute cost를 줄이기 위해 큰 receptive field를 encode 했으므로 decoder의 fine feature를

잘 잡기 위해 작은 receptive field를 쓰는 편.

이제는 자체 content를 보는 encoder와 transformer를 합친다면,

high-res를 합성할 수 있다고 한다.

* Approach.

Quadratic는 compute cost 때문이 Transformer를 활용한 Vision은 high-res로 다루고 싶었다.

기본상도 이미지는 전처리 맵을 이해해야 한다. 따라서 pixel 단위별로 인지적으로 흡사한 codebook을 사용한다.

codebook은 시청자 composition length를 두는 데 있고, transformer로 global한 interactive를 효율적으로 학습할 수 있다.

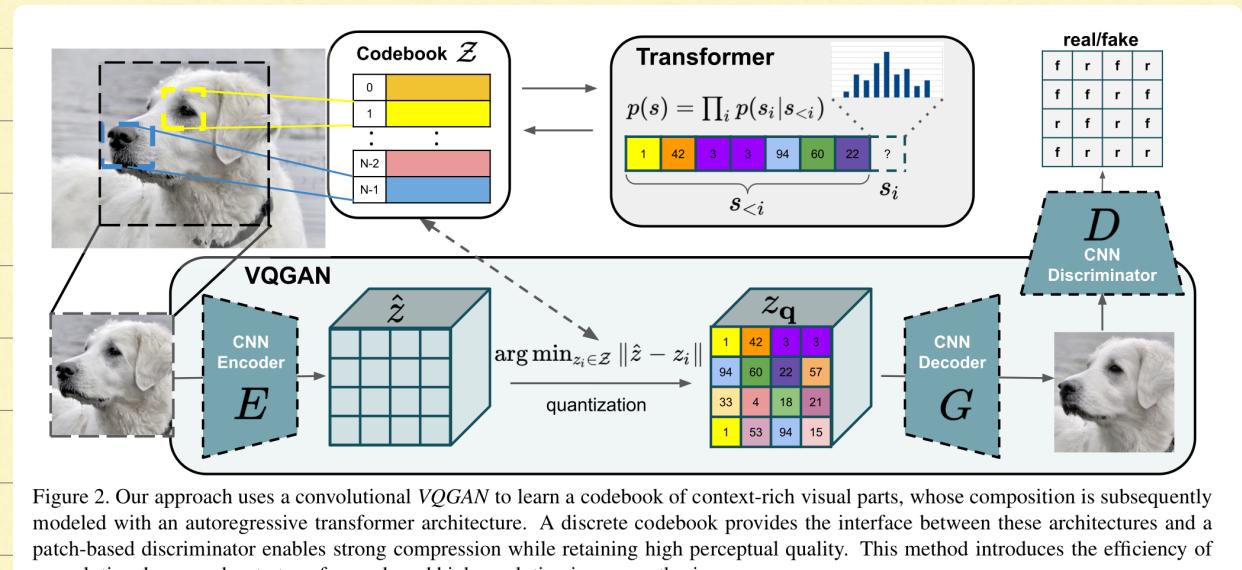


Figure 2. Our approach uses a convolutional VQGAN to learn a codebook of context-rich visual parts, whose composition is subsequently modeled with an autoregressive transformer architecture. A discrete codebook provides the interface between these architectures and a patch-based discriminator enables strong compression while retaining high perceptual quality. This method introduces the efficiency of convolutional approaches to transformer based high resolution image synthesis.

- Learning an Effective codebook of Image Constituents for Use in Transformer.

Transformer를 사용하면서 sequence를 input으로 받아야 한다.

pixel sequence를 놓칠까봐 받는 대신, codebook을 뽑는다. $Z_q \in \mathbb{R}^{h \times w \times n_z}$

\Rightarrow VQVAE 방식

단지, E와 G를 학습시킨다.

↳ E로 x를 만들고, G로부터 input x의 버전인 $\hat{x} = G(Z_q)$ 를 생성한다.

$$Z_q = q(\hat{z}) = \left(\arg \min_{z_k \in Z} \| \hat{z}_k - z_k \| \right) \in \mathbb{R}^{h \times w \times n_z}$$

$$\hat{x} = G(Z_q) = G(q(E(x)))$$

미분 불가능한 Q의 경계 backpropagation은 straight-through gradient estimate로 구현된다.

\hookrightarrow decoder enc로 gradient copy.

$$L_{VAE}(E, G, Z) = \|x - \hat{x}\|_1 + \|s_g[E(x)] - z_q\|_1 + \beta \|s_g[z_q] - E(x)\|_1$$

↑
Irec
↑
step-gradient.
↑
commitment loss

* Learning a perceptual rich codebook

latent의 구조로 image representation 분포에 transformer를 사용하면서, codebook 훈련이 잘 되어야 한다

따라서 VQVAE의 DCL perceptual loss를 쓴다.

$$L_{GAN}(E, G, Z, D) = [\log D(x) + \log(1 - D(\hat{x}))]$$

$Q = \{E, G, Z\}$ 의 경우 optimal compression을 찾기 위해

$$Q^* = \arg \min_{E, G, Z} \max_D [L_{VAE}(E, G, Z) + \lambda L_{GAN}(E, G, Z, D)]$$

$$\lambda = \frac{\nabla_{G_L}[L_{rec}]}{\nabla_{G_L}[L_{GAN}] + \delta}$$

↓
decel last layer의 디贶 gradient.

모든 pixel context를 통해 거창 (in-reason) 대화에 있는 attention 진행.

- Learning the Composition of Images with Transformer.

* Latent Transformer.

E, G 가 훈련되었을 때 G 의 codebook으로 image를 만들 수 있다

정확히는, 양자화된 image x 는 $z_q = q(E(x))$ 로 주어지고,

$$\sum_{ij} s_{ij} = k \text{ such that } (z_q)_{ij} = z_k$$

해당되는 codebooked mapping은 decoder로 복원된다. $\hat{x} = G(z_q)$

index set S_i 를 Transformer는 가능한한 index로 예측한다.

$$P(S_i | S_{<i})$$
 를 찾는다.

$$P(S) = \prod_i P(S_i | S_{<i})$$

$$\mathcal{L}_{\text{Tr}} = E_{x \sim p_{\text{tar}}} [-\log P(S)]$$

- Conditional Synthesis

Condition을 주면 image를 생성한다.

$$P(s|c) = \prod_i P(s_i | s_{<i}, c)$$

Condition c 는 codebook의 벡터로 사용된다. 단위는 s 를 더해준다.

$$P(s_i | s_{<i}, r)$$

- Generating high-resolution Image.

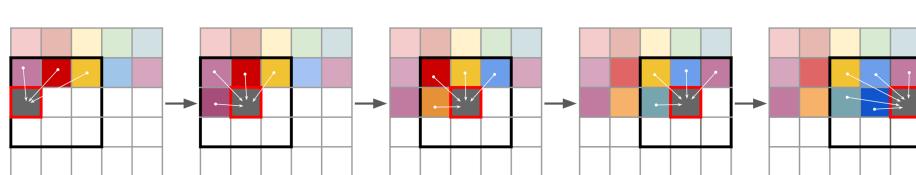


Figure 3. Sliding attention window.

Sequence h-w or limit을 준다.

mega-pixel을 gen하고면, patch로 잘라내고, s 를 최대로 한 후 생성된다.

이제 fig.3 같은 sliding window로 풀어보자.

VQGAN은 dataset의 spatial invariant 특성, spatial conditioning이 가능한데 model은 이를 faithful하게 한다.

실제로 업캐스팅은 아닙니다, unconditional 또는 coco-GAN처럼 condition을 넣을 수 있다.

* Experiments.

ConvNet Transformer 이점은 깊은 영역의 능력 평가.

codebook의 영향을 조사함.

$|\mathcal{X}| = 1024$ 이고, seq-len: 16x16으로 256

- Transformer in the Latent Space.

Transformer는 image의 autoregressive에 큰 발전, but low-res & shallow한 pixel은 고려하지 않았음.

→ (un-)conditional 어떤 Tfel conv 같은 차이를 조사함

Tfel pixelSNAIL의 차이를 맞춰 평가.

- Result

Negative Log-Likelihood (NLL)			
Data / # params	Transformer P-SNAIL steps	Transformer P-SNAIL time	PixelSNAIL fixed time
RIN / 85M	4.78	4.84	4.96
LSUN-CT / 310M	4.63	4.69	4.89
IN / 310M	4.78	4.83	4.96
D-RIN / 180 M	4.70	4.78	4.88
S-FLCKR / 310 M	4.49	4.57	4.64

Table 1. Comparing Transformer and PixelSNAIL architectures across different datasets and model sizes. For all settings, transformers outperform the state-of-the-art model from the PixelCNN family, PixelSNAIL in terms of NLL. This holds both when comparing NLL at fixed times (PixelSNAIL trains roughly 2 times faster) and when trained for a fixed number of steps. See Sec. 4.1 for the abbreviations.

→ Depth condition

→ semantic condition

⇒ Tfel은 낫다.

- Unified Model for Image Synthesis Tasks.

Transformer 모델로 이미지 합성 가능.

256x256 img, 16x16 latent size의 conditional synthesis의 경우

i) Semantic Image Synthesis

ii) structure-to-image.

iii) Pose-guided Synthesis

iv) Stochastic superresolution.

v) Class-conditional image Synthesis

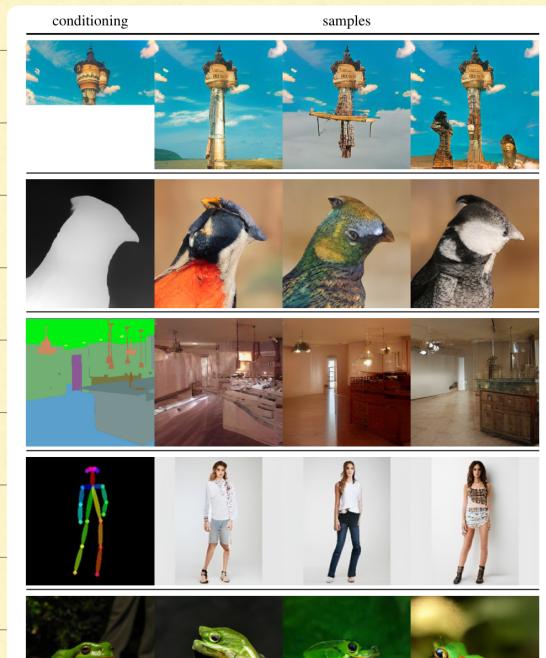


Figure 4. Transformers within our setting unify a wide range of image synthesis tasks. We show 256×256 synthesis results across different conditioning inputs and datasets, all obtained with the same approach to exploit inductive biases of effective CNN based *VQGAN* architectures in combination with the expressivity of transformer architectures. Top row: Completions from unconditional training on ImageNet. 2nd row: Depth-to-Image on RIN. 3rd row: Semantically guided synthesis on ADE20K. 4th row: Pose-guided person generation on DeepFashion. Bottom row: Class-conditional samples on RIN.

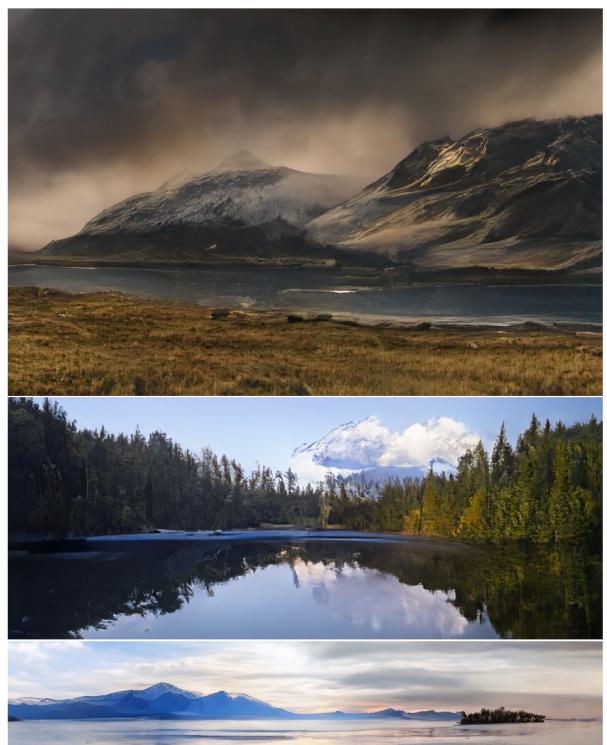


Figure 5. Samples generated from semantic layouts on S-FLCKR. Sizes from top-to-bottom: 1280×832 , 1024×416 and 1280×240 pixels. Best viewed zoomed in. A larger visualization can be found in the appendix, see Fig 29.

미시들은 같은 방법으로 학습된다.

Tf의 이점으로 task에 맞는 sequence를 찾을 수 있다.

다른 범위 mechanism으로 생각

- High-resolution Synthesis

Sliding Window 256x256 이미지를 기본으로

하지만 원래의 size의 ratio로 적용하는게 있다.

Dataset이 공급되는 대로 변화해야 한다.

256의 경우 $m=4$ (size $h = H/2^m$), block = $16 \times 16 \times 4$

S-FLCKR의 경우 $m=5$.

condition은 같이 Tf의 사용

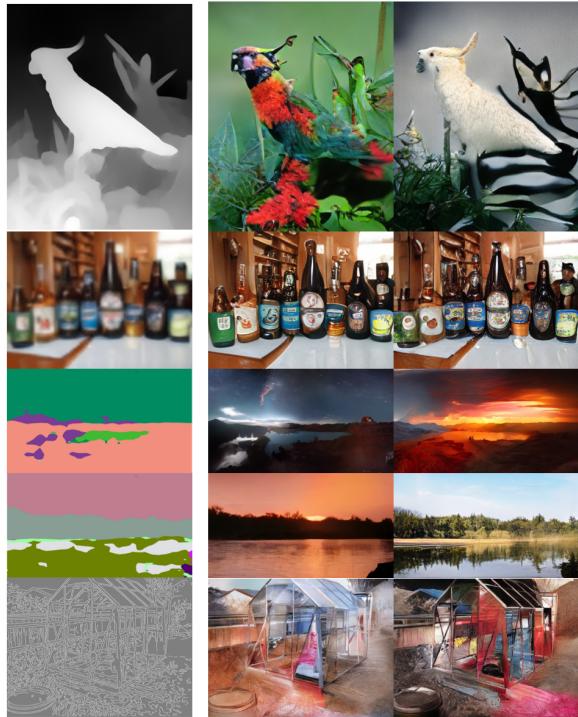


Figure 6. Applying the sliding attention window approach (Fig. 3) to various conditional image synthesis tasks. Top: Depth-to-image on RIN, 2nd row: Stochastic superresolution on IN, 3rd and 4th row: Semantic synthesis on S-FLCKR, bottom: Edge-guided synthesis on IN. The resulting images vary between 368×496 and 1024×576 , hence they are best viewed zoomed in.

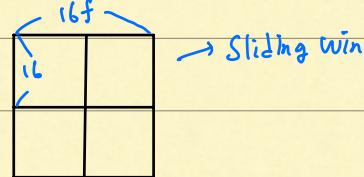
- Building Context-Rich Vocabularies

context-rich vocab의 중요성을 축소하기 위해 m 은 범위가 있는 Tf는 고정된 샘플을 진행한다.

image $H \times W$ 은 $H/f \times W/f$ 크기로 증가된다.
content size.

Tf의 경우 16×16 으로 훈련시키고, $16f$ 로 sample을 사용

↳ encoding size context의 양 줄임



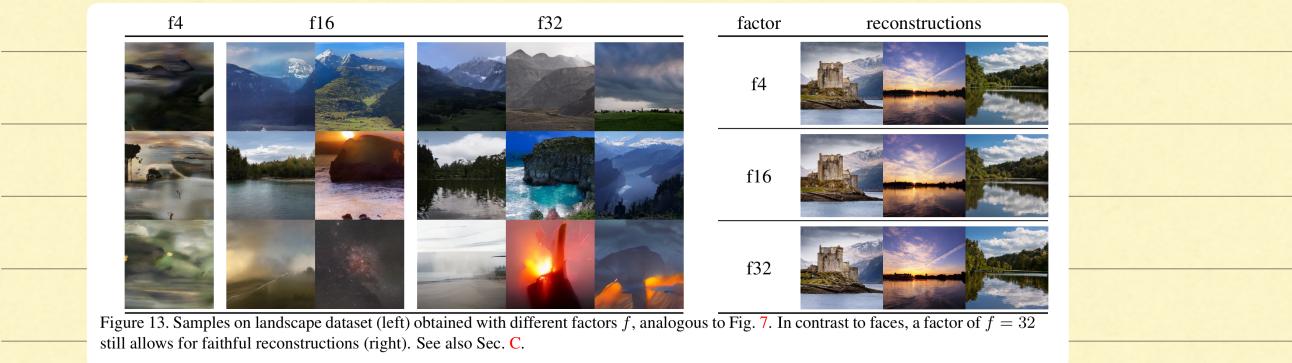
- Results

f1	f2	f8	f16	downsampling factor
1.0	3.86	65.81	280.68	
				speed-up

Figure 7. Evaluating the importance of effective codebook for HQ-Faces (CelebA-HQ and FFHQ) for a fixed sequence length $|s| = 16 \cdot 16 = 256$. Globally consistent structures can only be modeled with a context-rich vocabulary (right). All samples are generated with temperature $t = 1.0$ and top- k sampling with $k = 100$. Last row reports the speedup over the f1 baseline which operates directly on pixels and takes 7258 seconds to produce a sample on a NVIDIA GeForce GTX Titan X.

학습한 receptive field가 작으면 안된다.

$f=8$ 인 경우 물체가 생기고, $f=16$ 일 때 (full receptive) 동작이 잘 된다.



S-FL(learnable) 같은 결과.

정상적인 비교를 위해 image API 바로 학습하는 것과, codebook API 학습하는 것 사이의 차이를 봤다.

→ 18번 정확하고, 16번 사용

-Benchmarking Image Synthesis Results.

각 방법 비교.

Dataset	ours	SPADE [53]	Pix2PixHD (+aug) [75]	CRN [9]
COCO-Stuff	22.4	22.6/23.9(*)	111.5 (54.2)	70.4
ADE20K	35.5	33.9/35.7(*)	81.8 (41.5)	73.3

Table 2. FID score comparison for semantic image synthesis (256×256 pixels). (*): Recalculated with our evaluation protocol based on [50] on the validation splits of each dataset.

→ semantic synthesis

CelebA-HQ 256×256		FFHQ 256×256	
Method	FID ↓	Method	FID ↓
GLOW [37]	69.0	VDVAE ($t = 0.7$) [11]	38.8
NVAE [69]	40.3	VDVAE ($t = 1.0$)	33.5
PIONEER (B.) [23]	39.2 (25.3)	VDVAE ($t = 0.8$)	29.8
NCPVAE [1]	24.8	VDVAE ($t = 0.9$)	28.5
VAEBM [77]	20.4	VQGAN+P.SNAIL	21.9
Style ALAE [56]	19.2	BigGAN	12.4
DC-VAE [54]	15.8	ours (k=300)	9.6
ours (k=400)	10.2	U-Net GAN (+aug) [66]	10.9 (7.6)
PGGAN [31]	8.0	StyleGAN2 (+aug) [34]	3.8 (3.6)

⇒ unconditional face Synthesis

Table 3. FID score comparison for face image synthesis. CelebA-HQ results reproduced from [1, 54, 77, 24], FFHQ from [66, 32].

⇒ 투명한 method들도 네임드 있음.

⇒ TensorFlow top-k sampling 을 사용함

- Class-Conditional Synthesis on ImageNet.

VQGAN은 $f=16$, $Z=(6384 \times 128) \approx 800,000$

DFT VQGAN을 더 사실적으로 만든다.



Figure 8. Samples from our class-conditional ImageNet model trained on 256×256 images.

Model	acceptance rate	FID	IS
mixed k , $p = 1.0$	1.0	17.04	70.6 ± 1.8
$k = 973$, $p = 1.0$	1.0	29.20	47.3 ± 1.3
$k = 250$, $p = 1.0$	1.0	15.98	78.6 ± 1.1
$k = 973$, $p = 0.88$	1.0	15.78	74.3 ± 1.8
$k = 600$, $p = 1.0$	0.05	5.20	280.3 ± 5.5
mixed k , $p = 1.0$	0.5	10.26	125.5 ± 2.4
mixed k , $p = 1.0$	0.25	7.35	188.6 ± 3.3
mixed k , $p = 1.0$	0.05	5.88	304.8 ± 3.6
mixed k , $p = 1.0$	0.005	6.59	402.7 ± 2.9
DCTransformer [48]	1.0	36.5	n/a
VQVAE-2 [61]	1.0	~31	~45
VQVAE-2	n/a	~10	~330
BigGAN [4]	1.0	7.53	168.6 ± 2.5
BigGAN-deep	1.0	6.84	203.6 ± 2.6
IDDPM [49]	1.0	12.3	n/a
ADM-G, no guid. [15]	1.0	10.94	100.98
ADM-G, 1.0 guid.	1.0	4.59	186.7
ADM-G, 10.0 guid.	1.0	9.11	283.92
val. data	1.0	1.62	234.0 ± 3.9

Table 4. FID score comparison for class-conditional synthesis on 256×256 ImageNet, evaluated between 50k samples and the training split. Classifier-based rejection sampling as in VQVAE-2 uses a ResNet-101 [22] classifier. BigGAN(-deep) evaluated via <https://tfhub.dev/deepmind> truncated at 1.0. “Mixed” k refers to samples generated with different top-k values, here $k \in \{100, 200, 250, 300, 350, 400, 500, 600, 800, 973\}$.

⇒ autoregressive GAN 모델

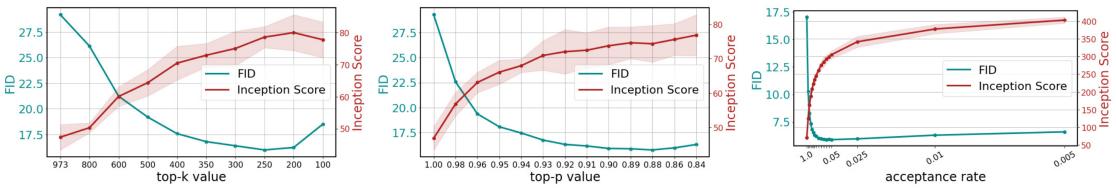


Figure 9. FID and Inception Score as a function of top-k, nucleus and rejection filtering.

- How good is VQGAN

GAN의 어떤 것(DDT 때문) 성능비교를 어떤 식으로

Model	Codebook Size	dim \mathcal{Z}	FID/val	FID/train
VQVAE-2	64×64 & 32×32	512	n/a	~ 10
DALL-E [59]	32×32	8192	32.01	33.88
VQGAN	16×16	1024	7.94	10.54
VQGAN	16×16	16384	4.98	7.41
VQGAN*	32×32	8192	1.49	3.24
VQGAN	64×64 & 32×32	512	1.45	2.78

Table 5. FID on ImageNet between reconstructed validation split and original validation (FID/val) and training (FID/train) splits.

*trained with Gumbel-Softmax reparameterization as in [59, 29].

= 역시 Codebook size LL 커가 훈련을 놓는다



Figure 12. Comparing reconstruction capabilities between VQVAEs and VQGANs. Numbers in parentheses denote compression factor and codebook size. With the same compression factor and codebook size, VQGANs produce more realistic reconstructions compared to blurry reconstructions of VQVAEs. This enables increased compression rates for VQGAN while retaining realistic reconstructions. See Sec. C.

* Conclusion.

T_f 를 사용해도 고비성도를 얻으되, quadratic cost를 해제

CNN과 T_f 의 장점을 합쳤고 general한 방식이라는 보기 얻어가 있다.

x Appendix

A. Changelog (이전 버전과의 차이점)

B. Implementation Details.

- Hyperparam

Experiment	n_{layer}	# params [M]	n_z	$ \mathcal{Z} $	dropout	length(s)	n_e	m
RIN	12	85	64	768	0.0	512	1024	4
c-RIN	18	128	64	768	0.0	257	768	4
D-RINV1	14	180	256	1024	0.0	512	768	4
D-RINV2	24	307	256	1024	0.0	512	1024	4
IN	24	307	256	1024	0.0	256	1024	4
c-IN	24	307	256	1024	0.0	257	1024	4
c-IN (big)	48	1400	256	16384	0.0	257	1536	4
IN-Edges	24	307	256	1024	0.0	512	1024	3
IN-SR	12	153	256	1024	0.0	512	1024	3
S-FLCKR, $f = 4$	24	307	256	1024	0.0	512	1024	2
S-FLCKR, $f = 16$	24	307	256	1024	0.0	512	1024	4
S-FLCKR, $f = 32$	24	307	256	1024	0.0	512	1024	5
(FacesHQ, $f = 1$) [*]	24	307	—	512	0.0	512	1024	—
FacesHQ, $f = 2$	24	307	256	1024	0.0	512	1024	1
FacesHQ, $f = 4$	24	307	256	1024	0.0	512	1024	2
FacesHQ, $f = 8$	24	307	256	1024	0.0	512	1024	3
FacesHQ ^{**} , $f = 16$	24	307	256	1024	0.0	512	1024	4
FFHQ ^{**} , $f = 16$	28	355	256	1024	0.0	256	1024	4
CelebA-HQ ^{**} , $f = 16$	28	355	256	1024	0.0	256	1024	4
FFHQ (big)	24	801	256	1024	0.0	256	1664	4
CelebA-HQ (big)	24	801	256	1024	0.0	256	1664	4
COCO-Stuff	32	651	256	8192	0.0	512	1280	4
ADE20K	28	405	256	4096	0.1	512	1024	4
DeepFashion	18	129	256	1024	0.0	340	768	4
LSUN-CT	24	307	256	1024	0.0	256	1024	4
CIFAR-10	24	307	256	1024	0.0	256	1024	1

Table 8. Hyperparameters. For every experiment, we set the number of attention heads in the transformer to $n_h = 16$. n_{layer} denotes the number of transformer blocks, # params the number of transformer parameters, n_z the dimensionality of codebook entries, $|\mathcal{Z}|$ the number of codebook entries, dropout the dropout rate for training the transformer, length(s) the total length of the sequence, n_e the embedding dimensionality and m the number of downsampling steps in the VQGAN. D-RINV1 is the experiment which compares to Pixel-SNAIL in Sec. 4.1. Note that the experiment (FacesHQ, $f = 1$)^{*} does not use a learned VQGAN but a fixed k-means clustering algorithm as in [8] with $K = 512$ centroids. A prefix “c” refers to a class-conditional model. The models marked with a ‘**’ are trained on the same VQGAN.

- Architecture (VQGAN)

Encoder	Decoder
$x \in \mathbb{R}^{H \times W \times C}$	$z_q \in \mathbb{R}^{h \times w \times n_z}$
Conv2D $\rightarrow \mathbb{R}^{H \times W \times C'}$	Conv2D $\rightarrow \mathbb{R}^{h \times w \times C''}$
$m \times \{ \text{Residual Block}, \text{Downsample Block} \} \rightarrow \mathbb{R}^{h \times w \times C''}$	Residual Block $\rightarrow \mathbb{R}^{h \times w \times C''}$
Residual Block $\rightarrow \mathbb{R}^{h \times w \times C''}$	Non-Local Block $\rightarrow \mathbb{R}^{h \times w \times C''}$
Non-Local Block $\rightarrow \mathbb{R}^{h \times w \times C''}$	Residual Block $\rightarrow \mathbb{R}^{h \times w \times C''}$
Residual Block $\rightarrow \mathbb{R}^{h \times w \times C''}$	$m \times \{ \text{Residual Block}, \text{Upsample Block} \} \rightarrow \mathbb{R}^{H \times W \times C'}$
GroupNorm, Swish, Conv2D $\rightarrow \mathbb{R}^{h \times w \times n_z}$	GroupNorm, Swish, Conv2D $\rightarrow \mathbb{R}^{H \times W \times C}$

Table 7. High-level architecture of the encoder and decoder of our VQGAN. The design of the networks follows the architecture presented in [25] with no skip-connections. For the discriminator, we use a patch-based model as in [28]. Note that $h = \frac{H}{2^m}$, $w = \frac{W}{2^m}$ and $f = 2^m$.

Transformers 구조는 각 블록으로 GPT-2의 capacity는 layer 수를 변경하여 조절

C. On Context-rich Vocabularies

factor m 은 encoding downsampling이 몇 번 이루어지는가? $f = \text{downsampling factor. } (f=2^m)$

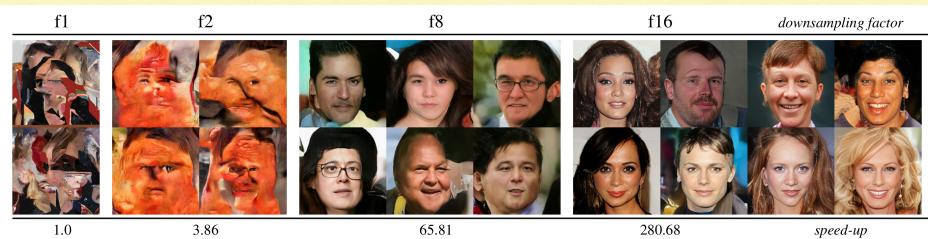


Figure 7. Evaluating the importance of effective codebook for HQ-Faces (CelebA-HQ and FFHQ) for a fixed sequence length $|s|=16$. Globally consistent structures can only be modeled with a context-rich vocabulary (right). All samples are generated with temperature $t=1.0$ and top- k sampling with $k=100$. Last row reports the speedup over the f1 baseline which operates directly on pixels and takes 7258 seconds to produce a sample on a NVIDIA GeForce GTX Titan X.

→ 성능에 놓은 (크면 좋음 → long-term interaction을 넓힐 가능)

→ 높지만 f가 크면 compression이 더 많이 되는데, 그래서 더욱 크면 recon의 성능↓

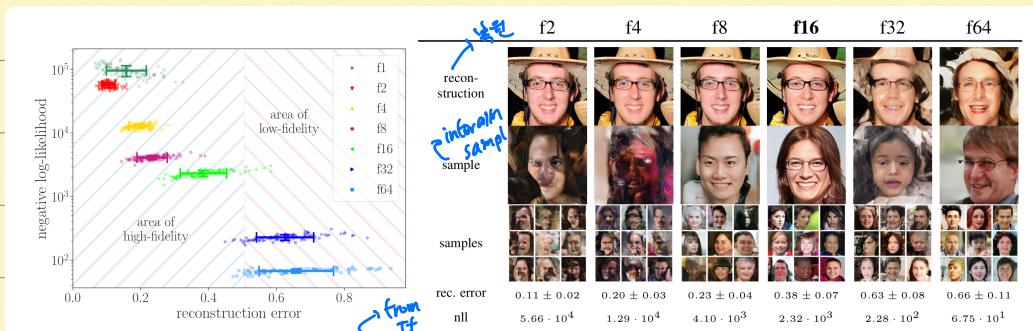


Figure 11. Trade-off between negative log-likelihood (nll) and reconstruction error. While context-rich encodings obtained with large factors f allow the transformer to effectively model long-range interactions, the reconstructions capabilities and hence quality of samples suffer after a critical value (here, $f=16$). For more details, see Sec. C.

NLL은 image의 representation distribution의 modeling 능력이 미친 나태변수.

f가 대로 크면 long-range의覃한 TFG가 잘되면서 NLL은 작지만 compression이 대로 많이 된다.

f가 상으로 recon을 잘되지만 sample quality↓

∴ 큰 f로 perceptually faithful한 recon을 encoding 선택된다.



Figure 12. Comparing reconstruction capabilities between VQVAEs and VQGANs. Numbers in parentheses denote compression factor and codebook size. With the same compression factor and codebook size, VQGANs produce more realistic reconstructions compared to blurry reconstructions of VQVAEs. This enables increased compression rates for VQGAN while retaining realistic reconstructions. See Sec. C.

↳ VQGAN, VQVAE, DALL-E 의 recon 능력 비교.

→ 그려진 부분적으로 떨어지면 그려진 그림, 성능↑

f의 크기는 dataset에 의존적이고.

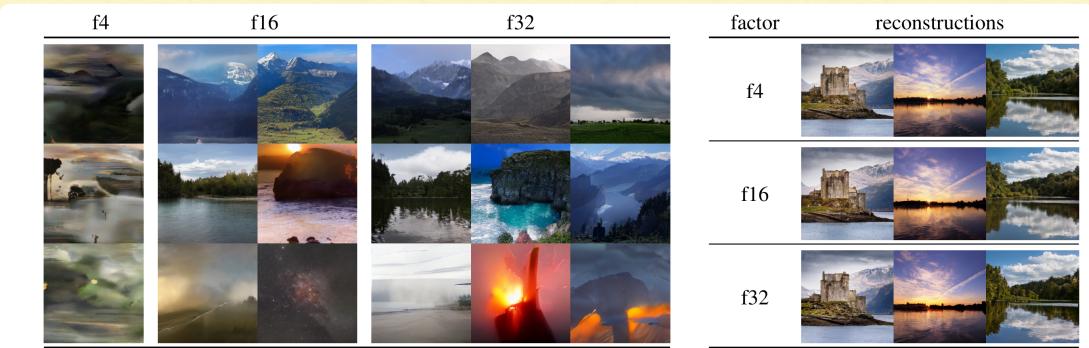


Figure 13. Samples on landscape dataset (left) obtained with different factors f , analogous to Fig. 7. In contrast to faces, a factor of $f = 32$ still allows for faithful reconstructions (right). See also Sec. C.

얼굴과 다르게 풍경에서는 perceptually 더 민감해서 f32가 성능이 좋다

* D. Additional Results.

- Qualitative Comparison.

Model	acceptance rate	FID	IS
mixed k , $p = 1.0$	1.0	17.04	70.6 ± 1.8
$k = 973$, $p = 1.0$	1.0	29.20	47.3 ± 1.3
$k = 250$, $p = 1.0$	1.0	15.98	78.6 ± 1.1
$k = 973$, $p = 0.88$	1.0	15.78	74.3 ± 1.8
$k = 600$, $p = 1.0$	0.05	5.20	280.3 ± 5.5
mixed k , $p = 1.0$	0.5	10.26	125.5 ± 2.4
mixed k , $p = 1.0$	0.25	7.35	188.6 ± 3.3
mixed k , $p = 1.0$	0.05	5.88	304.8 ± 3.6
mixed k , $p = 1.0$	0.005	6.59	402.7 ± 2.9
DCTransformer [48]	1.0	36.5	n/a
VQVAE-2 [61]	1.0	~31	~45
VQVAE-2	n/a	~10	~330
BigGAN [4]	1.0	7.53	168.6 ± 2.5
BigGAN-deep	1.0	6.84	203.6 ± 2.6
IDDPM [49]	1.0	12.3	n/a
ADM-G, no guid. [15]	1.0	10.94	100.98
ADM-G, 1.0 guid.	1.0	4.59	186.7
ADM-G, 10.0 guid.	1.0	9.11	283.92
val. data	1.0	1.62	234.0 ± 3.9

Table 4. FID score comparison for class-conditional synthesis on 256×256 ImageNet, evaluated between 50k samples and the training split. Classifier-based rejection sampling as in VQVAE-2 uses a ResNet-101 [22] classifier. BigGAN(-deep) evaluated via <https://tfhub.dev/deepmind> truncated at 1.0. “Mixed” k refers to samples generated with different top- k values, here $k \in \{100, 200, 250, 300, 350, 400, 500, 600, 800, 973\}$.

→ classifier 다른 성능

Dataset	ours-previous (+R)	BigGAN (-deep)	MSP	Dataset	ours-previous	ours-new
IN 256, 50K	19.8 (11.2)	7.1 (7.3)	n.a.	CelebA-HQ 256	10.7	10.2
IN 256, 18K	23.5	9.6 (9.7)	50.4	FFHQ 256	11.4	9.6

Table 6. Results from a previous version of this paper, see also Sec. A. Left: Previous results on class-conditional ImageNet synthesis with a slightly different implementation and evaluated against 50k and 18k training examples instead of the whole training split. See Tab. 4 for new, improved results evaluated against the whole training split. Right: Previous results on face-synthesis with a slightly different implementation compared to the new implementation. See also Tab. 3 for comparison with other methods.

⇒ 예시는 fig 14 ~ 17.

Semantic Synthesis

Dataset	ours	SPADE [53]	Pix2PixHD (+aug) [75]	CRN [9]
COCO-Stuff	22.4	22.6/23.9(*)	111.5 (54.2)	70.4
ADE20K	35.5	33.9/35.7(*)	81.8 (41.5)	73.3

Table 2. FID score comparison for semantic image synthesis (256×256 pixels). (*): Recalculated with our evaluation protocol based on [50] on the validation splits of each dataset.

⇒ 예시는 fig 40 ~ 41.

- Comparison to Image-GPT.

ImageGPT와 비교.

↳ TF의 image Synthesis 및 image representation 모듈이 훨씬 적이다.

↳ but (96x192 size) 구현

VQGAN의 high-res 가능

⇒ 예시: fig 27, 28

- Additional high-res results

Depth-to-Image, Edge-to-Image, Semantic Synthesis, Unconditional LSUN, Pose-to-image

예시 fig 29 ~ 44의 그림.

E. Nearest Neighbors of Sample.

GAN과 같이 likelihood-based gen model은 overfitting을 감지하는 징후가 없다.

이를 위해 큰 model (overfitted GDI 쪽)의 val NLL이 대비 NLL이 가장 큰 epoch

train NLL 대비 NLL이 가장 큰 epoch이 dataset에 유의.

각 epoch마다 sample을 뽑아, 그대로 사진 검색

train은 train set을 재구성, val은 새로운 data

∴ early-stopping을 고려해 봄지.

↳ 우리는 dataset 크기에 따른

FID는 고려해 봄지를 모니터링, val NLL은 품질을 보장 X

둘다 정기적인 matrix의 필요.

VQGAIN은 top-k로 다양성 확보.

F. On the Ordering of Image Representations.

NLP의 경우 이미지 순서를 정하지만, image는 명확 X

단지 Sliding의 경우 헤더 우선시 되기 때문에 5가지 순서의 생성 조사.

i) row major : 원쪽 위에서 오른쪽 아래.

ii) Spiral out : 중심부터 나선형.

iii) Z-curve

iv) subsample :

v) alternative : 헤더 평면은 방향을 고대로

vi) Spiral in : 끝부터 in으로

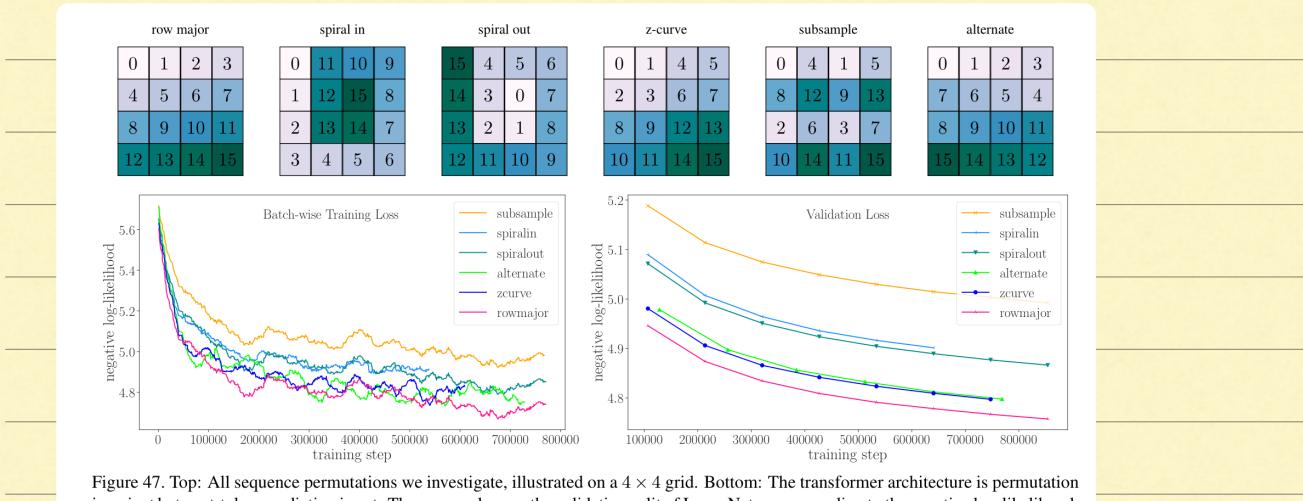


Figure 47. Top: All sequence permutations we investigate, illustrated on a 4×4 grid. Bottom: The transformer architecture is permutation invariant but next-token prediction is not: The average loss on the validation split of ImageNet, corresponding to the negative log-likelihood, differs significantly between different prediction orderings. Among our choices, the commonly used row-major order performs best.

Subsample와 spiral in이 경우 NLL이 가장 낮지만, 많은 질감을 생성했습니다.

다른 variant는 더 인식이 좋은 sample을 학습했습니다.

Subsample의 세부적인 bias와 도움이 되지는 않는 듯.

autoregressive codebook의 permutation-invariant이 아닙니다, major row가 가장 성능 좋다.

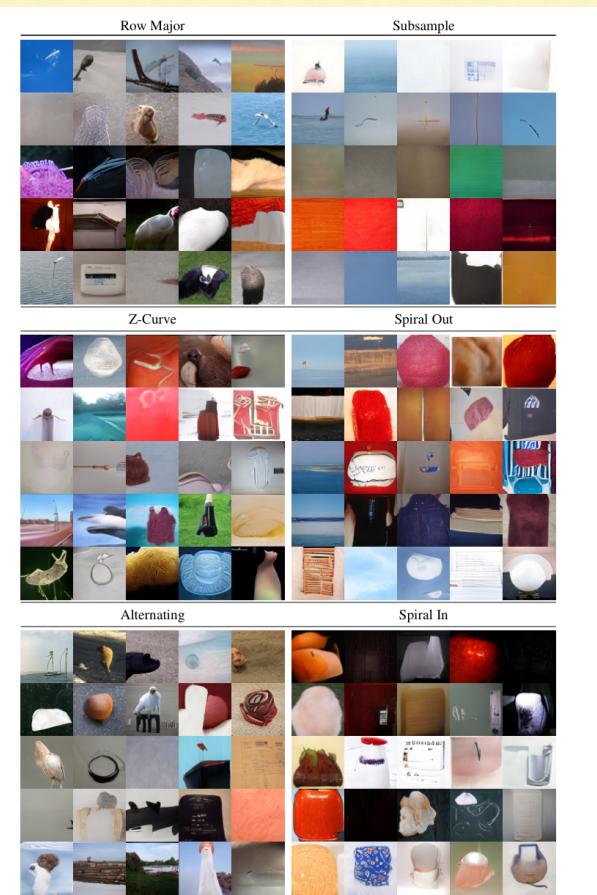


Figure 48. Random samples from transformer models trained with different orderings for autoregressive prediction as described in Sec. F.