# Linking enviPath to systems biology and sequencing data
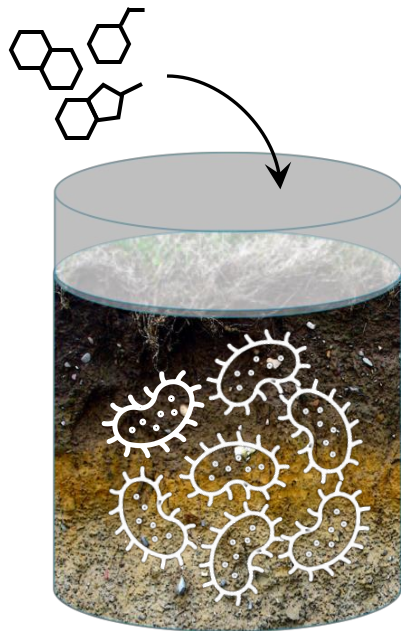
Jasmin Hafner, Albert Anguera Sempere,

Kathrin Fenner

11 May 2025

# Overview

# Why systems biology and sequencing data?
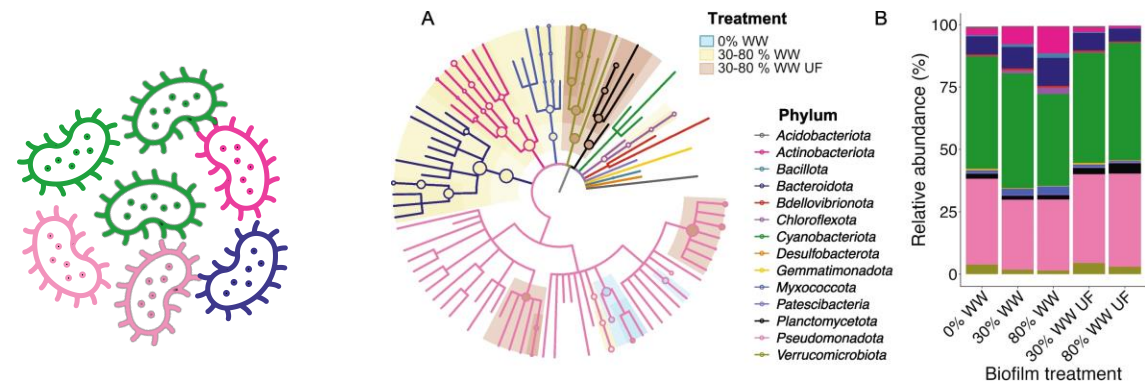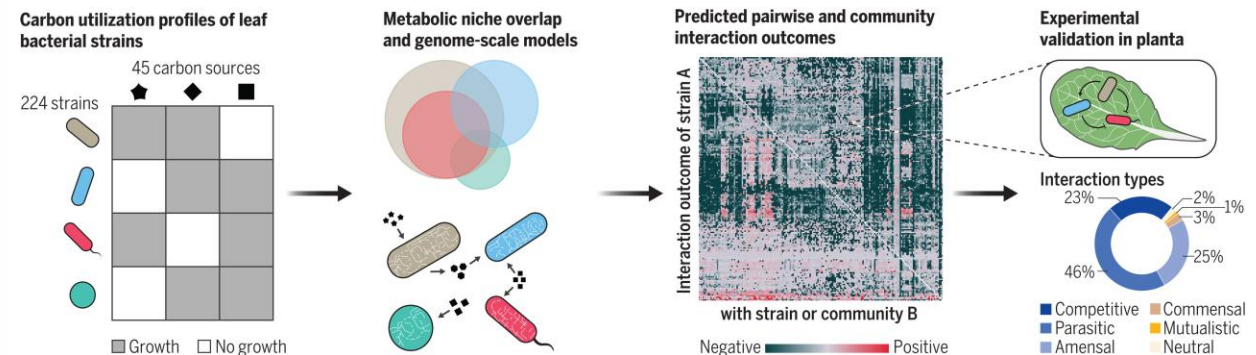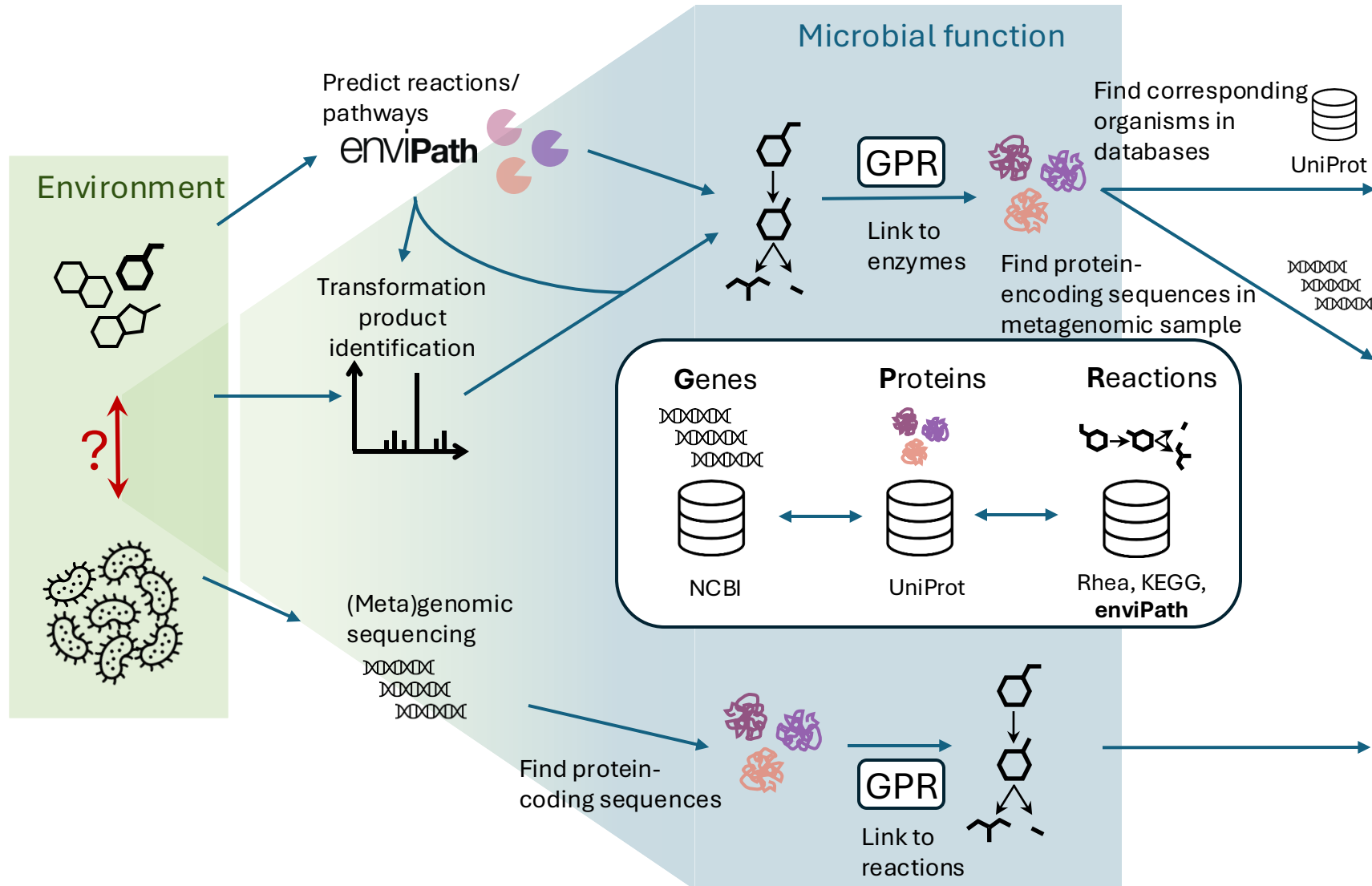
**Taxonomic composition** (sequencing data)

WHO?

**Enzymatic composition, community function** (metabolism / community interactions)

WHAT?
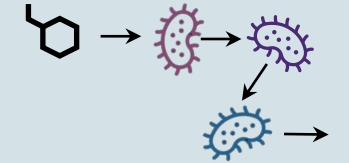
Microbial communities

Figures: Attrah et al., *Microbiome*, 2024., Schäfer et al., *Science,* 2023

# Different approaches – same challenge

Objectives

Microbial function

Environment

Predict reactions/pathways
enviPath

Find corresponding organisms in databases
UniProt

GPR

Link to enzymes

Find protein-encoding sequences in metagenomic sample

Transformation product identification

**Genes**          **Proteins**          **Reactions**

NCBI          UniProt          Rhea, KEGG, **enviPath**

?

(Meta)genomic sequencing

Find protein-coding sequences

GPR

Link to reactions
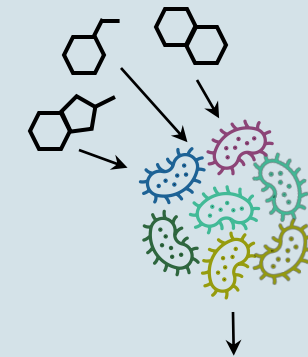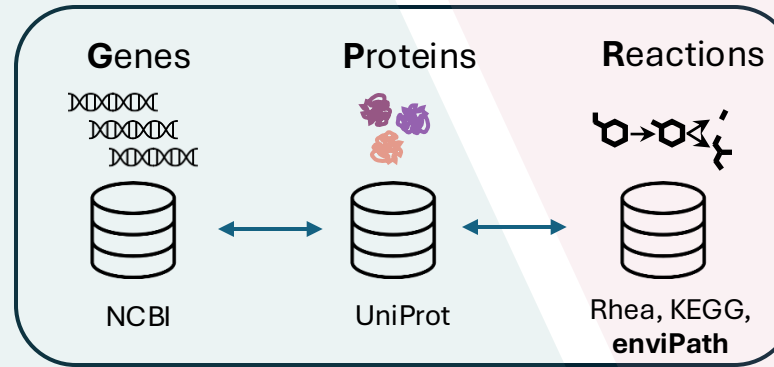
**Engineer** microbial communities

**Functional characterization** of microbial communities

4

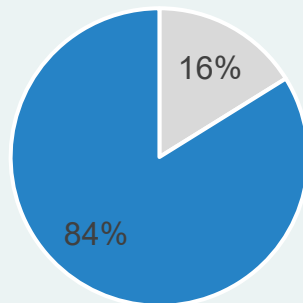# Challenge: Link function (reaction) to sequence (gene, protein)

- **Orphan gene / protein sequences** (no function associated)
- Promiscuous enzyme activity

**Genes**

**Proteins**

**Reactions**

NCBI ⟷ UniProt ⟷ Rhea, KEGG, **enviPath**

SwissProt (curated): ~ 500K proteins

UniProt

Trembl: > 200 M proteins

16%
84%

☐ Other
☐ Sequence linked to KEGG reaction

10%

**Unknown enzymes?**

**Orphan reactions** (no gene associated)

SIB Rhea

**Rhea:** 4,760 out of 17,614 reactions
→ **27% orphan**

**enviPath:** **?**
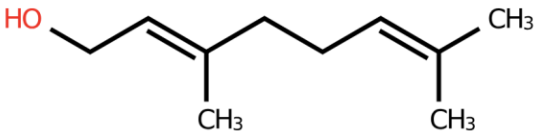
Rhea and UniProt consulted in April 2025

# Overview

1. Why **systems biology** and **sequencing** data?

2. Linking compounds and reactions to **external databases**

3. Briding the **Gene-Protein-Reaction** (GPR) knowledge gap

4. Outlook: Connecting to **metagenomes**

# Linking compounds to external databases



| Geraniol | 🔧 Actions ▾ |
|---|---|

**Image representation**

**Database look-up**
via SMILES and
InChiKey

**SMILES representation**

CC(=CCC/C(=C/CO)/C)C

**Canonical SMILES**

OCC=C(C)CCC=C(C)C

**InChIKey** — NEW

GLZPCOQZEFWAFX-JXMROGBWSA-N

**Description**

**Pathways**

**External Identifiers** — NEW

**PubChem Compound Identifier**

CID4458
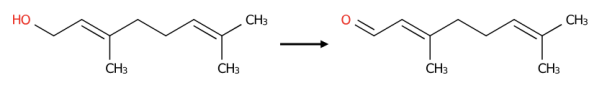CID637566

**ChEBI Identifier**

CHEBI:17447
CHEBI:24221

**KEGG Identifier**

C01500

# Linking reactions to external databases

**Eawag BBD reaction r1163** ★

**Image representation**



**Reaction Description**

| Geraniol | → | Geranial |

**SMIRKS representation**

CC(=CCC/C(=C/CO)/C)C>>CC(=CCC/C(=C/C=O)/C)C

---

Use **ChEBI** IDs to find **Rhea reactions**



Retrieve **UniProt** links from Rhea



---

**EC Numbers**

geraniol dehydrogenase (1.1.1.183)

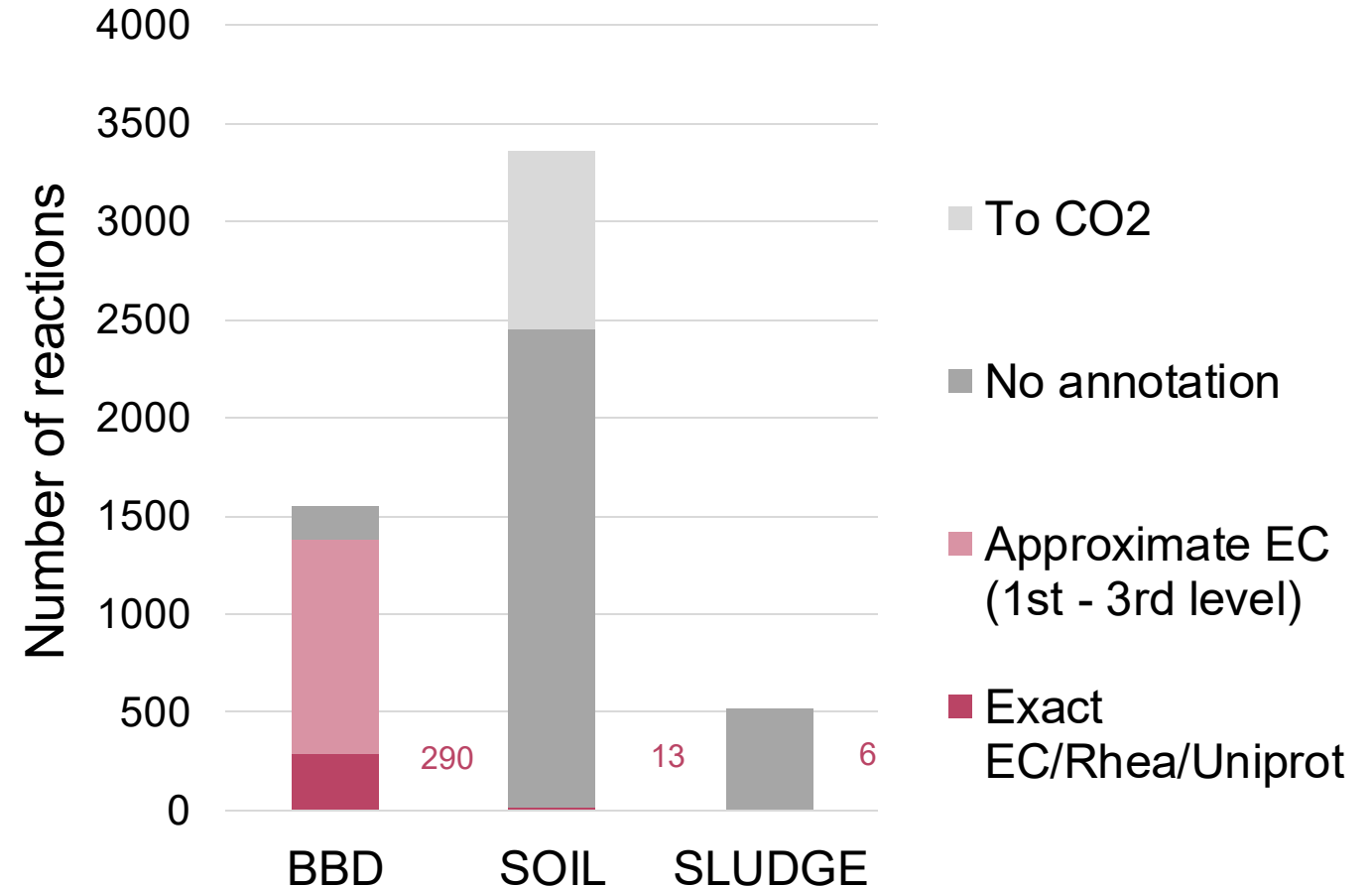**Pathways**
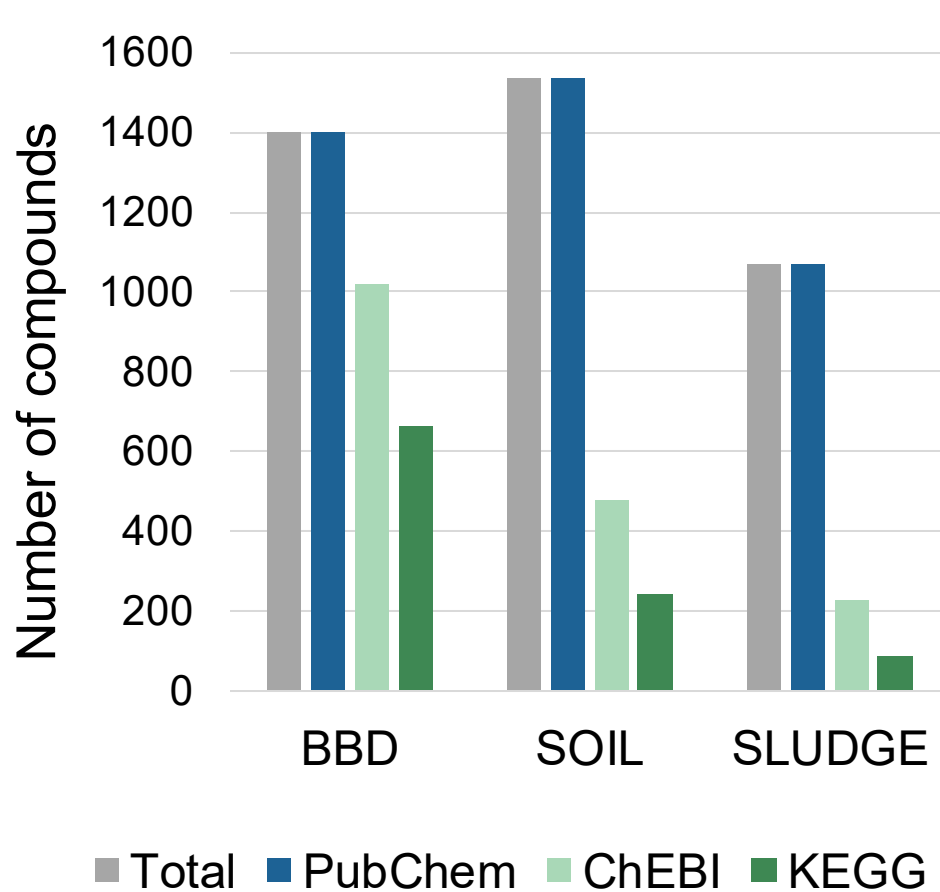
**References**

**NEW**

**Rhea**

34347
14521

**UniProt**

3 SwissProt entries (rhea:34347)
1 SwissProt entries (rhea:14521)

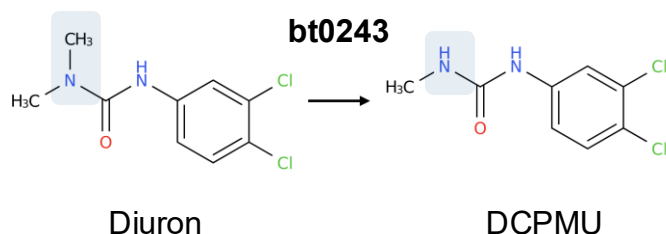# How many compounds and reactions could be linked?



In BBD, only **15% of reactions** (223) are linked to **UniProt** protein sequences – **85% are orphan**!

# Why most reactions don't have known enzymes?

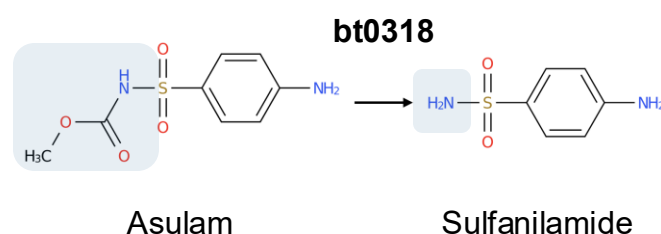Example: EAWAG-SOIL (excluded: reactions leading to $CO_2$)

| Known protein sequence | Mechanism known (bt rule) | Mechanism unknown (no rule) |
|---|---|---|



**bt0243**

Diuron → DCPMU

**bt0318**

Asulam → Sulfanilamide

Tetraconazole → M14360-DFA

**1%** of SOIL reactions

Linked to Rhea: 13
EC numbers: 8

**50%** of SOIL reactions

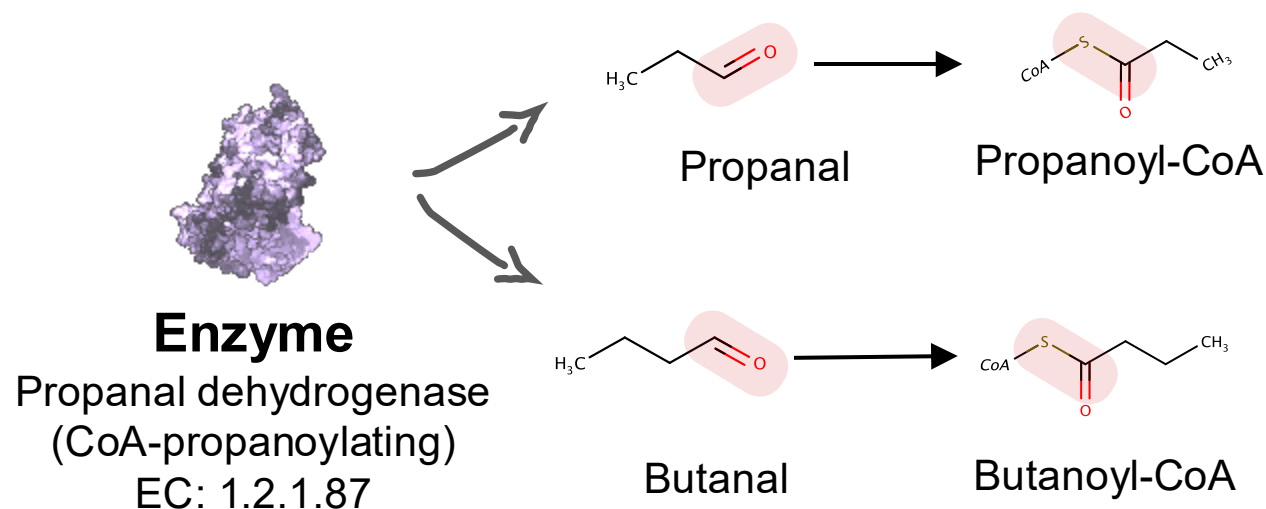**Potential to find sequences!**

**49%** of SOIL reactions

→ **Xenobiotic** metabolism not well covered by metabolic databases
→ **Co-metabolic** processes not well understood
→ Reported pathways in biodegradation studies **based on observed TPs** only

**Overview**

# Can we find potential protein sequences for orphan biotransformations?

**Assumption**: Similar reactions are catalyzed by similar/same enzymes

**Enzyme**
Propanal dehydrogenase
(CoA-propanoylating)
EC: 1.2.1.87

Propanal → Propanoyl-CoA

Butanal → Butanoyl-CoA

✔ Works well for **promiscuous** enzyme **activities**

✘ Difficult to find completely **new enzymes**

Orphan reaction

Create fingerprint
1001101001…

Compare to database of reactions linked to protein sequences

Find most similar reaction

Potentially catalyzing enzyme

# How can we link orphan reactions to enzymes?

Tools to **predict EC** numbers: Theia, BridgIT

Kunyang Zhang

Theia, BridgIT — **EC number** ← Rhea, KEGG

Direct link

**Reaction** ←→ **Protein**

|  | Theia | BridgIT | *Ideally* |
|---|---|---|---|
| Considers reactive site | no | yes | yes |
| Need for mass balance | no | yes | no |
| Open source | yes | no | yes |
| Bypass EC | no | no | yes |
| Xenobiotics metabolism | no | no | yes |

**Pretraining**
Chemical reactions

**Fine-tuning**
Biochemical & biodegradation reactions

AI-driven identifcation of relevant
**functional groups**

Important structures    Less important structures

Positive Contribution

Negative Contribution

**Similarity metric** for reaction comparison

13

# Overview

1. Why **systems biology** and **sequencing** data?

2. Linking compounds and reactions to **external databases**

3. Briding the **Gene-Protein-Reaction** (GPR) knowledge gap
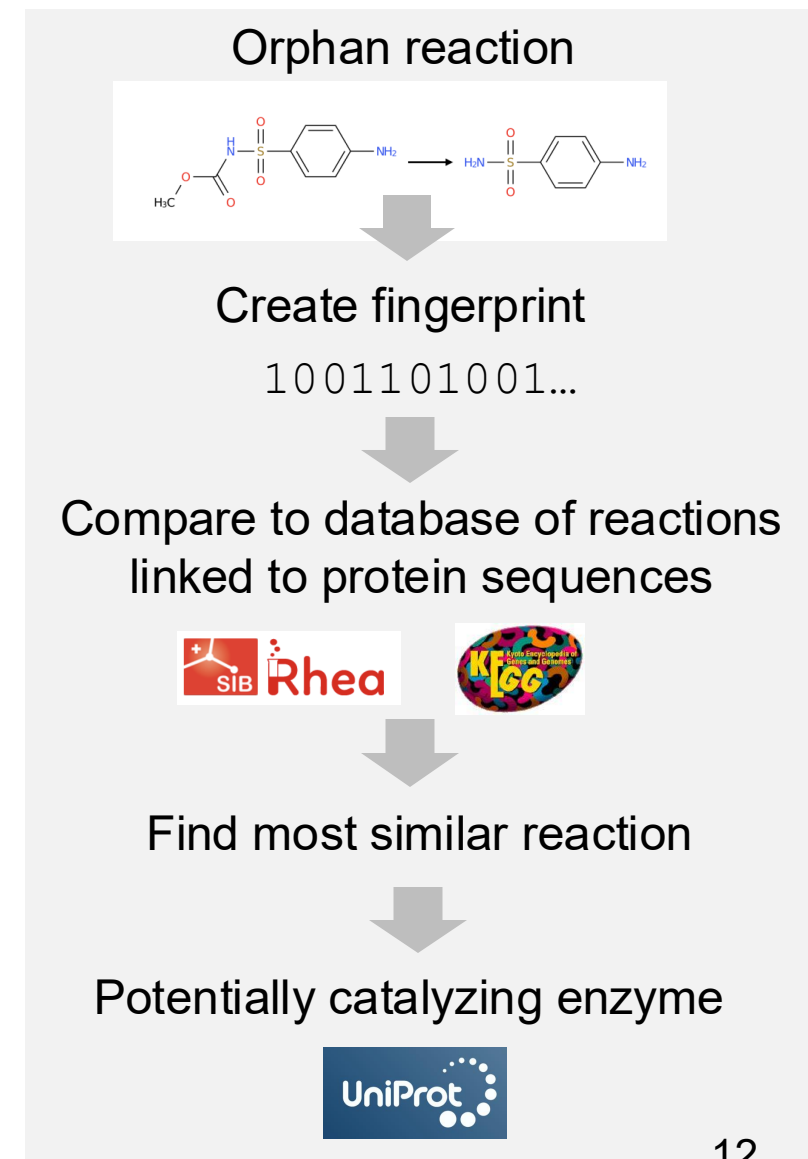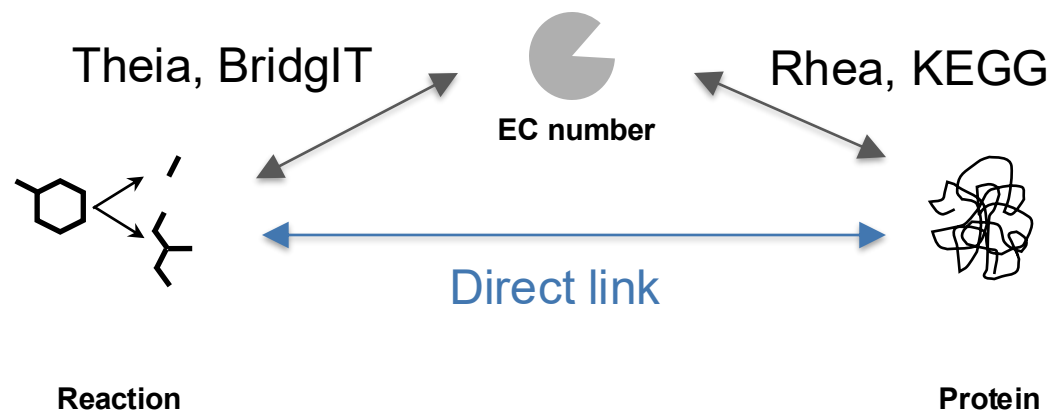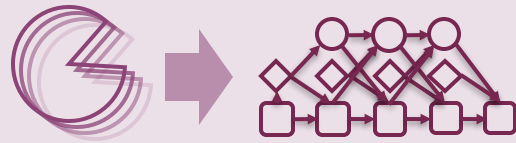
4. Outlook: Connecting to **metagenomes**

# New and ongoing developments

Repository of **HMM profiles** for **biotransformation** by

Dr. Serina Robinson, Victoria Poltorak

https://github.com/MSM-group/ContaHMM

Link to **metagenomic studies** in enviPath

Scenarios:
Add links to **ENA** and MG**nify** entries

Bridging **to Systems Biology**

Export **pathways** in Systems Biology Markup Language

https://envipath-python.readthedocs.io/en/develop/tutorials/download_pathway_SBML.html

# Take-home message

- Big **knowledge gaps** in linking enzymatic reactions to protein sequences, *especially* for xenobiotic biotransformations

- **Data curation** efforts needed to fill knowledge gaps

- In the future, **new ML/AI tools** may also help in filling those gaps

# Acknowledgements



Prof. Kathrin Fenner



Albert Anguera Sempere

## Funding sources



& the Fenner Team