



Budapesti Műszaki és Gazdaságtudományi Egyetem
Villamosmérnöki és Informatikai Kar
Méréstechnika és Információs Rendszerek Tanszék

Gráfmintaillesztő rendszerek tesztelése automatikus gráfgenerátorokkal

TDK dolgozat

Készítette:

Bekő Mária

Konzulens:

Semeráth Oszkár

2018

Tartalomjegyzék

Kivonat	i
Abstract	ii
1. Bevezetés	1
2. Előismeretek	2
2.1. Esettanulmány	2
2.2. Gráflekérdező rendszerek	2
2.2.1. Neo4j	3
2.3. Modellezés és metamodellezés	3
2.3.1. Cypher query-k metamodellje	3
2.3.2. xText	5
2.3.3. VIATRA jólformáltsági kényszerek	5
2.4. Gráfgenerálás	5
3. Áttekintés	6
3.1. Funkcionális áttekintés	6
3.2. Blokkdiagram	7
3.3. A tervezett alkalmazás felépítése	7
4. Gráflekérdezések automatikus generálása	8
4.1. Lekérdezések generálása	8
4.1.1. Nyelv minimalizálása	8
4.1.2. Jólformáltsági kényszerek betartása	8
4.1.3. Diverzitás biztosítása	8
4.2. Automatizálás	8
5. Értékelés	9
5.1. Benchmark eredmények	9
5.2. Feature fedés eredmények	9
5.3. Diverzitás mérése	9
6. Összefoglalás	10
6.1. Elméleti eredmények	10
6.2. Gyakorlati eredmények	10
6.3. Jövőbeli tervek	10
Köszönetnyilvánítás	11
Irodalomjegyzék	12

Kivonat

Napjainkban az adatokat többféle formátumban is tárolják. Ezek közé tartoznak a gráfadatbázisok, ahol csomópontok reprezentálják az entitásokat és az élek az entitások közötti kapcsolatokat. Az adatstruktúrához illeszkedve többféle gráflekérdező nyelv jött létre, amelyek képesek komplex struktúrák felírására.

A gráfmintaillesztő rendszerek tesztelése azonban komoly kihívást jelent, főképp automatizált megoldásokban nem bővelkedünk. A legnagyobb kihívást ebben az esetben a változatos modellek és lekérdezések automatikus és szisztematikus előállítása jelenti, melyek tesztbemenetként szolgálnak. Továbbá, gráfadatbázisok teljesítménymérését is nagyban segítené az automatikusan előállított modellkészlet.

Dolgozatom célja hogy ezekre a problémákra megoldást találjak. Kutatásom során megmutatom, hogy egy automatikusan előállított diverz modell halmazzal, amelynek modelljei lekérdezőként értelmezhetőek egy gráfmintaillesztő rendszerben (pl.: VIATRA vagy Neo4j), hogyan lehetséges az adott gráfmintaillesztő rendszer tesztelése.

Munkám során fejlett logikai következtetők alkalmazásával állítok elő modelleket, melyek diverzitását szomszédsági formákkal (neighborhood shape-ek) biztosítom. A logikai következtetők eredményeit lekérdezéseként, és adatbázisok tartalmaként egyaránt értelmezhetjük, amelyek eredményei különböző megvalósításokkal összehasonlíthatóvá válnak. A megoldásomat egy esettanulmány keretében prezentálom.

Ezzel a módszerrel lehetővé válik, nagyobb megbízhatóságú gráfmintaillesztő rendszerek fejlesztése olcsóbban. Illetve egy ekkora modell halmaz különböző gráfmintaillesztő rendszerekben lekérdezésekre fordítva és a válaszüzeneteket lemérve teljesítménymérésekre is használható.

Abstract

Nowadays the data is stored in multiple formats. One of this is the graph database, where entities are represented as graph nodes and the relations between entities as edges. Utilizing rich data structure, a variety of graph querying languages have been created in order to query complex structures.

Testing of graph query engines is a challenging task, especially in an automated way. The greatest challenge in this case is the automatic and systematic creation of a diverse set of models and queries, which serve as test inputs. In addition, the development of performance benchmarks for graph databases would be greatly aided.

The purpose of my thesis is to find solutions to these problems. In my thesis, I will show an approach to automatically generate a diverse set of models, that can be interpreted as queries of the graph query engine under test (e.g. VIATRA model query language or Neo4j graph database queries). Additionally, I propose a testing process.

In the course of my work, I produce models with the help of state-of-the-art logic solvers (SAT solvers and Graph Solvers), and use neighborhood shapes to ensure their diversity and create effective equivalence partitioning. The results of the logic solvers are interpreted as queries and as databases content, and the result of query evaluation can be compared to other implementations. I illustrate my solution in a case study.

This method makes it possible to develop more reliable graph query engines at a lower cost. And such a set of models translated to multiple graph querying languages can be used in those graph query engines for performance measurements by measuring response times.

1. fejezet

Bevezetés

Kontextus. Napjainkban

Problémafelvetés. Azonban ...

Célkitűzés. Dolgozatom célja ...

Kontribúció. Dolgozatomban bemutatok ...

Hozzáadott érték. Ezáltal ...

Dolgozat felépítése. A második fejezetben bemutatom a dolgozat megértéséhez szükséges háttérismereteket. blabla ...

2. fejezet

Előismeretek

2.1. Esettanulmány

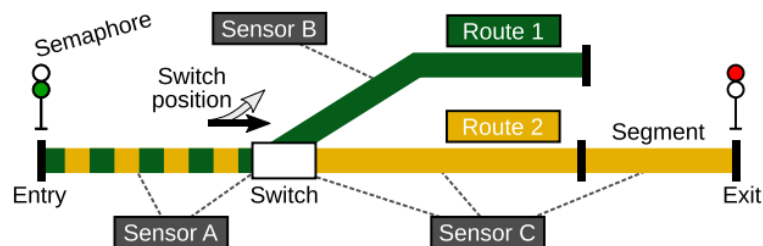
A dolgozatban elért eredményeket a Train Benchmark [8] esettanulmány segítségével fogom bemutatni. Ez a benchmark azért jött létre, hogy összehasonlítsuk különböző gráflekérdező rendszerek teljesítményét főleg időigény és memória felhasználás szempontjából. Ehhez többek között definiálja a vasút rendszer metamodeljét. Munkám során ezt a metamodelt felhasználva generáltam lekérdezéseimet, ezért bemutatom, hogy milyen elemekből áll.

A 2.1-es ábrán látható egy a trainbenchmark metamodeljére alapuló részlet. Ebben a kontextusban egy vasúti útvonal nem más mint szegmensek és váltók sorozata, illetve a belépést és a kilépést egy-egy szemafor jelzi. Ahhoz hogy biztonságos legyen a közlekedés szükség van szenzorokra, amelyek monitorozzák a különböző szegmensek és váltók kihasználtságát. Egy útvonal definiálásához a felsorolt elemeken kívül a váltók adott útvonalhoz tartozó pozícióját is el kell tárolnunk. Egy útvonal akkor aktív, ha a rendszerben a specifikációjának megfelelően állnak a váltók.

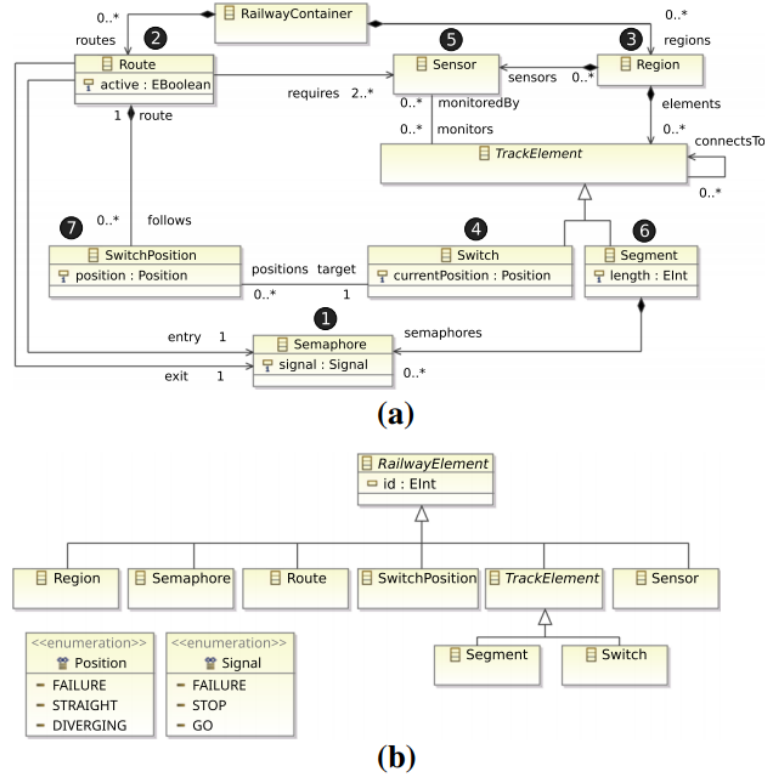
A metamodellezés egy technika arra, hogy definiáljunk modellező nyelveket, ahol a metamodell specifikálja a nyelv szintaktikáját. A trainbenchmark metamodelje a 2.2-es ábrán látható.

2.2. Gráflekérdező rendszerek

A gráfok intuitív formalizációt nyújtanak modellezési szempontból arra, hogy úgy írhatunk le a világot ahogy az ember gondolkozik róla. Tehát mint dolgok (csomópontok) és köztük lévő kapcsolatok (élek) [5]. A property gráf adatmodell kiterjeszti a gráfokat úgy, hogy címkéket/típusokat illetve tulajdonságokat ad a csúcsoknak és az éleknek. A gráf adatbázisok alkalmasak tulajdonság gráfok tárolására, és az abban lévő adatok lekérdezésére



2.1. ábra. Vasúti útvonal részlet (forrás: [8])



2.2. ábra. A train benchmark metamodellje. a, tartalmazási hierarchia és relációk b, öröklési relációk (forrás: [8])

komplex gráf minták használatával. Ilyen rendszerek például a Neo4j [1], OrientDB [2] és a SparkSee [4].

Ahhoz hogy jobban megérthessük mi is ez az adatmodell a 2.3 -es ábrán látható egy példa.

2.2.1. Neo4j

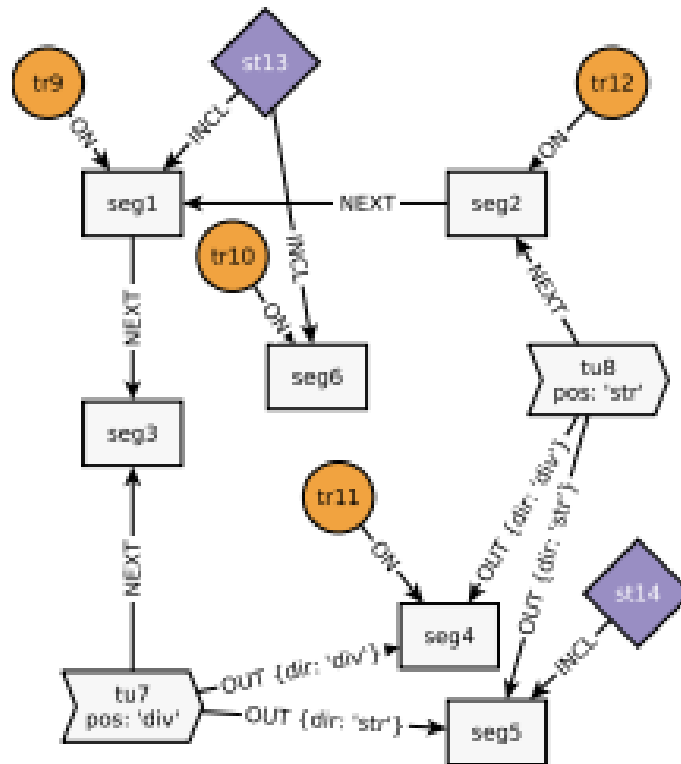
A Neo4J egy populáris NoSQL property gráf adatbázis és a Cypher lekérdező nyelvet kínálja lekérdezések írására. A Cypher egy magas szintű deklaratív lekérdező nyelv és mivel le van választva a lekérdező rendszerről, ezért az képes a Cypher nyelven írt lekérdezések optimalizálására. A Cypher szintaxisa olyan gráf minták megírását teszi lehetővé, amelyeknek megértése nagyon egyszerű. A 2.4 -es példán egy olyan lekérdezést láthatunk, ami az összes olyan vonat, szegmens párral tér vissza, ahol az adott vonat rajta van az adott szegmensen.

2.3. Modellezés és metamodellezés

A metamodellek definiálják a legfontosabb fogalmakat, relációkat és attribútumokat a cél domainen , hogy specifikálják modellek alap struktúráját [7]. Dolgozatomban az Eclipse Modeling Framework (EMF) [3] -öt használtam metamodellezésre.

2.3.1. Cypher query-k metamodellje

A 2.5 -ös ábrán látható a korábban említett Cypher nyelv egyszerűsített metamodellje. A [SinglePartQuery](#) elem reprezentálja a modell gyökerét. Egy ilyen elem két részből áll



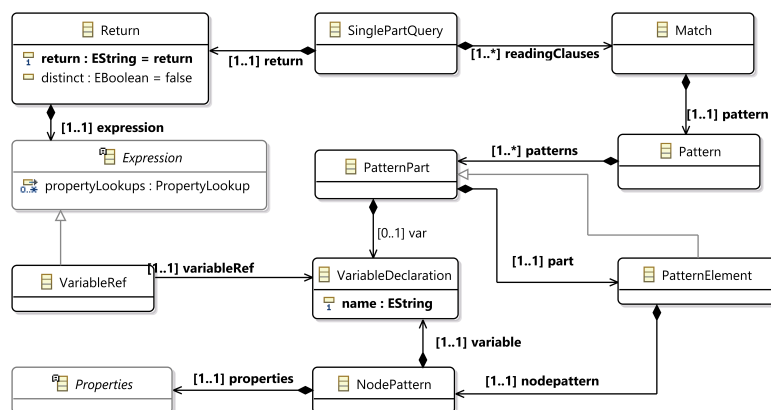
2.3. ábra. Property gráf példa(forrás: [5])

```

MATCH (tr:Train)-[:ON]->(seg:Segment)
RETURN tr, seg

```

2.4. ábra. Példalekérdezés(forrás: [5])



2.5. ábra. Cypher metamodel

: Egy **Match** és egy **Return** elemből. A **Match** elem minták összességéből áll (**Pattern**), amelyek részekre (**PatternPart**) bonthatóak. Egy ilyen részben pedig vagy tartalmaz változó deklarációt, vagy egy belső részt tartalmaz ami tartalmaz változó deklarációt. (**Variable Declaration**). Ezáltal az összes változót a **Match** elemen belül deklaráljuk. A **Return** elem pedig egy kifejezést (**Expression**) tartalmaz, amelyben mindenképp szerepel egy változó referencia (**VariableRef**) is, így összekötve egymással a **Match** és a **Return** elemet.

2.3.2. xText

2.3.3. Viatra jólformáltsági kényszerek

2.4. Gráfgenerálás

ábra milyen bemenetek milyen kimenetek

3. fejezet

Áttekintés

3.1. Funkcionális áttekintés

Munkám célja hogy mutasson egy olyan koncepciót, amelynek segítségével lehetséges gráf-mintaillesztő rendszerek tesztelése több aspektusból. Ezt úgy viszem véghez, hogy a rendszer nyelvén generálok mintákat/lekérdezéseket, amelyeket aztán a lekérdező motoron futtatok. Mivel generálás során nem akadály az hogy sok és diverz lekérdezést készítek el ezért a nagy lekérdezőhalmaz alkalmassá válik arra hogy teljesítményben tesztelje a lekérdező rendszereket. Illetve ha egyes lekérdezéseken nem tud futtatni a rendszer akkor kiderül az is hogy milyen funkcionalitásokat nem fed még le. Illetve egy referencia implementáció segítségével azt is ellenőrizni tudjuk, hogy a megfelelő válaszolakt adja-e a rendszer, hibás-e a működése. A koncepciómat egy Neo4J [1] gráf adatbázison mutatom be miközben a lekérdezéseket az ehhez kifejlesztet Cypher [6] nyelven generálom.

1. Nyelvi specifikáció
2. Szintaxis
3. Absztrakt Sintaxis

Gráf

4. Modell

A generálást gráfok segítségével végzem. Ezért van szükségem a nyelv szintaktikájának metamodelljére, aminek alapján példánymodelleket tudok generálni az adott nyelv felépítését betartva. Azonban vannak olyan szabályok amelyeket a metamodell nem fejez ki, betartásuk nélkül viszont a generált példánygráfok nem értelmezhetők Cypher nyelvű lekérdezésekként. Ezen szabályok betartására jólformáltsági kényszereket írok fel, amelyeket a generátor támpontként használ a lekérdezések készítése során. Ahhoz hogy mindezt véghez tudjam vinni választanom kellett egy gráf generáló szoftvert ami kényszerek és egy metamodell alapján képes hatékonyan példánymodellek generálására. Ehhez a ViatraSolver nevű alkalmazást választottam, ami eclipse fejlesztői környezetben működik, a jólformáltsági kényszereket pedig Viatra nyelven írtam fel.

Miután a példánygráfok létrejönnek már csak le kell őket fordítani Cypher nyelvű lekérdezésekre. A fordítás során ki kell tölteni a gráfgenerátor által üresen hagyott értékeket. Illetve át parszolni a gráfokat egy .cypher típusú fájlba.

3.2. Blokkdiagram

3.3. A tervezett alkalmazás felépítése

4. fejezet

Gráflekérdezések automatikus generálása

4.1. Lekérdezések generálása

4.1.1. Nyelv minimalizálása

4.1.2. Jólformáltsági kényszerek betartása

Ahhoz, hogy a nyelvtan alapján olyan modelleket tudjunk generálni, amelyeket vissza tudunk majd fordítani szintaktikailag helyes lekérdezésekre, ahhoz kényszereket kell felállítanunk, amelyek betartására kötelezzük a modellgenerátort.

Én a következő kényszereket írtam fel: `pattern x(s:State) State(s);`, ahol a `s` változó jelöli a...

```
pattern goodReferenceToVariable(q : SinglePartQuery, match : Match, vari : VariableDeclaration, variRef : VariableRef) {
    SinglePartQuery.readingClauses(q ,match);
    Match.^pattern.patterns.^var(match, vari);
    SinglePartQuery.^return.body.returnItems.items.expression(q , variRef);
    VariableRef.variableRef(variRef, vari);
}

@Constraint(severity = "error", key ={variRef}, message ="error")
pattern notGoodReferenceToVari(variRef : VariableRef){
    neg find goodReferenceToVariable(_, _, vari ,variRef );
    VariableRef.variableRef(variRef,vari);
}
```

A fenti két kényszer azt biztosítja, hogy a lekérdezésnek, amit generálunk mindenképpen meglegyen a két fő része. A match részben határozzuk meg azt a bizonyos mintát, amit keresünk, a return részben pedig visszatérünk az erre a mintára illeszkedő gráfrészletekkel.

Ezen belül ki kell kössük szintén, hogy a két fő rész ne csak csonkként vagy helytelenül generálódjon, hanem az egész elengedhetetlen részfa megszülessen.

Ezért van szükség a return ágon az alábbi kényszerekre:

4.1.3. Diverzitás biztosítása

4.2. Automatizálás

5. fejezet

Értékelés

5.1. Benchmark eredmények

5.2. Feature fedés eredmények

5.3. Diverzitás mérése

6. fejezet

Összefoglalás

6.1. Elméleti eredmények

6.2. Gyakorlati eredmények

6.3. Jövőbeli tervek

Köszönetnyilvánítás

Ez nem kötelező, akár törölhető is. Ha a szerző szükségét érzi, itt lehet köszönetet nyilvánítani azoknak, akik hozzájárultak munkájukkal ahhoz, hogy a hallgató a szakdolgozatban vagy diplomamunkában leírt feladatokat sikeresen elvégezze. A konzulensnek való köszönetnyilvánítás sem kötelező, a konzulensnek hivatalosan is dolga, hogy a hallgatót konzultálja.

Irodalomjegyzék

- [1] 2018 October 16–2018 October 15: The neo4j graph platform – the #1 platform for connected data. URL <https://neo4j.com/>.
- [2] Graph database | multi-model database. URL <https://orientdb.com/>.
- [3] Richard Gronback: Eclipse modeling framework (emf).
URL <https://www.eclipse.org/modeling/emf/>.
- [4] High-performance human solutions for extreme data.
URL <http://sparsity-technologies.com/#sparksee>.
- [5] József Marton–Gábor Szárnyas–Márton Búr: Model-driven engineering of an opencypher engine: Using graph queries to compile graph queries. In *International SDL Forum* (konferenciaanyag). 2017, Springer, 80–98. p.
- [6] Neo4j’s graph query language: An introduction to cypher.
URL <https://neo4j.com/developer/cypher-query-language/>.
- [7] Oszkár Semeráth–Ágnes Barta–Ákos Horváth–Zoltán Szatmári–Dániel Varró: Formal validation of domain-specific languages with derived features and well-formedness constraints. *Software & Systems Modeling*, 16. évf. (2017) 2. sz., 357–392. p.
- [8] Gábor Szárnyas–Benedek Izsó–István Ráth–Dániel Varró: The train benchmark: cross-technology performance evaluation of continuous model queries. *Software & Systems Modeling*, 17. évf. (2018) 4. sz., 1365–1393. p.