# TONY DAVIES COLUMN

# To add or not to add—that is the question for your reference spectra prediction databases

**Tony Davies,[a] Wolfgang Robien[b] and Jeff Seymour[c]**
[a]External Professor, University of Glamorgan, UK and Informatics Consultant
[b]Institute for Organic Chemistry, University of Vienna, Austria
[c]ModGraph Consultants, 1 Oakland View, Welwyn AL6 0RJ, UK

## Having a good idea

It is generally accepted that spectra prediction performance can be greatly enhanced if you add your own reference quality data to the database used for the predictions. This is mostly due to the fact that the reference databases that are generally available often represent only a small selection of data from across chemistry. What is probably worse, some prediction systems consist of data files that have been donated by one or two major industrial chemical manufacturers and as such are slanted towards the particular chemistry practiced within those companies—such as large organics or heterocyclic chemistry. They can therefore predict with excellent results within those fields but fall down badly when asked to predict spectra for structures whose chemistry is completely different.

So just buying bigger and bigger databases may not be the way forward to improving your prediction quality. Overall system performance is not really an issue anymore. Running Hierarchical Ordered description of the Substructure Environment (HOSE) code prediction data for medium sized structures C10, C20 against 100,000 in your reference database takes about 1–2s, using neural networks reduces this to a few ms. What is critical is to add reference data representative of your own chemistry. But, with limited resources available, which data should be added to really make a difference?

Recently, Wolfgang Robien and his colleagues were discussing this problem during a congress, with the example of a large pharmaceutical company who were looking to exploit their large collection of $^{13}$C-NMR spectra by making it available within their organisation for spectra prediction.

They came up with the idea of using the CSEARCH prediction engine[1] itself to locate the reference data sets of most value to improving the prediction quality of the knowledge base.

## Avoiding unnecessary laborious spectra/ structure assignments

Now if we look at the 100,000 spectra which represent the candidate data for adding to the knowledge base, we will usually find that most of the data sets have not been assigned. This means that there is no electronic assignment of the spectra features to the registered chemical structures. Were they all available as assigned data it would be possible to simply add them to the prediction set—and monitor statistically the prediction performance of the system.

But this is not how the real world works, so it is important to decide how to deploy limited resources to the best effect. The key question which needs to be answered is how do I select which of the candidate spectra and structures should be added to the database in order to maximise the benefit? In other words how much assignment work do I need to carry out to achieve the best return on my investment?

## HOSE codes

Now to understand how their solution works you need to understand something about how NMR prediction software does its job. Back in 1978 Wolfgang Bremser, then of BASF, published a paper "HOSE—A Novel Substructure Code" describing an encoding system for chemical structures which was ideally suited for use as a descriptor in NMR spectra prediction—the HOSE code.[2] HOSE codes describe the structural neighbours of the particular atom of interest, which in NMR essentially identifies those atoms within the molecule influencing the chemical shift of that atom. The neighbouring atoms are described in spheres—the nearest neighbours being in the first sphere. Typically, prediction databases only went out as far as five (or six) (CSEARCH, NMRPREDICT and NMRBENEFIT use up to five) spheres. Figure 1 shows how much of the testosterone molecule would be encoded for a central atom of interest.

In the reference database to be used as the knowledge-base for the predictions, the $^{13}$C NMR reference chemical shifts are stored along with HOSE codes
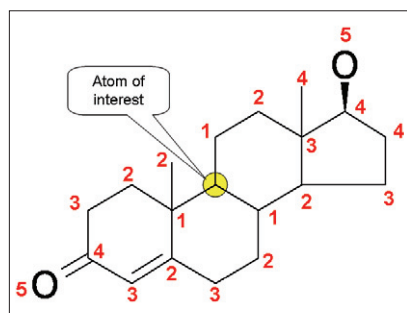
# TONY DAVIES COLUMN



**Figure 1.** HOSE codes describe nearest-neighbour relationships to atoms of interest. In this figure the atoms are labelled with the sphere number they have in relation to a particular atom of interest. The larger the number of a particular atom, the further away that atom lies from the atom whose NMR chemical is being predicted and probably the weaker the influence of that neighbour to the chemical shift.

for all spheres from one to five. A database with 100,000 reference spectra and their accompanying structures will yield (usually) more than 1 million structure codes. This is the knowledge base which the $^{13}$C-NMR prediction algorithms draw upon to calculate new chemical shift values for our chemical structures.

## Using HOSE codes for prediction

How does the prediction work? For our atom of interest, the HOSE code for the first sphere is generated and searched for in the reference database. If it is found then the HOSE code going out to the second sphere is searched for. Normally this continues, if possible, until the fifth sphere has been reached.

The prediction application then averages the chemical shift values it has mined from the reference database for the largest number of spheres and presents this as predicted value.

Obviously various software packages have numerous refinements on how the data is treated and the predictions weighted, but the basic principles are essentially the same.

For $^{13}$C-NMR, the predictions can be extremely reliable as the HOSE structure encoding methodology is exceptionally close to the real physical effects acting to change an atoms NMR shift value. Figure 2 shows an example of one such prediction package displaying the number of spheres used for the various chemical shift predictions so that the user can make their own evaluation about the quality of the predicted values.

## How best to improve predictions

As you can see, if your reference database contains different chemistries you will have little or no chance of locating
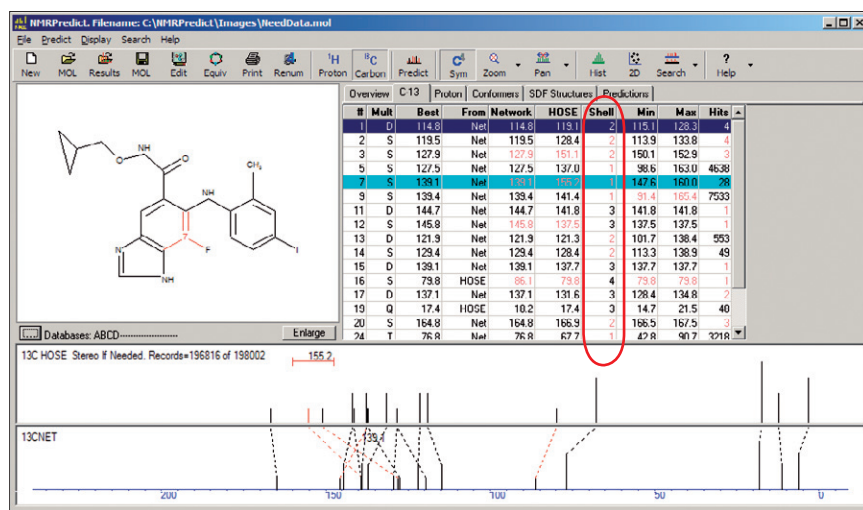
**Figure 2.** A screenshot from a [13]C-NMR prediction software package called NMRPredict showing (circled in red) how many spheres, or shells in this case, found in the reference database have contributed to the predictions.

reference shift values with the three, four or five spheres identical to the atom to be predicted. So now you know why it is better to have more, similar, reference data.

Robien's idea was beautiful in its simplicity and essentially turned the prediction route on its head. To identify which of the new data sets it would be worthwhile adding to the reference database he took the chemical structure of each of the new reference data sets and searched the database as if he was looking to predict its spectrum. By looking at the number of spheres found in the database it was possible to identify those structures which would have been well predicted by the current reference data and those structures which would have been poorly predicted, i.e. those structures whose encoded HOSE codes were not well represented in the database. By assessing which molecules would add the most new HOSE codes to the reference database you get a clear indication which reference data sets it would be worthwhile carrying out the time consuming assignment work on. More mathematically:

- you start with $N$ candidates for selection
- you do a prediction for all of them against your existing database holding $X$ entries

- you select the structure which is represented worst by your reference collection
- you "virtually" add this structure to your reference collection (has now $X + 1$ references) in order to benefit already from this dataset for the next prediction
- now you have $(N - 1)$ candidates and go back to step #1 until all candidates have been processed.

You end up with a detailed list in which sequence your internal data have to be assigned in order to get most benefit for future predictions.

So the end result is you mine your current reference database against your available additional reference data to come up with a sound basis for deciding how much effort will produce the best return on investment on improving NMR structure predictions.

For more information on NMRBenefit see Reference 3.
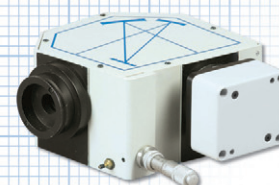
## References

1. H. Kalchhauser and W. Robien, *J. Chem. Inform. Comput. Sci.* **25**, 103–108 (1985).
2. W. Bremser, *Anal. Chim. Acta* **103**, 355–365 (1978).
3. http://www.modgraph.co.uk/product_nmr_benefit.htm