



Neptune Graph Database

0 to Production



Ohad Israeli
Platform Architect @ NI

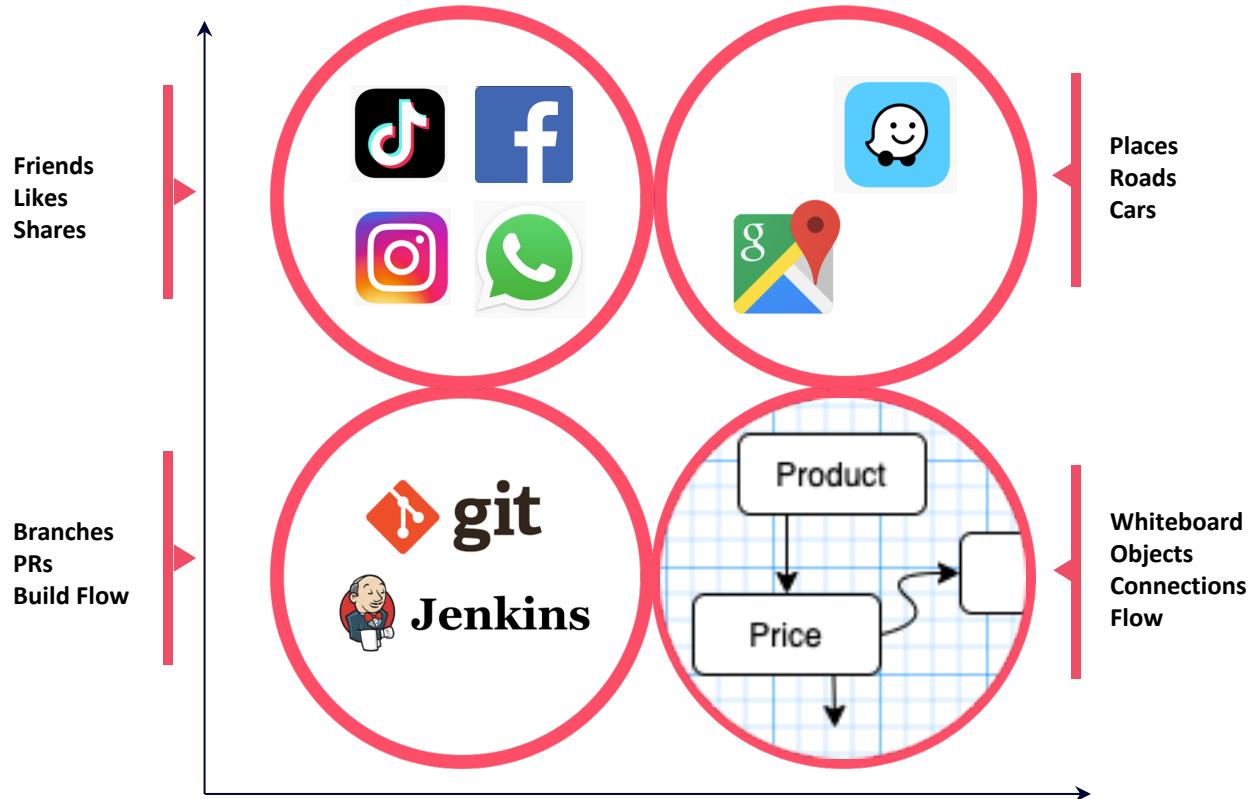
✉️ Ohad.Israeli@naturalint.com

Geektime
<CODE>



Jun 2021

Graphs are all around us





Ohad Israeli

Platform Architect @ NI

20+ Years of Software Architecture



What are we going to talk about today?

Sharing my experience of mapping
entities @ NI using AWS **Neptune** / **Kafka** / **NestJS**

90% 5% 5%

1. **Why** graph database ?
2. **How** did we use it ?
3. **What** we learned ?
4. **What** you need to know ?



setTimeout(presentationEnd, 30min, 'lets start');

Natural Intelligence in a Nutshell



Household



Shopping



Tech



Lifestyle



Family



Health & Wellness



Finance



Entertainment



Personal Growth



www.naturalint.com/jobs Full Stack, Frontend, Data Scientists, QA Lead, R&D Directors, Senior DevOps ...

Natural Intelligence in a Nutshell

6M

Leads generated
annually

9M

Monthly
Unique Visitors

130+

Countries
Reached

10

Years of
Expertise

320+

Talented
Employees



Natural Intelligence R&D in a Nutshell



NodeJS / NestJS

Serving platform + Backend
Microservices

KAFKA / Kinesis

Event Bus

meetup.com/CodeNaturally



www.naturalint.com/jobs Full Stack, Frontend, Data Scientists, QA Lead, R&D Directors, Senior DevOps ...

Java / Python

Data science / Machine learning

AWS

Fully invested



OK, let's get started

Chapter 1

Graph Data What ?



Databases Types

Relational Databases

Sql server

Oracle

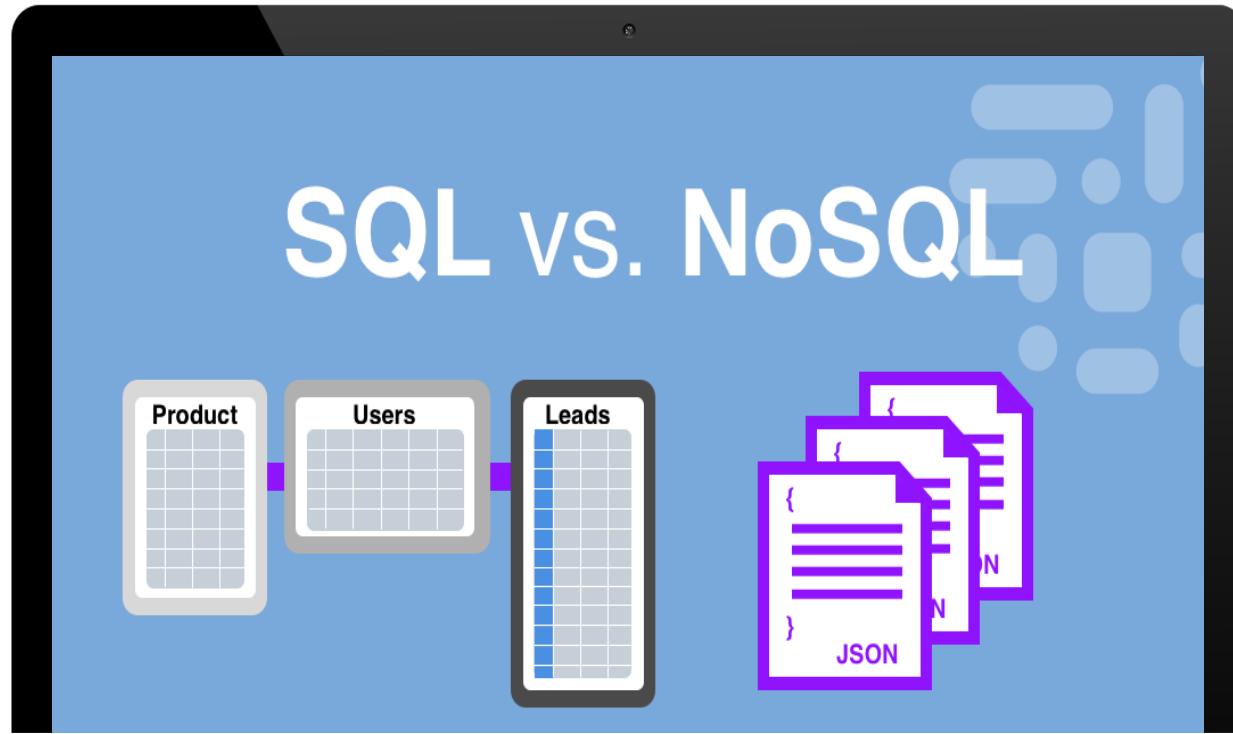
DB2

MySQL

Postgresql

...

...



Document Databases

Aerospike

CouchDB

MongoDB

Elasticsearch

Solr

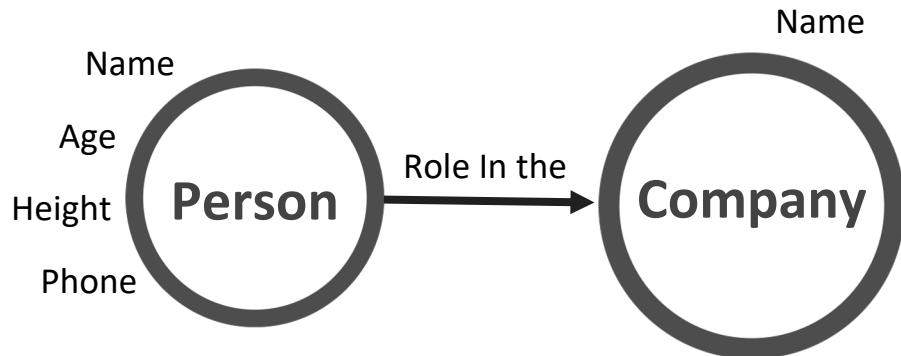
...

...



What is Graph Database ?

- No pre-defined model
- Data modeled like a drawing on a whiteboard
- Relationships between data as equally important to the data itself



Graph Databases

Graph Databases Query Language

AWS Neptune

NEO4J

DataStax

JanusGraph

OrientDB

ArangoDB

Dgraph

NabulaGraph

...

...

Gremlin

SparQL

Cypher

nGQL

Proprietary

...

...

Gremlin

```
g.V(person).valueMap();
```

SparQL

```
select ?p where {?p}
```

Cypher

```
MATCH (n:character {name:"person"}) RETURN properties(n)
```

nGQL

```
FETCH PROP ON character hash("person");
```





Coding time start in ...



SCAN ME



Amazon Neptune

[github.com / israelio / graphdb-intro](https://github.com/israelio/graphdb-intro)

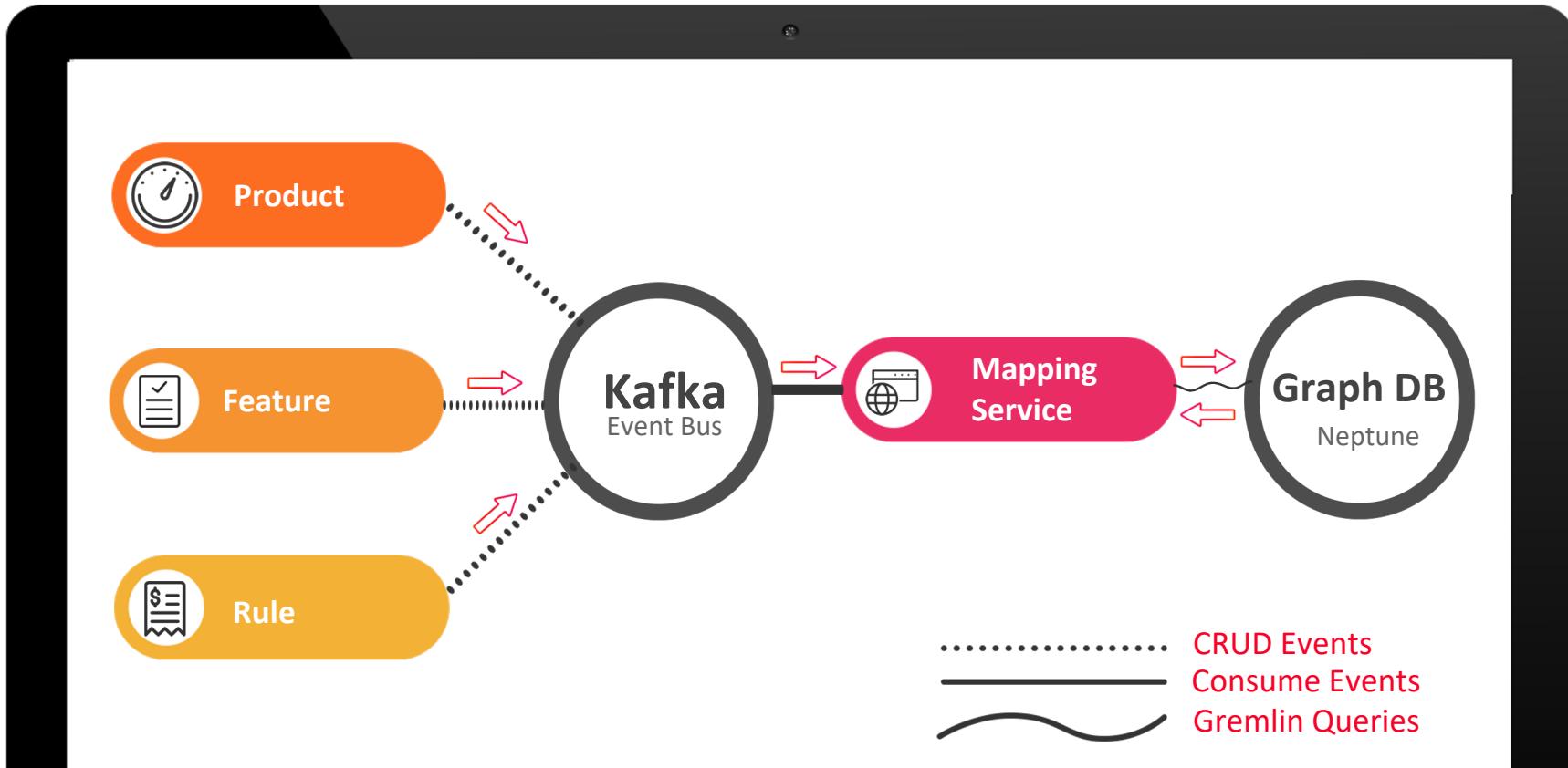


Chapter 2

How did we use it ?



Entity relations mapping solution architecture



Chapter 3

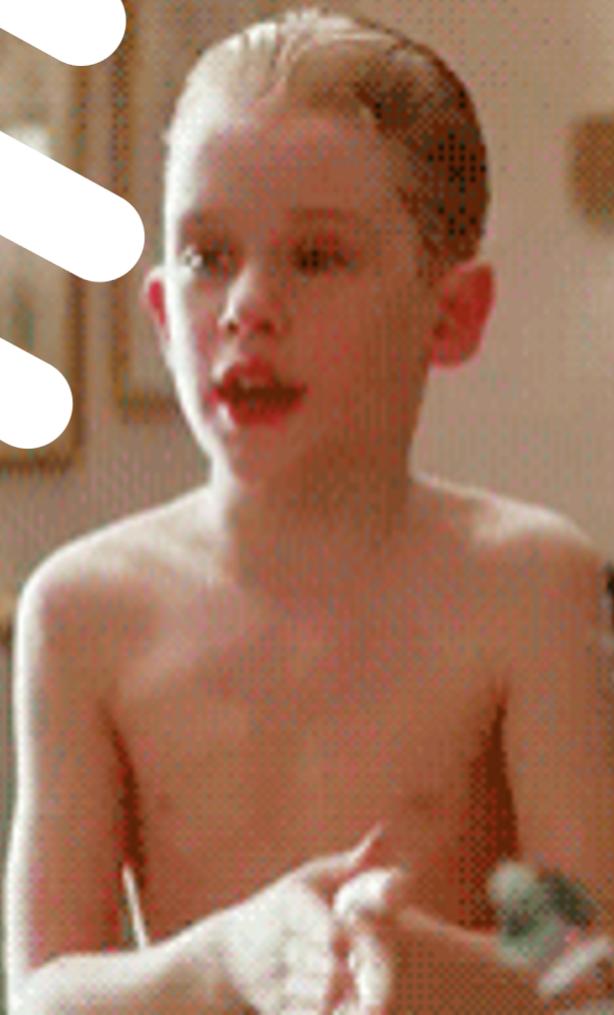
What did we learn ???

What do you need to know ?

DEVELOPMENT LIFECYCLE

Let's get into it

- Development
- Data import
- Data export
- Testing
- Deployment



Mapping Service



- NodeJS MicroServices Framework (nestjs.com)
 - Extensible
 - Promote Design patterns
 - Reusable services
 - Lots of integrations
 - Typescript enabled
-
- KafkaJS integration (already implemented – we extended it a bit)
 - AWS Neptune repository service (we wrote it internally)

NestJS + Kafka



```
npm i -g @nestjs/cli
```



```
nest new project-name
```

- Use **@EventPattern** to handle Kafka events



```
@EventPattern(topicName)
async handleEntityChangedEvent(payload: EntityRelationPayloadDto) {
    this.logger.info(`Event: ${topicName}, Payload: ${JSON.stringify(payload)}`);
    return this.statsDService.timer(
        () => this.handleEntityChanged(payload).catch((err) => this.getEntityProcessError(err, payload)),
        'handleKafkaEntityChanged_execution_time'
    );
}
```

Javascript + Gremlin



```
import * as gremlin from 'gremlin'
```

```
const {
  t: { id },
} = gremlin.process;
const {
  cardinality: { single },
} = gremlin.process;

const __ = gremlin.process.statics;
const __operators = gremlin.process.P;
```

- Use retry mechanism
- Use direct references
- Use UpSert pattern
- Add singularity

```
@Retryable({
  maxAttempts: RETRY_MAX_ATTEMPTS,
  backOffPolicy: BackOffPolicy.ExponentialBackOffPolicy,
  backOff: RETRY_BACKOFF,
  exponentialOption: { maxInterval: RETRY_EXP_MAX_INTERVAL, multiplier: RETRY_EXP_MULTIPLIER },
  doRetry: (e: Error) =>
    e.message.includes('ConcurrentModificationException') || e.message.includes('ReadOnlyViolationException'),
})
async addEntity(entity: Entity) {
  return this.graphs.writer
    .V(entity.nodeId)
    .fold()
    .coalesce(
      __.unfold()
        .property(single, ENTITY_FIELDS.LAST_UPDATE, new Date().toUTCString())
        .property(single, ENTITY_FIELDS.ATTRIBUTES, JSON.stringify(entity.filteredAttributes)),
      __.addV(entity.entityType)
        .property(id, entity.nodeId)
        .property(ENTITY_FIELDS.NODE_ID, entity.nodeId)
        .property(ENTITY_FIELDS.ENTITY_ID, entity.entityId)
        .property(ENTITY_FIELDS.ENTITY_TYPE, entity.entityType)
        .property(ENTITY_FIELDS.VAR_ID, entity.varId)
        .property(ENTITY_FIELDS.VAR_TYPE, entity.varType)
        .property(single, ENTITY_FIELDS.IS_PUBLISHED, !!entity.isPublished)
        .property(single, ENTITY_FIELDS.LAST_UPDATE, new Date().toUTCString())
        .property(single, ENTITY_FIELDS.VERSION, this.VERTEX_VERSION)
        .property(single, ENTITY_FIELDS.ATTRIBUTES, JSON.stringify(entity.filteredAttributes))
    )
    .id()
    .next();
}
```

1
2

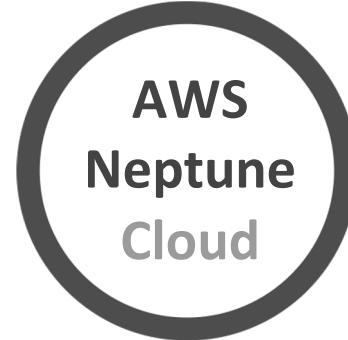
3

4

5



Cost explorer



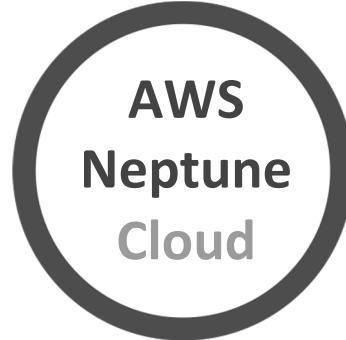
Production / Staging

- Price per instance / instance type



Development / Testing

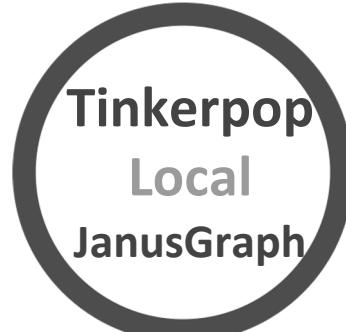
Cost explorer



Production / Staging

- Price per instance / instance type

- In memory
- Large database is an issue in load time and memory consumption

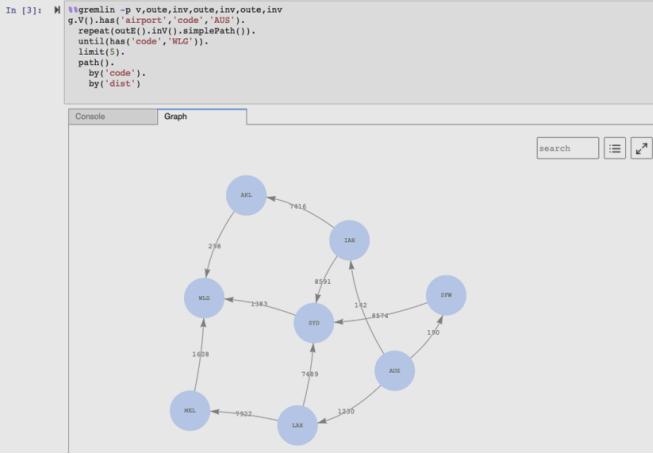


Development / Testing

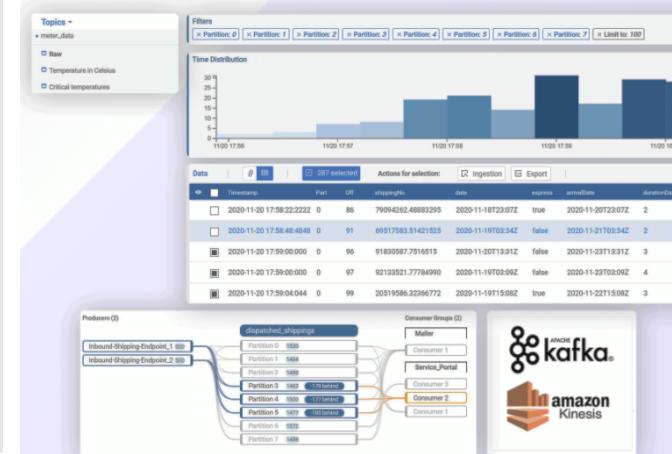
tinkerpop.apache.org / janusgraph.org

Dev/Test - Tinkerpop + Kafka

- Use dockers to host Tinkerpop and Kafka locally
- Configure Tinkerpop server yml to use persistency / Use JanusGraph
- Add all dockers to the same network
- Add dependencies between the dockers
- Add **Xeotek KaDeckWeb** docker – it will assist you in dev / test / prod with kafka
- Add **AWS graph notebook docker** - github.com/BarrMan/graph-notebook-docker



```
In [3]: g.tinkerPop->p.v.outE('inv_outE').inv_outE,inv
g.V().has('Airport','code','AUS'),
repeat(outE().inV(),simplePath()),
until(has('code','WLG')),
limit(5),
path,
by('code'),
by('dist')
```



Topics - **meter_data**

Filters: raw, Temperature in Celsius, Critical temperatures

Time Distribution

Timestamp	Part	Offset	Actions for selection:	Ingestion	Export		
2020-11-20 17:58:22.2222	0	86	7909426.48883295	2020-11-18T23:07Z	true	2020-11-29T23:07Z	2
2020-11-20 17:58:48.4848	0	91	69517582.51421325	2020-11-19T03:34Z	false	2020-11-21T03:34Z	2
2020-11-20 17:59:00.0000	0	96	91830592.79160115	2020-11-20T13:31Z	false	2020-11-23T13:31Z	3
2020-11-20 17:59:00.0000	0	97	92139521.77784990	2020-11-19T07:05Z	false	2020-11-23T07:05Z	4
2020-11-20 17:59:04.9444	0	99	20519586.32366772	2020-11-19T15:08Z	true	2020-11-22T15:08Z	3

Producers (2)

Inbound-Shipping-Endpoint.1:999
Inbound-Shipping-Endpoint.2:999

dispatched_shipping

Consumer Groups (2)

Mailer
Consumer 1
Service_Portal
Consumer 3
Consumer 2
Consumer 1

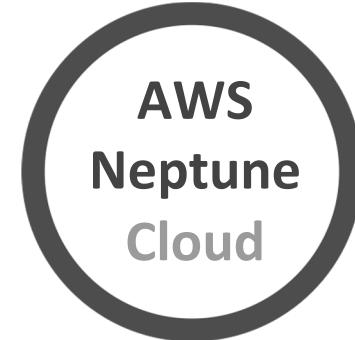
kafka
amazon Kinesis

Import / Export



Import data

- AWS Bulk data loader
- Write your own code / scripts
Java/NodeJS/Python



Export data

- **Neptune-Export.jar**
github.com/awslabs/amazon-neptune-tools

Production / Staging

- Write your own code / scripts
Java/NodeJS/Python



- Configure persistency

Development / Testing

Export



Export data

- **Neptune-Export.jar**

github.com/awslabs/amazon-neptune-tools

- **java -jar neptune-export.jar export-pg -e db-name.neptune.amazonaws.com**

-d ./Neptune

← Dump to local folder

--use-ssl

--serializer GRAPHBINARY_V1D0

← Compatibility with tinkerpop

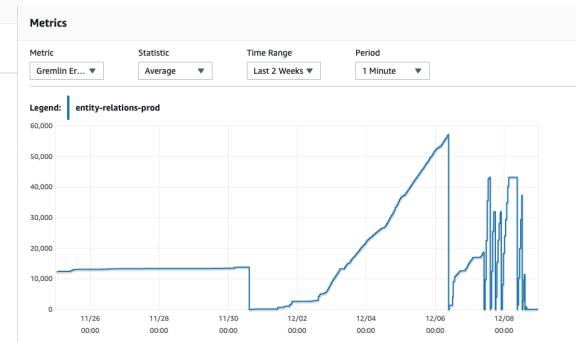
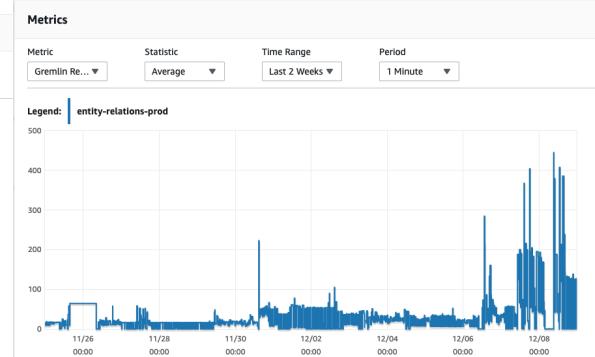
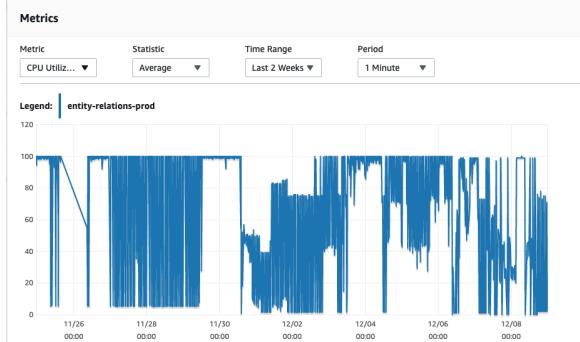
--clone-cluster

← Create new export cluster

--clone-cluster-instance-type db.r5.larg

Production

- Use AWS support if you need – they are great !
- We had performance issue with UpSert
 - It took 4.5 sec to complete
 - we got hotfix from AWS Support and it had been fixed in a later version
- We had high CPU for low numbers of gremlin transactions even on huge instances
 - Use direct references to nodes / edges when you can
- Do not forget to add retry on concurrency failures
- Monitor your queries, server load and gremlin errors



Summary

- Don't be afraid from using Graph Databases
- Use them only when / where it makes sense
 - (3 – 4 joins | traversal path)
- Working with NodeJS ? - Give NestJS a try
- Event sourcing is great for system decoupling
 - Kafka is awesome !

More info about this topic will be posted at:
medium.com/natural-intelligence



thank you

Ohad Israeli
Platform Architect
Natural Intelligence

ohad.israeli@naturalint.com



SCAN ME



<https://github.com/israelio/graphdb-intro>