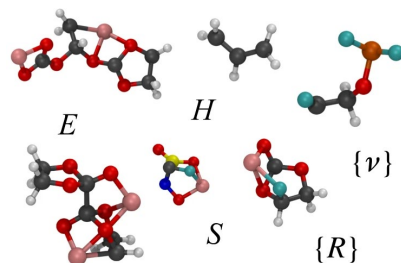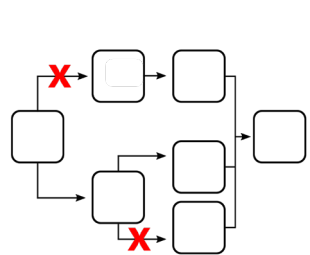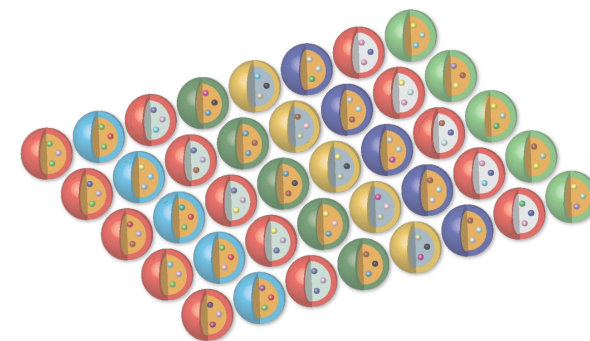# High-Throughput DFT and Monte Carlo for Reaction Networks and Machine Learning

Samuel M. Blau

Research Scientist

Lawrence Berkeley Lab

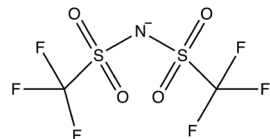# High-Throughput Molecular DFT Data Generation



Charged molecules

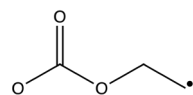Radical molecules

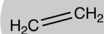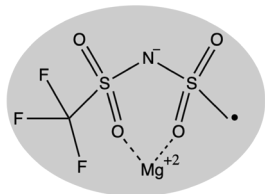Metal-coordinated molecules

Solvated molecules

Additional complexity

BERKELEY LAB

# High-Throughput Molecular DFT Data Generation



Typically, **<75% success**

With on-the-fly error correction:

**>98% success**

**S. M. Blau***, E. Spotte-Smith*, B. Wood, S. Dwaraknath, K. A. Persson, *ChemRxiv* 2020

# High-Throughput Molecular DFT Workflow Infrastructure



S. M. Blau*, E. Spotte-Smith*, B. Wood, S. Dwaraknath, K. A. Persson, *ChemRxiv* 2020

# We Use Workflows to Generate Unique Simulated Datasets



**LIBE**

**Lithium-Ion Battery Electrolyte**
**17,190 molecules**

E. W. C. Spotte-Smith*, **S. M. Blau***, et al., *Sci. Data* 2021

ωB97X-V/def2-TZVPPD/SMD

**MADEIRA**

**MA**gnesium **D**ataset of **E**lectrolyte and **I**nterphase **ReA**gents:
**11,502** molecules

E. W. C. Spotte-Smith, **S. M. Blau**, et al., *JACS (accepted)*

ωB97X-V/def2-TZVPPD/SMD

**RAPTER**

**ReA**ctants, **P**roducts, and **T**ransition-states of **E**lementary **R**eactions:
**>15,000** complex reactions

E. W. C. Spotte-Smith, **S. M. Blau**, et al., *In preparation*

ωB97X-D/def2-SVPD/PCM

Collaborators:

Evan Spotte-Smith    Kristin Persson

# We Use Workflows to Generate Unique Simulated Datasets



**ESCoMMS**

**E**lectronic **S**tructure of **Co**mplexes with **M**etals of **M**any **S**pins:
**>140,000** complexes

ωB97M-V/def2-SVPD

Michael Taylor

Ping Yang

**ORIONS**

**OR**bital **I**nteractions of **O**rga**N**ic **S**pecies:
**>230,000** molecules

D. Boiko et al., *ChemRxiv* 2022

Gabe Gomes

**SUNSET**

**S**imulated **U**pconverting **N**anoparticle **S**pectra for **E**missions **T**uning:
**>6,000** spectra (kMC, not DFT)

Eric Sivonxay       Emory Chan

BERKELEY LAB

4

# Machine Learning Atop Our DFT Datasets

BonDNet GNN



M. Wen, **S. M. Blau**, E. Spotte-Smith, S. Dwaraknath, K. A. Persson, *Chem. Sci.* 2021

**An orbital-based representation for accurate quantum machine learning**

Konstantin Karandashev[1,a] (iD) and O. Anatole von Lilienfeld[1,2,b] (iD)

"The LIBE dataset is of particular interest… [because] it contains species of different charge and spin states, enabling us to test [our model]'s ability to process them…"



D. Boiko*, T. Reschützegger*, B. Sanchez-Lengeling, **S. M. Blau**, G. d. P. Gomes, *In preparation*

# Introduction to Chemical Reaction Networks (CRNs)

△ = Unstable intermediate   ★ = Stable product

Initial
species

Chemical reaction network (CRN)

BERKELEY LAB

# Introduction to Chemical Reaction Networks (CRNs)



△ = Unstable intermediate    ★ = Stable product

Initial species

Chemical reaction network (CRN)

BERKELEY LAB

# Introduction to Chemical Reaction Networks (CRNs)



△ = Unstable intermediate    ★ = Stable product

Initial species

Chemical reaction network (CRN)

Important species

Reaction pathways

Free energy

Time dynamics

Concentration

Time

M. Wen, E. W. C. Spotte-Smith, **S. M. Blau**, M. J. McDermott, A. S. Krishnapriyan, K. A. Persson, *Nat. Comp. Sci.* 2023

# Background: Solid-Electrolyte-Interphase Formation



Negative electrode
("Anode")

Electrolyte

Positive electrode
("Cathode")

Salt

Solvent

e⁻

BERKELEY LAB

# Background: Solid-Electrolyte-Interphase Formation



**Goal:** enable next-generation batteries by **controlling SEI formation**

Big questions:

1. What species form?
   - Identify products
2. How do those species form?
   - Reaction mechanisms
3. How do individual species, pathways compete and interact?

AIMD, by-hand DFT investigations: limited insight

BERKELEY LAB

# A Data-Driven Approach to Understanding Reactivity



- Rational enumeration of possible species, reactions

- ΔG of each reaction in isolation via HT molecular DFT

- Network analysis: novel mechanistic insight

- Workflows necessary for data generation

# High-Throughput Molecular DFT Data Generation



pymatgen

Custodian

FireWorks

atomate

Q-CHEM
A QUANTUM LEAP INTO THE FUTURE OF CHEMISTRY

Principal Molecules

Solvent molecules          SEI products

Salts

Molecular Fragments

HIGH-THROUGHPUT DFT

LIBE

**Lithium-Ion Battery Electrolyte**
**17,190 molecules**

$E$     $H$

$\{\nu\}$

$S$     $\{R\}$

Selective Recombinant
Molecules

E. W. C. Spotte-Smith*, **S. M. Blau***, X. Xie, H. D. Patel, M. Wen,
B. Wood, S. Dwaraknath, K. A. Persson, *Sci. Data* 2021

BERKELEY LAB

9

# The Challenge of Reaction Generation

- Given e.g. 10k species, how to enumerate connecting reactions?

- Common approach – templates:



- Prescriptive templates are not well-suited to electron-driven chemistry

**Our solution:**

**filters**

Goals:

- Minimize prescriptive constraints in order to *facilitate discovery*

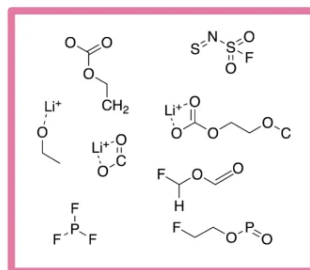- Want all reactions that:
  - Are likely to be single-step
  - May be kinetically viable

- Enable automated kinetic refinement

- Resolve complex competition

# High-Performance Reaction Generation: HiPRGen

$\mathbf{S}_{init}$

**Input:** initial species

LIBE-CHOLi = 8904 species

$\mathbf{S}_{filtered}$ $\mathbf{R}_{filtered}$

**Output:** species, reactions that compose network
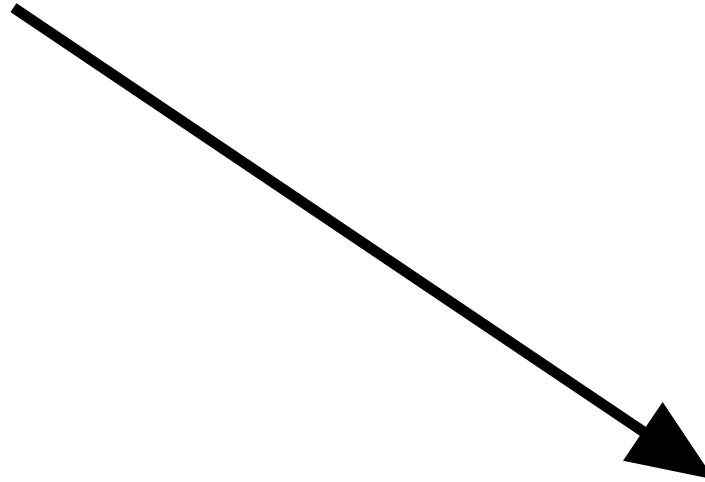
D. Barter*, E. W. C. Spotte-Smith*, N. Redkar, S. Dwaraknath, K. A. Persson, **S. M. Blau**, *Dig. Disc.* 2023

BERKELEY LAB

# High-Performance Reaction Generation: HiPRGen

**1.** Filter species

$\mathbf{S}_{init}$

**Input:** initial species

LIBE-CHOLi = 8904 species



H⁺ → $H^+$ KEEP

KEEP

DISCARD

KEEP

KEEP

KEEP

DISCARD

. . .

After filtering = 5193 species

- Metal-centric complexes
- $Li^0$-containing species

D. Barter*, E. W. C. Spotte-Smith*, N. Redkar, S. Dwaraknath, K. A. Persson, **S. M. Blau**, *Dig. Disc.* 2023

BERKELEY LAB

# High-Performance Reaction Generation: HiPRGen

**1. Filter species**

$S_{init}$

**Input:** initial species

LIBE-CHOLi = 8904 species

H[+]
KEEP

KEEP

DISCARD

KEEP

KEEP

KEEP

DISCARD

· · ·

- Metal-centric complexes
- $Li^0$-containing species
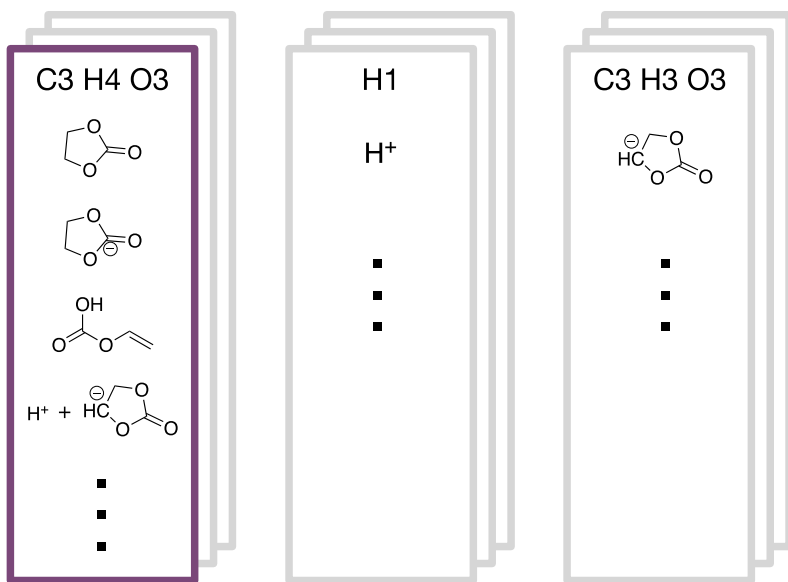
After filtering = 5193 species

**2. Bucket Species by Composition**

C3 H4 O3

H1

H[+]

C3 H3 O3

H[+] + 



D. Barter*, E. W. C. Spotte-Smith*, N. Redkar, S. Dwaraknath, K. A. Persson, **S. M. Blau**, *Dig. Disc.* 2023

# High-Performance Reaction Generation: HiPRGen

$S_{init}$

**Input:** initial species

LIBE-CHOLi = 8904 species

**1.** Filter species
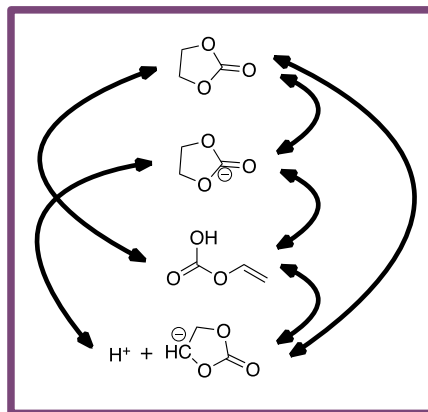
| $H^+$ | | | |
| KEEP | KEEP | DISCARD | KEEP |
| KEEP | KEEP | DISCARD | . . . |

After filtering = 5193 species

**2.** Bucket Species by Composition

| C3 H4 O3 | H1 | C3 H3 O3 |
| | $H^+$ | |

**3.** Generate reactions by stoichiometry

> 176 billion rxns

D. Barter*, E. W. C. Spotte-Smith*, N. Redkar, S. Dwaraknath, K. A. Persson, **S. M. Blau**, *Dig. Disc.* 2023

BERKELEY LAB

# High-Performance Reaction Generation: HiPRGen
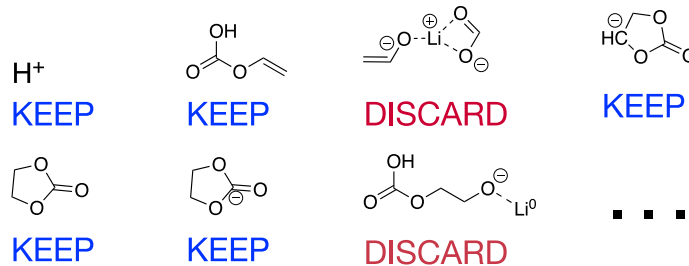
**1. Filter species**

H⁺ — KEEP
KEEP
DISCARD
KEEP
KEEP
KEEP
DISCARD

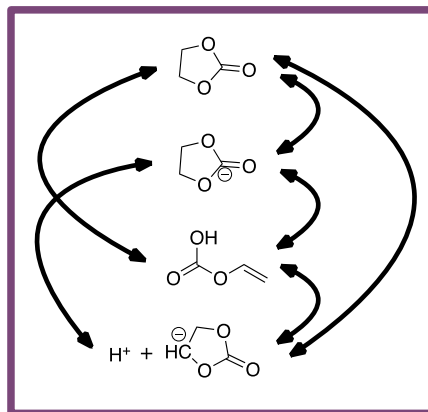$S_{init}$

**Input:** initial species

LIBE-CHOLi = 8904 species

After filtering = 5193 species

**2. Bucket Species by Composition**

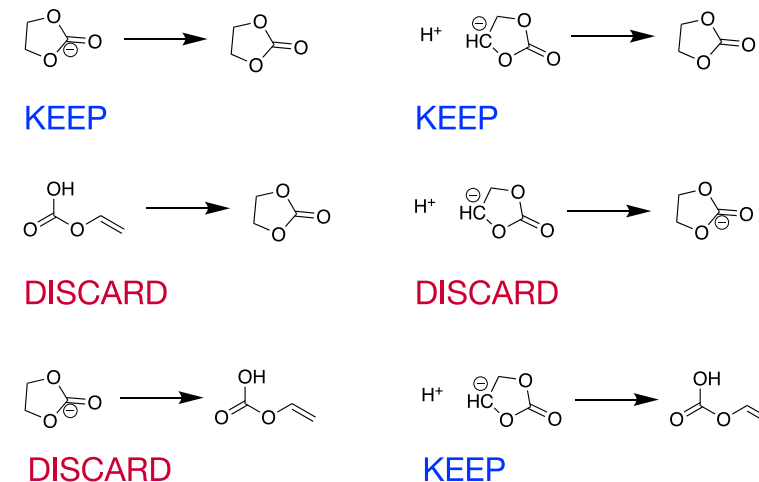C3 H4 O3    H1    C3 H3 O3

H⁺

**3. Generate reactions by stoichiometry**

> 176 billion rxns

**4. Filter reactions**

KEEP        KEEP

DISCARD     DISCARD

DISCARD     KEEP

- Too many bonds changing
- Bond change + redox
- Coordination + covalent bond change
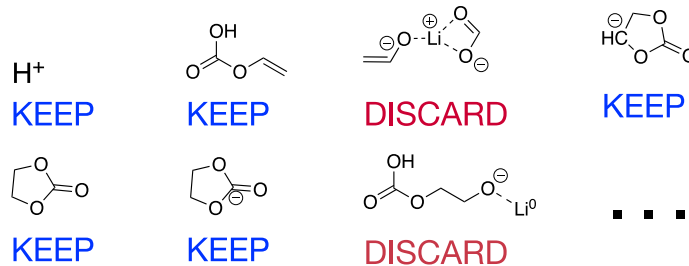
# High-Performance Reaction Generation: HiPRGen



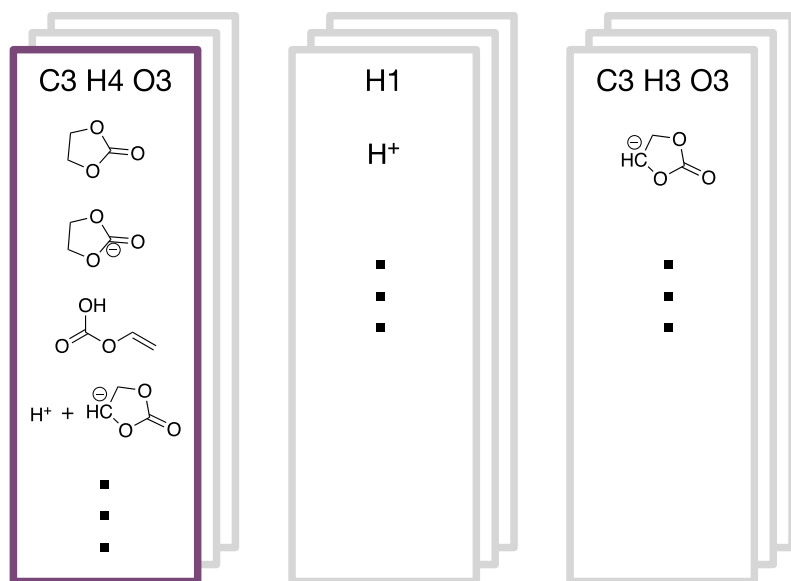**1. Filter species**

**4. Filter reactions**

$\mathbf{S}_{init}$

**Input:** initial species

LIBE-CHOLi = 8904 species

**2. Bucket Species by Composition**

**3. Generate reactions by stoichiometry**

5193 species

86 million rxns

$\mathbf{S}_{filtered}$ $\mathbf{R}_{filtered}$

**Output:** species, reactions that compose network

D. Barter*, E. W. C. Spotte-Smith*, N. Redkar, S. Dwaraknath, K. A. Persson, **S. M. Blau**, *Dig. Disc.* 2023

BERKELEY LAB

# Reaction Network Analysis: Graphs vs Kinetic Monte Carlo

## Molecule node

From reactions with A as a product → Specie A → To reactions with A as a reactant

## Reaction node

From reactant node(s) → Reactant(s) $\Phi$ Product(s) → To product node(s)

$\Phi$ : Cost

**S. M. Blau**, H. D. Patel, E. Spotte-Smith, X. Xie, S. Dwaraknath, K. A. Persson, *Chem. Sci.* 2021

- No concept of system state / concentrations
- Pathfinding to a given species scales as $O(N^2)$
- **Must know target of interest a priori**



- kMC... don't we need kinetics?
- No, all $\Delta G < 0$ rxns, all same rate
- "Thermodynamically bounded"

- Need initial state, evolve full system stepwise
- Stochastic sampling scales as $O(\log N)$ + parallelizable
- **Target prediction from full system exploration...?**

BERKELEY LAB

# Reaction Network Monte Carlo:  RNMC

Inputs:

$$\mathbf{S}_{filtered}$$

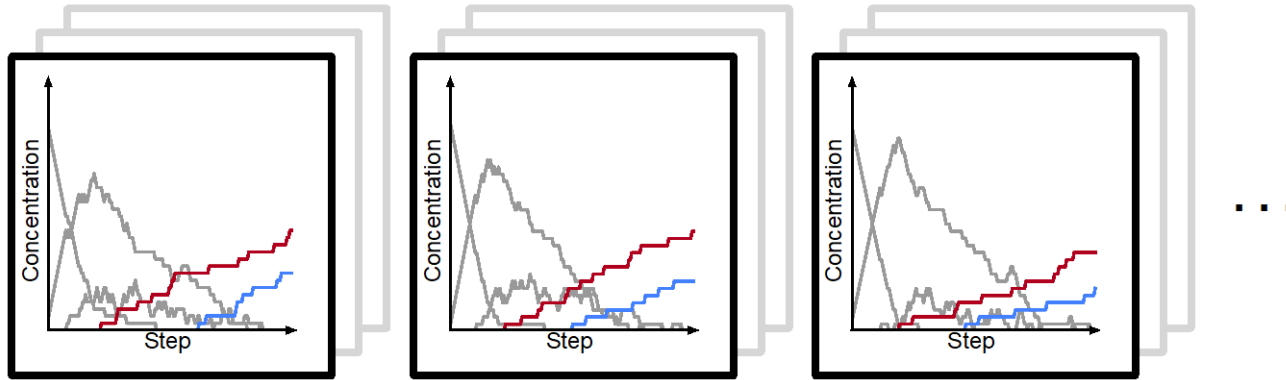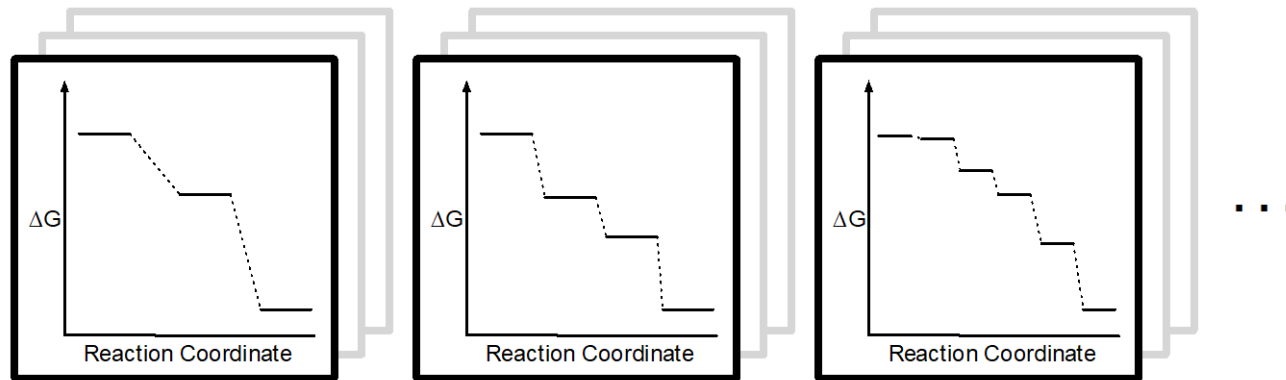$$\mathbf{R}_{filtered}$$

$$[x_i, x_j, ...]_o$$

Perform many thermodynamically bounded Monte Carlo trajectories



. . .

- 30 of each $x_i$
- All ΔG < 0: can run to completion
- 100k trajectories

D. Barter*, E. W. C. Spotte-Smith*, N. Redkar, S. Dwaraknath, K. A. Persson, **S. M. Blau**, *Dig. Disc.* 2023

BERKELEY LAB

# Reaction Network Monte Carlo: RNMC

Inputs:

$\mathbf{S}_{filtered}$

$\mathbf{R}_{filtered}$

$[x_i, x_j, \ldots]_o$

Perform many thermodynamically bounded Monte Carlo trajectories



- 30 of each $x_i$
- All ΔG < 0: can run to completion
- 100k trajectories

Extract shortest reaction pathways from each trajectory to each specie of interest



Can do pathfinding on up to approx. 300 million reactions

Group and order by cost
Φ ≈ # of reactions

github.com/BlauGroup/HiPRGen
github.com/BlauGroup/RNMC

D. Barter*, E. W. C. Spotte-Smith*, N. Redkar, S. Dwaraknath, K. A. Persson, **S. M. Blau**, *Dig. Disc.* 2023

# Converting RNMC and Identifying Network Products



Can the average trajectory identify network products?

D. Barter*, E. W. C. Spotte-Smith*, N. Redkar, S. Dwaraknath, K. A. Persson, **S. M. Blau**, *Dig. Disc.* 2023

# Converting RNMC and Identifying Network Products



- Totally heuristic
- **Network** products, not real products

High formation / consumption

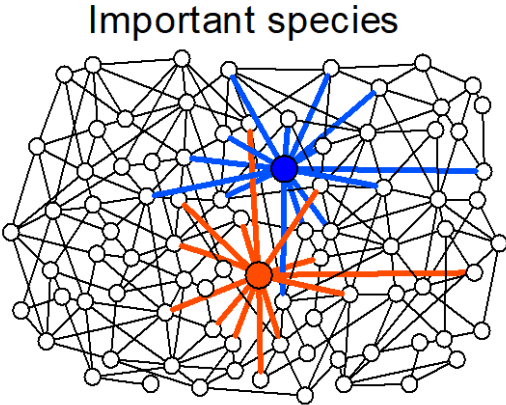Significant accumulation

Low-cost pathways available

**Network product?**

Can the average trajectory identify network products?

D. Barter*, E. W. C. Spotte-Smith*, N. Redkar, S. Dwaraknath, K. A. Persson, **S. M. Blau**, *Dig. Disc.* 2023

BERKELEY LAB

# Building Up and Picking Apart Complexity



**Species**

8,904

**I**: Filter species

**Reactions**

~176,000,000,000

**II**: Filter reactions

5,193

~86,000,000

**III**: Run MC trajectories*

3,182    ~312,000

**IV**: Find products in average trajectory

18

* 100,000 trajectories with Li⁺, EC, CO₂ at 0V vs. Li/Li⁺

Important species

Reaction pathways

Free energy

Time dynamics

Concentration

Time

D. Barter*, E. W. C. Spotte-Smith*, N. Redkar, S. Dwaraknath, K. A. Persson, **S. M. Blau**, *Dig. Disc.* 2023

BERKELEY LAB

# Predicted Battery Network Products: 36 out of 5139

**Small molecules/gases**



**Inorganics**



lithium carbonate     lithium oxalate

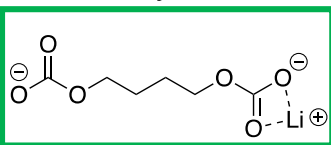[EC, Li$^+$] and [EC, Li$^+$, CO$_2$] at 0V and +0.5V vs. Li/Li$^+$

**Alkyl carbonates**



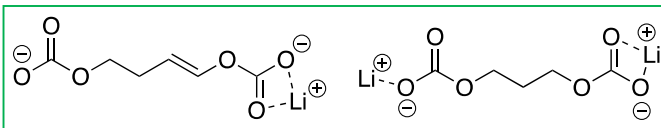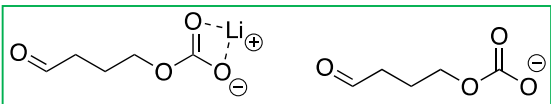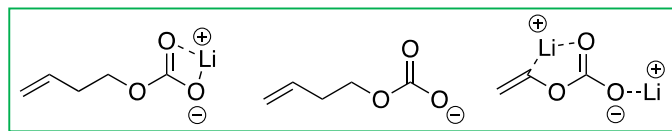LMC          vinyl carbonate          ethylene monocarbonate
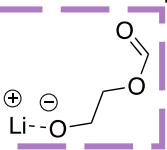


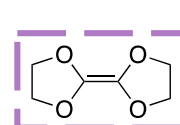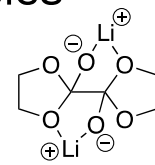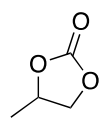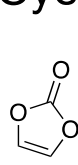LiEDC$^-$          LiBDC$^-$



**Carboxylates, esters, and oxides**



LFEO

**Cyclic species**



bi-dioxolylidene

BERKELEY LAB

# Predicted Battery Network Products: 36 out of 5139

**Small molecules/gases**



**Inorganics**



lithium carbonate      lithium oxalate

[EC, Li$^+$] and [EC, Li$^+$, CO$_2$] at 0V and +0.5V vs. Li/Li$^+$

**Alkyl carbonates**



LMC          vinyl carbonate          ethylene monocarbonate

LiEDC$^-$          LiBDC$^-$

- Recovered nearly all observed or proposed molecular SEI components
- Only thermodynamics – unexpectedly effective!
- So about those particularly weird molecules...

**Carboxylates, esters, and oxides**

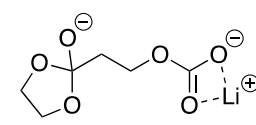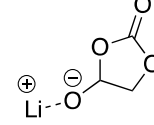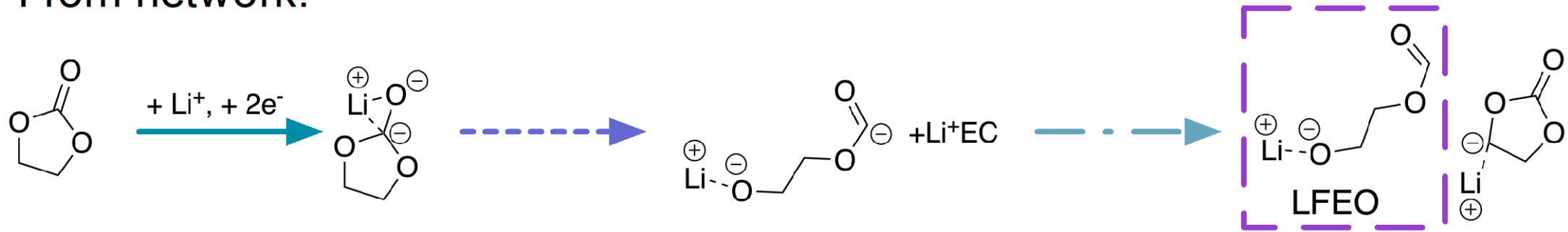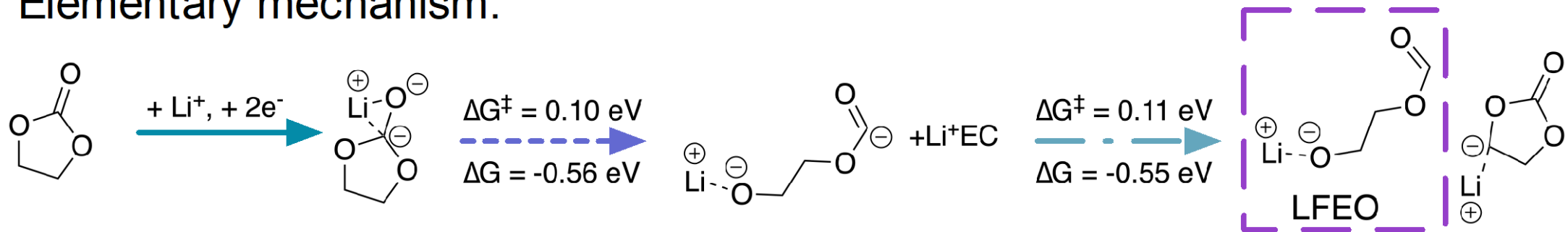

LFEO

**Cyclic species**



bi-dioxolylidene

BERKELEY LAB

# Network Path to Refined Mechanism: LFEO

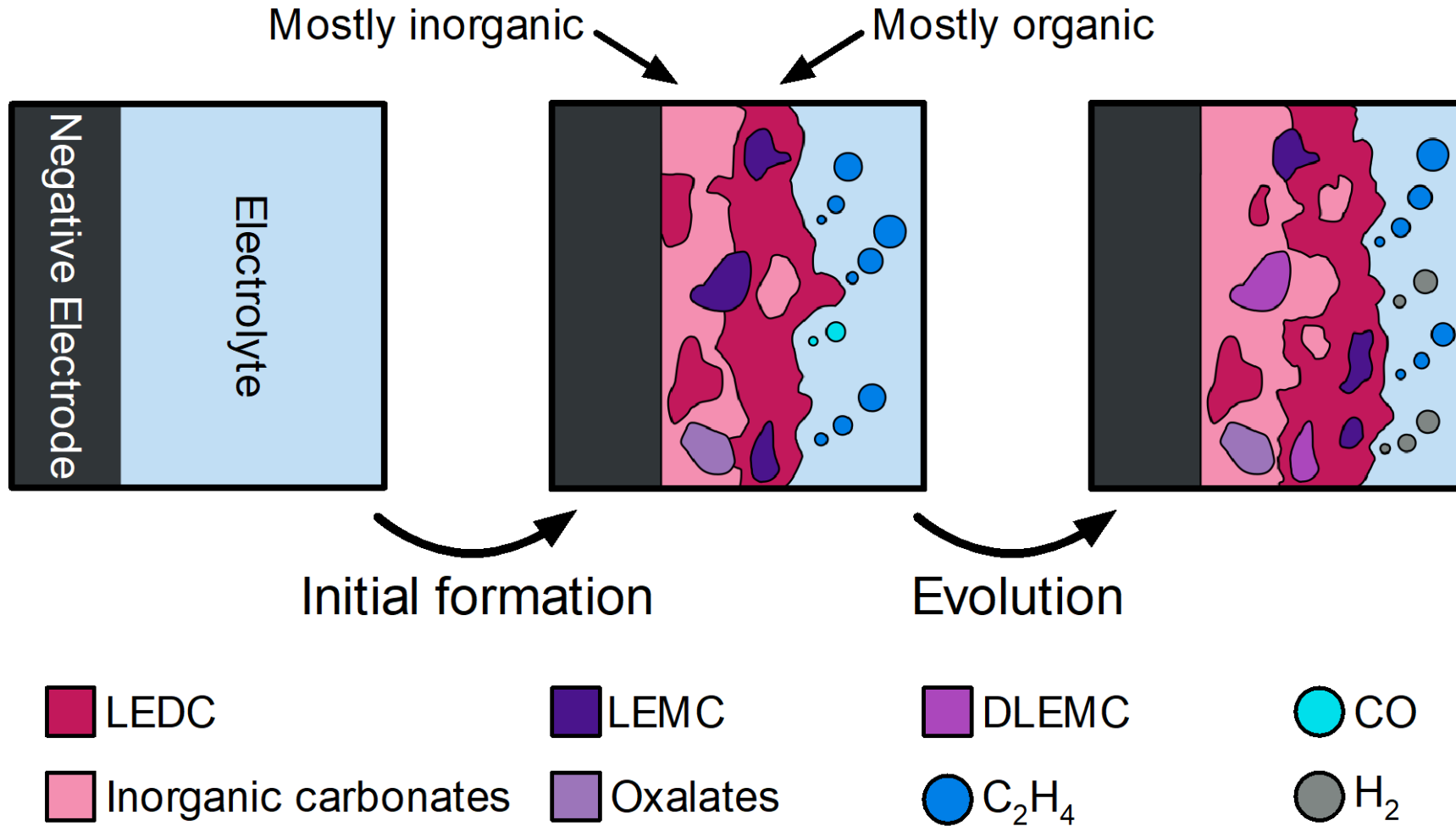Applied semi-automated TS procedure to 15 shortest thermo. paths – 12th shortest with [Li+, EC] at 0V vs Li/Li+:

From network:



Elementary mechanism:



$\Delta G^{\ddagger} = 0.10$ eV
$\Delta G = -0.56$ eV

$\Delta G^{\ddagger} = 0.11$ eV
$\Delta G = -0.55$ eV

D. Barter*, E. W. C. Spotte-Smith*, N. Redkar, S. Dwaraknath, K. A. Persson, **S. M. Blau**, *Dig. Disc.* 2023
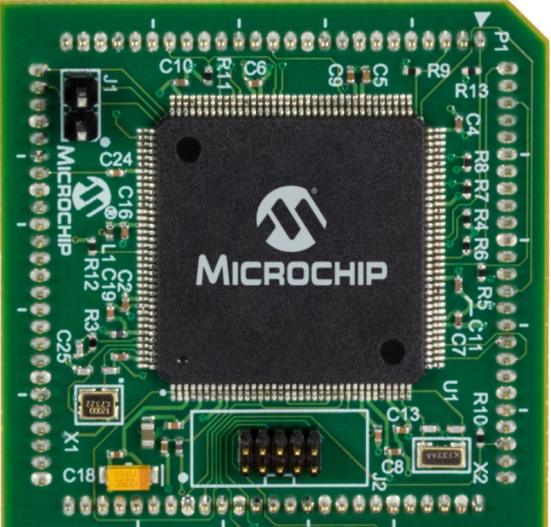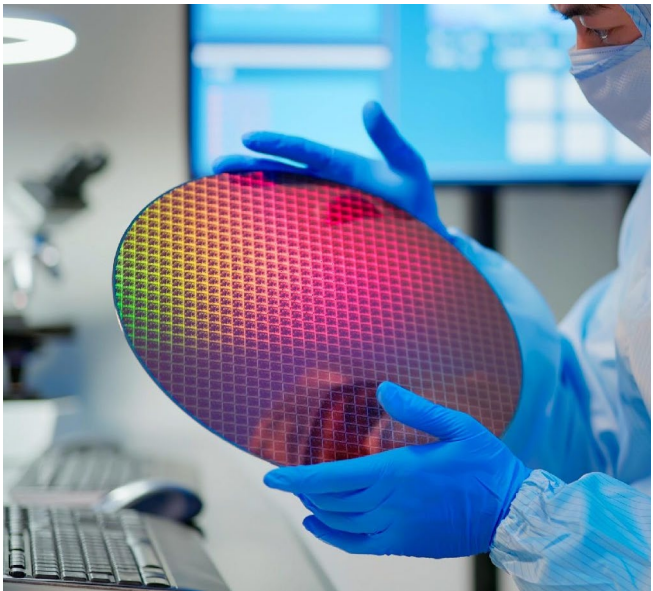
BERKELEY LAB

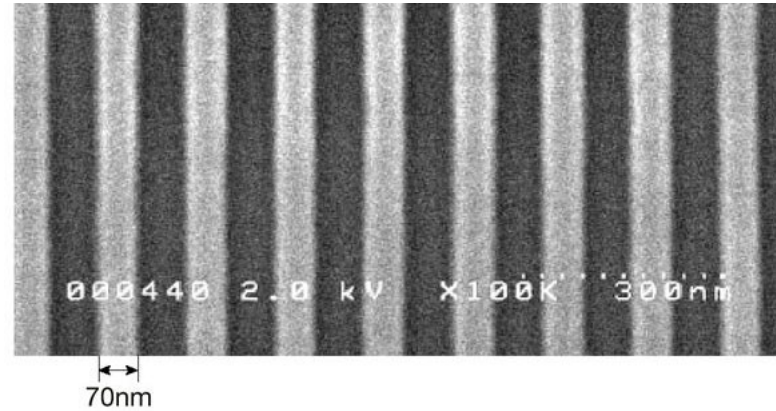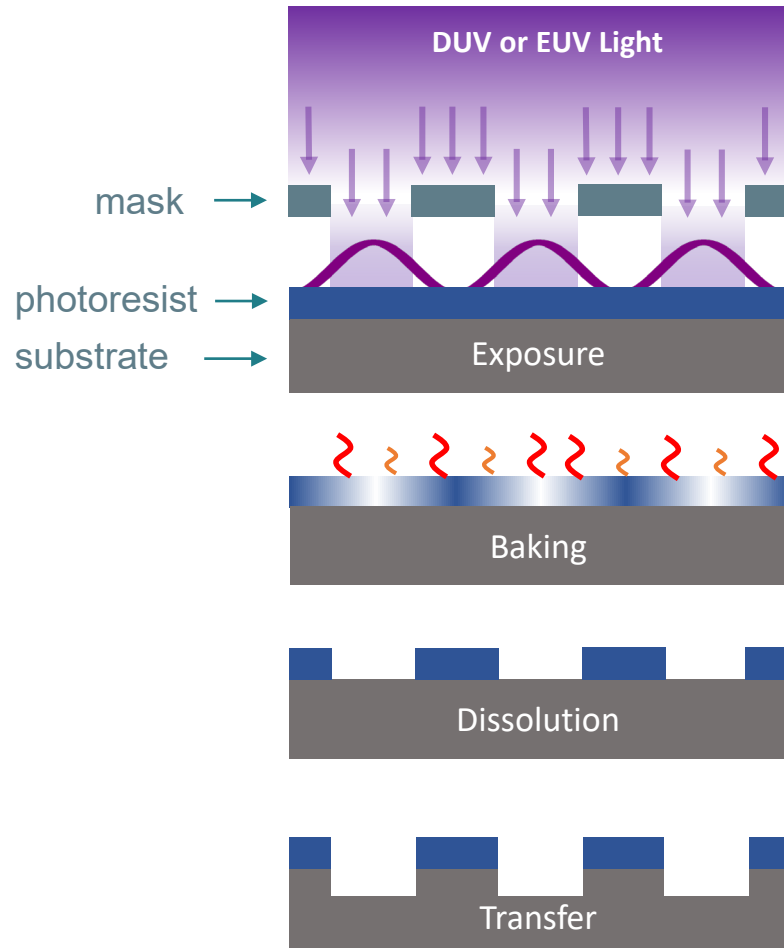# Mechanistic Model of SEI Formation Derived from CRN



- Pathways derived from CRN, semi-automated $\Delta G^{\ddagger}$ calcs

- Recovered bi-layer SEI from first principles for first time

- **Is this approach limited to just SEI formation? No!**

E. Spotte-Smith*, R. Kam*, D. Barter, X. Xie, T. Hou, S. Dwaraknath, **S. M. Blau**, K. A. Persson, *ACS Energy Lett.* 2022

BERKELEY LAB

# Background: Nanoscale Patterning with Photolithography

BERKELEY LAB

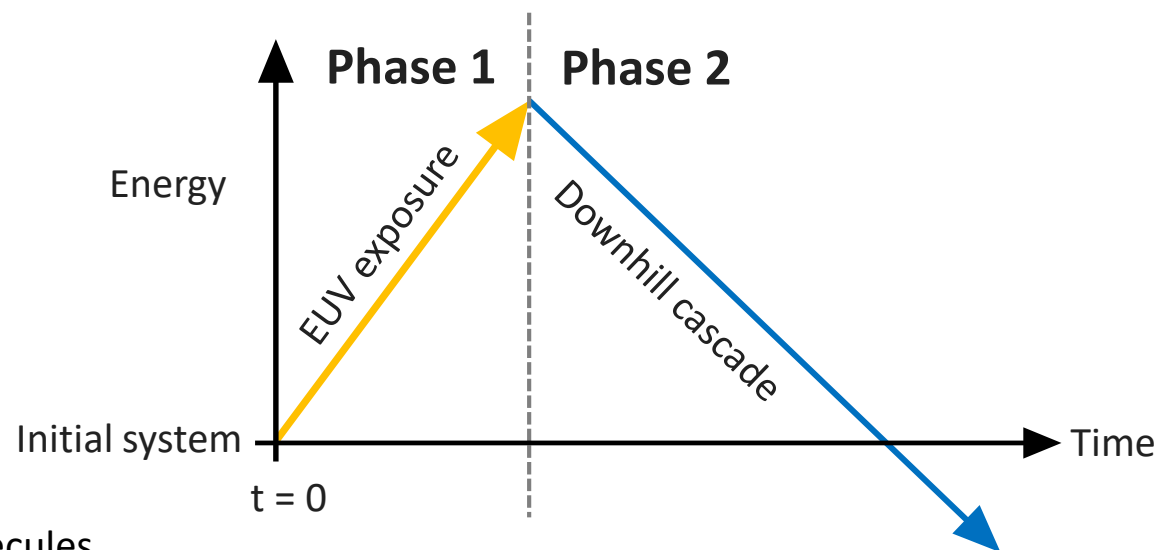# Background: Nanoscale Patterning with Photolithography



Chemical reactions
cause solubility switch

- 1994 to 2017: "deep" UV, 248 nm – 134 nm light
  - 5 eV – 9 eV photons
  - **Selective** resonant photochemistry
- Want smaller patterns? Need shorter wavelength!
- 2018 to now: "extreme" UV, 13.5 nm light
  - 92 eV photons
  - Stochastic photoionization yields **poorly understood** radical ion reaction cascade
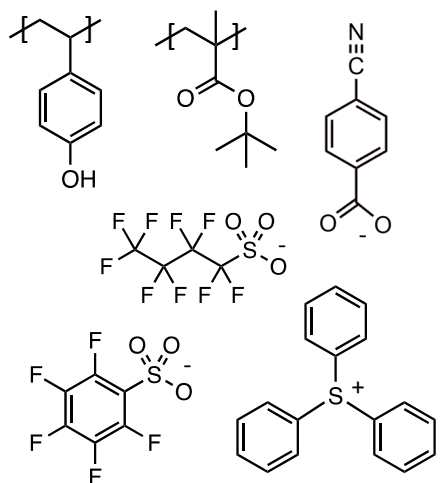
# EUV Lithography Reaction Network Construction



**Phase 1:** EUV exposure

*Keep* reactions that are:
- Electron detachment (all $\Delta G > 0$)
- Electron attachment (all $\Delta G < 0$)
- $\Delta G > 0$ one-bond fragmentation
- $\Delta G > 0$ $H^+$ or $H^0$ transfer

185,929 reactions

Principal molecules

Fragment / DFT opt. → 108 species → Recombine / DFT opt. → 3367 species → HiPRGen
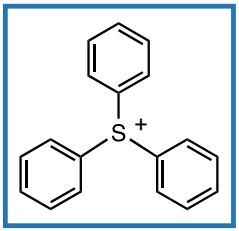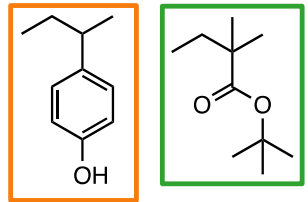
pymatgen   FireWorks

Custodian   atomate

**Phase 2:** Post-exposure cascade

*Remove* reactions with:
- $\Delta G > 0$
- Unbalanced redox
- >2 covalent bonds changing
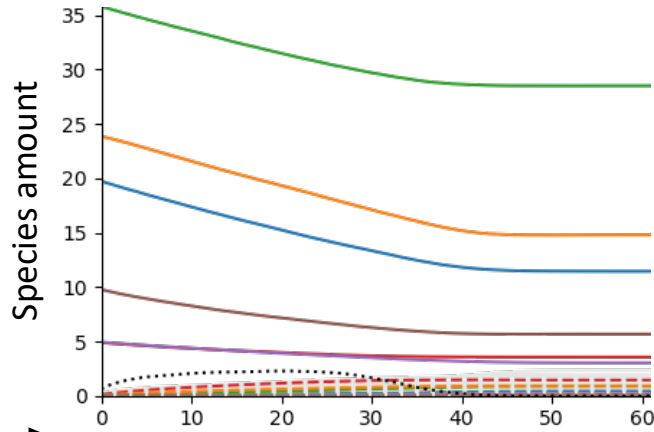- Sterically hindered reaction center

2,776,867 reactions
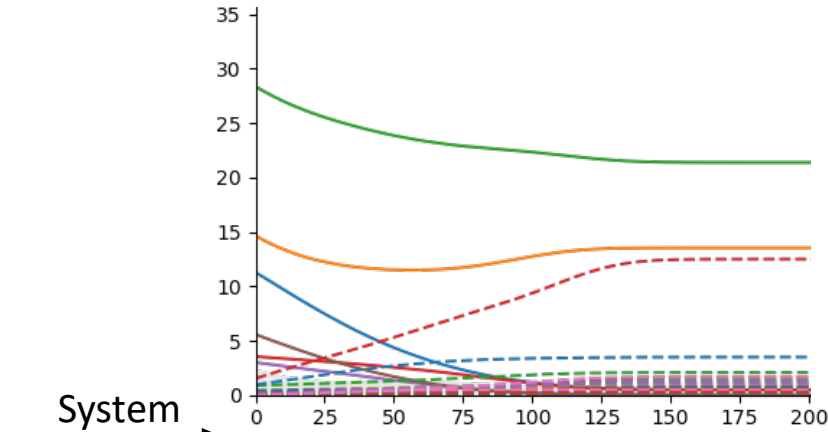
BERKELEY LAB

20

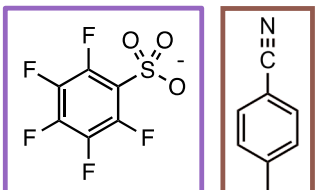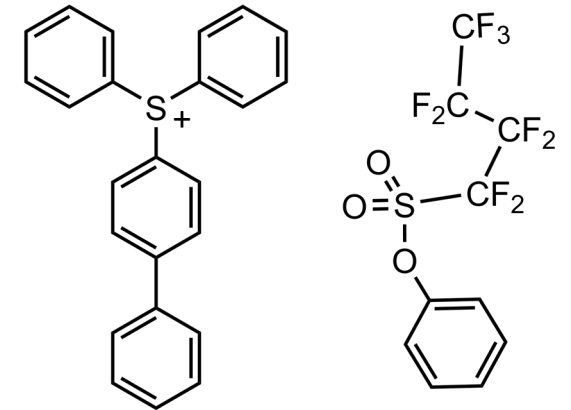# EUV Lithography Reaction Network Analysis



**Phase 1:** EUV exposure

**Phase 2:** Post-exposure cascade

Initial state

System state

Network products include:

$CO_2$

Free electron

Trajectory step

- +92 eV "energy budget"
- Explicit free electron species

RNMC

Normalized Intensity ($I/N_0$)

m/z

Cruz et al. *Proc. SPIE* 2022

BERKELEY LAB

21

# Recap: The Steps of Building and Analyzing a CRN

## 1. Species generation



## 2. Reaction generation



## 3. Pathway sampling



## 4. Identify Products



Transition state calcs
Build kinetic models

Discover novel important species and pathways

Under development: ML-assisted network expansion
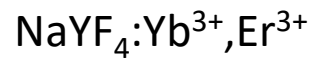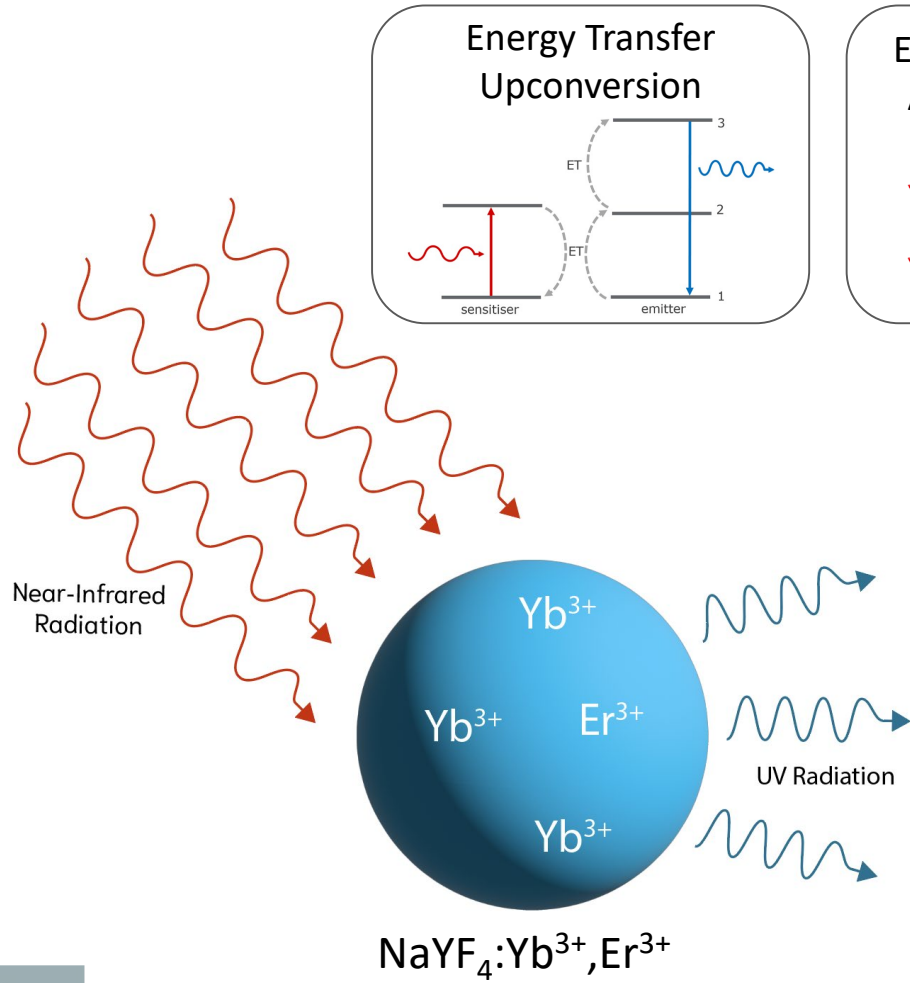
# Background: Upconverting Nanoparticles (UCNPs)



Energy Transfer Upconversion

Excited State Absorption

Near-Infrared Radiation

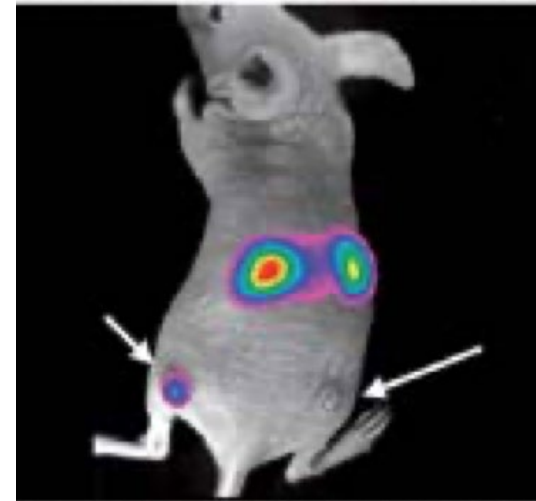Yb³⁺

Yb³⁺  Er³⁺

Yb³⁺
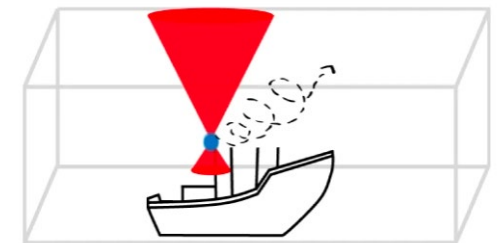
UV Radiation
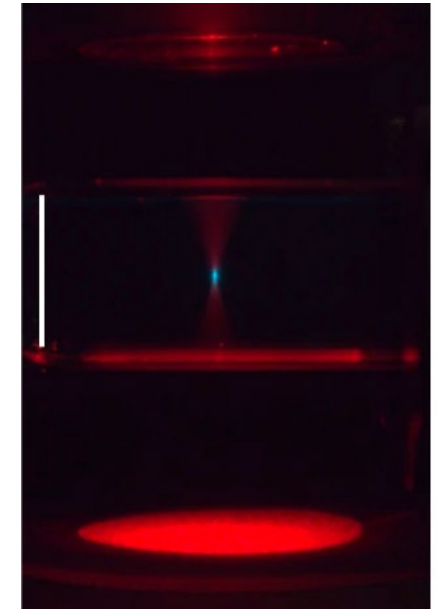
$NaYF_4:Yb^{3+},Er^{3+}$

Security Printing

Lu et al. *Nat. Photon.* 2014

Bio-imaging

Xiong et al, *Anal. Chem.* 2009

3D Printing

Sanders et al, *Nature* 2022

BERKELEY LAB

23

# UCNP Doping and Heterostructure

Host Material: $NaYF_4$



Dope $Y^{3+}$ sites with $Ln^{3+}$



$x_{Yb}$, $x_{Er}$, $x_{Nd}$

$x_{Yb}$, $x_{Er}$, $x_{Nd}$

Core    Shell 1    Shell 2    Shell 3

| 58 | 59 | 60 | 61 | 62 | 63 | 64 | 65 | 66 | 67 | 68 | 69 | 70 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|
| Ce | Pr | Nd | Pm | Sm | Eu | Gd | Tb | Dy | Ho | Er | Tm | Yb |

BERKELEY LAB

# UCNP Photophysics Can Be Simulated With kMC

**Transition rate constants**



Teitelboim et al. *J. Phys. Chem. C* 2019

**Generate ensemble of randomly doped structures**



Energy Transfer
Kinetic Monte Carlo (kMC)

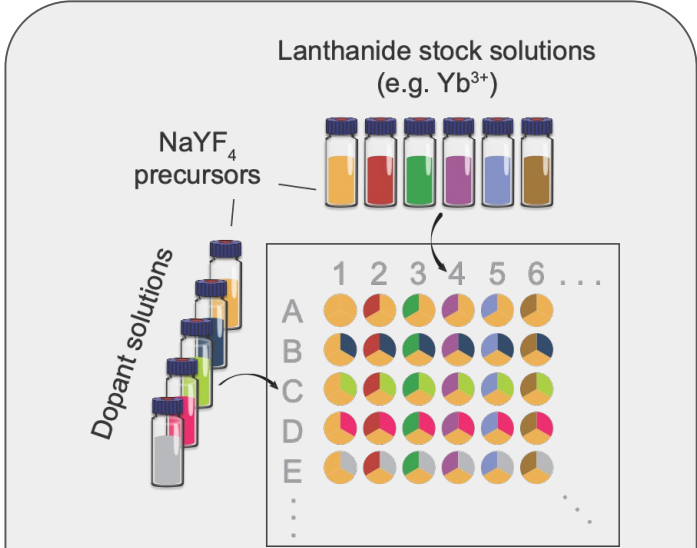github.com/BlauGroup/NanoParticleTools
github.com/BlauGroup/NPMC

- Simulations require 10-150+ hours on one CPU core
  - Cannot be parallelized

**Spectra**

BERKELEY LAB

# Large Search Space Necessitates Intelligent Searching



Lanthanide stock solutions (e.g. Yb$^{3+}$)

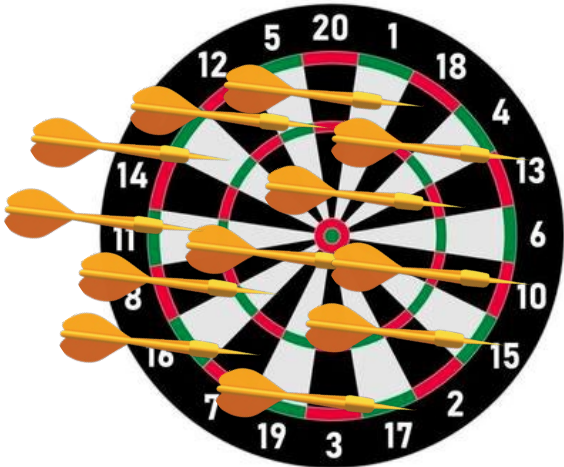NaYF$_4$ precursors

Dopant solutions

**Combinatorial/Robotic Synthesis**
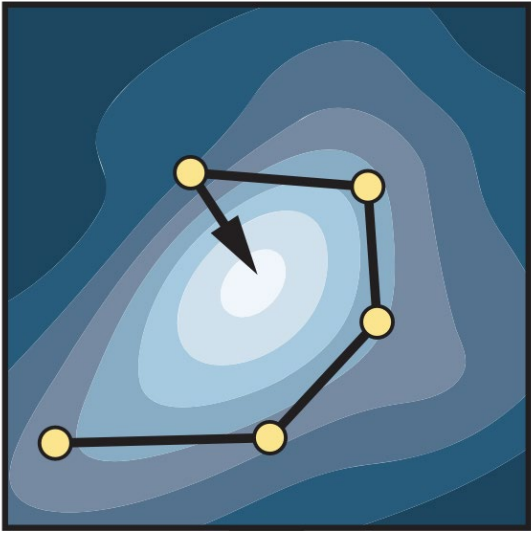
**Consider a simple spherical nanoparticle:**

- <u>Chose up to 4 dopants</u> (of 13 lanthanides)

  1,093 combinations

- <u>3 Dopant concentrations</u> - Low, Medium, High

  66,379 dopant configurations

- <u>5 particle sizes</u> - 4, 6, 8, 10, & 12nm

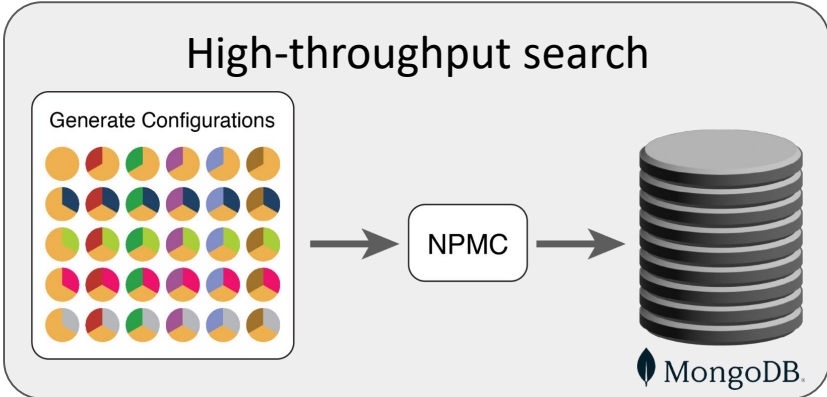  265,516 nanoparticle configurations

**Inverse Design**



Sanchez-Lengling et al. *Science* 2018



**High-throughput search**

Generate Configurations → NPMC → MongoDB

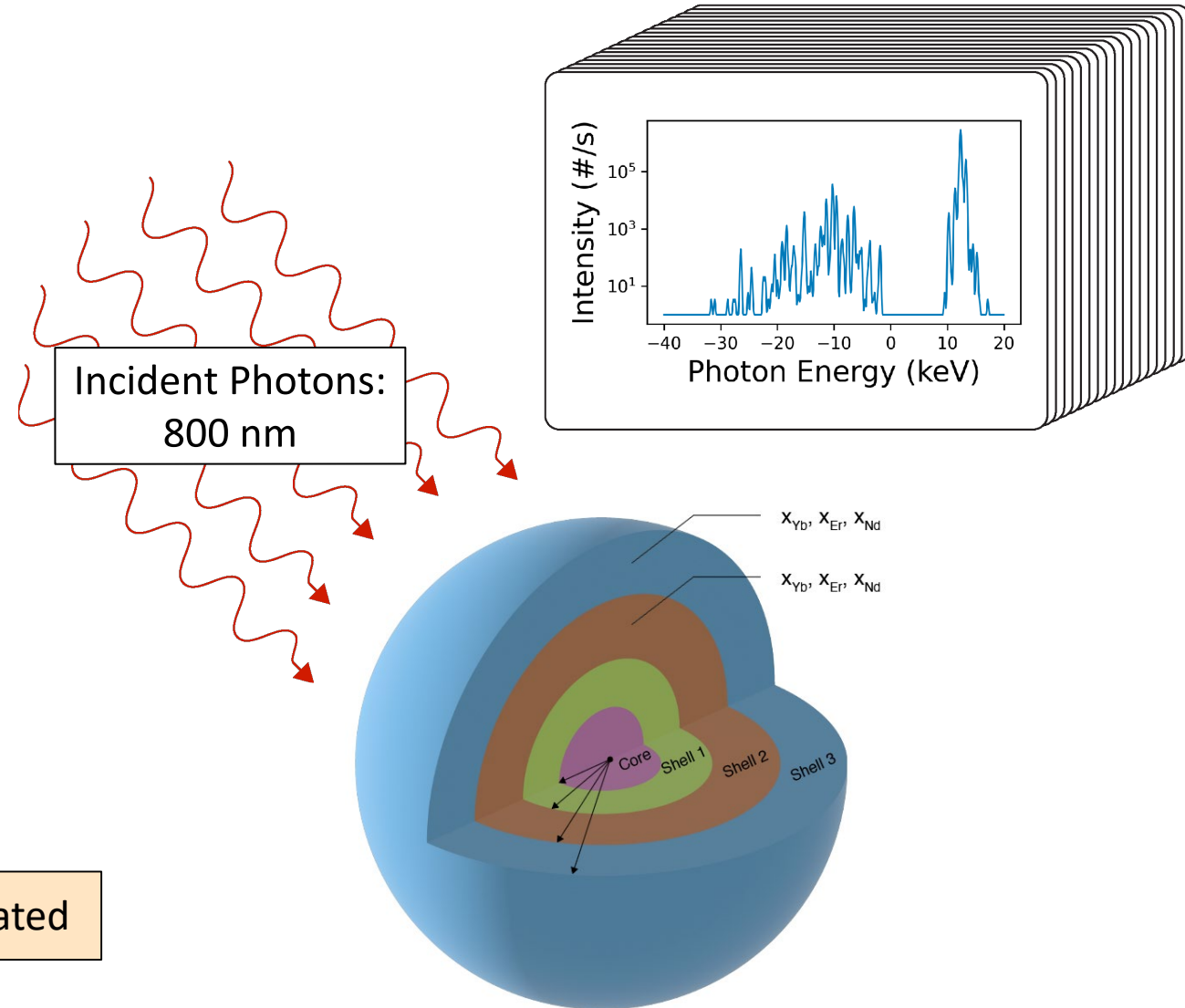BERKELEY LAB

# Generating a Dataset for Machine Learning

## IID Dataset:

- Up to 8 nm diameter core
- Up to 3 shells
  - Each shell is 1-2.5 nm thick
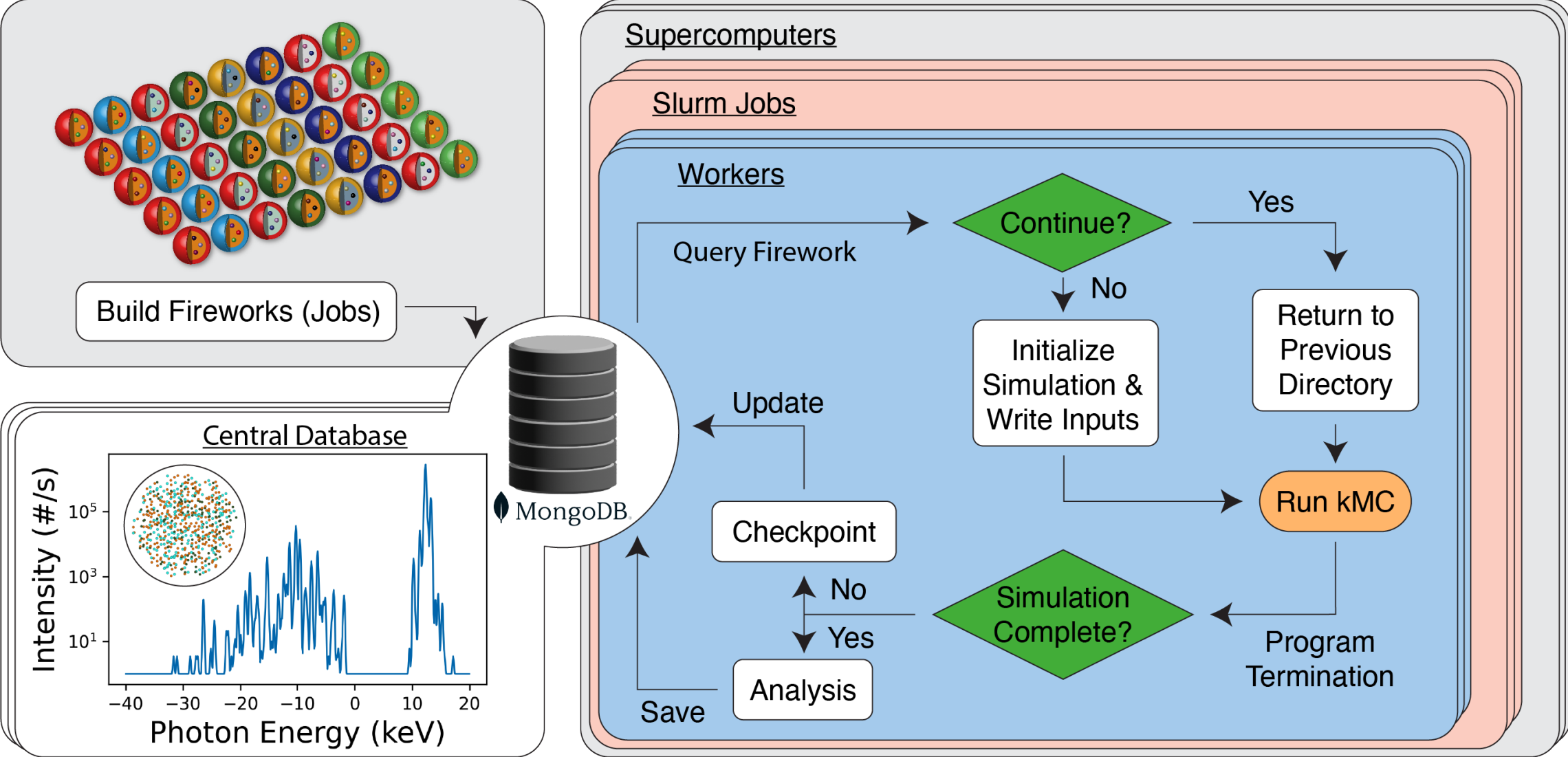- Consider only Yb, Er, and Nd dopants

## OOD Testing Dataset:

- Up to 8 nm diameter core
- 4 shells
  - Each shell is 1-2.5 nm thick
- Consider only Yb, Er, and Nd dopants

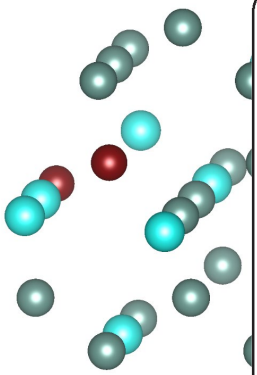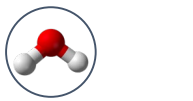>6,000 nanoparticle configurations/spectra simulated

Incident Photons: 800 nm



$x_{Yb}, x_{Er}, x_{Nd}$

$x_{Yb}, x_{Er}, x_{Nd}$

Core  Shell 1  Shell 2  Shell 3

BERKELEY LAB

# UCNP kinetic Monte Carlo Simulation Workflow

BERKELEY LAB

# Representations of Nanoparticles for Machine Learning

## Atomistic Representation

$l = 0, 1, 2, \ldots$

| ? | Prediction Accuracy and Generalizability |
| :-: | :-- |
| ✗ | Gradients w.r.t dopant conc. |
| ✗ | Gradients w.r.t layer radii |

## Image Representation

height

width

| ✓ | Prediction Accuracy and Generalizability |
| :-: | :-- |
| ✓ | Gradients w.r.t dopant conc. |
| ✗ | Gradients w.r.t layer radii |

height

width

Dopant Concentrations

## Tabular Representation

| $r_0$ | $V_0$ | $x_{0,Yb}$ | | | $x_{n,Yb}$ | $x_{n,Er}$ | $x_{n,Nd}$ |
| :-: | :-: | :-: | :-: | :-: | :-: | :-: | :-: |
| | | | | | | | |
| | Core features | | | | Layer features | | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

| ✗ | Prediction Accuracy and Generalizability |
| :-: | :-- |
| ✓ | Gradients w.r.t dopant conc. |
| ✓ | Gradients w.r.t layer radii |

BERKELEY LAB

# Developing a Physics-Infused Graph Representation



Dopant/Control Volume Graph

Dopant Interaction

One-Hot encoding of dopant

Dopant Concentration

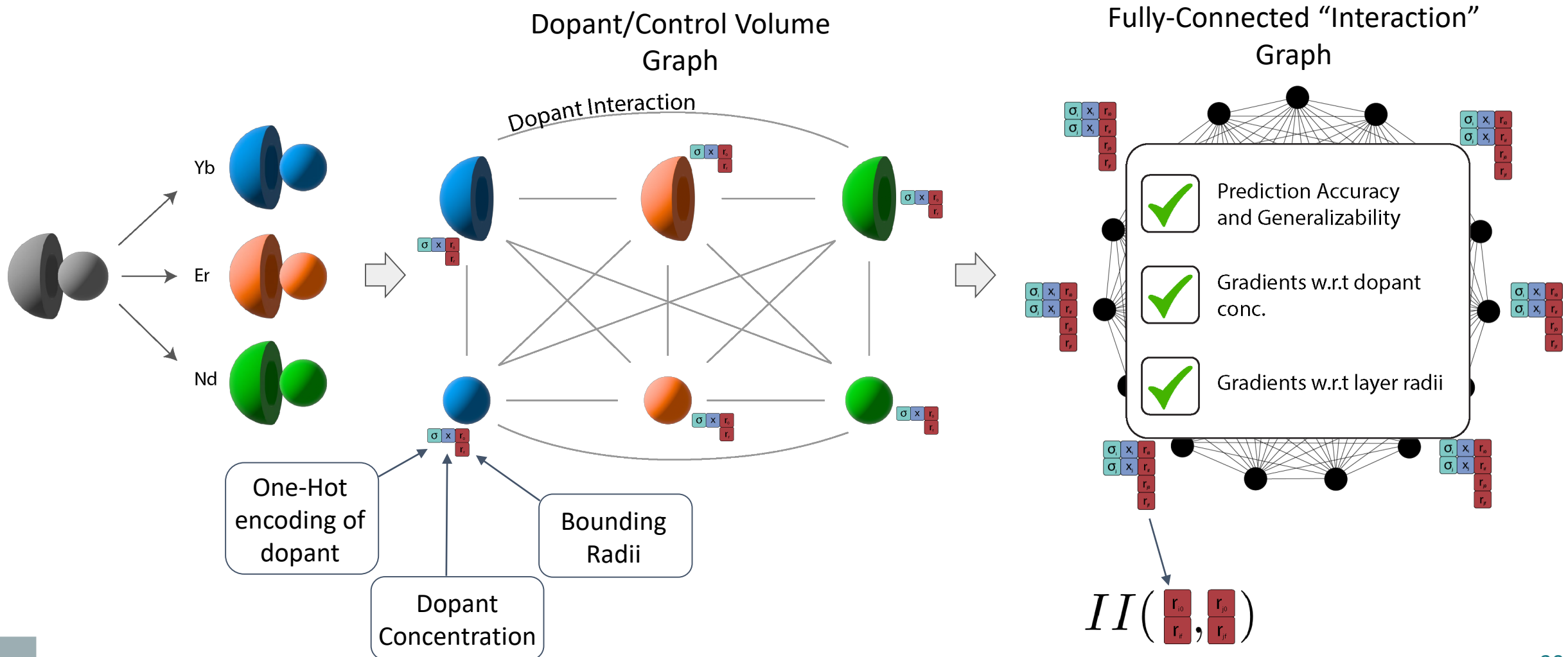Bounding Radii

Edge feature - Integrated Interaction

$$I_{ij} = \frac{1}{\sigma\sqrt{2\pi}} e^{-1/2 \frac{s_{ij}}{\sigma}^2}$$

Layer k

Layer l

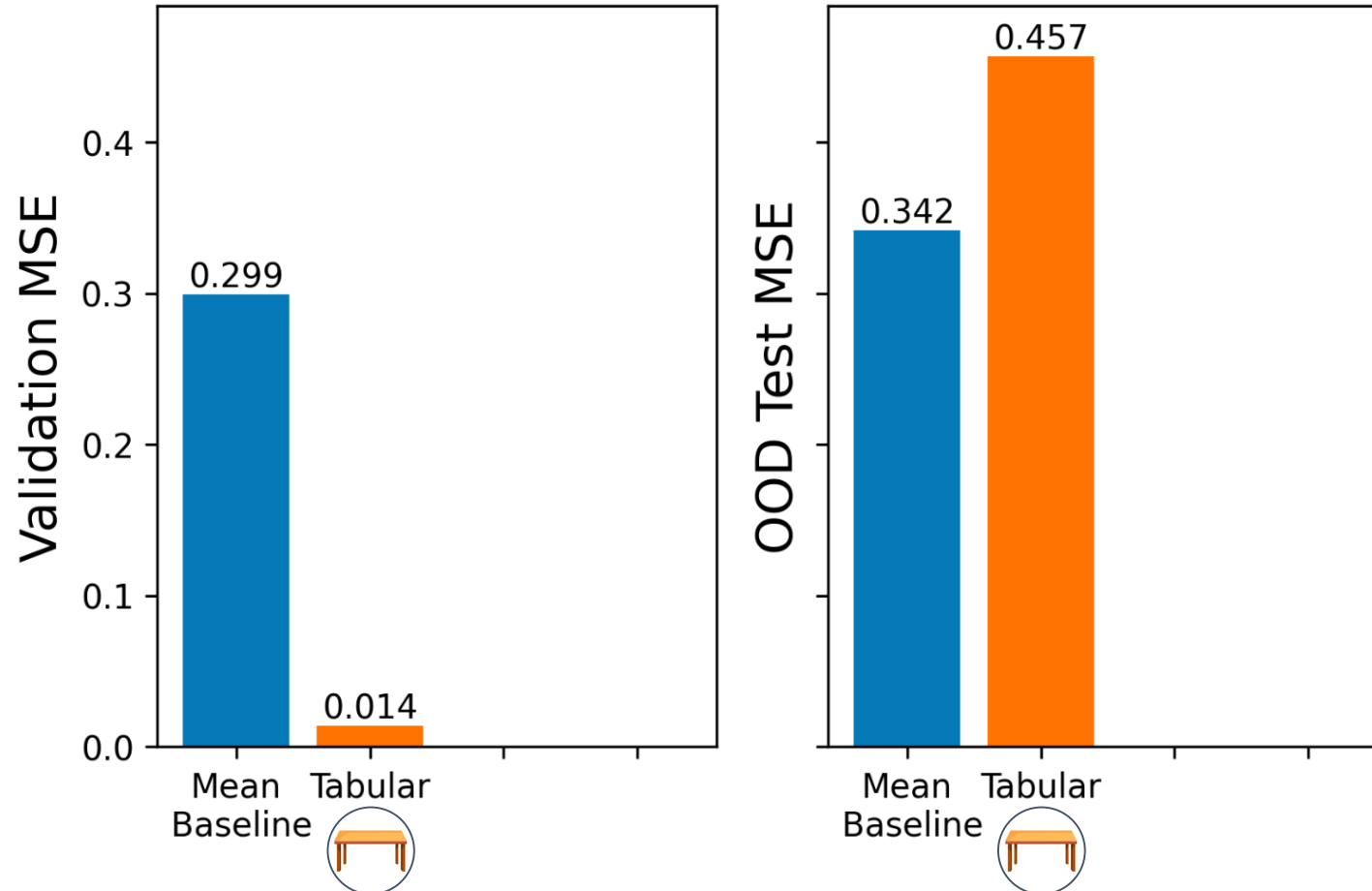$$II\left(\begin{array}{c} r_{i_0} \\ r_{i_f} \end{array}, \begin{array}{c} r_{j_0} \\ r_{j_f} \end{array}\right) = \int \int I_{ij} dV_k dV_l$$

# Developing a Physics-Infused Graph Representation



**Dopant/Control Volume Graph**

Dopant Interaction

One-Hot encoding of dopant

Dopant Concentration

Bounding Radii

**Fully-Connected "Interaction" Graph**

✓ Prediction Accuracy and Generalizability

✓ Gradients w.r.t dopant conc.

✓ Gradients w.r.t layer radii

BERKELEY LAB

# Comparing Tabular vs. Image vs. Graph Rep. Performance

BERKELEY LAB

# Comparing Tabular vs. Image vs. Graph Rep. Performance

BERKELEY LAB

# Comparing Tabular vs. Image vs. Graph Rep. Performance

BERKELEY LAB

# Inverse Design of Nanoparticles Via Gradient Ascent



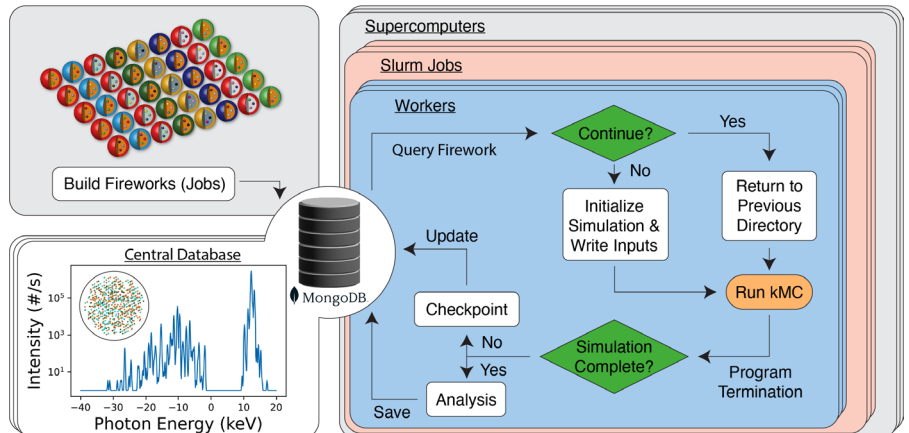E. Sivonxay, E. Chan, **S. M. Blau**, *In preparation*
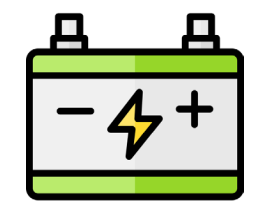
32

# Summary



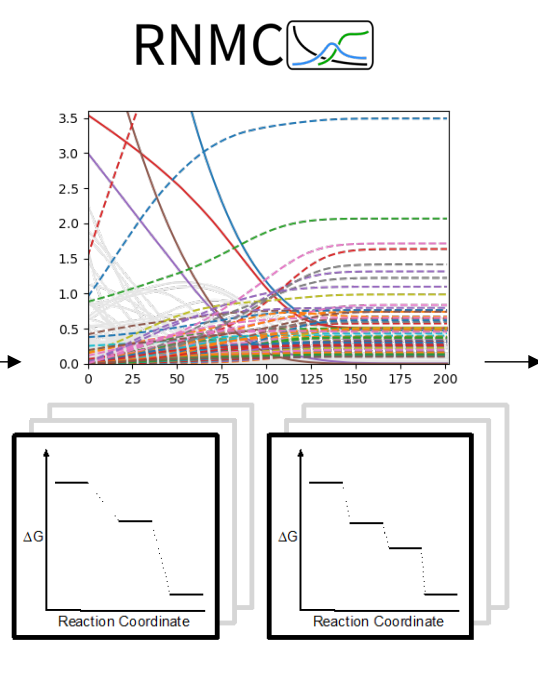Network with up to ~200,000,000 reactions

HiPRGen

RNMC

LFEO
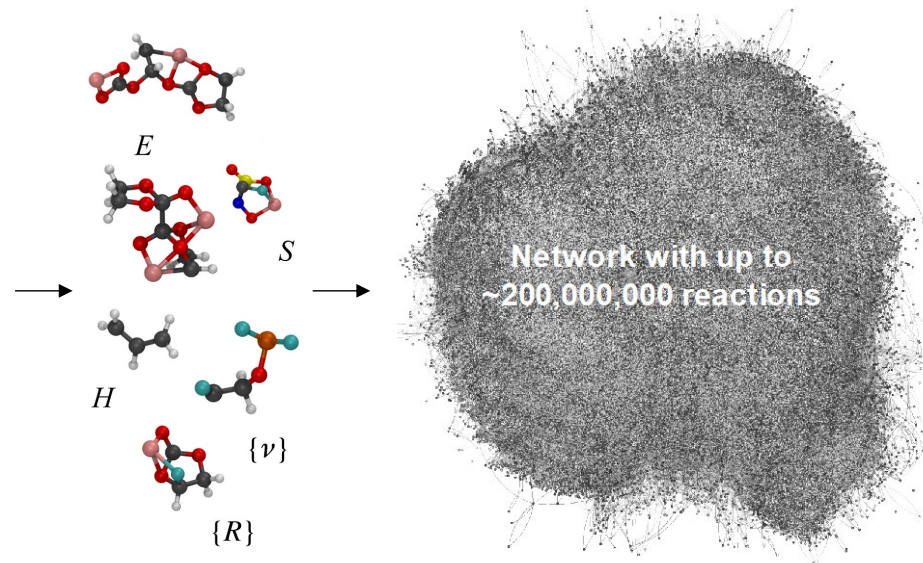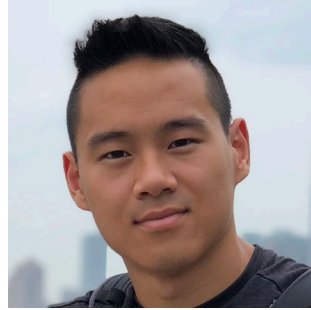
# Acknowledgements



Daniel Barter   Evan Spotte-Smith   Eric Sivonxay   Jacob Milton   Frances Houle   Kristin Persson   Emory Chan

BERKELEY LAB

34