

Kétsztályos WFQ kiszolgálás közelítő vizsgálata

HORVÁTH GÁBOR, TELEK MIKLÓS

Budapest Műszaki és Gazdaságtudományi Egyetem, Híradástechnikai Tanszék
{ghorvath, telek}@hit.bme.hu

Kulcsszavak: WFQ kiszolgálás, kétdimenziós Markov-lánc, várakozási idő várható értéke és szórása

A súlyozott igazságos kiszolgálási elvet (Weighted Fair Queueing – WFQ) régóta számos távközlési és számítástechnikai rendszerben alkalmazzák erőforrások megosztásának szabályozására. Látszólagos egyszerűsége ellenére a WFQ hatékony analitikus teljesítményvizsgálata még mindig nyitott kérdés. Az irodalomban megtalálható – numerikus, vagy komplex analízisen alapuló – algoritmusok gyakorlati alkalmazhatósága meglehetősen korlátozott. Ebben a cikkben egy egyszerű és gyors közelítő eljárást mutatunk be WFQ rendszerek vizsgálatára. Szimulációs eredményekkel igazoljuk, hogy egyszerűsége ellenére az ismertetett megközelítés megfelelően pontos.

1. Bevezetés

A súlyozott igazságos (WFQ) kiszolgálási elvet erőforrások megosztására alkalmazzák többosztályos környezetben, melyben az igények igényosztályokba sorolhatók. Minden igényosztályhoz tartozik egy súly. WFQ rendszerekben a teljes kiszolgálási kapacitás megosztását az igényosztályok között a pillanatnyilag jelenlévő igényosztályok súlyának aránya határozza meg. Ily módon az igények „fontossága” a súlyok segítségével szabályozható. Az igények torlódás esetén sorban állnak a kiszolgálóért, FCFS (first come first serve) elv szerint.

Ha az igények érkezési folyamatát Poisson folyamattal, a kiszolgálási időket pedig exponenciális eloszlással jellemezzük, akkor a WFQ rendszer egy kétdimenziós Markov-lánccal modellezhető.

Az irodalomban számos eljárás található ennek a kétdimenziós Markov-lánccal a megoldására. Először a numerikus eljárásokat vesszük sorra: [1,2]-ben a szerzők ugyanilyen működésű rendszer megoldásával foglalkoztak, bár *Coupled Processor Model*-nek hívták. A Markov-lánc egyensúlyi eloszlását a rendszer terhelésének a hatványsoraként fejezték ki, és adtak egy algoritmust a hatványsor együtthatóinak rekurzív kiszámolására. Ezzel a megközelítéssel csak néhány (2-3) igényosztályos rendszert tudtak kiszámolni, és ahogy a terhelés 1-hez közelít, túl sok együtthatót kell kiszámolni.

Egy másik megoldás ([3]) a végtelen Markov-lánccal egy véges Markov-lánccal közelíti, és egyfajta Gauss-eliminációt használ az egyensúlyi eloszlás kiszámítására. A Gauss-elimináció közben kihasználja a mátrix speciális struktúrája adta gyorsítási lehetőségeket. De ebben az esetben is, nagy terhelés mellett a végesített Markov-lánccal is túl sok állapota lesz, és a megoldás drasztikusan lelassul.

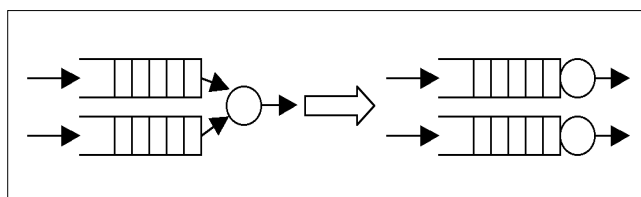
[4]-ben a szerzők felírják az egyensúlyi eloszlás generátorát (Laplace transzformáltját). Az eredmény egy

kétféle változós komplex (éppenséggel analitikus) függvény, ahol a problémát a peremeloszlások kiszámítása jelenti, ehhez ugyanis Wiener-Hopf faktorizációt kell alkalmazni.

Ebben a cikkben 2 osztályos WFQ rendszerrel foglalkozunk, de az ismertetett módszer egyszerűen kiterjeszhető többosztályos rendszerekre is. A fent ismertetett megoldásokkal ellentétben az igények érkezési és kiszolgálási idejét nem exponenciális eloszlásúnak tételezzük fel, hanem két momentumot veszünk figyelembe. Közelítést adunk az igények várakozási idejének várható értékére és szórására.

2. A közelítés elve

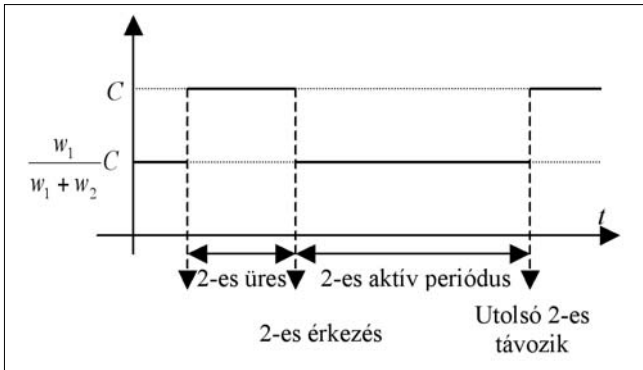
A közelítés lényege, hogy a két igényosztályt szeparáltan vizsgáljuk (1. ábra). Egy olyan kiszolgálási folyamatot konstruálunk mindkét igényosztály számára, ami az eredeti kiszolgáló viselkedését „imitálja”.



1. ábra Az igényosztályok szétválasztása

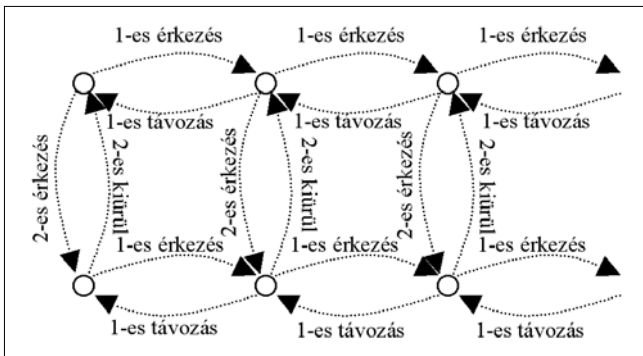
Például az 1. osztályt tekintve látható, hogy a kiszolgálási kapacitás a teljes kapacitás (C) és a súlyoknak megfelelően csökkentett kapacitás között változik attól függően, hogy van-e a rendszerben 2-es típusú igény.

Az ötlet egyszerű: jellemezzük a 2-es típusú igények aktív periódusának hosszát, és konstruáljunk egy olyan modulált kiszolgálási folyamatot, ahol a kapacitás modulációját a 2-es típusú igények aktív periódusa határozza meg (2. ábra).



2. ábra A modulált kiszolgálási folyamat

A két igénytípus szeparálása után az így már egyosztályossá váló sorbanállási rendszert kvázi születés-halálózási (QBD) folyamattal modellezzük, és mátrix geometrikus technikával oldjuk meg. Az így keletkező (egy igényosztályt modellező) Markov-láncon az állapotokat duplikáljuk: az egyik csoport a 2-es igények aktív, a másik csoport a passzív állapotához tartozik. A Markov-lánc makró szintű struktúráját a 3. ábra szemlélteti.



3. ábra A közelítő Markov-lánc szerkezete

A részletes tárgyalás előtt foglaljuk össze a megoldás menetét:

- Fázis-típusú eloszlást konstruálunk az igények érkezési és kiszolgálási idejének modellezésére. A fázis-típusú eloszlások használata teszi lehetővé a sormodellek mátrix-geometrikus megoldását.
- Kiszámoljuk a két igényosztály aktív periódusának a hosszát. Ez a lépés négy eredményt fog adni, az aktív periódus hosszát a két osztályra a két lehetséges kiszolgálási kapacitás (teljes, illetve súlyoknak megfelelően csökkentett kapacitás) mellett.
- A 3. ábrának megfelelően összeállítjuk és megoldjuk az igényosztályok sorbanállását egyenként jellemző Markov-láncon. A teljesítményjellemzőket ennek a Markov-láncnak a segítségével kapjuk meg.

2.1. Az érkezési és kiszolgálási folyamat

A továbbiakban az i igényosztályhoz tartozó mennyiségeket és jeleket az (i) index jelöli (ebben a cikkben két igényosztállyal foglalkozunk).

Az érkezési folyamatot két mennyiség jellemzi: az érkezési intenzitás $\lambda^{(i)}$, valamint az érkezési időközök relatív szórásnégyzete $c_A^{2(i)}$.

Az így adott két momentum alapján egy olyan másodrendű, aciklikus fázis típusú eloszlást készítünk (PH, [6]), melynek ugyanez az első két momentuma. A PH eloszlásoknak három paramétere van: a tranzien átmenetek generátor mátrixa $D^{(i)}$, a nyelőbe vezető átmenetek rátavektora $d^{(i)}$, valamint a kezdőállapot-valószínűség vektor $\delta^{(i)}$. Könnyen ellenőrizhető, hogy a következő PH eloszlás első két momentuma megegyezik az érkezési időközök első két momentumával:

$$D^{(i)} = \begin{bmatrix} -\lambda^{(i)} & \lambda^{(i)} \\ c_A^{2(i)} & -2\lambda^{(i)} \end{bmatrix}, \quad d^{(i)} = \begin{bmatrix} 0 \\ 2\lambda^{(i)} \end{bmatrix}$$

$$\delta^{(i)} = \begin{bmatrix} 1 \\ 2c_A^{2(i)} \end{bmatrix} 1 - \frac{1}{2c_A^{2(i)}}$$

Az igények által a rendszerbe hozott munka szintén két paraméterrel adott: a várható értékével $m_f^{(i)}$, és a relatív szórásnégyzetével $c_{sf}^{(i)}$. Ezekből a paraméterekből kiszámoljuk a kiszolgálási idő teljes kiszolgálási kapacitást feltételezve:

$$\mu_f^{(i)} = \frac{C}{m_f^{(i)}}, \quad c_{sf}^{2(i)} = c_f^{2(i)},$$

és csökkentett kiszolgálási kapacitást feltételezve:

$$\mu_r^{(i)} = \frac{w_i C}{\sum_i w_i m_f^{(i)}}, \quad c_{sr}^{2(i)} = c_f^{2(i)},$$

ahol C jelöli a kiszolgáló kapacitását, és w_i jelöli az i igényosztály súlyát. A fentiekhez hasonlóan PH eloszlást illesztünk a kiszolgálási időkre is:

$$S_f^{(i)} = \begin{bmatrix} -\mu_f^{(i)} & \mu_f^{(i)} \\ c_{sf}^{2(i)} & -2\mu_f^{(i)} \end{bmatrix}, \quad s_f^{(i)} = \begin{bmatrix} 0 \\ 2\mu_f^{(i)} \end{bmatrix}$$

$$\sigma_f^{(i)} = \begin{bmatrix} 1 \\ 2c_{sf}^{2(i)} \end{bmatrix} 1 - \frac{1}{2c_{sf}^{2(i)}}.$$

A csökkentett kapacitáshoz tartozó kiszolgálási idők PH leírása ($S_r^{(i)}$, $s_r^{(i)}$, $\sigma_r^{(i)}$) hasonlóan számítható.

2.2. Az aktív periódus

Vessünk újra egy pillantást a 3. ábrára. Az érkezésekkel, illetve kiszolgálási idővel kapcsolatos átmenetek az előző fejezet eredményei alapján már ismertek. Az egyetlen hiányzó átmenet a 2-es osztály kiürülési ideje, vagyis az aktív periódusának a hossza. Ebben a fejezetben az aktív periódus két momentumát számítjuk ki, teljes, valamint csökkentett kiszolgálási kapacitás mellett, a másik igényosztály hatását figyelmen kívül hagyva.

Mivel a sorok érkezési és kiszolgálási idejének eloszlása fázis-típusú, ezért az egyenként vizsgált sorok Markov-láncának szerkezete speciális mátrix-tridiagonál, vagyis kvázi születés-halálózási folyamat:

$$Q = \begin{bmatrix} A_1' & A_0' & & & \\ A_2' & A_1 & A_0 & & \\ & A_2 & A_1 & A_0 & \\ & & \ddots & \ddots & \ddots \end{bmatrix}$$

A generátor mátrix blokkjait az érkezési és kiszolgálási folyamat paramétereiből az alábbi módon kapjuk meg [6]:

$$\begin{aligned} A_{0f}^{(i)} &= d^{(i)} \delta^{(i)} \otimes I_{2 \times 2} \\ A_{1f}^{(i)} &= D^{(i)} \oplus S_f^{(i)} \\ A_{2f}^{(i)} &= I_{2 \times 2} \otimes s_f^{(i)} \sigma_f^{(i)} \\ A_{0f}^{(i)'} &= d^{(i)} \delta^{(i)} \otimes \sigma_f^{(i)} \\ A_{1f}^{(i)'} &= D^{(i)} \\ A_{2f}^{(i)'} &= I_{2 \times 2} \otimes s_f^{(i)} \end{aligned}$$

(A redukált kapacitás melletti mátrixok hasonlóképpen számíthatók.)

Egy érkező igény által indított aktív periódus k -adik momentuma ($m_{Bf}^k{}^{(i)}$) a következőképpen kapható meg:

$$m_{Bf}^k{}^{(i)} = (-1)^k (\delta^{(i)} \otimes \sigma_f^{(i)}) \frac{d^k}{ds^k} G_f^{(i)}(s) |_{s=0} h,$$

ahol $G_f^{(i)}(s)$ mátrix kielégíti a következő mátrix egyenletet:

$$s G_f^{(i)}(s) = A_{2f}^{(i)} + A_{1f}^{(i)} G_f^{(i)}(s) + A_{0f}^{(i)} (G_f^{(i)}(s))^2.$$

A $G_f^{(i)}(s)$ mátrix 0-dik deriváltja az $s=0$ helyen az úgynevezett *fundamental matrix*, melynek minél hatékonyabb megoldásával számos cikk és könyv foglalkozik [6]. Az első derivált kiszámításához fix-pont iterációt használtunk (csupa 0 elemű mátrixból kiindulva) az alábbi alakra hozott egyenlet segítségével:

$$\begin{aligned} \frac{d}{ds} G_f^{(i)}(s) |_{s=0} &= (A_{1f}^{(i)} - A_{0f}^{(i)} G_f^{(i)}(0))^{-1} \cdot \\ &\cdot \left(I - A_{0f}^{(i)} \frac{d}{ds} G_f^{(i)}(s) |_{s=0} \right) G_f^{(i)}(0) \end{aligned}$$

Közelítésünkben az aktív periódus hosszának két momentumát fogjuk felhasználni. Az eddigi eredményeket összegezve ez a két momentum az alábbi módon számítható ki:

$$\begin{aligned} m_{Bf}^1{}^{(i)} &= -(\delta^{(i)} \otimes \sigma_f^{(i)}) (A_{2f}^{(i)} + A_{0f}^{(i)} + A_{0f}^{(i)} G_f^{(i)}(0))^{-1} h \\ m_{Bf}^2{}^{(i)} &= 2(\delta^{(i)} \otimes \sigma_f^{(i)}) (A_{2f}^{(i)} + A_{0f}^{(i)} + A_{0f}^{(i)} G_f^{(i)}(0))^{-1} \cdot \\ &\cdot \left(A_{0f}^{(i)} \frac{d}{ds} G_f^{(i)}(s) |_{s=0} - I \right) \cdot m_{Bf}^1{}^{(i)} \end{aligned}$$

Ebből a két momentumból kiindulva fázis-típusú eloszlást konstruálunk, hasonlóan, ahogy azt az érkezési és kiszolgálási idők esetében tettük. A kapott PH eloszlás paramétereit jelölje ($B_f^{(i)}$, $b_f^{(i)}$, $\beta_f^{(i)}$).

A 3. ábra Markov-láncának felépítéséhez még egy dolog hiányzik. Tudnunk kell, hogy az aktív periódus befejeztével milyen fázisban lesz az adott sor érkezési folyamata. Ez a következő egyenlet segítségével kapható:

$$\alpha_f^{(i)} = (\delta^{(i)} \otimes \sigma_f^{(i)}) \cdot G_f^{(i)}(0) \cdot (h_2 \otimes I_{2 \times 2}).$$

2.3. A sorbanállási modell

Ebben a fejezetben a 3. ábrán felvázolt, a rendszert az i igényosztály szemszögéből leíró Markov-láncot konstruáljuk meg. A másik igényosztályt j -vel fogjuk jelölni (vagyis ha $i=1$, akkor $j=2$, és fordítva). Amint az az ábrán látható, az állapottér két részre osztható, az alsó, és a felső állapotsorra.

A felső részben, ahol nincsen a rendszerben másik típusú igény, az i igényosztály a kiszolgáló teljes kapacitását használhatja. Az állapottérnek ezen a részén követni kell (1) az i igényosztály érkezési folyamatának a fázisát, (2) az i igényosztály kiszolgálási idejének a fázisát, végül (3) a j igényosztály érkezési folyamatának a fázisát. Mivel mindhárom esetben 2 fázisú eloszlásról van szó, ez eddig 8 állapotot jelent.

Ha megjelenik a j típusú igény a rendszerben, a kiszolgáló kapacitása arányosan megoszlik, és az i igényosztálynak csökkentett kapacitás jut. A j típusú igény megjelenése ezért a Markov-láncot az állapottér alsó részébe viszi, ahol a kiszolgálás lassabb sebességgel történik. Ebben a részben követni kell az (1)-es és (2)-es fázisokat, de a harmadik nyomon követendő dolog a j igényosztály aktív periódus eloszlásának a fázisa lesz. Ez ismét 8 állapotot jelent. Így a teljes Markov-láncnak összesen 16 állapota lesz.

Már a 3. ábrából is látható, hogy ismét kvázi születési-halálzási folyamatot kapunk, melynek mátrixblokkjai – az eddig leírtak alapján – eképpen állnak össze:

$$\begin{aligned} C_0^{(i)} &= \begin{bmatrix} d^{(i)} \delta^{(i)} \otimes I \otimes I & 0 \\ 0 & d^{(i)} \delta^{(i)} \otimes I \otimes I \end{bmatrix} \\ C_1^{(i)} &= \begin{bmatrix} D^{(i)} \oplus D^{(j)} \oplus S_f^{(i)} & I \otimes d^{(j)} \beta_r^{(j)} \otimes I \\ I \otimes b_r^{(j)} \alpha_r^{(j)} \otimes I & D^{(i)} \oplus B_r^{(j)} \oplus S_r^{(i)} \end{bmatrix} \\ C_2^{(i)} &= \begin{bmatrix} I \otimes I \otimes s_f^{(i)} \sigma_f^{(i)} & 0 \\ 0 & I \otimes I \otimes s_r^{(i)} \sigma_r^{(i)} \end{bmatrix} \end{aligned}$$

Az irreguláris 0. szint mátrixai a következők:

$$\begin{aligned} C_0^{(i)'} &= \begin{bmatrix} d^{(i)} \delta^{(i)} \otimes I \otimes \sigma_f^{(i)} & 0 \\ 0 & d^{(i)} \delta^{(i)} \otimes I \otimes \sigma_r^{(i)} \end{bmatrix} \\ C_1^{(i)'} &= \begin{bmatrix} D^{(i)} \oplus D^{(j)} & I \otimes d^{(j)} \beta_f^{(j)} \\ I \otimes b_f^{(j)} \alpha_f^{(j)} & D^{(i)} \oplus B_f^{(j)} \end{bmatrix} \\ C_2^{(i)'} &= \begin{bmatrix} I \otimes I \otimes s_f^{(i)} & 0 \\ 0 & I \otimes I \otimes s_r^{(i)} \end{bmatrix} \end{aligned}$$

Mivel a fázisok száma kicsi (szintenként 16 állapot), a klasszikus QBD megoldó algoritmusok nagyon gyorsan (1 másodpercen belül) képesek kiszámítani a teljesítményjellemzőket, közöttük az esetünkben fontos várakozási idő momentumokat [6].

3. Numerikus eredmények

Hogy az eredmények használhatóságát bemutassuk, két példát állítottunk össze. Az első esetben a két igényosztály azonos mennyiségű munkát hoz a rendszerbe,

míg a másodikban a 2-es típusú igények tízszer annyit kiszolgálási időt igényelnek, mint az 1-es típusúak.

A példákban bemutatott görbék nem csak a cikkben ismertetett analitikus eljárás eredményeit tartalmazzák, hanem a szimulációval kapott eredményeket is, így lemérhető a közelítő eljárás pontossága.

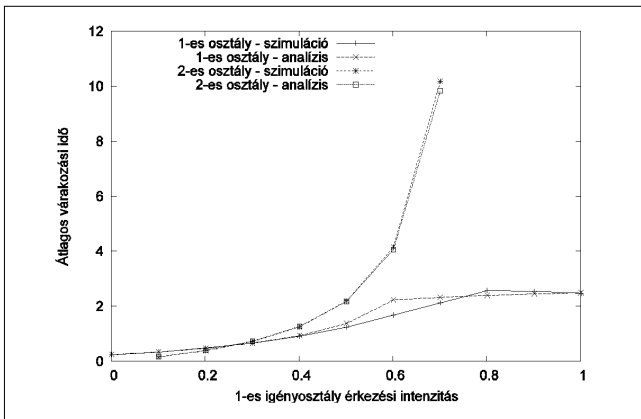
1. példa

A 4. és 5. ábrán a várakozási idő várható értékét és szórását ábrázoltuk, miközben az 1-es osztály terhelé-

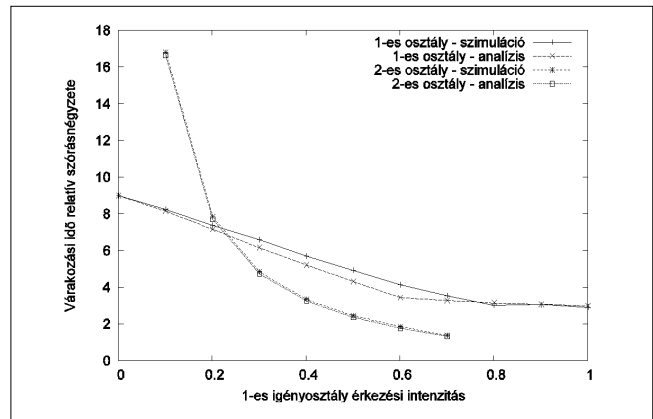
sét (forgalmi intenzitását) növeltük. Amint várható volt, a terhelés növekedésével a várható várakozási idő nő, a relatív szórásnégyzet csökken.

A 6.-9. ábrákon azt ábrázoltuk, hogy milyen hatással van a várakozási időre, ha a kiszolgálási idő, vagy az érkezési idő szórását növeljük.

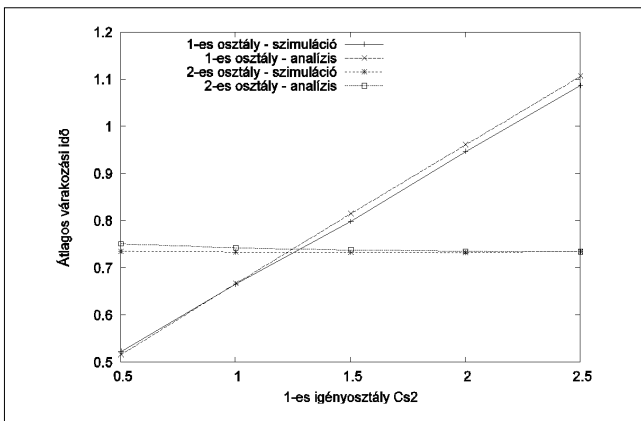
A bemutatott ábrák alapján elmondható, hogy a szórások növelése csak minimális mértékben van hatással a másik igényosztály vizsgált teljesítményjelzőire.



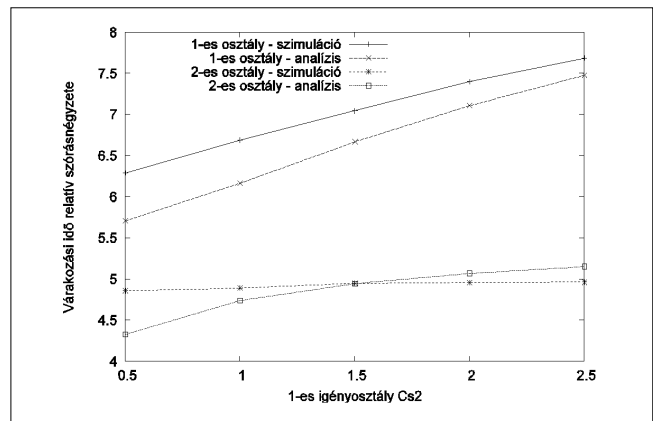
4. ábra Átlagos várakozási idő az érkezési intenzitás függvényében



5. ábra Várakozási idő relatív szórásnégyzete az érkezési intenzitás függvényében



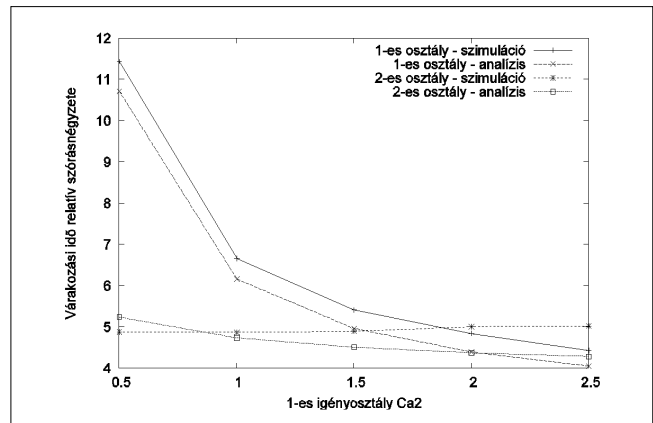
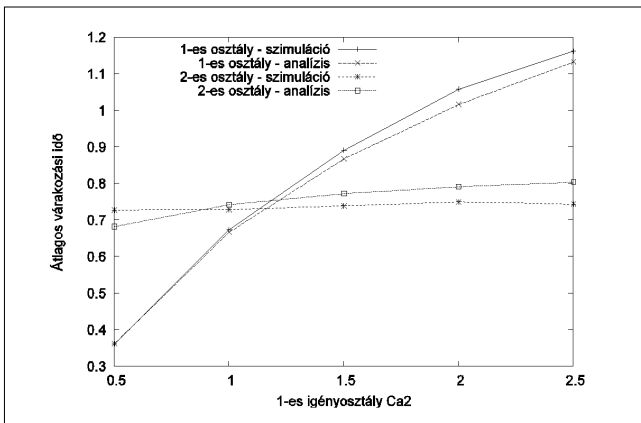
6. ábra Átlagos várakozási idő a kiszolgálási idő relatív szórásnégyzetének a függvényében

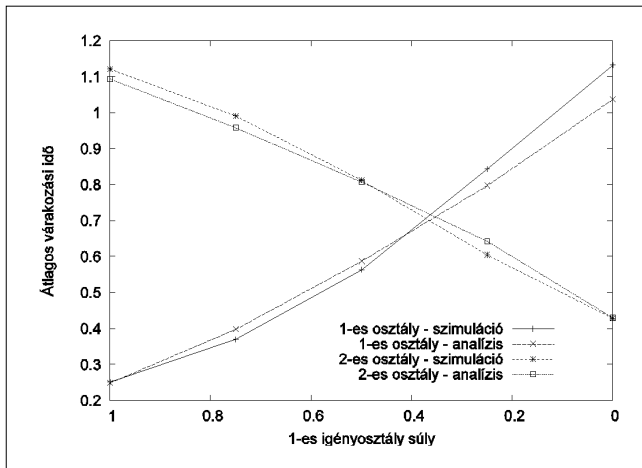


7. ábra Várakozási idő relatív szórásnégyzete a kiszolgálási idő relatív szórásnégyzetének a függvényében

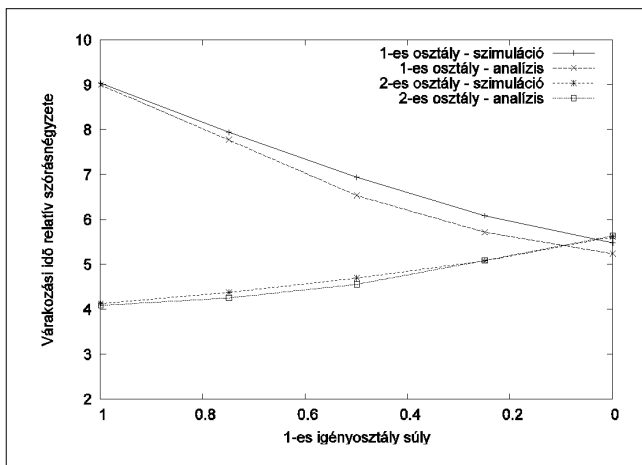
8. ábra Átlagos várakozási idő az érkezési időközők relatív szórásnégyzetének a függvényében

9. ábra Várakozási idő relatív szórásnégyzete az érkezési időközők relatív szórásnégyzetének a függvényében





10. ábra Átlagos várakozási idő a forgalmi osztály súlyának függvényében



11. ábra Várakozási idő relatív szórásnégyzete a forgalmi osztály súlyának függvényében

2. példa

Ebben a példában a 2-es típusú igények tízszer annyi kiszolgálási időt igényelnek, mint az 1-es típusúak.

A 12. és 13. ábra a várakozási idő alakulását mutatja a terhelés függvényében. Látható, hogy még a 2-es osztály túlterhelése esetén is az 1-es típusú igények megkapják a garantált, súlylal csökkentett kiszolgáló kapacitást, így a várakozási idő nem nő jelentősen.

A 14.-17. ábrák az érkezési és kiszolgálási idők szórásának hatását szemléltetik.

Ismét látható, hogy a 2-es típusú igények várakozási idejét nem befolyásolja az 1-es osztály érkezési és kiszolgálási idejének szórása. Igaz ugyan, hogy az analitikus megoldás mutat egy kis összefüggést, de a szimuláció görbéjétől való távolság még mindig elfogadható.

A 18. és 19. ábrán a várakozási idő két paraméterét a súlyok beállításának függvényében ábrázoltuk.

A várható várakozási idő közelítése jó, de a relatív szórásnégyzet esetén itt tapasztaltuk a legnagyobb különbséget a szimuláció és az analízis között (kb. 20%).

4. Összefoglalás

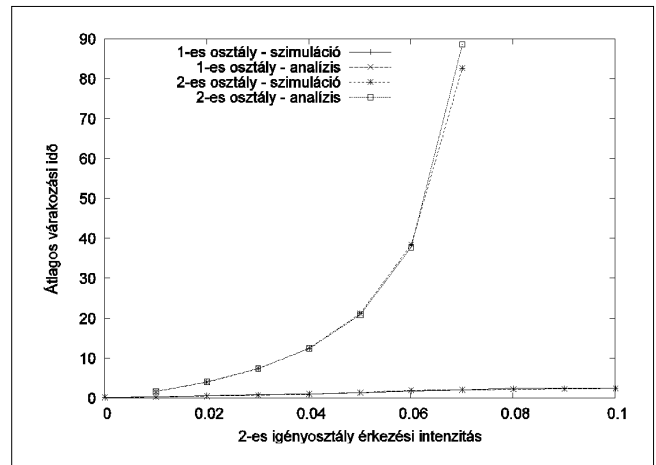
Ebben a cikkben egy közelítő eljárást ismertetünk WFQ rendszerek teljesítményvizsgálatára. A közelítés elvének bemutatása után klasszikus mátrix-geometriai eszközökkel oldottuk meg a felmerülő sorbanállás-elméleti problémákat.

Módszerünk előnye, hogy nagyon kicsi a számítási igénye, továbbá hogy az irodalomban látott korábbi megoldásoknál általánosabb, mivel az érkezési és kiszolgálási időközöknek nem csak a várható értékét, hanem a szórását is figyelembe veszi.

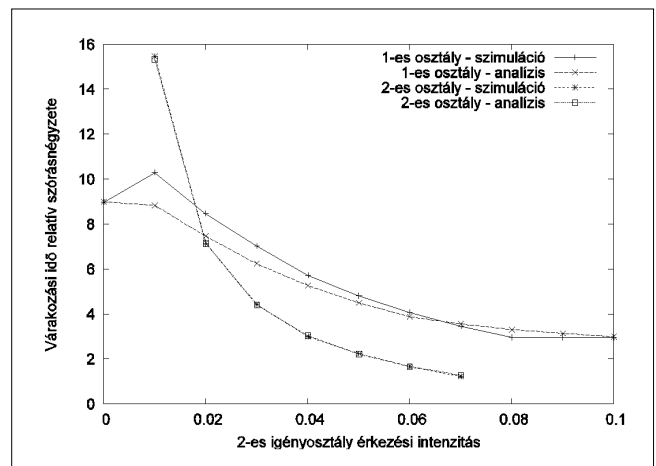
Két numerikus példán keresztül kimerítően vizsgáltuk az eljárás pontosságát, minden lehetséges paraméter függvényében. Az esetek nagy részében az eredmények nagyon jól közelítették a szimuláció eredményeit.

A legnagyobb eltérést (15-20%) a várakozási idő relatív szórásnégyzetében tapasztaltuk, a várakozási idő várható értékére azonban minden esetben jó közelítést kaptunk.

12. ábra Átlagos várakozási idő az érkezési intenzitás függvényében



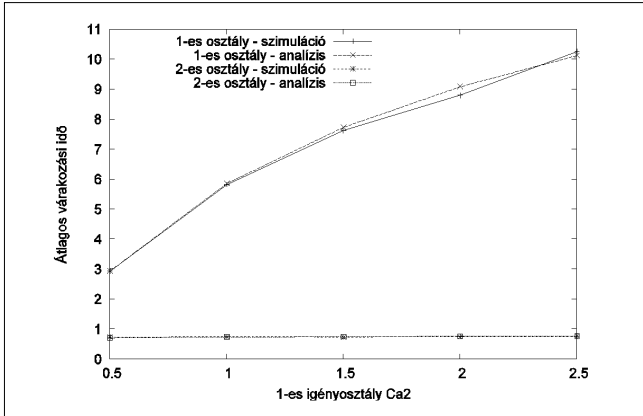
13. ábra Várakozási idő relatív szórásnégyzete az érkezési intenzitás függvényében



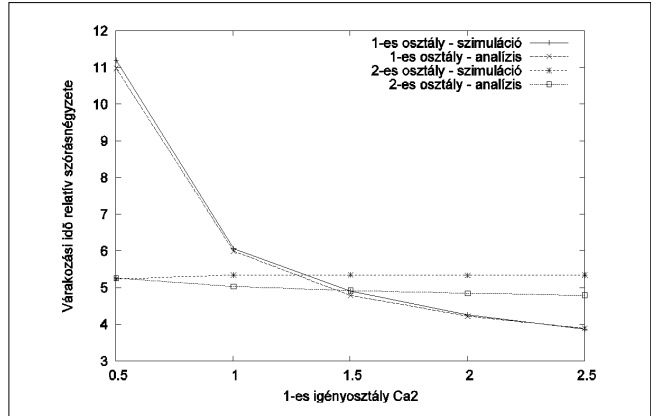
Irodalom

[1] G. Koole,
 "On the power series algorithm" (1994)
 Tech. Rep., Centrum voor Wiskunde en Informatica
 [2] J. P. C. Blanc (1988)
 "A numerical study of the coupled processor model",
 in Computer Performance and Reliability
 [3] Leslie D. Servi (2002)
 "Algorithmic solutions to two-dimensional birth-death
 processes with application to capacity planning",
 Telecom. Systems, Vol. 21, No.2, pp.205–212.

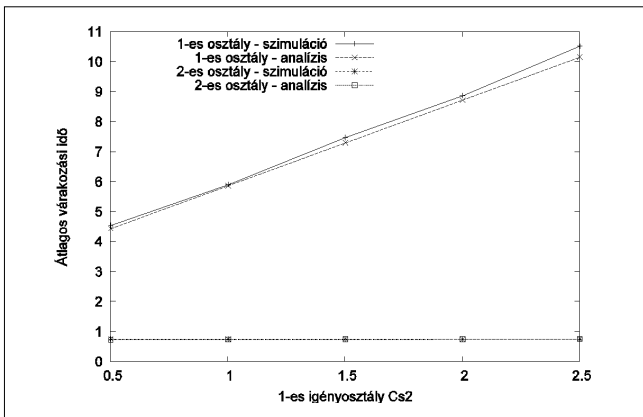
[4] F. Guillemin, R. Mazumdar, A. Dupuis, J. Boyer (2003)
 "Analysis of the fluid weighted fair queueing system",
 J. Appl. Probab., Vol. 40, No.1, pp.180–199.
 [5] G. Fayolle, R. Iasnogrodski, V. Malyshev (1999)
 Random Walks in the Quarter Plane,
 Springer-Verlag New York
 [6] G. Latouche, V. Ramaswami (1999)
 Introduction to Matrix Analytic Methods in
 Stochastic Modeling,
 American Statistical Association and the Society for
 Industrial and Applied Mathematics



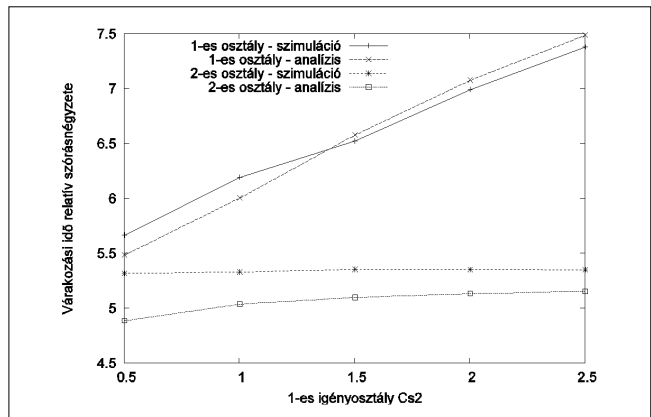
14. ábra Átlagos várakozási idő az érkezési időközök relatív szórásnégyzetének a függvényében



15. ábra Várakozási idő relatív szórásnégyzete az érkezési időközök relatív szórásnégyzetének a függvényében



16. ábra Átlagos várakozási idő a kiszolgálási idő relatív szórásnégyzetének a függvényében



17. ábra Várakozási idő rel. szórásnégyzete a kiszolgál. idő relatív szórásnégyzetének a függvényében

18. ábra Átlagos várakozási idő a forgalmi osztály súlyának függvényében

19. ábra Várakozási idő relatív szórásnégyzete a forgalmi osztály súlyának függvényében

