Argonne
NATIONAL LABORATORY

# Active Learning via Bayesian Optimization for Discovery of Energy Storage Materials

cm CHEMISTRY OF MATERIALS
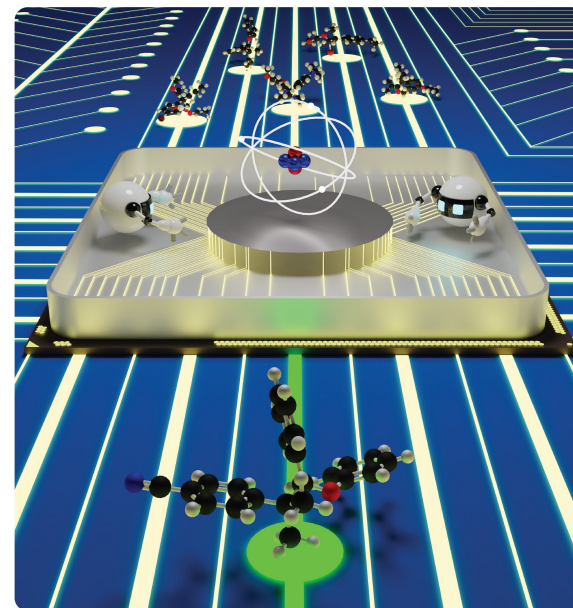
AUGUST 11, 2020 | VOLUME 32 | NUMBER 15 | pubs.acs.org/cm

Hieu A. Doan

Garvit Agarwal

Molecular Materials Group
Materials Science Division

Doan, Agarwal, Qian, Counihan, Rodríguez-López, Moore, & Assary.
(2020). https://doi.org/10.1021/acs.chemmater.0c00768

# Redox Flow Battery (RFB) as a Stationary Energy Storage System



Kowalski, Su, Milshtein, Brushett (2016) *Current Opinion in Chemical Engineering*
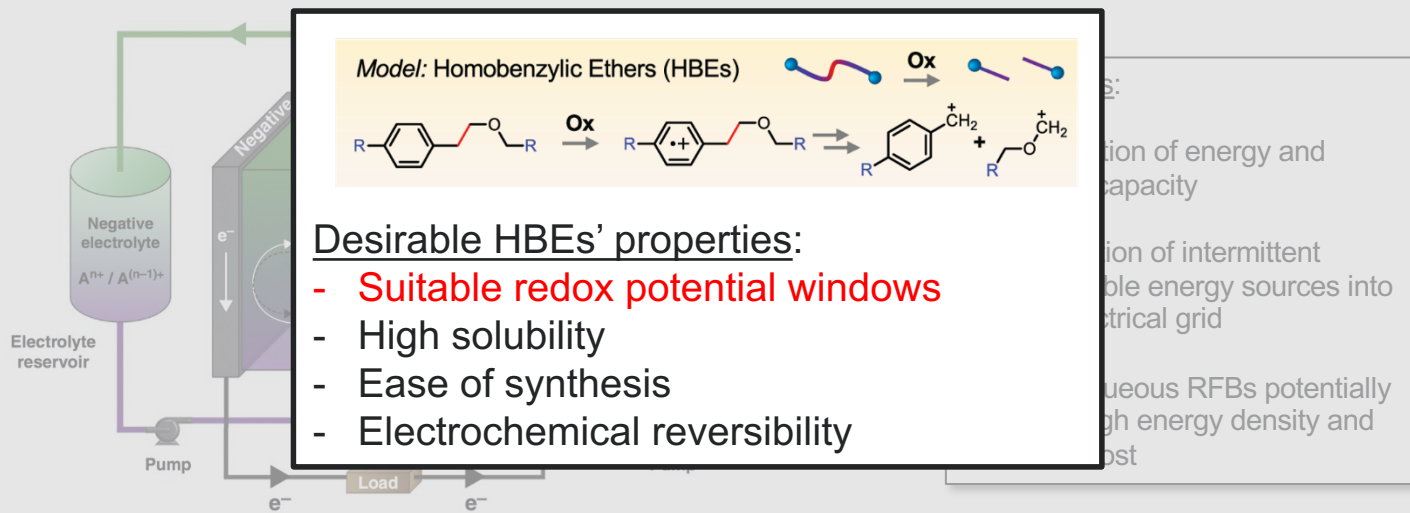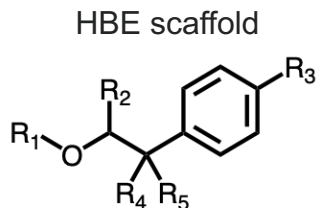
Advantages:

- ❑ Separation of energy and power capacity

- ❑ Integration of intermittent renewable energy sources into the electrical grid

- ❑ Non-aqueous RFBs potentially yield high energy density and lower cost

# Redox Flow Battery (RFB) as a Stationary Energy Storage System



Negative electrolyte

$A^{n+}$ / $A^{(n-1)+}$

Electrolyte reservoir

Pump

Load

e⁻        e⁻

**Model:** Homobenzylic Ethers (HBEs)

Desirable HBEs' properties:
- <span style="color:red">Suitable redox potential windows</span>
- High solubility
- Ease of synthesis
- Electrochemical reversibility

...tion of energy and ...capacity

...tion of intermittent ...ble energy sources into ...trical grid

...ueous RFBs potentially ...gh energy density and ...st

Kowalski, Su, Milshtein, Brushett (2016) *Current Opinion in Chemical Engineering*

Argonne NATIONAL LABORATORY

# The Challenges

## 1. Large space of molecular candidates

HBE scaffold



~

SMILES (Simplified Molecular-Input Line Entry System)

[**R3**]C1=CC=C(C([**R4**])([**R5**])C([**R2**])-[O][**R1**])C=C1

$R_1$ = -Me, -Et, -Pr, -Ph, -CN, -Eth, -COMe, -C(Me)Me, -CCOMe
$R_2$ = -N(Me)$_3^+$, -COMe, -Et, -OCMe, N(Me)$_2$, -NO$_2$, -C(=O), -Pr, -C(Me)Me, -CCOMe, -C(Me)OMe
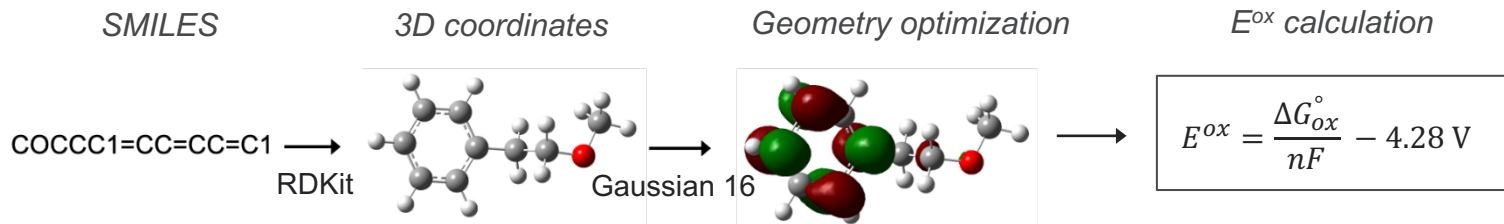$R_3$ = -N(Me)$_3^+$, -Me, -COMe, -Br, -C(=O), -OEth, -Pr, -C(Me)Me, -C(COMe), -C(Me)(OMe)
$R_4$ - $R_5$ = -N(Me)$_3^+$, -OMe, -COMe, -Br, -N(Me)Me, -Et, -OEt, -NO$_2$, -C(=O), -Pr, -C(Me)Me, -C(COMe), -C(Me)(OMe)

$> 10^5$ molecules

## 2. Expensive/time-consuming synthesis and characterization

Argonne NATIONAL LABORATORY

# The Opportunities

*1. Density Functional Theory (DFT) calculations starting from SMILES*



*SMILES*     *3D coordinates*     *Geometry optimization*     *$E^{ox}$ calculation*

COCCC1=CC=CC=C1 $\longrightarrow$ RDKit     Gaussian 16     $\longrightarrow$

$$E^{ox} = \frac{\Delta G^{\circ}_{ox}}{nF} - 4.28 \text{ V}$$

*2. Train machine learning (ML) models using DFT-computed $E^{ox}$*

Need ML models to not only make accurate predictions but also guide the selection of training data

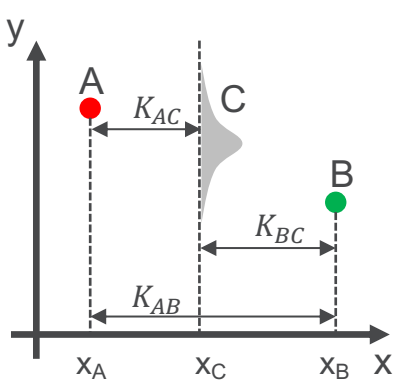$\longrightarrow$

Active learning/
Bayesian optimization

(Surrogate model + Acquisition function)

Argonne
NATIONAL LABORATORY

# Gaussian Process Regression (GPR) as a surrogate model

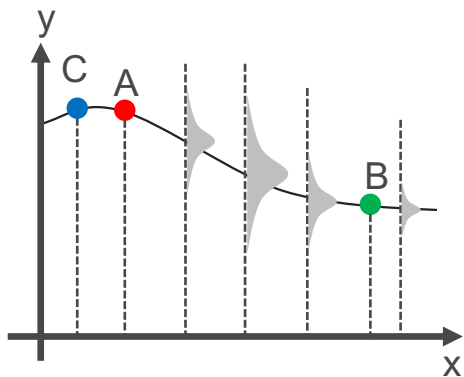<u>In a nutshell</u>: Predict properties/outputs based on feature/input differences (distances)

**Covariance** is calculated as a function of (feature) distances, e.g. $K(x_1, x_2) = exp\left[-\frac{1}{2}\frac{(x_1-x_2)^2}{l^2}\right]$
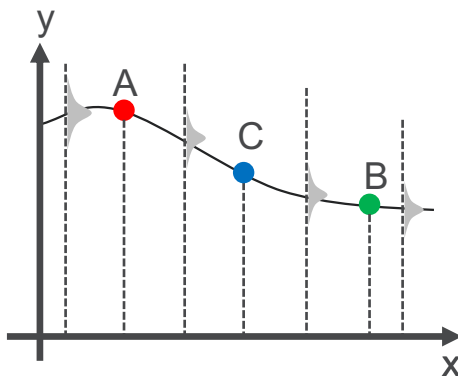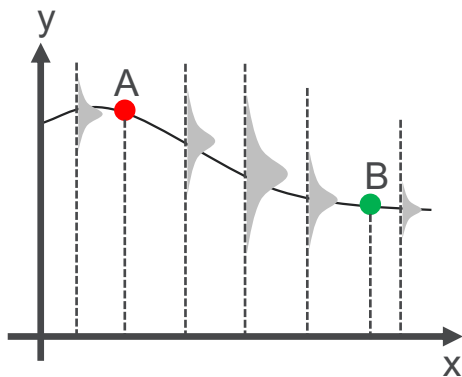
U.S. DEPARTMENT OF **ENERGY**  Argonne National Laboratory is a
U.S. Department of Energy laboratory
managed by UChicago Argonne, LLC.

Argonne
NATIONAL LABORATORY

# Use of GPR in Active learning/Bayesian optimization

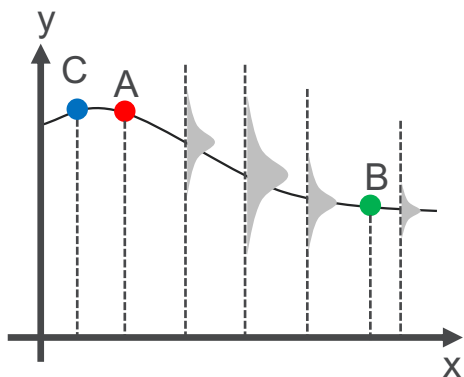GPR-predicted μ(x) and σ(x) enables **Active Learning/Bayesian optimization**
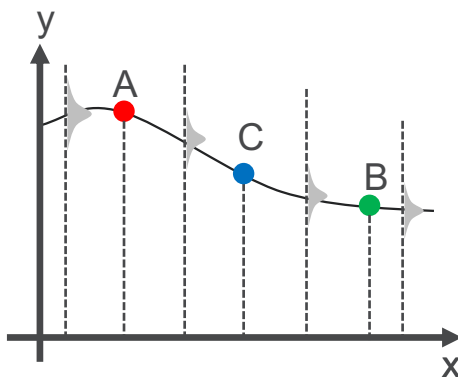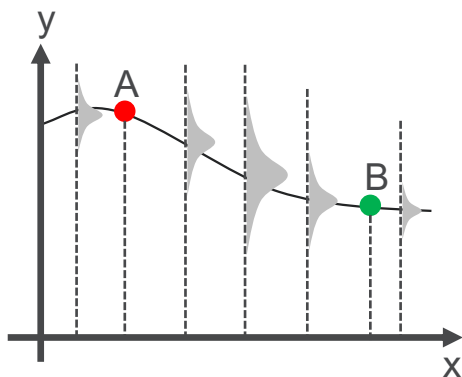


"Exploitation": follow $\mu(x)_{max}$

"Exploration": follow $\sigma(x)_{max}$

# Acquisition function

Acquisition function formulates an optimal strategy toward an objective by guiding the next evaluation



"Exploitation": follow $\mu(x)_{max}$
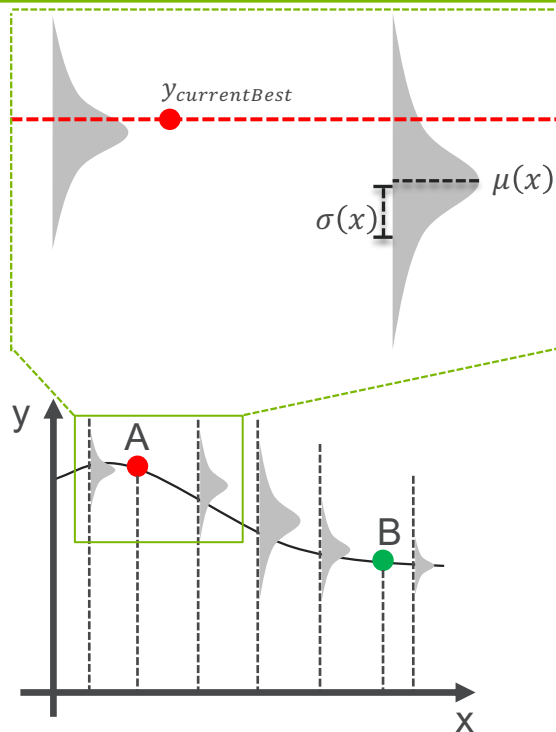
"Exploration": follow $\sigma(x)_{max}$

# Upper Confidence Bound (UCB)

Upper Confidence Bound (UCB)

$$UCB(x) = \mu(x) + \xi * \sigma(x)$$

$$x_{next} = argmax(UCB(x))$$



$\mu(x)$: Mean
$\sigma(x)$: Standard deviation

# Probability of improvement (PI) and Expected improvement (EI)

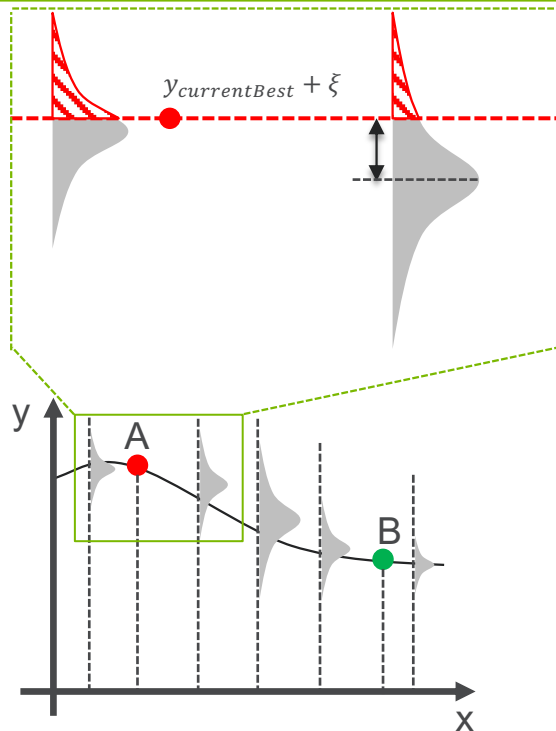Probability of Improvement (PI)

$$PI(x) = P\left(y'_x \geq y_{currentBest} + \xi\right)$$

$$= \Phi(Z)$$

$$Z = \frac{\mu(x) - y_{currentBest} - \xi}{\sigma(x)}$$

$$x_{next} = argmax(PI(x))$$

Expected Improvement(EI)

$$EI(x) =$$
$$\begin{cases}(\mu(x) - y_{currentBest} - \xi)\Phi(Z) + \sigma(x)\phi(x), & \sigma(x) > 0 \\ 0, & \sigma(x) = 0\end{cases}$$

$$x_{next} = argmax(EI(x))$$

$y_{currentBest} + \xi$

A

B

y

x

$\mu(x)$: Mean
$\sigma(x)$: Standard deviation

$\Phi$: Cummulative Distribution Function
$\phi$: Probability Density Function

Argonne
NATIONAL LABORATORY

# Problem definition for Bayesian optimization

### Input

Molecule candidates
(SMILES library)

*"SMILES.csv"*

RDKit

Feature vectors

### Output/Objective

Molecules w/ desired oxidation potential

E.g.



| MW | # of C | # of aromatic rings | … |
|-----|--------|---------------------|-----|
| 179 | 11 | 1 | ... |

*"features.csv"* 125-component vectors

Argonne
NATIONAL LABORATORY

# Bayesian optimization scheme

Objective: Given 1000 HBE molecules, find one with maximized $E^{ox}$ within N evaluations

"SMILES.csv"

| HBE Library |
| --- |
| SMILES strings |

"features.csv"

| Features |
| --- |
| MW, logP, TPSA… **Principle Component Analysis** |

Initial training data

10 random SMILES with computed $E^{ox}$

*Step 1*

| Surrogate Model Training |
| --- |
| Gaussian Process Regression |

Predicted μ and Σ

Computed $E^{ox}$

Bayesian optimization

| $E^{ox}$ Calculation |
| --- |
| DFT simulation |

Next SMILES to test

| Acquisition Function Evaluation |
| --- |
| E.g.: Expected-improvement |

*Step 3*

*Step 2*

Argonne NATIONAL LABORATORY

# Launch the nanoHUB tool

From your browser go to this link: https://nanohub.org/resources/bayesopt



Click to start the notebook