# Artificial Intelligence /Machine Learning for Secondary Metabolite Prediction
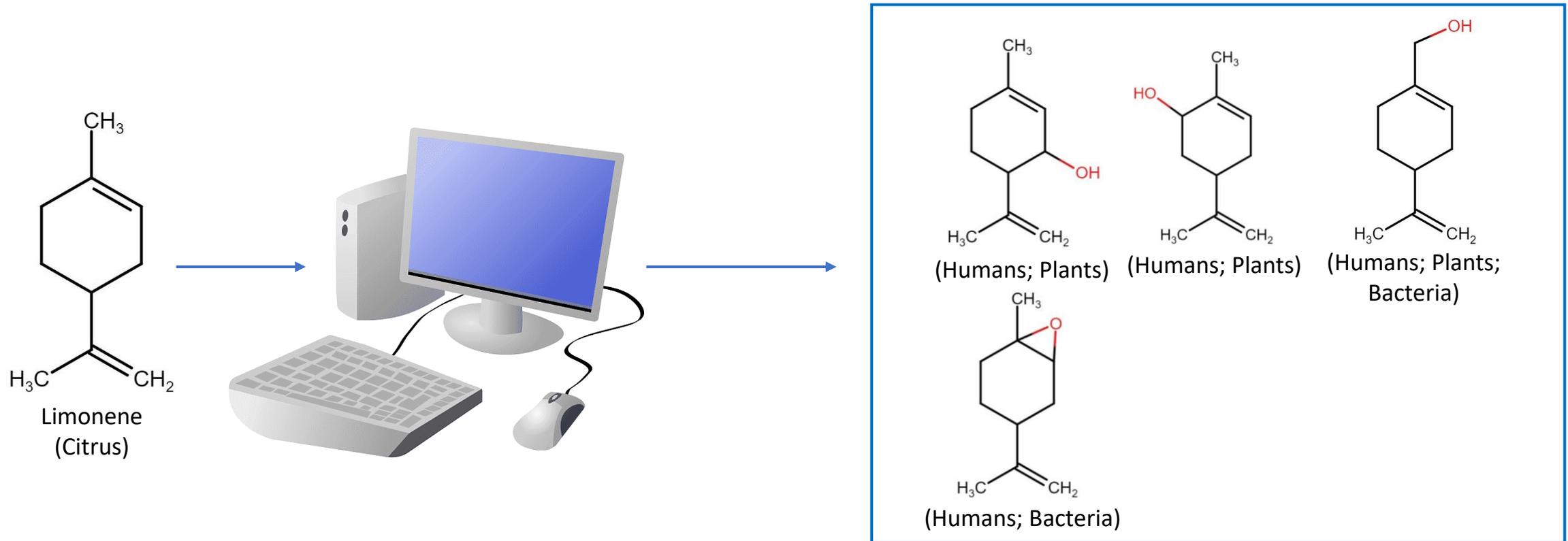
Yannick Djoumbou Feunang
Corteva Agriscience, Indianapolis, IN, US

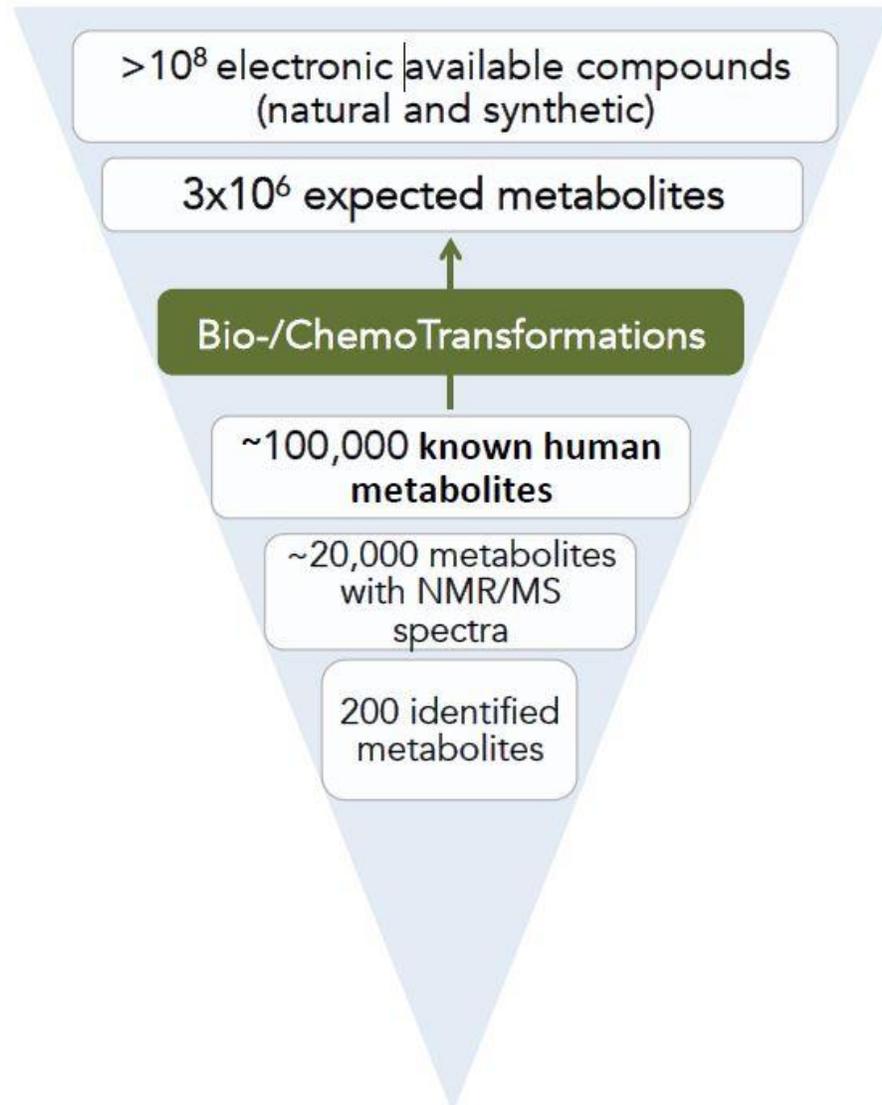Online workshop: Computational Applications in Secondary Metabolite Discovery
March 9, 2021

# *In Silico* Metabolism Prediction

**Task:** Given a small molecule, use <u>computational tools</u> to <u>predict</u> the outcome of its interactions with metabolic enzymes.

Here, we will focus on the <u>structural elucidation</u> of potential <u>metabolites</u>.
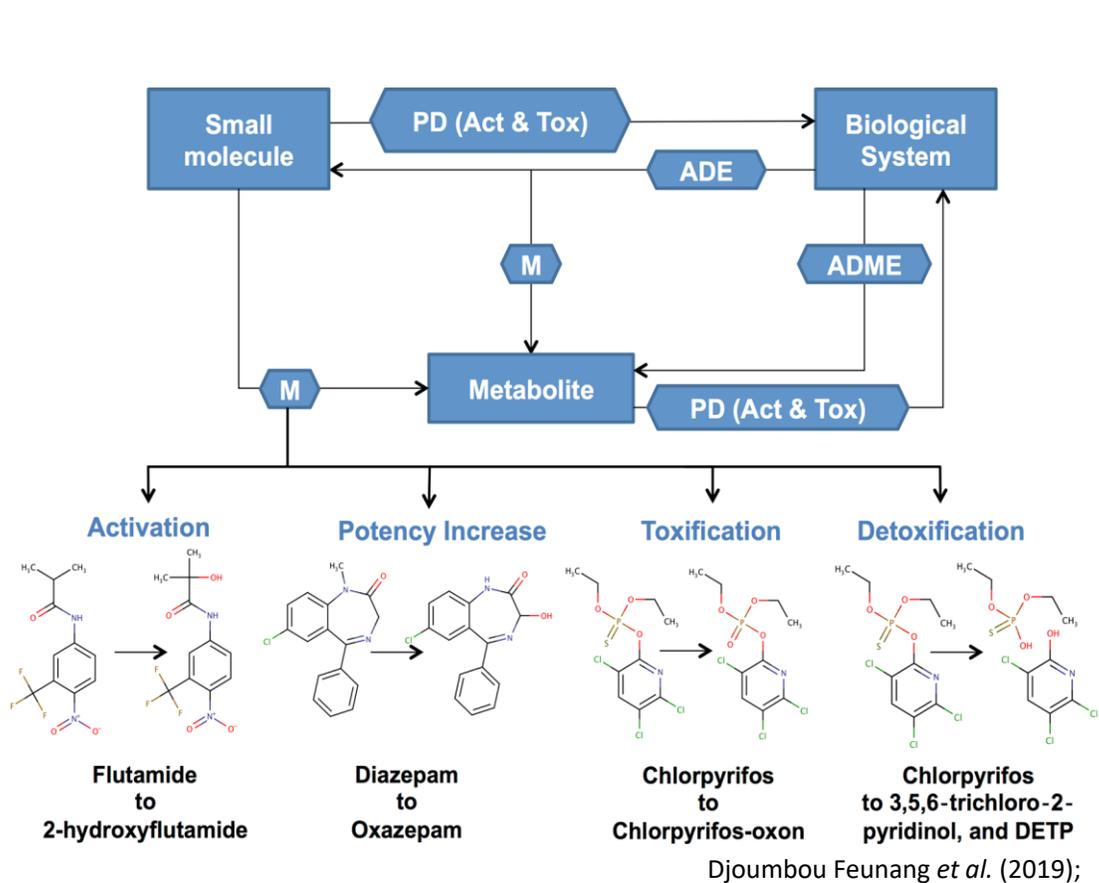
# Why is Metabolism so Important?



- Bio/chemo-transformations influence changes of our chemical exposome

- <2% of detectable peaks identifiable in large-scale untargeted Metabolomics

- What are these unknowns?

- How to determine their structures?

- How to determine their activities?

- Cheminformatics and AI can be used to:
  - Detect common patterns
  - Predict enzyme/ligand interactions
  - Understand metabolism
  - Generate biologically feasible structures through simulated reactions

The funnel diagram labels read:
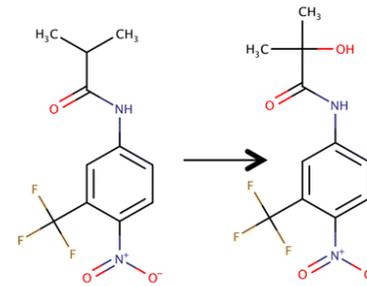- >$10^8$ electronic available compounds (natural and synthetic)
- $3 \times 10^6$ expected metabolites
- Bio-/ChemoTransformations
- ~100,000 known human metabolites
- ~20,000 metabolites with NMR/MS spectra
- 200 identified metabolites

# Why is Metabolism so Important?



**Pro/soft-drug design**

**Activation**

**Exposure science**

**ADMET profiling**

http://admet.scbdd.com/

**Microbial biosynthesis**

*Strep. griseus*

Artemisinin to
Artemisitone-9

Liu *et al.* (2006)

**Activation** — Flutamide to 2-hydroxyflutamide

**Potency Increase** — Diazepam to Oxazepam

**Toxification** — Chlorpyrifos to Chlorpyrifos-oxon

**Detoxification** — Chlorpyrifos to 3,5,6-trichloro-2-pyridinol, and DETP

Djoumbou Feunang *et al.* (2019);

**Nutrition science**

Curcumin/Polyphenols (C/P)

C/P Metabolites ↔ Altered Microbiota

Microbial Products/Metabolites

↓ Adipogenesis
↑ FGF21
↑ Energy expenditure
↑ Antioxidant action

Metabolic Regulation ← ↑ Gut Hormones; GLP-1/2

Jin *et al.* (2017)

Influence of metabolism on a xenobiotic's pharmacodynamics (PD), including pharmacological activity (Act), and toxicological effects (Tox).
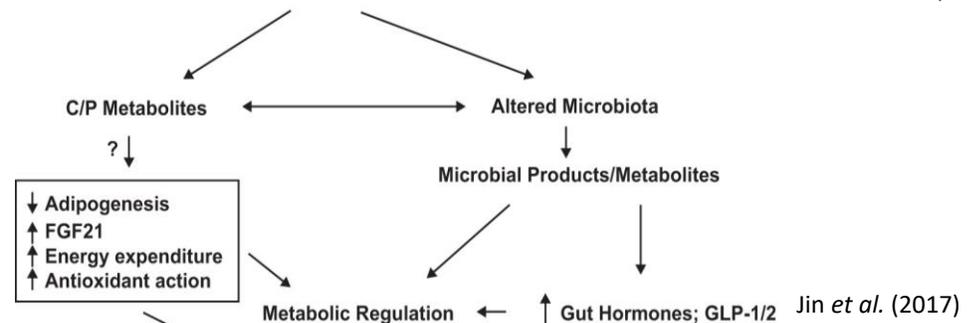
# Approaches for *In silico* Metabolism Prediction
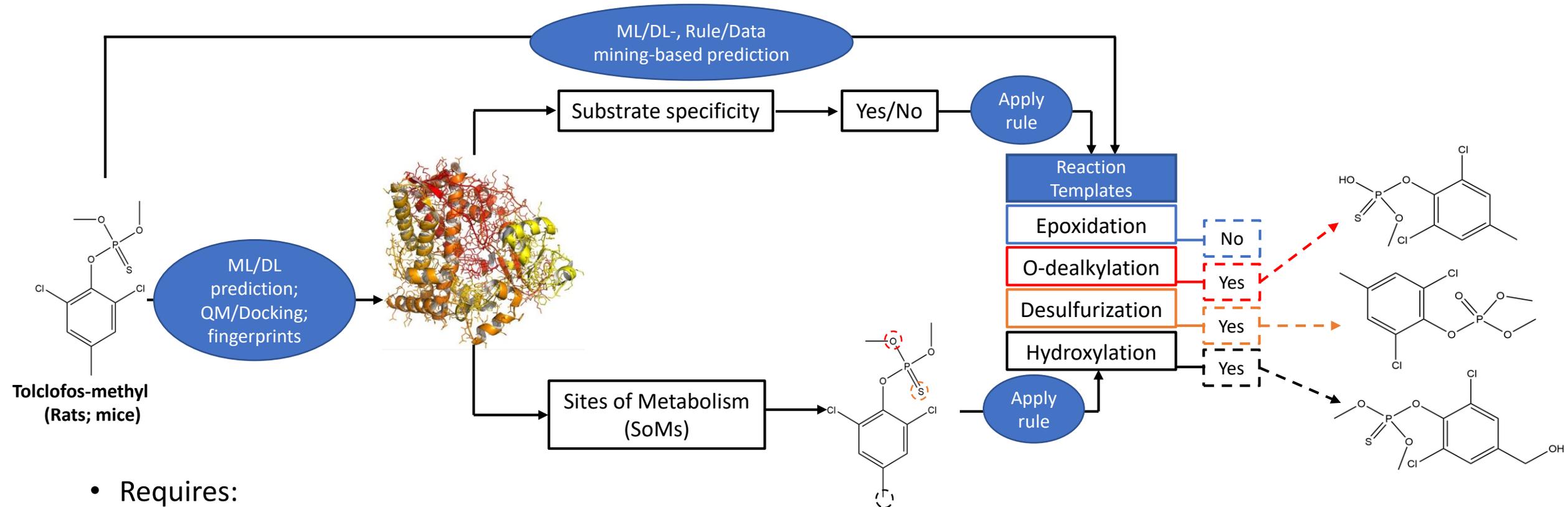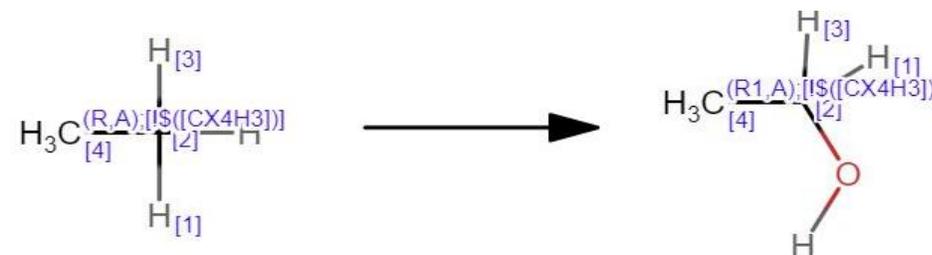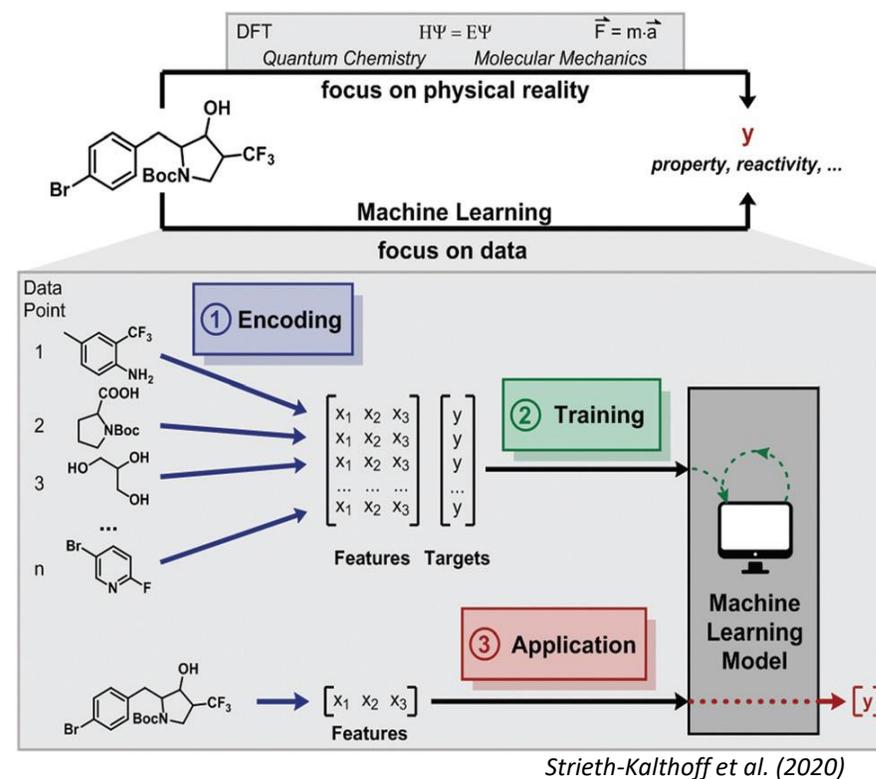


- Requires:
  - Module to predict (and rank/score) SoMs, enzyme-substrate selectivity, or reaction groups
  - Library of reaction templates to apply or select (via prediction) from
  - Modules are usually specific (chemical/enzyme classes), or comprehensive (whole species)
- Prediction approaches can be ligand- or structure-based
- Some tools include: ML/DL-based (MetaTrans, Glory), Rule-based (MetabolExpert), Hybrid (BioTransformer, Meteor Nexus)

# Ligand-based Prediction

- Use structures of ligands/non-ligands to predict specificity, accessibility to enzymes

- Some methods include, among others:
  - QSAR/QSMR (SoM, BoM, ESS)
  - Data mining/fingerprints (SoM)
  - Substructures/Rules (ESS, Reaction)

- Advantages:
  - Speed
  - Seem to perform as well as structure-based methods

- Disadvantages:
  - Approaches/models/rule bases tend to not work on novel chemistries
  - Data quantity and quality is a limiting factor



*Strieth-Kalthoff et al. (2020)*



Hydroxylation of methyl carbon adjacent to aliphatic ring

# Structure-based Prediction

- Explicitly model the interaction within the enzyme's binding pocket

- Help predicting sites of metabolism (SoMs)

- Some methods include, among others:
  - Protein-ligand docking
  - Molecular dynamics

- Advantages
  - Detailed study of the enzyme-substrate interaction

- Disadvantages
  - High computational power to model structural flexibility

Ménard *et al.* (2012)

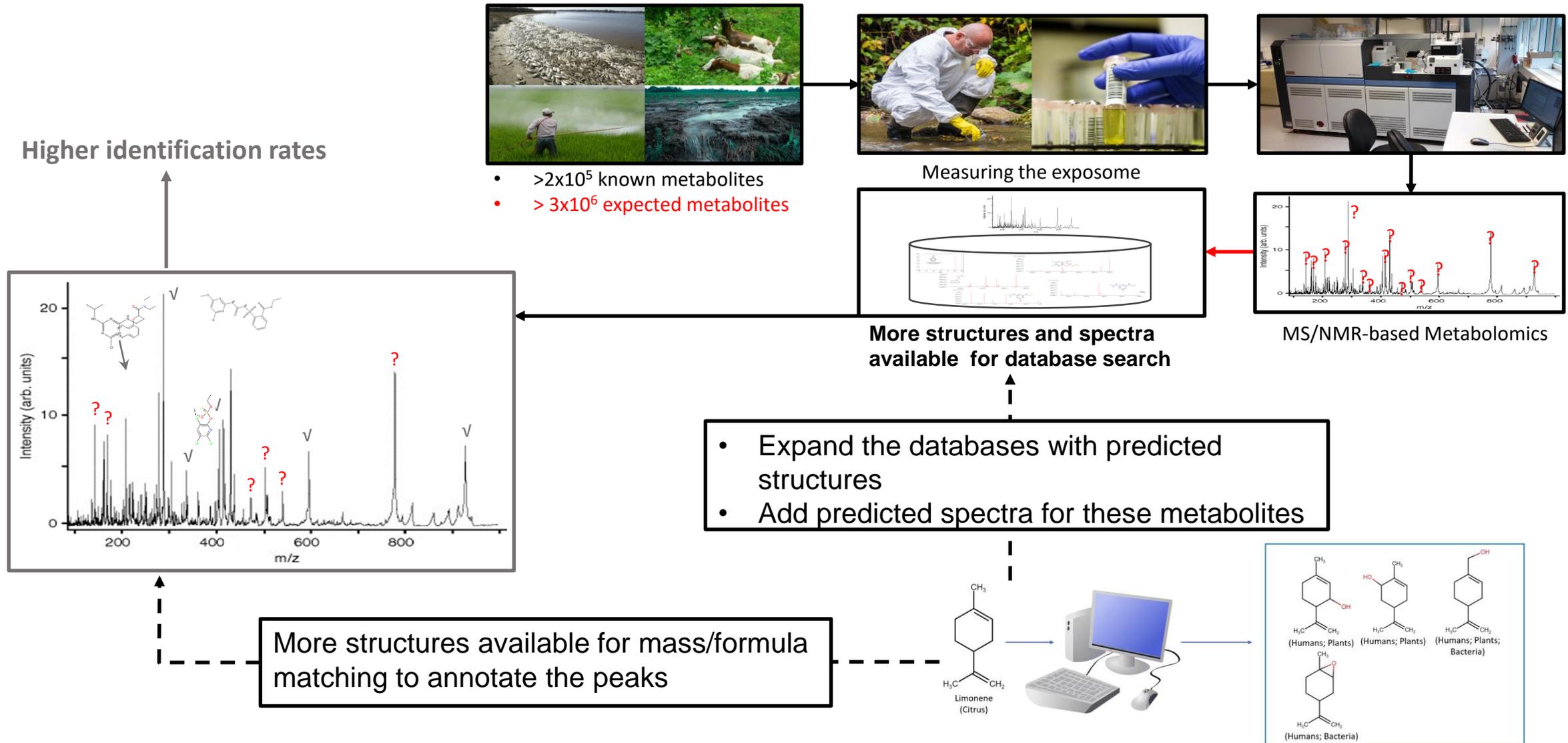# Examples of *In Silico* Metabolism Prediction Tools

Based on their coverage, *in silico* metabolism prediction tools can be classified as:

1.  Specific tools, which apply to simple biological systems (e.g. a small set of enzymes), and often, cover a rather small chemical space.

2.  Comprehensive tools, which apply to a more diverse and larger set of enzymes, species, and cover a larger chemical space.

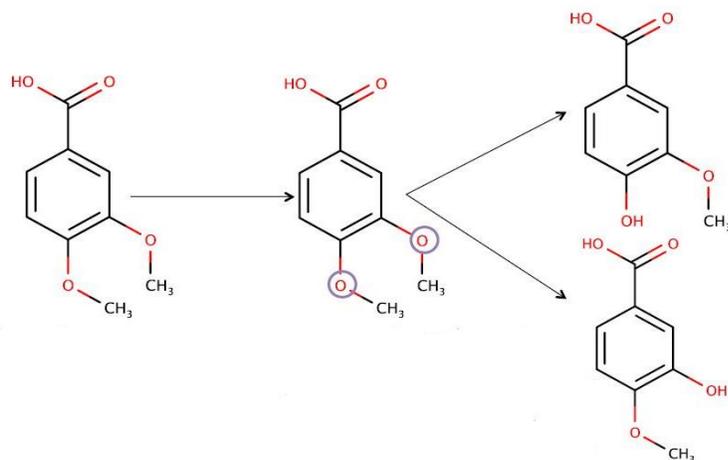| Prediction Tool | Accessibility | Spectrum | Approach | Methods | Availability |
|---|---|---|---|---|---|
| BioTransformer | Web server / Command line | Comprehensive | Ligand-based | Machine Learning Expert system | Open source |
| Meteor Nexus | GUI-based standalone | Comprehensive | Ligand-based | Machine Learning Expert system | Commercial |
| GLORYx | Web server | Comprehensive | Ligand-based | Machine Learning | Freely available |
| SmartCYP | Web server | Specific | Combined | Machine Learning | Freely available |
| MetaSite | GUI-based standalone | Specific | Combined | Machine Learning | Commercial |
| MetabolExpert | GUI-based standalone | Specific | Ligand-based | Expert system | Commercial |

*Examples of tools for metabolite prediction*

# Enhancing Metabolite Discovery and Identification



Measuring the exposome

- >2x10$^5$ known metabolites
- > 3x10$^6$ expected metabolites

MS/NMR-based Metabolomics

**More structures and spectra available for database search**

**Higher identification rates**

- Expand the databases with predicted structures
- Add predicted spectra for these metabolites

More structures available for mass/formula matching to annotate the peaks

Limonene (Citrus)

(Humans; Plants)   (Humans; Plants)   (Humans; Plants; Bacteria)

(Humans; Bacteria)

# BioTransformer: Open Source for Metabolite Identification



Djoumbou Feunang *et al.* (2019); J. Cheminf.; DOI 10.1186/s13321-018-0324-5

- **Djoumbou-Feunang *et al.* J Cheminform. 2019 Jan 5;11(1):2**
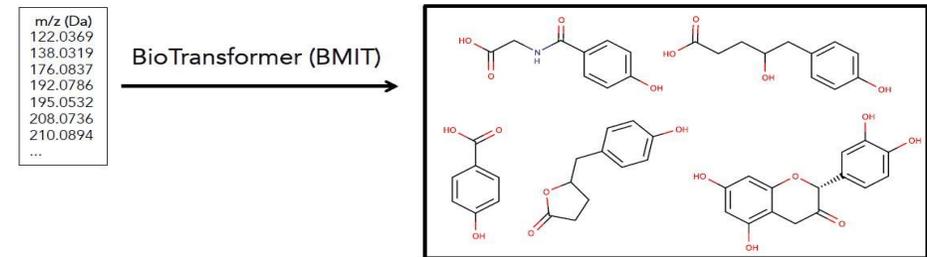- Open Source: bitbucket.org/djoumbou/biotransformer
- Docker: https://hub.docker.com/r/djoy2018/biotransformer-cl
- RESTful web server: www.biotransformer.ca

# BioTransformer: Open source for Metabolite Identification



Epicatechin

12-week old Wistar rats were fed epicatechin

Urine collected before and after 5-day treatment

???

Use BioTransformer to identify/suggest metabolites

Data Extraction (m/z and intensities)
260 neutral monoisotopic masses (53.48 – 969.86 Da)

Sample Analysis with UPHLC - QToF MS

In collaboration with Manach C., and Fiamoncini J., INRA, France

m/z (Da)
122.0369
138.0319
176.0837
192.0786
195.0532
208.0736
210.0894
...

BioTransformer (BMIT)

Metabolite never reported, but mass known

Metabolite and mass never reported before

- BioTransformer (BMIT) suggested 19 new compounds
  – 11 compounds match previously reported masses
  – 15 compounds match masses observed exclusively in our study
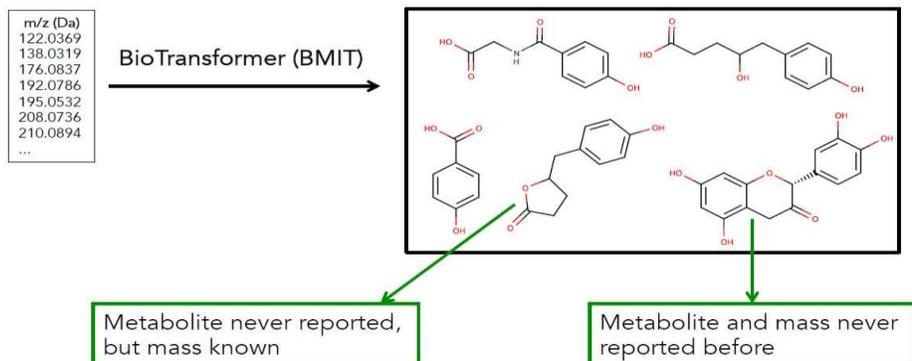  – Experimental validation needed

- Literature mining revealed 56 single- and multi-step epicatechin metabolites corresponding to 37 unique masses
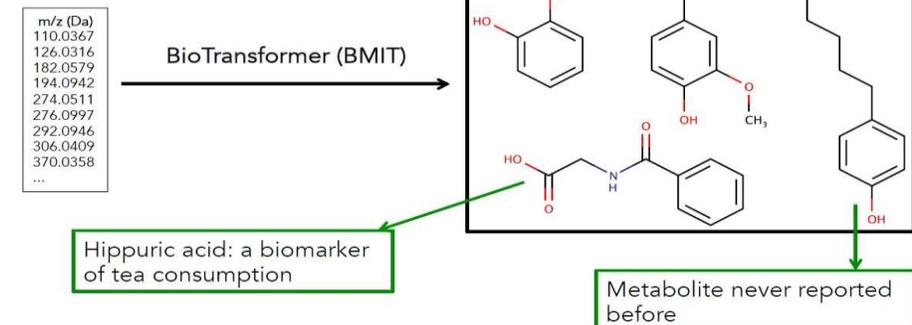  – 11 out of 37 masses were measured in our experimental study

m/z (Da)
122.0369
138.0319
176.0837
192.0786
195.0532
208.0736
210.0894
...

BioTransformer (BMIT)

- BioTransformer (BMIT) suggested 37 compounds (20 masses)
  – 22/37 matched the 11 unique and previously reported masses
  – 18/22 compounds were confirmed in previous reports

- BioTransformer (BMIT) was used to identify metabolites matching the 26 masses not measured in our study

m/z (Da)
110.0367
126.0316
182.0579
194.0942
274.0511
276.0997
292.0946
306.0409
370.0358
...

BioTransformer (BMIT)

Hippuric acid: a biomarker of tea consumption

Metabolite never reported before

- BioTransformer (BMIT) suggested 28 compounds (19 masses)
  – 21 previously reported metabolites matching 18 unique masses
  – Experimental validation needed for the 7 new compounds

# Challenges in Metabolism Prediction

- Challenges in understanding metabolism lead to lower prediction accuracy
  - ➤ ML tools can rank SoMs with high accuracy, but not the resulting biotransformations
  - ➤ Knowledge-driven approaches achieve higher recall, but suffer from combinatorial explosion

- Limited accessibility/availability of tools
  - ➤ Most tools developed in academia are available only via web-sever, raising confidentiality issues for most potential users
  - ➤ Restricted shareability of data generated by commercial tools

- The limited availability of high-quality data
  - ➤ Crucial for achieving better accuracy, higher coverage, and larger applicability domains
  - ➤ Needed to derive more reaction rules
  - ➤ However, data curation (Sites of metabolism, reaction annotation) is tedious and expensive

# Outlook

- Data sharing
  - ➢ Increase the amount of high-quality, publicly available, curated, and downloadable data (substrate-product, Sites of Metabolism). E.g.: MetXBioDB, PubChem, XMetDB, etc.
  - ➢ Develop publicly accessible databases of predicted metabolites (e.g.: BioTransformerDB)
  - ➢ Community-wide efforts needed

- Novel, and innovative AI approaches:
  - ➢ Seq2Seq transformer architectures have proven applicable for end-to-end learning-based method, with results comparable to existing tools (MetaTrans: Litsa *et al.* 2020)
    - ▪ Could improve accuracy, while bypassing manual rule design
  - ➢ Bonds of Metabolism (BoMs) seem to provide a better description of reaction centers, leading to higher accuracy, compared to SoMs (Upcoming - Tian, et al. (2021))

- Open source communities
  - ➢ Provide means for easier user feedback loops; valuable for improving prediction tools
  - ➢ Provide developers opportunities to improve software tools, in an agile, continuous manner
    - ▪ Successful projects include (BioTransformer, Chemistry Development Kit, Knime)

# Thank You

- The organizing committee:
  - Fidele Ntie-Kang
  - J. Ludwig Muller

- The Wishart Lab @UofAlberta, Canada

- Corteva Agriscience

- The BioTransformer community

- The listeners


- To learn more about BioTransformer:
  - *HS03: In Silico Prediction and Identification of Metabolites with BioTransformer : Enabling Secondary Metabolite Discovery; March 10, 2021*

https://unlockinglifescode.org/education-resource-profile/resource-month-ask-biologist