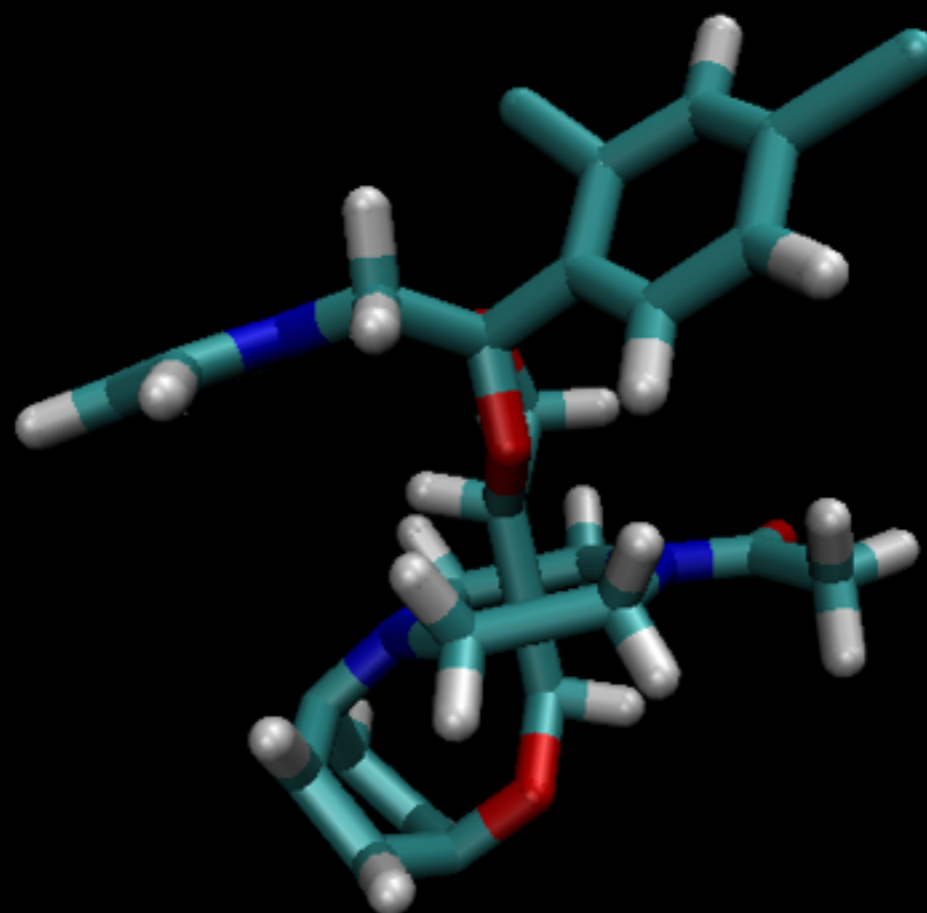# Difficult lessons learned from QM & p$K_a$ calculations in SAMPL5

Frank C. Pickard IV

D3R Meeting
11 March 2016

A "high energy" conformation
of molecule 92.

# SAMPL5 Description

A good proxy measurement for drug availability.

This is a "broad" definition of cosolvation, it includes ionic species, microsolvation, tautomers, snorkeling, *etc.*

$$\log D_{\text{chex/aq}}(x) = \log_{10}\left(\frac{[x]_{\text{chex}}}{[x]_{\text{aq}} + [x']_{\text{aq}} + [x'']_{\text{aq}} + ...}\right)$$

$$= (\Delta G_{\text{aq}} - \Delta G_{\text{chex}})\frac{\log_{10}(e)}{kT}$$

1 log D = 1.36 kcal/mol

# SAMPL5 Methods: log P calculation

- MM BAR

- QM Optimization

- QM NBB

- QM/MM Zwanzig

- Semi-Empirical NBB (in progress)

$$\log P_{\mathrm{chex/aq}}(x) = (\Delta G_{\mathrm{aq}} - \Delta G_{\mathrm{chex}}) \frac{\log_{10}(e)}{kT}$$

# SAMPL5 Methods: log P calculation

Non-Boltzmann Bennett Method

$$\Delta A = k_{\mathrm{B}} T \ln \left( \frac{\langle f(U_0 - U_1 + C) \rangle_1}{\langle f(U_1 - U_0 - C) \rangle_0} \right) + C \xrightarrow{V^b}$$

$$\langle X \rangle_{\mathrm{unbiased}} = \frac{\langle X \exp(V^b / k_{\mathrm{B}} T) \rangle_{\mathrm{biased}}}{\langle \exp(V^b / k_{\mathrm{B}} T) \rangle_{\mathrm{biased}}}$$

$$\Delta A = k_{\mathrm{B}} T \ln \left( \frac{\langle f(U_0 - U_1 + C) \exp(V_1^b / k_{\mathrm{B}} T) \rangle_1 \langle \exp(V_0^b / k_{\mathrm{B}} T) \rangle_0}{\langle f(U_1 - U_0 - C) \exp(V_0^b / k_{\mathrm{B}} T) \rangle_0 \langle \exp(V_1^b / k_{\mathrm{B}} T) \rangle_1} \right) + C$$

$$f(x) = \frac{1}{1 + \exp(x / k_{\mathrm{B}} T)}$$

König and Boresch J. Comp. Chem. (2010) 32, 1082.

# SAMPL5 Methods: log P calculation

QM/MM Non-Boltzmann Bennett

$$V^b = U_{\mathrm{MM}} - U_{\mathrm{QM}}$$

$$\Delta A = k_{\mathrm{B}} T \ln \left( \frac{\langle f(U_{0,\mathrm{QM}} - U_{1,\mathrm{QM}} + C) \rangle_{1,\mathrm{MM}}}{\langle f(U_{1,\mathrm{QM}} - U_{0,\mathrm{QM}} - C) \rangle_{0,\mathrm{MM}}} \right) + C$$

$$\Delta A = k_{\mathrm{B}} T \ln \left( \frac{\langle f(U_{0,\mathrm{QM}} - U_{1,\mathrm{QM}} + C) \exp(V_1^b/k_{\mathrm{B}}T) \rangle_{1,\mathrm{MM}} \langle \exp(V_0^b/k_{\mathrm{B}}T) \rangle_{0,\mathrm{MM}}}{\langle f(U_{1,\mathrm{QM}} - U_{0,\mathrm{QM}} - C) \exp(V_0^b/k_{\mathrm{B}}T) \rangle_{0,\mathrm{MM}} \langle \exp(V_1^b/k_{\mathrm{B}}T) \rangle_{1,\mathrm{MM}}} \right) + C$$

König, Pickard, Mei and Brooks J. Comput. Aided Mol. Des. (2014) 28, 245.

# SAMPL5 Methods: log P calculation

- CGenFF

- HREX Simulations, LD NVT

- 36 lambda points (6 electrostatic, 30 vdw)

- 1 fs timestep, 5 ns total

- 5000 QM or QM/MM calculations

# SAMPL5 Methods: log P calculation

- ## QM Optimization (our "control" submission)

  - w/ SMD Implicit Solvent (Vertical or Relaxed Solvation)

  - M06-2X/6-311++G**/6-31+G* with SMD

- ## QM NBB (optimized from SAMPL4 data)

  - w/ SMD Implicit Solvent

  - M06-2X/6-31+G* **or** OLYP/DZP

- ## QM/MM Zwanzig

  - w/ TIP3P Explicit Solvent

  - BLYP/6-31G*

# SAMPL5 Methods: log D correction

- p$K_a$ corrections

  - absolute/relative

  - vertical/relaxed solvation

- tautomerization (we only looked at aqueous)

- dimerization (in progress), trimerization, *etc.*

- wet cyclohexane

$$\log D_{\mathrm{chex/aq}}(x) = \log P_{\mathrm{chex/aq}}(x) + \mathbf{\Delta G}_{\mathrm{corr}} \frac{\log_{10}(e)}{kT}$$

$$\Delta G_{\mathrm{corr}} = \Delta G_{\mathrm{p}K_a} + \Delta G_{\mathrm{taut}} + \Delta G_{\mathrm{dimer}} + \Delta G_{\mu-\mathrm{solv}} + \ldots$$

# SAMPL5 Methods: p$K_a$ correction

Deprotonation Thermocycle:

$$\text{AH}^+_{(g)} \xrightarrow{\Delta G^\circ_{\text{gas}}(\text{AH}^+)} \text{A}_{(g)} \quad + \quad \text{H}^+_{(g)}$$

$$\downarrow \Delta G^*_{\text{solv}}(\text{AH}^+) \qquad \downarrow \Delta G^*_{\text{solv}}(\text{A}) \qquad \downarrow \Delta G^*_{\text{solv}}(\text{H}^+)$$

$$\text{AH}^+_{(aq)} \xrightarrow{\Delta G^\circ_{\text{aq}}(\text{AH}^+)} \text{A}_{(aq)} \quad + \quad \text{H}^+_{(aq)}$$

Convert to p$K_a$:

$$\Delta G^\circ_{\text{aq}}(AH^+) = \ln(10) RT \text{p}K_{\text{a}} + \Delta G^{\circ \rightarrow *}$$

Convert to populations at pH = 7.4:

$$\text{pH} = \text{p}K_a + \log_{10}\left(\frac{[\text{A}_{(aq)}]}{[\text{A}^+_{(aq)}]}\right)$$

Convert to free energy of protonation…

# SAMPL5 Methods: p$K_a$ correction

$$AH^+_{(g)} \xrightarrow{\Delta G^\circ_{gas}(AH^+)} A_{(g)} \quad + \quad H^+_{(g)}$$

$$\downarrow \Delta G^*_{solv}(AH^+) \qquad\qquad \downarrow \Delta G^*_{solv}(A) \qquad\qquad \downarrow \Delta G^*_{solv}(H^+)$$

$$AH^+_{(aq)} \xrightarrow{\Delta G^\circ_{aq}(AH^+)} A_{(aq)} \quad + \quad H^+_{(aq)}$$

Problematic Terms

## Absolute p$K_a$ calculations:

- Proton solvation free energy from experiment (265.9 kcal/mol)

  - Small experimental errors can yield **big** p$K_a$ errors!

- Robustly treat molecules with coupled protonation/tautomerization

- Shouldn't use with vertical solvation (expensive)

# SAMPL5 Methods: p$K_a$ correction

$$\mathrm{p}K_a = \widetilde{\mathrm{p}K_a} + \left[ G(\mathrm{A_{(aq)}}) - G(\mathrm{AH^+_{(aq)}}) - \underline{G(\mathrm{L_{(aq)}}) + G(\mathrm{LH^+_{(aq)}})} \right] / \left[\ln(10)RT\right]$$

Analogue
experiment

Analogue
calculation

## Relative p$K_a$ calculations:

- Experimental proton solvation free energy term drops out.

  - Uncertainty from analogue experiment

  - Results sensitive to analogue choice

- Can be more accurate than absolute calculations

- Can be used with vertical solvation (cheap)

Test Set: Cohort0

03

15

17

20

37

45

55

58

59

61

68

70

80

# Test Set: Cohort1

# Test Set: Cohort2



02

06

13

19

24

33

49

50

65

67

69

74

75

82

83

84
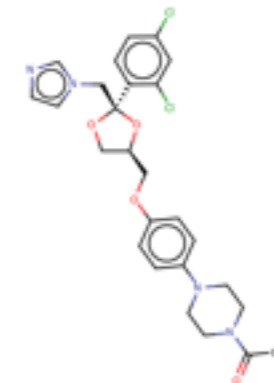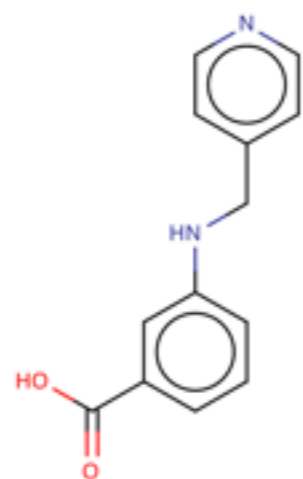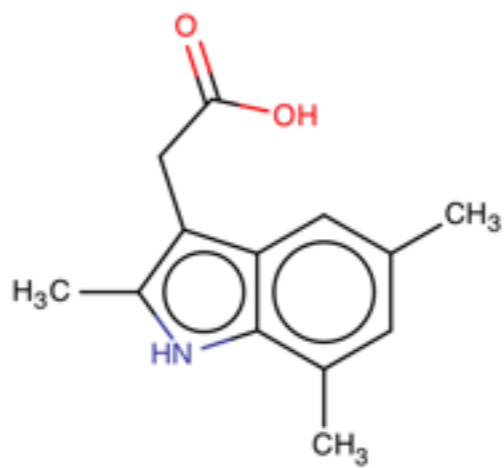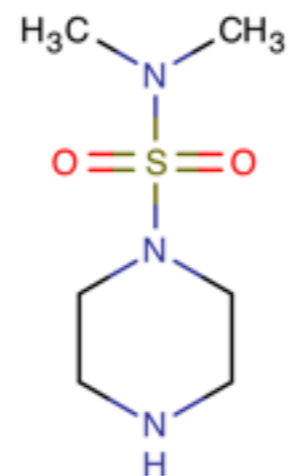
85

86

88

92

# Test Set: pKₐ Baddies (simple)



10

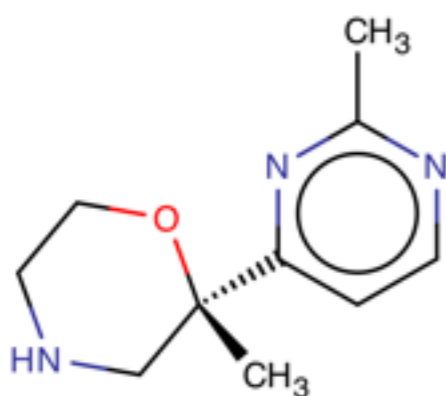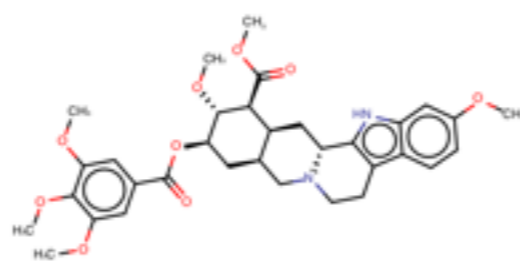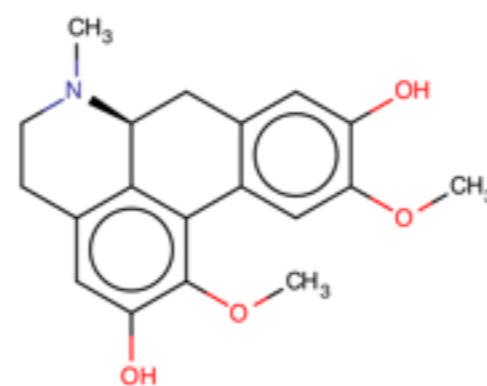26

37

47

61

65

69

70

84

85

92

11

15

17

50

56

60

63

83

# SAMPL5 Results



Rankings by RMS Error, Top 20, Cohort 0+1+2

QM Control: M06-2X/6-31+G*/SMD

QM | Other
NBB | MM

# SAMPL5 Results



Rankings by Pearson's R, Top 20, Cohort 0+1+2

# SAMPL5 Results



Rankings by Kendall's Tau, Top 20, Cohort 0+1+2

# SAMPL5 Results



Distribution Coefficients (Theory vs. Experiment)

# SAMPL5 Results



Distribution Coefficients (Theory vs. Experiment)

Legend:
- agreement
- $1\,k_\mathrm{B}T$ error
- $5\,k_\mathrm{B}T$ error

X-axis: Experimental ($\log_\mathrm{D}$)
Y-axis: Calculated ($\log_\mathrm{D}$)

# SAMPL5 Results



QM NBB OLYP/DZP

RMSE: 2.32
$\tau$: 0.46
LSq Slope: 1.08
LSq Shift: −0.45
LSq $R$: 0.63

Calculated ($\log_D$)

Experimental ($\log_D$)

— LSq Fit

Our predictions are significantly hydrophilic.

# SAMPL5 Results: The Baddies



13 of 53 predictions have errors > 5 kT.

# SAMPL5 Results: p$K_a$ corrections

# SAMPL5 Results: p$K_a$ corrections



QM NBB OLYP/DZP Abs. p$K_a$ Solv

# SAMPL5 Results



## QM Optimization M06-2X/6-31+G* (Naive QM)

RMSE: 2.58
$\tau$: 0.46
LSq Slope: 1.15
LSq Shift: −0.16
LSq $R$: 0.61

LSq Fit

Experimental ($\log_D$)

Calculated ($\log_D$)

# SAMPL5 Results



QM Optimization M06-2X/6-31+G* (Naive QM)

# SAMPL5 Results



QM Optimization M06-2X/6-31+G* Abs. $pK_a$ VSolv

# SAMPL5 Results



QM Optimization M06-2X/6-31+G* Abs. $pK_a$ VSolv

RMSE: 2.74
$\tau$: 0.55
LSq Slope: 1.51
LSq Shift: $-0.50$
LSq $R$: 0.72

Calculated ($\log_D$)

Experimental ($\log_D$)

LSq Fit
$pK_a \sim 7.4$

# SAMPL5 Results



MM BAR

RMSE: 5.57
$\tau$: 0.22
LSq Slope: 1.18
LSq Shift: 2.31
LSq $R$: 0.37

Calculated ($\log_D$)

Experimental ($\log_D$)

— LSq Fit

# SAMPL5 Results



MM BAR Abs. $pK_a$ Solv

RMSE: 3.14
$\tau$: 0.49
LSq Slope: 1.62
LSq Shift: 0.59
LSq $R$: 0.69

Calculated ($\log_D$)

Experimental ($\log_D$)

— LSq Fit
● $pK_a \sim 7.4$

70
69
26
83
50 85

# SAMPL5 Results



QM/MM Zwanzig BLYP/6-31G*

RMSE: $3.01$
$\tau$: $0.48$
LSq Slope: $1.48$
LSq Shift: $-0.42$
LSq $R$: $0.67$

— LSq Fit

Calculated ($\log_D$)

Experimental ($\log_D$)

# SAMPL5 Results



QM/MM Zwanzig BLYP/6-31G* Abs. $pK_a$ Solv

RMSE: $3.34$
$\tau$: $0.55$
LSq Slope: $1.78$
LSq Shift: $-0.83$
LSq $R$: $0.73$

LSq Fit
$pK_a \sim 7.4$

# SAMPL5 Results: pKa corrections



pKa and tautomerization correction to free energy (Relaxed Solvation)

Large differences between pKa methods might indicate tautomerization issues.

# SAMPL5 Results: Conclusions

- Predicting log D values is *difficult*.

- Lessons learned from SAMPL4 have carried over (choice of density functional and basis set).

- NBB QM calculations with implicit solvent are among the best options (RMSE rank 2nd) but have poor correlation.

- Our predictions are too hydrophilic, we ignored the wetness in cyclohexane

- Accounting for tautomers is very important, and universally improves our correlation (but reduces RMSE).