# Understanding Datasets

**Seeing the Unseen through Graph Automations**
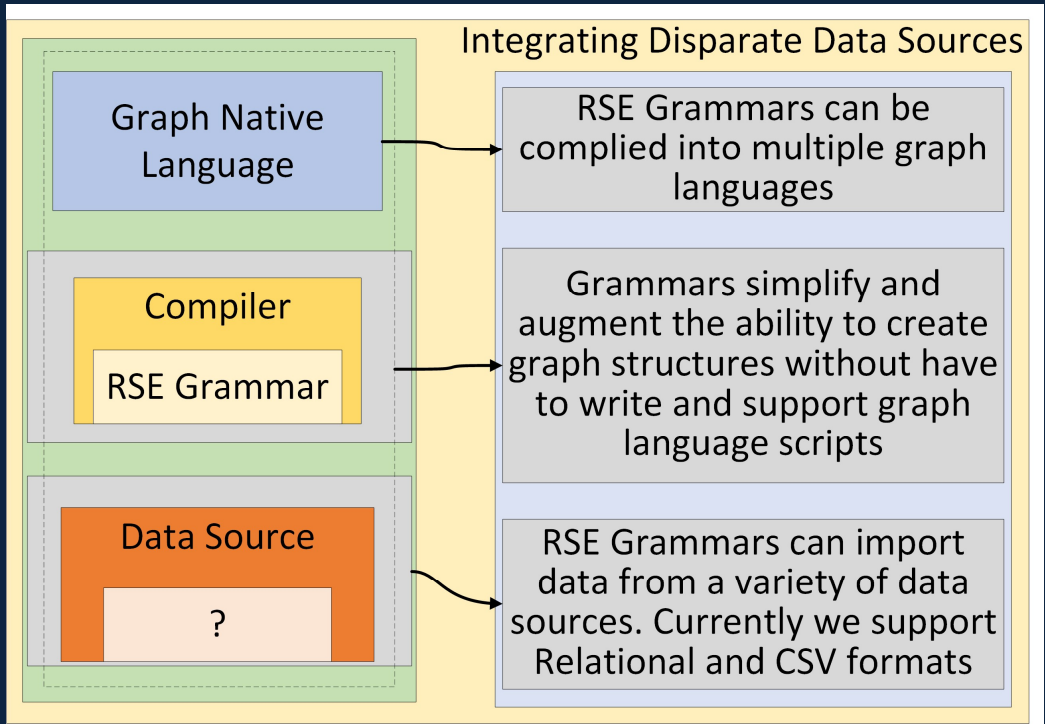
**September 2022**

**MITRE** | **SOLVING PROBLEMS FOR A SAFER WORLD™**

# The Need

- **Simplified mechanism for building and unifying disparate data sources into a normalize high performance Unified Data Fabric, which can be evaluated and understood for Analytics, Machine Learning, and Simulations**

- **Mechanisms to visually explore and understand both data entities and their relationships to one another**

- **Mechanisms to understand and traverse threads of augmented entities**

- **Mechanisms to project and understand through graph analytics the desired dataset**
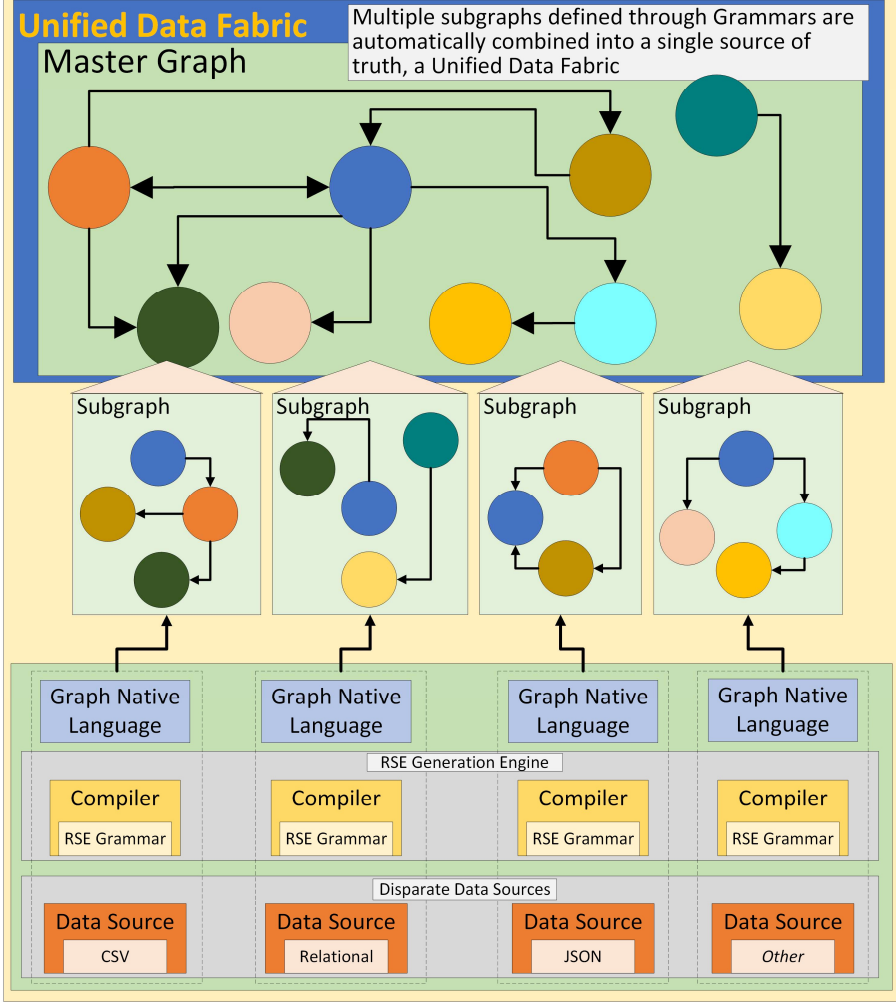
**The Rules Simulation Engine is the Answer, One mechanism for all these things, and more.**

# Why Graphs?

- High performance Entity<->Entity Relationship Management
- Any Entity can reference any other Entity. Graphs are highly Adaptable, and White Board Friendly (Schemaless)
- Graphs are the foundational element for Analytics, NLP, ML, and AI.
- Graphs are thousands of time faster than the best Relational Databases even on small datasets and the performance gap grows exponentially as dataset sizes increase
- Relationships are First Class Citizens
- Graphs are the cornerstone of our solution

## Find Reports and How Many are Managed Three Level Down

### Recursive Cypher

match (boss)-[:Manages*0..3]->(mgr)
where boss.name="Dave Green"
and (mgr-[Manages]->())

**Declarative**

Describe what you want
Not of how to do it

Graphs are schemaless and adaptable to changing requirements.

Relationships are first class citizens

### Recursive SQL

## Comparing Graph and Relational Structures

Dataset of 1000 people, each with 50 friends. Find all the paths from one person to another in four or less hops.

| | # of People | Query Time |
|---|---|---|
| RDBMS | 1,000 | 2000ms |
| Neo4J | 1,000 | 2ms |
| Neo4J | 1,000,000 | 2ms |

Dataset of 1,000,000 people each with 50 friends. Find friends of friends up to five levels deep.

| Depth | RDBMS (ms) | Neo4J (ms) | Records Returned |
|---|---|---|---|
| 2 | 0.016 | 0.010 | ~2,500 |
| 3 | 30.267 | 0.168 | ~110,000 |
| 4 | 1543.505 | 1.359 | ~600,000 |
| 5 | unfinished | 2.132 | ~800,000 |

rlukas

**MITRE**

Approved for Public Release; Distribution Unlimited. Public Release Case Number 22-2914

3

# Grammar Overview

**Grammars are complied into native graph languages, currently Cypher a Neo4J open standard, but could also generate other Graph languages.**
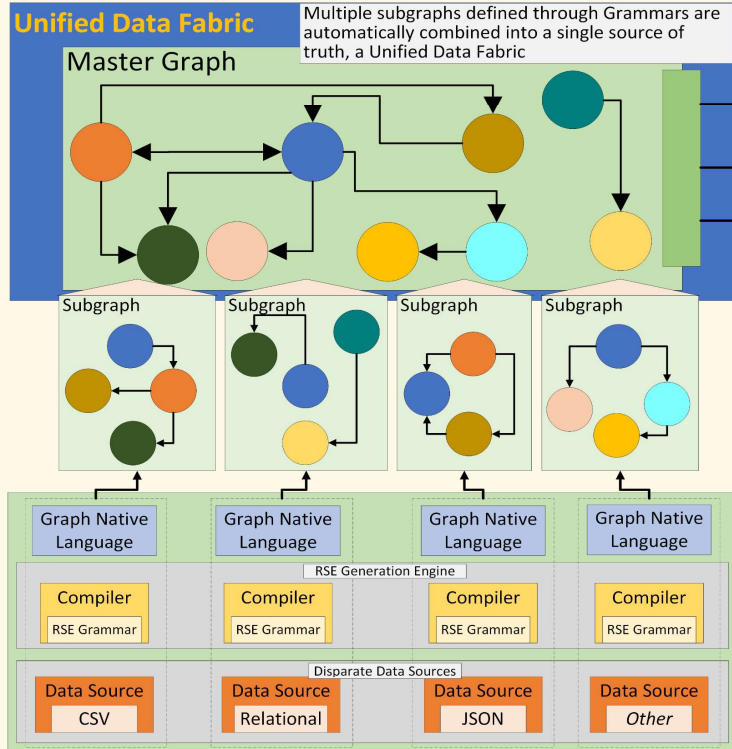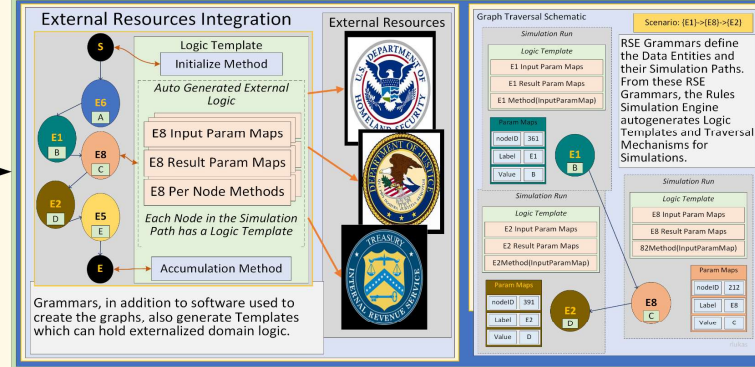**Grammars are able load data from an ever-growing set of data sources.**

## Integrating Disparate Data Sources

| Graph Native Language | RSE Grammars can be complied into multiple graph languages |

| Compiler<br>RSE Grammar | Grammars simplify and augment the ability to create graph structures without have to write and support graph language scripts |

| Data Source<br>? | RSE Grammars can import data from a variety of data sources. Currently we support Relational and CSV formats |

**Unified Data Fabric**
Master Graph

Multiple subgraphs defined through Grammars are automatically combined into a single source of truth, a Unified Data Fabric



Subgraph    Subgraph    Subgraph    Subgraph

| Graph Native Language | Graph Native Language | Graph Native Language | Graph Native Language |

RSE Generation Engine

| Compiler<br>RSE Grammar | Compiler<br>RSE Grammar | Compiler<br>RSE Grammar | Compiler<br>RSE Grammar |

Disparate Data Sources

| Data Source<br>CSV | Data Source<br>Relational | Data Source<br>JSON | Data Source<br>Other |

**MITRE**

Approved for Public Release; Distribution Unlimited. Public Release Case Number 22-2914

4

# Overview Schematic

## Rules Simulation Engine Schematic

Easily "See the Unseen" through an easily constructed, high performance, Unified Data Fabric composed from multiple subgraphs. Subgraphs are constructed, from disparate data sources, through our Grammar technologies.

Insights from the Unified Data Fabric are enhanced through high powered analytics, machine learning technologies, as well as our Rules Simulation Engine powered by autogenerated domain agnostic Templates.
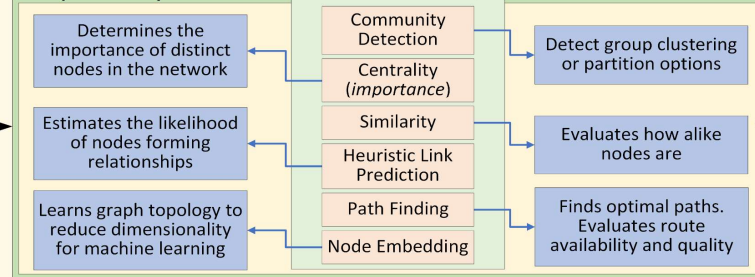
**Unified Data Fabric**

Multiple subgraphs defined through Grammars are automatically combined into a single source of truth, a Unified Data Fabric

Master Graph

Subgraph  Subgraph  Subgraph  Subgraph

Graph Native Language | Graph Native Language | Graph Native Language | Graph Native Language

RSE Generation Engine

Compiler | Compiler | Compiler | Compiler
RSE Grammar | RSE Grammar | RSE Grammar | RSE Grammar

Disparate Data Sources

Data Source | Data Source | Data Source | Data Source
CSV | Relational | JSON | *Other*

### Simulation Through Autogenerated Templates

External Resources Integration

Logic Template
Initialize Method

*Auto Generated External Logic*

E8 Input Param Maps
E8 Result Param Maps
E8 Per Node Methods

*Each Node in the Simulation Path has a Logic Template*

Accumulation Method

External Resources

Grammars, in addition to software used to create the graphs, also generate Templates which can hold externalized domain logic.

Graph Traversal Schematic

Scenario: {E1}->{E8}->{E2}

*Simulation Run*
*Logic Template*
E1 Input Param Maps
E1 Result Param Maps
E1 Method(InputParamMap)

Param Maps
nodeID | 361
Label | E1
Value | B

*Simulation Run*
*Logic Template*
E2 Input Param Maps
E2 Result Param Maps
E2Method(InputParamMap)

Param Maps
nodeID | 391
Label | E2
Value | D

*Logic Template*
E8 Input Param Maps
E8 Result Param Maps
82Method(InputParamMap)

Param Maps
nodeID | 212
Label | E8
Value | C

RSE Grammars define the Data Entities and their Simulation Paths. From these RSE Grammars, the Rules Simulation Engine autogenerates Logic Templates and Traversal Mechanisms for Simulations.

Traverse graph pathways invoking agnostic externalize domain logic. Grammars also generate these domain agnostic Templates

### Graph Analytics

Data Science Library

Determines the importance of distinct nodes in the network ← Centrality (*importance*)
Community Detection → Detect group clustering or partition options

Estimates the likelihood of nodes forming relationships ← Heuristic Link Prediction
Similarity → Evaluates how alike nodes are

Learns graph topology to reduce dimensionality for machine learning ← Node Embedding
Path Finding → Finds optimal paths. Evaluates route availability and quality

Perform rich graph analytics over a high performance, easily customizable, data structures
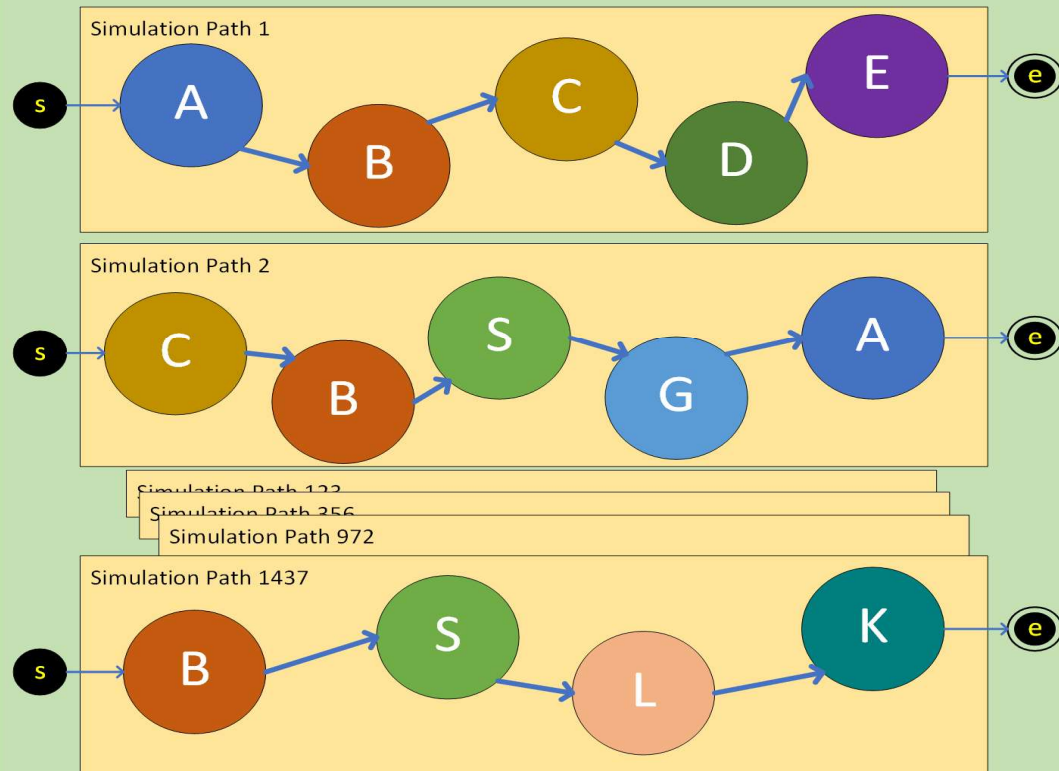
### Machine Learning and Knowledge Graphs

Combined Subgraphs

Relationships are First Class Citizens in a graph. Strong relationship understanding improves Machine Learning results. Our technology interfaces directly with Machine Learning services.

Data Science Library

| FastRP | Linear Algebra | Speed Accuracy |
| GraphSAGE | Neural Networks (GCNN) | Inductive Reasoning |
| Node2Vec | Neural Networks (SkipGram) | Interpretable Hyperparameters |

Embedded Vectors
1,0,1,0,1,0,1,1,0,0,0,1,1,0,0,0,1
1,0,0,1,1,1,0,0,0,.,10,1,1,0,0,0,1
1,1,1,1,0,0,1,1,1,1,1,0,0,1,0,1,1

Machine Learning Engine
Knowledge Graph
Descriptive Entity to Entity Relationships

Easily vectorize RSE structures for integration with Machine Learning Engines.

**MITRE**

Approved for Public Release; Distribution Unlimited. Public Release Case Number 22-2914

5

# Making Decisions Through Thread Pathways

## Create & Accumulate Graphs

### Variable Path, Variable Steps



Simulation Path 1

s → A → B → C → D → E → e

Simulation Path 2

s → C → B → S → G → A → e

Simulation Path 123
Simulation Path 256
Simulation Path 972

Simulation Path 1437

s → B → S → L → K → e

Most projects can identify *Data Entities* that might contribute to the eventual outcome, with new *Data Entities* constantly emerging.
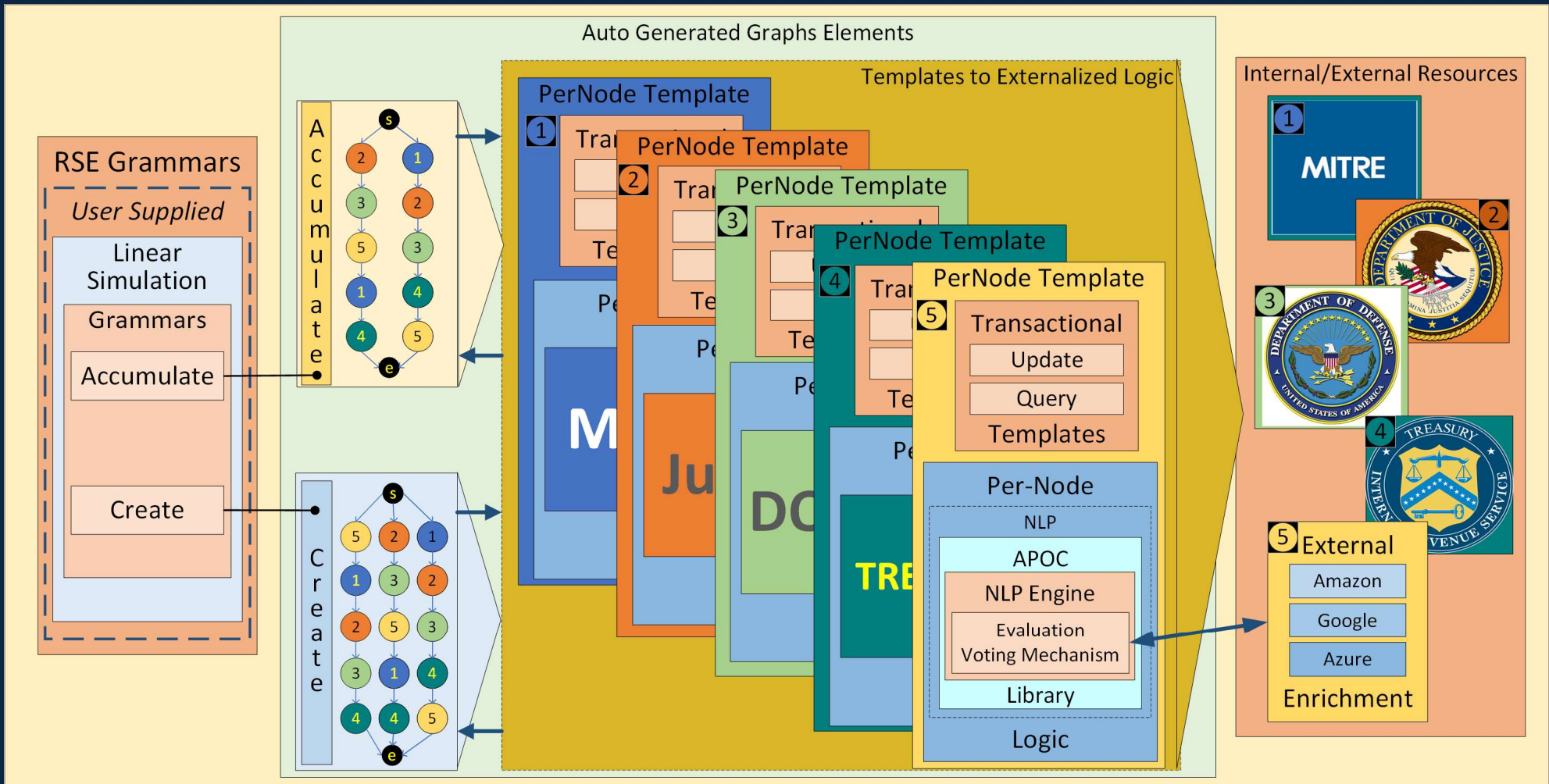
Our Solution provides a simple mechanism to define scenarios. These scenarios select, and order, *Data Entities* which can be included into variable length *Simulation Paths*.

{A}->{B}->{C}->{D}->{E},
{C}->{B}->{S}->{G}->{A},
{B}->{S}->{L}->{K}, ....

*Data Entity* values may have multiple settings. For example, Income Levels may be described as a set of ranges.

Because of the dynamic nature of projects, *Simulation Paths* must be easily, and dynamically, defined from external sources.

**MITRE**

Auto Generated Graphs Elements

Templates to Externalized Logic

Internal/External Resources

**RSE Grammars**

*User Supplied*

Linear Simulation

Grammars

Accumulate

Create

Accumulate

Create

PerNode Template

1 Transactional

PerNode Template

2 Transactional

PerNode Template

3 Transactional

PerNode Template

4 Transactional

PerNode Template

5 Transactional

Update

Query

Templates

Per-Node

NLP

APOC

NLP Engine

Evaluation Voting Mechanism

Library

Logic

1 MITRE

2 Department of Justice

3 Department of Defense

4 Treasury Internal Revenue Service

5 External

Amazon

Google

Azure

Enrichment

Approved for Public Release; Distribution Unlimited. Public Release Case Number 22-2914

7

# Data Science Library: A Global View

## Data Science Library

**Community Detection**

**Centrality (*importance*)**

**Similarity**

**Heuristic Link Prediction**

**Path Finding**

**Node Embedding**

Determines the importance of distinct nodes in the network

What are the important predictors that we should consider?

Detect group clustering or partition options

What do our clients have in common? How can we discover emerging group types?

Estimates the likelihood of nodes forming relationships

Clients with X probably also would choose Y.

Evaluates how alike nodes are

Maybe we don't need both nodes, simplify our models? Maybe Discover new trends?

Learns graph topology to reduce dimensionality for machine learning

Graphs are the base technologies for neural networks, so let us use them to feed our machine learning services.

Finds optimal paths. Evaluates route availability and quality

What are the most common pathways through a system, those are the ones we should optimize.

rlukas

**MITRE**

# Solution Entity Pallet Mechanisms (CI/CD)

- **Solution Entities are the cornerstone of our solution.**

- **Solution Entities undergo a Certification Process and normal CI/CD pipeline**

- **Solution Entities are combined provide perspectives into our high-performance datasets.**

# The AI Infrastructure Pipeline

Solution Entity Adapters are the core mechanism for integrating disparate components into a collective AI Infrastructure through Declarative Interfaces.

**Diagram Legend**
- Declarative Interface
- Solution Entity Mechanism
- Solution Entity Adapter

**Solution Entities are the next generation, domain agnostic Rule Engine Templates, residing in, and are selected from, a pallet library. Solution Entities expose and are connected through declarative interfaces**

**MITRE**

# Interface Mechanisms

## Current Interface Schema

Servers and Clients are Roles. A component often operates in both Roles

The need for an attribute set, generally requires acquisition through multiple interfaces and multiple requests from the Client.

The Client then has to parse the presented structures to obtain the desired results.

The Client is tightly bound to the structure of the Server. Changes to the Server interfaces (X -> X$_1$) directly affects all X dependent Clients.

**Client**

Requests
1
2  X$_1$

Response Processing
Extract
Corollate

1  2  x → X$_1$

Service
1  Request Processing
2  Request Processing
X$_1$  Request Processing

Each Service interface exposes one or more methods (functions with parameters and return types)

SOAP, REST, and Language Native, although different, are common interface types

Declarative Interfaces establish a common declarative request mechanism throughout the entire topology

Tightly bound components are contrary to Plug and Play topologies

Interface versioning complicates the topology and the exposing Services

Transferring unrequested attributes harms performance on multiple levels

## Declaritive Interface Schema

**Client**

Requested Attributes
A  Z  B

Response
A → ray
B → lukas
Z → MITRE

G

Service

Request Processing

Schema
A  SQL
B  Graph
Z  Object
MS Office

The Service Exposes a single interface to all Clients and Client versions.

Requests contain a the list of Requested Attributes. Services send only Requested Attributes simplifying Client processing.

The Schema contains pre-canned components which provide access to disparate data sources.

The Schema creates a standardized Key:Value Response object which is returned to the Requestor.

Declarative Interfaces simplifies both the micro and macro topologies

## Declarative interfaces promote a Plug and Play Infrastructure

**MITRE**

# Gaining and Evaluating Different Perspectives

## Perspectives

**P1**
- Solution Entity — *Unique Logic* — Mechanism 1 — Adapter
- Solution Entity — *Unique Logic* — Mechanism C — Adapter
- Solution Entity — *Unique Logic* — Mechanism X — Adapter

**P2**
- Solution Entity — *Unique Logic* — Mechanism C — Adapter
- Solution Entity — *Unique Logic* — Mechanism 3 — Adapter
- Solution Entity — *Unique Logic* — Mechanism Z — Adapter

**P3**
- Solution Entity — *Unique Logic* — Mechanism 1 — Adapter
- Solution Entity — *Unique Logic* — Mechanism Q — Adapter
- Solution Entity — *Unique Logic* — Mechanism Z — Adapter

## Evaluation Engine

- Measures
- Recommendations

| Ranking P3 | Ranking P2 | Ranking P1 |
|------------|------------|------------|
| 0.3 | 1.4 | 10.3 |
| Score | Score | Score |

## Ranked Recommendations

| Ranking P1 | Ranking P2 | Ranking P3 |
|------------|------------|------------|
| 10.3 | 1.4 | 0.3 |
| Score | Score | Score |

- **Solution Entities are combined in any order for multiple perspectives**
- **Perspectives are then fed into the Evaluation Engine which produce Ranking Scores.**

**MITRE**

# Observers Provide a Global View

**Each Solution Entity Adapter implements an Observer Pattern which reports Events are Relevant State into a common Mapping Service Repository**

MITRE

13

# General Mapping Service

Solution Entity States are preserved directly through their Adapters or the Observers into a high-performance graph structure which are easily be explored visually or through graph analytics.

# Wrap Up

- **Simplified mechanism for building and unifying disparate data sources into a normalize high performance Unified Data Fabric, which can be evaluated and understood**
  MITRE Grammars/Compiler provide a simple mechanism to create new graph structures to address emerging explorations through multiple services

- **Mechanisms to visually explore and understand both data entities and their relationships to one another**
  **Bloom:** Perspectives, Share Discoveries within teams, Near Natural Language processing
  **Browser:** Cypher developers interface simplified through our Grammars

- **Mechanisms to understand and traverse threads of augmented entities**
  Templates and Solution Entities provide domain agnostic Rules Engine thread processing

- **Mechanisms to project and understand through graph analytics the desired dataset**
  Analytic Projections and Machine Learning integration through Data Science Library

**The Rules Simulation Engine is the Answer, One mechanism for all these things, and more.**

**MITRE**

MITRE Labs inspires breakthroughs in applied science and advanced technology to transform the future of U.S. scientific and economic leadership. Our goal: Deliver disruptive innovation to support our mission of solving problems for a safer world.

**Please feel free to contact me with any questions**

**Thank You**

**Ray Lukas (rlukas@mitre.org)**
**Jim Lockett (jdlockett@mitre.org)**

**The Attached Appendix Supplies More Detailed Explanations of RSE Grammars and Subgraphs**

**rlukas@MITRE.org**

# Database Comparisons

| Database | Purpose | Attributes |
|---|---|---|
| **Relational** | | |
| Oracle | Table Based | Highly Refined |
| SQL Server | Rows and Columns | ACID Transactions |
| | | Schema Based |
| **No SQL** ⭐ | | |
| MongoDB | Key-Value Storage | Eventual Consistency |
| Cassandra | Document Storage | Great for Web Interactions |
| Spanner | | |
| **Graph** | | |
| Anzo | Model Simplicity | White Board Friendly |
| Neo4J | High Speed Entity-Entity Traversals | ACID Transactions |
| | | Schemaless |

⭐ Mongo DB contains ACID Transactions

## Relational

Strict schema and data normalization separating data into tables. To preserve data consistency ACID transactions are supported. This imposes limitations on how relationships can be queried. Translations from OO to Rel are difficult and expensive.

Relational model and other NoSQL database models link the data by implicit connections

Relationships are reunified at query time

## Graph

Whiteboard Friendly/Object Oriented: Native Graph databases have no pre-canned schema. Structures are directly mapped, any node can point to any other node. Unlimited query environment.

Relationships are first-class citizen in a graph database and can be labelled, directed, and assigned properties

★ Graphs and their Analytics Libraries are <u>Very Scalable</u>. Our graph technology can handle Trillions of nodes. Graphs with less than 10 million nodes and relationship are considered small.

**But Which Is The Best? It Depends On What You Are Doing**

**Performance: Connections are made a creation time, not at query time**

**MITRE**

# RSE General Terminology

## Basic Graph Terminologies

**Grammars**: Grammars provide a simplified abstraction mechanism from the underlying graph technologies (in our case the Cypher Language) used to describe, and create graph, structures. Grammars are much like a high level programming language (C, C++, Java) which abstract the underlying CPU capabilities (assembly/machine code) mechanisms. *Grammars are compiled into the underlying Cypher codes.*

**Pattern**: Patterns are basic text strings (called ASCII Art) which describes a Node's Relationship to another Node. (Node)-[Relationship]->(Node). *See Appendix for details.*

**Cypher**: Cypher is the language of Graphs. Our Grammars are abstractions of Cypher and are compiled into Cypher which can be easily deployed to create new high performant graph structures. Cypher's compliment in the relational data world would be SQL. Cypher is basic *SQL for Graphs, but much better.*

**Create and Merge Commands**: Two basic and commonly used Cypher commands. Create creates the presented Pattern, even if is already exists. Merge uses the presented Pattern if it already exists, else it creates a new pattern in the database. *This is the underlying mechanism for auto connecting subgraphs and/or sharing threads in the Accumulation graphs.*

## Rules Simulation Engine Terminologies

**Rules Engine**: A Rules Engine is a service that allows domain (system) logic to be defined and invoked externally from the engine. Since the logic of the system is externalized from the engine, this allows the engine to be easily adapted to, and used across, virtually any domain. *This is a cornerstone concept in our solution, and how we provide a domain agnostic solution.*

**Simulation Path**: A graph is a series of interconnected entities (nodes). Our Rules Simulation Engine is able to traverse these pathways and for each node invoke externalized domain logic living inside autogenerated templates. This is how our Rules Simulation Engine implements a high performance, domain agnostic, decision/simulation mechanism. *This technique is generally used with Create and Accumulate (threading type) graphs. Merge graphs are generally for graph analytics.*

**MITRE**

# RSE Grammar Node Types

## RSE Grammar Node Types

### Standard Element Types

Node (instanceVar:LabelList{AttributeNameValueList})

Relationship (node)-[instanceVar:RelationshipLabelList{AttributeNameValueList}]-(node)

### Enhanced Element Types

#### Enhanced Visualization

**Index Nodes**: These are special autogenerated nodes which group like nodes together. For example, all the states in a country, all the counties in a state, all cities in a county all the data centers for each city. Each of these topology nodes are index nodes.

#### Graph Debugging

**Data Source Nodes**: Special nodes created for each data source that connect to all the nodes from that data source. Data Source nodes allow you to visualize the data sources for each node in a subgraph.

**Error Graphs**: ErrorType nodes are generated when errors are found in the underlying data source. Each ErrorType node points to nodes which describe each occurrence of such errors. Currently null data elements are supported, with additional error detection schemas under development

**MITRE**

# Unified Data Fabric

**Subgraph Technology**: Each Data Source can define one or more subgraphs. Subgraphs can be loaded in any order, or not a at all.

**Proxy Nodes**: Proxy Nodes (DS1:2 and DS3:1) are Place holder (Proxy) nodes which define connection points that will be mapped to by other subgraphs. Proxy Nodes hold only the state information defined in that Data Source.



**Data Sources** — Three Independent Sub Graphs

Data Source 1
- DS1:1 [Atr1, Atr2, Atr3]
- Item 1: Augmented [Atr1, Atr2, Atr3]
- DS1:2 **P** — Item 2: Proxy Node

Data Source 2
- DS2:2 [Atr4, Atr5] — Item 2: Augmented [Atr4, Atr5]
- DS2:6 [Atr1, Atr3] — Item 6: Augmented [Atr1, Atr3]

Data Source 3
- DS3:5 [Atr4, Atr1] — Item 5: Augmented [Atr4, Atr1]
- DS3:4 [Atr1, Atr3]
- DS3:1 **P** — Item 1: Proxy Node
- Item 4: Augmented [Atr1, Atr3]

Creates →

**Unified Data Set/Fabric** — Load Order Agnostic
- DS1:1 [Atr1, Atr2, Atr3]
- DS2:2 [Atr4, Atr5]
- DS2:6 [Atr1, Atr3]
- DS3:5 [Atr4, Atr1]
- DS3:4 [Atr1, Atr3]

Subgraphs are Defined in the RSE V2.0 grammars.
Load Only the Graphs/Data that you Need and Want for a Particular purpose.
Proxy definitions are similar to Foreign Keys in a Relational Database, but Hyper Performant.
Each Node Definition can be augmented with multiple property values.
Nodes are still able to invoke external templates and external resources.
External Templates can also modify graph structures.

## Proxy Node Mappings

- DS1:2 is defined in DS2:2. When DS2 is loaded DS1:2 is mapped and augmented with Attribute 4 and Attribute 5
- DS3:1 is defined in DS1:1. Loading the Data Source 1 subgraph maps to DS3:1 and load Attribute 1, Attribute 2, and Attribute 3

**MITRE**

# Visual Exploration Mechanisms

## Bloom

GPU Accelerated Visualization

High performance physics and rendering

- Easily discover patterns yielding more questions and explorations.
- Near natural language visual explorations of graph structures without knowing SQL or Cypher. Bloom even can suggests entities to include in your visualizations.
- Provides mechanisms to easily update graph entities without knowing Cypher
- Graph customizations (millions of colors, property based styles, icons, auto sizing, etc.)
- Easily augmented capabilities through embedded Cypher enrichments
- Provides defined perspectives which are tailored for a specific role or context
- Subgraphs, Perspectives, Scenes, and other enrichments can be shared

## Browser

- Developer Tool requires Cypher, the Language of Graphs, Knowledge
- Limited display (colors, sizing of entities, no node icons, etc.) capabilities

**MITRE**

# Node Definition

## Grammar Components

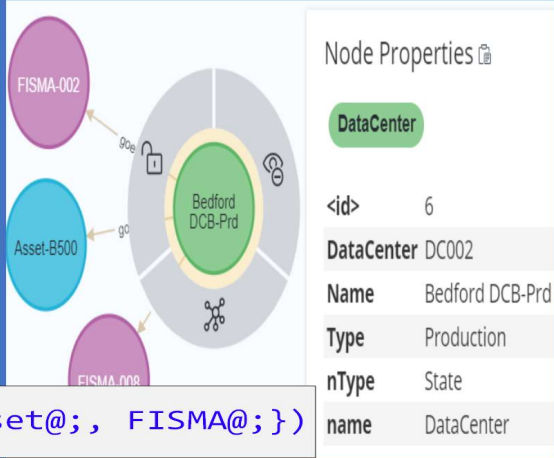| | |
|---|---|
| DS1 | `DataCntDS:DataCenter(Name{DataCenter, Type, Asset@;, FISMA@;})` |
| NodeName | `DataCntDS:DataCenter(Name{DataCenter, Type, Asset@;, FISMA@;})` |
| Identifier | `DataCntDS:DataCenter(Name{DataCenter, Type, Asset@;, FISMA@;})` |
| {Aug1, Aug2} | `DataCntDS:DataCenter(Name{DataCenter, Type, Asset@;, FISMA@;})` |

### MultiFields

MultiFields contain multiple delimited data elements embedded into a single column.
In this example these are separated by the semicolon character.

`[- - -]` Multifield Definition  Semicolon Delimiter

## Node Definition Grammar

DataSource

Unique Node Identifier. Attribute Name will be "value" in the subgraph

**DS1: NodeName( Identifier {Aug1, Aug2} )**

The Name of the Node in the Subgraph

Additional Data Source Resident Augmentations for this Node

`DataCntDS:DataCenter(Name{DataCenter, Type, Asset@;, FISMA@;})`

## Data Center Node



Node Properties
DataCenter
| <id> | 6 |
| DataCenter | DC002 |
| Name | Bedford DCB-Prd |
| Type | Production |
| nType | State |
| name | DataCenter |

FISMA-002
Asset-B500
Bedford DCB-Prd
FISMA-008

---

**Data Source**
Defines where this node should obtain its attributes. Current supported data sources include CSV files and Relational Database Tables.

**Node Name**
Node names are no longer restricted to the column name.

**Identifier**
List of the unique identifiers (primary key concept) for this node.

**Augmentations**
Defines additional attributes from the data source which should be included into this Node.

**Improved Autogenerated Graph Structures**
Accumulate, Create, and Merge graph

---

**MITRE**

# Subgraph Definition

## Subgraph Definition

### Root Node -> Destination Nodes

**Root Node**

```
DS1: NodeName( Identifier    {Aug1, Aug2}  )    >
```

**Destination Node List**

```
[  DS1: NodeName( Identifier    {Aug1, Aug2}  ) ,

   DS1: NodeName( Identifier    {Aug1, Aug2}  ) ]
```

### Data Center Subgraph



### Graph Definition Example

```
DataCntDS:Type@I(Type)  >  [DataCntDS:DataCenter(Name{DataCenter, Type, Asset@;, FISMA@;})]
```
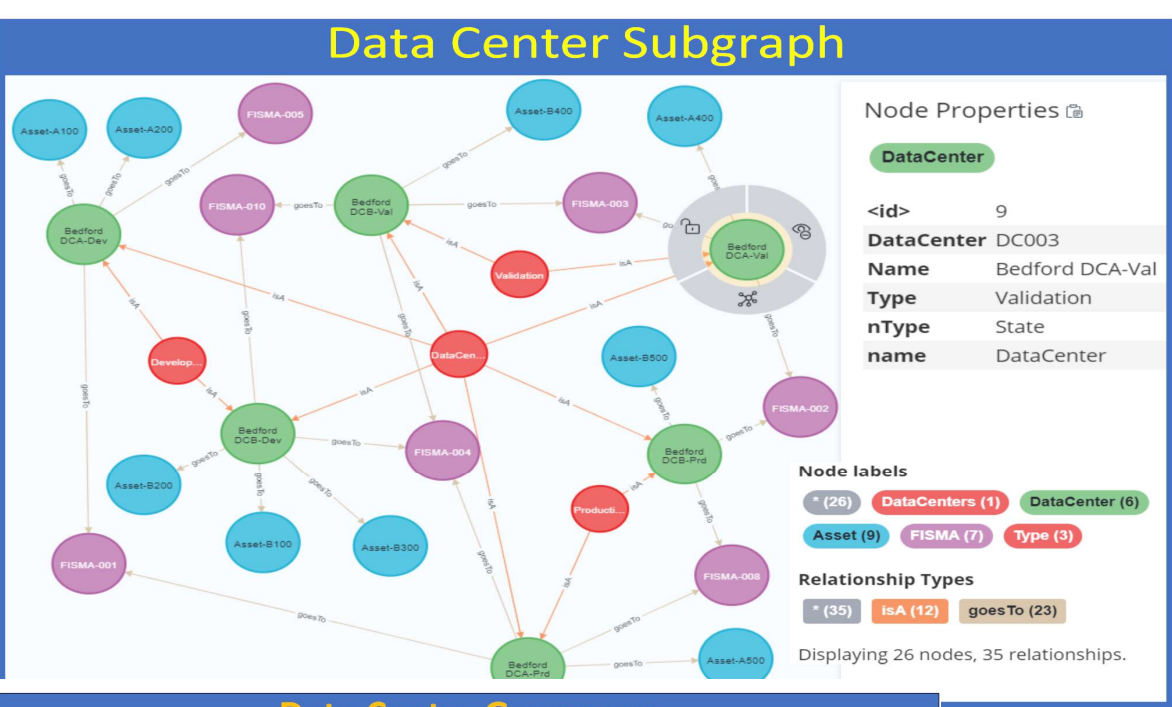
**Subgraph Technology:**

- **Each Data Source defines can a subgraph. The DataCenter data source defines the following Data Center subgraph.**

- **Multiple subgraphs can, through our simplifies grammars, be defined each from separate or shared data sources and types.**

- **Subgraphs can be loaded in any order, or not at all. Loading a subgraph automatically connects like nodes together into a Master Graph**

- **Only load the data you need for your specific explorations**

**MITRE**

23

# Data Center Subgraph

## Data Center Subgraph



**Node Properties**

**DataCenter**

| | |
|---|---|
| <id> | 9 |
| DataCenter | DC003 |
| Name | Bedford DCA-Val |
| Type | Validation |
| nType | State |
| name | DataCenter |

**Node labels**

* (26)  DataCenters (1)  DataCenter (6)  Asset (9)  FISMA (7)  Type (3)

**Relationship Types**

* (35)  isA (12)  goesTo (23)

Displaying 26 nodes, 35 relationships.

We can see that the Data Center Data Source contains two multifield elements, Assets and supported FISMA elements

There are two types of Index nodes, both shown in Red. DataCenter Index node in the center points to all Data Center nodes. The Type Index node (Production, Validation, and Development) connect to their respective Data Center types.

Asset nodes are Proxy Nodes, containing on the Asset ID.

## Data Center Grammars

### Define Data Center Index Node
DataCentersDS:DataCenters@I()

### Define Data Center Node with Asset and FISMA Multifields
DataCntDS:DataCenter(Name{DataCenter, Type, Asset@;, FISMA@;})

### Define Data Center Type Index Node
DataCentersDS:Type@I(Type)

## Data Center Data

| | DataCenter | Name | Type | Asset | FISMA |
|---|---|---|---|---|---|
| 1 | | | | | |
| 2 | DC001 | Bedford DCA-Prd | Production | Asset-A500 | FISMA-001; FISMA-004; FISMA-008 |
| 3 | DC002 | Bedford DCB-Prd | Production | Asset-B500 | FISMA-002;FISMA-008 |
| 4 | DC003 | Bedford DCA-Val | Validation | Asset-A400 | FISMA-003;FISMA-003;FISMA-002 |
| 5 | DC004 | Bedford DCB-Val | Validation | Asset-B400 | FISMA-004;FISMA-010;FISMA-003 |
| 6 | DC005 | Bedford DCA-Dev | Development | Asset-A100;Asset-A200 | FISMA-005;FISMA-001 |
| 7 | DC006 | Bedford DCB-Dev | Development | Asset-B100;Asset-B200;Asset-B300 | FISMA-004;FISMA-010 |

# Asset Subgraph

## Asset Graph



Node Properties

Asset

| | |
|---|---|
| <id> | 39 |
| Asset | Asset-B800 |
| AssetName | Unix OS |
| AssetType | Server |
| BiosGuid | S300-BG |
| Class | Virtual |
| DeviceType | Software |
| HostName | Host-B300 |
| IP_Address | 123.123.123.006 |
| IP_Address_Type | IPV4 |
| Vendor | DELL |
| Versions | U5.8 |
| nType | State |

## Asset Subgraph

### Define Asset Index Node

```
AssetsDS:Assets@I()
```

### Define Vendor, AssetType, ClassType, DeviceType Index Nodes

```
AssetsDS:Vendor@I(Vendor)
AssetsDS:AssetType@I(AssetType)
AssetsDS:ClassType@I(Class)
AssetsDS:DeviceType@I(DeviceType)
```

### Define Asset Node

```
AssetsDS:Asset(Asset{AssetName, Vendor, Class, AssetType, DeviceType,
          BiosGuid, Versions, HostName, IP_Address_Type, IP_Address})
```

## Asset Data

| | Asset | AssetName | Vendor | AssetType | Class | DeviceType | BiosGuid | Versions | HostName | IP_Address_Type | IP_Address |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Asset | AssetName | Vendor | AssetType | Class | DeviceType | BiosGuid | Versions | HostName | IP_Address_Type | IP_Address |
| 2 | Asset-A100 | Dell Server | Dell | Server | Virtual | Hardware | S100-BG | V3.2.4 | Host-A100 | IPV6 | 123.123.123.001 |
| 3 | Asset-A200 | Windows Office | Microsoft | Desktop | Physical | Software | D200-BG | C-2-5 | Host-A200 | IPV4 | 123.123.123.002 |
| 4 | Asset-A300 | Phone Systems | Plantronix | Phone | Physical | Com | PH100-BG | V8.3.1 | Host-A300 | IPV6 | 123.123.123.003 |
| 5 | Asset-A400 | Power Generator | PowerHouse | Power | Physical | Infra | PW100-BG | C2 | Host-A400 | IPV4 | 123.123.123.004 |
| 6 | Asset-A500 | Dell Server | Dell | Server | Virtual | Hardware | S200-BG | V3.4 | Host-A500 | IPV6 | 123.123.123.005 |
| 7 | Asset-B100 | Power Generator | PowerHouse | Power | Physical | Infra | PW200-BG | C2 | Host-B100 | IPV4 | 123.123.123.006 |

Assets have a single Asset Index node pointing to all Assets

Assets have multiple Index nodes each describing various Asset characteristics

# Vulnerability Subgraph

## Vulnerabilty Subgraph



### Node Properties

**Vulnerability**

| | |
|---|---|
| <id> | 4 |
| AssetName | Phone Syst |
| Description | Static on ex |
| DeviceType | Com |
| DiscoveryDate | 9/23/2020 |
| Effects | Degraded |
| Name | VName003 |
| RiskLevel | High |
| Vulnerability | Vuln003 |
| nType | State |
| name | Vulnerabili |

## Vulnerability Grammars

### Define Vulnerability Index Node

VulnerabilitiesDS:Vulnerabilities@I()

### Define RiskLevel, DeviceType, Effects, Vendor Index Nodes

VulnerabilitiesDS:RiskLevel@I(RiskLevel)
VulnerabilitiesDS:DeviceType@I(DeviceType)
VulnerabilitiesDS:Effects@I(Effects)
VulnerabilitiesDS:Vendor@I(Vendor)

### Define Vulnerabilty Node

VulnerabilitiesDS:Vulnerability(Vulnerability{Name, DiscoveryDate,
        Description, DeviceType, AssetName, Effects, RiskLevel,
        DeviceType, Versions@;, ImplemImpact@;})
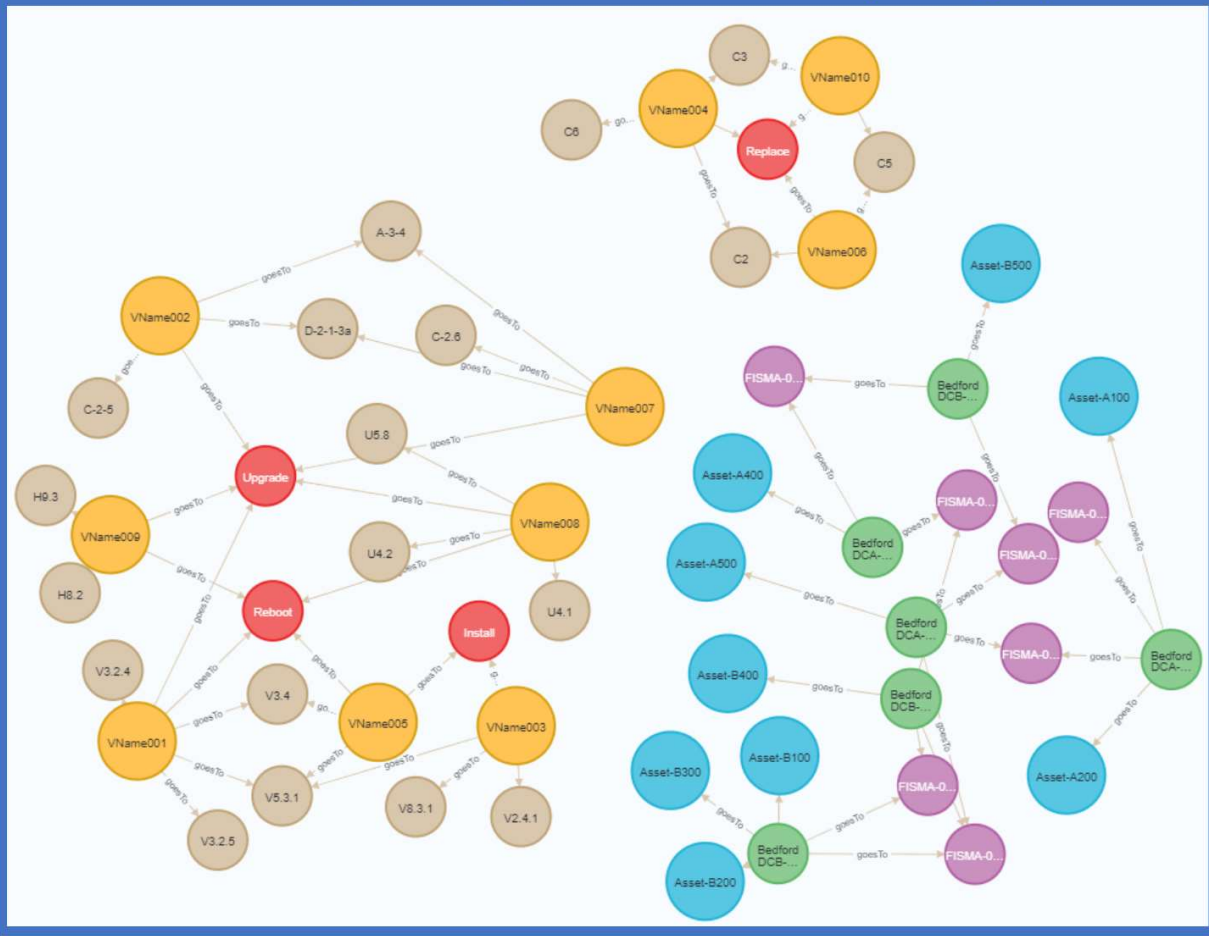
Vulnerability 003 and 002 are at High Risk

Vulnerability 001 and 002 require an Upgrade. 001 also requires a Reboot, etc.

## Vulnerability Data

| | A | B | C | D | E | F | G | H | I | J | K |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Vulnerabili | Name | DiscoveryDate | Description | DeviceType | AssetName | Effects | RiskLevel | Vendor | ImplemImpact | Versions |
| 2 | Vuln001 | VName001 | 3/15/1998 | Network Failures | Hardware | Dell Server | Disabled | High | Dell | Upgrade;Reboc | V3.2.4; V3.2.5; V3.4;V5.3.1 |
| 3 | Vuln002 | VName002 | 4/5/2021 | Performance Impact | Software | Windows Office | Degraded | Medium | Microsoft | Upgrade | A-3-4; C-2-5; D-2-1-3a |
| 4 | Vuln003 | VName003 | 9/23/2020 | Static on external conn | Com | Phone Systems | Degraded | Low | Photonix | Install | V5.3.1;V2.4.1;V8.3.1 |
| 5 | Vuln004 | VName004 | 3/6/2001 | Reduced Power Outpu | Infra | Power Generator | Degraded | Medium | Powerwerx | Replace | C6;C2;C3 |
| 6 | Vuln005 | VName005 | 5/10/2020 | Device access disabled | Hardware | Dell Server | Disabled | High | Dell | Install;Reboot | V3.4;V5.3.1 |
| 7 | Vuln006 | VName006 | 8/12/2019 | Faulty oputput voltage | Infra | Power Generator | Disabled | High | Powerwerx | Replace | C2;C5 |
| 8 | Vuln007 | VName007 | 7/25/2004 | Network Failures | Software | Windows Office | Disabled | Medium | Microsoft | Upgrade | A-3-4; C-2.6; D-2-1-3a |
| 9 | Vuln008 | VName008 | 7/10/2005 | Performance Impact | Software | Unix OS | Degraded | Low | IBM | Upgrade;Reboc | U4.2;U4.1;U5.8 |
| 10 | Vuln009 | VName009 | 9/3/2010 | Network Failures | Hardware | IBM Server | Disabled | Medium | IBM | Upgrade;Reboc | H9.3;H8.2 |
| 11 | Vuln010 | VName010 | 10/4/2014 | Voice stream interrupt | Com | Phone Systems | Degraded | Low | Photonix | Replace | C3;C5 |

# Combined Master Graph

## Data Center, Assets, and Vulnerability Subgraphs



The Master Graph contains all three subgraphs, and describes the Data Centers, Assets, and Vulnerabilities.

Assets (in Blue) are not connected to Vulnerabilities (in Gold) as shown in the Master Graph and the following Cypher Query.

### Show Asset-Vulnerability Connections

```
neo4j$ match (a:Asset)-[]→(b:Vulnerability) return a, b
```
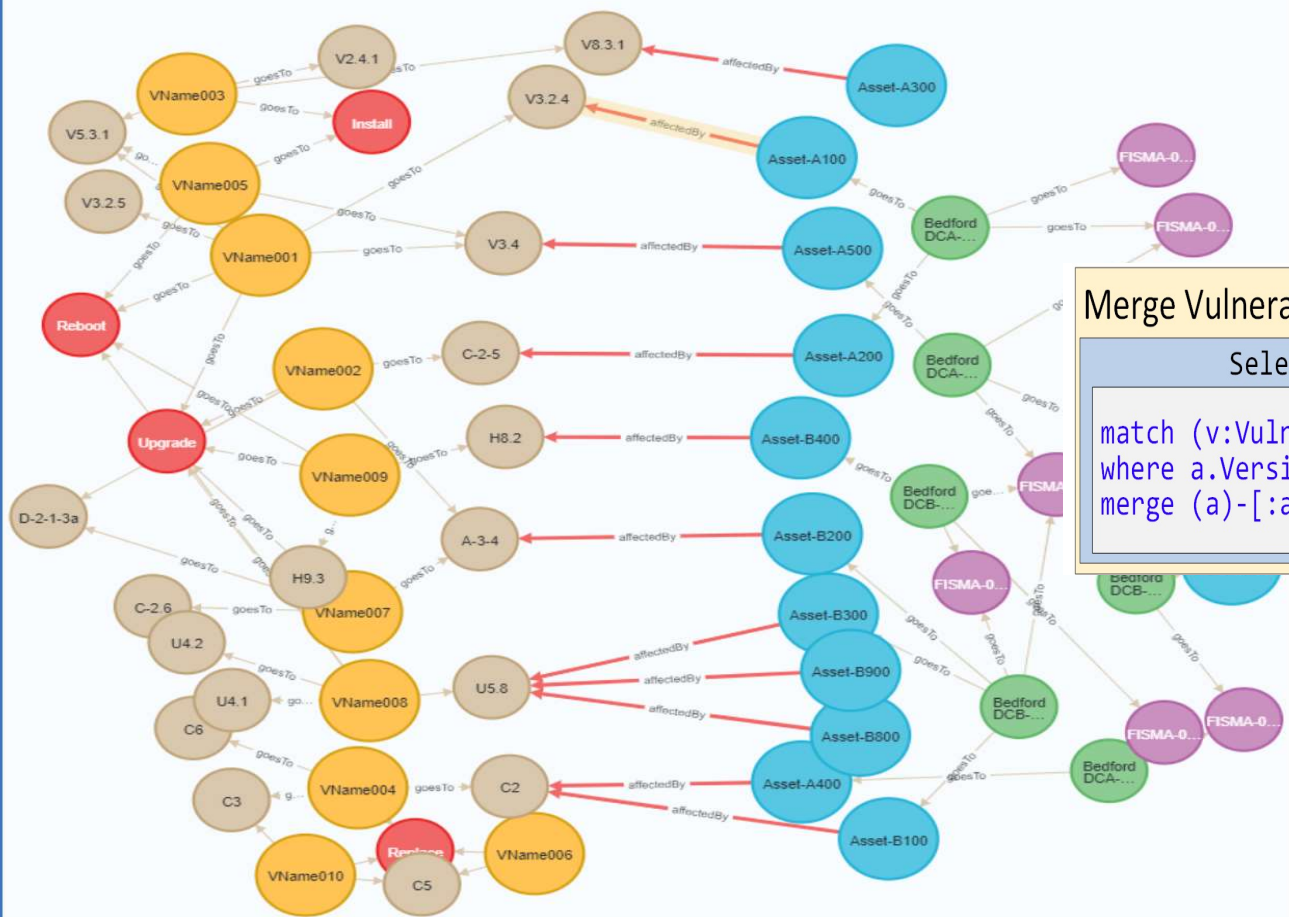
Table    (no changes, no records)

The underlying data sources have no such connections.

# Merging Vulnerability and Assets

## Merged Vulnerabilities and Assets



Standard Cypher, the language of graphs, as shown below, can be used to bind these entities together. The Bloom interface also provides this capability.

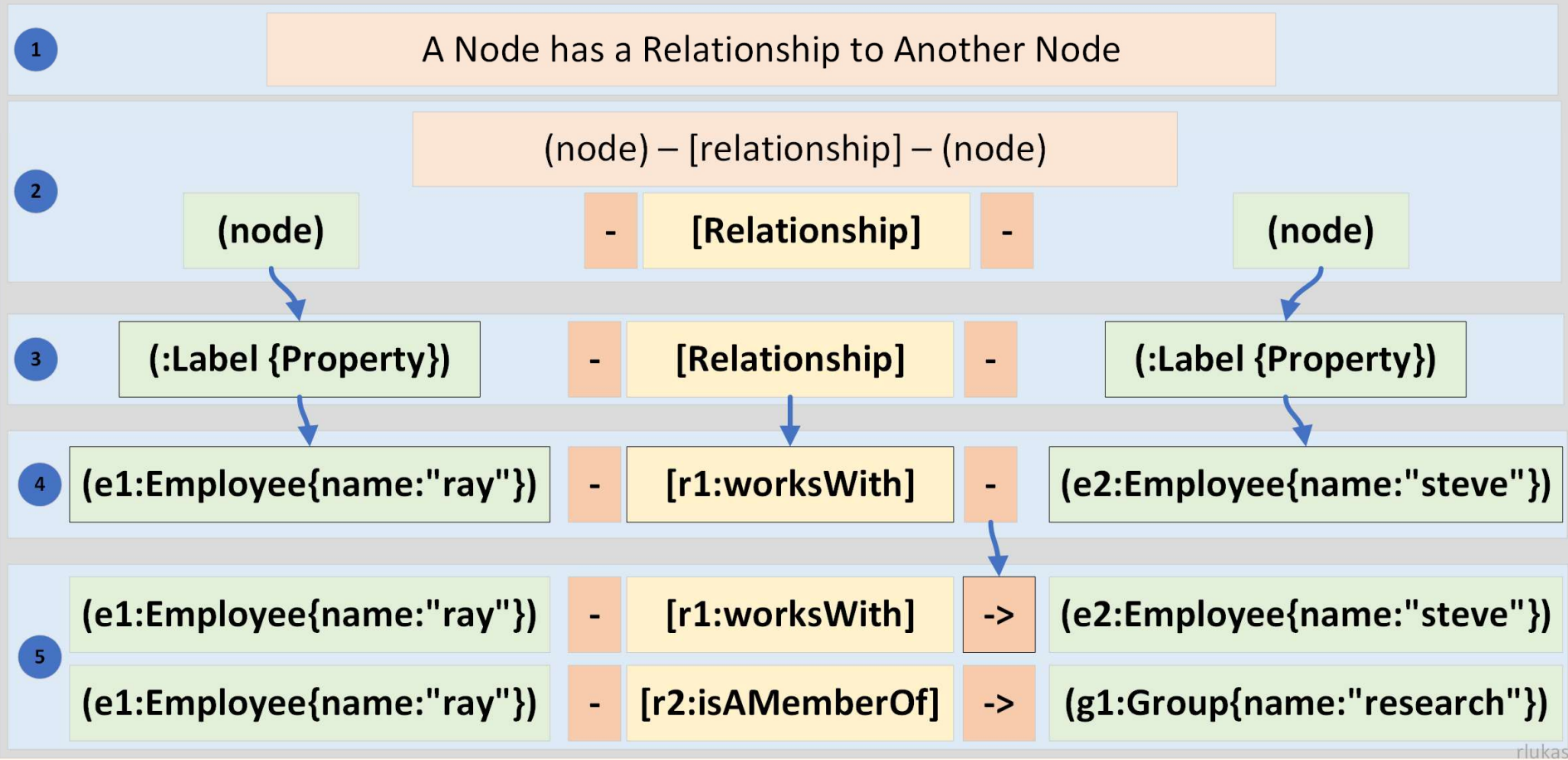### Merge Vulnerabilities and Affected Assets Together

Select matching nodes and build relationships

```
match (v:Vulnerability)-[gt:goesTo]->(ve:Versions), (a:Asset)
where a.Versions = ve.Versions and v.DeviceType = a.DeviceType
merge (a)-[:affectedBy]->(ve)
```

We can also see that some Assets, A330, B800, and B900 for example, are not bound in our Data Sets to a Data Center.
This rogue assets might require further exploration.

**MITRE**

Approved for Public Release; Distribution Unlimited. Public Release Case Number 22-2914

28

# Introduction to Cypher Patterns

## Understanding Patterns

ASCII Art

**1** A Node has a Relationship to Another Node

**2** (node) – [relationship] – (node)

(node)     -     [Relationship]     -     (node)

**3** (:Label {Property})     -     [Relationship]     -     (:Label {Property})

**4** (e1:Employee{name:"ray"})     -     [r1:worksWith]     -     (e2:Employee{name:"steve"})

**5**
(e1:Employee{name:"ray"})     -     [r1:worksWith]     ->     (e2:Employee{name:"steve"})

(e1:Employee{name:"ray"})     -     [r2:isAMemberOf]     ->     (g1:Group{name:"research"})

rlukas

**MITRE**

Approved for Public Release; Distribution Unlimited. Public Release Case Number 22-2914

29

# Generated Cypher Code

## Load Nodes and Linkages Mechanism

```
①:auto USING PERIODIC COMMIT 300   load ②csv with headers from ③"file:///dataSet.csv" as ④dataSetRow

   with ⑤dataSetRow

⑥ForEach(_ In Case When ((dataSetRow.dataCenterColName Is  not  Null)) Then [1] Else [] End|
       Merge (dataCenterInstance:DataCenter{name:"DataCenter",
   ⑦                              dataCenterID:trim(dataSetRow.datacenter_id_derived)}))

       set dataCenterInstance.nType = trim("State")⑧

           ForEach(_ In Case When ((dataSetRow.ServerID Is  not  Null)) Then [1] Else [] End|
               Merge (serverInstance:Server{name:"Server", ServerID:trim(dataSetRow.ServerID)})
               set serverInstance.macaddress = trim(dataSetRow.macAddress)
           ⑨ set serverInstance.dnsname = trim(dataSetRow.dnsname)
               set serverInstance.nType = trim("State")

               merge (dataCenterInstance)-[dataCenterToServerLink:contains]->(serverInstance)

           ) //end For ((dataSetRow.ServerID Is  not  Null)) command

⑥) //end For ((dataSetRow.dataCenterColName Is  not  Null)) command
```

**MITRE**

Approved for Public Release; Distribution Unlimited. Public Release Case Number 22-2914

30

# Generated Cypher Code

## Multifield Mechanism

```
:auto USING PERIODIC COMMIT 300  load csv with headers from "file:///newVersion.csv" as osRow
    with osRow

①  ForEach(_ In Case When ((osRow.ServerID Is  not  Null)) Then [1] Else [] End|
     Merge (serverInstance:Server{name:"Server", ServerID:trim(dataSetRow.ServerID)})

    ForEach(_ In Case When (osRow.os is not null) Then [1] Else [] End|
②    ForEach(osItem in split(osRow.os, ';') Then [1] Else [] End|
③      ForEach(_ IN CASE WHEN (osItem  <> '') Then [1] Else [] End|
④        merge (osInstance:OS{name:"OS", operSys:trim(osItem)})
⑤        set osInstance.nType = trim("Index")
          merge (serverInstance)-[:isA]->(osInstance)

        )

      )

    )

  )
```

**MITRE**

MITRE Labs inspires breakthroughs in applied science and advanced technology to transform the future of U.S. scientific and economic leadership. Our goal: Deliver disruptive innovation to support our mission of solving problems for a safer world.

**Please feel free to contact me with any questions**

**Thank You**

**Ray Lukas (rlukas@mitre.org)**
**Jim Lockett (jdlockett@mitre.org)**

**rlukas@MITRE.org**