# How to properly analyse spectra
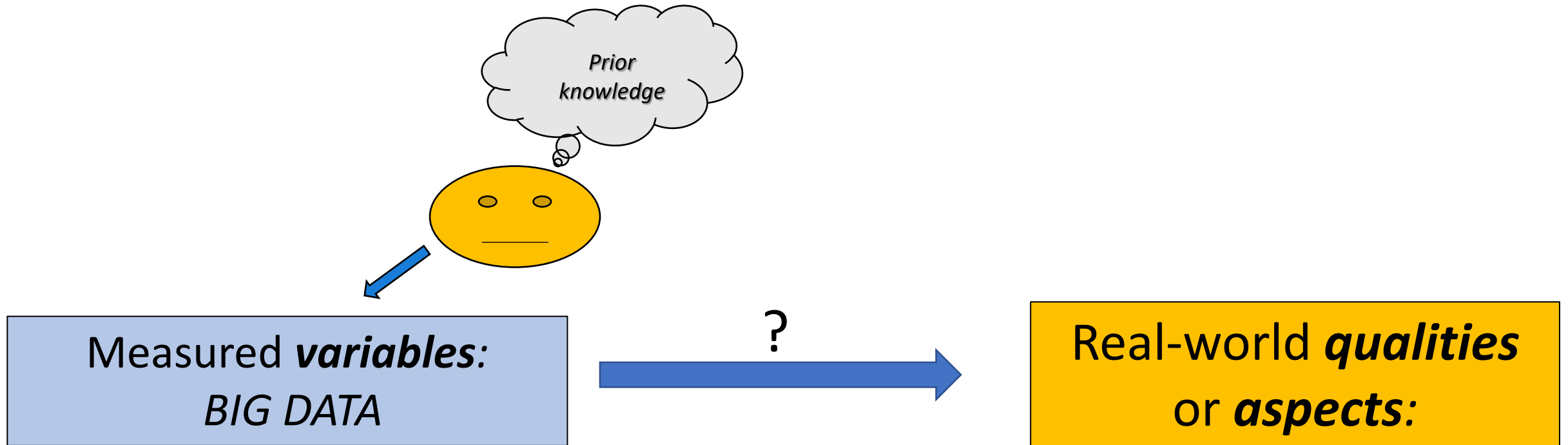
## Harald Martens, bio-chemometrician

Prof. emerit., DEPARTMENT OF ENGINEERING CYBERNETICS, Norw. U. of Sci. & Technol. NTNU, *Trondheim Norway*
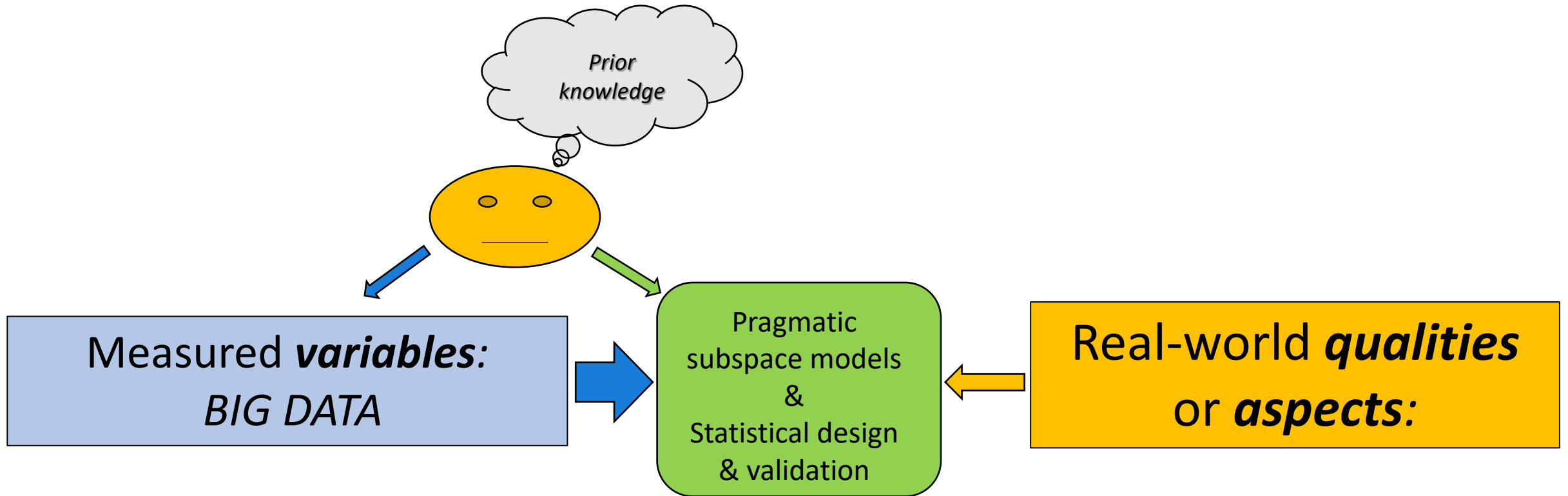
IDLETECHS AS (www.idletechs.com) *Trondheim Norway*

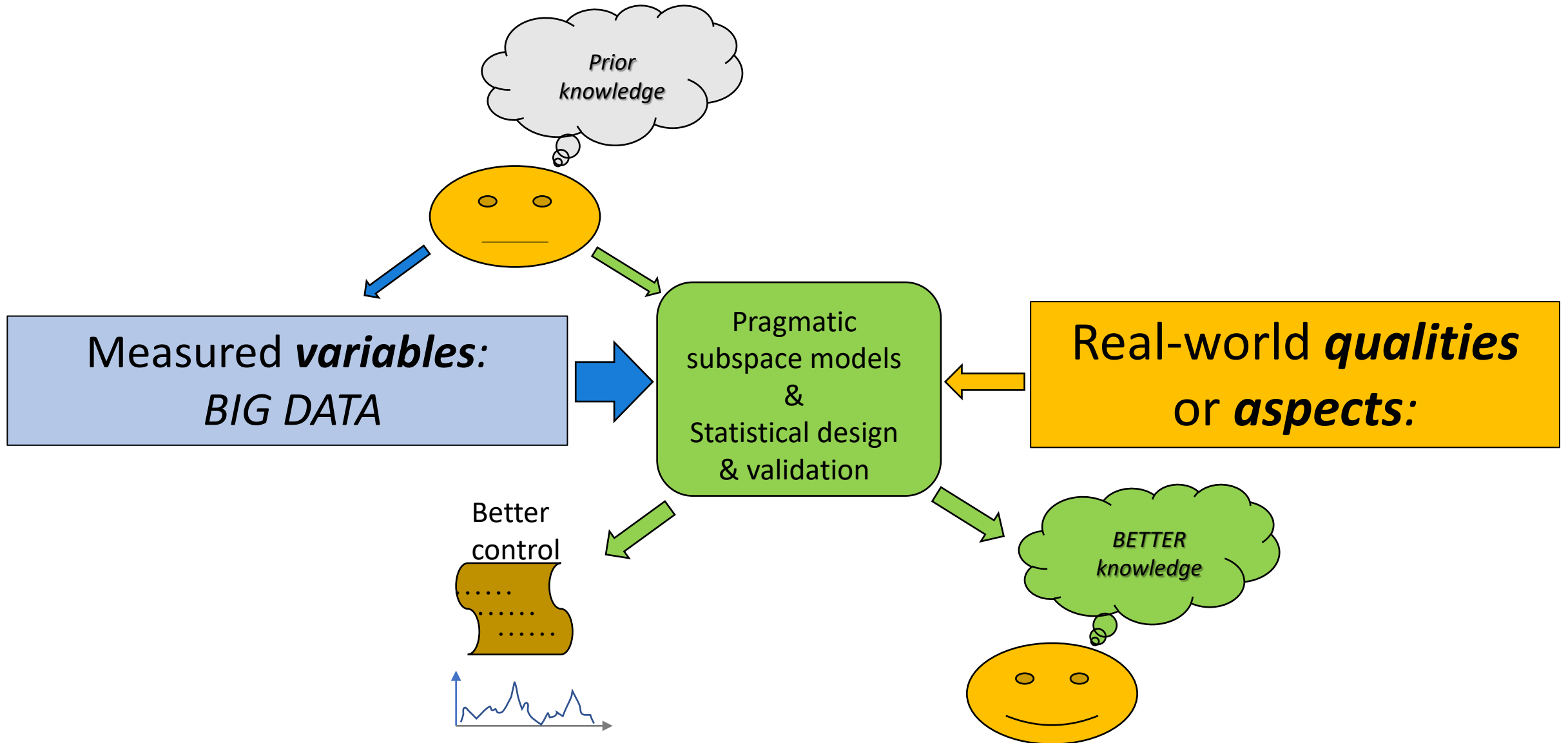# How to properly analyse spectra

## *My perspective after 50 years: A way to analyze spectra*

## Harald Martens, bio-chemometrician

Prof. emerit., DEPARTMENT OF ENGINEERING CYBERNETICS, Norw. U. of Sci. & Technol. NTNU, *Trondheim Norway*

IDLETECHS AS ([www.idletechs.com](www.idletechs.com)) *Trondheim Norway*

# Predicting something from many other things



*Prior knowledge*

Measured *variables*:
*BIG DATA*

?

Real-world *qualities*
or *aspects*:

# Predicting something from many other things



Prior knowledge

Measured *variables*:
*BIG DATA*

Pragmatic
subspace models
&
Statistical design
& validation

Real-world *qualities*
or *aspects*:

# Predicting something from many other things

BIG DATA

**Hyperspectral image of wood**

Multivariate data modelling

# Hybrid Chemometric Subspace Modelling

Principal components    #1    #2    #3

WHAT IS GOING ON HERE ???

BIG DATA

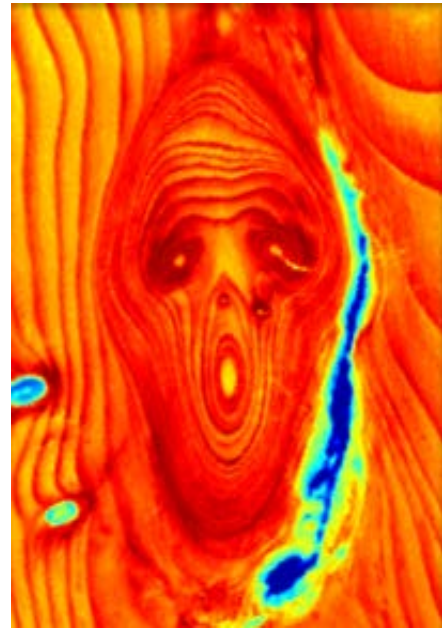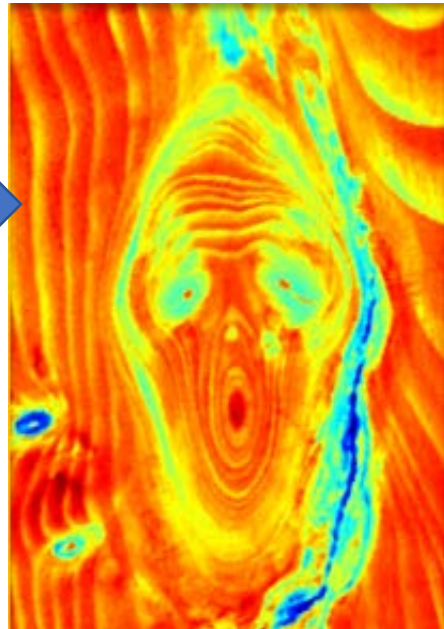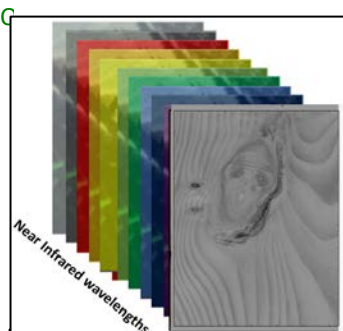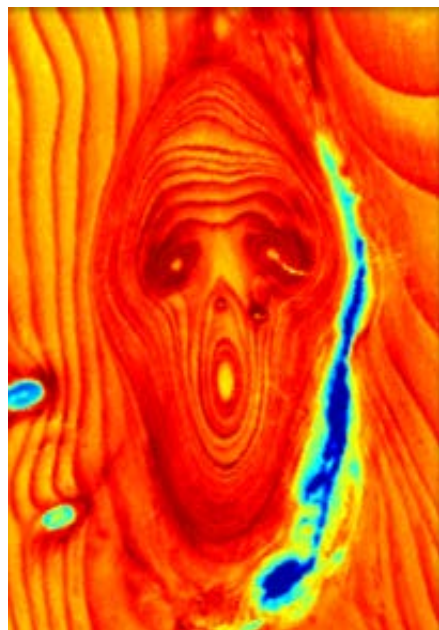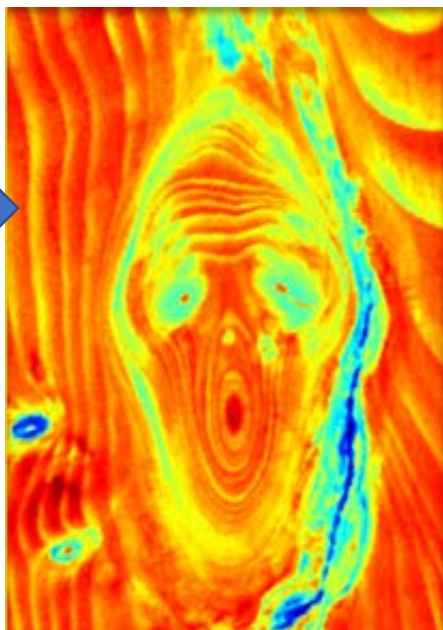# Hybrid Chemometric Subspace Modelling
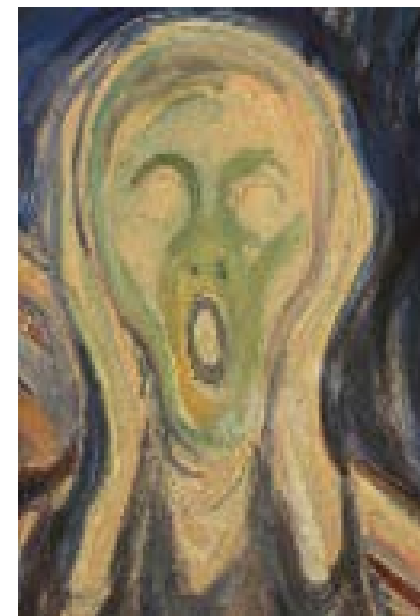
**Hyperspectral image of wood**

Multivariate data modelling

Edvard Munch: SCREAM
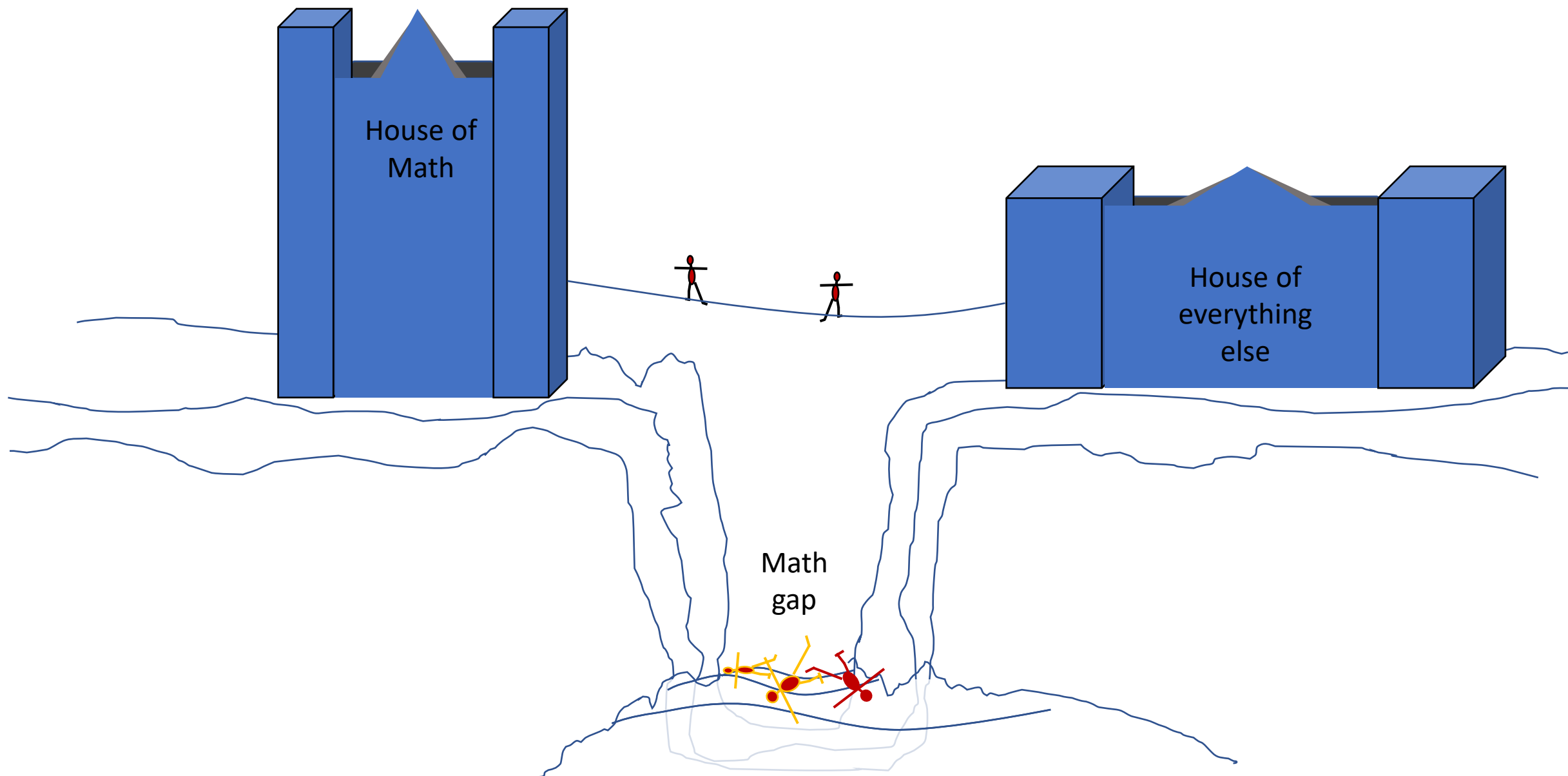


Principal components      #1                    #2                    #3

WHAT IS GOING ON HERE ???

MATHEMATICAL MODELLING ?

House of Math

House of everything else

Math gap

House of Math

House of everything else

Math gap

House of
Math
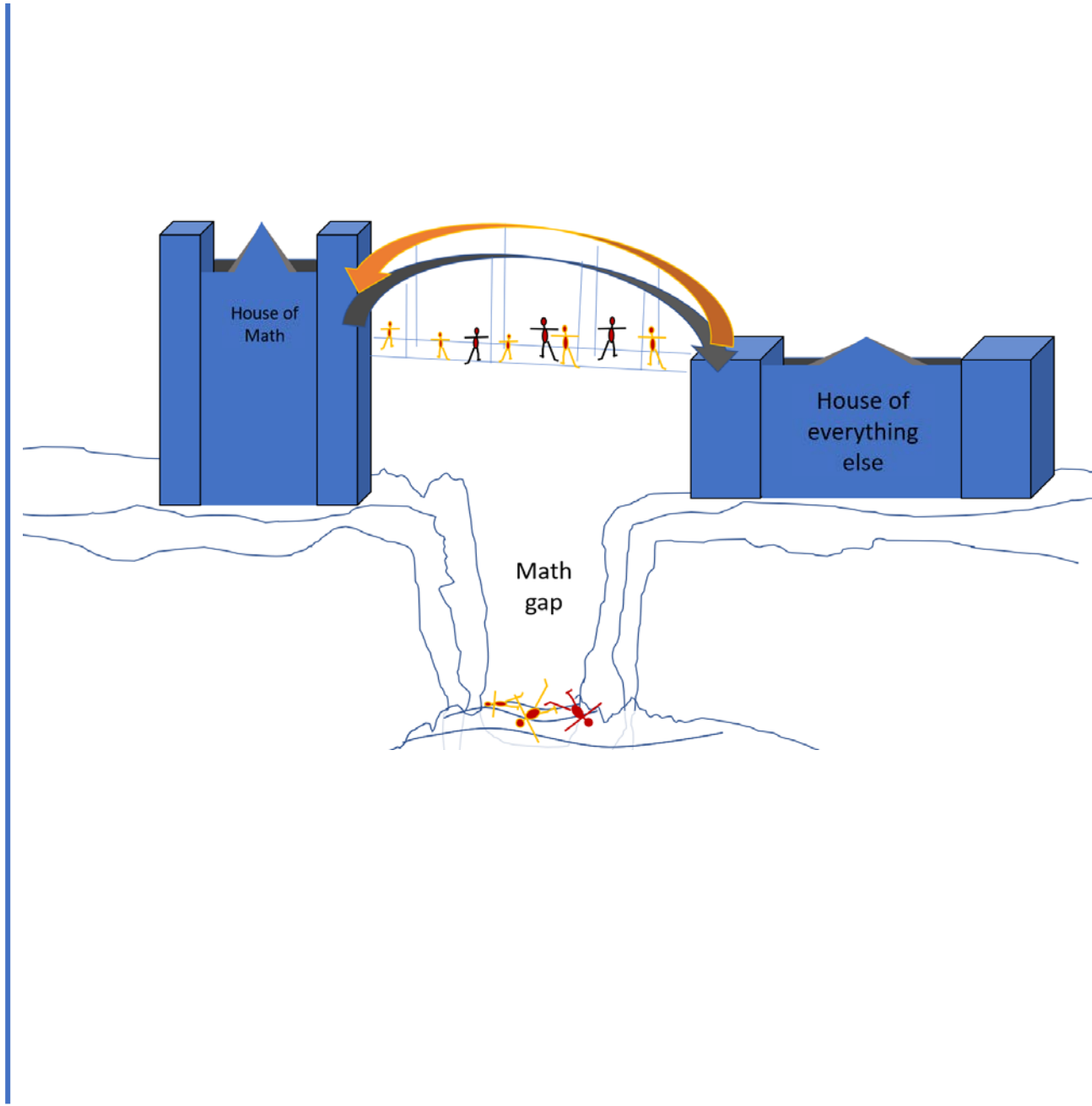
House of
everything
else

Math
gap

House of
Math

House of
everything
else

Math
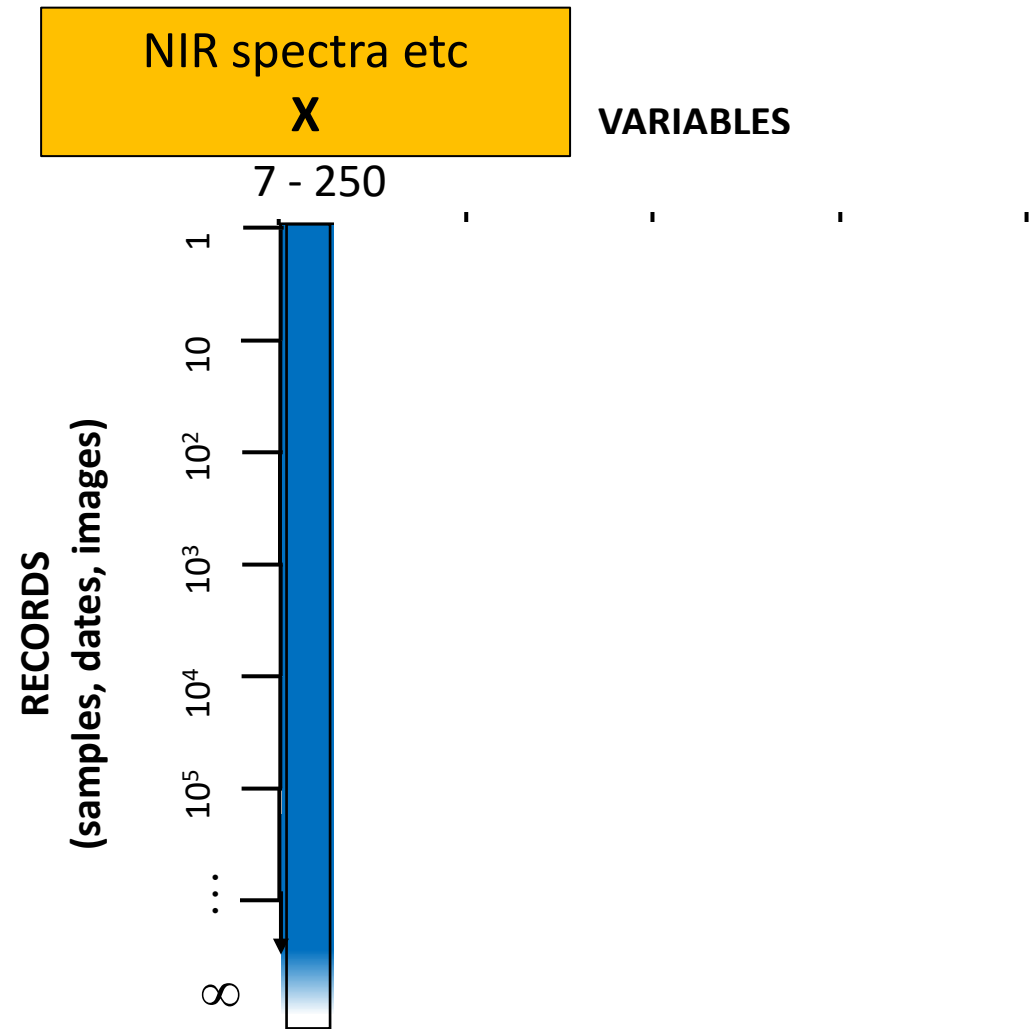gap

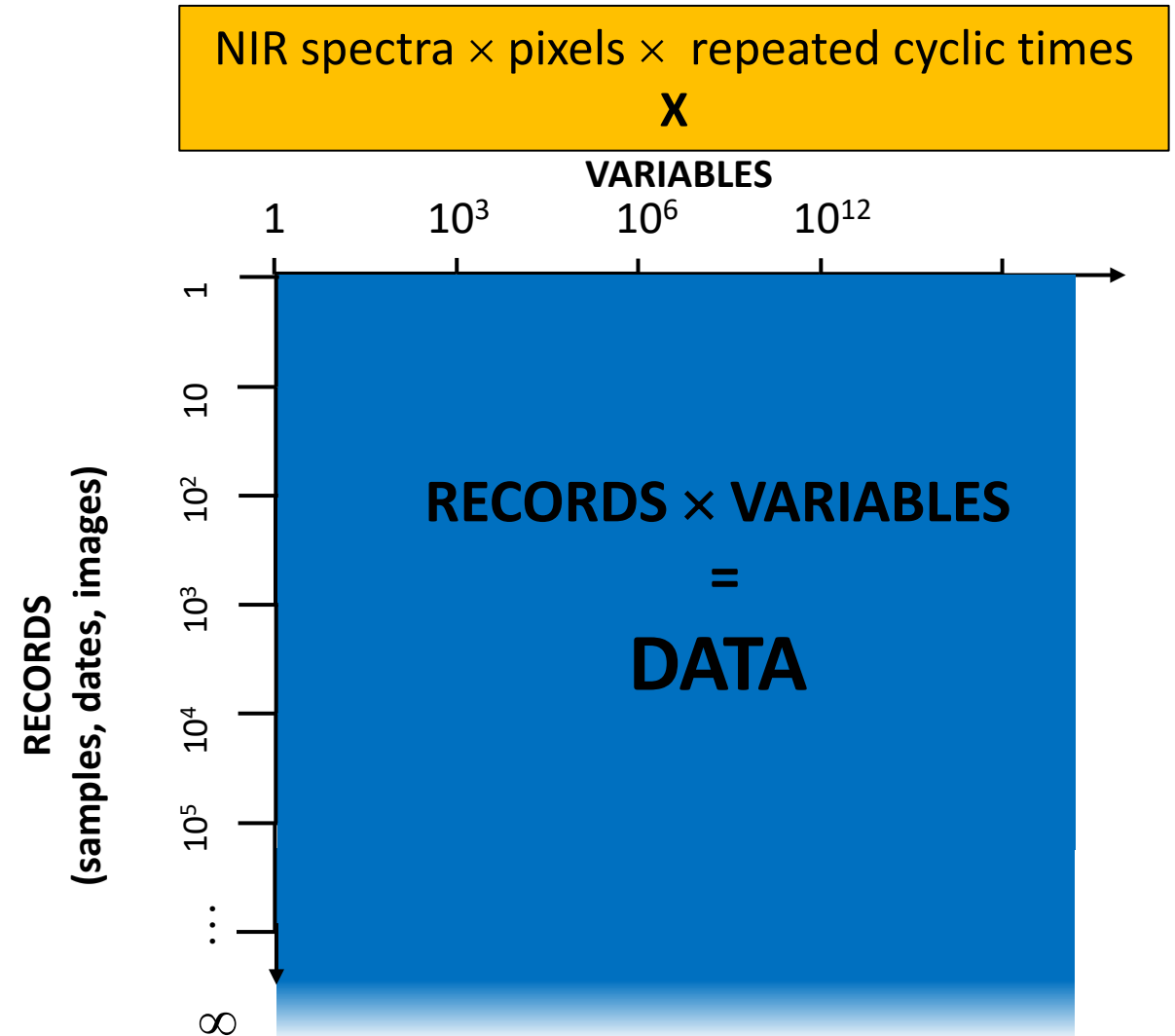# *A way to analyze spectra*

## Science in general :

What drives us?

I love
to do this

See how
good I am

Euro ⇒

This needs
to be done

# *A way to analyze spectra*

## NIR & Chemometrics:

What drives us?

*I love to do this*

*See how good I am*

**This needs to be done**

*Euro* ⇒

⇨ ***Euro***

# Input **DATA** = **RECORDS** $\times$ **VARIABLES**:

NIR spectra etc
**X**

**VARIABLES**

7 - 250

**RECORDS
(samples, dates, images)**

1

10

$10^2$

$10^3$

$10^4$

$10^5$

...

$\infty$

# Input **DATA = RECORDS** × **VARIABLES**:

NIR spectra × pixels × repeated cyclic times

**X**

**VARIABLES**

1      $10^3$    $10^6$    $10^{12}$

1

10

$10^2$

$10^3$

$10^4$

$10^5$

⋮

∞

**RECORDS (samples, dates, images)**

**RECORDS × VARIABLES**

**=**

**DATA**

# OTFP: Automatic modelling of continuous high-dimensional data streams
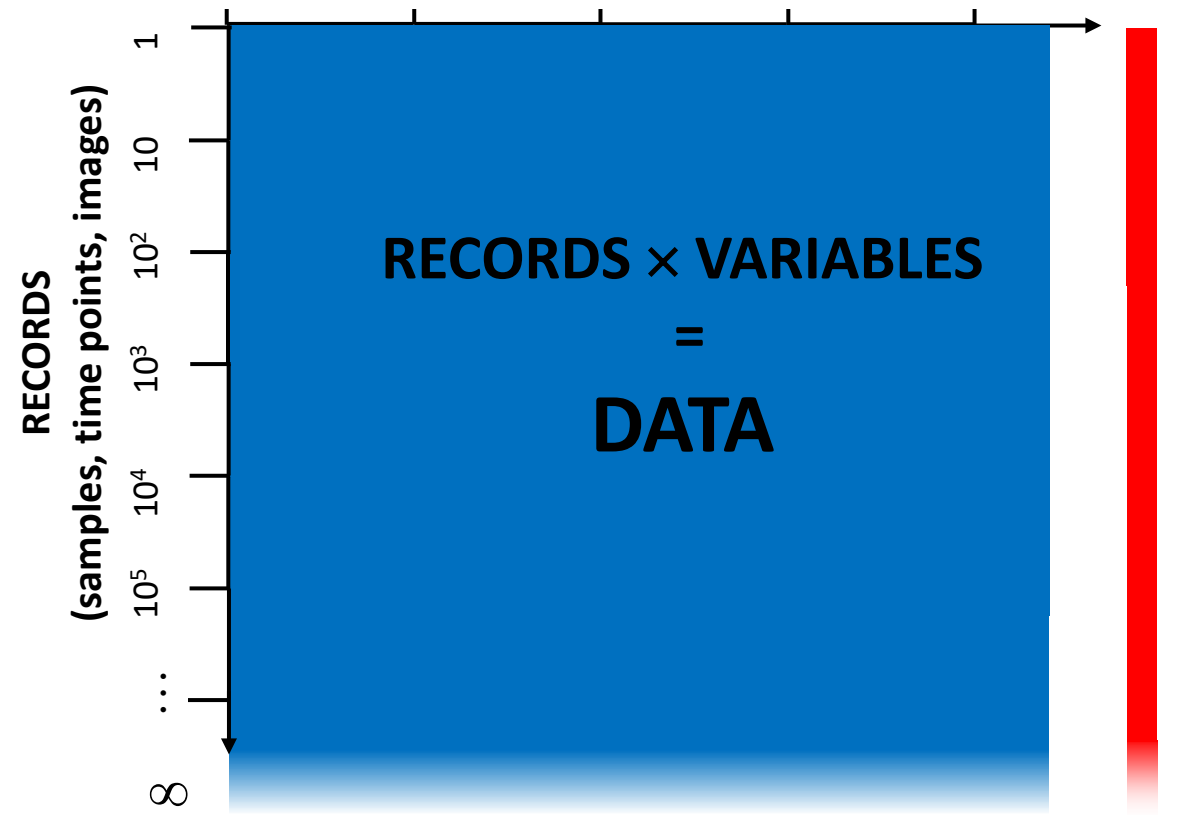


**On-The-Fly-Processing software for e.g. thermal – and hyperspectral video in industry (Vitale et al. 2017)**
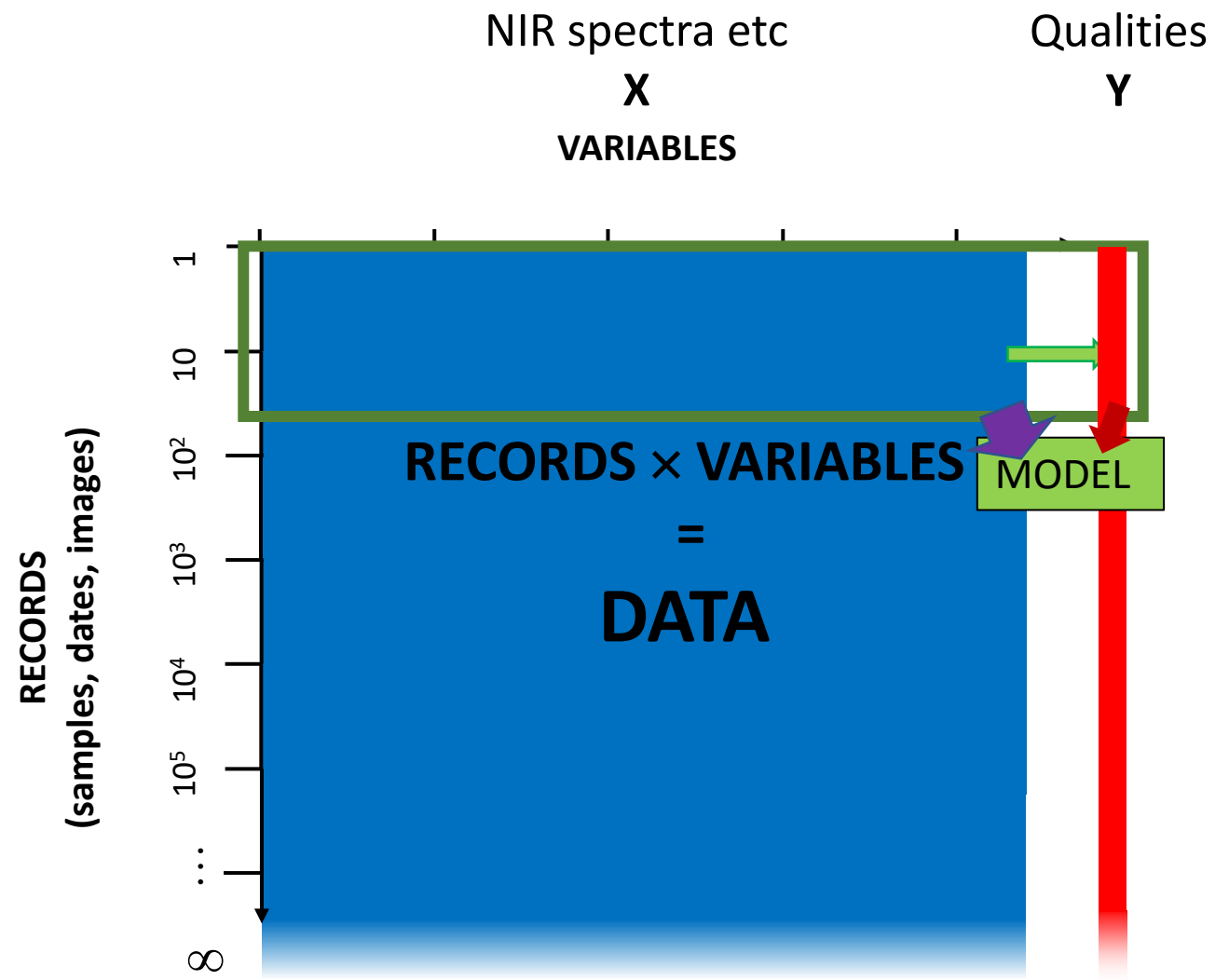
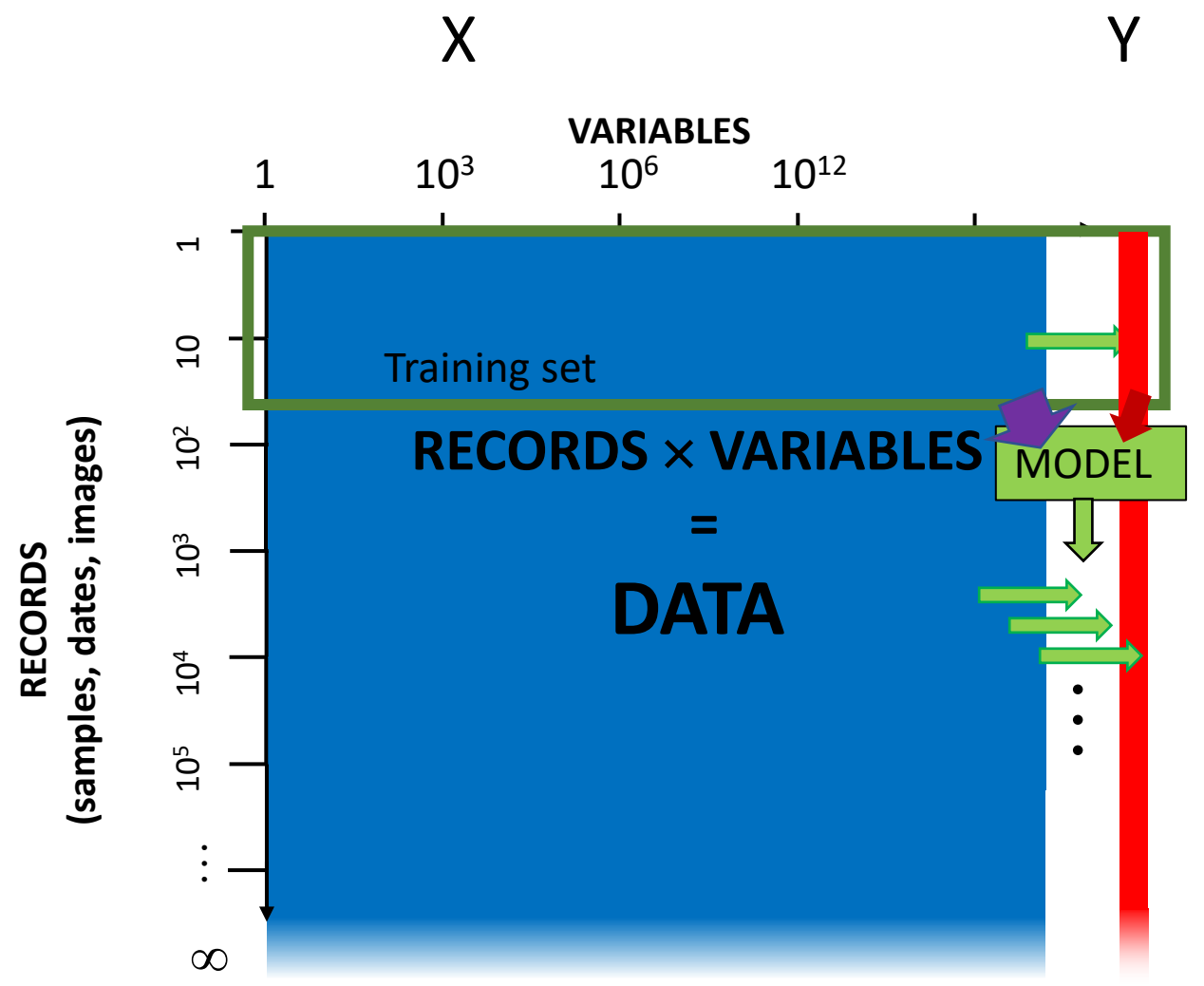# Input **DATA** = **RECORDS** × **VARIABLES**:



NIR spectra etc    →    Qualities
**X**           **Y**

**VARIABLES**

**RECORDS**
**(samples, time points, images)**

1
10
$10^2$
$10^3$
$10^4$
$10^5$
⋮
∞

**RECORDS × VARIABLES**
**=**
**DATA**

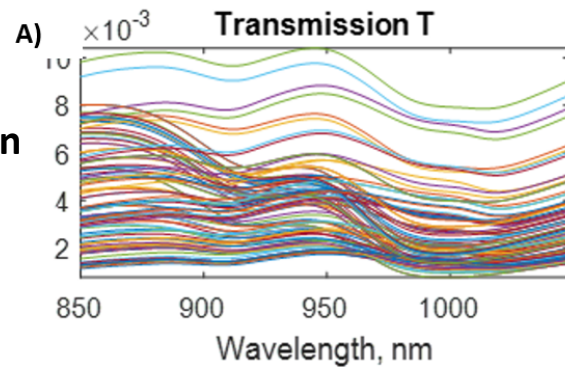# Input **DATA = RECORDS** × **VARIABLES:**
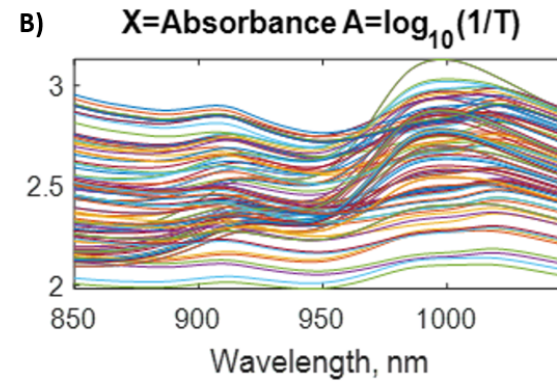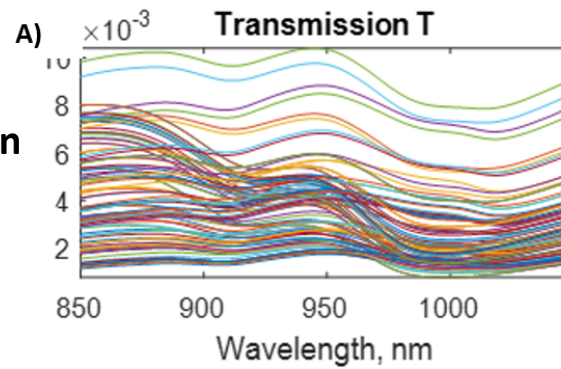
# Input **DATA** = **RECORDS** × **VARIABLES**:

# Five different powder mixtures
# measured by light transmission,
# each at varying sample - thickness and - compression

**Conventional linearization**
**+**
**multivariate calibration**
**(cross-validated PLSR)**



A) **Transmission T** $\times 10^{-3}$

B) **X=Absorbance A=$\log_{10}(1/T)$**

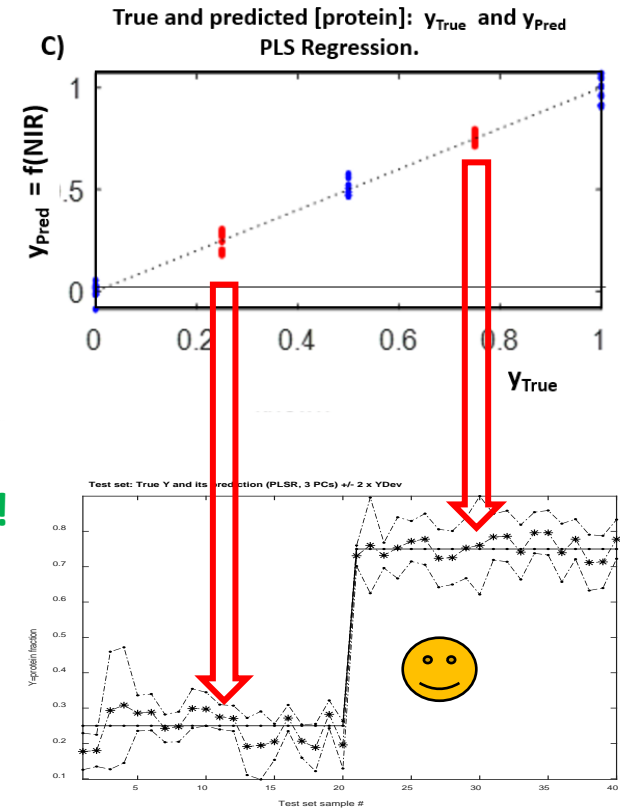C) True and predicted [protein]: $y_{True}$ and $y_{Pred}$
PLS Regression.

# Five different powder mixtures
# measured by light transmission,
# each at varying sample - thickness and - compression

**Conventional linearization**
**+**
**multivariate calibration**
**(cross-validated PLSR)**



A) Transmission T

B) X = Absorbance A = $\log_{10}(1/T)$

C) True and predicted [protein]: $y_{True}$ and $y_{Pred}$
PLS Regression.

**Predicted uncertainty: OK!**

Test set: True Y and its prediction (PLSR, 3 PCs) +/- 2 x YDev

# Five different powder mixtures
# measured by light transmission,
# each at varying sample - thickness and - compression

**Conventional linearization
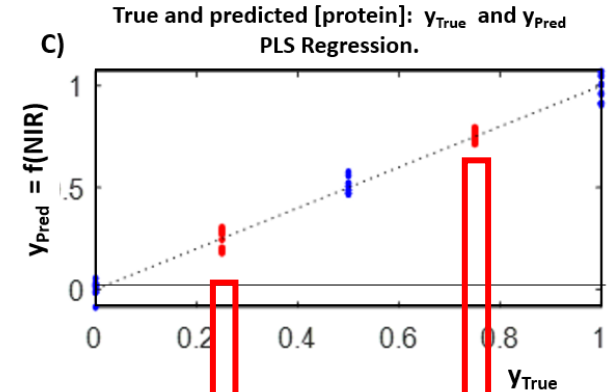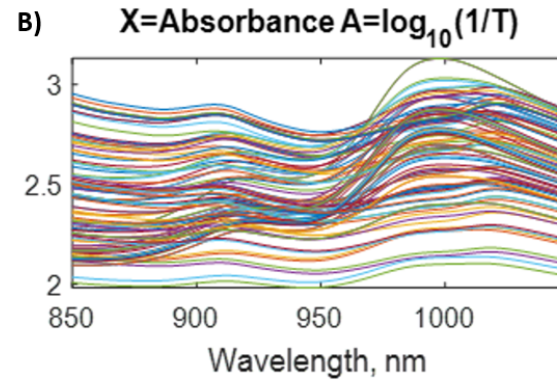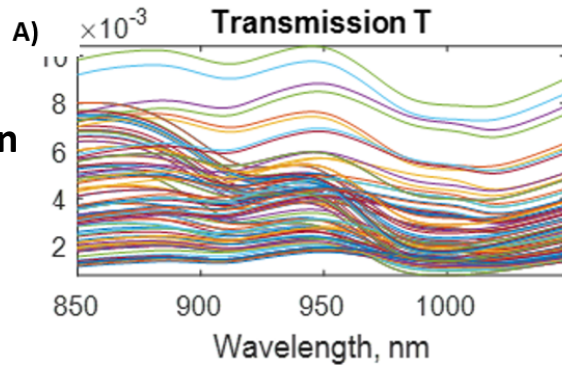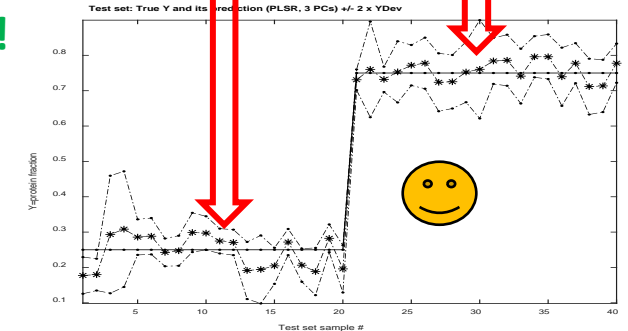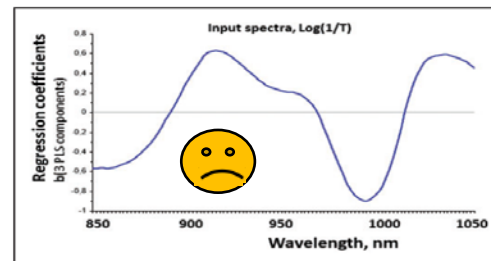+
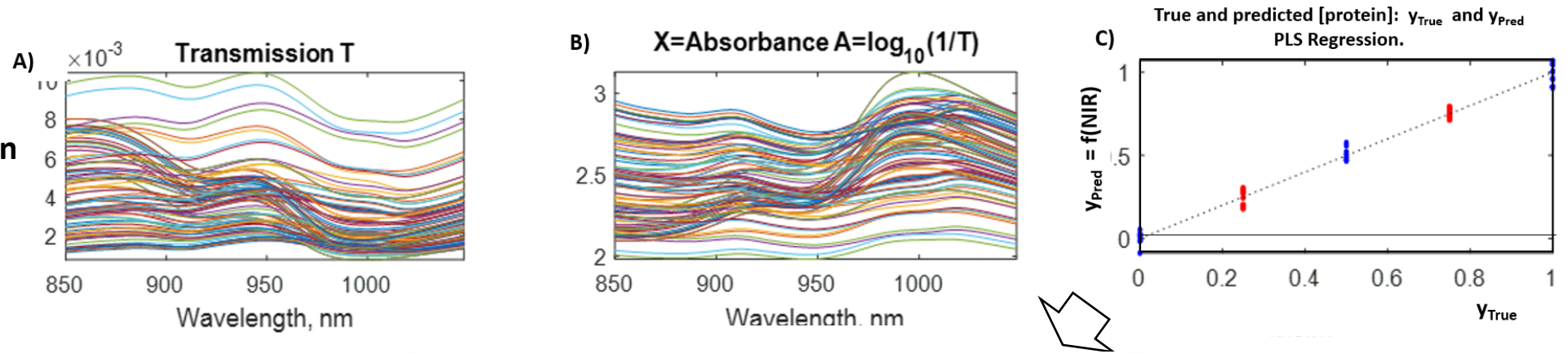multivariate calibration
(cross-validated PLSR)**



A) Transmission T

$\times 10^{-3}$

B) $X = \text{Absorbance } A = \log_{10}(1/T)$

C) True and predicted [protein]: $y_{True}$ and $y_{Pred}$ PLS Regression.

$y_{Pred} = f(NIR)$  vs  $y_{True}$

**Predicted uncertainty: OK!**

Test set: True Y and its prediction (PLSR, 3 PCs) +/- 2 x YDev

**Regression coefficient =
«Net Analyte Signal»:**

Input spectra, Log(1/T)

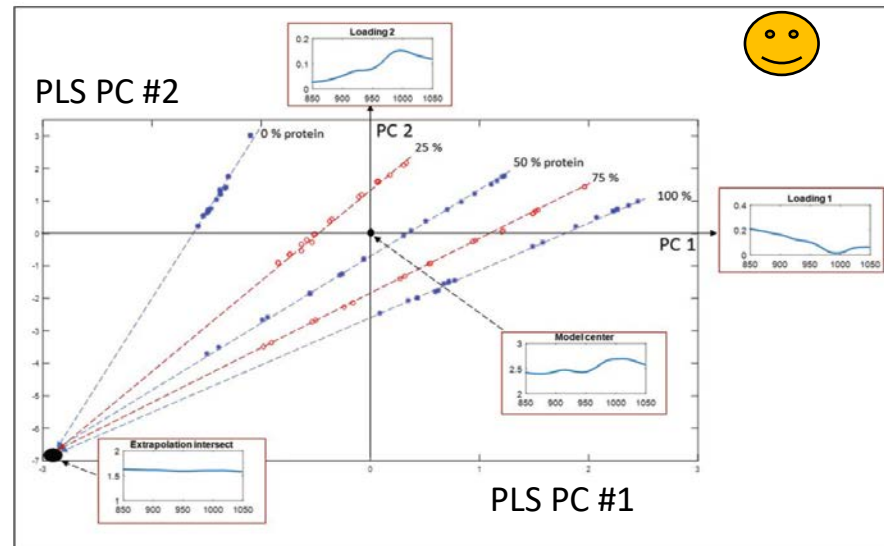Regression coefficients b(3 PLS components)

# Five different powder mixtures measured by light transmission, each at varying sample - thickness and - compression

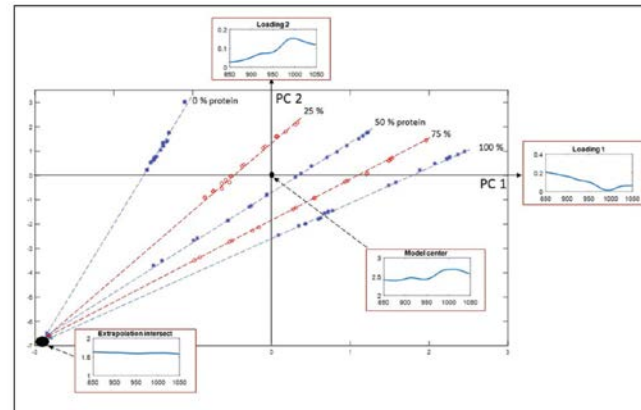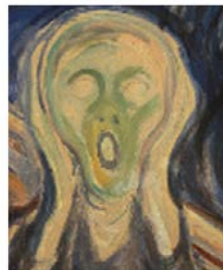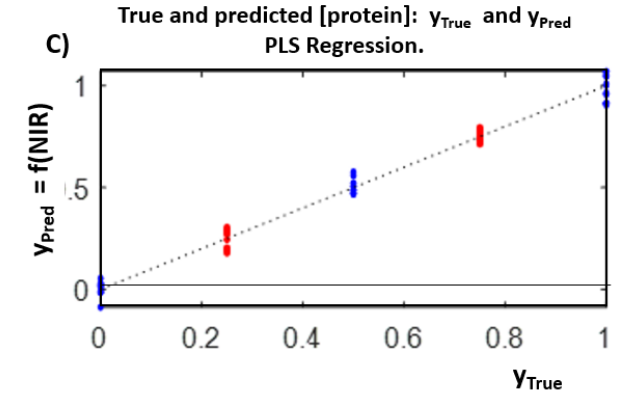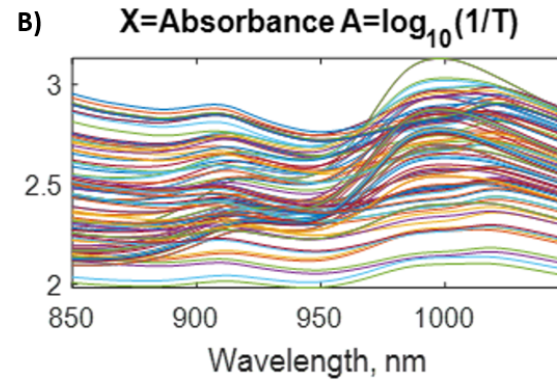**Conventional linearization + multivariate calibration (cross-validated PLSR)**

A) Transmission T

B) $X = $ Absorbance $A = \log_{10}(1/T)$

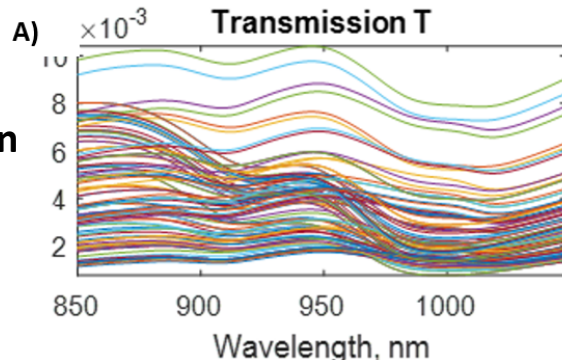C) True and predicted [protein]: $y_{True}$ and $y_{Pred}$
PLS Regression.

*PCA & PLS regression etc : Low-dimensional subspace ! graphic insight* $\Rightarrow$

# Five different powder mixtures
# measured by light transmission,
# each at varying sample - thickness and - compression

**Conventional linearization**
**+**
**multivariate calibration**
**(cross-validated PLSR)**



A) Transmission T

B) $X = Absorbance\ A = \log_{10}(1/T)$

C) True and predicted [protein]: $y_{True}$ and $y_{Pred}$
PLS Regression.

$y_{Pred} = f(NIR)$

$y_{True}$

**OEMSC linearization**
**+**
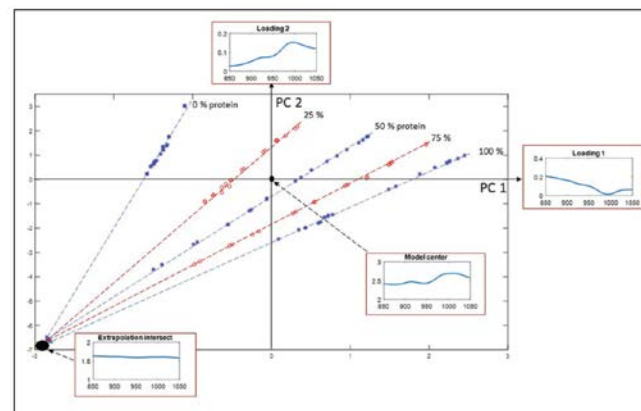**multivariate calibration**
**(cross-validated PLSR)**
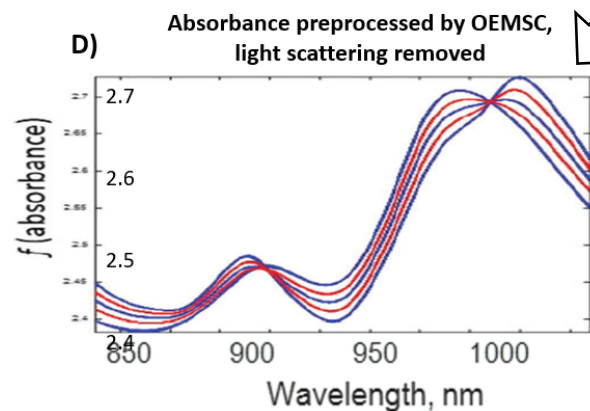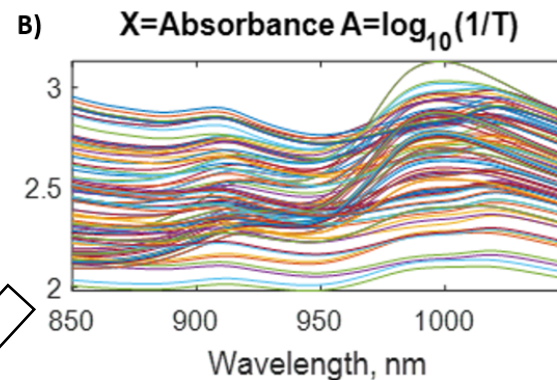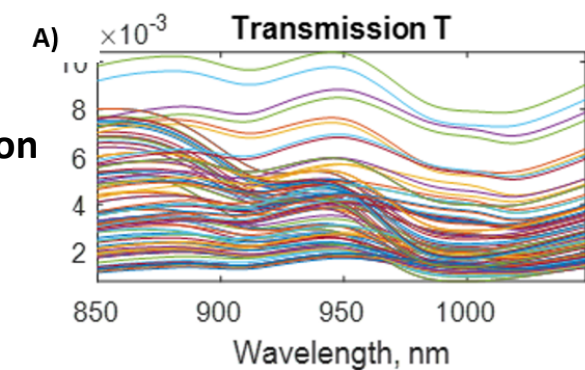


MATHEMATICAL MODELLING ?

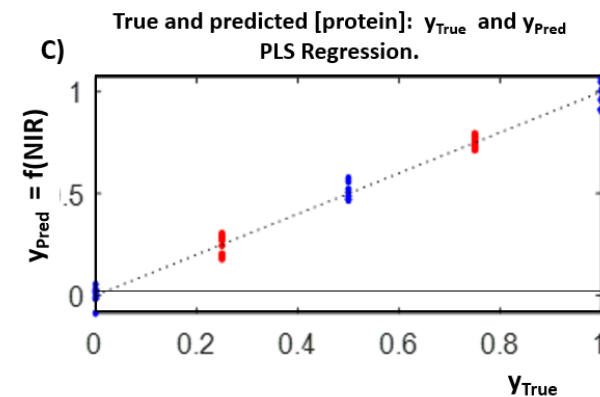$A \approx B \times C + D$

Subspace inspection
( two first PLS PCs)

# Five different powder mixtures measured by light transmission, each at varying sample - thickness and - compression

**Conventional linearization**
**+**
**multivariate calibration**
**(cross-validated PLSR)**



A) $\times 10^{-3}$ **Transmission T**

B) **X=Absorbance A=$\log_{10}(1/T)$**

C) True and predicted [protein]: $y_{True}$ and $y_{Pred}$ PLS Regression.

$y_{Pred} = f(NIR)$

$y_{True}$

**OEMSC linearization**
**+**
**multivariate calibration**
**(cross-validated PLSR)**

D) Absorbance preprocessed by OEMSC, light scattering removed
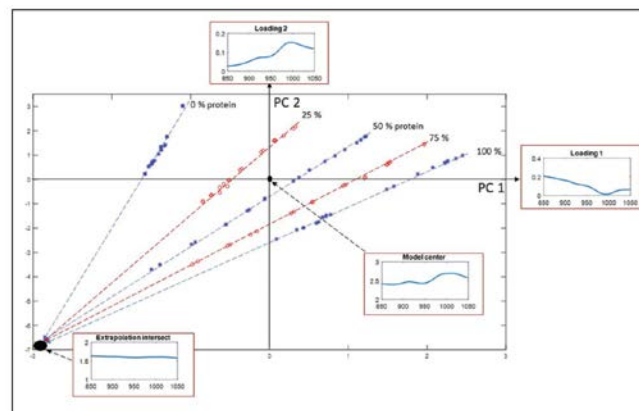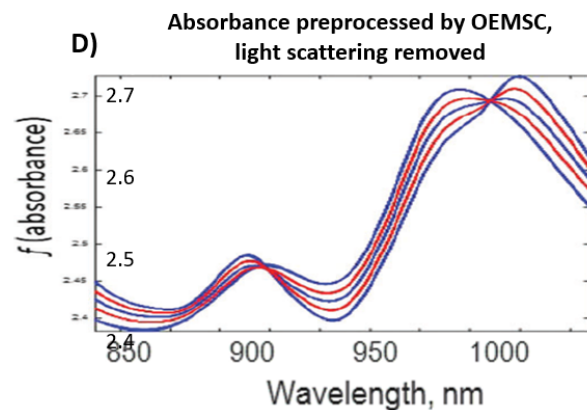
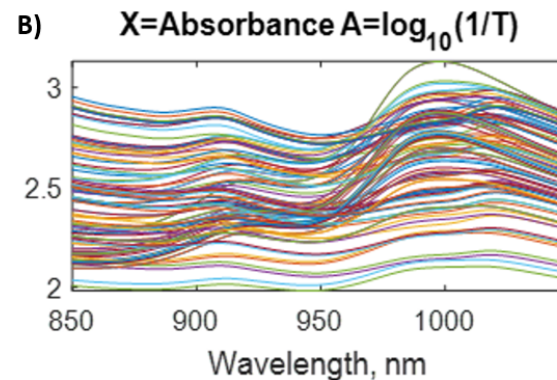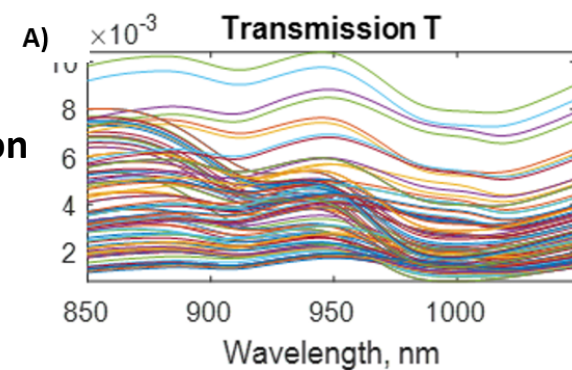$f$ (absorbance)

Retaining only
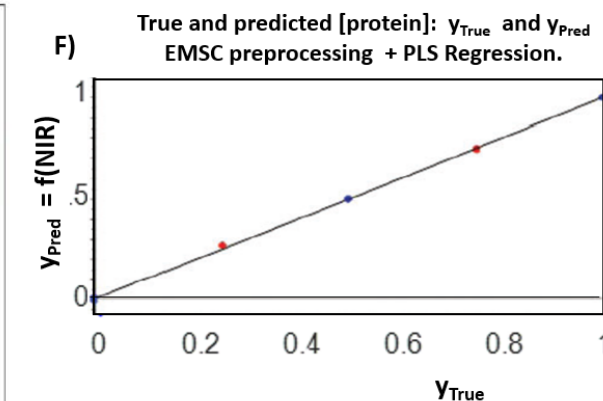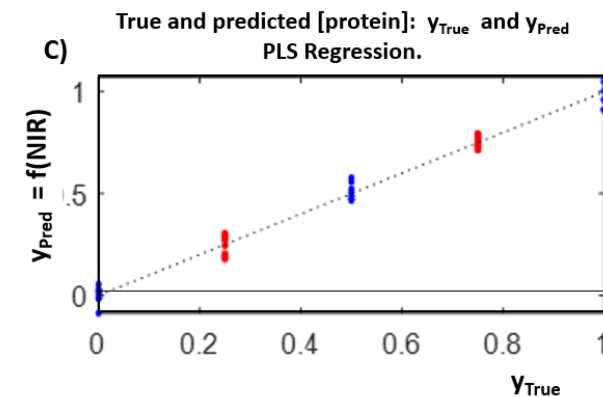chemical info

Subspace inspection
( two first PLS PCs)

# Five different powder mixtures measured by light transmission, each at varying sample - thickness and - compression

**Conventional linearization**
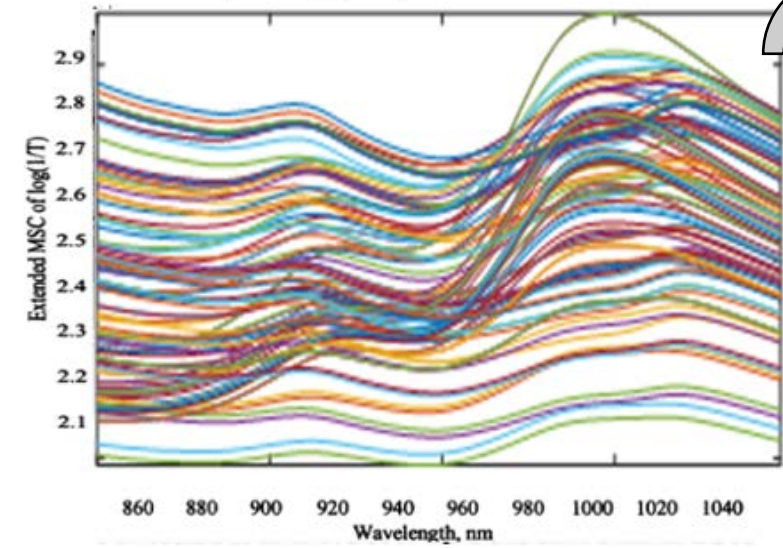**+**
**multivariate calibration (cross-validated PLSR)**



A) Transmission T

B) $X=\text{Absorbance } A=\log_{10}(1/T)$

C) True and predicted [protein]: $y_{True}$ and $y_{Pred}$ PLS Regression.

**OEMSC linearization**
**+**
**multivariate calibration (cross-validated PLSR)**



D) Absorbance preprocessed by OEMSC, light scattering removed

F) True and predicted [protein]: $y_{True}$ and $y_{Pred}$ EMSC preprocessing + PLS Regression.

Retaining only chemical info

Subspace inspection ( two first PLS PCs)

A: Input log(1/T)
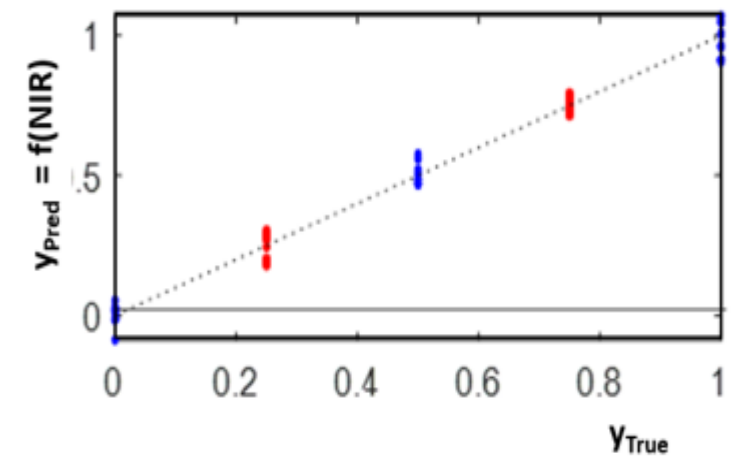


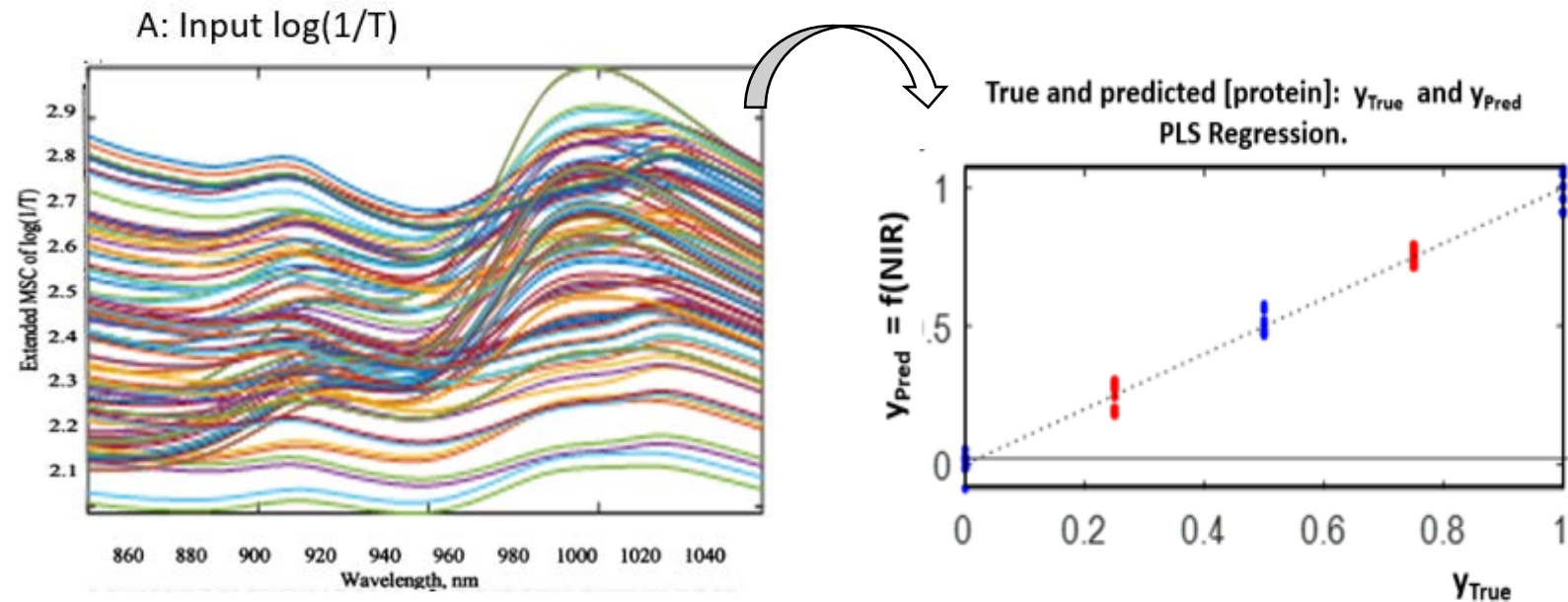True and predicted [protein]: $y_{True}$ and $y_{Pred}$
PLS Regression.

A: Input log(1/T)

True and predicted [protein]: $y_{True}$ and $y_{Pred}$
PLS Regression.

EMSC:   Extended Multiplicative Signal Correction
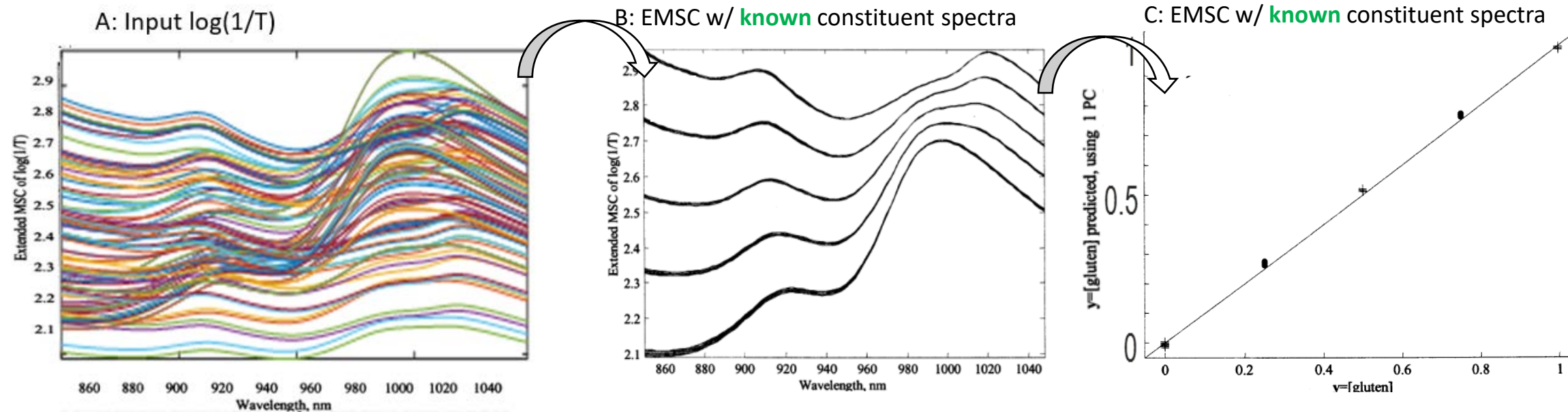
Simple linear model, using high-school algebra:

X=NIR absorbance log(1/T) or log(1/R)

A=Spectral knowledge

$X \approx A \times B + C$ :        Find B and C,         then $X_{corrected} = (X-C)/B$

A: Input log(1/T)

B: EMSC w/ **known** constituent spectra

C: EMSC w/ **known** constituent spectra

# EMSC:   Extended Multiplicative Signal Correction

Simple linear model, using high-school algebra:

X=NIR absorbance log(1/T) or log(1/R)
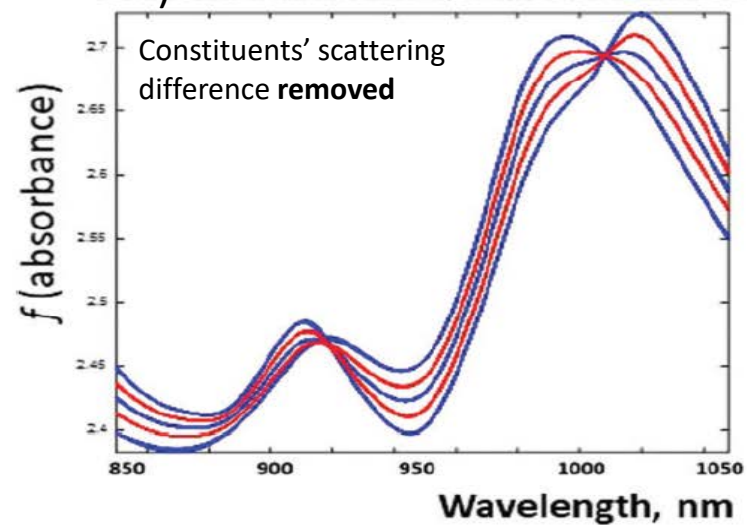
A=Spectral knowledge

$X \approx A \times B + C :$     Find B and C,        then $X_{corrected}=(X-C)/B$

A: Input log(1/T)

B: EMSC w/ **known** constituent spectra

C: EMSC w/ **known** constituent spectra

D: OEMSC w/ **unknown** constituent spectra

Constituents' scattering difference **removed**

E: OEMSC w/ **unknown** constituent spectra

Constituents' scattering difference **retained**
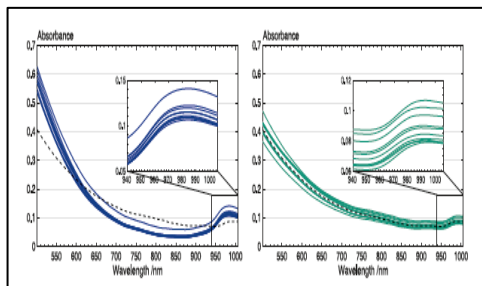
F: OEMSC w/ **unknown** constituent spectra

# Big Data: Hyperspectral «video»
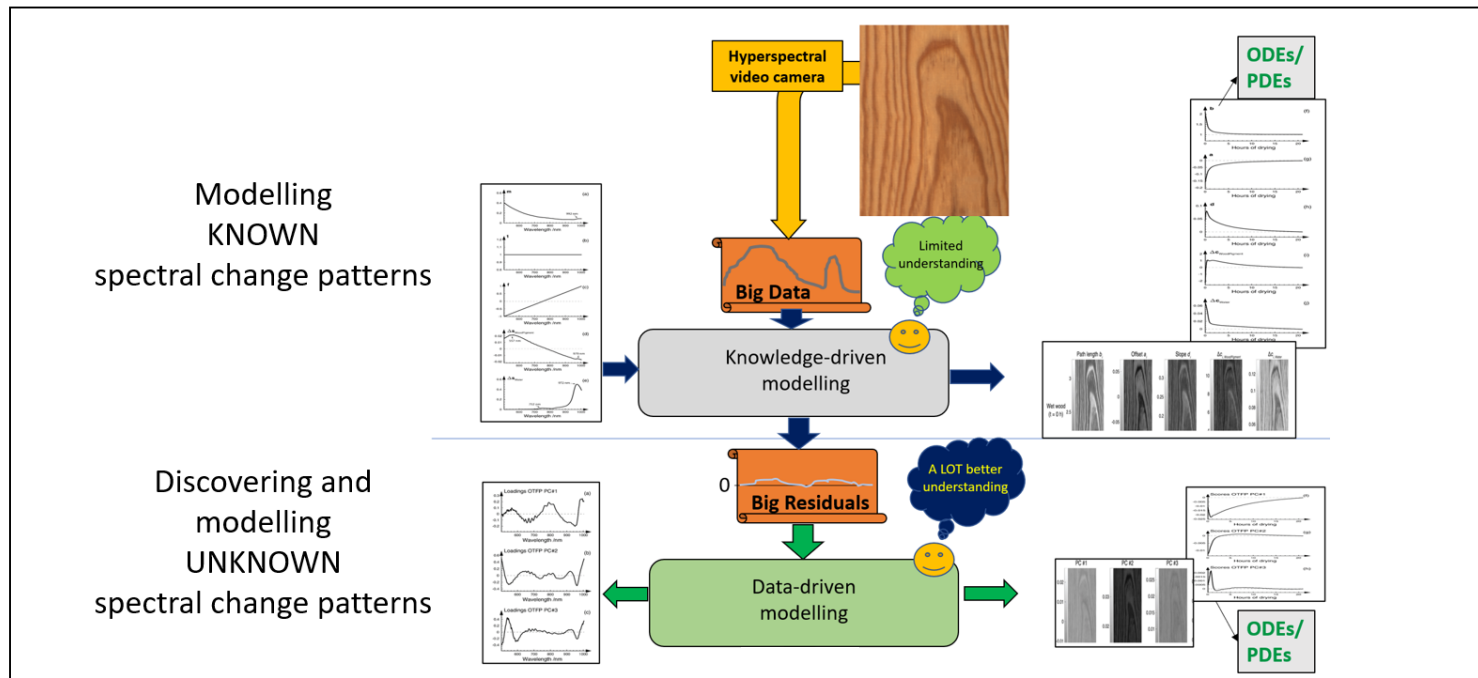 A single piece of drying wood:
>350 000 000 VNIR reflectance spectra,
≈200 channels each, measured at 150 consecutive times

VNIR;
400-1000 nm

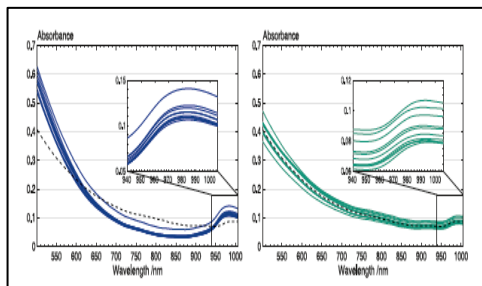**EMSC modelling KNOWN and UNKNOWN physics & chemistry:**

# Big Data: Hyperspectral «video»
 A single piece of drying wood:
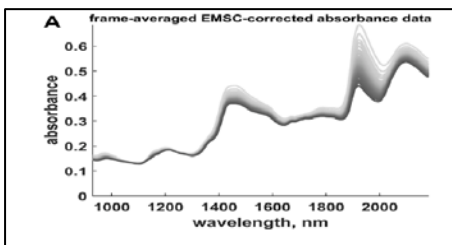>350 000 000 VNIR reflectance spectra,
≈200 channels each, measured at 150 consecutive times
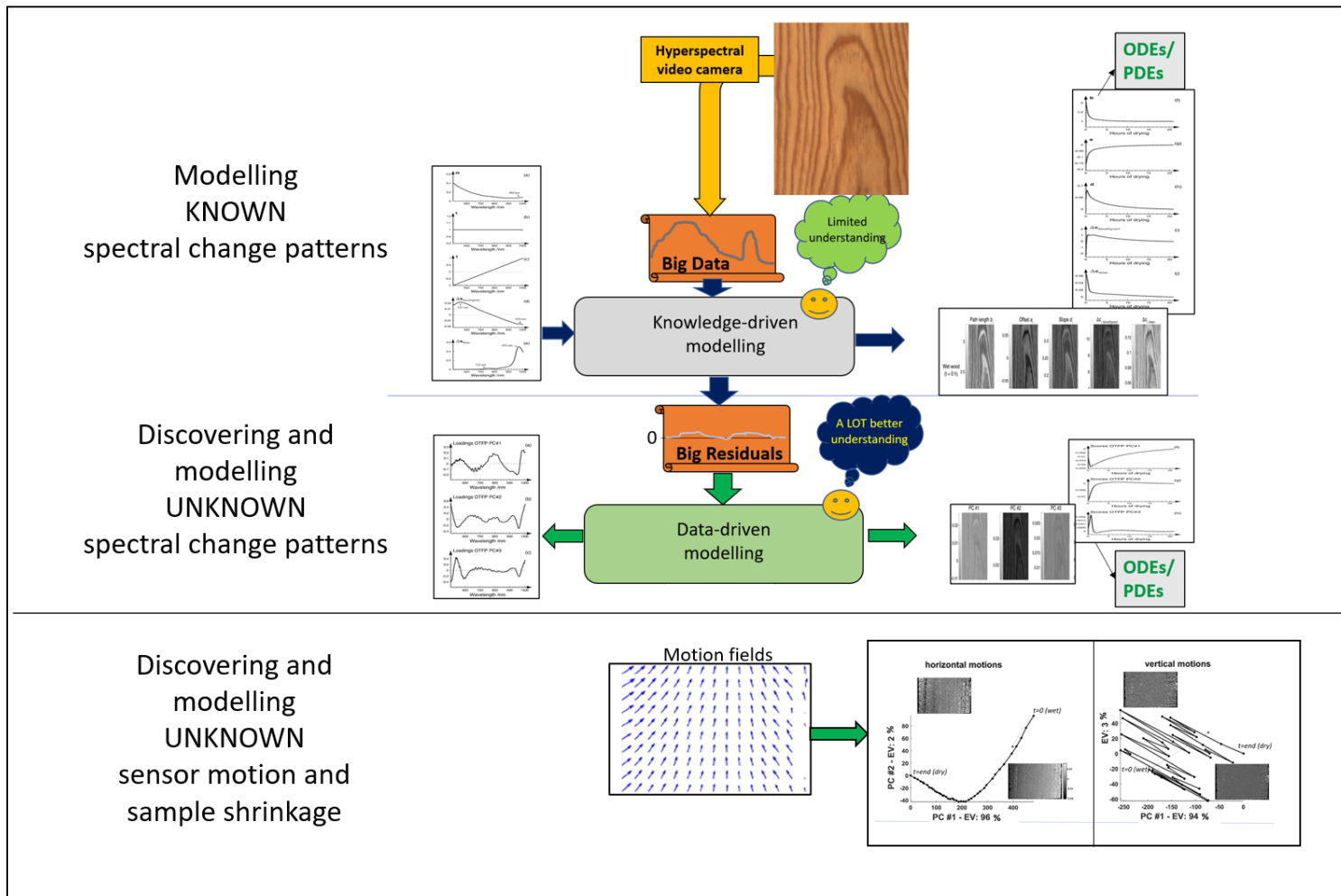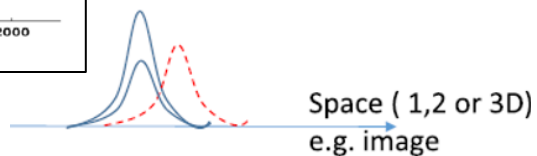
VNIR:
400-1000 nm

**EMSC modelling
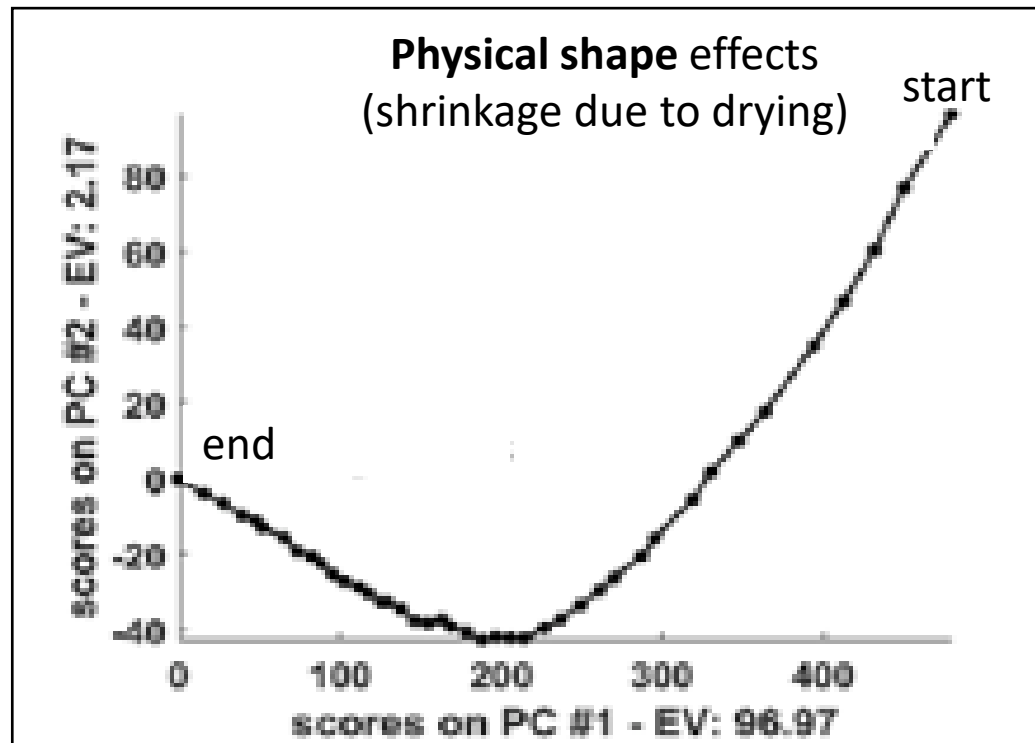KNOWN and
UNKNOWN
physics &
chemistry:**

SWIR:
900-2500 nm
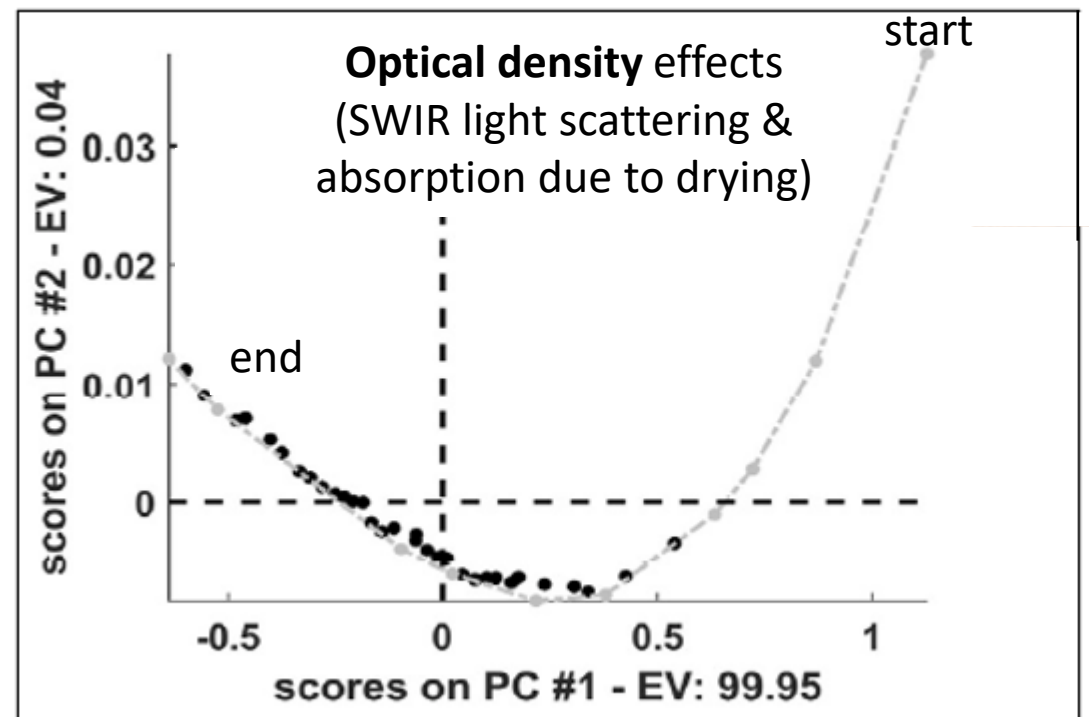
**Two-domain
IDLE modelling:**

# Drying wood in SWIR (900-2500 nm)

## Similar two phase-kinetics for
## physical shrinkage and chemical composition change



**Physical shape** effects
(shrinkage due to drying)

**Optical density** effects
(SWIR light scattering & absorption due to drying)

≈

House of Math

House of everything else

Math gap

House of Math

House of everything else

Math gap

# Thank you!

harald.martens@ntnu.no

# PLSR etc uses a linear method,
# but can often handle non-linear responses automatically

Many data points in a high-dimensional space
e.g. 100 wavelengths ( y, $x_1$, $x_2$,... , $x_{100}$),
happen to form a banana-shaped cloud :
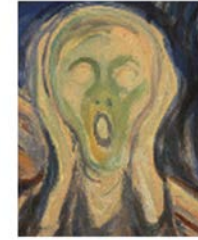
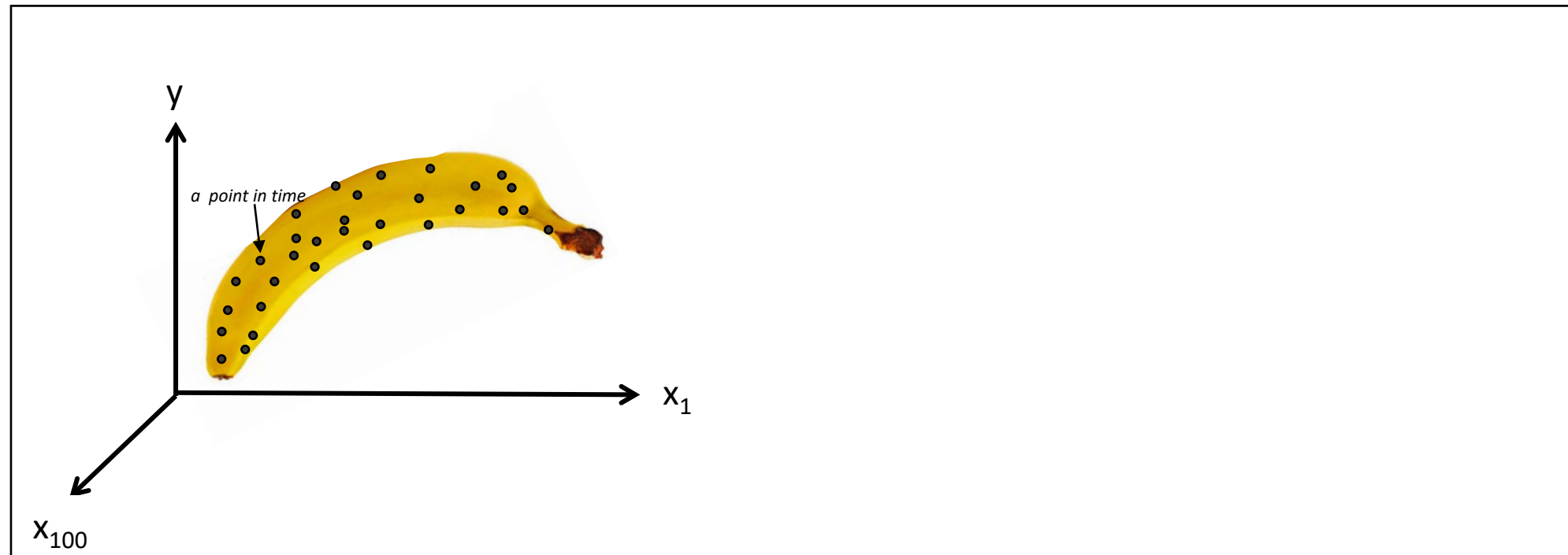MATHEMATICAL
MODELLING ?



y

*a point in time*

$x_1$

$x_{100}$

# PLSR etc uses a linear method,
# but can often handle non-linear responses automatically

Many data points in a high-dimensional space
e.g. 100 wavelengths ( y, $x_1$, $x_2$,... , $x_{100}$),
happen  to form a banana-shaped cloud :

MATHEMATICAL
MODELLING ?



y

*a  point in time*

$x_1$

1D   approximation

$x_{100}$

# PLSR etc uses a linear method,
# but can often handle non-linear responses automatically

Many data points in a high-dimensional space
e.g. 100 wavelengths ( $y$, $x_1$, $x_2$,… , $x_{100}$),
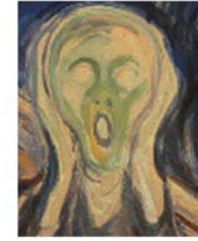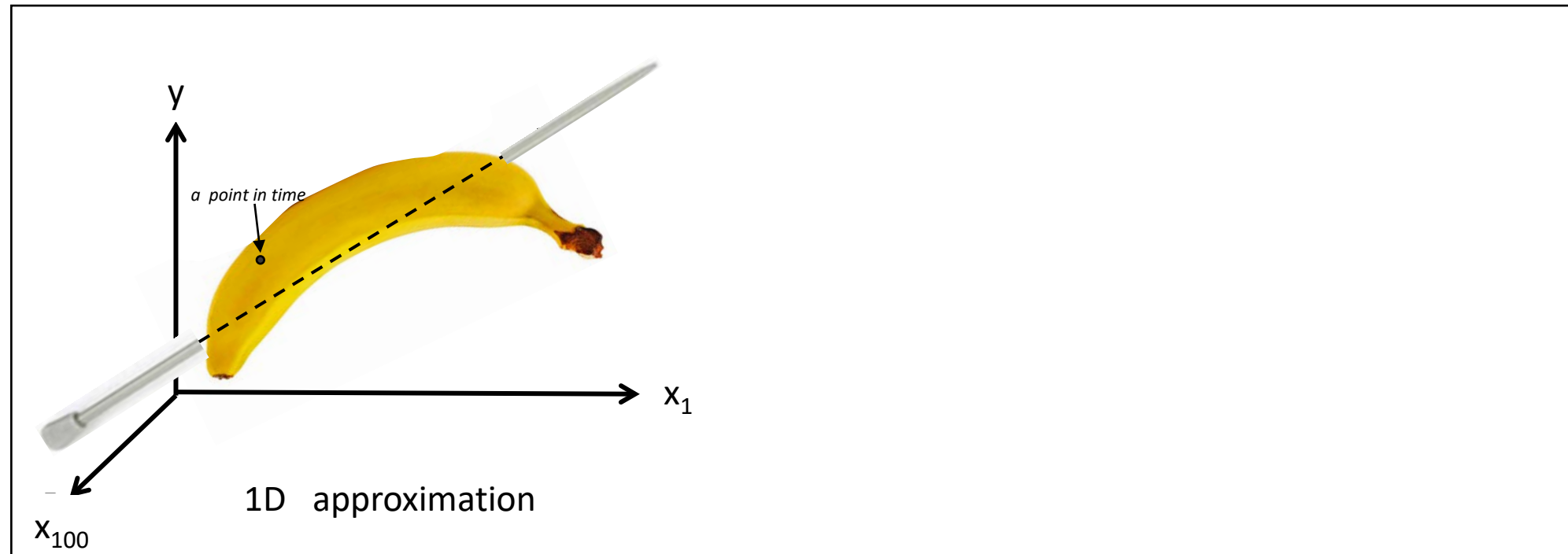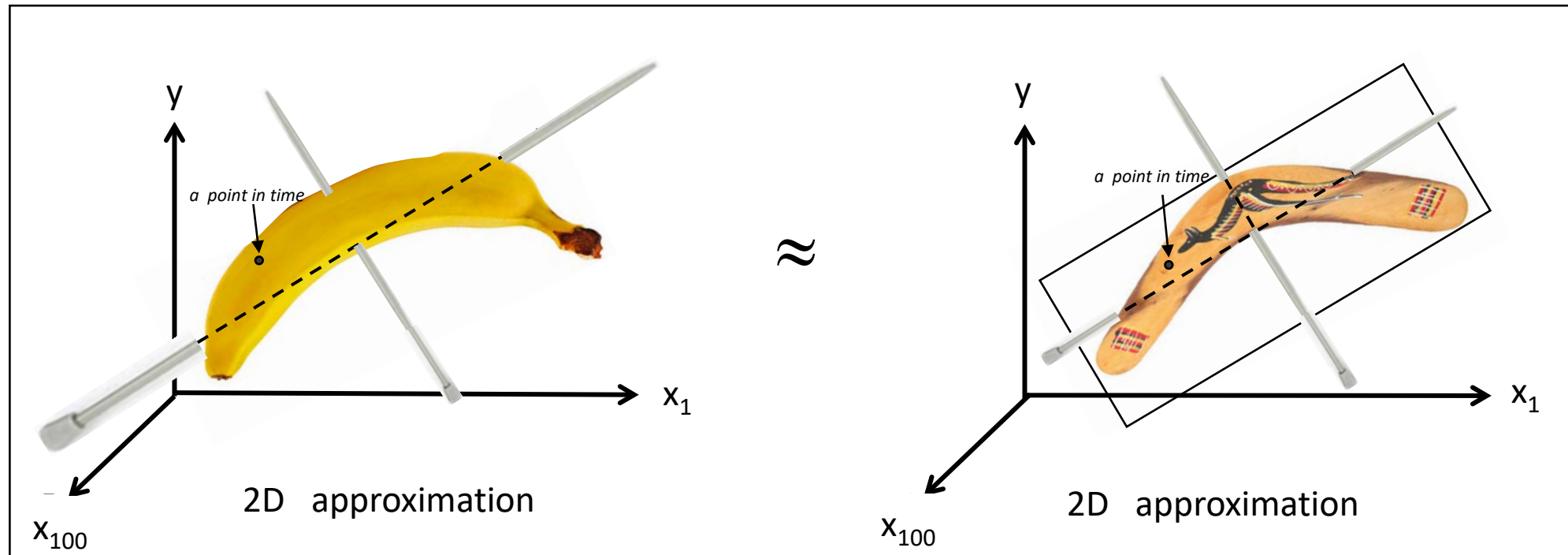happen to form a banana-shaped cloud :

MATHEMATICAL
MODELLING ?



$\approx$

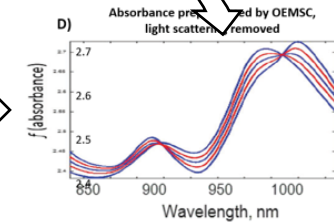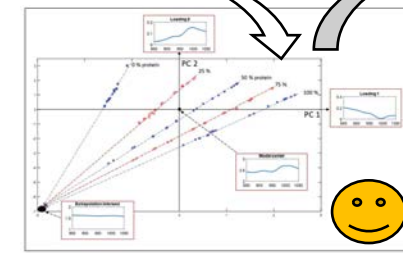2D approximation      2D approximation

# PLSR etc uses a linear method,
# but can often handle non-linear responses automatically

Many data points in a high-dimensional space
e.g. 100 wavelengths ( y, $x_1$, $x_2$,... , $x_{100}$),
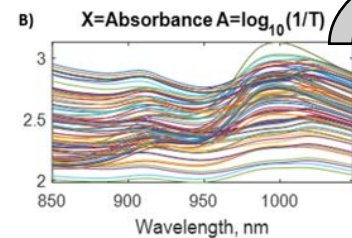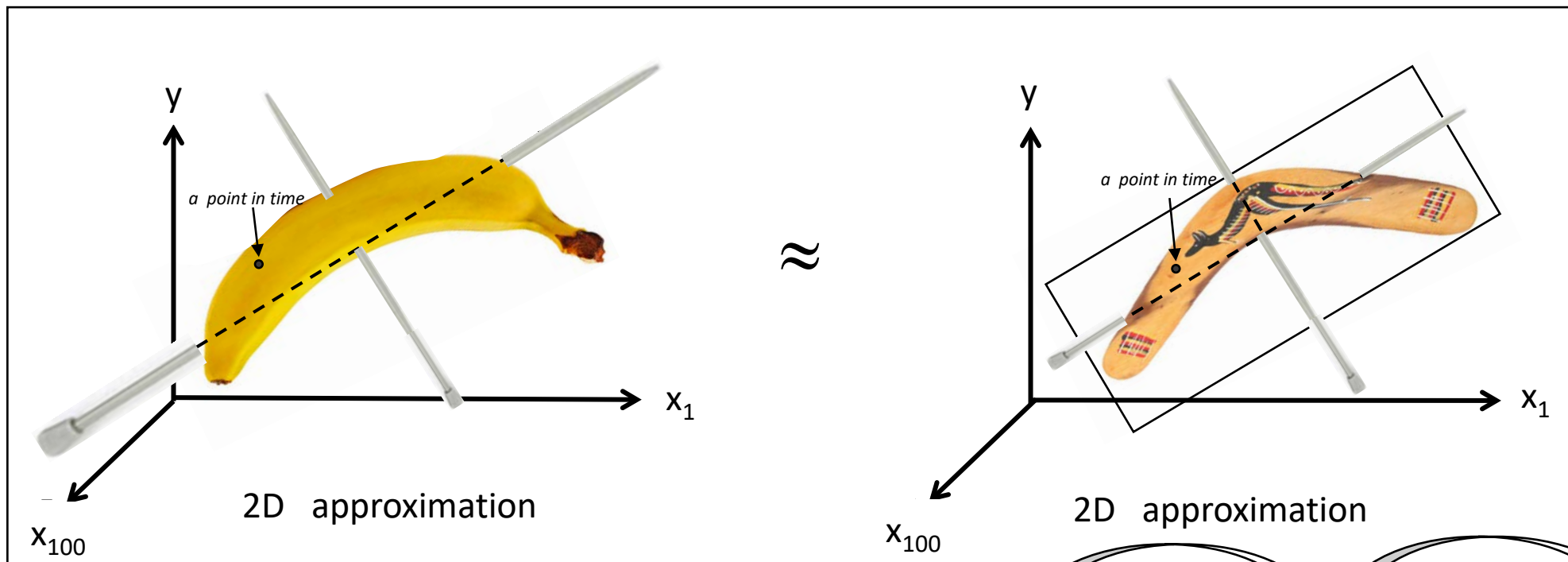happen  to form a banana-shaped cloud :

MATHEMATICAL
MODELLING ?



$\approx$

2D   approximation

2D   approximation

B)   X=Absorbance A=$\log_{10}$(1/T)

D)   Absorbance predicted by OEMSC, light scattering removed

# Linearizing the inputs:



**Where you see a PEAR,**  **I see a ...**

# Linearizing the inputs:

**Ways to combine Deep Learning (DL) and Hybrid Chemometrics (HC)**

1. DL as «hunting dog»:



*If DL did not find anyting, then find a better project !*

## Ways to combine Deep Learning (DL) and Hybrid Chemometrics (HC)

1. DL as «hunting dog»:



*If DL did not find anyting, then find a better project !*

*Find patterns in the hidden nodes of DL*

2. DL «hidden node» scores as data input:
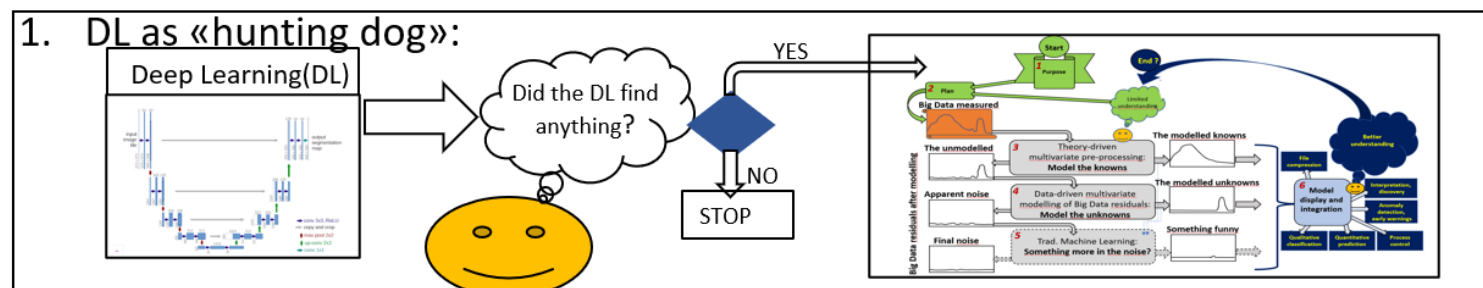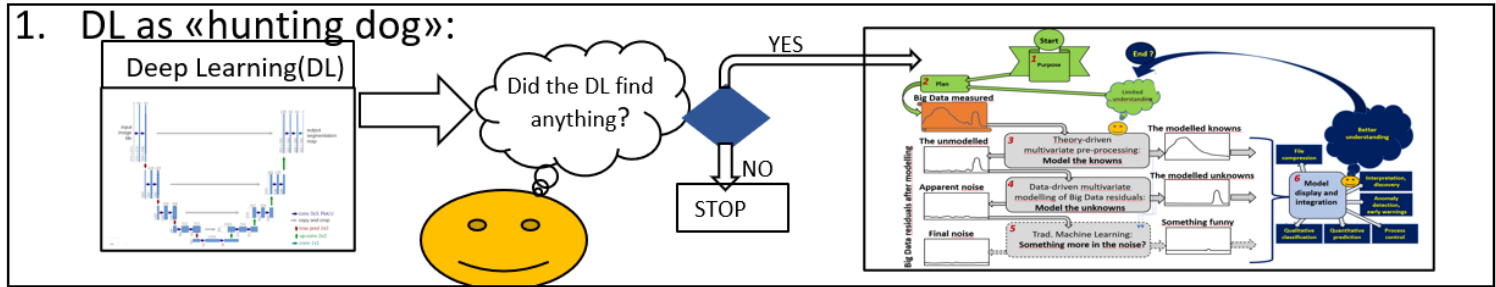
## Ways to combine Deep Learning (DL) and Hybrid Chemometrics (HC)

1.  DL as «hunting dog»:



*If DL did not find anyting, then find a better project !*

*Find patterns in the hidden nodes of DL*

2. DL «hidden node» scores as data input:



*Look for «funny» patterns in PCA residuals*

3. DL for clean-up:
Look for structures in residuals after PCA

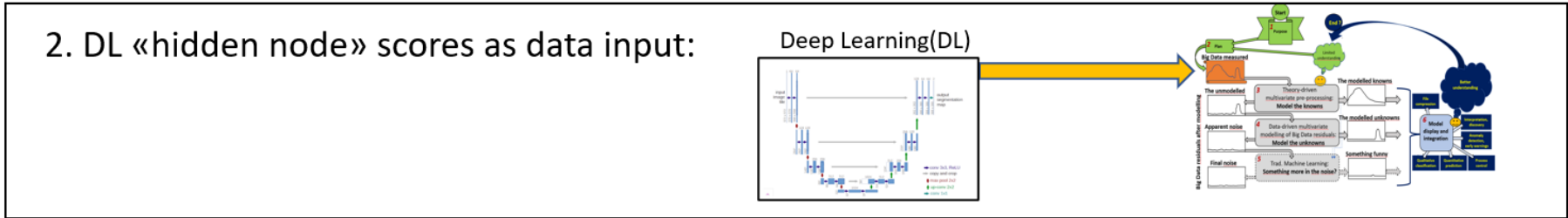## Ways to combine Deep Learning (DL) and Hybrid Chemometrics (HC)

1. DL as «hunting dog»:

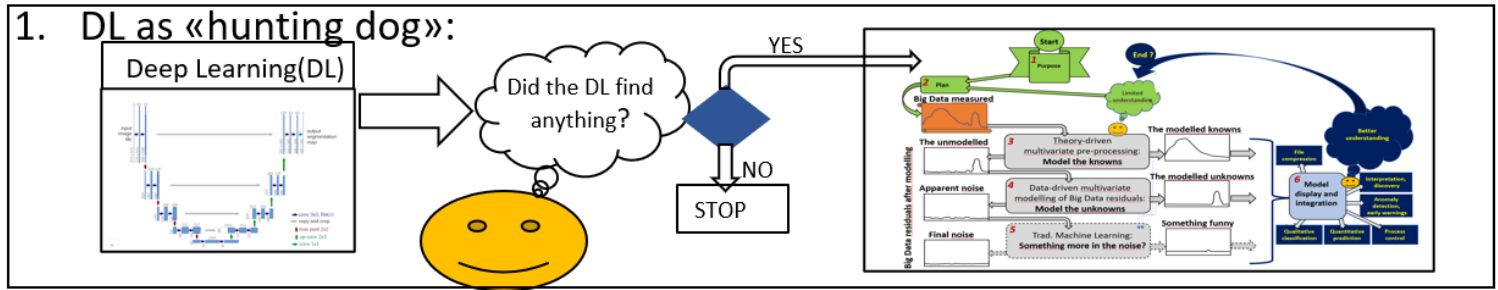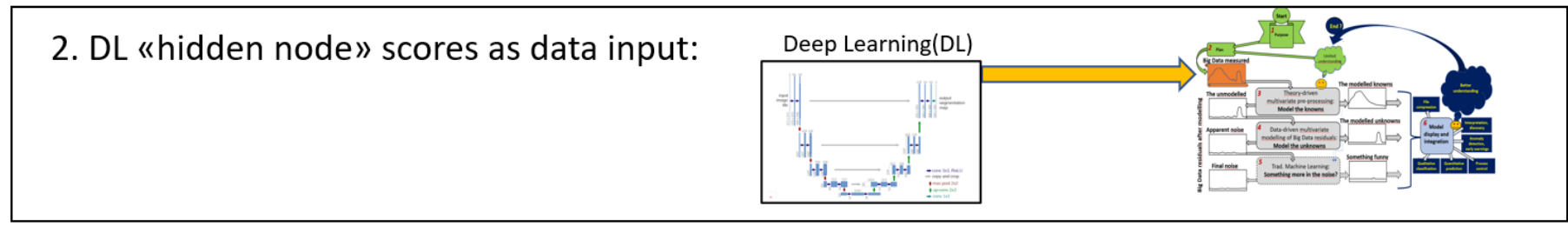

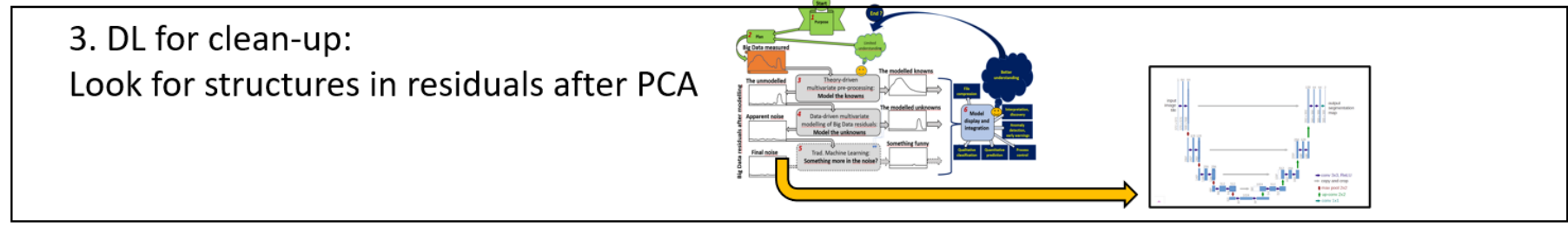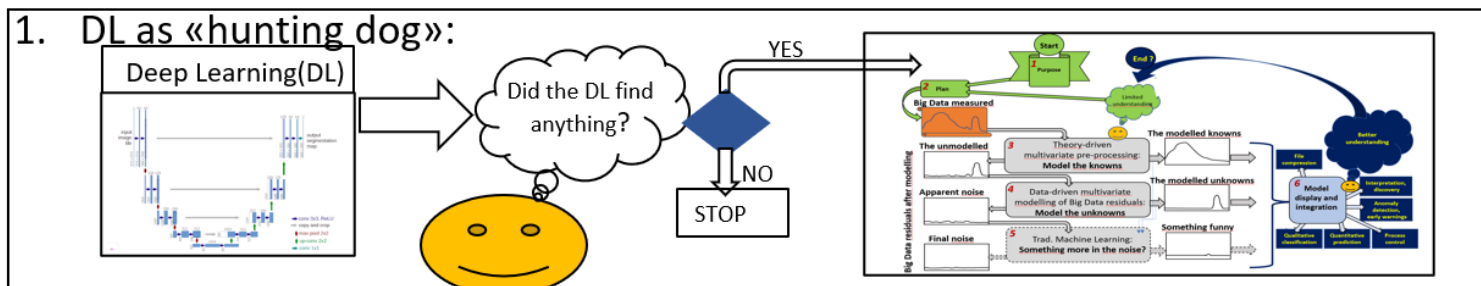*If DL did not find anyting, then find a better project !*

*Find patterns in the hidden nodes of DL*

2. DL «hidden node» scores as data input:



*Look for «funny» patterns in PCA residuals*

3. DL for clean-up:
Look for structures in residuals after PCA



*Look for Higher-order complexities in modelled subspace*

4. DL as post-processing:
Look for more complex structures in the combined output scores and residual statistics from EMSC & PCA

# BIG DATA in Science and Technology (S&T)



**Modern spectral instruments generate Big Data:**

High-speed scanners

Hyperspectral imaging

Hyperspectral video

# *Laws of nature, and other common causes*

Most technical scenes and samples change in systematic ways according to laws of nature, and so do most spectral instruments. A calibration MODEL is needed.

# Spectra from common causes show patterns

**Common causes generate *systematic change patterns* from spectrum to spectrum**



**The *systematic change pattern*s allows us to build a Calibration MODEL**

# *Some causes are expected*

**Many KNOWN causes give NICE,** *systematic change pattern*s



**Many KNOWN systematic change patterns
may be modelled by simple linear models
with additive and multiplicative model elements**

# *Other causes are unexpected*

**But UNKNOWN causes can still give NICE, *systematic change pattern*s**



Start

Purpose

Plan

End ?

**Big Data measured**

Limited understanding

+ *Constituent spectra*
+ *Math models of mechanisms*
+ *Noise levels*

**1** Deductive, theory-driven mechanistic modelling: **Quantify KNOWNS**

**2** Inductive, data-driven subspace modelling: **Discover and quantify UNKNOWNS**
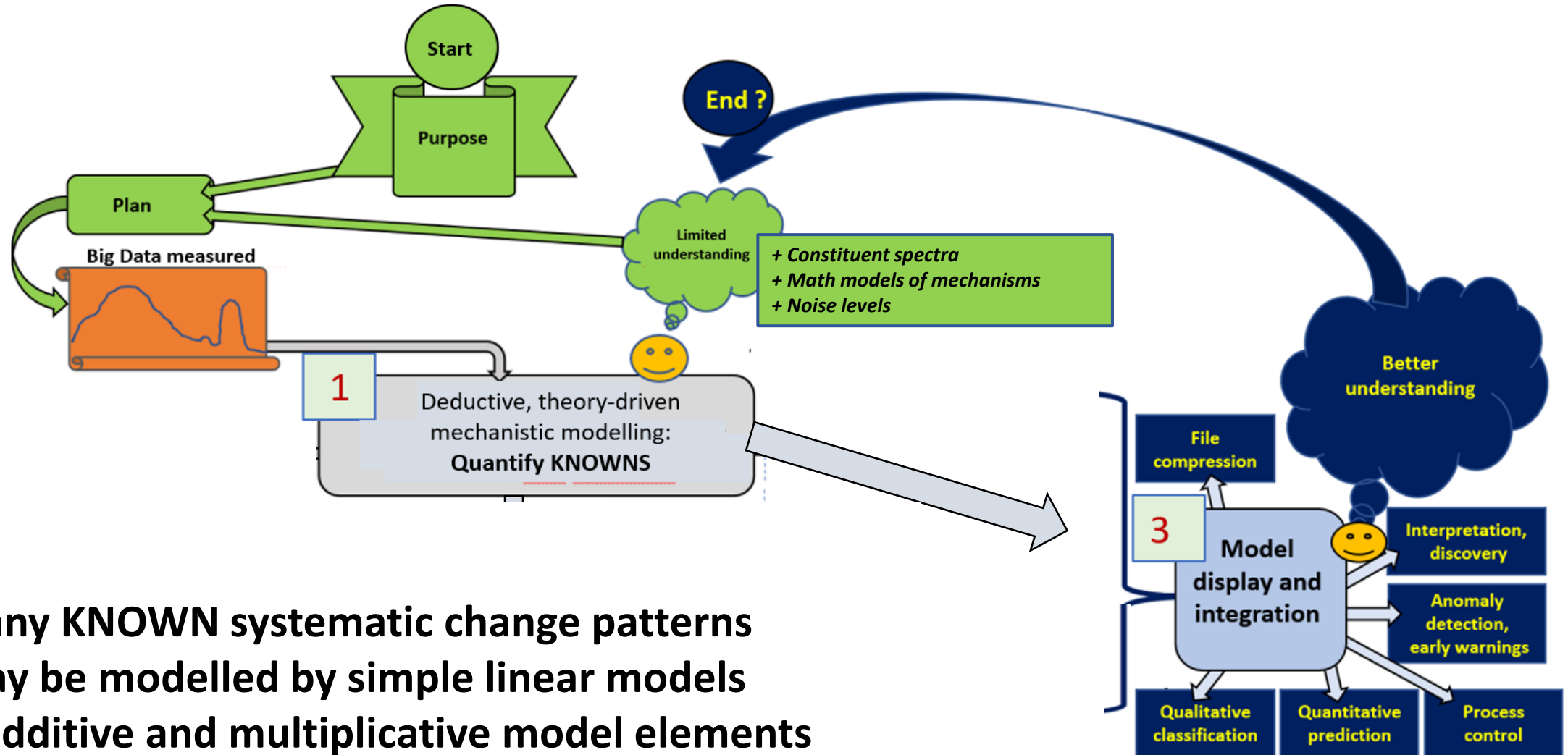
**Better understanding**

File compression

**3** Model display and integration

Interpretation, discovery

Anomaly detection, early warnings

Qualitative classification

Quantitative prediction

Process control

**Many UNKNOWN, but *systematic* change patterns may be modelled by purely additive elements**

Outliers and irrelevant anomalies

# Very strange behaviours

Open

UNEXPECTED, non-systematic PECULIARITIES may also have to be handled

Start

Purpose

Plan

End ?

Big Data measured

Limited understanding

+ Constituent spectra
+ Math models of mechanisms
+ Noise levels

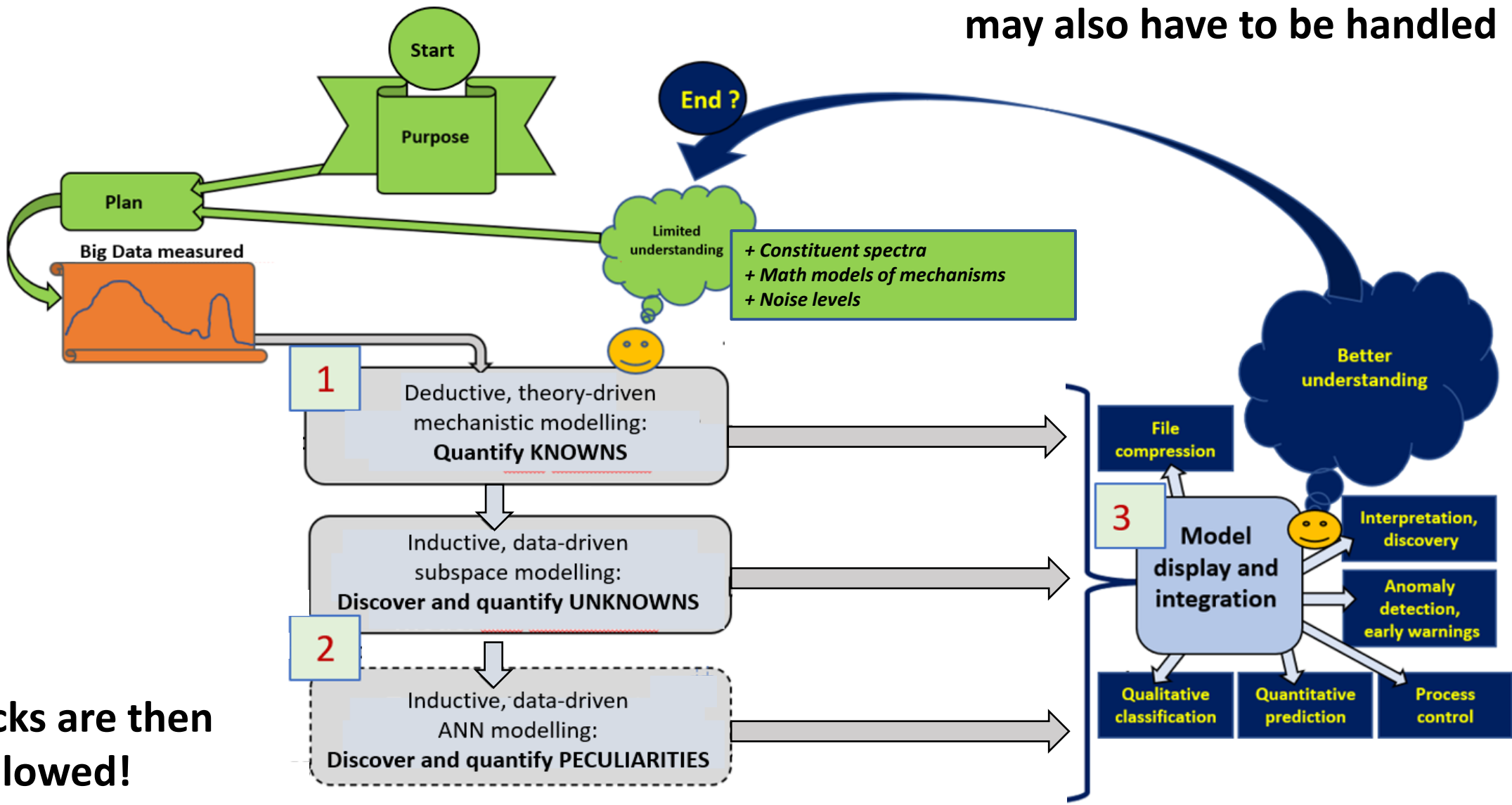**1** Deductive, theory-driven mechanistic modelling: **Quantify KNOWNS**

Inductive, data-driven subspace modelling: **Discover and quantify UNKNOWNS**

**2** Inductive, data-driven ANN modelling: **Discover and quantify PECULIARITIES**

**All tricks are then allowed!**

Better understanding

**3** Model display and integration

File compression

Interpretation, discovery

Anomaly detection, early warnings
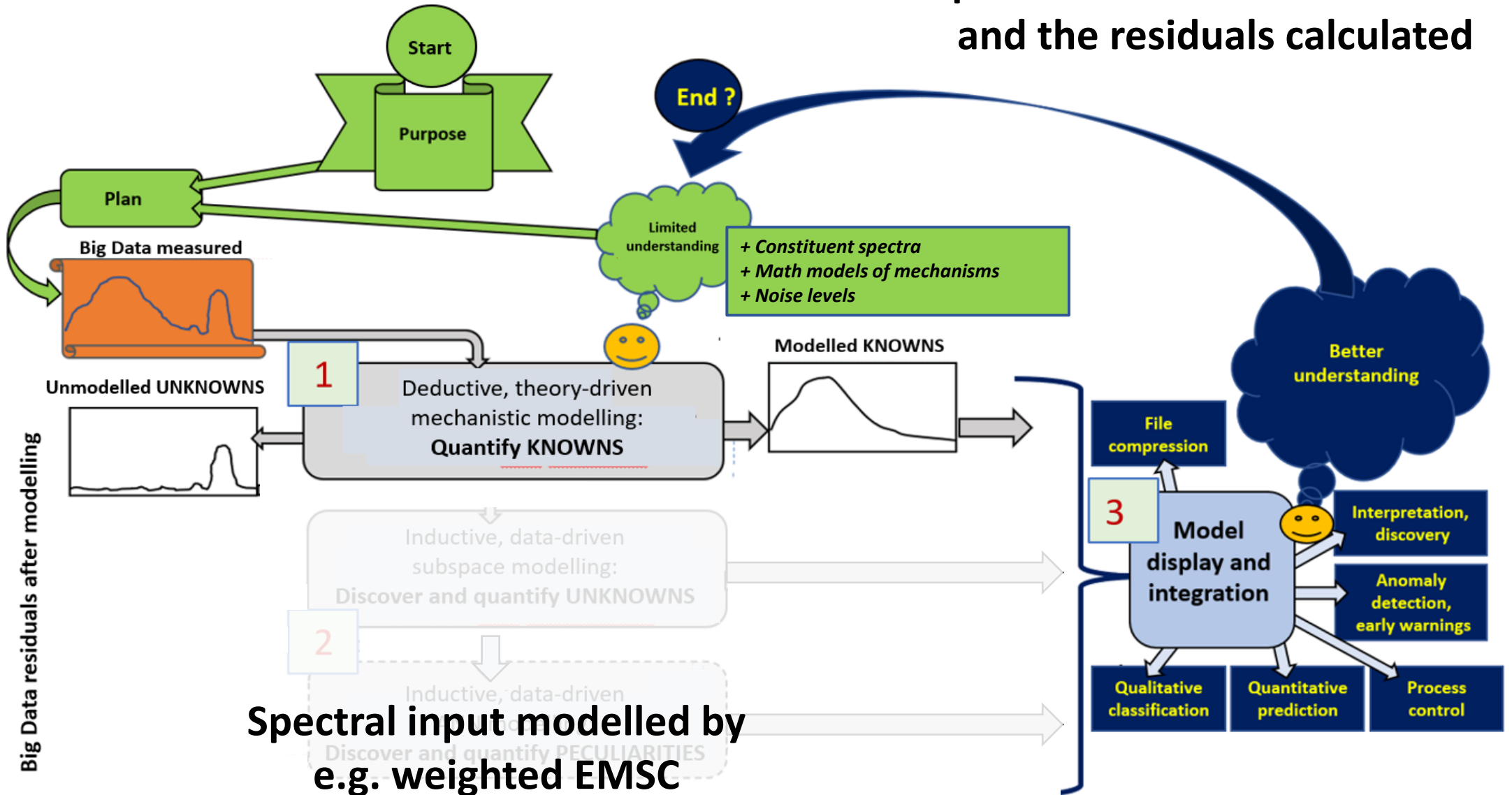
Qualitative classification

Quantitative prediction

Process control

# *Modelling KNOWN patterns*

**KNOWN**
*systematic change pattern*s
**are quantified in multivariate model,
and the residuals calculated**



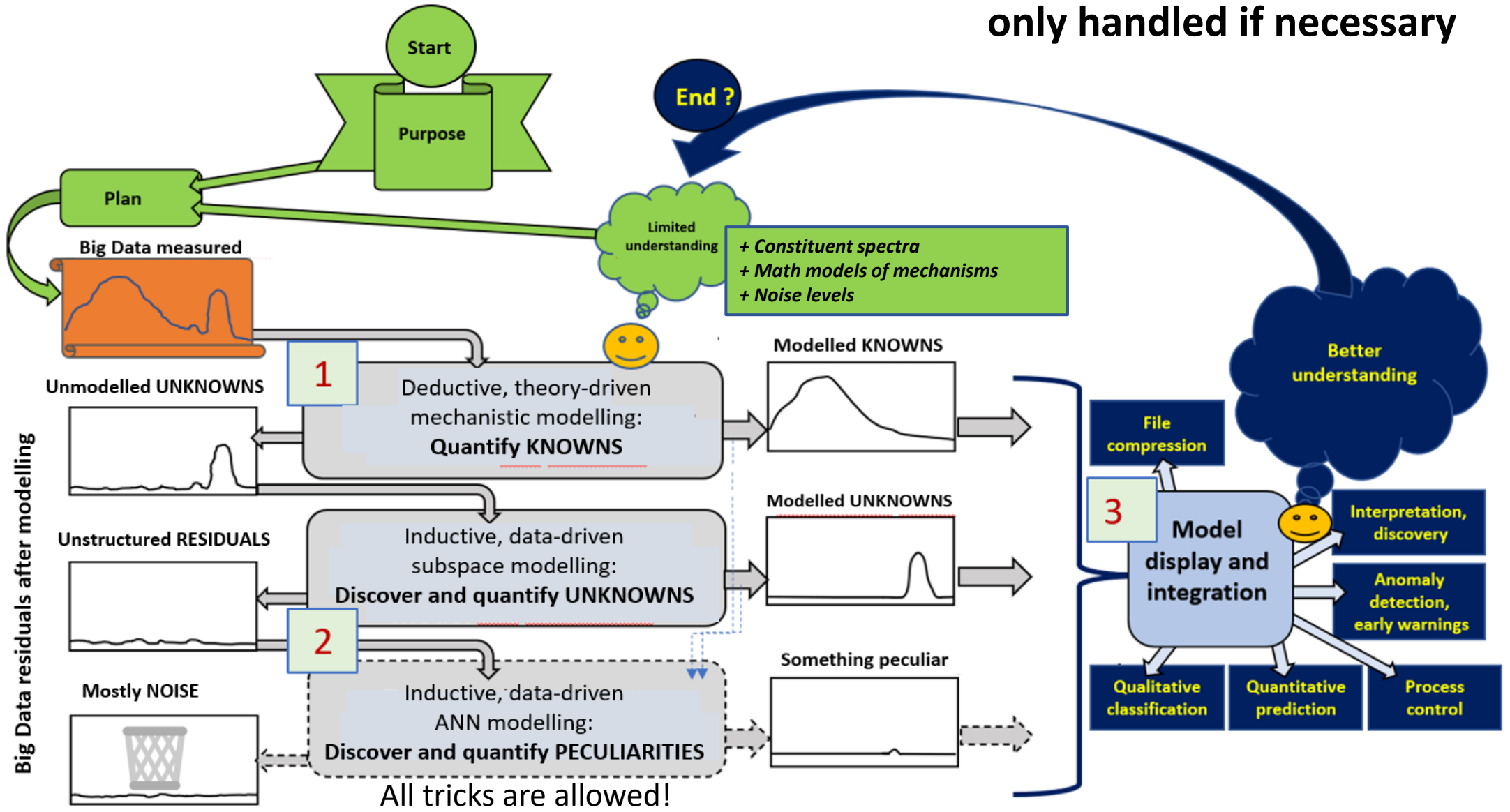**Spectral input modelled by
e.g. weighted EMSC**

# Modelling *UNKNOWN* patterns

**The UNKNOWN, but *systematic* change patterns Are discovered, profiled and quantified by multivariate "machine learning"**
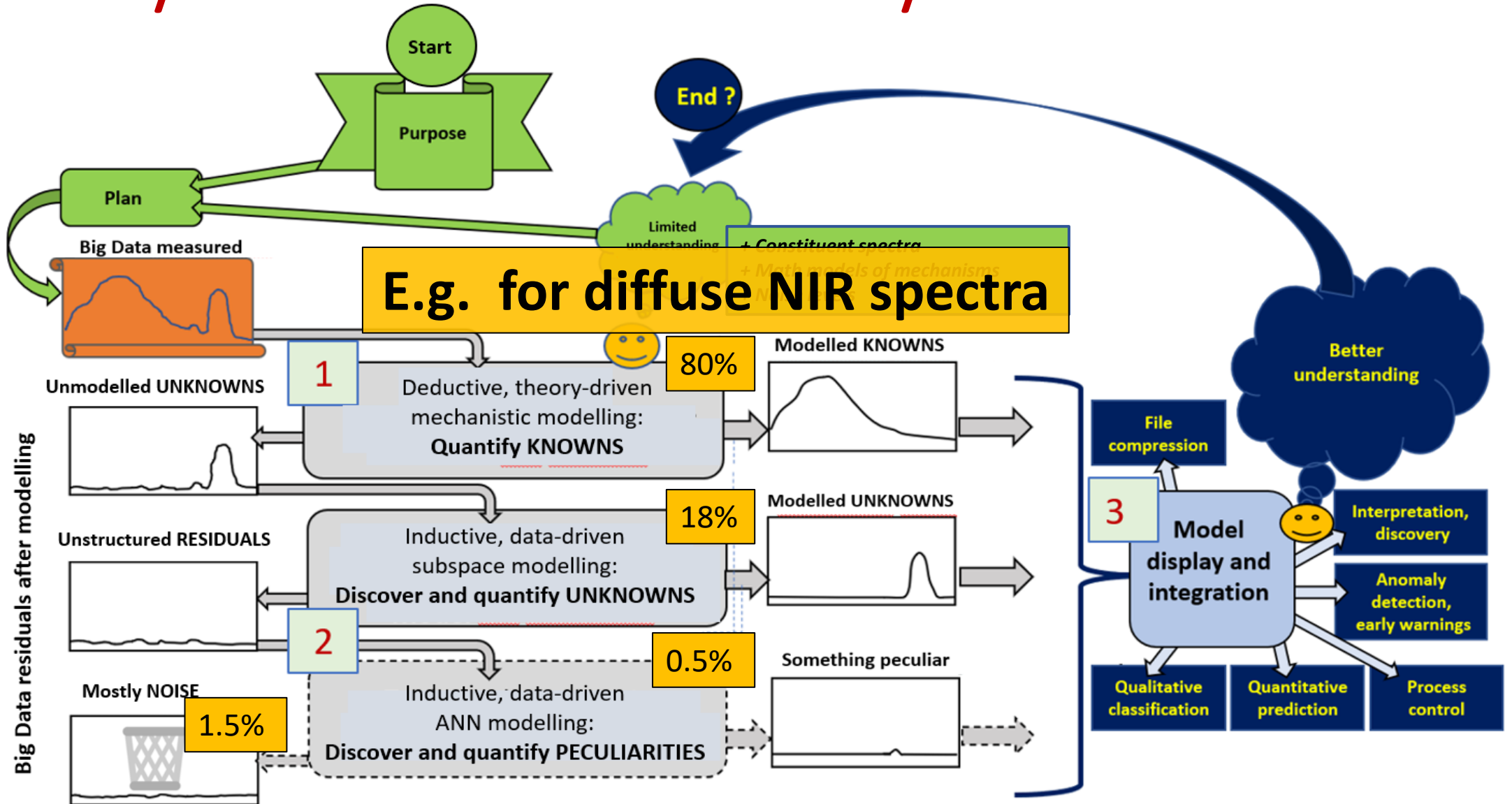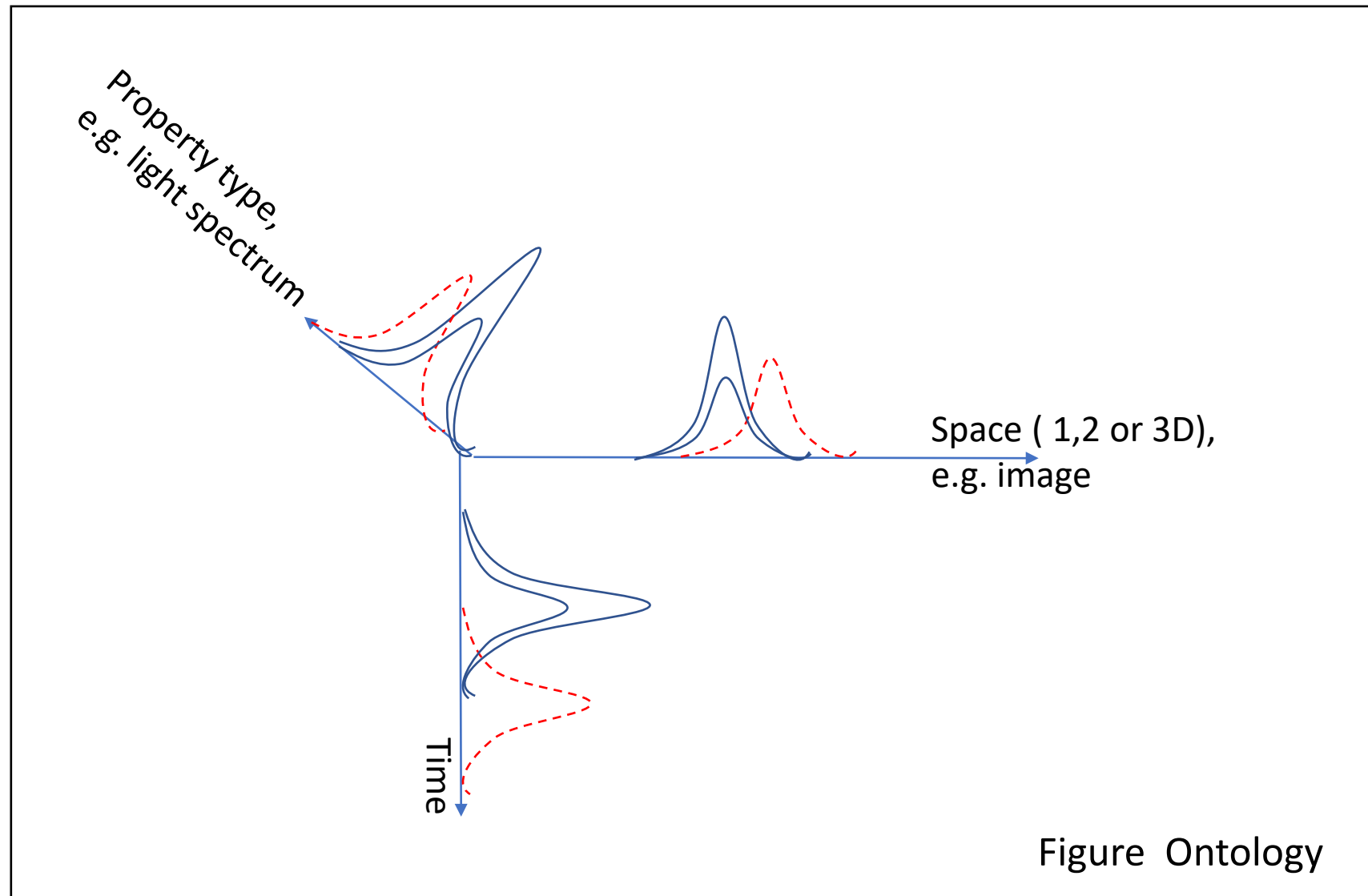
Start

Purpose

Plan

End ?

Limited understanding

+ *Constituent spectra*
+ *Math models of mechanisms*
+ *Noise levels*

**Big Data measured**

**Unmodelled UNKNOWNS**

**Unstructured RESIDUALS**

Big Data residuals after modelling

**1** Deductive, theory-driven mechanistic modelling: **Quantify KNOWNS**

**2** Inductive, data-driven subspace modelling: **Discover and quantify UNKNOWNS**

Inductive, data-driven Discover and quantify PECULIARITIES

**Modelled KNOWNS**

**Modelled UNKNOWNS**

Something peculiar

**Better understanding**

**3** Model display and integration

File compression

Interpretation, discovery

Anomaly detection, early warnings

Qualitative classification

Quantitative prediction

Process control

**Spectral residuals modelled by e.g. weighted PCA**

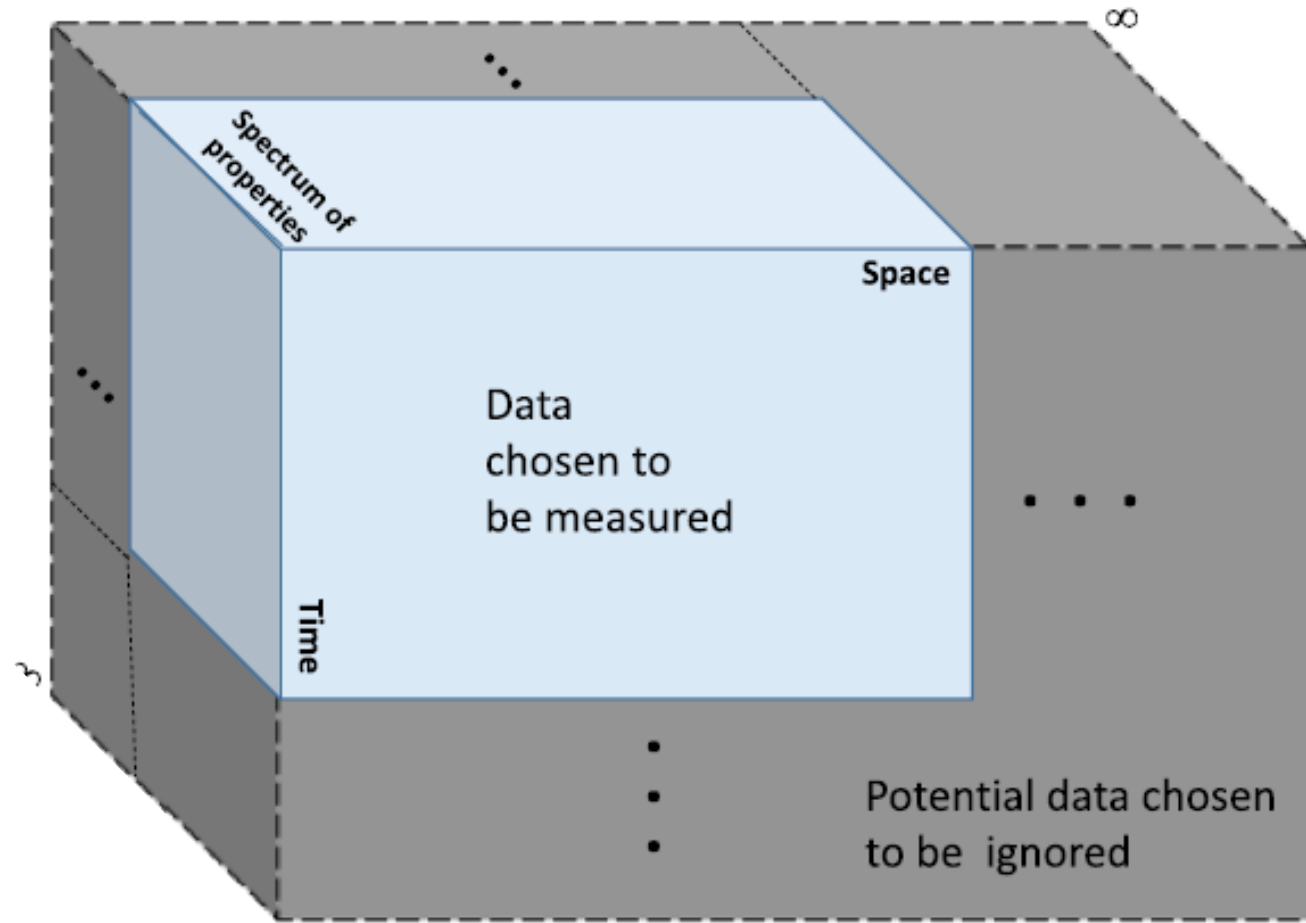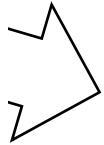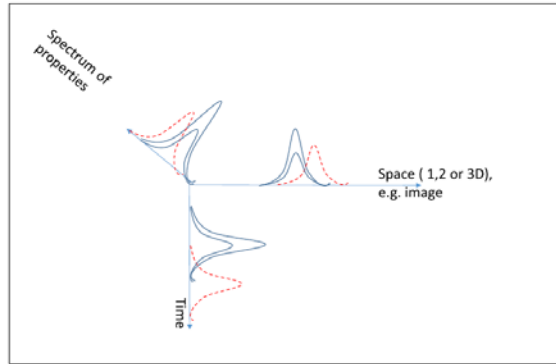*Variances explained in a representative set of samples*

Ontology: position and intensity variation in time, space and properties



Property type,
e.g. light spectrum

Space ( 1,2 or 3D),
e.g. image

Time

Figure  Ontology

# Which DATA are measured?

Epistemology: measure position and intensity variation in time, space and properties, and extract interpretable essence by data modelling
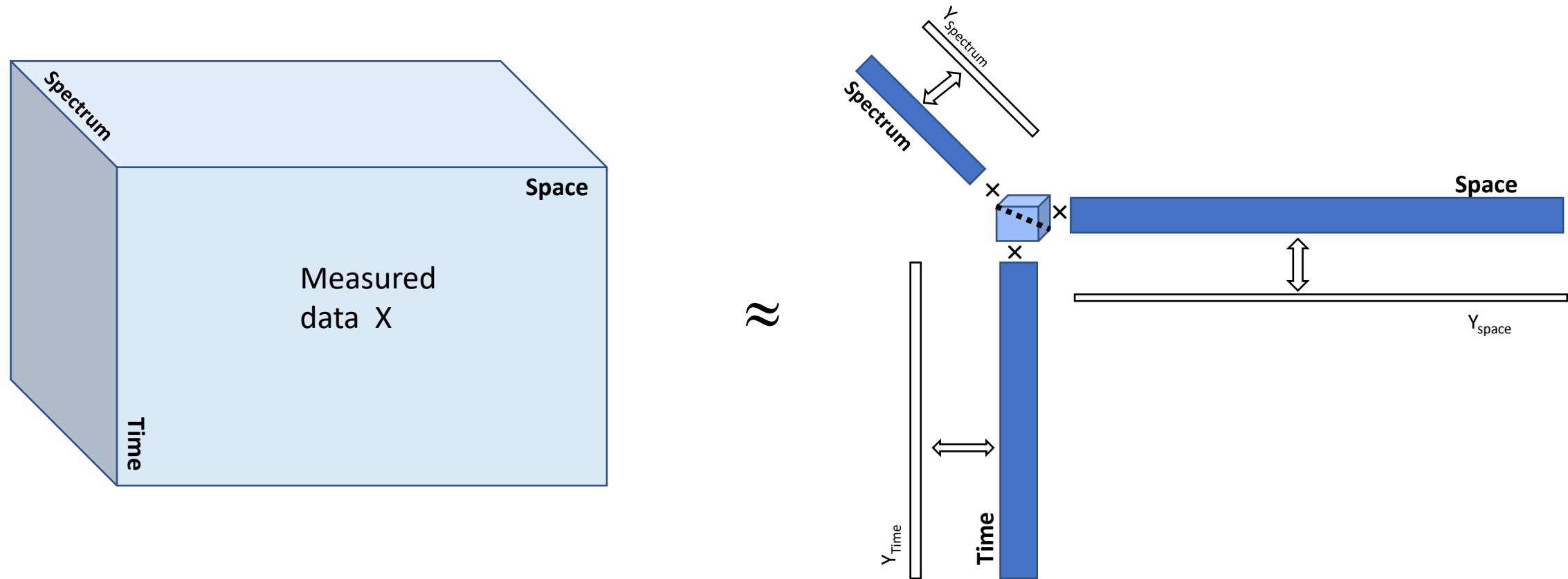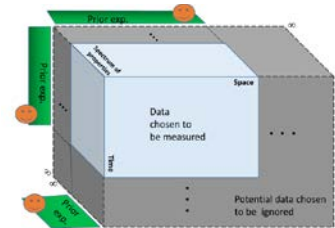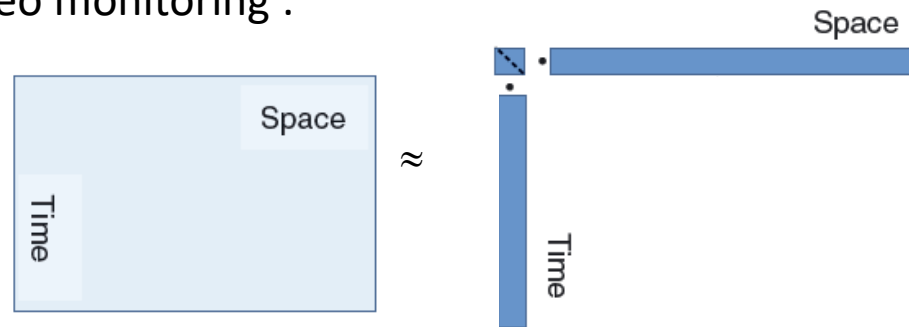


Figure N-linear
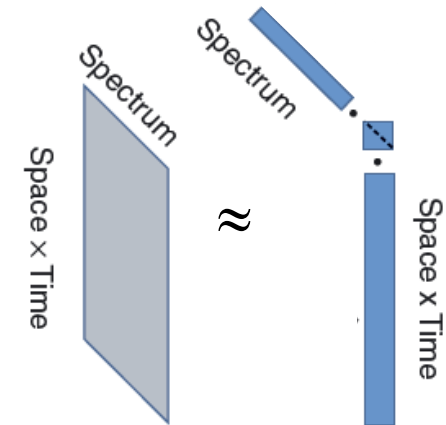
# Subspace autoencoder,  examples:



Pragmatic subspace models & Statistical design, validatio, graphics
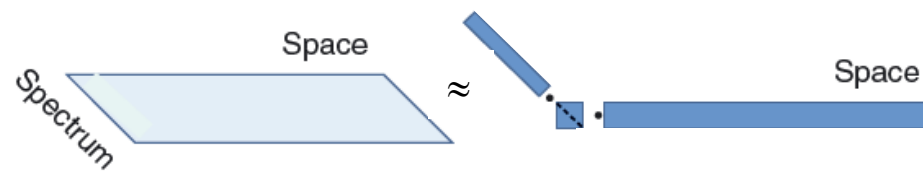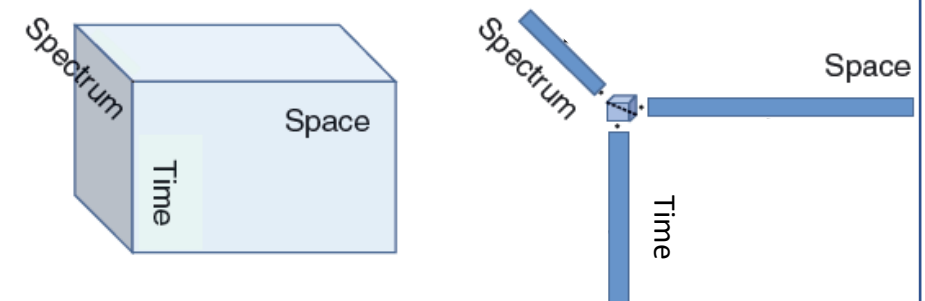
Process video monitoring :

Hyperspectral video monitoring :

Hyperspectral imaging :

Hyperspectral video monitoring :

# Subspace regressions, examples:



Pragmatic subspace models
&
Statistical design, validatio, graphics

Process video monitoring :

Hyperspectral video monitoring :

Hyperspectral imaging :

Hyperspectral video monitoring :