

# **Κεφάλαιο 6**

---

**Αποθήκευση και  
είσοδος/έξοδος**



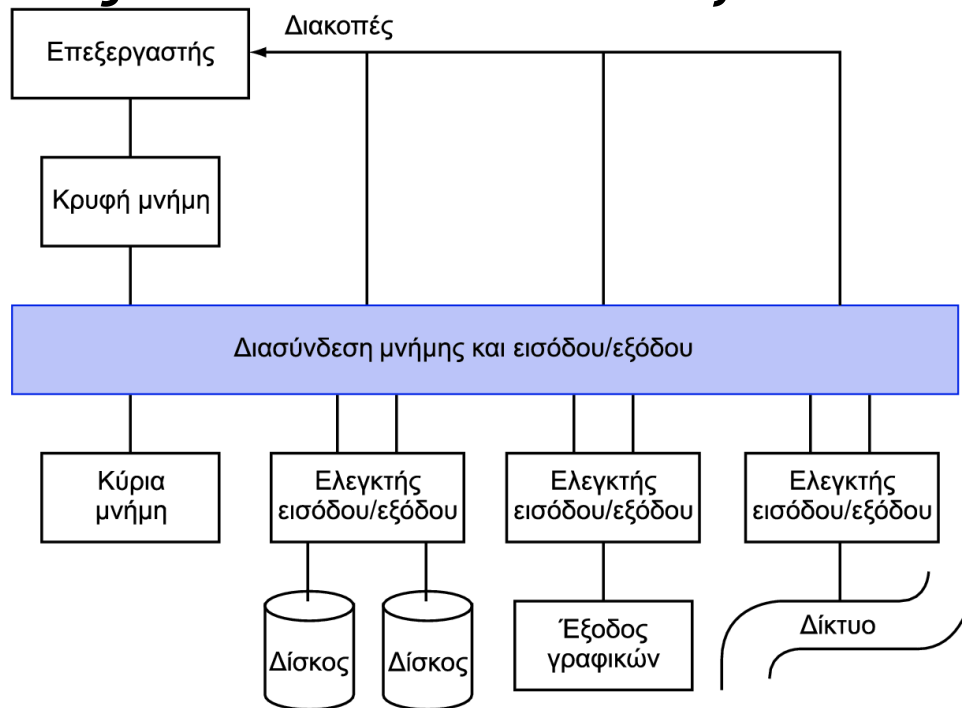
1956  
5 MB  
\$35.000

2015  
750 GB  
\$85



# Εισαγωγή

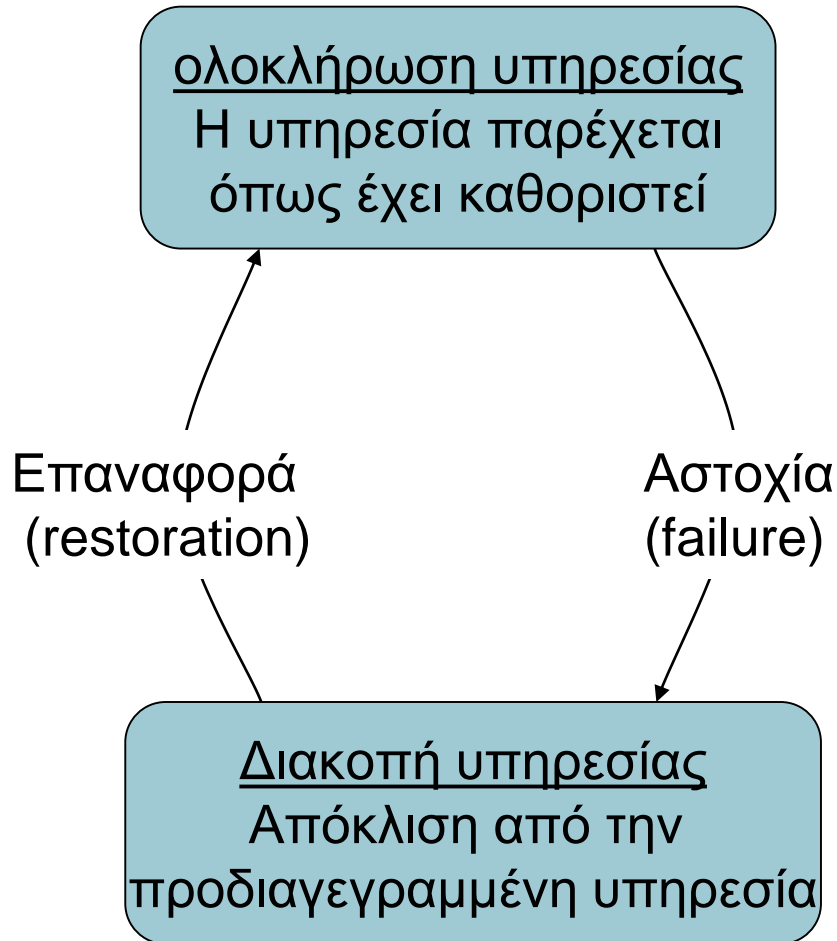
- οι συσκευές εισόδου/εξόδου χαρακτηρίζονται από
  - Συμπεριφορά: είσοδος, έξοδος, αποθήκευση
  - Εταίρο: άνθρωπος ή μηχανή
  - Ρυθμό δεδομένων: byte/sec, transfers/sec
- Συνδέσεις διαύλου εισόδου/εξόδου



# Χαρακτηριστικά συστήματος Ε/Ε

- Η φερεγγυότητα (dependability) είναι σημαντική
  - Ειδικά για συσκευές αποθήκευσης
- Μέτρα απόδοσης
  - Λανθάνων χρόνος (latency) ή χρόνος απόκρισης (response time)
  - Διεκπεραιωτική ικανότητα (throughput) ή εύρος ζώνης (bandwidth)
  - Επιτραπέζια και ενσωματωμένα συστήματα
    - Ενδιαφέρονται κυρίως για το χρόνο απόκρισης και την ποικιλομορφία των συσκευών
  - Διακομιστές
    - Ενδιαφέρονται κυρίως για τη διεκπεραιωτική ικανότητα και την επεκτασιμότητα των συσκευών

# Φερεγγυότητα (dependability)



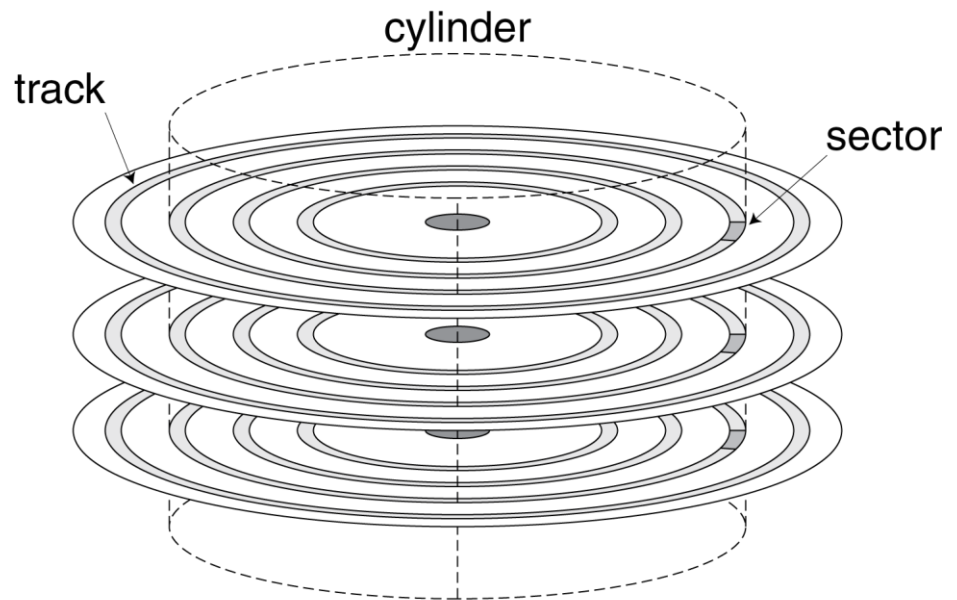
- Ελάττωμα (fault):  
αστοχία ενός  
συστατικού
  - Μπορεί να οδηγήσει ή  
να μην οδηγήσει σε  
αστοχία του  
συστήματος

# Μέτρα φερεγγυότητας

- **Αξιοπιστία** (reliability): μέσος χρόνος πρώτης αστοχίας (mean time to failure – MTTF)
- **Διακοπή υπηρεσίας** (service interruption): μέσος χρόνος επιδιόρθωσης (mean time to repair – MTTR)
- **Μέσος χρόνος μεταξύ αστοχιών** (mean time between failures - MTBF)
  - $MTBF = MTTF + MTTR$
- **Διαθεσιμότητα** (availability) =  $MTTF / (MTTF + MTTR)$
- Βελτίωση διαθεσιμότητας
  - Αύξηση MTTF: αποφυγή ελαττώματος (fault avoidance), ανοχή ελαττωμάτων (fault tolerance), πρόβλεψη ελαττωμάτων (fault forecasting)
  - Μείωση MTTR: βελτιωμένα εργαλεία και διαδικασίες διάγνωσης και επιδιόρθωσης

# Αποθήκευση στο δίσκο

- Μη πτητική (nonvolatile), περιστρεφόμενη μαγνητική αποθήκευση



# Τομείς δίσκου και προσπέλαση

- Κάθε τομέας (sector) καταγράφει
  - Την ταυτότητα τομέα (sector ID)
  - Δεδομένα (512 byte, 4096 byte προτεινόμενη τιμή)
  - Κώδικα διόρθωσης σφαλμάτων (error correcting code – ECC)
    - Για απόκρυψη ατελειών και καταγραφή σφαλμάτων
  - Πεδία συγχρονισμού και κενά (gaps)
- Η προσπέλαση ενός τομέα περιλαμβάνει
  - Καθυστέρηση αναμονής σε ουρά αν εκκρεμούν άλλες προσπελάσεις
  - Αναζήτηση (seek): μετακίνηση των κεφαλών
  - Λανθάνων χρόνος περιστροφής (rotational latency)
  - Μεταφορά δεδομένων
  - Επιβάρυνση ελεγκτή (controller)



# Παράδειγμα προσπέλασης δίσκου

- Δεδομένα
  - Τομέας των 512B, 15.000rpm (περιστροφές ανά λεπτό), μέσος χρόνος αναζήτησης 4ms, ρυθμός μεταφοράς 100MB/s, επιβάρυνση ελεγκτή 0.2ms, δίσκος αδρανής
- Μέσος χρόνος ανάγνωσης
  - 4ms χρόνος αναζήτησης  
+  $\frac{1}{2} / (15,000/60) = 2\text{ms}$  λανθάνων χρόνος περιστροφής  
+  $512 / 100\text{MB/s} = 0.005\text{ms}$  χρόνος μεταφοράς  
+ 0.2ms καθυστέρηση ελεγκτή  
= 6.2ms
- Αν ο πραγματικός μέσος χρόνος αναζήτησης ήταν 1ms:
  - Μέσος χρόνος ανάγνωσης = 3.2ms

# Ζητήματα απόδοσης δίσκου

- Οι κατασκευαστές αναφέρουν το **μέσο χρόνο αναζήτησης**, με βάση **όλες** τις πιθανές αναζητήσεις
- Η τοπικότητα και ο χρονοπρογραμματισμός του ΛΣ οδηγούν σε **μικρότερους** μέσους χρόνους αναζήτησης

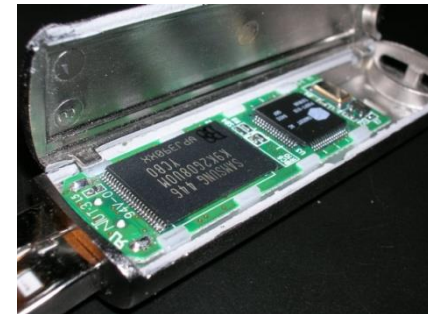
# Ελεγκτές δίσκων

- «Έξυπνος» ελεγκτής δίσκου κατανέμει τους φυσικούς τομείς του δίσκου
  - Εμφανίζει τη διασύνδεση των λογικών τομέων (logical sector interface) στον υπολογιστή
  - SCSI, ATA, SATA ελεγκτές
- Οι μονάδες δίσκου περιλαμβάνουν **και κρυφές μνήμες**
  - Εκ των προτέρων προσκόμιση (prefetch) τομέων με αναμονή προσπέλασής τους σύντομα
  - Αποφυγή αναζήτησης και καθυστέρησης περιστροφής



# Αποθήκευση σε μνήμη φλας

- Μη πτητική ημιαγωγική αποθήκευση
  - 100× – 1000× ταχύτερη από το δίσκο
  - Μικρότερο φυσικό μέγεθος,
  - Χαμηλότερη κατανάλωση ισχύος,
  - Πιο εύρωστη
  - Αλλά κοστίζει περισσότερα €/GB (ανάμεσα στο δίσκο και την DRAM)



# Τύποι μνήμης φλας

- NOR flash: κελί bit μοιάζει με πύλη NOR
  - Τυχαία προσπέλαση ανάγνωσης/εγγραφής
  - Χρησιμοποιείται για μνήμη εντολών σε ενσωματωμένα συστήματα
- NAND flash: κελί bit μοιάζει με πύλη NAND
  - Πιο πυκνή (bit/επιφάνεια), αλλά προσπέλαση ενός ολόκληρου μπλοκ τη φορά
  - Φθηνότερη ανά GB
  - Χρήση σε USB keys, αποθήκευση μέσω (ήχος, εικόνα), ...
- Τα bit της μνήμης φλας φθείρονται μετά από χιλιάδες προσπελάσεις
  - Δεν είναι κατάλληλη για να αντικαταστήσει τη RAM ή το δίσκο
  - Εξισορρόπηση φθοράς (wear leveling): επαναχαρτογράφηση δεδομένων σε λιγότερο χρησιμοποιημένα μπλοκ

# Συστατικά διασύνδεσης

- Ανάγκη διασύνδεσης μεταξύ
  - CPU, μνήμης, ελεγκτών E/E
- Δίαυλος: κοινόχρηστο κανάλι επικοινωνίας
  - Παράλληλο σύνολο αγωγών για δεδομένα και συγχρονισμό της μεταφοράς τους
  - Μπορεί να αποτελέσει σημείο συμφόρησης
- Η απόδοση περιορίζεται από φυσικούς παράγοντες
  - Μήκος αγωγού, αριθμός συνδέσεων
- Πιο πρόσφατη εναλλακτική: σειριακές συνδέσεις υψηλής ταχύτητας με μεταγωγούς
  - Όπως στα δίκτυα

# Τύποι διαύλου

- Δίαυλοι επεξεργαστή-μνήμης
  - Κοντοί, μεγάλη ταχύτητα
  - Η σχεδίαση ταιριάζει με την οργάνωση μνήμης
- Δίαυλοι εισόδου/εξόδου
  - Μακρύτεροι, επιτρέπουν πολλές συνδέσεις
  - Προδιαγράφονται με **πρότυπα** για λόγους διαλειτουργικότητας (interoperability)
  - Σύνδεση με το δίαυλο επεξεργαστή-μνήμης μέσω μιας γέφυρας (bridge)

# Σήματα διαύλου και συγχρονισμός

- Γραμμές δεδομένων
  - Μεταφέρουν διεύθυνση και δεδομένα
  - Με πολύπλεξη ή ξεχωριστά
- Γραμμές ελέγχου
  - Δείχνουν τον τύπο δεδομένων, συγχρονίζουν τις συναλλαγές (transactions)
- Σύγχρονη
  - Χρησιμοποιεί ρολόι διαύλου
- Ασύγχρονη
  - Χρησιμοποιεί γραμμές ελέγχου αίτησης/επιβεβαίωσης (request/acknowledge) για χειραψία (handshaking)



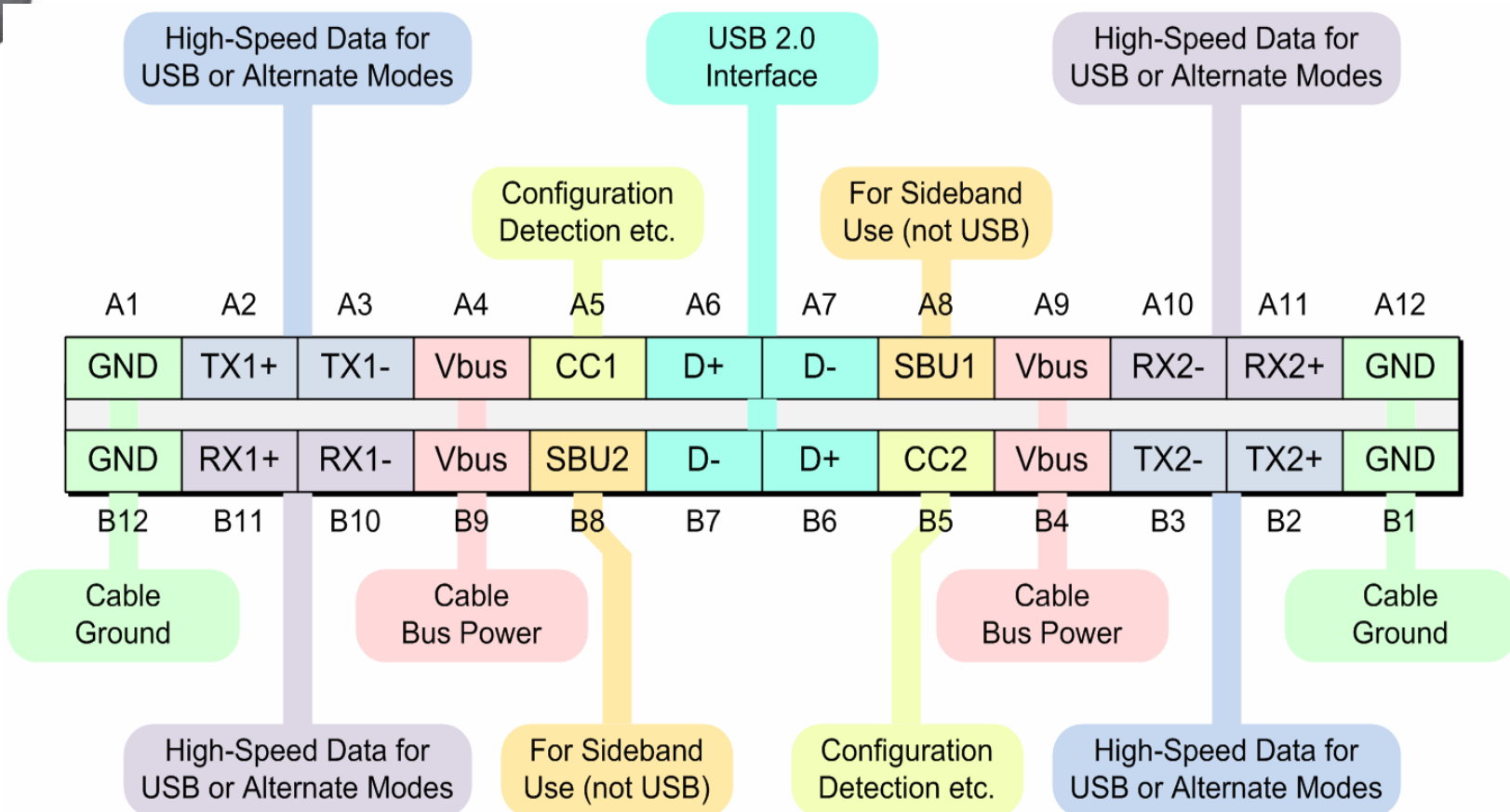
# Αρχαία Ιστορία



# Το Μέλλον (USB Type-C)



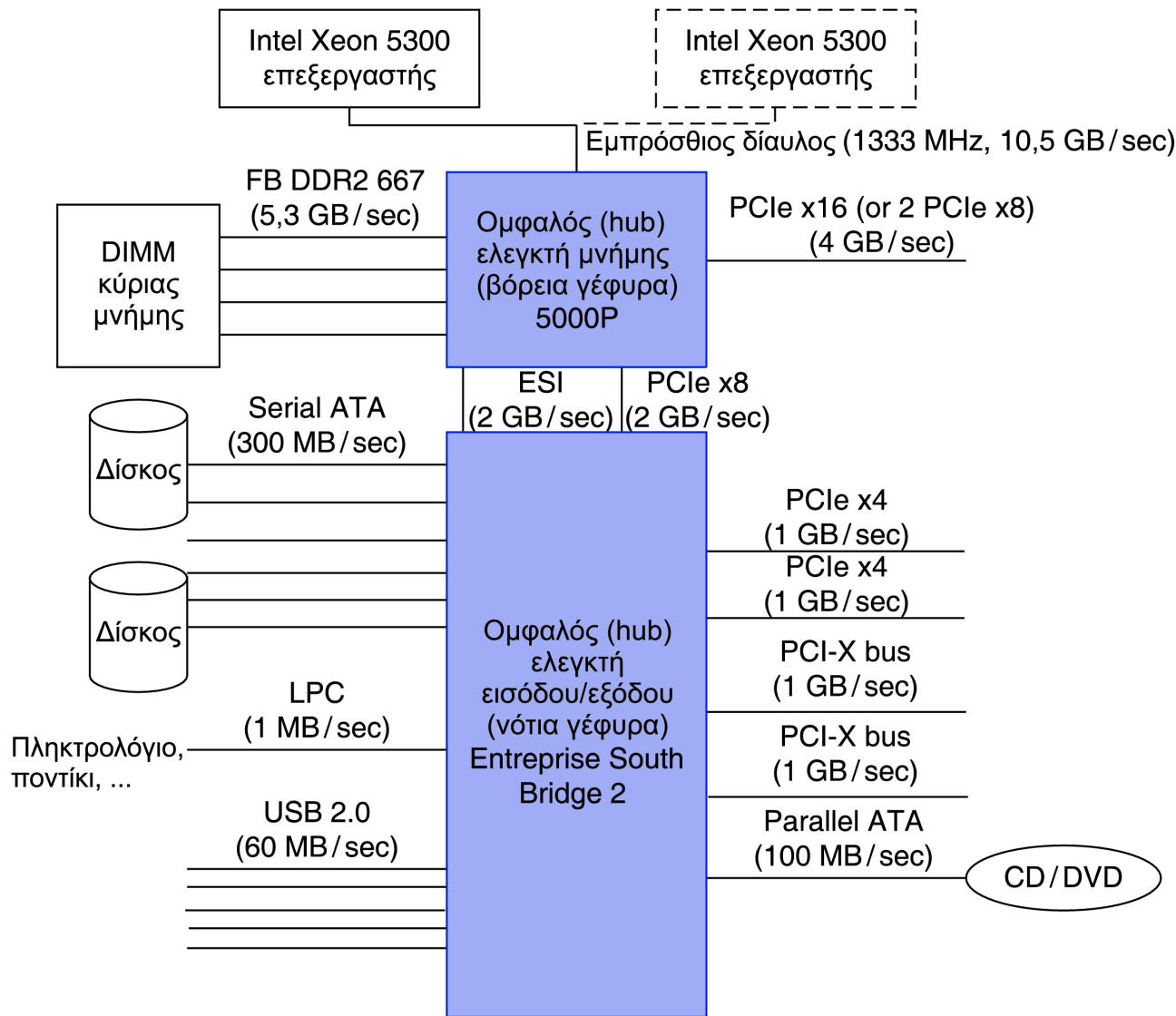
**20 Gbps, 100W**



# Παραδείγματα διαύλου E/E

	Firewire	USB 2.0	PCI Express	Serial ATA	Serial Attached SCSI
Πρόθεση χρήσης	Εξωτερική	Εξωτερική	Εσωτερική	Εσωτερική	Εξωτερική
Συσκευές ανά κανάλι	63	127	1	1	4
Εύρος δεδομένων	4	2	2/lane	4	4
Μέγιστο εύρος ζώνης	50MB/s ή 100MB/s	0.2MB/s, 1.5MB/s, ή 60MB/s	250MB/s/lane 1×, 2×, 4×, 8×, 16×, 32×	300MB/s	300MB/s
Σύνδεση «εν θερμώ»	Ναι	Ναι	Εξαρτάται	Ναι	Ναι
Μέγιστο μήκος	4.5m	5m	0.5m	1m	8m
Πρότυπο	IEEE 1394	USB Implementers Forum	PCI-SIG	SATA-IO	INCITS TC T10

# Τυπικό σύστημα E/E x86 PC



# Διαχείριση εισόδου/εξόδου

- Το ΛΣ είναι ο ενδιαμέσος για την Ε/Ε
  - Πολλά προγράμματα μοιράζονται πόρους εισόδου/εξόδου
    - Χρειάζεται προστασία και χρονοπρογραμματισμός
  - Η Ε/Ε προκαλεί ασύγχρονες διακοπές
    - Ίδιος μηχανισμός με τις εξαιρέσεις
  - ο προγραμματισμός Ε/Ε είναι δύσκολος
    - Το ΛΣ παρέχει αφαιρέσεις στα προγράμματα

# Διαταγές εισόδου/εξόδου

- Τις συσκευές Ε/Ε διαχειρίζεται το υλικό των ελεγκτών Ε/Ε
  - Μεταφέρουν δεδομένα από/προς τη συσκευή
  - Συγχρονίζουν τις λειτουργίες με λογισμικό
- Καταχωρητές διαταγών (command registers)
  - Αναγκάζουν τη συσκευή να κάνει κάτι
- Καταχωρητές κατάστασης (status registers)
  - Δείχνουν τι κάνει η συσκευή και την εμφάνιση σφαλμάτων
- Καταχωρητές δεδομένων (data registers)
  - Εγγραφής: μεταφέρουν δεδομένα σε μια συσκευή
  - Ανάγνωσης: μεταφέρουν δεδομένα από μια συσκευή

# Χαρτογράφηση καταχωρητών E/E

- E/E με χαρτογράφηση μνήμης (memory mapped I/O)
  - οι καταχωρητές προσπελάζονται στον ίδιο χώρο δ/νσεων με τη μνήμη
  - ο αποκωδικοποιητής δ/νσεων κάνει το διαχωρισμό
  - Το ΛΣ χρησιμοποιεί μηχανισμό μετάφρασης δ/νσεων ώστε να τους κάνει προσπελάσιμους μόνο στον πυρήνα (kernel) του ΛΣ
- Εντολές E/E
  - Ξεχωριστές εντολές για προσπέλαση καταχωρητών E/E
  - Μπορούν να εκτελεστούν μόνο σε κατάσταση πυρήνα
  - Παράδειγμα: x86

# Περίοδος (polling)

- Περιοδικός έλεγχος του καταχωρητή κατάστασης (status register) της E/E
  - Αν η συσκευή είναι έτοιμη, καμία λειτουργία
  - Αν υπάρχει σφάλμα, ανάληψη δράσης
- Συνήθης σε μικρά ή χαμηλών επιδόσεων ενσωματωμένα συστήματα πραγματικού χρόνου
  - Προβλέψιμος χρονισμός
  - Χαμηλό κόστος υλικού
- Σε άλλα συστήματα, σπατάλη χρόνου CPU



# Διακοπές (interrupts)

- Όταν μια συσκευή είναι έτοιμη ή όταν συμβεί σφάλμα
  - ο ελεγκτής διακόπτει τη CPU
- Η διακοπή είναι σαν εξαίρεση (exception)
  - Αλλά δε συγχρονίζεται με την εκτέλεση των εντολών
  - Μπορεί να καλέσει το χειριστή (handler) μεταξύ εντολών
  - Πληροφορία αιτίου (cause) προσδιορίζει συχνά τη συσκευή που προκαλεί διακοπή
- Διακοπές με προτεραιότητες
  - οι συσκευές που χρειάζονται πιο επείγουσα προσοχή λαμβάνουν υψηλότερη προτεραιότητα
  - Μπορούν να διακόψουν το χειριστή μιας διακοπής χαμηλότερης προτεραιότητας

# Μεταφορά δεδομένων E/E

- Περίοδευση και E/E οδηγούμενη από διακοπές
  - Η CPU μεταφέρει δεδομένα μεταξύ μνήμης και καταχωρητών δεδομένων E/E
  - Χρονοβόρα διαδικασία για συσκευές υψηλής ταχύτητας
- Άμεση προσπέλαση μνήμης (direct memory access – DMA)
  - Το ΛΣ παρέχει την αρχική δ/νση μνήμης
  - ο ελεγκτής E/E κάνει μεταφορά προς/από τη μνήμη αυτόνομα
  - ο ελεγκτής προκαλεί διακοπή όταν ολοκληρώσει τη μεταφορά ή σε περίπτωση σφάλματος

# Αλληλεπίδραση DMA/Cache

- Αν το DMA γράφει σε ένα μπλοκ μνήμης που βρίσκεται στην κρυφή μνήμη
  - Το αντίγραφο της κρυφής μνήμης γίνεται «παλιό»
- Αν η κρυφή μνήμη είναι ετερόχρονης εγγραφής και το μπλοκ είναι «ακάθαρτο», και το DMA διαβάζει το μπλοκ της μνήμης
  - Διαβάζει τα «παλιά» δεδομένα
- Πρέπει να εγγυηθούμε τη συνοχή (coherence) της κρυφής μνήμης
  - «Εκκένωση» (flush) των μπλοκ από τη κρυφή μνήμη αν πρόκειται να χρησιμοποιηθούν σε DMA
  - Ή χρήση θέσεων μνήμης που δεν αποθηκεύονται στη κρυφή μνήμη (non-cacheable) για τις λειτουργίες E/E

# Μέτρηση απόδοσης E/E

- Η απόδοση E/E εξαρτάται από
  - Υλικό: CPU, μνήμη, ελεγκτές, δίαυλοι
  - Λογισμικό: λειτουργικό σύστημα, σύστημα διαχείρισης βάσης δεδομένων, εφαρμογή
  - Φορτίο εργασίας: ρυθμοί και μοτίβα αιτήσεων
- Η σχεδίαση του συστήματος E/E μπορεί να κάνει συμβιβασμούς μεταξύ χρόνου απόκρισης και ρυθμού διεκπεραίωσης
  - οι μετρήσεις ρυθμού διεκπεραίωσης γίνονται συχνά με περιορισμένο χρόνο απόκρισης

# Μετροπρογράμματα επεξεργασίας συναλλαγών

- Συναλλαγές (Transactions)
  - Μικρές προσπελάσεις δεδομένων σε ένα σύστημα διαχείρισης βάσης δεδομένων (DBMS)
  - Το ενδιαφέρον είναι στο ρυθμό E/E, όχι το ρυθμό δεδομένων
- Μέτρηση ρυθμού διεκπεραίωσης (throughput)
  - Υπόκειται σε περιορισμούς χρόνου απόκρισης και χειρισμό αστοχιών
  - ACID (Atomicity/Ατομικότητα, Consistency/Συνέπεια, Isolation/Απομόνωση, Durability/Αντοχή)
  - Συνολικό κόστος ανά συναλλαγή
- Μετροπρογράμματα του Transaction Processing Council (TPC, [www.tpc.org](http://www.tpc.org))
  - TPC-APP: διακομιστής εφαρμογών και υπηρεσιών ιστού
  - TCP-C: περιβάλλον καταχώρισης παραγγελιών
  - TCP-E: επεξεργασία συναλλαγών μεσιτικού γραφείου
  - TPC-H: υποστήριξη αποφάσεων — κατά περίπτωση (ad-hoc) ερωτήματα με προσανατολισμό επιχειρήσεις

# Μετροπρογράμματα συστήματος αρχείων και Ιστού

- SPEC System File System (SFS)
  - Συνθετικό φορτίο εργασίας για διακομιστή NFS, με βάση παρακολούθηση πραγματικών συστημάτων
  - Αποτελέσματα
    - Ρυθμός διεκπεραίωσης, throughput (λειτουργίες/sec)
    - Χρόνος απόκρισης (μέσο ms/λειτουργία)
- SPEC Web Server benchmark
  - Μετράει τις ταυτόχρονες συνεδρίες (sessions) χρηστών, με βάση τον απαιτούμενο ρυθμό διεκπεραίωσης ανά συνεδρία
  - Τρία φορτία εργασίας: Τραπεζική, Ηλεκτρονικό εμπόριο, και Υποστήριξη

# Ε/Ε έναντι απόδοσης CPU

- Νόμος του Amdahl
  - Μην αγνοείς την απόδοση της Ε/Ε καθώς η παραλληλία αυξάνει την απόδοση των υπολογισμών
- Παράδειγμα
  - Το μετροπρόγραμμα διαρκεί 90s χρόνο CPU, 10s χρόνο Ε/Ε
  - Διπλάσιες CPU κάθε 2 χρόνια
    - Ε/Ε αμετάβλητη

Έτος	Χρόνος CPU	Χρόνος Ε/Ε	Παρελθών χρόνος	% Χρόνος Ε/Ε
Τώρα	90s	10s	100s	10%
+2	45s	10s	55s	18%
+4	23s	10s	33s	31%
+6	11s	10s	21s	47%

# RAID

- Πλεονασματικές συστοιχίες φθηνών (ανεξάρτητων) δίσκων – Redundant Array of Inexpensive (Independent) Disks
  - Χρήση πολλών μικρότερων δίσκων (σε σχέση με ένα μεγάλο )
  - Η παραλληλία βελτιώνει την απόδοση
  - Και πρόσθετοι δίσκοι για αποθήκευση πλεονασματικών δεδομένων
- Παρέχει σύστημα αποθήκευσης με ανοχή σε ελαττώματα (fault tolerant)
  - Ειδικά αν οι δίσκοι που αστοχούν δεν μπορούν να αντικατασταθούν «εν θερμώ»
- RAID 0
  - Χωρίς πλεονασμό (“AID”;)
    - Τα δεδομένα απλώς μοιράζονται σε πολλούς δίσκους
  - Αλλά βελτιώνει την απόδοση



# RAID 1 & 2

- RAID 1: Δημιουργία ειδώλων (mirroring)
  - $N + N$  δίσκοι, επανάληψη δεδομένων
    - Εγγραφή δεδομένων και στο δίσκο δεδομένων και στο δίσκο είδωλο
    - Σε περίπτωση αστοχίας δίσκου, ανάγνωση από το είδωλο
- RAID 2: Κώδικας διόρθωσης σφαλμάτων (Error correcting code – ECC)
  - $N + E$  δίσκοι (π.χ.,  $10 + 4$ )
  - Χωρισμός δεδομένων σε επίπεδο bit στους  $N$  δίσκους
  - Δημιουργία ECC των  $E$  bit
  - Υπερβολικά πολύπλοκο, δε χρησιμοποιείται στην πράξη

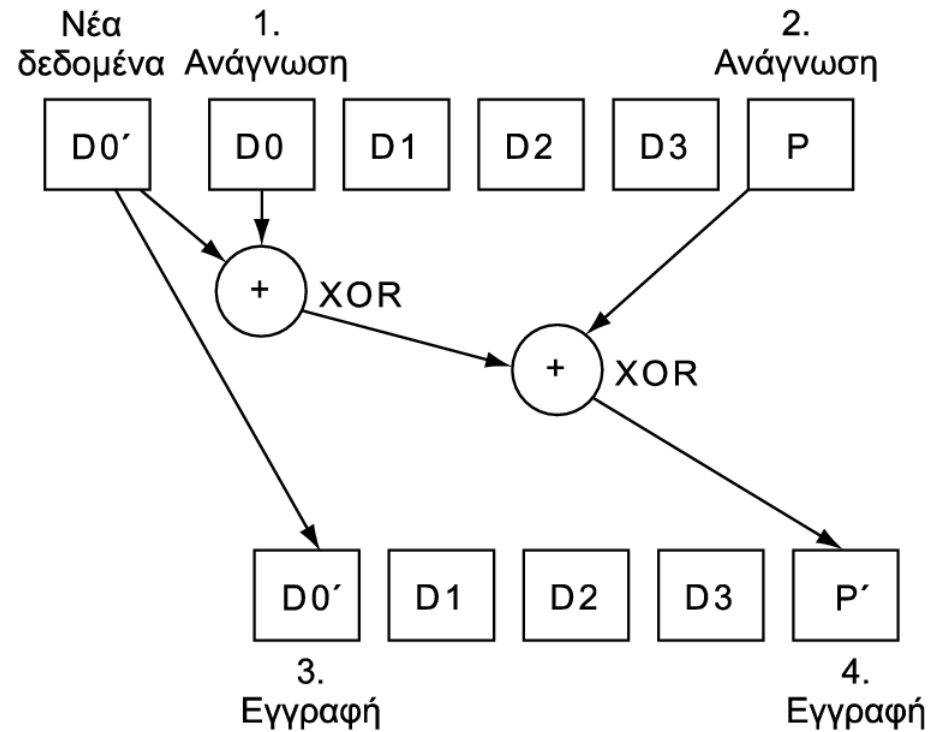
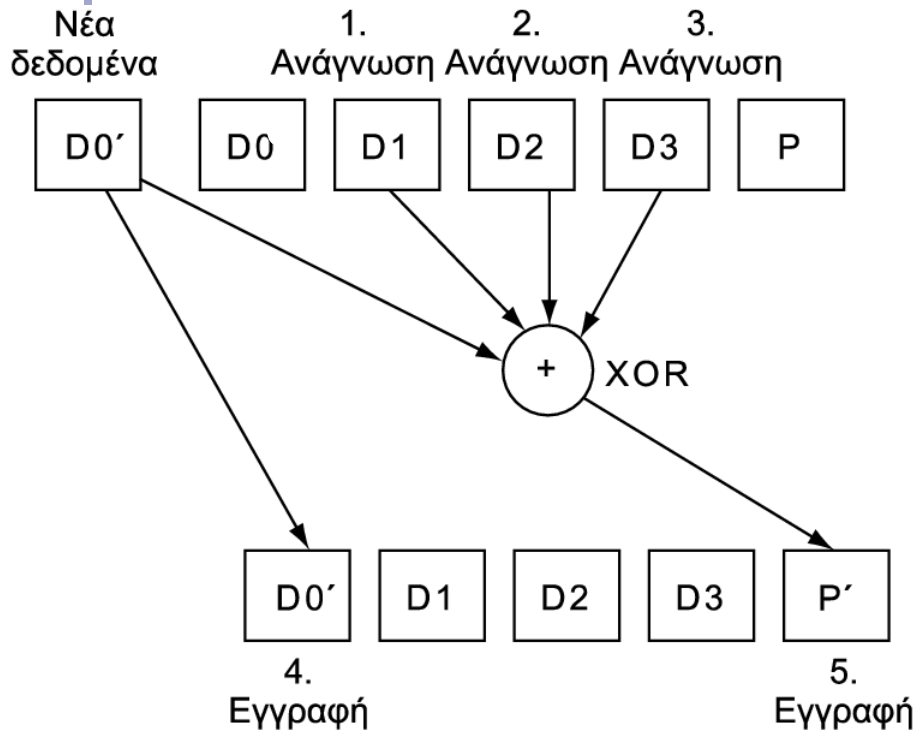
# RAID 3: Ισοτιμία πλέξης bit

- Bit-Interleaved Parity
- $N + 1$  δίσκοι
  - Δεδομένα μοιράζονται σε  $N$  δίσκους σε επίπεδο byte
  - Πλεονασματικός δίσκος αποθηκεύει την ισοτιμία
  - Προσπέλαση ανάγνωσης
    - Ανάγνωση όλων των δίσκων
  - Προσπέλαση εγγραφής
    - Δημιουργία νέας ισοτιμίας και ενημέρωση όλων των δίσκων
  - Σε περίπτωση αστοχίας
    - Χρήση ισοτιμίας για επανασύσταση των χαμένων δεδομένων
- Δε χρησιμοποιείται ευρέως

# RAID 4: Ισοτιμία πλέξης μπλοκ

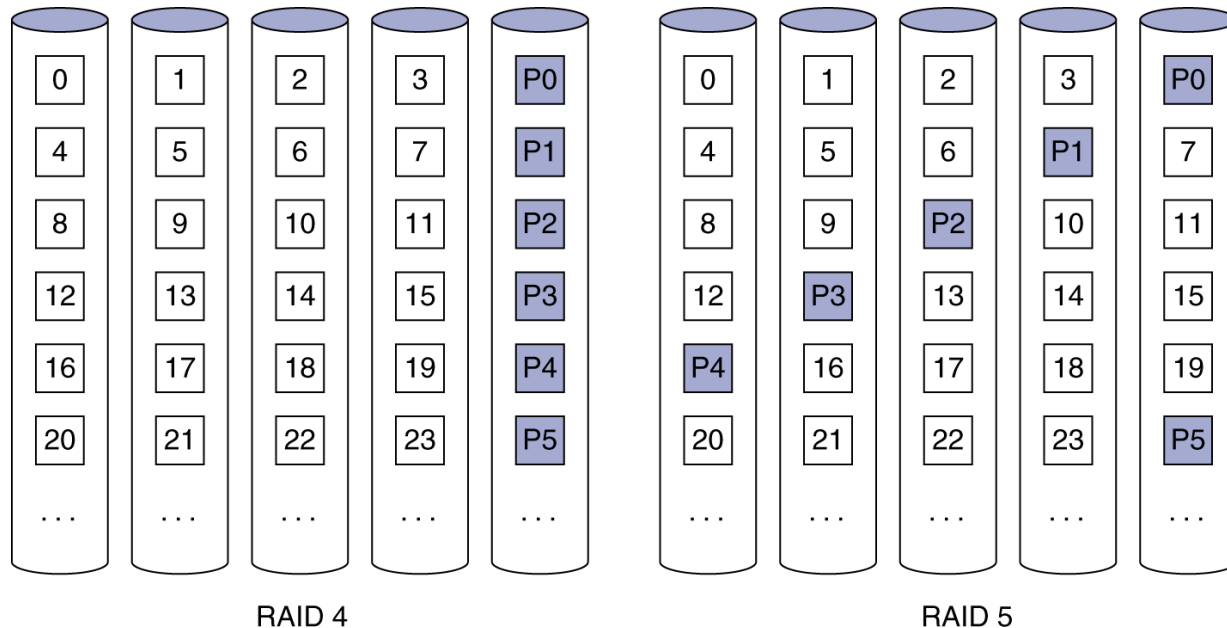
- Block-Interleaved Parity
- $N + 1$  δίσκοι
  - Τα δεδομένα μοιράζονται σε  $N$  δίσκους σε επίπεδο μπλοκ
  - Πλεονασματικός δίσκος αποθηκεύει την ισοτιμία για μια ομάδα μπλοκ
  - Προσπέλαση ανάγνωσης
    - Διαβάζει μόνο το δίσκο που περιέχει το ζητούμενο μπλοκ
  - Προσπέλαση εγγραφής
    - Απλώς διαβάζει το δίσκο οι οποίος περιέχει το μπλοκ που τροποποιείται, και το δίσκο ισοτιμίας
    - Υπολογισμός νέας ισοτιμίας, ενημέρωση δίσκου δεδομένων και δίσκου ισοτιμίας
  - Σε περίπτωση αστοχίας
    - Χρήση ισοτιμίας για την επανασύσταση των χαμένων δεδομένων
- Δε χρησιμοποιείται ευρέως

# RAID 3 έναντι RAID 4



# RAID 5: Κατανεμημένη ισοτιμία

- $N + 1$  δίσκοι
  - Όπως το RAID 4, αλλά τα μπλοκ ισοτιμίας κατανέμονται στους δίσκους
    - Αποφεύγει τη δημιουργία σημείου συμφόρησης (bottleneck) στο δίσκο ισοτιμίας
- Ευρεία χρήση



# RAID 6: Πλεονασμός P + Q

- P + Q Redundancy
- N + 2 δίσκοι
  - Σαν το RAID 5, αλλά με δύο «παρτίδες» ισοτιμίας
  - Μεγαλύτερη ανοχή σε ελαττώματα μέσω περισσότερου πλεονασμού
- Πολλαπλά RAID
  - Πιο προηγμένα συστήματα δίνουν παρόμοια ανοχή σε ελαττώματα με καλύτερη απόδοση

# Περίληψη RAID

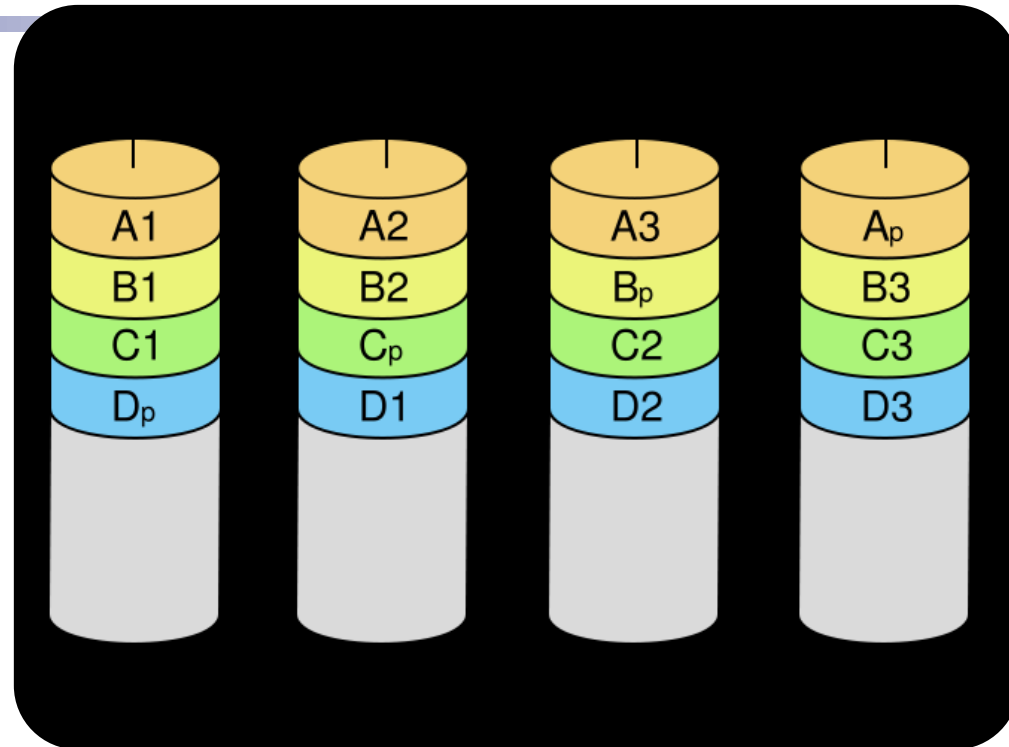
- Το RAID μπορεί να βελτιώσει την απόδοση και τη διαθεσιμότητα (availability)
  - Η υψηλή διαθεσιμότητα απαιτεί «εν θερμώ» εναλλαγή (hot swapping)
- Υποθέτει ότι οι αστοχίες δίσκων είναι ανεξάρτητες
  - Μεγάλο πρόβλημα αν καεί όλο το κτήριο!
- Δείτε το “Hard Disk Performance, Quality and Reliability”
  - <http://www.pcguides.com/ref/hdd/perf/index.htm>

# RAID 5 – Παράδειγμα

Συνεχόμενα blocks γράφονται εναλλάξ στους δίσκους, ενώ κατανέμεται σε αυτούς και ένα block ισοτιμίας.

Παρέχει υψηλή απόδοση στις αναγνώσεις, αφού αυτές μπορούν να γίνουν από πολλούς δίσκους εναλλάξ.

Παρέχει αξιοπιστία, αφού αν πάθει βλάβη ένας δίσκος, τα δεδομένα μπορούν να ανακτηθούν από τους υπόλοιπους .



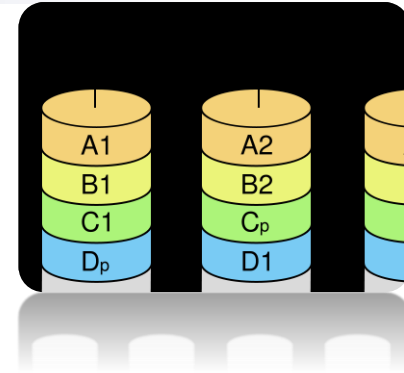


# RAID 5 - Παράδειγμα

Έστω ότι διαθέτουμε 4 δίσκους.  
Πώς δουλεύει το RAID 5;

Απάντηση: Ας θεωρήσουμε ότι οι 4 δίσκοι έχουν τα  
Παρακάτω δεδομένα (δυαδικό):

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	0010	
STRIPE1	0010	0000		0100
STRIPE2	0011		1010	1000
STRIPE3		0001	1101	1010



Στα κίτρινα σημεία, τοποθετούνται τα δεδομένα ιστοτιμίας. Η ιστοτιμία υπολογίζεται ως το Exclusive-OR (XOR) του ίδιου stripe όλων των δίσκων.

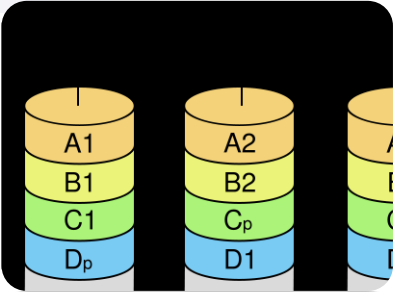
# RAID 5 - Παράδειγμα

Έστω ότι διαθέτουμε 4 δίσκους.  
 Πώς δουλεύει το RAID 5;

Απάντηση: Ας θεωρήσουμε ότι οι 4 δίσκοι έχουν τα Παρακάτω δεδομένα (δυναμικό):

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	0010	
STRIPE1	0010	0000		0100
STRIPE2	0011		1010	1000
STRIPE3		0001	1101	1010

Στα κίτρινα σημεία, τοποθετούνται τα δεδομένα ισότητας. Η ισότητα υπολογίζεται ως το Exclusive-OR (XOR) του ίδιου stripe όλων των δίσκων.



Για όσους δε θυμούνται της XOR ...

ο πίνακας αληθείας

XOR		
Είσοδος		Έξοδος
0	0	0
0	1	1
1	0	1
1	1	0

# RAID 5 - Παράδειγμα

	DISK0	DISK1	DISK2	DISK3
<b>STRIPE0</b>	<b>0100</b>	<b>0101</b>	<b>0010</b>	<b>0011</b>
STRIPE1	0010	0000		0100
STRIPE0	0011		1010	1000
STRIPE3		0001	1101	1010

STRIPE0,DISK3 = 0100 XOR 0101 XOR 0010 = 0011

# RAID 5 - Παράδειγμα

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	0010	0011
<b>STRIPE1</b>	<b>0010</b>	<b>0000</b>	<b>0110</b>	<b>0100</b>
STRIPE2	0011		1010	1000
STRIPE3		0001	1101	1010

STRIPE0,DISK3 = 0100 XOR 0101 XOR 0010 = 0011  
STRIPE1,DISK2 = 0010 XOR 0000 XOR 0100 = 0110

# RAID 5 - Παράδειγμα

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	0010	0011
STRIPE1	0010	0000	0110	0100
<b>STRIPE2</b>	<b>0011</b>	<b>0001</b>	<b>1010</b>	<b>1000</b>
STRIPE3		0001	1101	1010

STRIPE0,DISK3 = 0100 XOR 0101 XOR 0010 = 0011

STRIPE1,DISK2 = 0010 XOR 0000 XOR 0100 = 0110

STRIPE2,DISK1 = 0011 XOR 1010 XOR 1000 = 0001

# RAID 5 - Παράδειγμα

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	0010	0011
STRIPE1	0010	0000	0110	0100
STRIPE2	0011	0001	1010	1000
<b>STRIPE3</b>	<b>0110</b>	<b>0001</b>	<b>1101</b>	<b>1010</b>

STRIPE0,DISK3 = 0100 XOR 0101 XOR 0010 = 0011  
STRIPE1,DISK2 = 0010 XOR 0000 XOR 0100 = 0110  
STRIPE2,DISK1 = 0011 XOR 1010 XOR 1000 = 0001  
STRIPE3,DISK0 = 0001 XOR 1101 XOR 1010 = 0110

# RAID 5 - Παράδειγμα

Τελική Εικόνα της συστοιχίας ΔΙΣΚΩΝ  
με διάταξη RAID5

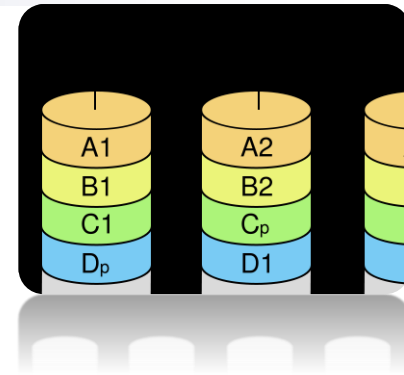
	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	0010	0011
STRIPE1	0010	0000	0110	0100
STRIPE2	0011	0001	1010	1000
STRIPE3	0110	0001	1101	1010

# RAID 5 - Παράδειγμα

Παράδειγμα: Τι γίνεται στις εγγραφές;

Απάντηση: Ας θεωρήσουμε ότι οι 4 δίσκοι έχουν τα Παρακάτω δεδομένα (δυναμικό):

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	0010	0011
STRIPE1	0010	0000	0110	0100
STRIPE2	0011	0001	1010	1000
STRIPE3	0110	0001	1101	1010



Έστω ότι γίνεται η εγγραφή του στοιχείου 1101 στο block 2 (αρίθμηση ξεκινάει από block 0).



# RAID 5 - Παράδειγμα

Έστω ότι γίνεται η εγγραφή του στοιχείου 1101 στο block 2  
(αρίθμηση ξεκινάει από block 0).

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	0010	0011
STRIPE1	0010	0000	0110	0100
STRIPE2	0011	0001	1010	1000
STRIPE3	0110	0001	1101	1010

# RAID 5 - Παράδειγμα

Έστω ότι γίνεται η εγγραφή του στοιχείου 1101 στο block 2 (αρίθμηση ξεκινάει από block 0).

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	<del>0010</del> 1101	0011
STRIPE1	0010	0000	0110	0100
STRIPE2	0011	0001	1010	1000
STRIPE3	0110	0001	1101	1010


ο ελεγκτής RAID κάνει την εγγραφή του στοιχείου στο αντίστοιχο block ...

# RAID 5 - Παράδειγμα

Έστω ότι γίνεται η εγγραφή του στοιχείου 1101 στο block 2 (αρίθμηση ξεκινάει από block 0).

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	<del>0010</del> 1101	<b>0011</b>
STRIPE1	0010	0000	0110	0100
STRIPE2	0011	0001	1010	1000
STRIPE3	0110	0001	1101	1010

ο ελεγκτής RAID κάνει την εγγραφή του στοιχείου στο αντίστοιχο block ... και ταυτόχρονα ξαναδημιουργεί την ισοτιμία για το συγκεκριμένο stripe, χρησιμοποιώντας παλιά τιμή, νέα τιμή και ισοτιμία

$$\text{STRIPE0,DISK3} = 0010 \text{ XOR } 1101 \text{ XOR } 0011 = 1100$$


# RAID 5 - Παράδειγμα

Έστω ότι γίνεται η εγγραφή του στοιχείου 1101 στο block 2  
(αρίθμηση ξεκινάει από block 0).

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	<del>0010</del> 1101	<del>0011</del> 1100
STRIPE1	0010	0000	0110	0100
STRIPE2	0011	0001	1010	1000
STRIPE3	0110	0001	1101	1010

ο ελεγκτής RAID κάνει την εγγραφή του στοιχείου στο αντίστοιχο block ... και ταυτόχρονα ξαναδημιουργεί την ισοτιμία για το συγκεκριμένο stripe, χρησιμοποιώντας παλιά τιμή, νέα τιμή και ισοτιμία  
 $STRIPE0, DISK3 = 0010 \text{ XOR } 1101 \text{ XOR } 0011 = 1100$

# RAID 5 - Παράδειγμα

Έστω ότι γίνεται η εγγραφή του στοιχείου 1101 στο block 2 (αρίθμηση ξεκινάει από block 0).

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	<b>1101</b>	<b>1100</b>
STRIPE1	0010	0000	<b>0110</b>	0100
STRIPE2	0011	<b>0001</b>	1010	1000
STRIPE3	<b>0110</b>	0001	1101	1010

ο ελεγκτής RAID κάνει την εγγραφή του στοιχείου στο αντίστοιχο block ... και ταυτόχρονα ξαναδημιουργεί την ισοτιμία για το συγκεκριμένο stripe, χρησιμοποιώντας παλιά τιμή, νέα τιμή και ισοτιμία  
 $STRIPE0, DISK3 = 0010 \text{ XOR } 1101 \text{ XOR } 0011 = 1100$

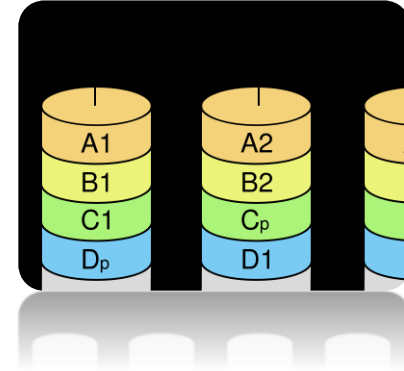
*Η εγγραφή στο RAID 5, ισοδυναμεί με 2 αναγνώσεις και 2 εγγραφές σε δίσκους.*

# RAID 5 - Παράδειγμα

Παράδειγμα: Τι γίνεται αν χαλάσει ένας δίσκος;

Απάντηση: Ας θεωρήσουμε ότι οι 4 δίσκοι έχουν τα Παρακάτω δεδομένα (δυναμικό):

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	1101	1100
STRIPE1	0010	0000	0110	0100
STRIPE2	0011	0001	1010	1000
STRIPE3	0110	0001	1101	1010



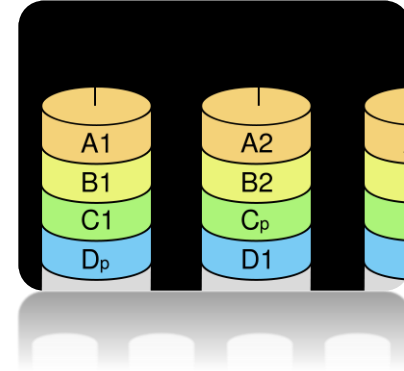
Έστω ότι χαλάει ο DISK2

# RAID 5 - Παράδειγμα

Παράδειγμα: Τι γίνεται αν χαλάσει ένας δίσκος;

Απάντηση: Ας θεωρήσουμε ότι οι 4 δίσκοι έχουν τα Παρακάτω δεδομένα (δυναμικό):

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	1101	1100
STRIPE1	0010	0000	0110	0100
STRIPE2	0011	0001	1010	1000
STRIPE3	0110	0001	1101	1010



Έστω ότι χαλάει ο DISK2

# RAID 5 - Παράδειγμα

Έστω ότι χαλάει ο DISK2

	DISK0	DISK1	DISK2	DISK3
<b>STRIPE0</b>	<b>0100</b>	<b>0101</b>	<del>1101</del>	<b>1100</b>
STRIPE1	0010	0000	<del>0110</del>	0100
STRIPE2	0011	<b>0001</b>	<del>1010</del>	1000
STRIPE3	<b>0110</b>	0001	<del>1101</del>	1010

ο ελεγκτής RAID εξυπηρετεί τις αιτήσεις για τις πληροφορίες που είχε ο DISK2, χρησιμοποιώντας όλους τους άλλους δίσκους + την ιστοιμιά.



# RAID 5 - Παράδειγμα

Έστω ότι χαλάει ο DISK2

	DISK0	DISK1	DISK2	DISK3
<b>STRIPE0</b>	<b>0100</b>	<b>0101</b>	<del>1101</del>	<b>1100</b>
STRIPE1	0010	0000	<del>0110</del>	0100
STRIPE2	0011	<b>0001</b>	<del>1010</del>	1000
STRIPE3	<b>0110</b>	0001	<del>1101</del>	1010

ο ελεγκτής RAID εξυπηρετεί τις αιτήσεις για τις πληροφορίες που είχε ο DISK2, χρησιμοποιώντας όλους τους άλλους δίσκους + την ισοτιμία.

Έτσι

$$\text{STRIPE0,DISK2} = 0100 \text{ XOR } 0101 \text{ XOR } 1100 = 1101$$

# RAID 5 - Παράδειγμα

Έστω ότι χαλάει ο DISK2

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	1101	1100
STRIPE1	0010	0000	0110	0100
<b>STRIPE2</b>	<b>0011</b>	<b>0001</b>	1010	<b>1000</b>
STRIPE3	0110	0001	1101	1010

ο ελεγκτής RAID εξυπηρετεί τις αιτήσεις για τις πληροφορίες που είχε ο DISK2, χρησιμοποιώντας όλους τους άλλους δίσκους + την ισοτιμία.

Έτσι

STRIPE0,DISK2 = 0100 XOR 0101 XOR 1100 = 1101

STRIPE2,DISK2 = 0011 XOR 0001 XOR 1000 = 1010

# RAID 5 - Παράδειγμα

Έστω ότι χαλάει ο DISK2

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	1101	1100
STRIPE1	0010	0000	0110	0100
STRIPE2	0011	0001	1010	1000
<b>STRIPE3</b>	<b>0110</b>	<b>0001</b>	1101	<b>1010</b>

ο ελεγκτής RAID εξυπηρετεί τις αιτήσεις για τις πληροφορίες που είχε ο DISK2, χρησιμοποιώντας όλους τους άλλους δίσκους + την ιστοιμιά.

Έτσι

STRIPE0,DISK2 = 0100 XOR 0101 XOR 1100 = 1101

STRIPE2,DISK2 = 0011 XOR 0001 XOR 1000 = 1010

STRIPE3,DISK2 = 0110 XOR 0001 XOR 1010 = 1101

# RAID 5 - Παράδειγμα

Έστω ότι χαλάει ο DISK2

	DISK0	DISK1	DISK2	DISK3
STRIPE0	0100	0101	1101	1100
STRIPE1	0010	0000	0110	0100
STRIPE2	0011	0001	1010	1000
STRIPE3	0110	0001	1101	1010

ο ελεγκτής RAID εξυπηρετεί τις αιτήσεις για τις πληροφορίες που είχε ο DISK2, χρησιμοποιώντας όλους τους άλλους δίσκους + την ισοτιμία.

Έτσι

STRIPE0,DISK2 = 0100 XOR 0101 XOR 1100 = 1101

STRIPE2,DISK2 = 0011 XOR 0001 XOR 1000 = 1010

STRIPE3,DISK2 = 0110 XOR 0001 XOR 1010 = 1101

*Κάθε ανάγνωση του χαλασμένου δίσκου, αντιστοιχεί σε αναγνώσεις σε όλους τους υπόλοιπους δίσκους. Καλό είναι να αντικαταστήσουμε το χαλασμένο δίσκο γρήγορα!*

# Διακομιστές

- οι εφαρμογές εκτελούνται όλο και περισσότερο σε διακομιστές (servers)
  - Αναζήτηση στον Ιστό, εφαρμογές γραφείου, εικονικοί κόσμοι, ...
- Απαιτούνται μεγάλοι διακομιστές κέντρων δεδομένων
  - Πολλοί επεξεργαστές, συνδέσεις δικτύου, μαζική αποθήκευση
  - Περιορισμοί χώρου και ηλεκτρικής ισχύος
- Εξοπλισμός διακομιστών για ικριώματα (racks) των 19 ιντσών
  - Ύψος σε πολλαπλάσια 1.75 ιντσών (1U)

# Διακομιστές για ικρίωμα

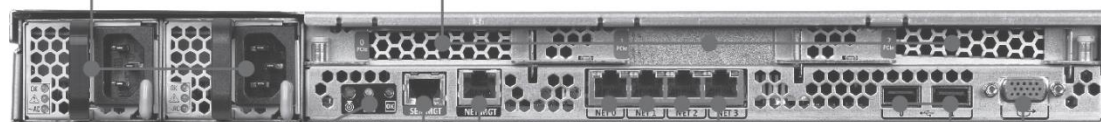


Διακομιστής Sun Fire x4150 1U



2 πλεονάζοντα  
τροφοδοτικά

3 υποδοχές PCI Express



LED κατάστασης συστήματος

Κάρτα διασύνδεσης  
δικτύου για διαχείριση

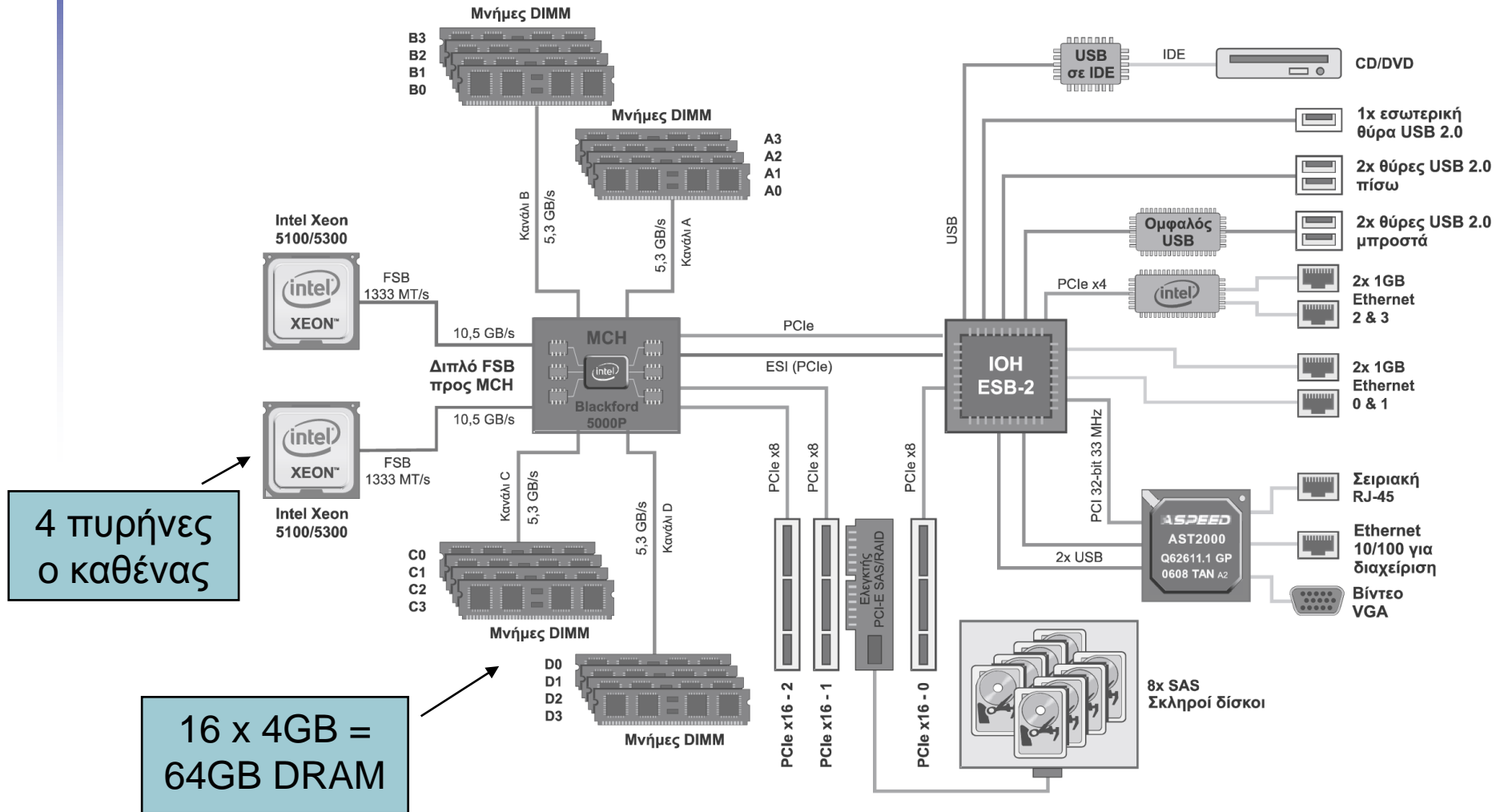
2 θύρες USB

Σειριακή θύρα  
για διαχείριση

Κάρτες διασύνδεσης  
δικτύου των 4 Gigabit

Εικόνα

# Διακομιστής Sun Fire x4150 1U



# Παράδειγμα σχεδίασης συστήματος E/E

- Δεδομένα: ένα σύστημα Sun Fire x4150 με
  - Φορτίο εργασία: αναγνώσεις δίσκου των 64KB
    - Κάθε λειτουργία E/E απαιτεί 200000 εντολές κώδικα χρήστη και 100000 εντολές του ΛΣ
  - Κάθε CPU:  $10^9$  εντολές/sec
  - FSB (Front Side Bus, Εμπρόσθιος δίαυλος): 10.6 GB/sec μέγιστο
  - DRAM DDR2 στα 667MHz: 5.336 GB/sec
  - PCI-E 8× δίαυλος:  $8 \times 250\text{MB/sec} = 2\text{GB/sec}$
  - Δίσκοι: 15000 rpm, 2.9ms μέσος χρόνος αναζήτησης, 112MB/sec ρυθμός μεταφοράς
- Ποιος είναι ο ρυθμός E/E που μπορεί να διατηρηθεί;
  - Για τυχαίες αναγνώσεις, και για ακολουθιακές αναγνώσεις



# Παράδειγμα σχεδίασης (συνέχεια)

- Ρυθμός E/E για τις CPU
  - Ανά πυρήνα:  $10^9 / (100000 + 200000) = 3333$
  - 8 πυρήνες: 26667 λειτουργίες/sec
- Τυχαίες αναγνώσεις, ρυθμός E/E για τους δίσκους
  - Υποθέστε ότι ο πραγματικός χρόνος αναζήτησης είναι το 1/4 μέσου
  - Χρόνος/λειτουργία = αναζήτηση + λανθάνων χρόνος + μεταφορά  
=  $2.9\text{ms}/4 + 4\text{ms}/2 + 64\text{KB}/(112\text{MB/s}) = 3.3\text{ms}$
  - 303 λειτουργίες/sec ανά δίσκο, 2424 λειτουργίες/sec για 8 δίσκους
- Ακολουθιακές αναγνώσεις
  - $112\text{MB/s} / 64\text{KB} = 1750$  λειτουργίες/sec ανά δίσκο
  - 14000 λειτουργίες/sec για 8 δίσκους

# Παράδειγμα σχεδίασης (συνέχεια)

- Ρυθμός E/E του PCI-E
  - $2\text{GB/sec} / 64\text{KB} = 31,250$  λειτουργίες/sec
- Ρυθμός E/E της DRAM
  - $5.336\text{ GB/sec} / 64\text{KB} = 83375$  λειτουργίες/sec
- Ρυθμός E/E του FSB
  - Υποθέστε ότι μπορούμε να διατηρήσουμε το μισό του μέγιστου ρυθμού
  - $5.3\text{ GB/sec} / 64\text{KB} = 81540$  λειτουργίες/sec ανά FSB
  - 163080 λειτουργίες/sec για 2 FSB
- ο πιο αδύναμος κρίκος: οι δίσκοι
  - 2424 λειτουργίες/sec τυχαίες, 14000 λειτουργίες/sec ακολουθιακές
  - Τα άλλα συστατικά έχουν άφθονο χώρο για να χωρέσουν αυτούς του ρυθμούς

# Πλάνη: Φερεγγυότητα δίσκων

- Αν ένας κατασκευαστής δίσκων δίνει ότι το MTTF είναι 1200000 ώρες (140 χρόνια)
  - Ένας δίσκος θα έχει τόσο μεγάλη διάρκεια
- Λάθος: αυτός είναι ο μέσος χρόνος πρώτης αστοχίας
  - Ποια είναι η κατανομή των αστοχιών;
  - Τι θα γίνει αν έχετε 1000 δίσκους;
    - Πόσοι θα αστοχούν κάθε χρόνο;

$$\text{Annual Failure Rate (AFR)} = \frac{1000 \text{ disks} \times 8760 \text{ hrs/disk}}{1200000 \text{ hrs/failure}} = 0.73\%$$

# Παγίδα: «ξεφόρτωμα» σε επεξεργαστές E/E

- Η επιβάρυνση της διαχείρισης των αιτήσεων του επεξεργαστή E/E μπορεί να κυριαρχεί
  - Ταχύτερο να γίνεται η μικρή λειτουργία στη CPU
  - Αλλά η αρχιτεκτονική E/E μπορεί να μην το επιτρέπει αυτό
- ο επεξεργαστής E/E μπορεί να είναι πιο αργός
  - Εφόσον υποτίθεται ότι είναι απλούστερος
- Αν γίνει ταχύτερος καθίσταται σημαντικό συστατικό του συστήματος
  - Μπορεί να χρειάζεται τους δικούς του συνεπεξεργαστές!

# Παγίδα: αντίγραφα ασφαλείας σε ταινία

- Η μαγνητική ταινία είχε πλεονεκτήματα
  - Φορητότητα, μεγάλη χωρητικότητα
- Τα πλεονεκτήματα άρχισαν να χάνονται με τις εξελίξεις της τεχνολογίας δίσκων
- Είναι πιο λογικό να επαναλαμβάνονται τα δεδομένα
  - Π.χ, RAID, απομακρυσμένη δημιουργία ειδώλων (remote mirroring)

# Παγίδα: μέγιστη απόδοση

- οι μέγιστοι ρυθμοί E/E είναι σχεδόν αδύνατον να επιτευχθούν
  - Συνήθως, κάποια άλλα συστατικά του συστήματος περιορίζουν την απόδοση
  - Π.χ., μεταφορές στη μνήμη μέσω ενός διαύλου
    - Σύγκρουση με την ανανέωση (refresh) της DRAM
    - Συναγωνισμός διαιτησίας με άλλου κύριους (masters) του διαύλου
  - Π.χ., δίαυλος PCI: μέγιστο εύρος ζώνης ~133 MB/sec
    - Στην πράξη, μπορεί να διατηρηθεί το 80MB/sec κατά μέγιστο

# Συμπερασματικές παρατηρήσεις

- Μέτρα απόδοσης E/E
  - Ρυθμός διεκπεραίωσης, χρόνος απόκρισης
  - Η φερεγγυότητα και επίσης το κόστος είναι σημαντικά
- Χρησιμοποιούνται δίαυλοι για τη σύνδεση CPU, μνήμη, ελεγκτές E/E
  - Περίοδευση, διακοπές, DMA
- Μετροπρογράμματα E/E
  - TPC, SPECint, SPECweb
- RAID
  - Βελτιώνει την απόδοση και τη φερεγγυότητα