

# **Graph Databases & High-Performance Python**

## **From Data Models to HPC Triangle Analytics**

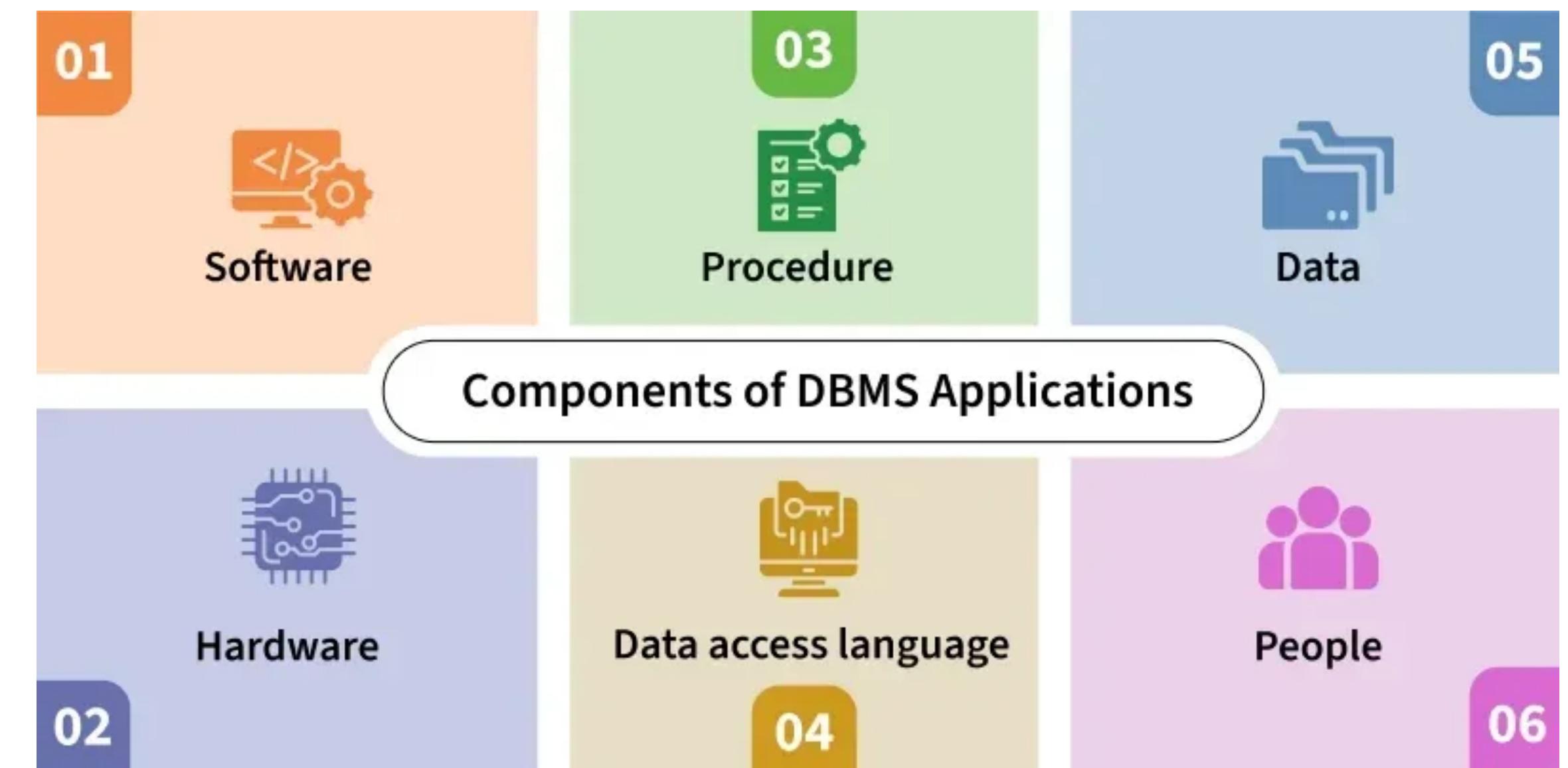
**Daria Kriuchkova**

# Agenda

- Graph Databases: Core Idea
- Relational vs Graph Models
- From ORM to OGM
- Python Objects to Graph Structures
- Limits of OGMS for Analytics
- HPC with Python on Graph Data
- Triangle Counting in  $G(n,p)$

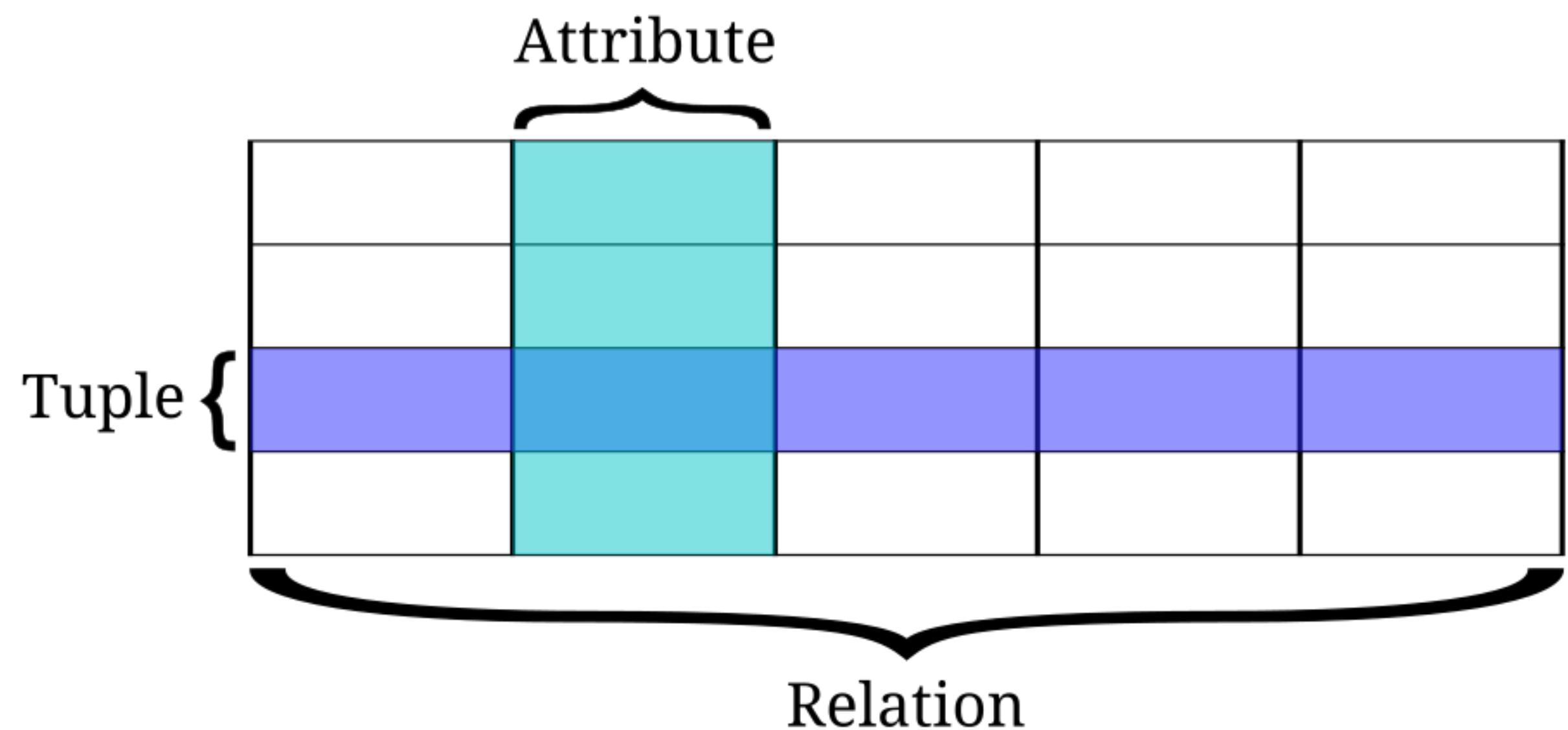
# What is a Database?

A **database** is an organised collection of data or a type of data store based on the use of a database management system (DBMS), the software that interacts with end users, applications, and the database itself to capture and analyse the data.



# Relational DBs

## RDBMS



DBMS engine:

- Implements **tables, rows, schemas**, and relational constraints.
- The query optimiser builds efficient execution plans for **joins**

# Relational Databases Lack Relationships

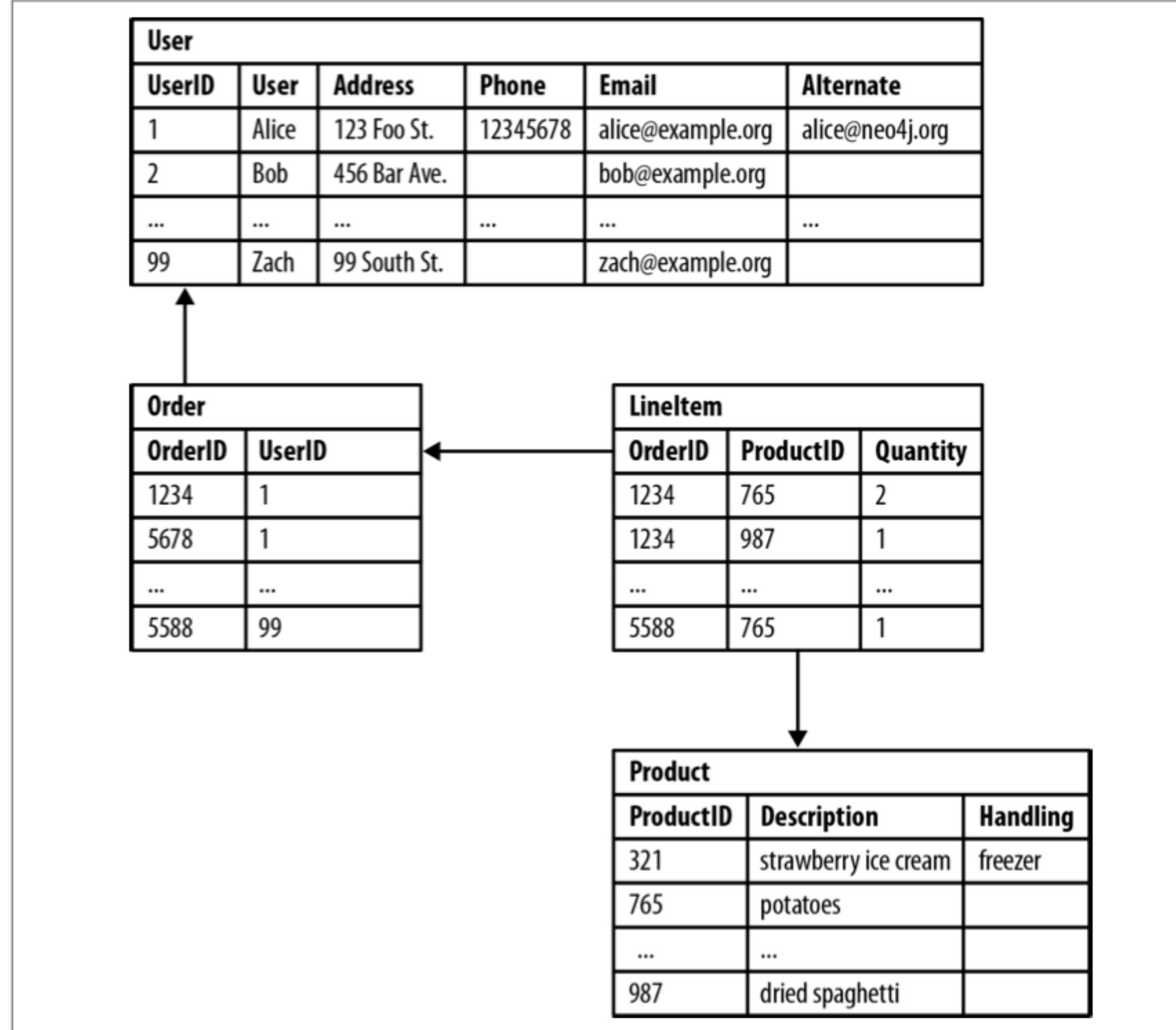


Figure 2-1. Semantic relationships are hidden in a relational database

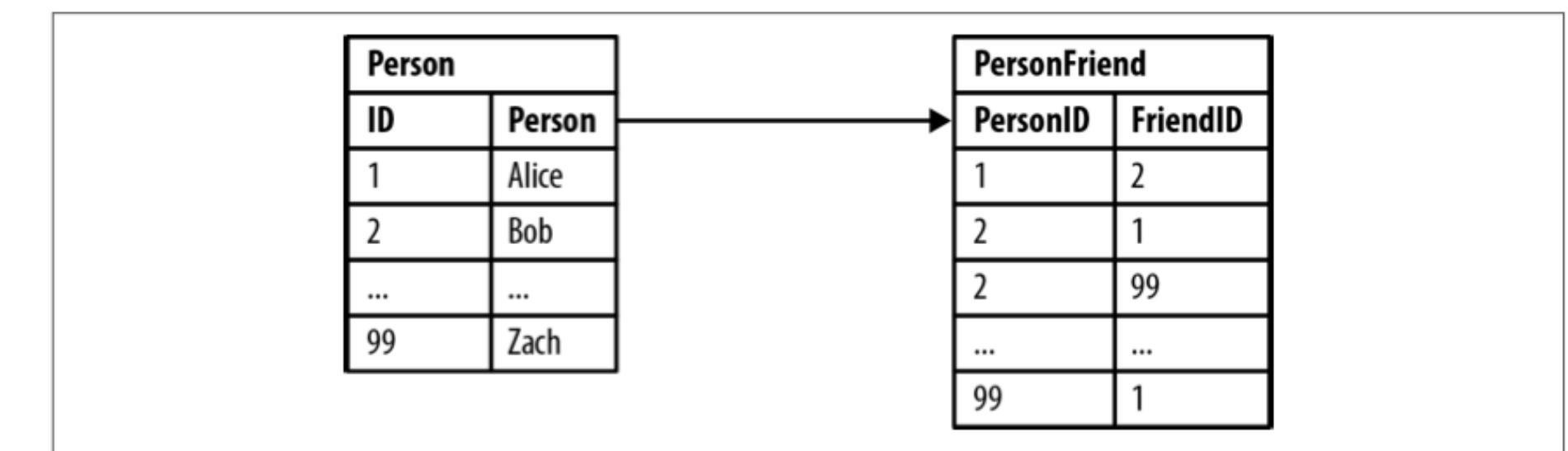


Figure 2-2. Modeling friends and friends-of-friends in a relational database

```
SELECT p1.Person
FROM Person p1 JOIN PersonFriend
  ON PersonFriend.FriendID = p1.ID
JOIN Person p2
  ON PersonFriend.PersonID = p2.ID
WHERE p2.Person = 'Bob'
```

# Graph Databases Relationships

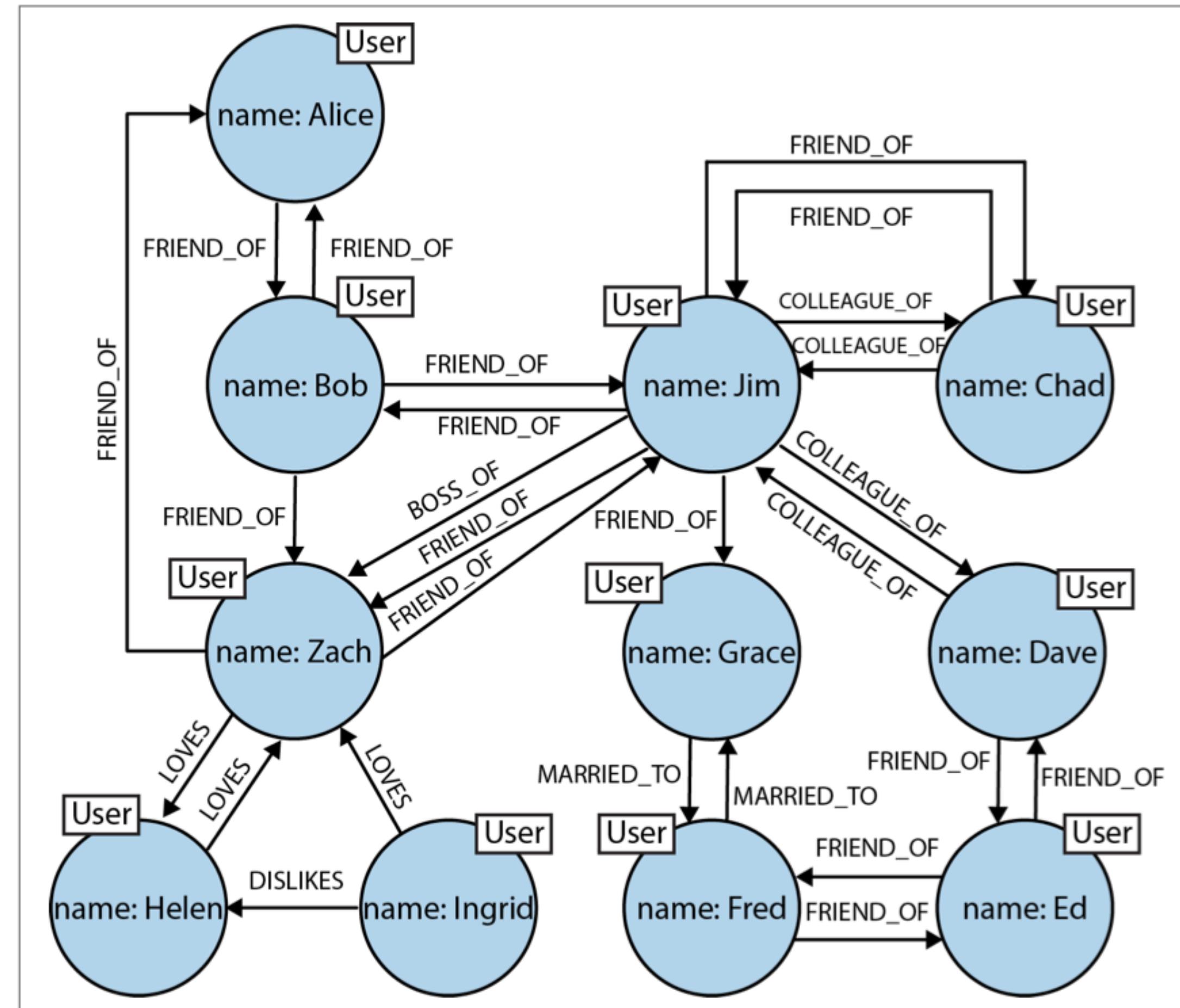
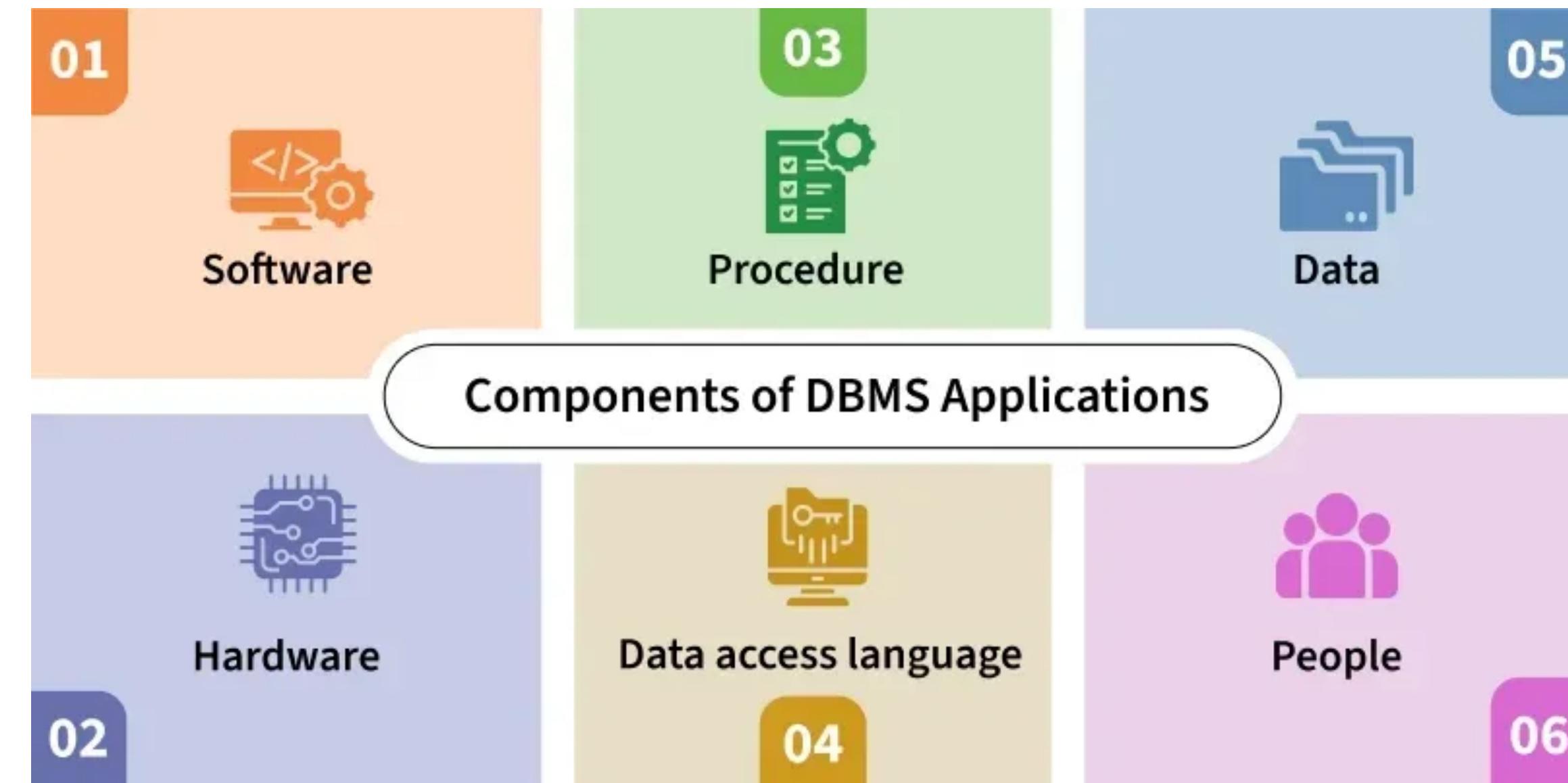
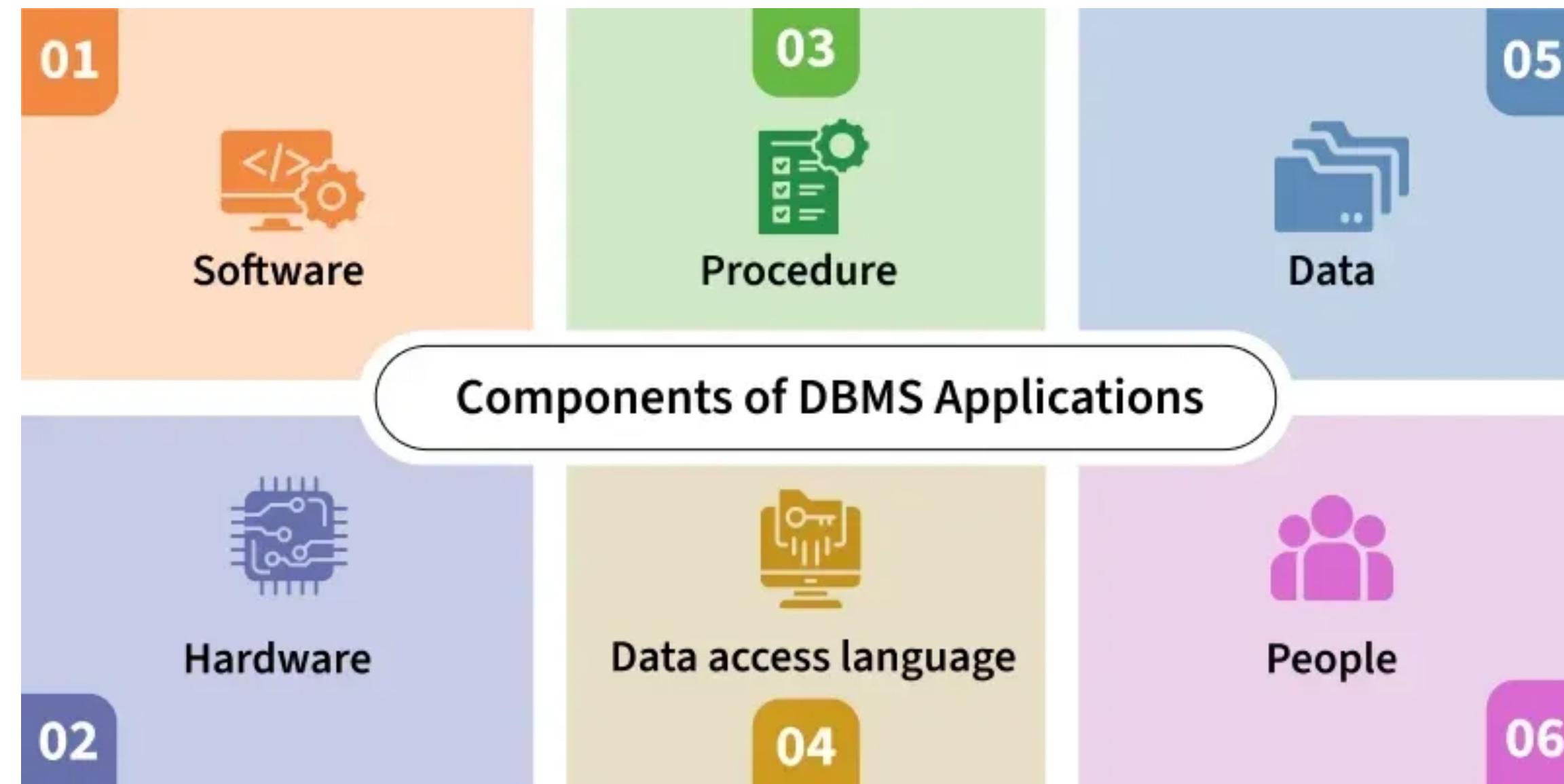


Figure 2-5. Easily modeling friends, colleagues, workers, and (unrequited) lovers in a graph



The database engine is built around **graph-native storage and processing**. It supports ACID transactions, persistence, recovery, and high availability, but is optimized for **connected data**, not tables.



The database engine is built around **graph-native storage and processing**. It supports ACID transactions, persistence, recovery, and high availability, but is optimized for **connected data**, not tables.

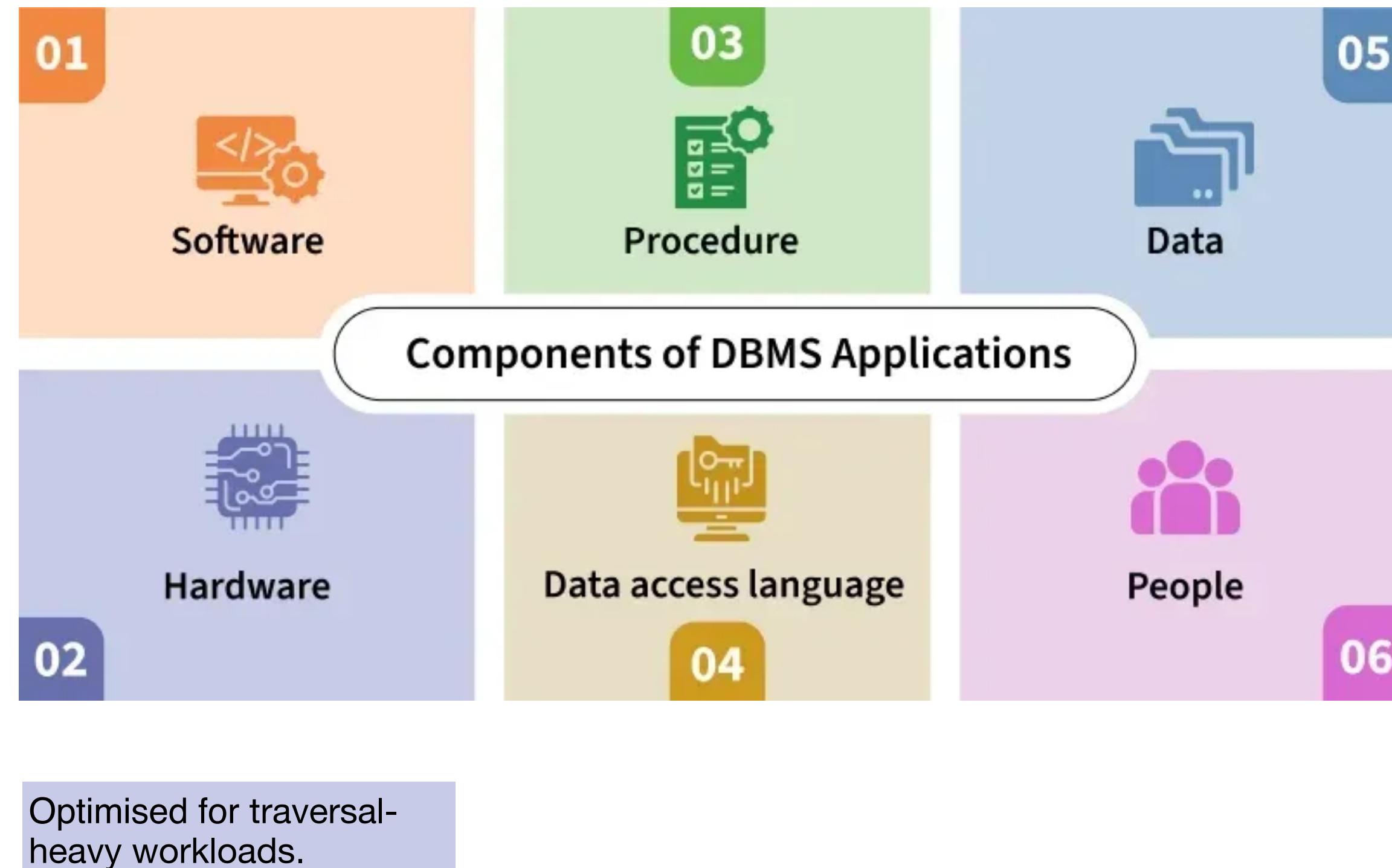


Figure 3-2. Simplified snapshot of application deployment within a data center

The database engine is built around **graph-native storage and processing**. It supports ACID transactions, persistence, recovery, and high availability, but is optimized for **connected data**, not tables.

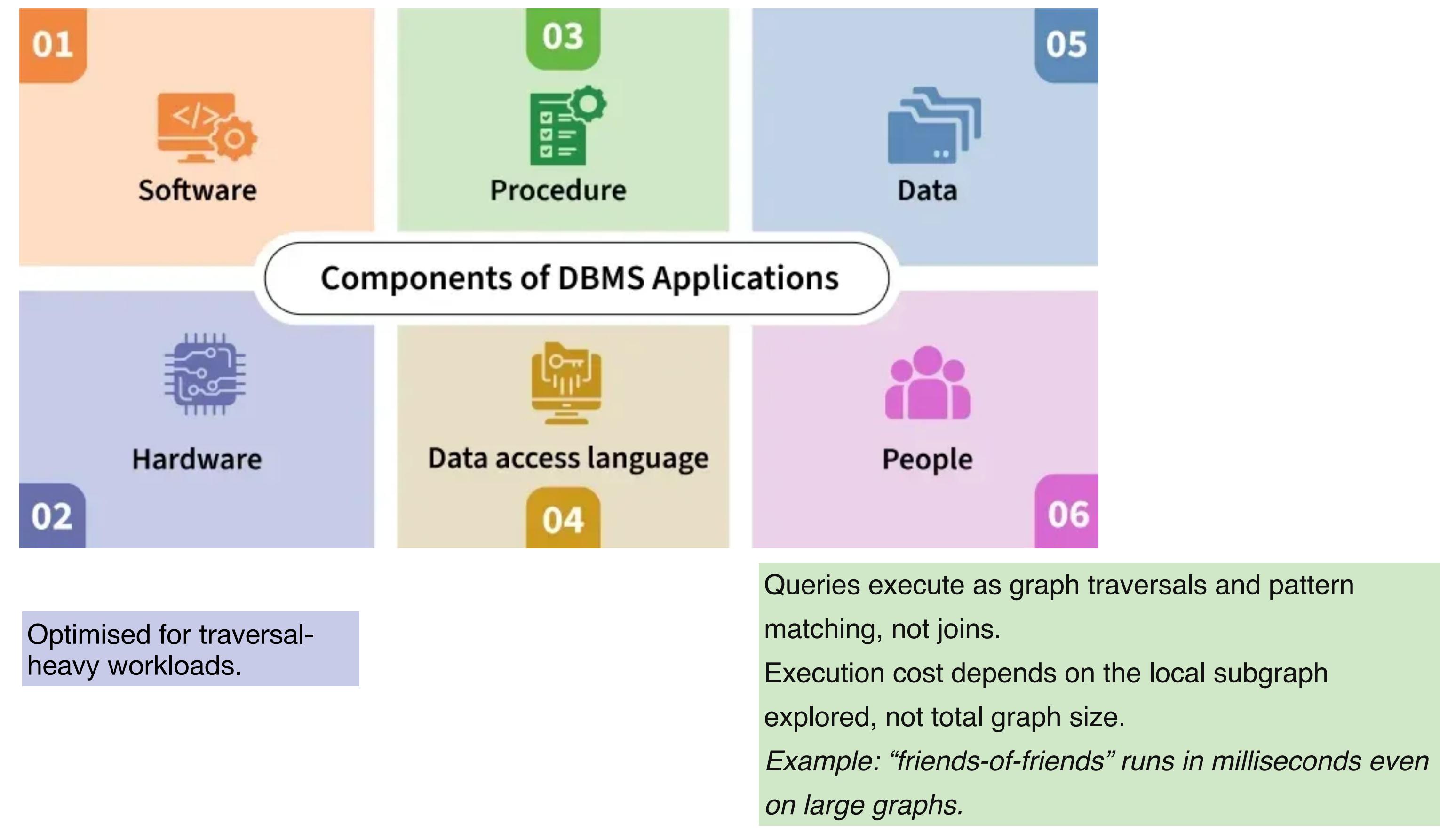
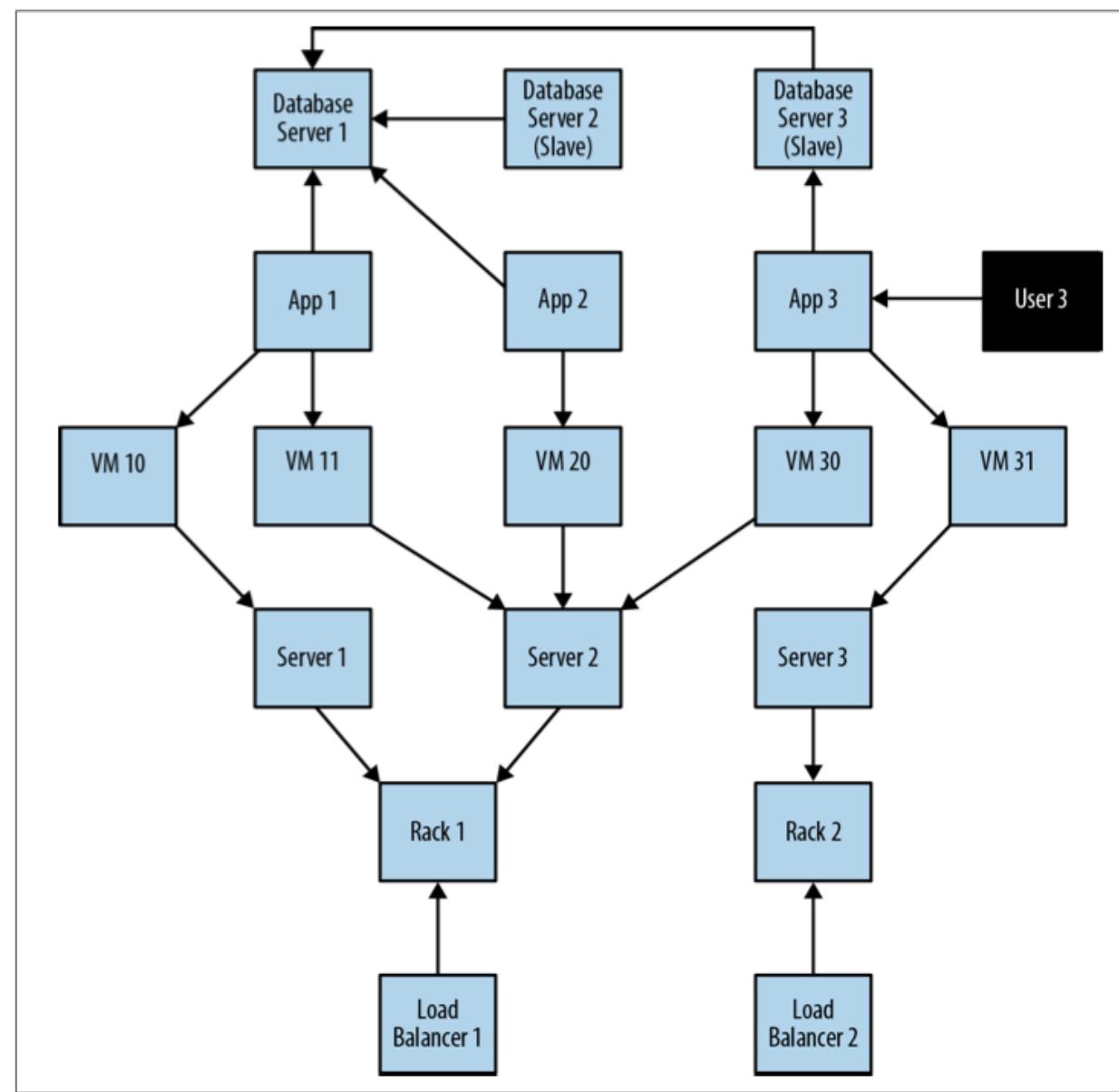
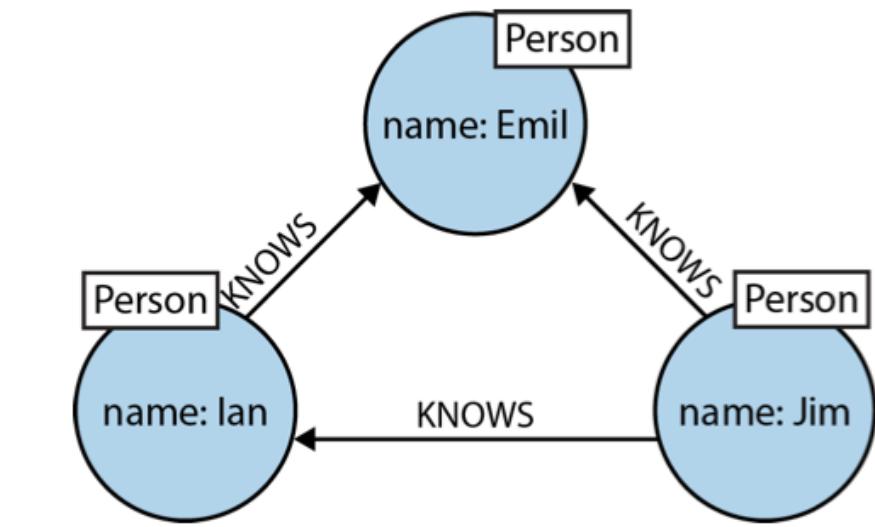


Figure 3-2. Simplified snapshot of application deployment within a data center

The database engine is built around **graph-native storage and processing**. It supports ACID transactions, persistence, recovery, and high availability, but is optimized for **connected data**, not tables.

A graph database stores data as a labeled property graph



Optimised for traversal-heavy workloads.

Queries execute as graph traversals and pattern matching, not joins.  
Execution cost depends on the local subgraph explored, not total graph size.  
*Example: “friends-of-friends” runs in milliseconds even on large graphs.*

Figure 3-2. Simplified snapshot of application deployment within a data center

# Graph Databases – Architectural Overview

storage model, query language, execution model, scalability



Neo4j



Memgraph



ArangoDB



JanusGraph

**native graph storage  
(nodes & relationships  
first-class)**

Cypher

disk-based,  
optimized traversals  
+ GDS (C++/native)

mature ecosystem, analytics +  
applications

**in-memory graph**

openCypher

real-time, streaming-  
friendly

fast subgraph extraction, tight  
Python workflows

**multi-model (graph +  
document + key-value)**

AQL

general-purpose  
engine

hybrid data, not a “pure” graph  
system

**distributed (on top of  
Cassandra / HBase /  
etc.)**

Gremlin

cluster-oriented

very large graphs, high  
operational complexity

# Object-Relational Mapping

solves “impedance mismatch”  
and converts Python objects into  
a database structure.

# Object-Relational Mapping

solves “impedance mismatch”  
and converts Python objects into  
a database structure.

**RDBMS:** Converts Python operations into SQL under the hood.

Python classes  $\leftrightarrow$  SQL tables  
Object attributes  $\leftrightarrow$  columns  
Object relationships  $\leftrightarrow$  foreign keys

```
class Person(Base):
    __tablename__ = "persons"
    id = Column(Integer, primary_key=True)
    name = Column(String)

class Friendship(Base):
    __tablename__ = "friendships"
    person_id = Column(ForeignKey("persons.id"), primary_key=True)
    friend_id = Column(ForeignKey("persons.id"), primary_key=True)
```

SQLAlchemy/Django

# Object-Graph Mapping

solves “impedance mismatch”  
and converts Python objects into  
a database structure.

Python classes ↔ graph nodes  
Object attributes ↔ node properties  
Python references ↔ graph edges

**RDBMS:** Converts Python operations into SQL under the hood.

Python classes ↔ SQL tables  
Object attributes ↔ columns  
Object relationships ↔ foreign keys

```
class Person(Base):  
    __tablename__ = "persons"  
    id = Column(Integer, primary_key=True)  
    name = Column(String)  
  
class Friendship(Base):  
    __tablename__ = "friendships"  
    person_id = Column(ForeignKey("persons.id"), primary_key=True)  
    friend_id = Column(ForeignKey("persons.id"), primary_key=True)
```

SQLAlchemy/Django

# Object-Graph Mapping

solves “impedance mismatch”  
and converts Python objects into  
a database structure.

Python classes ↔ graph nodes  
Object attributes ↔ node properties  
Python references ↔ graph edges

```
from neomodel import StructuredNode, StringProperty, RelationshipTo

class Person(StructuredNode):
    name = StringProperty()
    friends = RelationshipTo("Person", "FRIENDS_WITH")

alice = Person(name="Alice").save()
bob = Person(name="Bob").save()
alice.friends.connect(bob)

friends_of_alice = alice.friends.all()
```

Neomodel/Neo4

**RDBMS:** Converts Python operations into SQL under the hood.

Python classes ↔ SQL tables  
Object attributes ↔ columns  
Object relationships ↔ foreign keys

```
class Person(Base):
    __tablename__ = "persons"
    id = Column(Integer, primary_key=True)
    name = Column(String)

class Friendship(Base):
    __tablename__ = "friendships"
    person_id = Column(ForeignKey("persons.id"), primary_key=True)
    friend_id = Column(ForeignKey("persons.id"), primary_key=True)
```

SQLAlchemy/Django

# Object-Graph Mapping

solves “impedance mismatch”  
and converts Python objects into  
a database structure.

Python classes ↔ graph nodes  
Object attributes ↔ node properties  
Python references ↔ graph edges

```
from neomodel import StructuredNode, StringProperty, RelationshipTo

class Person(StructuredNode):
    name = StringProperty()
    friends = RelationshipTo("Person", "FRIENDS_WITH")

alice = Person(name="Alice").save()
bob = Person(name="Bob").save()
alice.friends.connect(bob)

friends_of_alice = alice.friends.all()
```

Neomodel/Neo4j

**RDBMS:** Converts Python operations into SQL under the hood.

Python classes ↔ SQL tables  
Object attributes ↔ columns  
Object relationships ↔ foreign keys

```
class Person(Base):
    __tablename__ = "persons"
    id = Column(Integer, primary_key=True)
    name = Column(String)

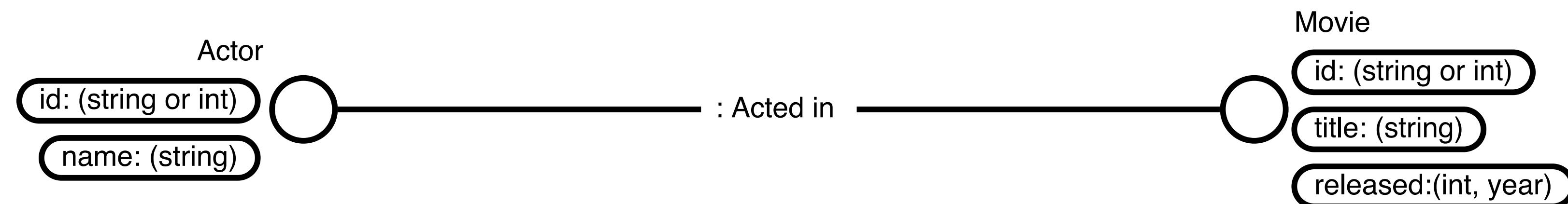
class Friendship(Base):
    __tablename__ = "friendships"
    person_id = Column(ForeignKey("persons.id"), primary_key=True)
    friend_id = Column(ForeignKey("persons.id"), primary_key=True)
```

SQLAlchemy/Django

- Relationships are **explicit objects**, not foreign-key fields.
- Traversal is done by walking edges, not joining tables.
- The graph model is much closer to the object model → *less mismatch*.

# How Python Maps Objects to Graph Structures

## Example Project: A Minimal Python Object–Graph Mapper



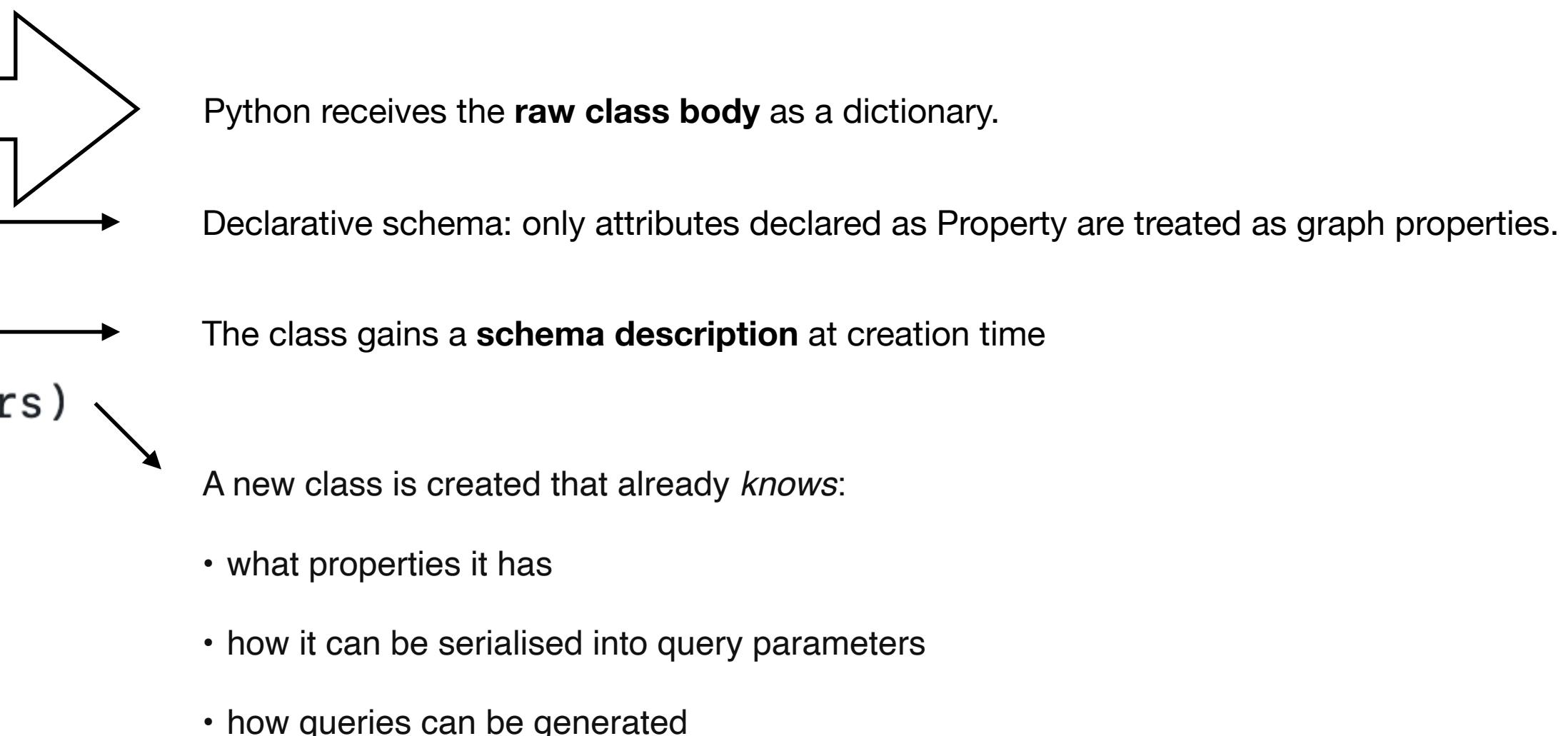
# How Python Maps Objects to Graph Structures

## Class Definition → Graph Schema

```
4  class NodeBase(type):
5      def __new__(mcs, name, bases, attrs):
6          properties = {}
7          for attr_name, attr in attrs.items():
8              if isinstance(attr, Property):
9                  properties[attr_name] = attr
10         attrs["_properties"] = properties
11         kls = type.__new__(mcs, name, bases, attrs)
12         return kls
13
14     def __init__(cls, name, bases, attrs):
15         super().__init__(name, bases, attrs)
16
```

```
4     class Movie:
5         def __init__(self, title: str, released: Optional[int] = None):
6             self.title = title
7             self.released = released
```

- **Node** → Movie (entity)
- **Properties** → title, released



# How Python Maps Objects to Graph Structures

## Python class → Metadata → Cypher

Collects values from \*\*kwargs and assigns them to declared properties

```
18 <  class Node(metaclass=NodeBase):
19
20 <      def __init__(self, /, **kwargs):
21          self.cached_properties = {}
22          for fn, prop in self._properties.items():
23              val = kwargs.pop(fn, None)
24              setattr(self, fn, val)
--
```

Validates input against the schema and extracts query parameters

```
42     @classmethod
43 <     def _get_params_from_kwargs(cls, **kwargs):
44         params = {}
45         for prop_name, prop in cls._properties.items():
46             params[prop_name] = kwargs.get(prop_name, None)
47         for k in kwargs:
48             if k not in cls._properties:
49                 raise AttributeError(f"'{k}' is not a valid property for model {cls}")
50     return params
```

# How Python Maps Objects to Graph Structures

## Python class → Metadata → Cypher

Compiles class metadata into a parameterized Cypher query

```
52     @classmethod
53     def build_create_query(cls, **params):
54         node_alias = "node"
55         label = cls.__name__  
The graph node label is derived directly from the Python class name. No manual mapping
56         query = f"CREATE {node_alias}:{label} "
Cypher is generated dynamically from Python metadata.
57
58         set_query = [
Only properties collected by the metaclass are used:
59             f"{node_alias}.{{k}}=${{k}}"
60             for k in cls._properties
61         ]
- schema consistency
62         set_query = ", ".join(set_query)
- no accidental fields
63         if set_query:
- safe query construction
64             set_query = "SET " + set_query
65
66         params = cls._get_params_from_kwargs(**params)  
Python objects → parameter dict
67
68         query = query + set_query
→ safe, parameterized Cypher queries
69         return query, params
```

# From Query Results to Python Objects

## Hydration

- Database returns nodes and relationships
- OGM instantiates Python objects from results
- Relationships reconnect objects into an in-memory graph

DB result → {"title": "Matrix", "released": 1999}

movie = Movie(title="Matrix", released=1999)

(Keanu)-[:ACTED\_IN]->(Matrix)

keanu.acted\_in.append(matrix)

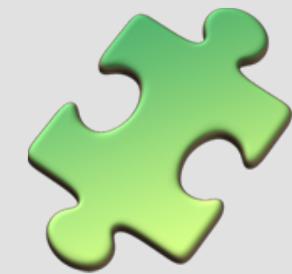
# From Graph Traversal to Graph Computation

Abstractions help us model graphs. Computation requires numeric engines.



## Graph Databases

- Role:** Data structure & storage
- nodes, edges, properties
  - persistence
  - indexing



## OGM

- Role:** Data access & modeling
- Python objects
  - schema abstraction
  - queries & hydration



## HPC / Numeric Stack

- Role:** Computation
- NumPy
  - Numba
  - vectorised loops

- ✓ triangles, centrality, simulations
- ✓ orders-of-magnitude speedups

# HPC with Python on Graph Data

## What Breaks and Why It Matters

### Problem 1 – CPython is slow for graph workloads

1) Python loops are slow → high per-iteration overhead

2) NetworkX is pure Python → collapses for  $n \gtrsim 10^5$

3) GIL → no true parallelism for CPU-bound graph algorithms

4) Poor cache locality:

- graphs cause random memory access

- CPU cache misses dominate runtime

# HPC with Python on Graph Data

## What Breaks and Why It Matters

### Problem 1 – CPython is slow for graph workloads

- 1) Python loops are slow → high per-iteration overhead
- 2) NetworkX is pure Python → collapses for  $n \gtrsim 10^5$
- 3) GIL → no true parallelism for CPU-bound graph algorithms
- 4) Poor cache locality:
  - graphs cause random memory access
  - CPU cache misses dominate runtime

### Problem 2 – Graph algorithms are computationally heavy

- 1) Triangle counting → worst case  $O(n^3)$
- 2) Betweenness centrality → extremely expensive (many shortest paths)
- 3) PageRank:
  - iterative matrix–vector multiplications
  - convergence takes many rounds

# How Python Escapes the Performance Trap

## HPC Stack

**What we  
already  
know**

# How Python Escapes the Performance Trap

## HPC Stack

What we  
already  
know

### NumPy – Vectorization

- Operates on arrays, not Python objects
- Inner loops run in C / Fortran
- Better cache locality
- Removes Python-level iteration

# How Python Escapes the Performance Trap

## HPC Stack

What we  
already  
know

### NumPy – Vectorization

- Operates on arrays, not Python objects
- Inner loops run in C / Fortran
- Better cache locality
- Removes Python-level iteration

### Numba – JIT Compilation

- Python → LLVM → machine code
- Removes interpreter overhead
- `@njit, parallel=True` → multi-core CPU
- Explicit control over loops

# How Python Escapes the Performance Trap

## HPC Stack

Scaling  
further

# How Python Escapes the Performance Trap

## HPC Stack

### ⚡ Parallel execution

- `parallel=True` → loop-level parallelism
- Bypasses GIL for numeric kernels
- Fits graph algorithms with independent iterations

Scaling  
further

# How Python Escapes the Performance Trap

## HPC Stack

### ⚡ Parallel execution

- `parallel=True` → loop-level parallelism
- Bypasses GIL for numeric kernels
- Fits graph algorithms with independent iterations

### 🧠 GPU (cuGraph)

- Thousands of cores
- Massive parallelism for graph primitives
- Same idea: move computation out of Python

Scaling  
further

# Triangle counting in random graphs $G(n,p)$

## Problem Solving

$G \sim G(n, p)$ ,  $X = \#\{\text{triangles in } G\}$ . Find  $\mathbb{E}[X]$ ,  $\text{Var}(X)$ .

$$\mathbb{E}[X] = \binom{n}{3} p^3$$

$$\text{Var}(X) = \binom{n}{3} p^3(1 - p^3) + 3 \binom{n}{4} p^5(1 - p)$$

# Triangle counting in random graphs $G(n,p)$

## Problem Solving

$G \sim G(n, p)$ ,  $X = \#\{\text{triangles in } G\}$ . Find  $\mathbb{E}[X]$ ,  $\text{Var}(X)$ .

### Why it matters

- Fundamental graph metric
- Computationally heavy
- Quickly breaks naïve Python implementations



### What I used

- NumPy → graph as numeric data
- Numba JIT → fast triangle counting
- CPU execution (no database, no OGM)

# Computation

# Sources

1. Introduction of DBMS — GeeksforGeeks: <https://www.geeksforgeeks.org/dbms/introduction-of-dbms-database-management-system-set-1/>
2. Overview of Database Types — DataCamp Blog: <https://www.datacamp.com/blog/types-of-databases-overview>
3. SQL vs NoSQL Databases — DataCamp Blog: <https://www.datacamp.com/blog/sql-vs-nosql-databases>
4. Relational Database — Wikipedia: [https://en.wikipedia.org/wiki/Relational\\_database](https://en.wikipedia.org/wiki/Relational_database)
5. Graph Traversal — Wikipedia: [https://en.wikipedia.org/wiki/Graph\\_traversal](https://en.wikipedia.org/wiki/Graph_traversal)
6. Getting Started with Graph Database — Neo4j Documentation: <https://neo4j.com/docs/getting-started/graph-database/>
7. Graph Databases: The Definitive Guide (O'Reilly) — <https://graphdatabases.com>
8. Building a Neo4j Python OGM (Nodes 2022) — Neo4j Conference Video - <https://neo4j.com/videos/086-building-a-neo4j-python-ogm-nodes2022-estelle-scifo/>
9. Object–Relational Mapping — Wikipedia [https://en.wikipedia.org/wiki/Object%20-%20relational\\_mapping](https://en.wikipedia.org/wiki/Object%20-%20relational_mapping)
10. Nodes 2022 GitHub Repository — Neo4j Python OGM example code - <https://github.com/stellasia/nodes2022>
11. A Comprehensive Guide to Cypher Map Projection — Medium (Neo4j) - <https://medium.com/neo4j/a-comprehensive-guide-to-cypher-map-projection-2622b879e886>

