

Algorithmic challenges in mass spectrometry

Anna Gambin

**Institute of Informatics,
University of Warsaw**



outline

I. modelling isotopic distribution

aggregated structure: BRAIN algorithm

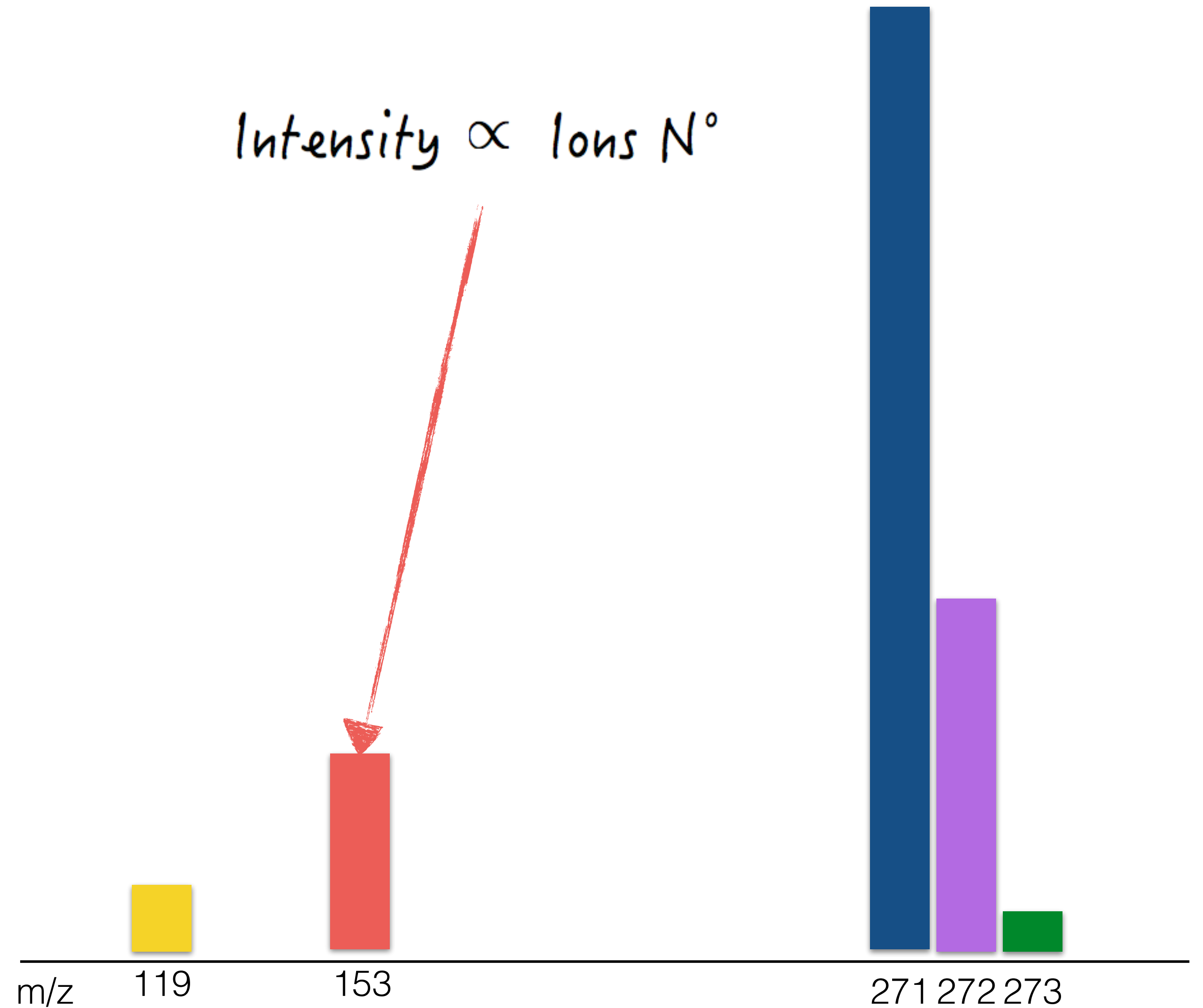
fine structure: ISOSPEC algorithm

II. Markov processes: modelling **fragmentation**

III. optimal transport in spectroscopy

Mass spectra

- mass spectrometer =
= highly precise
weighing scale
- compounds are
frequently
fragmented
- even a single
compound has a
complicated signal



Mass/charge [Th]

Chemical compounds are made of different isotopes

$$13.0033 - 12 = 1.0033 \text{ [Da]}$$

Element	Isotope	Extra Neutrons	Mass [Da]	Probability
Carbon	¹² C	0	12	0.9893
	¹³ C	1	13.0033	0.0107
Hydrogen	¹ H	0	1.0078	0.999885
	² H	1	2.0141	0.000115
Nitrogen	¹⁴ N	0	14.0031	0.99632
	¹⁵ N	1	15.0001	0.00368
Oxygen	¹⁶ O	0	15.9949	0.99757
	¹⁷ O	1	16.9991	0.00038
	¹⁸ O	2	17.9992	0.00205
Sulfur	³² S	0	31.9721	0.9493
	³³ S	1	32.9714	0.0076
	³⁴ S	2	33.9679	0.0429
	³⁶ S	4	35.9671	0.0002

elements have different numbers of stable isotopes

2

3

4

differences in frequencies of observation

isotopes of different elements differ in mass differences

$$32.9714 - 31.9721 = 0,9993 \text{ [Da]}$$

mathematical model of mass spectra

product of multinomial distributions

Assume

1) variants of isotopes of atoms are **independent**

2) elements **vary in abundances** of isotopes

$$\mathcal{P}({}^{12}\text{C}_{c_0} {}^{13}\text{C}_{c_1} {}^1\text{H}_{h_0} {}^2\text{H}_{h_1} {}^{14}\text{N}_{n_0} {}^{15}\text{N}_{n_1} {}^{16}\text{O}_{o_0} {}^{17}\text{O}_{o_1} {}^{18}\text{O}_{o_2} {}^{32}\text{S}_{s_0} {}^{33}\text{S}_{s_1} {}^{34}\text{S}_{s_2} {}^{36}\text{S}_{s_4}) =$$

$$\binom{c}{c_0, c_1} \mathcal{P}({}^{12}\text{C})^{c_0} \mathcal{P}({}^{13}\text{C})^{c_1} \binom{h}{h_0, h_1} \mathcal{P}({}^1\text{H})^{h_0} \mathcal{P}({}^2\text{H})^{h_1} \binom{n}{n_0, n_1} \mathcal{P}({}^{14}\text{N})^{n_0} \mathcal{P}({}^{15}\text{N})^{n_1} \times$$

$$\binom{n}{o_0, o_1, o_2} \mathcal{P}({}^{16}\text{O})^{o_0} \mathcal{P}({}^{17}\text{O})^{o_1} \mathcal{P}({}^{18}\text{O})^{o_2} \binom{s}{s_0, s_1, s_2, s_4} \mathcal{P}({}^{32}\text{S})^{s_0} \mathcal{P}({}^{33}\text{S})^{s_1} \mathcal{P}({}^{34}\text{S})^{s_2} \mathcal{P}({}^{36}\text{S})^{s_4}$$



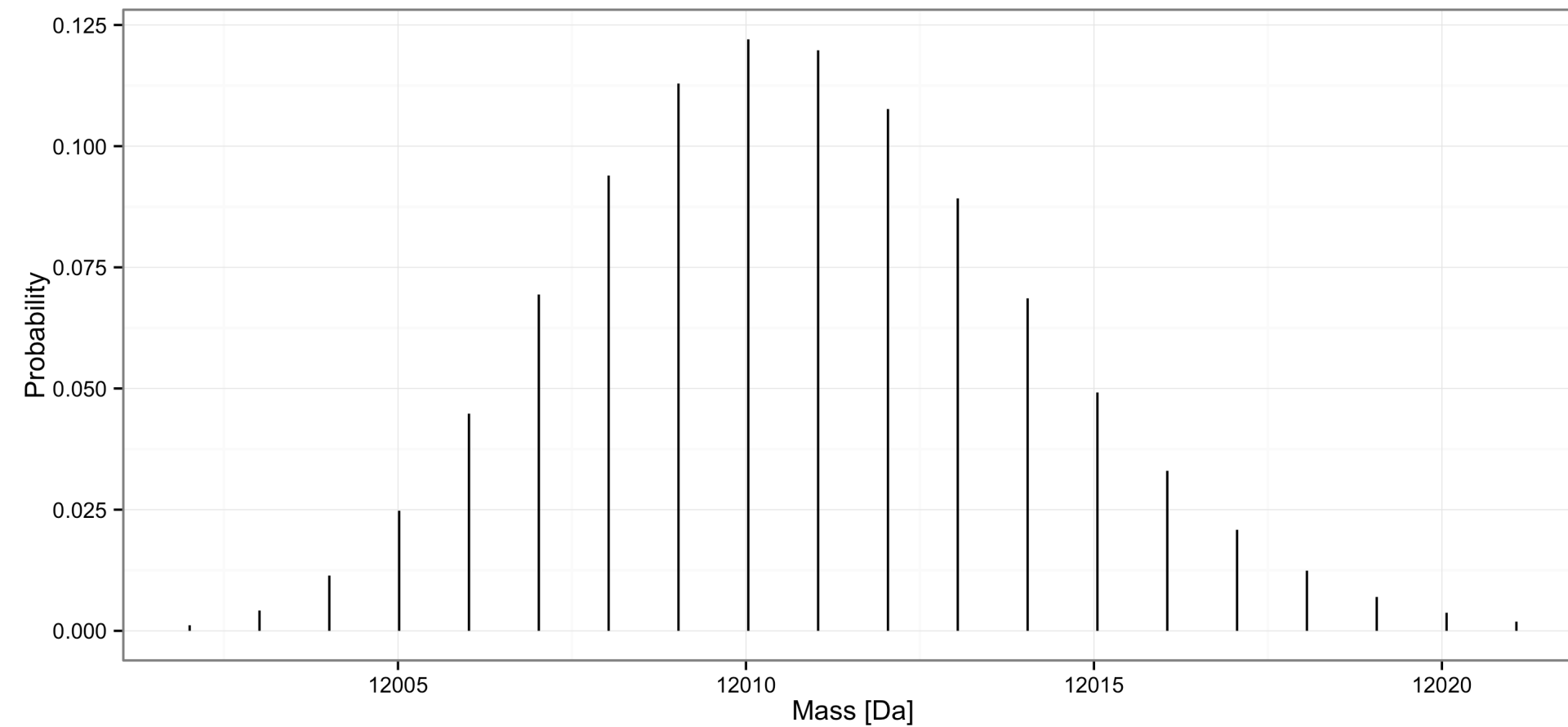
frequencies
of isotopes



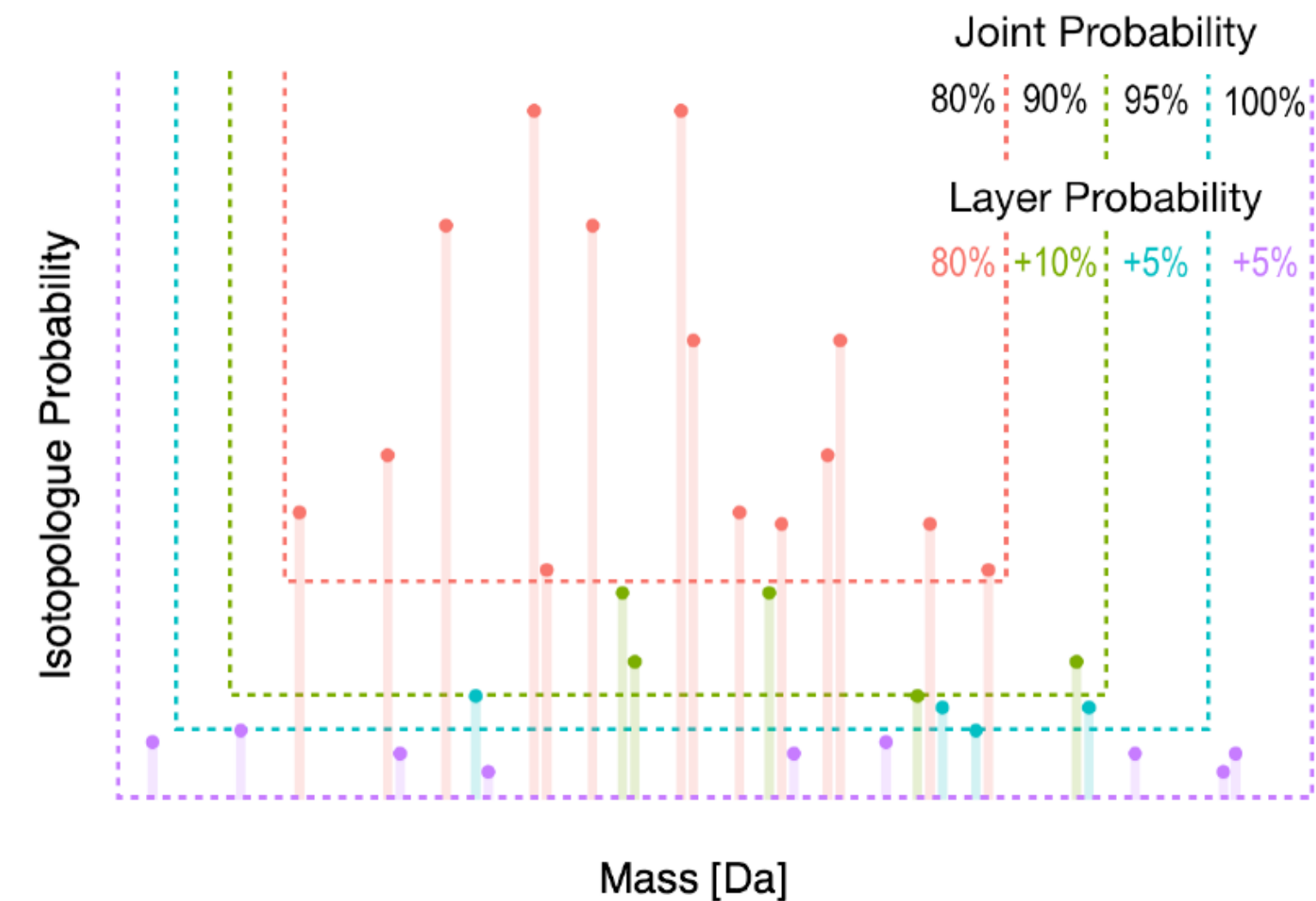
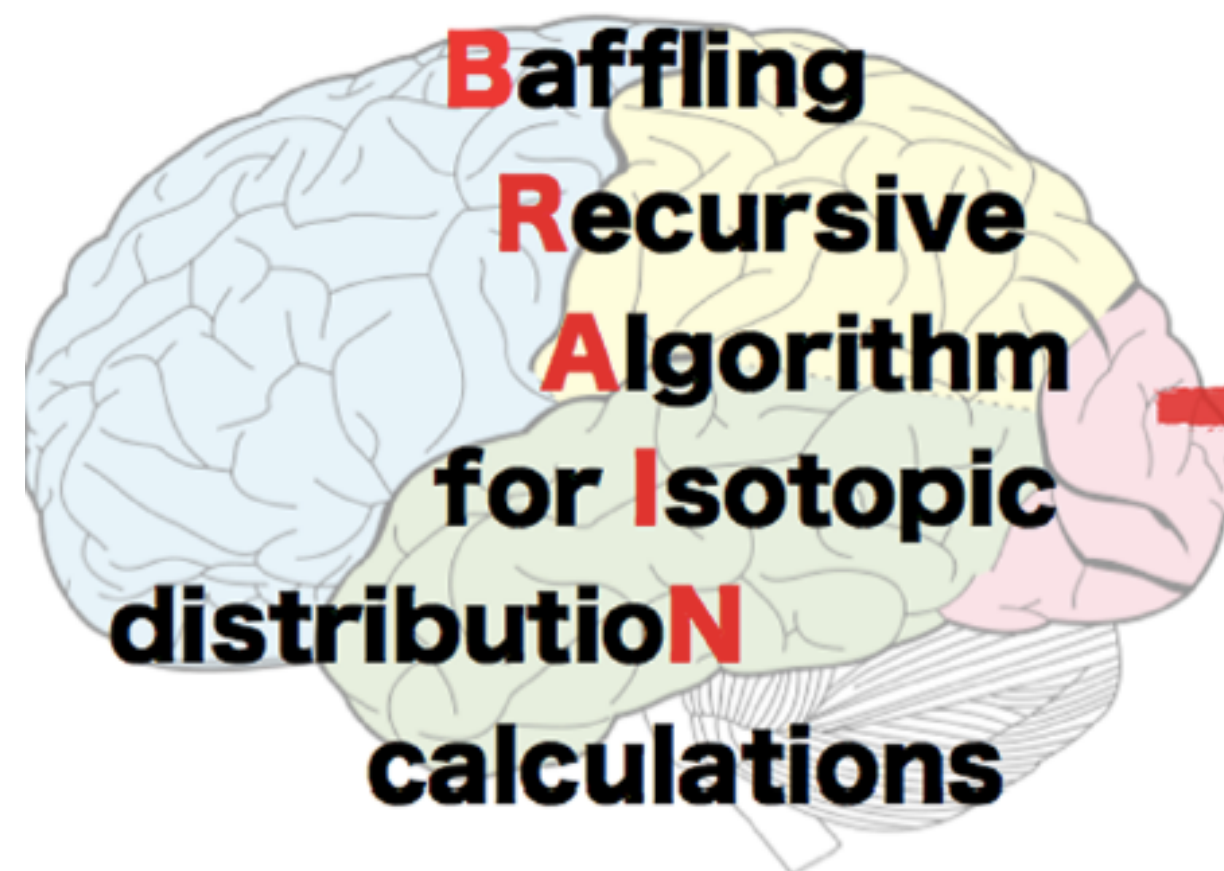
I U P A C

INTERNATIONAL UNION OF
PURE AND APPLIED CHEMISTRY

mathematical model of mass spectra



low resolution problem:
aggregated isotopic structure



high resolution problem:
fine isotopic structure



ISOSPEC

The Fine Isotopic Structure Calculator

aggregated isotopic distribution

we **group** together variants with **the same number of additional neutrons**

for molecular formula $C_v H_w N_x O_y S_z$ consider polynomial:

$$Q(I; v, w, x, y, z) = (P_{C_{12}} I^0 + P_{C_{13}} I^1)^v \times (P_{H_1} I^0 + P_{H_2} I^1)^w \\ \times (P_{N_{14}} I^0 + P_{N_{15}} I^1)^x \times (P_{O_{16}} I^0 + P_{O_{17}} I^1 + P_{O_{18}} I^2)^y \\ \times (P_{S_{32}} I^0 + P_{S_{33}} I^1 + P_{S_{34}} I^2 + P_{S_{36}} I^4)^z$$

$n = v + w + x + 2y + 4z$

frequencies
of isotopes, e.g.

$$P_{C_{12}} = 98.93\% \text{ and } P_{C_{13}} = 1.07\%$$

$$Q(I; v, w, x, y, z) \equiv \sum_{j=0}^n q_j I^j$$

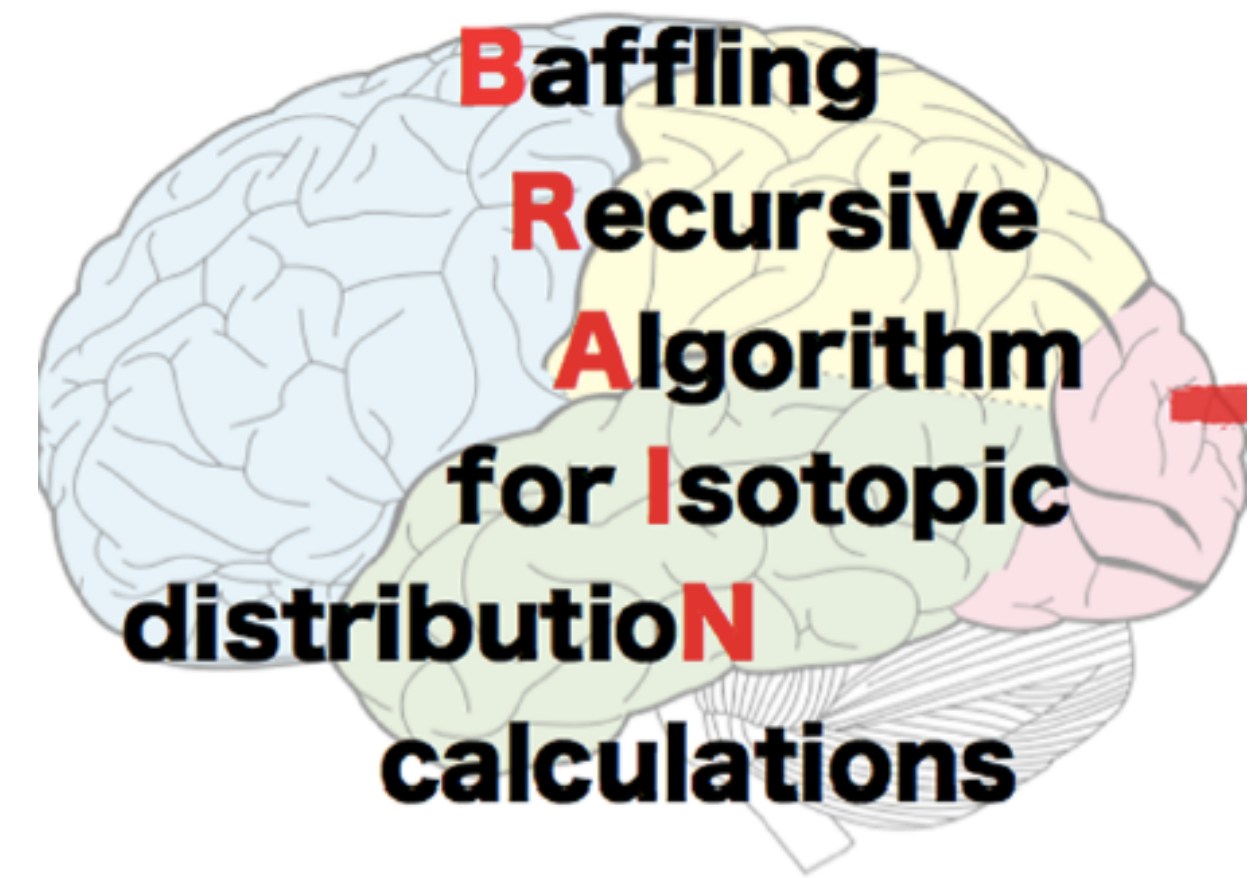
probability of peak
with j additional
neutrons

algorithm

to compute coefficients of polynomial:

$$Q(I; v, w, x, y, z) \equiv \sum_{j=0}^n q_j I^j$$

determine its roots: $r_C = -\frac{P_{C_{12}}}{P_{C_{13}}}$, $r_H = -\frac{P_{H_1}}{P_{H_2}}$, and $r_N = -\frac{P_{N_{14}}}{P_{N_{15}}}$. $r_O, \bar{r}_O = \frac{-P_{O_{17}} \pm \sqrt{P_{O_{17}}^2 - 4P_{O_{16}}P_{O_{18}}}}{2P_{O_{18}}}$



apply the recurrent formula (follows from Newton-Girard theorem and Viète's formulae)

$$q_j = -\frac{1}{j} \sum_{l=1}^j q_{j-l} \psi_l$$

where

$$\psi_l = v(r_C)^{-l} + w(r_H)^{-l} + x(r_N)^{-l} + (r_O)^{-j} + (\bar{r}_O)^{-j} + (r_{S,1})^{-j} + (\bar{r}_{S,1})^{-j} + (r_{S,2})^{-j} + (\bar{r}_{S,2})^{-j}$$

complexity: quadratic; **exact** values calculated

fine isotopic distribution

$$\mathcal{P}(^{16}\text{O}) = \frac{4}{9} \quad \mathcal{P}(^{17}\text{O}) = \frac{3}{9} \quad \mathcal{P}(^{18}\text{O}) = \frac{2}{9}$$

100 oxygen atoms

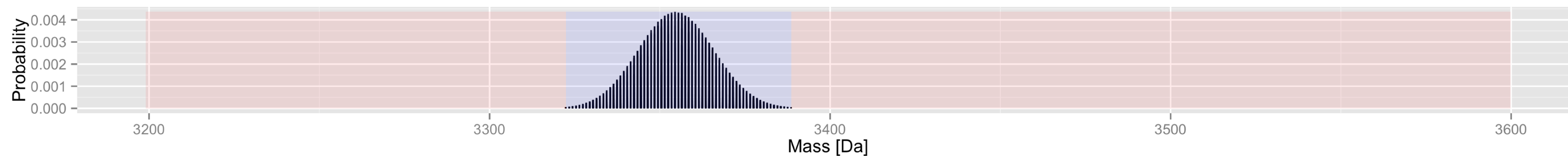
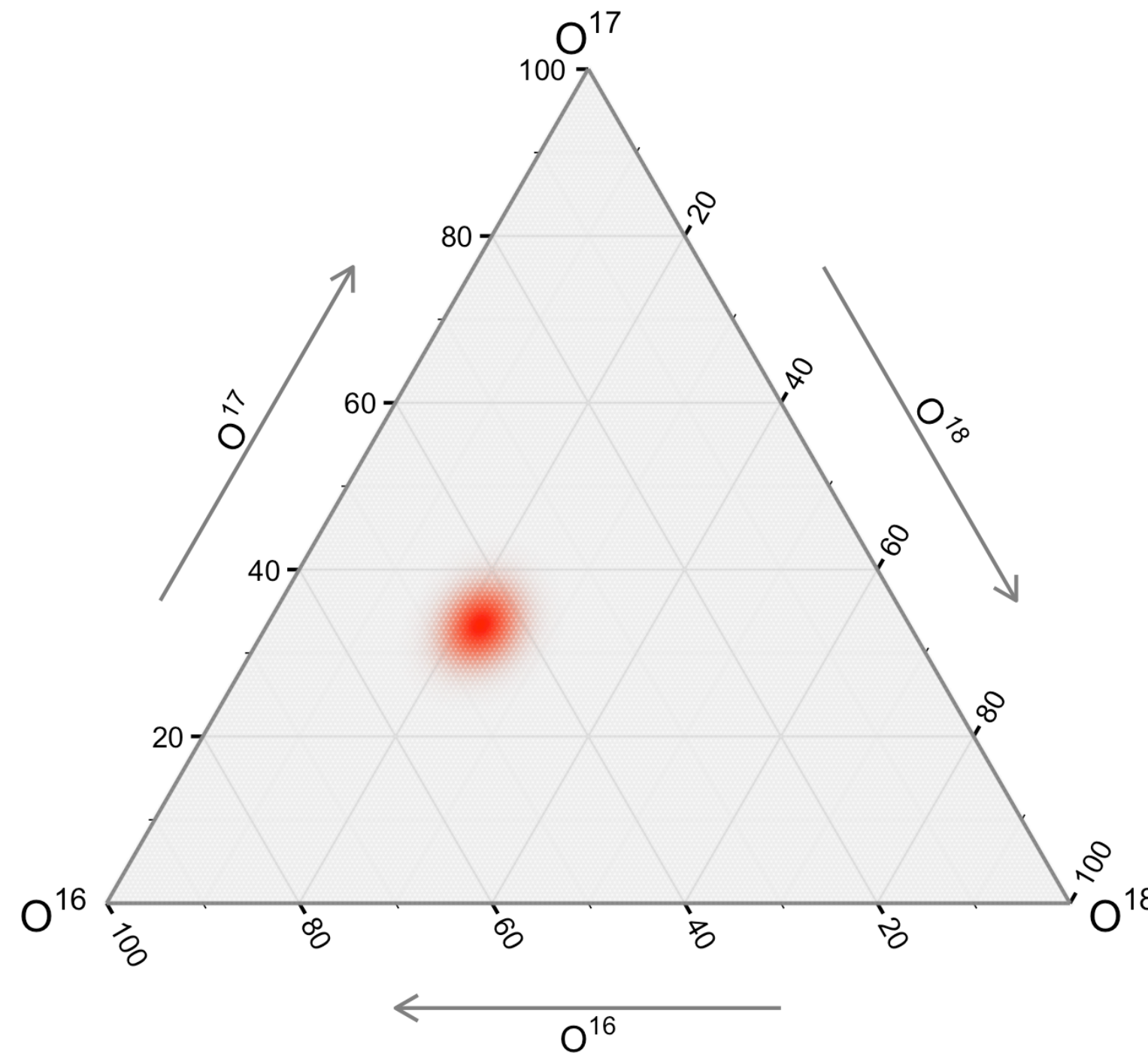
$$o_0 + o_1 + o_2 = 200$$

(not real world values!)

20301 variants

whereas 1043
bear 99% prob.

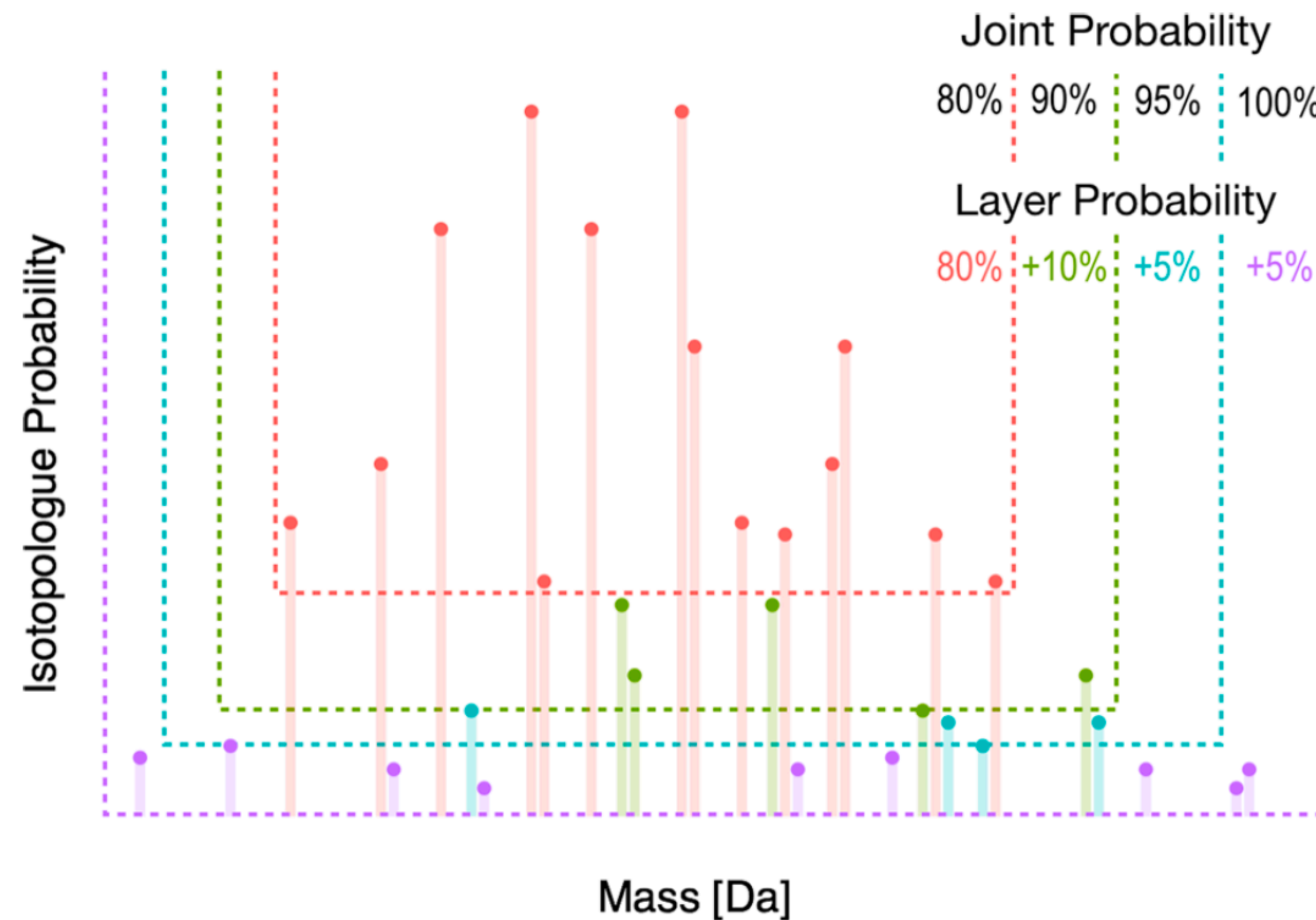
→ huge number
of isotopologues



fine isotopic distribution

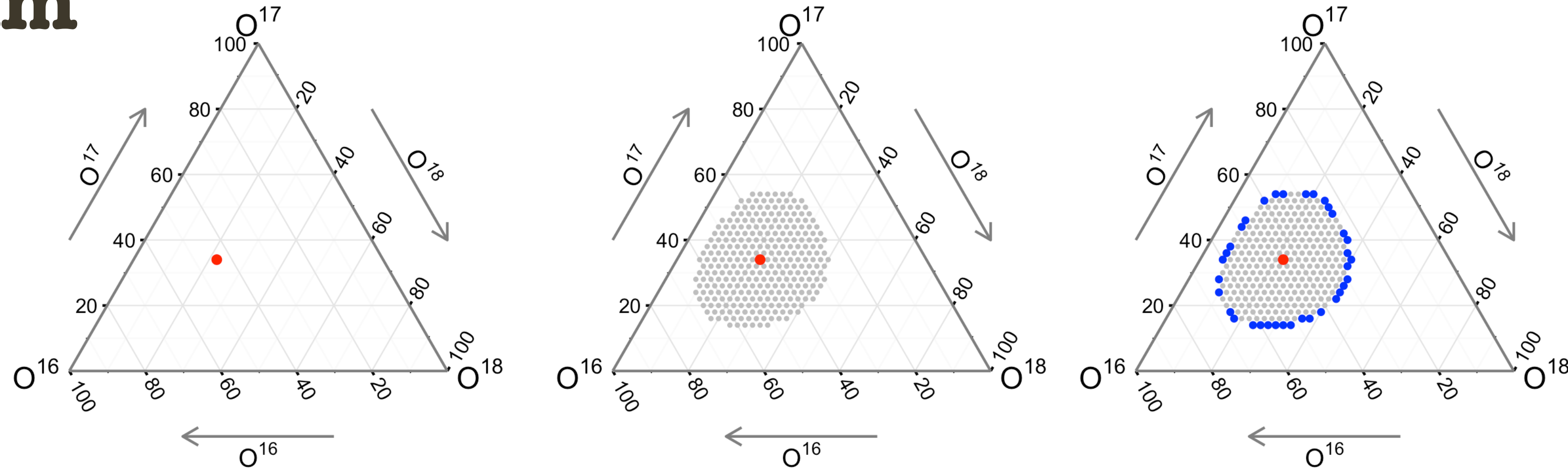
optimal p -set

smallest set of isotopologues that jointly surpass probability p



division of isotopic distribution into optimal p -sets: 80%, 90%, 95%, 100%

algorithm



To get the **optimal P-set**:

Find the **most probable variant**

while **Total Probability** < **P** :

Get **layer** of v so that $p > P(v) > 0,5 p$ where

$$p = P(v_{min \text{ previous layer}})$$

Trim the **least probable variants** from
the last layer so that **Total Probability** = **P**

complexity:

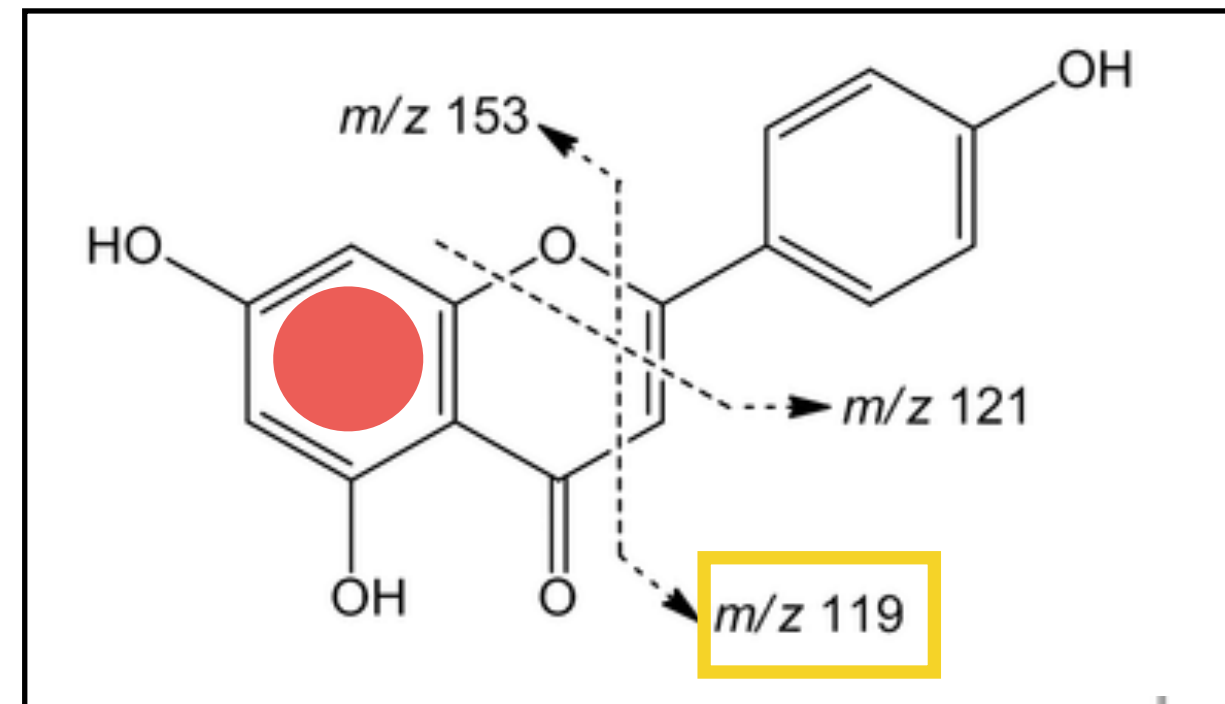
$O(n)$ in the total
number
of configurations



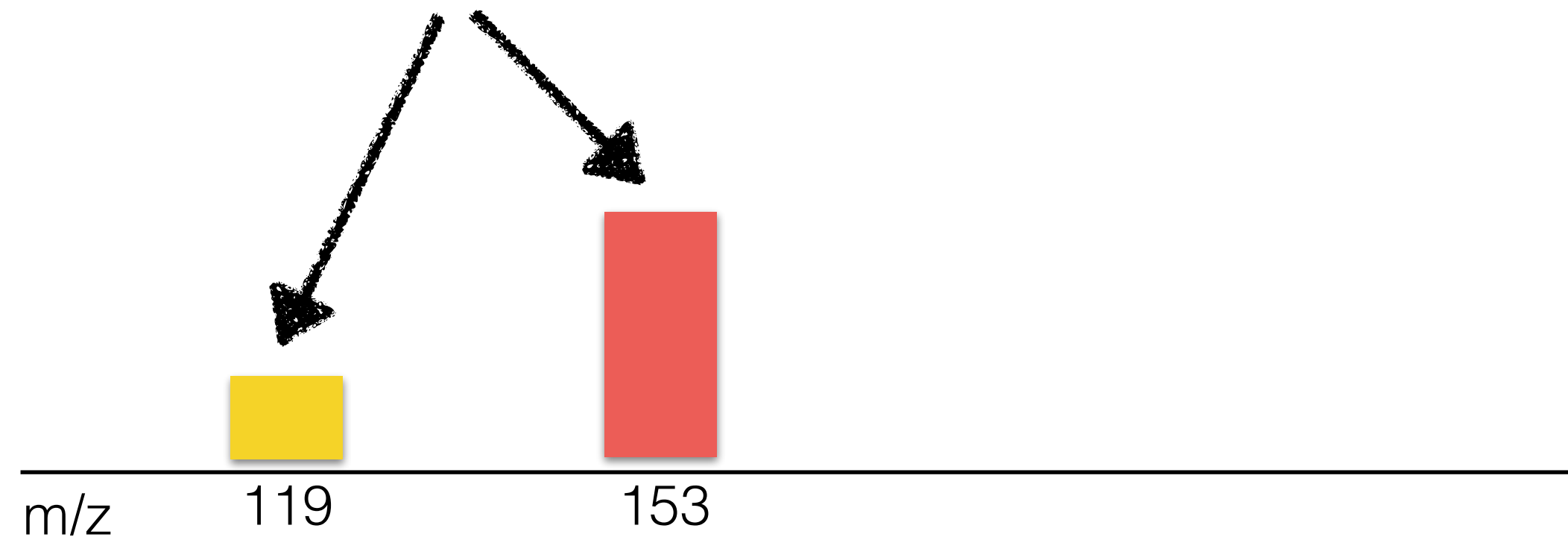
ISOSPEC

The Fine Isotopic Structure Calculator

II. Markov processes: modelling fragmentation



some **bonds** get easily broken

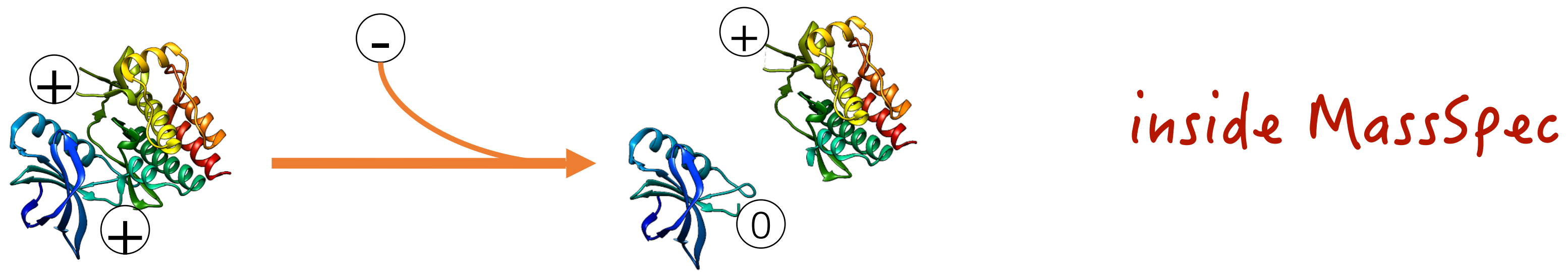


electron is transferred to the positively-charged protein or peptide, causing fragmentation along the peptide **backbone**

.. others **not**



modelling fragmentation: problem



Cleavage of protein backbone by a rapid neutralization of charge

ETD – main reaction, others = side reactions



problem: for the set of biochemical reactions

determine their intensitie having observed the substrates

modelling fragmentation: solution

model the phenomena as Markov process describing the flow of particles through the fragmentation graph

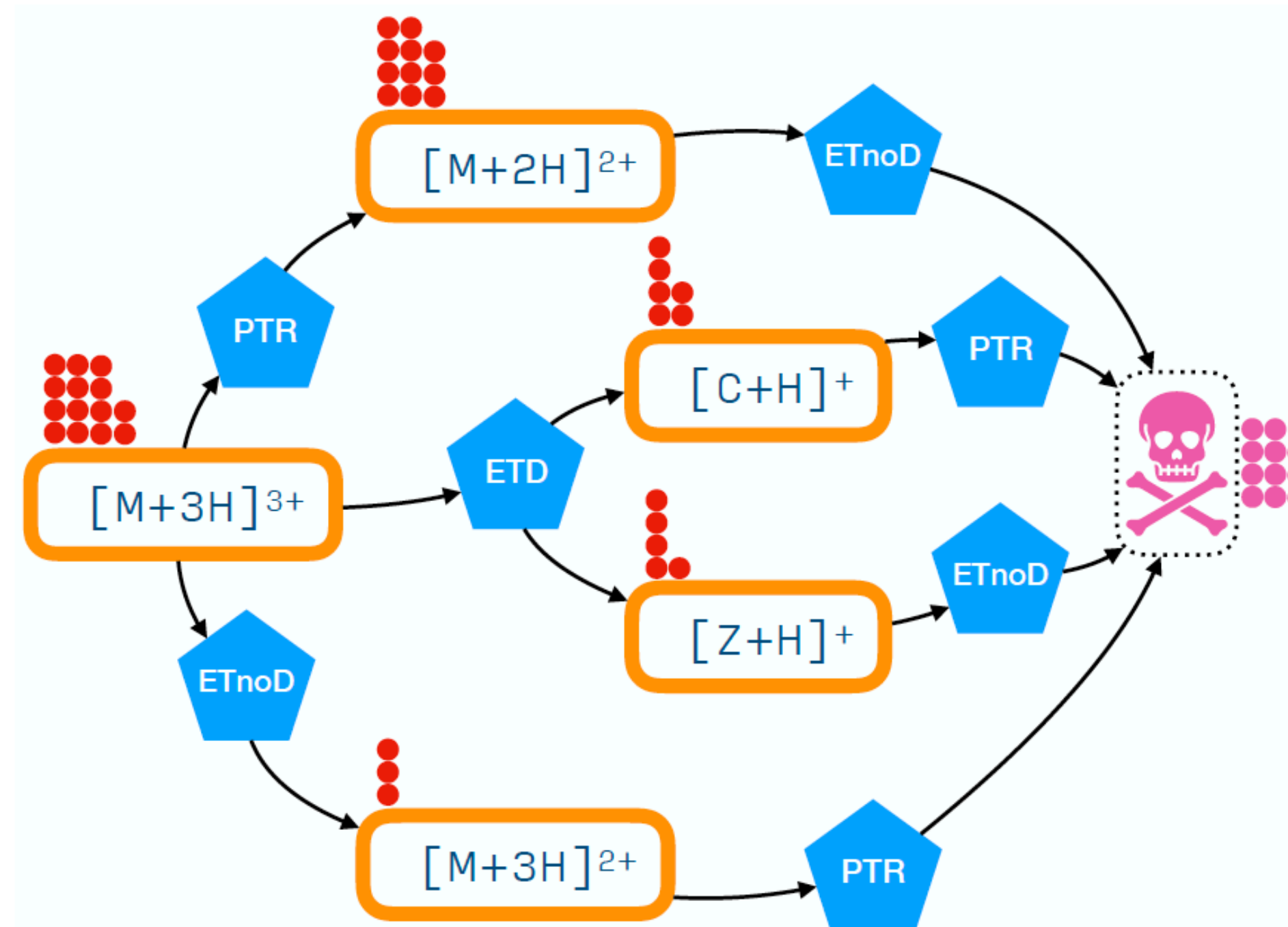
calculate expectance in the model:

use ODE description for big population of particles

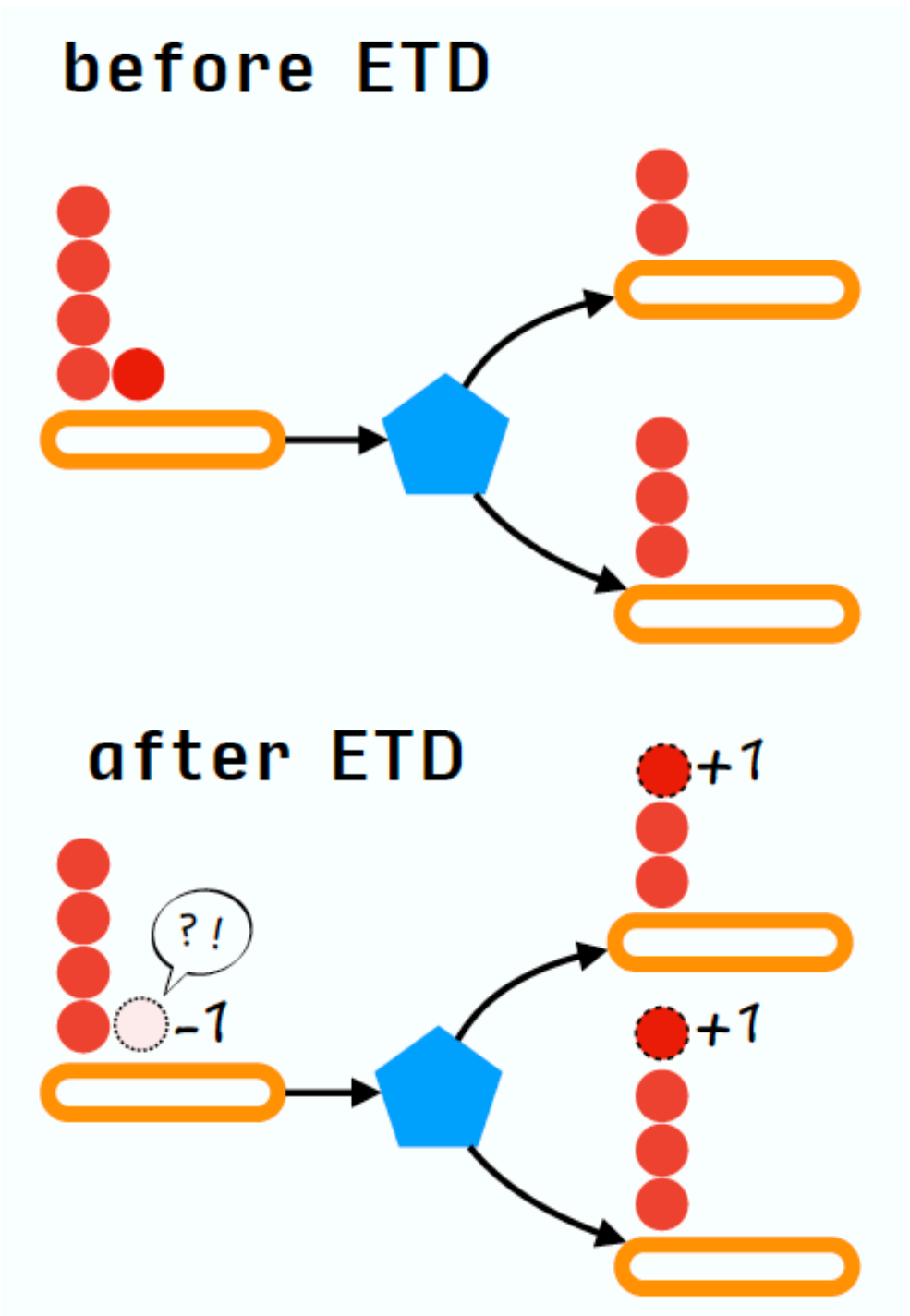
compare to observed data:

find intensities that best predict the observed data

by minimising the discrepancy (nonlinear optimisation)

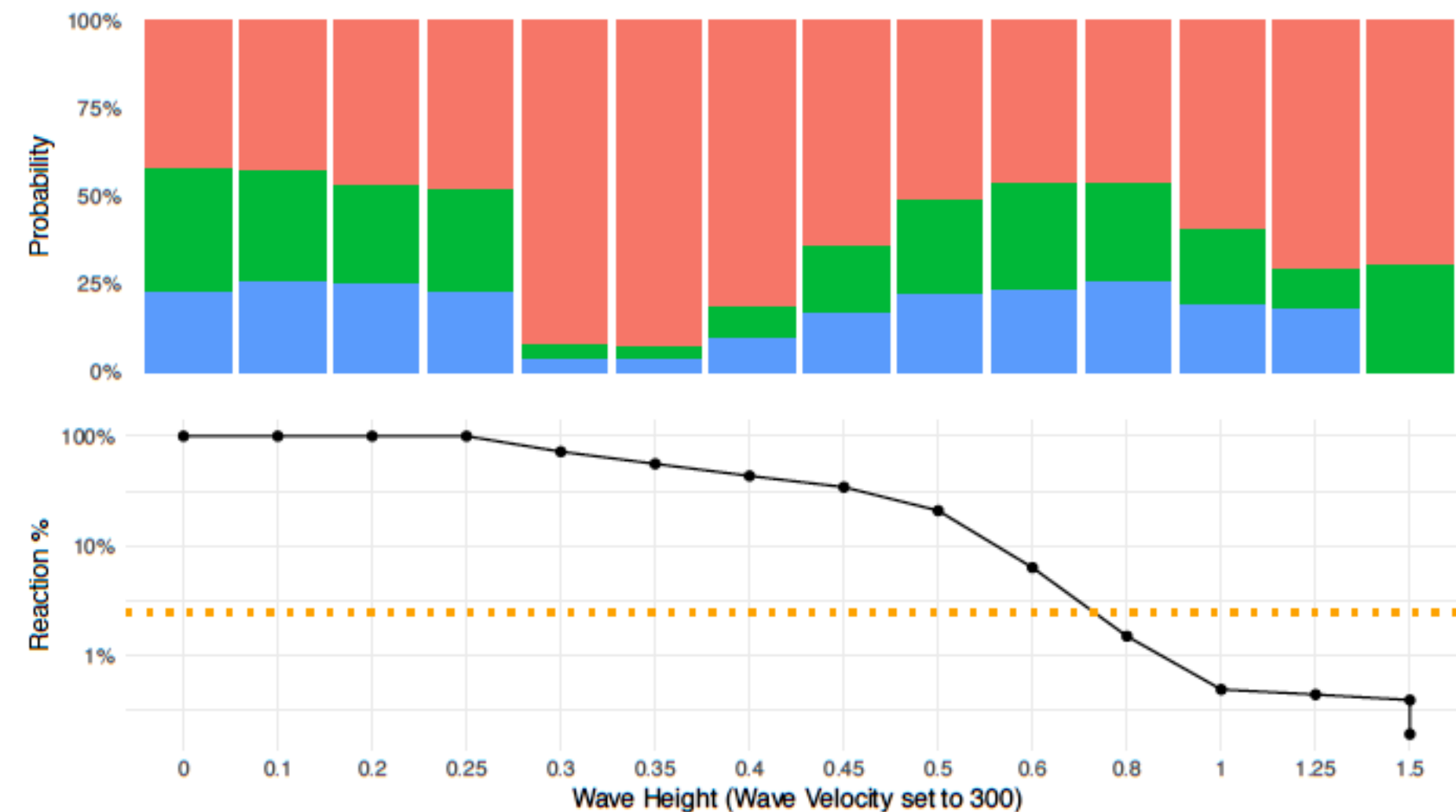
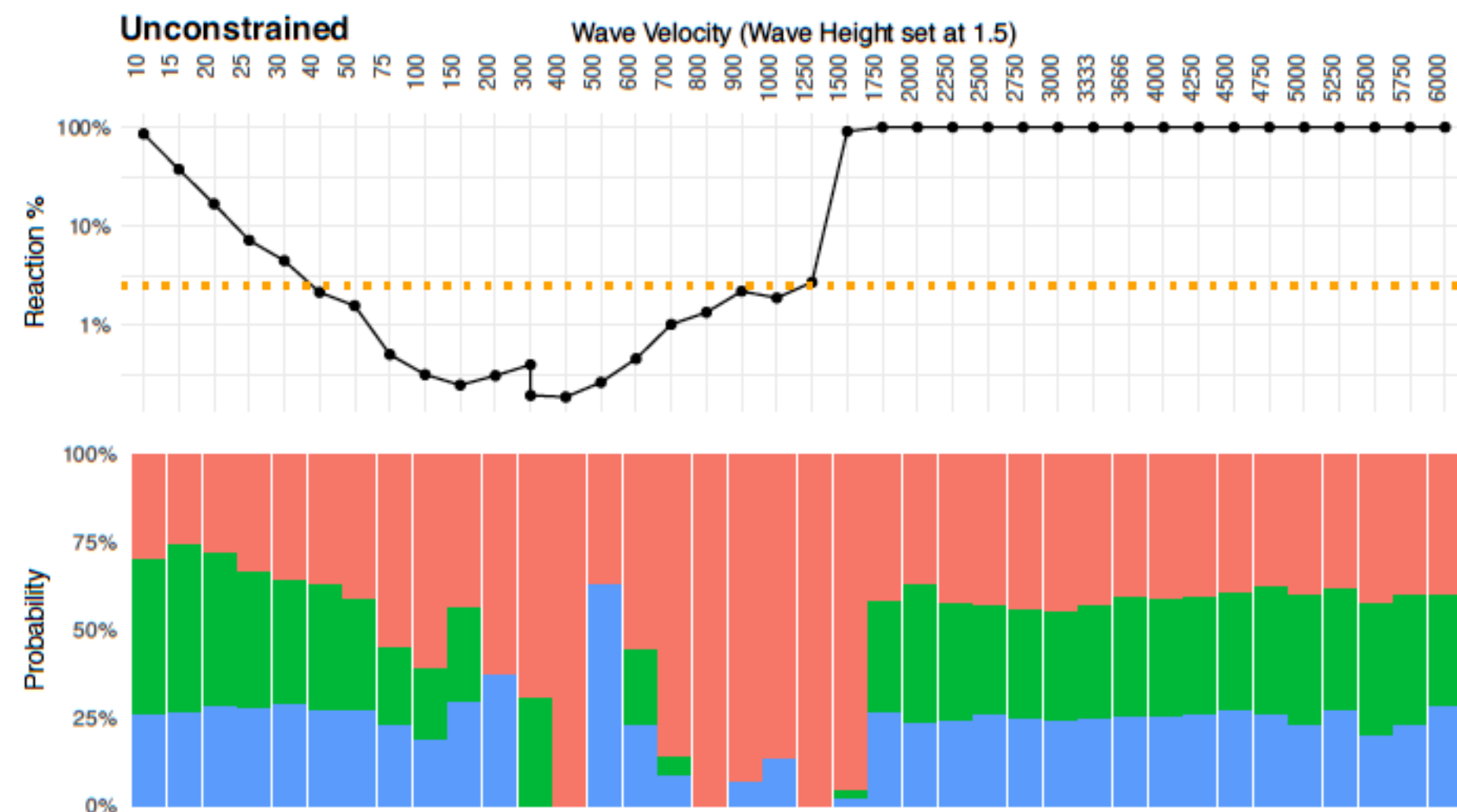


(a) A Petri Net

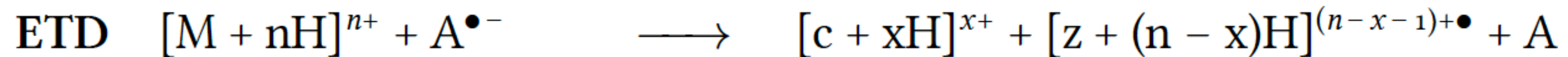


(b) ETD reaction

modelling fragmentation: results



proportions of **PTR**, **ETD**, **ETnoD** for different **MS** setup



III. optimal transport in spectroscopy



Eight tons of hope: world's strongest persistent magnet for NMR at ETH

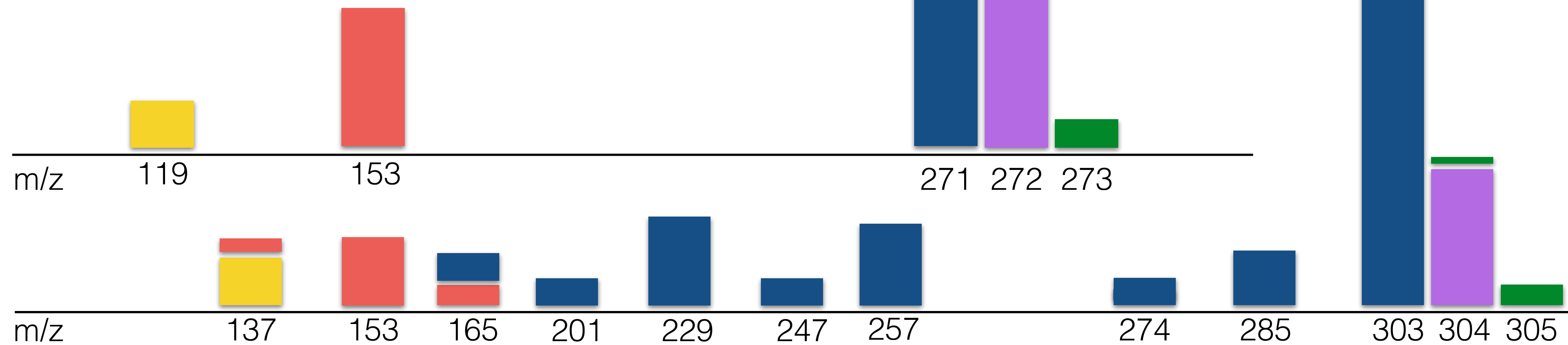
08.06.2020 by Julia Ecker

Wasserstein metric

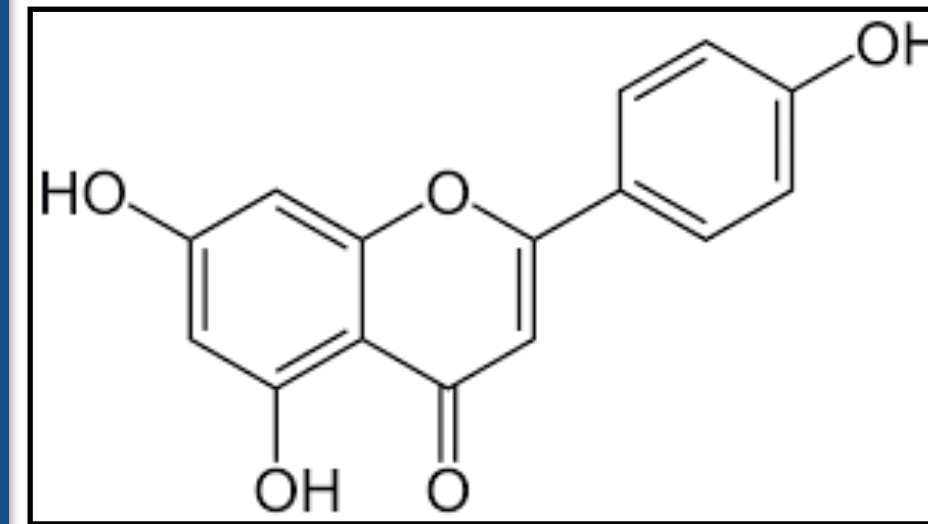
$$W(\mu, \nu) = \min_{\gamma} \sum_{x,y} |x - y| \gamma(x, y)$$

distance

amount of transported signal

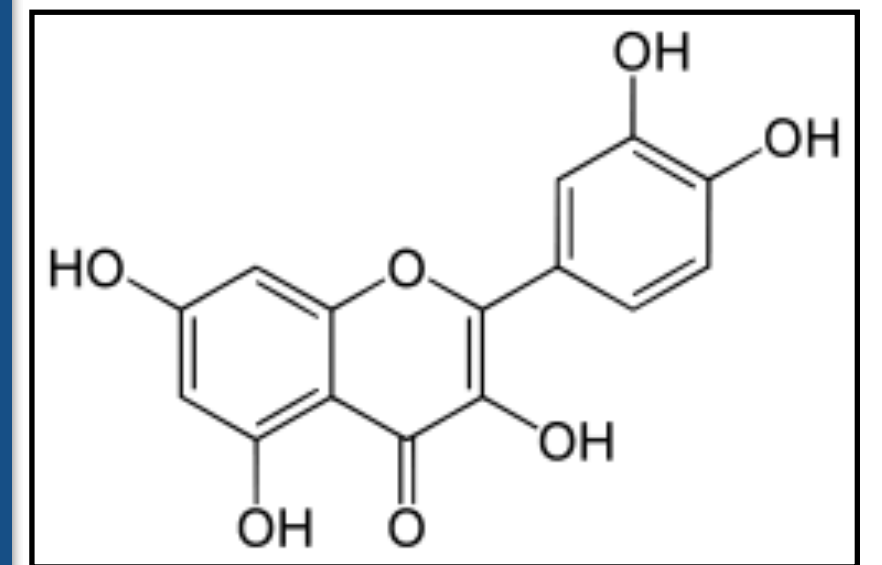


← Apigenin

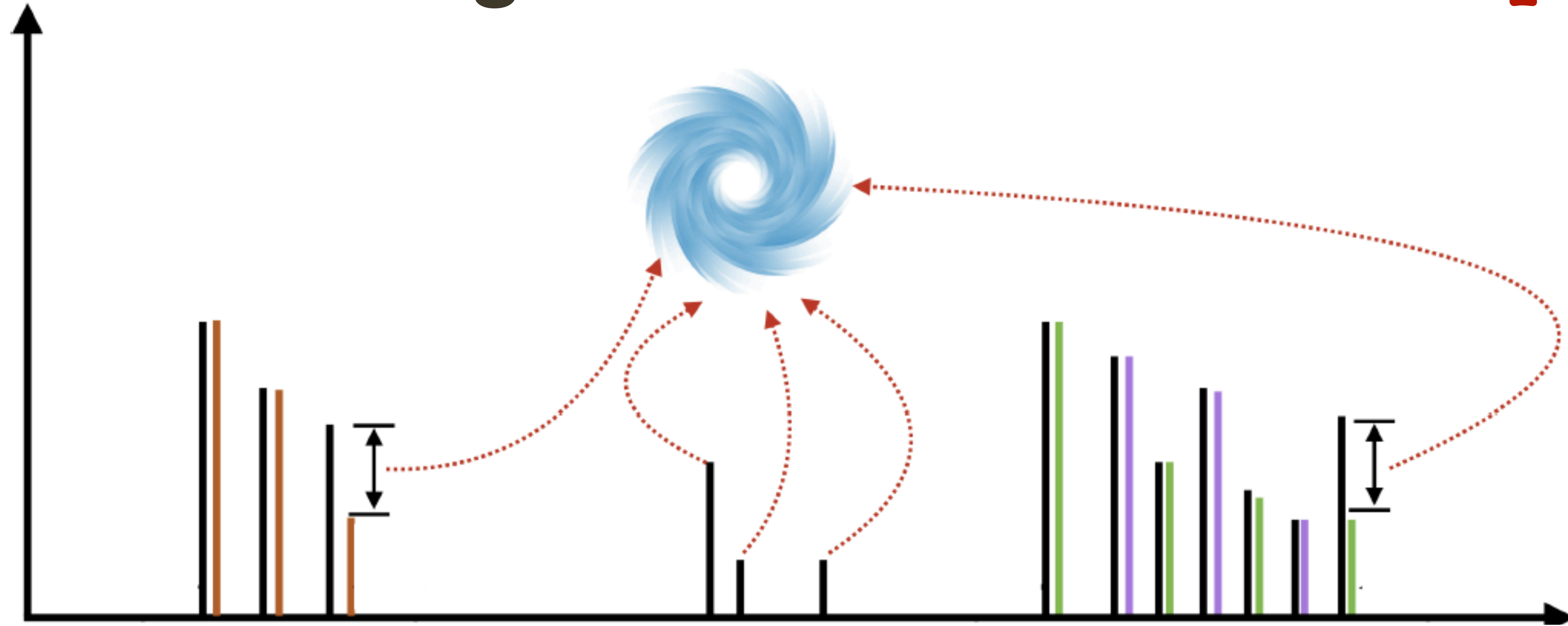


32 Da = 02

← Quercetin



Wasserstein regression: **deconvolution of spectra**



linear combination of **simulated = model** spectra



$$\nu_p = p_1\nu_1 + \dots + p_k\nu_k$$

observed spectrum = model + noise



$$\mu(x) = \nu_p(x) + \varepsilon(x)$$

Wasserstein **regression**: find proportions



$$p^* = \arg \min_p W(\mu, \nu_p)$$

finally add **vortex** to transport noise

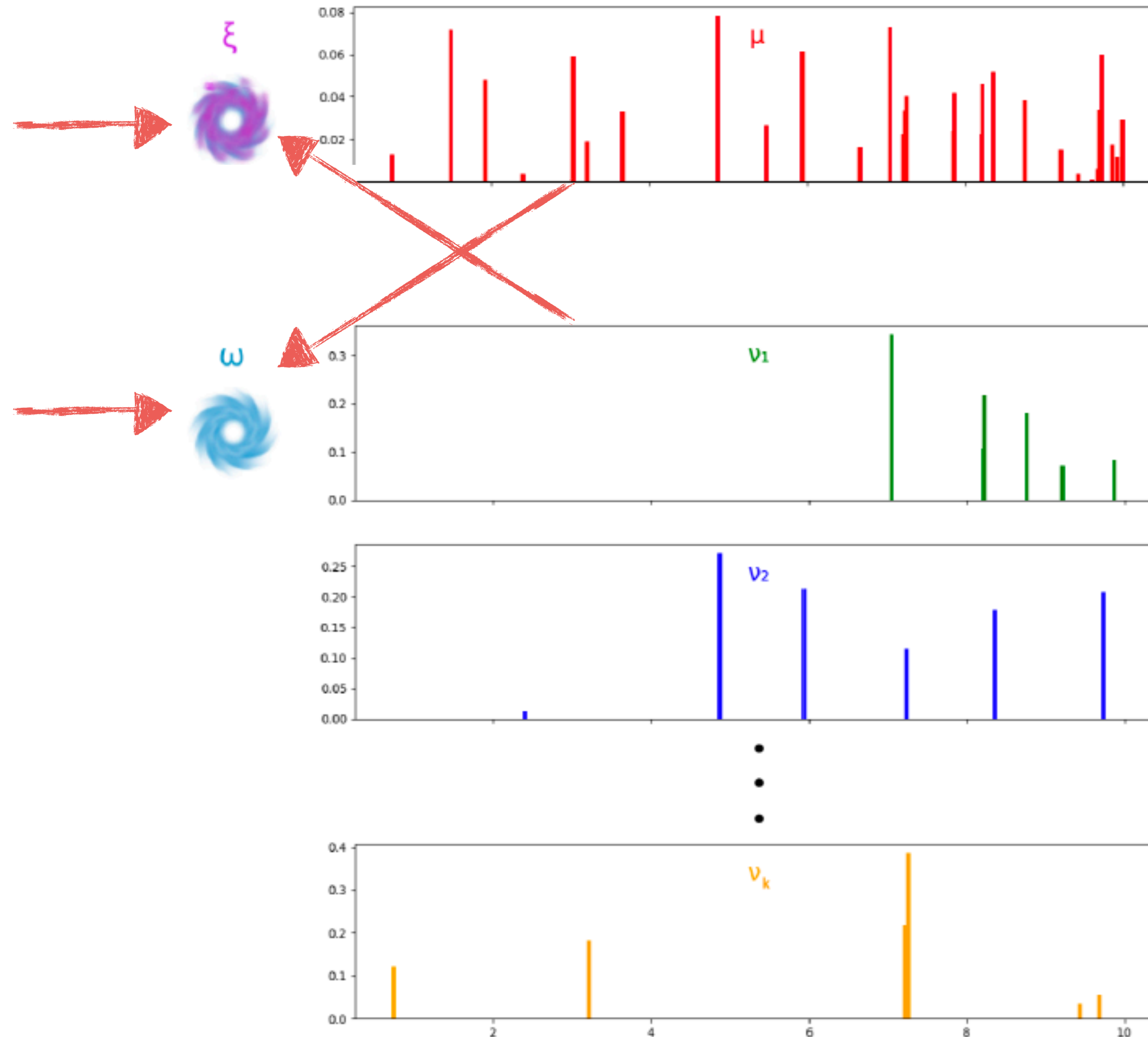


$$p^* = \arg \min_p W(\mu, p_0\omega + \nu_p)$$

working with two vortices

$$p^* = \min_{p=(p_1, p_2, \dots, p_k)} W(p_0 \omega + p_1 \nu_1 + p_2 \nu_2 + \dots + p_k \nu_k, (1 - p'_0) \mu + p'_0 \xi),$$

to remove **excess** of hypothetical spectra in the model



to remove **noise** from experimental data

1 dimensional case is easy to calculate:

$$W(\mu, \nu) = \int_{\mathbb{R}} |M(t) - N(t)| dt$$

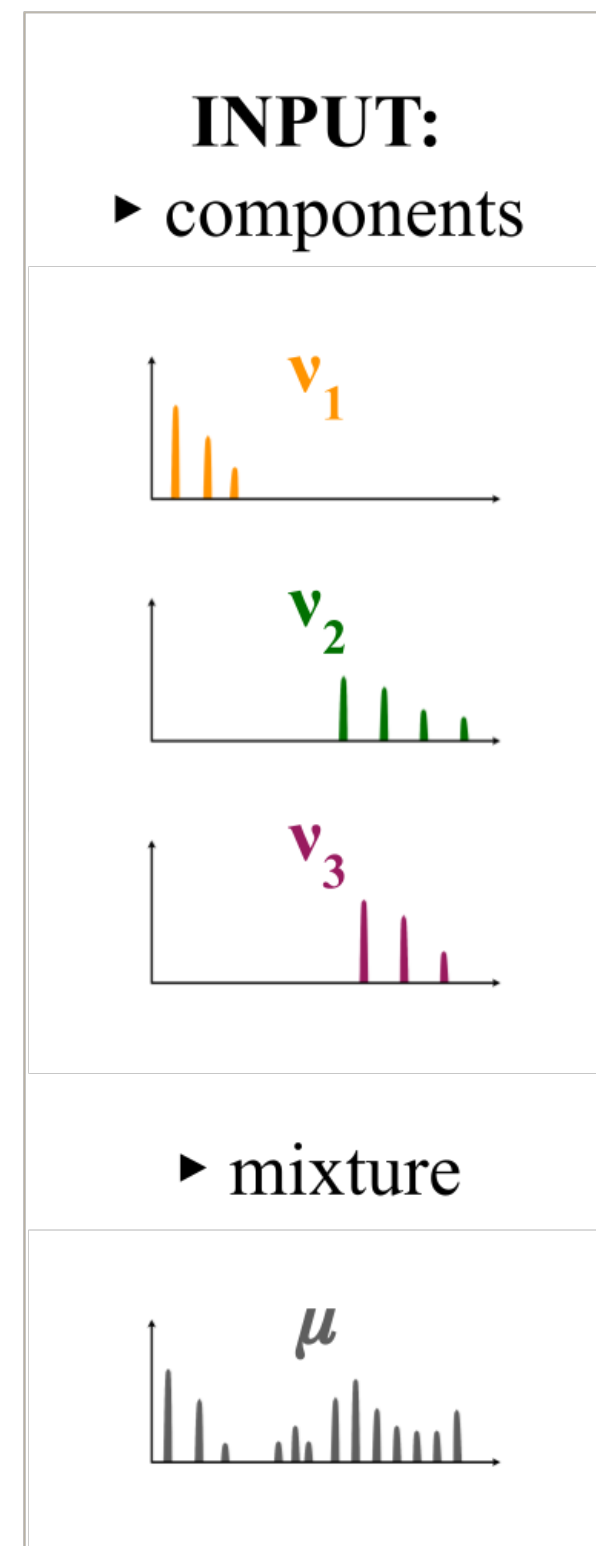
using

$$N(t) = \sum_{j=1}^k p_j N_j(t)$$

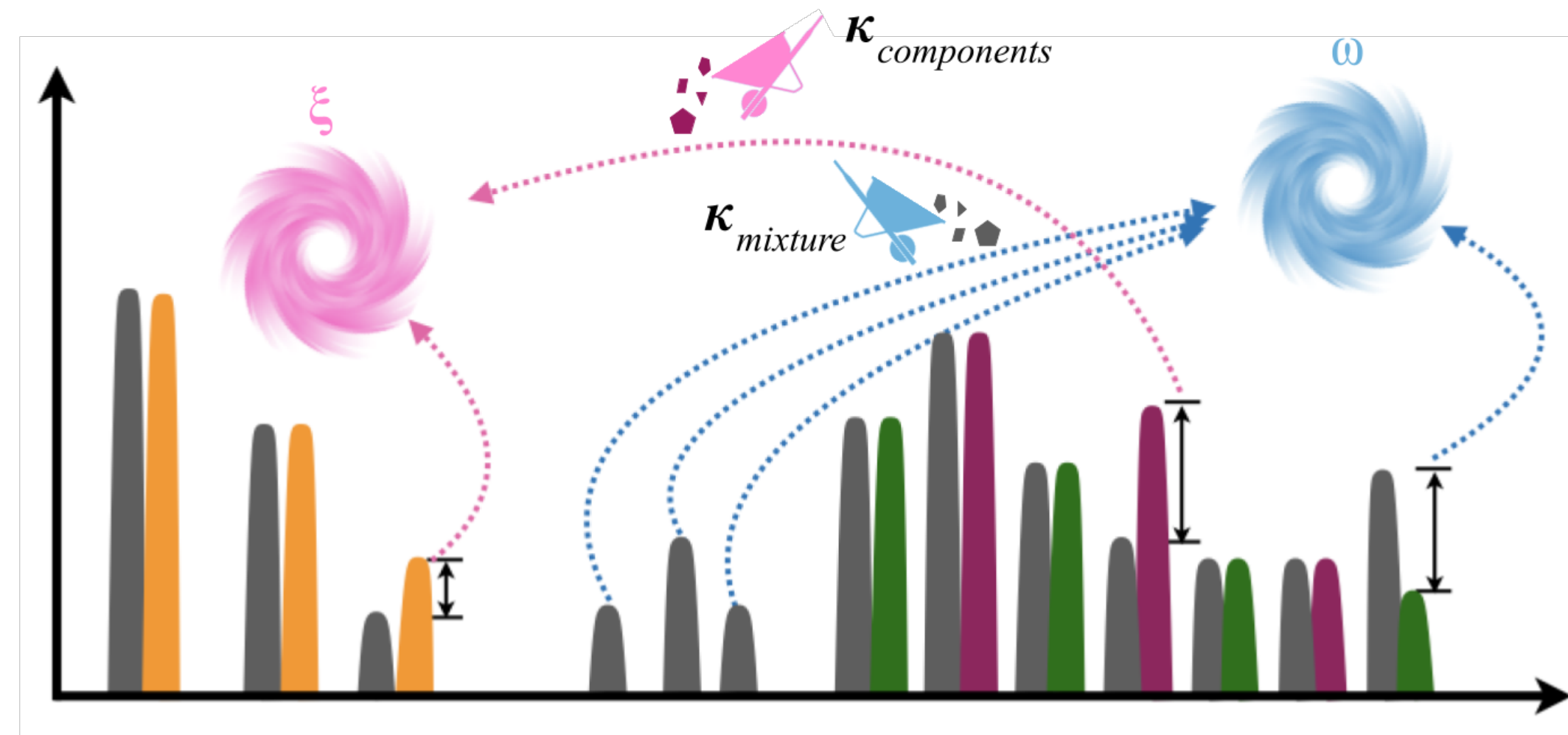
$$M(t)$$

cumulative distribution functions

Magnetstein for NMR analysis

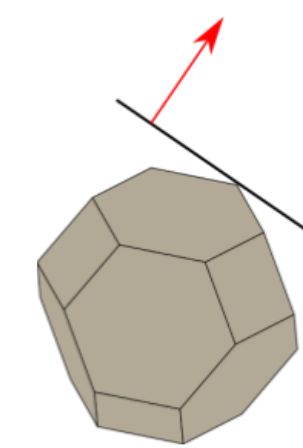
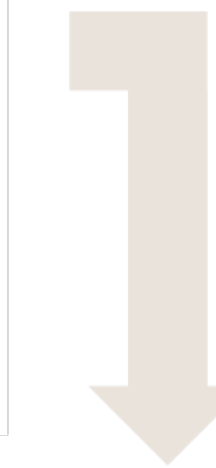


normalization



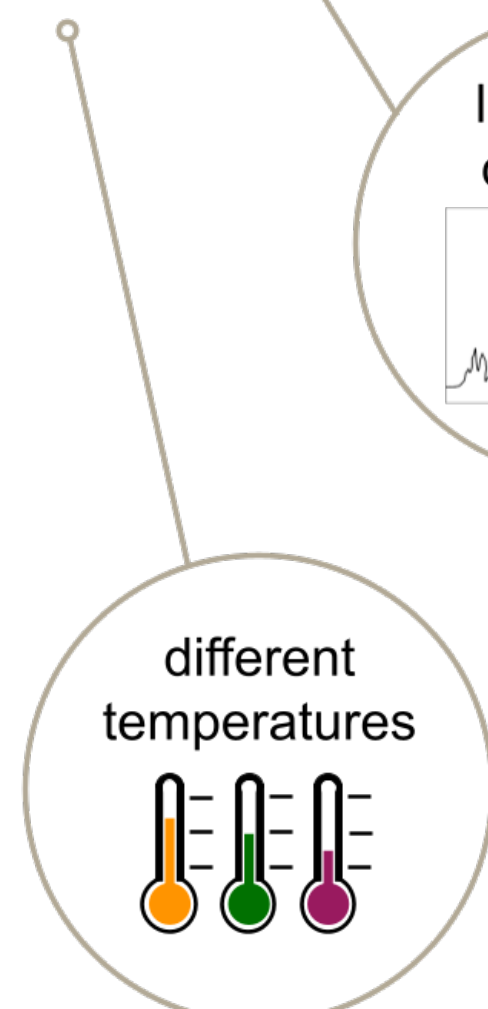
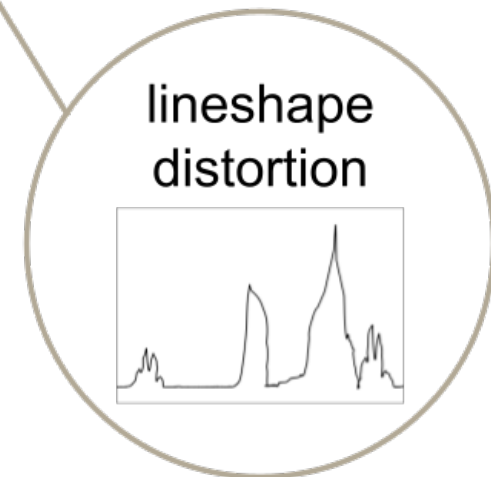
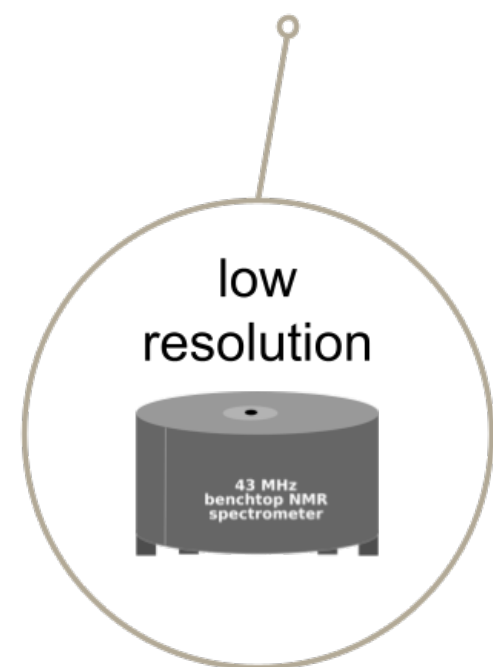
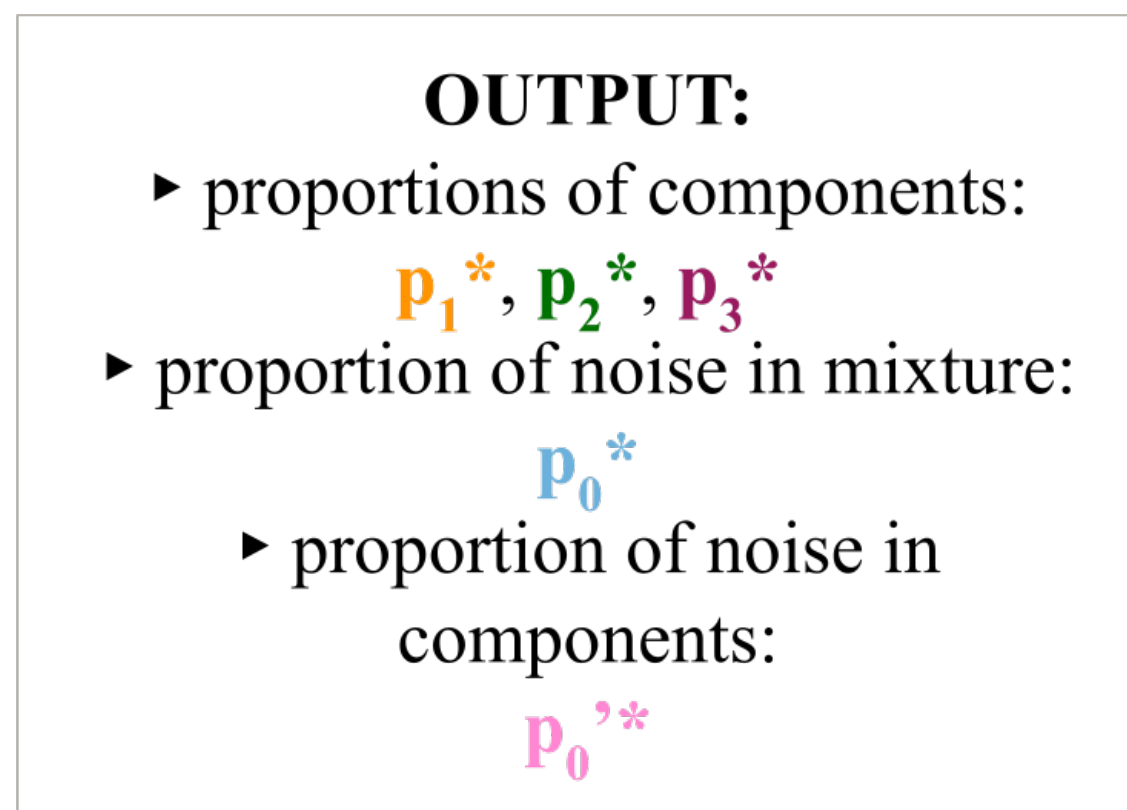
Regression with the Wasserstein distance:

$$\min_{(p_0, p_0', p_1, p_2, p_3)} \mathbf{W}^{\kappa}(p_0 \omega + p_1 \mathbf{v}_1 + p_2 \mathbf{v}_2 + p_3 \mathbf{v}_3, (1 - p_0') \mu + p_0' \xi)$$



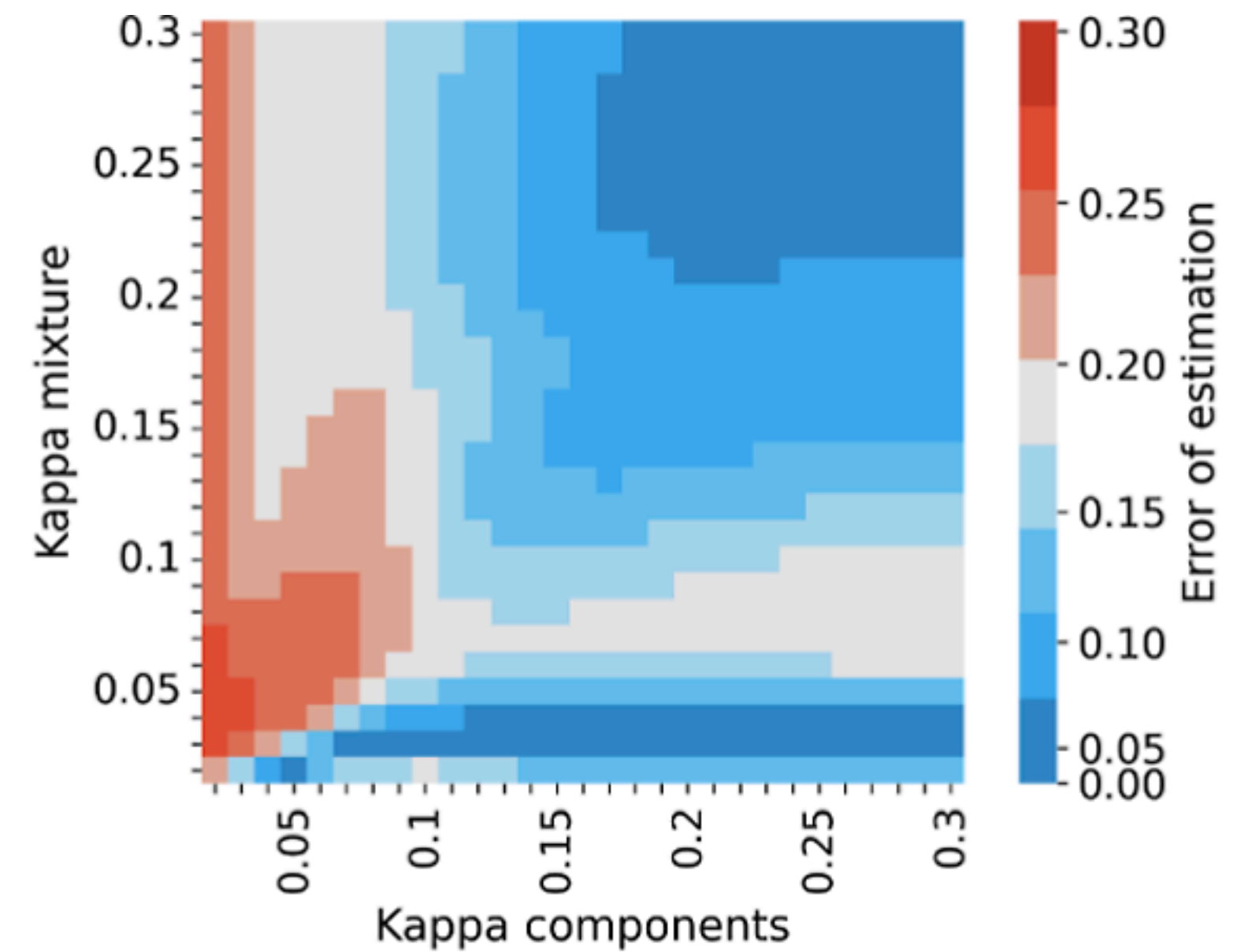
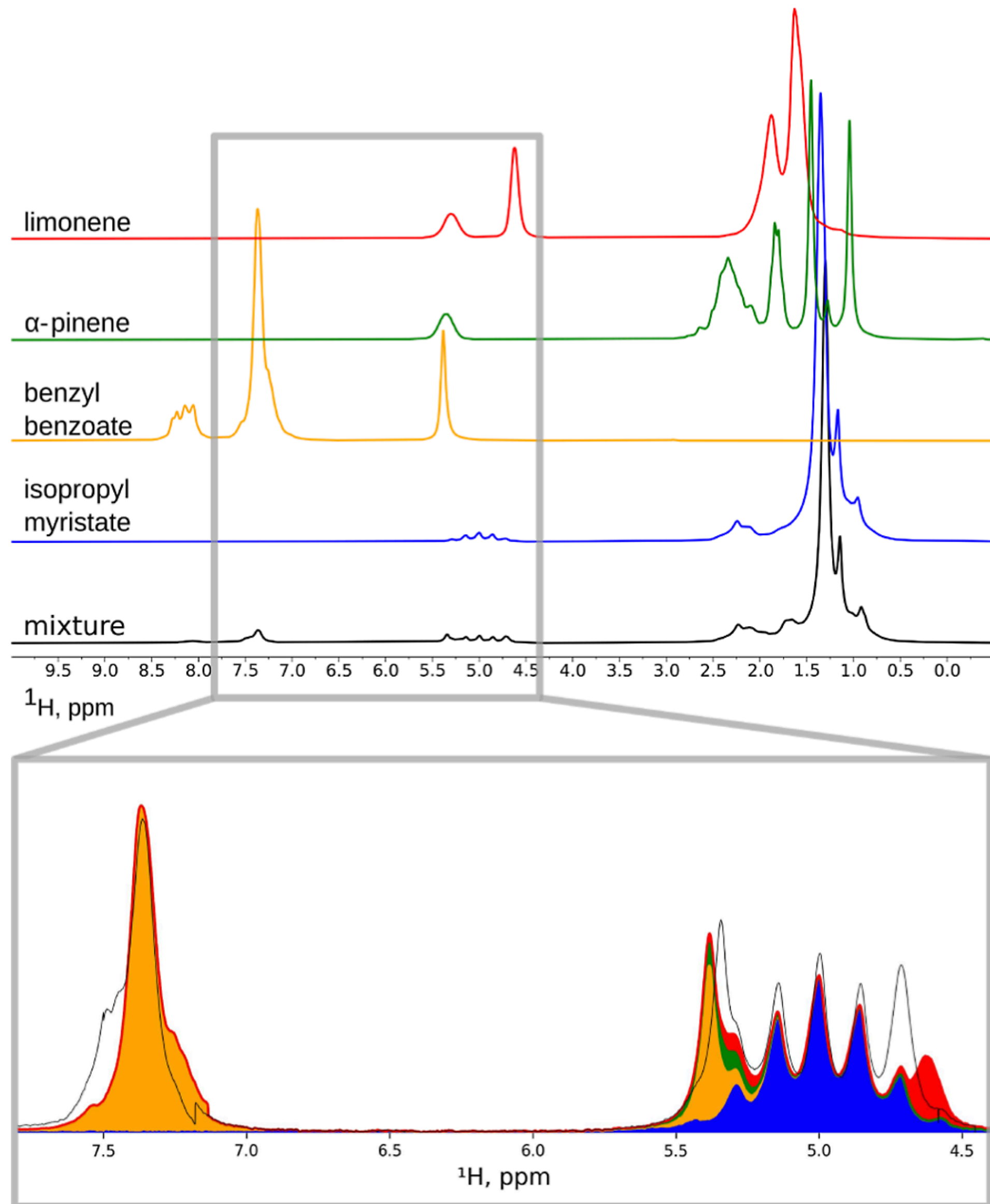
linear program

regression with Wasserstein distance can be formulated as linear program and solved efficiently



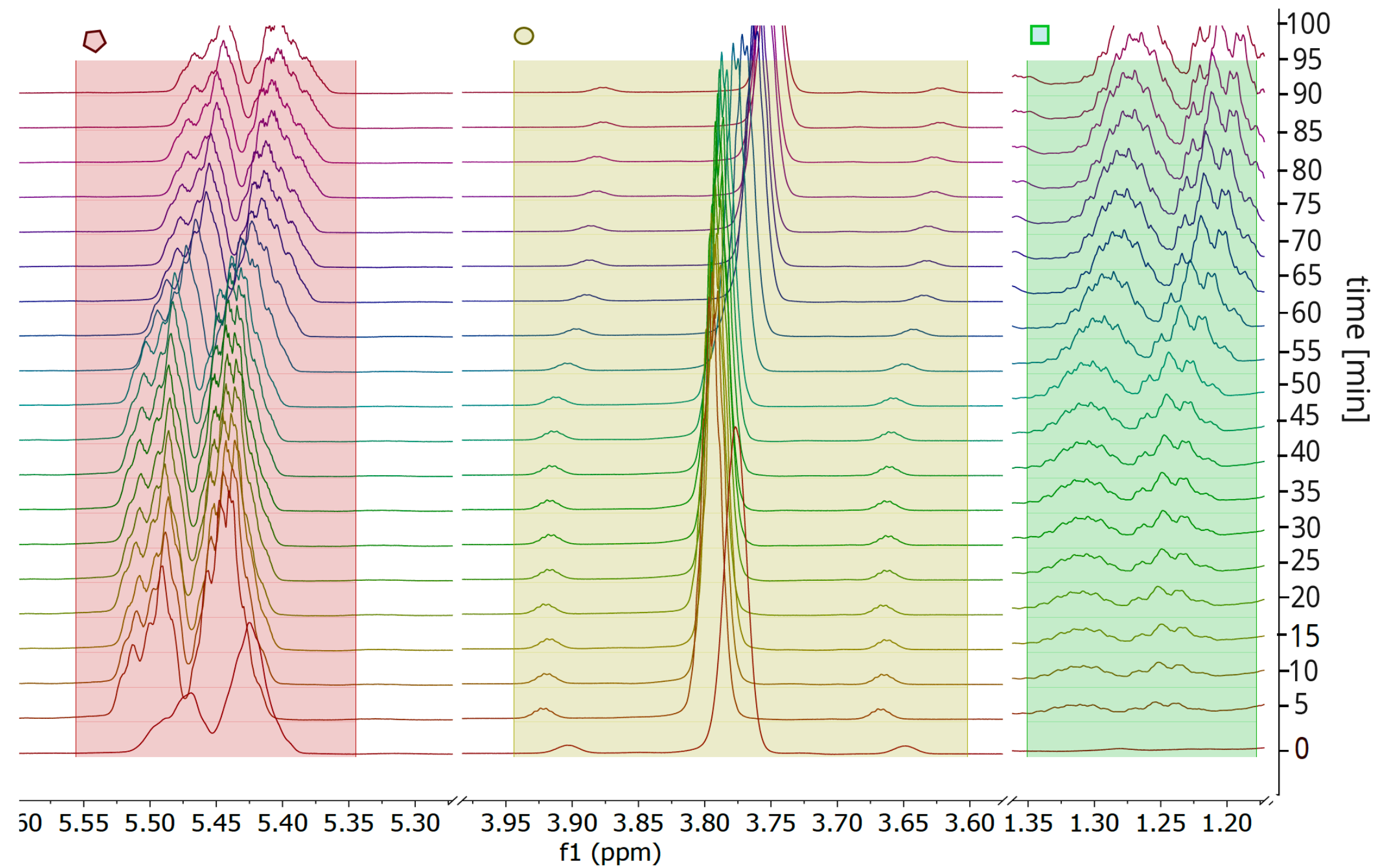
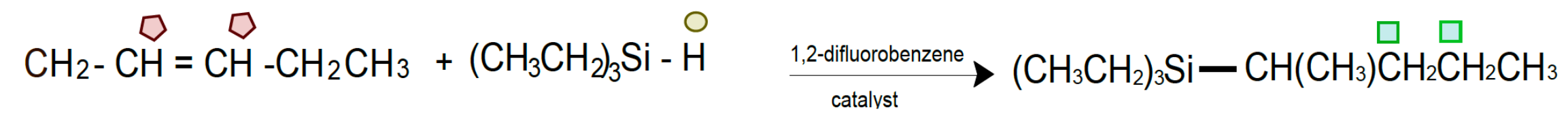
Magnetstein in action

Magnetstein can **quantitatively** analyze difficult spectra with the **estimation trueness** an order of magnitude **higher** than that of commercial tools...



... having **only two parameters** with default values applicable to a broad range of experiments...

chemical reactions revisited

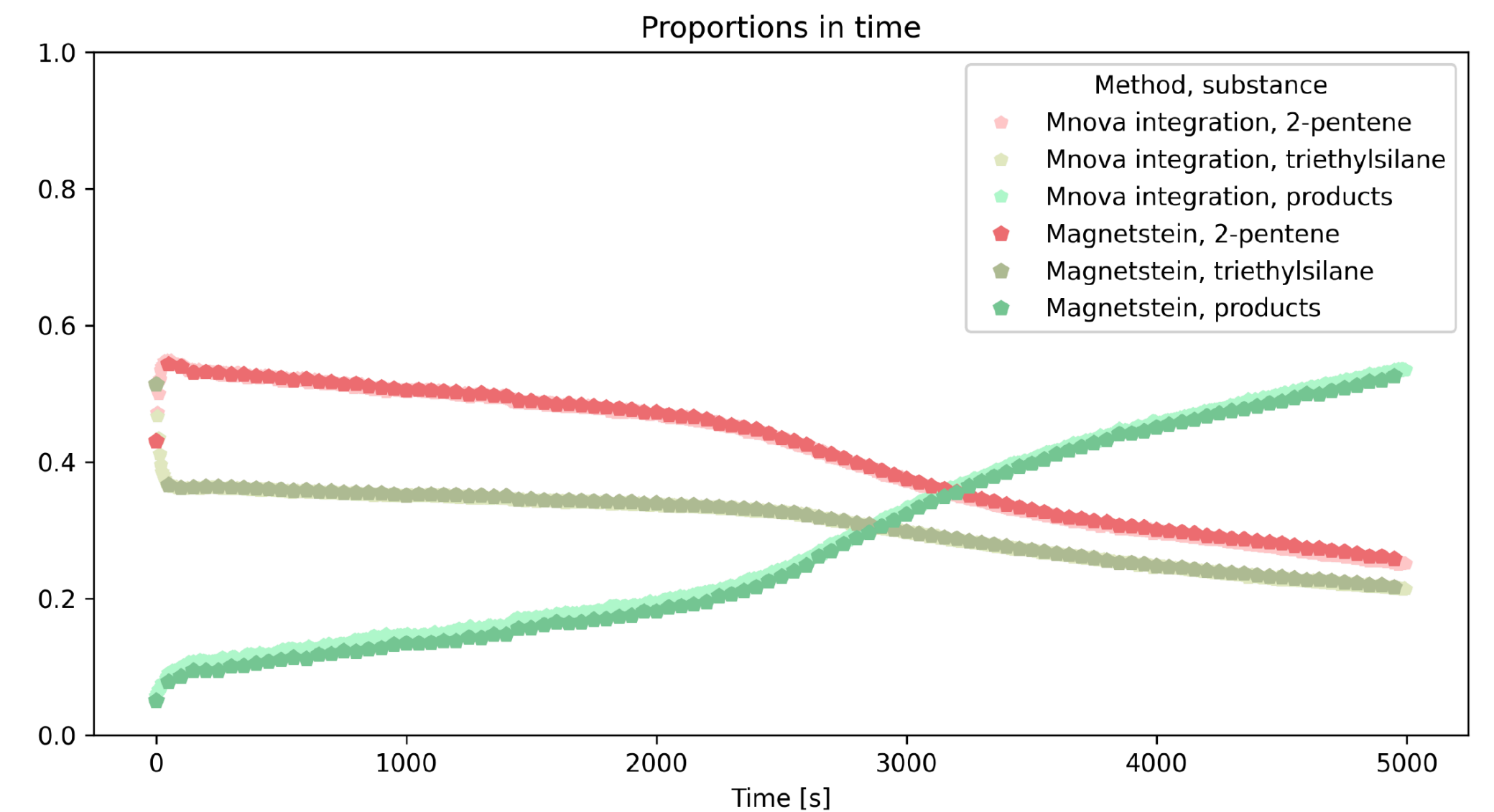


we can effectively quantify the components of a reacting mixture without a need for peak-picking

1. solve Wasserstein regression

2. get a sequence of proportions in consecutive timepoints

3. infer about kinetics of the monitored reaction



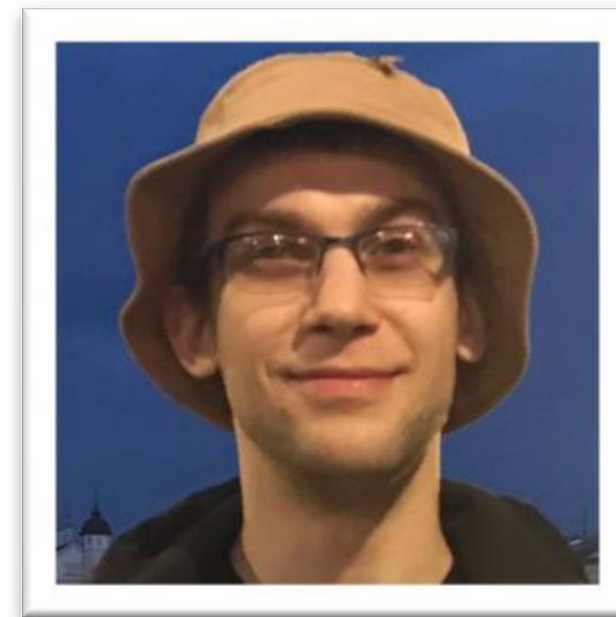
Many thanks to collaborators



**Michał
Startek**



**Dirk
Valkenborg**



**Grzegorz
Skoraczyński**



**Błażej
Miasojedow**



Frederik Lermyte



Mateusz Łacki



Krzysztof Kazimierczuk



Piotr Dittwald



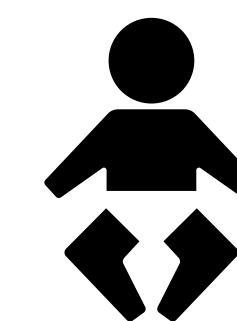
Frank Sobott



Michał Ciach



Barbara Domżał



Alan Rockwood

