

Simple Guidance Mechanisms for Discrete Diffusion Models



**Yair
Schiff***



**Subham
Sahoo***



**Hao
Phung***



**Guanghan
Wang***



**Sam
Boshar**



**Hugo
Dalla-torre**



**Bernardo P
de Almeida**



**Alexander
Rush**

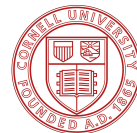


**Thomas
Pierrot**



**Volodymyr
Kuleshov**

**Equal contribution*



Write a function for LLM infer

Iterations

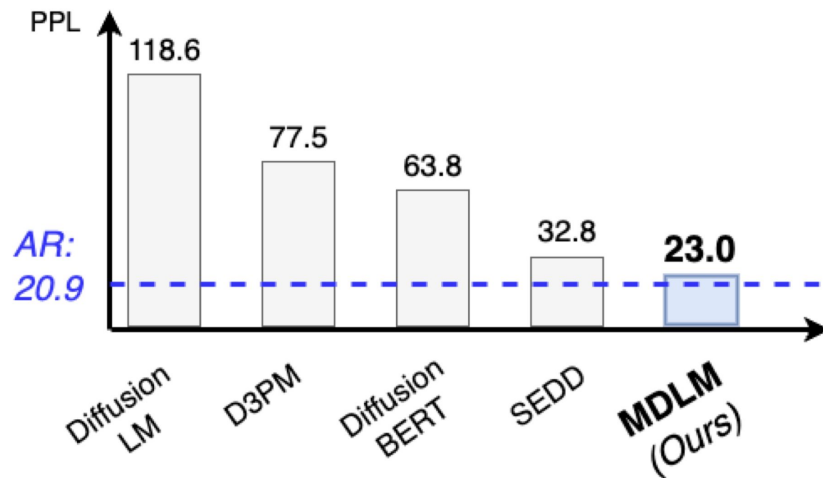
0

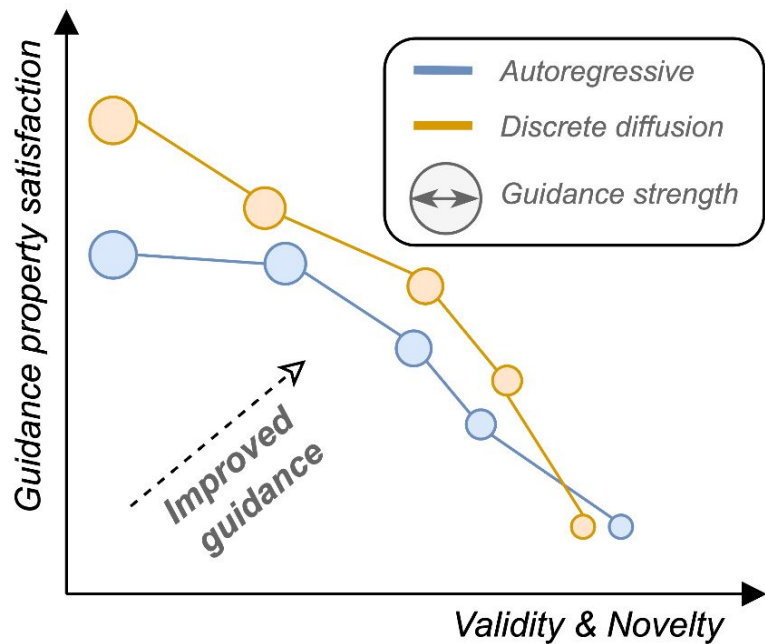
AUTOREGRESSIVE LLM
LEFT-TO-RIGHT GENERATION

Iterations

0

INCEPTION DIFFUSION LLM
COARSE-TO-FINE GENERATION





Autoregressive models:
make “**local**” predictions



Diffusion models: make
“**global**” refinements



An ink sketch style illustration of a small hedgehog holding a piece of watermelon with its tiny paws, taking little bites with its eyes closed in delight.



The 'bling zoo' shop in new york city is both a jewelry store and zoo. sabertooth tigers with diamond and gold adornments...

Notation

s, t timesteps, $s < t \in [0, 1]$

\mathbf{z}_t latent vector

Continuous: $\mathbf{z}_t \in \mathbb{R}^d$

Discrete: $\mathbf{z}_t^{(1:L)} \in |\mathcal{V}|^L$ (vocab \mathcal{V} , seq len L)

Guidance Background

Sample from **conditional**, **tempered** distribution

$$\mathbf{z}_s \sim p^{\boxed{\gamma}}(\mathbf{z}_s \mid \mathbf{z}_t, \boxed{y})$$

Sample from **conditional**, tempered distribution

$$\mathbf{z}_s \sim p(\mathbf{z}_s \mid \mathbf{z}_t, y)$$



$$p(\mathbf{z}_s \mid \mathbf{z}_t, y) \propto p(y \mid \mathbf{z}_s, \mathbf{z}_t) \cdot p(\mathbf{z}_s \mid \mathbf{z}_t)$$



Sample from conditional, **tempered** distribution

$$\mathbf{z}_s \sim p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y)$$



$$p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) \propto p(y \mid \mathbf{z}_s, \mathbf{z}_t)^\gamma \cdot p(\mathbf{z}_s \mid \mathbf{z}_t)$$

Continuous: Follow log-prob gradient with respect to input

$$p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) \propto p(y \mid \mathbf{z}_s, \mathbf{z}_t)^\gamma \cdot p(\mathbf{z}_s \mid \mathbf{z}_t)$$



$$\nabla_{\mathbf{z}_s} \log p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) = \gamma \nabla_{\mathbf{z}_s} \log p(y \mid \mathbf{z}_s, \mathbf{z}_t) + \nabla_{\mathbf{z}_s} \log p(\mathbf{z}_s \mid \mathbf{z}_t)$$

Continuous: Follow log-prob gradient with respect to input

$$p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) \propto p(y \mid \mathbf{z}_s, \mathbf{z}_t)^\gamma \cdot p(\mathbf{z}_s \mid \mathbf{z}_t)$$



$$\nabla_{\mathbf{z}_s} \log p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) = \gamma \nabla_{\mathbf{z}_s} \log p(y \mid \mathbf{z}_s, \mathbf{z}_t) + \nabla_{\mathbf{z}_s} \log p(\mathbf{z}_s \mid \mathbf{z}_t)$$

**Classifier-based
guidance**

External classifier

Diffusion model

Continuous: Follow log-prob gradient with respect to input

$$\nabla_{\mathbf{z}_s} \log p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) = \gamma \nabla_{\mathbf{z}_s} \log p(y \mid \mathbf{z}_s, \mathbf{z}_t) + \nabla_{\mathbf{z}_s} \log p(\mathbf{z}_s \mid \mathbf{z}_t)$$

$$\nabla_{\mathbf{z}_s} \log p(\mathbf{z}_s \mid \mathbf{z}_t, y) - \nabla_{\mathbf{z}_s} \log p(\mathbf{z}_s \mid \mathbf{z}_t)$$



Continuous: Follow log-prob gradient with respect to input

$$\nabla_{\mathbf{z}_s} \log p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) = \gamma \nabla_{\mathbf{z}_s} \log p(y \mid \mathbf{z}_s, \mathbf{z}_t) + \nabla_{\mathbf{z}_s} \log p(\mathbf{z}_s \mid \mathbf{z}_t)$$



$$\nabla_{\mathbf{z}_s} \log p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) = \gamma \nabla_{\mathbf{z}_s} \log p(\mathbf{z}_s \mid \mathbf{z}_t, y) + (1 - \gamma) \nabla_{\mathbf{z}_s} \log p(\mathbf{z}_s \mid \mathbf{z}_t)$$

Continuous: Follow log-prob gradient with respect to input

$$\nabla_{\mathbf{z}_s} \log p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) = \gamma \nabla_{\mathbf{z}_s} \log p(y \mid \mathbf{z}_s, \mathbf{z}_t) + \nabla_{\mathbf{z}_s} \log p(\mathbf{z}_s \mid \mathbf{z}_t)$$



$$\nabla_{\mathbf{z}_s} \log p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) = \gamma \nabla_{\mathbf{z}_s} \log p(\mathbf{z}_s \mid \mathbf{z}_t, y) + (1 - \gamma) \nabla_{\mathbf{z}_s} \log p(\mathbf{z}_s \mid \mathbf{z}_t)$$

**Classifier-free
guidance**

Conditional
diffusion model

Unconditional
diffusion model

Discrete Guidance

Discrete:



\mathbf{z}_s



Discrete: Back to the drawing board

$$\log p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) = \gamma \log p(y \mid \mathbf{z}_s, \mathbf{z}_t) + \log p(\mathbf{z}_s \mid \mathbf{z}_t) + \textit{const}.$$

Discrete: Preview of where we'll end up

$$p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) \propto \\ (p_\theta) \cdot (\text{some other scaling distribution})$$

Discrete Classifier-Free Guidance

Discrete: Classifier-free guidance

$$\log p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) = \boxed{\gamma \log p(y \mid \mathbf{z}_s, \mathbf{z}_t)} + \log p(\mathbf{z}_s \mid \mathbf{z}_t) + \text{const.}$$



$$\gamma \log p(\mathbf{z}_s \mid \mathbf{z}_t, y) - \gamma \log p(\mathbf{z}_s \mid \mathbf{z}_t) + \underbrace{\gamma \log p(\mathbf{y} \mid \mathbf{z}_t)}_{\text{constant}}$$



Discrete: Classifier-free guidance

$$\log p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) = \gamma \log p(y \mid \mathbf{z}_s, \mathbf{z}_t) + \log p(\mathbf{z}_s \mid \mathbf{z}_t) + \text{const.}$$



$$\log p^\gamma(\mathbf{z}_s \mid \mathbf{z}_t, y) = \gamma \log p(\mathbf{z}_s \mid \mathbf{z}_t, y) + (1 - \gamma) \log p(\mathbf{z}_s \mid \mathbf{z}_t) + \text{const.}$$

Conditional
diffusion model

Unconditional
diffusion model

Discrete Classifier-Based Guidance

Discrete: Classifier-based guidance

$$p^{\gamma}(\mathbf{z}_s \mid \mathbf{z}_t, y) \propto p(y \mid \mathbf{z}_s, \mathbf{z}_t)^{\gamma} \cdot p(\mathbf{z}_s \mid \mathbf{z}_t)$$

Discrete: Classifier-based guidance

$$p^{\gamma}(\mathbf{z}_s \mid \mathbf{z}_t, y) \propto p(y \mid \mathbf{z}_s, \mathbf{z}_t)^{\gamma} \cdot p(\mathbf{z}_s \mid \mathbf{z}_t)$$



In practice, we model **sequences of text**, not just individual tokens

$$\mathbf{z} \rightarrow \mathbf{z}^{(1:L)}$$



$$p(y \mid \mathbf{z}_s, \mathbf{z}_t) \rightarrow p(y \mid \mathbf{z}_s^{(1:L)}, \mathbf{z}_t^{(1:L)})$$



$$\mathcal{O}(|V|^L) \text{ terms}$$



Assume $p^\gamma(\mathbf{z}_s^{(1:L)} \mid \mathbf{z}_t^{(1:L)}, y) = \prod_{\ell=1}^L p^\gamma(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)}, y)$



We will model
this per token
distribution
instead

$$p(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)}, y) \propto p(y \mid \mathbf{z}_s^{(\ell)}, \mathbf{z}_t^{(1:L)}) \cdot p(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)})$$

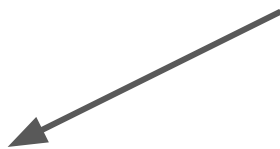
$$p(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)}, y) \propto p(y \mid \mathbf{z}_s^{(\ell)}, \mathbf{z}_t^{(1:L)}) \cdot p(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)})$$



Model with
discrete diffusion

p_θ

$$p(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)}, y) \propto p(y \mid \mathbf{z}_s^{(\ell)}, \mathbf{z}_t^{(1:L)}) \cdot p(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)})$$



Train classifier on
noisy sequences:

$$p_\phi(\mathbf{z}_t^{(1:L)})$$



Evaluate it on sequences of
the form:

$$p_\phi(y \mid [\mathbf{z}_t^{(1:\ell-1)}, \mathbf{z}_s^{(\ell)}, \mathbf{z}_t^{(\ell+1:L)}])$$

Putting it all together

Discrete Classifier-free guidance

$$p^\gamma(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)}, y) \propto p(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)}, y)^\gamma \cdot p(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)})^{(1-\gamma)}$$

Discrete Classifier-based guidance

$$p^\gamma(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)}, y) \propto p_\phi(y \mid [\mathbf{z}_t^{(1:\ell-1)}, \mathbf{z}_s^{(\ell)}, \mathbf{z}_t^{(\ell+1:L)}])^\gamma \cdot p(\mathbf{z}_s^{(\ell)} \mid \mathbf{z}_t^{(1:L)})$$

Uniform Diffusion Language Models

For Masked Diffusion Models, generated tokens are ‘fixed’

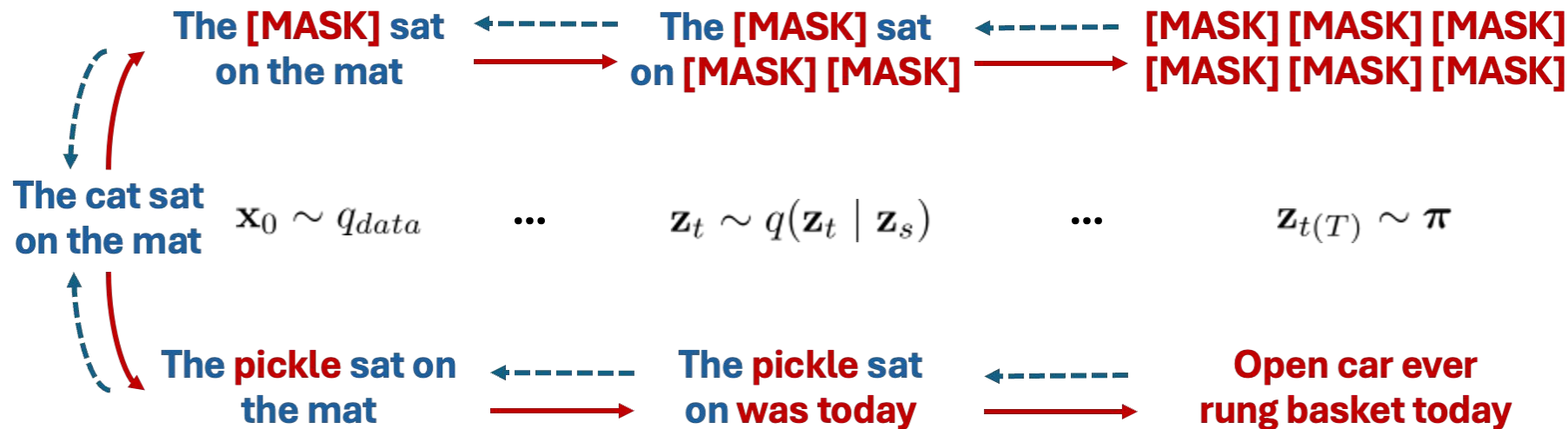
$$q(\mathbf{z}_s \mid \mathbf{z}_t, \mathbf{x}) = \begin{cases} \text{Cat}(\mathbf{z}_s; \mathbf{z}_t), & \mathbf{z}_t \neq \mathbf{m} \\ \text{Cat}(\mathbf{z}_s; \frac{\alpha_s - \alpha_t}{1 - \alpha_t} \mathbf{x} + \frac{1 - \alpha_s}{1 - \alpha_t} \mathbf{m}), & \mathbf{z}_t = \mathbf{m} \end{cases}$$



Unmasked tokens are
‘locked-in’ (even if incorrect!)

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

π = absorbing state



$\leftarrow \text{---} p_{\theta}(\mathbf{z}_s | \mathbf{z}_t)$

$\rightarrow q(\mathbf{z}_t | \mathbf{z}_s)$

Experiments

Table 1: UDLM performs best with smaller vocabs. Perplexity (\downarrow) on various datasets. Best values are **bolded**. * indicates values reported from early stopping on the validation set; otherwise validation performance at the end of training is used. [†]From [Sahoo et al. \(2024a\)](#). ^{\$}From [Lou et al. \(2023\)](#).

	Vocab.	AR	MDLM	UDLM
Species10	12	2.88	3.17 _≤	3.15 _≤
QM9*	40	2.19	2.12 _≤	2.02 _≤
CIFAR10	256	-	9.14 _≤	11.21 _≤
text8	35	2.35 ^{\$}	2.62 _≤	2.71 _≤
Amazon*	30,522	21.67	24.93 _≤	27.27 _≤
LM1B	30,522	22.32 [†]	27.04 _≤ [†]	31.28 _≤

Table 2: UDLM outperforms other uniform discrete diffusion on text8. Best value is **bolded** & best uniform diffusion value is underlined.
[†]From Lou et al. (2023). *From Shi et al. (2024).

Method	BPC (\downarrow)
<i>Discrete Uniform Diffusion</i>	
D3PM Uniform [†] (Austin et al., 2021)	1.61 \leq
SEDD Uniform [†] (Lou et al., 2023)	1.47 \leq
UDLM (<i>Ours</i>)	<u>1.44\leq</u>

Table 3: UDLM outperforms other uniform discrete diffusion on LM1B. Best value is **bolded** & best discrete uniform diffusion value is underlined.
[†]From Sahoo et al. (2024a). *From Lou et al. (2023). [§]From Austin et al. (2021).

Method	PPL (\downarrow)
<i>Discrete Uniform Diffusion</i>	
D3PM Uniform [§] (Austin et al., 2021)	137.9 \leq
SEDD Uniform* (Lou et al., 2023)	40.25 \leq
UDLM (<i>Ours</i>)	<u>31.28\leq</u>

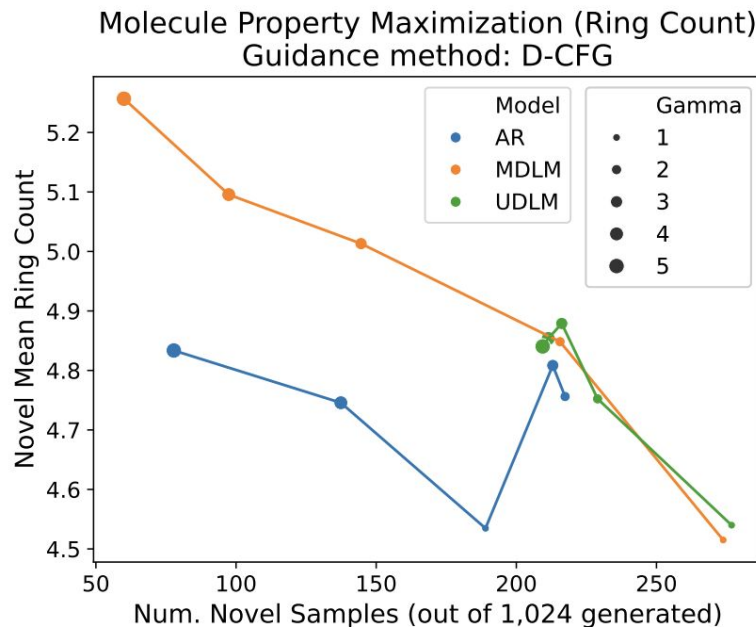
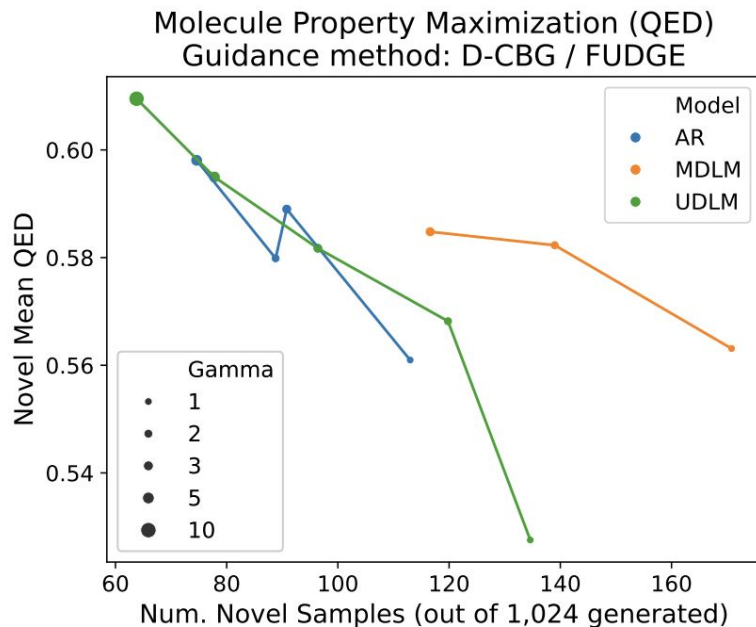
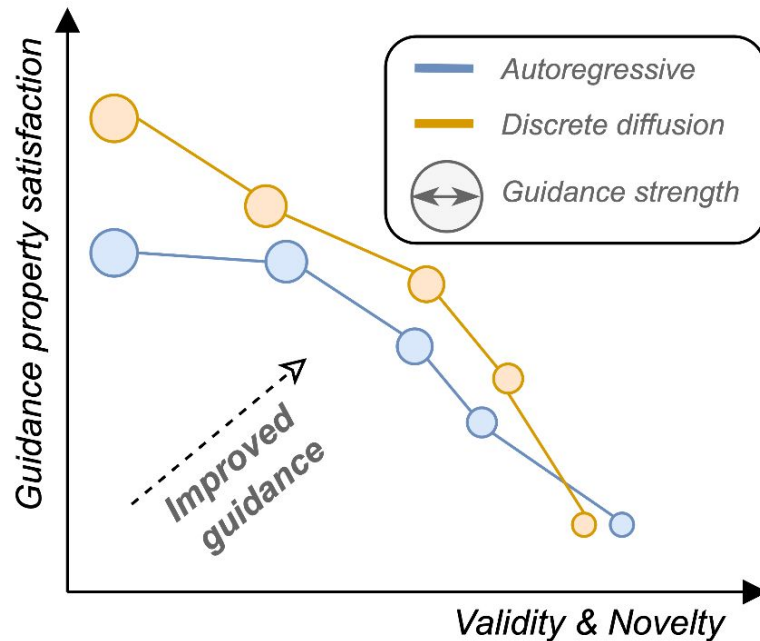
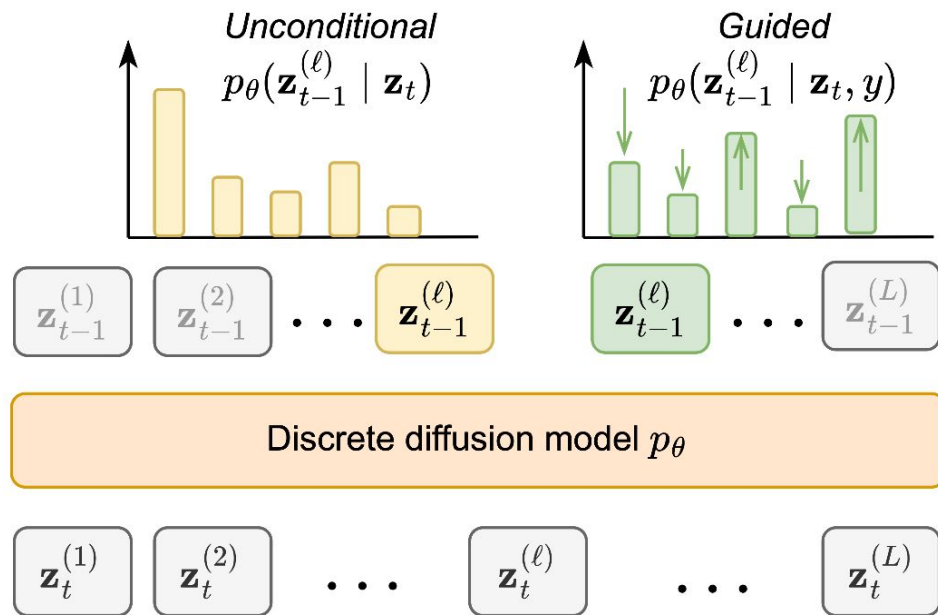


Figure 3: Diffusion models extend the steer-ability Pareto frontier. (*Left*) D-CBG outperforms FUDGE classifier guidance when maximizing drug-likeness (QED). (*Right*) D-CFG with diffusion better trades-off novel generation and ring-count maximization compared to AR.

Conclusion



<https://arxiv.org/abs/2412.10193>



<https://github.com/kuleshov-group/discrete-diffusion-guidance>

Thank you!

Thursday April 24
Hall 3 + Hall 2B
Poster session 1

