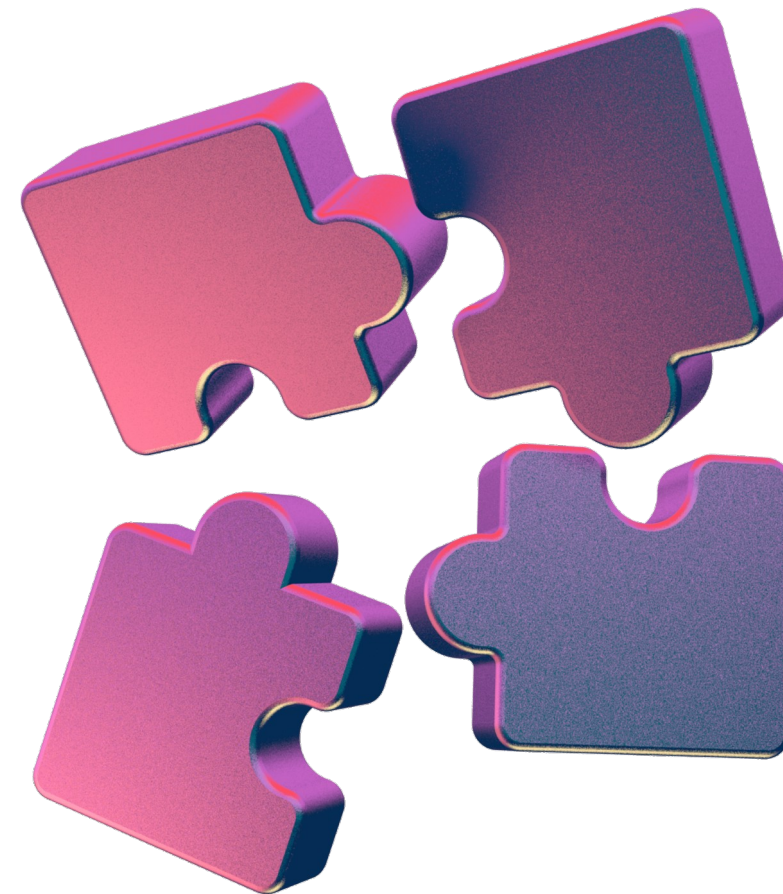


**ML for Drug Discovery
Summer School 2024
@ Mila
A brief summary**



Motivation: What's the current state of ML for DD?

Motivation

Recent headlines (2018-2020)

SPOTLIGHT · 30 MAY 2018

How artificial intelligence is changing drug discovery

World first breakthrough in AI drug discovery

By Emma Morriss · January 30, 2020

RAPID GROWTH IN PUBLISHED RESEARCH USING AI FOR DRUG DISCOVERY

More papers since 2010 than in all prior years combined

AI 2020:

THE FUTURE OF DRUG DISCOVERY

72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 00 01 02 03 04 05 06 07 08 09 10 11 12 13 14 15 16 17

Source: PubMed, July 11, 2018, using this query ("artificial intelligence" or "machine learning" or "deep learning" or "neural network") and (drug or drugs). 1972-2017.



Eli Lilly And Co

NYSE: LLY

Overview

Financials

Compare

Market Summary > Eli Lilly And Co

918.00

 USD

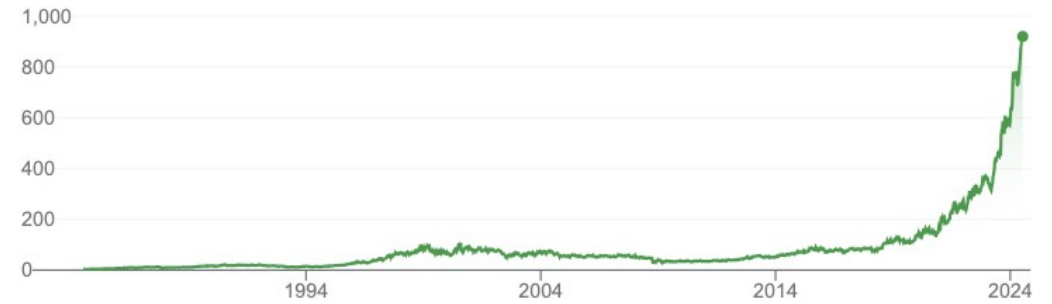
+ Follow

+914.23 (24,250.13%) ↑ all time

Closed: Jul 8, 7:58 p.m. EDT · Disclaimer

After hours 920.75 +2.75 (0.30%)

1D | 5D | 1M | 6M | YTD | 1Y | 5Y | Max



Open	917.89	Mkt cap	872.47B	GDP score	A-
High	928.60	P/E ratio	135.12	52-wk high	928.60
Low	912.00	Div yield	0.57%	52-wk low	434.34

Feedback

More about Eli Lilly And Co →

Open Questions

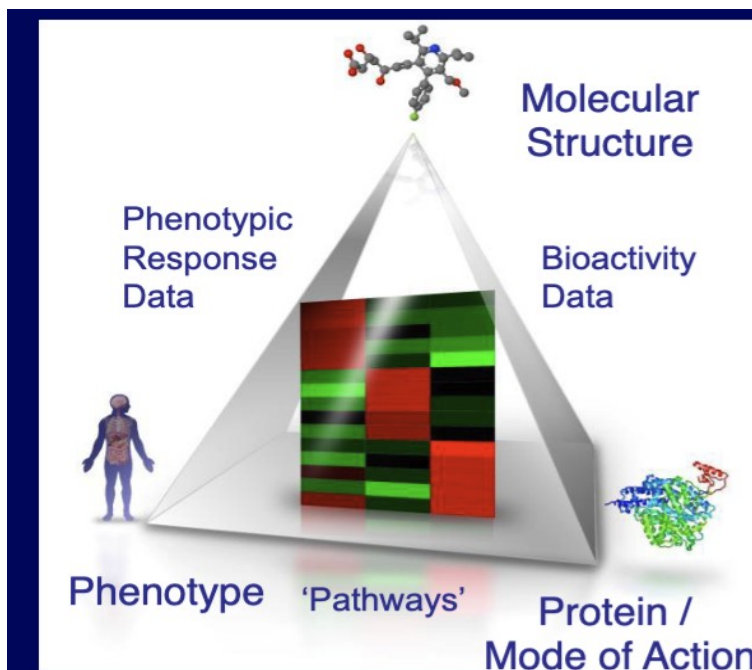
Are we moving from the 'all you need is small algorithms' to the 'data is all you need' paradigm?
(so what data do we need, then?) (Computer is tremendously powerful... but is our data?)

what is the correct data modality to key into my new setting'
(for example, for cellular biology should it be ATACseq, or RNAseq, or proteomics or all of the above)?"

What are the theoretical and scientific limitations of AIDD (AI for Drug Discovery)?
(Causality, confidence, quantification)

Do LLMs have use in AIDD?

Are we reliving 'The Billion Dollar Molecule',
but with an AI twist?



a.k.a.
"The world is flat"

= "We believe our labels"

(which are often
insufficiently quantified, not
directed, unconditional,
don't have time/
concentration/biological
setup dependence, *etc.*)

Open Questions

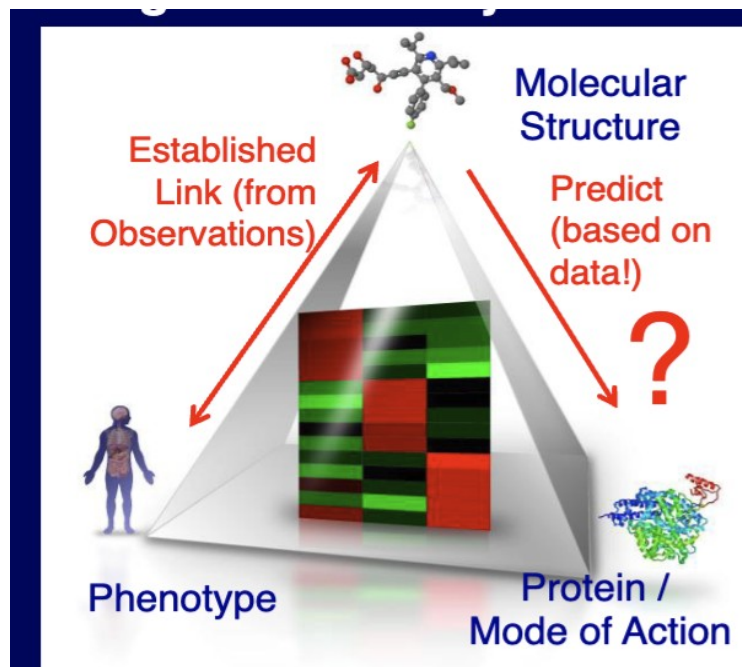
Are we moving from the 'all you need is small algorithms' to the 'data is all you need' paradigm?
(so what data do we need, then?) (Computer is tremendously powerful... but is our data?)

what is the correct data modality to key into my new setting'
(for example, for cellular biology should it be ATACseq, or RNAseq, or proteomics or all of the above)?"

What are the theoretical and scientific limitations of AIDD (AI for Drug Discovery)?
(Causality, confidence, quantification)

Do LLMs have use in AIDD?

Are we reliving 'The Billion Dollar Molecule',
but with an AI twist?

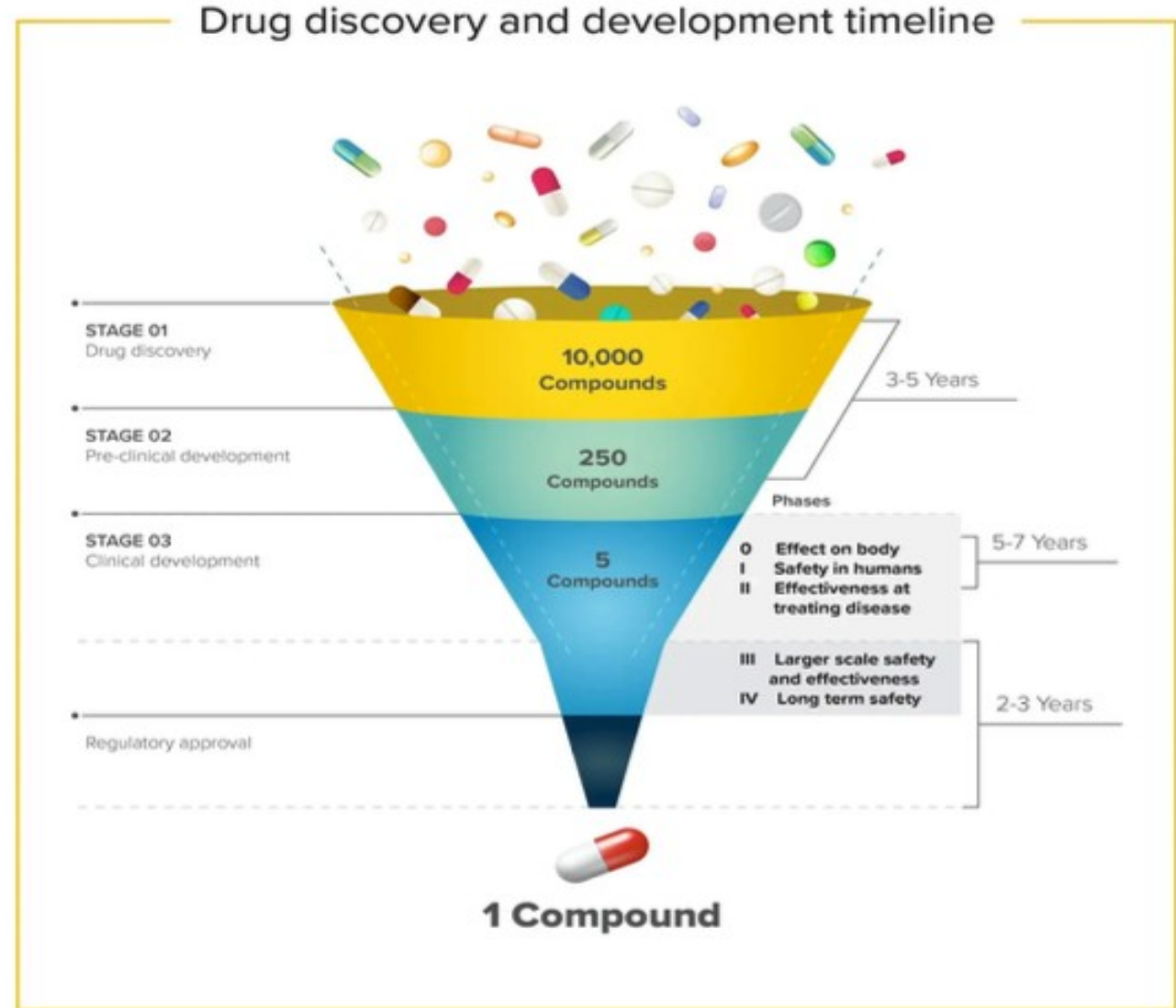
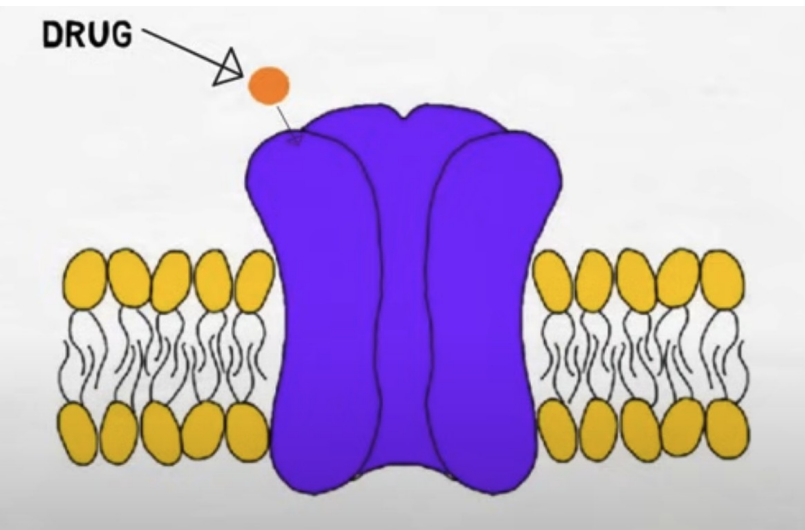
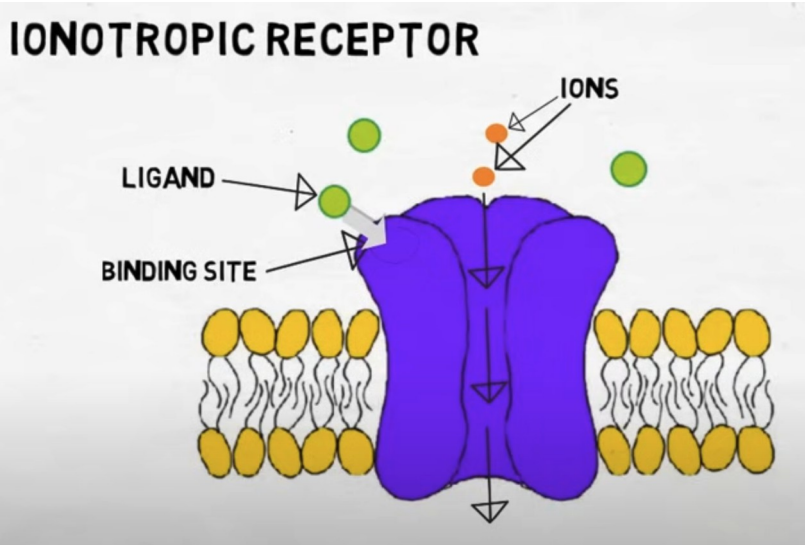


a.k.a.
"The world is flat"

= "We believe our labels"

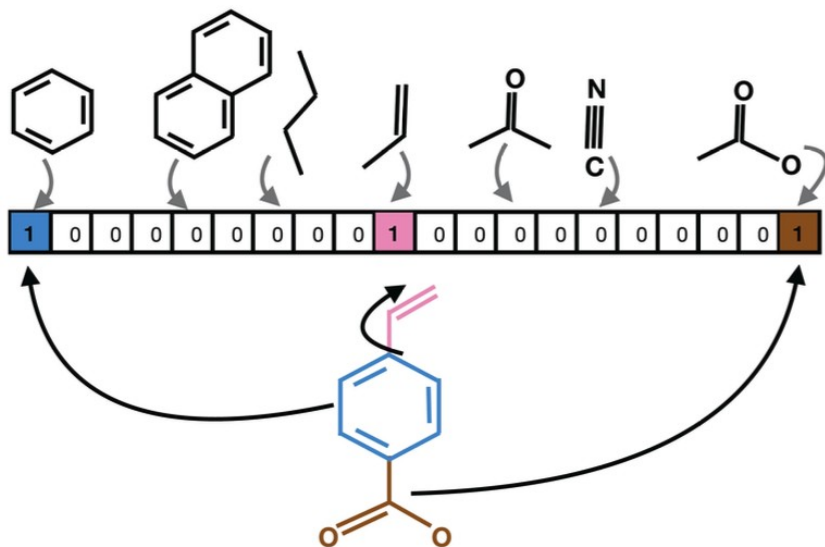
(which are often
insufficiently quantified, not
directed, unconditional,
don't have time/
concentration/biological
setup dependence, *etc.*)

THE LONG FUNNEL TOWARDS A DRUG



Refresher

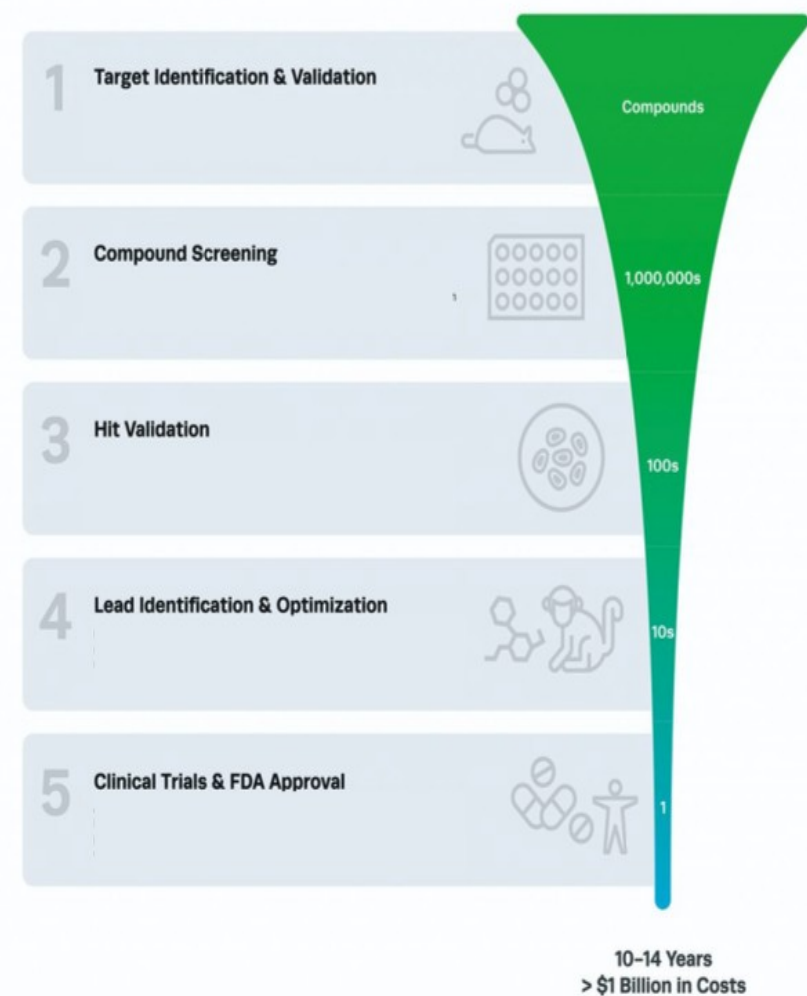
MACHINE LEARNING 101: TURN MOLECULES INTO VECTORS



Turn into Linear Algebra !

Gateway to ML !

LOW DATA IS A FUNDAMENTAL CHALLENGE IN DRUG DISCOVERY



4 secret sauce ingredients for successful applications of ML in biology:

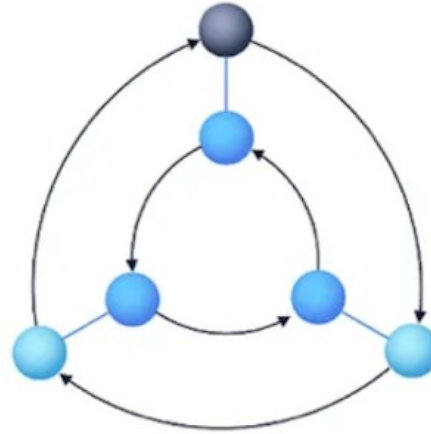
Ingredient 1:
Large and diverse data



Ingredient 2:
Clear performance criterion



Ingredient 3:
Internal structure



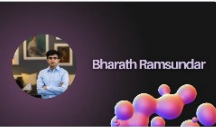
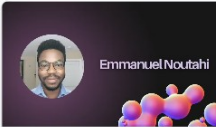
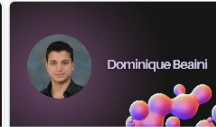
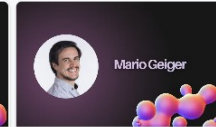
Ingredient 4:
Degenerate solution space



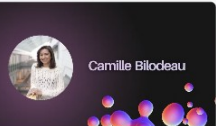
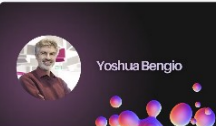
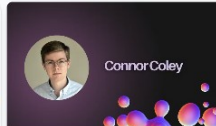

Four “ingredients” necessary to make biological problems amenable for machine learning methods: 1) Large and, most importantly, diverse data; 2) clearly defined performance criterion that can be used as a loss function; 3) useful internal structures generated through human language or biological evolution or coming from physical constraints, e.g., a symmetry group shown here; and 4) a degenerate solution space, e.g., manifested in the form of “manifold hypothesis.”

Schedule & Outline


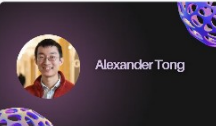

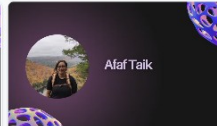
Day 1: Foundations and ML in Ligand-Based Modeling

 <p>1.1 Intro to ML in Drug Discovery: Principles & Applications</p> <p>Professor: Bharath Ramsundar Lecture Recording: Coming soon... Slides:</p> <p>📄 0 📅 06/12/2024</p>	 <p>1.2 Molecular Representation & Scoring</p> <p>Professor: Emmanuel Noutahi Lecture Recording: Coming soon... Lecture Slides:</p> <p>📄 5 📅 06/12/2024</p>	 <p>1.3 Graph Neural Networks for Chemistry</p> <p>Professor: Dominique Beaini Lecture Recording: Coming soon... Lecture Slides:</p> <p>📄 2 📅 06/12/2024</p>	 <p>1.4 Learning Geometry & 3D Symmetries</p> <p>Professor: Mario Geiger Lecture Recording: Coming soon... Lecture Slides: https://slides.com/mariogeiger/learning-geometry-3d-...</p> <p>📄 0 📅 06/12/2024</p>
---	--	--	---

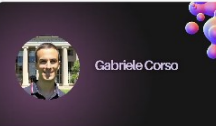
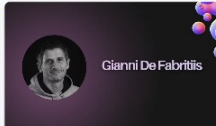
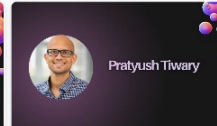
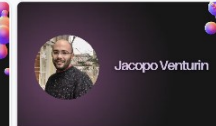
Day 3: Generative Models and Molecular Design

 <p>3.1 Generative Models of Molecular Structures</p> <p>Professor: Camille Bilodeau Lecture Recording: Lecture Slides:</p> <p>📄 5 📅 06/14/2024</p>	 <p>3.2 Exploring Molecular Space & Active Learning</p> <p>Professor: Yoshua Bengio Lecture Recording: Lecture Slides:</p> <p>📄 0 📅 06/14/2024</p>	 <p>3.3 Synthesizability & Molecular Synthesis</p> <p>Professor: Connor Coley Lecture Recording: Lecture Slides:</p> <p>📄 0 📅 06/14/2024</p>	 <p>3.4 Sampling Physical & 3D Systems</p> <p>Professor: Michael Bronstein Lecture Recording: Lecture Slides:</p> <p>📄 1 📅 06/14/2024</p>
---	--	---	---

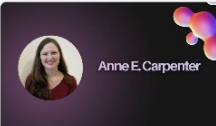
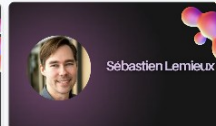
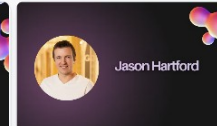
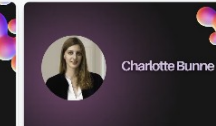
Day 5: Frontiers in AI for Drug Discovery

 <p>5.1 LLMs in Drug Discovery</p> <p>Professor: Andres M Bran Lecture Recording: Coming soon... Slides: Notebook: Colab</p> <p>📄 1 📅 06/18/2024</p>	 <p>5.2 Protein Folding & Design</p> <p>Professor: Alex Tong Lecture Recording: Coming soon... Slides:</p> <p>📄 4 📅 06/18/2024</p>	 <p>5.3 Open-Source Initiatives & Benchmarking Efforts</p> <p>Professor: Karmen Condic-Jurkic Lecture Recording: Coming soon... Slides:</p> <p>📄 2 📅 06/18/2024</p>	 <p>5.4 Ethical & Bias Concerns</p> <p>Professor: Afaf Taik Lecture Recording: Coming soon... Slides:</p> <p>📄 1 📅 06/18/2024</p>
--	--	--	---

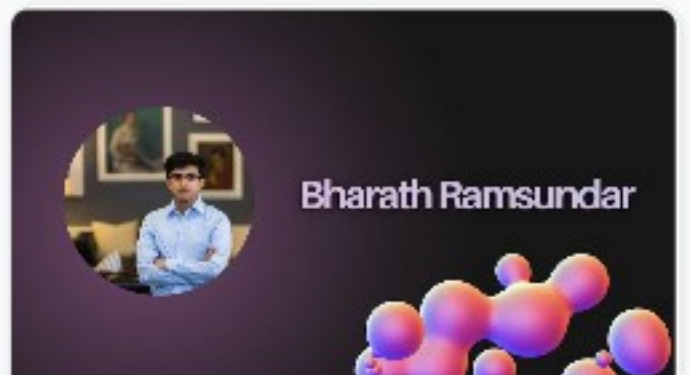
Day 2: ML in Structure-Based Drug Discovery

 <p>2.1 ML in Structure-Based Drug Discovery</p> <p>Professor: Gabriele Corso Lecture Recording: Lecture Slides:</p> <p>📄 10 📅 06/13/2024</p>	 <p>2.2 Learning ML Interatomic Potentials</p> <p>Professor: Gianni De Fabritiis Lecture Recording: Coming soon... Lecture Slides:</p> <p>📄 0 📅 06/13/2024</p>	 <p>2.3 Accelerate Atomistic Simulations, Sampling, and Dynamics</p> <p>Professor: Pratyush Tiwary Lecture Recording: Coming soon... Lecture Slides: PPT https://docs.google.com/presentation/d/1sODOQ8r9eGmQZbx-...</p> <p>📄 0 📅 06/13/2024</p>	 <p>2.4 Coarse-Grained Biological Systems</p> <p>Professor: Jacopo Venturin Lecture Recording: Coming soon... Lecture Slides:</p> <p>📄 3 📅 06/13/2024</p>
---	--	--	---

Day 4: Target Discovery and Deconvolution

 <p>4.1 Phenomics in Drug Discovery</p> <p>Professor: Anne Carpenter Lecture Recording: Lecture Slides:</p> <p>📄 0 📅 06/17/2024</p>	 <p>4.2 Multi-Modal Omics & AI</p> <p>Professor: Sébastien Lemieux Lecture Recording: Lecture Slides:</p> <p>📄 0 📅 06/17/2024</p>	 <p>4.3 Causal Discovery & Representation Learning</p> <p>Professor: Jason Hartford Lecture Recording: Lecture Slides:</p> <p>📄 0 📅 06/17/2024</p>	 <p>4.4 Modeling Population Dynamics</p> <p>Professor: Charlotte Bunne Lecture Recording: Lecture Slides:</p> <p>📄 6 📅 06/17/2024</p>
---	---	--	---

1.1 Bharath



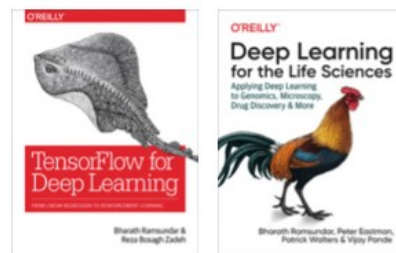
1.1 Intro to ML in Drug Discovery: Principles & Applications

Professor: Bharath Ramsundar
Lecture Recording: Coming soon... Slides:

0 06/12/2024

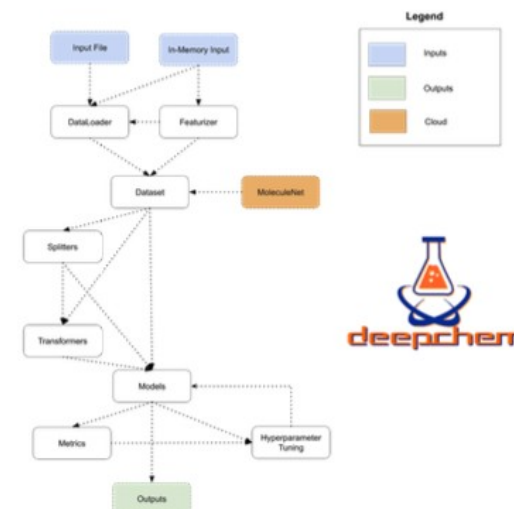
BRIEF INTRO: BHARATH RAMSUNDAR

- ▶ Stanford PhD, Pande Group, UCB EECS/Math Alum
- ▶ Lead Developer of DeepChem
- ▶ Founder/CEO of Deep Forest Sciences

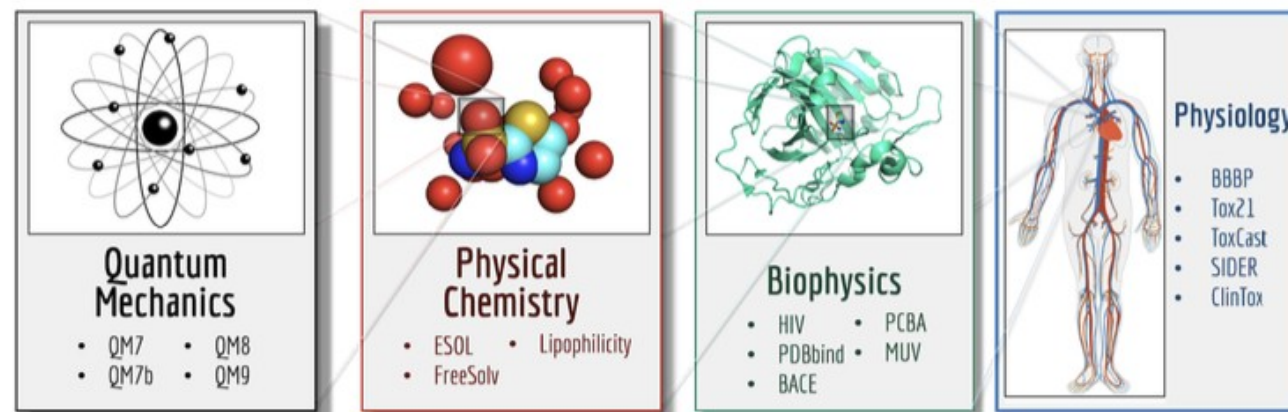


SCIENTIFIC MACHINE LEARNING WITH DEEPCHEM

- ▶ DeepChem is a framework to apply AI to drug discovery and other scientific problems.
- ▶ DeepChem offers composable re-usable pipelines for scientific workflows
- ▶ DeepChem increasingly supports non-molecular applications for general scientific computing.



MOLECULENET PROVIDES CONVENIENT BENCHMARKS



1.1 Summary

1. Despite the hype, AI/ML is not “that useful” in Drug Discovery, e.g. because it cannot understand long time scales that drug discovery requires, e.g. similar to civil engineering and building of bridge and castles

- Don't trust machine learning too much.
- Physics > Datasets (most of the time)
- Understand the biology and chemistry of the system.
- AI is not a substitute for basic science. AI/ML does not generalize well; science does.

2. Low Data is the fundamental challenge in Drug Discovery, e.g. as we take out biology predictions from the test tube, they break down in the real world, because out of distribution generalization is a really hard problem (and unsolvable with the wrong data)

3. Strategies to tackle low data are

- Physical Simulations: to generate synthetic data
- Metric + Meta Learning: to compare prediction vectors
- Multitask-Learning: to combine datasets
- Transfer Learning: from assay readouts

1.2 Emmanuel

This talk had overlap with Bharat 's talk

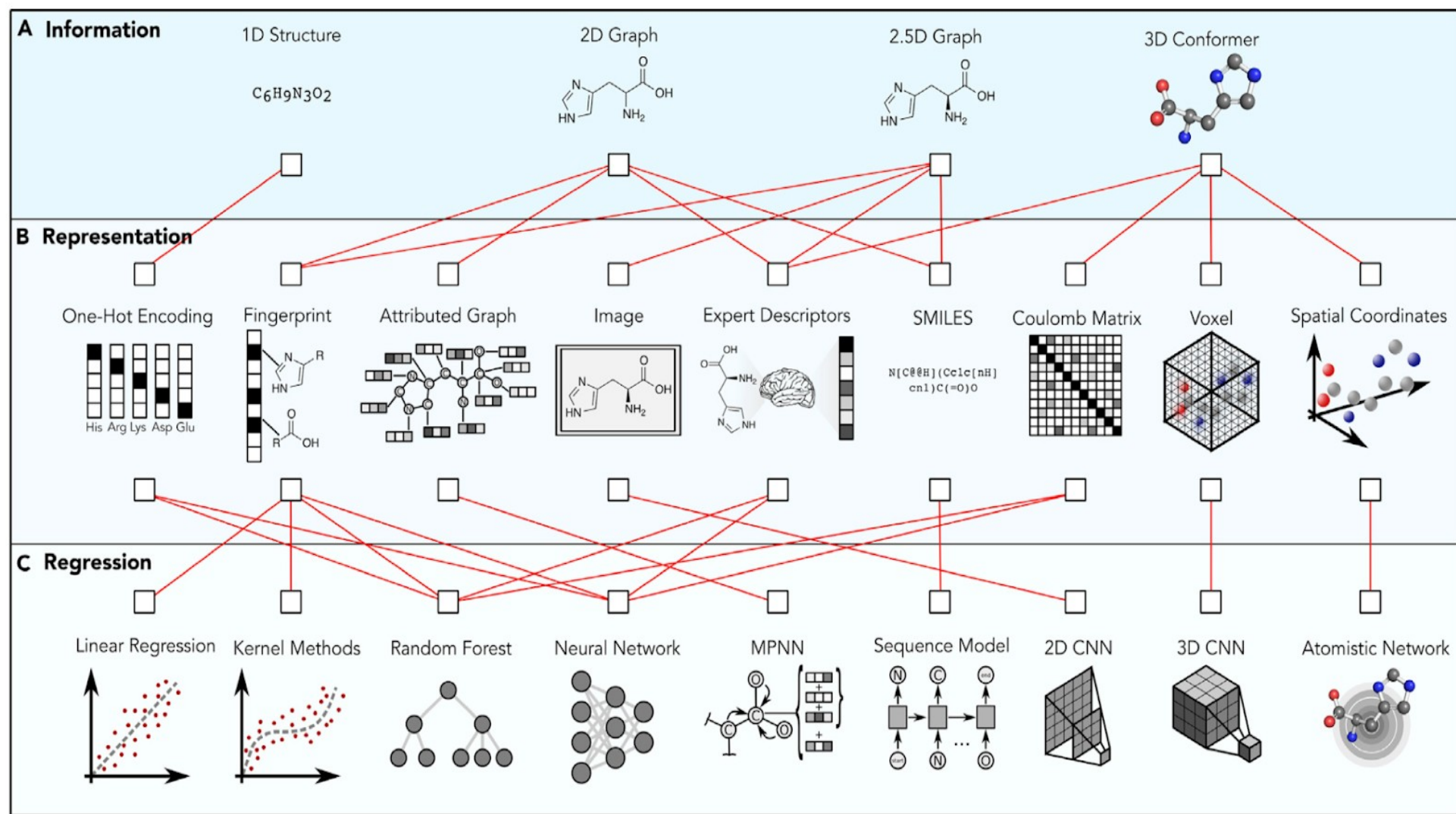


Emmanuel Noutahi

1.2 Molecular Representation & Scoring

Professor: Emmanuel Noutahi
Lecture Recording: Coming soon... Lecture Slides:

5 06/12/2024




1.2 Summary

This talk had overlap with Bharat 's talk


Good representations are critical: “A method cannot save an unsuitable representation, which cannot remedy irrelevant data, for an ill thought-through question.”

Random Forest + ECFP is a strong base line,
(as shown in “Random forest classification for predicting lifespan-extending chemical compounds” by Kapsiani et al.)

1.3 Dominique



Dominique Beaini



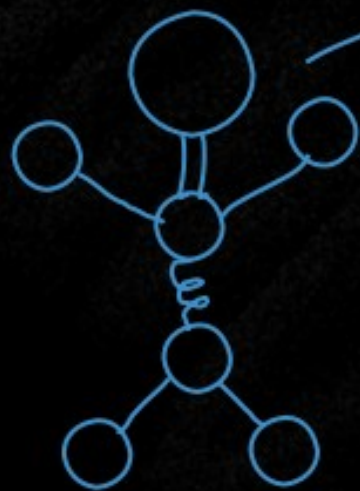
1.3 Graph Neural Networks for Chemistry

Professor: Dominique Beaini
Lecture Recording: Coming soon... Lecture Slides:

2 06/12/2024



Meeting Graphy 🖐️



Hello everyone 🖐️! I'm Dom's assistant for today!

Let's visit the **molecular graph world** together!

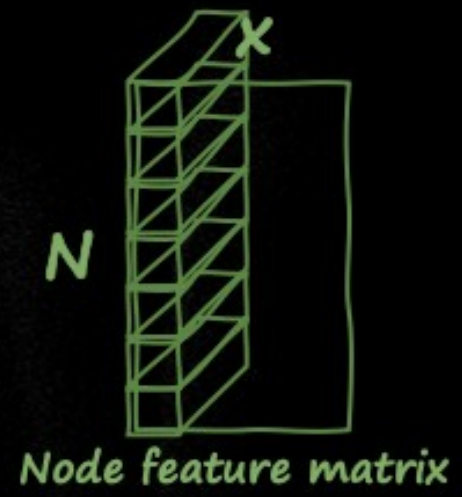
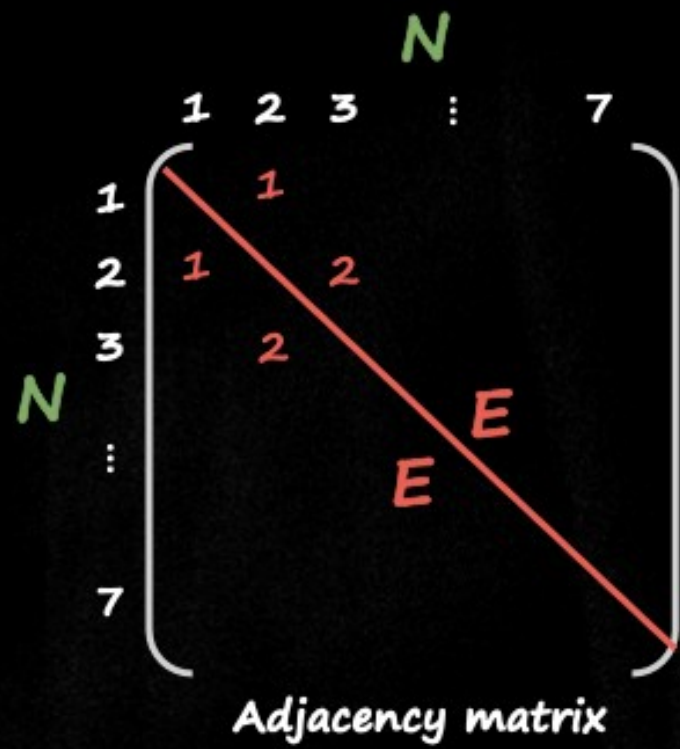
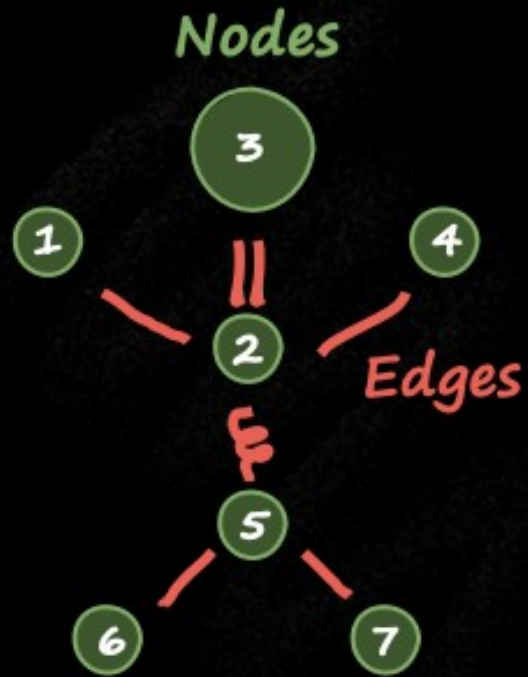
We'll first learn **what are graphs** and how to manipulate them

Then we'll look into **standard GNNs** graph neural networks

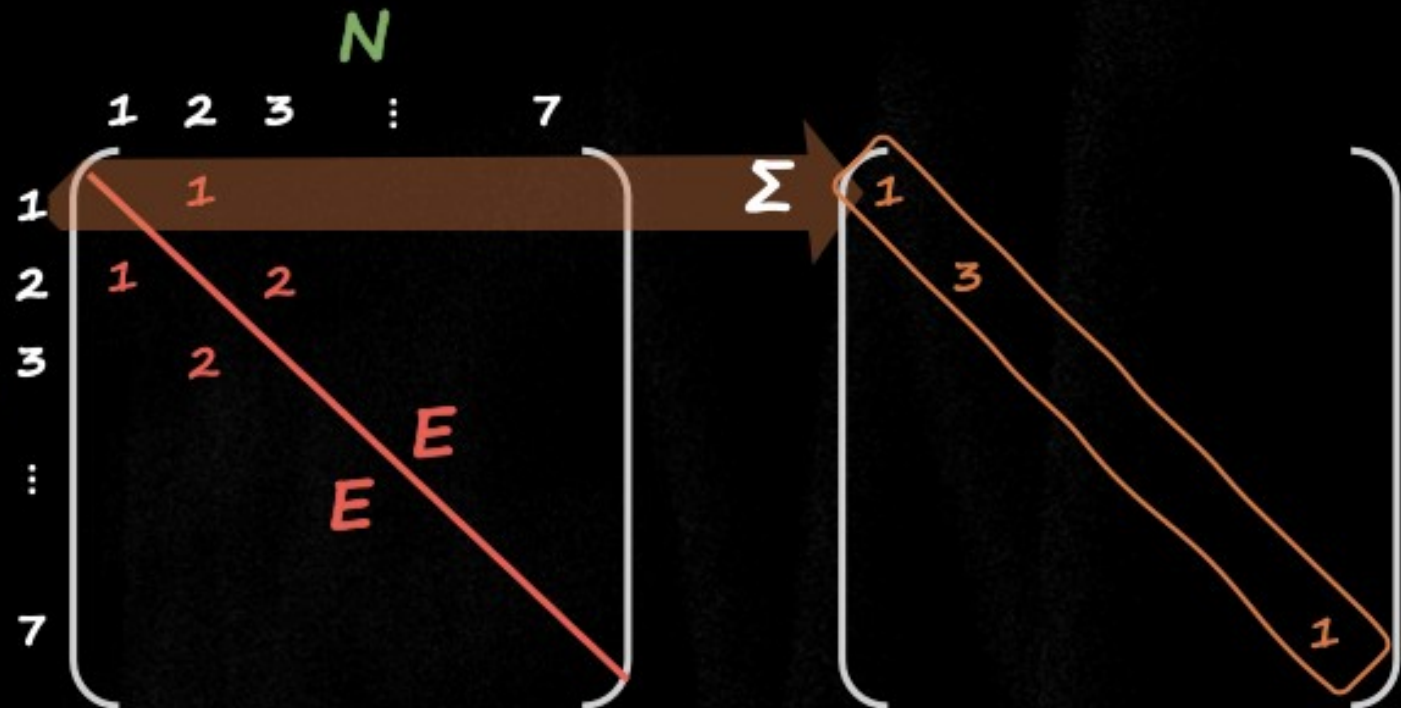
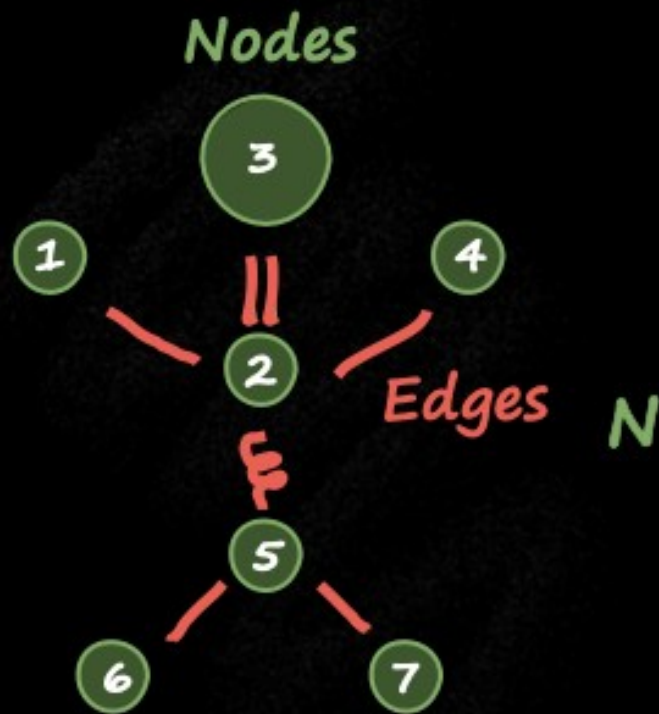
And how to build more **expressive GNNs** for molecules

Finally, we will scale a **Graph Transformer** together

Anatomy of Graphy



Laplacian matrix



Adjacency matrix

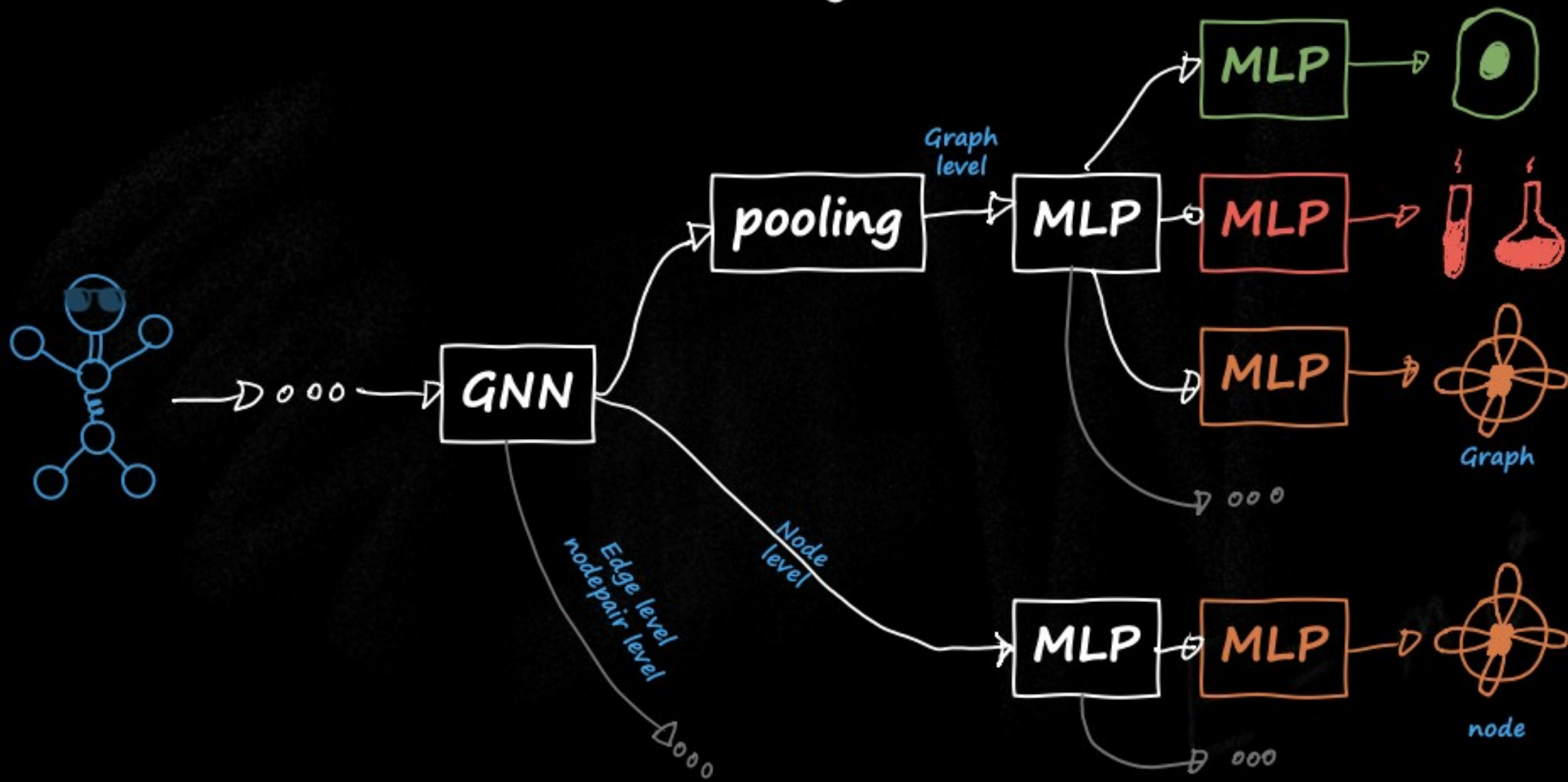
Degree Matrix

A

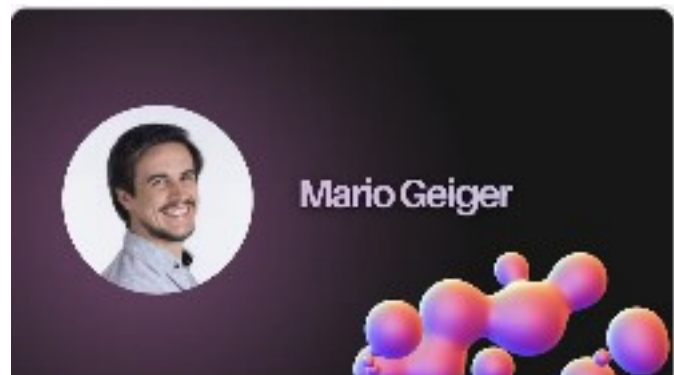
D

$$L = D - A$$

Multi-level multi-tasking



1.4 Mario



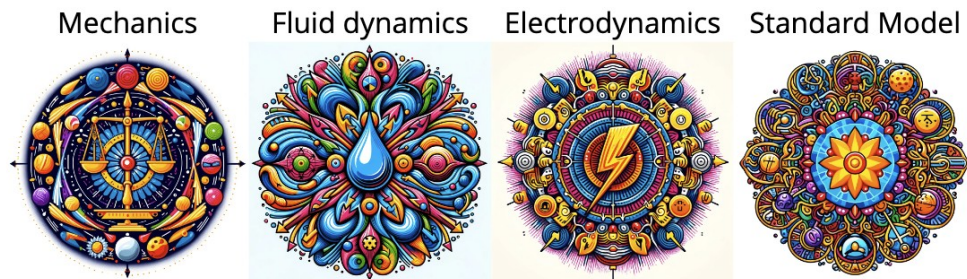
1.4 Learning Geometry & 3D Symmetries

Professor: Mario Geiger Lecture Recording: Coming soon...

Lecture Slides:

<https://slides.com/mariogeiger/learning-geometry-3d-...>

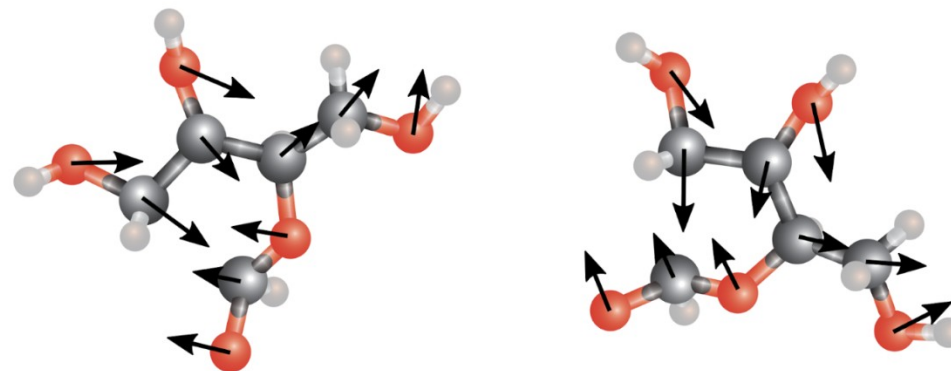
0 06/12/2024



	Mechanics	Fluid dynamics	Electrodynamics	Standard Model
Rotation	✓	✓	✓	✓
Translation	✓	✓	✓	✓
Time translation	✓	✓	✓	✓
Boosts (Galilean or Lorentz)	✓	✓	✓	✓

Example of Equivariance

$$f : \text{positions} \rightarrow \text{forces}$$



Where is equivariance used in AI?

Protein Folding

EquiFold Jae Hyeon Lee et al.

Protein Docking

DIFFDOCK Gabriele Corso et al.

Molecular Dynamics

Nequip S. Batzner et al.

MACE I. Batatia et al.

Allegro A. Musaelian et al.

Open Catalyst Project.

Solid State Physics

Prediction of Phonon Density Z. Chen et al.

Molecular Electron Densities

Cracking the Quantum Scaling Limit with Machine Learned Electron Densities J. Rackers

Cosmology, Medical Images and others



44M atoms while taking advantage of up
to 5120 GPUs
Albert Musaelian

Group and Representations

"what are the operations" "how they compose"

"vector spaces on which the action of the group is defined"

Representation $D(g, x)$

- $g \in G, x \in V$
- Linear $D(g, x + y) = D(g, x) + D(g, y)$
- Follow the structure of the group
 $D(gh, x) = D(g, D(h, x))$

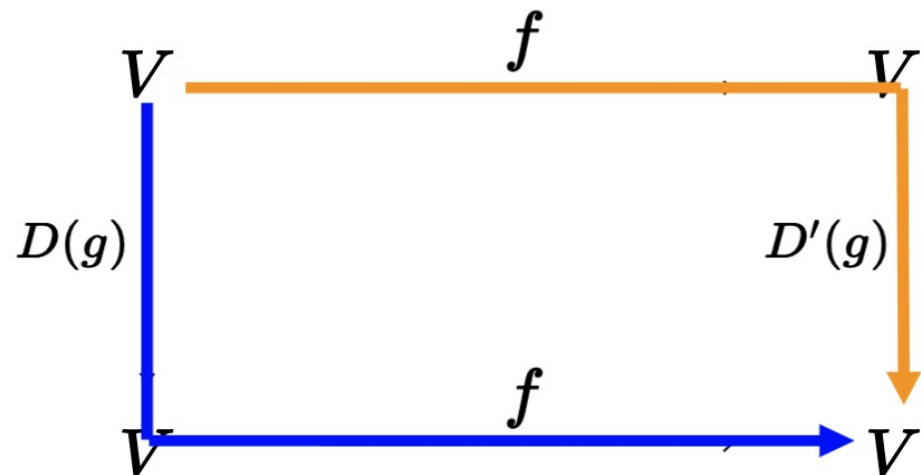
Equivalent notation $D(g)x$

- $D(g) : V \rightarrow V$
- $D(g) \in \mathbb{R}^{d \times d}$
- $D(gh) = D(g)D(h)$

Group G

- identity $\in G$
- associativity $g(hk) = (gh)k$
- inverse $g^{-1} \in G$

Equivariance



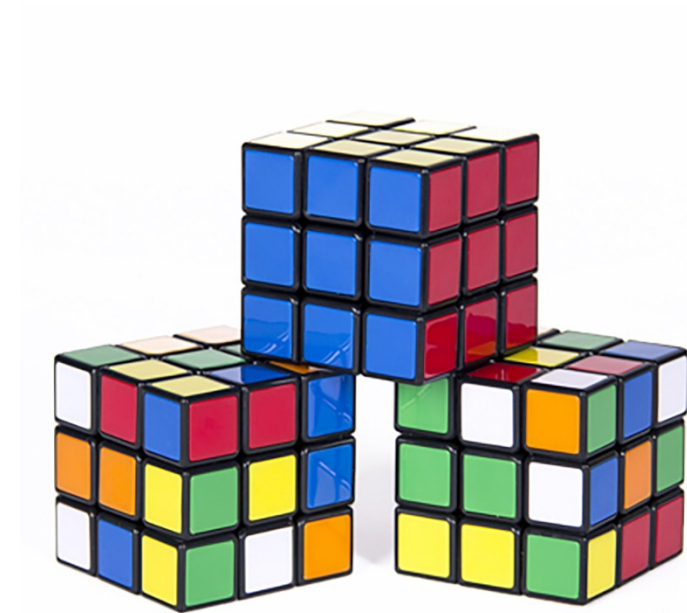
$$f(D(g)x) = D'(g)f(x)$$

1.4 Summary

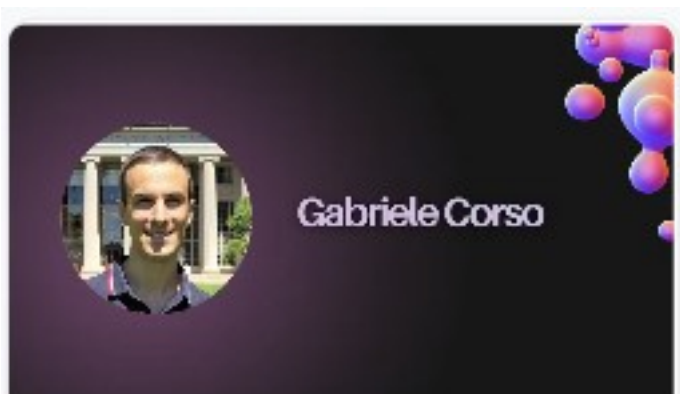
Equivariance is largely motivated by equivariance observed in Physics, e.g. Mechanics, Fluid Dynamics, Electrodynamics and the Standard Model

There is some evidence that equivariance leads to more data efficiency.
(There is also evidence against that)

Equivariance tends to introduce additional computational complexity, but that can be tackled with **reducible representations**.



2.1 Gabriele

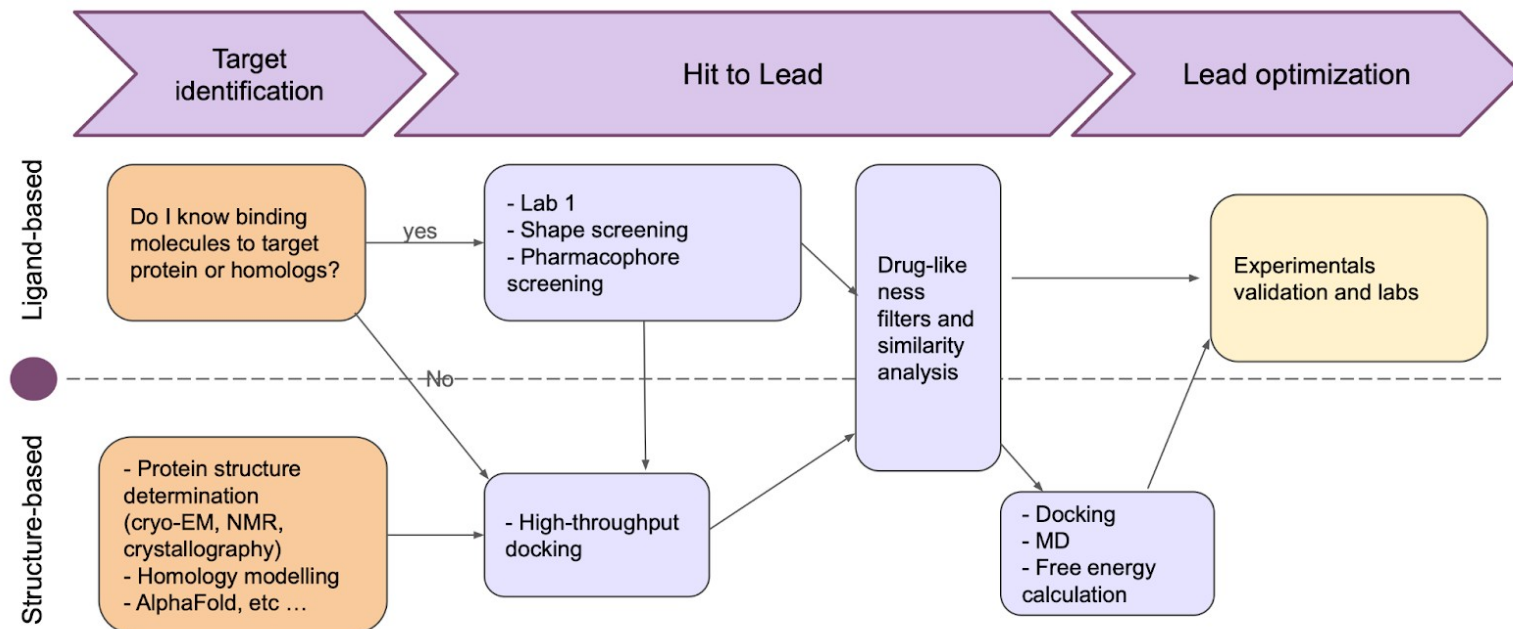


2.1 ML in Structure-Based Drug Discovery

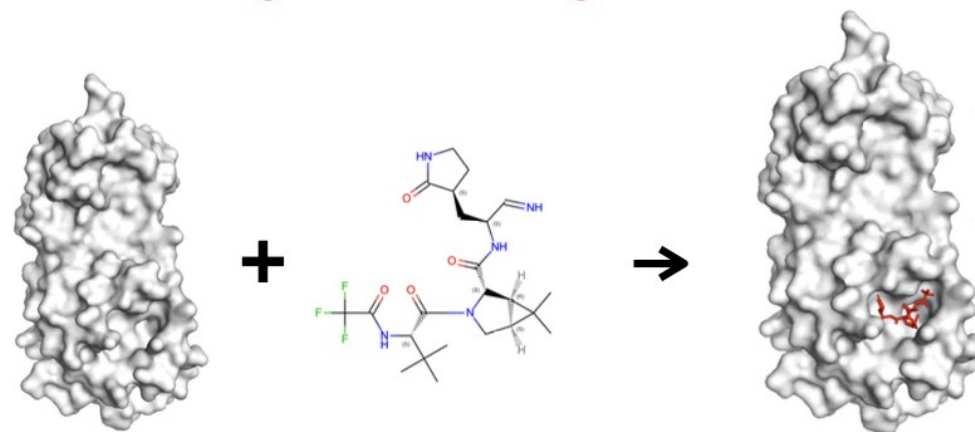
Professor: Gabriele Corso Lecture
Recording: Lecture Slides:

10 06/13/2024

Ligand-based and Structure-based are interconnected



Protein-Ligand Docking



Input: protein + molecule

Output: bound structure

2.1 Summary

There are three approaches to docking: Search-based methods, Regression models and Generative models (see diagram)

Generative models are preferred (at least in his talk) due to **aleatoric uncertainty** introducing an averaging effect as well as issue with model uncertainty in comparison to regression.

There are three components to Diffusion Generative Models: forward diffusion, scoring and reverse diffusion. (see slide)

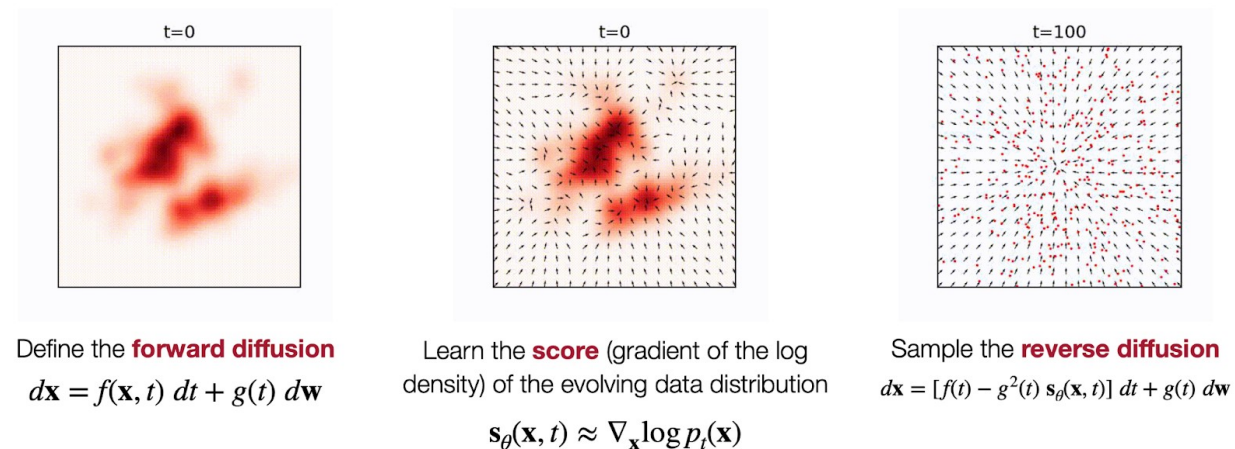
Docking Screening.

Different Approaches to Docking

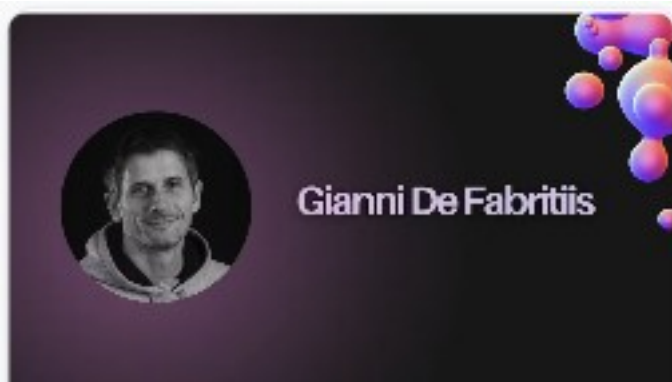


Corso et al., ICLR 2023

Diffusion Generative Models



2.2 Gianni

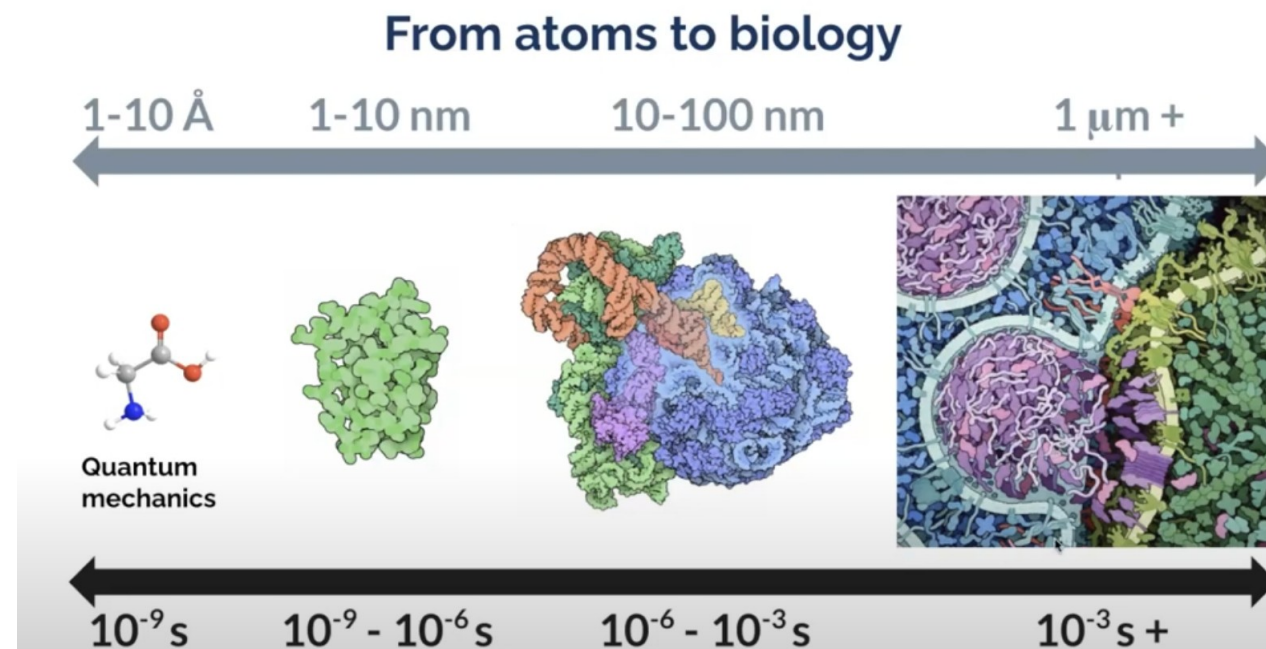


2.2 Learning ML Interatomic Potentials

Professor: Gianni De Fabritiis
Lecture Recording: Coming soon...
Lecture Slides:

0 06/13/2024

Acellera: Accelerating Drug Discovery



PlayMolecule™

Molecular discovery

MembraneBuilder

DeepSite

Parameterize

SystemBuilder

AceProfiler

AcePrep

2.2 Summary

His talk is a mixed bag (at least for me as an outsider) of his recent research results on:

Machine Learning Potential Architectures: All about learning energy functions which otherwise are not computable, see slide

- There are different type of atomic interactions: bonded and non-bonded.
- Energy functions describe the different types of interactions.
- The parameters of the functions constitute a molecular force-field .

$$U = \sum_{\text{All Bonds}} \frac{1}{2} K_b (b - b_0)^2 + \sum_{\text{All Angles}} \frac{1}{2} K_\theta (\theta - \theta_0)^2$$
$$+ \sum_{\text{All Torsion Angles}} K_\phi [1 - \cos(n\phi + \delta)]$$
$$+ \sum_{\text{All nonbonded pairs}} \epsilon \left[\left(\frac{r_0}{r}\right)^{12} - 2\left(\frac{r_0}{r}\right)^6 \right]$$
$$+ \sum_{\text{All partial charges}} \frac{332 q_i q_j}{r}$$


Cartesian Tensor Representations for Efficient Learning of Molecular Potentials:

All about “Incorporating **transformation properties** of physical quantities into NNPs (Neural network potentials for chemistry)

Coarse-grained Potentials:

All about the “need to guarantee physical equivariiances/invariances: translation, rotational, order” in coarse grained representations

2.3 Pratyush



Pratyush Tiwary

2.3 Accelerate Atomistic Simulations, Sampling, and Dynamics

Professor: Pratyush Tiwary
Lecture Recording: Coming soon... Lecture Slides: PPT
<https://docs.google.com/presentation/d/1sODOQ8r9eGmQZbx-...>

0 06/13/2024

With infinite compute, all we need is a single equation,
but compute is not infinite, so :0

The 2nd most important formula in thermodynamics & statistical mechanics.

In an ideal world with infinite compute power,
this is the only equation one needs

$$Z = \sum_i e^{-\beta E_i}$$

But the world isn't ideal and compute power is finite

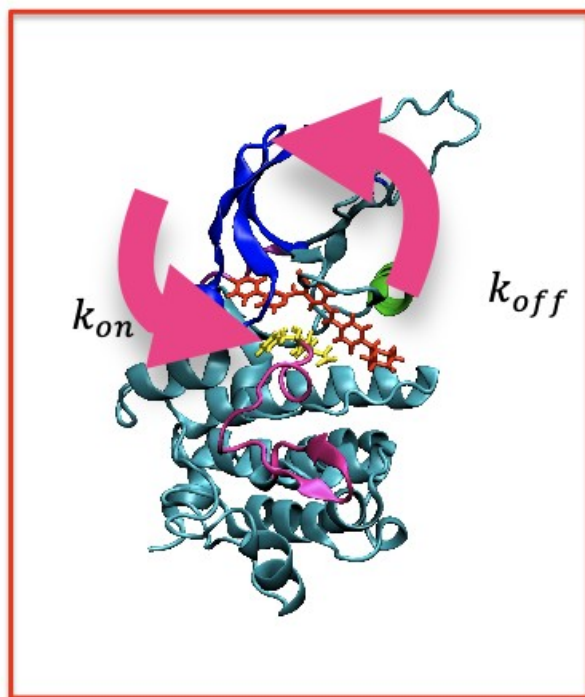
Much of **statistical mechanics** (120+ years)
and **machine learning** (?? years) have been working around finding
tricks to sample the partition function Z ,
or more generally, avoid sampling Z .

The 2 disciplines have much to teach other

2.3 Summary

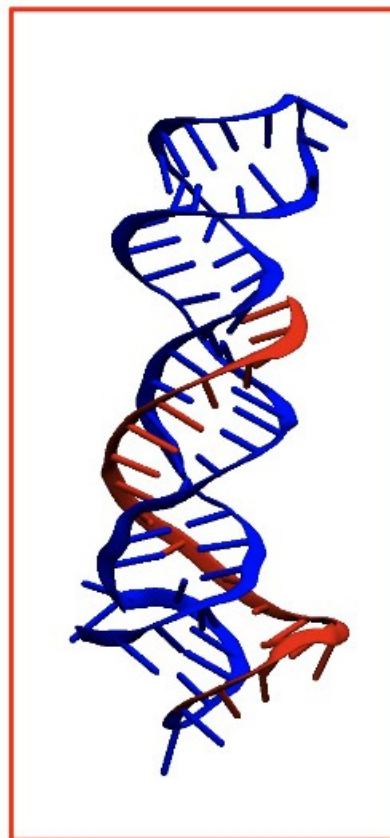
My lab combines statistical mechanics and generative AI to tackle problems of human health and energy relevance

guided by structure & dynamics

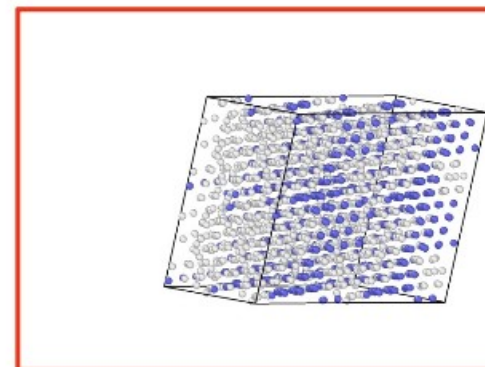


**Protein-
small molecules
drug discovery**

github.com/tiwaryl原因



**RNA
therapeutics**



**Finite-temperature
crystal polymorphs &
phase transitions**

2.3 Summary

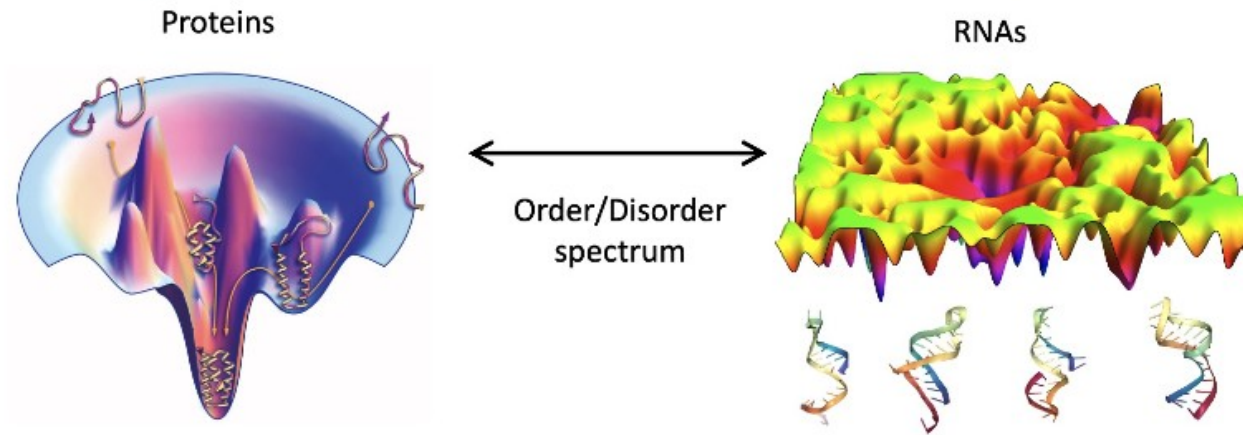
**AI can now be used to easily, routinely
predict structure and dynamics**

Molecular dynamics is dead

Not quite ...

2.3 Summary

Life is not about a single structure but an ensemble of structures
with **just the right fluctuations across different length & timescales**



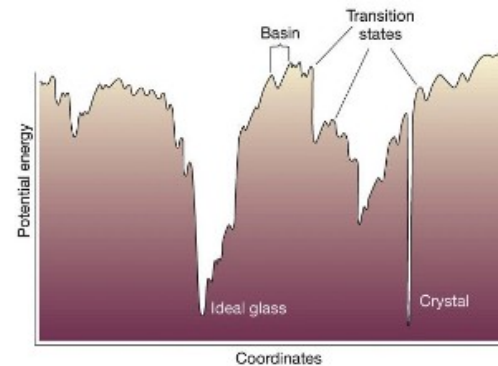
K. Dill, J. McCallum, *Science* (2012).
The Protein-Folding problem, 50 years on

**Fluctuations between metastable states
described by few slow modes & timescale separation**

D.J. Wales, *Ann. Rev. Phys. Chem.* (2017).
Exploring Energy Landscapes

**No obvious timescale separation &
no dominant driving fluctuations**

Somewhere on the order/disorder spectrum

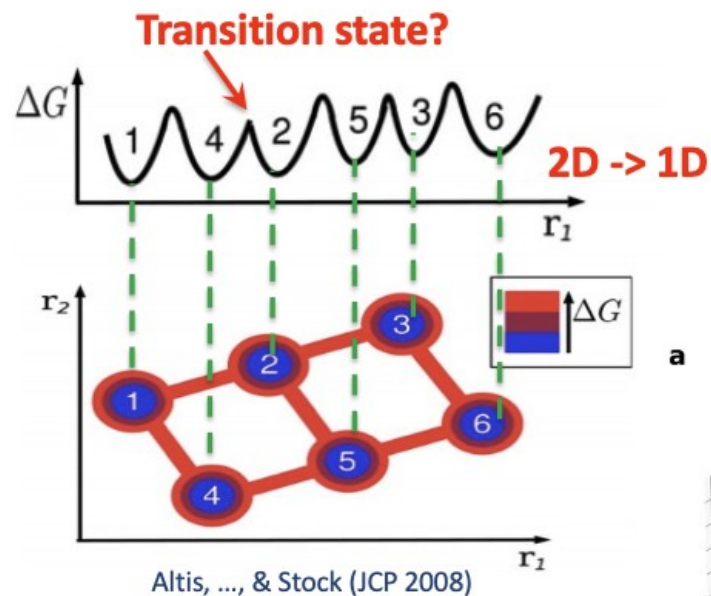


**Beyerle, Zou, Tsai, Tiwary
PNAS 2023**

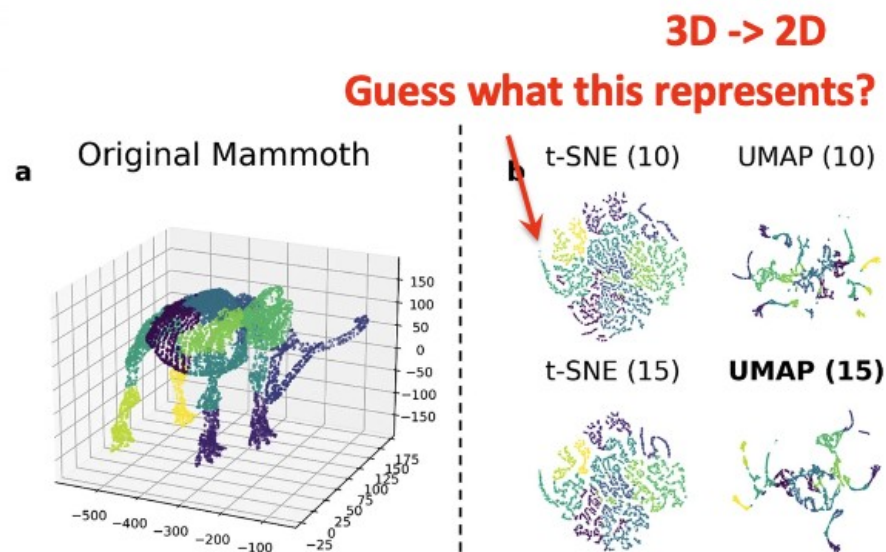
Debenedetti, Stillinger, Nature (2001).
Supercooled Liquids & the Glass Transition

2.3 Summary

All reduced representations are correct
Some are useful
Few are meaningful



Tribello & Gasparotto (Frontiers in Mol. Bio. 2019)

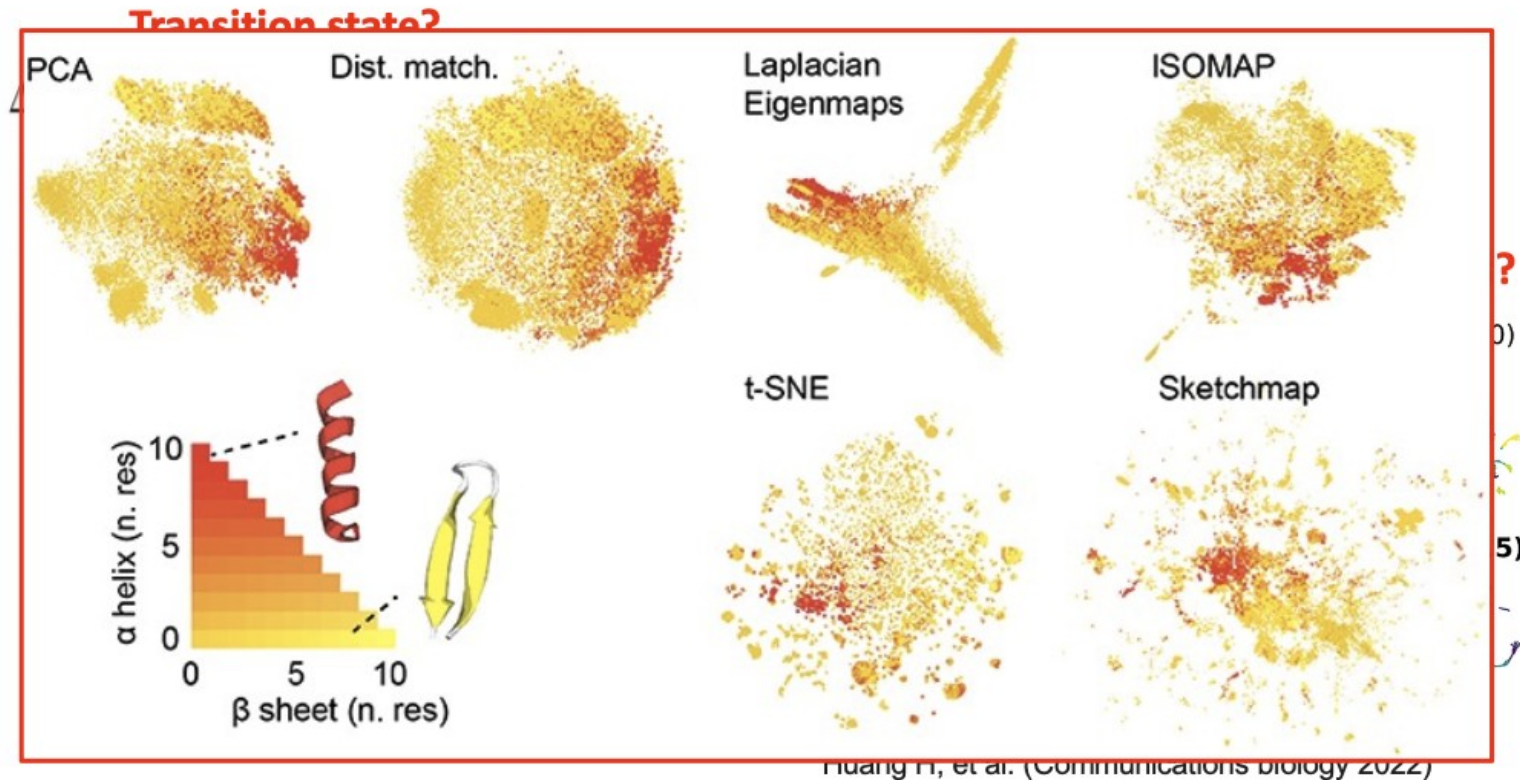


Huang H, et al. (Communications biology 2022)

Why should you trust these projections?
N-D \rightarrow 2D

2.3 Summary

All reduced representations are correct
Some are useful
Few are meaningful



Tribello & Gasparotto (Frontiers in Mol. Bio. 2019)

Why should you trust these projections?
N-D -> 2D

2.3 Summary

Take-home:

Molecular Dynamics – 60 years old and still very relevant

- **AI tools are efficient hypothesis generators** that when combined with physics based simulators – especially MD - in a rigorous feedback loop might have predictive power
- A 100ns MD trajectory is not always 10x more data than 10ns. Important to understand **metastable state to state emergent dynamics**
- **Careful enhanced MD generates correct fluctuations across length and timescales**, critical for training AI models that can mimic reliable molecular conformational diversity and dynamics, & also for next gen forcefields
- Showed 2 methods today – **RAVE and path-sampling LSTM**
- Demo codes open-source @ github.com/tiwarylalab

2.4 Jacopo



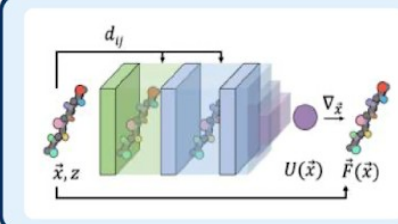
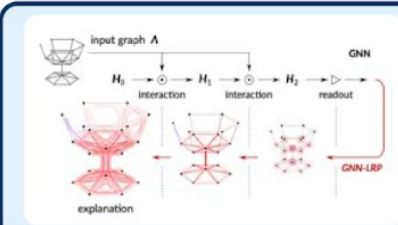
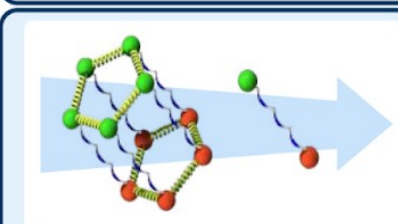
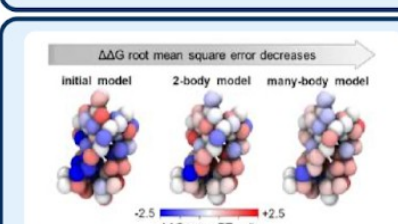
Jacopo Venturin

2.4 Coarse-Grained Biological Systems

Professor: Jacopo Venturin
Lecture Recording: Coming soon... Lecture Slides:

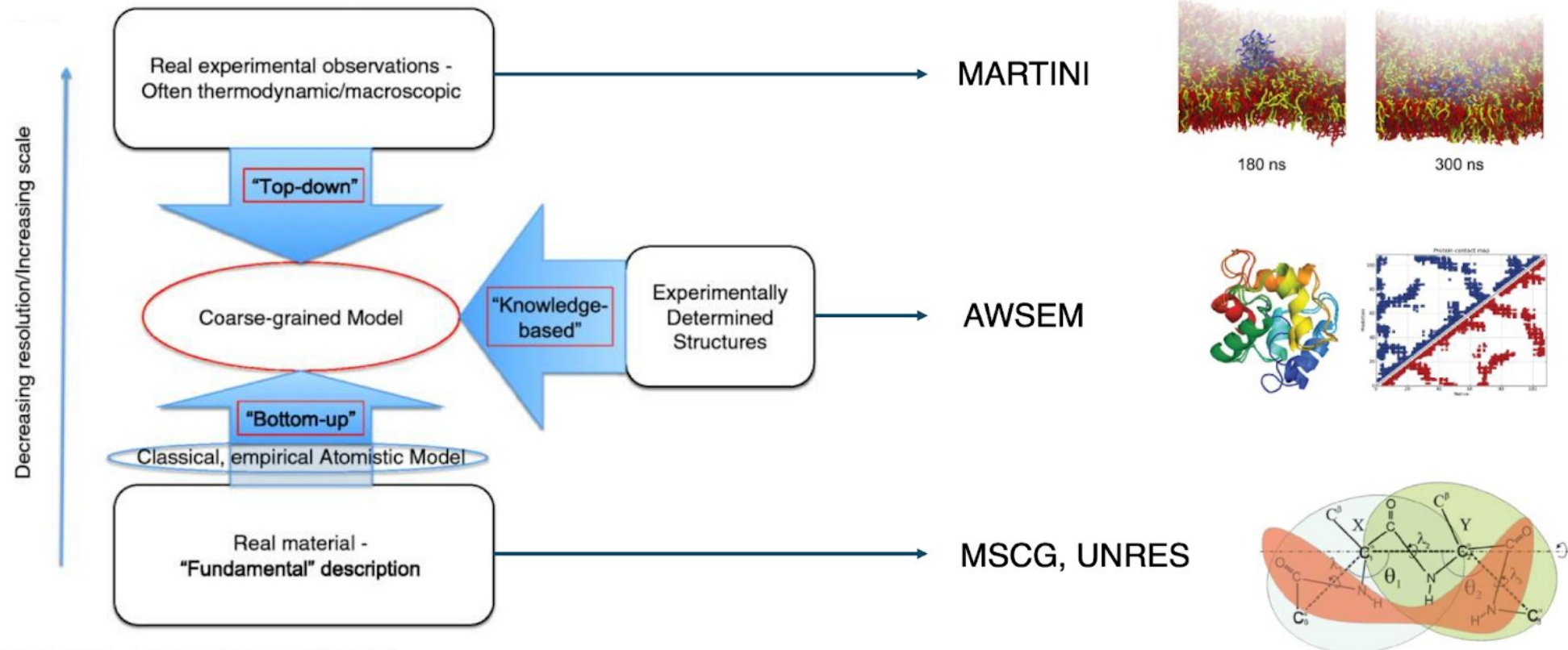
3 06/13/2024

Outline

	Transferable ML-CG force fields for proteins
	Interpreting ML-CG force-fields
	ML-CG force-fields for Path Integral Molecular Dynamics
	Incorporating experimental data in CG force-fields

2.4 Summary

Coarse-graining philosophies and approaches



Noid, W. G. (2013)., *J. Chem. Phys.* 139(9), 90901 (2005)
 Marrink, S. J. *et al*, *J. Chem. Phys. B*, 111(27), 7812-7824 (2007)
 Davtyan, A. *et al*, *J. Chem. Phys. B*, 116(29), 8494-8503 (2012)
 Liwo, A. *et al*, *J. Mol. Model.* 20, 2306 (2014)
 Thorpe, I. F., Zhou, J., & Voth, G. A., *J. Phys. Chem. B*, 112(41), 13079–13090 (2008).

3.1 Camille

Today seems to be a special day with many high-profile speakers.



Camille Bilodeau

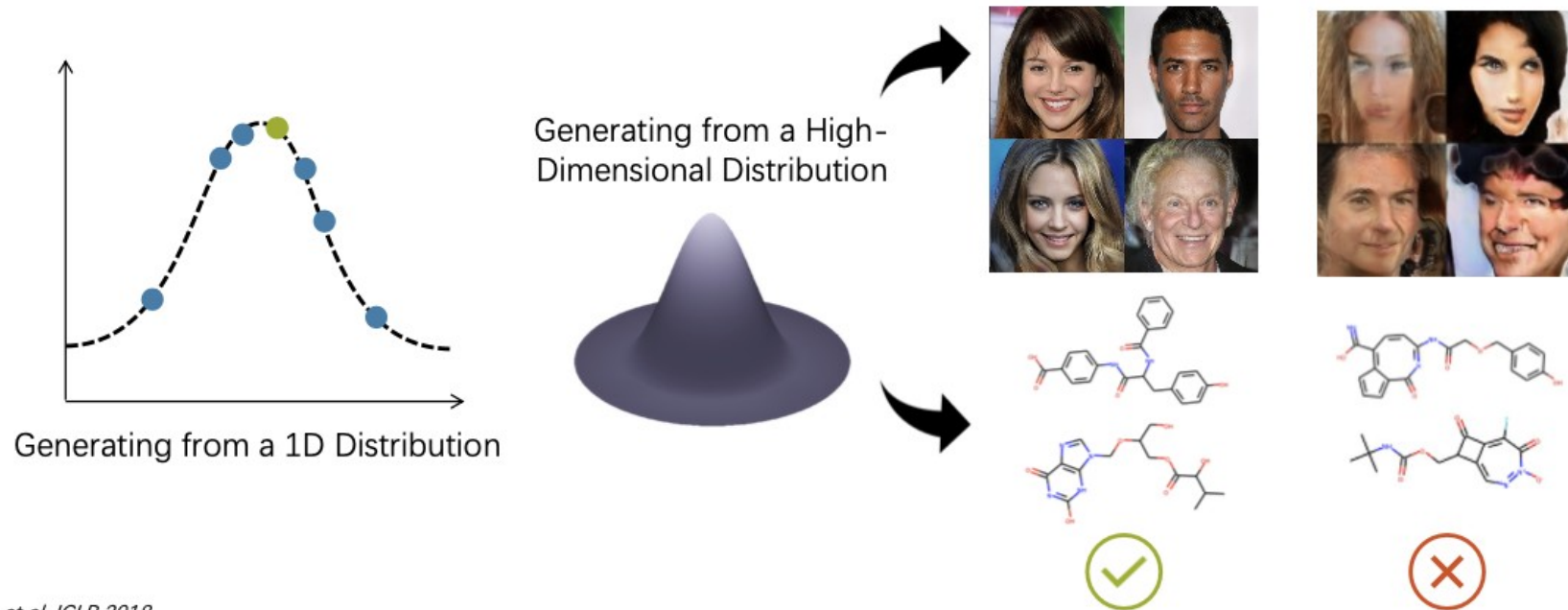
3.1 Generative Models of Molecular Structures

Professor: Camille Bilodeau
Lecture Recording: Lecture Slides:

5 06/14/2024

What is a generative model?

A generative model **learns a distribution** from training examples at training time and **draws new examples** from this learned distribution at test time



arras et al. ICLR 2018

ómez-Bombarelli et al. ACS Cent. Sci. 2018

3.1 Summary

The mechanics of how to generate molecules

How do we train a neural network to learn a distribution?

Objective: We want a model that generates data such that the distribution of the generated data **matches the distribution of the input data as closely as possible.**

Variational Autoencoder

Learn the parameters of a normalized function that approximates $P(x)$



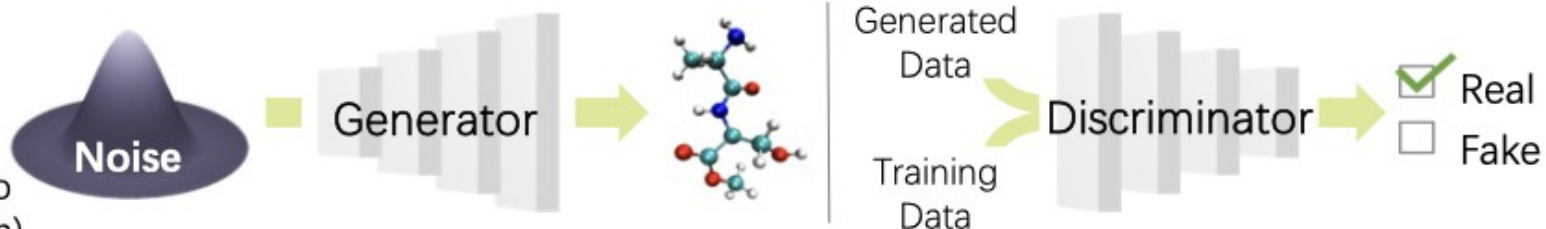
Normalizing Flow Model

Learn the parameters of a series of invertible, normalized transformations



Generative Adversarial Network

Discriminator Imposes Distribution Matching (no explicit $P(x)$ approximation)



We can cast molecular discovery as an **optimization problem**

3.1 Summary

Human Version:



Navigating the Molecular Design Space

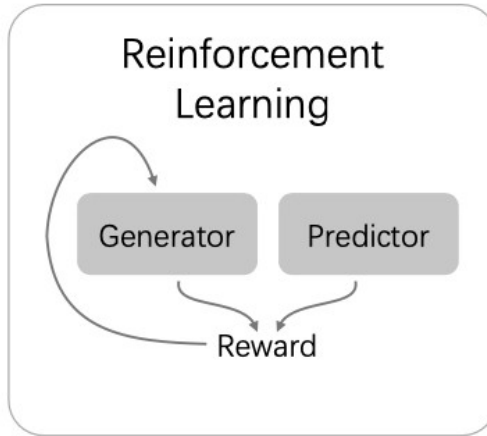
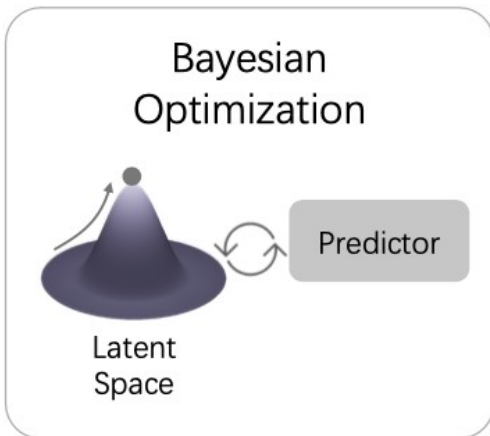
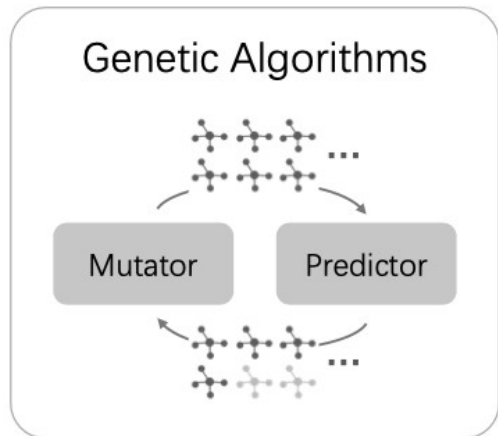
To turn this into an algorithmically well-defined problem we need:

- 1) A design space
- 2) Algorithms for making moves in that design space
- 3) An objective function to guide those algorithms

The molecular design space is inherently discontinuous, but can be mapped to a continuous space using generative models:



Some Optimization Methods:



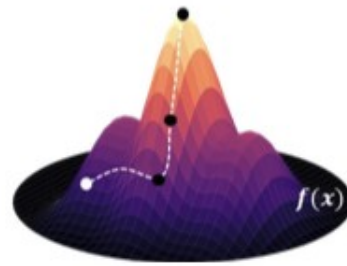
3.1 Summary

Generative modeling workflows rely on predictive models with limited certainty

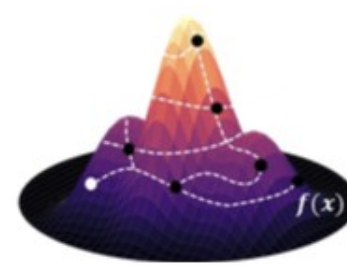
Solution: Incorporate information about uncertainty into your objective function

- Bayesian optimization methods yield probability distributions instead of discrete predictions making enabling straightforward uncertainty characterization
- Other uncertainty characterization methods can be incorporated directly into the objective function
 - For neural networks these can include methods such as Ensemble Variance or Monte Carlo Dropout

Embrace
uncertainty!



Exploitation (Risk Avoiding)



Exploration (Risk Seeking)

Note: Optimizing with a bad predictive model can limit exploration; it is better to avoid explicitly optimizing with bad models

3.1 Summary

The science of how to discover molecules

a bit technical again~

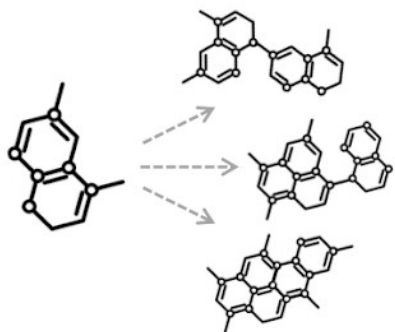


Structure-Constrained Molecular Generation

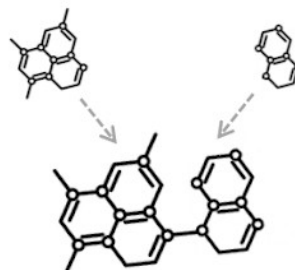
Goal: Take an input scaffold or structure and build on it to generate a new example.

Structure-constrained molecular generation **can be unconstrained OR property constrained.**

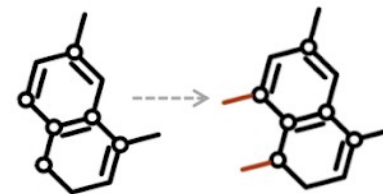
Scaffold/Fragment
Expansion



Scaffold/Fragment
Merging



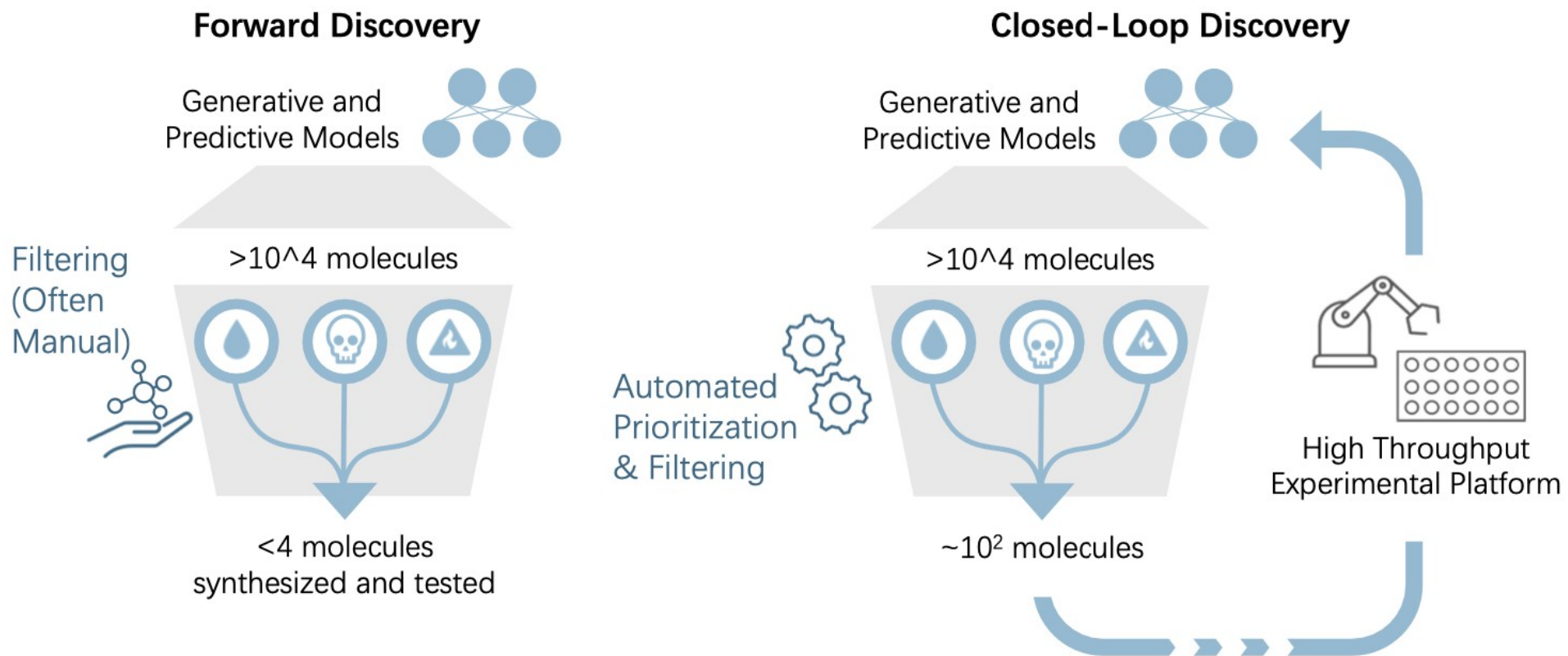
Molecular
Translation



3.1 Summary

The art of how to discover molecules

Molecular Generation in Real Applications



3.1 Summary

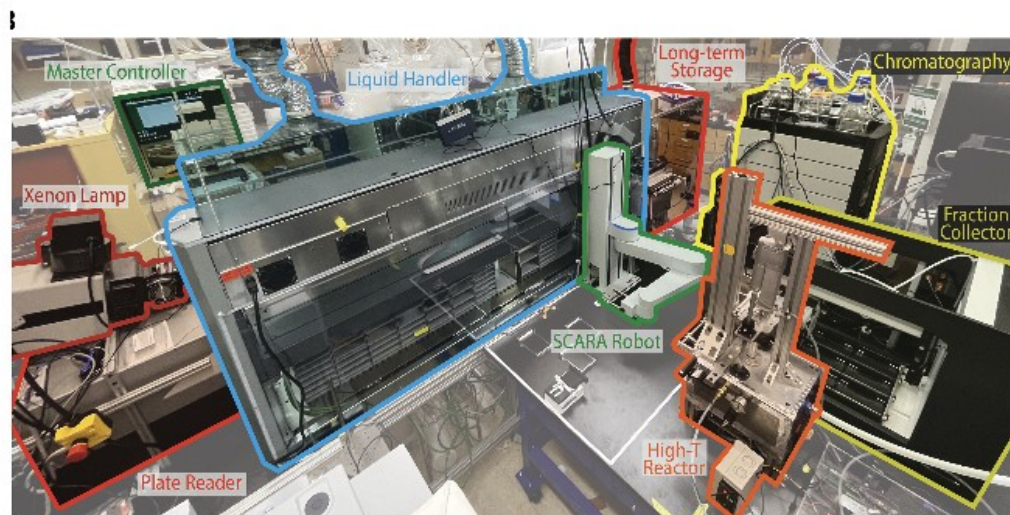
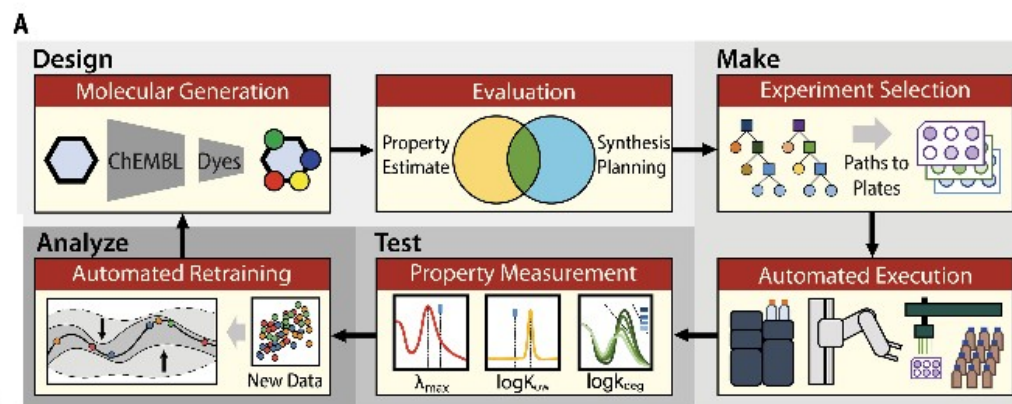
The art of how to discover molecules

Real Application: Closed Loop Discovery

Science

Autonomous, multiproperty-driven molecular discovery: From predictions to measurements and back

BRENT A. KOSCHER , RICHARD B. CANTY , MATTHEW A. McDONALD , KEVIN P. GREENMAN , CHARLES J. MCGILL, CAMILLE L. BILODEAU , WENGONG JIN, HAQYANG WU , FLORENCE H. VERMEIRE , BROOKE JIN , TRAVIS HART , TIMOTHY KULESZA, SHIH-CHENG LI , TOMMI S. JAAKKOLA , REGINA BARZILAY , RAFAEL GÓMEZ-BOMBARELLI , WILLIAM H. GREEN , AND KLAUS F. JENSEN  [fewer](#) [Authors Info & Affiliations](#)



3.2 Yoshua



Yoshua Bengio

3.2 Exploring Molecular Space & Active Learning

Professor: Yoshua Bengio Lecture
Recording: Lecture Slides:

0 06/14/2024

The slide features a circular portrait of Yoshua Bengio on the left and a decorative graphic of colorful spheres on the right. The background is dark purple.

1) It's a needle in the haystack problem

Discovery: Searching Needles in a Haystack

- *The space of experiments and the space of theories are huge*
- *Can't enumerate*
- ***Tiny fraction will work***
- ***How to get all good ones?***
- ***Amortized inference can help us to efficiently represent posteriors over theories, sample candidate experiments and reason about outcomes***



3.2 Summary

1) It's a needle in the haystack problem

In-Silico Design - Desiderata

We aim to combine several metrics:

- **Performance:** Measure of the desired effect of the generated candidates
- **Diversity:** Measure of the variety in generated candidates
- **Novelty:** Measure of novelty with respect to known candidates

$$\text{Mean}(\mathcal{D}) = \frac{\sum_{(x_i, y_i) \in \mathcal{D}} y_i}{|\mathcal{D}|}$$

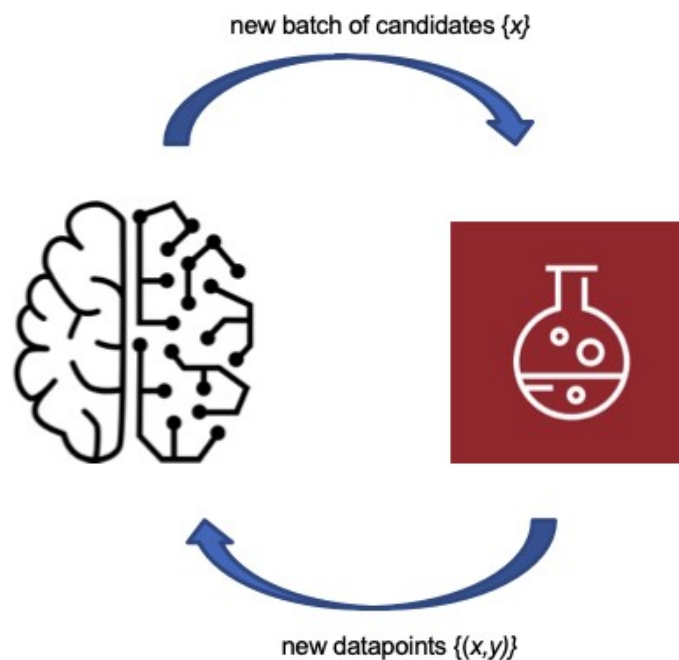
$$\text{Diversity}(\mathcal{D}) = \frac{\sum_{(x_i, y_i) \in \mathcal{D}} \sum_{(x_j, y_j) \in \mathcal{D} \setminus \{(x_i, y_i)\}} d(x_i, x_j)}{|\mathcal{D}|(|\mathcal{D}| - 1)}$$

$$\text{Novelty}(\mathcal{D}) = \frac{\sum_{(x_i, y_i) \in \mathcal{D}} \min_{s_j \in \mathcal{D}_0} d(x_i, s_j)}{|\mathcal{D}|}$$

3.2 Summary

Active Learning, Growing Dataset, ML and Assays Feedback Loop

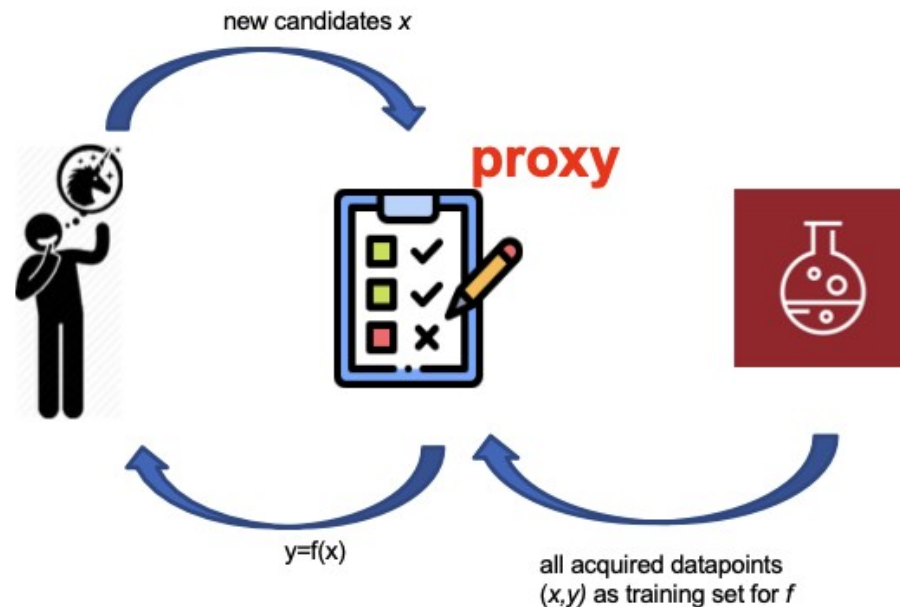
- ML system proposes next batch of candidates $\{x\}$
- Next-stage "oracle" provides more or less precise answers y per x
- Dataset is augmented with new set $\{(x,y)\}$ of pairs
- ML system retrained



3.2 Summary

Generating Candidates & Evaluating Them

Generating candidates can be obtained with a **generative model** (that we can draw samples from) trained to produce a diverse set of good candidates, where "good" is defined by an **evaluation function** trained with the data collected from a downstream "oracle", based on a **proxy** for the true property of interest



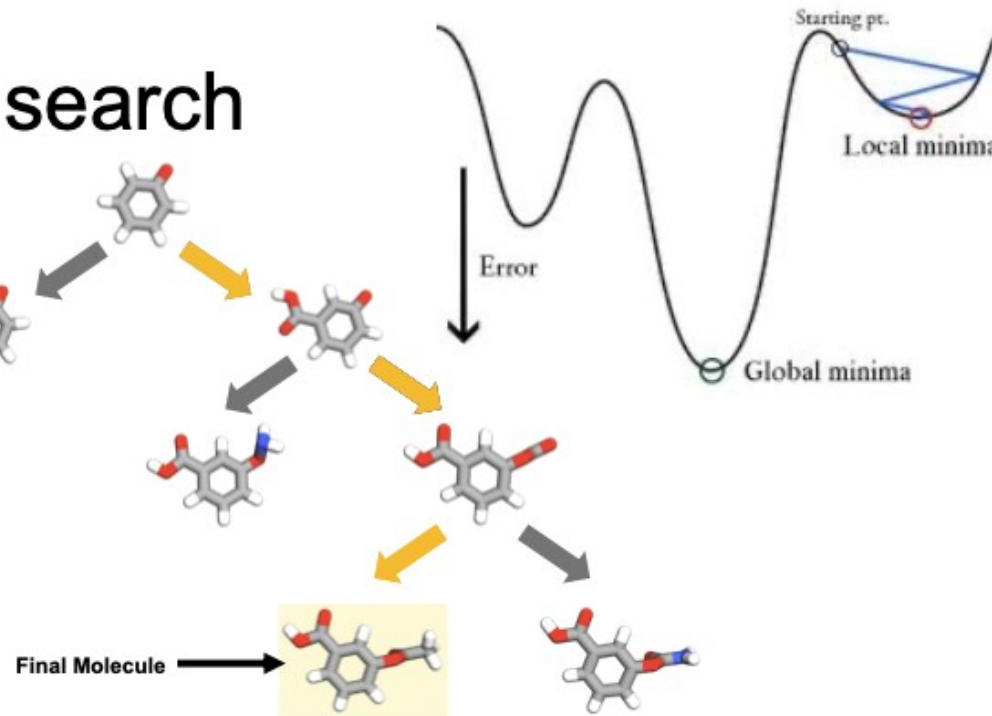
3.2 Summary

<https://yoshuabengio.org/gflownet-tutorial>

Papers at NeurIPS 2021, ICML 2022, UAI 2022, NeurIPS 2022, ICLR 2023, ICML 2023, NeurIPS 2023, ICML 2024 + ArXiv
Mila

Local discrete search

We can create a discrete action space which corresponds to transforming a molecule into another by adding or removing an atom or a frequently occurring rigid block of atoms.



Similarly we can search in the space of peptides by adding one amino-acid at a time

Black-Box ~~Optimization~~ Exploration & Scientific Discovery

- We have access to an expensive oracle (real-world), outputs $y=f(x)$
- We can only call it a few times, each round with batched queries B
- We look for a set of x 's whose $f(x)$ is large (modes of f)

- Examples:

- Discovering new drugs
- Discovering new materials
- Experimental design of experiments
- Reasoning & causal discovery:

discovering good explanations & causal models: $-\log p = \text{energy fn}$



3.3 Connor



3.3 Synthesizability & Molecular Synthesis

Professor: Connor Coley Lecture
Recording: Lecture Slides:

0 06/14/2024

The formulation of molecular optimization

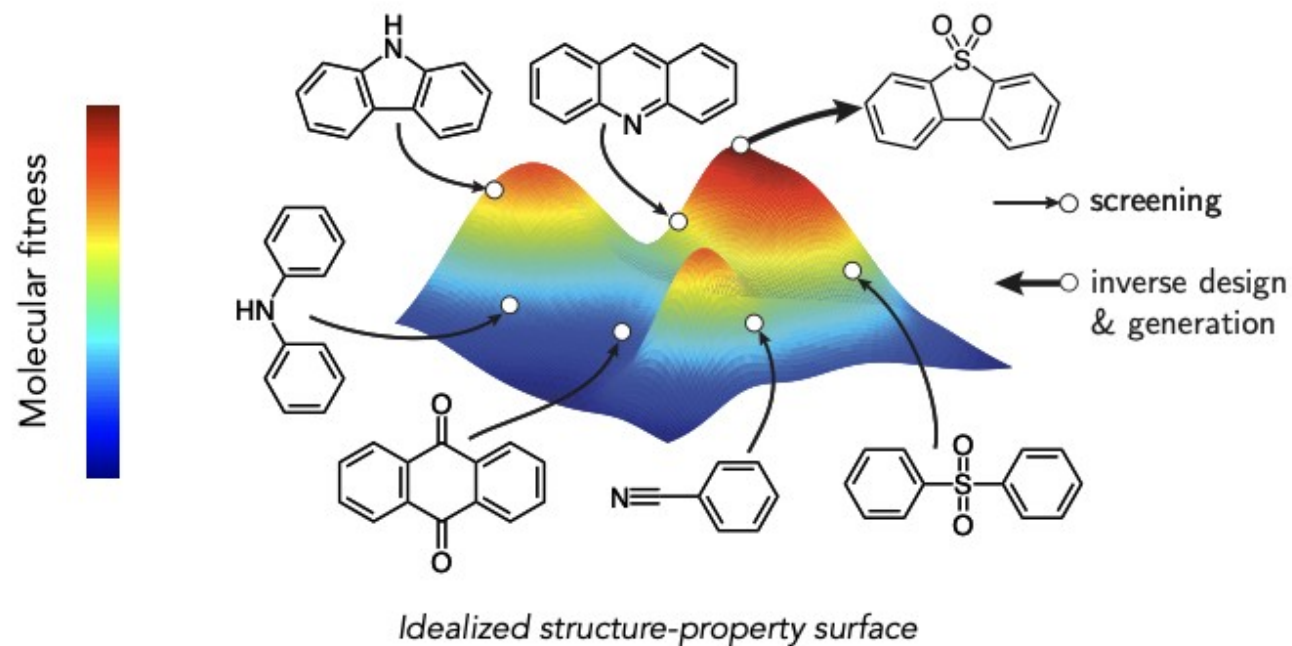
best design

oracle

$$x^* = \arg \max_{x \in \mathcal{X}} f(x)$$

all possible molecules

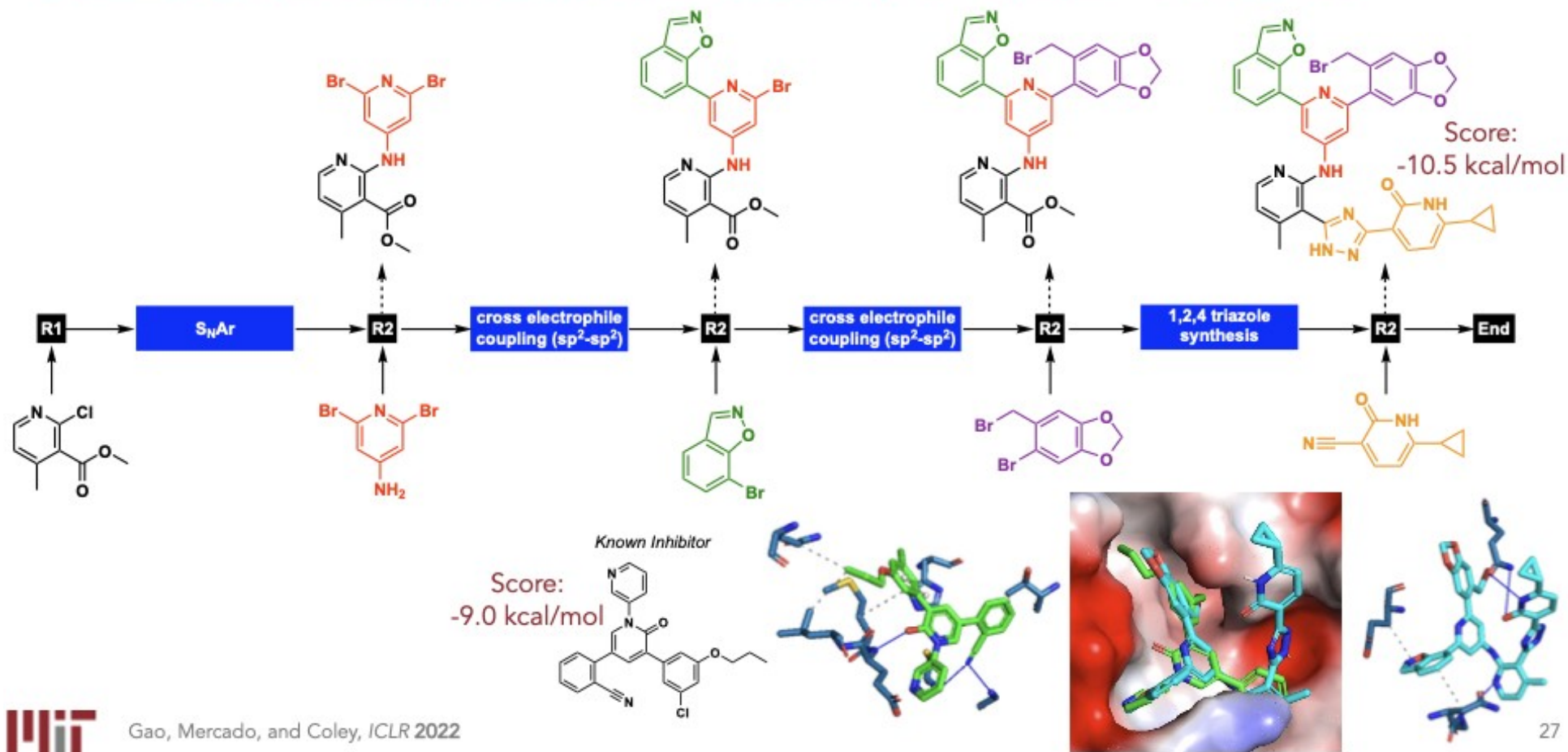
Molecular design and the ideation of structures



3.3 Summary

Building molecules by building synthetic trees

Optimization objective: design a molecule with a good docking score against M^{PRO} of SARS-CoV-2



3.3 Summary

What is the chemical space that these models search?

$O(100,000)$ commercial building blocks

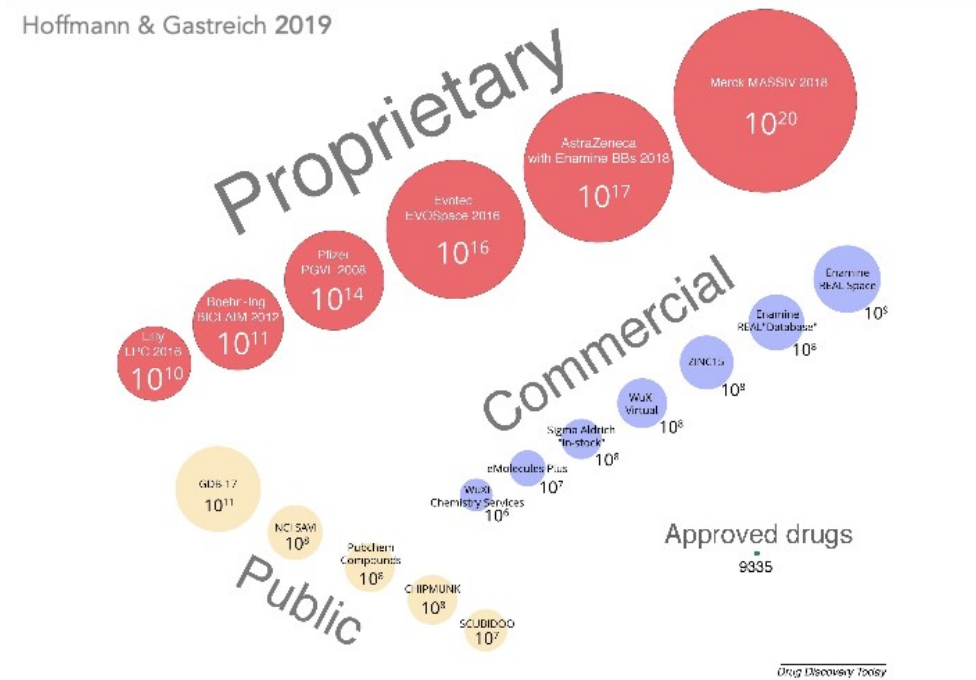
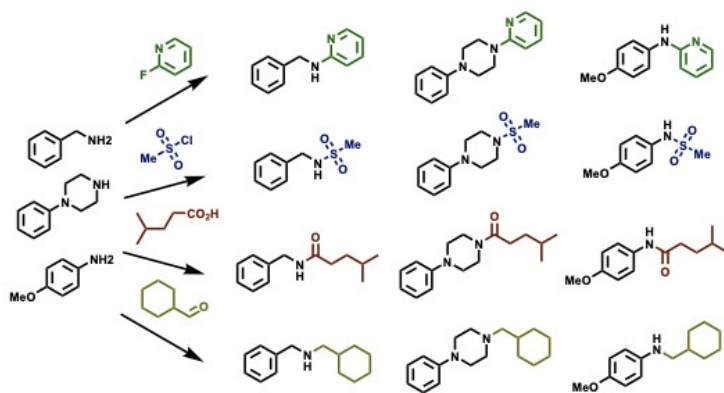
x

$O(100)$ expert-defined reaction templates

x

1-10 reaction steps

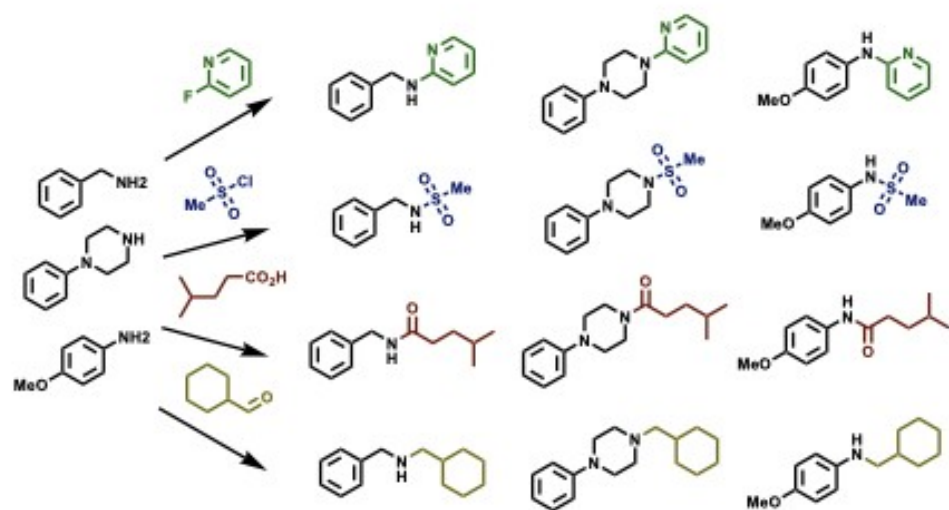
These types of models search the same kinds of chemical spaces as "make-on-demand" virtual libraries but avoid explicit enumeration



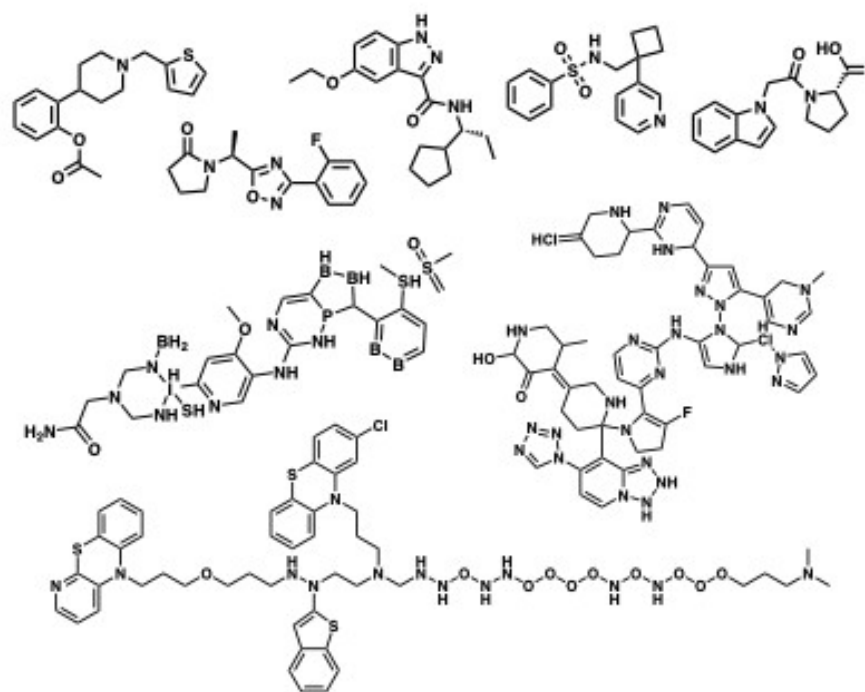
3.3 Summary

Summary: don't overlook synthesis in molecular design

Virtual libraries are often “make-on-demand” libraries enumerated using chemical transformation rules we believe to be robust



Generative models produce new compounds for which we must plan synthetic routes unless already constrained by synthesis



3.4 Michael



Michael Bronstein

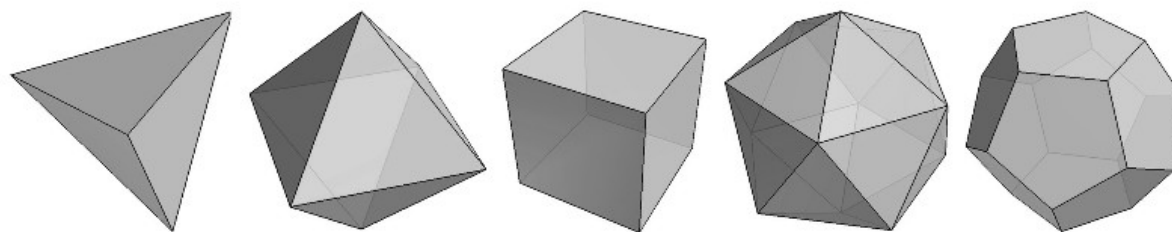
3.4 Sampling Physical & 3D Systems

Professor: Michael Bronstein
Lecture Recording: Lecture
Slides:

1 06/14/2024

geometric deep learning
graph representation learning

“Symmetry, as wide or as narrow as you may define its meaning, is one idea by which man through the ages has tried to comprehend and create order, beauty, and perfection”



“Platonic solids”

“It is only slightly overstating the case to say that Physics is the study of symmetry”

— More is different



H. Weyl



Plato

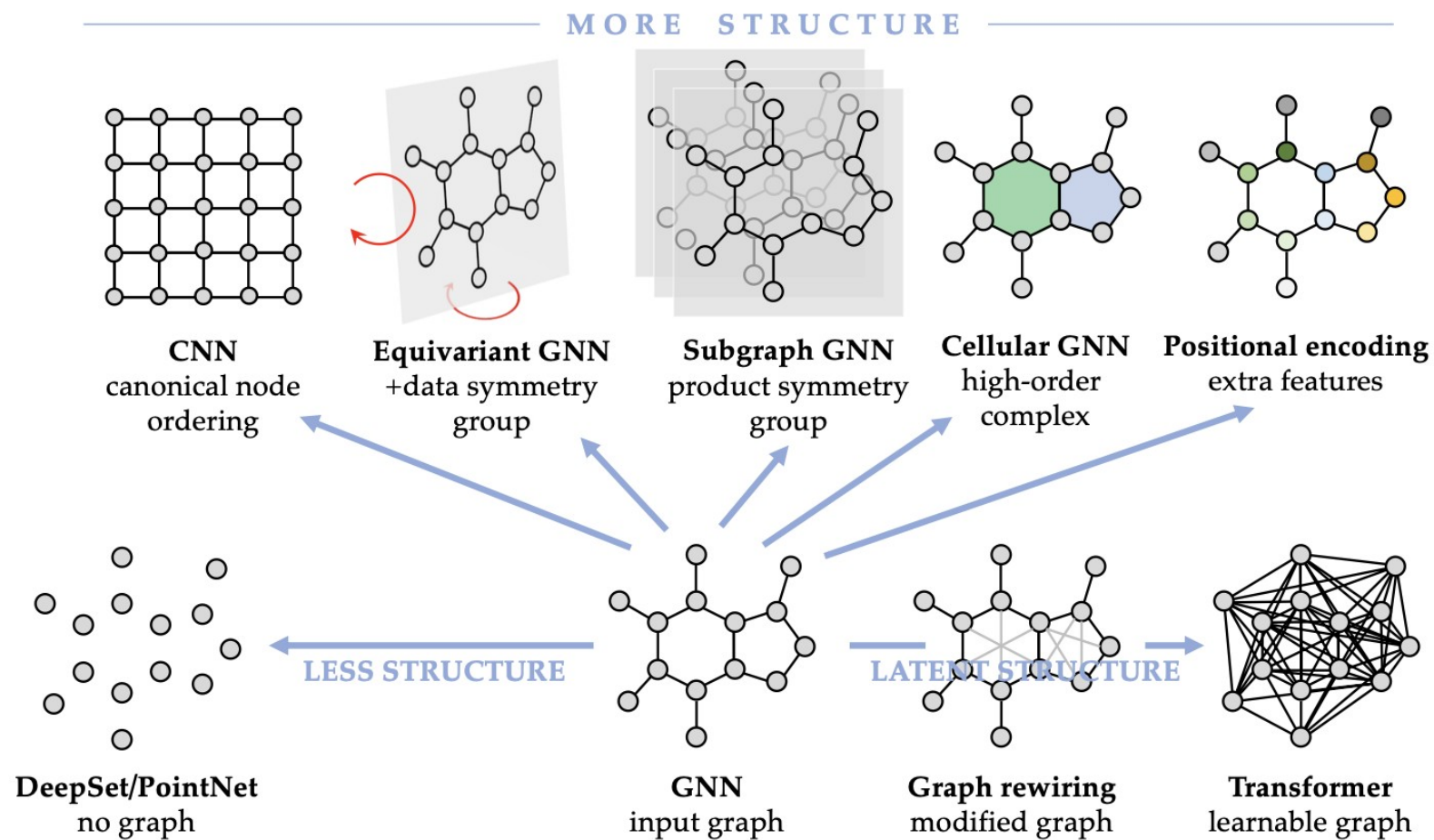
370 BC



P. Anderson

3.4 Summary

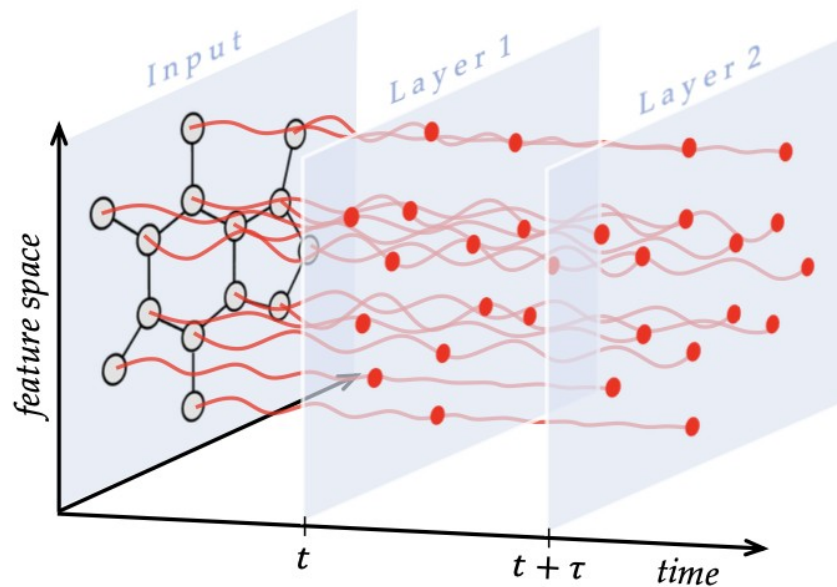
GNNs is (possibly) the parent class of all NNs,
developed around the idea of geometric symmetry :



3.4 Summary

Physics-inspired GNNs are also cool :

Physical metaphor of Graph ML



- GNN = dynamic system
- layers = discretisation of time
- graph = coupling function (discretisation of space)

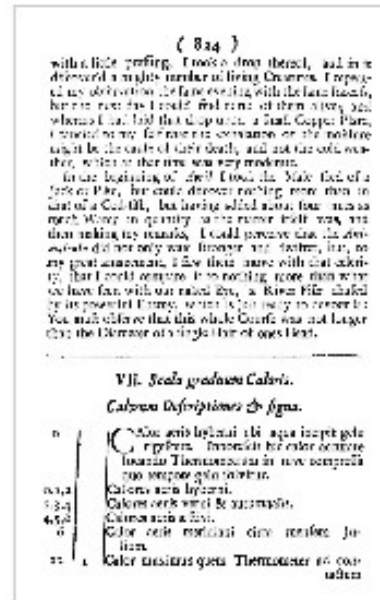
$$\mathbf{X}(t + \tau) = \mathbf{X}(t) + \tau \mathbf{F}_{\theta(t)}(\mathbf{X}(t), \mathcal{G})$$

3.4 Summary

Heat Diffusion

Newton Law of Cooling:
“the [temperature] a hot body loses in a given time is proportional to the temperature difference between the object and the environment”

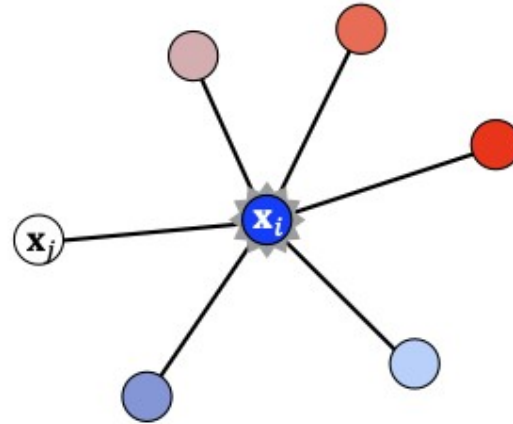
Anonymous 1701



I. Newton

3.4 Summary

Heat Diffusion Equation on Graphs

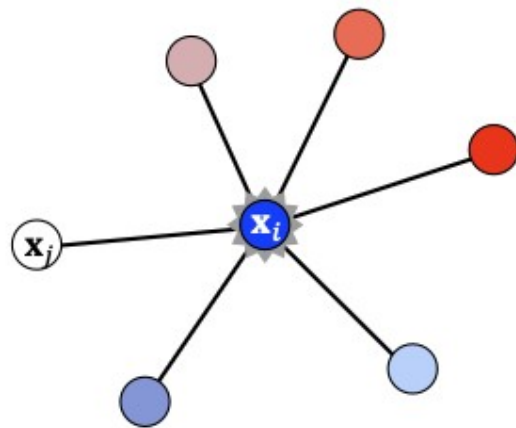


$$\dot{\mathbf{x}}_i(t) = \mathbf{x}_i(t) - \frac{1}{d_i} \sum_{j \in \mathcal{N}_i} a_{ij} \mathbf{x}_j(t)$$

rate of temperature change self temperature temperature of the environment

3.4 Summary

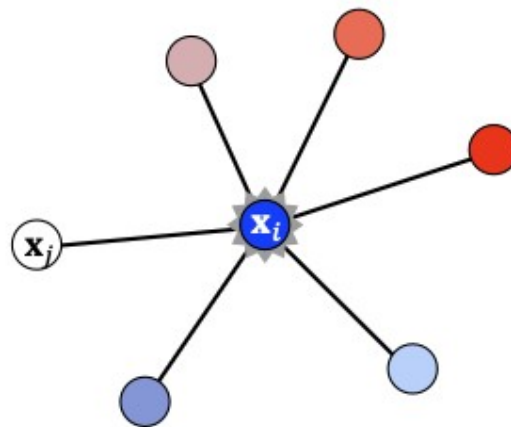
Heat Diffusion Equation on Graphs



$$\dot{\mathbf{x}}_i(t) = \underbrace{\frac{1}{d_i} \sum_{j \in \mathcal{N}_i} a_{ij}}_{\substack{\text{divergence} \\ \text{div}}} \underbrace{\left(\mathbf{x}_i(t) - \mathbf{x}_j(t) \right)}_{\substack{\text{gradient} \\ - (\nabla \mathbf{X})_{ij}}}$$

3.4 Summary

Heat Diffusion Equation on Graphs



$$\dot{\mathbf{x}}_i(t) = \underbrace{\frac{1}{d_i} \sum_{j \in \mathcal{N}_i} a_{ij}}_{\text{divergence div}} \underbrace{(\mathbf{x}_i(t) - \mathbf{x}_j(t))}_{\text{gradient } -(\nabla \mathbf{X})_{ij}}$$

$$\dot{\mathbf{X}}(t) = \Delta \mathbf{X}(t)$$

3.4 Summary

Michael concludes:


Conclusions

- Geometric deep learning is a powerful tool for molecular modelling
- Rootes in fundamental mathematical principles of invariance & symmetry
- Geometric generative models models
- Experimental validation!

4.1 Anne

But how do we know what protein to target??

It's... messy

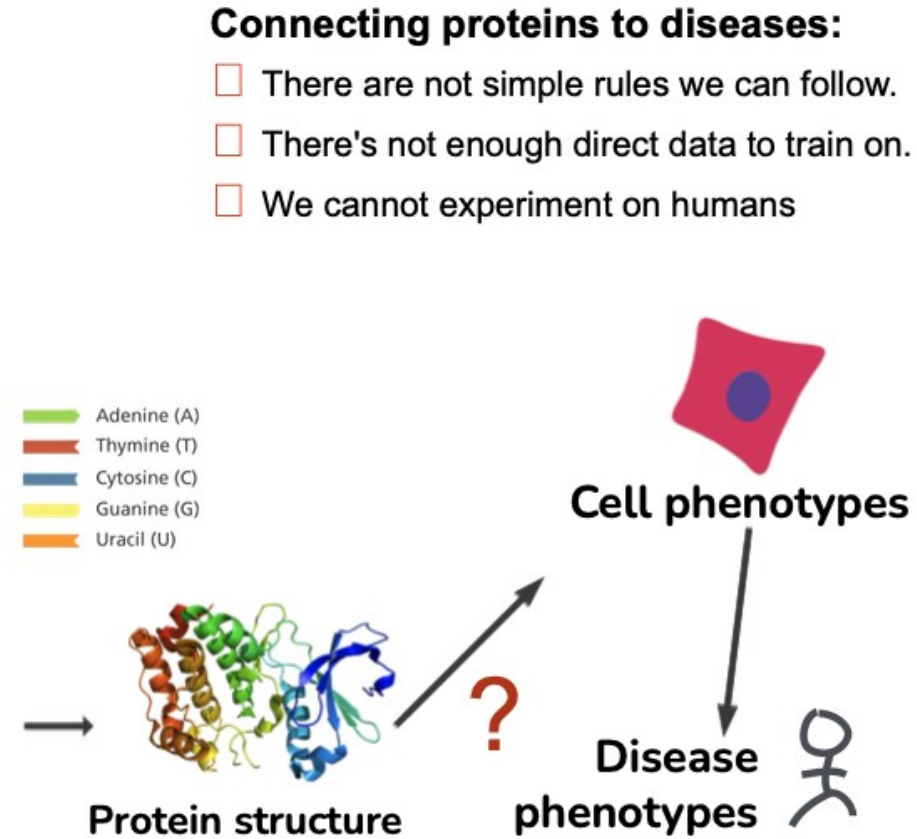
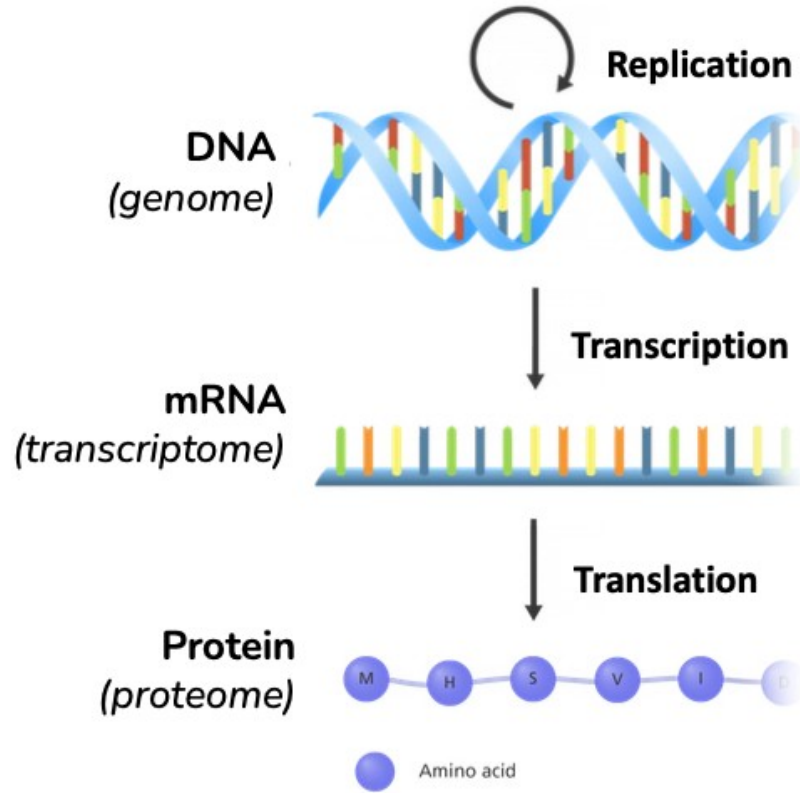


Anne E. Carpenter

4.1 Phenomics in Drug Discovery

Professor: Anne Carpenter
Lecture Recording: Lecture Slides:

0 06/17/2024



Connecting proteins to diseases:

- There are not simple rules we can follow.
- There's not enough direct data to train on.
- We cannot experiment on humans

4.1 Summary

Cell painting assays can be used for drug discovery:

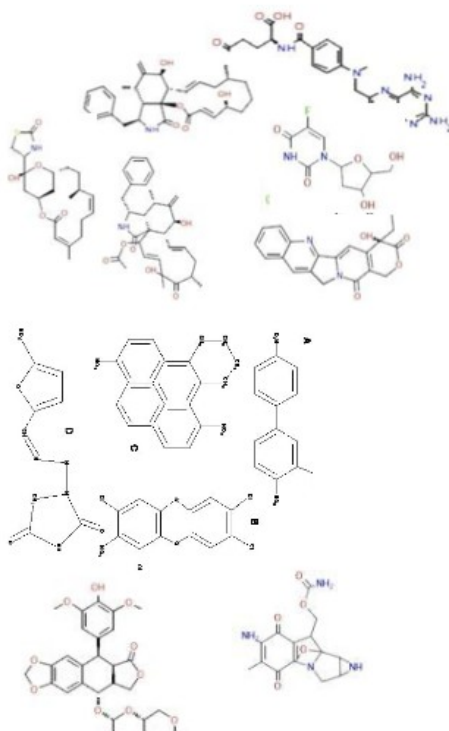
Drug discovery using cells as proxies



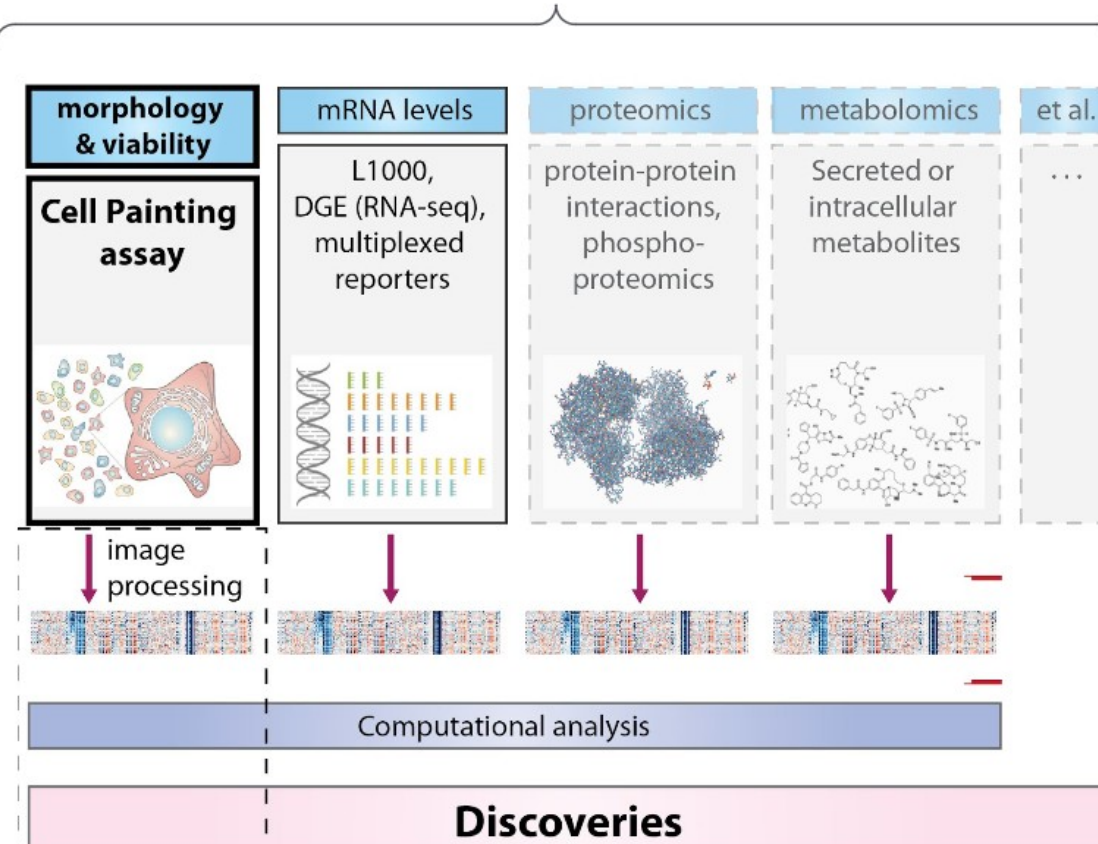
thousands of genetic perturbations



or millions of chemicals



profiling assays

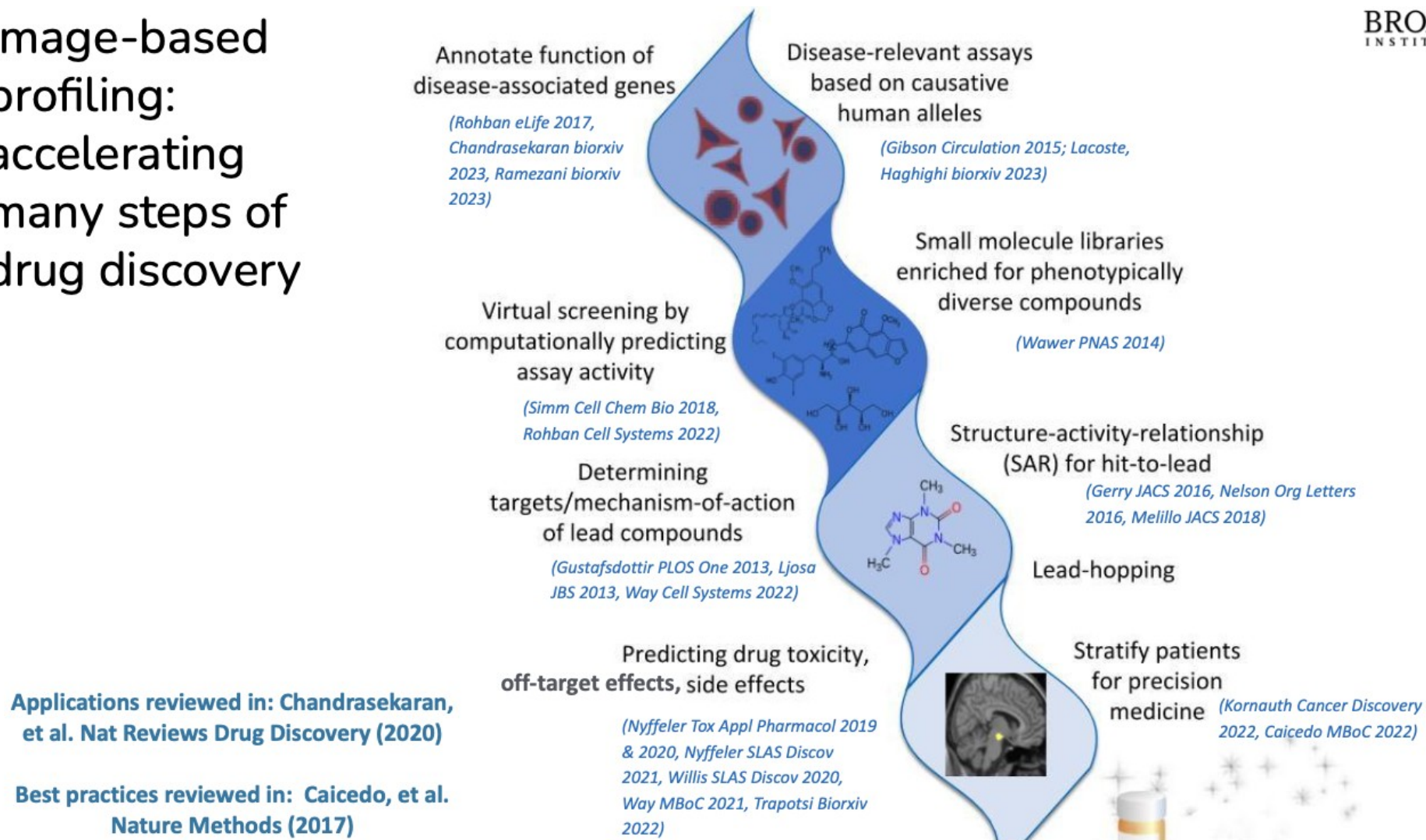


4.1 Summary

Image-based profiling can be applied to many different steps:

Image-based profiling:
accelerating
many steps of
drug discovery

BROAD IMAGING
INSTITUTE PLATFORM



4.1 Summary

There are substantial success of image based profiling such as:

More successes of image-based profiling

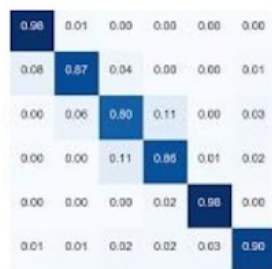


Applications reviewed in: Chandrasekaran, Ceulemans, Boyd, Carpenter Nat Reviews Drug Discovery (2020)

Best practices reviewed in: Caicedo, et al. Nature Methods (2017)

Annual conferences: CytoData, SBI2, Images 2 Knowledge (I2K)

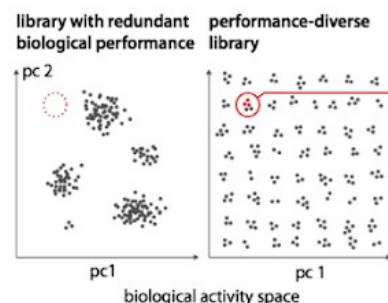
Label-free imaging flow cytometry diagnostics



classify cell cycle, red blood cell lesions, leukemic progression and response

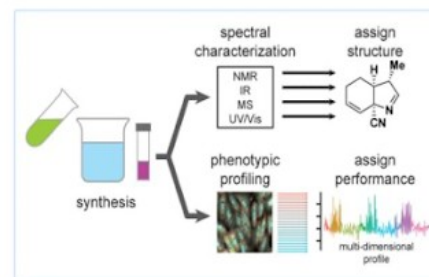
Blasi, et al. Nature Comm. 2016
Hennig, et al. Methods 2016
Eulenberg, et al. Nature Comm. 2017
Doan, et al. Trends Biotech 2018

Create performance-diverse screening libraries



improved hit rate
50% increase vs randomly chosen compounds

Test bioactivity for newly synthesized compounds




improved turnaround time
from years to days

Gerry, et al. JACS 2016
Nelson, et al. Org Lett 2016

Predict compound toxicity

OASIS Consortium launched!

4.2 Sebastien



Sébastien Lemieux

4.2 Multi-Modal Omics & AI

Professor: Sébastien Lemieux
Lecture Recording: Lecture
Slides:

0 06/17/2024

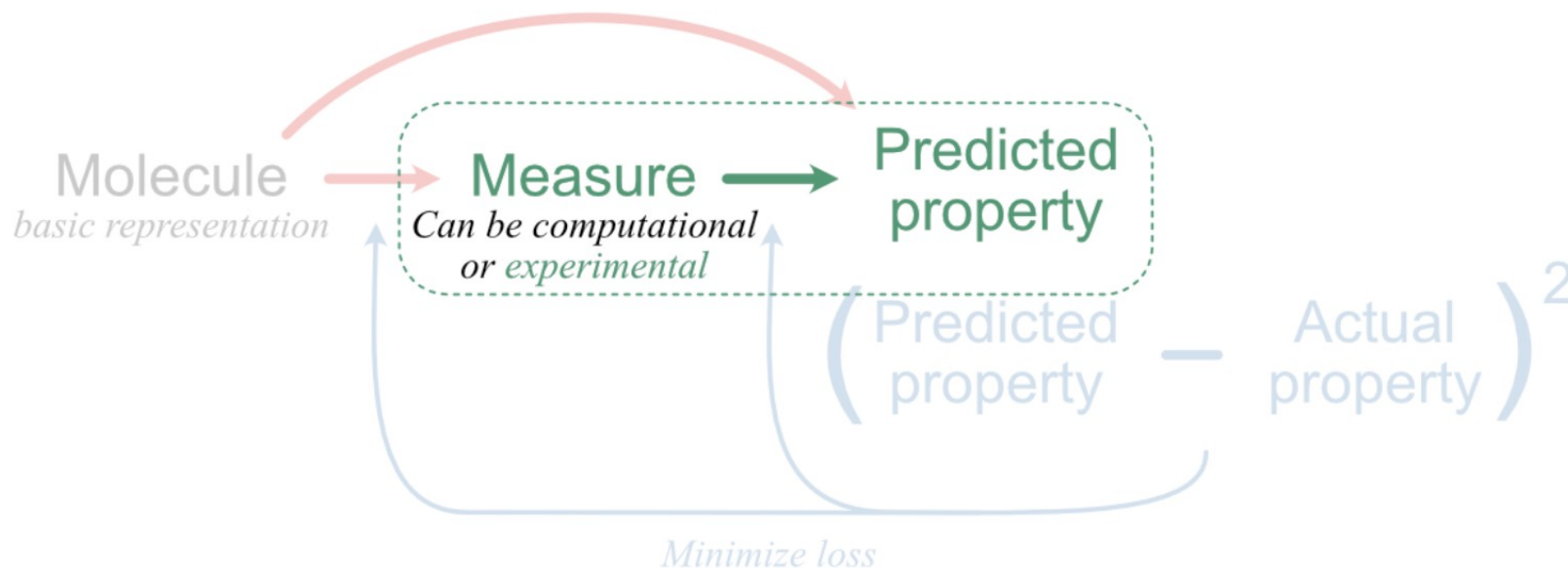
What's omics?

- Emerges in the mid-90s:
 - with microarrays, shotgun sequencing, proteome-wide mass spectrometry;
 - with advances closely tied with bioinformatics / computational biology.
- Large-scale datasets generated by high-throughput technologies.
 - Typically, 10^3 to 10^6 measures per biological sample.
- This data encompasses a broad range of biological fields:
 - Genomics (DNA)
 - Transcriptomics (mRNA)
 - Proteomics (proteins)
 - Metagenomics (microbial communities), metabolomics (metabolites), epigenomics (non sequence-based DNA modifications), etc.
- Such data aims to provide a holistic view of the molecular components to facilitate a deeper *understanding* of complex biological processes.

4.2 Summary

What are the tasks of multimodal learning in omics >

Tasks



We'll focus on experimental "measures"

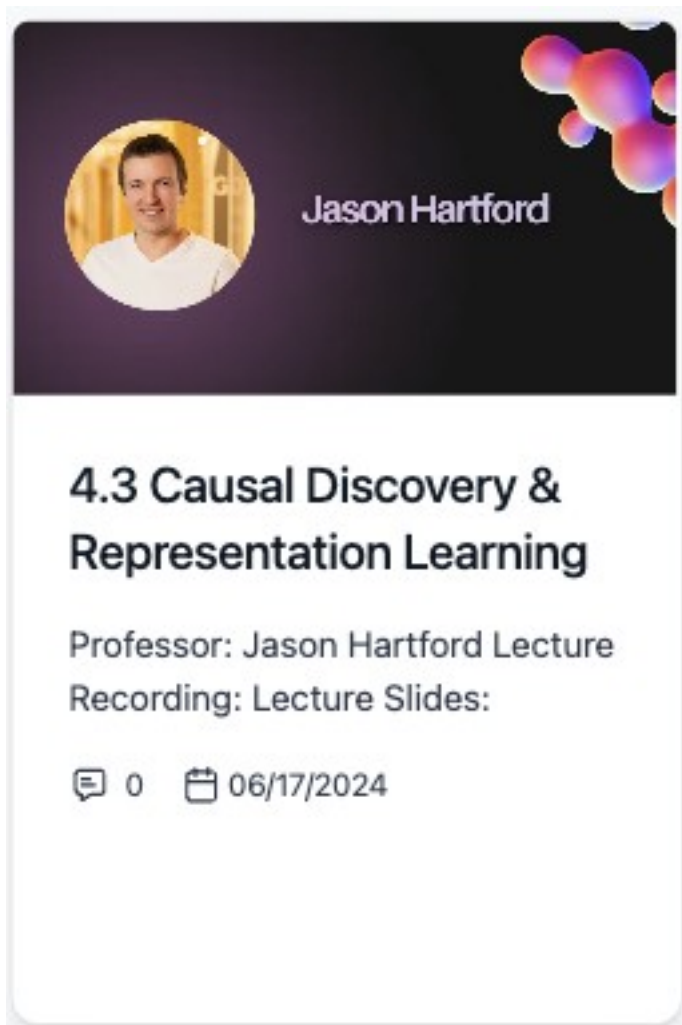
The molecule must be synthesized and assayed

Promises a much richer characterization of the molecule

Take-home message

- Biological prediction tasks need to account for an infinite number of input dimensions, showing an important level of correlation.
 - Selection is always made, *ad hoc*, at the selection of experimental data generation, for availability, cost or feasibility reasons.
- Use of multimodal approaches in omics data usage is still poorly studied, despite being desired for at least 20 years.¹
 - Still today, the field is dominated by late fusion and discretization, reusing modality-specific analyses, often designed for other uses.
- Current, direct application to drug discovery are still few.
 - Target deconvolution.^{2,3}
 - Drug activity prediction from L1000 expression profiles.⁴
- Use of factorized embeddings offers clear opportunity in omics data integration.^{5,6}

4.3 Jason



Jason Hartford

4.3 Causal Discovery & Representation Learning

Professor: Jason Hartford Lecture Recording: Lecture Slides:

0 06/17/2024

A Tutorial on Causal Representation Learning | Jason Hartford & Dhanya Sridhar



Valence Labs
7.64K subscribers

Subscribe



Valence Labs

@valence_labs · 7.64K subscribers · 241 videos

Harnessing computation to radically improve lives. ...more

Subscribe

Home Videos Live Playlists Search

Created playlists



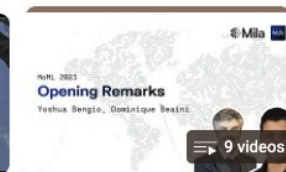
CARE: Causality, Abstraction, Reasoning & Extrapolation...

View full playlist



Valence Labs Launch

View full playlist



2023 Molecular Machine Learning Conference

View full playlist



M2D2: Molecular Modeling and Drug Discovery

View full playlist

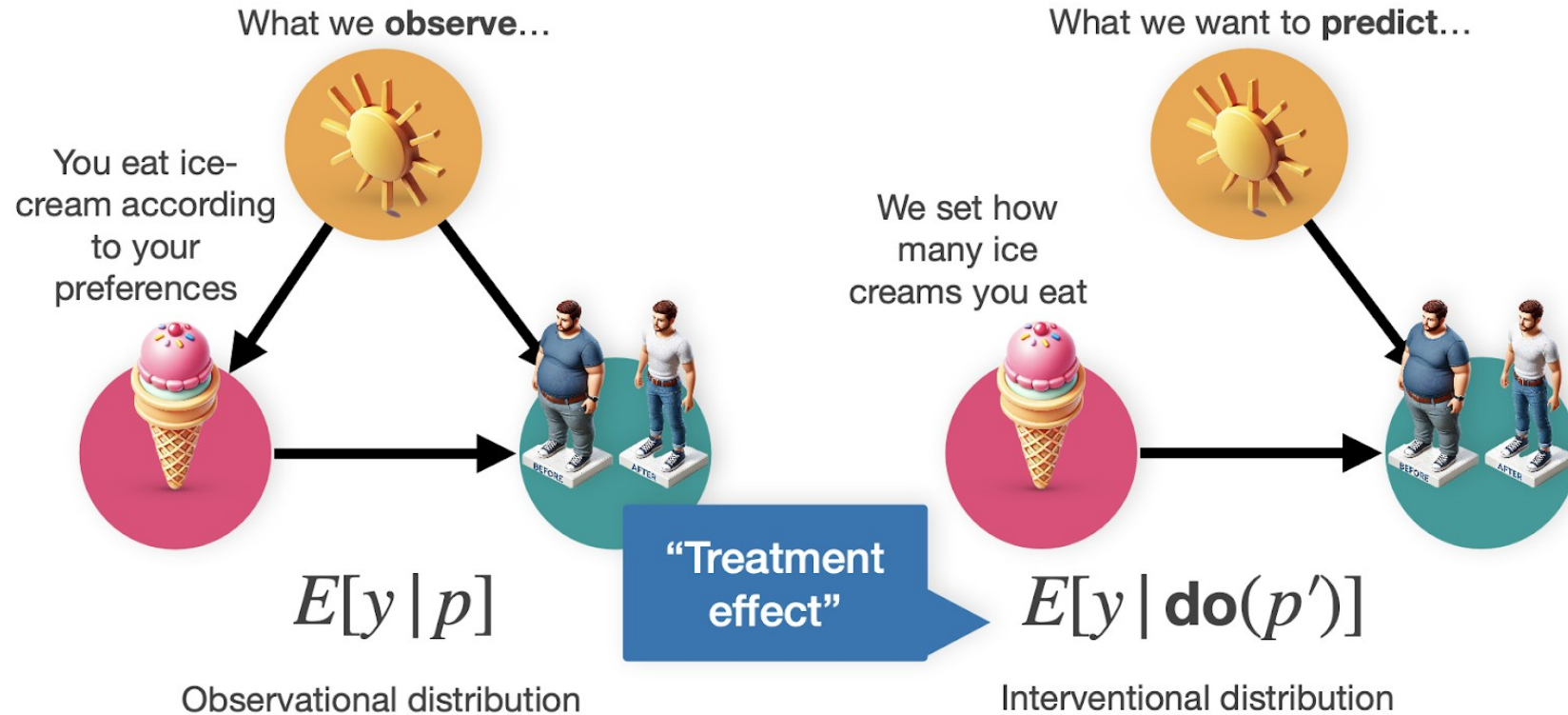


Graphs and Geometry Reading Group

Updated yesterday
View full playlist

4.3 Summary

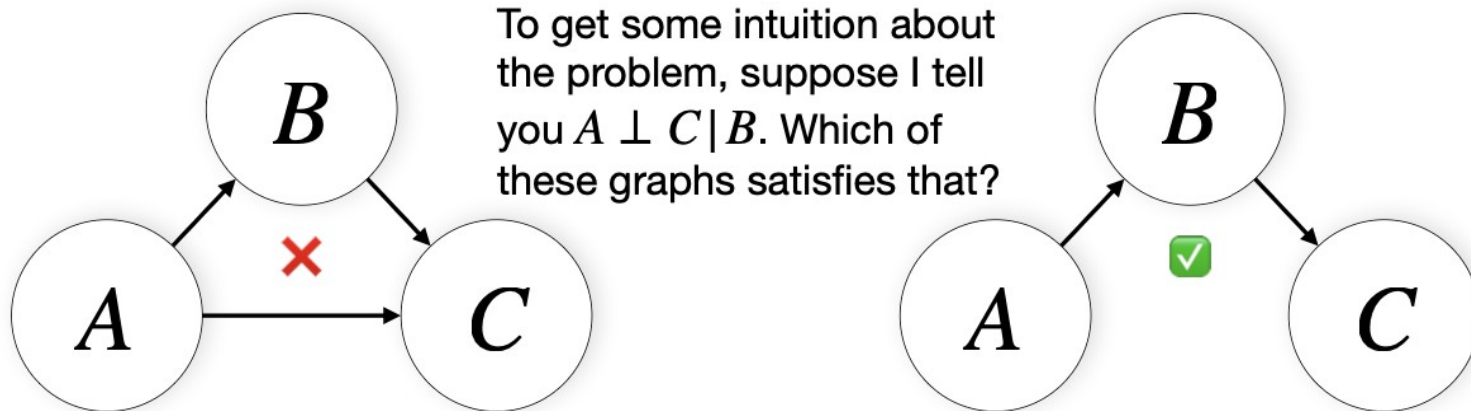
Observations vs interventions



Causal discovery

d -separation allows to infer conditional independencies from the graphical structure implied by our assumptions.

Causal discovery goes in the opposite direction: given a set of conditional independencies, what does that imply about the underlying graph?

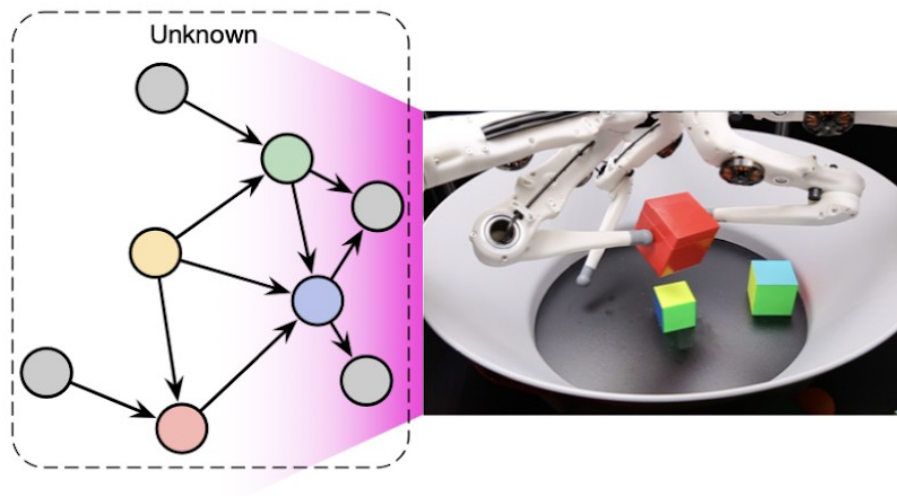


Causal representation learning

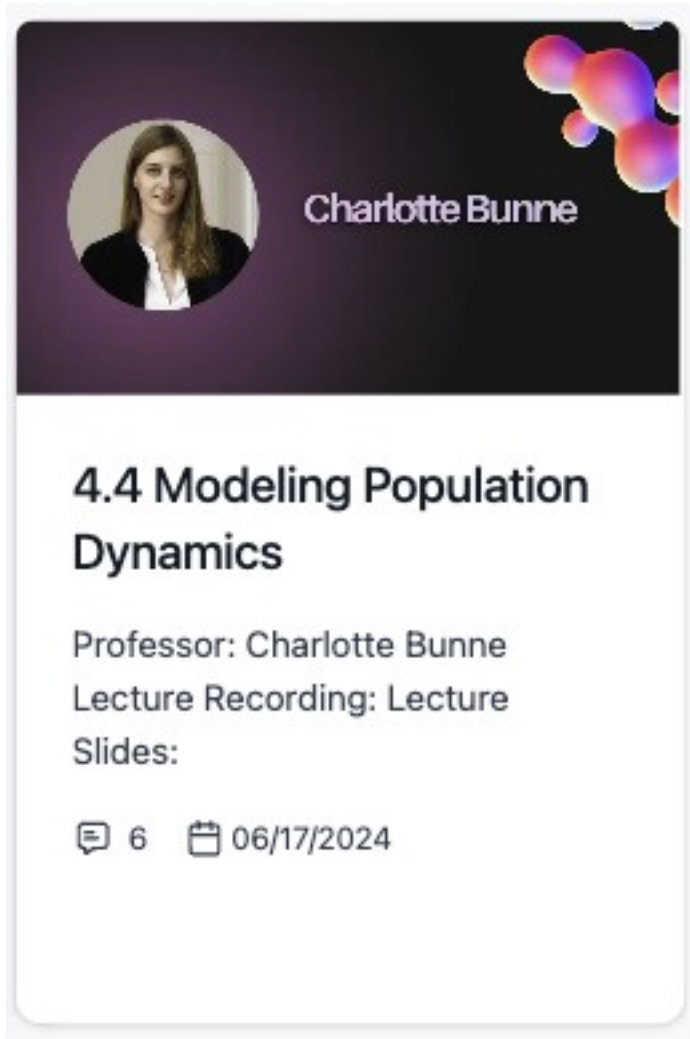
- Everything we've seen so far assumed semantically meaningful variables are observed. What do we do if we just have unstructured data like the pixels in an image?
- **Today:** why is this problem hard? What do existing techniques do?

Towards Causal Representation Learning

Bernhard Schölkopf [†], Francesco Locatello [†], Stefan Bauer ^{*}, Nan Rosemary Ke ^{*}, Nal Kalchbrenner
Anirudh Goyal, Yoshua Bengio



4.4 Charlotte



A lecture slide thumbnail for Charlotte Bunne. The top half features a dark purple background with a circular portrait of Charlotte Bunne on the left and a cluster of colorful, glowing spheres on the right. The name 'Charlotte Bunne' is written in white text to the right of the portrait. Below the image, the title '4.4 Modeling Population Dynamics' is displayed in a large, bold, black font. Underneath the title, the text 'Professor: Charlotte Bunne' and 'Lecture Recording: Lecture Slides:' is shown in a smaller black font. At the bottom left, there is a speech bubble icon with the number '6' and a calendar icon with the date '06/17/2024'.

4.4 Modeling Population Dynamics

Professor: Charlotte Bunne
Lecture Recording: Lecture Slides:

6 06/17/2024

Static Optimal Transport as Inductive Bias for Reconstructing Molecular Responses in Biomedicine

Dynamic Optimal Transport and Connections to Diffusion Models and Flow Matching

Diffusion Models and Flow Matching for Reconstructing Dynamic Processes in Molecular Medicine

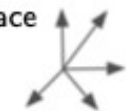
Optimal Transport

4.4 Summary

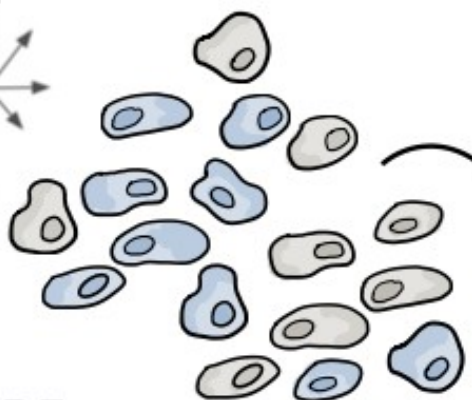
A Dynamic Perspective on Optimal Transport

Next: **Dynamic Optimal Transport**

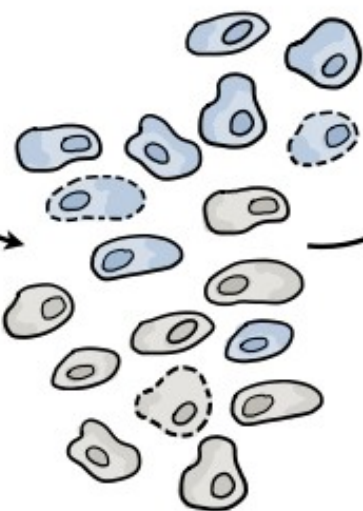
cell data space



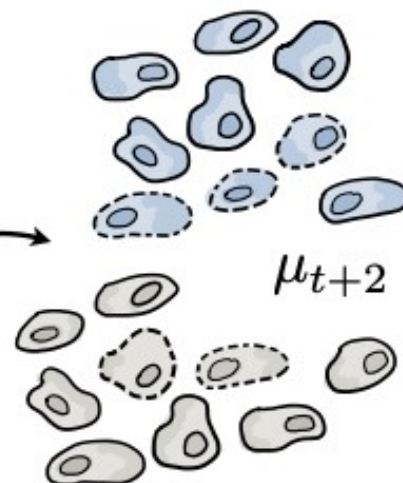
μ_t



μ_{t+1}

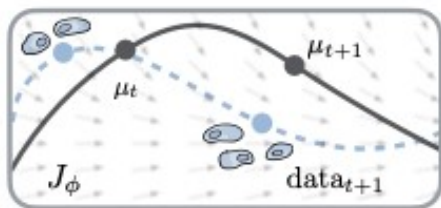


μ_{t+2}



via PDEs

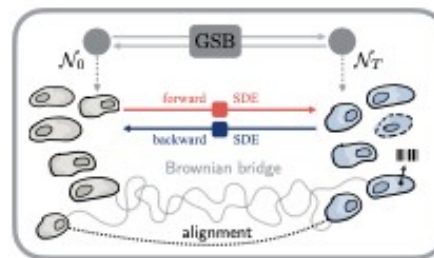
JKONET
AISTATS'22



via SDEs

GSBFLOW
AISTATS'23

SBALIGN
UAI'23



time

4.4 Summary

Solving Schrödinger Bridges

Schrödinger bridge optimality given by

(Léonard 2013)

$$* \begin{cases} \frac{\partial \Phi}{\partial t} = -\nabla \Phi^\top f - \frac{1}{2} \sigma^2 g^2 \Delta \Phi \\ \frac{\partial \hat{\Phi}}{\partial t} = -\nabla \cdot (\hat{\Phi} f) + \frac{1}{2} \sigma^2 g^2 \Delta \hat{\Phi} \end{cases} \quad \text{s.t.} \quad \begin{cases} \Phi(0, \cdot) \hat{\Phi}(0, \cdot) = \mu_0 \\ \Phi(1, \cdot) \hat{\Phi}(1, \cdot) = \mu_1 \end{cases}$$

Solution of $*$ can be expressed via two SDEs of the form

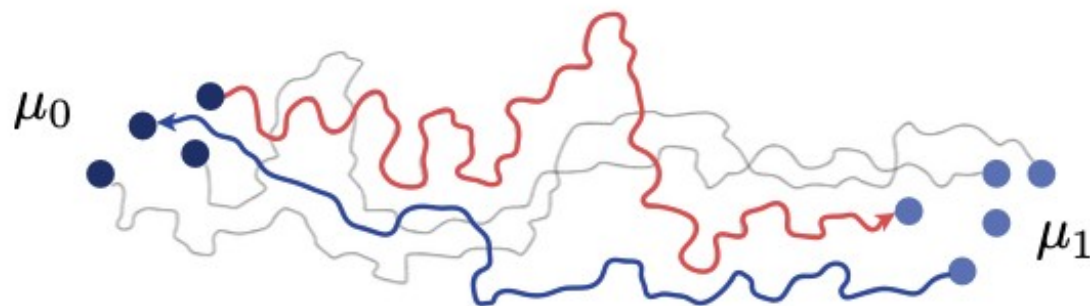
$$\begin{aligned} dX_t &= \left[f + g^2 \overbrace{\nabla \log \Phi(t, X_t)}^{\text{forward policy}} \right] dt + g d\mathbb{W}_t, & X_0 &\sim \mu_0 \\ dX_t &= \left[f - g^2 \underbrace{\nabla \log \hat{\Phi}(t, X_t)}_{\text{backward policy}} \right] dt + g d\mathbb{W}_t, & X_1 &\sim \mu_1 \end{aligned}$$

4.4 Summary

Solving Schrödinger Bridges

The solution of the SB is a system of a forward and backward SDE

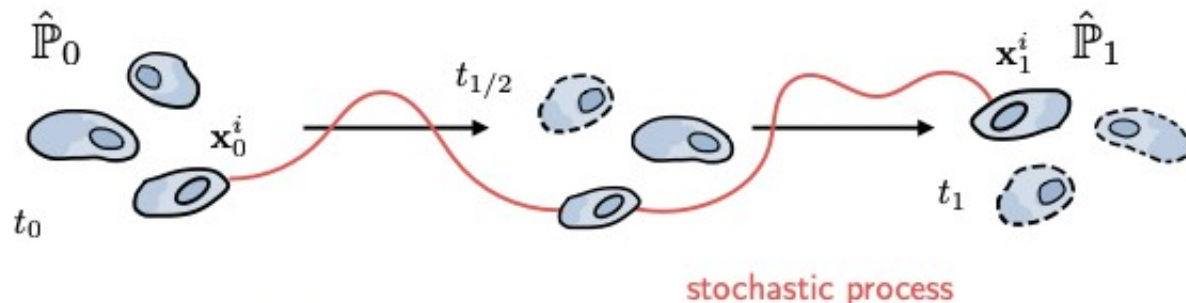
$$X_0 \sim \mu_0 \quad dX_t = [f + g^2 \overset{\text{forward policy}}{\nabla \log \Phi(t, X_t)}] dt + g dW_t$$



$$\text{reverse SDE} \quad dX_t = [f - g^2 \underset{\text{backward policy}}{\nabla \log \hat{\Phi}(t, X_t)}] dt + g dW_t,$$

(Anderson 1982)

Reconstructing Dynamics from Partial Alignments



$$\text{couplings over } \hat{\mathbb{P}}_0 \text{ and } \hat{\mathbb{P}}_1 \quad \pi^* := \underset{\mathbb{P}_0 = \hat{\mathbb{P}}_0, \mathbb{P}_1 = \hat{\mathbb{P}}_1}{\operatorname{argmin}} D_{\text{KL}}(\mathbb{P}_{0,1} \parallel \mathbb{Q}_{0,1})$$

reference process

Schrödinger's Bridge

We can learn the drift describing cellular dynamics by minimizing the resulting loss function

$$\min_{\theta} \mathbb{E} \left[\int_0^1 \left\| \frac{\mathbf{x}_1 - X_t}{\beta_1 - \beta_t} - (b_t^\theta + \nabla \log h_t^\theta(X_t)) \right\|^2 dt \right]$$

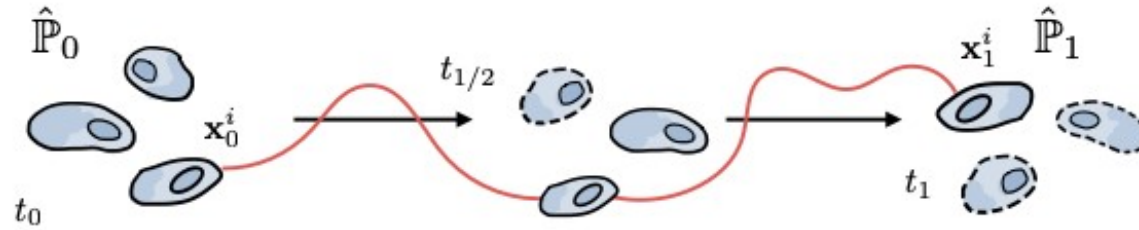
drift from Brownian bridge drift describing cellular dynamics

.. similar to score-matching objective successfully employed in previous works (Song et al., 2021)

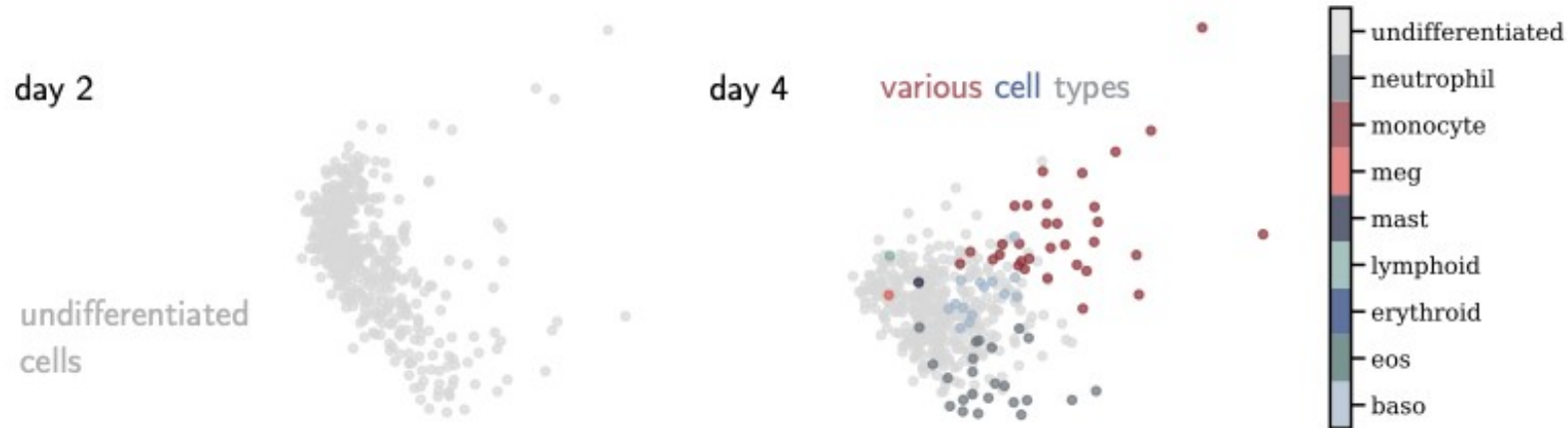
4.4 Summary

Recovering *Continuous* Dynamics of Cell Differentiation

UAI'23





Application: Fate determination in hematopoiesis.



.. dataset by Weinreb et al., Science (2020)

5.1 Andres

The only talk with a Colab Demo!



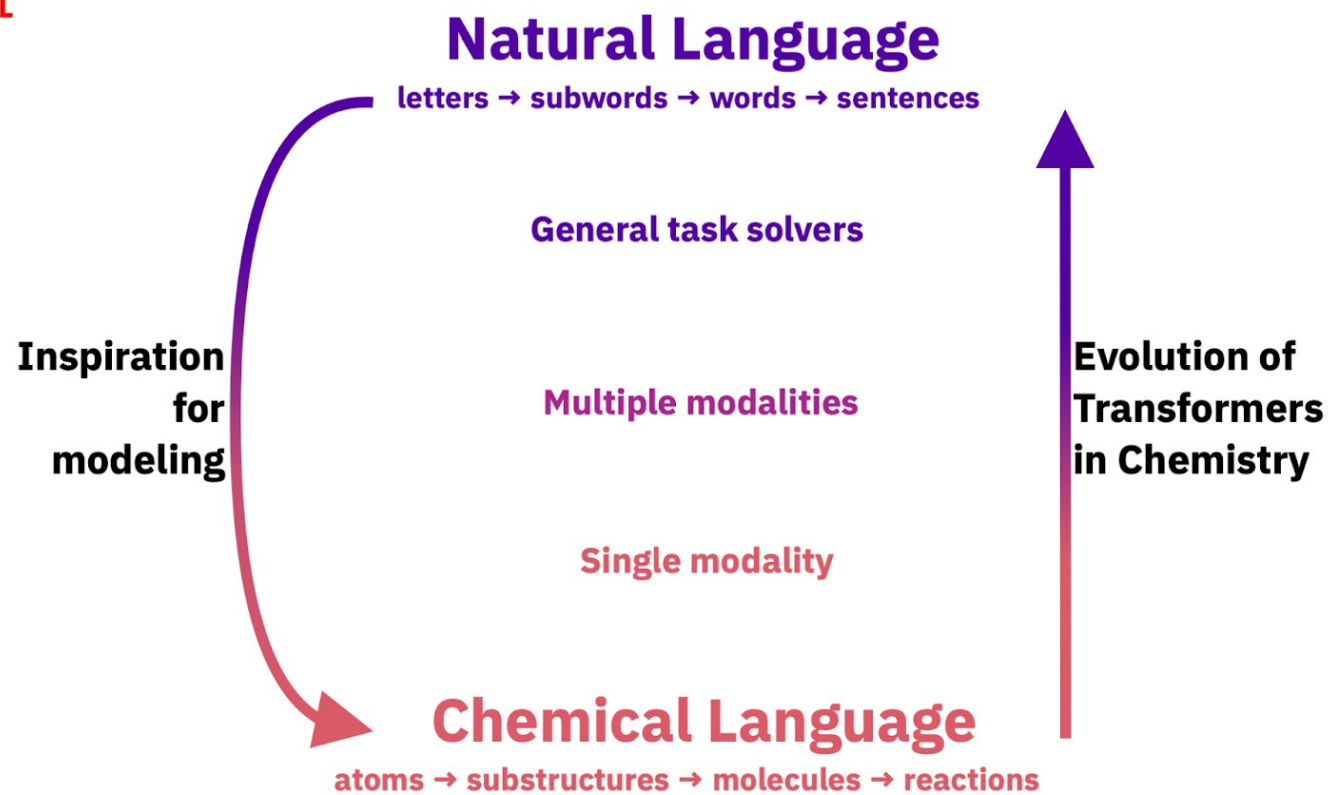
Andres M Bran

5.1 LLMs in Drug Discovery

Professor: Andres M Bran
Lecture Recording: Coming soon...
Slides: Notebook: Colab


1 06/18/2024

EPFL



Bran, A. and Schwaller, P. "Transformers and Large Language Models for Chemistry and Drug Discovery." *ArXiv abs/2310.06083* (2023)



5.4 Afaf



Afaf Taik

5.4 Ethical & Bias Concerns

Professor: Afaf Taik Lecture
Recording: Coming soon... Slides:

 1  06/18/2024

The End