

Molecular Structure Elucidation: More **accurate** or more **diverse** models?

Guido Falk von Rudorff, University of Kassel



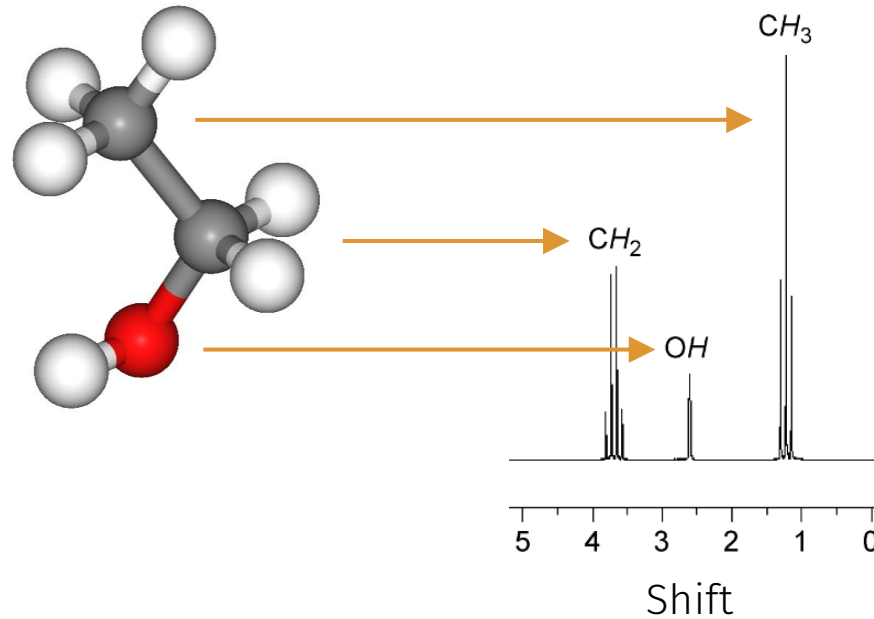
ferchault



nablachem.org/talks

What we have:

^1H -NMR, ^{13}C -NMR, UV-VIS, stoichiometry, ...

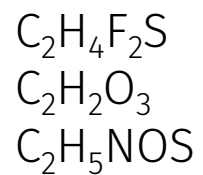


What we want:

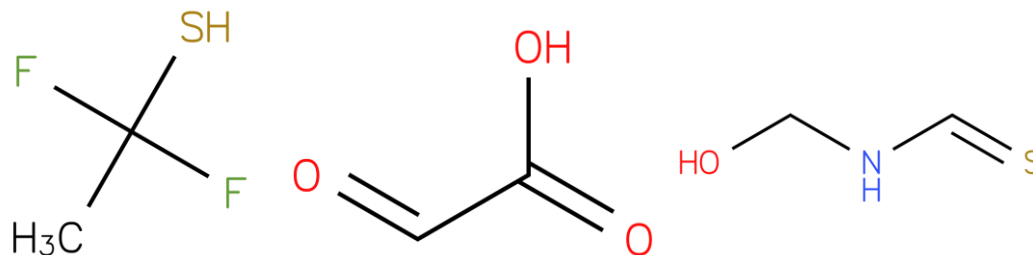
(Ranked) list of candidate geometries

Take some atoms from {C, O, N, F, S}, H-saturated

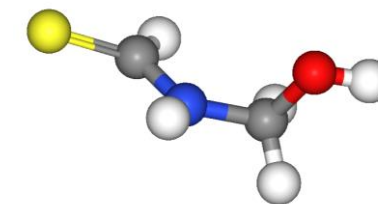
349:



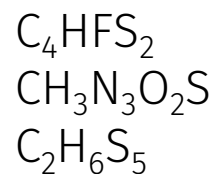
9,917:



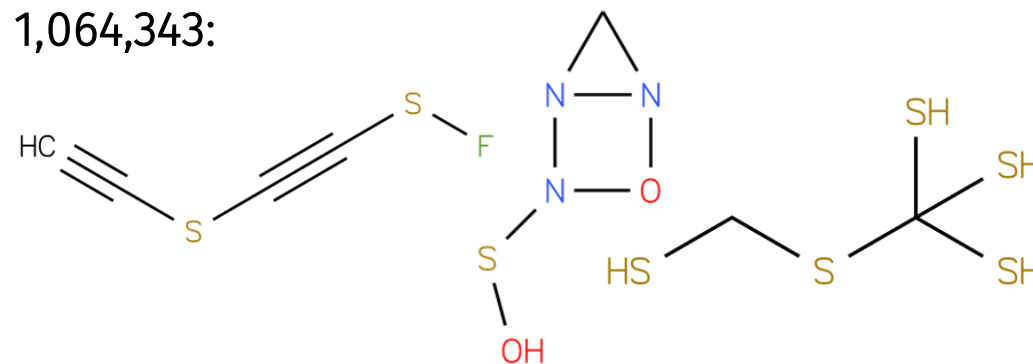
~52k:



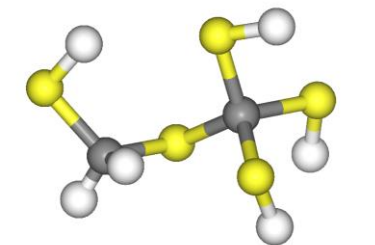
1,050:



1,064,343:



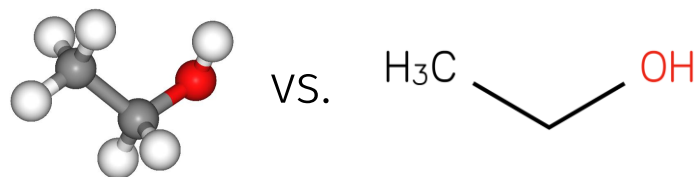
~23M:



Ibuprofen: ~120B

Cost: Too many candidates for forward estimates^[1-4]

Chemistry: Models often lack 3D information



hindrance and strain

Accuracy: Geometries nearly indistinguishable
Experimentally: ~0.2 ppm^[5]



Estimation by adding noise to computed NMR shifts

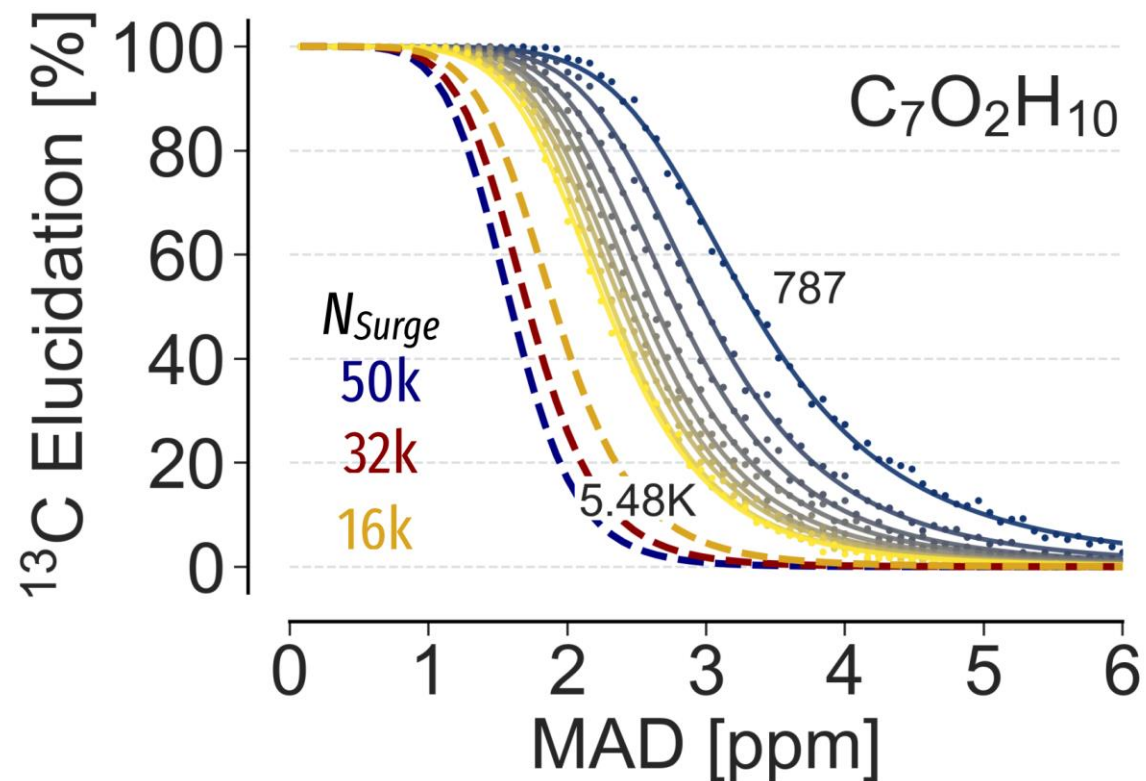
[1] E Jonas, *NeurIPS*, 2019. [2] Z Huang et al, *Chem Sci*, 2021. [3] B Sridharan et al, *J Phys Chem Lett*, 2022.

[4] A. Howarth et al, *Chem Sci*, 2022. [5] F. Fardus-Reid et al, *Anal. Methods*, 2016.

How accurate do we need to be?

Database: ~5.5k spectra from QM9-NMR + noise

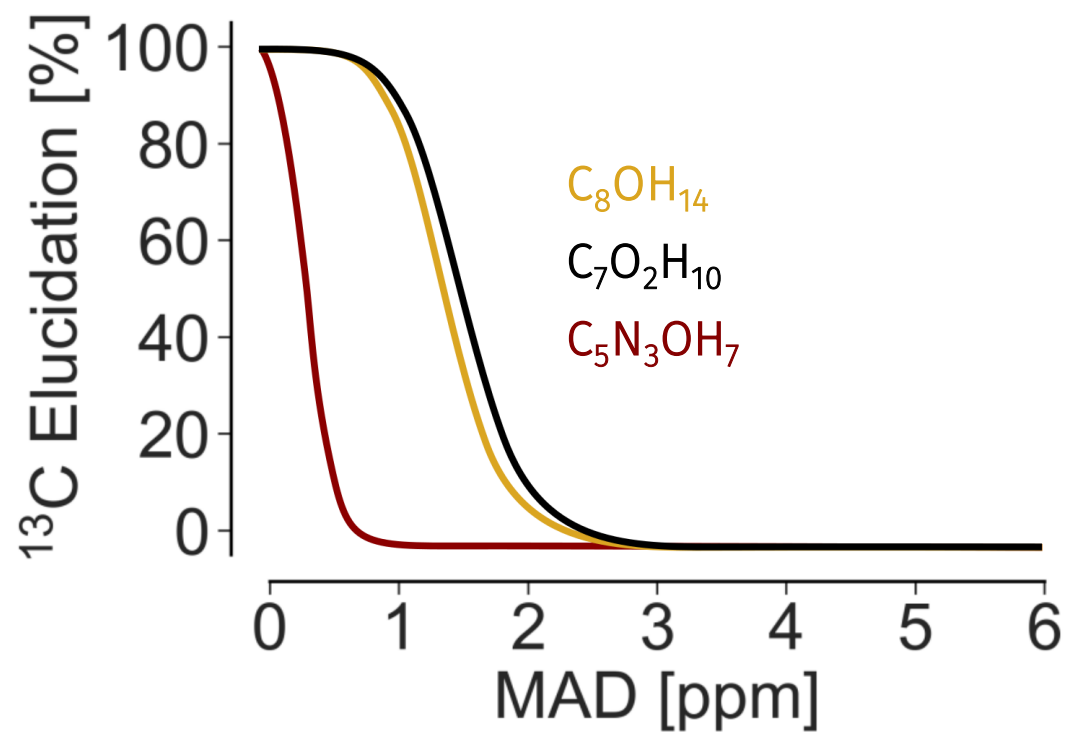
Possible graphs: ~50k



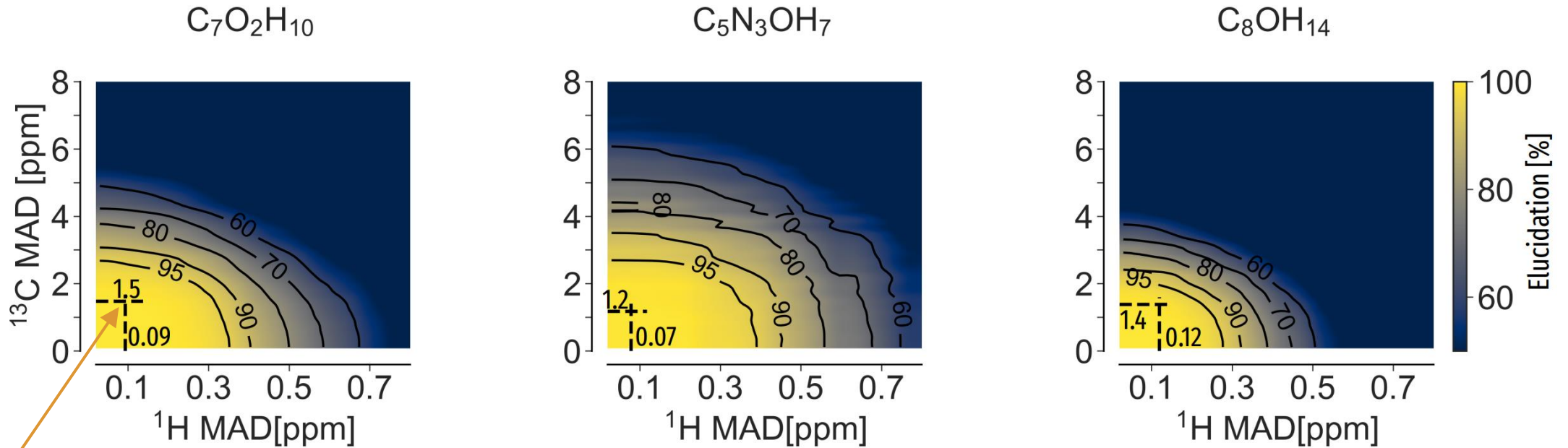
Optimistic: only one geometry/graph

Depends on

- number of graphs
- self-similarity

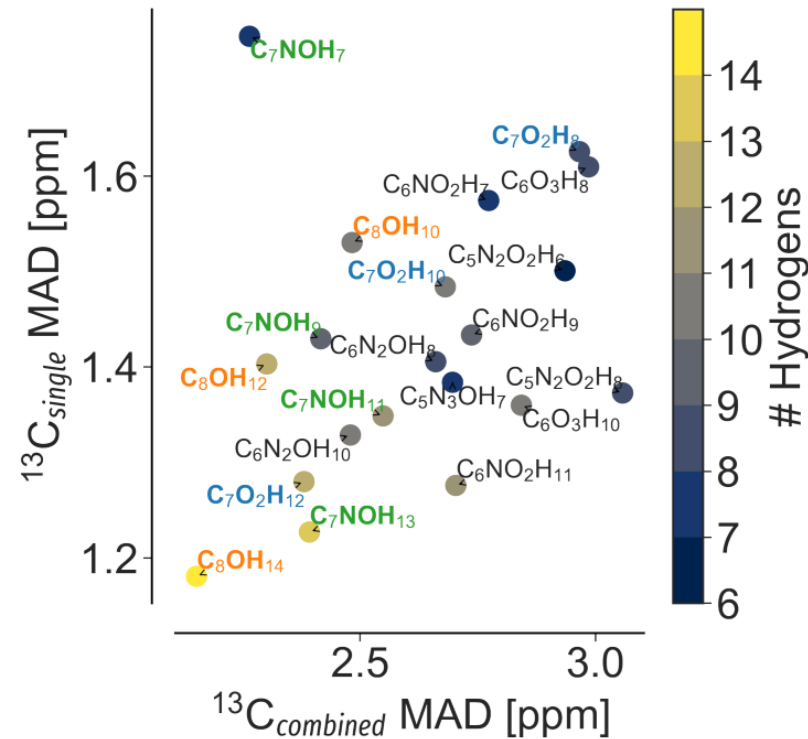


Combining ^{13}C and ^1H



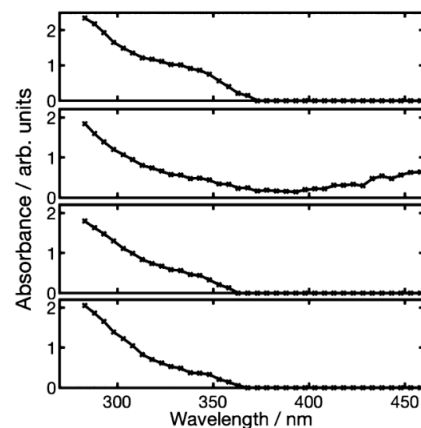
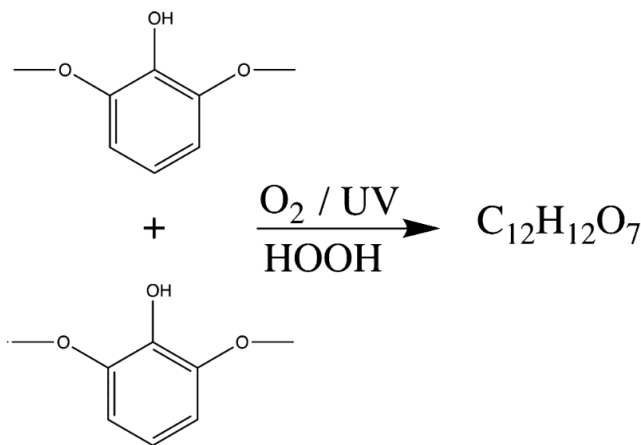
95% for only one spectrum

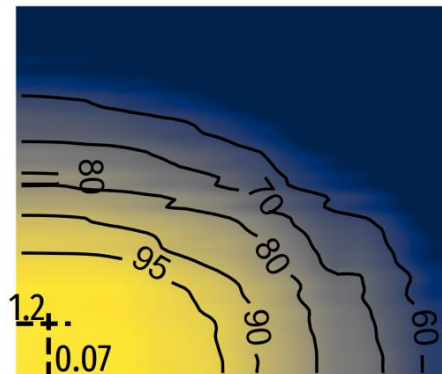
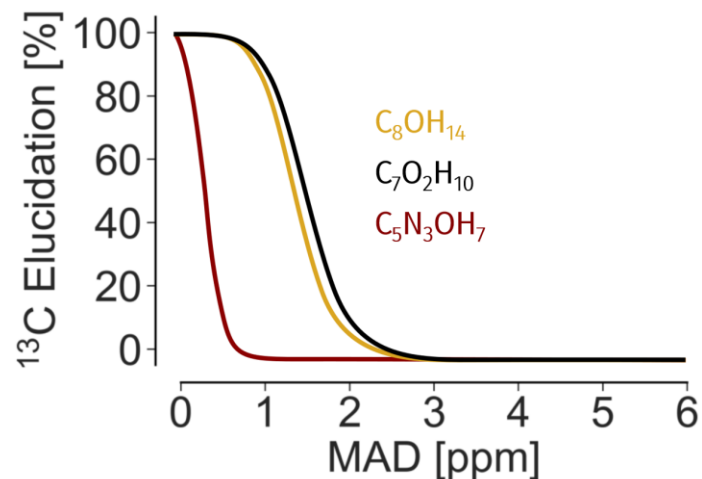
- Twice the permissible error: 1/100th of data points needed
- More hydrogens / fewer rings: harder problem despite more information





total molecules with two C ₆ rings	263,917,411
and which have OH groups	263,917,411
and which have no oxygen chain longer than 2	161,160,394
and which have an oxygen connecting the carbon rings	115,715,458
and which have one aromatic ring	134,944
and which are stable	64,121
and which match spectrum 1	36,518





Tradeoff | Multiple properties might be more valuable than accuracy

Data | Small-data regime potentially sufficient

NMR | Accuracy and gain quantified

NMR | ^{13}C alone insufficient even for small molecules