

# KNOWLEDGE GRAPHS

## Lecture 10: Property Graphs

Markus Krötzsch  
Knowledge-Based Systems

TU Dresden, 17th Dec 2019

## Review

### Datalog is a general language for recursive, relational queries

- Easy to adopt to graphs (edges = special relations)
- Plain Datalog is a “pure” paradigm without the technical extensions of real query languages (esp. data types, filters)
- Adding negation is useful, but the interplay with recursion must be limited
- Adding aggregation must consider similar issues

### VLog4j can be used to issue Datalog queries:

- Data can be loaded from various sources, including SPARQL endpoints
- Stratified negation is supported (~)
- Syntax inspired by RDF and SPARQL to work with IRIs and datatype literals

## Review

## Property Graphs

# What is a Property Graph?

**Property Graph** is a type of graph, which can be described as follows:

- **directed** (edges have source and target vertices)
- **vertex-labelled** (for some kind of “label”)
- **edge-labelled** (for some kind of “label”)
- **multi-graph** (several versions of the exact same edge may exist)
- **with self-loops** (vertices can have edges to themselves), and
- **with sets of attribute-value pairs** associated with any vertex or edge

## Obvious questions

**Many issues require further specification:**

- What are the ids for vertices and edges?
- What are “attributes”?
- What are “values”? Which datatypes are supported? How are they defined?
- What are those “labels” that one can use for edges?
- What are those “labels” that one can use for vertices?
- If vertices and edge can have arbitrary attribute-value pairs, why do we also need labels?

Unfortunately, Property Graph as such is not an answer to all such questions:

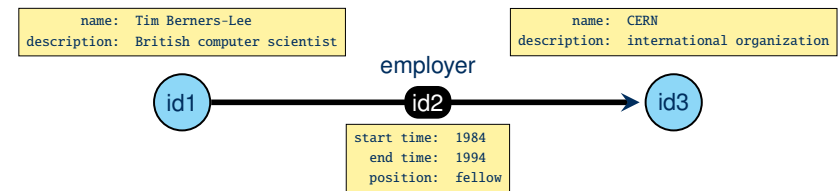
The name “Property Graph” refers to a **broad class of enriched graph structures**, allowing for **many technical interpretations** in different software systems. These interpretations are often **incompatible** and **based on different assumptions**.

# Example

## Tim Berners-Lee (Q80)

British computer scientist [edit](#)  
TimBL | Sir Tim Berners-Lee | Timothy John Berners-Lee | TBL | Tim Berners Lee | T. Berners-Lee | T Berners-Lee | Tim Berners-Lee | T.J. Berners-Lee

employer	CERN	<a href="#">edit</a>
	start time	1984
	end time	1994
	position held	Fellow
0 references		<a href="#">+ add reference</a>



## Types of “Property Graphs”: Object Model

The name “Property Graph” primarily hints at the attribute-value pairs (called “properties”) that can be associated with nodes and edges. There are different ways to interpret this model when designing actual data structures.

### View 1: Property Graph as an Object Model

“If you have ever worked with an object model or an entity-relationship diagram, the labeled property graph model will seem familiar.”  
– Neo4j, <https://neo4j.com/developer/guide-data-modeling/>

- Graphs viewed as data-modelling API in a programming language (often Java)
- “Values” are could be any other objects that represent data in programming
- Programmatic data access approaches are preferred over query language services
- Examples: Apache TinkerPop/Gremlin<sup>1</sup>, Neo4j/Cypher, multi-model object databases (e.g., Azure Cosmos DB, OrientDB, Oracle Spatial and Graph),

<sup>1</sup> Many graph DBMS have TinkerPop bindings, but TinkerPop’s native view is object-based.

## Types of “Property Graphs”: Relational Database Extension

The name “Property Graph” primarily hints at the attribute-value pairs (called “properties”) that can be associated with nodes and edges.  
There are different ways to interpret this model when designing actual data structures.

### View 2: Property Graph as Extended Relational Model

“Vertex attributes match to columns of the vertex table. Edge attributes match to columns of the edge table. The maximum number of attributes is bound by the maximum number of columns for the underlying tables”

– SAP HANA Graph Reference, v1.0 – 2016-11-30

- Graphs viewed as indexing and access layer on top of RDBMS
- “Values” can be any values in standard or proprietary SQL datatypes
- Various data access paradigms: in-DB-code (e.g., SAP Hana GraphScript) or query language (e.g., Tigergraph GSQL)
- Examples: SAP Hanah Graph, Tigergraph<sup>1</sup>

<sup>1</sup>Not based on a full RDBMS, but strongly using RDB concepts, rigid schema.

## Types of “Property Graphs”: RDF-Database Extension

The name “Property Graph” primarily hints at the attribute-value pairs (called “properties”) that can be associated with nodes and edges.  
There are different ways to interpret this model when designing actual data structures.

### View 3: Property Graph as Access Layer on Top of RDF

“property graph data can be loaded and accessed via the TinkerPop3 API, but underneath the hood the data will be stored as RDF” – Blazegraph TinkerPop3

- Property Graphs stored internally as RDF graphs
- “Values” can be any values supported in RDF (XML Schema + proprietary extensions)
- Multiple access paradigms: Apache Tinkerpop Gremlin or SPARQL
- Examples: Amazon Neptune, Stardog, BlazeGraph

## Types of “Property Graphs”: Summary

### Main approaches:

- **Object databases:** flexible, schema-less; often multi-model; includes many graph extensions of noSQL DBMS; varying datatypes and formats (e.g., JSON for many object DBs)
- **Relational databases:** rather rigid, schema-based; graph extensions of classical RDBMS; SQL datatypes
- **RDF databases:** flexible, schema-less, highly normalised (data atomised into triples); property graph extensions of RDFDBs; RDF datatypes

### Other types of graph databases:

- Simpler graph models (neither RDF nor property graph); mostly for network analysis; e.g., Apache Giraph
- Based on other paradigms; e.g., AllegroGraph (RDF database with Prolog support)
- Combinations and specialised components/frameworks; e.g., some data stored in Lucene or Solr (for text search), exchangeable storage back-end

## Property graph: data access

Methods for data access are as diverse as the data structures that are used.

### Programmatic access:

- Proprietary or common APIs (mostly Tinkerpop)
- Scripting and processing languages (e.g., Apache Gremlin, SAP GraphScript)
- MapReduce, Spark, and other processing frameworks

### Query languages:

- Neo4j Cypher
- Oracle PGQL
- Tigergraph GSQL
- ...

## Property Graph: What to teach?

The current implementation chaos has prompted some activities:

- **OpenCypher** publishes a specification for some parts of Neo4j's Cypher query language
- The Linked Data Benchmark Council has proposed a merger of PGQL, Cypher, and some Gremlin features, called **G-CORE** (not implemented yet)
- A current push towards a unified "(Property) Graph Query Language" **GQL** is underway (ISO standardisation has been started; optimistic ETA: late 2021)

→ In this course: focus on (Open) Cypher + discussion of proposed changes

## Identity and labelling

**Nodes** and **relationships** have independent identity

- Especially: two **relationships** can be indistinguishable in terms of data (same source, target, label, **properties**) and yet be different
- Identity conferred by identifiers, which are implementation specific

Labels are based on strings:

- Vertex labels are sets of unicode strings, each called **label**
- Edge labels are single unicode strings, called **relationship types**
- Both vertices and edges may have no labelling

In practice, **labels** are used to take the role of types or classes, e.g., one may have a **label** person used for all **nodes** representing people.

**Relationship types** play the role of RDF properties, denoting the type of relationship that an edge expresses.

## A note on terminology

Unfortunately, the OpenCypher/Neo4j world uses completely different names for concepts than the RDF world:

<b>Graph-theoretic concept</b>	<b>OpenCypher terminology</b>
vertex	<b>node</b>
edge	<b>relationship</b>
vertex label	a set of strings, each of which is called <b>label</b>
edge label	<b>relationship type</b>
key-value pair	<b>property</b>

Terms are shown in a **distinct style** to clarify this special meaning

## Properties

Sets of **properties** can be assigned to vertices and edges.

**Properties** are key-value pairs:

- **Property keys** are unicode strings
- **Property values** are either concrete values of some datatype, or lists of values of the same datatype

**Property keys** must be unique in a set of **properties** used on some **node** or **relationship**, but lists can be used to encode several values.

## Datatypes

**Supported datatypes** for **property values** in OpenCypher:

- INTEGER: “exact numbers without decimals” (apparently of arbitrary magnitude)
- FLOAT: double precision (64bit) floating point numbers
- STRING: unicode strings
- BOOLEAN: true or false
- lists of values of the above

**Missing datatypes** of much practical importance

- dates and times
- geographic coordinates
- fine-grained numeric types

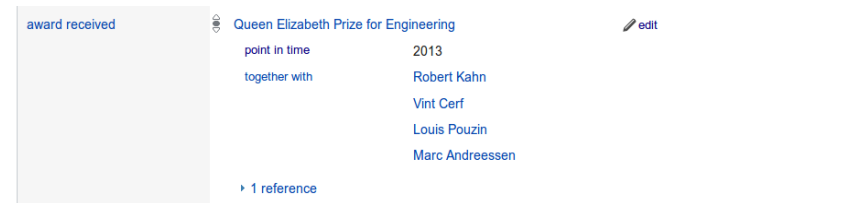
↪ might be supported as proprietary extensions in implementations

## Anything missing?

Another major omission and one of the biggest shortcomings of Property Graph models:

**Property values** cannot be references to vertices or edges.

Property Graph in this form is therefore not suitable to model, e.g., Wikidata statements:



This information could be captured in a Property Graph that looks like the RDF graph we used before (with seven **relations**), but not using any **properties** at all.

## From Property Graph to RDF

As expected, every **Property Graph** can be expressed as an **RDF graph**:

- Use reification to represent edges by auxiliary nodes
- Use special RDF properties to encode source and target of edges, and the **node-label** and **relationship-relationship type** association
- Use application-based RDF properties to encode **property keys**
- Depending on the exact Property graph implementation, use some appropriate datatype-to-RDF mapping (e.g., based on RDF2RDB mappings from SQL datatypes to RDF)

**Remarks:**

- This scheme is implemented natively in existing DBMS (Blazegraph, Amazon Neptune), and scales there despite of the increase size of the graph data that the system has to work with
- The use of reification can be avoided for **relationships** without **properties**
- One can also use both (reified and direct statements) for flexibility (similar to Wikidata's RDF encoding)

## From RDF to Property Graph

As expected, every **RDF Graph** can be expressed as a **Property Graph**:

- Represent all RDF resources by **nodes**
- Use **property keys** to associate resources with IRIs and/or datatype literal information
- Represent all RDF triples with auxiliary **nodes**
- Use **relationships** with special **relationship types** to associate auxiliary **nodes** with triple subject, predicate, and object

**Remarks:**

- There does not seem to be any simpler way of capturing the full power of RDF in Property Graph, due to the restrictions of the latter (no reference in **relationship types** or **property values** to **nodes**)
- In some cases, certain triples could be represented as **properties** (if their datatype is supported in Property Graph and we do not need their RDF property to be addressable in the graph)
- Many Property Graph implementations will have performance problems in handling graphs with so many edges (part of the motivation for moving data into **properties** is to reduce the graph size)

## (No) Schema Modelling

In RDF, properties were identified by IRIs and could be subject of triples

- to define labels and descriptions in several languages
- to specify the datatype for the property
- to relate it to other properties, e.g., to their inverse

**Example 10.1:** Wikidata describes properties on own pages, and allows them to be used in statements, references, or statement qualifiers.

In Property Graph, **labels**, **relationship types**, and **property keys** are plain strings

- They cannot occur in the graph
- They can have neither **relationships** nor **properties**

### Workarounds:

- One can create nodes that refer to a string token through a **property value**, and encode the knowledge that this is meant as a reference in application software
- Some database management system may support the declaration of **constraints** that restrict the usage of **labels**, **relationship types**, or **property keys**

## Cypher

## Cypher overview

Cypher is a (family of) query languages for Property Graph:

- Proprietary query language of the Neo4j graph database
- Subset supported by other tools as well: openCypher
- Might be an important input to future graph query language standards

openCypher supports two main functions:

- A query language
- An update language

Current specification of openCypher 9:

<https://s3.amazonaws.com/artifacts.opencypher.org/openCypher9.pdf>

Currently **not defined in openCypher v9** and depending on implementation:

- Parts of the query language (e.g., comparison and ordering of some values)
- Result formats for sending query results
- Protocol for sending queries and receiving answers

## Cypher queries

The heart of Cypher is its query language.

**Example 10.2:** The following Cypher query asks for a list of all nodes that are in an EMPLOYER relationship:

```
MATCH (person)-[:EMPLOYER]->(company)
RETURN person, company
```

This corresponds to the following SPARQL query:

```
SELECT ?person ?company
WHERE { ?person :EMPLOYER ?company }
```

### Basic concepts:

- Cypher uses **variables**, marked by their context of use
- The core of a query is the **query condition** within **MATCH** { ... }
- Conditions can be simple patterns based on graph edges in a custom syntax
- **RETURN** specifies how results are produced from query matches

## Basic Cypher by example

**Example 10.3:** Find up to ten people whose daughter is a professor:

```
MATCH
  (parent)-[:HAS_DAUGHTER]->(child {occupation:'Professor'})
RETURN parent
LIMIT 10
```

**Example 10.4:** Count all **relationships** in the database:

```
MATCH ()-[relationship]->()
RETURN count(relationship) AS count
```

**Example 10.5:** Count all **relationship types** in the database:

```
MATCH ()-[relationship]->()
RETURN count(DISTINCT type(relationship)) AS count
```

## Basic Cypher by example (2)

**Example 10.6:** Find the person with most friends:

```
MATCH (person)-[:HAS_FRIEND]->(friend)
RETURN person, count(DISTINCT friend) AS friendCount
ORDER BY friendCount DESC
LIMIT 1
```

**Example 10.7:** Find pairs of siblings:

```
MATCH
  (parent)-[:HAS_CHILD]->(child1),
  (parent)-[:HAS_CHILD]->(child2)
WHERE id(child1) <> id(child2)
RETURN child1, child2
```

## Basic Cypher by example (3): properties and labels

Queries can also access

- **Labels** (the additional strings used on **nodes**)
- **Properties** (of **nodes** and **relationships**)

**Example 10.8:** Find friends of all people with name Paul Erdős, and return their name and the start date of the friendship:

```
MATCH
  (:Human {name: 'Paul Erdős'})-[rel:HAS_FRIEND]->(friend:Human)
RETURN friend.name, rel.startDate
```

Here Human is a **label**, and name and startDate are **property keys**.

## Cypher: Outlook

**Cypher has many further features:**

- Filter expressions (similar to SPARQL)
- Some regular path expressions (less than SPARQL)
- Returning of paths and shortest paths (more than SPARQL)
- Groups and aggregates (largely as in SPARQL)
- Subqueries
- Unions (more limited than in SPARQL)
- Optional (similar to SPARQL)

See next lecture for further details ...

## Summary

Property Graph is a general concept for organising graph data in two layers: a primary graph layer and a sub-ordinate key-value-set layer

Property graph has many different, incompatible implementations, based on several database paradigms (object, relational, RDF graph)

Cypher is an influential query language for property graph

### What's next?

- Holidays
- Detailed overview of Cypher
- Quality control in knowledge graphs

