# The Benefits of Spectral Databasing within Chemistry and Natural Products

## Matthew Reid, PMAS

syngenta

# Introduction

- **Syngenta Jealott's Hill**

- **Structural Chemistry Group at Syngenta Jealott's Hill**

  - **Who We Are and What We Do**

  - **Spectroscopic Capabilities**

  - **Data Archiving and Prediction Training**

  - **What for?**

  - **What Software We Use and How We Do It**

  - **Future Developments**
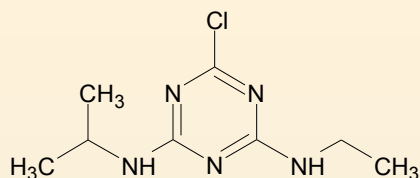
**syngenta**

# Syngenta Jealotts Hill



Jealott's Hill is our largest site for R&D and product support with ca.800 employees

Jealott's Hill Research Centre in the UK is our centre for Crop Protection Discovery, Bioscience, Weed Control Research, Seeds Research, bio-performance enhancement and Product Safety Research.
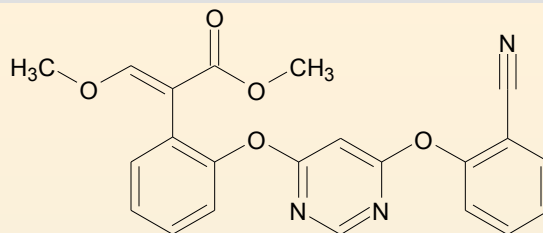
Key activities on site include research into discovery of new active ingredients (AIs), new formulation technologies to develop products from existing AIs and technical support of our product range.
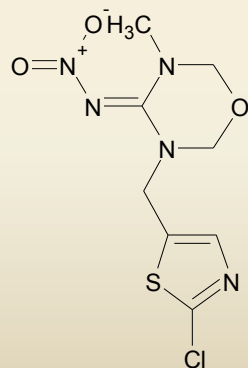
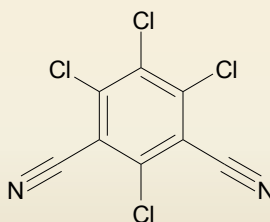**syngenta**

# Syngenta Jealotts Hill
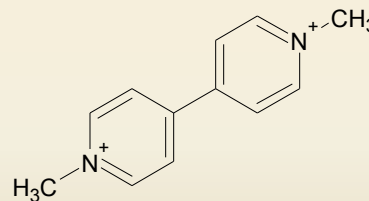
## Key Products



Atrazine - Herbicide

Azoxystrobin - Broad Spectrum Fungicide

Thiamethoxam - Broad Spectrum Insecticide

Chlorothalonil - Broad Spectrum Fungicide

Paraquat - Broad Spectrum Herbicide

syngenta

# Structural Chemistry Group at JH

## Who We Are and What We Do

- **The group is part of Product Metabolism and Analytical Science**
- **Provide qualitative and quantitative services and technical expertise to global projects in Research, Development and beyond.**

### NMR Spectroscopy

Structure elucidation for Research Chemistry (including natural products, impurity ID and physical chemistry),

Solid-state NMR.

Process studies, mode of action.

Competitor product analysis

Formulation studies, product safety, biochemistry (including protein NMR), quantitative NMR and Stage 1 metabolism

### Mass Spectrometry

Open access LC and GC-MS support for Research Chemistry, Formulation and Process Science

Structure elucidation studies for Research Chemistry, Formulation and Process Science

Structure elucidation for Stage 1 and Stage 2 metabolism

syngenta

# Structural Chemistry Group at JH

- **Spectroscopic Capabilities**



**Varian Inova 400MHz**

- Routine Open Access instrument

**Bruker AVANCE III 500MHz**

- Research instrument - 1D/2D experiments
- Primarily used for structure elucidation

**Varian Inova 600MHz**

- Research instrument – high sensitivity for $^1H$ and $^{13}C$
- Primarily used for structure elucidation of low level samples



**2 x Waters ZQ Open Access LC-MS systems, Shimadzu GC-MS, Jeol GC-MS**

- **Routine Open Access instruments**

**LTQ Velos**

- **Accurate mass instrument used for structure elucidation**
- **Applications within Chemistry, Process Science, Metabolism, Formulation**

**LTQ Orbitrap XL**

- **Accurate mass instrument used for metabolite ID**

**syngenta**

# Structural Chemistry Group at JH

- **Spectroscopic Capabilities**

**Sample Throughput – 2011**

- NMR

  - Full Structure Elucidations – 349

  - Open Access NMR – 20639

- MS

  - Open Access MS – 38152

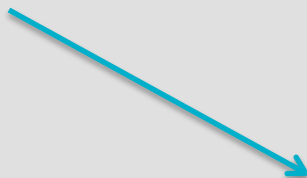**Sample Throughput – 2012 Jan - Apr**

- NMR

  - Full Structure Elucidations – 316

  - Open Access NMR – 8930

  - Batch Analysis - 10166

- MS

  - Full Analysis – 62

  - Open Access MS – 11587

  - Batch Analysis - 11782

syngenta

# Data Archiving

The raw data acquired on our analytical instruments is archived onto a server and can be accessed by chemists and analysts for reference.



**ADA Storage Server**

This process is automatic and the system can be searched using a number of variables.

syngenta

# Data Archiving

**•From a Spectroscopists/Analysts Point of View….**

- A raw data or paper archive isn't enough – we need something more specific:-

  - A spectral database that contains only what is useful and relevant

  - Something that's searchable by variables I want to search by

  - A database that is quick, powerful and can be customised

  - Something that stops duplication – very important in not wasting time

  - A system that can accept spectral data in multiple vendor formats

**syngenta**

# Data Archiving

- **Why do I need this?**



**Structure elucidation takes time - everything we've done in the past is useful. A spectral database gives us valuable legacy data from a range compounds and projects:-**

- Have we seen this compound before? If so where and when.
- What does the spectral information look like – is it similar to what we are currently working on?
- Are my raw materials the right compounds and of sufficient purity?
- Are the polymers present in my formulation similar to those in our database?

syngenta

# Data Archiving

- **What do we want to put in our database?**

  - **As much data from Chemistry projects as possible**

    - **The more assigned compounds we have the more it can aid in structure elucidation of compounds from the same projects**



**$^1$H, $^{13}$C, COSY, HSQC, HMBC, NOESY, LR or HRMS**

**$^{15}$N has become very important!!!**

    - **For natural products in particular, a comprehensive record should be kept, especially if the compound is novel**

syngenta

# Data Archiving

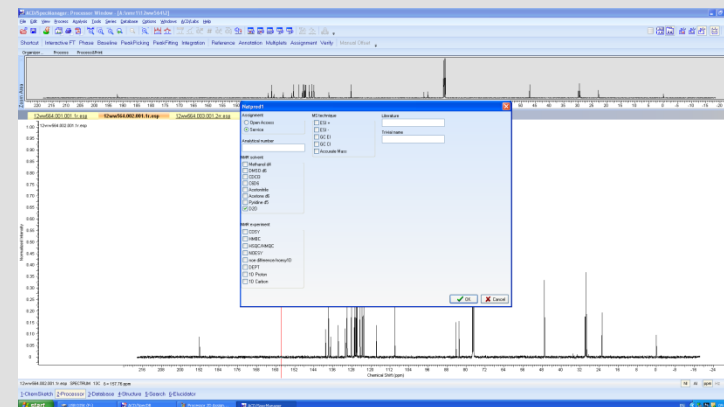- **The software of choice…ACD/Labs SpecDB Enterprise**



**Chemistry and natural product database with fully assigned NMR and MS data – Over 300 records.**

**Includes 1D and 2D data for 1H, 13C and 15N.**
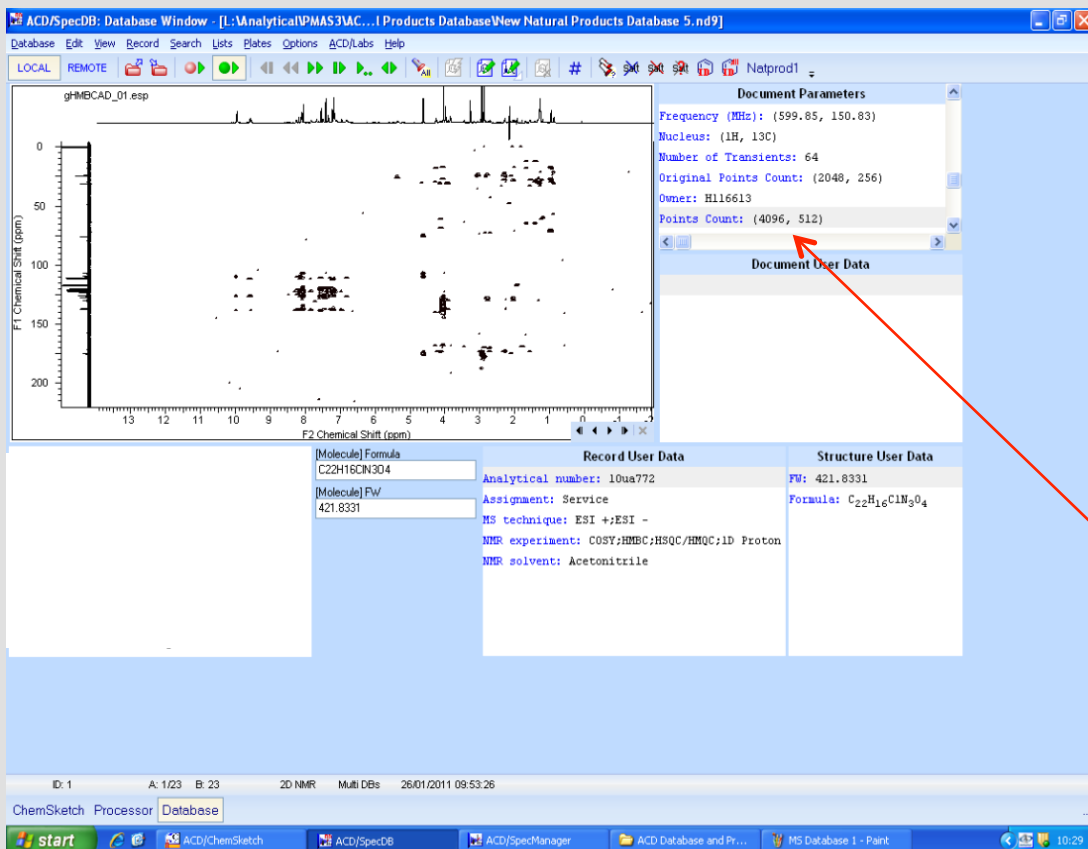
**Polymer and surfactant database – 209 compounds**

# Data Archiving

- **From Sample to Spectrum to Database**

syngenta

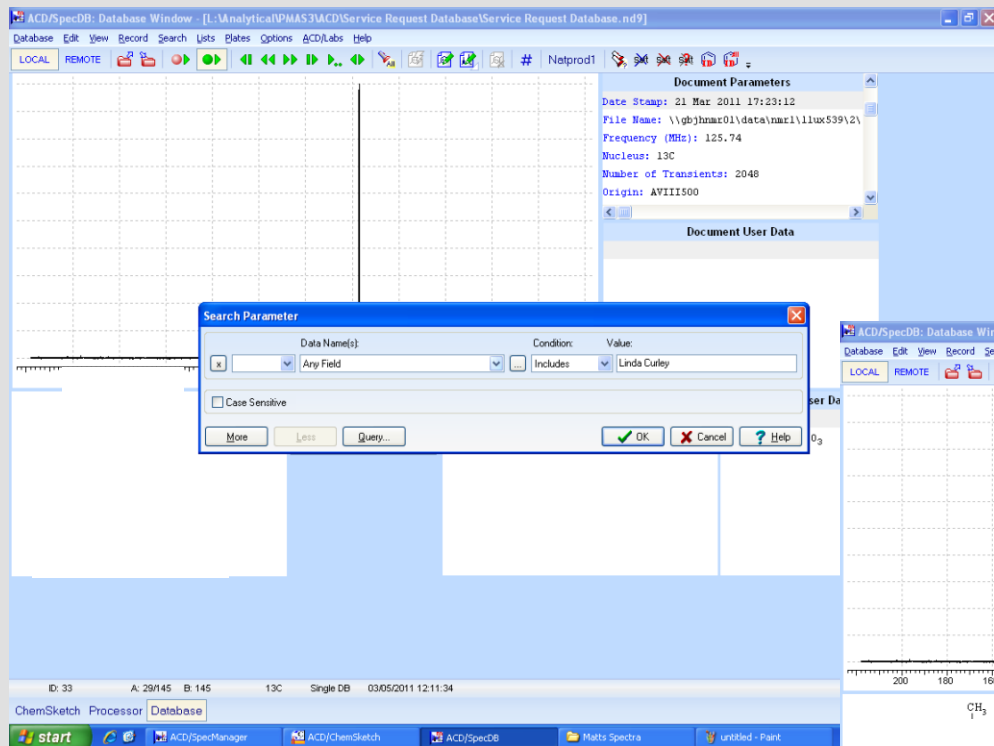# Data Archiving

- ## The Database



The format of the database is fully customisable. You only see what you want to see.
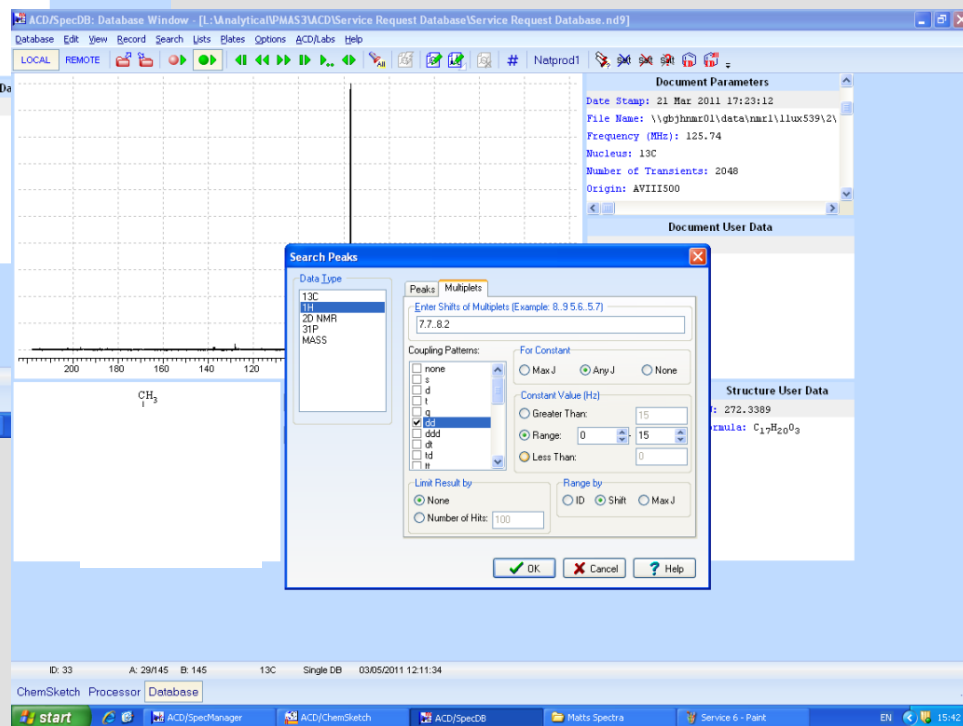
This customisation can be done at database creation or later on.

All the experimental/user/ record data is taken from the spectrometer automatically.
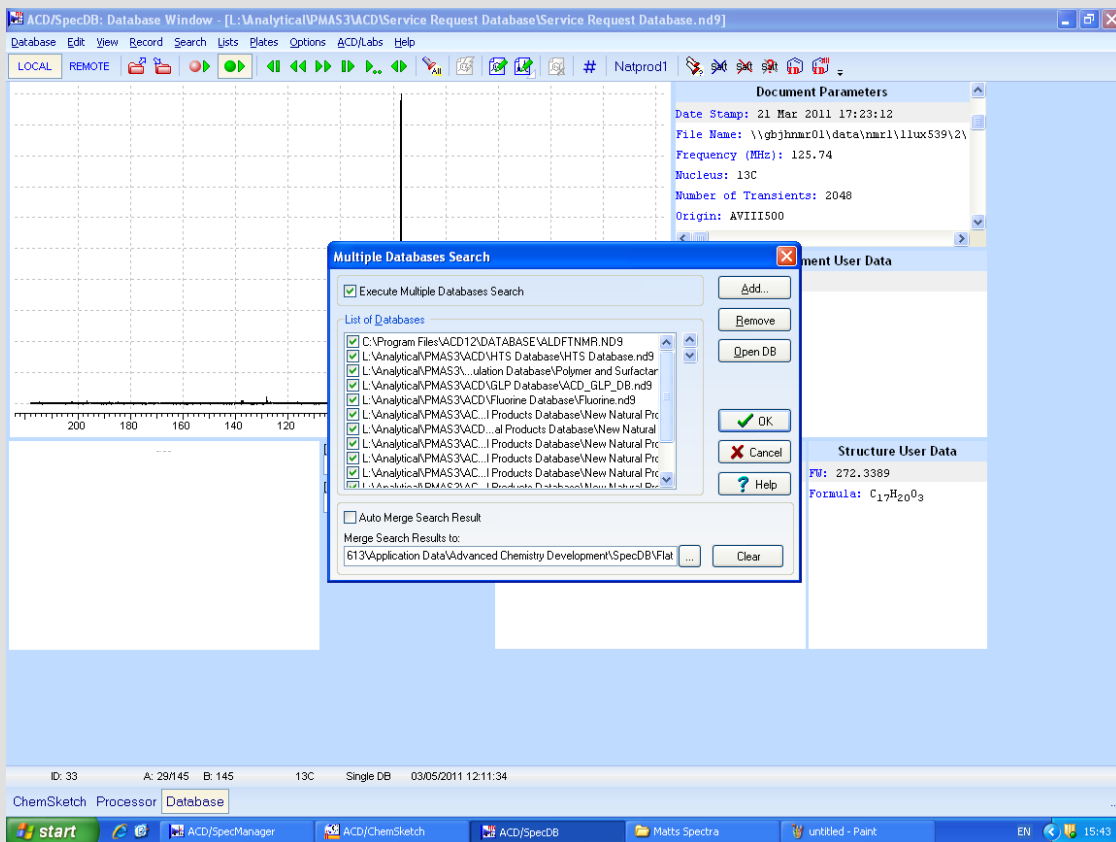
syngenta

# Data Archiving



The searching options within the database are too numerous to mention. Examples are to search by username or by NMR chemical shift and coupling constant

# Data Archiving

- ## The Database



Saving time…search multiple databases at once

Search for a structure, m/z, NMR chemical shift a range of database

Select the database you wish to search
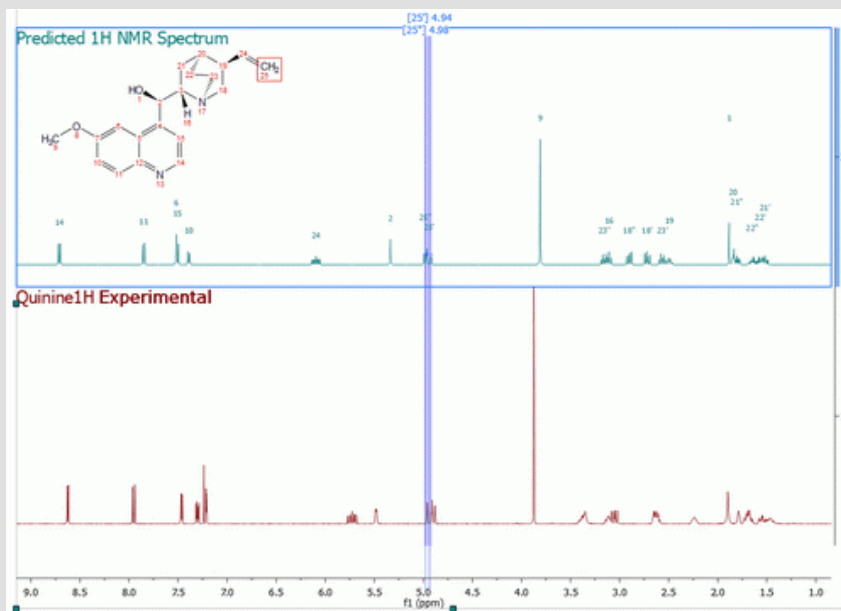
syngenta

# Data Archiving

•**Conclusions**

    ❑**ACD/Labs database software allows us to archive all our data From MS/NMR/UV/IR instrumentation.**

    ❑**The software is vendor neutral so has no bias to any particular data format**

    ❑**The databases are searchable in a huge number of ways and each database can be presented as the user sees fit**

    ❑**The data is only as good as the user inputs – careful consideration should be given as to what is submitted**

**syngenta**

# NMR Prediction

**The ability to predict what the NMR spectrum of a small organic molecule may look like is a valuable tool for the organic chemist and analyst.**

**It allows the chemist to gain some idea as to whether they have synthesised the correct compound by comparison of experimental and calculated $^1$H and $^{13}$C spectra.**



**This can save the chemist valuable time - giving the confidence the right compound had been made - they can move on to the next stage of synthesis**

# NMR Prediction

## There are a number of key software vendors who offer NMR chemical shift prediction packages:

**ACD/Labs**  For both $^1$H and $^{13}$C NMR use a combination of HOSE code and neural works. Prediction training is possible with HOSE code, hence the      value mental chemical shift databases

**Modgraph**  $^1$H NMR uses additivity rules (all the substituents have a contribution to the shift of the atom in question) and partial atomic charges and steric effects.
$^{13}$C NMR uses HOSE code/NN – dependent on the quality of shift database

**Perch**  Deterministic methods – high computational time to determine      ions

**Quantum Chemistry Predictors (e.g. Gaussian)**
Predicted from first principals – computationally very heavy

syngenta

# NMR Prediction

The software uses the HOSE code approach in its prediction (Hierarchical Organisation of Spherical Environments).
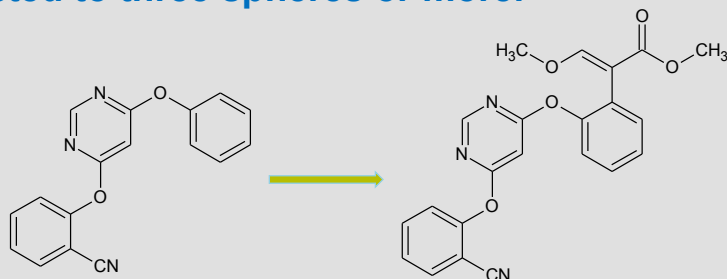


The HOSE starts at the atom whose shift is to be predicted.

- It looks one bond away from the atom and tries to find this environment in its database.

- If it is successful it moves two bonds away, tries again and so on until it comes across something not represented in the database, reaches the boundary of the molecule or the HOSE code limits.
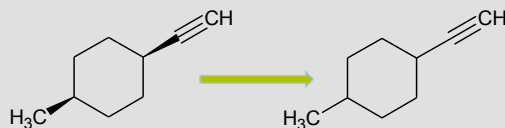
# NMR Prediction

**HOSE Code works well if:-**

•Structures are well represented in the reference collection.

•If atoms can be predicted to three spheres or more.

**HOSE Code doesn't work well if:-**

•Query structure is not well represented in the database

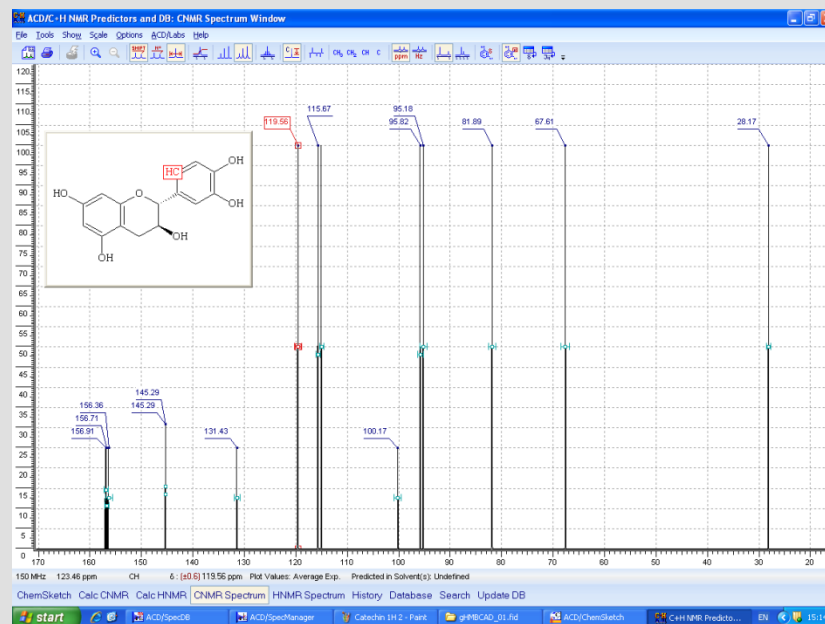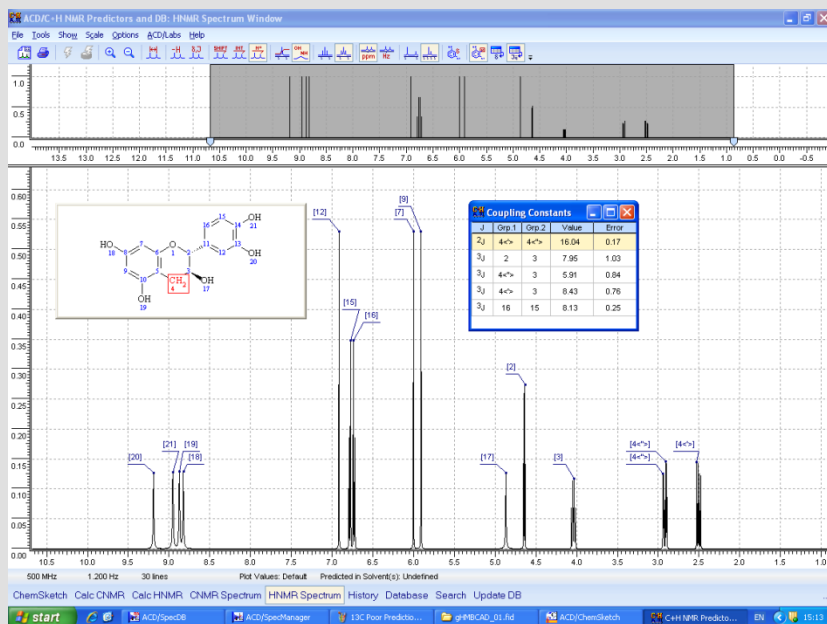•The atom in question can only be predicted to one or two spheres.

HOSE code approach exactly reproduces the contents of the reference database including every error within the reference database.

# NMR Prediction

## So how do we use the predictors?

**Very simply - just draw the structure and predict the $^1H/^{13}C$ spectrum. The structure atoms and spectrum peaks will be correlated and peaks will light up when atoms are hovered over with the mouse and vice versa.**
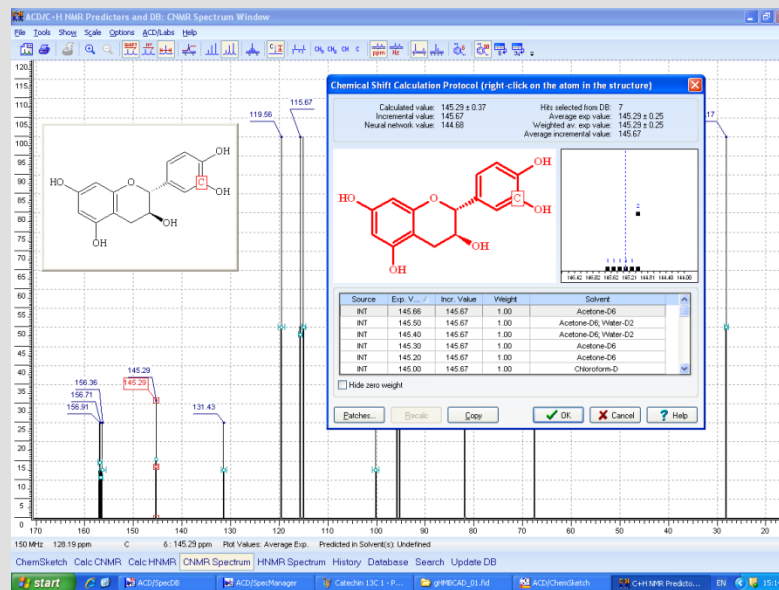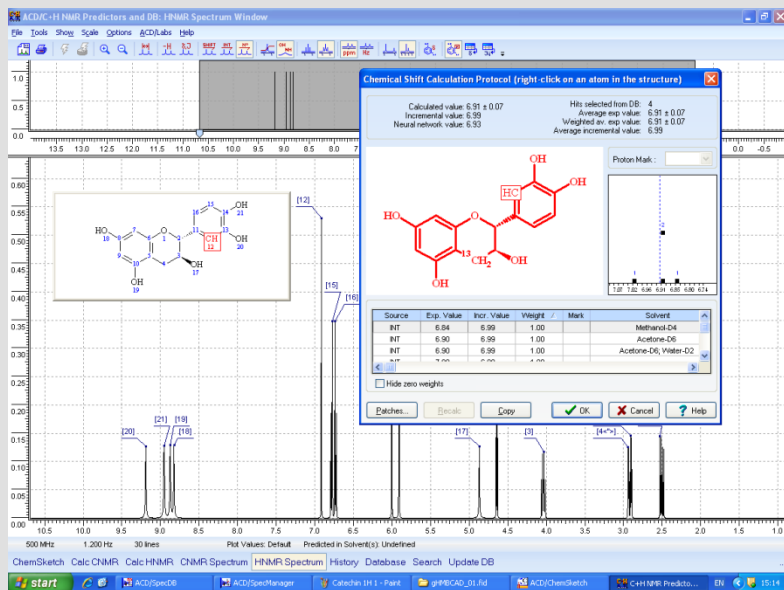


**Tables of chemical shifts, coupling constants and confidence limits can be displayed and couplings to other atoms/exchangeable protons can also be turned on or off.**

syngenta

# NMR Prediction

**By right clicking on an appropriate atom on the structure, its possible to see if this structure has good representation in the internal chemical shift database that the HOSE code uses.**

**Within the [1]H database are over 210,700 structures and within the [13]C database over 200,100 structures.**
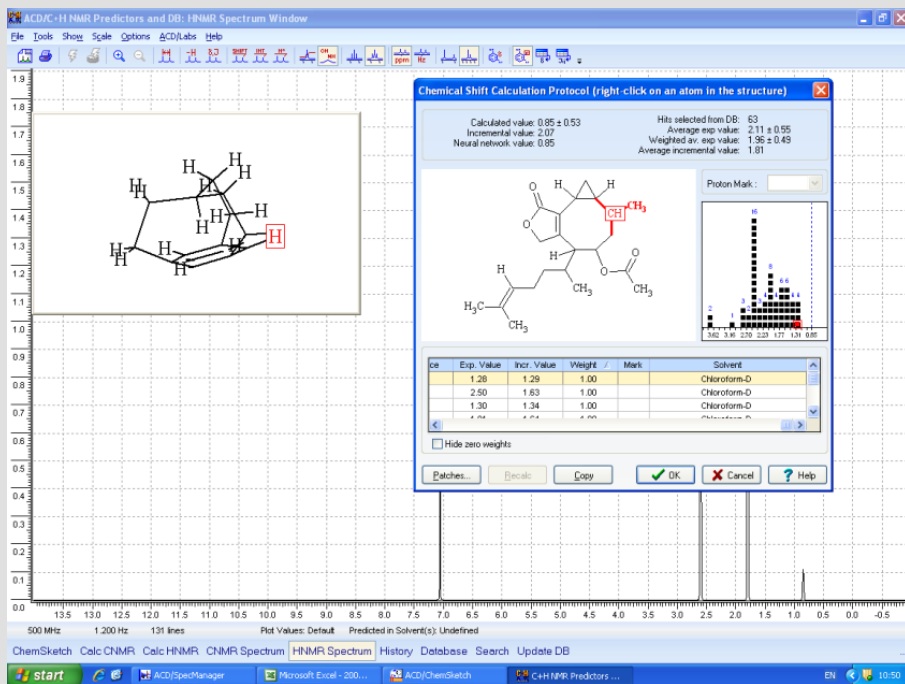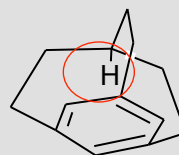


**The predictions can also be solvent specific so if you recorded your spectrum in $CDCl_3$ then use only database compound shift recorded in $CDCl_3$.**

syngenta

# NMR Prediction

## The problem with HOSE code....

Because we work with novel chemical entities, they are not represented well in the internal ACD databases. When this is the case the predictions can fall down.

A primary example is [10] cyclophane:





The proton highlighted resonates at ~-4ppm. The best ACD can do is ~0.7ppm for the proton highlighted.
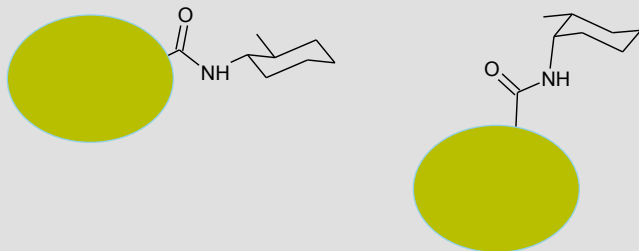
This is because there are no similar structures at all in the internal database.
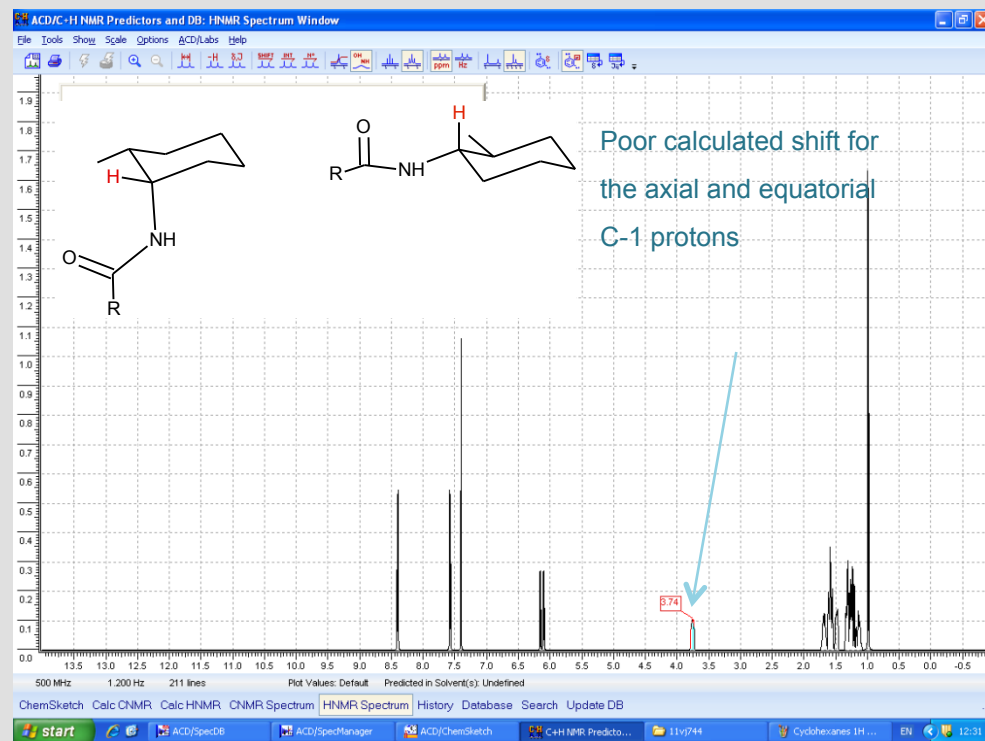
# NMR Prediction

## How do we overcome this problem?

To overcome the problem of poor representation in ACD's internal databases is to simply add the experimental data yourself and ask the software to use it. We **train the predictor**.

A good example is for the cyclohexanes below:-

The ${}^1$H prediction for these two molecules is poor and they cannot be differentiated.

By creating a ACD User HNMRDB we can add this experimental data in as we have recorded the ${}^1$H and ${}^{13}$C spectra ourselves.



Poor calculated shift for the axial and equatorial C-1 protons

# NMR Prediction

## How do we overcome this problem?

**We assign our experimental spectra and push these chemical shifts to the HNMRDB**



**The database contains all the assignments/coupling constants we measured**

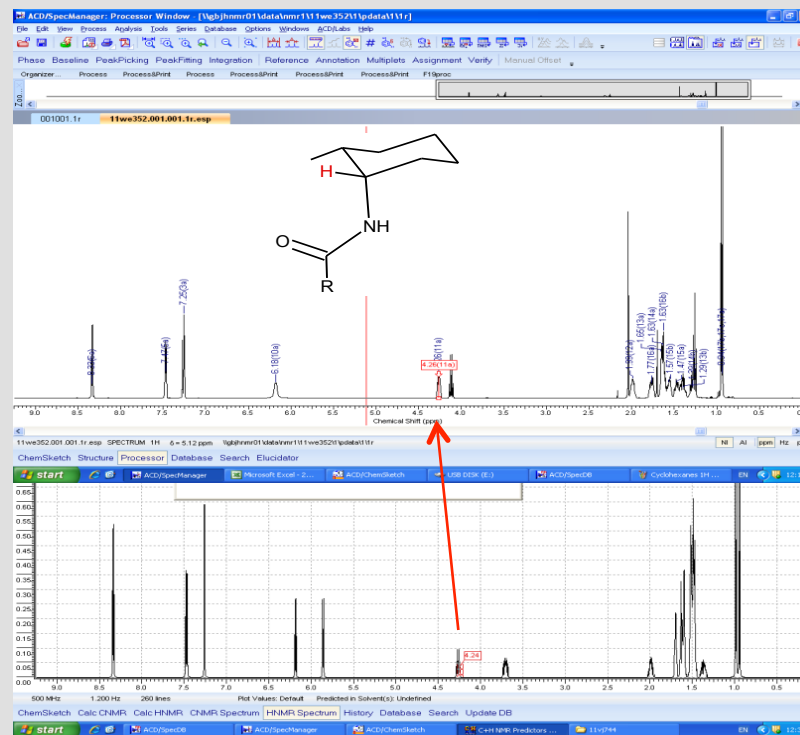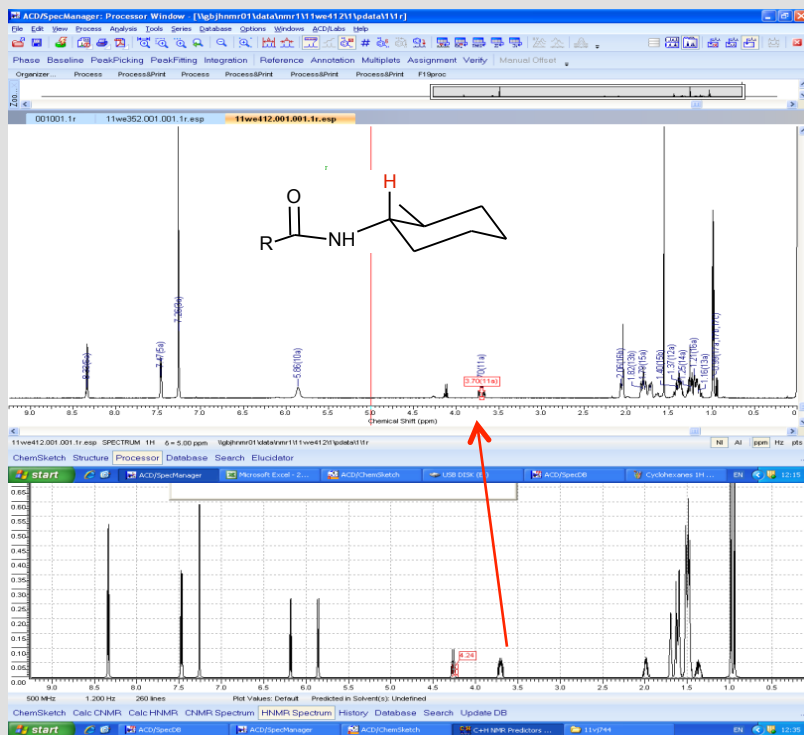**Same navigation and search functions as a spectral database described earlier.**

**We can also:**

•**Use multiple databases for prediction training.**

•**Edit shifts at any time**

•**Set the HOSE code depth to what we want – loose or tight**

syngenta

# NMR Prediction

## The result...

We predict again and now obtain a accurate prediction. We force the software to use our own chemical shift database - the predictions below should be a 100% match with the experimental data.
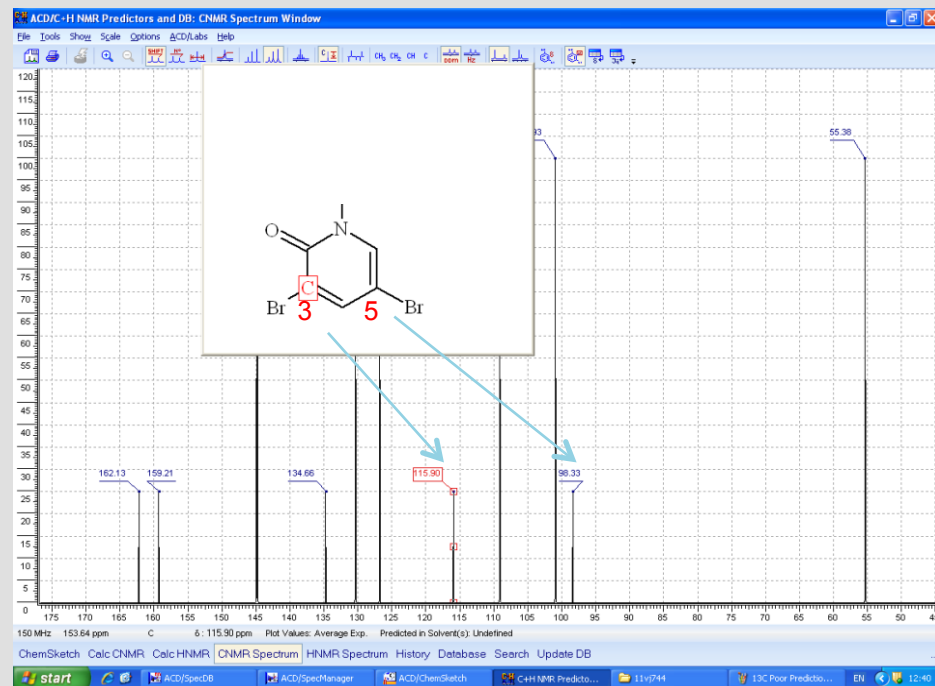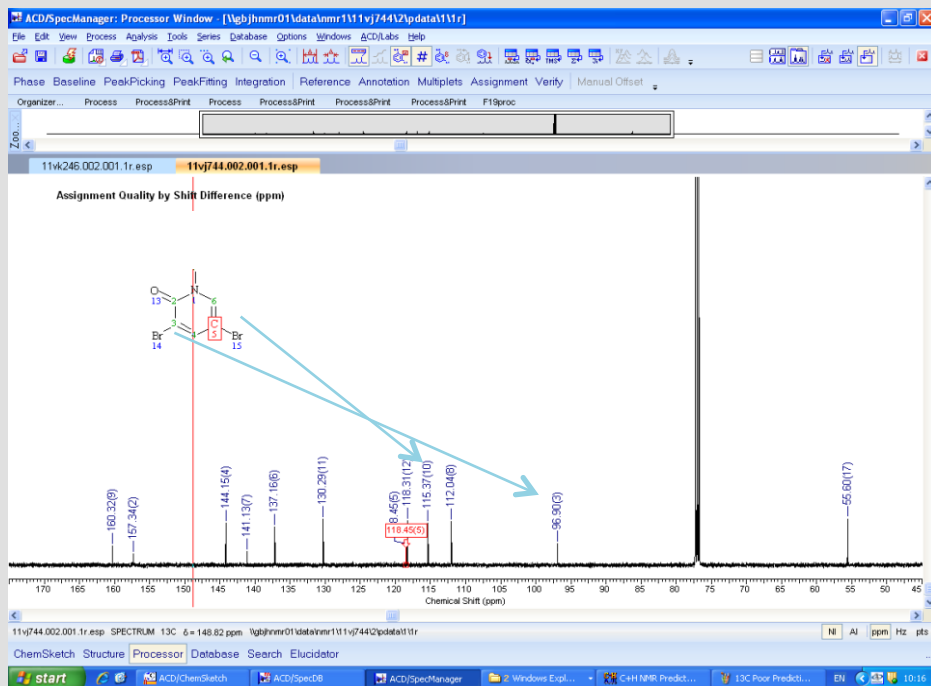
The value we get is further down the line when we predict spectra for similar compounds – we now have a representative compound in our user HNMRDB.

# NMR Prediction

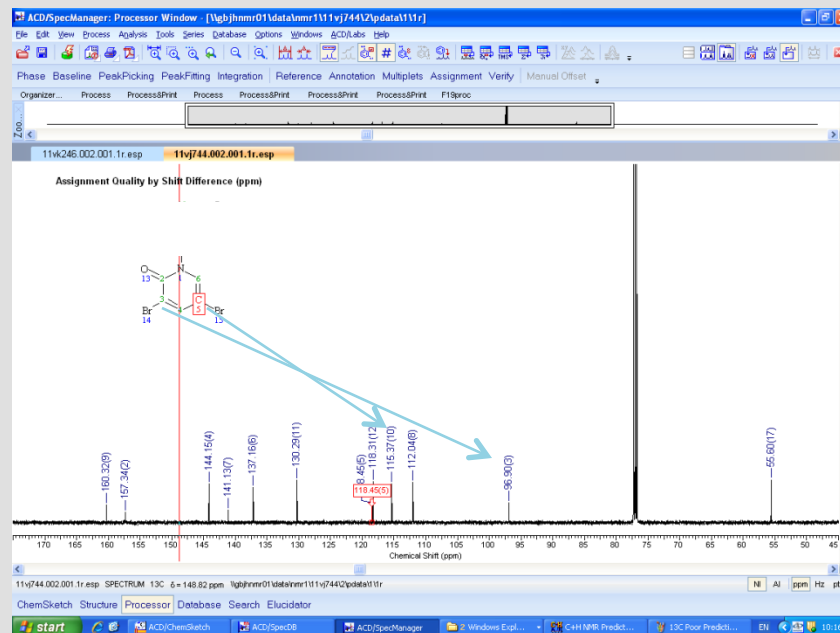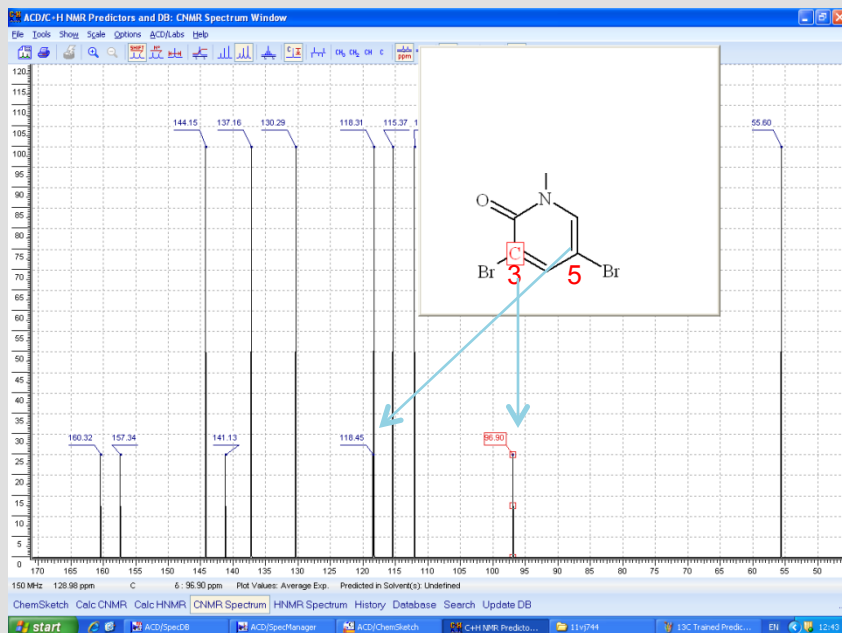## The result...

**Same principal for any nuclei – in this case $^{13}$C. Atoms 3 and 5 are predicted completely the wrong way around. In reality atom 3 is at much lower ppm than atom 5:-**

# NMR Prediction

## The result...

**With training, the predicted shifts are corrected and for subsequent compounds that contain the dibromopyridinone fragment the predictions will be a great deal better.**

# NMR Prediction

## Our Internal Databases:-

**We have multiple user prediction databases available, all created in house.**

Natural Products database contains ~140 compounds with $^1$H and $^{13}$C data

Chemistry database contains over 160 compounds

Natural product $^{13}$C database with over 6100 compounds

Coupling constant database with over 120 $^1$H-$^{13}$C coupling constants

**These databases are updated regularly. The desire is to start populating $^{19}$F and $^{31}$P databases**

syngenta

# NMR Prediction

## Conclusions

**Using HOSE code based ACD predictors has numerous advantages over other prediction programs:-**

- **Not computationally expensive**

- **They can incorporate neural network algorithms to compliment each other**

- **Multiple nuclei can be predicted quickly and easily**

- **The internal database contains a huge amount of chemical shift and coupling information**

- **The databases can be trained – this ability makes database based predictions very effective**

**syngenta**

# NMR Prediction

## Conclusions

**However, the disadvantages:**

- **Its very time consuming to get the data into databases. Full and accurate assignment of the experimental data must be performed**

- **The predictions without training are only as good as the data in the ACD internal databases.**
  - How do we know where this data came from?
  - Is it accurate – who recorded it, did they mis-assign the spectrum, etc?

syngenta