# Computational Tools for the Identification of Unknowns

**David Wishart, University of Alberta**

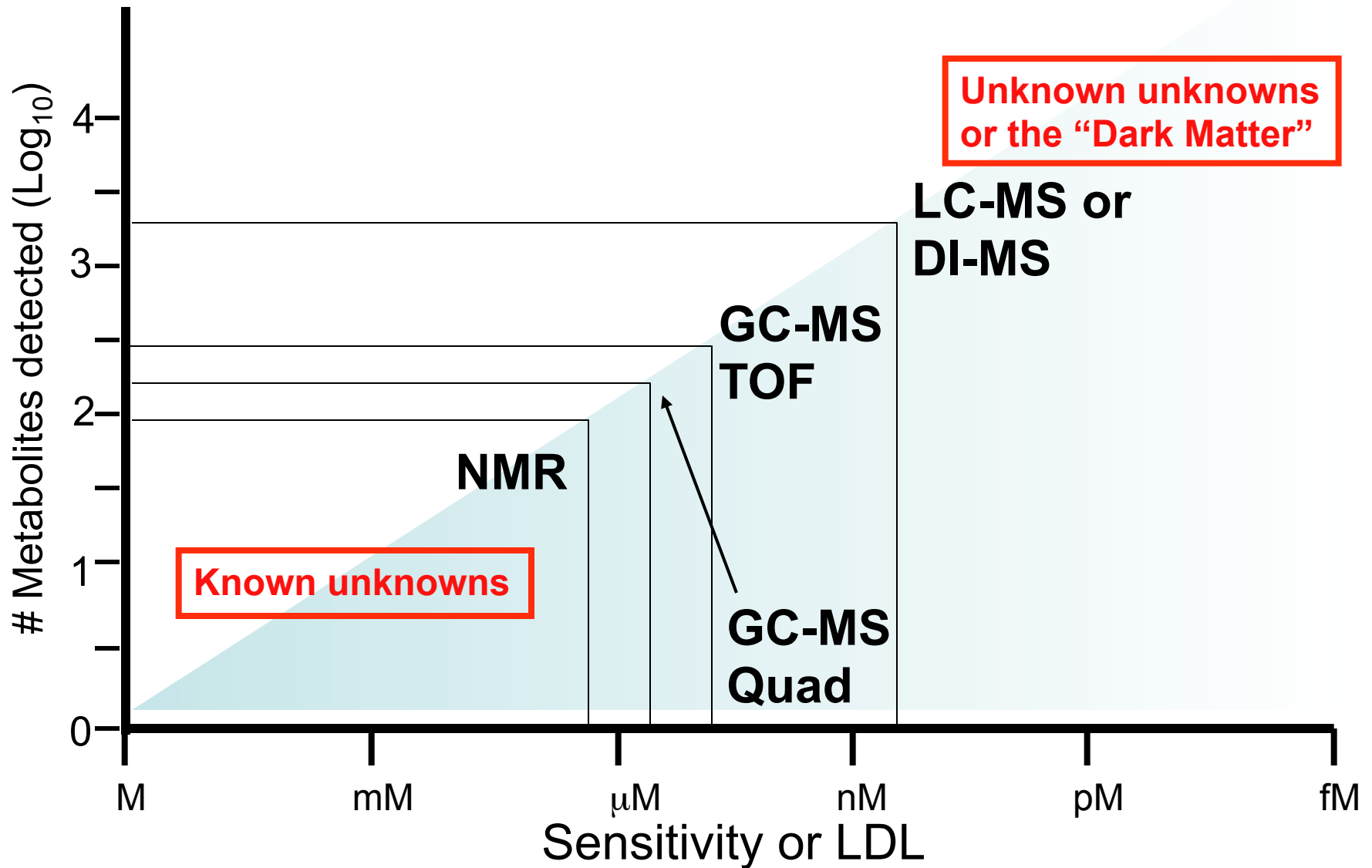**3rd International Electronic Conference on Metabolomics**

**Nov. 15-30, 2018**

# What We Don't Know

- "*…there are known unknowns; that is to say we know there are some things we do not know. But there are also unknown unknowns – the ones we don't know we don't know.*"
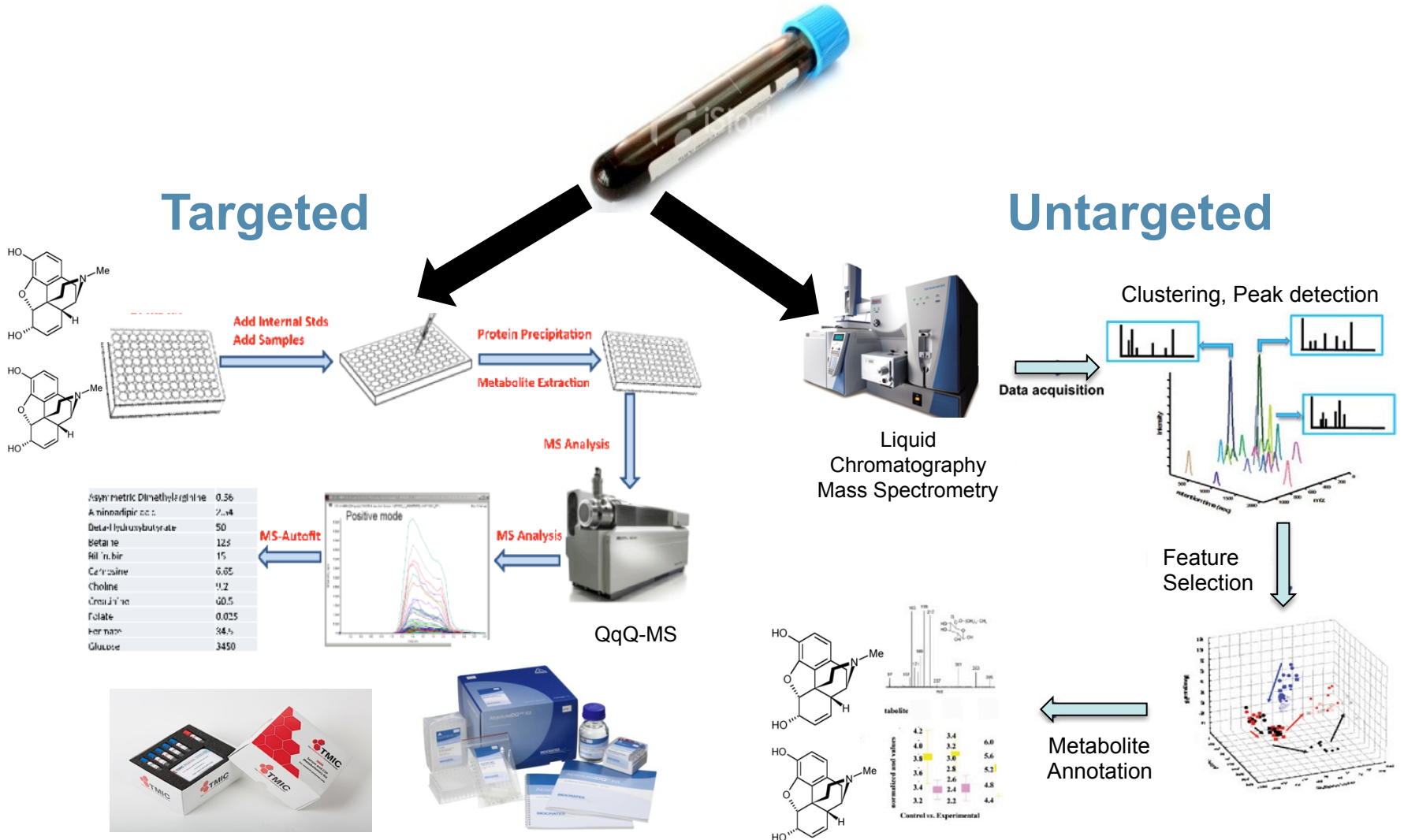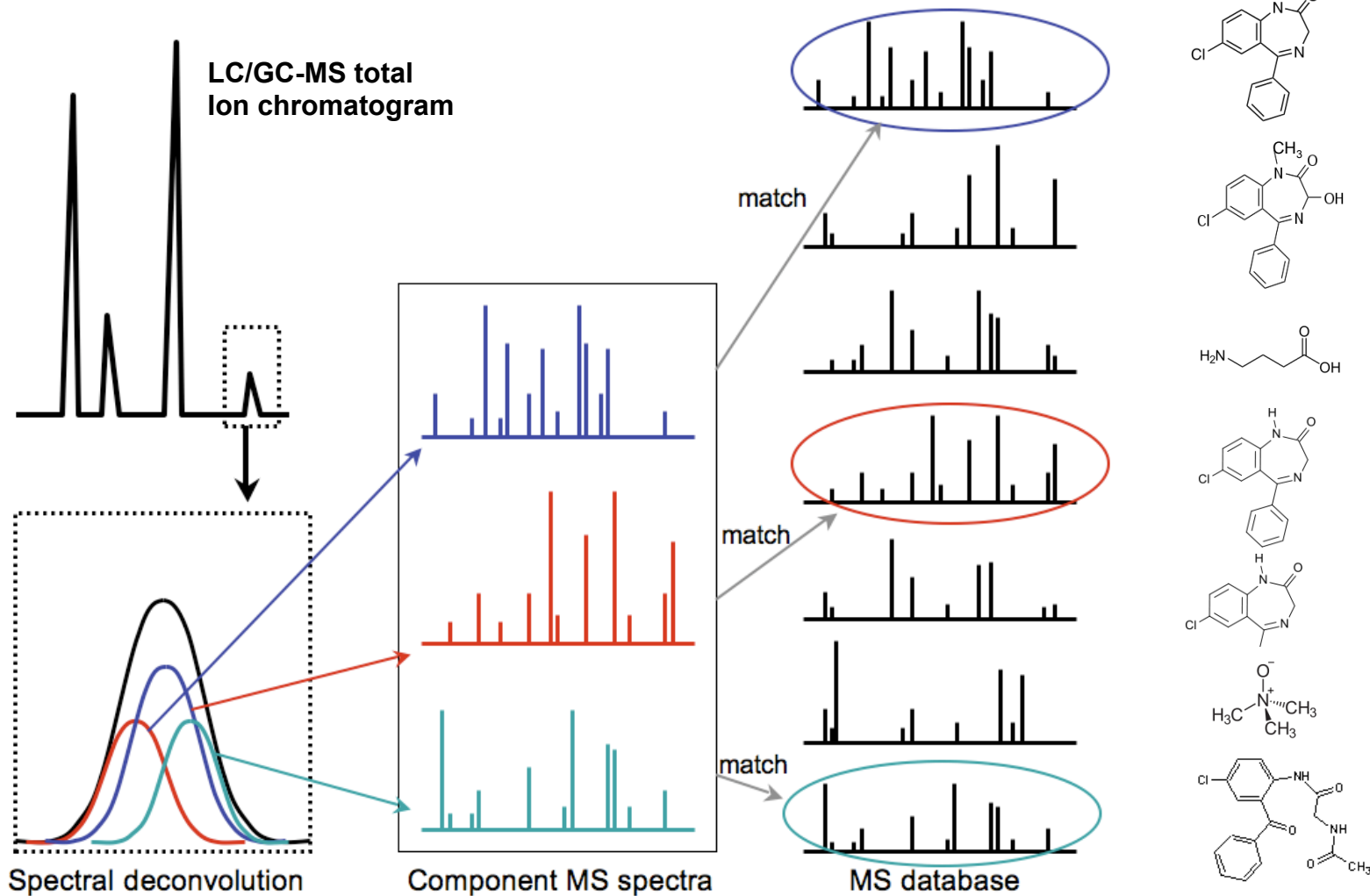


Donald Rumsfeld, US Secretary of Defense - Nov. 2001

# Technology & Sensitivity

# 2 Routes to Metabolomics

**Targeted**

**Untargeted**

Add Internal Stds
Add Samples

Protein Precipitation

Metabolite Extraction

MS Analysis

Liquid
Chromatography
Mass Spectrometry

Data acquisition

Clustering, Peak detection

| Asymmetric Dimethylarginine | 0.36 |
| Aminoadipic acid | 2.14 |
| Beta-hydroxybutyrate | 50 |
| Betaine | 125 |
| Bilirubin | 15 |
| Carnosine | 6.65 |
| Choline | 9.2 |
| Creatinine | 60.5 |
| Folate | 0.025 |
| Formate | 34.5 |
| Glucose | 3450 |

MS-Autofit

Positive mode

MS Analysis

QqQ-MS

Feature
Selection

Metabolite
Annotation

# Untargeted MS Compound Identification



**LC/GC-MS total Ion chromatogram**

match

match

match

Spectral deconvolution

Component MS spectra

MS database

# Levels of Metabolite ID for Untargeted Metabolomics

- **4 levels of metabolite identification**

- **Level 1 - Positively identified compounds**
  - Confirmed by MS/MS match and RT match to an actual/authentic standard

- **Level 2 - Putatively identified compounds**
  - Match to EI-MS + RT or MS/MS + RT from a reference database

- **Level 3 - Compounds putatively identified via molecular formula or m/z matching**
  - Match to high resolution m/z and nothing else

- **Level 4 - Unknown compounds**

# Compound Identification (Formula Matching – Level 3 ID)



**68 million chemicals**

**96 million chemicals**

# Compound Identification (Spectral Matching – Level 2 ID)



**14,009 "real" compounds**
**72,036 "real" spectra**
**~150,000 total compounds**

**213,019 MS spectra**
**75,270 compounds**

# Other Resources for MS/MS Spectral Matching

- **HMDB** – 302,219 spectra, 114,100 cmpds

- **mzCloud** – 191,722 spectra, 8304 cmpds

- **NIST17 MS/MS** – 652,475 spectra, 14,351 cmpds

- **MassBank** – 28,185 spectra, 11,500 cmpds

- **Wiley LC-MS$^n$** – >10,000 spectra, 4500 poisons

- **ReSpect** – 9017 spectra, 3595 cmpds

- **GNPS** – 221,000 spectra, 18,163 cmpds

# How Well Do We Do?

- **Untargeted LC-MS of human biofluids** – 100-250 compounds positively ID'd (level 2) out of 10,000+ features (1%), 700-1000 tentatively ID'd (level 3) out of 10,000+ features (8%)

- **Untargeted LC-MS of river water** – 649 compounds identified (level 1 or 2) out of 8535 features (8%) *(Schymanski et al. Anal Bioanal Chem (2015) 407:6237–6255)*

*Overall we are doing pretty badly*
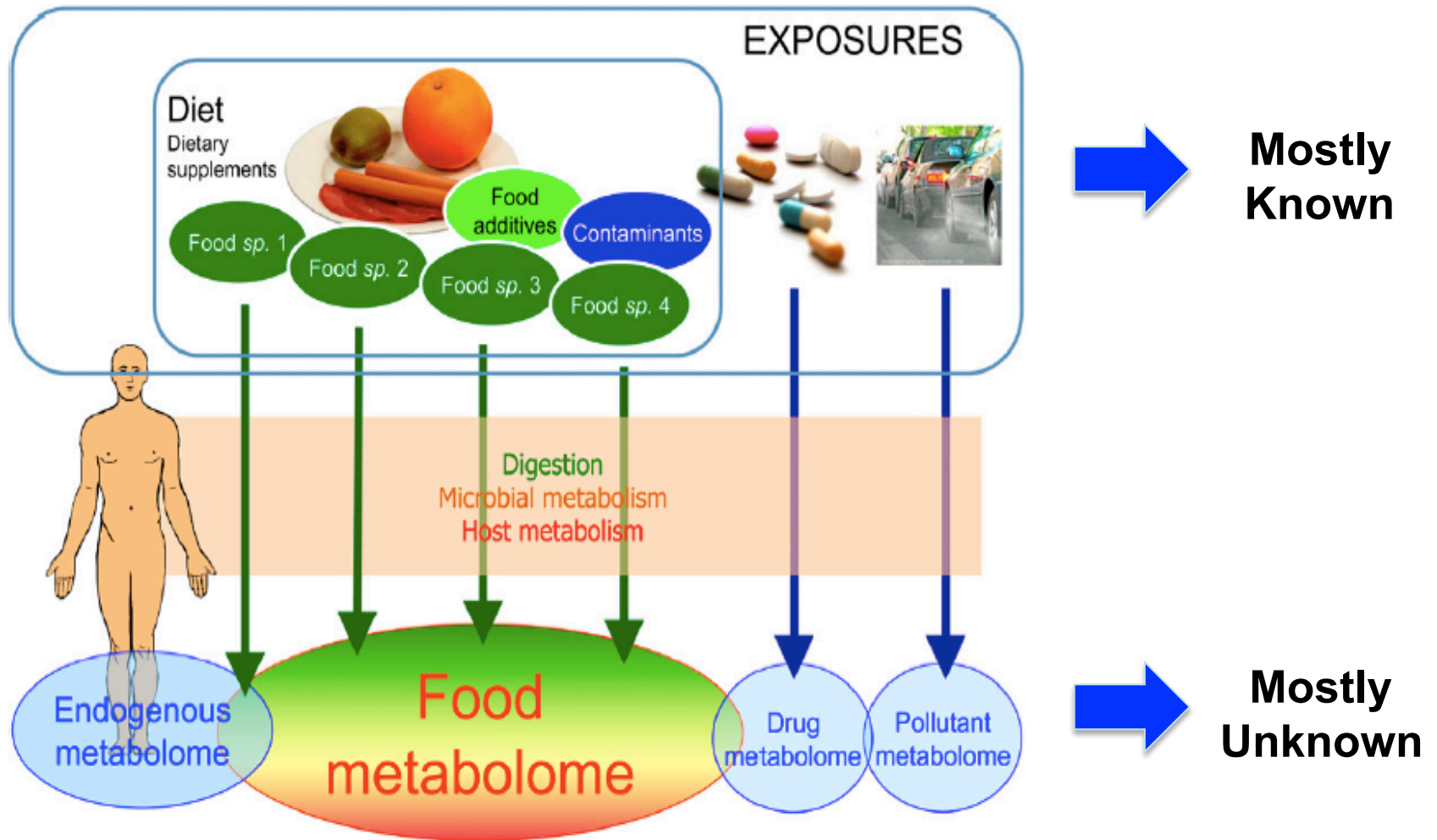
# Why Are We Doing So Badly?

- **32,000,000,000** chemical formulae (<2000 Da)

- **2,400,000,000** chemically feasible formulae

- **96,500,000** chemicals in PubChem

- **1,500,000** LC/GC-MS spectra collected on ~15 different platforms

- **80,000** chemicals with EI-MS spectra

- **~20,000** chemicals with high resolution MS spectra

- **~1500** chemicals that are biologically relevant

# Why Are We Doing So Badly?

- **Using larger databases (PubChem, ChemSpider) and m/z matching is leading to many, many false positives**

- **<0.2% of compounds in PubChem or ChemSpider have ever left the laboratory or are likely to be found in humans or in the environment (Bigger isn't better)**

- **Only 1500 "meaningful" chemicals are routinely available from vendors (Synthetic chemistry is hard)**

- **Enormous resources going to collect lots of MS spectra on a tiny number of chemicals (Measuring the same thing over doesn't make the problem go away)**

- **Most of the unknowns are not in PubChem or Chemspider, or anywhere else (Where are they?)**

# What Are These Unknowns?

# What To Do?

- **Obtain or synthesize all commonly available xenobiotics (HPVs, drugs, pollutants, foods, etc.), prepare or synthesize their metabolites and collect their NMR, LC-MS and GC-MS spectra COST: 5 million cmpds X $1000/cmpd = $5 billion**

- **OR**

- **Do this entire exercise computationally**

  **COST - 5,000,000 cmpds X $0.10/cmpd = $500,000**

# *In Silico* Metabolomics

- **An emerging concept to facilitate metabolite ID of unknown unknowns**

- **Realization that all metabolites will never be synthesized or isolated and most will never have reference MS/MS or NMR spectra collected**

- **Based on *in silico* prediction of biologically feasible metabolites**

- **Based on *in silico* prediction of observables (RI, RT, NMR spectra, IR, CCS, MW, MS/MS spectra)**

# *In Silico* Metabolomics

HMDB/DrugBank/T3DB

BioTransformer



Known compounds (250,000)

Predicted biotransformations
(250,000 --> 5,000,000)

CFM-ID/NMRPred

Match observed spectra
to predicted spectra to ID compounds

Predicted MS/MS, NMR, GC-MS
spectra of knowns + biotransformed

# What Is Known?

# UofA Metabolomics Databases



www.hmdb.ca

www.drugbank.ca

www.ymdb.ca

www.phenol-explorer.eu

www.ecmdb.ca

www.foodb.ca

www.cowmetdb.ca

www.t3db.ca

www.smpdb.ca

www.csfmetabolome.ca

www.serummetabolome.ca

www.urinemetabolome.ca

# The New Human Metabolome Database (HMDB)



**http://www.hmdb.ca**

- HMDB 4.0 has 114,100 "quantified", "detected", "expected" and "predicted" metabolites (3X more than version 3.0)

- HMDB 3.0 had 442 biological pathways, HMDB 4.1 has 48,627 (100X more than version 3.0)

- New version has >500,000 MS/MS & GC-MS spectra, 3900 NMR spectra

- New version has 6800 metabolite-SNP associations, 2500 metabolite-drug associations and 2900 metabolite-age/gender associations

- 78,000 new lipids/peptides to be added in late 2018 – total = 192,000

# The **New** Drug Database (DrugBank v. 5.0)



- A comprehensive database of drugs, drug actions and drug targets
- 2533 small molecule drugs
- >5700 experimental drugs
- Detailed ADMET, MOA and pharmacokinetic data
- >3850 drugs with metabolizing enzyme data
- >1360 drug metabolites
- >6000 MS+NMR spectra
- >5130 unique drug targets
- 215 data fields/drug
- *Published on Jan. 1, 2018*

**http://www.drugbank.ca**

# The Food Database (FooDB)



www.foodb.ca

- **A comprehensive food composition database (more than polyphenols)**
- **28,771 compounds**
- **718,405 concentration values for 722 raw/processed foods**
- **31,791 references**
- **1435 cmpds with health effects**
- **2692 cmpds w flavour attributes**
- **2000+ reference MS/NMR spectra**
- **Structure & text searches**
- **>100 data fields/compound**
- ***Publicly released on Jan. 1, 2018, manuscript being prepared***

# The Toxic Exposome Database (T3DB)



**http://www.t3db.ca**

- **Comprehensive data on toxic compounds (drugs, pesticides, herbicides, endocrine disruptors, drugs, solvents, carcinogens, etc.)**

- **Detailed mechanisms, binding constants, target info, lots of ToxCast data**

- **>3600 toxic compounds**

- **>1900 reference spectra**

- **~2100 toxic targets**

- **Supports sequence, spectral, structure, text searches as well as compound browsing**

- **Full data downloads**

# ContaminantDB



www.contaminantdb.ca

- **Data on 54,249 probable or known chemical contaminants**
- **Expected to grow to 80,000+ by Sept. 2018**
- **Exp. MS data for 5000+ cmpds**
- **Pred. 54,000 EI-MS spectra, 150,000 ESI-MS/MS spectra**
- **Source or role information for most compounds**
- **>40% of the compounds in ContaminantDB are not found in PubChem or ChemSpider**
- **Supports spectral, structure and text searches as well as compound browsing**

# PhytoBank



- **179,729 plant-derived compounds from more than 23,700 plant species including >8,318 food/crop plants and >2,439 medicinal plants**

- **>33% of the compounds in PhytoBank are <u>not</u> found in PubChem or ChemSpider**

- **Will offer same resources as HMDB, DrugBank, etc.**

# What Can We Predict?

# *In Silico* Metabolomics

HMDB/DrugBank/T3DB

BioTransformer



Unknown Unknowns

Known compounds (250,000)

Predicted biotransformations
(250,000 --> 5,000,000)

CFM-ID/NMRPred

Match observed spectra
to predicted spectra to ID compounds

Predicted MS/MS, NMR, GC-MS
Spectra of knowns + biotransformed

# Examples of Biotransformation



Tempazepam

Oxazepam

Nordazepam

Diazepam

N-(2-Benzoyl-4-chlorophenyl)-2-acetamidoacetamide

# Commercial Tools



Meteor-Nexus



ADMET Predictor



Metabolizer



MetabolExpert

# BioTransformer (Free)



**Input**
Prediction: Structure (SMILES, MOL, SDF)
Identification: Structure, masses and mass tolerance

(1) Parsing and first filtering

(2) Standardization

(3) Prediction and structure generation

(4) Metabolic tree reconstruction and metabolite annotation

**Output in SDF**
Prediction: InChI, ID, InChI key, mass, LogP, formula, biosystem, enzyme, reaction type, synonyms, PubChem CID.
Identification: InChI, InChI key, mass, synonyms, PubChem CID, Metabolic trees

Reasoning Engine

Knowledgebase System

Human Super Transformer

EC-based

CYP450

Phase II

Human Gut Microbial

Envir. Microbial

CypPred

Machine Learning-Based Module

- **A comprehensive tool for predicting metabolite structures that have been biotransformed through phase I, phase II, microbial, promiscuous and environmental processes**

- **Uses a large knowledgebase and a large set of heuristic (and machine-learned) biotransformation rules**

- **Performs much better than well-known commercial tools**

- *Publicly released, manuscript submitted*

# BioTransformer

**Parathion**

Phase I (163)

Env. Microbial (301)

**Glycitein**

Human Gut (201)

Phase II (74)

**DG(16:0/16:0/0:0)**

EC-based (408)

# BioTransformer

**Oxidation of p-substituted anilides (BTMR1018)**
**Human (Liver) Phase I Metabolism**

**Hydrolysis of secondary amide (BTMR0704)**
**Environmental microbial (EAWAG-BBD)**

| | No. of enzymes | No. of biotransformation rules | No. of enzyme-rule associations | No. of covered biosystems |
|---|---|---|---|---|
| **EC-based** | 285 | 408 | 459 | 2 |
| **CYP450** | 9 | 163 | 712 | 1 |
| **Human gut micro.** | 53 | 201 | 204 | 2 |
| **Phase II** | 9 | 74 | 81 | 2 |
| **Envir. microbial** | 1 | 301 | 301 | 1 |

**Overall, 1,240 biotransformation rules were manually designed and tested. Overall, 2,150 reaction preference rules were implemented.**

**Quercetin**

FLAVONE_REDUCTION TO_FLAVANONE

**Taxifolin**

FLAVANONE_C_RING_FISSION

**Hydrokampferol chalcone**

TAXIFOLIN_C_RING_CONTRACTION

CHALCONE_DIHYDROGENATION_PATTERN1

ARYL_O_METHYLATION

**3'-methoxyquercetin**

**Aphitonin**

ALPHITONIN_DEGRADATION

ARYL_OH_GLUCURONIDATION

**4'-methoxyquercetin**

FLAVONOL_3_O_GLUCOSYLATION

**Quercetin 3-O Glucuronide**

**Phoroglucinol**

**3,4-dihydroxyphenylacetic acid**

**Quercetin 3-O Glucoside**

# Quercertin BioTransformartion

# BioTransformer Evaluation

| | BioTransformer | Meteor |
|---|---|---|
| **True Positives** | 188 | 153 |
| **False Positives** | 198 | 281 |
| **False Negatives** | 35 | 70 |
| **Total No. of Predictions** | 386 | 431 |
| **Precision** | 0.49 | 0.35 |
| **Recall** | 0.84 | 0.69 |
| **# of reported metabolites** | 223 | |

1. Test set: 40 compounds (incl. drugs and pesticides)
2. Metabolism data was retrieved from >60 references
3. BioTransformer (v. 1.0.4) and Meteor Nexus (v. 3.0.1) were used for 1- step prediction of **human metabolism**

# BioTransformer Evaluation

1. Test set: 60 compounds (incl. drugs, pesticides, food compounds, steroids)
2. Metabolism data was retrieved from 60+ references
3. BioTransformer (v. 1.0.4) and ADMET Predictor (v. 8.5.11) were used for 1- step prediction of **human CYP450 metabolism**

|  | BioTransformer | ADMET Predictor |
|---|---|---|
| True Positives | 162 | 110 |
| False Positives | 188 | 122 |
| False Negatives | 18 | 70 |
| Total No. Predictions | 350 | 232 |
| Precision | 0.46 | 0.47 |
| Recall | 0.9 | 0.61 |
| No. of Reported Metabolites | 180 | |

1. Test set: 20 compounds (incl. endogenous metabolites, phytochemicals, and other xenobiotics)
2. Metabolism data was retrieved from >50 references
3. BioTransformer (v. 1.0.4) was used for 1-step prediction of **human and gut microbial metabolism**

| BioTransformer | |
|---|---|
| True Positives | 111 |
| False Positives | 49 |
| False Negatives | 17 |
| Total No. Predictions | 160 |
| Precision | 0.69 |
| Recall | 0.87 |
| No. of Reported Metabolites | 128 |

# BioTransformer Updates

- **Added/updated:**
  - Available as a web server and a program, added metabolite ID option via m/z or molecular formula

- **Number of compounds expected**
  - 1-step reactions generate 5X as many compounds while 2-step reactions generate 20X as many compounds
  - Expect HMDB+BioTransformer will generate 2.2 million new compounds, all DBs+BioTransformer = 5 million cpds

- **Benchmark for computing time**
  - 1,000 FooDB compounds generate 5,071 human and gut microbial metabolites in 1 step (all enzymes)
  - 1h 29 mins (~5.35 s/compound per processor)

If you want compounds processed now, send your queries to
Yannick Djoumbou Feunang --- djoumbou@ualberta.ca

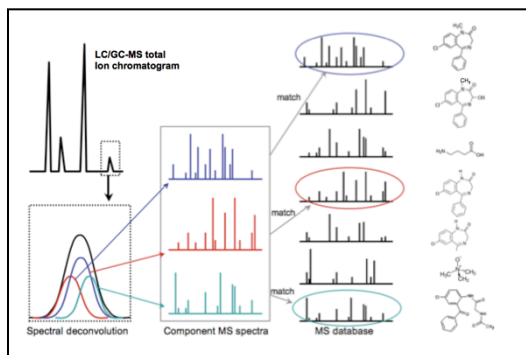# *In Silico* Metabolomics

HMDB/DrugBank/T3DB

BioTransformer



The Dark Matter

Known compounds (250,000)

Predicted biotransformations
(250,000 --> 2,500,000)

CFM-ID/NMRPred

Match observed spectra
to predicted spectra to ID compounds

Predicted MS/MS, NMR, GC-MS
Spectra of knowns + biotransformed

# Competitive Fragment Modeling (CFM-ID)



F. Allen et al. Metabolomics 11: 98-110 (2015).

# CFM-ID

- **Uses a large training set of high resolution MS/MS data of known compounds at low (10 eV) medium (20 eV) and high (40 eV) collision energies**

- **Employs an initially naïve chemical fragmenter that generates potential fragments and the coresponding MS/MS spectra**

- **The fragmenter slowly learns from its training data (via an HMM)**

- *The more training data, the better the overall performance*

# CFM-ID Performance



F. Allen et al. Metabolomics 11: 98-110 (2015).

# Performance

- **Significant performance improvement in CFM-ID and all other fragment or structure predictors if the database being searched is smaller or more targeted**

- **Significant improvement if multiple collision energies (10, 20, 40 eV) are used rather than a single collision energy**

- **80% correct for DB ~30,000 cmpds**

- **50% correct for DB ~1,000,000 cmpds**

- **20% correct for DB ~50,000,000 cmpds**

# The CFM-ID Server



http://cfmid.wishartlab.com

- **A web server that predicts MS/MS spectra, annotates input MS/MS spectra and permits compound identification from input MS/MS spectra**
- **Matches predicted MS/MS spectra (from HMDB or KEGG) to input MS/MS spectra**
- **1st and 2nd in the 2014 CASMI competition, used by winners of 2016 CASMI competition**

# CFM-ID Example Output

# CFM-ID Updates

- **Version 3.0 completed (to be released in Dec):**
  - Significantly improved (5X) performance with respect to lipid MS/MS spectral prediction
  - Supports matches to known MS/MS spectra for better compound ID, improved scoring based on citations

- **Version 4.0 in progress (to be finished by 2019)**
  - Much larger training data set (4X larger) covering QTOF and OrbiTrap MS/MS spectra at multiple collision energies
  - Improved generative function, improved chemical and bond descriptors boosts spectral prediction performance by 30-40% over previous CFM-ID version
  - Combined improvements should increase overall performance by at least 50% (still not perfect)

# CFM-ID 3.0 for Predicting 70,000 Lipid Spectra in the HMDB

# "Observables" Prediction

- **MS/MS spectral prediction alone will not be sufficient to ID all unknown compounds**

- **Other observables need to be included for confirmation such as RT (retention time), RI (retention index), CCS (collisional cross section), and gas phase IR or IR ion spectra**

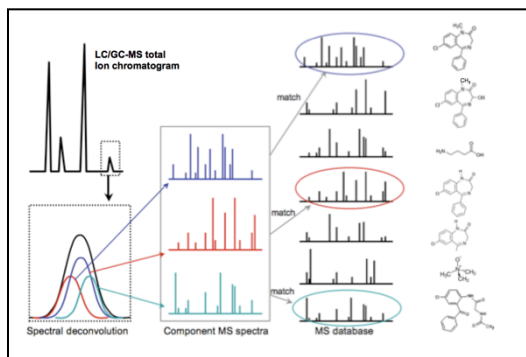Martens et al., *J. Inherit Metabol. Dis.* 2018 41(3):367-377

# *In Silico* Metabolomics

HMDB/DrugBank/T3DB

BioTransformer



Known compounds (250,000)

Predicted biotransformations
(250,000 --> 5,000,000)

CFM-ID/NMRPred
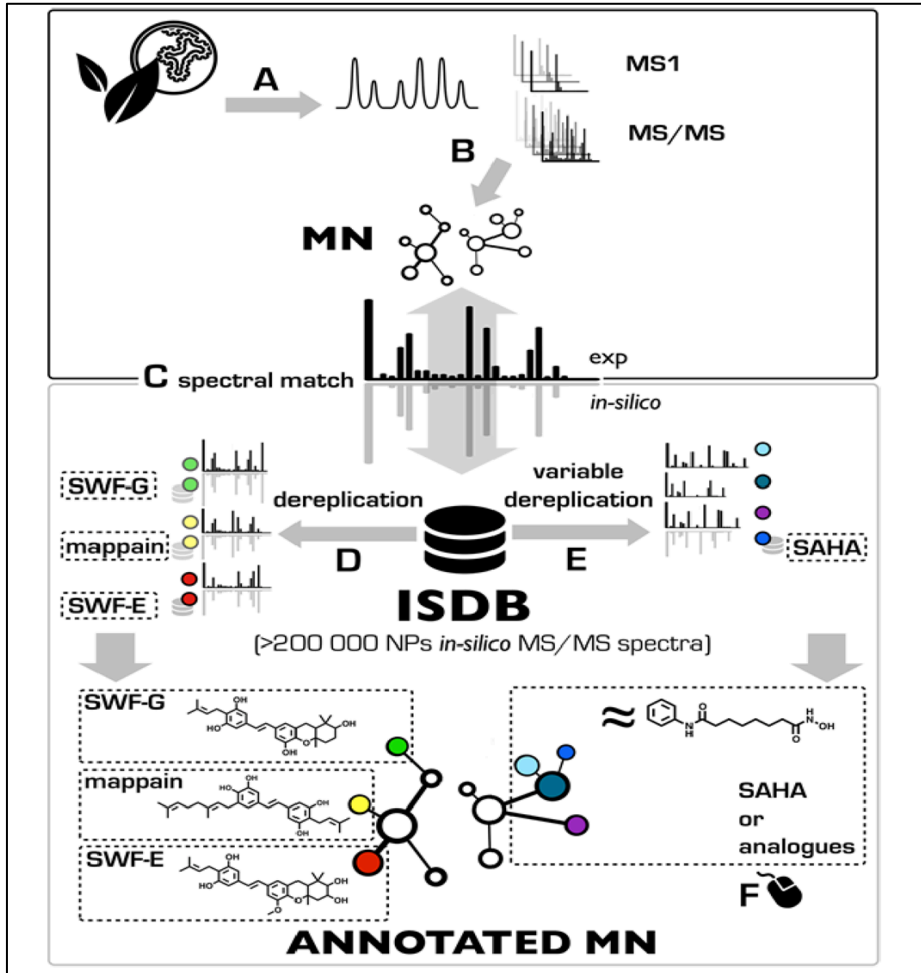
Match observed spectra
to predicted spectra to ID compounds

Predicted MS/MS, NMR, GC-MS
Spectra of knowns + biotransformed

# Who Is Doing This?

PM Allard et al. *Anal. Chem*. 2016 88(6) 3317-3323

F Qiu et al. *Anal. Chem*. 2016 88(23) 11373-11383

# Conclusions

- **Compound identification of unknown compounds by MS/MS analysis is hard**

- **We have insufficient MS/MS and compound resources now and the foreseeable future**

- ***In silico* methods offer a possible solution to the problem of inadequate spectral libraries and inadequate collections of compounds**

- **Predicted compound libraries and and predicted MS/MS spectra are still imperfect, but they are getting better every year**

- **These *in silico* methods are already being used by several groups in natural product analysis**

# Thanks To…

- **Yannick Djoumbou Feunang**
- **Ana Marcu**
- **AnChi Guo**
- **Kevin Liang**
- **Rosa Vazquez-Fresno**
- **Tanvir Sajed**
- **Daniel Johnson**
- **Carin Li**
- **Naama Karu**
- **Zinat Sayeeda**
- **Elvis Lo**
- **Nazanin Assempour**
- **Augustin Scalbert**

- **Sandeep Singhal**
- **David Arndt**
- **Yongjie Liang**
- **Hasan Badran**
- **Jason Grant**
- **Arnau Serra-Cayuela**
- **Yifeng Liu**
- **Rupasri Mandal**
- **Vaness Neveu**
- **Allison Pon**
- **Craig Knox**
- **Mike Wilson**
- **Claudine Manach**