# ASPIS

*ASPIS – A shield used in ancient Greece; revolutionary defensive tool enabling advancing legions to join forces for better protection*

# Enabling innovation: from data science research to regulatory application

Philippe Rocca-Serra | University of Oxford

email: philippe.rocca-serra@oerc.ox.ac.uk
linkedin: /philipperoccaserra/
github: @proccaserra
orcid: 0000-0001-9853-5668

JRC meeting

- 31.01.2024

# The ⬤ASPIS Cluster

## PRECISION TOX

Leveraging evolutionary diversity to reveal the molecular basis of toxicity
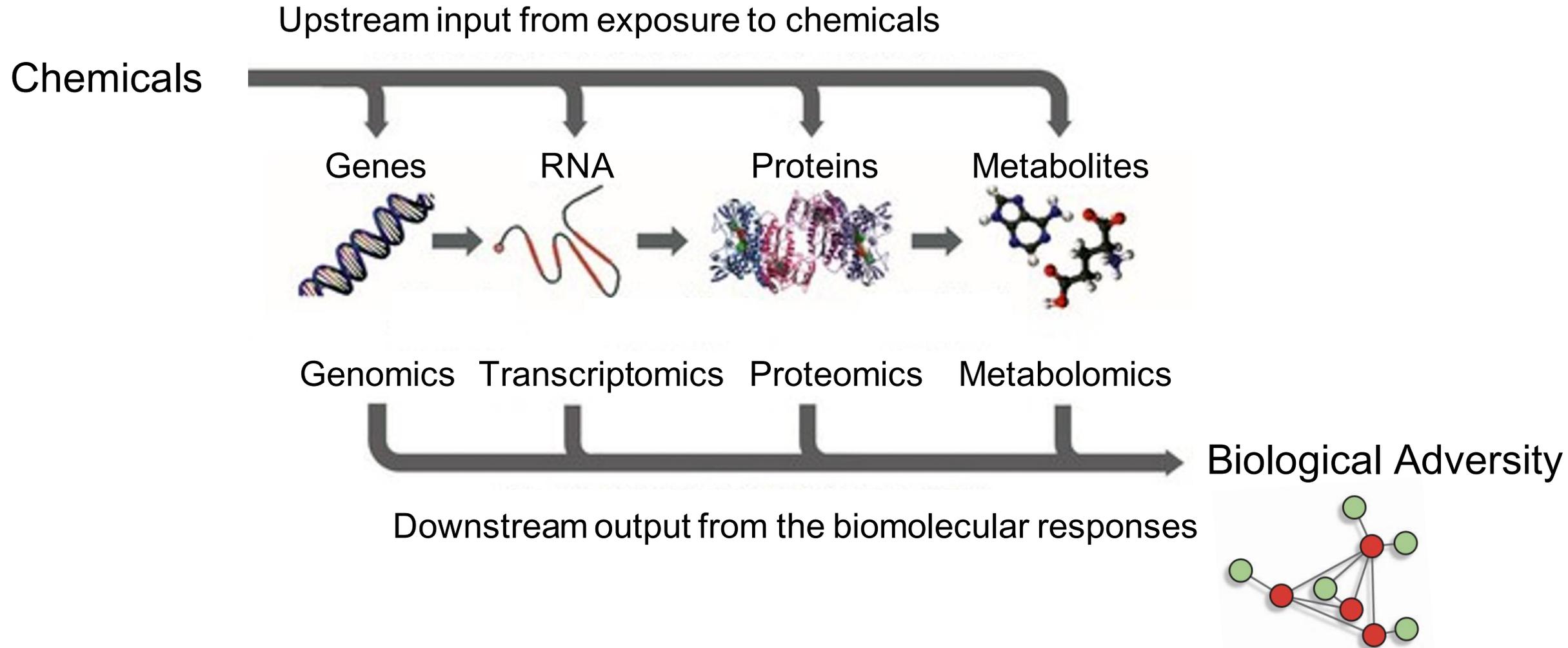
## ⬤NTOX

Synthesising toxicology knowledge to support next generation risk assessment
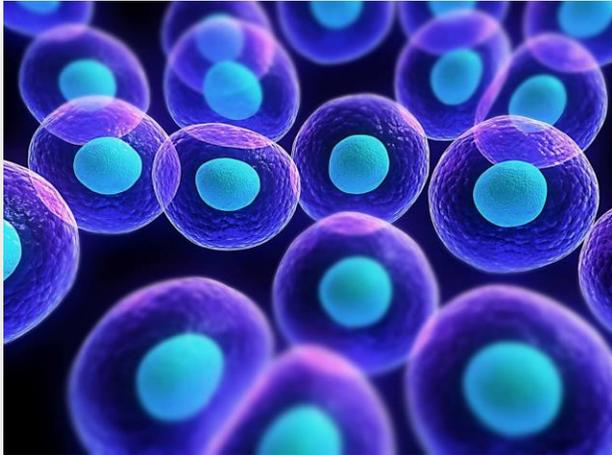
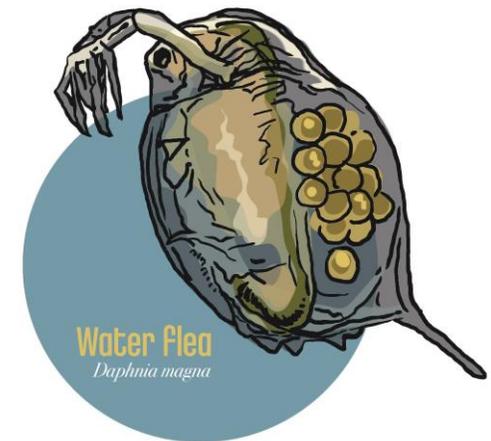## RISK [∷∷] HUNT3R

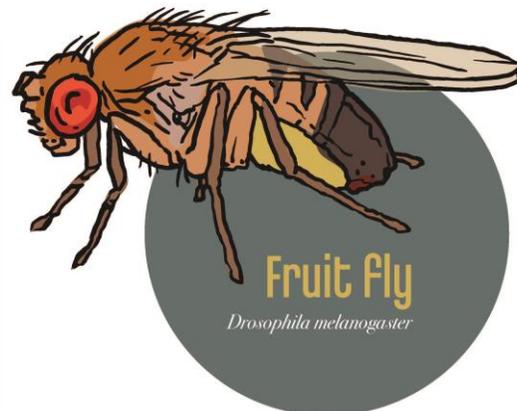Human-centric chemical safety assessment utilising systems toxicology
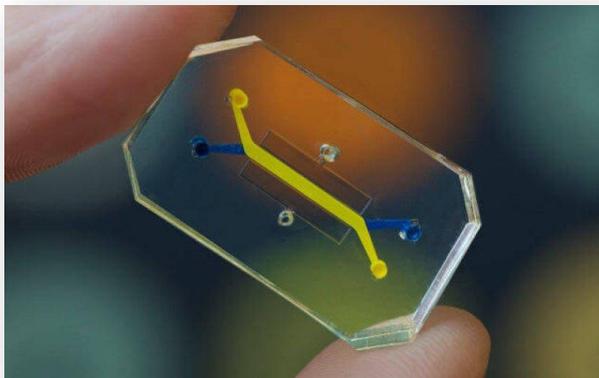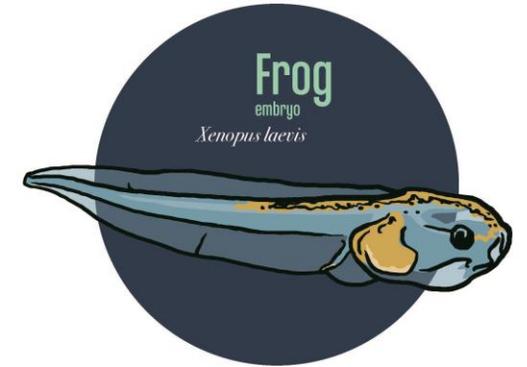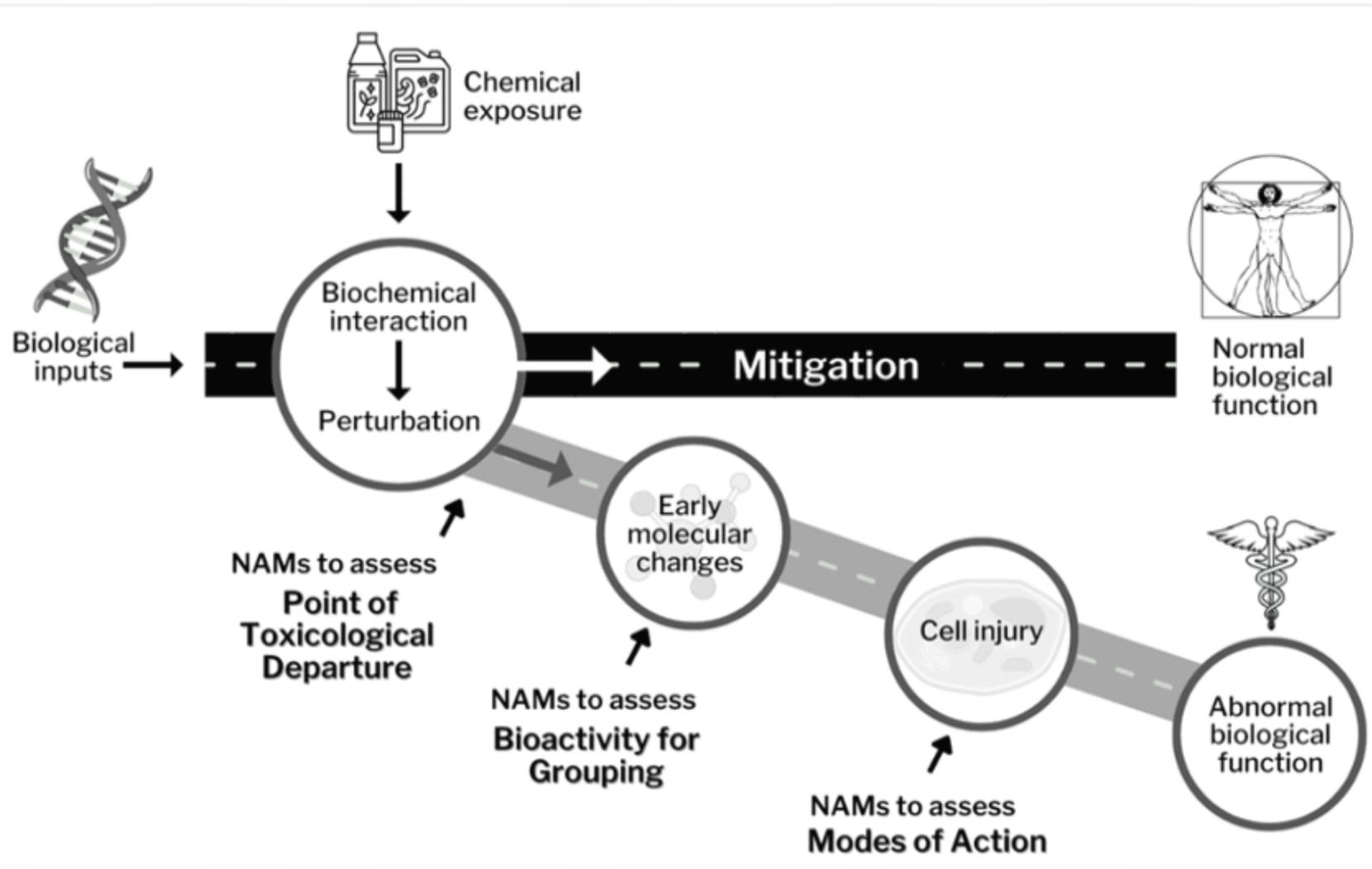
# Biomolecular Research Data

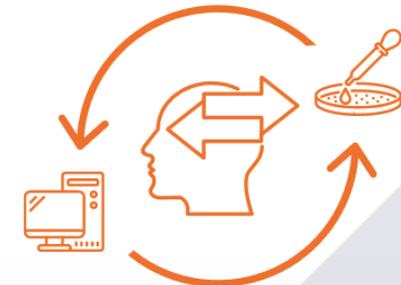# Experimental Data from Variety of Testing Platforms



Existing Data

# The New Paradigm: Activation of Toxicity Pathways



Each NAM will consist of an ontology-driven and artificial intelligence-based computational system linked with a battery of *in vitro* assays and *in silico* tools for hazard prediction combined with customized exposure assessment.

# Why is F.A.I.R. important for innovation and regulatory applications?

- FAIR means self-describing data

- FAIR means explicit semantics (commitment)

- FAIR means availability of data to software agents ("machine actionability")

- By ensuring sufficient metadata always accompany data, FAIR enables trust

- FAIR practices protect investments by enhancing data reuse potential

PRECISION
—TOX—

# Planning for FAIR from the onset

- Survey of the art - catalogues of standards (FAIRsharing) - Promoting reuse
  - Domain-specific data repositories & associated requirements
  - Data deposition syntax
  - Controlled vocabularies / terminologies
- Semantic models and identifiers patterns
  - Preparing for Extraction Transform Load (ETL) to RDF & KG from data sources
- Licensing and access rights around in machine actionable form
- Data release and cataloguing for optimum discoverability and reuse
- Consolidation in an extensive and detailed Data Management Plan

# Testing process during pilot phases

- As with all projects, a central issue is that of phasing:
  - Delivering the right tool at the right moment
  - Developing lean processes as mitigating plans
    - PrecisionTox Spreadsheet Templates with naming conventions
    - Collect signal for refine software specifications

- Taking advantage of projects pilots
  - Coordination user requirements for data collection
  - Defining annotation requirements with subject matter experts
  - Developing software solution prototypes
  - Following software engineering best practice standards
    - Continuous integration, Code review, test driven development, automation
    - Code documentation, user documentation…

PRECISION
TOX

# Key concept: prospective FAIR approach

- Implementation of **Metadata Manager Tool -** PrecisionTox example

- Globus based file transfer process between data producing centers

- Data publication / dissemination procedure

- Embedding FAIR practice in computational workflows

- Ready for knowledge graph representations

# Developing software to make data FAIR by design

# Developing software to make data FAIR at design

## Precision Toxicology Metadata Manager: sample exposure and collection



Batista Dominique (0000-0002-2109-489X), Data Readiness Group

# Developing software to make data FAIR by design



{"sample":{"batch":"AA","compound":{"name":"Ethoprophos","ptox_id":"PTX002"},"dose":"BMD10","ptox_id":"RAA002LA1","replicate":1,"timepoint_hours":10,"vehicle":"Water"}}

# Integrated Data and Knowledge Management

# Data Management and Analysis

## Behind the scenes

| Data structure | Data import | Data use |
|---|---|---|

**Data structure**
- Excel template file (EU-ToxRisk)
- Required fields for all the datasets
- Custom fields if needed

| Column Name | Constant Values |
|---|---|
| Sample ID | |
| Method name | method name indicated in the UKN3b_NeuroTox_LUH_neuri |
| Toxicity domain | Neuro |
| Information domain | Cytotoxicity |
| Date | |
| Experiment ID | |
| Organization abbreviation | |

**Data import**
- Each dataset is validated
- Published on Biostudies
- Automatically imported to EdelweissData
- Data becomes accessible through web requests (URL)

EdelweissData™
Convenient publishing of scientific data with proper versioning, rich metadata support and a powerful API

**Data use**
- Access data & metadata directly from the database
- Target specific version of dataset (e.g. latest)
- Consume data directly into your data analysis tool (R, Python, Excel, Jupyter, Colab, Observable)

Observable

# Data Management and Analysis Workflows based on Biostudies and EdelweissData



Version 1.0 of EdelweissData (harmonised data management product of Edelweiss Connect) completed by 2020 and in commercial use.
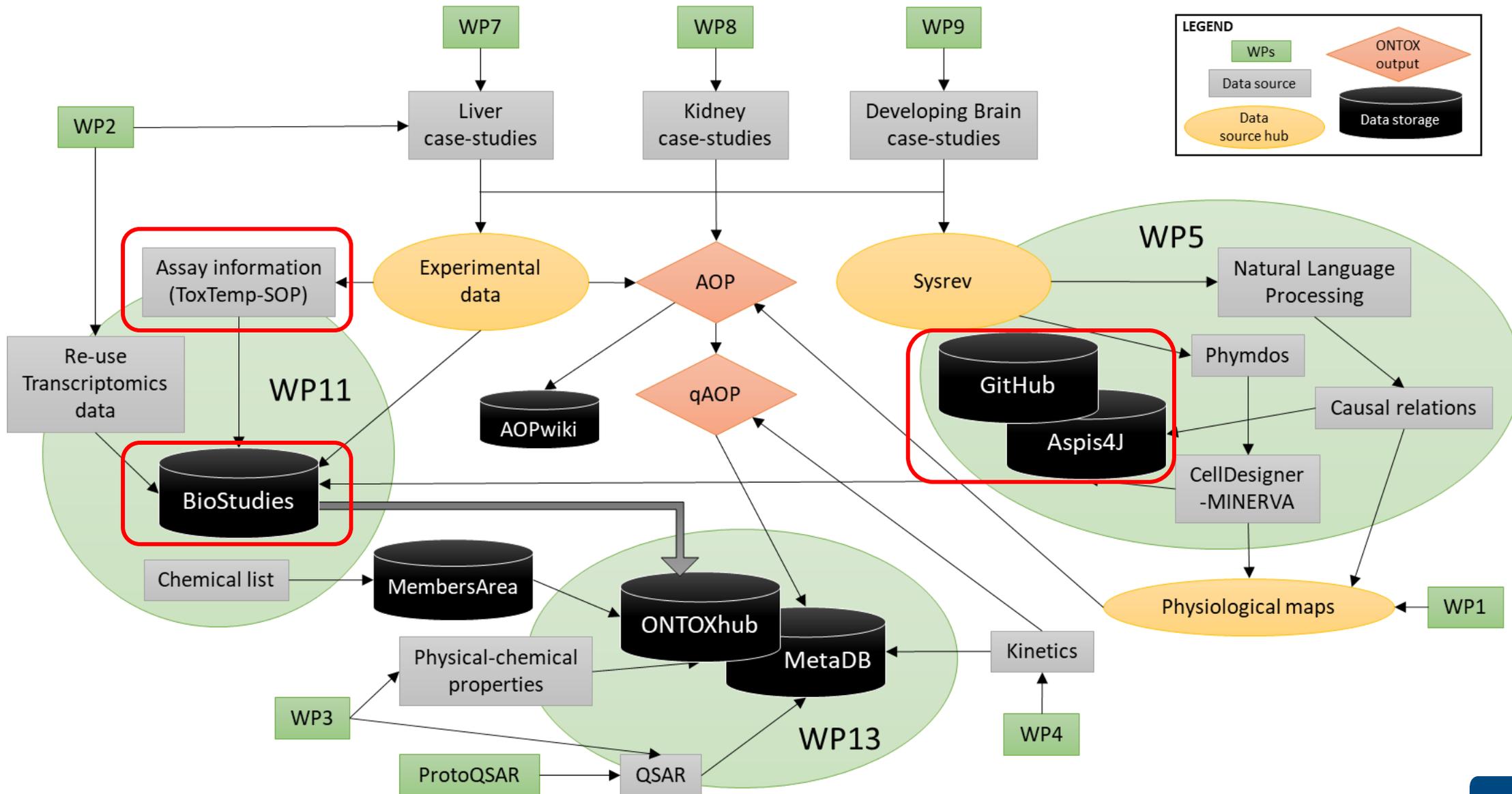
**RISK HUNT3R**

Supporting data life cycle from generation through to consumption.

https://edelweissdata.com/

Credits:
Barry Hardy (Edelweiss Connect)
Email: Barry.Hardy@edelweissconnect.com
LinkedIn: https://www.linkedin.com/in/barryhardy/

# ONTOX data flowchart

# PrecisionTox and ELIXIR collaboration





Convergence & Alignment

Packaging Metadata, Data and Computational Worflow as a Research Object

# Promoting adoption and convergence

# Knowledge Graph (EU-ToxRisk data being extended to ASPIS)

# Acknowledgments

Pr. John Colbourne (UoB)

Dr Ralf Weber, Dr Martin Jones, Marianne Barnard (UoB)

D Batista (PhD candidate) (UOXF)

Dr Juan Carlos Gonzales-Sanchez, Prof Rob Russel (HEI)

Dr Nadine Cistiakova (MGI)

Dr Anna Vogt (CRG)

All our colleagues from PrecisionTox

Dr. Barry Hardy (CEO Edelweiss Connect)

All our colleagues from RISK-HUNT3R

Dr. Danyel Jennen

All our colleagues from Ontox