



ΠΑΝΕΠΙΣΤΗΜΙΟ
ΠΑΤΡΩΝ
UNIVERSITY OF PATRAS

Ανάλυση Απόδοσης Πληροφοριακών Συστημάτων

Διαλέξεις 6-9: Λειτουργική Ανάλυση

Καθ. Γιάννης Γαροφαλάκης

ΜΔΕ Επιστήμης και Τεχνολογίας Υπολογιστών

Τμήμα Μηχανικών Η/Υ & Πληροφορικής

2016-2017

Εισαγωγή στη Λειτουργική Ανάλυση I

- **Λειτουργικοί νόμοι** : Απλές σχέσεις που δεν απαιτούν κατανομή χρόνων μεταξύ αφίξεων ή εξυπηρέτησης (Buzen 1976, Denning & Buzen 1978)
- **"Λειτουργική"**: απευθείας μετρήσιμο.

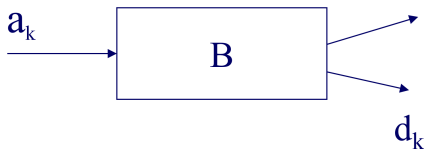
Υποθέσεις :

1. **Job flow balance** = Εργοδικότητα (Μετρήσιμο)
2. Σε συγκεκριμένο χρόνο T , **αριθμός αφίξεων = αριθμός αναχωρήσεων**. (Μετρήσιμο)
3. **Ανεξαρτησία**. Η ακολουθία χρόνων εξυπηρέτησης δεν μπορεί να εντοπιστεί με μετρήσεις αν είναι ανεξάρτητες ή όχι. (Μη Μετρήσιμο)



Εισαγωγή στη Λειτουργική Ανάλυση II

- **Λειτουργικές ποσότητες** : ποσότητες που μπορούν να μετρηθούν κατά τη διάρκεια πεπερασμένης περιόδου παρατήρησης $\underline{\quad}$
Παράδειγμα λειτουργικού νόμου : Νόμος του Little: $\overline{N} = \lambda T$



- a_k : βλέπει \overline{N} στο B.
- d_k : αφήνει \overline{N} στο B.
- Άρα, οι \overline{N} που βλέπει φεύγοντας, ήρθαν όλοι, όσο ο πελάτης βρισκόταν μέσα (T). Κατά την παρουσία του ήρθαν λT κατά μέσο όρο. Άρα, $\overline{N} = \lambda T$.



Λειτουργικές ακολουθίες συμπεριφοράς και ιδιότητες vs. Στοχαστικές διαδικασίες

- Ίδια αποτελέσματα με τις στοχαστικές διαδικασίες αλλά περισσότερο "δαισθητικά" και εφαρμόσιμα σε μη στοχαστικά συστήματα.
- 3 ιδιότητες που ενδιαφέρουν:
 - Ομογενής συμπεριφορά.

Αφίξεις Εξυπηρέτησεις Δρομολόγηση	}	Ανεξάρτητα από τον αριθμό πελατών στο σύστημα. (Εξαίρεση: $\mu = 0$ όταν δεν υπάρχουν ή $\lambda = 0$ όταν έχουμε \bar{N} στο σύστημα)
---	---	--

- Flow – Balance
 - Συμπεριφορά ενός βήματος : Μόνο ένα γεγονός τη φορά.
- **Σημαντικό** : Οι λειτουργικές ιδιότητες ορίζονται για μια συγκεκριμένη ακολουθία συμπεριφοράς :
 - "Μέσο" = στατιστικό μέσο.
 - "Πιθανότητα" = σχετική συχνότητα.



Βασικές λειτουργικές σχέσεις I



- Χρόνος Παρατήρησης T .
- Μπορούμε να μετρήσουμε τις λειτουργικές ποσότητες:
 - A_i : Αριθμός αφίξεων.
 - C_i : Αριθμός Εξυπηρετήσεων (completions).
 - B_i : Busy time κατά το T .



Βασικές λειτουργικές σχέσεις II

- Επιπλέον, μπορούμε να πάρουμε και τις εξής λειτουργικές ποσότητες :
 - Μέσος Ρυθμός Αφίξεων : $\lambda_i = \frac{A_i}{T}$
 - Throughput : $X_i = \frac{C_i}{T}$
 - Utilization : $U_i = \frac{B_i}{T}$
 - Μέσος Χρόνος Εξυπηρέτησης : $S_i = \frac{B_i}{C_i}$
- Οι λειτουργικές ποσότητες είναι μεταβλητές που μπορεί να αλλάζουν από τη μία περίοδο παρατήρησης στην επόμενη. Όσες παραμένουν σταθερές, ονομάζονται *λειτουργικοί νόμοι*.



Νόμος Χρησιμοποίησης

- ▣ Έστω C_i ο αριθμός εξυπηρετήσεων, B_i το busy time της device i κατά μία περίοδο T .
- ▣ Νόμος χρησιμοποίησης :

$$U_i = \frac{B_i}{T} = \frac{C_i B_i}{T C_i} = X_i S_i$$

$$X_i = \frac{C_i}{T} = \frac{C_i B_i}{B_i T} = \frac{1}{S_i} U_i = \frac{U_i}{S_i}$$

$$S_i = \frac{B_i}{C_i} = \frac{B_i T}{T C_i} = U_i \frac{1}{X_i} = \frac{U_i}{X_i}$$



Νόμος Εξαναγκασμένης Ροής I

- Συνδέει το throughput του συστήματος με το throughput ενός device.
- Σε ένα ανοιχτό δίκτυο το system throughput είναι ο αριθμός των jobs που φεύγουν από το σύστημα ανά μονάδα χρόνου.
- Σε ένα κλειστό δίκτυο, καμία job δεν μπορεί να το εγκαταλείψει. Συνεπώς, ορίζεται ένα link ώστε ότι αποχωρεί από το δίκτυο, ταυτόχρονα, να εισέρχεται. Όσες διέρχονται από εκεί ορίζουν το system throughput.



Νόμος Εξαναγκασμένης Ροής II

- Αν το T είναι τέτοιο ώστε $A_i = C_i$, τότε το device ικανοποιεί την υπόθεση Job Flow Balance.
- Έστω ότι κάθε job κάνει V_i αιτήσεις για το i -στο device. Ισχύει:

$$C_i = C_0 V_i \text{ ή } V_i = \frac{C_i}{C_0}$$

- όπου C_0 , ο αριθμός των jobs που διατρέχουν το εξωτερικό link.
- Το V_i ονομάζεται *visit ratio* ή *relative throughput*, (λόγος επισκέψεων το i device προς επισκέψεις στο out link).
- System Throughput : $X = \frac{C_0}{T}$ και

$$X_i = \frac{C_i}{T} = \frac{C_i}{C_0} \frac{C_0}{T} \Rightarrow \boxed{X_i = X V_i}$$



Νόμος Εξαναγκασμένης Ροής III

- Ο Νόμος Εξαναγκασμένης Ροής ισχύει όποτε ισχύει το Job Flow Balance.
- Συνδυάζοντας το Νόμο Εξαναγκασμένης Ροής και το Νόμο Χρησιμοποίησης, προκύπτει:

$$\left. \begin{array}{l} X_i = XV_i \\ U_i = X_i S_i \end{array} \right\} U_i = XV_i S_i \Rightarrow U_i = XD_i$$

όπου, $D_i = V_i S_i$ είναι η συνολική απαίτηση service στο device- i για όλες τις επισκέψεις μίας job.

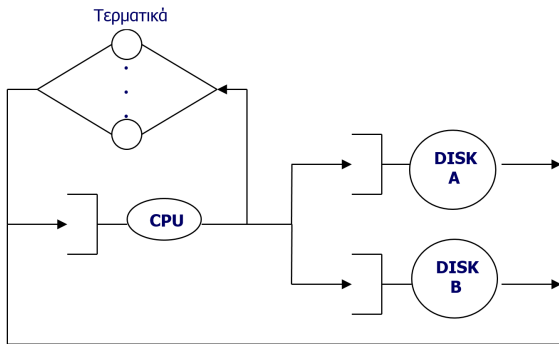
- Δηλαδή, το device με το μεγαλύτερο service demand D_i έχει το μεγαλύτερο Utilization και είναι bottleneck device.



Νόμος Εξαναγκασμένης Ροής IV

Example

- Διαθέτουμε ένα central server model ενός timesharing συστήματος.



Νόμος Εξαναγκασμένης Ροής V

- Μετρώντας τα log data προκύπτουν τα εξής:
 - Κάθε πρόγραμμα απαιτεί 5 sec CPU time και κάνει 80 I/O request στο δίσκο A και 100 στο δίσκο B.
 - Ο μέσος χρόνος σκέψης των χρηστών είναι 18 sec.
 - Ο δίσκος A απαιτεί 50 msec για μία I/O request ενώ ο δίσκος B, 30 msec.
 - Με 17 active terminals το throughput του A μετρήθηκε ίσο με 15.7 I/O requests/sec.
- Θέλουμε να βρούμε το system throughput και το utilization των devices.



Νόμος Εξαναγκασμένης Ροής VI

Δηλαδή: $D_{CPU} = 5 \text{ sec}$, $V_A = 80$, $V_B = 100$, $Z = 18 \text{ sec}$, $S_A = 0.05 \text{ sec}$, $S_B = 0.03 \text{ sec}$, $N = 17$, $X_A = 15.7 \text{ jobs/sec}$.

Επειδή οι jobs πρέπει να επισκεφθούν τη CPU πριν τους δίσκους ή τα τερματικά, το visit ratio της CPU είναι:

$$V_{CPU} = V_A + V_B + 1 = 181$$

Το πρώτο βήμα στην λειτουργική ανάλυση, γενικά, είναι να καθοριστούν τα D_i όλων των devices:

$$\begin{aligned} D_{CPU} &= 5 \text{ sec}, \\ D_A &= S_A V_A = 4 \text{ sec}, \\ D_B &= S_B V_B = 3 \text{ sec} \end{aligned}$$



Νόμος Εξαναγκασμένης Ροής VII

Χρησιμοποιώντας το νόμο εξαναγκασμένης ροής, τα throughputs προκύπτουν:

$$\begin{aligned}X &= \frac{X_A}{V_A} = 0.1963 \text{ jobs/sec,} \\X_{CPU} &= X V_{CPU} = 35.48 \text{ requests/sec,} \\X_B &= X V_B = 19.6 \text{ requests/sec}\end{aligned}$$

Χρησιμοποιώντας το νόμο χρησιμοποίησης, τα utilizations προκύπτουν :

$$\begin{aligned}U_{CPU} &= X D_{CPU} = 98\%, \\U_A &= X D_A = 78.4\%, \\U_B &= X D_B = 58.8\%\end{aligned}$$



Νόμος Εξαναγκασμένης Ροής VIII

- ▣ Άλλος τρόπος περιγραφής της δρομολόγησης jobs σε ένα queuing network είναι οι πιθανότητες μετάβασης (transition probabilities), p_{ij} . Είναι ισοδύναμη περιγραφή με τα V_i .
- ▣ Σε δίκτυο με Job Flow Balance:

$$C_j = \sum_{i=0}^M C_i p_{ij}$$

- ▣ Το 0 επισημαίνει τις επισκέψεις στο εξωτερικό link.



Νόμος Εξανασκασμένης Ροής IX

- Ισχύει,

$$C_j = \sum_{i=0}^M C_i p_{ij} \Rightarrow \frac{C_j}{C_0} = \sum_{i=0}^M \frac{C_i}{C_0} p_{ij} \Rightarrow V_j = \sum_{i=0}^M V_i p_{ij}, \quad V_0 = 1$$

- $V_0 = 1$, διότι κάθε επίσκεψη στο εξωτερικό link ορίζεται σαν την ολοκλήρωση της job.
- Οι παραπάνω εξισώσεις είναι γνωστές σαν *visit ratio equations*. Ισχύουν εφόσον το δίκτυο είναι συνδεδεμένο λειτουργικά : Κάθε συσκευή επισκέπτεται από κάθε job τουλάχιστον 1 φορά.



Νόμος του Little

- Είναι λειτουργικός νόμος.
- Εφαρμόζουμε το Νόμο του Little για να σχετίσουμε το μήκος της ουράς Q_i , με το χρόνο απόκρισης R_i , στη device- i :

$$\begin{array}{l} \text{Μέσος Αριθμός jobs στην } i \\ Q_i \end{array} = \begin{array}{l} \text{ρυθμός αφίξεων} \\ \lambda_i \end{array} \cdot \begin{array}{l} \text{μέσο χρόνο στην device} \\ R_i \end{array}$$

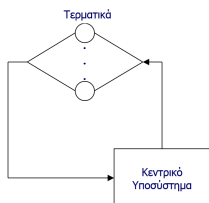
- Αν το job flow είναι σε ισορροπία, ο ρυθμός άφιξης είναι ίσος με το throughput ($\lambda_i = X_i$) και ισχύει

$$Q_i = X_i R_i$$



Γενικός Νόμος Χρόνου Απόκρισης I

Timesharing System



terminal υποσύστημα
κεντρικό υποσύστημα

I user \longrightarrow I terminal.

- Ο Νόμος του Little εφαρμόζεται σε όλα τα τμήματα του συστήματος, αρκεί να υπάρχει job flow balance.
- Σε κεντρικό υποσύστημα, ισχύει :

$$\begin{aligned} \text{Συνολικός Αριθμός jobs στο υποσύστημα} &= \text{Throughput} \cdot \text{Response Time} \\ Q &= X \cdot R \end{aligned}$$



Γενικός Νόμος Χρόνου Απόκρισης II

Για M devices έχουμε

$$Q = Q_1 + Q_2 + \dots + Q_M$$

$$XR = X_1R_1 + X_2R_2 + \dots + X_MR_M$$

$$\frac{XR}{X} = \frac{X_1R_1 + X_2R_2 + \dots + X_MR_M}{X} \xrightarrow{X_i = XV_i}$$

$$R = V_1R_1 + V_2R_2 + \dots + V_MR_M$$

$$R = \sum_{i=0}^M R_i V_i$$

Γενικός Νόμος Χρόνου Απόκρισης:

$$R = \sum_{i=0}^M R_i V_i$$



Γενικός Νόμος Χρόνου Απόκρισης III

- Ισχύει και σε περιπτώσεις που δεν ισχύει το Job Flow Balance.
- Δηλαδή, ο συνολικός χρόνος που καταναλώνει μία job σε ένα server είναι το γινόμενο του χρόνου ανά επίσκεψη, επί τον αριθμό των επισκέψεων, στο server.
- Ο συνολικός χρόνος συστήματος είναι το άθροισμα των συνολικών χρόνων στους διάφορους servers.



Interactive Response Time Law

- Έστω Z ο χρόνος think time στο terminal.
- Συνολικός χρονικός κύκλος αιτήσεων: $R + Z$.
- Κάθε χρήστης δημιουργεί περίπου $\frac{T}{R+Z}$ αιτήσεις στο χρονικό διάστημα T .
- Αν έχουμε N χρήστες (τερματικά):

$$\begin{aligned} \text{System Throughput: } X &= \frac{\text{Συνολικός Αριθμός Αιτήσεων}}{\text{Χρόνος}} = \\ &= \frac{N \frac{T}{R+Z}}{T} = \frac{N}{R+Z} \Rightarrow R = \frac{N}{X} - Z \end{aligned}$$

- **Interactive Response Time Law:**

$$R = \frac{N}{X} - Z$$



Bottleneck analysis I

Νόμος εξαναγκασμένης ροής: $U_i = XD_i \Rightarrow U_i \propto D_i$. Το device με το υψηλότερο φορτίο αιτήσεων για service D_i , έχει το μεγαλύτερο utilization και είναι το bottleneck device. Η βελτίωση του θα βελτιώσει το σύστημα.

Βήμα 1: Προσδιορισμός του device για μελέτη βελτίωσης απόδοσης.

- Έστω ότι βρίσκουμε ότι το device b είναι το bottleneck. Δηλαδή,

$$D_b = D_{max} = \max(D_1, D_2, \dots, D_M)$$

Τότε, το throughput και οι χρόνοι απόκρισης του συστήματος περιορίζονται από τις παρακάτω σχέσεις:

$$X^{(N)} \leq \min\left(\frac{1}{D_{max}}, \frac{N}{D + Z}\right)$$
$$R^{(N)} \geq \max(D, ND_{max} - Z)$$



Το $D = \sum_i D_i$ είναι το άθροισμα των service devices εκτός των terminals.

Bottleneck analysis II

Απόδειξη:

Παρατηρήσεις :

1. Το utilization οποιουδήποτε device δεν μπορεί να είναι μεγαλύτερο από 1. Με αυτόν τον τρόπο τίθεται ένα όριο για το μέγιστο δυνατό throughput. (A)
2. Το response time του συστήματος με N users δεν μπορεί να είναι μικρότερο από αυτό του συστήματος με 1 user. (B)
3. Η σχέση $R = \frac{N}{X} - Z$ μπορεί να χρησιμοποιηθεί για τη μετατροπή του ορίου του throughput σε όριο για το response time και αντίστροφα.



Bottleneck analysis III

Βασισμένοι στις προηγούμενες παρατηρήσεις, έχουμε, για το bottleneck device:

$$\begin{aligned}U_b &= XD_{max} \stackrel{(A)}{\Rightarrow} \\U_b &\leq 1 \Rightarrow \\X &\leq \frac{1}{D_{max}}\end{aligned}$$

Με 1 χρήστη στο σύστημα, δεν υπάρχει αναμονή στην ουρά:

$$R^{(1)} = D_1 + D_2 + \dots + D_M = D \quad (1)$$

Με περισσότερους χρήστες, μπορεί να υπάρξει αναμονή:

$$R^{(N)} \geq D \quad (2)$$



Bottleneck analysis IV

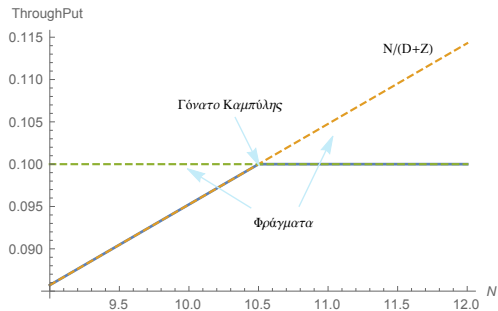
Ισχύει από την (A):

$$R^{(N)} = \frac{N}{X^{(N)}} - Z \geq ND_{max} - Z \xrightarrow{(B)}$$
$$X^{(N)} = \frac{N}{R^{(N)} + Z} \leq \frac{N}{D + Z}$$

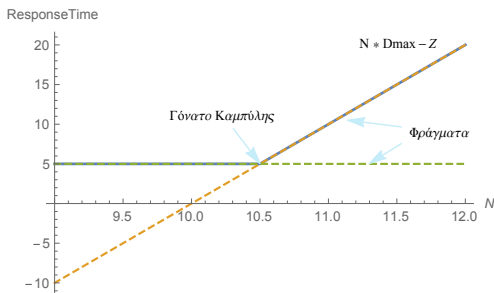
Συνδυάζοντας τα παραπάνω όρια με τις (B), (1) προκύπτει το όριο για το $X^{(N)}$ και με τις (A), (2) προκύπτει το όριο για το $R^{(N)}$.



Bottleneck analysis V




Bottleneck analysis VI



Η προηγούμενες γραφικές αναπαριστούν τα ασυμπτωτικά όρια.

Παρατηρείται, ότι το k_{nue} εμφανίζεται για την ίδια τιμή του αριθμού χρηστών :

$$D = N^* D_{max} - Z \Rightarrow N^* = \frac{D + Z}{D_{max}}$$

 Αν ο αριθμός χρηστών είναι μεγαλύτερος από N^* μπορούμε με βεβαιότητα να πούμε ότι υπάρχει αναμονή στο σύστημα.

Bottleneck analysis VII

Συνέχεια παραδείγματος 1 :

Έχουμε $D_{CPU} = 5, D_A = 4, D_B = 3, Z = 18$. Άρα,

$$D = D_{CPU} + D_A + D_B = 12 \Rightarrow D_{max} = D_{CPU} = 5$$

Τα όρια είναι:

$$X^{(N)} \leq \min \left\{ \frac{N}{30}, \frac{1}{5} \right\}$$

$$R^{(N)} \geq \max \{12, 5N - 18\}$$

Το κπee εμφανίζεται για N^* : $12 = 5N^* - 18 \Rightarrow N^* = 6$. Δηλαδή, αν υπάρχουν περισσότεροι από 6 χρήστες σίγουρα θα υπάρχει αναμονή κάπου. υπάρχει αναμονή κάπου.



Bottleneck analysis VIII

Πόσα τερματικά μπορούν να υποστηριχθούν αν το response time πρέπει να μην ξεπερνά τα 100 sec;

Με τη χρησιμοποίηση των ασυμπτωτικών ορίων, έχουμε:

$$R^{(N)} \geq \max \{12, 5N - 18\}$$

Το response time θα είναι μεγαλύτερο από 100 αν

$$\max \{12, 5N - 18\} \geq 100 \Rightarrow N \geq 23.6$$

Συνεπώς, το σύστημα δεν μπορεί να υποστηρίξει περισσότερους από 23 χρήστες ώστε να είναι $R^{(N)} \leq 100$.



Operational Laws

- Utilization Law : $U_i = X_i S_i = X D_i$
- Forced Flow Law : $X_i = X V_i$
- Little's Law : $Q_i = X_i R_i$
- General Response Time Law : $R = \sum_{i=1}^M R_i V_i$
- Interactive Response Time Law : $R = \frac{N}{X} - Z$
- Asymptotic bounds :

$$R \geq \max \{D, N D_{max} - Z\}$$
$$X \leq \min \left\{ \frac{1}{D_{max}}, \frac{N}{D + Z} \right\}$$



Mean Value Analysis: Ανάλυση ανοικτών δικτύων

- Μοντέλα για *Transaction Processing* συστήματα (ρυθμός αφίξεων ανεξάρτητος από φορτίο συστήματος).
- Αφίξεις *Poisson* με μέσο ρυθμό λ .
- Τα service centers είναι:
 - fixed capacity (single server εκθετικοί) ή
 - delay centers (infinite server εκθετικοί)



MVA Ανάλυση ανοικτών δικτύων: Fixed Capacity Centers

Σε όλα τα **fixed capacity service centers**:

□ Response Time: $R_i = S_i(1 + Q_i) = \underbrace{S_i}_{\text{service}} + \underbrace{S_i Q_i}_{\substack{\text{service των jobs} \\ \text{που βλέπει στην } i}}$

□ Δεν είναι operational law διότι υποθέτει *memoryless service*, κάτι που δεν είναι λειτουργικά δοκιμαζόμενο.

□ Λόγω του job flow balance, ισχύει: $X = \lambda$.

□ Σε κάθε device i :

□ Forced flow law : $X_i = X V_i$.

□ Νόμος χρησιμοποίησης: $U_i = X_i S_i = X V_i S_i = \lambda D_i$.

□ Με τη χρήση του N. Little:

$$\begin{aligned} Q_i &= X_i R_i = X_i S_i (1 + Q_i) = U_i (1 + Q_i) \Rightarrow \\ Q_i &= \frac{U_i}{1 - U_i} \Rightarrow R_i = \frac{S_i}{1 - U_i} \end{aligned} \quad (3)$$



MVA Ανάλυση ανοικτών δικτύων: Delay Centers

Σε όλα τα delay centers:

- Ισχύει,

$$R_i = S_i$$

- Και,

$$Q_i = X_i R_i = X V_i S_i = X D_i = U_i$$



MVA - Ανάλυση Κλειστών Δικτύων I

- Γενικά εφαρμόζεται για πολλά είδη service κατανομών και τρόπων εξυπηρέτησης.
- Θεωρούμε *fixed capacity service centers*.
 - Αν έχουμε κλειστό δίκτυο με N εργασίες, οι Reiser και Lavenberg (1980), έδειξαν:

$$R_i^{(N)} = S_i (1 + Q_i^{(N-1)})$$

- όπου, $Q_i^{(N-1)}$, ο μέσος αριθμός εργασιών στο i -οστό device, όταν είναι $N - 1$ jobs στο δίκτυο.
- Δηλαδή, όταν προσθέσουμε ένα πελάτη στο δίκτυο, ενώ υπάρχουν ήδη $N - 1$ πελάτες, μόλις φθάσει στο i -οστό device, θα βρει $Q_i^{(N-1)}$ jobs μέσα.
- Θα περιμένει χρόνο $S_i Q_i^{(N-1)}$ μέχρι να εξυπηρετηθεί και θα ξοδέψει χρόνο S_i κατά την εξυπηρέτησή του.



MVA - Ανάλυση Κλειστών Δικτύων II

- System Response Time (Γενικός Ν. Χρόνου Απόκρισης):

$$R^{(N)} = \sum_{i=1}^M V_i R_i$$

- System Throughput (Interactive Response Time Law):

$$X^{(N)} = \frac{N}{R^{(N)} + Z}$$

- Device throughputs βάσει των jobs:

$$X_i^{(N)} = X^{(N)} V_i$$

- Μήκη ουρών με τη χρήση του N. Little:

$$Q_i^{(N)} = X_i^{(N)} R_i^{(N)} = X^{(N)} V_i R_i^{(N)}$$



MVA - Ανάλυση Κλειστών Δικτύων III

Θεωρούμε **delay centers**.

- Ισχύει, $R_i^{(N)} = S_i$
- Device throughputs βάσει των jobs:

$$X_i^{(N)} = X^{(N)} V_i$$

- Μήκη ουρών με τη χρήση του N. Little:

$$Q_i^{(N)} = X_i^{(N)} R_i^{(N)} = X^{(N)} V_i R_i^{(N)}$$

- Η MVA ισχύει, αν είναι *Product Form Network*:
 - Job Flow Balance.
 - One-Step behavior.
 - Device homogeneity.
 - Εκθετικοί χρόνοι εξυπηρέτησης.



Αλγόριθμος MVA

Είσοδοι: N, Z, M, S_i, V_i

Έξοδοι: X, Q_i, R_i, R, U_i

for $i = 0$ to M **do** $Q_i \leftarrow 0$

end for

for $n = 0$ to N **do**

for $i = 0$ to M **do**

$$R_i = \begin{cases} S_i(1 + Q_i) \\ S_i \end{cases} \quad \begin{array}{l} \text{fixed capacity} \\ \text{delay centers} \end{array}$$

end for

$$R = \sum_{i=1}^M R_i V_i$$

$$X = \frac{Z+R}{Z+R}$$

for $i = 0$ to M **do** $Q_i \leftarrow X V_i R_i$

end for

end for

□ Αρχικοποίηση

□ Device Throughputs: $X_i = X V_i$

□ Device Utilizations: $U_i = X S_i V_i$



Approximate MVA I

- Η MVA είναι αναδρομικός αλγόριθμος. Για μεγάλο N , εμφανίζει μεγάλο υπολογιστικό κόστος.
- Η προσέγγιση Schweitzer (Schweitzer's Approximation) (1979) αποφεύγει την αναδρομή του MVA.
- Schweitzer's Approximation:
 1. Εκτίμηση του Q_i με N jobs.
 2. Υπολογισμός των R_i, X_i .
 3. Υπολογισμός του Q_i . Αν η διαφορά με το προηγούμενο Q_i είναι μικρή, η νέα τιμή είναι καλή εκτίμηση.



Approximate MVA II

- ▣ **Υπόθεση:** Όσο μεγαλώνει ο αριθμός των jobs στο δίκτυο, το μήκος ουράς κάθε device αυξάνεται ανάλογα. Δηλαδή,

$$\frac{Q_i^{(N)}}{N} = a_i = \text{σταθερό } \forall N. \text{ Πιο συγκεκριμένα,}$$

$$\begin{aligned}\frac{Q_i^{(N-1)}}{N-1} &= \frac{Q_i^{(N)}}{N} \Rightarrow \\ Q_i^{(N-1)} &= \frac{N-1}{N} Q_i^{(N)}\end{aligned}$$

- ▣ Συνεπώς οι εξισώσεις MVA αλλάζουν ως εξής:

$$R_i = \begin{cases} S_i(1 + \frac{N-1}{N} Q_i) & \text{fixed capacity} \\ S_i & \text{delay centers} \end{cases}$$

$$R = \sum_{i=1}^M R_i V_i \quad X = \frac{n}{Z + R}$$



Approximate MVA III

Είσοδοι: $N, Z, M, S_i, V_i, \epsilon$: μέγιστο επιτρεπόμενο λάθος στα $Q_i^{(N)}$

Έξοδοι: X, Q_i, R_i, R, U_i

for $i = 0$ to M **do** $Q_i = N/M$

□ Αρχικοποίηση

end for

while $\max_i \{Q_i - X R_i V_i\} > \epsilon$ **do**

for $i = 0$ to M **do**

$$R_i = \begin{cases} S_i(1 + \frac{N-1}{N}Q_i) & \text{fixed capacity} \\ S_i & \text{delay centers} \end{cases}$$

end for

$$R = \sum_{i=1}^M R_i V_i$$

$$X = \frac{n}{Z+R}$$

for $i = 0$ to M **do** $Q_i = X V_i R_i$

end for

end while



Approximate MVA IV

- Συγκλίνει;
 - Οι αρχικές τιμές των $Q_i^{(N)}$ επηρεάζουν τον αριθμό των iterations και όχι το αποτέλεσμα.
- Αρχική τιμή: $Q_i^{(N)} = \frac{N}{M}$
- Device Throughputs: $X_i = X V_i$
- Device Utilizations: $U_i = X S_i V_i$



Balanced Job Bounds I

- Zahorjan, Sevcik, Eager, Galler (1982).
- Άνω και κάτω όρια βασισμένα στην παρατήρηση ότι ένα σύστημα σε ισορροπία έχει καλύτερη απόδοση από ένα σύστημα που δεν βρίσκεται σε ισορροπία.
- Σύστημα σε ισορροπία : σύστημα χωρίς bottleneck. Οι συνολικοί χρόνοι εξυπηρέτησης των αιτήσεων είναι ίσοι σε όλα τα devices.
- Η απόδοση ενός unbalanced συστήματος μπορεί να βελτιωθεί με την αντικατάσταση του bottleneck από ένα πιο γρήγορο device, μέχρι να δημιουργηθεί κάποιο άλλο bottleneck.



Balanced Job Bounds II

Ο χρόνος απόκρισης και το throughput ενός timesharing συστήματος φράσσονται ως εξής:

$$\max \left\{ ND_{max} - Z, D + (N - 1)D_{avg} \frac{D}{D + Z} \right\} \leq R^{(N)} \leq D + (N - 1)D_{max} \frac{(N - 1)D}{(N - 1)D + Z}$$
$$\frac{N}{Z + D + (N - 1)D_{max} \frac{(N - 1)D}{(N - 1)D + Z}} \leq X^{(N)} \leq \min \left\{ \frac{1}{D_{max}}, \frac{N}{Z + D + (N - 1)D_{avg} \frac{D}{D + Z}} \right\}$$

όπου, $D_{avg} = D/M$. Τα όρια είναι αυστηρά.



Balanced Job Bounds III

- Υποθέσεις: Όλα τα *service centers* είναι *fixed capacity*. Τα τερματικά είναι *delay centers*. Τα *balanced job bounds* προκύπτουν από τα εξής βήματα:
 1. Βρίσκουμε έκφραση για τα X, R ενός *balanced* συστήματος.
 2. Με δεδομένο ένα *unbalanced* σύστημα, κατασκευάζουμε ένα αντίστοιχο "*best case*" *balanced*, με ίδιο αριθμό *devices* και άθροισμα απαιτήσεων *service*. Από αυτό προκύπτουν τα άνω όρια.
 3. Κατασκευάζουμε ένα αντίστοιχο "*worst case*" *balanced*, ώστε κάθε *device* έχει *demand* ίσο με *bottleneck* και ο αριθμός των *devices* προσαρμόζεται ώστε το άθροισμα των απαιτήσεων να είναι ίδιο σε *balanced* και *unbalanced* συστήματα. Από αυτό προκύπτουν τα κάτω όρια.



Balanced Job Bounds IV

Απόδειξη:

Βήμα 1

Έστω σύστημα με κεντρικό υποσύστημα balanced:

$$D_i = \frac{D}{M}$$

Οι χρόνοι απόκρισης κάθε device, με MVA είναι: (εκθετικοί servers)

$$R_i^{(N)} = S_i(1 + Q_i^{(N-1)}) \quad i = 1, 2, \dots, M$$

Το σύστημα είναι balanced, άρα:

$$Q_i^{(N-1)} = \frac{Q^{(N-1)}}{M}$$

όπου $Q^{(j)}$ είναι ο συνολικός αριθμός jobs στο subsystem όταν υπάρχουν j jobs στο σύστημα. Ο αριθμός των jobs στα terminals είναι $j - Q^{(j)}$.



Balanced Job Bounds V

Η απόκριση του συστήματος δίνεται:

$$R^{(N)} = \sum_{i=1}^M V_i R_i^{(N)} = \sum_{i=1}^M \frac{D}{M} \left(1 + \frac{Q^{(N-1)}}{M} \right) \Rightarrow$$
$$R^{(N)} = D \left(1 + \frac{Q^{(N-1)}}{M} \right) = D + \frac{D}{M} Q^{(N-1)}$$

Διαδικασία να βάλουμε bound στο $Q^{(N)}$:

1. Αν αντικαταστήσουμε το σύστημα με N workstations έτσι ώστε κάθε χρήστης να έχει το δικό του workstation και κάθε workstation είναι ίδιο με το αρχικό σύστημα, το νέο σύστημα θα έχει καλύτερο R και X . Το νέο περιβάλλον έχει N single user systems. Κάθε χρήστης κάνει κύκλους με Z χρόνο σκέψης και D χρόνο computing. Κάθε job έχει πιθανότητα $\frac{D}{D+Z}$ να είναι στο κεντρικό υποσύστημα και συνεπώς:

$$\frac{Q^{(N)}}{N} \geq \frac{D}{D+Z}$$

(4)



Balanced Job Bounds VI

2. Τώρα, θεωρούμε ένα άλλο περιβάλλον, ίδιο με το προηγούμενο, με τη διαφορά ότι κάθε χρήστης έχει workstation N φορές αργότερο από το αρχικό σύστημα. Το νέο περιβάλλον έχει συνολική υπολογιστική ισχύ ίδια με το αρχικό, αλλά δεν υπάρχει διαμοιρασμός (sharing). Δηλαδή, οι χρήστες περνούν περισσότερο χρόνο στο κεντρικό σύστημα:

$$\frac{Q^{(N)}}{N} \geq \frac{ND}{ND+Z} \quad (5)$$

- Από τις 4 και 5 προκύπτει το εξής όριο για τον αριθμό στις devices:

$$\begin{aligned} \frac{D}{D+Z} &\leq \frac{Q^{(N)}}{N} \leq \frac{ND}{ND+Z} \Rightarrow \\ (N-1) \frac{D}{D+Z} &\leq Q^{(N-1)} \leq (N-1) \frac{(N-1)D}{(N-1)D+Z} \Rightarrow \\ D + \frac{D}{M} (N-1) \frac{D}{D+Z} &\leq R^{(N)} \leq D + \frac{D}{M} (N-1) \frac{(N-1)D}{(N-1)D+Z} \end{aligned}$$



Balanced Job Bounds VII

Βήμα 2

- Έστω ότι έχουμε *unbalanced* σύστημα τέτοιο ώστε, οι απαιτήσεις εξυπηρέτησης στο i -οστό device να είναι D_i . Θεωρούμε την εξής διάταξη:

$$D_1 \leq D_2 \leq \dots \leq D_M$$

- Κατ' αυτόν τον τρόπο το M -στο device είναι το αργότερο (*bottleneck device*) και το πρώτο το γρηγορότερο. Εκτελούμε το παρακάτω πείραμα:
 1. Κάνουμε το bottleneck device λίγο γρηγορότερο και το γρηγορότερο device, λίγο αργότερο. Η performance θα βελτιωθεί.
 2. Συνεχίζουμε τη διαδικασία μέχρι το bottleneck να γίνει ίσο με το bottleneck κάποιου άλλου device. Συνεπώς, διαθέτουμε 2 bottleneck devices.



Balanced Job Bounds VIII

3. Παίρνουμε τα δύο πιο αργά και τα δύο πιο γρήγορα devices στο σύστημα και εφαρμόζουμε το βήμα 1. Αυτό θα βελτιώσει την performance μέχρι και το τρίτο device να μπει στο bottleneck group. Την ίδια στιγμή προσθέτουμε ένα device στο group των γρήγορων devices.
4. Στο τέλος, θα έχουμε ένα balanced σύστημα όπου κάθε device θα έχει:

$$D_i = D_{avg} = \frac{1}{M} \sum_{i=1}^M D_i = \frac{D}{M}$$

- Το νέο balanced σύστημα έχει καλύτερη performance. Η παρατήρηση αυτή, οδηγεί στα ακόλουθα όρια της απόδοσης του unbalanced συστήματος:

$$R^{(N)} \geq D + (N - 1)D_{avg} \frac{D}{D + Z} \quad (6)$$

$$X^{(N)} \leq \frac{N}{Z + D + (N - 1)D_{avg} \frac{D}{D + Z}} \quad (7)$$



Balanced Job Bounds IX

Βήμα 3

Ένα άλλο πείραμα στο unbalanced σύστημα:

1. Τα devices με D_i ίσο με D_{max} (bottleneck) απομακρύνονται από το σύνολο των devices που πρόκειται να δεχθούν αλλαγές. Από τα υπόλοιπα, παίρνουμε το γρηγορότερο και το αργότερο. Έστω ότι αυτά είναι τα devices 1 και k , αντίστοιχα. Ισχύει, $k = M - 1$, εκτός κι αν, $D_M = D_{M-1}$, οπότε και παίρνουμε το μεγαλύτερο k με $D_k \neq D_M$. Μειώνουμε κατά ΔD το γρηγορότερο και αυξάνουμε κατά ΔD το k .
2. Συνεχίζουμε την διαδικασία, μέχρι το γρηγορότερο να γίνει ίσο με 0 και το αργό να γίνει ίσο με D_{max} . Όποιο από τα devices φτάσει πρώτο το όριο της, απομακρύνεται από το σύνολο των devices που πρόκειται να δεχθούν αλλαγές.
3. Επαναλαμβάνουμε το βήμα 1 για το νέο σύνολο devices, κ.ο.κ.



Balanced Job Bounds X

4. Στο τέλος θα έχουμε ένα balanced σύστημα όπου $M' = \frac{D}{D_{max}}$ devices έχουν demands ίσα με D_{max} . Τα υπόλοιπα έχουν 0 και απομακρύνονται από το σύστημα.
5. Αφού κάθε αλλαγή χειροτερεύει την performance, τελικά θα έχουμε χειρότερη performance από το unbalanced. Δηλαδή,

$$D + (N - 1)D_{max} \frac{(N - 1)D}{(N - 1)D + Z} \geq R^{(N)}$$
$$\frac{N}{Z + D + (N - 1)D_{max} \frac{(N - 1)D}{(N - 1)D + Z}} \leq X^{(N)}$$

Συνδυάζοντας τις 6, 7 με τις παραπάνω προκύπτουν τα ασυμπτωτικά όρια.



Σημείωμα Αναφοράς

Copyright Πανεπιστήμιο Πατρών, Γιάννης Γαροφαλάκης, Αθανάσιος Ν. Νικολακόπουλος. «Ανάλυση Απόδοσης Πληροφοριακών Συστημάτων. Λειτουργική Ανάλυση». Έκδοση: 1.1. Πάτρα 2016. Διαθέσιμο από τη δικτυακή διεύθυνση:
<https://eclass.upatras.gr/courses/CEID1094/>.



Σημείωμα Αδειοδότησης

Το παρόν υλικό διατίθεται με τους όρους της άδειας χρήσης Creative Commons Αναφορά, Μη Εμπορική Χρήση Παρόμοια Διανομή 4.0 (1) ή μεταγενέστερη, Διεθνής Έκδοση. Εξαιρούνται τα αυτοτελή έργα τρίτων π.χ. φωτογραφίες, διαγράμματα κ.λ.π., τα οποία εμπεριέχονται σε αυτό και τα οποία αναφέρονται μαζί με τους όρους χρήσης τους στο «Σημείωμα Χρήσης Έργων Τρίτων».



Ως **Μη Εμπορική** ορίζεται η χρήση:

- ▣ που δεν περιλαμβάνει άμεσο ή έμμεσο οικονομικό όφελος από την χρήση του έργου, για το διανομέα του έργου και αδειοδόχο
- ▣ που δεν περιλαμβάνει οικονομική συναλλαγή ως προϋπόθεση για τη χρήση ή πρόσβαση στο έργο
- ▣ που δεν προσπορίζει στο διανομέα του έργου και αδειοδόχο έμμεσο οικονομικό όφελος (π.χ. διαφημίσεις) από την προβολή του έργου σε διαδικτυακό τόπο



Ο δικαιούχος μπορεί να παρέχει στον αδειοδόχο ξεχωριστή άδεια να χρησιμοποιεί το έργο για εμπορική χρήση, εφόσον αυτό του ζητηθεί.