# Amazon Neptune

**Ian Robinson, DBS Customer Advisory, EMEA | 2018**

**ianrob@amazon.com**

@awscloud

#AWSInnovate

# Amazon Neptune – A Fast, Reliable Graph Database



Optimized for storing and querying highly connected data
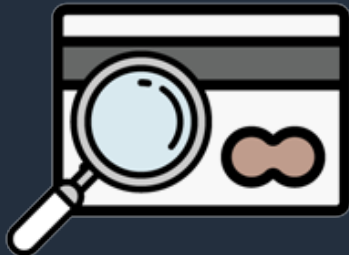
aws

# Neptune Use Cases

**Social Networking**

**Recommendations**

**Knowledge Graphs**

**Fraud Detection**

**Life Sciences**

**Network & IT Operations**

aws

# Characteristics of a Graph Workload

- Complex domain model
- Variable schema (per entity)
    - Optional attributes
- Variable structure (across domain)
    - Highly connected data
- Queries require joining entities or navigating relationships
    - Need to understand that things are connected
    - Need to understand *strength*, *weight* or *quality* of relationships

aws

# Navigate a Web of Global Tax Policies



THOMSON REUTERS

*"Our customers are increasingly required to **navigate a complex web of global tax policies** and regulations. We need an approach to **model the sophisticated corporate structures** of our largest clients and deliver an end-to-end tax solution. We use a microservices architecture approach for our platforms and are beginning to leverage Amazon Neptune as a graph-based system to quickly **create links within the data**."*

Tim Vanderham, Chief Technology Officer, Thomson Reuters Tax & Accounting

aws

# Amazon Neptune – A Fully-Managed Graph Database

**Fast**

Query billions of relationships with millisecond latency

**Reliable**

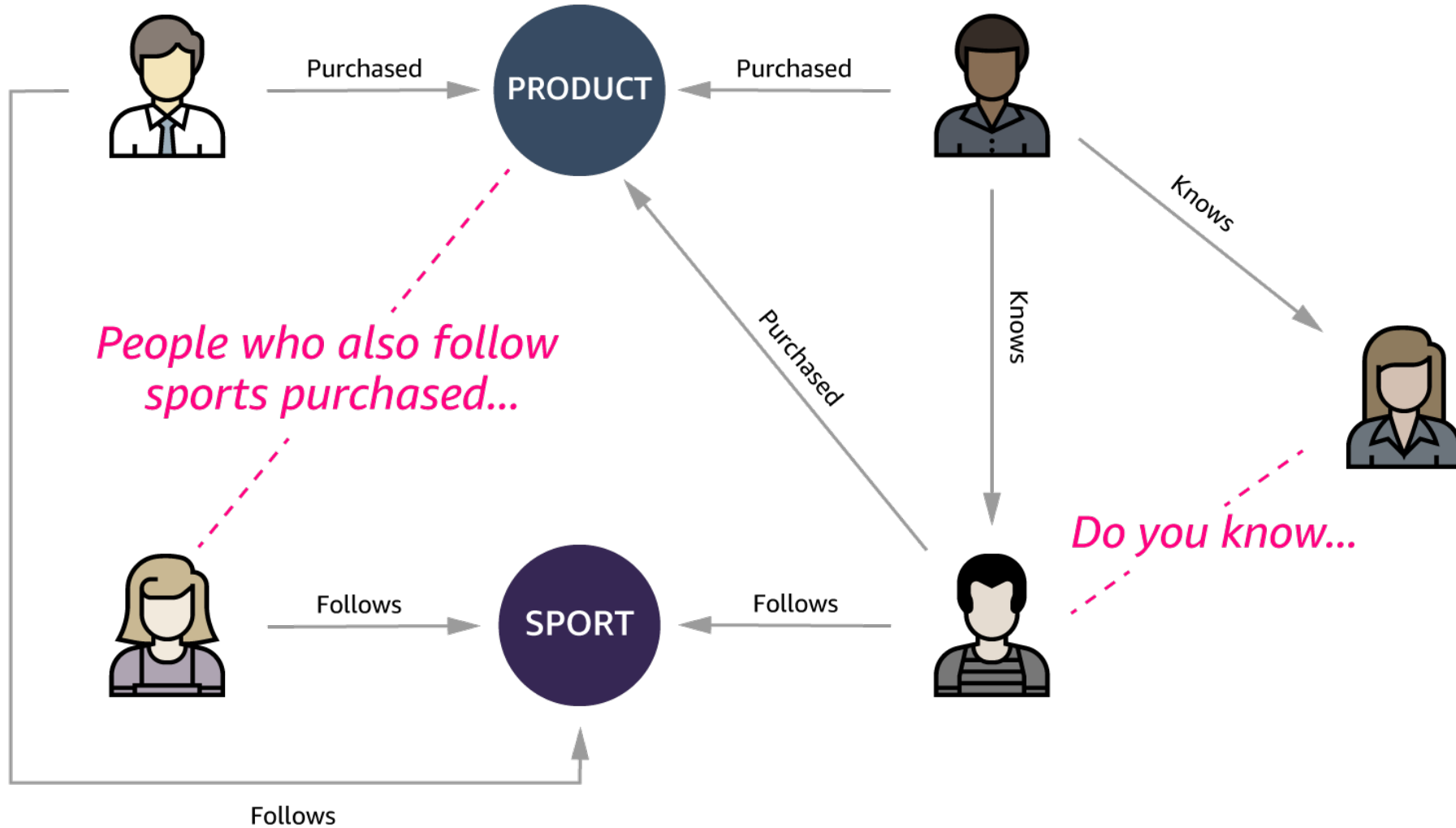6 replicas of your data across 3 AZs with full backup and restore

**Easy**

Build powerful queries easily with Gremlin and SPARQL

**Open**

Supports Apache TinkerPop & W3C RDF graph models
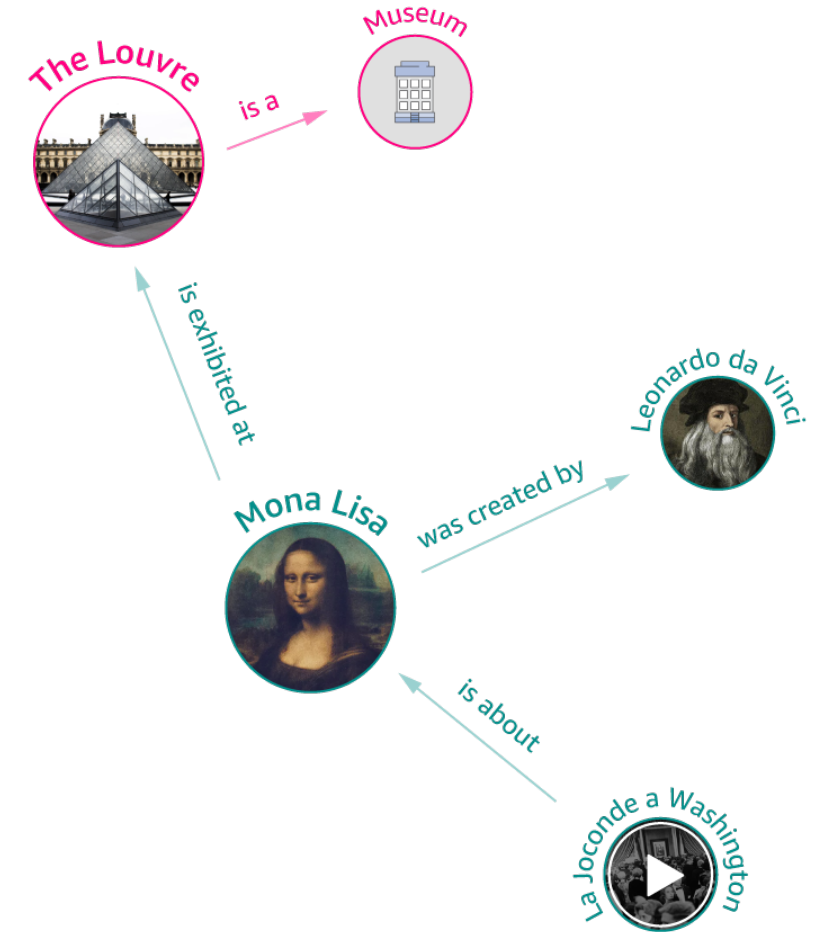
aws

# Recommendations Based on Relationships

# Knowledge Graph Applications

Who painted the
Mona Lisa?

What museums should
Alice visit while in Paris?

What artists have
paintings in The Louvre?

# Leading Graph Models and Frameworks

## Property Graph

- Vertices and edges (nodes and relationships) with properties
- Both record-like items

## Resource Description Framework (Rdf)

- Triples
- subject – predicate – object

aws

# Creating a Property Graph

Add 2 vertices

```
g.addV('User').property('name','Bill');
g.addV('User').property('name','Sarah');

...

g.V().hasLabel('User').has('name','Sarah').as('a').
  V().hasLabel('User').has('name','Bill').
  addE('FRIEND').to('a');
```
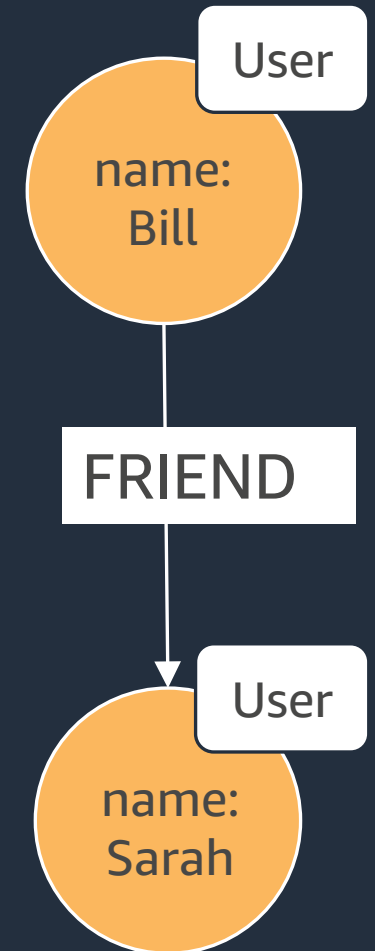
Connect with an edge

User
name: Bill

FRIEND →

User
name: Sarah

Gremlin (Apache TinkerPop 3.x)

aws

# Creating an RDF Graph

```
@prefix contacts: <http://www.socialnetwork.com/people#>.


<http://www.socialnetwork.com/person#1>
    rdf:type contacts:User;
    contact:name: "Bill" .


<http://www.socialnetwork.com/person#1>
    contacts:friend <http://www.socialnetwork.com/person#2> .


<http://www.socialnetwork.com/person#2>
    rdf:type contacts:User;
    contact:name: "Sarah" .
```
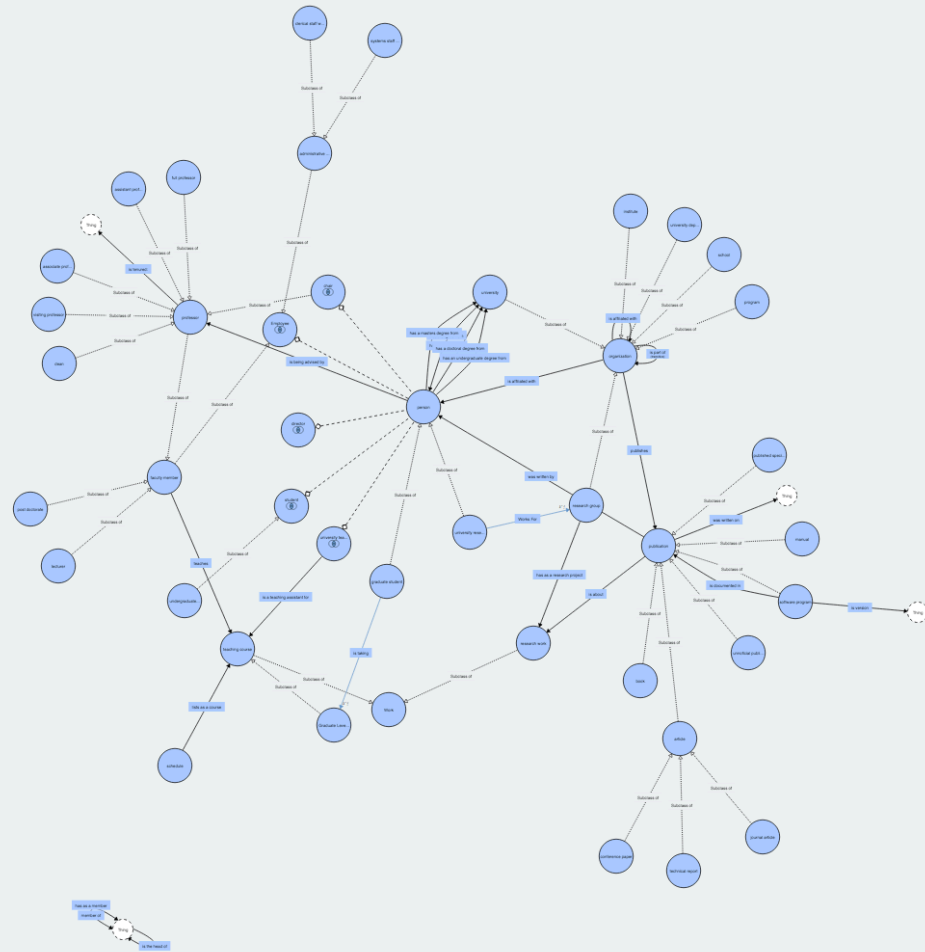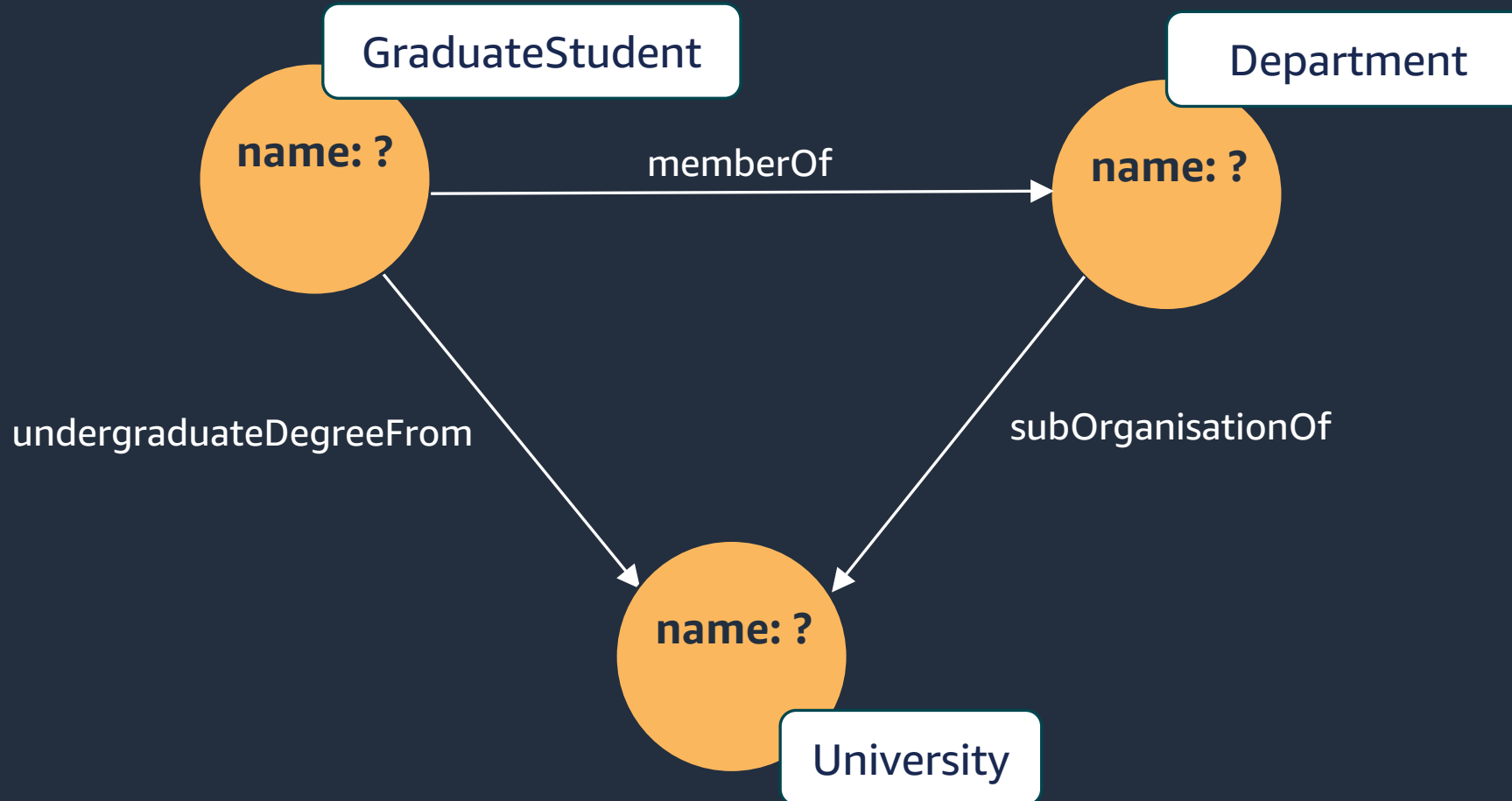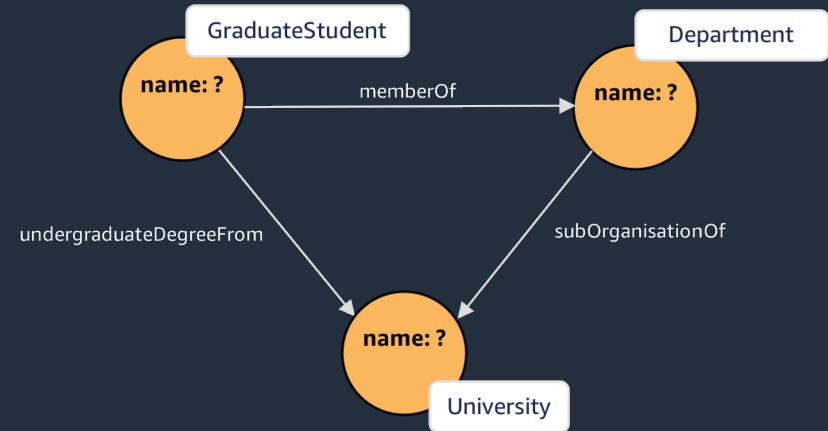
# Example: Universities



- Institutions

- Organisational structure

- Roles and job titles

- Research, teaching and administrative staff

- Subjects and course catalogs

- Timetables

- Undergraduate and graduate populations

# Find the graduate students who received an undergraduate degree from the same university
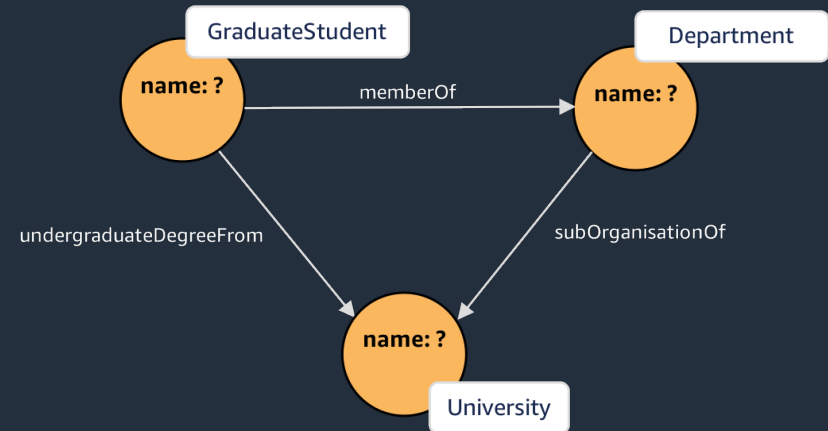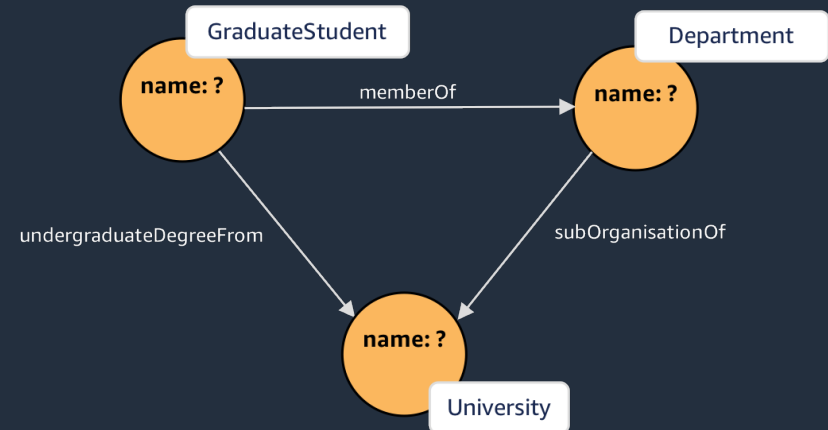
# Gremlin

```
g.V().hasLabel('GraduateStudent').as('student').
  out('memberOf').
  out('subOrganisationOf').
  in('undergraduateDegreeFrom').
  where(eq('student'))
```
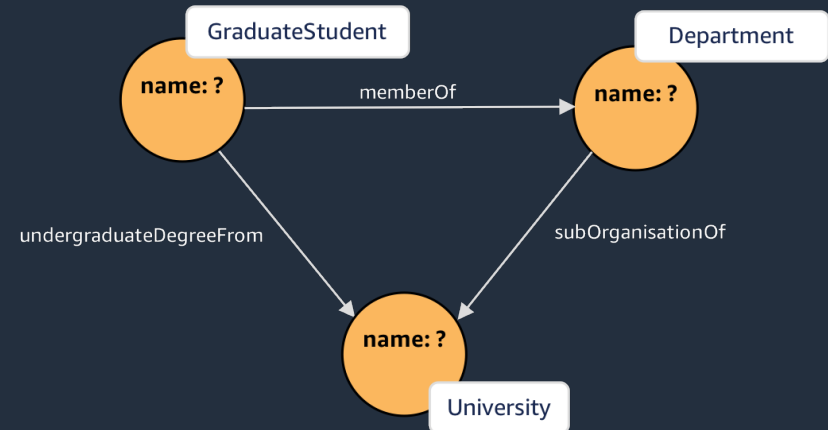
# Gremlin

```
g.V().hasLabel('GraduateStudent').as('student').
  out('memberOf').
  out('subOrganisationOf').
  in('undergraduateDegreeFrom').
  where(eq('student'))
```
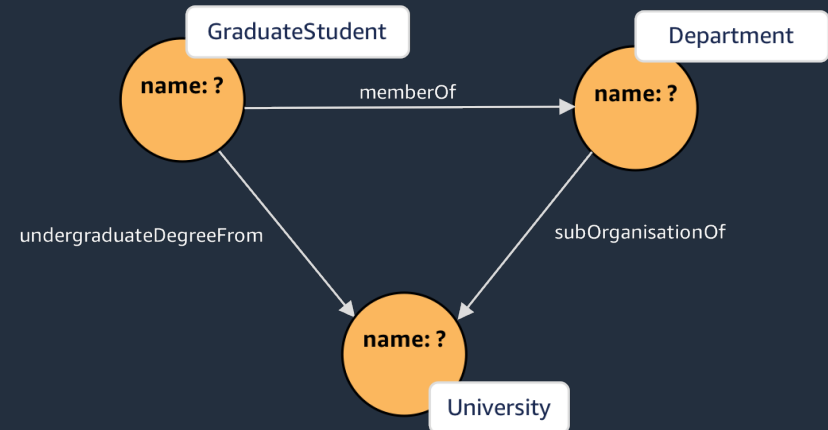
# Gremlin

```
g.V().hasLabel('GraduateStudent').as('student').
    out('memberOf').
    out('subOrganisationOf').
    in('undergraduateDegreeFrom').
    where(eq('student'))
```

# Gremlin

```
g.V().hasLabel('GraduateStudent').as('student').
  out('memberOf').
  out('subOrganisationOf').
  in('undergraduateDegreeFrom').
  where(eq('student'))
```
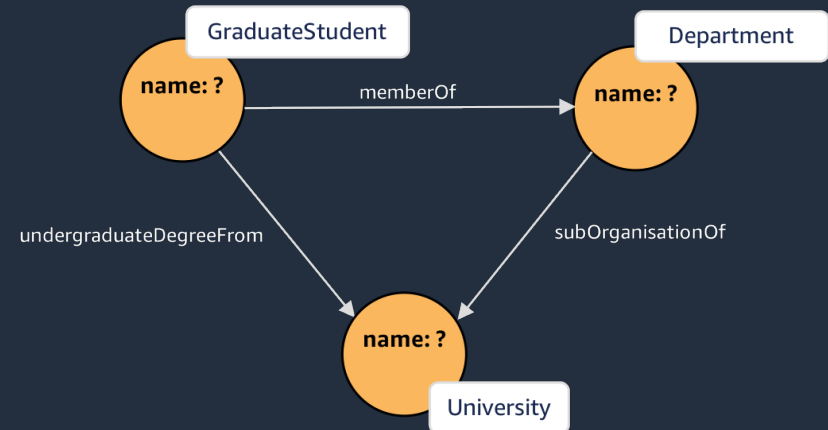
# Gremlin

```
g.V().hasLabel('GraduateStudent').as('student').
  out('memberOf').
  out('subOrganisationOf').
  in('undergraduateDegreeFrom').
  where(eq('student'))
```
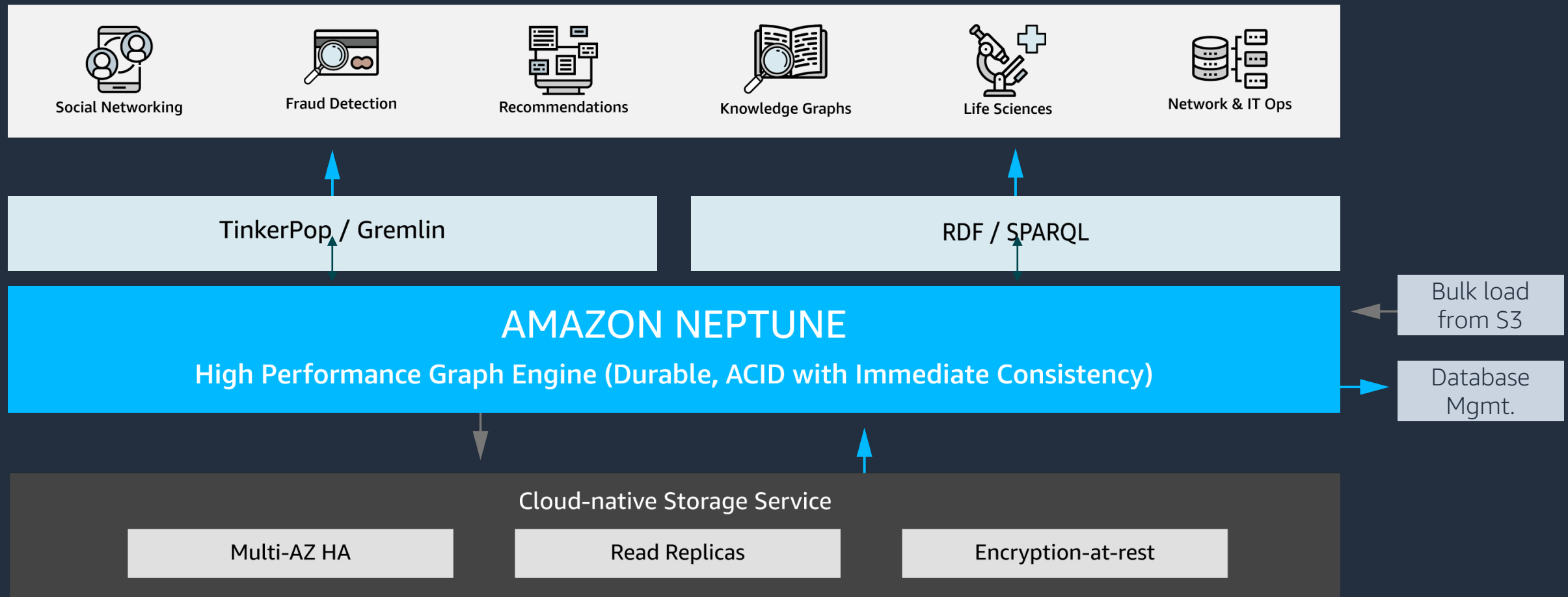
# Gremlin

```
g.V().hasLabel('GraduateStudent').as('student').
   out('memberOf').
   out('subOrganisationOf').
   in('undergraduateDegreeFrom').
   where(eq('student'))
```

# SPARQL

```sparql
PREFIX rdf:http://www.w3.org/1999/02/22-rdf-syntax-ns#
PREFIX ub:http://www.lehigh.edu/~zhp2/2004/0401/univ-bench.owl#

SELECT ?student WHERE {
        ?student rdf:type ub:GraduateStudent .
        ?univ rdf:type ub:University .
        ?dept rdf:type ub:Department .
        ?student ub:memberOf ?dept .
        ?dept ub:subOrganizationOf ?univ .
        ?student ub:undergraduateDegreeFrom ?univ
}
```
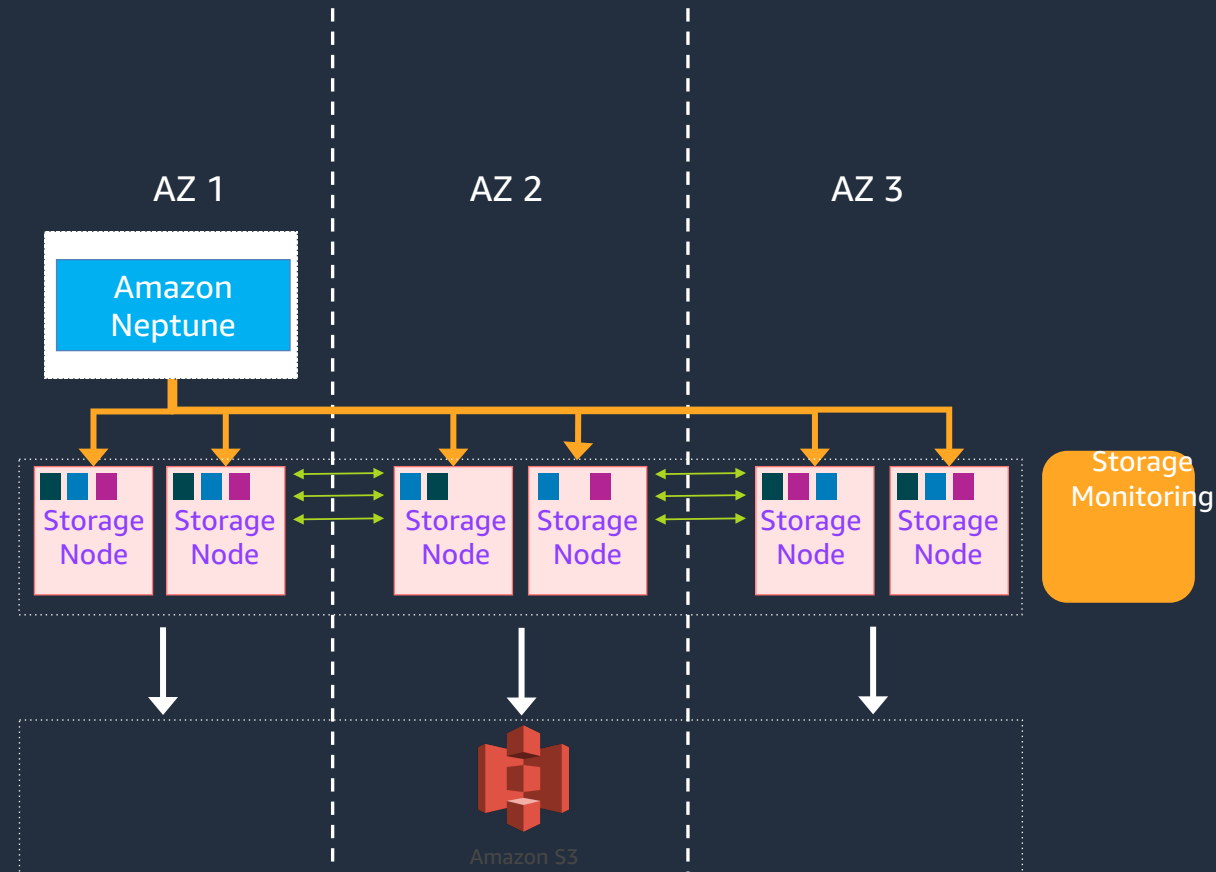
# Amazon Neptune High Level Architecture

Social Networking  Fraud Detection  Recommendations  Knowledge Graphs  Life Sciences  Network & IT Ops

TinkerPop / Gremlin

RDF / SPARQL

## AMAZON NEPTUNE

### High Performance Graph Engine (Durable, ACID with Immediate Consistency)

Bulk load from S3

Database Mgmt.

Cloud-native Storage Service

Multi-AZ HA

Read Replicas

Encryption-at-rest

aws

# Cloud-Native Storage

- Data is replicated 6 times across 3 AZs

- Continuous backup to Amazon S3
  - built for 11 9s durability

- Continuous monitoring of nodes and disks

- 10 GB segments as unit of repair or hotspot rebalance

- Quorum system for read/write; latency tolerant

- Quorum membership changes do not stall writes

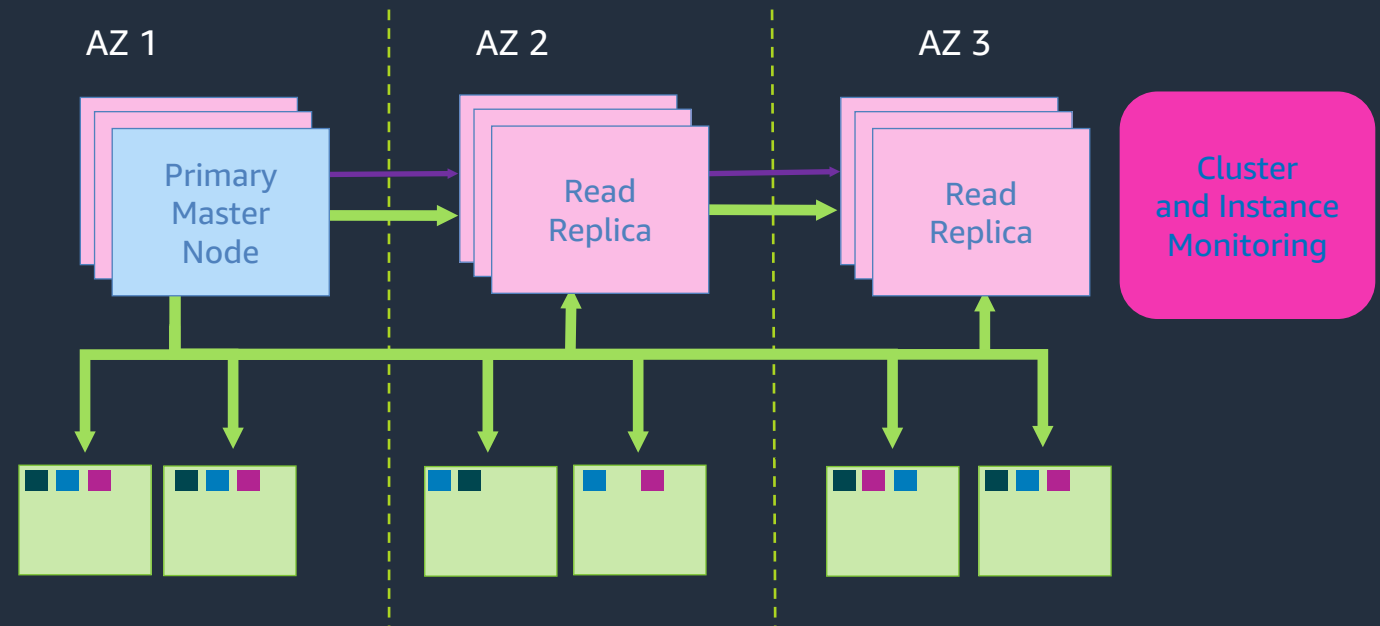- Storage volume automatically grows up to 64 TB

# Amazon Neptune Read Replicas

## Availability

- Failing database nodes are automatically detected and replaced
- Failing database processes are automatically detected and recycled
- Replicas are automatically promoted to primary if needed (failover)
  - Customer specifiable failover order

## Performance

- Customer applications can scale out read traffic across read replicas
- Read balancing across read replicas
  - Use reader endpoint

# Security



- Network isolation via Virtual Private Cloud
  - Use security groups to control ingress
- Encryption at rest using AWS Key Management Service (KMS)
  - underlying storage, automated backups, snapshots, and replicas in the same cluster
- IAM Policies to secure resources
- IAM Database Authentication
  - IAM Policies for database access
  - Each request must be signed using AWS Signature Version 4
  - Libraries for Gremlin and SPARQL clients

aws

# Loading Data into Neptune

**Bulk Load API**

- UTF-8 encoded files in S3

- Supports `gzip` compression of single files

- Requires a VPC endpoint for Amazon S3

- Add IAM Role to Neptune allowing `s3::Get*` and `s3::List*` permissions for S3 bucket

- CSV formatted files for Gremlin; 4 standard formats for RDF (N-Triples, N-Quads, RDF/XML, Turtle)

- Load, get status and cancel job via `/loader` HTTP endpoint

- Append-only

aws

# Backup and Restore

## Automated Backups

- Daily automated backups during backup window
- Full storage volume snapshot
- Taken from read replica if possible
- Retention period: 1-35 days

## Manual Snapshots

- Backup entire DB instance
- Can be shared with other AWS accounts

# Backup and Restore

## Restoring from a DB Snapshot (Automated or Manual)

- Creates new DB instance

- Apply custom parameter groups and security groups after restore using **Modify** command

## Point-in-Time Restore

- Restore from DB instance (not snapshot)

- Creates new DB instance

- Choose "Latest restorable time" or specify custom data and time

- 1 second granularity

aws

# Monitoring

**AWS CloudTrail**

- Log all Neptune API calls to S3 bucket

**Event Notifications**

- Create SNS subscription via CLI or SDK
- Sources: db-instance | db-cluster | db-parameter-group | db-security-group | db-snapshot | db-cluster-snapshot

**Amazon CloudWatch**



| CPUUtilization | GremlinRequestsPerSec | Http429 | SparqlErrors |
|---|---|---|---|
| ClusterReplicaLag | Http100 | Http500 | SparqlRequests |
| ClusterReplicaLagMaximum | Http101 | Http501 | SparqlRequestsPerSec |
| ClusterReplicaLagMinimum | Http200 | LoaderErrors | StatusErrors |
| EngineUptime | Http400 | LoaderRequests | StatusRequests |
| FreeableMemory | Http403 | NetworkReceiveThroughput | VolumeBytesUsed |
| GremlinErrors | Http405 | NetworkThroughput | VolumeReadIOPs |
| GremlinRequests | Http413 | NetworkTransmitThroughput | VolumeWriteIOPs |

# Audit Logs

`neptune_enable_audit_log` DB cluster parameter



- Timestamp
- ServerHost
- ClientHost
- ConnectionType
- RequestMesssage