

---

# Generating images with sparse representations

---

Charlie Nash<sup>1</sup> Jacob Menick<sup>1</sup> Sander Dieleman<sup>1</sup> Peter W. Battaglia<sup>1</sup>

## Abstract

The high dimensionality of images presents architecture and sampling-efficiency challenges for likelihood-based generative models. Previous approaches such as VQ-VAE use deep autoencoders to obtain compact representations, which are more practical as inputs for likelihood-based models. We present an alternative approach, inspired by common image compression methods like JPEG, and convert images to quantized discrete cosine transform (DCT) blocks, which are represented sparsely as a sequence of DCT channel, spatial location, and DCT coefficient triples. We propose a Transformer-based autoregressive architecture, which is trained to sequentially predict the conditional distribution of the next element in such sequences, and which scales effectively to high resolution images. On a range of image datasets, we demonstrate that our approach can generate high quality, diverse images, with sample metric scores competitive with state of the art methods. We additionally show that simple modifications to our method yield effective image colorization and super-resolution models.

## 1. Introduction

Deep generative models of images are neural networks trained to output synthetic imagery. Current models generate sample images that are difficult for humans to distinguish from real images, and have found applications in image super-resolution (NVIDIA), colorization (Antic), and text-guided generation (Ramesh et al.). Such models fall broadly into three categories: generative adversarial networks (GANs, Goodfellow et al. 2014), likelihood-based models, and energy-based models. GANs use discriminator networks that are trained to distinguish samples from generator networks and real examples. Likelihood-based models, including variational autoencoders (VAEs, Kingma & Welling 2014; Rezende et al. 2014), normalizing flows

<sup>1</sup>DeepMind, London, United Kingdom. Correspondence to: Charlie Nash <charlienash@google.com>.



Figure 1. Images often have sparse structure, with salient content distributed unevenly across locations. Our model DCTransformer predicts *where* to place content, as well as *what* content to add. (left) Images generated by DCTransformer, and (right) associated heatmaps of image locations selected in the generation process.

(Rezende & Mohamed, 2015), and autoregressive models (Van Oord et al., 2016), directly optimize the model log-likelihood or the evidence lower bound. Energy-based models estimate a scalar energy for each example that corresponds to an unnormalized log-probability, and can be trained with a variety of objectives (Du & Mordatch, 2019).

Likelihood-based models offer several important advantages. The training objective incentivizes learning the full data distribution, rather than only a subset of the modes (termed “mode-dropping”), which is a common downside of GANs. Likelihood-based training tends to be more stable than adversarial alternatives, and the objective can be used to detect overfitting using a held-out set. The dramatic recent advances in modelling natural language made by GPT-3 (Brown et al., 2020) demonstrate that optimizing a

simple log-likelihood objective, using expressive models and large datasets, is an effective approach for modelling complex data.

Optimizing the likelihood of pixel-based images can be problematic due to their complexity and high dimensionality, however. To our knowledge, no likelihood-based model operating on raw pixels has demonstrated competitive sample quality on the ImageNet dataset (Russakovsky et al., 2015) at a resolution of 256x256 or higher. For autoregressive models, conditioning on, and sequentially sampling, the hundreds of thousands of pixels in a typical image can be prohibitive. Some likelihood-based approaches address this by reducing the dimensionality using, e.g., low-precision or quantized color spaces and images (Kingma & Dhariwal, 2018; Chen et al., 2020; Van Oord et al., 2016). VQ-VAE (Oord et al., 2017; Razavi et al., 2019) uses a vector-quantized autoencoder network to first perform neural lossy compression before downstream generative modelling, which reduces representation size while maintaining quality.

Beyond generative models, the dimensionality of images poses a challenge whenever data storage, transmission, and processing budgets are at a premium. Fortunately, natural images have tremendous redundancy<sup>1</sup>, which all modern image compression methods exploit: even the images in this paper’s PDF file are stored in a compressed format.

JPEG (Wallace, 1992) and other popular lossy image compression methods use the discrete cosine transform (DCT) (Ahmed et al., 1974; Rao & Yip, 1990) to separate spatial frequencies of an image, and encode them with controllable resource budgets (e.g., the “quality” parameter). This takes advantage of the statistical structure of natural images (e.g., smooth signals with high frequency noise) and human vision (e.g., low frequencies tend to be more perceptually salient) by dropping high frequency information, to strike favorable efficiency/quality trade-offs (Wang et al., 2004).

To capitalize on the decades of engineering behind modern compression tools, we propose a generative model over DCT representations of images, rather than pixels. We convert an image to a 3D tensor of quantized DCT coefficients, and represent them sparsely as a sequence of 3-tuples that encode the DCT channel, spatial location, and DCT coefficient for each non-zero element in the tensor. We present a novel Transformer-based autoregressive architecture (Vaswani et al., 2017), called “DCTransformer”, which

<sup>1</sup>Kersten (1987) estimated 4-bit grayscale pixel images have (an upper bound of) 1.42 median bits of information per pixel, which means they can generally be compressed to  $\frac{1.42}{4} = 35.5\%$  their original size. They used human vision as a model, analogous to how Shannon (1951) measured the redundancy of English by having people guess the next character in a text sequence.

is trained to sequentially predict the conditional distribution over the next element in the sequence, resulting in a model that predicts both where to add content in an image, and what content to add (Figure 1). We find that the DCTransformer’s sample quality is competitive with state-of-the-art GANs, but with better sample diversity. We find sparse DCT-based representations help mitigate the inference time, memory, and compute costs of traditional pixel-based autoregressive models, as well as the time-consuming training of neural lossy compression embedding functions used in models such as VQ-VAE (Oord et al., 2017).

## 2. DCT-based sparse image representations

### 2.1. Block DCT

The DCT projects an image into a collection of cosine components at differing 2D frequencies. The two-dimensional DCT is typically applied to zero-centered  $B \times B$  pixel blocks  $\mathbf{P}$  to obtain a  $B \times B$  DCT block  $\mathbf{D}$ :

$$D_{uv} = \frac{1}{4} \alpha(u) \alpha(v) \times \sum_{x=0}^{B-1} \sum_{y=0}^{B-1} P_{xy} \cos \left[ \frac{(2x+1)u\pi}{2B} \right] \cos \left[ \frac{(2y+1)v\pi}{2B} \right] \quad (1)$$

$$\alpha(u) = \begin{cases} \frac{1}{\sqrt{2}}, & \text{if } u = 0 \\ 1, & \text{otherwise} \end{cases},$$

where  $x$  and  $y$  represent horizontal and vertical pixel coordinates,  $u$  and  $v$  index the horizontal and vertical spatial frequencies, and  $\alpha$  is a normalizing scale factor to enforce orthonormality. We follow the JPEG codec and split images into the YCbCr color space, containing a brightness component Y (luma) and two color components Cb and Cr (chroma). We perform 2x downsampling in both the horizontal and vertical dimensions for the two chroma channels, and then apply the DCT independently to blocks in each channel. Chroma downsampling in this way substantially reduces color information, at a minimal perceptual cost.

### 2.2. Quantization

In the JPEG codec, DCT blocks are quantized by dividing elementwise by a quantization matrix, and then rounding to the nearest integer. The quantization matrix is typically structured so that higher frequency components are squashed to a larger extent, as high frequency variation is typically harder to detect than low frequency variation. We follow the same approach, and use quality-parameterized quantization matrices for  $8 \times 8$  blocks as defined by the Independent JPEG Group in all our experiments. For block sizes other than 8, we interpolate the size 8 quality matrix. Appendix A has more details.

## Generating images with sparse representations

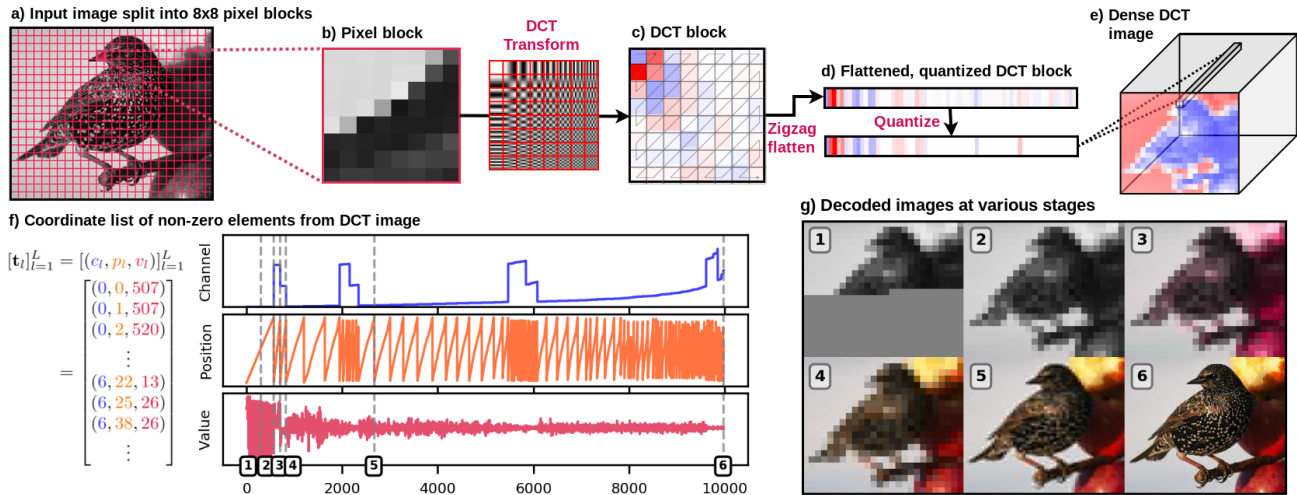


Figure 2. a) The input image is split into 8x8 blocks of pixels. b-c) The symmetric 2D DCT is used to transform each block into frequency coefficients. d) The block is flattened using the indicated zigzag ordering, and then quantized using a quality-parameterized quantization vector. e) The collection of quantized DCT vectors can be reformed into a DCT image, with  $8 \times 8 = 64$  channels. f) The DCT image is converted to a coordinate list of non-zero channel-position-value tuples, as a sparsified representation of the image. g) Images decoded at intermediate steps indicated on the coordinate list in f).

### 2.3. Sparsity

Sparsity is induced in DCT blocks through the quantization process, which squashes high frequency components and rounds low values to zero. The JPEG codec takes advantage of this sparsity by flattening the DCT blocks using a zigzag ordering from low to high frequency components, and applying a run-length encoding to compactly represent the strings of zeros that typically appear at the end of the DCT vectors.

We use an alternative sparse representation, where zigzag-flattened DCT vectors are reassembled in their corresponding spatial positions into a DCT image, of size  $H/B \times W/B \times B^2$ , where  $H$  and  $W$  are the height and width of the image, respectively, and  $B$  is the block size (Figure 2e). The DCT image’s non-zero elements are then serialized to a sparse list consisting of (DCT channel, spatial-position, value) tuples. This is repeated for the Y, Cb, and Cr image channels, with the Y channel occupying the first  $B^2$  DCT bands, and the Cb and Cr channels occupying the second and third  $B^2$  DCT bands, respectively. The spatial positions of the  $2 \times$  downsampled chroma channels are multiplied by 2 to take them into correspondence with the Y channel positions. The resulting lists are concatenated, and a stopping token is added in order to indicate the end of the variable length sequences.

For unconditional image generation we sort from low to high frequencies, with color channels interleaved at intervals with the luminance channel. This results in a natural upsampling structure, where low frequency content is represented first, and high frequency content is progressively added, as shown

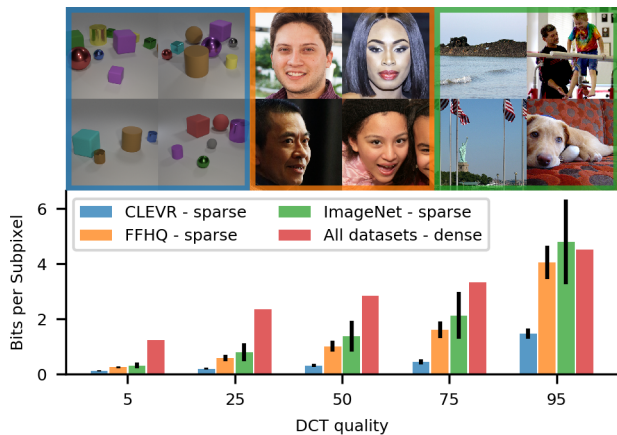


Figure 3. Bits per image-subpixel for dense and sparse block-DCT representations, quantized at various quality settings, reported for 1000 images from CLEVR (Johnson et al., 2017), FFHQ (Karras et al., 2019) and ImageNet (Russakovsky et al., 2015). Dense representations use the same amount of bits in each image location, whereas sparse representations use more bits on regions of greater detail, resulting in variable length codes. For spatially sparse image datasets like CLEVR, sparse DCT representations are substantially more compact, even at higher DCT quality settings.

in Figure 2g. Note, this ordering is not the only option: An alternative that we discuss in Section 4.3 places the luma data before the chroma data, resulting in an image colorization scheme.

Figure 3 compares the size of dense and sparse representations as a function of the DCT quality setting, and shows that for all but the highest quality setting, sparse representations are substantially smaller.

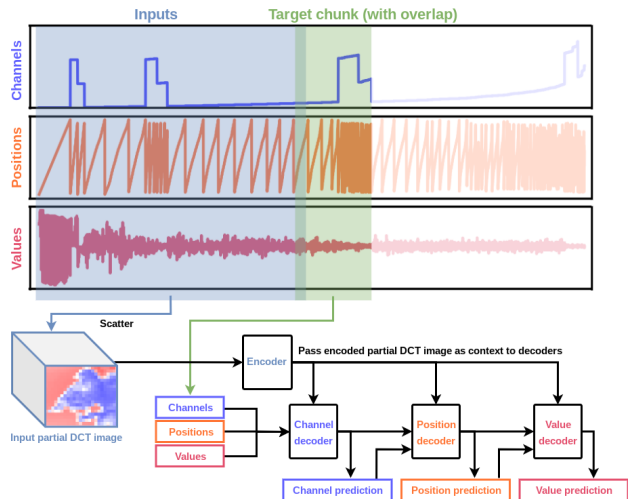


Figure 4. Chunk-based training and stacked Transformer architecture. A target chunk is selected during training, and previous elements are collected in an input slice. The input slice is gathered into a 3D DCT image Tensor. The DCT-image is flattened and embedded using a Transformer encoder, and the target slice is passed through a series of right-masked Transformer decoders, each of which performs cross-attention into the encoded DCT image, in order to predict channels, positions and values.

### 3. DCTransformer

We model the sparse DCT sequence autoregressively, by predicting the distribution of the next element conditioned on the previous elements. Samples are generated by iteratively sampling from the predictive distributions, and decoding the resulting sparse DCT sequence to images. For a sequence consisting of channel, position, value tuples  $[t_i]_{i=1}^L = [(c_i, p_i, v_i)]_{i=1}^L$  the joint distribution factors as:

$$\prod_{l=1}^L p(c_l | t_{<l}; \theta) p(p_l | c_l, t_{<l}; \theta) p(v_l | c_l, p_l, t_{<l}; \theta), \quad (2)$$

where  $\theta$  represents model parameters. While pixel-based autoregressive models predict values at every spatial location, DCTransformer predicts the channel of the next value, then the spatial position, and finally the quantized DCT value itself. We use categorical predictive distributions for each of the variable types, picking upper and lower bounds for the quantized DCT values.

#### 3.1. Chunked training for long sequences

Although sparse DCT representations are more compact than raw pixels, the length of DCT sequences depends on the resolution, dataset and quality setting of the quantization matrix used. For the highest resolution datasets we consider, DCT sequences can contain upwards of 100k tuples. This

presents a challenge for standard Transformer-based architectures, as the memory requirements of self attention layers scale quadratically with the sequence length. In practice, memory-issues are common in Transformers when operating on sequences with more than a few thousand elements. While memory-efficient Transformer variants (Child et al., 2019; Kitaev et al., 2020; Choromanski et al., 2021; Dhariwal et al., 2020; Zaheer et al., 2021; Roy et al., 2020) make long-sequence training more practical, care is still required to optimize model hyperparameters against available memory when training across datasets with varying sequence lengths.

Our approach is to train on fixed-size target chunks, while conditioning on prior context in a memory and compute-efficient way (Figure 4). During training a fixed-size target chunk is randomly selected from the DCT sequence. The sequence elements prior to the target chunk are treated as inputs, and using the inverse of the sparsification process shown in Figure 2e-f, are converted to a dense DCT image. The DCT image is optionally downsampled, then spatially flattened and embedded using a Transformer encoder. The target slice is processed with a Transformer decoder, that alternates between right-masked self-attention layers, cross attention layers that attend to the embedded DCT image, and dense layers. We use a target chunk size of 896 in all our experiments, and add an overlap of size 128 into the input chunk, so that the first predictions in the target chunk can attend directly into their preceding elements.

The resulting architecture, shown in Figure 4 uses constant memory and compute with respect to the size of the input slice, enabling training on large sequences, as well as simplifying the practical application of the model to datasets of varying resolutions. For more information about our strategy for selecting target chunks see Appendix C.

#### 3.2. Stacked channel, position, and value decoders

In order to model the three distinct variable types in DCT coordinate lists, our architecture must yield autoregressive predictions for each variable, and should condition on all the available context. One possible approach is to flatten the tuples, and to pass the resulting values through a single Transformer decoder, projecting the resulting states to logits as appropriate for each variable. To increase the amount of content modelled in a particular target sequence, we instead opt to use three distinct Transformer decoders, one for each of the channel, position and value predictions. The three decoders are stacked on top of one another, yielding sequence lengths three times smaller than flattened sequences.

Let  $\mathbf{D}$  be the input DCT image associated with a given input slice, then DCT image embeddings are obtained by passing

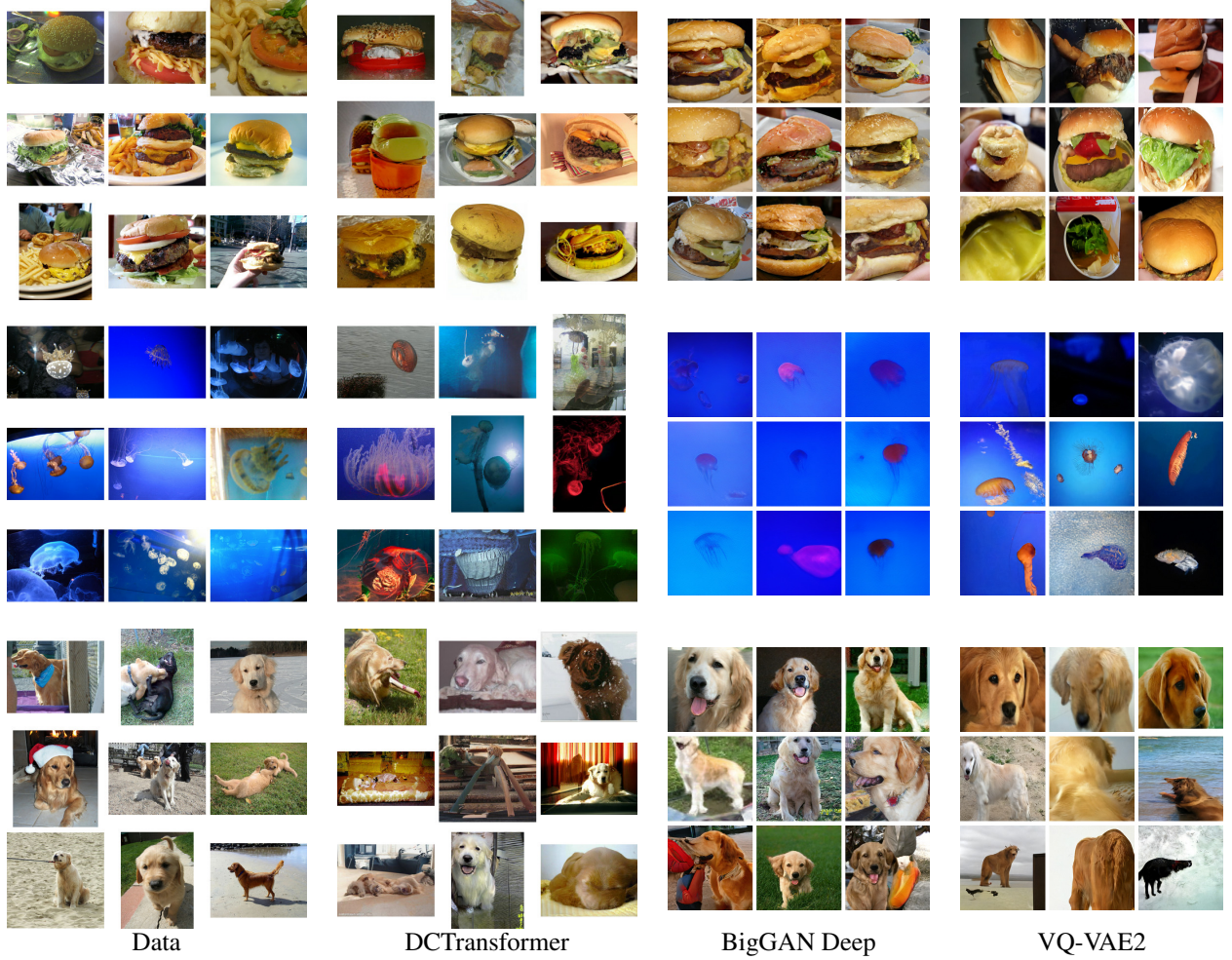


Figure 5. Comparison between ImageNet images and model samples for cheeseburger (label-933), jellyfish (label-107) and golden retriever (label-207) classes. DCTransformer produces variable aspect ratio samples with long-side resolution 384. BigGAN and VQ-VAE are trained and sample at a fixed 256x256 resolution corresponding to a resized long-side crop of the input images. BigGAN samples use truncation 1.0 to yield maximum diversity.

the flattened DCT image through a Transformer encoder:

$$\mathbf{E}_{\text{input}} = \text{encode}(\mathbf{D}_{\text{flat}}). \quad (3)$$

Predictive distributions for the channels are obtained by passing the encoded DCT image along with summed-channel  $\mathbf{C}$ , image position  $\mathbf{P}$ , value  $\mathbf{V}$ , and chunk-position  $\mathbf{P}_{\text{chunk}}$  embeddings into the channel decoder:

$$\mathbf{E}_{\text{channel}} = \mathbf{C}_{1:S-1} + \mathbf{P}_{1:S-1} + \mathbf{V}_{1:S-1} + \mathbf{P}_{\text{chunk}} \quad (4)$$

$$\mathbf{H}_{\text{channel}} = \text{decode}_{\text{channel}}(\mathbf{E}_{\text{channel}}; \mathbf{E}_{\text{input}}). \quad (5)$$

We distinguish here between inputs that are processed with masked self attention (left), and those that are processed with cross-attention (right). The final hidden state  $\mathbf{H}_{\text{channel}}$  is passed along with updated channel embeddings to the position decoder:

$$\mathbf{E}_{\text{position}} = \mathbf{H}_{\text{channel}} + \mathbf{C}_{2:S} \quad (6)$$

$$\mathbf{H}_{\text{position}} = \text{decode}_{\text{position}}(\mathbf{E}_{\text{position}}; \mathbf{E}_{\text{input}}) \quad (7)$$

When predicting a particular DCT value, we know its associated channel, and spatial position (Equation 2). We use this information to gather values from the embedded DCT inputs that are in the same spatial position as the target DCT value. This allows the value decoder to access all of the prior DCT values in the spatial position where it is making a prediction. The gathered DCT embeddings are added to the prior hidden state  $\mathbf{H}_{\text{position}}$  and passed to the value decoder:

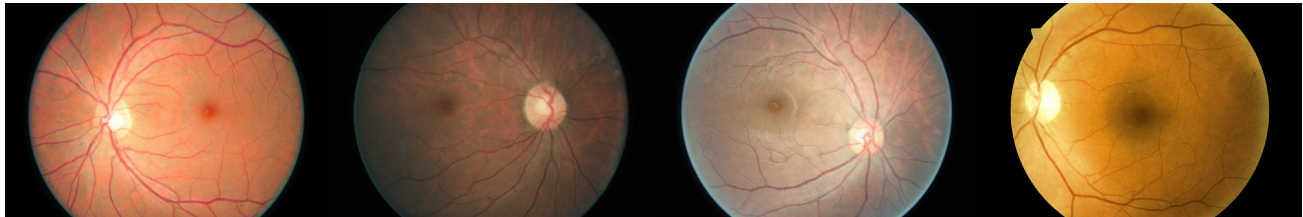
$$\mathbf{E}_{\text{value}} = \mathbf{H}_{\text{position}} + \text{gather}(\mathbf{E}_{\text{inputs}}, \mathbf{P}_{2:S}) \quad (8)$$

$$\mathbf{H}_{\text{value}} = \text{decode}_{\text{value}}(\mathbf{E}_{\text{value}}; \mathbf{E}_{\text{input}}) \quad (9)$$

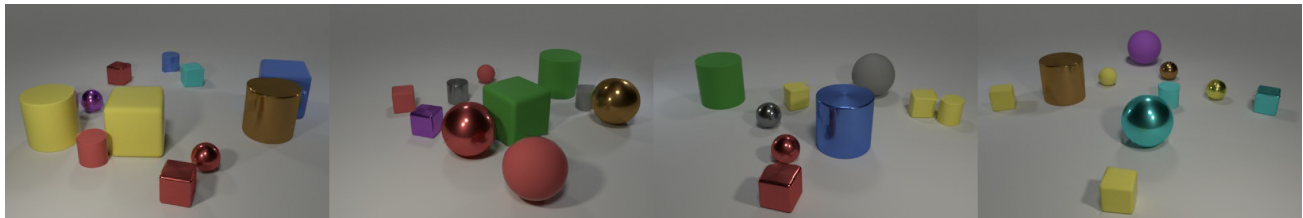
The resulting hidden states  $\mathbf{H}_{\text{channel}}$ ,  $\mathbf{H}_{\text{position}}$  and  $\mathbf{H}_{\text{value}}$  are each passed through linear layers to obtain the logits of their associated predictive distributions.



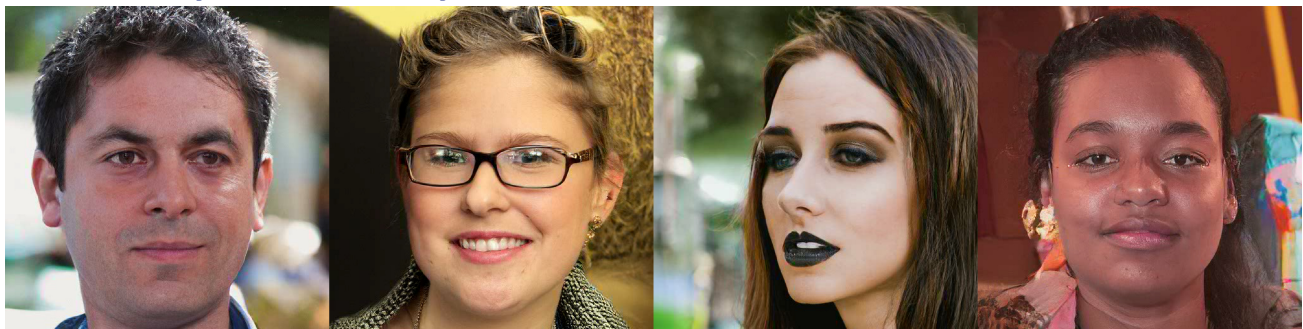
(a) **Plant leaves** (Chouhan et al., 2019) Long-side resolution 2048, block-size 32 and DCT quality 50. Available at [https://www.tensorflow.org/datasets/catalog/plant\\_leaves](https://www.tensorflow.org/datasets/catalog/plant_leaves)



(b) **Diabetic retinopathy** (Kaggle & EyePacs, 2015) Long-side resolution 1024, block-size 16 and DCT quality 75. Available at [https://www.tensorflow.org/datasets/catalog/diabetic\\_retinopathy\\_detection](https://www.tensorflow.org/datasets/catalog/diabetic_retinopathy_detection)



(c) **CLEVR** (Johnson et al., 2017) Long-side resolution 480, block-size 8 and DCT quality 90. Available at <https://www.tensorflow.org/datasets/catalog/clevr>



(d) **FFHQ** (Karras et al., 2019) Long-side resolution 1024, block-size 16 and DCT quality 35 Available at <https://github.com/NVlabs/ffhq-dataset>

Figure 6. Selected image samples on a range of varied image datasets. Each set of 4 images is selected from a sampled batch of 50 to showcase the image quality achievable by DCTransformer.

### 3.3. Sampling

As with training, sampling operates in a chunked fashion. Given a sequence of observed values from a DCT coordinate list, the model constructs an input dense DCT block, and proceeds to autoregressively sample a fixed size continuation. Once the chunk is complete its values are added to the input DCT image, and the process is repeated until a maximum number of chunks have been produced, or the stopping token is sampled. As with training, the sampling process requires constant computation and memory for each

block. This contrasts with standard Transformers where these factors expand with the sequence length.

## 4. Experiments

Quantitative comparison of generative image models is challenging, as not all model classes define normalized likelihoods (e.g., GANs), and even for models that do, the underlying data representation may not be the same, and so likelihoods are not directly comparable. In VQ-VAE and re-

Model	Precision	Recall	FID	sFID
<b>LSUN Bedrooms</b>				
StyleGAN	<b>0.55</b>	0.48	<b>2.45</b>	6.68
ProGAN	0.43	0.40	8.35	9.46
DCTransformer	0.44	<b>0.56</b>	6.40	<b>6.66</b>
<b>LSUN Towers</b>				
ProGAN	0.51	0.33	10.24	10.02
DCTransformer	<b>0.54</b>	<b>0.54</b>	<b>8.78</b>	<b>7.70</b>
<b>LSUN Churches</b>				
StyleGAN2	<b>0.60</b>	0.43	<b>4.01</b>	<b>9.28</b>
ProGAN	0.61	0.38	6.42	10.47
DCTransformer	0.60	<b>0.48</b>	7.56	10.71
<b>FFHQ</b>				
StyleGAN	<b>0.71</b>	<b>0.41</b>	<b>4.39</b>	11.57
DCTransformer	0.51	0.40	13.06	<b>9.44</b>
<b>ImageNet (class conditional)</b>				
BigGAN-deep	<b>0.78</b>	0.35	<b>6.59</b>	<b>6.66</b>
VQ-VAE2	0.36	0.57	31.11	17.38
DCTransformer	0.36	<b>0.67</b>	36.51	8.24

Table 1. Precision, recall and FID metrics comparison. sFID is equivalent to FID but uses intermediate spatial features in the inception network rather than the spatially-pooled features used in standard FID.

lated models for example, likelihood scores are meaningful only in relation to the latent codes associated with a particular neural encoder. We therefore use the sample-based Frechet inception distance (FID) (Heusel et al., 2017) and precision and recall (Kynkäänniemi et al., 2019) metrics, which enable comparison between any model that produces samples. We report FID using both the standard `pool_3` inception features, and the first 7 channels from the intermediate `mixed_6/conv` feature maps, which we refer to sFID (Szegedy et al., 2016). `pool_3` features compress spatial information to a large extent, making them less sensitive to spatial variability. We include the intermediate `mixed_6/conv` features to provide a picture of the spatial distributional similarity between models. We use the first 7 channels in order to obtain a feature space of size  $17 \times 17 \times 17 = 2023$ , which is comparable size to the size 2048 `pool_3` feature vector. We report Precision and Recall scores based on `pool_3` features. For all FID scores we follow Brock et al. (2019) and compare 50k model samples to features computed on the entire training set.

Model	Precision	Recall	FID	sFID
<b>ImageNet (class conditional 8x upsampling)</b>				
Val. set (low res. )	0.02	0.09	163.13	314.27
Val. set	<b>0.69</b>	0.59	<b>5.72</b>	14.06
DCTransformer	0.61	<b>0.61</b>	9.98	<b>11.06</b>
<b>OpenImagesV4 colorization</b>				
Val. set (greyscale)	0.61	0.28	34.19	38.09
Val. set	<b>0.75</b>	0.40	26.95	27.29
DCTransformer	0.72	<b>0.41</b>	<b>22.52</b>	<b>22.83</b>

Table 2. Precision, recall, and FID metrics for upsampling and colorization models.

#### 4.1. Image generation benchmarks

Table 1 shows FID, precision, recall scores on LSUN datasets (Yu et al., 2015), Flickr faces HQ (FFHQ, Karras et al. (2019)) as well as class-conditional ImageNet (Rusakovsky et al., 2015). On the LSUN datasets we compare to ProGAN (Karras et al., 2018) and StyleGAN (Karras et al., 2019; 2020) baselines, using publically available repositories of samples from pre-trained models. For LSUN towers samples from a pre-trained StyleGAN model are not publically available. On Imagenet we compare to BigGAN (Brock et al., 2019), and VQ-VAE2 (Razavi et al., 2019), using sample repositories sent to us by the VQ-VAE2 authors.

We find that overall the picture is mixed, with GANs achieving the strongest precision and FID scores, but DCTransformer typically achieving the best recall scores, reflecting the likelihood training metric that emphasises data coverage. The FID gap is closed to a large extent when using the spatial sFID, with DCTransformer achieving the best scores on three of the five datasets studied. This is consistent with our observation that DCTransformer samples tend to be diverse, and to exhibit realistic textures and structures, but are somewhat less reliably coherent than GAN samples. Figure 5 shows uncurated samples from three ImageNet classes, and Figures 10 and 11 in appendix E show uncurated samples on FFHQ and LSUN datasets respectively. For model hyperparameters and training details see appendix C.

#### 4.2. Image generation on diverse datasets

Due to the resolution independence of our chosen image representation and architecture, we can apply the same model to images that vary in resolution. Unlike GANs, our models are very stable in training, and enable us to detect overfitting by evaluating log-likelihood scores. As such DCTransformer is straightforward to apply to new image datasets. To illustrate the flexibility of our model, we train on a diverse set of image datasets, ranging from synthetic 3D objects



Figure 7. Image colorization. (left) grayscale image, (middle) three samples generated by DCTransformer, (right) original image.



Figure 8. 8x image upsampling. (left) downsampled image, (middle) three samples generated by DCTransformer, (right) original image.

to medical images of the retina. Figure 6 shows curated samples selected from a batch of 50 examples. We note the high quality results obtained on the plant leaves dataset at long-side resolution 2048. To our knowledge these are the first high-quality unconditional generations at this resolution in the literature using neural networks.

### 4.3. Image colorization and upsampling

The ordering we use for image generation naturally results in a model that upsamples from low frequency DCT components to high frequency (see Section 2.3). We investigate the effectiveness of DCTransformer as a class-conditional upsampling model by providing the first DCT channel for the luminance and chrominance components as a context, and sampling sequence continuations using a model trained on ImageNet for class-conditional generation. The first DCT component represents the average intensity in a pixel block, and so this procedure is equivalent to upsampling by a factor of the block size. Figure 8 shows some example results, and Figure 12 in appendix E shows additional uncurated results. Table 2 shows sample metrics for the upsampling model, where we compare DCTransformer-upsampled validation set images to the original validation set. The table shows that with the exception of the precision metric, DCTransformer-upsampled images obtain scores similar or better than to the validation set images.

In our standard DCT sequence ordering, color channels and luminance channels are interleaved, which means color

information is distributed throughout the sequence. If we instead place the color channels at the end of the sequence and condition on all the luminance information, we can train DCTransformer as a colorization model. We train a model on OpenImages V4 (Kuznetsova et al., 2018), and unlike the upsampling model, we train only on target slices from the color channels at the end of the sequence. Figure 7 shows some example results, and Figure 13 in Appendix E shows additional uncurated results. Table 2 shows sample metric scores for DCTransformer, compared to ground truth, as well as the greyscale components of the ground truth images. The sample metrics are very close to the validation set scores, suggesting a close distributional match to the training distribution.

## 5. Related Work

Our method builds on prior autoregressive generative models of images. Perhaps the most similar work is VQ-VAE2 (Razavi et al., 2019), which is also an autoregressive model of an image representation that has undergone lossy compression. In the case of VQ-VAE2, the image representation has fewer dimensions than the raw pixels but always a fixed size, whereas our sparse representation has a dynamic size depending on the content. VQ-VAE2 uses a hierarchy of latent codes, each modelled with a PixelSnail-style network with self-attention components. For the low-to-high frequency ordering we use in most experiments, a similar coarse-to-fine representation is obtained.



Other work has also found it advantageous to model a representation of the image which has lost information deemed less perceptually less relevant, e.g. representing RGB values with lower bit-depth (Kingma & Dhariwal, 2018; Menick & Kalchbrenner, 2018). The Subscale Pixel Network (SPN) also performs autoregressive spatial upsampling, and makes use of Transformers (Menick & Kalchbrenner, 2018). Similar to Section 3.1, SPN also subsamples a target sequence during training and uses a separate encoder network to condition on past context in order to obtain constant memory/compute costs with respect to image size.

## 6. Discussion

We proposed using sparse DCT-based image representations for generative modelling, and introduce a Transformer-based autoregressive architecture for modelling sparse image sequences which overcomes challenges in modelling long sequences with a chunked training regime. Our DC-Transformer achieves strong performance on sample quality and diversity benchmarks, and easily supports super-resolution upsampling, as well as colorization tasks. We believe there is much to be gained from exploring the data representations used in classical data compression methods as a basis for neural generative modelling, and plan to explore related representations in audio and video domains.

There are still some challenges: For complex and high resolution datasets, good results require large models and substantial computational resources (Appendix C). This contrasts with GANs in particular, which achieve high quality results with less computational resources. However, we believe the computational disparity is primarily due to the challenging likelihood maximization objective, which strongly incentivizes full coverage of the data distribution, and enables us to produce highly diverse outputs.

## References

- Ahmed, N., Natarajan, T., and Rao, K. R. Discrete cosine transform. *IEEE transactions on Computers*, 100(1):90–93, 1974.
- Antic, J. DeOldify. <https://github.com/jantic/DeOldify>. Accessed: 2020-02-03.
- Bachlechner, T., Majumder, B. P., Mao, H. H., Cottrell, G. W., and McAuley, J. J. Rezero is all you need: Fast convergence at large depth. *CoRR*, abs/2003.04887, 2020.
- Brock, A., Donahue, J., and Simonyan, K. Large scale GAN training for high fidelity natural image synthesis. In *ICLR*. OpenReview.net, 2019.
- Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In *NeurIPS*, 2020.
- Chen, M., Radford, A., Child, R., Wu, J., Jun, H., Luan, D., and Sutskever, I. Generative pretraining from pixels. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pp. 1691–1703. PMLR, 2020.
- Child, R., Gray, S., Radford, A., and Sutskever, I. Generating long sequences with sparse transformers. *CoRR*, abs/1904.10509, 2019.
- Choromanski, K., Likhoshesterov, V., Dohan, D., Song, X., Gane, A., Sarlós, T., Hawkins, P., Davis, J., Mohiuddin, A., Kaiser, L., Belanger, D., Colwell, L., and Weller, A. Rethinking attention with performers. In *ICLR*, 2021.
- Chouhan, S. S., Kaul, A., Singh, U. P., and Jain, S. A database of leaf images: Practice towards plant conservation with plant pathology. Mendeley Data, 2019.
- De, S. and Smith, S. L. Batch normalization biases residual blocks towards the identity function in deep networks. In *NeurIPS*, 2020.
- Dhariwal, P., Jun, H., Payne, C., Kim, J. W., Radford, A., and Sutskever, I. Jukebox: A generative model for music, 2020.
- Du, Y. and Mordatch, I. Implicit generation and modeling with energy based models. In *NeurIPS*, 2019.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. C., and Bengio, Y. Generative adversarial nets. In *NIPS*, pp. 2672–2680, 2014.
- Google, 2018. Cloud tpu. <https://cloud.google.com/tpu/>. Accessed: 2021.
- Henighan, T., Kaplan, J., Katz, M., Chen, M., Hesse, C., Jackson, J., Jun, H., Brown, T. B., Dhariwal, P., Gray, S., Hallacy, C., Mann, B., Radford, A., Ramesh, A., Ryder, N., Ziegler, D. M., Schulman, J., Amodei, D., and McCandlish, S. Scaling laws for autoregressive generative modeling, 2020.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NIPS*, pp. 6626–6637, 2017.
- Independent JPEG Group. Independent JPEG Group. <https://ijg.org/>. Accessed: 2020-02-03.

- Johnson, J., Hariharan, B., van der Maaten, L., Fei-Fei, L., Zitnick, C. L., and Girshick, R. B. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In *CVPR*, pp. 1988–1997. IEEE Computer Society, 2017.
- Kaggle and EyePacs. Kaggle diabetic retinopathy detection, jul 2015. URL <https://www.kaggle.com/c/diabetic-retinopathy-detection/data>.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. Scaling laws for neural language models, 2020.
- Karras, T., Aila, T., Laine, S., and Lehtinen, J. Progressive growing of gans for improved quality, stability, and variation. In *ICLR*. OpenReview.net, 2018.
- Karras, T., Laine, S., and Aila, T. A style-based generator architecture for generative adversarial networks. In *CVPR*, pp. 4401–4410. Computer Vision Foundation / IEEE, 2019.
- Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., and Aila, T. Analyzing and improving the image quality of stylegan. In *CVPR*, pp. 8107–8116. IEEE, 2020.
- Kersten, D. Predictability and redundancy of natural images. *JOSA A*, 4(12):2395–2400, 1987.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *ICLR (Poster)*, 2015.
- Kingma, D. P. and Dhariwal, P. Glow: Generative flow with invertible 1x1 convolutions, 2018.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. In *ICLR*, 2014.
- Kitaev, N., Kaiser, L., and Levskaya, A. Reformer: The efficient transformer. In *ICLR*. OpenReview.net, 2020.
- Kuznetsova, A., Rom, H., Alldrin, N., Uijlings, J. R. R., Krasin, I., Pont-Tuset, J., Kamali, S., Popov, S., Mallocci, M., Duerig, T., and Ferrari, V. The open images dataset V4: unified image classification, object detection, and visual relationship detection at scale. *CoRR*, abs/1811.00982, 2018.
- Kynkäänniemi, T., Karras, T., Laine, S., Lehtinen, J., and Aila, T. Improved precision and recall metric for assessing generative models. In *NeurIPS*, pp. 3929–3938, 2019.
- Mandava, S., Migacz, S., and Florea, A. F. Pay attention when required. *CoRR*, abs/2009.04534, 2020.
- Menick, J. and Kalchbrenner, N. Generating high fidelity images with subscale pixel networks and multidimensional upscaling. *CoRR*, abs/1812.01608, 2018. URL <http://arxiv.org/abs/1812.01608>.
- NVIDIA. NVIDIA DLSS 2.0: A Big Leap In AI Rendering. <https://www.nvidia.com/en-gb/geforce/news/nvidia-dlss-2-0-a-big-leap-in-ai-rendering/>. Accessed: 2020-02-03.
- Oord, A. v. d., Vinyals, O., and Kavukcuoglu, K. Neural discrete representation learning. *arXiv preprint arXiv:1711.00937*, 2017.
- Parisotto, E., Song, H. F., Rae, J. W., Pascanu, R., Gülçehre, Ç., Jayakumar, S. M., Jaderberg, M., Kaufman, R. L., Clark, A., Noury, S., Botvinick, M., Heess, N., and Hadsell, R. Stabilizing transformers for reinforcement learning. In *ICML*, volume 119 of *Proceedings of Machine Learning Research*, pp. 7487–7498. PMLR, 2020.
- Ramesh, A., Pavlov, M., Goh, G., Gray, S., Chen, M., Child, R., Misra, V., Mishkin, P., Krueger, G., Agarwal, S., and Sutskever, I. DALL·E: Creating Images from Text. <https://openai.com/blog/dall-e/>. Accessed: 2020-02-03.
- Rao, K. R. and Yip, P. C. *Discrete Cosine Transform - Algorithms, Advantages, Applications*. 1990.
- Razavi, A., van den Oord, A., and Vinyals, O. Generating diverse high-fidelity images with VQ-VAE-2. In *NeurIPS*, pp. 14837–14847, 2019.
- Rezende, D. J. and Mohamed, S. Variational inference with normalizing flows. In *ICML*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 1530–1538. JMLR.org, 2015.
- Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, volume 32 of *JMLR Workshop and Conference Proceedings*, pp. 1278–1286. JMLR.org, 2014.
- Roy, A., Saffar, M., Vaswani, A., and Grangier, D. Efficient content-based sparse attention with routing transformers, 2020.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M. S., Berg, A. C., and Li, F. Imagenet large scale visual recognition challenge. *Int. J. Comput. Vis.*, 115(3):211–252, 2015.
- Shannon, C. E. Prediction and entropy of printed english. *Bell system technical journal*, 30(1):50–64, 1951.

- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. Rethinking the inception architecture for computer vision. In *CVPR*, pp. 2818–2826. IEEE Computer Society, 2016.
- Van Oord, A., Kalchbrenner, N., and Kavukcuoglu, K. Pixel recurrent neural networks. In *International Conference on Machine Learning*, pp. 1747–1756. PMLR, 2016.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention is all you need. *arXiv preprint arXiv:1706.03762*, 2017.
- Wallace, G. K. The jpeg still picture compression standard. *IEEE transactions on consumer electronics*, 38(1):xxviii–xxxiv, 1992.
- Wang, Z., Bovik, A. C., Sheikh, H. R., and Simoncelli, E. P. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004.
- Yu, F., Zhang, Y., Song, S., Seff, A., and Xiao, J. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.
- Zaheer, M., Guruganesh, G., Dubey, A., Ainslie, J., Alberti, C., Ontanon, S., Pham, P., Ravula, A., Wang, Q., Yang, L., and Ahmed, A. Big bird: Transformers for longer sequences, 2021.

## A. Quantization matrices

Following the [Independent JPEG Group](#) standard we use the following quality parameterized quantization matrix  $\mathbf{Q}$  to quantize DCT pixel blocks. For quality  $q$  in  $[1, 100]$  the matrix is defined as:

$$\mathbf{T}_{\text{luma}} = \begin{bmatrix} 16 & 11 & 10 & 16 & 24 & 40 & 51 & 61 \\ 12 & 12 & 14 & 19 & 26 & 58 & 60 & 55 \\ 14 & 13 & 16 & 24 & 40 & 57 & 69 & 56 \\ 14 & 17 & 22 & 29 & 51 & 87 & 80 & 62 \\ 18 & 22 & 37 & 56 & 68 & 109 & 103 & 77 \\ 24 & 35 & 55 & 64 & 81 & 104 & 113 & 92 \\ 49 & 64 & 78 & 87 & 103 & 121 & 120 & 101 \\ 72 & 92 & 95 & 98 & 112 & 100 & 103 & 99 \end{bmatrix}$$

$$s(q) = \begin{cases} 5000/q, & \text{if } q < 50 \\ 200 - 2q, & \text{otherwise} \end{cases}$$

$$\mathbf{Q}_{\text{luma}}(q) = \text{floor} \left( \frac{s(q)\mathbf{T}_{\text{luma}} + 50}{100} \right)$$

For the chrominance components we replace  $\mathbf{T}_{\text{luma}}$  with an alternative base matrix  $\mathbf{T}_{\text{chroma}}$  which provides stronger quantization:

$$\mathbf{T}_{\text{chroma}} = \begin{bmatrix} 17 & 18 & 24 & 47 & 99 & 99 & 99 & 99 \\ 18 & 21 & 26 & 66 & 99 & 99 & 99 & 99 \\ 24 & 26 & 56 & 99 & 99 & 99 & 99 & 99 \\ 47 & 66 & 99 & 99 & 99 & 99 & 99 & 99 \\ 99 & 99 & 99 & 99 & 99 & 99 & 99 & 99 \\ 99 & 99 & 99 & 99 & 99 & 99 & 99 & 99 \\ 99 & 99 & 99 & 99 & 99 & 99 & 99 & 99 \\ 99 & 99 & 99 & 99 & 99 & 99 & 99 & 99 \end{bmatrix}.$$

When using block sizes other than 8, we use nearest neighbour interpolation to resize the base matrices  $\mathbf{T}_{\text{luma}}$  and  $\mathbf{T}_{\text{chroma}}$  to the target size.

## B. Architecture details

DCTransformer consists of a Transformer encoder that processes partial DCT images, and three stacked Transformer decoders that process slices from the DCT co-ordinate list (Section 3.2). We use a number of modifications to the original architecture that we found to improve stability, training speed, and memory consumption:

**Layer norm placement** Following [Child et al. \(2019\)](#) and [Parisotto et al. \(2020\)](#) we use Transformer blocks with layer norm placed inside the residual path, rather than applying layer norm

**ReZero** We use ReZero ([Bachlechner et al., 2020](#); [De & Smith, 2020](#)), multiplying each residual connection with a

zero-initialized scalar value, which is optimized jointly with the model parameters. In our experiments we found this to improve training speed and stability to a small degree. The combination of ReZero and our chosen layer norm placement results in residual connections of the following form:

$$\mathbf{H}_l = \mathbf{H}_{l-1} + \alpha_l f_l(\text{layernorm}(\mathbf{H}_{l-1})), \quad (10)$$

where  $\mathbf{H}_l$  is a sequence of activations at layer  $l$ , and  $f$  is the residual function.

**PAR Transformer** [Mandava et al. \(2020\)](#) showed through Neural architecture search that the default alternation of fully connected and self attention layers in Transformer blocks is sub-optimal with respect to performance-speed trade offs. Based on the results of the search they proposed the PAR Transformer, which applies a series of fully connected layers after each self-attention layer, resulting in improved inference speed, and memory savings. We use a PAR Transformer style architecture in DCTransformer encoder and decoders, and while we didn't experiment rigorously with the ratio of fully-connected to self-attention layers, we found that architectures using 2-4 fully connected layers per self-attention layer helped to boost the parameter count at a given memory budget.

Table 3 details the architecture configurations used for our main experiments.

## C. Training details

**Optimization** We train our models using Adam Optimizer ([Kingma & Ba, 2015](#)) and train for a fixed number of tokens, where the number of tokens processed in a batch is the number of elements in the target chunk that we apply a loss to. We use a linear warmup of 1000 steps up to a maximum learning rate, and use a cosine decay over the course of training. We train all models using Google Cloud TPUv3 ([Google, 2018](#)), and detail the number of cores used during training in Table 3.

**Chunk selection policy** As described in Section 3.1, we train DCTransformer on input and target chunks selected from larger sequences. If the target chunk is sampled uniformly from the sequence representation of an image, the sample gradient is unbiased because it is a Monte Carlo estimate of the gradient computed on the full sequence. As discussed in Section 2, lossy image compression codecs such as JPEG preferentially discard high frequency data when allocating limited storage capacity. It stands to reason that we may benefit from making an analogous decision when allocating limited model capacity. We found it advantageous for sample quality to bias the selection of target chunks toward the beginning of the sequence, which contains more low frequency information (see Figure 2).

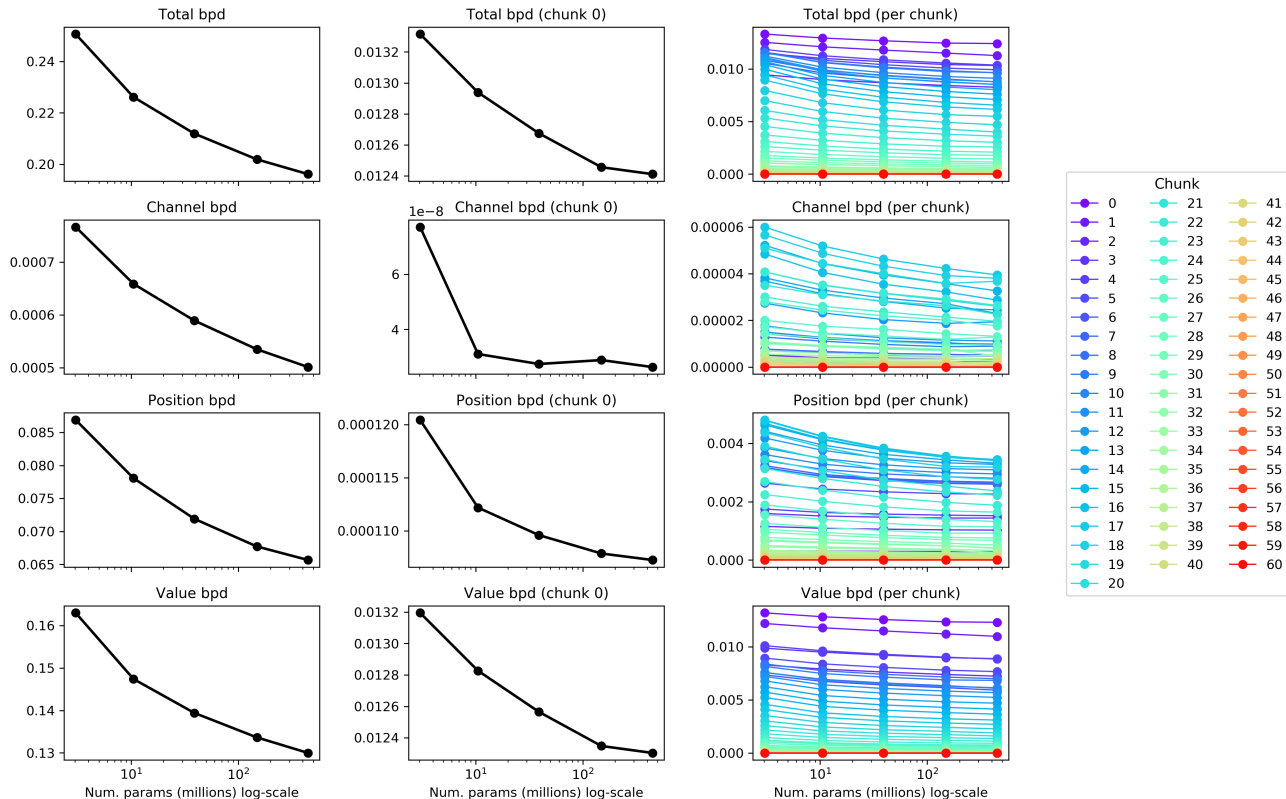


Figure 9. Model performance, reported in bits per image subpixel (bpd) as a function of model size for DCTransformer trained on LSUN bedrooms. We report the total bpd, as well as the contributions from the channel, position and value distributions. We additionally report bpd broken down by sequence chunk, where each chunk corresponds to 896 sequence elements. The model size is reported in terms of the number of total parameters, ranging from roughly 3 million to 448 million.

We select chunks using the following process: A sequence of length  $L$  is split into chunks of size  $C$ , with the final chunk containing  $L \bmod C$  elements. The probability of selecting a chunk with start position  $l$  decays to a lower limit  $p_{\min}$ , with probability proportional to a polynomial in  $l$ . By default we decay the probability of selecting a chunk beginning at position  $l$  proportional to  $l^{-3}$  down to a minimum of  $p_{\min} = 0.1$ , a hyperparameter choice we fixed early in model development and found no need to adjust.

**Sequence length bias adjustment** Chunk-based training introduces an issue in unconditional generative modelling: It biases the model towards chunks from shorter sequences. Consider the first chunk, that occurs at the very start of the sequence. For long sequences, this chunk will be selected relatively infrequently compared to short sequences. The model will therefore assign greater probability to initial chunks from short sequences, than initial chunks from long sequences.

We counter the bias by randomly filtering out sequences with a probability that is inversely proportional to the sequence length. We pick a maximum filtering sequence length  $L_{\max}$

and filter our sequences with probability:

$$p_{\text{filter}}([\mathbf{t}_l]_{l=1}^L) = \text{maximum} \left( \frac{L}{L_{\max}}, 1 \right). \quad (11)$$

## D. Model scaling properties

Kaplan et al. (2020) show that for Transformer language models, performance improves reliably subject to constraints on model size, compute and dataset size. In particular, they show that if compute and data are not bottlenecks, then test-loss performance improves log-linearly with model size. Follow-up work Henighan et al. (2020) has shown that this phenomenon applies more generally beyond just language data. We investigate the extent to which DCTransformer scales as a function of model size by training models of varying sizes on LSUN bedrooms. Figure 9 shows the results broken down by channel, position and value prediction contributions, and chunk position. We find that the total bits-per-dimension (bpd) roughly matches the expected log-linear fit, with the exception of the smallest model, which performs worse than the expected trend. Another exception is chunk 0, where the performance improvement is less than expected for the largest model.

We also find that the value predictions account for the largest portion of the total bpd, followed by positions, and finally channels, which accounts for a very small portion of the total bpd. For the total bpd, the contribution per chunk decreases with the chunk position. This is likely because the total bpd is dominated by the value bpd, and we expect the value bpd to decrease as we transition from lightly quantized low frequency components, to more heavily compressed high frequency components.

### **E. Additional samples**

Figures 10 and 11 compare uncurated samples from DC-Transformer and baselines to a random selection of real data on the FFHQ and LSUN datasets respectively. Figures 12 and 13 show uncurated upsampling and colorization results on ImageNet and OpenImagesV4 respectively.

## Generating images with sparse representations



Figure 10. Comparison between FFHQ images and uncurated model samples, all at 1024x1024 resolution.

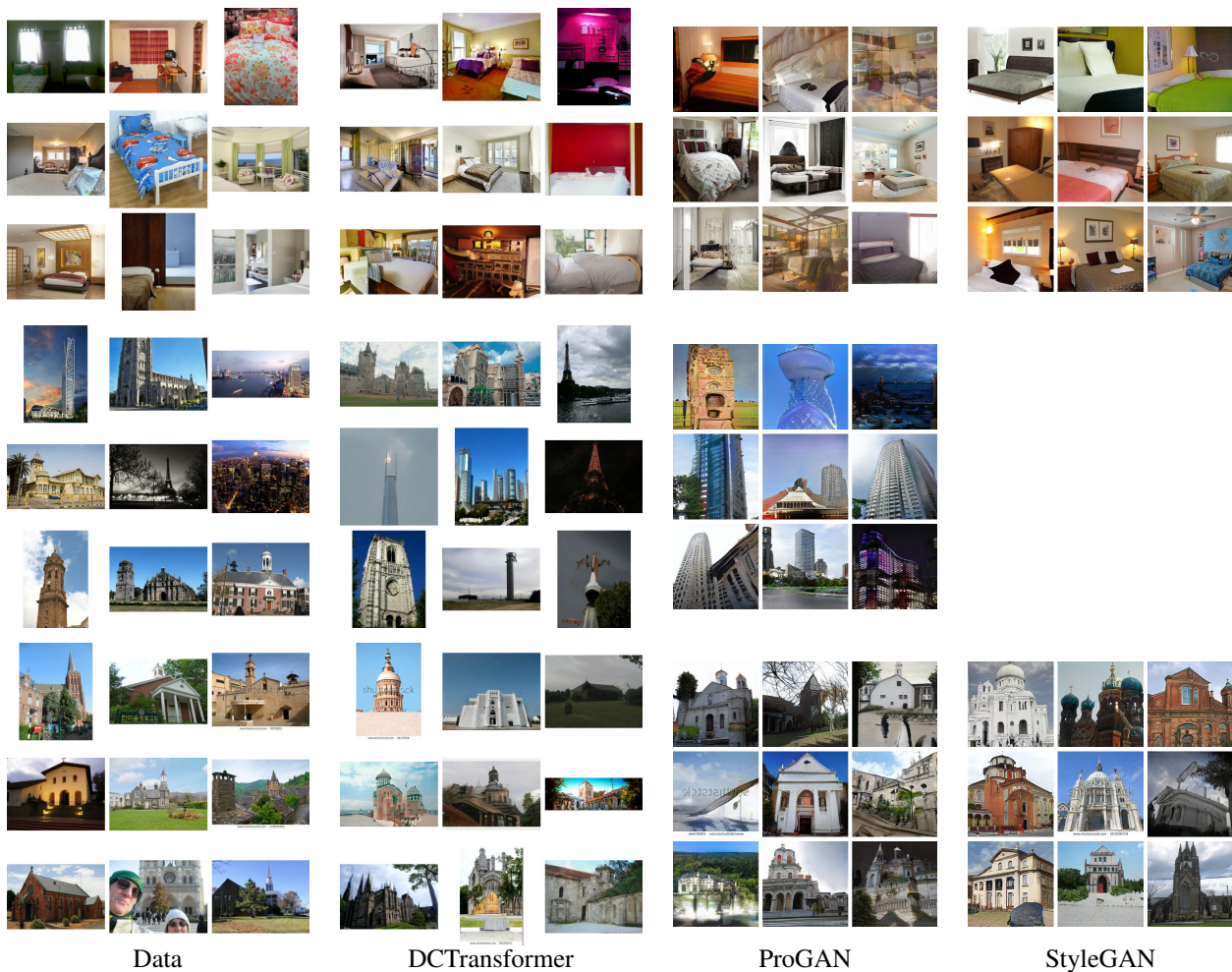


Figure 11. Comparison between LSUN images and uncurated model samples for bedroom, tower and church-outdoor subsets. StyleGAN refers to StyleGAN1 for bedrooms, and StyleGAN2 for church-outdoor samples. DCTransformer produces variable aspect ratio samples with long-side resolution 384. BigGAN and VQ-VAE are trained and sample at a fixed 256x256 resolution corresponding to a resized long-side crop of the input images. ProGAN and StyleGAN samples use truncation 1.0 to yield maximum diversity.

Generating images with sparse representations

	LSUN (all)	FFHQ	ImageNet	OpenImagesV4
Image resolution	384	1024	384	640
DCT block size	8	16	8	8
DCT quality	75	35	75	50
DCT clip value	1200	3200	1200	1200
Target chunk size	896	896	896	896
Target chunk overlap	128	128	128	128
Hidden units	896	896	1152	896
Self-attention heads	14	14	14	18
Layer spec (encoder)	[(1,2)] * 4	[(1, 2)] * 4	[(1,2)] * 4	[(1,2)] * 4
Layer spec (channel decoder)	[(1, 2)] * 3 + [(1, 4)]	[(1, 2)] * 3 + [(1, 4)]	[(1, 2)] * 3 + [(1, 4)]	[(1, 2)] * 3 + [(1, 4)]
Layer spec (position decoder)	[(1, 2)] * 3 + [(1, 4)]	[(1, 2)] * 3 + [(1, 4)]	[(1, 2)] * 3 + [(1, 4)]	[(1, 2)] * 3 + [(1, 4)]
Layer spec (value decoder)	[(1, 2)] * 5 + [(1, 7)]	[(1, 2)] * 6 + [(1, 7)]	[(1, 2)] * 6 + [(1, 7)]	[(1, 2)] * 6 + [(1, 7)]
DCT image downsampling kernel	kernel size 4, stride 2	kernel size 6, stride 3	kernel size 8, stride 4	kernel size 6, stride 3
Batch size	512	448	512	512
Dropout rate	0.1	0.1	0.01	0.01
Learning Rate Start	5e-4	5e-4	5e-4	5e-4
Tokens processed	300e9	250e9	1000e9	1000e9
Parameters	448e6	473e6	738e6	533e6
TPUv3 cores	64	64	128	64

	Plant Leaves	Retinopathy	CLEVR
Image resolution	2048	1024	480
DCT block size	32	16	8
DCT quality	50	75	90
DCT clip value	4000	3200	1200
Target chunk size	896	896	896
Target chunk overlap	128	128	128
Hidden units	768	768	896
Self-attention heads	12	12	14
Layer spec (encoder)	[(1,2)] * 4	[(1, 2)] * 4	[(1,2)] * 4
Layer spec (channel decoder)	[(1, 2)] * 3 + [(1, 3)]	[(1, 2)] * 3 + [(1, 3)]	[(1, 2)] * 3 + [(1, 4)]
Layer spec (position decoder)	[(1, 2)] * 3 + [(1, 4)]	[(1, 2)] * 3 + [(1, 4)]	[(1, 2)] * 3 + [(1, 4)]
Layer spec (value decoder)	[(1, 2), (1, 2), (1, 3), (1, 3), (1, 7)]	[(1, 2), (1, 2), (1, 3), (1, 3), (1, 7)]	[(1, 2)] * 6 + [(1, 7)]
DCT image downsampling kernel	kernel size 4, stride 2	kernel size 4, stride 2	kernel size 6, stride 3
Batch size	256	512	128
Dropout rate	0.4	0.1	0.5
Learning Rate Start	5e-4	5e-4	5e-4
Tokens processed	100e9	200e9	50e9
Parameters	325e6	318e6	483e6
TPUv3 cores	64	64	32

Table 3. Model and training hyperparameters. The layer spec for the Transformer encoder and decoders is a list of tuples, where each tuple describes the number of self-attention layers, followed by the number of fully-connected layers in a Transformer block. For example [(1, 2)] \* 3 + [(1, 4)] expands to [(1,2),(1,2),(1,2),(1,4)], and corresponds to four Transformer blocks, where the first three blocks consist of a single self-attention layer, followed by two fully-connected layers. The final block has one self-attention layer followed by four fully-connected layers.



## Generating images with sparse representations

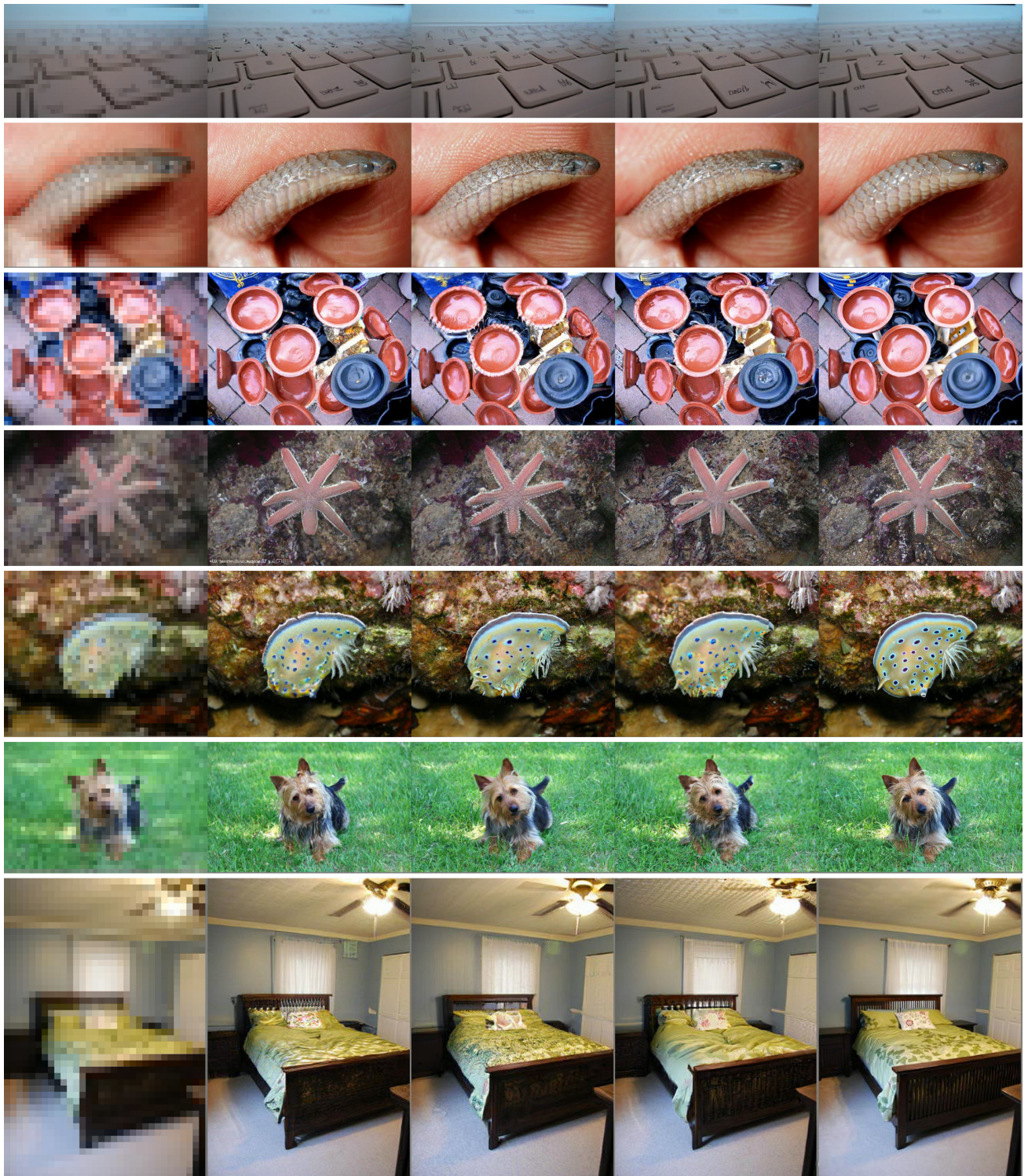


Figure 12. Uncurated 8x image upsampling results on ImageNet validation set. (left) input downsampled image, (middle) three samples generated by DCTransformer, (right) original image.

## Generating images with sparse representations



Figure 13. Uncurated image colorization results on OpenImagesV4 validation set. (left) input grayscale image, (middle) three samples generated by DCTransformer, (right) original image.