

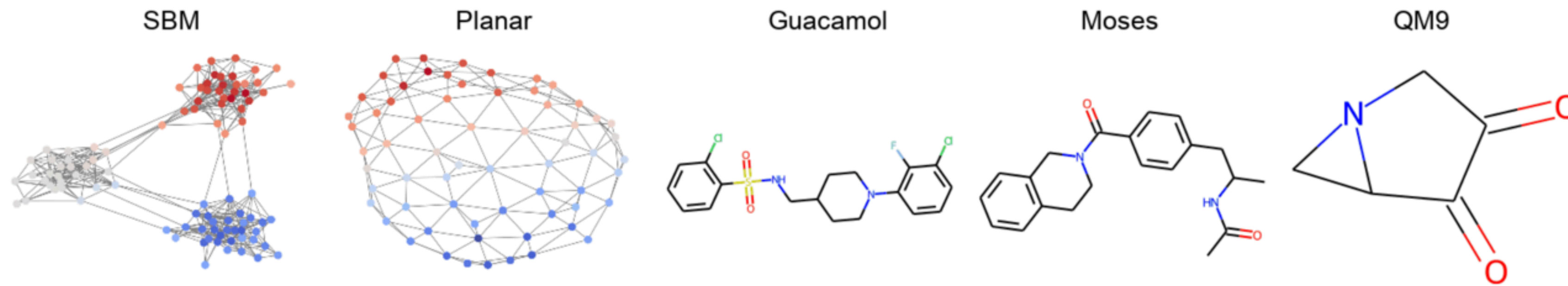
# DiGress: Discrete denoising diffusion for graph generation

Clément Vignac

Joint work with Igor Krawczuk (co-1st author), Antoine Siraudin, Bohan Wang, Volkan Cevher & Pascal Frossard (EPFL)

## GOAL

Generate graphs with categorical node and edge attributes that *look like* graphs in a training set



# Outline

1. Key properties of diffusion models
2. DiGress: a discrete diffusion model for graphs
3. Improvements over the vanilla model
4. Results
5. What's next?

# I Denoising diffusion models

## 01 DENOISING DIFFUSION MODELS



DALLE-2

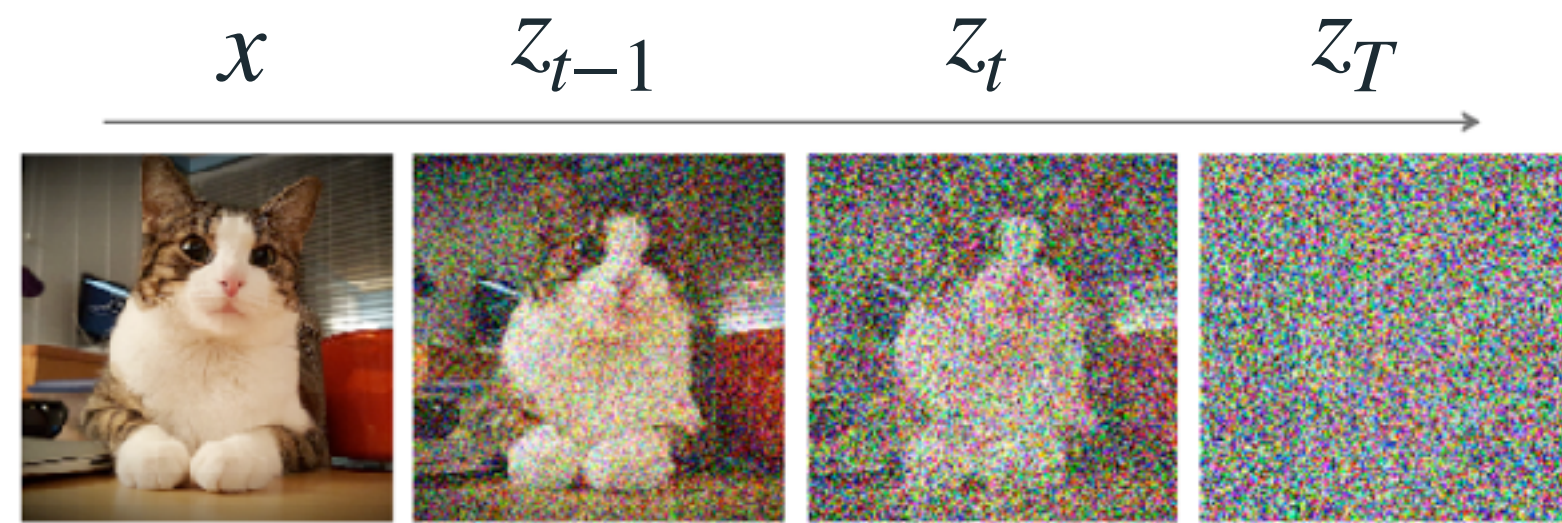


Stable diffusion

+ 3d point clouds, text, audio...

2 equivalent perspectives: score-matching and denoising diffusion

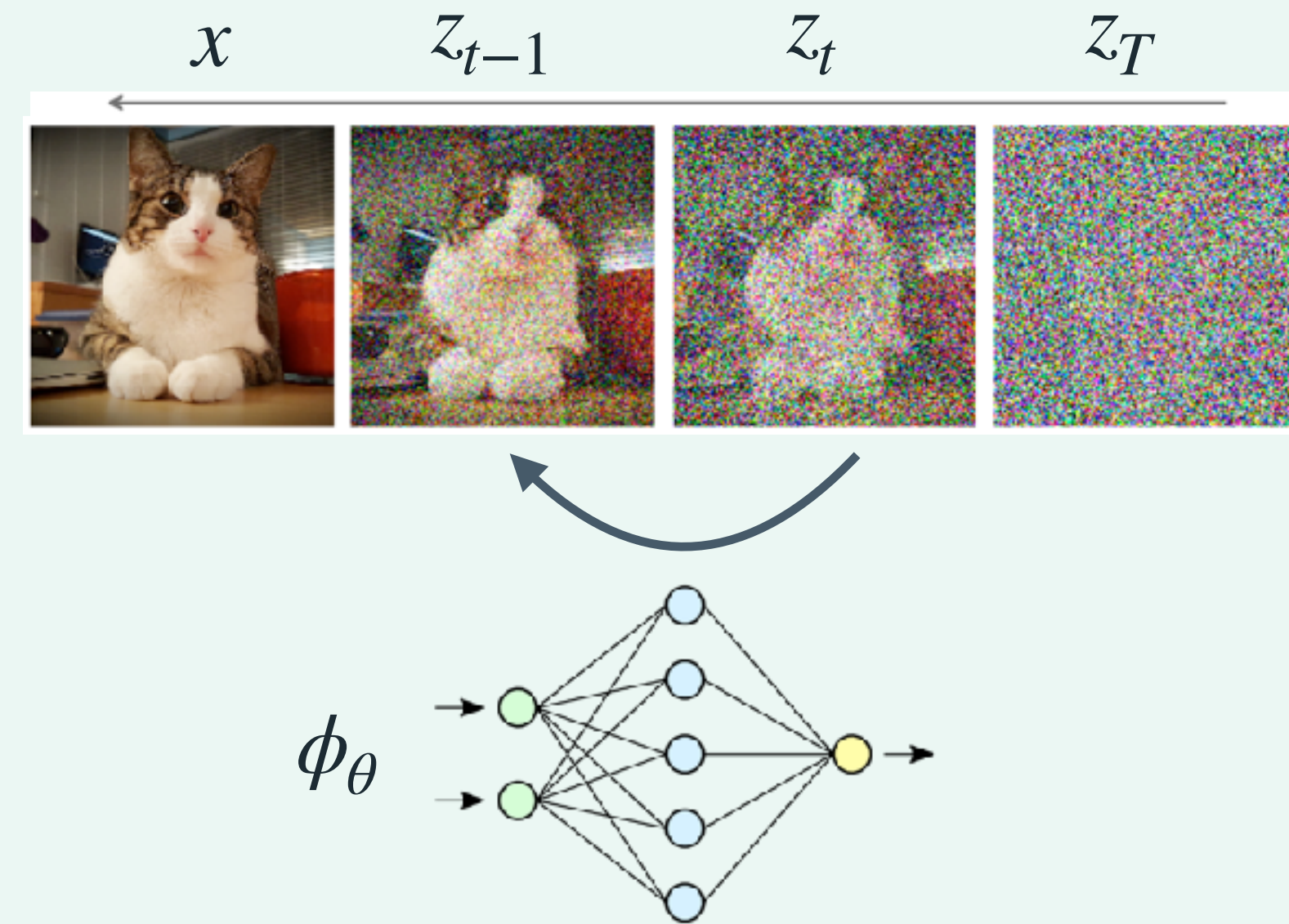
## 01 DENOISING DIFFUSION MODELS: THE NAIVE MODEL



$$q(z_1, \dots, z_T | x) = q(z_1 | x) \prod_{t=2}^T q(z_t | z_{t-1})$$

### Noise model

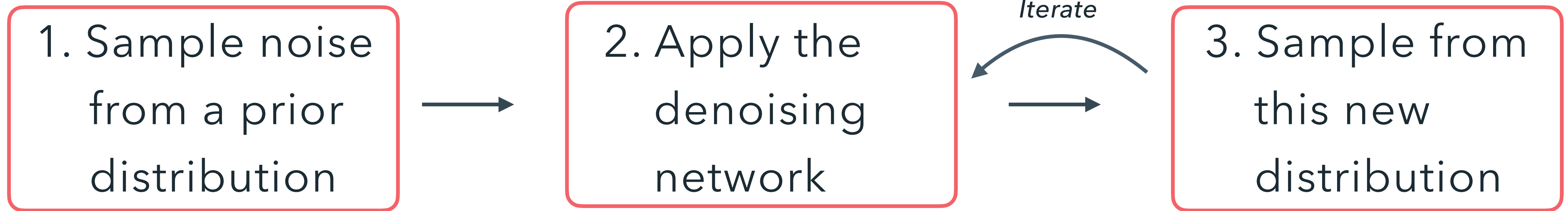
Generate **diffusion trajectories** by recursively adding noise to a data point.



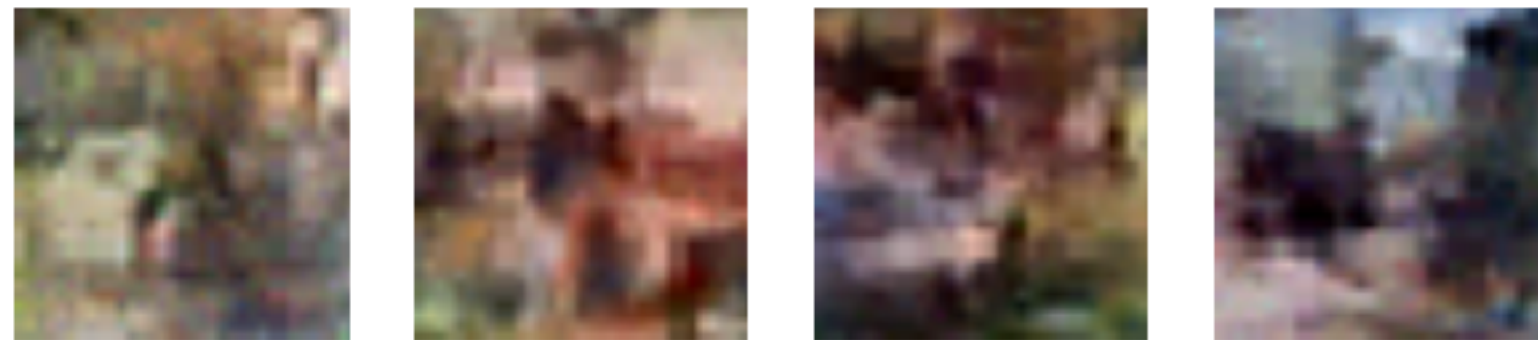
### Denoising network

Train a network to predict the **inverse** diffusion iterations.

## 01 DENOISING DIFFUSION MODELS: SAMPLING



to predict the distribution  
**at previous step.**



(Sohl-Dickstein et al. 15)



LDSM (Ho et al. 20)

What is missing?

---

**Algorithm** : Training

---

**Input:** Data point  $x$

Sample  $t \sim \mathcal{U}(1, \dots, T)$

Sample  $z_t \sim \text{noise}(x, t)$

$\hat{p}_x \leftarrow f_\theta(z_t, t)$

optimizer.step(loss( $x, \hat{p}_x$ ))

---

---

**Algorithm** : Sampling

---

Sample  $z_t$  from a prior distribution

**for**  $t = T$  **to**  $1$  **do**

$\hat{p}_x \leftarrow f_\theta(z_t, t)$

$\hat{p}_{t-1} \leftarrow \int_x q(z_{t-1} | z_t, x) d\hat{p}_x$

$z_{t-1} \sim \hat{p}_{t-1}$

**end**

**return**  $z_0$

---

### *Differences with the naive model*

- The model is not trained on trajectories
- $\hat{p}_{t-1}$  is deduced from  $\hat{p}_x$  or  $\hat{p}_\epsilon$

## 01 DENOISING DIFFUSION MODELS: MODERN MODELS

---

### Algorithm : Training

---

**Input:** Data point  $x$

Sample  $t \sim \mathcal{U}(1, \dots, T)$

Sample  $z_t \sim \text{noise}(x, t)$

$\hat{p}_x \leftarrow f_\theta(z_t, t)$

optimizer.step(loss( $x, \hat{p}_x$ ))

---

- $z^{t-1}$  (= the target of early diffusion models) is **very noisy**
- There are two sources of noise:
  - The data sample  $x \sim p_{\text{data}}$
  - The trajectory sampled from  $x$
- If we can predict  $x$  directly we remove some stochasticity

### *Differences with the naive model*

- The model is not trained on trajectories
- $\hat{p}_{t-1}$  is deduced from  $\hat{p}_x$  or  $\hat{p}_\epsilon$

## 01 DENOISING DIFFUSION MODELS: THE PRICE OF EFFICIENCY

### Property 1

$q(z^t | x)$  should have a closed-form formula

Avoid adding noise recursively during training to generate trajectories.

### Property 2

$\int q(z^{t-1} | z^t, x) d\hat{p}_x$   
should be tractable

Use  $x$  as the target of the neural network instead of  $z^{t-1}$ .

### Property 3

$\lim_{T \rightarrow \infty} q(z^T | x)$  should not depend on  $x$

Prior distribution =  
Limit noise distribution

2 main frameworks: Gaussian noise (95% of the papers) and discrete diffusion

## 01 DENOISING DIFFUSION MODELS: DISCRETE DIFFUSION

- Noise model = Markov process over a discrete state-space

Defined by transition matrices  $(Q^1, \dots, Q^T)$ :

$$q(z^t | z^{t-1}) = z^{t-1} Q^t$$

- Multiple steps:  $q(z^t | x) = x \bar{Q}^t$  for  $\bar{Q}^t = Q^1 \dots Q^t$

Property 1

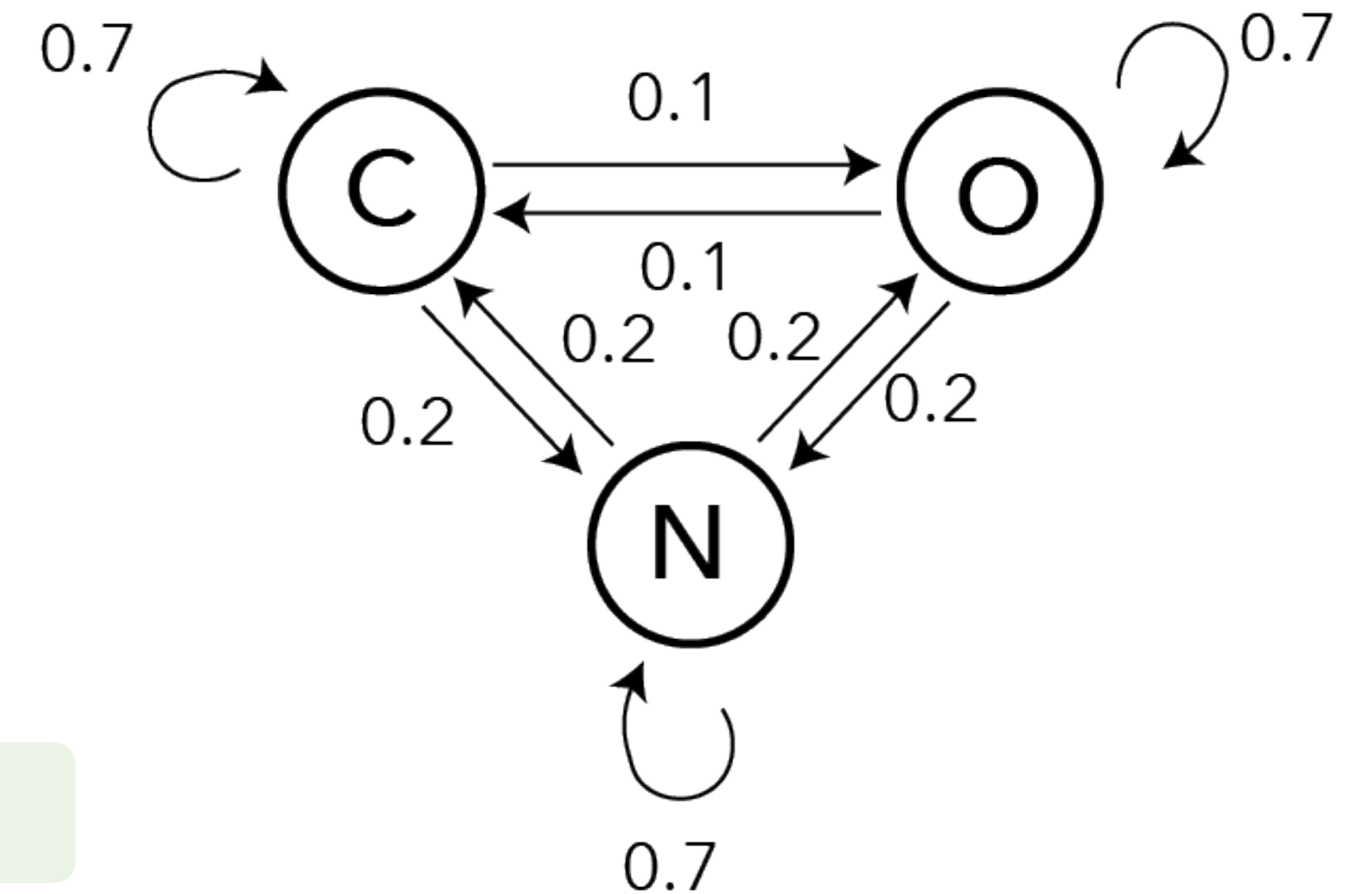
- Bayes rule:  $q(z^{t-1} | z^t, x) \propto z^t (Q^t)' \odot x \bar{Q}^{t-1}$

Property 2

- Most common model: Uniform transitions

Property 3

$$\begin{bmatrix} 0.8 & 0.1 & 0.1 \\ 0.1 & 0.8 & 0.1 \\ 0.1 & 0.1 & 0.8 \end{bmatrix} Q^t = \alpha^t I + \beta^t \mathbf{1}_d \mathbf{1}'_d / d \quad \lim_{t \rightarrow \infty} Q^t x = \mathbf{1}_d / d$$



Applications:

Text, image, audio...

# II DiGress: Discrete Graph Denoising Diffusion

# Challenges of learning on graphs

- Varying graph sizes
- Limited representation power

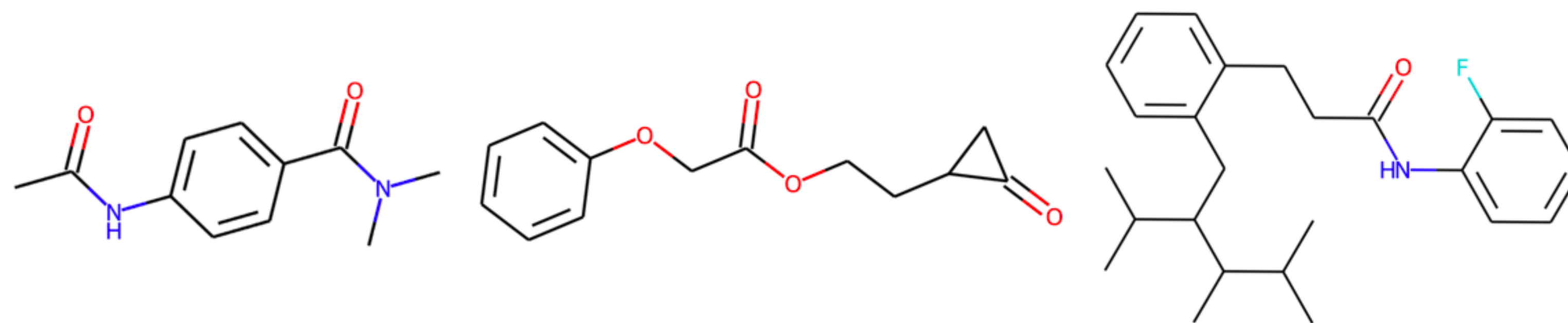
# Challenges of graph generation

- Unsupervised learning task that might require graph matching
- Sparsity
- Evaluation

## 02 DIGRESS: EXISTING DIFFUSION MODELS FOR GRAPHS

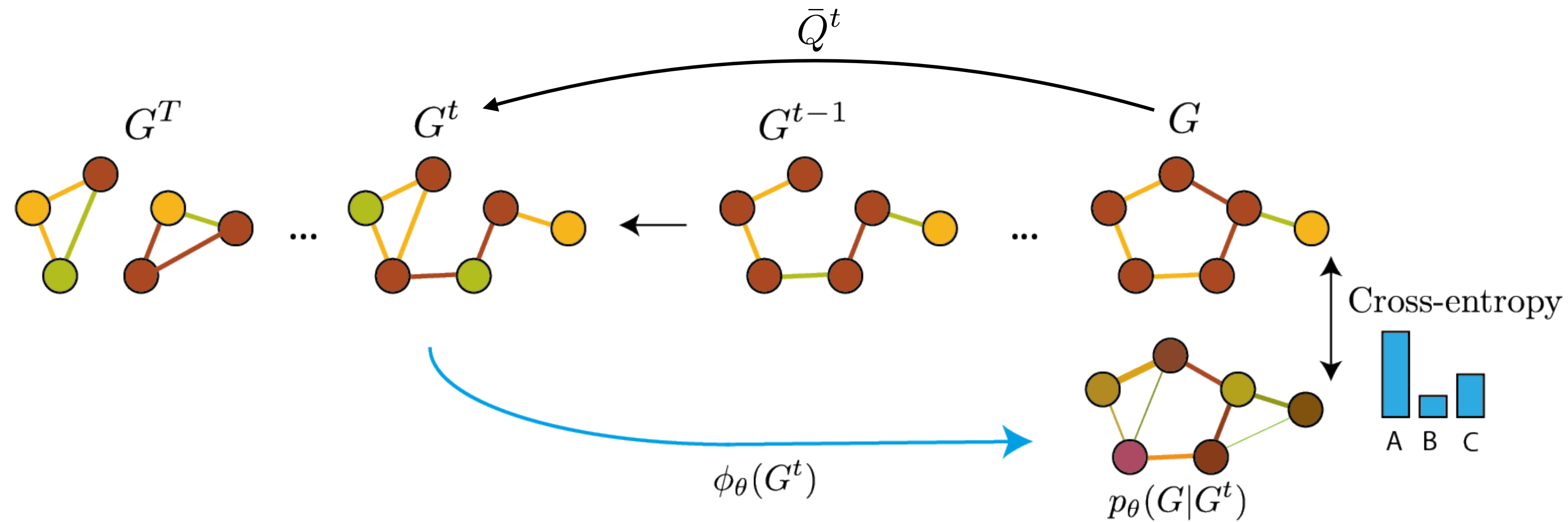
### Score-based Generative Modeling of Graphs via the System of SDEs

Jo et al. 2022



- Add **Gaussian noise** to a graph adjacency matrix, and the node and edge features
- Use a graph attention network to denoise this process
- **Scales** significantly better than previous work
- Does not respect the **discrete nature** of graphs
- The noisy graphs are **complete**: structural information (connectivity, cycles) is not defined

## 02 DIGRESS: METHOD (TRAINING)



Corrupt *independently* on each node and edge

Absence of edge = a particular edge type

$\mathcal{X}$ : space of node types

$\mathcal{E}$ : space of edge types

---

### Algorithm : Training DiGress

---

**Input:** A graph  $G = (\mathbf{X}, \mathbf{E})$

Sample  $t \sim \mathcal{U}(1, \dots, T)$

Sample  $G^t \sim \mathbf{X}\bar{Q}_X^t \times \mathbf{E}\bar{Q}_E^t$

$\hat{p}^X, \hat{p}^E \leftarrow \phi_\theta(G_t, t)$

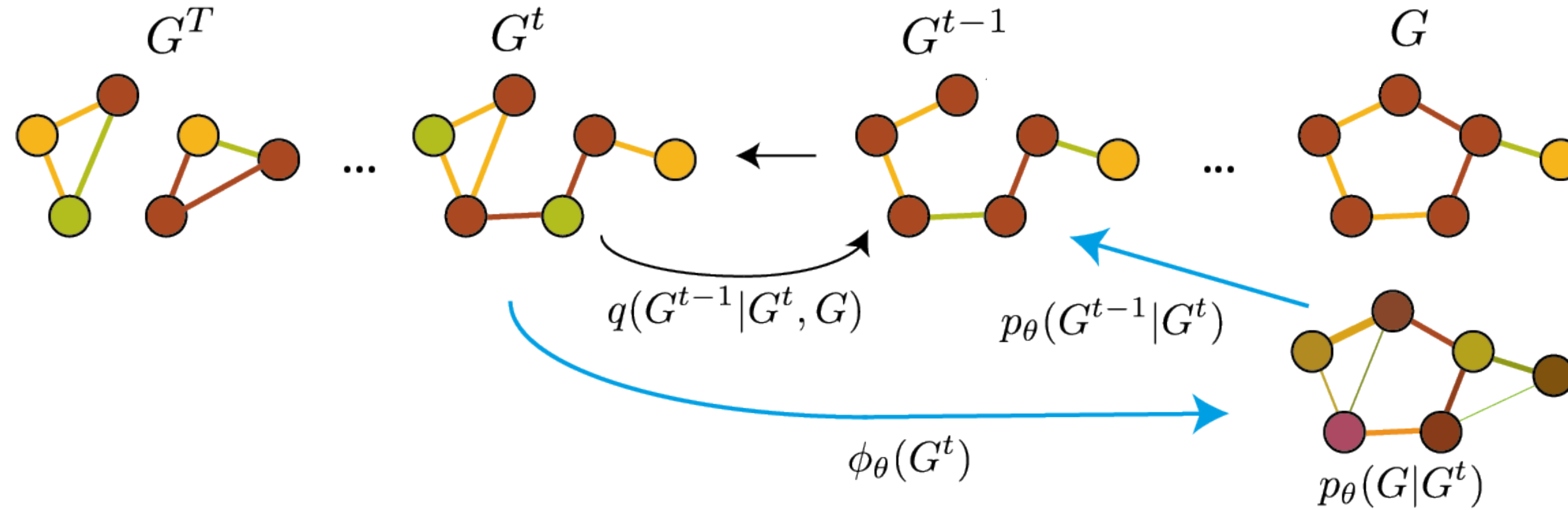
$loss \leftarrow l_{CE}(\hat{p}^X, \mathbf{X}) + \lambda l_{CE}(\hat{p}^E, \mathbf{E})$

optimizer.step(loss)

---

Graph generation = sequence of node and edge classification tasks

## 02 DIGRESS: METHOD (SAMPLING)



Model:

$$p_{\theta}(G^{t-1} | G^t) = \prod_i p_{\theta}(x_i^{t-1} | G^t) \prod_{i,j} p_{\theta}(e_{ij}^{t-1} | G^t)$$

$$p_{\theta}(x_i^{t-1} | G^t) = \int_{x_i} p_{\theta}(x_i^{t-1} | x_i, G^t) dp_{\theta}(x_i | G^t) = \sum_{x \in \mathcal{X}} q(x_i^{t-1} | x_i = x, x_i^t) \hat{p}_i^X(x)$$

---

### Algorithm : Sampling from DiGress

---

Sample  $n$  from the training data distribution

Sample  $G^T \sim q_X(n) \times q_E(n)$

**for**  $t = T$  **to**  $1$  **do**

$\hat{p}^X, \hat{p}^E \leftarrow \phi_{\theta}(G^t, t)$

    Sample  $G^{t-1} \sim \prod_i p_{\theta}(x_i^{t-1} | G^t) \times \prod_{i,j} p_{\theta}(e_{ij}^{t-1} | G^t)$

**end**

**return**  $G^0$

---

Graph generation = sequence of node and edge classification tasks




## 02 DIGRESS: PROPERTIES

DiGress is  
permutation  
equivariant

It is made of permutation  
equivariant blocks

Its loss is  
permutation  
invariant

$$l(\hat{p}^G, G) = \sum_i \text{CE}(x_i, \hat{p}_i^X) + \sum_{ij} \text{CE}(e_{ij}, \hat{p}_{ij}^E)$$


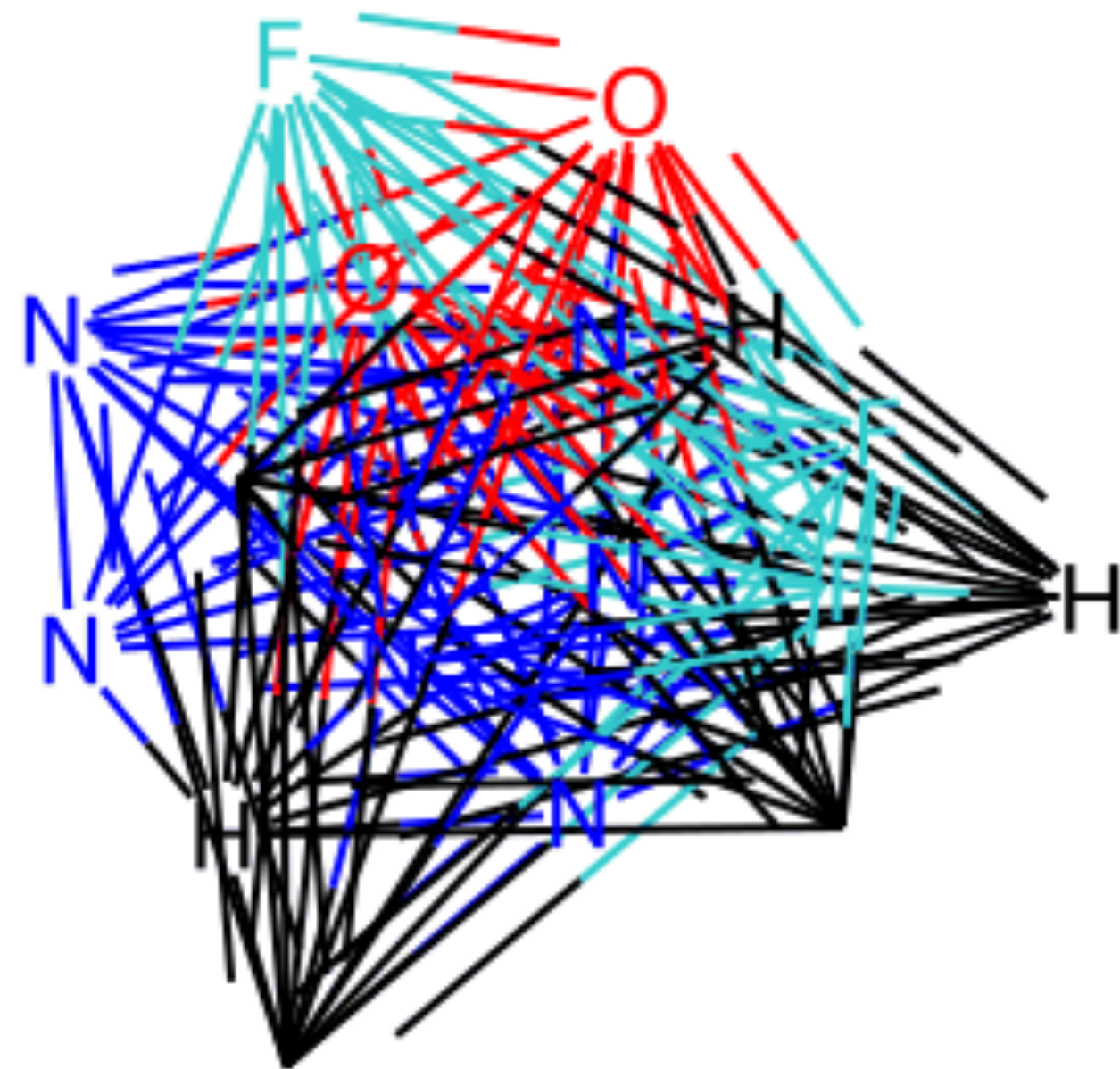
Do not depend on  $i, j$

The model is  
exchangeable

**All permutations** of a generated  
graph **are equally likely**

Needed for likelihood computation

# First results



Frame 0

## Have we addressed the specificities of learning on graphs?

- Varying sizes 👍
- Equivariance properties 👍
- Sparsity 🤔
- Limited expressive power 🤔

III Improving DiGress with marginal transitions and structural features

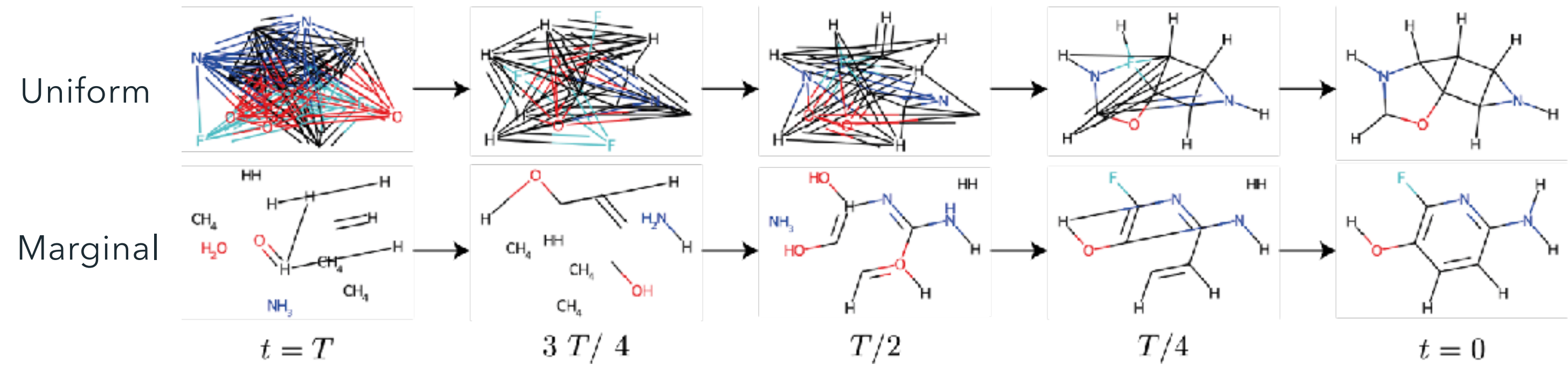
### 03 IMPROVING DIGRESS

# I. Choice of the noise model

Rather than **uniform transitions** over node and edge classes, choose a noise model that **preserves the marginal distribution of node and edge types**.

$$Q_X^t = \alpha^t I + \beta^t 1_a m'_X \quad \begin{bmatrix} 0.8 & 0.1 & 0.1 \\ 0.2 & 0.7 & 0.1 \\ 0.2 & 0.1 & 0.7 \end{bmatrix}$$

$$Q_E^t = \alpha^t I + \beta^t 1_b m'_E$$



# II. Structural and spectral features

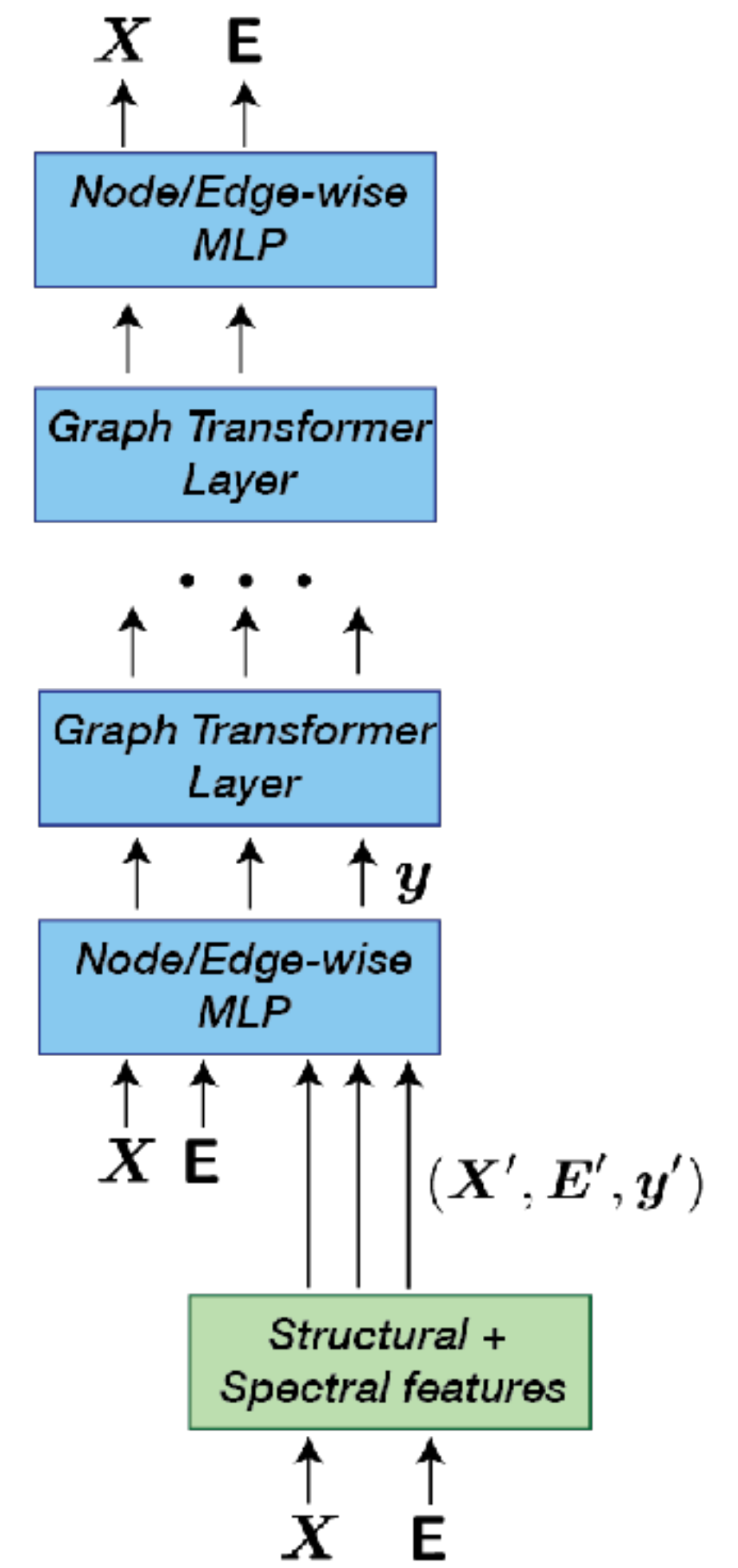
Message-passing neural networks have a **limited representation power** and cannot **detect substructures** such as cycles.

Solution 1: use a more expressive network

Very costly in practice

Solution 2: add features that the network cannot compute by itself

- Cycle counts
- Spectral features (# 0 eigenvalues, first non zero eigenvalues, first eigenvectors)



Xu et al. 18, How powerful are graph neural networks?

Bouritsas et al. 22, Improving graph neural network expressivity via subgraph isomorphism counting

Beaini et al. 21, Directional graph networks

# Algorithms

---

**Algorithm 1: Training DiGress**


---

**Input:** A graph  $G = (\mathbf{X}, \mathbf{E})$

Sample  $t \sim \mathcal{U}(1, \dots, T)$

Sample  $G^t \sim \mathbf{X} \bar{\mathbf{Q}}_X^t \times \mathbf{E} \bar{\mathbf{Q}}_E^t$

▷ Add noise

$z \leftarrow f(G^t, t)$

▷ Structural and spectral features

$\hat{p}^X, \hat{p}^E \leftarrow \phi_\theta(G^t, z)$

▷ Forward pass

$loss \leftarrow l_{CE}(\hat{p}^X, \mathbf{X}) + \lambda l_{CE}(\hat{p}^E, \mathbf{E})$

optimizer.step(loss)

---

**Algorithm 2: Sampling from DiGress**


---

Sample  $n$  from the training data distribution

Sample  $G^T \sim q_X(n) \times q_E(n)$

▷ Random graph

**for**  $t = T$  **to**  $1$  **do**

$z \leftarrow f(G^t, t)$

    ▷ Structural and spectral features

$\hat{p}^X, \hat{p}^E \leftarrow \phi_\theta(G^t, z)$

    ▷ Forward pass

    Sample  $G^{t-1} \sim \prod_i p_\theta(x_i^{t-1} | G^t) \times \prod_{ij} p_\theta(e_{ij}^{t-1} | G^t)$

    ▷ Reverse process

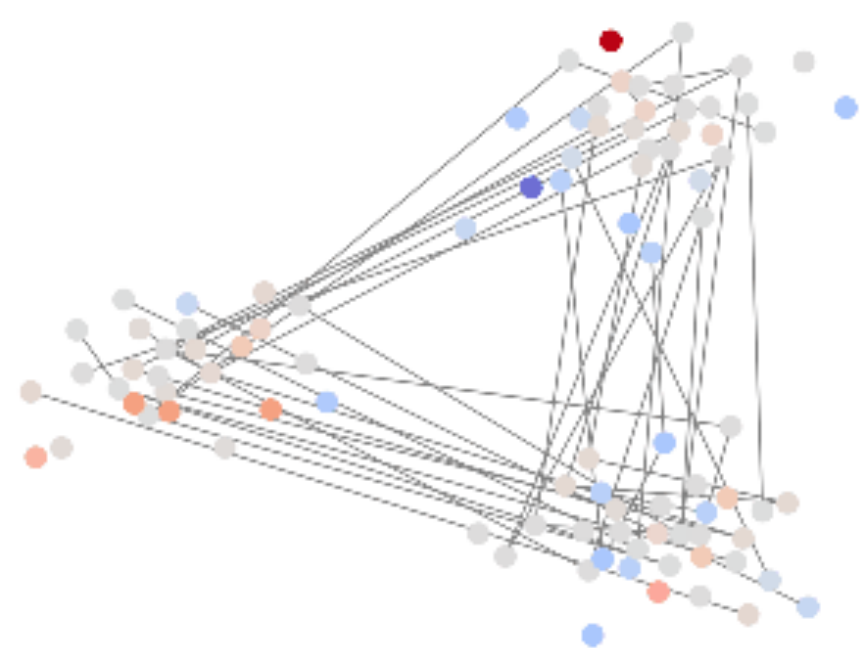
**end**

**return**  $G^0$

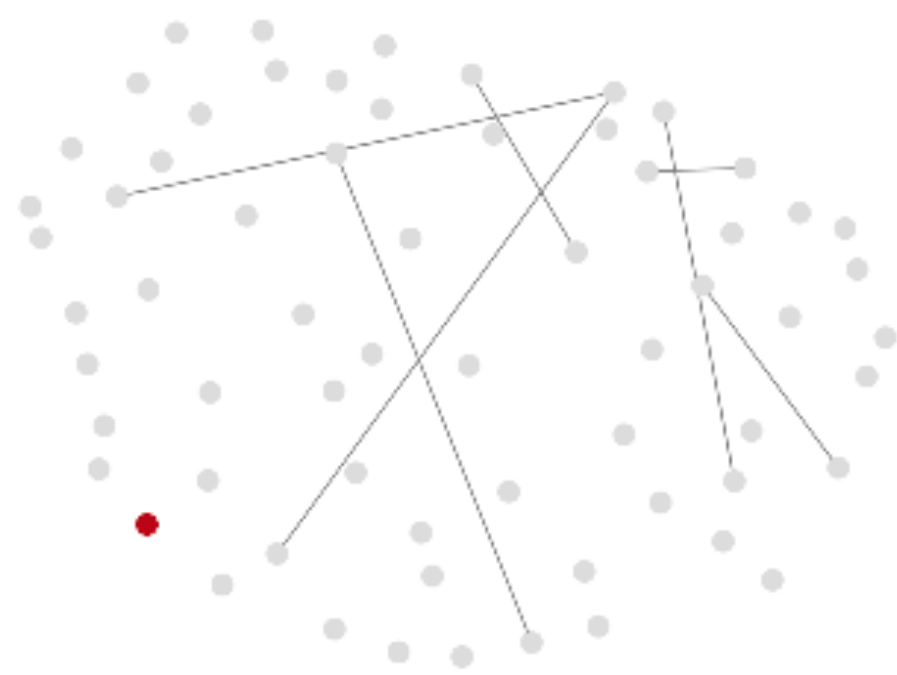
---

# Results

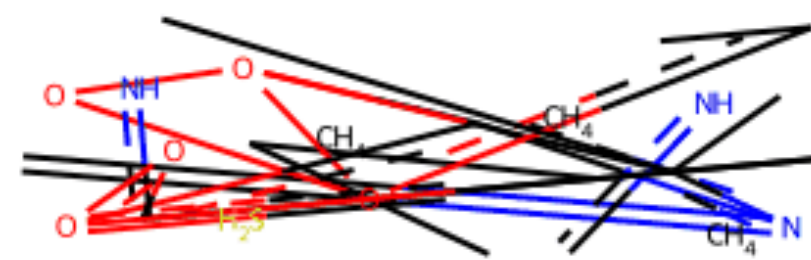
SBM



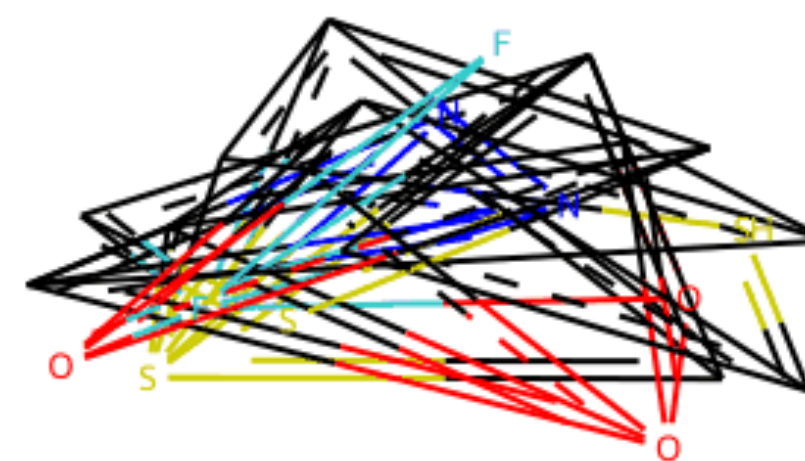
Planar



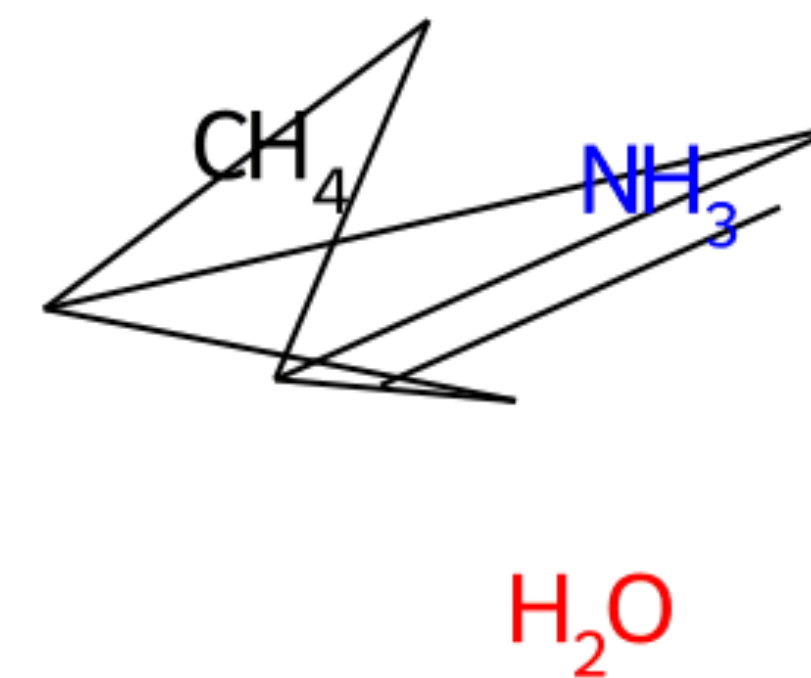
Guacamol



Moses



QM9



## 04 RESULTS

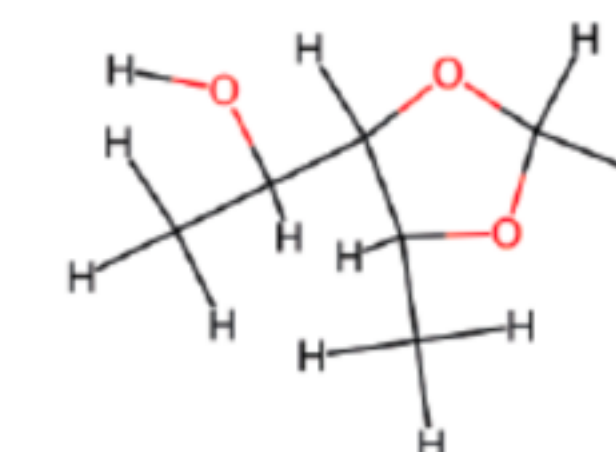
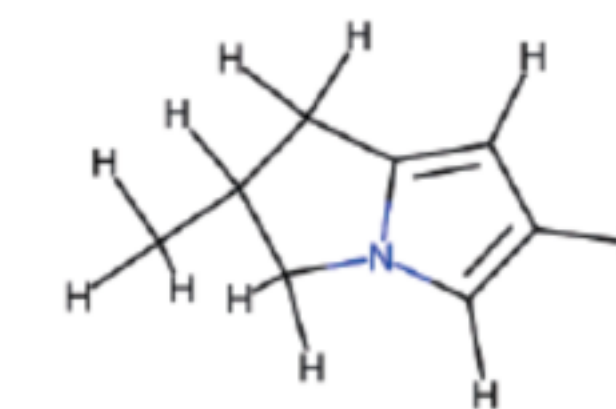
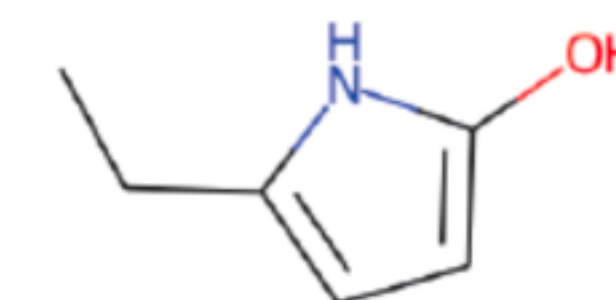
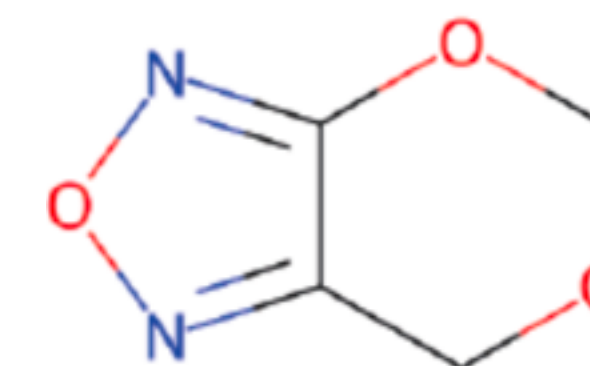
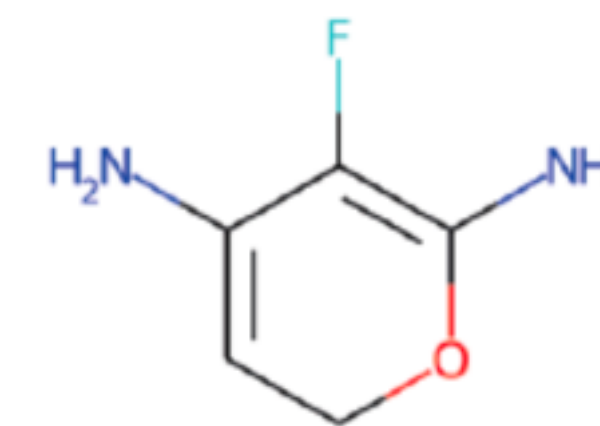
# Qm9

### Implicit hydrogens

Method	NLL	Valid	Unique	Training time (h)
Dataset	–	99.3	100	–
Set2GraphVAE	–	59.9	93.8	–
SPECTRE	–	87.3	35.7	–
GraphNVP	–	83.1	99.2	–
GDSS	–	95.7	<b>98.5</b>	–
ConGress (ours)	–	98.9 $\pm$ 1.1	96.8 $\pm$ 2.2	7.2
DiGress (ours)	23.2 $\pm$ 0	<b>99.0<math>\pm</math>1.1</b>	96.2 $\pm$ 1.1	<b>1.0</b>

### Explicit hydrogens

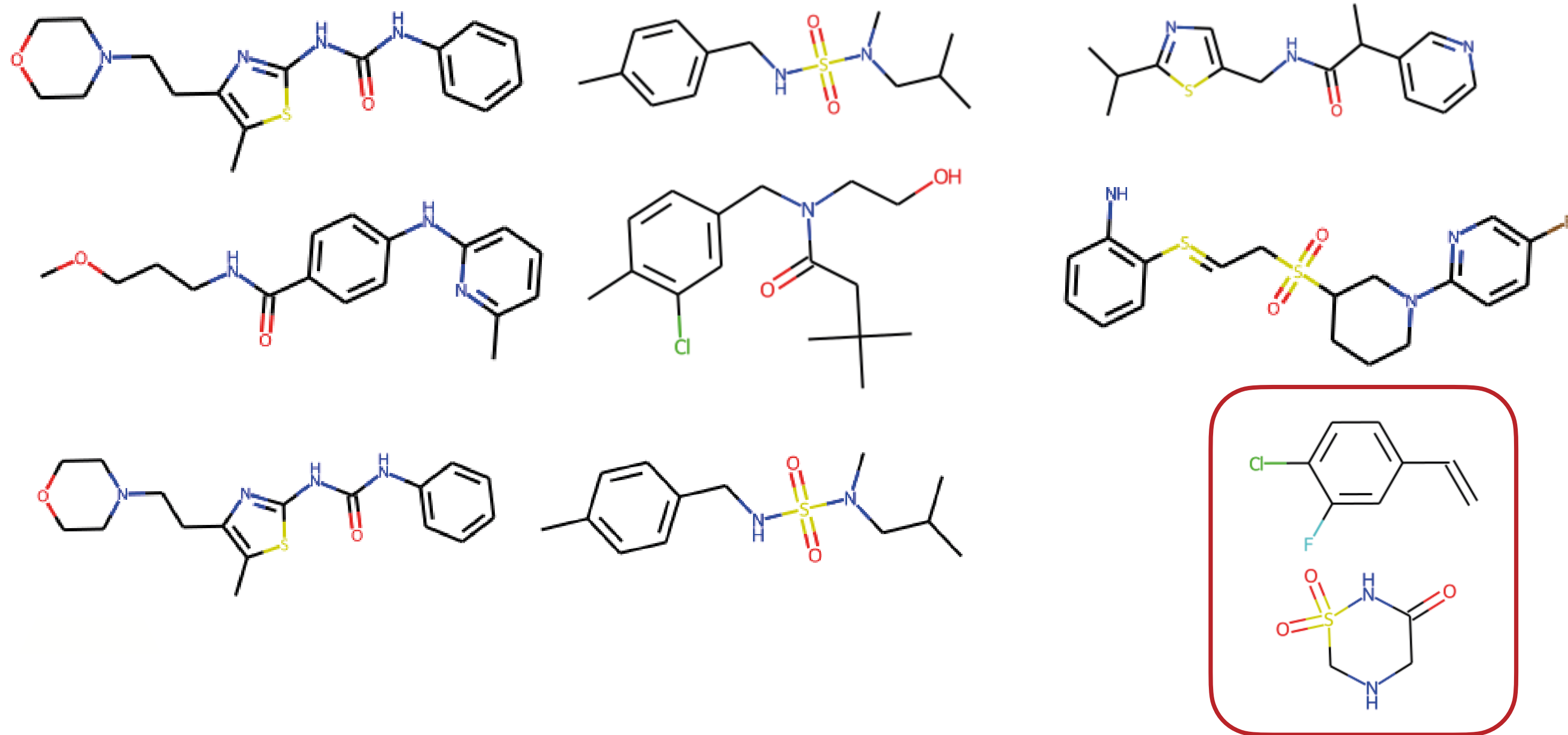
Model	Valid $\uparrow$	Unique $\uparrow$	Atom stable $\uparrow$	Mol stable $\uparrow$
Dataset	97.8	100	98.5	87.0
ConGress	86.7 $\pm$ 1.8	<b>98.4<math>\pm</math>0.1</b>	97.2 $\pm$ 0.2	69.5 $\pm$ 1.6
DiGress (uniform)	89.8 $\pm$ 1.2	97.8 $\pm$ 0.2	97.3 $\pm$ 0.1	70.5 $\pm$ 2.1
DiGress (marginal)	92.3 $\pm$ 2.5	97.9 $\pm$ 0.2	<b>97.3<math>\pm</math>0.8</b>	66.8 $\pm$ 11.8
DiGress (marg. + additional features)	<b>95.4<math>\pm</math>1.1</b>	97.6 $\pm$ 0.4	<b>98.1<math>\pm</math>0.3</b>	<b>79.8<math>\pm</math>5.6</b>



## 04 RESULTS

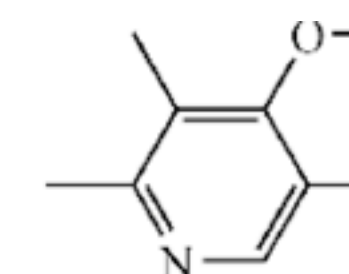
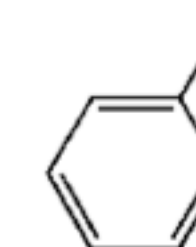
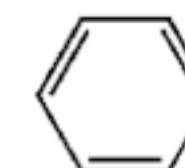
# MOSES

- 1.9M molecules filtered from Zinc Clean Leads
- Results are reported on a test set of separate scaffolds



Model	Class	Val ↑	Unique ↑	Novel ↑	Filters ↑	FCD ↓	SNN ↑	Scaf ↑
VAE	SMILES	97.7	99.8	69.5	99.7	0.57	0.58	5.9
JT-VAE	Fragment	100	100	99.9	97.8	1.00	0.53	10
GraphINVENT	Autoreg.	96.4	99.8	—	95.0	1.22	0.54	12.7
ConGress (ours)	One-shot	83.4	99.9	96.4	94.8	1.48	0.50	16.4
DiGress (ours)	One-shot	85.7	100	95.0	97.1	1.19	0.52	14.8

Example of SMILES string:  
Cn1cnc2n(C)c(=O)n(C)c(=O)c12

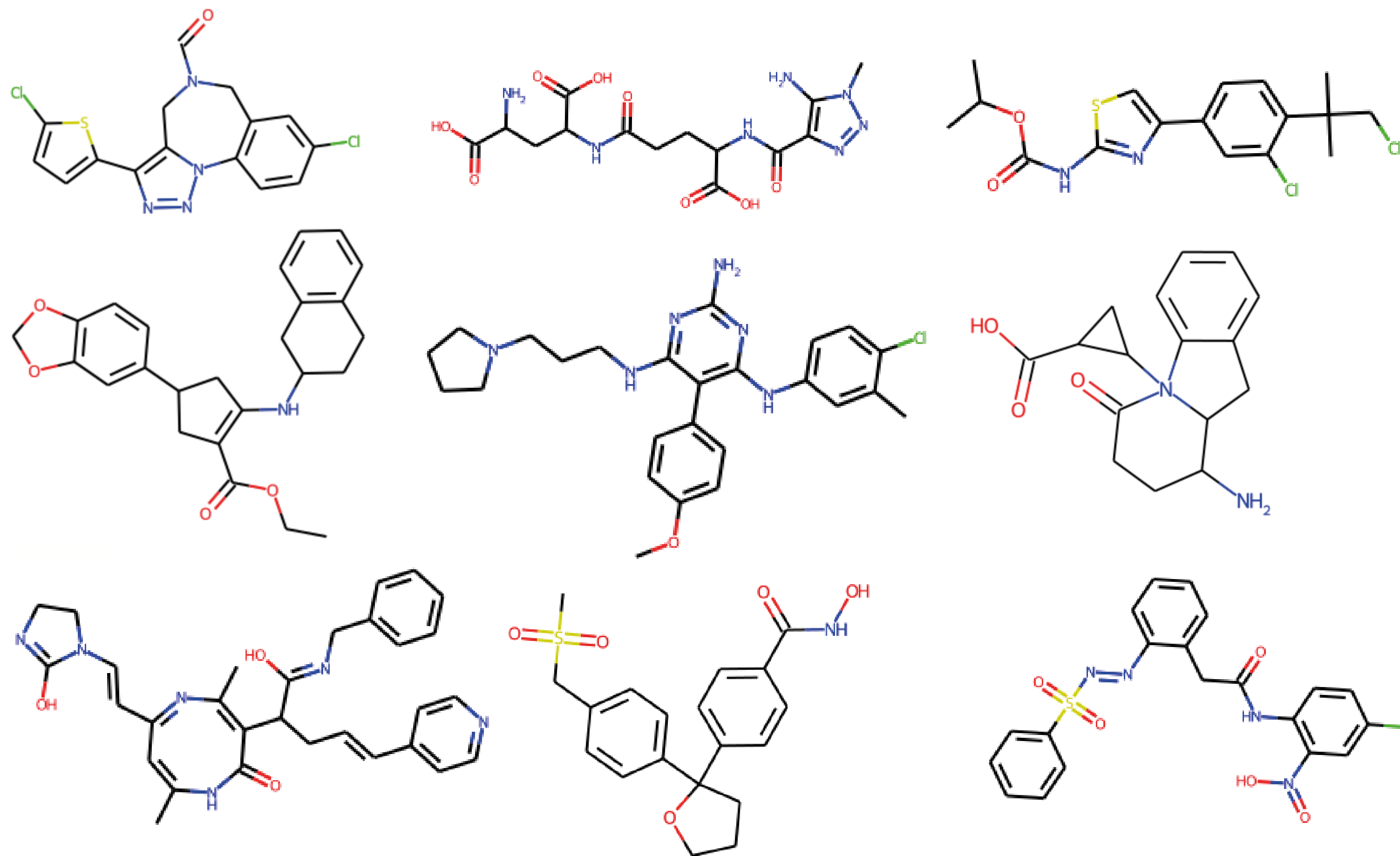


Examples of fragments:

## 04 RESULTS

# GuacaMol

- 1.3M molecules from ChEMBL 24
- Our model currently does not handle partial charges, which are present in the dataset



Model	Class	Valid $\uparrow$	Unique $\uparrow$	Novel $\uparrow$	KL div $\uparrow$	FCD $\uparrow$
LSTM	Smiles	95.9	100	91.2	99.1	91.3
NAGVAE	One-shot	92.9	95.5	100	38.4	0.9
MCTS	One-shot	100	100	95.4	82.2	1.5
ConGress (ours)	One-shot	0.1	100	100	36.1	0.0
DiGress (ours)	One-shot	73.9	99.9	99.9	89.8	53.1

## 04 RESULTS

# Abstract graphs

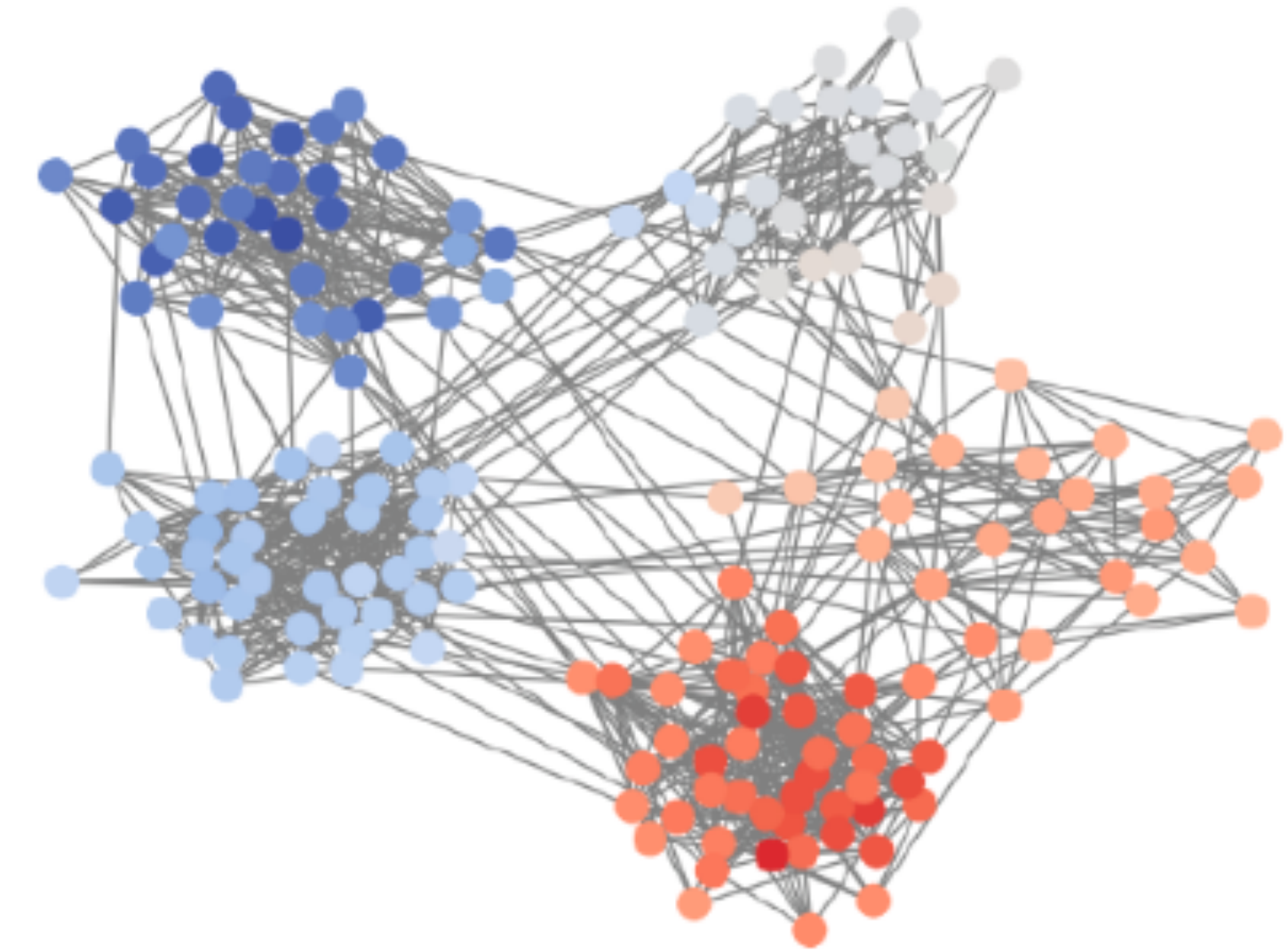
V.U.N: Valid, unique and novel graphs

Model	Deg ↓	Clus ↓	Orb ↓	V.U.N. ↑
<i>Stochastic block model</i>				
GraphRNN	6.9	1.7	3.1	5%
GRAN	14.1	1.7	2.1	25%
GG-GAN	4.4	2.1	2.3	25%
SPECTRE	1.9	1.6	1.6	53%
ConGress	34.1	3.1	4.5	0%
DiGress	<b>1.6</b>	<b>1.5</b>	1.7	<b>74%</b>
<i>Planar graphs</i>				
GraphRNN	24.5	9.0	2508	0%
GRAN	3.5	1.4	1.8	0%
SPECTRE	2.5	2.5	2.4	25%
ConGress	23.8	8.8	2590	0%
DiGress	<b>1.4</b>	<b>1.2</b>	<b>1.7</b>	<b>75%</b>

Could it be a graph drawn from SBM?

Is it a planar graph?

To measure generation quality, compare statistics from the generated graphs with statistics from a test set



# Summary

- DiGress solves graph generation as a sequence of node and edge classification task
- Diffusion models for graphs significantly outperform previous methods – opens the way to many applications
- Discrete diffusion helps scaling to large graphs
- Limitation:  $O(n^2)$  complexity

V What's next?

## 5: WHAT'S NEXT? CONDITIONAL GENERATION

# Discrete regression guidance (proof-of-concept)

**Goal:** condition the generation on a vector  $y$  representing **graph-level properties**

- Train an unconditional model  $\phi_\theta$
- Train a regressor  $g_\eta$  to predict  $y$  from  $G^t$
- At each sampling step, use the gradients of  $g_\eta$  to push the generation towards graph that have property  $y$

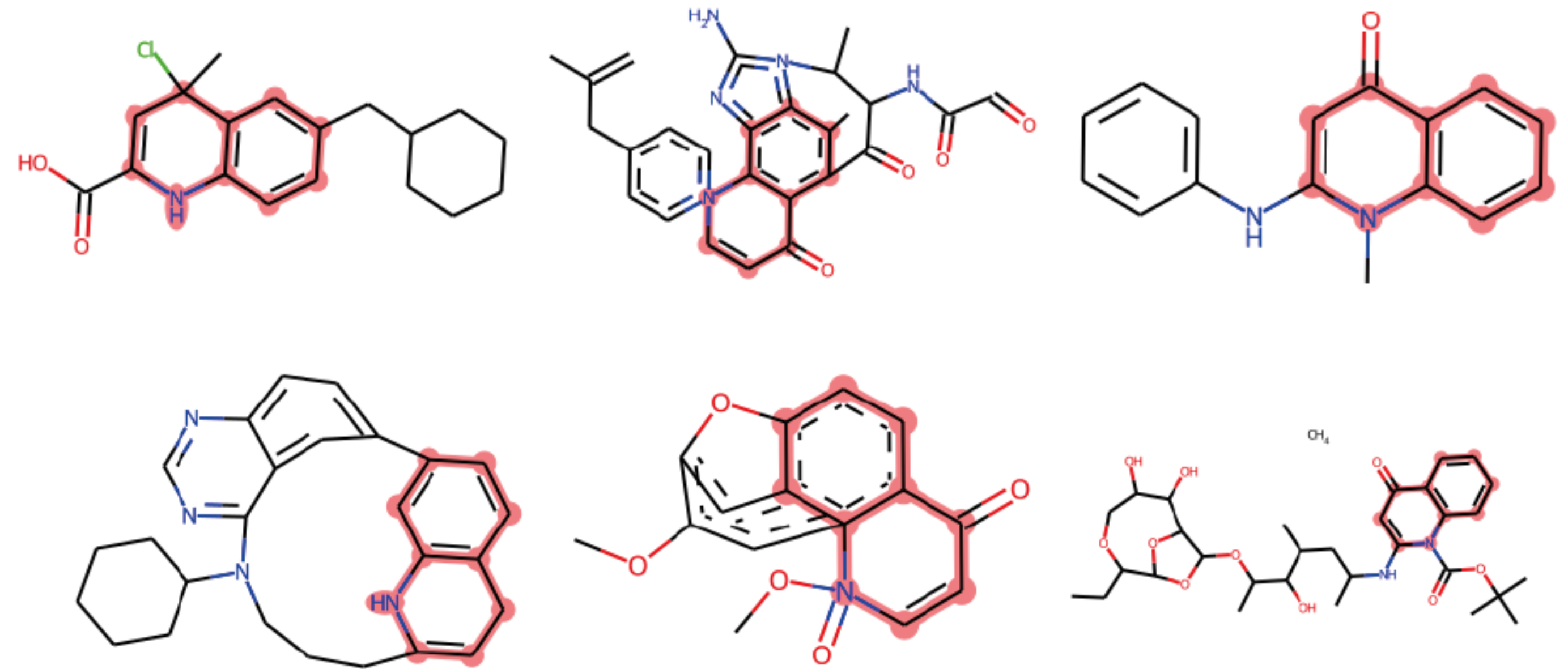
*MAE between target and obtained property*

Target	$\mu$	HOMO	$\mu$ & HOMO
Uncondit.	1.71 $\pm$ .04	0.93 $\pm$ .01	1.34 $\pm$ .01
Guidance	0.81 $\pm$ .04	0.56 $\pm$ .01	0.87 $\pm$ .03

Guidance helps, but there is room for improvement

## 5: WHAT'S NEXT?

# Subgraph Extension



**Long-scale** consistency issues

A resampling strategy similar to **RePaint** could be used

Also...

- Mixed continuous-discrete generation
- Better tools for conditional generation, property optimization, fine-tuning...
- Is it possible to keep permutation equivariance without the  $O(n^2)$  complexity?

# Questions?

Clément Vignac

[clement.vignac@epfl.ch](mailto:clement.vignac@epfl.ch)



**EPFL**

# Molecule generation in 3D

## **Unconditional generation in 3d**

Example: EDM

Generate the 3d coordinates of atoms

Bonds are deduced from the atomic distances and the atom types  
(with a lookup table) -> not learnt

## **Conformer generation**

Examples: Geodiff, Torsional diffusion

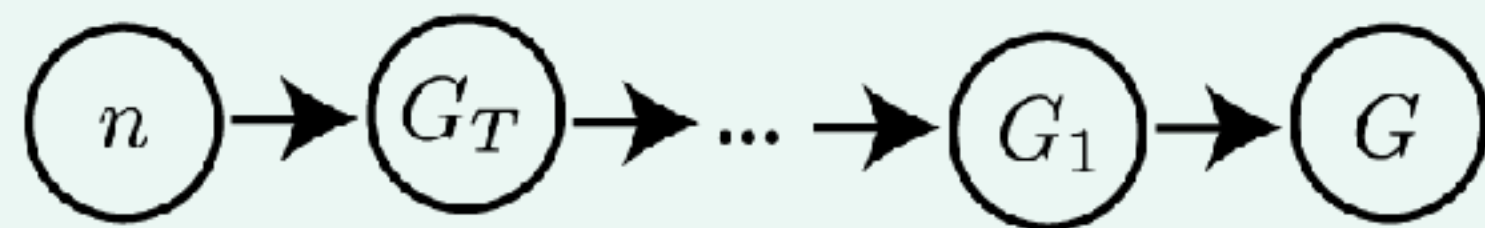
Take as input a graph structure, output 3d coordinates

Hoogeboom et al. 22, Equivariant diffusion for molecule generation in 3d

Xu et al. 22, GeoDiff: a Geometric Diffusion Model for Molecular Conformation Generation

Jing et al. 22, Torsional diffusion for molecular conformer generation

# Likelihood computation



$$\begin{aligned} \log p_\theta(G) &= \log \sum_{n \in \mathbb{N}} p(n) \int p(G^T | n) p_\theta(G^{t-1}, \dots, G^1 | G^T) p_\theta(G | G^1) d(G^1, \dots, G^T) \\ &= \log p(n_G) + \log \int p(G^T | n_G) \prod_{t=2}^T p_\theta(G^{t-1} | G^t) p_\theta(G | G^1) d(G^1, \dots, G^T) \end{aligned}$$

$$\log p_\theta(G) \geq \log p(n_G) + \underbrace{D_{\text{KL}}[q(G^T | G) || q_X(n_G) \times q_E(n_G)]}_{\text{Prior loss}} + \underbrace{\sum_{t=2}^T L_t(x)}_{\text{Diffusion loss}} + \underbrace{\mathbb{E}_{q(G^1 | G)}[\log p_\theta(G | G^1)]}_{\text{Reconstruction loss}}$$

with  $L_t(G) = \mathbb{E}_{q(G^t | G)}[\text{KL}[q(G^{t-1} | G^t, G) || p_\theta(G^{t-1} | G^t)]]$