# Going Meta #26
## Unpicking the data.world benchmark on the role of KGs in LLM QA over enterprise data

**Semantic Layer as the Data Interface for LLMs**

A new benchmark for natural language questions against databases dropped last week.
What does it mean and how does the dbt Semantic Layer stack up?

JASON GANZ
26 NOV 2023

*layering structured Semantic Knowledge on top of your data leads to much str... ...ility to correctly answer ad-ho... questions about y... ...niza... ...rge Language Models...*

Raw Data

Data Platform          dbt Semantic Layer

Output

BI Tools

Metric
Data

Notebooks

ML Models

Semantic
Layer
Queries

Catalogs

Observability

Data Apps

https://roundup.getdbt.com/p/semantic-layer-as-the-data-interface

# The two systems

# The Unpicking... The ontology



**ACME_Insurance/ontology/insurance.ttl**

```
293    ###  http://data.world/schema/insurance/PolicyHolder
294    in:PolicyHolder rdf:type owl:Class ;
295                rdfs:isDefinedBy <http://data.world/schema/insurance/> ;
296                rdfs:label "Policy Holder" .


75     ###  http://data.world/schema/insurance/hasPolicyHolder
76     in:hasPolicyHolder rdf:type owl:ObjectProperty ;
77                rdfs:domain in:Policy ;
78                rdfs:range in:PolicyHolder ;
79                rdfs:isDefinedBy <http://data.world/schema/insurance/> ;
80                rdfs:label "has policy holder" .
81
```

https://github.com/datadotworld/cwd-benchmark-data

# The Unpicking... The mapping onto-source data



**ACME_Insurance/data/data.world_P&C_Insurance_Ontology_V1.r2rml**

```
41    map:TripleMap_PolicyHolderID_12 a rr:TriplesMap ;
42        rr:predicateObjectMap  [ rr:objectMap [ rr:column "party_identifier" ] ;
43                                 rr:predicate <http://data.world/schema/insurance/policyHolderId> ] ;
44        rr:subjectMap          [ rr:template "https://myinsurancecompany.linked.data.world/d/omg-pc-database/Policy-Holder-{party_
45        rr:logicalTable        [ rr:sqlQuery """select distinct party_identifier
46    from agreement_party_role
47    join policy on agreement_party_role.agreement_identifier = policy.policy_identifier
48    where agreement_party_role.party_role_code = 'PH'""" ;
```
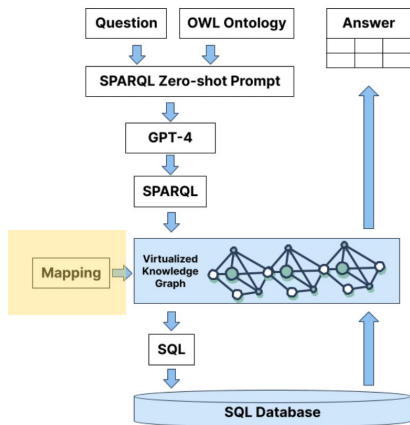
Agreement

Policy — Agreement_Party_role — Party_role — Party

https://github.com/datadotworld/cwd-benchmark-data

# The Unpicking… The mapping engine





https://github.com/chrdebru/r2rml

# The Unpicking... Alternatives to formalise the semantic layer



models/omg_semantics/agreement_party_role.yaml



https://github.com/dbt-labs/semantic-layer-llm-benchmarking

# A Neo4j based version



1. Build the property graph building the semantic layer into the data
   a. Use the ontology as target model
   b. Data-Import tool (or code based alternatives) check GM#5
2. Get LLM to retrieve the ontology from the graph in Neo4j and generate the cypher query

https://github.com/jbarrasa/goingmeta

# The completions API (or manually in ChatGPT)

```python
question = "How many policies have agents sold by agent id"

system = "You are an assistant capable of creating cypher queries from

prompt=f"""

Given a property graph with the structure described by the following li

GRAPH MODEL: {graph_struct}

Write a Cypher query that answer

QUESTION: {question}
"""
```

```python
from openai import OpenA

client = OpenAI()
completion = client.chat
    model="gpt-4",
    temperature=0,
    messages=[
        {"role": "sys
        {"role": "use
    ])
```

**You**

Given a property graph with the structure described by the following list of node and relationship types and their properties:

GRAPH MODEL: Label:Policy. Properties:policyEffectiveDate, policyExpirationDate, policyNumber, policyId
Label:PolicyHolder. Properties:policyHolderId
Label:Agent. Properties:agentId
relationship type: hasPolicyHolder, from: Policy, to: PolicyHolder
relationship type: soldByAgent, from: Policy, to: Agent

Write a Cypher query that answers the question. Do not explain the query. Return just the query, so it can be run verbatim from your response.

QUESTION: How many policies have agents sold by agent id?

**ChatGPT**

```cypher
cypher                                          Copy code

MATCH (:Policy)-[:soldByAgent]->(a:Agent)
RETURN a.agentId AS AgentID, COUNT(*) AS PoliciesSold
GROUP BY a.agentId
ORDER BY PoliciesSold DESC
```