

Modern Approaches to Methods Development



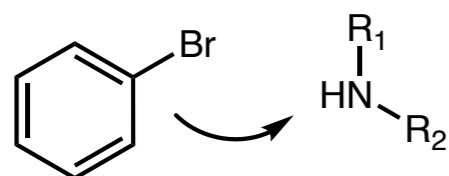
Johnny Wang

Group Meeting

Feb. 14, 2025

The path to developing a useful method

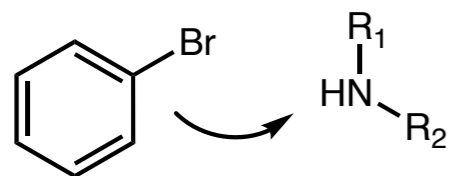
1. Identify a problem



C-N bonds are prevalent but
hard to form

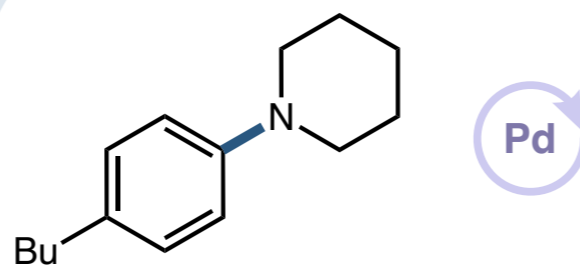
The path to developing a useful method

1. Identify a problem



C-N bonds are prevalent but hard to form

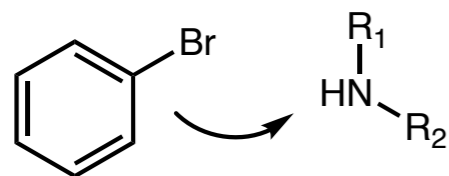
2. Initial hit



1983 - migita

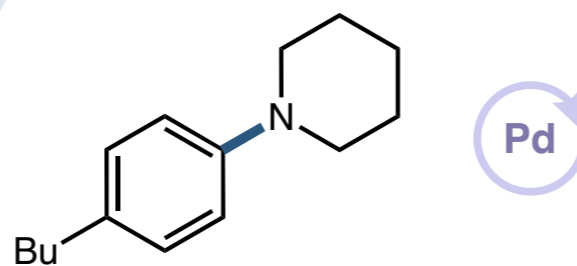
The path to developing a useful method

1. Identify a problem



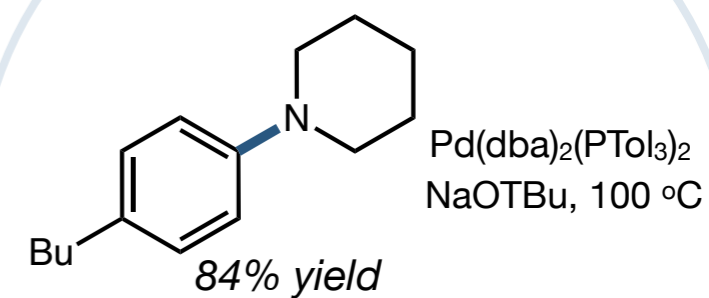
C-N bonds are prevalent but hard to form

2. Initial hit



1983 - migita

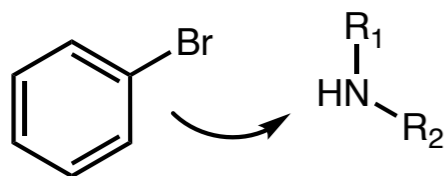
3. Optimization



1994 - Buchwald/Hartwig

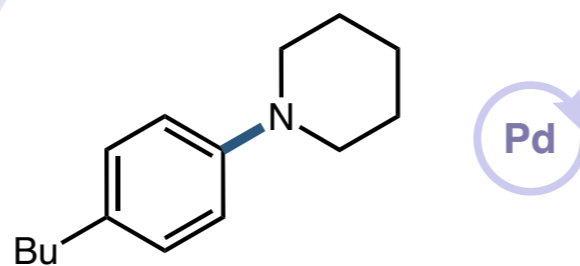
The path to developing a useful method

1. Identify a problem



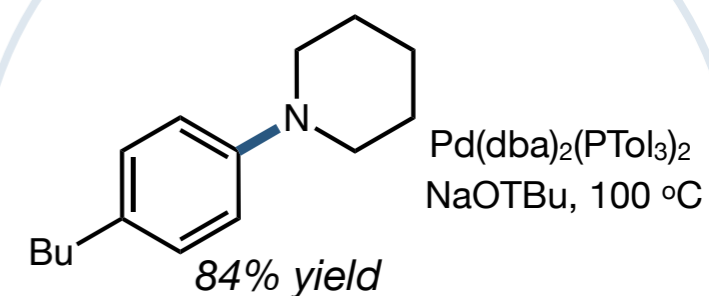
C-N bonds are prevalent but hard to form

2. Initial hit



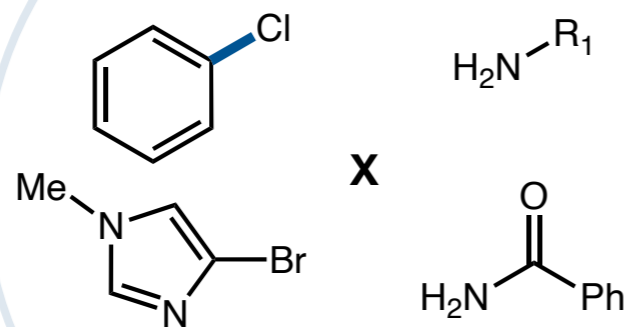
1983 - migita

3. Optimization



1994 - Buchwald/Hartwig

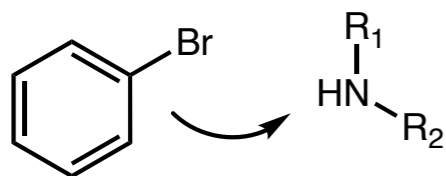
4. Generalize



~ongoing (but mostly done)

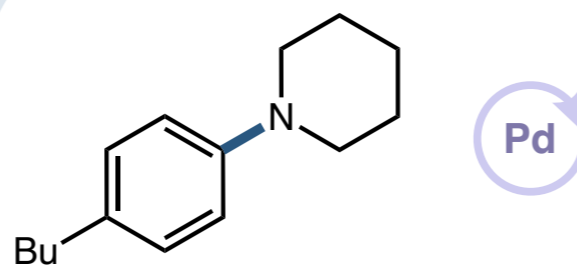
The path to developing a useful method

1. Identify a problem



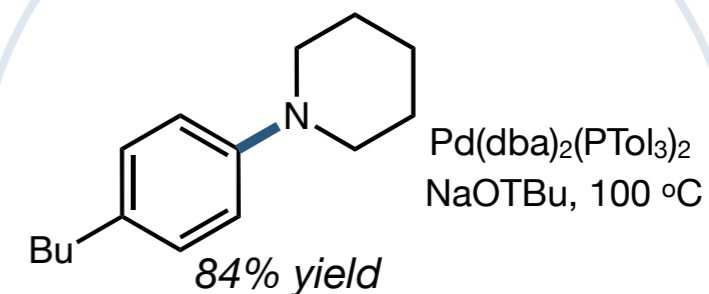
C-N bonds are prevalent but hard to form

2. Initial hit



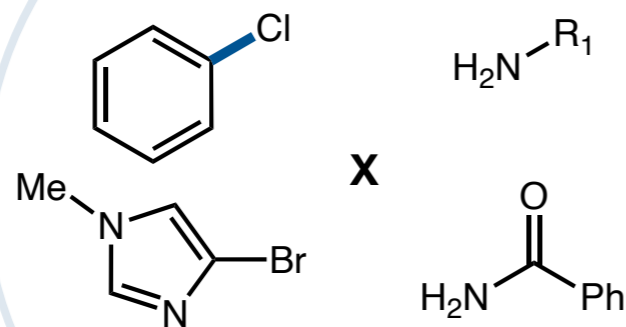
1983 - migita

3. Optimization



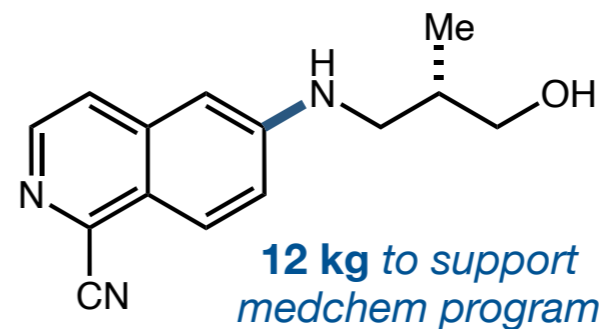
1994 - Buchwald/Hartwig

4. Generalize



~ongoing (but mostly done)

5. General adoption

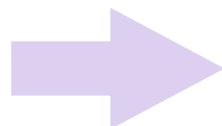


Widely adopted

Needs to be highly predictable

Traditional approach to methods development

**Initial discovery
1983**

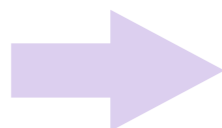


Today

- *Thousands of publication*
- *100's of ligands*
- *100's of substrate combinations*
- *Among most used reactions in industry*

Traditional approach to methods development

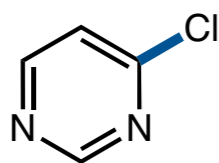
Initial discovery
1983



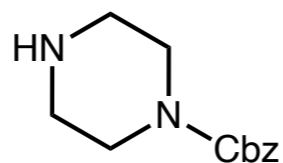
Today

- Thousands of publications
- 100's of ligands
- 100's of substrate combinations
- Among most used reactions in industry

Conditions and ligands known for almost any given combination

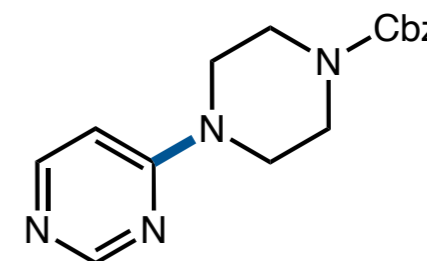


heteroaryl chloride



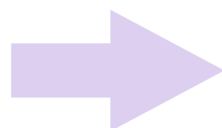
dialkyl amine

Conditions?



Traditional approach to methods development

Initial discovery
1983



Today

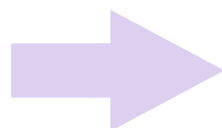
- Thousands of publications
- 100's of ligands
- 100's of substrate combinations
- Among most used reactions in industry

I_HAR	BINAP, Cs ₂ CO ₃	78%			Xantphos, NaOtBu	61%				Xantphos, Cs ₂ CO ₃	74%					
	Xantphos, Cs ₂ CO ₃	60%			BINAP, NaOtBu	52%										
	BINAP, K ₂ CO ₃	54%			Xantphos, Cs ₂ CO ₃	36%										
I_ARY	DPEphos, NaOtBu	88%	Xantphos, NaOtBu	81%		BINAP, NaOtBu	58%		P(Ph) ₃ , NaOtBu	80%	Xantphos, Cs ₂ CO ₃	54%				
	SPhos, NaOtBu	74%	P(tBu) ₃ , KOtBu	77%		P(o-tol) ₃ , NaOtBu	58%		P(tBu) ₃ , NaOtBu	72%						
	dppf, NaOtBu	74%	P(tBu) ₃ , NaOtBu	72%		BINAP, Cs ₂ CO ₃	57%									
Cl_HAR	XPhos, KOtBu	90%	P(tBu) ₃ , NaOtBu	70%	BrettPhos, LiHMDS	78%	DavePhos, KOtBu	80%	BINAP, Cs ₂ CO ₃	87%	P(tBu) ₃ , KOtBu	80%	BINAP, Cs ₂ CO ₃	62%	dCyph, Cs ₂ CO ₃	96%
	142691-72-3, NaOtBu	76%			BINAP, K ₂ CO ₃	73%	Trisobutylphosphatane, NaOtBu	80%	DavePhos, K ₂ CO ₃	84%	P(tBu) ₃ , K ₃ PO ₄	76%			dppf, K ₃ PO ₄	85%
	CyJohnPhos, NaOtBu	70%			Cy-tBu-Josiphos, NaOtBu	68%	DavePhos, NaOtBu	76%	Xantphos, Cs ₂ CO ₃	40%	P(tBu) ₃ , NaOtBu	70%			BINAP, Cs ₂ CO ₃	60%
Cl_ARY	SiPr, KOtBu	98%	P(tBu) ₃ , NaOtBu	75%	Cy-tBu-Josiphos, NaOtBu	94%	Trisobutylphosphatane, NaOtBu	83%			cBRIDP, NaOtBu	78%			JohnPhos, Cs ₂ CO ₃	94%
	BrettPhos, NaOtBu	94%	SPhos, NaOtBu	60%	Ad-BippyPhos, KOPh	90%	XPhos, NaOtBu	71%			P(tBu) ₃ , NaOtBu	68%			tBuBrettPhos, K ₃ PO ₄	92%
	P(Ph) ₃ , NaOtBu	88%			Cy-tBu-Josiphos, LiHMDS	85%	JohnPhos, NaOtBu	60%							BrettPhos, Cs ₂ CO ₃	51%
Br_HAR	DPEphos, KOtBu	83%	P(tBu) ₃ , NaOtBu	68%	Cy-tBu-Josiphos, NaOtBu	72%	Xantphos, NaOtBu	73%			Xantphos, Cs ₂ CO ₃	72%	BINAP, Cs ₂ CO ₃	74%	BINAP, Cs ₂ CO ₃	72%
	BINAP, KOtBu	79%	dppf, NaOtBu	60%	BINAP, NaOtBu	62%	RuPhos, LiHMDS	71%			P(tBu) ₃ , NaOtBu	65%	BINAP, NaOtBu	66%	Xantphos, Cs ₂ CO ₃	54%
	tBuXPhos, NaOtBu	70%			Xantphos, NaOtBu	52%	RuPhos, Cs ₂ CO ₃	70%			P(Ph) ₃ , NaOtBu	64%			XPhos, Cs ₂ CO ₃	46%
Br_ARY	SPhos, Cs ₂ CO ₃	95%	RuPhos, LiHMDS	94%	Ad-BippyPhos, KOPh	95%	XPhos, NaOPent	83%	P(tBu) ₃ , NaOtBu	76%	tBuXPhos, NaOtBu	90%	dppf, NaOtBu	88%	dppf, NaOtBu	81%
	QUINAP, NaOtBu	91%	P(tBu) ₃ , NoBase	81%	Cy-tBu-Josiphos, NaOtBu	94%	Trisobutylphosphatane, NaOtBu	78%	JohnPhos, NaOtBu	75%	P(tBu) ₃ , K ₃ PO ₄	83%	BINAP, NaOtBu	85%	cBRIDP, NaOtBu	80%
	QUINAP, Cs ₂ CO ₃	90%	P(tBu) ₃ , KOtBu	79%	Xantphos, KOtBu	86%	RuPhos, LiHMDS	77%	BINAP, Cs ₂ CO ₃	74%	P(o-tol) ₃ , NaOtBu	80%	Xantphos, Cs ₂ CO ₃	64%	Xantphos, K ₃ PO ₄	80%
	Aryl		DiAryl		Alkyl		DiAlkyl		Alkyl-Aryl		aromN		Ketimine		Amide	

Electrophile Type

Traditional approach to methods development

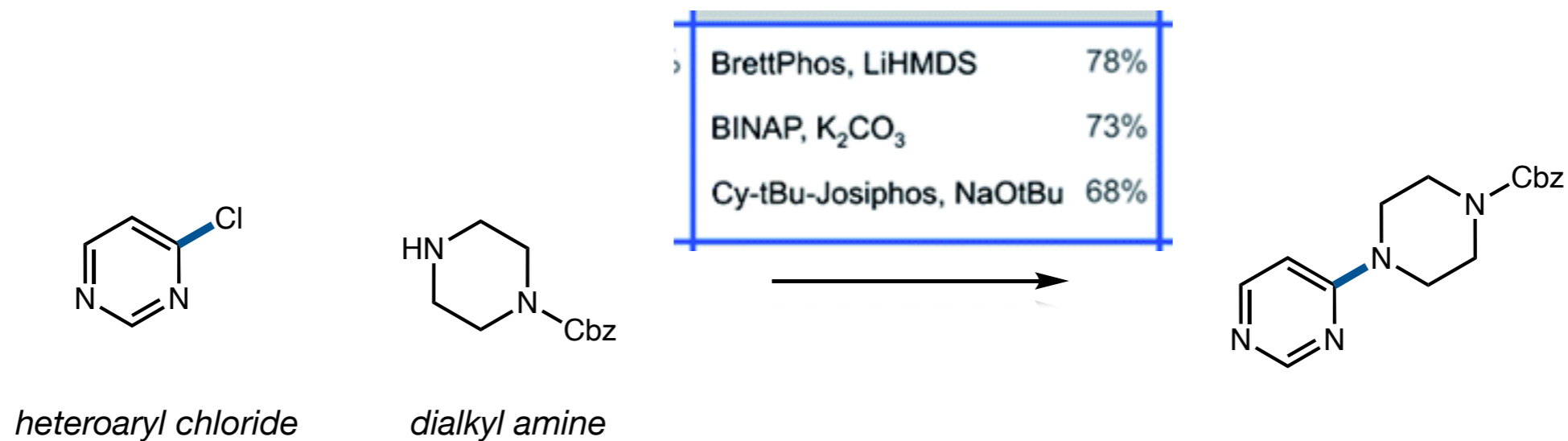
Initial discovery
1994



Today

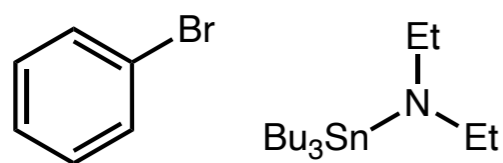
- Thousands of publications
- 100's of ligands
- 100's of substrate combinations
- Among most used reactions in industry

Conditions and ligands known for almost any given combination

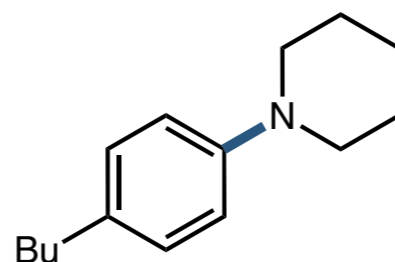


Traditional approach to methods development-belaboring the point

Initial hit -1983



First practical system -1994



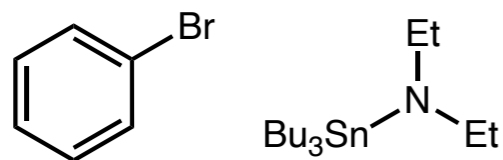
Today: Mostly generalized

40 years, 1390 publications

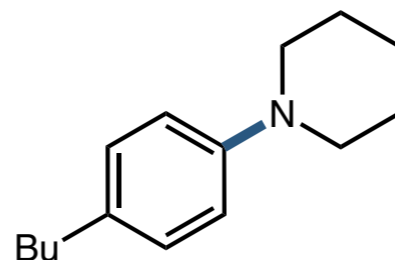
**New combinations still
being optimized**

Traditional approach to methods development-belaboring the point

Initial hit -1983



First practical system -1994

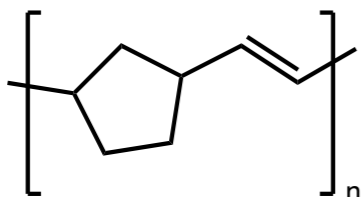


Today: Mostly generalized

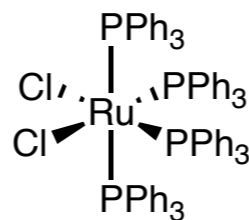
40 years, 1390 publications

**New combinations still
being optimized**

Initial hit -1955



First practical system -1992



Today: Mostly generalized

70 years, 15,622 publications

**New combinations still
being optimized**

Development of a general, useful reaction is slow

Traditional approach to methods development

How can the “generalization” and development of new reactions be greatly accelerated?

Modern approaches to methods development

Modern Paradigms in Screening

- *Reaction generalization*
- *“Accelerate” Serendipity*
- *Miniaturization of unique reaction set ups*

Data Science

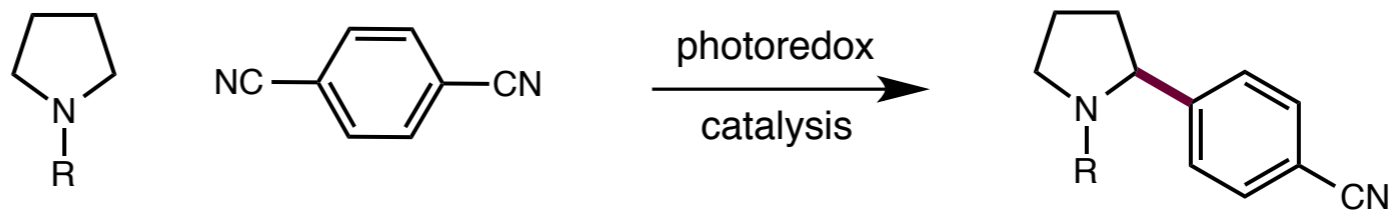
- *Catalyst optimization*
- *Predicting selectivity*
- *Discovery of new catalysts*

Machine Learning

- *What is machine learning*
- *Prediction of optimal conditions*
- *Selectivity prediction for complex systems*
- *Catalyst Discovery*

Accelerated Serendipity

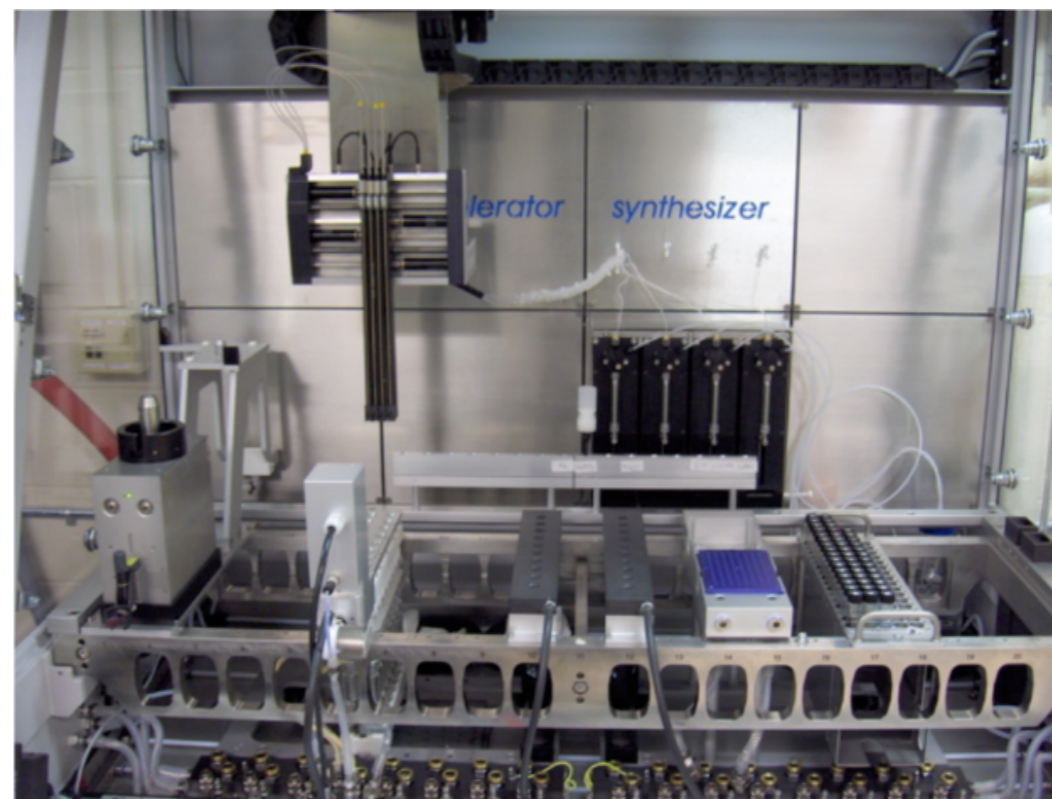
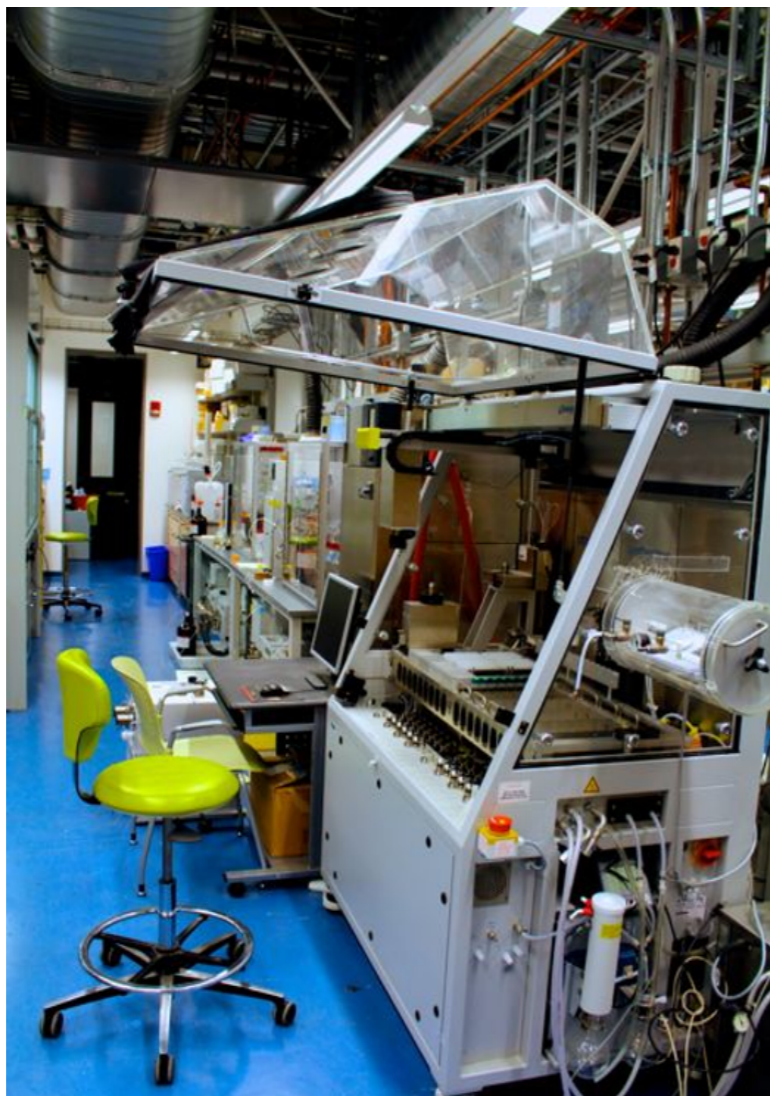
D,E	D,F	E,F	3 x
B,F	C,D	C,E	C,F
A,F	B,C	B,D	B,E
A,B	A,C	A,D	A,E



Discovery of new amine C-H functionalization through accelerated serendipity

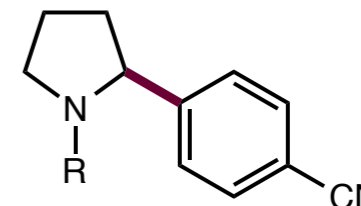
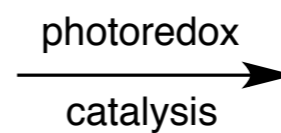
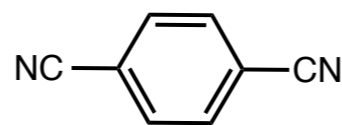
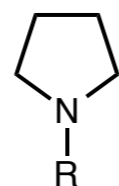
Accelerated Serendipity

Puts together 500 - 1,000 reactions per day in a highly controlled fashion



- Parallel processing of reactions at multiple temperatures in multiple solvents
- Multi channel filtration and vacuum capabilities, solid weighing and work ups performed
- 96 well LED plates with easy installation of LEDs of variable wavelengths

Accelerated Serendipity

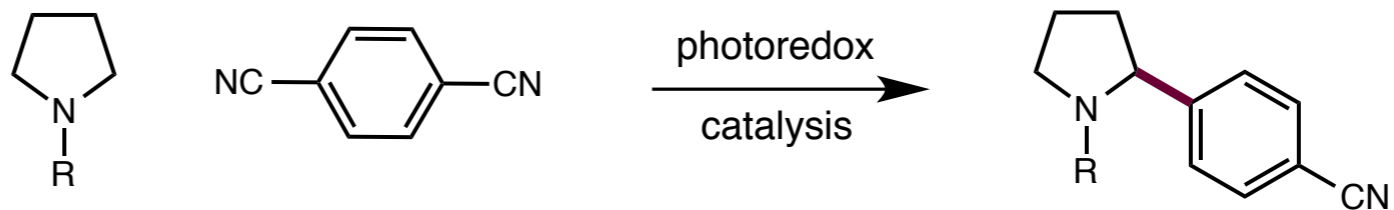


Discovery of new amine C-H functionalization through accelerated serendipity

19 x 19 matrix

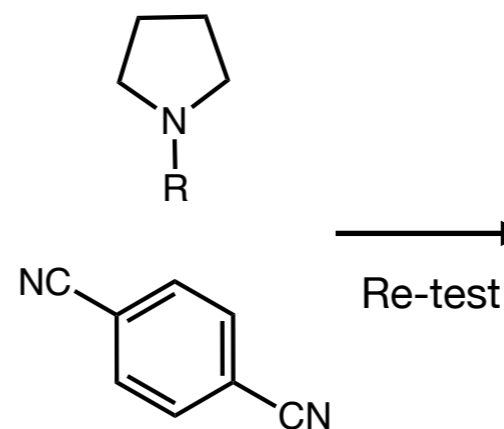
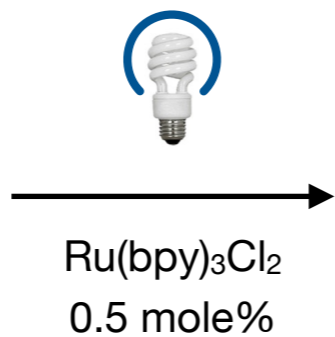


Accelerated Serendipity



Discovery of new amine C-H functionalization through accelerated serendipity

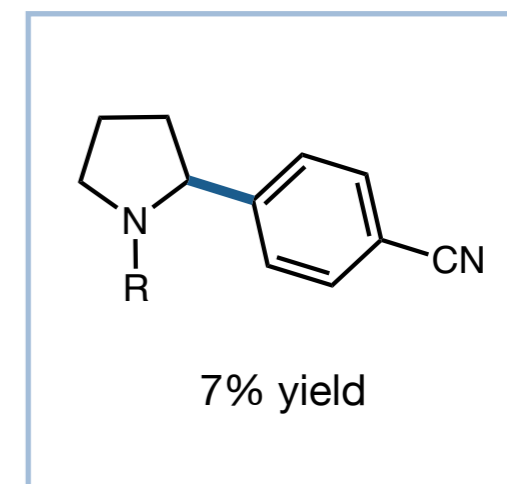
19 x 19 matrix



**Reaction
Evaluation**

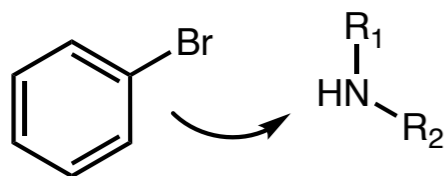


**Initial
Result**



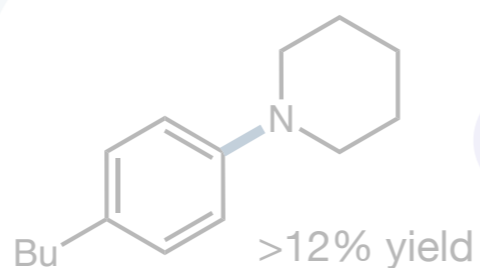
Accelerated Serendipity

1. Identify a problem



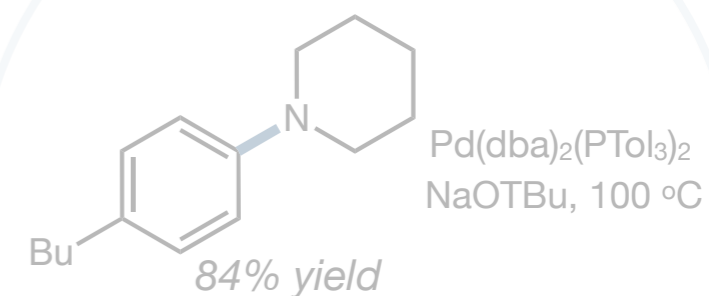
C-N bonds are prevalent but hard to form

2. Initial hit



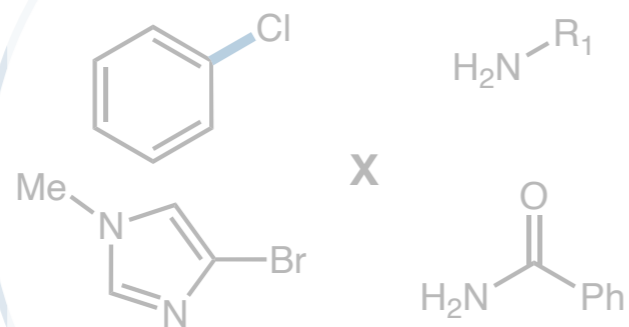
Bond formation at low yields

3. Optimization



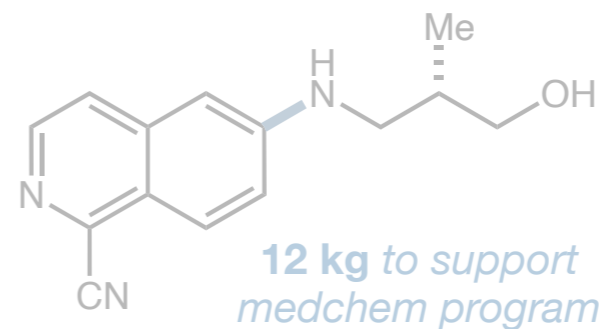
Conditions give high yields for initial system

4. Generalize



Works for most substrate combinations

5. General adoption

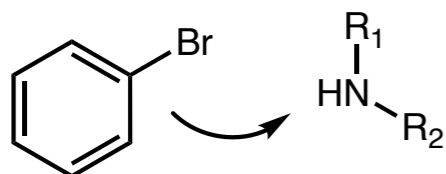


Wide adoption by community

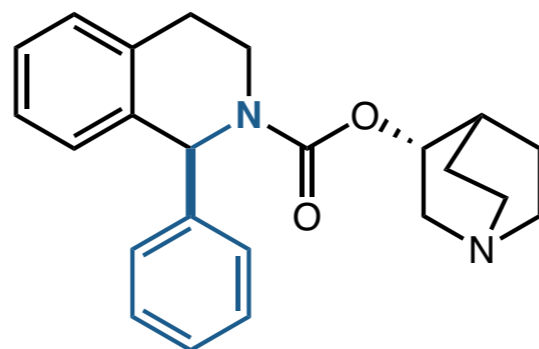
Needs to have easily predictable results

Accelerated Serendipity

1. Identify a problem

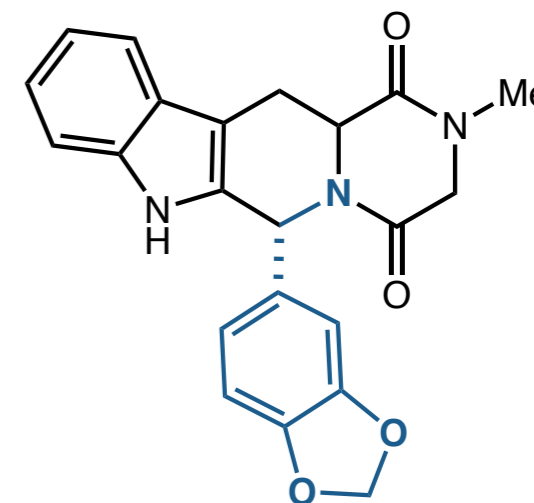


C-N bonds are prevalent but hard to form



Vesicare (GSK) No. 115/200

Muscarinic receptor antagonist

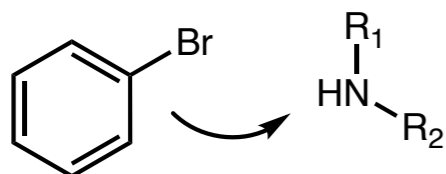


Cialis (Lilly) No. 66/200

Phosphodiesterase inhibitor

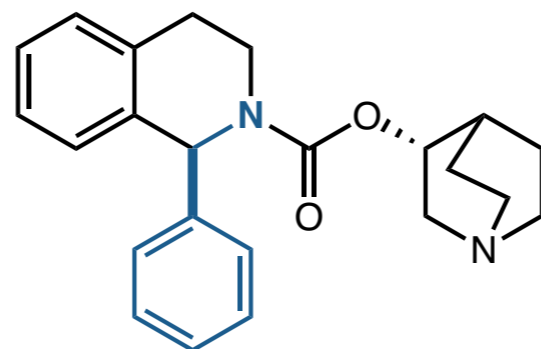
Accelerated Serendipity

1. Identify a problem

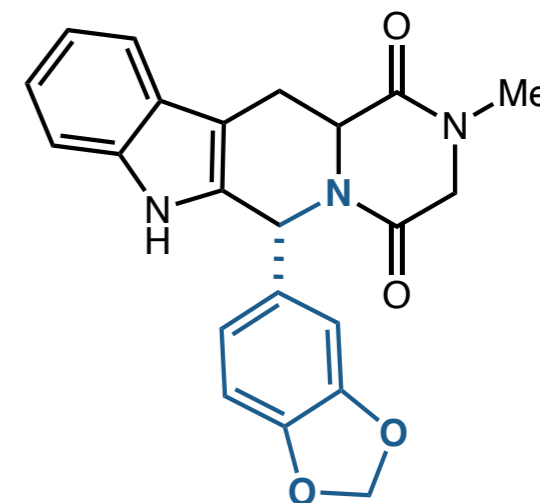


C-N bonds are prevalent but hard to form

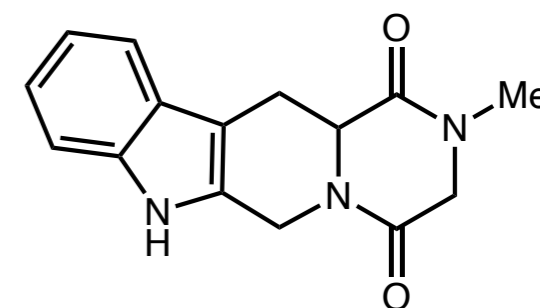
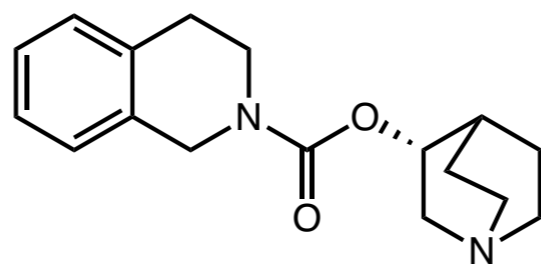
C-H arylation is an extremely desirable reaction



Vesicare (GSK) No. 115/200
Muscarinic receptor antagonist

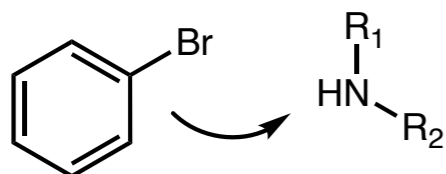


Cialis (Lilly) No. 66/200
Phosphodiesterase inhibitor



Accelerated Serendipity

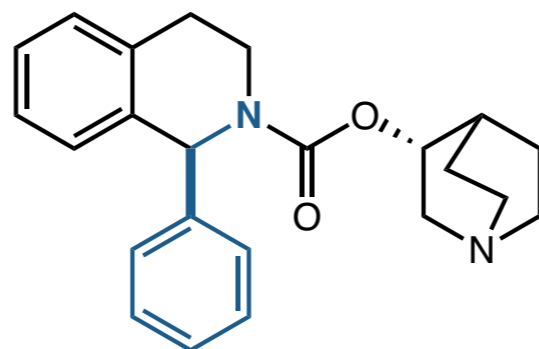
1. Identify a problem



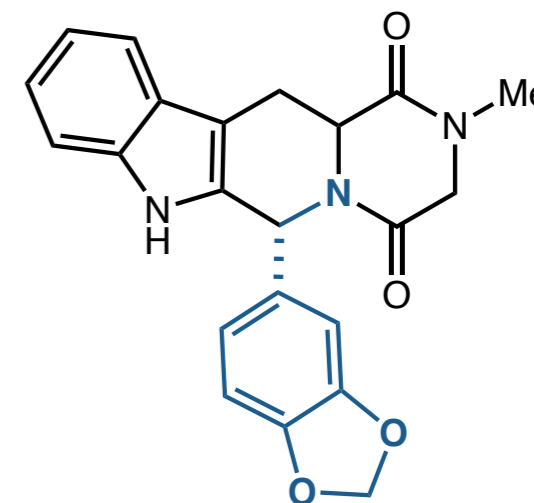
C-N bonds are prevalent but hard to form

C-H arylation is an extremely desirable reaction

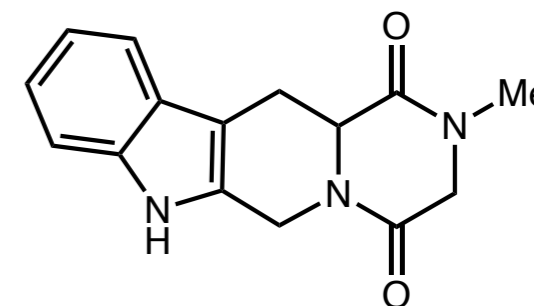
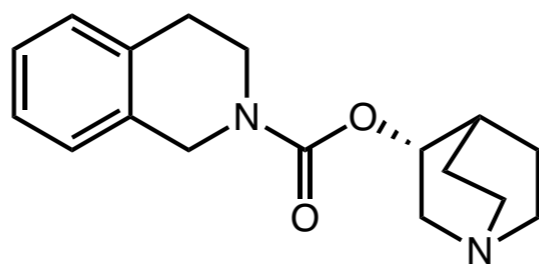
Not compatible with the synthetic logic of the time!



Vesicare (GSK) No. 115/200
Muscarinic receptor antagonist

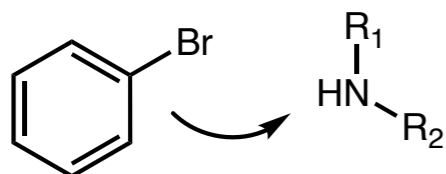


Cialis (Lilly) No. 66/200
Phosphodiesterase inhibitor



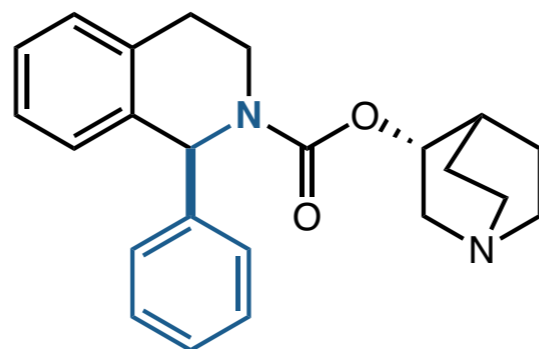
Accelerated Serendipity

1. Identify a problem

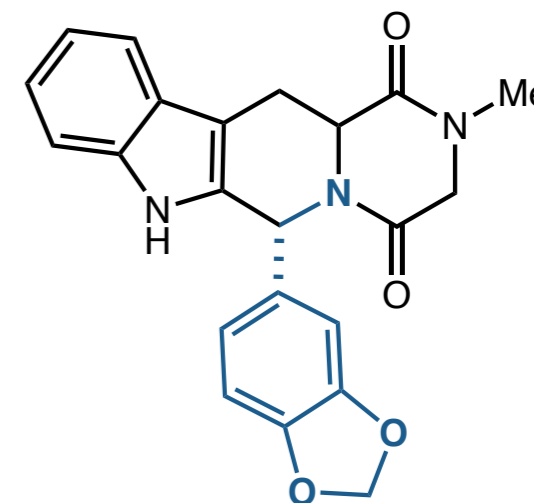


C-N bonds are prevalent but hard to form

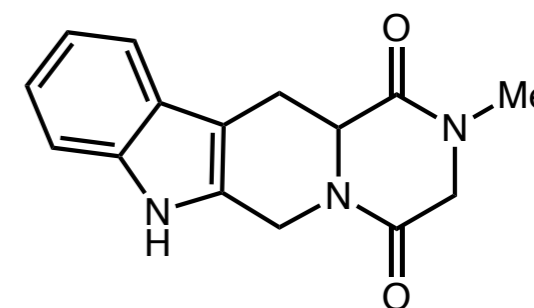
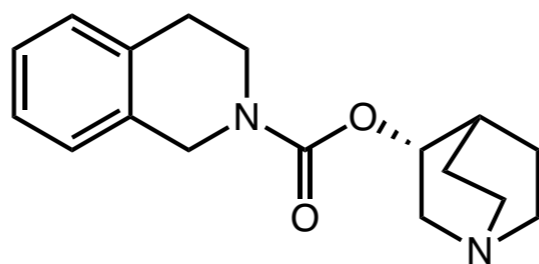
“Predict” a problem and solve it at the same time



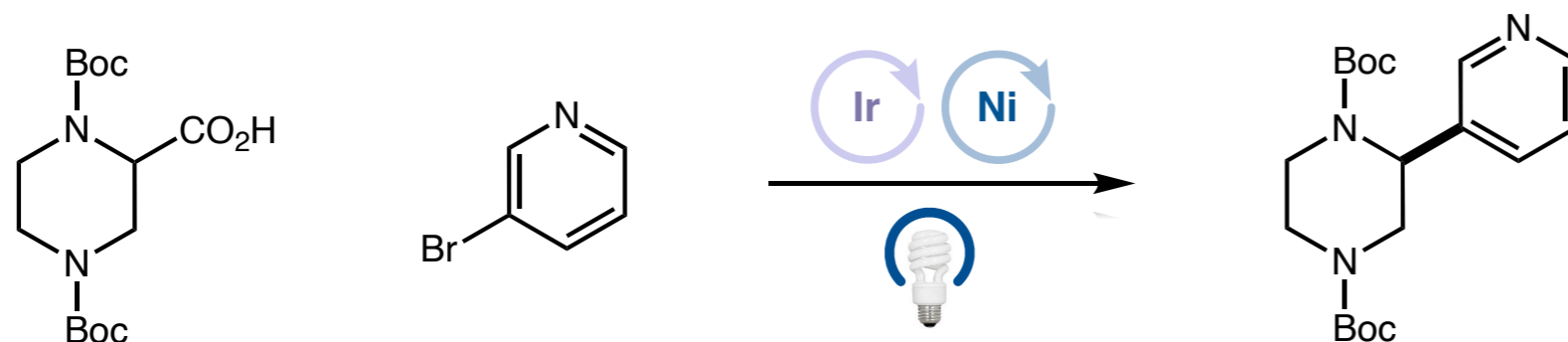
Vesicare (GSK) No. 115/200
Muscarinic receptor antagonist



Cialis (Lilly) No. 66/200
Phosphodiesterase inhibitor

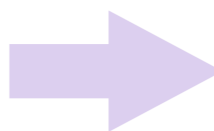


Additive Screening



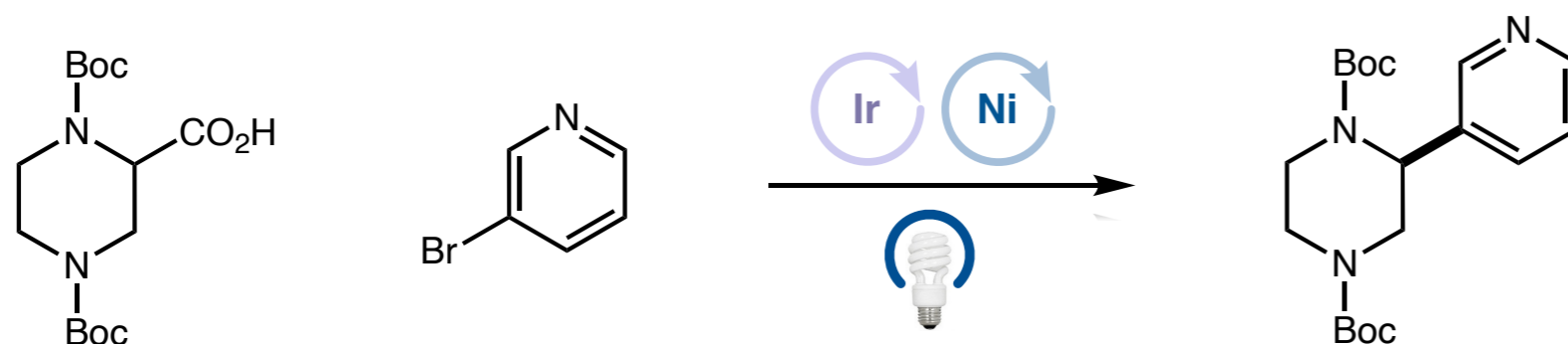
Can this reaction be generalized faster than historical approaches?

Design a new catalyst for every problematic combination

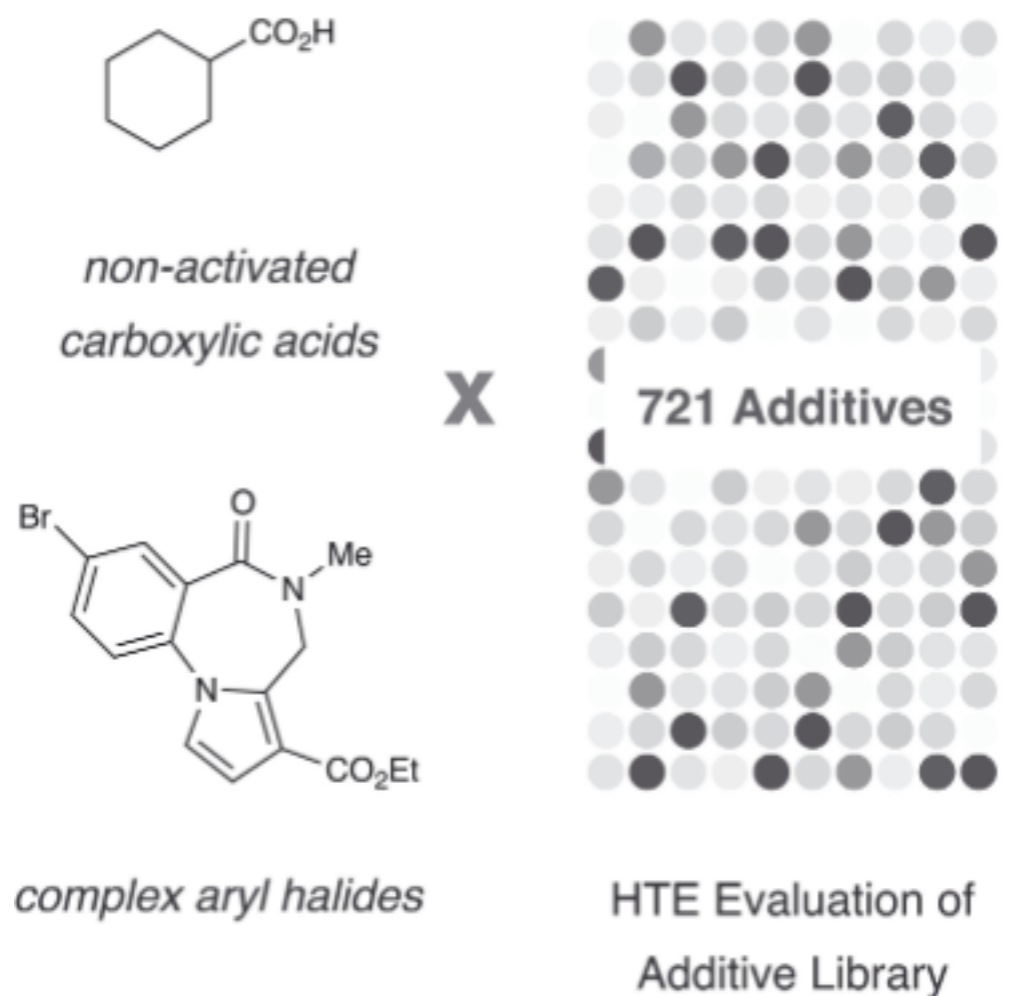


Universal additive that helps all substrate combinations?

Additive Screening

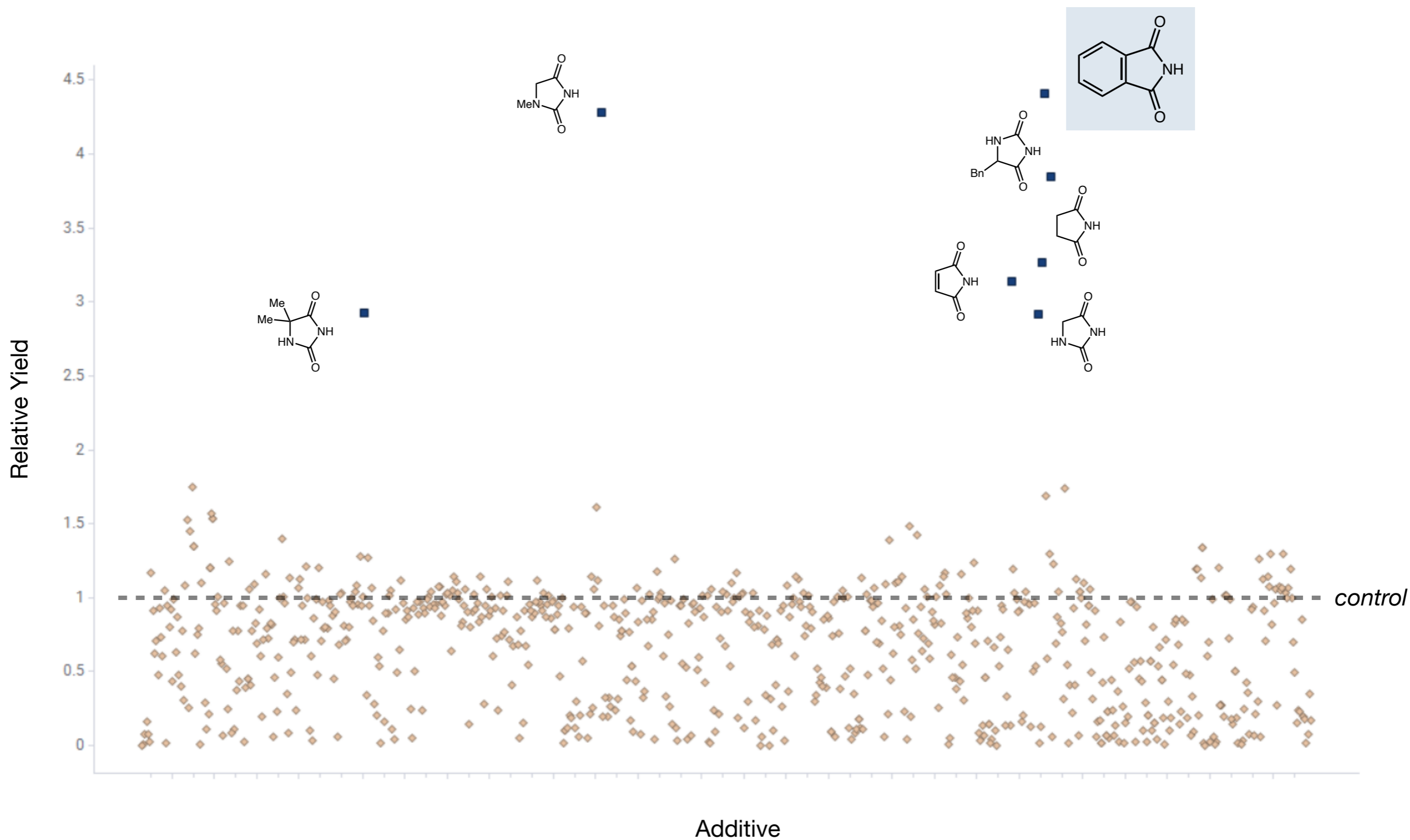


Can this reaction be generalized faster than historical approaches?

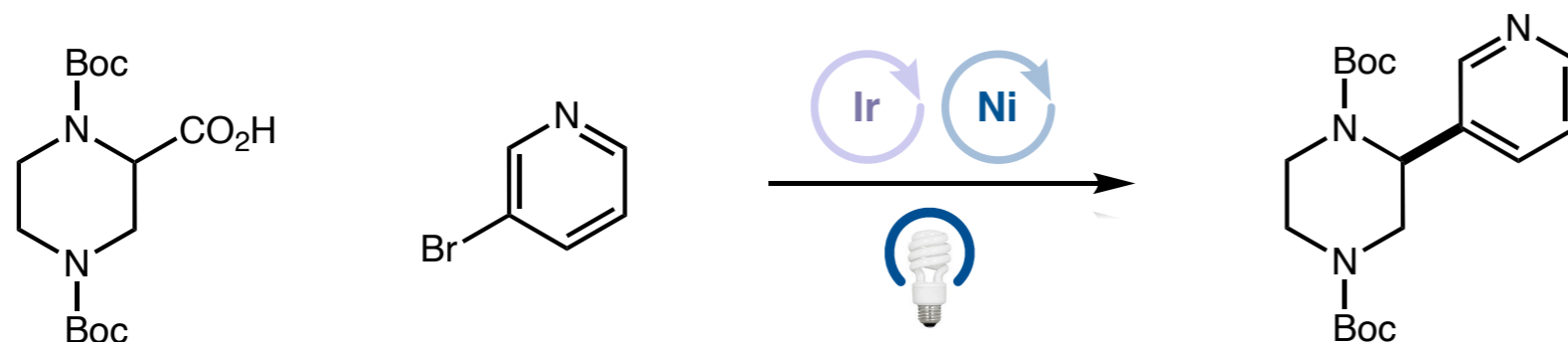


Universal additive that helps all substrate combinations?

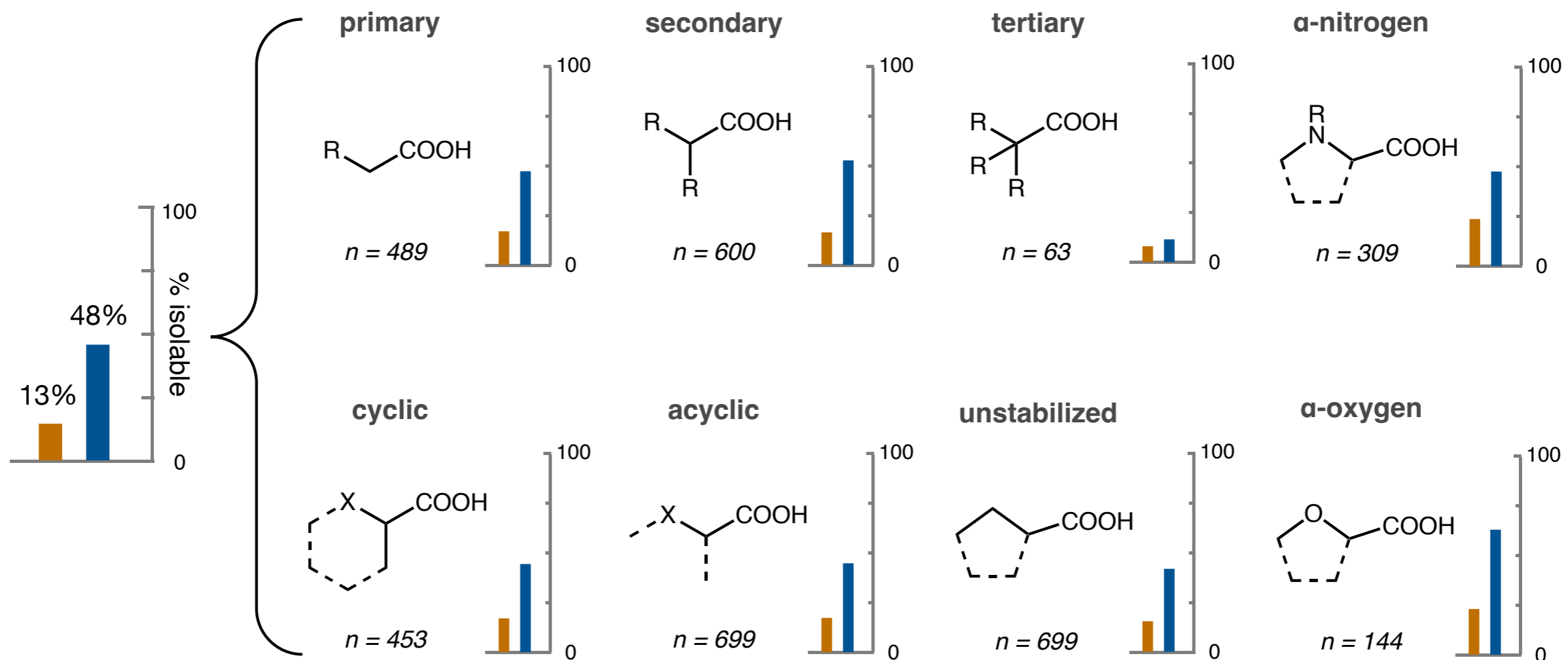
Additive Screening



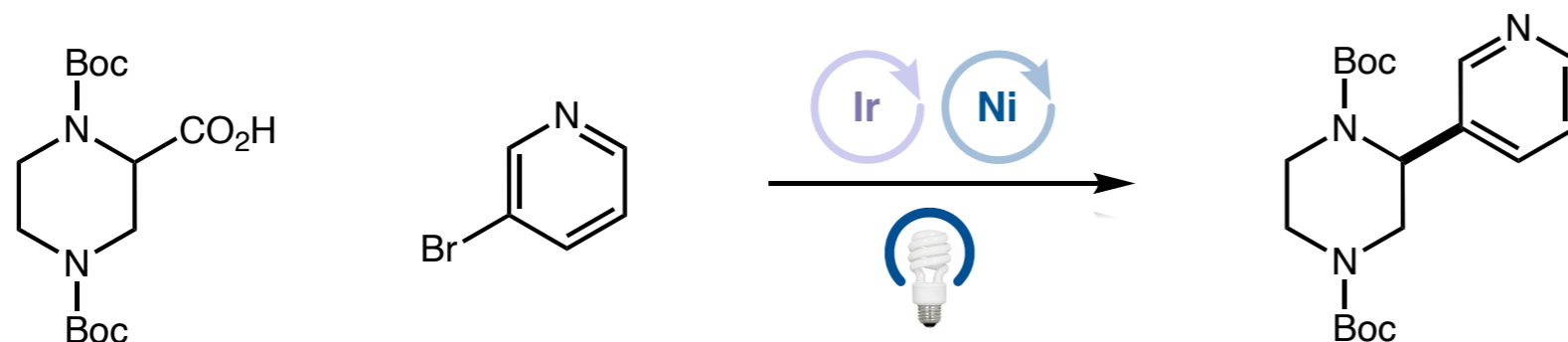
Additive Screening



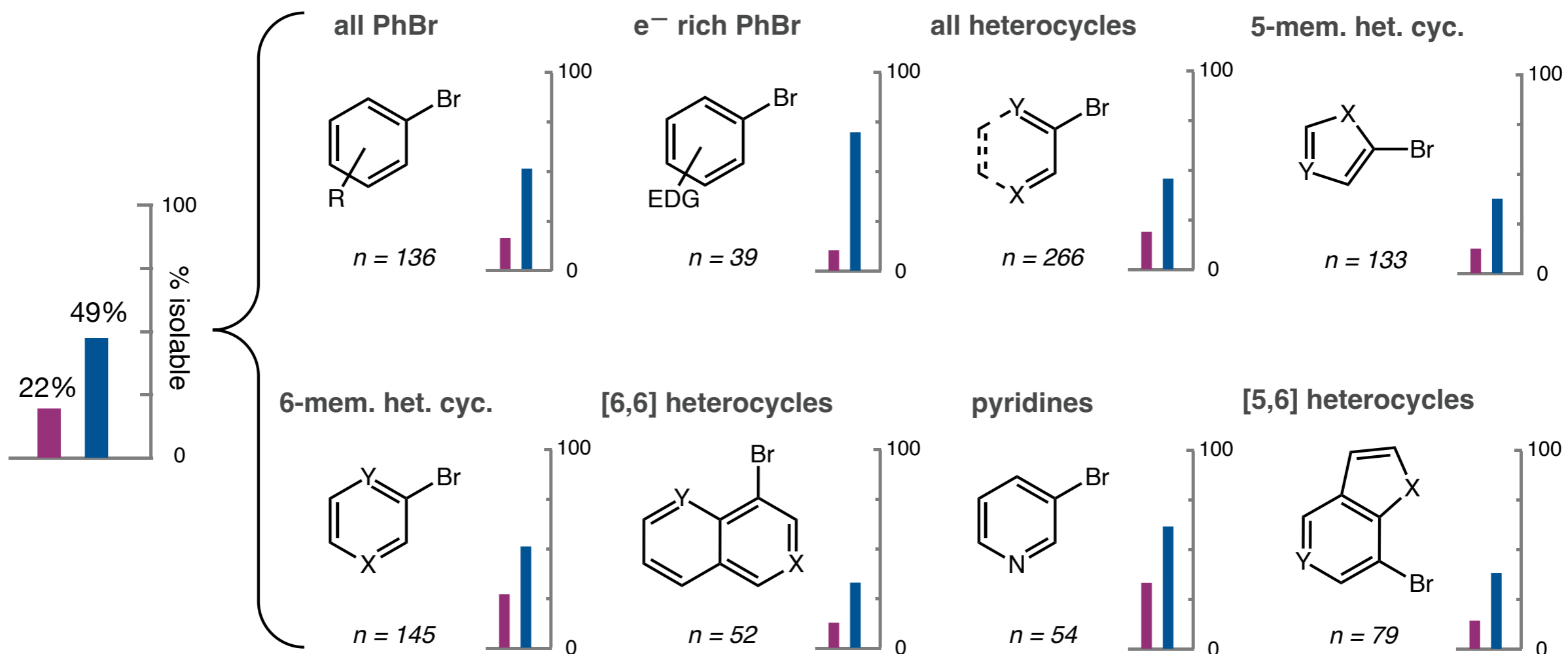
Can this reaction be generalized faster than historical approaches?



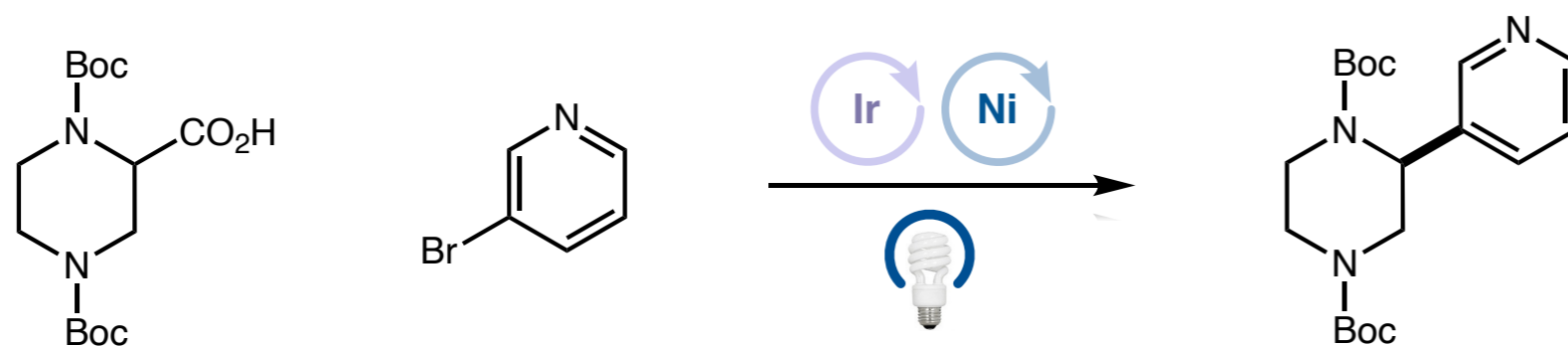
Additive Screening



Can this reaction be generalized faster than historical approaches?

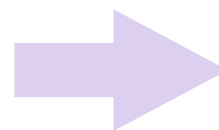
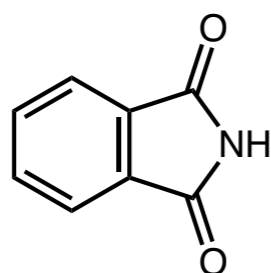


Additive Screening



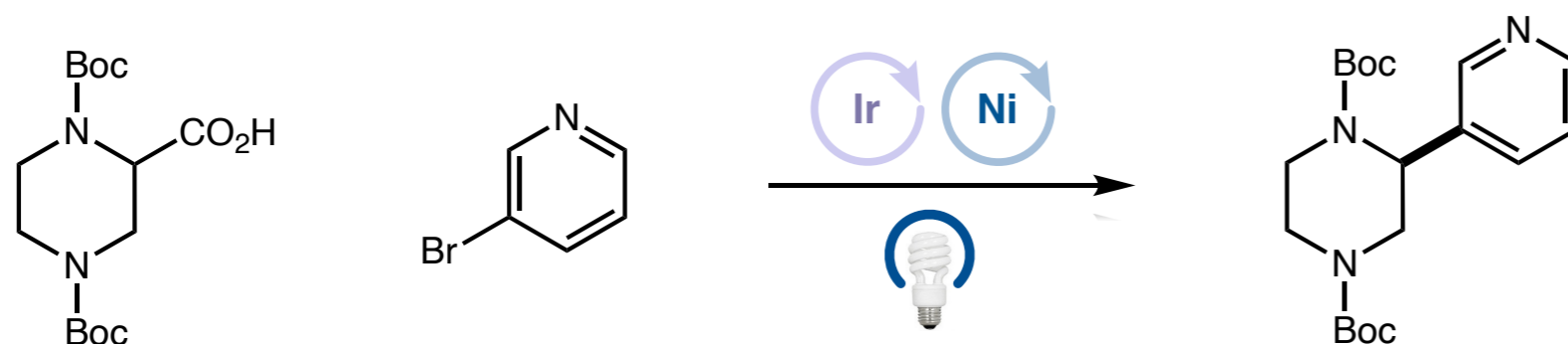
Can this reaction be generalized faster than historical approaches?

just add phthalimide

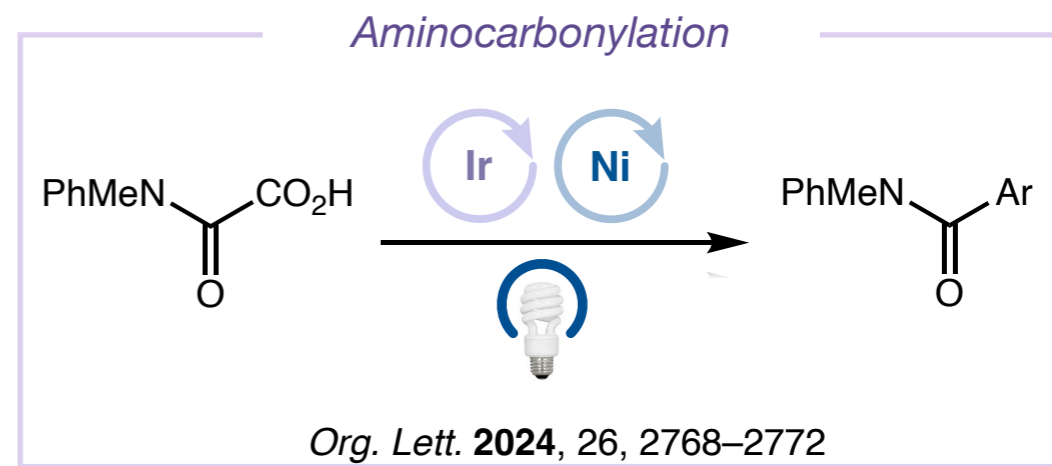
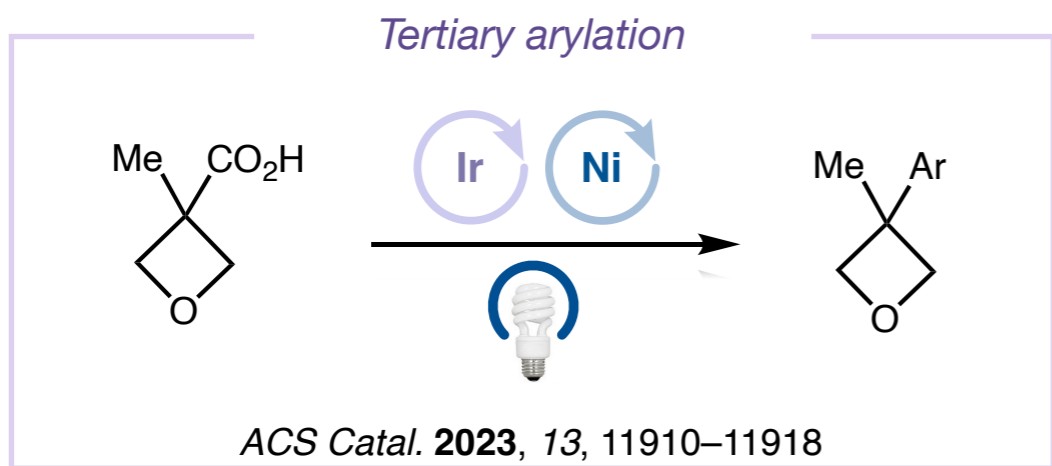
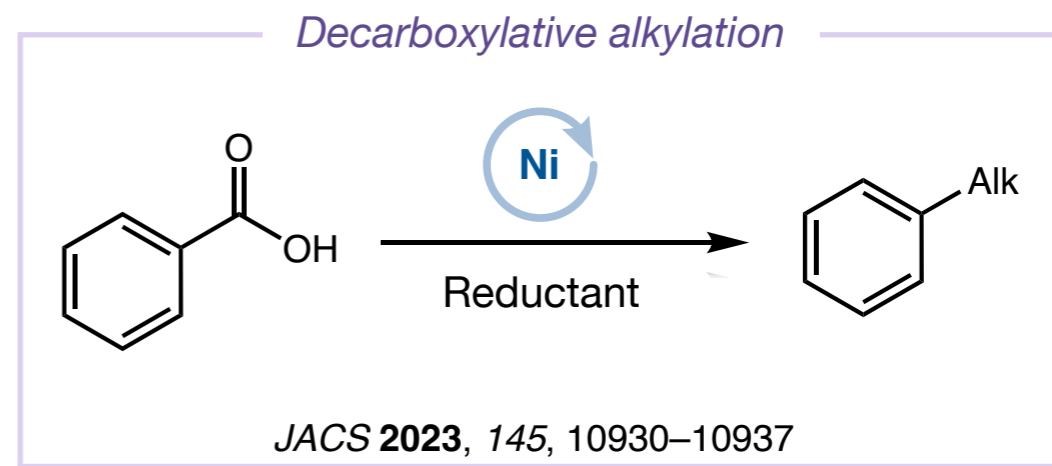
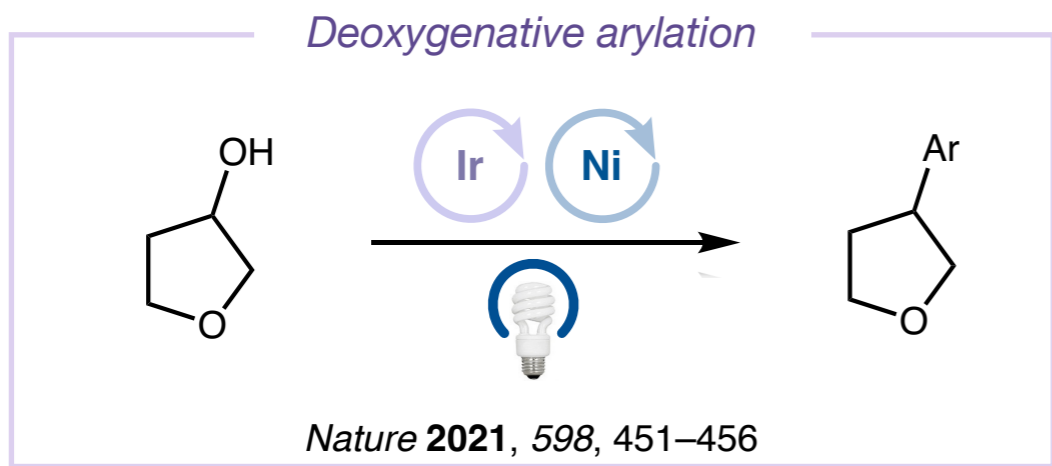


Greatly more general reaction

Additive Screening

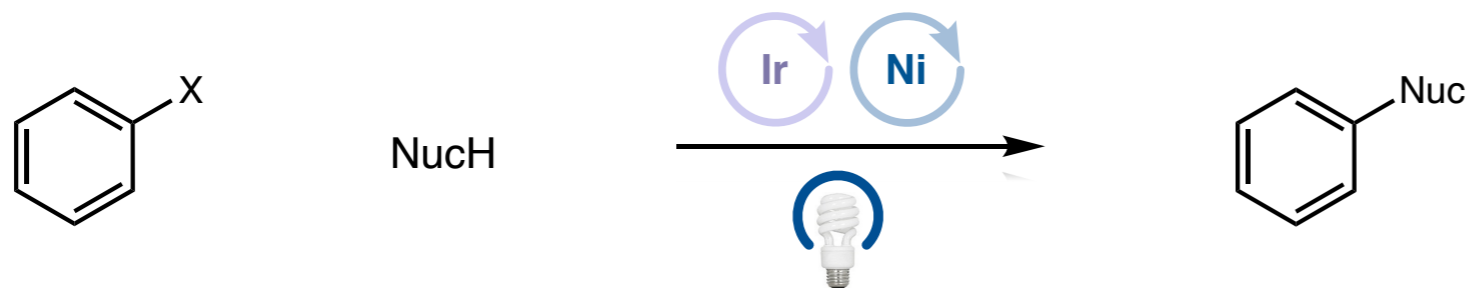


Can this reaction be generalized faster than historical approaches?



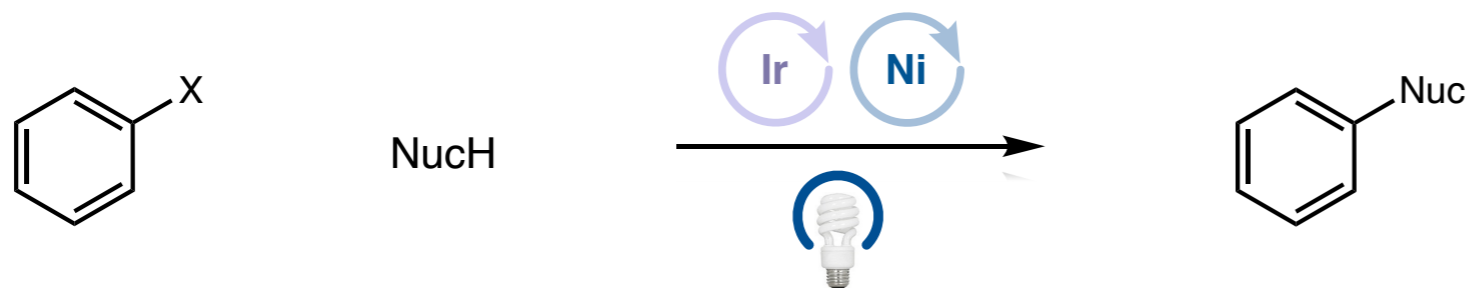
Some of many examples where phthalamide has helped with optimizing and generalizing novel methods

Additive Screening

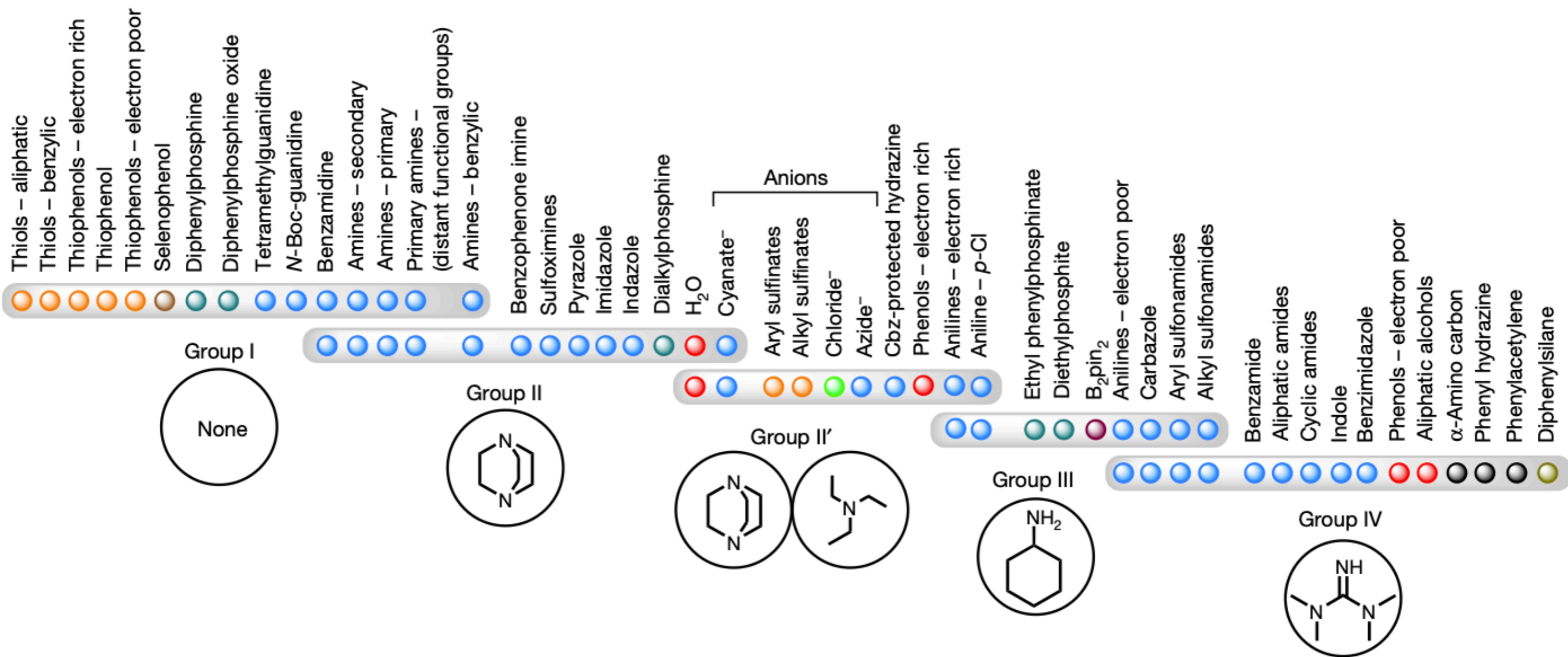


**General conditions using
different additives for
different Nuc**

Additive Screening

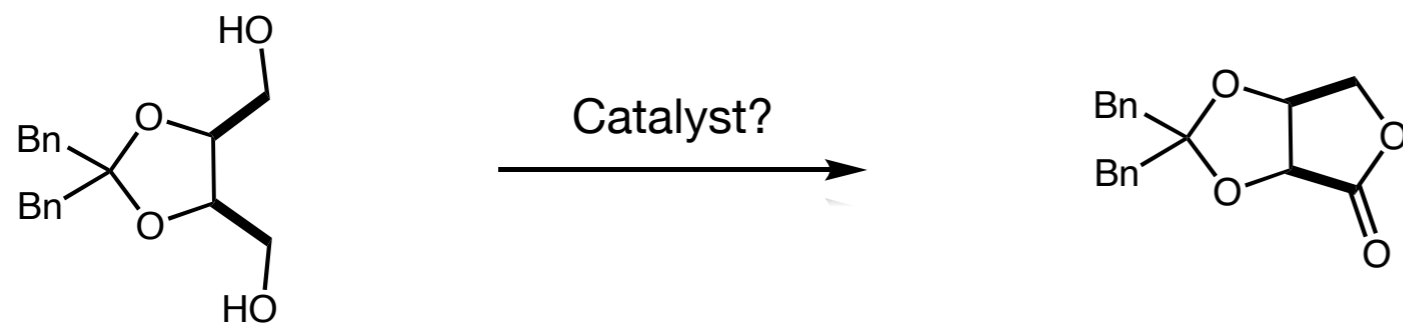


General conditions using different additives for different Nuc

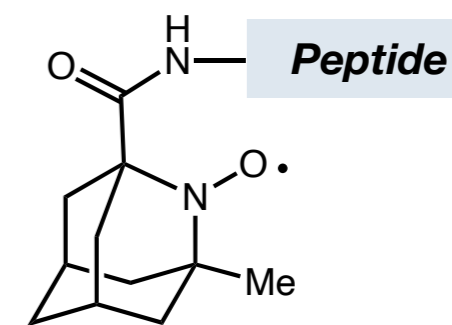


Generalization Through Multi-Substrate Optimization in Enantioselective Catalysis

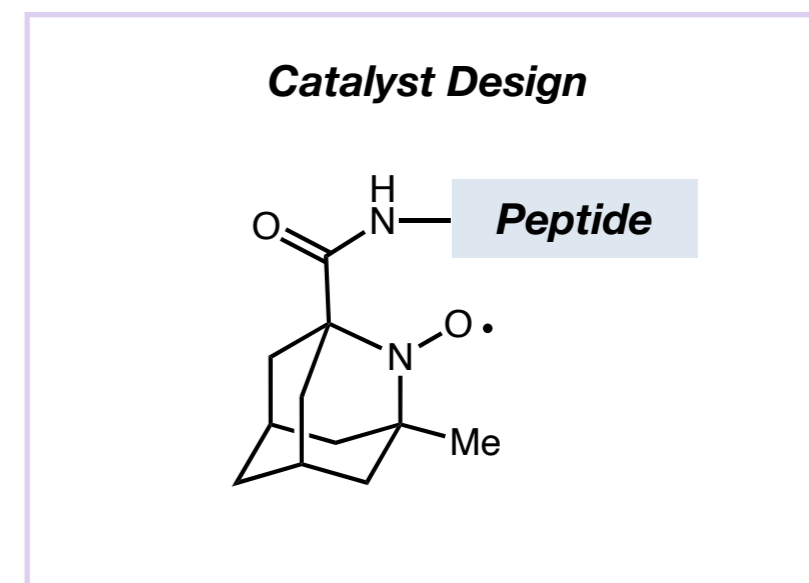
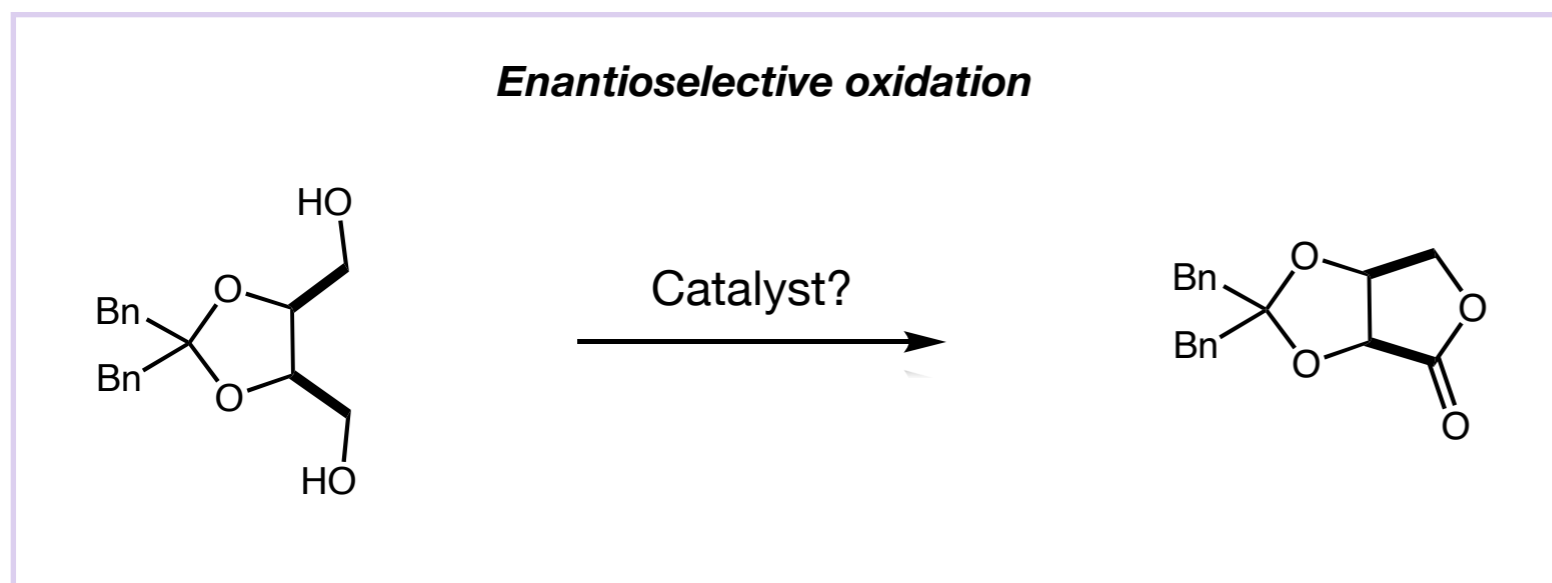
Enantioselective oxidation



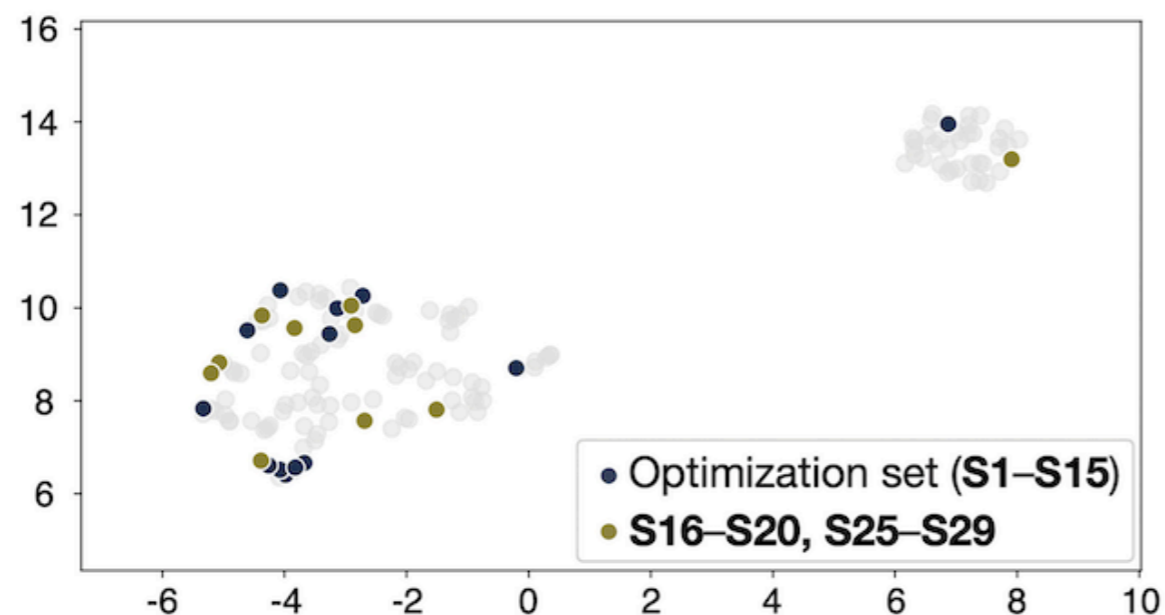
Catalyst Design



Generalization Through Multi-Substrate Optimization in Enantioselective Catalysis

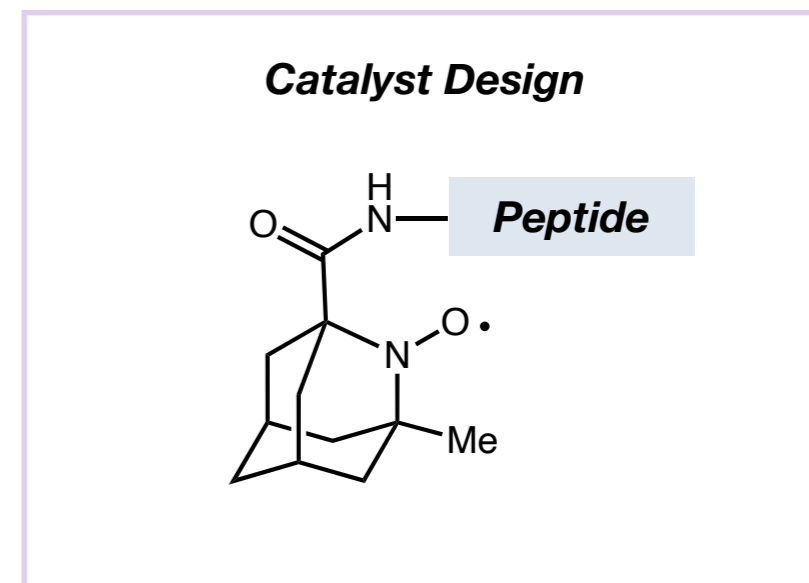
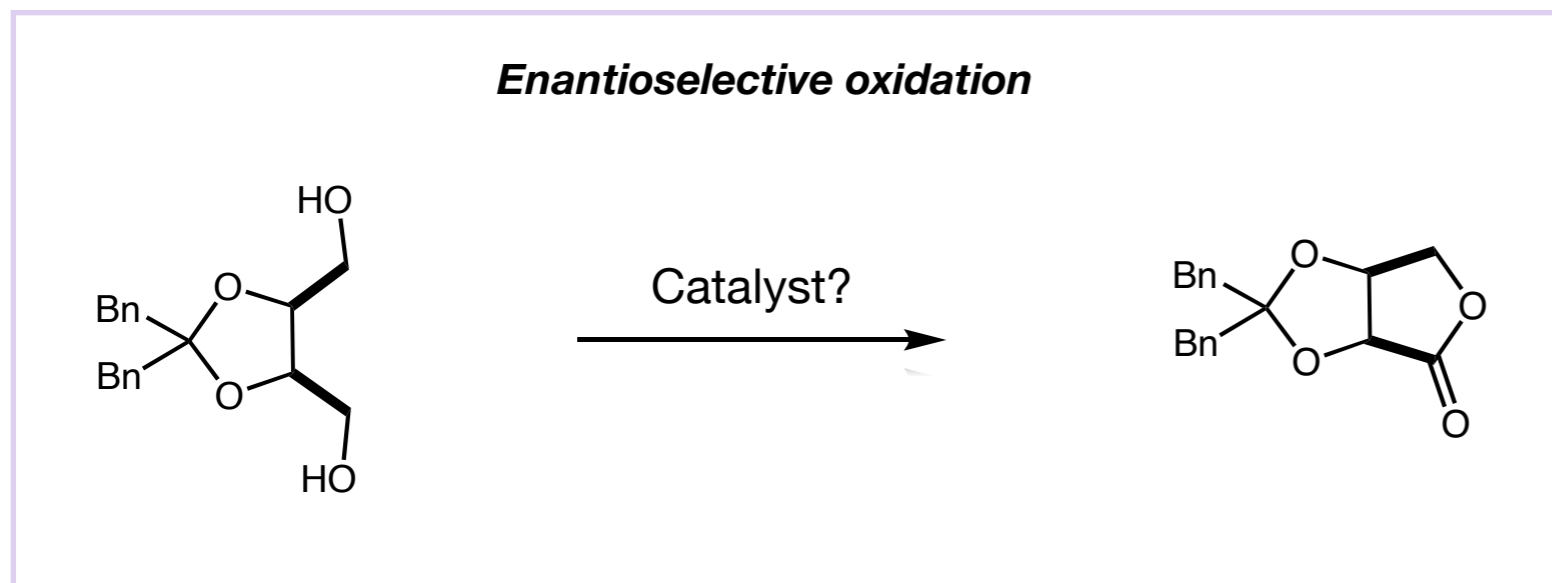


2D visualization of 1,4 diol chemical space

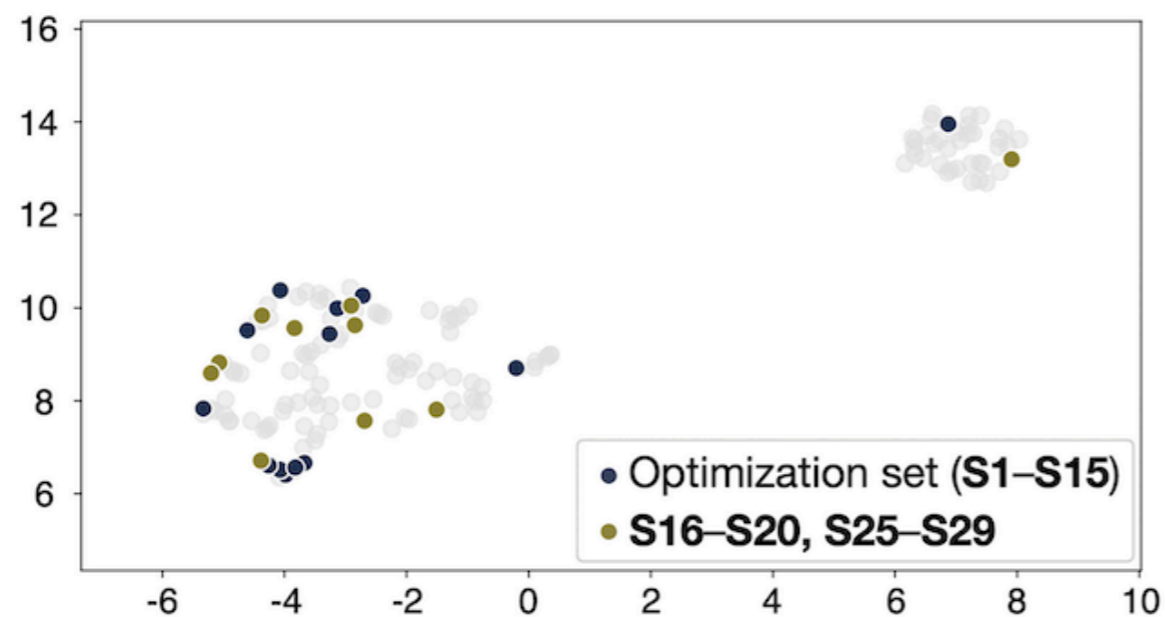


For more details about scope parameterization:
Kariofillis, S. K.; Jiang, S.; Zuranski, A. M.; Gandhi, S. S.; Alvarado, J. I. M.; Doyle, A. G. *J. Am. Chem. Soc.* **2022**, *144*, 1045–1055

Generalization Through Multi-Substrate Optimization in Enantioselective Catalysis



2D visualization of 1,4 diol chemical space



15 optimization substrates

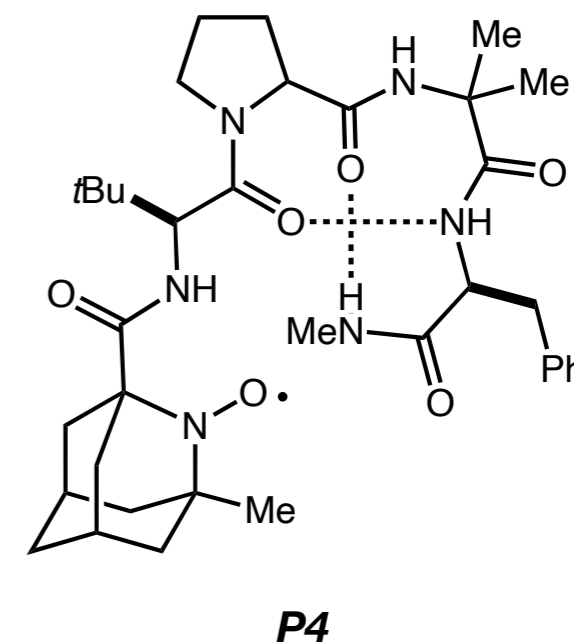
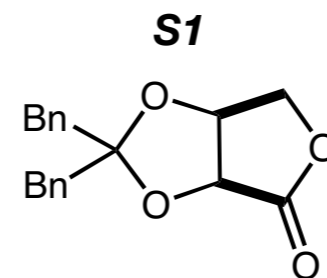
X

7 rounds of catalyst design

Generalization Through Multi-Substrate Optimization in Enantioselective Catalysis

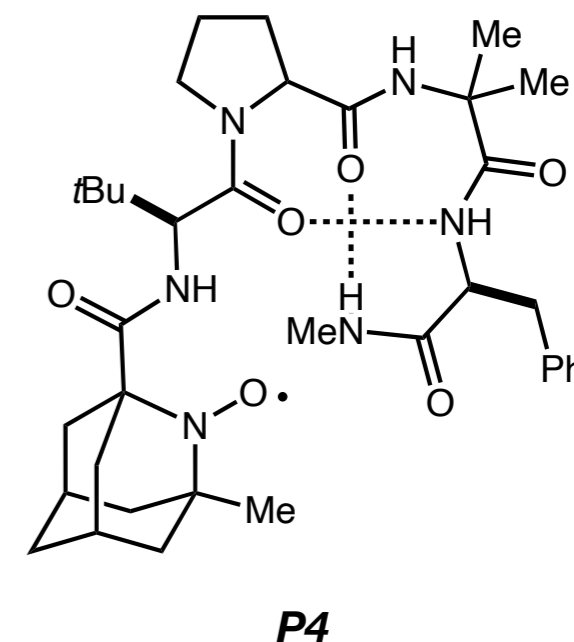
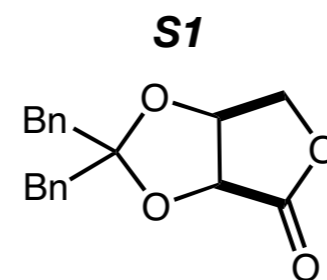
		Catalyst			
%ee		P1	P2	P3	P4
Substrate	S1	75	75	85	94

Under a one substrate optimization mode P4 is an “optimal catalyst”



Generalization Through Multi-Substrate Optimization in Enantioselective Catalysis

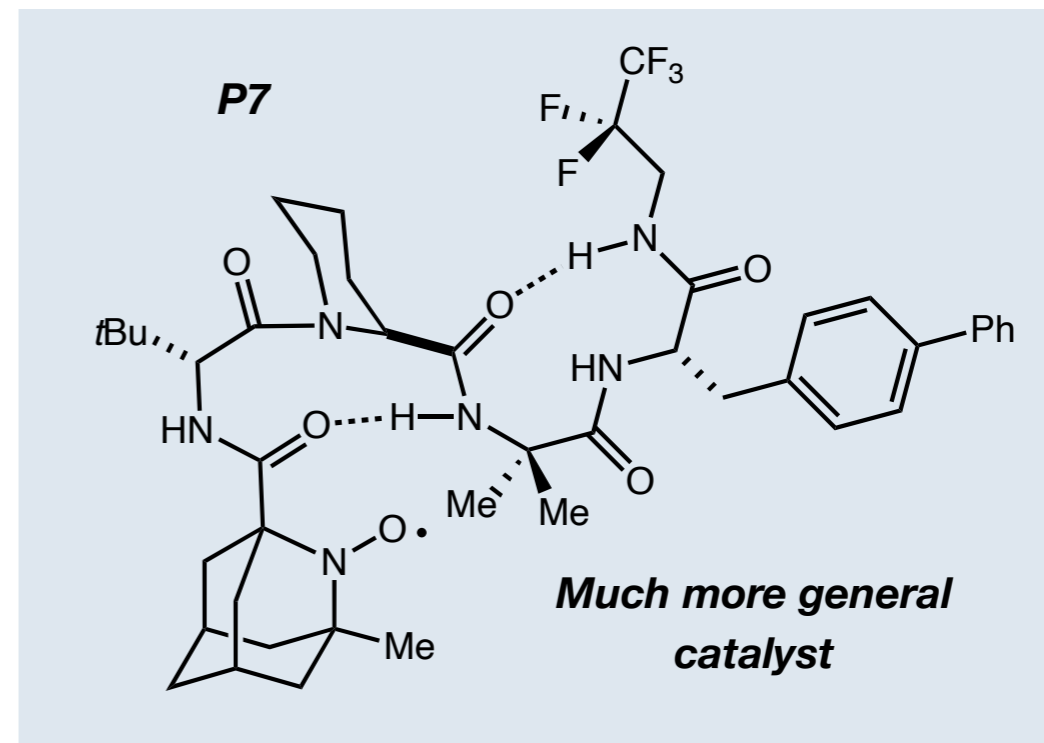
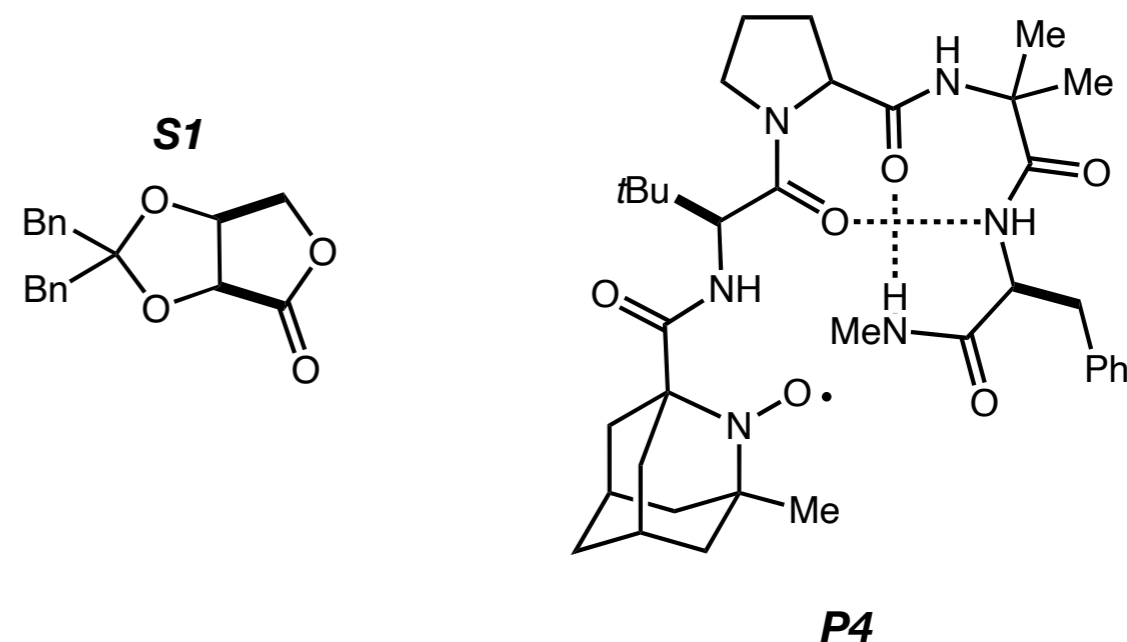
Substrate	Catalyst			
	P1	P2	P3	P4
S1	75	75	85	94
S2	6	2	11	69
S3	19	30	26	75
S4	23	33	1	77
S5	68	72	76	86
S6	36	37	69	87
S7	27	40	4	63
S8	2	13	14	71
S9	8	0	29	69
S10	27	40	28	48
S11	3	2	8	46
S12	7	5	17	22
S13	13	29	6	4
S14	9	34	15	1
S15	16	31	9	1
%<i>ee</i>_{med}	16	31	15	69



Not an optimal catalyst in a multi-substrate optimization paradigm

Generalization Through Multi-Substrate Optimization in Enantioselective Catalysis

		Catalyst						
		P1	P2	P3	P4	P5	P6	P7
Substrate	%ee							
	S1	75	75	85	94	95	95	97
	S2	6	2	11	69	83	86	97
	S3	19	30	26	75	72	93	96
	S4	23	33	1	77	91	90	96
	S5	68	72	76	86	86	95	95
	S6	36	37	69	87	85	96	94
	S7	27	40	4	63	82	83	93
	S8	2	13	14	71	74	89	93
	S9	8	0	29	69	48	79	93
	S10	27	40	28	48	81	68	90
	S11	3	2	8	46	36	76	79
	S12	7	5	17	22	46	51	79
	S13	13	29	6	4	15	40	65
	S14	9	34	15	1	15	4	63
S15	16	31	9	1	20	1	26	
%ee _{med}		16	31	15	69	74	83	93



Additional Reading

Reaction Generalization

Kim, H.; Gerosa, G.; Aronow, J.; Kasaplar, P.; Ouyang J.; Lingau, J .B.; Guerry, P.; Fares, C.; List, B. *Nat. Comm.* **2019**, *10*.

Betinol, I. O.; Lai, J.; Thakur, S.; Reid, J. P. *J. Am. Chem. Soc.* **2023**, *145*, 12870–12883

Rana, D.; Pfluger, P. M.; Holter, N. P.; Tan, G.; Glorius, F. *ACS Cent. Sci.* **2024**, *10*, 899–906

Wagen, C. C.; McMinn, S. E.; Kwan, E. E.; Jacobsen, E. N. *Nature* **2022**, *610*, 680–686

Miniaturization

Gorski, B.; Rein, J.; Norris, S.; Ji, Y.; McEuen, P. L.; Lin, S. *Nature*, **2025**, *637*, 354–361

Gesmundo, N.; Dykstra, K.; Douthwaite, J. L.; Kao, Y.; Zhao, R.; Mahjour, B.; Ferguson, R.; Dreher, S. Sauvagnat, B.; Sauri, J.; Cernak, T. *Nature Synthesis*, **2023**, *2*, 1082–1091

Esguevillas, M.; Fernandez, D. F.; Rincon, J. A.; Barberis, M.; Frutos, O. D.; Mateos, C.; Cerredá, S.; Agejas, J.; MacMillan, D. W. C. *J. Am. Chem. Soc.* **2021**, *7*, 1126–1134

Chemputer/Accelerated Serendipity

Robbins, D. W.; Hartwig, J. F. *Science*, **2011**, *333*, 1423–1427

Granda, J M.; Donina, L.; Dragone, V.; Long, D. Cronin, L. *Nature*, **2018**, *559*, 377–381

Collins, K. D.; Mensch, T.; Glorius, F. *Nat. Chem.* **2014**, *6*, 859–871

Modern approaches to methods development

Modern Paradigms in Screening

- *Reaction generalization*
- *“Accelerate” Serendipity*
- *Miniaturization of unique reaction set ups*

Data Science

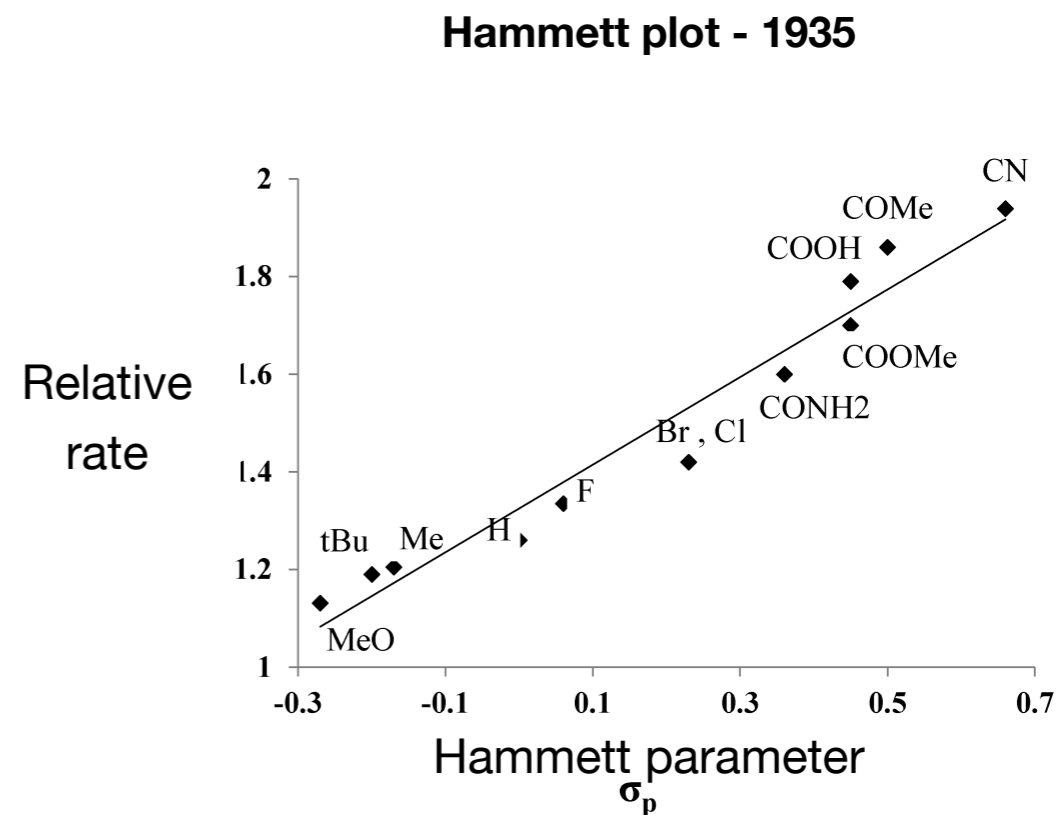
- *Catalyst optimization*
- *Predicting selectivity*
- *Discovery of new catalysts*

Machine Learning

- *What is machine learning*
- *Prediction of optimal conditions*
- *Selectivity prediction for complex systems*
- *Catalyst Discovery*

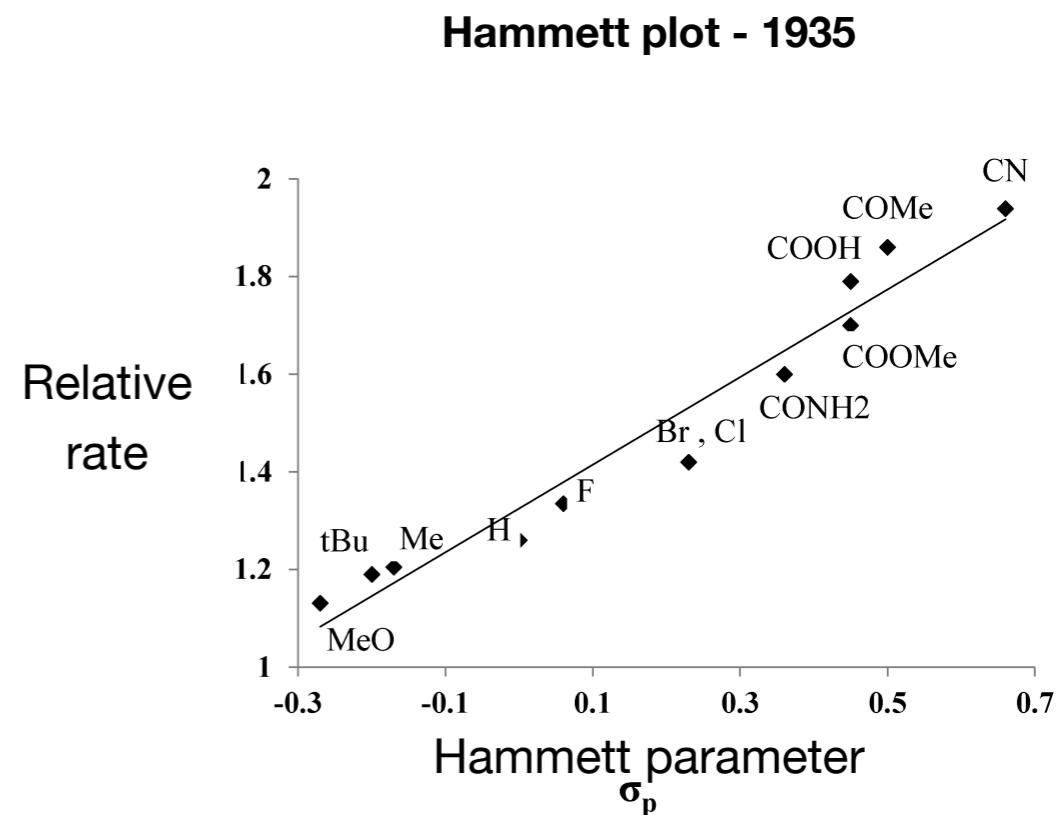
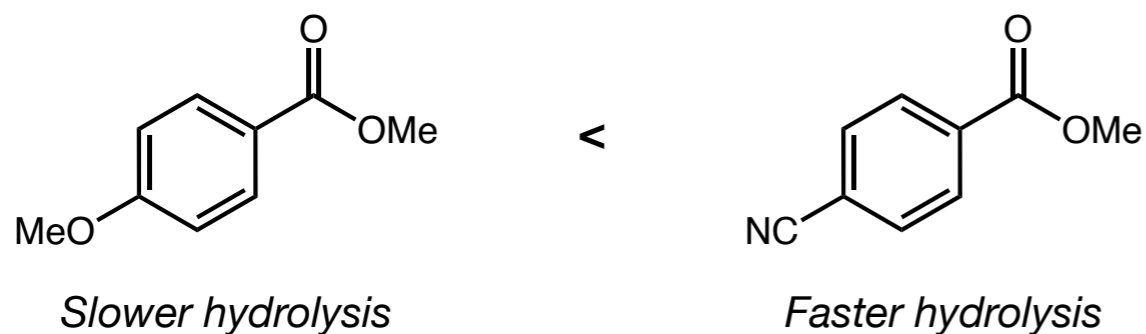
Linear Free Energy Relationships

Correlation of data to make
mechanistic conclusions



Linear Free Energy Relationships

Correlation of data to make
mechanistic conclusions



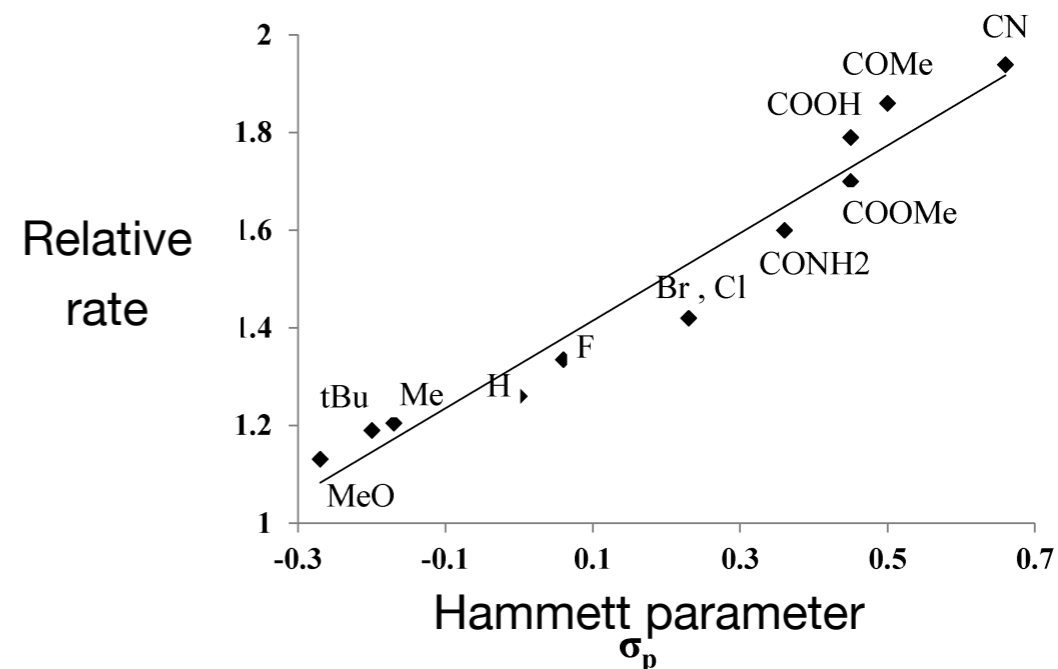
One of the most widely used “mechanistic probes” with **~66,600**
publications mentioning **Hammett plots**

Linear Free Energy Relationships

Correlation of data to make
mechanistic conclusions



Hammett plot - 1935



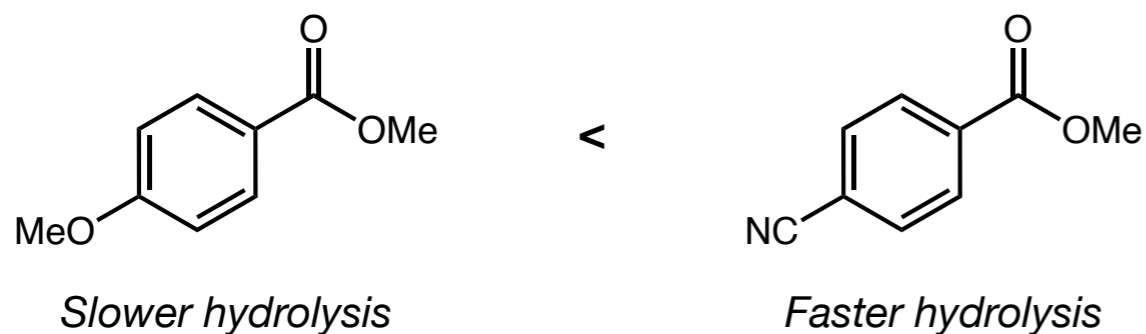
Basic Hammett equation

$$\log\left(\frac{K}{K_0}\right) = \sigma\rho$$

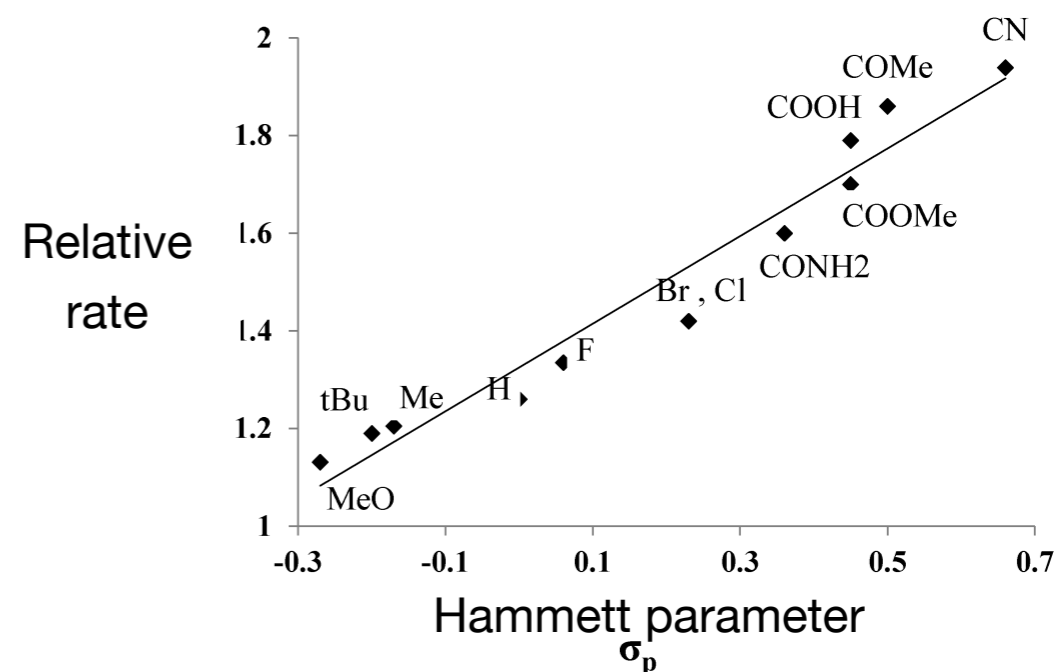
Linear free-energy relation
ship between relative rate
and hammett parameter

Linear Free Energy Relationships

Correlation of data to make
mechanistic conclusions



Hammett plot - 1935



Basic Hammett equation

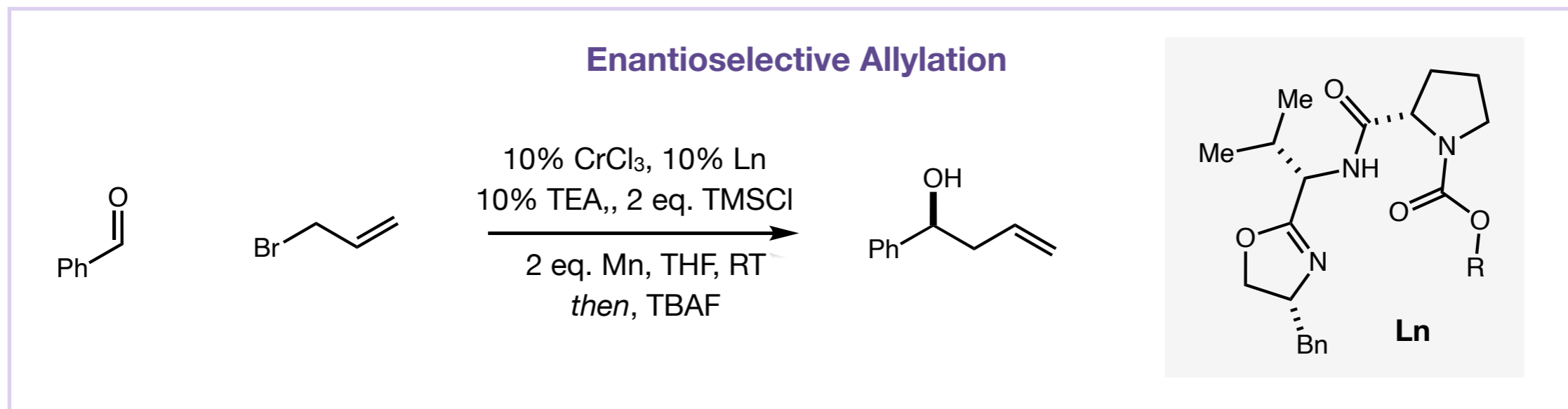
$$\log\left(\frac{K}{K_0}\right) = \sigma\rho$$

Linear free-energy relation
ship between relative rate
and hammett parameter

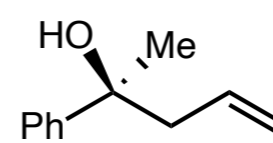
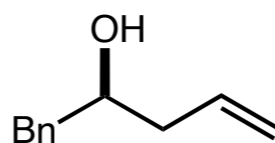
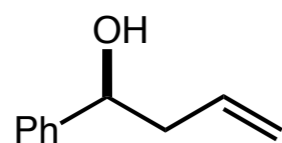
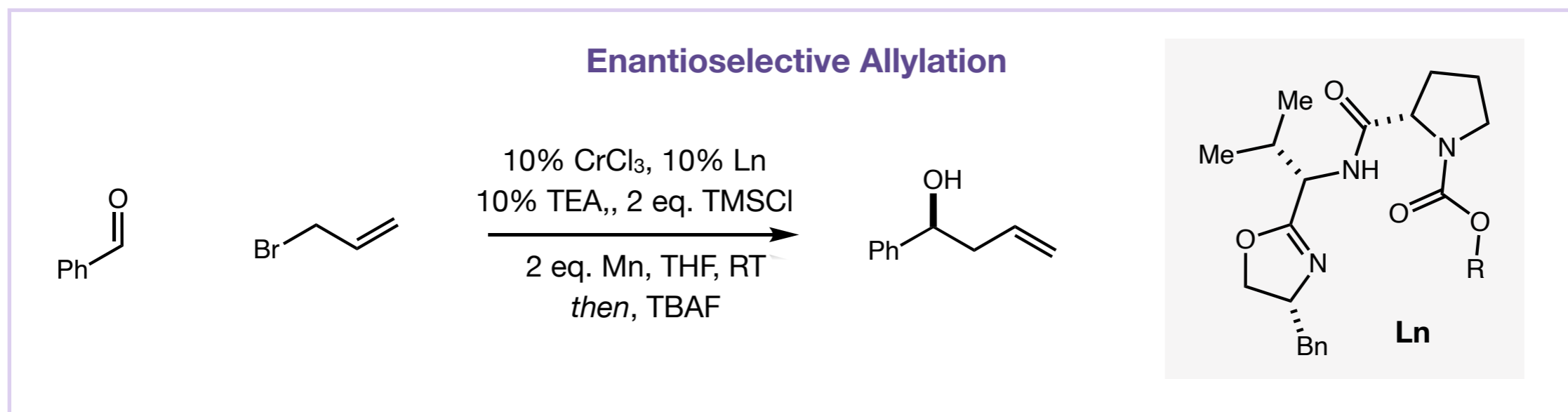
Many other parameters

Sterimol
Charton
Taft
A-values
electric field
Ect.

Linear Free Energy Relationships - Exceptions



Linear Free Energy Relationships - Exceptions



R = Me

1.5 e.r.

1 e.r.

0.3 e.r.

R = Et

1.9 e.r.

1.1 e.r.

0.3 e.r.

R = *i*-Pr

3.5 e.r.

1.3 e.r.

0.6 e.r.

R = *t*-Bu

23 e.r.

2.8 e.r.

2.2 e.r.

R = Ad

23 e.r.

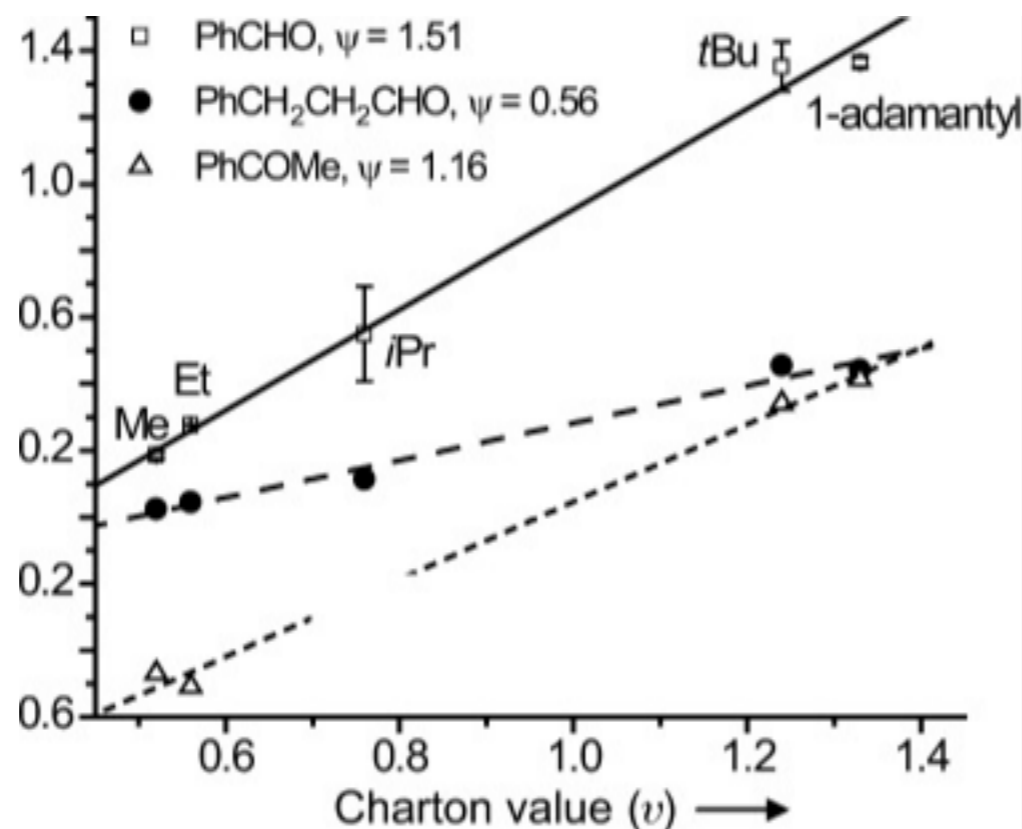
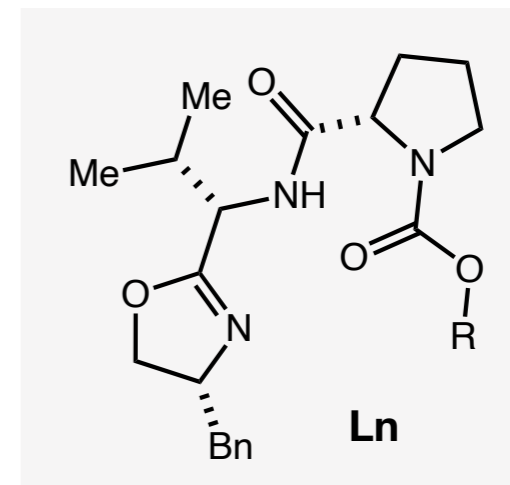
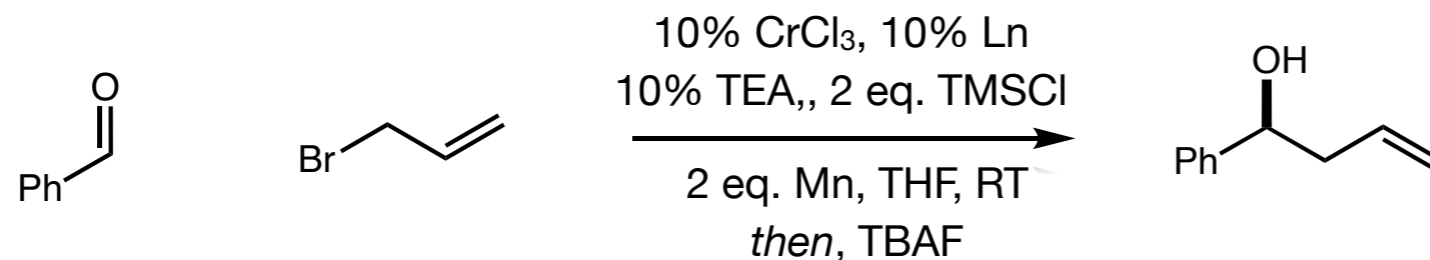
2.8 e.r.

2.6 e.r.

*Linear free-energy relationship
between log (e.r.) and Steric
parameters*

Linear Free Energy Relationships - Exceptions

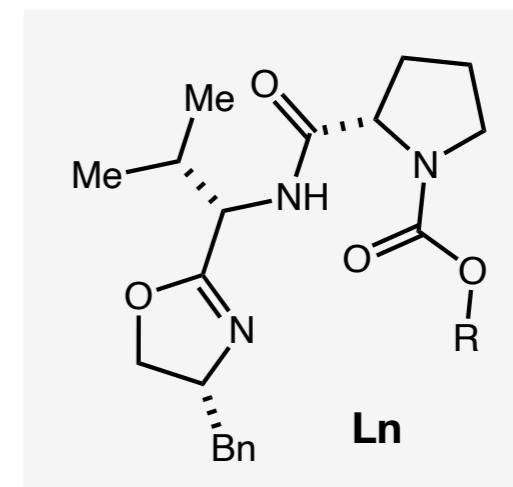
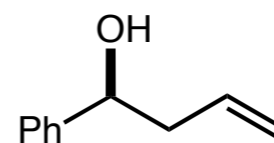
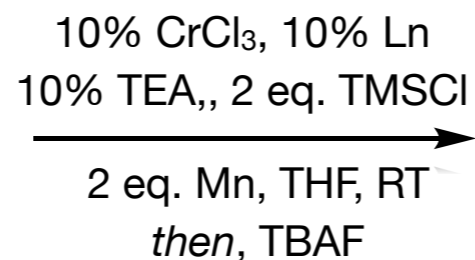
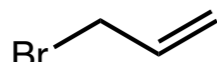
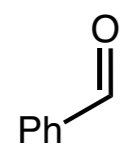
Enantioselective Allylation



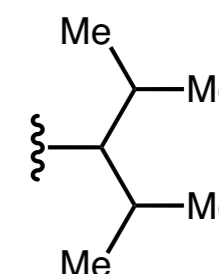
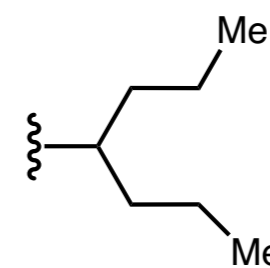
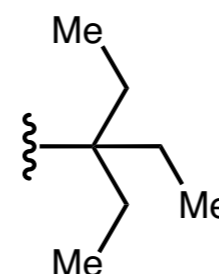
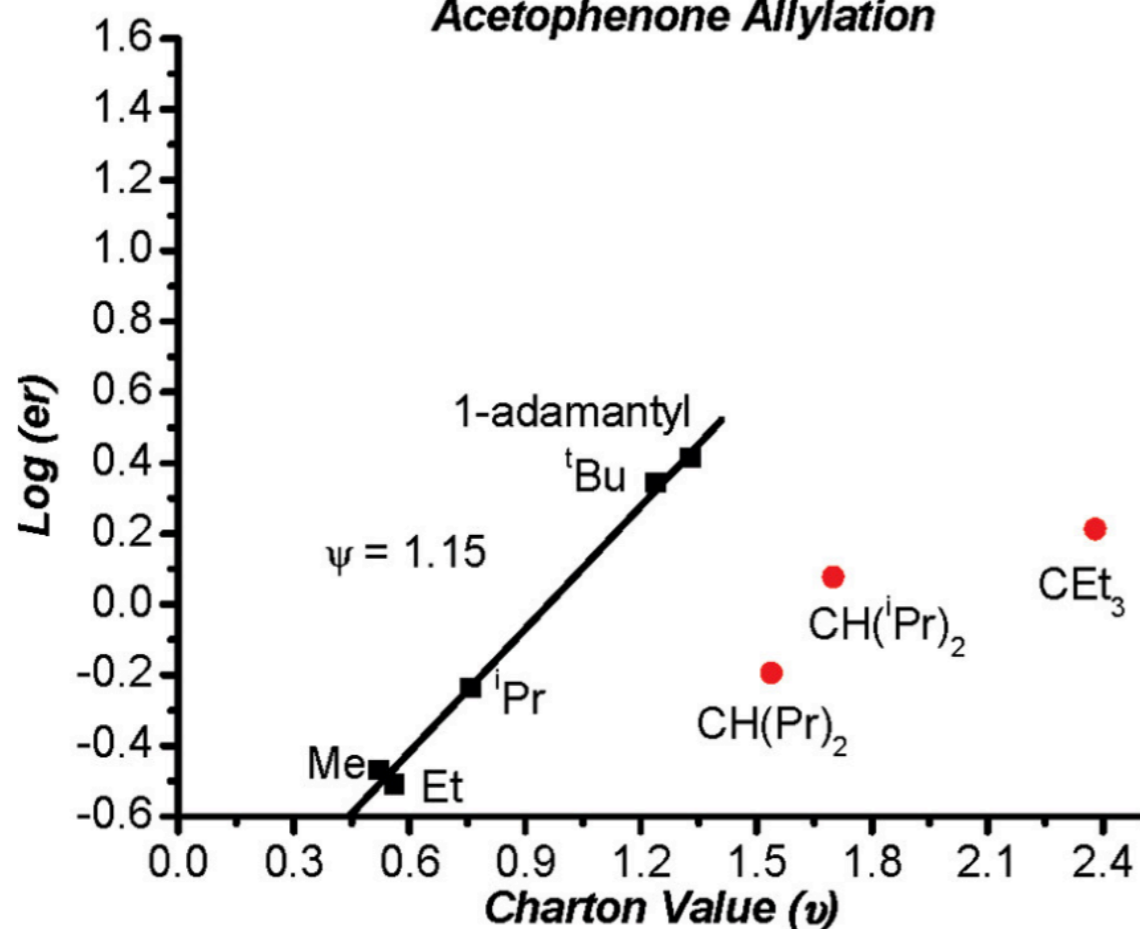
Linear free-energy relationship
between $\log(e.r.)$ and Steric
parameters

Linear Free Energy Relationships - Exceptions

Enantioselective Allylation

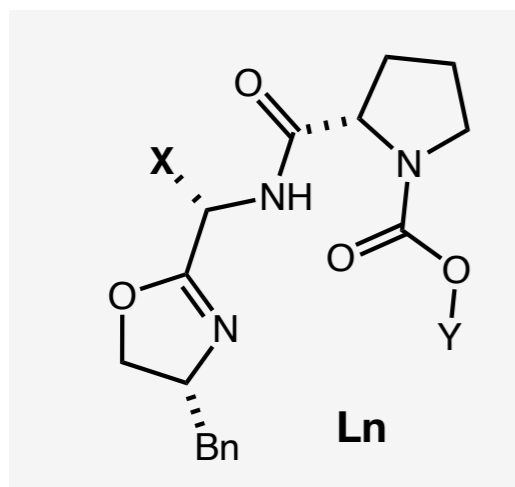


Acetophenone Allylation



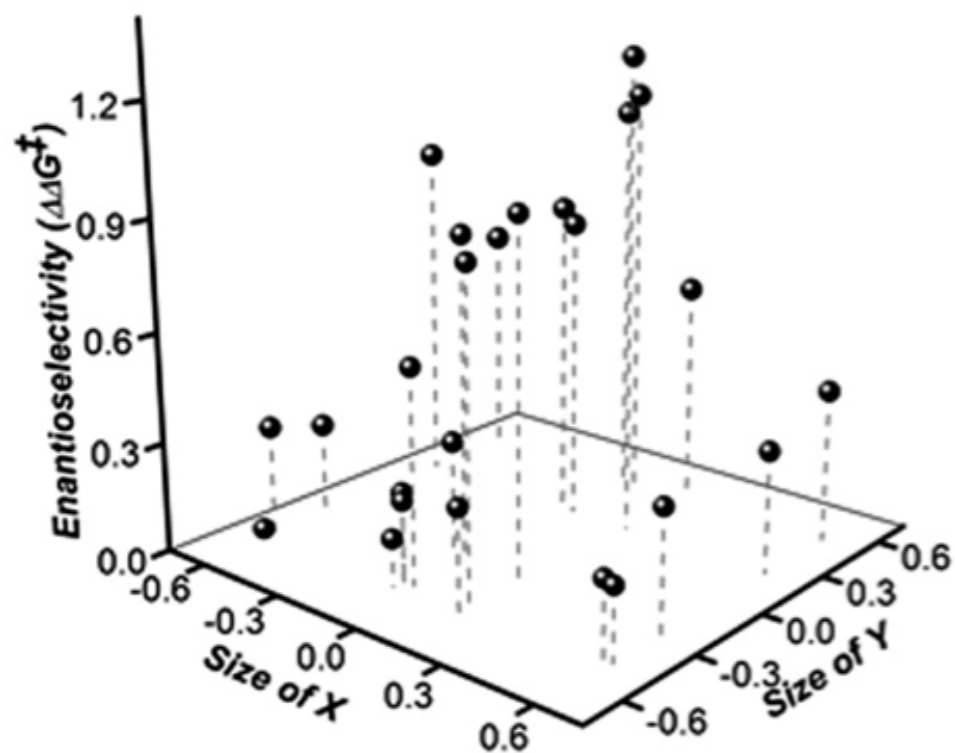
Even bigger Charton Values
give worse e.r., breaks the
linear relationship

Multivariate Linear Free Energy Relationships

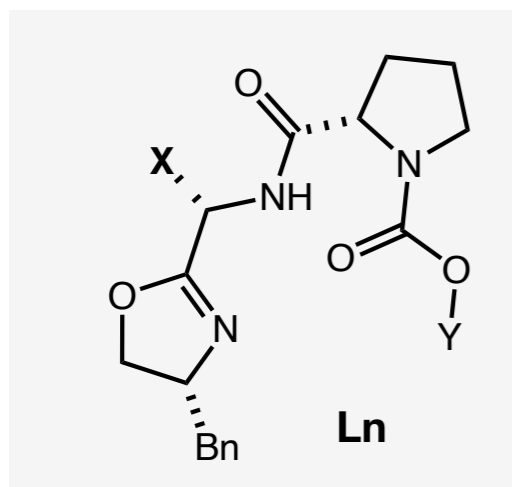


X =	Y =
H	H
Me	Me
Et	iPr
iPr	tBu
tBu	CH(Pr) ₂

25 ligand library to determine relationship
between steric bulk at X and Y

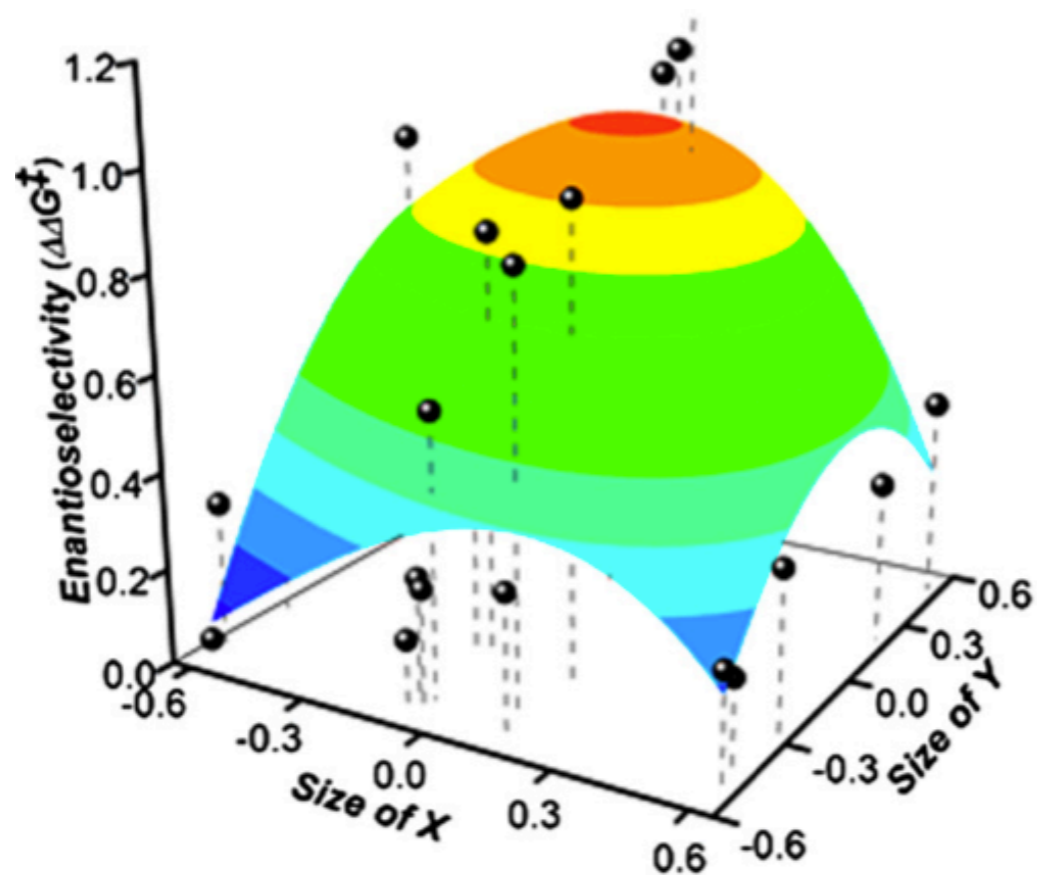


Multivariate Linear Free Energy Relationships

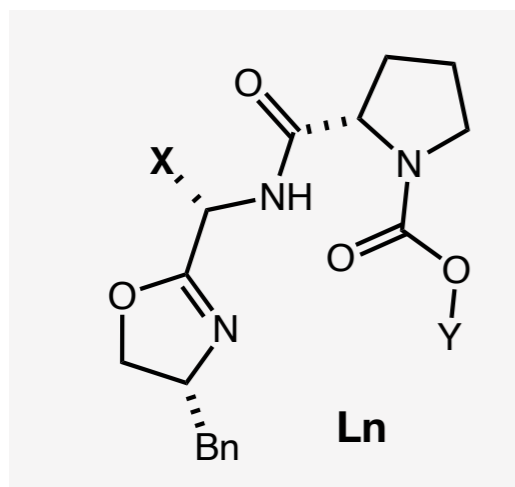


X =	Y =
H	H
Me	Me
Et	iPr
iPr	tBu
tBu	CH(Pr) ₂

25 ligand library to determine relationship
between steric bulk at X and Y

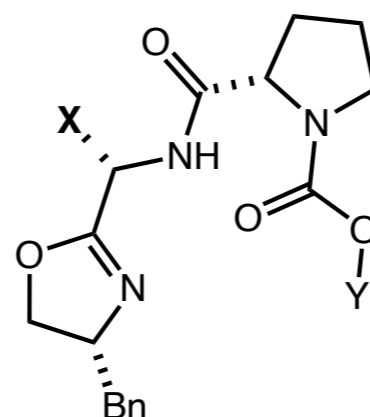
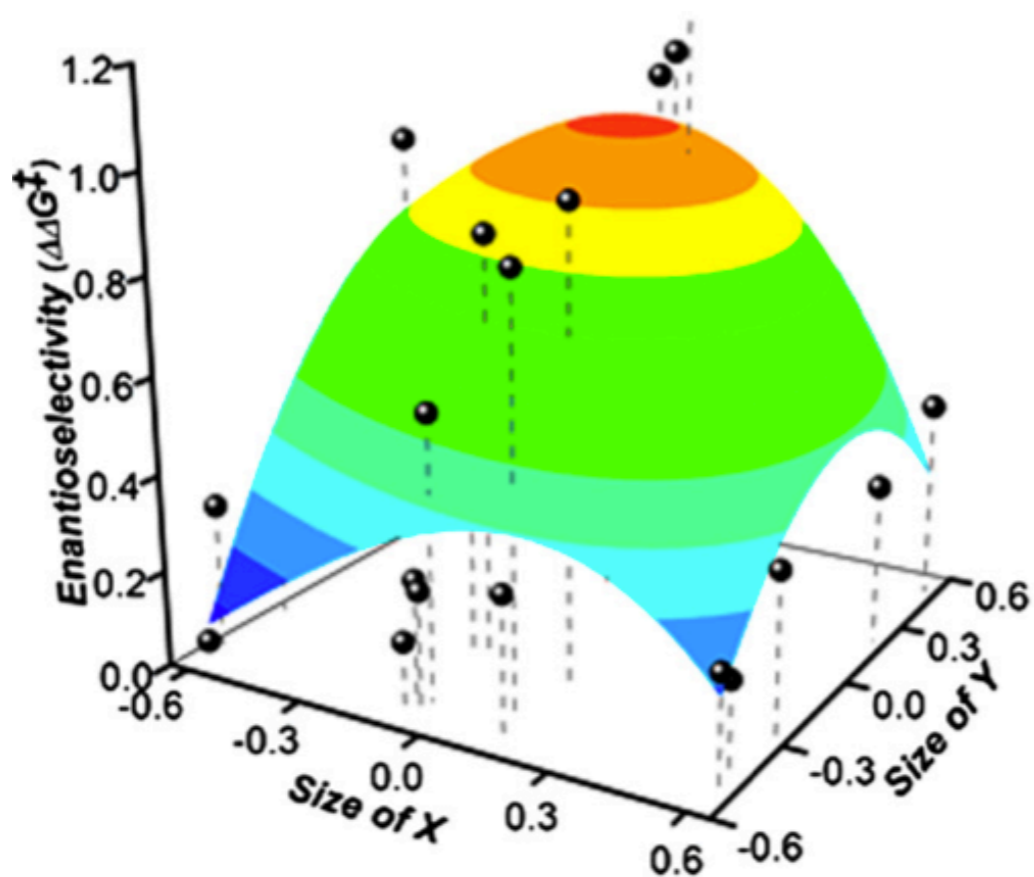


Multivariate Linear Free Energy Relationships



X =	Y =
H	H
Me	Me
Et	iPr
iPr	tBu
tBu	CH(Pr) ₂

25 ligand library to determine relationship between steric bulk at X and Y

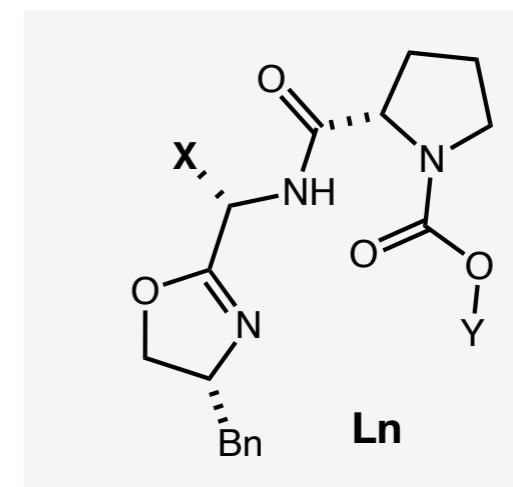
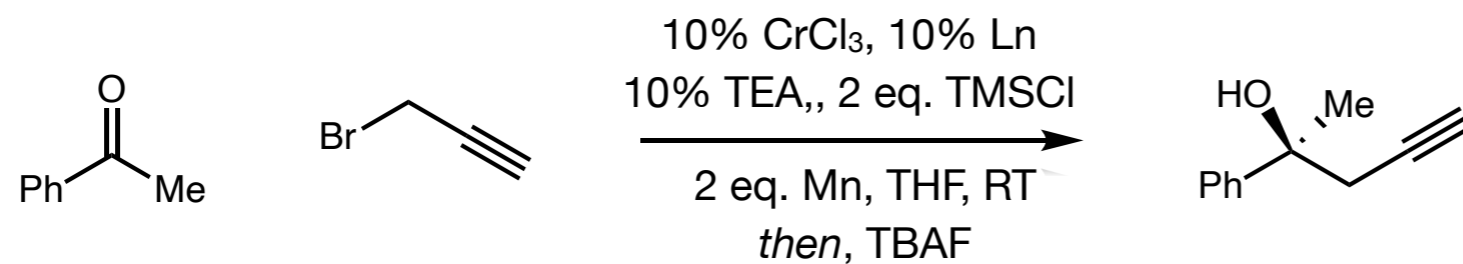


Relatively accurate predictor even for much bigger Y groups

X =	Y =	Predicted e.r.	Measured e.r.
H	Me	49:51	49:51
Me	tBu	44:56	46:54
t-Bu	C(Et) ₃	40:60	42:58
iPr	C(Et) ₃	49:51	54:46

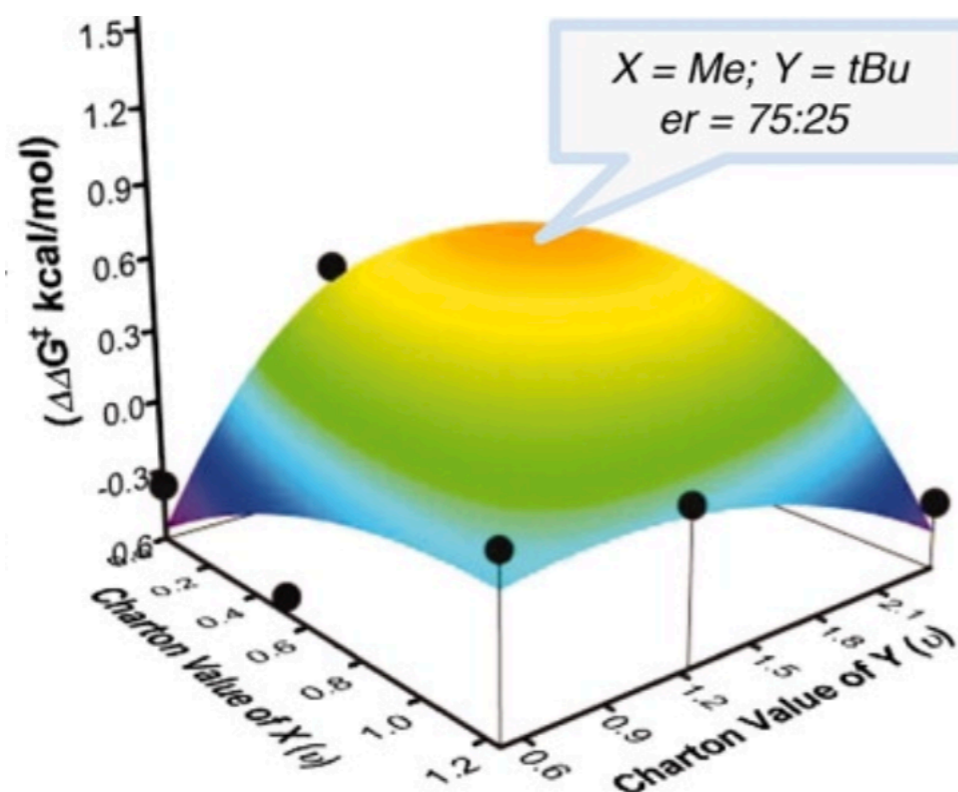
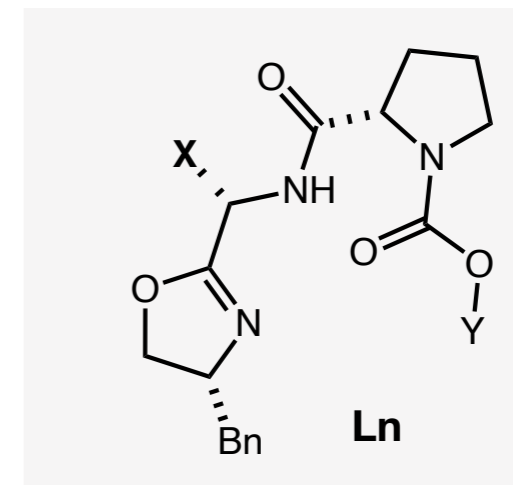
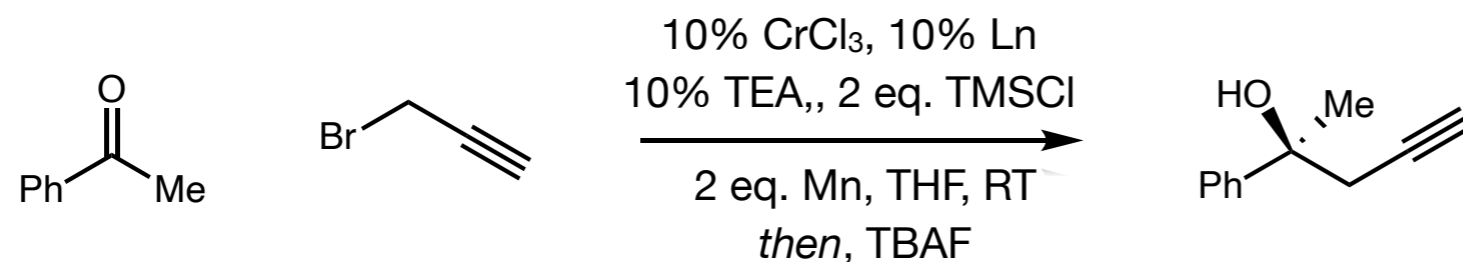
Multivariate Linear Free Energy Relationships

Enantioselective Propargylation



Multivariate Linear Free Energy Relationships

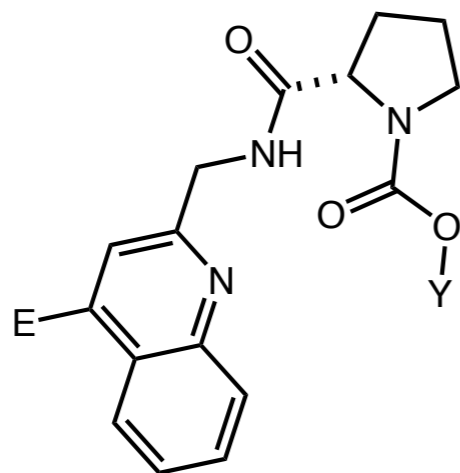
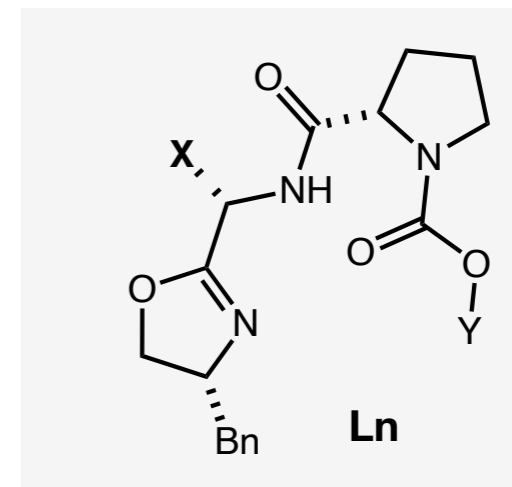
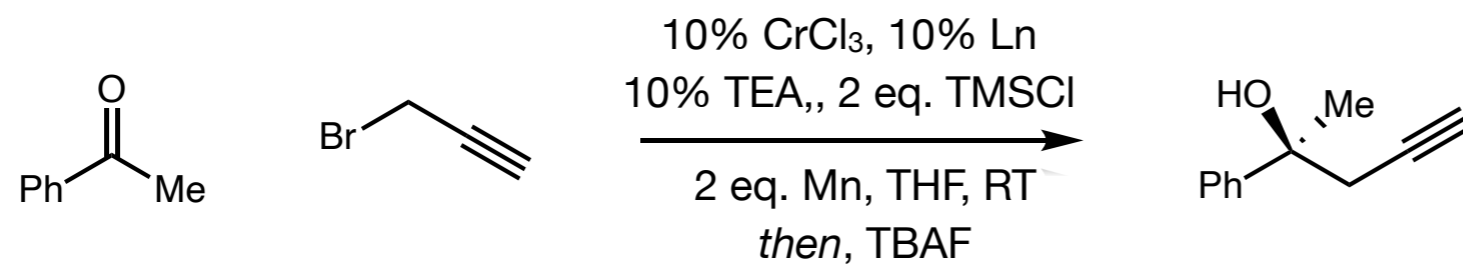
Enantioselective Propargylation



Low e.r. predicted and overall small correlation between e.r. and steric parameters

Multivariate Linear Free Energy Relationships

Enantioselective Propargylation

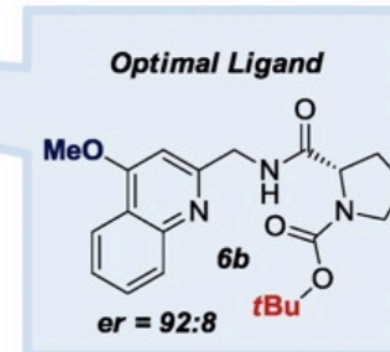
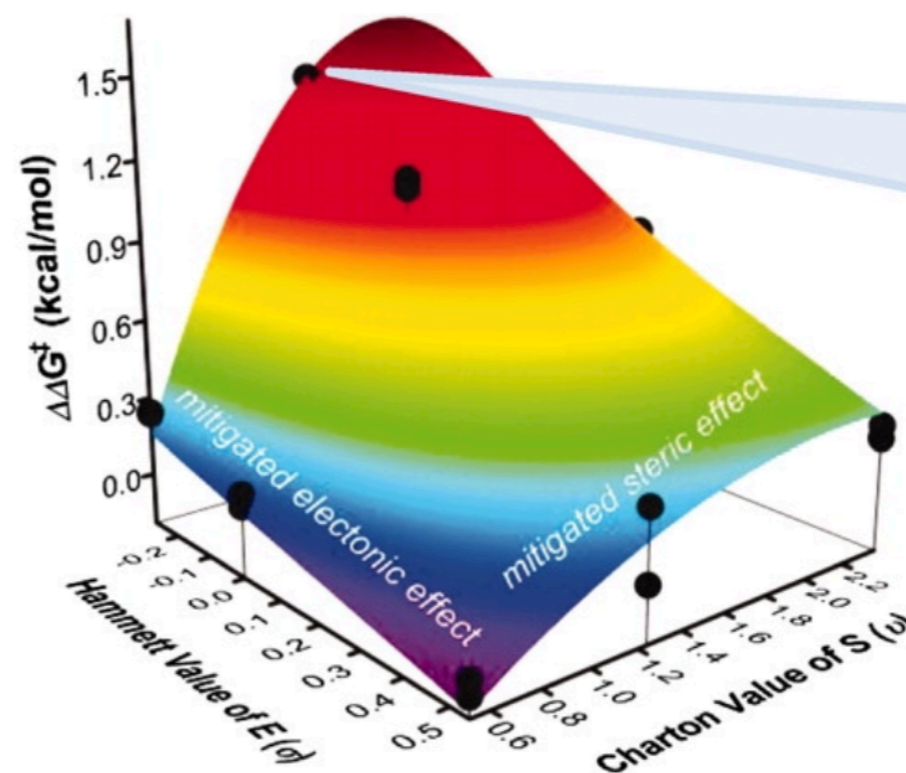
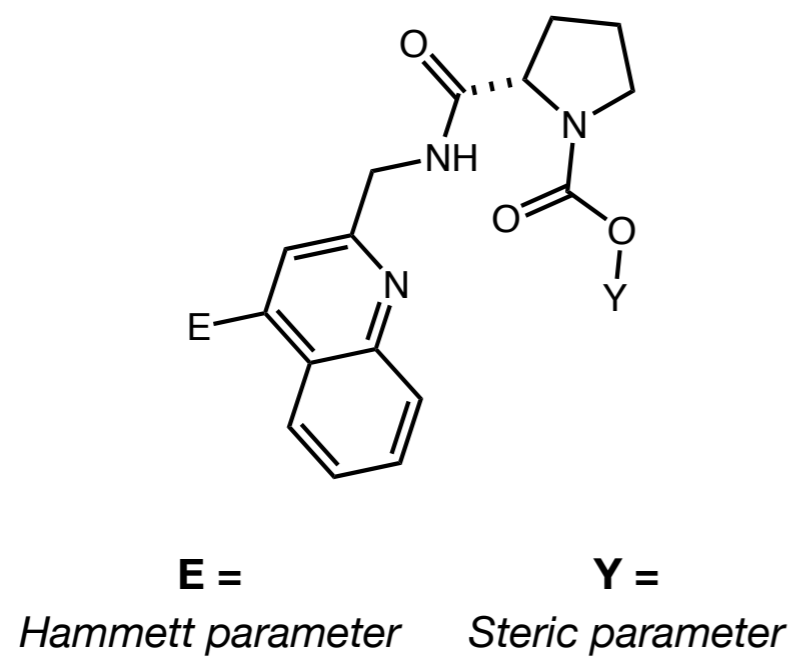
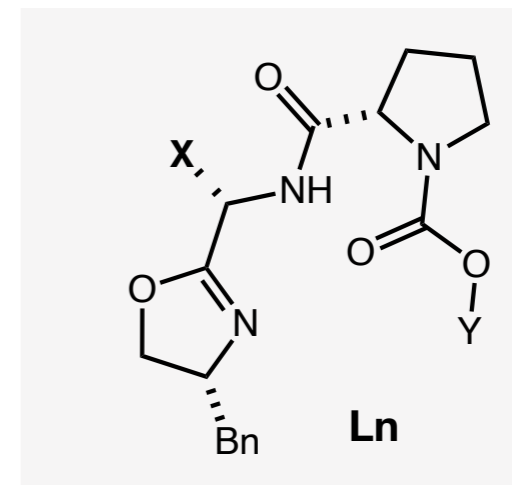
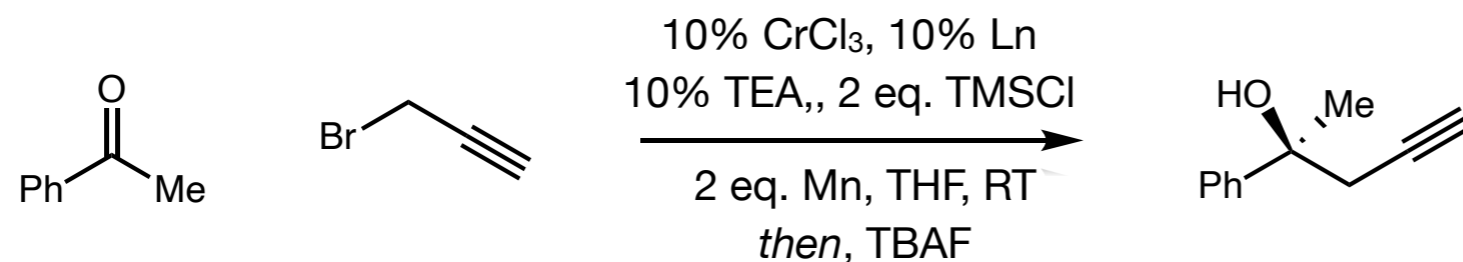


E =
Hammett parameter

Y =
Steric parameter

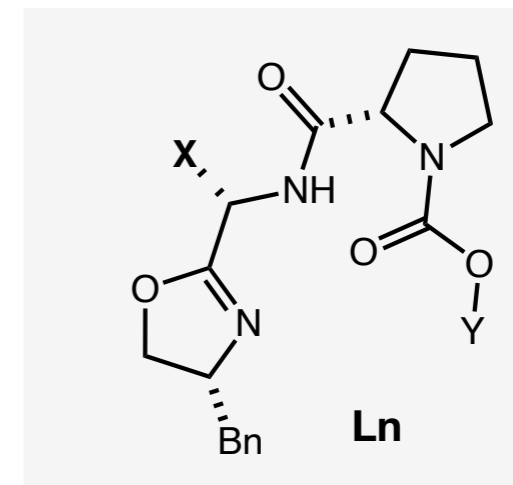
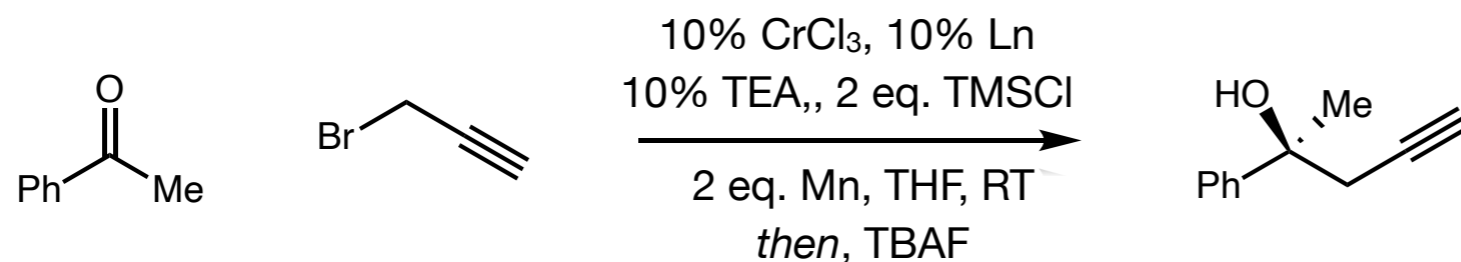
Multivariate Linear Free Energy Relationships

Enantioselective Propargylation

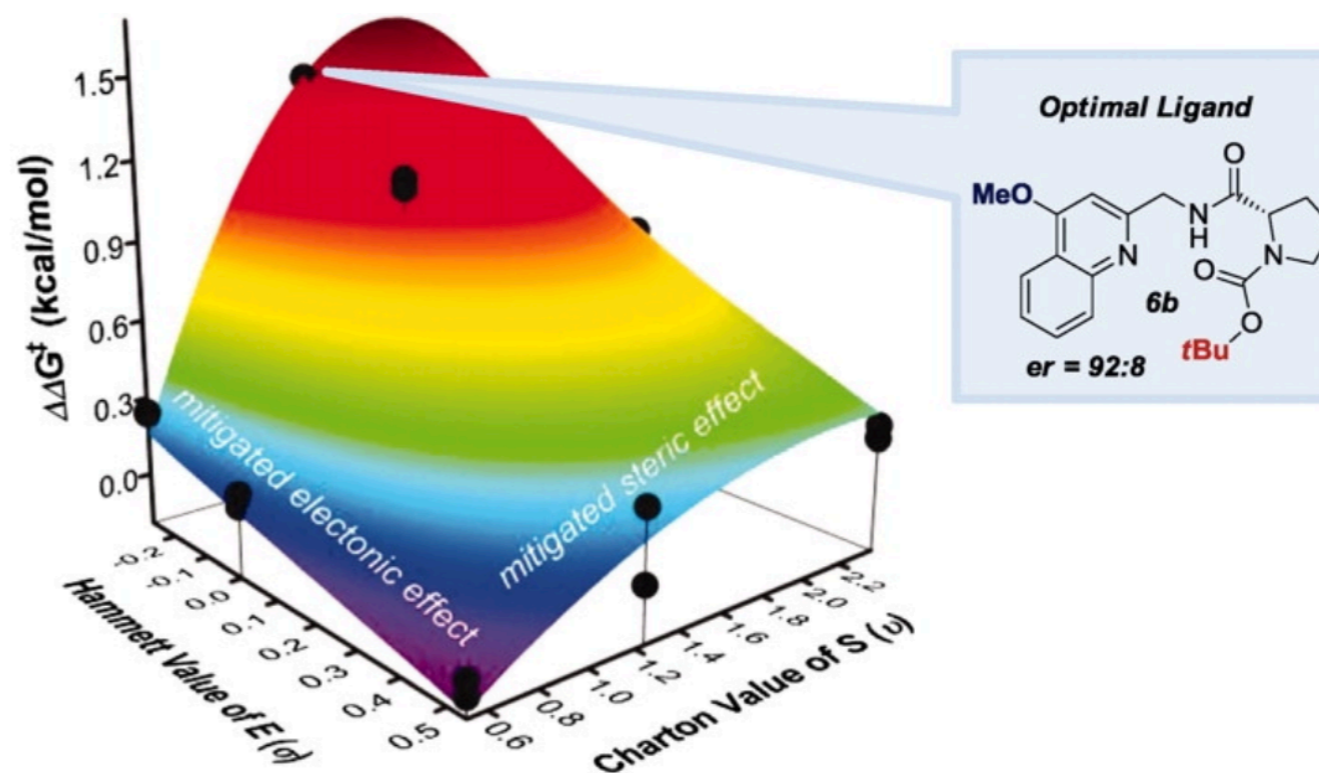
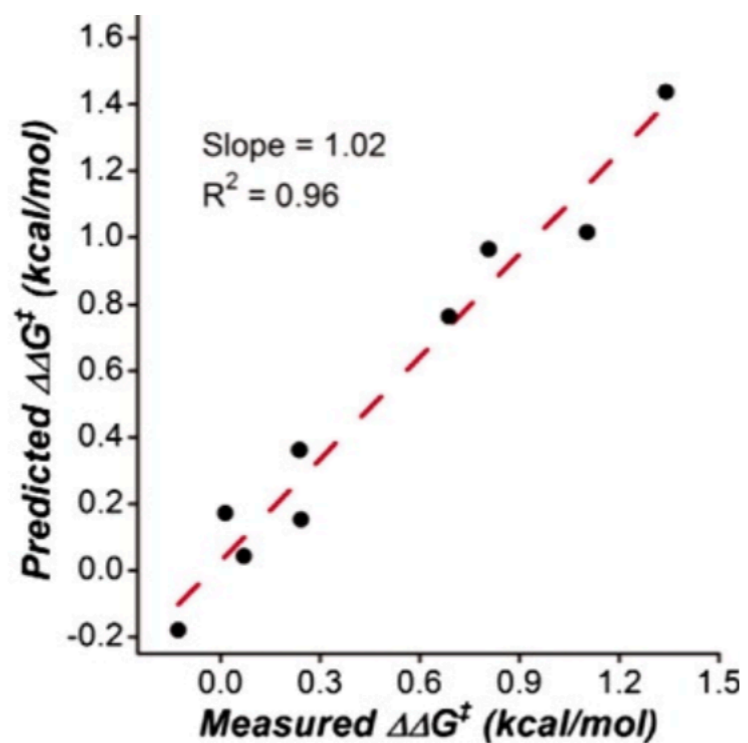


Multivariate Linear Free Energy Relationships

Enantioselective Propargylation

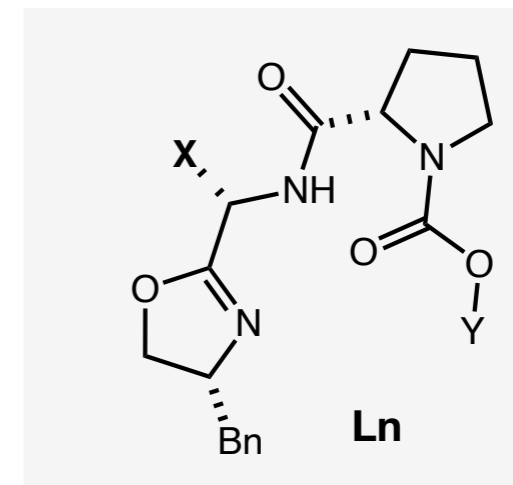
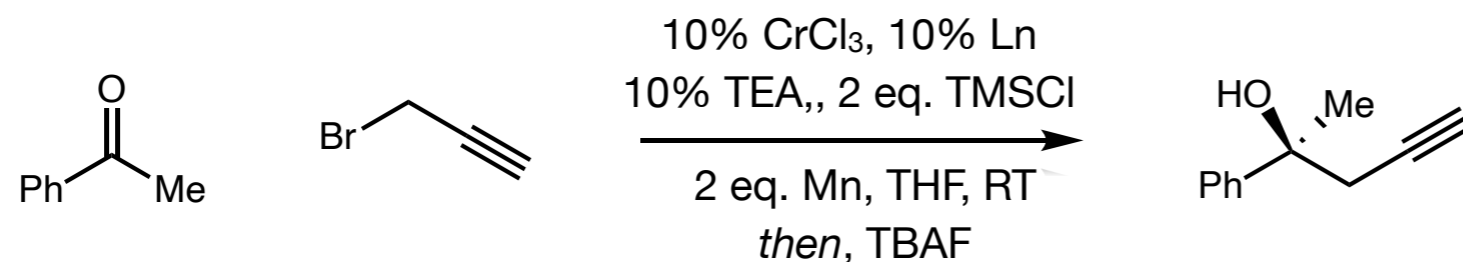


Model gives accurate predictions

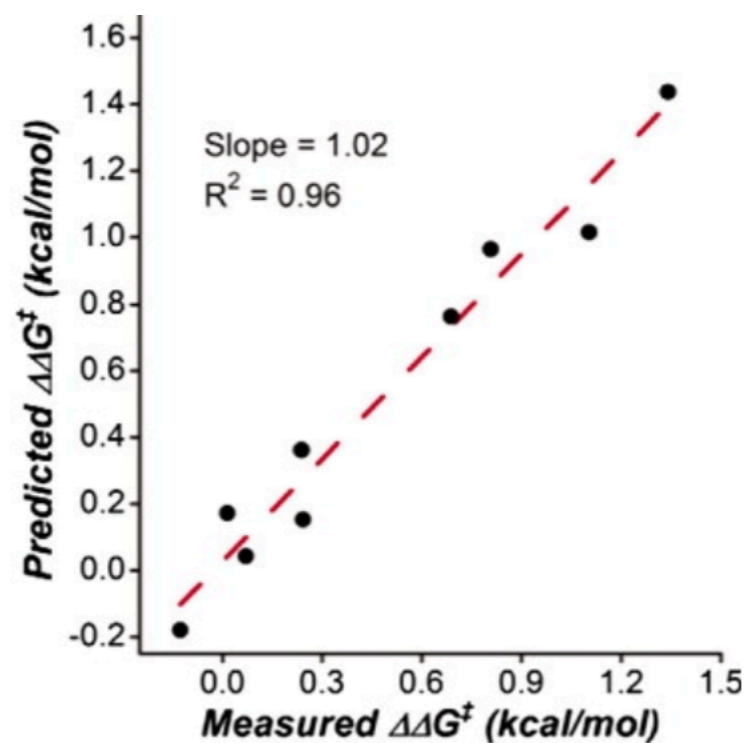


Multivariate Linear Free Energy Relationships

Enantioselective Propargylation



Model gives accurate predictions

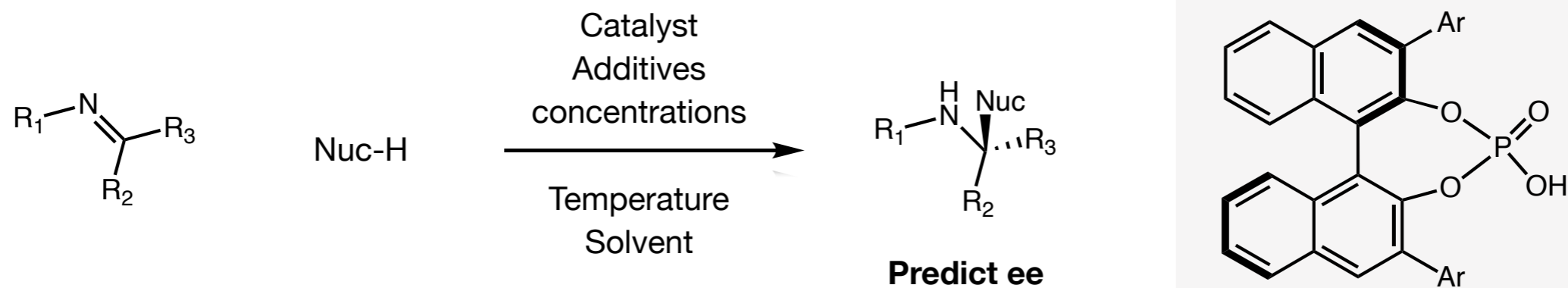


Mathematical formula

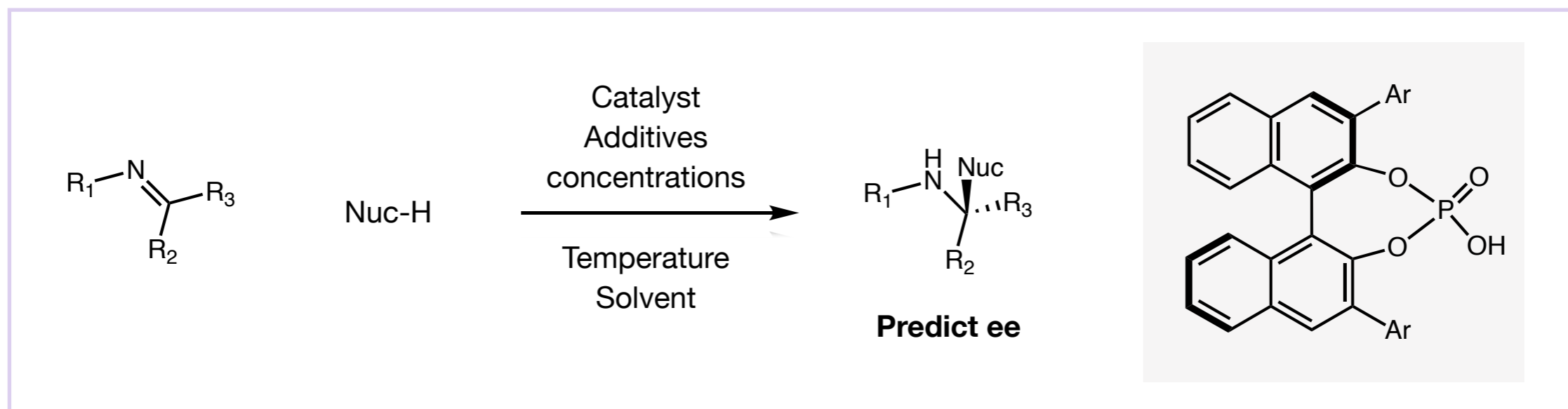
$$\Delta\Delta G^\ddagger = -1.20 + 1.22E + 2.84S - 0.85S^2 - 3.79ES + 1.25ES^2$$

Modeling can be done on more than just 2 parameters

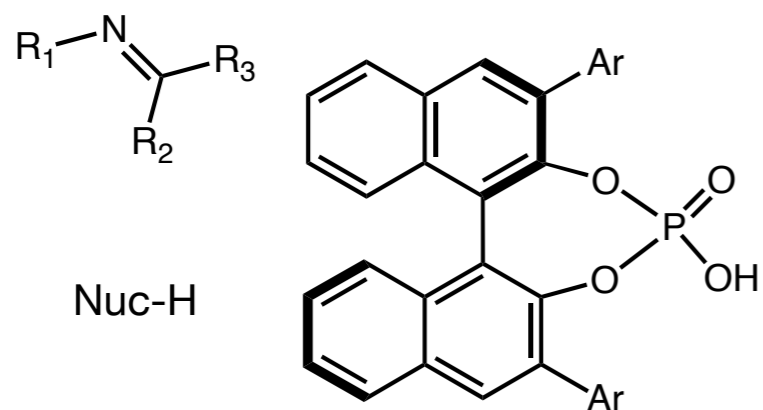
Multivariate Linear Free Energy Relationships - Predicting ee



Multivariate Linear Free Energy Relationships - Predicting ee

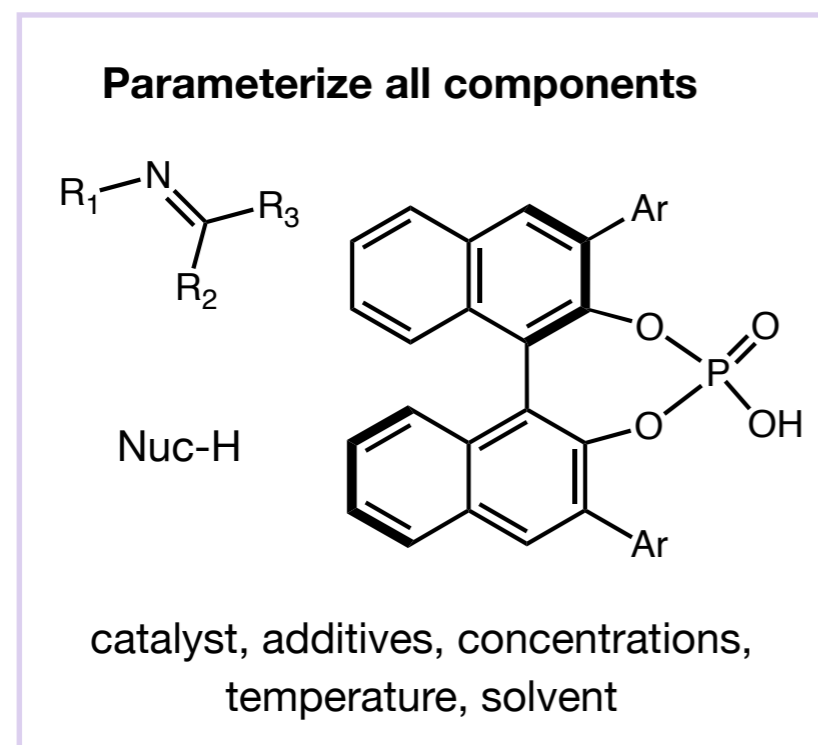
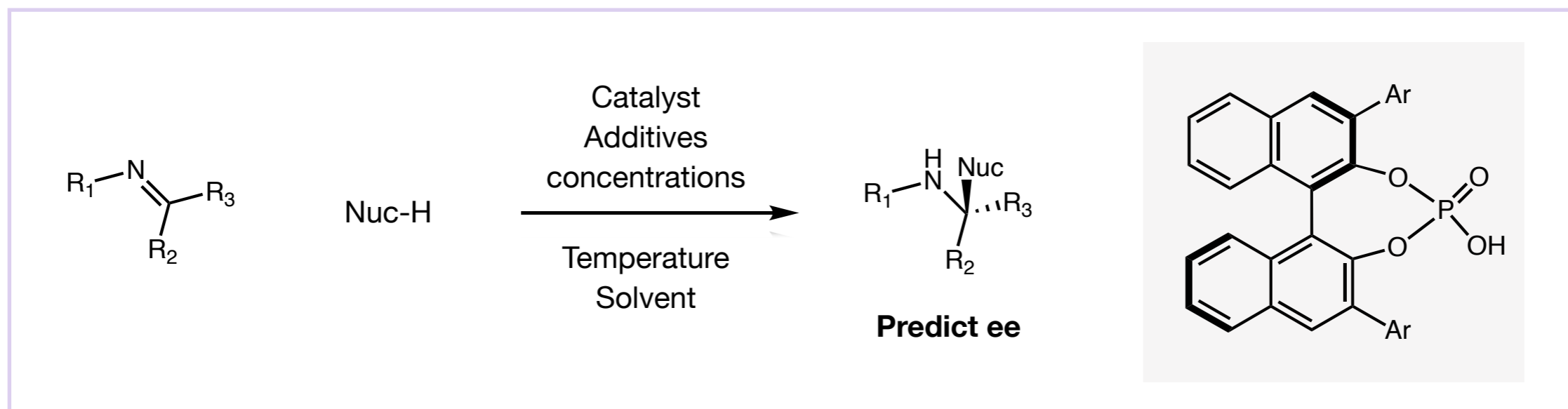


Parameterize all components



catalyst, additives, concentrations,
temperature, solvent

Multivariate Linear Free Energy Relationships - Predicting ee

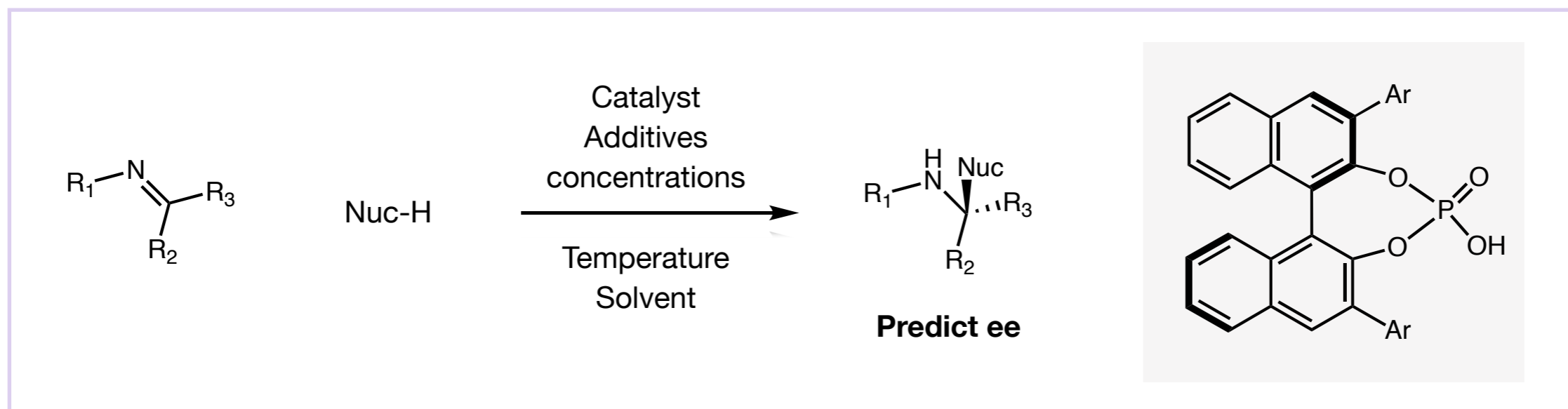


Gather data on a
training set

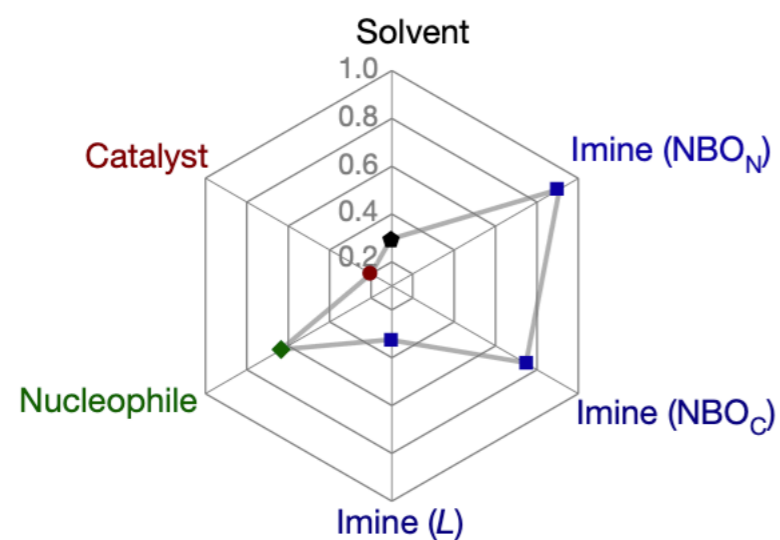
Analyze and model

1. Determine the parameters that have the highest effect on ee
2. Derive a mathematical formula that predicts ee based on the important parameters

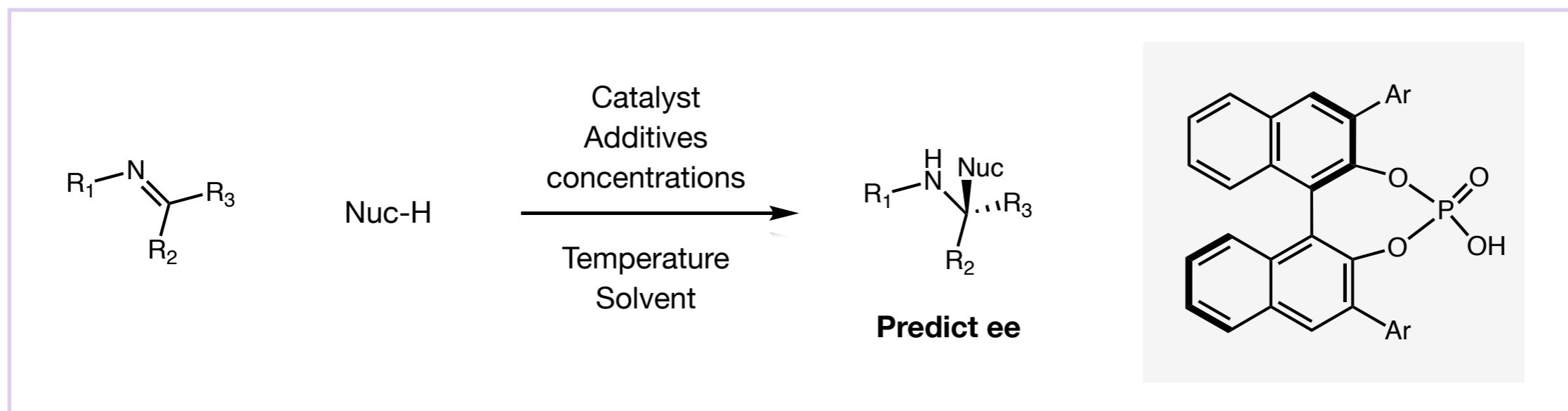
Multivariate Linear Free Energy Relationships - Predicting ee



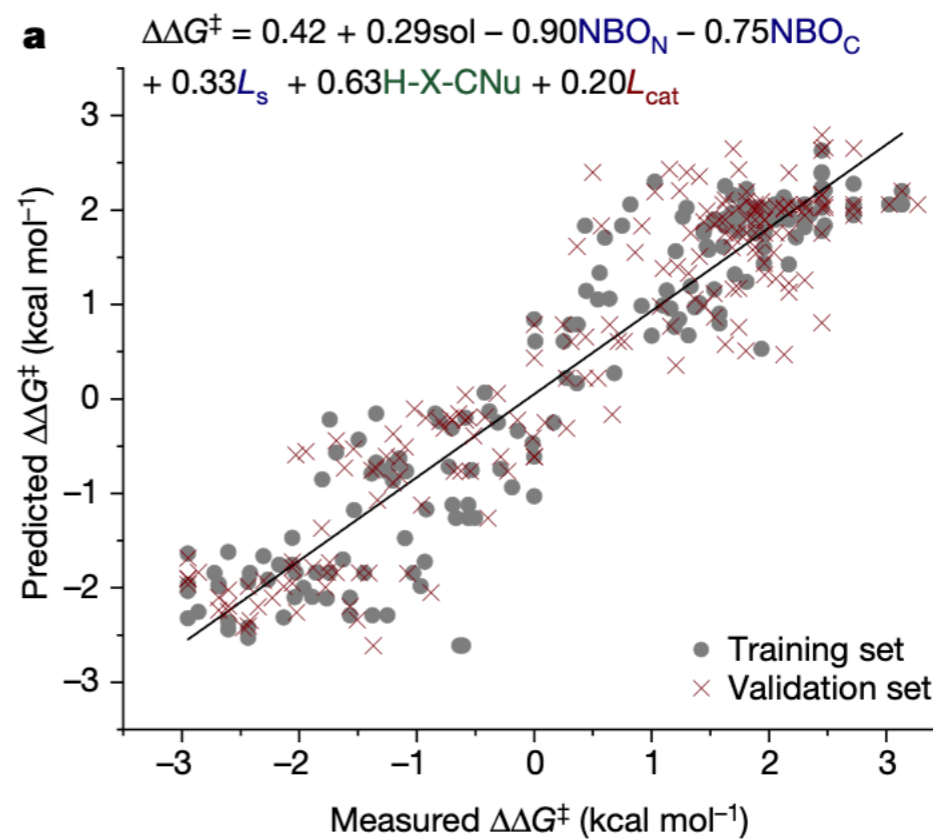
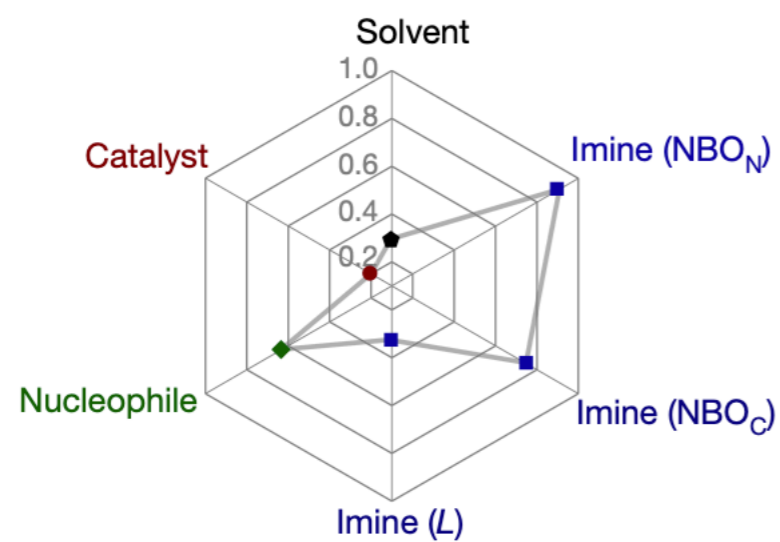
Determine and weigh
important parameters



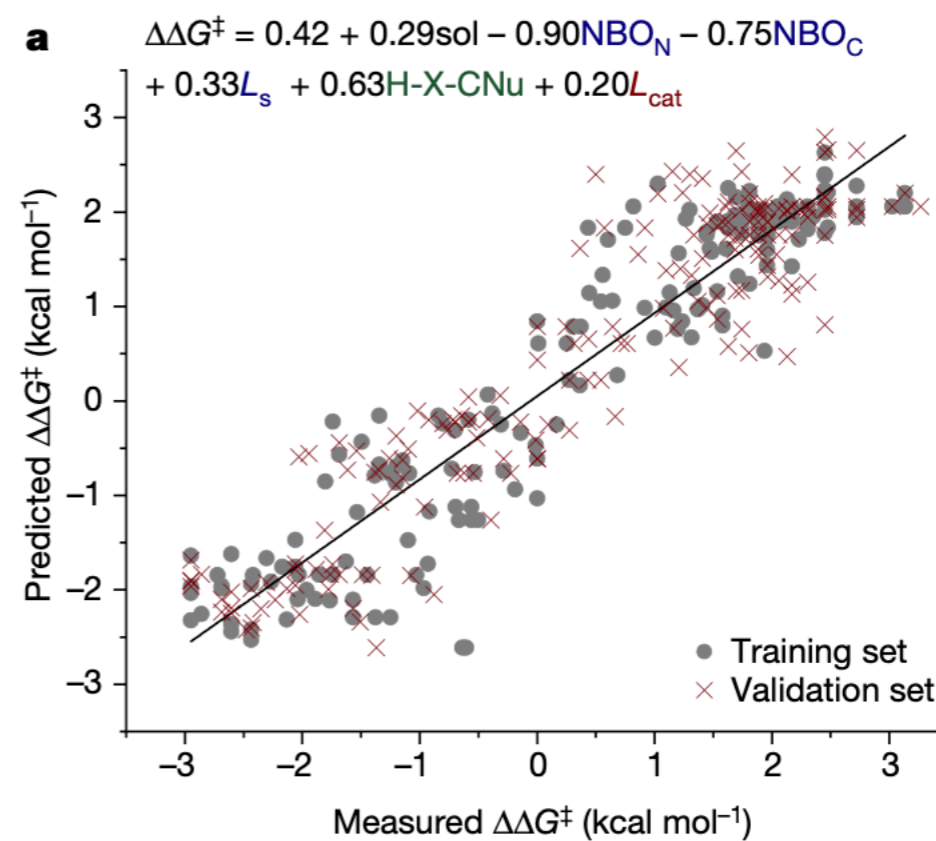
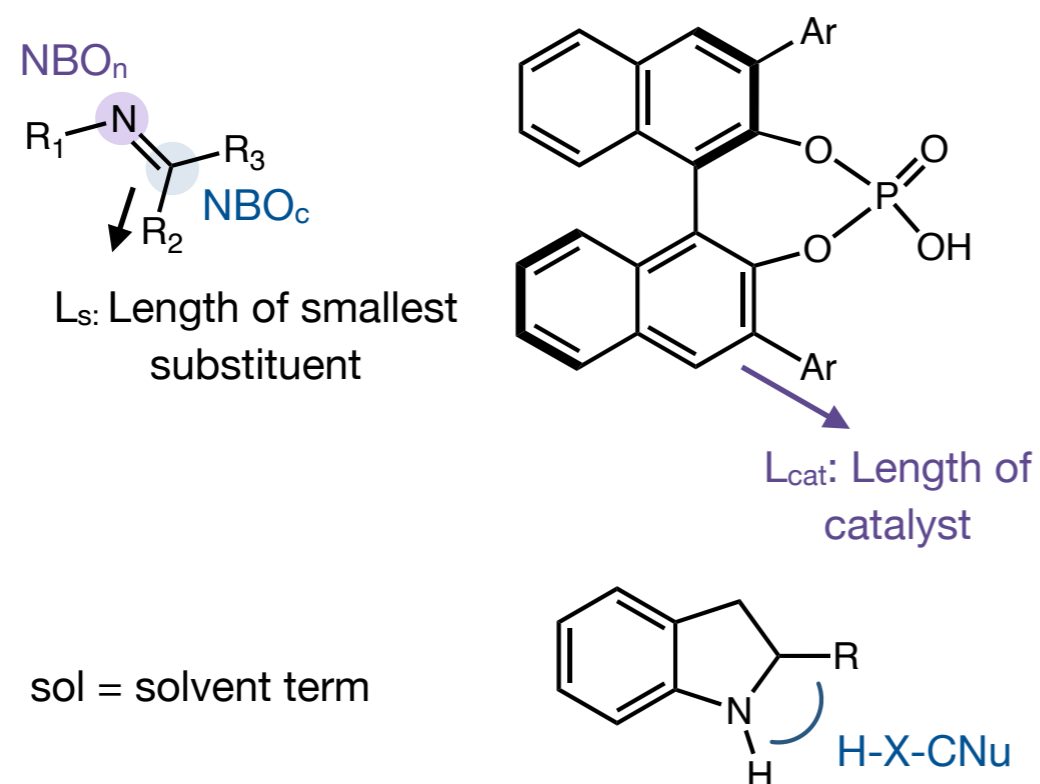
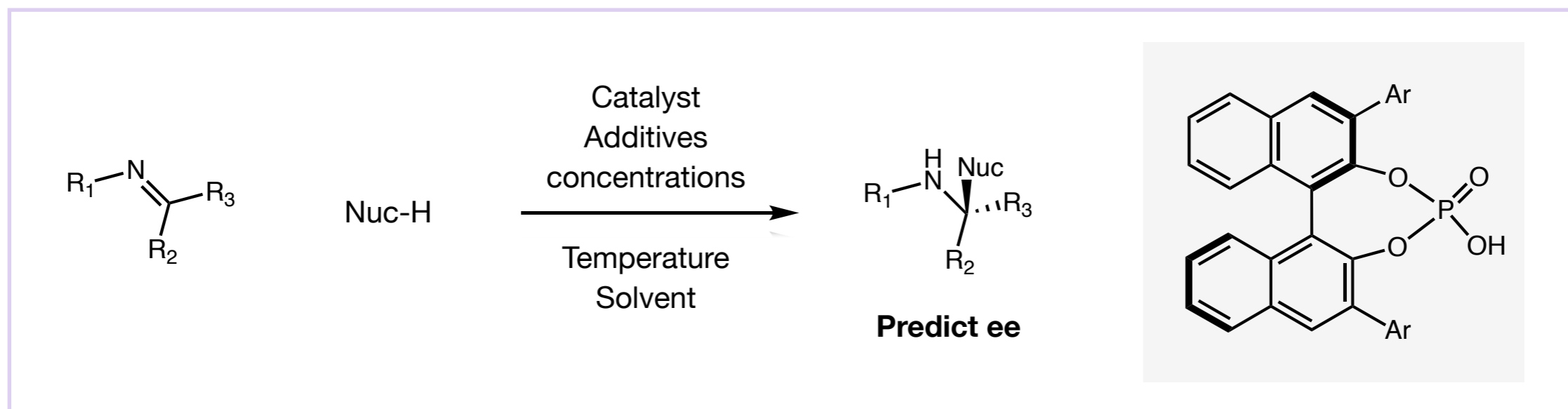
Multivariate Linear Free Energy Relationships - Predicting ee



Determine and weigh
important parameters

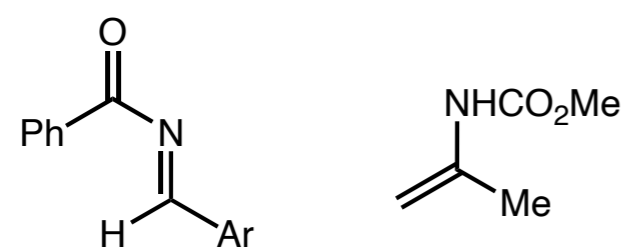


Multivariate Linear Free Energy Relationships - Predicting ee

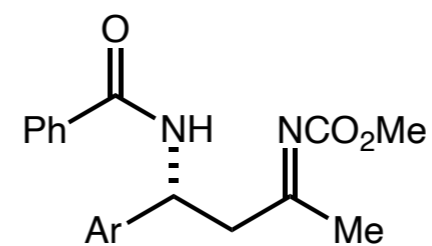
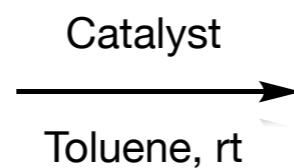


Multivariate Linear Free Energy Relationships - Predicting ee

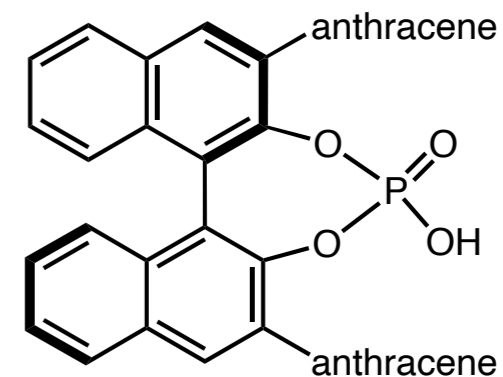
Testing the accuracy of the model with out-of-sample test substrates



15 examples

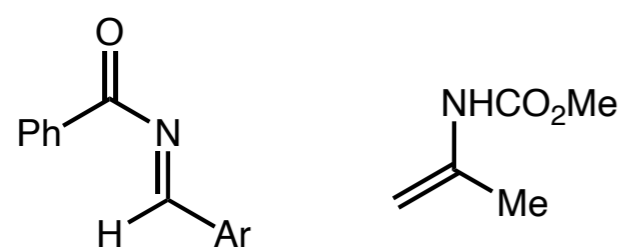


All predicted within 5% ee

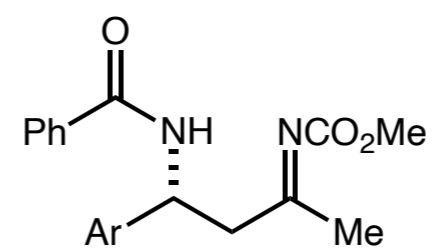


Multivariate Linear Free Energy Relationships - Predicting ee

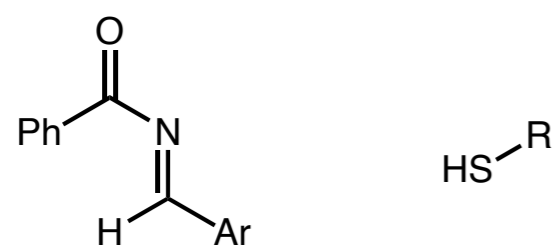
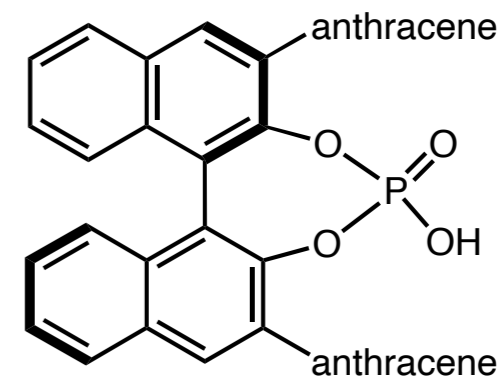
Testing the accuracy of the model with out-of-sample test substrates



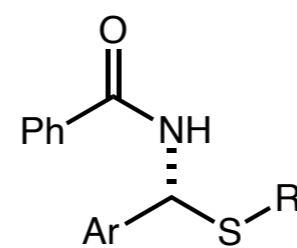
15 examples



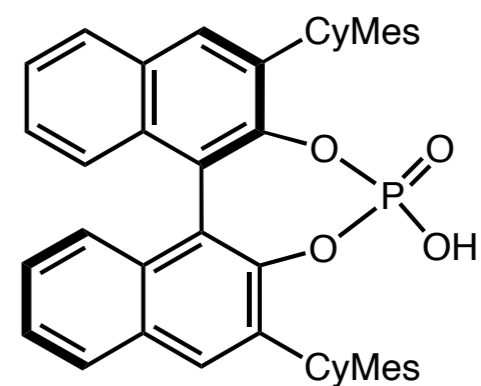
All predicted within 5% ee



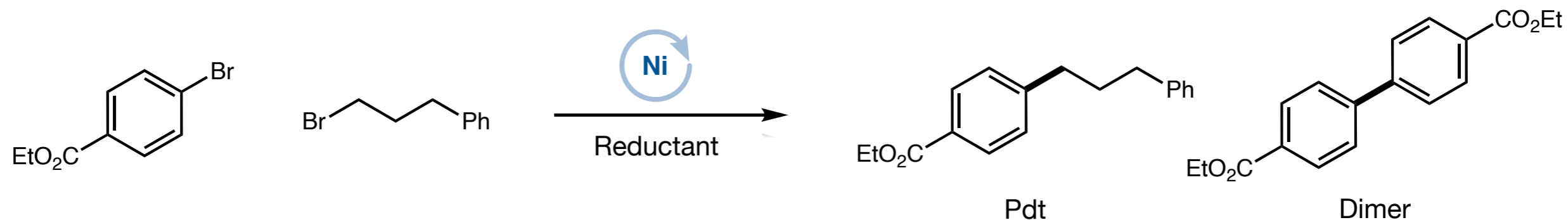
34 examples



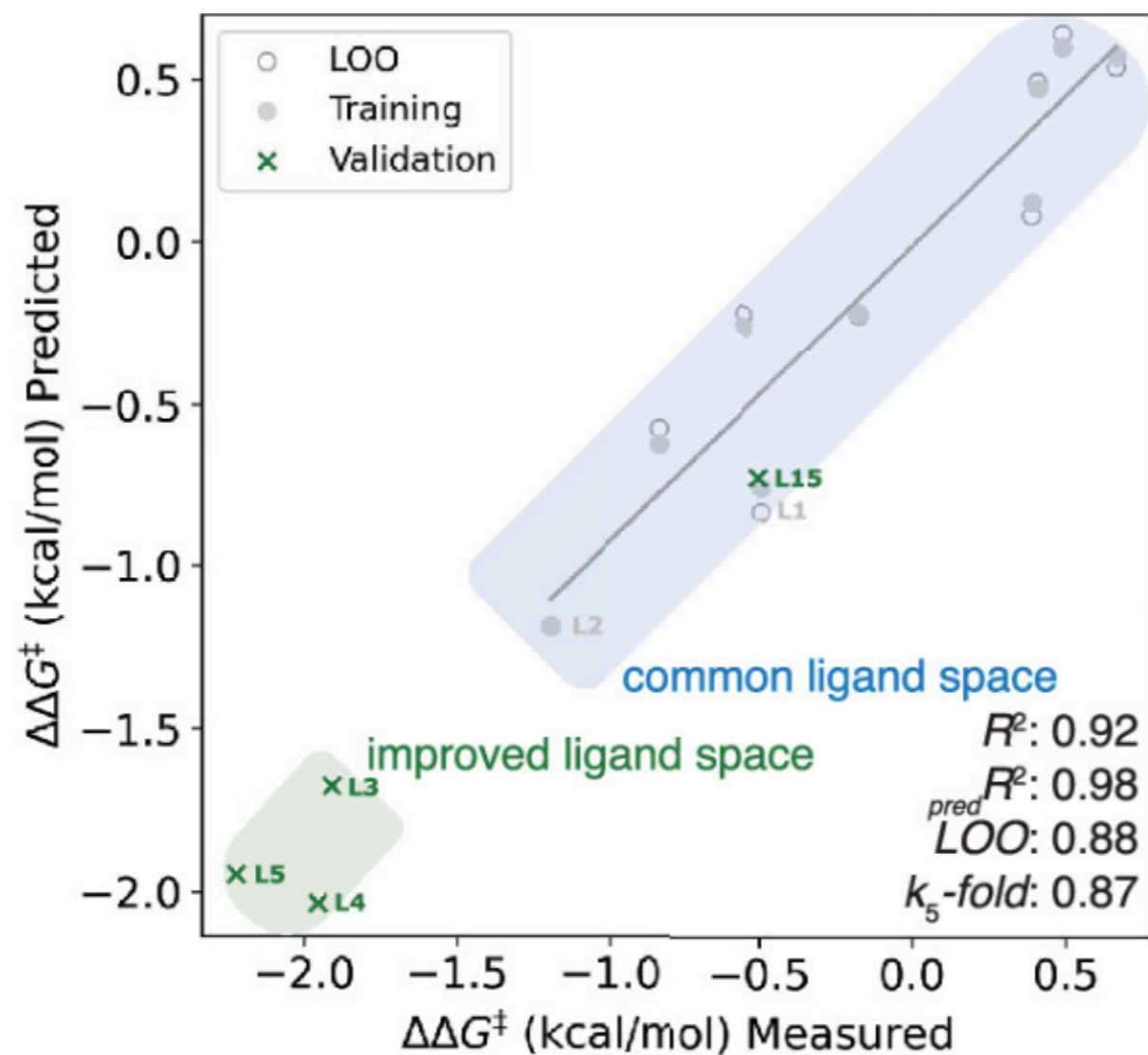
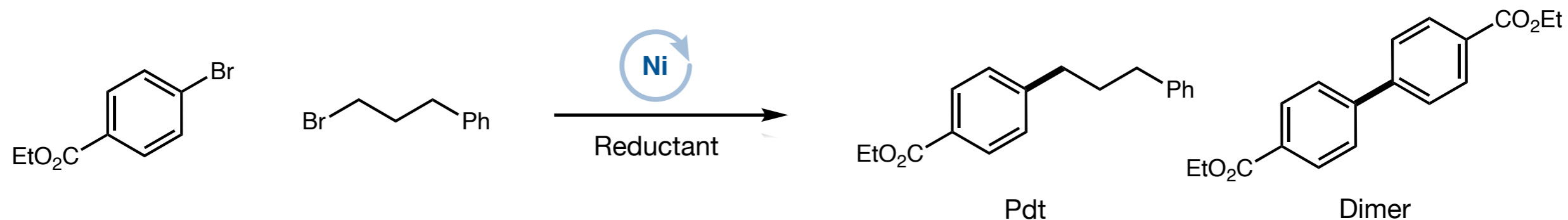
26 examples within 5% ee



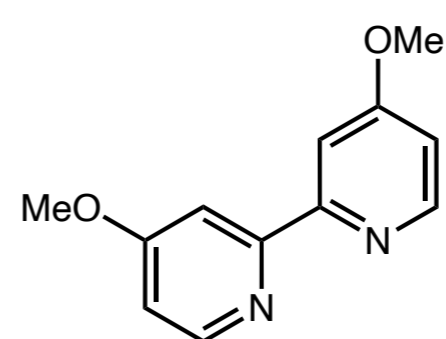
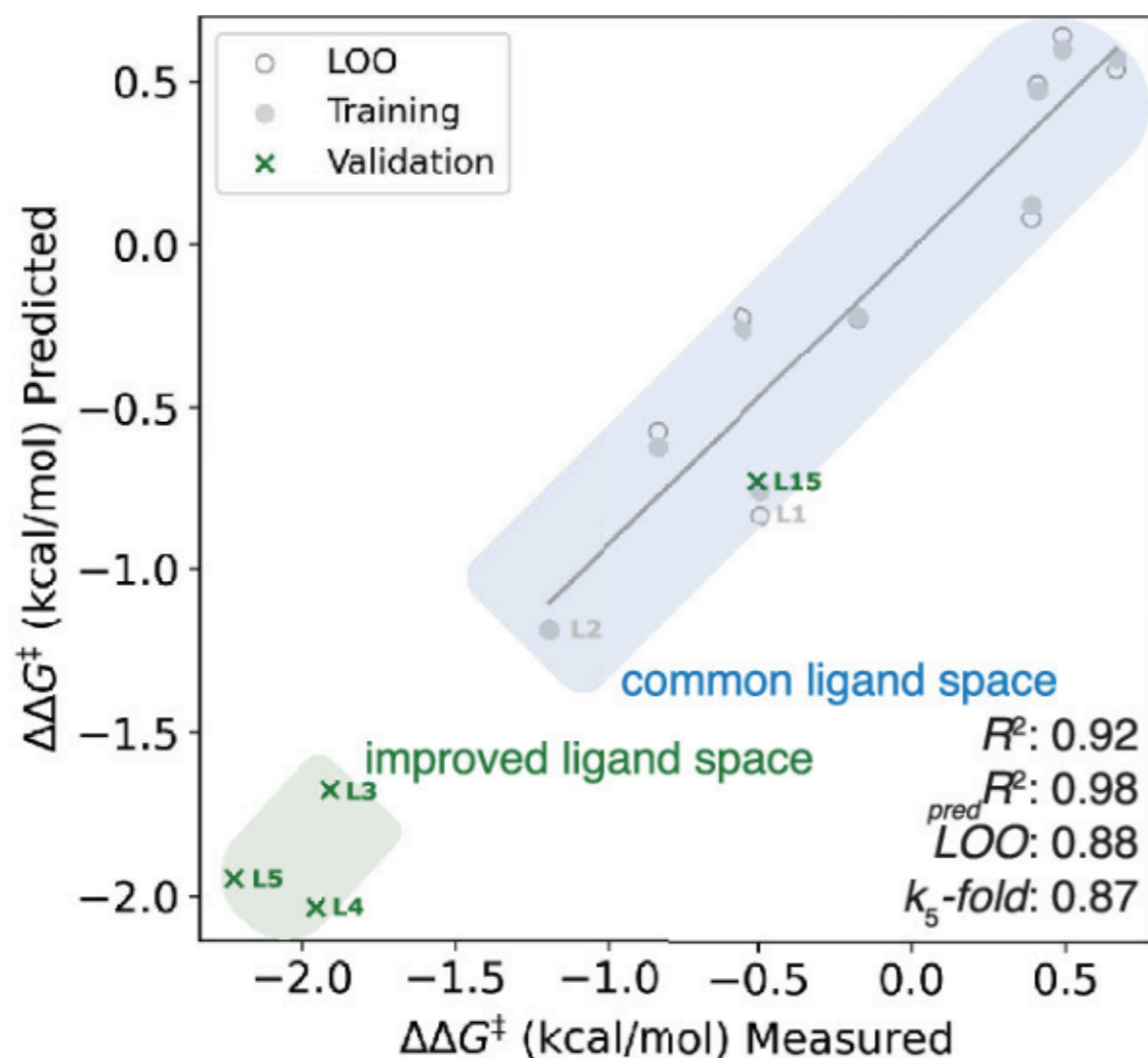
Multivariate Linear Free Energy Relationships - Finding New Catalysts



Multivariate Linear Free Energy Relationships - Finding New Catalysts

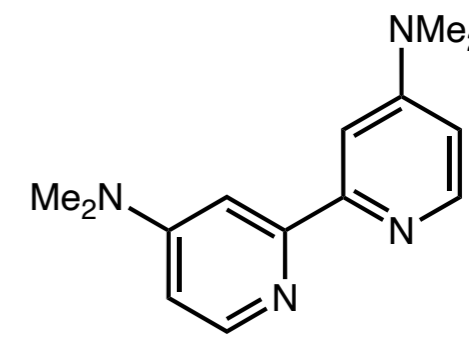


Multivariate Linear Free Energy Relationships - Finding New Catalysts



Previous best

61% yield
7:1 Pdt:dimer

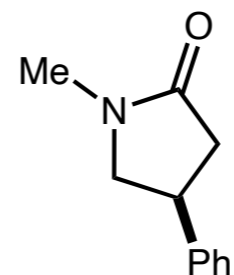
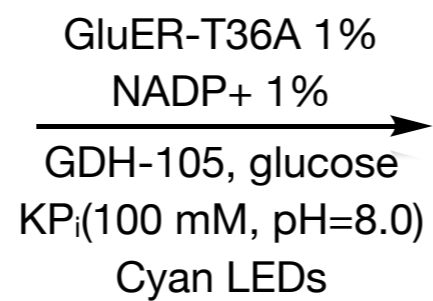
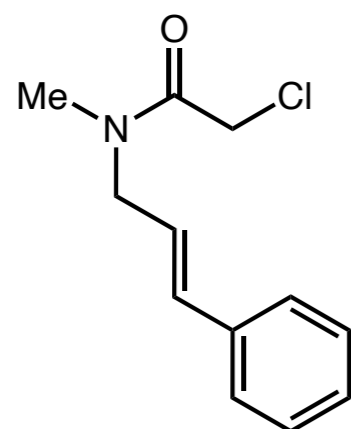


Predicted best

92% yield
187:1 Pdt:dimer

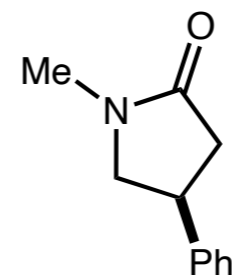
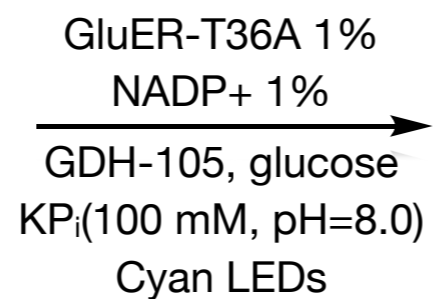
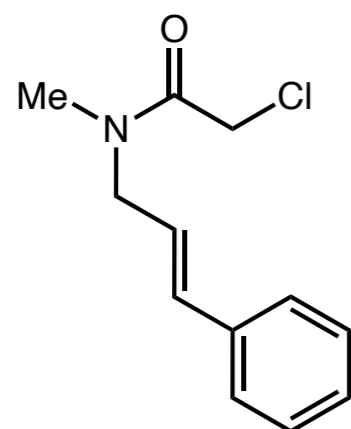
Poor understanding of the origin of the dimer, but model provides a hypothesis free solution to the problem

Multivariate Linear Free Energy Relationships - Finding New Catalysts

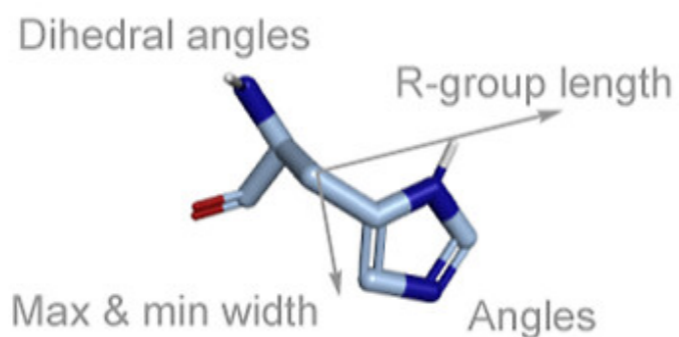


GluER-T36A

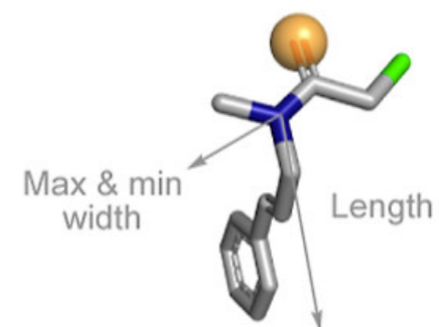
Multivariate Linear Free Energy Relationships - Finding New Catalysts



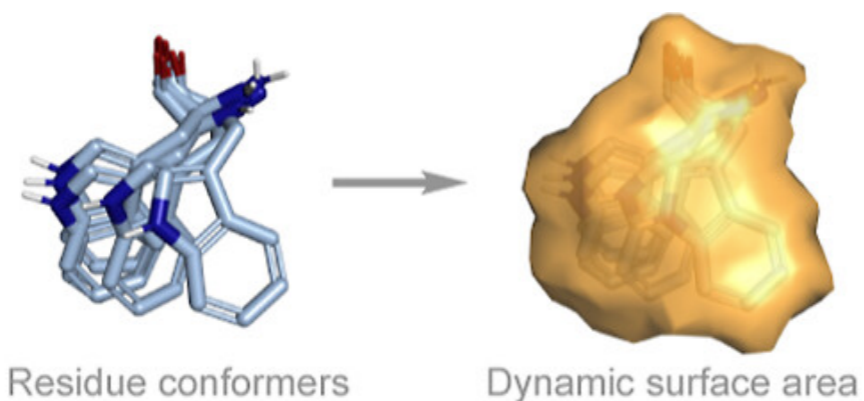
GluER-T36A



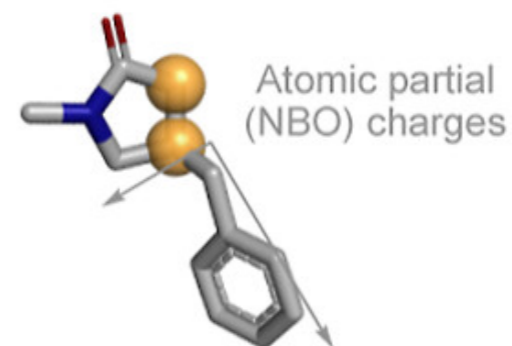
Geometric descriptors of relevant residues



Substrate sterics

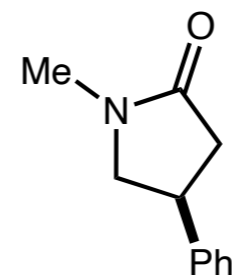
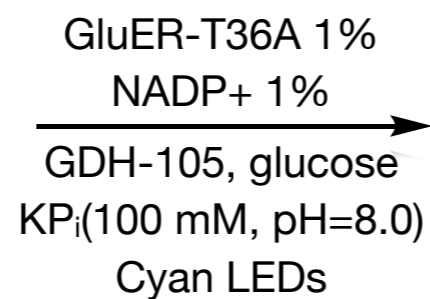
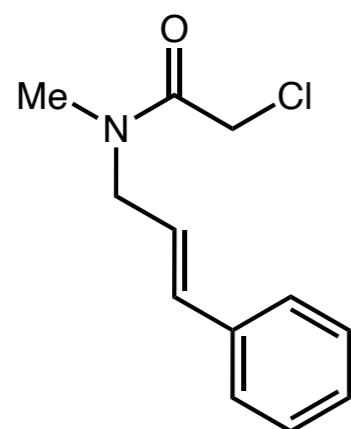


Enzyme dynamics

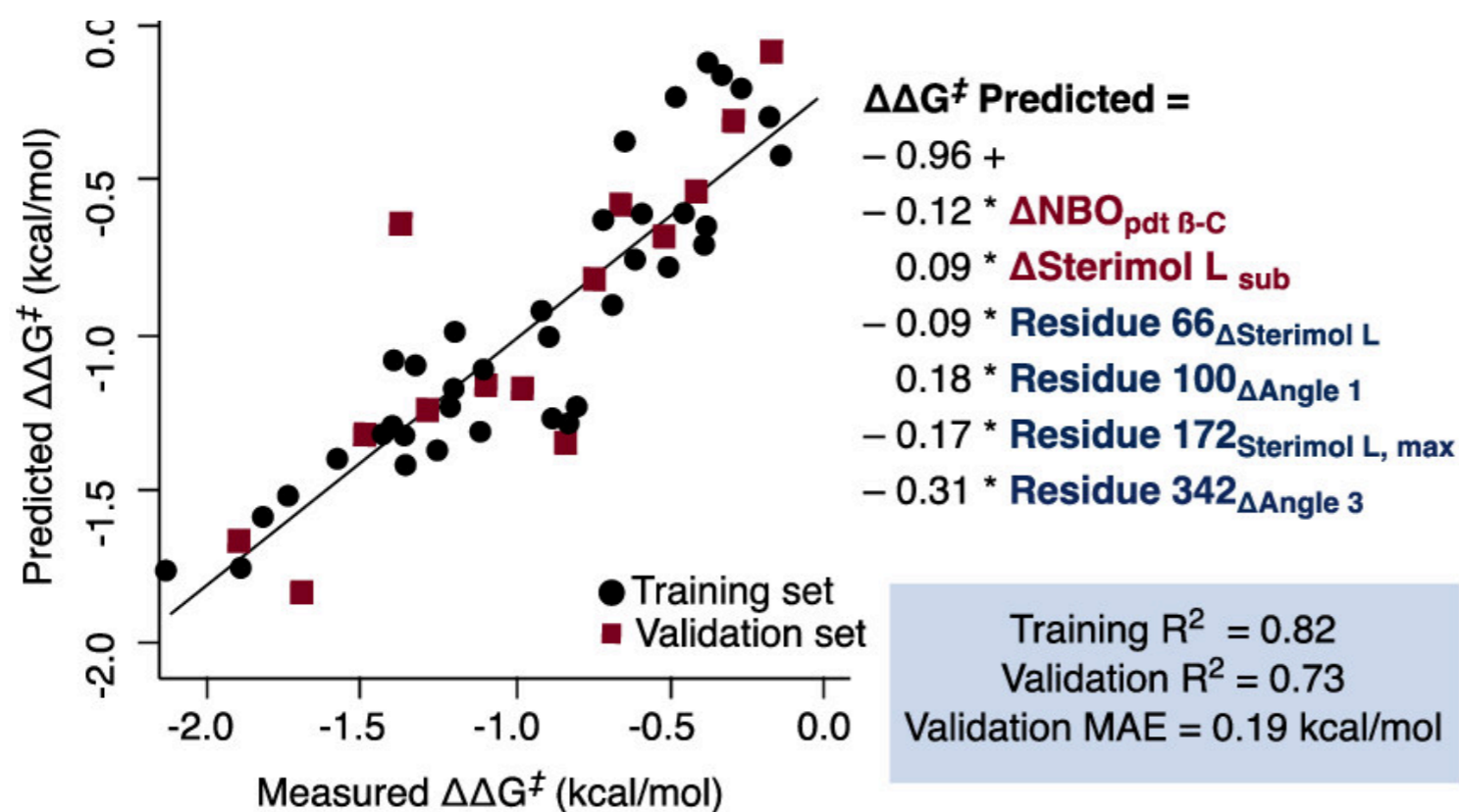


Substrate electronics

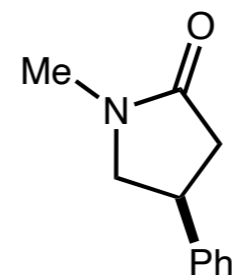
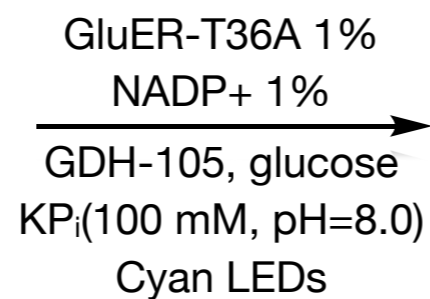
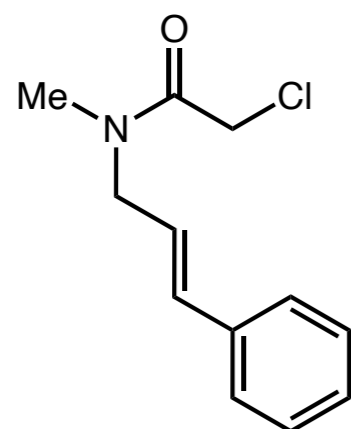
Multivariate Linear Free Energy Relationships - Finding New Catalysts



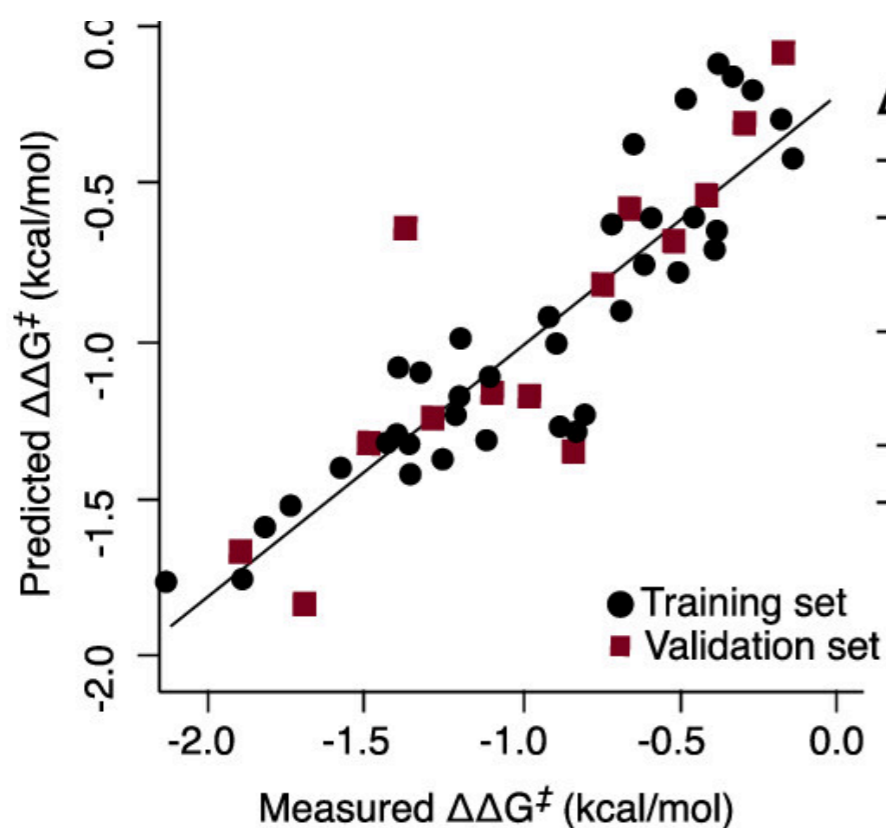
GluER-T36A



Multivariate Linear Free Energy Relationships - Finding New Catalysts



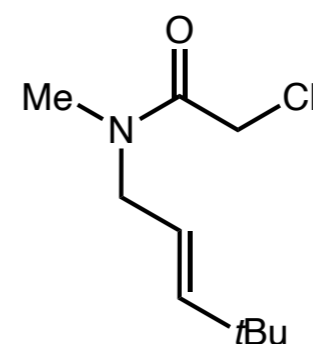
GluER-T36A



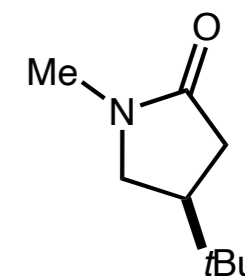
$\Delta\Delta G^\ddagger$ Predicted =

$$\begin{aligned} & -0.96 + \\ & -0.12 * \Delta\text{NBO}_{\text{pdt } \beta\text{-C}} \\ & 0.09 * \Delta\text{Sterimol } L_{\text{sub}} \\ & -0.09 * \text{Residue } 66_{\Delta\text{Sterimol } L} \\ & 0.18 * \text{Residue } 100_{\Delta\text{Angle } 1} \\ & -0.17 * \text{Residue } 172_{\text{Sterimol } L, \text{ max}} \\ & -0.31 * \text{Residue } 342_{\Delta\text{Angle } 3} \end{aligned}$$

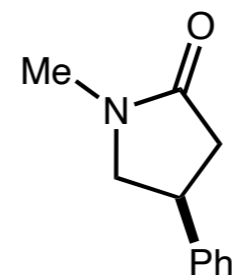
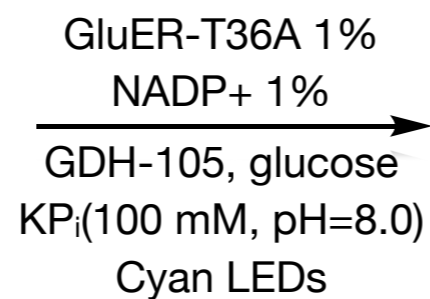
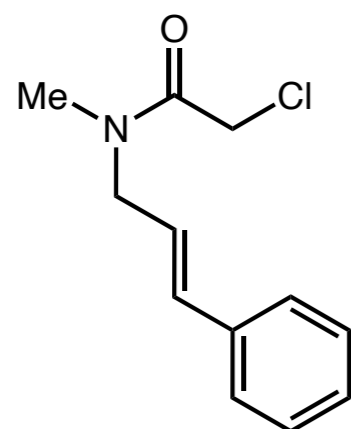
Training $R^2 = 0.82$
Validation $R^2 = 0.73$
Validation MAE = 0.19 kcal/mol



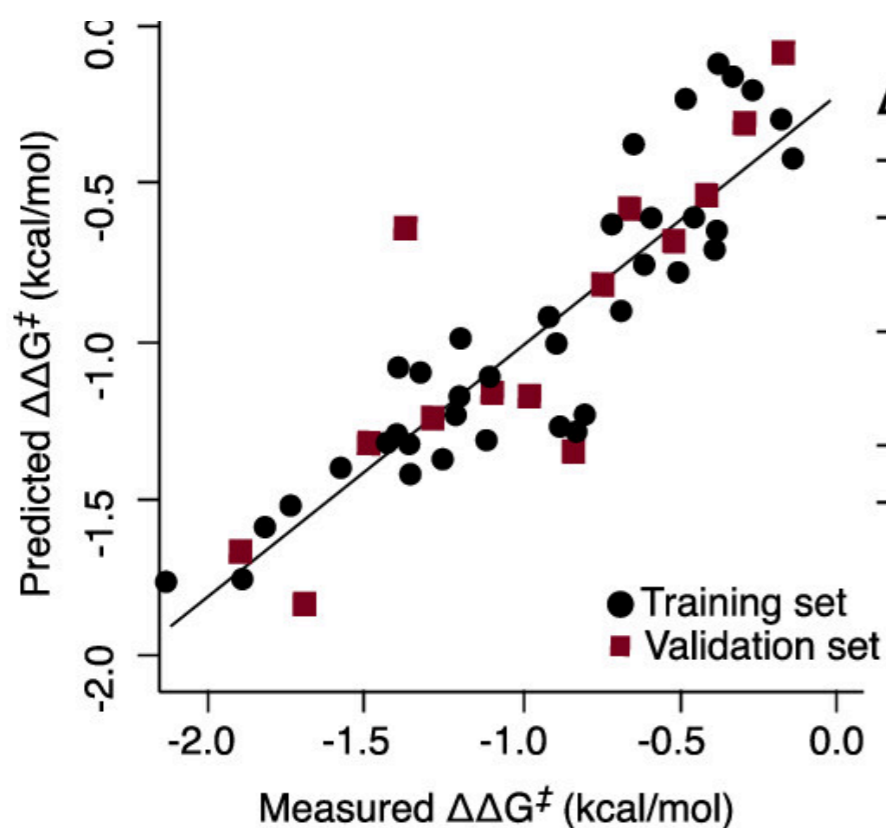
Mutant prediction



Multivariate Linear Free Energy Relationships - Finding New Catalysts

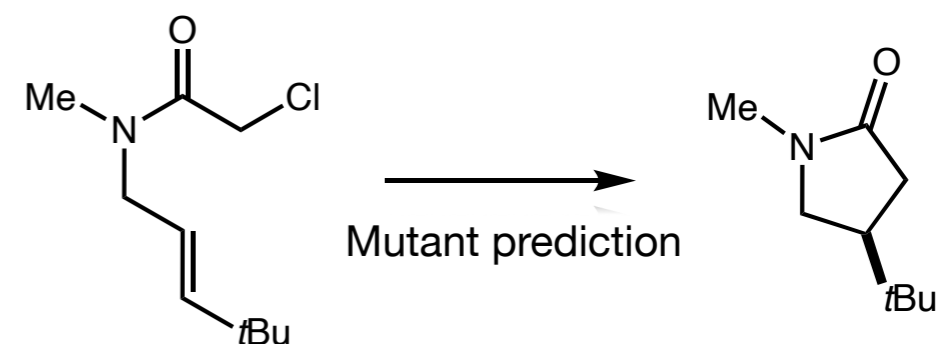


GluER-T36A



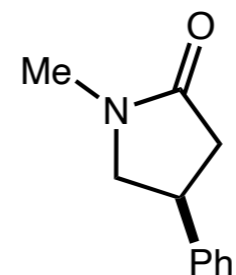
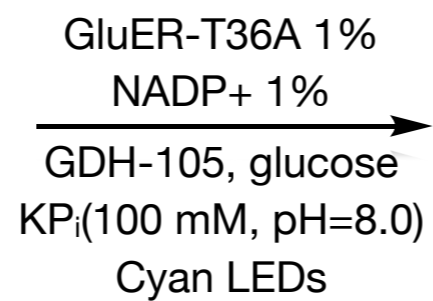
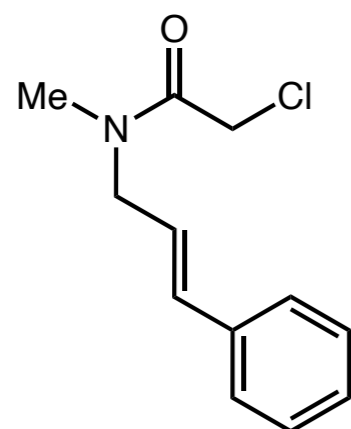
$\Delta\Delta G^\ddagger$ Predicted =
 - 0.96 +
 - 0.12 * $\Delta\text{NBO}_{\text{pdt } \beta\text{-C}}$
 0.09 * $\Delta\text{Sterimol L}_{\text{sub}}$
 - 0.09 * $\text{Residue 66}_{\Delta\text{Sterimol L}}$
 0.18 * $\text{Residue 100}_{\Delta\text{Angle 1}}$
 - 0.17 * $\text{Residue 172}_{\text{Sterimol L, max}}$
 - 0.31 * $\text{Residue 342}_{\Delta\text{Angle 3}}$

Training $R^2 = 0.82$
 Validation $R^2 = 0.73$
 Validation MAE = 0.19 kcal/mol

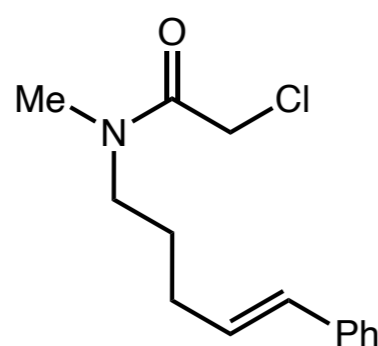


Entry	Pred. %ee ^a	Meas. %ee
5-W66A	34% ee	21% ee
5-W66L	34% ee	9% ee
5-Y177F	27% ee	36% ee
5-Q232F	44% ee	3% ee
5-Y343A	23% ee	29% ee
5-Y343F	80% ee	72% ee
5-Y343W	44% ee	38% ee

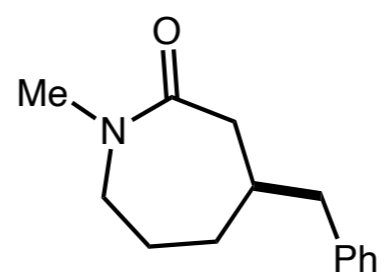
Multivariate Linear Free Energy Relationships - Finding New Catalysts



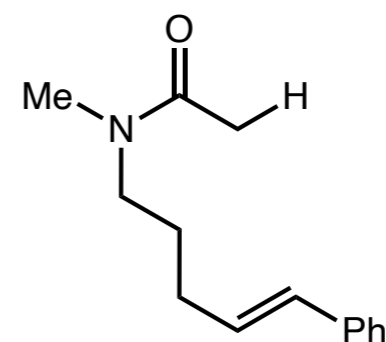
GluER-T36A



Mutant prediction

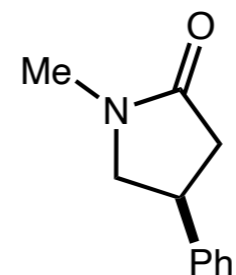
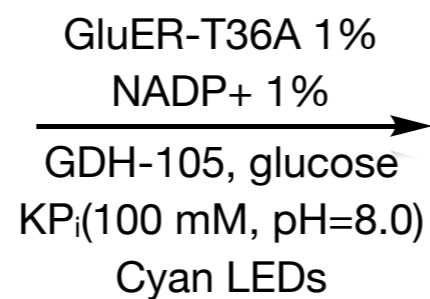
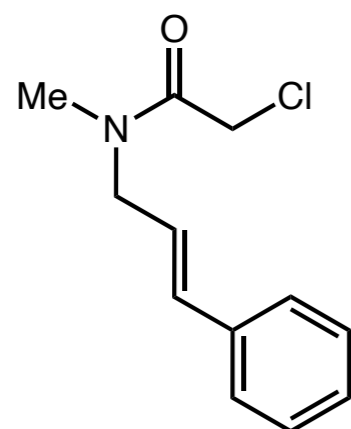


6b

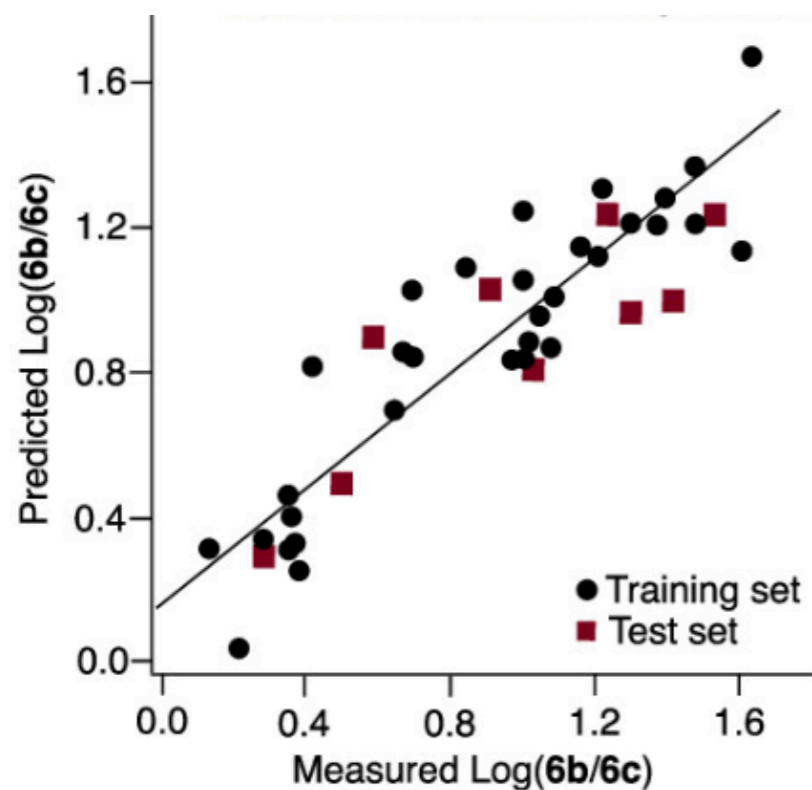


6c

Multivariate Linear Free Energy Relationships - Finding New Catalysts



GluER-T36A

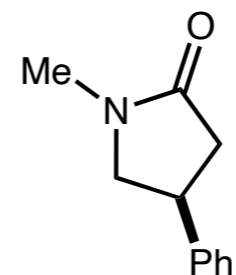
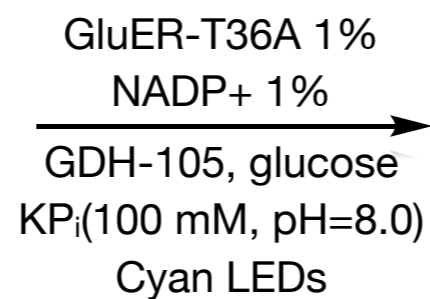
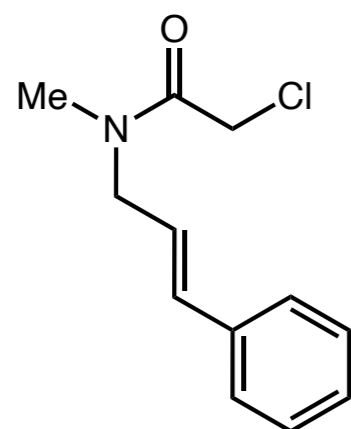


$\Delta\Delta G^\ddagger$ Predicted =

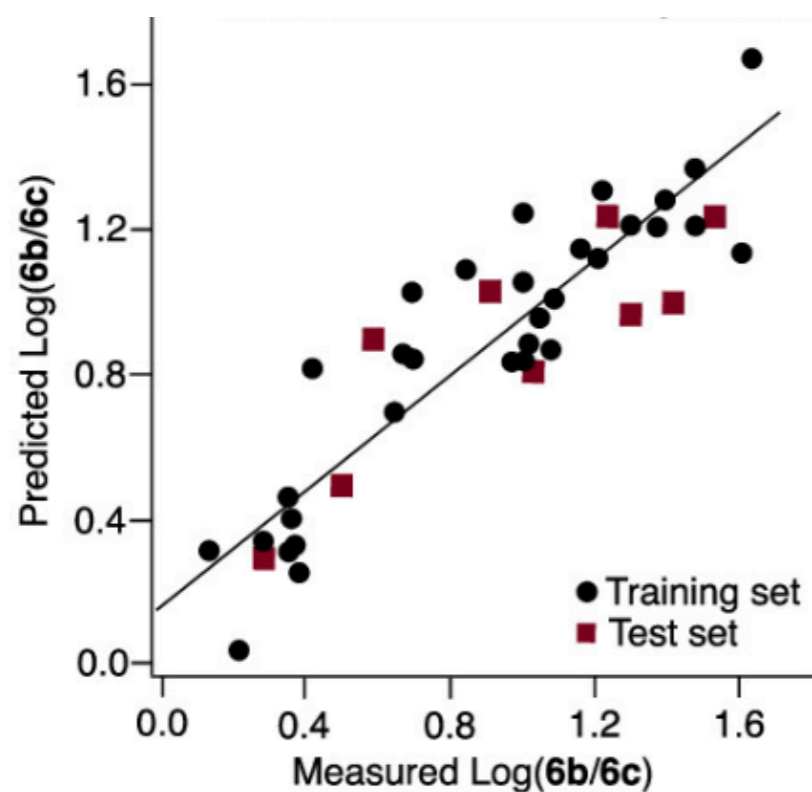
$$\begin{aligned} &0.90 + \\ &-0.33 * \text{NBO}_{\text{pdt carbonyl C, GS}} \\ &-0.23 * \text{NBO}_{\text{pdt } \beta\text{-H, GS}} \\ &-0.29 * \text{Sterimol } L_{\text{pdt, min}} \\ &0.13 * \text{Residue 269}_{\text{Sterimol B5 sub, GS}} \end{aligned}$$

Training $R^2 = 0.82$
Validation $R^2 = 0.70$
Validation MAE = 0.19 kcal/mol

Multivariate Linear Free Energy Relationships - Finding New Catalysts



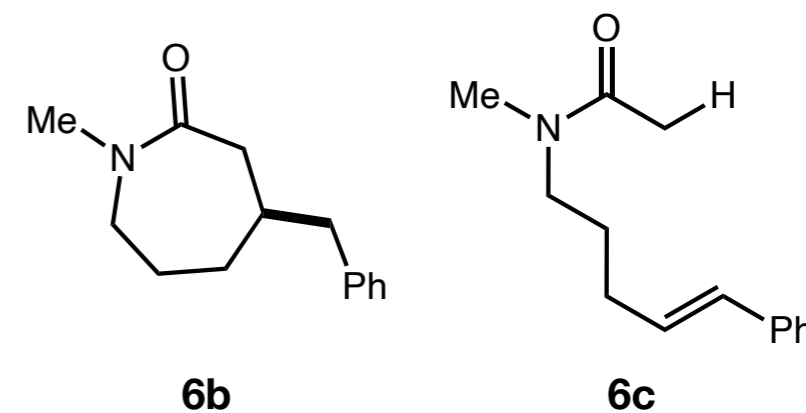
GluER-T36A



$\Delta\Delta G^\ddagger$ Predicted =

$$\begin{aligned}
 &0.90 + \\
 &-0.33 * \text{NBO}_{\text{pdt carbonyl C, GS}} \\
 &-0.23 * \text{NBO}_{\text{pdt B-H, GS}} \\
 &-0.29 * \text{Sterimol } L_{\text{pdt, min}} \\
 &0.13 * \text{Residue 269}_{\text{Sterimol B5 sub, GS}}
 \end{aligned}$$

Training $R^2 = 0.82$
 Validation $R^2 = 0.70$
 Validation MAE = 0.19 kcal/mol



Entry	Pred. (6b/6c)	Meas. (6b/6c)
6-W66A	4.0	4.6
6-W66L	3.9	9.2
6-Y177F	1.9	7.9
6-Q232F	2.8	4.3
6-Y343A	4.7	6.6
6-Y343F	2.7	2.8
6-Y343W	2.6	n.r.

Additional Reading

Data Science in Chemistry

Reviews

Williams, W. L. Zeng, L.; Gensch, T.; Sigman, M. S.; Doyle, A. G., **Anslyn, E. V.** *ACS Cent. Sci.* **2021**, *21*, 1622–1637

Raghavan, P.; Haas, B. C.; Ruos, M. E.; Schleinitz, J.; Doyle, A. G.; Reisman, S. E.; Sigman, M. S.; Coley, C. W. *ACS Cent. Sci.* **2023**, *9*, 2196–2204

Crawford, J. M.; Kingston, C.; **Toste, D. F.**; Sigman, M. S. *Acc. Chem. Res.* **2021**, *54*, 3136–3148

Additional Examples

Morak, T.; Myers, T. E.; Karas, L. J.; Hardy, M. A.; Mercado, B. Q.; Sigman, M. S.; **Miller, S. J.** *J. Am. Chem. Soc.* **2023**, *145*, 22322–22328

Nistanaki, S. K.; Williams, C. G.; Wigman, B.; Wong, J. J.; Haas, B. C.; Popov, S.; Werth, J.; Sigman, M. S.; Houk, K. N.; **Nelson, H. M.** *Science*, **2022**, *378*, 1085–1091

Sasha, M. H.; Wahlman, J. L. H.; Read, J. A.; Werth, J.; **Jacobsen, E. N.**; Sigman, M. S. *ACS Catal.* **2022**, *12*, 14836–14845

Boni, Y. T.; Cammarota, R. C.; Liao, K.; Sigman, M. S.; **Davies, H. M. L.** *J. Am. Chem. Soc.* **2022**, *144*, 15549–15591

Modern approaches to methods development

Modern Paradigms in Screening

- *Reaction generalization*
- *“Accelerate” Serendipity*
- *Miniaturization of unique reaction set ups*

Data Science

- *Catalyst optimization*
- *Predicting selectivity*
- *Discovery of new catalysts*

Machine Learning

- *What is machine learning*
- *Prediction of optimal conditions*
- *Selectivity prediction for complex systems*
- *Catalyst Discovery*

What is Machine Learning?

General use of algorithms and data to create autonomous or semi-autonomous tasks

Define a task



identify apples?

What is Machine Learning?

General use of algorithms and data to create autonomous or semi-autonomous tasks

Define a task



identify apples?

Parameterize Data



Pictures of apples

What is Machine Learning?

General use of algorithms and data to create autonomous or semi-autonomous tasks

Define a task



identify apples?

Parameterize Data



Pictures of apples

Curate Data



Is an object an apple?

What is Machine Learning?

General use of algorithms and data to create autonomous or semi-autonomous tasks

Define a task



identify apples?

Parameterize Data



Pictures of apples

Curate Data



Is an object an apple?

Choose Model

Extract patterns from data



101010
010101
101010



Train the AI

What is Machine Learning?

General use of algorithms and data to create autonomous or semi-autonomous tasks

Define a task



identify apples?

Parameterize Data



Pictures of apples

Curate Data



Is an object an apple?

Choose Model

Extract patterns from data



101010
010101
101010



Train the AI

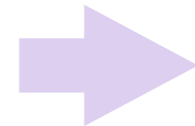
Validate Model

Input



Apples

Output



What is Machine Learning?

General use of algorithms and data to create autonomous or semi-autonomous tasks

Define a task



identify apples?

Parameterize Data



Pictures of apples

Curate Data



Is an object an apple?

Choose Model

Extract patterns from data



Train the AI

Application



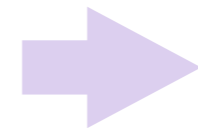
Input

Validate Model

Output



Apples



What is Machine Learning?

General use of algorithms and data to create autonomous or semi-autonomous tasks

Define a task



identify apples?

Parameterize Data



Pictures of apples

Curate Data



Is an object an apple?

Choose Model

Extract patterns from data



Train the AI

Application

Same algorithm to search through as many pictures as you want

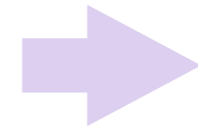
Input



Apples

Validate Model

Output



Machine Learning Applied to Chemistry

General use of algorithms and data to create autonomous or semi-autonomous tasks

Define a task



identify apples?

Parameterize Data



Pictures of apples

Curate Data



Is an object an apple?

Choose Model

Extract patterns from data



101010
010101
101010



Train the AI

Application

Same algorithm to search through as many pictures as you want

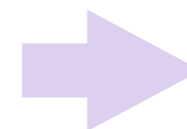
Input



Apples

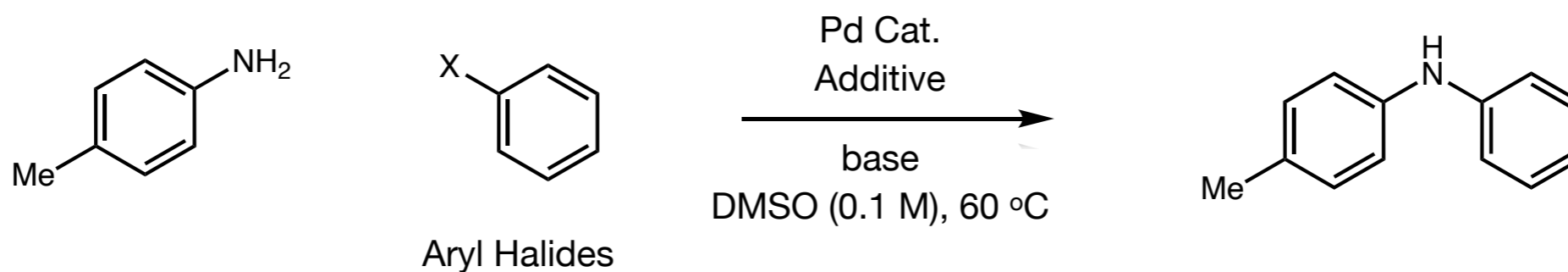
Validate Model

Output



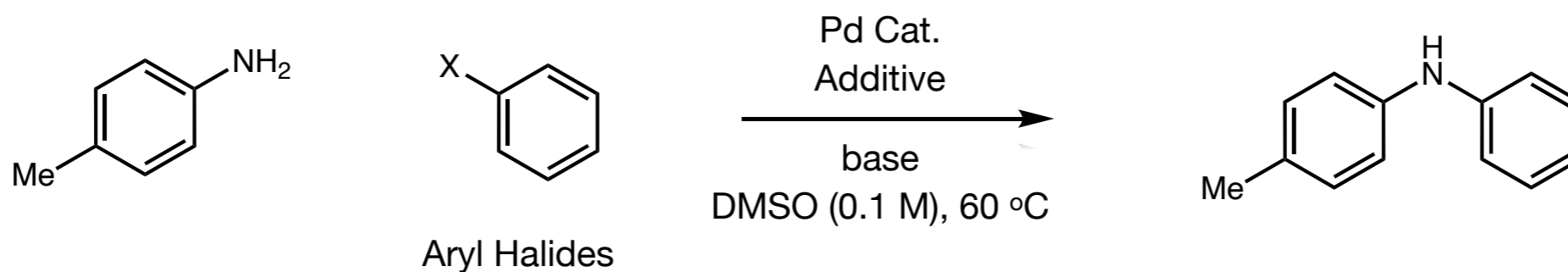
Machine Learning Applied to Chemistry

Task: Predict the yield of the Buchwald-Hartwig reaction under varying reaction conditions

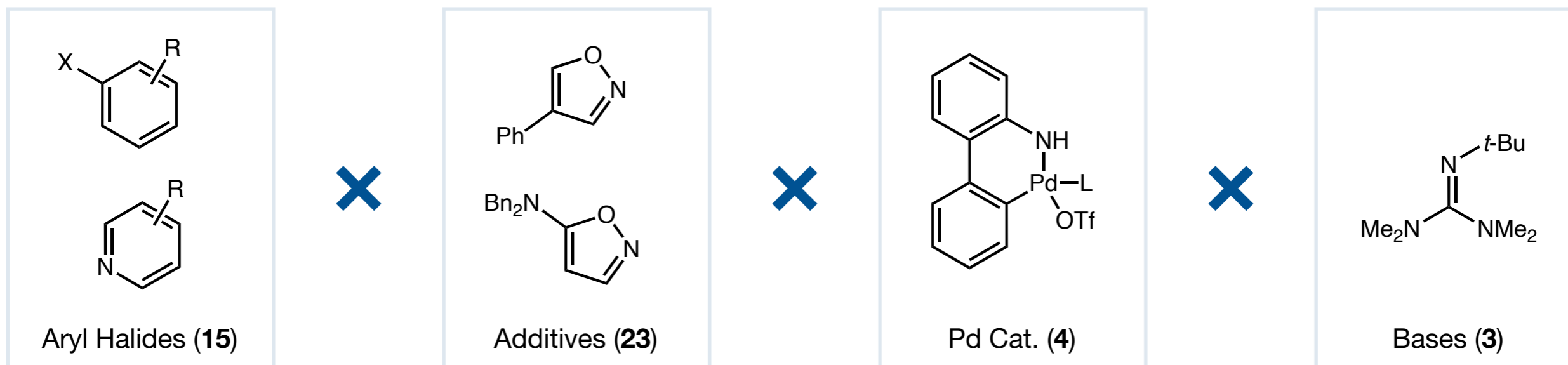


Machine Learning Applied to Chemistry

Task: Predict the yield of the Buchwald-Hartwig reaction under varying reaction conditions

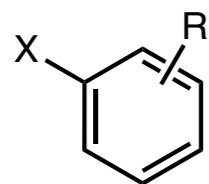


Variables

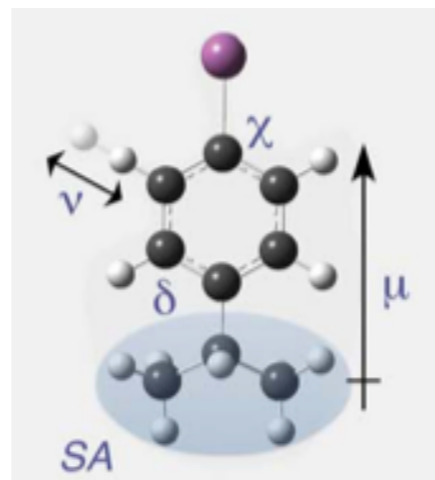
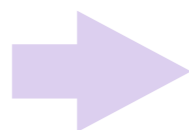


How can we feed data to the algorithm about each variable?

Machine Learning Applied to Chemistry



Spartan
(software)



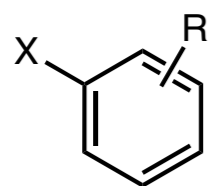
Parameterized Data

- *Vibrational properties*
- *Surface area*
- *Dipole moment*
- *Atomic electrostatic charge*

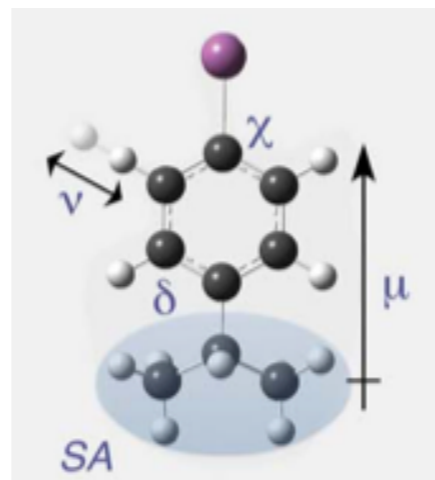
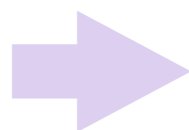
Total of 120 descriptors obtained for each individual reaction

What is the data input?

Machine Learning Applied to Chemistry



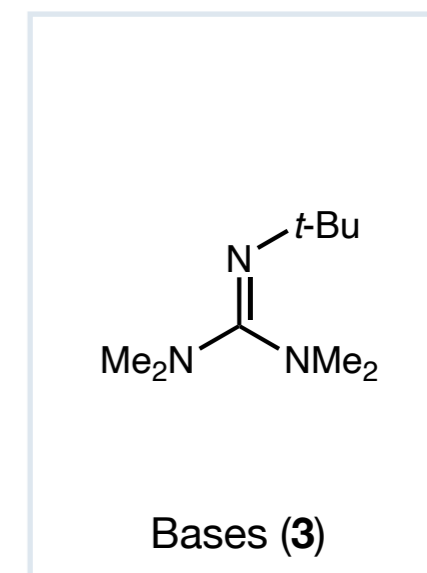
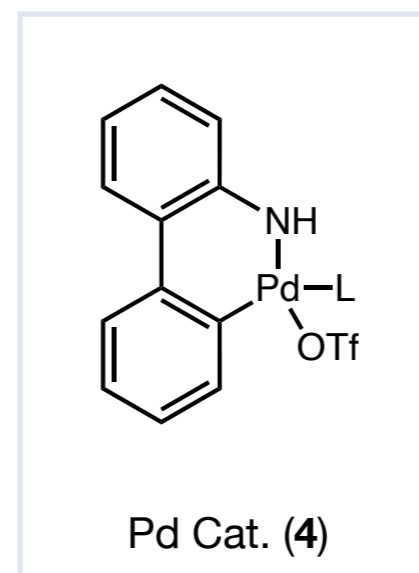
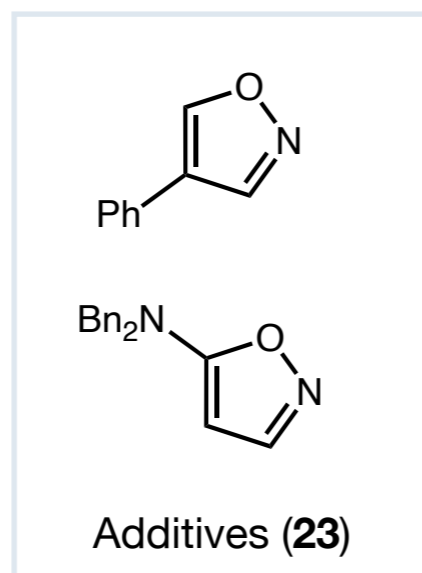
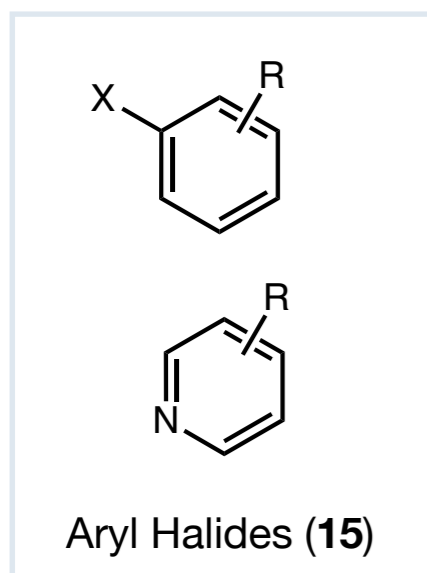
Spartan
(software)



Parameterized Data

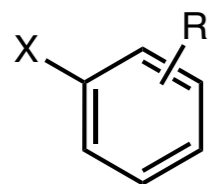
- *Vibrational properties*
- *Surface area*
- *Dipole moment*
- *Atomic electrostatic charge*

Total of 120 descriptors obtained for each individual reaction

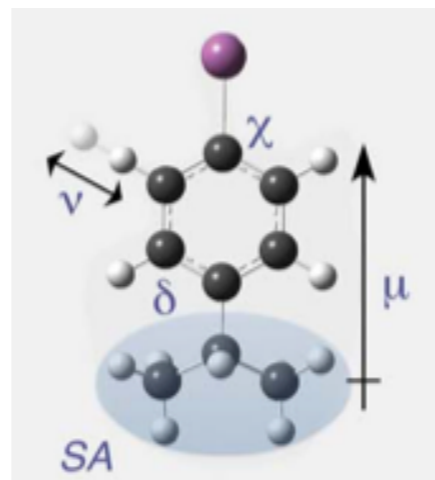
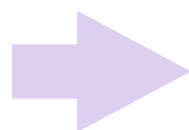


Run 4140 reactions through HTE and input yield along with parameterized reaction conditions

Machine Learning Applied to Chemistry



Spartan
(software)



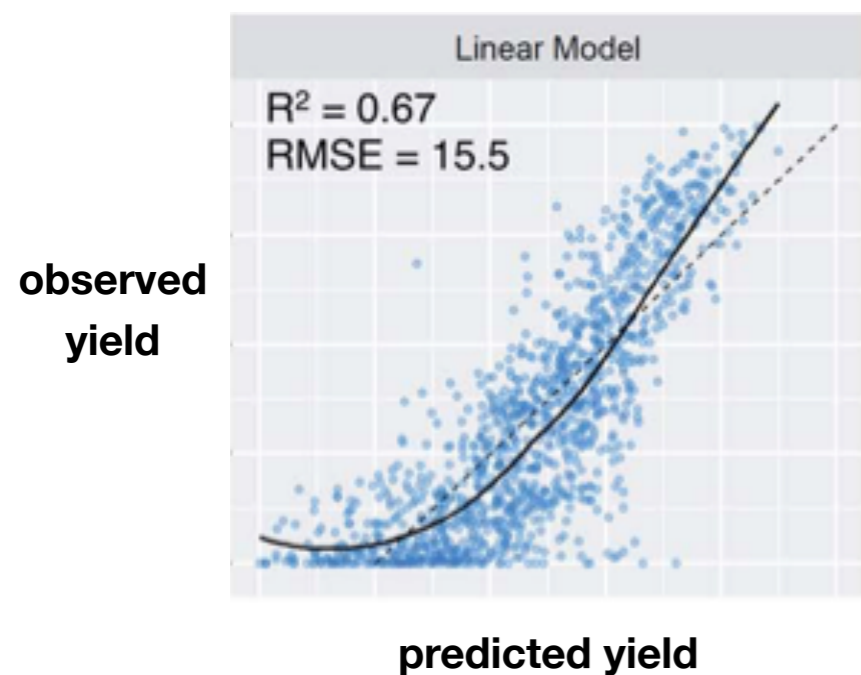
Parameterized Data

- *Vibrational properties*
- *Surface area*
- *Dipole moment*
- *Atomic electrostatic charge*

Total of 120 descriptors obtained for each individual reaction

**Unlike MLFER all 120 descriptors can be used
in the same model in a null-hypothesis manner**

Machine Learning Applied to Chemistry



RMSE: root mean square error



Accuracy

R²: Coefficient of determination

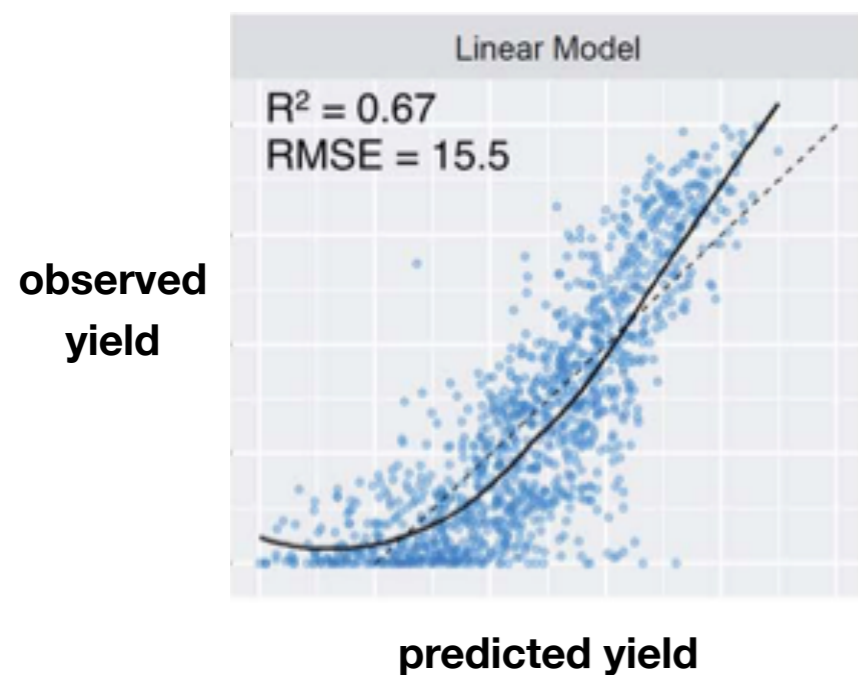


Consistency

straight, well correlated fits are more accurate models

Linear Regression Model is inaccurate and inconsistent

Machine Learning Applied to Chemistry



RMSE: root mean square error



Accuracy

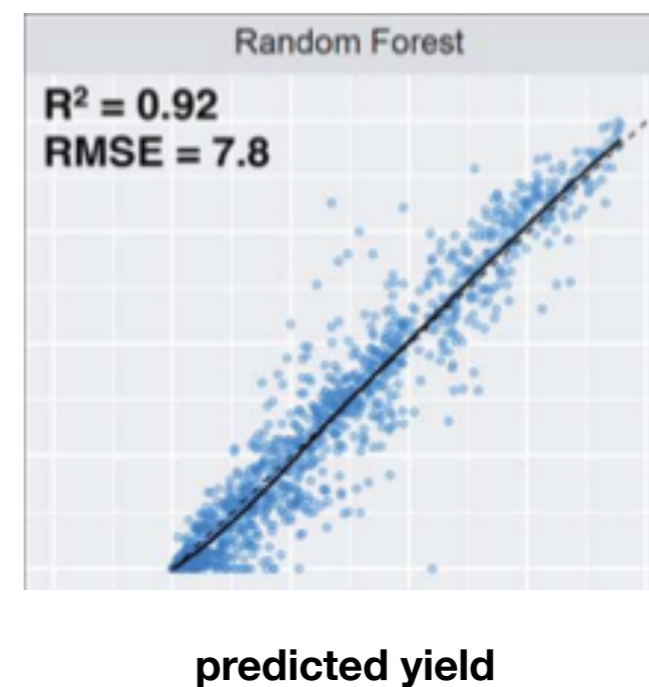
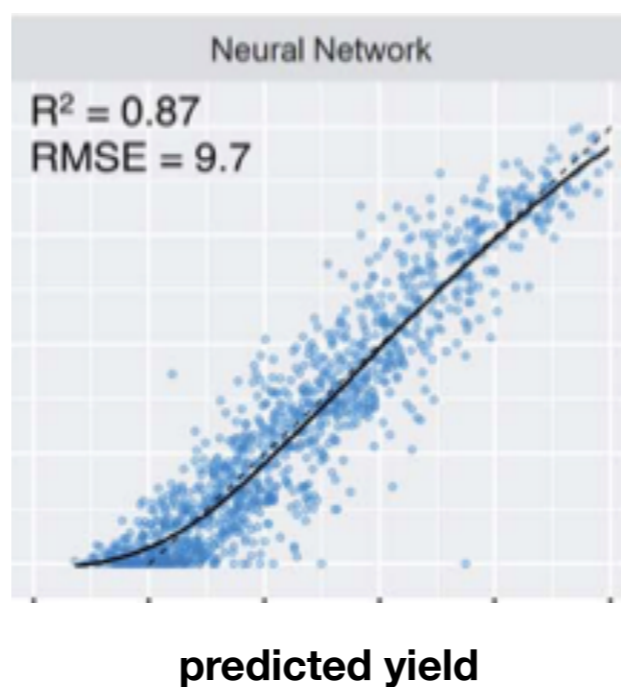
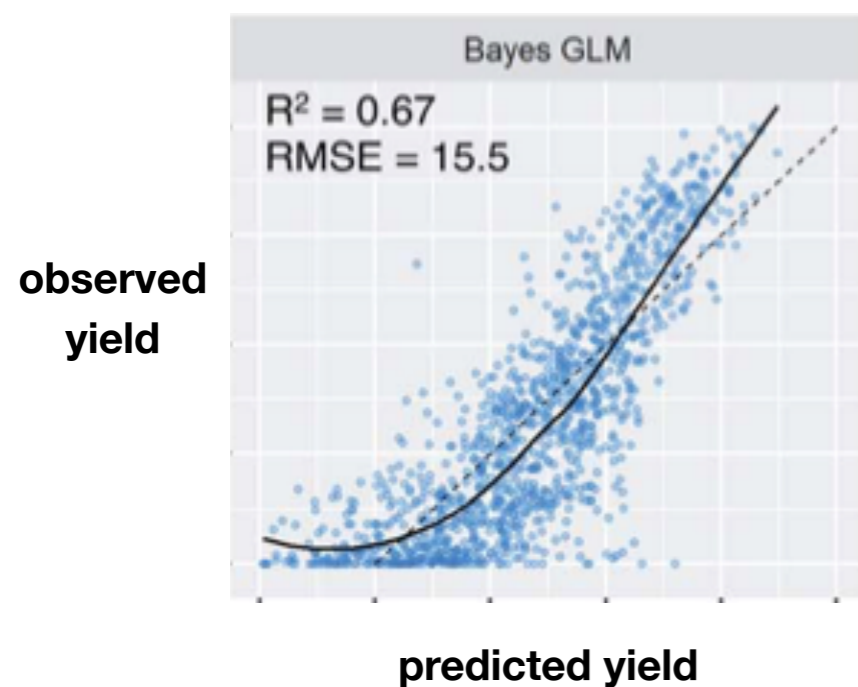
R²: Coefficient of determination



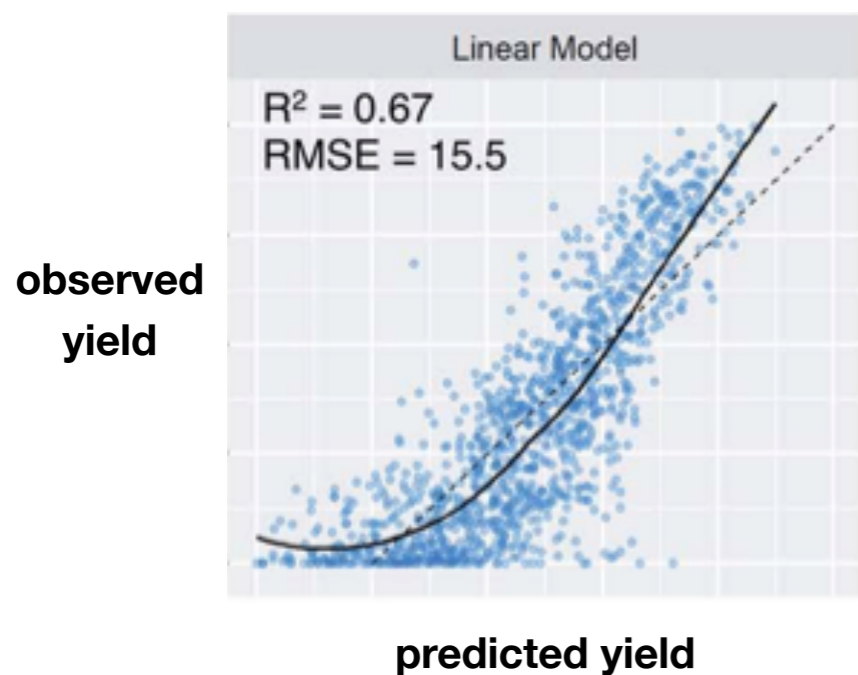
Consistency

straight, well correlated fits are more accurate models

supervised machine learning models



Machine Learning Applied to Chemistry



RMSE: root mean square error



Accuracy

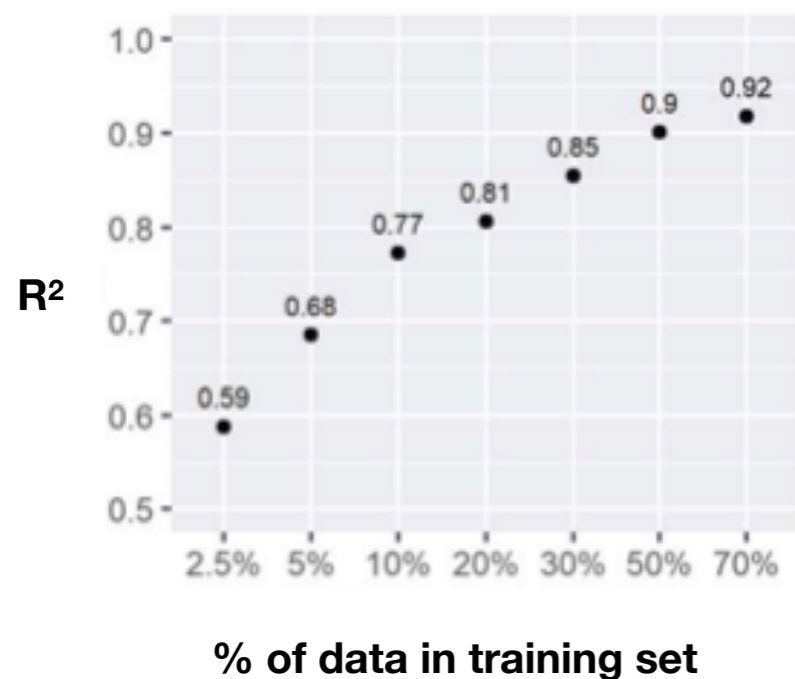
R²: Coefficient of determination



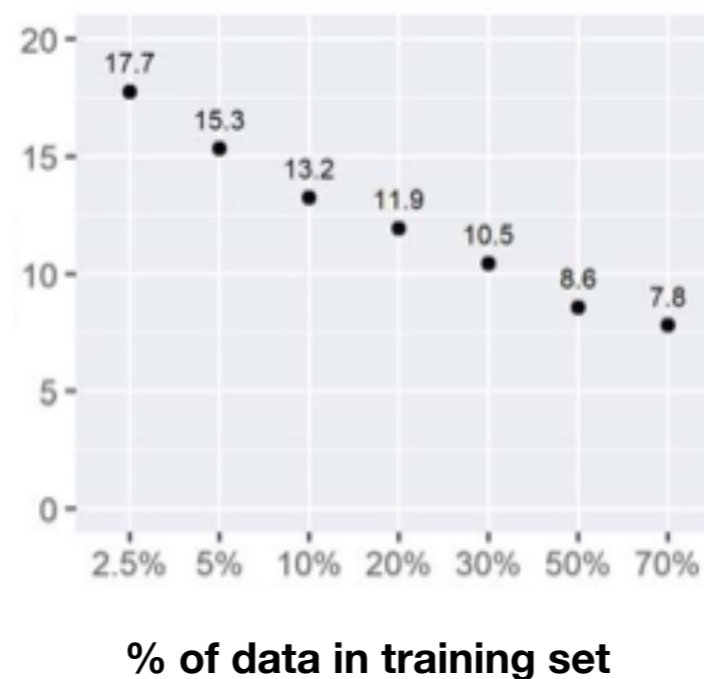
Consistency

straight, well correlated fits are more accurate models

supervised machine learning models



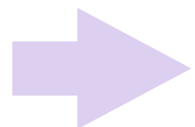
RMSE



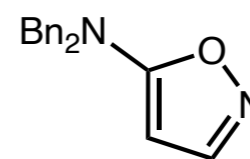
Bigger training set makes predictions more consistent and more accurate

Machine Learning Applied to Chemistry

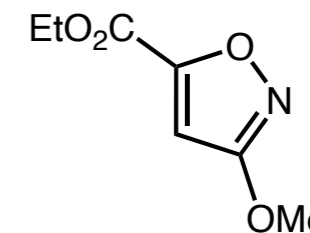
70% of data used
to train model



30% of data used
to validate model

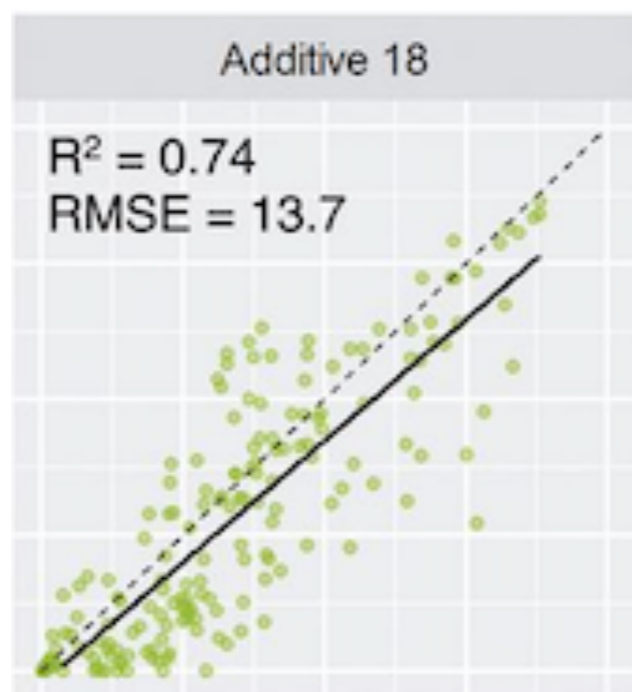


Additive 18



Additive 23

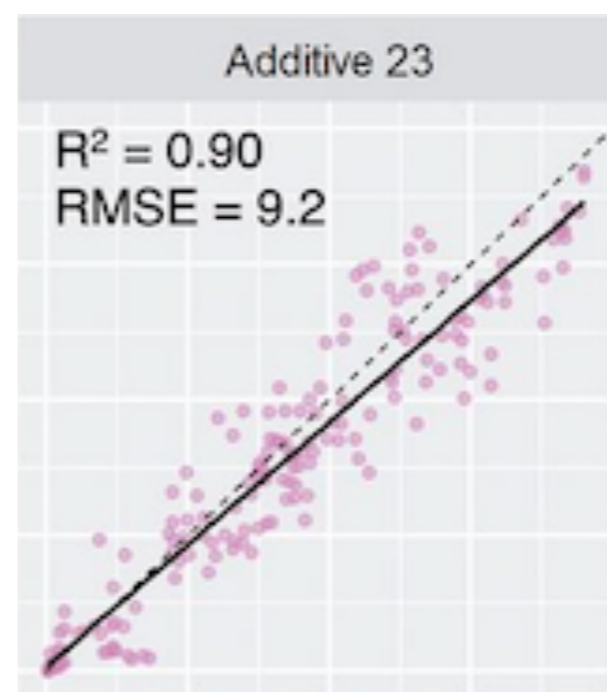
Poorly predictive



observed
yield

predicted yield

predictive

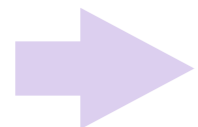


predicted yield

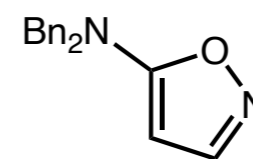
What causes the difference in accuracy between additive 18 and 23?

Machine Learning Applied to Chemistry

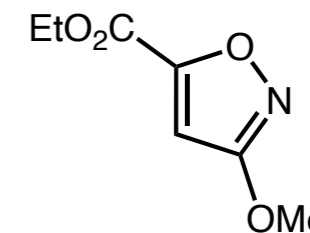
70% of data used
to train model



30% of data used
to validate model

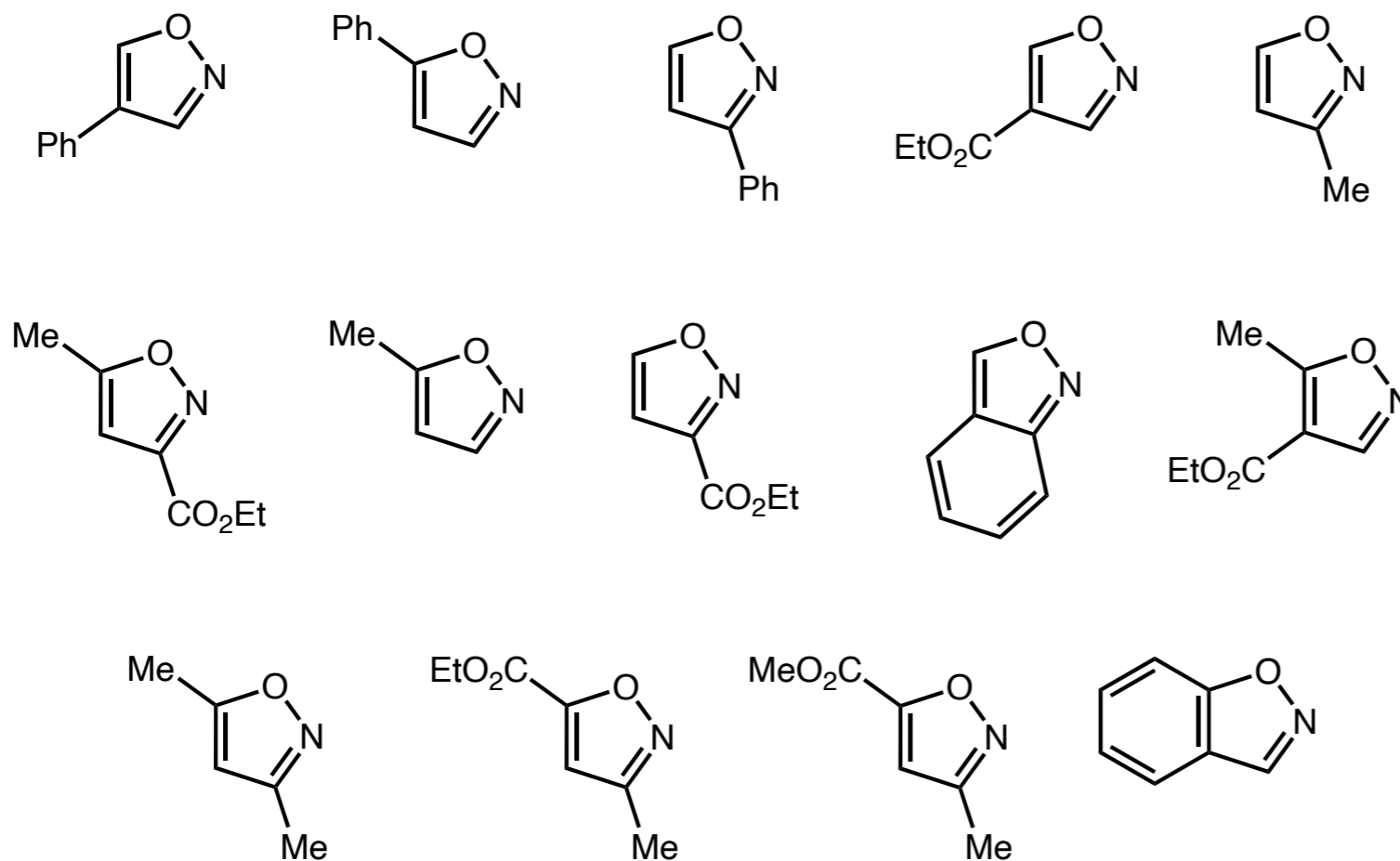


Additive 18



Additive 23

training set additives

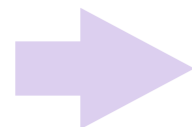


All “electron poor” to
neutral isoxazoles

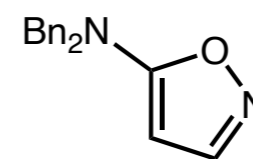
Model poorly predicts yields for
additive 18 because it falls out
of the scope of the training set

Machine Learning Applied to Chemistry

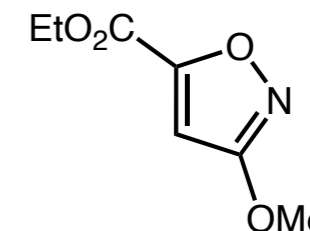
70% of data used
to train model



70% of data used
to train model



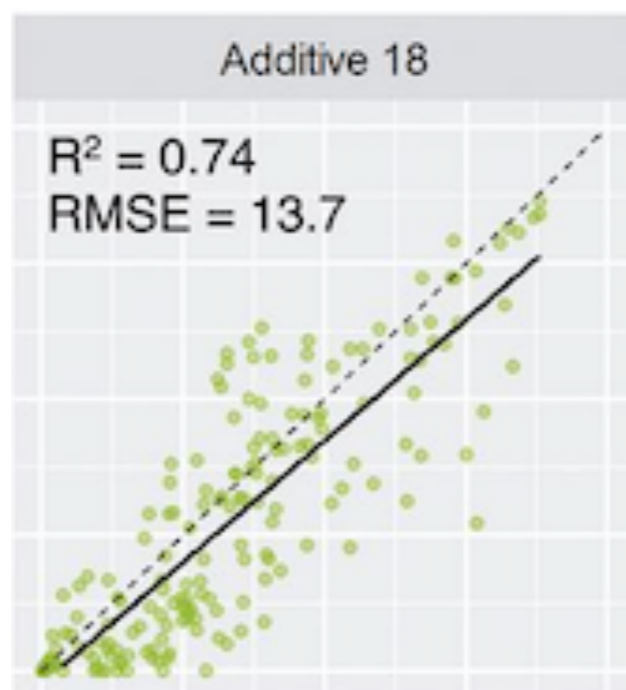
Additive 18



Additive 23

Poorly predictive

observed
yield

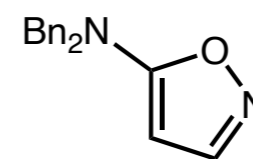


predicted yield

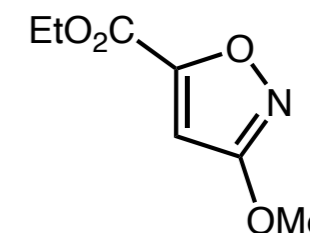
Inaccurate for parameters
outside of training data set

Machine Learning Applied to Chemistry

Why does the identity of the additive matter?



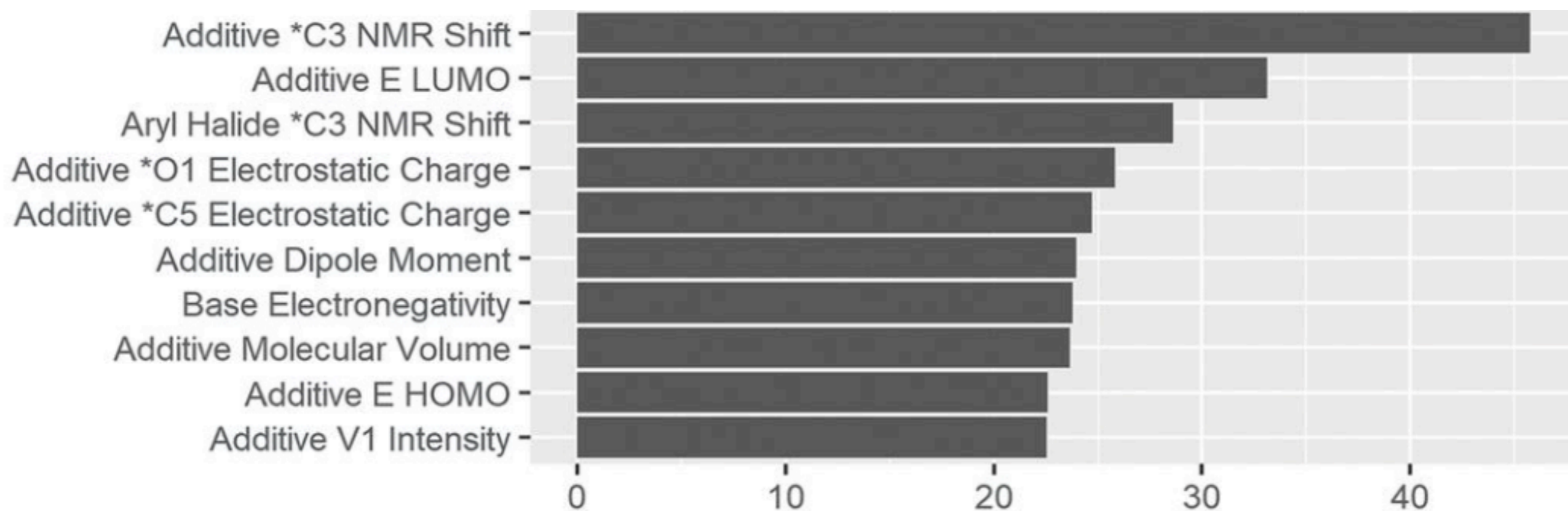
Additive 18



Additive 23

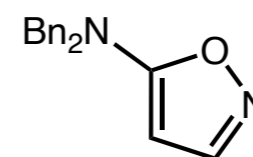
Descriptor

Increase in RMSE when excluded from model

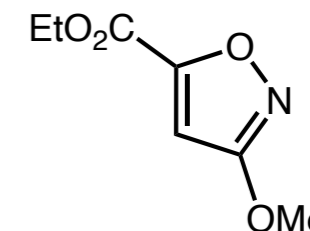


Machine Learning Applied to Chemistry

Why does the identity of the additive matter?



Additive 18



Additive 23

Descriptor

Additive *C3 NMR Shift

Additive E LUMO

Aryl Halide *C3 NMR Shift

Additive *O1 Electrostatic Charge

Additive *C5 Electrostatic Charge

Additive Dipole Moment

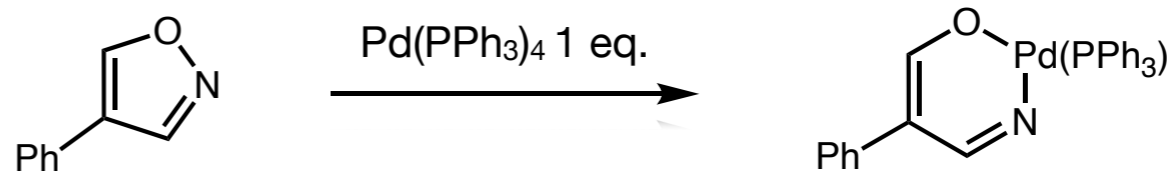
Base Electronegativity

Additive Molecular Volume

Additive E HOMO

Additive V1 Intensity

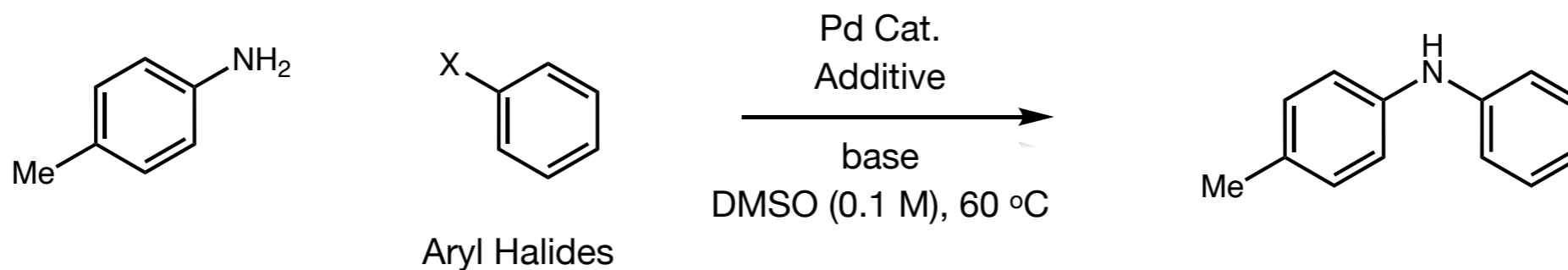
Top descriptors suggest additive electrophilicity influences yield



known for nickel but not for palladium

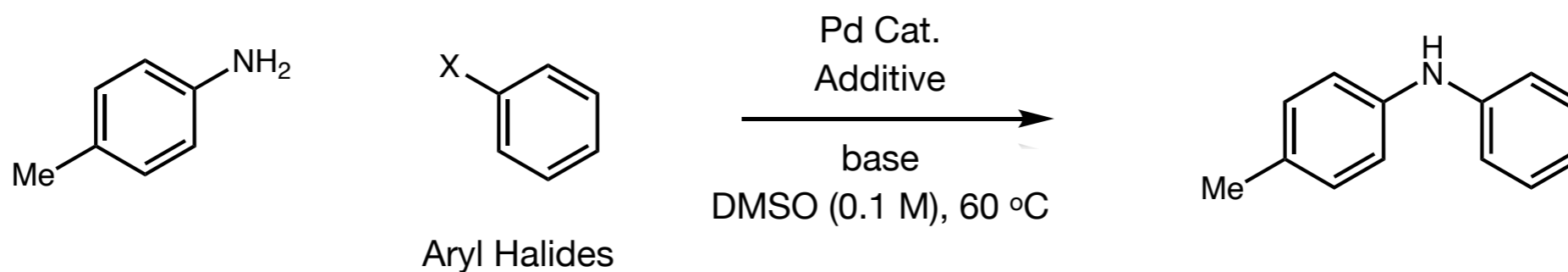
Machine Learning Applied to Chemistry

Task: Predict the yield of the Buchwald-Hartwig reaction under varying reaction conditions

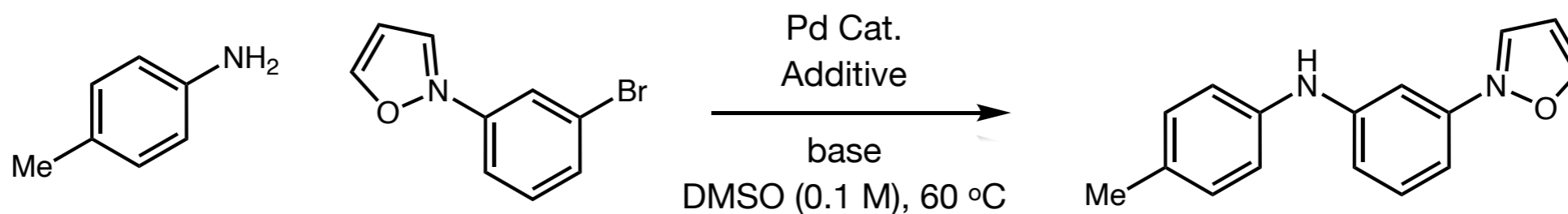


Machine Learning Applied to Chemistry

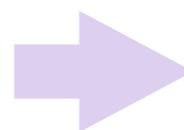
Task: Predict the yield of the Buchwald-Hartwig reaction under varying reaction conditions



Result: Model that can roughly predict the yield of a Buchwald Hartwig reaction in which isoxazoles are present



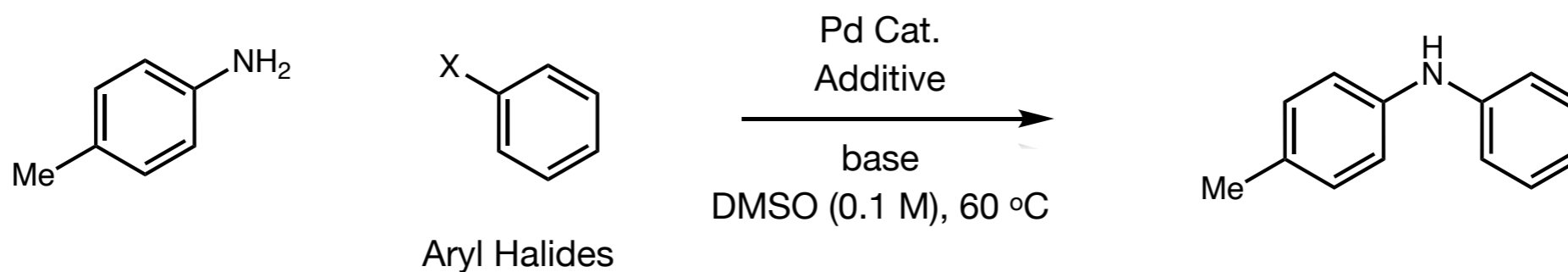
relatively specific
task



Needed 4,140
experimental yields

Machine Learning Applied to Chemistry

Task: Predict the yield of the Buchwald-Hartwig reaction under varying reaction conditions

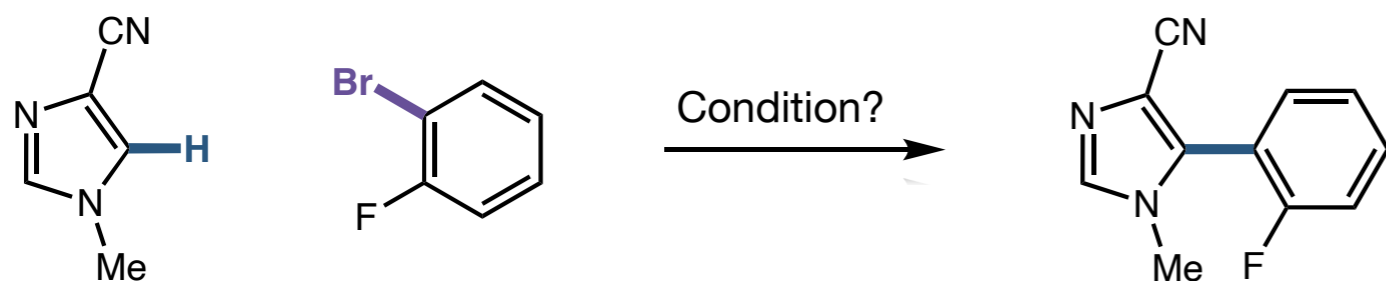


Result: Model that can roughly predict the yield of a Buchwald Hartwig reaction in which isoxazoles are present

strategy not fit for
reaction optimization

Machine Learning Applied to Chemistry - Reaction Optimization

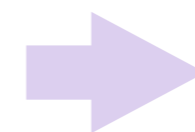
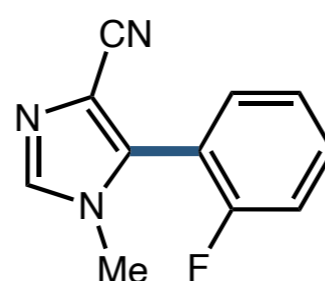
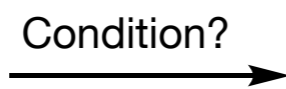
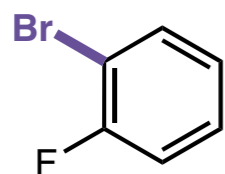
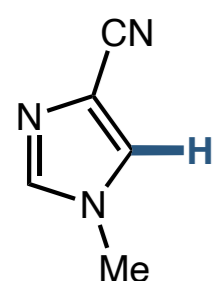
Task: Optimize a reaction where 1000's of experimental yields are not available



No previous data

Machine Learning Applied to Chemistry - Reaction Optimization

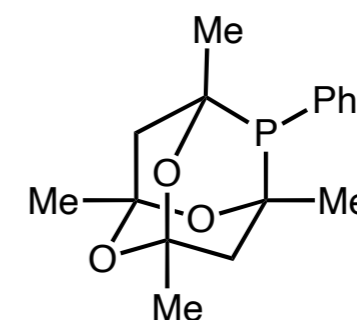
Task: Optimize a reaction where 1000's of experimental yields are not available



Optimizing
Algorithm

No previous data

Optimal Conditions

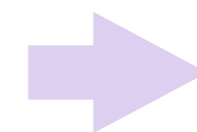
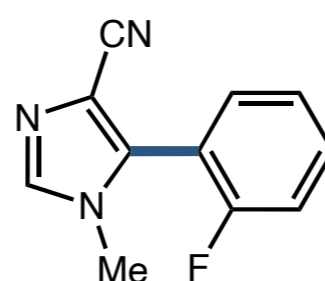
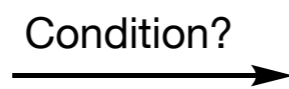
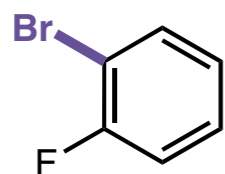
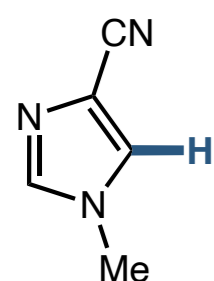


PdCl₂(allyl)₂
CsOPiv, 105 °C
DMA (0.153 M),

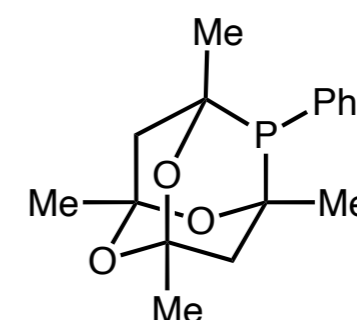
100% yield
Less than 50 experiments

Machine Learning Applied to Chemistry - Reaction Optimization

Task: Optimize a reaction where 1000's of experimental yields are not available



Optimizing
Algorithm



Optimal Conditions
PdCl₂(allyl)₂
CsOPiv, 105 °C
DMA (0.153 M),

100% yield
Less than 50 experiments

No previous data

Choose Model

Extract patterns from data



101010
010101
101010

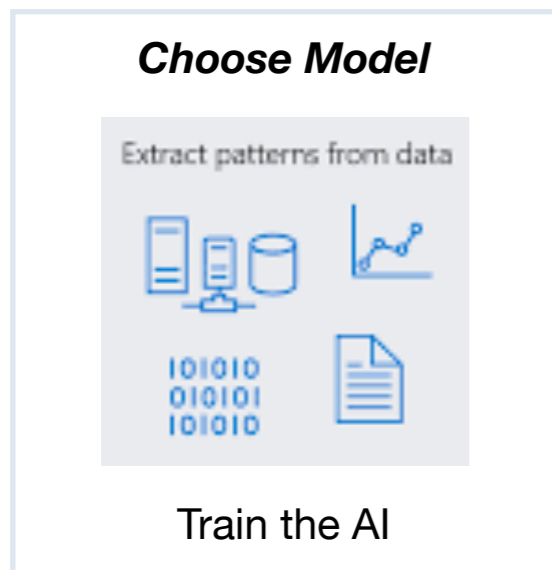


Train the AI

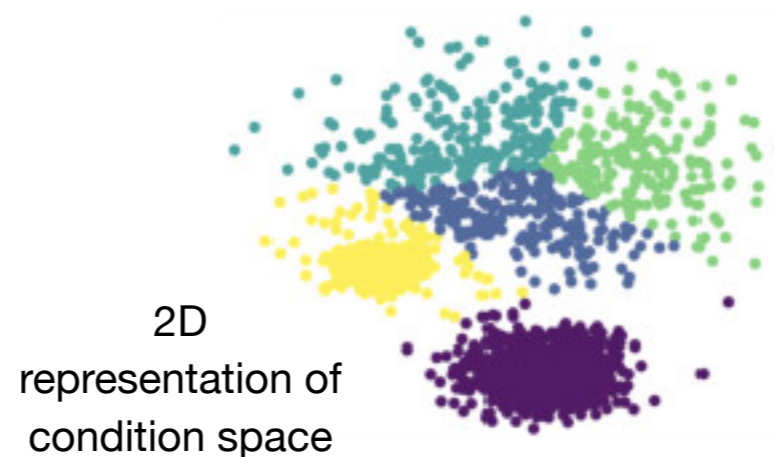
How would you train a model that
trains itself for new systems?

What would the model need to do to
quickly find optimal conditions?

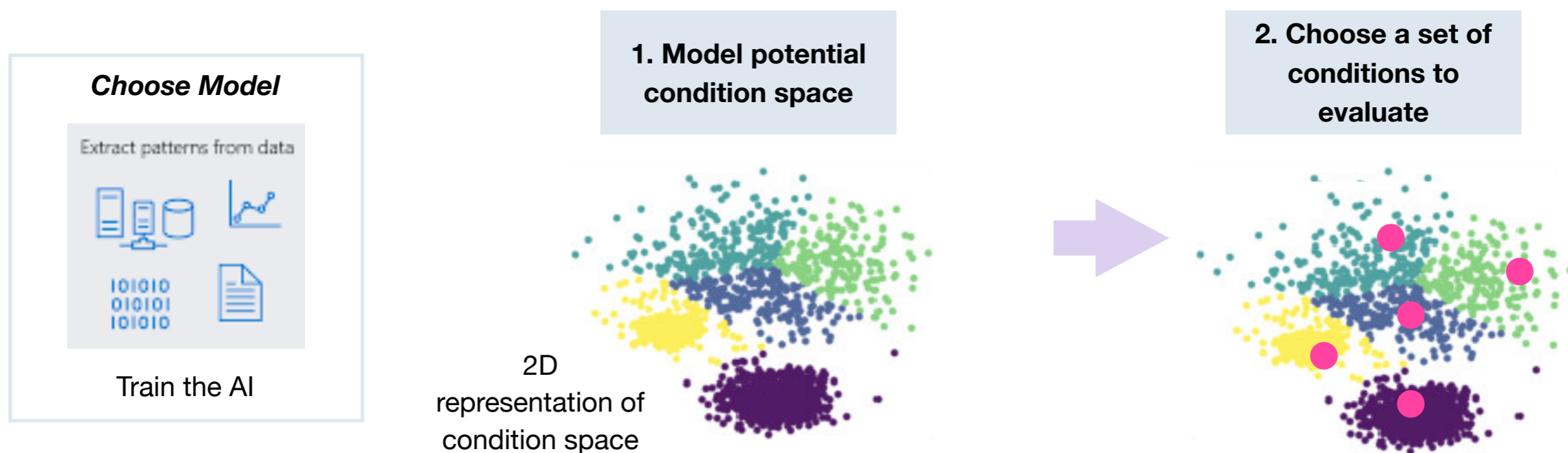
Machine Learning Applied to Chemistry - Reaction Optimization



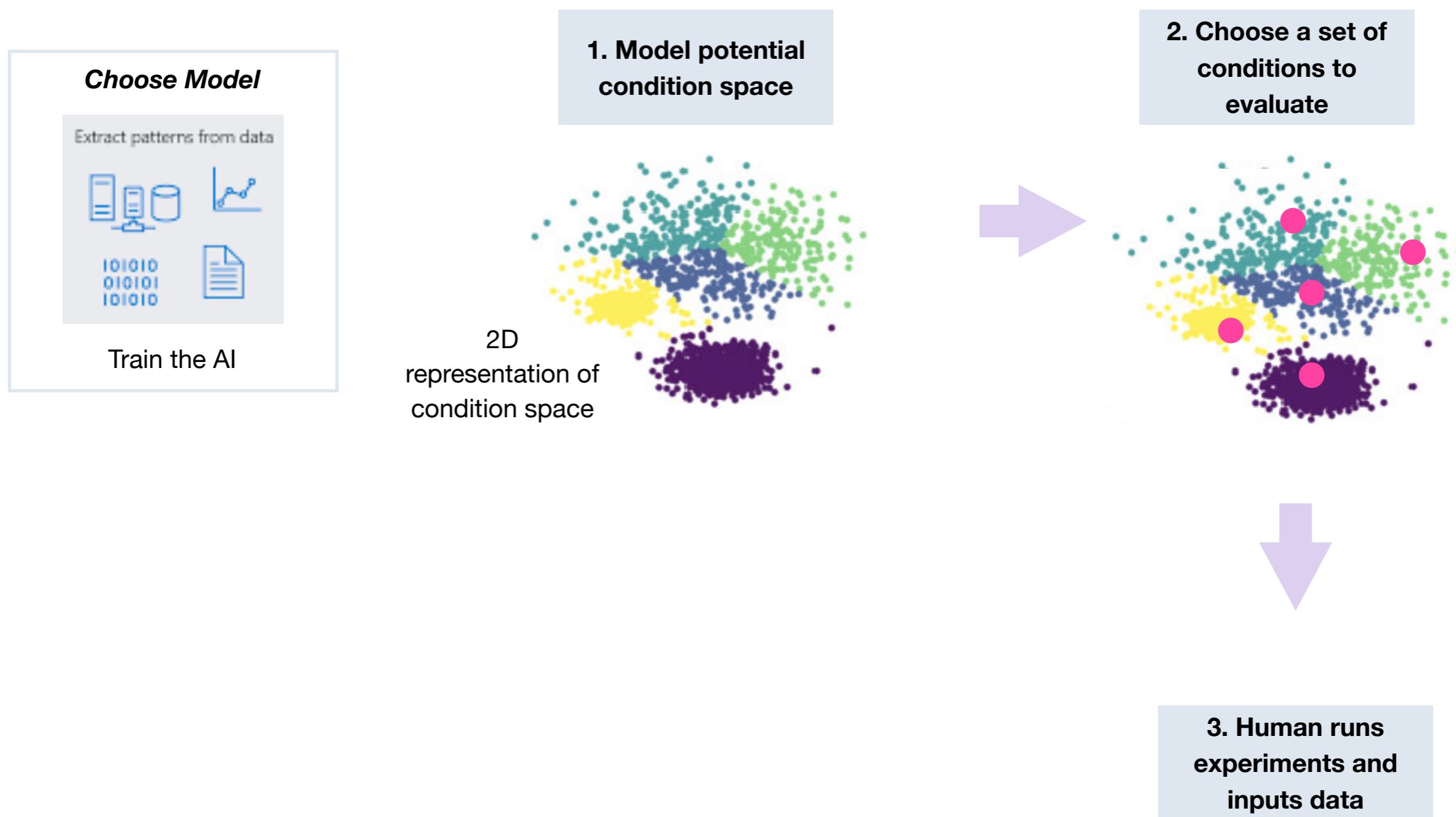
1. Model potential condition space



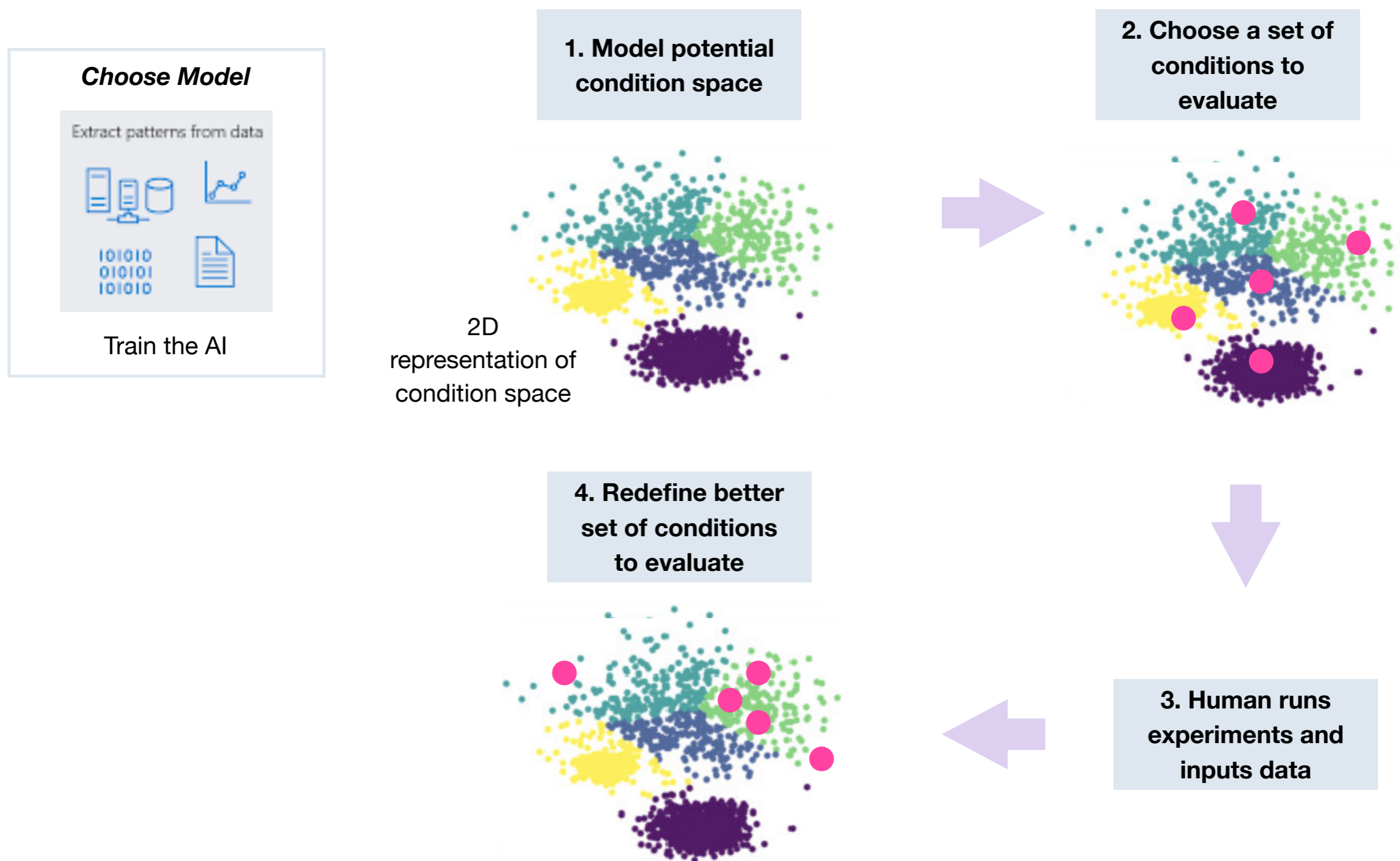
Machine Learning Applied to Chemistry - Reaction Optimization



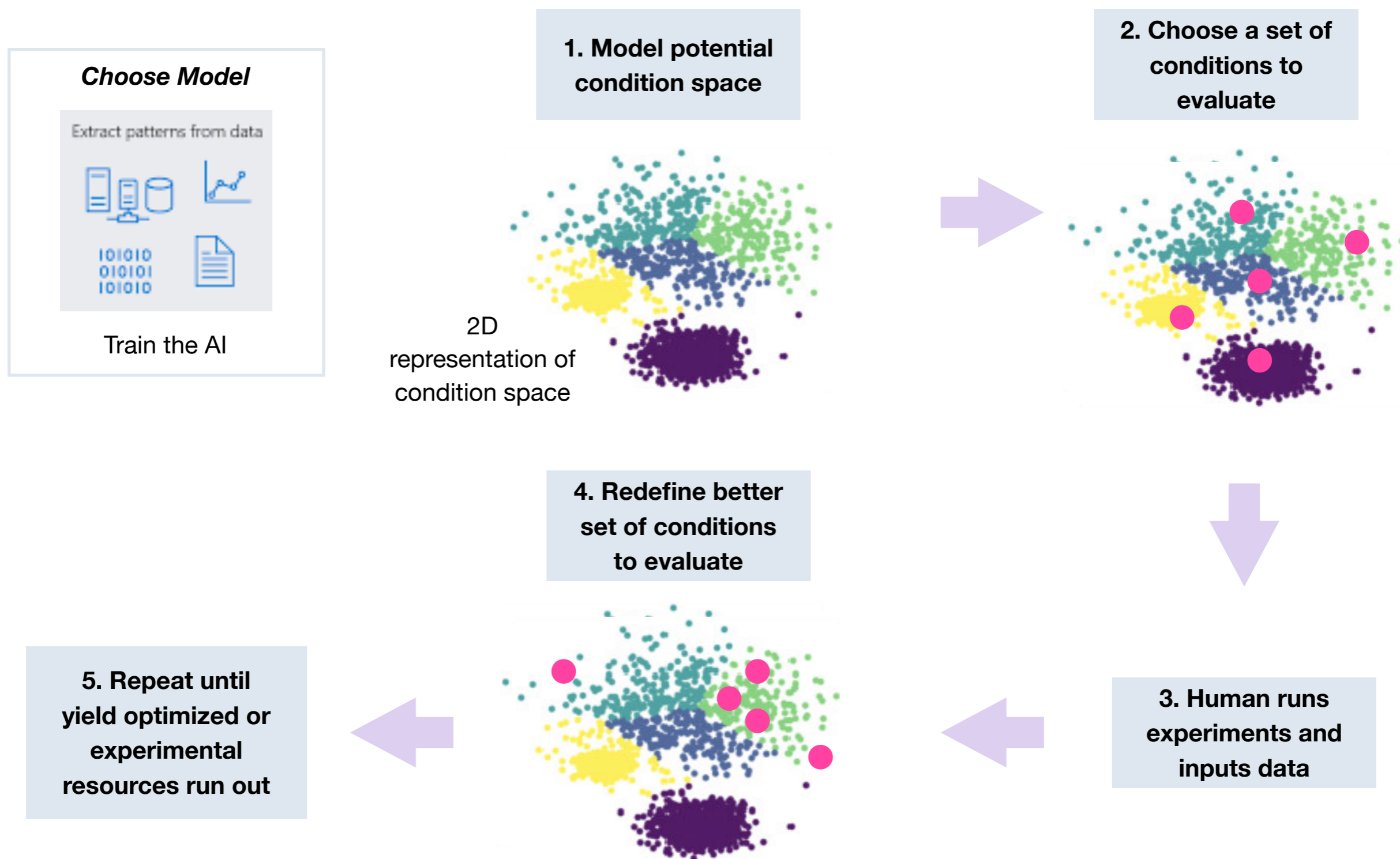
Machine Learning Applied to Chemistry - Reaction Optimization



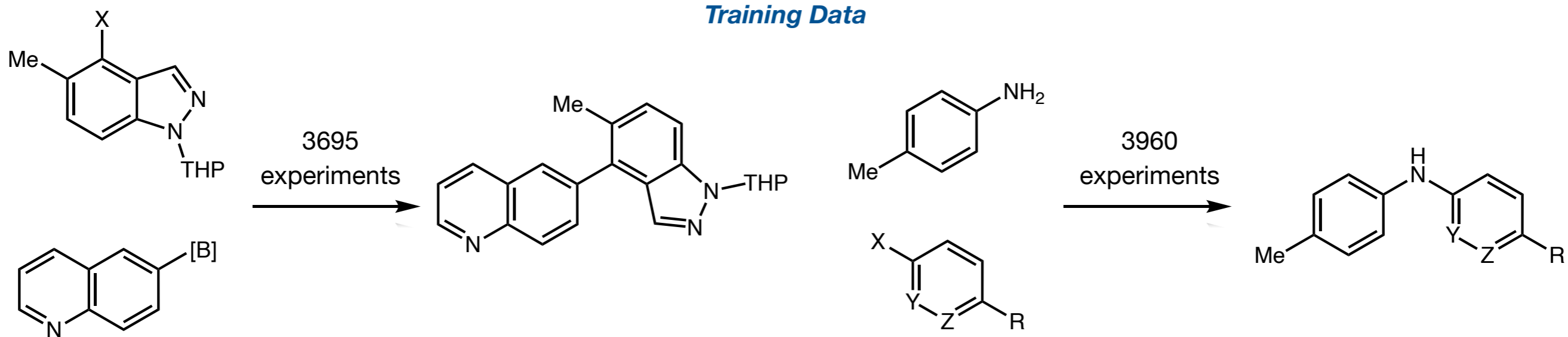
Machine Learning Applied to Chemistry - Reaction Optimization



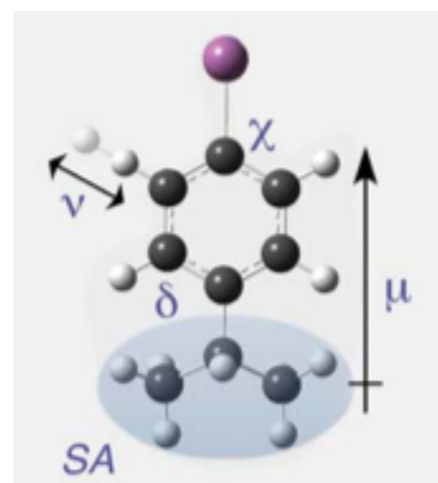
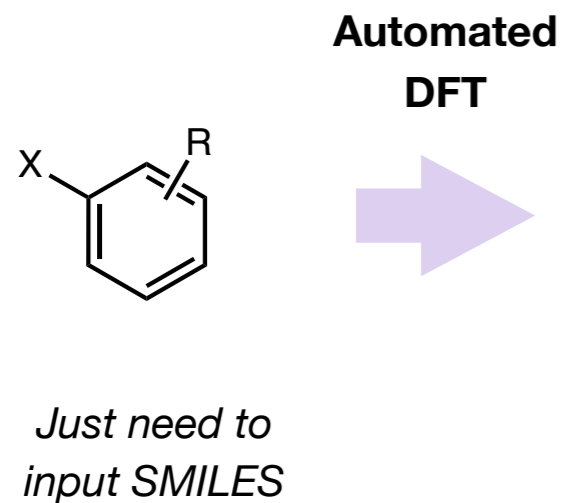
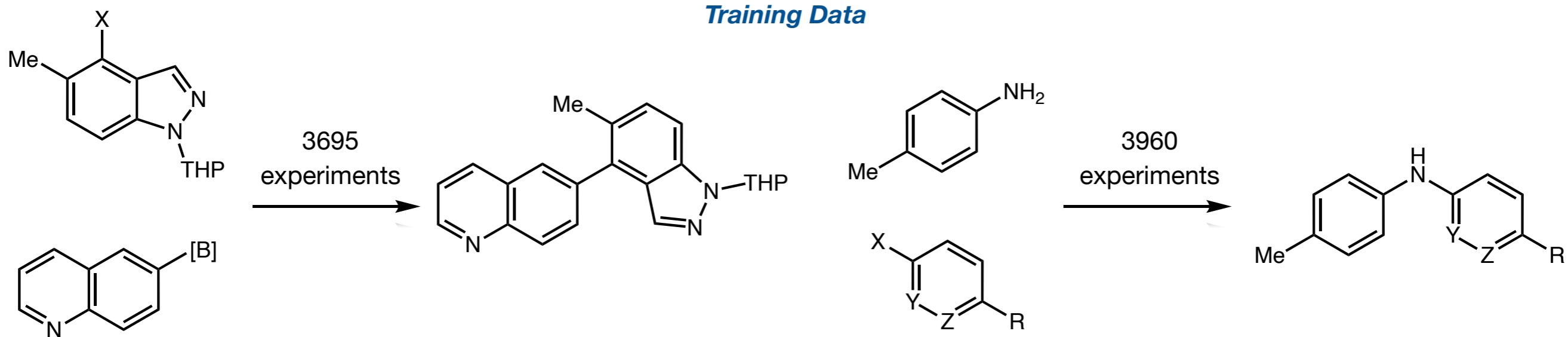
Machine Learning Applied to Chemistry - Reaction Optimization



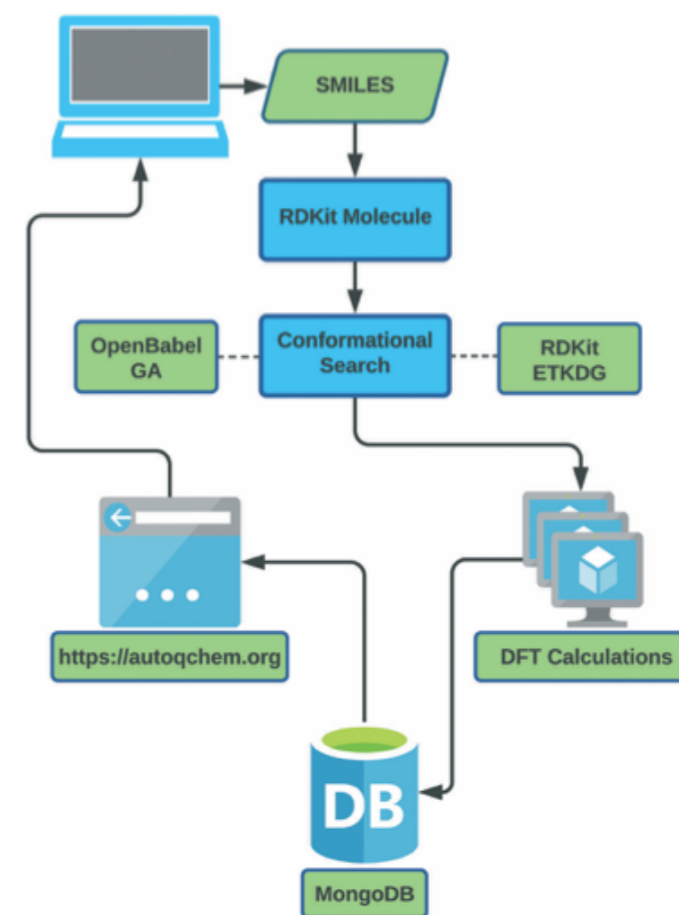
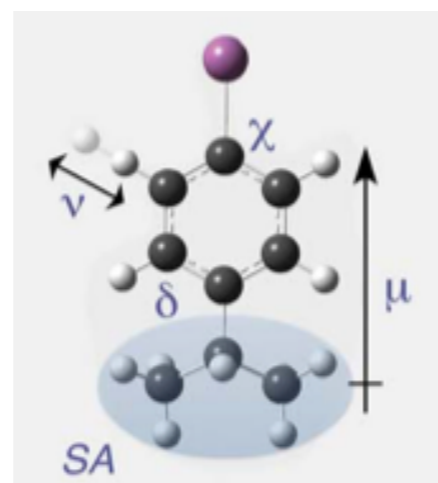
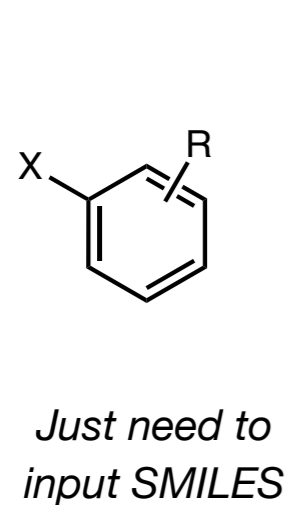
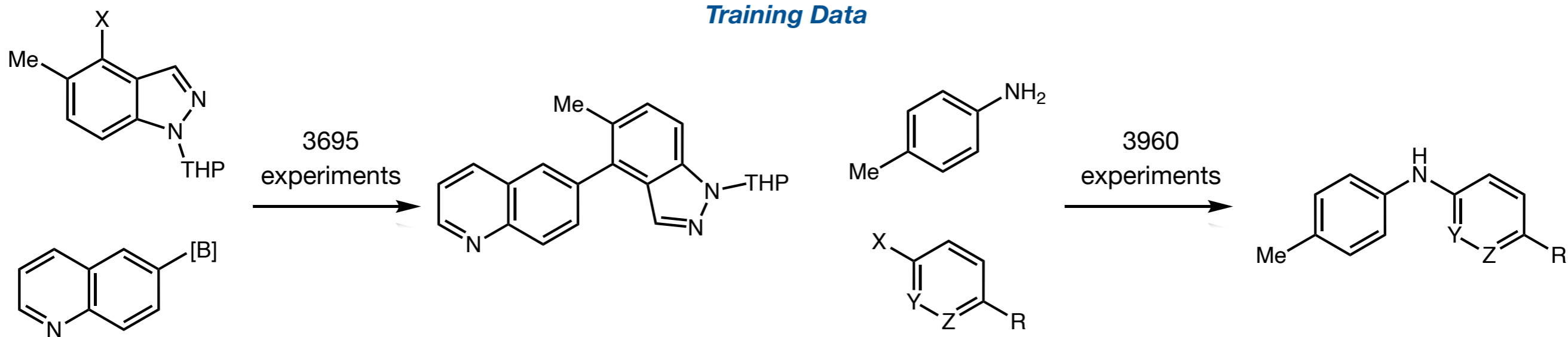
Machine Learning Applied to Chemistry - Reaction Optimization



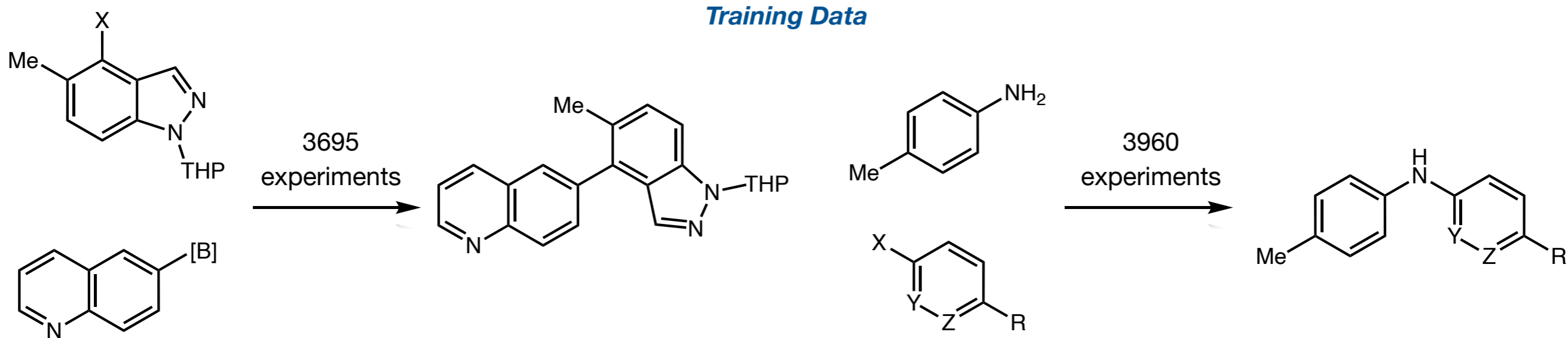
Machine Learning Applied to Chemistry - Reaction Optimization



Machine Learning Applied to Chemistry - Reaction Optimization



Machine Learning Applied to Chemistry - Reaction Optimization



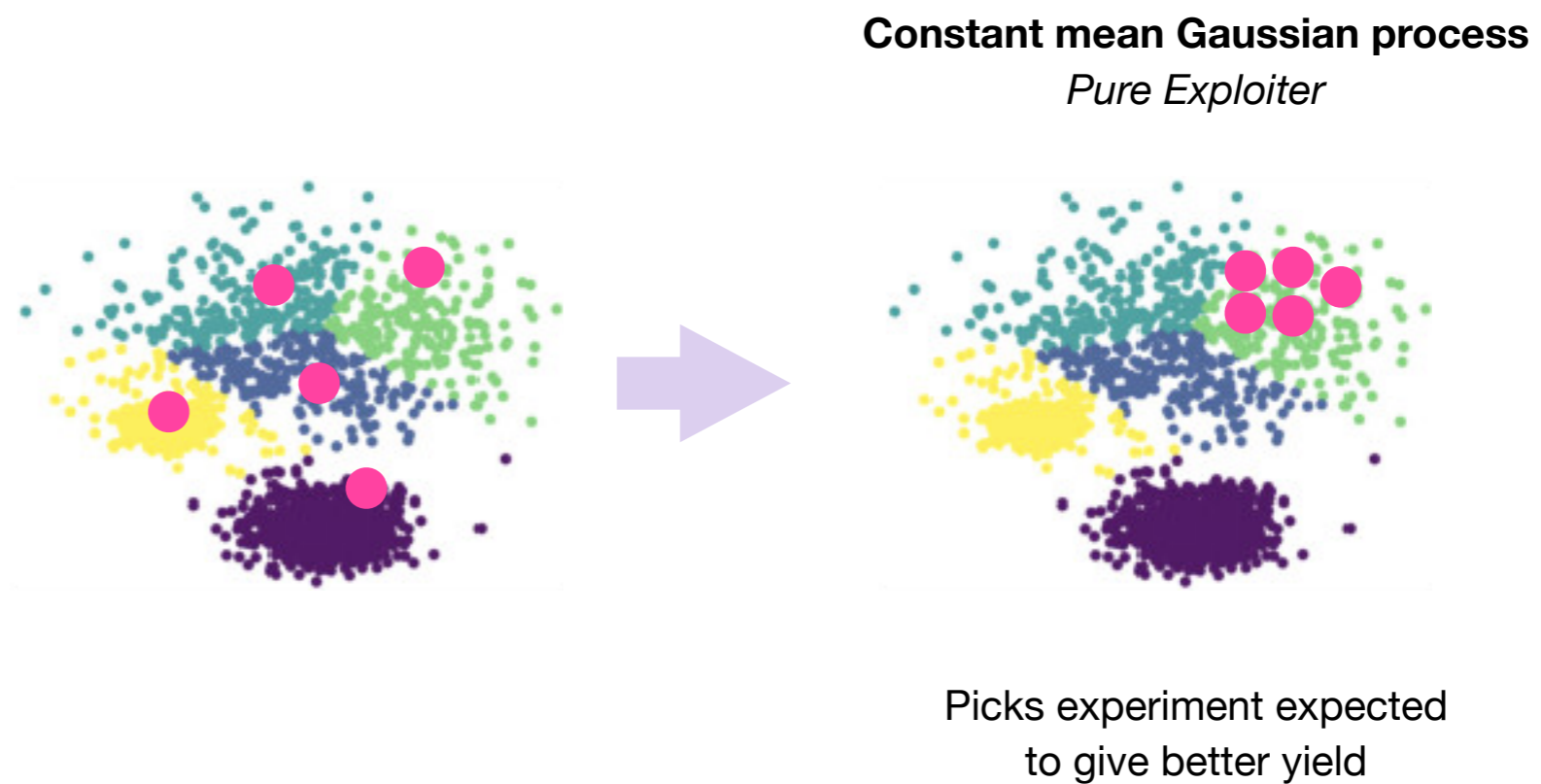
For more details about Auto-QChem and other applications see:

Zuranski, A. M.; Wang, J. Y.; Shields, B. J.; Doyle, A. G. *React. Chem. Eng.* **2022**, 7, 1276–1284

Machine Learning Applied to Chemistry - Reaction Optimization



Machine Learning Applied to Chemistry - Reaction Optimization



Machine Learning Applied to Chemistry - Reaction Optimization

Pioneering acquisition function
Pure Explorer



Picks experiment of greatest predictive uncertainty



Constant mean Gaussian process
Pure Exploiter



Picks experiment expected to give better yield

Machine Learning Applied to Chemistry - Reaction Optimization

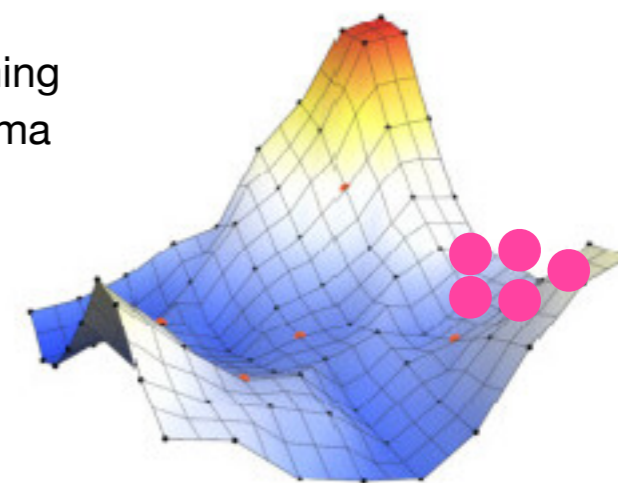
Pioneering acquisition function
Pure Explorer



Constant mean Gaussian process
Pure Exploiter



Susceptible to becoming trapped in local maxima



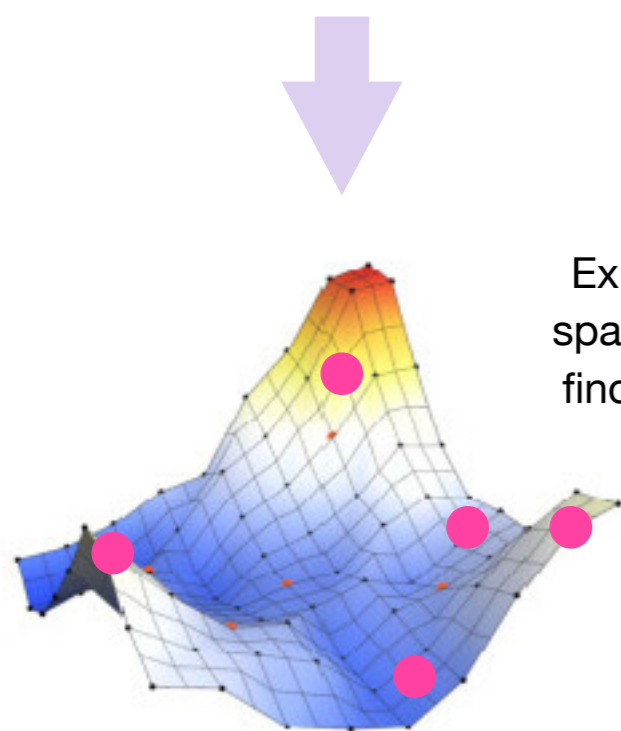
3D representation

Machine Learning Applied to Chemistry - Reaction Optimization

Pioneering acquisition function
Pure Explorer



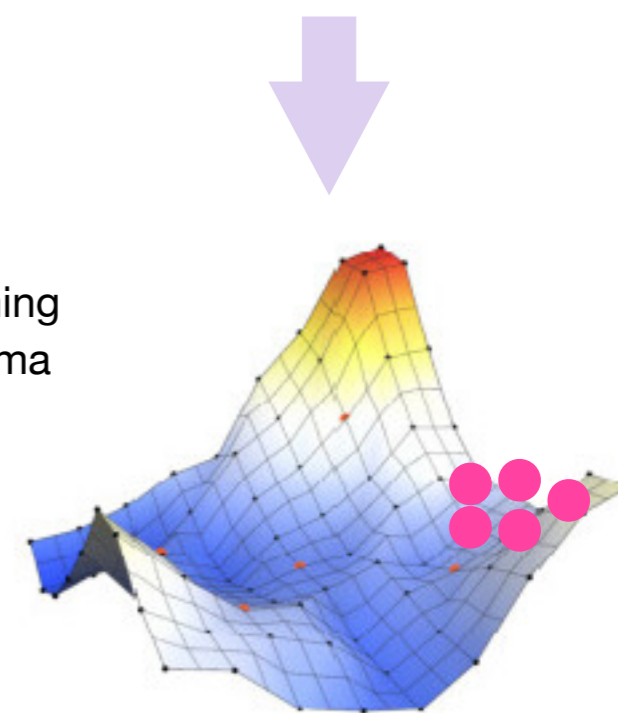
Constant mean Gaussian process
Pure Exploiter



Explores more condition space but not designed to find the absolute maxima

3D representation

Susceptible to becoming trapped in local maxima



3D representation

Machine Learning Applied to Chemistry - Reaction Optimization

Pioneering acquisition function
Pure Explorer



Constant mean Gaussian process
Pure Exploiter



Bayesian optimization with expected improvement
Mix of both



Machine Learning Applied to Chemistry - Reaction Optimization

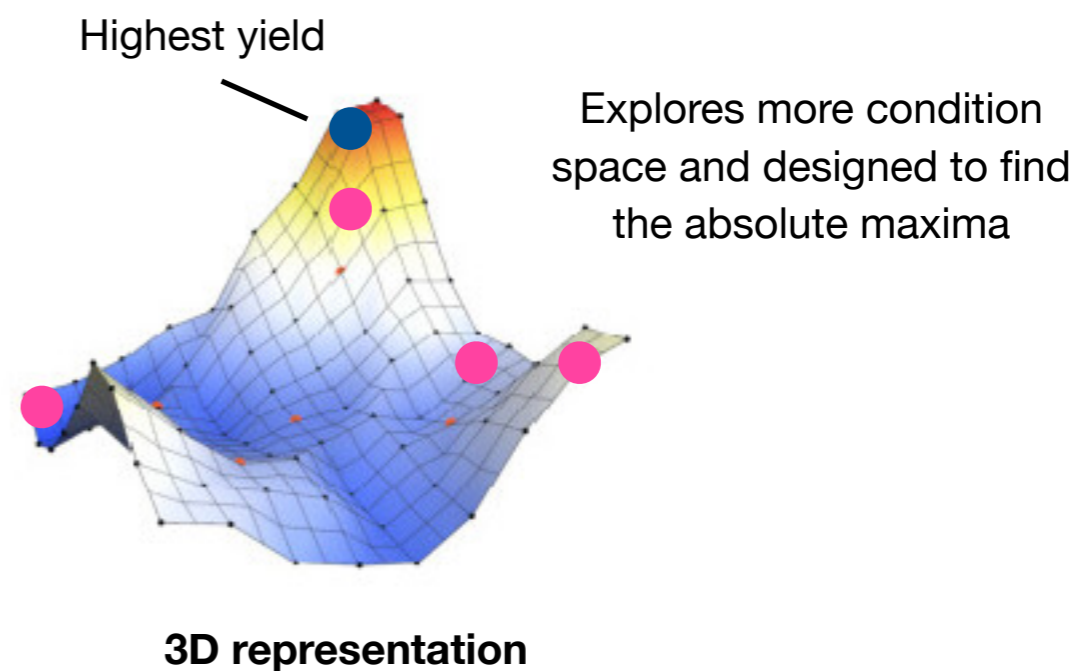
Pioneering acquisition function
Pure Explorer



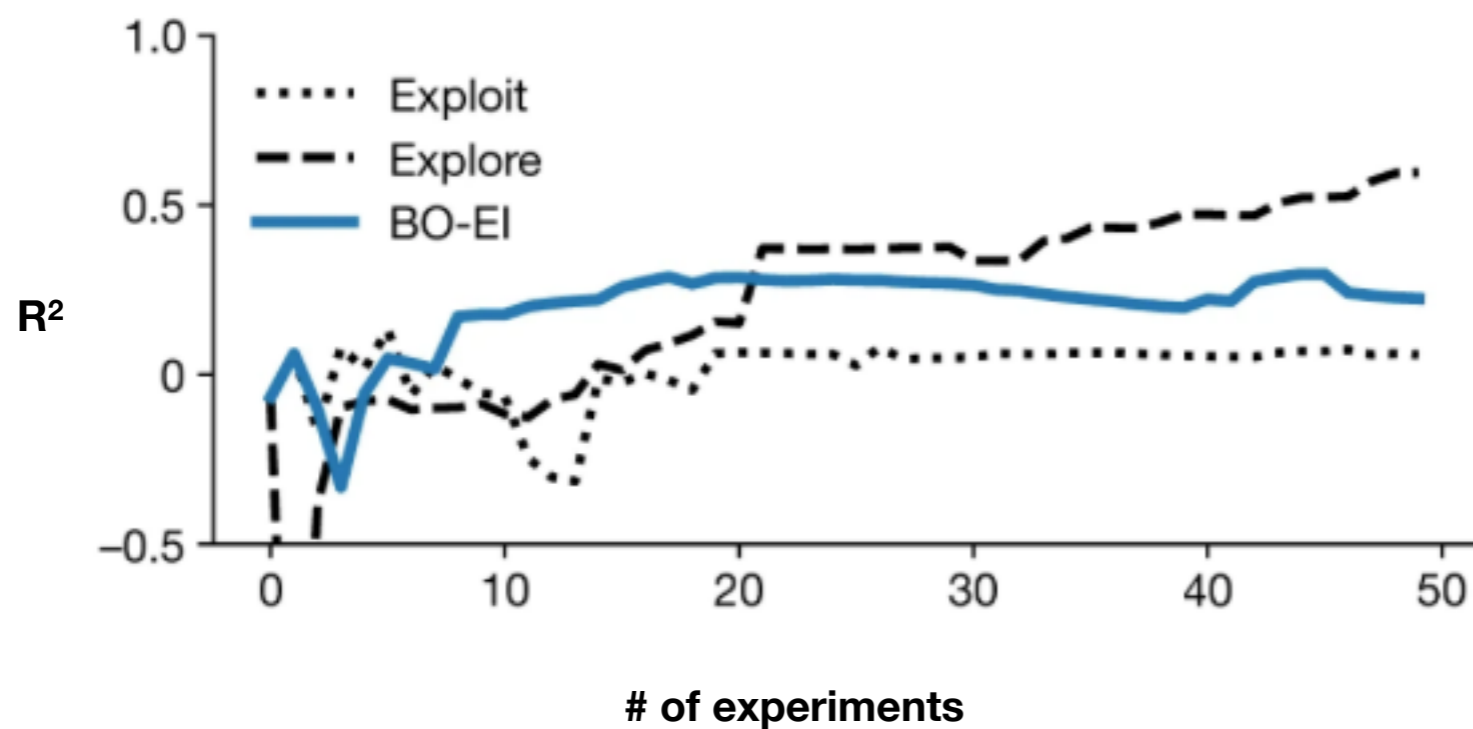
Constant mean Gaussian process
Pure Exploiter



Bayesian optimization with expected improvement
Mix of both

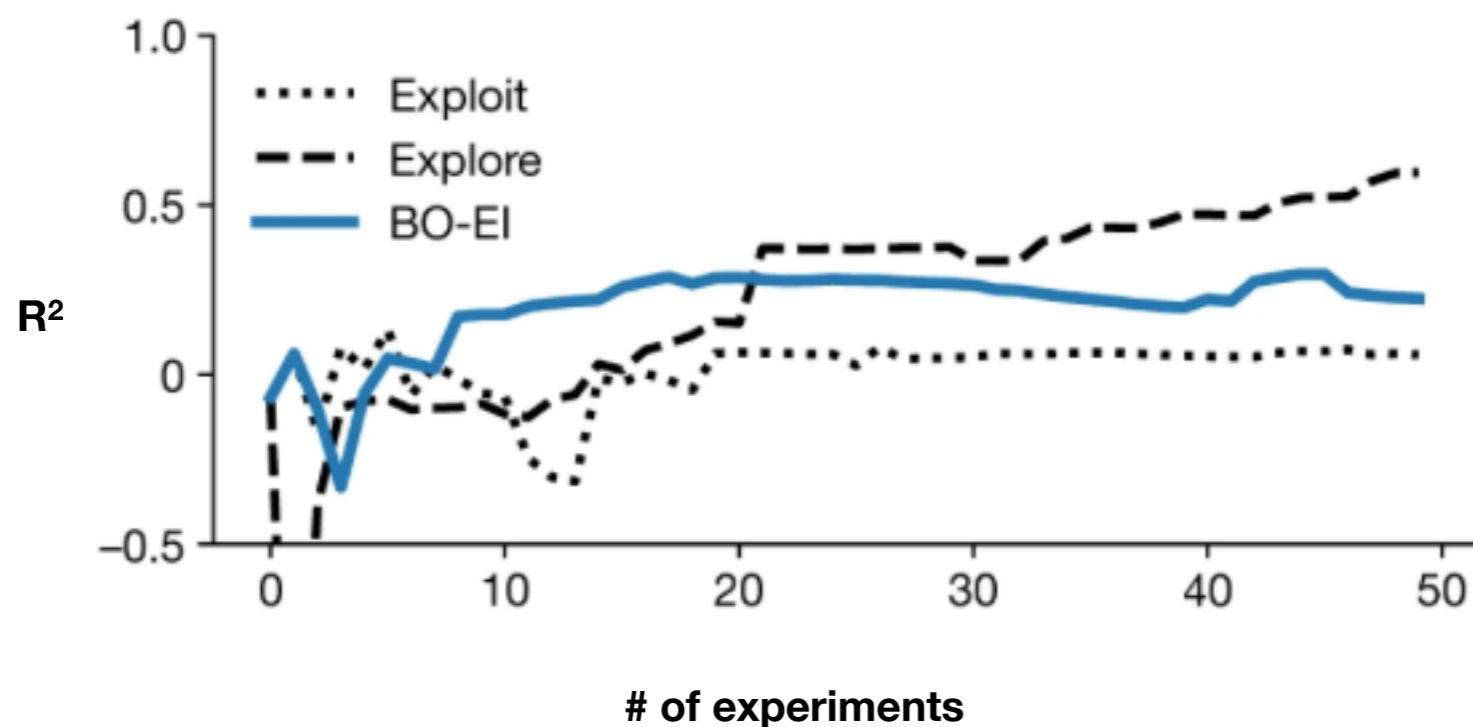


Machine Learning Applied to Chemistry - Reaction Optimization

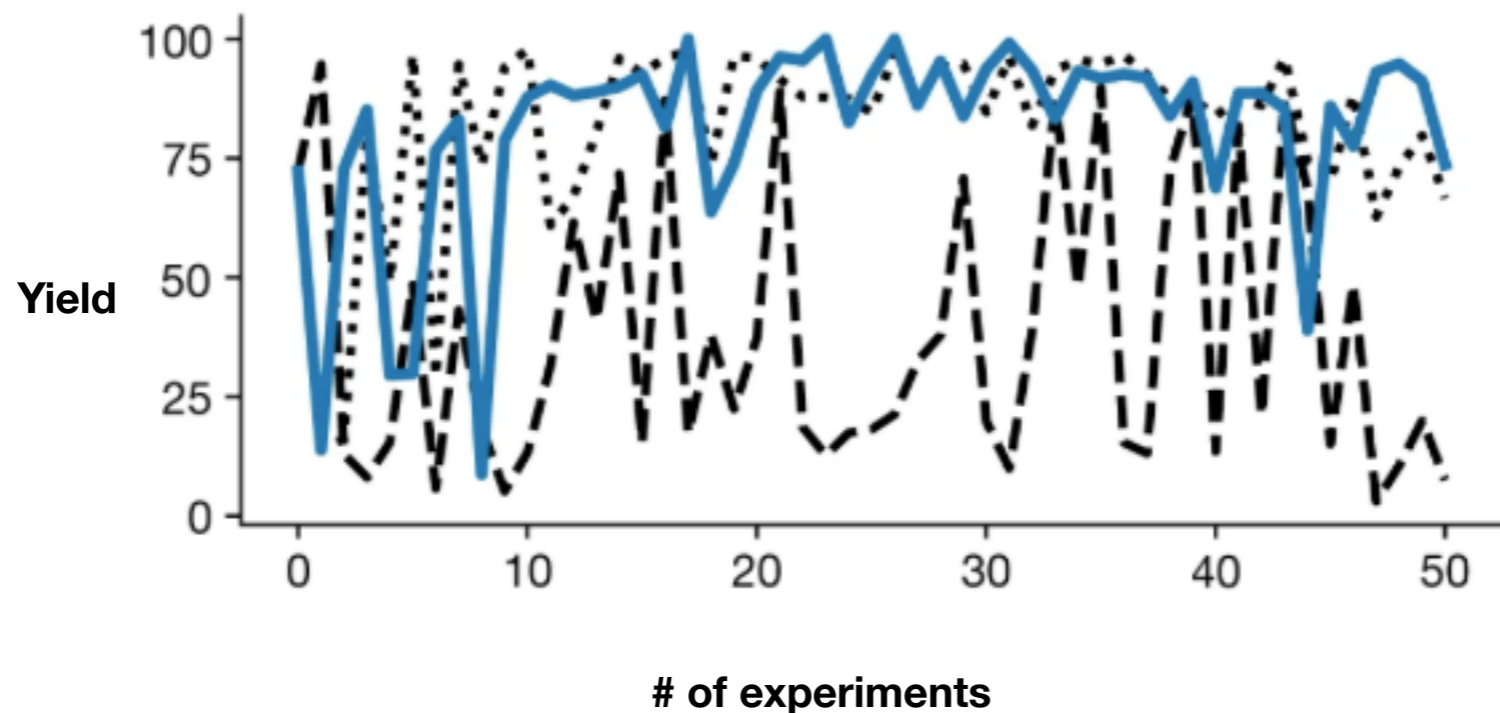


Explore algorithm and mixed algorithm have a much more accurate “understanding” of the reaction space

Machine Learning Applied to Chemistry - Reaction Optimization

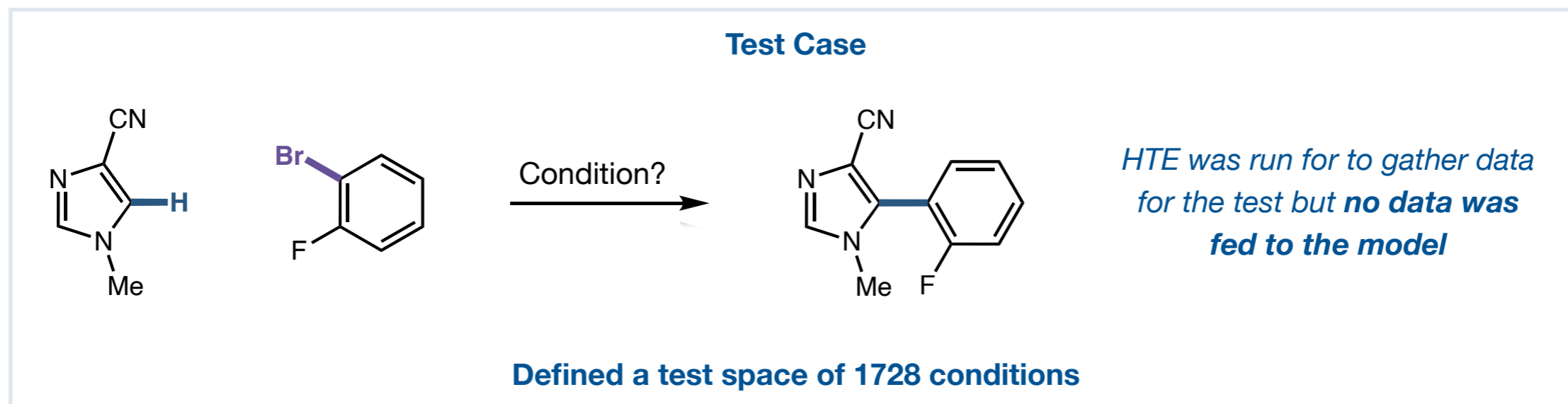


Explore algorithm and mixed algorithm have a much more accurate “understanding” of the reaction space

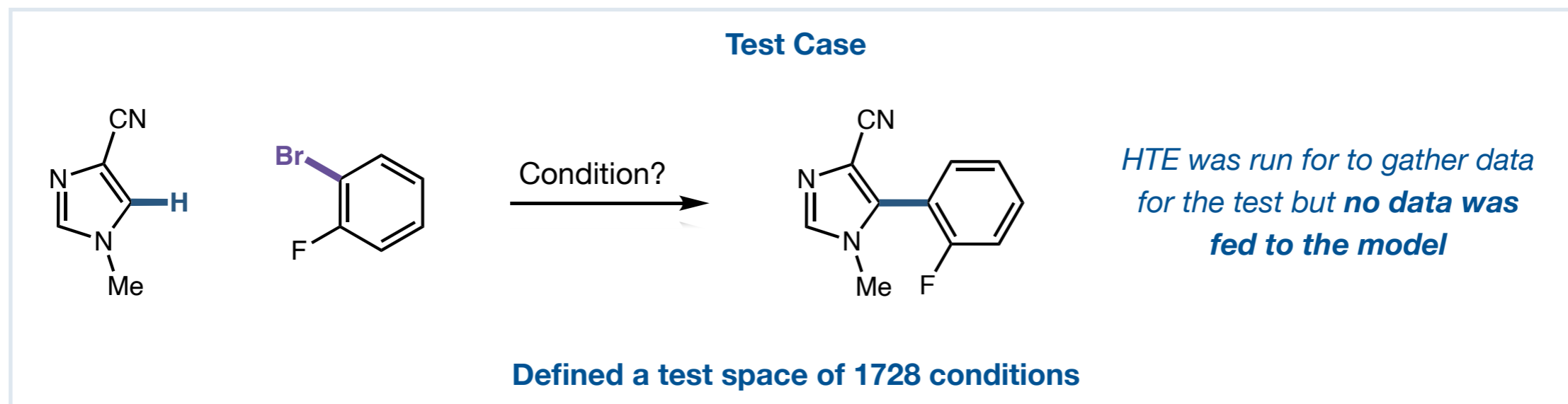


Explore algorithm doesn't care about highest yield and thus doesn't converge on optimal conditions as well as BO-EI

Machine Learning Applied to Chemistry - Reaction Optimization



Machine Learning Applied to Chemistry - Reaction Optimization



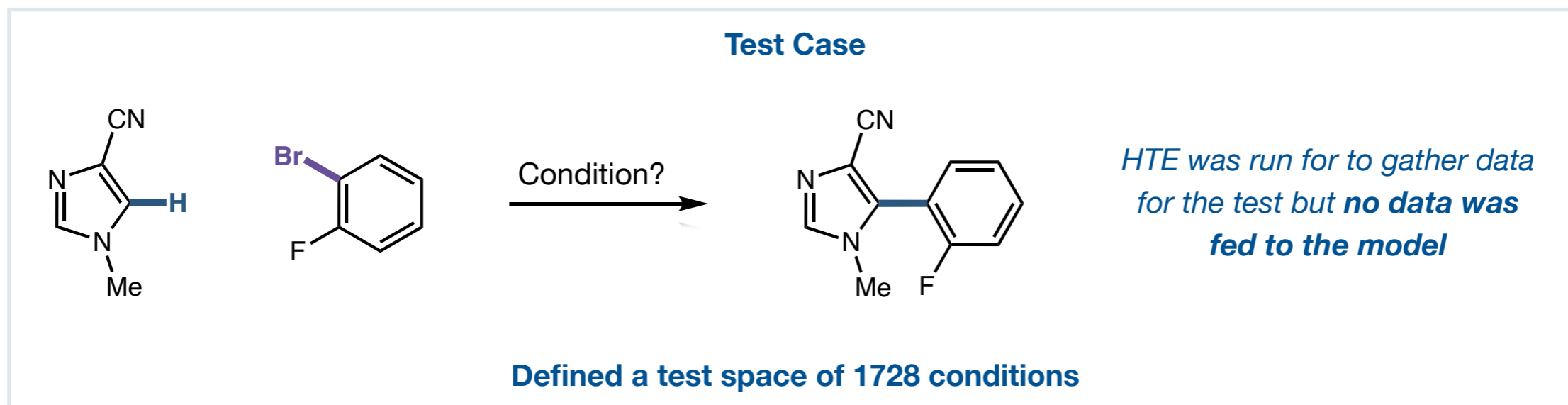
50 expert chemists

VS

50 run of model

Each player could submit 20 batches of 5 experiments, and between each batch get the results from the HTE data

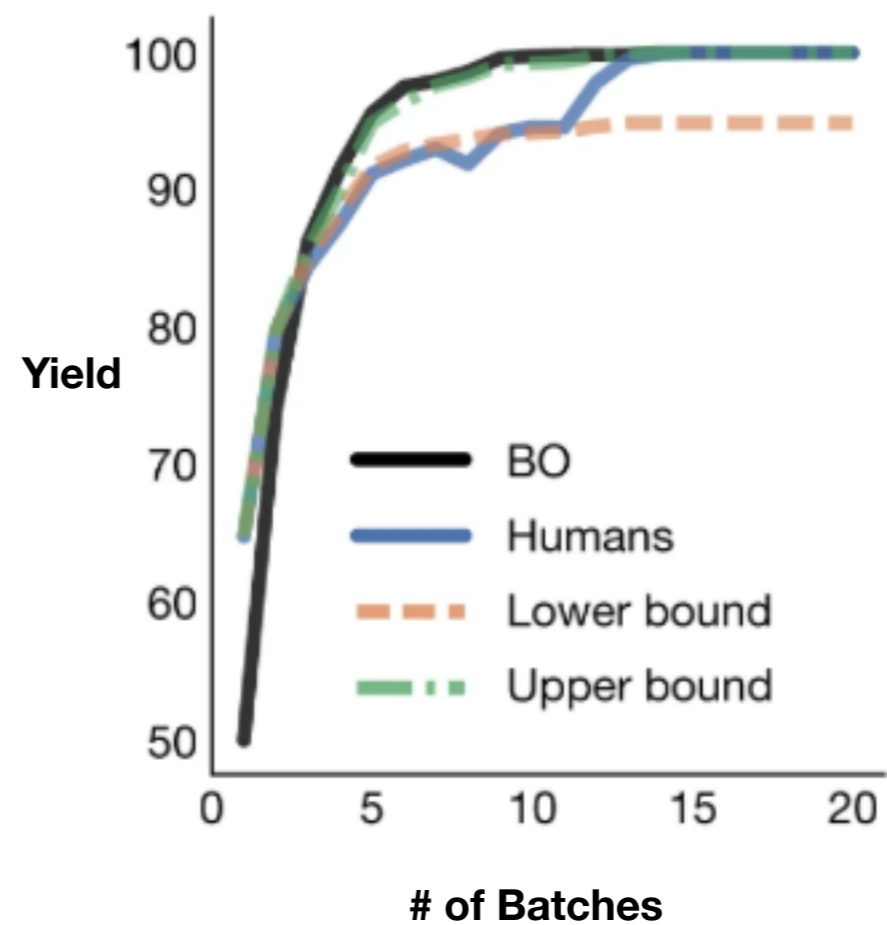
Machine Learning Applied to Chemistry - Reaction Optimization



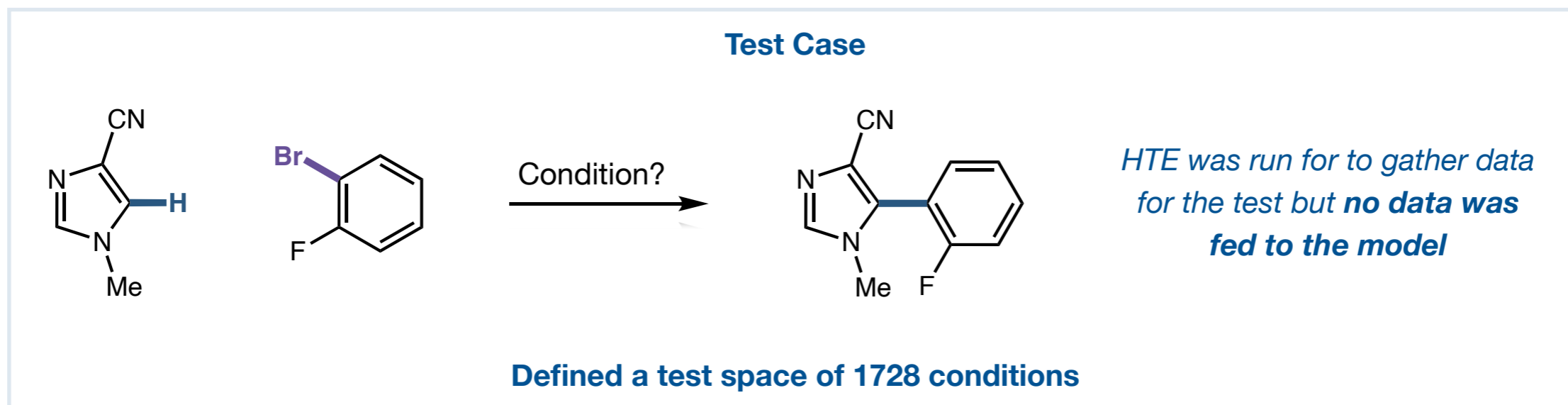
50 expert chemists

VS

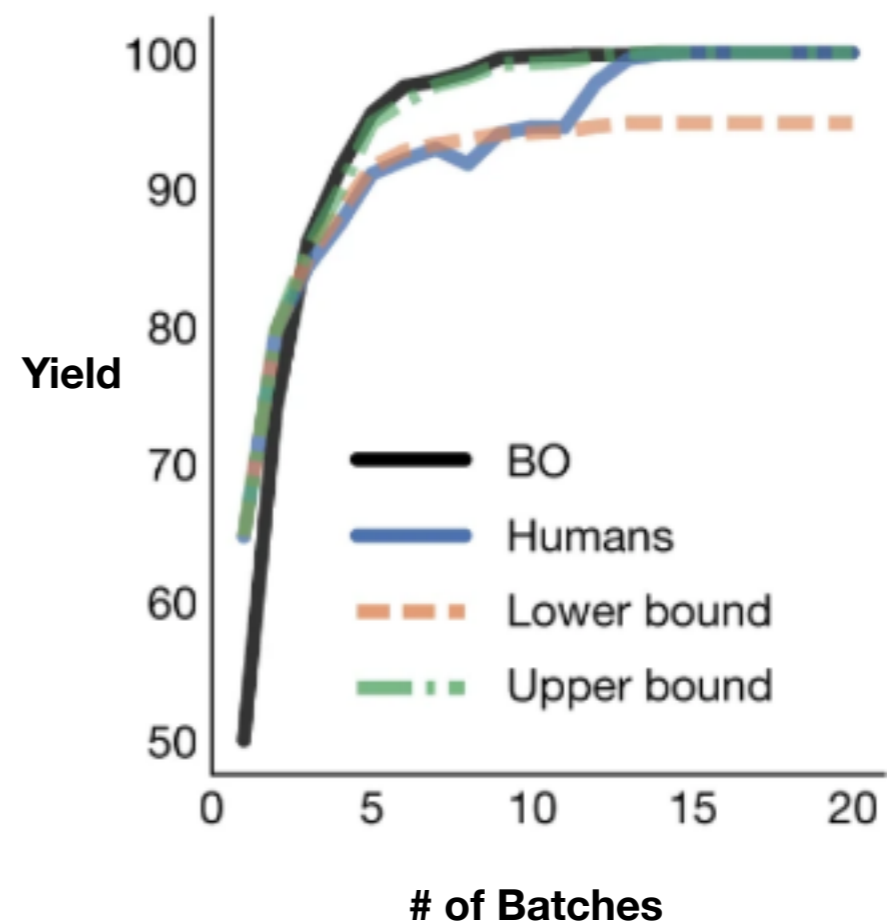
50 run of model



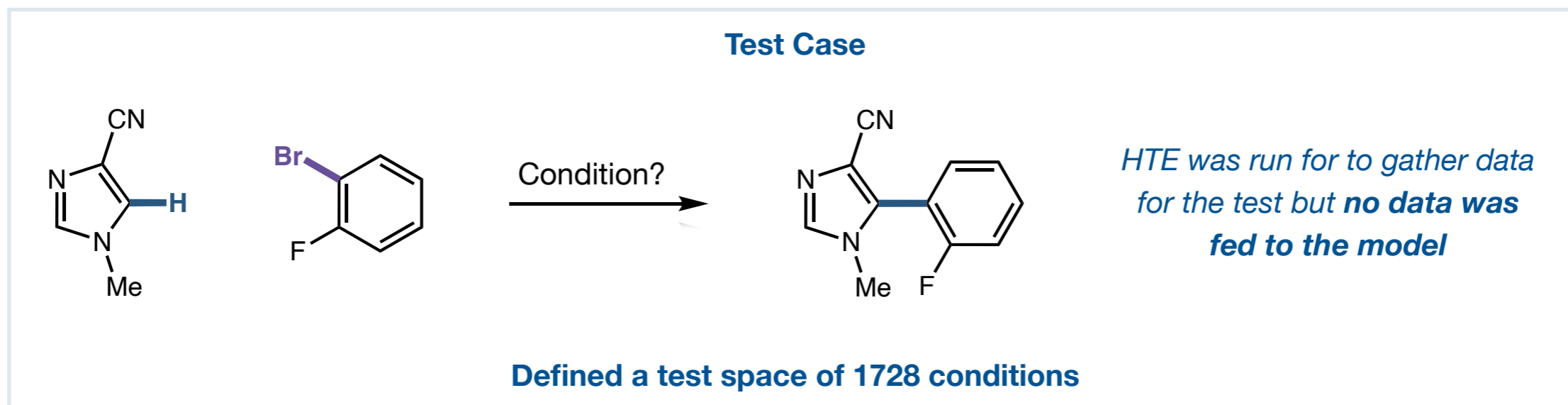
Machine Learning Applied to Chemistry - Reaction Optimization



On average the model outperformed experts by the 5th batch of experiments



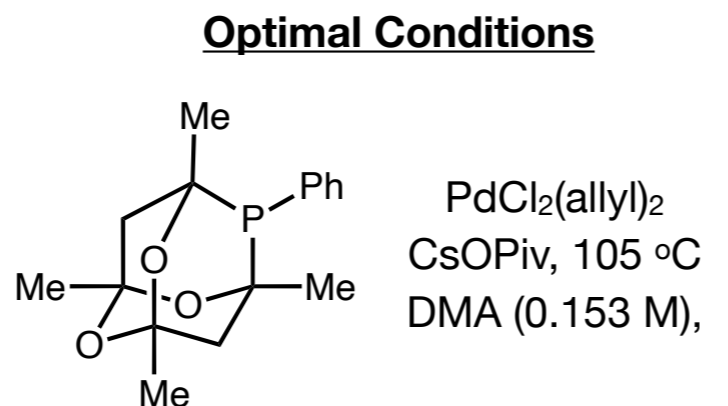
Machine Learning Applied to Chemistry - Reaction Optimization



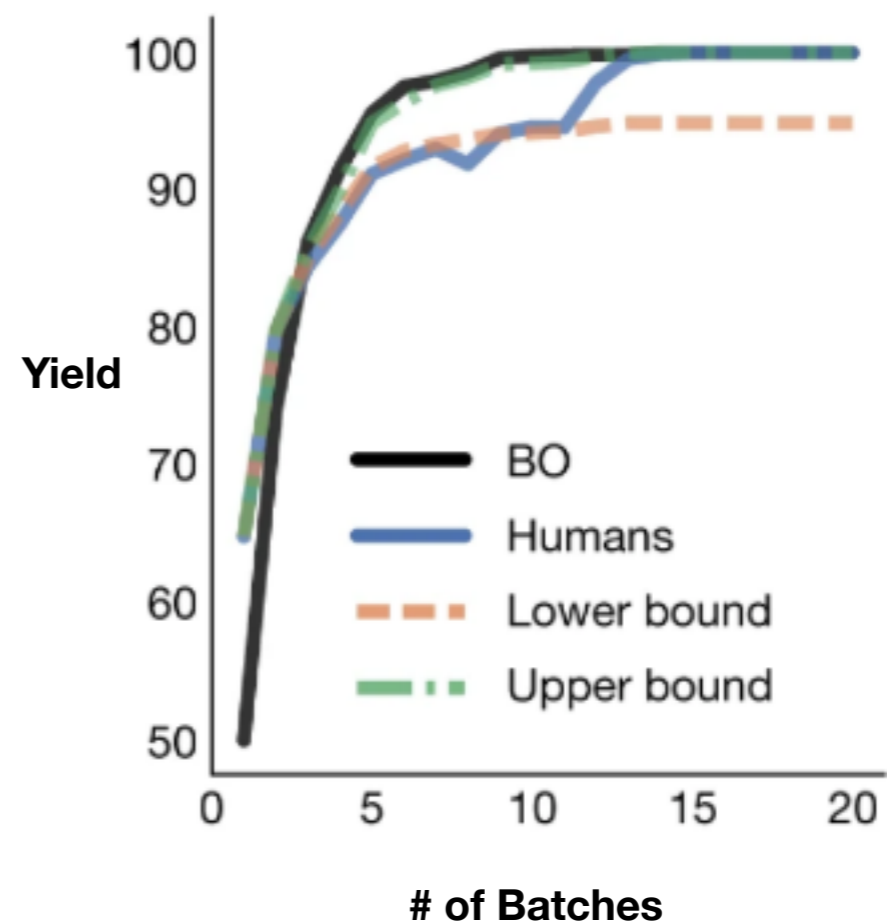
Unprecedented ligand for this class of reaction



Difficult for the experts to “find” this ligand

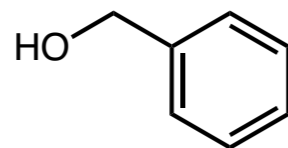
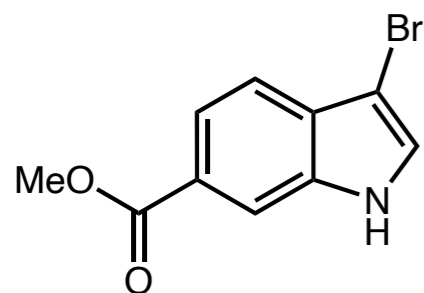


100% yield
Less than 50 experiments for all 50 different starting points

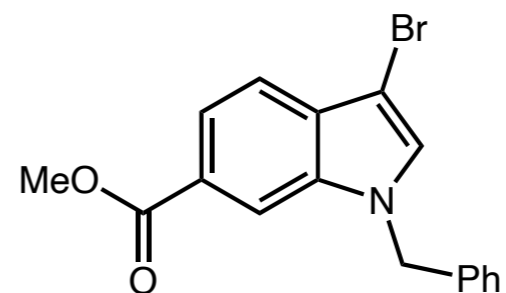


Machine Learning Applied to Chemistry - Reaction Optimization

Application to Mitsunobu reaction



Condition?
→
180,000 conditions

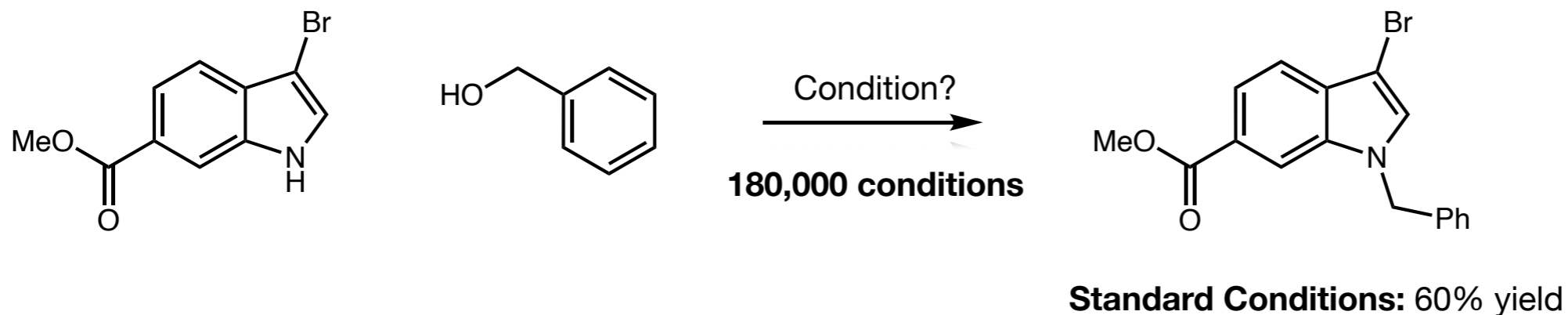


Standard Conditions: 60% yield

30 experiments
99% yield

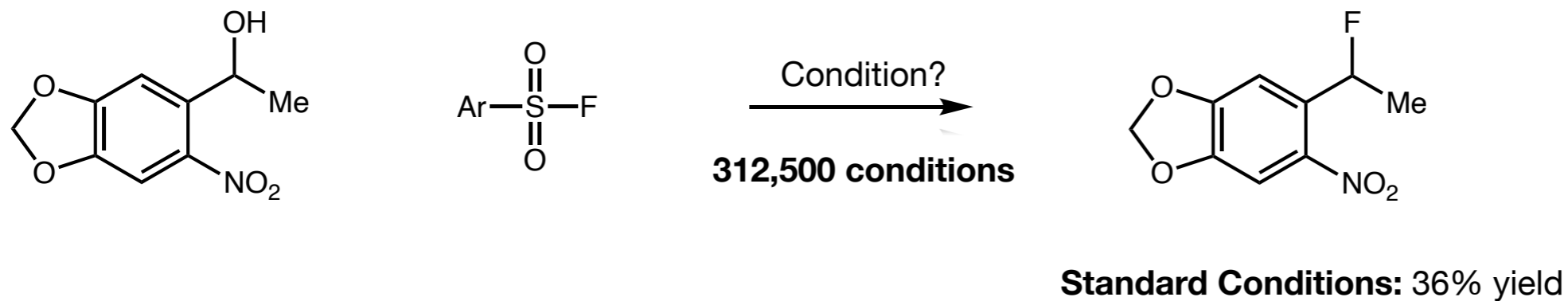
Machine Learning Applied to Chemistry - Reaction Optimization

Application to Mitsunobu reaction



30 experiments
99% yield

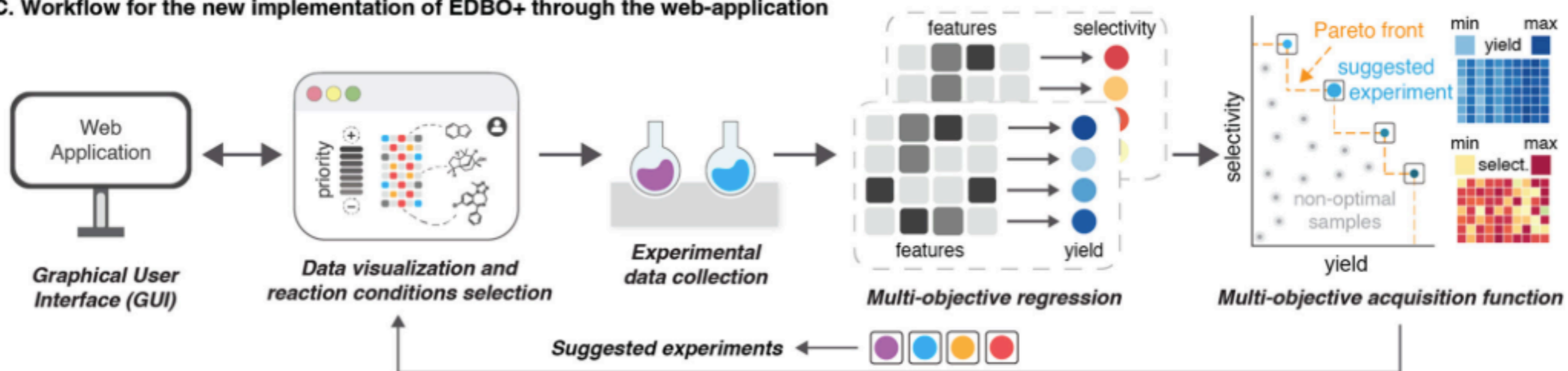
Application to deoxyfluorination reaction



40 experiments
70% yield

Machine Learning Applied to Chemistry - Reaction Optimization

C. Workflow for the new implementation of EDBO+ through the web-application

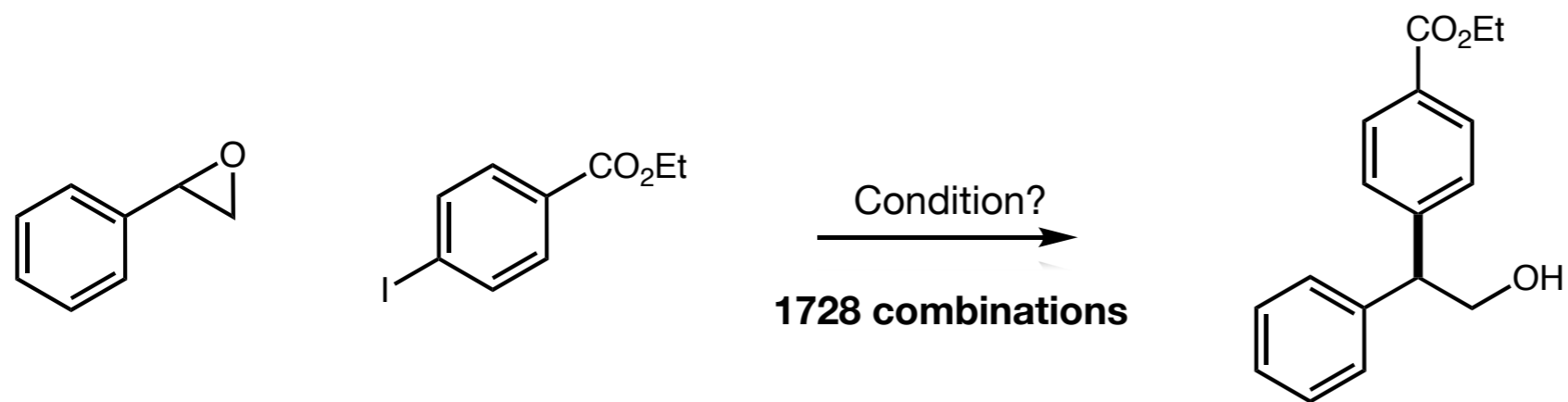


Website to use EDBO:

<https://www.edbowebapp.com>

Machine Learning Applied to Chemistry - Reaction Optimization

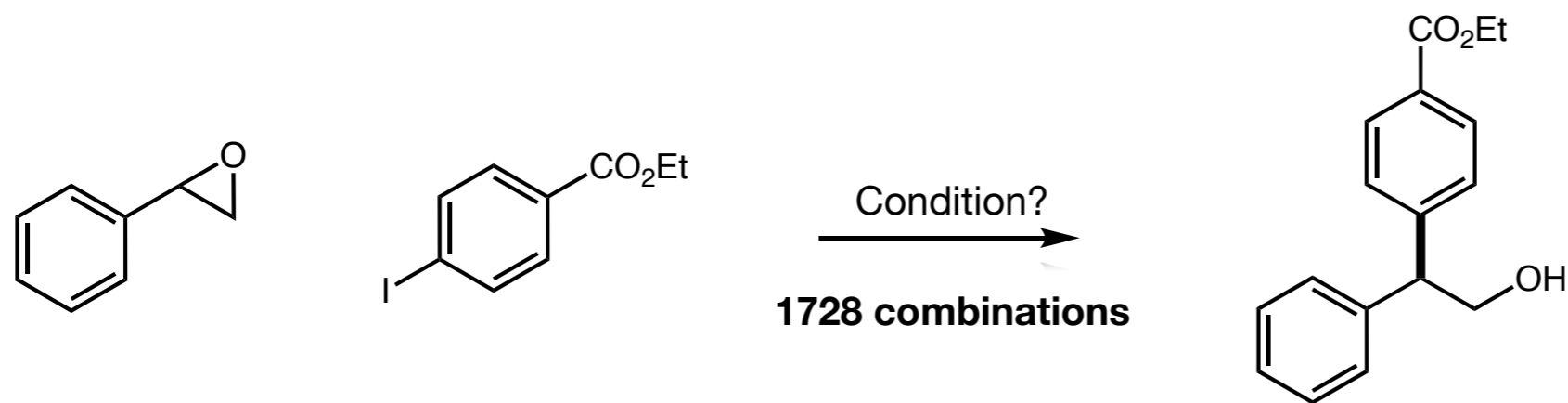
Bayesian Optimization applied to enantioselective photocatalytic systems



Human optimization: ~500 reactions for 70% yield, 80% ee

Machine Learning Applied to Chemistry - Reaction Optimization

Bayesian Optimization applied to enantioselective photocatalytic systems

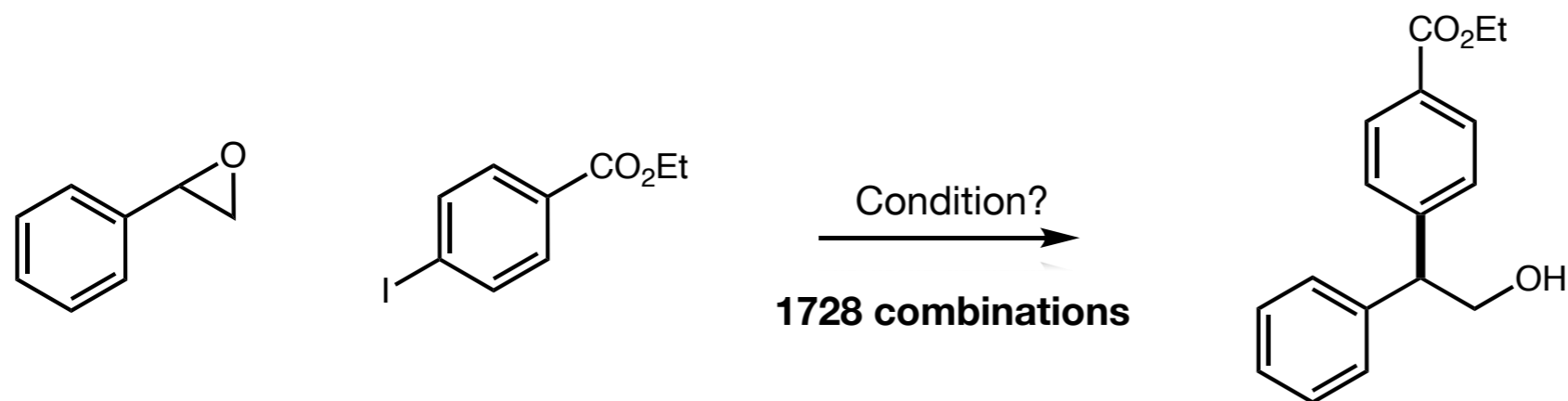


24 experiments
80% yield, 91% ee

Human optimization: ~500 reactions for 70% yield, 80% ee

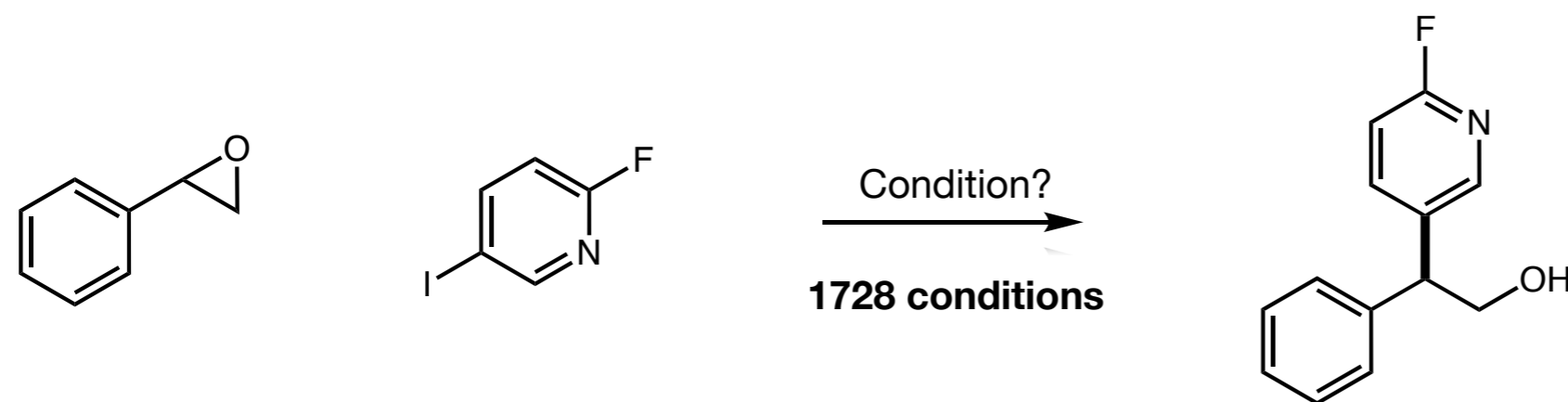
Machine Learning Applied to Chemistry - Reaction Optimization

Bayesian Optimization applied to enantioselective photocatalytic systems



24 experiments
80% yield, 91% ee

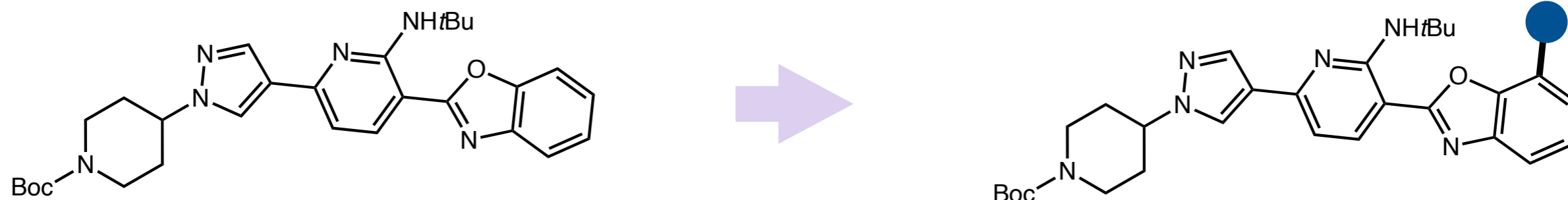
Human optimization: ~500 reactions for 70% yield, 80% ee



15 experiments
59% yield, 77% ee

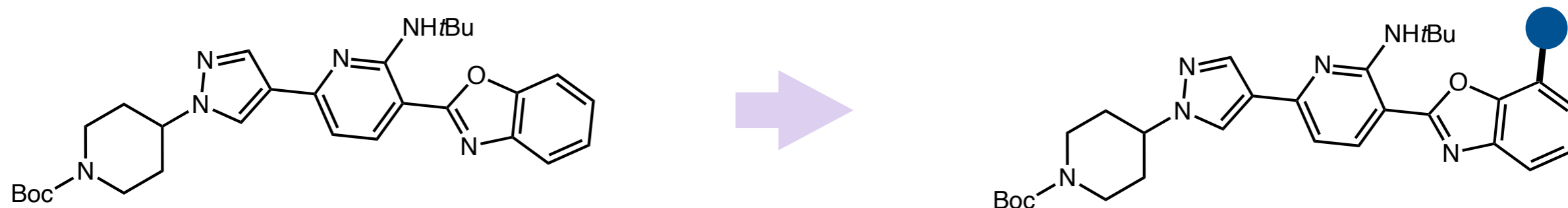
Human optimization: 49% yield, 76% ee

Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

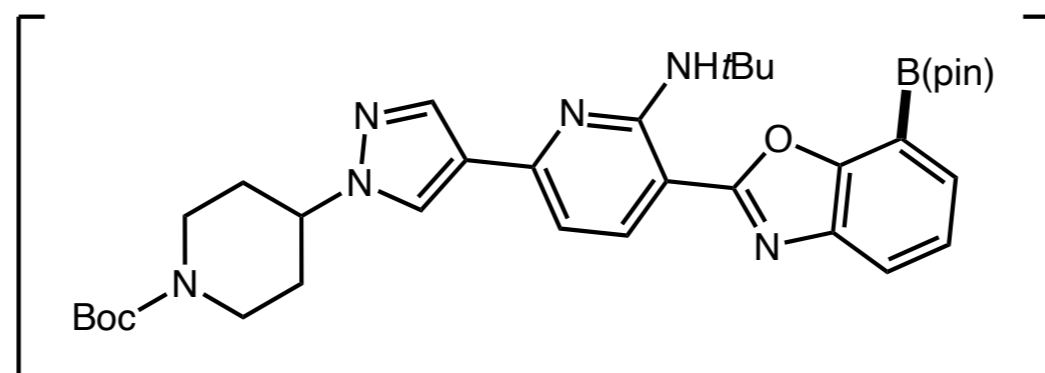


Theoretical Question: How can we access substitution at this position to fill a binding pocket?

Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

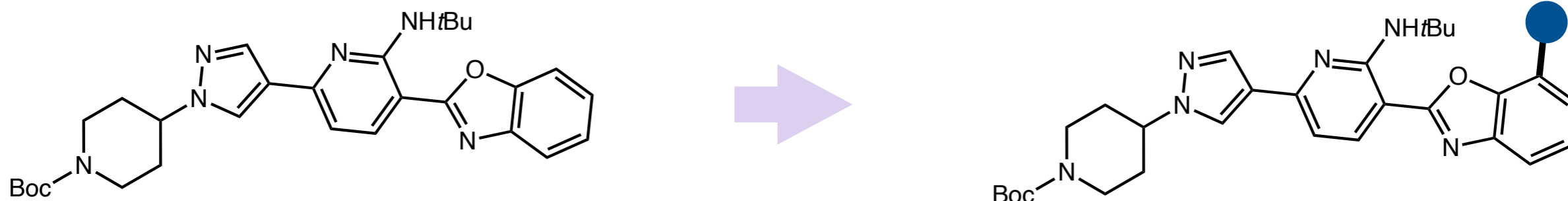


Theoretical Question: How can we access substitution at this position to fill a binding pocket?



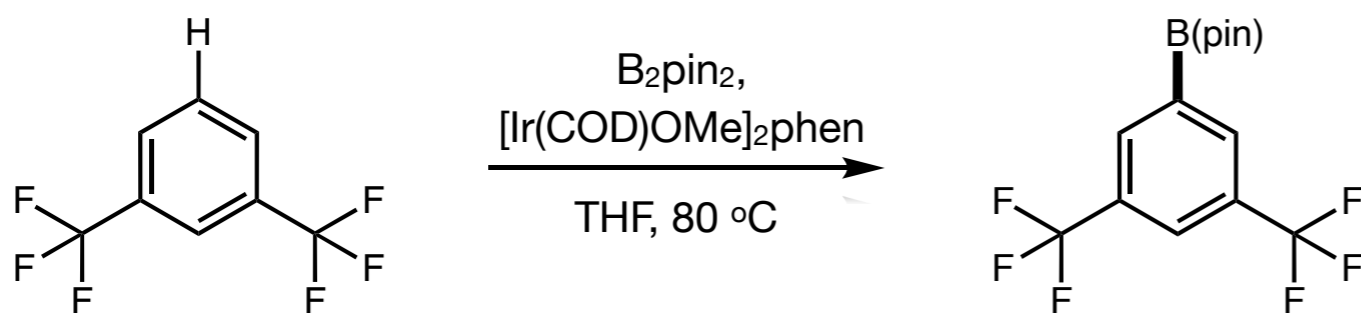
Late-stage functionalization to lynchpin intermediate

Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings



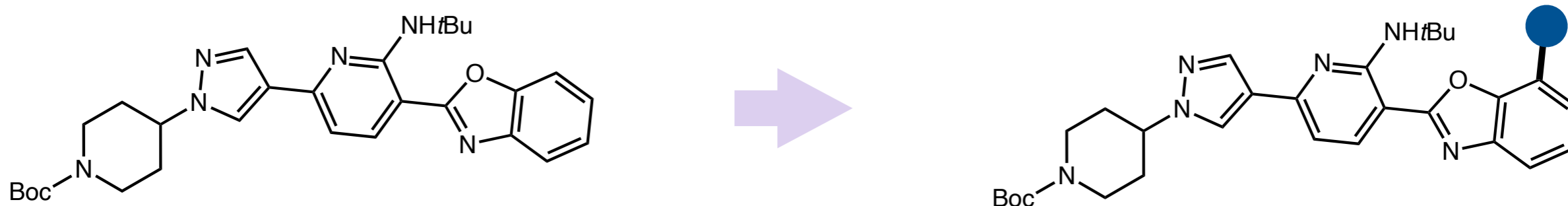
Theoretical Question: How can we access substitution at this position to fill a binding pocket?

Hartwig C-H Borylation

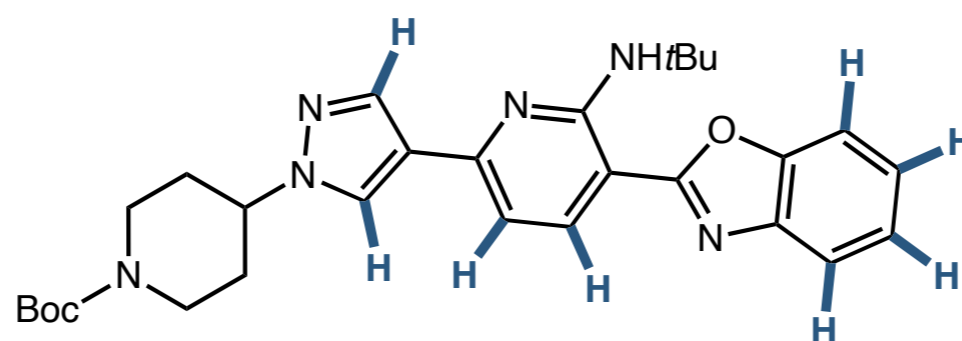


For simple systems selectivity can be a priori predicted by empirical knowledge of selectivity from other substrates

Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

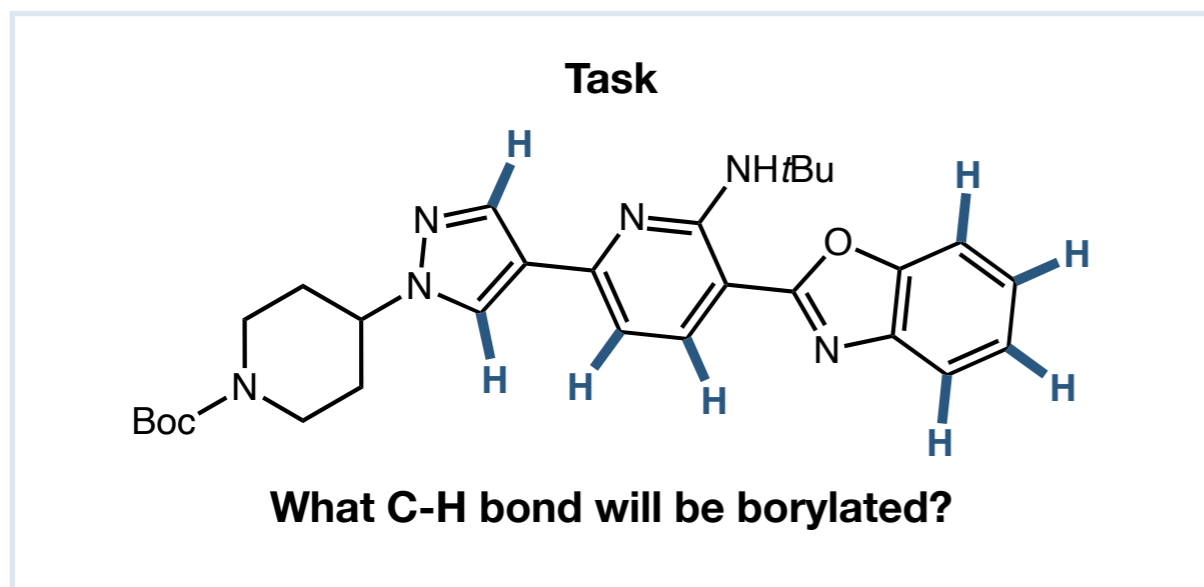


Theoretical Question: How can we access substitution at this position to fill a binding pocket?

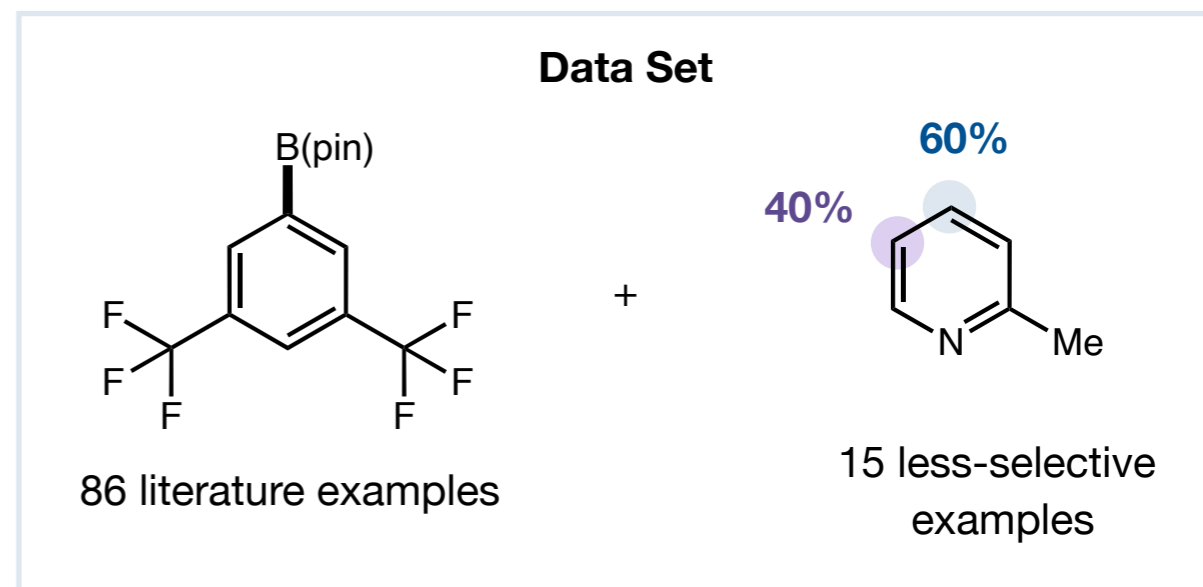
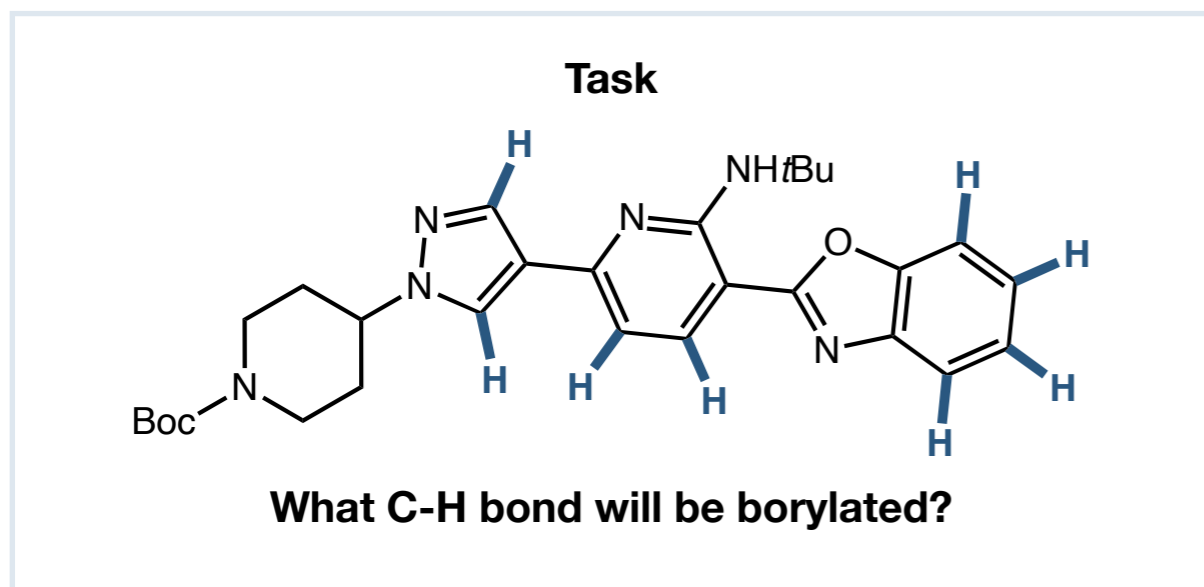


Many C-H bonds: predicting selectivity is not trivial

Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

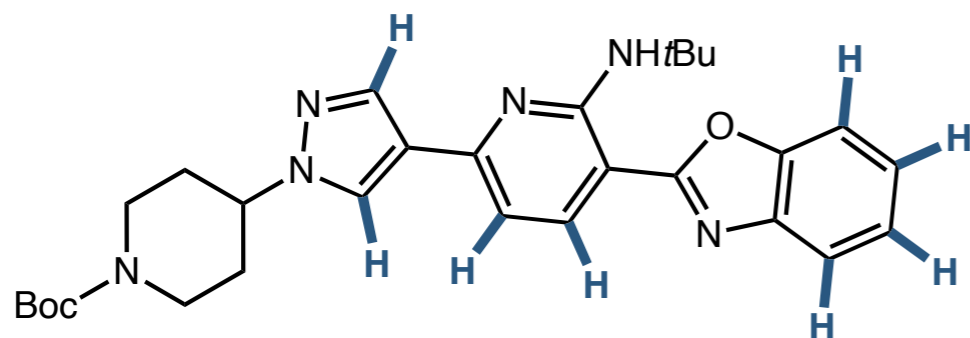


Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings



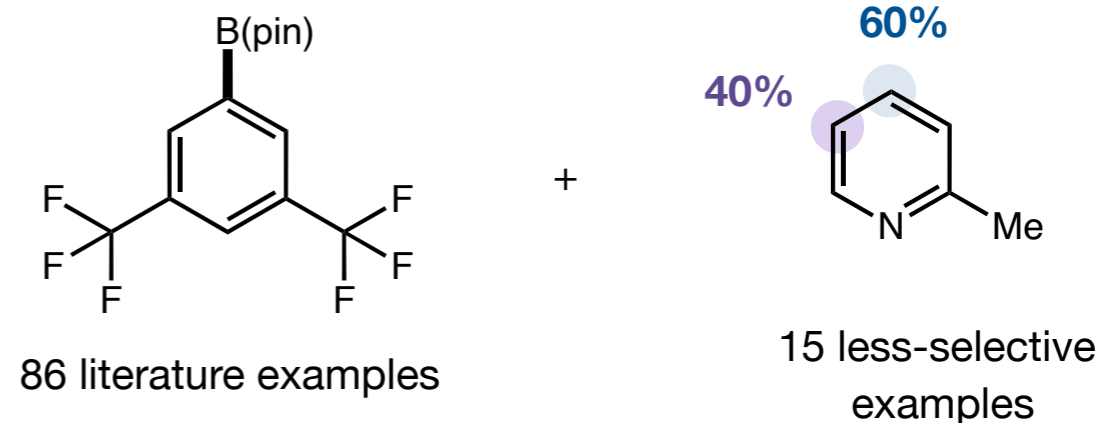
Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

Task



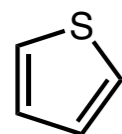
What C-H bond will be borylated?

Data Set

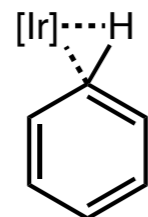
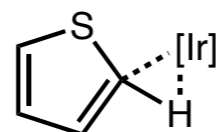
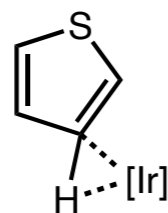


Data parameterization

Input SMILES



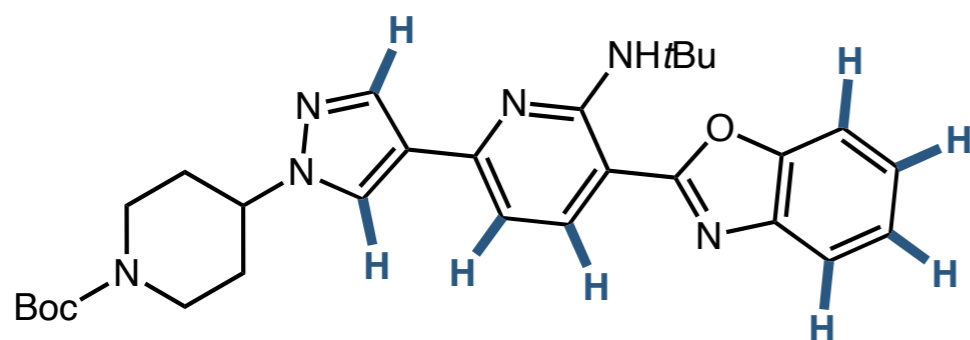
DFT TS for each C-H bond



Ph used as reference

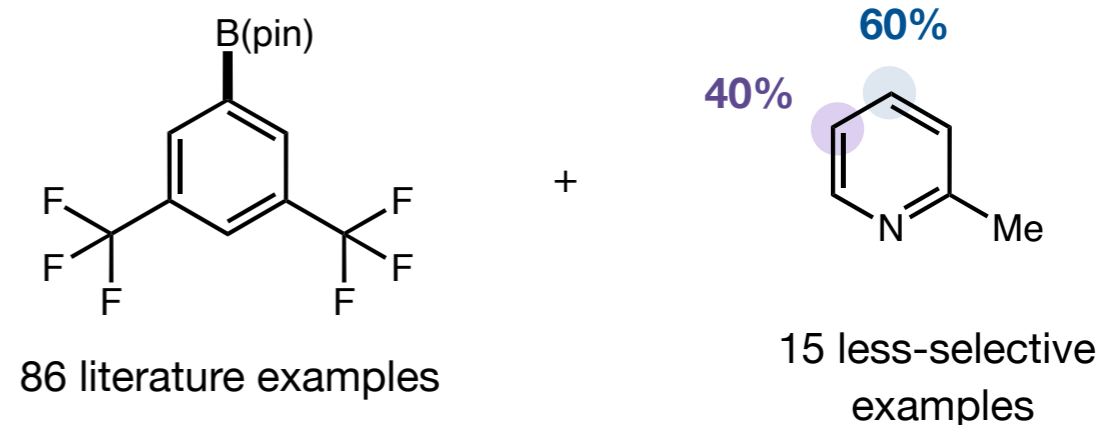
Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

Task



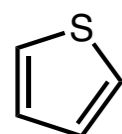
What C-H bond will be borylated?

Data Set

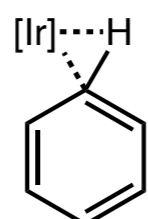
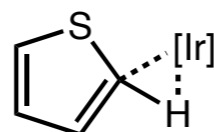
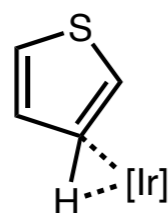


Data parameterization

Input SMILES



DFT TS for each C-H bond



Ph used as reference

Modeling

Partial Least Squares algorithm (electronic)

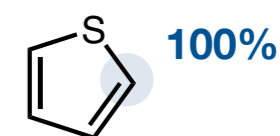
+

Sterimol Parameters (steric)

=

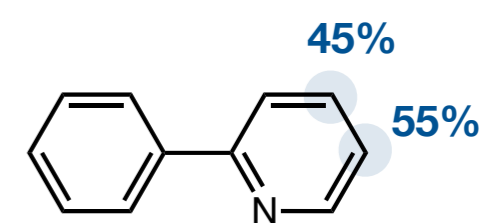
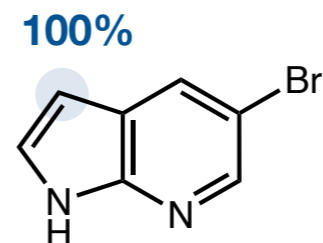
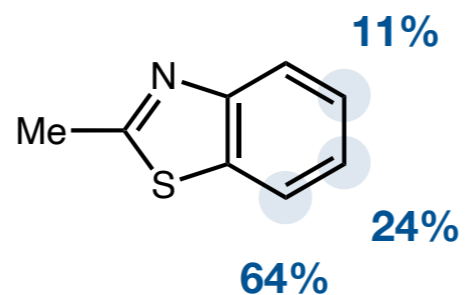
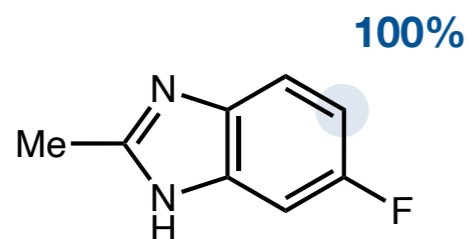
Predicted relative energy barrier

Boltzmann weights give predicted product distribution

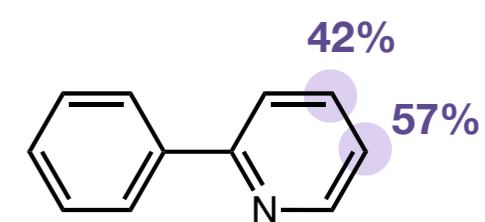
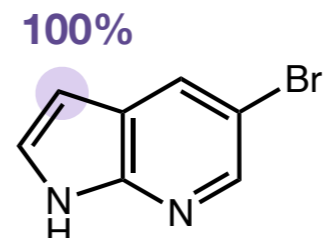
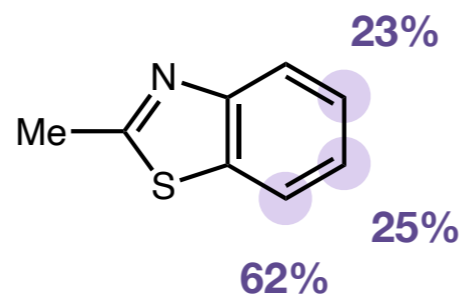
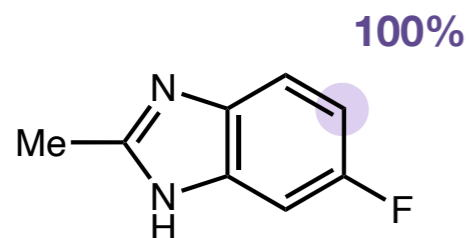


Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

Experimental Selectivities



Predicted Selectivities



Model is quite accurate in predicting selectivity for simple substrates

Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

**predict the major site
of borylation**

***15 expert
chemists***

VS

Model

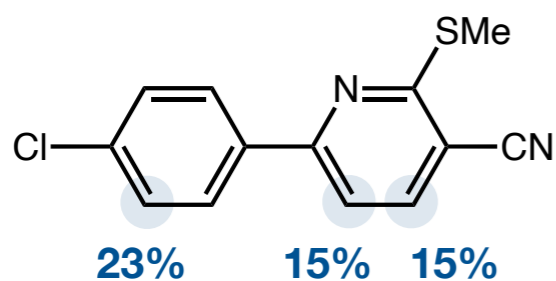
Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

predict the major site
of borylation

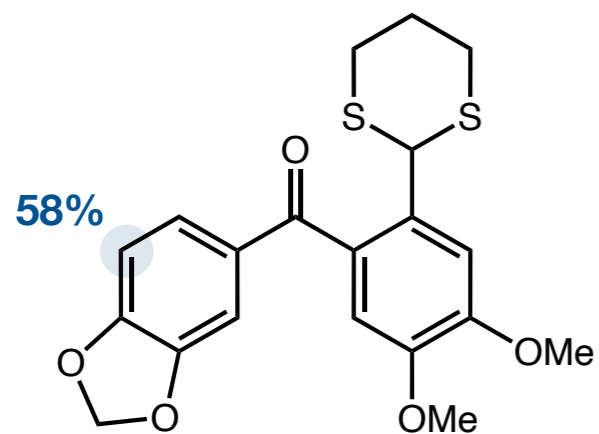
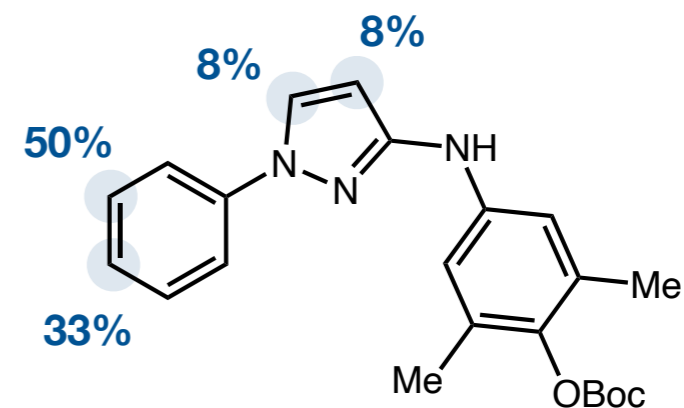
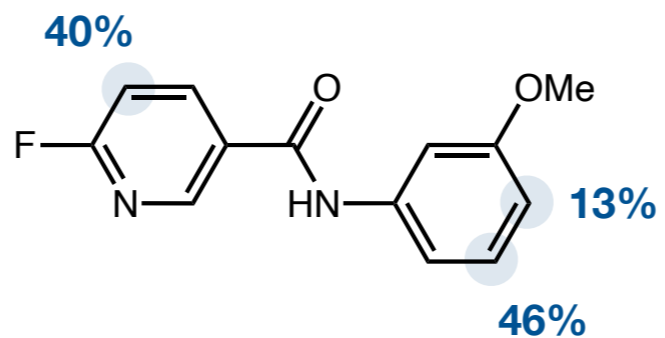
15 expert
chemists

VS

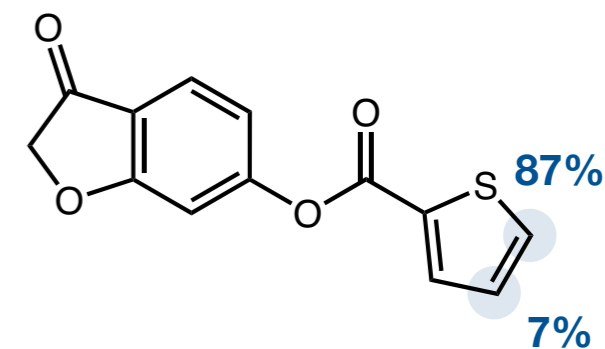
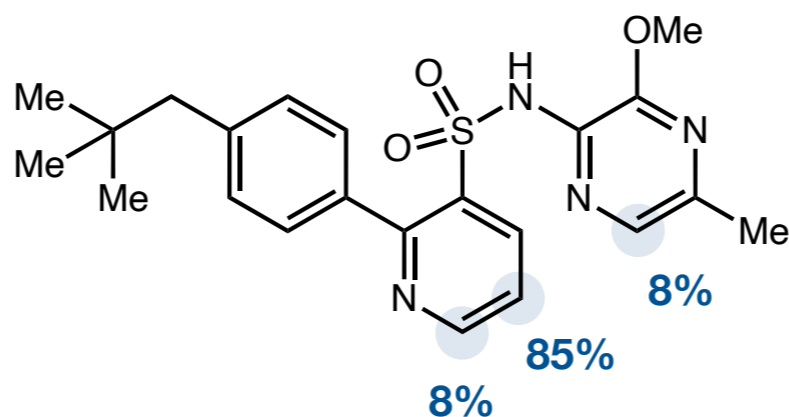
Model



No reaction: 46%



No reaction: 42%



No reaction: 7%

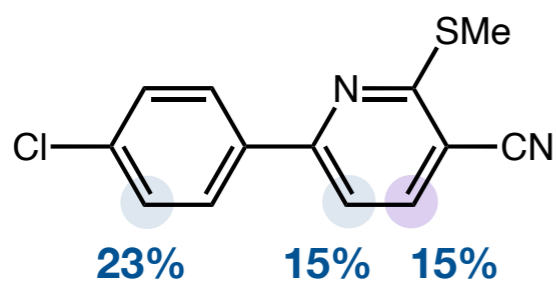
Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

predict the major site
of borylation

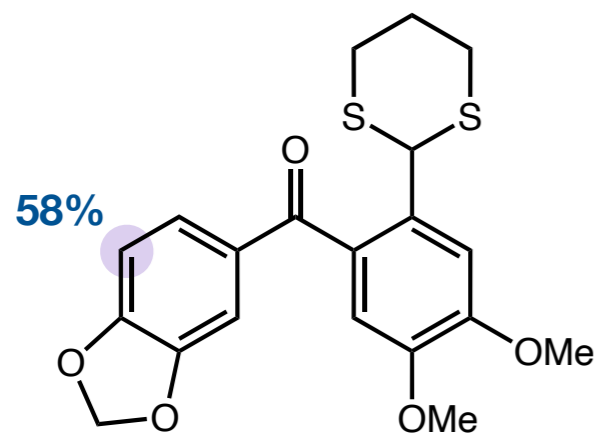
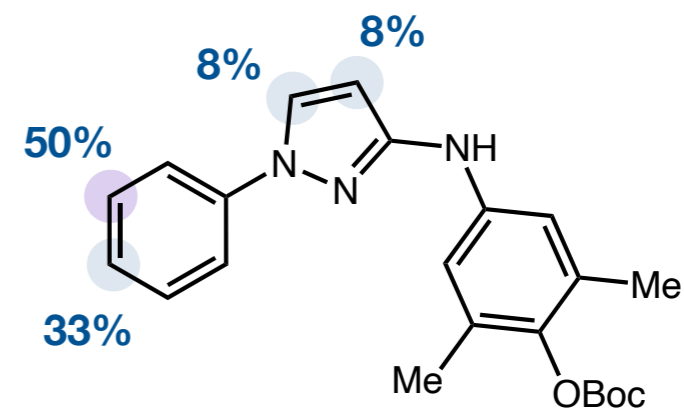
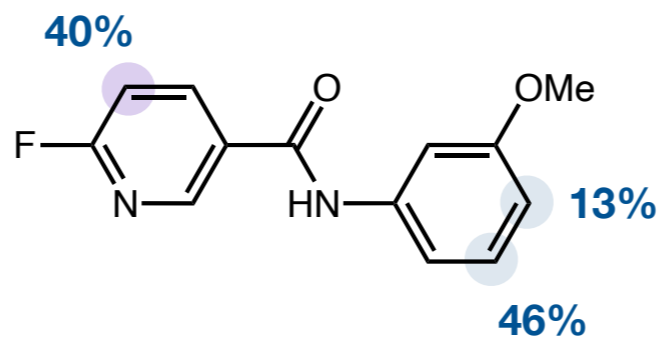
15 expert
chemists

VS

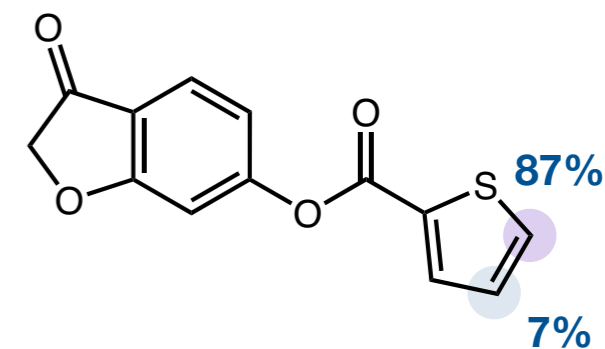
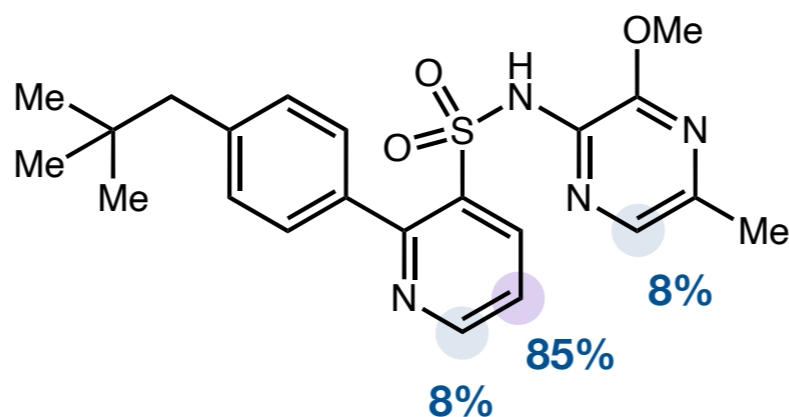
Model



No reaction: 46%



No reaction: 42%



No reaction: 7%

Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

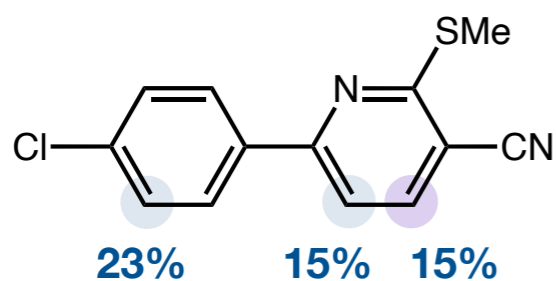
predict the major site
of borylation

15 expert
chemists

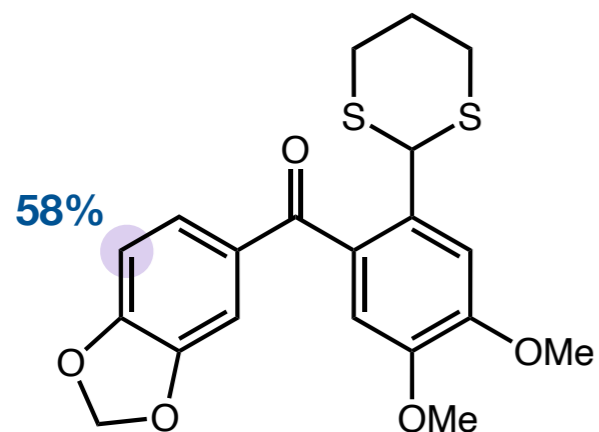
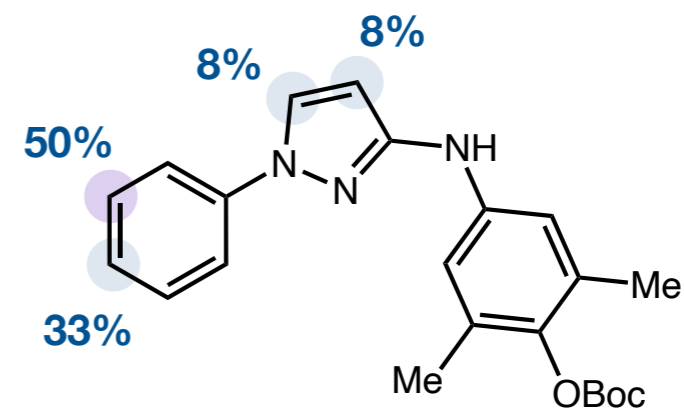
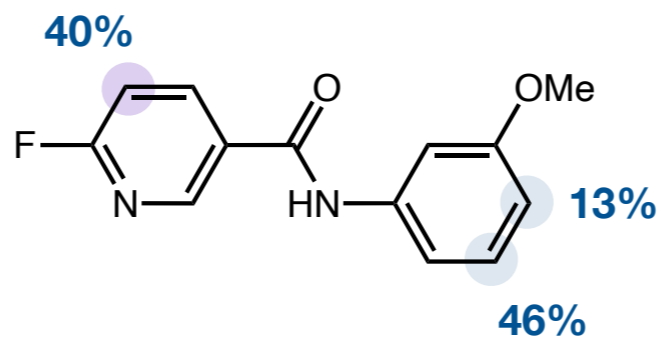
VS

Model

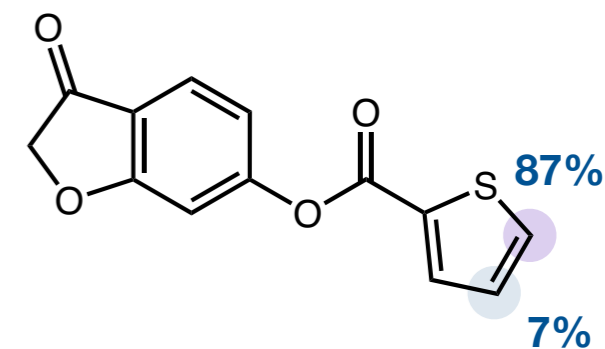
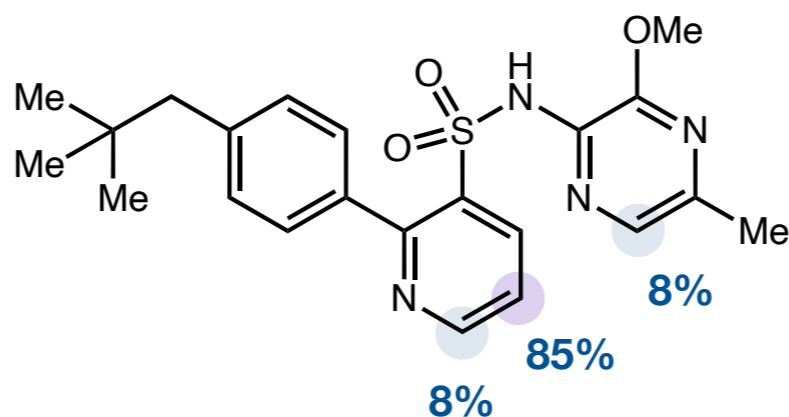
Correct answer in
purple



No reaction: 46%



No reaction: 42%



No reaction: 7%

Machine Learning Applied to Chemistry - Selectivity Predictions in Complex Settings

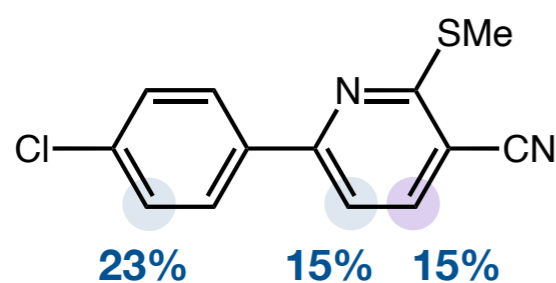
Overall: 56% accurate

15 expert chemists

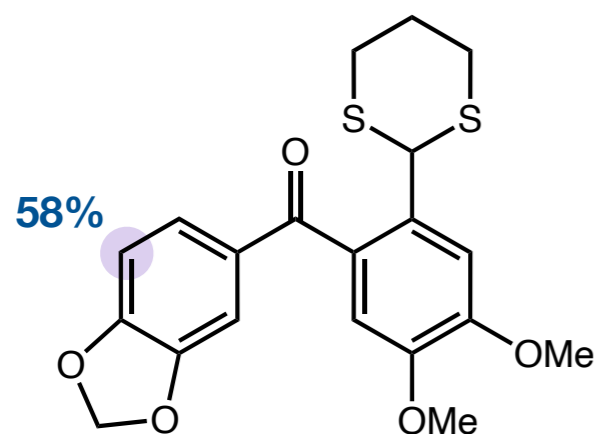
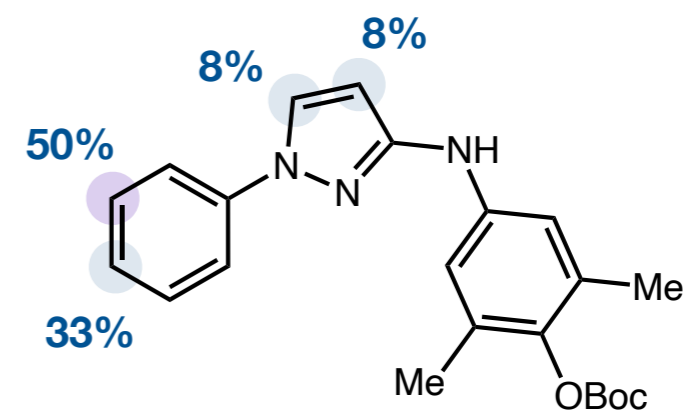
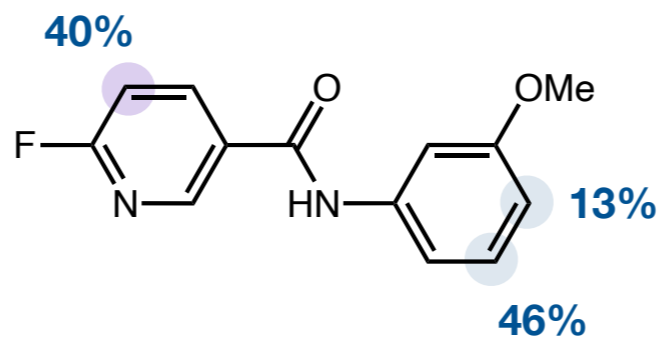
VS

Model

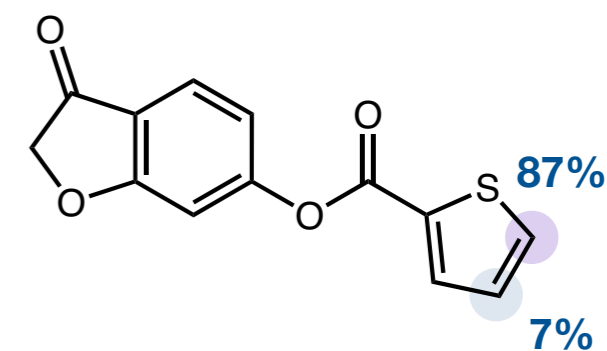
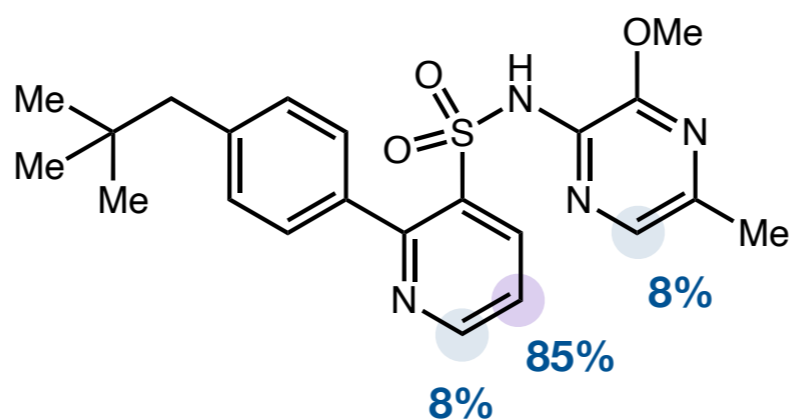
Overall: 100% accurate



No reaction: 46%



No reaction: 42%



No reaction: 7%

Machine Learning in Chemistry

Reviews

Zuranski, A. M.; Martinez Alvarado, J. I.; Shields, B. J.; Doyle, A. G. *Acc. Chem. Res.* **2021**, *54*, 1856–1865

Meuwly, M. *Chem. Rev.* **2021**, *121*, 10218–10239

Tu, Z.; Stuyver, T.; Coley, C. W. *Chem Sci*, **2022**, *14*, 226–244

Aal E Ali, R. S.; Meng, J.; Khan, M. E. I.; Jiang, X. *Artificial Intelligence Chemistry*, **2024**, *2*

Additional Examples

Romer, N. P.; Min, D. S.; Wang, J. Y.; Walroth, R. C.; Mack, K. A.; Sirois, L. E.; Gosselin, F.; Zell, D.; Doyle, A. G.; Sigman, M. S. *ACS Catal.* **2024**, *14*, 4699–4708

Wang, J. Y. [...] Doyle, A. G. *Nature*, **2024**, *626*, 1025–1033

Zuranski, A. M.; Gandhim S. S.; Doyle, A. G. *J. Am. Chem. Soc.* **2023**, *145*, 7898–7909

Zahrt, A. F.; Henle, J. J.; Rose, B. T.; Wang, Y.; Darrow, W. T.; Denmark, S. E. *Science*, **2019**, *247*

Rinehart, N. I.; Saunthwal, R. K.; Wellauer, J.; Zahrt, A. F.; Schlemper, L.; Shved, A. S.; Bigler, R.; Fantasia, S.; Denmark, S. E. **2023**, *381*, 965–972

Baczewska, P.; Kulczykowski, M.; Zambron, B.; Adamczak, J.; Pakulski, Z.; Roszak, R.; Grzybowski, B. A.; Młynarski, J. *Angew. Chem. Int. Ed.* **2024**, *136*

Gao, H.; Struble, T. J.; Coley, C. W.; Wang, Y.; Green, W. H.; Jensen, K. F. *Acs Cent. Sci.* **2018**, *4*, 1465–1476

Angello, N. H.; Rathore, V.; Beker, W.; Wolos, A.; Jira, E. R.; Roszak, R.; Wu, T. C.; Schroder, C. M.; Guzik, A. A.;

Grzybowski, B. A.; Burke, M. D. *Science*, **2022**, *378*, 399–405

Limitations

Schnitzer, T.; Schnurr, M.; Zahrt, A. F.; Sakhaee, N.; Denmark, S. E.; Wennemers, H. *Acs Cent. Sci.* **2024**, *10*, 367–373

Beker, W.; Roszak, R.; Wolos, A.; Angello, N. H.; Rathore, V.; Burke, M. D.; Grybowski, B. A. *J. Am. Chem. Soc.* **2022**, *144*, 4819–1827

The path to developing a useful method

1. Identify a problem

Accelerated serendipity

2. Initial hit

*Accelerated serendipity
Modern HTE*

3. Optimization

*Bayesian optimization
ML based virtual screening
Multidimensional LFER*

4. Generalize

*ML prediction of optimal catalyst
Additive screening strategies
General HTE
Data science substrate mapping*

5. General adoption

*ML prediction of selectivity
ML prediction of yield*

The path to developing a useful method

1. Identify a problem

Accelerated serendipity

2. Initial hit

Accelerated serendipity
Modern HTE

3. Optimization

Bayesian optimization
ML based virtual screening
Multidimensional LFER

Questions?

4. Generalize

ML prediction of optimal catalyst
Additive screening strategies
General HTE
Data science substrate mapping

5. General adoption

ML prediction of selectivity
ML prediction of yield