



Coarse graining autoencoders and evolutionary learning of atomistic potentials

IPAM, Los Angeles, September 24 2019 Rafael Gómez-Bombarelli rafagb@mit.edu http://gomezbombarelli.mit.edu/

Computational spectrum - virtuous cycle

There is essentially a continuum of higher parametrization and statistical learning connecting first principles (theory-based simulations) to black-box statistical learning over experiments.



High-throughput virtual screening



Successful applications

Organic Light Emitting Diodes

- High end displays, potentially lighting.
- Lightweight, flexible, transparent, high contrast, low power





Organic Flow battery electrolytes

- High-scale energy storage
- Emerging technology, promising low-cost





Charting chemical space



5

GCNN for interpolating

Cheap surrogate function to save useless computations

Atoms represented as nodes and bonds as distance-labeled edges

Node and edge features updated iteratively based on learned neighborhood mappings

Message, update, and embedding functions are neural networks



 r_1

GCNN and Neural Potentials

 Typically achieve state of the art performance over topological QSPR regression problems in chemistry

• Require 10⁴ or more to be truly effective

As interatomic potentials < 1 kcal/mol energies and 1 kcal/mol/A forces (as low as 0.1)

• Can overfit in chemical and configurational space



Smith, J. S. et. al.. *Sci. Data* **4**, 170193 (2017); Smith, J. S *Chem. Sci.* **8**, 3192–3203 (2017). Smith, J. S. et. al. *J. Chem. Phys.* **148**, 241733 (2018); Schüt et al. *J. Chem. Phys.* **148**, 241722 (2018); Hansen, K. et al. *J. Phys. Chem. Lett.* **6**, 2326–2331 (2015); Chmiela, S. et al. *Sci. Adv.* **3**, e1603015 (2017); Schütt, K. T. *Nat. Commun.* **8**, 13890 (2017); Chmiela, S. *Nat. Commun.* **9**, 3887 (2018).

Inverse design

Progress in predicting performance given candidate



Can we generate candidate based on design targets?



Inverse design



Semisupervised Molecular VAE



The latent representation now encodes mapping to one or more properties.



Chemical space is different

IN THE PIPELINE

Derek Lowe's commentary on drug discovery and the pharma industry. An editorially independent blog from the publishers of *Science Translational Medicine*. All content is Derek's own, and he does not in any way speak for his employer.

CHEMICAL NEWS

Calculating A Few Too Many New Compounds

By Derek Lowe 8 November, 2016

Is there a distribution we want to draw from?

In chemistry, one's ideas, however beautiful, logical, elegant, imaginative they may be in their own right, are simply without value unless they are actually applicable to the one physical environment. (Woodward)





Coarse-graining MD

TOWARDS INVERSE DESIGN IN 3D

Coarse Grained Methods

Coarse Graining MD simulates coarse grained variables that represents **slow collective atomistic motions** derived from full atomistic simulations

Coarse Graining methods find the "effective" coarse grained potential **given** a predetermined coarse graining mapping



Coarse Grained Methods

Extensively studied how to find the coarse graining potentials that reproduces *equilibrium* structural correlation function from atomistic simulations given a pre-determined CG mapping

Methods to approximate Coarse Grained Force Fields: Relative Entropy, Force Matching, g-YBG (implemented in VOTCA, BOCS, etc.)

Systematic Coarse Grained force fields for Biomolecules: MARTINI







[2]



[3]

Learning to Coarse-Grain





A learning problem

15

Coarse-graining framework



 \mathbf{x} : atomistic coordinates $\mathbf{V}(\mathbf{x})$: All-Atom Potential \boldsymbol{z} : coarse grainedcoordinates $V_{CG}(\boldsymbol{z})$: coarse grainedPotential

```
Instead of

given x, V(x), E(x) \rightarrow \text{find } V_{CG}(z)

We propose

given x, V(x) \rightarrow \text{find } E(x) \text{ and } V_{CG}(z = E(x))
```

 $V_{CG}(z)$ can have an arbitrary functional form

Classical

$$V_{CG}(z) = V_{bonded}(z) + V_{non-bonded}(z)$$

• Neural (using MPNN)

Force matching



The encoding function



We propose a neural network like encoding with the following constraints

- 1) $z_I = E_I(x) = \sum_j E_{Ij} x_j$ [1] 2) $\sum_j E_{Ij} = 1$ and $E_{Ij} \ge 0$ 3) $M_I = (\nabla E^{-1})^T M_j \nabla E^{-1}$
- 4) We use Gumbel-softmax during training to enforce the learning of **discrete** coarse graining variables to ensure that each atom only contributes to one CG atom

"Coarse Graining" Forces



- We also need a function that variationally determines *F*
- $F = \langle -b \cdot \nabla V(x) \rangle_{E(x)=z}$ where *b* is the force coarse graining function
- A consistent choice for *b* from statistical mechanics: $b = \frac{\nabla E(x)^{T}}{\nabla E(x)^{T} \cdot \nabla E(x)}$
- Computing the mean force *F* requires constrained dynamics.
- However, we want a one-shot optimization stack without running extra MD simulations.

Stochastic Force Matching



• We compute the instantaneous stochastic "coarse grained" force *F*_{ins} [1]

•
$$F_{ins} = -b \cdot \nabla V(x)$$
 (instantaneous force)

•
$$F = \langle F_{ins} \rangle_{E(x)=z}$$
 (mean force)

$$L = |F + \nabla_z V_{CG}(E(x))|^2$$



 $L_{ins} = |F_{ins} + \nabla_z V_{CG}(E(x))|^2$



• Supervised force matching conditioned on E(x) learns coarse graining mapping field and the potential in CG.



- AutoEncoder automatically coarse-grains atomistic coordinates to CG coordinates in a data-driven way
- Force matching also helps to shape the learning of CG and obtain $V_{CG}(z = E(x))$ for CG simulations



- AutoEncoder automatically coarse-grains atomistic coordinates to CG coordinates in a data-driven way
- Force matching also helps to shape the learning of CG and obtain $V_{CG}(z = E(x))$ for CG simulations

Pipeline

Run short all-atom MD trajectory on small system.

Shuffle frames

Decide degree of compression, functional form of V_{CG}

Train E and V_{CG} simultaneously

Deploy CG trajectory on larger system

Compare statistics

• Encoded AA vs. CG MD sampling



Alanine Dipeptide fully unsupervised

Automatic CG for small molecules

Compression vs. fluctuations

A pre-training step that is unsupervised (i.e. no V_{CG} is learnt) followed by adding supervision

The difference between *F*_{ins} and *F* is related to the friction, memory and fluctuations in the GLE.

We choose to add additional regularization to minimize $(F_{ins})^2$

Pipeline

Split training

Α

- into pre-train with force regularization and reconstruction
- $\,\circ\,$ fully supervised training of E and V_{CG}

ρ	0.0
Force regularization	4498.0
Reconstruction	0.36

ρ	0.01
Force regularization	1008.0
Reconstruction	0.36

CG of liquid ethane

(a) (b) bond probability density cg1 cg1 cg2 Ĥ. - Ĥ 1.90 2.15 2.20 1.95 2.10 2.00 2.05 r (Â) 0.20 0.05 0.10 0.15 0.25 (c) cg1-cg1 (d) ethane-ethane 1.75 — CG 1.50 --- atomistic 2.0 1.25 1.5 1.00 (L) **b** g (r) 0.75 0.50 0.5 0.25 0.00 0.0 6 10 12 14 0 10 12 r (Å) r (Å)

Classical potential

MPNN neural potential

CG of liquid ethane

MPNN neural potential

CG of polyethylene

120 atoms of polyethylene

C24 in bulk

1:2 and 1:3 compression

Structure correlations are recovered correctly.

Diffusion kinetics are modified (fortuitous agreements)

Evolutionary Neural Potentials

Lithium chelation

Lithium ion batteries require electrolytes to shuffle Li cations between electrodes.

Organic electrolytes

- have a great design space
- Liquids have high conductivity can catch fire.

• Polymer solids are safer but lower conductivity. Polyethylene oxides

Explore the design space of lithium-binding moieties

Chemistries and configurations

We are interested in Li-organic binding enthalpies (and free energies too, ...)

Global energy minimization requires good sampling of supramolecular clusters.

• A sampling problem!

• Can access with MD with some degree.

First-principles is too costly, and the chemistry is fairly constrained: create a potential.

Evolutionary NNs

- Initialize with normal—mode data only and using neural network to sample new configurations via MD.
- Obtain DFT energies on neural network sampled data and validate on accuracy
- Train on <= 30000 points BP86/SVP-D3. (E and forces) across 14 chemistries with <= 20 heavy atoms, validated on 40 larger ones
- Decouple sampling for data acquisition.
- **Evolve** by incorporating neural network sampled configurations into the next generation of training data

Evolutionary NNs

Validation

- Validated on seen and unseen chemicals, of up to 2x the size.
- It achieves good accuracy when validated on converged geometries on species that are not trained by the model -> good transferability over ether chemical space.

Error: 3.0 kcal/mol

Trained species prediction MAE: 0.1 kcal/mol
 New species prediction MAE: 3.1 kcal/mol

Validating Hessians

Exploring ether space

Predicting binding affinity:

- Run neural NVE MD •
- Gradient descent at some frames •

<2ms second per

O(N^2) and

Thanks!

Wujie Wang (CG, Evolving NN) Wil Harris (Graph NN)

Daniel Schwalbe Koda Somesh Mohapatra James Damewood Shi Jun Ang

