

NVIDIA: GTC 2019 Generative Molecular DLNN

Ellen Du, Joey Storer, and Abe Stern*

*NVIDIA

The Dow Chemical Company





DOW TEAM: Ellen Du, Joey Storer Sukrit Mukhopadhyay Matthew Christianson Hein Koelman, Ryan Marson William Edsall, Bart Rijksen Jonathan Moore Christopher Roth, Peter Margl Clark Cummins, Dave Magley



NVIDIA TEAM/Alumni: **Abe Stern**, Michelle Gill, John Ashley, Alex Volkov



Outline

Problem statement

Efforts in generative molecular deep learning methods

Our approach

- Hardware/software
- Tooling
- Data curation
- Model Training and convergence
- Latent space analysis and inference
- Generative capability evaluation



Can a molecular generative deep learning system be trained to deliver new molecular designs relevant to our research needs?



Introduction: Generative Molecular Systems

Challenges:

- Molecular encoding (Canonical SMILES)
- Molecular descriptors (100's)
- Vastness of chemical search space (10⁶⁰)
- Unknown structure/property relationships f(n)
- Promise of the latent space dimensionality (32-bit)
- Limits on data set used for training (ChEMBL, ZINC)
- Organization of target properties within the latent space (AlogP)
- Molecule discovery workflow (post-filtering)



Attraction of Molecular VAE/GANs

Convert discrete molecules to continuous latent representations

- Molecules are discrete entities
- Subtle molecular transformations have large differences in performance

Undocumented benefit to using negative data in ml/dl

- Availability of a molecular structure axis in DL that is not generally available to ML
- Tendency in science to "move on" relative to negative or poor results



Gomez-Bombarelli, et al., ACS Cent. Sci., 2018, 4 (2), pp 268–276



General intro on methods: VAEs

Generally there are numerous methods appearing in the open literature:

- Chemical VAE
- Grammar VAE
- Junction Tree
- ATNC RL
- FC-NN (NVIDIA-Dow)

The best way to go is not entirely clear.

Junction Tree – may be best because of the more natural graph representation – but it may constrain diversity

FC-NN is potentially more efficient.



Inferencing Comparison to Literature

Method	Reconstruction	Validity
	Knowns	Inferenced(unknown)
Chem-VAE	44 %	1 % ^{lit.}
Dow-Chem-VAE	94 %	10 %
Grammar-VAE	54 %	7 % ^{lit.}
SD-VAE	76 %	44 % ^{lit.}
Graph-VAE	-	14 % ^{lit.}
JT-VAE	77 %	100 % ^{lit.}
Dow-FC-NN	90 %	%
ATNC-RL	-	71 % ^{lit.}



Models and Training Details



Model Details: Architectures Explored



Model Details: Differences In Inputs





Model Details: Differences In Sequence Modeling





Model Training Details: Data Compilation







Model Training Details: Hardware

NVIDIA DGX-1





ATTIDIE DOTE 1	CI
NVIDIA DGX-I	Specifications
It i Duit Don 1	opeometation

CPUs	2x Intel Xeon E5-2698 v3 (16 core, Haswell-EP)
GPUs	8x NVIDIA Tesla P100 (3584 CUDA Cores)
System Memory	512GB DDR4-2133 (LRDIMM)
GPU Memory	128GB HBM2 (8x 16GB)
Storage	4x Samsung PM863 1.92TB SSDs
Networking	4x Infiniband EDR 2x 10GigE
Power	3200W
Size	3U Rackmount
GPU Throughput	FP16: 170 TFLOPs FP32: 85 TFLOPs FP64: 42.5 TFLOPs



Model Training Details: Software Environment

Container: Docker container

Standard Lightweight Secure



Packages

Chemistry: RDKit, DeepChem *Data Processing:* Numpy, Pandas, Rapids *ML/DL:* SciKitLearn, Keras, Tensorflow, Pytorch, XGBoost *Tuning/Scaling Up:* Hyperopt, Horovod



Model Training Details: Hyperparameter Optimization







Model Training Details: Distributed Model Training

Uber Horovod





- Data Parallelism
- Network Optimal
- User friendly



Model Training Details: Latent Space Organization



Dow

Generative Capability Evaluation



Hit Rate Analysis (> 0 hits/1000 attempts)

ChEMBL TEST = 11800 test molecules inferenced (1000 attempts)

Model	Hit Rate
C-VAE	
55550	94.4 %
JT-NN	
7	100 %
FC-NN	
14587*	94 %



VAE Hit rate: Molecules that never decoded

Analysis of molecules from ChEMBL-TEST (655) that did not decode with 1000 attempts:

- 1. SMILES string length distribution for the non-decoding molecules
- 2. Inference study increased to 10,000 attempts/molecule
 - a. 549/655 still never decoded
 - b. 16 % successful decoded at least once on 10,000 additional attempts
 - c. One molecule decoded an additional 44 times



Distribution is not remarkable compared to ChEMBL





Distribution of SMILES string lengths



Dow

Distribution of SMILES string length w.r.t model type



E. Putin, et al., *Mol. Pharmaceutics 2018, 15, 4386-4397*

22

Distance calculation and performance

<u>GPU enabled–distance matrix calculation:</u>

- 1. Characterizing latent space
- 2. Support inferencing
 - a. Nearest neighbor analysis
 - b. Gaussian process support

Rough method comparison: (30,000 molecules, 900 x 10⁶ distances) Python (Simple, non-vectorized) 5 x 10⁵ (DGX-1) Scipy.spatial.distance.euclidean 10⁴ (DGX-1) Numba/CUDA 1 (DGX-1)



Latent Space Vectors (Kernel Density Est) C-VAE, JT-NN, FC-NN







How far apart are the molecules in the Latent Space?



25

Interpolation from the Latent Space

Linear interpolation

 Stepping through training set and linearly interpolating between endpoints chosen from the training set

Spherical-linear interpolation

 Stepping through training set and spherical-linearly interpolating between endpoints chosen from the training set

Hyperspheres

 Utilizing the distance matrix to select point for expanding hyperspheres





Molecular Interpolation in a Continuous Design Space

LERP/SLERP

Algorithm only chose points across the whole of the training set (118,000 molecules) and then interpolated between points in ranges to ensure that, at a minimum, each molecule became an end-point for interpolation





Inferencing followed by molecular filtering





Synthetic Accessibility Score



The SAScore across:

INPUT: ChEMBL (118,000) OUTPUT: Inferenced_e55500 TEST: Dow

- Ertl, P.; Schuffenhauer, A., J. Cheminf. 2009, 1:8
- Ertl, P.; Landrum, G. <u>https://github.com/rdkit/rdkit/tree/master/Contrib/SA_Score</u>





C-VAE

Chem-VAE modeled after Bombarelli works better than reported and delivers good molecules. The time/epoch is high and the number of epochs needed is ~ 50,000.

JT-NN

Junction Tree converges faster, is a more natural representation of molecules, and delivers good molecules.

FC-NN

Fully Convolutional works well, converges faster than C-VAE, and delivers good molecules.





END

The Dow Chemical Company