

Di[M]O: Distilling Masked Diffusion Models into One-step Generator

Yuanzhi Zhu, Xi Wang, Stéphane Lathuilière, Vicky Kalogeiton

LIX, École Polytechnique, CNRS, IPP, Inria, Univ. Grenoble Alpes, CNRS, LJK



Content

- Preliminaries
- Methods
- Results

Text

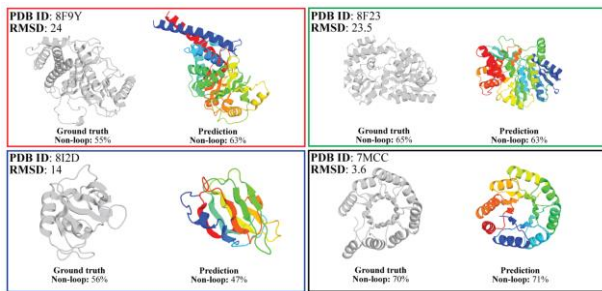
Multi-round Dialogue Case

User	Could you please share with me what the initial two lines of the renowned poem "The Road Not Taken"?
LLaDA	Certainly! The initial two lines of the renowned poem "The Road Not Taken" by Robert Frost are: "Two roads diverged in a yellow wood, And sorry I could not travel both." These lines set the stage for the poem, introducing the idea of a speaker facing a choice between two paths.
User	Please help me translate into Chinese.
LLaDA	Sure! The first two lines of "The Road Not Taken" by Robert Frost can be translated into Chinese as: "两条路分岔在黄色的树林中, 遗憾我不能同时走"
User	Please translate into German.
LLaDA	Sure! The first two lines of "The Road Not Taken" by Robert Frost can be translated into German as: "Zwei Wege trennten sich im gelben Wald, und leider konnte ich nicht beide tun."
User	Please also write a poem about life choices. I require 5 sentences, and the first word of each sentence starts with C as the first letter.
LLaDA	Certainly! Here's a poem about life choices, with each sentence starting with the letter "C": Choices cascade through life's journey, Creating connections with each decision. Careful consideration leads to better paths, Courageous choices shape our destiny, Choosing wisely guides us through the years.

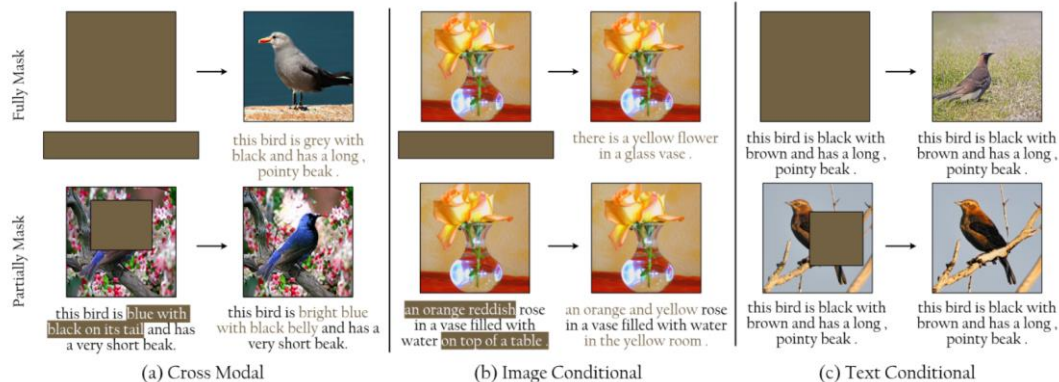
Image



Protein



Multi-modality



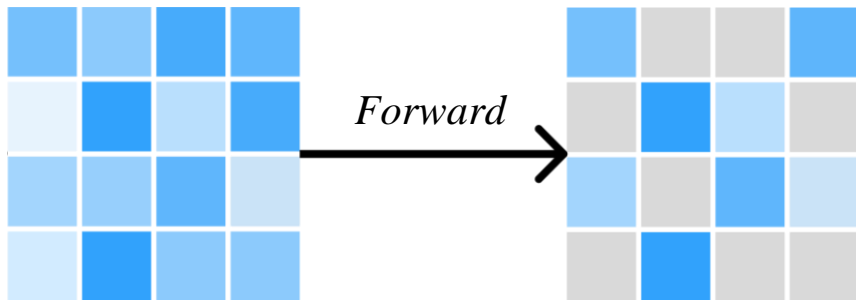
(a) Cross Modal

(b) Image Conditional

(c) Text Conditional

Mask Diffusion Models

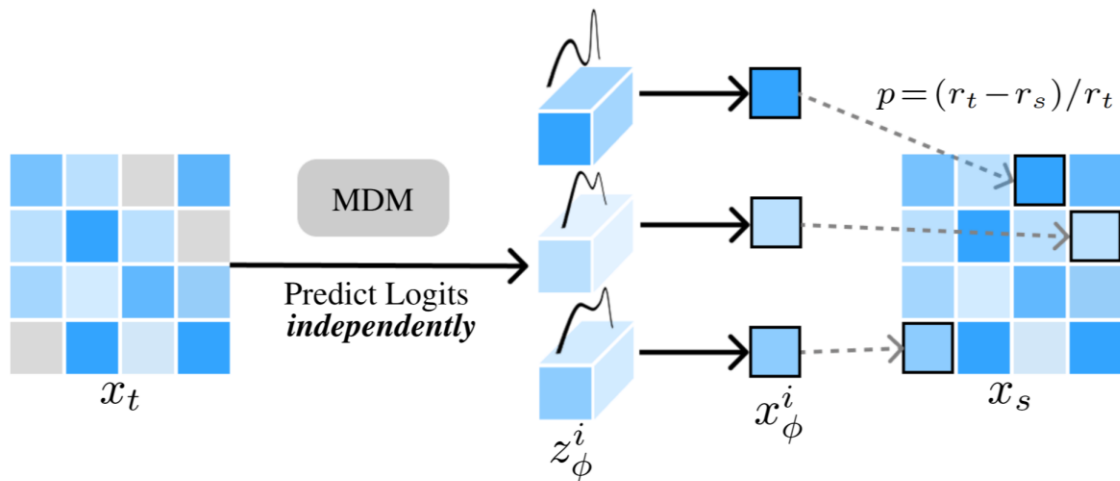
Forward Process: mask *independently* each token at position i



$$q_{t|0}(x_t|x_0) = \prod_{i=0}^{L-1} \text{Cat}(x_t^i; (1 - r_t)\delta(x_0^i) + r_t\delta([\mathbf{M}])))$$

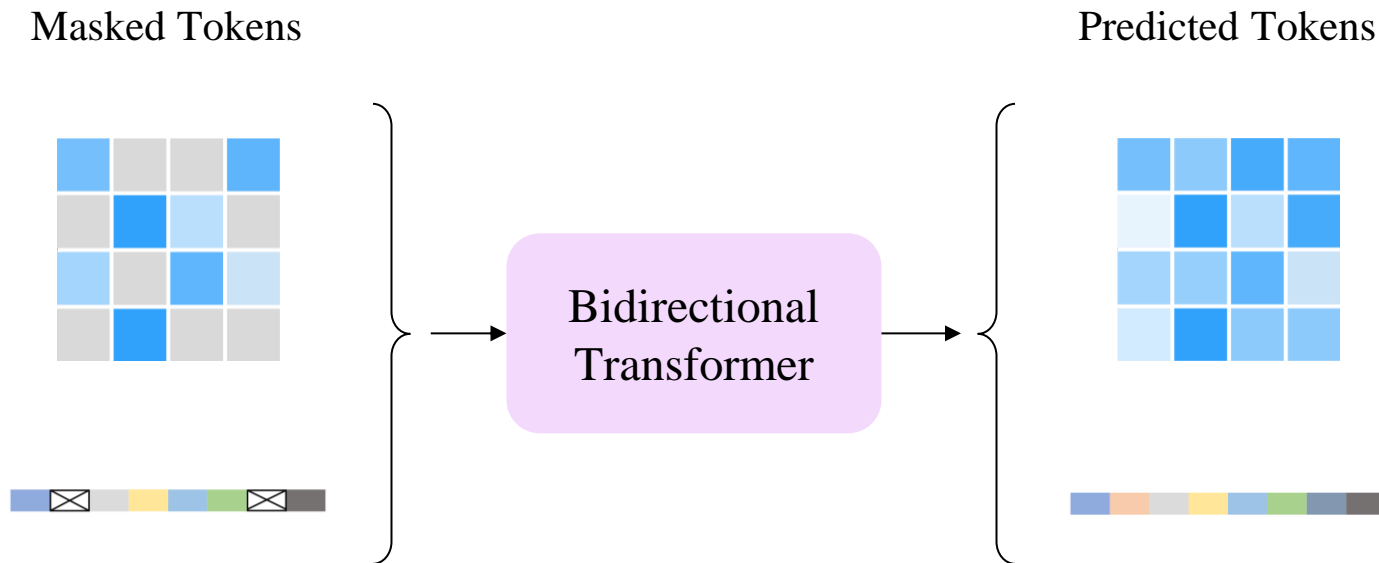
Mask Diffusion Models

Reverse Process: fill mask conditional *independently*



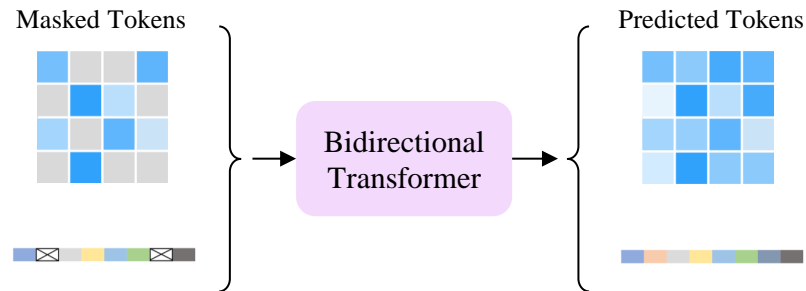
$$\text{Predicted Probability: } p_\phi(x_0^i | x_t) := \text{softmax}(z_\phi^i(x_t) / \tau)$$

Mask Diffusion Models

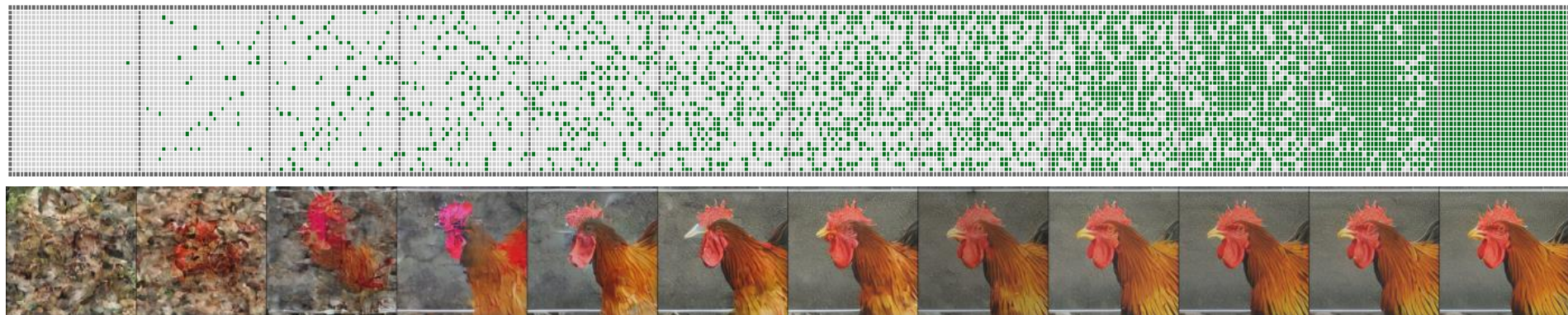


$$\mathcal{L}_{\text{MDM}} = \mathbb{E}_{x_0, t} \left[\left(\mathbb{E}_{q_{t|0}} \left[-\log p_{0|t}(x_0 | x_t, \phi) \right] \right) \right]$$

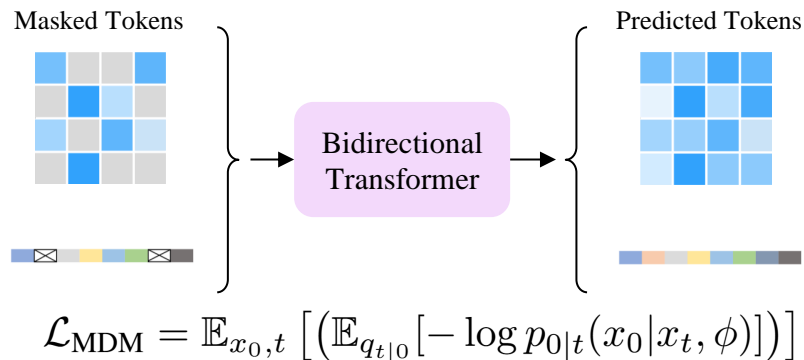
Mask Diffusion Models



Sampling Process



Promise of MDM



Theorem: Minimizing the *expected* cross-entropy loss is equivalent to maximizing the upper bound on the negative log-likelihood, i.e.:

$$-\mathbb{E}_{p_{\text{data}}(x_0)} [\log p_{\theta}(x_0)] \leq \mathcal{L}_{\text{MDM}}(\theta)$$

MDM loss (maximize likelihood) + Transformer (less inductive bias) + Data \rightarrow Scalability ?

Limitations of MDM

- Incorrect joint distribution for multi-token prediction 😞
- Need sufficient steps for better performance 😞

Ideas in Inference-time Scaling can Benefit Generative Pre-training Algorithms

Jiaming Song, Linqi Zhou

Luma AI

Theoretical Benefit and Limitation of Diffusion Language Model

Guhao Feng^{*1} Yihan Geng^{*1} Jian Guan² Wei Wu² Liwei Wang¹ Di He¹

Multi-token prediction (MTP) is of great interest to the language modeling community because of its potential to achieve faster inference [GIR⁺24], which allows efficient inference-time scaling. However, the current multi-token prediction models often predict the softmax values of multiple tokens in parallel, which is a naive conditional independence assumption (*i.e.*, naïve Bayes). We argue that this inference design greatly limits the capacity of the model distribution and more efforts should be spent resolving this fundamental issue.

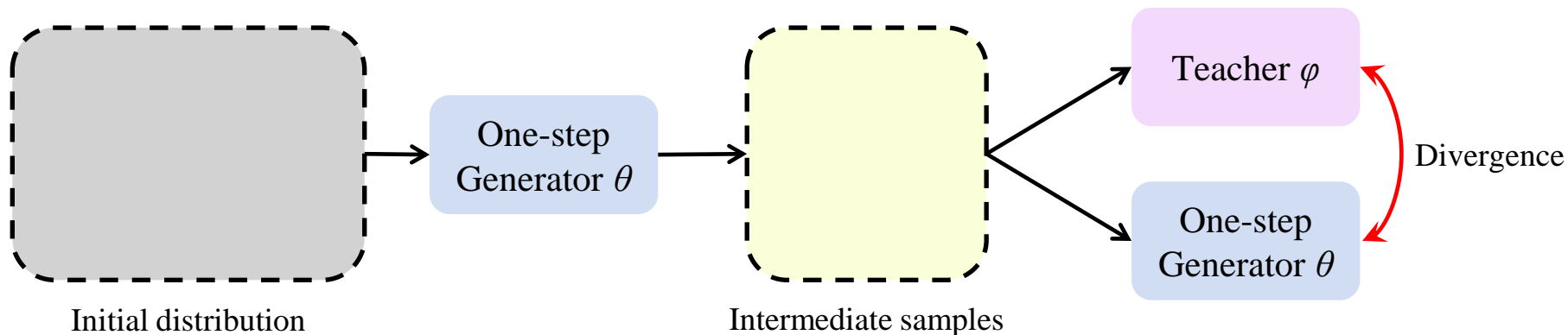
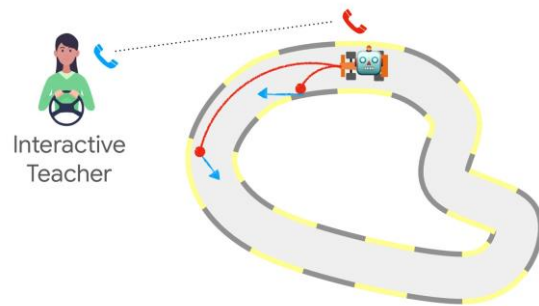
such as a reasoning chain—we show that the required sampling steps must scale linearly with sequence length to obtain “correct” sequences, thereby eliminating MDM’s efficiency advantage over autoregressive models. Our analysis estab-

Content

- Preliminaries
- Methods
- Results

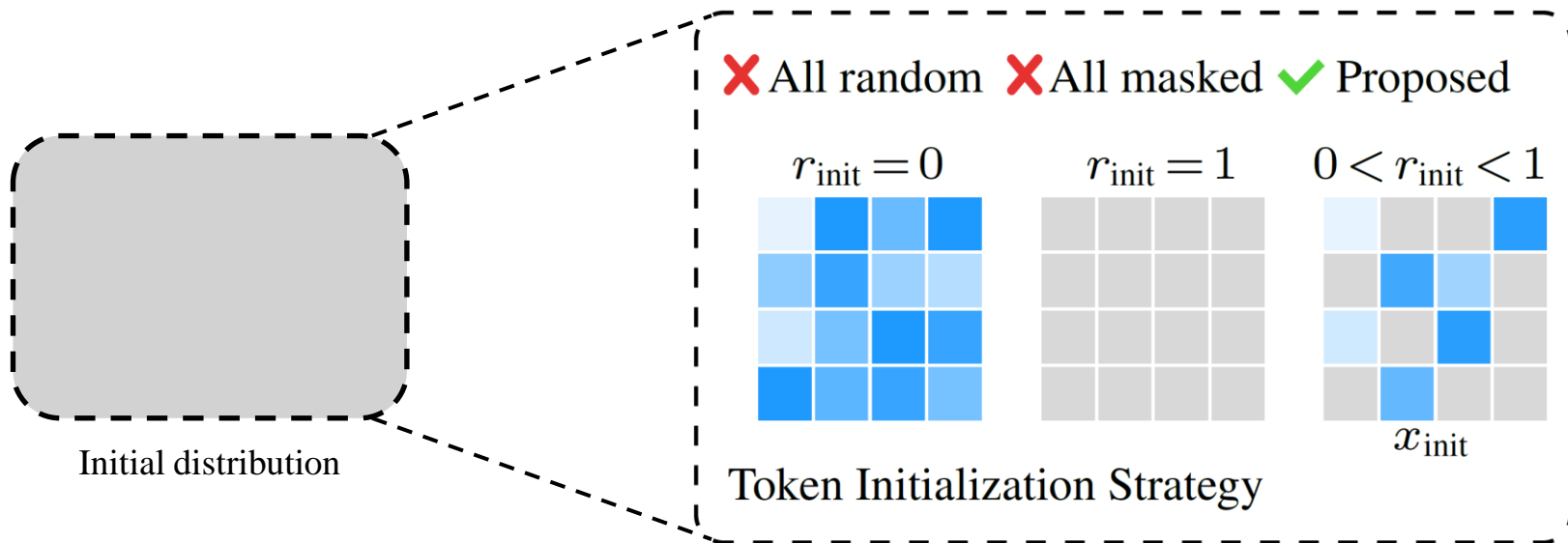
Main Idea: Distillation in an *On-Policy* fashion

Match teacher and generator for all possible intermediate samples



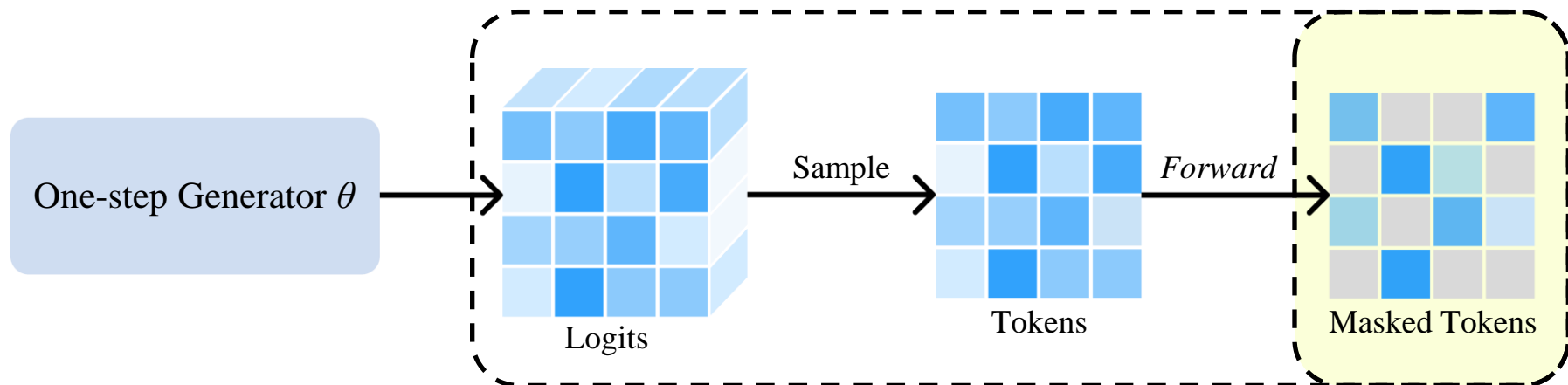
$$\mathcal{L}_{\text{Di[M]O}}(\theta) := \mathbb{E}_{x_{\text{init}}, t} \left[w(t) \left(\mathbb{E}_{q_t|o} [D(p_\phi || p_\theta)(\tilde{x}_t)] \right) \right]$$

Initial Distribution



Intermediate Samples

Via *forward mask diffusion process*



Consider loss on each intermediate state \tilde{x}_t

$$\mathcal{L}_{\text{Di}[M]\text{O}}(\theta) := \mathbb{E}_{x_{\text{init}}, t} \left[w(t) \left(\mathbb{E}_{q_t | o} \left[D(p_\phi || p_\theta)(\tilde{x}_t) \right] \right) \right]$$

Further decompose into *Token-level Divergence*

$$D(p_\phi || p_\theta)(\tilde{x}_t) := \frac{1}{L_M} \sum_{\substack{i=1 \\ \tilde{x}_t^i = [M]}}^L D(p_\phi(x_0^i | \tilde{x}_t) || p_\theta(x_0^i | \tilde{x}_t))$$

Gradient of Divergence:

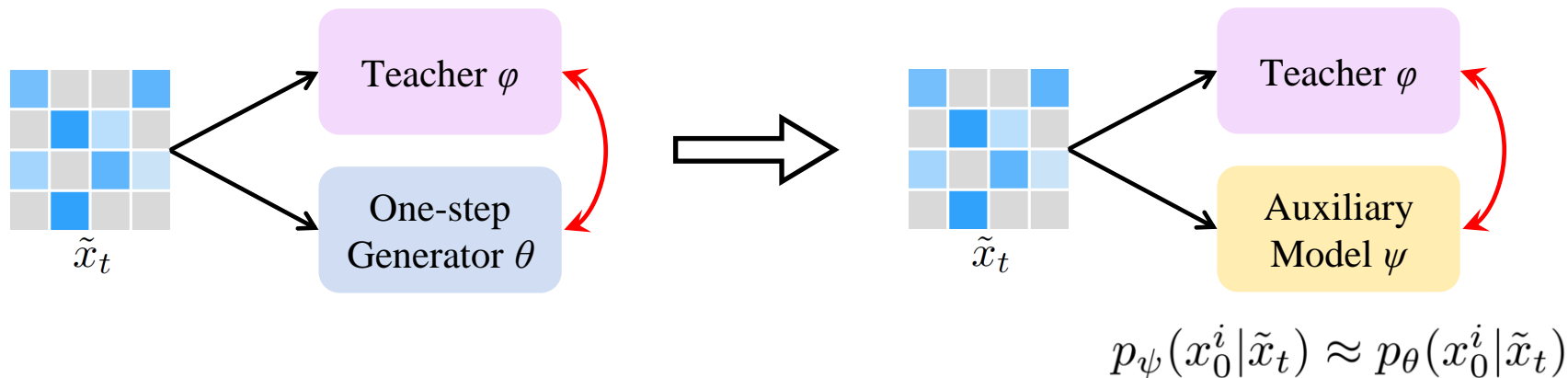
$$\nabla_{\theta} D((p_{\phi} || p_{\theta})(\tilde{x}_t)) = \nabla_{z_{\theta}} D(p_{\phi} || p_{\theta})(\tilde{x}_t) \frac{dz_{\theta}(\tilde{x}_t)}{d\theta}$$

Intractable terms

$$p_{\theta}(x_0^i | \tilde{x}_t) = \text{softmax}(z_{\theta}(\tilde{x}_t))$$

Approximation: $p_{\theta}(x_0^i | \tilde{x}_t)$

Introduce an *auxiliary model* to approximate the student output on intermediate samples

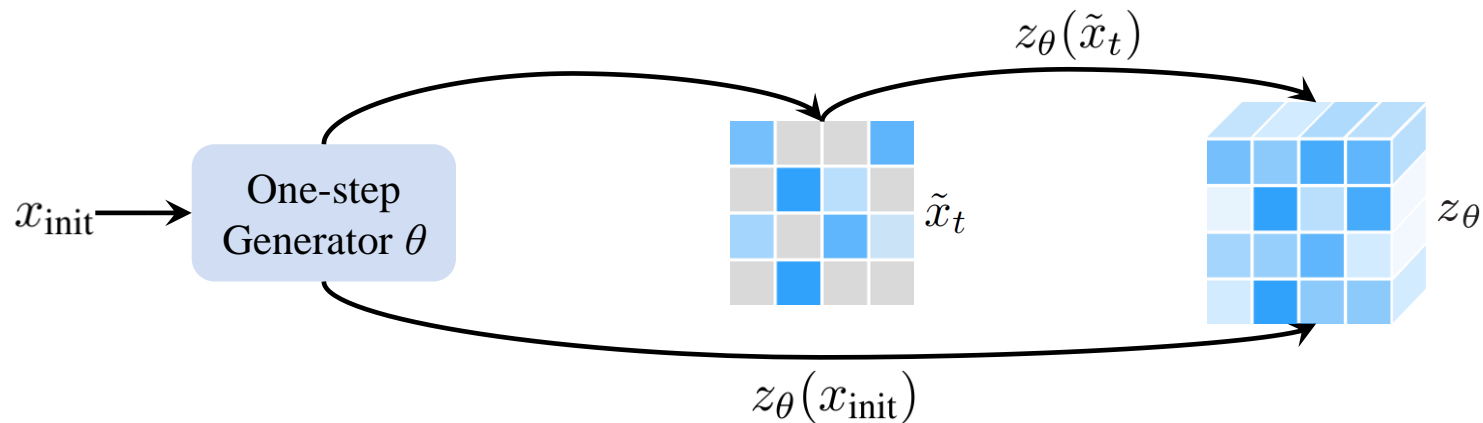


$$\nabla_{z_{\theta}} D(p_{\phi} || p_{\theta})(\tilde{x}_t) \quad \Rightarrow \quad \nabla_{z_{\psi}} D(p_{\phi} || p_{\psi})(\tilde{x}_t)$$

Approximation: $z_\theta(\tilde{x}_t)$

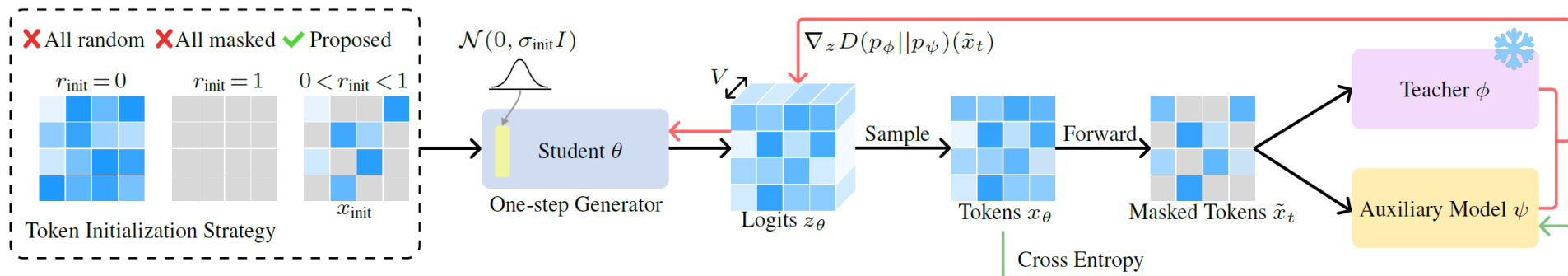
We approximate the output logits based on the *consistency assumption*

Also, we require this term as effective gradient for the generator



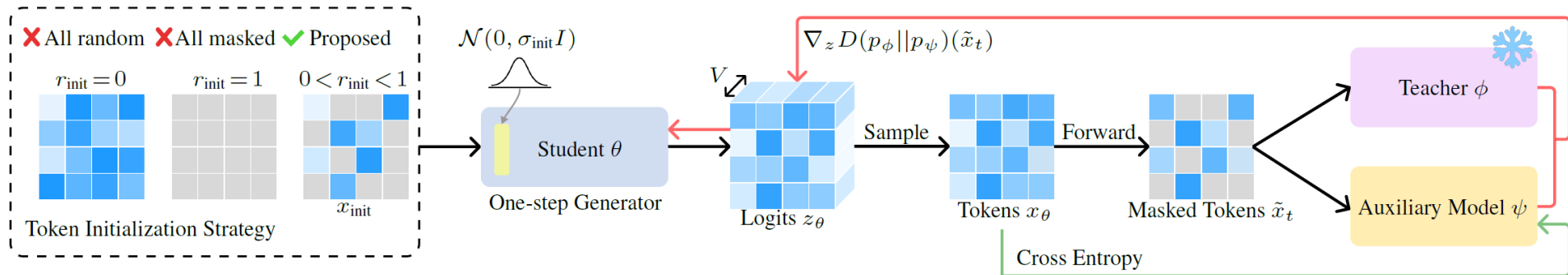
$$\frac{dz_\theta(\tilde{x}_t)}{d\theta} \Rightarrow \frac{dz_\theta(x_{\text{init}})}{d\theta}$$

Overview of the Method

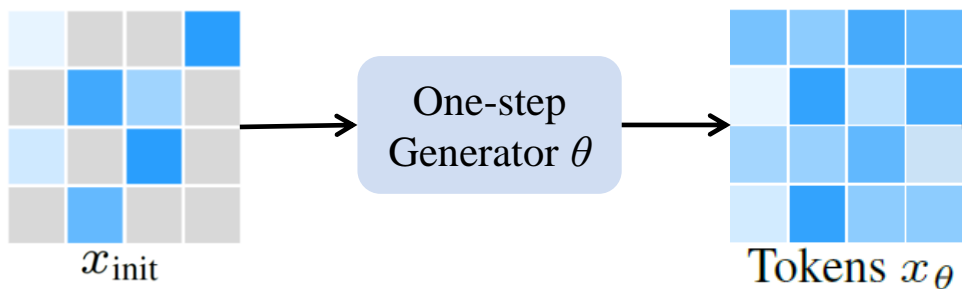


$$\nabla_\theta \mathcal{L}_{\text{Di[M]O}} \approx \mathbb{E}_{x_{\text{init}}, t} \left[w(t) \left(\mathbb{E}_{q_t | 0} \left[\nabla_{z_\psi} D(p_\phi || p_\psi)(\tilde{x}_t) \quad \frac{dz_\theta(x_{\text{init}})}{d\theta} \right] \right) \right]$$

Overview of the Method



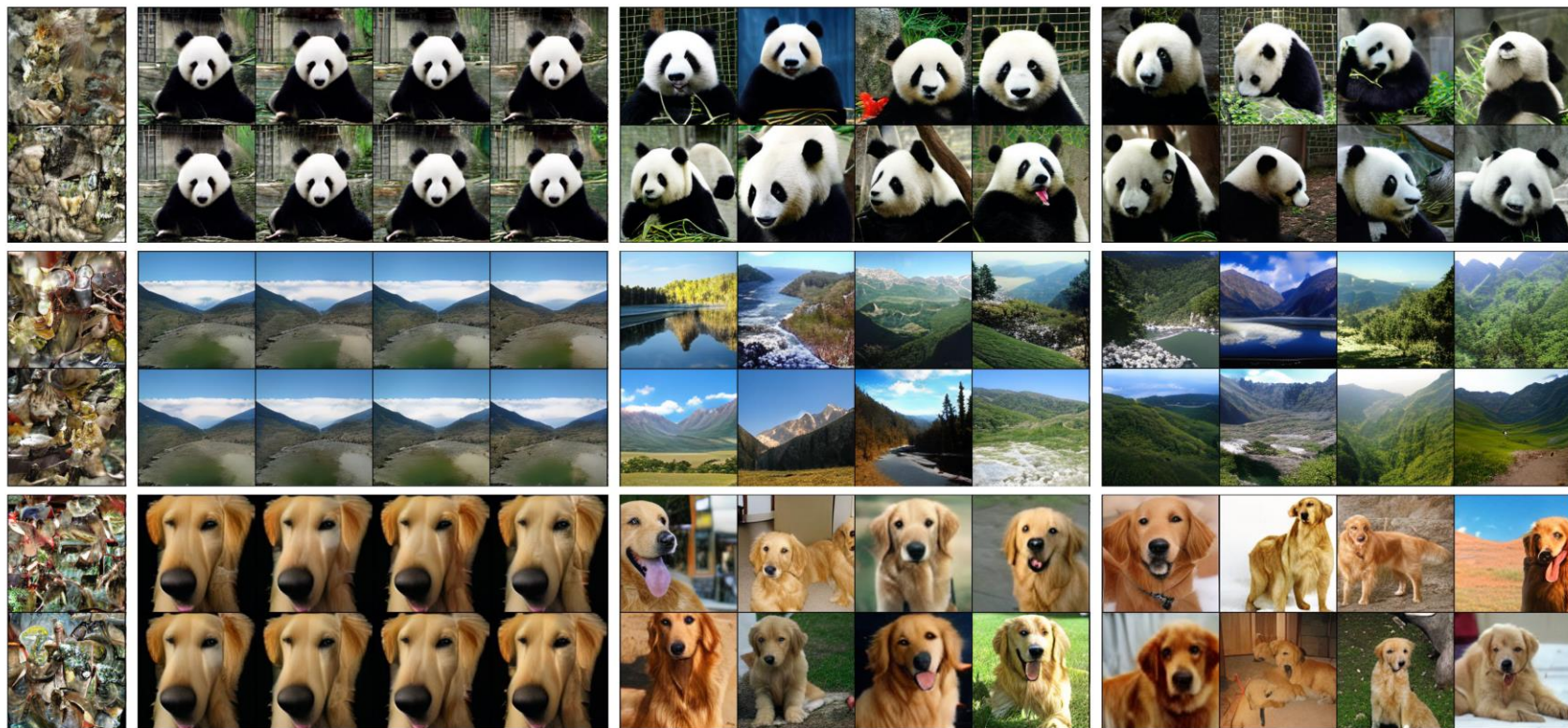
- From random initial states to samples from *correct joint distribution* 😊
- Greatly reduce the sampling steps \rightarrow One step generation 😊



Content

- Preliminaries
- Methods
- Results

Visual Results



$r_{\text{init}} = 0$

Di[M]O ($r_{\text{init}} = 1$) one step

Di[M]O ($r_{\text{init}} = 0.6$) one step

Teacher 16 steps

Visual Results

64
Steps



16
Steps



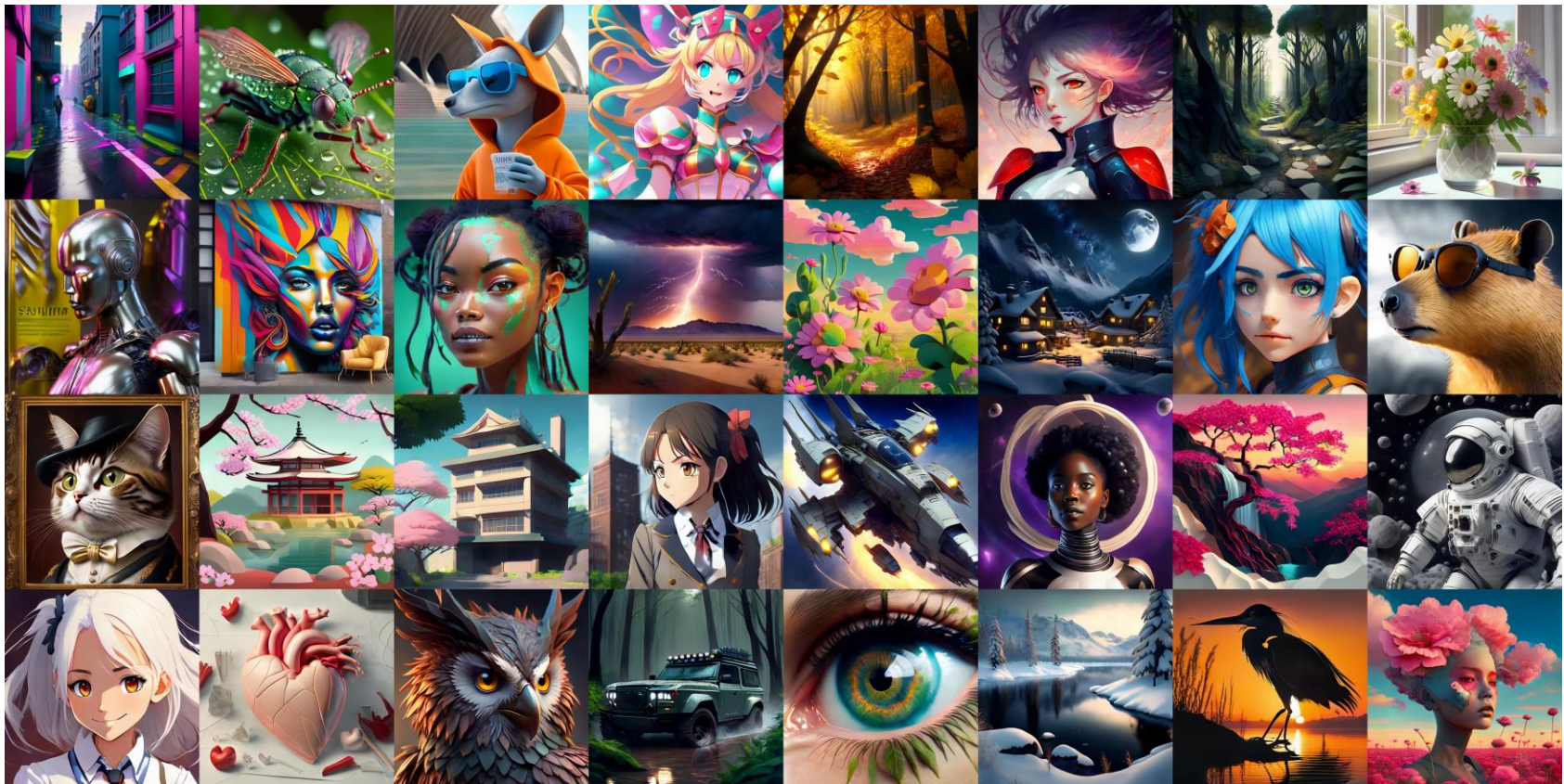
4
Steps



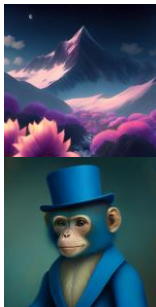
1
Step
(ours)



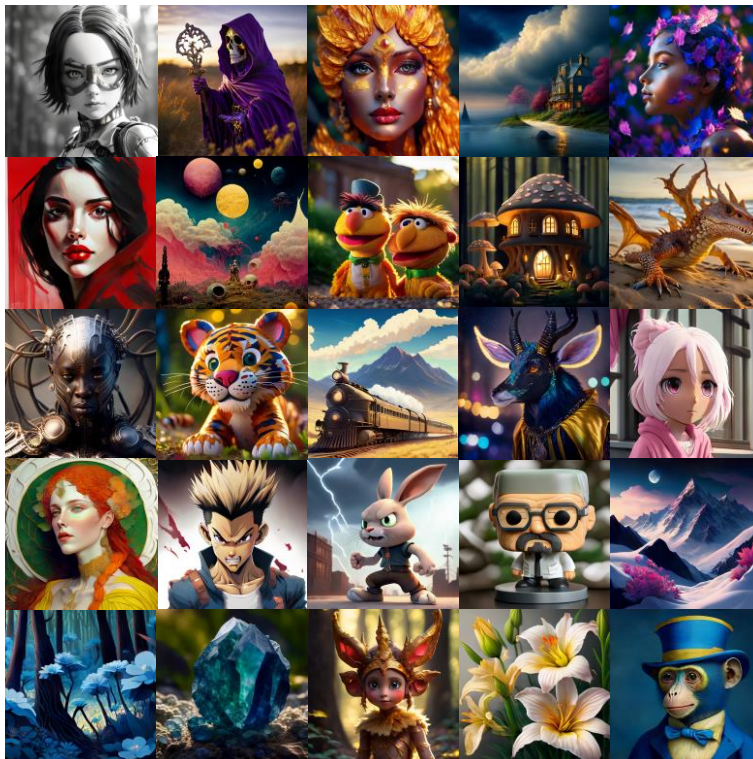
Visual Results



Runtime Comparison



10.5 s/image



0.80 s/image

Quantitative Results (MaskGit Teacher)

ImageNet	Method	Step (\downarrow)	FID (\downarrow)	IS (\uparrow)	Precision (\uparrow)	Recall (\uparrow)	Density (\uparrow)	Coverage (\uparrow)
Teacher	MaskGit [7]	16	6.60	224.07	0.831	0.402	1.246	0.977
	MaskGit [7]	8	6.66	221.57	0.827	0.397	1.233	0.974
	MaskGit [7]	4	10.73	192.29	0.748	0.313	1.011	0.920
	MaskGit [7]	2	91.35	13.37	0.178	0.164	0.091	0.122
Sampler	θ -trapezoidal [81]*	64	6.7	-	-	-	-	-
	θ -trapezoidal [81]*	32	7.1	-	-	-	-	-
Distillation	di4c [31]	4	6.79	209.2	-	-	-	-
	di4c-d [31]	4	6.57	213.6	-	-	-	-
	Di[M]O	1	6.91	214.0	0.828	0.377	1.255	0.967

Quantitative Results (Meissonic Teacher)

HPS v2.0

Model	Step	Anim.	Concept-art	Painting	Photo	Averaged
Latent Diffusion [83]	25	25.73	25.15	25.25	26.97	25.78
DALL-E 2 [79]	-	27.34	26.54	26.68	27.24	26.95
Stable Diffusion v1.4 [83]	50	27.26	26.61	26.66	27.27	26.95
Stable Diffusion v2.0 [83]	50	27.48	26.89	26.86	27.46	27.17
DeepFloyd-XL [19]	25	27.64	26.83	26.86	27.75	27.27
SDXL Base 1.0 [76]	50	28.88	27.88	27.92	28.31	28.25
SDXL Refiner 1.0 [76]	50	28.93	27.89	27.90	28.38	28.27
InstaFlow [55]	1	25.98	25.79	25.93	26.32	26.01
SD Turbo [90]	1	27.98	27.59	27.16	27.19	27.48
SwiftBrush v2 [18]	1	27.25	27.62	26.86	26.77	27.15
Meissonic [4]	48	29.57	28.58	28.72	28.45	28.83
	32	29.18	28.32	28.28	27.96	28.44
	16	28.61	27.82	27.84	27.32	27.90
	8	25.62	26.49	26.67	27.07	26.46
	4	25.01	24.95	24.87	23.80	24.66
	2	23.06	23.28	23.22	22.38	22.98
Di[M]O	1	28.64	27.91	27.99	27.92	28.11

GenEval

Model	Step	Overall	Objects		Counting	Colors	Position	Color Attribution
			Single	Two				
SD v1.5 [83]	50	0.43	0.97	0.38	0.35	0.76	0.04	0.06
SD v2.1 [83]	50	0.50	0.98	0.51	0.44	0.85	0.07	0.17
DALL-E2 [79]	-	0.52	0.94	0.66	0.49	0.77	0.10	0.19
SDXL [76]	50	0.55	0.98	0.74	0.39	0.85	0.15	0.23
Meissonic [4]	48	0.54	0.99	0.66	0.42	0.86	0.10	0.22
	32	0.46	0.92	0.53	0.33	0.80	0.08	0.13
	16	0.37	0.82	0.39	0.20	0.70	0.05	0.08
	8	0.20	0.58	0.12	0.05	0.40	0.02	0.04
	4	0.09	0.31	0.02	0.01	0.18	0.01	0.01
	2	0.03	0.14	0.01	0.00	0.05	0.00	0.00
Di[M]O	1	0.43	0.91	0.53	0.22	0.75	0.07	0.11

Thank You!