# Data Explosion in Chemistry: What Are We Going to Do With All the Data, and What Will It Do to Us

Caroline Lynn Kamerlin

kamerlinlab.com

# Transformation in How We Do Chemistry



How it started: "Never trust a computational chemist"

How it's going: even "lab chemists" need to code!

Image © Joppe van der Spoel, Studio de Wilde Muis

Georgia Tech

# Transformation in How We Do Chemistry



Breakdown of traditional barriers between experimentalists and theorists – digitalization is unavoidable!

Image © Joppe van der Spoel, Studio de Wilde Muis

# Data Explosion Revolutionizes Chemistry

- Large language models: turning words into structures, molecules and chemical reactions.

- AI changing drug discovery, human health, chemical synthesis.

- Digital chemistry in materials science: design of new polymers, functional materials.

- Machine learning in enzyme design: sustainable chemistry.

- More fundamentally, digital chemistry also changes how we interact and communicate with each other as scientists.

Georgia Tech

# What Does Digital Chemistry Mean To You?

"Digital chemistry" means different things to different people.

- Computational biophysicist – molecular simulations?

- Synthetic organic chemistry – AI in synthesis?

- Pharmaceutical chemist – AI in drug discovery?

- Structural biology – databases? Protein structure prediction?

  Asked 4 diverse colleagues the same question, got 4 answers.

Georgia Tech

# Digital Chemistry Means New Data Pipelines



"We've decided to take big data to the next level - humongous data" (from vox.com cartoon)

# Paradigm Shift for Chemistry Data

- How we obtain data?
    - experimental / computational workflows.

- How we store data?
    - infrastructure challenges as data gets huge.

- How we curate data?
    - how to make it understandable to others.

- How we share it with others.
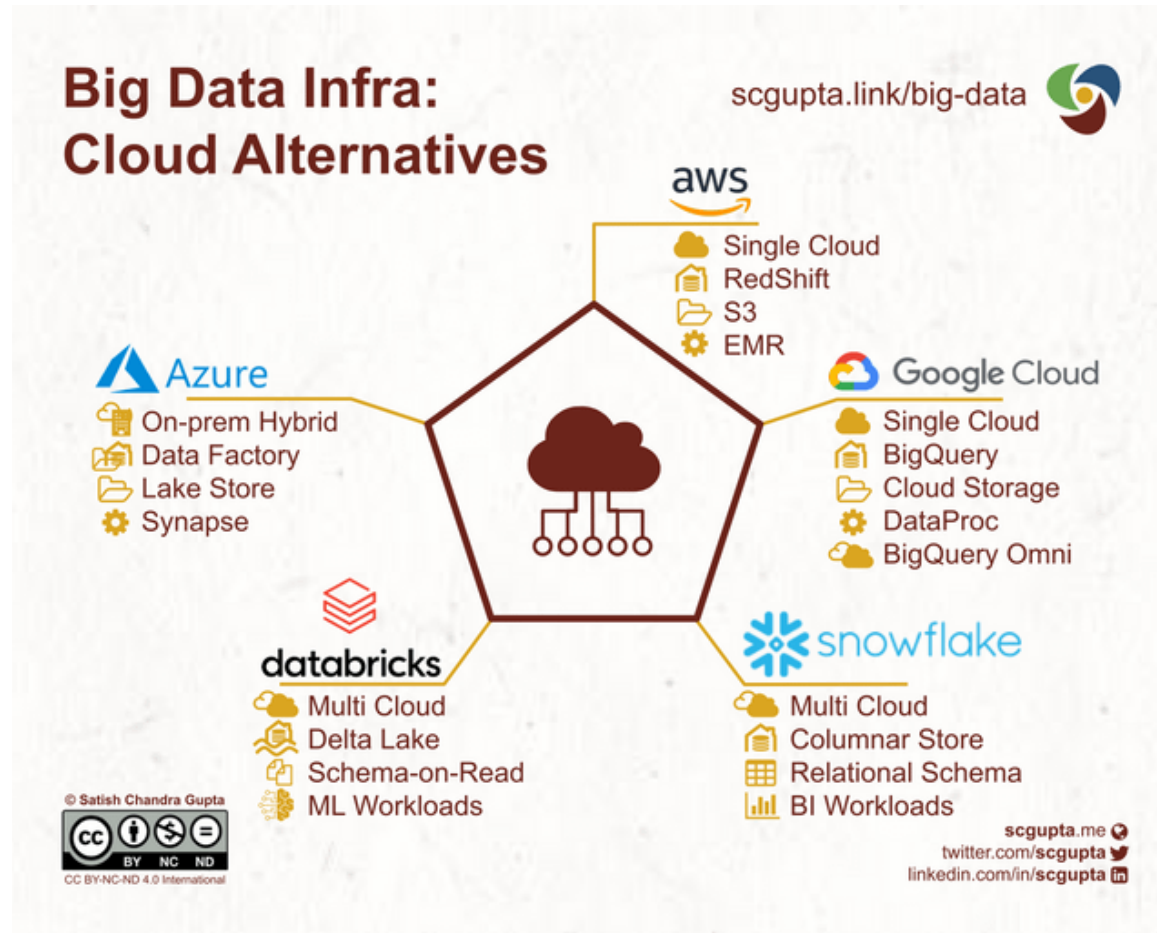    - infrastructure + curation + what/when?

New ways of working with data, needs new skills, approaches, *etc*.

Georgia Tech

# Paradigm Shift for Chemistry Data

- How we obtain data?
    - experimental / computational workflows.

- How we store data?
    - infrastructure challenges as data gets huge.

- How we curate data?
    - how to make it understandable to others.

- How we share it with others.
    - infrastructure + curation + what/when?

Different disciplines have different data needs (Excel spreadsheets to PB of data) and different expectations of digital chemistry.

Georgia Tech

Give us more infrastructure! Thank you for coming to my talk.

Big Data Infra: Cloud Alternatives

Need an "infrastructure revolution" – for analyzing the data, storing the data, disseminating the data.

# Infrastructure – Handling All The Data



Skills training – how do the next generation of chemists interface with this infrastructure?

# Paradigm Shift – Transition to Open Data



Not just papers need to be open access!

# Paradigm Shift – Transition to Open Data



Training to identify reputable / safe repositories to store data.

# FAIR and Digital Standards for Chemistry



[http://go-fair.org](http://go-fair.org)

GO-FAIR + IUPAC working on fair standards for chemistry.

Image from Foster Open Science

# GO-FAIR Chemistry Implementation Network

GO-FAIR ChIN guiding principles:

- Findable chemistry data.

- Reusable code and data – validation, compilation/aggregation, incorporation into future work, data mining.

- The use of standards at source and throughout the information cycle.

- Availability and accessibility of tools and infrastructure.

http://go-fair.org

# GO-FAIR Chemistry Implementation Network

- The use of persistent identifiers and machine readability at the core.

- Management and oversight of standards.

- Use of general data standards outside of chemistry where appropriate and FAIR in their implementation.

- Enable and promote use of chemical data standards in other disciplines that work with chemical data.

http://go-fair.org

# IUPAC Push for Digital Standards

- The IUPAC International Chemical Identifier (InCHI).

- ThermoML (XML based IUPAC standard for thermodynamic property data).

- JCAMP-DX (IUPAC standard format family for spectral data exchange).

- In progress: IUPAC SMILES+, machine accessible periodic table, standard for FAIR data management of spectroscopic data, making IUPAC assets FAIR, and more!

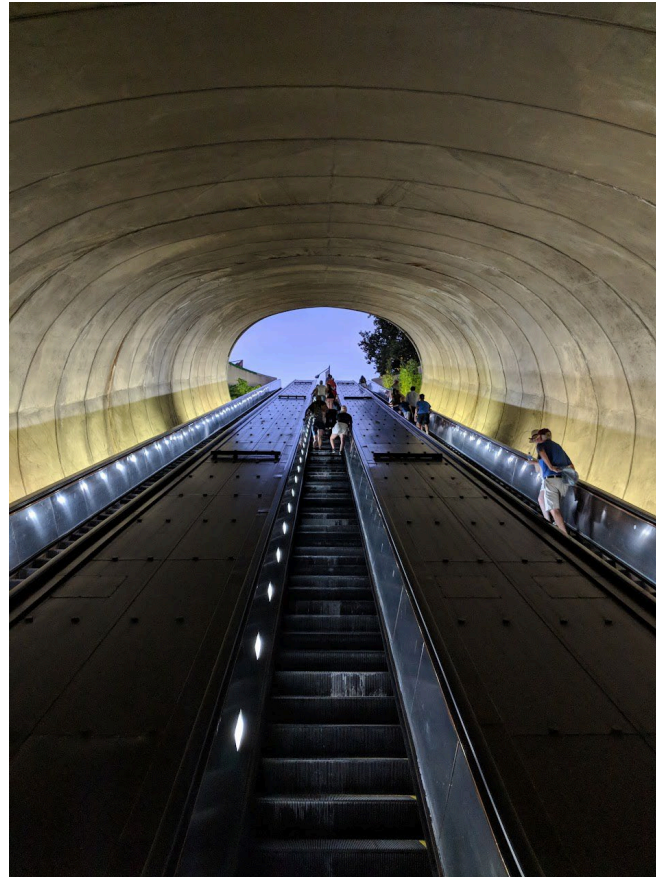https://iupac.org/what-we-do/digital-standards/

# Training: Chemistry Degrees Are Changing!

- Chemistry degree 1998-2002: minimal comp chem, only in final year (enough to make me be a theorist though!)

- Revolution in skills training for chemists – so many new skills you need to learn.

- EU "Traineeship in Digital Skills" as part of mobility schemes.

- Dedicated courses in "Digital Chemistry" at many institutions.

Next generation of chemists are going to live with data, computers and workflows that would be unrecognizable 20 years ago!

# What Does the Future Hold?



Interdisciplinarity | Big Data | Complexity | Integration