

# **Leverage the Wealth of Internal and External Information to Drive Collaboration and Project-Centricity into Your Research Informatics Platform for Drug Discovery and Development; A Strategic Imperative**

**BioIT World Conference 2016**

April 5-7, 2016

**James Connelly**

Head of Global Research Data Management,  
**Sanofi**

# Synopsis



Every pharmaceutical company holds a wealth of information collected for numerous drug candidates during research, pre-clinical and clinical development and for marketed drugs.

Much of this information is trapped, undiscoverable and not optimally useable for drug discovery and development or effectively combined with the massive amount of externally available data.

Two pilot studies with IBM Watson successfully extracted high quality information from toxicology reports and enabled researchers to discover critical insight for drug-repurposing proposals.

Cloud-based SAR data platforms will further stimulate the use of broad data sources for research and create a project-centric and collaborative environment.

Increasing reliance of Pharma on a diverse range of drug discovery collaborations has caused severe challenges to adapt Big Pharma data warehouse-centric SAR data platforms.

Evolution and innovation of cloud-based SAR data platforms has revolutionized our approach to SAR data integration, sharing and analysis.

We will describe our current approach to collaboration data and the evolutive roadmap towards a fully capable SAR data platform in the cloud that will utilize cloud-based services and big data technology.

This will reduce the cost of internally supported systems and create a scalable external SAR informatics system.

These are major trends in the pharma industry and are driven by the acceptance, evolution and scalability of cloud-based and computational services and for collaborative SAR data sharing and analysis.

There is a strong value-added in utilization of this wealth of information for R&D and strategic imperative to build informatics systems using Big Data technology to access and utilize all relevant information sources.

- **IBM Watson Drug Phenotype “Information Cloud”**
  - **Drug Safety and Toxicology “Cloud” Pilot**
    - Concept and Execution
    - Data Access and Visualization
- **Project-Centric, SAR Data Platforms “in the Cloud”**
  - **Traditional approaches to harmonize SAR Data**
    - Compounds, Bioassays and ISIS DB’s
    - Integrated SAR Data Platform and DataWarehouse-centricity
    - High Cost of big RED DP integration and support
  - **Back to the future with project-centric DB’s**
    - Drug Discovery Collaboration Data Exchange “in the cloud”
    - Evolution of full single-vendor SAR data platforms
  - **Leverage “Big Data” Technologies to build a distributed SAR data environment**

**IBM Watson**

# Drug Phenotype Information “Cloud”

**April 5, 2016**

## **Sanofi**

James Connelly, Research Director  
David Aldous, Senior Vice President  
Richard Brennan, Director Scientific Advisor

## **IBM Watson**

Louisa Roberts, Scott Spangler and Ying Chen



Create an advanced, next generation R&D Informatics capability that revolutionizes the use of internal and external data content for an information-driven, iterative approach to Drug Discovery and Development

## Information-Driven Approach to

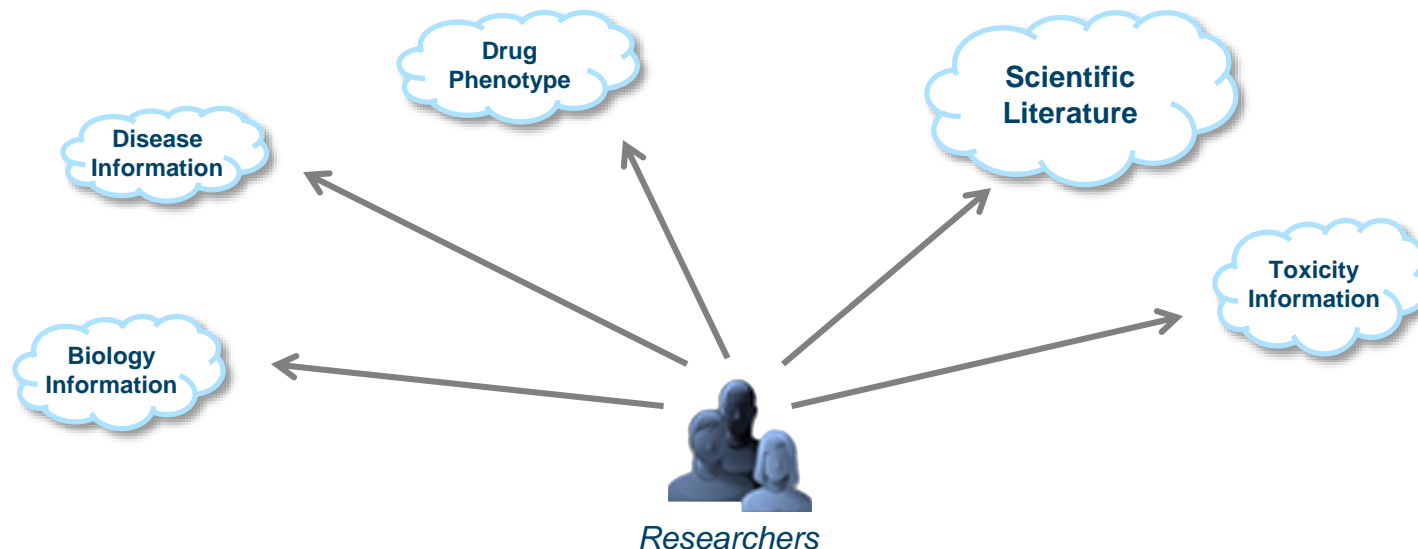
1. Learn from legacy projects & information?
2. Detect liabilities for early resolution or fast fail?
3. Uncover and analyze information in many different silos?
4. Make the right/best decisions around projects?
5. Accelerate programs, decrease attrition and lower costs?
6. Find new indications for existing Sanofi compounds?



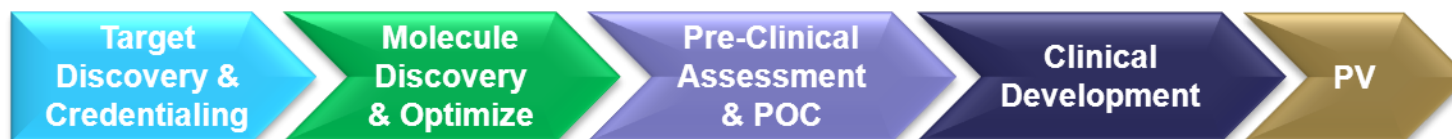
## To Create the Anticipated Impact

- ✓ Shape the future Drug R&D paradigm and profiling strategies to advance superior candidates to the clinic
- ✓ Identify new therapeutic targets, anti-targets, biomarkers, MOA, liabilities and correlation with responders/non-responders in patient populations
- ✓ Reduce R&D cost and accelerate the delivery of innovative therapeutic options for patients.

Sanofi has a wealth of information trapped in reports, presentations and siloed databases that is unavailable or undiscoverable for scientific queries, analysis or prospective insight and hypothesis generation/testing



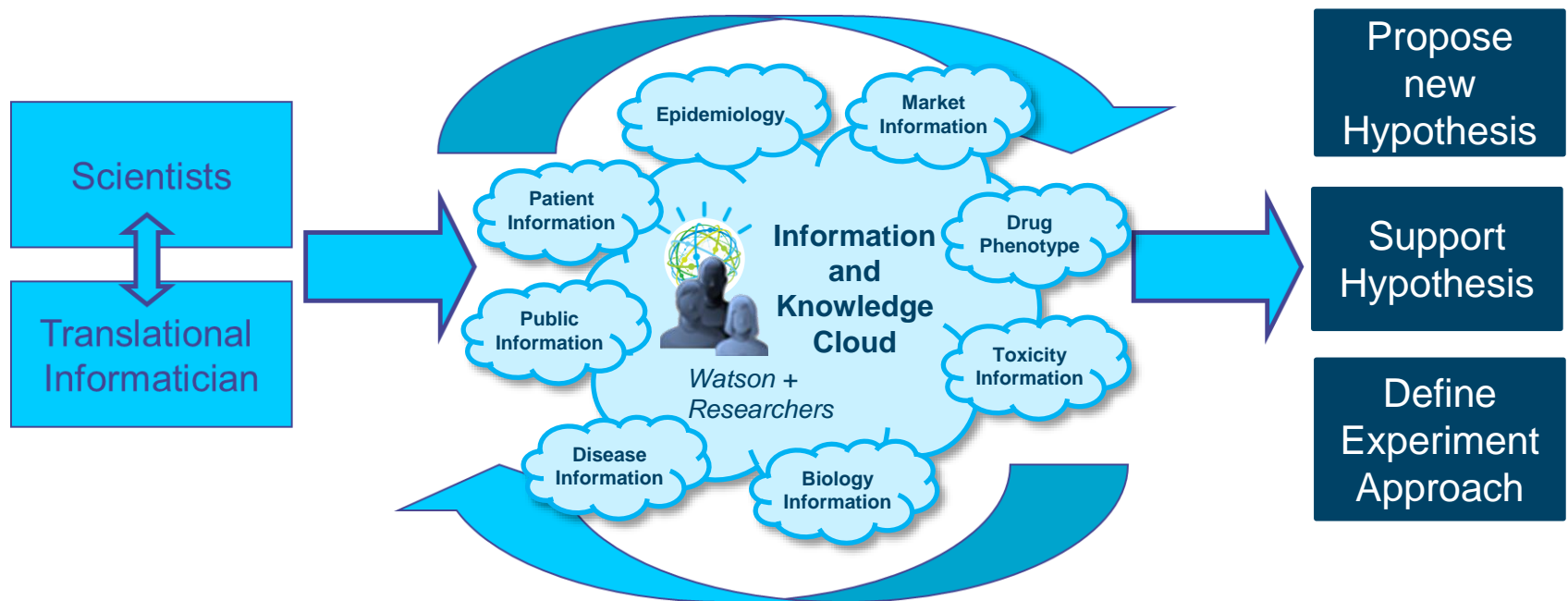
**Drug Discovery has historically worked in a linear model (Shots on Goal)**



- Programs end somewhere on this continuum and new projects are initiated
- Relatively little of the information collected is re-used for learning and insight

## Next-generation R&D Informatics for an Information-driven Era of Drug Discovery Research

*The future is an iterative information-driven model*



The key is information availability, interoperability and usability

## Background

- Information critical to drug discovery research is typically siloed, fragmented, and poorly accessible
- Sanofi partnered with IBM Watson to demonstrate that text analytics, data integration, and machine-learning could enable the extraction of scientific findings from unstructured data
- A pilot study was undertaken to extract and structure preclinical safety data from historical repeat dose toxicology study reports
- And to integrate and visualize these data alongside discovery pharmacology data to facilitate the resolution of compound → target → toxicity relationships



# The Information and Knowledge Cloud will drive value across three key areas

## Efficiency

- Speed up the deep exploration of scientific literature, less reliant on experts
- De-risking programs before significant investment to allow early targeted resolution or fast-fail. Information driven = reduce redundant/expensive experimentation

## Effectiveness

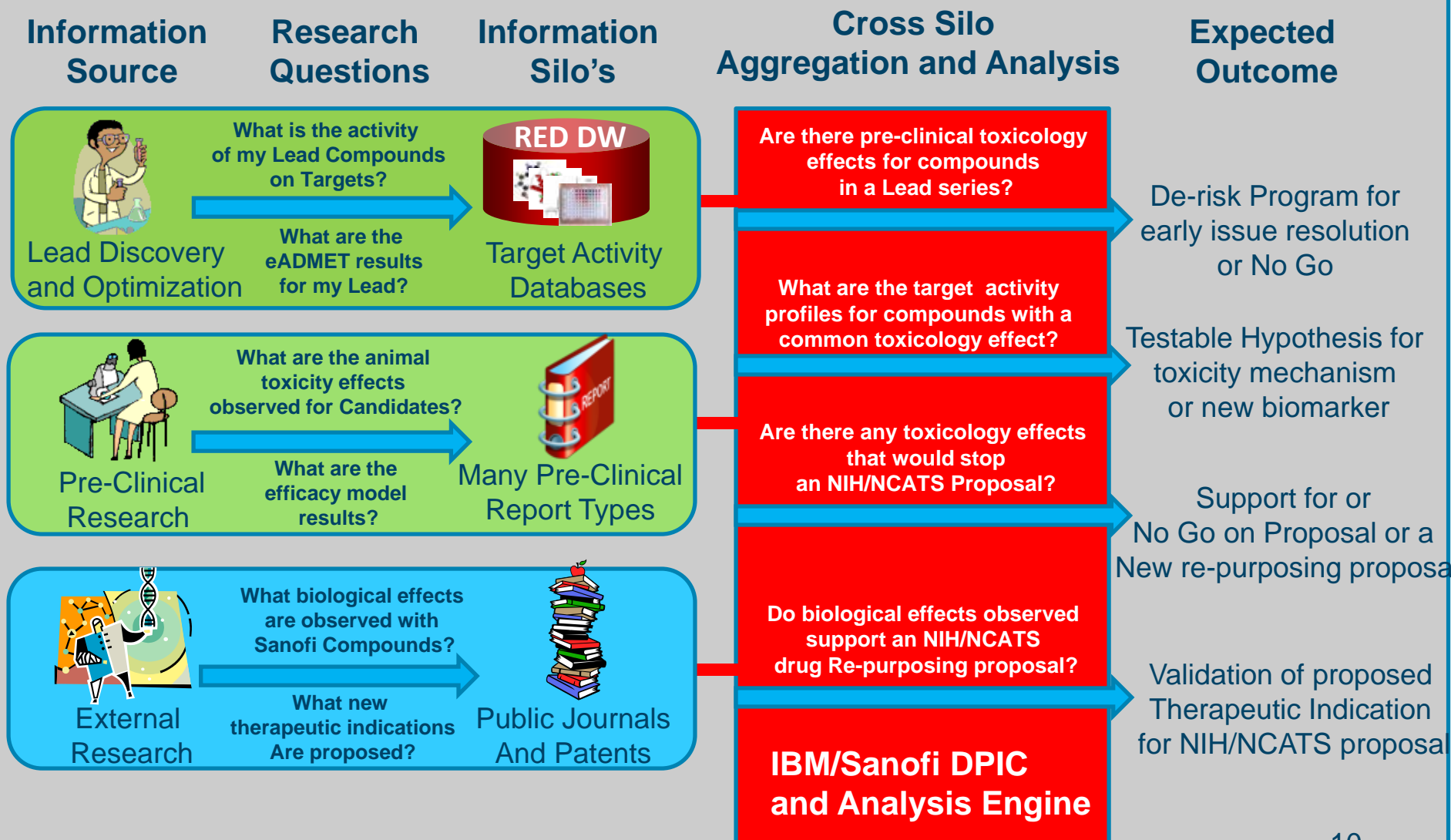
- Strategic advantage through immediate and effective leveraging of pre-competitive information or new information from collaborations
- Increase quality of decision making at Target, Lead and Development Candidate, Choose the best programs.
- Shape the future profiling and disease area strategies to advance superior candidates to the next stage of clinical trials

## Enablement

- Enable every researcher to be as good as the best with access to accumulated knowledge

Potential to **decrease attrition**, **reduce R&D costs**, and **increase success rate** for selected candidates to accelerate the delivery of new therapeutic options to patients

# Pilot Exercise Concept Illustration



# Safety & Toxicology Use Case – Pilot Overview

- Focus on a prototype Toxicology Information Cloud that will extract information from a Rat Toxicology report set used in the IMI eTOX initiative and known toxicity test-cases for benchmarking.

## ***Key Questions:***

1. Were there any pre-clinical rat toxicology effects observed for compounds similar to my lead series?
2. Can we observe a pattern of biological target activity among a series of compounds that show a common toxicological effect?
3. Are there any toxicity issues that impact decisions on NIH/NVCATS initiative?

### **Evaluation Criteria 1:**

Queries return at least 80% of known effects in the eTOX data in our VITIC database benchmark

### **Evaluation Criteria 2:**

Detect a pattern of target-based activity for compounds with a common toxicity effect and leads to an experimentally testable hypothesis

### **Evaluation Criteria 3:**

Successful detection of known toxicity effects on the NIH/NCATS compounds

## Safety & Toxicology Benchmarking

1. Find compounds with specific toxicity outcomes
  2. Identify toxic effects caused by similar compounds
  3. Link adverse effects to pharmacology
  4. Demonstrate the breadth/richness of the data extracted and more exploratory uses
- The above searches could be linked to additional information such as:
    - Which of these compounds also caused ALT increase?
    - What was the NOAEL for these compounds?
    - What were their indications?
    - What published compounds cause liver toxicity and also have structural similarity to Sanofi compounds withdrawn from development?

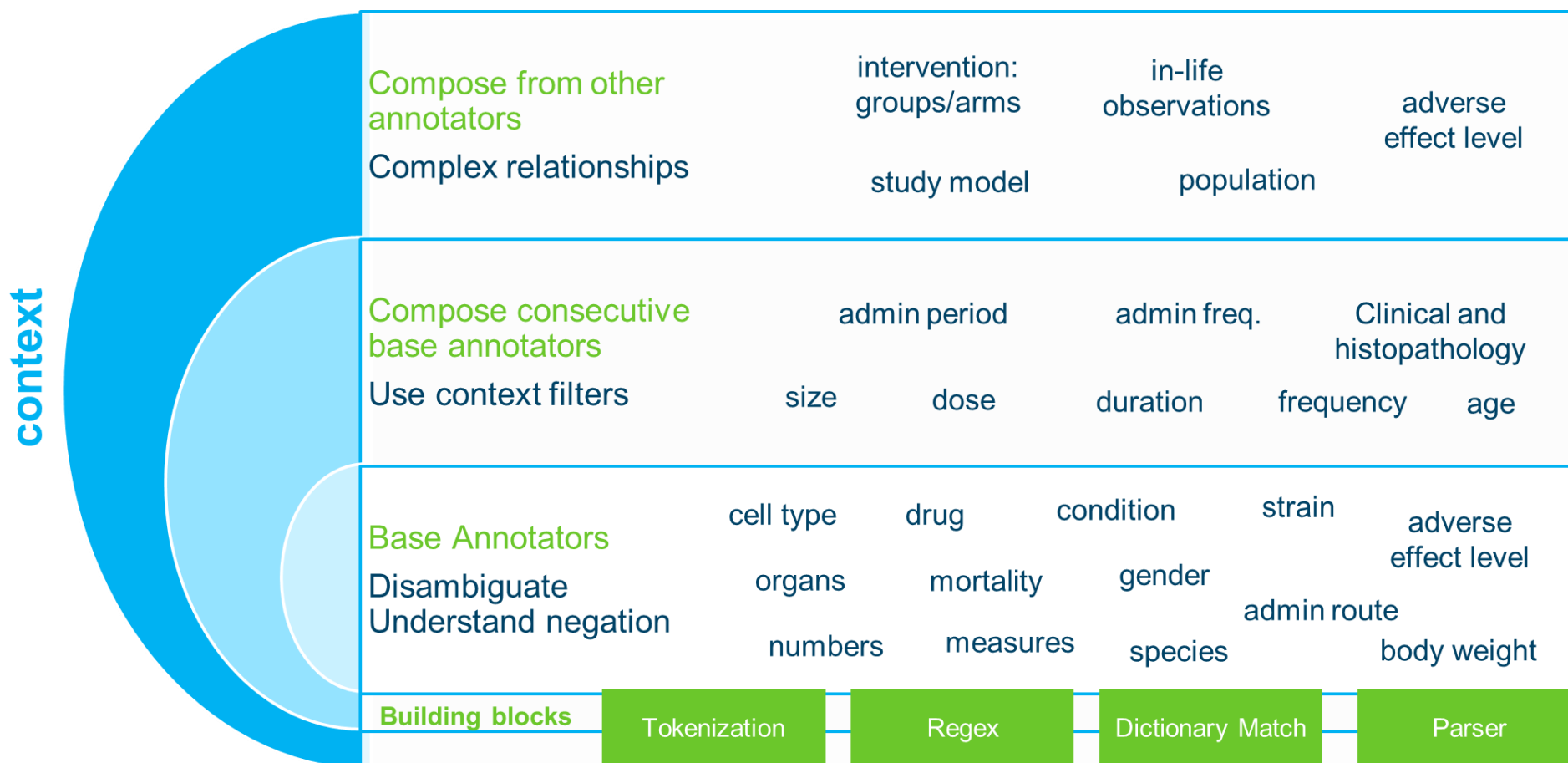
## Execution

The configuration of complex text annotators was required to automatically extract highly contextual information, and accurately recognize experimental design concepts such as test-article, dose levels, and administration frequency and period.


Standardized terminologies were adapted or built for multiple concepts such as organ, cell type, histopathological finding, and severity.

Intelligent recognition and reconstruction of experimental design concepts from free text enabled the association of outcomes to specific study groups

# Configuration of complex annotators required to extract highly contextual information



In the liver a **minimal diffuse hepatocellular hypertrophy** was seen in 9 of 10 males and all 10 females at 0.03 mg/kg/day, minimal to mild diffuse hepatocellular hypertrophy occurred in all 10 males and 10 females at 0.3 mg/kg/day, and mild to moderate diffuse hepatocellular hypertrophy was noted in all 10 males and females at 3 mg/kg/day.

- 
1. Identify structure of the document: e.g. sentence belongs to “Results” section
  2. Identify sub-entities and sub-elements. May require use of specifically-developed dictionaries and/or pattern based discovery
    - Detect cell type via a mix of dictionary and discovery
    - Imprecise severity term “mild to moderate” is mapped to “moderate”
    - Detect modifier term “diffuse” for histological outcome “hypertrophy”
    - Organ type may be recognized or inferred: “hepatocellular” maps to cell type (hepatocyte) AND organ (liver)
  3. Map to canonical terms for compatibility across studies/pathologists/authors/CROs etc.
  4. Relate entities within sentence to form facts (“tuples”)

Organ	Cell Type	Severity	Modifier	Outcome	Section
LIVER	HEPATOCYTE	MINIMAL	DIFFUSE	HYPERTROPHY	Results

## Challenges that have been overcome

- Build relevant dictionaries for entities, including information for mapping synonyms to canonical version (required iterations with SME)
- Configure text extractors to understand:
  - Contextual disambiguation
  - Negation
  - Variations
- Identify relationships between entities to surface complex facts
- Produce generic text extractors while using only a limited training set



## Study Intervention

Text analytics result, Number of rows: 53

Showing page 1 of 1

gro...	gender (S...	drug (TEXT)	adminPeriod (...)	adminFrequen...	adminRoute (TEXT)	dose (TEXT)	spans (STRING)	headers (STRI...	Input Document
0	BOTH	Cpd1	1 month	1 dose per 1 day	Administration, Oral	0 mg/kg	[68582-68583];...	FIGURES	Sanofi_
1	BOTH	Cpd1	1 month	1 dose per 1 day	Administration, Oral	0.03 mg/kg	[6856-16860];...	SUMMARY;M...	Sanofi_
2	BOTH	Cpd1	1 month	1 dose per 1 day	Administration, Oral	0.3 mg/kg	[6862-16865];...	SUMMARY;M...	Sanofi_
3	BOTH	Cpd1	1 month	1 dose per 1 day	Administration, Oral	10 mg/kg	[6870-16872];...	FIGURES;SUM...	Sanofi_
0	BOTH	Cpd2	1 month	1 dose per 1 day	Administration, Oral	0 ppm	[9153-19154];...	REFERENCES;...	Sanofi_
1	BOTH	Cpd2	1 month	1 dose per 1 day	Administration, Oral	10000 ppm	[1195-11201];...	REFERENCES;S...	Sanofi_
2	BOTH	Cpd2	1 month	1 dose per 1 day	Administration, Oral	50000 ppm	[1206-11212];...	REFERENCES;S...	Sanofi_
1	BOTH	Cpd3	1 month	1 dose per 1 day	Administration, Oral	1 mg/kg	[30025-30026];...	FIGURES;MAT...	Sanofi_
2	BOTH	Cpd3	1 month	1 dose per 1 day	Administration, Oral	15 mg/kg	[30086-30088];...	FIGURES;MAT...	Sanofi_
3	BOTH	Cpd3	1 month	1 dose per 1 day	Administration, Oral	35 mg/kg	[6230-16232];...	FIGURES;SUM...	Sanofi_
0	BOTH	Cpd3	1 month	1 dose per 1 day	Administration, Oral	0 mg/kg	[6815-56816];...	FIGURES	Sanofi_
1	BOTH	Cpd3	1 month	1 dose per 1 day	Administration, Oral	0.2 mg/kg	[6323-16326];...	FIGURES;SUM...	Sanofi_
2	BOTH	Cpd3	1 month	1 dose per 1 day	Administration, Oral	2 mg/kg	[6328-16329];...	FIGURES;SUM...	Sanofi_
3	BOTH	Cpd3	1 month	1 dose per 1 day	Administration, Oral	5 mg/kg	[6335-16336];...	FIGURES;SUM...	Sanofi_
0	BOTH	Cpd4	2 weeks	1 dose per 1 day	Administration, Oral	0 mg/kg	[41966-41969];...	EXPERIMENTA...	Sanofi_
1	BOTH	Cpd4	2 weeks	1 dose per 1 day	Administration, Oral	10 mg/kg	[5205-25207]	INTRODUCTI...	Sanofi_
2	BOTH	Cpd4	2 weeks	1 dose per 1 day	Administration, Oral	150 mg/kg	[6538-16541];...	LIST OF ABBR...	Sanofi_
3	BOTH	Cpd4	2 weeks	1 dose per 1 day	Administration, Oral	25 mg/kg	[6563-16565];...	LIST OF ABBR...	Sanofi_
0	BOTH	Cpd4	1 month	1 dose per 1 day	Administration, Oral	0 mg/kg	[29811-29812]	LIST OF ABBR...	Sanofi_
1	BOTH	Cpd4	1 month	1 dose per 1 day	Administration, Oral	10 mg/kg	[29114-29116];...	LIST OF ABBR...	Sanofi_
2	BOTH	Cpd4	1 month	1 dose per 1 day	Administration, Oral	100 mg/kg	[30798-30801];...	LIST OF ABBR...	Sanofi_
3	BOTH	Cpd4	1 month	1 dose per 1 day	Administration, Oral	30 mg/kg	[29121-29123];...	LIST OF ABBR...	Sanofi_

Gender

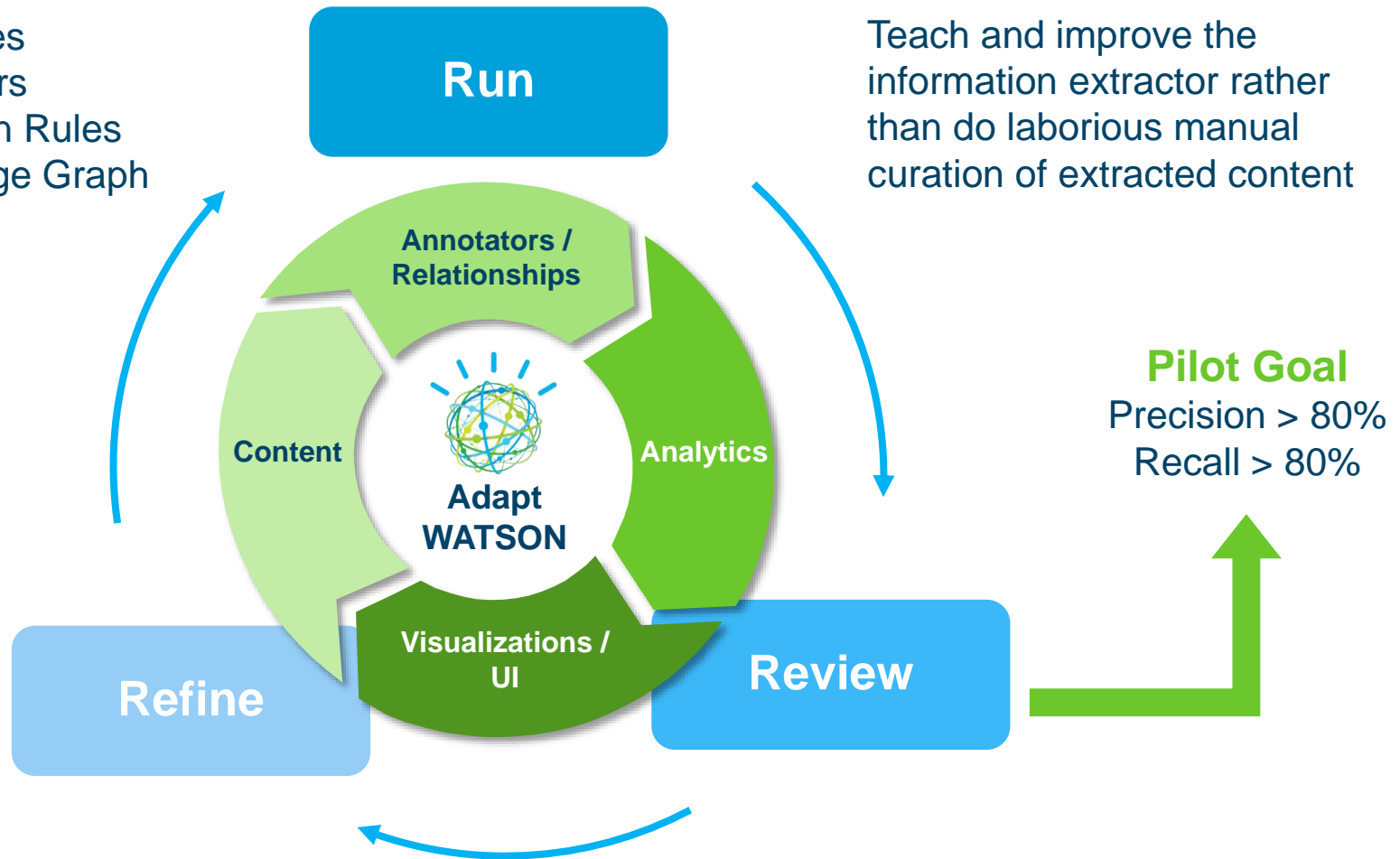
Drug

Admin  
PeriodAdmin  
FrequencyAdmin  
Route

Dose

# Watson Training and Evaluation Process

1. Ontologies
2. Annotators
3. Extraction Rules
4. Knowledge Graph



## Data Access & Visualization

- A flexible, web-based search and retrieval interface, designed by and for toxicologists, provides access to detailed study report data from multiple query points (compound, outcome, target organ, MOA...).
- Real-time filtering of results on multiple criteria, and complex multi-parameter searches are possible.
- Extracted findings are associated to the text from which they were extracted and access to the full study report is provided if required.
- Incorporation of target activity data from an internal SAR database and matching to study via unique compound ID, coupled with innovative visualization tools facilitates the identification of novel relationships between different data domains, such as target-toxicity relationships, co-occurrence of different pathologies.

Find compounds with specific toxicity outcomes

### Example: Liver Necrosis

Text Search:

**Histopathology:**

**Clinical Pathology:**

**Gross Pathology:**

**Compounds:**

**In Life Observation:**

**Drug Name:**

**Mode of Action:**

→ Hepatocellular necrosis was observed in one male (4M37, minimal) and one female (4F36, mild) at 200 mg/kg/day, and minimal to moderate karyorrhexic or apoptotic cells were observed in 4/6 males (2M10, 2M11, 2M12 and 2M13) at 600 mg/kg/day. [52178 - 52418]

# Watson extracts a list of drugs associated with liver necrosis by reading through its corpus of information

Sanofi Internal Data

Text Search: e.g. liver and heart Go

Query Builder

**Histopathology:**

LIVER Select Cell Type

HYPERTROPHY Select Modifier

Select Severity

**Clinical Pathology:**

Select Category

Select Parameter

Select Direction

**Gross Pathology:**

Select Organ or Tissue

Select Observation

Select Direction

Select Modifier

Select Severity

**Compounds:**

Enter list of drug names or SMILES.

[Populate From Working Set](#)

**In Life Observation:**

Select Category

Select Direction

**Drug Name:**

Select Drug Name

**Mode of Action:**

Select Mode of Action

Found 55 Sanofi Internal Document

Search Filter

Empty

Refine Search

**Species**

rat 55

**Tissue and Organs**

LIVER 55

BODY CAVIT... 53

HEAD 53

KIDNEY 52

STOMACH 51

BRAIN 50

SPLEEN 50

ARTERY, AORTA 48

SMALL INTE... 47

THYMUS 47

**Genders**

female male

**Admin Routes**

Administration, Oral 52

Infusions, Intraveno... 1

Injections, Intraper... 1

**Table** **Targets - Drugs Table** **Visualization**

**Drugs**

Cpd1

Cpd2

Cpd3

Cpd4

Cpd5

Cpd6

Cpd7

Cpd8; Cpd8metabolite

**# of Documents**

3

2

2

2

2

2

2

2

**Drug Info**

[View Drug Info](#)

[View Drug Info](#)

[View Drug Info](#)

[View Drug Info](#)

[View Drug Info](#)

[View Drug Info](#)

[View Drug Info](#)

[View Drug Info](#)

**List of Drugs Associated with Liver Necrosis**

Export Selected Compounds Add to Compounds

Identified  
Sanofi Tox  
Reports related  
to Liver  
Hypertrophy

Dashboard  
Categorizes  
documents by  
species, tissue,  
gender, etc.

Watson enables the scientist to drill down into the source documents to investigate further

The toxicologist can view the related documents and drug information for AVE-ABCD...

Table Targets - Drugs Table Visualization

Drugs	# of Documents	Drug Info
<input type="checkbox"/>	2	<a href="#">View Drug Info</a>
<input type="checkbox"/>	1	<a href="#">View Drug Info</a>
<input type="checkbox"/>	1	<a href="#">View Drug Info</a>
<input type="checkbox"/>	1	<a href="#">View Drug Info</a>
<input type="checkbox"/>	1	<a href="#">View Drug Info</a>
<input type="checkbox"/>	1	<a href="#">View Drug Info</a>
<input type="checkbox"/>	1	<a href="#">View Drug Info</a>
<input type="checkbox"/>	1	<a href="#">View Drug Info</a>
<input type="checkbox"/>	1	<a href="#">View Drug Info</a>
<input type="checkbox"/>	1	<a href="#">View Drug Info</a>

Export Selected Compounds Add to Compounds

Found 2 Sanofi Internal Documents for ((sanofi\_histopathology\_instances THRU CONTENT instances))  
x Filtered by ave

sanofi\_histopathology\_annotator\_v1\_0.canonicaloutcome:"NECROSIS") WITHIN (CONTENT

List of Documents Related to "AVE-ABCD"

ID	Title	Drug Name
GGA_Sanofi_2005-0871.txt	TOXICOLOGY STUDY REPORT DSE 2002-0871   : 2-WEEK ORAL EXPLORATORY TOXICITY STUDY IN RATS Author(s):	
GGA_Sanofi_2003-0148.txt	TOXICOLOGY STUDY REPORT DSE 2003-0148   : 4-WEEK ORAL TOXICITY STUDY IN RATS WITH A 2-WEEK RECOVERY	



# Watson understands, extracts and organizes toxicological information from a source document

## Histopathology information extracted and organized from a Sanofi Tox Report

TOXICOLOGY STUDY REPORT DSE 2002-0871 : 2-WEEK ORAL EXPLORATORY TOXICITY STUDY IN RATS Author(s): GGA_Sanofi_2005-0871.txt					
Report	Study Intervention	Study Population	In Life Observation	Histopathology	Gross Pathology
Organ/Tissue	Modifier	Cell Type	Outcome	Outcome Mention	Sentence
HAIR	N/A	N/A	ULCER	ulceration	Minimal skin hair follicle atrophy or moderate epidermal ulceration were noted in 1/5 and 2/5 female rats treated with 600 and 1200 mg/kg/day, respectively.
HAIR	N/A	N/A	ULCER	ulceration	Minimal skin hair follicle atrophy or moderate epidermal ulceration were noted in 1/5 and 2/5 female rats treated with 600 and 1200 mg/kg/day, respectively.
	N/A	HEPATOCTYE	HYPERTROPHY	hypertrophy	Compound—related microscopic findings consisted of centrilobular hepatocellular hypertrophy and centrilobular hepatocellular Vacuolation at all doses with the highest incidence and/or severity at 300 and 600 mg/kg/day.
	N/A	HEPATOCTYE	ENLARGE	enlarged	Minimal to mild centrilobular hypertrophy was characterized by enlarged hepatocytes with ground glass cytoplasm.
	N/A	HEPATOCTYE	ENLARGED	enlarged	Minimal to mild centrilobular hypertrophy was characterized by enlarged hepatocytes with ground glass cytoplasm.
LIVER	CENTRIOBLULAR	HEPATOCTYE	VACUOLES	Vacuoles	Centrilobular Vacuolation was characterized by intra—cytoplasmic round Vacuoles in the centrilobular hepatocytes; this finding was minimal to moderate, more pronounced in males than in females and occurred at all doses with the highest incidence at 300 and 600 mg/kg/day.
LIVER	N/A	HEPATOCTYE	NECROSIS	necrotic	Mononuclear cell infiltrates were characterized by small aggregates of lymphocytes and/or macrophages occasionally centered upon single necrotic hepatocytes.
LIVER	N/A	N/A	CONTENT	Contents	3.4.3 Microscopic observations (Table 14 and 23- see Table of Contents) Liver
LIVER	N/A	N/A	CONTENT(S)	Contents	3.4.3 Microscopic observations (Table 14 and 23- see Table of Contents) Liver
LIVER	N/A	N/A	CONTENTS	Contents	3.4.3 Microscopic observations (Table 14 and 23- see Table of Contents) Liver
LIVER	N/A	N/A	NECROSIS	necrosis	1074) dosed at 1200 mg/kg/day had microscopic focal subcapsular necrosis of the liver.

Histopathology Table  
Specific entry related to  
Liver Necrosis and  
sentence from the source  
document that supports it

## The scientist can view the full tox report, which highlights toxicological information

The toxicologist can use this information to:

- Obtain more context around liver necrosis within the source document
- Decide what to explore next

Full Tox Report  
*Paragraph describing an instance of liver necrosis*

Minimal to mild centrilobular hypertrophy was characterized by **enlarged** hepatocytes with ground glass cytoplasm. This change was present at all doses and was more pronounced in females. This microscopic finding was the histological correlate of the increases in liver weights.

Centrilobular Vacuolation was characterized by intra—cytoplasmic round **Vacuoles** in the centrilobular hepatocytes; this finding was minimal to moderate, more pronounced in males than in females and occurred at all doses with the highest incidence at 300 and 600 mg/kg/day. The content of the Vacuoles was consistent with neutral lipids (Oil Red O stain positive).

Mononuclear cell infiltrates were characterized by small aggregates of lymphocytes and/or macrophages occasionally centered upon single **necrotic** hepatocytes. This change increased in incidence and severity in treated groups as compared with controls.

Aventis Pharma — Confidential

Toxicology Study Report

DSE 2002-0871

One female animal (no. 1074) dosed at 1200 mg/kg/day had microscopic focal subcapsular **necrosis** of the liver. The relationship of this finding to the treatment is considered doubtful based on its focal and isolated nature and the occasional spontaneous occurrence of such change in untreated animals.

Stomach

Treatment—related microscopic changes consisted of submucosal edema and epithelial erosion/ulcer of the forestomach that were noted in all groups with an increase in incidence in animals treated with \_\_\_\_\_ when compared to controls.

Non— glandular stomach microscopic changes are presented in the table below.



# Demonstrate the breadth/richness of the data extracted and more exploratory uses

Example: Liver necrosis in rats + ALT increase

IBM Accelerated Discovery Solutions: Toxicity

Sanofi Internal Data

Text Search:  Go

Query Builder

Histopathology: LIVER, Select Modifier, Select Cell Type, NECROSIS

Clinical Pathology: Select Category, Alanine transaminase, INCREASE

Gross Pathology: Select Organ or Tissue, Select Organ Weight, Select Direction, Select Severity

Compounds:

In Life Observation: Select Category, Select Direction

Drug Name: Select Drug Name

Mode of Action: Select Mode of Action

Populate From Working Set

Species: rat

Tissue and Organs: ABDOMEN, ADRENAL GLAND, ARTERY, AORTA, BODY CAVITY, BRAIN, CERVIX, DUCT, BILE, EPICARDIUM, EPIDIDYMIS

Genders: male, female

Admin Routes: Administration, Oral

Table Targets - Drugs Table Visualization

atorvastatin:ave5530

STUDY REPORT DDO0984 ATORVASTATIN: Exploratory 14-day oral toxicity study in rats Authors: Dr. Kerstin Wase, Dr. GGA\_Sanofi\_Atorvastin\_DDO0984.txt

Report Study Intervention Study Population In Life Observation Histopathology Clinical Pathology Gross Pathology

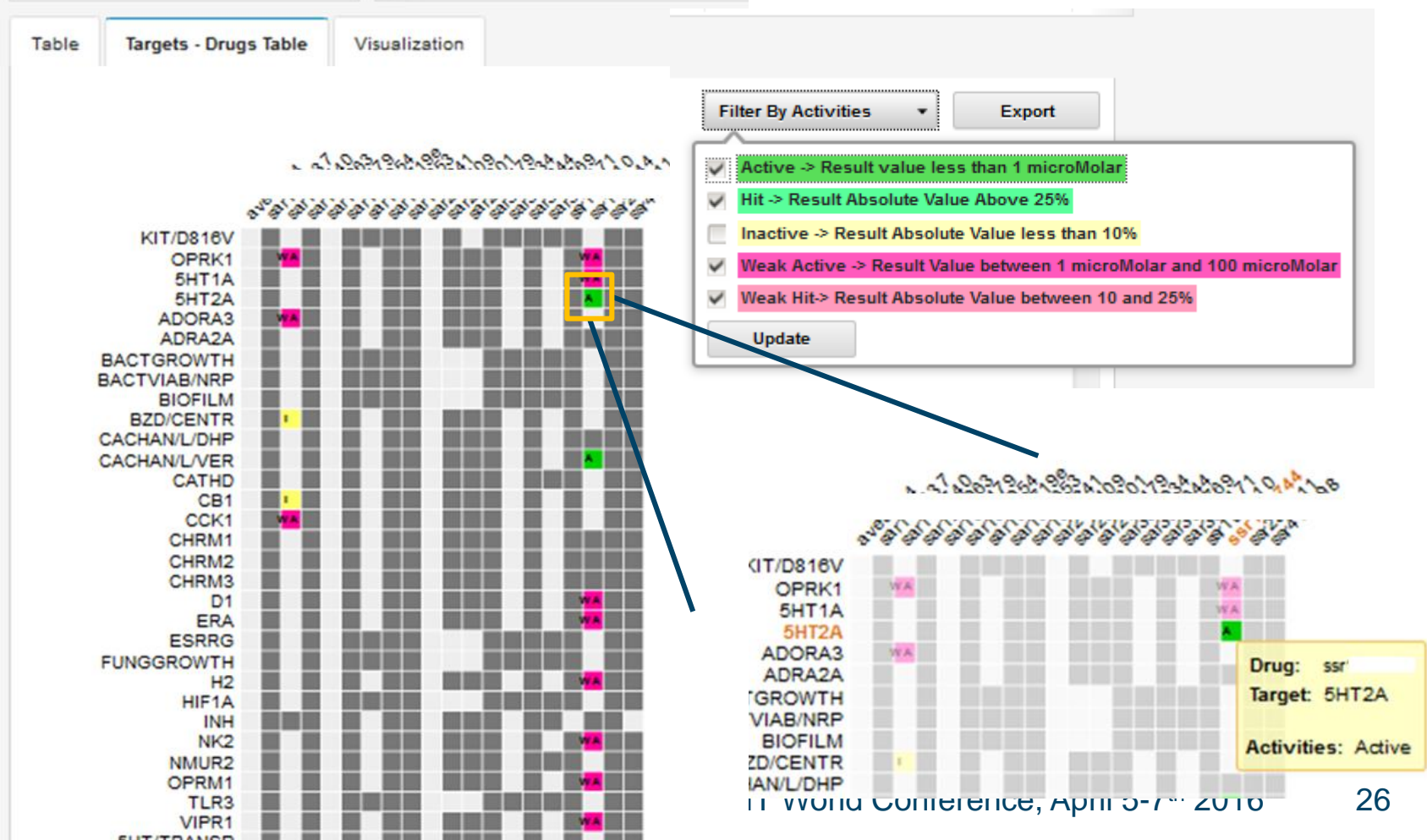
Parameter Category	Parameter	Parameter Mention	Changed Direction	Severity	Sentence
Clinical Chemistry	Alanine transaminase	ALT	INCREASE	N/A	ALT was increased 3.0- and 3.2-fold in males and 3.1- and 3.3-fold in females, respectively.
Clinical Chemistry	Alkaline phosphatase	ALP	INCREASE	N/A	ALP was increased 1.5- and 1.8 fold in males and 1.4- and 1.8-fold in females.

Found 1 compound (atorvastatin) that includes both attributes

List of Drugs Associated with Liver Necrosis + ALT Increase

Clinical Pathology Table Specific entry related to ALT increase and sentence from the source document that supports it

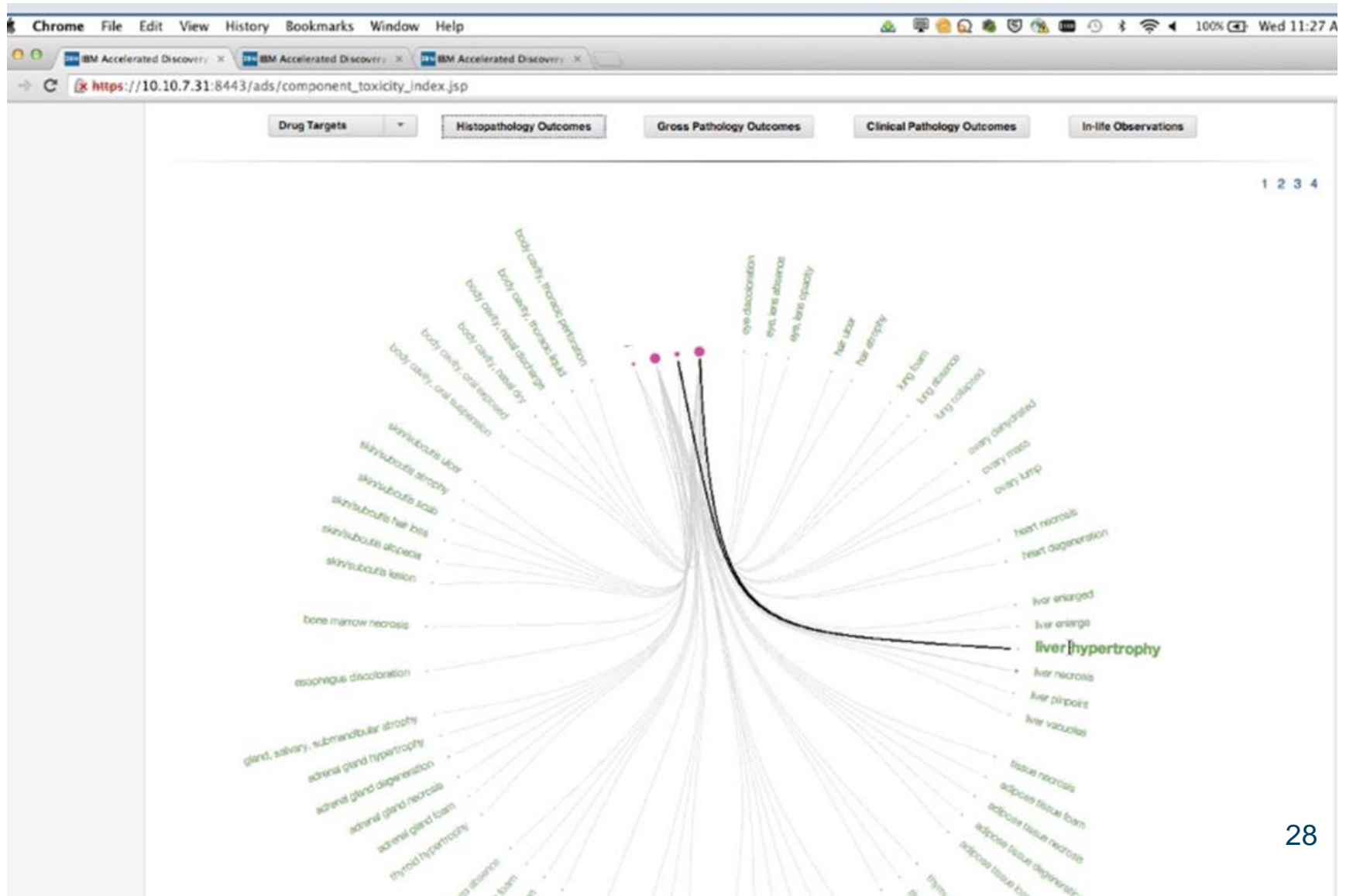
SSR-ABCDEF shows some activity towards target 5HT2A



AVEABC is linked to liver necrosis and about 30 other toxicity effects such as liver hypertrophy, skin lesions, etc. (**bold cords**)



# Identify toxic effects caused by similar compounds





### 3. Link adverse effects to pharmacology

#### Example: Liver Weight Increase

“SSR-ABCDEF” is associated with liver weight increase.

Visually, the toxicologist can see the drug-target relationships

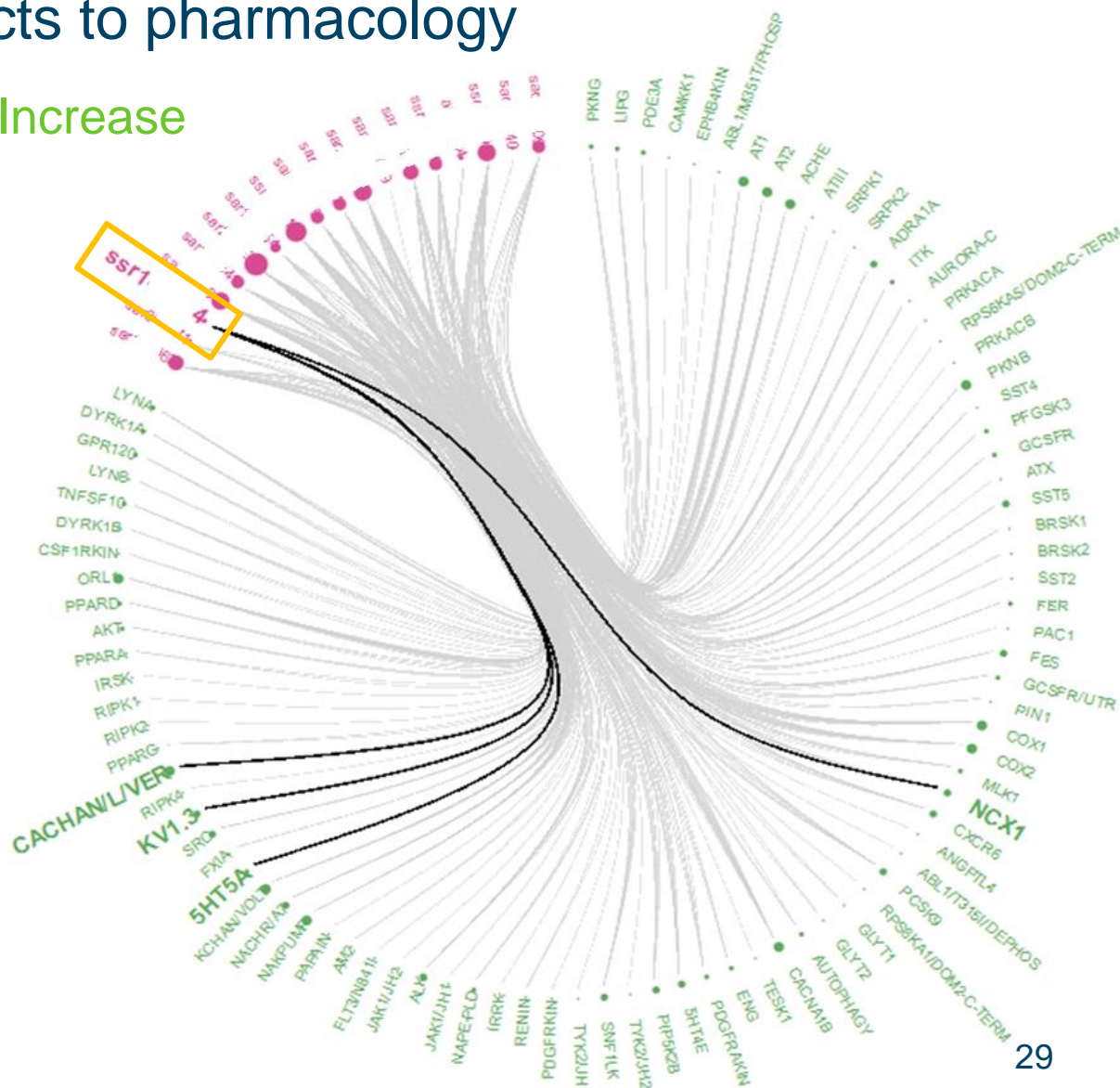
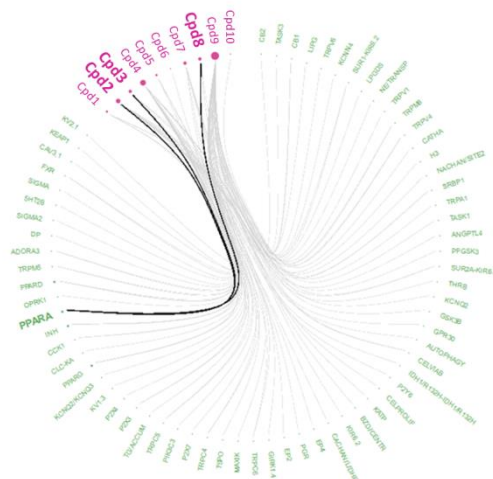


Table Targets - Drugs Table Visualization

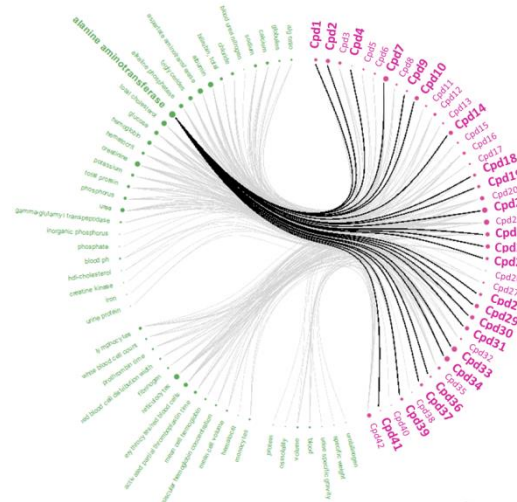
Found 46 targets for drugs in result set

	Cpd1	Cpd2	Cpd3	Cpd4	Cpd5	Cpd6	Cpd7	Cpd8	Cpd9	Cpd10
INH										
PPARA										
PPARD										
PPARG										
ADORA3										
CCK1										
CELVIAB										
SRBP1										
SUR2A-KIR6.2										
TRPC4										
TRPC5										
TRPC6										
TRPM5										
TRPV4										
5HT2B										
ANGPTL4										
AUTOPHAGY										
CACHAN/LDHP										
CATHA										
CB2										
CELPROLIF										
DP										
EP2										
EP4										
GPR30										
GSK3B										
H3										
DH1/R132H-IDH1/R132H										
KEAP1										
LIPG										
LPGDS										
NACHAN/SITE2										
NE/TRANSP										
OPRK1										
P2X3										
P2Y6										

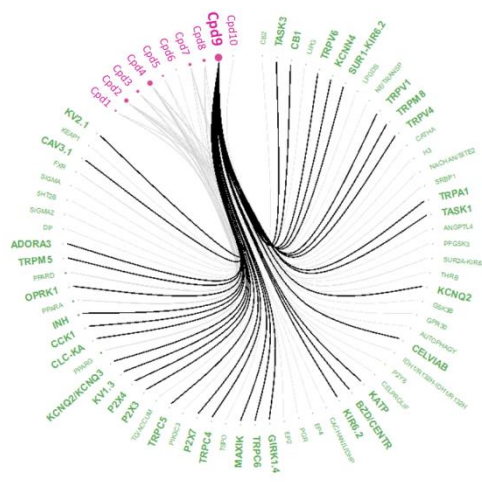
## Promiscuous Target



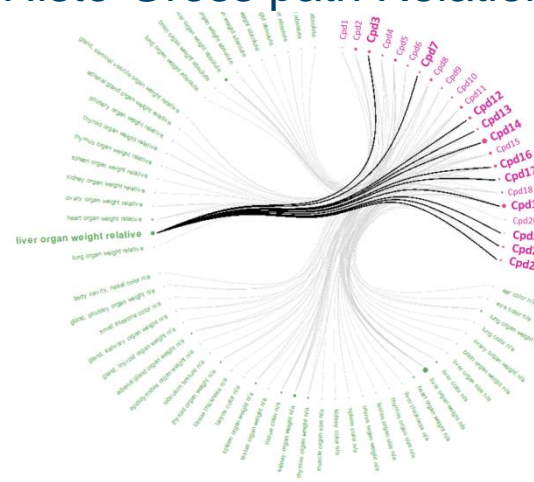
## Promiscuous Compound



## Histo-Clin path Relationships



## Histo-Gross path Relationships



## Safety & Toxicology Results Executive Summary – *Unlocking scientific information from unstructured data*

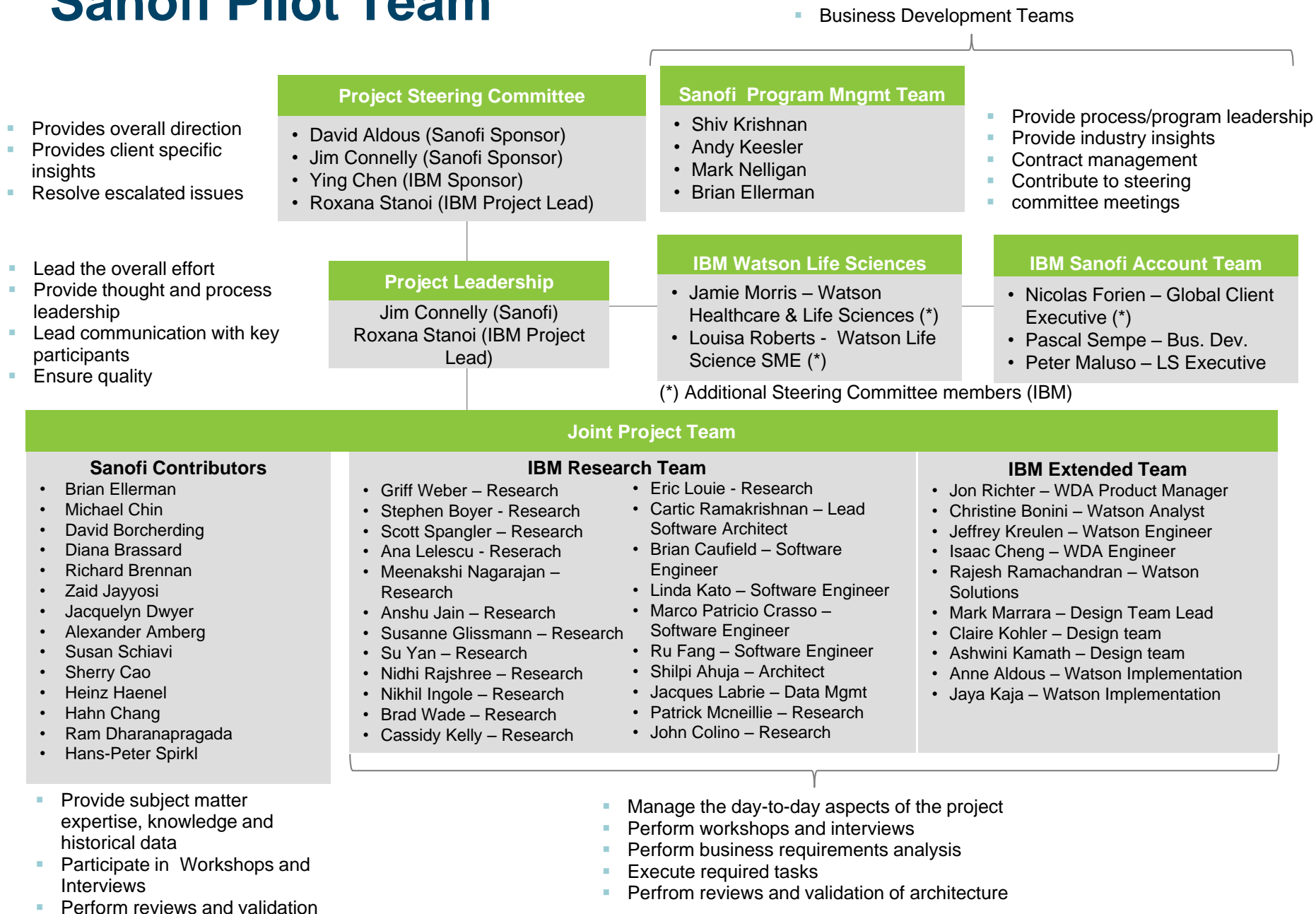
### Watson was able to exceed all pilot objectives and metrics

- ✓ Build Toxicology relevant ontology, annotators and Knowledge Graph
- ✓ Extract High Quality and detailed information from Toxicology reports
- ✓ Find compounds with specific toxicity outcomes
- ✓ Identify toxic effects caused by similar compounds
- ✓ Link toxicology outcomes with pharmacology
- ✓ Extract rich and deep information that can be further explored
- ✓ Precision and recall of information extraction ranged from 80-95%



**WATSON** has demonstrated considerable capabilities to help identify safety and toxicity signals, and we plan to expand the Watson corpus with additional internal toxicity and safety reports to further drive value

# Sanofi Pilot Team





# Single Vendor SAR Data Platforms “in the Cloud”

**James Connelly**

Head of Global Research Data Management,  
**Sanofi**

# Outline

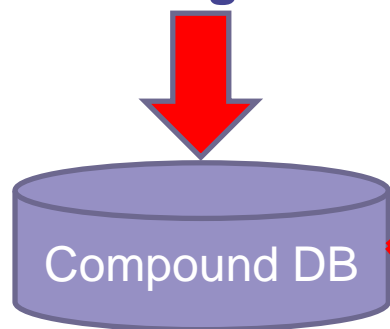
---

- Traditional approaches to harmonize SAR Data
  - Compounds, Bioassays and ISIS DB's
  - Integrated SAR Data Platform and DataWarehouse-centricity
  - Evolution of RED DP and challenges/gaps
  - High Cost of big RED DP integration and support
- Back to the future with project-centric DB's
  - Drug Discovery Collaboration “in the cloud”
  - Local project-centric DM's
  - Evolution of full single-vendor SAR data platforms
- Leverage new “Big Data” Technologies to build a distributed SAR data environment

# Original “Project Databases”

## Project Databases In ISIS

Compound/Substance  
Batch/Logistics



Data Transfer  
By Excel & eMail

Proj 1

Proj 2

Proj 3

Proj 4

Proj 5

Proj 6

Proj 7

Proj 8

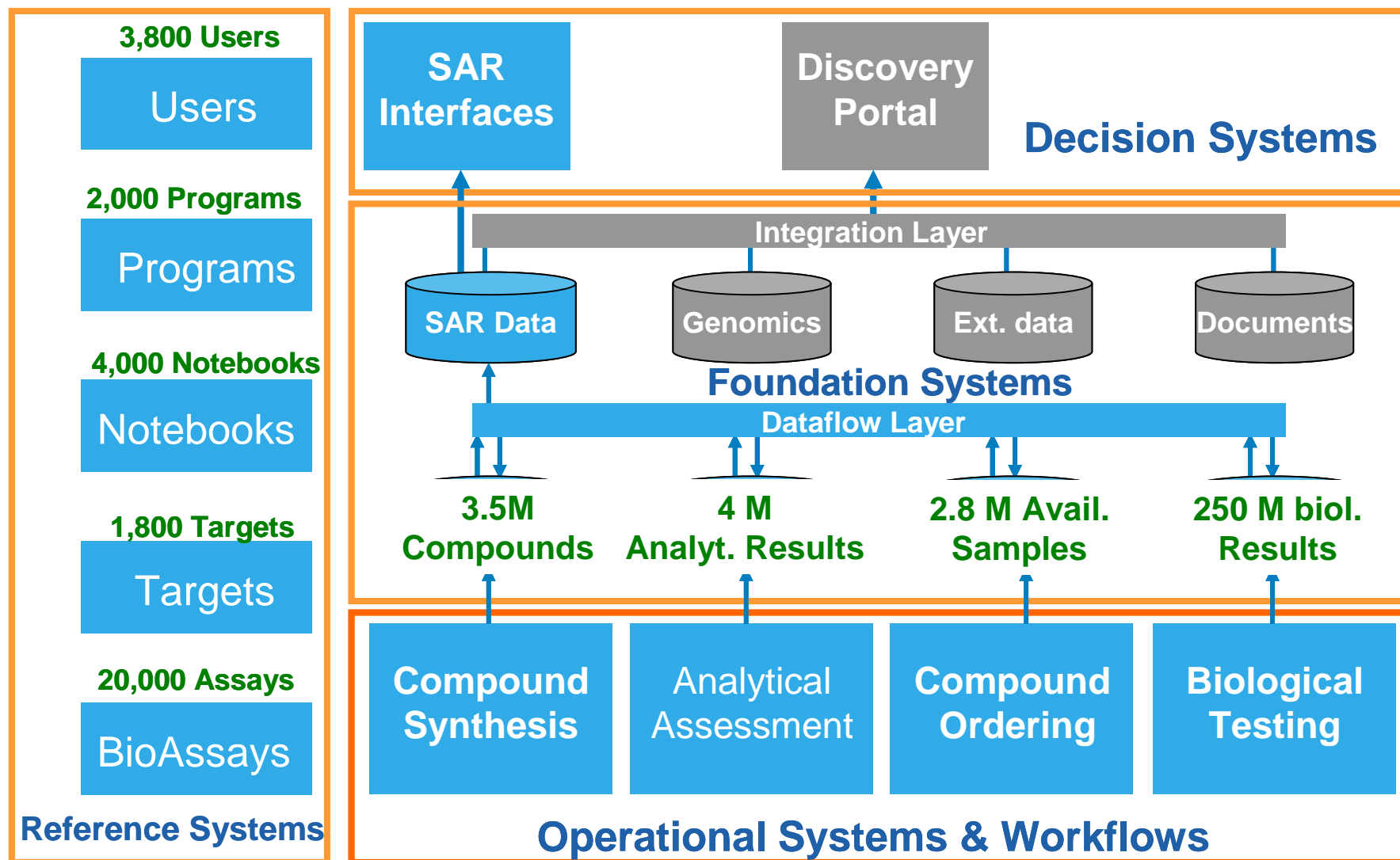
...

Analytical Results  
&  
BioAssay Results

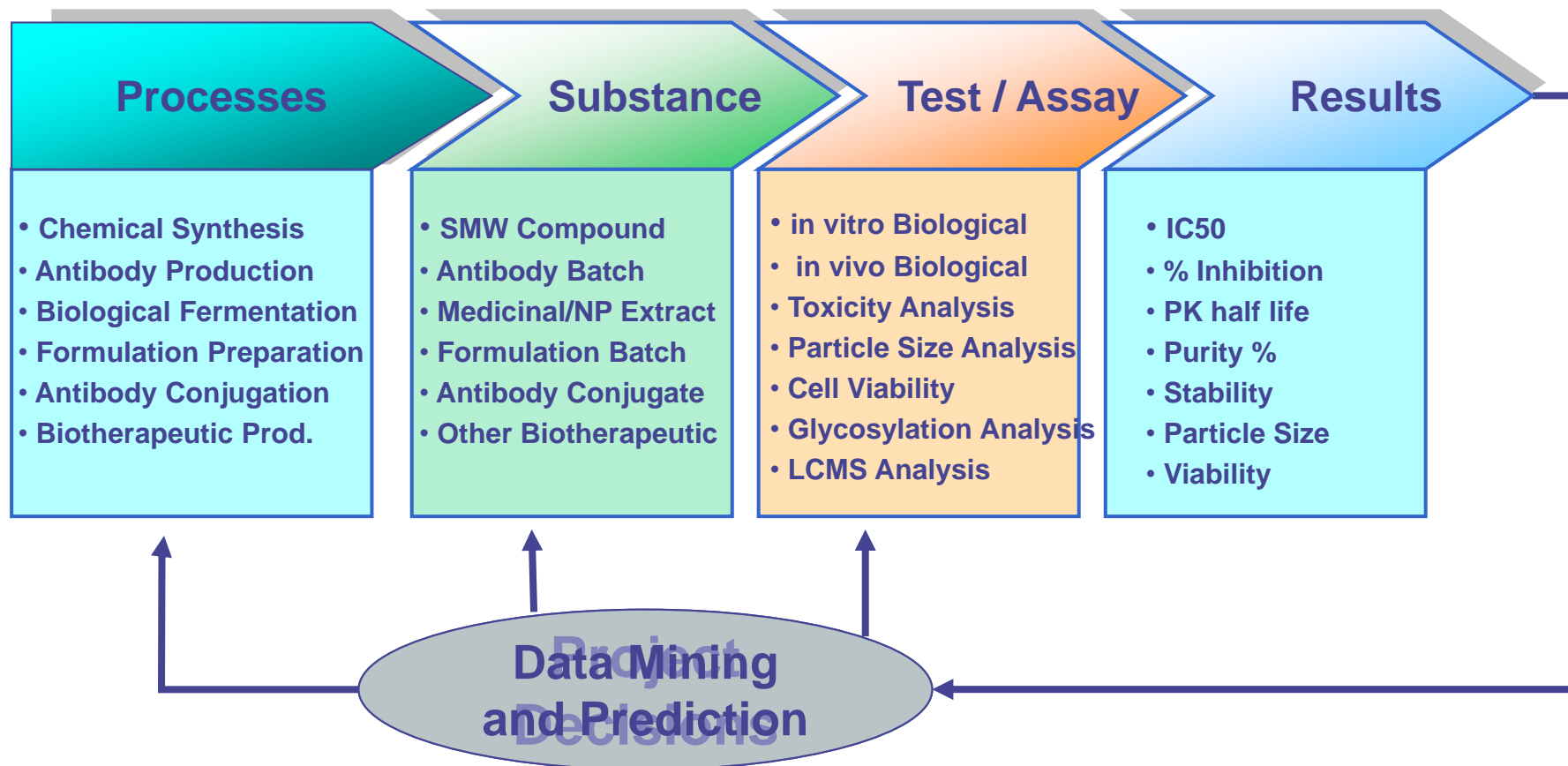
Data Transfer  
By SD File & eMail

- Good for Projects
  - Form View preferred by MedChemists
- However,
- Created “Data Islands”
  - Data-Mining Impossible

# Discovery SAR Data System

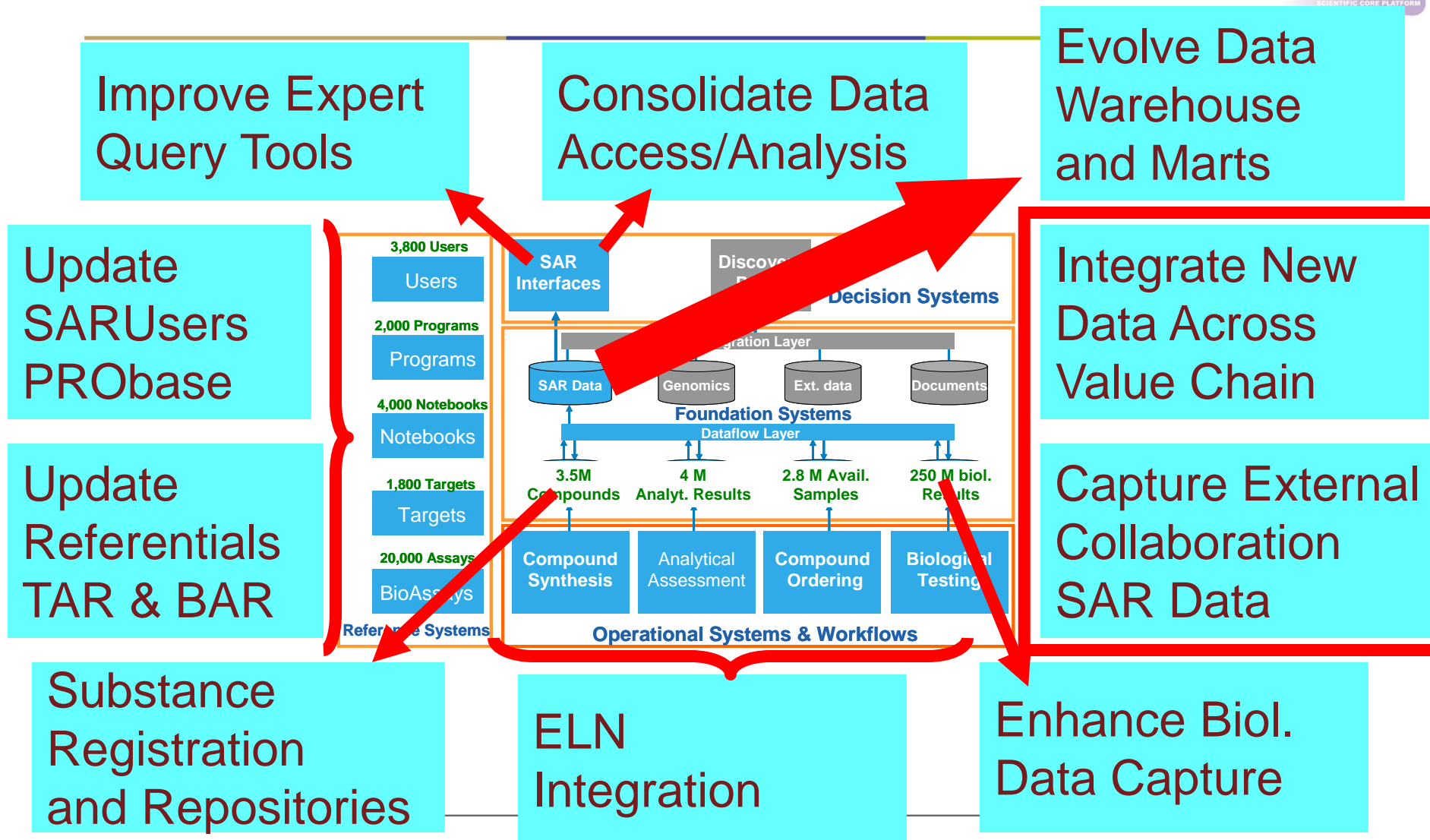


# From **Structure** to **Substance** Activity Relationship



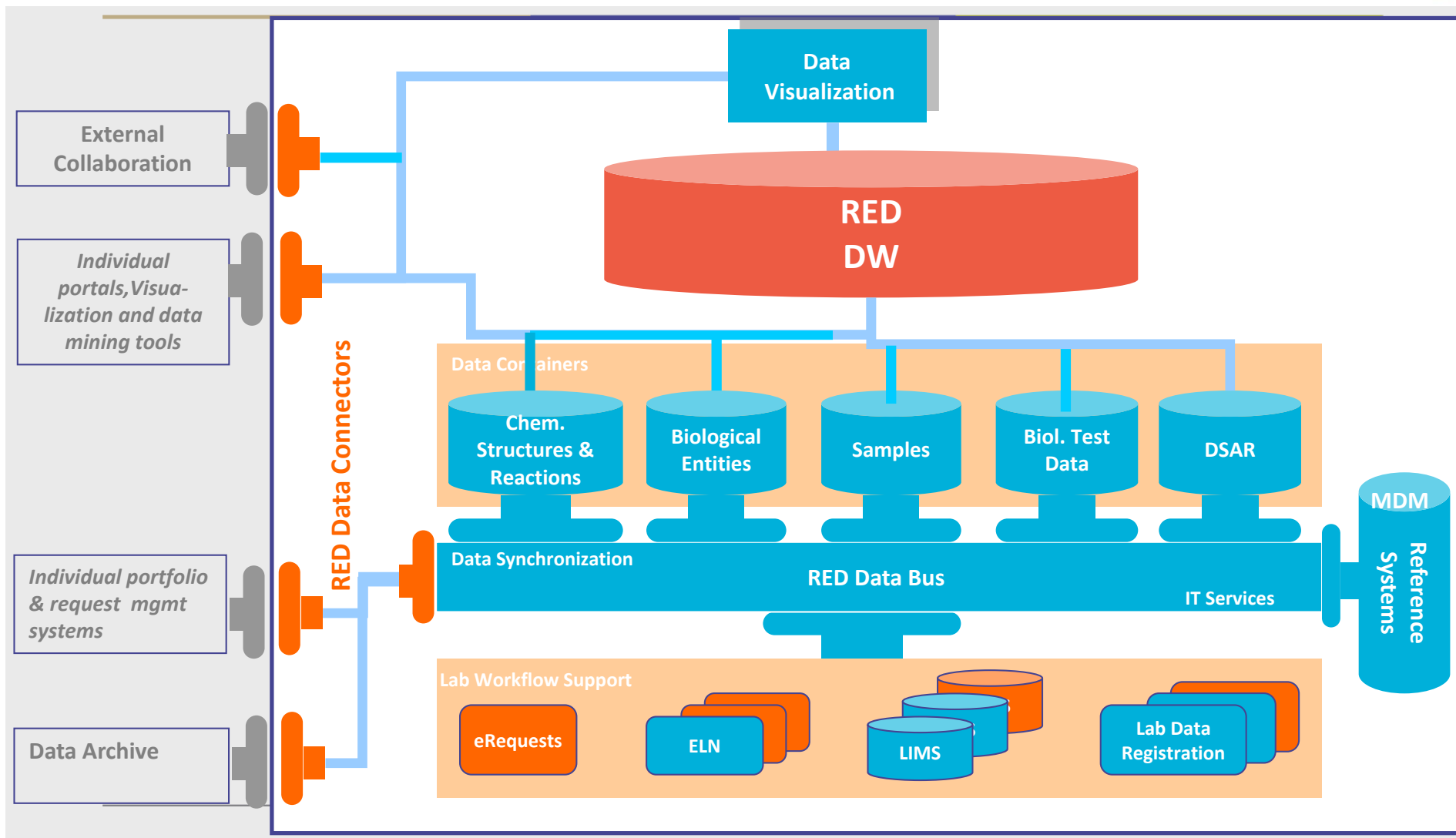
**Need to view results in context of processes, substance metadata and test/method**

# Evolution of Discovery SAR Data Systems





# Data Warehouse-Centric Platform



# Large, Integrated DataWarehouse Platforms

- Chose “Best in Class” applications and systems
- Enabled data-mining, modeling and drug design
- Provided harmonized, high quality data

## However,

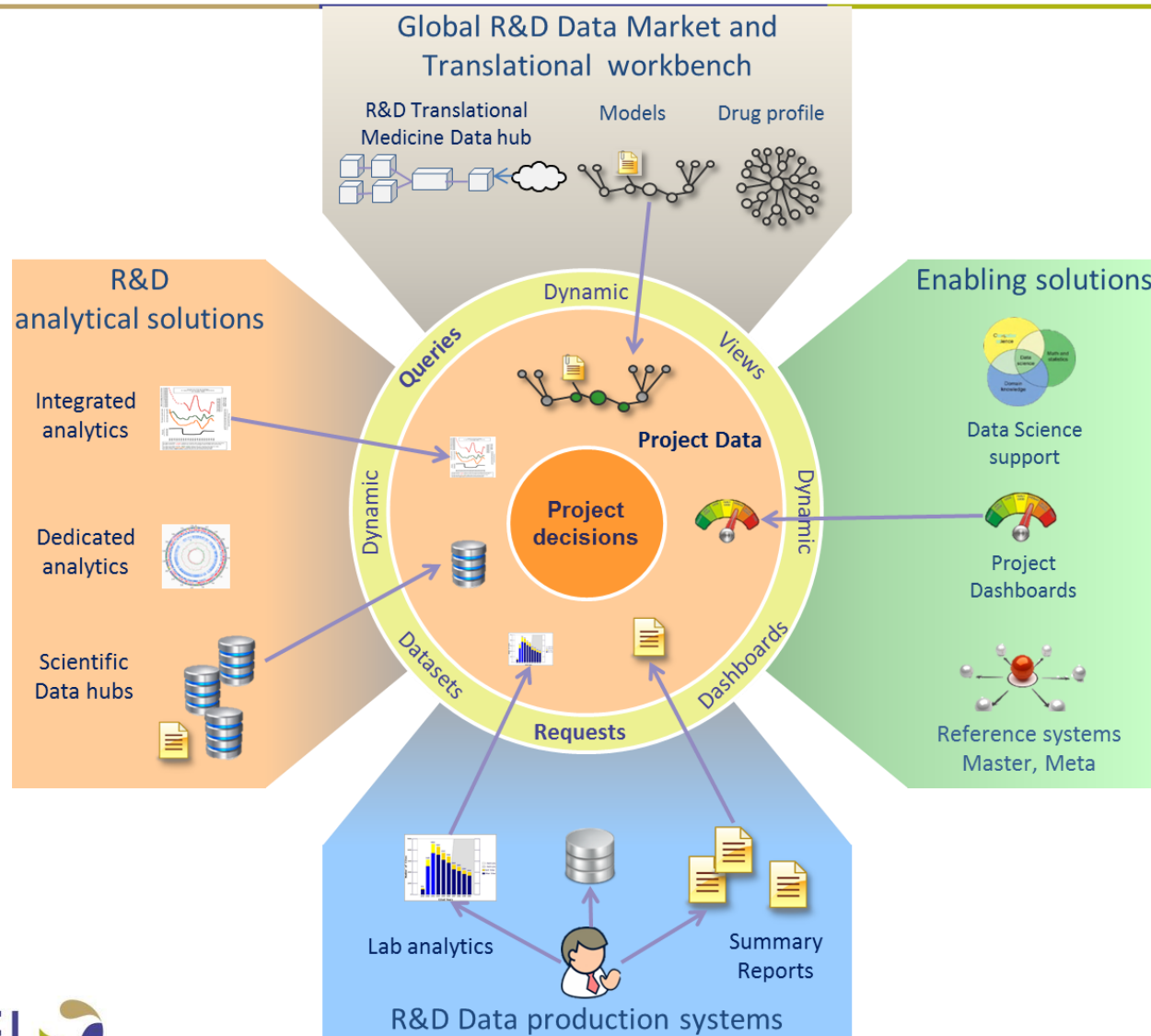
- Difficult to introduce new data types and complex assay results
- Performance and capability gap issues that delayed project progression
  - Data availability lag, query performance, “one size fits all”
- Complex data integration required high level of support, high cost
- Platform was rigid, relatively static and did not evolve with availability of new technologies and capabilities.
- Vendors raised prices because we were locked in

## Related to collaborations

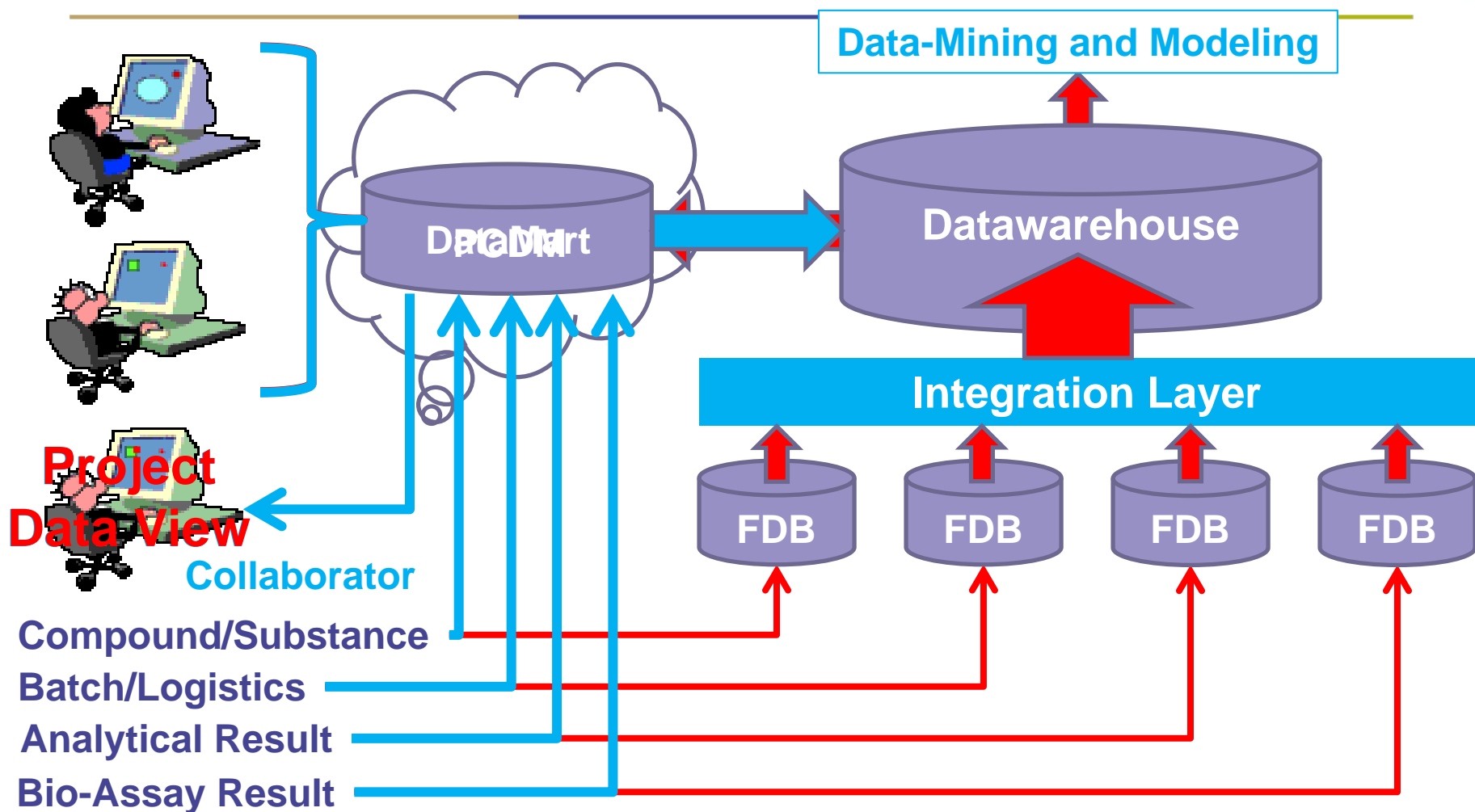
- IP tracking related issues and data access control issues
- Difficult to remove data from DW at end of collaboration
- Registration without structures and subsequent data visualization issues

# The Vision

## A Project-Centric Platform



# Reverse the Flow of Data



“Set up a **service** to support drug discovery **data exchange** with **external partners**”

## Previous Data Exchange Practices and Issues

- Excel Sheet data exchange by email or “SharePoint Dumps”
- Confusion on registration in corporate databases
- Difficulty to integrate and aggregate with internal data
- Limited data visualization and analysis options
- Intellectual property and contractual considerations
- Diversity of data standards and conventions
- Diversity of partner types
  - (program purchase, non-profit venture, CRO and true collaboration)
- Large workload on scientists

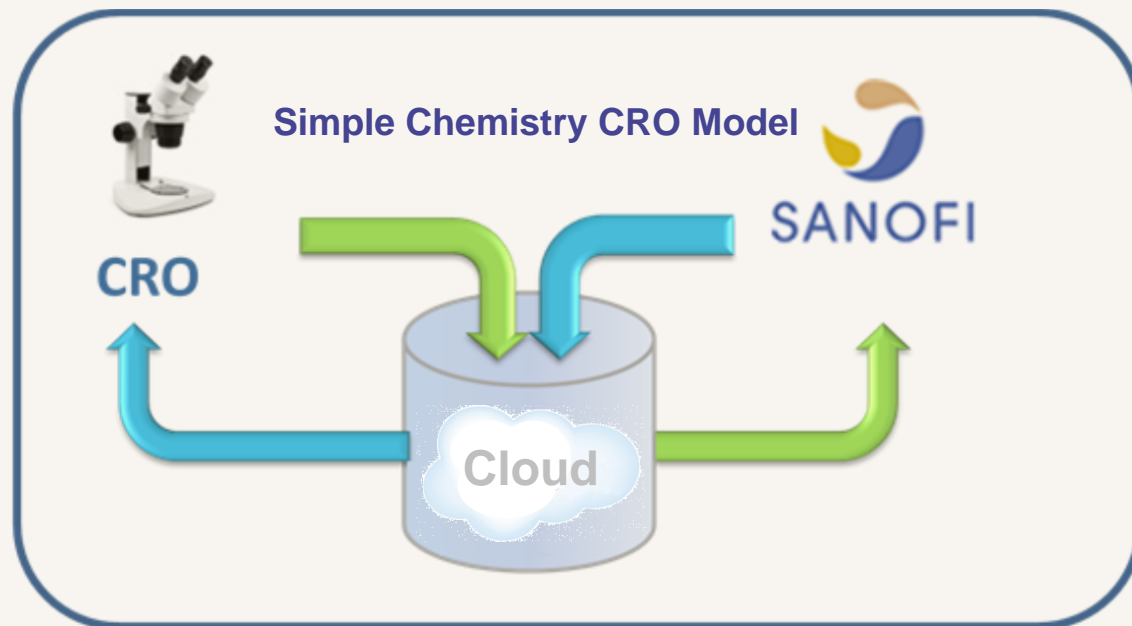
# Basic CRO Model



- To share, exchange and register drug-discovery data with external partners (CRO, academics...)
- To collaborate in an external cloud



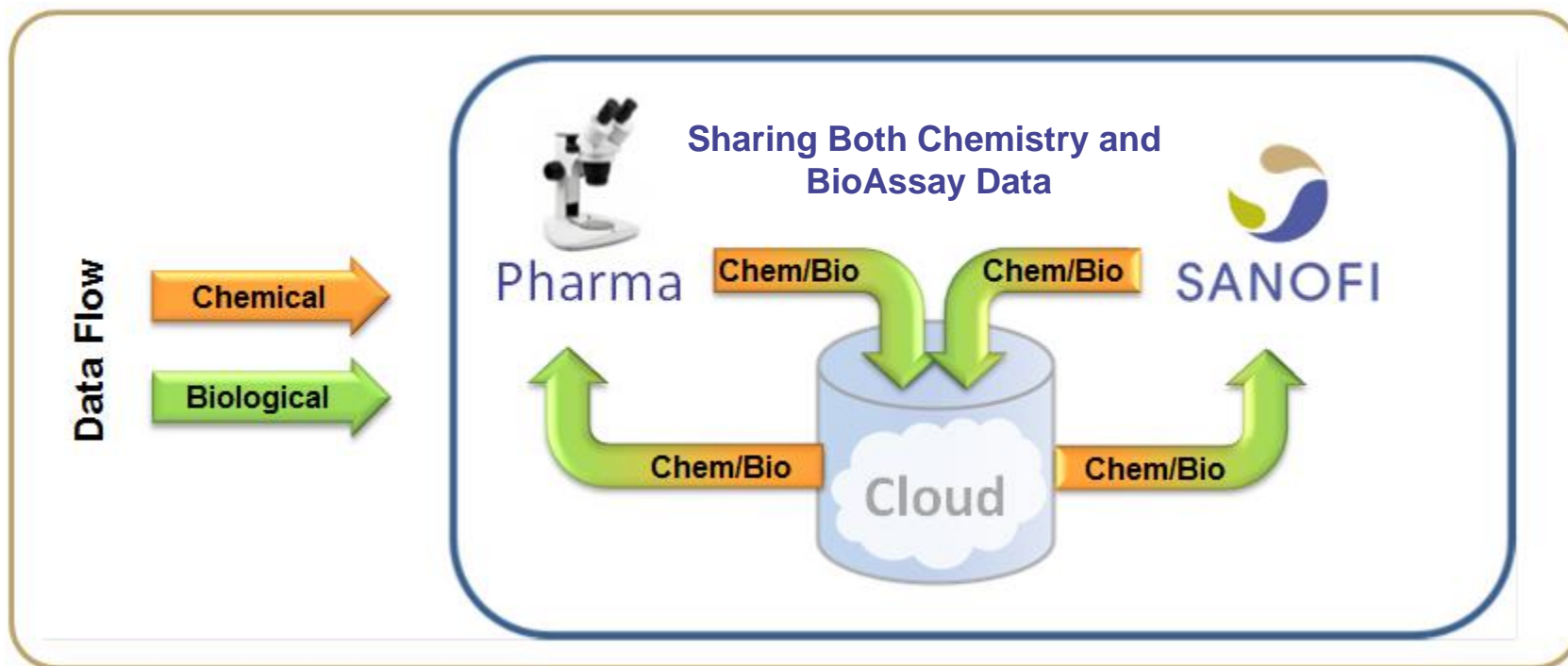
Data Flow





# Partnership Model

- In the context of a partnership with a big Pharma or Biotech Partner, the requirements are more complex...



## **Flexibility** to allow configuration for **many Collaboration Models**

- Outsourcing to **Chemistry CROs** (smw, peptides and soon biologics)
- Outsourcing to **Bio-profiling assay providers**
  - Integration with internal request system
  - automated results upload by providers
- Collaboration with other **Pharmas, Biotechs and Academics**, with or without **Data Ownership**
  - Harmonization and conversion of chemical representations and biological data terminologies and models
  - Automatic import from DW with mapping of partner model
  - Applications provided to partner and used by all scientists
  - Fast and clean project closure and IP tracking
- **Consortiums, PPPs**
  - Adjust automatic transfer with on-premise systems to comply consortiums business rules
- **Evolution to Fully Capable SAR Data Platform “in the Cloud”**



Corporate DB

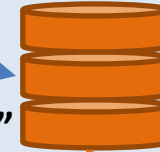
Only Collaboration projects  
data replicated

Internal project  
New data



data  
files

**"Collaboration Specific"**  
Data Mart



Complete SAR Data Platform

- Data Registration
- Data Integration
- Visualization and Analysis

Project Team  
Works in  
Cloud via  
Web Access

Sharepoint (file  
repository)

Option 1 or 2

Option 3



data  
files

Partner database?



**EXTERNAL  
PARTNER**



External Partner data registration options:  
1. direct data registration to cloud platform  
2. csv file upload to cloud platform  
3. data replication from collaborator database

# Single Vendor SAR Data Platforms “in the Cloud”; a Sampling



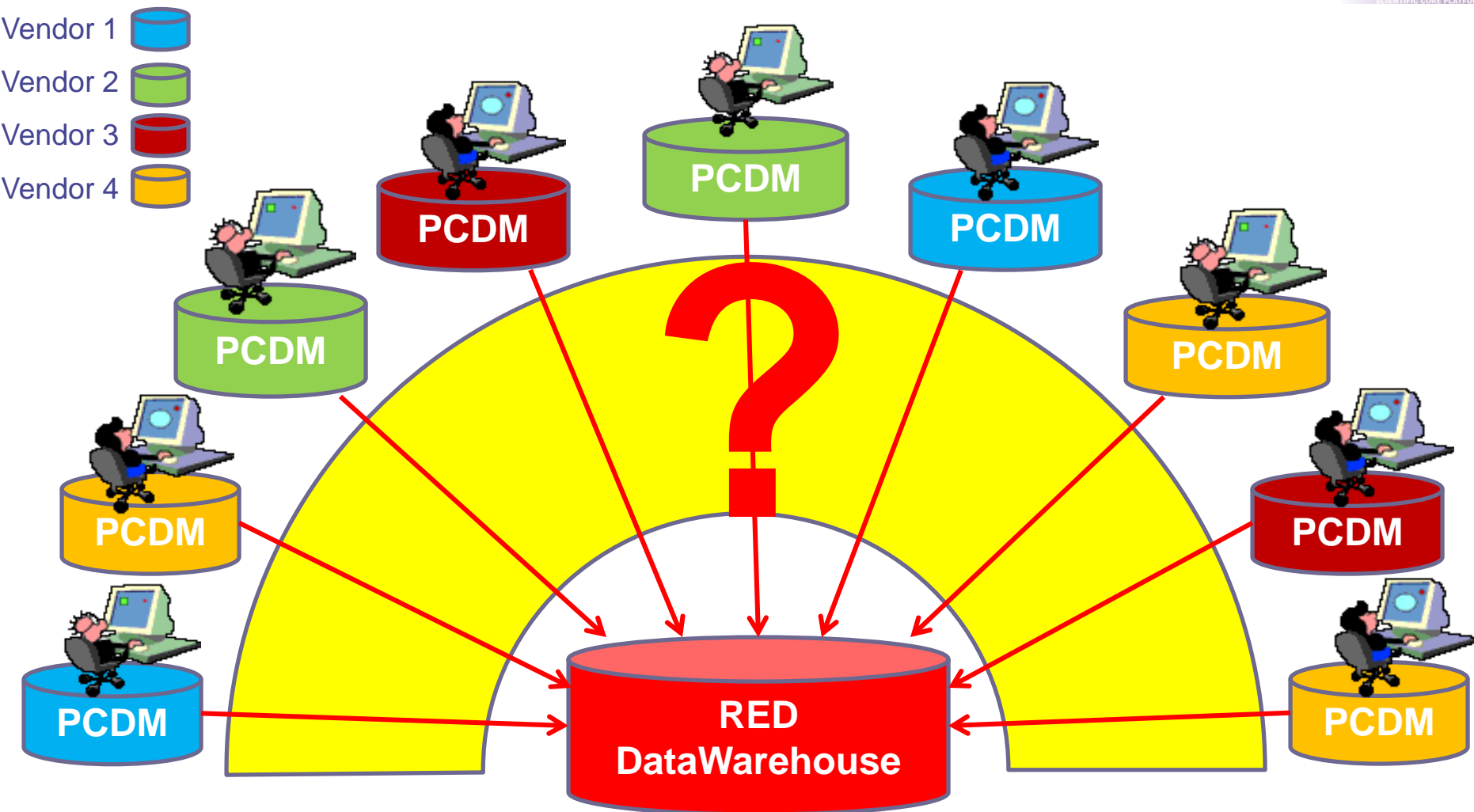
- Vendor 1 and 2; Platform designed for collaborative SAR data exchange, evolving into powerful full SAR Data Platform in the Cloud
- Vendor 3; Hosted Complete SAR Data Platform designed for rapid, inexpensive deployment to Biotech companies. (Also for collaboration)
- Vendor 4; Complete, fully integrated SAR Data Platform moved into the Cloud
- Vendor 5; Easily configurable SAR data platform that can be adapted to support collaborations
- Vendor 6; Powerful analysis and modeling platform available “in the Cloud” that can evolve into a full SAR data platform (with big data)
- Vendor 7; Data capture platform offering web access “in the Cloud” and integration with existing on-premise systems

## Commonalities

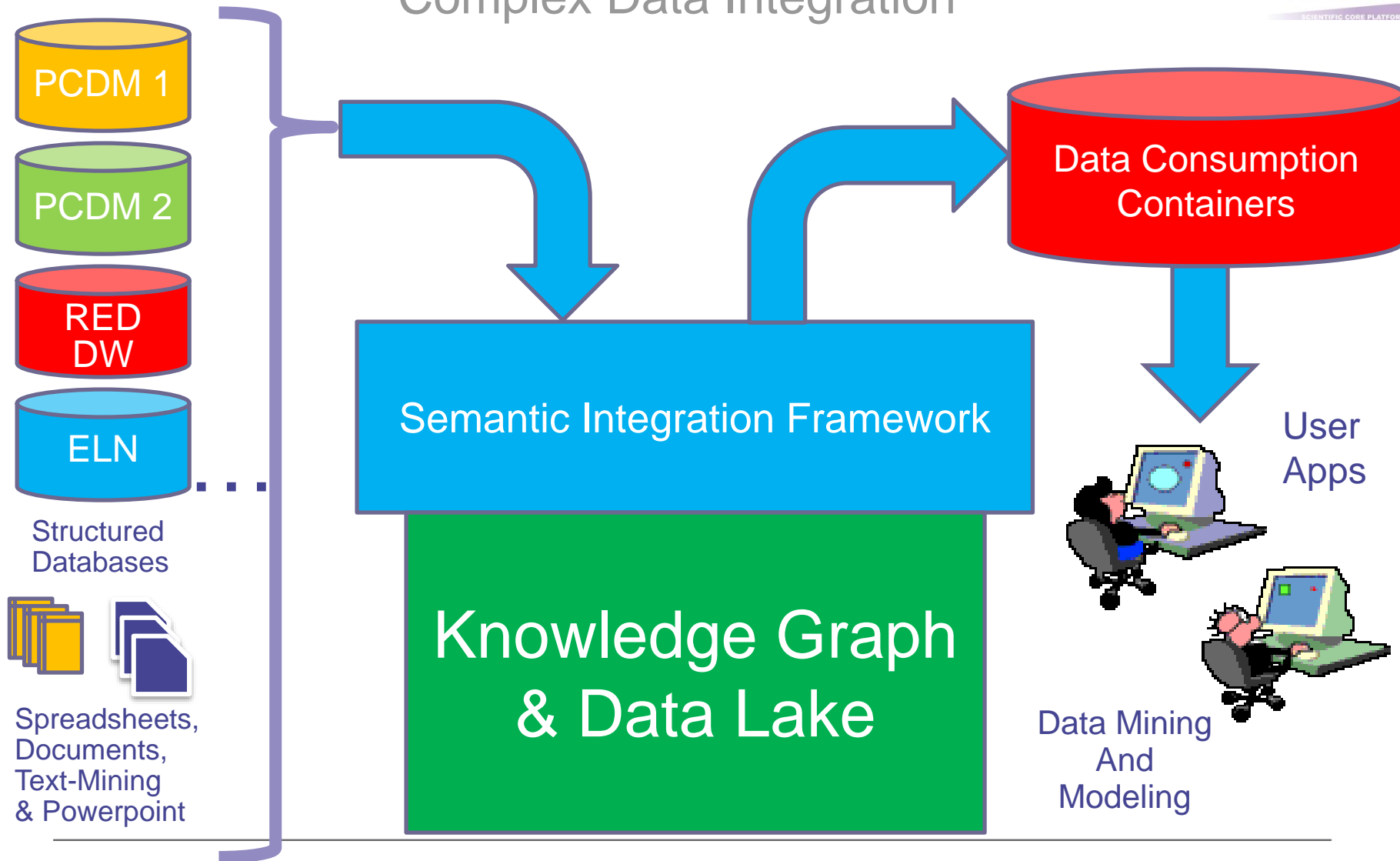
- Web services access
- Easy configuration and access control
- Off-prem or using Cloud resources

Vendor 8; Creating data integration platform that is application agnostic

# How do we Integrate Data Islands?



# “Big Data” Technology used for Complex Data Integration



# Advantages

- Scalable to Fit Demand
- Changes Funding Model
- Flexible Configuration
- Rapid Evolution of Capabilities and Technology
- Efficient for Project Teamwork and Progression
- Rapid Deployment
- Configurable to Needs of Projects
- Vendors Compete to Deliver Best Platform/Cost
- If these systems, designed for external collaboration, work well, why would we not adopt them for internal projects?
- Vendors should Drive Data Standardization



---

# THANKS FOR YOUR ATTENTION!