Large Language Diffusion Models

Shen Nie, Fengqi Zhu, Zebin You, Xiaolu Zhang, Jingyang Ou, Jun Hu, Jun Zhou, Yankai Lin, Ji-Rong Wen, Chongxuan Li

Presenter: Ye YUAN



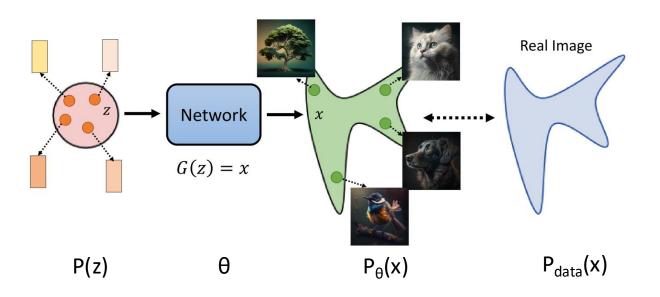




Contents

- Preliminary Knowledge
- Two Paradigms of Diffusion Models for Language Modeling
- Large Language Diffusion Models
- Discussion and Conclusions
- Questions and Answers

Preliminary – Likelihood-based Generative Models



- P_{data}(x): probability distribution of the data
- $P_{\theta}(x)$: approximate probability distribution of the data
- P(z): probability distribution of the latent variable, usually a Gaussian distribution

$$\begin{split} &\operatorname{Sample}\,\{x^1,x^2,...,x^m\}\,\operatorname{from}\,P_{data}(x) \\ &\theta^* = arg\,\max_{\theta}\prod_{i=1}^{m}P_{\theta}\big(x^i\big) = arg\,\max_{\theta}\log\prod_{i=1}^{m}P_{\theta}\big(x^i\big) \\ &= arg\,\max_{\theta}\sum_{i=1}^{m}\log P_{\theta}\big(x^i\big) \approx arg\,\max_{\theta}E_{x\sim P_{data}}[\log P_{\theta}(x)] \\ &= arg\,\max_{\theta}\int_{x}P_{data}(x)logP_{\theta}(x)dx - \int_{x}P_{data}(x)logP_{data}(x)dx \\ &= arg\,\max_{\theta}\int_{x}P_{data}(x)log\frac{P_{\theta}(x)}{P_{data}(x)}dx = \underset{\theta}{\operatorname{Difference between}\,P_{data}\,\operatorname{and}\,P_{\theta}} \end{split}$$

Maximum Likelihood = Minimize KL Divergence

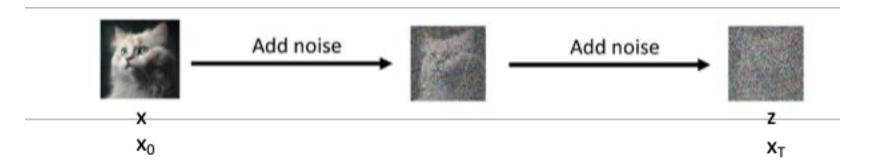
Preliminary – Evidence Lower Bound (ELBO)

- In practice, we usually maximize the log likelihood.
- Sometimes, the log likelihood is also called evidence.
- We can maximize the ELBO instead.

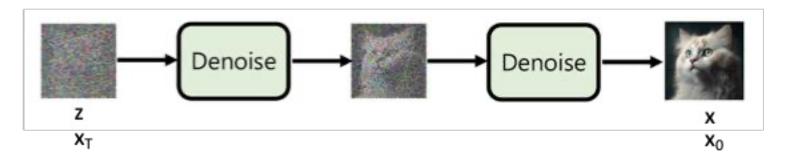
$$\begin{split} \log p(x) &= \log p(x) \int q_{\phi}(z|x) dz \\ &= \int q_{\phi}(z|x) (\log p(x)) dz \\ &= \underset{q_{\phi}(z|x)}{\mathbb{E}} \left[\log p(x) \right] \\ &= \underset{q_{\phi}(z|x)}{\mathbb{E}} \left[\log \frac{p(x,z)}{p(z|x)} \right] \\ &= \underset{q_{\phi}(z|x)}{\mathbb{E}} \left[\log \frac{p(x,z)}{p(z|x)} \right] \\ &= \underset{q_{\phi}(z|x)}{\mathbb{E}} \left[\log \frac{p(x,z)q_{\phi}(z|x)}{p(z|x)q_{\phi}(z|x)} \right] \\ &= \underset{q_{\phi}(z|x)}{\mathbb{E}} \left[\log \frac{p(x,z)q_{\phi}(z|x)}{p(z|x)} \right] + \underset{q_{\phi}(z|x)}{\mathbb{E}} \left[\log \frac{q_{\phi}(z|x)}{p(z|x)} \right] \\ &= \underset{q_{\phi}(z|x)}{\mathbb{E}} \left[\log \frac{p(x,z)}{q_{\phi}(z|x)} \right] + D_{KL}(q_{\phi}(z|x)||p(z|x)) \\ &\geq \underset{q_{\phi}(z|x)}{\mathbb{E}} \left[\log \frac{p(x,z)}{q_{\phi}(z|x)} \right]. \end{split}$$

Preliminary – Denoising Diffusion Models

Diffusion Process (Forward Process)



Denoising Process (Backward Process)



Preliminary – ELBO of Diffusion Models

 $\mathbb{E}_{q(x_1|x_0)}[\log p_{\theta}(x_0|x_1)]$: a reconstruction term, predicting the log probability of the original data sample given the first-step latent.

 $D_{KL}(q(x_T|x_0)||p(x_T))$: represents how close the distribution of the final noisified input is to the standard Gaussian prior.

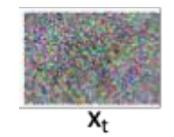
 $\sum_{t=2}^{T} \mathbb{E}_{q(\mathbf{x}_{t}|\mathbf{x}_{0})} \left[D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_{t},\mathbf{x}_{0})||p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_{t})) \right] : a$ denoising matching term. The q(x_{t-1}|x_t, x₀) defines how to denoise a noisy image x_t with access to what the final, completely denoised image x₀ should be.

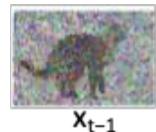
$$\log p(\mathbf{x}) = \log \int p(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T}$$

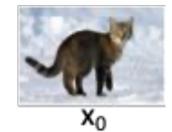
$$\geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_{0})} \left[\log \frac{p(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_{0})} \right]$$

$$= \mathbb{E}_{q(\mathbf{x}_{1}|\mathbf{x}_{0})} \left[\log p_{\theta}(\mathbf{x}_{0}|\mathbf{x}_{1}) \right] - D_{KL}(q(\mathbf{x}_{T}|\mathbf{x}_{0})||p(\mathbf{x}_{T}))$$

$$- \sum_{t=2}^{T} \mathbb{E}_{q(\mathbf{x}_{t}|\mathbf{x}_{0})} \left[D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_{t},\mathbf{x}_{0})||p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_{t})) \right].$$



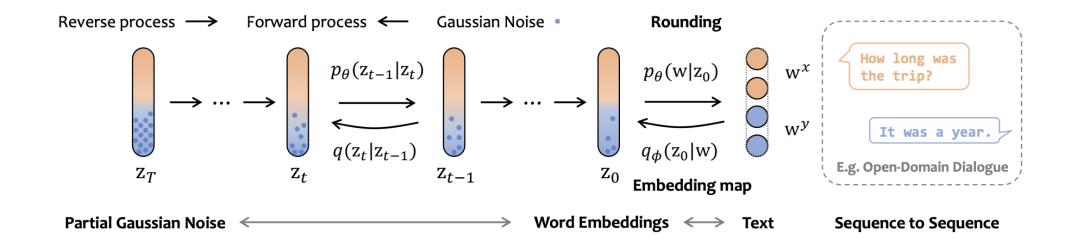




^[1] Luo, C. Understanding diffusion models: A unified perspective, 2022.

Diffusion for Language Modeling (Continuous)

Convert to Continuous State Space:



Diffusion for Language Modeling (Discrete)

Reformulate Two Key Equations in Discrete State Space:

Suppose we have a vocabulary table of size K=4: "apple", "banana", "carrot", "[MASK]".

Define a transition matrix $Q_t \in \mathbb{R}^{K \times K}$ at time t that govern the corruption process:

$$Q_t = egin{bmatrix} 0.1 & 0.1 & 0.1 & 0.7 \ 0.1 & 0.1 & 0.1 & 0.7 \ 0.1 & 0.1 & 0.1 & 0.7 \ 0.0 & 0.0 & 0.0 & 1.0 \end{bmatrix}$$

This means any normal token (e.g., "banana") has a 70% chance of turning into [MASK] at time t. Once a token becomes [MASK], it stays that way.

Diffusion for Language Modeling (Discrete)

"apple", "banana", "carrot", "[MASK]".

Forward Process:

$$q(x_t \mid x_{t-1}) = \operatorname{Cat}(x_t; \, p = x_{t-1}Q_t)$$

$$Q_t = egin{bmatrix} 0.1 & 0.1 & 0.1 & 0.7 \ 0.1 & 0.1 & 0.1 & 0.7 \ 0.1 & 0.1 & 0.1 & 0.7 \ 0.0 & 0.0 & 0.0 & 1.0 \end{bmatrix}$$

 x_{t-1} is a one-hot row vector. For "banana", $x_{t-1} = [0, 1, 0, 0]$.

 $x_{t-1}Q_t = [0.1, 0.1, 0.1, 0.7]$: give the probability of transitioning to other tokens.

•
$$q(x_t \mid x_0) = \text{Cat}(x_t; p = x_0 \overline{Q}_t)$$
, with $\overline{Q}_t = Q_1 Q_2 \dots Q_t$

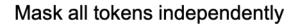
how likely each candidate x_{t-1} would result in the current x_t

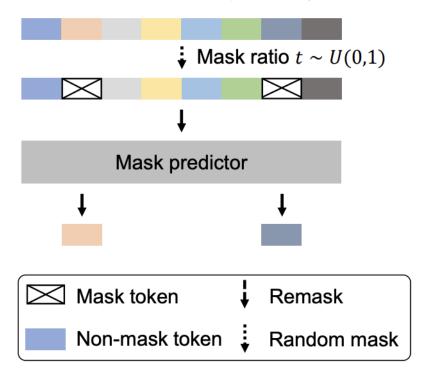
from the original input x_0

how likely each x_{t-1} was generated from the original input x_0

Large Language Diffusion Models - Pretraining

• This paper introduces large language diffusion model (LLaDA), with discrete states formulation. Each token can only be corrupted to [MASK] token in this paper.





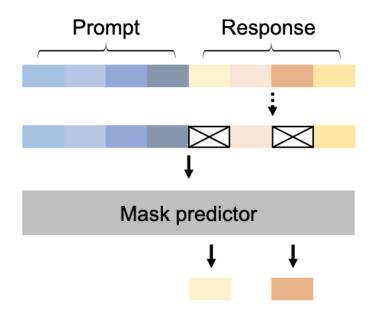
$$x_0 \sim p_{\mathrm{data}}, t \sim \mathrm{U}(0, 1]$$

 $x_t \sim q_{t|0}(x_t|x_0)$

This is like the BERT pretraining phase but with random mask ratio between 0 to 1.

Large Language Diffusion Models – SFT

 To enhance the instruction following capability of LLaDA, the authors perform supervised fine-tuning, similar to other autoregressive LLMs like Llama.



$$p_0, r_0 \sim p_{\text{data}}, t \sim \text{U}(0, 1]$$

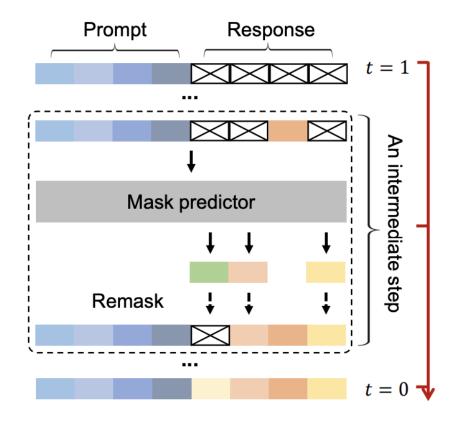
 $r_t \sim q_{t|0}(r_t|r_0)$

 p_0 is a prompt, r_0 is a response. Only apply masking to response.



Large Language Diffusion Models - Sampling

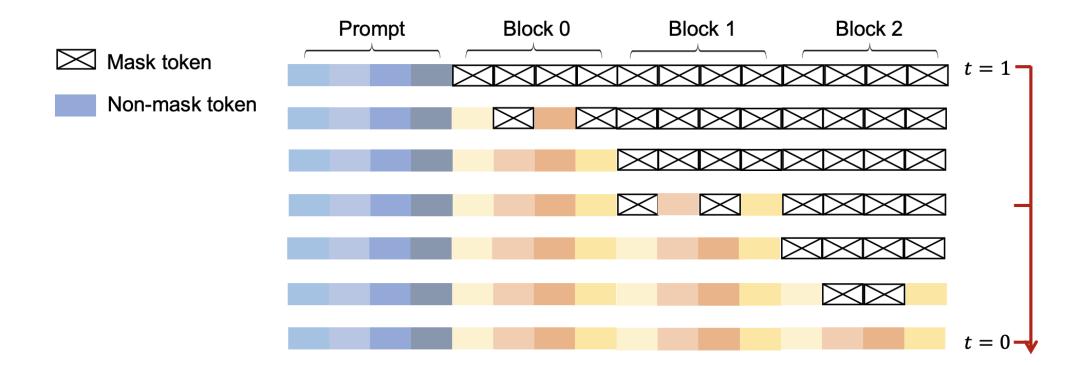
• LLaDA predicts all [MASK] tokens at once, then remasks some tokens and predicts them again. By doing this, the model can iteratively generate a response.



- Initialize the response as a sequence of [MASK] tokens with length s. The total steps for generation is T. The time steps could be [0, 1/T, ..., 1], and we remask $s \times i/T$ tokens at discrete time step i.
- Three different remask strategies:
 - Random: randomly remask tokens.
 - Low-Confidence: remask lowest confidence tokens.
 - **Semi-Autoregressive**: divide the sequence into several blocks and generate them from left to right.



Semi-Autoregressive Remasking



Discussion and Conclusion

- Advantages of LLaDA:
 - Faster sampling than autoregressive model.
 - Better reversal reasoning ability.

Table 3. Comparison in the Poem Completion Task.

	Forward	Reversal
GPT-4o (2024-08-06)	82.7	34.3
Qwen2.5 7B Instruct	75.9	38.0
LLaDA 8B Instruct	48.8	42.4

- Disadvantages of LLaDA:
 - Harder to calculate the probability of a generated sequence. (Need Monte-Carlo Estimations)
 - Generated sequence length and total generation steps are hyperparameters.
 - Still lower performance on general tasks than autoregressive LLMs.



Thanks for your attention!



References

- [1] Luo, C. Understanding diffusion models: A unified perspective, 2022.
- [2] Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. In Advances in Neural Information Processing Systems (2020).
- [3] S. Gong, M. Li, J. Feng, Z. Wu, and L. Kong, DiffuSeq: Sequence to Sequence Text Generation with Diffusion Models. International Conference on Learning Representations (2023).
- [4] J. Austin, D. D. Johnson, J. Ho, D. Tarlow, and van, Structured Denoising Diffusion Models in Discrete State-Spaces. In Advances in Neural Information Processing Systems (2021).
- [5] S. Nie et al., "Large Language Diffusion Models," 2025.