

Integrated Health Care Datasets for Knowledge Discovery

Huawei Global Technology Summit 2022
July 7th, 2022

Christophe Guéret
christophe.gueret@accenture.com



Agenda

We'll go through topics around making, and then reasoning over a complex Knowledge Graph

- 1 Assembling and using a graph
- 2 -> Make
- 3 -> Use
- 4 Can we do things differently?
- 5 Take away messages

Why integrate healthcare data?

Three factors

- Increasing amount of data available
- Increasing data processing capabilities
- Improved outcomes using the “big picture”

One “omic” to study each “ome”

- Genome: data about genes
- Exposome: data about exposition
- Proteome: data about proteins
- Etc...

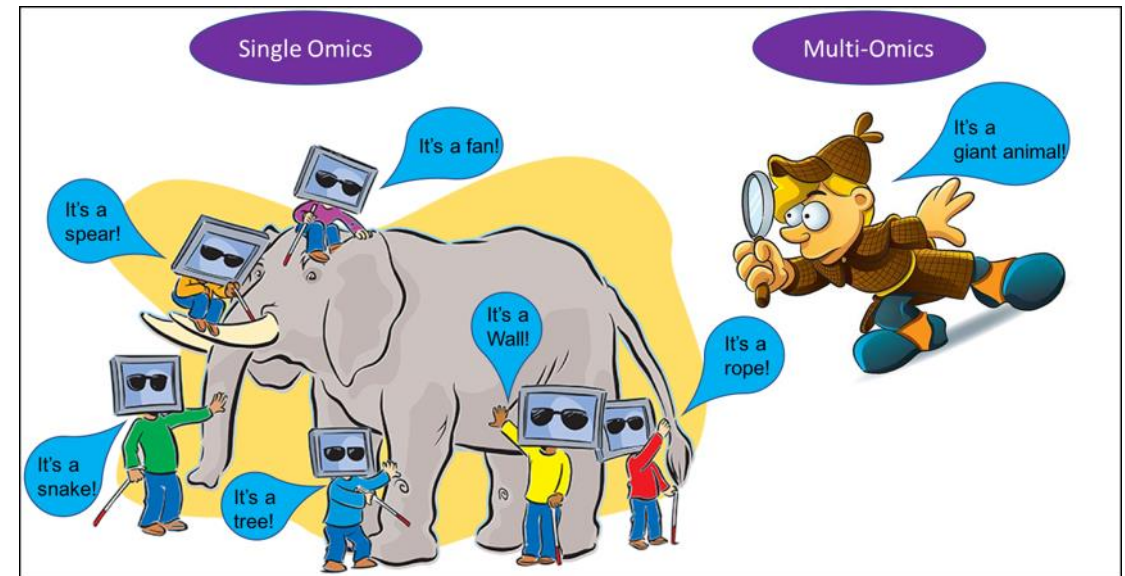
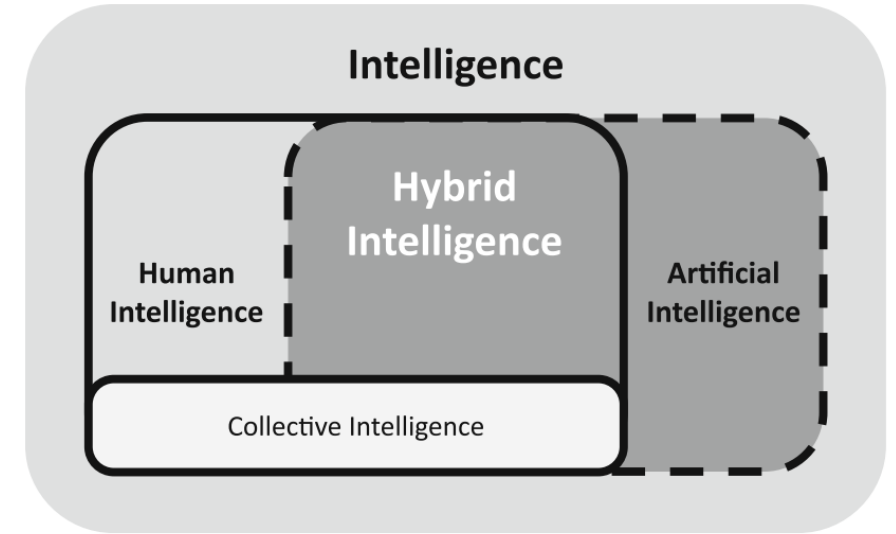


Image from [Multi-Omics: a Revolutionary Approach to Data Analysis](#)

KG-driven Knowledge Discovery

- Knowledge Graphs are an established way to connect data coming from different silos into so-called “360 views”
- Knowledge Discovery is the process of extracting useful information from this data.
- In our work, **we aim at building AI systems to work with humans on knowledge discovery tasks**. This ranges from exploring the data together to validating ideas using the data.

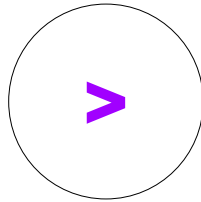


Positioning of Hybrid Intelligence [Dellermann 2021]

The story of the client and the Knowledge Graph geeks



"I need an AI system to help me work on a knowledge discovery task"

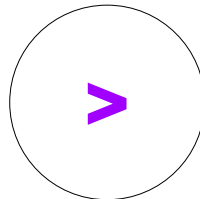


"We can make a graph and do graph machine learning on it"



"Cool! Here is all the data I have, feel free to enrich it with more stuff"

(sometime later...)



"Here you go; predictions! What do you think?"

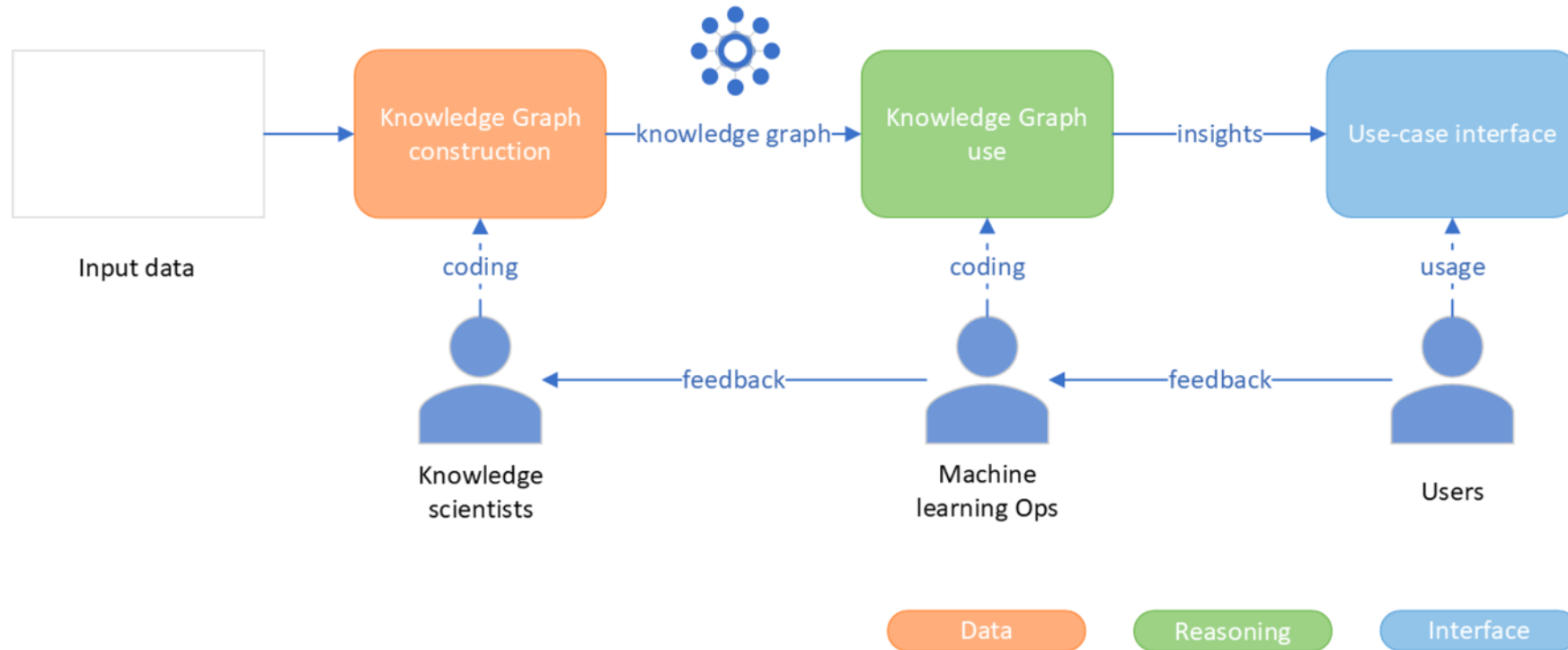
Assembling and using a graph

Getting from raw data to insights



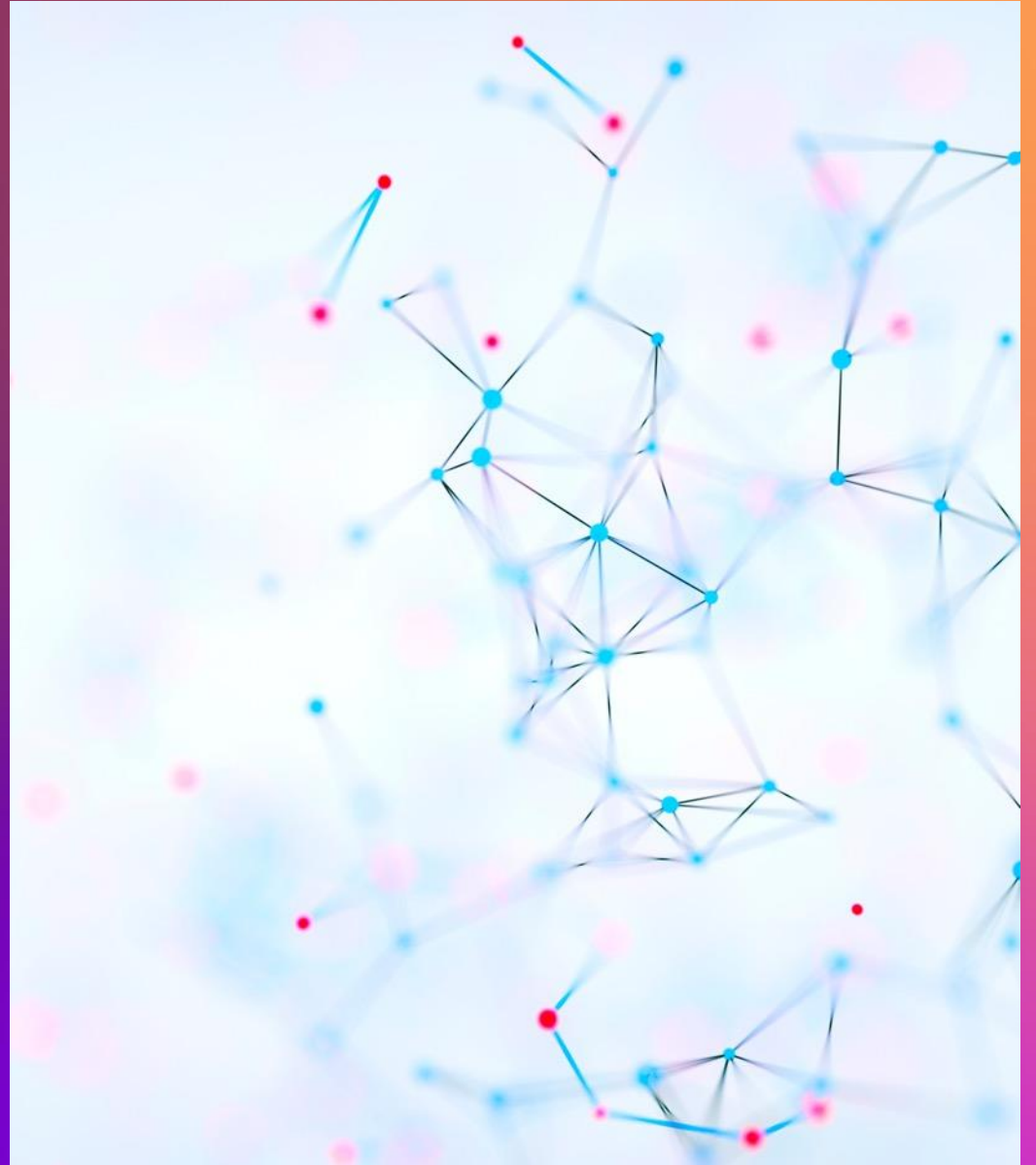
Two-step pipeline

- State of the art is to: build a graph, ship it to someone using it, and then ship the outcome of the AI part. Then back-track and repeat as needed



Make

Assemble a nice, big, Knowledge Graph



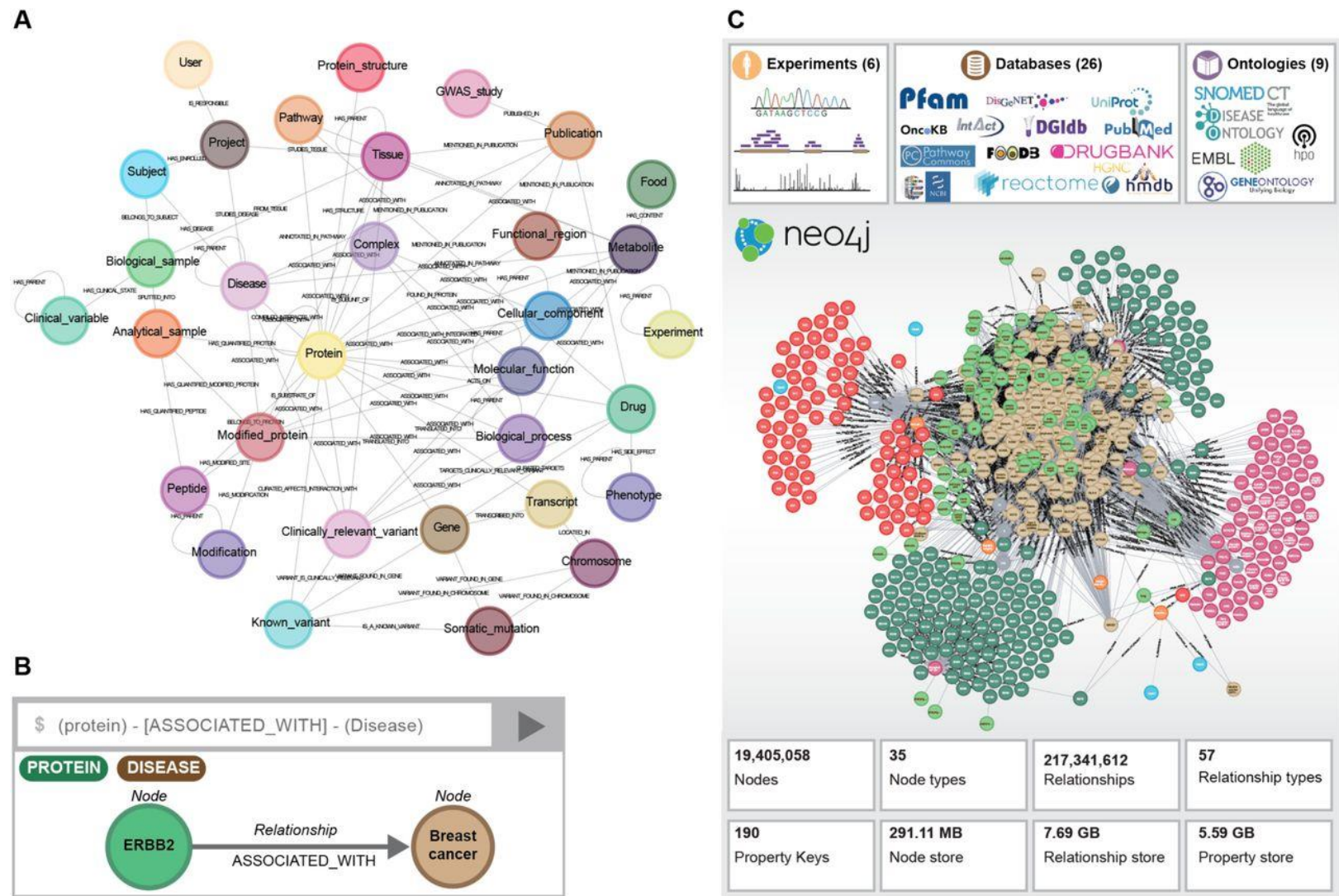
Graphs as take away items

- There are a lot of different integrated medical KGs out there
 - [PubChem](#)
 - [DISCOVER](#)
 - [DisGeNET - a database of gene-disease associations](#)
 - [Clinical Knowledge Graph \(CKG\)](#)
 - [Hetionet - An integrative network of biomedical knowledge](#)
 - [Open Pharmacological Space \(openphacts.org\)](#)
 - ...
- Collaboration networks such as Elixir are also interesting to study as a source of data and tools
- However, all those graphs are created with this one-size-fits all approach and share the other shortcomings of any “in house” graph constructed with SoTA approaches



This Photo by Unknown Author is licensed under CC BY-SA-NC

One example graph: CKG

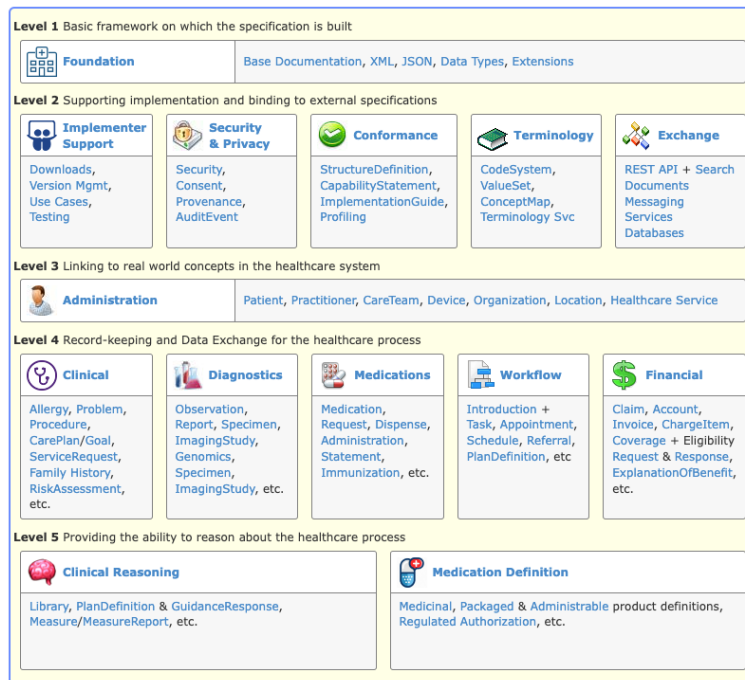


Some challenges of make/use approach

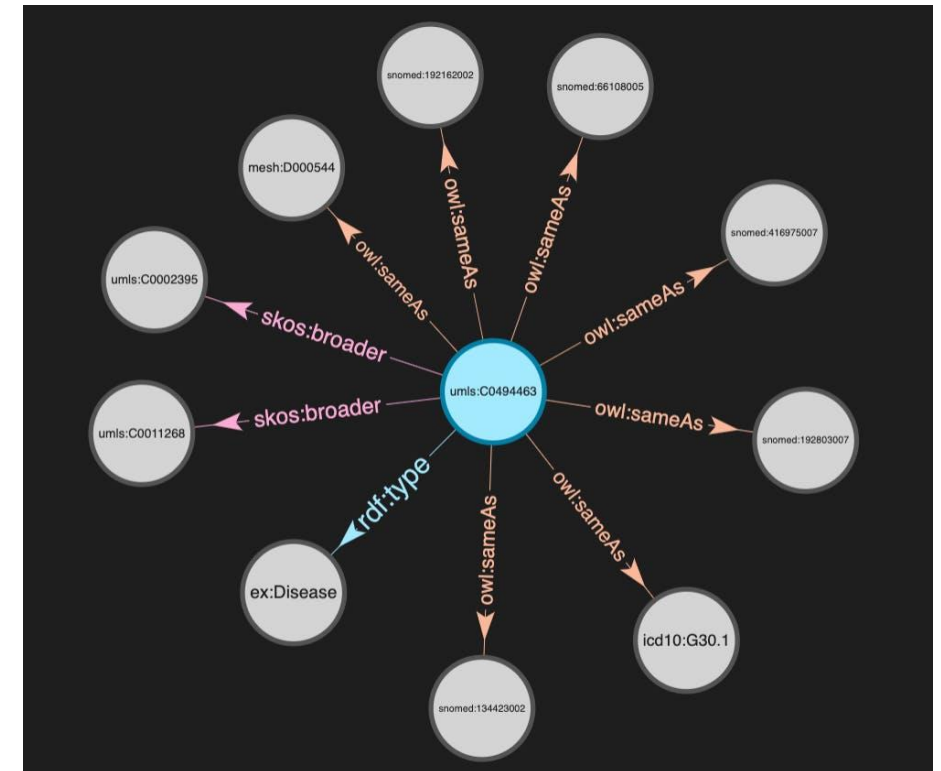
- **Data refresh:** it becomes challenging to release a new integrated KG to match a refresh of a single data source. Changes must be “big enough” to warrant a new release;
- **Data provenance:** the graph construction processes being decoupled from the graph consumption processes there is an information gap between the two;
- **Data uniformity:** performing data integration is, by nature, about fitting a source world conceptualisation into a target one. The assumption is that the integrated graph is a one-size-fits-all one.

Data integration

- We need to select a target ontology / model and an approach to deal with the many identifier schemes used across all the fields



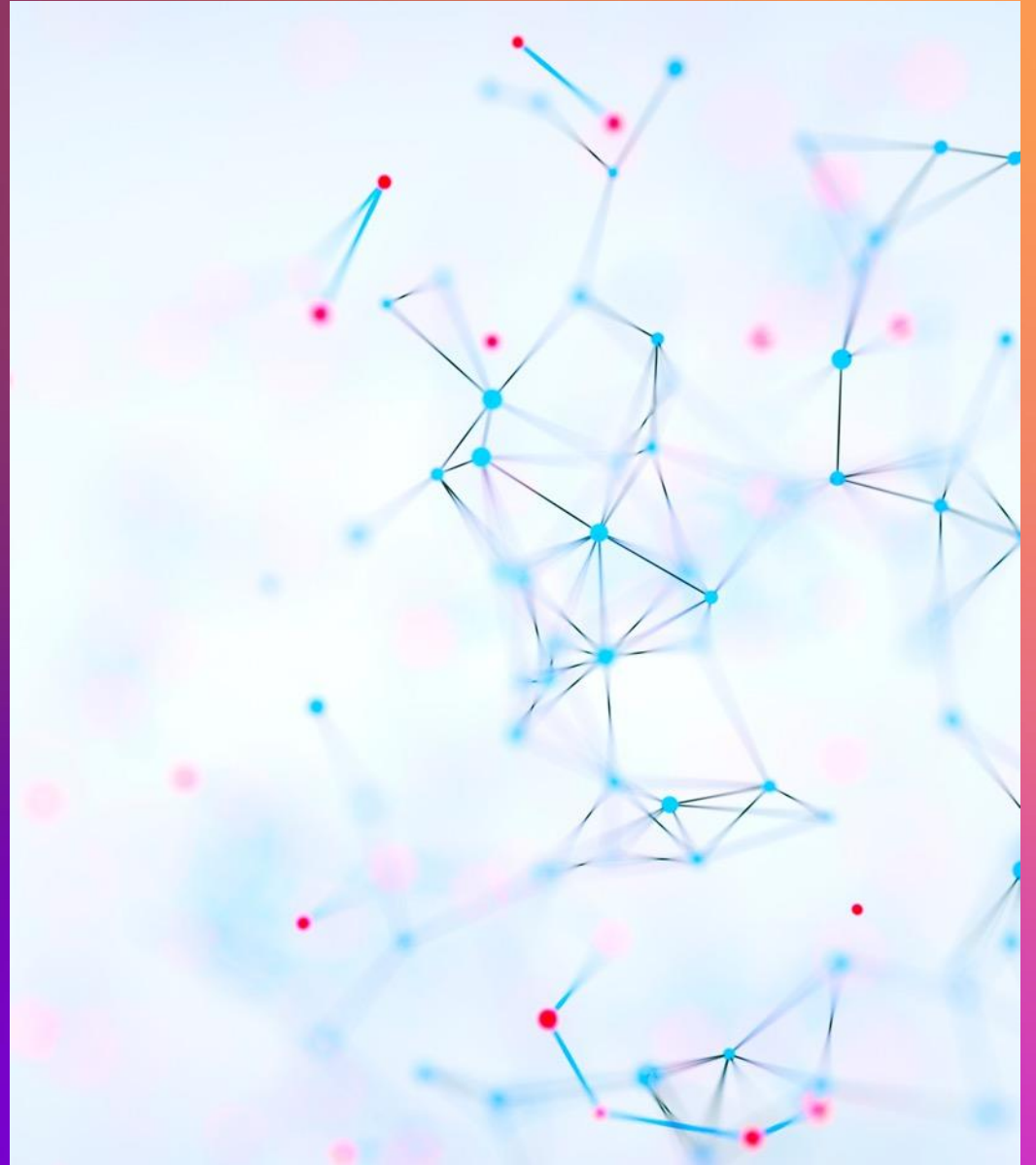
FHIR is an established data model to integrate data



`sameAs` reasoning can be applied to help with identifiers

Use

Now query the graph and do some machine learning with it



Example of questions

- The integrated KG provides exploration capabilities spanning over the whole spectrum of multi-omics data
- Users can formulate this kind of queries over the graph
 - “What is the gene encoding protein X?”
 - “What are the drugs containing a compound acting on target Y?”
- We can add IF/THEN rules to infer some statements: “IF compound X acts on gene Y which encodes protein Z, THEN compound X acts on protein Z”
- The challenge for answering all the above queries is to align the dataset semantics in a target ontology and reconcile identifiers

Interactive exploration and query

- There is no lack of options! Picking one depends on the target audience and the interaction pattern(s)

Yasgui

Query

```
1 - PREFIX rdfs: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
2 - PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
3 - SELECT * WHERE {
4   ?sub rdfs:sub ?obj .
5 } LIMIT 10
```

Yasqe

Table

Response

10 results in 0.061 seconds

Filter query results

Page size: 50

	sub	pred	obj
1	http://www.openlinksw.com/virtrdf-data-formats#default-t-id	rdfs:type	http://www.openlinksw.com/schemas/virtrdf#QuadMapFormat
2	http://www.openlinksw.com/virtrdf-data-formats#default-t-id-nulable	rdfs:type	http://www.openlinksw.com/schemas/virtrdf#QuadMapFormat
3	http://www.openlinksw.com/virtrdf-data-formats#default-t-id-nonblank	rdfs:type	http://www.openlinksw.com/schemas/virtrdf#QuadMapFormat
4	http://www.openlinksw.com/virtrdf-data-formats#default-t-id-nonblank-nulable	rdfs:type	http://www.openlinksw.com/schemas/virtrdf#QuadMapFormat
5	http://www.openlinksw.com/virtrdf-data-formats#default-1	rdfs:type	http://www.openlinksw.com/schemas/virtrdf#QuadMapFormat
6	http://www.openlinksw.com/virtrdf-data-formats#default-t-nulable	rdfs:type	http://www.openlinksw.com/schemas/virtrdf#QuadMapFormat

Yasr

[Yasgui API Reference - Docs - Triply](#)



[Sparnatural - Javascript SPARQL query builder](#)



[15 Best Graph Visualization Tools for Your Neo4j Graph Database](#)

Entry Point APIs

Open PHACTS
Open Pharmacological Space

Map free text to a concept URL

/search/freetext GET

Chemical Structure Exact Search

/structure/exact GET

InchiKey to URL

/structure GET

Inchi to URL

/structure GET

Chemical Structure Similarity Search

/structure/similarity GET

SMILES to URL

/structure GET

Chemical Structure Substructure Search

/structure/substructure GET

< 18 of 47 >

Open PHACTS API Walkthrough

Jul. 31, 2013 • 4 likes • 2,904 views

Download Now

Download to read offline

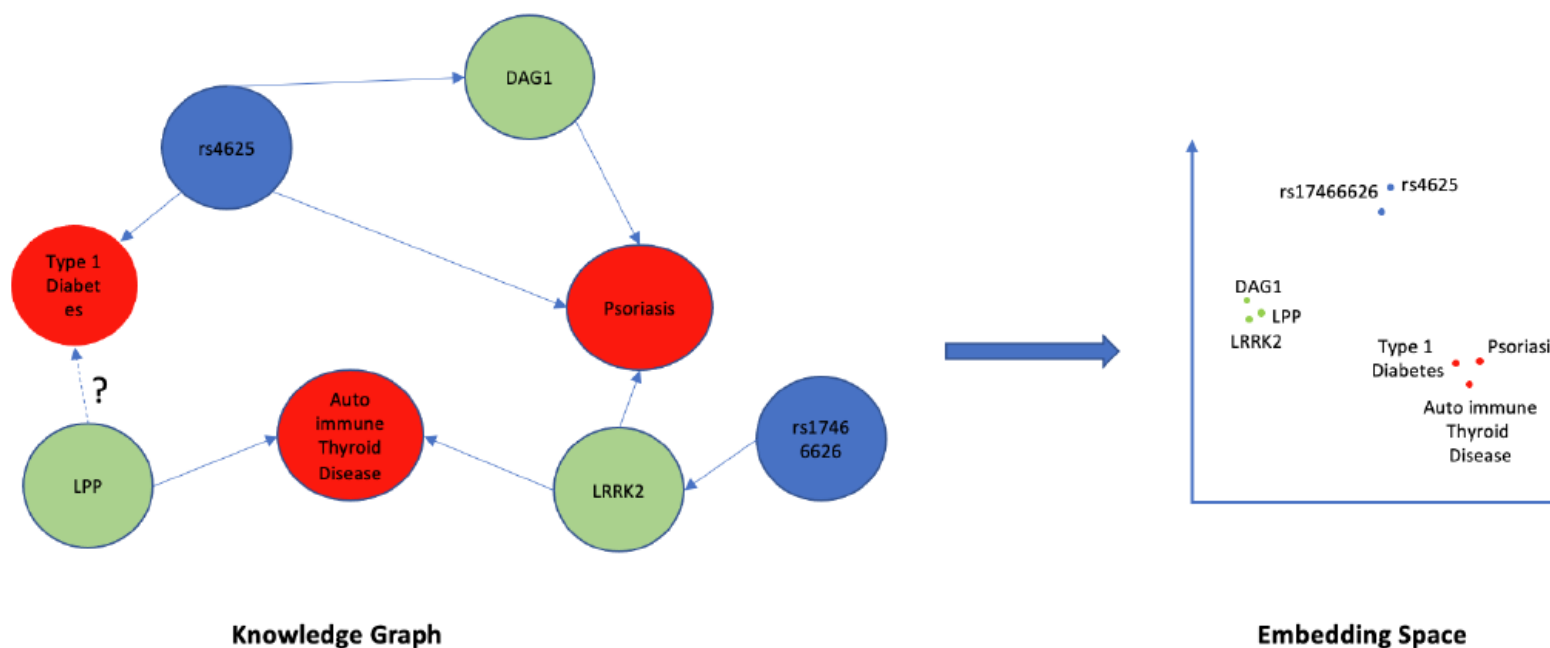
Paul Groth

[Open PHACTS API](#)



Graph embeddings

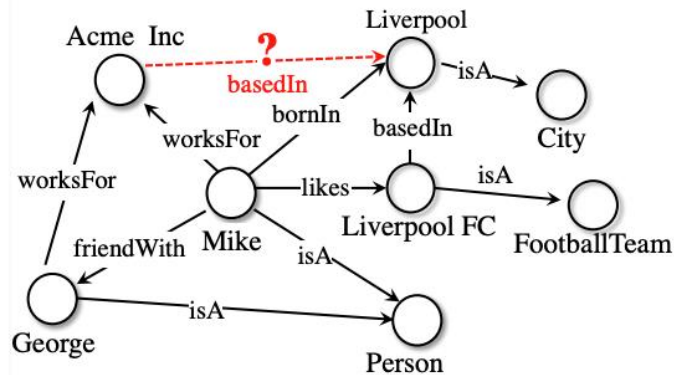
- The AI needs to learn the graph in order to reason over it. We do that by mapping the content of the graph into a vector space



Then, three possible type of investigations

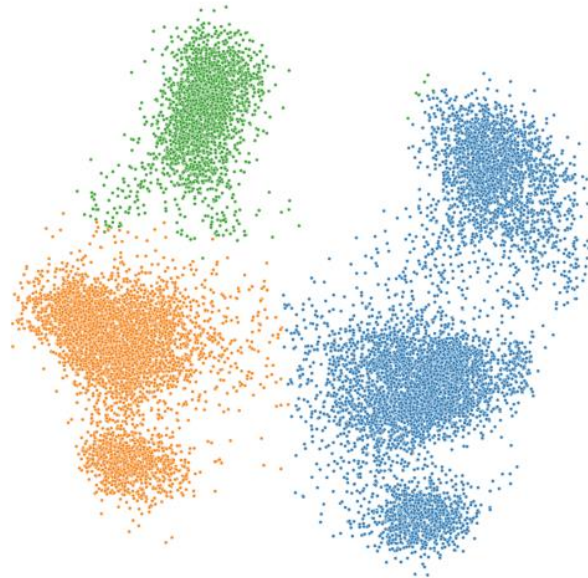
LINK PREDICTION / TRIPLE CLASSIFICATION

- Knowledge graph completion
- Content recommendation
- Question answering



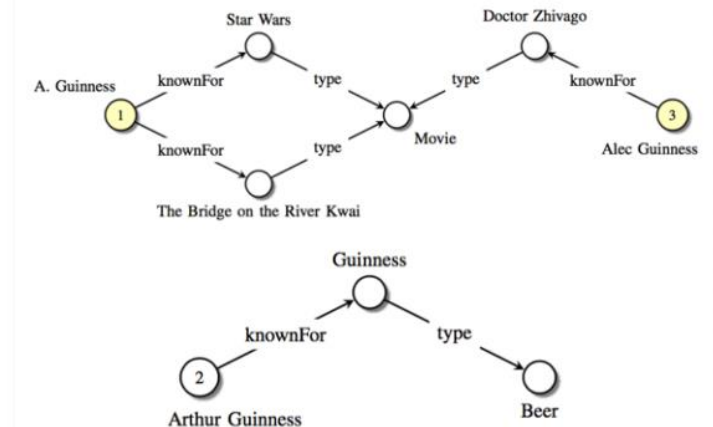
COLLECTIVE NODE CLASSIFICATION / LINK-BASED CLUSTERING

- Customer segmentation



ENTITY MATCHING

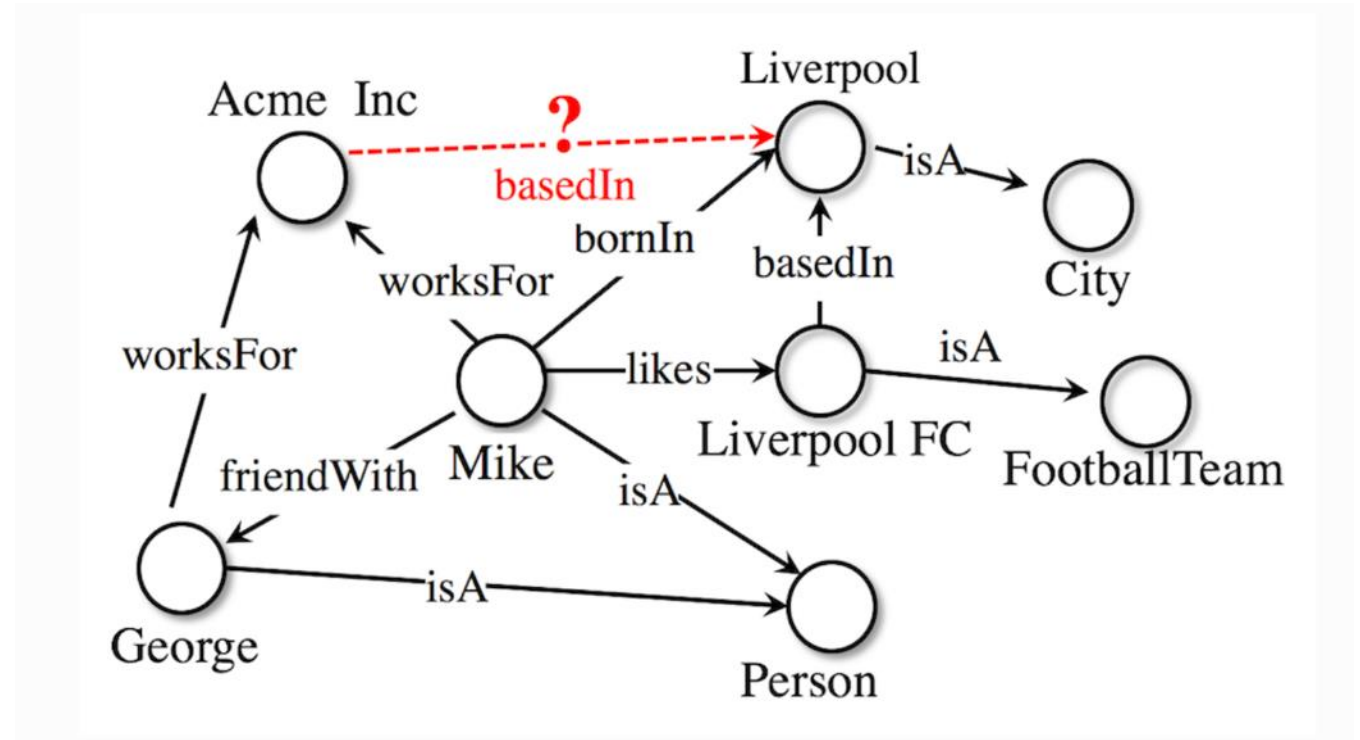
- Duplicate detection
- Inventory items deduplication



Pic from [Nickel et al. 2016a]

Predicting links

- The goal is to predict what is the likelihood of a link not present in the graph
- The outcome depends on all the content in the graph, at any distance from the target nodes
- We have a sub-graph explanation sub system able to state which nodes were most influential in the scoring



The task here is to predict if the link in red could be a statement in the graph

Training with Uncertainties/Importance of links

- As an extension to state of the art graph machine learning we developed an approach to weigh the links based on importance

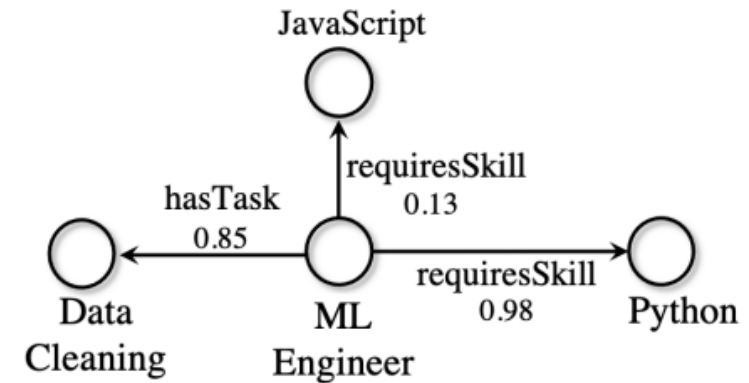
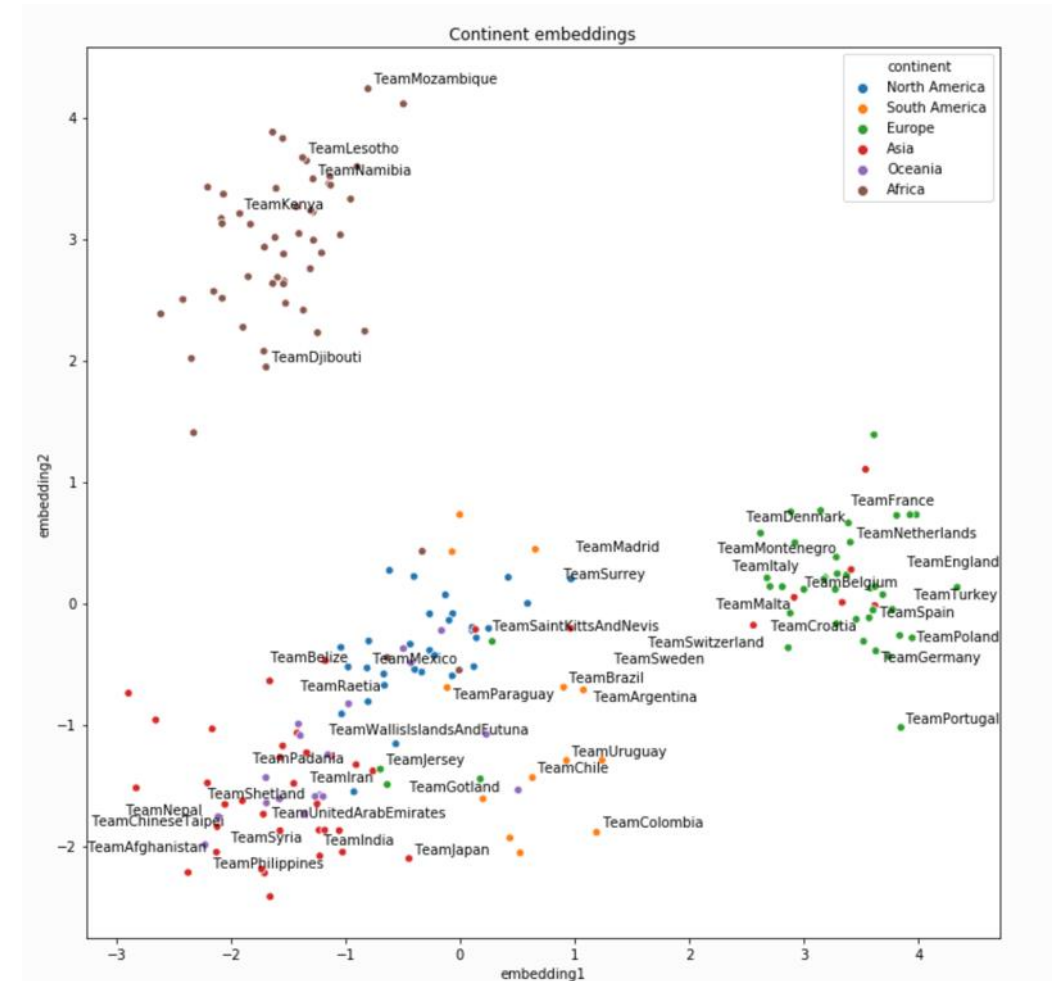


Figure 1: A Knowledge graph with numeric attributes associated to triples.

Clustering

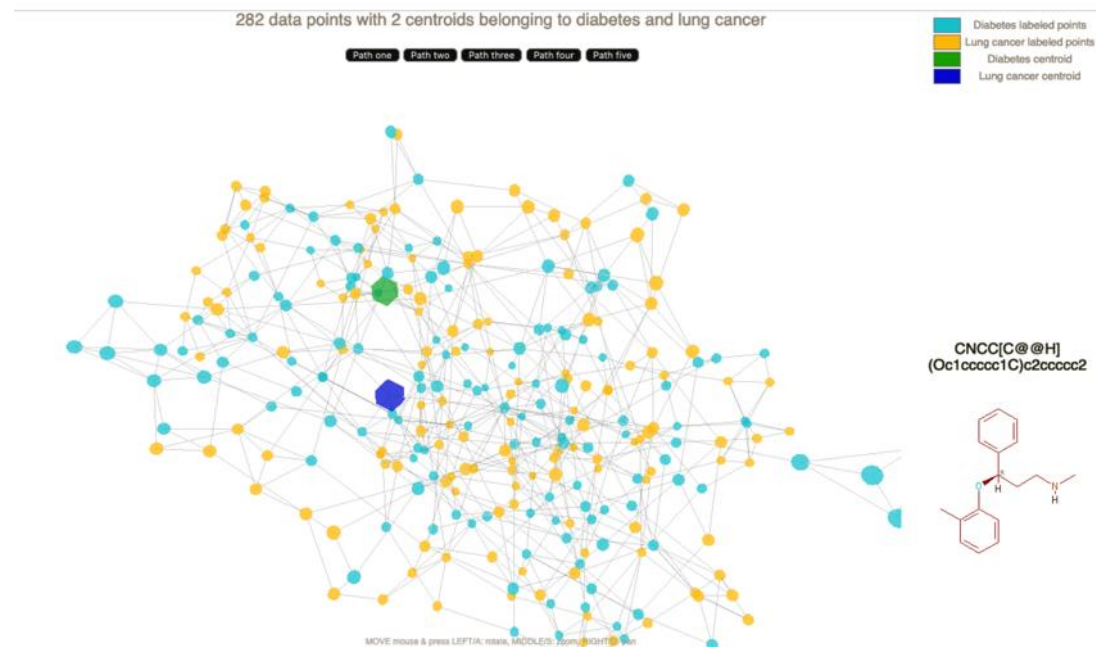
- If we cluster entities based on their vector representation, we see emerging features
- In this example, the continents are not in the graph but emerge from other type of edges



Clustering the nodes based on continents (image from [Examples — AmpliGraph 1.4.0 documentation](#))

Using embedding space for generation

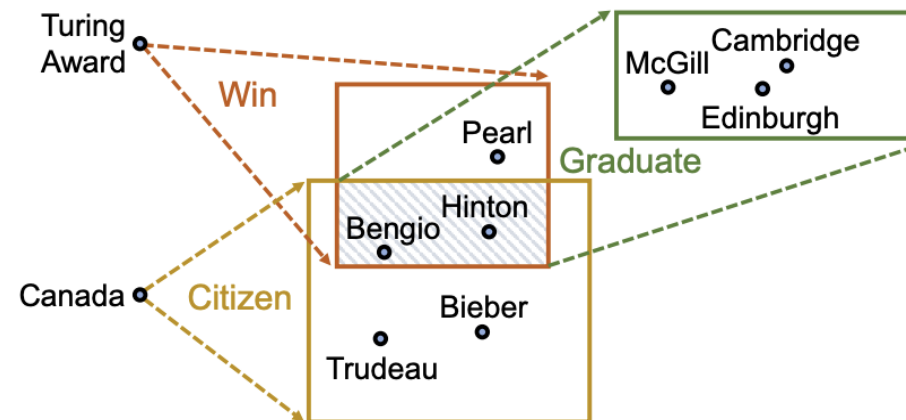
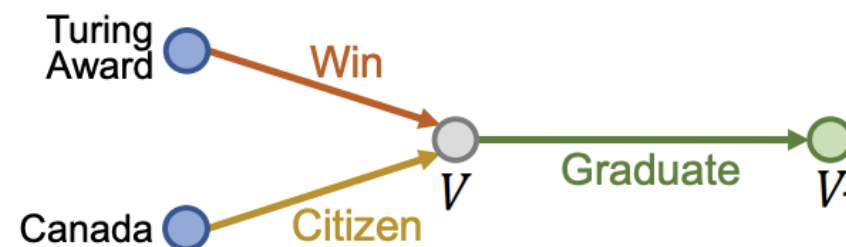
- In ChemoVerse we learn to represent and generate structures of small molecules from an embedding space (or latent space).
- We use a manifold traversal with heuristic search to explore a latent space created from knowledge on this chemical space.
- Different heuristics and scores such as the Tanimoto coefficient, synthetic accessibility, binding activity, and QED drug-likeness can be incorporated to increase the validity and proximity for desired molecular properties of the generated molecules.
- With this novel traversal method, we are able to find more unseen compounds and more specific regions to mine in the latent space.



Searching for complex combinations

- For future work, we are considering doing query answering using the embedding space. Eventually using an approach like Query2Box

$$q = V_? . \exists V : \text{Win}(\text{TuringAward}, V) \wedge \text{Citizen}(\text{Canada}, V) \wedge \text{Graduate}(V, V_?)$$

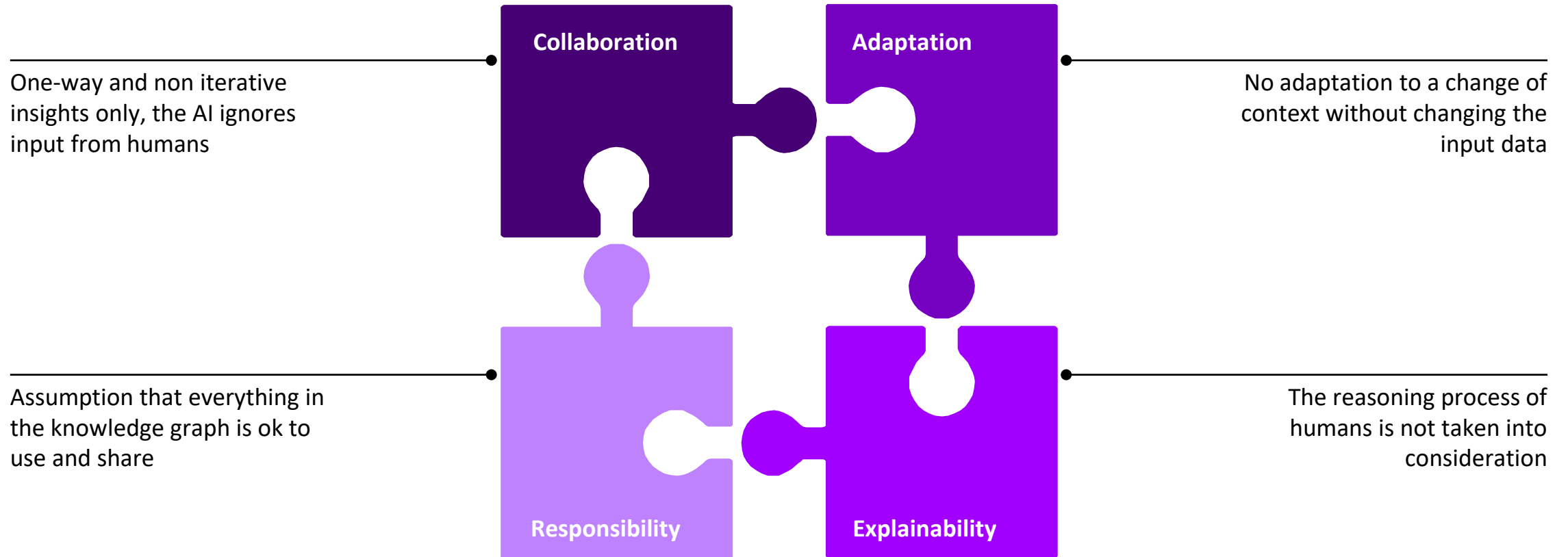


What can we do better?

Exploring what could be in the graph



Is this ok for Hybrid Intelligence?



Grid derived from: Akata, Z., Balliet, D., de Rijke, M., Dignum, F., Dignum, V., Eiben, G., Fokkens, A., Grossi, D., Hindriks, K., Hoos, H., Hung, H., Jonker, C., Monz, C., Neerincx, M., Oliehoek, F., Prakken, H., Schlobach, S., van der Gaag, L., van Harmelen, F., ... Welling, M. (2020). A Research Agenda for Hybrid Intelligence: Augmenting Human Intellect With Collaborative, Adaptive, Responsible, and Explainable Artificial Intelligence. *Computer*, 53(8), 18–28.

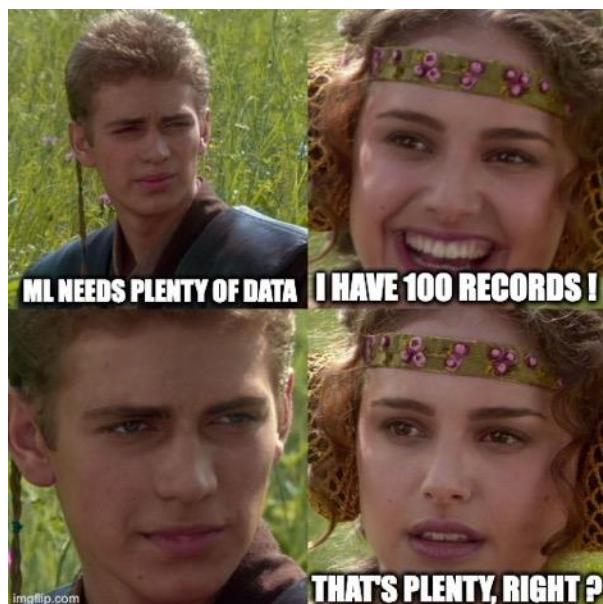
<https://doi.org/10.1109/MC.2020.2996587>

Opportunities

3 keys things we can leverage to improve on our pipeline

Reasoning rules

Lack of data compensated by expert knowledge



Rise of data fabrics

KG constructed on demand rather than pre-assembled



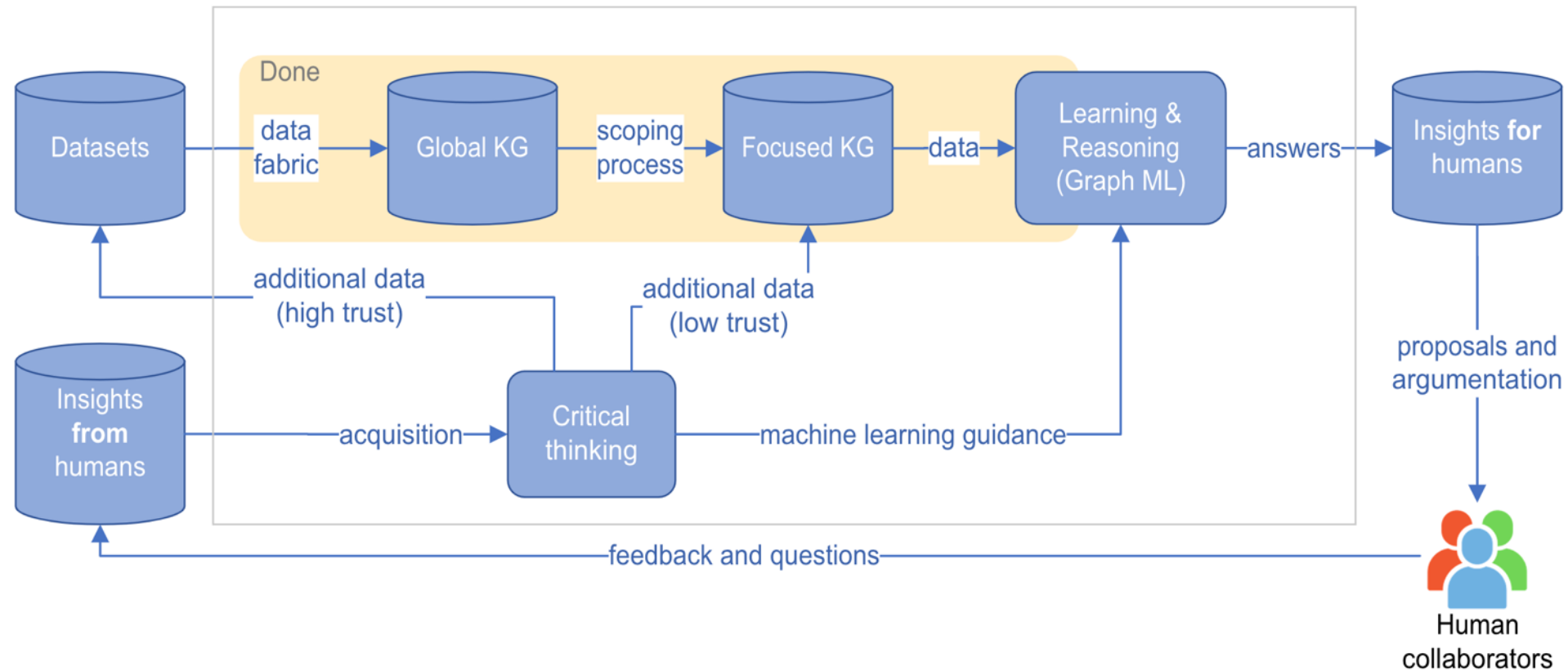
Shared understanding

Graphs used as is by all stakeholders, human and AI



Our proposed approach

- Add scoping and critical thinking elements. The goal is to incorporate information from the interaction with the users down to the KG construction



Take-away

Knowledge Graphs can be a key back-end component when introducing AI collaborators in a team

They enable:

- Having all stake-holders use the same conceptual model
- Reason and discuss over this model
- Put up a bidirectional data-to-insight pipeline

To chat more please reach out at:
christophe.gueret@accenture.com

