

图数据库解谜与数据质量应用

吉思为

DEVELOPER ADVOCATE @  vesoft



国际软件质量工程
International Software Quality Engineering

古思为

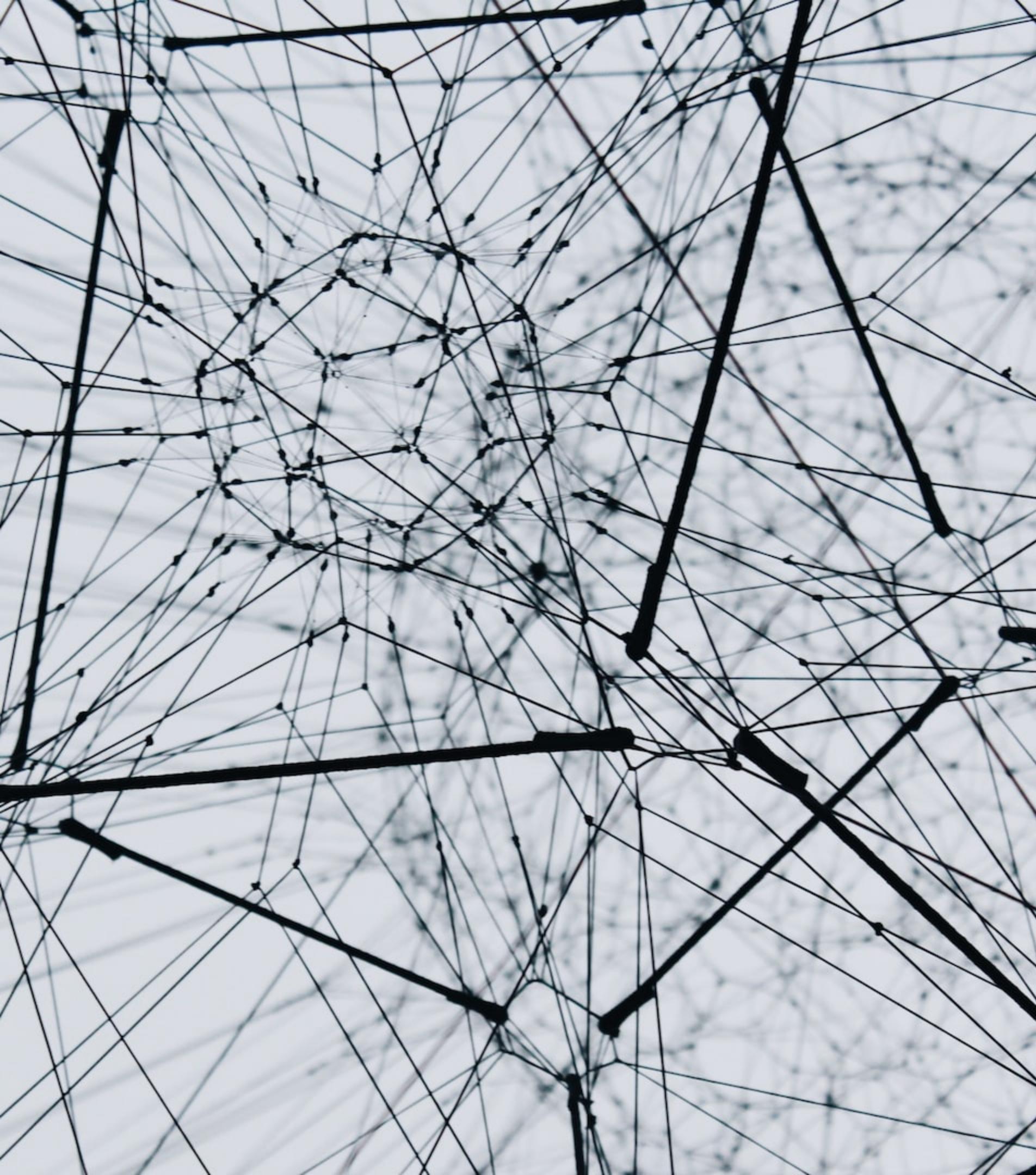
- Software Engineer @ Shanghai
- Open Source Believer
- Developer Advocate of NebulaGraph @vesoft

❏ [wey-gu](#)

❏ [wey_gu](#)

❏ [siwei.io/about](#)



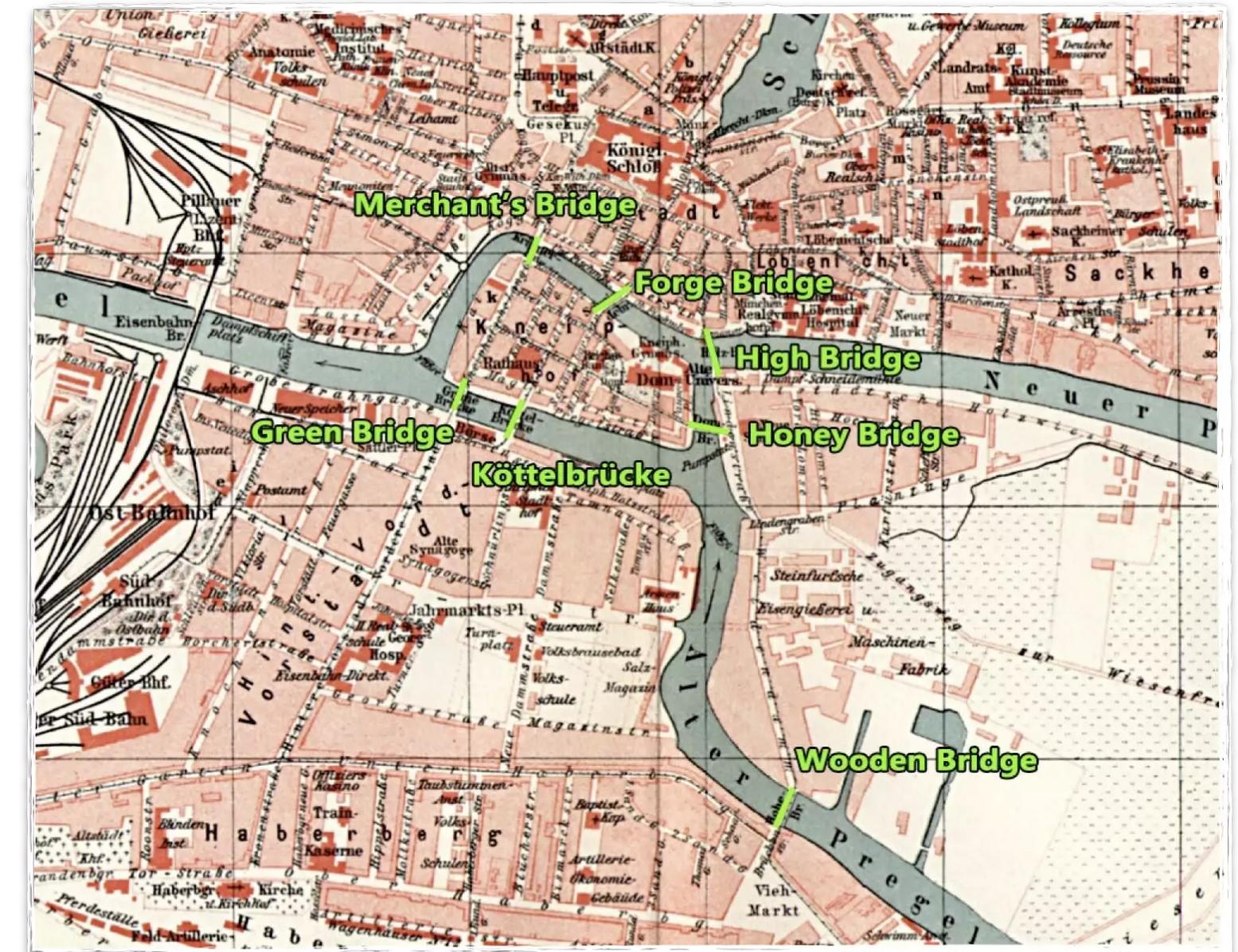
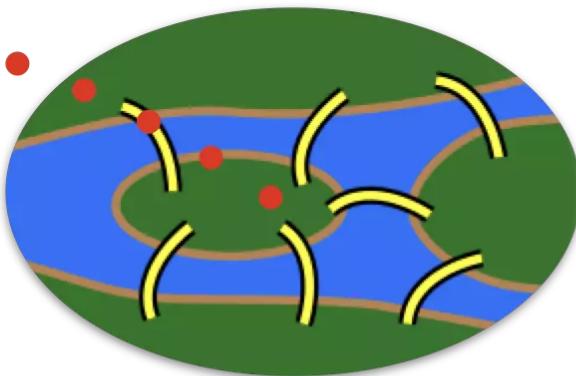
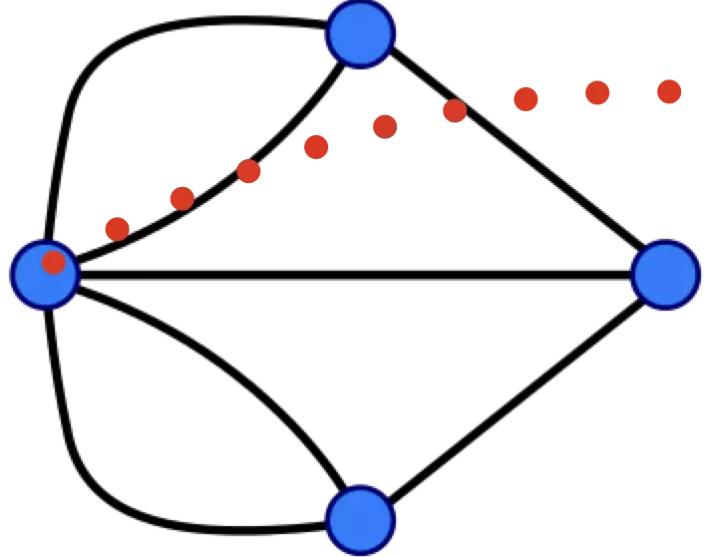


Overview

- Graph Database Explained
- Graph Tech Application: Data Lineage

Graph Database

What is Graph? What is Graph DB? Why yet another DB?



Map of Königsberg with the seven bridges labeled, circa 1905

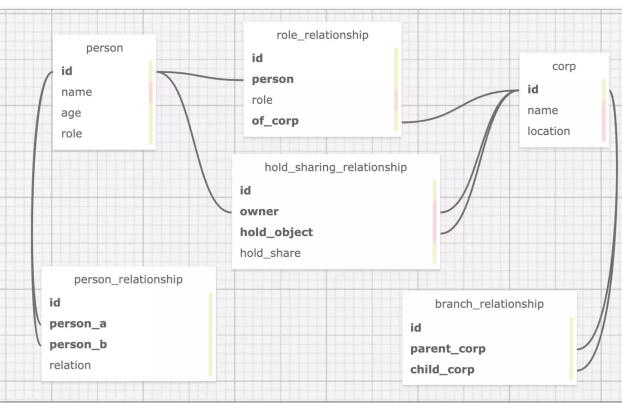
"A database that uses graph structures for semantic queries with nodes, edges, and properties to represent and store data

[wikipedia.org/wiki/graph_database](https://en.wikipedia.org/wiki/Graph_database)

More on what a GDB is

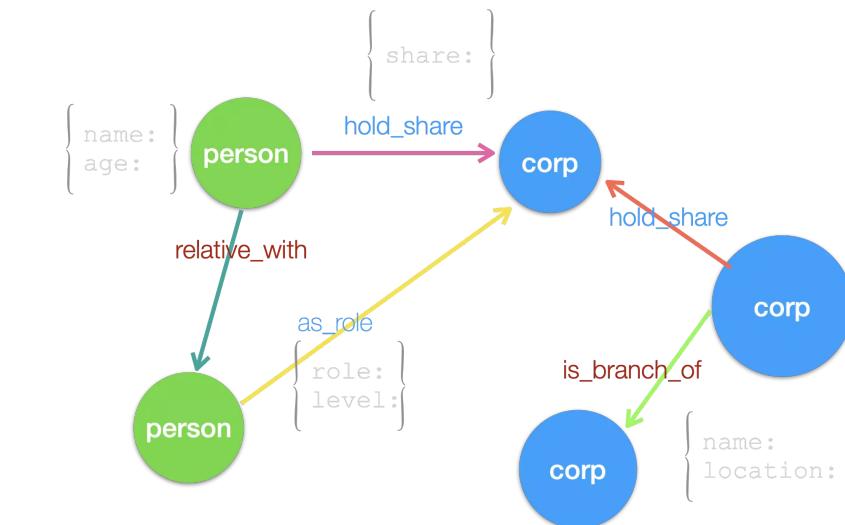
Why Yet Another DB?

Relational DB



Graph Schema

Graph DB



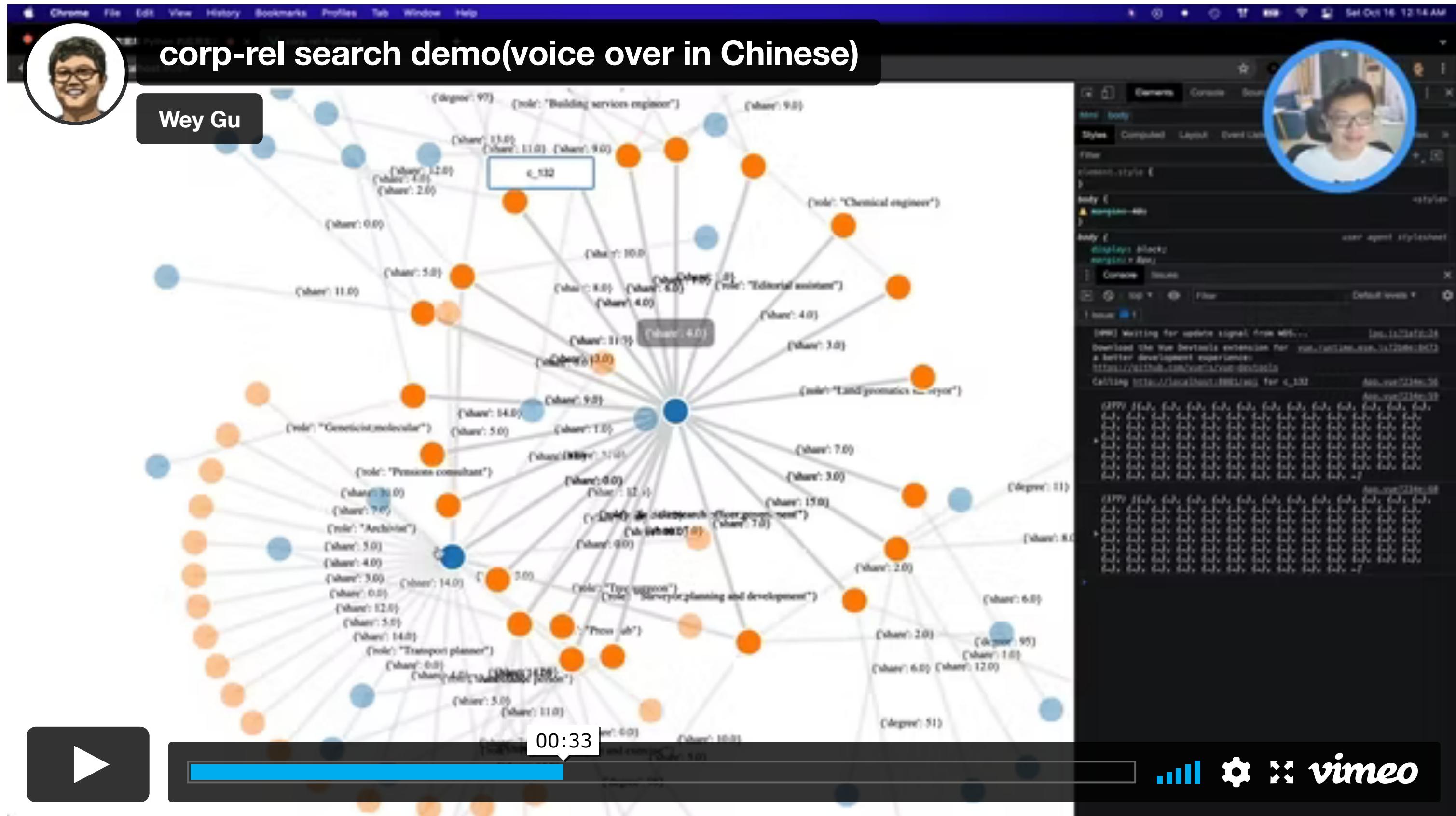
Graph Semantic Queries

```
SELECT a.id, a.name, c.name
FROM player a
JOIN serve b ON a.id=b.player_id
JOIN team c ON c.id=b.team_id
WHERE c.name IN (SELECT c.name
FROM player a
JOIN serve b ON a.id=b.player_id
JOIN team c ON c.id=b.team_id
WHERE a.name = 'Tim Duncan')
```

```
GO FROM 100 OVER serve YIELD serve._dst AS Team | \
GO FROM $-.Team OVER serve REVERSELY YIELD $$ .player.name;
```

Performance

	Designed Scenario	2-hop latency (~2.5K)	3-hop latency (~110K)	4-hop latency (~600K)
Graph DB	Relationship Walk	0.01 sec	0.168 sec	1.36 sec
SQL DB	Information retrieval	0.016 sec	30 sec	1544 sec



github.com/wey-gu/nebula-corp-rel-search

iSQE

A screenshot of a Katacoda profile page for Wey Gu. The top navigation bar shows "siwi demo in katacoda(voice over in Chinese)". The profile section features a circular profile picture of Wey Gu, the name "Wey Gu", and the O'Reilly Katacoda logo. Below this is a large blue banner with the heading "Wey Gu, Scale the Magic to others!" and a sub-headline "@wey". A bio text states: "I am a developer @vesoft working as Developer Advocate of Nebula Graph, the open-source distributed Graph Database I create toolings and content for Nebula Graph Database to help Developers in the open-source community. I am working in open source and consider it is a privilege 1. It took me a couple of my early career years to figure out that my passion lies in helping others with my thoughts & the tech/magic I have learned." To the right is a "Share Your Success" section with "Share" buttons for LinkedIn and Twitter. Below the banner are two scenario cards: "Shareholding Ownership Analysis with Nebula Graph Database" and "Siwi the Knowledge Graph Dialog System with Nebula Graph". Each card has a "Start Scenario" button. At the bottom is a video player showing a play button, a progress bar at 05:53, and a Vimeo logo.

https://katacoda.com/wey/scenarios/siwi-kgqa

katacoda.com/wey/scenarios/siwi-kgqa

iSQE

挑战：



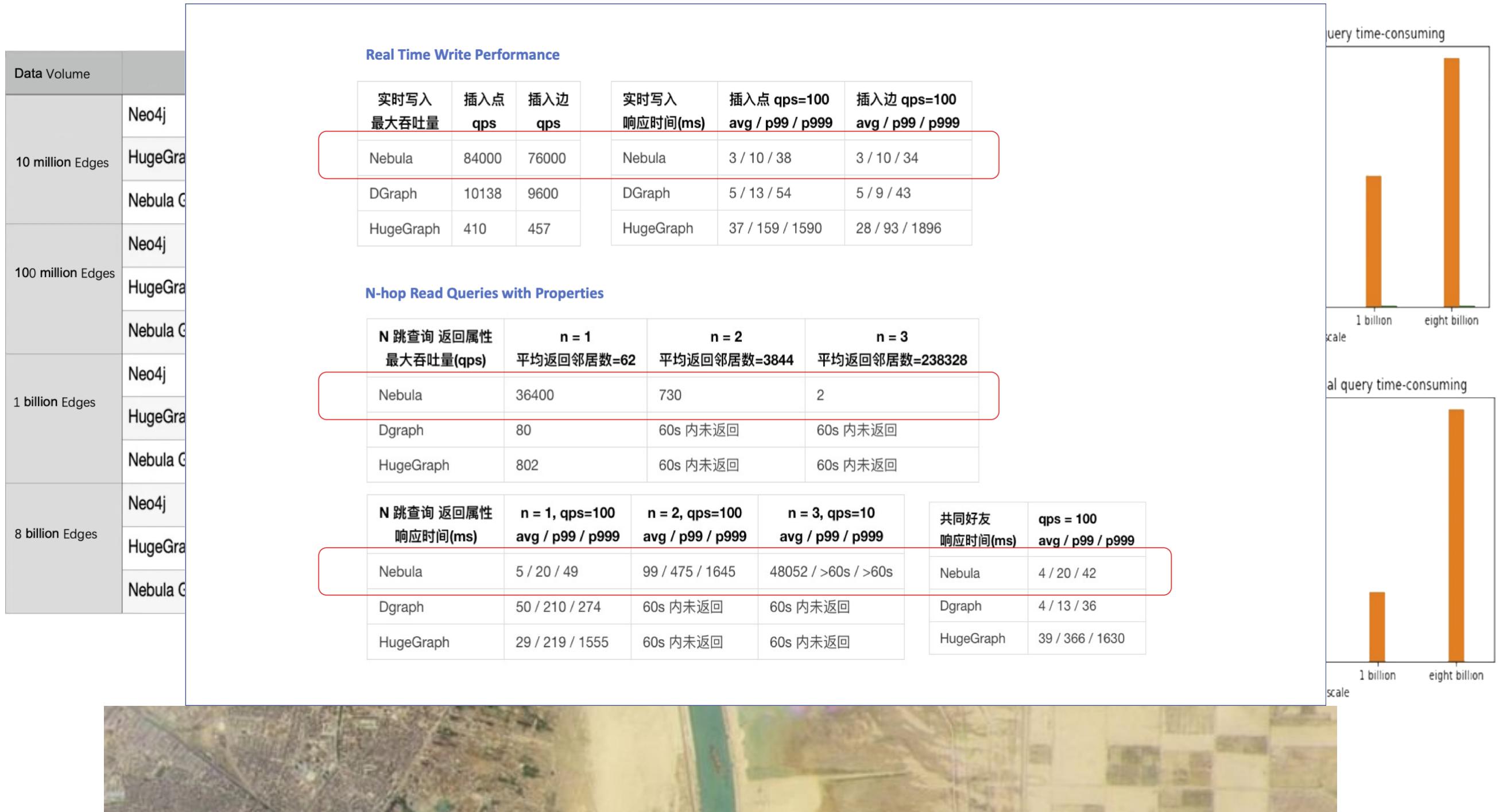
Graph DB nGQL



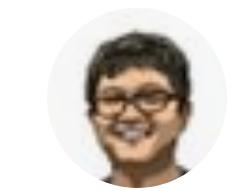
iSQE

挑战：数据规模

NebulaGraph is **highly performant yet linearly scalable** as it's designed **shared**



① nebula-graph.io/cases/



Wey Gu 古思为
@wey_gu · [Follow](#)



Scale makes differences
[@NebulaGraph](#)



Martin Beeby @thebeebs

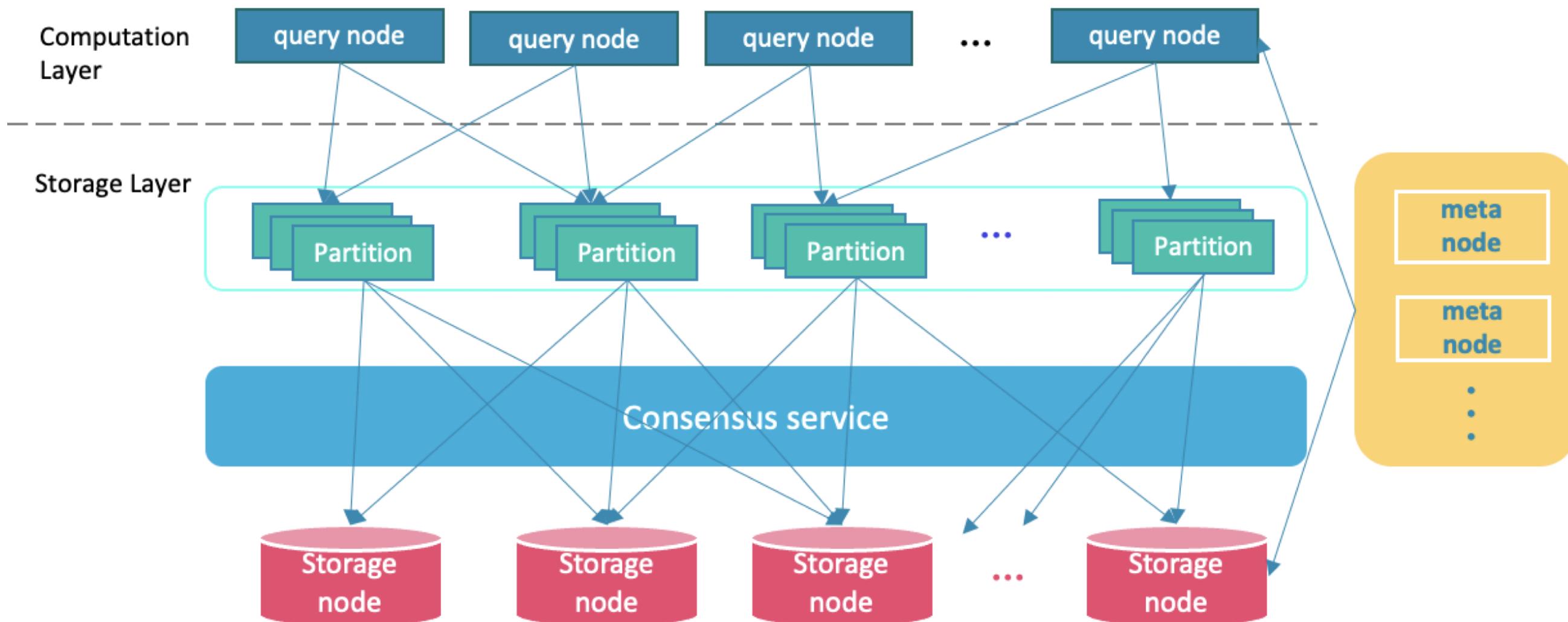
Small and medium businesses may have similar problems to big tech companies. But should never assume they need to use the same solution. The answer is not always a shiny new library, a complicated frontend build process, or the latest orchestrator. It is all a matter of scale.



iSQE

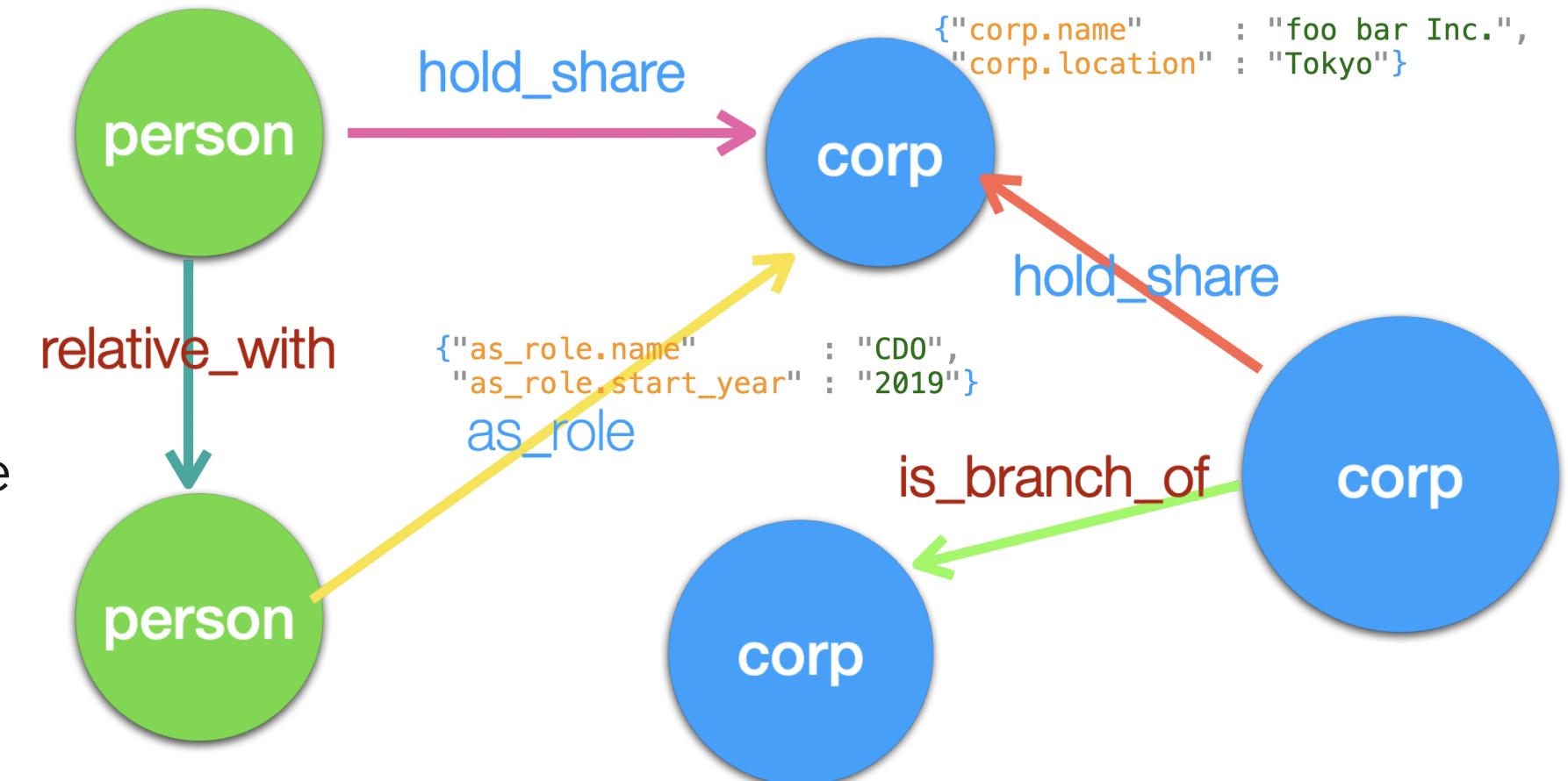
Distributed GraphDB Architecture

- Separated Compute and Storage: Scale Independently
- Shared-Nothing Architecture: Scale Easily and Fault Tolerance
- Auto Sharding(Hash Key Range, shards are way more than nodes)
- Raft(For both Meta and Storage nodes)



Graph Database Data Model

- Property Graph Data Model
- Schema:
 - Tag defines a set of properties for a vertex
 - Edge Type defines a set of properties for a edge
- Vertex
 - A vertex is a node in the graph, with a Vertex ID
 - 0, or more Tags could be attached to a vertex
- Edge
 - An edge is a 4-tuple of (Vertex ID, Edge Type, Vertex ID, Rank)
 - Edge connected two vertices in given type and rank



Graph Query Language

NGQL

Traversing graph

```
GO 3 Steps FROM "player102" OVER follow YIELD dst(edge);
```

Index based

```
LOOKUP ON player WHERE player.name == "Tony Parker" \
YIELD id(vertex);
```

Directly retrieve properties:

```
FETCH PROP ON player "player100" YIELD properties(vertex);
```

```
FIND SHORTEST PATH FROM "player102" TO "team204" OVER * \
YIELD path AS p;
```

```
GET SUBGRAPH 5 STEPS FROM "player101" \
YIELD VERTICES AS nodes, EDGES AS relationships;
```

Connected With Pipeline

```
GO FROM "player100" OVER follow \
YIELD dst(edge) AS did, properties($$).name AS Name | \
GO FROM $-.did OVER follow YIELD dst(edge);
```

OPENCYCER

```
MATCH (v:player{name:"Tim Duncan"})-->(v2)<--(v3) \
RETURN v3.player.name AS Name;
```

```
MATCH (v:player) \
WHERE not (v)--() \
RETURN v;
```

```
MATCH (v:player)--(v2) \
WHERE id(v2) IN ["player101", "player102"] \
RETURN v;
```

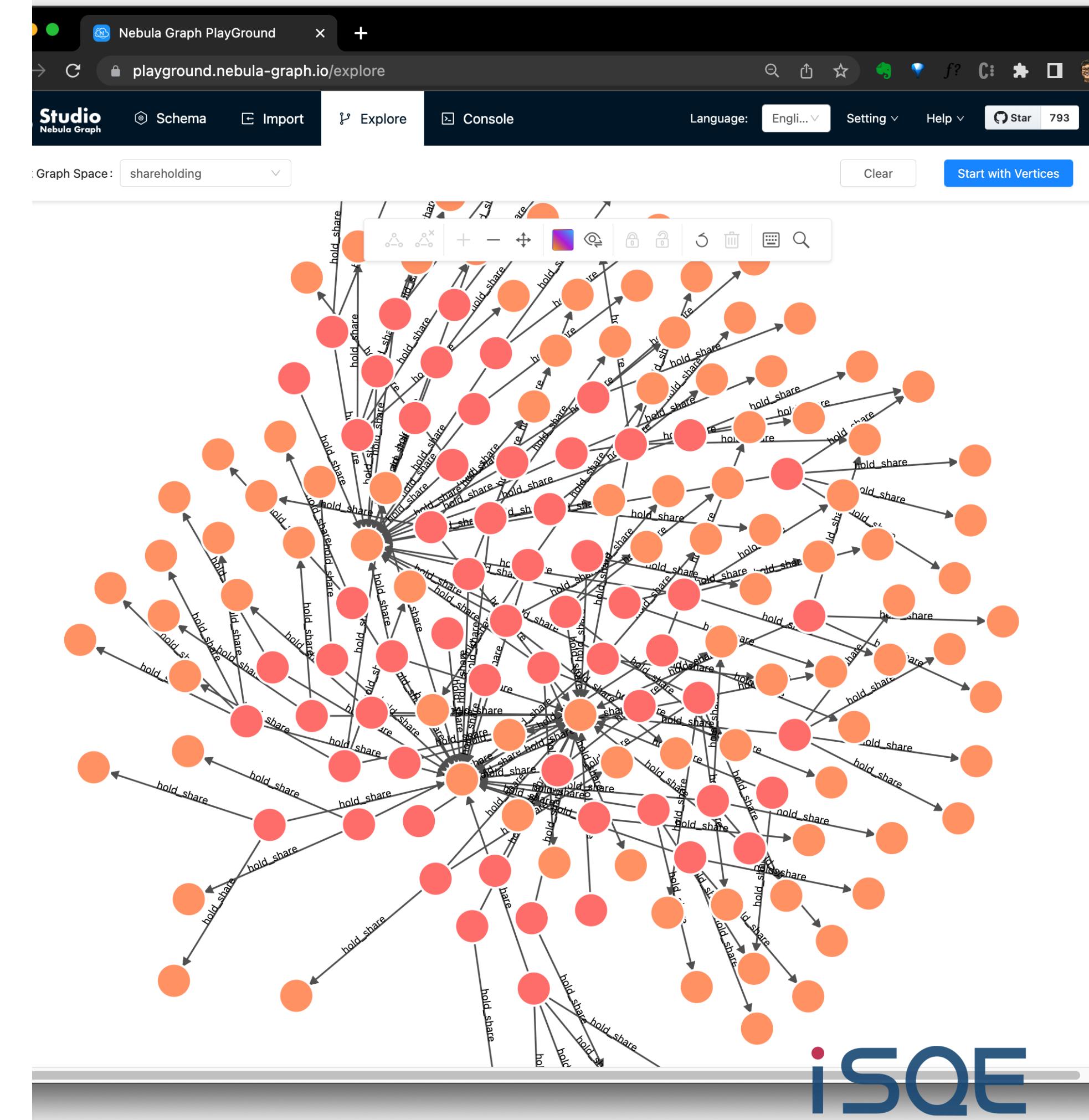
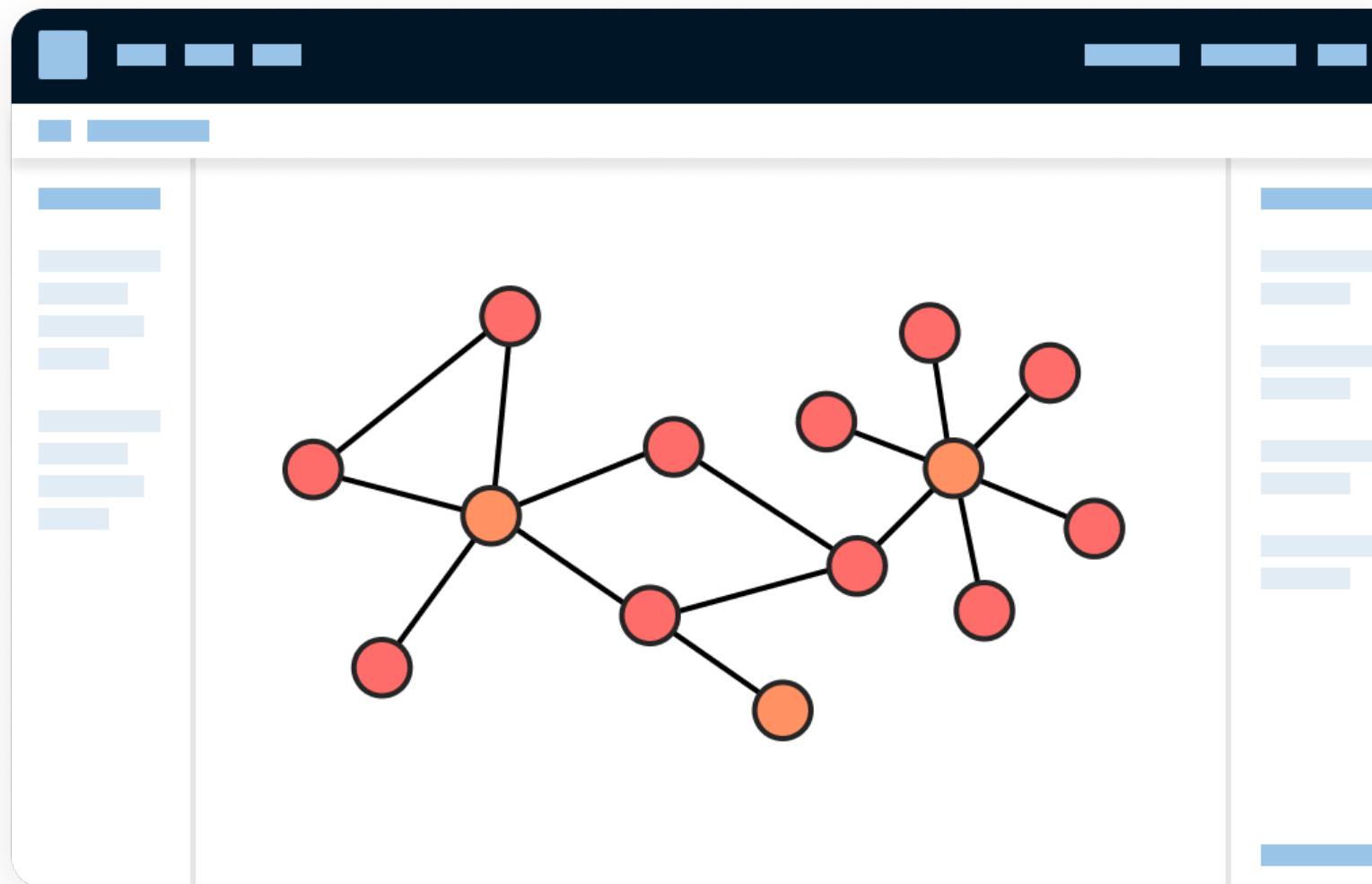
```
MATCH (m)-[]→(n) WHERE id(m)=="player100" \
OPTIONAL MATCH (n)-[]→(l) WHERE id(n)=="player125" \
RETURN id(m),id(n),id(l);
```

WILL BE FULLY COMPATIBLE WITH FUTURE ISO GQL STANDARD

优势：Data Visualization

Get one's 3-depth equity relationship with 1+% share

```
MATCH p=(:person{name: "Debra Ingram"})-[e*1..3]-(x)
WHERE ALL(e_ in e WHERE e_.share > 1.0)
RETURN p,e
```



iSQE

The screenshot shows the NebulaGraph Explorer Demo 4 interface. At the top, there's a navigation bar with tabs: Explorer, Visual Query (highlighted in red), and Workflow. Below the navigation bar, a sidebar on the left lists various components and features: MAT, WCC, CWC, CCW, Node importance, Graph feature, Community discovery, Clustering, and a Recent used components section. The main workspace contains a visual workflow editor. A 'Query' component is connected to a 'LoeVain' component. The 'Query' component has two outputs: 'output1' and 'output2'. The 'LoeVain' component has three inputs: 'src' (from 'output1'), 'dist' (from 'output2'), and 'weight' (from 'output1'). On the right side of the interface, there's a 'Query' panel with a search bar, an 'Input' section, a 'Query language' dropdown set to 'basketballplayer', and a 'Results' section. Below the input section, there's a code editor with a single line of Neo4j-like query: `match ()-[e:follow]->v return src(e) as s, dist(e) as d`. A play button and a progress bar at 00:01 are visible at the bottom of the video player area.

```
WITH CASE WHEN ALL(up_len IN upstream_len WHERE up_len IS NULL) THEN COLLECT(NULL)
ELSE COLLECT(DISTINCT {level:SIZE(upstream_len), source:split(id(upstream_entity),'://')[0],
key:id(upstream_entity), badges:upstream_badges, usage:upstream_read_count, parent:id(nodes(upath)[-2])})
END AS upstream_entities CASE WHEN ALL(down_len IN downstream_len WHERE down_len IS NULL) THEN collect(NULL)
```

iSQE

优势：SRE, Cloud Native

- Nebula K8s Operator
- Ansible/Docker-Compose
- Grafana
- Prometheus
- Stats Exporter
- Backup & Restore

Nebula enriches the CNCF Landscape

CLOUD NATIVE LANDSCAPE

You are viewing 1,147 cards with a total of 3,319,280 stars, market cap of \$19.6T and funding of \$53.8B.

Landscape Card Mode Members Serverless Wasm

100% +

Database Streaming & Messaging Application

App Definition and Development

The Cloud Native Computing Foundation Landscape visualization displays a grid of 1,147 cards representing various open-source projects. The projects are categorized into three main horizontal sections: Database, Streaming & Messaging, and Application. The Database section includes projects like KV, Vitess, CarbonData, Apache Hadoop, Ignite, ArangoDB, BIGCHAINDB, Cassandra, Cockroach Labs, Couchbase, Crate.io, crunchydata, Hazelcast IMDG, GraphScope, IBM DB2, Iguazio, Infinispan, InterSystems IRIS Data Platform, Kube, MariaDB, MySQL, NebulaGraph, Neo4j, Beam, CD Events, Apache Nifi, Apache Spark, Apache RocketMQ, Apache Storm, Azure Event Hubs, Beam, CD Events, Deepstream, EMQ, Flink, Fluvio, Google Cloud Dataflow, Hazelcast Jet, Kafka, KubeMQ, Lightbend, OpenMessaging, PostgreSQL, Presto, Qubole, SingleStore, Snowflake, ShardingSphere, TiDB, Timescale, VERTICA, VOLTDB, RabbitMQ, Redpanda, SeaTunnel, Siddhi, StreamSets, ServiceComb, Shipwright, Skaffold, Squash, and many others. The Application section includes Helm, Backstage, Amazon Kinesis, Heron, Bitnami, CARVEL, Eclipse Che, Gitpod, Gradle, KubeOrbit, KubeVela, Kui, Kudo, On-Prem, Open Application Model, OpenServiceBrokerAPI, and OPENAPI Initiative. The App Definition and Development section includes Nebula, CRUX, Databend, Dgraph, druid, FoundationDB, GraphScope, Hazelcast IMDG, IBM DB2, Iguazio, Infinispan, InterSystems IRIS Data Platform, Kube, MariaDB, MySQL, NebulaGraph, Neo4j, Beam, CD Events, Apache Nifi, Apache Spark, Apache RocketMQ, Apache Storm, Azure Event Hubs, Beam, CD Events, Deepstream, EMQ, Flink, Fluvio, Google Cloud Dataflow, Hazelcast Jet, Kafka, KubeMQ, Lightbend, OpenMessaging, PostgreSQL, Presto, Qubole, SingleStore, Snowflake, ShardingSphere, TiDB, Timescale, VERTICA, VOLTDB, RabbitMQ, Redpanda, SeaTunnel, Siddhi, StreamSets, ServiceComb, Shipwright, Skaffold, Squash, and many others. The interface also features navigation tabs (Landscape, Card Mode, Members, Serverless, Wasm), zoom controls (100%, +, -), and a large 'iSQE' watermark at the bottom right.

优势：数据管道

- Importer(`.csv`)/Console(`.ngql`)
- Exchange
 - Source: Postgres, Kafka, Clickhouse etc.
 - Target: NebulaGraph Server/SST File
- Connector/Clients
 - Spark/Flink
 - Python/Java/Go/CPP
 - NodeJS/PHP/.NET/JDBC
 - HTTP-Gateway
 - ORM: GO/Java*



What is Nebula Exchange

[Nebula Exchange](#) (Exchange) is an Apache Spark™ application for bulk migration of cluster data to Nebula Graph in a distributed environment, supporting batch and streaming data migration in a variety of formats.

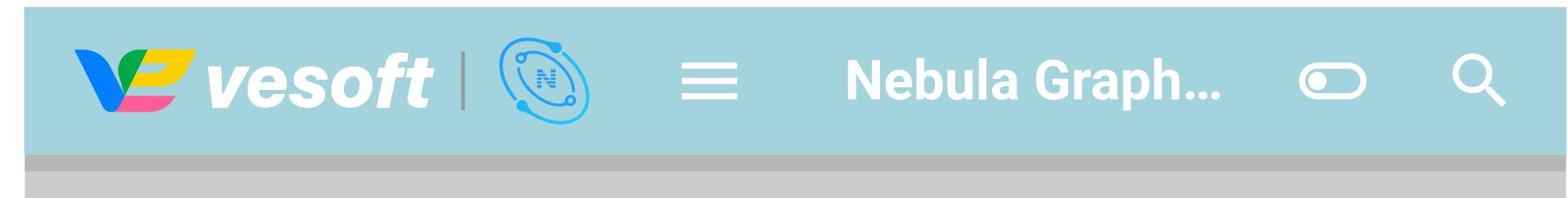
Exchange consists of Reader, Processor, and Writer. After Reader reads data from different sources and returns a DataFrame, the Processor iterates through each row of the DataFrame and obtains the corresponding value based on the mapping between `fields` in the configuration file. After iterating through the number of rows in the specified batch, Writer writes the captured data to the Nebula Graph at once. The following figure illustrates the process by which Exchange completes the data conversion

Graph Algorithm

Open Source state-of-the-art Graph Algorithm Platform on Spark consuming NebulaGraph.

- PageRank/ Louvain/ KCore
- LabelPropagation/ Hanp
- ConnectedComponent/ Strongly ConnectedComponent
- ShortestPath/ BFS
- TriangleCount/ GraphTriangleCount
- BetweennessCentrality/ Closeness
- DegreeStatic
- ClusteringCoefficient
- Jaccard
- Node2Vec

Nebula Analytics



Nebula Analytics

Nebula Analytics is a high-performance graph computing framework tool that performs graph analysis of data in the Nebula Graph database.

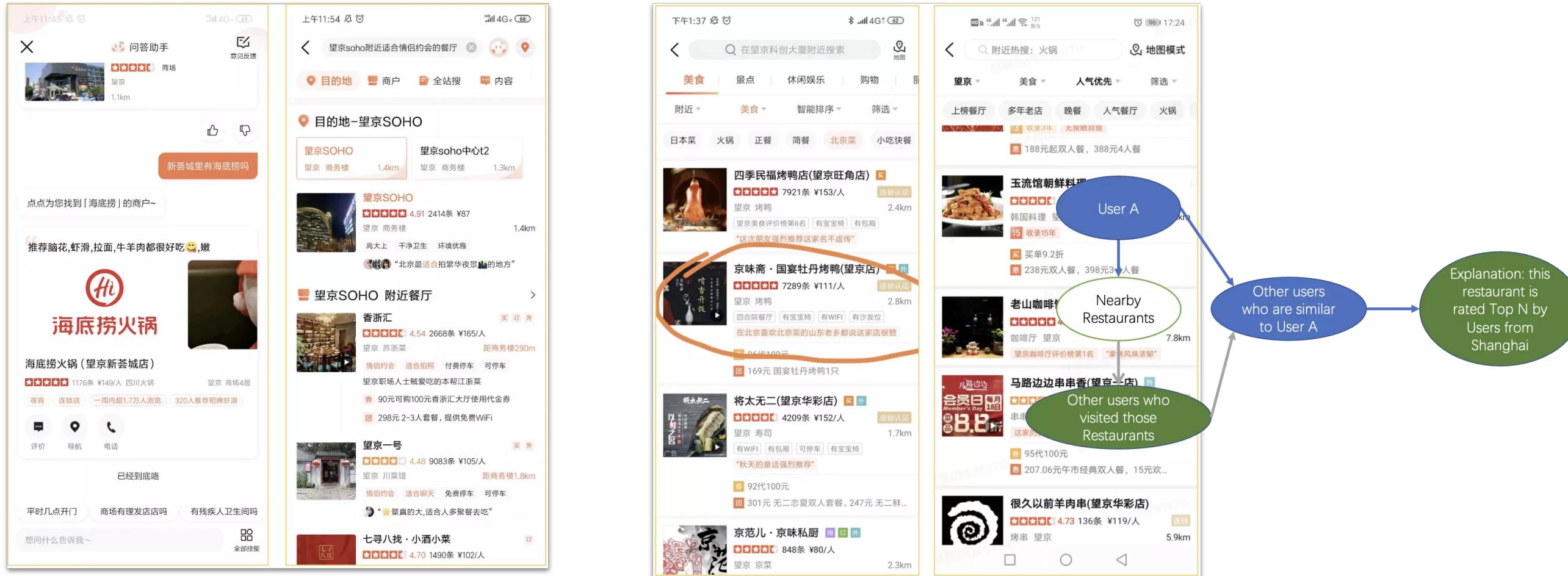
ⓘ Enterpriseonly

Only available for the Nebula Graph Enterprise Edition.

Scenarios

广泛的应用

- " Is there any McDonald's nearby street Foofar?
- " Those who from Beijing enjoying the Korean cuisine said this restaurant is great.



SNS

Risk Control

Public Security

Knowledge Graph

ML

Biopharmacy

IoT

Blockchain

Data Lineage

AI Ops

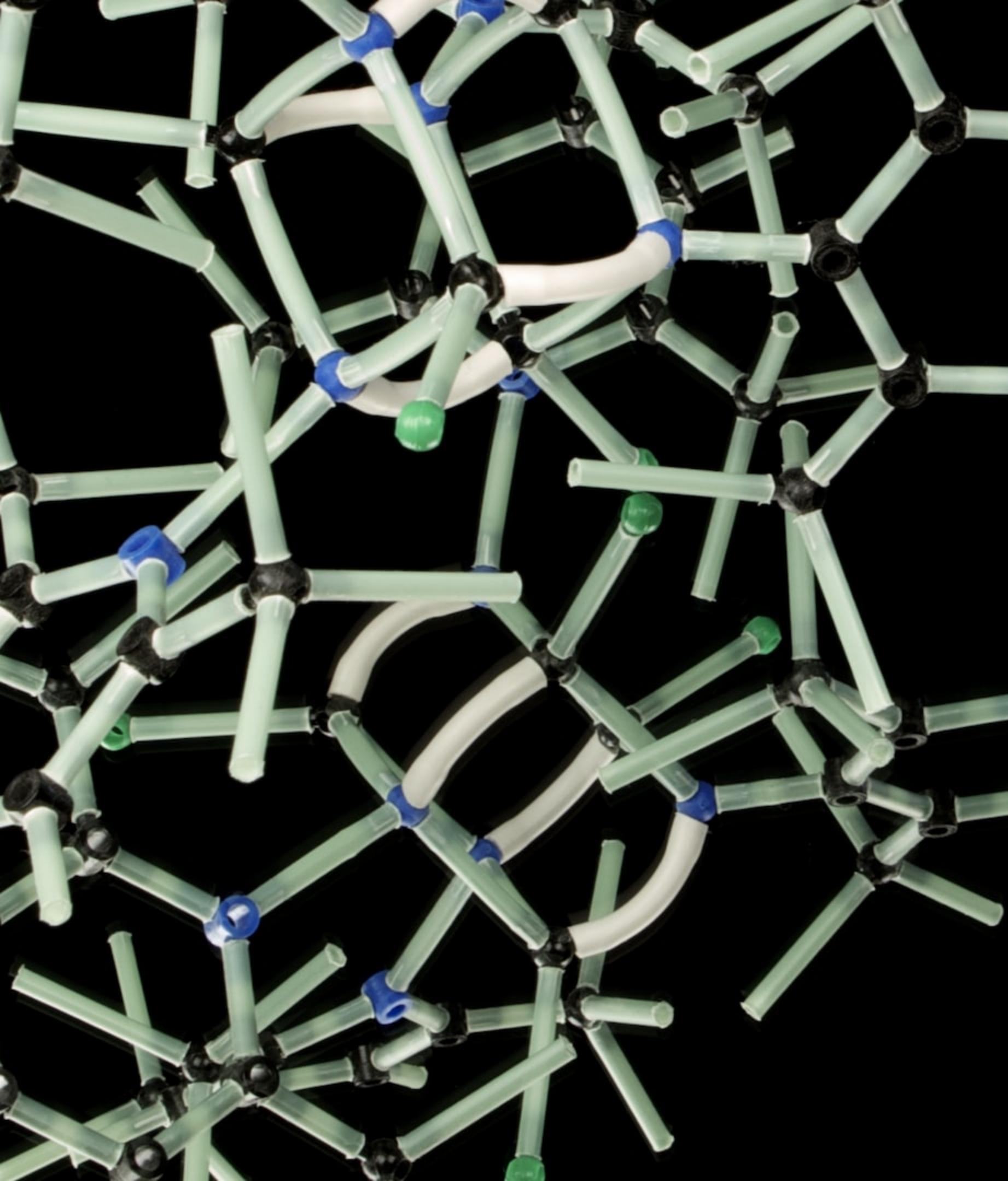
iSQE

数据质量工程

Accessibility, Correctness, Consistency, Uniqueness, Reasonableness, Data integrity

Data Governance is used to ensure capability that enables an organization to ensure that high data quality exists throughout the complete lifecycle of the data, and data controls are implemented that support business objectives.

① [wikipedia.org/wiki/Data_quality](https://en.wikipedia.org/wiki/Data_quality)



Data Governance

Data Quality Engineering is about:
Engineering process & tooling to build trust of
data, system and team, where the critical parts
are:

- Data Ownership
 - Shall we send mails to all or shout in WeChat groups?
- Metadata
 - Manage and discover data schema.
 - Maintain documentation of the schema
- Data Lineage
 - Traceability of data, service, pipelines etc.



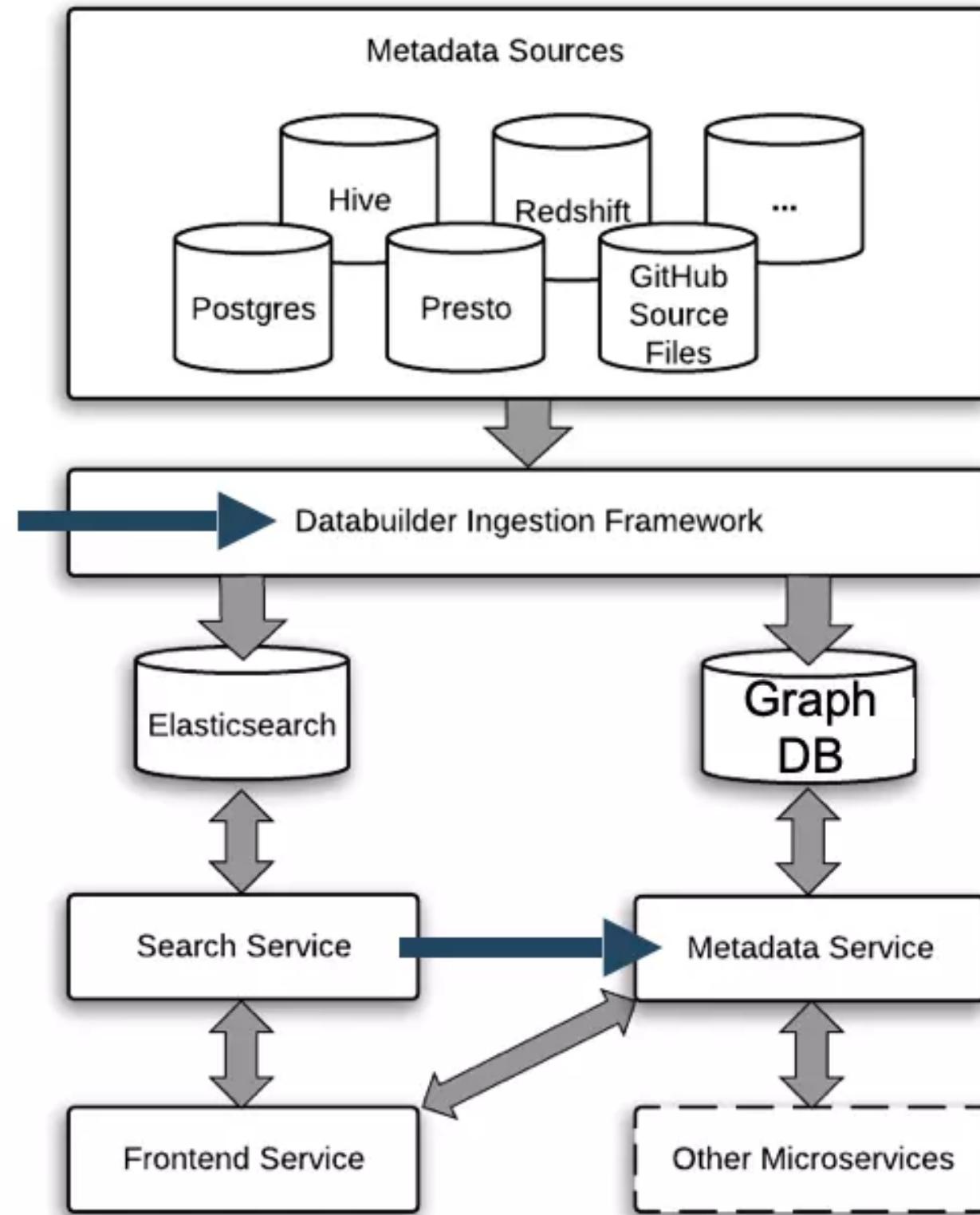
Open source data discovery and metadata engine



github.com/amundsen-io/amundsen

iSQE

Arch



Code

```
.
├── common
│   └── amundsen_common
└── databuilder
    ├── clients
    ├── extractor
    ├── loader
    ├── models
    ├── publisher
    ├── rest_api
    └── utils
├── docs
└── frontend
    └── amundsen_application
└── metadata
    └── metadata_service
        ├── api
        ├── cli
        ├── client
        ├── proxy
        └── util.py
└── search
    └── search_service
```



Chrome File Edit View History Bookmarks Profiles Tab Window Help

26 Sun Not Secure | 192.168.8.127:5000 Sun Dec 26 3:02 AM

AMUNDSEN Browse

Search for data resources... Advanced Search >

Available Badges

(Beta) (Fk) (Json) (Npl) (Pd) (Pk)

Popular Tags

cheap 1 delta 1 expensive 1 low_quality 1 needs_documentation 2 recommended 1 tag1 2 tag2 1

Browse all tags

My Bookmarks

Datasets (0)

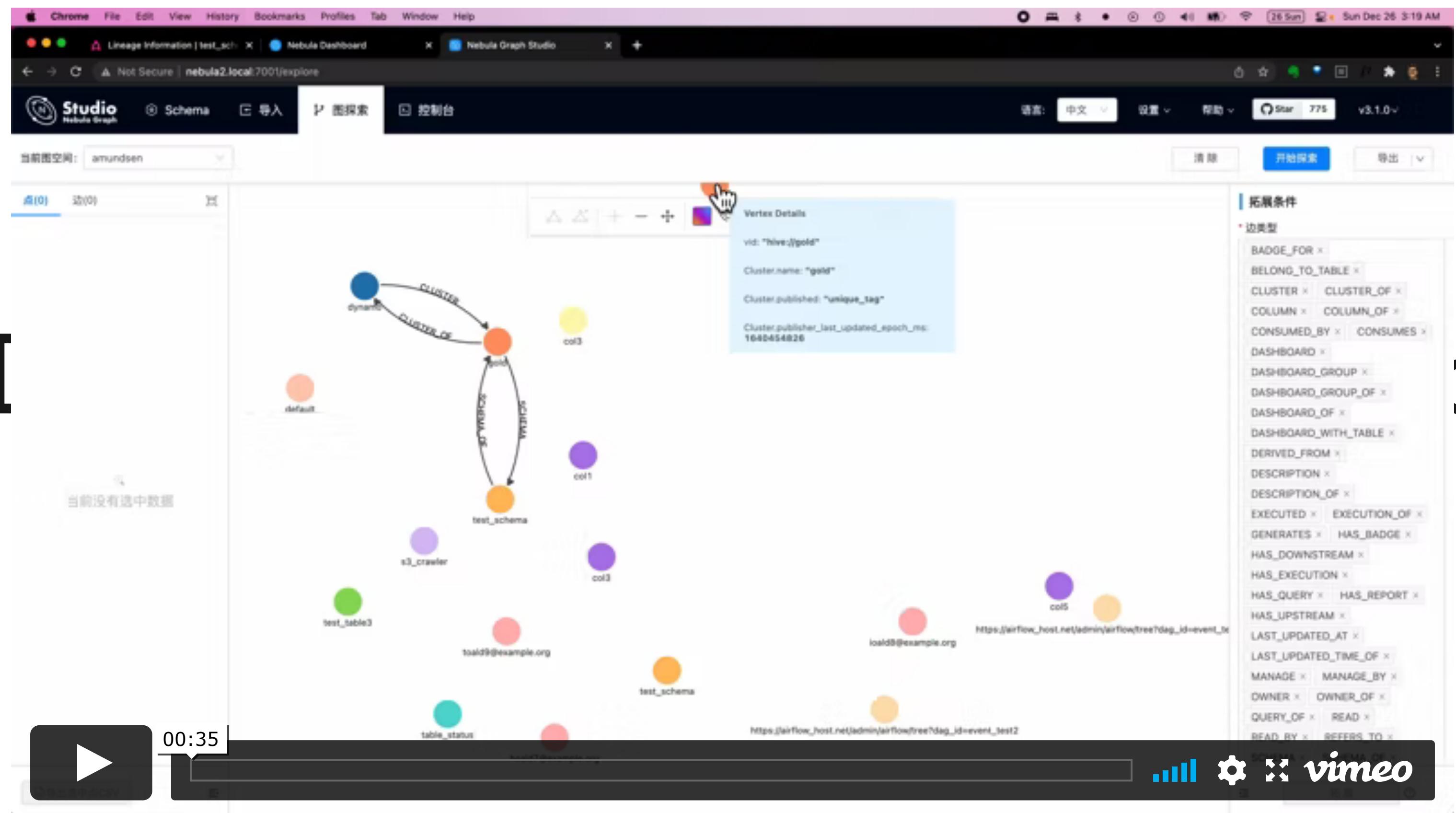
You don't have any bookmarks. Use the star icon to save a bookmark.

Popular Resources ⓘ

00:44 Datasets (0)

Amundsen was last indexed on December 25th 2021 at 11:50:00 am

vimeo



iSQE

D

```
Alacritty
~/.dev/nebula-console
~/dev/nebula-console remotes/origin/master~1*
> bat ~/Downloads/demo.ngql

File: /Users/weyl/Downloads/demo.ngql

1 USE amundsen;
2 MATCH (tbl:Table) \
3   WHERE id(tbl) == "hive://gold.test_schema/test_table1" \
4     OPTIONAL MATCH (wmk:Watermark)-[:BELONG_TO_TABLE]-(tbl) \
5     OPTIONAL MATCH (app_producer:Application)-[:GENERATES]-(tbl) \
6     OPTIONAL MATCH (app_consumer:Application)-[:CONSUMES]-(tbl) \
7     OPTIONAL MATCH (tbl)-[:LAST_UPDATED_AT]-(t:'Timestamp') \
8     OPTIONAL MATCH (owner:'User')<-[:OWNER]-(tbl) \
9     OPTIONAL MATCH (tbl)-[:TAGGED_BY]->(`tag`:'Tag') \
10    OPTIONAL MATCH (tbl)-[:HAS_BADGE]->(badge:Badge) \
11    OPTIONAL MATCH (tbl)-[:SOURCE]->(src:Source) \
12    OPTIONAL MATCH (tbl)-[:DESCRIPTION]->(prog_descriptions:Programmatic_Description) \
13    OPTIONAL MATCH (tbl)-[:HAS_REPORT]->(resource_reports:Report) \
14    RETURN collect(distinct wmk) as wmk_records, \
15      collect(distinct app_producer) as producing_apps, \
16      collect(distinct app_consumer) as consuming_apps, \
17      t.last_updated_timestamp as last_updated_timestamp, \
18      collect(distinct owner) as owner_records, \
19      collect(distinct `tag`) as tag_records, \
20      collect(distinct badge) as badge_records, \
21      src, \
22      collect(distinct prog_descriptions) as prog_descriptions, \
23      collect(distinct resource_reports) as resource_reports

~/dev/nebula-console remotes/origin/master~1*
>
```



00:09



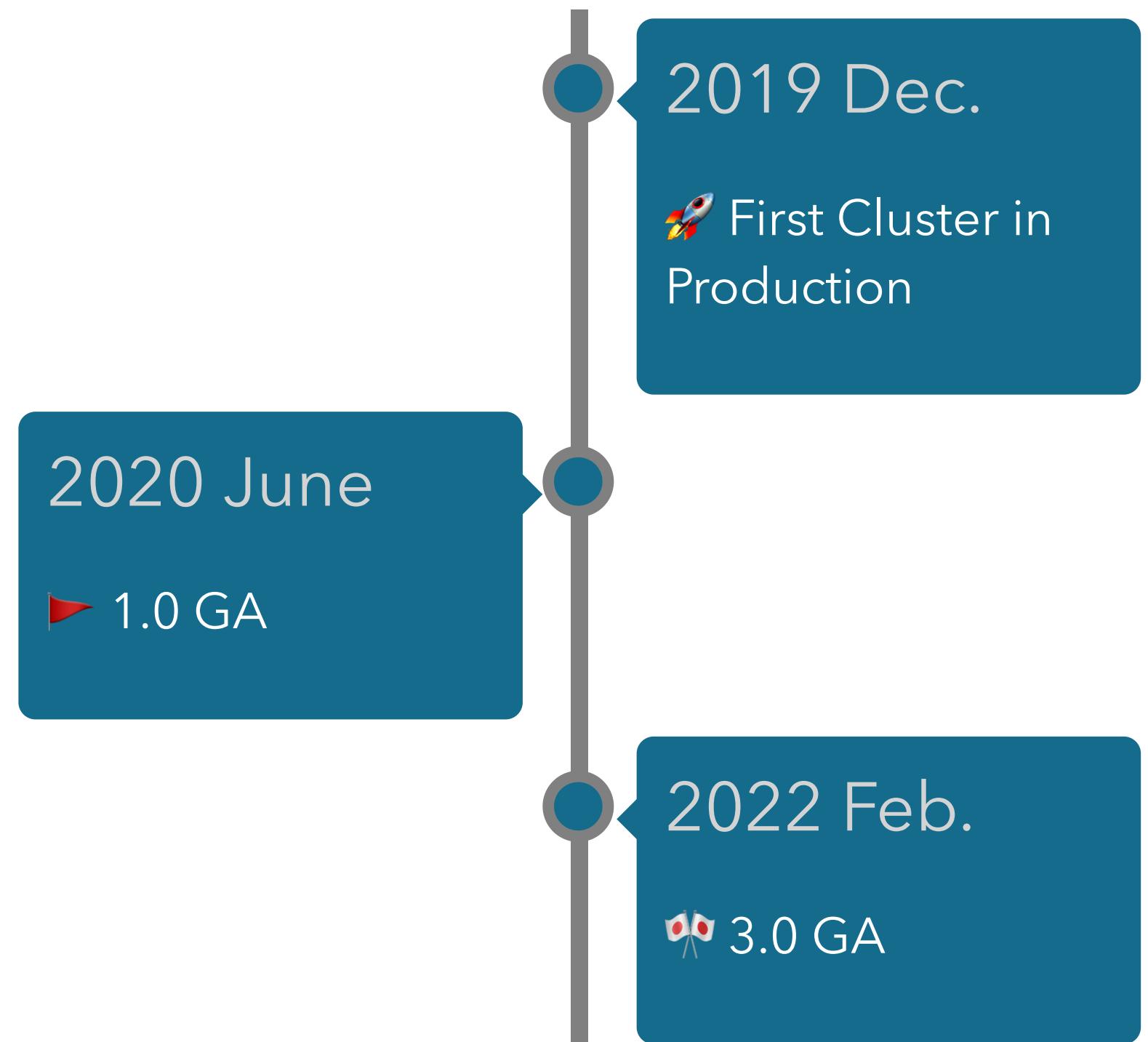
iSQE

NebulaGraph + Amundsen

👉 github.com/amundsen-io/amundsen/issues/1816

Open Source

Power users and contributors kept NebulaGraph evolving to the next level in many industries.

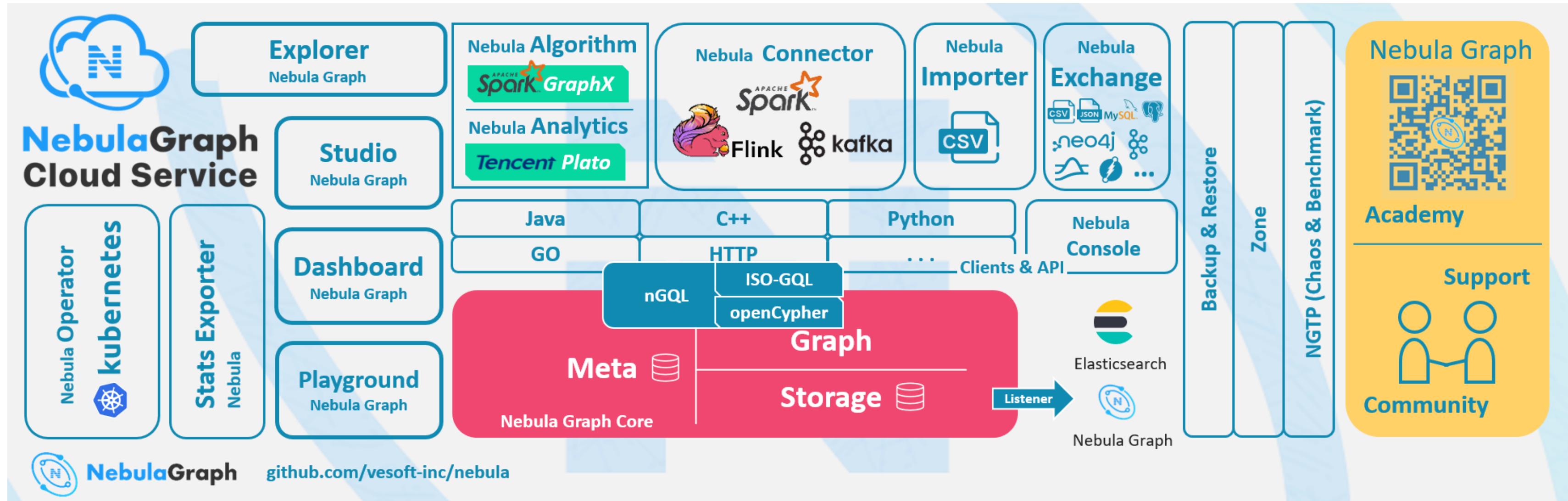


- ① nebula-graph.io/cases/
- ① db-engines.com/en/ranking/graph+dbms

NebulaGraph Ecosystem

Nebula Community is rich in ecology and still expanding and exploring, welcome to join and contribute!

- Deployment, Monitoring
- Data Visualization
- Algorithm, Analytic
- Clients, Connectors, ETL



Thank YOU!

<https://github.com/vesoft-inc/nebula>

