David Ryan Koes June 25, 2025 **BioInformatics in Torun**



Deep Learning for Structure-Based Drug Discovery: From Scoring to Generative Design

University of
Dittsburgh



Structure Based Drug Discovery





GNINA 1.0

https://github.com/gnina/gnina





SOFTWARE

Open Access

GNINA 1.0: molecular docking with deep learning

Andrew T. McNutt¹, Paul Francoeur¹, Rishal Aggarwal², Tomohide Masuda¹, Rocco Meli³, Matthew Ragoza¹, Jocelyn Sunseri¹ and David Ryan Koes^{1*}







Protein-Ligand Scoring



Pose Prediction

Binding Discrimination

Affinity Prediction



Protein-Ligand Scoring



Pose Prediction

Binding Discrimination

Affinity Prediction



Protein-Ligand Representation



(R,G,B) pixel



Protein-Ligand Representation



Matthew Ragoza,^{†,‡} Joshua Hochuli,^{‡,¶} Elisa Idrobo,[§] Jocelyn Sunseri,[∥] and David Ryan Koes^{*,∥}©

(R,G,B) pixel \rightarrow

(Carbon, Nitrogen, Oxygen,...) **voxel**

The only parameters for this representation are the choice of grid resolution, atom density, and atom types.

JOURNAL OF CHEMICAL INFORMATION AND MODELING

Article

pubs.acs.org/jcim

Protein–Ligand Scoring with Convolutional Neural Networks







Convolutional Neural Networks



Convolutional Filters



-1	-1	-1
0	0	0
1	1	1

-1	0	1	-1	-1	
-1	0	1	-1	8	
-1	0	1	-1	-1	





Protein Ligand Scoring





684,640 Parameters



Default2017

Default2018

HiRes Affinity

0.58

9

し

onvolutio

U

m

3x3x

BatchN

0.56

HiRes Pose





Cross-Docked Protein Ligand Scoring

Three-Dimensional Convolutional Neural Networks and a Cross-**Docked Data Set for Structure-Based Drug Design**

Paul G. Francoeur, Tomohide Masuda, Jocelyn Sunseri, Andrew Jia, Richard B. Iovanisci, Ian Snyder, and David R. Koes*





Clustered Cross-validation





18,450 complexes 22.6 million poses









Docking Performance







)



DUD-E Virtual Screening Performance



Open Access Article

Virtual Screening with GNINA 1.0



by 😫 Jocelyn Sunseri 🖂 问 and 😫 David Ryan Koes * 🖂 问

DUD-E			LIT-PCBA		
AUC	NEF1%	EF1%	AUC	NEF1%	EF1
0.683	0.0514	3.02	0.6	0.013	1.2
0.963	0.857	51.9	0.542	0.00733	0.7
0.745	0.118	7.05	0.581	0.011	1.
0.764	0.187	11.4	0.577	0.0103	0.9
0.756	0.179	11.6	0.579	0.037	2.0
0.702	0.156	10.3	0.498	0.0147	1.
0.795	0.27	17.7	0.616	0.037	2.5
0.767	0.313	20.4	0.514	0.0238	1.8
0.795	0.258	15.6	0.611	0.0238	1.8
0.744	0.241	15.8	0.512	0.0147	1.4
	AUC 0.683 0.963 0.745 0.764 0.756 0.702 0.795 0.767 0.795 0.795 0.744	DUD-EAUCNEF1%0.6830.05140.9630.8570.7450.1180.7640.1870.7560.1790.7020.1560.7950.270.7670.3130.7950.2580.7440.241	DUD-EAUCNEF1%EF1%0.6830.05143.020.9630.85751.90.7450.1187.050.7640.18711.40.7560.17911.60.7020.15610.30.7950.2717.70.7670.31320.40.7950.25815.60.7440.24115.8	DUD-EAUCNEF1%EF1%AUC0.6830.05143.020.60.9630.85751.90.5420.7450.1187.050.5810.7640.18711.40.5770.7560.17911.60.5790.7020.15610.30.4980.7670.2717.70.6160.7950.25815.60.6110.7440.24115.80.512	DUD-ELIT-PCBAAUCNEF1%EF1%AUCNEF1%0.6830.05143.020.60.0130.9630.85751.90.5420.007330.7450.1187.050.5810.0110.7640.18711.40.5770.01030.7560.17911.60.5790.0370.7020.15610.30.4980.01470.7670.31320.40.5140.02380.7950.25815.60.6110.02380.7440.24115.80.5120.0147







GNINA 1.3 https://github.com/gnina/gnina

GNINA 1.3: the next increment in molecular docking with deep learning

Andrew T. McNutt, Yanjing Li, Rocco Meli, Rishal Aggarwal & David Ryan Koes

Journal of Cheminformatics 17, Article number: 28 (2025) **Cite this article**



Caffe \rightarrow Torch easy covalent docking retrained models

GNINA 1.3 Performance

GNINA vs End-to-end Deep Docking

GNINA DynamicBind Boltz1

GNINA vs End-to-end Deep Docking

GNINA DynamicBind Boltz1

Generative Modeling

Discriminative Model

Features X

Generative Model

Features X

Computational and Systems Biology

Generative Model

Computational and Systems Biology

Generative Model

Computational and Systems Biology

Features **X**

Learning a Continuous Representation of 3D Molecular Structures with Deep Generative Models

Matthew Ragoza*

Comp. & Systems Biology University of Pittsburgh Pittsburgh, PA 15213 mtr22@pitt.edu

Tomohide Masuda* Comp. & Systems Biology University of Pittsburgh Pittsburgh, PA 15213 tmasuda@pitt.edu

David Ryan Koes dkoes@pitt.edu

Comp. & Systems Biology University of Pittsburgh Pittsburgh, PA 15213

NeurIPS 2020 Workshop Machine Learning for Structural Biology

Variational Autoencoding Examples

Atom Fitting

Computational and Systems Biology

Computational and Systems Biology

Generative Adversarial Networks

True Examples

Generator

Generative Adversarial Networks

True Examples

Generator

4.5 years of GAN progress on face generation. arxiv.org/abs/1406.2661 arxiv.org/abs/1511.06434 arxiv.org/abs/1606.07536 arxiv.org/abs/1710.10196 arxiv.org/abs/1812.04948

Computational and Systems Biology

University of Pittsburgh

0.5

1.0

2.0 Variability factor

Diffusion Models

STRUCTURE-BASED DRUG DESIGN WITH EQUIVARIANT DIFFUSION MODELS

Arne Schneuing^{1*}, Yuanqi Du^{2*}, Charles Harris³, Arian Jamasb³, Ilia Igashov¹, Weitao Du⁴, Tom Blundell³, Pietro Lió³, Carla Gomes², Max Welling⁵, Michael Bronstein⁶ & Bruno Correia¹

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

Quantitative Biology > Biomolecules

[Submitted on 4 Oct 2022 (v1), last revised 11 Feb 2023 (this version, v2)]

DiffDock: Diffusion Steps, Twists, and Turns for Molecular Docking

Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, Tommi Jaakkola

Diffusion Models

STRUCTURE-BASED DRUG DESIGN WITH EQUIVARIANT DIFFUSION MODELS

Arne Schneuing^{1*}, Yuanqi Du^{2*}, Charles Harris³, Arian Jamasb³, Ilia Igashov¹, Weitao Du⁴, Tom Blundell³, Pietro Lió³, Carla Gomes², Max Welling⁵, Michael Bronstein⁶ & Bruno Correia¹

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

$\exists r \times lV > q$ -bio > arXiv:2210.01776

Quantitative Biology > Biomolecules

[Submitted on 4 Oct 2022 (v1), last revised 11 Feb 2023 (this version, v2)]

DiffDock: Diffusion Steps, Twists, and Turns for Molecular Docking

Gabriele Corso, Hannes Stärk, Bowen Jing, Regina Barzilay, Tommi Jaakkola

Keypoint Conditioned Diffusion

Ian Dunn, David Ryan Koes

[Submitted on 22 Nov 2023 (v1), last revised 8 May 2024 (this version, v2)]

Accelerating Inference in Molecular Diffusion Models with Latent **Representations of Protein Structure**































- GVP all-atom
- EGNN all-atom
- GVP keypoints
- EGNN keypoints
- GVP C_α
- EGNN C_{α}



Generated molecules are not synthetically accessible or physically plausible

Reference Molecules







Generated Molecules















Generated molecules are not synthetically accessible or physically plausible

Reference Molecules







Generated Molecules



"epoxide, good luck with stability"











Generated molecules are not synthetically accessible or physically plausible

Reference Molecules









"certainly unstable"



Generated molecules are not synthetically accessible or physically plausible

Reference Molecules









"certainly unstable"



Generated molecules are not synthetically accessible or physically plausible

Reference Molecules









"certainly unstable"



Practical Measures of Molecule Quality

Existing literature primarily focuses on validity/valency; necessary but insufficient dimensions of molecule quality



We propose to evaluate molecule quality at the level of functional groups and ring systems





Practical Measures of Molecule Quality



ate of OOD Ring Systems		
P2M Data Source	DiffSBDD	



Unconditional Generation with FlowMol











ΗΝΟ

Ian Dunn **David Ryan Koes** Dept. of Computational & Systems Biology Dept. of Computational & Systems Biology University of Pittsburgh University of Pittsburgh Pittsburgh, PA 15260 Pittsburgh, PA 15260 ian.dunn@pitt.edu dkoes@pitt.edu

Exploring Discrete Flow Matching for 3D De Novo Molecule Generation



FlowMol v3



- State of the art validity \bullet
- Improves chemical plausibility and synthetic accessibility ullet



FlowMol

Mixed continous/categorical flow-matching model for de novo molecule generation.

🔵 Python 🛣 105 🖌 5







FlowMol v3



- State of the art validity \bullet
- Improves chemical plausibility and synthetic accessibility ullet



FlowMol

Mixed continous/categorical flow-matching model for de novo molecule generation.

🔵 Python 🛣 105 🖌 5







OMTRA: A Multi-Task Generative Model for SBDD





Individual Modalities









OMTRA: A Multi-Task Generative Model for SBDD

Individual Modalities





Composed into Many Tasks









Generating *Descriptions* of Molecules









22 May 2025, Version 1

Emma Flynn 💿, Riya Shah 💿, Ian Dunn 💿, Rishal Aggarwal, David Koes 💿







22 May 2025, Version 1

Emma Flynn 💿, Riya Shah 💿, Ian Dunn 💿, Rishal Aggarwal, David Koes 💿

DeepFrag: a deep convolutional neural network for fragment-based lead optimization[†]

Harrison Green, (D^a David R. Koes^b and Jacob D. Durrant (D^{*a})

https://durrantlab.pitt.edu/deepfrag/



Chemical Science



Prospective Evaluation





CACHE Challenge #1





Ian Dunn

CRITICAL ASSESSMENT OF COMPUTATIONAL HIT-FINDING EXPERIMENTS





Large-Scale Docking with GNINA



A Tale of Two Methods

Pharmacophore Screening with Pharmit





Large-Scale Docking with GNINA



A Tale of Two Methods

Pharmacophore Screening with Pharmit





molport



2 screening methods 2 scoring methods

1k ligands gnina scores vina scores

Round 1 Submission



- ZINC20: 20 mil molecules
- MCULE: 45 mil molecules
- MCULE-ULTIMATE: 126 mil molecules

Molecule Libraries



3.5k ligands gnina scores vina scores





Round 1 Results

- Selection limited/ skewed by database availability
- 84 ligands tested
 - 59 from docking
 - 24 from pharm screen





Round 1 Results

- Selection limited/ skewed by database availability
- 84 ligands tested
 - 59 from docking
 - 24 from pharm screen





Round 1 Results

- 2/84 were hits
 - Both from docking •



Round 2: Hit Optimization



Hit Optimization Pipeline

Parent Compound

Similarity screen against Enamine REAL Return 5000 most similar ligands by tanimoto score



Crystal Structure





Hit Optimization Results



Parent Compound











Participant

David Koes, University of Pittsburgh

Olexandr Isayev & Maria Kurnikova, Carne Cherkasov, University of British Columbia

Christina Schindler, Merck KGaA

Dmitri Kireev, University of Missouri

Christoph Gorgulla, St. Jude Children's Res University

Didier Rognan, Université Strasbourg

Pavel Polishchuk, Palacky University

Kam Zhang, Centre for Biosystems Dynan

Shuangjia Zheng, Shanghai Jiao Tong Uni

Carlos Zepeda, Treventis/UHN

Fabian Liessmann, Leipzig University

Final Results

	Participant ID	Aggregated score
	1181	18
egie Mellon University & Artem	1209	18
	1193	17
	1183	16
search Hospital and Harvard	1195	16
	1202	16
	1210	16
nic Research, RIKEN	1188	15
iversity (previously Galixir)	1187	14
	1200	14
	1201	14
	1179	13



CACHE Challenge #2

- RNA binding site of SARS-COV2 NSP13
- "Deep Docking" of Enamine (4B)



Train XGBoost model to predict docking scores from ligand fingerprint

Predict docking scores for entire database





_
- 1
_
- 1
_ 1
 - 1
I
_

CACHE Challenge #2

- RNA binding site of SARS-COV2 NSP13
- "Deep Docking" of Enamine (4B)



Train XGBoost model to predict docking scores from ligand fingerprint

Predict docking scores for entire database





_ 1
- 1
_
- 1
_ 1
 - 1
I
_
5/50 compounds identified as potential hits >2x the average hit rate 4/5 hits from last round of active learning Highest affinity round 1 hit in the competition





Spill







CACHE #2 Results











5/50 compounds identified as potential hits >2x the average hit rate 4/5 hits from last round of active learning Highest affinity round 1 hit in the competition





Spill







CACHE #2 Results











Scalable Screening of Ultra Large Libraries





Quantitative Biology > Biomolecules

[Submitted on 23 Nov 2024 (v1), last revised 20 Jan 2025 (this version, v2)]

Andrew T. McNutt, Abhinav K. Adduri, Caleb N. Ellington, Monica T. Dayao, Eric P. Xing, Hosein Mohimani, David R. Koes

Learn co-embedding of ligands and proteins

Ligands are "close" to proteins they bind to

Scaling Structure Aware Virtual Screening to Billions of Molecules with SPRINT



Enhanced Deep Docking with SPRINT



- Cannabinoid







https://github.com/abhinadduri/panspecies-dti



https://github.com/abhinadduri/panspecies-dti





















Rigorous evaluation is essential







No replacement for prospective evaluation



Rigorous evaluation is essential





Compounds



Acknowledgements





Current Ian Dunn Emma Flynn Riya Shah Rishal Aggarwal **Drew McNutt** Daniel Penaherrera Jacky Chen Somayeh Pirhadi Fareeda Abu-Juam Ben Krummenacher

Previous

Jocelyn Sunseri Matthew Ragoza Tomohide Masuda Paul Francoeur Jonathan King Rocco Meli Josh Hochuli Elisa Idrobo Lily Turner Alec Helbling Andrew Jia Rich Iovanisci Ian Snyder Nick Rego

R01GM108340

National Institute of General Medical Sciences R35GM140753

CHE-1800435



Shameless Plugs



gnina

https://github.com/gnina

Pharmit



https://pharmit.csb.pitt.edu

FlowMol





https://github.com/Dunni3/FlowMol

@dkoes.compstruct.org

https://gina.github.io/libmolgrid/

py3Dmol

[] import py3Dmol

[] p = py3Dmol.view(query='mmtf:lycr') p.setStyle({'cartoon': {'color':'spectrum'}})





Shameless Plugs



gnina

(pronounced NEE-na)

https://github.com/gnina

Pharmit



https://pharmit.csb.pitt.edu







FlowMol

https://github.com/Dunni3/FlowMol

@dkoes.compstruct.org

https://gina.github.io/libmolgrid/

py3Dmol

[] import py3Dmol

[] p = py3Dmol.view(query='mmtf:lycr') p.setStyle({'cartoon': {'color':'spectrum'}})



