

Are Self-driving Labs for Chemistry Next?

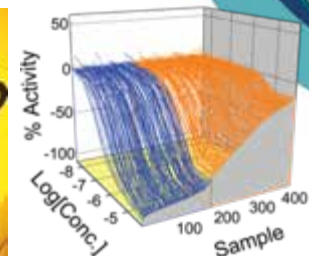
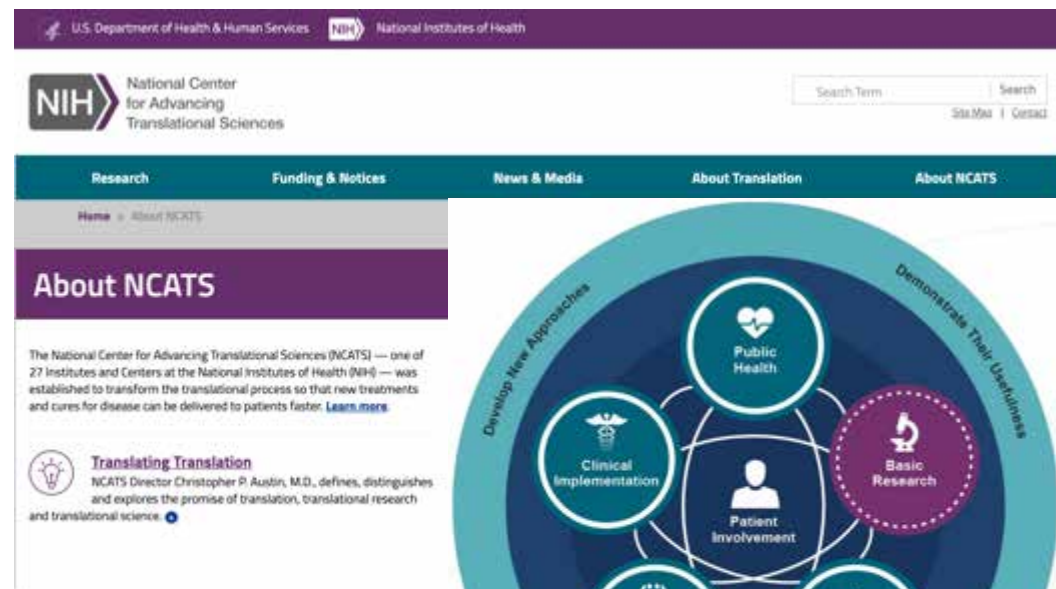
Alex Godfrey
NAS Space Science Week
CBPSS Meeting
March 31, 2020

National Center for Advancing Translational Sciences



National Center for Advancing Translational Sciences

- Founded 2004 as part of NIH Molecular Libraries Roadmap Initiative
- NIH Chemical Genomics Center (NCGC)
- MLPCN (screening & chemical synthesis; compound repository; PubChem database; funding for assay, library and technology development)
 - Develop new chemical probes for basic research and leads for therapeutic development, particularly for rare/neglected diseases
 - New paradigms & applications of HTS for chemical biology / chemical genomics
- Incorporated into NCATS in January 2012



ncats.nih.gov

NCATS Mission



To catalyze the generation of **innovative methods and technologies** that will enhance the development, testing and implementation of diagnostics and therapeutics across a wide range of human diseases and conditions.

ASPIRE

A Specialized Platform for Innovative Research Exploration



**Re-thinking and Re-designing the
Research Laboratory Ecosystem**

Modern
Lab Bench



Integrated
Automated
Solutions

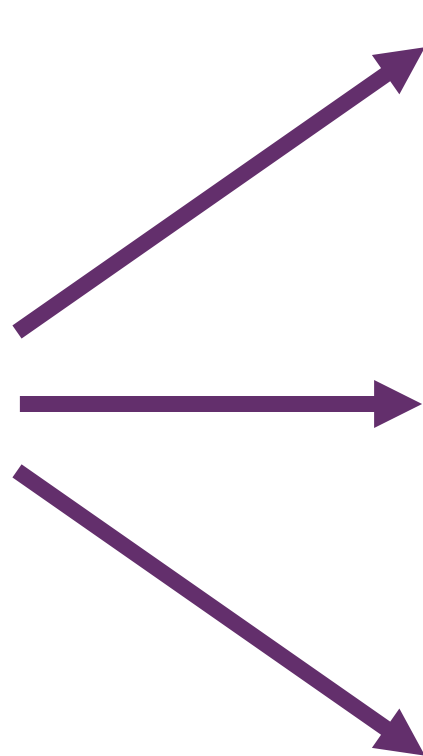
“Self-Driving Lab”

Implies a system that can autonomously learn and improve its function by analyzing and generating its own data.



Schneider, G. (2019). Mind and machine in drug design. *Nature Machine Intelligence* 1(3), 128-130.
doi: 10.1038/s42256-019-0030-7.

ASPIRE



Modern Lab Bench



Integrated Automation



Space Chemistry??



Many Challenges in Automating Chemistry

Methodology

Robust, efficient, safe, mild conditions and high yielding

Physical Operations

Reagent delivery
Product Isolation (workup)
Evaporation (solvent removal or exchange)

Analysis

Invasive (sampling, chemical derivation, e.g. TLC, HPLC)
Non-invasive (spectral analysis, UV-vis, FTIR, Raman, NMR)

Informatics

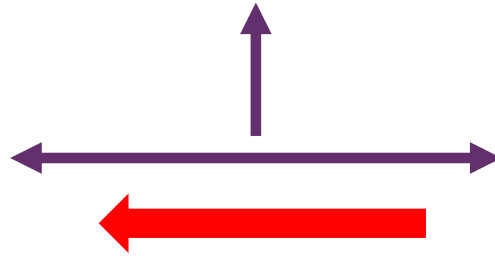
Data mining, curation, and feature engineering
Contextual representation

ASPIRE's Technology Ecosystem



ASPIRE

Translational



Next Gen Medchem

Collaboration-Oriented

Advanced Technology Development Lab

Technology-Driven



ASPIRE Next Gen MedChem Dashboard

“Complete Situational Awareness”

ASPIRE Modules Search

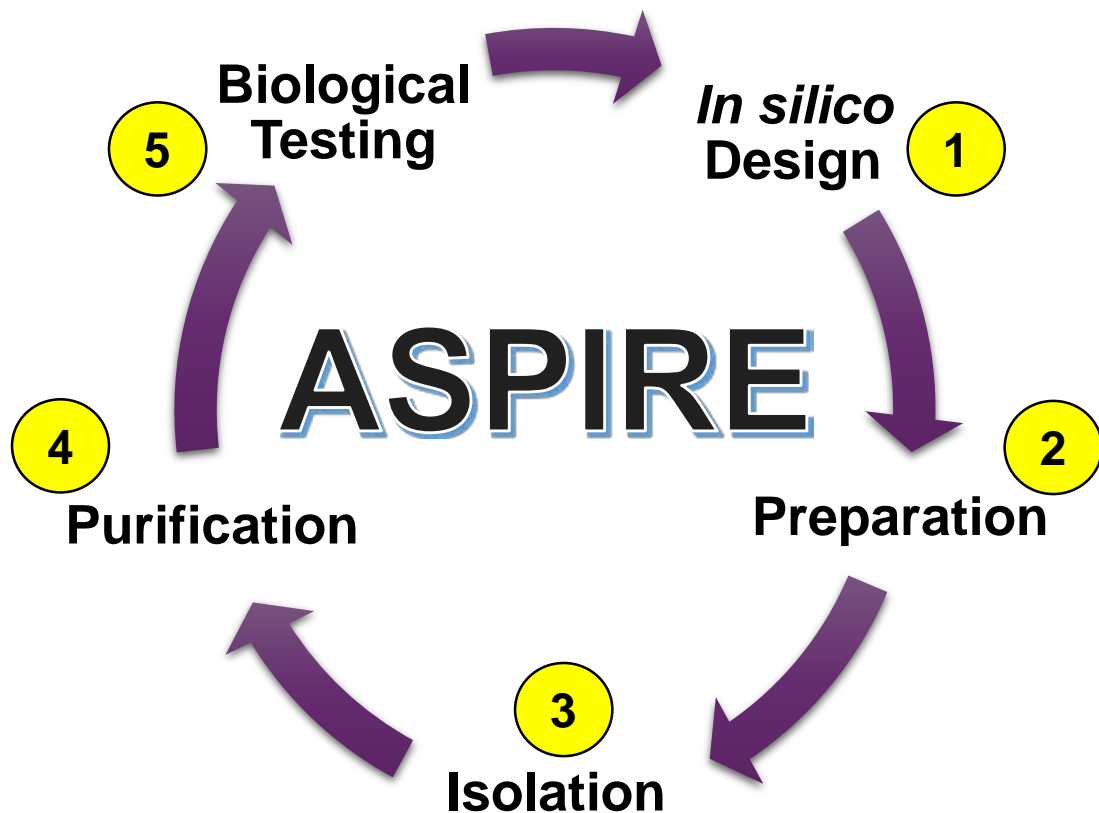
H HCube Pro Location : C2 Hood 24 Availability: ★★★★★ status: ●	Uw Microwave Initiator Location : C2 Hood 11 Availability: ☆☆☆☆☆ status: ●	Gm Weigh Station Location : C2 Hood 3 Availability: FREE ☆☆☆☆☆ status: ●
Sp SPE Station Location : C2 Hood 14 Availability: ★★★★★ status: ●	Zp Zaiput Extractor Location : C2 Hood 22 Availability: ★★★★★ status: ●	Ly Lyophilizer Location : C2 Hood 6 Availability: ★★★★★ status: ●
Ev Batch Evaporator Location : C2 Hood 5 Availability: ☆☆☆☆☆ status: ●	Nm NMR Location : C2 Hood 14 Availability: ☆☆☆☆☆ status: ●	Po V-10 Pooling Location : C2 Hood 16 Availability: ★★★★★ status: ●

Facilitated Development & Deployment Through Azure IOT Platform



An Evolution from Research in Silos to
Research in Real-time!

ASPIRE – Automated, Closed Loop Processing



Process Detail

① In silico Design

- Target Design
- Scored Synthetic Target(s)
- Retrosynthetic Analysis
- Route Selection
- Reaction Templating

② Preparation

- Reaction Setup
- Reaction Control
- Reaction Monitoring
- Reaction Completion

③ Isolation

- Workup Selection
- Workup Platform(s)
- Purification Submission

④ Purification

- Method Selection
- Purification Platform(s)
- Fraction Handling
- Characterization
- Delivery (Storage Logistics)
- Reaction Outcome Reporting

⑤ Biological Testing

- Assay Selection
- Assay Platform(s)
- Data Collection / Analysis
- Knowledge Assessment
- New Hypothesis (next iteration)

A New Renaissance in Chemical Synthesis Automation (In a Research Setting)

**ASL
2009**

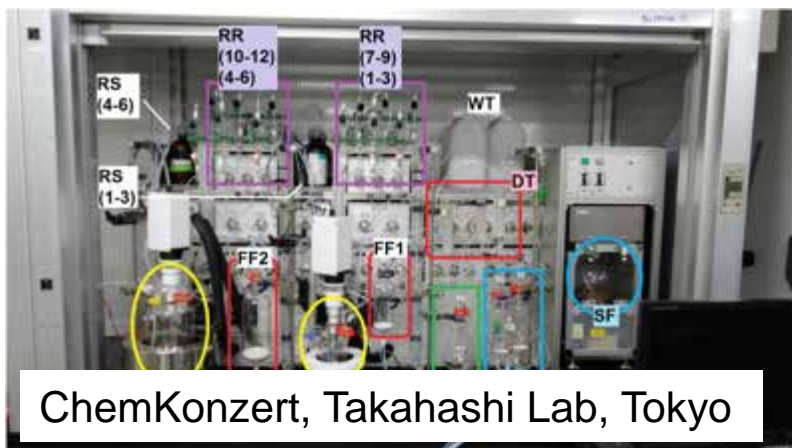
**Lilly Life Sciences Studio (L2S2)
2020**



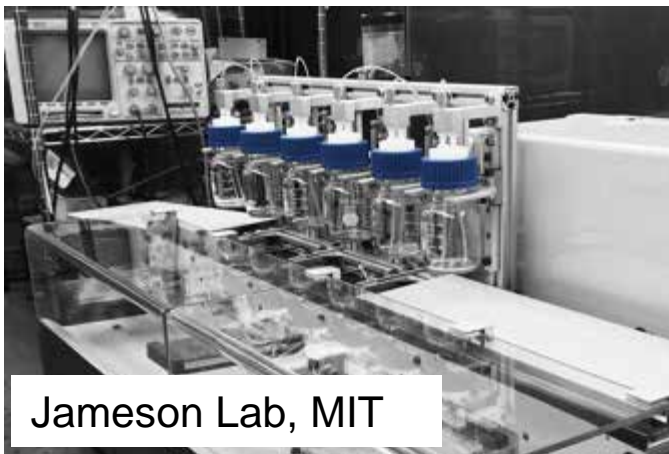
Eli Lilly, Indianapolis, IN

Eli Lilly / Strateos, San Diego, CA

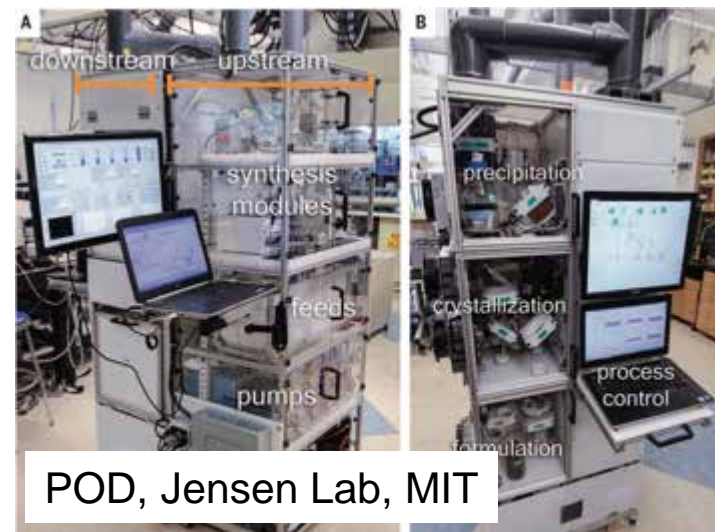
Additional Contemporary Examples



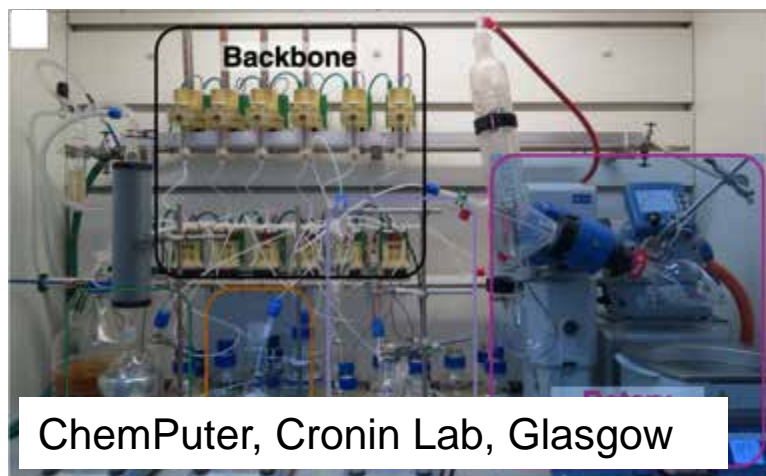
ChemKonzert, Takahashi Lab, Tokyo



Jameson Lab, MIT



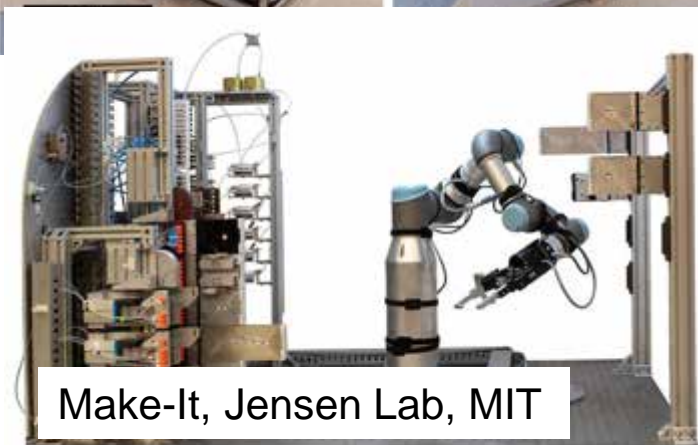
POD, Jensen Lab, MIT



ChemPuter, Cronin Lab, Glasgow



Radial Synthesis, Gilmore Lab, Max Planck Inst.



Make-It, Jensen Lab, MIT

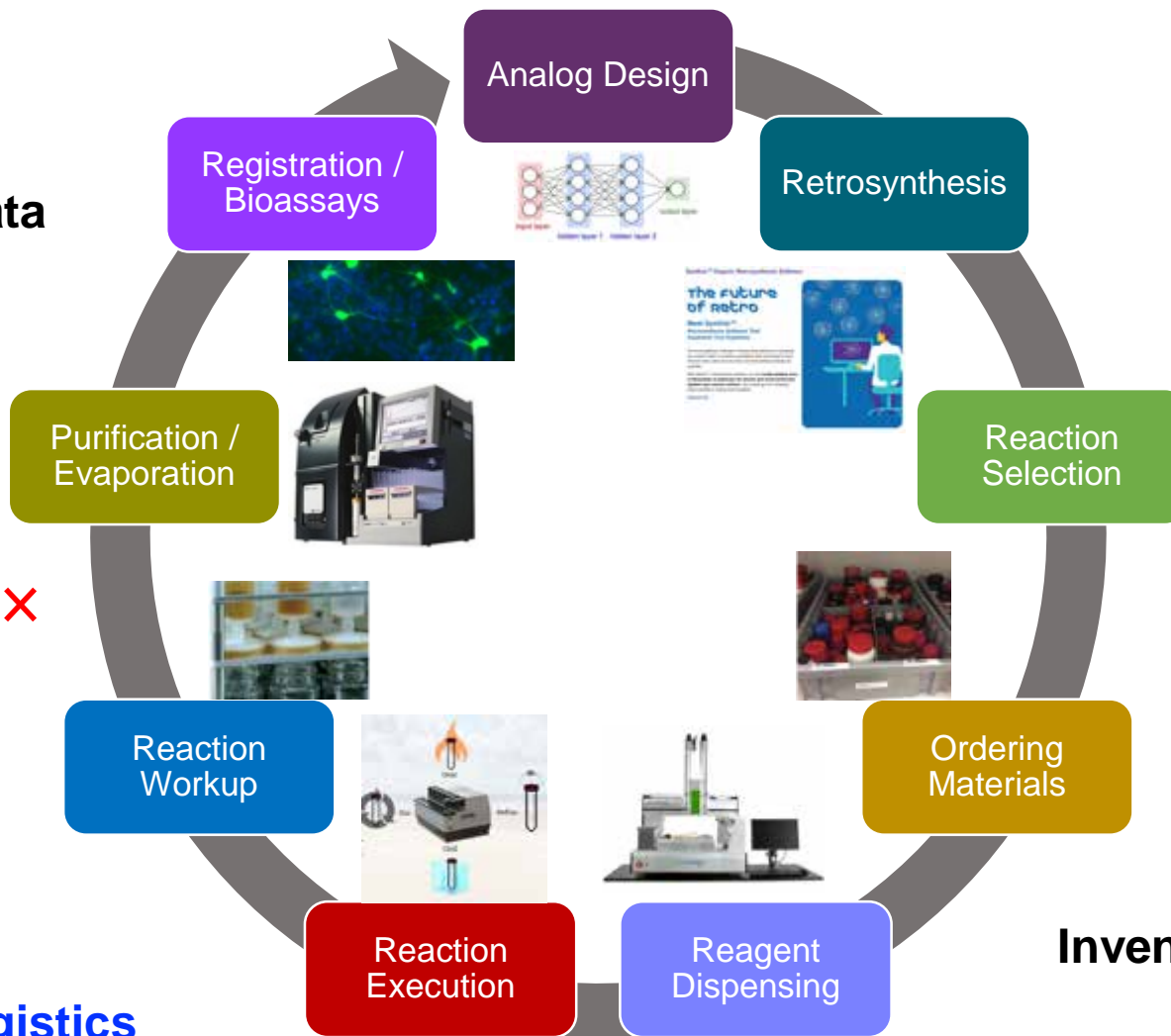
The Value of Integration

✓ **Biological Data**

Analytical Data ✗

ELN Data ✗

Inventory Data ✗



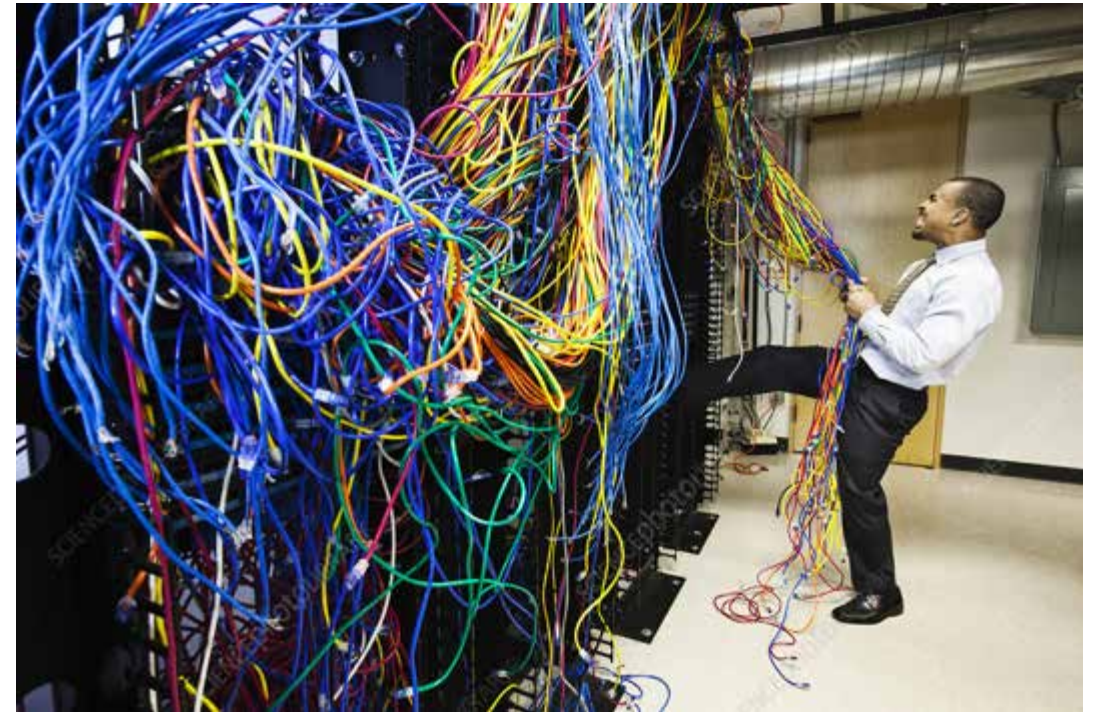
✓ **Value in improving efficiency in logistics**

✓ **Value in providing greater scientific connectivity & context**

Our Love-Hate Relationship with the Electronic Notebook



What we wish we had



What we really have

Lots of other folks have been looking at this problem as well

eLNs



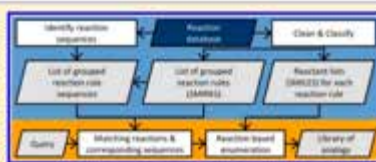
Mining Electronic Laboratory Notebooks: Analysis, Retrosynthesis, and Reaction Based Enumeration

Clara D. Christ,^{*,2} Matthias Zentgraf,³ and Jan M. Kriegl¹

¹Department of Lead Identification and Optimization Support, Boehringer Ingelheim Pharma GmbH & Co. KG, Birkendferstrasse 65, 88397 Biberach an der Riss, Germany

ABSTRACT: An approach to automatically analyze and use the knowledge contained in electronic laboratory notebooks (eLNs) has been developed. Reactions were reduced to their reactive center and converted to a string representation (SMIRKS) which formed the basis for reaction classification and *in silico* (retro-)synthesis. Of the SMIRKS that occurred at least five times, 98% successfully regenerated the original product. The extracted reaction rules (SMIRKS) and corresponding reactants span a virtual chemical space which showed a strong dependence on the size of the reactive center.

Whereas relatively few robust reaction types were sufficient to describe a large part of all reactions, considerably more reaction rules were necessary to cover all reactions in the eLN. Furthermore, reaction sequences were extracted to identify frequent combinations and diversifying reaction steps. Based on the extracted knowledge a (retro-)synthesis tool was built allowing for *de novo* design of compounds which have a high chance of being synthetically accessible. In an example application of the *de novo* design tool, various feasible retrosynthetic routes to the query molecule were obtained. Reaction based enumeration along the top ranked route yielded a library of 29 920 compounds with diverse properties, 99.9% of which are novel in the sense that they are unknown to the public domain.



Christ, C. D.; Zentgraf, M.; Kriegl, J. M., *J. Chem. Inf. Model.* **2012**, 52 (7), 1745-1756.

Patents



Big Data from Pharmaceutical Patents: A Computational Analysis of Medicinal Chemists' Bread and Butter

Nadine Schneider,^{*,1} Daniel M. Lowe,² Roger A. Sayle,³ Michael A. Tarselli,² and Gregory A. Landrum^{1,2}

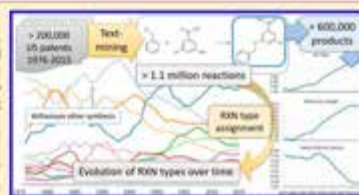
¹Novartis Institutes for BioMedical Research, Novartis Pharma AG, Novartis Campus, 4002 Basel, Switzerland

²Novartis Institutes for BioMedical Research, 186 Massachusetts Avenue, Cambridge, Massachusetts 02139, United States

³NextMove Software Ltd., Innovation Centre, Unit 23, Science Park, Milton Road, Cambridge CB4 0EY, U.K.

Supporting Information

ABSTRACT: Multiple recent studies have focused on unraveling the content of the medicinal chemist's toolbox. Here, we present an investigation of chemical reactions and molecules retrieved from U.S. patents over the past 40 years (1976–2015). We used a sophisticated text-mining pipeline to extract 1.15 million unique whole reaction schemes, including reaction roles and yields, from pharmaceutical patents. The reactions were assigned to well-known reaction types such as Wittig olefination or Buchwald–Hartwig amination using an expert system. Analyzing the evolution of reaction types over time, we observe the previously reported bias toward reaction classes like amide bond formations or Suzuki couplings. Our study also shows a steady increase in the number of different reaction types used in pharmaceutical patents but a trend toward lower median yield for some of the reaction classes. Finally, we found that today's typical product molecule is larger, more hydrophobic, and more rigid than 40 years ago.



Schneider, N.; Lowe, D. M.; Sayle, R. A.; Tarselli, M. A.; Landrum, G. A., *J. Med. Chem.* **2016**, 59 (9), 4385-4402.

Literature



Information Retrieval and Text Mining Technologies for Chemistry

Martin Krallinger,^{1,2} Obedula Rabal,^{1,3} Anália Lourenço,^{1,4,5} Julen Oyarzabal,^{1,6} and Alfonso Valencia^{*,4,5,7}

¹Structural Computational Biology Group, Structural Biology and BioComputing Programme, Spanish National Cancer Research Centre, C/Melchor Fernández Almagro 3, Madrid E-28029, Spain

²Small Molecule Discovery Platform, Molecular Therapeutics Program, Center for Applied Medical Research (CIMA), University of Navarra, Avenida Pio XII 55, Pamplona E-31008, Spain

³ISEI - Department of Computer Science, University of Vigo, Edificio Politécnico, Campus Universitario As Lagoas s/n, Ourense E-32004, Spain

⁴Centro de Investigaciónes Biomédicas (Centro Singular de Investigación de Galicia), Campus Universitario Lugo-Marcoende, Vigo E-36130, Spain

⁵CEB-Centre of Biological Engineering, University of Minho, Campus de Guizara, Braga 4710-057, Portugal

⁶Life Science Department, Barcelona Supercomputing Centre (BSC-CNS), C/Jordi Girona, 29-31, Barcelona E-08034, Spain

⁷Joint RSC-IBB-CRG Program in Computational Biology, Parc Científic de Barcelona, C/ Baldiri Reixac 10, Barcelona E-08028, Spain

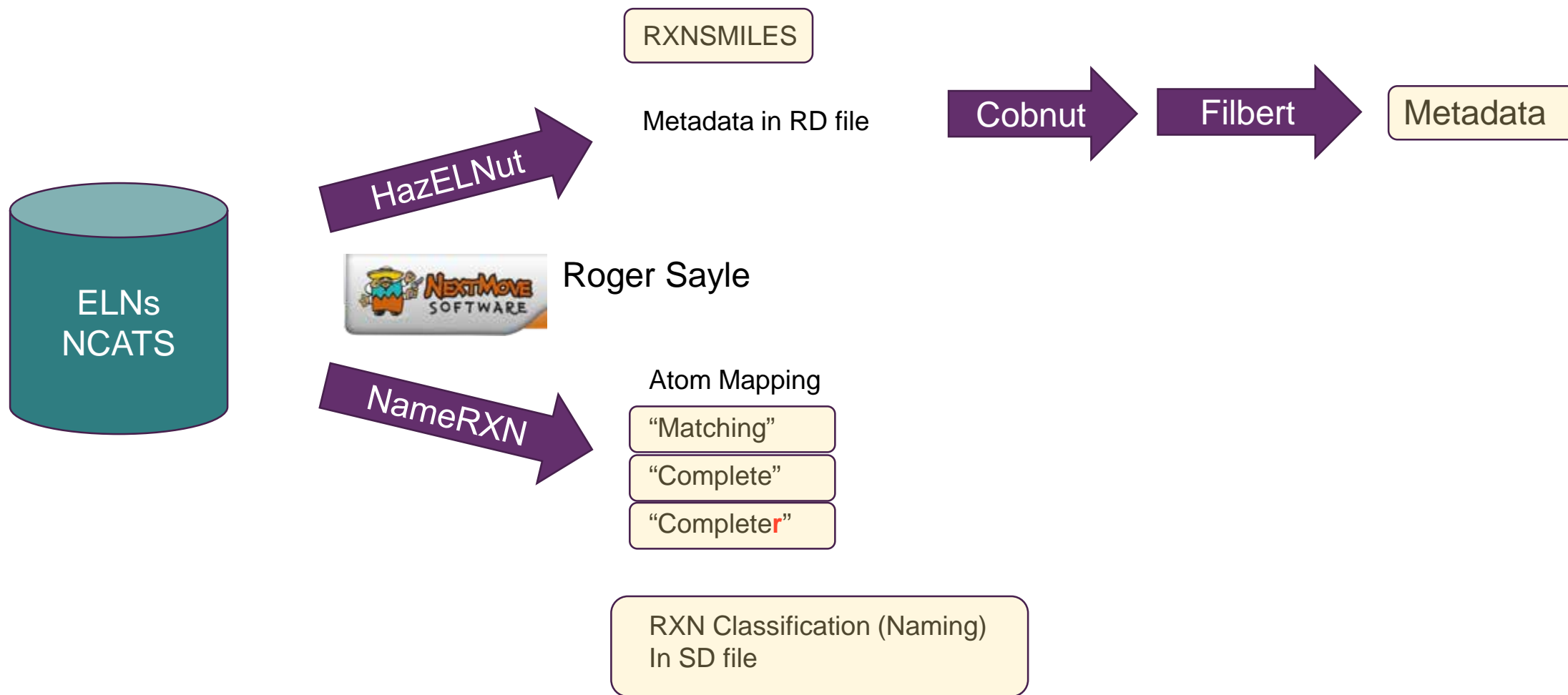
^{*}Institució Catalana de Recerca i Estudis Avançats (ICREA), Passatge de Lluís Companys 23, Barcelona E-08010, Spain

ABSTRACT: Efficient access to chemical information contained in scientific literature, patents, technical reports, or the web is a pressing need shared by researchers and patent attorneys from different chemical disciplines. Retrieval of important chemical information in most cases starts with finding relevant documents for a particular chemical compound or family. Targeted retrieval of chemical documents is closely connected to the automatic recognition of chemical entities in the text, which commonly involves the extraction of the entire list of chemicals mentioned in a document, including any associated information. In this Review, we provide a comprehensive and in-depth description of fundamental concepts, technical implementations, and current technologies for meeting these information demands. A strong focus is placed on community challenges addressing systems performance, more particularly CHEMDNER and CHEMDNER patents tasks of BioCreative IV and V, respectively. Considering the growing interest in the construction of automatically annotated chemical knowledge bases that integrate chemical information and biological data, cheminformatics approaches for mapping the extracted chemical names into chemical structures and their subsequent annotation together with text mining applications for linking chemistry with biological information are also presented. Finally, future trends and current challenges are highlighted as a roadmap proposal for research in this emerging field.



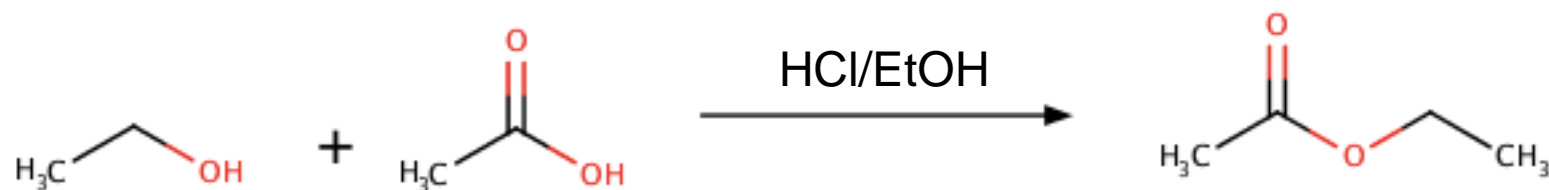
Krallinger, M.; Rabal, O.; Lourenço, A.; Oyarzabal, J.; Valencia, A., *Chem. Rev.* **2017**, 117 (12), 7673-7761.

ELN Extraction - Raw Output



Reaction Representation

- Reaction SMILES



OCC.CC(=O)O>[H+].[Cl-].OCC>CC(=O)OCC |f:2.3|

Ref: <https://www.daylight.com/meetings/summerschool01/course/basics/smirks.html>

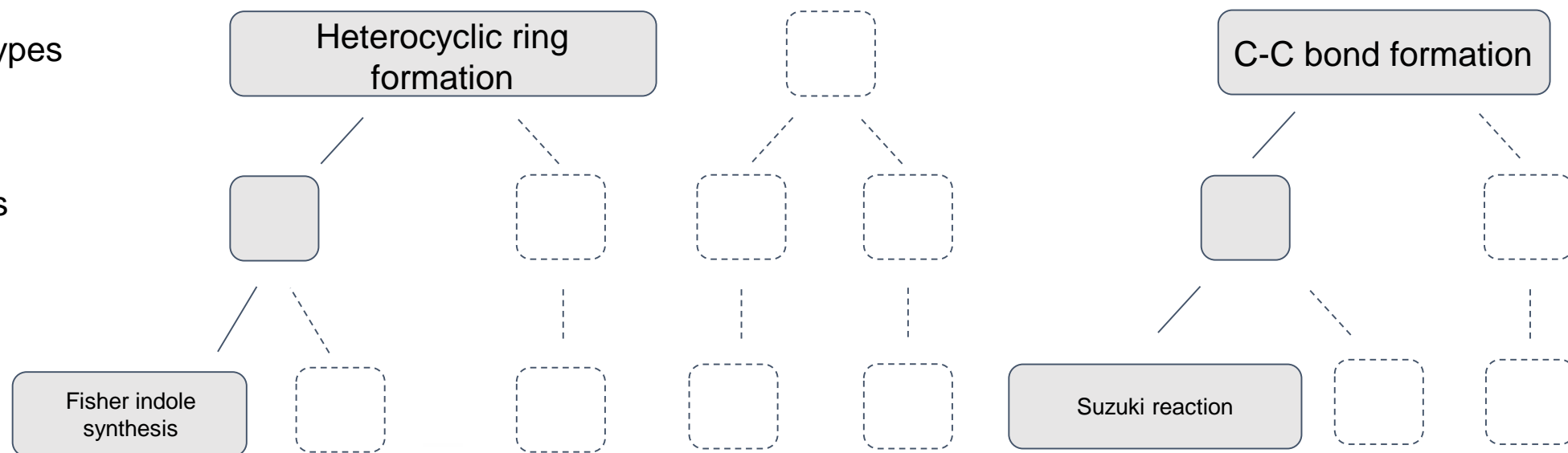
Reaction Ontology

- RXNO reaction ontology invented by Royal Society of Chemistry
- Ontology available at: <https://github.com/rsc-ontologies/rxno>
- Implemented by NextMove
- Can be utilized in curation and data mining

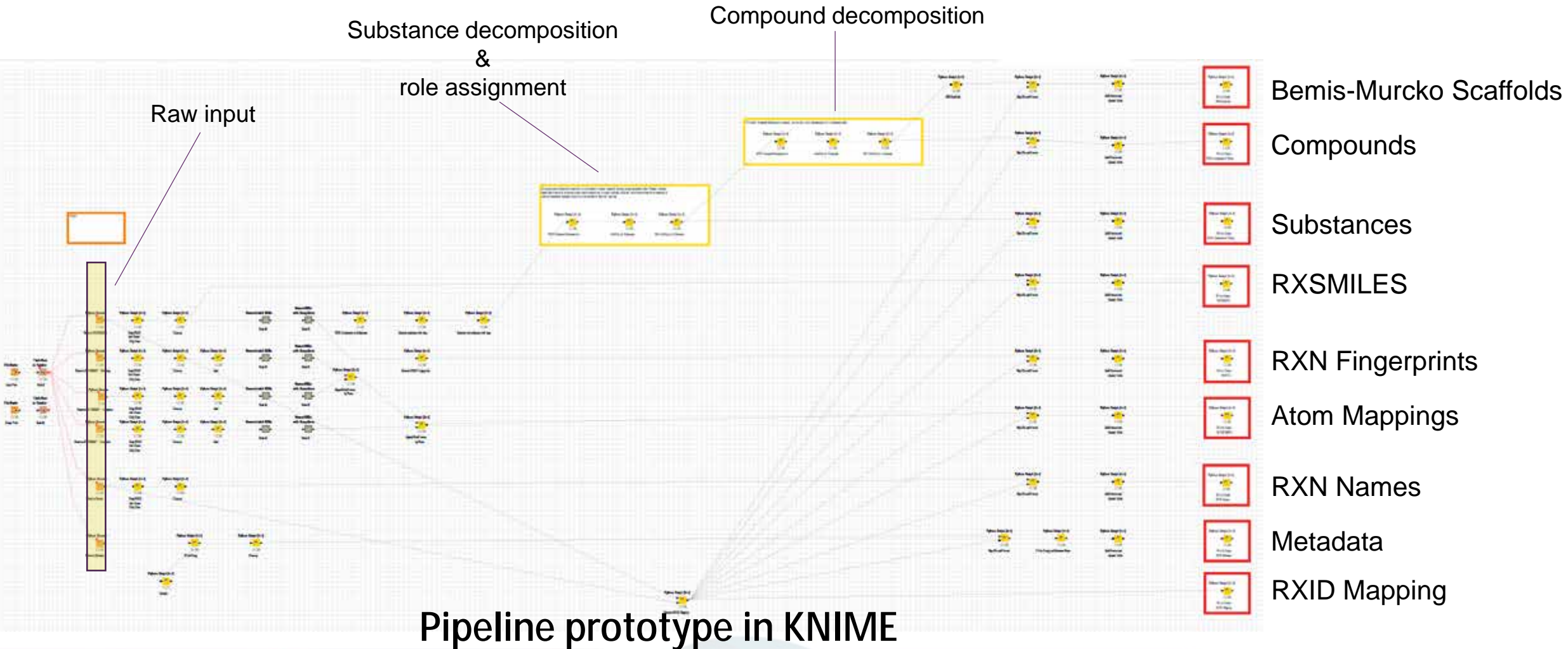
- 12 basic types

- 84 classes

- 500+ named reactions



Data Post-processing



Overall Process

ELN Extraction (~90K expts)



Data Curation (~16K expts)



Data Processing



Database Upload



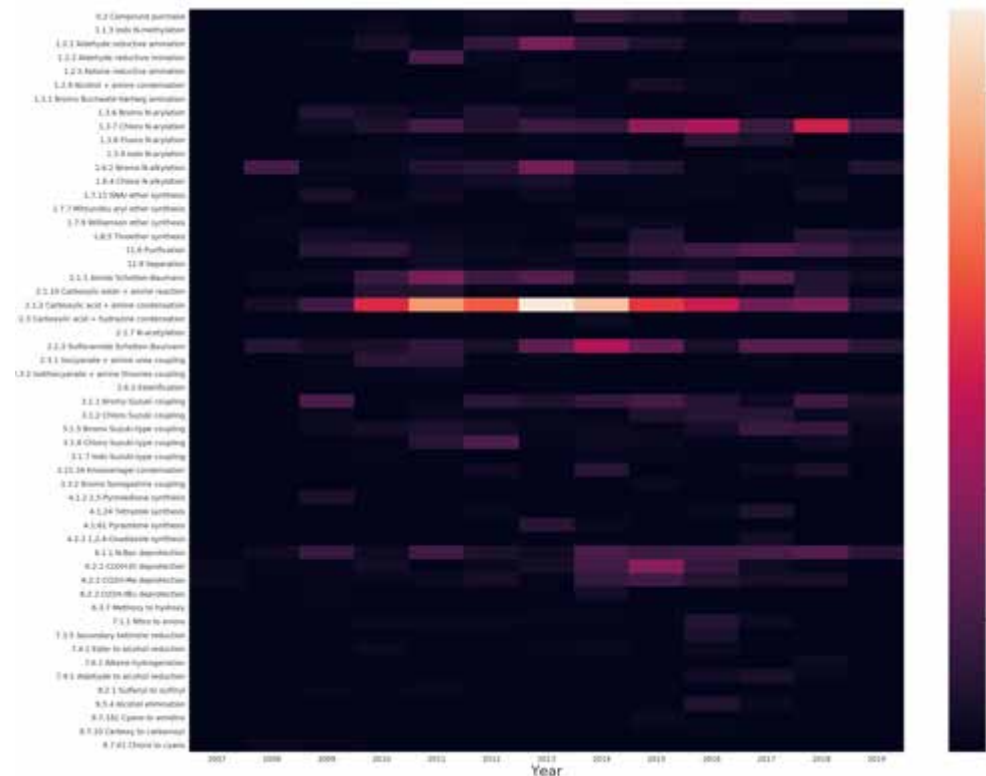
Reaction Analytics Dashboard



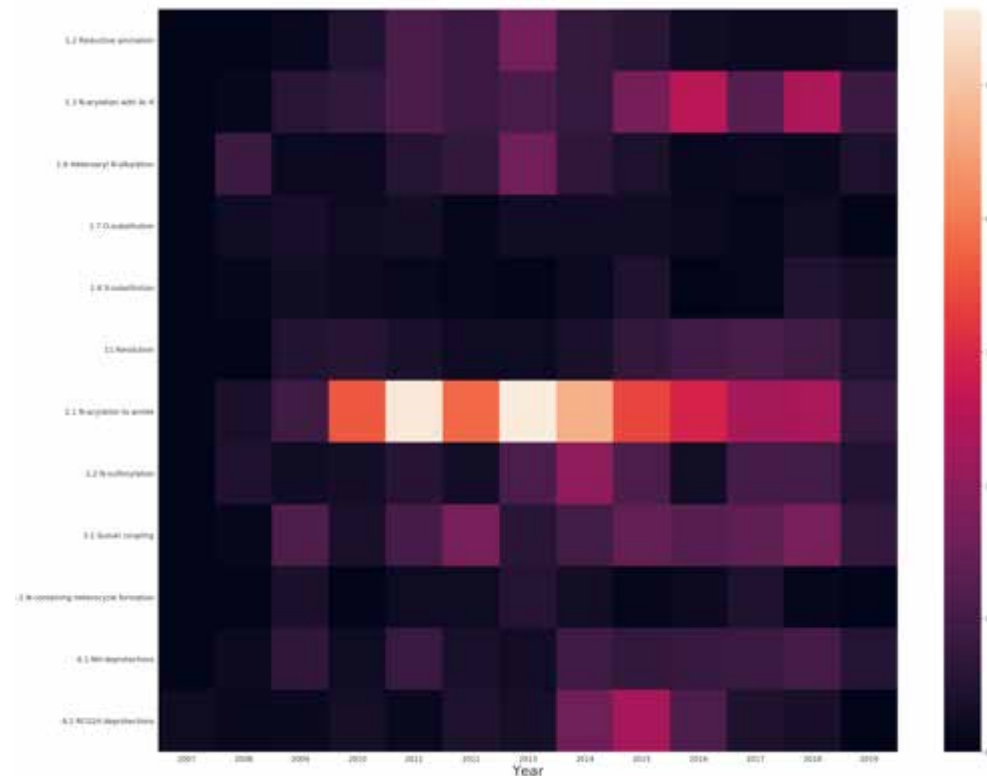
Some Preliminary Observations

Most Frequently Utilized Reactions per Year

Reaction Name (Top 25% of 216)

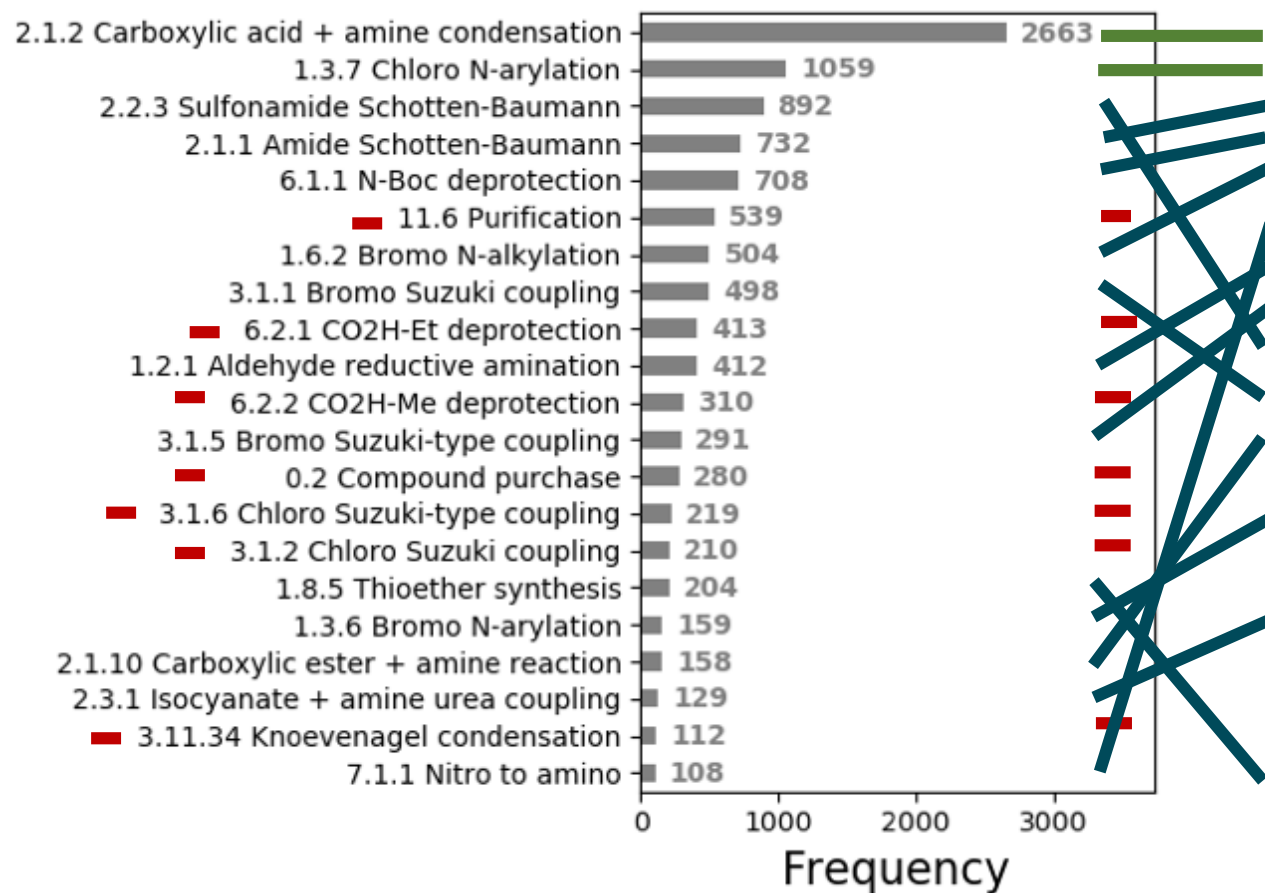


Reaction Class (Top 25% of 50)

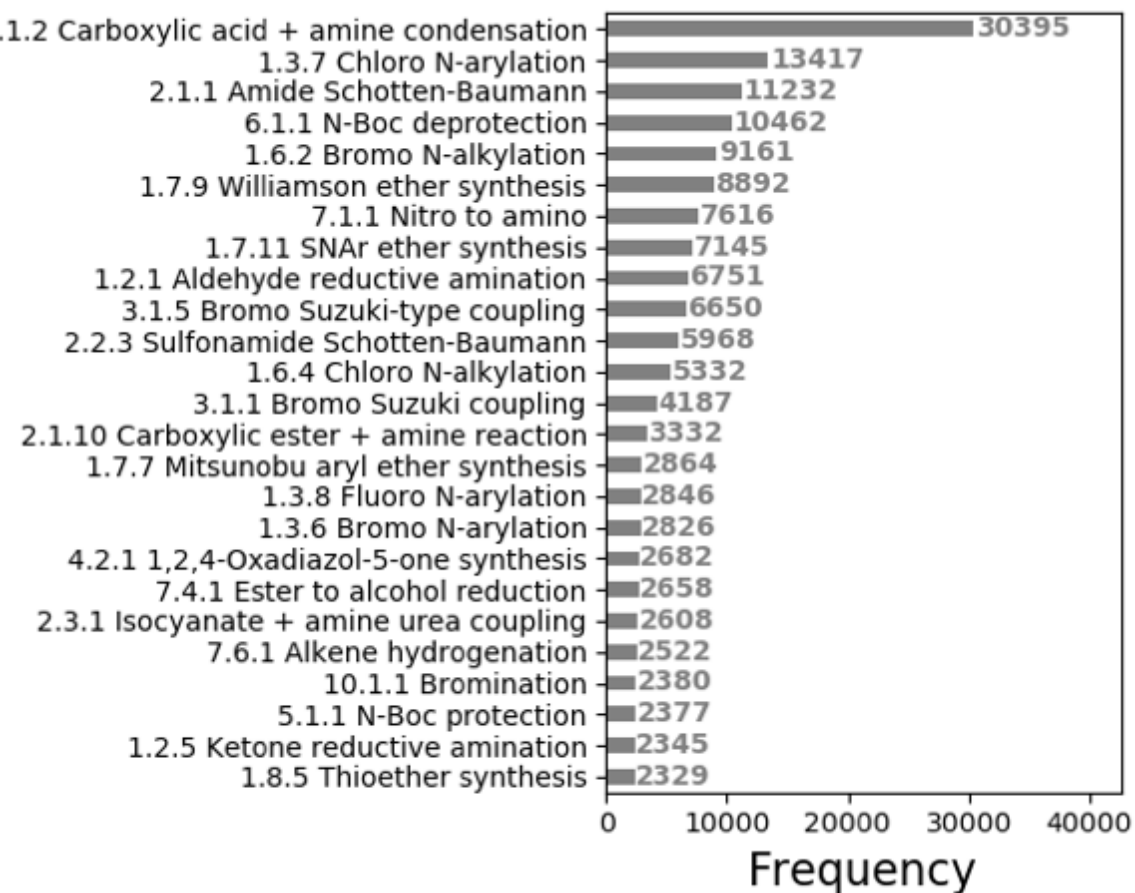


Comparison of Name Reaction Frequencies

Top 10% of 210 Reactions - ASPIRE



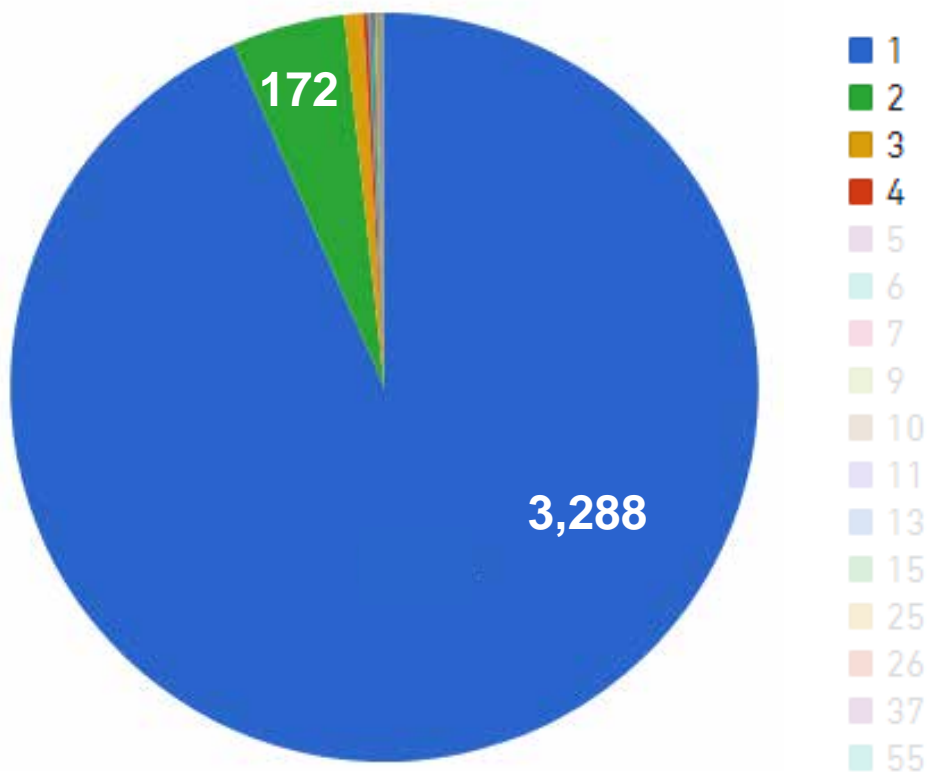
Top 5% of 505 Reactions - USPTO



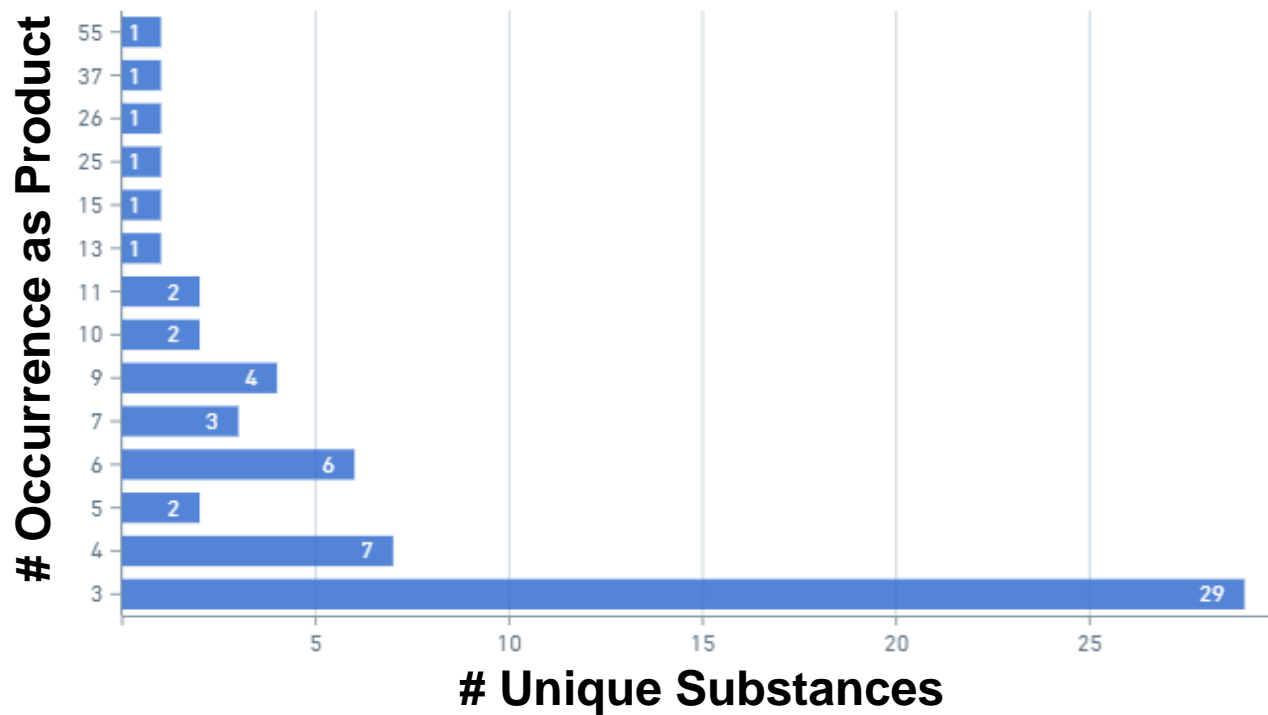
Substance Synthesis Frequency Distribution

Closed Reactions Only (~4K expts)

Occurrence as product ≥ 1



Occurrence as product ≥ 3



Reaction Analytics Dashboard Web App Prototype

← → ↻ ⓘ Not secure | 165.112.226.175:8050

CYMFHCDONPWFL-UHFFF

RXN ID | SMILES | InChi-Key

FIND

Graph 1

Count

Compound

Graph 2

Count

RX Type

Compound

OH

rx1dm	yes
Filter: data...	
NCGCRX-0004044	2018 5
NCGCRX-000782	
NCGCRX-001841	
NCGCRX-001957	
NCGCRX-0050396	2017 Secondary ketimine reduction [Cyano or imi

PREVIOUS NEXT

Timeline: 2007, 2010, 2012, 2015, 2017, 2019

Marvin JS

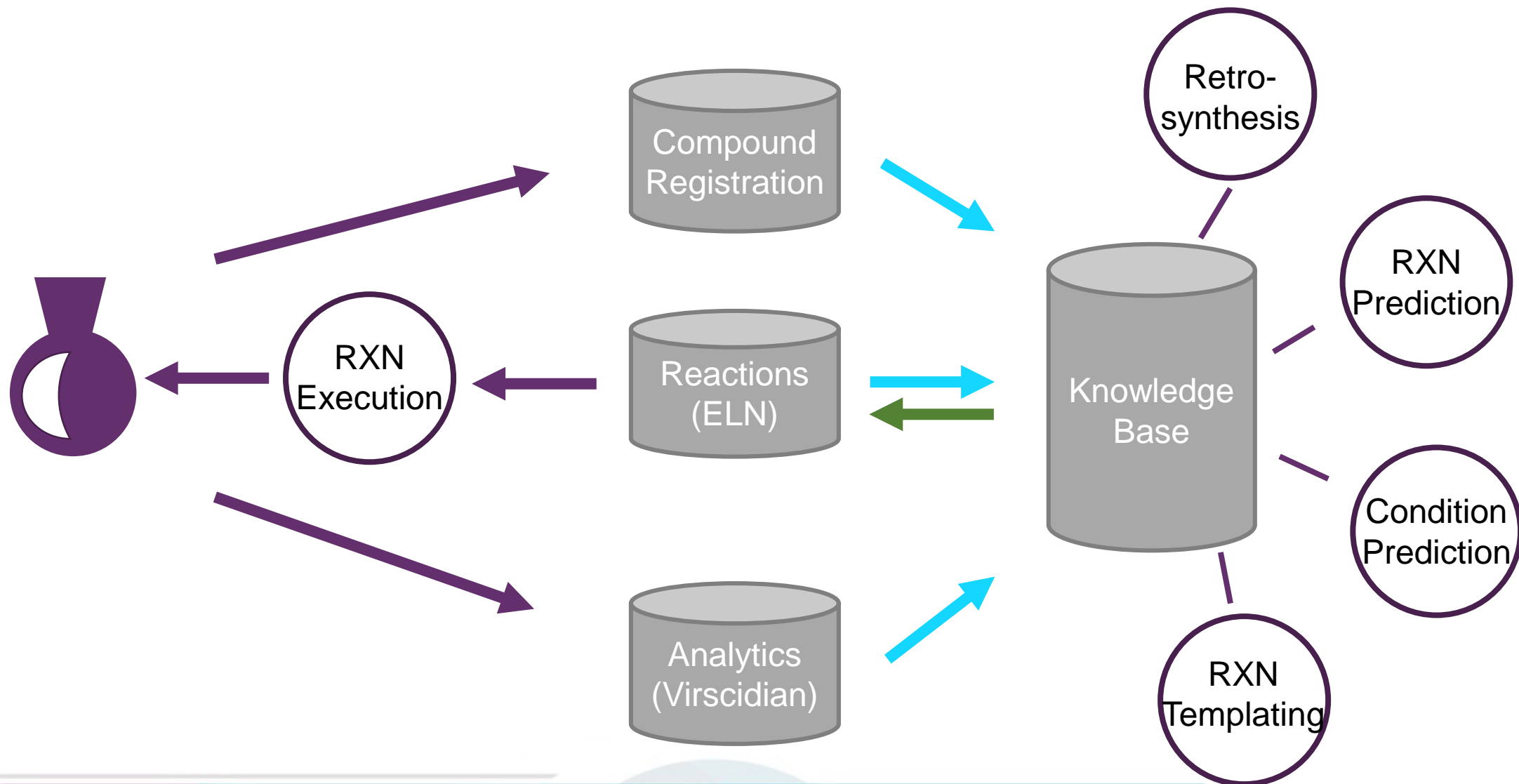
Chemical network graph showing reaction types and their relationships.

Navigation: < >



Gergely Zahoránszky-Köha

Computational Platform to Support Chemistry Automation



Reaction Analytics – From Information to Actionable Knowledge

Frequency Distributions:

- Reaction Type (overall or by year)
- Reagent use by name or type
- Reagent use by quantity
- Catalyst use (most common vs least common)
- Catalyst use (by yield)
- Reactant use

Cross Comparisons with Inventory:

- Reactants infrequently used or never used
- Use of rare (long lead time) reactants

Reaction Networks:

- Average or median length of synthesis
- Longest synthetic sequence
- Reaction types represented / NOT represented
- Reactant diversity analysis
 - Total count & frequency
 - Performance by functionality

Productivity Metrics:

- Product production cycle times
- Cycle time by reaction type
- Yield Distributions

Reaction Analytics by Platform (i.e. Target Space)

Reaction Networks:

- Average of
- Longest
- Reaction
- Reacta
 - Total
 - Perform

How is our reaction space evolving over time?

How would we like to see it evolve?

Where can we intervene most effectively?

What opportunities might we be overlooking?

Project /
Platform B

Project /
Platform D

Reaction Templating

“How do we get reactions to automatically set themselves up?”

Think: Object + [Role] + Properties → Robotic Process

Features:

- 1) Reaction Ontology
- 2) Complete Reaction Annotation (“No atom left behind”)
- 3) Standardized Reaction Equations (scale-independent)
- 4) Machine-Interpretable Temporal Instructions

The Work Ahead

- Clearly defined annotation strategy for benchmark reactions
- Test how discretized scale factors will work in practice
- Complete the reaction template lexicon / graph
- Build and test the hardware specification map
- Detail module design and construction
- Integrate-test, integrate-test, integrate-test...

Thank you!

Work presented here was supported in part by the Intramural/Extramural research program of the NCATS, NIH

