# Exploring chemical space: Molecular Space Shuttle

Alán Aspuru-Guzik
Professor
Harvard University

http://aspuru.chem.harvard.edu
Twitter: A_Aspuru_Guzik
aspuru@chemistry.harvard.edu

How large is space?

$10^{82}$ atoms in the observable universe

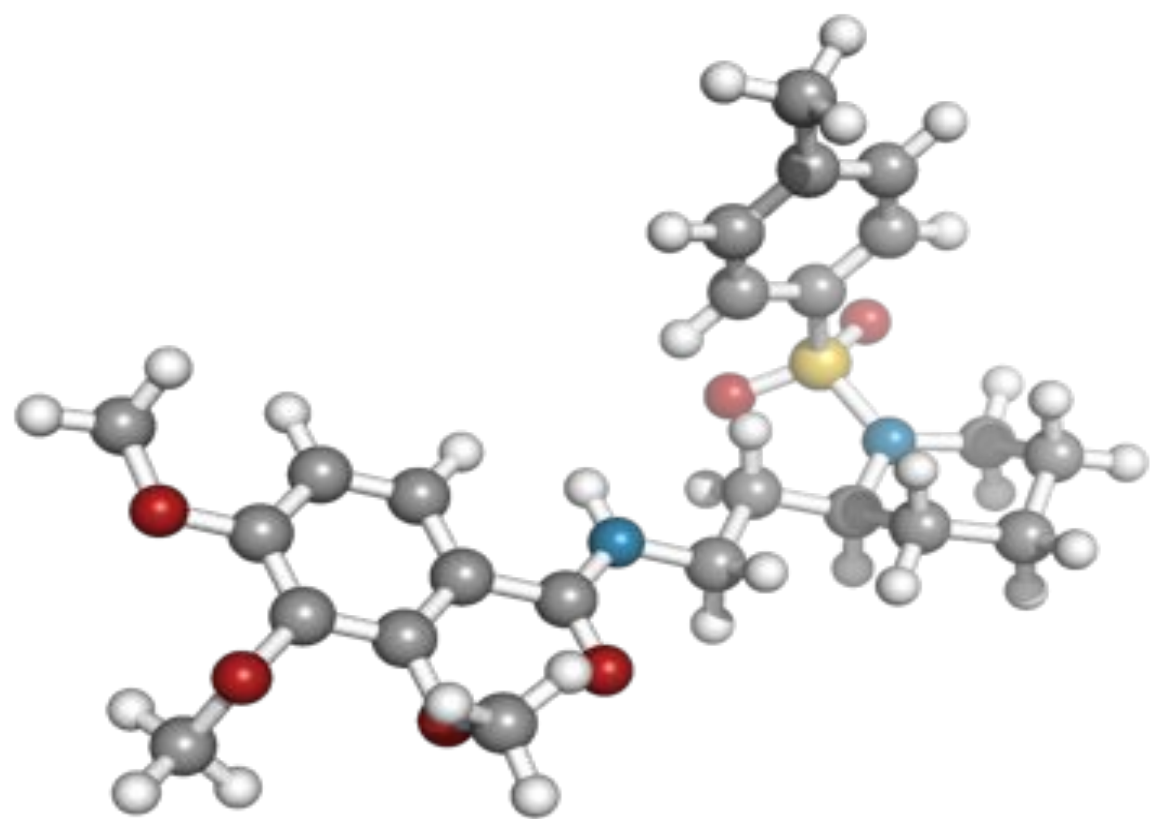How large is chemical space?

$10^{60} - 10^{180}$ medium-size molecules

# Molecular screening for organic materials

How good is this molecule as a battery material?

Quantum Mechanics

Machine Learning

# From $10^{60}$ to $10^6$ to 10...



Initial library

Computational
screening

Synthesis
and testing

**Computational
cost**

**Molecules most likely
to be of interest**

Organic materials in the larger context

US Materials Genome Initiative

Pyzer-Knapp, et al. Ann Rev. Mat Sci. (2015)

Materials Innovation Infrastructure

- Human Welfare
- Clean Energy
- National Security
- Next Generation Workforce
- Computational Tools
- Experimental Tools
- Digital Data

Organic Materials

Inorganic Materials

Organic Pharmaceuticals

**Shared Features**
- Timescale is important
- Automated techniques
- Data-driven discovery
- Computational funnel

- Number of descriptors
- Size of search space
- Level of approximation

Large
Medium
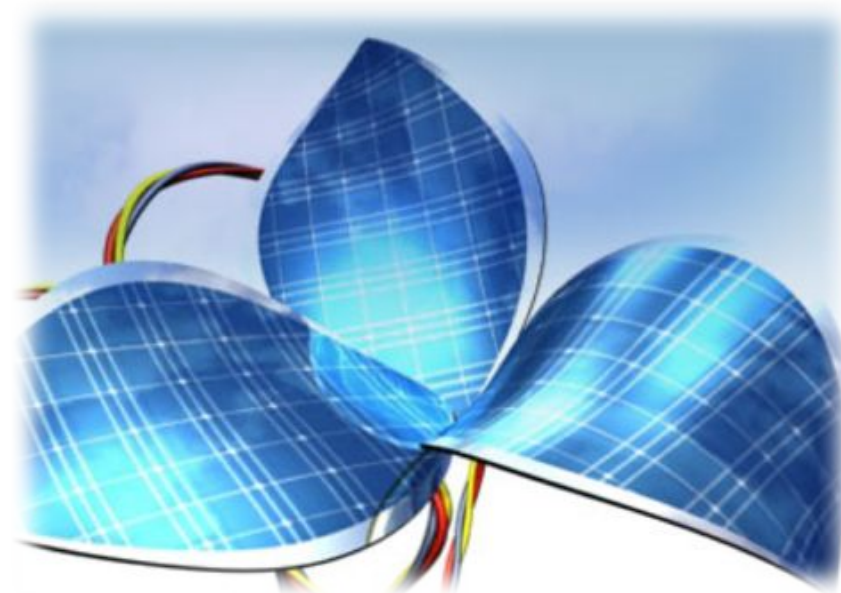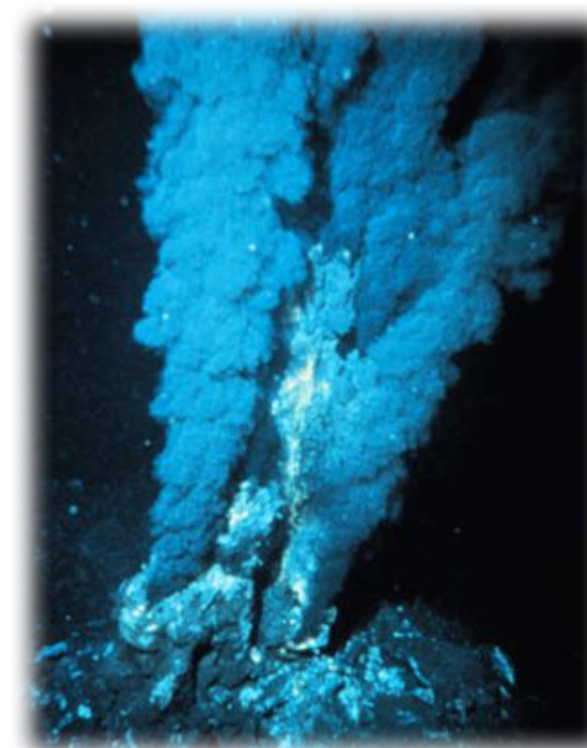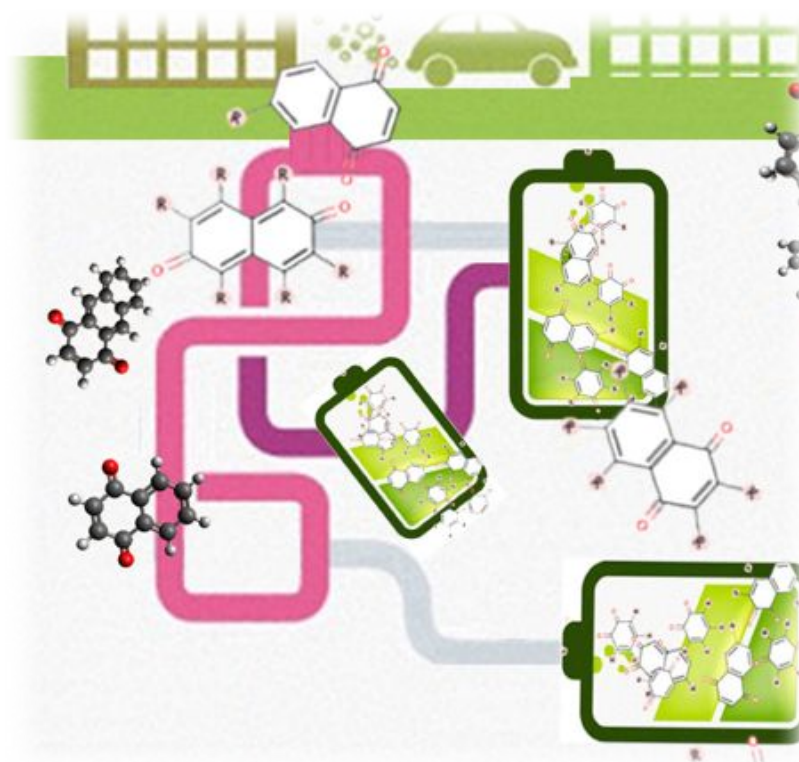Small

# My research group's explorations of chemical space



**The Harvard Clean Energy Project**
Generating renewable energy

**Blue Organic LED**
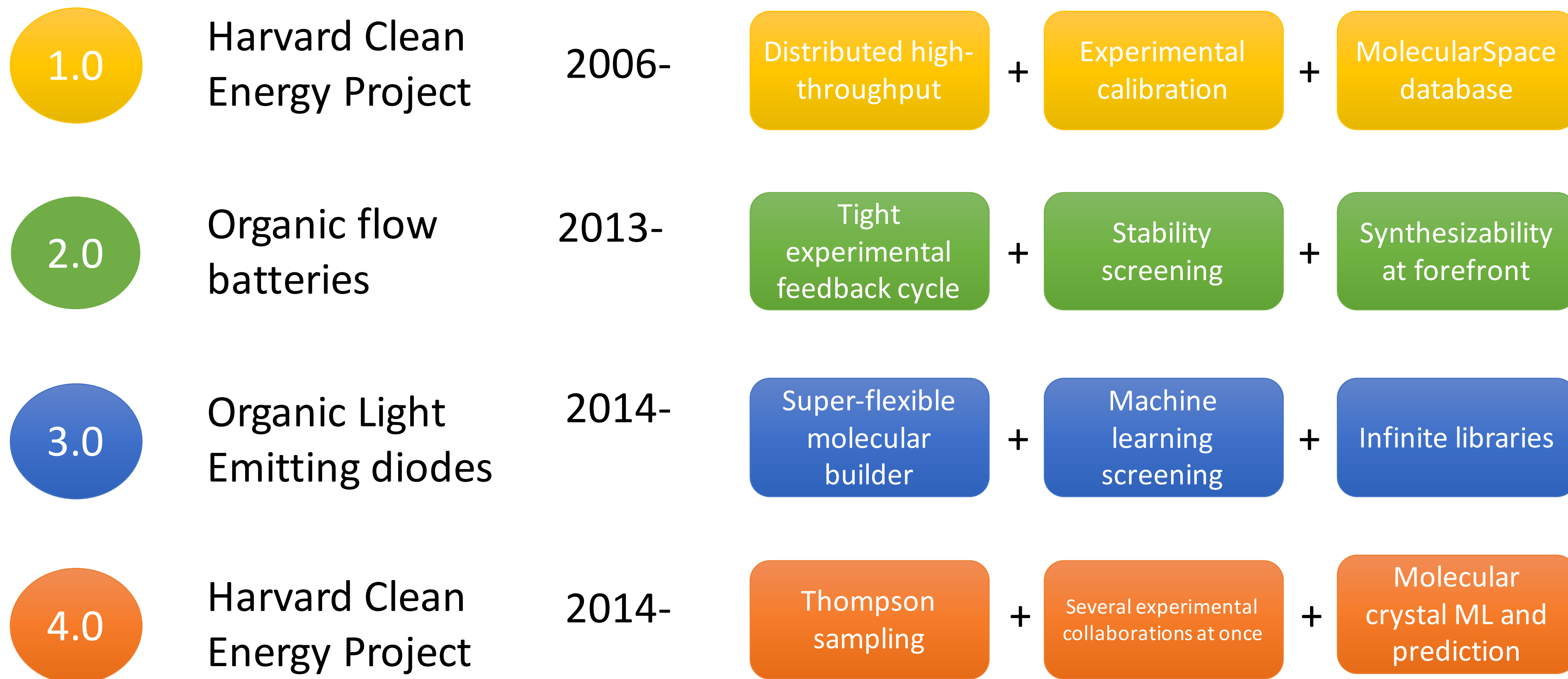For your next gadget or TV

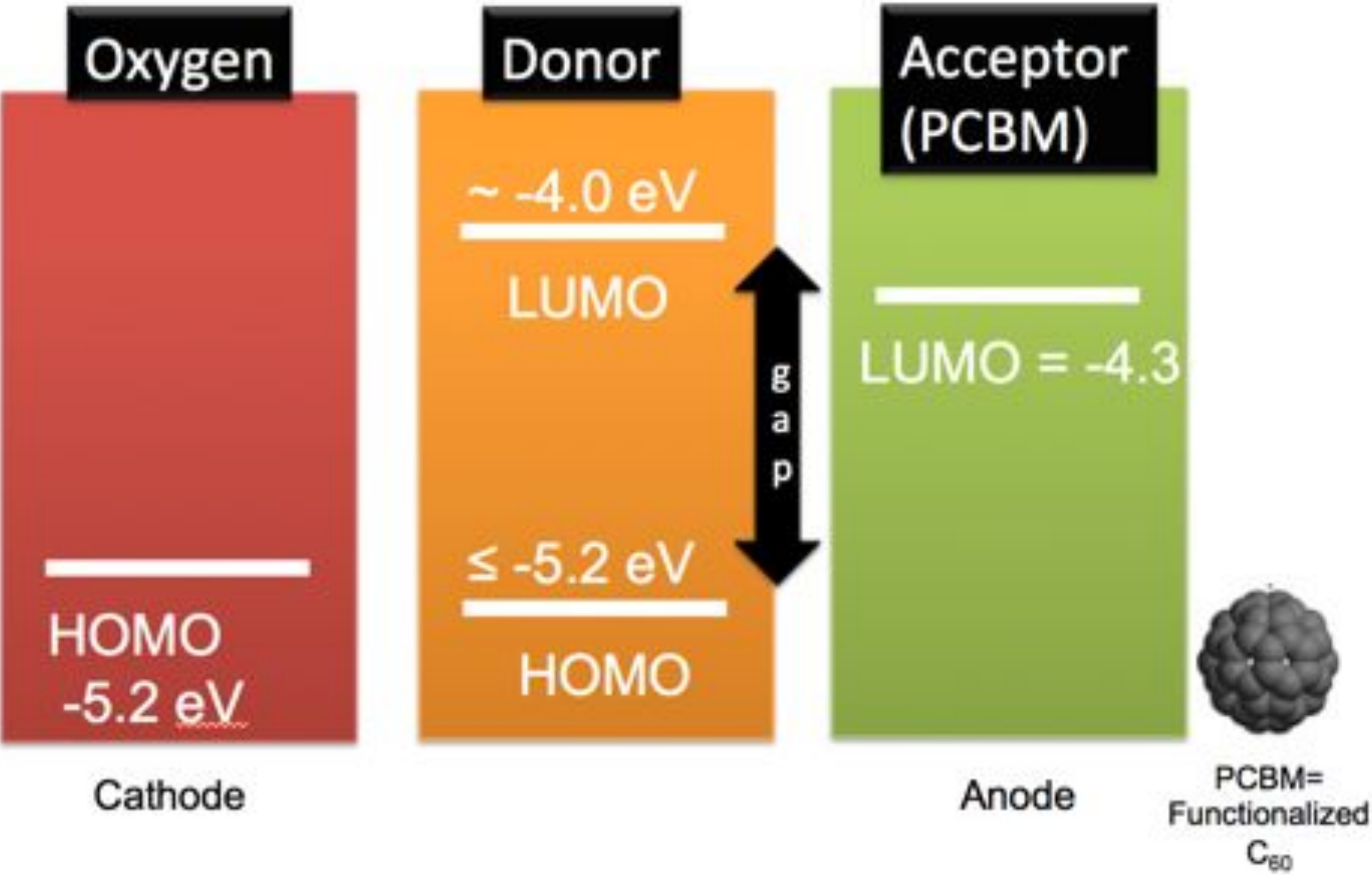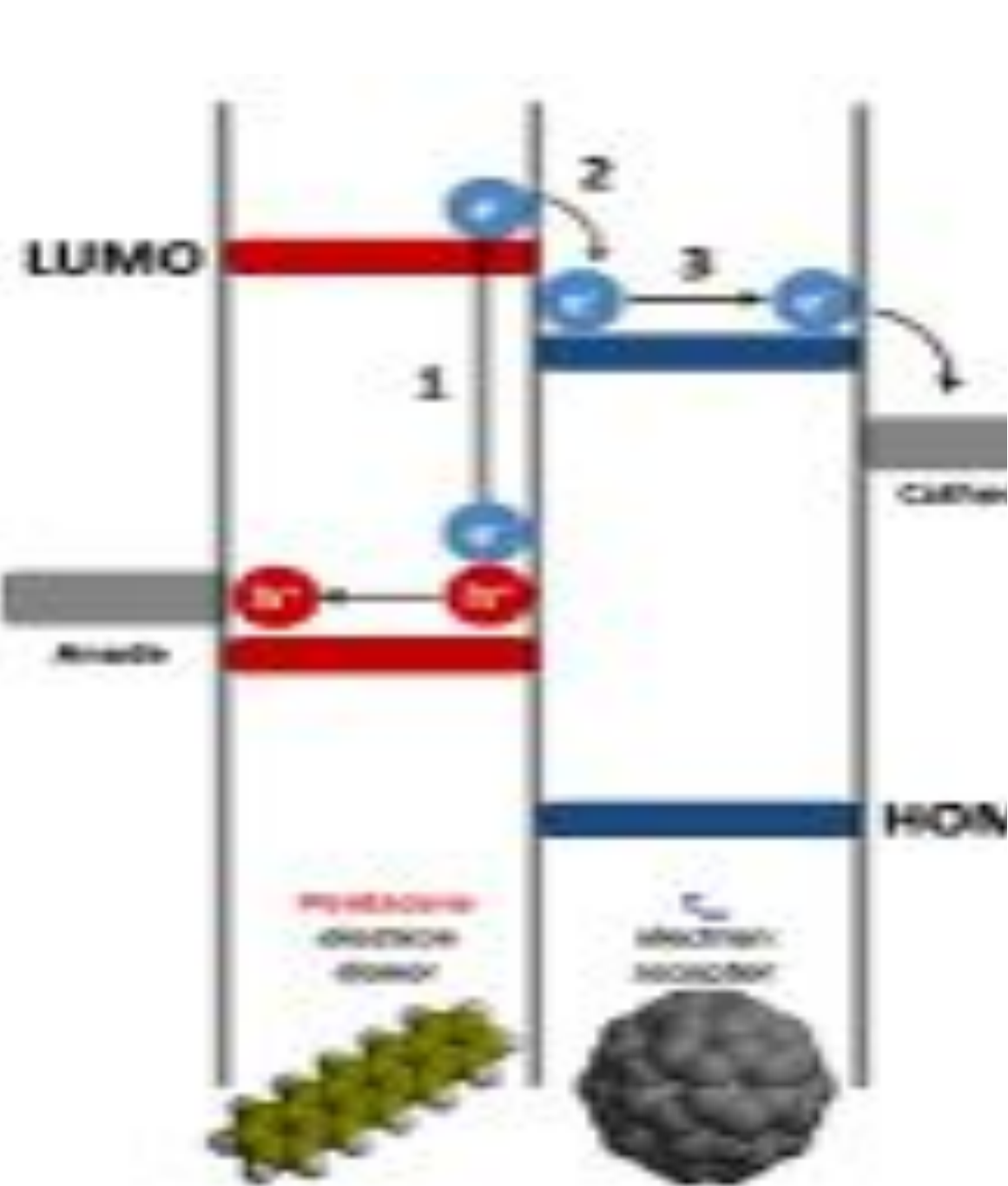**Organic flow batteries**
Storing renewable energy

**Origins of life**
How life may have come about?

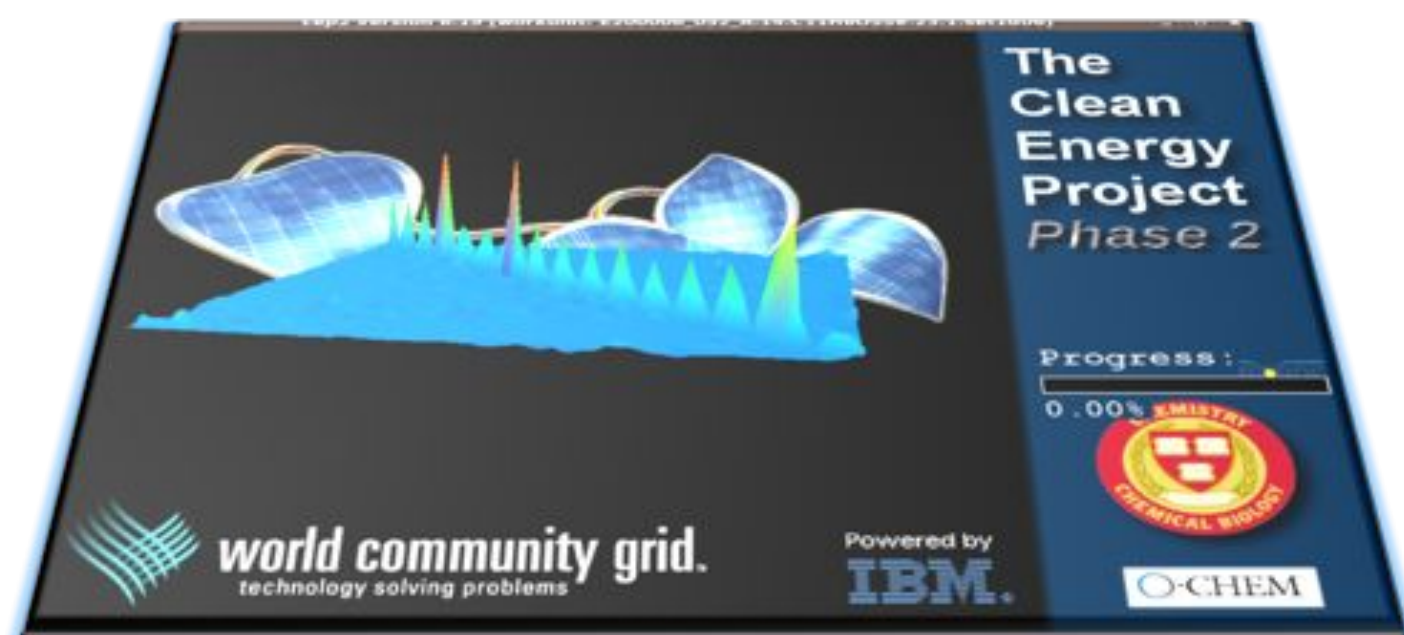# Project chronology and screening methodology improvements

**1.0** — Harvard Clean Energy Project — 2006-

Distributed high-throughput + Experimental calibration + MolecularSpace database

**2.0** — Organic flow batteries — 2013-

Tight experimental feedback cycle + Stability screening + Synthesizability at forefront

**3.0** — Organic Light Emitting diodes — 2014-

Super-flexible molecular builder + Machine learning screening + Infinite libraries

**4.0** — Harvard Clean Energy Project — 2014-

Thompson sampling + Several experimental collaborations at once + Molecular crystal ML and prediction

# Harvard Clean Energy Project:Organic solar cells



Oxygen

Donor

Acceptor (PCBM)

~ -4.0 eV
LUMO

LUMO = -4.3

gap

≤ -5.2 eV
HOMO

HOMO
-5.2 eV

Cathode

Anode

PCBM= Functionalized $C_{60}$

LUMO

HOMO

# The Harvard Clean Energy Project

Part IV: Clean Energy Project

30,000+ CPU years have led to more than 35,000 high-performance organic photovoltaic candidates.

Collaborators: Juan Hindo (IBM) Zhenan Bao (Stanford), Johannes Hachmann (Buffalo), Alejandro Briseño (UMass), Carlos Amador (UNAM), and …

Hachmann, et. Al, J. Phys. Chem Lett. (2011), Energ. Env. Sci. (2014)

# Energy Levels and Efficiency



Scharber M, Mühlbacher D, Koppe M, et al., Advanced Materials (2006)

# Sifting through 2.3 million molecules

BP86/SVP//PBE0/TZVP

10%

~35000 molecules
**(1.5% of sample space)**

Energy and Environmental Science, 7, 698 (2014)

# The Clean Energy Project gets an artificial intelligence boost!

4.0

Machine learning



Bayesian calibration

Calculate

Learn

Prioritize

Easy to synthesize libraries



Smart Screening Using machine learning



Neural Fingerprints

E. O. Pyzer-Knapp, et al. Advanced Functional Materials  2015

E. O. Pyzer-Knapp, et al. arXiV:1510.00388 2015

D. Duvenaud, arXiV:1509.09292 NIPS 2015

Collaborator: Ryan Adams (Harvard)

# Organic Flow Batteries

**The Harvard Clean Energy Project**
Generating renewable energy

**Blue Organic LED**
For your next gadget or TV

**Organic flow batteries**
Storing renewable energy

**Chemical networks**
Origins of life
Organic reactions
Chemical autoencoders

# Organics for storing clean energy

Organic flow batteries



Suh, et al., Chem. Sci., 6, 2015, p. 885

Huskinson, et al., Nature, 505, 2014, p. 195

Lin, et al., Science, 349, 2015, p. 1529

Collaborators: Mike Aziz and Roy Gordon (Harvard)

# Search space for redox potentials



Oxidations | Reductions

Aromatics
Amines
Hydrocarbon
Amides
Amines
Phenols
Quinones
Quinolines
Catechols
Halogens

Olefins
Esters
Ketones
Aldehydes
Conjugated Esters
Ethers
Disulfides
Nitros
Diazos

+2.0    +1.0    0    -1.0    -2.0

E (V vs. Saturated Calomel Electrode)
Estimated potential range of organic functional group@ 25°C

*Handbook of Electrochemistry Ed. C.G. Zoski*

Choice for combinatorial library:
1R and fully substituted cases only

1. N(CH$_3$)$_2$
2. NH$_2$
3. OCH$_3$
4. OH
5. SH
6. CH$_3$
7. SiH$_3$
8. F
9. Cl
10. C$_2$H$_3$
11. CHO
12. COOCH$_3$
13. CF$_3$
14. CN
15. COOH
16. PO$_3$H$_2$
17. SO$_3$H
18. NO$_2$

2.0

*1,4-BenzoQuinones*

*1,2-BenzoQuinones*

Naphtoquinones

Anthraquinones

# Theoretical calibration of quinone redox potentials



**GEN-4 model**
with inclusion of GEN-4 molecules (ADS and PDS) and NQPS

# > **300** new candidate quinones predicted



S. Er., C. Suh, M. P. Marshak, A. Aspuru-Guzik,
Chemical Science (2015)

# Our metal-free aqueous flow battery

Computational screening of
10,000 quinone molecules

**Intense design cycle**

Synthesize molecules
Test in flow battery

**Selected molecule**

# Theory-experiment collaboration



**LETTER**
doi:10.1038/nature12909

# A metal–free organic–inorganic aqueous flow battery

Brian Huskinson[1]*, Michael P. Marshak[1,2]*, Changwon Suh[2], Süleyman Er[2,3], Michael R. Gerhardt[1], Cooper J. Galvin[2], Xudong Chen[2], Alán Aspuru-Guzik[2], Roy G. Gordon[1,2] & Michael J. Aziz[1]

**Michael Aziz**
Engineering

**Roy Gordon**
Chemistry

**Alán Aspuru-Guzik**
Chemistry

Nature, 505, 2014, p. 195

# Molecular Flow Battery Data View

Blue: Stable molecule
Red: Unstable molecule

X axis: Redox Potential
Y axis: Free energy of Solvation

~ 100,000 molecules shown

Molecular Flow Battery Data View

Filtering the data view

# Molecular Flow Battery Data View

Baseball card view

**Moelcular Space Shuttle**: advanced molecular discovery platform

High-throughput materials discovery process and tools

Info Cards

Voting Interface

Detailed Tables

Bubble Plot

Web tools and critically enable partner communication and successful molecular discovery

# Feedback tool

Database-backed web system tracks:
- ~1,000,000 machine-generated molecules
- ~1,500 (8000 including oxidation, decomposition and dissociation products)

# Complex quinone redox pathways

Highly reduced

Highly oxidized

# Additional current-theory work 1: quinone stability

## Screening procedures excluding potential Michael addition



1st screening with
Michael addition
9,866 couples

22,364 couples

2nd screening:
Fewer than two R-groups
2,290 couples

3rd screening:
Fewer than two R-groups
and good solubility
($\Delta G^0_{solv} < -0.75$ eV )
2,052 couples

**Molecules most likely
to be of interest**

- *X-axis: Quinone redox potential ($E^0_1$) / Y-axis: Stability ($K_{hyd}$)*
- Warmer colors represent higher density of molecules

# Additional current-theory work 2: Second-oxidation quinones

## Screening Procedures with consideration of Michael addition



1st screening with Michael addition

762 couples

2nd screening: Fewer than two R-groups

98 couples

3rd screening: Fewer than two R-groups and good solubility ($\Delta G^0_{solv} < -0.75$ eV )

84 couples

8,252 couples

Molecules most likely to be of interest

# Molecular baseball cards including stability



| | | |
|---|---|---|
| Potential (eV) | 0.19 | ● ○ |
| Log K hyd | -4.00 | ● |
| Mike Energy (eV) | -0.17 | ◗ |
| Solvation | -1.06 | |
| Weight (amu) | 370.00 | |
| Reduced CAS | No CAS found | |
| Oxidized CAS | No CAS found | |
| Reduced SA | 3.2 | |
| Reduced IKEY | CLRJBXWOLKGZSY-UHFFFAOYSA-N | |
| view family graph | | |

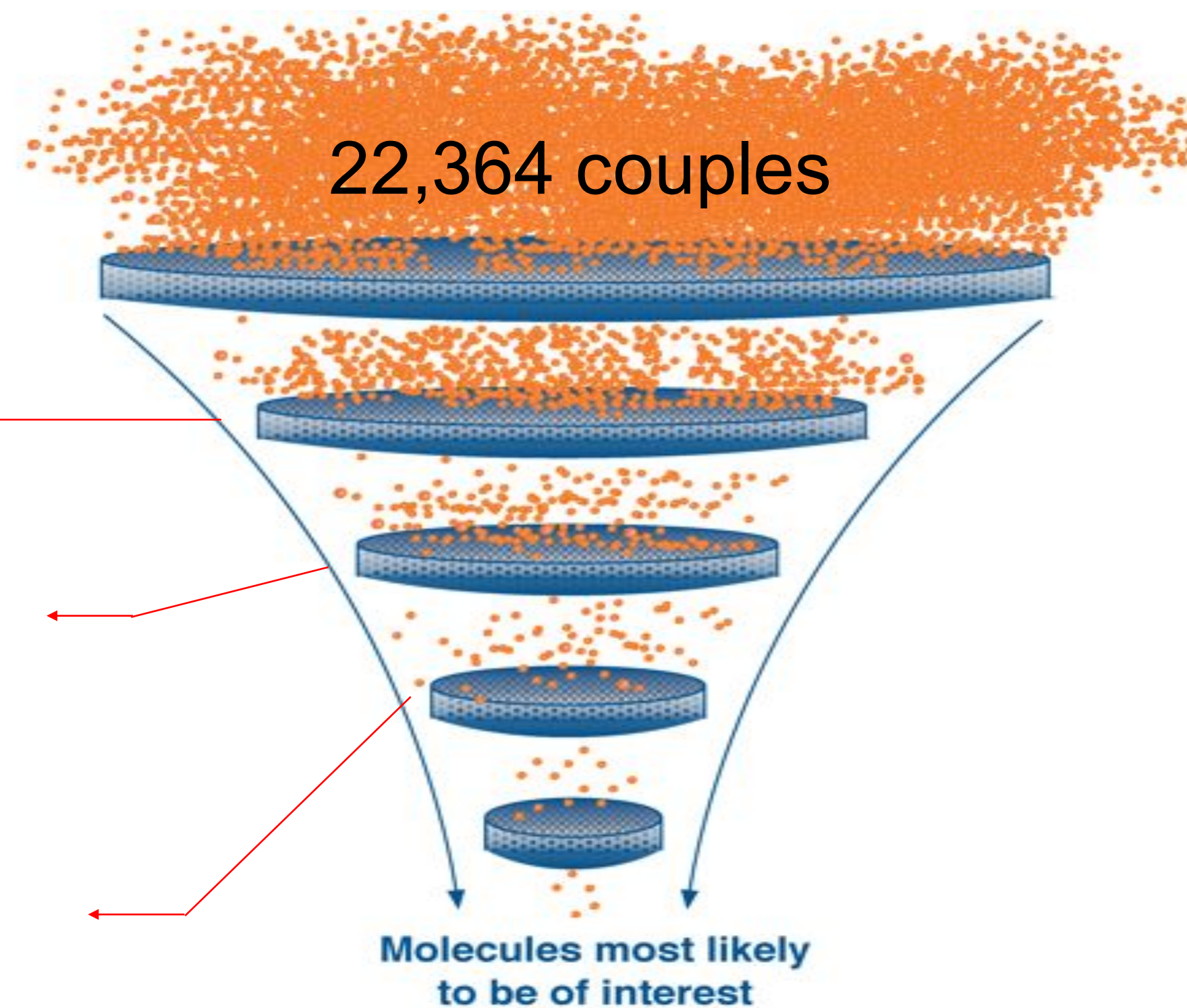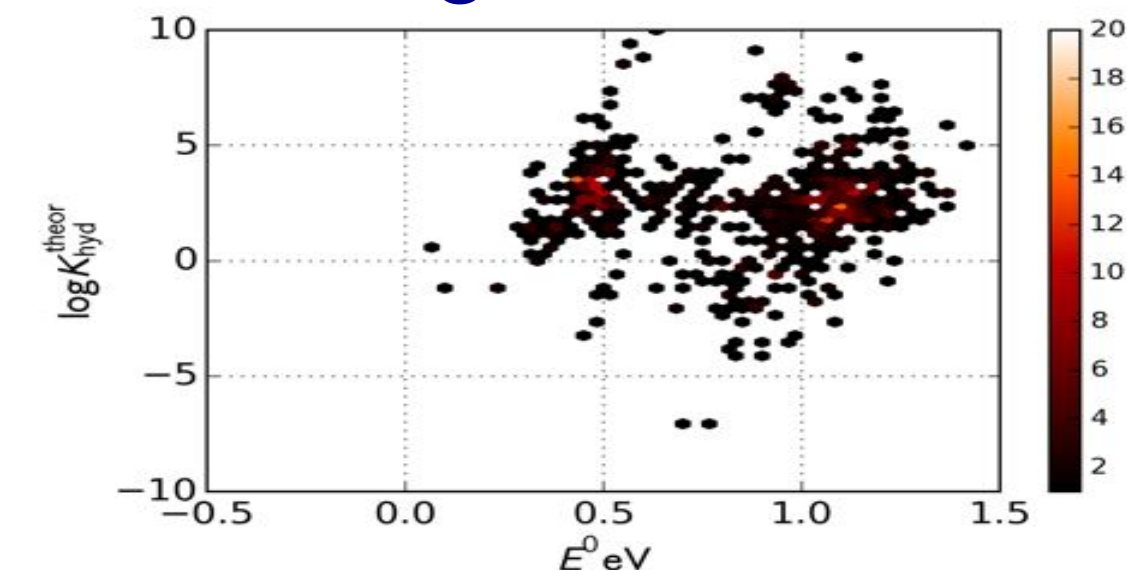| | | |
|---|---|---|
| Potential (eV) | 0.70 | ○ ○ |
| Log K hyd | 2.53 | ◗ |
| Mike Energy (eV) | -0.55 | ○ |
| Solvation | -0.84 | |
| Weight (amu) | 298.05 | |
| Reduced CAS | No CAS found | |
| Oxidized CAS | No CAS found | |
| Reduced SA | 2.4 | |
| Reduced IKEY | ZJXDHPQZYHGAJA-UHFFFAOYSA-N | |
| view family graph | | |

| | | |
|---|---|---|
| Potential (eV) | 0.20 | ● ○ |
| Log K hyd | -2.76 | ● |
| Mike Energy (eV) | -0.00 | ● |
| Solvation | -1.05 | |
| Weight (amu) | 370.00 | |
| Reduced CAS | No CAS found | |
| Oxidized CAS | No CAS found | |
| Reduced SA | 3.3 | |
| Reduced IKEY | LGNXOEAQMPUSPK-UHFFFAOYSA-N | |
| view family graph | | |

# Beyond quinones



*Sulfolobus* archaebacteria

Pineda-Flores, et al. J. Phys. Chem. C 119 21800 (2015)

# Long-lasting blue organic LED

**The Harvard Clean Energy Project**
Generating renewable energy

**Blue Organic LED**
For your next gadget or TV

**Organic flow batteries**
Storing renewable energy

**Origins of life**
How life may have come about?

# Harvard-MIT collaboration

Harvard



Ryan Adams
**Machine Learning**

Alán Aspuru-Guzik
**High-throughput
quantum chemistry**

MIT



Tim Swager
Stephen Buchwald
**Synthetic Chemistry**
Marc Baldo
**Device Engineering**
Troy Van Voorhis
**Microscopic theory**

>450,000 molecules screened so far!  ~25 synthesized and tested

# Speedy screening: Machine Learning

**1** Sample
Sample random molecules in the DB

**2** Calculate
Run a few thousand calculations on those molecules

**3** Train
Train ML system on calcualtion results

**4** Prioritize
Use predictions to rank candidates and prioritized the calculations

Enhanced discovery rate

# Machine Learning

- Supervised learning algorithms
  - Neural networks for ultrafast predictions leveraging thousands of data-points.
  - Result in 10x speedup by discarding poor candidates

- Role of dimensionality
  - Chemical space is sparse but libraries are dense. Powerful interpolation

- Explore-exploit strategy

Selecting molecules is like dating.

tinder

It's how people meet

Download the App

▶ Watch **Tinder Plus**

Organic LED Screening

Synthetic accessibility voting tool

# Neural Net Training Workflow

# Data mining 500,000 quantum calculations



R. Gomez-Bombarelli, et al. Submitted (2015)

# Batches

## Batches

Selection of 100-200 molecules for experimentalists to browse in a contained way.

Usually explore some chemical *family*, using ancestry from database.

Need to confirm novelty *post hoc*: sometimes re-discover known molecules.

22.5%

Rafael Gómez-Bombarelli, Jorge Aguilera-Iparaguirre, Tim Hirzel
Martin BloodAdams, Baldo, Swager groups, Samsung IT

# Key breakthroughs in efficiency: Strength

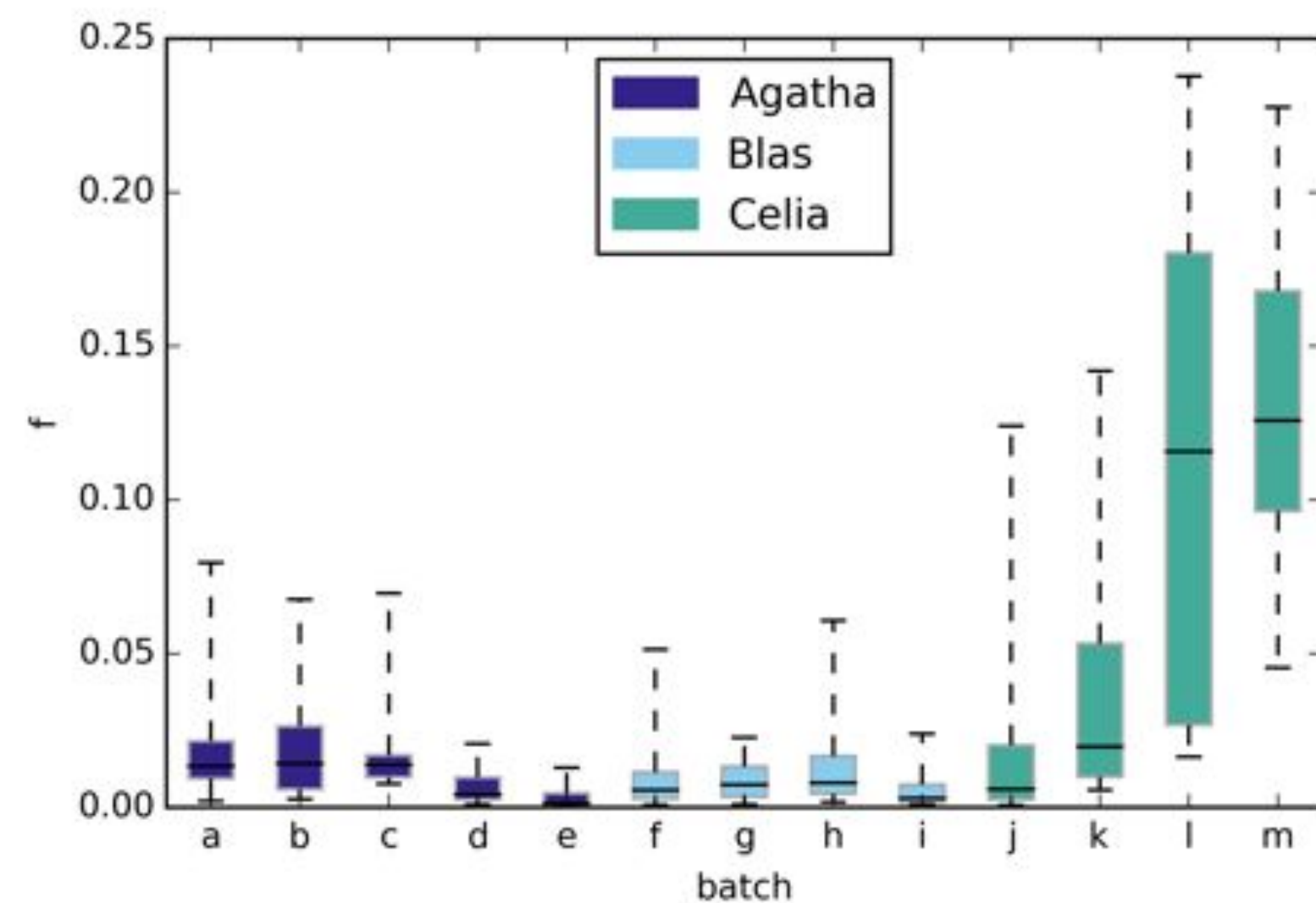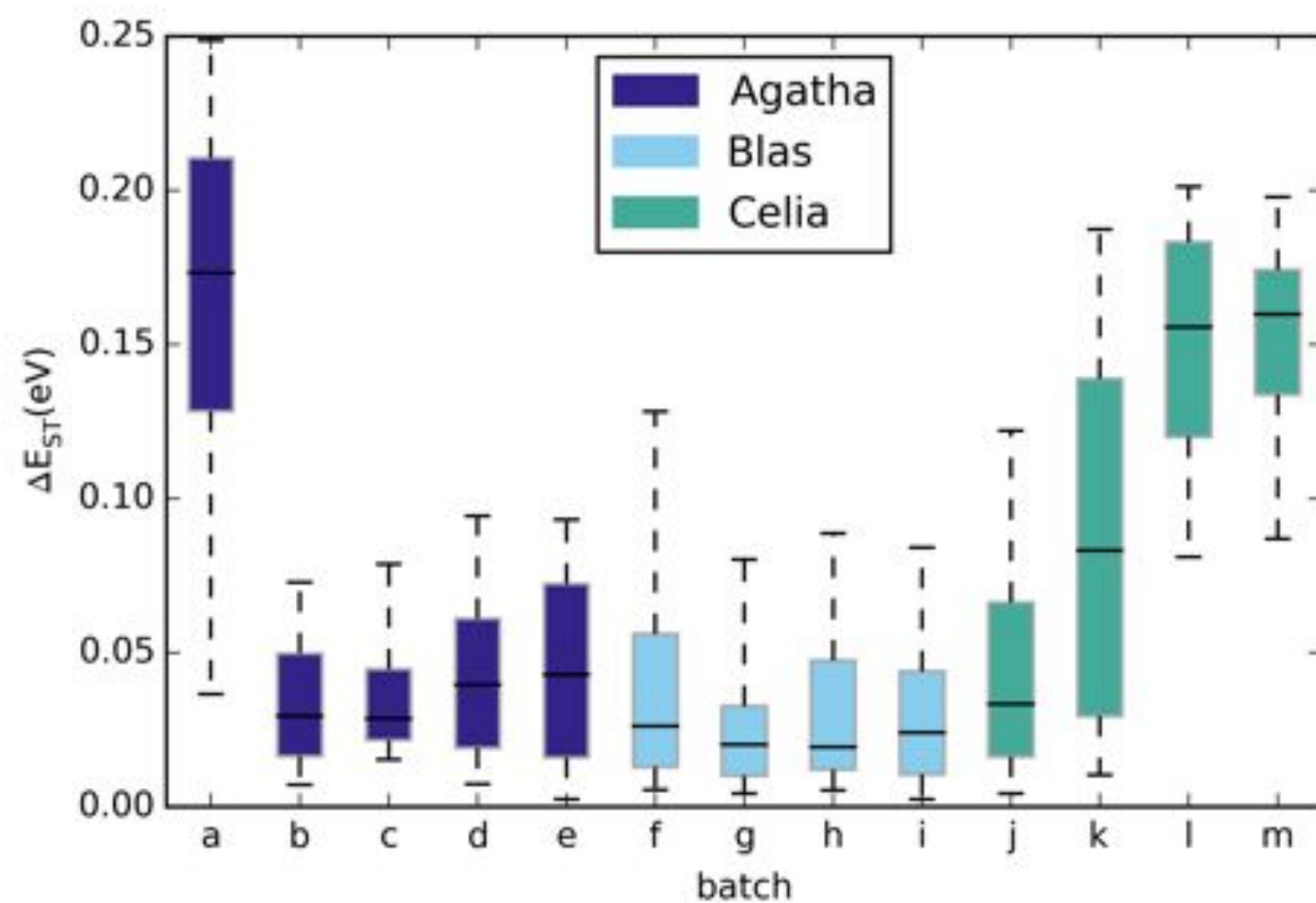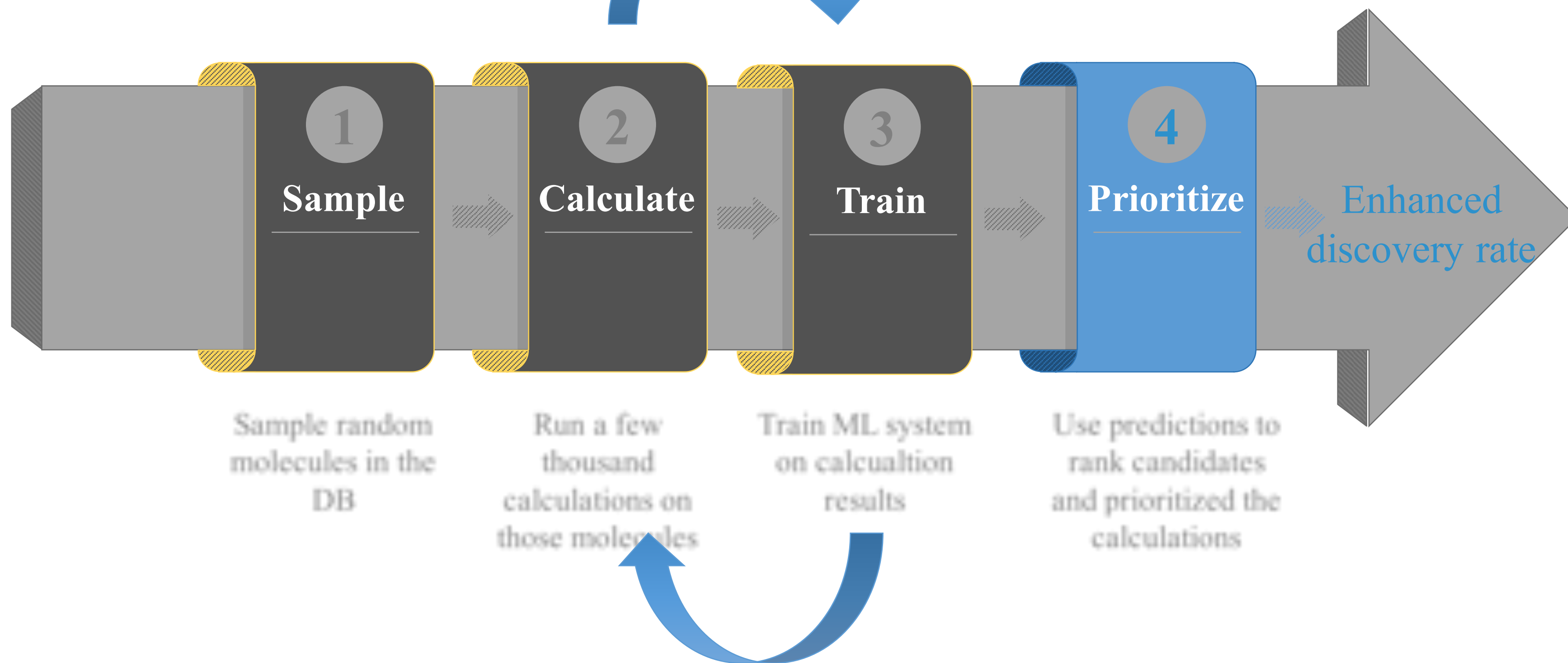| Name | $S_0$ splitting | $T_1$ splitting | $S_0$ strength | $T_1$ strength | EQE(%) |
|---|---|---|---|---|---|
| 4CzIPN | 0.124 | 0.101 | 0.063 | 0.049 | 20 |
| Foxtrot1-21 | 0.015 | 0.031 | 0.003 | 0.000 | 20 |
| Hotel1-38 | 0.017 | 0.046 | 0.008 | 0.012 | 7 |
| Julie2-16-1 | 0.104 | 0.145 | 0.124 | 0.186 | 22 |
| Lima17-36 | 0.179 | 0.187 | 0.257 | 0.240 | 17* |

- A small gap is crucial for TADF behavior

- We need also need a big fluorescence

- We have managed to control both for great overall efficiency

# Screening billions of molecules: Machine learning takes the driver's seat

To design something really well you have to get it. You have to really grok what it's all about. It takes a passionate commitment to really thoroughly understand something. Chew it up, not quickly swallow it. Most people don't take time to do that.

# Aspuru-Guzik group

http://aspuru.chem.harvard.edu
Twitter: A_Aspuru_Guzik
aspuru@chemistry.harvard.edu