

Machine Learning and Chemistry

QMMM Study Group

Dakota Folmsbee

2018/08/31

What is Machine Learning?

Ability to learn without rules-based programming

Learning:

- Supervised

- Unsupervised

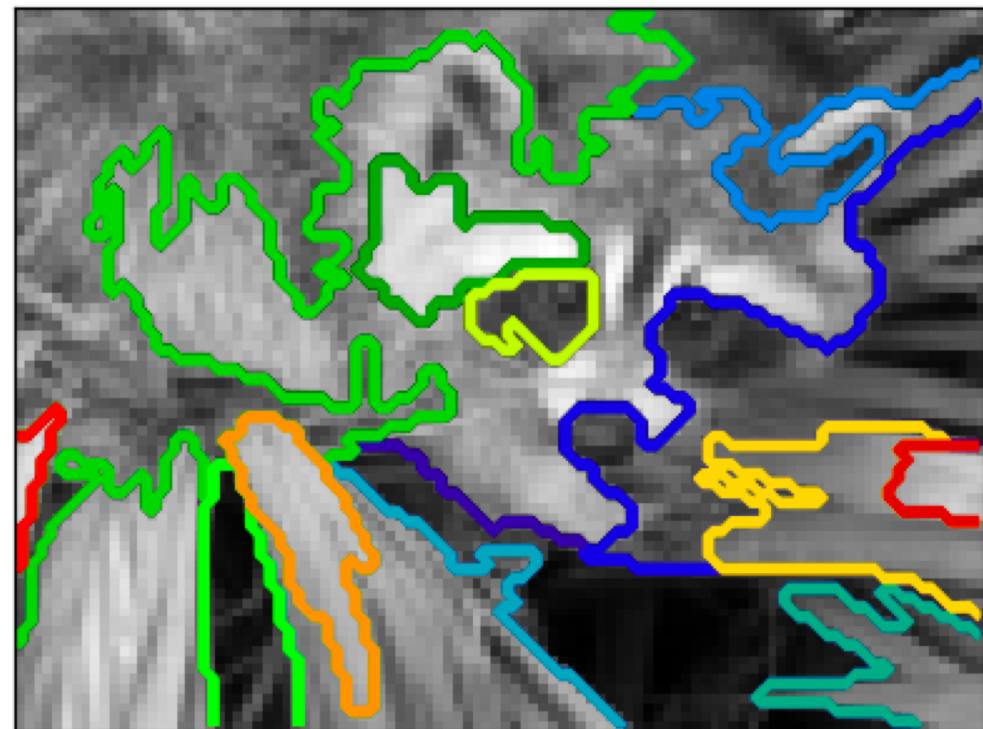
- Reinforcement

Types:

- Classification

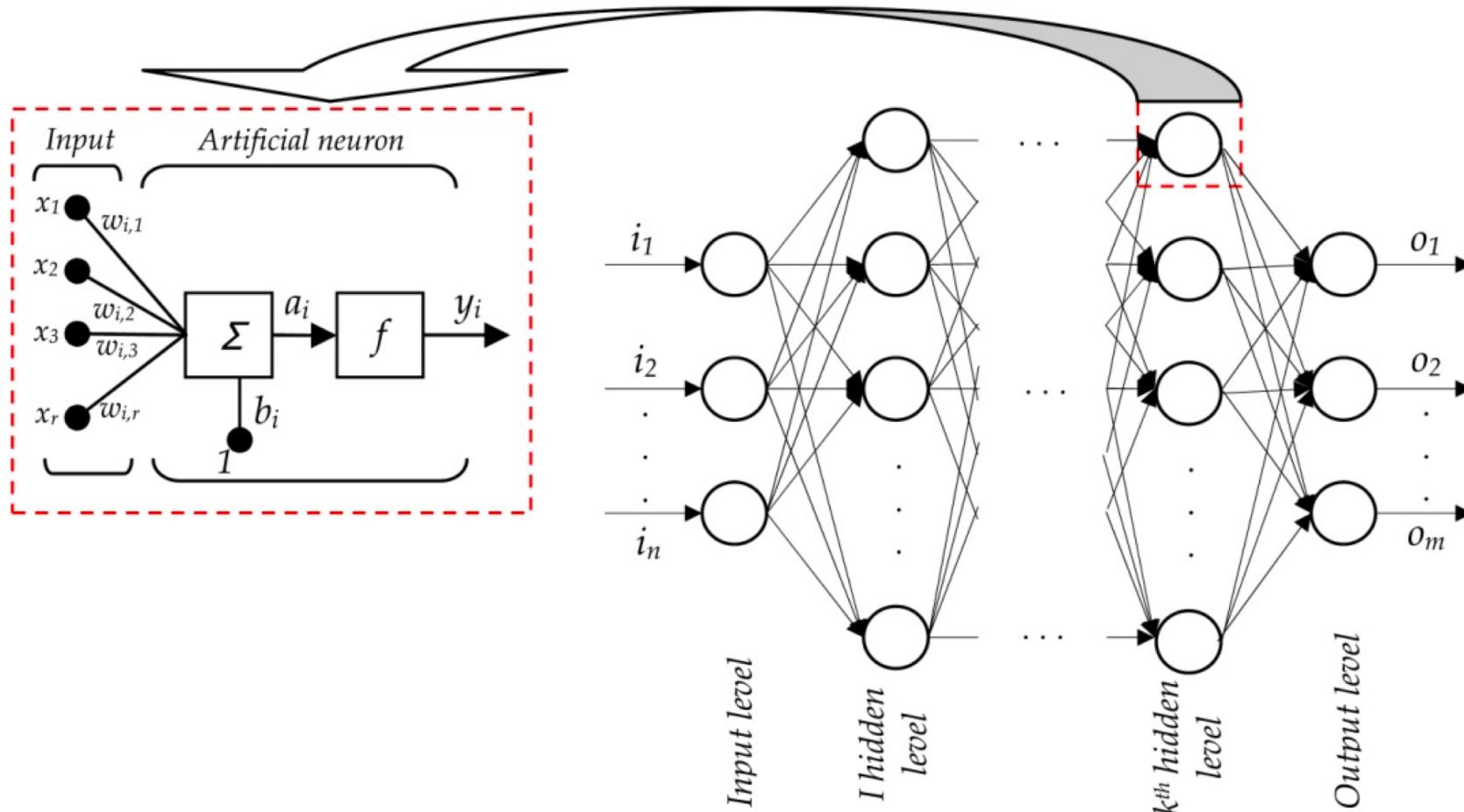
- Regression

- Clustering



http://scikit-learn.org/stable/auto_examples/cluster/plot_face_ward_segmentation.html

What is Machine Learning?



Why use ML in Chemistry?

Chemical space $\sim 10^{60}$ stable molecules

Time requirement

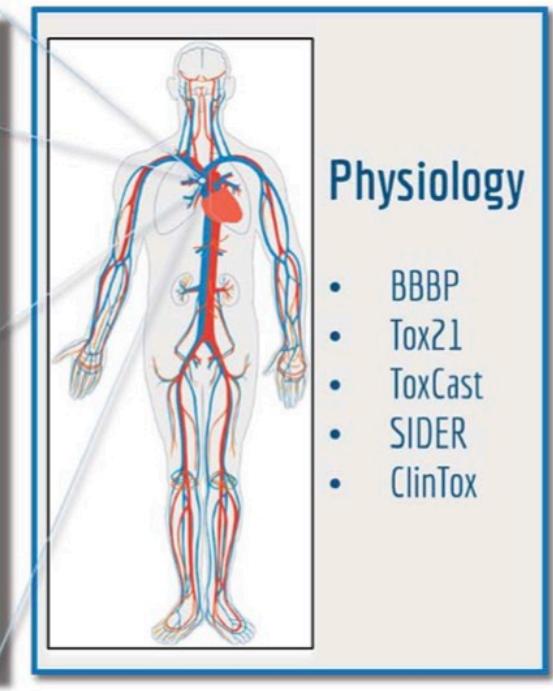
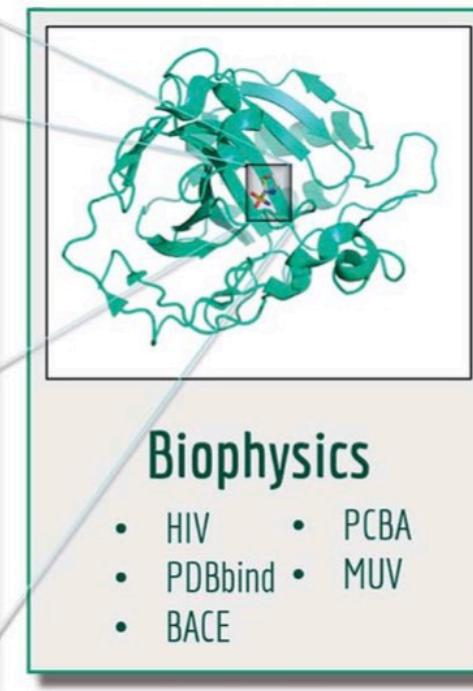
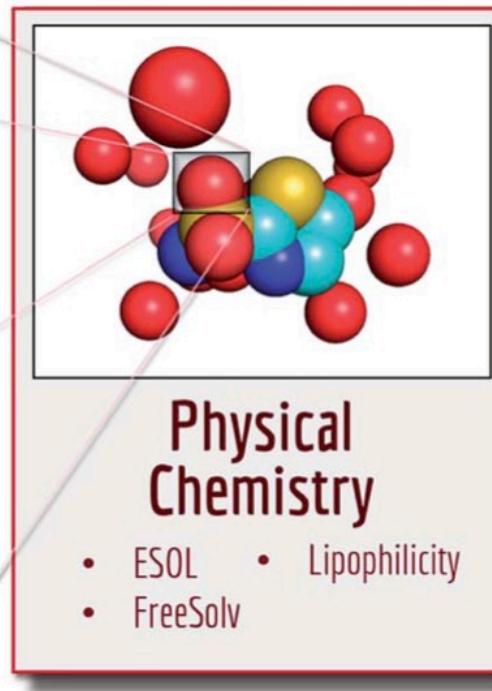
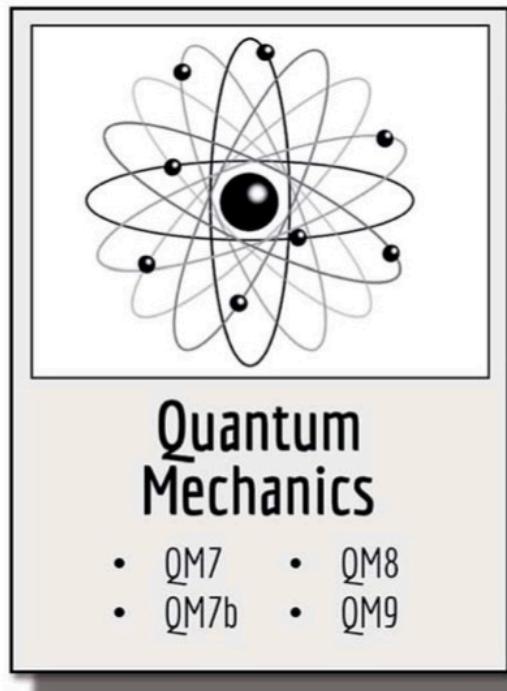
Multi-property training

Multi-level theory training

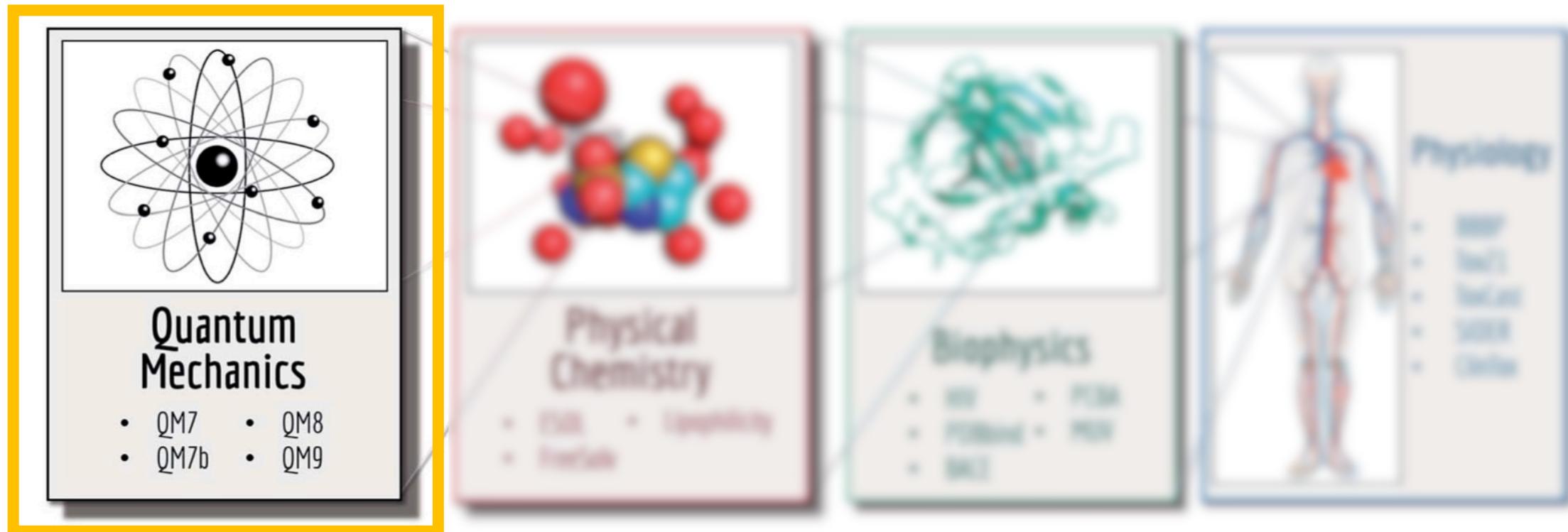
What is out there?

- Von Lilenfeld Group
 - Work with QM datasets
 - Bayesian Ridge, Kernel Ridge, and more
 - Coulomb Matrix, Bag of Bonds, BA Representations
- MoleculeNet and DeepChem
 - Benchmarks for machine learning models
 - Implementation of Scikit-Learn and Tensorflow
- And Much More!

Datasets



Datasets



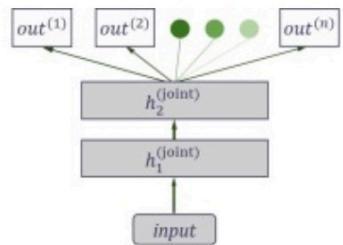
Current Approaches

- Machine Learning Methods
- Representations
- Benchmarks

Current Approaches

- Machine Learning Methods
 - Conventional
 - Graph
- Representations
- Benchmarks
- Packages

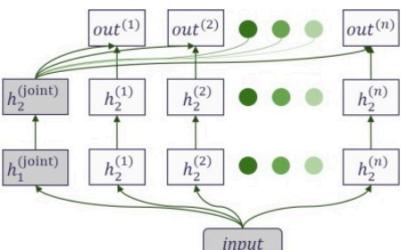
Conventional Models



Multitask Network

Standard neural network prediction method designed for multitask settings. Input features are processed through multiple shared fully-connected layers and then fed into separate linear classifiers/regressors for each different task. In case of single task dataset, it will become vanilla neural network model.

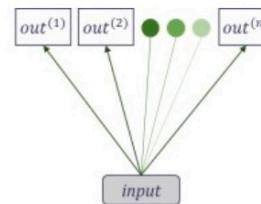
Regression
Classification
Deterministic



Bypass Network

A modified version of multitask network designed for uncorrelated tasks. Based on the structure of multitask network, it adds "bypass" layers directly connecting input features and each individual task, hence increasing explanatory power in case of unrelated variations in the sample.

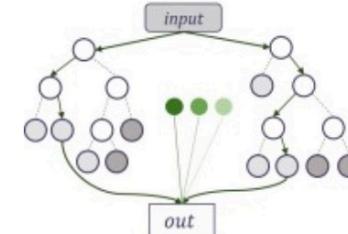
Regression
Classification
Deterministic



Logistic Regression

Standard classification model by applying logistic function to weighted linear combination of input features.

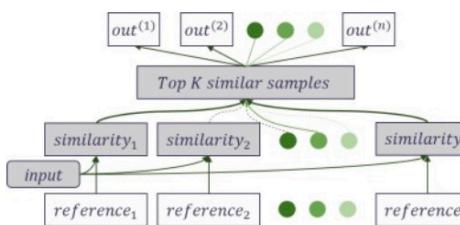
Classification
Deterministic



Random Forest

Standard classification and regression method based on an ensemble of decision trees, each trained on a different subsampled version of the original dataset.

Regression
Classification
Deterministic

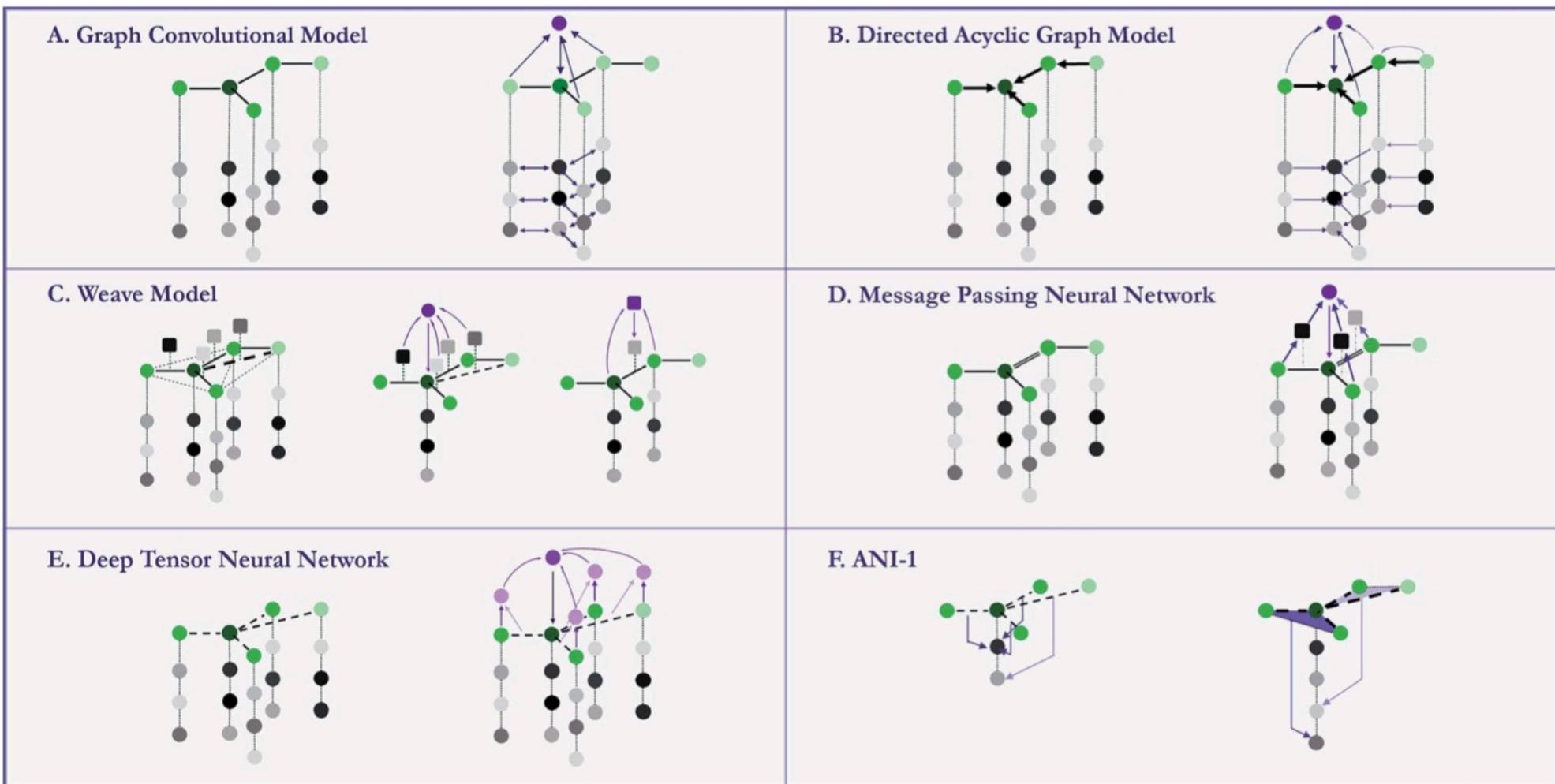


Influence Relevance Vector

Refined K-nearest neighbour classifier. Using the hypothesis that compounds with similar substructures have similar functionality, it makes prediction by combining labels from the top-K compounds most similar to the sample.

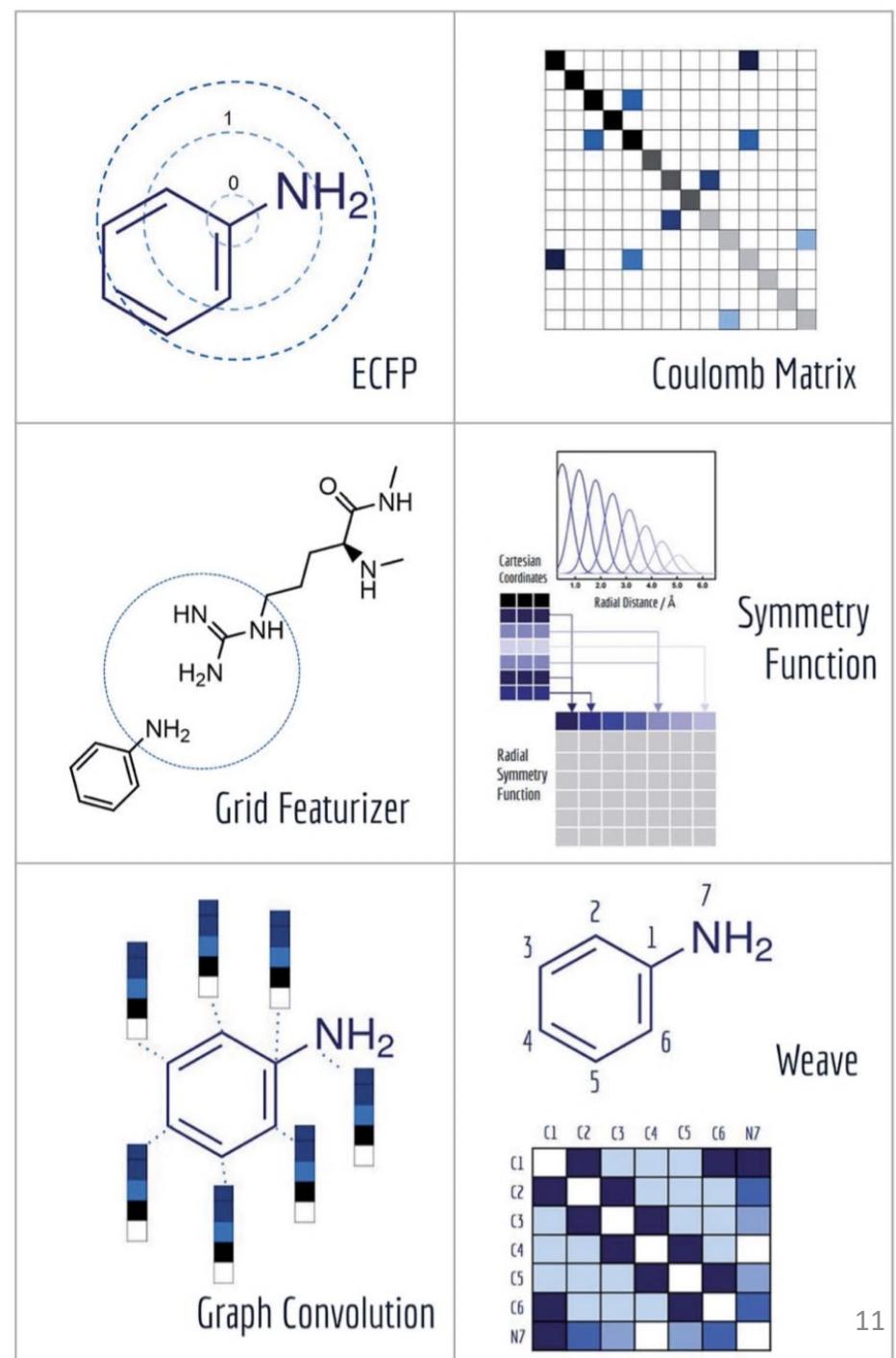
Classification
Deterministic

Graph Models



Current Approaches

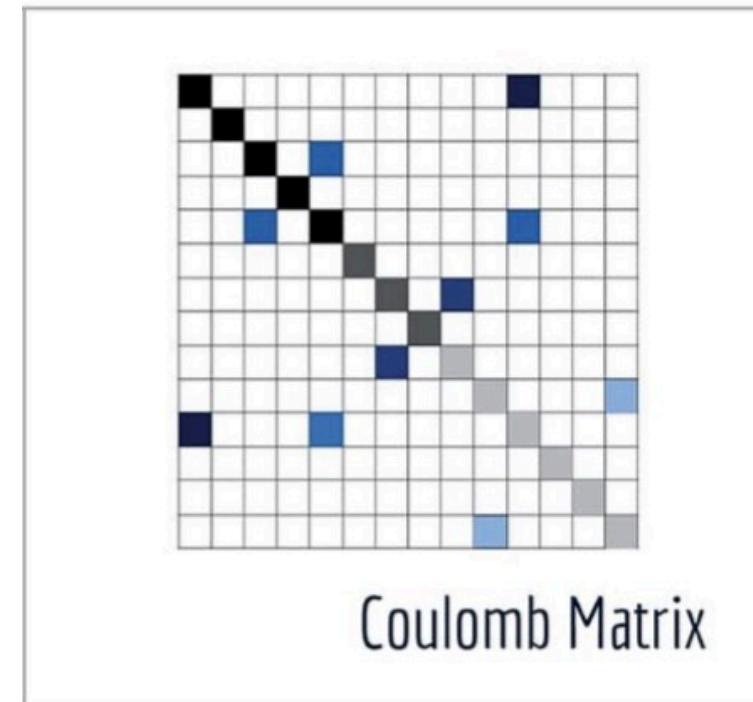
- Machine Learning Methods
- **Representations**
 - Coulomb Matrix
 - Bag of Bonds
 - BAML
 - HD/HDAD
 - Molecular Graphs
- Benchmarks
- Packages



Representations

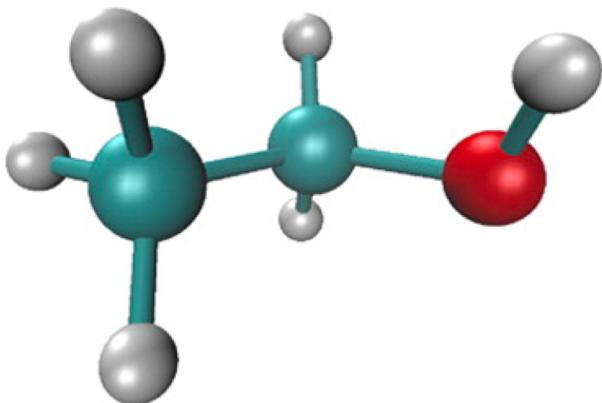
- Coulomb Matrix

$$M_{IJ} = \begin{cases} 0.5Z_I^{2.4} & \text{for } I = J \\ \frac{Z_I Z_J}{|\mathbf{R}_I - \mathbf{R}_J|} & \text{for } I \neq J \end{cases}$$



Representations

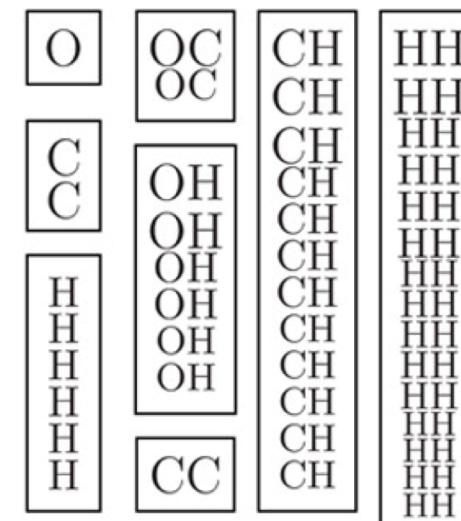
- Bag of Bonds



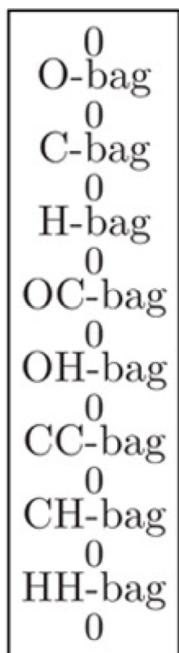
(a)

	O	C	C	H	H	H	H	H	H
O	o	oc	oc	oh	oh	oh	oh	oh	oh
C	oc	c	cc	ch	ch	ch	ch	ch	ch
C	oc	cc	c	ch	ch	ch	ch	ch	ch
H	oh	ch	ch	h	hh	hh	hh	hh	hh
H	oh	ch	ch	hh	h	hh	hh	hh	hh
H	oh	ch	ch	hh	hh	h	hh	hh	hh
H	oh	ch	ch	hh	hh	hh	h	hh	hh
H	oh	ch	ch	hh	hh	hh	hh	h	hh
H	oh	ch	ch	hh	hh	hh	hh	hh	h

(b)

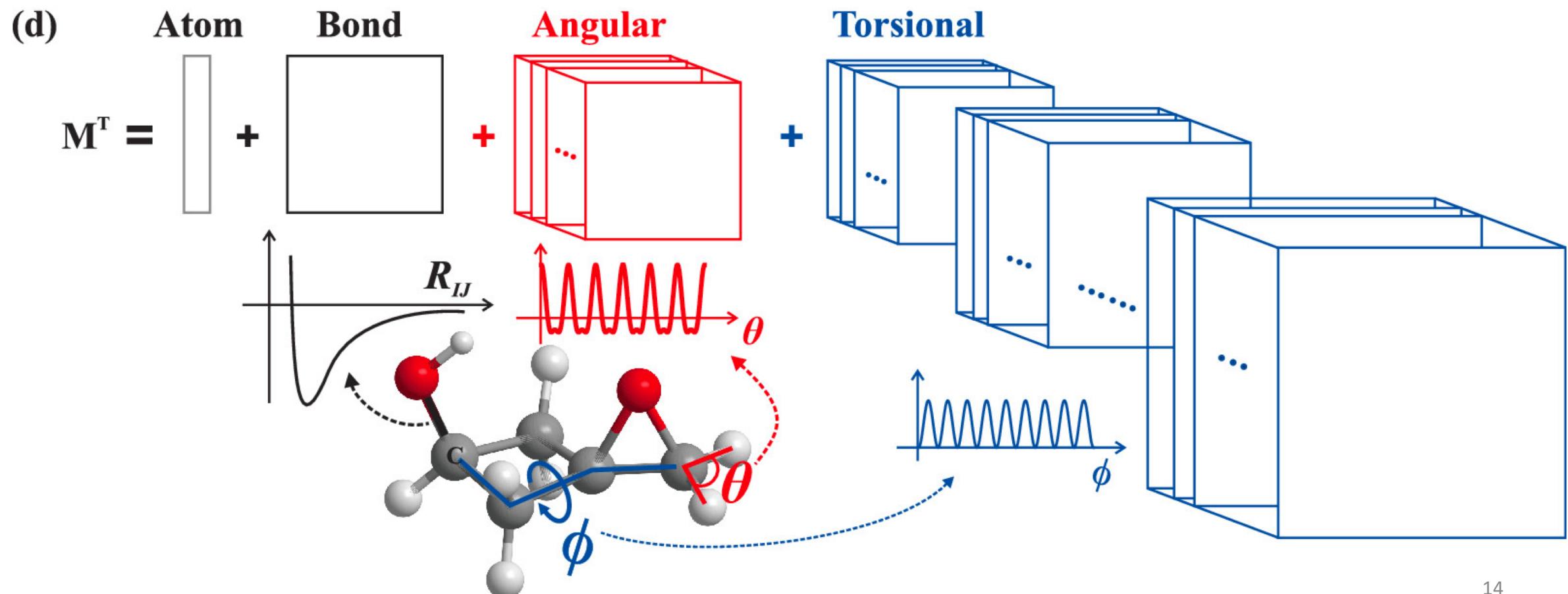


(c)



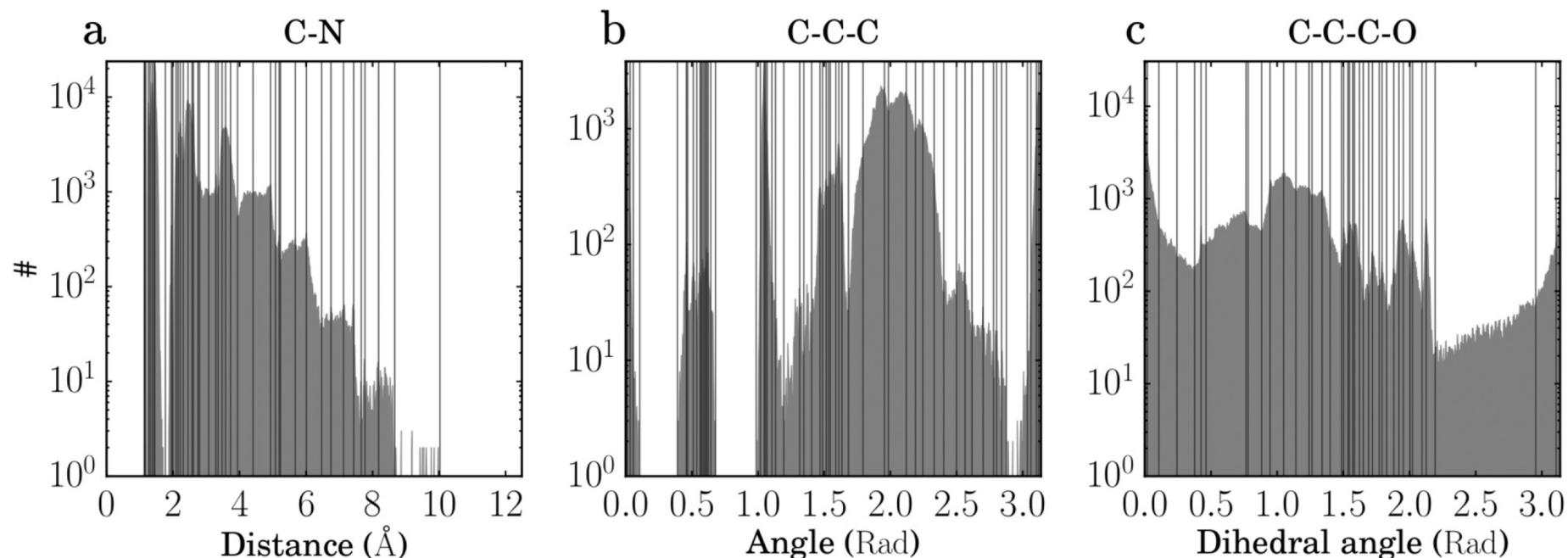
Representations

- BAML



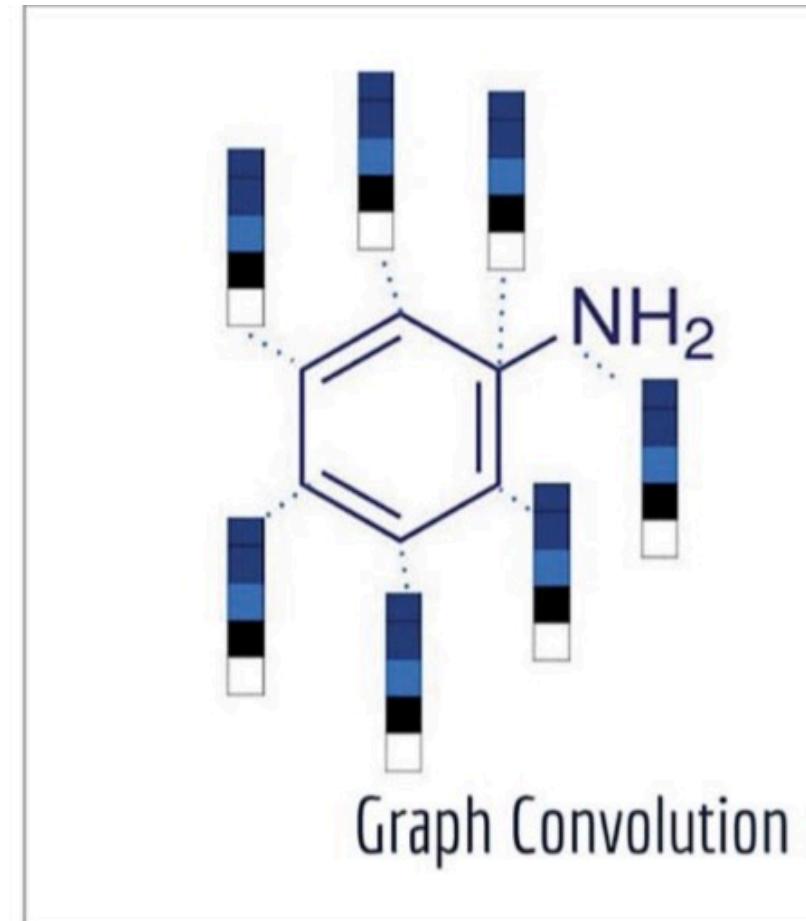
Representations

- Histogram of Distances (HD)
- Histogram of Distances, Angles, and Dihedrals (HDAD)



Representations

- Molecular Graphs
- Contains feature vectors:
 - Local chemical environment
 - Atom types
 - Hybridization types
 - Valence structures



Current Approaches

- Machine Learning Methods
- Representations
- **Benchmarks**
 - **MoleculeNet**
 - **von Lilienfeld group**
- Packages

Benchmarks

Method	Representation	U_0 (ev)	ϵ_{HOMO} (ev)	ϵ_{LUMO} (ev)	$\Delta\epsilon$ (ev)	μ (D)	α (bohr ³)	ZPVE (ev)	Cv (cal/molK)
KRR [†]	CM	0.128	0.133	0.183	0.229	0.449	0.433	0.0048	0.118
	BOB	0.0667	0.0948	0.122	0.148	0.423	0.298	0.00364	0.0917
	BAML	0.0519	0.0946	0.121	0.152	0.46	0.301	0.00331	0.082
	HDAD	0.0251	0.0662	0.0842	0.107	0.334	0.175	0.00191	0.0441
	HD	0.0644	0.0874	0.113	0.143	0.364	0.299	0.00316	0.0844
GG [†]	MG	0.0421	0.0567	0.0628	0.0877	0.247	0.161	0.00431	0.0837
GC [‡]	MG	0.15	0.0549	0.062	0.0869	0.101	0.232	0.00966	0.097
Multitask [‡]	CM	0.0984	0.00506	0.00645	0.0086	0.519	0.85	0.00207	0.39
GC [‡]	MG	0.1479	0.00716	0.00921	0.0112	0.583	1.37	0.00299	0.65
DTNN [‡]	MG	0.1054	0.00388	0.00513	0.0066	0.244	0.95	0.00172*	0.27
MPNN [‡]	MG	0.0889	0.00541	0.00623	0.0082	0.358	0.89	0.00216	0.42

Current Approaches

- Machine Learning Methods
- Representations
- Benchmarks
- **Packages**
 - **Methods**
 - **Representations**

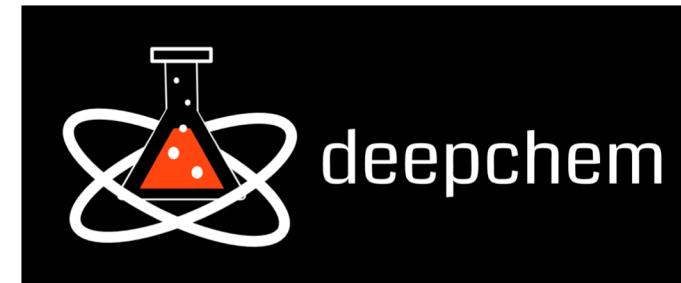
Machine Learning Libraries

- Scikit-Learn
 - <http://scikit-learn.org/stable/>
- Tensorflow
 - <https://www.tensorflow.org/>
- DeepChem
 - <https://deepchem.io/>
 - <https://github.com/deepchem/deepchem>



Representation Libraries

- DeepChem
 - <https://deepchem.io/>
 - <https://github.com/deepchem/deepchem>
- QML
 - <http://www.qmlcode.org/>
 - <https://github.com/qmlcode/tutorial>
- MolML
 - <https://github.com/crcollins/molml>



Find this interesting?

Want to recreate representations from the literature?

Check out chemreps and help Amanda, Shiv, and I do just that!

<https://github.com/dlf57/chemreps>

You are welcome to fork, submit issues, and add feature requests!

Useful Resources

- **MoleculeNet** <http://moleculenet.ai/>
- **MoleculeNet: a benchmark for molecular machine learning**
Wu, Zhenqin and Ramsundar, Bharath and Feinberg, Evan N. and Gomes, Joseph and Geniesse, Caleb and Pappu, Aneesh S. and Leswing, Karl and Pande, Vijay
Chem. Sci. **2018** 9 (2), 513-530
DOI: 10.1039/C7SC02664A
- **Prediction Errors of Molecular Machine Learning Models Lower than Hybrid DFT Error**
Felix A. Faber, Luke Hutchison, Bing Huang, Justin Gilmer, Samuel S. Schoenholz, George E. Dahl, Oriol Vinyals, Steven Kearnes, Patrick F. Riley, and O. Anatole von Lilienfeld
Journal of Chemical Theory and Computation **2017** 13 (11), 5255-5264
DOI: 10.1021/acs.jctc.7b00577