# THE IMPORTANCE OF DATA QUALITY IN SCIENTIFIC AI

Jacob Al-Saleem, Ph.D.

CAS
A division of the
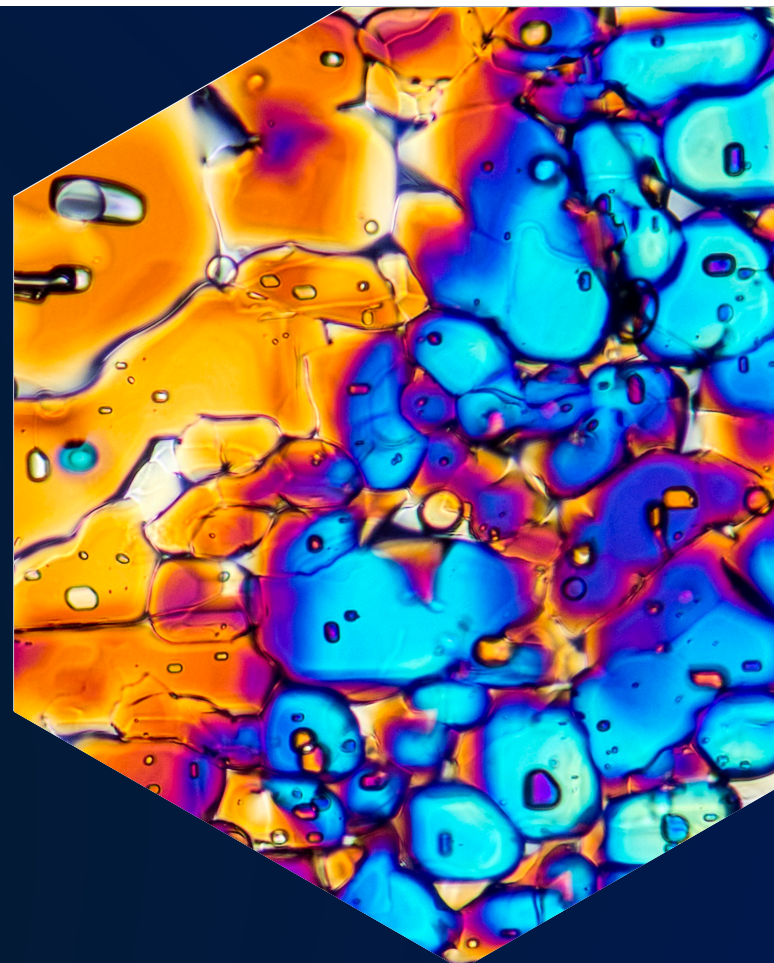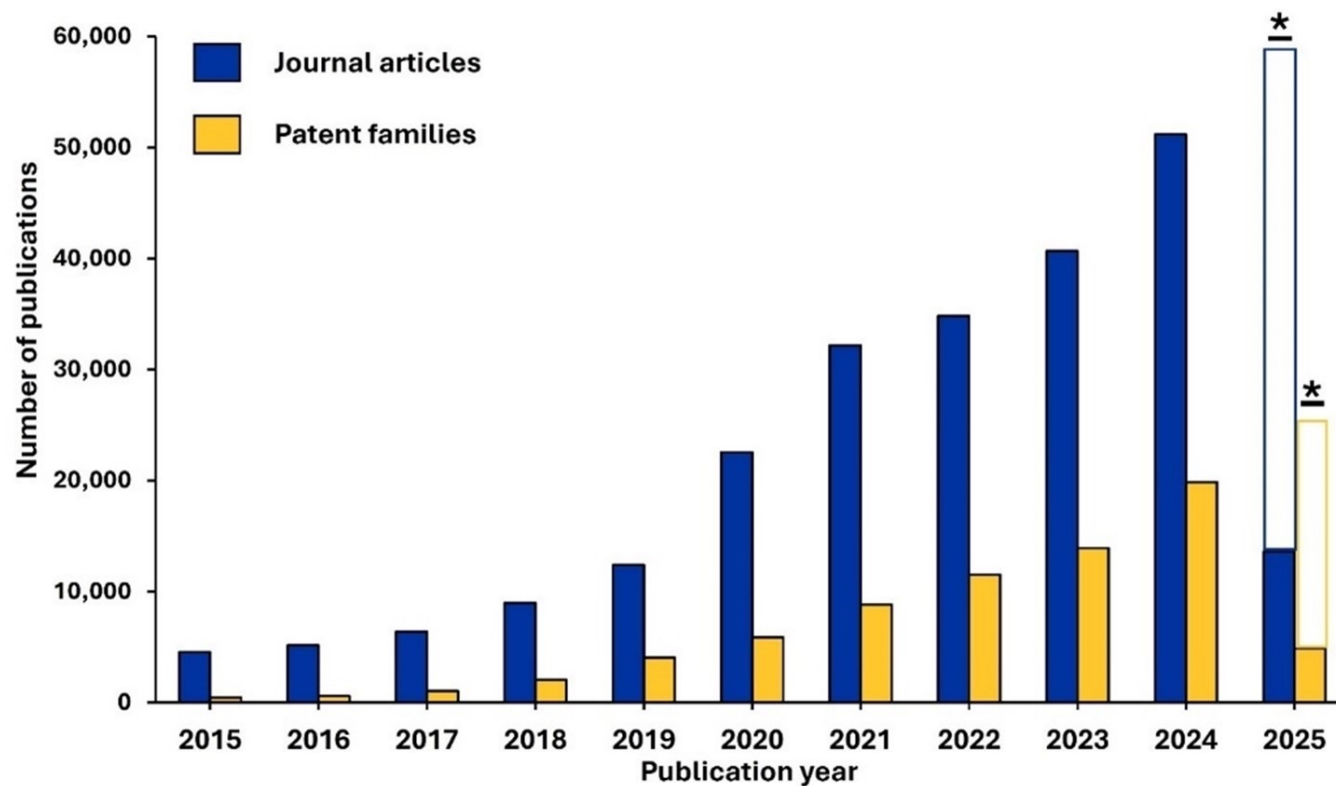American Chemical Society

# About CAS

CAS connects the world's scientific knowledge to accelerate breakthroughs that improve lives.

As a specialist in scientific knowledge management, we empower global innovators across industries with the essential data, solutions, services, and expertise needed to navigate today's complex digital information landscape and make more confident decisions.
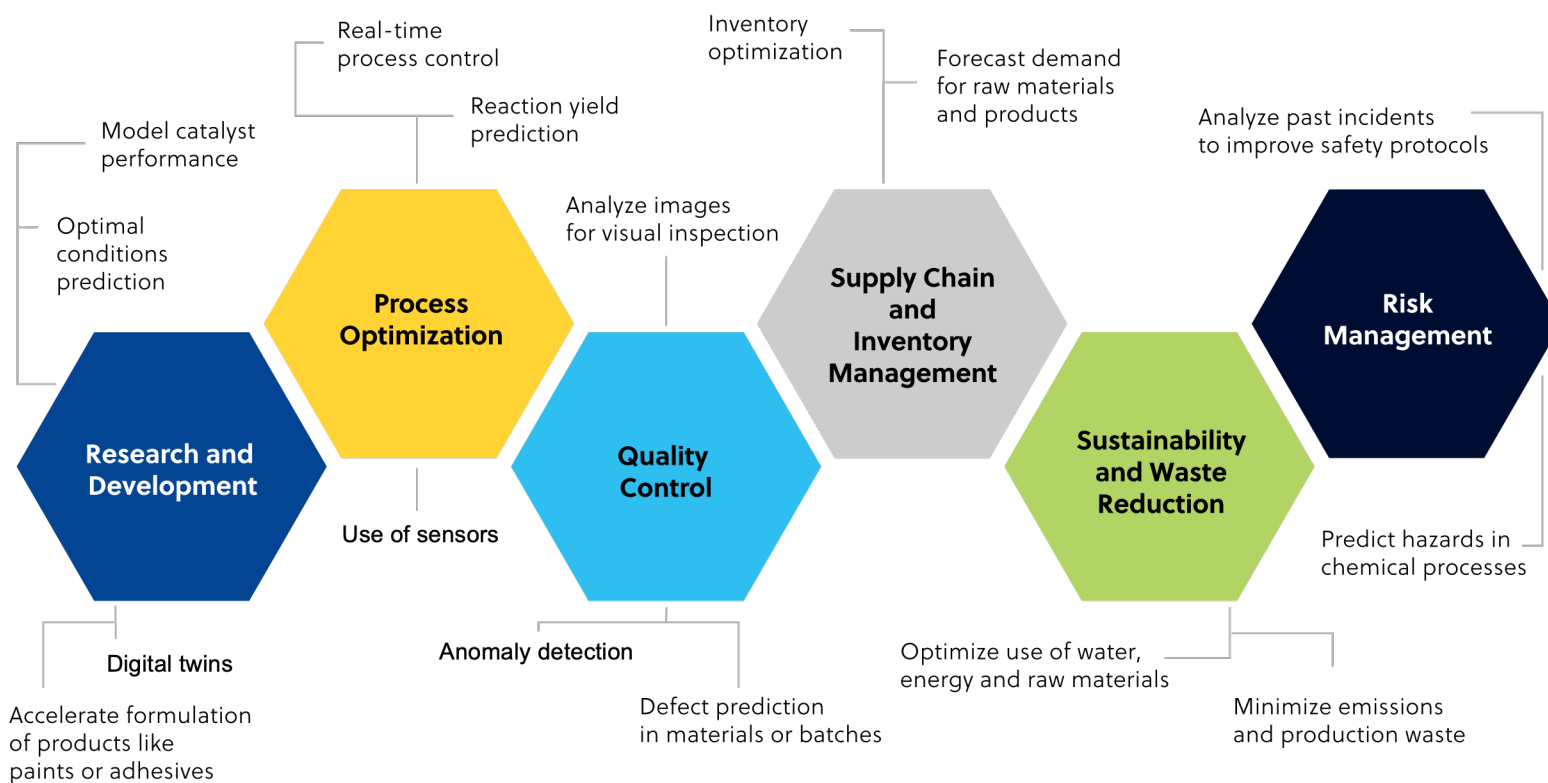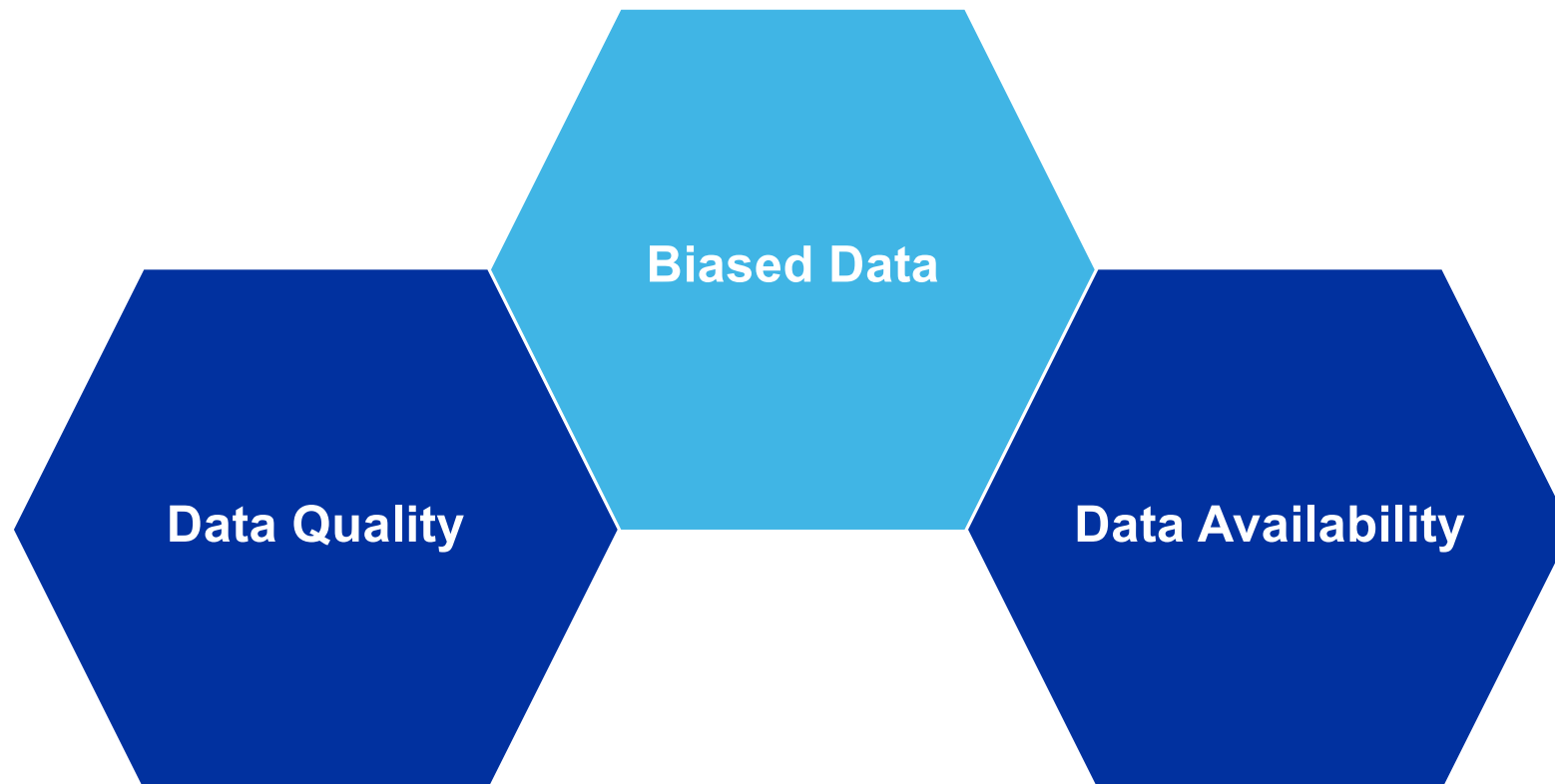
CAS is a division of the American Chemical Society.

**CAS**
A division of the
American Chemical Society

# Publications featuring AI grows year over year

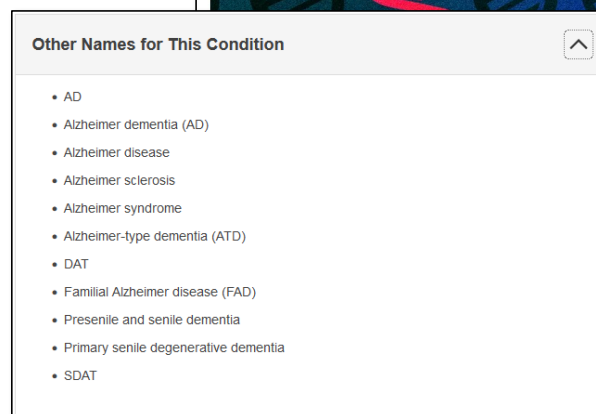# Emerging applications touch every phase of innovation



Model catalyst performance

Optimal conditions prediction

Real-time process control

Reaction yield prediction

**Process Optimization**

Analyze images for visual inspection

Inventory optimization

Forecast demand for raw materials and products

**Supply Chain and Inventory Management**

Analyze past incidents to improve safety protocols

**Risk Management**

**Research and Development**

Use of sensors

**Quality Control**

**Sustainability and Waste Reduction**

Predict hazards in chemical processes

Digital twins

Anomaly detection

Accelerate formulation of products like paints or adhesives

Defect prediction in materials or batches

Optimize use of water, energy and raw materials

Minimize emissions and production waste

CAS
A division of the
American Chemical Society

# What challenges exist in scientific AI?



Biased Data

Data Quality

Data Availability

# What Challenges exist in scientific AI?



Data Quality

Biased Data

Data Availability

CAS
A division of the
American Chemical Society
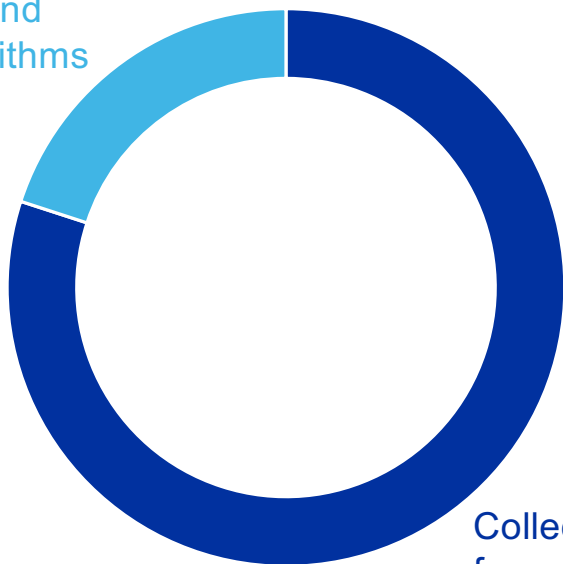
# Inconsistency in data

## Leading to problems downstream

– Many names for the same entity

– Different labs prefer different names

– Author errors

– Technology errors

– Different databases and IDs



**Genome Biology**

Home  About  Articles  Submission Guidelines  Submit manuscript

Comment | Open access | Published: 23 August 2016

Gene name errors are widespread in the scientific literature

Mark Ziemann, Yotam Eren & Assam El-Osta

SCIENCE / TECH / MICROSOFT

Scientists rename human genes to stop [spreadsheets] from misreading them as dates

/ Sometimes it's easier to rewrite genetics than update [software]

By James Vincent, a senior reporter who has covered AI, robotics, and more for eight years at The Verge.

Aug 6, 2020, 8:44 AM EDT

0  Comments (0 New)

If you buy something from a Verge link, Vox Media may earn a commission. See our ethics statement.

**Other Names for This Condition**

- AD
- Alzheimer dementia (AD)
- Alzheimer disease
- Alzheimer sclerosis
- Alzheimer syndrome
- Alzheimer-type dementia (ATD)
- DAT
- Familial Alzheimer disease (FAD)
- Presenile and senile dementia
- Primary senile degenerative dementia
- SDAT

CAS
A division of the
American Chemical Society

# The 80:20 data science rule

Messy data consumes time

Developing and
refining algorithms



Collecting, cleaning,
formatting, and
prepping data

## 90%



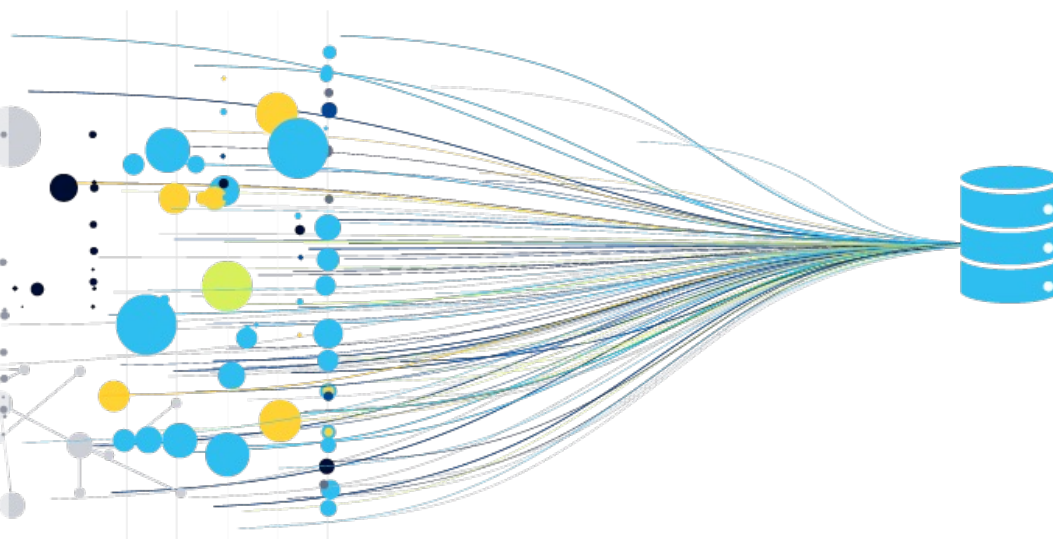**Data generated by
humans created in
the last 2 years
(est.)**

*Stobierski T. Harvard Business School Online. 2021*

CAS
A division of the
American Chemical Society

# Harmonization creates consistency

## From chaos to order

– Corrects inconsistencies and errors

– Combines different types of data

– Connects across data sources

CAS
A division of the
American Chemical Society
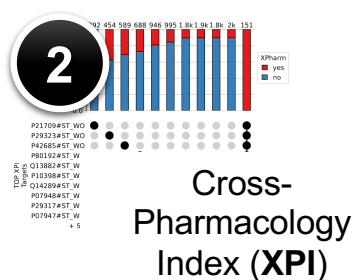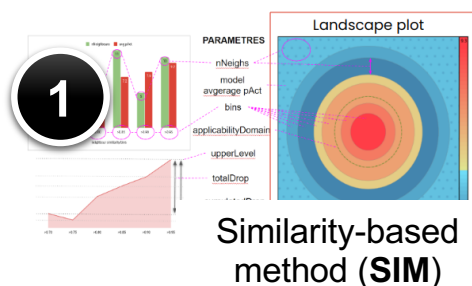
# Good harmonization is human driven

## Human effort prior to AI efficiency

– Human input ensures alignment to authority constructs

– Human assisted curation creates training data at the standards set by data scientists

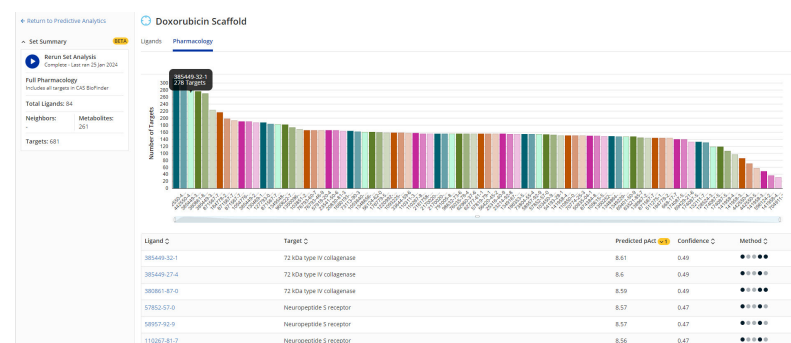– CAS employs a global team of hundreds of scientists and technologists to accomplish this

CAS
A division of the
American Chemical Society

# Drug-target activity prediction

## Estimating the pActivity of ligands towards protein targets



Similarity-based method (**SIM**)

Cross-Pharmacology Index (**XPI**)

Machine Learning Method (**MLM**)

Similarity ensemble approach (**SEA**)

Simplest Active Substructure (**SAS**)

Estimated pActivity for a given ligand-target pair obtained by aggregating the outputs of an ensemble of individual models into a **single consensus value**

# The difference harmonized data makes

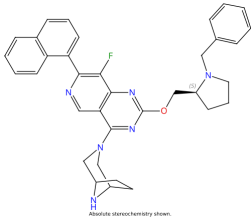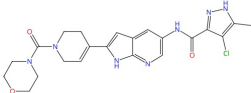## Higher quality data lead to higher quality results

- CAS has an extensive set of harmonized bioactivity data

- This harmonized data was used to retrain existing models

**56%**

Reduction of difference in Predicted vs. Experimental

**23%**

Reduction of Std. Dev. of Predicted vs. Experimental

| Ligand | Target | Experimental pActivity | Original Model | Retrained Model |
|---|---|---|---|---|
|  Absolute stereochemistry shown. | KRas | 4.2 | 11.1 | 4.8 |
|  | KIT | 10.0 | 4.4 | 9.9 |

# Effective partnering for scientific AI solutions

## Requires the "Triangle for Success"



**Domain Expert**
Deep understanding of the science and workflows
– Brings expertise as well as real-world data

**Tech/Algorithm Expert**
Mathematical and computational capabilities
– Brings expertise as well as unique technologies

**Content Expert**
Expertise in data modeling, curation, harmonization
– Brings expertise, data, and information management tech

# Thank you

Connect with us at cas.org

in linkedin.com/company/cas          𝕏 @CASchemistry



**Jacob Al-Saleem PhD**
JAl-Saleem@cas.org

CAS
A division of the
American Chemical Society