

# You're your own best teacher: A Self-Supervised Learning Approach For Expressive Representations

Johan Vik Mathisen

May 22, 2024

**TODO:** Introduce this section. Why are these models presented?

## 0.1 MaskGIT

The Masked Generative Image Transformer (MaskGIT)[3] is a generative transformer model for image synthesis developed by Google Research. The novelty of the model lies in the token generation. Unlike popular autoregressive generative transformers, who treat images as a sequence of tokens, MaskGIT introduces an image synthesis paradigm using a bi-directional transformer. This means that during training MaskGIT learns to predict tokens in all directions, an intuitively more natural way to consider images. At inference time MaskGIT starts out with a blank canvas and predicts the entire image, and iteratively keeps and conditions on the most confident pixels.

MaskGIT assumes a tokenization procedure for stage 1. In the original paper [3] VQGAN [2] was used and the actual contribution of the work revolved around improving stage 2, hence we present that part only.

### 0.1.1 Masked Visual Token Modeling (Prior learning)

For prior learning the codebook learned in the tokenization procedure is provided with a masking vector, which is the embedding of the special masking token, which we denote by  $\mathbf{M}$ . The input embedding in the bidirectional transformer is initialized with this expanded codebook. For some image  $X$  in the dataset  $\mathcal{D}$ , let  $z = \{z_{k_i}\}_{i=1}^N$  denote the sequence of codewords obtained by passing  $X$  through the VQ-Encoder. Such a sequence can equivalently be described as a sequence of indices  $s = \{k_i\}_{i=1}^N$ . The prior learning amounts to masking such a sequence and training the bidirectional transformer to predict the masked indices.

Let  $s = \{k_i\}_{i=1}^N$  be the sequence of indices described above and denote the corresponding binary mask by  $M = \{m_i\}_{i=1}^N$ . During training a subset of  $s$  is replaced by the masking token  $\mathbf{M}$  according to the binary mask  $M$ . This is done by

$$s_{\text{Mask}} = s \odot (1_N - M) + M \cdot \mathbf{M}, \quad (1)$$

where  $\odot$  is the Hadamard product, i.e point wise multiplication, and  $1_N$  is a vector with the same shape as  $M$  and  $s$ .

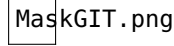
The sampling procedure, or choice number of tokens to mask, is parameterized by a mask scheduling function  $\gamma$ . The sampling can be summarized as follows

- Sample  $r \sim U(0, 1]$ .
- Sample  $\lceil \gamma(r) \cdot N \rceil$  indices  $I$  uniformly from  $\{0, \dots, N-1\}$  without replacement.
- Create  $M$  by setting  $m_i = 1$  if  $i \in I$ , and  $m_i = 0$  otherwise.

The training objective is to minimize the negative log likelihood of the masked tokens, conditional on the unmasked.

$$\mathcal{L}_{\text{Mask}} = -\mathbb{E}_{s \in \mathcal{D}} \left[ \sum_{i \in I} p(s_i | s_{\text{Mask}}) \right] \quad (2)$$

The bidirectional transformer is used to predict the probabilities  $p(s_i | s_{\text{Mask}})$  of each masked token, and  $\mathcal{L}_{\text{Mask}}$  is computed as the cross entropy between the ground truth one-hot token and the predicted token probabilities.



**Figure 1:** MaskGIT forward computation.

### 0.1.2 Iterative decoding (Image generation)

The bi-directional transformer could in principle predict all masked tokens and generate a sample in a single pass by simply sampling from the predicted probabilities  $p(\hat{s}_i | s_{\text{Mask}})$  from a forward pass of an all masked sequence. However, there are challenges with this approach. In their original article [3] proposes a novel non-autoregressive decoding method to synthesize samples in a constant number of steps.

The decoding process goes from  $t = 0$  to  $T$ . To generate a sample at inference time one starts out with an all masked sequence which we denote by  $s_{\text{Mask}}^{(0)}$ . At iteration  $t$  the model predicts the probabilities for all the mask tokens,  $p(\hat{s}_i | s_{\text{Mask}}^{(t)})$ , in parallel. At each masked index  $i$  a token  $s_i^{(t)}$  is sampled according to the predicted distribution, and the corresponding probability  $c_i^{(t)}$  is used as a measure of the confidence in the sample. For the unmasked tokens a confidence of 1 is assigned to the true position. The number of  $s_i^{(t)}$  with highest confidence kept for the next iteration is determined by the mask scheduling function. We mask  $n = \lceil \gamma(t/T) \cdot N \rceil$  of the lower confidence tokens by calculating  $M^{(t+1)}$  by

$$m_i^{(t+1)} = \begin{cases} 1, & \text{if } c_i < \text{Sort}([c_1^{(t)}, \dots, c_N^{(t)}])[n] \\ 0, & \text{otherwise} \end{cases} \quad (3)$$



**Figure 2:** Illustration of first pass of the iterative decoding algorithm.

The algorithm synthesizes a full image in  $T$  steps. For image generation, cosine scheduling function proved best across all experiments in the original paper.

## 0.2 TimeVQVAE

TimeVQVAE is a time series generation model based on VQVAE and MaskGIT. It is the first to our and the authors knowledge that utilizes vector quantization (VQ) to address the TSG problem. It leverages a two stage approach similar to VQVAE and uses a bidirectional transformer akin to MaskGIT for prior learning. Additionally, the authors propose VQ modeling in time-frequency domain, separating data into high and low frequency components to better retain temporal consistencies and generate higher quality samples.

The contributions TimeVQVAE presents is, in addition to VQ-modeling in time-frequency domain, a process of sampling jointly from high and low frequency latent spaces and guided class-conditional sampling. By appending a class token, similarly to [dosovitskiy2021image], the prior is learned such that the model can generate synthetic samples both conditionally and unconditionally.

Our work consists of extending a variation of the TimeVQVAE model without the high-low frequency split. This reduces the prior learning method to MaskGIT, with the addition of guided class-conditional sampling. Hence we present only the tokenization stage and refer the reader to [VQVAE] for the prior model training.

### 0.2.1 Tokenization

The tokenization stage is similar to VQVAE presented in section ?? except for the frequency split. An overview of the model is presented in figure ?. First a time series is mapped to time-frequency domain using the Short-time Fourier Transform (STFT). Then the time-frequency representation is separated into a two branches, one zero-padding the HF region and the other zero-padding the LF region. From here the two branches follow the VQVAE architecture, with separate encoders, decoders and codebooks denoted by  $E_{LF}$ ,  $E_{HF}$ ,  $D_{HF}$ ,  $D_{LF}$  and  $Z_{LF}$ ,  $Z_{HF}$  respectively. The output of the decoders are again zero-padded giving  $\hat{u}_{LF}$  and  $\hat{u}_{HF}$ , before being mapped back to time domain by the Inverse Short-time Fourier Transform (ISTFT) to produce the reconstructed HF and LF components,  $\hat{x}_{LF}$  and  $\hat{x}_{HF}$ , of the time series.

TimeVQVAE1.png

**Figure 3:** Stage 1: Tokenization. Figure taken with permission from [1]

## Loss

The codebook loss of TimeVQVAE is similar to codebook loss presented section ?? equation ?? but reflects the HF-LF split

$$\begin{aligned}\mathcal{L}_{\text{codebook}} = & \| \text{sg}[E_{LF}(\mathcal{P}_{LF}(\text{STFT}(x)))] - z_q^{\text{LF}} \|_2^2 \\ & + \| \text{sg}[E_{HF}(\mathcal{P}_{HF}(\text{STFT}(x)))] - z_q^{\text{HF}} \|_2^2 \\ & + \beta \| E_{LF}(\mathcal{P}_{LF}(\text{STFT}(x))) - \text{sg}[z_q^{\text{LF}}] \|_2^2 \\ & + \beta \| E_{HF}(\mathcal{P}_{HF}(\text{STFT}(x))) - \text{sg}[z_q^{\text{HF}}] \|_2^2,\end{aligned}\tag{4}$$

The reconstruction loss is performed both on time and time-frequency reconstructions, and is given by

$$\begin{aligned}\mathcal{L}_{\text{recons}} = & \| x_{\text{LF}} - \hat{x}_{\text{LF}} \|_2^2 + \| x_{\text{HF}} - \hat{x}_{\text{HF}} \|_2^2 \\ & + \| u_{\text{LF}} - \hat{u}_{\text{LF}} \|_2^2 + \| u_{\text{HF}} - \hat{u}_{\text{HF}} \|_2^2.\end{aligned}\tag{5}$$

The total loss is given by

$$\mathcal{L}_{\text{VQ}} = \mathcal{L}_{\text{codebook}} + \mathcal{L}_{\text{recons}}.\tag{6}$$

In order to update the codebooks TimeVQVAE uses an exponential moving average presented in appendix A.1 of [VQVAE].

## 0.3 SSL

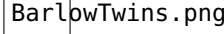
Our model leverages SSL algorithms in order to learn more expressive latent representations. Here we present the relevant algorithms for our work, Barlow Twins and VibCReg.

### 0.3.1 Barlow Twins

Barlow Twins is a non-contrastive SSL method based on applying the *redundancy-reduction principle* (or efficient coding hypothesis) [5] from the neuroscientist H. Barlow to a pair of identical networks.

In essence the model wants to encourage representations of similar samples to be similar, while simultaneously reducing the amount of redundancy between the components of the vectors. This is enforced by producing two augmented views of each sample and projecting the their representations onto a vast feature space, in such a way that their cross-correlation is close to the identity.

The Barlow Twins algorithm starts out by creating two different augmented views for each datapoint in a batch  $D$ . The augmentations are selected by sampling from a collection of augmentations  $\mathcal{T}$ . We denote the batches of augmented views



**Figure 4:** Overview of the Barlow Twins architecture. Figure inspired by [6]

$T(D) = X$  and  $T'(D) = X'$ , for augmentations  $T, T' \sim \mathcal{T}$ . The batches are then passed through an encoder (give representations  $Y$  and  $Y'$ ) and a *projector* to produce batches of embeddings  $Z$  and  $Z'$ . The embeddings are assumed to be mean centered across the batch dimension.

The loss function is calculated using the cross correlation matrix  $\mathcal{C}$  between  $Z$  and  $Z'$ , and measuring its deviance from the identity. In particular the Barlow Twins loss is defined as

$$\mathcal{L}_{\text{BT}} = \overbrace{\sum_i (1 - \mathcal{C}_{ii})^2}^{\text{Invariance}} + \lambda \overbrace{\sum_i \sum_{j \neq i} \mathcal{C}_{ij}^2}^{\text{Redundancy reduction}}, \quad (7)$$

where

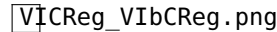
$$\mathcal{C}_{ij} = \frac{\sum_b z_{b,i} z'_{b,j}}{\sqrt{\sum_b (z_{b,i})^2} \sqrt{\sum_b (z'_{b,j})^2}}. \quad (8)$$

The *invariance term* assists in making the embedding invariant to the distortions introduced by the augmentations, hence pushes the representations closer together. The *redundancy reduction term* decorrelates the different vector components, which reduces the information redundancy.

### 0.3.2 VlbCReg

VlbCReg [7] is a non-contrastive SSL model with siamese architecture based on VICReg [8]. It can be seen as VICReg with better covariance regularization and IterNorm [9]. Overall the architecture is similar to Barlow Twins, but a key difference is that variance/covariance regularization is done in each branch individually.

As before a batch  $D$  is augmented to create two views and passed through an encoder and projector. The embedding  $Z$  and  $Z'$  are *whitened* using IterNorm [9].



**Figure 5:** Overview of VlbCReg, and comparison with VICReg. Taken with permission from [lee2024computer]

The loss consists of a similarity loss between the branches, and feature decoration (FD) loss together with a feature component expressiveness (FCE) term

define or elaborate on this

at each branch. Input data is processed in batches. Let  $Z \in \mathbb{R}^{B \times F}$  where  $B$  and  $F$  denotes the batch and feature sizes respectively. We denote a row in  $Z$  by  $Z_b$  and column by  $Z_f$ , and similarly for  $Z'$ .

The similarity loss is defined as the MSE of the two embeddings

$$s(Z, Z') = \frac{1}{B} \sum_{b=1}^B \|Z_b - Z'_b\|_2^2, \quad (9)$$

which encourages them to be similar. The FcE term acts on each branch separately and encourages the variation across a batch to stay at a specified level  $\gamma$ . It is defined as

$$v(Z) = \frac{1}{F} \sum_{f=1}^F \max(0, \gamma - \sqrt{\text{Var}(Z_f) + \epsilon}), \quad (10)$$

where  $\text{Var}()$  is a variance estimator,  $\gamma$  is a target value for the standard deviation, which both in VIBReg and VICReg is set to 1.  $\epsilon$  is a small scalar preventing numerical instabilities.

For the FD loss we first mean shift and normalize along the batch dimension

$$\hat{Z}_b = \frac{Z_b - \bar{Z}}{\|Z_b - \bar{Z}\|_2} \text{ where } \bar{Z} = \frac{1}{B} \sum_{b=1}^B Z_b, \quad (11)$$

$$\hat{Z} = [\hat{Z}_1, \dots, \hat{Z}_B]^T, \quad (12)$$

compute the normalized covariance matrix

$$C(Z) = \frac{1}{B-1} \hat{Z}^T \hat{Z}, \quad (13)$$

and take the mean square across all off-diagonal elements to obtain the FD loss

$$c(Z) = \frac{1}{F^2} \sum_{i \neq j} C(Z)_{ij}^2. \quad (14)$$

The total loss is then given by

$$\mathcal{L}_{\text{VIBReg}} = \lambda s(Z, Z') + \mu [v(Z) + v(Z')] + \nu [c(Z) + c(Z')] \quad (15)$$

where  $\lambda, \mu$  and  $\nu$  are hyperparameters determining the importance of each term. The normalization of the covariance matrix keeps the range of the FD loss small, independent of data, and eases hyperparameter tuning across datasets.