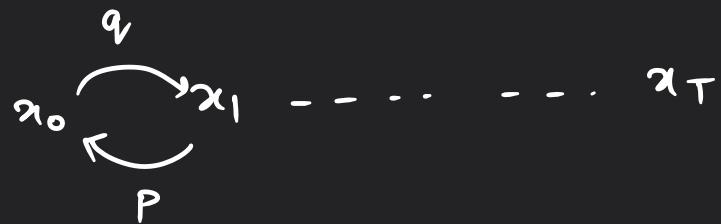


# Diffusion Models and their use in Language Modelling

cumulating loss problem  
fast inference  
control over generation

↳ Main observation points

## # DDPM:



$$q(\gamma_t | \gamma_{t-1}) = N(\gamma_t | \sqrt{1-\beta_t} \gamma_{t-1}, \beta_t I)$$

$$\gamma_t = \sqrt{1-\beta_t} \gamma_{t-1} + \sqrt{\beta_t} \varepsilon \quad \varepsilon \sim N(0, I)$$

$$q(\gamma_t | \gamma_0) = N(\gamma_t | \sqrt{\alpha_t} \gamma_0, (1-\alpha_t) I)$$

$$\gamma_t = \sqrt{\alpha_t} \gamma_0 + \sqrt{1-\alpha_t} \varepsilon \quad \varepsilon \sim N(0, I)$$

objective:  $\max_{\theta} E_{\gamma_0 \sim p^*} [\log p_{\theta}(\gamma_0)]$

$$p_{\theta}(\gamma_0) = p_{\theta}(\gamma_0, z) \quad z = \gamma_1: T$$

$$\begin{aligned} \log p_{\theta}(\gamma_0) &= \log \int p_{\theta}(\gamma, z) dz \\ &= \log \int \frac{q(z|\gamma)}{q(z|\gamma)} p_{\theta}(\gamma, z) dz \end{aligned}$$

$$= \log E_q \left( \frac{p_{\theta}(\gamma, z)}{q(z|\gamma)} \right)$$

$$\geq E_q \left( \log \frac{p_{\theta}(\gamma, z)}{q(z|\gamma)} \right) \rightarrow ELBO$$

$$p_{\theta}(\gamma_0, z) = p_{\theta}(\gamma_0, \gamma_1, \dots, \gamma_T) = p(\gamma_T) \prod_{t=1}^T p(\gamma_{t-1} | \gamma_t)$$

$$q(z|\gamma) = q(\gamma_1, \dots, \gamma_T | \gamma_0) = \prod_{t=1}^T q(\gamma_t | \gamma_{t-1})$$

$$Eq \left( \log \frac{\prod_{t=1}^T P(x_{t-1}|x_t)}{\prod_{t=1}^T q(x_t|x_{t-1})} \right)$$

$$= Eq \left( \sum_{t=1}^T \log \frac{P(x_{t-1}|x_t)}{q(x_t|x_{t-1})} \right)$$

$$= Eq \left( \log \frac{P(x_0|x_1)}{q(x_1|x_0)} + \sum_{t=2}^{\infty} \log \frac{P(x_{t-1}|x_t)}{q(x_t|x_{t-1})} \right)$$

$$q(x_t|x_{t-1}) = q(x_t|x_{t-1}, x_0)$$

$$= \frac{q(x_{t-1}|x_t, x_0) q(x_t|x_0)}{q(x_{t-1}|x_0)} \quad \begin{matrix} q(x_t|x_0) \\ \curvearrowleft \\ q(x_T|x_0) \end{matrix}$$

$$= Eq \left( \sum_{t=2}^{\infty} \log \frac{P(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} + \log \prod_{t=2}^T \frac{q(x_{t-1}|x_0)}{q(x_t|x_0)} \right. \\ \left. + \log \frac{P(x_0|x_1)}{q(x_1|x_0)} \right)$$

$$= Eq \left( \sum_{t=2}^T \log \frac{P(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} + \log \frac{P(x_0|x_1)}{q(x_T|x_0)} \right)$$

$$= \sum_{t=2}^T E_{x_t \sim q_t, x_{t-1} \sim q_{t-1}} \left( \log \frac{P(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)} \right)$$

$$+ E_{x_1 \sim q_1} \log P(x_0|x_1)$$

$$= \sum_{t=2}^T E_{\alpha_t \sim q(\alpha_t | \alpha_0)} \left( -D_{KL}(q(\alpha_{t-1} | \alpha_t, \alpha_0) || P(\alpha_{t-1} | \alpha_t)) \right. \\ \left. + E_{\alpha_1 \sim q_1} \log P(\alpha_0 | \alpha_1) \right)$$

$$\text{Obj} \min_{\theta} L_{\theta}$$

$$L(\theta) = \sum_{t=2}^T E_{q(\alpha_t | \alpha_0)} \left( D_{KL}(q(\alpha_{t-1} | \alpha_t, \alpha_0) || P(\alpha_{t-1} | \alpha_t)) \right. \\ \left. - E_{q_1(\alpha_1 | \alpha_0)} \log P(\alpha_0 | \alpha_1) \right)$$

$$q(\alpha_{t-1} | \alpha_t, \alpha_0) = \frac{q(\alpha_t | \alpha_{t-1}, \alpha_0) q(\alpha_{t-1} | \alpha_0)}{q(\alpha_t | \alpha_0)}$$

$$= \frac{N(\sqrt{\alpha_t} \alpha_{t-1}, (1-\alpha_t) I) N(\sqrt{\alpha_{t-1}} \alpha_0, (1-\alpha_{t-1}) I)}{N(\sqrt{\alpha_t} \alpha_0, (1-\alpha_t) I)}$$

$$\propto \exp \left( - \frac{(\alpha_t - \sqrt{\alpha_t} \alpha_{t-1})^2}{2(1-\alpha_t)} - \frac{(\alpha_{t-1} - \sqrt{\alpha_{t-1}} \alpha_0)^2}{2(1-\alpha_{t-1})} \right.$$

$$\left. + \frac{(\alpha_t - \sqrt{\alpha_t} \alpha_0)^2}{2(1-\alpha_t)} \right)$$

$$\propto \exp \left( \frac{1}{2(1-\alpha_t)} (-\alpha_t \alpha_{t-1}^2 + 2\sqrt{\alpha_t} \alpha_t \alpha_{t-1}) \right. \\ \left. - \frac{1}{2(1-\alpha_{t-1})} (\alpha_{t-1}^2 - 2\sqrt{\alpha_{t-1}} \alpha_0 \alpha_{t-1}) \right)$$

$$\propto \exp \left( -\frac{1}{2} \left[ \left( \frac{\alpha_t}{(1-\alpha_t)} + \frac{1}{(1-\bar{\alpha}_{t-1})} \right) \alpha_{t-1}^2 - 2 \left( \frac{\sqrt{\alpha_t} \alpha_t}{1-\alpha_t} + \frac{\sqrt{\alpha_{t-1}} \alpha_0}{1-\bar{\alpha}_{t-1}} \right) \alpha_{t-1} \right] \right)$$

$$\propto \exp \left( -\frac{1}{2} \left[ \frac{\alpha_t - \bar{\alpha}_t + 1 - \bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \alpha_{t-1}^2 - 2 \left( \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\alpha_t + \sqrt{\alpha_{t-1}}(1-\alpha_t)\alpha_0}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \alpha_{t-1} \right] \right)$$

$$\propto \exp \left( -\frac{1}{2} \left( \frac{1-\bar{\alpha}_t}{(1-\alpha_t)(1-\bar{\alpha}_{t-1})} \right) \alpha_{t-1}^2 - \frac{2\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{(1-\bar{\alpha}_t)} \alpha_{t-1} \right)$$

$$\Sigma = \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} I =$$

$$\tilde{M} = \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})\alpha_t + \sqrt{\alpha_{t-1}}(1-\alpha_t)\alpha_0}{1-\bar{\alpha}_t}$$

$$\tilde{N} = \frac{\sqrt{\alpha_{t-1}}(1-\alpha_t)}{1-\bar{\alpha}_t} \alpha_0 + \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} \alpha_t$$

$$\alpha_t = \sqrt{\alpha_t} \alpha_0 + \sqrt{1-\bar{\alpha}_t} \varepsilon \quad \alpha_0 = \frac{\alpha_t - \sqrt{1-\bar{\alpha}_t} \varepsilon}{\sqrt{\alpha_t}}$$

$$\begin{aligned}
\tilde{\mu}^2 &= \frac{\left(\alpha_t - \sqrt{1-\bar{\alpha}_t} \varepsilon\right)}{\sqrt{\bar{\alpha}_t}} \frac{\sqrt{\alpha_{t-1}(1-\alpha_t)}}{1-\bar{\alpha}_t} + \frac{\sqrt{\alpha_t(1-\bar{\alpha}_{t-1})}}{1-\bar{\alpha}_t} \alpha_t \\
&= \frac{1}{\sqrt{\alpha_t}} \left( -\frac{(1-\alpha_t) \varepsilon}{\sqrt{1-\bar{\alpha}_t}} \right) + \frac{1}{\sqrt{\alpha_t}} \left( \frac{\alpha_t(1-\alpha_t)}{1-\bar{\alpha}_t} + \frac{\alpha_t(1-\bar{\alpha}_{t-1})}{(1-\bar{\alpha}_t)} \alpha_t \right) \\
&\quad \frac{\alpha_t - \cancel{\alpha_t} \alpha_t + \cancel{\alpha_t} \alpha_t}{-\bar{\alpha}_t \alpha_t} = \frac{\alpha_t}{(1-\bar{\alpha}_t)} = \tilde{\alpha}_t
\end{aligned}$$

$$\tilde{\mu} = \frac{1}{\sqrt{\alpha_t}} \left( \alpha_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}} \varepsilon \right) \quad \tilde{\Sigma} = \frac{\beta_t(1-\bar{\alpha}_{t-1})}{(1-\bar{\alpha}_t)} I$$

$$\text{obj} = E_{\alpha_0 \sim P^2, t \sim U[0, T], \alpha_t \sim q_t} [ D_{KL}(q_t || p_t) ]$$

$$q_t \sim N(\tilde{\mu}, \tilde{\Sigma}) \quad p_t \sim N(\mu_\theta, \Sigma_\theta)$$

decoder model

# #MDLM:

obj:

$$\begin{aligned} & \mathbb{E}_{\alpha_0 \sim P_d} \left[ \sum_{t=2}^{\infty} \mathbb{E}_{q(\alpha_t | \alpha_0)} \left[ D_{KL} \left( q(\alpha_{t-1} | \alpha_{t-1}, \alpha_0) \| P(\alpha_{t-1} | \alpha_t) \right) \right. \right. \\ & \quad \left. \left. + \mathbb{E}_{q(\alpha_t | \alpha_0)} [\log P(\alpha_0 | \alpha_t)] \right] \right] \\ & = \mathbb{E}_{t \in \{1, 2, \dots\}} \mathbb{E}_{\alpha_0 \sim P_d} \mathbb{E}_{\alpha_t \sim q(\alpha_t | \alpha_0)} [D_{KL}(q || P)] \end{aligned}$$

instead of Gaussian, Categorical

$$q(\alpha_t | \alpha_{t-1}) = \text{cat}(\Theta_t^\top \alpha_{t-1})$$

$$\begin{aligned} q(\alpha_t | \alpha_0) &= \text{cat}((\Theta_1, \dots, \Theta_t)^\top \alpha_0) \\ &= \text{cat}(\bar{\Theta}_t^\top \alpha_0) \end{aligned}$$

$$q(\alpha_t | \alpha_0) = \text{cat}(\alpha_t \alpha_0 + (1-\alpha_t)m) \quad m \rightarrow \text{mask}$$

$$q(\alpha_t | \alpha_{t-1}) = \text{cat}(\alpha_t \alpha_{t-1} + (1-\alpha_t)m)$$

$$s < t \quad q(\alpha_t | \alpha_s) = \text{cat}(\alpha_t \alpha_s + (1-\alpha_t)m)$$

$$\Theta_t|s = \alpha_t|s I + (1-\alpha_t|s) 1m^\top$$

$$q(\alpha_t | \alpha_s) = \text{cat}(\Theta_t|s^\top \alpha_s)$$

$$q(\alpha_s | \alpha_t, \alpha_0) = \frac{q(\alpha_t | \alpha_s) q(\alpha_s | \alpha_0)}{q(\alpha_t | \alpha_0)}$$

$$= \frac{\text{cat}(\Theta_{t|S} \alpha_t) \text{cat}(\Theta_S^\top \alpha_0)}{\text{cat}(\Theta_t^\top \alpha_0)}$$

$$= \text{cat}\left(\frac{\Theta_{t|S} \alpha_t \odot \Theta_S^\top \alpha_0}{\alpha_t^\top \Theta_t^\top \alpha_0}\right)$$

Model: given  $\alpha_t$  it predicts  $\alpha_0$  which is put in forward process-

$$\alpha_0 \rightarrow f_\theta(\alpha_t, t)$$

$$\underline{\text{obj}}: DKL\left(q(\alpha_s | \alpha_t, \alpha_0) \parallel P(\alpha_s | \alpha_t)\right) \quad s < t$$

$$P(\alpha_s | \alpha_t) = q(\alpha_s | \alpha_t, f_\theta(\alpha_t, t))$$

$$P(\alpha_s | \alpha_t = m) = \text{cat}(\alpha_t)$$

$$P(\alpha_s | \alpha_t = m) = \text{cat}\left(\frac{\alpha_t m \odot \alpha_0 + (\alpha_s - \alpha_t) \alpha_0 + (1 - \alpha_s) m}{\alpha_t \prec \alpha_0, m \succ + 1 - \alpha_t}\right)$$

$$L(\theta) = \sum_{\alpha_s} q(\alpha_s | \alpha_t, \alpha_0) \log \frac{q(\alpha_s | \alpha_t, \alpha_0)}{P(\alpha_s | \alpha_t)}$$

$$\text{if } \alpha_t = x \quad L(\theta) = 0 \quad \because P(\alpha_s | \alpha_t) = q(\alpha_s | \alpha_t, \alpha_0) = \text{cat}(\alpha_t)$$

if  $\alpha_t = m$

$$L_2(\theta) = \sum_{\alpha_s} q(\alpha_s | \alpha_t, \alpha_0) \log \frac{q(\alpha_s | \alpha_t, \alpha_0)}{P(\alpha_s | \alpha_t)}$$

$$= \frac{\alpha_s - \alpha_t}{1 - \alpha_t} \log \frac{\alpha_t < \alpha_0, m > + (1 - \alpha_t)}{(1 - \alpha_t) < \alpha_0, \alpha >}$$

$$+ \frac{1 - \alpha_s}{1 - \alpha_t} \log \frac{(1 - \alpha_s)(\alpha_t < \alpha_0, m >) + (1 - \alpha_t)}{(1 - \alpha_t)(\alpha_s < \alpha_0, m > + (1 - \alpha_s))}$$

$$L(\theta) = E_t E_q(\alpha_t | \alpha) \left( L_1(\theta) \xrightarrow{0} \alpha_t < \alpha_t, \alpha > + L_2(\theta) < \alpha_t, m > \right)$$
$$= E_{t \sim \{1, \dots, T\}} E_q(\alpha_t | \alpha) [ \nabla L_2(\theta) < \alpha_t, m > ]$$

SUBS parameterization

① zero-masking prob:

observing  $< \alpha, m > = 0$  we put  $< \alpha_0, m > = 0$

② carry-over unmasking:

$\pi_0(\alpha_t, t) = \alpha_t$  if  $\alpha_t \neq m$ .

(i.e predict only for the masked tokens

subs. in  $L(\theta)$  to get

$$L(\theta) = E_{t \sim \{1, \dots, T\}} \sum_{\alpha_t \sim q_t} \left[ \nabla \left[ \frac{\alpha_s - \alpha_t}{1 - \alpha_t} \log < \alpha_0, \alpha > \right] \right]$$



Rao-Blackwellized bound

C.E over  
masked  $\alpha_t$ s

## #Flex-MDLM:

unmasking posterior  $\rightarrow P(\alpha_s | \alpha_t = m)$

Insertion expectation  $\rightarrow$  exp. number of tokens between  $\alpha_{i-1}$  &  $\alpha_i$  + i

$$L_\theta = L_{MDLM} + L_{\text{insertion}}$$

$\downarrow$   
unmasking loss

use Bregman divergence for ins. loss

$$L_{MDLM} = E_t E_{q_p(\alpha_t | \alpha)} \left[ \frac{\hat{\beta}_t}{1 - \hat{\beta}_t} \sum_{\substack{i \\ \in \{\alpha_t^i = m\}}} \log (f_\theta(\alpha_{t,i}), \alpha_0^i) \right]$$

$$L_{\text{ins}} = E_t E_{q_p(\alpha_t | \alpha)} \left[ \frac{\hat{\alpha}_t}{1 - \hat{\alpha}_t} \sum_{i=0}^{\text{len}(\alpha_t)} \phi(s_t(i) - s_t(i-1)^{-1}, g_\theta^i(\alpha_{t,i})) \right]$$

