

Navigating the Design Space of Equivariant Diffusion-Based Generative Models for De Novo 3D Molecule Generation

TUAN LE^{*1,2}, JULIAN CREMER^{*1,3}, FRANK NOÉ^{2,4}, DJORK-ARNÉ CLEVERT¹ AND KRISTOF T. SCHÜTT¹

¹ Pfizer Research & Development

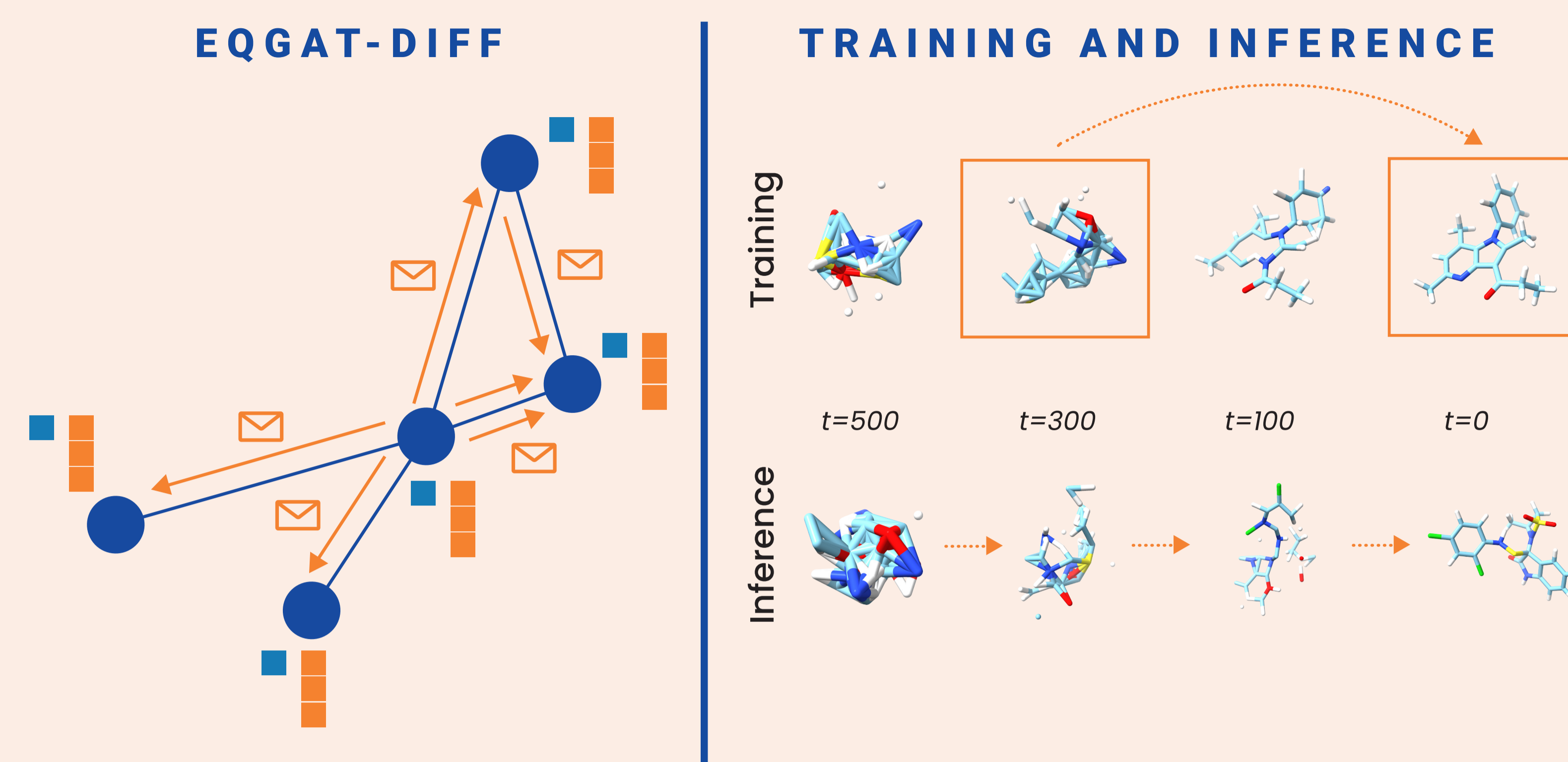
³ Universitat Pompeu Fabra, Barcelona Biomedical Research Park (PRBB)

² Freie Universität Berlin

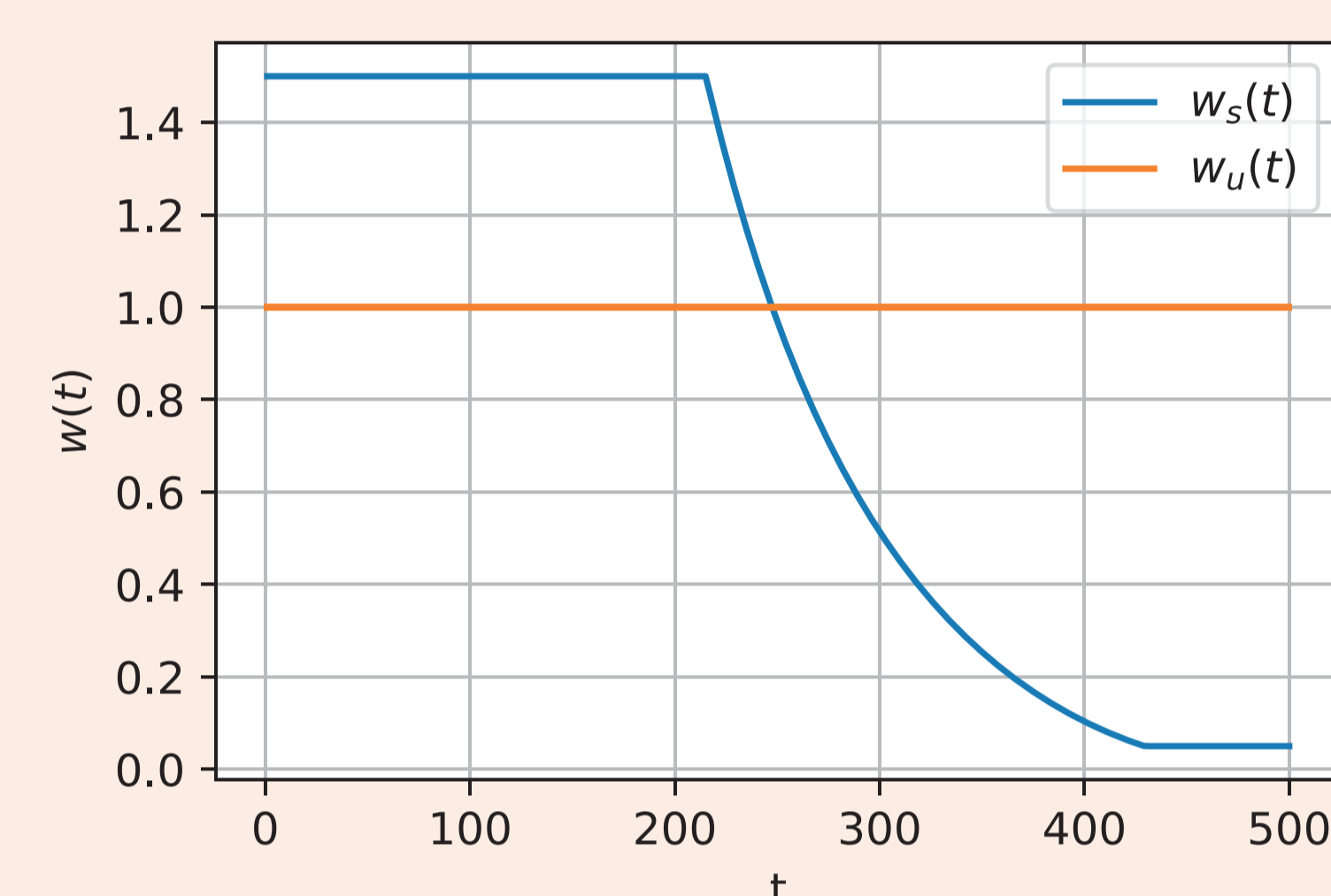
⁴ Microsoft Research

*Equal contribution

Overview



TIMESTEP-DEPENDENT LOSS WEIGHTING



$$L_{t,\epsilon} = w(t) \|\epsilon_t - \hat{\epsilon}_\theta(x_t, t)\|^2$$

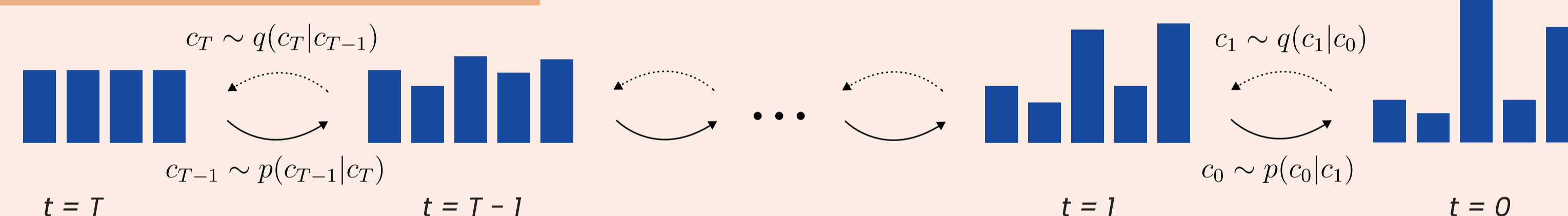
$$L_{t,x_0} = w(t) \cdot l_d(x_0, \hat{x}_\theta(x_t, t); \lambda_m)$$

$$w_s(t) = \max(0.05, \min(1.5, \text{SNR}(t)))$$

CATEGORICAL AND CONTINUOUS DIFFUSION

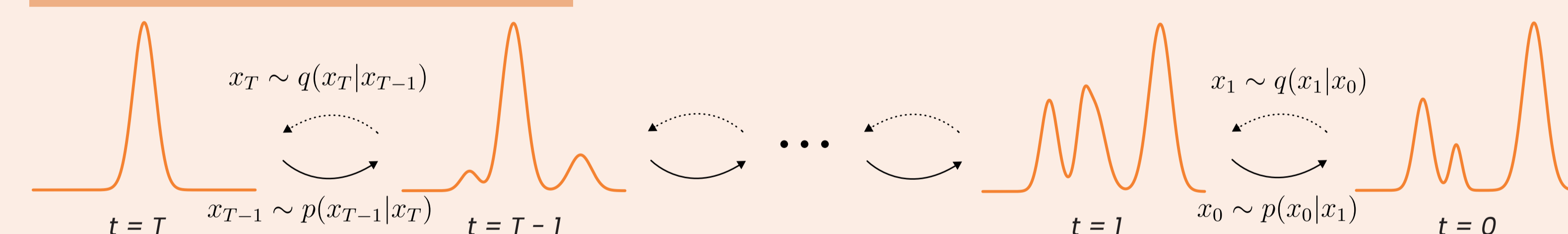
Categorical diffusion of atom types, charges and bond types

$$q(c_t|c_0) = C(c_t|\bar{a}_t c_0 + (1 - \bar{a}_t) \bar{c})$$

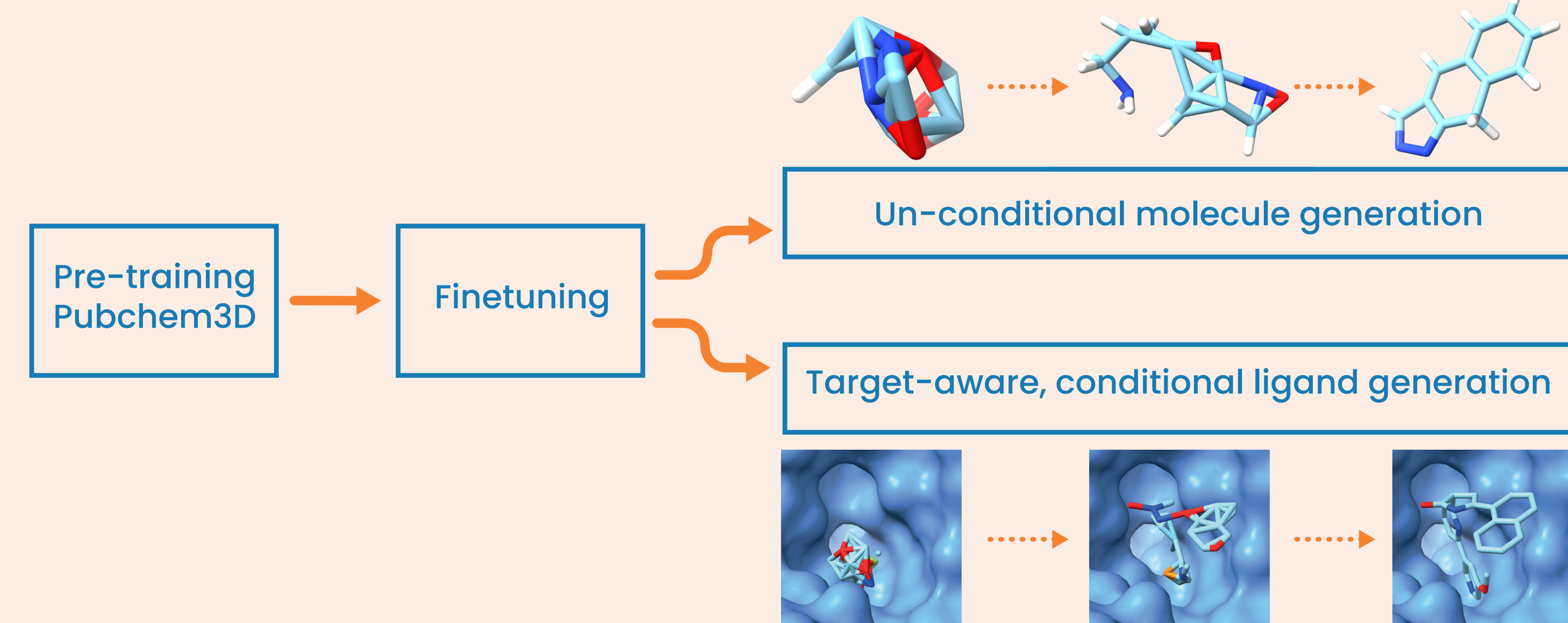


Continuous diffusion of atom positions

$$q(x_t|x_0) = \mathcal{N}(x_t|\sqrt{\bar{a}_t}x_0, (1 - \bar{a}_t)\mathbf{I})$$



TRANSFERABILITY OF DIFFUSION MODELS



MOTIVATION AND BACKGROUND

We propose EQGAT-diff, a generative model p_θ for 3D molecules $x = (\mathbf{H}, \mathbf{X}, \mathbf{E})$ with $\mathbf{H} \in \{0, 1\}^{N \times K_a}$, $\mathbf{E} \in \{0, 1\}^{N \times N \times K_b}$, $\mathbf{X} \in \mathbb{R}^{N \times 3}$, operating on continuous and discrete variables invariant to the permutation of atoms and roto-translations, i.e., $p_\theta(x) = p_\theta(g \cdot x)$

Research Questions:

1. Prior works using equivariant autoregressive or diffusion models for molecule design showed suboptimal performance on more complex data distributions, like GEOM-Drugs. Are those fundamental architectural limitations or can we improve the performance with better design choices and which are those?
2. Should we respect the data modalities and apply data-dependent noising and de-noising steps? That is, for continuous data, Gaussian diffusion and for discrete data, categorical diffusion.
3. For continuous variables, should the noise or the data parameterization be used, i.e., ϵ or \hat{x}_θ ?
4. Can we increase performance and generality by pre-training the diffusion model on large corpora of molecular data derived at a low level of theory? Can this be transferred to structure-based drug design?

TIMESTEP-DEPENDENT LOSS WEIGHTING

We hypothesize that denoising requires high accuracy close to the data distribution for generating valid molecules, while errors close to the noise distribution are negligible. We propose using the time-dependent weighting:

$$w_s(t) = \max(0.05, \min(1.5, \text{SNR}(t))). \quad (1)$$

	QM9				GEOM-Drugs			
Weighting	Mol. Stability \uparrow	Validity \uparrow	Connect. Comp. \uparrow		Mol. Stability \uparrow	Validity \uparrow	Connect. Comp. \uparrow	
w_u	97.39 \pm 0.23	97.99 \pm 0.20	99.70 \pm 0.03		87.59 \pm 0.19	71.44 \pm 0.22	86.57 \pm 0.33	
$w_s(t)$	98.68 \pm 0.11	98.96 \pm 0.07	99.94 \pm 0.03		91.60 \pm 0.14	84.02 \pm 0.19	95.08 \pm 0.12	

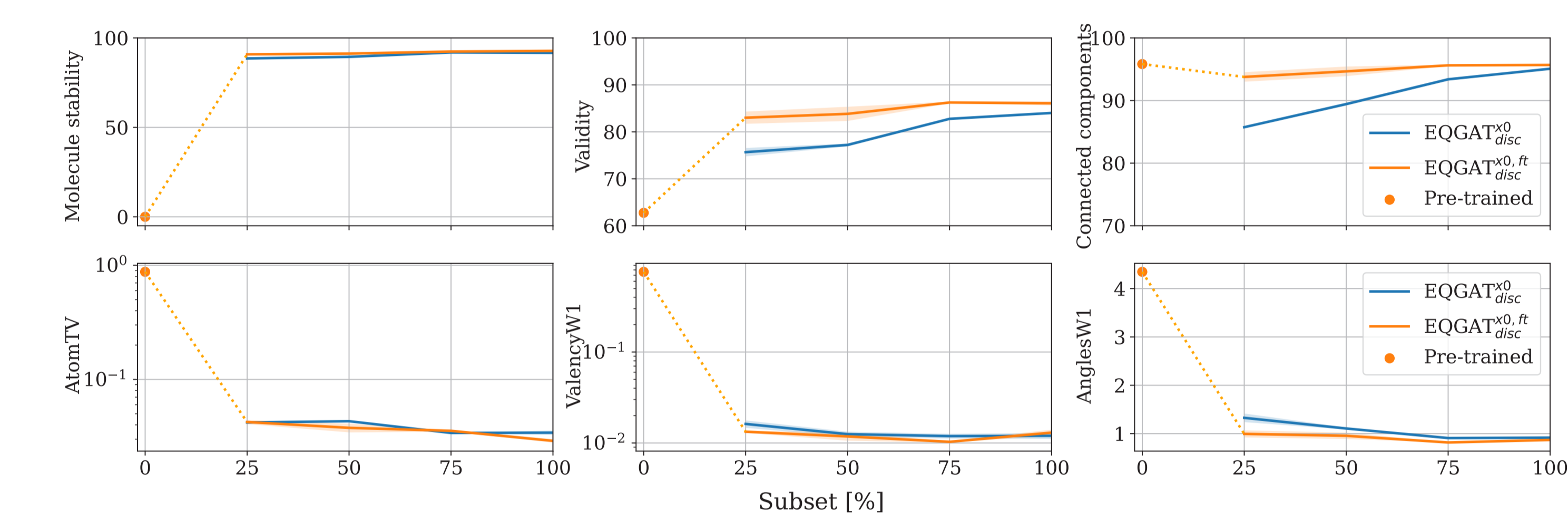
STATE-OF-THE-ART MOLECULE GENERATION

Dataset	GEOM-Drugs					
	EQGAT _{disc} ^{x0}	EQGAT _{disc} ^{x0,ft}	EQGAT _{disc} ^{x0,af}	EQGAT _{disc} ^{x0,af,ft}	EDM	MiDi
Mol. Stab. \uparrow	93.11 \pm 0.31	93.92 \pm 0.13	94.51 \pm 0.18	95.01 \pm 0.37	40.3	89.7 \pm 0.60
Atom. Stab \uparrow	99.79 \pm 0.01	99.81 \pm 0.01	99.83 \pm 0.01	99.84 \pm 0.00	97.8	99.7 \pm 0.01
Validity \uparrow	85.86 \pm 0.33	88.04 \pm 0.17	87.89 \pm 0.31	88.42 \pm 0.26	87.8	70.5 \pm 0.41
Connect. Comp. \uparrow	96.32 \pm 0.25	96.57 \pm 0.18	96.36 \pm 0.25	96.71 \pm 0.20	41.4	88.76 \pm 0.55
Novelty \uparrow	99.82 \pm 0.05	99.84 \pm 0.02	99.82 \pm 0.05	99.82 \pm 0.03	100.00	100.00 \pm 0.00
Diversity \uparrow	89.03 \pm 0.03	89.05 \pm 0.05	88.98 \pm 0.02	88.96 \pm 0.01	-	-
KL Divergence \uparrow	87.66 \pm 0.31	87.58 \pm 0.56	88.38 \pm 0.25	87.62 \pm 0.19	-	-
Train Similarity \downarrow	0.114 \pm 0.0	0.113 \pm 0.0	0.114 \pm 0.0	0.114 \pm 0.0	-	-
AtomsTV [10 ⁻²] \downarrow	3.02 \pm 0.08	3.02 \pm 0.10	2.88 \pm 0.10	2.91 \pm 0.10	21.2	5.11 \pm 0.19
BondsTV [10 ⁻²] \downarrow	2.44 \pm 0.01	2.40 \pm 0.00	2.42 \pm 0.00	2.40 \pm 0.00	4.8	2.44 \pm 0.00
ValencyW ₁ [10 ⁻²] \downarrow	1.18 \pm 0.09	1.20 \pm 0.00	0.85 \pm 0.12	0.90 \pm 0.10	28.5	2.48 \pm 0.52
BondLengthsW ₁ [10 ⁻²] \downarrow	0.56 \pm 0.38	0.10 \pm 0.00	0.50 \pm 0.51	0.20 \pm 0.10	0.2	0.2 \pm 0.10
BondAnglesW ₁ \downarrow	0.83 \pm 0.03	0.79 \pm 0.02	0.65 \pm 0.01	0.62 \pm 0.01	6.23	1.73 \pm 0.32

Experiments

TRANSFERABILITY OF DIFFUSION MODELS

Diffusion models are effective in learning vast data distributions but require large datasets to be trained on. We explore the effect of pre-training on PubChem3d for generalization on the Geom-Drugs dataset. We observe that the fine-tuned model performs better than models trained from scratch, even if the fine-tuning was done on a subset of data.



STRUCTURE-BASED LIGAND DESIGN

Here, we train a conditional diffusion model $p_\theta(x|P)$ where P is a protein pocket. Particularly, a pre-trained model that learned on a vast, yet unconditional chemical space, generalizes better when fine-tuned on CrossDocked2020 than a model trained from scratch.

Model	Validity \uparrow	Connect. Comp. \uparrow	BondLengths W ₁ [10 ⁻²] \downarrow	BondAngles W ₁ \downarrow
EQGAT _{disc} ^{x0} (w_u)	85.51 \pm 0.09	95.15 \pm 0.14	0.20 \pm 0.0	4.37 \pm 0.20
EQGAT _{disc} ^{x0} ($w_s(t)$)	89.62 \pm 0.08	97.65 \pm 0.11	0.12 \pm 0.0	2.12 \pm 0.26
EQGAT _{disc} ^{x0,ft} ($w_s(t)$)	95.65 \pm 0.12	99.66 \pm 0.10	0.11 \pm 0.0	1.55 \pm 0.21

Model	Vina (All) \downarrow	Vina (Top-10%) \downarrow	QED \uparrow	SA \uparrow	Lipinski \uparrow	Diversity \uparrow
EQGAT _{disc} ^{x0,ft} ($w_s(t)$)	-7.423 \pm 2.33	-9.571 \pm 2.14	0.522 \pm 0.18	0.697 \pm 0.20	4.66 \pm 0.72	0.742 \pm 0.07
TargetDiff	-7.318 \pm 2.47	-9.669 \pm 2.55	0.483 \pm 0.20	0.584 \pm 0.13	4.594 \pm 0.83	0.718 \pm 0.09
DiffSBDD-cond	-6.950 \pm 2.06	-9.120 \pm 2.16	0.469 \pm 0.21	0.578 \pm 0.13	4.562 \pm 0.89	0.728 \pm 0.07

