# Open PHACTS
## Open Pharmacological Space

# Incorporating Commercial and Private Data into an Open Linked Data Platform for Drug Discovery

Carole Goble, **Alasdair J G Gray**, Lee Harland, Karen Karapetyan, Antonis Loizou, Ivan Mikhailov, Yrjänä Rankka, Stefan Senger, Valery Tkachenko, Antony J Williams, and Egon L Willighagen
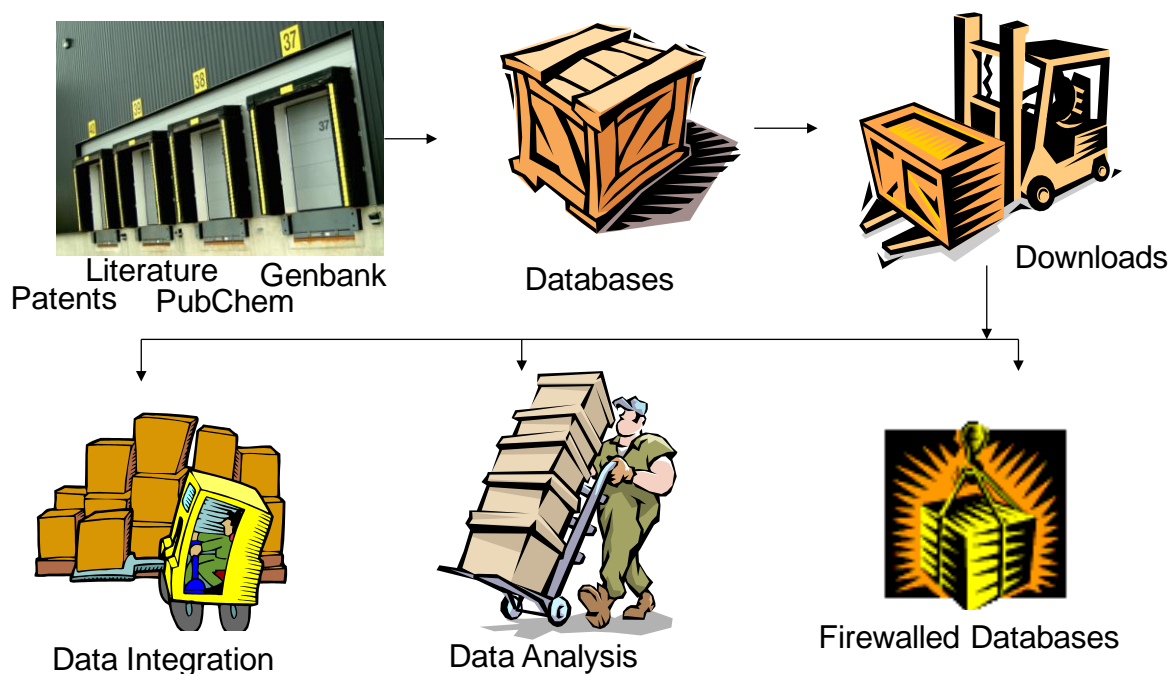
[www.openphacts.org](www.openphacts.org)          [A.J.G.Gray@hw.ac.uk](A.J.G.Gray@hw.ac.uk)

@open_phacts          @gray_alasdair

# Pre-competitive Informatics

Pharmaceutical companies are all accessing, processing, storing & re-processing external research data

Patents  Literature  Genbank  Databases  Downloads  X  **Repeat @ each company**
PubChem

Data Integration  Data Analysis  Firewalled Databases

# Open PHACTS objective



Open
Standards

Apps

Interactive
responses

Domain API

Provenance of
data

Drug Discovery Platform

Production
quality

# Drug Discovery Data

Pathways

Proteins

Interactions

Genes

Pharmacological Activities

Transcripts

Clinical Drug Applications

Biological Processes

Drugs

Indications

Pathological Processes

Compounds

Diseases

# Public Data

# Real Business Questions

# OPS Discovery Platform



Open PHACTS
Open Pharmacological Space

**Apps**

**Core Platform**

Identity Resolution Service

*"Adenosine receptor 2a"*

**Linked Data API** **(RDF/XML, TTL, JSON)**

**Domain Specific Services**

Identifier Management Service

P12374
EC2.43.4
CS4532

**Semantic Workflow Engine**

**Data Cache**
**(Virtuoso Triple Store)**

Chemistry Registration Normalisation & Q/C

**Indexing**

**Public Ontologies**

| VoID | VoID | VoID | VoID | VoID |
| --- | --- | --- | --- | --- |
| | Nanopub | | Nanopub | Nanopub |
| **RDF** | **RDF** | **RDF** | **RDF** | **RDF** |
| Db | Db | Db | Db | |

**User Annotations**

Public Content

Commercial

# Present Content: Public Data

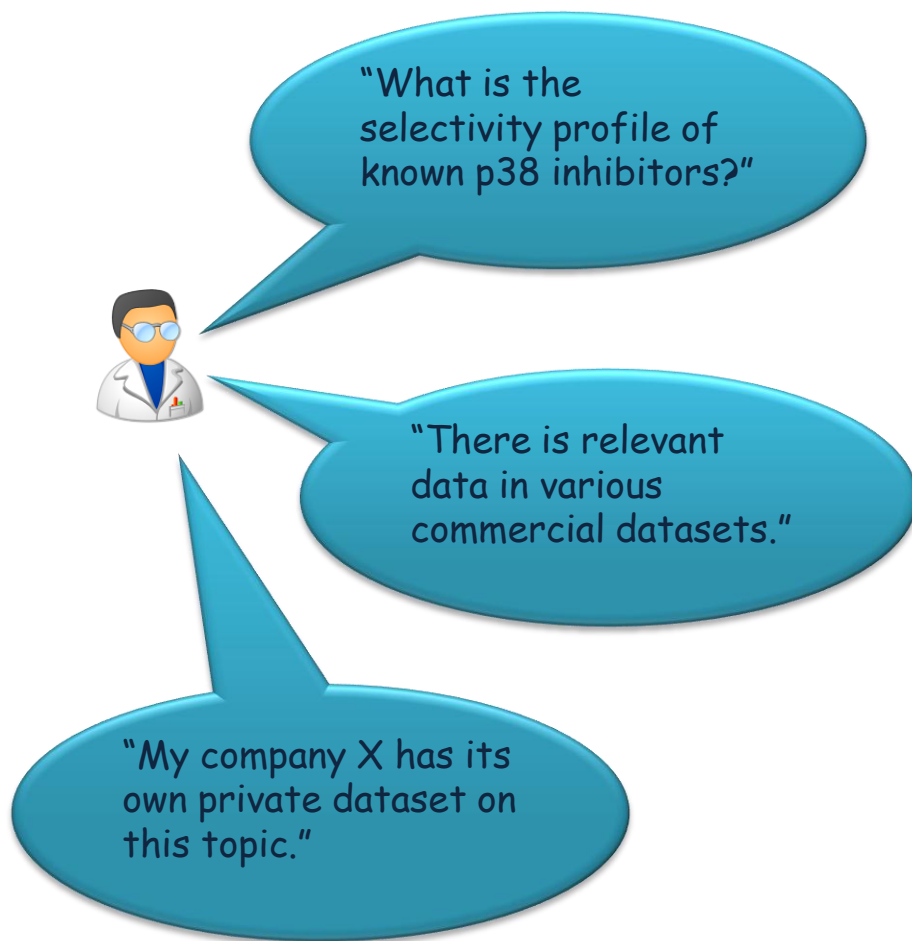| Source | Initial Records | Triples | Properties |
|---|---|---|---|
| ChEMBL | 1,247,403 | 305,419,649 | 77 |
| DrugBank | 19,628 | 517,584 | 74 |
| UniProt | ? | 533,394,147 | 82 |
| ENZYME | | 73,838 | 2 |
| ChEBI | | 40,575 | 2 |
| GeneOntology | 38,137 | | 26 |
| GOA | ? | | 15 |
| ChemSpider | 1,194,437 | 161,336,857 | 26 |
| ConceptWiki | 2,828,966 | 3,739,884 | 1 |
| WikiPathways | 946 | 1,449,981 | 34 |

Over a billion triples

# Semantic Integration Methodology

1. Define use cases
2. Identify Data
   – Create RDF
   – VoID dataset descriptions
3. Create mappings
   – between data set and known data sets (instance level)
   – index for text to URL conversion
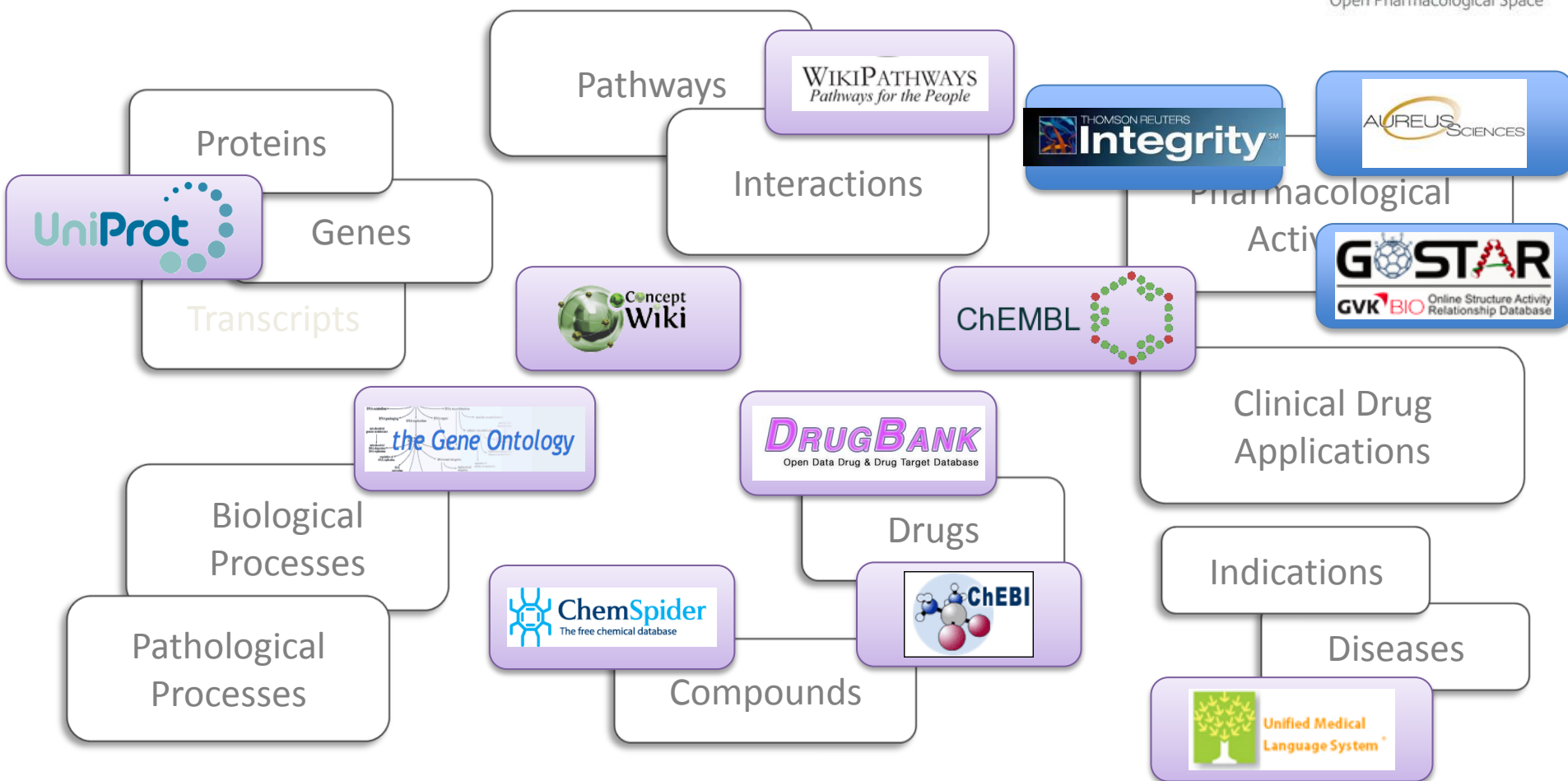
# Semantic Integration Methodology

4. Ingest RDF into data cache
   (i.e. triple store)
5. Define access paths to core concepts in data
6. Extend or create SPARQL queries for API calls
7. Publish API calls

# Commercial Data Use Case

"What is the selectivity profile of known p38 inhibitors?"

"There is relevant data in various commercial datasets."

"My company X has its own private dataset on this topic."

- Comprehensive data coverage
  - Commercial data collections
  - Extensive private collections

- Control data responses
  - Only authorised data

# Commercial Data Sets Pilot

# Linked Open Data

★           make your stuff available on the Web (whatever format) under an open license

★★        make it available as structured data (e.g. Excel instead of image scan of a table)

★★★     use non-proprietary formats (e.g. CSV instead of Excel)

★★★★    use URIs to denote things, so that people can point at your stuff

★★★★★ link your data to other data to provide context

http://5stardata.info/

# Commercial Linked Data

- Same conversion challenges as Open Data!
  - Goal to have 5⭐ linked data
  - www.openphacts.org/specs/rdfguide/
- Pilot (sample) data provided as data dumps
  - XML
  - CSV
  - RDF
- Structurally similar to ChEMBL
- Converted to interoperable RDF

# Data Modelling Challenges

- Contain private terminologies
  - Mapped to public equivalents
  - On going work

- Units represented as strings
  - Not always consistent, e.g. IC50, IC_50, IC-50
  - QUDT extended, e.g. IC$_{50}$
  - www.openphacts.org/specs/units/

# Dataset Descriptions

www.openphacts.org/specs/datadesc/

**Enable**

- Discovery
  - Name
  - Description
  - Coverage
- Access control
  - License
  - File locations
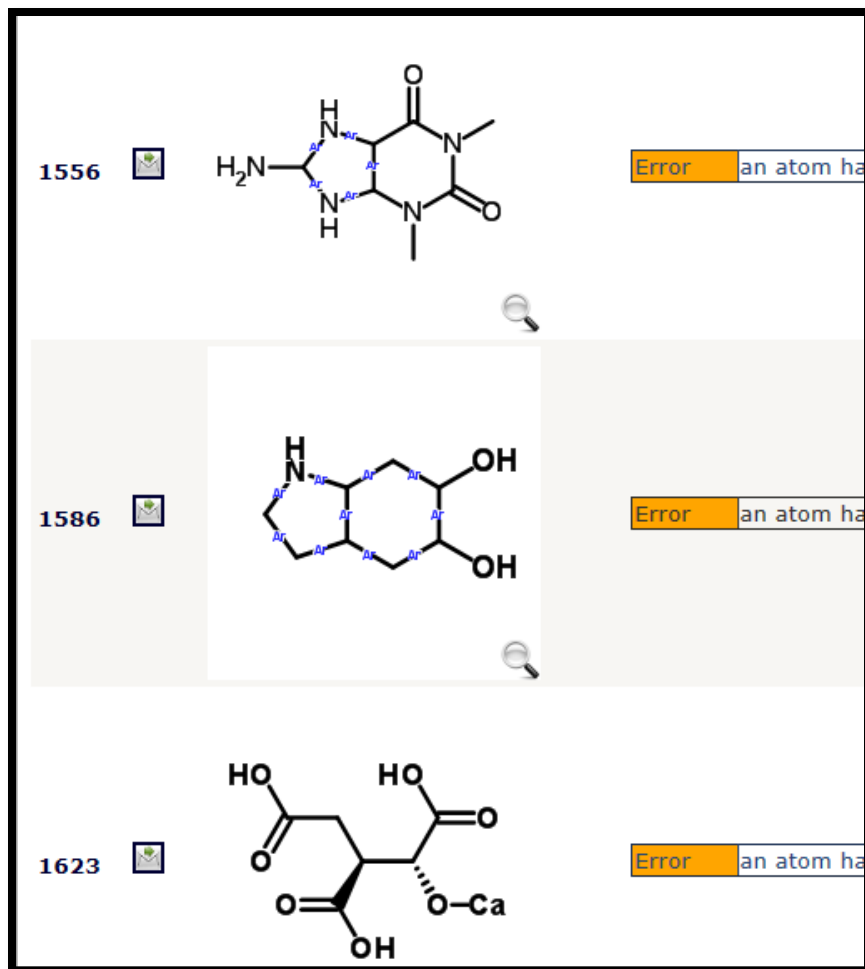- Answer Provenance
  - Returned data links to description

**Commercial Data Description**

- Publicly discoverable
  - Advertisement for data
  - Bring in more customers
- Restricted access by license

**Private Data Description**

- Hidden to all but authorised
- Restricted access

# Chemical mappings



- Data is messy!
- Identify common problems:
  - Charge imbalance
  - Stereochemistry
- Link based on structure
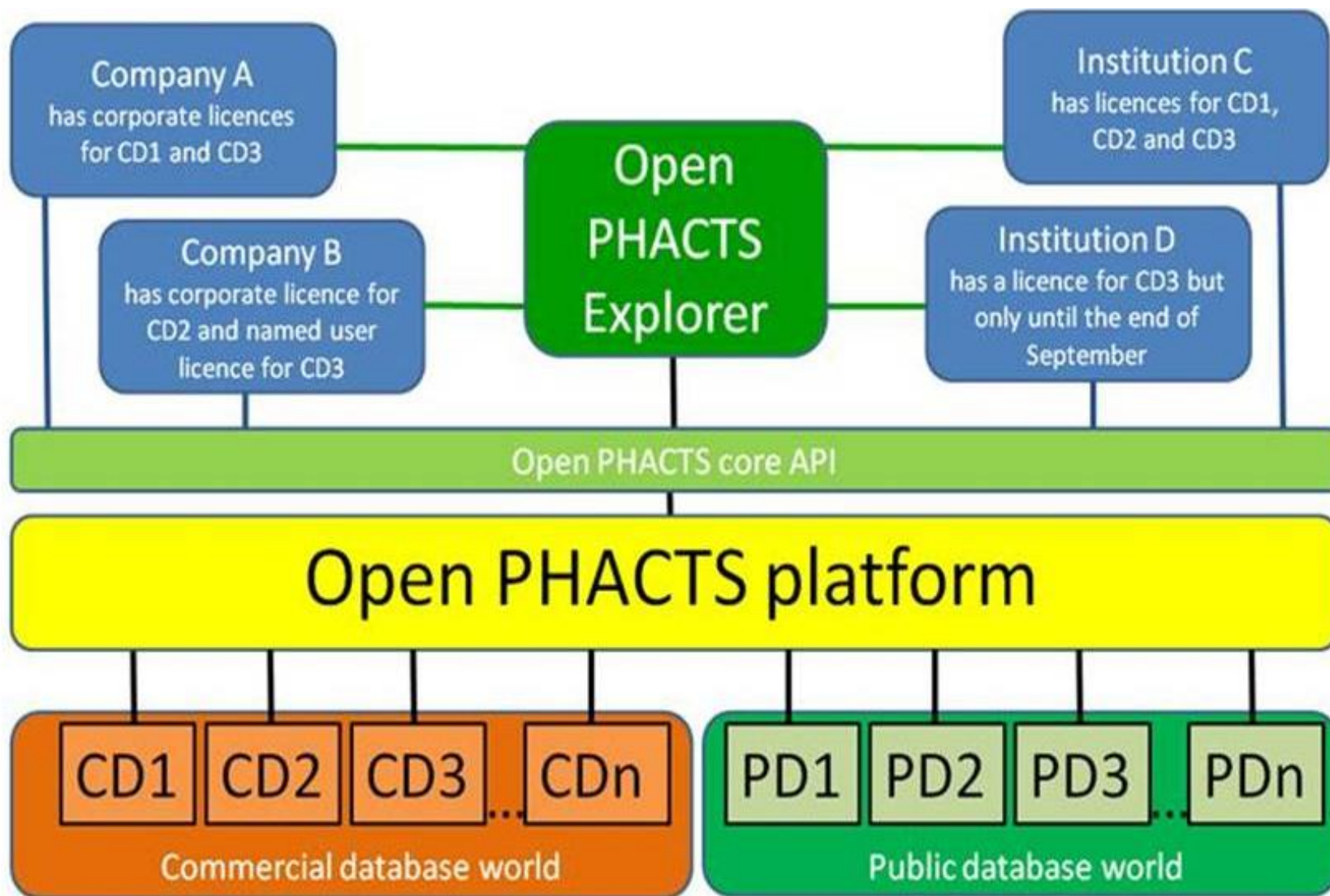
# Chemistry Registration

**ChemSpider Service**

- Validates and standardizes chemical representations
- Manual curation by RSC staff
- Data loaded in ChemSpider
- Open data: unsuitable for
  - Commercial data
  - Private data

**Chemical Registration Service**

- Utilizes ChemSpider Validation and Standardization platform
- Utilizes FDA rule set as basis for standardization
- Generates OPSID for chemicals
- Computes properties

# Access Requirements

# Data Access

- Each data set loaded into separate graph in cache

- Pilot data same form as open ChEMBL data
  – Extend queries with sub-queries for each set

- Restricted access
  – Virtuoso offers graph-based access restriction
  – Commercial data sets turned on/off

# Conclusions

- Drug discovery requires full data coverage
  - Public/open data
    - Open description
    - Open data
  - Commercial data
    - Open description
    - Restricted data
  - Private data
    - Restricted description
    - Restricted data
- Pilot study with three commercial datasets

# Conclusions

- Data Modelling
  - Similar challenges as public data

- Access restriction
  - Provided by standard mechanisms
  - Graph-based access

- Open PHACTS Discovery Platform
  - Releasing version 1.3 (late 2013)
  - Version 1.4 will contain commercial data (2014)

# Acknowledgements



- GVK Bio
  GOSTAR
  gostardb.com



- Thomson Reuters
  Integrity
  integrity.thomson-pharma.com



- Aureus Sciences Elsevier
  AurSCOPE
  www.aureus-sciences.com

# Questions

A.J.G.Gray@hw.ac.uk

www.macs.hw.ac.uk/~ajg33

@gray_alasdair

pmu@openphacts.org

www.openphacts.org

@open_phacts

## Open PHACTS Project