

The Symbiotic Relationship between Data Spaces and AI. Challenges and Trends

Prof. Edward Curry
University of Galway

Insight SFI research Centre for Data Analytics

Insight



SFI RESEARCH CENTRE FOR DATA ANALYTICS

HOST INSTITUTIONS



PARTNER INSTITUTIONS

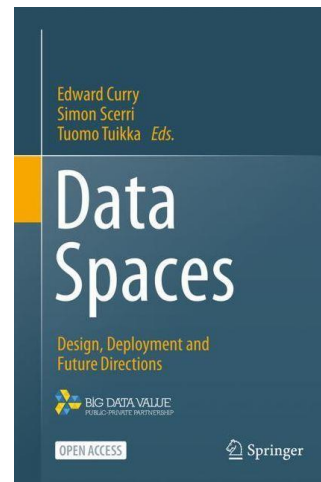
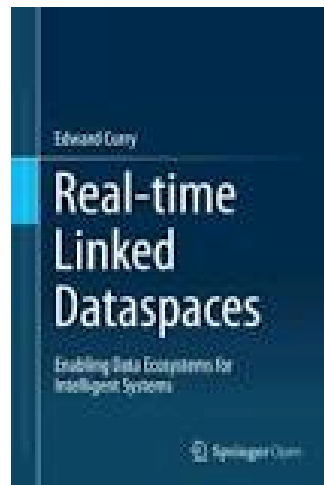
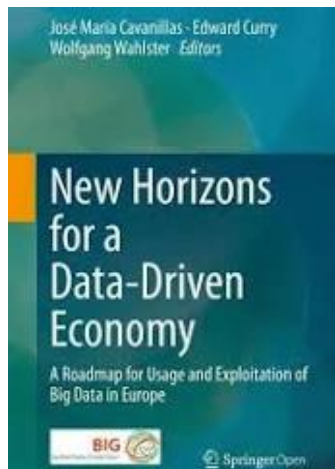


FUNDED BY:



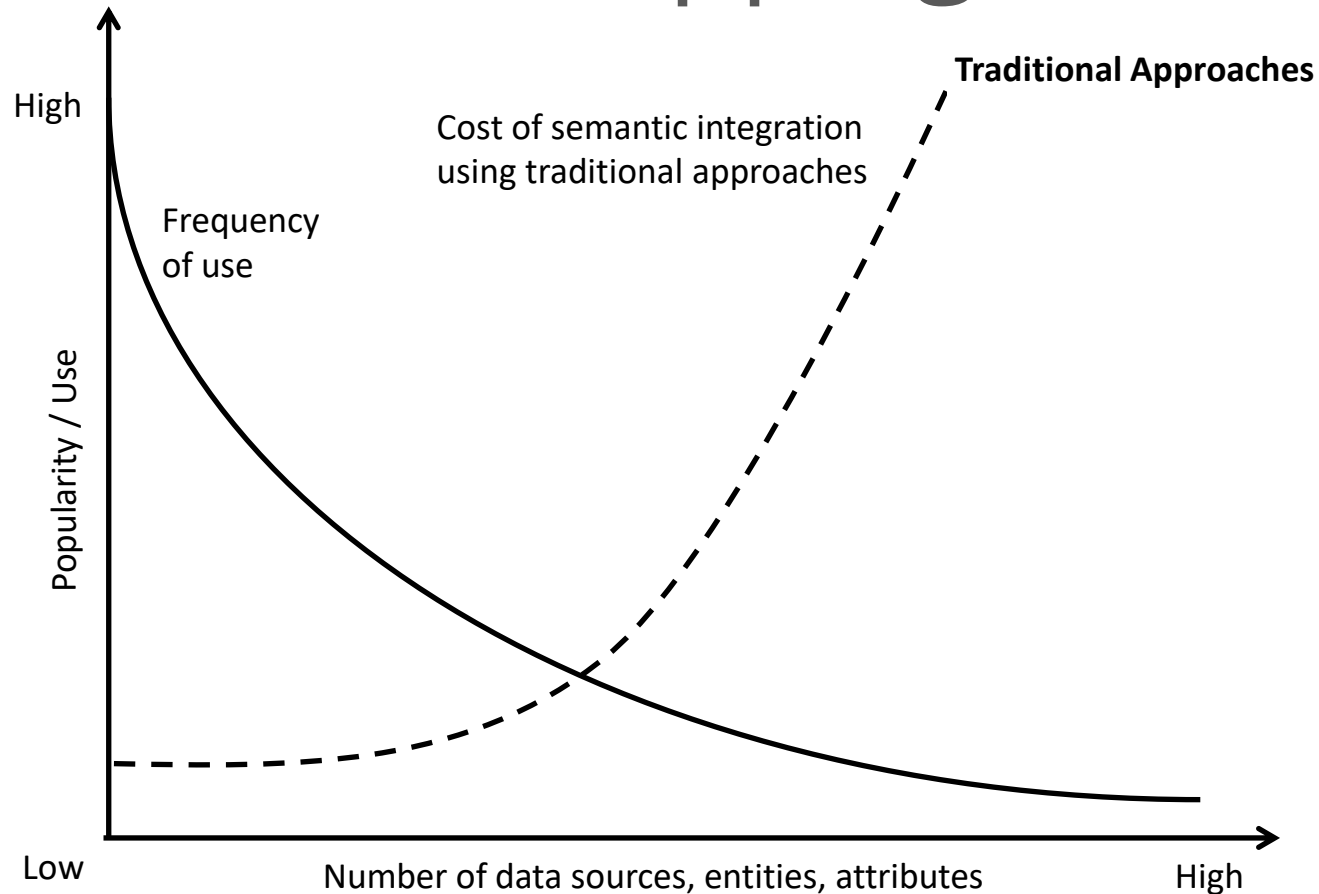
Edward Curry

I have been researching the underlying technology for data spaces for the last decade...



Data Spaces need a data co-existence
approach for “Good Enough”
Interoperability

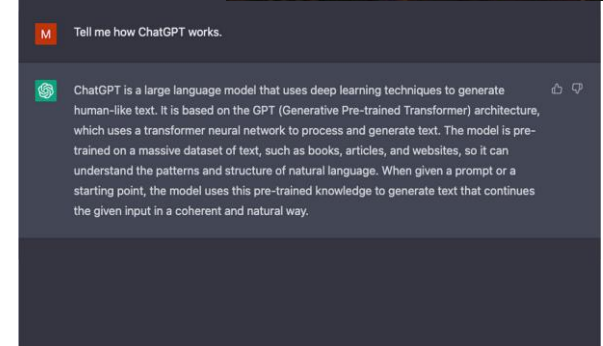
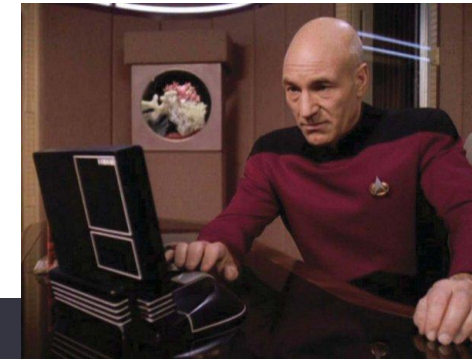
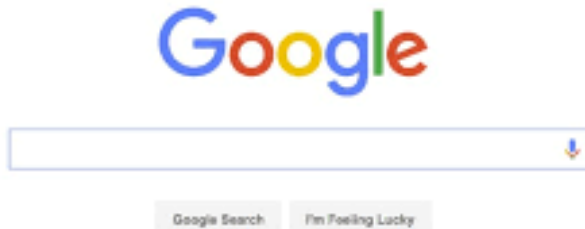
Long Tail Semantics continues to grow... and the cost of mapping increases....



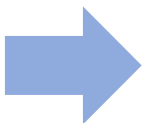
The Long Tail of Data

<http://dataspaces.info>

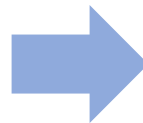
A new paradigm for human-data Interaction.....



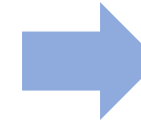
From Structure



to Search



to Knowledge Graph



to Conversations ?

~1995
~100K Websites
Exact Results
Human Curated

~1998
~2.4M Websites
Approximate Results
Computed

~2012
~700M
Approximate Results + Exact
Computed + Crowd

~2023
~2B
Approximate Results ?
Content creation?

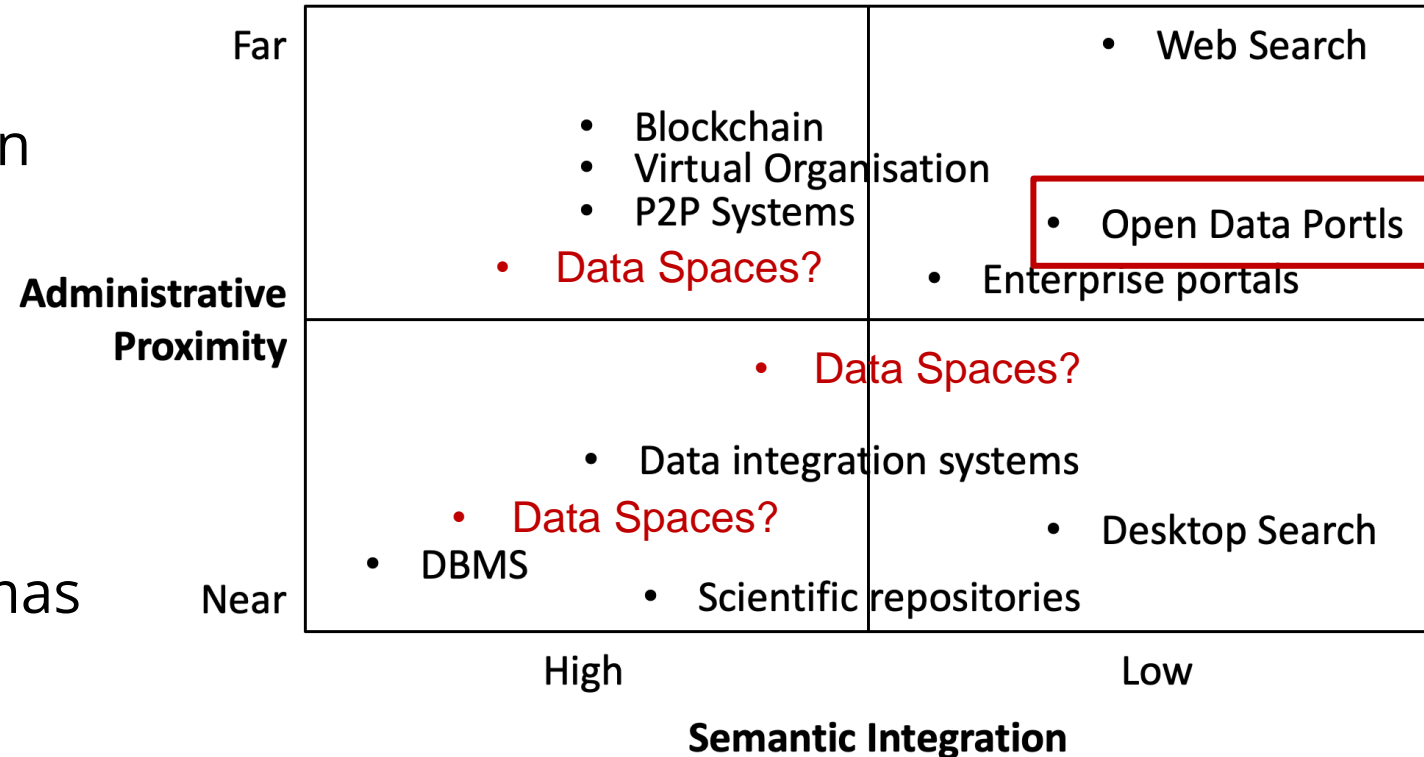
Control and Coordination

Administrative Proximity

- Close vs. Loose Coordination
- Assumptions concerning guarantees such as data, access, quality, and consistency,

Semantic Integration

- Degree to which data schemas are matched up (types, attributes, and names).



Halevy, A., Franklin, M. and Maier, D. 2006. Principles of dataspace systems. *25th ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems - PODS '06* (New York, New York, USA, 2006), 1–9.

In today's humongous database systems, clarity may be relaxed, but business needs can still be met.

BY PAT HELLAND

If You Have Too Much Data, then 'Good Enough' Is Good Enough

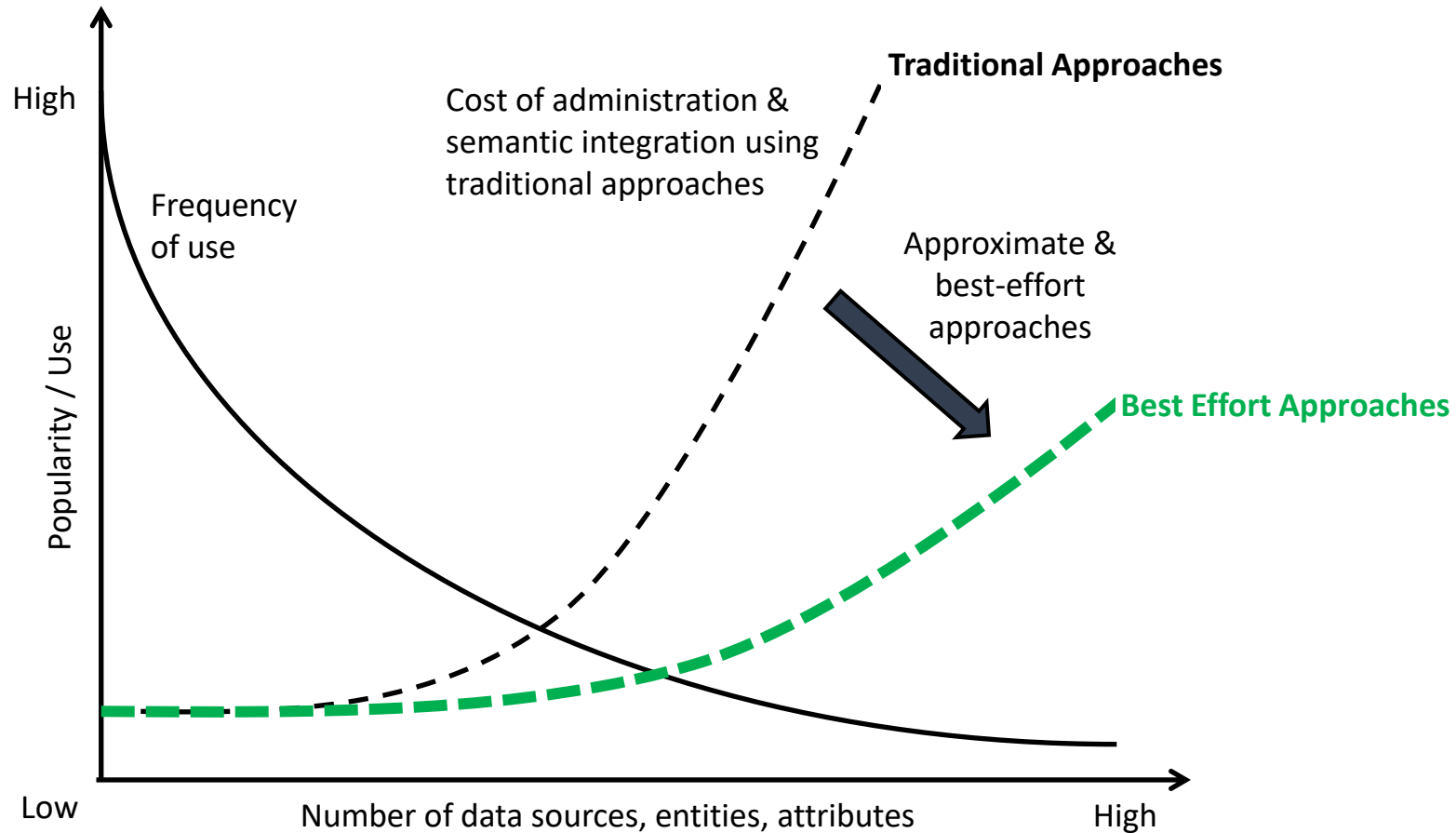
"We can no longer pretend to live in a clean world. SQL and its Data Definition Language (DDL) assume a crisp and clear definition of the data, but that is a subset of the business examples we see in the world around us. It's OK if we have lossy answers—that's frequently what business needs."

What is a Dataspace? (2006)

“Dataspaces are not a data integration approach; rather, they are more of a **data co-existence approach**. The goal of dataspace support is to provide base functionality over all data sources, regardless of how integrated they are.” (Halevy, A., Franklin, M. and Maier, D. 2006.)

Incrementalism, Approximate, Interactive

Approximate and Best Effort Approaches



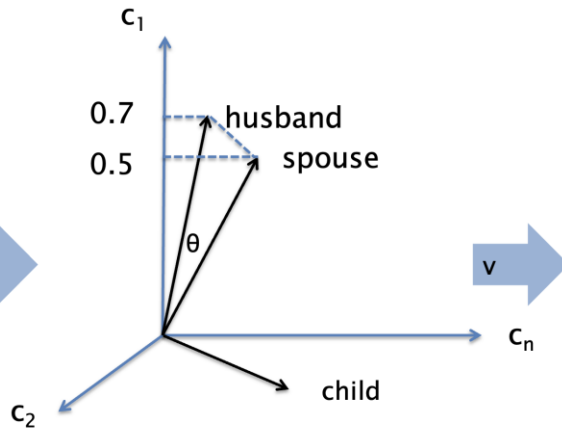
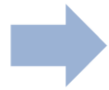
The Long Tail of Data

<http://dataspaces.info>

Creating Approximate Services is a Key Challenge

Investigate techniques to enable approximate and best-effort support services for loose administrative proximity and semantic integration

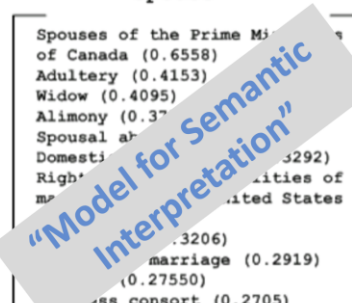
Distributional semantic model:
 Semantic statistical knowledge extracted from large Web corpora



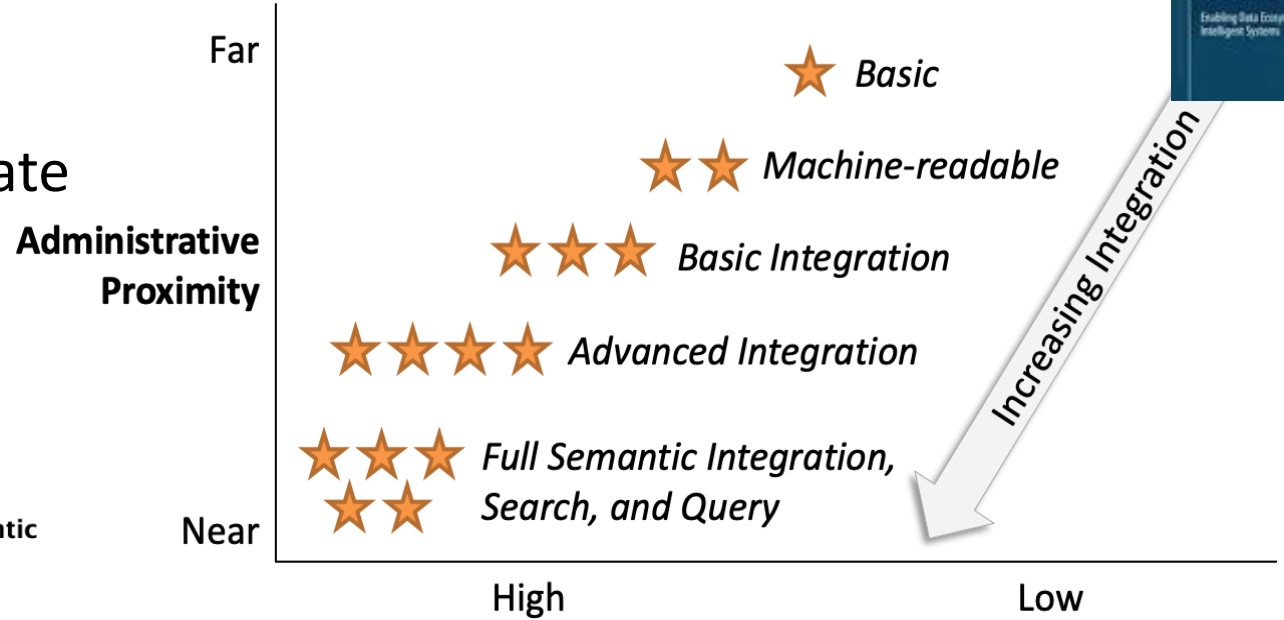
Works as a semantic ranking function

spouse

```
Spouses of the Prime M...
of Canada (0.6558)
Adultery (0.4153)
Widow (0.4095)
Alimony (0.37...
Spousal ab... (0.3292)
Right... ities of
... ited States
... (0.3206)
... marriage (0.2919)
... (0.27550)
... consort (0.2705)
D...orce (0.2383)
...
```



E.g. $dsm(spouse, widow) = 0.4095$
 E.g. $dsm(spouse, car) = 0.009$



Incremental support services

- Catalog
- Entity management
- Query and search
- Data discovery
- Human tasks
- Quality of service
- Complex event processing
- Streams dissemination
- Approximate semantic event matching

LLMs for "Good Enough" Interoperability

User: As a Semantic Web expert, please extend the SHACL shapes given below such that I can also describe digital versions of paintings with them, using properties from the Integrated Authority File, the GND. For the painter, please use a 'gndo:firstArtist' that is a 'gndo:DifferentiatedPerson' which has its GND id as a property. For the year of creation, please use 'gndo:dateOfProduction'.

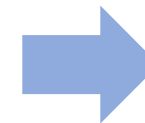
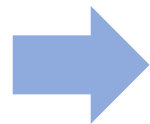
```
```turtle ...
```

```

1 @prefix sh: <http://www.w3.org/ns/shacl#> .
2 @prefix dcat: <http://www.w3.org/ns/dcat#> .
3 @prefix dct: <http://purl.org/dc/terms/> .
4 @prefix odrl: <http://www.w3.org/ns/odrl/2/> .
5
6 [] a sh:NodeShape ;
7 sh:targetClass dcat:Dataset ;
8 sh:property [
9 sh:path dct:title ;
10 sh:maxCount 1 ; ↔ sh:minCount 1 ;
11 sh:message "A dcat:Dataset must have exactly one
12 ↳ dct:title." ;
13] ;
14 sh:property [
15 sh:path odrl:hasPolicy ;
16 sh:maxCount 1 ; ↔ sh:minCount 1 ;
17 sh:message "A dcat:Dataset must have exactly one
18 ↳ odrl:hasPolicy." ;
19] .

```

Listing 1: The base SHACL shapes.



```

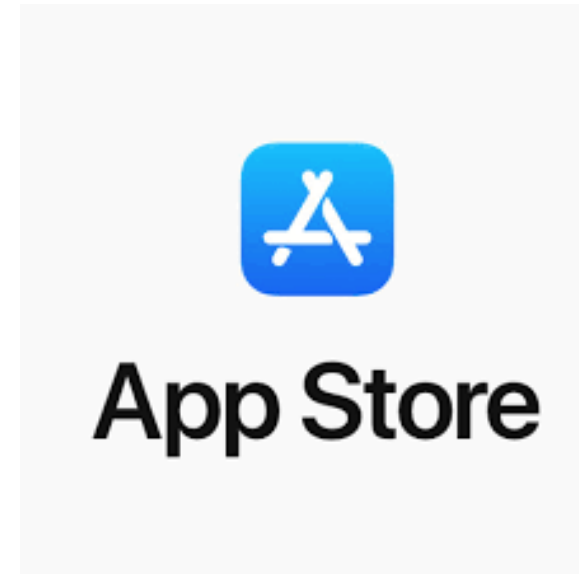
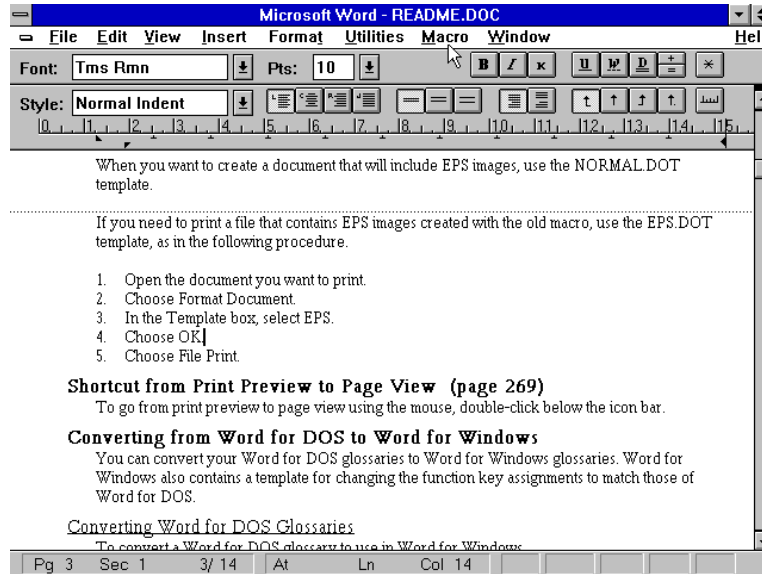
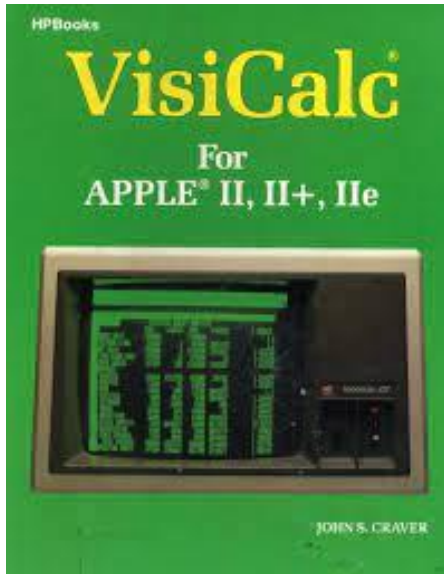
1 @prefix sh: <http://www.w3.org/ns/shacl#> .
2 @prefix dcat: <http://www.w3.org/ns/dcat#> .
3 @prefix dct: <http://purl.org/dc/terms/> .
4 @prefix odrl: <http://www.w3.org/ns/odrl/2/> .
5 @prefix gndo:
6 ↳ <http://d-nb.info/standards/elementset/gnd#> .
7
8 [] a sh:NodeShape ;
9 sh:targetClass dcat:Dataset ;
10
11 ... omitted for conciseness ...
12
13 sh:property [
14 sh:path gndo:firstArtist ;
15 sh:class gndo:DifferentiatedPerson ;
16 sh:maxCount 1 ; ↔ sh:minCount 1 ;
17 sh:message "A digital representation of a painting
18 ↳ must have exactly one gndo:firstArtist." ;
19 sh:node [
20 a sh:PropertyShape ;
21 sh:path gndo:preferredNameForThePerson ;
22 sh:maxCount 1 ; ↔ sh:minCount 1 ;
23 sh:message "Each gndo:DifferentiatedPerson must
24 ↳ have exactly one
25 ↳ gndo:preferredNameForThePerson." ;
26] ;
27 sh:node [
28 a sh:PropertyShape ;
29 sh:path gndo:gndIdentifier ;
30 sh:maxCount 1 ; ↔ sh:minCount 1 ;
31 sh:message "Each gndo:DifferentiatedPerson must
32 ↳ have exactly one gndo:gndIdentifier." ;
33] ;
34 sh:property [
35 sh:path gndo:dateOfProduction ;
36 sh:maxCount 1 ; ↔ sh:minCount 1 ;
37 sh:datatype xsd:date ;
38 sh:message "A digital representation of a painting
39 ↳ must have exactly one gndo:dateOfProduction." ;
40] .

```

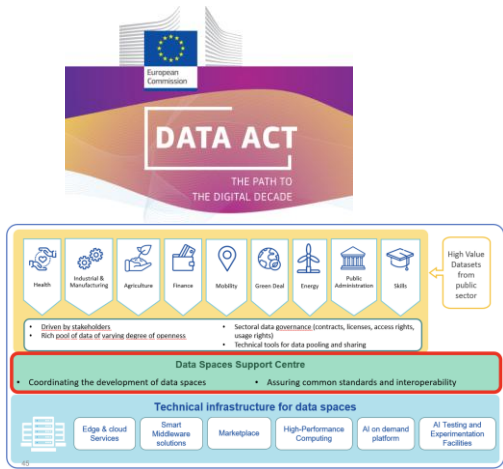
Listing 2: The SHACL shapes with extensions related to paintings, by GPT-4.

Data Spaces will require a unified Data and AI Lifecycle if we are to maximise the potential of Generative AI

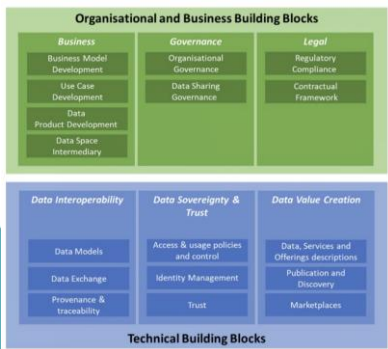
# Generative AI and Foundation models will be the Killer App for Data Spaces....



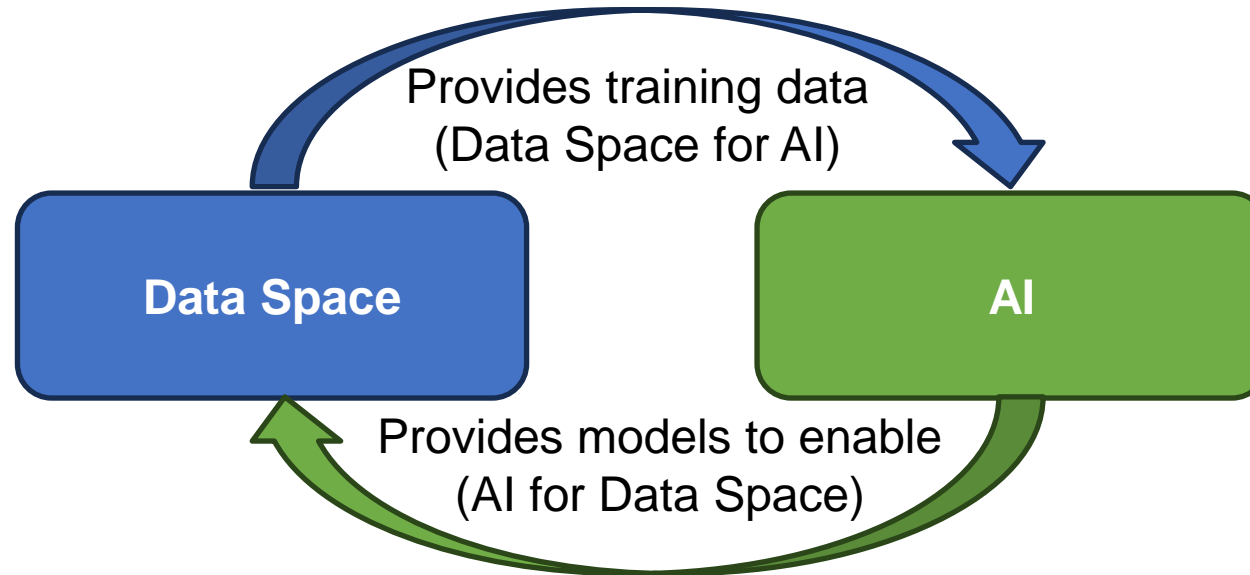
# Symbiotic Relationship between Data Spaces and AI....



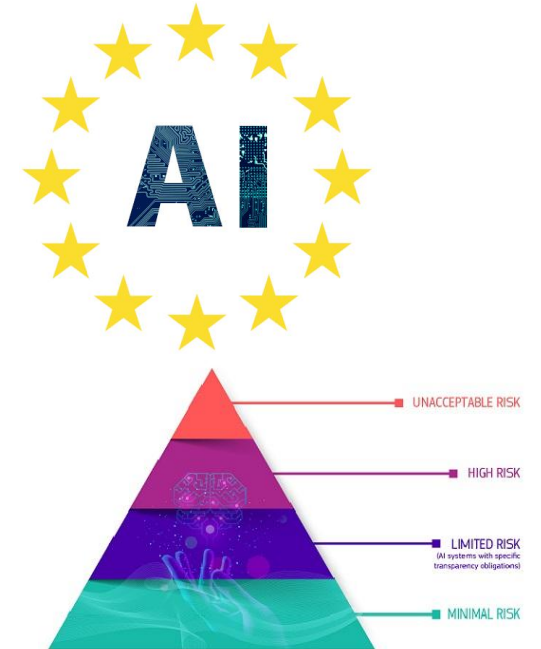
**DATA SPACES SUPPORT CENTRE**



High Quality Training Data  
Community of Data Users/Owners  
Data Governance/compliance



Tools for Diverse Data Space Tasks  
Faster Value Extraction from Multimodal Data  
Knowledge Foundation for Long-Tail Semantics

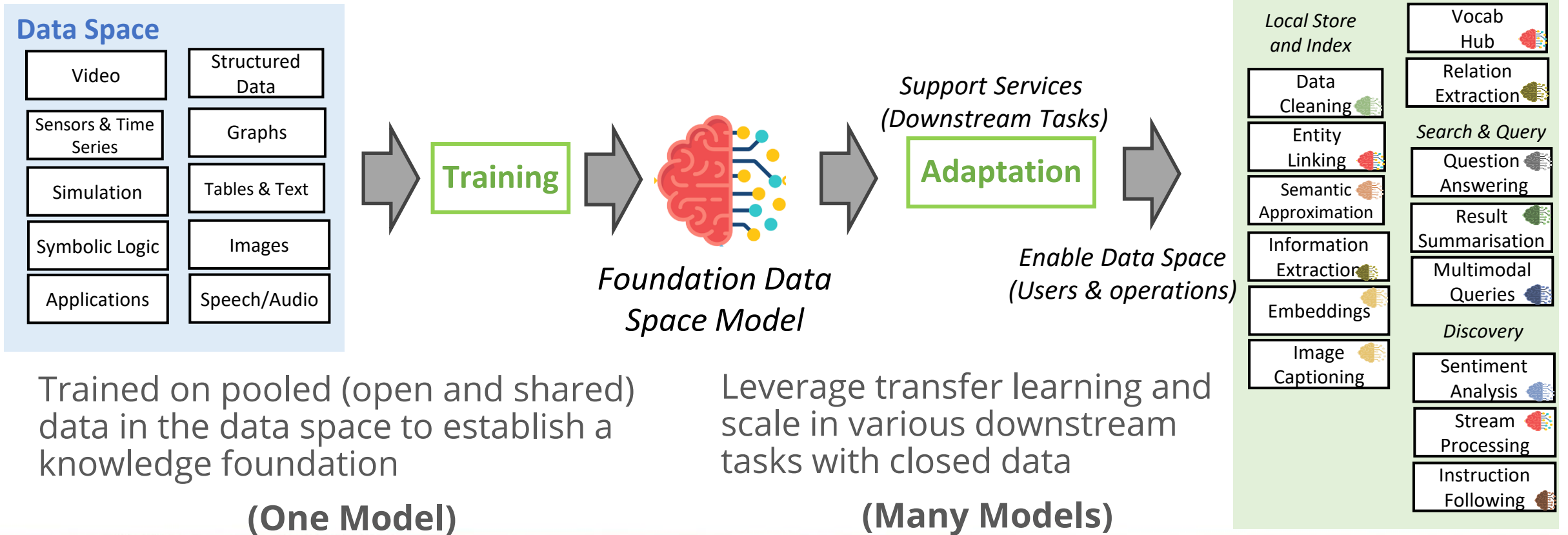


**Adra**  
The AI Data Robotics Association

E. Curry, M. Timilsina, T. Zaarour, M. Al-QATF, R. Haque, "Foundational Data Space Models: Bridging the AI and Data Ecosystems (Vision Paper), Proceedings of the 2023 IEEE International Conference on Big Data (Big Data), Sorrento, Italy, 2023



# Data Spaces for GenAI: Centralise information from the data space, then adapt to a wide range of downstream tasks....



Trained on pooled (open and shared) data in the data space to establish a knowledge foundation

**(One Model)**

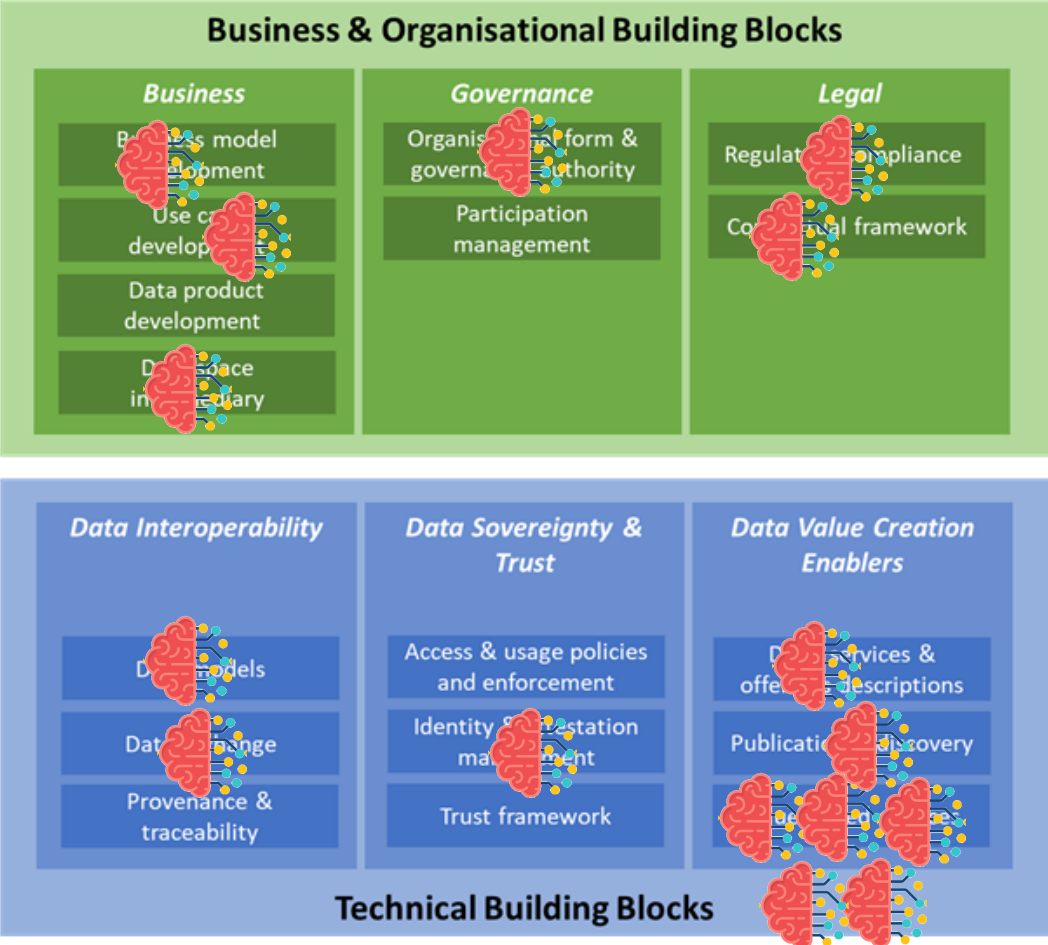
Leverage transfer learning and scale in various downstream tasks with closed data

**(Many Models)**

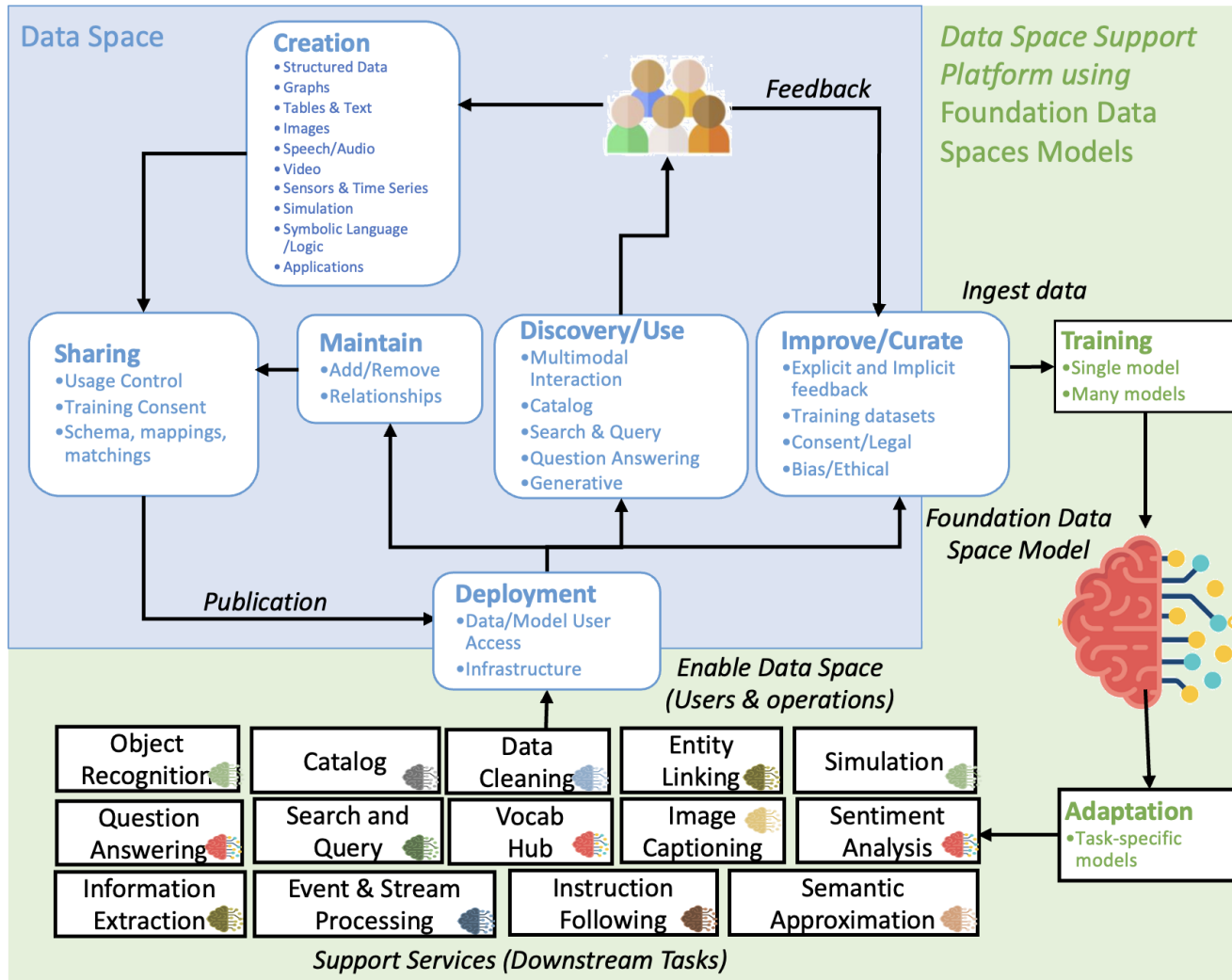
# GenAI for Data Spaces: *Supports tasks across the life cycle of the data space...*



**DATA SPACES  
SUPPORT CENTRE**



# We need Unified Data and AI Lifecycles



A unified lifecycle for data and AI models recognizes the **symbiotic relationship between both ecosystems** and can serve as the basis to simplify the development, operation, and use of data-intensive AI systems.