

# COSMO-RS based predictions for the SAMPL6 logP challenge

Christoph Loschen, Jens Reinisch, Andreas Klamt  
COSMOlogic GmbH & Co. KG – Dassault Systèmes

[loschen@cosmologic.de](mailto:loschen@cosmologic.de)

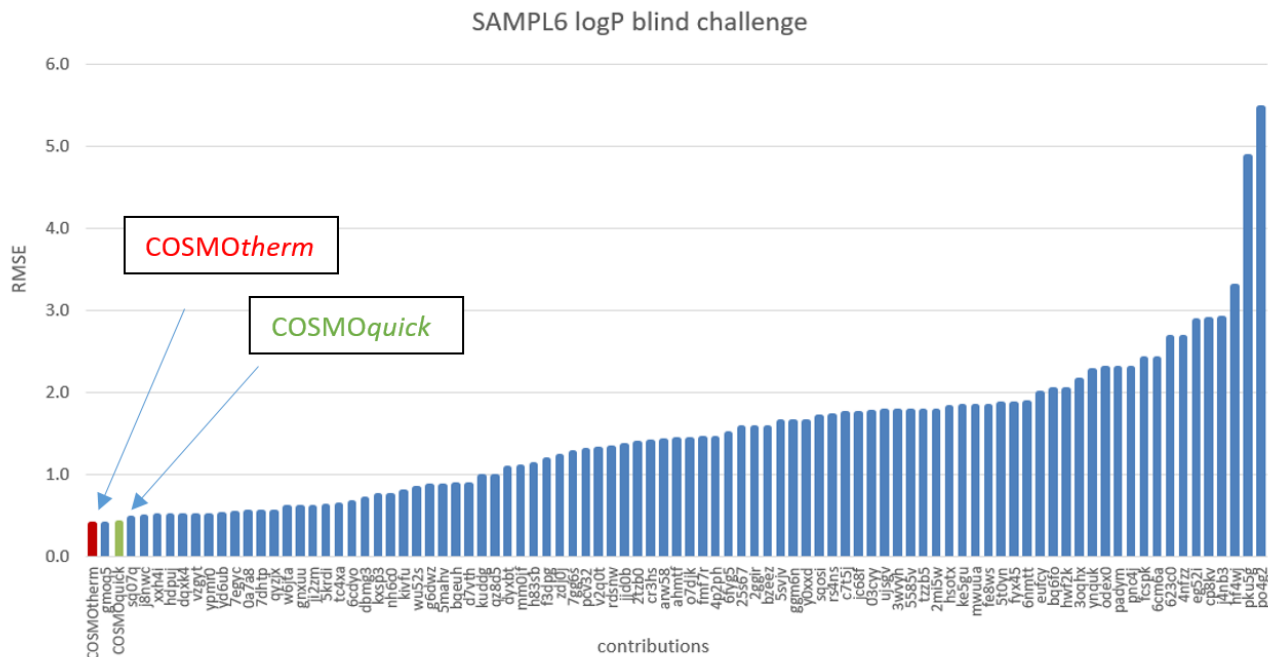
# SAMPL6 logP blind challenge

## Methods used

- QSPR/machine learning
- Deep Learning
- Molecular Dynamics
- Implicit solvation models (SMD)
- 3D-RISM
- COSMO-RS

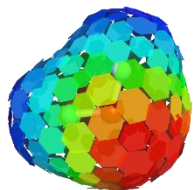
# Results

## COSMO-RS based methods

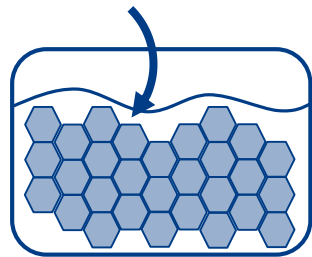


# COSMO-RS

In a nutshell

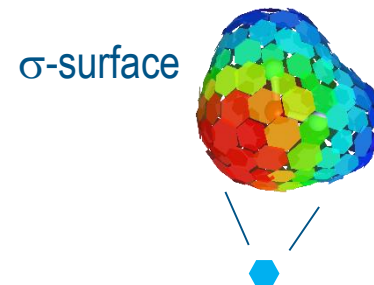


➤ COSMO: implicit solvation model via DFT



➤ COSMO-RS: statistical thermodynamics

- Electrostatic:
- Hydrogen bonds:
- Van der Waals:
- Combinatorial term:



Intermolecular interactions  
are based on  
 $\sigma$  surface segments

$$E \sim (\sigma + \sigma')^2$$

$$E \sim (\sigma \times \sigma')$$

$$E \sim \text{area}$$

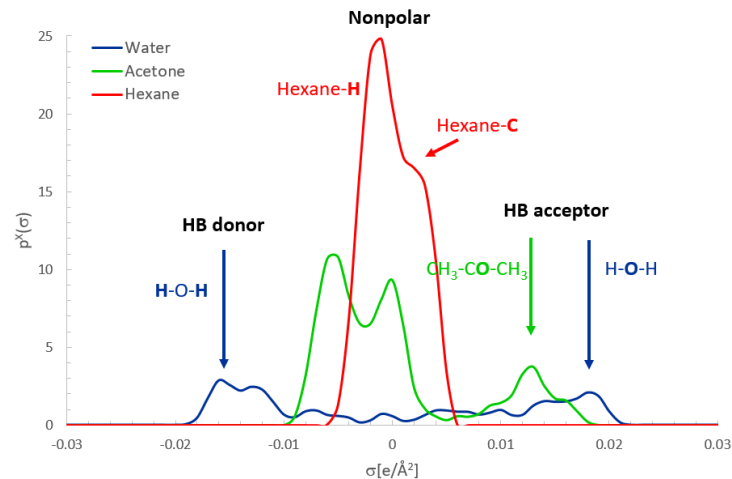
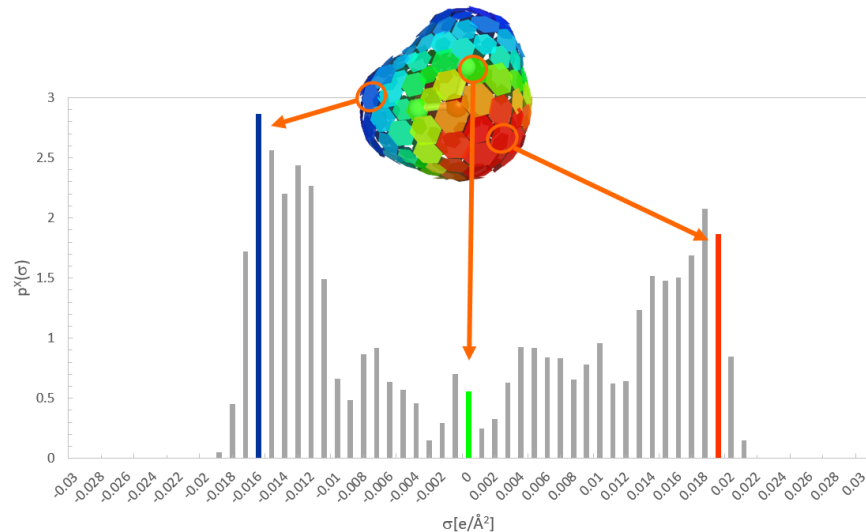
$$E (\text{Shape})$$

Klamt, A. *J. Phys. Chem.* **1995**, 99, 2224-2235.

Klamt, A. *WIREs: Comput. Mol. Sci.* **2011**, 1, 699-70

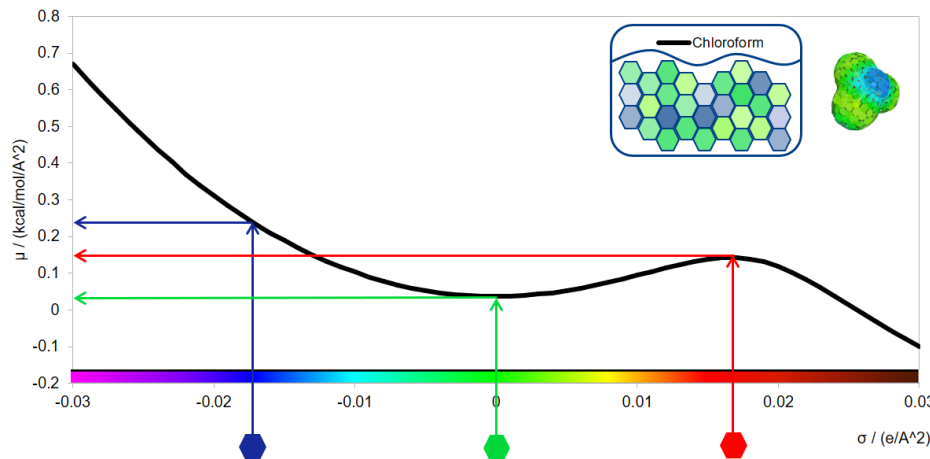
# COSMO-RS

$\sigma$ -profile  $p(\sigma)$ : a histogram of charged surface segments of a molecule



# COSMO-RS

The  $\sigma$ -potential  $\mu_s(\sigma)$  is a characteristic function of a system at a given T.



From  $\mu_s(\sigma)$  one obtains the chemical potential of a substance in solution  $\mu_s$  and **all** related properties:

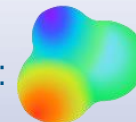
$$\mu_s = \int d\sigma p(\sigma) \mu_s(\sigma) + RT \ln \{x \gamma_{comb.,s}\}$$

# COSMOtherm workflow

Conformational sampling  
in liquid phase:  
COSMOconf v4.3



Liquid Phase Statistical  
Thermodynamics COSMO-RS:  
COSMOtherm v19



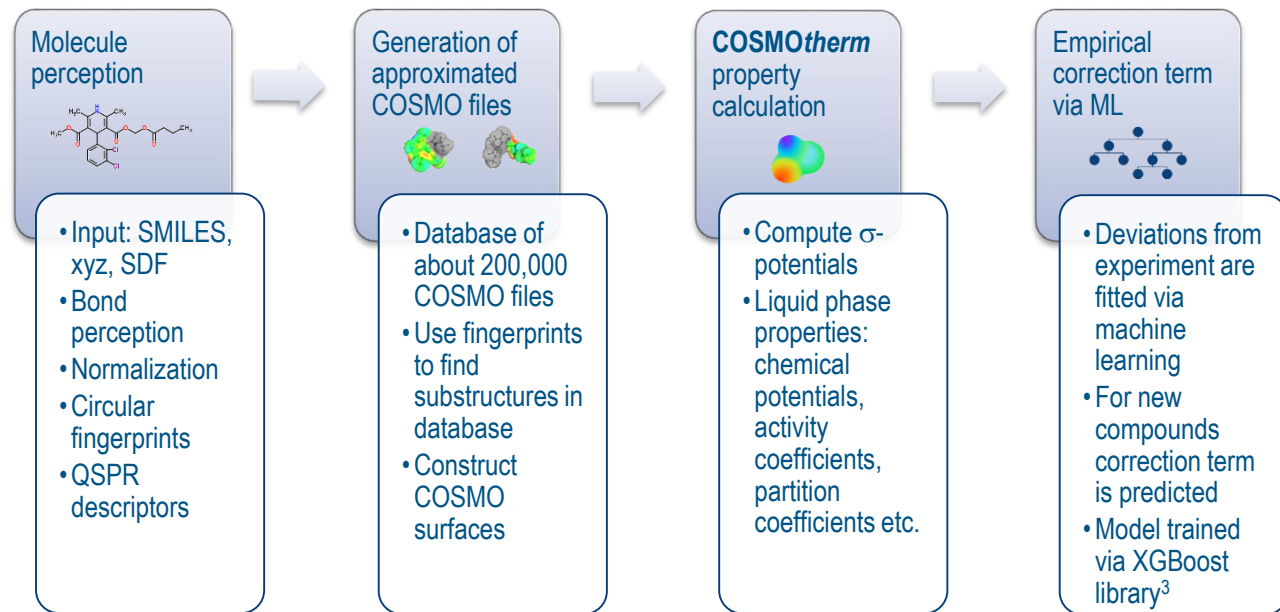
- Structure generation (Balloon & RDKit)
- Iterative conformer reduction according to energy and clustering of structures and (liquid) chemical potentials
- COSMO-Levels used: BP/SV(P) , BP/TZVP & BP/TZVPD (Turbomole v7.3)<sup>3</sup>
- Identification of optimal conformational set

- COSMO-RS based computation of chemical potentials in water & octanol<sup>4-6</sup>
- Consideration of conformational effects
- Assuming wet octanol (27.4% mf water)
- Parameterization: BP\_TZVPD\_FINE\_19

$$\log_{10}(P) = \log_{10}\left(\frac{c_1}{c_2}\right) = (\mu_2^x - \mu_1^x) / RT \ln(10) + \log_{10}(V_2 / V_1)$$

1. COSMOconf 4.3; COSMOlogic GmbH & Co. KG; <http://www.cosmologic.de>; Leverkusen, Germany, 2018.
2. Klamt, A.; Eckert, F.; Didenhofen, J. Phys. Chem. B 2009, 113 (14), 4508–4510.
3. TURBOMOLE V7.3; University of Karlsruhe and Forschungszentrum Karlsruhe GmbH, 1989-2007, TURBOMOLE GmbH, since 2007; available from <http://www.turbomole.com>; Karlsruhe, Germany, 2018.
4. Klamt, A. J. Phys. Chem. 1995, 99 (7), 2224–2235.
5. Klamt, A.; Jonas, V.; Bürger, T.; Lohrenz, J. C. J. Phys. Chem. A 1998, 102 (26), 5074–5085.
6. COSMOtherm, Release 19; COSMOlogic GmbH & Co. KG; <http://www.cosmologic.de>; Leverkusen, Germany, 2019.

# COSMOquick workflow



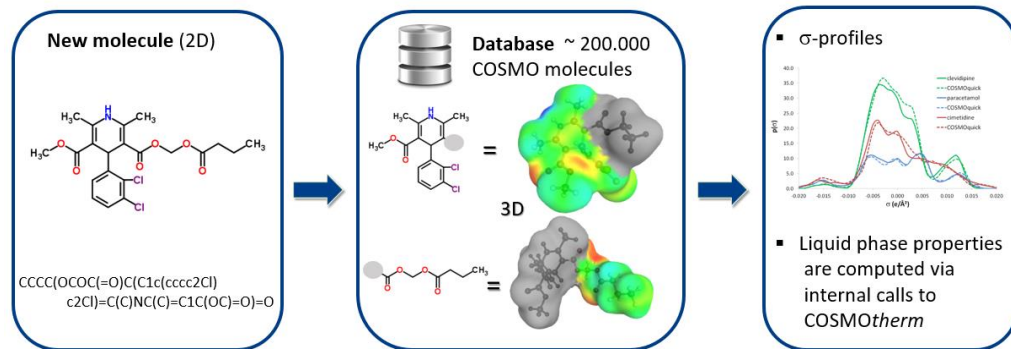
1. COSMOquick 1.7; COSMOlogic GmbH & Co. KG; <http://www.cosmologic.de>; Leverkusen, Germany, 2018.  
2. Chen, T.; Guestrin, C. Xgboost: A Scalable Tree Boosting System. In Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining; ACM, 2016; pp 785–794.



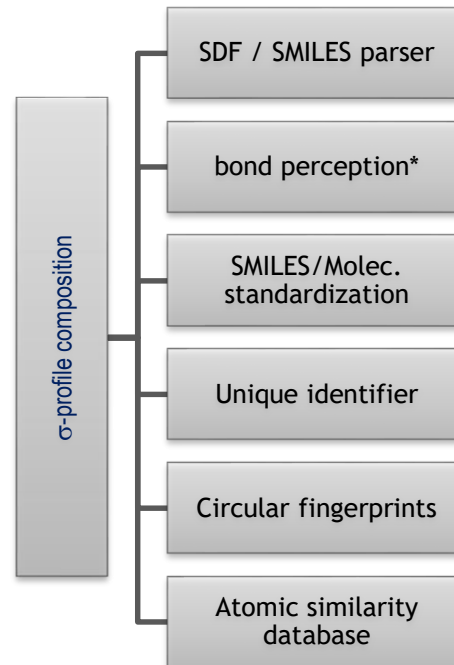
# COSMOquick: instant $\sigma$ -profile composition

Idea: Compose larger molecules from a database of pre-calculated molecules

Assumption:  $\sigma$ -profiles of compounds are somewhat additive!



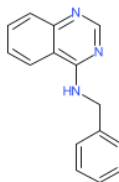
cheminformatics backend



# COSMOquick: instant $\sigma$ -profile composition

Compound	N,fragments	<similarity>
SM02	5	4.9
SM04	4	5.2
SM07	3	6.2
SM08	4	3.5
SM09	6	4.9
SM11	0.0	9.0
SM12	4	6.3
SM13	2	6.5
SM14	3	4.3
SM15	4	4.0
SM16	2	5.6

Example: SM07

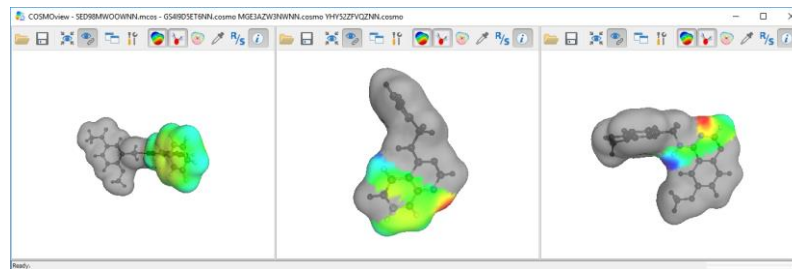


Atomic weight strings:

w={00000000000000001000001111111000000000001001111111000}

w={1111000000000011011000000010}

w={00000001111000000000000000001100000000}



1: lowest similarity (single atom match)

9: highest similarity (full match)

- significant fragmentation effect expected!
- Only a fraction of a second needed for  $\sigma$ -profile

# COSMOquick

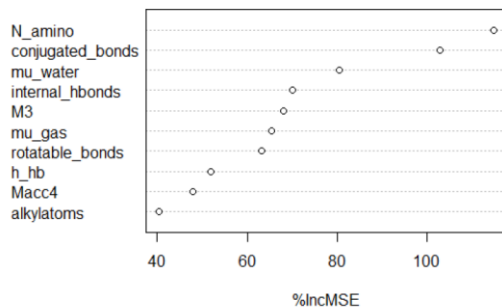
ML correction term:

$$\log P = \log P_{TZVP} + \Delta \log P_{ML-corr}$$

COSMOtherm  
& approx.  $\sigma$

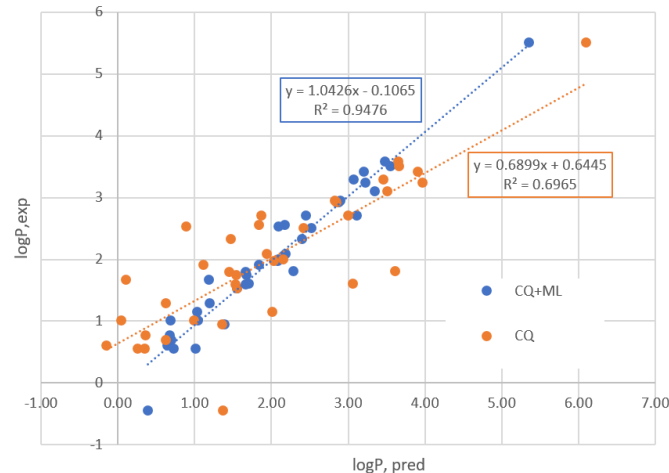
XGBoost

variable(i.e. descriptor) importance



data type	n	source
Training & crossval.	10964	PHYSPROP
Test	37	PHYSPROP subset via substructure search

test set results (RMSE=0.46)

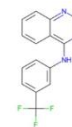


# Results

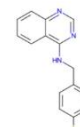
## Comparison of COSMO $\text{therm}$ submissions

compound	logP <sub>exp</sub>	COSMO $\text{therm}$	COSMOquick
SM02	4.09	4.42	3.98
SM04	3.98	3.86	3.74
SM07	3.21	3.48	3.19
SM08	3.1	2.85	2.74
SM09	3.03	3.44	3.38
SM11	2.1	2.00	2.64
SM12	3.83	3.82	4.29
SM13	2.92	<b>3.84</b>	3.41
SM14	1.95	2.21	2.24
SM15	3.07	2.77	<b>2.28</b>
SM16	2.62	3.05	2.88

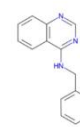
id	$\sigma$ -surface from	Method	n , fragments	level	RMSE
hmz0n	COSMO files	Turbomole	-	FINE19	0.38
3vqbi	SMILES	COSMOfrag	37	TZVP+ML	0.41
	COSMO files	Turbomole	-	TZVP+ML	<b>0.35</b>
0.5 * hmz0n+0.5 * 3vqbi	COSMO files / SMILES	consensus	-	FINE19 TZVP+ML	<b>0.34</b>



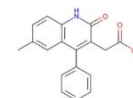
SM02



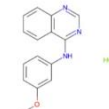
SM04



SM07



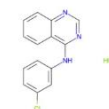
SM08



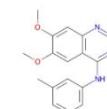
SM09



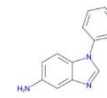
SM11



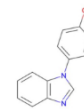
SM12



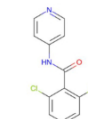
SM13



SM14



SM15

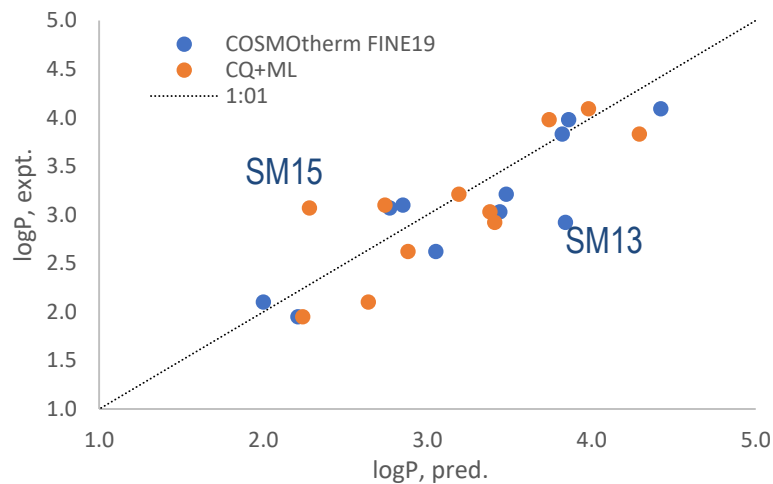


SM16

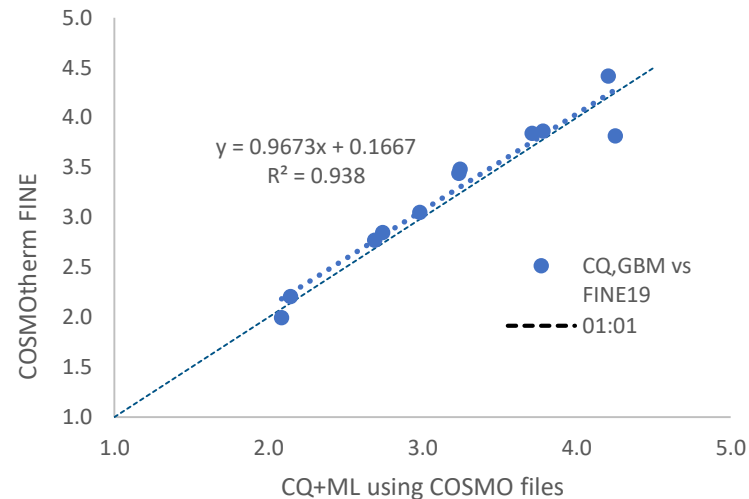
# Results

## Comparison of COSMOtherm Submissions

versus experiment



with each other without fragmentation



# SAMPL6

## Learnings

- ▶ Need for systematic tautomer & deprotonation workflow (pKa part)
- ▶ TZVP (intermediate level) shows systematic problems with this compounds class
- ▶ Probably more potential for ML based approaches
- ▶ Monitor fragmentation error (COSMO*quick*)

Many thanks to the SAMPL6 organizers for setting up the challenge!