

The anatomy of a chemical reaction: Dissection by machine learning algorithms

Alex M. Clark, Ph.D.

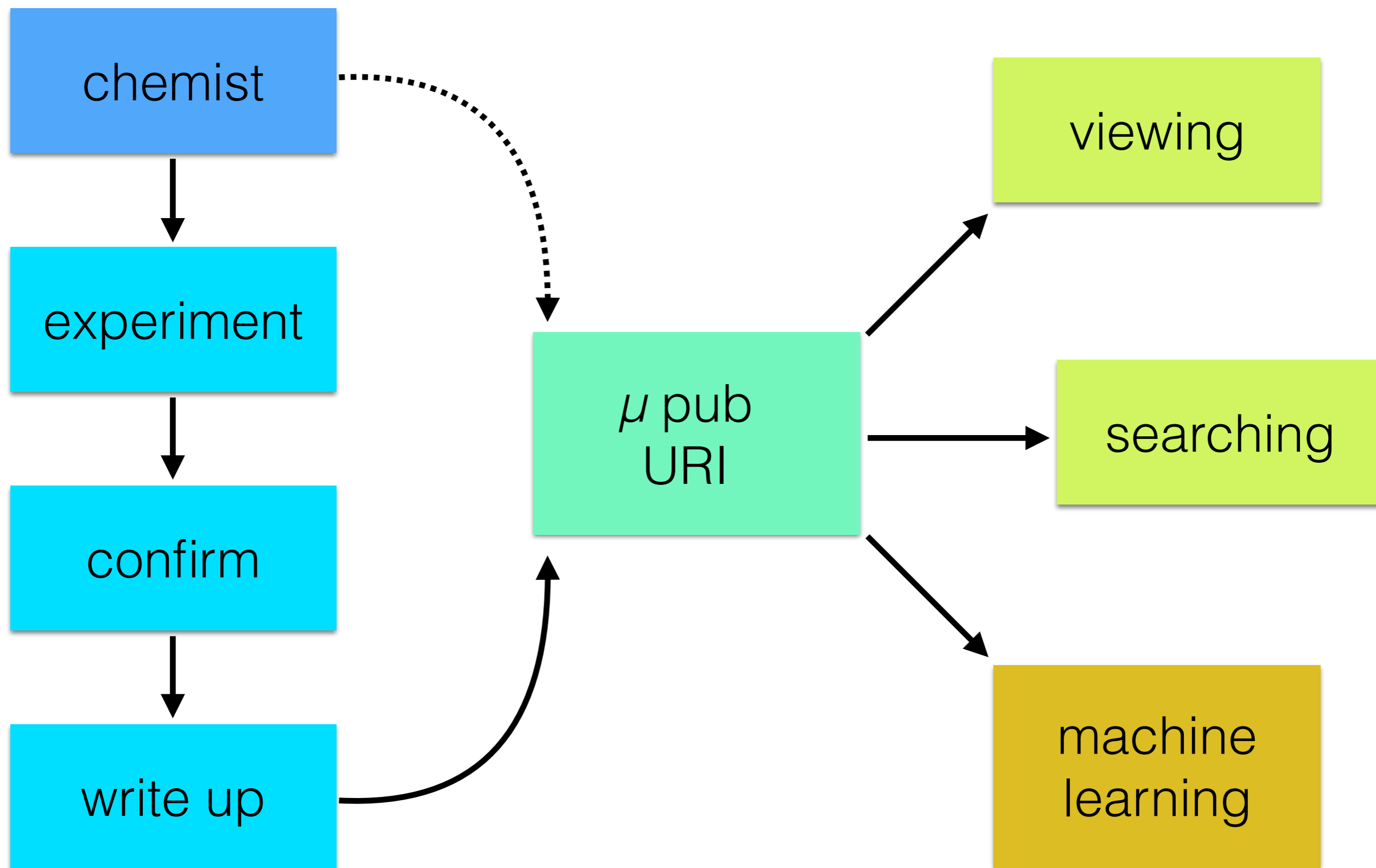
August 2014



© 2015 Molecular Materials Informatics, Inc.

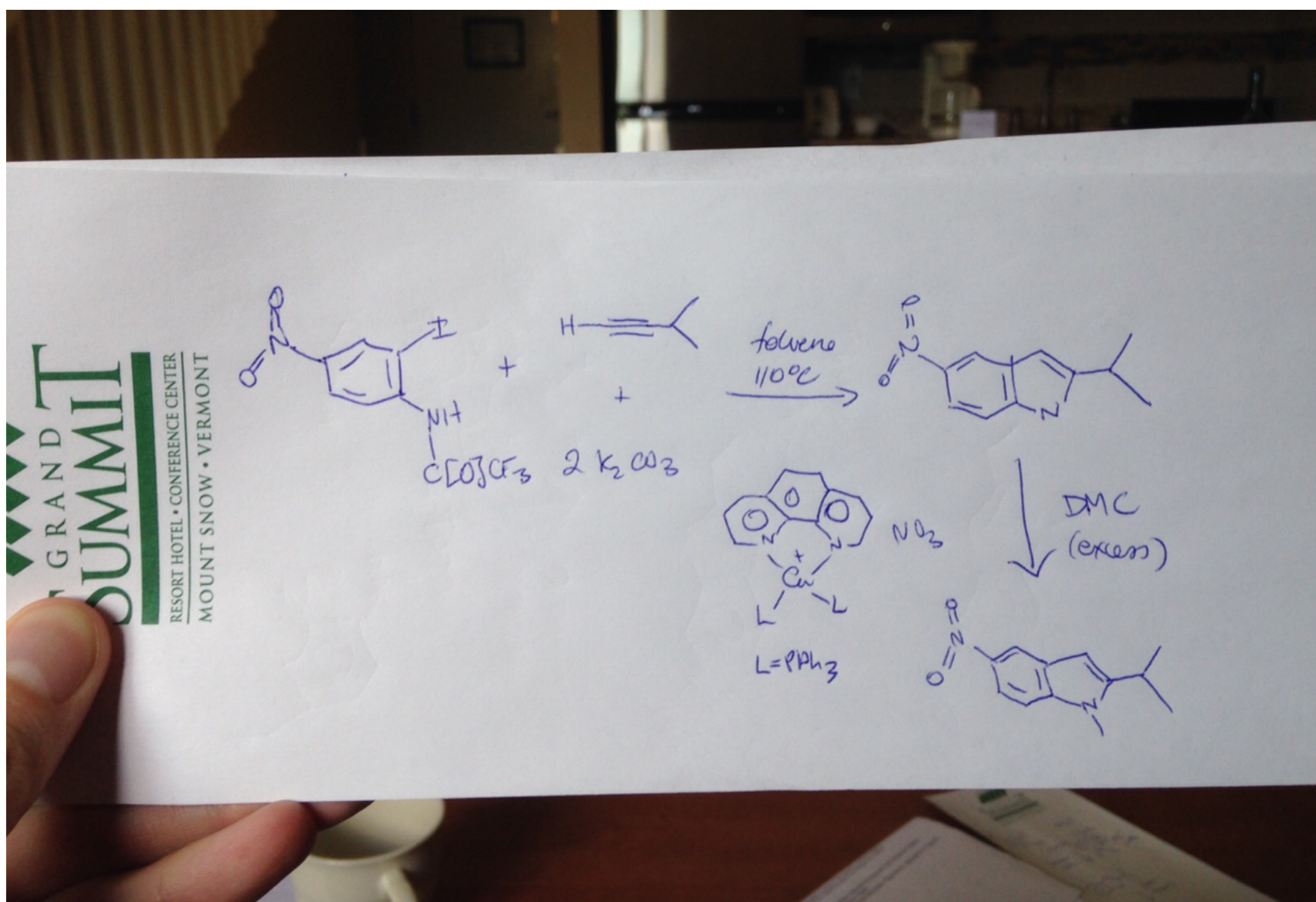
<http://molmatinf.com>

21st Century Publishing



All your byte are belong to us

- Just because a reaction scheme is digital...



- ... doesn't mean it's of any use to a computer.

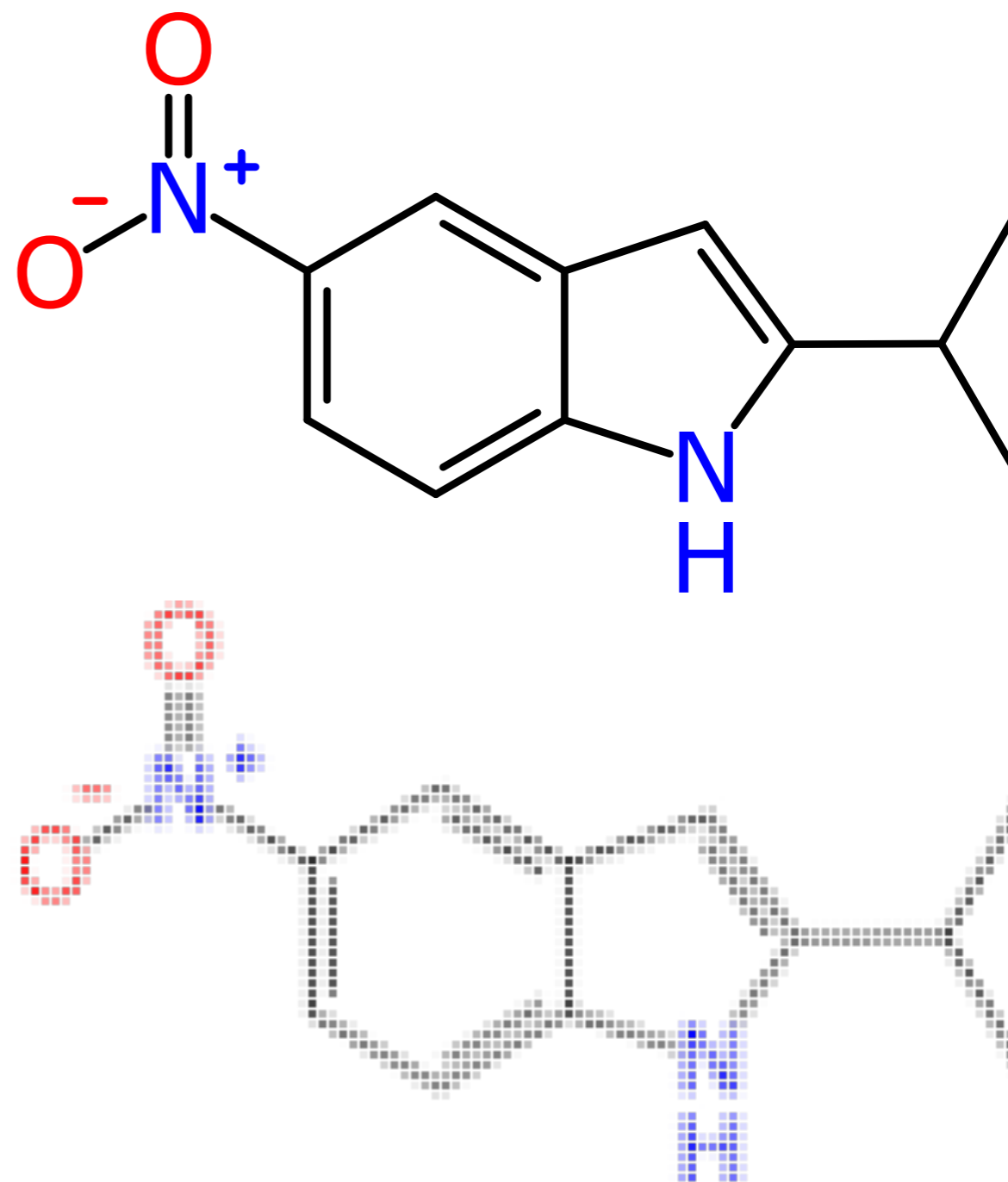
Production Raster Graphics

Generic molfile

```

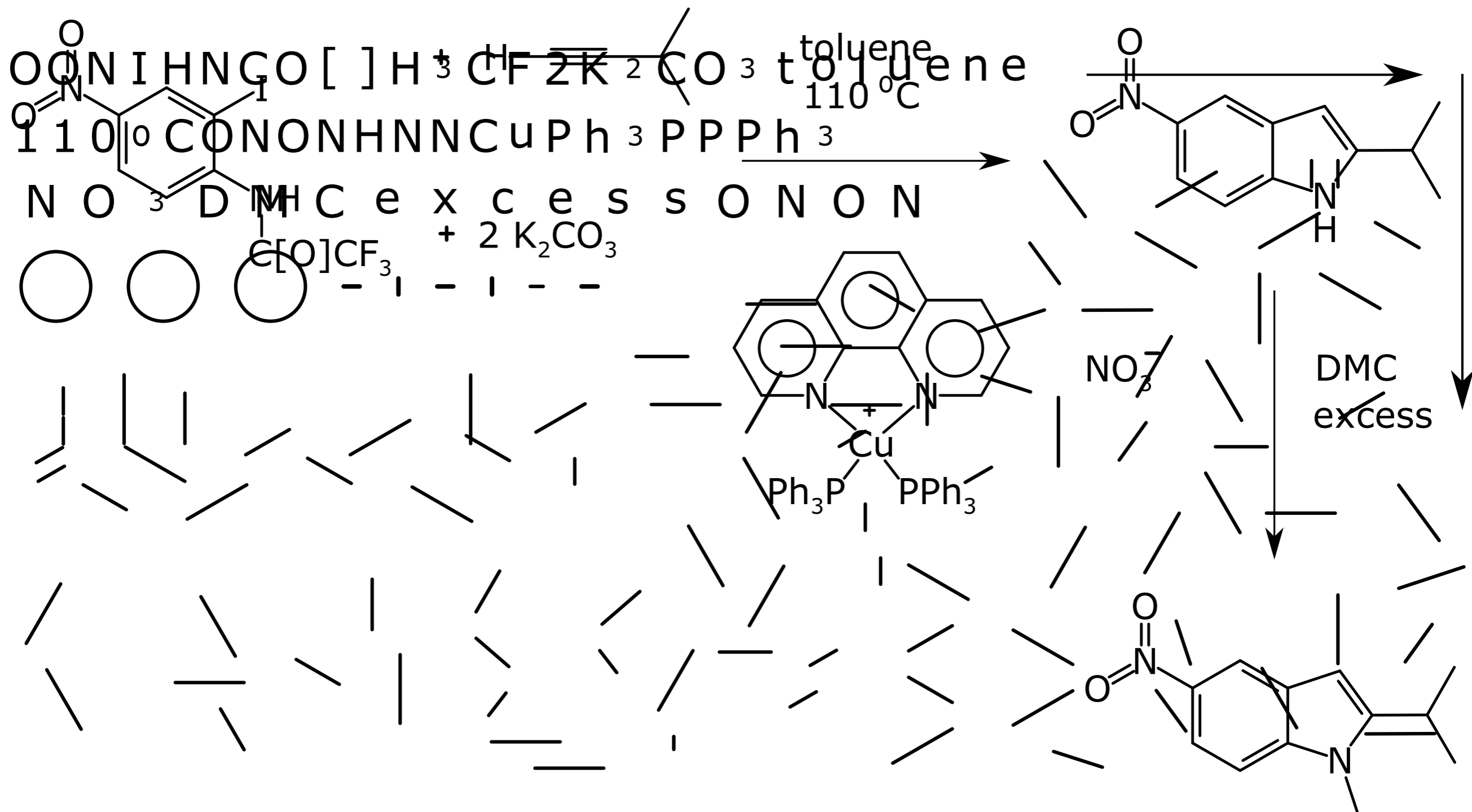
15 16 0 0 0 0 0 0 0 0999 v2000
-3.9510 4.0500 0.0000 C 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
-5.2500 3.3000 0.0000 C 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
-2.6519 3.3000 0.0000 C 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
-5.2500 1.8000 0.0000 C 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
-3.9510 1.0500 0.0000 C 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
-2.6519 1.8000 0.0000 C 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
-1.2306 3.7694 0.0000 C 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
-0.3482 2.5611 0.0000 C 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
-1.2240 1.3407 0.0000 N 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
-6.5490 4.0500 0.0000 N 0 3 0 0 0 0 0 0 0 0 0 0 0 0 0
-7.8481 3.3000 0.0000 O 0 5 0 0 0 0 0 0 0 0 0 0 0 0 0
-6.5490 5.5500 0.0000 O 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 1.1518 2.5673 0.0000 C 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 1.9072 1.2714 0.0000 C 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 1.8964 3.8695 0.0000 C 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
 1 2 1 0 0 0 0
 1 3 2 0 0 0 0
 2 4 2 0 0 0 0
 4 5 1 0 0 0 0
 5 6 2 0 0 0 0
 6 3 1 0 0 0 0
 3 7 1 0 0 0 0
 7 8 2 0 0 0 0
 8 9 1 0 0 0 0
 9 6 1 0 0 0 0
 2 10 1 0 0 0 0
10 11 1 0 0 0 0
10 12 2 0 0 0 0
 8 13 1 0 0 0 0
13 14 1 0 0 0 0
13 15 1 0 0 0 0
M CHG 2 10 1 11 -1
M END

```



Production Vector Graphics

- Manuscripts usually delivered as PDFs:



Spreadsheets

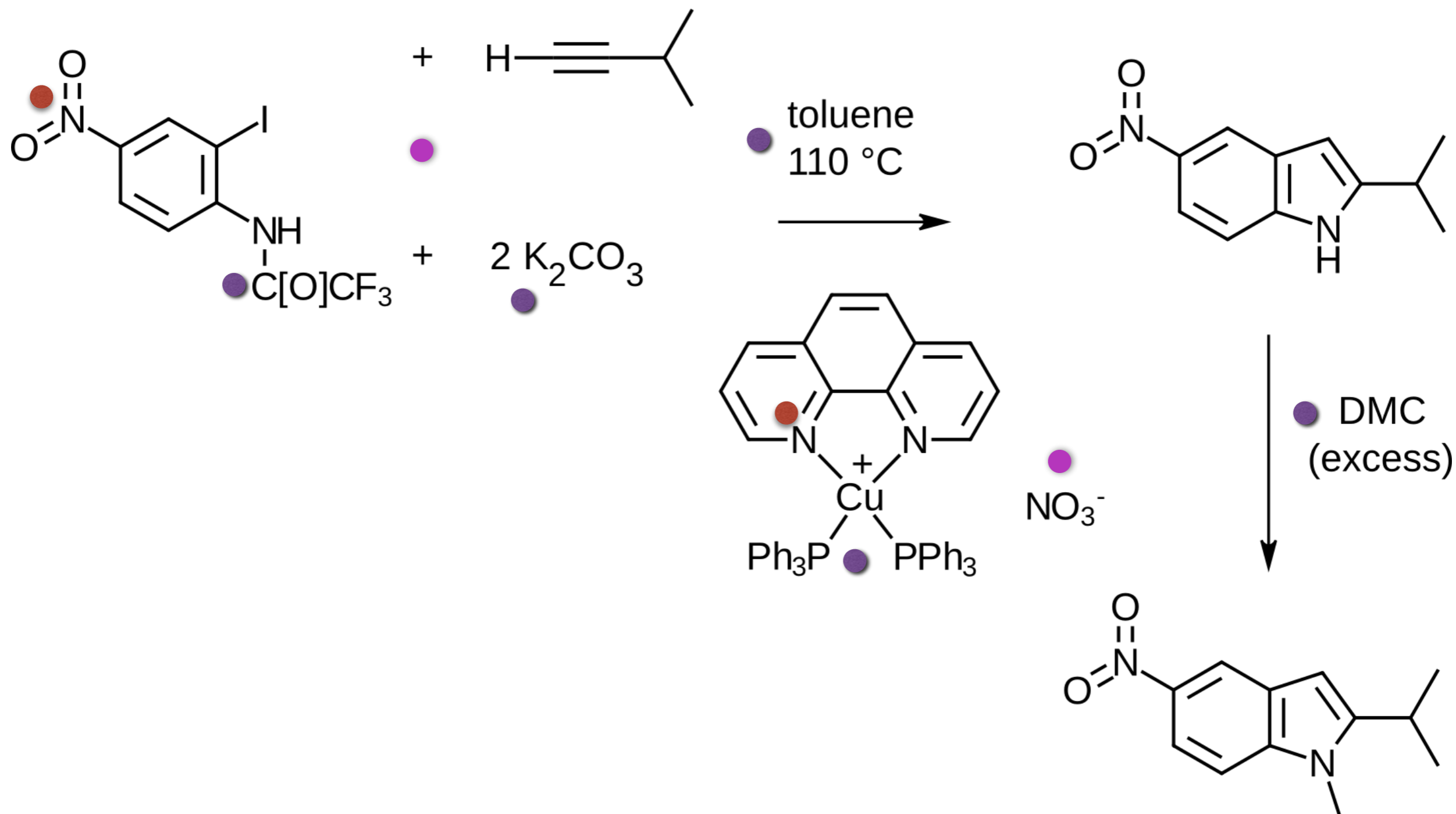
The screenshot shows a Microsoft Excel spreadsheet titled "ssreaction.xlsx - Microsoft Excel non-commercial use". The spreadsheet is organized into columns for reaction data. The active cell is B2, containing the date "10-08-2015".

	<u>Reactants</u>	<u>Reagents</u>	<u>Solvents</u>	<u>Catalysts</u>	<u>Product</u>
Step 1	<chem>[O-][N+](=O)c1cc(l)c(cc1)NC(=O)C(F)(F)F</chem>	<chem>[K+].[K+].[O-]C([O-])=O</chem>	<chem>Cc1ccccc1</chem>	<chem>[O-][N+](O)c1ccc2ccc3ccc[n+]4c3c2[n+]1[Cu-3]4([P+](c1ccccc1)(c1ccccc1)c1ccccc1)[P+](c1ccccc1)(c1ccccc1)c1ccccc1</chem>	<chem>CC(C)c1cc2cc(ccc2[nH]1)[N+](O)=O</chem>
	<chem>CC(C)C#C</chem>				
Step 2		<chem>COC(=O)OC</chem>			<chem>CC(C)c1cc2cc(ccc2[n]1C)[N+](O)=O</chem>

- Data gives the impression of organisation
- Very high degrees of freedom, nothing for structures

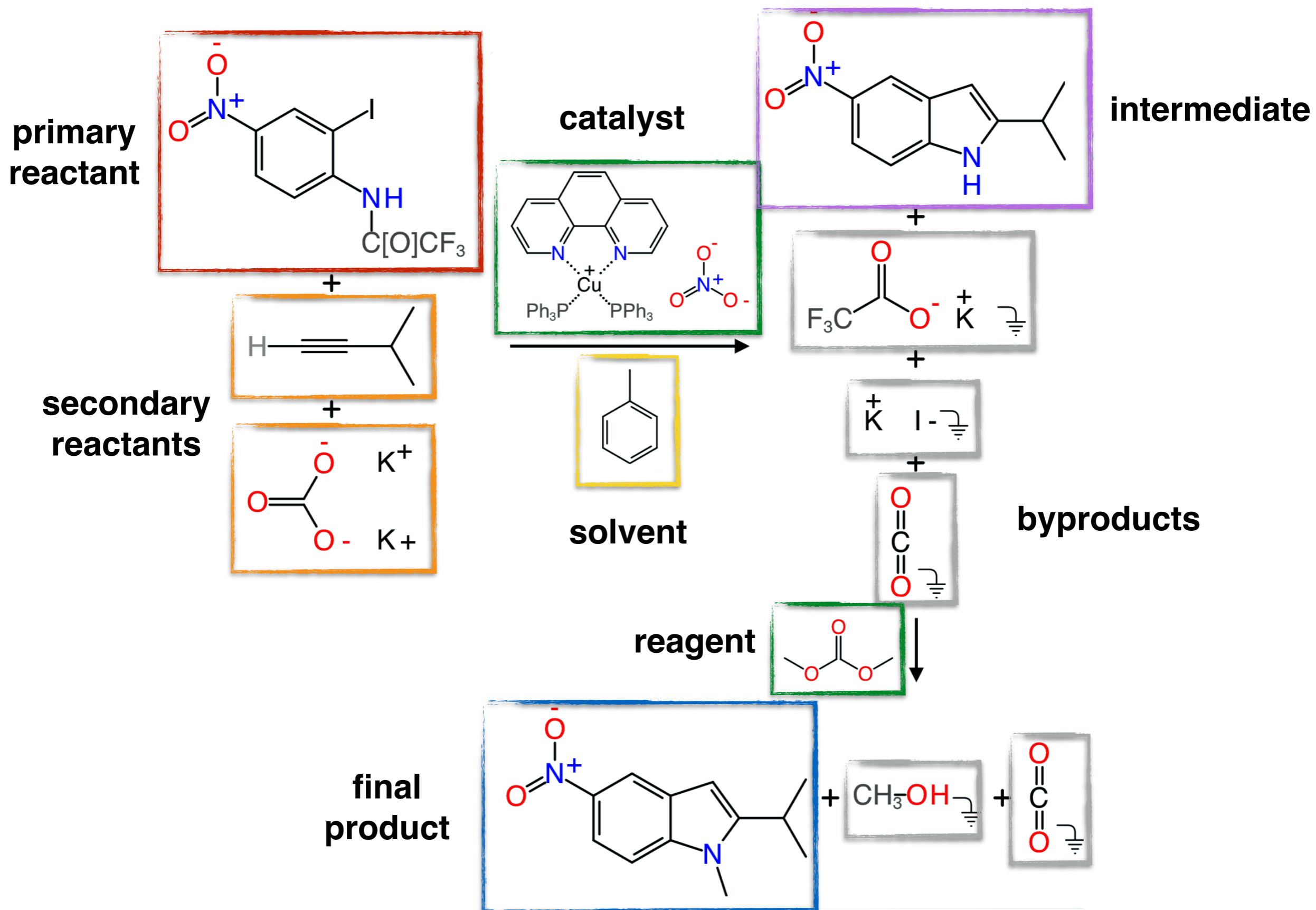


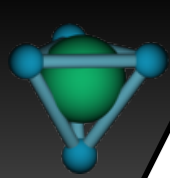
Common Scheme





Digitally Friendly





Representation

Step	Stoich.	Structure	Role
1	1		Reactant
1	1		Reactant
1	1		Reagent
1	1		Reagent
1	1		Reagent
1	1		Product
2	1		Reagent
2	1		Product

- For **machines**: representation must be very rigidly defined
- For **humans**: can generate diagram programmatically
- MDL RXN/RDfile ~50% there
- DataSheet XML with Experiment aspect

<http://molmatinf.com/fmtaspect.html>



Balancing

Balanced
Aug 11, 2015

Step 1


$$C_{14}H_{12}F_3IK_2N_2O_6 \rightarrow C_{11}H_{12}N_2O_2 + C_3F_3IK_2O_4$$


Step 2

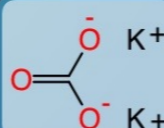
$$C_{14}H_{18}N_2O_5 \rightarrow C_{12}H_{14}N_2O_2 + C_2H_4O_3$$

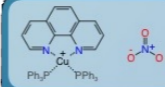


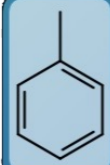
Quantities

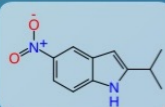
	Equiv:	1
	MW:	360.029 g/mol
	Mass:	1 g
	Volume:	
	Moles:	0.00277756 mol
	Density:	
	Conc:	

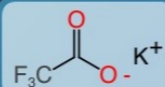
	Equiv:	1
	MW:	68.117 g/mol
	Mass:	0.189199 g
	Volume:	
	Moles:	0.00277756 mol
	Density:	
	Conc:	

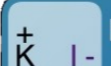
	Equiv:	1
	MW:	138.206 g/mol
	Mass:	0.383874 g
	Volume:	
	Moles:	0.00277756 mol
	Density:	
	Conc:	

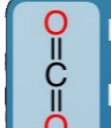
	Equiv:	
	MW:	830.327 g/mol
	Mass:	
	Volume:	
	Moles:	
	Density:	
	Conc:	

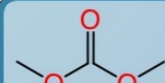
	Equiv:	
	MW:	92.1384 g/mol
	Mass:	
	Volume:	
	Moles:	
	Density:	0.865 g/mL
	Conc:	

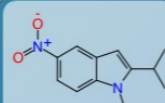
	Equiv:	1
	MW:	204.225 g/mol
	Mass:	0.567247 g
	Volume:	
	Moles:	0.00277756 mol
	Density:	
	Conc:	
	Yield:	100 %


	Equiv:	1
	MW:	152.114 g/mol
	Mass:	0.422504 g
	Volume:	
	Moles:	0.00277756 mol
	Density:	
	Conc:	
	Yield:	

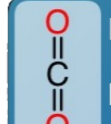
	Equiv:	1
	MW:	166.003 g/mol
	Mass:	0.461082 g
	Volume:	
	Moles:	0.00277756 mol
	Density:	
	Conc:	
	Yield:	

	Equiv:	1
	MW:	44.0095 g/mol
	Mass:	0.122239 g
	Volume:	
	Moles:	0.00277756 mol
	Density:	
	Conc:	
	Yield:	

	Equiv:	
	MW:	90.0779 g/mol
	Mass:	
	Volume:	
	Moles:	
	Density:	
	Conc:	

	Equiv:	1
	MW:	218.252 g/mol
	Mass:	0.606206 g
	Volume:	
	Moles:	0.00277756 mol
	Density:	
	Conc:	
	Yield:	100 %

	Equiv:	1
	MW:	32.0419 g/mol
	Mass:	0.0889981 g
	Volume:	
	Moles:	0.00277756 mol
	Density:	
	Conc:	
	Yield:	

	Equiv:	1
	MW:	44.0095 g/mol
	Mass:	0.122239 g
	Volume:	
	Moles:	0.00277756 mol
	Density:	
	Conc:	
	Yield:	

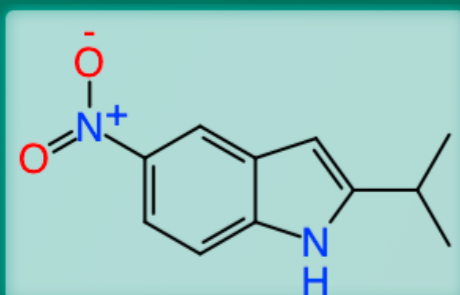


Green Metrics

$$\Sigma \text{ reactants} = 1 \text{ g} + 0.189199 \text{ g} + 0.829233 \text{ g} + 0.115314 \text{ g} + 17.34 \text{ g} + 0.5 \text{ g} = 19.9737 \text{ g}$$

$$\Sigma \text{ products} = 0.567247 \text{ g} + 0.606206 \text{ g} = 1.17345 \text{ g}$$

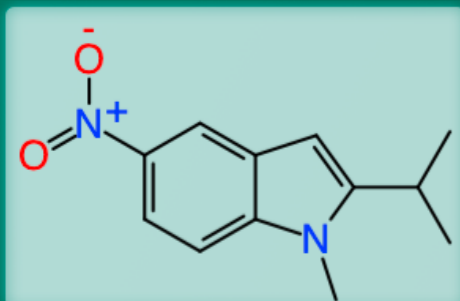
$$\Sigma \text{ waste} = 0.422504 \text{ g} + 0.461082 \text{ g} + 0.122239 \text{ g} + 0.0889981 \text{ g} + 0.122239 \text{ g} = 1.21706 \text{ g}$$



$$\text{PMI} = \frac{1 \text{ g} + 0.189199 \text{ g} + 0.829233 \text{ g} + 0.115314 \text{ g} + 17.34 \text{ g}}{0.567247 \text{ g}} = 34.3303$$

$$\text{E-factor} = \frac{0.422504 \text{ g} + 0.461082 \text{ g} + 0.122239 \text{ g}}{0.567247 \text{ g}} = 1.77317$$

$$\text{Atom-E} = \frac{204.225}{360.029 + 68.117 + 138.206} = 36.0598 \%$$



$$\text{PMI} = \frac{1 \text{ g} + 0.189199 \text{ g} + 0.829233 \text{ g} + 0.115314 \text{ g} + 17.34 \text{ g} + 0.5 \text{ g}}{0.606206 \text{ g}} = 32.9487$$

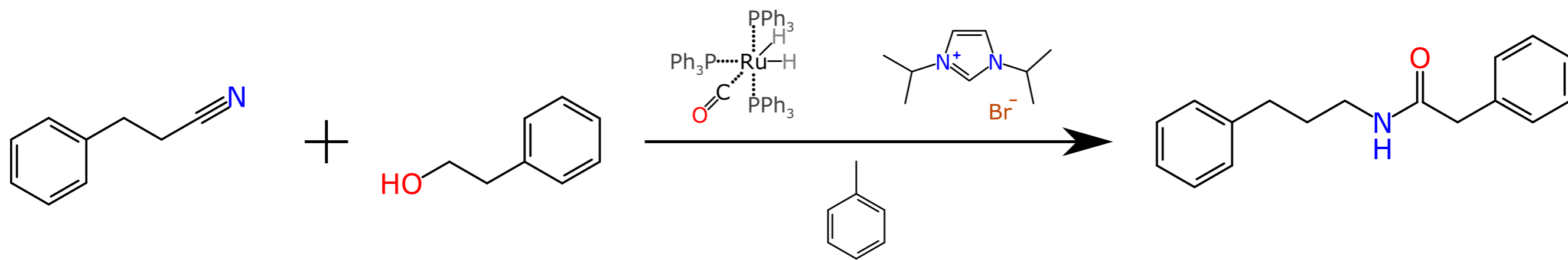
$$\text{E-factor} = \frac{0.422504 \text{ g} + 0.461082 \text{ g} + 0.122239 \text{ g} + 0.0889981 \text{ g} + 0.122239 \text{ g}}{0.606206 \text{ g}} = 2.00767$$

$$\text{Atom-E} = \frac{218.252}{204.225 + 90.0779} = 74.1588 \%$$

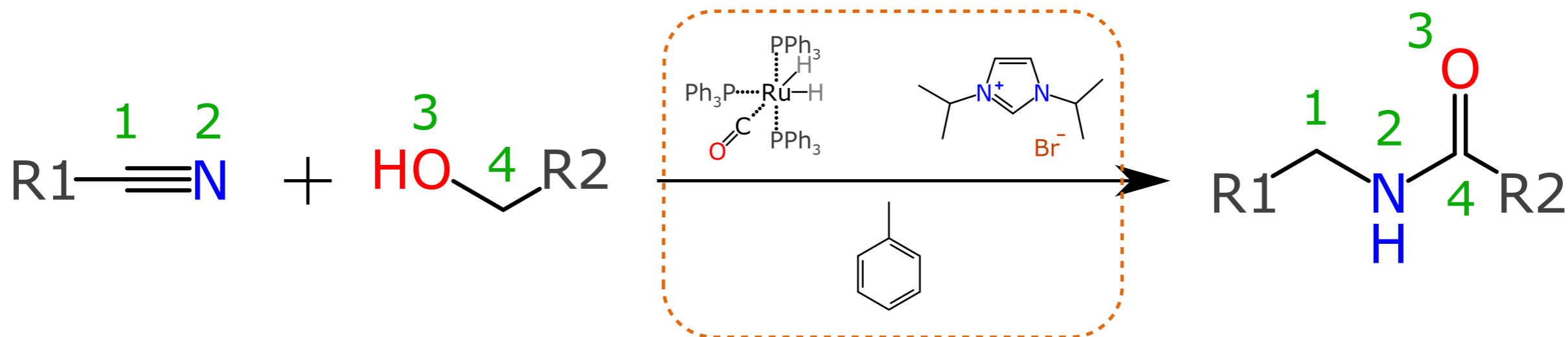
- Totals for reactants, products & waste
- For each non-waste product: **yield**, **PMI**, **E-factor**, **Atom-E**... always calculated, always recorded

Reaction Transforms

- **Reaction** = specific description of experiment



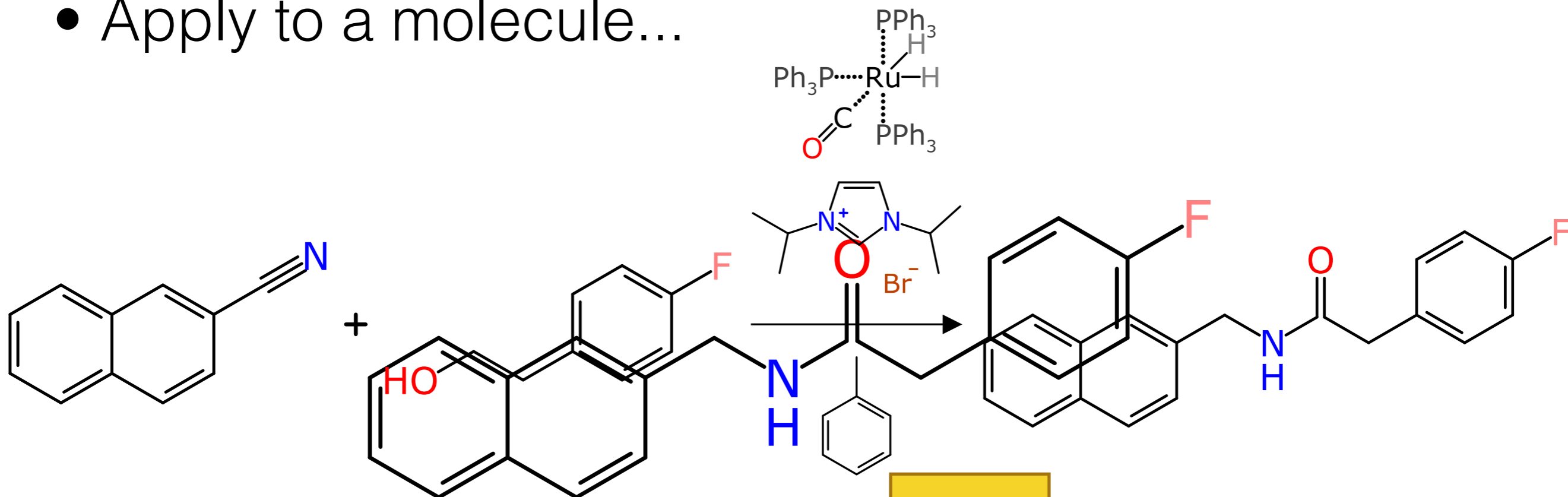
- **Transform** = the generic form of a reaction





Convenience

- Apply to a molecule...



	Equiv:	1
	MW:	153.18 g/mol
	Mass:	5.22202 g
	Volume:	
	Moles:	0.0340907 mol
	Density:	
	Conc:	

	Equiv:	1
	MW:	140.155 g/mol
	Mass:	4.77798 g
	Volume:	
	Moles:	0.0340907 mol
	Density:	
	Conc:	

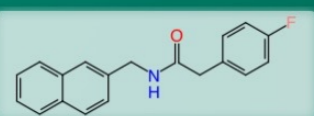
	Equiv:	0.1
	MW:	917.952 g/mol
	Mass:	3.12936 g
	Volume:	
	Moles:	0.00340907 mol
	Density:	
	Conc:	

	Equiv:	0.1
	MW:	233.149 g/mol
	Mass:	0.794821 g
	Volume:	
	Moles:	0.00340907 mol
	Density:	
	Conc:	

	Equiv:	137.692
	MW:	92.1384 g/mol
	Mass:	432.5 g
	Volume:	500 mL
	Moles:	4.69402 mol
	Density:	0.865 g/mL
	Conc:	

	Equiv:	1
	MW:	293.335 g/mol
	Mass:	10 g
	Volume:	
	Moles:	0.0340907 mol
	Density:	
	Conc:	
	Yield:	100 %

Σ reactants = 5.22202 g + 4.77798 g + 3.12936 g + 0.794821 g + 432.5 g = 446.424 g Σ products = 10 g = 10 g Σ waste = ? = 0



$$\text{PMI} = \frac{5.22202 \text{ g} + 4.77798 \text{ g} + 3.12936 \text{ g} + 0.794821 \text{ g} + 432.5 \text{ g}}{10 \text{ g}} = 44.6424$$

$$\text{E-factor} = \frac{436.424 \text{ g}}{10 \text{ g}} = 43.6424$$

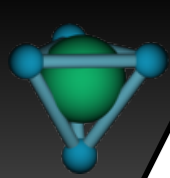
$$\text{Atom-E} = \frac{293.335}{153.18 + 140.155} = 100 \%$$



Decision Making

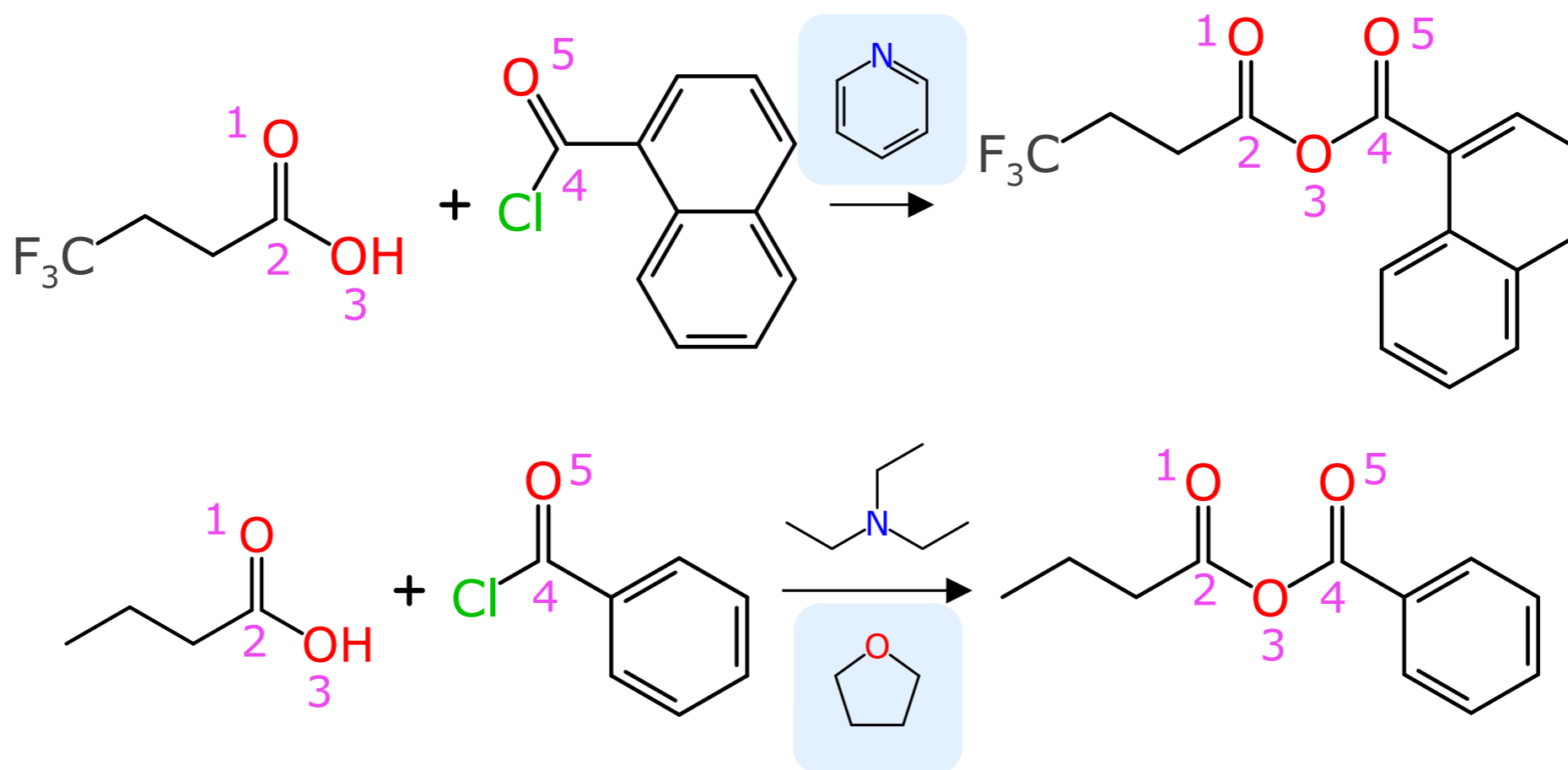
Product Search Results

	Yield	PMI	E-factor	Atom Economy
	100%	2.18	1.18	100%
	84%	12.49	11.49	93.3%
	82%	19.17	18.17	87.4%
	100%	8.26	7.26	56.3%
	63%	8.93	7.93	73.3%

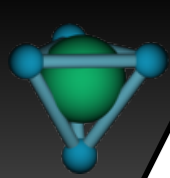


Model Building

- Most reaction data is noisy and incomplete
- Imagine opportunities with quantity & quality...



- For example: model solvent substitution



Conclusions & Future

- Most published reactions intractable to machines
- Most reaction informatics formats 50% complete
- Full description has immediate benefits...
- ... eventual large scale machine learning.
- μ Publications with provenance: the path to open repositories - **but** requires attention to content

Acknowledgments

- Antony Williams
- Sean Ekins
- Leah McEwen
- Open data advocates
- Inquiries to info@molmatinf.com

**MOLECULAR
MATERIALS
INFORMATICS**

<http://molmatinf.com>

<http://molsync.com>

<http://cheminf20.org>

@aclarkxyz

