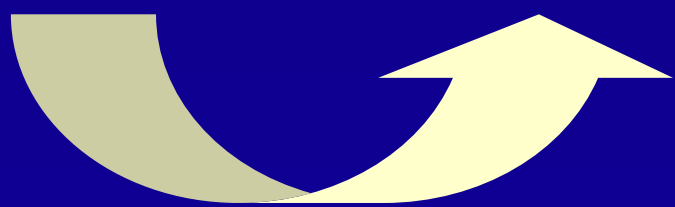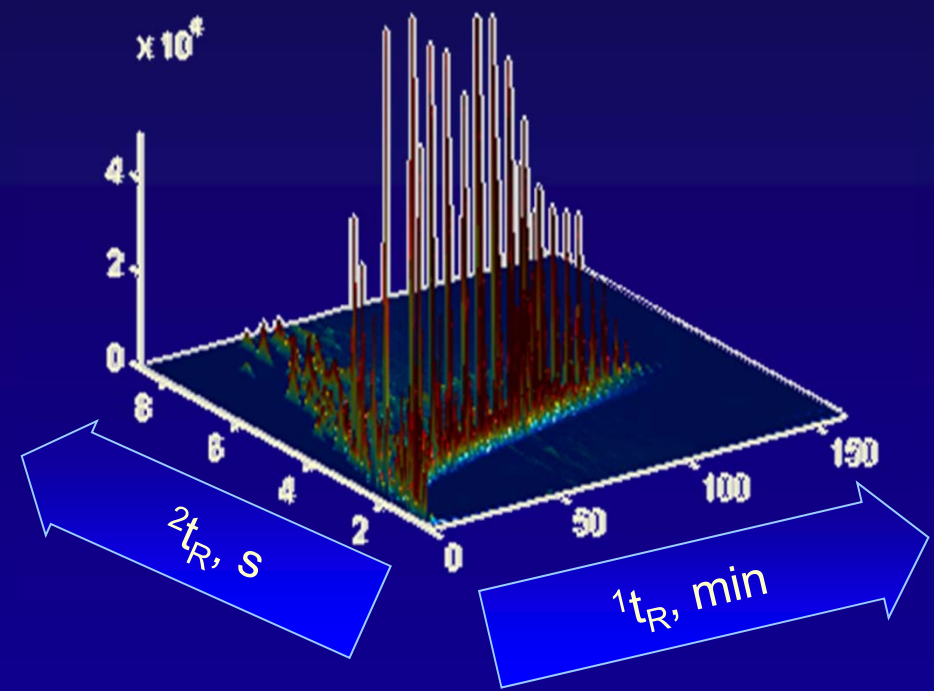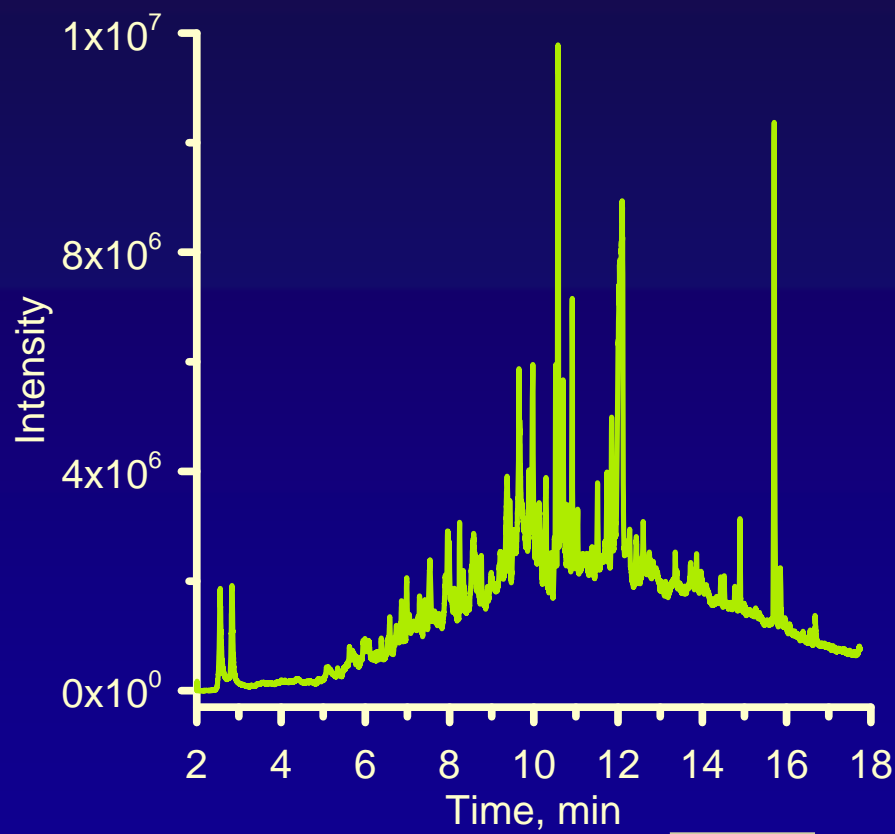# *Data analysis in GCxGC*

*Gabriel Vivó Truyols*
*Analytical-chemistry group*
*Van 't Hoff Institute for Molecular Sciences*
*University of Amsterdam*
*g.vivotruyols@uva.nl*

**From 1D chromatography to 2D chromatography: what does change?**

## The complexity of the data changes

Instruments can be classified according to the order of the tensor of data used to represent a single experiment:

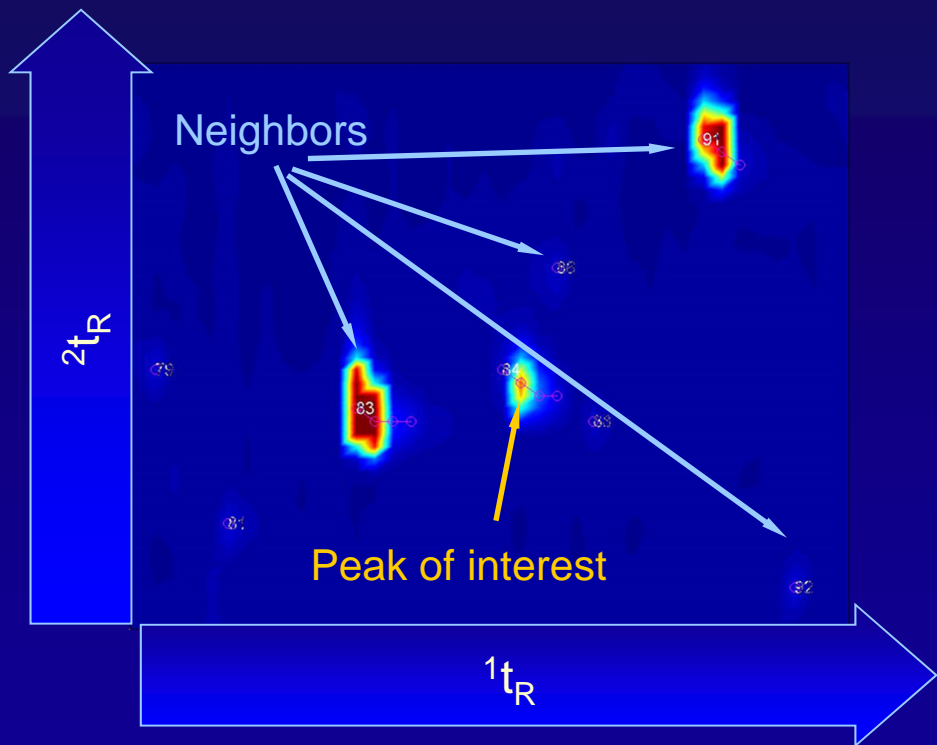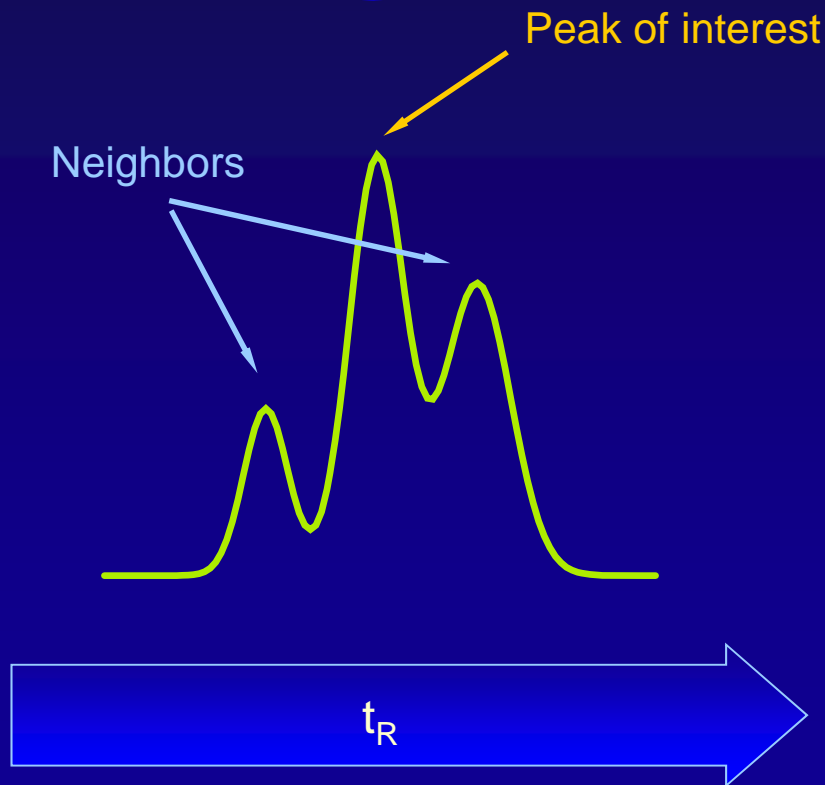| | | | | |
|---|---|---|---|---|
| Zero-order instruments | Produce | Zero-order tensor (e.g. a number) | Example | pH - meter |
| First-order instruments | Produce | First-order tensor (e.g. a vector) | Example | UV-VIS spectrometer |
| Second-order instruments | Produce | Second-order tensor (e.g. a matrix) | Example | LC-MS, GCxGC |
| Third-order instrument | Produce | Third-order tensor (e.g. a "cube" of data) | Example | GCxGC-MS |

$n^{th}$-order instrument exist, but they are rare

S. Peters, G. Vivó-Truyols, P. Marriott, P.J. Schoenmakers, J. Chromatogr. A. 1146 (2007), 232-241.

S. Peters, G. Vivó-Truyols, P. Marriott, P.J. Schoenmakers, J. Chromatogr. A. 1146 (2007), 232-241.

# *First step: visualization*

**Step 1**

**View**

- Folding
- Phasing

**Step 2**

**Pre-process**

- Base-line correction
- Noise filtering
- Spike filtering
- Alignment
- … etc.

**Step 3**

**Measure**

- Peak detection / integration
- Calibration
- Deconvolution
- Pattern recognition
- Class separation

## Raw data in 2D chromatography

Chromatogram of Glycine preparation (254nm)
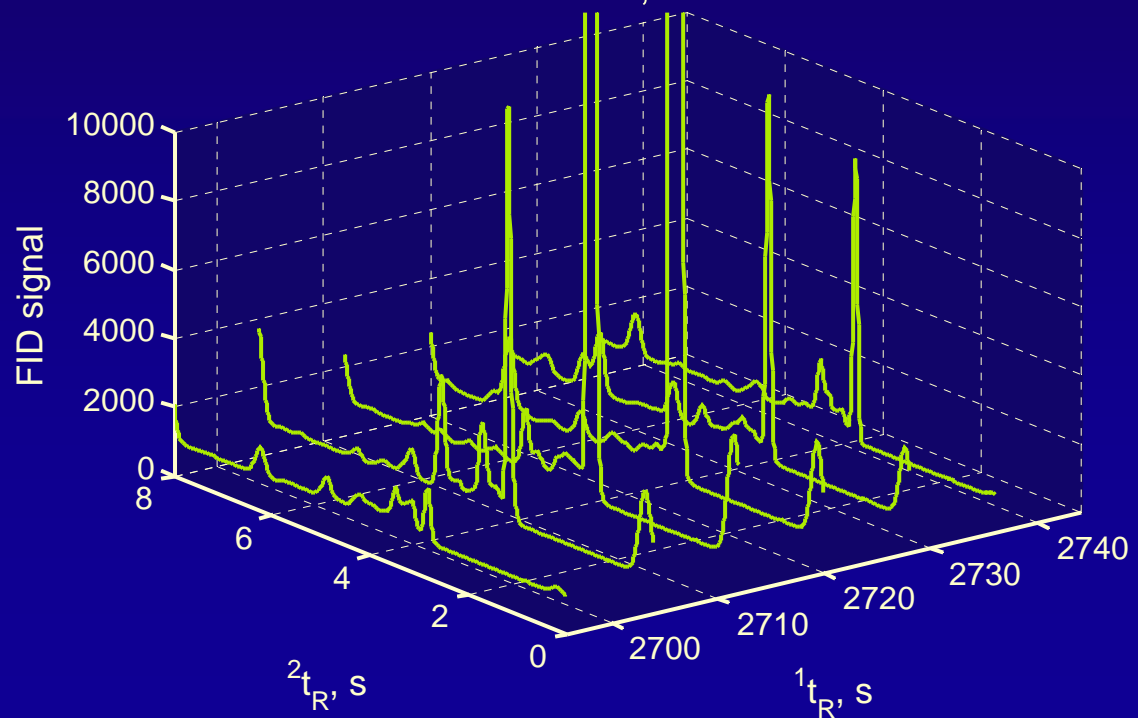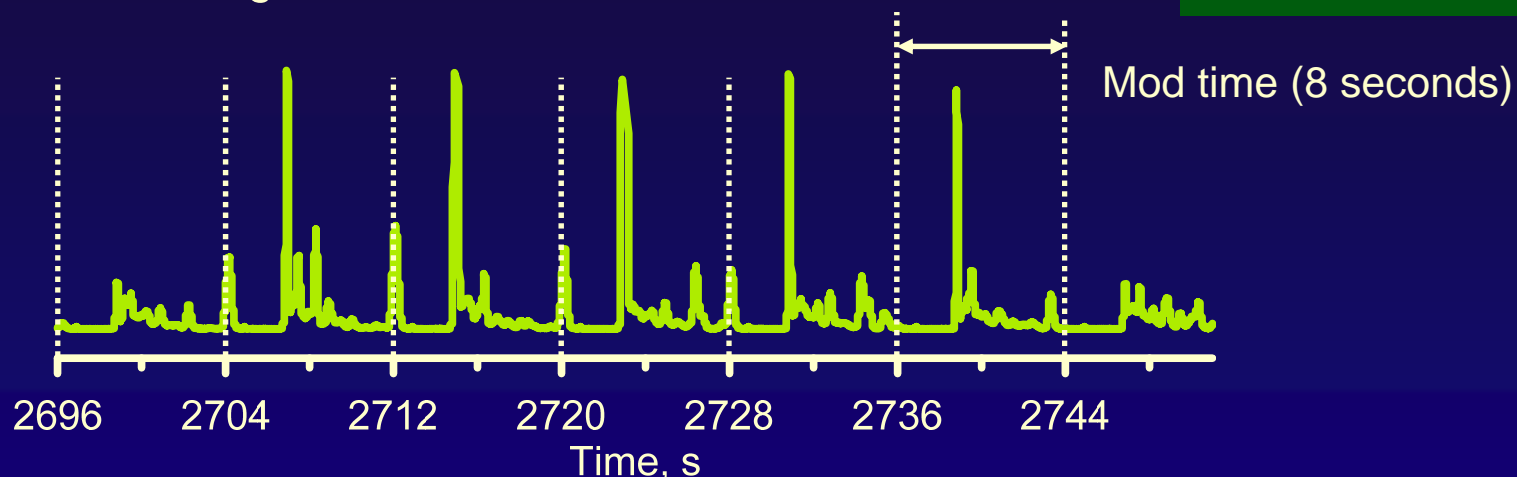(data courtesy of University of Valencia)



## Two-dimensional chromatography

GCxGC chromatogram of diesel (FID detector)
First dimension: non-polar; Second dimension: polar



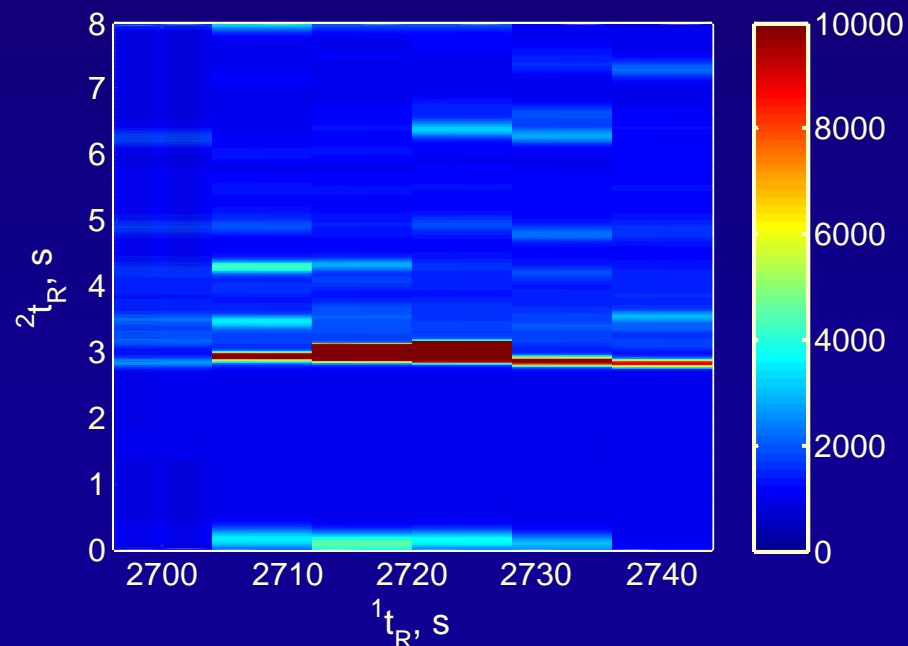*GCxGC, Riva - 2014*
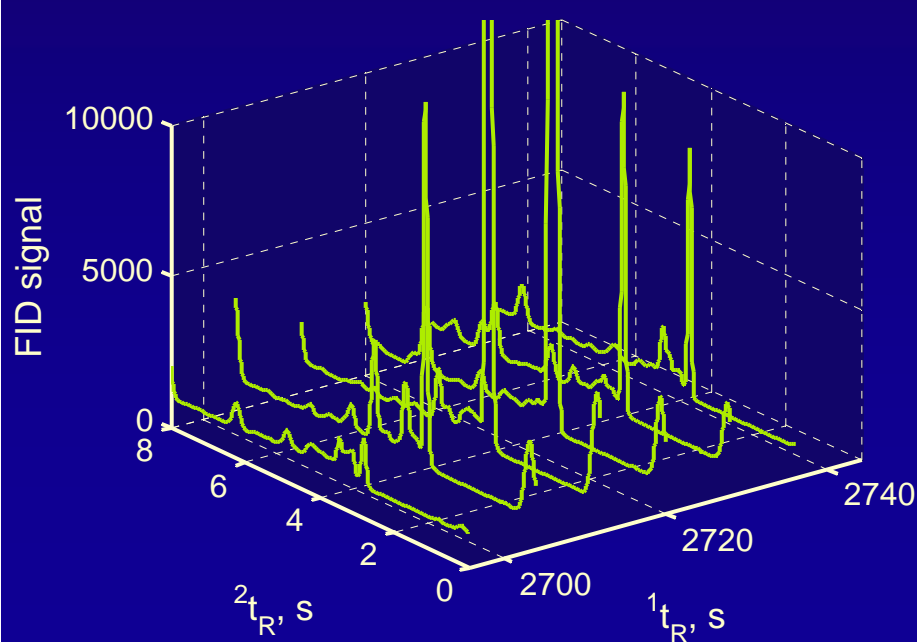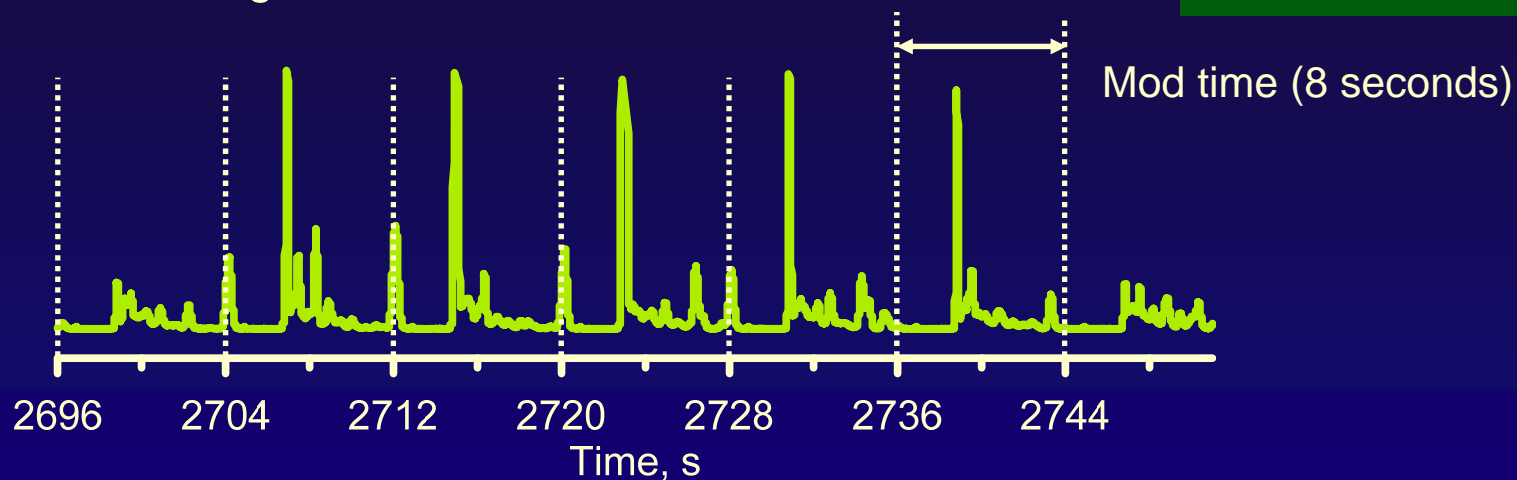
"Folding" the chromatogram

*Visualisation*

Mod time (8 seconds)

Time, s

FID signal

$^2t_R$, s

$^1t_R$, s

GCxGC, Riva - 2014

Cylindrical coordinates. An alternative way to represent the data.

J.J.A.M. Weusten, E.P.P.A. Derks , J.H.M. Mommers, S. van der Wal, Anal. Chim. Acta 726 (2012), 9

Interpolation

*Visualisation*

Not interpolated

Interpolated

FID signal

$^2t_R$, s

$^1t_R$, s

$^1t_R$, s

$^1t_R$, s

*GCxGC, Riva - 2014*

Interpolation

## Welcome to the magic world of chemometrics!

## "Folding" the chromatogram: final result

## Conclusions

- Visualising is simple, and gives a lot of information.

- Folding (one-dimensional) data into (2D) image introduces discontinuities in the edges. Other visualization methods (cylindrical coordinates) possibloe.

- Phasing can be of great help.

- Careful with "cosmetic" effects!

## Typical problems: base-line drifts and noise

**Base-line drifts**

**Noise**

Van 't Hoff Institute for Molecular Sciences    University of Amsterdam

Base-line drifts.

Pre-processing

Weighted least squares fitting

Corrected chromatogr.

Fitted base-line correction

Original chromatogram

GCxGC, Riva - 2014

# Base-line drifts.

Original
chromatogram

*Weighted least
squares fitting*

Fitted base-line
correction

Corrected
chromatogram

*GCxGC, Riva - 2014*

## Base-line drifts.

Original chromatogram

*Other (more sophisticated) options:*

*- Use splines*
*- Base-line correction coupled to peak detection*
*- Fourier-transform based approaches*
*- Wavelet-based approaches*

*Weighted least squares fitting*

Fitted base-line correction

Corrected chromatogram

FID signal

Time, s

*GCxGC, Riva - 2014*

## Base-line drifts.

Base-line is reached at the (half) upper part

Base-line is reached at the (half) bottom part

S.E. Reichenbach, M. Ni, D. Zhang, E.B. Ledford Jr., J. Chromatogr. A, 985 (2003) 47 - 56

Base-line drifts.

Consider the positions with the smallest values in each half

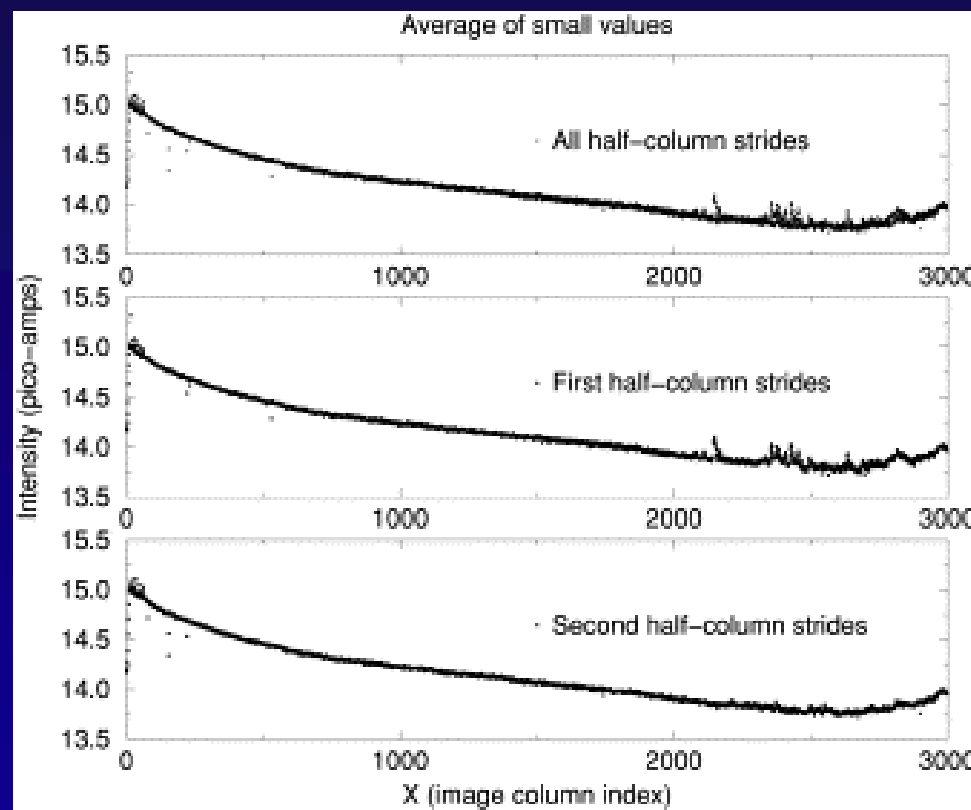Estimate local background parameters using neighboring values

Interpolate the main background trend and substract it



Average of small values

S.E. Reichenbach, M. Ni, D. Zhang, E.B. Ledford Jr., J. Chromatogr. A, 985 (2003) 47 - 56

Noise removal. Smoothing and derivatives.        *Pre-processing*

Savitzky-Golay filter is the most common method

Two parameters should be optmisized
- Window size
- Polynomial degree

These parameters govern how much correlated noise is removed

- Large window sizes and low polynomial degree

*Too much noise is removed (chromatograms appear deformed)*

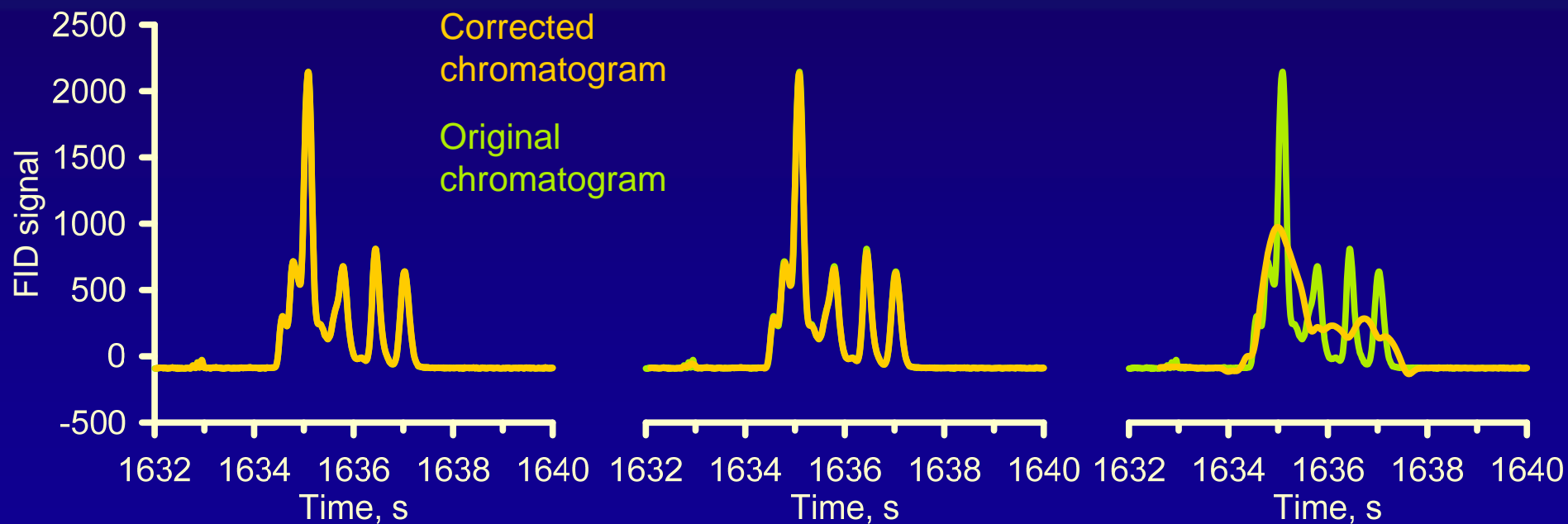- Small window sizes and large polynomial degrees

*Too much noise remains*

## Alignment. Using (truly) 2D algorithms

*Score alignment in GCxGC-MS*

S. Castillo, I. Mattila, J. Miettinen, M. Orešič, T.Hyötyläinen, Anal. Chem. 83 ( 2011) 3058–3067

*COW-adapted GCxGC-MS (using single channel)*

D. Zhang, X. Huang, F.E. Regnier, M. Zhang, Anal. Chem., 80 (2008) 2664–2671



a. The Partitioned Chromatograph     b. The Warped Chromatograph

## Conclusions

- Pre-processing methods are almost the same: one-dimensional = two-dimensional. Normally done in the (pre-folded) raw data.

- Every case needs a particular solution (it always exists, but some care should be taken!)
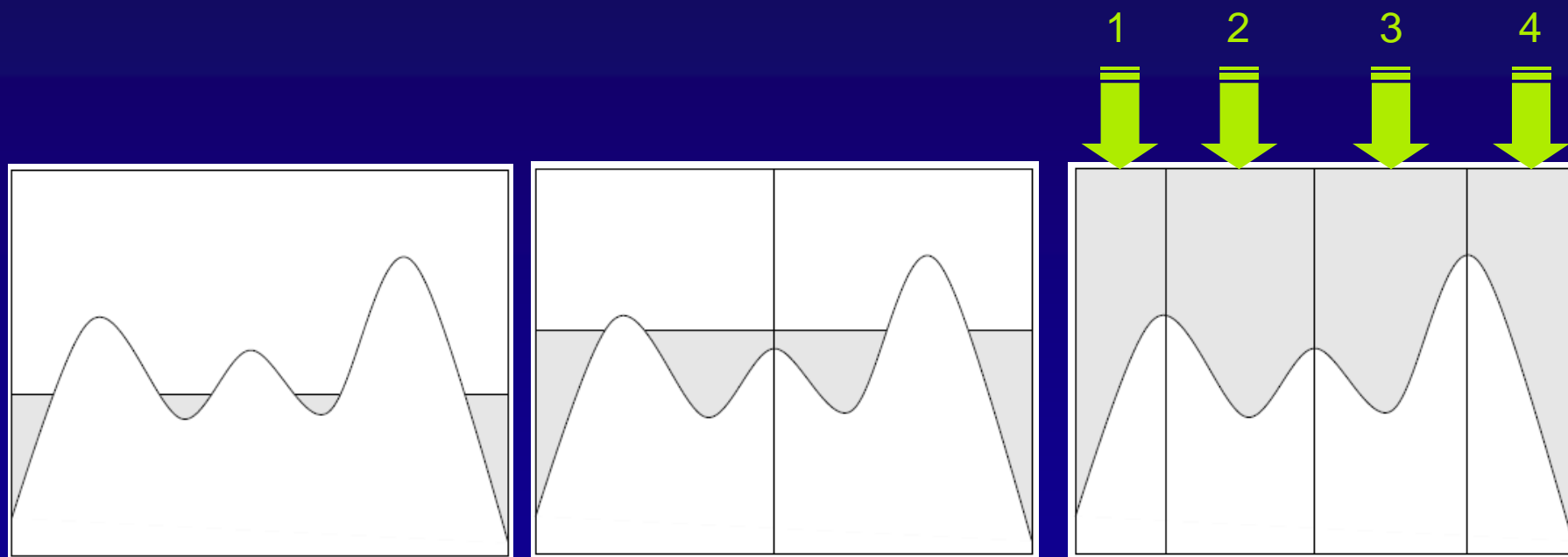
# *Third step: measure*

**Step 1**

**View**

- Folding
- Phasing

**Step 2**

**Pre-process**

- Base-line correction
- Noise filtering
- Spike filtering
- Alignment
- … etc.

**Step 3**

**Measure**

- Peak detection / integration
- Calibration
- Deconvolution
- Pattern recognition
- Class separation

Peak detection in one step: the watershed algorithm

Most common software programs use the watershed algorithm to detect peaks in 2D chromatography:

1    2    3    4

J. De Bock et al., doi 10.1007/11558484

Peak detection in one step: the watershed algorithm

*Peak detection*

Single catchment basin?

????

G. Vivó-Truyols, H.G. Janssen, J. Chromatogr. A, doi:10.1016/j.chroma.2009.12.063

Peak detection in one step: the watershed algorithm

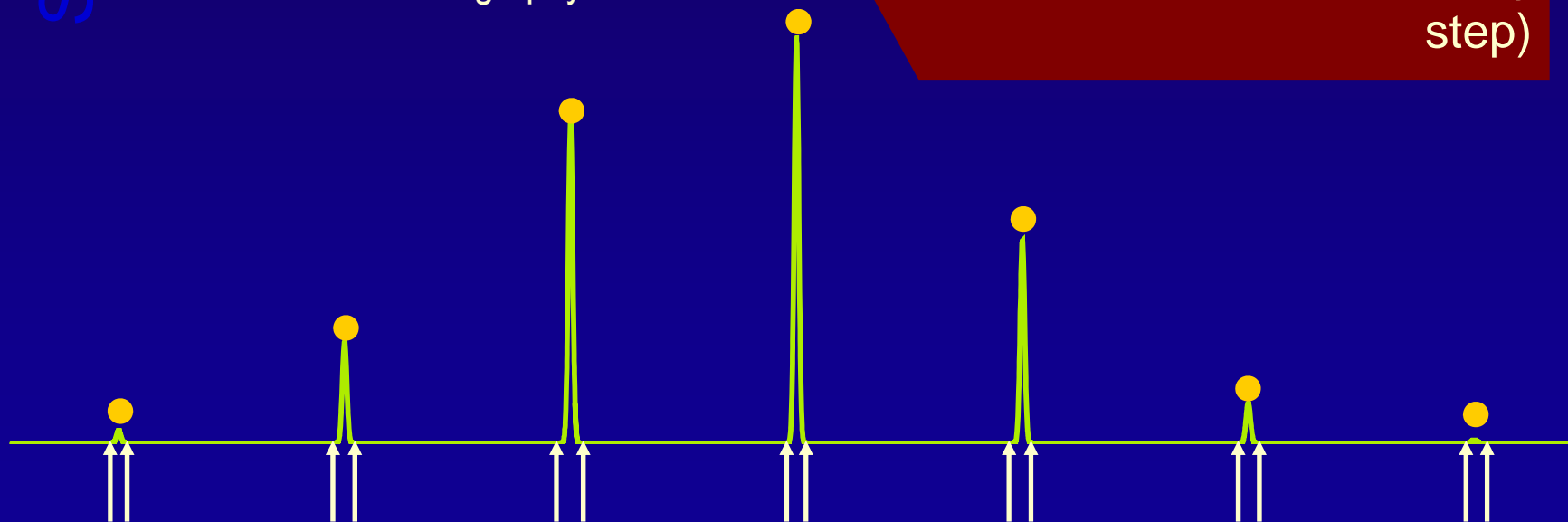The first problem using the watershed algorithm.

Peak detectioni in two steps.

S. Peters, G. Vivó-Truyols, P.J. Marriott and P.J. Schoenmakers, J. Chromatogr. A 1156 (2007) 14.

E.J.C. van der Klift, G. Vivó-Truyols, F.W. Claassen, F.L. van Holthoon, T.A. van Beek, J. Chromatogr. A, 1178 (2008) 43.

Step 1  Detect peaks as in one-dimensional chromatography

Use information from derivatives (pre-processing step)
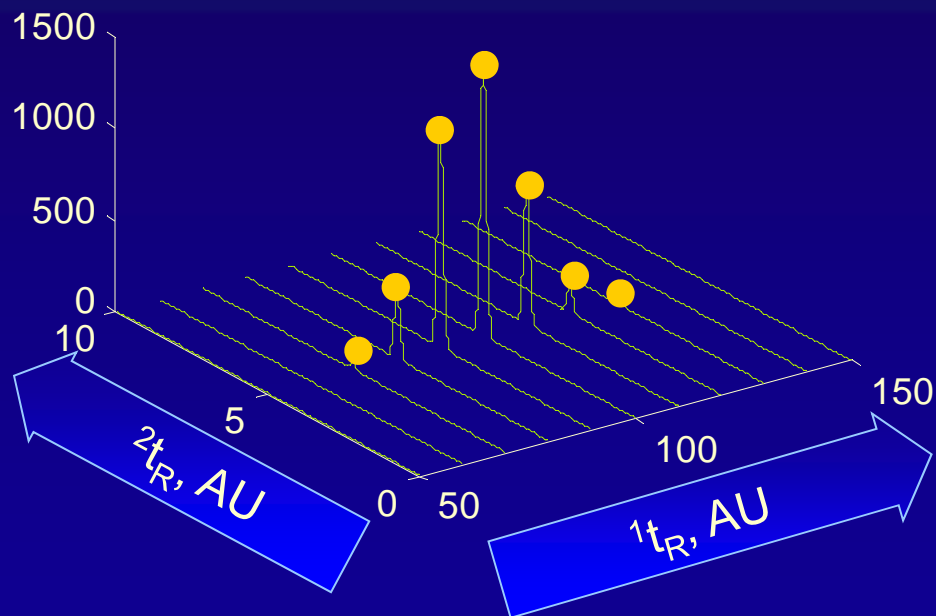
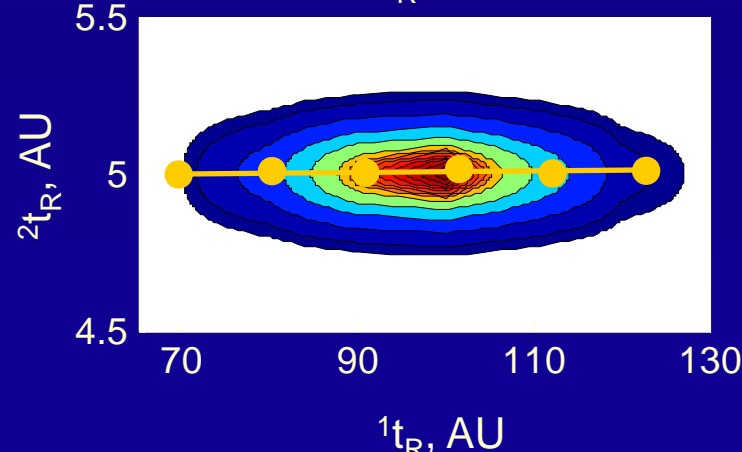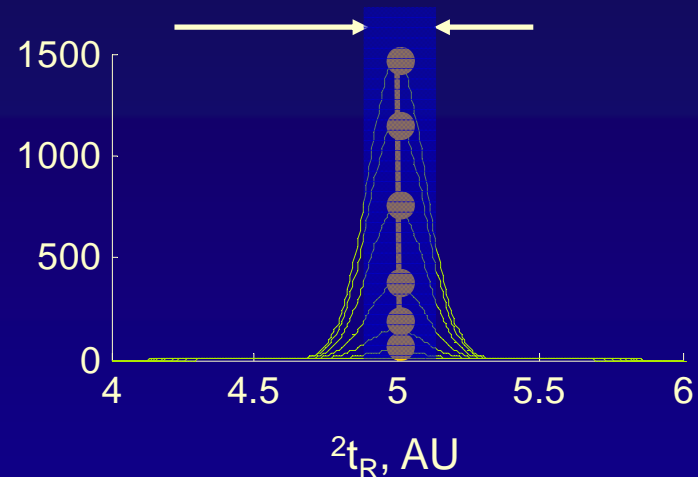Time, arbitrary units (AU)

GCxC

Peak detectioni in two steps.

**2** step

Merge peaks that belong to the same compound according to 2nd-dimension retention time differences



T: Tolerance criterion

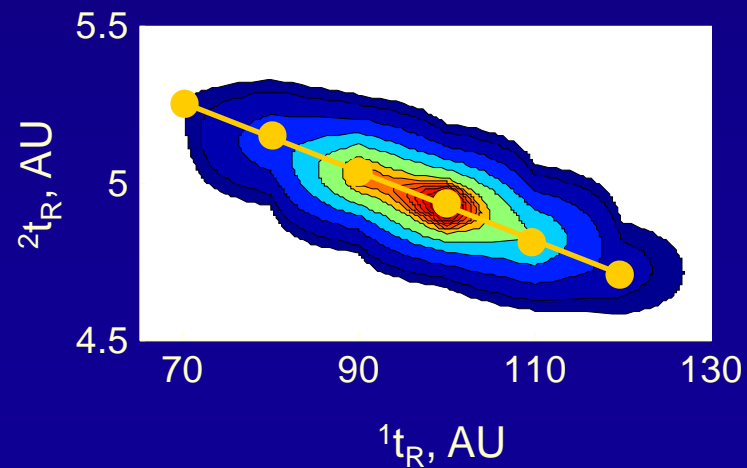Peak detectioni in two steps.
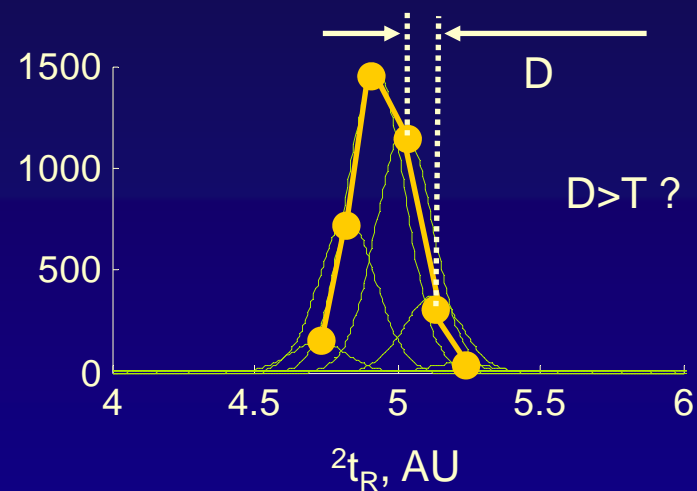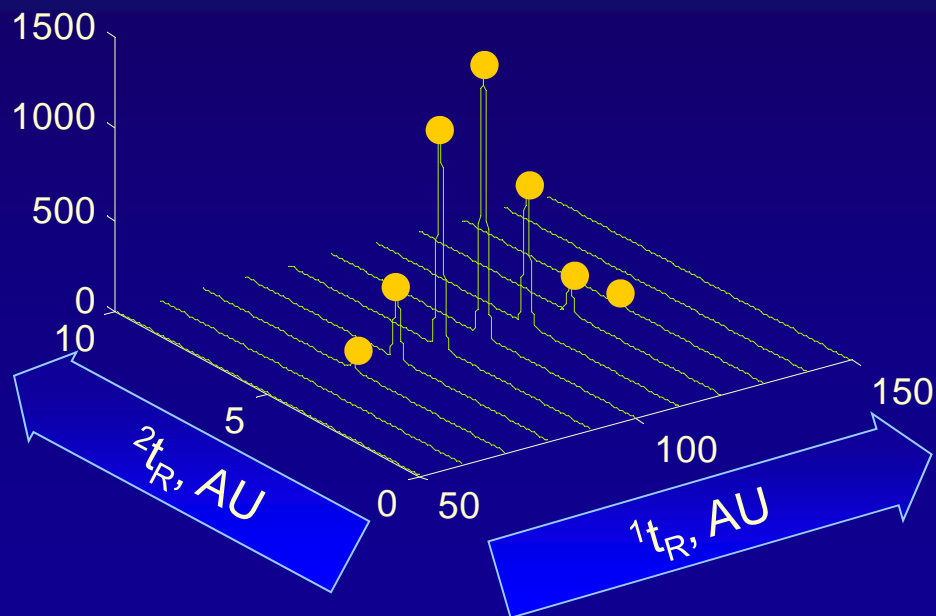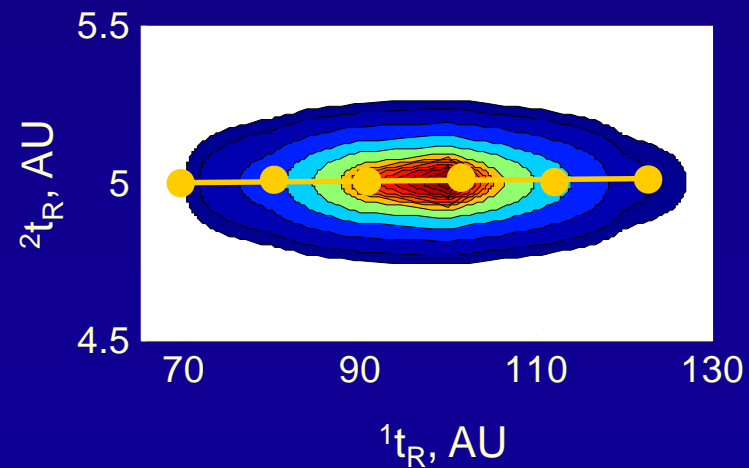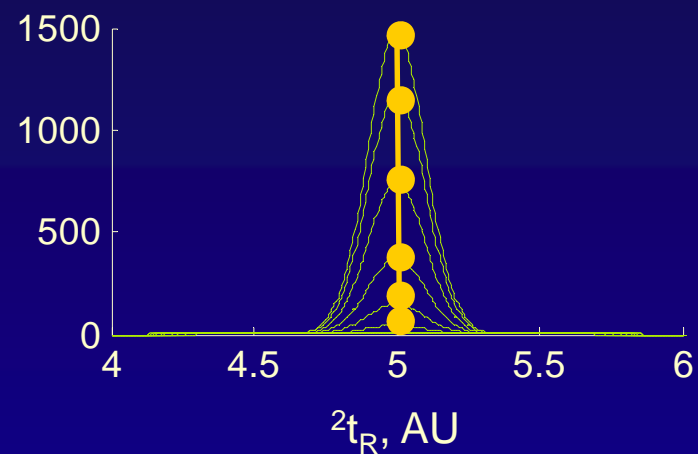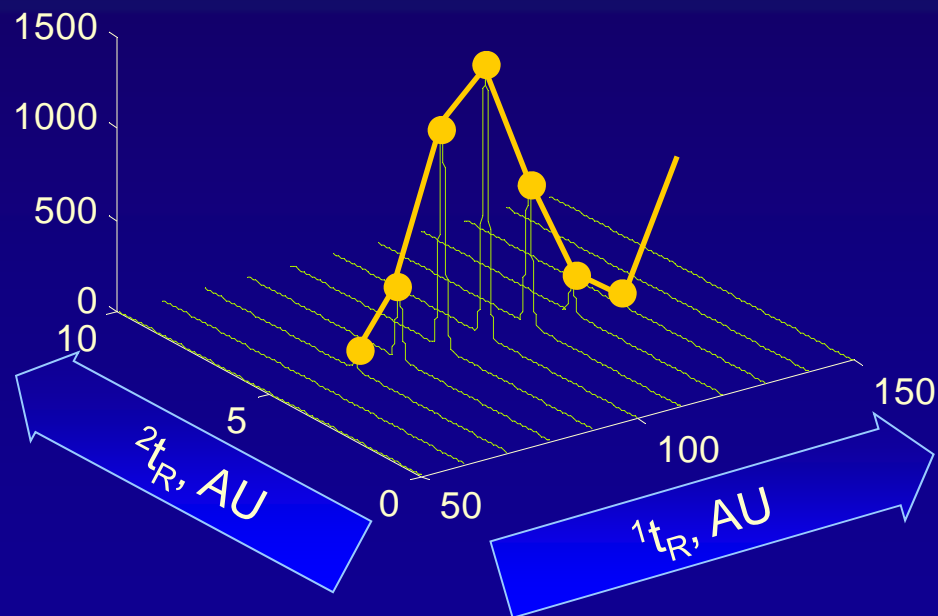
**3** Check unimodality

## Conclusions

- Two methods available: (inverse) watershed, and two-step peak detection process.

- Peak detection seems to be still a subject of discussion.

## Deconvolution mehtods.

### An example of PARAFAC



A.E. Sinha, J.L. Hope, B.J. Prazen, C.G. Fraga, E.J. Nilsson, R.E. Synovec, J. Chromatogr. A, 1056 (2004) 145 - 154

# *Third step: measure*

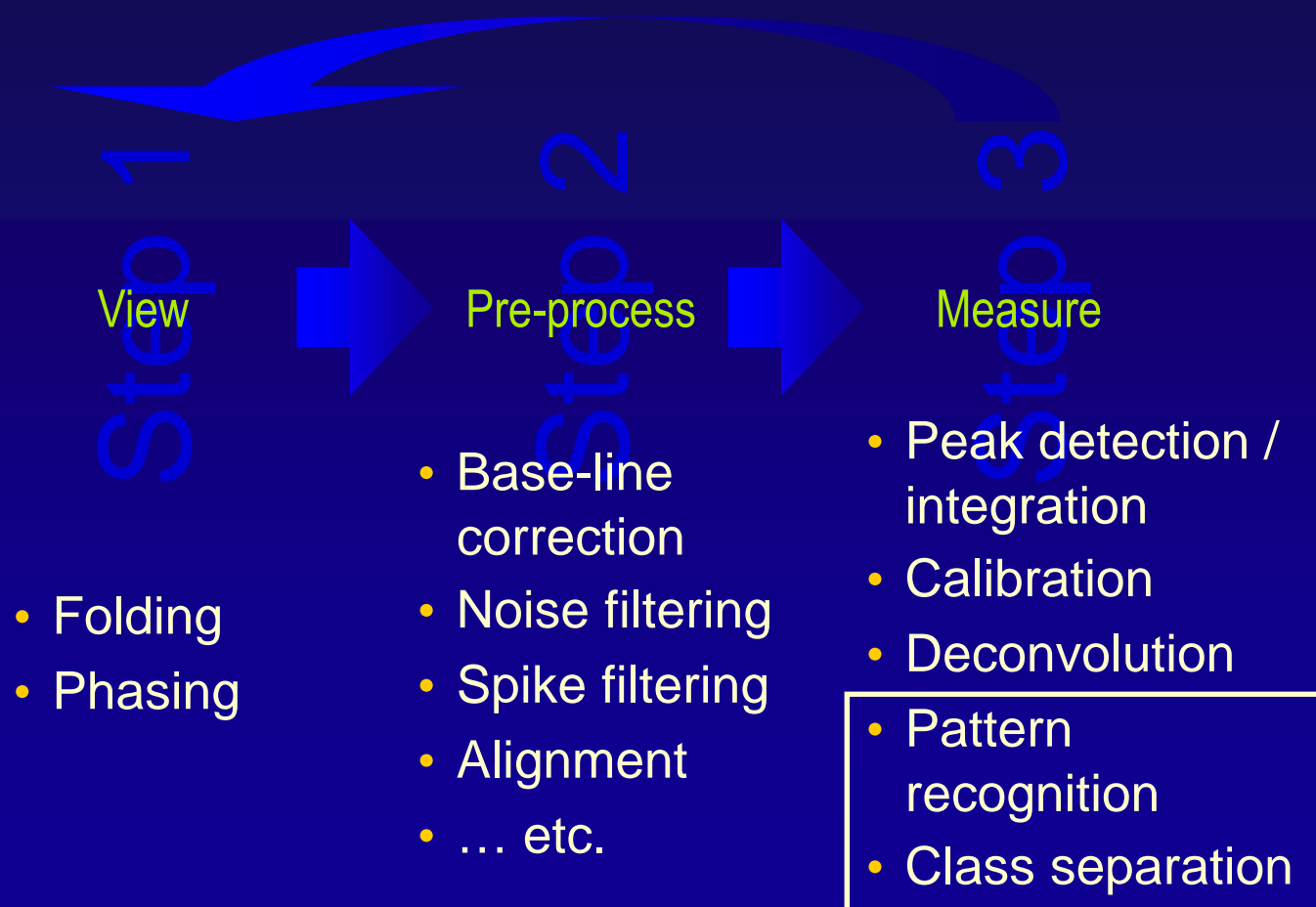Step 1

**View**

- Folding
- Phasing

Step 2

**Pre-process**

- Base-line correction
- Noise filtering
- Spike filtering
- Alignment
- … etc.

Step 3

**Measure**

- Peak detection / integration
- Calibration
- Deconvolution
- Pattern recognition
- Class separation

Pattern recognition in GCxGC. Supervised methods.

*Pattern recognition*

In supervised pattern recognition of GCxGC, a tremendous reduction of variables is performed (form millions to a few tens/hundreds)

*Any method will be prone to overfitting*

Any variable pre-reduction (e.g. using Fisher ratios) should be done within a cross-validation loop

Otherwise the results will be optimistic (a method that seems to work, when in fact it only works for that data)

## Pattern recognition in GCxGC. Example of a wrong strategy    *Pattern recognition*

Objective: discovering metabolites responsible for cancer tumor

Obtain GCxGC chromatograms for sick (50) and healthy (50)

→

Variable selection: Fisher ratio on the raw data

→

Supervised pattern recognition: PLS-DA to separate sick from healthy

→

Consider the coefficients from PLS-DA as indicators of potential metabolites

*Probably align GCxGC data*

*Keep only variables with a FR>threshold*

*Use only the selected variables from step 2*

*Hurrah! I have a collection of interesting metabolites!*

**Aren't you Overfitting?**

→ No, I've been cross-validating the PLS-DA

… but the variable pre-selection has been done with the full data set!!

## Conclusions

- Deconvolution: normally done with the unfolded data (less problems with between-modulation alignment)

- Deconvolution: problem to establish the number of compounds (normally done in a manual way)

- Two ways for pattern recognition: with raw data (normally preferred) or with peak table.

- Careful with validation of supervised pattern recognition. Variable pre-selection should be included in the validation loop.