

Partition Generative Modeling

Deschenaux et al. 2025

Reviewed & Presented by Joon Hyeok Kim

Contents

1. Recap of **D3PM**
2. Recap of **MDLM** based on **D3PM**
3. **PGM**'s Idea based on **MDLM**
4. Implementation Details

D3PM : General but heavy framework!

Forward transition probabilities matrix

- $[Q_t]_{ij} = q(x_t = j \mid x_{t-1} = i)$

$$Q_t = \begin{bmatrix} q(x_t = 1 \mid x_{t-1} = 1) & q(x_t = 2 \mid x_{t-1} = 1) & \cdots & q(x_t = N \mid x_{t-1} = 1) \\ q(x_t = 1 \mid x_{t-1} = 2) & q(x_t = 2 \mid x_{t-1} = 2) & \cdots & q(x_t = N \mid x_{t-1} = 2) \\ \vdots & \vdots & \ddots & \vdots \\ q(x_t = 1 \mid x_{t-1} = N) & q(x_t = 2 \mid x_{t-1} = N) & \cdots & q(x_t = N \mid x_{t-1} = N) \end{bmatrix}$$

- Props.)
 - Massive $Q_t \in N^2$ where N is the size of the vocabulary

D3PM continues

Question) What Q should we choose?

Answer) **Absorbing State** Transition Matrix

- Def.)

- For the absorbing state m ,

$$[Q_t]_{ij} = \begin{cases} 1 & \text{if } i = j = m \\ 1 - \beta_t & \text{if } i = j \neq m \\ \beta_t & \text{if } j = m, i \neq m \end{cases}$$

- Prop.)

- This is equivalent to the interpolation notation : $Q_t = (1 - \beta_t)\mathbf{I} + \beta_t \mathbf{1}e_m^\top$
 - where e_m is a vector with a 1 on the absorbing state and 0s elsewhere.
 - Use this for the **MASK** token!

- e.g.) When the absorbing token $m = 2$

$$\begin{aligned}
 & (1 - \beta_t) \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} + \beta_t \begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \underbrace{1}_{\text{abs}} & 0 & \dots & 0 \end{bmatrix} \\
 &= \begin{bmatrix} 1 - \beta_t & \beta_t & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & \beta_t & 1 - \beta_t & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \underbrace{\beta_t}_{\text{abs}} & 0 & \dots & 1 - \beta_t \end{bmatrix}
 \end{aligned}$$

D3PM continues

Forward Process

One Step Marginal

$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \text{Cat}(\mathbf{x}_t; \mathbf{p} = \mathbf{x}_{t-1} \mathbf{Q}_t)$$

- i.e.) Sample \mathbf{x}_t from the categorical distribution $\mathbf{x}_{t-1} \mathbf{Q}_t$

One Shot Marginal

$$q(\mathbf{x}_t \mid \mathbf{x}_0) = \text{Cat}(\mathbf{x}_t; \mathbf{p} = \mathbf{x}_0 \overline{\mathbf{Q}}_t) : \text{the } t\text{-step marginal}$$

- where $\overline{\mathbf{Q}}_t = \mathbf{Q}_1 \mathbf{Q}_2 \cdots \mathbf{Q}_t$

Posterior

$$q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \mathbf{x}_0) q(\mathbf{x}_{t-1} \mid \mathbf{x}_0)}{q(\mathbf{x}_t \mid \mathbf{x}_0)} = \text{Cat} \left(\mathbf{x}_{t-1}; \mathbf{p} = \frac{\mathbf{x}_t \mathbf{Q}_t^\top \odot \mathbf{x}_0 \overline{\mathbf{Q}}_{t-1}}{\mathbf{x}_0 \overline{\mathbf{Q}}_t \mathbf{x}_t^\top} \right)$$

D3PM continues

Reverse Process

$$p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) \propto \sum_{\tilde{\mathbf{x}}_0} \underbrace{q(\mathbf{x}_{t-1}, \mathbf{x}_t \mid \tilde{\mathbf{x}}_0)}_{\text{term 0}} \underbrace{\tilde{p}_{\theta}(\tilde{\mathbf{x}}_0 \mid \mathbf{x}_t)}_{\text{nn}_{\theta}(\mathbf{x}_t) \text{ predicts this!}}$$

$$\underbrace{q(\mathbf{x}_{t-1}, \mathbf{x}_t \mid \tilde{\mathbf{x}}_0)}_{\text{term 0}} = \underbrace{q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \tilde{\mathbf{x}}_0)}_{\text{term 1}} \cdot \underbrace{q(\mathbf{x}_{t-1} \mid \tilde{\mathbf{x}}_0)}_{\text{term 2}}$$

$$\underbrace{q(\mathbf{x}_t \mid \mathbf{x}_{t-1}, \tilde{\mathbf{x}}_0)}_{\text{term 1}} = q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$$

$$\underbrace{q(\mathbf{x}_{t-1} \mid \tilde{\mathbf{x}}_0)}_{\text{term 2}} = \tilde{\mathbf{x}}_0 \overline{Q}_{t-1} = \tilde{\mathbf{x}}_0 Q_1 \cdots Q_{t-1} \quad \leftarrow \text{Expensive } \overline{Q}$$

D3PM continues

Loss

$$L_\lambda = L_{\text{vb}} + \lambda \mathbb{E}_{q(\mathbf{x}_0)} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[-\log \tilde{p}_\theta(\mathbf{x}_0 | \mathbf{x}_t) \right]$$

$$L_{\text{vb}} = \mathbb{E}_{q(\mathbf{x}_0)} \left[\underbrace{D_{\text{KL}}[q(\mathbf{x}_T|\mathbf{x}_0)||p(\mathbf{x}_T)]}_{L_T} + \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}[q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)]]}_{L_{t-1}} \right. \\ \left. \underbrace{-\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{L_0} \right]. \quad (1)$$

forward posterior reverse denoiser

Drawback

- q and p_θ are N dimensional distributions, which are heavy.
- Parameterization is complicated (Consider all possible \tilde{x}_0 s)

MDLM : Single Token Problem

Idea

Keep the **absorbing state transition** to use the MASK token!

Instead of considering N dimensional distribution Q_t , focus only on **two** tokens!

1. original token \mathbf{x}
2. mask token \mathbf{m}

Inheriting the **interpolation notation** in the forward process

- $q(\mathbf{z}_t \mid \mathbf{x}) = \text{Cat}(\mathbf{z}_t ; \alpha_t \mathbf{x} + (1 - \alpha_t) \mathbf{m})$

i.e.) Sample the intermediate state \mathbf{z}_t from the distribution $\alpha_t \mathbf{x} + (1 - \alpha_t) \mathbf{m}$

MDLM : Single Token Problem continues

Forward Process

Marginal

$$q(\mathbf{z}_t \mid \mathbf{x}) = \text{Cat}(\mathbf{z}_t ; \alpha_t \mathbf{x} + (1 - \alpha_t) \mathbf{m})$$

- Then, \mathbf{z}_t is sampled from $\alpha_t \mathbf{x} + (1 - \alpha_t) \mathbf{m}$.
 - i.e.) $\mathbf{z}_t = \begin{cases} \mathbf{x} & \text{with } p = \alpha_t \\ [\text{MASK}] & \text{with } p = 1 - \alpha_t \end{cases}$

Posterior

$$q(\mathbf{z}_s \mid \mathbf{z}_t, \mathbf{x}) = \begin{cases} \text{Cat}(\mathbf{z}_s ; \mathbf{z}_t) & \mathbf{z}_t \neq \mathbf{m} \\ \text{Cat} \left(\mathbf{z}_s ; \frac{(1 - \alpha_s) \mathbf{m} + (\alpha_s - \alpha_t) \mathbf{x}}{1 - \alpha_t} \right) & \mathbf{z}_t = \mathbf{m} \end{cases}$$

- where $(s < t)$

MDLM : Single Token Problem continues

Reverse Process : SUBS Parameterization- Model)

$$p_{\theta}(\mathbf{z}_s \mid \mathbf{z}_t) = q(\mathbf{z}_s \mid \mathbf{z}_t, \mathbf{x} = \mathbf{x}_{\theta}(\mathbf{z}_t, t)) = \begin{cases} \text{Cat}(\mathbf{z}_s ; \mathbf{z}_t) \\ \text{Cat} \left(\mathbf{z}_s ; \frac{(1 - \alpha_s)\mathbf{m} + (\alpha_s - \alpha_t)\mathbf{x}_{\theta}(\mathbf{z}_t, t)}{1 - \alpha_t} \right) \end{cases}$$

- Desc.)
 - Similar to the forward posterior.
 - Only difference is that \mathbf{x} is parameterized as $\mathbf{x}_{\theta}(\mathbf{z}_t, t)$
 - Cases
 - If \mathbf{z}_t is masked (initial state),
 - we draw \mathbf{z}_s from the distribution $\frac{(1 - \alpha_s)\mathbf{m} + (\alpha_s - \alpha_t)\mathbf{x}_{\theta}(\mathbf{z}_t, t)}{1 - \alpha_t}$
 - Else (i.e. already unmasked),
 - $\mathbf{z}_s = \mathbf{z}_t$ definitively.

MDLM : Single Token Problem continues

Loss

Start with ELBO

$$L_{vb} = \mathbb{E}_{q(\mathbf{x}_0)} \left[\underbrace{D_{\text{KL}}[q(\mathbf{x}_T|\mathbf{x}_0)||p(\mathbf{x}_T)]}_{L_T} + \sum_{t=2}^T \underbrace{\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{\text{KL}}[q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)||p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)]]}_{L_{t-1}} \right. \\ \left. \underbrace{-\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_{\theta}(\mathbf{x}_0|\mathbf{x}_1)]}_{L_0} \right]. \quad (1)$$

forward posterior reverse denoiser

- Put $L_{T-1} = \mathcal{L}_{\text{diffusion}}$

MDLM : Single Token Problem continues

$$\begin{aligned}\mathcal{L}_{\text{diffusion}} &= \mathbb{E}_q \left(\sum_{i=1}^T D_{\text{KL}} \left[q(\mathbf{z}_{s(i)} \mid \mathbf{z}_{t(i)}, \mathbf{x}) \parallel p_{\theta}(\mathbf{z}_{s(i)} \mid \mathbf{z}_{t(i)}) \right] \right) \\ &= \sum_{i=1}^T \mathbb{E}_q \left(D_{\text{KL}} \left[\underbrace{q(\mathbf{z}_s \mid \mathbf{z}_t, \mathbf{x})}_{\text{Forward Posterior}} \parallel \underbrace{p_{\theta}(\mathbf{z}_s \mid \mathbf{z}_t, \mathbf{x}_{\theta}(\mathbf{z}_t, t))}_{\text{SUBS para'd Reverse}} \right] \right)\end{aligned}$$

where $s =$
 $\because \mathbf{z}_t \in \{\mathbf{x},$

⋮

Rao-Blackwellized Likelihood Bounds backs up Simplification

⋮

$$= \sum_{i=1}^T \mathbb{E}_q \left[\frac{\alpha_t - \alpha_s}{1 - \alpha_t} \underbrace{\log \langle \mathbf{x}_{\theta}(\mathbf{z}_t, t), \mathbf{x} \rangle}_{\text{Cross Entropy Loss!}} \right]$$

cf.) $\begin{cases} \mathbf{z}_t = \\ \mathbf{z}_t = \end{cases}$

MDLM : Single Token Problem continues

Continuous Time Loss

Making the number of steps $T \rightarrow \infty$, we may get

$$\mathcal{L}_{\text{NELBO}}^{\infty} = \mathbb{E}_q \int_{t=0}^{t=1} \frac{\alpha'_t}{1 - \alpha_t} \underbrace{\log \langle \mathbf{x}_{\theta}(\mathbf{z}_t, t), \mathbf{x} \rangle}_{\text{Cross Entropy Loss!}} dt$$

- where $\alpha'_t = \lim_{T \rightarrow \infty} T(\alpha_t - \alpha_s) \quad (\because T \rightarrow \infty \Rightarrow s \rightarrow t)$

MDLM dealing with length L sequence of tokens

Forward noising process is applied **independently** across a sequence.

- $q(\mathbf{z}_t^{1:L} \mid \mathbf{x}^{1:L}) = \prod_{\ell=1}^L q(\mathbf{z}_t^\ell \mid \mathbf{x}^\ell)$

Denoising process factorizes **independently** across tokens, conditioned on a sequence of latents $\mathbf{z}_t^{1:L}$

- $p_\theta(\mathbf{z}_s^{1:L} \mid \mathbf{z}_t^{1:L}) = \prod_{\ell=1}^L p_\theta(\mathbf{z}_s^\ell \mid \mathbf{z}_t^{1:L})$

Optimization)

- $\mathcal{L}_{\text{NELBO}}^\infty = \mathbb{E}_q \int_{t=0}^{t=1} \frac{\alpha'_t}{1 - \alpha_t} \sum_{\ell=1}^L \log \langle \mathbf{x}_\theta^\ell(\mathbf{z}_t^{1:L}, t), \mathbf{x}^\ell \rangle dt$

PGM also solves **Mask Generative Modeling** problem

$$\mathcal{L}_{\text{MGM}} := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, t \sim \mathcal{U}[0,1]} [w(t) \text{CE}(\mathbf{x}_{\theta}(\mathbf{z}_t, t), \mathbf{x})]$$

According to them **MDLM**'s **NELBO** loss is just a MGM loss with special **weight**

$$\mathcal{L}_{\text{NELBO}}^{\infty} = \mathbb{E}_q \int_{t=0}^{t=1} \frac{\alpha'_t}{1 - \alpha_t} \sum_{\ell=1}^L \underbrace{\log \langle \mathbf{x}_{\theta}(\mathbf{z}_t, t), \mathbf{x} \rangle}_{\text{Cross Entropy Loss!}} dt$$

- where $w(t) = \frac{\alpha'_t}{1 - \alpha_t}$

PGM's Idea : Based on the MDLM framework...

- Treat masked and unmasked tokens as **group 0** and **group 1**
- Let each group learn each other in the opposite time schedule

Objective

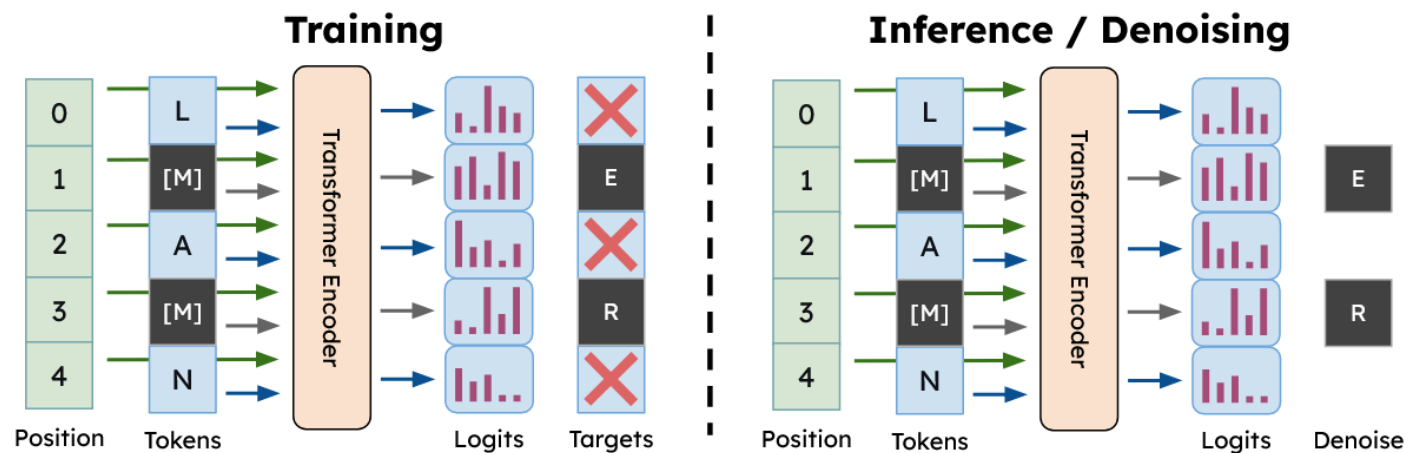
$$\mathcal{L}_{\text{PGM}} = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}, t \sim \mathcal{U}[0,1]} \left[w^{\text{PGM}}(\mathbf{g}, t) \text{CE}(\mathbf{x}_{\theta}(\mathbf{x}; \mathbf{g}; t), \mathbf{x}) \right]$$

- where $w^{\text{PGM}}(\mathbf{g}, t)_i = \begin{cases} w(t) & \text{if } \mathbf{g}_i = 0 \\ w(1 - t) & \text{if } \mathbf{g}_i = 1 \end{cases}$

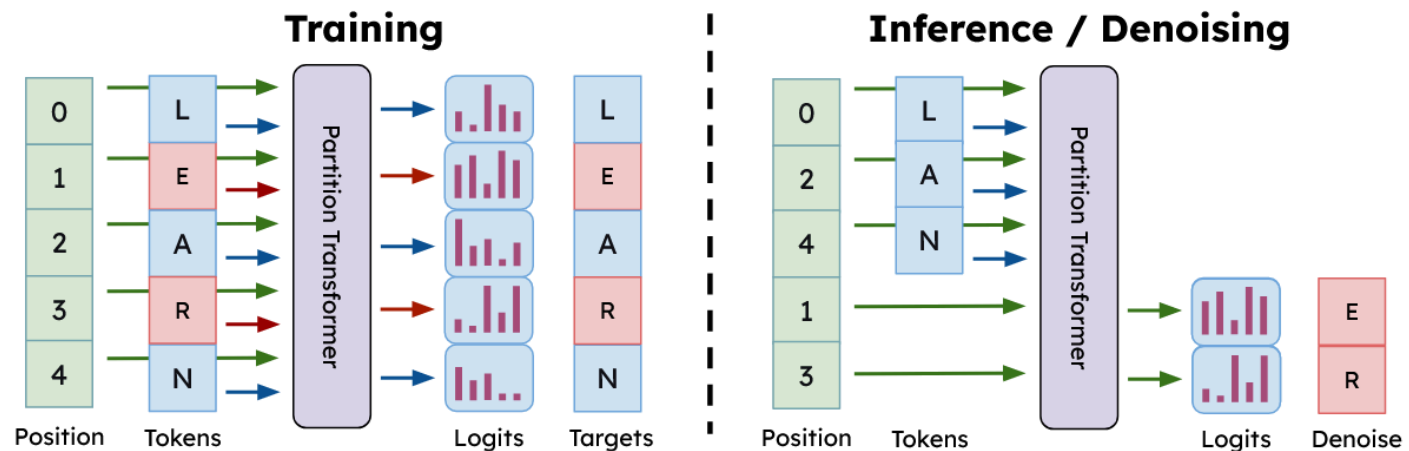
How?

- Enable this using the architecture

MDLM vs PGM

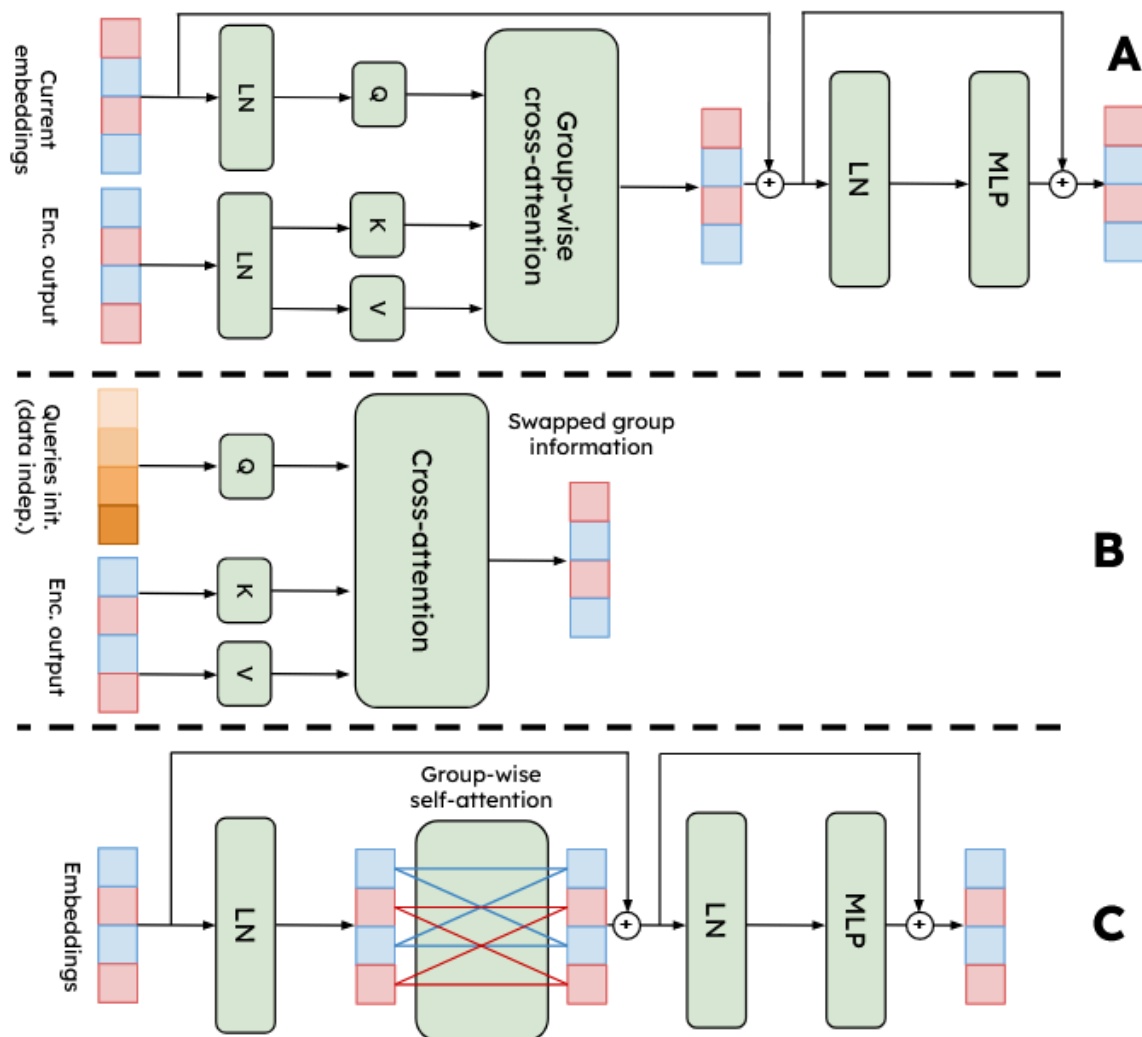
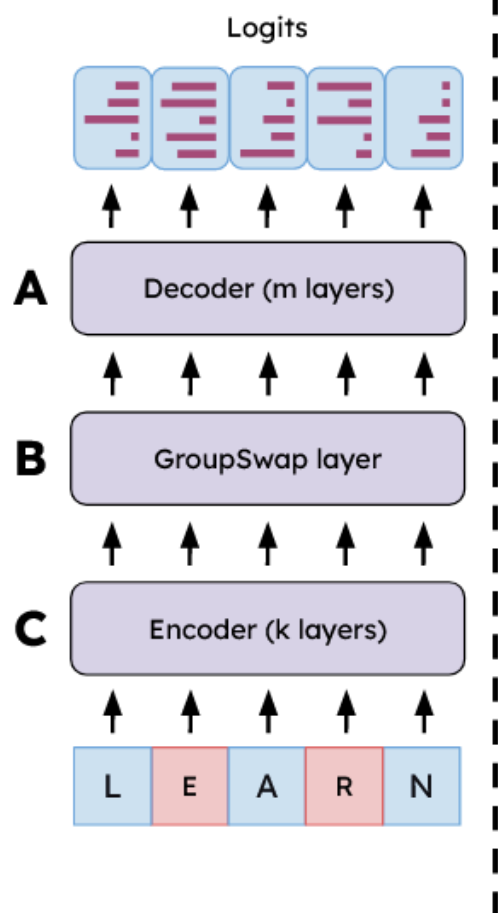


(a) Masked Generative Modeling (MGM)



(b) Partition Generative Modeling (PGM)

Partition Transformer



Encoder

- Desc.)
 - Partition-wise **self-attention** blocks
- Prop.)
 - Similar to standard bidirectional transformer blocks
 - Different on that separate groups do not attend each other

Decoder

- Desc.)
 - Cross-attention layer
 - Keys and Values are computed based on the output of the **encoder**
 - Queries are computed using either...
 - the output of the **GroupSwap layer** (Current Implementation)
 - the output of the previous decoder block (Future work?)
- Prop.)
 - No self-attention layer
 - Why?)
 - Compute predictions solely at the positions that we will decode.
 - Paper argues that this allows efficient generation.

GroupSwap Layer

- Goal)
 - Recall that encoder localized each group information.
 - This layer allows information exchange between groups.
 - Why?)
 - Prediction on group 1 should be based on group 0
 - and vice versa
- Implementation)
 - Cross Entropy
 - Queries
 - Data-Independent
 - Data-Dependent
 - Keys and Values
 - Inputted from the [encoder](#)
 - No specified details provided.

Experiments

Language Modeling : PGM vs MDLM

- PGM achieved....
 - lower validation **perplexity**
 - 5 times higher sampling **throughput**
 - faster inference without sacrificing **downstream performance**
 - i.e. tasks never trained
 - higher accuracy after **distillation**

Image Modeling : PGM vs MaskGIT (VQGAN on ImageNet)

- PGM (64 steps) achieved....
 - improved FID score
 - 3.9 times faster

Isolating the effect of Complementary Masking

- Setting)
 - Use only MDLM
 - Input complementary masked sequences
 - e.g.) "L[M]A[M]N", "[M]E[M]R[M]"
 - Compare it with the normal MDLM.
- Result

Model ↓	#Params	Val. PPL	Latency (sec) ↓	TP (tok/sec) ↑
<i>LM1B (ctx len. 128)</i>				
MDLM	170M	27.67	3.78	1'081.57
MDLM [†] (Compl. masking)	170M	25.72	3.78	1'081.57
PGM 6 / 6	171M	<u>26.80</u>	2.12	1'930.93
<i>OpenWebText (ctx len. 1024)</i>				
MDLM	170M	23.07	31.41	1'043.22
MDLM [†] (Compl. masking)	170M	22.98	31.41	1'043.22
PGM 8 / 8	203M	<u>22.61</u>	5.86	5'585.57
PGM 6 / 6 (dim. 1024)	268M	21.43	5.93	5'518.09