



Al-Driven Drug Discovery & Graph Generative Model

MoFlow: An Invertible Flow Model for Generating

Molecular Graphs (KDD 2020 Research)



Chengxi Zang and Fei Wang

Weill Cornell Medicine

www.calvinzang.com

Drug Discovery and Development



Background: Drug Discovery





1. Lengthy, costly, & with high failure rate

\$2.6 billion, ≥ 10 years in total, clinical success ~12%, poor translation in patients
 Our focus: Drug discovery (lead discovery and optimization) ~ 5 years and 33% of total cost

□How to accelerate the process, reduce its cost, and increase the success rate?

Nature 2010, Proteomes 2016

Background: Drug Discovery





2. Big chemical space but largely unexplored

o The scale of drug-like small molecules: $10^{33} \sim 10^{60}$

 $_{\rm \sim 10^6}^{\rm o}$ Existing chemical database to (linearly) screen:

•A huge gap! Exhaustive enumeration is impossible!

How to efficiently explore such a big chemical space?

Nature 2017

KDD 2020 -- MOFLOW

Background: Drug Discovery



Design



Evaluate

Image from Sygnaturediscovery

3. Difficult to optimize molecules

• Difficult to design novel & better molecules:

- High-throughput virtual screen, or
- Medicinal chemists' knowledge

• Difficult to **evaluate**:

\$\$ expensive experiments, In-vitro, in-vivo, in-silico

How to efficiently optimize molecules guided by the targeted properties?

Our Vision: Al for Drug Discovery

- Driven by <u>AI</u> and <u>Big Chemical Data</u>
- Ito reduce time, cost and failure rate of drug discovery process o3-5 years → 3-5 months
- to efficiently explore big chemical space
 <u>~ 10⁶⁰</u> drug-like chemical space
- to efficiently and automatically design novel molecules with optimized properties
 <u>automatic</u>, <u>in silico</u>, <u>learning</u> from data and human knowledge

Problem Definition

Goal: To generate novel molecular graphs with optimized properties

Data Input:

oDiscrete 2D molecular graphs, etc.

o{ $G_1, G_2, ..., G_N$ }: Molecular graph data samples o{ $y_{1,k}, y_{2,k}, ..., y_{N,i}$ }_{k=1...K}Some properties of molecules

Output:

•**Novel** molecular graphs $\{G_{N+1}, G_{N+2}, ...\}$ with **optimized** properties.

Metformin (二甲双胍) CN(C)C(=N)N=C(N)N



Discrete molecular graph data and its combinatorial complexity

oNodes/atoms and edges/bonds can have multiple types

- Node types: C, H, O, etc., Edge types: single, double, triple bond.
- •Combinatorial Complexity
 - \clubsuit the scale of small molecular graphs $\sim 10^{60}$

Deep models are majorly designed for regular grid structures (image or text)

VS.

Complex molecular graph optimization task: \circ Graph generation: G \sim P(G) • Graph property prediction: f(G)• Graph optimization: $G \rightarrow G'$ and maximizing f(G') - f(G)Continuous **Property Space** Direct nverse Inverse Desired properties (redox potential, solubility, toxicity) Optimization. Experiment or High-throughput virtual evolutionary strategies, Discrete simulation (Schrödinger screening (e.g., with 3 generative models (VAE, equation) filtering stages) GAN, RL) **Graph Space** rug-like, photovoltaics, Image adapted from Science 2018

10

Why Is It Hard?

Encoding graph is hard, Decoding graph is much harder

Encoding, embedding, inference with graph input



Related Works

Classified by Data:

Sequence: SMILESGraph: molecular graphs

Classified by Deep Generative Models:

Autoregressive Models (AR)
Variational Autoencoders (VAE)
Generative Adversarial Networks (GAN)
Normalizing Flow Models (Flow)

Classified by Search & Optimization

Gradient ascend

Reinforcement learning

<u>KDD2020-</u>Tutorial Recent Advances on Graph Analytics and Its Applications in Healthcare <u>http://www.calvinzang.com/kdd2020_tutorial_medical_graph_analytics.html</u>

Our Choice

Classified by Data:

oSequence: SMILES

OGraph: molecular graphs

Classified by Deep Generative Models:

Autoregressive Models (AR)
Variational Autoencoders (VAE)
Generative Adversarial Networks (GAN)
Normalizing Flow Models (Flow)

Classified by Search & Optimization • Gradient ascend

oReinforcement learning

<u>KDD2020-</u>Tutorial Recent Advances on Graph Analytics and Its Applications in Healthcare <u>http://www.calvinzang.com/kdd2020_tutorial_medical_graph_analytics.html</u>

Basics of Normalizing Flow

An invertible generative model

• Goal: $X \sim P(X)$, by leveraging <u>an invertible mapping</u> $f_{\theta}(X)$

□ Inference: $Z = f_{\theta}(X)$ ◦ From complex to simple, e.g. Z is Gaussian

Generation: $X = f_{\theta}^{-1}(Z)$ • Generate complex by invertible mapping

Exact Maximum Likelihood Training

• Change of variable $\log P(X) = \log P(Z) + \log |\det(\frac{\partial f_{\theta}}{\partial Z})|$ • $\arg \max_{\theta} E_{M \sim P_{data}}[\log P_M(M; \theta)]$

Constraints of network structures:

• f_{θ} : invertible DNNs, each layer is invertible • Computing det $(\frac{\partial f_{\theta}}{\partial z})$ should be efficient

Image from: Dinh et al. 2017. Density Estimation using Real NVP. ICLR.



P(X):

Related works: RealNVP Model

- RealNVP: Real-valued Non-Volume Preserving flow
- Invertible layers: splitting dimensions + affine updated alternately



Dinh et al. 2014. <u>Nice: Non-linear independent components estimation</u> Dinh et al. 2017. <u>Density Estimation using Real NVP.</u> *ICLR*. Chen et al. 2019. <u>Neural Ordinary</u> <u>Differential Equations.</u> *NeurIPS*.

Why Flow Frameworks

Invertible mappings

•Potentials to generate more novel molecules

oVAE, GAN, AR are not invertible, see diagrams below

•Flow learns a strict superset and explores chemical space better





•VAE,GAN are not

Efficient one-shot inference and generation

 Capturing molecular structures in a holistic way v.s. AR's step-bystep way.

Better performance shown later

■ Molecular Graph: Molecule = (Atom, Bond) • Atoms → Nodes, Atom ∈ $\{0,1\}^{n \times k}$, n Nodes in k (atom) types • Bonds → Edges, Bond ∈ $\{0,1\}^{c \times n \times n}$, Edges in c (bond) types



Idea of our MoFlow

•Molecule=(Atom, Bond) How to model discrete atom-bond structures of molecule?

- $\circ P_M(M) = P_M((A,B)) \approx P_B(B) P_{A|B}(A|B)$
- **1.** Any flow model $f_B(B)$ for bonds $P_B(B)$

Conversion Generating graph skeleton by $P_B(B)$

<u>2. Graph conditional flow</u> $f_{A|B}(A|B)$ for atoms given bonds $P_{A|B}(A|B)$

Contracting nodes given graph skeleton by $P_{A|B}(A|B)$

3. Assembling atom and bonds with validity correction

The Generative Framework



KDD 2020 -- MOFLOW

A variant of Glow for Bond/Edge





Graph Conditional Flow For Atoms Given Bonds

Actnorm2D: • Stable dynamics • $B = \frac{B-\mu}{\sqrt{\sigma^2 + c}}$ each row over batch Split: • Discretization of Hamiltonian system on Graphs • $A = (A_1, A_2)$ by each row $\circ \mathbf{Z} = (\mathbf{Z}_{\mathbf{A1}|\mathbf{B}}, \mathbf{Z}_{\mathbf{A2}|\mathbf{B}})$ Graphnorm • $\widehat{B}_i = D^{-1}B_i$, $D = \sum_{c,i} B_{c,i,j}$ in-degree over all channels GraphConv(A|B), multi-channel $\circ \sum_{i=1}^{c} \widehat{B_i}(M \odot A) W_i + (M \odot A) W_0$ o update each row by the remaining rows Affine coupling: Stable (batchnorm, Sigmoid) and expressive power (Affine) $\circ Z_{A1|B} = A_1$

 $\circ Z_{A2|B} = A_2 \odot Sigmoid(S_{\theta}(A_1|B)) + T_{\theta}(A_1|B)$

Deep: alternating update in next layer



Molecular Graph Generation



Graph Property Prediction



Molecular Graph Optimization



□Valid molecules: valency constraints

 $o\sum_{c,j} c * B(c,i,j) ≤ Valency(Atom_i) + Formal_Charge$ oC: 4, O:2, O+:3

Validity Correction

•While checking valency constraints:

if follows constraints:

• Return the greatest connected component

else:

• Delete unnecessary bond or add charge to invalid atoms according to chemical rules

Experiments

- **1.** Molecular Generation & Reconstruction
- **2.** Visualization of Continuous Latent Space
- **3.** Property Optimization
- 4. Constrained Property Optimization

EXP1: Molecular Generation & Reconstruction

The Problem:

oInput: { G_1 , G_2 , ... } molecular graphs oModel

*Learned molecular generative model P_M , and its invertible mapping f

- **Contraction:** $G = f^{-1}(Z)$, where Z follows isotropic Gaussian
- *****Reconstruction: $G = f^{-1}(Z)$ where Z = f(G)

•Goal: To generate valid & unique & novel molecular graphs

Datasets:		#Graphs	#Nodes	#Node/Atom Types	#Edge/Bond Types
0	QM9	134K	9	4	3
	ZINC	250K	38	9	3

EXP1: Molecular Generation & Reconstruction

Evaluation metrics:

- 1. <u>Validity</u>: %chemically valid molecules in all the generated molecules
- 2. Validity without check/correction
- 3. <u>Uniqueness</u>: %chemically valid and unique molecules in all the generated molecules
- 4. <u>Novelty</u>: %generated valid molecules not in training dataset
- 5. Reconstruction rate: % training dataset which can be reconstructed from their latent representations
- 6. N.U.V.: %novel, unique and valid molecules in all the generated molecules

EXP1: Molecular Generation & Reconstruction

More novel & unique & valid	Table 1: Generative performance on QM9							
molecules		% Validity	% Validity w/o check	% Uniqueness	% Novelty	% N.U.V.	% Reconstruct	
	GraphNVP	83.1 ± 0.5	-	99.2 ± 0.3	58.2 ± 1.9	47.97	100	
	GRF	84.5 ± 0.70	-	66.0 ± 1.15	58.6 ± 0.82	32.68	100	
Reconstruction	GraphAF	100	67	94.51	88.83	83.95	100	
 Strict superset of training 	MoFlow	100.00 ± 0.00	95.74 ± 0.65	99.48 ± 0.33	98.69 ± 0.39	98.18 ± 0.53	100.00 ± 0.00	
dataset								

Table 2: Generative performance on Zinc250k

Better validity without check

 Than AR models. Oneshot models, a holistic way

Our MoFlow explores the big chemical space further and better!

	% Validity	% Validity w/o check	% Uniqueness	% Novelty	% N.U.V.	% Reconstruct
JT-VAE	100	-	100	100	100	76.7
GCPN	100	20	99.97	100	99.97	-
MRNN	100	65	99.89	100	99.89	-
GraphNVP	42.6 ± 1.6	-	94.8 ± 0.6	100	40.38	100
GRF	73.4 ± 0.62	-	53.7 ± 2.13	100	39.42	100
GraphAF	100	68	99.10	100	99.10	100
MoFlow	100.00 ± 0.00	81.94 ± 0.45	99.94 ± 0.05	100.00 ± 0.00	99.94 ± 0.05	100.00 ± 0.00

EXP2: Visualization of latent space

Encode & decode between discrete graph space and continuous latent space!

- Grid interpolation around the latent representation of one molecular graph, and decode its neighbors
 - Smooth latent space ← → Similar graph structures (Tanimoto similarity)

Linear interpolation between two molecules

Changing trajectory from one graph to another one.



EXP3: Property Optimization

■ To Generate <u>Novel</u> Molecules with the best Quantitative Estimate of <u>Druglikeness</u> (QED) scores as many as possible

 Searching latent space by gradient ascend

Our MoFlow generates more novel molecules with top QED scores!

Table 3: Discovered novel molecules with top QED score. Our MoFlow finds more molecules with the best QED score. More results in

Method	1st	2nd	3rd	4th
ZINC (Dataset)	0.948	0.948	0.948	0.948
JT-VAE	0.925	0.911	0.910	-
GCPN	0.948	0.947	0.946	-
MRNN	0.948	0.948	0.947	-
GraphAF	0.948	0.948	0.947	0.946
MoFlow	0.948	0.948	0.98	0.948

EXP3: Property Optimization



KDD 2020 -- MOFLOW

EXP4: Constrained Property Optimization

Find a new molecular graph G' from a seed molecular graph G

- To maximize: similarity(\mathbf{G} , \mathbf{G}') and property $Y(\mathbf{G}') Y(\mathbf{G})$
 - Tanimoto similarity of Morgan fingerprint
 - Target property Y: penalized logP (plogP), which is the octanol-water partition coefficients (logP) penalized by the synthetic accessibility (SA) score and number of long cycles.

EXP4: Constrained Property Optimization

Best similarity

Second best improvement

□More realistic

 AR+RL model tends to generate long chains



Table 4: Constrained optimization on Penalized-logP

		JT-VAE			GCPN	
δ	Improvement	Similarity	Success	Improvement	Similarity	Success
0.0	1.91 ± 2.04	0.28 ± 0.15	97.5%	4.20 ± 1.28	$\textbf{0.32} \pm \textbf{0.12}$	100%
0.2	1.68 ± 1.85	0.33 ± 0.13	97.1%	4.12 ± 1.19	0.34 ± 0.11	100%
0.4	0.84 ± 1.45	0.51 ± 0.10	83.6%	2.49 ± 1.30	0.48 ± 0.08	100%
0.6	0.21 ± 0.71	0.69 ± 0.06	46.4%	0.79 ± 0.63	0.68 ± 0.08	100%
		GraphAF		r — — — —	MoFlow	
δ	Improvement	Similarity	Success	Improvement	Similarity	Success
0.0	13.13 ± 6.89	0.29 ± 0.15	100%	8.61 ± 5.44	0.30 ± 0.20	98.88%
0.2	11.90 ± 6.86	0.33 ± 0.12	100%	7.06 ± 5.04	$\textbf{0.43} \pm \textbf{0.20}$	96.75%
0.4	8.21 ± 6.51	0.49 ± 0.09	99.88%	4.71 ± 4.55	0.61 ± 0.18	85.75%
0.6	4.98 ± 6.49	0.66 ± 0.05	96.88%	2.10 ± 2.86	$\boldsymbol{0.79 \pm 0.14}$	58.25%
						.0
			R			
			\longrightarrow			
	╵ͺͺ╶║╶╵		+16.48			

EXP4: Constrained Property Optimization-Visualization.gif

VizBoard
for
Molecule
Generation

Visu

sele

Sel sel

/isualization tasks:	
random generation	
Linear Interpolation	n
Explore Latent Space	ce
Optimization	
elect Model	
Select	
select Dataset	
zinc250k	× +

Optimiz	zatio	n			
Type mething					
8399316623 Cc1ccc2c(c1)c1c3n2	2CC[NH+](C	:CC)[C@I	H]3CCC1		
choose a property:					
Plogp Optimize	× •				
About	Viev	N			
Use presets Rendering style					
Default/reset		× •			
Atom style					
Select		¥.			

COc1ccccc1C(=O)Oc1cc2c3c(c1) C(C)=CC(C)(C)N3C(=O)C2=O



COc1cccc1C(=O)OC1=CC=C2 C(C1)C(C)=CC(C)(C)N2C(=O)C=0

Summary

Novel MoFlow model for molecular graph generation

•A variant of Glow for bonds

A novel Graph conditional flow for atoms given bonds
Novel validity correction
Invertible, fast inference and generation at one shot

□The state-of-the-art results

•Best results for generation and reconstruction

*w.r.t. novelty, uniqueness, validity, and reconstruction rate

•Best results for QED property optimization

More drug-like molecules

 Best similarity scores for constraint optimization and second best improvement scores for plogP





Al-Driven Drug Discovery & Graph Generative Model

MoFlow: An Invertible Flow Model for Generating

Molecular Graphs (KDD 2020 Research)



Chengxi Zang and Fei Wang

Weill Cornell Medicine

www.calvinzang.com