# DiT / Diffusion Transformer Architecture Notes

## Left diagram (DiT architecture)

- Linear and reshape
- layer norm — nn.LayerNorm()
- dit-blocks
- position embed — nn.parameter(torch.randn(1, self.patch_count**2, emb_size))
- patchify
  - nn.Linear(in_features=self.patch_size**2*channel, out_features=emb_size)
  - view
  - permute(0, 2, 3, 1)
  - conv2d(inchannel, outchannel 3, 3*4*4)
  - patch_size = 4 = stride
  - patch_count = img_size//pathc_size
- Noised Latent
- input x — batch, channel, img_size, img_size / batch,3,256,256

- x: batch, 256/4 * 256/4, 3*4*4
- reshape (x.size(0), self.channel, self.patch_count*patch_size, self.patch_count*patch_size) — permute — view — emb_size, self.patch_size**2*channel
- x: batch, 256, 256
- x: batch,256/4 * 256/4, emb_size
- x: batch,256/4 * 256/4, 3*4*4

- time step T
  - T: freqexp(-log(max_period*i)/half) i ∈[0, half-t] embedding i(t) = [cos(i*freq), sin(t, freq i)], t ∈[0, half-t]
  - T: batch,emb_size
  - nn.Sequential( TimeEmbedding(emb_size), nn.Linear(emb_size,emb_size), nn.ReLU(), nn.Linear(emb_size,emb_size))
- label y
  - Y: batch, emb_size
  - nn.Embedding(num_embeddings=self.label_num, embedding_dim=emb_size)

## Middle-right diagram (attention block)

- Y: (batch,seq_len,emb_size)
- nn.Linear(nhead*emb_size,emb_size)
- (batch,seq_len,nhead*emb_size)
- reshape
- permute
- Y=attn*v — (batch,nhead,seq_len,emb_size)
- attn=torch.softmax(q@k/math.sqrt(q.size(2)), dim=-1) (batch,nhead,seq_len,seq_len)
- q:(batch,nhead,seq_len, emb_size) / k:(batch,nhead,emb_size,seq_len) / v:(batch,nhead,seq_len, emb_size)
- view
- q:(batch,seq_len,nhead*emb_size) / k:(batch,seq_len,nhead*emb_size) / v:(batch,seq_len,nhead*emb_size)
- nn.Linear (emb_size, nhead*emb_size)

## Right diagram (DiT block with adaLN)

- Z=Z+Y  +
- sacle
- points wise feedforward — nn.Sequential( nn.Linear(emb_size, emb_size*4), nn.ReLU(), nn.Linear(emb_size*4, emb_size) )  α2*Z
- scale shift  (1+γ2)Z+β2
- layer norm  Z
- Y=X+Y  +
- scale  α1*Y
- multi-head self-attention
- scale shift  (1+γ1)Y+β1
- layer norm  Y
- MLP — nn.Linear(emb_size, emb_size)
- Input tokens — x: batch,256/4*256/4, emb_size (batch, self.patch_count**2, emb_size)
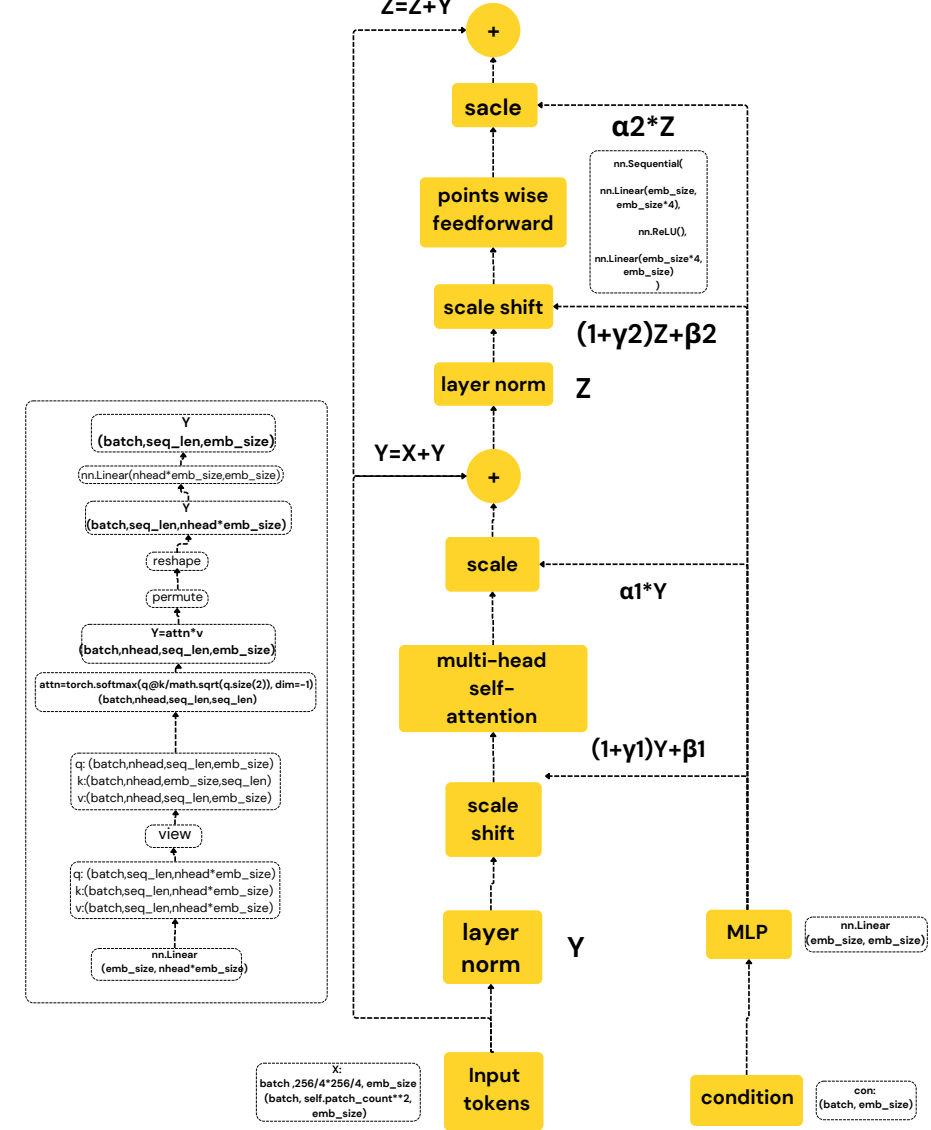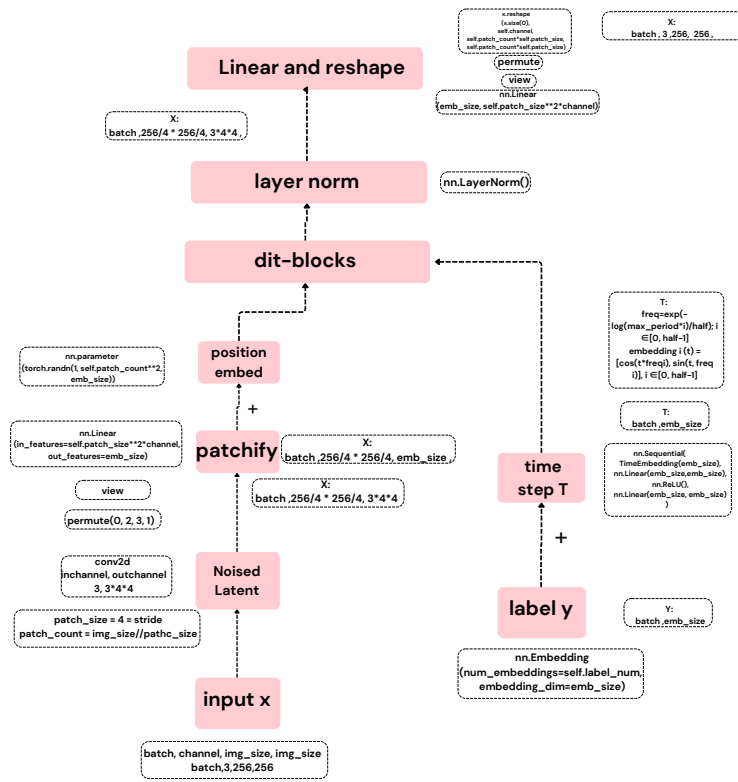- condition — con: (batch, emb_size)

Figure 3. The Diffusion Transformer (DiT) architecture. Left: We train conditional latent DiT models. The input latent is decomposed into patches and processed by several DiT blocks. Right: Details of our DiT blocks. We experiment with variants of standard transformer blocks that incorporate conditioning via adaptive layer norm, cross-attention and extra input tokens. Adaptive layer norm works best.

Latent Diffusion Transformer / DiT Block with adaLN-Zero / DiT Block with Cross-Attention / DiT Block with In-Context Conditioning

## Left text (code notes)

torch.cumprod(self.alpha, dim=-1)

self.alpha_cumprod[t].view(x.size(0), 1, 1, 1)

noise = torch.randn_like(x)

alphas_cumprod_prev=torch.cat((torch.tensor([1.0]),alphas_cumprod[:-1]),dim=-1) # alpha_t-1累乘 (T,), [1,a1,a1*a2,a1*a2*a3,.....]

variance=(1-alphas)*(1-alphas_cumprod_prev)/(1-alphas_cumprod) # denoise用的方差 (T,)

## For each training step:

**1. Randomly select a time step & encode it**

$t = 14$ —encode→ Time step embedding: 0.12 0.31 0.34 ⋯ 0.02

**2. Add noise to image**

Noisy image = Original image + Gaussian noise

$\varepsilon \sim \mathcal{N}(0,1)$

$\alpha_t = 1 - \beta_t$

$\bar{\alpha}_t = \prod_{i=1}^{t} \alpha_i$

$x_t = \sqrt{\bar{a}_t}\, x_0 + \sqrt{1-\bar{a}_t}\,\varepsilon$

Adjust the amount of noise according to the time step $t$

**3. Train the UNet**

Noisy image → UNet → Predicted noise --loss-- True noise $\varepsilon$
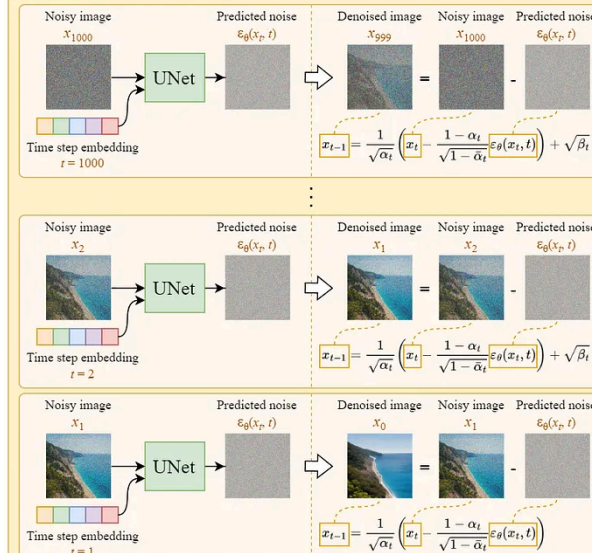
Time step embedding

## Reverse Diffusion / Denoising / Sampling

**1. Sample a Gaussian noise**

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

E.g. $T = 1000$, $x_{1000} \sim \mathcal{N}(0, I)$

**2. Iteratively denoise the image**

Noisy image $x_{1000}$ → UNet → Predicted noise $\varepsilon_\theta(x_T, t)$ ⇒ Denoised image $x_{999}$ = Noisy image $x_{1000}$ - Predicted noise $\varepsilon_\theta(x_T, t)$

Time step embedding $t = 1000$

$x_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon_\theta(x_t,t)\right) + \sqrt{\beta_t}\,\varepsilon$

Noisy image $x_2$ → UNet → Predicted noise $\varepsilon_\theta(x_T, t)$ ⇒ Denoised image $x_1$

Time step embedding $t = 2$

$x_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon_\theta(x_t,t)\right) + \sqrt{\beta_t}\,\varepsilon$

Noisy image $x_1$ → UNet → Predicted noise $\varepsilon_\theta(x_T, t)$ ⇒ Denoised image $x_0$

Time step embedding $t = 1$

$x_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\varepsilon_\theta(x_t,t)\right)$

**3. Output the denoised image**

Denoised image $x_0$