

# Bírálat

Toka László

## Erőforrás-kezelés felhőrendszerekben (Resource Provisioning in Cloud Systems)

### c. MTA doktori értekezéséről

**Bíráló: Jelasity Márk**

#### Bevezető

Toka László disszertációja a felhő infrastruktúrákban fellépő menedzsment feladatok, mint pl. skálázás vagy erőforrás elosztás, vizsgálatával foglalkozik, és ennek számos területén mutat be értékes kontribúciókat: algoritmusokat és elméleti eredményeket egyaránt. A téma jelenleg is aktuális, jól motivált. A felhő, és a hozzá kapcsolódó perem rendszerek, technológiák (IoT, kód, mobil hálózatok) továbbra is kiemelt szerepet játszanak a jelenkori informatikai infrastruktúrában, amelyeknek az optimalizálása rendkívül fontos.

Bár nem tartozik szorosan a doktori műhöz, kiemelem, hogy Toka László, minőségi újságcikkei mellett, rendelkezik számos értékes konferencia publikációval, közelebbről több INFOCOM cikk szerzője. Ezen kívül is rendszeresen jár jó minőségű konferenciákra. Ennek a motiválása a doktori szempontrendszerben nagyobb szerepet kellene hogy kapjon, mivel a láthatóság és presztizs szempontjából ez kulcsfontosságú. Idézettsége kiváló, és egyértelműen felszálló ágban van.

A továbbiakban az egyes tézisekkel kapcsolatos megjegyzéseimet és kérdéseimet ismertetem.

#### 1. Téziscsoport: A felhők közötti üzleti kapcsolatok

Ebben a téziscsoportban a szerző elméleti modelleket állít fel a felhő szolgáltatások komplex piaci mechanizmusaira, amiből egyfelől a felhő szolgáltatások által indukált hálózati struktúrákat lehet megmagyarázni, másfelől a szolgáltatások árazását lehet elemezni.

A tézisek első részcsoportja arról szól, hogy a szolgáltatók közötti kapcsolatok kialakulásának milyen feltételei vannak. Ehhez a szerző játékelmélet modellt ír fel, amelyben a lehetséges stratégiákat a kapcsolatok kialakítása jelenti, és a jutalomfüggvény pedig a kapcsolatok által indukált előnyöket (közvetítő szerep révén bevétel) és költségeket jelenít meg. A modell újdonságtartalma a közvetítő szerepből származó bevételek figyelembevétele, ami miatt nemtriviális hálózati struktúrák jelenhetnek meg.

**Kérdés 1.1:** *Jól értem-e, hogy a játék egyensúlyi helyzete a levelek felé haladva egyre kedvezőtlenebb lehet, mert a leveleknek nincs beleszólása abba, hogy számukra ki szolgáltat kapcsolatot? Ez azért van, mert a legrövidebb úton levő költségek vannak figyelembe véve, azonban a legrövidebb út csak a súlyozatlan topológia alapján kerül kiszámításra, és ennek kialakítására nincs egy csomópontnak hatása.*

**Kérdés 1.2:** *Nem jelent a homogén szolgáltatás feltevése túl nagy egyszerűsítést? A valóságban a kapcsolatok megbízhatósága, ebből adódó bizalmi viszonyok, a kapcsolatok kialakításának/felbontásának költségei (nem csak a használatuk költségei) is fontosak lehetnek.*

**Kérdés 1.3:** *Gyakorlati, valós rendszerekben vajon meg lehet-e figyelni az elmélet által megjósolt viselkedést?*

A tézisek második csoportja szinten a felhő szolgáltatás topológiájához kapcsolódik, de ezúttal az árazási kérdéseket veszi szemügyre. Egy komplex modellt állít fel amely a felhasználói igényeket, pénzügyi keretet, és a szolgáltatói kapacitásokat (topológiával együtt) figyelembe veszi. Megállapítható, hogy a modell exact megoldása túl költséges (NP-teljes problémára vezet) ezért egy sztochasztikus modellt vezet be a felhasználók viselkedésére, és ezt felhasználva elemzi az árképzés alakulását. A sztochasztikus modellben a felhasználók vakon döntenek, abban az értelemben, hogy nem ismernek árakat, csak véletlenszerű büdzséjük van, és ha a szolgáltató ennél többre kerül, a szolgáltatás egyszerűen nem valósul meg. Ebben a modellben számos egyszerű topológiában elemzi a szerző a kialakuló árakat például a szolgáltatói kapacitás függvényében.

**Kérdés 1.4:** *A véletlenszerű (vak) felhasználói viselkedés nem rugaszkodik-e el a valóságtól túlságosan? Sok esetben majdnem triviális eredményre vezet, pl ha a szolgáltató kapacitása végtelen, akkor az ár a nullához tarthat, mert érdekében áll mindenkint kiszolgálni. Azonban ez teljesen máshogy alakulna, ha a felhasználók tudatában lennének az aktuális áraknak és a kapacitásnak is, hiszen akkor a büdzsé is tartana a nullához.*

## 2. Téziscsoport: nagyléptékű peremfelhők erőforráskezelése

Ebben a téziscsoportban a szerző olyan algoritmusokat javasol, amelyek peremfelhőkben tesznek lehetővé hatékony erőforráskiosztást, különös tekintettel arra a problémára, hogy a helyőrzők (placeholders) kijelölése optimálisan történjen (azaz a legkevésbé redundáns módon, azonos robusztusság mellett).

Elsősorban egy olyan online módszert javasol, amellyel a placeholderek kiosztása nincs túl messze az optimálistól egy olyan egyszerű heurisztikát használva, hogy ha a pod környezetében már van placeholder, akkor azt fogja használni. Jelentősen leegyszerűsített feltevések mellett belátható, hogy a módszer 3-approximációs.

**Kérdés 2.1:** *A jelentősen egyszerűsítő 3.1 feltevések elhagyása milyen módon nehezíti az elméleti elemzést? Milyen általánosabb eset lenne vizsgálható? A gyakorlatban mennyire teljesülnek az itteni feltevések?*

Ezt követően egy offline módszert javasol a hozzárendelések további optimalizálására.

**Kérdés 2.2:** *Az offline módszer esetén nem lenne lehetséges az optimalizálási problémát valamilyen standard alakra hozni és egy hatékony solverrel megoldani? A gyakorlatban előforduló problémaméretek ezt nem teszik lehetővé?*

A probléma méretének a csökkentésére a szerző klaszterező algoritmust fejleszt ki, ami a csomópontokat olyan klaszterekbe sorolja, amelyek nem átfedőek, és tetszőleges pár között a késleltetés adott érték alatt van. Megmutatja, hogy a probléma NP-teljes, majd egy egyszerű heurisztikát ad, amely ugyanakkor nem garantált, hogy megfelelő megoldást ad vissza.

Ennek az eredménynek a bemutatása kissé nehezen követhető, a 3.3 tétel tautológikusnak tűnik, amennyiben azt állítja, hogy ha a polinomiális idejű heurisztika talál megoldást, akkor a megoldási idő polinomiális. A heurisztika megadásánál a fogalmak használata furcsa, a „deterministic” jelző használatát nem értem. A 3.11 lemma pedig egyenesen helytelennek tűnik: a bizonyítás vége az hogy „If all of the disconnected components are complete subgraphs (cliques), the output is positive (yes), G can be clustered based on d deterministically, otherwise it cannot.” Világos, hogy létezik determinisztikus algoritmus, amely megfelelő klaszterezést ad, akkor is, amikor a heurisztika nem működik (pl további éleket kell törölni valami determinisztikus módon, amíg a maradék összefüggő részgráfok mind klikkek, ami bekövetkezik, legkésőbb akkor, amikor már nincs él).

**Kérdés 2.3:** *Léteznek olyan szimulátorok vagy trace-ek amiken az eddigiekben megadott algoritmusokat a state-of-the-art algoritmusokkal empirikusan össze lehet hasonlítani? Voltak ilyen kutatások?*

A téziscsoport következő kontribúciója egy modell amely a valóságban megfigyelhető skálázást képes modellezni abból a célból, hogy a skálázó módszerek szimulációban kiértékelhetők legyenek.

**Kérdés 2.4:** *A 3.5 ábrán az osztályozó alkalmazás esetén a lineárisnál jobban nő a hatékonyság. Ez hogyan lehet?*

Ezt követi egy olyan elemzés, amelynek során a szerző több prediktív módszert hasonlít össze abból a szempontból, hogy a skálázást ezekkel milyen mértékben lehet javítani az alapértelmezett kubernetes algoritmushoz képest. A módszerek amelyekről szó van: autóregresszió, LSTM, HTM, és Q-tanulás. A javasolt módszer végül mindezen módszereket kombinálja HPA+ néven, mivel más időpontokban más és más módszerek bizonyulnak jobbnak.

**Kérdés 2.5:** *A 3.12 ábrán nem látszik, hogy az „SLA violation” és a „POD minutes” metrikák külön külön hogy változnak a HPA és HPA+ tekintetében. De mintha arról lenne szó, hogy a HPA+ növeli a POD perceket és csökkenti az SLA megsértések számát? Ugyanez a hatás nem érhető el a sima HPA paraméterezésével?*

Végül egy módszert mutat be a szerző, amellyel az alkalmazások modulokra darabolása végezhető el a hatékonyságot maximalizálva.

**Kérdés 2.6:** *A konklúzió kitér arra, hogy a módszer nem veszi figyelembe a modulok közötti kommunikációt. Ez a gyakorlatban mekkora limitáció? Van esetleg konkrét példa rá, amikor megtehető, hogy nem a kommunikációs gráf alapján csoportosítjuk egy konténerbe a komponenseket?*

### **3. Téziscsoport: sávszélesség megosztás**

Ebben a téziscsoportban a szerző olyan módszerekkel foglalkozik, amik a sávszélesség megosztásának a problémájával foglalkoznak. Az első részben egy játékelméleti megközelítést mutat be, amelynek során definiál egy iterált játékot, ahol a játékosok a sávszélességért licitálnak. Több stratégia lehetséges, amelyeknek az összehasonlítását empirikusan végzi el. Érdekes eredmény, hogy a legsikeresebb stratégia annak a függvénye, hogy mekkora terhelés alatt van a rendszer.

**Kérdés 3.1:** *Nem volt világos a kutatás motivációja: egy szolgáltatót mi ösztönözheti arra, hogy megengedjen egy viszonylag előrejelezhetetlen dinamikájú játékelméleti felállást ahelyett, hogy valamilyen fair algoritmus alapján osztana sávszélességet?*

A második részben a szerző definiál egy problémát, melynek alapfeltevése, hogy kamerák osztoznak limitált feltöltési sávszélességen, és a felhőben ezeknek a kameráknak a segítségével egy alkalmazás fut, amelynek feladata, hogy valamilyen eseményeket detektáljon. A feladat az, hogy rendeljük hozzá a limitált sávszélességet a kamerákhoz annak a függvényében, hogy a kérdéses események detekciója minél jobban megvalósuljon. A szerző definiál egy dinamikus programozási feladatot bizonyos feltevések mellett, és ezt oldja meg.

**Kérdés 3.2:** *A feladat nagyban épít arra, hogy meg tudjuk jósolni, hogy egyes kamerák mennyire hasznosak lesznek detekció szempontjából. Vajon ez mennyire realisztikus feltevés? Nekem nehezebb feladatnak tűnik ennek a megjósolása, mint aztán ez alapján egy optimális sávszélesség kiosztás. Ha pedig nem jók a jóslatok, akkor nem biztos hogy érdemes szofisztikált módszereket ezekre építeni.*

## Összegzés

A disszertáció összességében jól érthetően van megfogalmazva. Nem találtam optimálisnak, hogy a vonatkozó irodalom teljesen izoláltan van tárgyalva az eredményektől, fontos lett volna integráltabban motiválni az egyes feladatokat, illetve a megközelítések újszerűségét.

A disszertációban megfogalmazott tézisek mindegyikét új tudományos eredménynek fogadom el, amely a szerző saját munkája.

A doktori művet az MTA doktora cím elnyeréséhez elegendőnek, egyúttal a nyilvános vitára alkalmasnak tartom.

2023. 09. 30. Szeged



Jelasity Márk

egyetemi tanár

Szegedi Tudományegyetem