

# CS6208 : Advanced Topics in Artificial Intelligence

## Graph Machine Learning

### Lecture 9 : Molecular Science and Graph Generative Models

Semester 2 2022/23

Xavier Bresson

<https://twitter.com/xbresson>



Department of Computer Science  
National University of Singapore (NUS)



# Outline

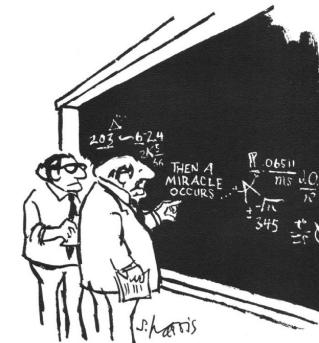
- Deep learning for molecular science
- Graph generative model with VAE
  - Auto-encoder
  - Encoder and decoder
  - System description
  - Numerical experiments
- Conclusion

# Outline

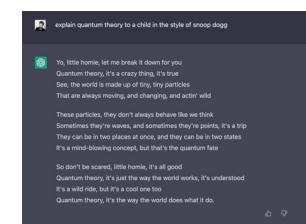
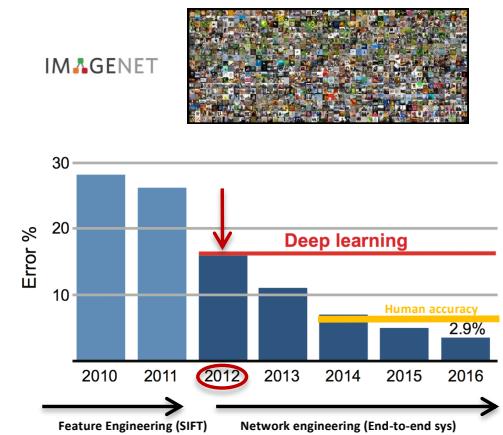
- Deep learning for molecular science
- Graph generative model with VAE
  - Auto-encoder
  - Encoder and decoder
  - System description
  - Numerical experiments
- Conclusion

# Deep Learning

- The DL revolution – 2012 ImageNet
  - First breakthrough in computer vision (CV)
  - Then speech (SR) and language processing (NLP)
- DL works very well.. in practice but why?
  - High-dimensional data/curse of dimensionality
  - Characterization of energy landscapes/solutions
  - Generalization (feature learning, inductive bias, expressivity)
- DL products
  - Computer Vision : Face recognition, video surveillance, autonomous driving
  - Natural Language Processing : Machine translation, chatbot (ChatGPT)
  - Speech Recognition : Virtual assistants (Alexa/Siri/Google/Cortana)
- DL beyond CV/NLP/SR for scientific discovery

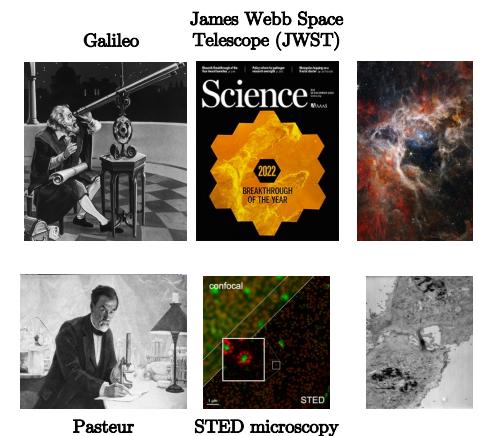
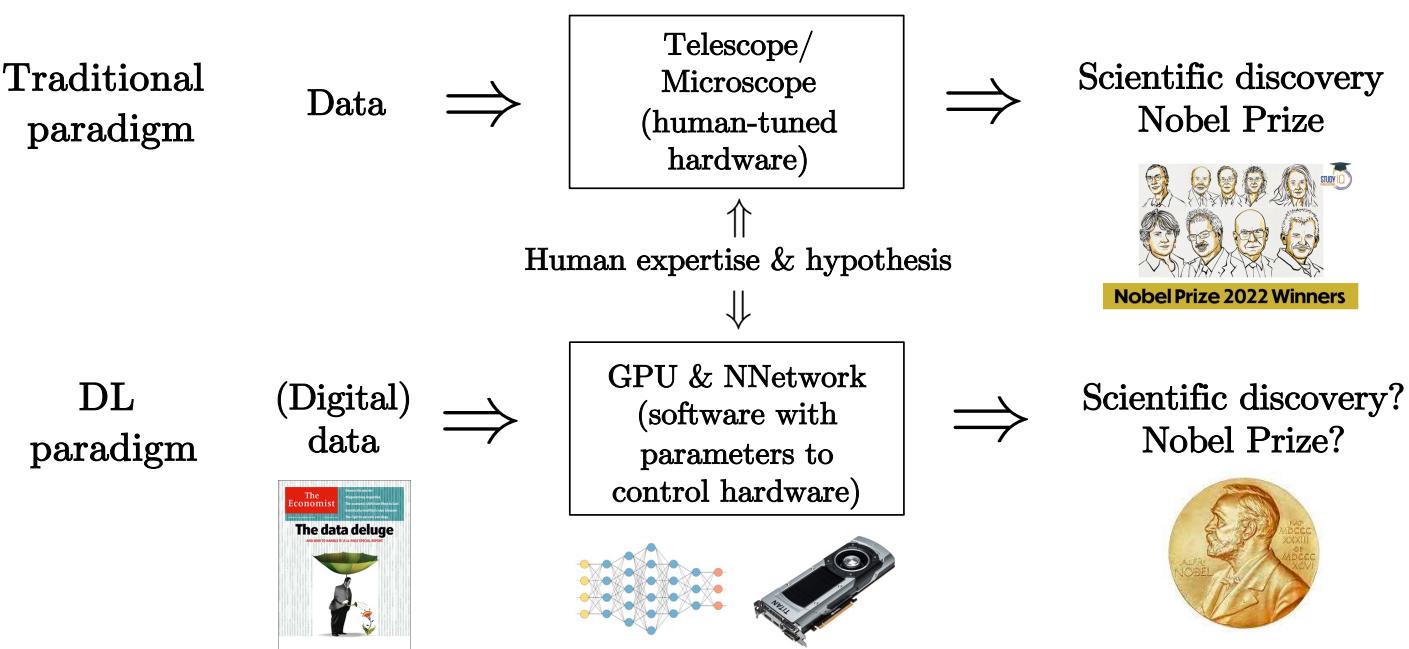


"I THINK YOU SHOULD BE MORE EXPLICIT HERE IN STEP TWO."



# DL for science

- DL has been applied to biology, physics, chemistry, material, medicine, engineering, etc
- Scientific tools s.a. telescope/microscope have served physics (Galileo's proof that earth is round and revolves around the sun, study of universe) and biology (Pasteur's discovery of microbes, study of cells, molecules).
- Can DL be the new telescope/microscope for scientific discovery?



# DL for molecular science

- A promising application of DL and potentially breakthrough is in molecular science.
- Case studies
  - Drug discovery (halicin antibiotic)
  - De novo drug design (new molecules w/ optimized chemical property)
  - Protein folding (3D structure)
  - Protein-drug interaction (drugs with strong binding)
  - Generate protein with text prompt (ChatGPT for biology)

# GNNs for DFT estimation

- Density Functional Theory (DFT) is one of the main workhorses to estimate molecular properties s.a. energy (Walter Kohn awarded Nobel Prize in Chemistry in '98 for DFT).
- DFT is  $O(N^3)$ , N being the number of atoms vs. GNNs being  $O(N)$ .

---

## Neural Message Passing for Quantum Chemistry

---

Justin Gilmer<sup>1</sup> Samuel S. Schoenholz<sup>1</sup> Patrick F. Riley<sup>2</sup> Oriol Vinyals<sup>3</sup> George E. Dahl<sup>1</sup>

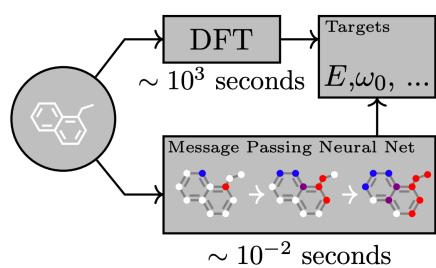
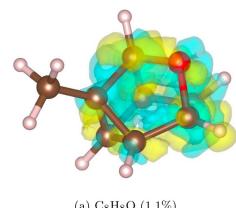
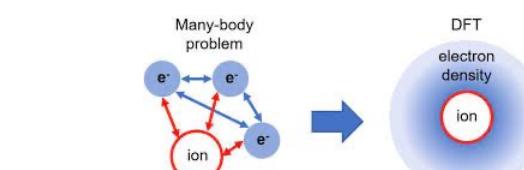


Figure 1. A Message Passing Neural Network predicts quantum properties of an organic molecule by modeling a computationally expensive DFT calculation.



Xavier Bresson



npj | computational materials

Explore content ▾ About the journal ▾ Publish with us ▾

nature > npj computational materials > articles > article

Article | Open Access | Published: 23 August 2022

### Equivariant graph neural networks for fast electron density estimation of molecules, liquids, and solids

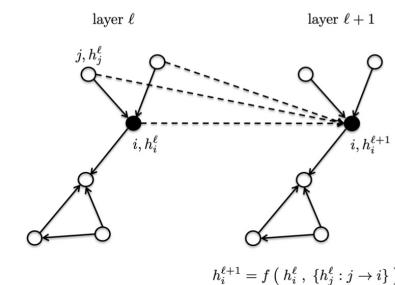
Peter Bjørn Jørgensen & Arghya Bhownik

npj Computational Materials 8, Article number: 183 (2022) | Cite this article

3382 Accesses | 13 Altmetric | Metrics

#### Abstract

Electron density  $\rho(\vec{r})$  is the fundamental variable in the calculation of ground state energy with density functional theory (DFT). Beyond total energy, features and changes in  $\rho(\vec{r})$  distributions are often used to capture critical physicochemical phenomena in functional materials. We present a machine learning framework for the prediction of  $\rho(\vec{r})$ . The model is based on equivariant graph neural networks and the electron density is predicted at special query point vertices that are part of the message-passing graph, but only receive messages. The model is tested across multiple datasets of molecules (QM9), liquid ethylene carbonate electrolyte (EC) and Li<sub>x</sub>Ni<sub>y</sub>Mn<sub>z</sub>Co<sub>(1-y-z)/2</sub>O<sub>2</sub> lithium ion battery cathodes (NMC). For QM9 molecules, the accuracy of the proposed model exceeds typical variability in  $\rho(\vec{r})$  obtained from DFT done with different exchange-correlation functionals. The accuracy on all three datasets is beyond state of the art and the computation time is orders of magnitude faster than DFT.



npj | computational materials

Explore content ▾ About the journal ▾ Publish with us ▾

nature > npj computational materials > articles > article

Article | Open Access | Published: 21 May 2021

### Accurate and scalable graph neural network force field and molecular dynamics with direct force architecture

Cheol Woo Park, Mordechai Kornblith, Jonathan Vandermuse, Chris Wolverton, Boris Kozinsky & Jonathan P. Malloia

npj Computational Materials 7, Article number: 73 (2021) | Cite this article

9838 Accesses | 36 Citations | 1 Altmetric | Metrics

#### Abstract

Recently, machine learning (ML) has been used to address the computational cost that has been limiting ab initio molecular dynamics (AIMD). Here, we present GNNFF, a graph neural network framework to directly predict atomic forces from automatically extracted features of the local atomic environment that are translationally-invariant, but rotationally-covariant to the coordinate of the atoms. We demonstrate that GNNFF not only achieves high performance in terms of force prediction accuracy and computational speed on various materials systems, but also accurately predicts the forces of a large MD system after being trained on forces obtained from a smaller system. Finally, we use our framework to perform an MD simulation of Li<sub>2</sub>P<sub>3</sub>S<sub>11</sub>, a superionic conductor, and show that resulting Li diffusion coefficient is within 14% of that obtained directly from AIMD. The high performance exhibited by GNNFF can be easily generalized to study atomistic level dynamics of other material systems.

# GNNs for drug discovery

Cell

Volume 180, Issue 4, 20 February 2020, Pages 688-702.e13



Article

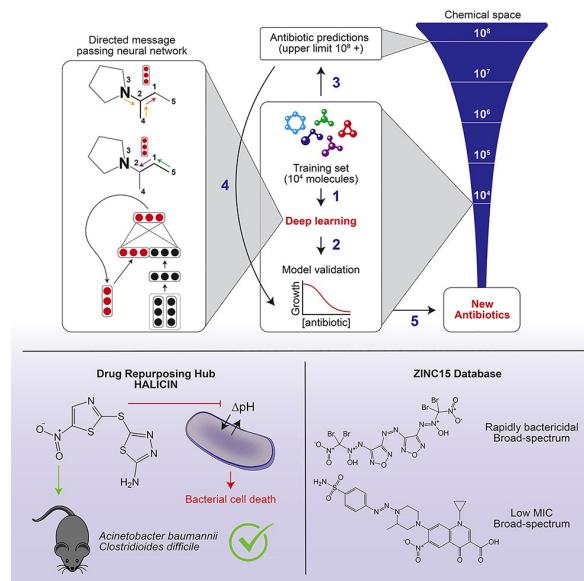
## A Deep Learning Approach to Antibiotic Discovery

Jonathan M. Stokes<sup>1, 2, 3</sup>, Kevin Yang<sup>3, 4, 10</sup>, Kyle Swanson<sup>3, 4, 10</sup>, Wengong Jin<sup>3, 4</sup>, Andres Cubillos-Ruiz<sup>1, 2, 5</sup>, Nina M. Donghia<sup>1, 5</sup>, Craig R. MacNair<sup>6</sup>, Shawn French<sup>6</sup>, Lindsey A. Carfrae<sup>6</sup>, Zohar Bloom-Ackermann<sup>2, 7</sup>, Victoria M. Tran<sup>2</sup>, Anush Chiappino-Pepe<sup>5, 7</sup>, Ahmed H. Badran<sup>2</sup>, Ian W. Andrews<sup>1, 2, 5</sup>, Emma J. Chory<sup>1, 2</sup>, George M. Church<sup>5, 7, 8</sup>, Eric D. Brown<sup>6</sup>, Tommi S. Jaakkola<sup>3, 4</sup> ... James J. Collins<sup>1, 2, 5, 8, 9, 11</sup>  

### Highlights

- A deep learning model is trained to predict antibiotics based on structure
- Halicin is predicted as an antibacterial molecule from the Drug Repurposing Hub
- Halicin shows broad-spectrum antibiotic activities in mice
- More antibiotics with distinct structures are predicted from the ZINC15 database

<https://www.sciencedirect.com/science/article/pii/S0092867420301021>



Halicin was previously developed for anti-diabetic treatment.

MIT  
Technology  
Review

77 Mass Ave

AI vs. bacteria

nature

Explore content ▾ Journal information ▾ Publish with us ▾ Subscribe Sign up

nature > news > article

NEWS · 20 FEBRUARY 2020

### Powerful antibiotics discovered using AI

Machine learning spots molecules that work even against 'untreatable' strains of bacteria.

FINANCIAL TIMES

Artificial intelligence 

AI discovers antibiotics to treat drug-resistant diseases

Machine learning uncovers potent new drug able to kill 35 powerful bacteria

Quanta magazine Physics Mathematics Biology Computer Science All Articles

ARTIFICIAL INTELLIGENCE  
Machine Learning Takes On Antibiotic Resistance

To combat resistant bacteria and refill the trickling antibiotic pipeline, scientists are getting help from deep learning networks.



# GNNs for de novo drug design

## Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules

Rafael Gómez-Bombarelli,<sup>†,‡,§,¶</sup> Jennifer N. Wei,<sup>‡,§,¶</sup> David Duvenaud,<sup>¶,§,¶</sup> José Miguel Hernández-Lobato,<sup>§,¶</sup> Benjamin Sánchez-Lengeling,<sup>§,¶</sup> Dennis Sheberla,<sup>‡,§,¶</sup> Jorge Aguilera-Iparraguirre,<sup>¶</sup> Timothy D. Hirzel,<sup>¶</sup> Ryan P. Adams,<sup>¶,||</sup> and Alán Aspuru-Guzik<sup>¶,‡,§,¶,||</sup>

<sup>†</sup>Kylia North America Inc., 10 Post Office Square, Suite 800, Boston, Massachusetts 02109, United States

<sup>‡</sup>Department of Chemistry and Chemical Biology, Harvard University, Cambridge, Massachusetts 02138, United States

<sup>§</sup>Department of Computer Science, University of Toronto, 6 King's College Road, Toronto, Ontario M5S 3H5, Canada

<sup>¶</sup>Department of Engineering, University of Cambridge, Trumpington Street, Cambridge CB2 1FZ, U.K.

<sup>||</sup>Google Brain, Mountain View, California, United States

<sup>||</sup>Princeton University, Princeton, New Jersey, United States

<sup>¶</sup>Bioinspired Solar Energy Program, Canadian Institute for Advanced Research (CIFAR), Toronto, Ontario M5S 1M1, Canada

### Supporting Information

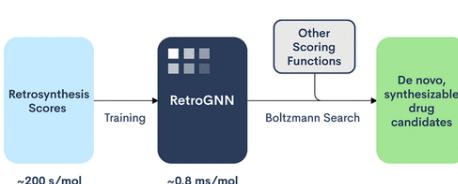
**ABSTRACT:** We report a method to convert discrete representations of molecules to and from a multidimensional continuous representation. This model allows us to generate new molecules for efficient exploration and optimization through open-ended spaces of chemical compounds. A deep neural network was trained on hundreds of thousands of existing chemical structures to construct three coupled functions: an encoder, a decoder, and a predictor. The encoder converts the discrete representation of a molecule into a real-valued continuous vector, and the decoder converts these continuous vectors back to discrete molecular representations. The predictor estimates chemical properties from the latent continuous vector representation of the molecule. Continuous representations of molecules allow us to automatically generate novel chemical structures by performing simple operations in the latent space, perturbing known chemical structures, or interpolating between molecules. Continuous representations also allow the use of powerful gradient-based optimization to efficiently guide the search for optimized functional compounds. We demonstrate our method in the domain of drug-like molecules and also in a set of molecules with fewer than nine heavy atoms.



## RetroGNN: Fast Estimation of Synthesizability for Virtual Screening and De Novo Design by Learning from Slow Retrosynthesis Software

Cheng-Hao Liu<sup>1\*</sup>, Maksym Korablyov<sup>1</sup>, Stanisław Jastrzębski<sup>1</sup>, Paweł Włodarczyk-Pruszyński<sup>1</sup>, Yoshua Bengio<sup>1</sup>, and Marwin Segler<sup>1\*</sup>

Cite this: *J. Chem. Inf. Model.* 2022, 62, 10, 2293–2300  
Publication Date: April 22, 2022  
<https://doi.org/10.1021/acs.jcim.1c01476>  
Copyright © 2022 American Chemical Society



Xavier Bresson

~200 s/mol

~0.8 ms/mol

## scientific reports

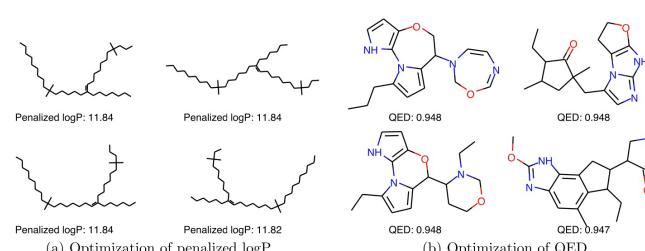
Explore content ▾ About the journal ▾ Publish with us ▾

nature > scientific reports > articles > article

Article | Open Access | Published: 24 July 2019

## Optimization of Molecules via Deep Reinforcement Learning

Zhenpeng Zhou, Steven Kearnes, Li Li, Richard N. Zare & Patrick Riley



**Figure 2.** Sample molecules in the property optimization task. (a) Optimization of penalized logP from MoldQNN-bootstrap; note that the generated molecules are obviously not drug-like due to the use of a single-objective reward. (b) Optimization of QED from MoldQNN-twosteps.

**logP:** molecular solubility

**QED:** Quantitative Estimate of Druglikeness

Sundin et al. *Journal of Cheminformatics* (2022) 14:86  
<https://doi.org/10.1186/s13321-022-00667-8>

Journal of Cheminformatics

## RESEARCH

## Open Access

## Human-in-the-loop assisted de novo molecular design

Iiris Sundin<sup>1\*</sup>, Alexey Voronov<sup>2</sup>, Haoping Xiao<sup>1</sup>, Kostas Papadopoulos<sup>2,5</sup>, Esben Jannik Bjerrum<sup>2,5</sup>, Markus Heinonen<sup>1</sup>, Atanas Patronov<sup>2,5</sup>, Samuel Kaski<sup>1,3</sup> and Ola Engkvist<sup>2,4</sup>

REVIEWS  
Drug Discovery Today • Volume 26, Number 6 • June 2021  
Teaser This paper highlights the progress and challenges in de novo drug design using graph neural network technology.  
Check for updates

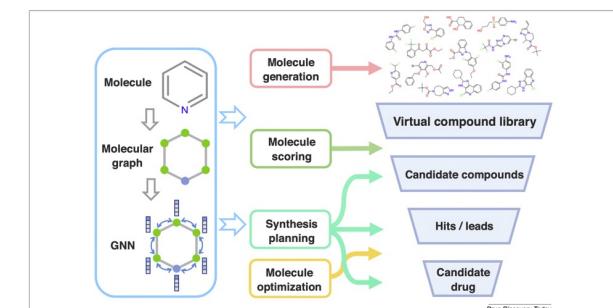
## Graph neural networks for automated de novo drug design

Jiacheng Xiong<sup>1,2</sup>, Zhaoping Xiong<sup>1,3</sup>, Kaixian Chen<sup>1,2</sup>, Hualiang Jiang<sup>1,2</sup> and Mingyue Zheng<sup>1,2</sup>

<sup>1</sup>Drug Discovery and Design Center, State Key Laboratory of Drug Research, Shanghai Institute of Materia Medica, Chinese Academy of Sciences, 555 Zuchongzhi Road, Shanghai 201203, China

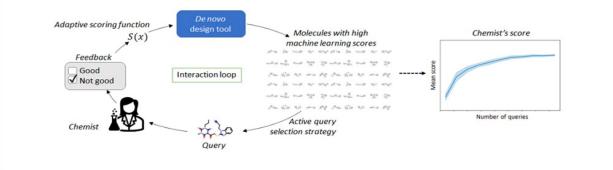
<sup>2</sup>University of Chinese Academy of Sciences, No. 19A Yuquan Road, Beijing 100049, China

<sup>3</sup>Shanghai Institute for Advanced Immunological Studies, and School of Life Science and Technology, ShanghaiTech University, Shanghai 200031, China



**FIGURE 1**  
The applications of graph neural networks (GNNs) in all stages of automated de novo drug design.

## Graphical Abstract



# GNNs for molecular generation

- Categories of graph generative models : VAE, GAN, Diffusion Model

## A Two-Step Graph Convolutional Decoder for Molecule Generation

Xavier Bresson  
 School of Computer Science and Engineering  
 NTU, Singapore  
`xbresson@ntu.edu.sg`

Thomas Laurent  
 Department of Mathematics  
 Loyola Marymount University  
`tlaurent@lmu.edu`

### Abstract

We propose a simple auto-encoder framework for molecule generation. The molecular graph is first encoded into a continuous latent representation  $z$ , which is then decoded back to a molecule. The encoding process is easy, but the decoding process remains challenging. In this work, we introduce a simple two-step decoding process. In a first step, a fully connected neural network uses the latent vector  $z$  to produce a molecular formula, for example  $\text{CO}_2$  (one carbon and two oxygen atoms). In a second step, a graph convolutional neural network uses the same latent vector  $z$  to place bonds between the atoms that were produced in the first step (for example a double bond will be placed between the carbon and each of the oxygens). This two-step process, in which a bag of atoms is first generated, and then assembled, provides a simple framework that allows us to develop an efficient molecule auto-encoder. Numerical experiments on basic tasks such as novelty, uniqueness, validity and optimized chemical property for the 250k ZINC molecules demonstrate the performances of the proposed system. Particularly, we achieve the highest reconstruction rate of 90.5%, improving the previous rate of 76.7%. We also report the best property improvement results when optimization is constrained by the molecular distance between the original and generated molecules.

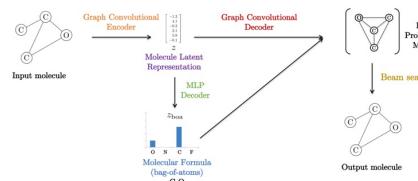


Figure 1: Proposed auto-encoder. The encoder reduces the molecular graph to a latent vector  $z$ . The decoder uses a MLP to produce a molecular formula and a graph convolutional network classifies each bond between the atoms given by the molecular formula. Finally, a beam search generates a valid molecule.

## Graph Convolutional Policy Network for Goal-Directed Molecular Graph Generation

Jiaxuan You<sup>1\*</sup>  
`jiaxuan@stanford.edu`

Bowen Liu<sup>2\*</sup>  
`liubowen@stanford.edu`

Rex Ying<sup>1</sup> Vijay Pande<sup>3</sup> Jure Leskovec<sup>1</sup>  
`rexying@stanford.edu` `pande@stanford.edu` `jure@cs.stanford.edu`

<sup>1</sup>Department of Computer Science, <sup>2</sup>Department of Chemistry, <sup>3</sup>Department of Bioengineering  
 Stanford University  
 Stanford, CA, 94305

### Abstract

Generating novel graph structures that optimize given objectives while obeying some given underlying rules is fundamental for chemistry, biology and social science research. This is especially important in the task of molecular graph generation, whose goal is to discover novel molecules with desired properties such as drug-likeness and synthetic accessibility, while obeying physical laws such as chemical valency. However, designing models to find molecules that optimize desired properties while incorporating highly complex and non-differentiable rules remains to be a challenging task. Here we propose Graph Convolutional Policy Network (GCPN), a general graph convolutional network based model for goal-directed graph generation through reinforcement learning. The model is trained to optimize domain-specific rewards and adversarial loss through policy gradient, and acts in an environment that incorporates domain-specific rules. Experimental results show that GCPN can achieve 61% improvement on chemical property optimization over state-of-the-art baselines while resembling known molecules, and achieve 184% improvement on the constrained property optimization task.

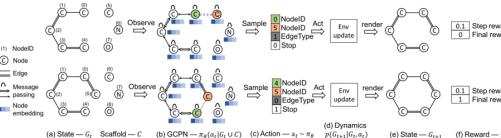


Figure 1: An overview of the proposed iterative graph generation method. Each row corresponds to one step in the generation process. (a) The state is defined as the intermediate graph  $G_t$ , and the set of scaffold subgraphs defined as  $C$  is appended for GCPN calculation. (b) GCPN conducts message passing to encode the state as node embeddings then produce a policy  $\pi_\theta$ . (c) An action  $a_t$  with 4 components is sampled from the policy. (d) The environment performs a chemical valency check on the intermediate state, and then returns (e) the next state  $G_{t+1}$  and (f) the associated reward  $r_t$ .

## Equivariant Diffusion for Molecule Generation in 3D

Emiel Hoogeboom<sup>\*1</sup> Victor Garcia Satorras<sup>\*1</sup> Clément Vignac<sup>\*2</sup> Max Welling<sup>1</sup>

### Abstract

This work introduces a diffusion model for molecule generation in 3D that is equivariant to Euclidean transformations. Our Equivariant Diffusion Model (EDM) learns to denoise a diffusion process with an equivariant network that jointly operates on both continuous (atom coordinates) and categorical features (atom types). In addition, we provide a probabilistic analysis which admits likelihood computation of molecules using our model. Experimentally, the proposed method significantly outperforms previous 3D molecular generative methods regarding the quality of generated samples and efficiency at training time.

### 1. Introduction

Modern deep learning methods are starting to make an important impact on molecular sciences. Behind the success of AlphaFold in protein folding prediction (AlQuraishi, 2019), an increasing body of literature develops deep learning models to analyze or synthesize (in silico) molecules (Simonovsky & Komodakis, 2018; Gebauer et al., 2019; Klicpera et al., 2020; Simm et al., 2021).

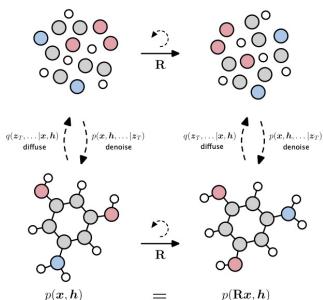


Figure 1. Overview of the EDM. To generate a molecule, a normal distributed set of points is denoised into a molecule consisting of atom coordinates  $\mathbf{x}$  and atom types  $\mathbf{h}$ . As the model is rotated equivariant, the likelihood is preserved when a molecule is rotated by  $\mathbf{R}$ .

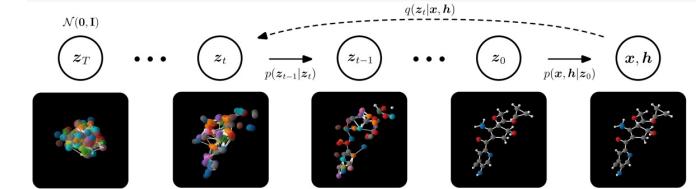


Figure 2. Overview of the Equivariant Diffusion Model. To generate molecules, coordinates  $\mathbf{x}$  and features  $\mathbf{h}$  are generated by denoising variables  $\mathbf{z}_t$  starting from standard normal noise  $\mathbf{z}_T$ . This is achieved by sampling from the distributions  $p(\mathbf{z}_{t-1}|\mathbf{z}_t)$  iteratively. To train the model, noise is added to a datapoint  $\mathbf{x}, \mathbf{h}$  using  $q(\mathbf{z}_t|\mathbf{x}, \mathbf{h})$  for the step  $t$  of interest, which the network then learns to denoise.

# GNNs for protein folding/structure

**AlphaFold: a solution to a 50-year-old grand challenge in biology**

SHARE



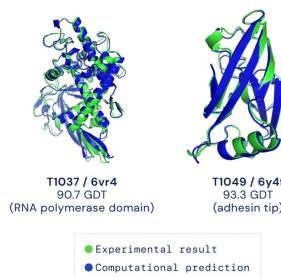
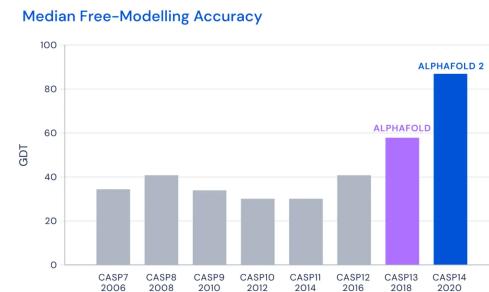
AUTHORS

TAt The AlphaFold team

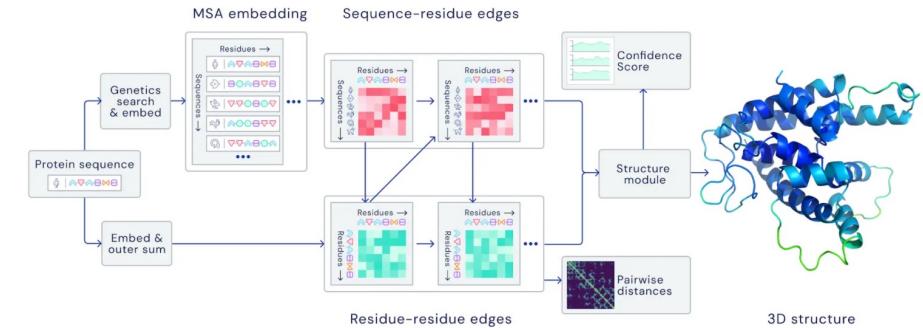


Proteins are essential to life, supporting practically all its functions. They are large complex molecules, made up of chains of amino acids, and [what a protein does largely depends on its unique 3D structure](#). Figuring out what shapes proteins fold into is known as the “[protein folding problem](#)”, and has stood as a grand challenge in biology for the past 50 years. In a major scientific advance, the latest version of our AI system [AlphaFold](#) has been recognised as a solution to this grand challenge by the organisers of the biennial Critical Assessment of protein Structure Prediction ([CASP](#)). This breakthrough demonstrates the impact AI can have on scientific discovery and its potential to dramatically accelerate progress in some of the most fundamental fields that explain and shape our world.

<https://deepmind.com/blog/article/alphafold-a-solution-to-a-50-year-old-grand-challenge-in-biology>



Xavier Bresson



Graph Transformer

Science

RESEARCH ARTICLES  
Cite as: M. Back et al., *Science* 10.1126/science.abj8754 (2021).

**Accurate prediction of protein structures and interactions using a three-track neural network**

Minkyung Back<sup>1,2</sup>, Frank DiMaio<sup>1,2</sup>, Ivan Anishchenko<sup>1,2</sup>, Justas Dauparas<sup>1,2</sup>, Sergey Ovchinnikov<sup>3,4</sup>, Gyu Rie Lee<sup>1,2</sup>, Jue Wang<sup>1,2</sup>, Qian Cong<sup>1,2</sup>, Lisa N. Kinch<sup>1</sup>, R. Dustin Schaeffer<sup>1</sup>, Claudia Millán<sup>5</sup>, Hahnbeom Park<sup>1,2</sup>, Carson Adams<sup>1,2</sup>, Caleb R. Glassman<sup>1,10</sup>, Andy DeGiovanni<sup>1,2</sup>, Jose H. Pereira<sup>12</sup>, Andria V. Rodrigues<sup>12</sup>, Alberdina A. van Dijk<sup>12</sup>, Ana C. Ebrecht<sup>12</sup>, Diederik J. Opperman<sup>14</sup>, Theo Sagmeister<sup>15</sup>, Christoph Buhpheller<sup>15,16</sup>, Tea Pavkov-Keller<sup>15,17</sup>, Manoj K. Rathinaswamy<sup>18</sup>, Udit Dalwadi<sup>19</sup>, Calvin K. Yip<sup>19</sup>, John E. Burke<sup>19</sup>, K. Christopher Garcia<sup>10,11,20</sup>, Nick V. Grishin<sup>6,21,27</sup>, Paul D. Adams<sup>12,22</sup>, Randy J. Read<sup>24</sup>, David Baker<sup>1,2,23\*</sup>

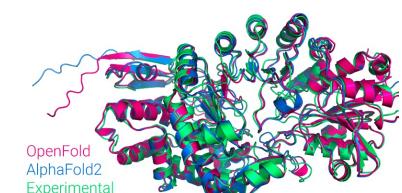


RoseTTAFold

bioRxiv

THE PREPRINT SERVER FOR BIOLOGY  
OpenFold: Retraining AlphaFold2 yields new insights into its learning mechanisms and capacity for generalization

Gustaf Ahrlitz, Nazim Bouatta, Sachin Kadyan, Qinghai Xia, William Gerecke, Timothy J. O'Donnell, Daniel Berenberg, Ian Fisk, Nicollo Zanichelli, Bo Zhang, Arkadiusz Novaczynski, Bei Wang, Marta M. Stepienwska-Dziubinska, Shang Zhang, Adegoke Ojewole, Murat Efe Guney, Stella Biderman, Andrew M Watkins, Stephen Ra, Pablo Ribalta Lorenzo, Lucas Nixon, Brian Weitzner, Yih-En Andrew Ban, Peter K Sorger, Emad Mostaque, Zhao Zhang, Richard Bonneau, Mohammed AlQuraishi



OpenFold

# GNNs for protein-drug discovery & design

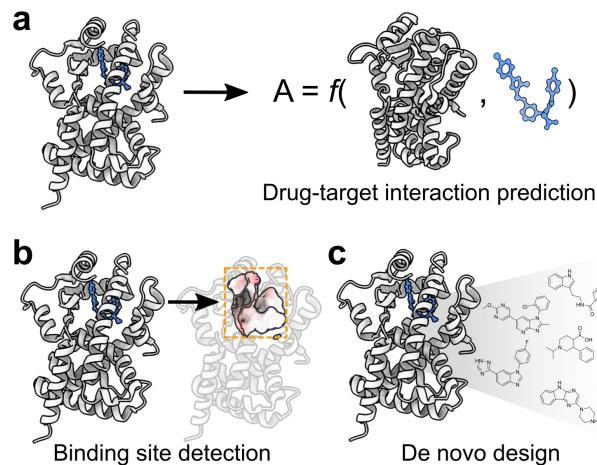
## Structure-based drug discovery with deep learning

Riza Özcelik<sup>1,2\*</sup>, Derek van Tilborg<sup>1,2\*</sup>, José Jiménez-Luna<sup>3</sup>, and Francesca Grisoni<sup>1,2\*</sup>

<sup>1</sup>Eindhoven University of Technology, Institute for Complex Molecular Systems and Dept. Biomedical Engineering, Eindhoven, Netherlands.

<sup>2</sup>Centre for Living Technologies, Alliance TU/e, WUR, UU, UMC Utrecht, Netherlands.

<sup>3</sup>Microsoft Research Cambridge, Cambridge, United Kingdom.

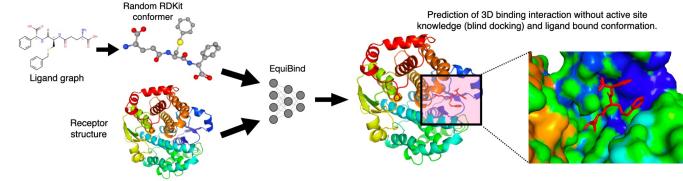


**Figure 1. Structure-based drug discovery tasks discussed in this review:** (a) *drug-target interaction prediction*, which aims to predict the affinity between a protein and a ligand, using the structural information of both molecular entities; (b) *binding site detection*, which aims to identify ‘druggable’ cavities in the protein structure; (c) *de novo design*, aiming to design bioactive molecules from scratch using the information of a protein target.



**EquiBind: Geometric Deep Learning for Drug Binding Structure Prediction**  
Hannes Stärk, Octavian-Eugen Ganea, Lagnajit Pattaik, Regina Barzilay, Tommi Jaakkola

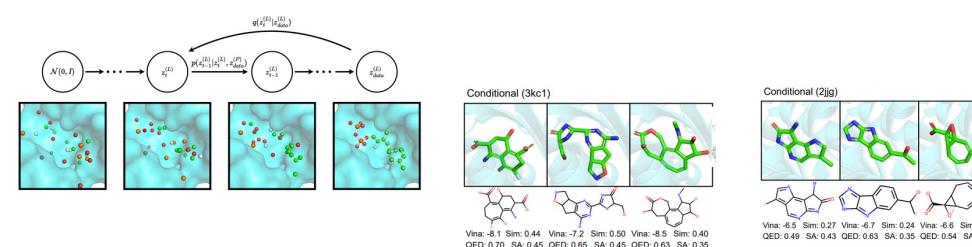
Predicting how a drug-like molecule binds to a specific protein target is a core problem in drug discovery. An extremely fast computational binding method would enable key applications such as fast virtual screening or drug engineering. Existing methods are computationally expensive as they rely on heavy candidate sampling coupled with scoring, ranking, and fine-tuning steps. We challenge this paradigm with EquiBind, an SE(3)-equivariant geometric deep learning model performing direct-shot prediction of both i) the receptor binding location (blind docking) and ii) the ligand’s bound pose and orientation. EquiBind achieves significant speed-ups and better quality compared to traditional and recent baselines. Further, we show extra improvements when coupling it with existing fine-tuning techniques at the cost of increased running time. Finally, we propose a novel and fast fine-tuning model that adjusts torsion angles of a ligand’s rotatable bonds based on closed-form global minima of the von Mises angular distance to a given input atomic point cloud, avoiding previous expensive differential evolution strategies for energy minimization.



## STRUCTURE-BASED DRUG DESIGN WITH EQUIVARIANT DIFFUSION MODELS

Arne Schneuring<sup>1\*</sup>, Yuanqi Du<sup>2\*</sup>, Charles Harris<sup>3</sup>, Arian Jamash<sup>3</sup>, Ilya Igashov<sup>1</sup>, Weitao Du<sup>4</sup>, Tom Blundell<sup>3</sup>, Pietro Lió<sup>3</sup>, Carla Gomes<sup>2</sup>, Max Welling<sup>5</sup>, Michael Bronstein<sup>6</sup> & Bruno Correia<sup>1</sup>

<sup>1</sup>École Polytechnique Fédérale de Lausanne, <sup>2</sup>Cornell University, <sup>3</sup>University of Cambridge, <sup>4</sup>USTC, <sup>5</sup>Microsoft Research AI4Science, <sup>6</sup>University of Oxford



# Generate protein with text prompt

- Inspired by breakthrough multi-modal image-language models s.a. Midjourney, OpenAI's DALL-E, Google Brain's Imagen and Stable Diffusion:



Midjourney image from the prompt "swimming pool filled with a galaxy on a moonlit night" 



An image generated with DALL-E 2 based on the text prompt "Teddy bears working on new AI research underwater with 1990s technology"



An image generated by Stable Diffusion based on the text prompt "a photograph of an astronaut riding a horse"

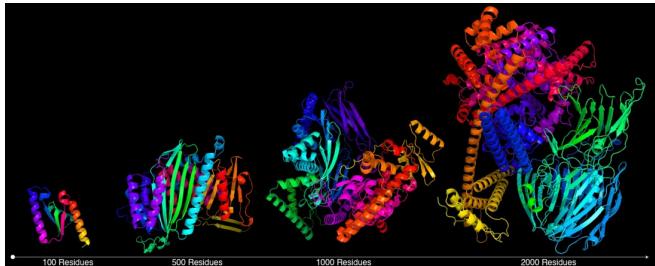
# Generate protein with text prompt

Illuminating protein space  
with a programmable generative model

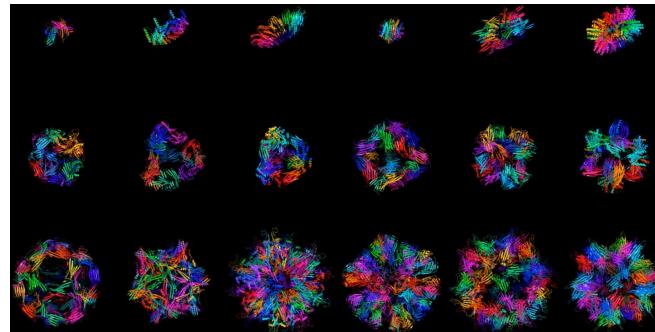
John Ingraham, Max Baranov, Zak Costello, Vincent Frappier,  
Ahmed Ismail, Shan Tie, Wujie Wang, Vincent Xue, Fritz Obermeyer,  
Andrew Beam, Gevorg Grigoryan

<https://www.generatebiomedicines.com/chroma>

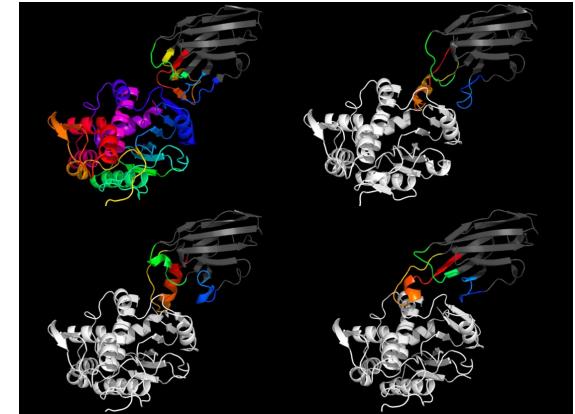
Programming proteins with Chroma



Making giants



Symmetry groups



Protein infilling

# Generate protein with text prompt

## Broadly applicable and accurate protein design by integrating structure prediction networks and diffusion generative models

Joseph L. Watson<sup>#1,2</sup>, David Juergens<sup>#1,2,3</sup>, Nathaniel R. Bennett<sup>#1,2,3</sup>, Brian L. Trippe<sup>#2,4</sup>, Jason Yim<sup>#2,6</sup>, Helen E. Eisenach<sup>#1,2</sup>, Woody Ahern<sup>#1,2,7</sup>, Andrew J. Borst<sup>1,2</sup>, Robert J. Ragotte<sup>1,2</sup>, Lukas F. Milles<sup>1,2</sup>, Basile I. M. Wicky<sup>1,2</sup>, Nikita Hanikel<sup>1,2</sup>, Samuel J. Pellock<sup>1,2</sup>, Alexis Courbet<sup>1,2,9</sup>, William Sheffler<sup>1,2</sup>, Jue Wang<sup>1,2</sup>, Preetham Venkatesh<sup>1,2,8</sup>, Isaac Sappington<sup>1,2,8</sup>, Susana Vázquez Torres<sup>1,2,8</sup>, Anna Lauko<sup>1,2,8</sup>, Valentin De Bortoli<sup>9</sup>, Emile Mathieu<sup>10</sup>, Regina Barzilay<sup>6</sup>, Tommi S. Jaakkola<sup>6</sup>, Frank DiMaio<sup>1,2</sup>, Minkyung Baek<sup>12</sup>, David Baker<sup>\*1,2,11</sup>

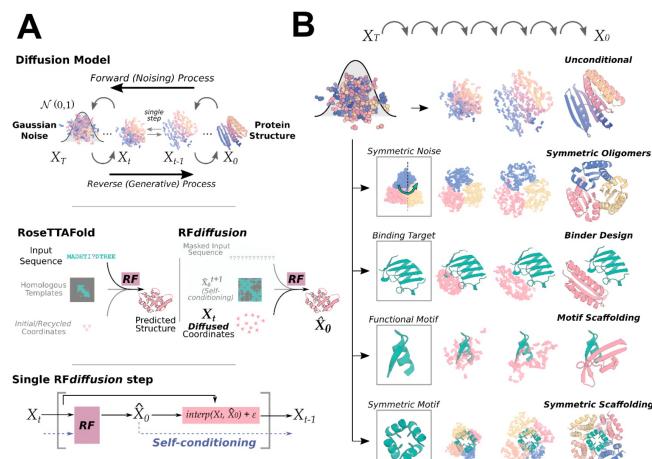


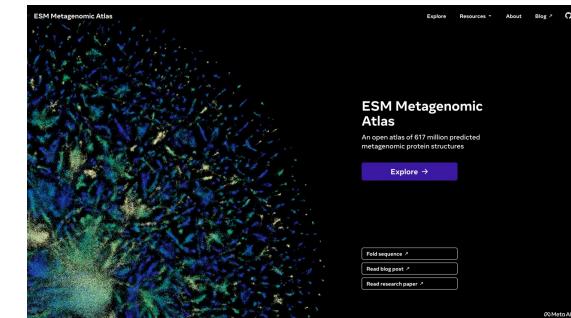
Figure 1: RFdiffusion is a denoising diffusion probabilistic model with RoseTTAFold fine-tuned as the denoising network. A) Top panel: Diffusion models for proteins are trained

RoseTTAFold Diffusion: “Generate a protein that binds to X”

Xavier Bresson

## Evolutionary-scale prediction of atomic level protein structure with a language model

Zeming Lin<sup>1,2\*</sup> Halil Akin<sup>1\*</sup> Roshan Rao<sup>1\*</sup> Brian Hie<sup>1,3\*</sup> Zhongkai Zhu<sup>1</sup> Wenting Lu<sup>1</sup> Nikita Smetanin<sup>1</sup> Robert Verkuil<sup>1</sup> Ori Kabeli<sup>1</sup> Yaniv Shmueli<sup>1</sup> Allan dos Santos Costa<sup>4</sup> Maryam Fazel-Zarandi<sup>1</sup> Tom Sercu<sup>1†</sup> Salvatore Candido<sup>1†</sup> Alexander Rives<sup>1,5‡</sup>



Metagenomic Atlas  
<https://esmatlas.com>

# UCLA workshop on GNNs and molecular science

The screenshot shows the IPAM website's workshop page. At the top, there's a navigation bar with links for Programs, Videos, News, People, Your Visit, About IPAM, Donate, and Contact Us. Below the navigation is a search bar with the placeholder "ENHANCED BY Google". The main content area has a blue header with the text "Workshops" and a sub-header "Learning and Emergence in Molecular Systems". Below the header, it says "JANUARY 23 - 27, 2023". There are four tabs at the top of the content area: "OVERVIEW" (selected), "SPEAKER LIST", "LOGGING", and "APPLICATION & REGISTRATION". The "OVERVIEW" tab contains text about the workshop's purpose and a figure showing a network of atoms and molecules. The "SPEAKER LIST" tab is currently inactive.

<http://www.ipam.ucla.edu/programs/workshops/learning-and-emergence-in-molecular-systems>

## Speaker List

Mohammed AlQuraishi (Harvard Medical School)  
**Xavier Bresson** (National University of Singapore)  
**Steve Brunton** (University of Washington)  
Stefan Chmiela (Technische Universität Berlin)  
Kyunghyun Cho (New York University)  
**Cecilia Clementi** (Freie Universität Berlin)  
Bruno Correia (École Polytechnique Fédérale de Lausanne (EPFL))  
Kyle Cranmer (University of Wisconsin-Madison)  
Payel Das (IBM Research)  
Ron Dror (Stanford University)  
Rafael Gomez-Bombarelli (Massachusetts Institute of Technology)

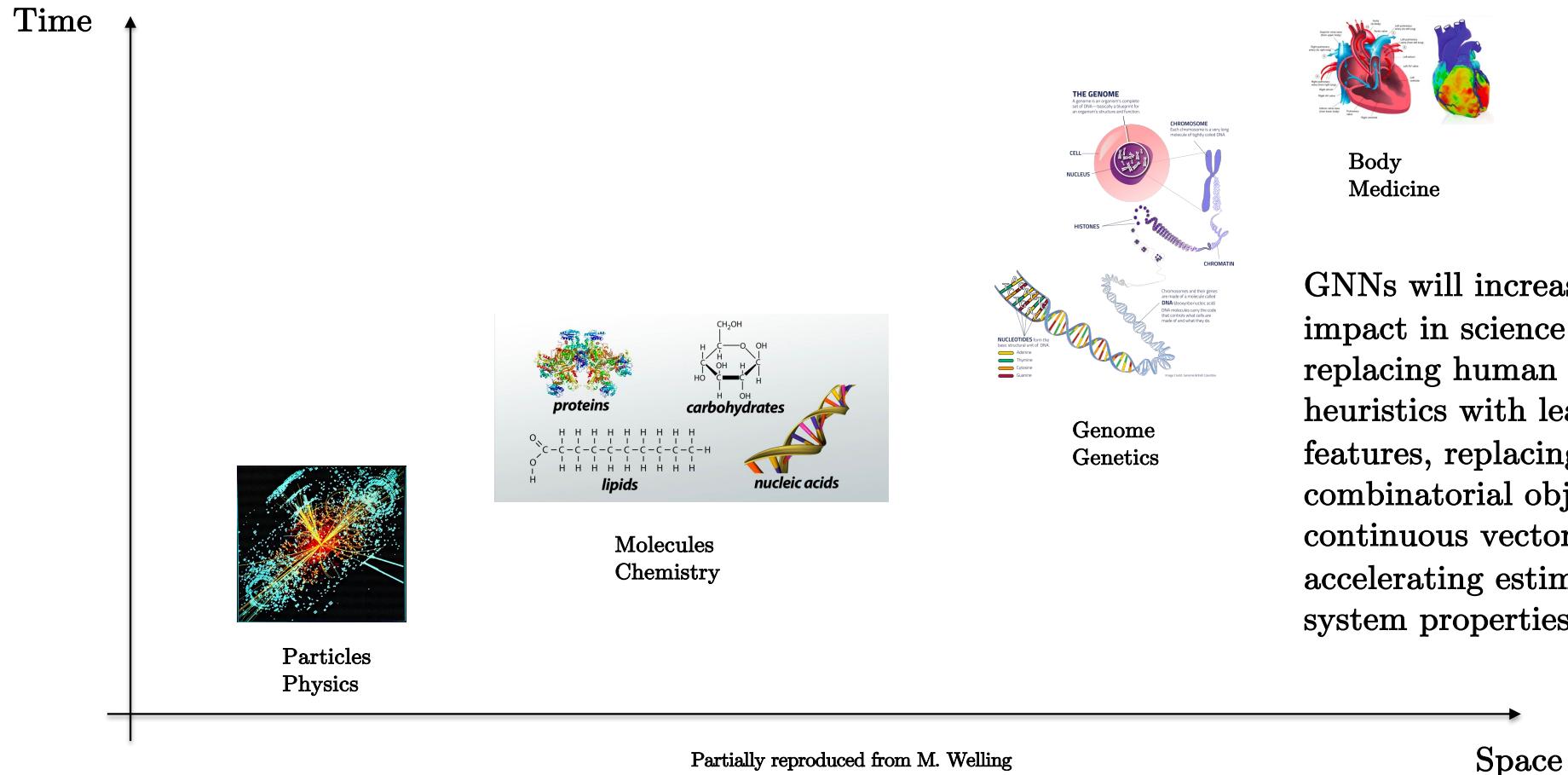
Tommi Jaakkola (Massachusetts Institute of Technology)  
**Frank Noe** (Freie Universität Berlin)  
Kristin Persson (University of California, Berkeley (UC Berkeley))  
**Patrick Riley** (Relay Therapeutics)  
Jutta Rogal (New York University)  
Tess Smidt (Massachusetts Institute of Technology)  
**Alexandre Tkatchenko** (University of Luxembourg)  
**Mark Tuckerman** (New York University)  
Anatole von Lilienfeld (University of Toronto)  
**Max Welling** (Microsoft Research)  
**Rose Yu** (University of California, San Diego (UCSD))  
Maxim Zlatdinov (Oak Ridge National Laboratory)

The screenshot shows a YouTube playlist titled "2023 Learning and Emergence in Molecular Systems" from the channel "Institute for Pure & Applied Mathematics (IPAM)". The playlist contains 21 videos with 743 views and was last updated on Jan 27, 2023. The thumbnails for the first seven videos are displayed:

1. Steve Brunton - Machine Learning for Scientific Discovery, with Examples in Fluid Mechanics
2. Rose Yu - Incorporating Symmetry for Learning Spatiotemporal Dynamics - IPAM at UCLA
3. Tess Smidt - Learning how to break symmetry with symmetry-preserving neural networks - IPAM at UCLA
4. Frank Noe - Advancing molecular simulation with deep learning - IPAM at UCLA
5. Kyunghyun Cho - Lab-in-the-loop de novo antibody design - what are we missing from machine learning?
6. Tommi Jaakola - Diffusion based distributional modeling of conformers, blind docking and proteins
7. Anatole von Lilienfeld - Go EAST, young scientist - First principles view on chemical compound space

<https://www.youtube.com/playlist?list=PLHyI3Fbmv0SeEiaFaO2XHMuHdZnZZiqa8>

# GNNs for science/biology



# Outline

- Deep learning for molecular science
- Graph generative model with VAE
  - Auto-encoder
  - Encoder and decoder
  - System description
  - Numerical experiments
- Conclusion

# Molecule generation

- Goal is to design a neural network that can
  - Auto-encode molecules
  - Generate novel molecules
  - Produce molecules with optimized chemical property

---

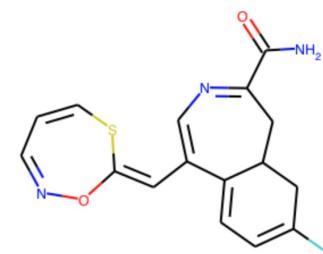
## A Two-Step Graph Convolutional Decoder for Molecule Generation

---

Xavier Bresson  
School of Computer Science and Engineering  
NTU, Singapore  
[xbresson@ntu.edu.sg](mailto:xbresson@ntu.edu.sg)

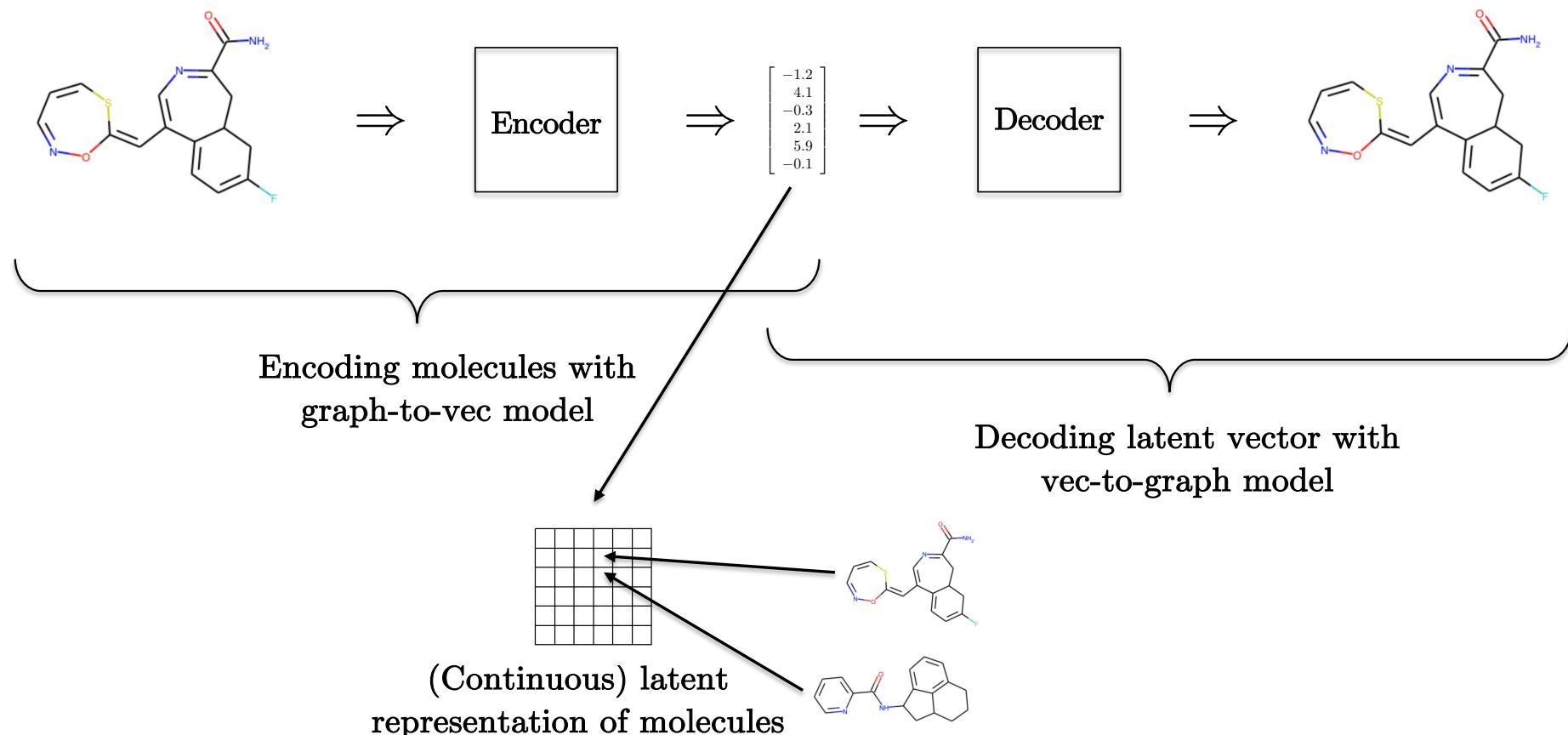
Thomas Laurent  
Department of Mathematics  
Loyola Marymount University  
[tlaurent@lmu.edu](mailto:tlaurent@lmu.edu)

ArXiv : <https://arxiv.org/pdf/1906.03412.pdf>



# Graph auto-encoder

- Graph-to-graph Model :

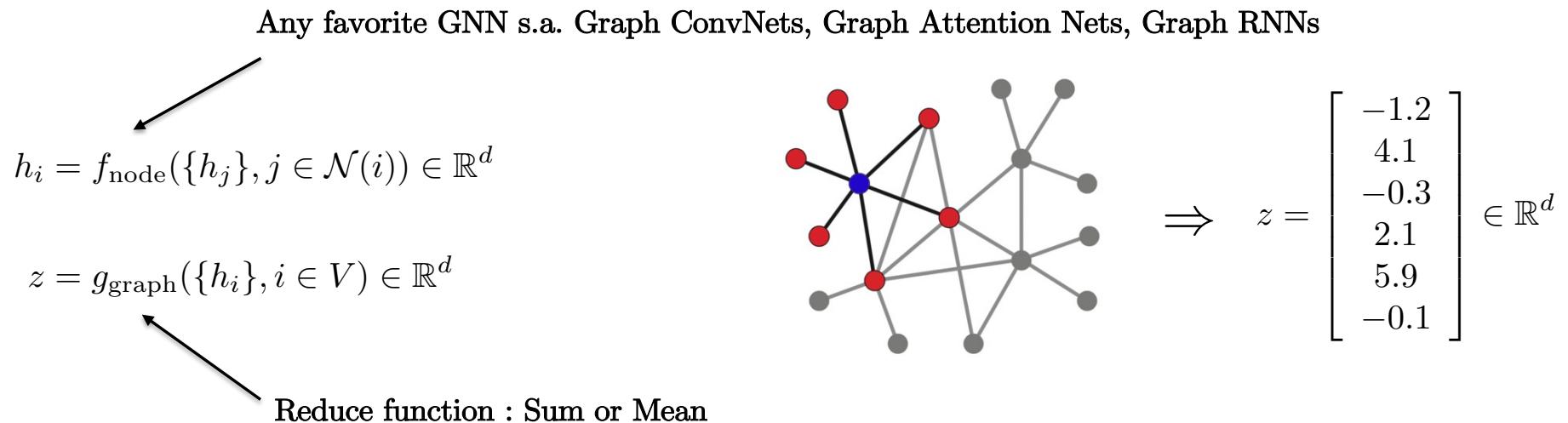


# Outline

- Deep learning for molecular science
- Graph generative model with VAE
  - Auto-encoder
  - Encoder and decoder
  - System description
  - Numerical experiments
- Conclusion

# Graph encoder

- GNNs have been used to encode molecules into a continuous vectorial space.
  - GNNs used for regression<sup>[1,2]</sup> to predict molecular properties (1-2 orders of magnitude faster than solving Schrodinger equation w/ DFT).

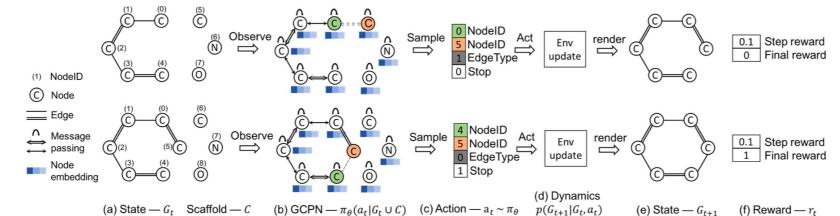


[1] Duvenaud, Maclaurin, Iparraguirre, Bombarell, Hirzel, Aspuru-Guzik, Adams, Convolutional networks on graphs for learning molecular fingerprints, 2015

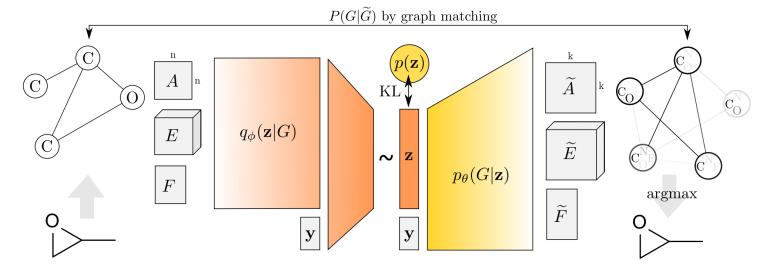
[2] Gilmer, Schoenholz, Riley, Vinyals, Dahl, Neural message passing for quantum chemistry, 2017

# Decoder for graph generation

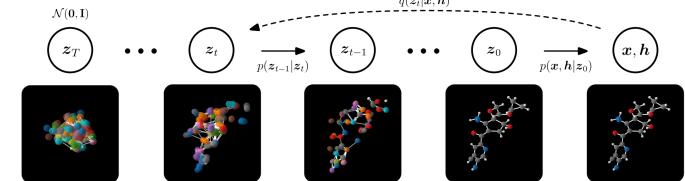
- Encoding is easy. Decoding is the challenging part!
- Three approaches :
  - Step-by-step auto-regressive models : Sequential generation of molecules (atom-by-atom)
    - Jin-et.al, 2018<sup>[1]</sup>, You-Leskovec-et.al, 2018<sup>[2]</sup>, etc
  - One-shot models : Generation of all atoms and bonds in a single pass
    - Simonovsky, Komodakis, 2018<sup>[3]</sup>, De Cao, Kipf, 2018<sup>[4]</sup>, etc
  - Step-by-step generative models : Sequential denoising of molecules
    - Hoogeboom-et.al, 2022<sup>[5]</sup>, Xu-et.al, 2022<sup>[6]</sup>



You-Leskovec-et.al, 2018



Simonovsky, Komodakis, 2018

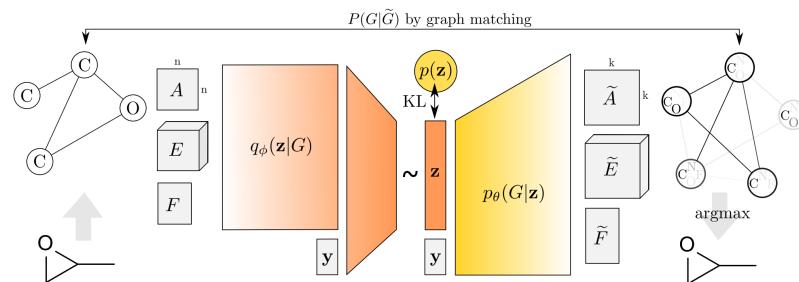


Hoogeboom-et.al, 2022

- [1] Jin, Barzilay, Jaakkola, Junction Tree Variational Autoencoder for Molecular Graph Generation, 2018
- [2] You, Liu, Ying, Pande, Leskovec, Graph convolutional policy network for goal-directed molecular graph generation, 2018
- [3] Simonovsky, Komodakis, GraphVAE: Towards generation of small graphs using variational autoencoders, 2018
- [4] De Cao, Kipf, MolGAN: An implicit generative model for small molecular graphs, 2018
- [5] Hoogeboom, Satorras, Vignac, Welling, Equivariant diffusion for molecule generation in 3d, 2022
- [6] Xu, Yu, Song, Shi, Ermon, Tang, Geodiff: A geometric diffusion model for molecular conformation generation, 2022

# One-shot decoder

- A challenge with one-shot decoder is to generate molecules of different sizes.
  - It is hard to generate simultaneously
    - The number of atoms
    - The type of atoms
    - The bond structures between the atoms
  - Techniques<sup>[1,2]</sup> generated molecules with a fixed size (the size of the largest molecule).

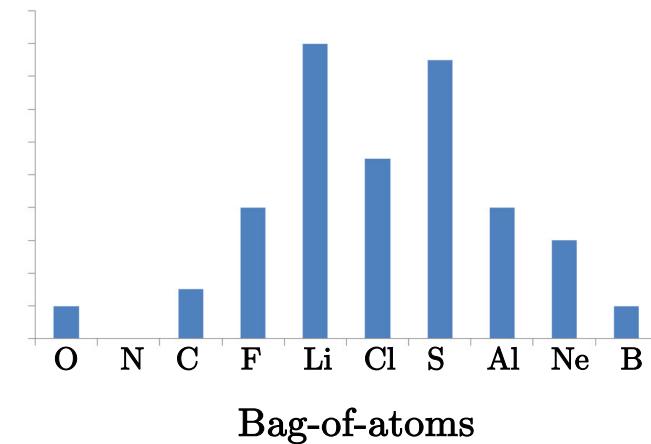
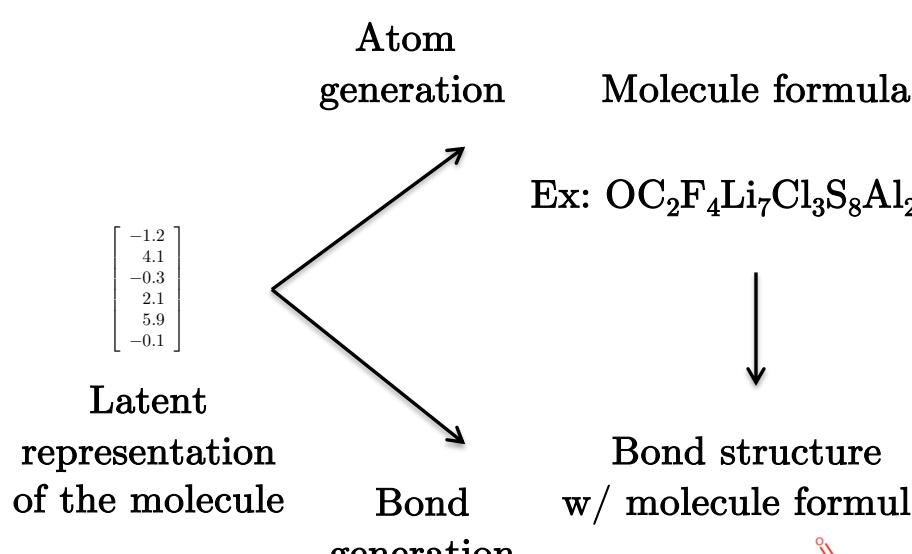


[1] Simonovsky, Komodakis, GraphVAE: Towards generation of small graphs using variational autoencoders, 2018

[2] De Cao, Kipf, MolGAN: An implicit generative model for small molecular graphs, 2018

# Proposed decoder

- We propose to disentangle these 3 problems :



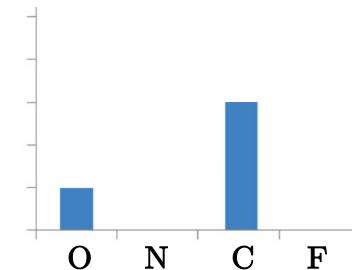
Sum of histogram is the number  
of atoms in a molecule

# Atom decoder

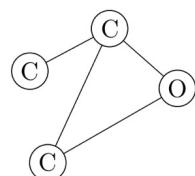
- We decode the latent representation of the molecule with a Multi-Layer Perceptron (MLP) to produce the histogram over the atoms in a molecule :

$$\begin{bmatrix} -1.2 \\ 4.1 \\ -0.3 \\ 2.1 \\ 5.9 \\ -0.1 \end{bmatrix}$$

MLP



Latent  
representation of  
the molecule

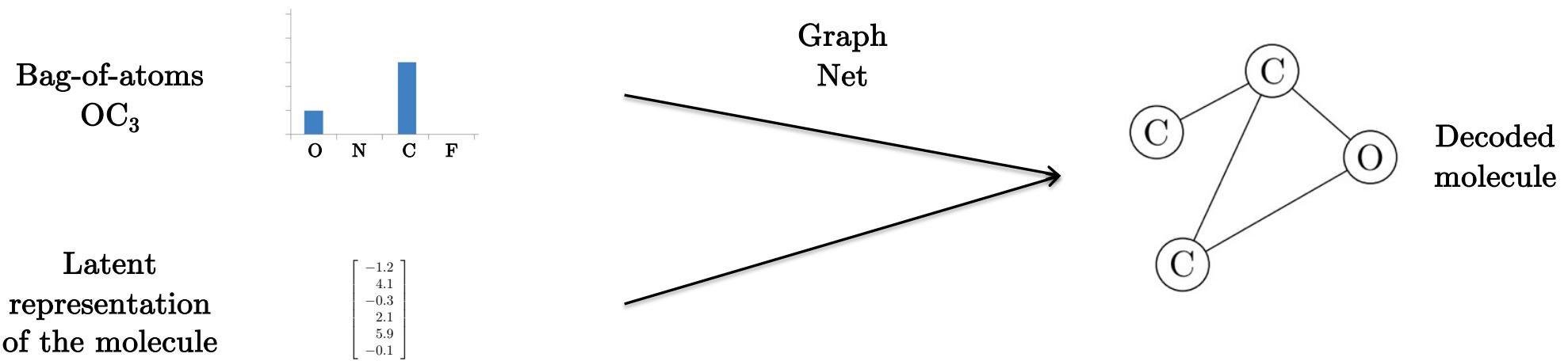


Bag-of-atoms

$\text{OC}_3$

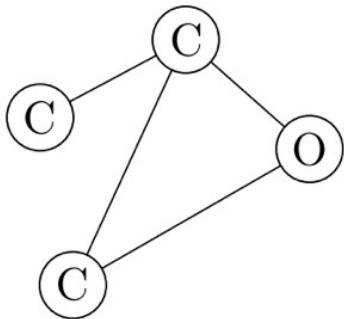
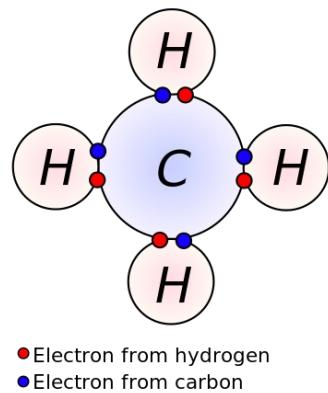
# Bond decoder

- The “IKEA” model :
  - The bag-of-atoms indicates what atoms are in the molecule (IKEA pieces),
  - The atoms are assembled with a GNN (IKEA assembly instructions).



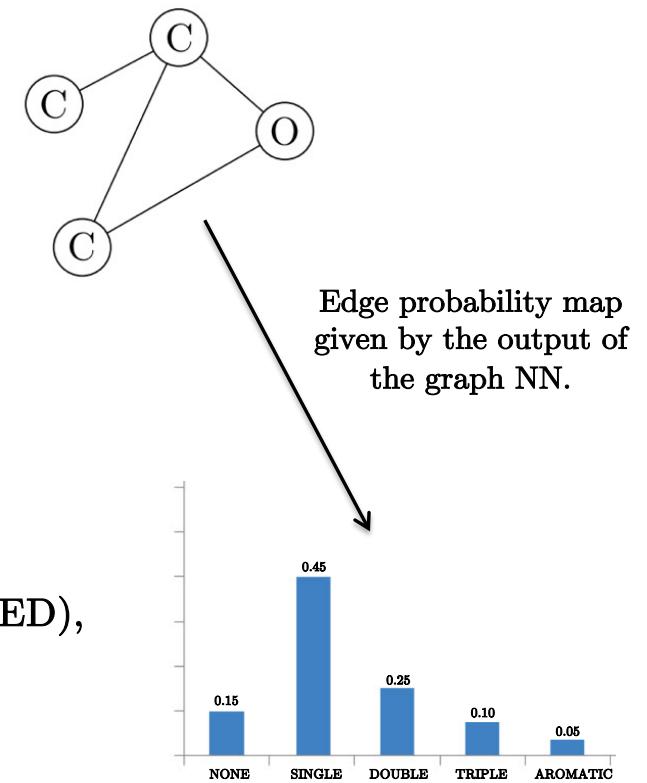
## Beam search for chemical valency

- The one-shot model may not produce a chemically valid molecule.
  - Violation of atom valency (maximum number of electrons in the outer shell of the atom that can participate of a chemical bond).
- We use beam search to produce valid molecules.



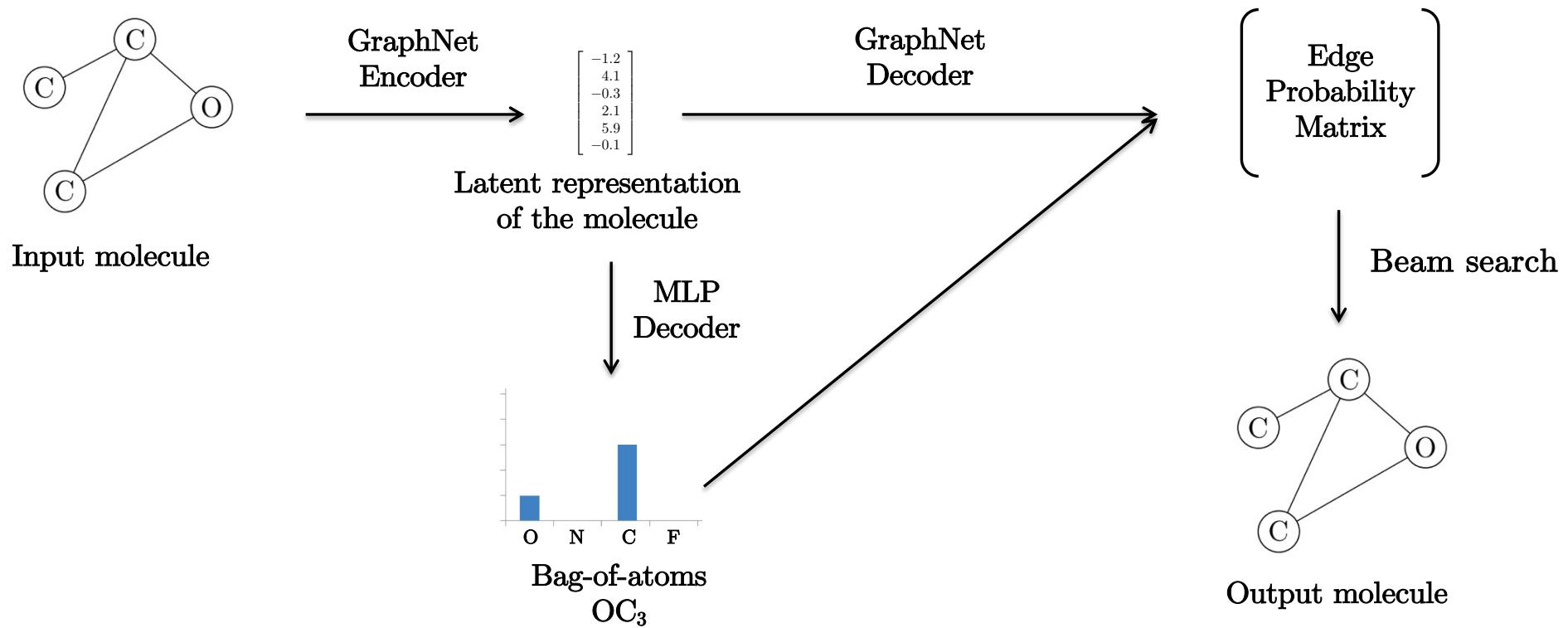
# Beam search for chemical valency

- Beam search :
  - Start with a random edge.
  - Select the next edges that
    - have the largest probability (or Bernouilli sampling),
    - are connected to selected edges,
    - do not violate valency.
- Repeat for a number of different initializations.
- Select the molecule that maximizes
  - The product of edge probabilities or,
  - The chemical property to be optimized s.a. druglikeness (QED), constrained solubility ( $\log P$ ), etc.



## Proposed auto-encoder

- Molecular auto-encoder system :

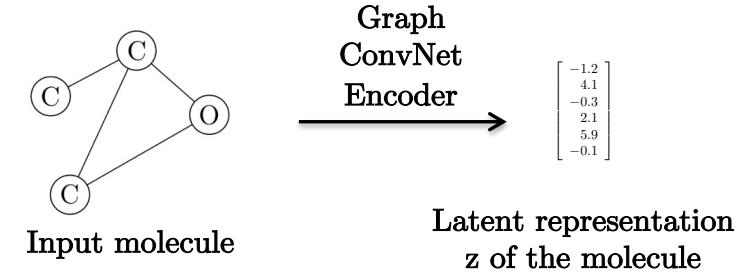


# Outline

- Deep learning for molecular science
- Graph generative model with VAE
  - Auto-encoder
  - Encoder and decoder
  - System description
  - Numerical experiments
- Conclusion

# Encoder description

- We use graph ConvNet<sup>[1]</sup> :



**Node and edge representations**

$$\left\{ \begin{array}{l} h_i^{\ell+1} = h_i^\ell + \text{ReLU}\left(\text{BN}\left(W_1^\ell h_i^\ell + \sum_{j \sim i} \eta_{ij}^\ell \odot W_2^\ell h_j^\ell\right)\right) \quad \text{with} \quad \eta_{ij}^\ell = \frac{\sigma(e_{ij}^\ell)}{\sum_{j' \sim i} \sigma(e_{ij'}^\ell) + \varepsilon} \\ e_{ij}^{\ell+1} = e_{ij}^\ell + \text{ReLU}\left(\text{BN}\left(V_1^\ell e_{ij}^\ell + V_2^\ell h_i^\ell + V_3^\ell h_j^\ell\right)\right) \end{array} \right.$$

**Dense attention**

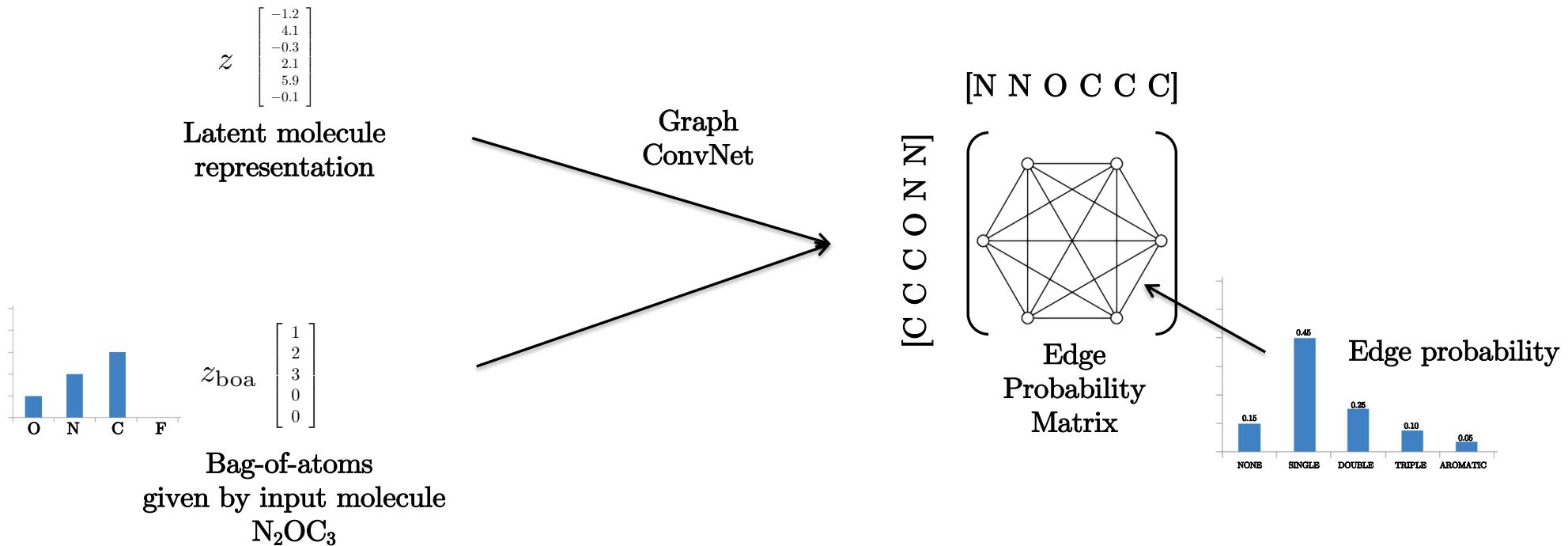
**Graph representation**

$$z = \sum_{i,j=1}^N \sigma(A e_{ij}^L + B h_i^L + C h_j^L) \odot W e_{ij}^L$$

[1] Bresson, Laurent, Residual gated graph convnets, 2017

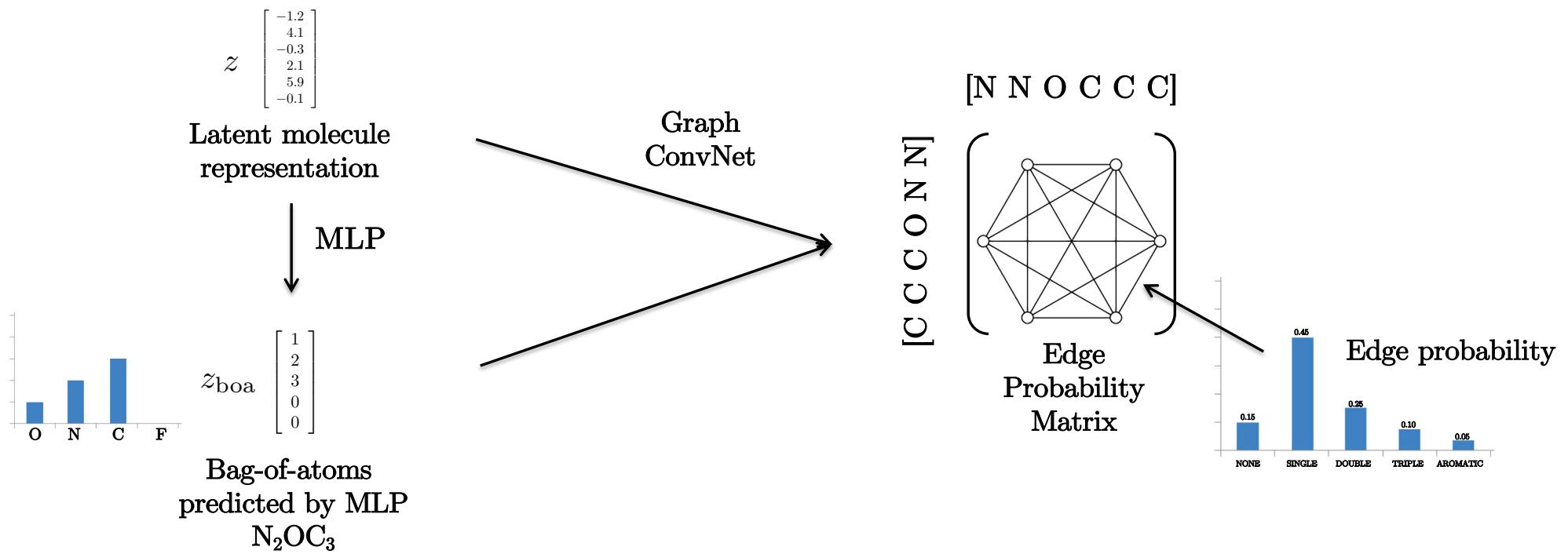
# Bond decoding during training

- Given the latent encoding  $z$  of the molecule and the bag-of-atoms  $z_{boa}$ , we use another graph ConvNet to decode the bonds between the atoms :



## Bond decoding at test time

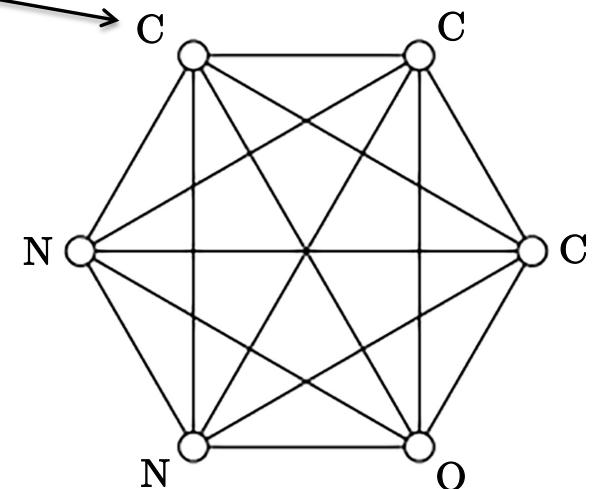
- The bag-of-atoms of the input molecule is predicted by a MLP :



## Breaking atom symmetry

- The bond decoder starts with a fully connected graph with the atom type  $z_{\text{ato}}$  on each node.
- This is not enough for the GNN to be able to differentiate the 3 atoms of Carbon and the 2 atoms of Nitrogen!
  - We break this symmetry by introducing positional features  $z_{\text{pos}}$ , which will differentiate several atoms of the same type.
  - We concatenate this positional feature with the atom type  $z_{\text{ato}}$  to form the input node feature of the decoder.

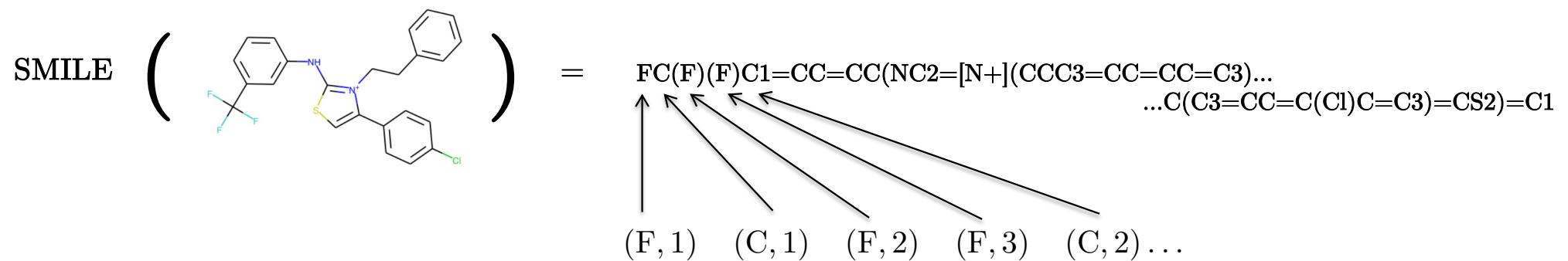
$$h_i^{\ell=0} = [ z_{\text{ato}_i} ]$$



$$h_i^{\ell=0} = \begin{bmatrix} z_{\text{ato}_i} \\ z_{\text{pos}_i} \end{bmatrix}$$

# Atom positional encoding

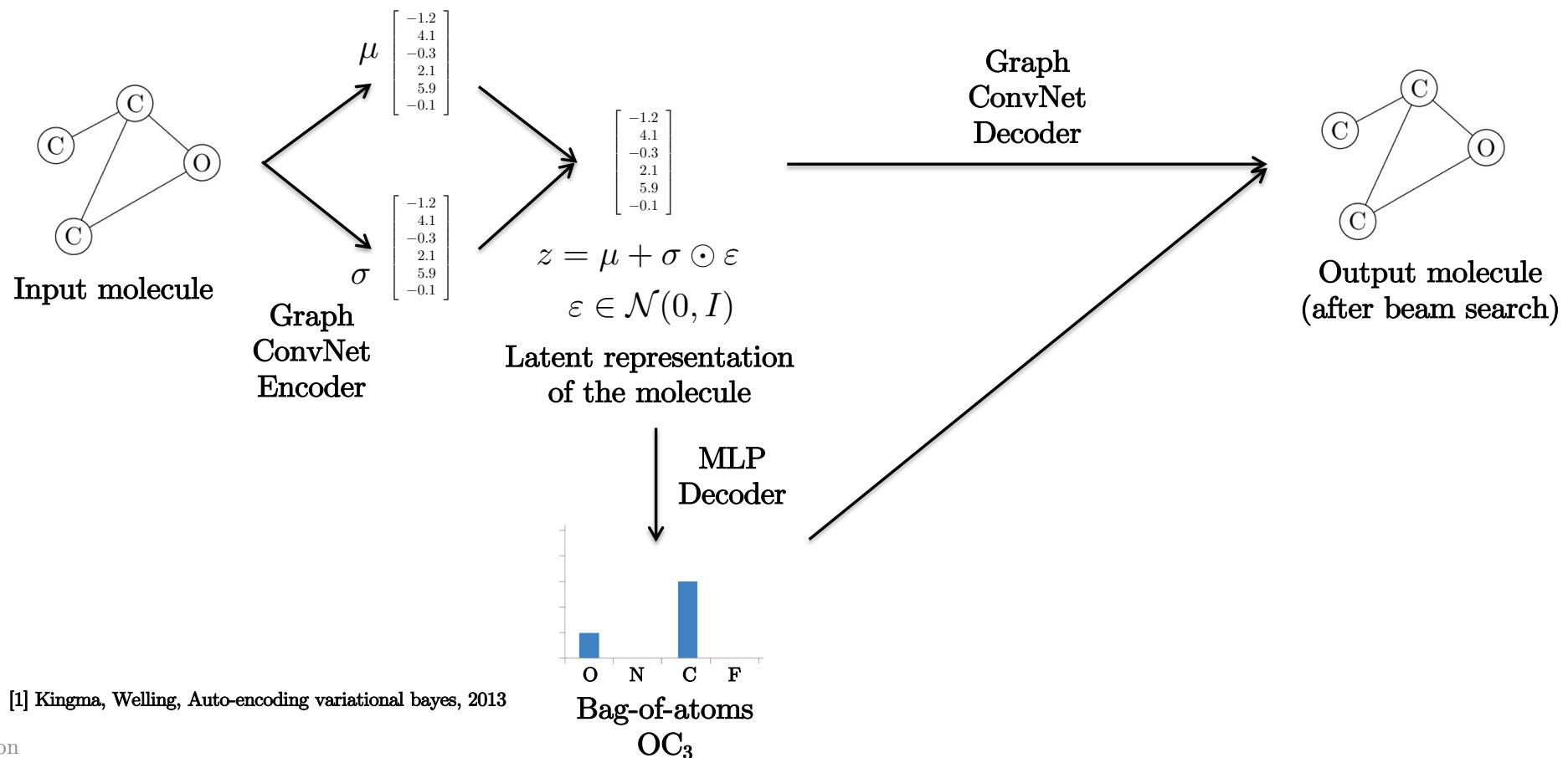
- We need to order the atoms.
- We use the SMILE<sup>[1]</sup> representation of molecules to order the atoms.
  - A SMILE is a sequence of string characters that encodes atoms and bonds of a molecule.



[1] Weininger, SMILES, a chemical language and information system, 1988

# Variational auto-encoder

- We use the VAE formulation<sup>[1]</sup> to improve the molecule generation “by filling the latent space” :



## Molecular generative loss

- The generative loss is composed of
  - Cross-entropy loss for edge probability,
  - Cross-entropy loss for bag-of-atoms probability,
  - Kullback–Leibler divergence for the VAE Gaussian distribution.

$$L = \lambda_e \sum_e \hat{p}_e \log p_e + \lambda_a \sum_a \hat{p}_a \log p_a - \frac{\lambda_{vae}}{2} \sum_k (1 + \log \sigma_k^2 - \mu_k^2 - \sigma_k^2)$$

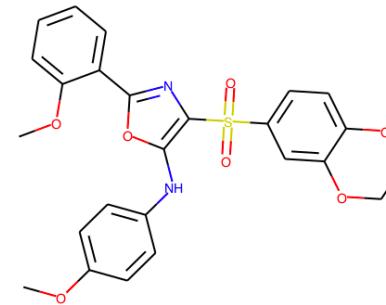
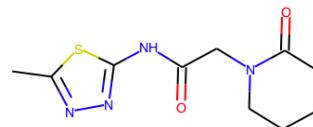
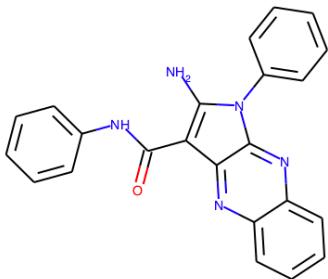
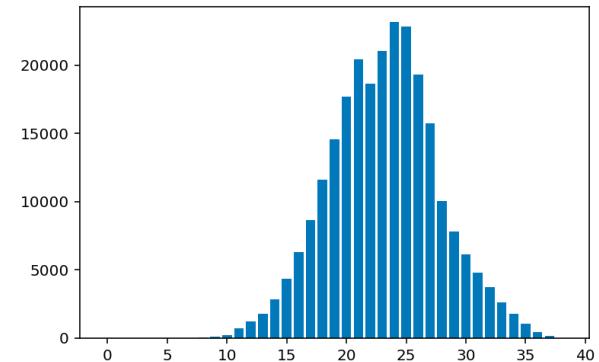
- No matching process is necessary between input and output molecules because the same atom ordering is used (with the SMILE representation).

# Outline

- Deep learning for molecular science
- Graph generative model with VAE
  - Auto-encoder
  - Encoder and decoder
  - System description
  - Numerical experiments
- Conclusion

# Dataset

- ZINC dataset
  - 250k drug-like molecules
  - Up to 38 heavy atoms (excluded hydrogen)



# Training

- Training process
  - Mini-batch of 50 molecules
  - Learning rate is decreased by 1.25 after each epoch if training loss does not decrease by 1%.
  - Learning stops when LR is less than  $10^{-6}$ .
  - Training takes 28 hours on a single Nvidia 1080Ti GPU.

# Numerical experiments

- Molecule reconstruction
  - How many molecules are correctly decoded?
- Molecule novelty
  - Beyond memorization – how many molecules sampled from the learned distribution are not in the training set?
- Molecule optimization
  - How much property improvement can we obtain when optimizing in the latent space?
  - The chemical property is here the constrained solubility of molecules.

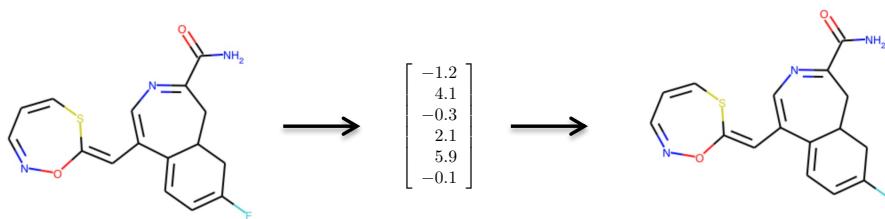
# Baseline techniques

- VAE + SL + AR :
  - JT-VAE : Jin, Barzilay, Jaakkola, Junction Tree Variational Autoencoder for Molecular Graph Generation, 2018
- GAN + RL + AR :
  - GCPN : You, Liu, Ying, Pande, Leskovec, Graph convolutional policy network for goal-directed molecular graph generation, 2018

# Molecule reconstruction

Method	Reconstruction	Validity
CVAE [Gomez-Bombarelli et al., 2016]	44.6%	0.7%
GVAE [Kusner et al., 2017]	53.7%	7.2%
SD-VAE [Dai et al, 2018]	76.2%	43.5%
GraphVAE [Simonovsky, Komodakis, 2018]	-	13.5%
JT-VAE (SL) [Jin et al, 2018]	76.7%	<b>100.0%</b>
GCPN (GAN+RL) [You et al, 2018]	-	-
OURS (VAE+SL)	<b>90.5%</b>	<b>100.0%</b>

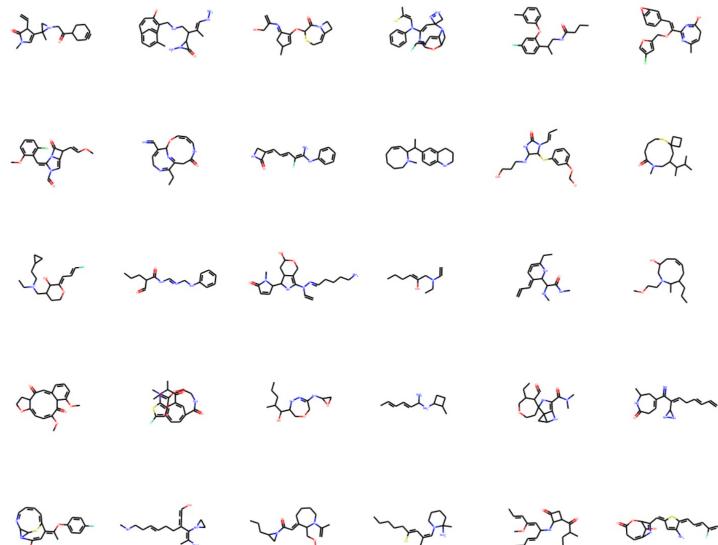
Table 1: Percentage of successful reconstruction of 250k ZINC molecules.



# Molecule novelty

Method	Novelty	Uniqueness
JT-VAE (SL) [Jin et al, 2018]	<b>100.0%</b>	<b>100.0%</b>
GCPN (GAN+RL) [You et al, 2018]	-	-
OURS (VAE+SL)	<b>100.0%</b>	<b>100.0%</b>

Table 2: Sample 5000 molecules from learned prior distribution.



# Molecule optimization task #1

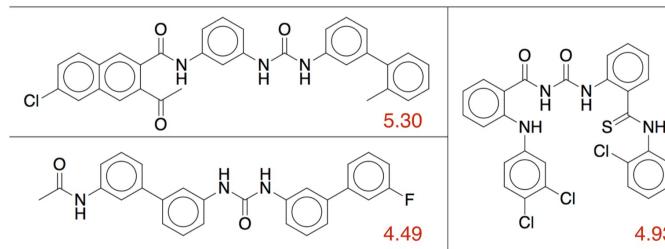
- Molecule optimization :
  - Goal is to maximize the constrained solubility of the training molecules.
  - Optimization is done by gradient ascent in the latent space of molecules.
  - Following JT-VAE, we report the top 3 optimized molecules :

Method	1st	2nd	3rd	Mean
ZINC	4.52	4.30	4.23	4.35
CVAE, Gómez-Bombarelli et al. [2018]	1.98	1.42	1.19	1.53
GVAE, Kusner et al. [2017]	2.94	2.89	2.80	2.87
SD-VAE, Dai et al. [2018]	4.04	3.50	2.96	3.50
JT-VAE, Jin et al. [2018]	5.30	4.93	4.49	4.90
OURS (VAE+SL)	5.24	5.10	5.06	5.14
GCPN (GAN+RL), You et al. [2018]	<b>7.98</b>	<b>7.85</b>	<b>7.80</b>	<b>7.88</b>

# Molecule optimization task #1

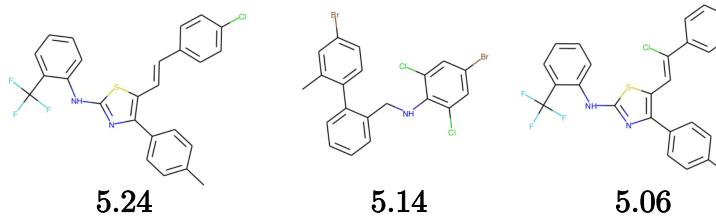
- Top 3 optimized molecules :

JT-VAE (VAE+SL)



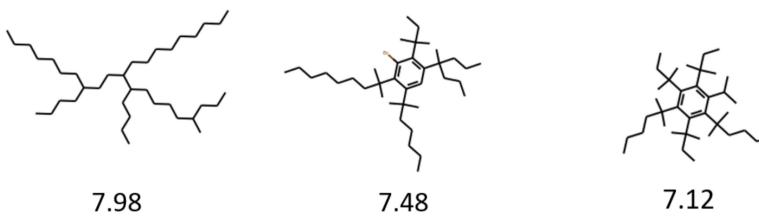
Mean is 4.90

OURS (VAE+SL)



Mean is 5.14

GCPN (GAN+RL)



Mean is 7.52

## Molecule optimization task #2

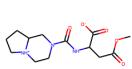
- Constrained optimization :
  - Goal is to maximize the constrained solubility of the 800 test molecules with the lowest value.
  - The optimization of the chemical property is constrained by the similarity between the original molecule and the new generated molecule.
  - Following JT-VAE, we report property improvements w.r.t. molecule similarity  $\delta$  :

$\delta$	JT-VAE [Jin et al, 2018] (SL)			GCPN [You et al, 2018] (GAN+RL)			OURS (VAE+SL)		
	Improvement	Similarity	Success	Improvement	Similarity	Success	Improvement	Similarity	Success
0.0	$1.91 \pm 2.04$	$0.28 \pm 0.15$	97.5%	$4.20 \pm 1.28$	<b><math>0.32 \pm 0.12</math></b>	<b>100.0%</b>	<b><math>5.24 \pm 1.55</math></b>	$0.18 \pm 0.12$	<b>100.0%</b>
0.2	$1.68 \pm 1.85$	$0.33 \pm 0.13$	97.1%	$4.12 \pm 1.19$	<b><math>0.34 \pm 0.11</math></b>	<b>100.0%</b>	<b><math>4.29 \pm 1.57</math></b>	$0.31 \pm 0.12$	98.6%
0.4	$0.84 \pm 1.45$	<b><math>0.51 \pm 0.10</math></b>	83.6%	$2.49 \pm 1.30$	$0.47 \pm 0.08$	<b>100.0%</b>	<b><math>3.05 \pm 1.46</math></b>	<b><math>0.51 \pm 0.10</math></b>	84.0%
0.6	$0.21 \pm 0.71$	<b><math>0.69 \pm 0.06</math></b>	46.4%	$0.79 \pm 0.63$	$0.68 \pm 0.08$	<b>100.0%</b>	<b><math>2.46 \pm 1.27</math></b>	$0.67 \pm 0.05$	40.1%

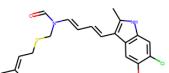
Table 7: Molecule optimization results.

## Molecule optimization task #2

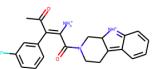
Molecule  
similarity 0.0



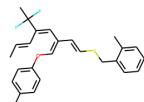
4: -8.38



4: 2.19



77: -5.81

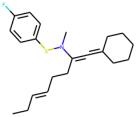


77: 4.75

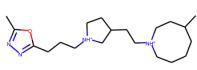
Molecule  
similarity 0.2



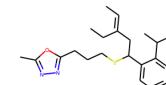
143: -5.01



143: 3.54

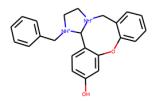


136: -5.06

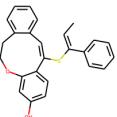


136: 3.10

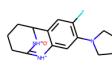
Molecule  
similarity 0.4



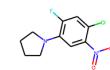
604: -2.94



604: 3.39

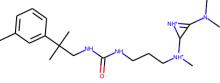


103: -5.40

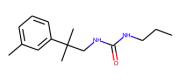


103: 0.88

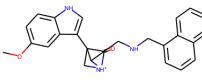
Molecule  
similarity 0.6



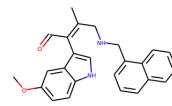
89: -5.64



89: 0.94



782: -2.57



782: 2.44

# Outline

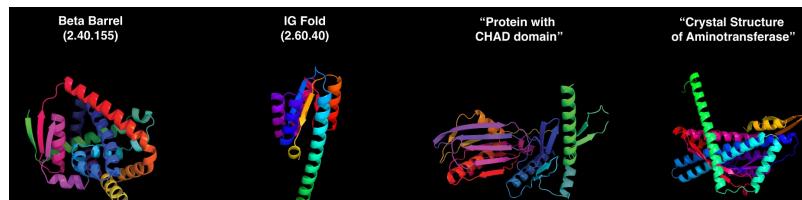
- Deep learning for molecular science
- Graph generative model with VAE
  - Auto-encoder
  - Encoder and decoder
  - System description
  - Numerical experiments
- Conclusion

# Conclusion

- We propose an efficient VAE for atoms and bonds decoding.
  - Single-shot reconstruction + beam search
  - Simple/fast (GPU parallelizable)
  - Alternative to step-by-step techniques
- We report highest VAE accuracy on ZINC dataset for
  - Molecule reconstruction
  - Molecule optimization of constrained solubility property
- Comparing VAE+SL vs GAN+RL :
  - GAN+RL generates better molecules (outside the training statistics)
  - VAE+SL generates better optimized molecules similar to the original ones
  - GAN+RL generates optimized molecules with 100% success
- Limitation
  - Large molecular generation

# Conclusion

- Graph generative models are still in their early development.
- Like for image/text, there is a huge potential to develop pre-trained graph molecular networks from all available datasets and then prompt the generation of new molecules for
  - Drug-like molecules, ligands, proteins, RNAs with optimized chemical, biological, potency, toxicity, pharmacokinetics properties
  - Material science with optimized mechanical, thermal, electrical, optical, chemical and environmental properties



Prompt: Generate a protein with CHAD domain

The end / selfie time



