# Too Much Data or Too Little Cooperation?

**Tom Plasterer, PhD.**
**Research & Development Information (RDI)**
**Director, US Cross-Science**

AstraZeneca

BigDiP USA 2015
Analytics & Insights for Pharma

# Approaching a Pharma Big Data Problem:
## Requirements of the CI Informatics Landscape

**Must span the entire drug development lifecycle**
- and back (post-market surveillance to discovery)

**Must support large and very heterogeneous da**
- single nucleotide polymorphisms to countries

**Will change with new science & new regulatio**
- Medline just under 1M articles/year

**Must work with multiple, international regulato**
- Emerging markets

**Partners, customers and collaborators will cha**
- and will have divergent technical aptitudes

**Must be work with precompetitive consortia**
- Can they perform common tasks for the community
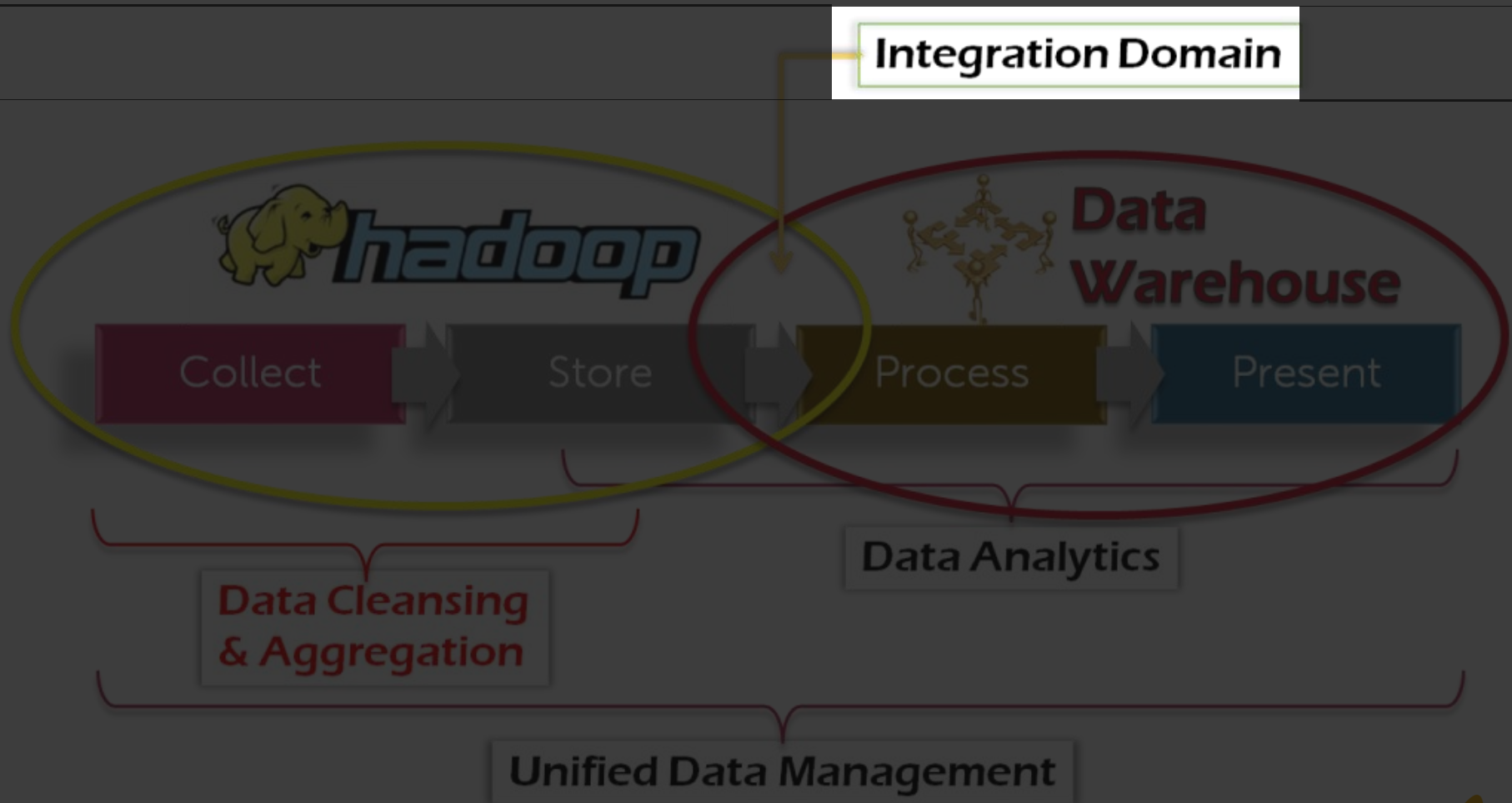
**Must be able to work with legacy data**
- Lots of unmined gems here!
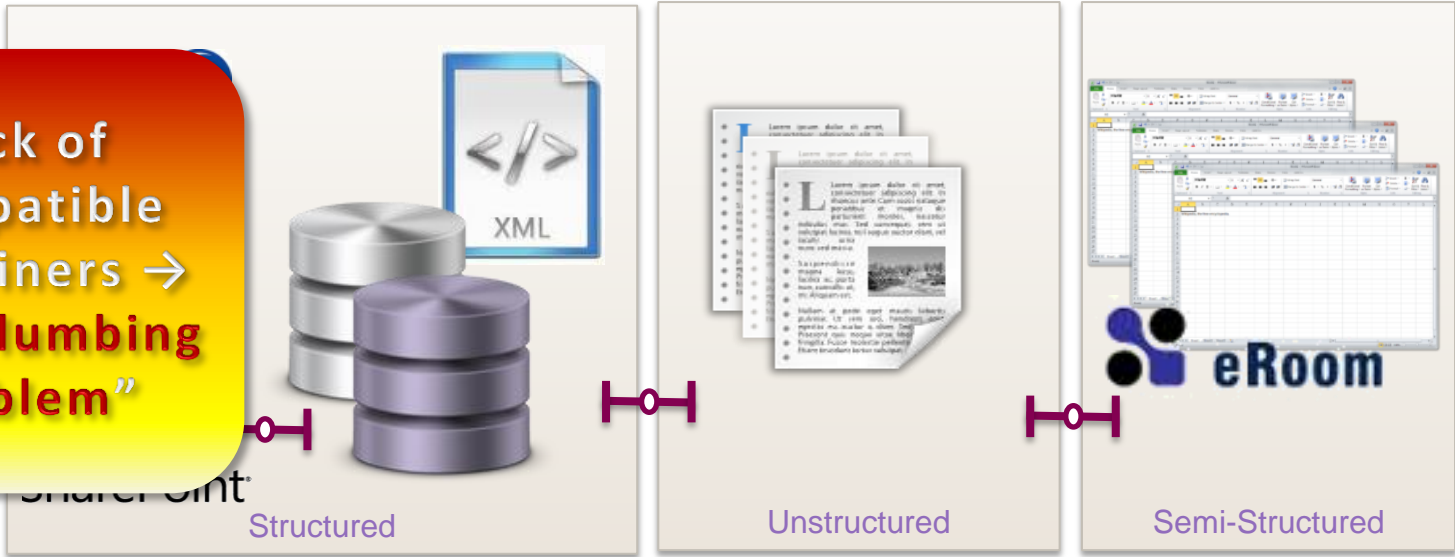
These are Big Data **Variety** and **Veracity** Challenges

# Typical Big Data Integration Process
## Document-Centric model



Integration Domain

Collect · Store · Process · Present

Data Cleansing & Aggregation

Data Analytics

Unified Data Management

R&D | RDI

# Integration Quandary:
## Content Does Not Combine Easily



**Lack of Compatible Containers →** the "**Plumbing Problem**"

**Lack of Compatible Semantics →** the "**Meaning Problem**"

Content

Structured

Unstructured

Semi-Structured

XML

eRoom

Models

HL7

MedDRA

ICD-9 & ICD-10

Fit-for-Purpose to "Standards"

# What's Needed?
# Linked Data!



LOD Cloud 2014

R&D | RDI

# Linked Data Demystified

## Addresses *Plumbing* and *Meaning* Challenges

### A Graph is the fundamental data model
•Not a table or a hierarchy or a document

### This model uses RDF* and is the web data model
•The underlying content need not be RDF only its published interface

### Web uniform resource identifiers (URIs) name things
•Resolving the URI gives a useful description

### URIs link data and integrate with other Linked Data
•Two things sharing the same URI are the same thing

### Fit-for-Purpose yet Scalable Applications are enabled
•Easy to mash-up and scales as the web scales

### Features Flexible & Adaptable Information Models
•You can change the data model without breaking downstream applications

### Encourages Shared Understanding via shared vocabularies
•Communities build out as needed to support their business questions

R&D | RDI

# Solution: IMI Open PHACTS Project

**Open PHACTS Mission:**
**Integrate Multiple Research**
**Biomedical Data Resources**
**Into A Single Open, Sustainable &**
**Free**
**Access Point**

# The Open PHACTS Discovery Platform

- **Cloud-Based "Production" Level System. Secure & Private**

- **Guided By Business Questions**

- **Uses Semantic Web Technology And provides a simple REST-ful API for the everyone else**

Open PHACTS has built a cutting edge, unique, flexible and powerful infrastructure for semantic data in life sciences
We have the platform, data, services, experience and capabilities to tackle big data challenges

- Cloud-based '**production**' level secure and private system
- We are expanding into new data areas
- Pathways, patents, disease
- Comprehensive workflow components for advanced use cases
- Open PHACTS is big data



**Future Concepts & Data Sources to support questions**

Disease · Gene · Protein · Pathway · Compound · Compound · Patent

http://www.openphactsfoundation.org/

- Delivered on all of the IMI project deliverables on or ahead of time
- Open PHACTS Foundation is a UK-based member-owned non-profit company founded to sustain and develop the Open PHACTS infrastructure
- OPF is now a <u>registered charity</u> with aims to further the public understanding of science through research
- OPF operations are distributed across funded partners to maintain and further develop public:private partnership
- A sought-after partner for Horizon 2020 projects
- Bring <u>industry perspective</u> to new academic partners
- Leverage the heritage of open data services and public:private partnership
- Significant academic research and professional network of subject matter experts



**BIG DATA EUROPE**
Empowering Communities
with Data Technologies

*OPF is a funded partner in the BDE project, another project approved and completing grant agreement, and 4 other pending submissions*

http://www.openphactsfoundation.org/

# The Open PHACTS Foundation

*OPF is a not-for-profit membership organisation, supporting the Open PHACTS Discovery Platform:*
*A sustainable, open, vibrant and interoperable information infrastructure for applied life science research and development.*

To reduce the barriers to drug discovery in industry, academia and for small businesses, the Open PHACTS Discovery Platform provides tools and services to interact with multiple integrated and publicly available data sources. To integrate this data, extensive cross-referencing of scientific concepts is needed across all databases.

The Open PHACTS Foundation ensures the sustainability of the Open PHACTS Discovery Platform infrastructure and acts as a hub for relevant scientific research and development.



## Key Resources

⚗ Open PHACTS API

🐙 Open PHACTS Repository

## Subscribe to the Foundation Newsletter

email address

**Subscribe**

## Contact us

✉ Email:
info@openphactsfoundation.org

🐦 Twitter: @Open PHACTS

# Solution: Emerging Public Solutions

http://bio2rdf.org

# EMBL-EBI RDF Platform

http://www.ebi.ac.uk/rdf/

# ChEMBL 18 RDF Model

# National, International Health Systems
## Data.gov, Data.gov.uk, WHO, datahub.io



https://catalog.data.gov/organiza...

http://data.gov.u...

http://gho.aksw.org/

http://datahub.io/dataset?q=health

# Best Practices

# Node and Edge Informatics
## Interfaces within the Drug Development Process

| Target Discovery | Lead Discovery | Lead Optimization | Pre-Clinical Development | Clinical Development | Registration | Marketing & Sales |
|---|---|---|---|---|---|---|
| NGS Exome analysis | RNAi | SAR | GLP Tox | IND | NDA/BLA | PMR REMS |
| Pathway Analysis | Assay Development | In vivo non-human testing | Formulation | Safety, Tolerability | MAA | PSUR |
| Structure Analysis | HTS | Exploratory PK | ADME | Phase I-III | | Observational Research |
| | | Exploratory Tox | PK | | | |
| Pathway Enrichment | | | Efficacy | | | |
| Disease Contextualization | | | | | | |

Seamless information connectivity (an EDGE) needed across domain NODEs

R&D | RDI

# It's the Data…
## …not the app, not the container

# Disposable Applications
## Questions, Answers, Insights Persist



Capture Business Questions and Sources

Domain Expert Concept Map

Build Formal Ontology

Challenge with Linked Data

Examine with a Faceted Browser

Share insights with a Knowledge Base

# Cooperation…
## …without Coordination

# Take-Aways

## Get your plumbing right

- And you won't be stuck in a silo

## Leverage working public solutions

- No need to reinvent the wheel

## Use Edge Informatics

- Consider handoffs—you don't know how your data will be used in the future

# Thanks

Big Data in Pharma 2015
Conference Organizers

AZ Linked Data Community

Key Influencers
David Wood
Toby Segaran
Tim Berners-Lee
Lee Harland
Bryn Williams-Jones
Eric Neumann
Dean Allemang
Barend Mons
Bernadette Hyland
Bob Stanley