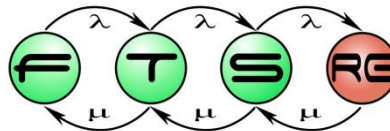


# Hatékony gráflekérdezési technikák

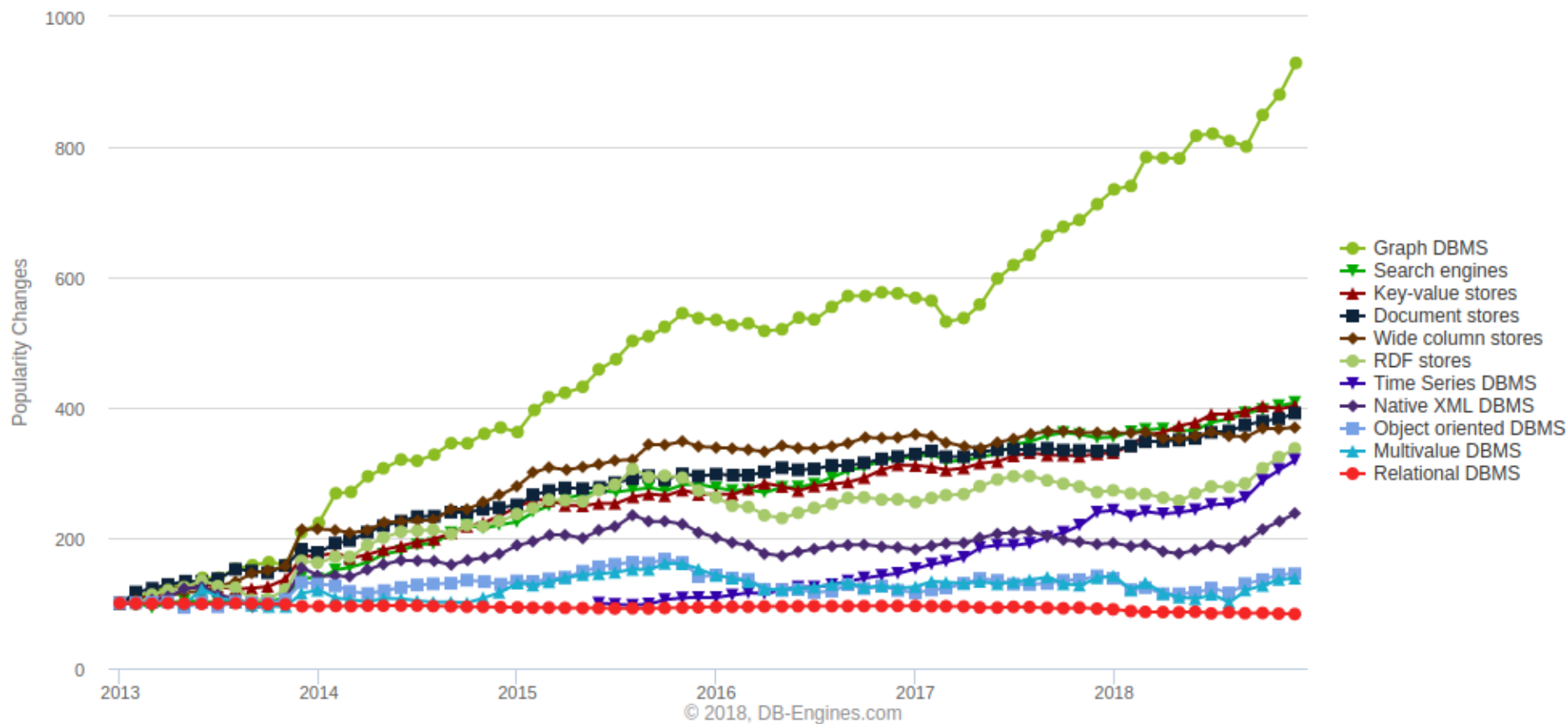
**Antal János Benjamin**

**Konzulens: Szárnyas Gábor**

**Budapesti Műszaki és Gazdaságtudományi Egyetem  
Hibatűrő Rendszerek Kutatócsoport**



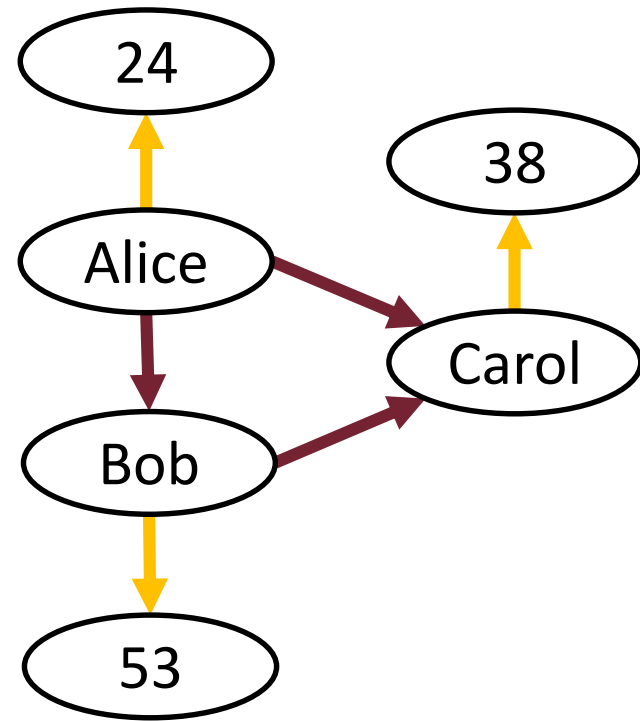
# Motiváció



- Gráfadatbázisok előnyei
  - Intuitív adatmodell
  - Olvashatóság
  - Tömörség
  - Gyors prototipizálás
- Hatékony megoldások?
- Optimalizációk?

# Szemantikus gráf

- Egy gráfban van a meta- és példánymodell
- Resource Description Framework
- Alany, állítmány, tárgy hármassok
- Csúcs- és élcímkezett gráfok

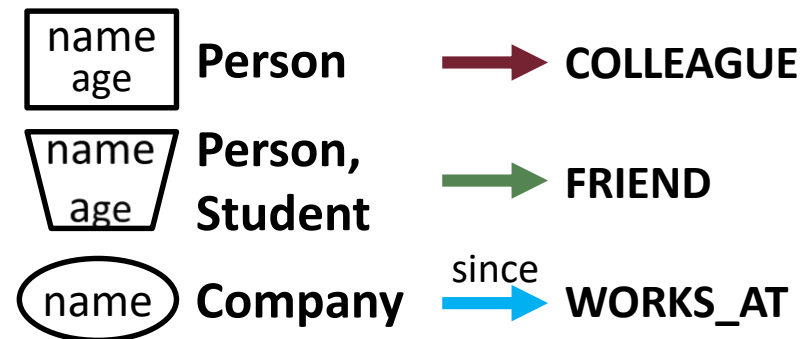
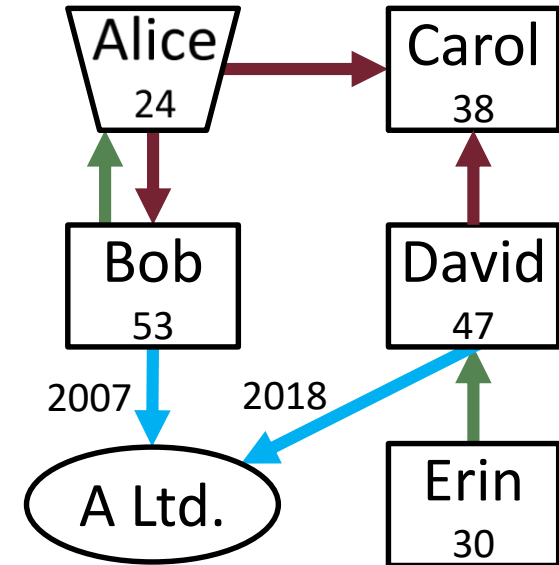


→ COLLEAGUE

→ AGE

# Tulajdonsággráf

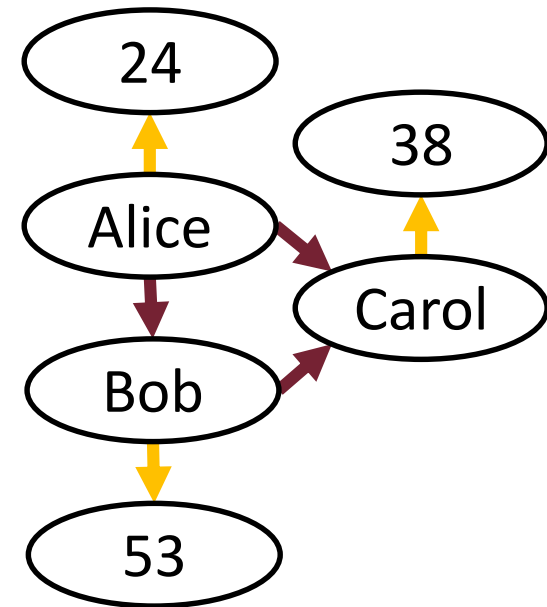
- Leggyakoribb gráf-alapú adatmodell
- Tankönyvi gráf:  $G = (V, E)$
- Kiterjesztve:
  - Címkézett csúcsok
  - Típusos élek
  - Csúcs- és éltulajdonságok
- Implicit séma
  - Adatok felépítése adja
  - Bármikor bővíthető



# SPARQL

- Szemantikus gráfokhoz használható
- $H \leftarrow B$  alakú lekérdezés
  - $B$ : RDF alapú komplex gráfminta (**OPTIONAL**, **UNION** stb.)
  - $H$ : eredmény összeálítása (projekció, szelekció)
- Alice kollégái és az ő kollégáik:

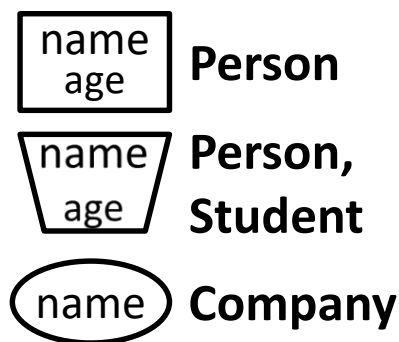
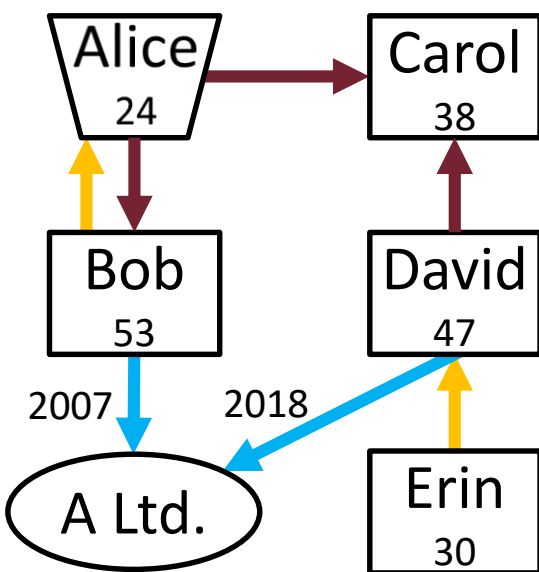
```
SELECT ?pName WHERE {  
  'Alice' COLLEAGUE/COLLEAGUE? ?pName  
}
```



# Cypher

- Tulajdonsággráfokhoz használható
- Legelterjedtebb gráflekérdező nyelv
- Alice kollégái és az ő kollégáik:

```
MATCH (p1:Person {name: 'Alice'})-[c:COLL*1..2]-(p2:Person)  
RETURN p2.name
```



# TELJESÍTMÉNYMÉRÉS



# Teljesítménymérési keretrendszer

- Linked Data Benchmark Council **LDBC** 
  - Teljesítménymérési keretrendszerek és munkafolyamatok definiálása
  - Gráf alapú adatbáziskezelőkhöz
  - Auditált eredmények publikálása
- Social Network Benchmark – LDBC SNB
  - Szintetizált közösségi háló gráf alapú adathalmaz
  - Lekérdezések specifikációja
    - *Business Intelligence* terhelési profil
    - *Interactive* terhelési profil
  - Méréshez szükséges szoftveres keretrendszer

# LDBC SNB bővítése

- Szükséges modulok implementálása több adatbázis-kezelőhöz
- Specifikáció frissítése, konzisztens állapotra hozása
- Hibajavítások
- Business Intelligence terhelési profil
  - 24 darab lekérdezés implementálása SPARQL nyelven
  - A lekérdezések validációja
- Interactive terhelési profil
  - 29 darab lekérdezés implementálása Cypher és SPARQL nyelven
  - A lekérdezések validációja
- Teljesítménymérés

# LDBC SNB bővítése

- Szükséges modulok implementálása több adatbázis-kezelőhöz
- Specifikáció frissítése, konzisztens állapotra hozása
- Hibajavítások
- Business Intelligence
  - 24 darab lekérdezés implementálása
  - A lekérdezések validációja
- Interactive SQL
  - 29 darab lekérdezés implementálása Cypher és SPARQL nyelven
  - A lekérdezések validációja
- Teljesítménymérés



# LDBC SNB bővítése

- Szükséges modulok implementálása több adatbázis-kezelőhöz
- Specifikáció frissítése, konzisztens állapotra hozása
- Hibajavítások
- Business Intelligence
  - 24 darab lekérdezés
  - A lekérdezések vandergera
- Interactive
  - 29 darab lekérdezés
  - A lekérdezések vandergera
- Teljesítménymérés



A lekérdezések fehérszóközök nélkül  
**94 986** karakterből állnak

SQL nyelven

# Teljesítménymérés

- Eszközök
  - Sparksee (referencia implementáció)
  - PostgreSQL (már meglévő implementáció)
  - SDB1: szemantikus adatbázis-kezelő
  - SDB2: szemantikus adatbázis-kezelő
- Lekérdezésenként legalább 20 db különböző behelyettesítési paraméterrel

# Teljesítménymérés

## ■ Eszközök

- Spa Anonimizált eredmények, mert nem auditáltak

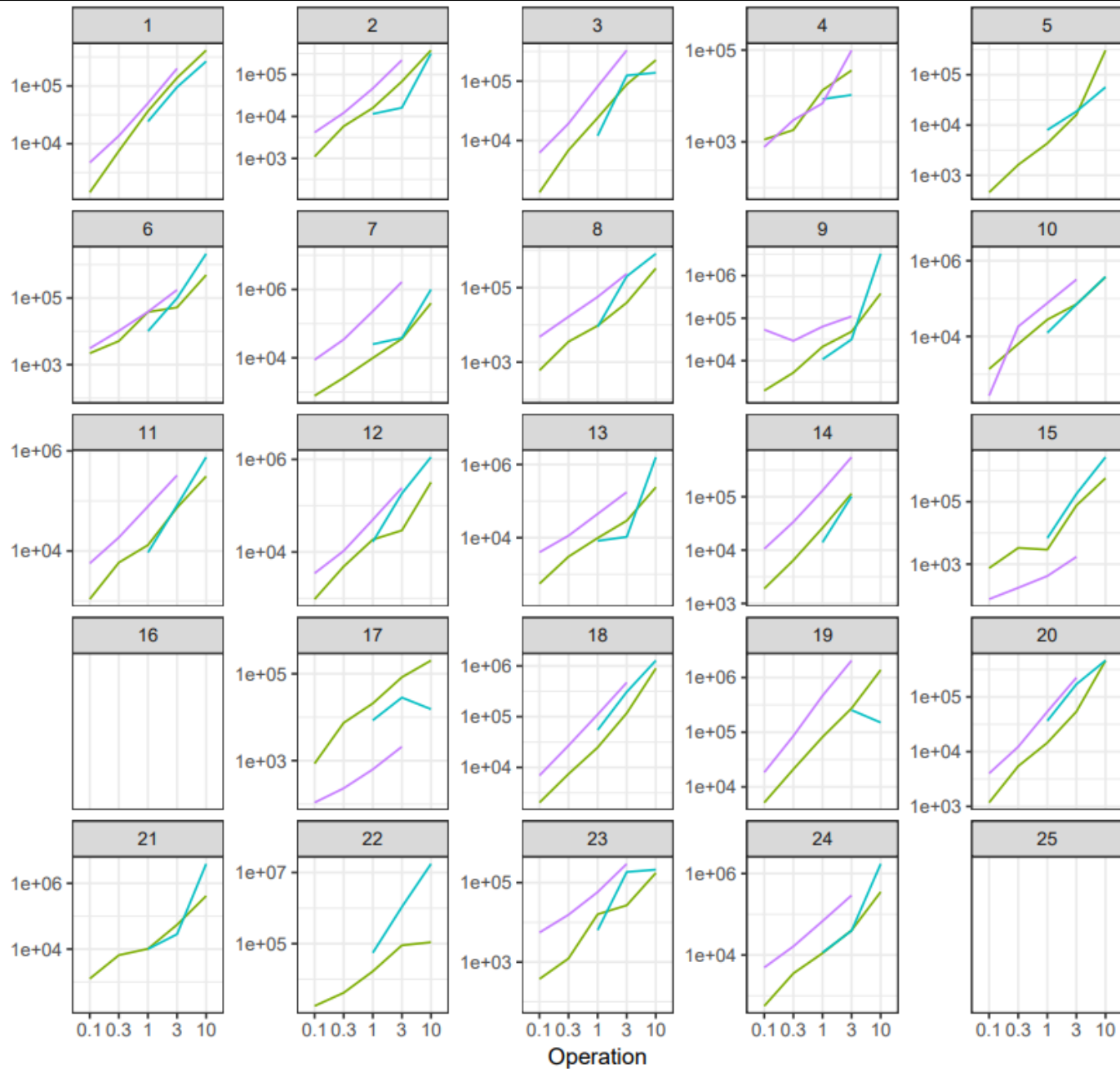
- PostgreSQL (nem teljeskörű implementáció)

- SDB1: szemantikus adatbázis-kezelő

- SDB2: szemantikus adatbázis-kezelő

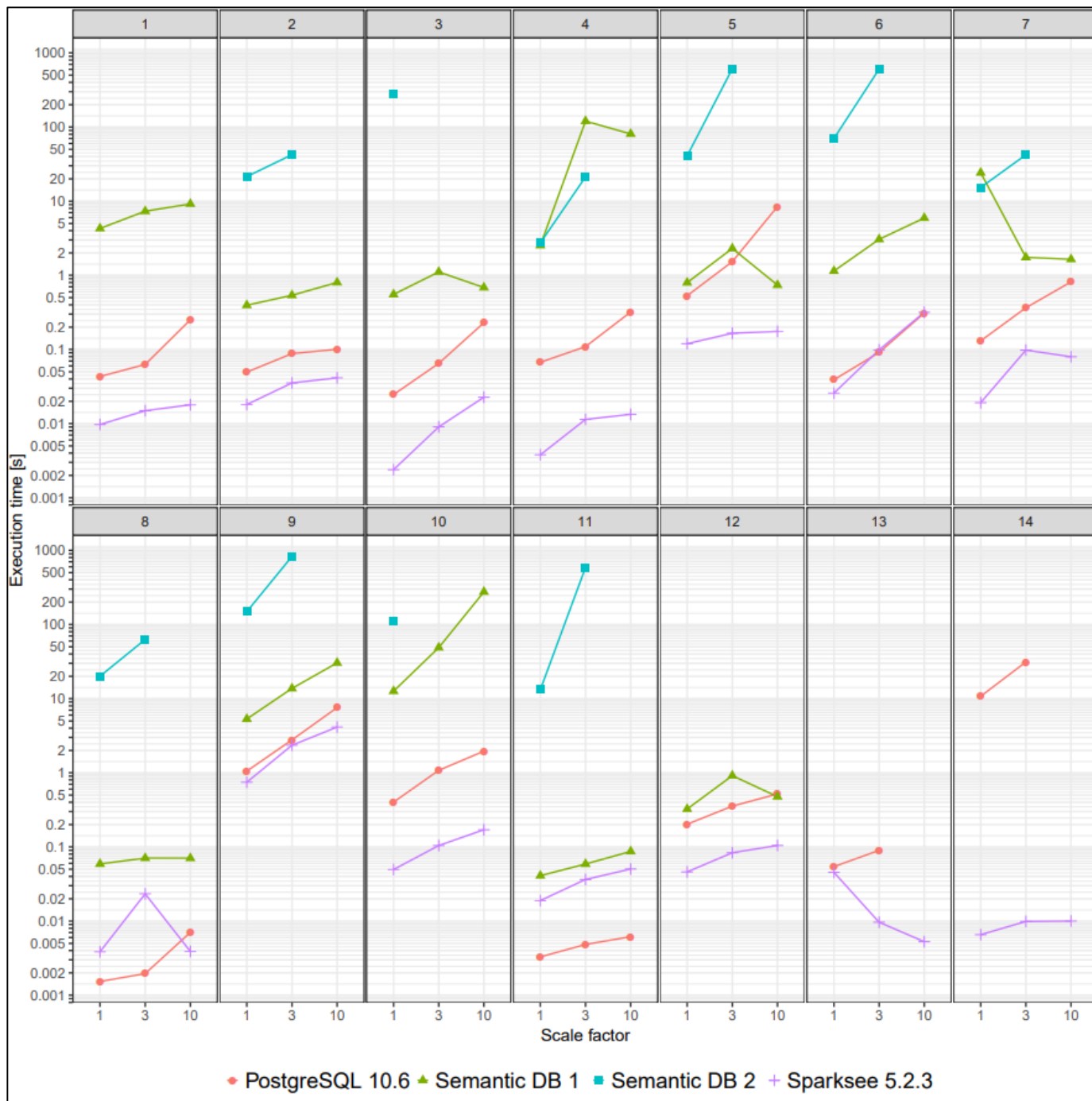
## ■ Lekérdezésenként legalább 20 db különböző behelyettesítési paraméterrel

Execution time [ms]



tool

PostgreSQL Sparksee SDB2





# INKREMENTÁLIS NÉZETKARBANTARTÁS

# Adatfolyam alapú számítási modellek

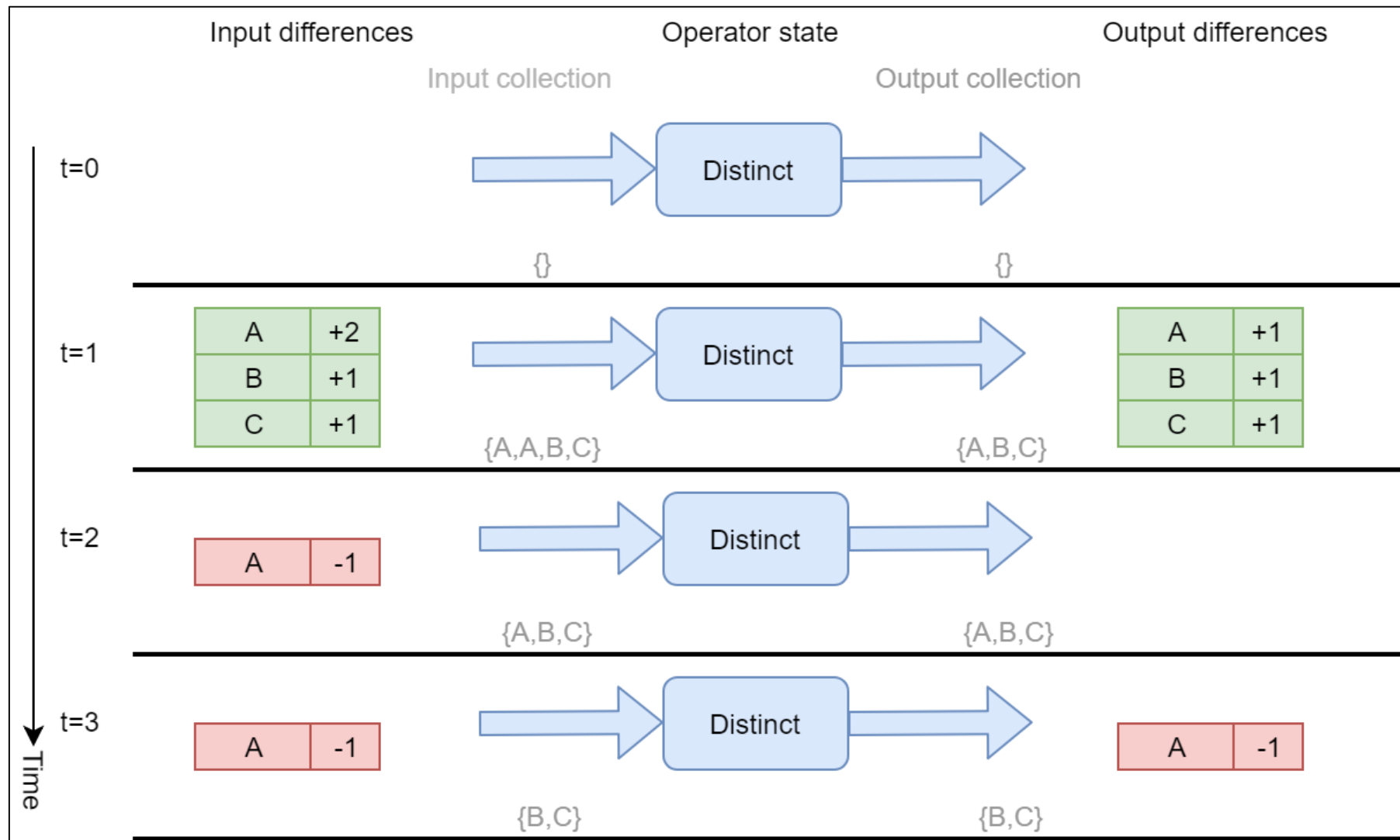
## ■ Időzített adatfolyam

- Az adatfolyamban minden adatrekordhoz egy virtuális időbélyeget rendelünk
- Nagy áteresztőképességű, elosztott számítási modell
- Támogat iteratív számításokat is

## ■ Differenciális adatfolyam

- Időzített adatfolyamon felett implementált
- Inkrementális nézetkarbantartáshoz
- Csak a változás továbbítódik az adatfolyamban
- Hagyományos operátorok: **SELECT, JOIN, DISTINCT, COUNT**
- **FIXEDPOINT, ENTERLOOP**

# Distinct operátor működése

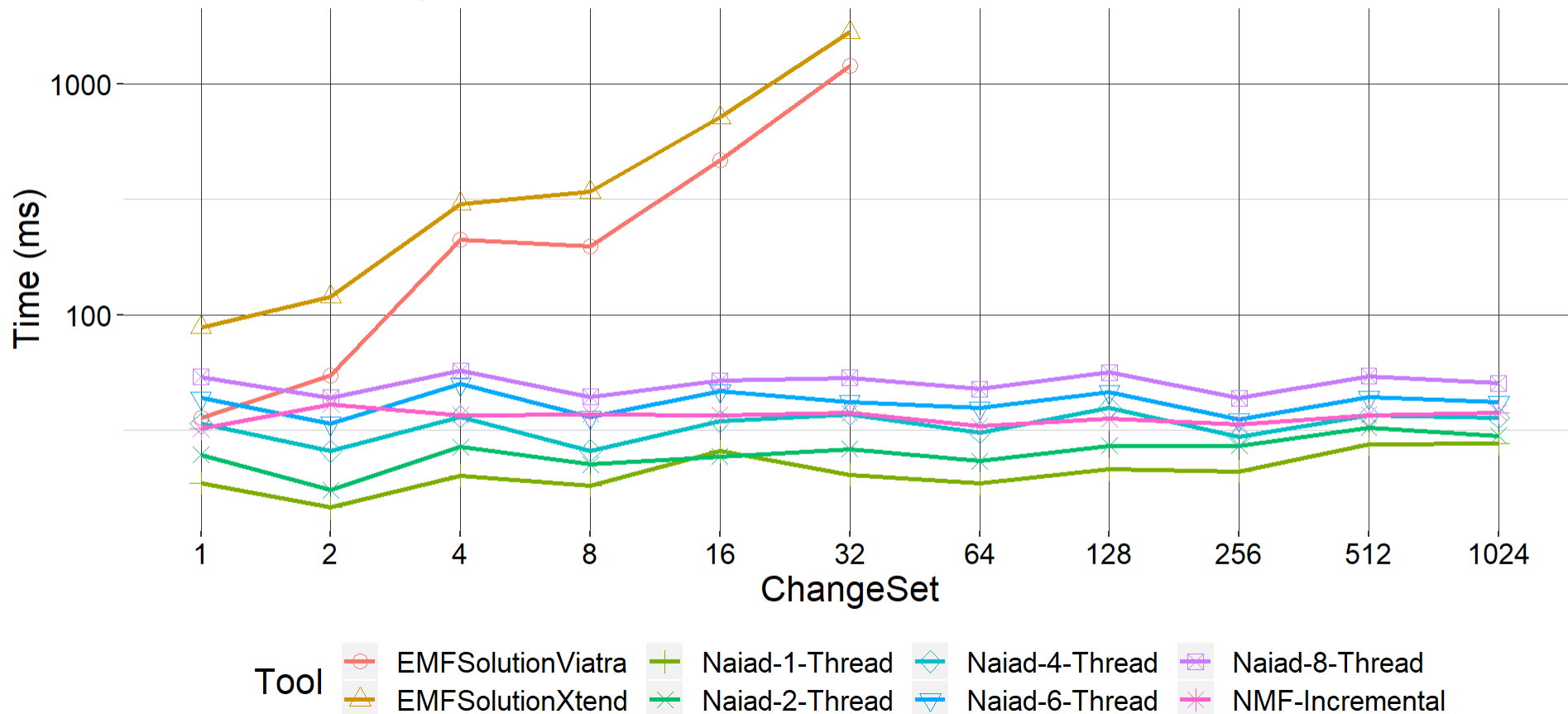


# Transformation Tool Contest

- Modelltranszformációs eszközök összehasonlítása
  - kifejezőerő, használhatóság és teljesítmény
- 2018-as feladat az LDBC SNB egyszerűsített sémája feletti inkrementális nézetkarbantartás
- Kettő konkrét lekérdezés megoldása volt a feladat
- Több, már meglévő megoldás
  - Xtend
  - Viatra
  - .NET Modelling Framework (NMF)

# TTC 2018 Q1 mérési eredmények

Q1, Function: Update



# Összefoglalás

- Gráf-alapú adatmodellek és lekérdezőnyelvek
  - Szemantikus gráf, tulajdonsággráf, SPARQL , Cypher
- LDBC SNB bővítése
  - Business Intelligence: SPARQL implementáció (**32 417** karakter)
  - Interactive: SPARQL és Cypher implementáció (**62 569** karakter)
  - Szükséges szoftvermodulok
  - Teljesítménymérés (25 + 14 darab lekérdezés)
- Inkrementális nézetkarbantartás
  - Időzített és differenciális adatfolyamok megismerése
  - TTC 2018 megoldása differenciális adatfolyammal
  - Teljesítménymérés

# Bíráloi kérdések 1.

- A háttérismeretek fejezetben szerepel, hogy „[Az inkrementális nézetkarbantartás] két csoportját különböztetjük meg: az algebrai és a procedurális megközelítéseket”. A differenciális adatfolyamokat használó technikát hova sorolja, és miért?
- Válasz:
  - A kategorizálást alapvetően deklaratív lekérdezőnyelvekhez találták ki
  - Általánosítani lehet
  - **Procedurális**

# Bíráloi kérdések 2.

- A Cypher lekérdezés Gremlinre fordított változata szuboptimális és több technikai probléma is felmerült. Van-e valós gyakorlati igény a Cypher-Gremlin transzformációra?
- Válasz:
  - Microsoft CosmosDB támogatja a Gremlin lekérdezőnyelvet
    - 367 szavazat a Cypher támogatásra is [1]
  - Az motiváció a több gráfadatbázis teljesítménymérése a lekérdezések újabb implementációja nélkül („low hanging fruit”)

[1] <https://feedback.azure.com/forums/263030-azure-cosmos-db/suggestions/19275547-support-cypher-as-a-query-language-for-graph-data>



# Bírálóí kérdések 3./I

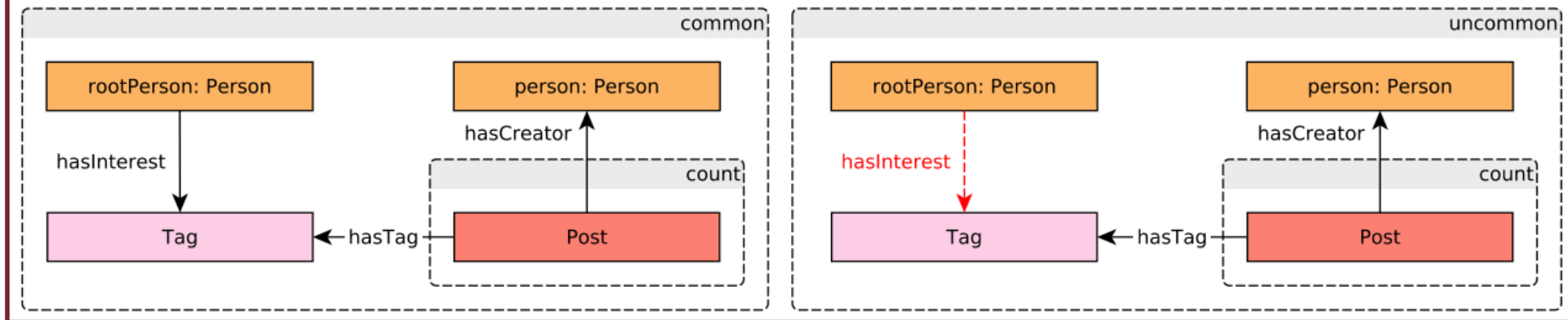
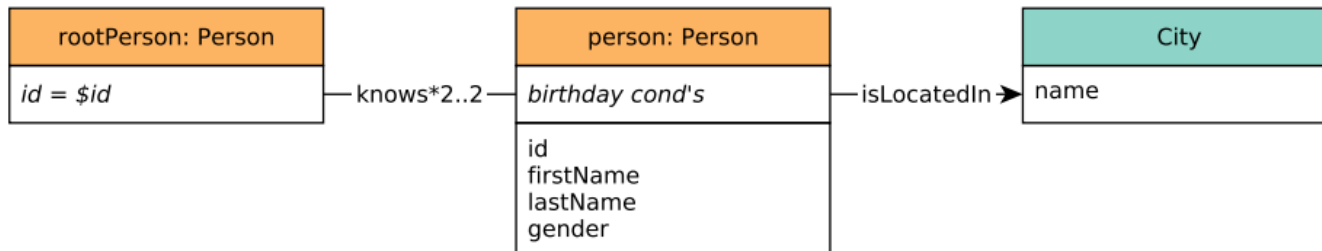
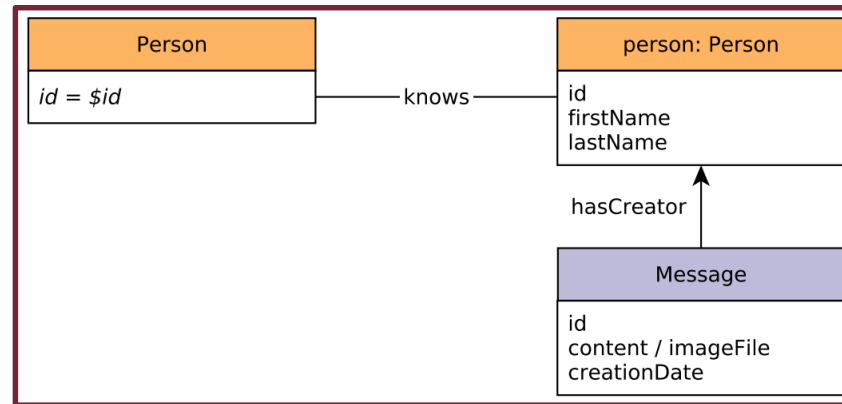
- A 2. és 10. (és esetleg még a 11.) lekérdezések jól látható teljesítménybeli sorrendet állítanak fel az egyes rendszerek között. Ki tud-e valami közös jellemzőt emelni ezen lekérdezések esetén, amik ezt a jelenséget okozhatják?

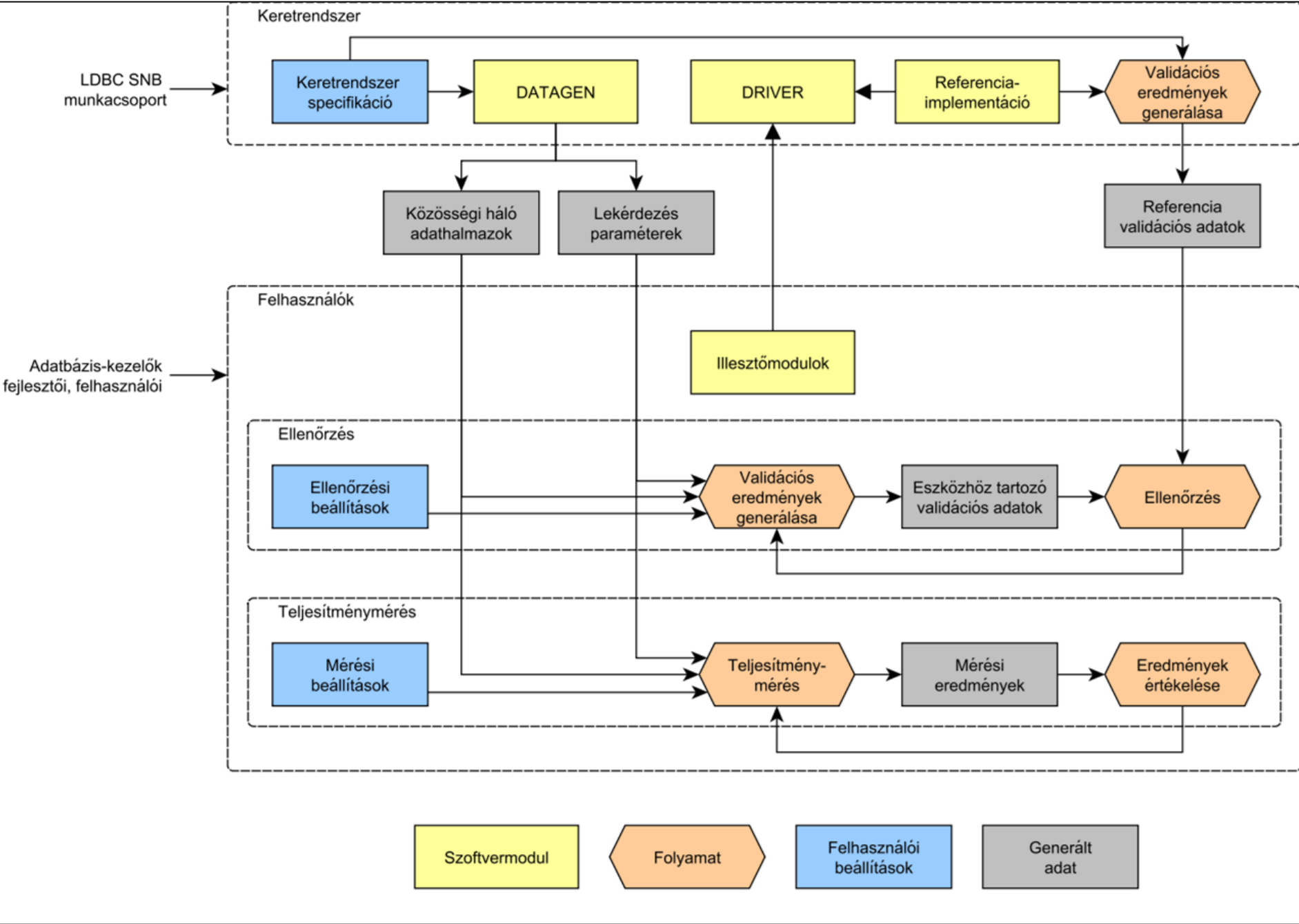
# Bírálóí kérdések 3./I

- A 2. és 10. (és esetleg még a 11.) lekérdezések jól látható teljesítménybeli sorrendet állítanak fel az egyes rendszerek között. Ki tud-e valami közös jellemzőt emelni ezen lekérdezések esetén, amik

	1.1	1.2	1.3	1.4	2.1	2.2	2.3	2.4	3.1	3.2	3.3	4.1	4.2	4.3	5.1	5.2	5.3	6.1	7.1	7.2	7.3	7.4	8.1	8.2	8.3	8.4	8.5	8.6
IC 1					●												●							●				
IC 2	●					●	●			●																	●	
IC 3					●				●						●									●			●	
IC 4							●																	●			●	
IC 5							●				●													●			●	
IC 6															●									●				
IC 7						●	●				●				●								●		●			
IC 8								●		●	●						●											
IC 9	●	●				●	●			●	●																●	
IC 10							●				●	●	●		●	●		●	●									●
IC 11			●				●	●			●																	
IC 12											●									●	●			●				
IC 13											●									●	●		●				●	
IC 14											●										●	●		●				●

# Bírálói kérdések 3./II





# TTC 2018 Q1

