



# The Great Escape: Liberating 20+ Years of Legacy Enterprise Data

Demo Day  
July 19, 2023



# About Today's Speaker



## **Eliud Polanco**

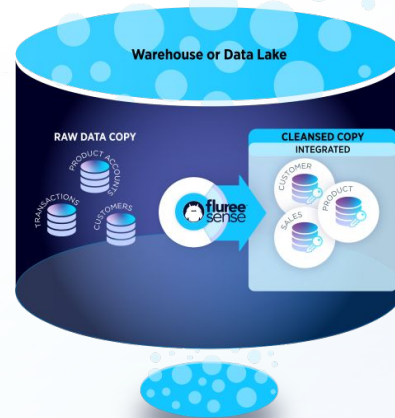
- Former Head of Analytics at HSBC, DeutscheBank, ScotiaBank
- Former Head of Data Innovation @ Citi
- Current President at Fluree PBC, a start-up focused on secure data-sharing and collaboration at scale using Zero-Trust and Semantics

# What Fluree Offers



## Intelligent Knowledge Graph

**Fluree** combines cryptography and semantic data standards to enable secure data collaboration



## AI/ML Data Cleansing Pipeline

With **Fluree** it is faster and easier to clean raw data and make it consumable for users and algorithms via semantic Knowledge Graph

# Agenda

- I. **The Knowledge Graph opportunity.** Why Knowledge Graphs are emerging as more critical for Enterprise use cases.
- II. **The Legacy Data Problem.** Why Knowledge Graphs have historically been difficult to adopt at scale.
- III. **The Innovations.** How JSON-LD and AI/ML are helping bridge the gap between relational and semantic worlds.
- IV. **The Demo.** An How Fluree converts scale legacy data into data ready for semantic Knowledge Graphs.
- V. **Getting Started With Knowledge Graphs In Your Organization.** A BluePrint for socializing KG use cases in your company.



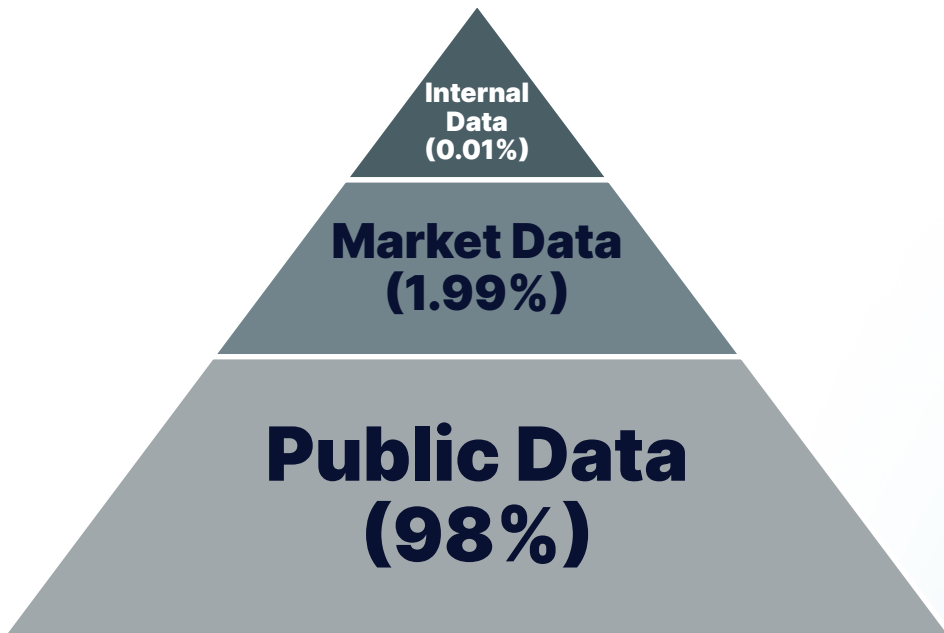
# I. The Knowledge Graph Opportunity

# Why Look At Knowledge Graphs?

- **Data sharing** is becoming an imperative for competitive differentiation
  - Inside a company's four walls
  - Between companies and regulatory agencies who have adopted semantic standards for how they want to receive data
  - Between companies and suppliers, partners, even competitors using common industry ontologies
- **AI, Algos and Machines** are becoming the **primary consumer** of vast stores of data inside a company
  1. Generative AI and **Large Language Models** (LLM) shook the world
  2. **Internet of Things**: devices create data, and devices consume data

**Semantic interoperability is no longer a niche use case... it is an economic imperative**

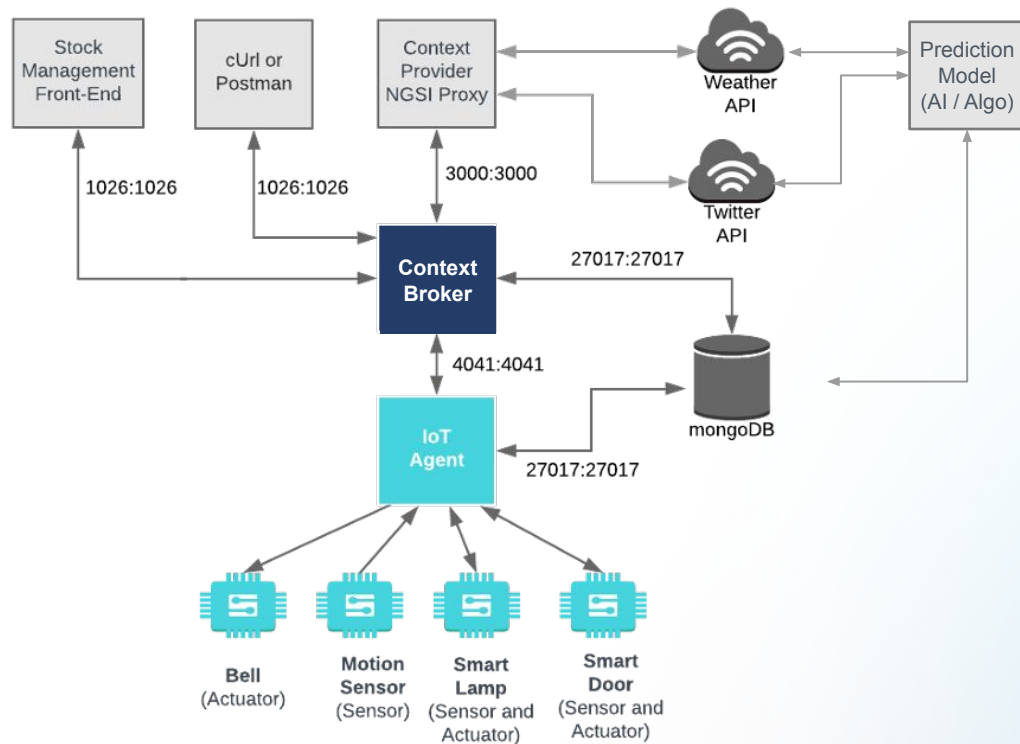
# 1. Large Language Models (LLMs)



- The vast quantity of data consumed into LLMs are public data that is either unstructured, or using semantic Web standards (RDF)
- Some market data (data available for a price) is opportunistically being integrated into LLMs (at a price)
- Most internal data is out of the reach, because:
  - It is proprietary.
  - It lacks semantic context.
- Large enterprises today are experimenting in how to **build their own internal LLMs** using a mix of their own structured and unstructured corporate internal data.
- This will drive **competitive differentiation inside companies** in how they create products and optimize their own business processes



# Internet of Things (IoT)



- Data being created by sensors and actuators...
- ...and being distributed via APIs...
- ...to be consumed by algorithmic models (e.g., predictive maintenance, or personalization recommendation)...
- ...with little man-in-the-middle intervention.
- Many IoT devices use JavaScript Object Notation (JSON) as the defacto lightweight data interchange format. It is easy for machines to parse and generate.
- This means, that over time, a vast amount of Operational Technology (OT) data will be created in JSON.
- There is a desire to link Information Technology (IT) or relational world data to OT data to create even better models.



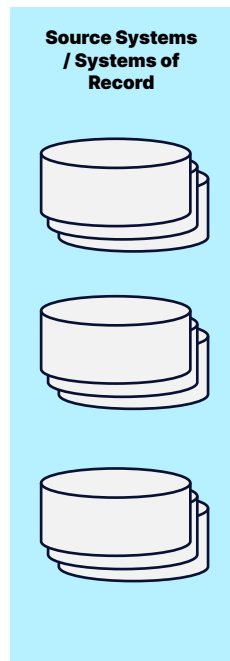


## II. The Legacy Data Problem

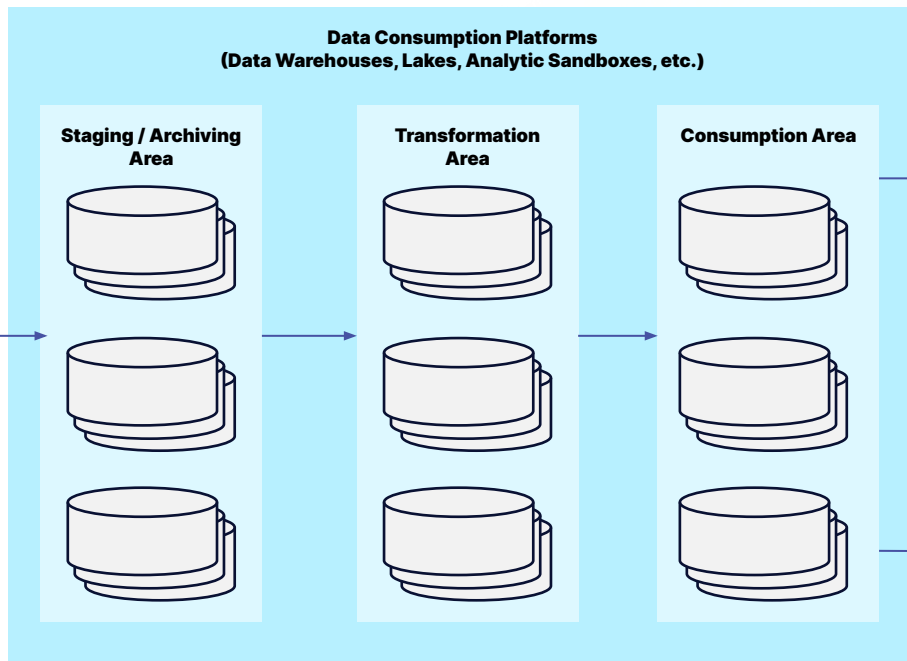
# The Problem In A Nutshell

## Common Conceptual Data Architecture Pattern in Most Large Enterprises

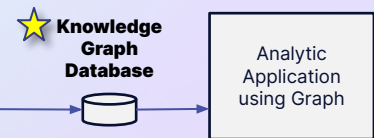
**Operational Data**  
GB / TB scale



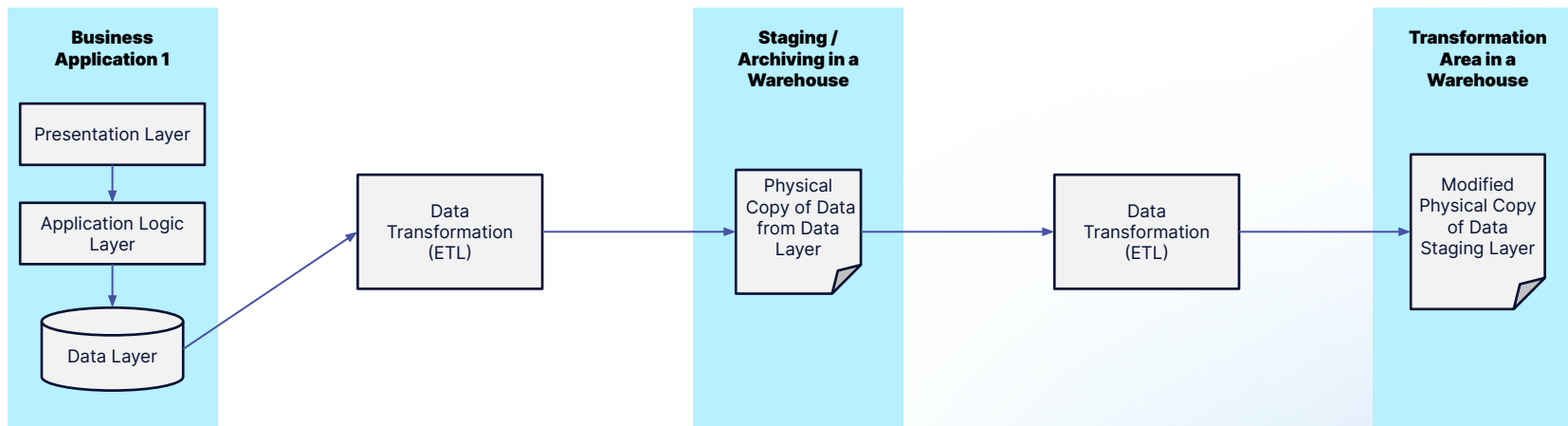
**Analytic Data**  
TB / PB scale



**Graph Analytic Data**  
KB / MB scale!!!



# Barrier 1: Contextless Data Proliferation



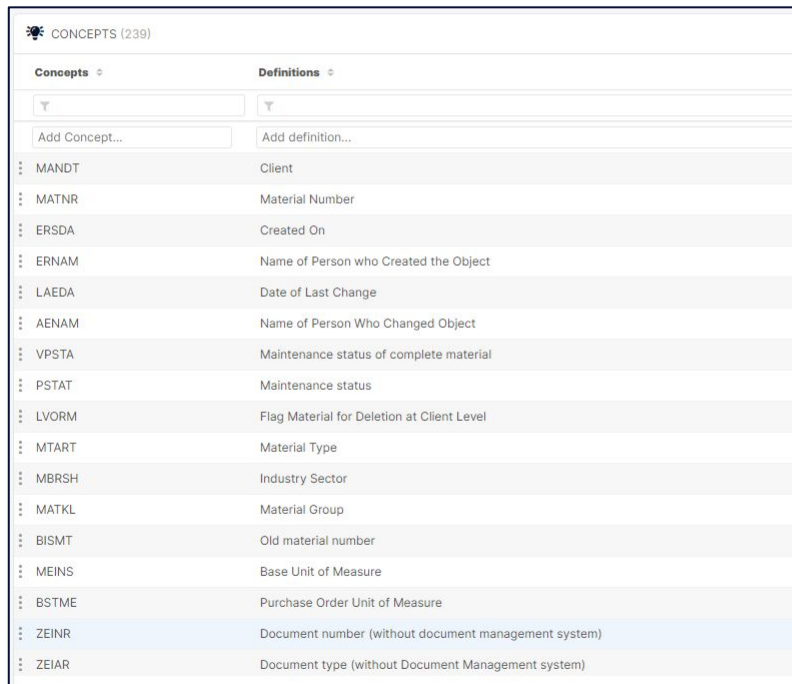
- Data created by app, for the app. Mostly for Operational use cases
- Semantic context for Data Layer embedded inside the application logic layer
- Data saved for a consumer other than the originating app
- But data is flat and has no context!!
- The copy has a different schema than the original
- Any consumer can create their own consumption schema on a use case by use case basis!

# Example: Data from an ERP System

mandt	matnr	ersda	ernam	laeda	aenam	vpsta	pstat	lvorm	mtart	mbrsh	matkl	bism
400	0000000000000000...	2001-06-04T00:00:0...	DNOBLE	2020-05-21T00:00:0...	PHOU	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1005	223
400	0000000000000000...	2001-06-04T00:00:0...	DNOBLE	2019-06-20T00:00:0...	MSMITH3	KDELBSVQGCZX	KDELBSVQGC	X	Z1H	P	1000	103
400	0000000000000000013	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1000	101
400	0000000000000000072	2001-06-04T00:00:0...	DNOBLE	2020-02-04T00:00:0...	RJACKAN	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1000	106
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2019-11-29T00:00:0...	LSANTACR	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1000	107
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2019-01-25T00:00:0...	MSMITH3	KDELBSVQGCZX	KDELBSVQGC	X	Z1H	P	1000	110
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1005	223
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2020-05-21T00:00:0...	PHOU	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1005	223
400	0000000000000000001	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1000	101
400	0000000000000000019	2001-06-04T00:00:0...	DNOBLE	2019-06-20T00:00:0...	MSMITH3	KDELBSVQGCZX	KDELBSVQGC	X	Z1H	P	1000	104
400	0000000000000000061	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1000	110
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2020-05-21T00:00:0...	PHOU	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1005	223
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1000	105
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1000	110
400	0000000000000000103	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1005	223
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1005	223
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1000	107
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1000	106
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	X	Z0H	P	1100	103
400	0000000000000000000	2001-06-04T00:00:0...	DNOBLE	2019-07-03T00:00:0...	BSAFINEJ	KDELBSVQGCZX	KDELBSVQGC	null	Z1H	P	1005	223

- What is this...?!!
- What do these columns mean...?

# ERP System Data Dictionary



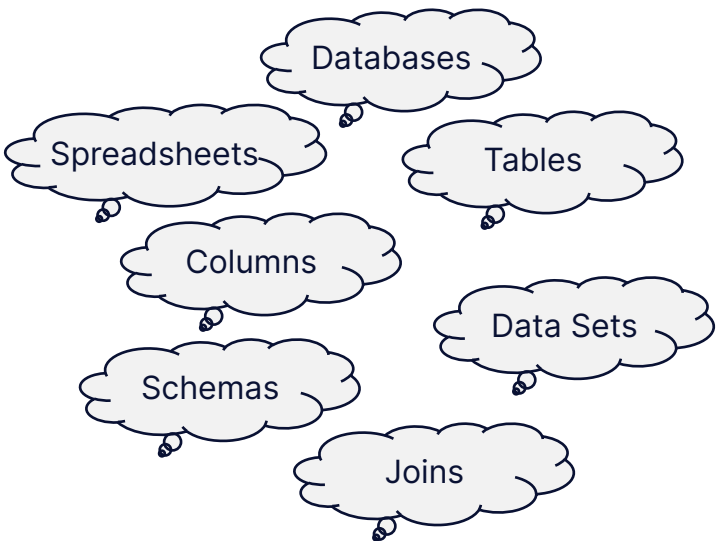
CONCEPTS (239)

Concepts	Definitions
<input type="text" value="Add Concept..."/>	<input type="text" value="Add definition..."/>
MANDT	Client
MATNR	Material Number
ERSDA	Created On
ERNAM	Name of Person who Created the Object
LAEDA	Date of Last Change
AENAM	Name of Person Who Changed Object
VPSTA	Maintenance status of complete material
PSTAT	Maintenance status
LVORM	Flag Material for Deletion at Client Level
MTART	Material Type
MBRSH	Industry Sector
MATKL	Material Group
BISMT	Old material number
MEINS	Base Unit of Measure
BSTME	Purchase Order Unit of Measure
ZEINR	Document number (without document management system)
ZEIAR	Document type (without Document Management system)

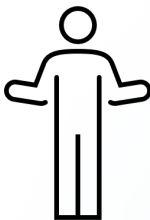
- OK, this table does have a meaning...
- It's just that the context resides in some document somewhere else **apart from the data** (e.g., some catalog system or document repository)
- Now multiply this problem for hundreds of apps and tens of thousands of database tables

# Barrier 2: Mental Model and Syntax

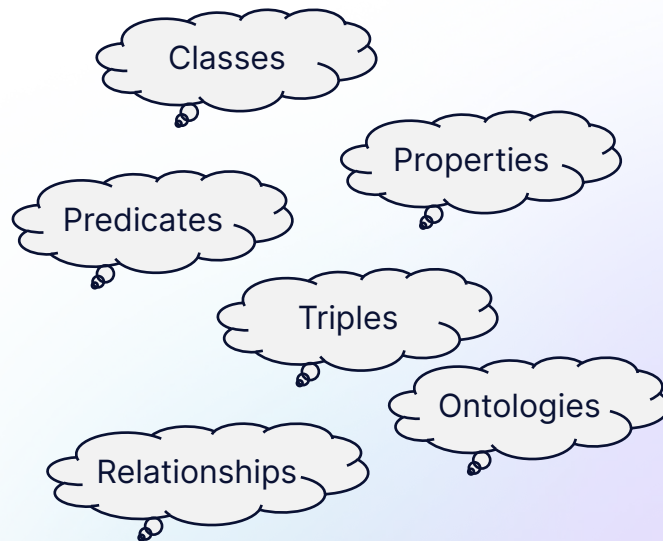
## Relational World



?????



## Semantic World



- Two completely different worlds
- Very few people in an organization can operate proficiently in both
- Years of legacy investment, time, training in Relational World



## III. The Innovations



# Two Technical Innovations

- **JSON Linked Data Standard (JSON-LD)**: a nice go-between Relational World and Semantic World
- **Robust data classification AI Models**: a method for layering semantics to flat data at the massive scale required

# What is JSON Linked Data? (JSON-LD)

- A method for encoding linked data ("semantic context") using JSON
- Easy to link ontologies to JSON documents and map to RDF

```
{ "@context": {  
  "name": "http://xmlns.com/foaf/0.1/name",  
  "homepage": {  
    "@id": "http://xmlns.com/foaf/0.1/workplaceHomepage",  
    "@type": "@id"  
  },  
  "Person": http://xmlns.com/foaf/0.1/Person  
},  
"@id": "https://me.example.com",  
"@type": "Person",  
"name": "John Smith",  
"homepage": "https://www.example.com/"  
}
```

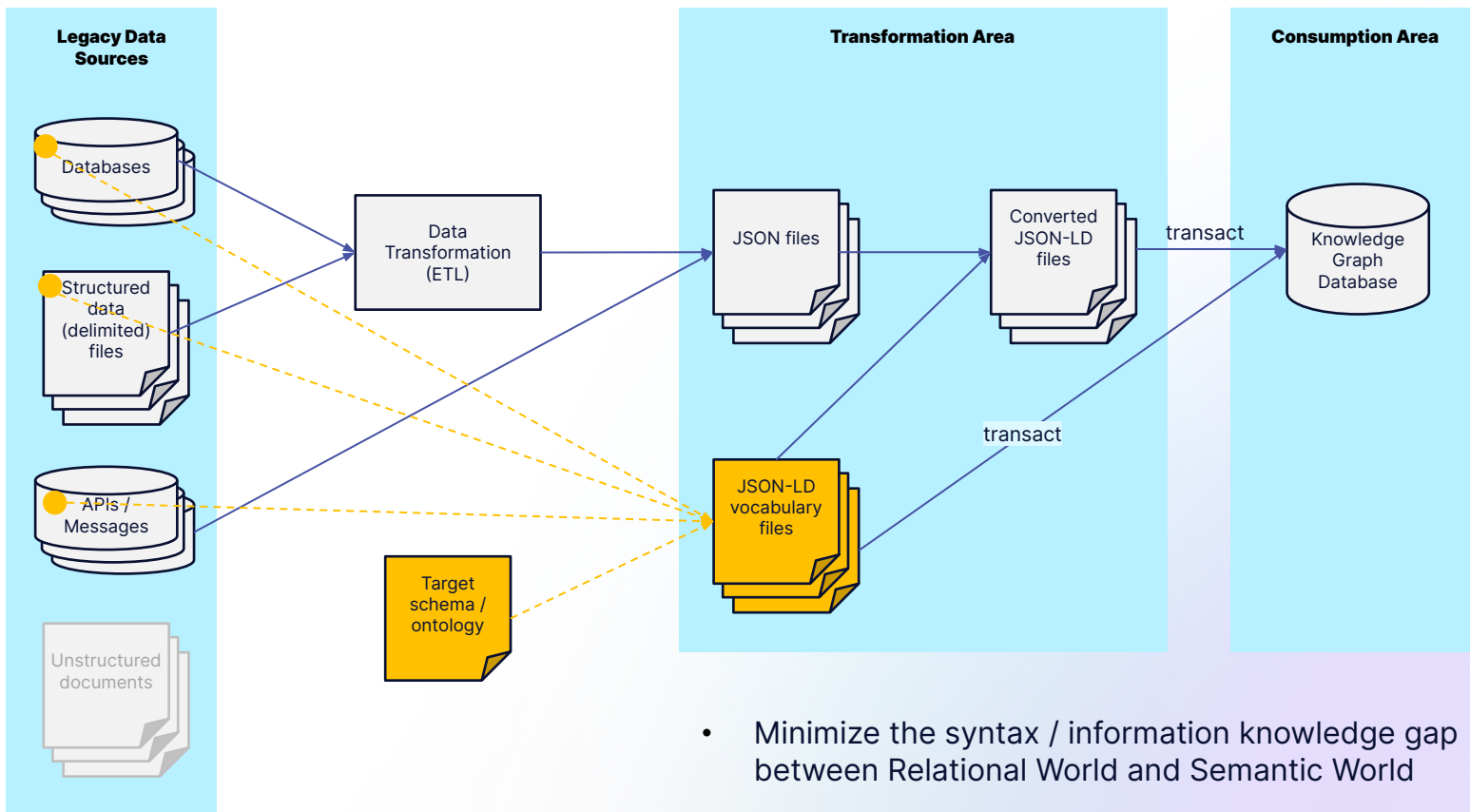
Describes the vocabulary  
being used to describe a  
Person entity

# Why does JSON-LD matter?

- IoT transacts in JSON
- Lots of APIs transact in JSON
- Very easy to ETL from CSV to JSON
- JSON-LD is not a huge leap from JSON
- Lots of websites use JSON-LD for Google Search Engine Optimization
- Generative AI uses it a lot!!! Easy fit for LLMs

# JSON-LD as a bridge to RDF / Triple Stores

- Represents a DTD, DDL, schema or data dictionary that describes the file located somewhere (inside the app logic layer, catalog, functional requirements doc, etc.)



- Minimize the syntax / information knowledge gap between Relational World and Semantic World

# Using Data Classification Models to add semantics at scale

- Machine learning models can be trained how to tag columns to vocabulary terms and then infer vocabulary to vocabulary synonym terms
- Once the model is trained, you can run an automated process to scan all columns and then classify them to their respective terms.
- Once complete, you can convert the discovered mappings and semantic relationships into json-ld files that can be transacted into a Knowledge Graph.
- There are many ML classification models; for this talk we will go through an example using a Neural Network

# AI/ML Data Classification Methodology

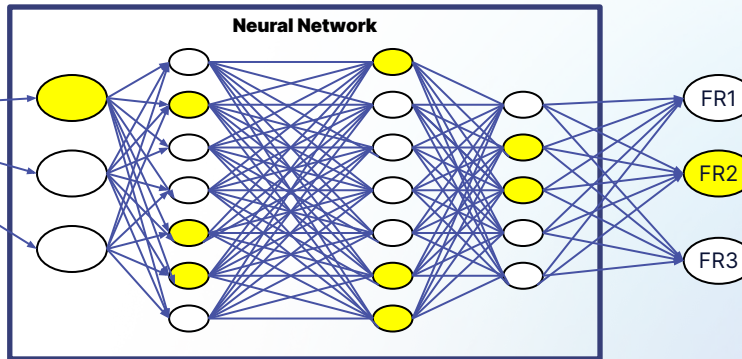
## Training Data

Col A	Col B	Col C	Col D
00001	Apple	Red	Shiny
00002	Banana	Yellow	Long
00003	Grape	Purple	Squishy

## Classifiers (Vocabulary for an Entity Called "Fruit")

ID	Term	Definition
FR1	Fruit Name	The name of the fruit
FR2	Fruit Color	The hue or color of the fruit
FR3	Characteristic	Defining feature whether shape or texture used to best describe the fruit

Col C
Red
Yellow
Purple



Col	Mapping
Col B	FR1
Col C	FR2
Col D	FR3

# Scaling the Process Through Crowdsourcing

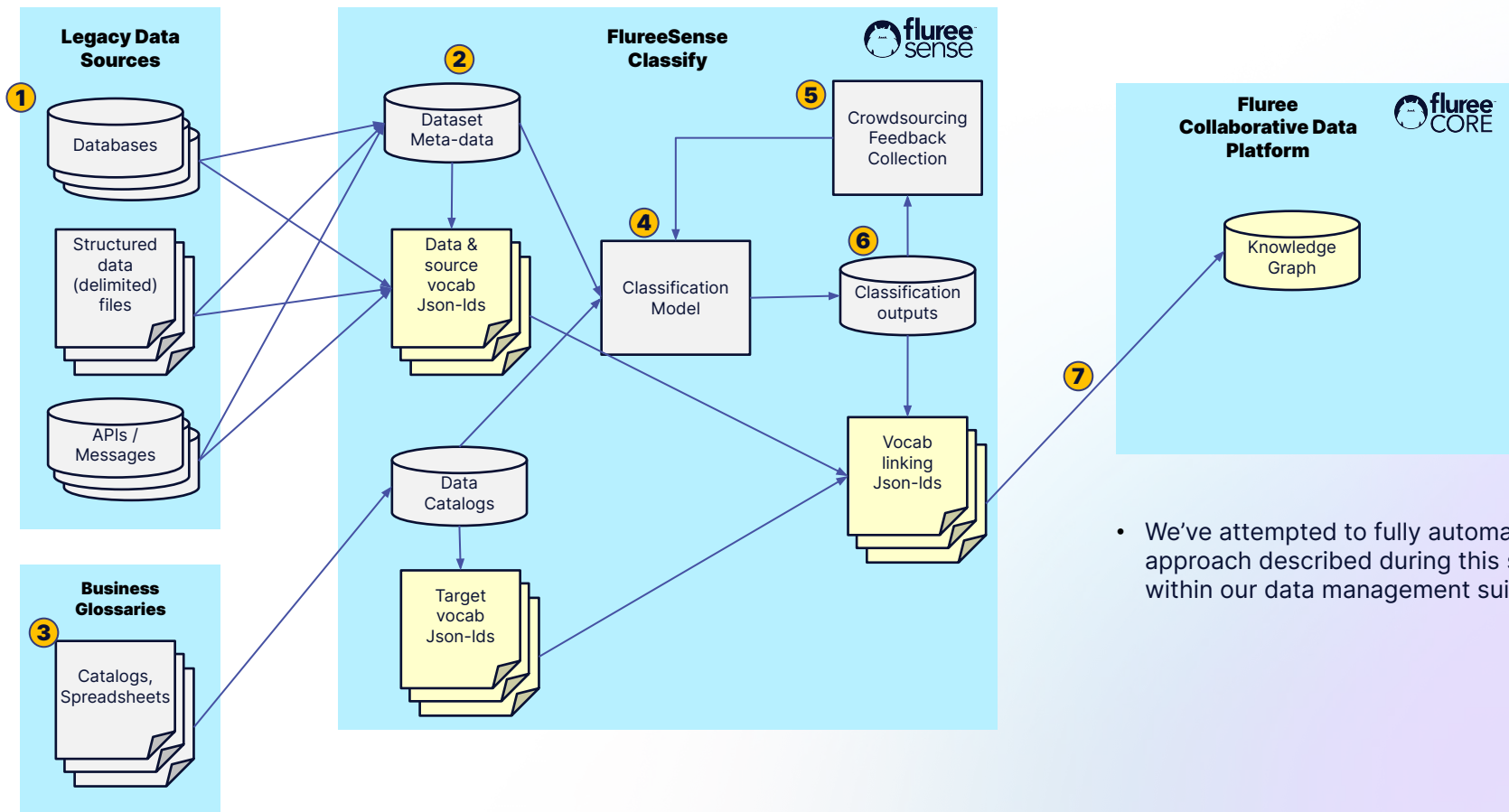
- The power of the approach is to distribute the training to as many people as possible
- Tens of people can individually classify 1-2 tables, or validate the machine prediction and give it reinforcement
- This will generate thousands or tens of thousands of new features to further refine the model





## IV. The Demo

# Fluree's Approach to Legacy Data Conversion



- We've attempted to fully automate the approach described during this session within our data management suite

# Step 1. Connect to Data Sources

fluree sense Classify Start by Default search

Pharma Demo User Sign Out

**Data Sources**

Create New Data Source +

DATA SOURCES (10)

Data Source Name	Data Source Type	Created By	Last Modified	# of Datasets	# of Group Entitlements	# of User Entitlements
IDMP Mapping Full	Hadoop	pharma_user	04/11/2023	10	1	6
IDMP Mapping	Hadoop	pharma_user	04/11/2023	14	1	6
SAP ECC	Hadoop	pharma_user	03/29/2023	2	1	6
GSRS GSRS registration system for the ingredients in medicinal products	Hadoop	pharma_user	04/7/2023	2	1	6
IQVIA Market intelligence of pharmaceutical customers, products and clinical inf...	Hadoop	pharma_user	04/11/2023	5	1	6
GEN Pharma Default Datasource Default Datasource created while creating a tenant	Hadoop	59_user	04/7/2023	13	1	6

fluree | Copyright © 2023 Fluree - All rights reserved | Version: v0.2.1 Environment: fspa

# Step 2. Generate Data Set Profiles

Classify Start by Default

PDU Pharma Demo User  
Sign Out

Data Sets > **MARA**

Data Set Description

**SAP Material data used to retain pharma drug mate...**

Object Tags

MARA
Medicinal Products
MedicinalProduct ...

Data Source

**SAP ECC**

Data Set Type

**Parquet**

Total Records

**273,105**

Total of Columns

**22**

Last Modified On

**03/29/2023**

6a

Overall quality score

**100%**

**Data Set Attributes**
Data Set Sample
Data Set Relationships
Data Refresh
Data Entitlements
Related Projects
Data Quality

Attribute	Data Type	Classification
ManufacturerCode	STRING	
ManufacturerDesc	STRING	
MAKTX	STRING	
mara_ZZ_PLANT	LONG	
STRU_DESC	STRING	
FAMIL_DESC	STRING	
SUBF_DESC	STRING	
STRE_DESC	STRING	
PAFM_DESC	STRING	
PASI_DESC	LONG	
MATNR	LONG	
mara_MATNR	LONG	
WRKST	STRING	
ZZ_STRUCTURE	LONG	
ZZ_PFAMGRP	STRING	
mara_ZZ_PFAMGRP	STRING	
ZZ_PRODFM	STRING	
MTART	STRING	
UMREN	STRING	
LAEDA	DATE	

**Profiler Result**

Total Size for Profiling: 273,105 Rows  
Total Rows: 273,105

Total Columns: 22  
Total Sampled Rows: 273,105

Last Extract on  
03/29/2023

**Results**

uniques 13 (0.0%)

Nulls 252974.0 (92.63%)

**Distribution Of Values**

Min. Value 2

Max. Value 99

**Data Types**

NUMBER 100.0%

**Regex Patterns**

9 72.89%

99 27.11%

**Top 10 Most Frequent**

252974	15
8784	96
4694	105
2408	136
1240	168
1166	357
536	396
396	536
357	1186
168	1240

**Top 10 Least Frequent**

94	15
12	96
97	105
96	136
11	168
16	357
98	396
95	536
9	1186
1240	1240

**Classification Result**

Run Model

Copyright © 2023 Fluree - All rights reserved

Version: v0.2.1 Environment: fsqa

# Step 3. Import Ontologies

fluree Classify Start by Default search

Pharma Demo User Sign Out

Data Catalogs

Create New Data Catalog

DATA CATALOGS (6)

Data Catalog	Data Catalog Description	Created By	Last Modified On	Semantic Objects	Concepts	Mapped Data Sources	Mapped Data Sets	Mapped Data Set Columns	# of Data Quality Rules	Measured Data Quality
Automation_catalogtest	null	Pharma Demo User	05/2/2023	0	0	0	0	0	Add	n/a
GEN Pharma Catalog	Business information model for Generics Pharmaceutical company	Pharma Demo User	03/30/2023	3	32	3	3	51	Add	n/a
Global Substance Registration System (GSRS)	Registration system for the ingredients in medicinal products. It makes it easier for regulators and other stakeholders.	Pharma Demo User	02/27/2023	34	786	2	2	25	2	100%
Identification of Medicinal Products (IDMP)	IDMP is a suite of five standards developed within the ISO to facilitate the unique identification of medicinal products.	Pharma Demo User	03/9/2023	3	90	3	3	32	2	100%
SAP ECC Data Dictionary	Data glossary for SAP ECC R3 system	Pharma Demo User	04/7/2023	75	4784	1	1	11	Add	n/a
zFlureeSense System Catalog	Out of the box system catalog and semantic entities	Fluree Sense	02/27/2023	2	262	1	1	2	Add	n/a

fluree Copyright © 2023 Fluree - All rights reserved. Version: v0.2.1 Environment: fsqa

# Step 3. Import Ontologies (cont'd)

fluree sense Classify Start by Default search

Data Catalogs > Identification of Medicinal Products (IDMP) > Pharmaceutical Products

Technical View Business View

**Pharmaceutical Products**  
Data elements and structures for unique identification and exchange of regulated pharmaceutical product information

3 Data Sources As of 5/8/2023

4 Data Sets As of 5/8/2023

0 Data Quality Rules As of 5/8/2023

n/a Measured Data Quality As of 5/8/2023

CONCEPTS (19)

Concepts	Definitions	Tags	Actions	Synonyms	Mapped Data Columns	Train Model	# of Data Quality Rules	Measured Data Quality
Add Concept...	Add definition...							
Reference Strength	strength of an active substance(s) and/or specified substance(s) used as a reference from which the strength of an inv...	+ Add Tag		1	0		Add	n/a
Presentation Strength	quantity or range of quantities of the substance/specified substance present per unitary volume (or mass) expressed in...	+ Add Tag		1	1		Add	n/a
Pharmaceutical Product Identifier	unique identifier for a pharmaceutical product	+ Add Tag		2	3		Add	n/a
Placebo	inactive substance, treatment or procedure that is intended to provide baseline measurements for the experimental pro...	+ Add Tag		1	0		Add	n/a
Packaged Medicinal Product	medicinal product in a container being part of a package, representing the entirety that has been packaged for sale or s...	+ Add Tag		1	0		Add	n/a
Investigational Medicinal Product I...	unique identifier allocated to an investigational medicinal product supplementary to any existing identifier as ascribed b...	+ Add Tag		0	0		Add	n/a
Clinical Trial	investigation in human subjects intended to discover or verify the clinical, pharmacological and/or other pharmacodyna...	+ Add Tag		2	0		Add	n/a
Package Item		+ Add Tag		1	1		Add	n/a
Outer Packaging		+ Add Tag		1	0		Add	n/a
Intermediate Packaging		+ Add Tag		1	0		Add	n/a
Immediate Container		+ Add Tag		1	0		Add	n/a
Investigational Medicinal Product		+ Add Tag		0	0		Add	n/a
Medicinal Product Identifier	unique identifier allocated to a medicinal product supplementary to any existing authorization number as ascribed by a ...	+ Add Tag		0	0		Add	n/a
Medicines Regulatory Agency	institutional body that, according to the legal system under which it has been established, is responsible for the grantin...	+ Add Tag		1	0		Add	n/a
Mass Based Strength		+ Add Tag		1	0		Add	n/a
Concentration		+ Add Tag		1	0		Add	n/a
Activity Based Strength		+ Add Tag		1	0		Add	n/a

fluree | Copyright © 2023 Fluree - All rights reserved. | Version: v0.2.1 Environment: fscg

# Step 4. Train Classification Model

fluree sense Classify Start by Default search

Data Catalogs > Identification of Medicinal Products (IDMP) > Pharmaceutical Products > Pharmaceutical Product Identifier > Train Model

Review Prior Training Data Select Data Sets Map Data Columns Finalize Training Data

Map the column from the data sets selected that accurately represent the Concept that is being trained.

**TM TOTAL DATA ATTRIBUTES (8)**

Data set	Attribute	Mapping(0)
IDMP Medicinal Products - IMS	Manufacturer	
IDMP Medicinal Products - IMS	MedicinalProductIdentifier	
IDMP Medicinal Products - IMS	MedicinalProductName	
IDMP Medicinal Products - IMS	Packagetem	
IDMP Medicinal Products - IMS	PharmaceuticalProductIdentifier	
IDMP Medicinal Products - IMS	PresentationStrength	
IDMP Medicinal Products - IMS	RegistrationNumber	
IDMP Medicinal Products - IMS	id	

**PHARMACEUTICAL PRODUCTS SEMANTIC OBJECT MODEL MAPPINGS (0 FIELDS MAPPED)**

Concept	Mapping
Pharmaceutical Product Identifier	Drag & drop data set attribute here.

IDMP Medicinal Products - IMS  
PharmaceuticalProductIdentifier

Cancel Previous Step Next Step

fluree | Copyright © 2023 Fluree - All rights reserved | Version: v0.2.1 Environment: fsqa



# Step 5. Review results and calibrate model

Classify

Start by Default

Data Sets > **MARA**

Data Set Description  
**SAP Material data used to retain pharma drug mate...**

Object Tags  
MARA Medicinal Products MedicinalProduct ...

Data Source  
**SAP ECC**

Data Set Type  
**Parquet**

Total Records  
**273,105**

Total of Columns  
**22**

Last Modified On  
**03/29/2023**

Overall quality score **100%**

**Data Set Attributes**

Data Set Sample

Data Set Relationships

Data Refresh

Data Entitlements

Related Projects

Data Quality

Attribute	Data Type	Classification
ManufacturerCode	STRING	ManufacturerId
ManufacturerDesc	STRING	Manufacturer ManufacturerName
MAKTX	STRING	Medicinal Product Name ...
mara_ZZ_PLANT	LONG	ZZ_PLANT
STRU_DESC	STRING	
FAMIL_DESC	STRING	Agent Substance Name ...
SUBF_DESC	STRING	Substance Name Classifier...
STRE_DESC	STRING	Strength
PAFM_DESC	STRING	
PASI_DESC	LONG	
MATNR	LONG	MATNR ...
mara_MATNR	LONG	MATNR
WRKST	STRING	Medicinal Product Name ...
ZZ_STRUCTURE	LONG	ZZ_STRUCTURE
ZZ_PFMGRP	STRING	Medicinal Product Name ...
mara_ZZ_PFMGRP	STRING	Medicinal Product Name ...
ZZ_PRODFM	STRING	MedicinalProductHierarchyG...
MTART	STRING	MTART
UMREN	STRING	
LAEDA	DATE	LAEDA

**Profiler Result**

Total Size for Profiling: 273,105 Rows  
Total Rows: 273,105

**Results**

uniques	13 (0.0%)
Nulls	252974.0 (92.63%)

**Distribution Of Values**

Min. Value	2
Max. Value	99

Total Columns: 22  
Total Sampled Rows: 273,105  
Last Extract on: 03/29/2023

**Data Types**

NUMBER	100.0%
--------	--------

**Regex Patterns**

9	72.89%
99	27.11%

**Top 10 Most Frequent**

252974	94
8794	12
4694	97
2408	96
1240	11
1166	16
536	98
396	95
357	9
168	99

**Top 10 Least Frequent**

15	96
105	105
136	168
357	357
396	396
536	536
1186	1186
1240	1240

**Classification Result**

Semantic Object	Concept	Mapping Confidence Level
MARA	ZZ_PLANT	<span style="color: green; font-weight: bold;">HIGH</span> 100%  1
MedicinalProduct	Pharmaceutical.	<span style="color: red; font-weight: bold;">LOW</span> 35%
Medicinal Products	Pharmaceutical.	<span style="color: red; font-weight: bold;">LOW</span> 31%
MedicinalProduct	MedicinalProdu.	<span style="color: red; font-weight: bold;">LOW</span> 0%   1
MedicinalProduct	MedicinalProdu.	<span style="color: red; font-weight: bold;">LOW</span> 0%   1

Run Model

Copyright © 2023 Fluree - All rights reserved.

Version: v0.2.1 Environment: fsqa

# Step 5. Review results and calibrate model (cont'd)

fluree sense Classify Start by Default search

Data Sets > MARA

Data Set Description: SAP Material data used to retain pharma drug mate...  
 Object Tags: MARA Medicinal Products MedicinalProduct ...

Data Source: SAP ECC  
 Data Set Type: Parquet  
 Total Records: 273,105  
 Total of Columns: 22  
 Last Modified On: 03/29/2023  
 Overall quality score: 100%

**Edit Classification Tag : ManufacturerId**

You've given this mapping a thumbs down.  
 Keep my selected mapping  
 - Or select a better mapping:

CLASSIFICATION TAGS

Catalog Name	Semantic Object	Concept	Mapping Confidence Level
GEN Pharma Catalog	MedicinalProduct	Strength	HIGH 100%
Identification of Medicinal Products	Pharmaceutical Products	Presentation Strength	LOW 47%
GEN Pharma Catalog	PharmaceuticalProduct	SubstanceAmount	LOW 47%
GEN Pharma Catalog	MedicinalProduct	Manufacturer	LOW 42%
GEN Pharma Catalog	MedicinalProduct	MedicinalProductHierarchyGrou	LOW 41%
GEN Pharma Catalog	MedicinalProduct	ManufacturerId	LOW 35%
GEN Pharma Catalog	MedicinalProduct	Manufacturer	LOW 33%
GEN Pharma Catalog	MedicinalProduct	Manufacturer	LOW 0%

I'm Not Sure

Buttons: Cancel Save & Close

Mapping Confidence Level table:

Concept	Mapping Confidence Level	Score	Feedback
Strength	HIGH	100%	1 thumbs up
Presentation St.	LOW	47%	1 thumbs down
SubstanceAmo.	LOW	47%	1 thumbs down
Manufacturer	LOW	42%	1 thumbs down
MedicinalProdu.	LOW	41%	1 thumbs down
Manufacturerid	LOW	35%	1 thumbs down
MedicinalProdu.	LOW	33%	1 thumbs down
Presentation St.	LOW	0%	1 thumbs down

Run Model

fluree Copyright © 2023 Fluree - All rights reserved. Version: v0.2.1 Environment: foga

# Step 5. Review results and calibrate model (cont'd)

Classify
Start by Default

Pharma Demo User  
sign out

Data Sets > **MARA**

Data Set Description

**SAP Material data used to retain pharma drug mate...**

Object Tags

MARA Medicinal Products MedicinalProduct ...

Data Source

**SAP ECC**

Data Set Type

**Parquet**

Total Records

**273,105**

Total of Columns

**22**

Last Modified On

**03/29/2023**

Overall quality score

**100%**

Classify Set

Data Set Attributes    Data Set Sample    **Data Set Relationships**    Data Refresh    Data Entitlements    Related Projects    Data Quality

**Data Set Relationships:** Please pinch or zoom using your mouse or Trackpad to zoom in or zoom out the diagram

Technical View    **Business View**

The diagram illustrates the relationships between four data sources:

- MARA:** ManufacturerCode, ManufacturerDesc, MATX, mara\_zz\_PLANT, STRU\_DESC, FAM1\_DESC, SUBP\_DESC, STRE\_DESC, PAFM\_DESC, PAS1\_DESC, MATNR, mara\_MATNR, WRKST, ZZ\_STRUCTURE, ZZ\_PPAMORP, mara\_ZZ\_PPAMORP, ZZ\_PRODFM, MTART, UMREN, LAEDA, mara\_NORMT, NORMT.
- GEN Pharma Catalog:** PharmaceuticalProduct (Manufacturer, PharmaceuticalProductid), MedicinalProduct (Strength, MedicinalProductid, Manufacturerid, ManufacturerName, MedicinalProductName, MedicinalProductHierarchyGroupName, PharmaceuticalProductid), Substance (SubstanceName, SubstanceOfficialName).
- SAP ECC Data Dictionary:** MARA, MATNR, LAEDA, MTART, NORMT, ZZ\_PPAMORP, ZZ\_PLANT, ZZ\_STRUCTURE, ZZ\_PRODFM.
- Identification of Medicinal Products (IDMP):** Pharmaceutical Products (Pharmaceutical Product Identifier, Manufacturer), Substance Identification (Substance Name Classifier-Official Name).

Relationships are shown as lines connecting fields between these categories.

Copyright © 2023 Fluree - All rights reserved

Version: v0.2.1    Environment: fscg

# Step 6. Validate Discovered Synonyms

fluree sense Classify Start by Default search

Pharma Demo User

## Synonyms

SYNONYM NAME: SUBSTANCE IDENTIFICATION:(15)

From Original Catalog: Identification of Medicinal Products (IDMP)

From Other Catalogs : Global Substance Registration System (GSRS)

Catalog	Semantic Object	Concept	Synonym Catalog	Synonym Object	Synonym Concept	Synonym Confidence
Identification of Me...	Substance Identification	Moiety Identifier	Global Substance Registration System (GSRS)	Substance	Approval ID	HIGH 100%  1
Identification of Me...	Substance Identification	Molecular Structure	Global Substance Registration System (GSRS)	Structure	Moiety Formula	HIGH 100%  1
Identification of Me...	Substance Identification	Substance Code	Global Substance Registration System (GSRS)	Substance	Approval ID	HIGH 100%  1
Identification of Me...	Substance Identification	Nucleic Acid Substance	Global Substance Registration System (GSRS)	NucleicAcid	Nucleic Acid UUID	HIGH 100%  1
Identification of Me...	Substance Identification	Polymer Substance	Global Substance Registration System (GSRS)	Polymer	Polymer UUID	HIGH 100%  1
Identification of Me...	Substance Identification	Protein Substance	Global Substance Registration System (GSRS)	Protein	Protein UUID	HIGH 100%  1
Identification of Me...	Substance Identification	Substance Name Clas...	Global Substance Registration System (GSRS)	SubstanceReference	Agent Substance Name	HIGH 100%  1
Identification of Me...	Substance Identification	Excipient	GEN Pharma Catalog	Substance	Substanceld	HIGH 100%  1
Identification of Me...	Substance Identification	Moiety	GEN Pharma Catalog	Substance	Substanceld	HIGH 100%  1
Identification of Me...	Substance Identification	Substance Code	GEN Pharma Catalog	Substance	Substanceld	HIGH 100%  1
Identification of Me...	Substance Identification	Substance Name Clas...	GEN Pharma Catalog	Substance	SubstanceName	HIGH 100%  1
Identification of Me...	Substance Identification	Substance Name Clas...	GEN Pharma Catalog	Substance	SubstanceSynonymName	HIGH 100%  1
Identification of Me...	Substance Identification	Substance Code	GEN Pharma Catalog	PharmaceuticalProduct	Substanceld	HIGH 100%  1
Identification of Me...	Substance Identification	Substance Name Clas...	GEN Pharma Catalog	Substance	SubstanceOfficialName	HIGH 100%  1
Identification of Me...	Substance Identification	Molecular Structure	GEN Pharma Catalog	Substance	SubstanceMolecularStructureId	HIGH 100%  1

Save & Close Save & Run

# Step 7. Publish Semantic Data Sets

The screenshot displays the 'Object Model Diagram' for 'Medicinal Products' within the fluree Sense application. The interface includes a top navigation bar with 'fluree Sense', 'Classify', 'Start by Default', and a search bar. The breadcrumb trail shows 'Data Catalogs > Identification of Medicinal Products (IDMP) > Medicinal Products > Medicinal Products Object Model'. On the left, a sidebar contains navigation icons. The main area shows three tables: 'SAP ECC' (with fields: MARA, Manufacturer\_Desc, MAKTX, MATNR, WKST, ZZ\_PFAMGRP, mara\_ZZ\_PFAMGRP, NORMT), 'IQVIA' (with fields: IQVIA Product List, DIN, Product\_Name, Product\_Strength, Product\_Pack, Manufacturer\_Desc, CMF10, Product\_Key), and 'Medicinal Products' (with fields: Active Ingredient Entire Substance Basis of Strength, Medicinal Product Identifier, Batch Manufacturing Process, Marketing Authorization Number, Mass Based Strength, Active Ingredient Reference Substance Basis of Strength, Batch Number, Active Ingredient Active Moiety Basis of Strength, Medicinal Product Name, Clinical Trial Identifier, Lot Number, Reference Strength, Clinical Trial, Activity Based Strength, Registration Number, Investigational Medicinal Product, Continuous Manufacturing Process, Process Step Identifier, Outer Packaging, Process Identifier, Intermediate Packaging, Packaged Medicinal Product, Package Item, Active Ingredient Without Basis of Strength, Placebo, Product Constituency, Medicinal Product Batch Identifier for Outer Packaging, Medicinal Product Manufacturing Process, Discrete Manufacturing Process, Pharmaceutical Product Batch Number, Immediate Container, Manufacturer, Medicinal Product Batch Identifier for Immediate Packaging, Concentration, Clinical Trial Authorization, Investigational Medicinal Product Identifier, Pharmaceutical Product Identifier, Presentation Strength, Authorized Medicinal Product). Lines connect the source tables to the target table. At the bottom right, a red circle highlights the 'Publish Semantic Data Sets' button, with a red text annotation: 'This exports the selected source tables to full json-Id'. Other buttons include 'Back to Catalog' and 'Pharma Demo User sign out'.

This exports the selected source tables to full json-Id



## V. Getting Started In Your Organization

# Define Your Use Case

- Which business use case?
- Who are the users?
- Which ontology?
- Which data sources to start with?

Driver	Common Use Cases	Users	Ontology	Data Sources
LLM	<ul style="list-style-type: none"> <li>• Internal knowledge mgmt</li> <li>• Business process optimization</li> </ul>	<ul style="list-style-type: none"> <li>• Internal transformation teams</li> </ul>	<ul style="list-style-type: none"> <li>• Business process models</li> </ul>	<ul style="list-style-type: none"> <li>• KM systems</li> <li>• Project management systems</li> <li>• Accounting systems</li> </ul>
IT / OT Data Integration	<ul style="list-style-type: none"> <li>• Supply Chain Optimization</li> <li>• Smart Factory Enablement</li> </ul>	<ul style="list-style-type: none"> <li>• Factory Operations</li> </ul>	<ul style="list-style-type: none"> <li>• Standardized Bill of Materials</li> </ul>	<ul style="list-style-type: none"> <li>• Sensors</li> <li>• ERP</li> <li>• SCADA</li> <li>• PLCs</li> </ul>
Internal Data Sharing	<ul style="list-style-type: none"> <li>• Marketing (Upsell / Cross-sell)</li> <li>• Corporate Functions data sharing (Risk, Finance, Compliance)</li> </ul>	<ul style="list-style-type: none"> <li>• Business End Users</li> <li>• Corporate Functions End Users</li> </ul>	<ul style="list-style-type: none"> <li>• Internal Business Information Models</li> <li>• Upper Ontologies</li> </ul>	<ul style="list-style-type: none"> <li>• Customer</li> <li>• Product</li> <li>• Account</li> <li>• Transaction</li> </ul>
External Data Sharing	<ul style="list-style-type: none"> <li>• Supply Chain Optimization</li> <li>• Industry Data Sharing</li> </ul>	<ul style="list-style-type: none"> <li>• Operations across external Partnerships</li> </ul>	<ul style="list-style-type: none"> <li>• Industry Standard Ontologies</li> </ul>	<ul style="list-style-type: none"> <li>• Manufacturing data</li> </ul>



# Run a FailFast PoC

Phase I: Data Acquisition	Phase II: Classification Training	Phase III: Data Consumption
<ul style="list-style-type: none"> <li>• Connect to existing systems that you can easily access</li> <li>• Start with existing data warehouses and data lakes</li> <li>• Stand-up the analytic infrastructure</li> </ul>	<ul style="list-style-type: none"> <li>• Collect business information models and ontologies</li> <li>• Collect as much Source to Target mapping information, data catalogs, documented ETL logic</li> <li>• Import as training material for classification model</li> <li>• Conduct focused 2-week training with Human Subject Matter Experts</li> </ul>	<ul style="list-style-type: none"> <li>• Ingest data into Knowledge Graph</li> <li>• Generate test queries and results that demonstrate the business use case objectives</li> </ul>
<b>3-6 weeks</b>	<b>4 weeks</b>	<b>3 weeks</b>

# In Closing

- We're at this interesting juxtaposition where the world is both waking up to the need for semantics, AND the tools and capabilities for creating it at scale are better than ever
- We hope we've shared practical concepts on how it can be done
- Please do not hesitate to reach out to me if you are interested further in pursuing this journey!!
- **Eliud Polanco** – [epolanco@flur.ee](mailto:epolanco@flur.ee)

# Questions?



Keep in touch with us  
[flur.ee/sense](https://flur.ee/sense)