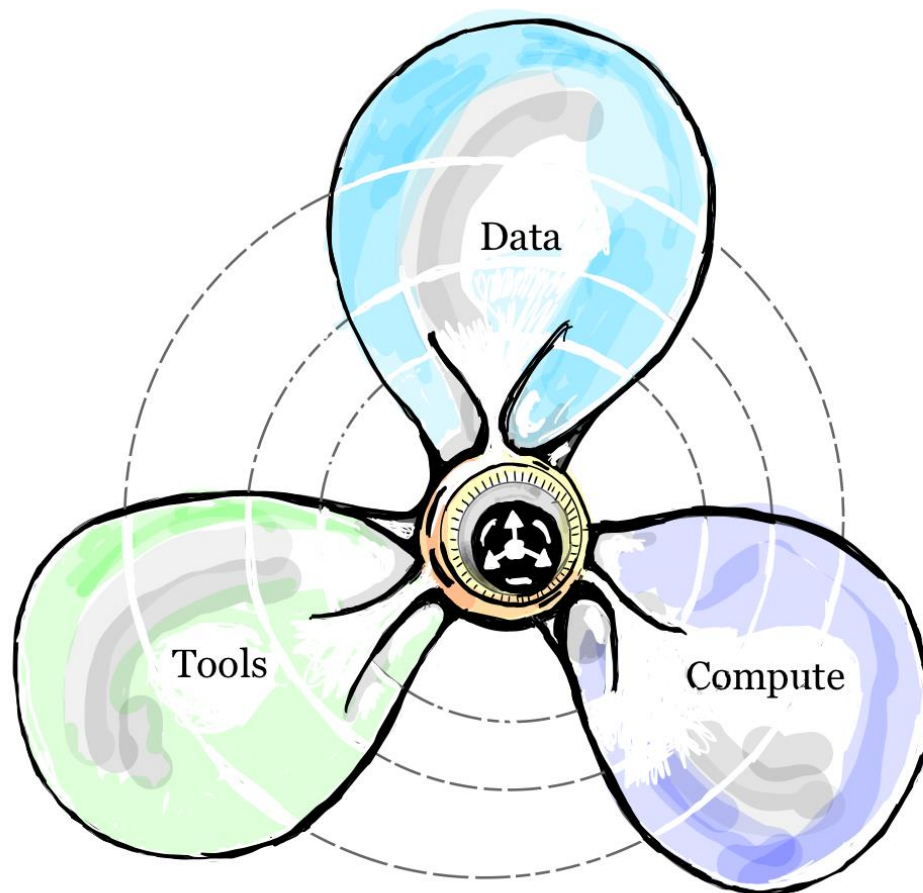




Open PHACTS lessons learned in pioneering semantic technology for drug discovery

Barend Mons
Stockholm, PSWC

What we need: The Internet of FAIR data and Services





<https://vimeo.com/162062013>

Pre-history

BMC Bioinformatics

HOME

ABOUT

ARTICLES

SUBMISSION GUIDELINES

COMMENTARY | OPEN ACCESS

Which gene did you mean?

Barend Mons 

BMC Bioinformatics 2005 6:142 | DOI: 10.1186/1471-2105-6-142 | © Mons; licensee BioMed Central Ltd. 2005

Received: 26 May 2005 | Accepted: 07 June 2005 | Published: 07 June 2005

Abstract

Computational Biology needs computer-readable information records. Increasingly, meta-analysed and pre-digested information is being used in the follow up of high throughput experiments and other investigations that yield massive data sets. Semantic enrichment of plain text is crucial for computer aided analysis. In

Download PDF

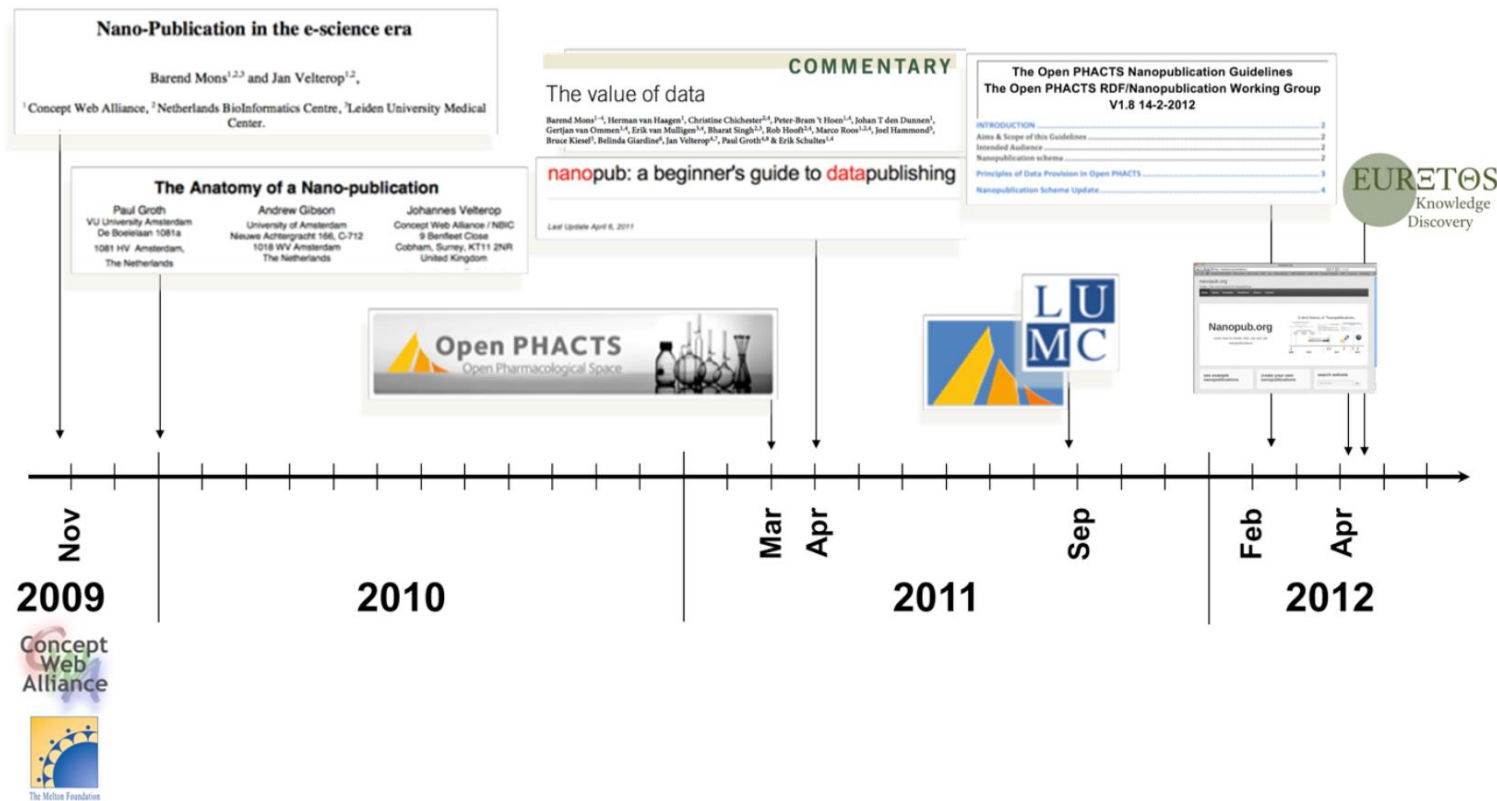
Export citations 

Table of Contents

- Abstract
- Text mining?Why bury it first and then mine it again?
- Too much to read
- Semantics: a crucial addition to text
- Databases and pathway tools are not enough
- However, text is a nightmare for

Made it to be my most hated statement

History



The FAIR principles became the unifying force

World 2017



Realising the European Open Science Cloud

First report and recommendations of the Commission High Level Expert Group on the European Open Science Cloud



European Open Science Cloud
Latest news 19 April 2016 – European Open Science Cloud

Giving a major boost to Open Science in Europe, the Commission today presented its blueprint for cloud-based services and world-class data infrastructure to ensure science, business and public services reap benefits of big data revolution.

By bolstering and interconnecting existing research infrastructure, the Commission plans to create a new European Open Science Cloud that will offer Europe's 1.7 million researchers and 70 million science and technology professionals a virtual environment to store, share and re-use their data across disciplines and borders. This will be underpinned by the European Data Infrastructure, deploying the high-bandwidth networks, large scale storage facilities and super-computer capacity necessary to effectively access and process large datasets stored in the cloud.

GO FAIR

EURETOS Knowledge Platform

OUTOFORCE

The Open PHACTS Foundation



Jointly Designing a Data FAIRPORT
Workshop: 13 - 16 January 2016, Leiden, The Netherlands

The FAIR Guiding Principles
The FAIR Guiding Principles for Findable, Accessible, Interoperable and Reusable data

DTL Data Science Community
Home Current issue archive

nature genetic
Home Current issue archive

FORCE11
The State of Research Communications and e-Scholarship

The Commons
Biomedical discovery enabling sharing objects

European Open Science Cloud
Home Open Access

The FAIR data principles are simple guidelines for ensuring that machines can find and use data, supporting data reuse by individuals. More—and better—research can be generated by designing data and algorithms to be findable, accessible, interoperable and reusable, with the tools and workflows that led to these data.



something to refer to

The FAIR Guiding Principles for scientific data management and stewardship

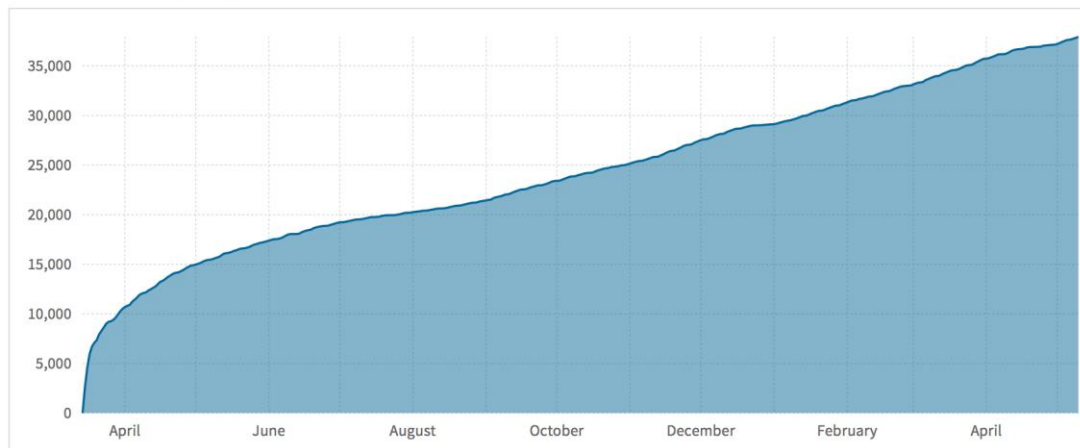
MD Wilkinson, M Dumontier, IJJ Aalbersberg, G Appleton, M Axton, ...
Scientific data 3

114

2016

Page views (37,928)

View as: Cumulative | Line | Table

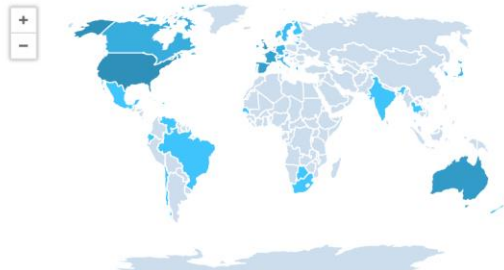


Mentions in news, blogs & Google+

News articles (1) Scientific blogs (19) Google+ posts (5)

Guidelines for best practices in the publication of scientific data
Phys.org

Twitter demographics



Online attention



Altmetric score (what's this?)

- Tweeted by 756
- Blogged by 19
- On 13 Facebook pages
- Mentioned in 5 Google+ posts
- Picked up by 1 news outlets

Show more

This Altmetric score means that the article is:

- in the 99th percentile (ranked 615th) of the 278,235 tracked articles of a similar age in all journals
- in the 95th percentile (ranked 1st) of the 23 tracked articles of a similar age in *Scientific Data*



G20 HANGZHOU SUMMIT

HANGZHOU, CHINA 4-5 SEPTEMBER

'We support appropriate efforts to promote open science and facilitate appropriate access to publicly funded research results on findable, accessible, interoperable and reusable (FAIR)'



Open Science



Open Innovation

FAIR data from a pharma perspective

from: Herman van Vlijmen (Janssen and chair Open FACTS Foundation)



- ✦ A lot of data in certain domains
- ✦ A lot of it is “unFAIR”, “triple-F (Far From FAIR)”
- ✦ Extensive use of CROs makes data standards essential
- ✦ Internal development of FAIR data is in early stages
- ✦ Strong interest in **FAIR data concept** because of opportunities for answering complex scientific questions and multidomain use cases
 - Open PHACTS (IMI)
 - Open Targets application (EBI)
 - FAIRification of IMI/EFPIA data (new IMI call topic)
- ✦ Key driver of effort should be USE CASES



Affiliate Members



Open PHACTS (IMI project 2011-2016)

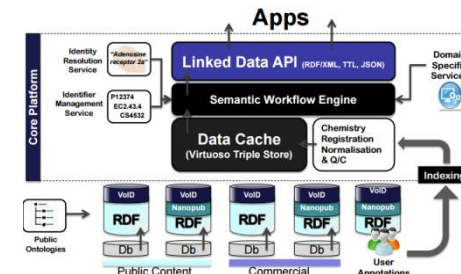
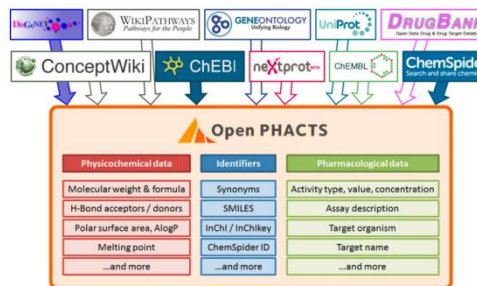


Objective: Integrate multiple research biomedical data resources into a single open, sustainable and free access point

Pharma partners: Pfizer, GSK, Lundbeck, Lilly, Janssen, AZ, Novartis, Esteve

Pharma needs:

- **Data Integration:** public data is heavily underused and needs processing to make it usable with internal data
- **Tools** to answer more complex questions that go across scientific domains



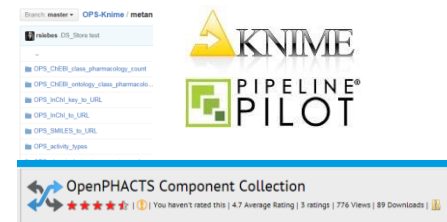
Integrated databases

Available to the world

Scientific competency questions as the basis for semantically enriched open pharmacological space development

Scientific questions

Kamal Azzaoui¹, Edgar Jacoby¹⁴, Stefan Senger², Emiliano Cuadrado Rodríguez³, Mabel Loza⁵, Barbara Zdrzil⁶, Marta Pinto⁶, Antony J. Williams⁵, Victor de la Torre⁶, Jordi Mestres⁷, Manuel Pastor⁷, Olivier Taboureau⁸, Matthias Rarey⁹, Christine Chichester¹⁰, Steve Pettifer¹¹, Niklas Blomberg^{12,a}, Lee Harland¹³, Bryn Williams-Jones¹³ and Gerhard F. Ecker⁴



Answerable with query tools and workflows

Priority use cases collected in Open PHACTS project

TABLE 1

The top 20 research questions

Question number	Question
Cluster I	
Q1	Give me all oxidoreductase inhibitors active <100 nM in human and mouse
Q2	Given compound X, what is its predicted secondary pharmacology? What are the on- and off-target safety concerns for a compound? What is the evidence and how reliable is that evidence (journal impact factor, KOL) for findings associated with a compound?
Q3	Given a target, find me all actives against that target. Find/predict polypharmacology of actives. Determine ADMET profile of actives
Q4	For a given interaction profile – give me similar compounds
Q5	The current Factor Xa lead series is characterized by substructure X. Retrieve all bioactivity data in serine protease assays for molecules that contain substructure X
Q6	A project is considering protein kinase C alpha (PRKCA) as a target. What are all the compounds known to modulate the target directly? What are the compounds that could modulate the target directly? I.e. return all compounds active in assays where the resolution is at least at the level of the target family (i.e. PKC) from structured assay databases and the literature
Q7	Give me all active compounds on a given target with the relevant assay data
Q8	Identify all known protein-protein interaction inhibitors
Q9	For a given compound, give me the interaction profile with targets
Q10	For a given compound, summarize all 'similar compounds' and their activities
Q11	Retrieve all experimental and clinical data for a given list of compounds defined by their chemical structure (with options to match stereochemistry or not)
Cluster II	
Q12	For my given compound, which targets have been patented in the context of Alzheimer's disease?
Q13	Which ligands have been described for a particular target associated with transthyretin-related amyloidosis, what is their affinity for that target and how far are they advanced into preclinical/clinical phases, with links to publications/patents describing these interactions?
Q14	Target druggability: compounds directed against target X have been tested in which indications? Which new targets have appeared recently in the patent literature for a disease? Has the target been screened against in AZ before? What

Scientific competency questions as the basis for semantically enriched open pharmacological space development

Kamal Azzaoui¹, Edgar Jacoby¹⁴, Stefan Senger², Emiliano Cuadrado Rodríguez³, Mabel Loza³, Barbara Zdrzil⁴, Marta Pinto⁴, Antony J. Williams⁵, Victor de la Torre⁶, Jordi Mestres⁷, Manuel Pastor⁷, Olivier Taboureau⁸, Matthias Rarey⁹, Christine Chichester¹⁰, Steve Pettifer¹¹, Niklas Blomberg^{12,a}, Lee Harland¹³, Bryn Williams-Jones¹³ and Gerhard F. Ecker⁴

Drug Discovery Today • Volume 18, Numbers 17/18 • September 2013

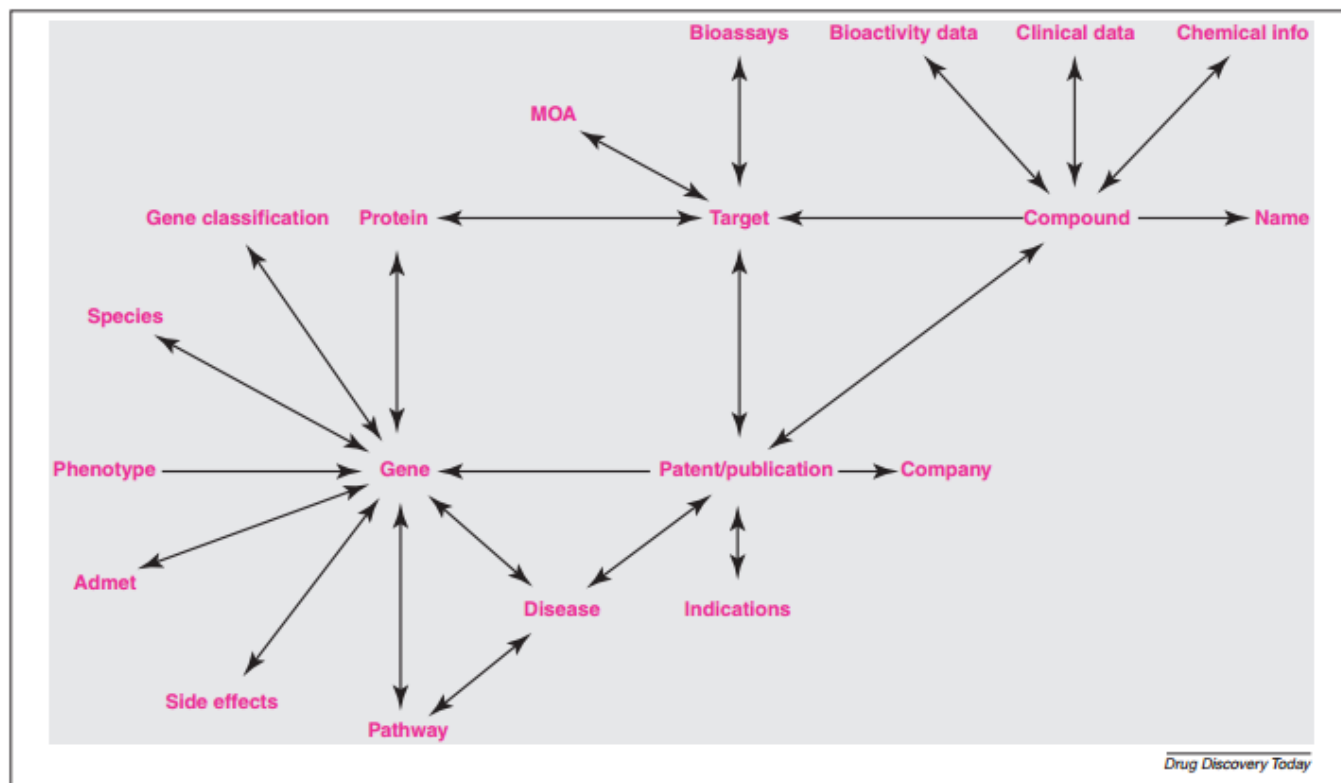


FIGURE 2

Network of data associations needed to answer the top-ranked scientific competency questions. The network reflects a cartoon that summarizes the data associations that are needed to target the top 20 research questions.

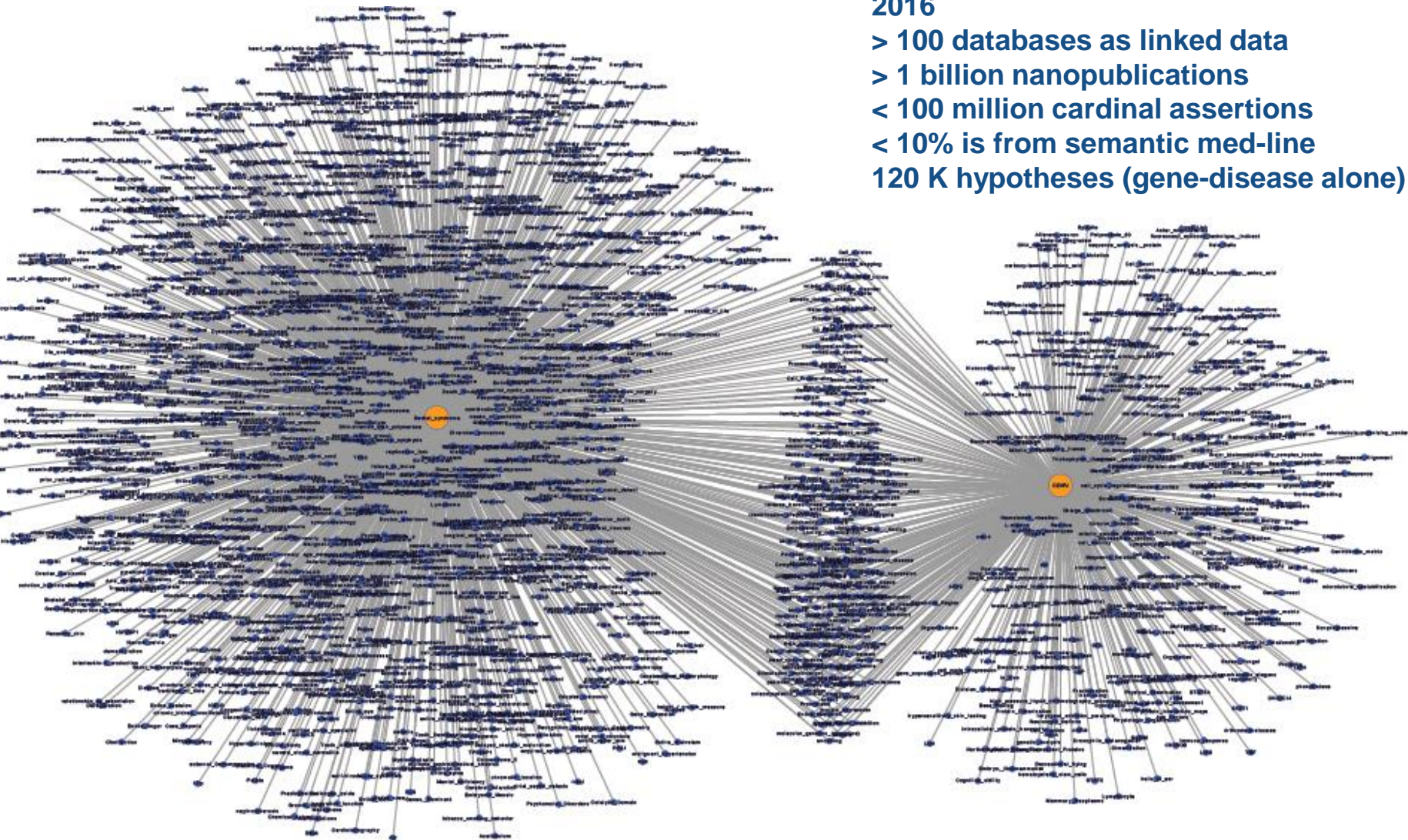
ROC RDT 1010 101010
AUC 111010 1111
HTTP://xzwvwx
 $\sqrt{e/2} \rightarrow 200d$

ATCG
TCAGAAAT
GCAATTCGTAAT
Phenotype drug TT
mitochondria





We publish about less than a million LSConcepts !



2016

> 100 databases as linked data

> 1 billion nanopublications

< 100 million cardinal assertions

< 10% is from semantic med-line

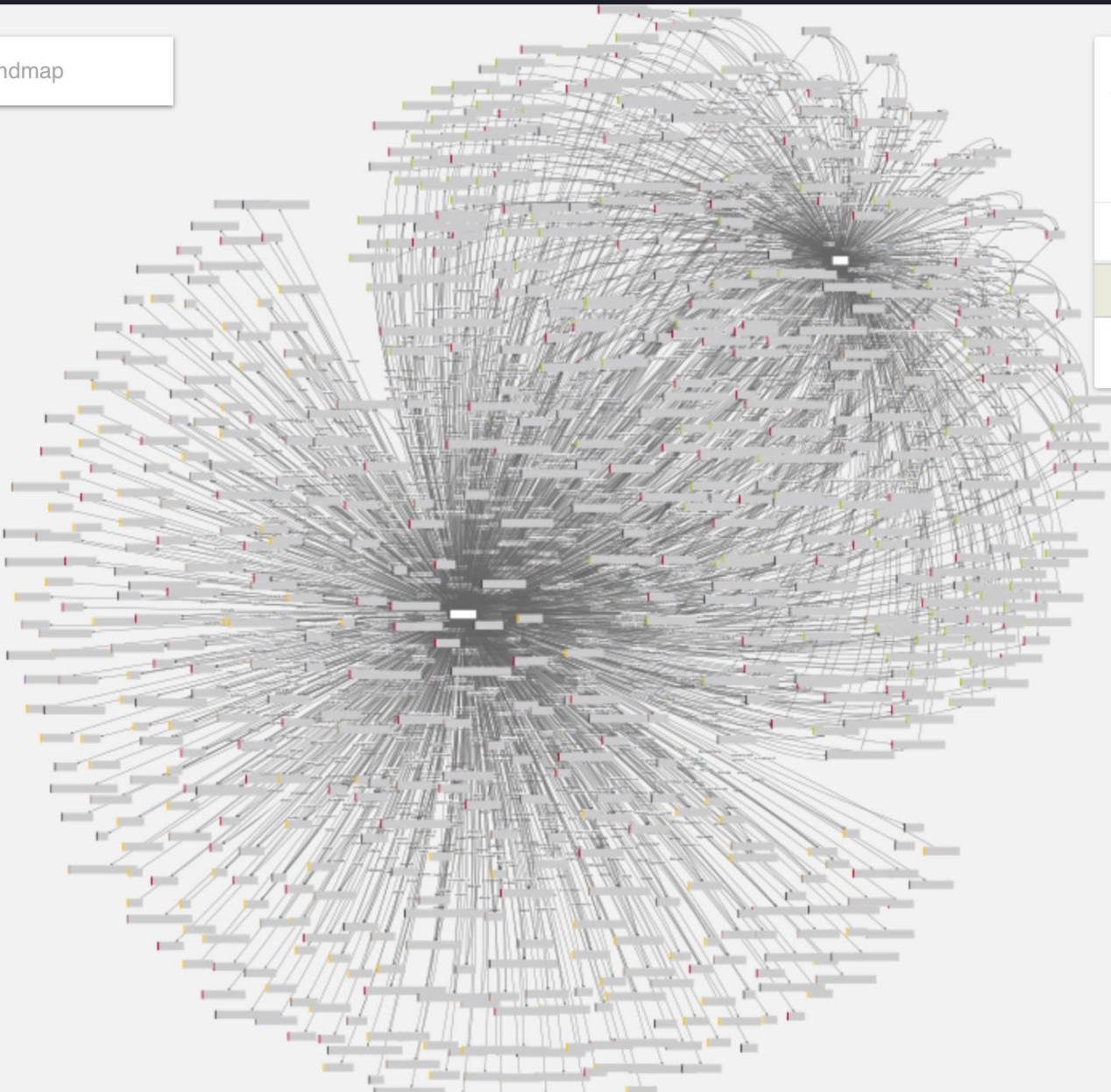
120 K hypotheses (gene-disease alone)

Made during Nir's talk on Metformin

< metformin Applications >

? v barendmons@gmail.com

mindmap



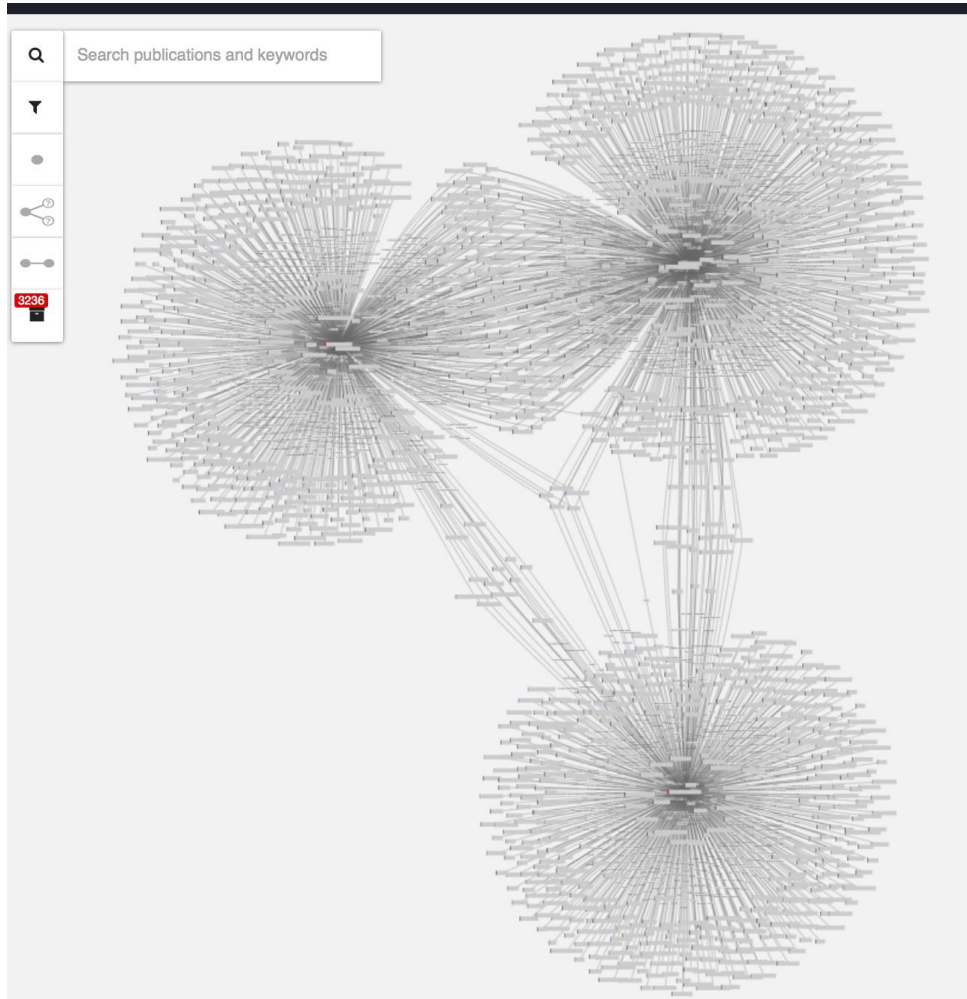
MAP DETAILS

🎯	Keywords	2 / 656
↔	Relations	0 / 2287

2 keywords

aging **i** metformin **i**

A real knowlet image (Euretos Knowledge Platform)



share: in this case 5 objects are shared between all three knowlets (in this case: metabolic syndrome



Summary Article

Genes predicted for [disease] based on known genes.

Abstract

Many genetic diseases, even those that are thought to be caused by mutations in a single gene do show a wide variation in disease phenotypes, severity of symptoms and for instance response to drugs and other interventions. These variation may be partly due to variations in other genetic factors. Not only variants in other associated genes, but also regulation factors and differences in metabolism and general physiology, which in turn may have a strong genetic basis. Also age factors and environmental factors may play a role.

Here, we have used [x] genes that are previously reported to be involved in the [disease] and closely related phenotypic disorders such as [subdisorder 1, sub disorder 2, closely related phenotype] to predict candidate genes that may have a direct or indirect influence in the phenotypic manifestation of [disease] patients.

[x] candidate genes were found with a prediction score ranging between [x] and [y] [reference to mothership paper to be written about methods used in workflows]

For each of those genes we have recovered all relevant direct and indirect associations between the gene, other genes in the cluster and relevant physiological and molecular processes.

Method

The detailed methods used for recovery and analysis of direct and indirect associations found in [x] pre-mined and annotated, interoperable data resources are described in [ref. 1]. Briefly, all databases were etc.

Results:

Summary (table comparable to table in seed article)

main disease	parkinson disease	more information	viewer
sub-phenotypes included	x,y,z	HPO	
number of genes used as reference set	17	nature genetics (DOI)	
predicted genes with direct connections	x	show	show
predicted genes	y	show	show
most frequent connecting semantic categories			
other genes and proteins	x	show	show
gene inhibiting molecules	x	show	show
pathways	x	show	show
cellular processes	x	show	show
chemicals, metabolites and drugs	x	show	show
tissues	x	show	show
cell types	x	show	show
organelles	x,y,z	show	show

Upload file

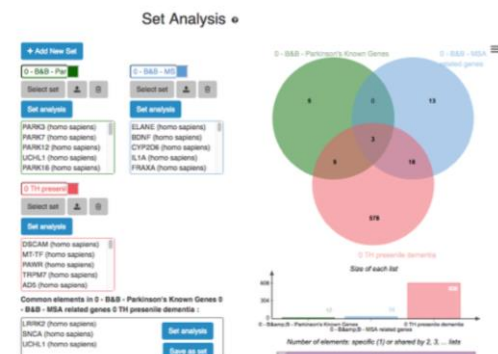
Get file Choose File PARKknowngenes.csv

Identified: 17
Conflicted: 0
Not identified: 0

Source Identifier	Name
atp13a2	ATP13A2 (homo sapiens)
elif4g1	EIF4G1 (homo sapiens)
fbxo7	FBXO7 (homo sapiens)
htra2	HTRA2 (homo sapiens)
lrrk2	LRRK2 (homo sapiens)
park10	PARK10 (homo sapiens)
park11	GIGYF2 (homo sapiens)
park12	PARK12 (homo sapiens)
park16	PARK16 (homo sapiens)
park3	PARK3 (homo sapiens)
park7	PARK7 (homo sapiens)
pink1	PINK1 (homo sapiens)
pla2g6	PLA2G6 (homo sapiens)
prkn	PARK2 (homo sapiens)
snca	SNCA (homo sapiens)
uchl1	UCHL1 (homo sapiens)
vps35	VPS35 (homo sapiens)

Invert selection Remove selected

Continue Cancel



H-factor
Journal impact factor
Tenure track
Nature, science....
Text, Tables, Figures



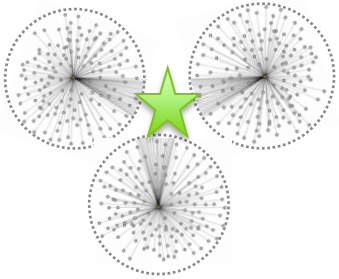
Discount on:
Spurious correlations
Ridiculograms

Not in stock:
Reproducible results
Actionable knowledge

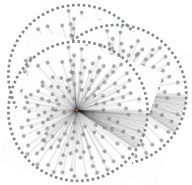




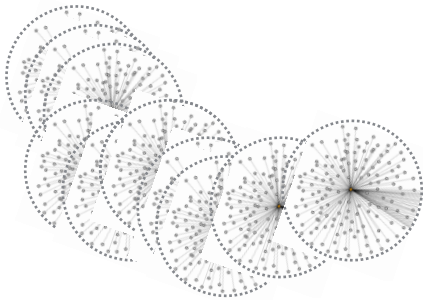
near-sameness and semantic drift of concepts represented by Knowlets



1. The implicitome > new hypothesis generated by the computer

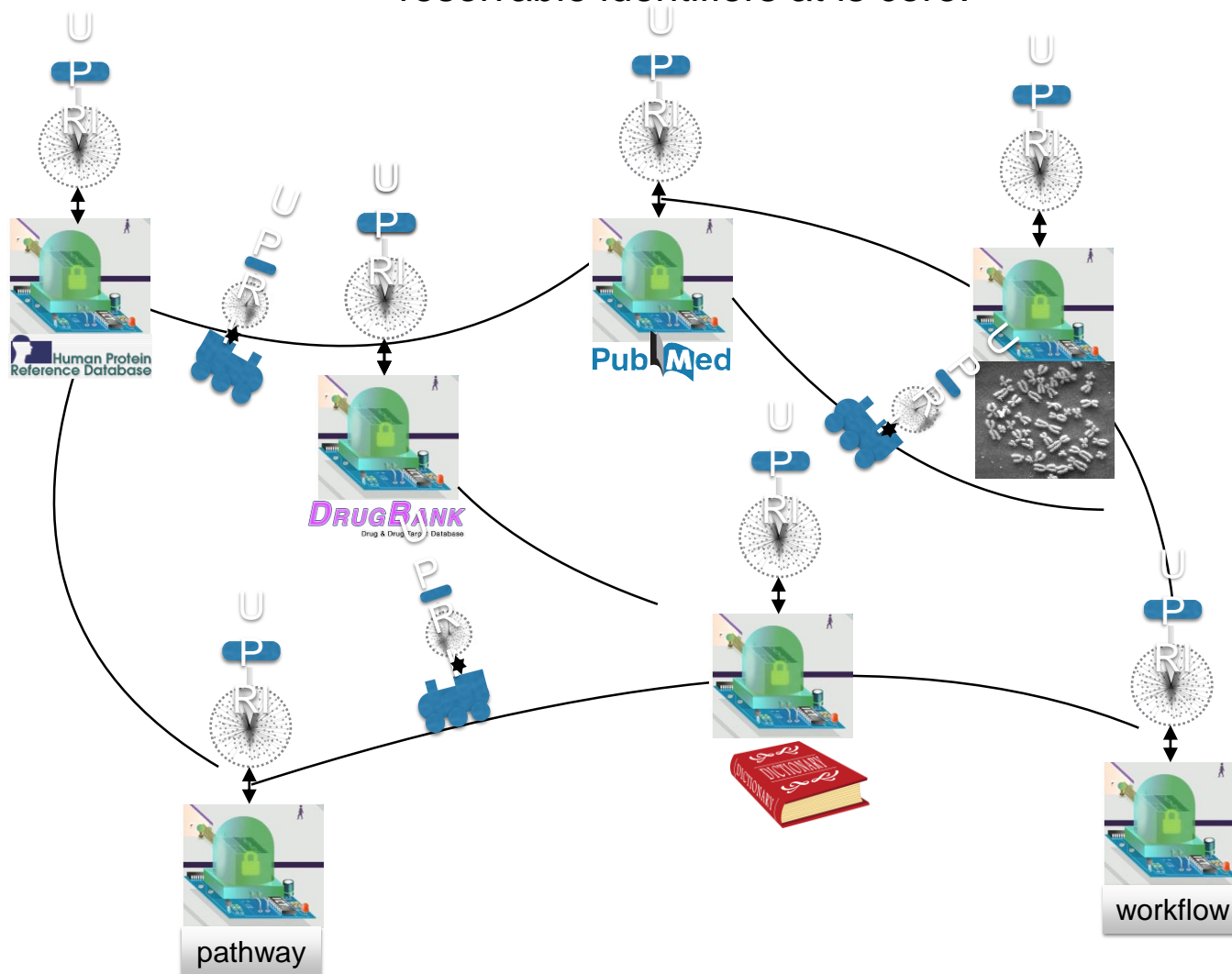


2. dealing with near-sameness (human versus mouse gene etc.)



2. dealing with semantic and conceptual drift)

The Internet of Data and Services (PHT model)
controlled at two levels: UPRIs and a web of associative Knowlets
representing abstract concepts, digital objects (including distributed VMs) and
physical objects, devices & things in a universal way with scalable, unique,
resolvable identifiers at its core.





Scholarly practice is changing profoundly as we embrace new methods of digital research and engage society.

Our centuries-old research communication practices that underpin scholarship are to be celebrated — but are they still fit for their purpose?

PSWC2022, Stockholm October 18th , in combination with FORCE 2020

Pre-symposium event:

The Nobel Prize in NON-Literature

The Swedish Academy (Svenska Akademien), Börssalen, Källargränd 4, Stockholm

More information: <http://svenskaakademien.se/>

The first Nobel prize ever awarded to a Social Machine !! (join us)

Reception Co-sponsor: The Innovative medicines Initiative