

**Néhány gyakoribb várakozósoros modell
rendszertervezéshez.**

Dr. Gyarmati G. Péter

1976. július.

Tartalomjegyzék

Néhány gyakoribb várakozósoros modell rendszertervezéshez.

1. Bevezetés	5
2. A várakozósor leírása	6
A forrás	7
Az igények beérkezése	7
A kiszolgálási idő	8
Egy paraméter	8
A várakozósoros rendszer kapacitása	9
A kiszolgáló egységek száma	9
A sorbanállási rend	9
3. A modell leírása	11
A rendszerek rövidített jelölése	11
A rendszer forgalma	12
A kiszolgálási egység foglaltsága	12
Egy valószínűség: t időben n igény van a rendszerben	12
További jellemzők	13
4. Modellek egy kiszolgáló egységgel.	15
Az M/M/1 modell	15
Az M/M/1/K modell	19
Az M/G/1 modell	21
Az M/D/1 modell	22
Az M/E _k /1 modell	23
5. Elsobbégi, prioritásos modellek	24
Kivárásos rendszer, non-preemptive	25
Megszakításos rendszer, preemptive	25
Az M/G/1 prioritásos modell főbb összefüggései	25
6. Összefoglalás	27
7. Mellékletek	28
A. Az alkalmazott jelölések és fogalmak	28
B. Irodalomjegyzék	31

1. Bevezetés

A sorbanállási elmélet, --amelyet telefon hálózatok tervezéséhez A. K. Erlang dán matematikus dolgozott ki-- a számítógépes rendszerek tervezésének és vizsgálatának hasznos és nélkülözhetetlen eszközévé vált a ráfordított adaptációs és fejlesztési munka eredményeként.

Az elmélet segítségével előre meghatározhatók a számítógép teljesítőképességének jellemző adatai, például a várakozási idő egy on-line terminálnál, vagy a tárolóigények üzenetközvetítő rendszerekben, vagy a prioritás hozzárendelések hatása, stb.

Ez a tanulmány ilyen és más hasonló kérdések tárgyalásához szükséges elméleti eszközök bemutatásával --a gyakorlati alkalmazás síkján-- kíván foglalkozni.

A sor a kiszolgálásra várakozó igények tárolója például egy kommunikációs-, vagy számítástechnikai rendszerben, míg a sorbanállási elmélet ezek kezelésének a tanulmányozására, tervezésére alkalmas eszköz.

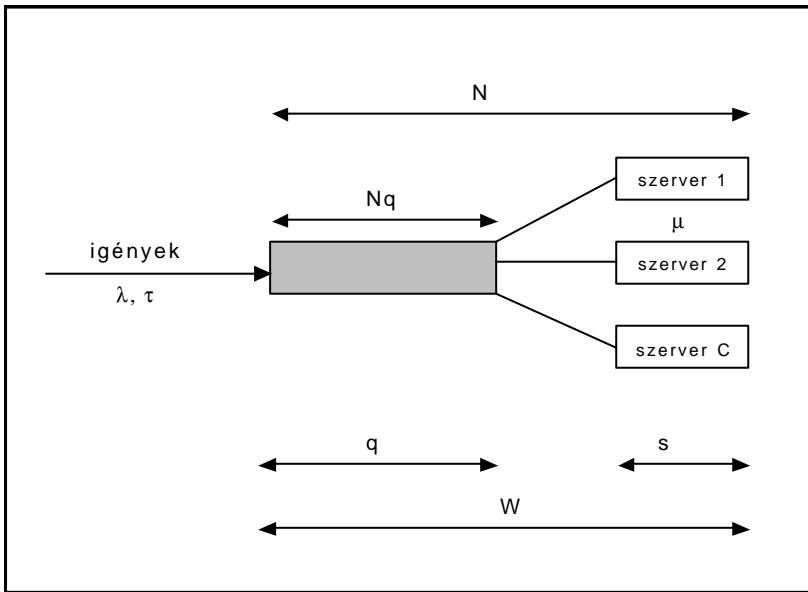
A feldolgozásra várakozó igények egy rendszerben sorok alkotnak. Az ilyen kiszolgálásra várakozó sorok további sorokat alkothatnak, mint például az input-output igények sora egy feldolgozási programban.

2. A várakozósor leírása

Az elméletben az alapelem az igénylo, aki, vagy ami igényli és élvezi a kiszolgálást, ezért az ilyen rendszereket --legyen az kommunikációs, számítástechnikai, informatikai, vagy kereskedelmi, pénzügyi stb. természetű-- kiszolgálási rendszereknek nevezzük és a tervezéshez modellezzük.

Az igénylo által keltett igénylések sorban állnak és várakoznak a kiszolgálásra, ezért várakozósoros modellnek nevezzük.

Az informatikai rendszerekben az igénylo lehet például egy átviendő üzenet, egy feldolgozási program, egy választ váró kérdés, egy adatkérés, egy eredményközlés, vagy ezek bármilyen ismétlése, kombinációja, stb.



1. számú ábra

Az igénylo kiszolgálást kér --üzenet átvitelét egy csatornán, egy program utasításainak végrehajtását, egy kifejezés értelmezését, egy adat megkeresését, adat kérését, átvételét más forrásokból, eredmények közlését, stb.-- egy erre alkalmas kiszolgáló-, szolgáltató forrásból.

A kiszolgáló forrás egy vagy több kiszolgáló egységből áll, amelyek a kívánt szolgáltatást nyújtják az igényloknak.

Ha az összes kiszolgáló egység foglalt, amikor az igénylo belép a rendszerbe, akkor az igény egy várakozó sorba kerül és várakozik addig, ameddig egy kiszolgáló egység a számára fel nem szabadul.

Az 1.sz ábrán látható az általános modell a használatos kifejezésekkel, a változókkal és jelentésükkel.

A sorbanállási rendszert, vagy modellt matematikai tanulmányozásához az alábbi módon írhatjuk le.

A forrás

A forrás, ahonnan az igények származnak, illetve maguk az igények. Az igények forrása lehet véges, vagy végtelen.

Véges forrású rendszerben nem lehet a kiszolgálásra várakozó sor tetszőlegesen hosszú és az igények száma a rendszerben befolyásolja a beérkezési arányt.

Extrém esetben, ha minden igény vár, vagy kiszolgálás alatt van, akkor a beérkezési arány nulla értékű lesz.

Ha a források száma véges, de nagy, akkor végtelen forrást feltételezhetünk az alkalmazott matematikai módszer kezelhetősége érdekében.

Az igények beérkezése

Feltételezzük, hogy az igények $t_0 < t_1 < t_2 < t_3 \dots < t_k \dots$ időben érkeznek a rendszerbe. A beérkezések közötti időt a $\mathbf{t} \approx t_k - t_{k-1}$ (ahol $k \geq 1$) valószínűségi változóval írjuk le, amelyről feltételezzük hogy független és azonos eloszlású valószínűségi változók sorát képezi.

Tetszőleges beérkezések közötti idő esetén a τ jelet használjuk. A beérkezéseket a beérkezési közötti idő eloszlásfüggvényével $A(\mathbf{t}) = P(\mathbf{t} \leq t)$ jellemezzük.

Ha a beérkezési eloszlásfüggvény exponenciális, --a következő igény beérkezésének valószínűsége az idő múlásával exponenciálisan növekszik-- tehát minden \mathbf{t} -ra $P(\mathbf{t} \leq t) = 1 - e^{-\lambda t}$, akkor annak valószínűsége, hogy bármely t időben n igény érkezik be: $e^{-\lambda t} (\lambda t)^n / n!$, ahol $n=0, 1, 2, \dots$ és λ

az átlagos beérkezési arány. Ebben az esetben a beérkezések a Poisson eloszlást követik.

Gyakran használatos még az Erlang- k és a konstans eloszlás. A használatos eloszlásfüggvények leírásai megtalálhatók az [1., 2.] irodalomban.

A kiszolgálási idő

Legyen s_k a k -adik igény kiszolgálásához szükséges idő és feltételezzük, hogy az s_k egy független, azonos eloszlású valószínűségi változó. (Tetszőleges kiszolgálási idő esetén az s jelölést használjuk.)

A kiszolgálási idő eloszlásfüggvénye legyen $W_s(t) = P(s \leq t)$.

Az átlagos kiszolgálási arányt \mathbf{m} -vel jelöljük.

Véletlen kiszolgálás esetén --exponenciális eloszlás-- az eloszlásfüggvény: $W_s(t) = 1 - e^{-\mathbf{m}t}$, ha $t \geq 0$.

A gyakoribb eloszlásfüggvények még az Erlang- k és a konstans eloszlás.

Leírásuk az [1., 2.] irodalomban található.

Egy paraméter

Egy hasznos statisztikai paraméter az $A(t)$ és $W(t)$ eloszlásfüggvények jellemzésére az u.n. négyzetes variációs

együttható: $C_x^2 = \frac{\text{var}(x)}{E(x^2)}$.

Ha x egy konstans valószínűségi változó, akkor $C_x^2 = 0$,

ha x exponenciális eloszlású, akkor $C_x^2 = 1$,

ha x Erlang- k eloszlású, akkor $C_x^2 = 1/k$.

Következtetések:

- a. Ha C_t^2 megközelítőleg nulla, akkor a beérkezéseknek szabályos rendje van, ha közel van az egyhez, akkor a beérkezések véletlenszerűnek jellemezhetők és végül, ha nagyobb egynél, akkor a beérkezések csoportosulására lehet következtetni.

- b. Ha C_s^2 megközelítőleg nulla, akkor a kiszolgálási idő megközelítőleg állandó, nagy értékei pedig a kiszolgálási idő nagy mértékű változásait jelzik

A várakozósoros rendszer kapacitása

Egyes rendszerekben a sor kapacitását végtelennek feltételezik, vagyis minden igény addig várakozhat, amíg kiszolgáláshoz nem jut.

Más rendszerekben a sor kapacitása nulla, vagyis, ha minden kiszolgáló egység foglalt amikor egy igény beérkezik, akkor az igény visszautasításra kerül.

Megint más rendszerek véges kapacitással rendelkeznek, vagyis bizonyos mennyiségű igény felvételére képesek.

A kiszolgáló egységek száma

A legegyszerűbb várakozósoros rendszer egy kiszolgáló egységgel rendelkezik, egyidőben egy igényt szolgál ki.

A több kiszolgáló egységű rendszerek C számú azonos kiszolgáló egységgel rendelkeznek és egyidőben C számú igényt elégítenek ki.

Végtelen számú kiszolgálási egységű rendszerben minden igény azonnal kiszolgálásra kerül.

A sorbanállási rend

A sorbanállási rend, vagy kiszolgálási rend az a szabály, amely szerint kiválasztásra kerül a kiszolgálásban sorrakerülő következő igény.

A leggyakrabban alkalmazott ilyen szabály az „első be - első ki”, a *FIFO* (first-in, first-out), vagy másnéven *FCFS* (first-come, first-served). Ez nem más, mint a mindennap emlegetett sorszám.

Egy másik másik gyakori rend az „utolsó be - első ki”, a *LIFO* (last-in, first-out), vagy más néven *LCFS* (last-come, first-served). Ilyen modell például a közismert veremautomata.

Megemlíthető még a „véletlen kiválasztás”, az *RSS* (Random Selection for Service, más néven *SIRO* (Service In Random Order) és a „fontossági sorrend”, a *PRI* (priority service).

A rövidítések kombinálásával további modellek is létrehozhatók, amelyeknek a mi vizsgálataink szempontjából kevésbé van jelentősége.

Mellékesen megemlítjük a fogyasztói társadalomban kialakult „közgazdasági modellt, a *GIGO*-t (garbage-in, garbage-out)”. Ennek tárgyalása nem tárgya jelen tanulmányunknak.

3. A modell leírása

A továbbiakban az eddig bemutatott alapfogalmak szerinti sorbanállási, vagy várakozósoros rendszerek gyakorlati használatra is megfelelő modelljét mutatjuk be.

A rendszerek rövidített jelölése

A sorbanállásos rendszereket az u.n. *Kendall jelöléssel* szokás rövidített formában leírni az alábbi módon:

$A/B/c/K/m/Z$, ahol

- A a beérkezések közötti idő eloszlása,
- B a kiszolgálási idő eloszlása,
- C a kiszolgáló egységek száma,
- K a rendszer kapacitása (a tárolható igények száma),
- M a források száma,
- Z a sorbanállási rend.

Gyakran csak az $A/B/c$ jelölést használjuk, ilyenkor a várakozósornak nincs felső korlátja, a források száma végtelen és a várakozási rend *FIFO*.

Az A és B helyén a következő eloszlásokat használjuk a leggyakrabban:

- D konstans eloszlás (kiszolgálásra, beérkezésre),
- $E-k$ Erlang- k eloszlás ((kiszolgálásra, beérkezésre),
- G_i általános, független beérkezések közötti idő,
- G általános kiszolgálási idő a függetlenség feltételezésével,
- M exponenciális eloszlás (kiszolgálásra, beérkezésre).

Például az $M/E4/3/20/8/SIRO$ rendszer exponenciális eloszlású beérkezések közötti idővel, három kiszolgáló

egységgel, azonos Erlang-k kiszolgálási idő eloszlással, 20 igény (ebből 3 kiszolgálás alatt, 17 a várakozó sorban) kapacitással, végtelen számú igényforrással és véletlen kiválasztásos (minden várakozó igény u.a. valószínűséggel kerül sorra) kiszolgálási renddel jellemezhető.

A rendszer forgalma

A rendszer forgalma a forgalom intenzitásával u -jellemezhető, ami a kiszolgálás várható értékének $E(s)$ - és a beérkezések közötti idő várható értékének $E(t)$ - az aránya. Ez az arány a sorbanállásos rendszer egyik legfontosabb paramétere és az alábbi formulával írható le:

$$u = \frac{E(s)}{E(t)} = \lambda E(s) = \frac{\lambda}{\mu}$$

A forgalom intenzitása (u) határozza meg a kiszolgálási egységek minimális számát, amely szükséges a rendszer egyensúlyához a beérkező igények számát tekintve.

Ha például $E(t) = 10$ sec és $E(s) = 15$ sec, akkor legalább két kiszolgáló szükséges, mivel: $\frac{E(s)}{E(t)} = \frac{15}{10} = 1,5 \rightarrow 2$.

Az u egységét *erlang*nak nevezték el A.K. Erlang után.

A kiszolgálási egység foglaltsága

Egy másik fontos paraméter a kiszolgálási egységenkénti intenzitás: r , ami u/c , ha a forgalom egyenletesen oszlik meg az egységek között.

A kiszolgálási egységek foglaltsága annak a valószínűsége, hogy az adott egység foglalt. A nagy számok törvénye értelmében a r annak az idonek a megközelítő része, amikor minden egység foglalt.

Egy valószínűség: t időben n igény van a rendszerben

Ez a valószínűség $p_n(t)$ nemcsak t -tol függ, hanem a rendszer kezdeti állapotától is, vagyis a kiszolgálás indulásakor már fennálló igények számától.

A legtöbb működő rendszer esetében t növekedésével $p_n(t)$ közelíti és eléri a p_n értéket, ami független t -tol és így az iniciális állapottól. Az ilyen rendszert kiegyensúlyozott rendszernek nevezzük.

Ebben a tanulmányban csak kiegyensúlyozott rendszerekkel foglalkozunk, mivel az időfüggő és tranzien্স megoldások általában túl bonyolultak a gyakorlati használatra. De a legfontosabb, hogy a legtöbb esetben -- amelyekről itt szót ejtünk -- kiegyensúlyozott rendszer kialakítása a célunk. A nem kiegyensúlyozott, vagyis az időfüggő, illetve tranzien্স rendszerek tervezésével, illetve az ilyen jelenségek felismerésére, elkerülésére irányuló módszerekkel például az Adaptív irányítások (a [6., 9.] irodalom) foglalkoznak.

További jellemzők

A sorbanállási elmélet a rendszer teljesítményének statisztikus mérését teszi lehetővé és így segíti a rendszertechnikust, hogy az igényelt kiszolgálási szinthez szükséges minimális ráfordítású rendszert tervezhessen, alakíthasson ki.

Ezek a statisztikus mértékek és változataik többek között az alábbiak:

- a várakozósorban eltöltött idő várható értéke: Wq ,
- a rendszerben eltöltött idő várható értéke: W ,
- a sorban álló igények várható értéke: Lq ,
- a rendszerben lévő igények várható értéke: L .

Ezek a mértékek nem függetlenek egymástól, ezért valamelyik ismeretében a többi már kiszámítható. A kapcsolatok az alábbi formulákban fejezhetők ki:

$$1. \quad u = \frac{E(s)}{E(t)} = \mathbf{l}E(s) = \frac{\mathbf{l}}{\mathbf{m}}$$

$$2. \quad \mathbf{r} = \frac{u}{c}$$

$$3. \quad W = q + s$$

$$4. \quad W = E(W) = E(q) + E(s) = W_q + W_s$$

$$5. \quad N(t) = N_q(t) + N_s(t)$$

$$6. \quad N = N_q + N_s$$

$$7. \quad L = E(N) = IW = E(N_q) + E(N_s)$$

$$8. \quad L_q = E(N_q) = IW_q$$

Ha például W_q -t ismerjük, akkor a többi már számolható:

$$L_q = IW_q, \quad W = W_q + W_s, \quad L = IW.$$

További hasznos teljesítési mérték a rendszer válaszadási idejének értéke, illetve statisztikusan az u.n. 90%-hoz tartozó érték, a $\mathbf{p}_w(90)$. Ez az az idő, amelynél nem többet tölt el a rendszerben a beérkező igények 90%-a. Formálisan a $p(w \leq \mathbf{p}_w(90)) = 0,9$ egyenlettel fejezhető ki. A 90%-os időérték a várakozósorra is definiálható ugyanígy, jele $\mathbf{p}_q(90)$.

4. Modellek egy kiszolgáló egységgel.

Ebben a fejezetben néhány gyakrabban előforduló modell leírását és formuláit --a levezetések mellozésével-- adjuk meg gyakorlati használatra.

Az M/M/1 modell

A leggyakrabban alkalmazott modell az M/M/1 típus, az egyszerű forma és a fontosabb változók eloszlásának pontos meghatározhatósága miatt.

A modell exponenciális eloszlású beérkezések közötti- és kiszolgálási idővel rendelkezik.

Jelentősen különbözik ez a modell azoktól, amelyeknél --a legtöbb modell ilyen-- a fontosabb változók átlagos, vagy várható értéke, esetleg még a szórása számítható „csak” ki. Egy másik előnye, hogy gyakran célszerű véletlenszerű beérkezésekkel számolni, ugyanakkor a véletlenszerű kiszolgálás feltételezése a gyakorlati szokásokkal találkozódik.

A CPU típusú kiszolgálási idő eloszlásoknál a szórás sokkal nagyobb lehet a várható értéknél és a modell esetenként túl optimisztikus következtetésekre vezethet.

Az alábbiakban az M/M/1 modell kiegyensúlyozott rendszerekre vonatkozó formuláit adjuk meg a levezetések mellozésével:

-a rendszer foglaltsága

$$\mathbf{r} = \rho = \frac{E(s)}{E(t)} = \lambda E(s) = \frac{\lambda}{\mu} \quad \text{mert } c = 1$$

-annak valószínűsége, hogy a rendszerben n darab igény van

$$P_n = P(N = n) = (1 - \mathbf{r}) \mathbf{r}^n \quad n = 0, 1, 2, \dots$$

-annak valószínűsége, hogy a rendszerben n -nél több igény van

$$P(N \geq n) = \mathbf{r}^n \quad n = 1, 2, \dots$$

-a rendszerben lévő igények várható száma

$$L = E(N) = \frac{\mathbf{r}}{1 - \mathbf{r}}$$

-a rendszerben lévo igények szórásnégyzete

$$\mathbf{s}_N^2 = \frac{\mathbf{r}}{(1 - \mathbf{r}^2)}$$

-a rendszerben tartózkodási ido eloszlásfüggvénye (annak a valószínűsége, hogy az igény legfeljebb t idot tolt el a rendszerben)

$$W(t) = P(w \leq t) = 1 - e^{-m(1-\mathbf{r})t} = 1 - e^{-t/E(w)}$$

-a rendszerben eltöltött ido várható értéke

$$W = E(w) = \frac{E(s)}{1 - \mathbf{r}} = \frac{1}{m(1 - \mathbf{r})}$$

-a rendszerben eltöltött ido szórásnégyzete

$$\mathbf{s}_w^2 = E(w)^2$$

-a várakozósor hosszának várható értéke

$$\begin{aligned} L_q = E(N_q) &= \frac{\mathbf{r}^2}{1 - \mathbf{r}} \quad \text{ha } N_q \geq 0 \\ &= \frac{1}{1 - \mathbf{r}} \quad \text{ha } N_q > 0 \quad (\text{nem üres sor}) \end{aligned}$$

-a várakozósor hosszának szórásnégyzete

$$\begin{aligned} \mathbf{s}_{N_q}^2 &= \frac{\mathbf{r}^2(1 + \mathbf{r} - \mathbf{r}^2)}{(1 - \mathbf{r})^2} \quad \text{ha } N_q \geq 0 \\ &= \frac{\mathbf{r}}{(1 - \mathbf{r})^2} \quad \text{ha } N_q > 0 \quad (\text{nem üres sor}) \end{aligned}$$

-a várakozósorban eltöltött idő eloszlásfüggvénye (annak a valószínűsége, hogy az igény legfeljebb t időt tölt el a várakozósorban)

$$W_q(t) = P(q \leq t) = 1 - r e^{-m(1-r)t} = 1 - r e^{-t/E(w)}$$

-a várakozósorban eltöltött idő várható értéke

$$\begin{aligned} W_q = E(q) &= \frac{rE(s)}{1-r} = \frac{r}{m(1-r)} \quad \text{ha } q \geq 0 \\ &= \frac{E(s)}{1-r} = \frac{1}{m(1-r)} \quad \text{ha } q = 0 \quad (\text{nem üres sor}) \end{aligned}$$

-a várakozósorban eltöltött idő szórásnégyzete

$$\begin{aligned} s_q^2 &= \frac{r(2-r)(E(s))^2}{(1-r)^2} = \frac{r(2-r)}{m^2(1-r)^2} \quad \text{ha } N_q \geq 0 \\ &= \left(\frac{E(s)}{1-r} \right)^2 = \left(\frac{1}{m(1-r)} \right)^2 \quad \text{ha } N_q > 0 \quad (\text{nem üres sor}) \end{aligned}$$

-az r százalékos idő a rendszerben

$$p_w(r) = \frac{E(s)}{1-r} \log\left(\frac{100}{100-r}\right) = E(w) \log\left(\frac{100}{100-r}\right)$$

-a 90%-os idő

$$p_w(90) = 2,3E(w)$$

-a 95%-os idő

$$p_w(95) = 3E(w)$$

-az r százalékos idő a várakozósorban

$$p_q(r) = E(w) \log\left(\frac{100r}{100-r}\right) = \frac{E(q)}{r} \log\left(\frac{100r}{100-r}\right)$$

-a 90%-os ido

$$p_q(90) = E(w) \log(10r)$$

-a 95%-os ido

$$p_q(95) = E(w) \log(20r)$$

A valószínűségi változók többsége az M/M/1 modellben közismert formájú:

- az igények száma a rendszerben (N) geometrikus eloszlású,
- a rendszer válaszadási ideje (w) exponenciális eloszlású,
- a várakozósorban eltöltött ido (q) vegyes eloszlású
- (diszkrét belépéskor: $P(q=0)=1-r$, folytonos később).

Kiegyensúlyozott rendszerben a várakozósorban eltöltött ido eloszlása az alábbi: $W_q(t) = 1 - r^{-t/E(w)}$ ha $t \geq 0$

A modell formuláiból kitunik néhány változó erosen nem-lineáris függése a rendszer foglaltságától (ρ). Ilyen a rendszer különböző részeiben eltöltött ido. A non-linearitás azt mutatja, hogy magas rendszerfoglaltság esetén a várakozási ido nagyon megnövekszik és minden határon túl no $r \approx 1$ esetén. Mintegy $r=0,8$ értéknél a foglaltság egészen kis növekedése dramatikusan lerontja a rendszer átbocsájto képességét.

Ennek elkerülésére az igények prioritásának bevezetésére van szükség. Egy alkalmasan megtervezett prioritásos rendszer magas foglaltság esetén is jól muködik. Erre a későbbiekben még visszatérünk.

Az M/M/1 modell alkalmas az u.n. léptékhatás bemutatására: Ha az átlagos beérkezési ido aránya I és az átlagos kiszolgálási ido arány m megkétszerezodnek (változtatlan $r=I/m$ mellett !), az eltöltött ido várható értéke úgy a várakozósorban $E(q)$, mint a rendszerben $E(W)$ egyaránt megfelezoedik. Az igények számának várható értéke mindkét helyen változatlan marad.

Általában ha az új rendszerben I helyett nI és m helyett nm áll, akkor a lépték-viszony az alábbi lesz:

$$\frac{E(q)_{új}}{E(q)_{rég}} = \frac{\left(\frac{r}{(1-r)nm}\right)}{\left(\frac{r}{(1-r)m}\right)} = \frac{1}{n}$$

és

$$\frac{E(w)_{új}}{E(w)_{rég}} = \frac{\left(\frac{1}{(1-r)nm}\right)}{\left(\frac{1}{(1-r)m}\right)} = \frac{1}{n}$$

A léptékhatás egy következtetésre ad lehetőséget. Ha egy nagyszámítógép terhelését egyenlően elosztjuk n darab olyan kisebb számítógép között, amelyek sebessége a nagy gép $1/n$ -szerese, akkor a válaszadási idő nem változik és a felhasználók alkalmasabban elhelyezhető és üzembiztosabb rendszerhez jutnak.

Errol részletesebben a [12., 14.] irodalomban olvashatnak.

Az M/M/1/K modell

($K \geq I$ és az igények száma $n \neq K$)

Ez a modell az M/M/1 modelltől annyiban különbözik, hogy a rendszer igényfelvételi kapacitása véges és legfeljebb K lehet ($n \neq K$).

Ilyen rendszereknél az igények visszautasításra kerülnek, ameddig $n=K$ -val, vagyis a rendszer teljes igénytároló kapacitása foglalt.

A modell az alábbi formulákkal jellemezhető a levezetések mellozésével:

-a rendszer foglaltsága

$$r = (1 - p_k)u$$

-a valószínűség, hogy a rendszerben n darab igény van

$$P_n = \frac{(1 - \mathbf{m})u^n}{1 - u^{K+1}} \quad \text{ha } \mathbf{l} \neq \mathbf{m} \quad \text{és} \quad n = 0, 1, 2, \dots, K$$

-a valószínűség, hogy a beérkező igény visszautasításra kerül

$$P_k = P(n = K)$$

-a tényleges beérkezési arány

$$I_a = (1 - p_k)I$$

-a rendszerben lévő igények várható száma

$$L = E(n) = \frac{u(1 - (K + 1)u^K + Ku^{K+1})}{(1 - u)(1 - u^{K+1})} \quad \text{ha } \mathbf{l} \neq \mathbf{m}$$
$$= \frac{K}{2} \quad \text{ha } \mathbf{l} = \mathbf{m}$$

-a rendszerben eltöltött idő várható értéke

$$W = E(w) = \frac{L}{I_a} = \frac{E(n)}{(1 - p_k)I}$$

-a várakozósor hosszának várható értéke

$$L_q = E(N_q) = L - (1 - p_0)$$

-a várakozósorban eltöltött idő várható értéke

$$W_q = E(q) = \frac{L_q}{I_a} \quad \text{ha } q \geq 0$$
$$= \frac{W_q}{1 - p_0} \quad \text{ha } q > 0$$

Az M/M/1/K modell bizonyos szempontból jobb, mint az M/M/1 modell, amint az a formulákból is látszik a rendszer stabil marad még akkor is, ha a beérkezési arány \mathbf{l} meghaladja a kiszolgálási arányt \mathbf{m} mivel a rendszer telítettsége esetén a beérkező igények visszautasításra kerülnek.

Az M/G/1 modell

Ennél a modellnél általában nem nyerhetők az N , N_q , W , q eloszlási függvényei úgy, mint az M/M/1 esetében és így csak a várható értékek kiszámítására szoritkozhatunk.

Ha a kiszolgálási idő várható értékének első három momentuma -- $E(s)$, $E(s^2)$, $E(s^3)$ -- ismert, akkor több valószínűségi változó várható értéke és szórása számítható. A modell az alábbi formulákkal jellemezhető:

-a várakozósor hosszának várható értéke

$$L_q = E(N_q) = \frac{\mathbf{I}^2 E(s^2)}{2(1-\mathbf{r})} = \frac{\mathbf{I}^2 \mathbf{s}_s^2 + \mathbf{r}^2}{2(1-\mathbf{r})}$$

-a várakozósorban eltöltött idő várható értéke

$$\begin{aligned} W_q &= E(q) = \frac{L_q}{\mathbf{I}} && \text{ha } q \geq 0 \\ &= E(q | q > 0) = \frac{W_q}{\mathbf{r}} && \text{ha } q > 0 \\ &= E(q^2) = \frac{\mathbf{I}E(s^3)}{3(1-\mathbf{r})} + \frac{1}{2} \left(\frac{\mathbf{I}E(s^2)}{1-\mathbf{r}} \right)^2 \end{aligned}$$

a várakozósorban eltöltött idő szórása

$$\mathbf{s}_q^2 = E(q^2) - W_q^2$$

-a rendszerben lévő igények számának várható értéke

$$L = E(N) = L_q + \mathbf{r}$$

-a rendszerben eltöltött idő várható értéke

$$W = E(w) = \frac{L}{I}$$

$$E(w^2) = E(q^2) + \frac{E(s^2)}{1-r}$$

-a rendszerben eltöltött idő szórása

$$\mathbf{s}_w^2 = E(w^2) - W^2$$

-a rendszerben lévő igények számának a szórása

$$\mathbf{s}_N^2 = \frac{I^3 E(s^3)}{3(1-r)} + \left(\frac{I^2 E(s^2)}{2(1-r)} \right)^2 + \frac{I^2 (3-2r) E(s^2)}{2(1-r)} + r(1-r)$$

Az $E(q^2)$, \mathbf{s}_q^2 , $E(w^2)$, \mathbf{s}_w^2 , \mathbf{s}_N^2 csak akkor számítható, ha a kiszolgálási idő első három momentuma szerinti értékek ismertek.

A kiegyensúlyozott rendszerre vonatkozó L és W értékek egyenletei az irodalomban Pollacsek-Khintchine egyenletekként ismertek.

Adott $E(s)$ esetén az L , Lq , W és Wq értékek akkor minimálisak, ha a kiszolgálási idő szórása nulla. (\mathbf{s}_s^2), azaz a kiszolgálási idő *konstans*. Ilyenkor a momentumok kiszámíthatók és ez az M/D/1 modellhez vezet

Az M/D/1 modell

Az M/D/1 modell az M/G/1 speciális esetének tekinthető, amikor

- a kiszolgálási idő: $E(s) = konstans$,
- a momentumok: $E(s^2) = E(s)^2$ illetve $E(s^3) = E(s)^3$,
- a szórás természetesen nulla ($\mathbf{s}_s^2 = 0$).

A modell formulái az alábbiakra egyszerűsödnek:

-a várakozó sor hosszának várható értéke

$$L_q = E(N_q) = \frac{r^2}{2(1-r)}$$

-a várakozósorban eltöltött idő várható értéke

$$W_q = E(q) = \frac{L_q}{I} = \frac{IE(s)}{2(1-r)} \quad \text{ha } q \geq 0$$
$$= E(q | q > 0) = \frac{W_q}{r} = \frac{E(s)}{2(1-r)} \quad \text{ha } q > 0$$

a várakozósorban eltöltött idő szórása

$$s_q^2 = \frac{rE(s)^2}{3(1-r)} + \frac{r^2E(s)^2}{4(1-r)^2}$$

-a rendszerben lévő igények számának várható értéke

$$L = E(N) = L_q + r = \frac{r^2}{2(1-r)} + r$$

-a rendszerben eltöltött idő várható értéke

$$W = E(N) = \frac{L}{I} = W_q + E(s) = \frac{rE(s)}{2(1-r)} + E(s)$$

-a rendszerben eltöltött idő szórása

$$s_w^2 = s_q^2$$

-a rendszerben lévő igények számának a szórása

$$s_N^2 = \frac{r^4}{4(1-r)^2} + \frac{r^3}{3(1-r)} + \frac{r^2(3-2r)}{2(1-r)} + r(1-r)$$

Az M/Ek/1 modell

Az M/Ek/1 modell egy fontos modell, mert számos esetben a kiszolgálási idő jól közelíthető az Erlang eloszlással. További elonye az Erlang eloszlás alkalmazásának, hogy a

momentumok ismertek és így a q , w , N változók várható értékei és a szórás számíthatók. Ezek ismeretében a 90%-os és 95%-os értékek is becsülhetők. Ez a modell végülis az M/G/1 modell speciális eseteként is tekinthető. A kiszolgálási idő eloszlási függvényének a momentumai az alábbi módon számíthatók ki:

$$E(s^2) = \frac{(k+1)}{k} (E(s))^2$$

$$E(s^3) = \frac{(k+1)(k+2)}{k^2} (E(s))^3$$

Ezek behelyettesítésével az M/G/1 modell formuláiba az $E(q^2)$, L_q , S_q^2 értékek már kiszámolhatóak.

A q , w , N százalékos valószínűségei becslésére gamma-disztribúciós táblázat alkalmazását ajánljuk, mivel az Erlang- k a gamma-disztribúció speciális esete, ahol a k paraméter diszkrét egész értékű. A táblázatok az u.n. négyzetes variációs

együttható: $C_x^2 = \frac{\text{var}(x)}{E(x^2)}$ különböző értékeihez adják meg a p

t.

Ha x Erlang- k eloszlású, akkor $C_x^2 = 1/k$.

5. Elsobbségi, prioritásos modellek

Az elsobbségi, vagy prioritásos modellek lényege, hogy az igényeket, vagy forrásaikat számozott osztályokba sorolják 1-től n -ig. Az alacsonyabb sorszámúnak elsobbsége van a magasabb számúval szemben, azaz az i osztályba tartozónak elsobbsége van a j osztálybeli előtt, ha $i < j$, továbbá az 1 osztályba tartozó az összes előtt van. Az azonos osztályba tartozók kiszolgálása FIFO, beérkezési sorrend szerint történik.

Az i igény a beérkezésekor a már várakozó i igények mögé kerül, ha ilyen nincs, akkor a j igény elé kerül ahol $i < j$.

A j kiszolgálásának befolyásolásától függoen kétféle taktika lehetséges:

Kivárasos rendszer, non-preemptive

vagy HOL (head of line)

i csak akkor kap kiszolgálást, amikor j kiszolgálása befejeződött.

Megszakításos rendszer, preemptive

j kiszolgálása azonnal megszakad és i kiszolgálása megkezdodik.

Ezután a visszatérés a j kiszolgálásához kétféle módon lehetséges:

- megszakítás újraindítással, preemptive-repeat:
a j kiszolgálása teljes egészében újraindul,

- megszakítás folytatással, preemptive-resume:
a j kiszolgálása a megszakítástól folytatódik.

Az M/G/1 prioritásos modell fobb összefüggései

Általános összefüggések:

- az átlagos beérkezési arány

$$I = I_1 + I_2 + \dots + I_n$$

- a kiszolgálásban eltöltött ido várható értéke

$$E(s) = \frac{I_1}{I} E(s_1) + \frac{I_2}{I} E(s_2) + \dots + \frac{I_n}{I} E(s_n)$$

$$E(s^2) = \frac{I_1}{I} E(s_1^2) + \frac{I_2}{I} E(s_2^2) + \dots + \frac{I_n}{I} E(s_n^2)$$

- a rendszer forgalma

$$u_j = I_1 E(s_1) + I_2 E(s_2) + \dots + I_j E(s_j) \quad \text{ahol } j = 1, 2, \dots, n$$

$$u_n = u = I E(s)$$

Nonpreemptive esetben:

- a várakozósorban eltöltött ido várható értéke

$$E(q_j) = \frac{I E(s^2)}{2(1-u_{j-1})(1-u_j)} \quad \text{ahol } j = 1, 2, \dots, n \text{ és } u_0 = 0$$

$$E[q] = \frac{I_1}{I} E[q_1] + \frac{I_2}{I} E[q_2] + \dots + \frac{I_n}{I} E[q_n]$$

- a rendszerben eltöltött ido várható értéke

$$E[w_j] = E[q_j] + E[s_j] \quad \text{ahol } j = 1, 2, \dots, n$$

$$E[w] = E[q] + E[s]$$

- a várakozósor hosszának várható értéke

$$L_q = E[N_q] = I E[q]$$

$$L = E[N] = I E[w]$$

Preemptive -resume esetben:

- a várakozósorban eltöltött ido várható értéke

$$E[q_j] = E[w_j] - E[s_j] \quad \text{ahol } j = 1, 2, \dots, n$$

$$w_q = E[q] = \frac{I_1}{I} E[q_1] + \frac{I_2}{I} E[q_2] + \dots + \frac{I_n}{I} E[q_n]$$

- a rendszerben eltöltött ido várható értéke

$$E[w_j] = \frac{1}{1-u_{j-1}} \left[E[s_j] + \frac{\sum_{i=1}^j I_i E[s_i^2]}{2(1-u_j)} \right] \quad \text{ahol } j = 1, 2, \dots, n; u_0 = 0$$

$$W = E[w] = E[q] + E[s]$$

- a várakozósor hosszának várható értéke

$$L_{q_j} = E[N_{q_j}] = I_j E[q_j] + E[s_j] \quad \text{ahol } j = 1, 2, \dots, n$$

$$L_q = E[N_q] = I E[q] = I w_q$$

$$L = E[N] = I E[w] = I w$$

6. Összefoglalás

Az itt összefoglalt technika a számítóközpontok üzemének tervezéséhez használatos módszerek egyik készlete. Segítségével tervezési és döntési értékű válasz nyerhető --a számítógépekkel kapcsolatos elvárások alapján-- a rendszer konfigurációs adatainak meghatározására.

Ez az összeállítás megjelent a KSH Rendszertechnikai Közlemények sorozatában. Az ott ismertetett módon került sor ennek a technikának alkalmazására KSH számítóközpontja legújabb gépének konfigurálására és az Üzemeltetési Kézikönyv és -Szabályzat elkészítésére.

A szerző köszönetet mond munkatársainak a tanulmány elkészítéséhez nyújtott segítségükért, értékes javaslataikért.
Budapest, 1975.

7. Mellékletek

A. Az alkalmazott jelölések és fogalmak

- $A(t)$ -a beérkezések között eltelt idő eloszlásfüggvénye,
vagyis $A(t) = P(t \leq t)$
- C -az azonos kiszolgálók száma
- D -állandó eloszlás (beérkezés, vagy kiszolgálás)
- E_k -Erlang-k eloszlás (beérkezés, vagy kiszolgálás)
- $E(N_q | N_q > 0)$ -a nem-üres várakozósor hosszának várható értéke, vagy átlaga
- $E(q | q > 0)$ -a nem-üres várakozósorban eltöltött idő várható értéke, vagy átlaga
- $FCFS$ -first-come, first-served, lásd még FIFO
- $FIFO$ -first-in, first-out, beérkezési sorrend szerint
- G -általános eloszlás (kiszolgálás)
- GI -általános, független eloszlás (beérkezés, vagy kiszolgálás)
- $GIGO$ -garbage-in, garbage-out
- K -a rendszer kapacitása: legfeljebb ennyi igény lehet a várakozósorokban és a kiszolgálókban együttesen
- L - $E(N)$, a rendszerben lévő igények számának várható értéke
- L_q - $E(N_q)$, a várakozósorban lévő igények várható száma, vagy átlaga

- LCFS* -last-come, first-served, lásd még LIFO
- LIFO* -last-in, first-out, az utoljára beérkezett következik
(verem automata)
- I* -az átlagos beérkezési idő
- M* -exponenciális eloszlás (beérkezés, vagy kiszolgálás)
- m* -az átlagos kiszolgálási idő
- N(t)* -valószínűségi változó, a rendszerben lévő igények várható száma t időpontban
- N* -valószínűségi változó, a rendszerben lévő igények számának várható értéke
- Nq(t)* -valószínűségi változó, a várakozósorban lévő igények várható száma t időpontban
- Nq* -valószínűségi változó, a várakozósorban lévő igények számának várható értéke
- Ns(t)* -valószínűségi változó, a kiszolgálásban lévő igények várható száma t időpontban
- Ns* -valószínűségi változó, a kiszolgálásban lévő igények számának várható értéke
- Pn(t)* -annak a valószínűsége, hogy t időpontban n igény van a rendszerben
- Pn* -annak a valószínűsége, hogy n igény van a rendszerben
- PRI* -elsobbbségi kiszolgálási szabály
- q* -valószínűségi változó, az igénynek a várakozósorban eltöltött ideje a kiszolgálás megkezdése előtt
- RSS* -Random Selection for Service, véletlenszerű kiszolgálás lásd még SIRO

- \mathbf{r} -a kiszolgáló kihasználtsága: $\mathbf{r} = \frac{\mathbf{l}}{c\mathbf{m}}$
- s -valószínűségi változó, a kiszolgálási idő leírására
- SIRO* -Service In Random Order, minden igény ugyanolyan valószínűséggel kerül kiszolgálásra, lásd még RSS
- \mathbf{t} -valószínűségi változó, a beérkezések között eltelt idő leírására
- u -rendszer forgalom: $u = \frac{E(s)}{E(t)} = \mathbf{l}E(s) = \frac{\mathbf{l}}{\mathbf{m}}$
- w -valószínűségi változó, a rendszerben eltöltött idő leírására
- $W(t)$ -a w (a rendszerben eltöltött idő) eloszlásfüggvénye:
 $W(t) = P(w \leq t)$
- W - $E(w)$, az igénynek a rendszerben eltöltött ideje
- $Wq(t)$ -a w_q (a várakozósorban eltöltött idő) eloszlásfüggvénye:
- Wq - $E(q)$, a várakozósorban eltöltött idő várható értéke
- $Ws(t)$ -a w_s (a kiszolgálásban eltöltött idő) eloszlásfüggvénye:
- Ws - $E(s)$, a kiszolgálásban eltöltött idő várható értéke.

B. Irodalomjegyzék

1. G.A. Korn: Matematikai kézikönyv muszakiaknak. Muszaki, 1975.
2. Prékopa A.: Valószínűségelmélet, Muszaki, 1974.
3. P. Denning: The working set model of program behaviour. Comm. ACM., 1968.
4. E.G. Coffman: Analysis of two time-sharing algorithms designed for limited swapping. Journal of ACM., 1968.
5. L. Kleinrock: A continuum of time-sharing algorithms. Proc of AFIPS., 1970.
6. Gyarmati P.: Dinamikus erőforrás elosztás vegyes üzem esetén. SZTAKI Közl., 1975
7. A. Brandwajn: A model of time-sharing virtual memory system using equivalence and decomposition method. Acta info., 1974.
8. T. Beretvas: A simulation model representing OS/VS2 CP. Proc of OS., 1974.
9. P. Gyarmati: On the adaptiv control of Operating Systems., 1975.
10. KSH Rendszertechnikai Közlemények /3. Szerk.: Gyarmati Péter., 1974.
11. G.S. Shedler: A queuing model of multiprogrammed computer with a two level storage system. Comm. ACM., 1973.
12. L. Kleinrock: Queuing systems vol.1 theory., 1975.
13. W. Chang: Single-server queuing processes in computer systems. IBM Sytems Journal., 1970.
14. D. Gross and C. M. Harris: Fundamentals of Queuing Theory. -The McMillan Co., 1974.

