# Interoperable and Scalable Metabolomics Data Analysis with Microservices

Christoph Steinbeck
and the
PhenoMeNal consortium

**European Bioinformatics Institute, Hinxton**

**Friedrich-Schiller-University, Jena**

Genomics Education Programme

NHS Health Education England

## Traditional Medicine

## Personalised Medicine

@genomicsedu    www.genomicseducation.hee.nhs.uk    f /genomicsedu

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# Phenome - Exposome

# Reaction times following external change



Metabolism (Seconds)
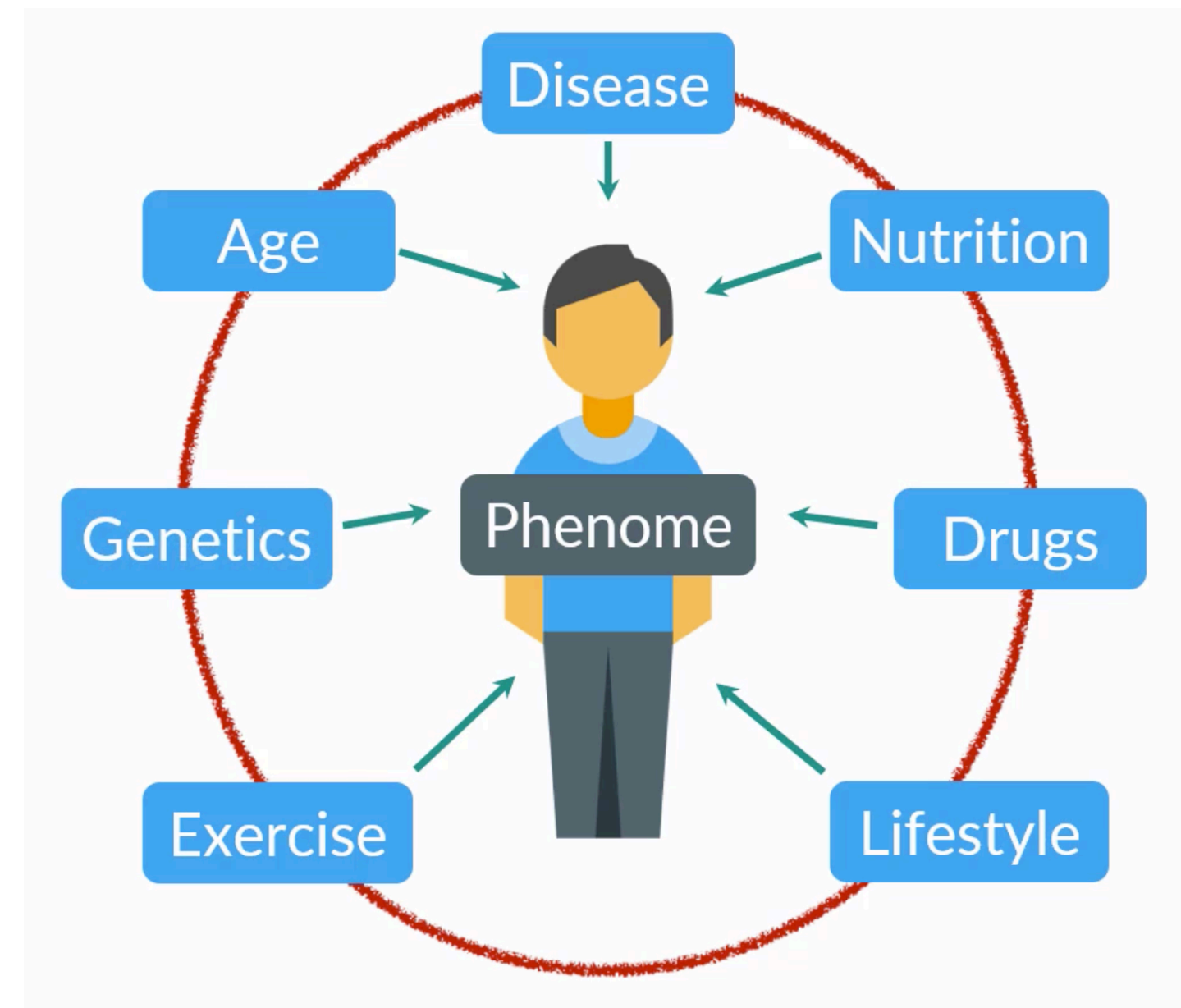
Genetics (Decades, Centuries...)
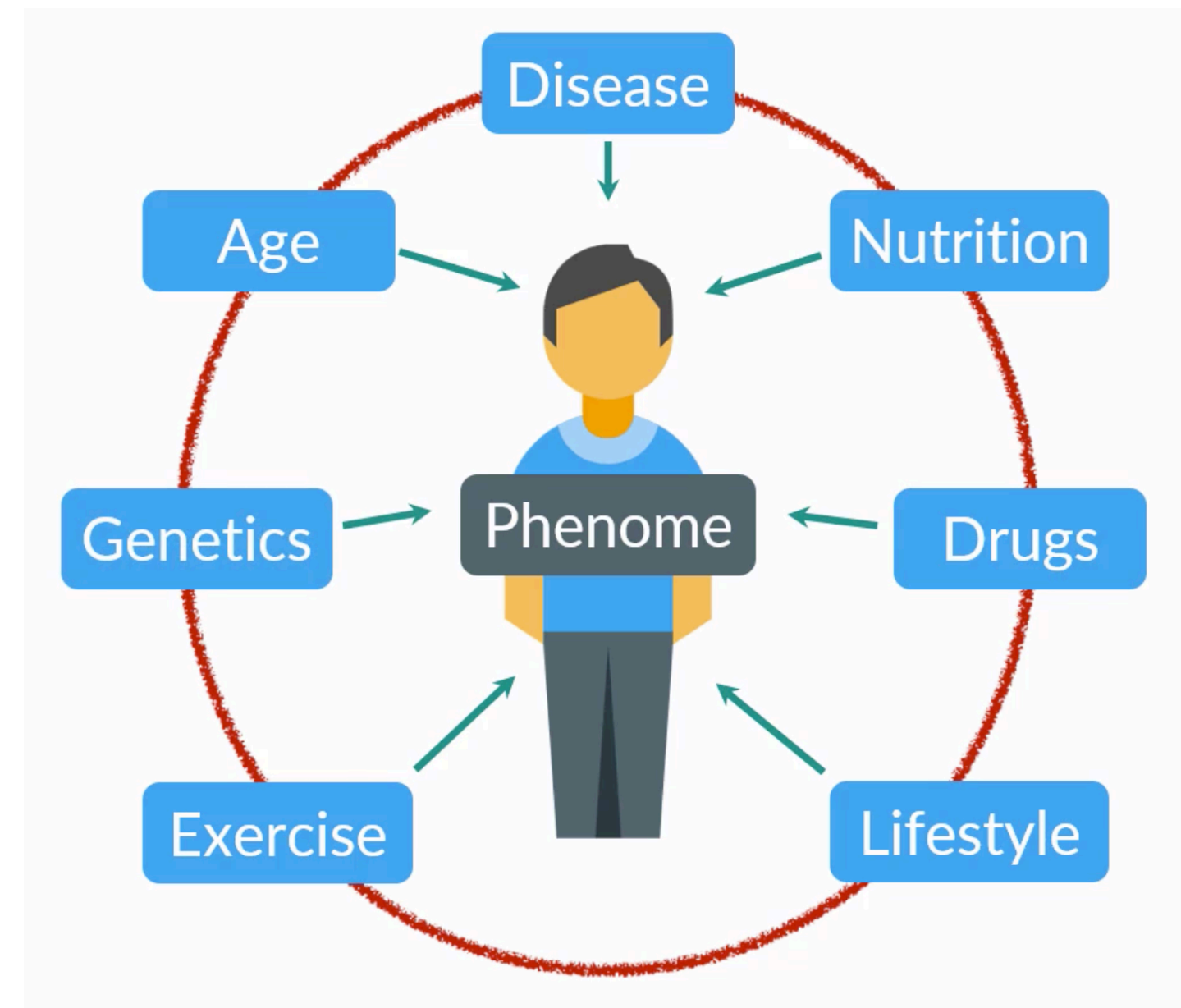
Epigenetics (Days, Months, Years...)

Gene Expression (Hours)

# The Metabolome
# is an easily accessible and dynamically changing Molecular Phenotype
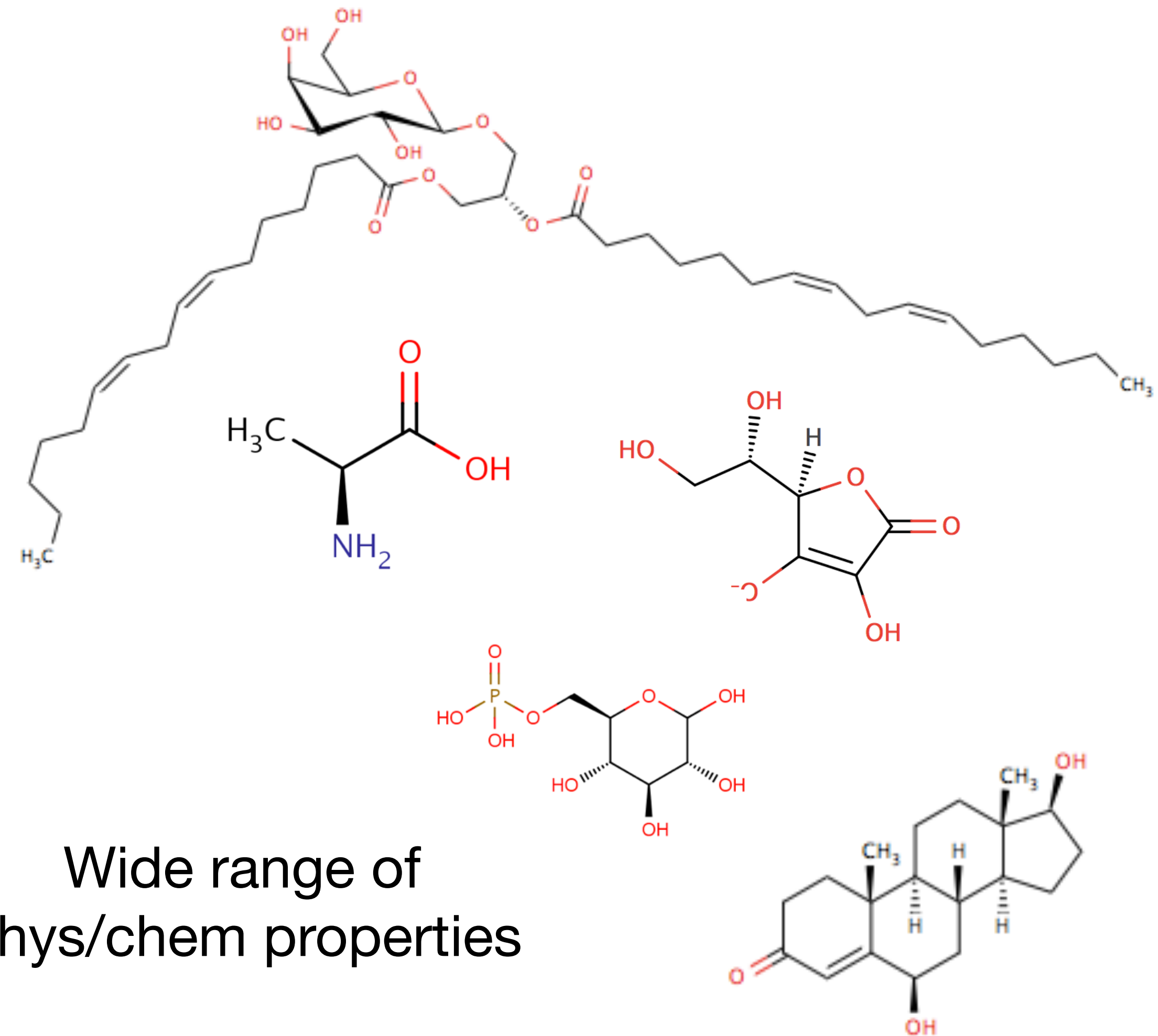
# Phenome - Exposome

# Phenome - Exposome



**Big, well-annotated** metabolomics **data** required to statistically link individual components of the **exposome** to effects in the **molecular phenotype**

# Metabolomics

Measures **occurrence** and **concentrations** of many small molecules (**metabolites**) in an organism at once.

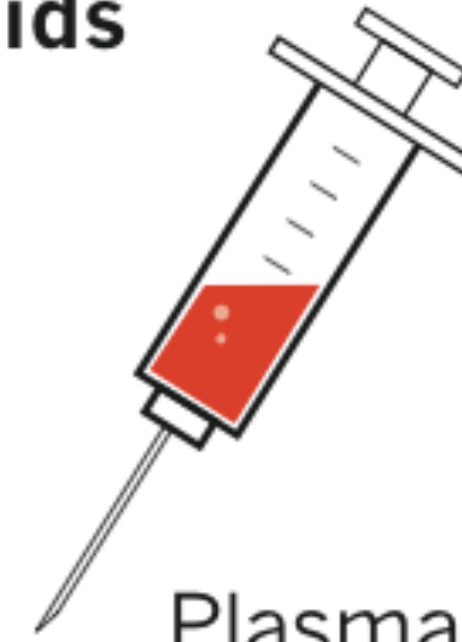**Metabolites**: (Endogenous) small molecules in biological organisms



Wide range of phys/chem properties

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

**Diagnostic fluids**

Urine
(time-averaged data)

Plasma
(snap-shot data)

**Other accessible
analytical compartments**

Pathological fluids

Specialized fluids and biopsies
(selected fluids)

Artificial fluids

Nicholson et al., Nature, 491(7424), 384–392
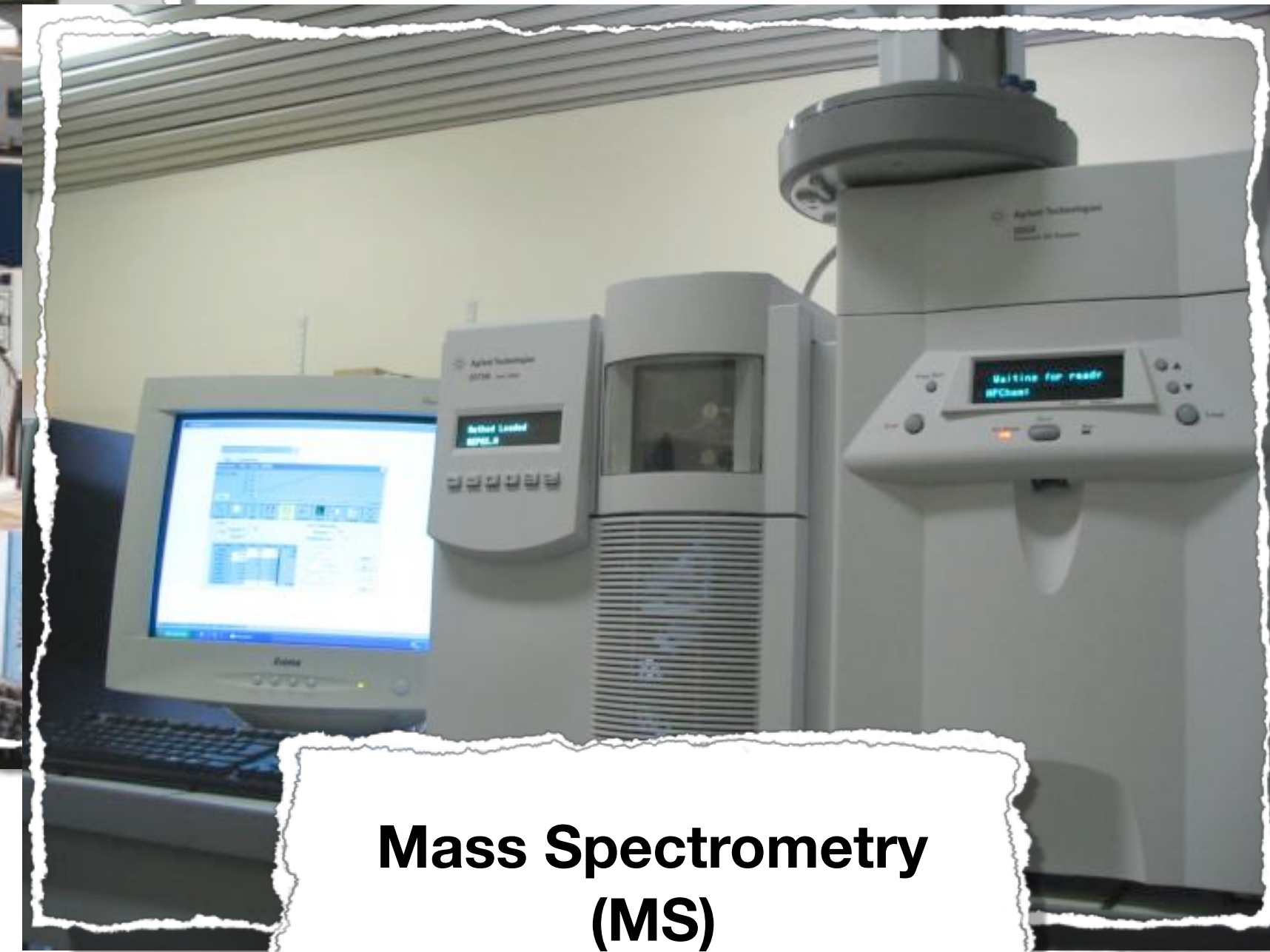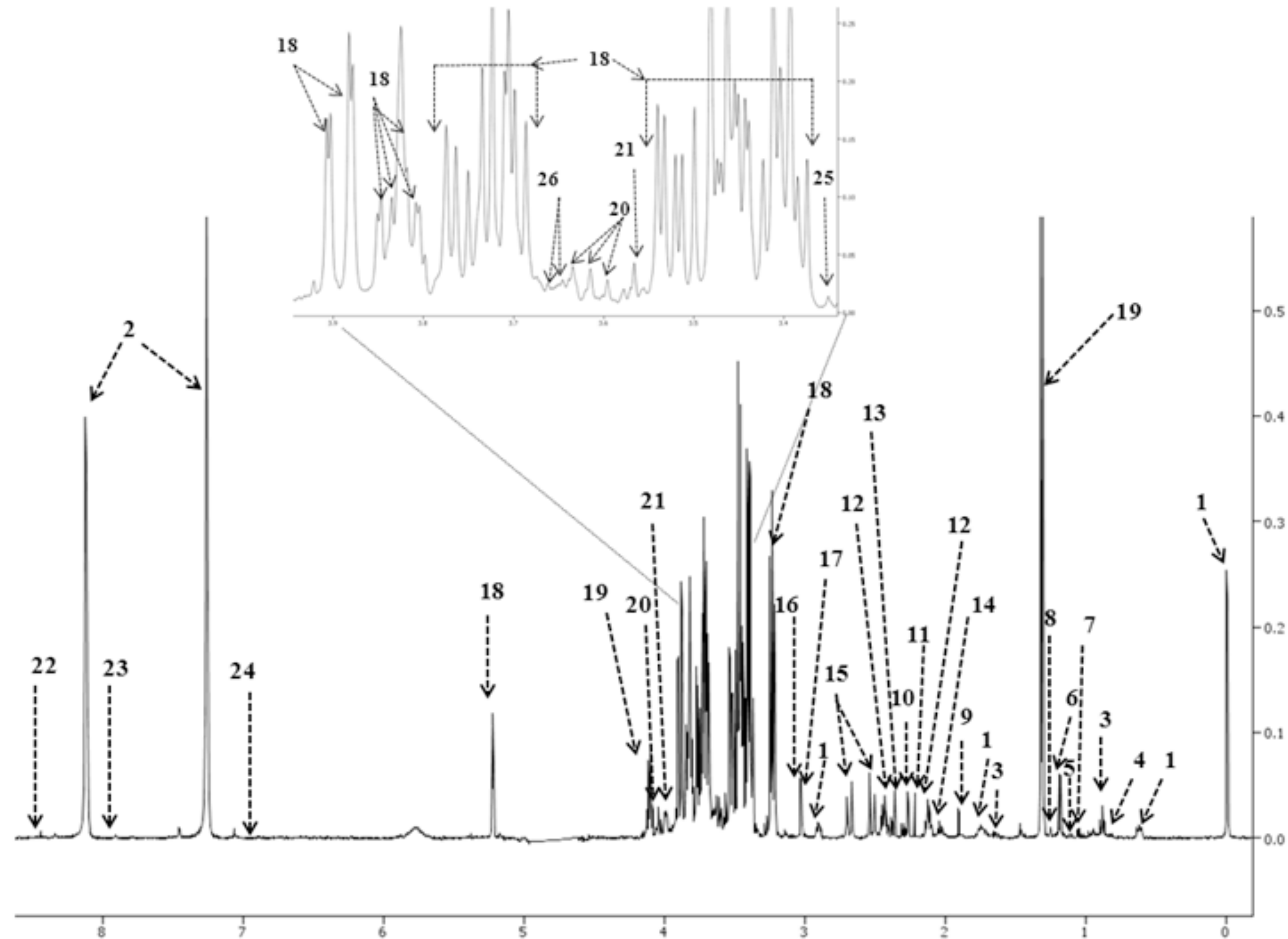
# Metabolomics uses a wide-range of analytical techniques



Nuclear Magnetic Resonance (NMR)

Mass Spectrometry (MS)

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA
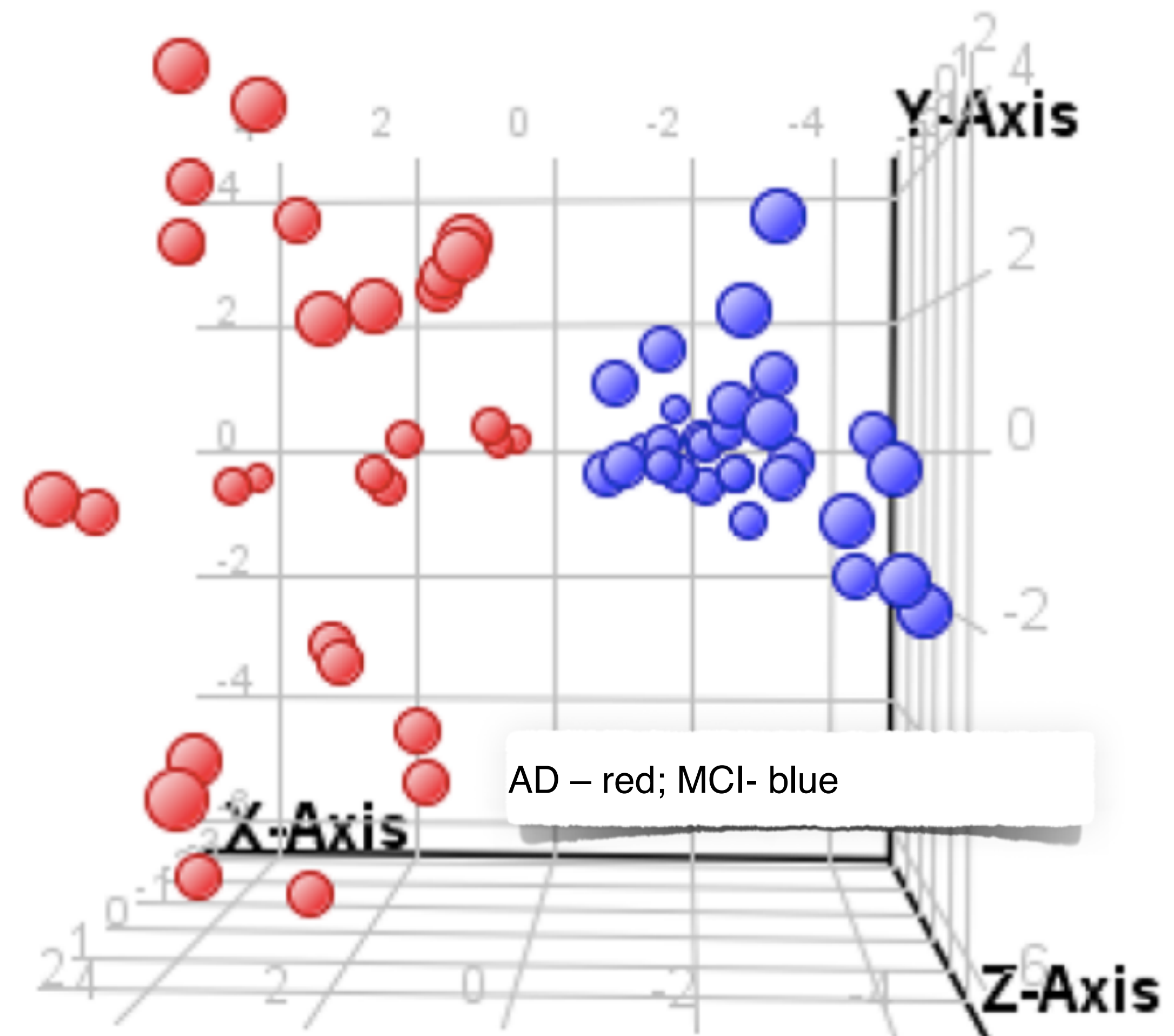
Typical 500 MHz ¹H-NMR spectrum of human cerebrospinal fluid. Numbers indicate the following metabolites:

1. DSS, 2. imidazole, 3. 2-hydroxybutyric acid, 4. 2-hydroxyisovaleric acid, 5. 2-oxoisovaleric acid, 6. 3-hydroxybutyric acid, 7. 3-hydroxyisobutyric acid, 8. 3-hydroxyisovaleric acid, 9. acetic acid, 10. acetoacetic acid, 11. acetone, 12. L-glutamine, 13. pyruvic acid, 14. L-glutamic acid, 15. citric acid, 16. creatinine, 17. creatine, 18. D-glucose, 19. lactic acid, 20. myo-inositol, 21. D-fructose, 22. formic acid, 23. L-histidine, 24. L-tyrosine, 25. methanol, 26. glycerol

source: http://www.csfmetabolome.ca

# Example: Untargeted Metabolomics of Cerebrospinal Fluid



AD – red; CN- blue

AD – red; MCI- blue

Alzheimer's Disease (AD), Mild Cognitive Impairment (MCI), Cognitive Normal (CN)

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

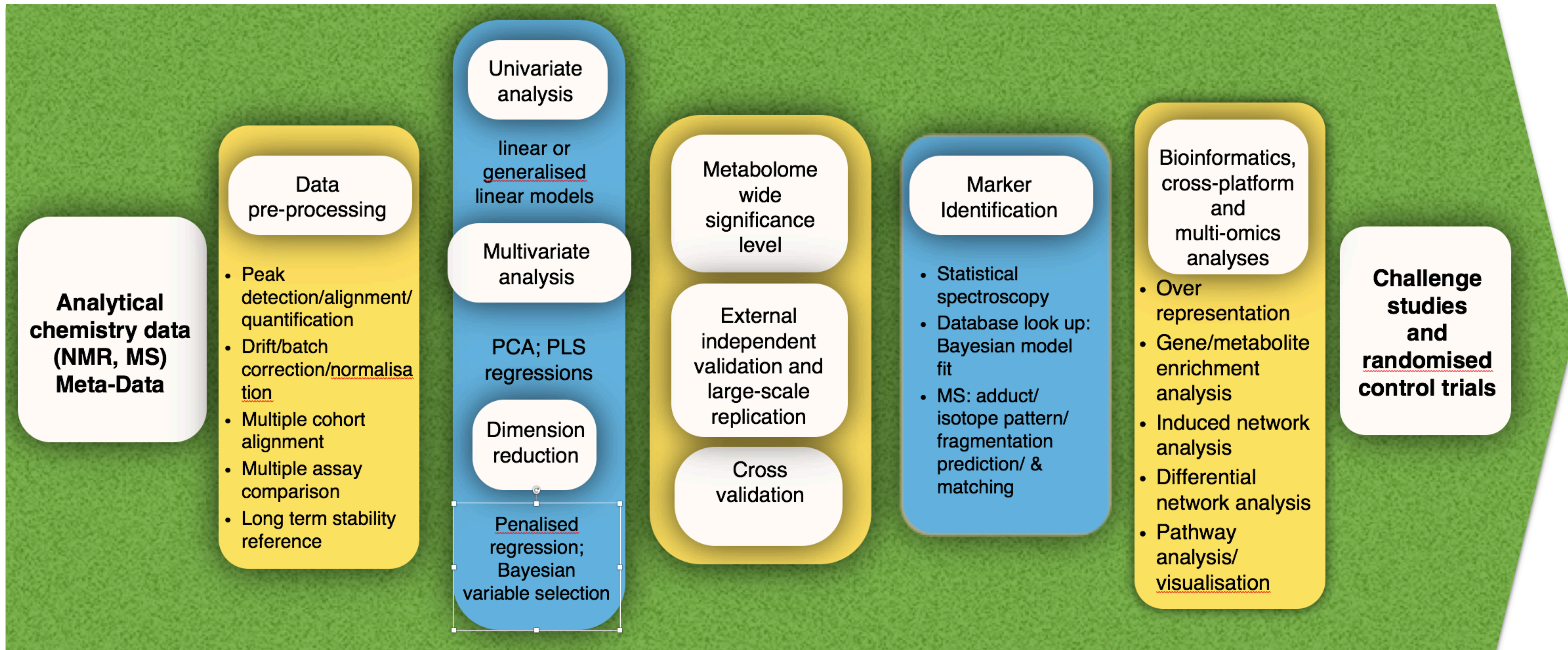# Large Scale Computing with Medical Metabolomics Data



- H2020 e-infra
- 3 Years
- 13 Partners
- 8 Mio €
- 830 PM
- Kick-off 9/15

**PhenoMeNal**

An e-infrastructure for

# Large Scale Computing with Medical Metabolomics Data

**Analytical chemistry data (NMR, MS) Meta-Data**

**Data pre-processing**

- Peak detection/alignment/quantification
- Drift/batch correction/normalisation
- Multiple cohort alignment
- Multiple assay comparison
- Long term stability reference

**Univariate analysis**

linear or generalised linear models

**Multivariate analysis**

PCA; PLS regressions

**Dimension reduction**

Penalised regression; Bayesian variable selection

**Metabolome wide significance level**

**External independent validation and large-scale replication**

**Cross validation**

**Marker Identification**

- Statistical spectroscopy
- Database look up: Bayesian model fit
- MS: adduct/isotope pattern/fragmentation prediction/ & matching

**Bioinformatics, cross-platform and multi-omics analyses**

- Over representation
- Gene/metabolite enrichment analysis
- Induced network analysis
- Differential network analysis
- Pathway analysis/visualisation

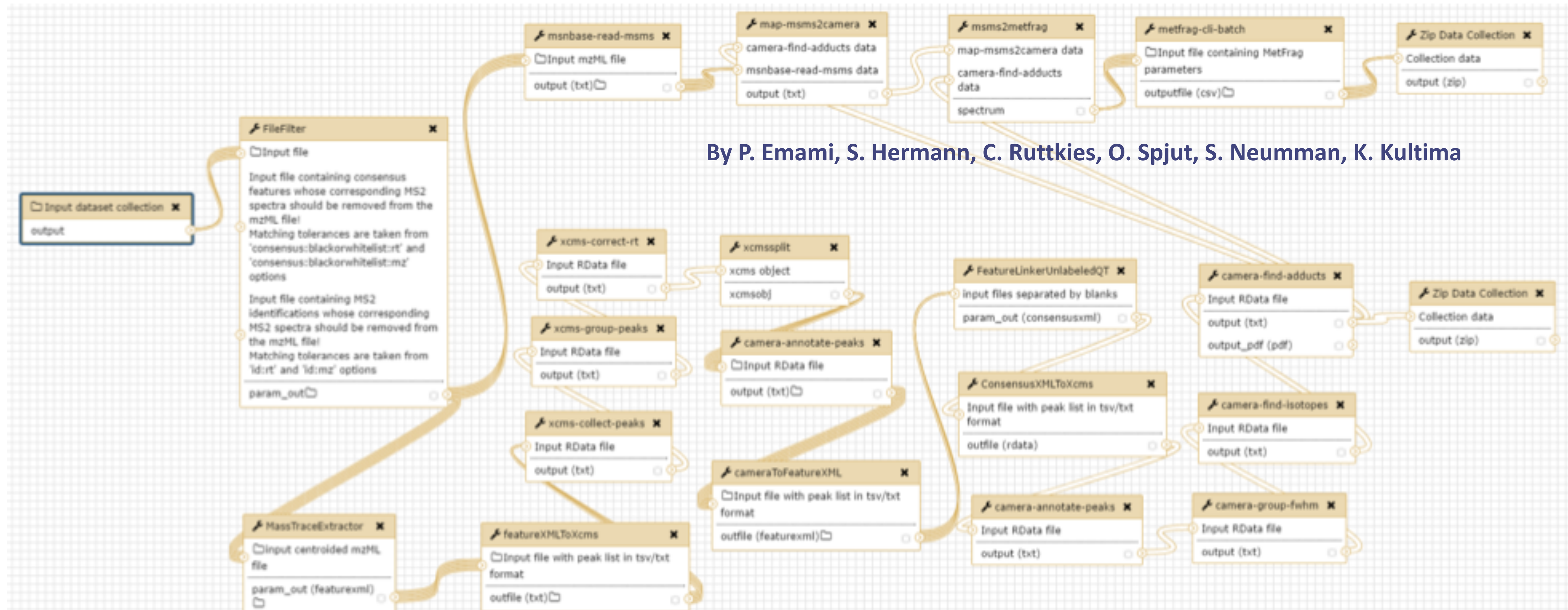**Challenge studies and randomised control trials**
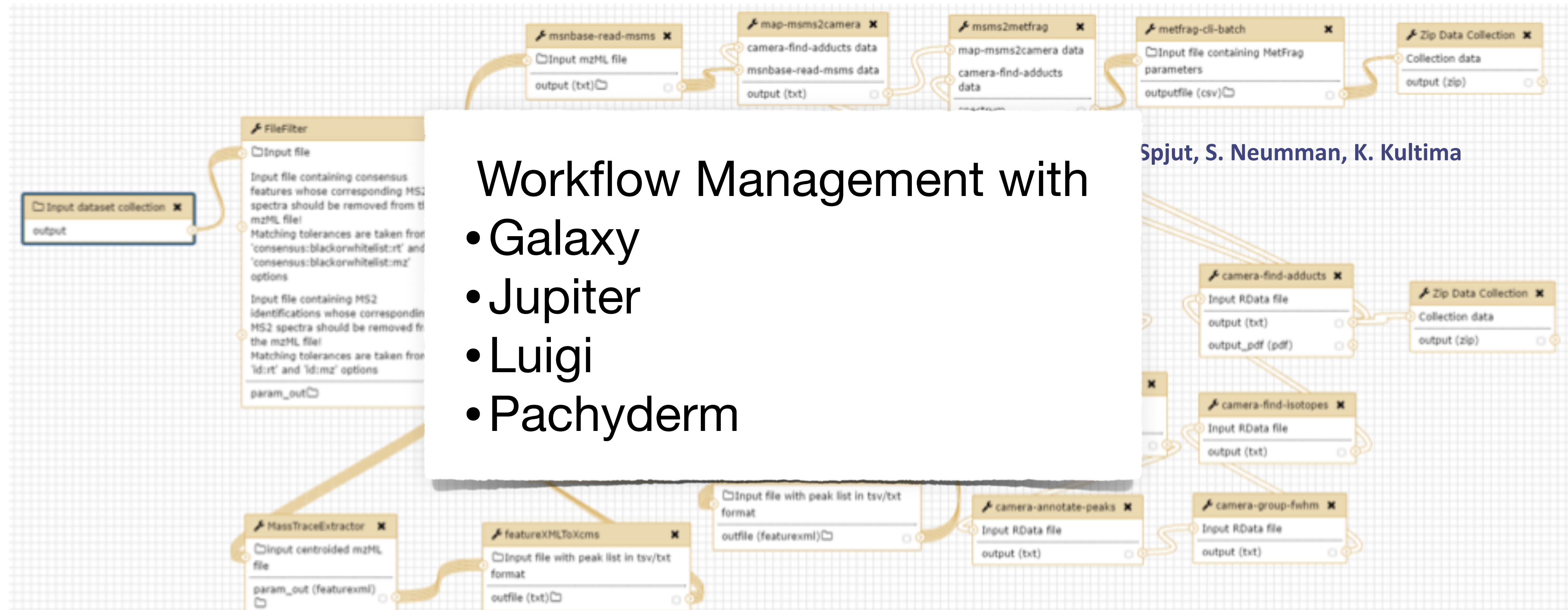
# Microservices



- Architectural design pattern

- Independent, potentially distributed processes

- Language-independent interfaces.

- Services decoupled. Perform a small task ("Do one thing and do it well").

- Individual service should be easy replaceable.

# MS1/MS2 – XCMS-OpenMS



By P. Emami, S. Hermann, C. Ruttkies, O. Spjut, S. Neumman, K. Kultima

# MS1/MS2 – XCMS-OpenMS



Workflow Management with
- Galaxy
- Jupiter
- Luigi
- Pachyderm

Spjut, S. Neumman, K. Kultima

100s of common metabolomics tools ready to run as Galaxy workflows in the cloud of your choice

# Large Scale Computing with Medical Metabolomics Data



**AWS**

**Azure**

**Literally a click of a button**

**PhenoMeNal**

A ready-made, well-tested, best-practice, Virtual Research Environment

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

# Managing large sets of machines in the cloud

# Generic Cloudification

# Economies of Scale

# Economies of Scale



- Assume 24.000 samples. Each takes 1 minute to process. That is 400 CPU core hours, or 16 days, or 2 weeks.

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

PhenoMeNal
*Large-Scale Computing for Medical Metabolomics*

# Economies of Scale

- Assume 24.000 samples. Each takes 1 minute to process. That is 400 CPU core hours, or 16 days, or 2 weeks.
- At $0.04 per core hour, plus 1TB storage, that's $16 in the cloud.

# Economies of Scale



- Assume 24.000 samples. Each takes 1 minute to process. That is 400 CPU core hours, or 16 days, or 2 weeks.
- At $0.04 per core hour, plus 1TB storage, that's $16 in the cloud.
- It is also $16 if you rent 400 cloud servers and you'd be done in 1 hour.

FRIEDRICH-SCHILLER-
UNIVERSITÄT
JENA

PhenoMeNal
*Large-Scale Computing for Medical Metabolomics*
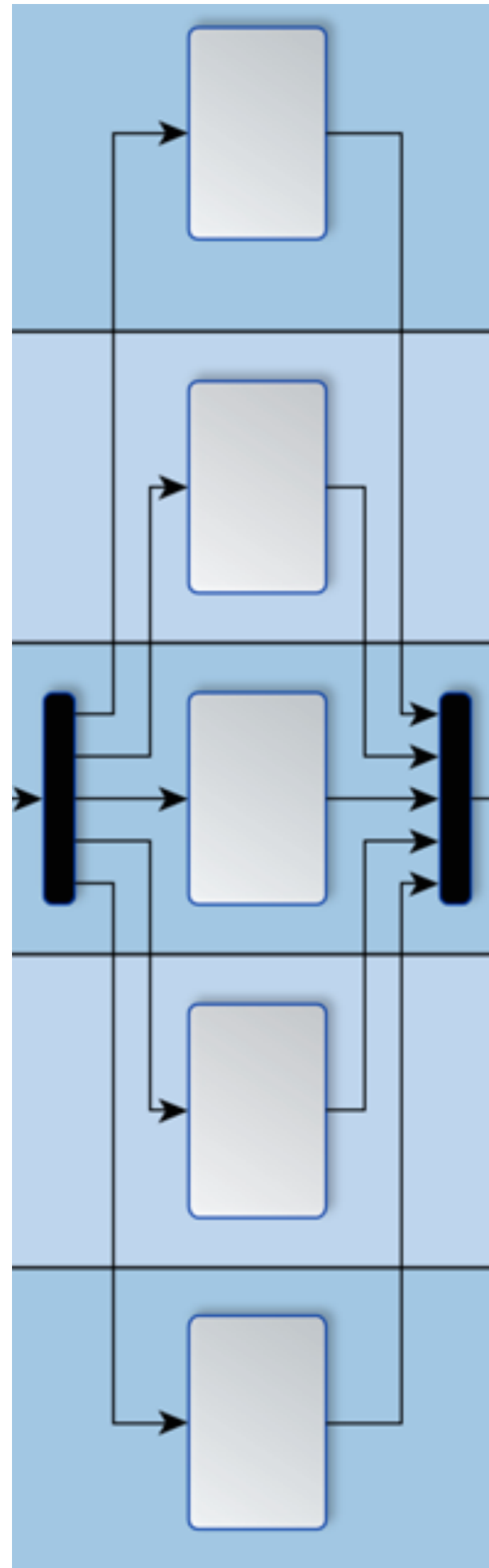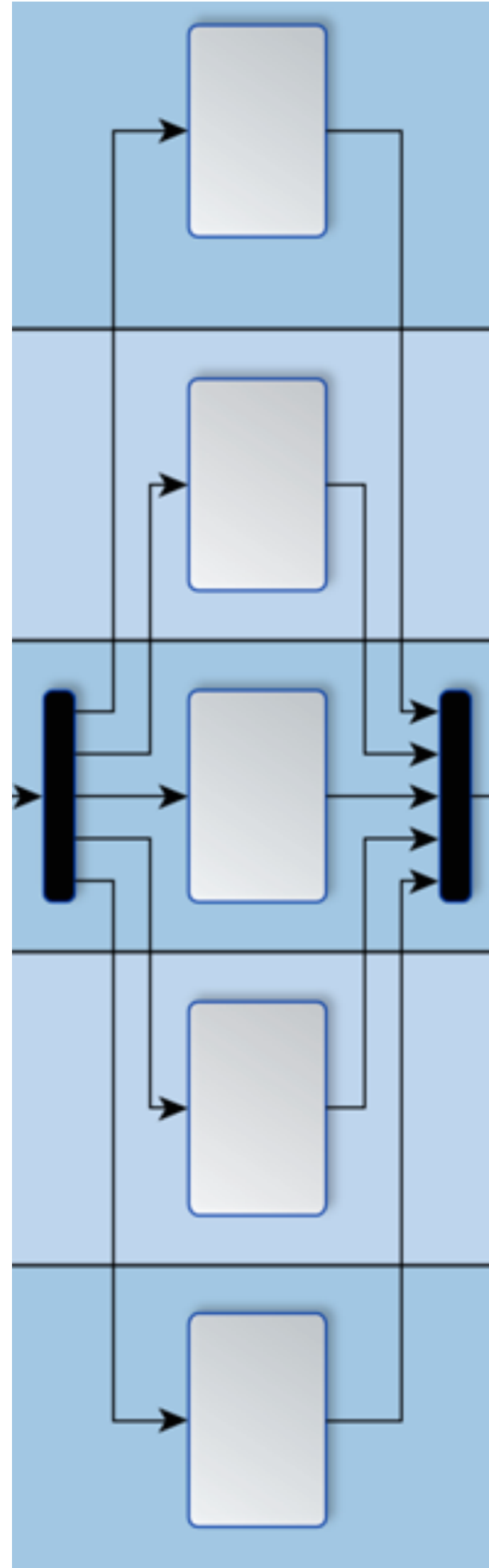
# Economies of Scale

- Assume 24.000 samples. Each takes 1 minute to process. That is 400 CPU core hours, or 16 days, or 2 weeks.
- At $0.04 per core hour, plus 1TB storage, that's $16 in the cloud.
- It is also $16 if you rent 400 cloud servers and you'd be done in 1 hour.
- It is also $16 if you rent 24.000 cloud servers, and be done in 1 minute (ignoring overhead …)

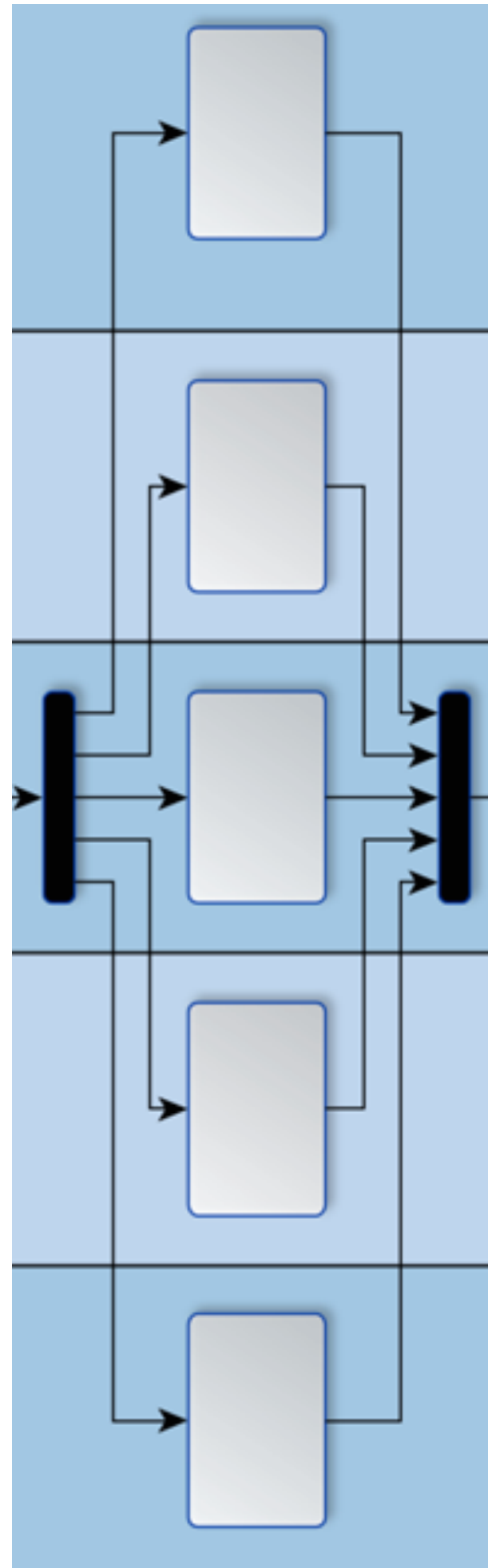# Economies of Scale
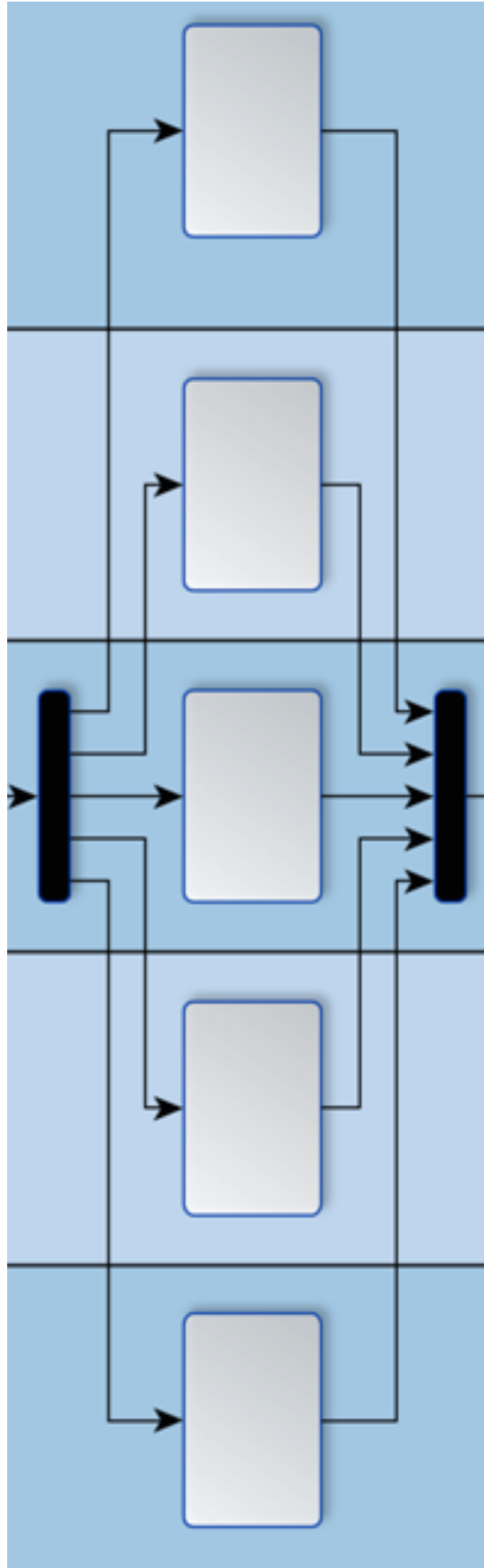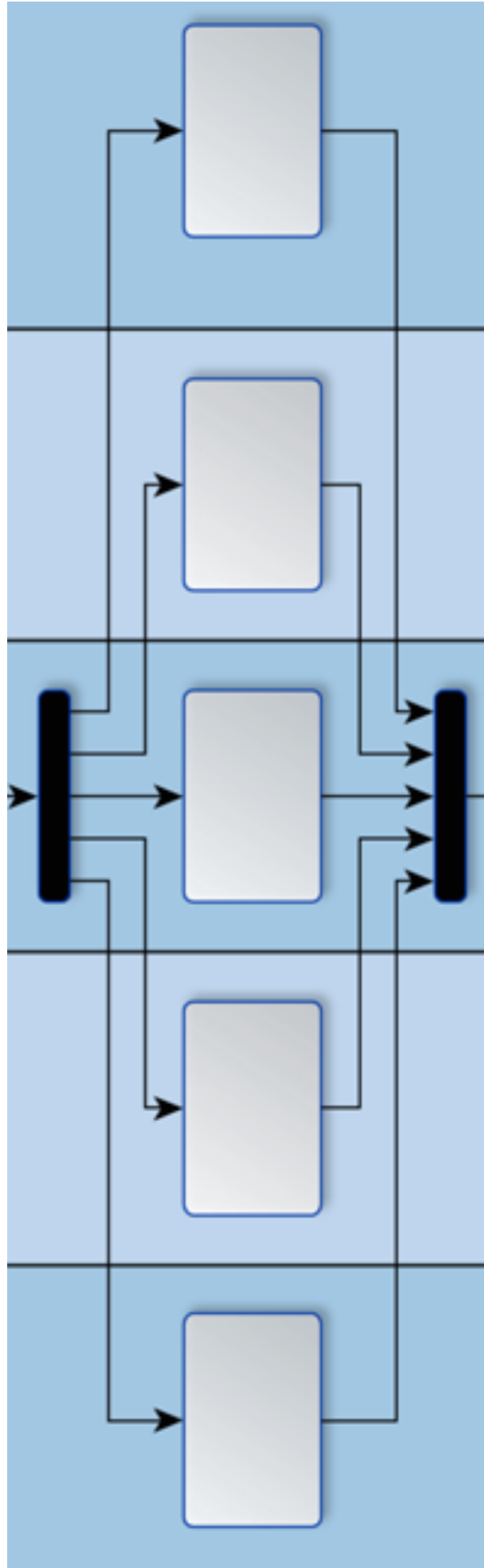
- Assume 24.000 samples. Each takes 1 minute to process. That is 400 CPU core hours, or 16 days, or 2 weeks.
- At $0.04 per core hour, plus 1TB storage, that's $16 in the cloud.
- It is also $16 if you rent 400 cloud servers and you'd be done in 1 hour.
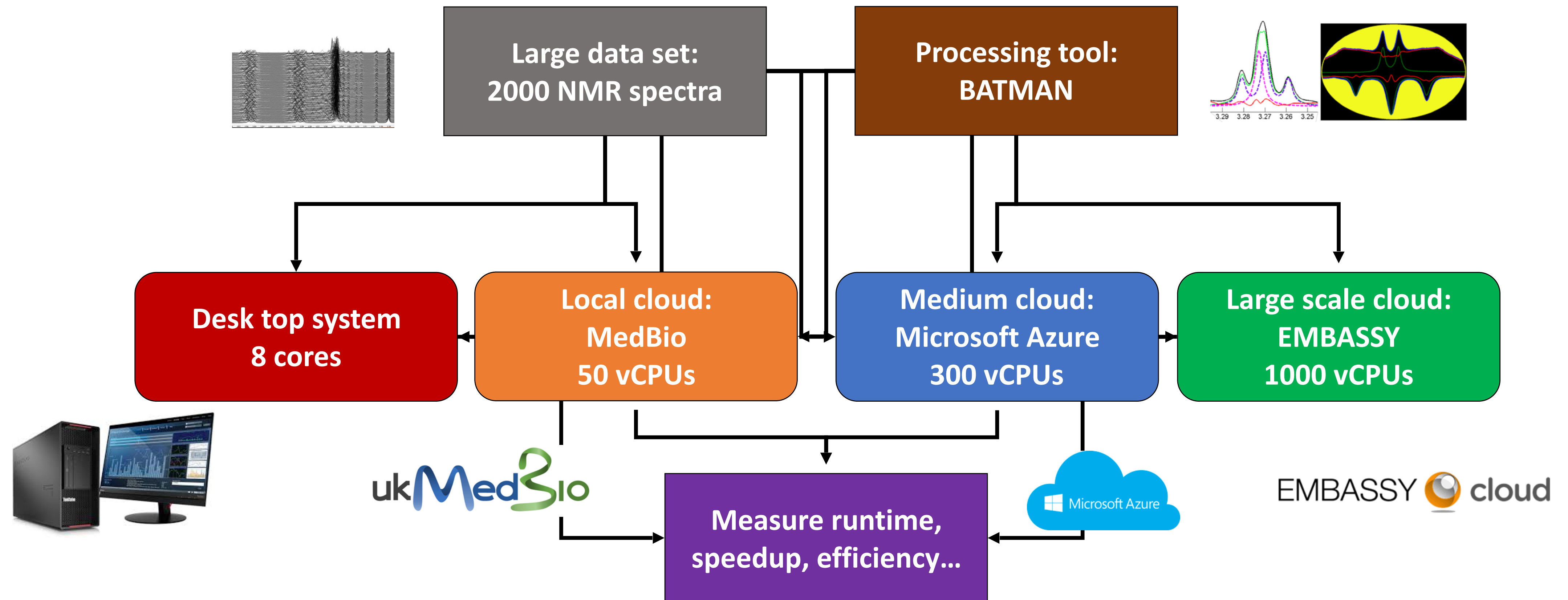- It is also $16 if you rent 24.000 cloud servers, and be done in 1 minute (ignoring overhead …)

FRIEDRICH-SCHILLER-UNIVERSITÄT JENA

**PhenoMeNal**
*Large-Scale Computing for Medical Metabolomics*

# Testing scalability

# PhenoMeNal scalability – BATMAN NMR processing



- 4 systems:
  - Desktop (8 cores)
  - MedBio (50 vCPUs)
  - Azure (300 vCPUs)
  - EMBASSY (1000 vCPUs)

- Considerable speed up possible
  - Running time for 2000 spectra down from **3 days** (1 core) to **10 mins** (1000 vCPUs)

FRIEDRICH-SCHILLER-UNIVERSITÄT JENA

PhenoMeNal
*Large-Scale Computing for Medical Metabolomics*

# Summary



- Large analytical chemistry data produced in medical metabolomics

- Computational analysis can take days or weeks on a single node

- **Industry standard tools** such as Kubernetes, Terraform and Ansible allow us to

    - seamlessly **deploy analysis pipelines**

        - composed of **microservices**

        - on **1000s of cloud machines**

    - at **low costs** and **no upfront investment**

- The free and open PhenoMeNal infrastructure encapsulates those **orchestration tools** together with **hundreds of tools in computational metabolomics**

# Funding and Collaborators

European Commission

Grant # 654241

PhenoMeNal Consortium

Netherlands Metabolomics Centre

4 Wm Workflow4metabolomics

METABOHUB

EMBL-EBI

Imperial College London

isasoftwaresuite

MRC Human Nutrition Research
Improving health through nutrition research

UNIVERSITY OF CAMBRIDGE

FRIEDRICH-SCHILLER-UNIVERSITÄT JENA