# How We Handle Mass Spectra
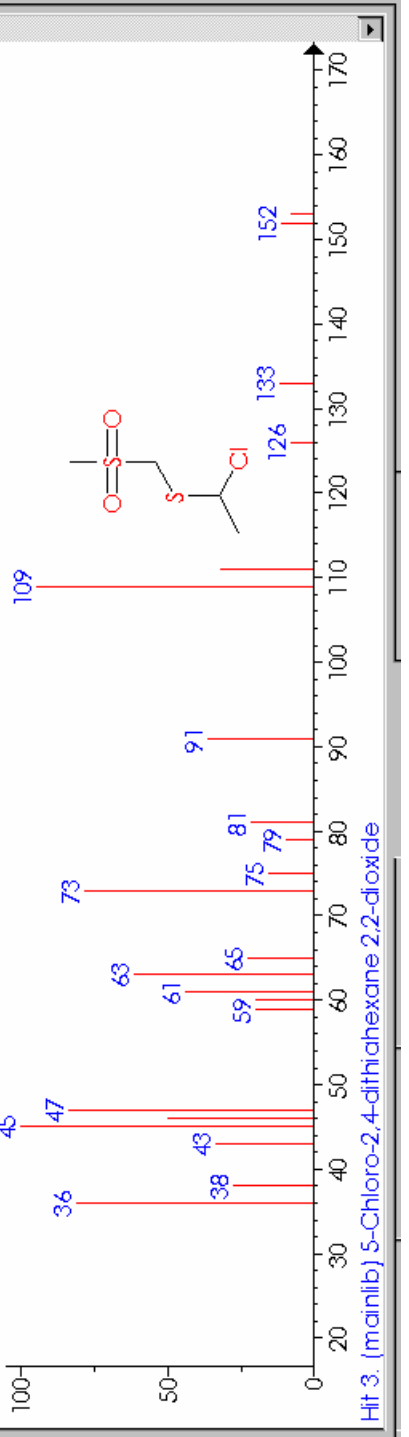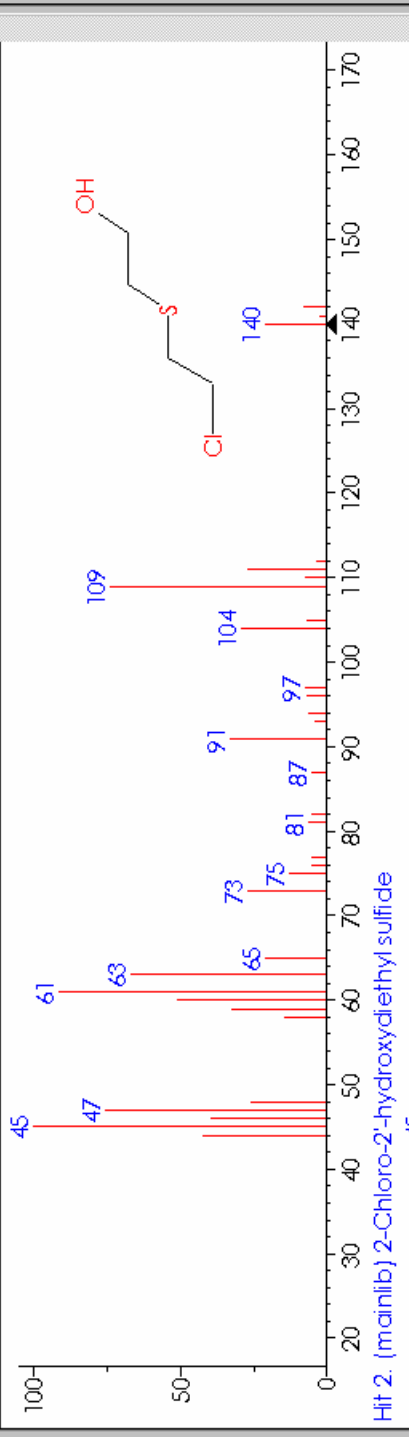
## NIST Mass Spectrometry Data Center
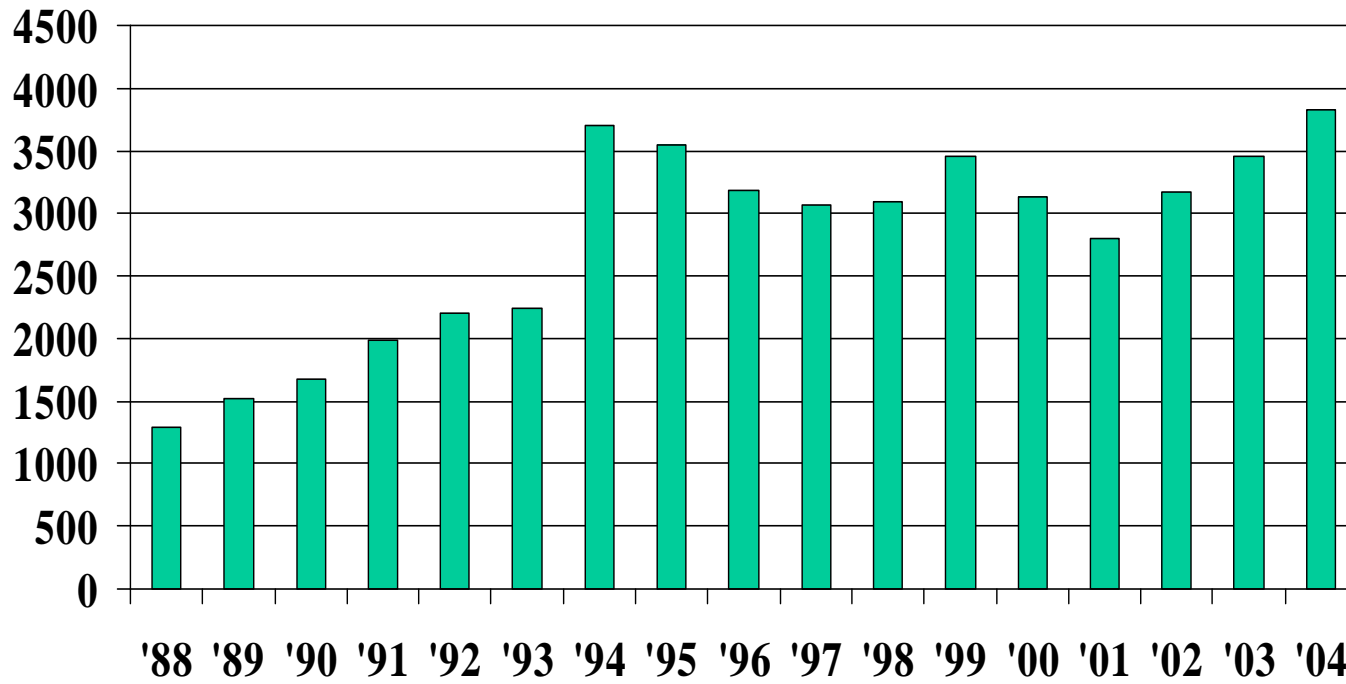
Hit 1. (mainlib) Mustard Gas

Hit 2. (mainlib) 2-Chloro-2'-hydroxydiethyl sulfide

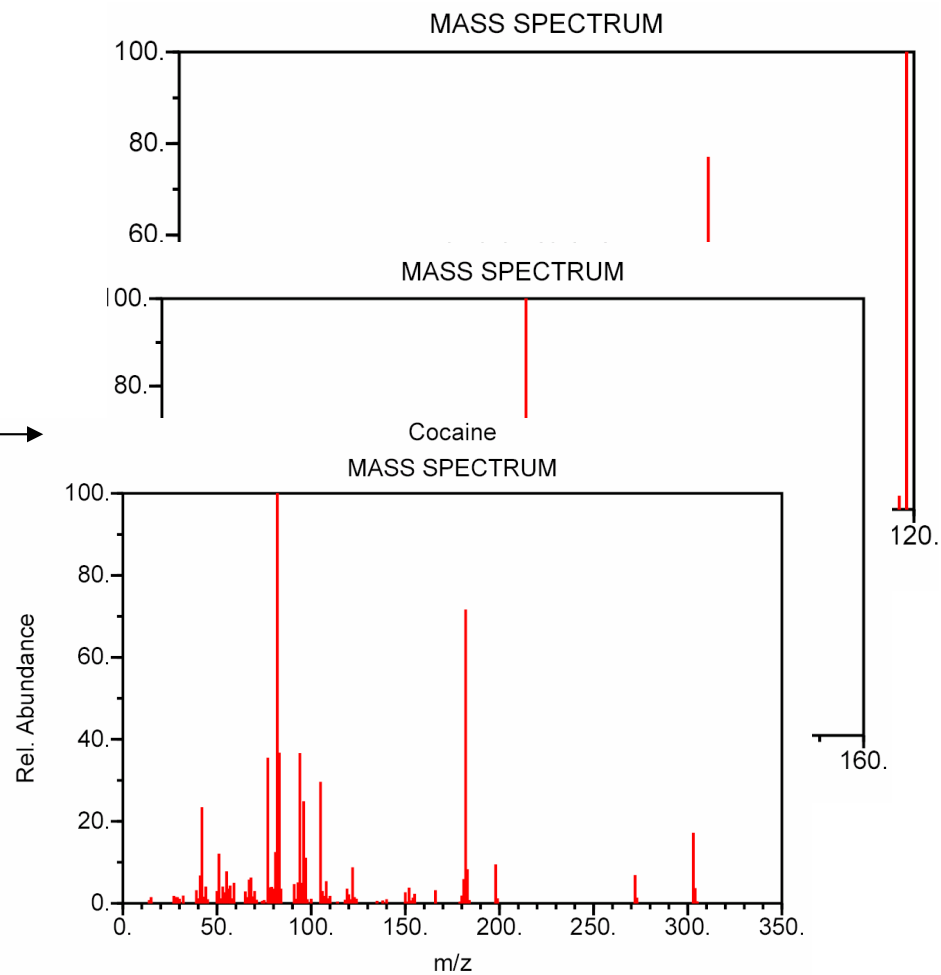Hit 3. (mainlib) 5-Chloro-2,4-dithiahexane 2,2-dioxide
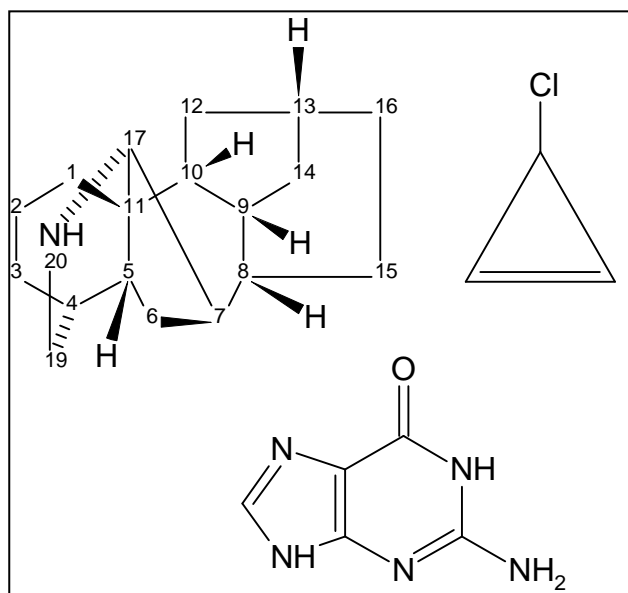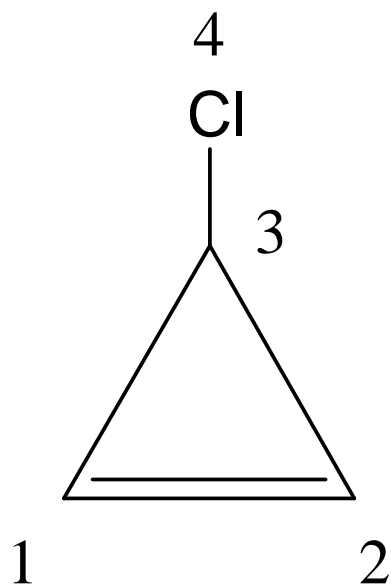
# NIST/EPA/NIH Mass Spectral Library

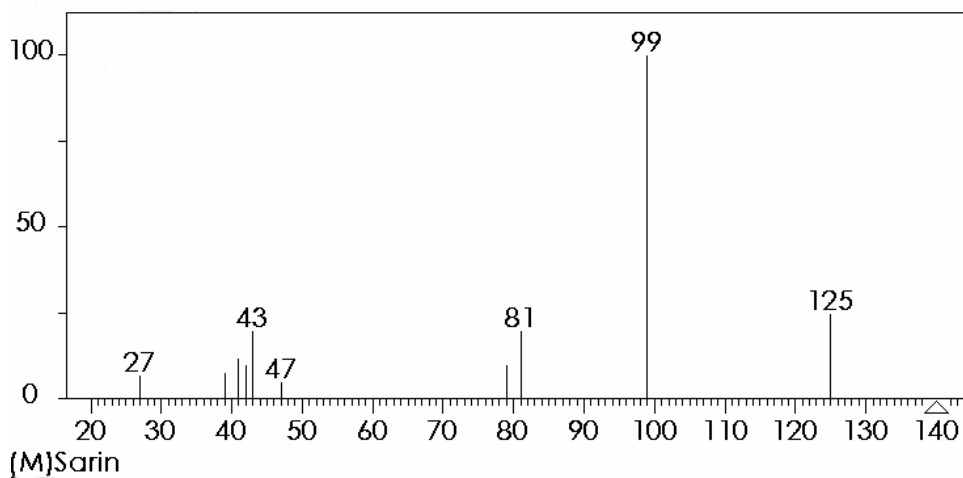## Numbers of Spectra
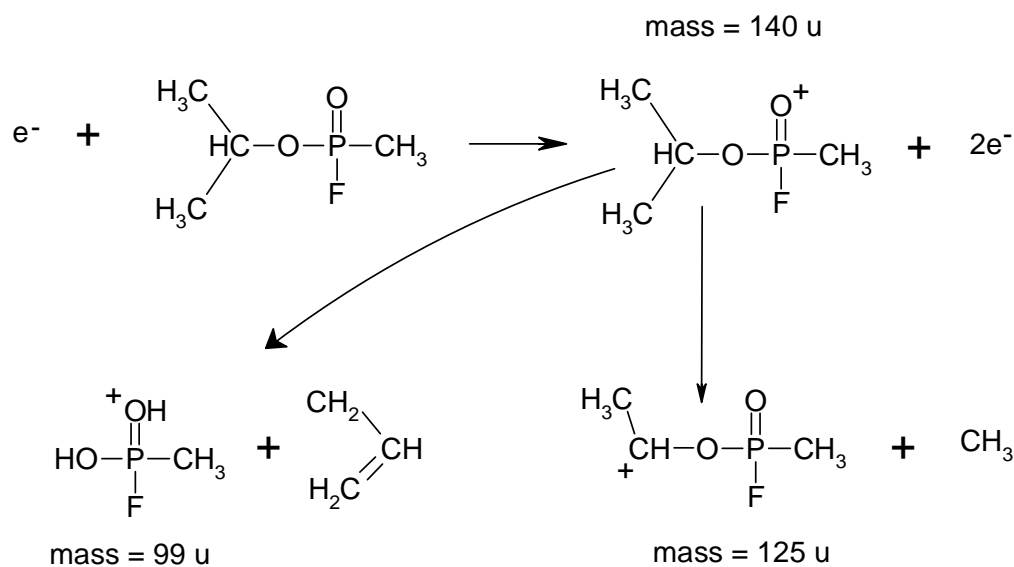
# Libraries Distributed/Year

# The Data

# Connection Table



|   | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 |   | D | S |   |
| 2 | D |   | S |   |
| 3 | S | S |   | S |
| 4 |   |   |   | S |

# From Structure to Spectrum: A Mass "Fragmentogram"

mass = 140 u



mass = 99 u                    mass = 125 u



(M)Sarin

# Molecular Fingerprints



(M)Phosphonothioic acid, methyl-, S-[2-[bis(1-methylethyl)amino]ethyl] o-ethyl ester

(M)Bis(2-chloroethyl) sulphide

(M)Sarin

# I will discuss

- Library Searching
  - Full and Partial Spectra
- Spectrum Purification
- Chemical Structure Representation
- Peptide Spectra Libraries

# Instrument 'Noise Signature'

## 250 Hexachlorobenzene Spectra
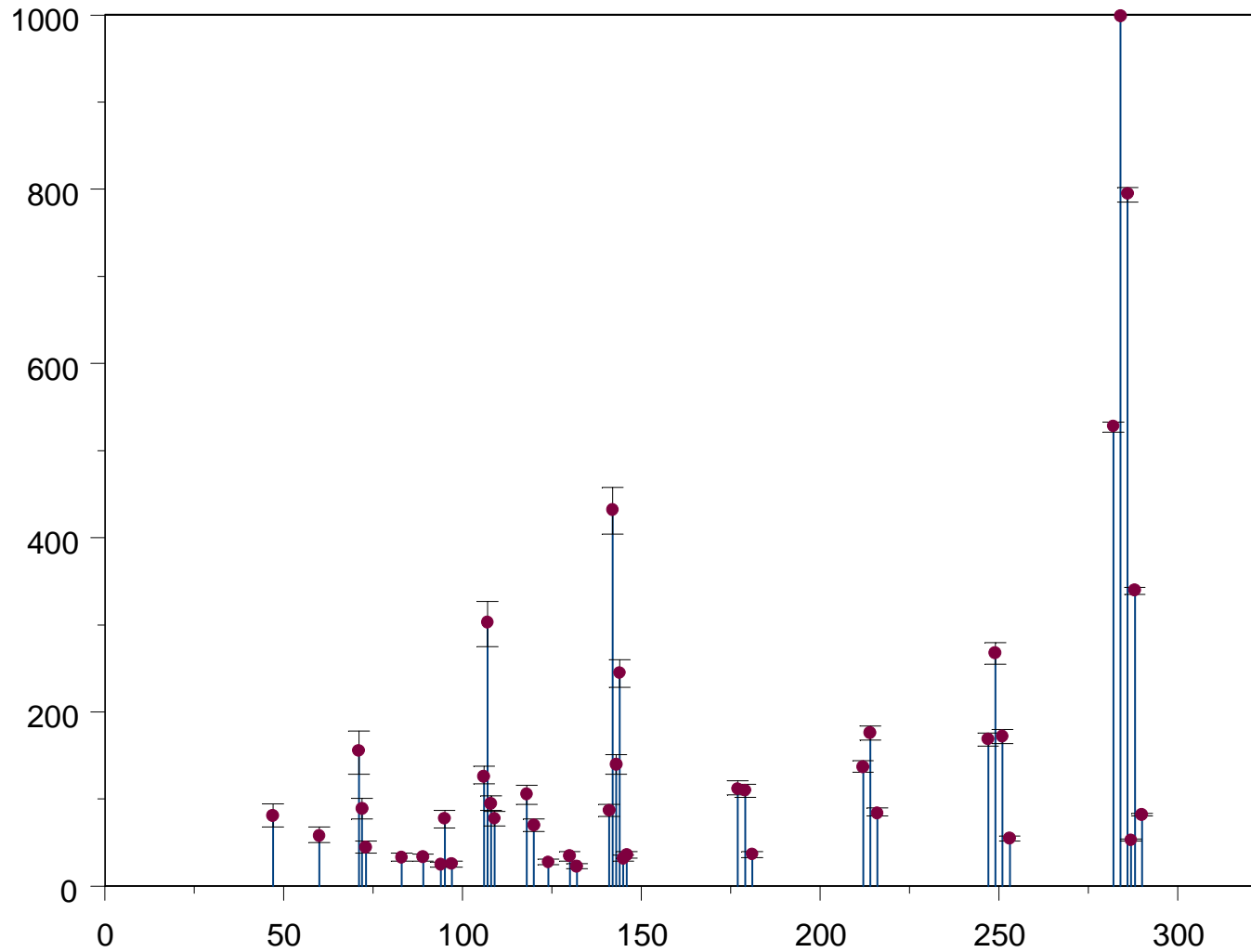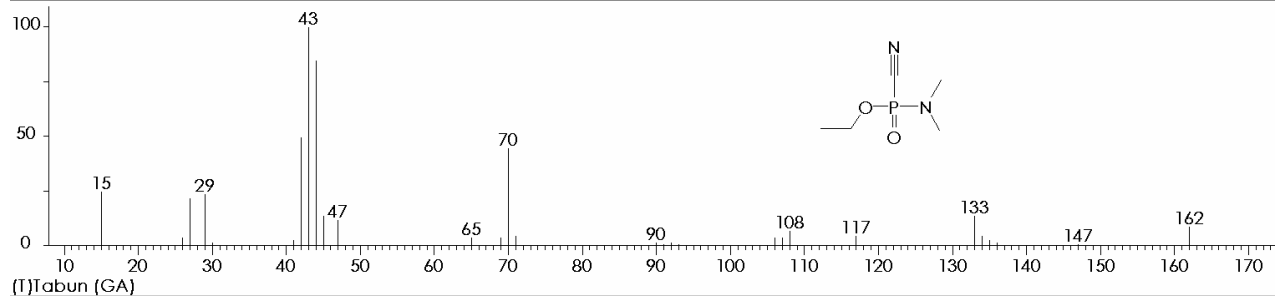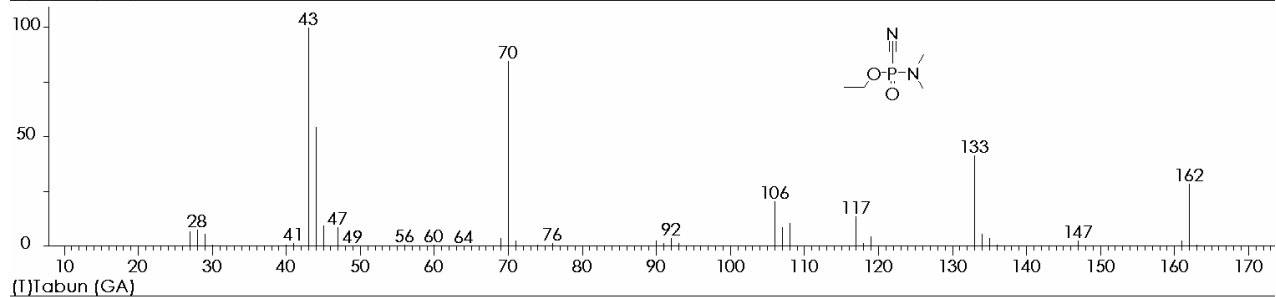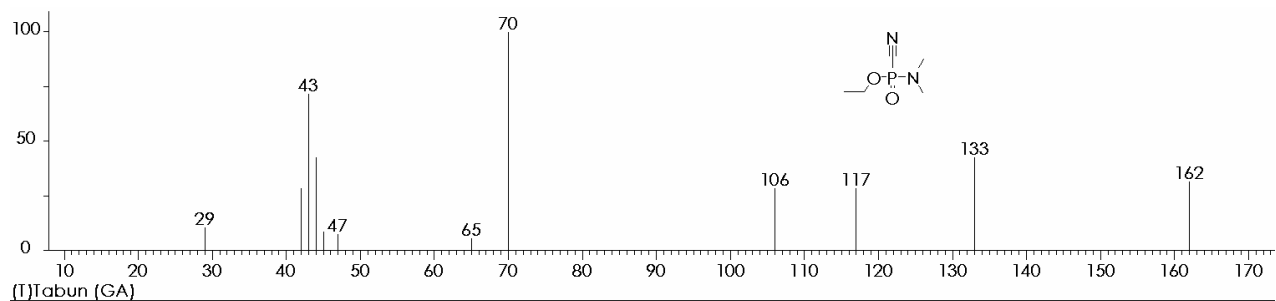## same instrument, calibration mix

Bars show
quartiles

# Instrument Effects



(T)Tabun (GA)

(T)Tabun (GA)

(T)Tabun (GA)

# Library Search

# Spectral Similarity

$$\frac{\sum \sqrt{MR}}{\sqrt{\sum M \sum R}}$$
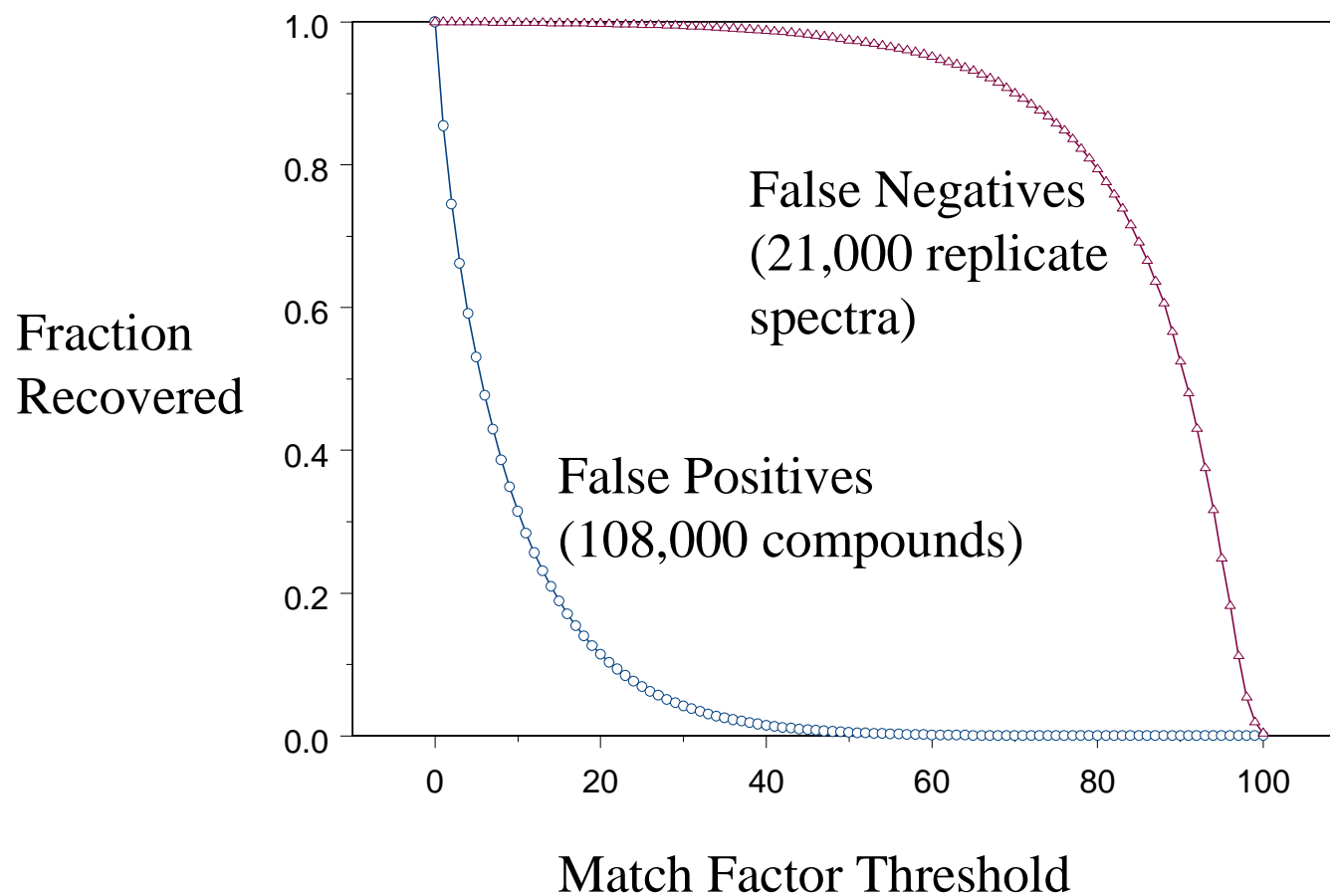
- M = $f$(Abundance) Peak in Measured Spectrum
- R = $f$(Abundance) Peak in Reference Spectrum
- Sum over all peaks
- $f$(Abundance)
  - Abundance
  - Abundance * m/z
  - Certainty

# Algorithm Performance

12,592 Replicate Spectra against NIST Library

| Model | Percent Correct | | |
|---|---|---|---|
| | **Top Hit** | **Top 2 Hits** | **Top 3 Hits** |
| **Correlation – Weighted** | 74.9 | 86.9 | 91.7 |
| **Correlation** | 72.9 | 85.9 | 90.8 |
| **Euclidean Distance** | 71.9 | 83.9 | 88.9 |
| **Absolute Distance** | 67.9 | 80.3 | 85.5 |
| **PBM - Published** | 64.7 | 78.4 | 84.8 |
| **Hites/Hertz/Biemann** | 64.4 | 77.2 | 83.2 |

FP/FP Above Given Match Factor
for NIST Library Spectra

FP/FN
Expanded View

m/z weighting

Fraction
Recovered

no weighting

FN

FP x 10,000

Match Factor

# FP Depends on Spectrum Uniqueness



HCB

DMPB

TMB    decalin    decane

malathion    sarin

FP

Match Factor

HCB = hexachlorobenzene
DMPB = dimethylpenobarbital
TMB = 1,2,3-trimethylbenzene
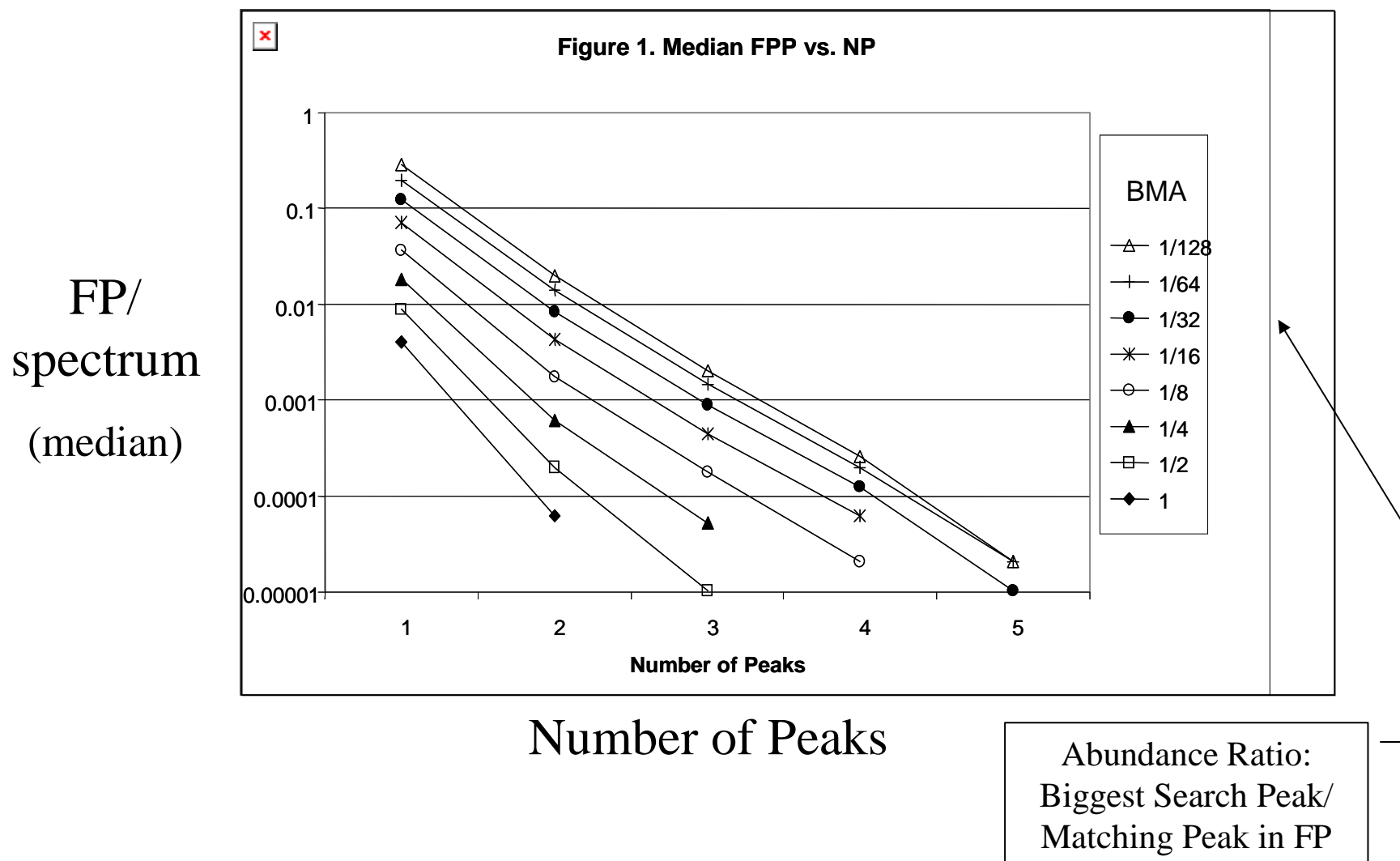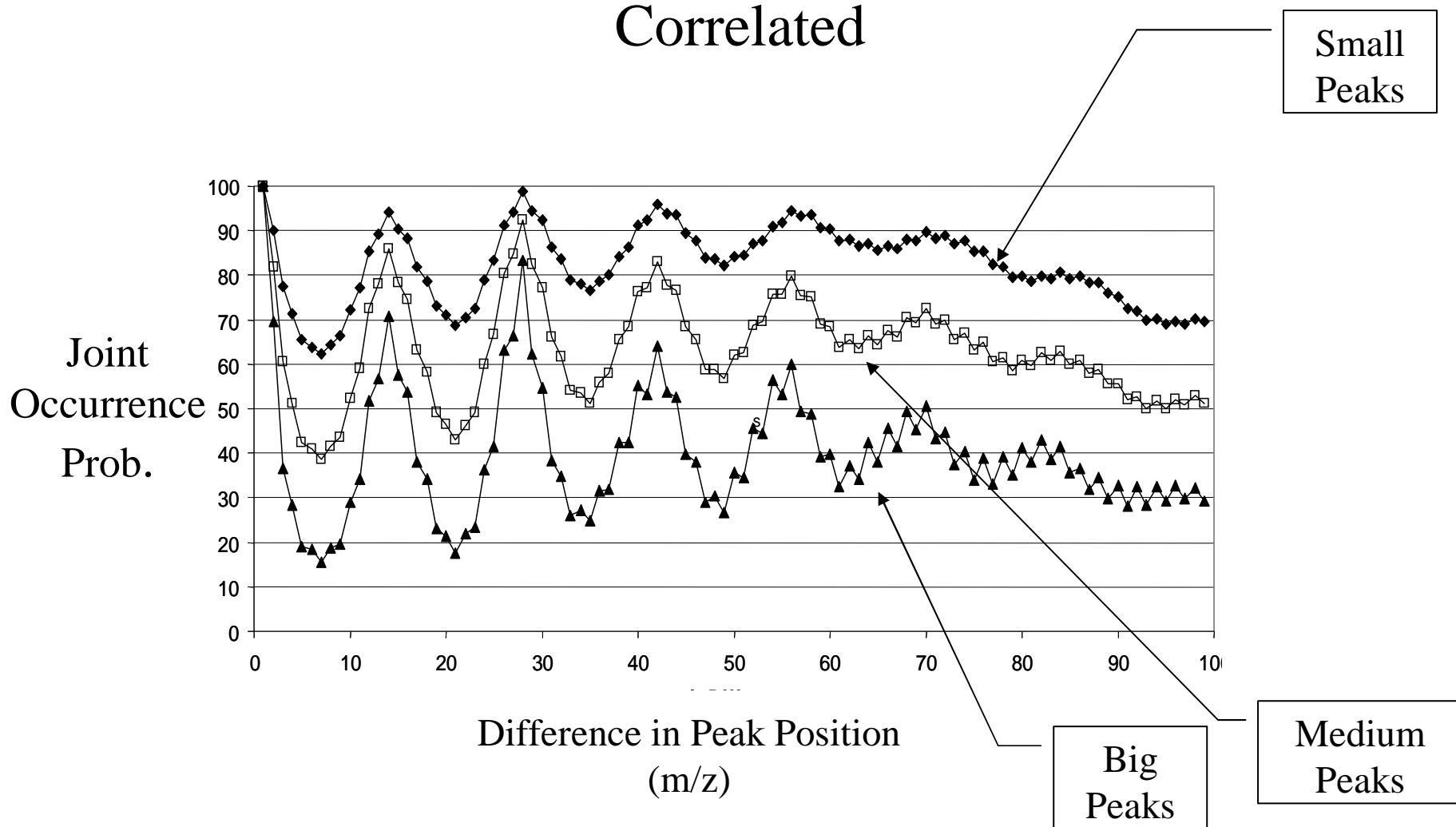
# Multiple Ion Monitoring

- ## What is is?
  - Use 2-5 Major Peaks in Spectrum of Target
    - 10 – 100 more sensitive

- ## What's the problem?
  - Can match major Target peaks with Minor Sample Peaks

- ## What we have done:
  - Examine risk using library as source of potential false positive IDs
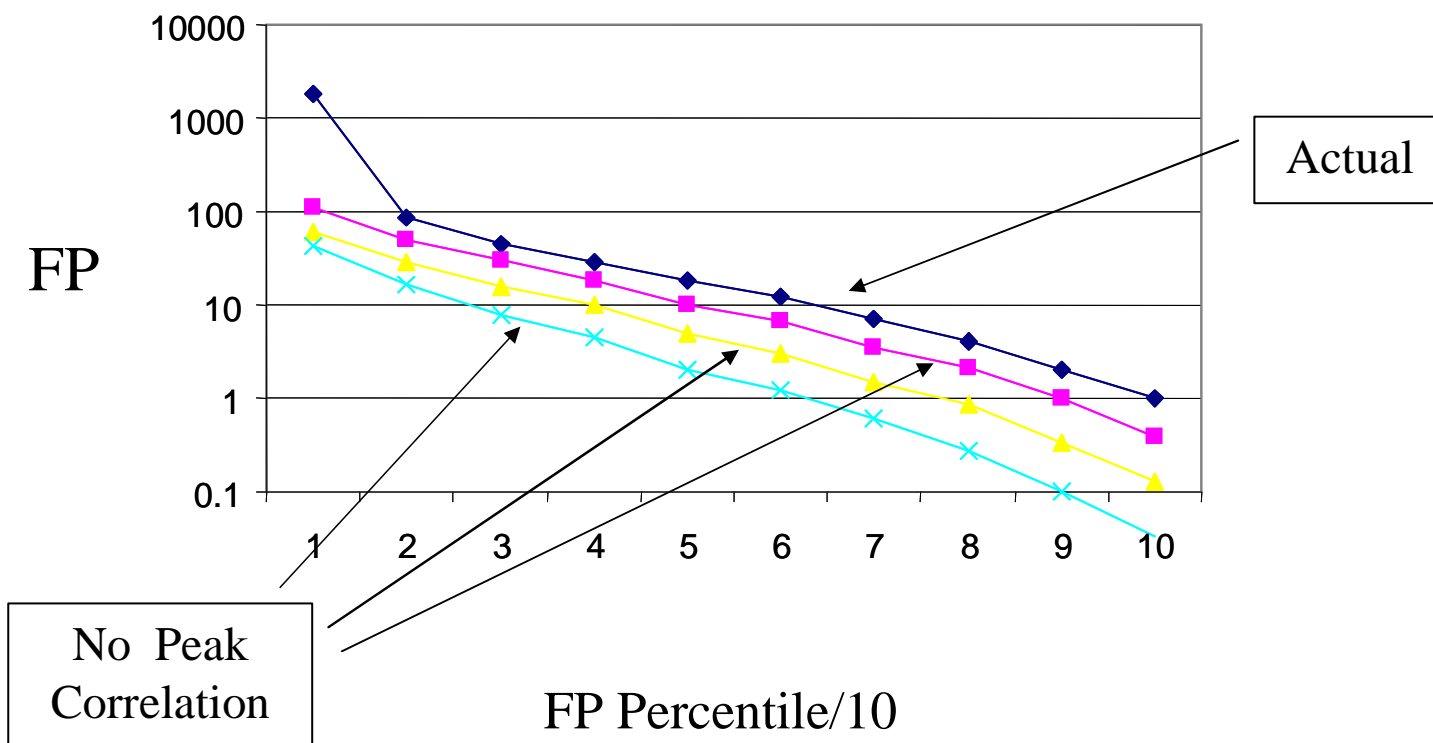
# False Positive Risk vs Number of Peaks Used



Figure 1. Median FPP vs. NP

FP/ spectrum (median)

Number of Peaks

BMA
- 1/128
- 1/64
- 1/32
- 1/16
- 1/8
- 1/4
- 1/2
- 1

Abundance Ratio:
Biggest Search Peak/
Matching Peak in FP
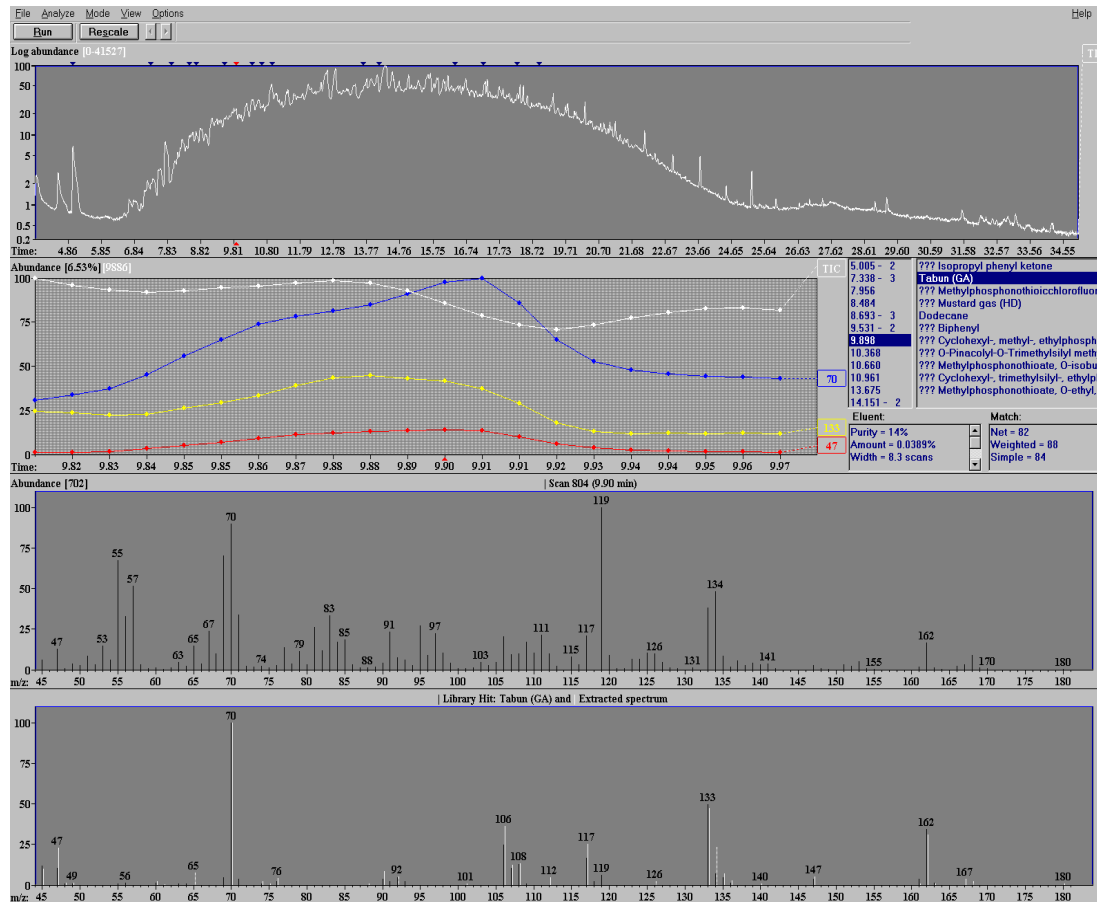
Mass Spectral Peak Occurrences are Correlated

# FP Observed and Computed
## (from individual peak probabilities)
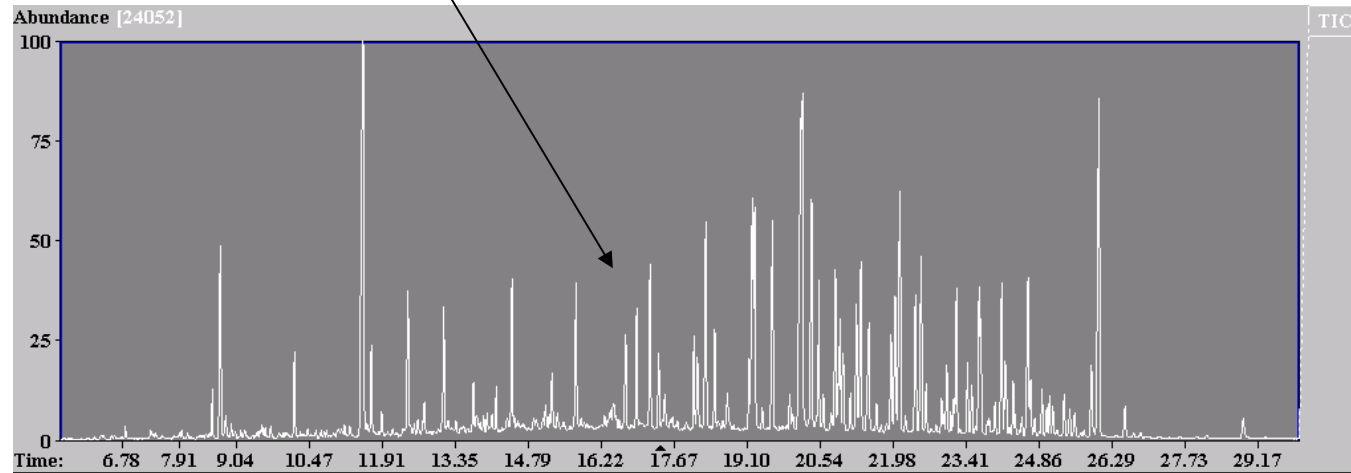
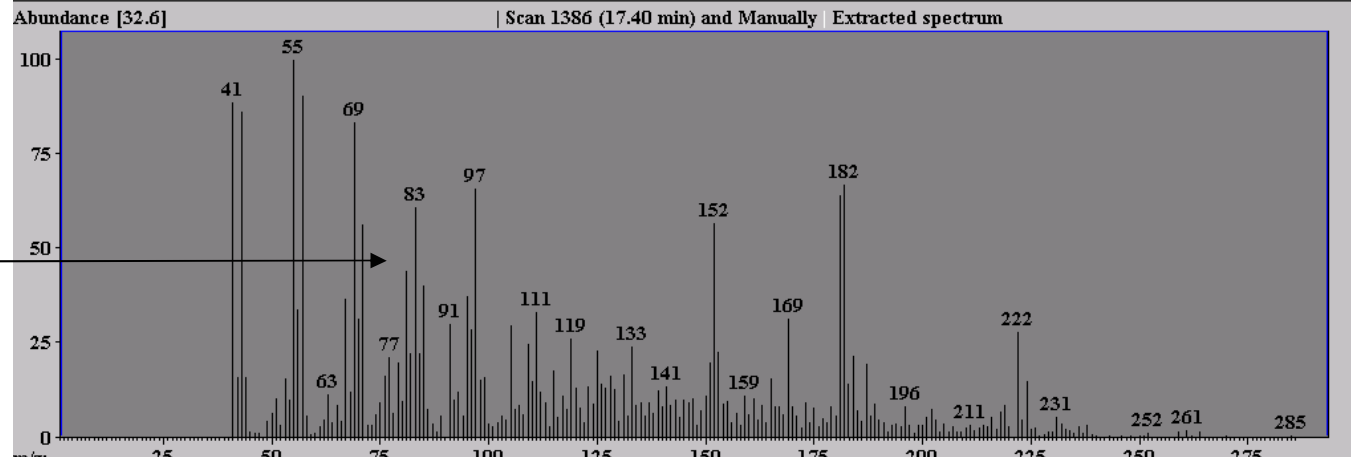# Search Results Depend on Search Spectrum Quality
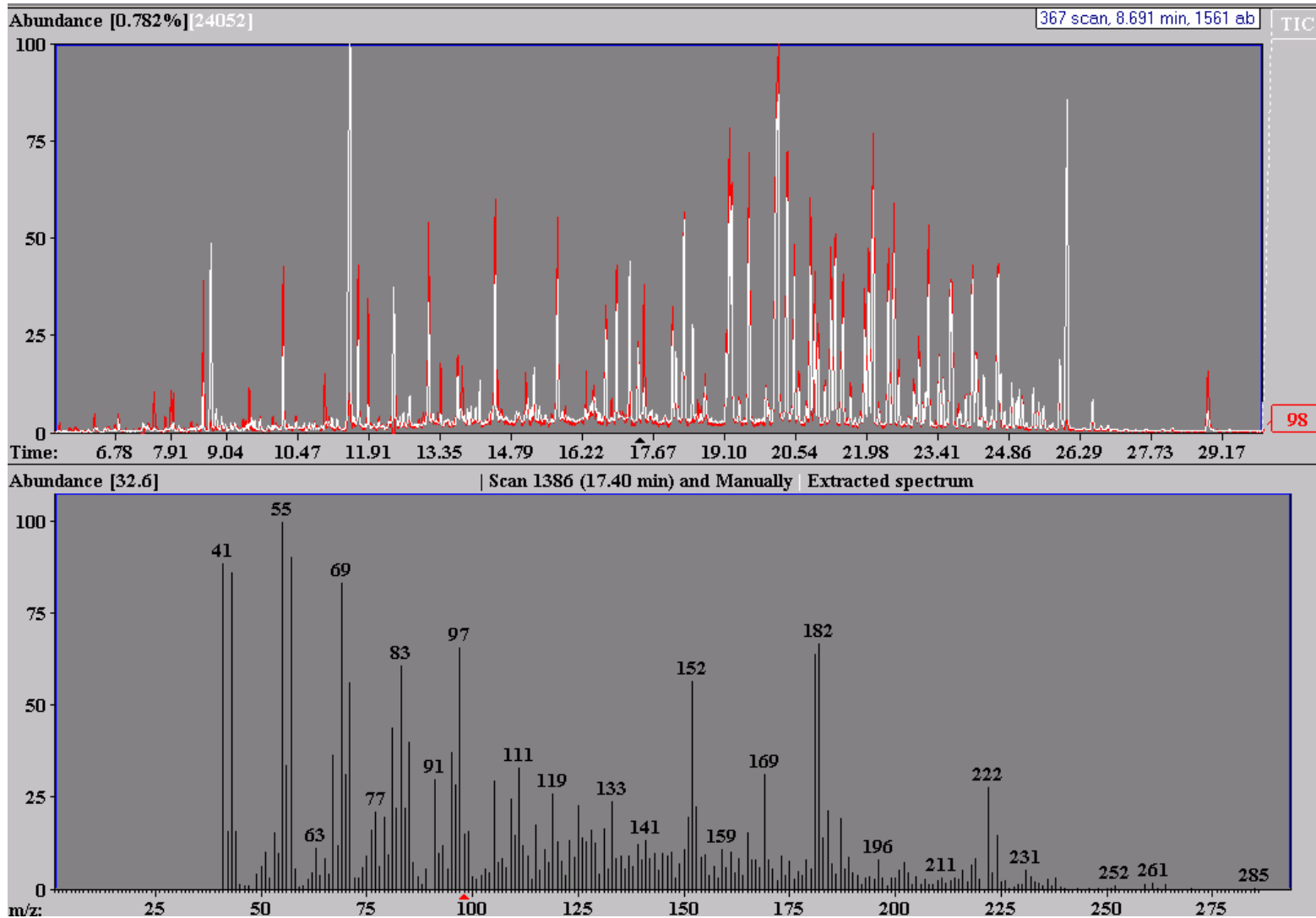


AMDIS:
http://chemdata.nist.gov
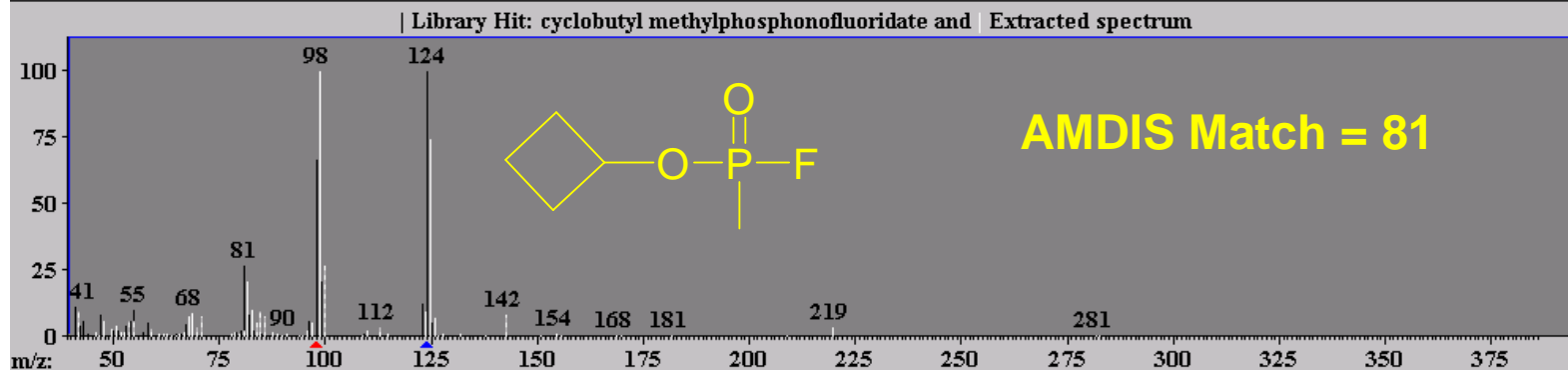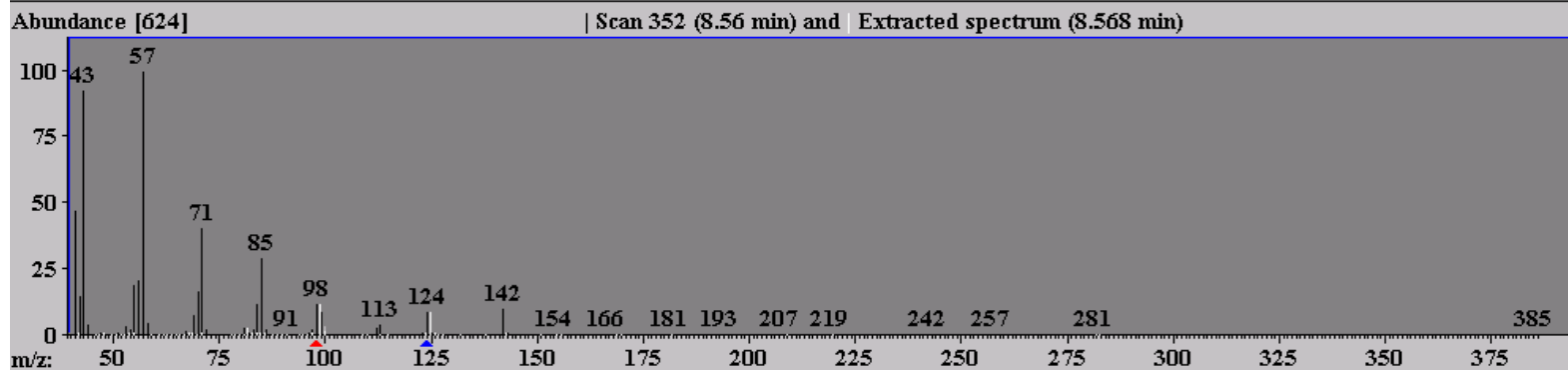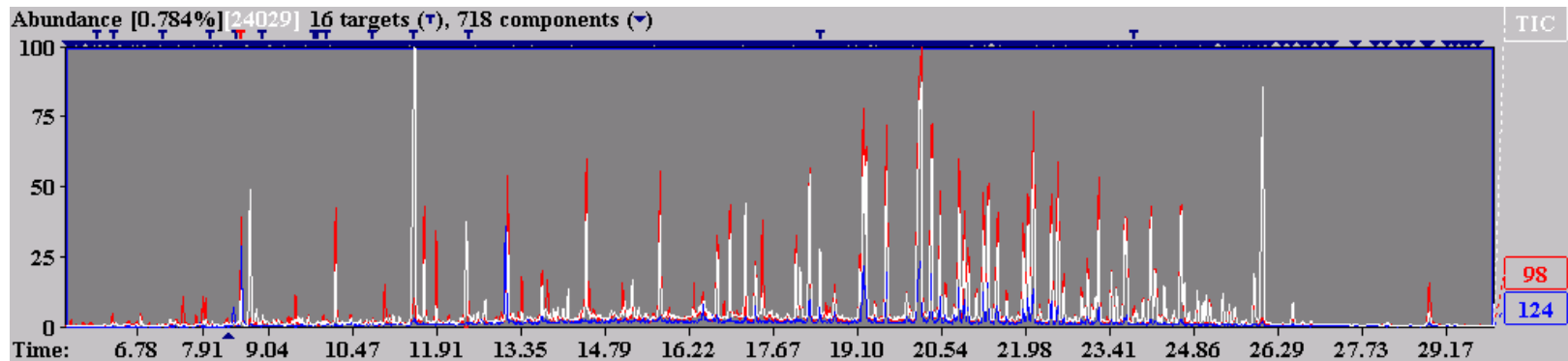
# Real Data

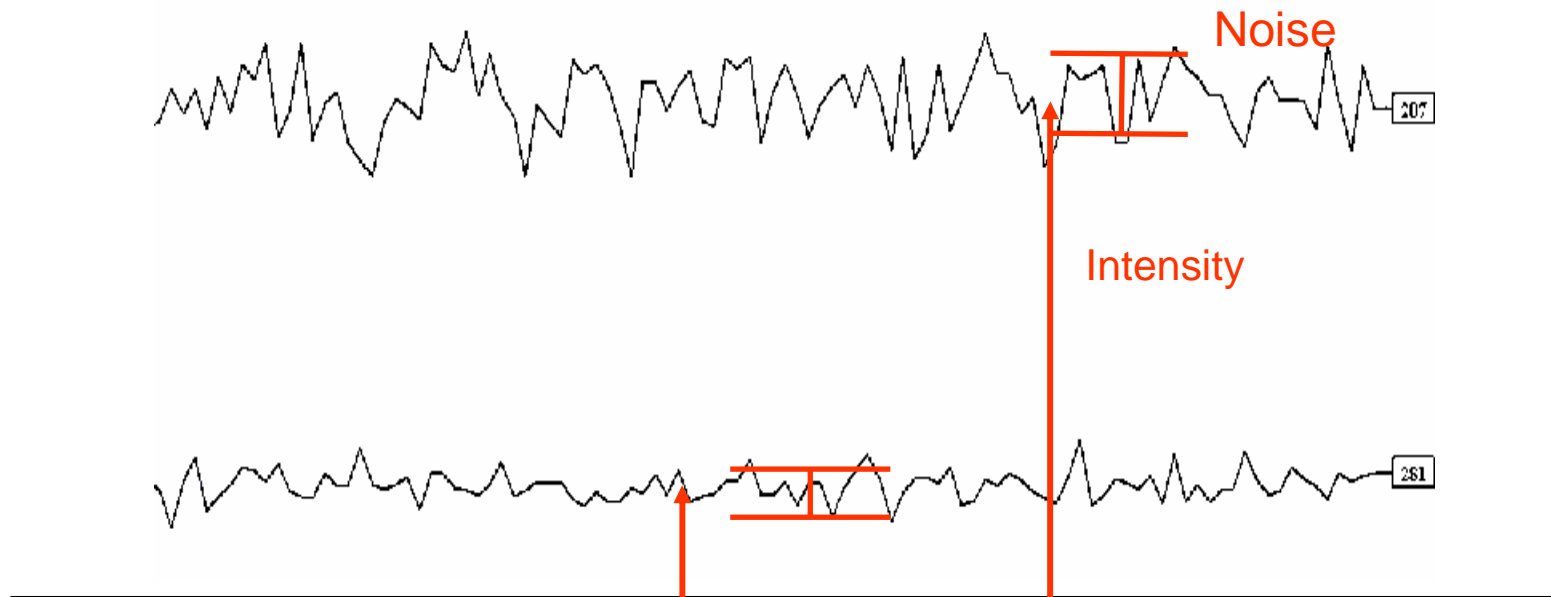Total ion chromatogram

A mass spectrum (scan)

# Chromatogram with single ion

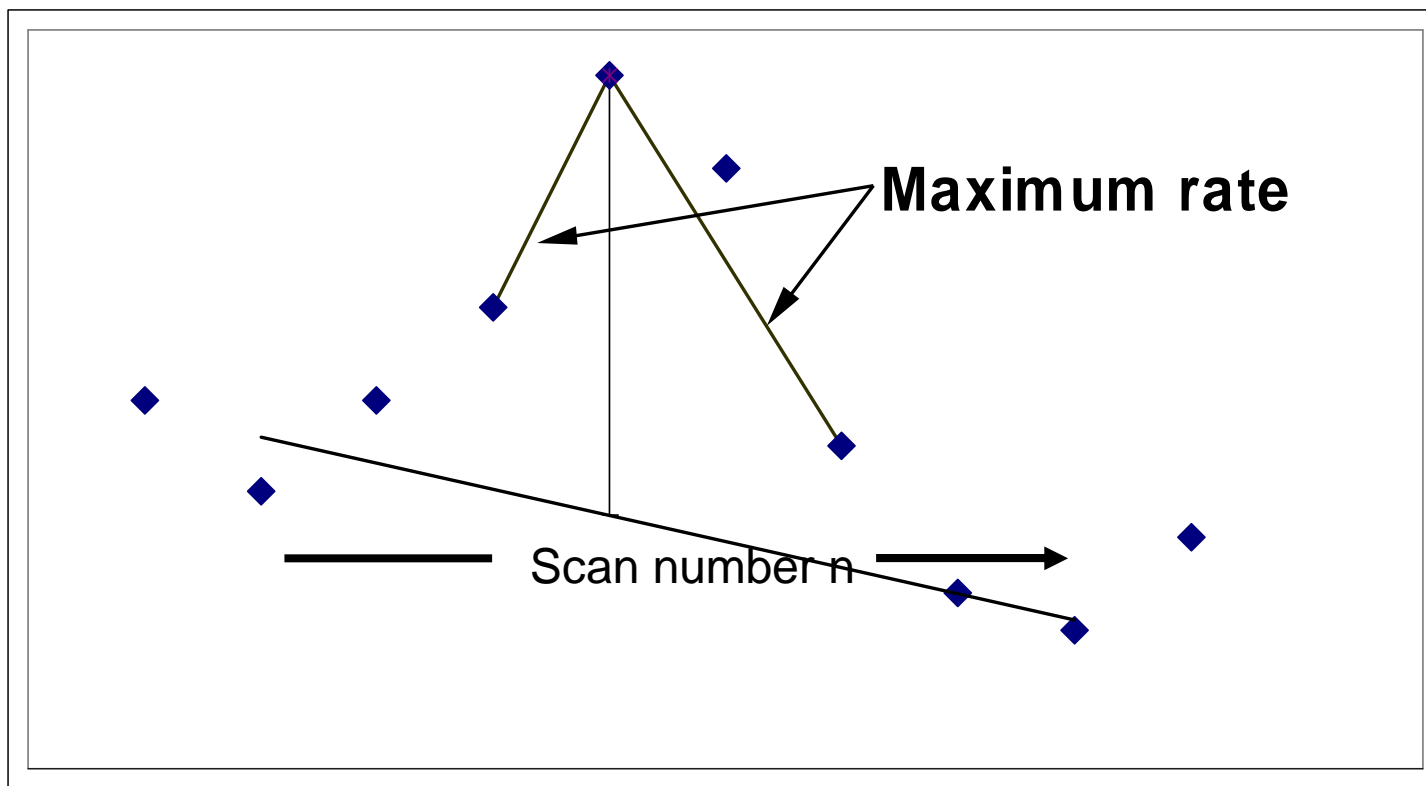# AMDIS Analysis of Data

# Order of Analysis

- Noise Analysis – find 'Noise Factor'
- Find and quantify maximizing ions
- Combine to create 'Model Peak'
- Use Model Peak shape (intensity vs time) to purify spectra
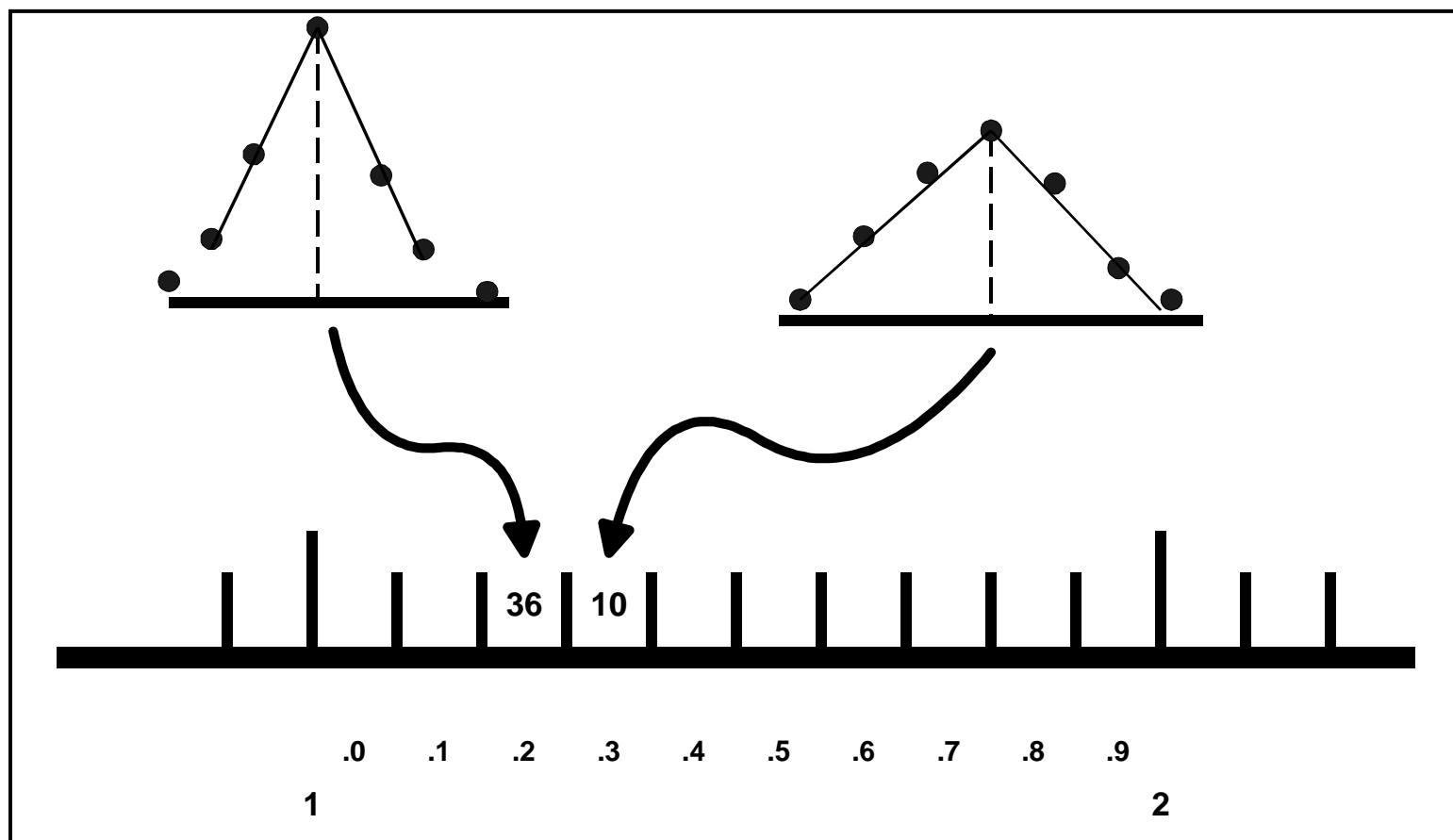- Find best matching library spectrum

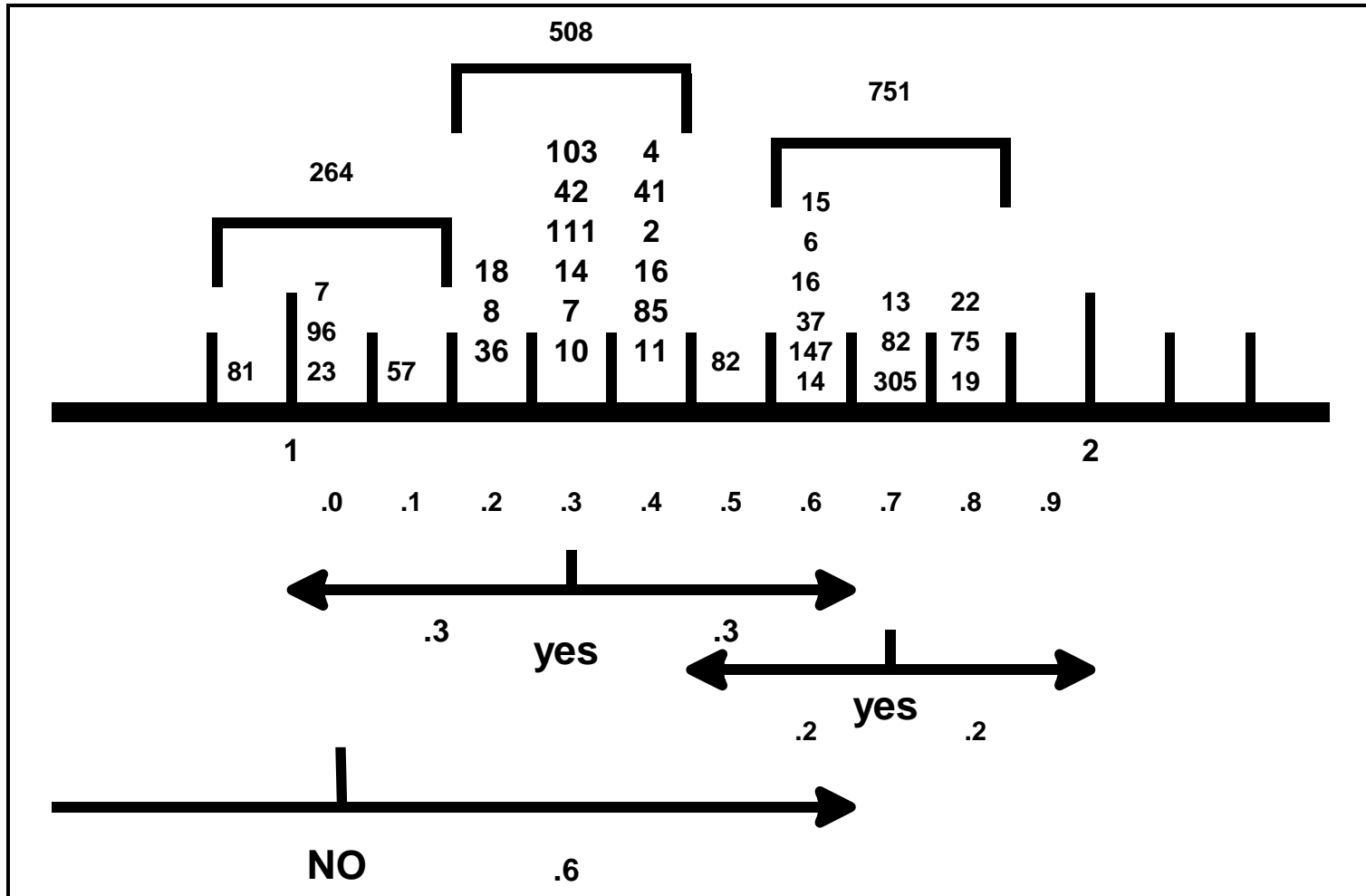# Derive Noise Factor



$$Noise = K_{noise} \sqrt{Intensity}$$

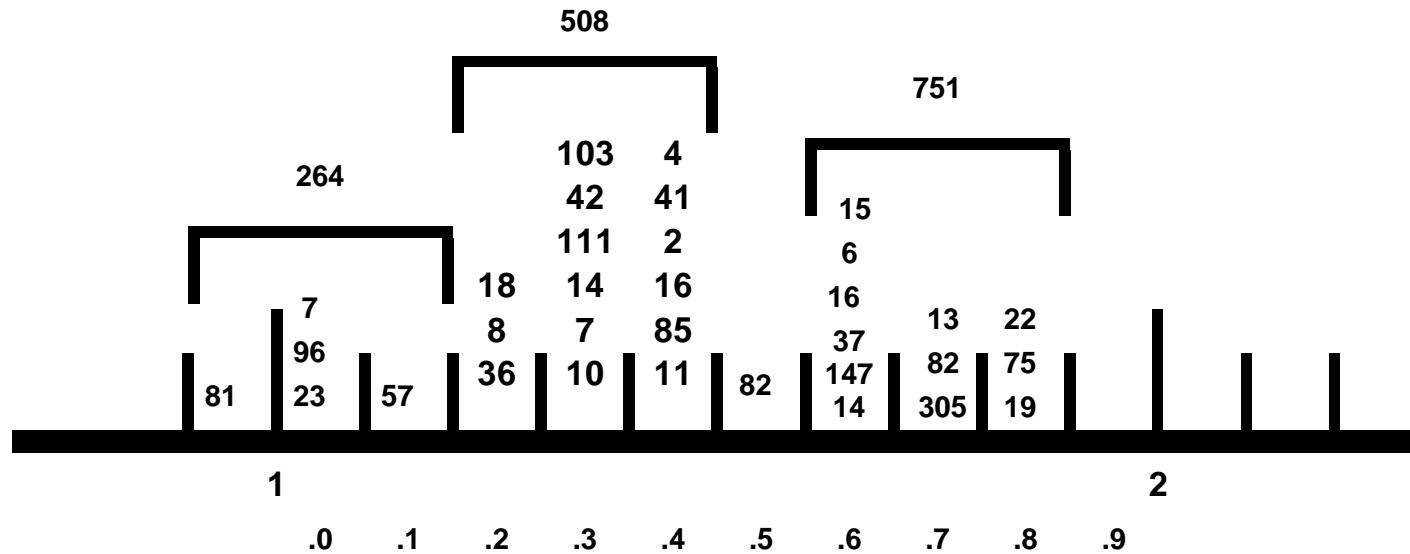# Finding Possible Peaks for Each m/z

# Find Possible Compounds:
# Do Ions Maximize at Same Time?
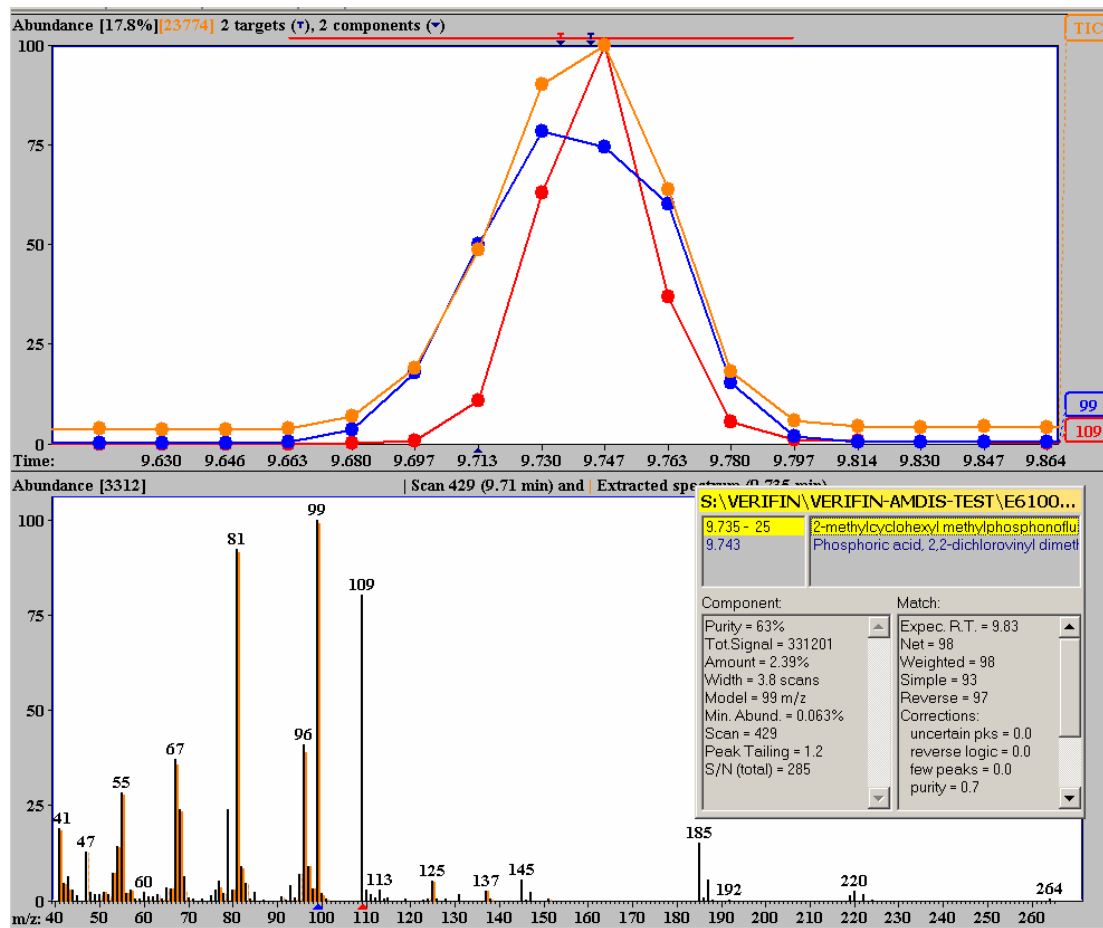
# Separate the Components

# A 'Model Peak' Provides Shape



The model shape is defined as the sum of all of the ion chromatograms that maximize within the range and have a sharpness value within 75% of the maximum.

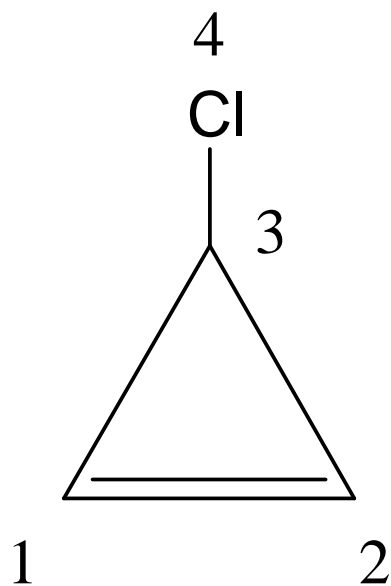# AMDIS Testing – Closely Eluting Components

# Representing Chemical Identity

- Visual: 2D Structure
- Text: IUPAC Name
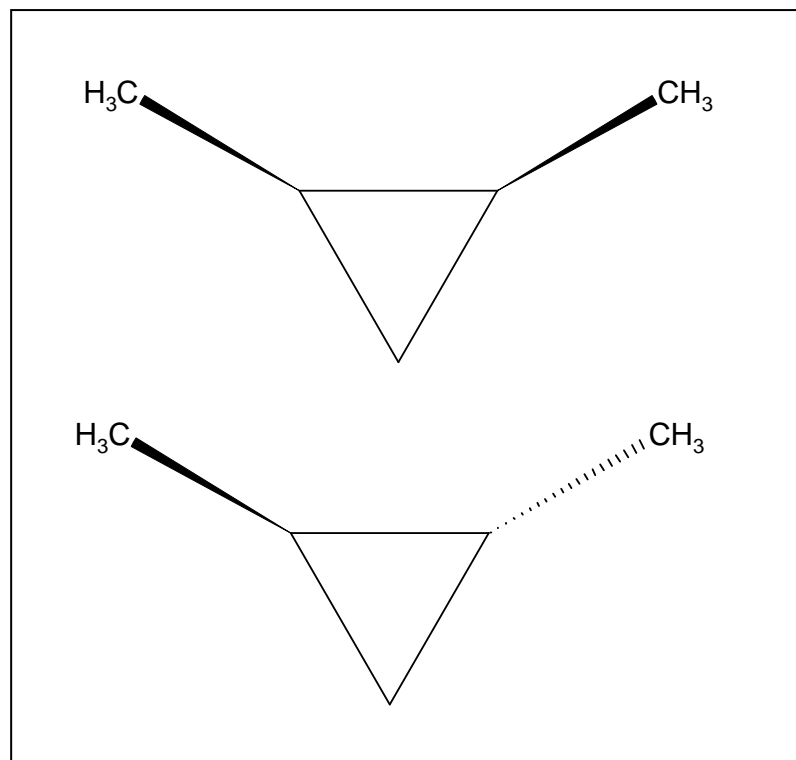- Digital: No Accepted, Open Method

- Solution:

The IUPAC/NIST Chemical Identifier

# Connection Table



|   | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 |   | D | S |   |
| 2 | D |   | S |   |
| 3 | S | S |   | S |
| 4 |   |   |   | S |

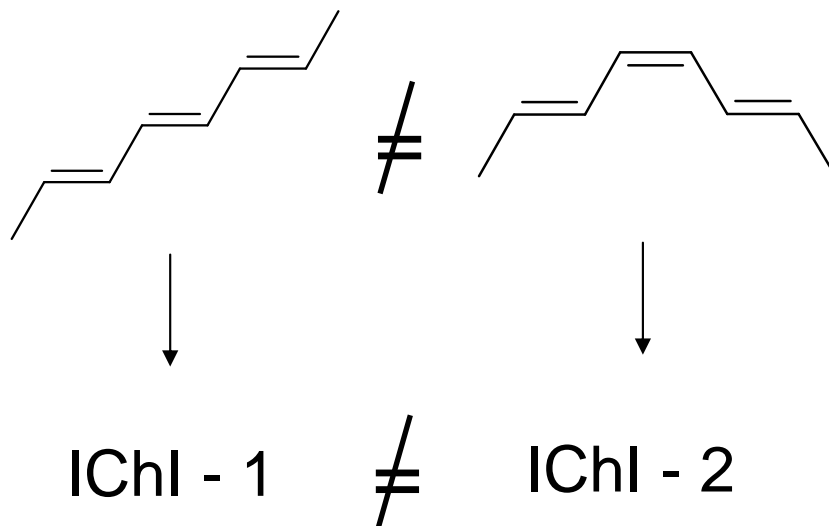# Chemical Identity Problems



Registry Number possible for each exact form,
mixture, unknown, unspecified

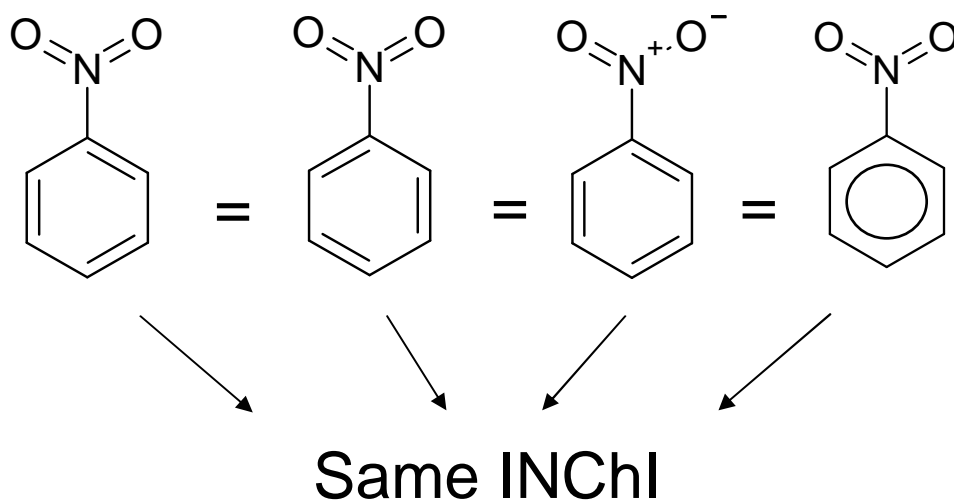Experts required

Expensive, ambiguous and error prone

# Requirements

- Different compounds have different identifiers
  - Keep all distinguishing structural information



IChI - 1 ≠ IChI - 2

# Requirements

- One compound has only one identifier
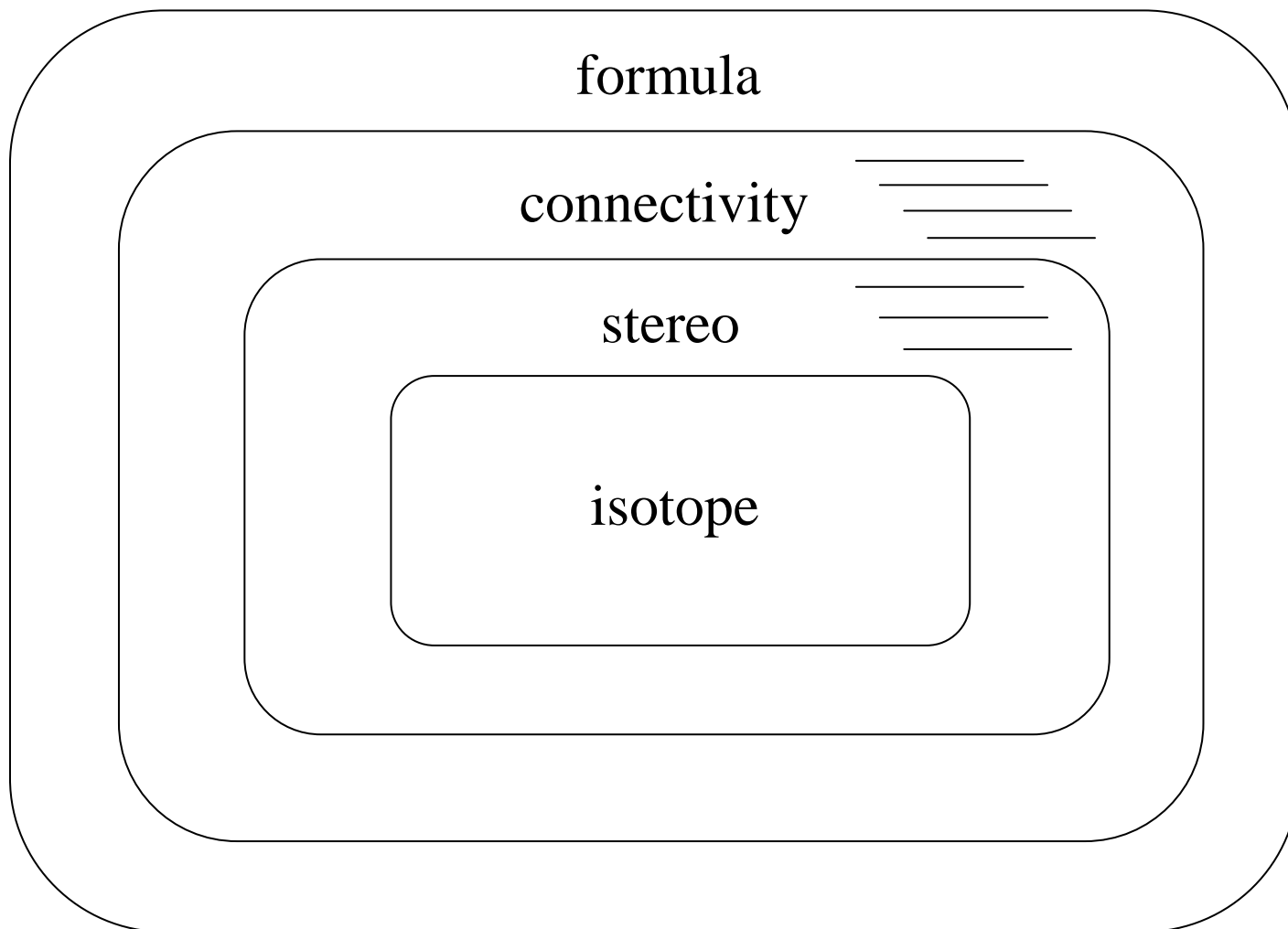  - Omit unnecessary information



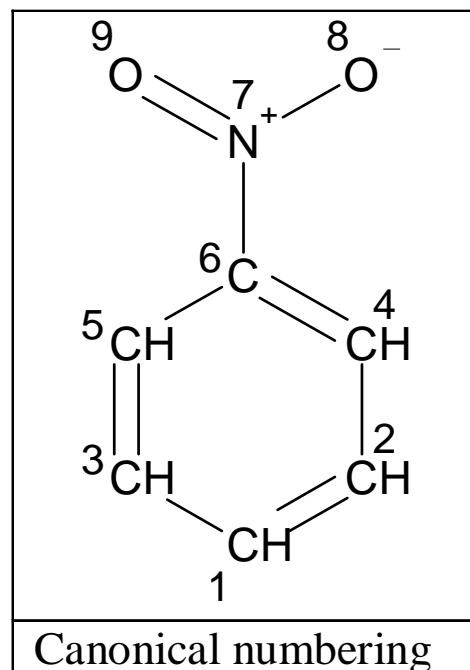Same INChI

# 3 Steps to INChI

- Chemistry
  - 'Normalize' Input Structure
    - Implement chemical rules

- Math
  - 'Canonicalize' (label the atoms)
    - Equivalent atoms get the same label

- Format
  - 'Serialize' Labeled Structure
    - Output as character string ('name')
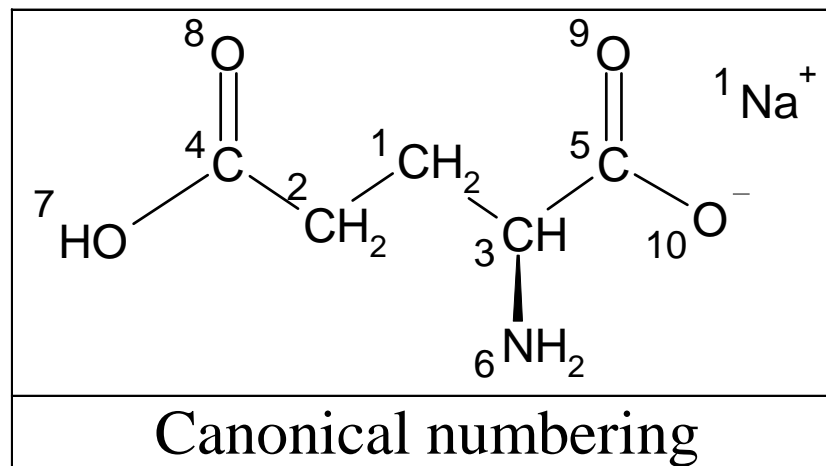
"**Layers**"   Chemical Substances

formula

connectivity

stereo

isotope

# Nitrobenzene



Canonical numbering

| Description | Layers |
|---:|:---|
| formula | C6H5NO2 |
| connectivity | 8-7(9)6-4-2-1-3-5-6 |
| H-atoms | 1-5H |
| charges | |

# MSG



8 O 9 O ¹Na⁺
4 C 1 CH₂ 5 C
7 HO 2 CH₂ 3 CH 10 O⁻
6 NH₂

Canonical numbering

| Description | Layers |
|---|---|
| formula | `C5H8NO4.Na` |
| connectivity | `6-3(5(9)10)1-2-4(7)8;` |
| H-atoms | `1-2H2,3H,6H2(H-,7,8,9,10);` |
| stereo sp$^3$ | `3-;` |
| | |
| charges | `-1;+1` |

C5H9NO4.Na/c6-3(5(9)10)1-2-4(7)8;/h1-2H2,3H,6H2,(H,7,8)(H,9,10);/q;+1/p-1/t3-;/m1./s1

**Input/ Result**

**Mobile H On/Off**
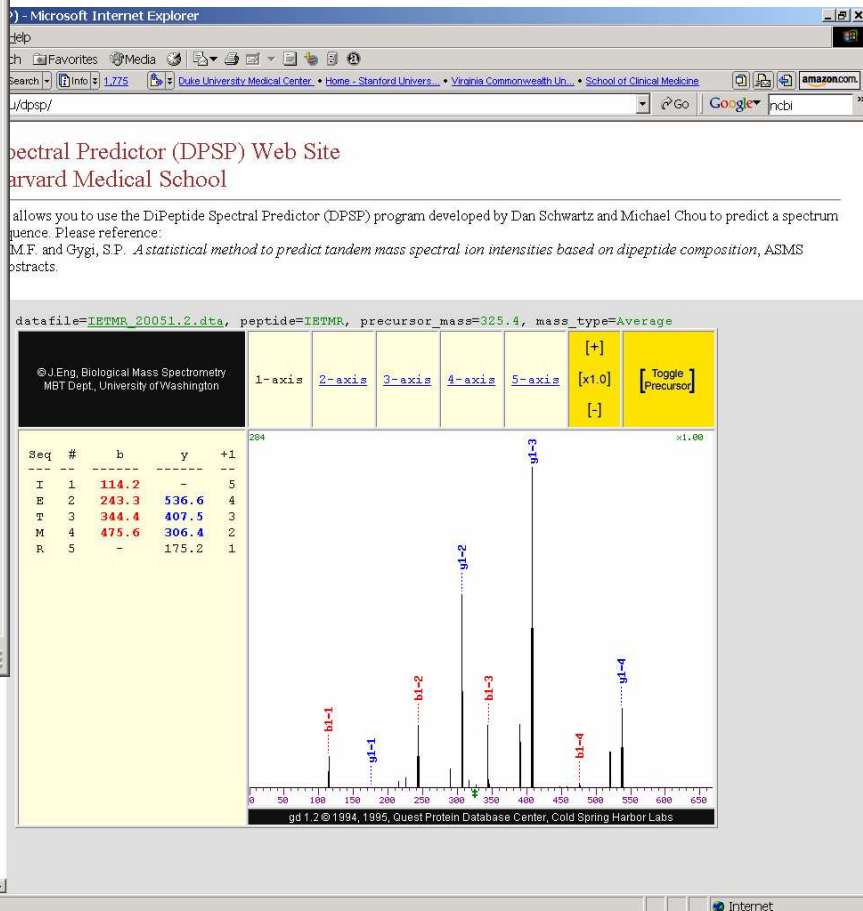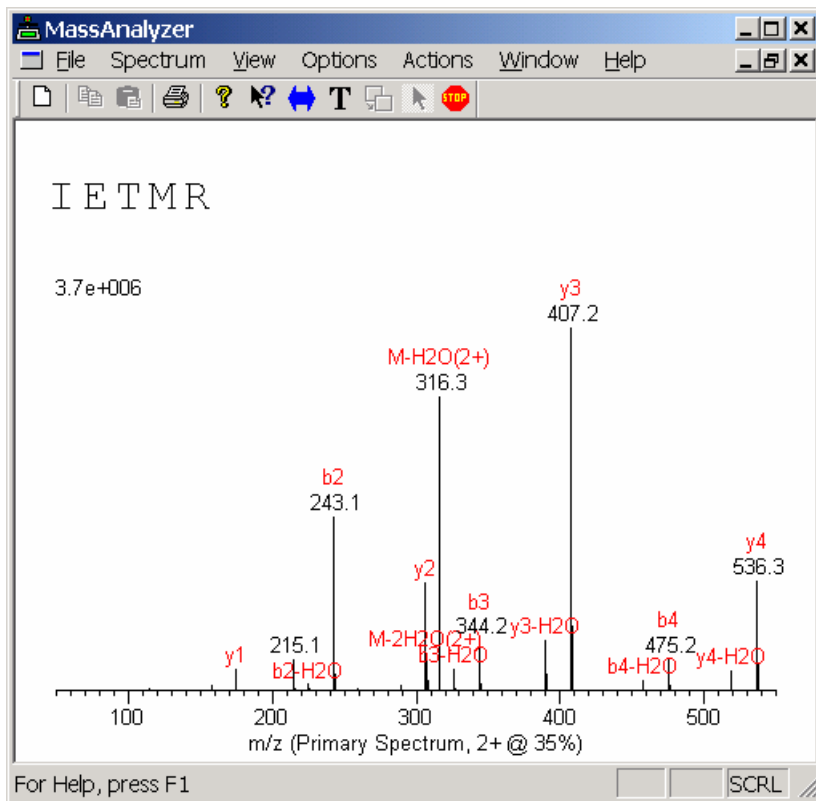
**Include Org-Metal Bonds**

**INChI Test Version**
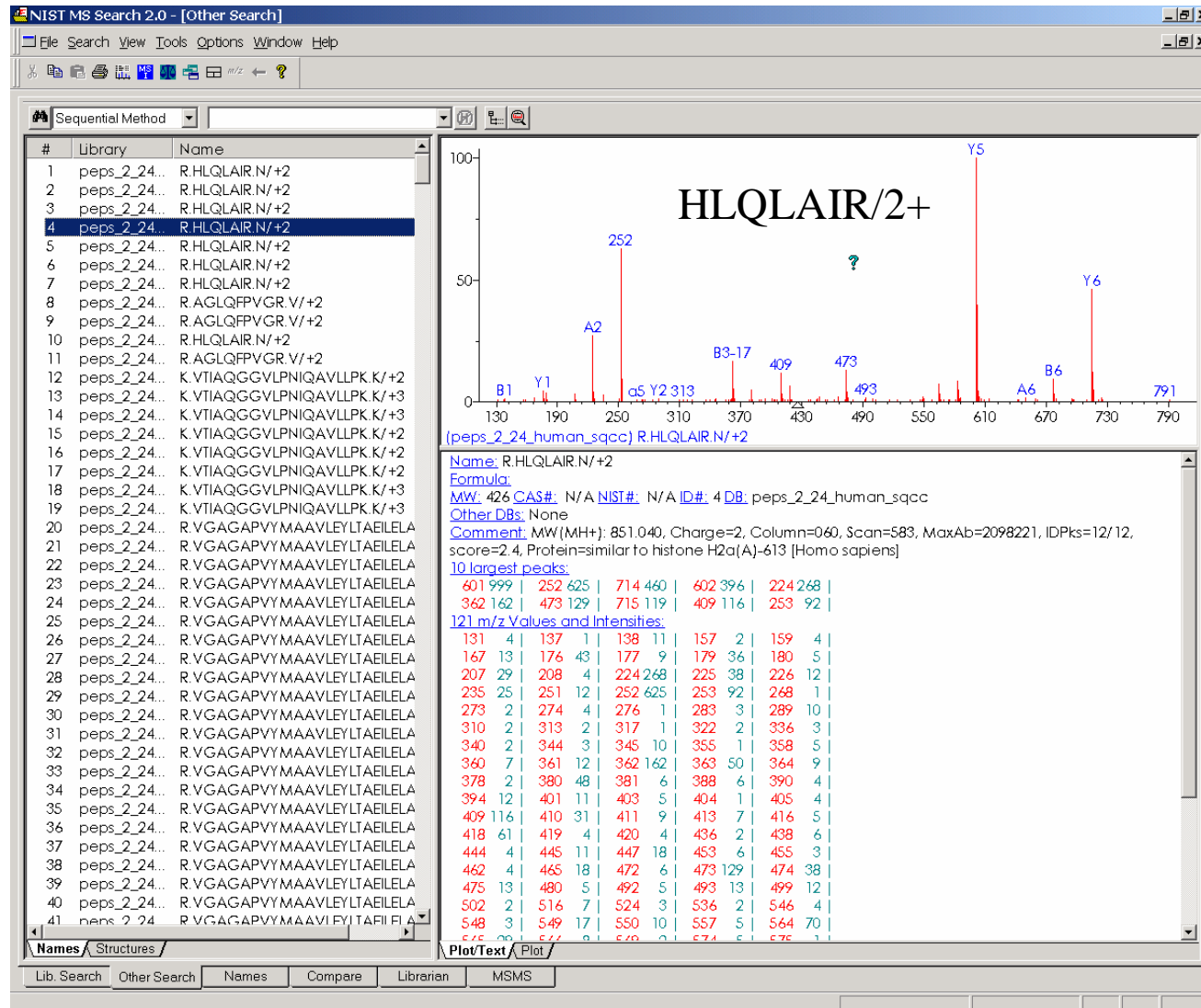
# Peptide Mass Spectra: Libraries for Organisms

- Proteins are linear sequences of amino acids
  - characteristic of Genome (organism)
- Peptides are 'digested' fragments of proteins
- MS 'sequences' peptides to reveal source Protein
- Peptides fragmentation spectra are not quite predictable
- Peptide fragmentation spectra for a 'genome' can be contained in one Library.
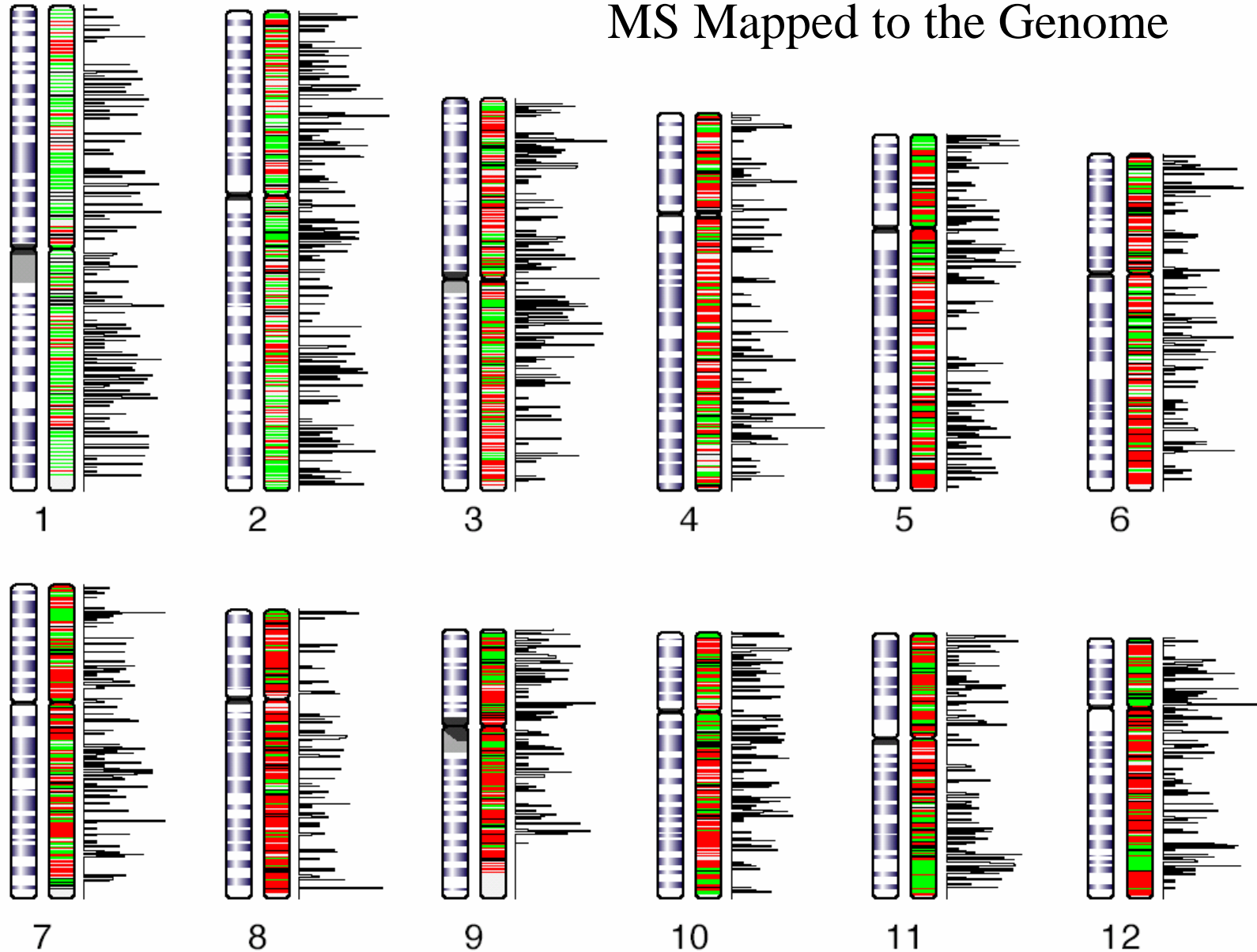
# Spectrum Prediction Programs

# Peptide Spectra Reference Library
## (multiple measurements each of 10,000 peptides)

MS Mapped to the Genome



From Eric Deutsch, ISB, 6/2004