

Workshop on Synergizing Biomedical Ontologies 2021

CDD Annotator and Perspectives from the Data FAIRy Initiative

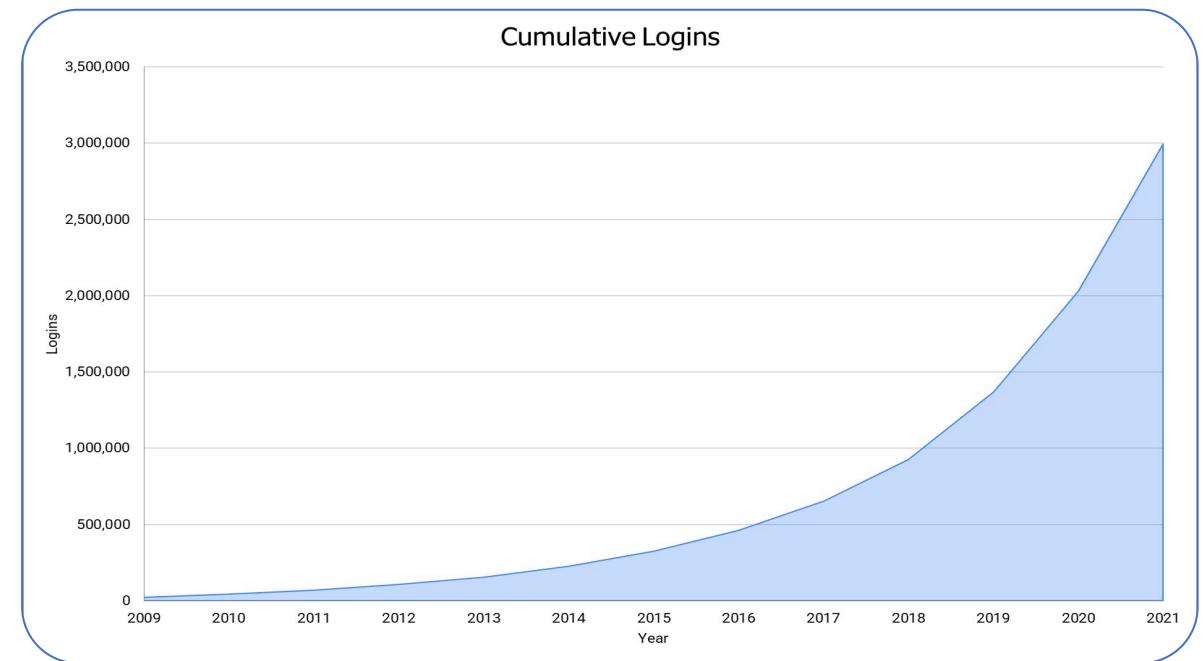
Samantha Jeschonek, Ph.D.

Collaborative Drug Discovery

15 July 2021

Collaborative Drug Discovery

Started in San Francisco in 2004, CDD provides cloud-based data management and discovery platforms to the scientific community worldwide. Our software engineers and application scientists have a rich variety of backgrounds and expertise in many fields, which they leverage to support your modern drug research informatics needs.



CDD's platforms for (meta)data management





CDD.VAULT®
Complexity Simplified™

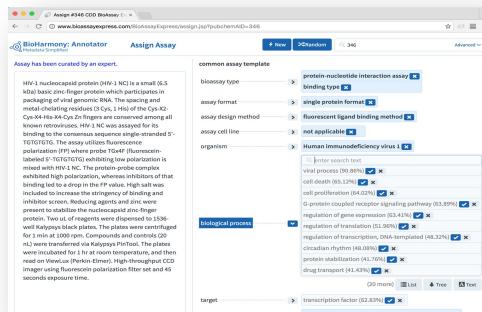
The industry's most trusted cloud-based drug discovery data management platform





BioHarmony
DATA SIMPLIFIED

Your centralized source for semantic drug data



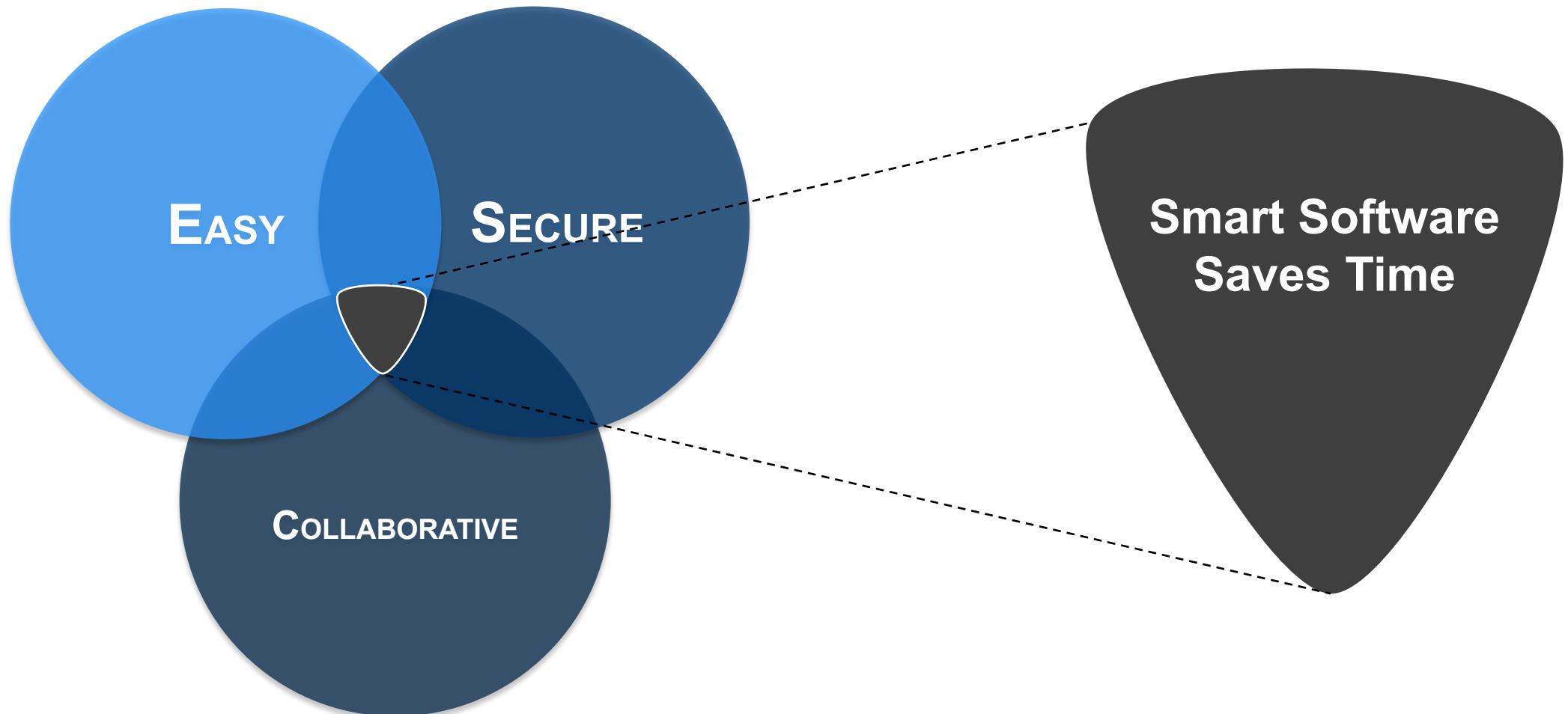


Annotator
METADATA SIMPLIFIED

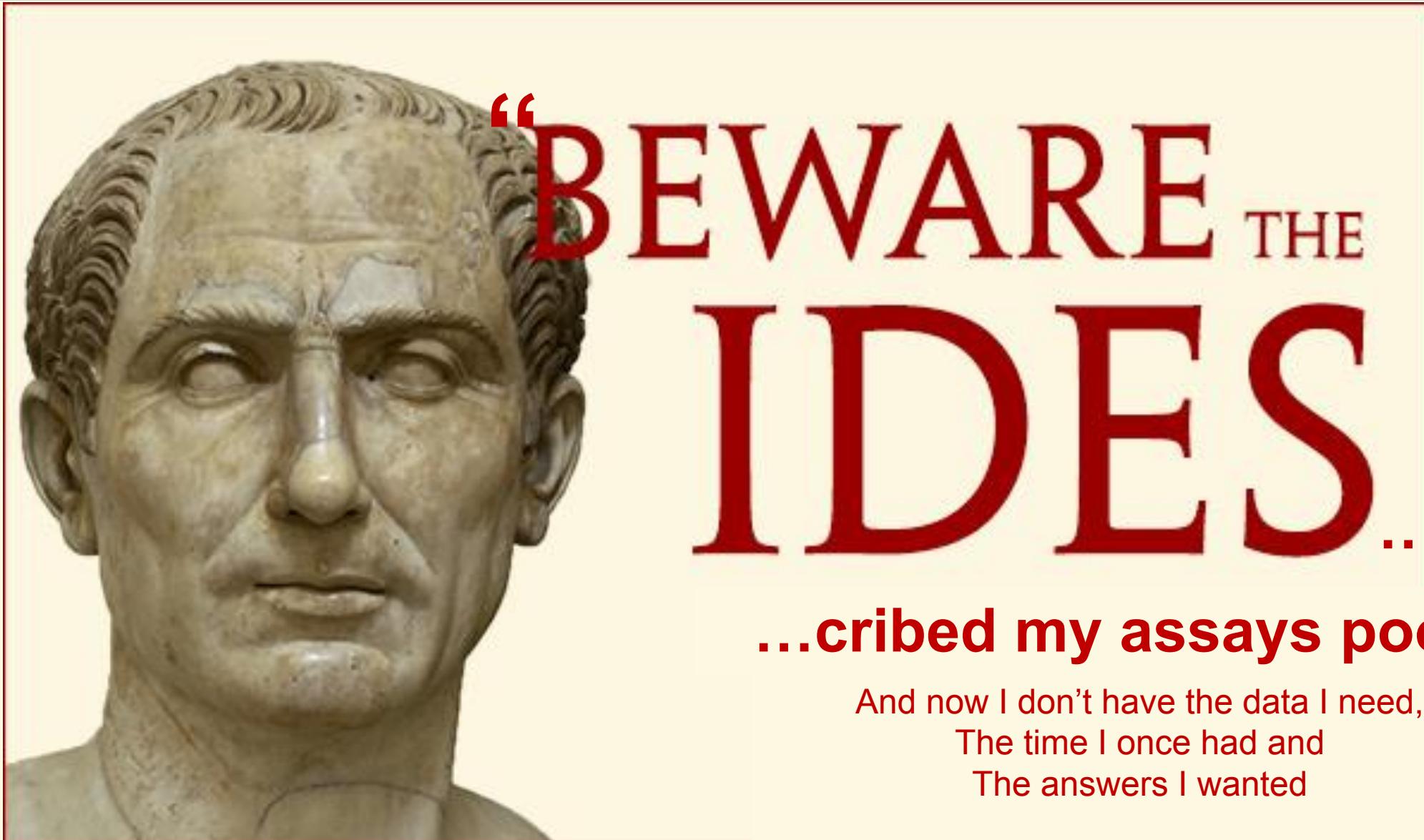
Automatically converts your assay metadata into semantic content

Management
Discovery
Annotation

CDD's principles to save you *TIME!*



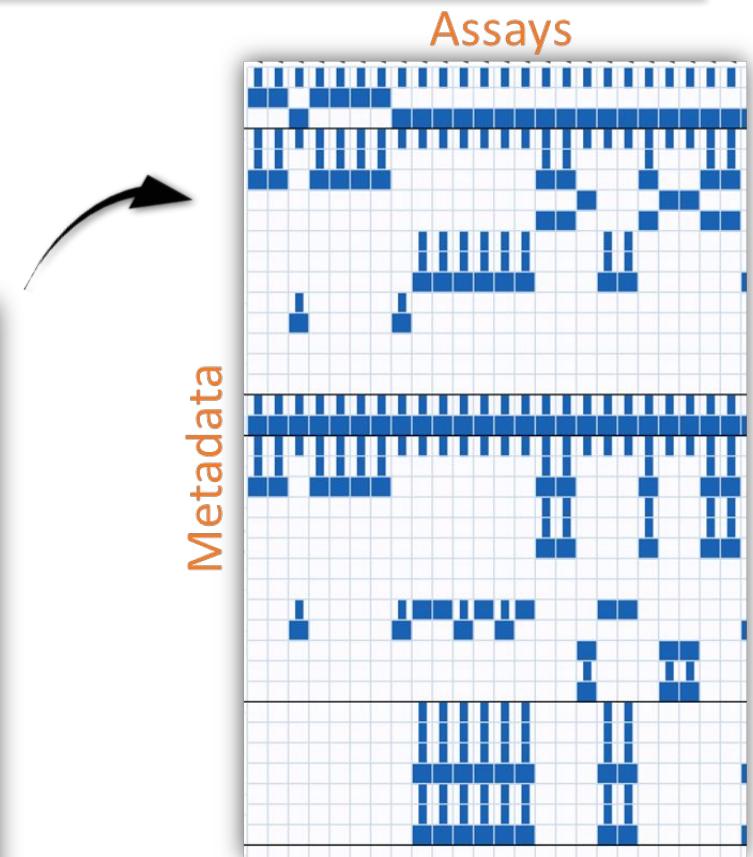
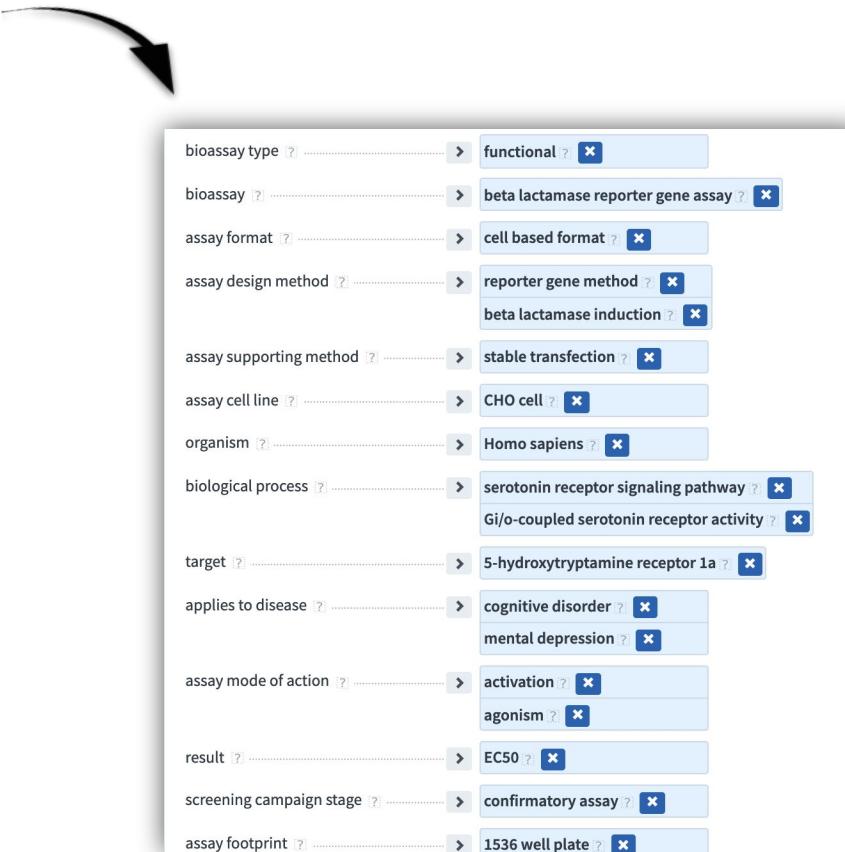
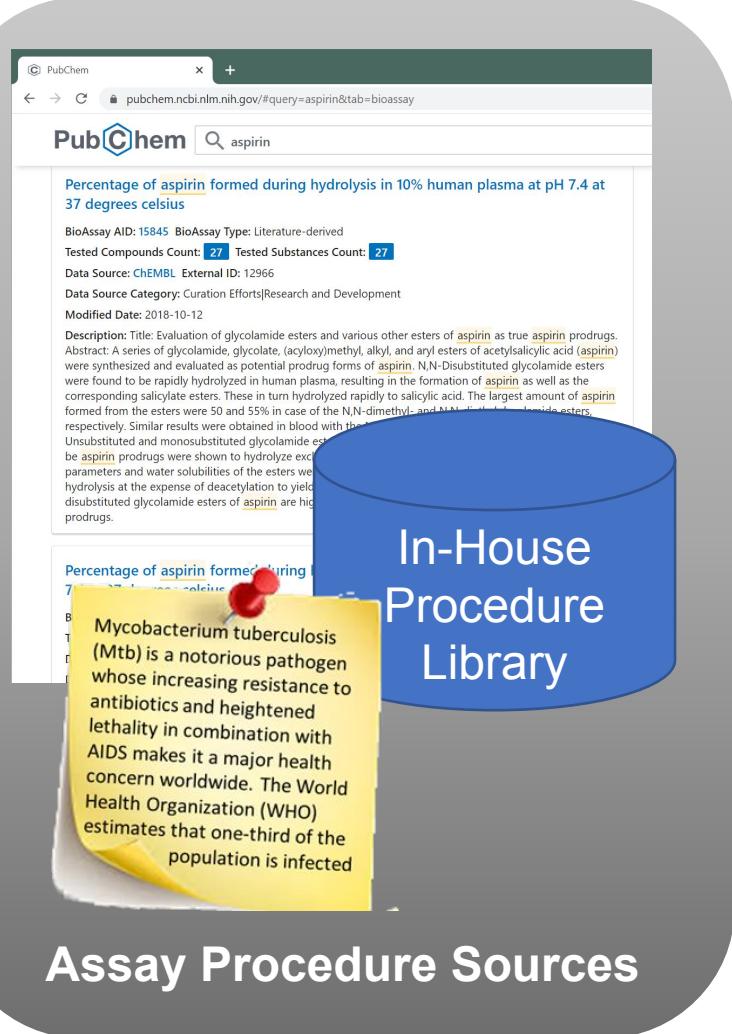
Are you capturing enough data today to answer *tomorrow's* questions?



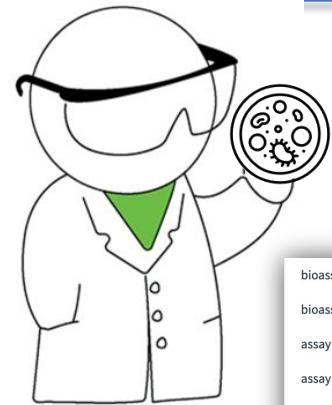
And now I don't have the data I need,
The time I once had and
The answers I wanted

"

Annotator (previously BioAssay Express) enables true assay informatics



Annotator also enables new assay registration



bioassay type ? > functional ?

bioassay ? > beta lactamase reporter gene assay ?

assay format ? > cell based format ?

assay design method ? > reporter gene method ? beta lactamase induction ?

assay supporting method ? > stable transfection ?

assay cell line ? > CHO cell ?

organism ? > Homo sapiens ?

biological process ? > serotonin receptor signaling pathway ? Gi/o-coupled serotonin receptor activity ?

target ? > 5-hydroxytryptamine receptor 1a ?

applies to disease ? > cognitive disorder ? mental depression ?

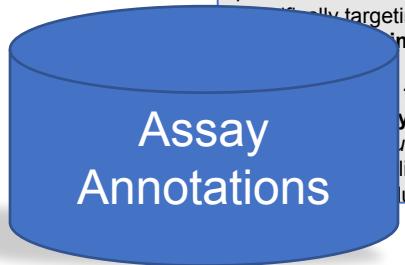
assay mode of action ? > activation ? agonism ?

result ? > EC50 ?

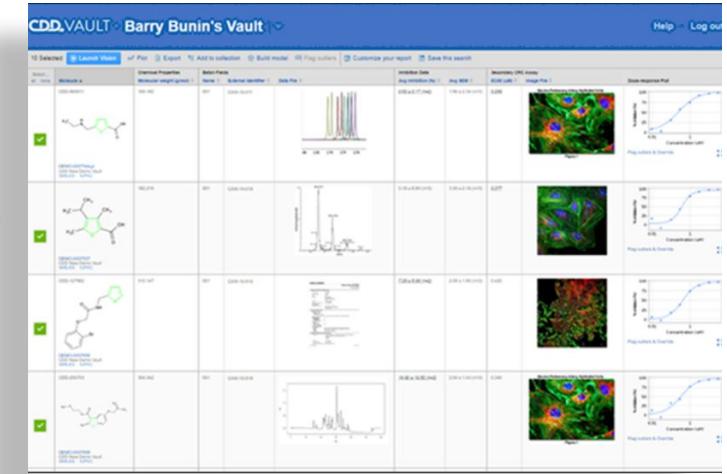
screening campaign stage ? > confirmatory assay ?

assay footprint ? > 1536 well plate ?

Assay Registration



Human &
Machine-Readable
Content



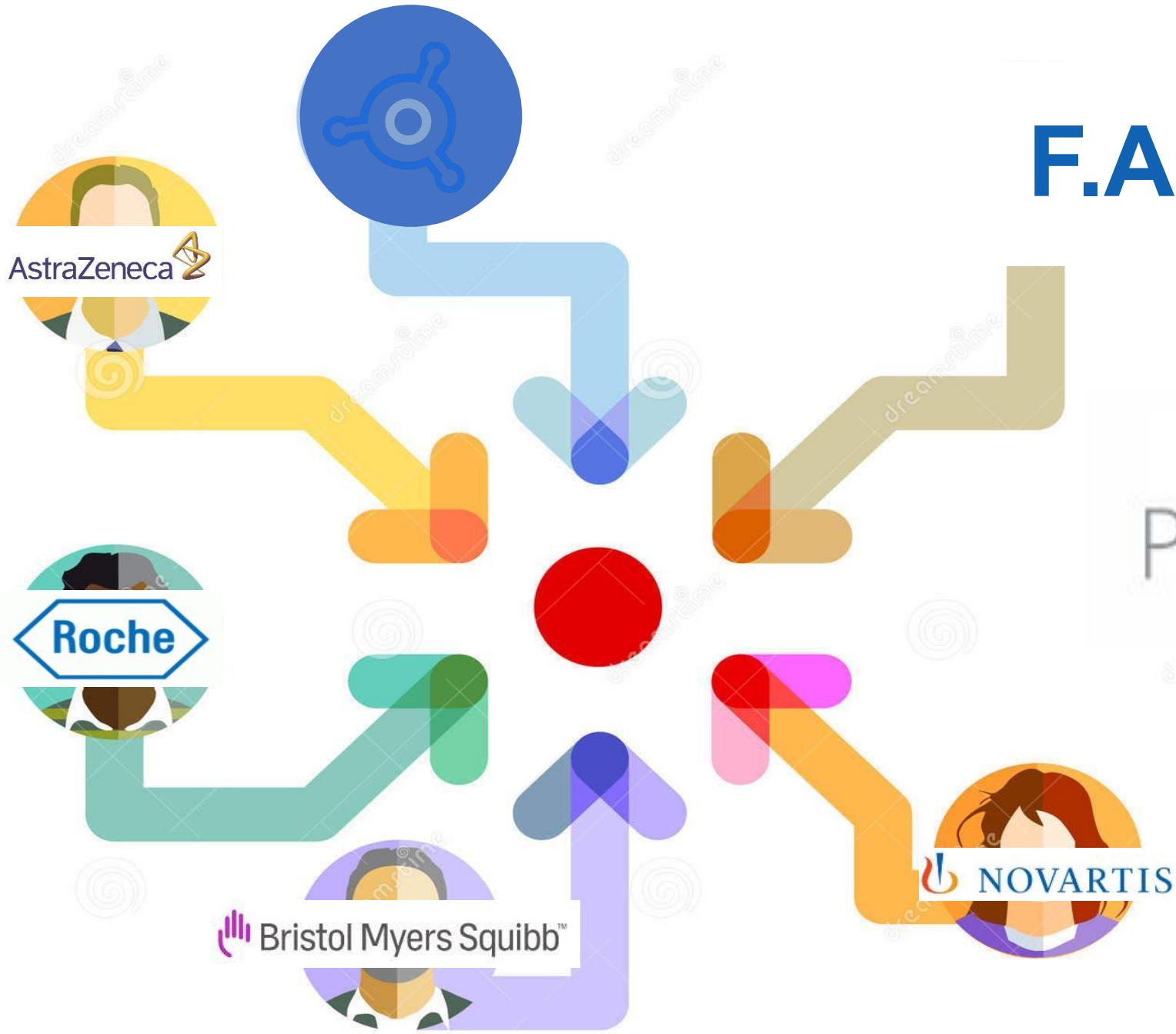
Project Data Analysis

Auto-Generate Text Description

This is a **confirmatory assay** to identify potential treatments for **human immunodeficiency virus infectious disease**, by investigating the biological process of **viral RNA genome packaging**, specifically targeting **see Gene human immunodeficiency virus**

type/protein-nucleotide binding in a single protein **fluorescent ligand binding** line not applicable was conducted in **1536 well plates**

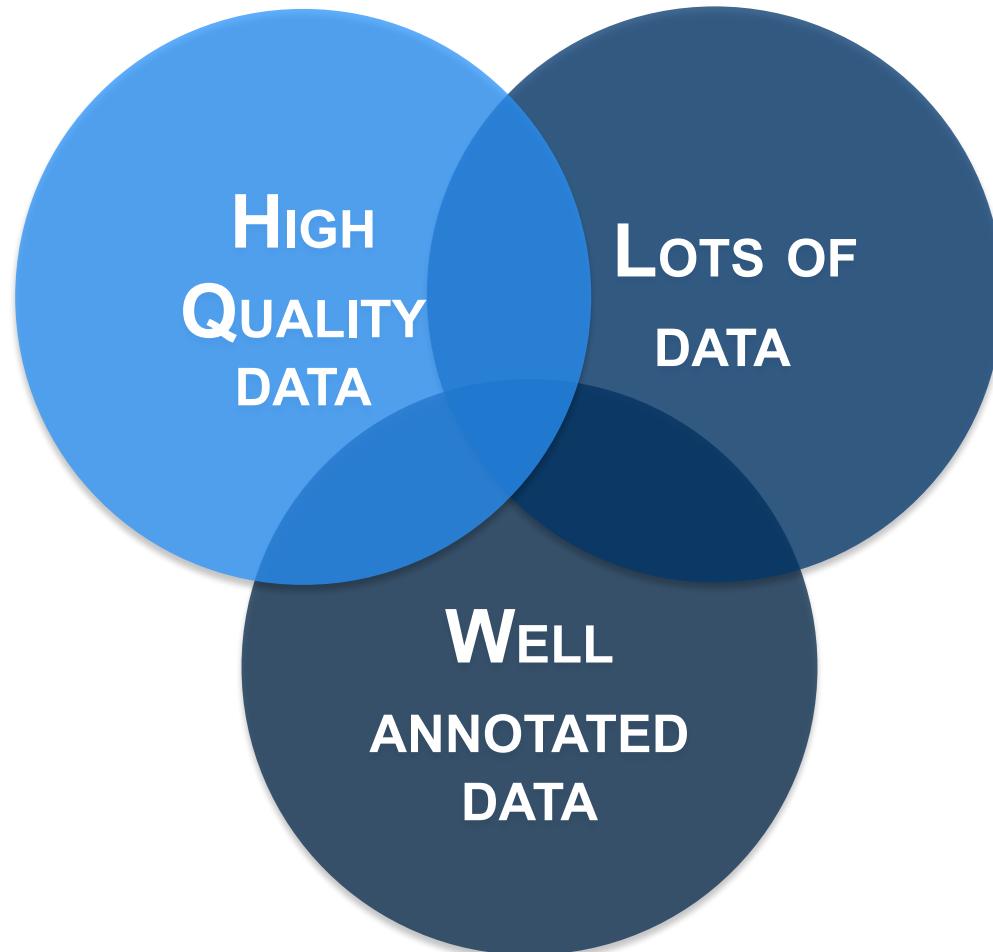
F.A.I.R. Mindsets



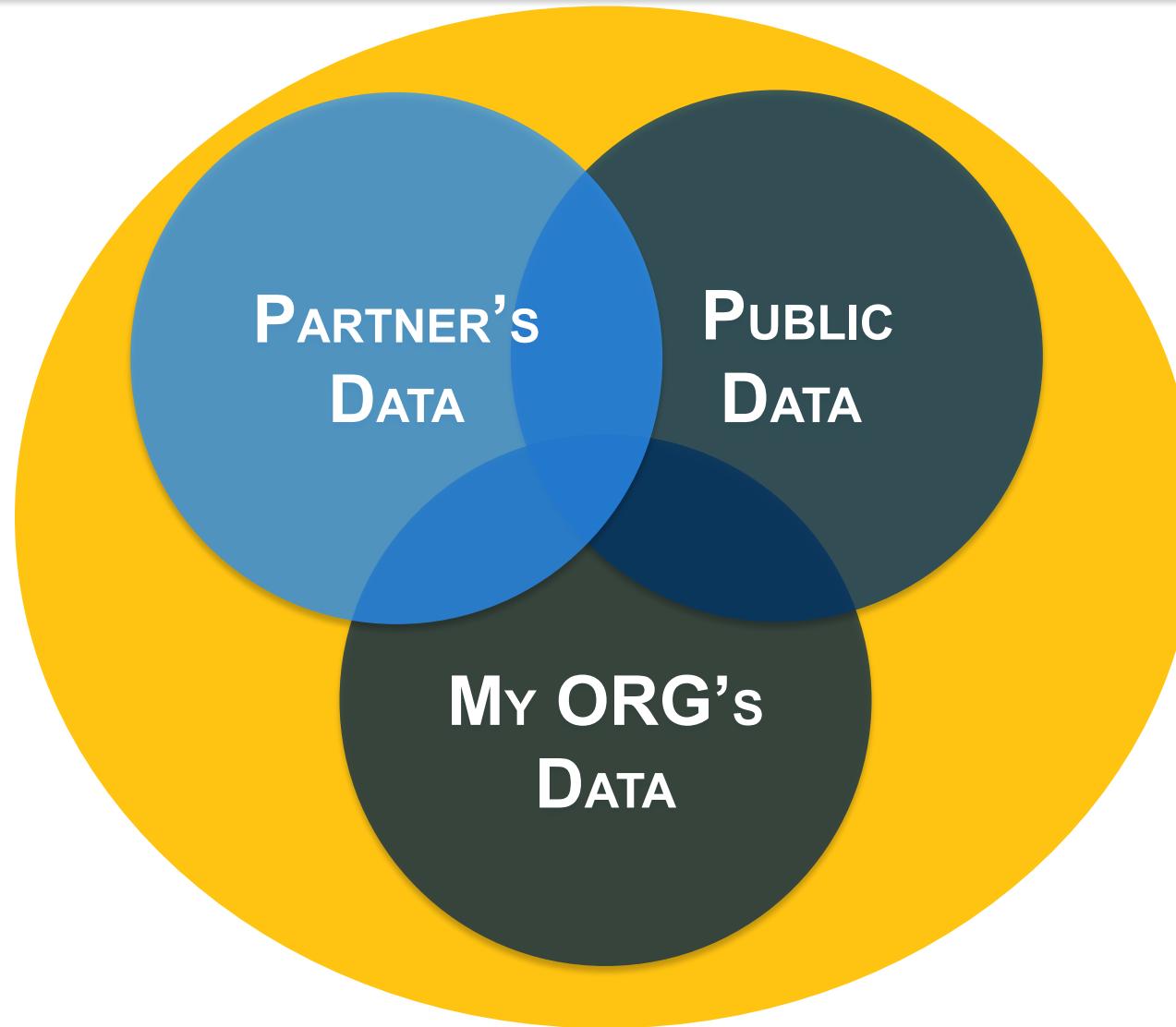
Pistoia Alliance



Community Needs



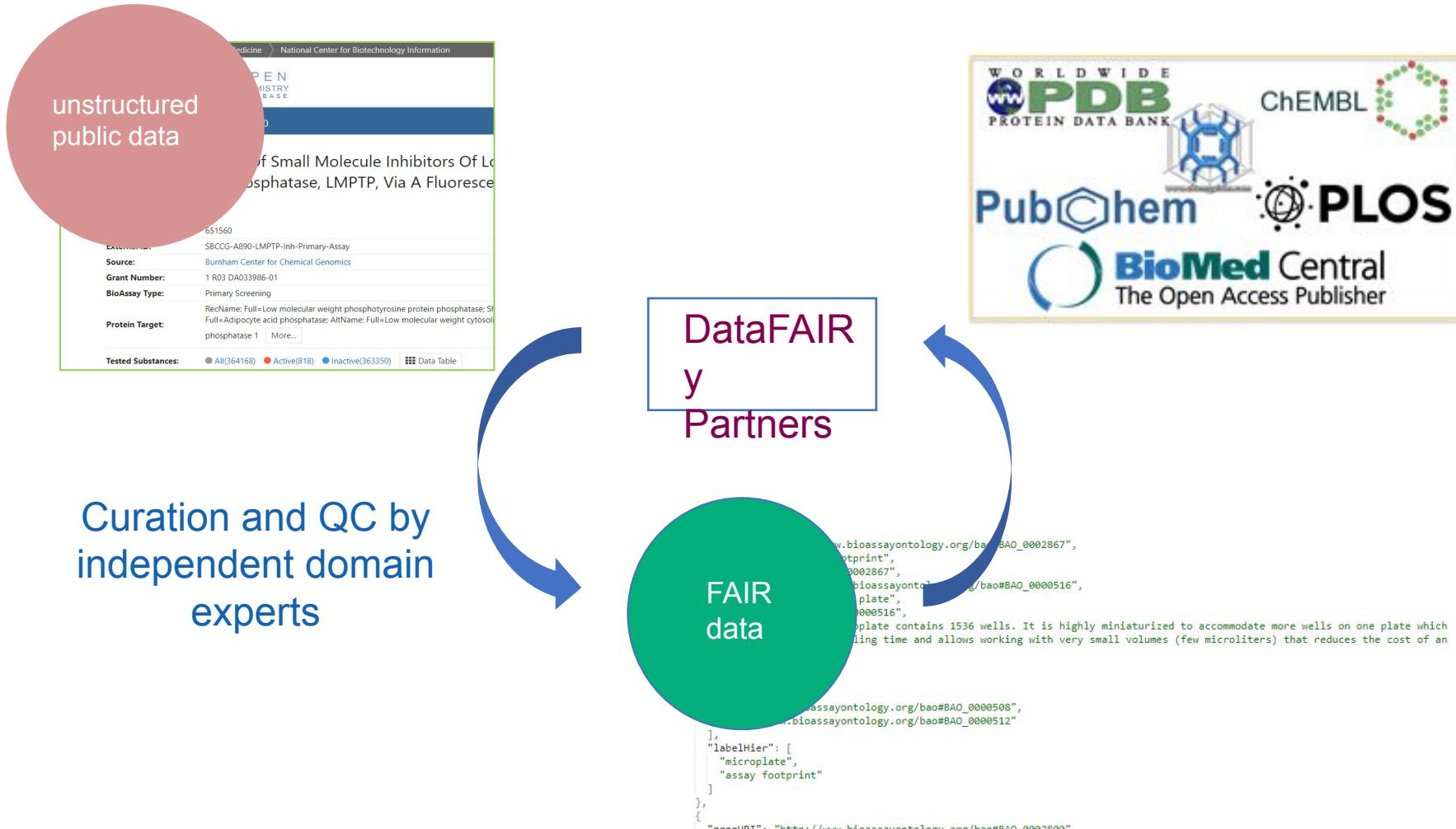
Siloed Data is Not Helpful



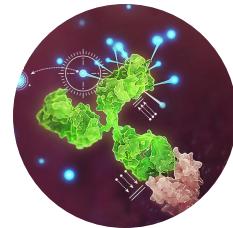
The proposed DataFAIRy operational model (2018)

A vision by Isabella Feierberg...

Cost-shared annotation of public domain bioassay descriptions with high quality, using an agreed data model, making data FAIR

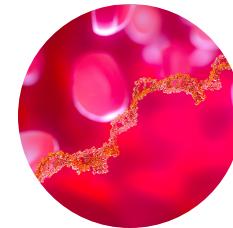


Project team



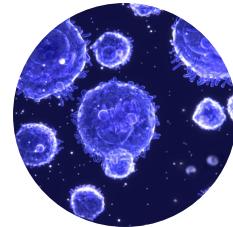
Novartis

Anosha Siripala
Gabriel Backiananthan



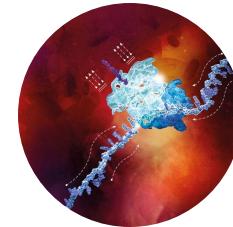
AstraZeneca

Tim Ikeda
Isabella Feierberg*



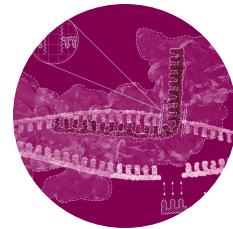
BMS

Dana Vanderwall*



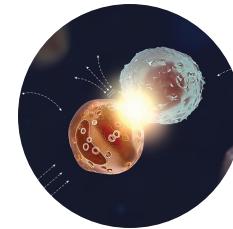
Collaborative Drug Discovery

Samantha Jeschonek
Jason Harris
Whitney Smith



Roche

Rama Balakrishnan
Martin Romacker



Pistoia Alliance

Vladimir Makarov
Thomas Liener





Create FAIR data
with intention



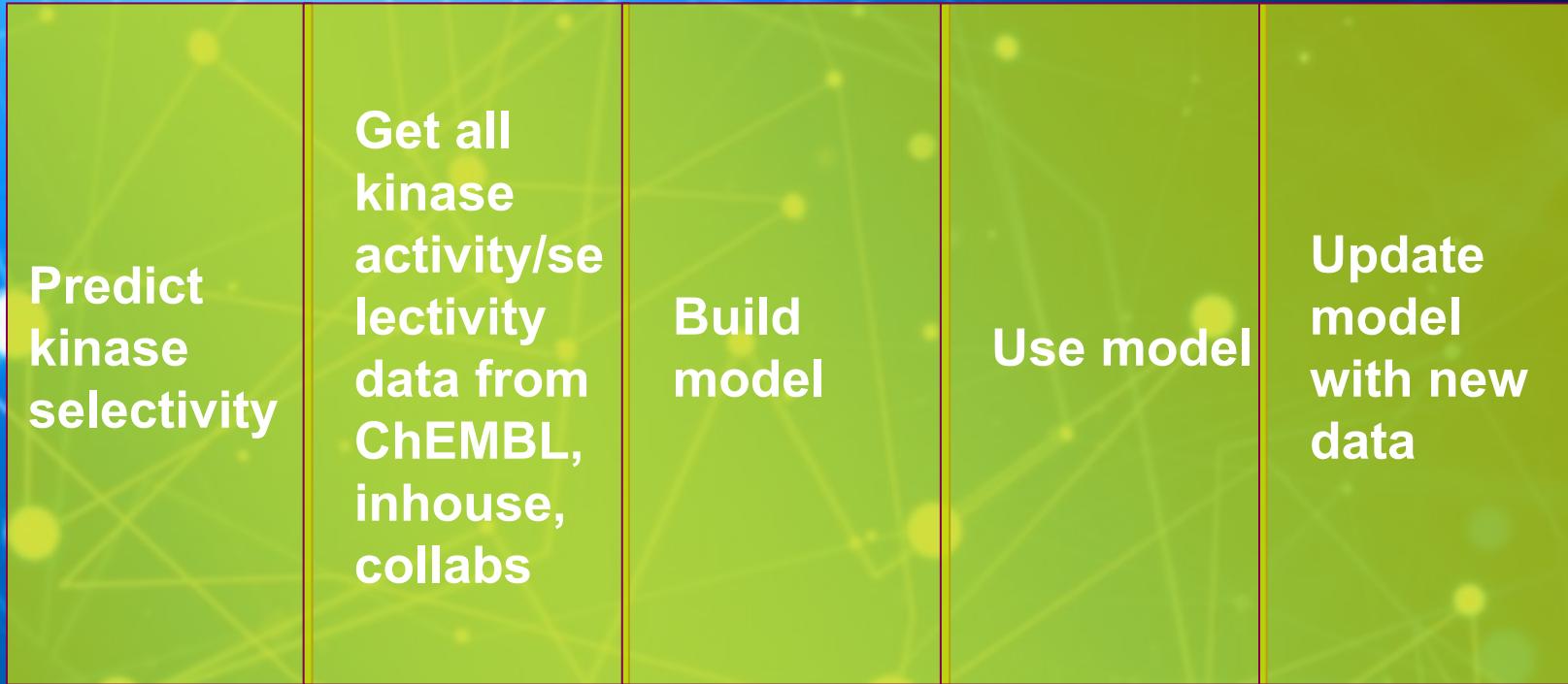
Pilot Project - Business questions

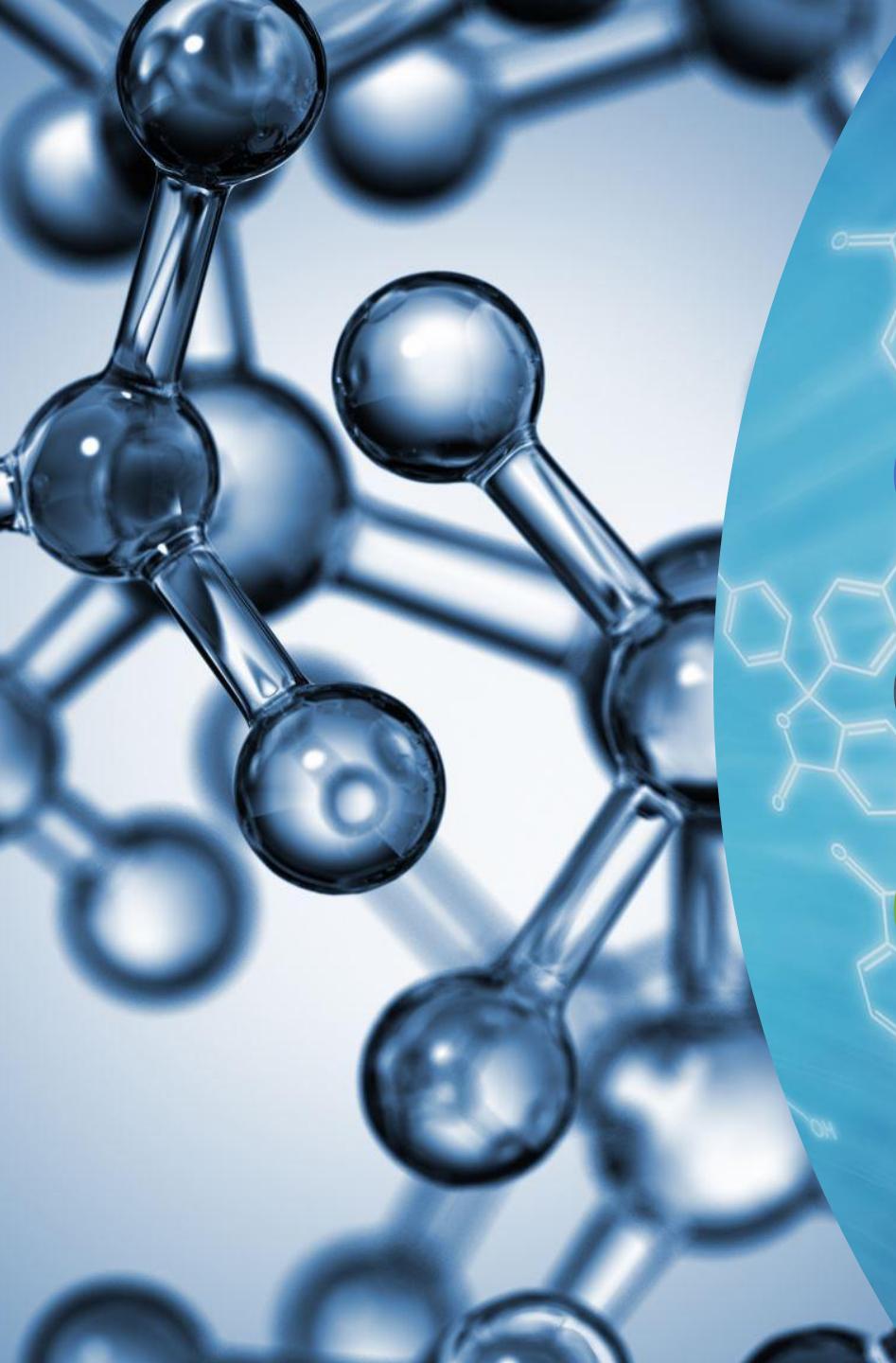
26 initial questions, pruned down to 15, across 5 main categories

- 1 Biology oriented literature mining for discovery project planning
- 2 Assay technology oriented
- 3 Chemistry/tool compound oriented
- 4 Specific assay conditions
- 5 Computational chemogenomic modelling
(e.g., target activity, "PAINS")



Example





Pilot project (2020) – Summary

- Feasibility study, guidance for a larger initiative, example creation
- Curation of 496 public domain assay descriptions were converted into FAIR information objects **using an agreed data model**, which was guided by jointly defined business questions. Upload of the metadata to PubChem.
- Learning points were captured along with recommendations for future endeavors



Bristol Myers Squibb™



AstraZeneca

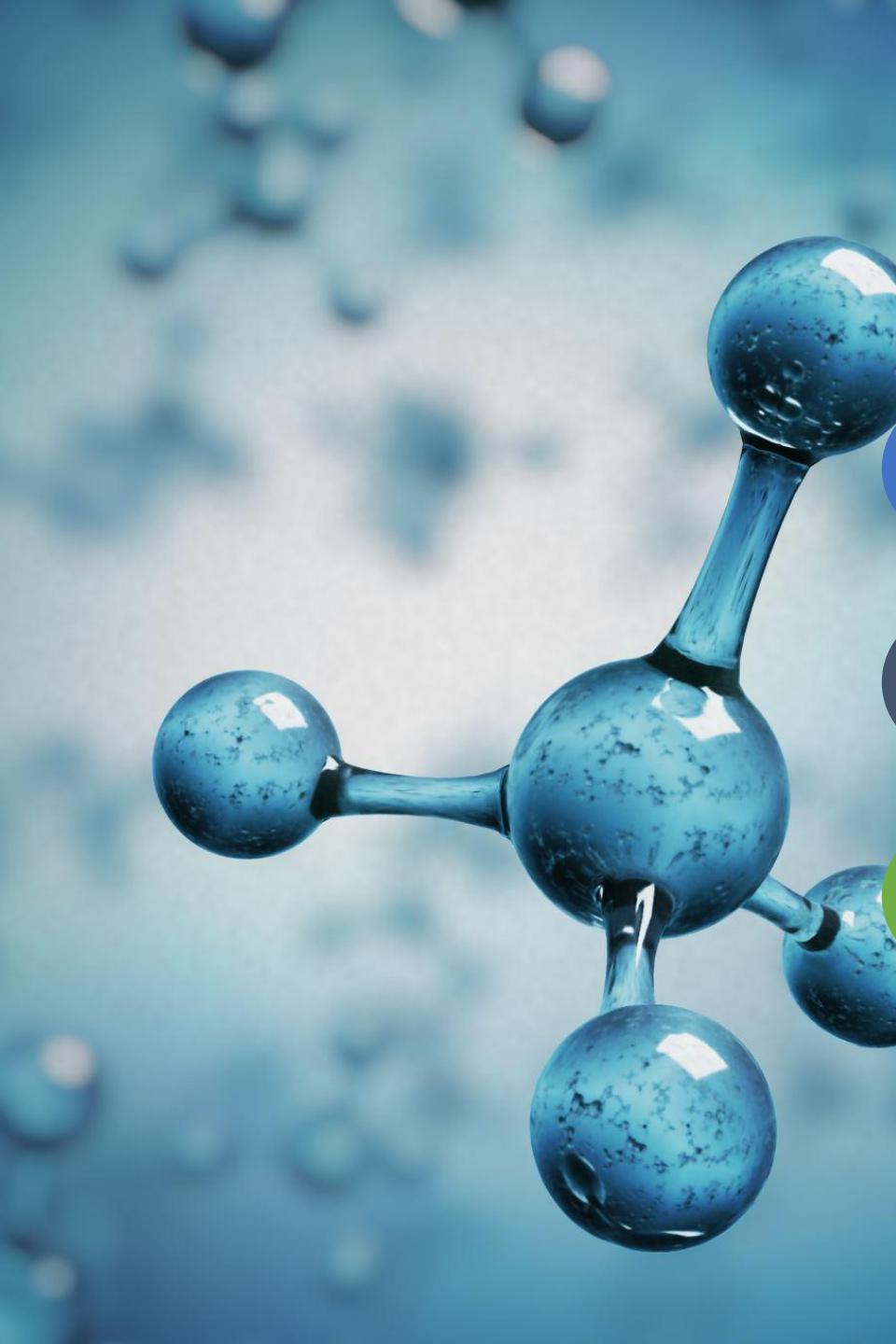


NOVARTIS



Roche





Pilot Project – Assay Selection

1

245 Commercial panel assays: ThermoFisher's kinase selectivity Z'-lyte panel

2

42 PubChem NCATS assays – qHTS, large datasets

3

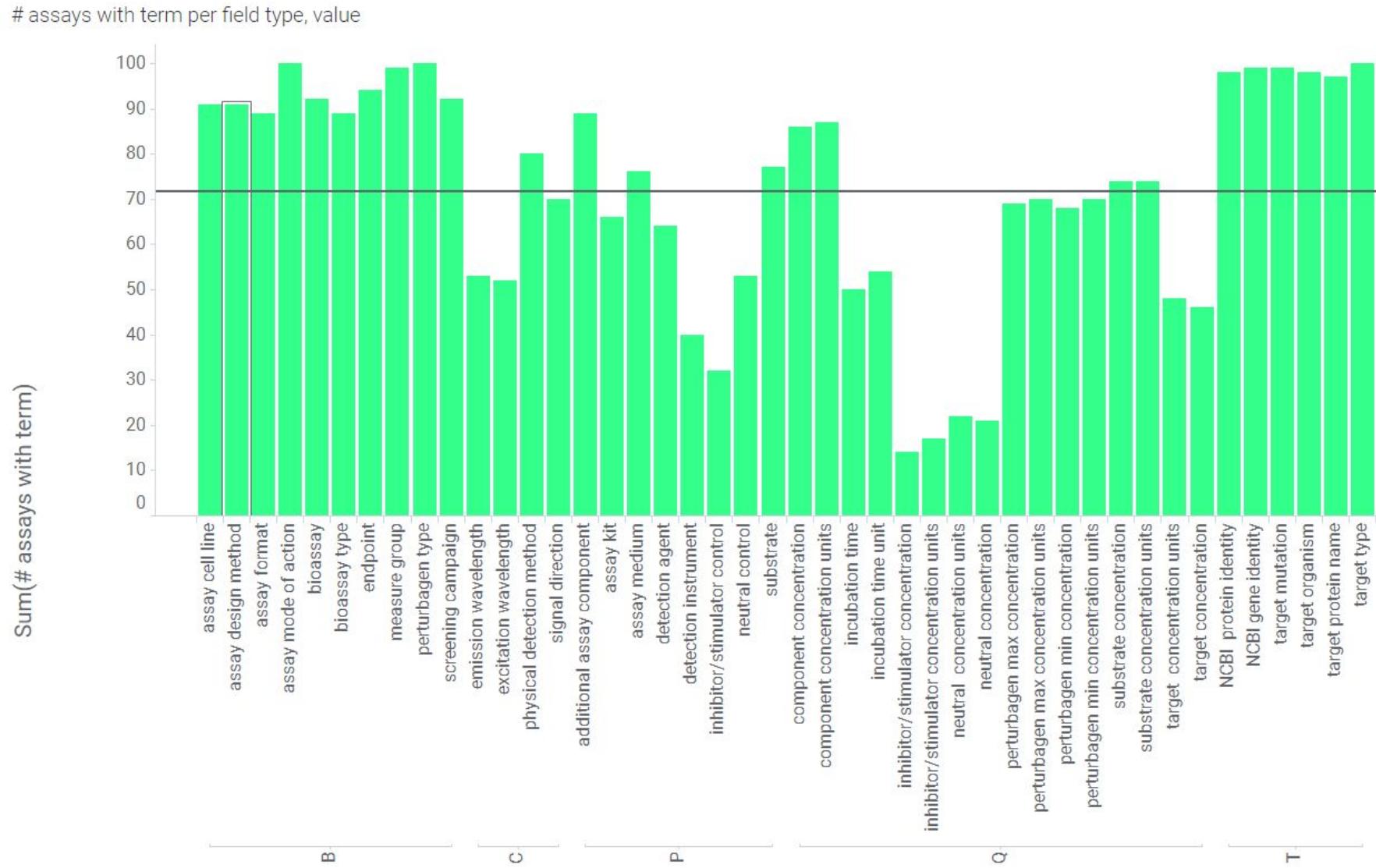
210 publication assays: ChEMBL assays where the target is EGFR, and the reference is Open Access

100 were subjected to manual QC by project team members



Annotator
METADATA SIMPLIFIED

How well did the pilot assays get annotated?



Pilot Project – Learnings

1

Review of supplements & citations = High cost. Choose assays wisely.

2

Commercial assay panels were the easiest to annotate (low-hanging fruit)

3

No persistent links exist for commercial assay panel protocols

4

Errors propagate between papers

5

Fully automated is not fully accurate: Benefit from good work practices: audit trail, versioning, iterative QC

6

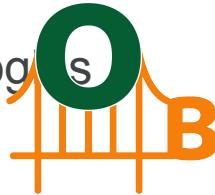
Need for published assay protocols to be well-annotated in public databanks and linked to the publication

7

Need for a common community data standard for future assay publications.

8

Need for flexible & synergized Ontologies (e.g.: OntoloBridge)





Next steps:

1

Scale up (x 10-100) in next steps. Having more partners lowers cost per partner per assay and overhead cost

2

Optimize process, data sources, tools, QC within quality constraints. Define quality metrics.

3

Define and promote a **community standard** for assay reporting and publishing -- **align** with vendors, publishers, government agencies.

4

Attract new project members and sufficient funding to start the next phase

Thanks to

AstraZeneca

Nigel Green
David Hayes
Tom Plasterer

BMS

Rick Bishop

Janssen

Herman van Vlijmen

Novartis

Fabien Pernot

MMV

Jeremy Burrows

PubChem

Evan Bolton

ChEMBL

Anna Gaulton
Andrew Leach

Roche

Olivier Roche

Medicines Discovery

Catapult

John Overington
Mark Davies

Pangeadata.ai

Vibhor Gupta

University of Miami

Stephan Schürer

BioSci Consulting

Scott Wagers

Collaborative Drug

Discovery

Barry Bunin
Frank Cole
Alex Clark
Hande Küçük McGinty
(now Univ. Of Ohio)

Pistoia Alliance

Carmen Nitsche (now CCDC)
Nick Lynch (now Curlew Research)

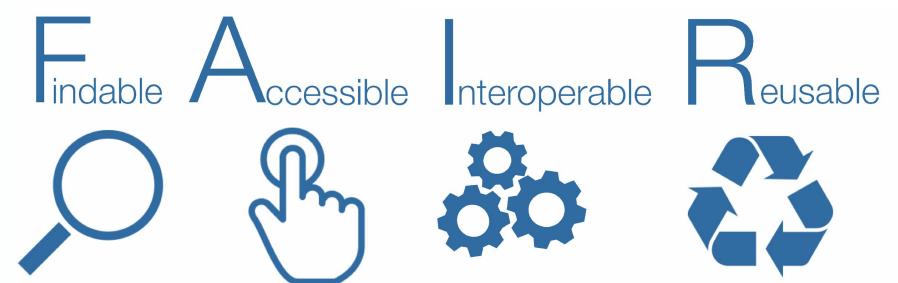
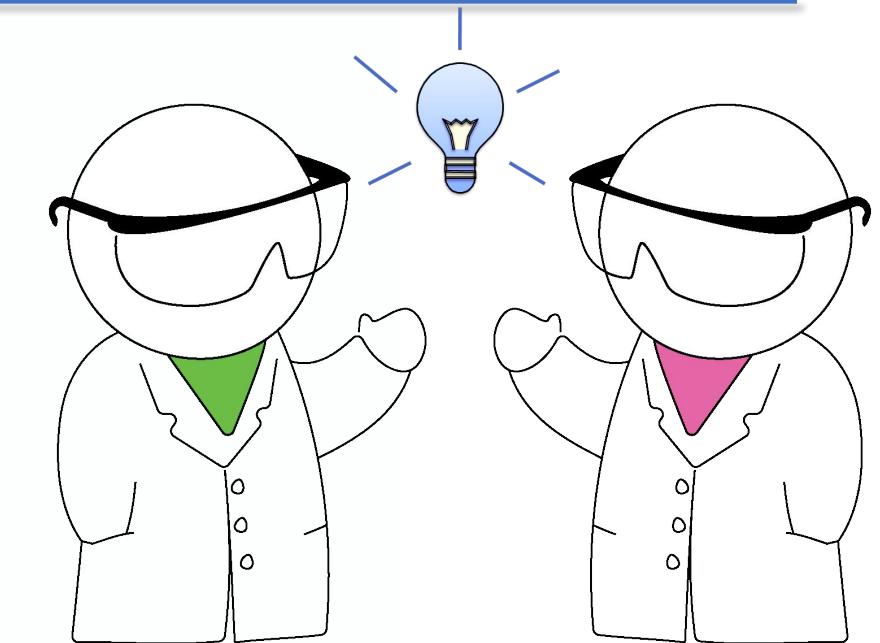
Metadata Managed! – Annotator In Action

Unique ID ▾ assay ID common assay template Template Sections new content

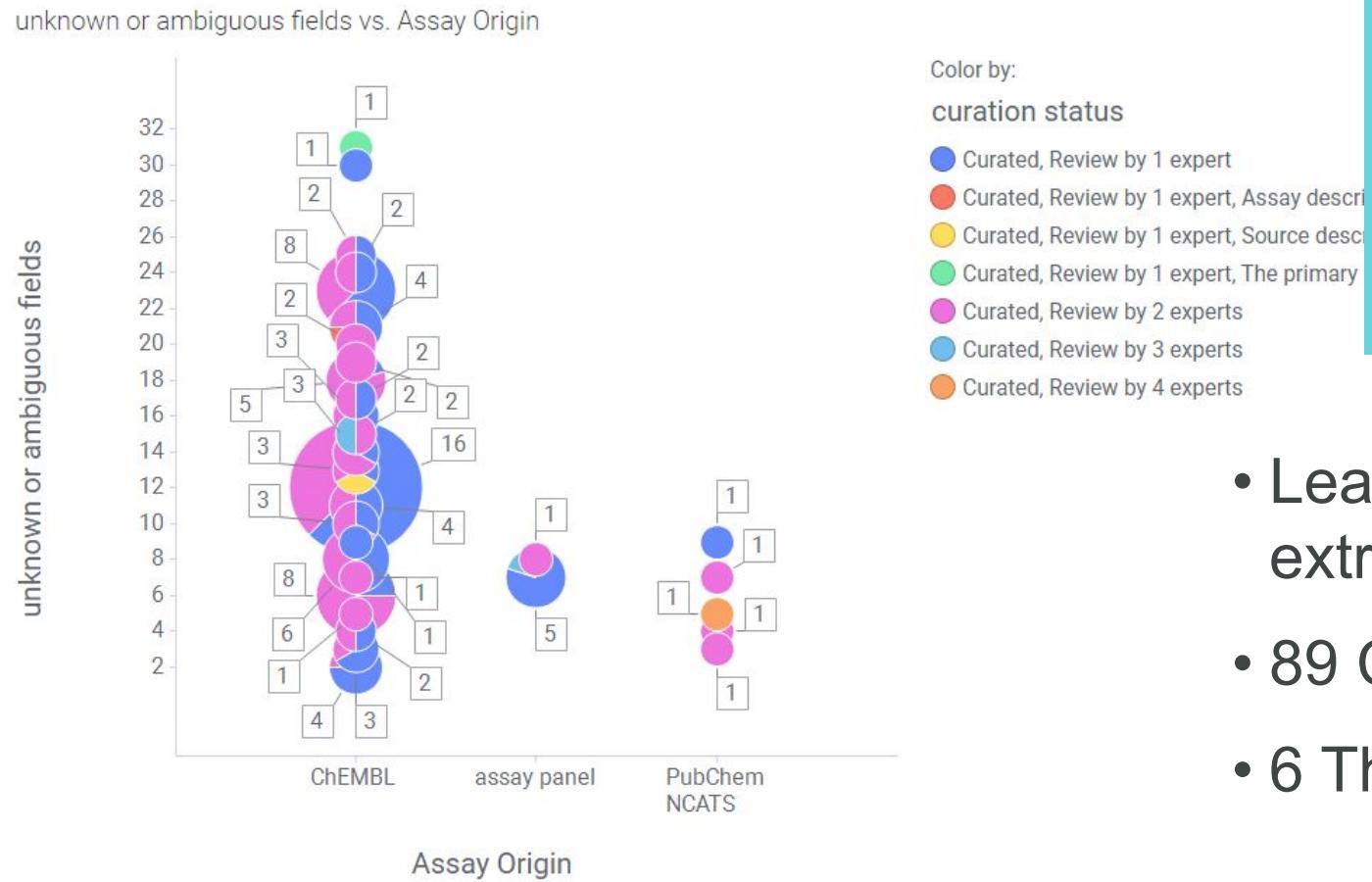
Protocol Full Text

I

	assay title	?	>	<input type="text"/>
	bioassay type	?	>	<input type="text"/>
	bioassay	?	>	<input type="text"/>
	assay format	?	>	<input type="text"/>
	assay design method	?	>	<input type="text"/>
	assay supporting method	?	>	<input type="text"/>
	assay cell line	?	>	<input type="text"/>
	organism	?	>	<input type="text"/>
	biological process	?	>	<input type="text"/>
	target	?	>	<input type="text"/>
	applies to disease	?	>	<input type="text"/>
	assay mode of action	?	>	<input type="text"/>
	result	?	>	<input type="text"/>
	screening campaign stage	?	>	<input type="text"/>
	assay footprint	?	>	<input type="text"/>

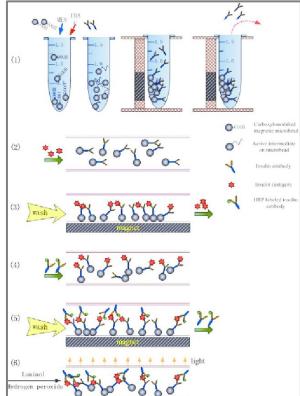


Pilot: 100 QC:d assays (~20%)



- Learning points are largely extrapolating on the 100 QC:d assays
 - 89 ChEMBL assays, 5 NCATS assays
 - 6 ThermoFisher panel assays QC:d

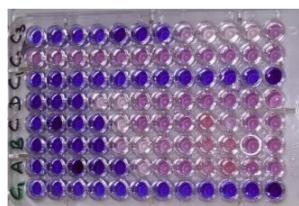
Are you capturing enough data today to answer tomorrow's questions?



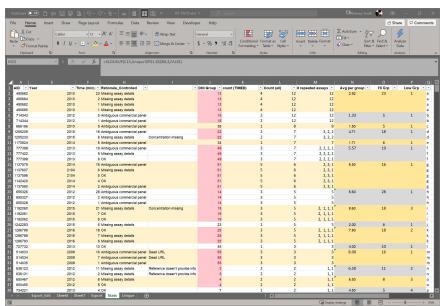
Procedures



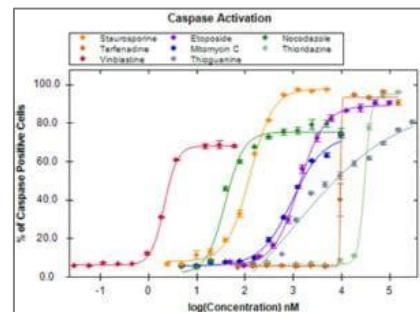
Instrumentation



Samples & Results

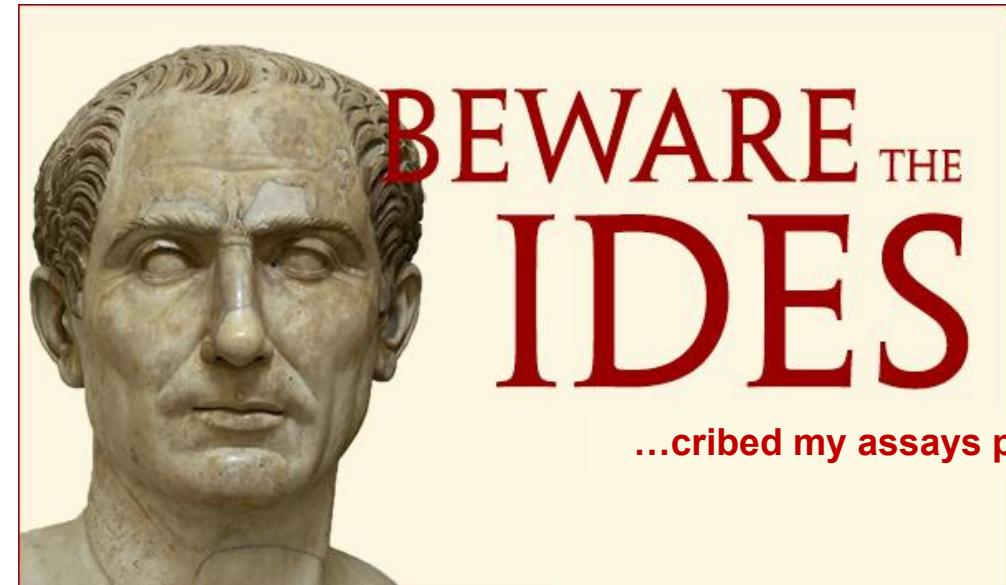


Raw Data & Observations



Data Analysis & Assumptions

Compound	Result (mM)	Note
XYZ-0001234	12	Added too much Buffer 1
XYZ-0001235	105	
XYZ-0001236	>200	Didn't look good
XYZ-0001237	<5	
XYZ-0001238	<5	
XYZ-0001239	>200	Looked bad
XYZ-0001240	55	
XYZ-0001241	56	
XYZ-0001242	20	
XYZ-0001243	25	Sent to Fred for follow-up (New Fred, not Old Fred)





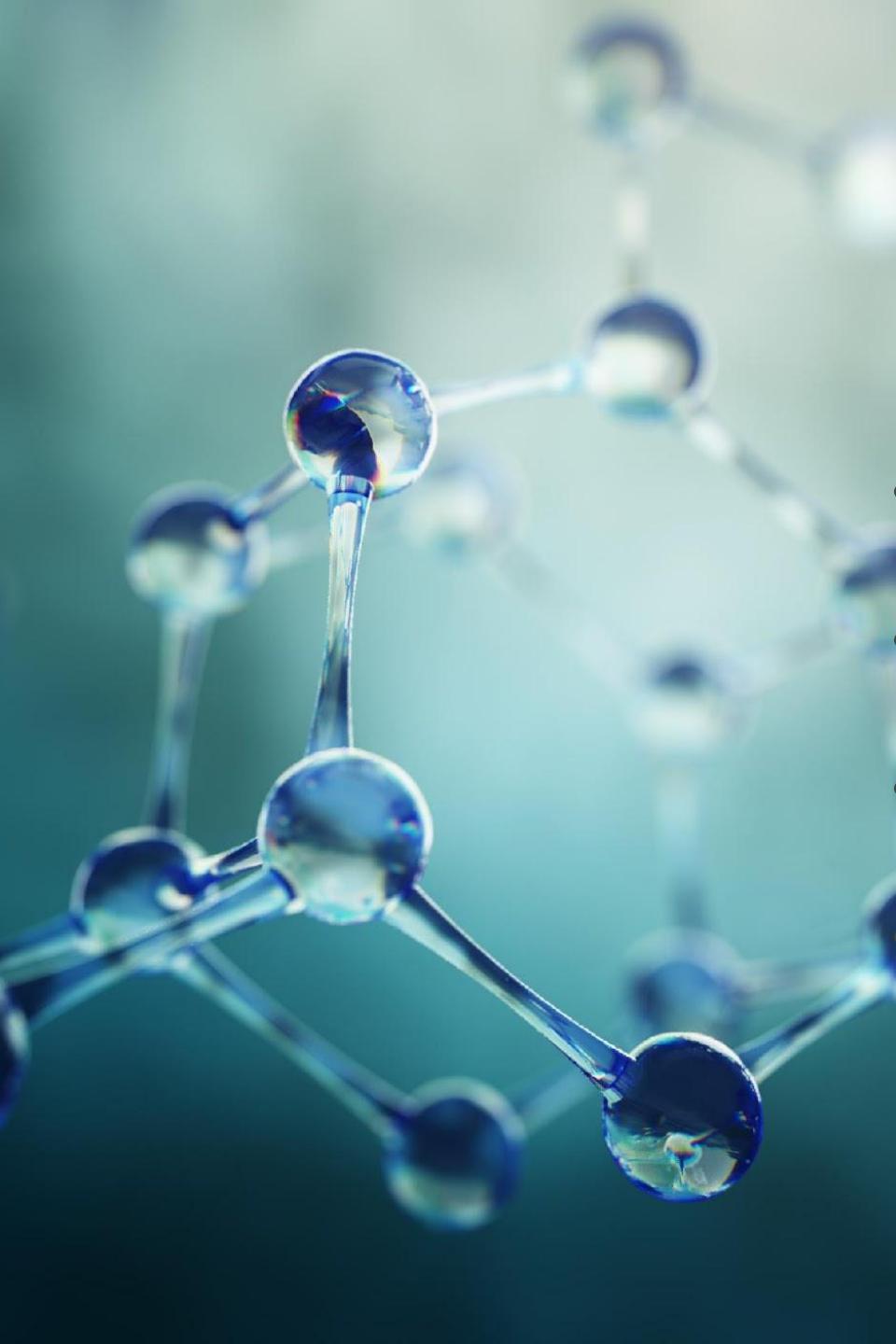
Why DataFAIRy?

There is a need for FAIR public domain data with high quality annotations using public ontologies and a common data model

Available metadata in (public domain) data repositories is often insufficient for answering current and future business questions

Substantial investments are being made in AI, ML and FAIR data across life science industry and academia

Pharma companies already pay for curation of partially overlapping public domain data (e.g., ChEMBL, papers, chemistry patents)



Small molecule bioassays make up a good pilot case

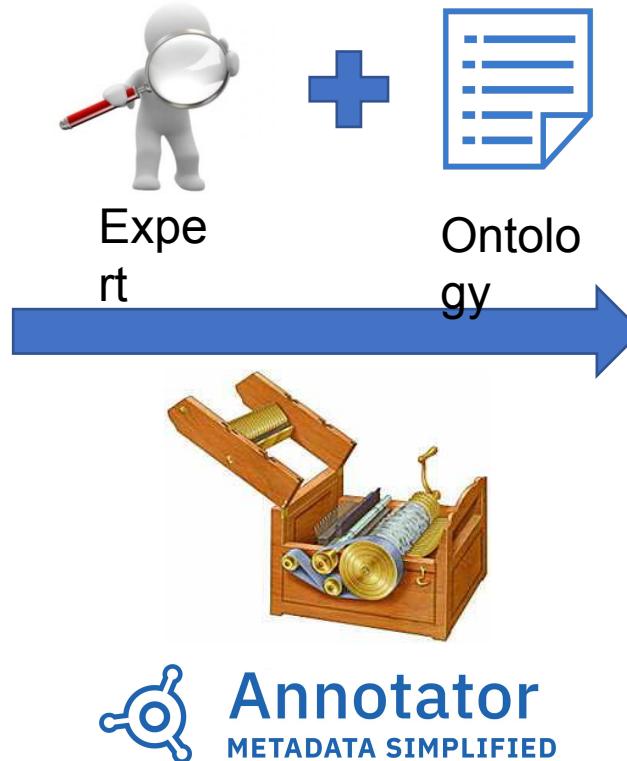
- Project planning – what is available in the public domain?
- Assay development, e.g., assay conditions and tool compounds
- Chemogenomic model building
- Enriching public chemogenomics data with FAIR metadata will show impact across the cheminformatics domain

What we want



A Generic Ontology-Driven Annotation Tool

Ambiguous, Unstructured Data



Unambiguous, Structured Data



Example annotation use cases: Biochemical, cell, phenotypic assays, pharmacology assays, clinical assay data, bicycles...