# PXRD with RAMAN SPECTROSCOPY, DSC and IR DATA

**Chris Gilmore, Gordon Barr and Gordon Cunningham**

# This document was presented at PPXRD - Pharmaceutical Powder X-ray Diffraction Symposium

## Sponsored by The International Centre for Diffraction Data

This presentation is provided by the International Centre for Diffraction Data in cooperation with the authors and presenters of the PPXRD symposia for the express purpose of educating the scientific community.

*All copyrights for the presentation are retained by the original authors.*

The ICDD has received permission from the authors to post this material on our website and make the material available for viewing. Usage is restricted for the purposes of education and scientific research.

PPXRD Website – www.icdd.com/ppxrd          ICDD Website - www.icdd.com

# Classifying Data

- **Statistics/data mining problem: put patterns (PXRD, spectra) into clusters where each cluster contains patterns which are most similar to each other.**

- **Not always a unique solution.**

- **Problems with:**

  - **Data quality.**

  - **Sample quality.**

  - **Data quantity.**

  - **Need for automation, and speed.**

# Correlation

- **Forget peaks.**
- **Match *every* data point.**
- **Use correlation coefficients:**
  - **Pearson correlation coefficient (parametric).**
  - **Spearman correlation coefficient (non-parametric).**

# Pearson correlation coefficient

$$r_{Pearson} = \frac{\displaystyle\sum_{\text{all points}} \left( x_i - \overline{x} \right)\left( y_i - \overline{y} \right)}{\sqrt{\displaystyle\sum_{\text{all points}} \left( x_i - \overline{x} \right)^2} \sqrt{\displaystyle\sum_{\text{all points}} \left( y_i - \overline{y} \right)^2}}$$

# Spearman correlation coefficient

$$r_{Spearman} = \frac{\sum\limits_{\text{all points}} \left( R_i - \overline{R} \right) \left( S_i - \overline{S} \right)}{\sqrt{\sum\limits_{\text{all points}} \left( R_i - \overline{R} \right)^2} \sqrt{\sum\limits_{\text{all points}} \left( S_i - \overline{S} \right)^2}}$$
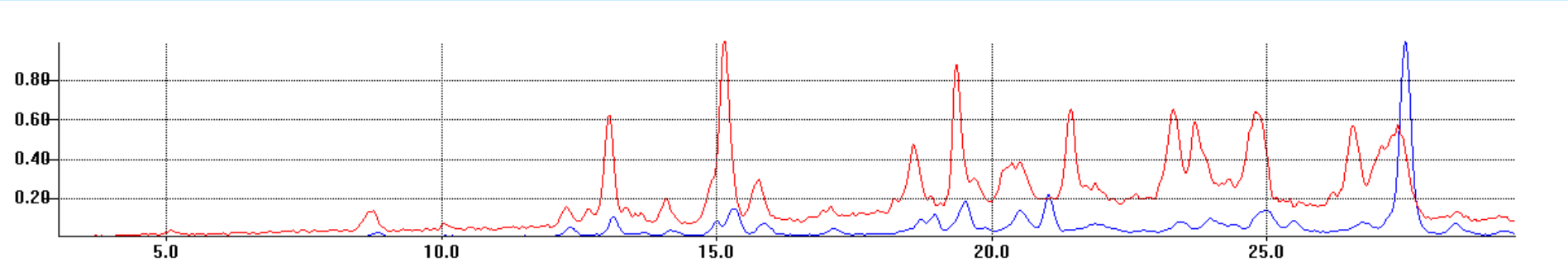
**Ranks**

# Combine them

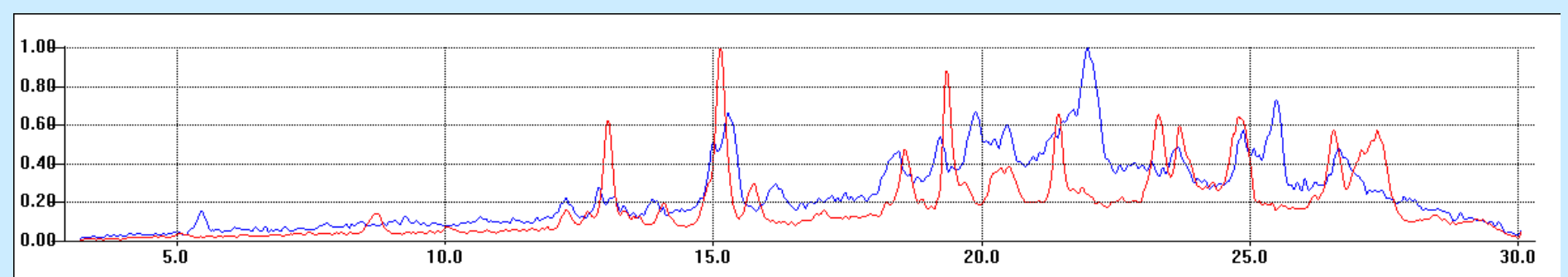$$r = w_1 r_{Pearson} + w_2 r_{Spearman}$$

**Use Fischer transforms:**

$$r = \tanh\left(\frac{\tanh^{-1}(r_{Pearson}) + \tanh^{-1}(r_{Spearman})}{2}\right)$$

# Correlation

## Mean correlation coefficient = 0.673



## Mean correlation coefficient = 0.771

# Correlation matrix

| Rank: | 01_s4.raw | 02_s3.raw | 03_s2.raw | 04_c1.raw | 05_c3.raw | l6_s4+3.raw | l7_s3+2.raw |
|---|---|---|---|---|---|---|---|
| 01_s4.raw | 1.000 | 0.853 | 0.810 | 0.615 | 0.590 | 0.753 | 0.679 |
| 02_s3.raw | 0.853 | 1.000 | 0.725 | 0.595 | 0.546 | 0.710 | 0.673 |
| 03_s2.raw | 0.810 | 0.725 | 1.000 | 0.715 | 0.623 | 0.707 | 0.674 |
| 04_c1.raw | 0.615 | 0.595 | 0.715 | 1.000 | 0.691 | 0.525 | 0.472 |
| 05_c3.raw | 0.590 | 0.546 | 0.623 | 0.691 | 1.000 | 0.531 | 0.476 |
| 06_s4+3.raw | 0.753 | 0.710 | 0.707 | 0.525 | 0.531 | 1.000 | 0.902 |
| 07_s3+2.raw | 0.679 | 0.673 | 0.674 | 0.472 | 0.476 | 0.902 | 1.000 |
| 08_s4+2.raw | 0.778 | 0.660 | 0.706 | 0.520 | 0.492 | 0.692 | 0.562 |
| 09_c1+3.raw | 0.509 | 0.496 | 0.554 | 0.719 | 0.855 | 0.441 | 0.413 |
| 10_s2+3+4.raw | 0.832 | 0.889 | 0.768 | 0.657 | 0.646 | 0.754 | 0.674 |
| 11_s2+c1.raw | 0.646 | 0.592 | 0.618 | 0.532 | 0.515 | 0.574 | 0.562 |
| 12_s3+c1.raw | 0.692 | 0.679 | 0.747 | 0.539 | 0.592 | 0.798 | 0.853 |
| 13_s4+c1.raw | 0.837 | 0.814 | 0.868 | 0.796 | 0.771 | 0.754 | 0.683 |
| 14_s2+c3.raw | 0.842 | 0.725 | 0.754 | 0.606 | 0.612 | 0.702 | 0.585 |
| 15_s3+c3.raw | 0.510 | 0.471 | 0.520 | 0.510 | 0.673 | 0.462 | 0.448 |
| 16_s4+c3.raw | 0.563 | 0.494 | 0.515 | 0.501 | 0.546 | 0.474 | 0.456 |

# Correlation, distance and Similarity

$$-1 \leq \rho \leq 1$$

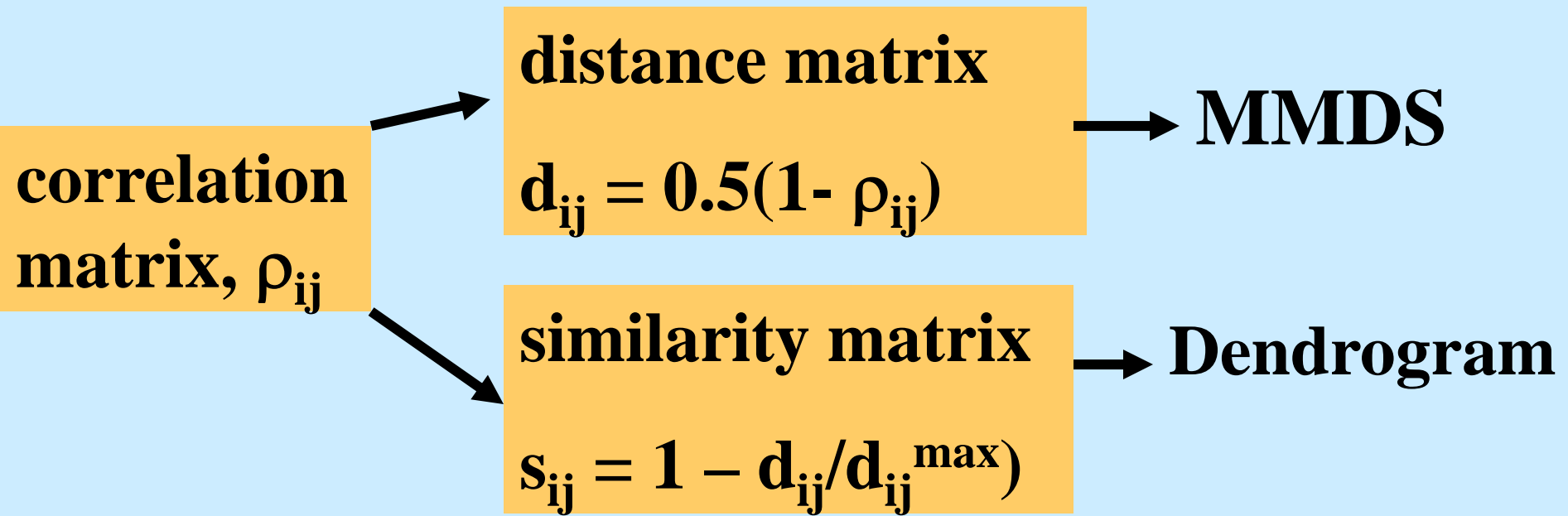**Larger** **the value of $\rho$ the closer the match.**

$$0 \leq d \leq 1$$

**Smaller** **the value of d the closer the match.**

$$0 \leq s \leq 1$$

**Larger** **the value of s the closer the match.**

# Correlation, distance and Similarity

**correlation matrix, $\rho_{ij}$**

→ **distance matrix**

$$d_{ij} = 0.5(1 - \rho_{ij})$$

→ **MMDS**

→ **similarity matrix**

$$s_{ij} = 1 - d_{ij}/d_{ij}^{max})$$

→ **Dendrogram**

# Example: Phase Transitions in Ammonium Nitrate

**Work with**

- **Michael Hermannn (Fraunhofer-Institut Chemische Technologie)**

- **Karsten Knorr (Bruker AXS, Karlsruhe)**

**See Hermann & Engel *Propellants, Explosives, Pyrotechnics* 22**, 143-147 (1997).

# Ammonium Nitrate

- **5 phases I – V**

- **Start with a mixture of IV + V that transforms to IV, II and I.**

Ammoniumnitrat; Pfinztal 15. März 2007, M. Herrmann

# Ammonium Nitrate

- **75 PXRD data sets.**

- **How to visualize and interpret these data?**

# Dendrogram
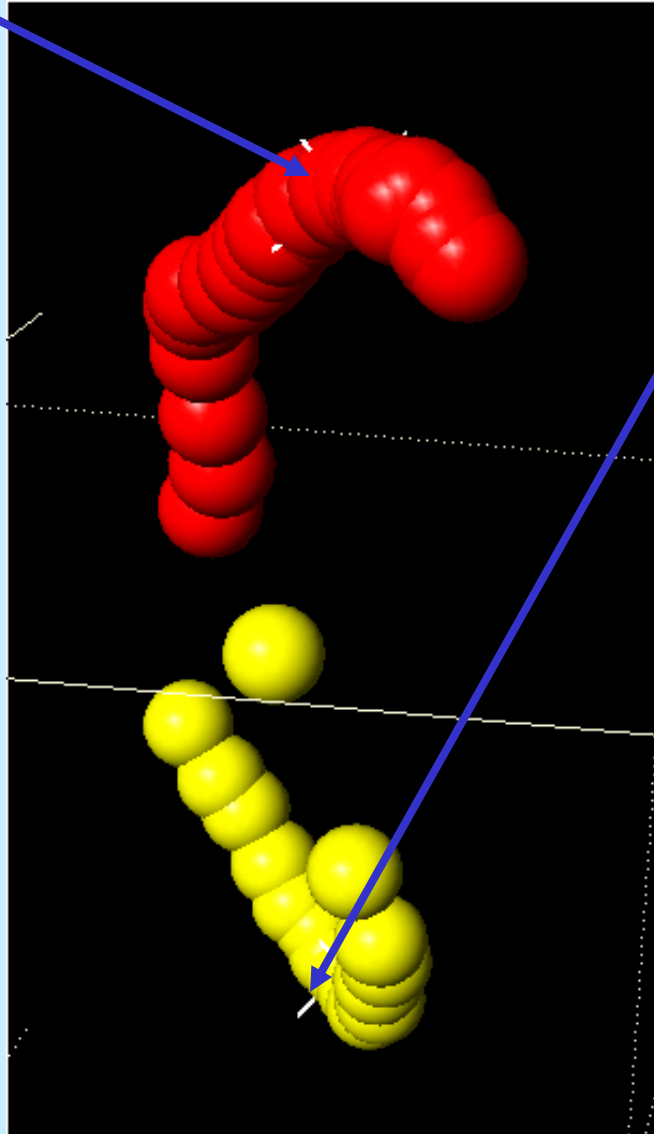
# Numbering

# Metric Multidimensional Scaling (MMDS)
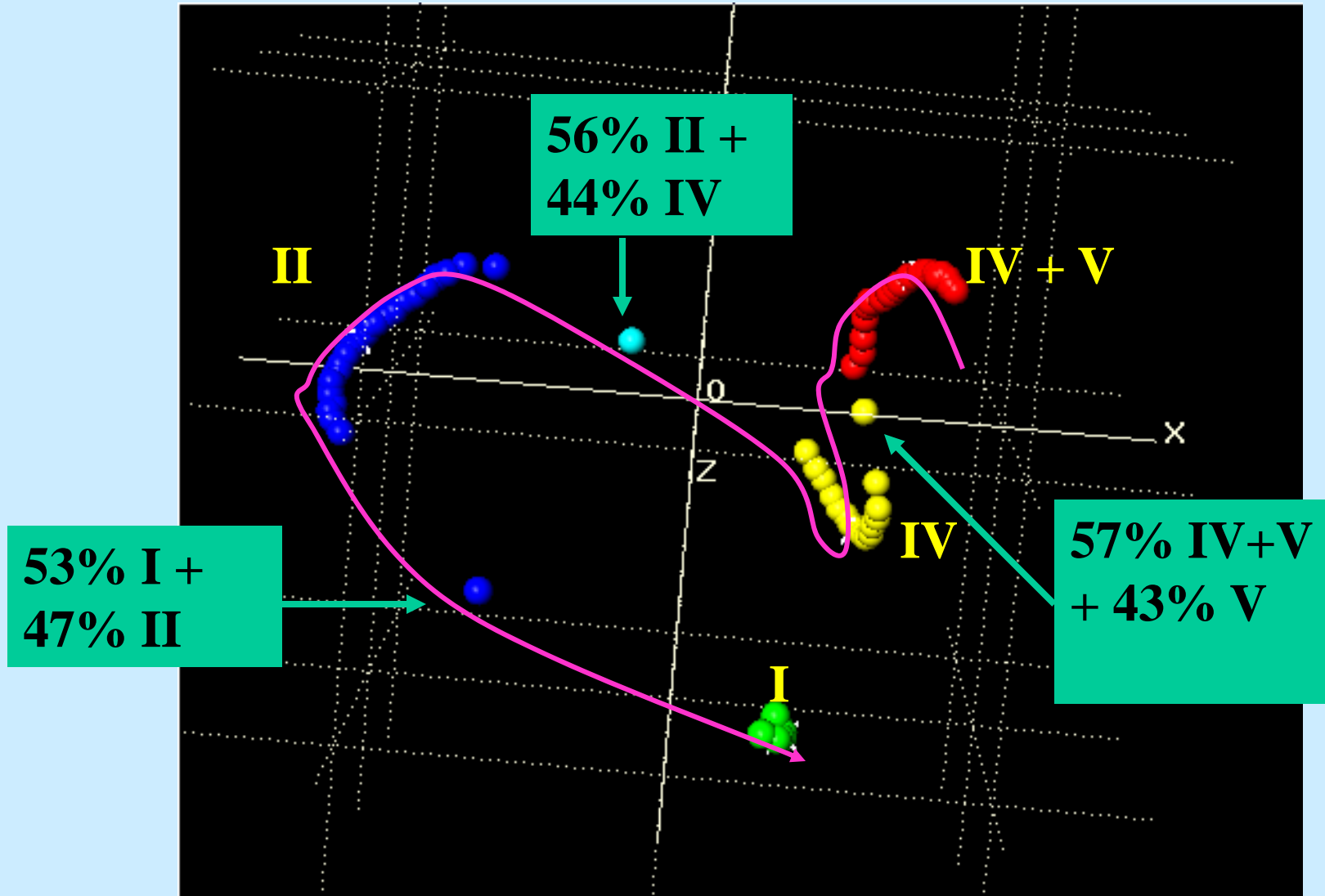
# Metric Multidimensional Scaling

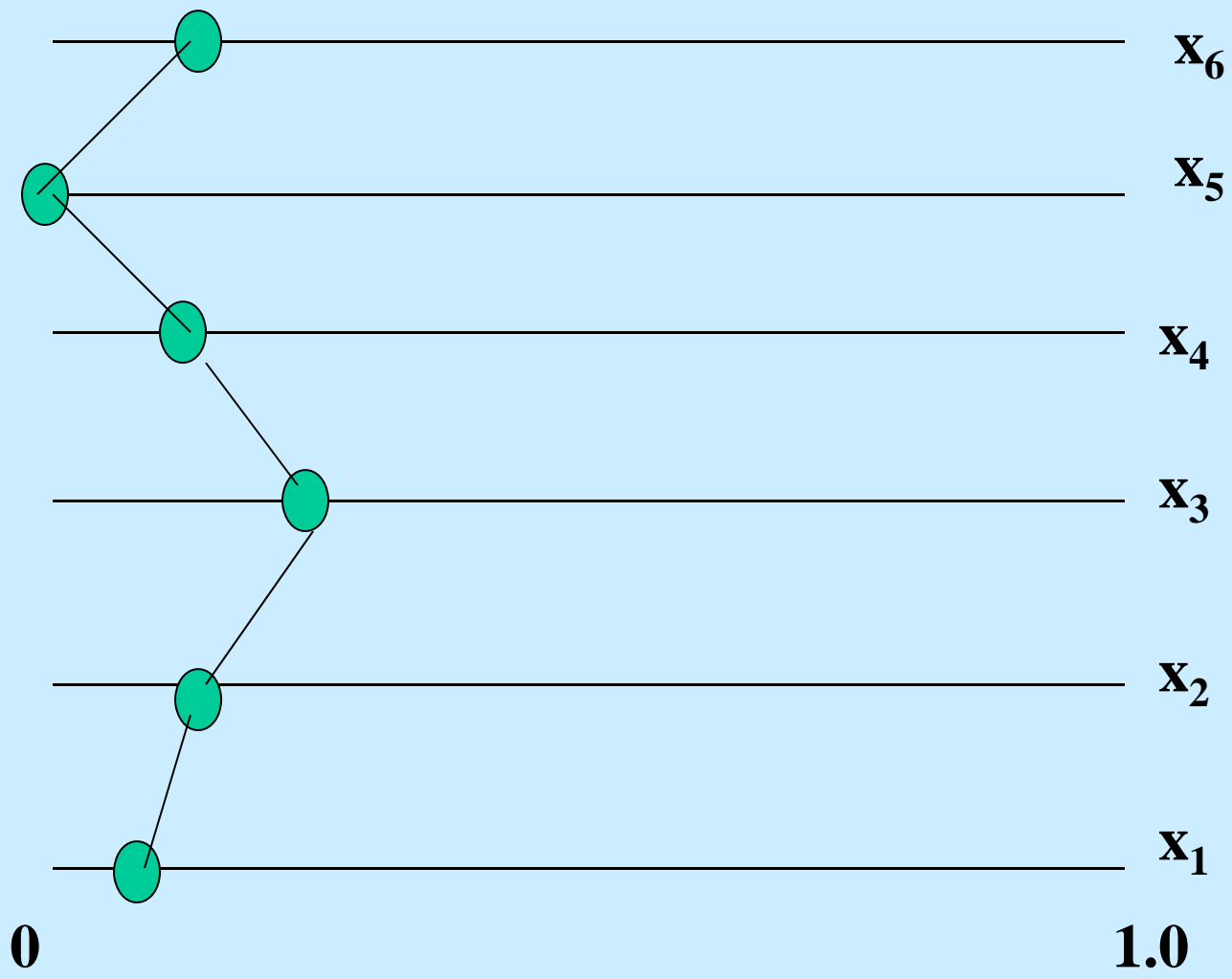**Most Representative Sample (MRS)**

# Simple Quantitative Analysis

- **We know the sequence IV + V $\rightarrow$ IV $\rightarrow$ II $\rightarrow$ I**

- **Take MRS of IV+V, IV, II, I and use in quantitative analysis.**
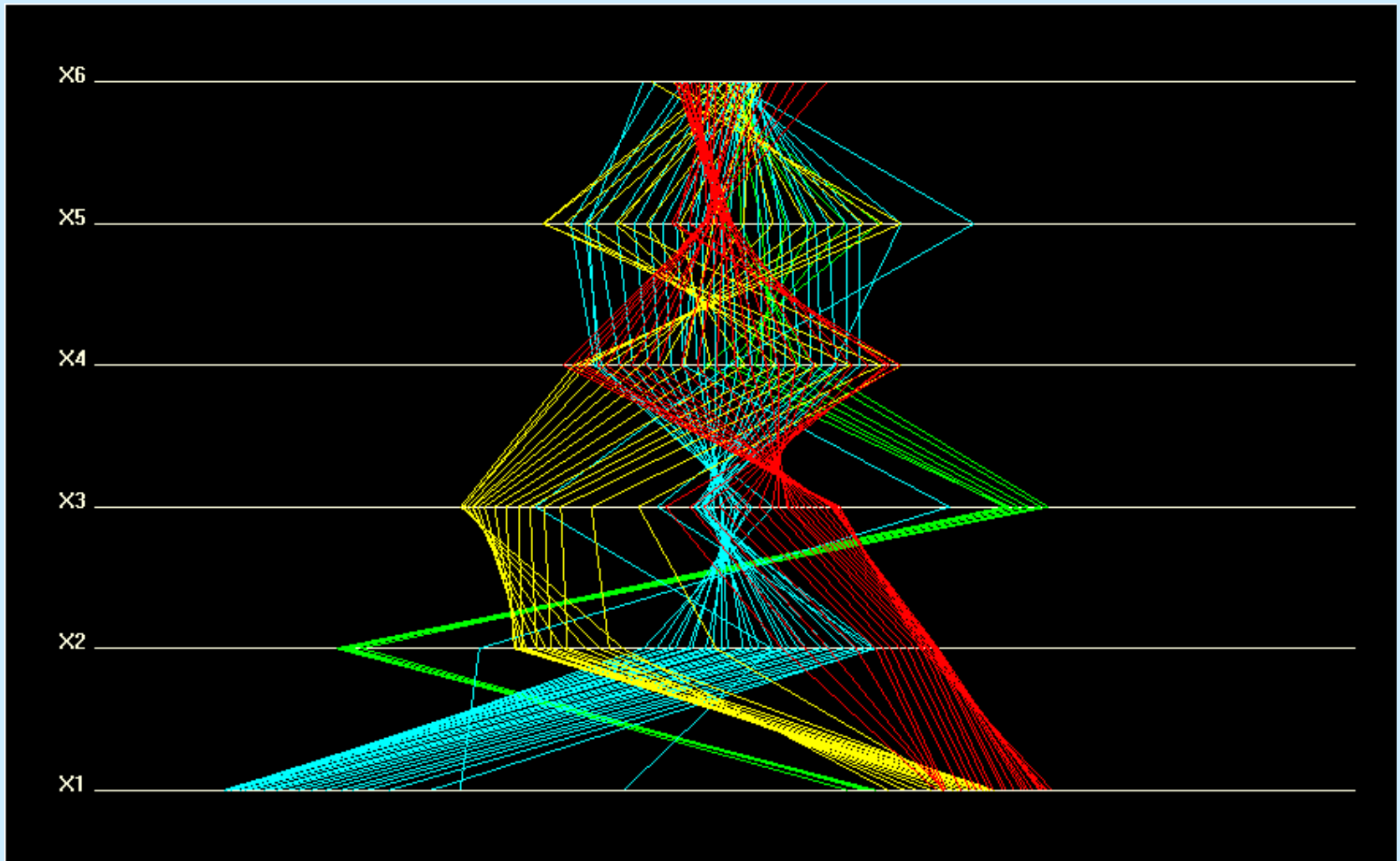
# Simple Quantitative Analysis

# Parallel Coordinate Plots

- **The MMDS (or PCA) calculation gives results in > 3 dimensions, and we merely select the 1$^{st}$ 3 (e.g. the first 3 eigenvectors for each eigenvalue)**

- **Consider the 6 dimensional coordinates**

**$(0.1, 0.2, 0.3, 0.2, 0.0, 0.2) = (x_1, x_2, x_3, x_4, x_5, x_6)$**

$x_6$

$x_5$

$x_4$

$x_3$

$x_2$

$x_1$

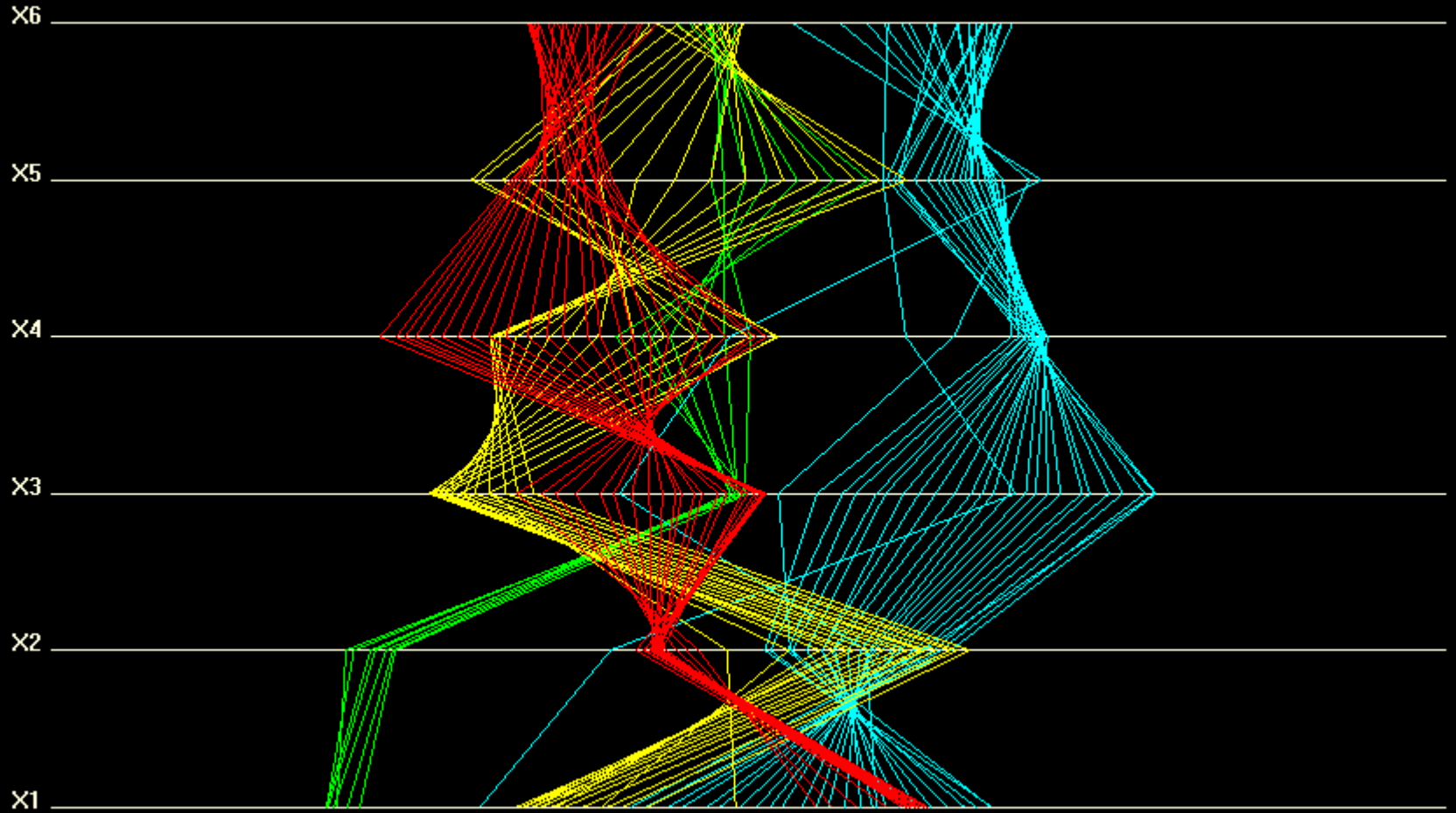0                                                                          1.0

# Other Visualization Tools: Parallel Coordinate Plots – More than 3 Dimensions

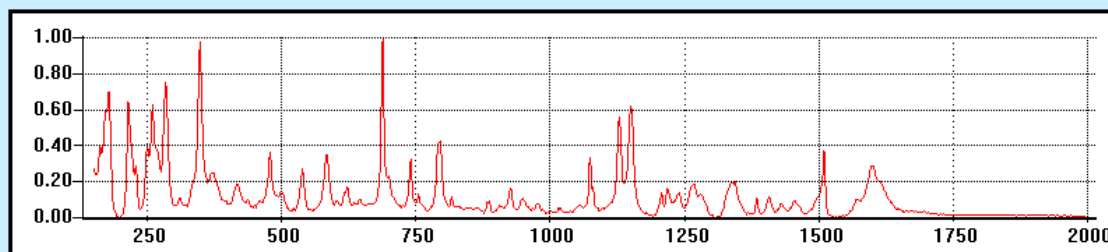# Another View

# Is PXRD the Gold Standard?

**Powder X-ray diffraction is usually considered the gold standard in high throughput studies designed to identify polymorphs, salts, co-crystals *etc.*, but other techniques such as Raman and IR spectroscopy, or differential scanning calorimetry (DSC) can have a major role to play also especially with poorly characterised samples.**
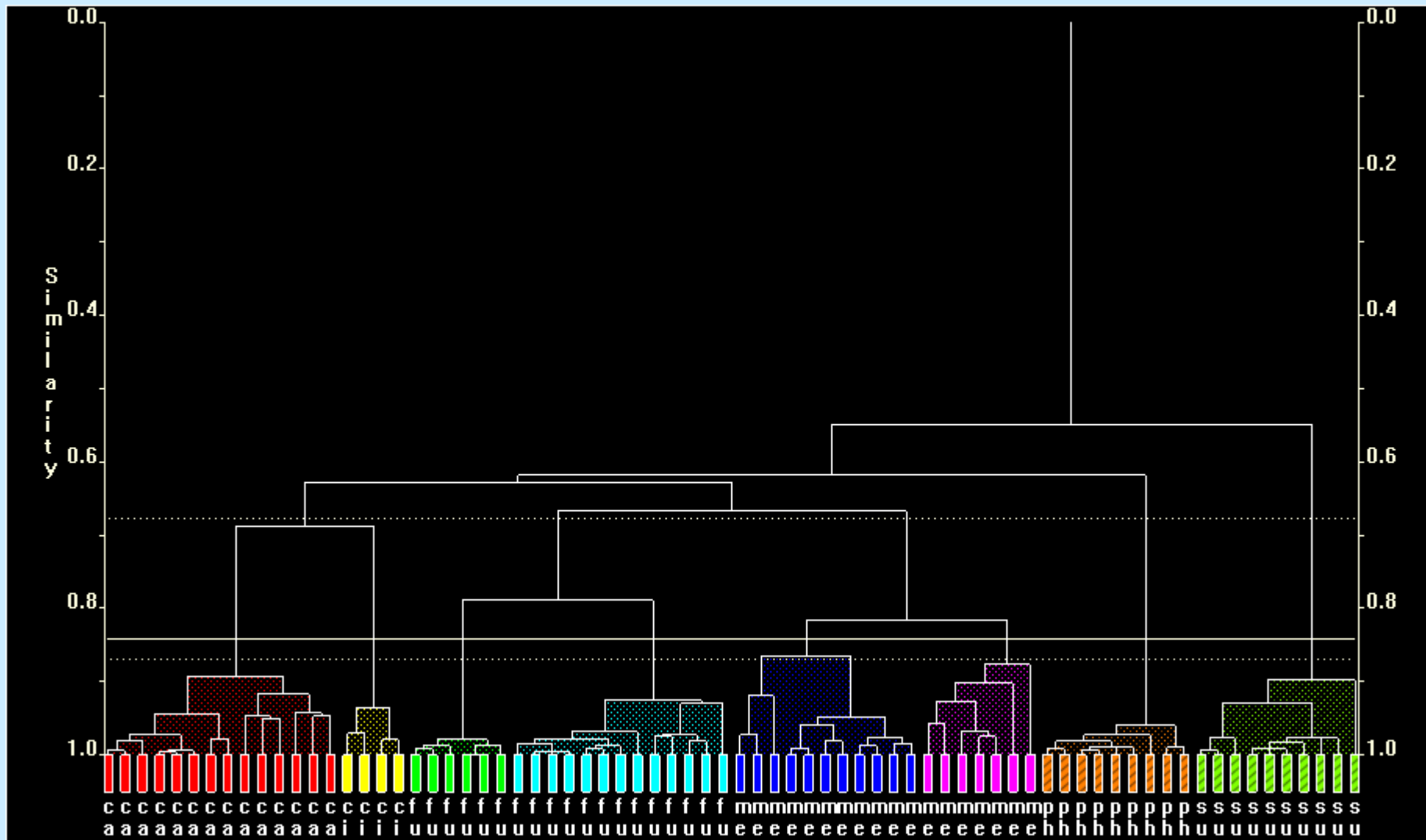
# Raman Data

- **Raman spectroscopy is well suited to screening: good quality spectra can be collected in a few minutes, and sample preparation is straightforward and flexible, although the resulting spectra are not always as distinct as the PXRD equivalent**

*High-throughput powder diffraction V: the use of Raman spectroscopy with and without X-ray powder diffraction data*
**Barr, Cunningham, Dong, Gilmore & Kojima** *J. Appl. Cryst.* **(2009). 42, 706–714**
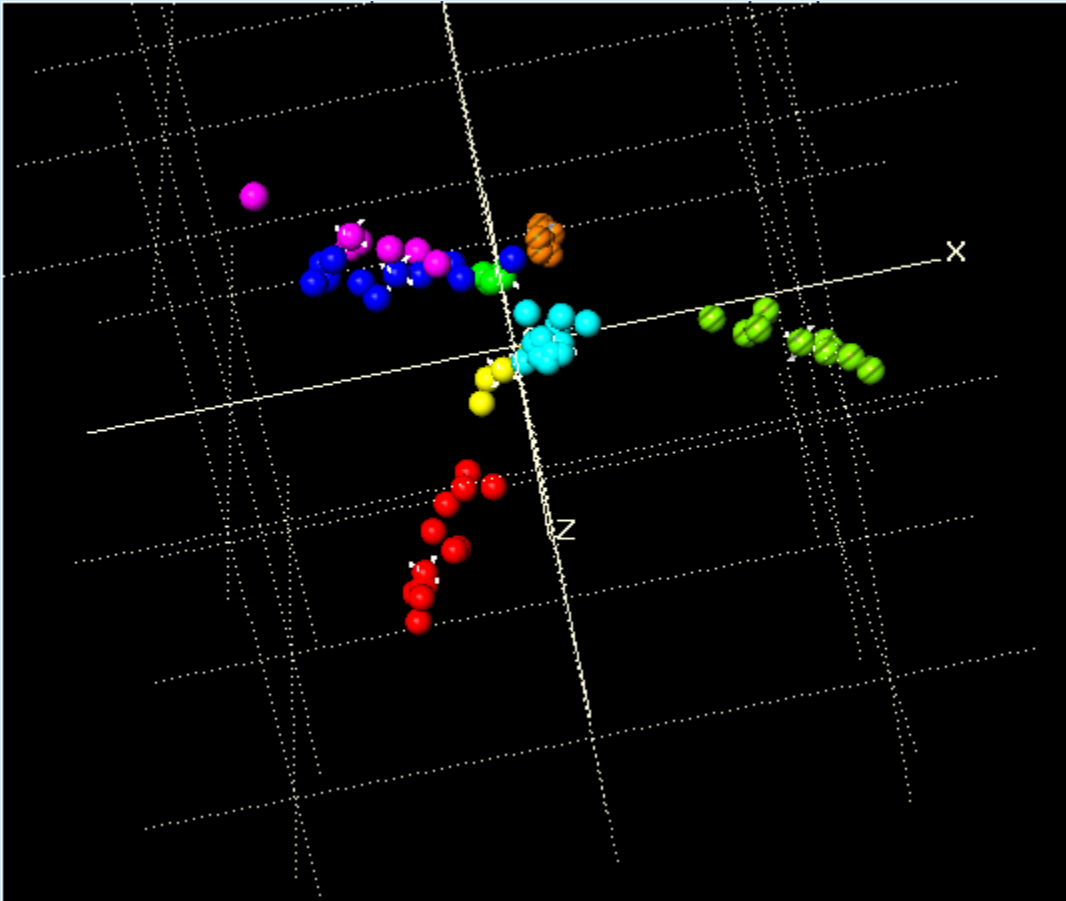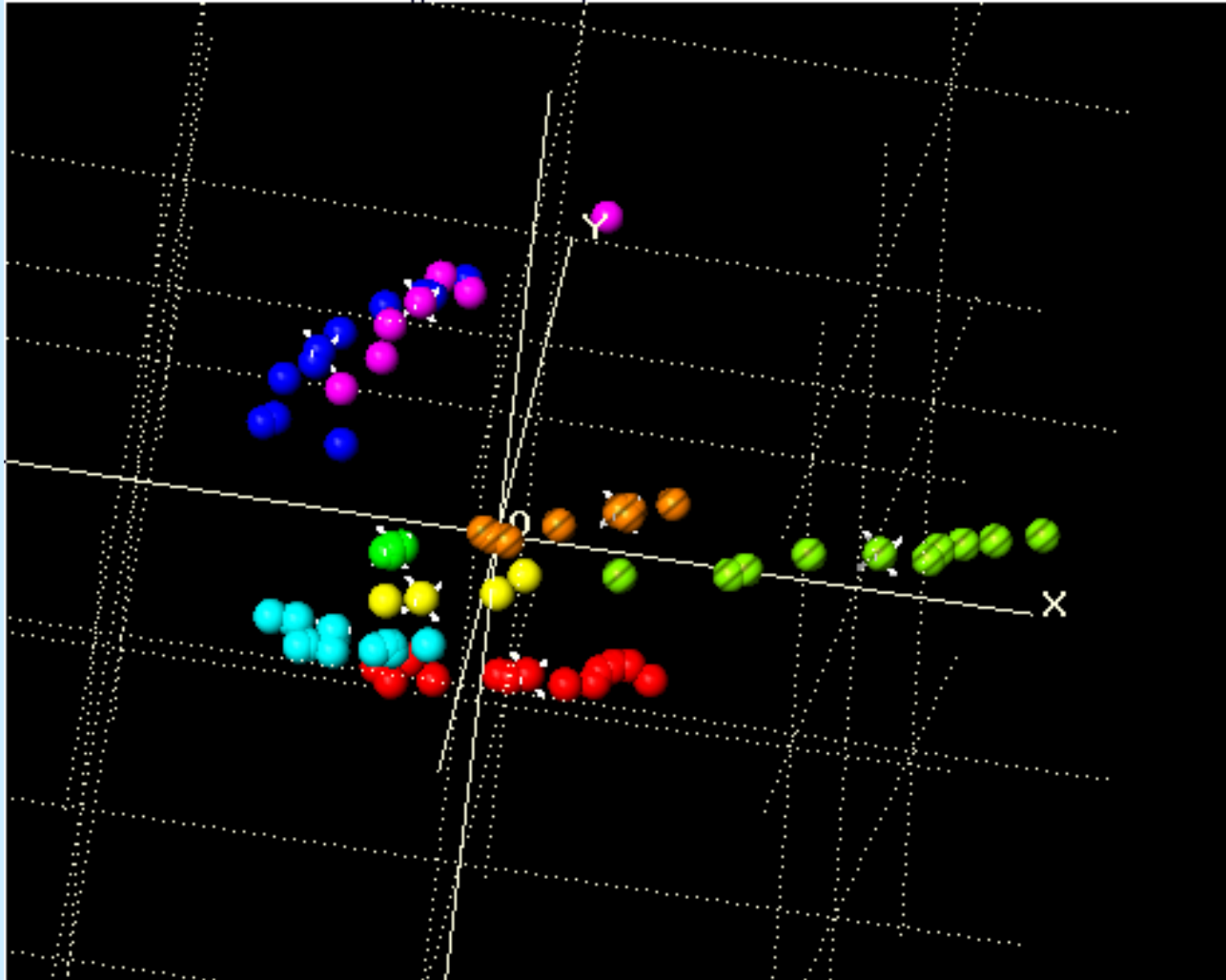
carbamazepine

furosemide

phenylbutazone

cimetidene

mefenamicacid

sulfamerazine
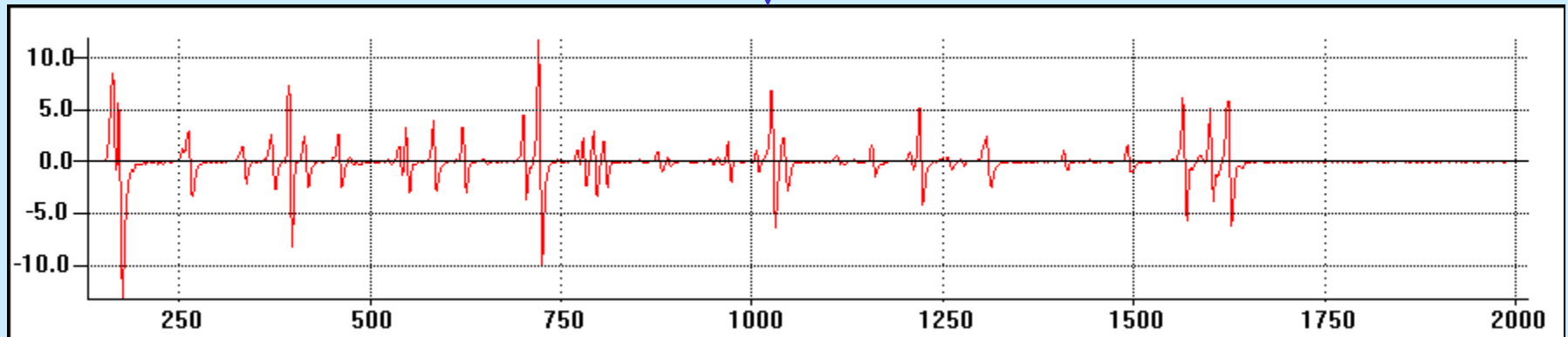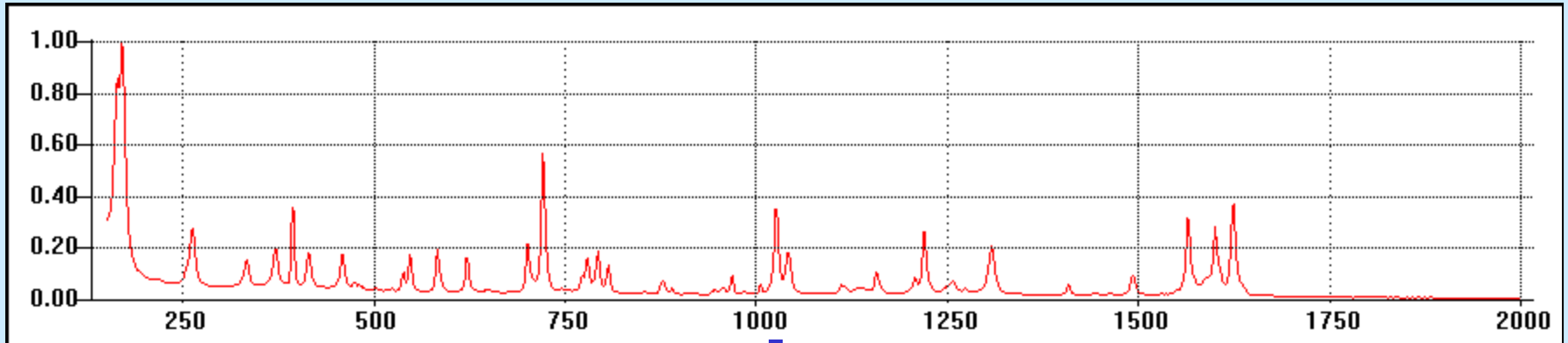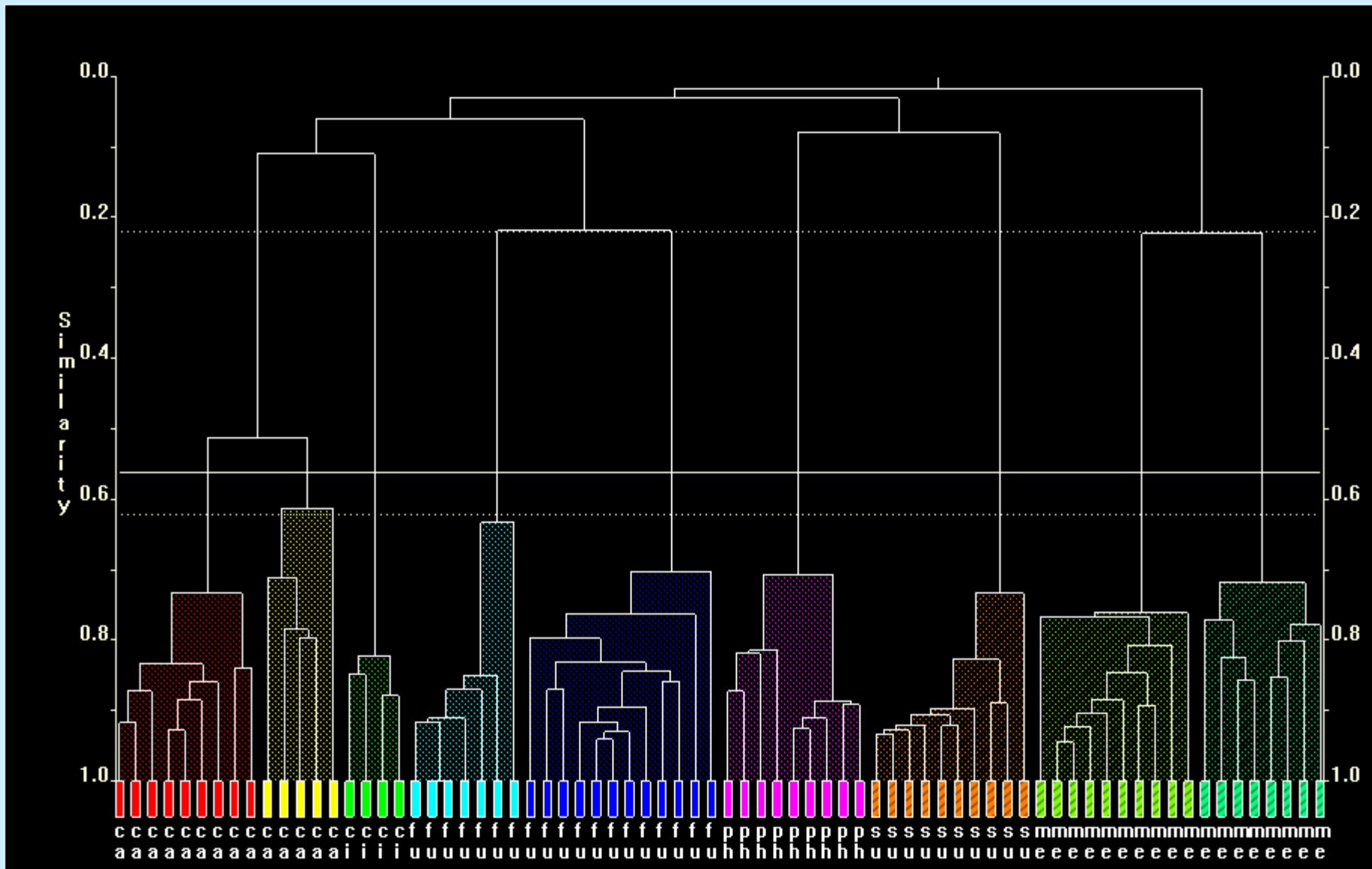
# MMDS

# PCA – Be Wary!
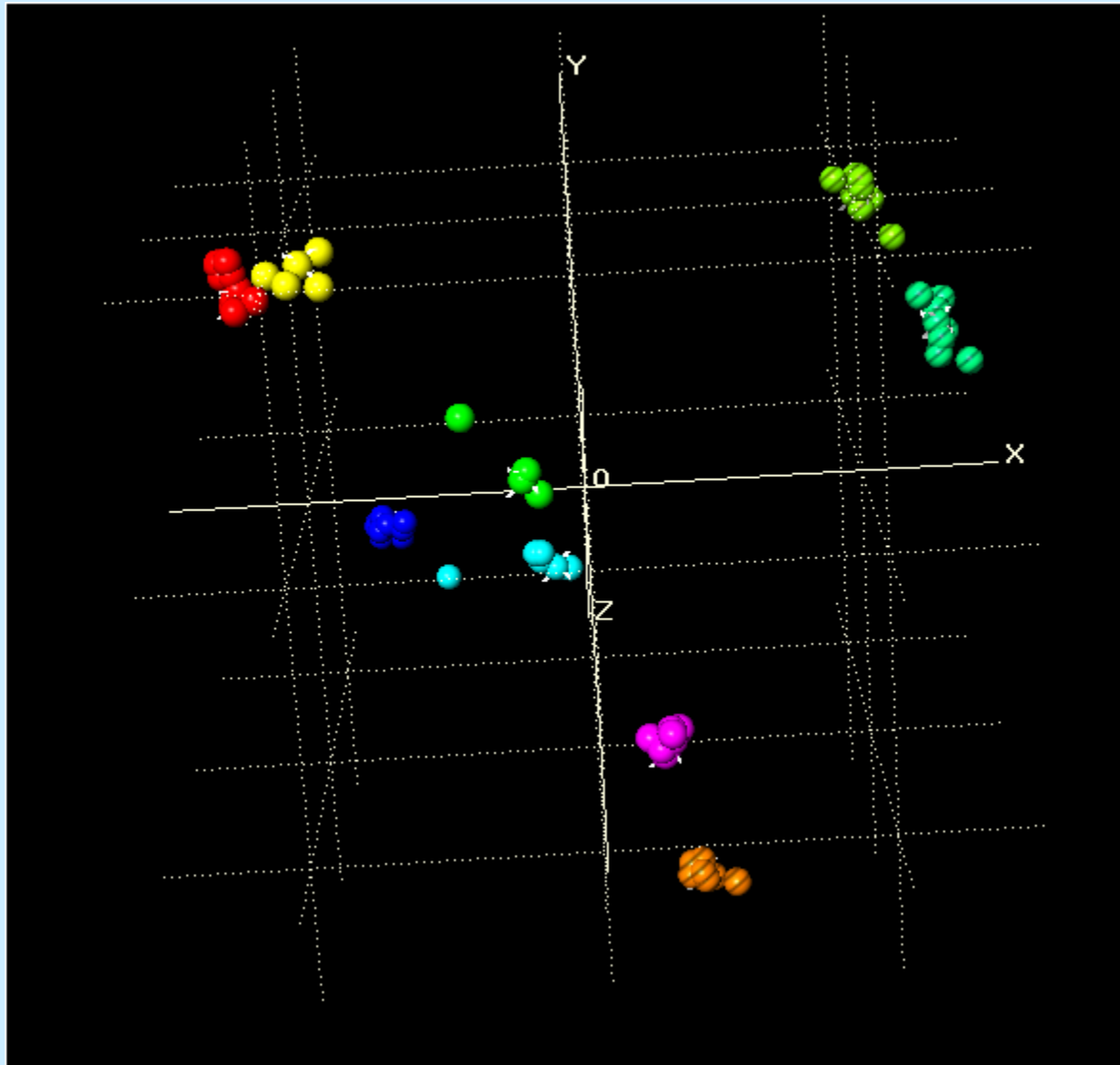
# Use 1st Derivative Data

# Dendrogram on Derivative Data

# MMDS on Derivative Data

# How Do You Combine Different Data Types?

- **Combined XRPD + Raman instruments now available**

- **Applying multiple techniques to the same samples gives additional information**

- **How would we actually combine results from two (or more) such different techniques ?**

# Combining Data Types

- **Manual weighting:**
  - **Give a single weight to each dataset *e.g.* Powder 0.8, Raman 0.2**
  - **Use Fischer transforms.**

- **Dynamic weighting:**
  - **Automatically calculate optimal weighting for each entry in each dataset**
  - **Unbiased solution that scales the differences between individual distance matrices**

# Dynamic weighting

- **Dynamic Weighting using INDSCAL:**

  **Independent Scaling of Differences**
  **Carroll & Chang, (1970) *Psychometrica* 35, 283-319**

- **Each data set has a 2-D distance matrix $d$**

- **$D_k$ is squared ($n \times n$) distance matrix for dataset $k$**

  ***e.g.* we have Raman and XRPD data on 20 samples, so $k = 2$, $n=20$.**

- **We want a Group Average Matrix, G, to optimally describe our data**

- **Specify diagonal weight matrices $W_k$ which can vary over the $k$ datasets**

# Dynamic Weighting

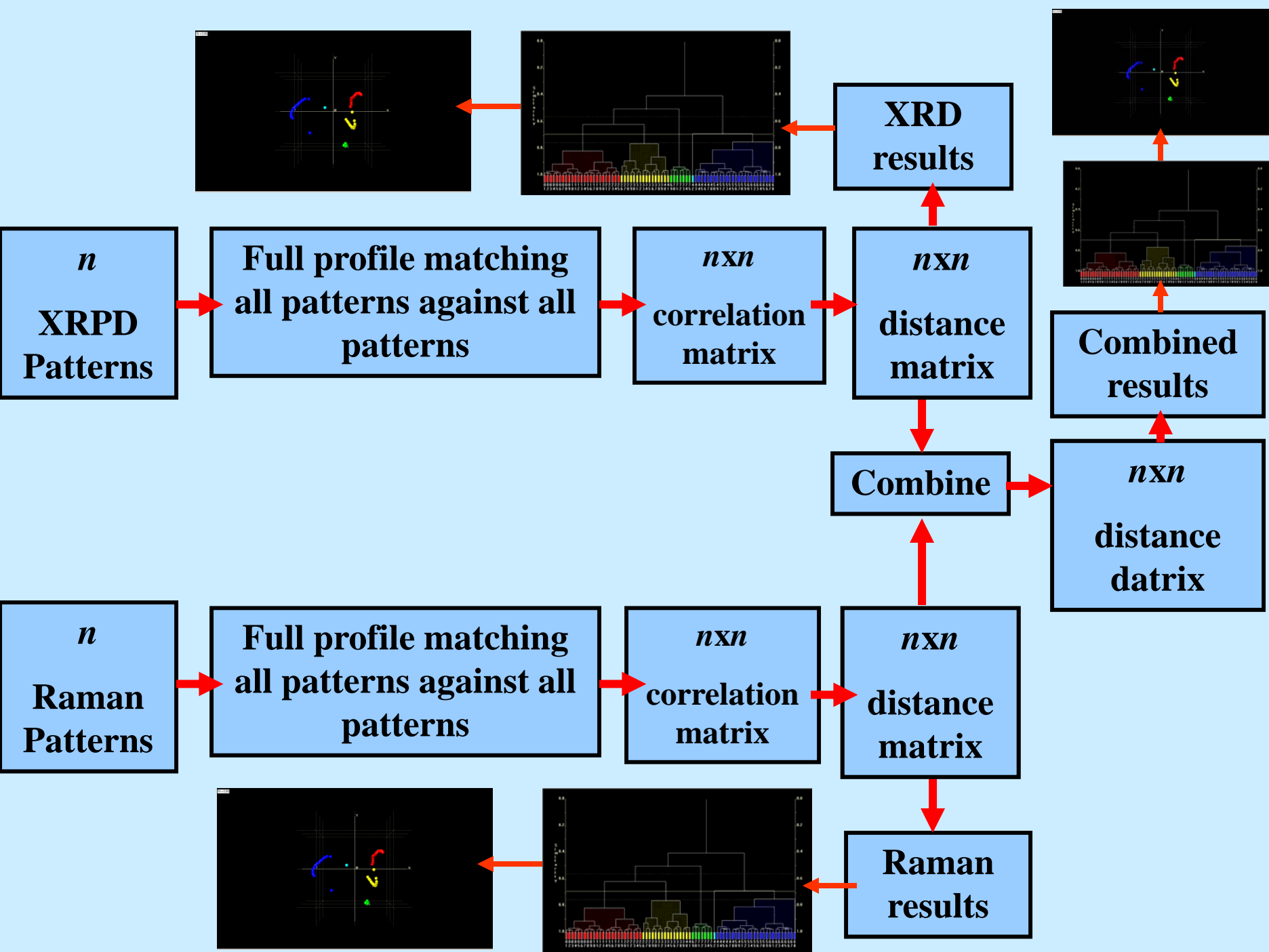**Matrices are matched to weighted form of G by minimising**

$$\sum_{k=1}^{K} \left\| \mathbf{B}_k - \mathbf{G}\mathbf{W}_k^2\mathbf{G}' \right\|$$

**(1)**

**Where**

$$\mathbf{B}_k = -\frac{1}{2}(\mathbf{I} - \mathbf{N})\mathbf{D}_k(\mathbf{I} - \mathbf{N})$$

**(a double-centering operation on D), and solve (1) to get best values for G and W**

**The resulting G matrix is then used as before**

**XRD results**



| **$n$** | | |
|---|---|---|
| **XRPD Patterns** | **Full profile matching all patterns against all patterns** | **$n$x$n$ correlation matrix** |

**$n$x$n$ distance matrix**

**$n$x$n$ distance matrix**

**Combine**

**Combined results**

**$n$x$n$ distance datrix**

| **$n$** | | |
|---|---|---|
| **Raman Patterns** | **Full profile matching all patterns against all patterns** | **$n$x$n$ correlation matrix** |

**$n$x$n$ distance matrix**





**Raman results**

# PXRD + Raman

- **Forms 2,3 and 4 of sulfathizole**
- **48 samples, no mixtures.**

**Run PolySNAP on...**

Analyse: ○ **Single Dataset**  ◉ **Multiple Datasets**

**Dataset 1:** [Powder XRD ▼]

C:\Cluster Analysis\Gordon Data Sets\gordon multiple datatype data\sulfathiazole-carbamazepine\xray\

[Folder...]  [File...]

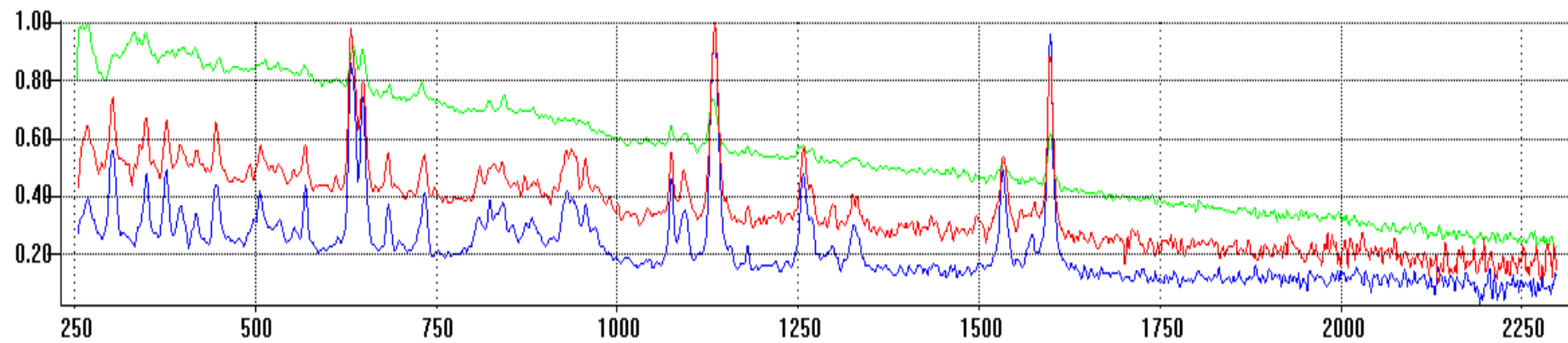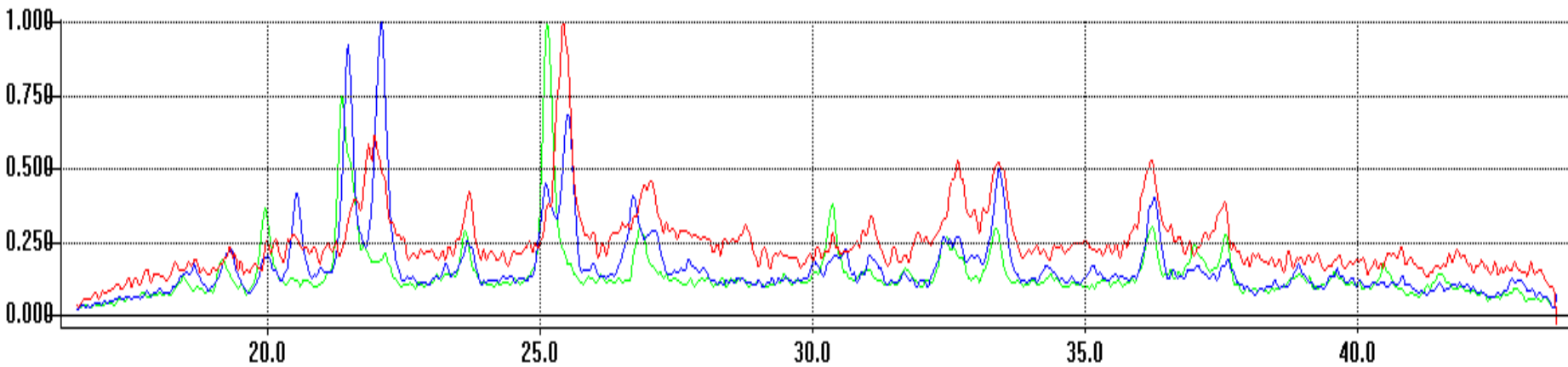☐ Use folder or database containing Known/Reference data files:
C:\Cluster Analysis\Gordon Data Sets\gordon multiple datatype data\sulfathiazole-carbamazepine\Pure PXRD phases\

[Folder...]  [File...]

☐ Load sample image files from separate folder:
<None>

[Folder...]

**Advanced Options**

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| Allow x-shift calculation (sin theta) for datasets | ☐ | ☐ | ☐ | ☐ |
| Denoise Patterns for datasets | ☐ | ☑ | ☐ | ☐ |
| Subtract Background from datasets | ☐ | ☑ | ☐ | ☐ |
| Check for amorphous samples in datasets | ☑ | ☐ | ☐ | ☐ |
| Remove cosmic ray spikes from datasets | ☐ | ☐ | ☐ | ☐ |
| Mask specified regions in datasets | ☐ | ☐ | ☐ | ☐ |
| Set matching range subset in datasets | ☐ | ☐ | ☐ | ☐ |
| Apply signal transform to datasets | ☐ | ☐ | ☐ | ☐ |

☐ Include reference files in main calculation    ☐ Hide results similar to references

**Output options:**
Combine the multiple datasets using weights: ◉ Automatic ○ None ○ Manual:
Dataset 1: [1.0]  Dataset 2: [1.0]  Dataset 3: [1.0]  Dataset 4: [1.0]

**Dataset 2:** [Raman ▼]

C:\Cluster Analysis\Papers\Snap Paper 5\ps2demo_suthaz\raman\

[Folder...]  [File...]

☐ Use folder or database containing Known/Reference data files:
C:\

[Folder...]  [File...]

**Dataset 3:** [IR ▼]

C:\Cluster Analysis\Gordon Data Sets\gordon multiple datatype data\sulfathiazole-carbamazepine\IR\

[Folder...]  [File...]

☐ Use folder or database containing Known/Reference data files:
C:\

[Folder...]  [File...]

**Dataset 4:** [DSC ▼]

C:\Cluster Analysis\Gordon Data Sets\gordon multiple datatype data\sulfathiazole-carbamazepine\dsc\

[Folder...]  [File...]

☐ Use folder or database containing Known/Reference data files:
C:\

[Folder...]  [File...]

[Cancel]  [OK]

# PXRD

**Raman**

# Raman and PXRD Averaged

INDSCAL

# Other Combinations

- **Raman + Raman derivative data.**

- **Different data collection protocols/apparatus on the same samples.**

- **You can include numeric data as a data type *e.g.* image analysis data.**

# Numeric data – get the distance matrix directly.

**Sample 1:  113.431 58.531  155.845 … {x$_{11}$, x$_{12}$ , x$_{13}$…….}**

**Sample 2:  113.44   58.328  153.602  …  {x$_{21}$, x$_{22}$ , x$_{23}$…….}**

**Sample 3: 117.873  60.117  93.686    … {x$_{31}$, x$_{32}$ , x$_{33}$…….}**

$$d_{ij} = \left( \sum_{k=1}^{m} w_k \left| x_{ik} - x_{jk} \right|^2 \right)^{\frac{1}{2}} \Rightarrow \mathbf{d}$$

# Raman + Derivative Data

# PXRD, Raman, IR + DSC

**16 samples containing 3 forms of sulfathiazole and carbamazepine + mixtures:**

- **PXRD:  Bruker C2 GADDS,**

- **IR: JASCO FT/IR 4100,**

- **DSC: TA instruments Q100.**

- **Raman: Renishaw inVia Reflex Spectrometer System.**

# PXRD, Raman, IR + DSC

- **15 different data sets and combinations!**

- **You need good software to explore all these options.**

# PXRD + IR + Raman + DSC

# PXRD + IR + Raman + DSC

All sulfathiazole mixtures

All carbamazepine form 3 mixtures

# PXRD + XRF



Clustering of XRF subset with XRD

Lowe-Ma, et al.

Ford Motor Company

# PolySNAP 3

**You can do this with PolySNAP3**

http://www.chem.gla.ac.uk/snap/
   PolySNAP_index.html

**See:**

**Barr, Dong & Gilmore,** *J. Appl. Cryst.,* **(2009),** <u>42</u>**, 965-974.**

**or**

**Use MATLAB/R/Sage**

# In progress

- **Missing data.**
- **DSC processing**
- **Non numeric data**

# macroSNAP

## dSNAP for proteins – watch this space!

# Acknowledgements

- **Michael Hermannn (Fraunhofer-Institut Chemische Technologie)**

- **Karsten Knorr & Arnt Kern  (Bruker AXS, Karlsruhe)**

- **Chris Frampton & Susie Buttar, Pharmorphix**