

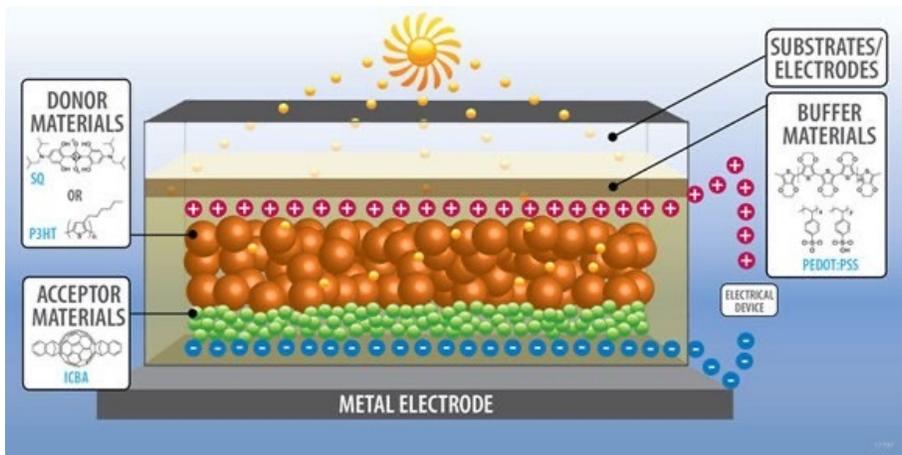
A Unified Active Learning Framework for Designing Energy-Relevant Molecules

Shomik Verma¹, Jiali Li², Kevin Greenman¹,
Rafael Gomez-Bombarelli¹, Xiaonan Wang^{2,3}, Aron Walsh⁴

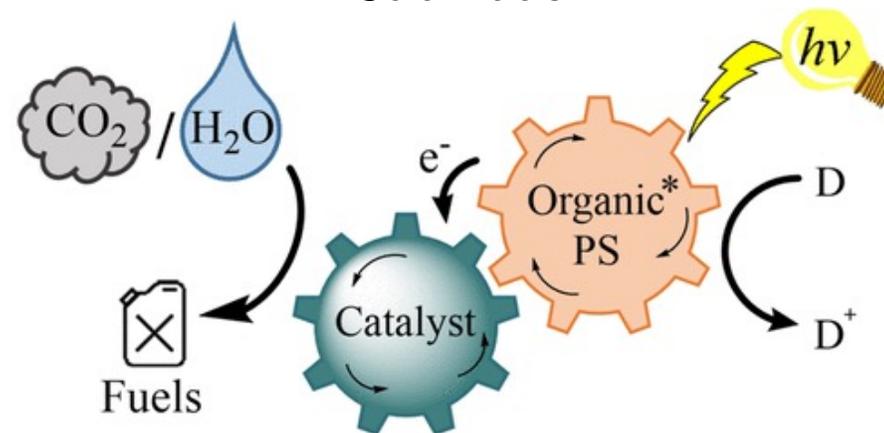
¹Massachusetts Institute of Technology, ²National University of Singapore,
³Tsinghua University, ⁴Imperial College London

Molecules relevant to energy applications

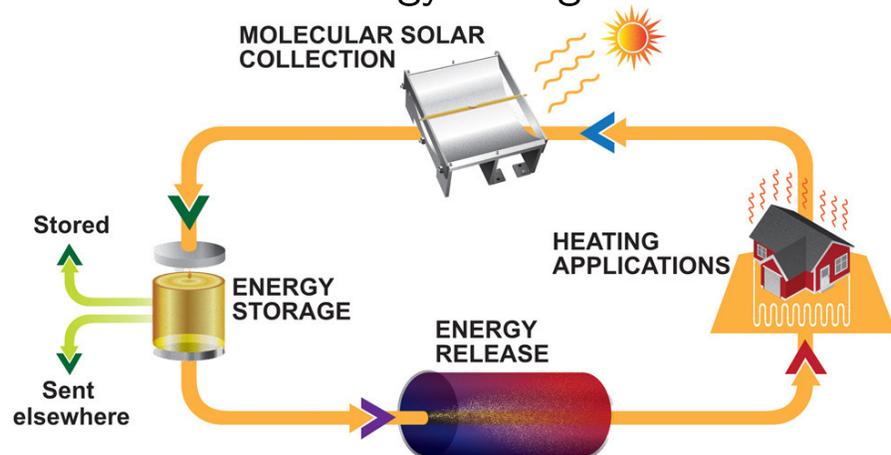
Organic photovoltaics



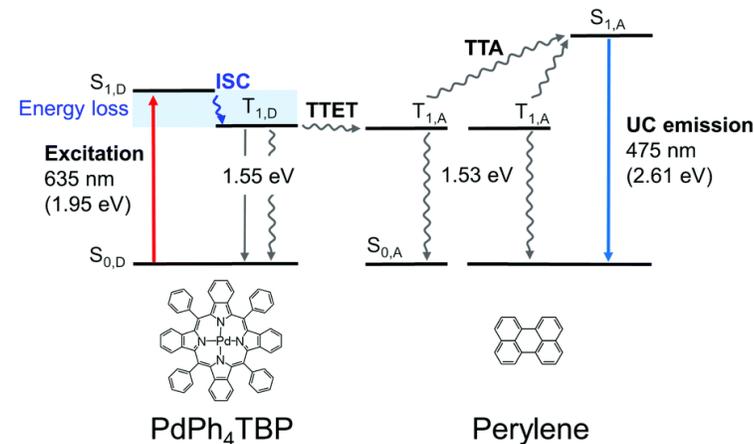
Solar fuels



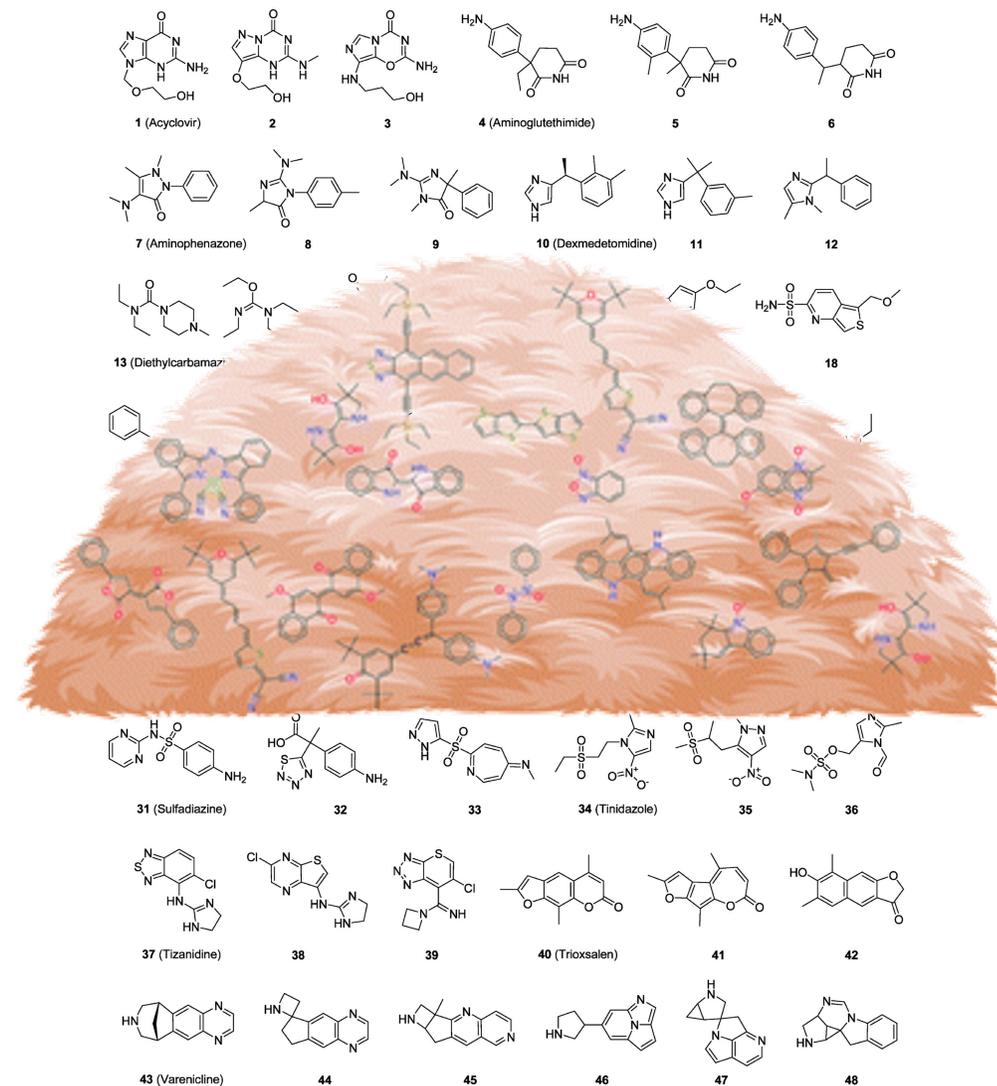
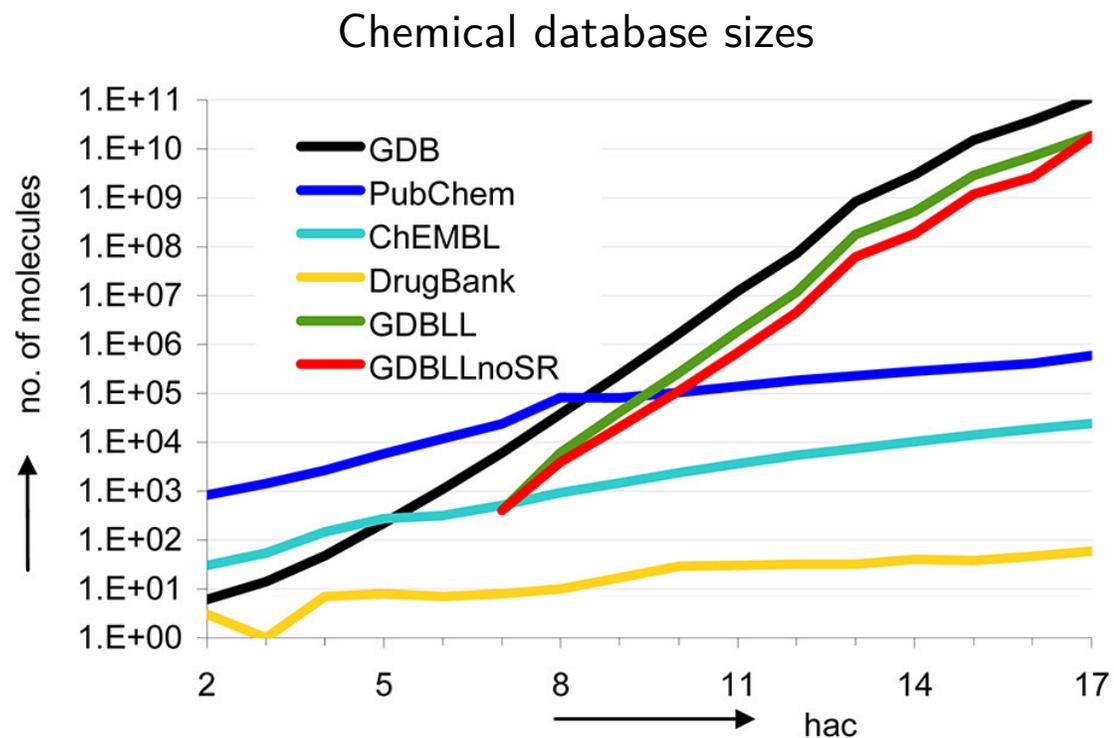
Energy storage



Photon conversion

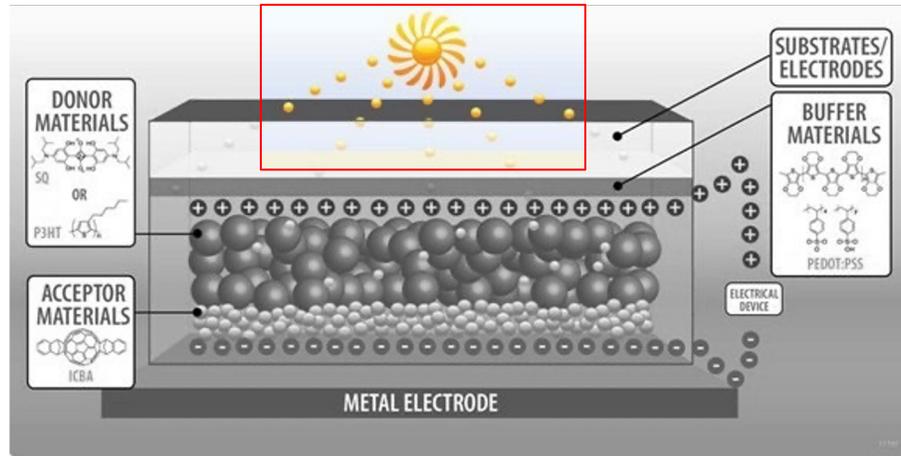


Molecular space is massive

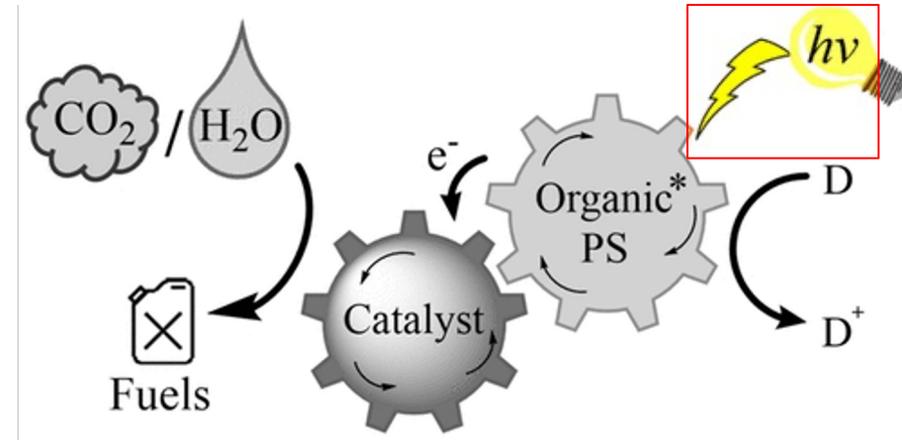


Light interactions = excited state energies required

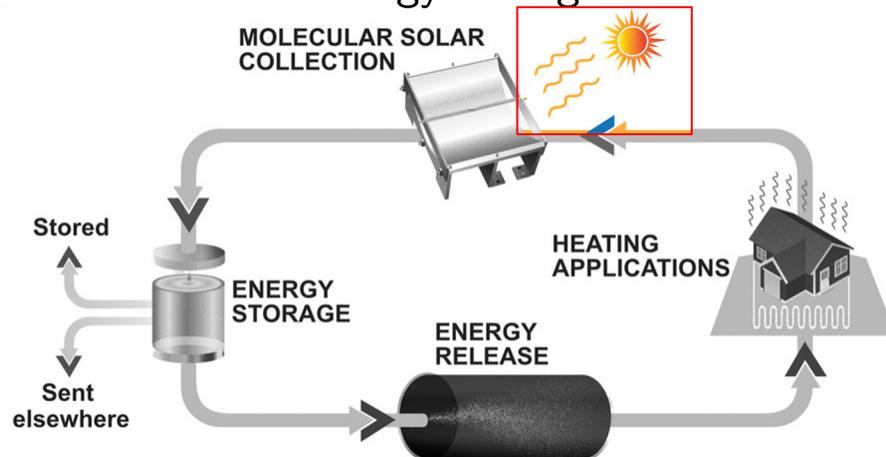
Organic photovoltaics



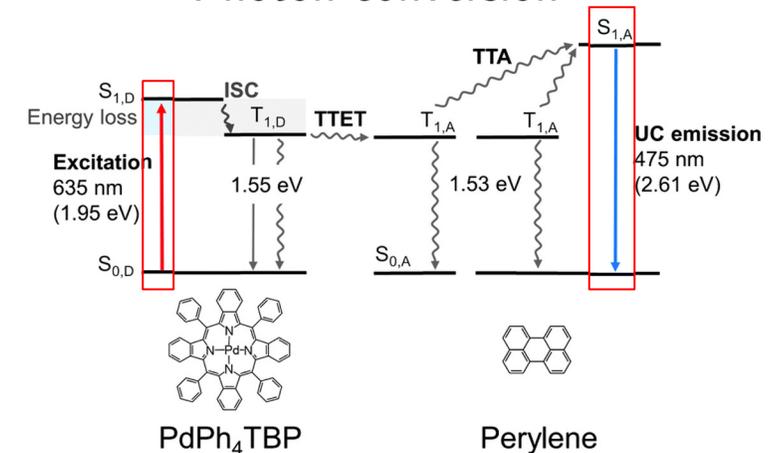
Solar fuels



Energy storage

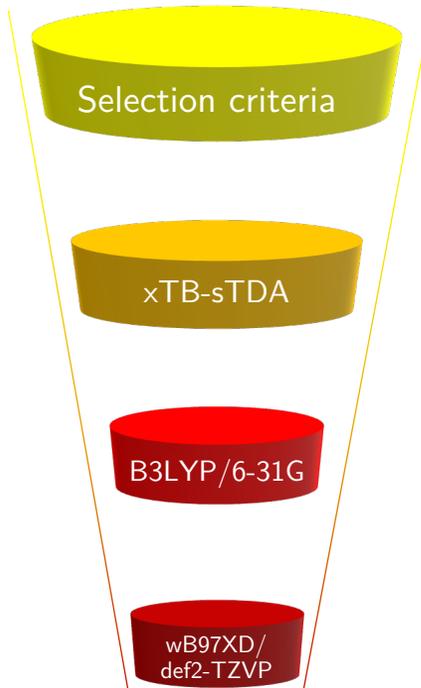


Photon conversion



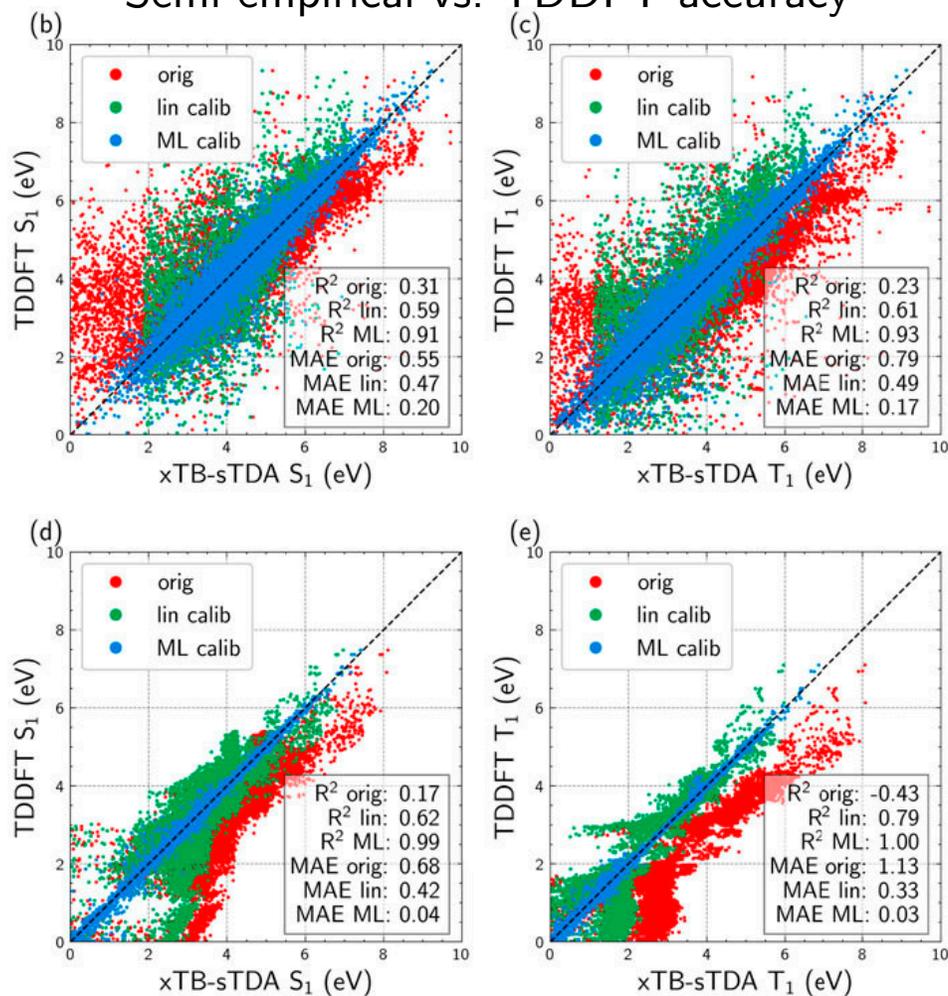
Fast and accurate chemical screening

Molecular dataset



Filtered molecules

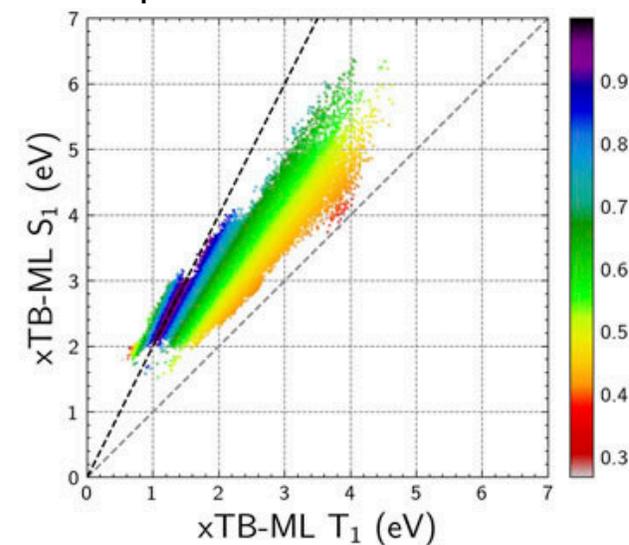
Semi-empirical vs. TDDFT accuracy



Calibration results

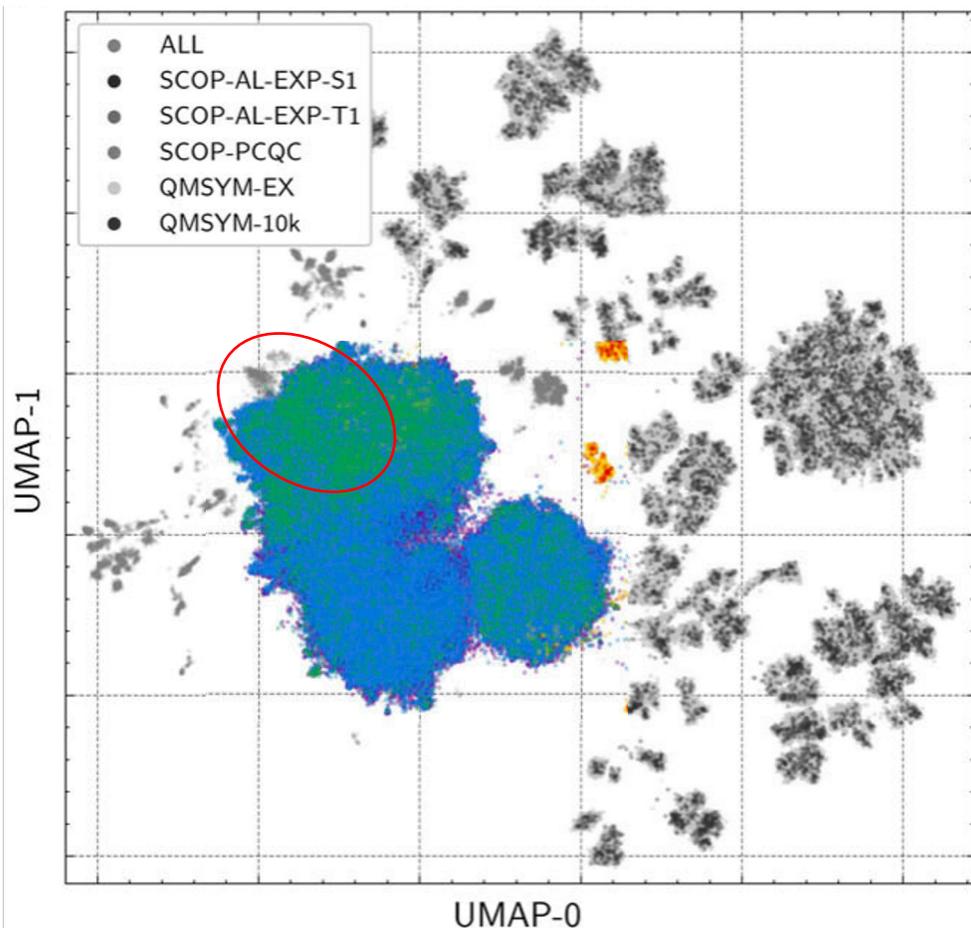
	xTB-sTDA T_1	xTB-ML-20k Class 1 T_1	xTB-ML-300k Class 1 T_1	xTB-ML-20k Class 2 T_1	xTB-ML-300k Class 2 T_1
MOPSSAM (143)	0.59	0.17	0.13	0.10	0.08
MOPSSAM (1k)	0.39	0.19	0.15	0.14	0.12
INDT (10k)	0.40	0.13	0.15	0.60	0.38
Verde (1.5k)	0.48	0.19	0.22	0.26	0.29

HTVS for photon conversion molecules



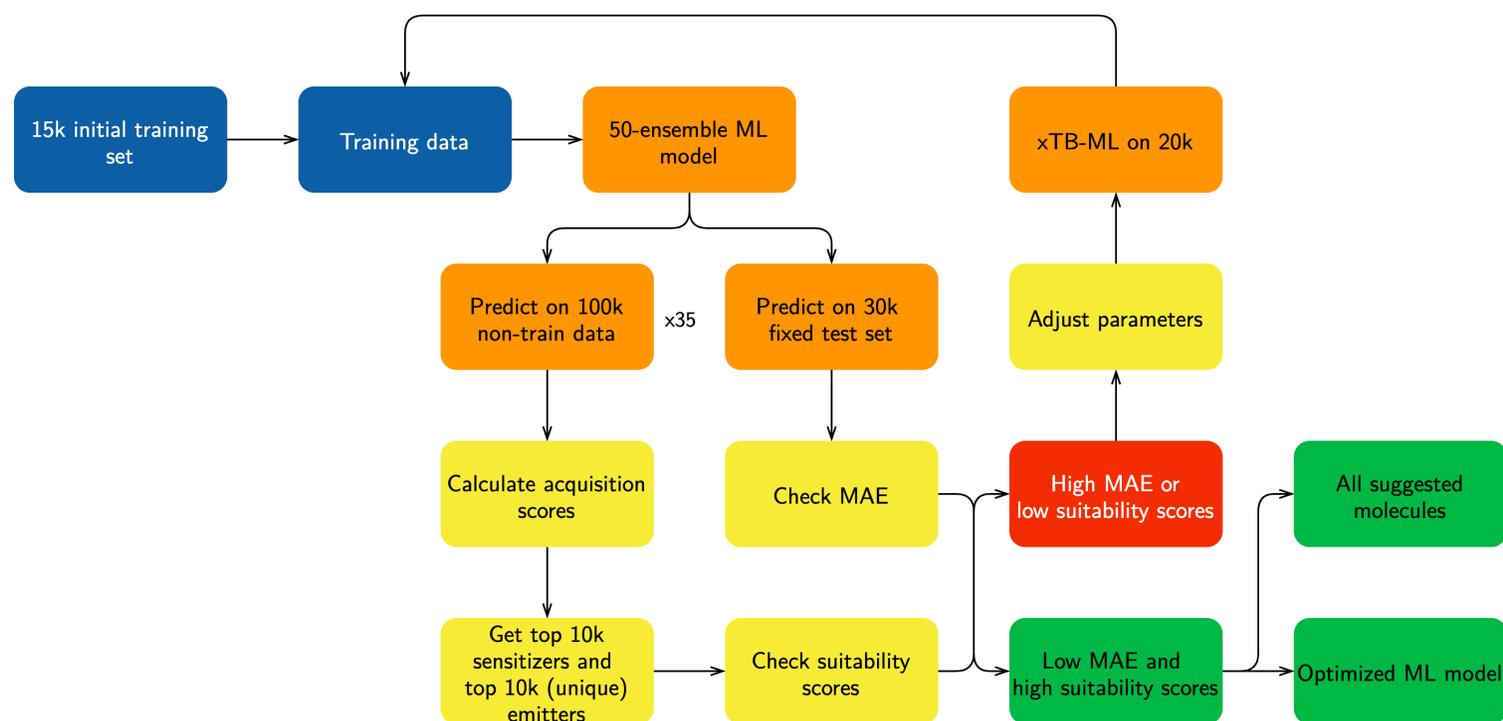
Efficient exploration with active learning

Global chemical space map



○ 250k HTVS space

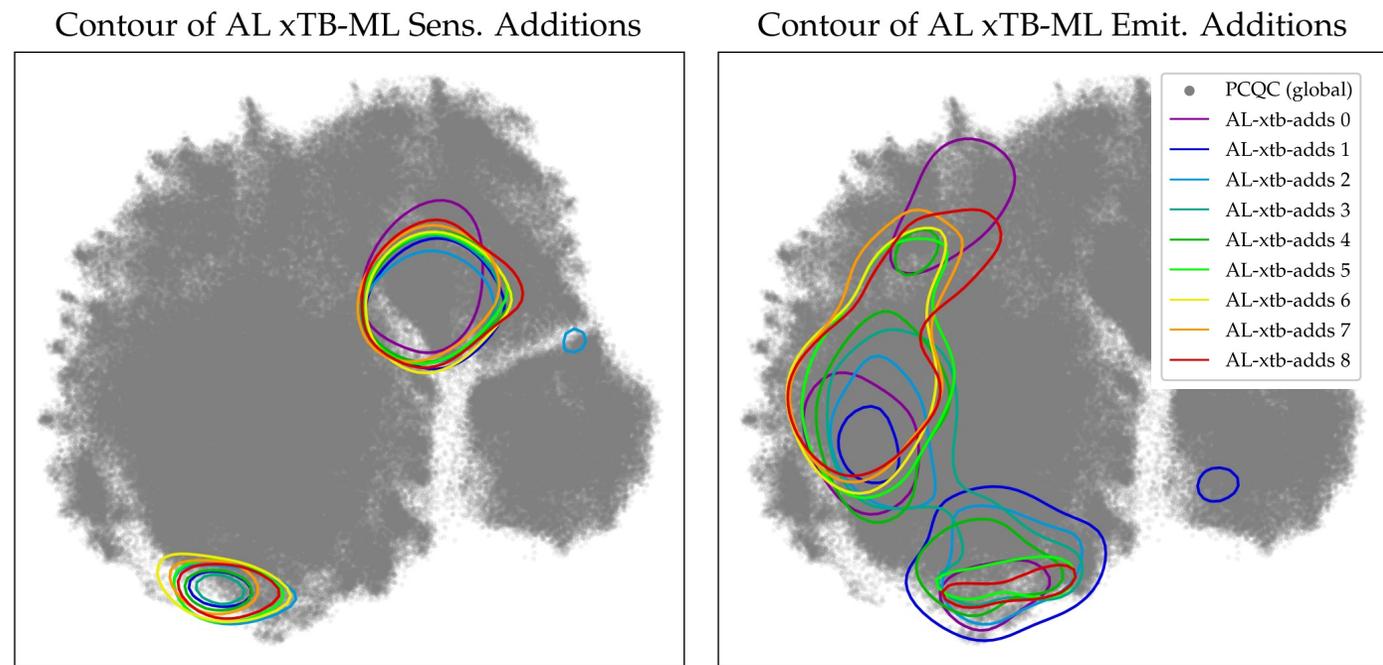
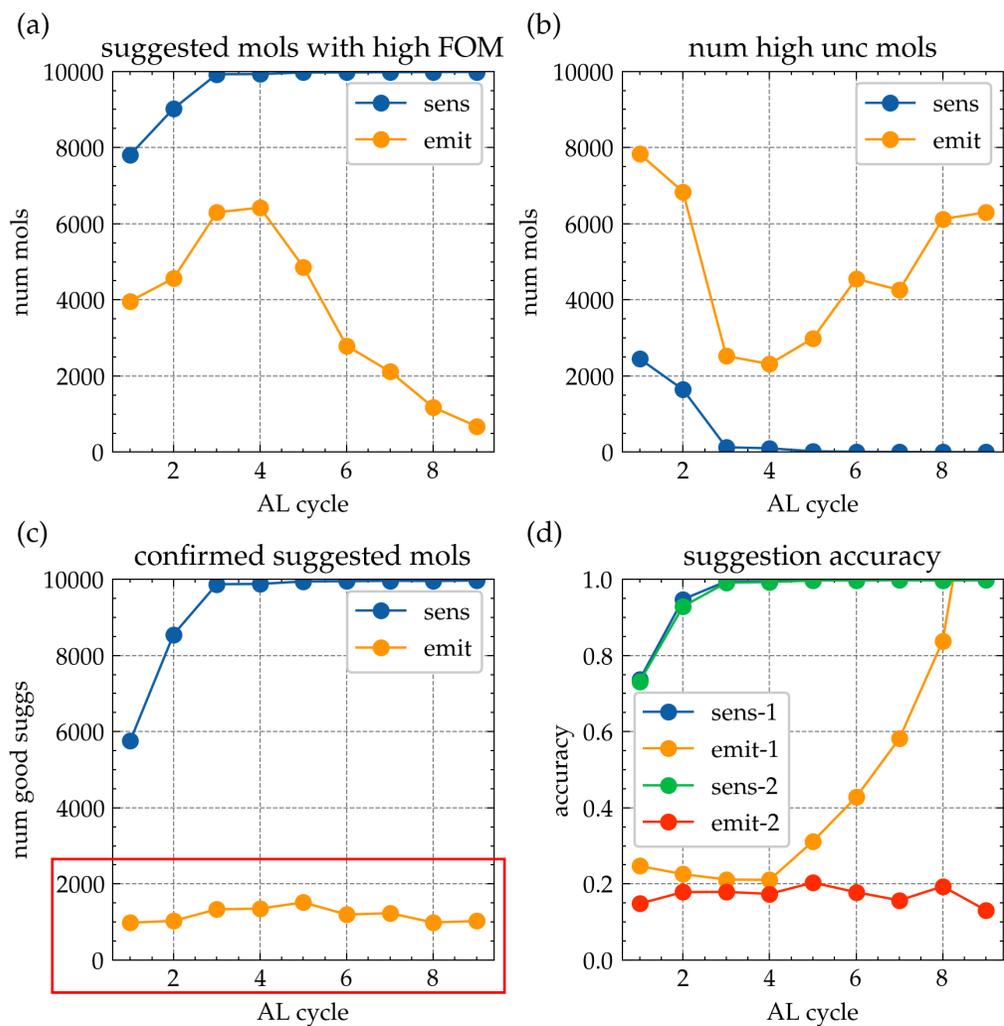
AL workflow



$$\alpha_{sens} = B \cdot \epsilon_{sens} + \sigma_{T1} + \sigma_{S1}$$

$$\alpha_{emit} = B \cdot \epsilon_{emit} + \sigma_{T1} + \sigma_{S1}$$

AL results on small molecule database



Low AL accuracy for emitters

How to improve results?

1

Expand database

2

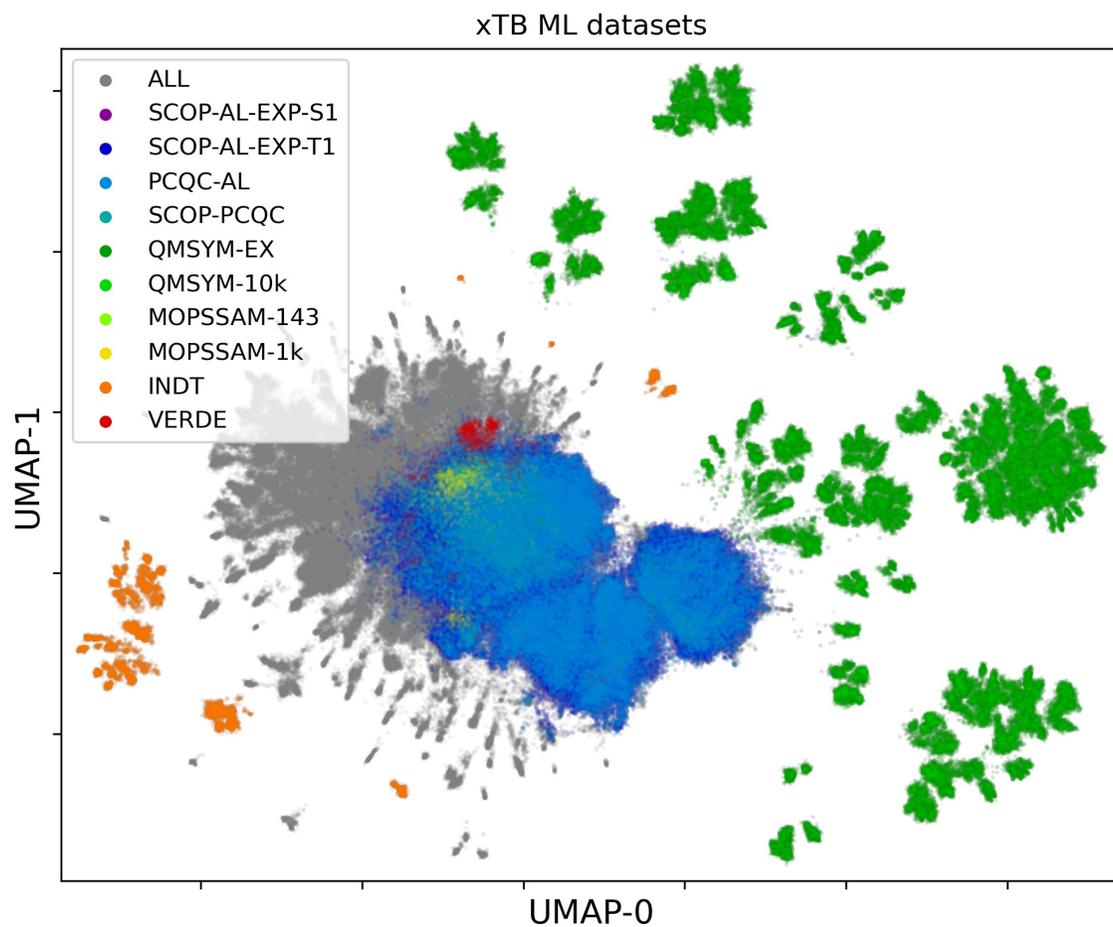
Improve AL
framework

3

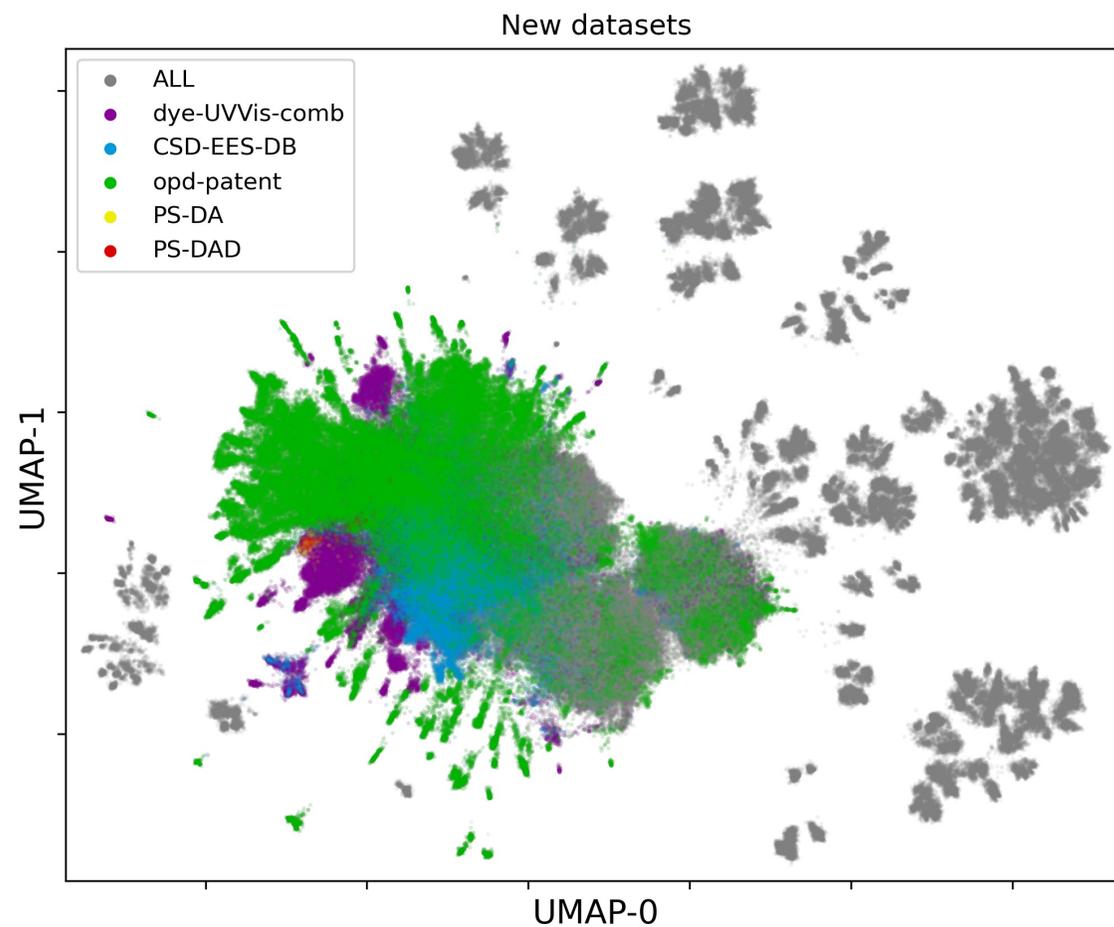
Implement
generative ML

Expanding to large molecules

(a)

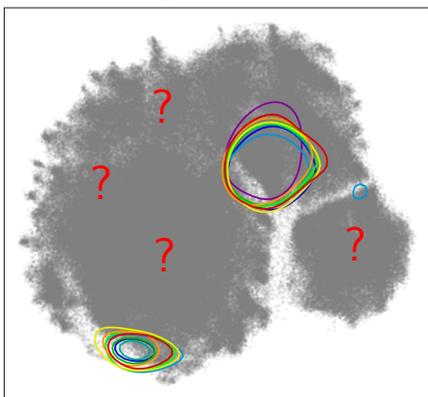


(b)



Improving the AL framework

Training data locations

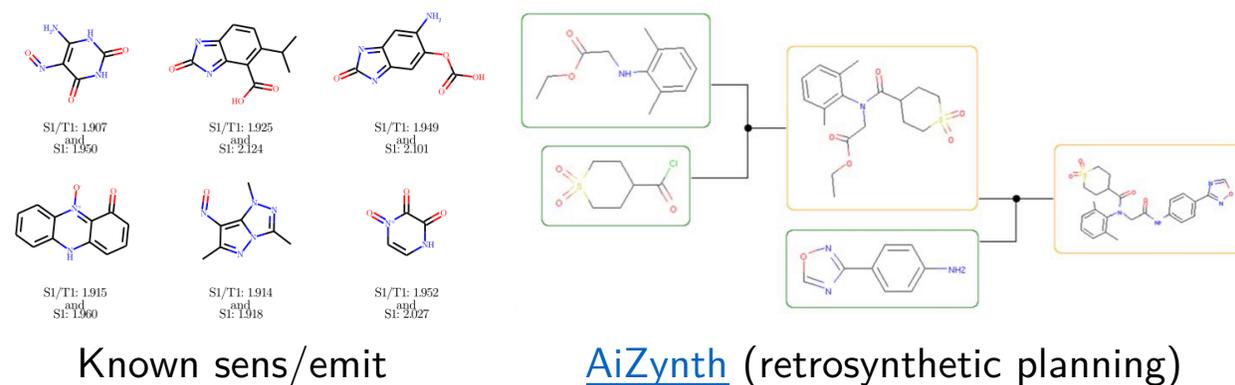


Uncertainty

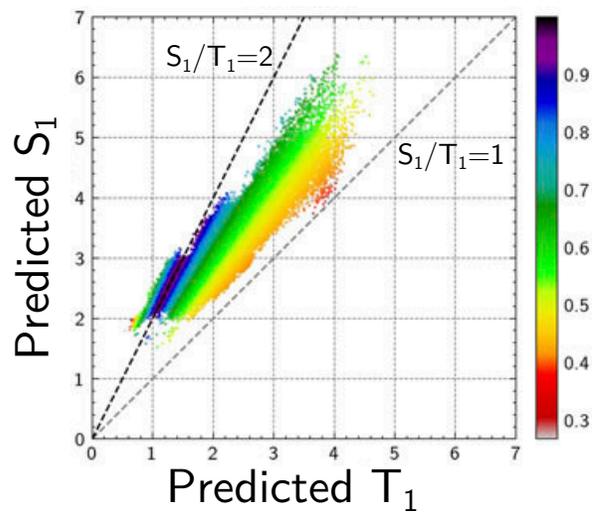
$$\sigma_{T_1} = \frac{\sum_{i=1}^n (T_{1i} - T_{1mean})^2}{n}$$

$$\sigma_{S_1} = \frac{\sum_{i=1}^n (S_{1i} - S_{1mean})^2}{n}$$

Domain Knowledge / Synthesizability



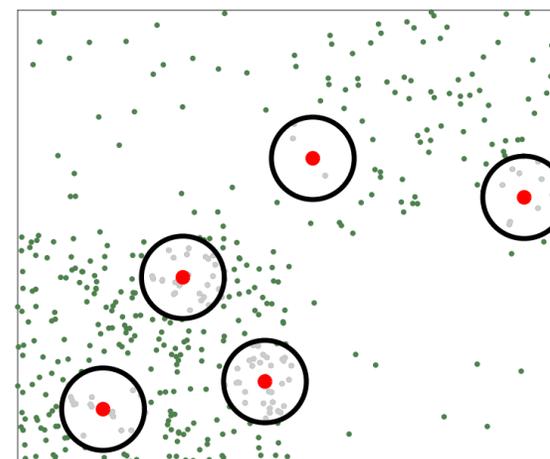
Target Property



$$\epsilon_{sens} = e^{-A\left(1 - \frac{S_1}{T_1}\right)}$$

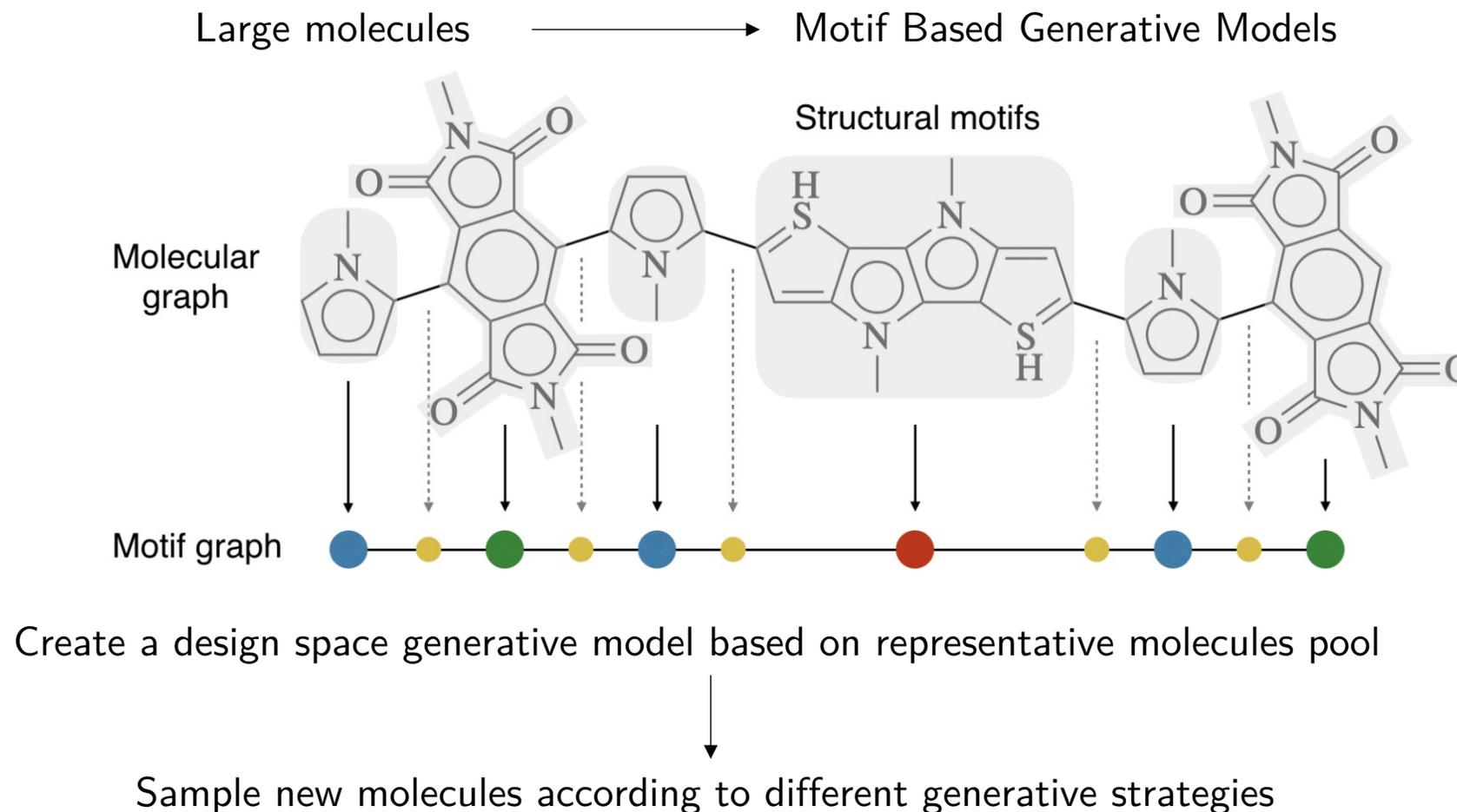
$$\epsilon_{emit} = e^{-A\left(2 - \frac{S_1}{T_1}\right)}$$

Batch diversification

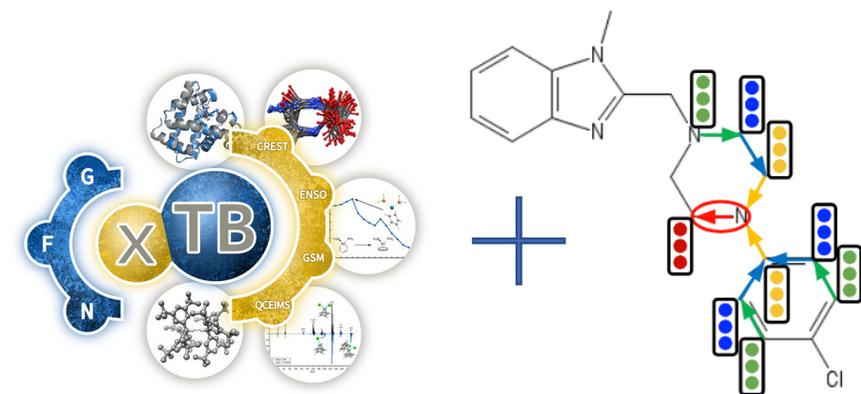


$$r_{new} - r_i > \theta, \forall i$$

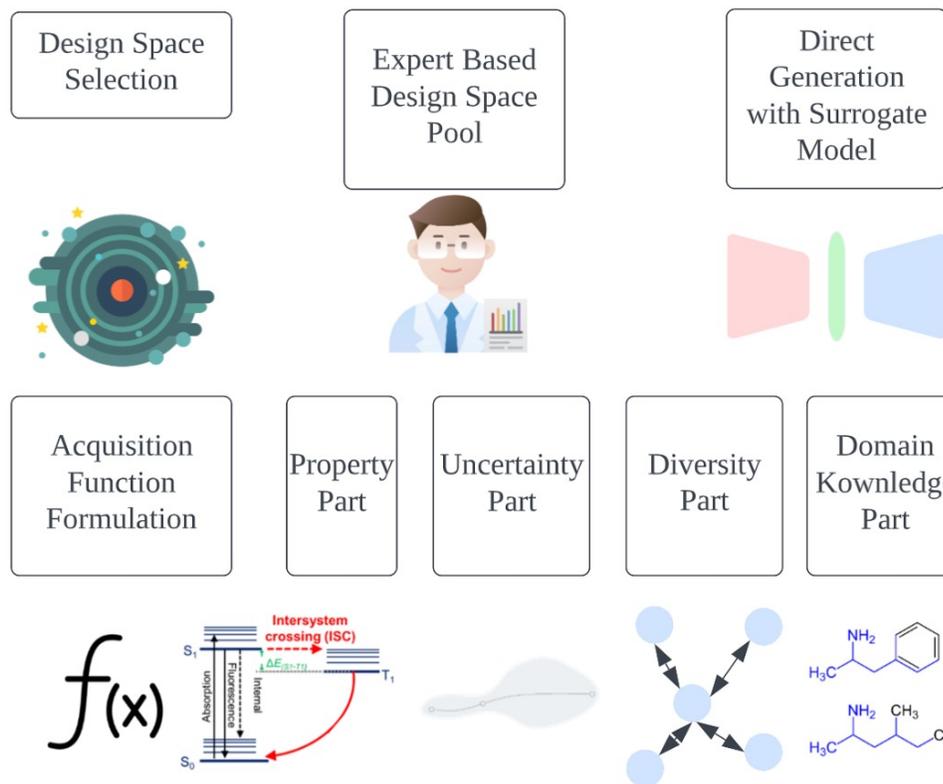
Generative models



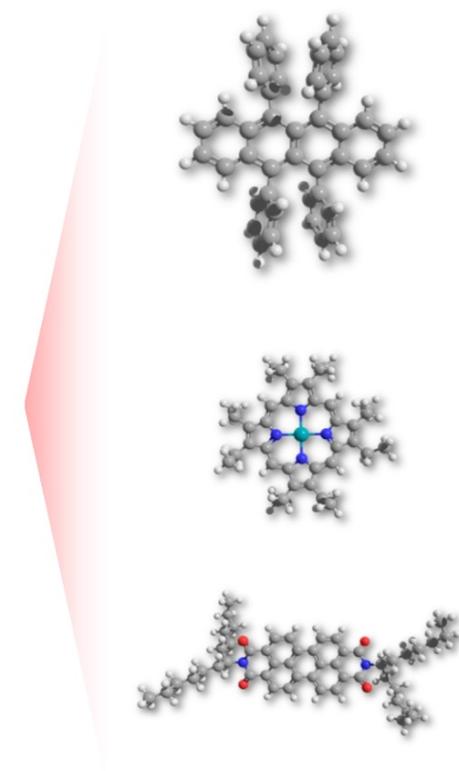
Summary



Method 1: xTB+ML



Method 2: Unified AL



Thanks!

Questions?



Massachusetts
Institute of
Technology

Imperial College
London



NUS
National University
of Singapore