# Escape from 'availability bias' in compound design

ANDRÁS STRÁCZ, ÁKOS TARCSAY, IVÁN SOLT, GÁBOR IMRE

**ABSTRACT** Small molecule design is an information demanding activity, since all relevant knowledge is to be accessible within a single space and requires synchronized application of computational models to assist decision making on synthesis candidates. Our study aims to evaluate a software platform coping with this complexity (Marvin Live). The tool provides central management of innovative ideas and a framework that helps triage these. Decision support is based on predicted properties that span phys-chem descriptors, combined metrics like CNS MPO score, 3D overlay and modelling results conducted with KNIME on the one hand, and cross-checking with available knowledge collected from a variety of sources to do rapid freedom to operate analysis and SAR by catalog by ultra-fast searching of patent and compound catalog databases, respectively on the other. This hypothetic study shows the potential evolution of a compound idea from the reference compound to a synthesis candidate through an example discovery project of 5-hydroxytryptamine receptor 6 (HTR6) ligands.

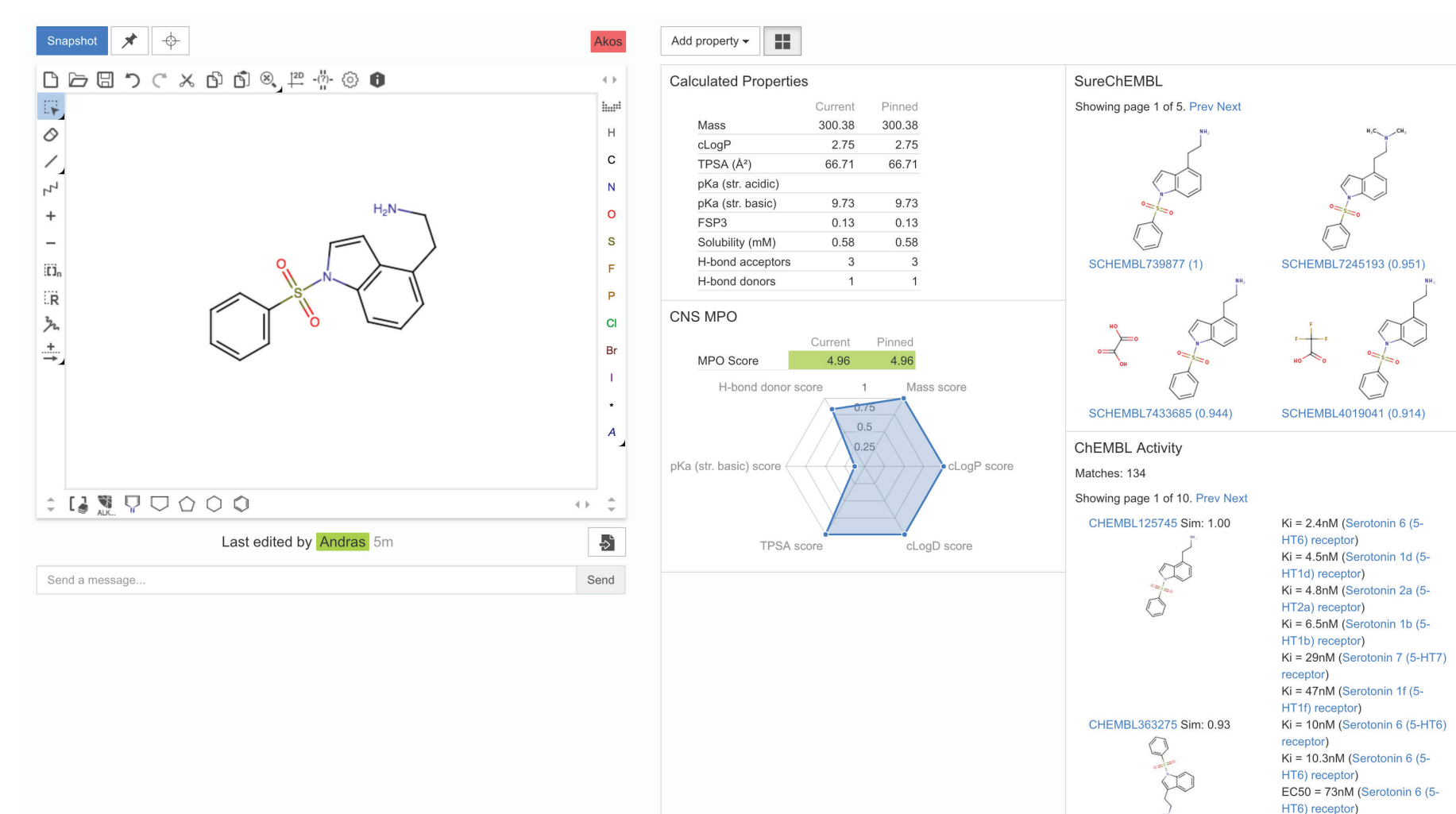## 1. Characterization of the reference compound

Pros:
- Balanced phys-chem properties (**SCHEME 1.** Calculated properties)
- pKa limits the CNS MPO[1] score (**SCH. 1.** CNS MPO)
- Body of public activity data available (**SCH. 1.** ChEMBL Activity)
- High affinity (2.4 nM) on the primer target (HTR6) (**SCH. 1.** ChEMBL Activity)
- SAR can be deduced from literature data
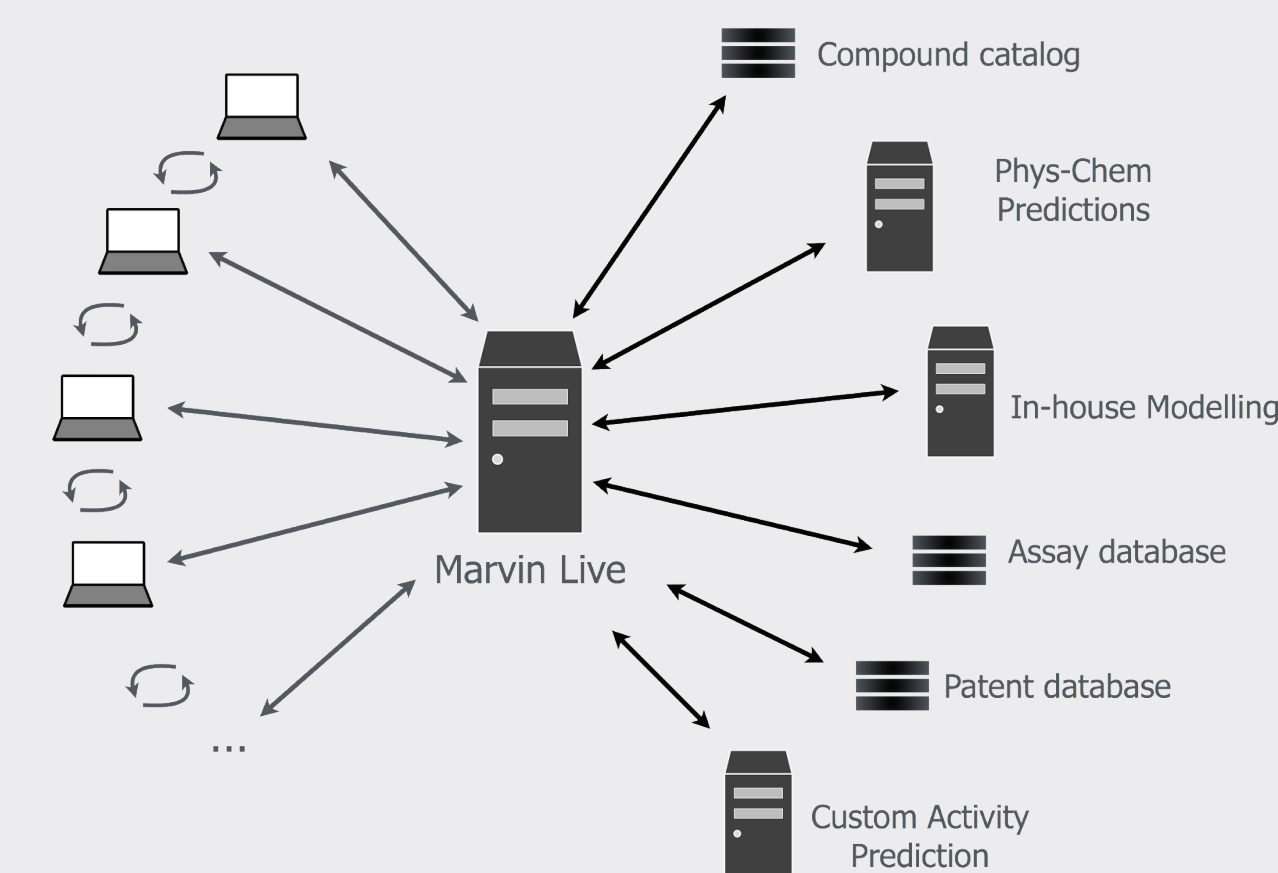- 134 ChEMBL analogs are available within 0.8 Tanimoto similarity threshold

Cons:
- Patented (**SCH. 1.** SureChEMBL[2])
- Selectivity might be an issue

**SCHEME 1.** Interface to create chemical structures and to display results of plugins. Calculated phys-chem properties, most similar hits from SureChEMBL, and most similar compound from ChEMBL with public assay data are shown. Active links allow drill down for more information from the original source. ▶
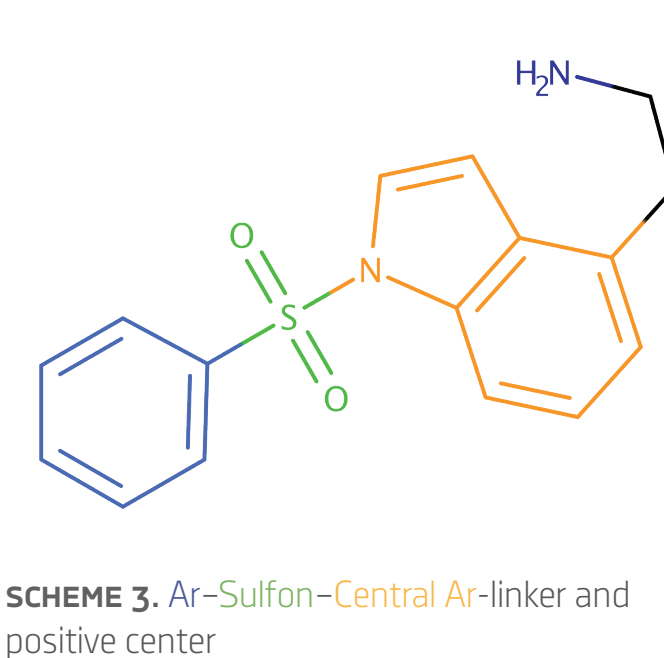


### Marvin Live

A NodeJS[3] based plugin system was created to communicate with available services. The plugin code sets up the connection to databases, web services, etc., handles the conversion of the molecule to the desired chemical structure format and formats the query to match the provided APIs. When the request completes rich HTML and reportable key-value pairs are sent to the connected browsers. This way Marvin Live and its plugin system acts as a central communication channel (**SCHEME 2.**) through which any internal or public resource is available to chemists for ideation and optimization.



**SCHEME 2.** A schematic representation of how Marvin Live acts as an interface between chemists and diverse set of cheminformatics and modelling capabilities.

## 2. Scaffold hopping to IP-free space

- Generation of ideas within the space defined by the pharmacophore elements (**SCHEME 3**)
- Monitor freedom to operate (FTO) by instant similarity search against SureChEMBL database using MadFast Similarity Search[4]
- Preferred structures are located in the chemical space below 0.8 Tanimoto similarity threshold from exemplified structures collected from patent literature
- Top 4 excerpt of the closest analogs are shown in **TABLE 1**.
- 3 out of 5 exemplified structures can be considered in the next round of idea triage (**TABLE 1.**)
- Please note that this rapid similarity search does not include search in Markush structure, a more detailed IP analysis can be achieved by consulting with IP-experts
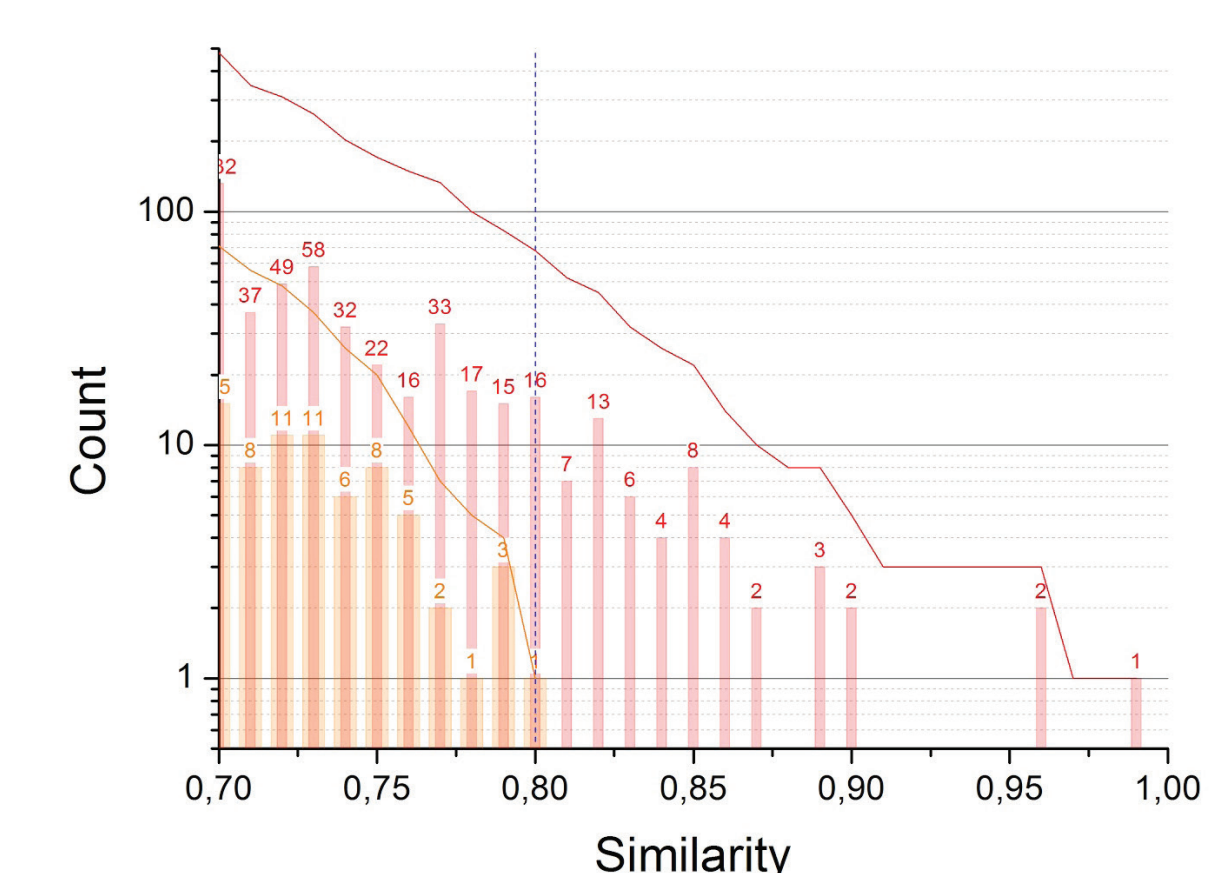- Comparison of abundantly (HTR6 H2L - 9) and less represented (HTR6 H2L - 8) structures with regard to IP space are shown on **SCHEME 4.**



**SCHEME 3.** Ar–Sulfon–Central Ar–linker and positive center



◀ **TABLE 1.** Overlap analysis of the designed analogs versus the chemical space of exemplified structures from patents (SureChEMBL). Heat coloring of the Tanimoto similarity values reveals the IP risk.



**SCHEME 4.** Similarity based distribution of SureChEMBL analogs. The number of analogs are shown with binning size of 0.01 Tanimoto similarity values. ▶
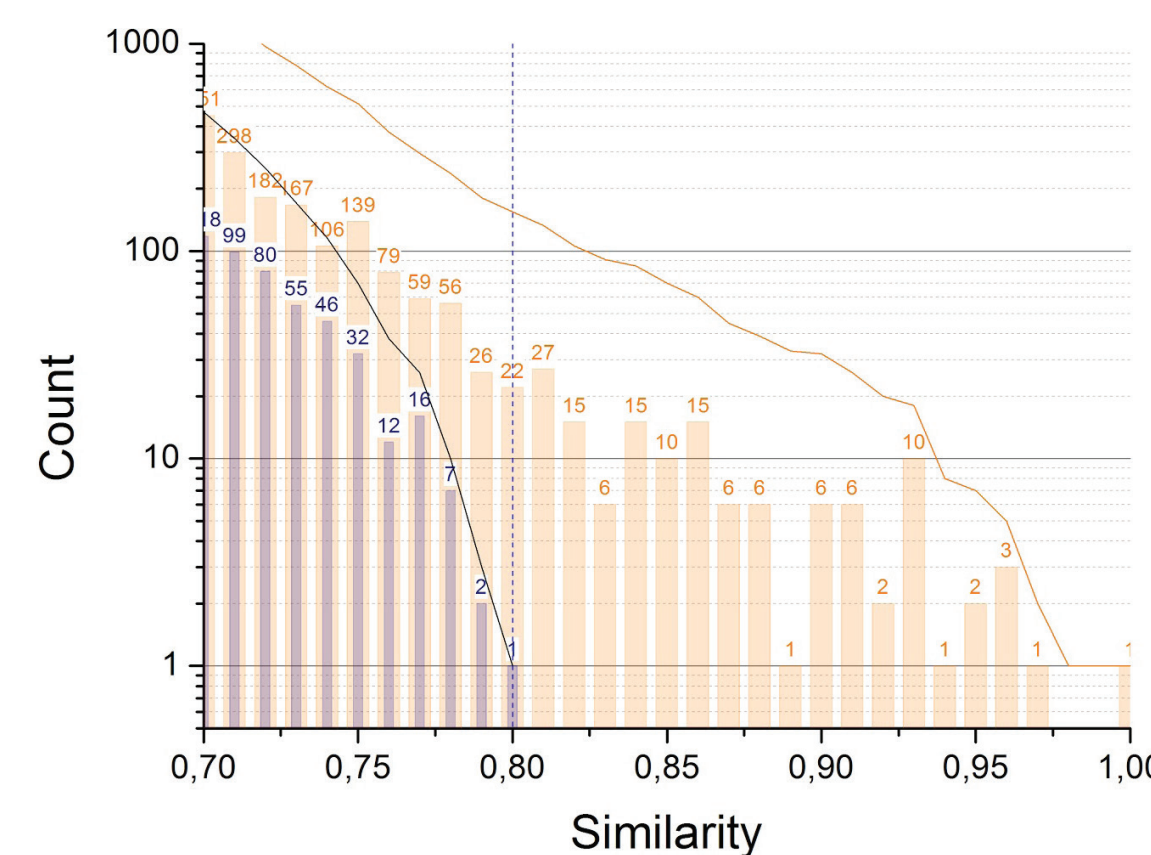
## 3. Scaffold availability

- Assess feasibility by similarity based exploration of the availability in Enamine REAL database[5, 6]
- Top 4 excerpt of the closest analogs are shown in **TABLE 2.**
- Check opportunity to explore early SAR based on most similar analogs present in Enamine REAL. Example of well represented and purly represented structures are shown on **SCHEME 5**
- Two derivatives (HTR6 H2L - 6 and - 8) are available in the vendor DB



**TABLE 2.** Overlap analysis of the top ranked scaffolds with Enamine REAL ▶



◀ **SCHEME 5.** Similarity based distribution of Enamine REAL analogs. The number of analogs are shown with binning size of 0.01 Tanimoto similarity values.

### MadFast Similarity Search

The Enamine REAL and SureChEMBL databases searched using 1024 bit Chemical Hashed Finger-prints[7] for the 171.5 and 17 million molecules, respectively. These were loaded for in-memory similarity search into a single virtual machine having 24 cores and 64 GB RAM. The high performance search engine, MadFast Similarity Search was accessed through a Marvin Live plugin, to automatically run similarity searches on the selected molecule with a maximum hit count of 20 sorted by similarity. The measured search times are displayed in **TABLE 1** and **TABLE 2.**

## 4. 3D alignment of the scaffolds to the reference

- Evaluation of steric complementarity with the pharmacophore elements of the reference compound revealed that HTR6 H2L - 8 fits, while HTR6 H2L - 6 displays larger deviations.
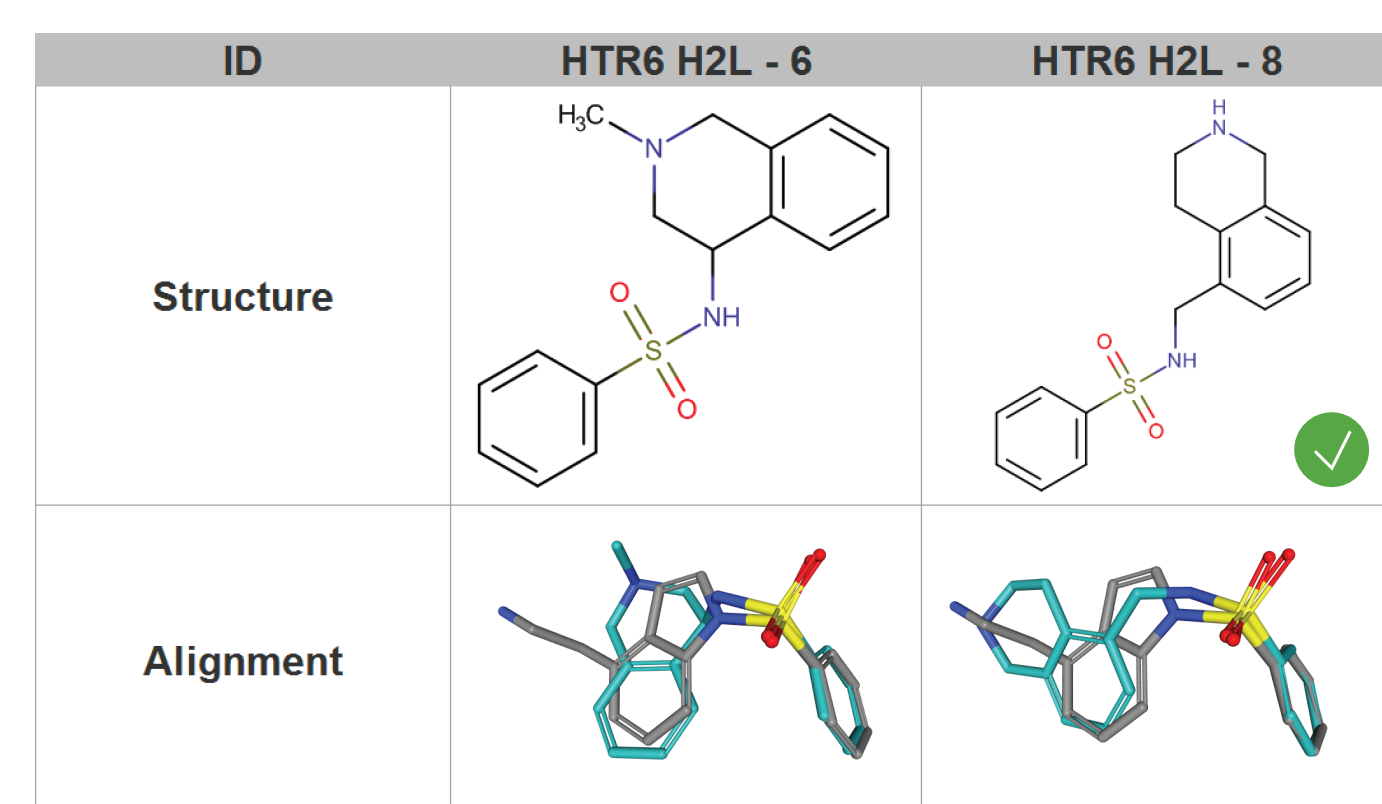


**TABLE 3.** MCS based 3D alignment of the two selected compounds. Alignment was done by 3D align tool from ChemAxon.[8]

## 5. SAR by catalog

- The *design mode* was used to find close structural analogs in the Enamine REAL catalog
- Halo- and methyl scan on the phenyl ring were both available from catalog as well as ethyl substituted derivatives
- The collection of close analogs was filtered based on similarity to patented structures from SureChEMBL in the *overview mode* (**SCHEME 6**)
- Sorting on values provided by plugins, like *Mass, cLogP* or *Solubility* aids early triage and decision making based on compound quality
- Status and task assignment facilitates the creation of a cherry picked library as a distinguished subgroup.



**SCHEME 6.** Overview mode of the finally selected custom library around HTR6 H2L -8 with selected calculated properties. ▶

## 6. Comprehend one-by-one

- Assessment of calculated properties compared to the reference compound (Pinned structure)
- Radar chart with visualization of CNS MPO increments for detailed comparison
- Cross-check using 3D alignment plugin
- Predict hERG liability by a utilizing complex KNIME workflow as a plugin



**SCHEME 7.** Comparison of one selected analog to the reference compound in Design mode. ▶

### Conclusions

Our results reveal that ligand design can be fostered by providing instant access to extra large compound collections, knowledge bases and predicted properties. Therefore, it facilitates widening the scope of the chemical space and helps escape an availability bias that would constrain the series of ideas and increase the likelihood of project failure due to insufficient freedom to operate.

The rich visualization in design mode and the idea management in overview mode help prioritizing the compound series and decision making in Marvin Live.

[1] Wager, Travis T., et al. ACS Chem. Neurosci., **2010**, 1 (6), pp 435–449
[2] Papadatos, George, et al. Nucleic Acids Res. **2016** Jan 4; 44(Database issue): D1220–D1228. Database downloaded 01 May 2017
[3] https://nodejs.org/ Accessed 04 June 2017
[4] https://www.chemaxon.com/products/madfast/ Accessed 04 June 2017
[5] Shivanyuk, A. N., et al. Chemistry Today **2007**, 25, 58-59. Database downloaded 01 May 2017
[6] Tolmachev, Andrey, et al. ACS Combinatorial Science **2016**, 18.10, 616–624.
[7] https://docs.chemaxon.com/display/docs/Chemical+Hashed+Fingerprint Accessed 04 June 2017
[8] https://www.chemaxon.com/products/calculator-plugins/molecular-modelling/ Accessed 04 June 2017
[9] Galambos J., et al. ACS J. Med. Chem., **2017**, 60 (6), pp 2470–2484