



Πανεπιστήμιο Δυτικής Μακεδονίας  
Τμήμα Μηχανικών Πληροφορικής & Τηλεπικοινωνιών

# Τα κέντρα δεδομένων ως υπολογιστής

*Μια εισαγωγή στο σχεδιασμό των μηχανών αποθήκευσης*

Επιβλέπων καθηγητής  
Δρ. Μηνάς Δασυγένης

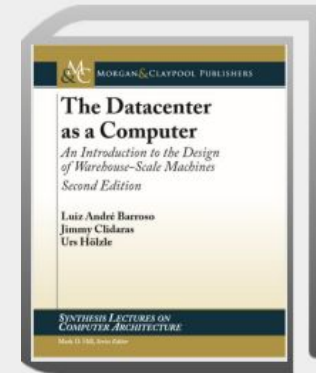
Επιμέλεια παρουσίασης  
Χριστιάνα Τσίρκα

Κοζάνη, 2018

Εργαστήριο Ψηφιακών Συστημάτων και Αρχιτεκτονικής Υπολογιστών

# ΒΑΣΙΚΟ ΣΤΟΙΧΕΙΟ ΤΗΣ ΠΑΡΟΥΣΙΑΣΗΣ

Η παρουσίαση αυτή βασίζεται στο βιβλίο με τίτλο “The Datacenter as a Computer”, 2η έκδοση, με συγγραφείς τους: Luiz André Barroso, Jimmy Clidaras, Urs Hölzle.



Το βιβλίο είναι γραμμένο στα Αγγλικά και είναι διαθέσιμο σε ηλεκτρονική μορφή στο διαδίκτυο.



- Φόρτος εργασίας και δομή Λογισμικού
- Δομή του Υλικού
- Βασικά στοιχεία των Κέντρων Δεδομένων
- Απόδοση σε ενέργεια και ισχύ
- Κόστη των διάφορων σχεδιασμών
- Χειρισμός βλαβών και επιδιορθώσεις



- Η εμφάνιση δημοφιλών υπηρεσιών Διαδικτύου ως e-mail, αναζήτησης και κοινωνικών δικτύων βασισμένων στο Web, καθώς και η αύξηση της παγκόσμιας διαθεσιμότητας σύνδεσης υψηλής ταχύτητας έχουν επιταχύνει την τάση προς υπολογιστές από πλευράς διακομιστή ή “cloud”.
- Ο υπολογιστής και ο αποθηκευτικός χώρος μετακινούνται από πελάτες τύπου PC (Personal Computer) σε μικρότερες, συχνά κινητές συσκευές, σε συνδυασμό με μεγάλες υπηρεσίες Διαδικτύου.
- Το λογισμικό ως υπηρεσία επιτρέπει ταχύτερη ανάπτυξη εφαρμογών, επειδή είναι πιο εύκολο για τους πωλητές λογισμικού να κάνουν αλλαγές και βελτιώσεις.



## ΕΙΣΑΓΩΓΗ ΣΤΑ ΚΕΝΤΡΑ ΔΕΔΟΜΕΝΩΝ (2/3)

- Τα οικονομικά δεδομένα του κέντρου δεδομένων επιτρέπουν σε πολλές υπηρεσίες εφαρμογών να εκτελούνται με χαμηλό κόστος ανά χρήστη.
- Ο ίδιος ο υπολογισμός μπορεί να γίνει φθηνότερος σε μια κοινόχρηστη υπηρεσία (π.χ. ένα συνημμένο ηλεκτρονικό ταχυδρομείο που λαμβάνεται από πολλούς χρήστες μπορεί να αποθηκευτεί μία φορά και όχι πολλές φορές).
- Οι διακομιστές και οι αποθηκευτικοί χώροι σε ένα κέντρο δεδομένων μπορούν να διαχειρίζονται ευκολότερα από την επιφάνειας εργασίας ή από ένα φορητό υπολογιστή.
- Οι υπηρεσίες αναζήτησης (Web, εικόνες, κλπ.) είναι ένα πρωταρχικό παράδειγμα αυτής της τάξης φόρτου εργασίας, αλλά εφαρμογές όπως η γλωσσική μετάφραση μπορούν επίσης να λειτουργήσουν πιο αποτελεσματικά σε μεγάλες κοινές υπολογιστικές εγκαταστάσεις λόγω της εξάρτησης από τα γλωσσικά μοντέλα μεγάλης κλίμακας.



## ΕΙΣΑΓΩΓΗ ΣΤΑ ΚΕΝΤΡΑ ΔΕΔΟΜΕΝΩΝ (3/3)

- Η τάση προς την πλευρά υπολογιστών από την πλευρά του διακομιστή και η εκρηκτική δημοτικότητα των υπηρεσιών του Διαδικτύου δημιούργησε μια νέα τάξη υπολογιστικών συστημάτων που έχουμε ονομάσει υπολογιστές αποθήκης ή WSCs (Warehouse-Scale Computers).
- Το όνομα προορίζεται να επιστήσει την προσοχή στο πιο χαρακτηριστικό γνώρισμα αυτών των μηχανών: την τεράστια κλίμακα της υποδομής του λογισμικού, των αποθετηρίων δεδομένων και της πλατφόρμας υλικού.
- Το υλικό για μια τέτοια πλατφόρμα αποτελείται από χιλιάδες ατομικούς κόμβους υπολογιστών με τα αντίστοιχα υποσυστήματα δικτύωσης και αποθήκευσης, εξοπλισμό διανομής ισχύος και κλιματισμού και εκτεταμένα συστήματα ψύξης. Το περίβλημα αυτών των συστημάτων είναι στην πραγματικότητα ένα κτίριο.



## WSCs ΥΠΟΛΟΓΙΣΤΕΣ ΑΠΟΘΗΚΕΥΣΗΣ (1/2)

---

- Τα κέντρα δεδομένων είναι κτίρια όπου συναποτελούνται πολλοί εξυπηρετητές και μέσα επικοινωνίας λόγω των κοινών περιβαλλοντικών απαιτήσεων και φυσικών αναγκών ασφαλείας και για ευκολία συντήρησης.
- Ένα WSC είναι ένας τύπος κέντρου δεδομένων.
- Τα WSCs τροφοδοτούν σήμερα τις υπηρεσίες που προσφέρουν εταιρείες όπως το Google, το Amazon, το Facebook και το διαδικτυακό τμήμα της Microsoft.



## WSCs ΥΠΟΛΟΓΙΣΤΕΣ ΑΠΟΘΗΚΕΥΣΗΣ (2/2)

- Διαφέρουν σημαντικά από τα παραδοσιακά κέντρα δεδομένων: ανήκουν σε έναν μόνο οργανισμό, χρησιμοποιούν μια σχετικά ομοιογενή πλατφόρμα λογισμικού υλικού και συστημάτων και μοιράζονται ένα κοινό επίπεδο διαχείρισης συστημάτων.
- Το πιο σημαντικό είναι ότι τα WSCs διαχειρίζονται μικρότερο αριθμό πολύ μεγάλων εφαρμογών (ή υπηρεσιών Διαδικτύου) και η κοινή υποδομή διαχείρισης πόρων επιτρέπει σημαντική ευελιξία ανάπτυξης.
- Η επίτευξη λειτουργίας χωρίς σφάλματα σε ένα τέτοιο σύστημα είναι δύσκολη και δυσχεραίνεται από τον μεγάλο αριθμό εμπλεκόμενων διακομιστών.





# ΑΠΟΤΕΛΕΣΜΑΤΙΚΟΤΗΤΑ ΚΟΣΤΟΥΣ

- Η κατασκευή και η λειτουργία μιας μεγάλης πλατφόρμας υπολογιστών είναι δαπανηρή και η ποιότητα μιας υπηρεσίας που παρέχεται μπορεί να εξαρτάται από τη συνολική διαθέσιμη χωρητικότητα επεξεργασίας και αποθήκευσης, την περαιτέρω αύξηση του κόστους οδήγησης και την ανάγκη επικέντρωσης στην αποδοτικότητα του κόστους.
- Η αναζήτηση στο Web, η αύξηση των αναγκών σε υπολογιστές οφείλεται σε τρεις κύριους παράγοντες:
  - Η αυξημένη δημοτικότητα των υπηρεσιών μεταφράζεται σε υψηλότερα φορτία αιτήσεων.
  - Το μέγεθος του προβλήματος συνεχίζει να αυξάνεται - ο ιστός αυξάνεται κατά εκατομμύρια σελίδες την ημέρα, γεγονός που αυξάνει το κόστος κατασκευής και προβολής ενός ευρετηρίου ιστού.
  - Βελτιώσεις ποιότητας με άλλους αλγόριθμους θα επιφέρει μεγαλύτερο κόστος.



# ΤΑ ΚΕΝΤΡΑ ΔΕΔΟΜΕΝΩΝ ΔΕΝ ΕΙΝΑΙ ΑΠΛΑ ΜΙΑ ΣΥΛΛΟΓΗ ΔΙΑΚΟΜΙΣΤΩΝ

- Τα κέντρα δεδομένων δεν είναι πλέον απλά μια συλλογή μηχανών τοποθετημένων σε μια εγκατάσταση και ενσύρματα.
- Είναι μια νέα κατηγορία μηχανών μεγάλης κλίμακας που οδηγούνται από μια νέα και ταχέως εξελισσόμενη σειρά φόρτου εργασίας.
- Η συμπεριφορά σφαλμάτων και οι εκτιμήσεις της ισχύος και της ενέργειας έχουν σημαντικότατο αντίκτυπο στον σχεδιασμό των WSCs.
- Τα WSCs έχουν ένα επιπλέον επίπεδο πολυπλοκότητας πέραν των συστημάτων που αποτελούνται από μεμονωμένους διακομιστές ή μικρές ομάδες διακομιστών.



# ΕΝΑ ΚΕΝΤΡΟ ΔΕΔΟΜΕΝΩΝ vs ΔΙΑΦΟΡΑ ΚΕΝΤΡΑ ΔΕΔΟΜΕΝΩΝ

---

- Σχεδιασμός ενός υπολογιστή ως κέντρο δεδομένων και όχι πολλών κέντρων δεδομένων που βρίσκονται μακριά.
- Χαρακτηριστικά παραδείγματα είναι οι υπηρεσίες που αφορούν τις σταθερές ενημερώσεις δεδομένων χρήστη και συνεπώς απαιτούν πολλαπλά αντίγραφα για λόγους προστασίας από την καταστροφή.
- Για τέτοιους υπολογισμούς, ένα σύνολο κέντρων δεδομένων μπορεί να είναι το πιο κατάλληλο σύστημα. Αλλά επιλέξαμε να σκεφτόμαστε το σενάριο πολλαπλών κέντρων δεδομένων ως περισσότερο ανάλογο με ένα δίκτυο υπολογιστών.



# Η ΣΗΜΑΣΙΑ ΤΩΝ WSCs ΣΕ ΕΜΑΣ

- Τα WSCs μπορεί να θεωρηθούν ως κάτι εξειδικευμένο, διότι το μέγεθος και το κόστος τους καθιστούν δύσκολη τη δημιουργία τους εκτός από μερικές μεγάλες εταιρίες Διαδικτύου που μπορούν να το πράξουν.
- Η ελκυστική οικονομία της φόρμουλας των διακομιστών χαμηλών επιπέδων τοποθετεί εκατοντάδες κόμβους σε ένα ευρύ φάσμα εταιρειών και ερευνητικών ιδρυμάτων.
- Όταν συνδυάζεται με τις τάσεις προς δημιουργία επεξεργαστών με μεγάλο αριθμό πυρήνων σε μία μοναδική μητρική, ένα μόνο ερμάριο διακομιστή μπορεί σύντομα να έχει τόσα ή περισσότερα υλικά από πολλά σύγχρονα κέντρα δεδομένων.



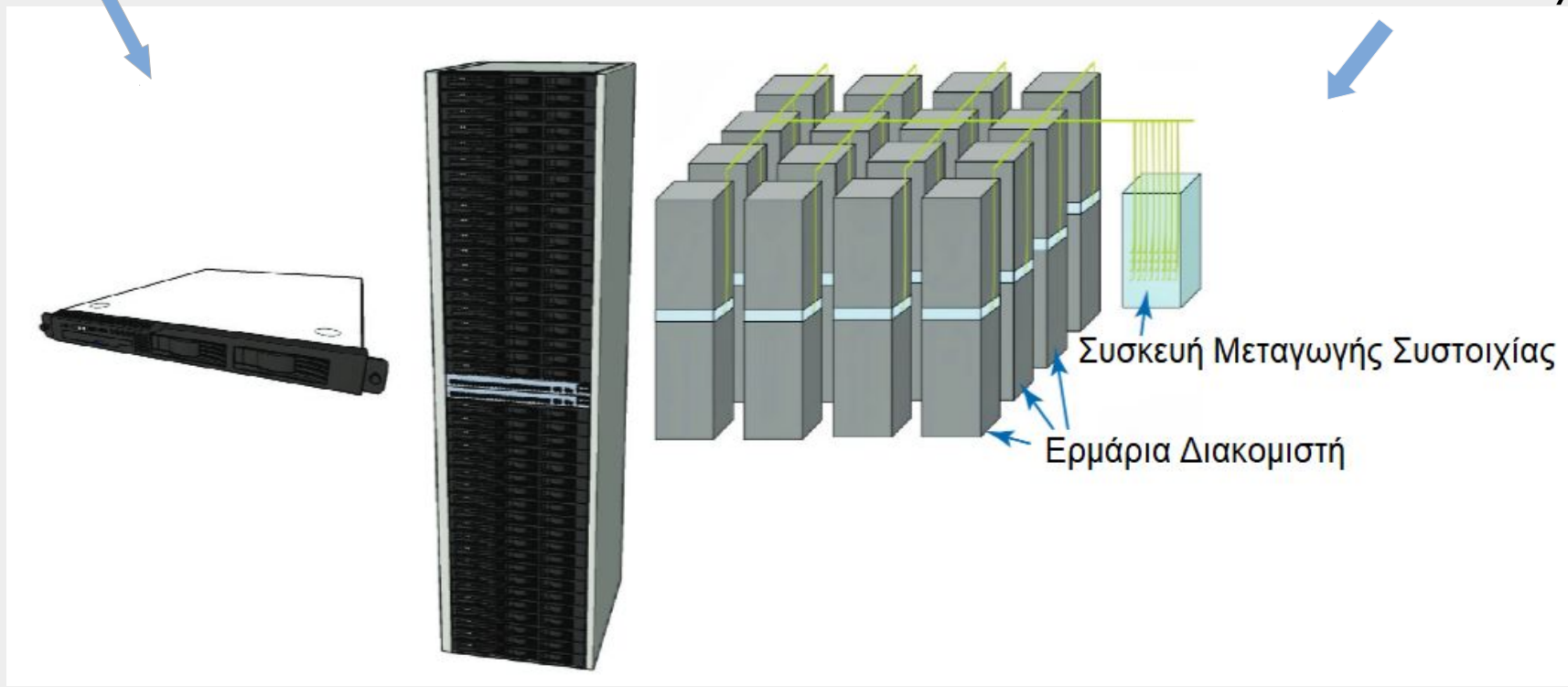
# ΑΡΧΙΤΕΚΤΟΝΙΚΗ ΤΩΝ WSCs (1/2)

Σχέδιο των χαρακτηριστικών στοιχείων σε συστήματα WSCs:

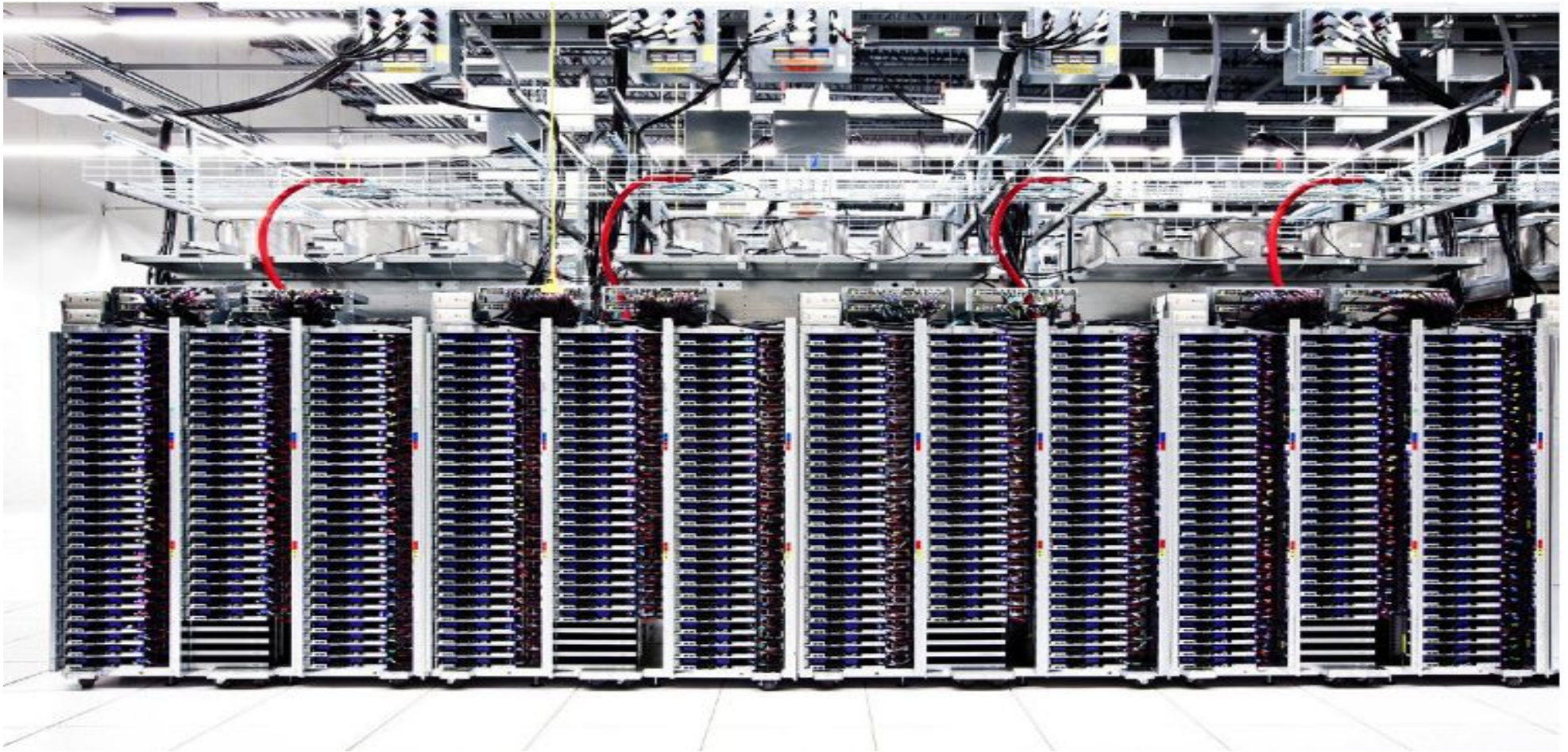
7' ερμάριο με μεταγωγέα Ethernet

1U διακομιστής

διάγραμμα ενός μικρού συμπλέγματος με Ethernet μεταγωγέα/δρομολογητή σε επίπεδο συστοιχίας



## ΑΡΧΙΤΕΚΤΟΝΙΚΗ ΤΩΝ WSCs (2/2)



Εικόνα μιας σειράς διακομιστών σε ένα Google WSC, 2012



## ΑΠΟΘΗΚΕΥΣΗ ΣΕ ΚΕΝΤΡΑ ΔΕΔΟΜΕΝΩΝ (1/2)

- Οι μονάδες δίσκου ή οι συσκευές Flash συνδέονται απευθείας σε κάθε μεμονωμένο διακομιστή και διοικούνται από ένα παγκόσμιο κατακευμαμένο σύστημα αρχείων (όπως το GFS (Google File System) της Google) ή μπορούν να αποτελούν μέρος των συσκευών αποθήκευσης συνδεδεμένων στο δίκτυο (NAS), Network Attached Storage, που είναι άμεσα συνδεδεμένες με τη μεταγωγή σε επίπεδο δομής συστοιχιών.
- Ένα NAS τείνει να είναι μια απλούστερη λύση για την ανάπτυξη, αρχικά επειδή επιτρέπει ορισμένες από τις αρμοδιότητες διαχείρισης δεδομένων να ανατεθεί σε έναν προμηθευτή συσκευών NAS.
- Τα συστήματα που μοιάζουν με GFS είναι σε θέση να διατηρούν τα δεδομένα διαθέσιμα ακόμη και μετά την απώλεια ενός ολόκληρου διακομιστή ή κομμάτι του και μπορεί να επιτρέπουν μεγαλύτερο συνολικό εύρος ζώνης ανάγνωσης επειδή τα ίδια δεδομένα μπορούν να προέρχονται από πολλαπλά αντίγραφα.



## ΑΠΟΘΗΚΕΥΣΗ ΣΕ ΚΕΝΤΡΑ ΔΕΔΟΜΕΝΩΝ (2/2)

- Οι μονάδες Nearline, που δημιουργήθηκαν αρχικά για διακομιστές δημιουργίας αντιγράφων που βασίζονται σε δίσκους, προσθέτουν κοινές ισχυρές λειτουργίες, όπως αυξημένη ανεκτικότητα σε κραδασμούς και είναι κατάλληλες για συνεχή λειτουργία. Τέλος, οι δίσκοι αυτοί προσφέρουν τις υψηλότερες επιδόσεις με το υψηλότερο κόστος.
- Η τεχνολογία NAND Flash έχει κάνει τους SSDs (Solid State Drives) προσιτούς για μια αυξανόμενη τάξη αποθηκευτικών αναγκών σε WSCs. Ενώ το κόστος ανά byte που είναι αποθηκευμένο σε SSD θα παραμείνει πολύ υψηλότερο από ότι στους δίσκους για το προβλέψιμο μέλλον, πολλές υπηρεσίες Web έχουν συντελεστές εισόδου/εξόδου που δεν μπορούν να επιτευχθούν εύκολα με συστήματα δίσκων.





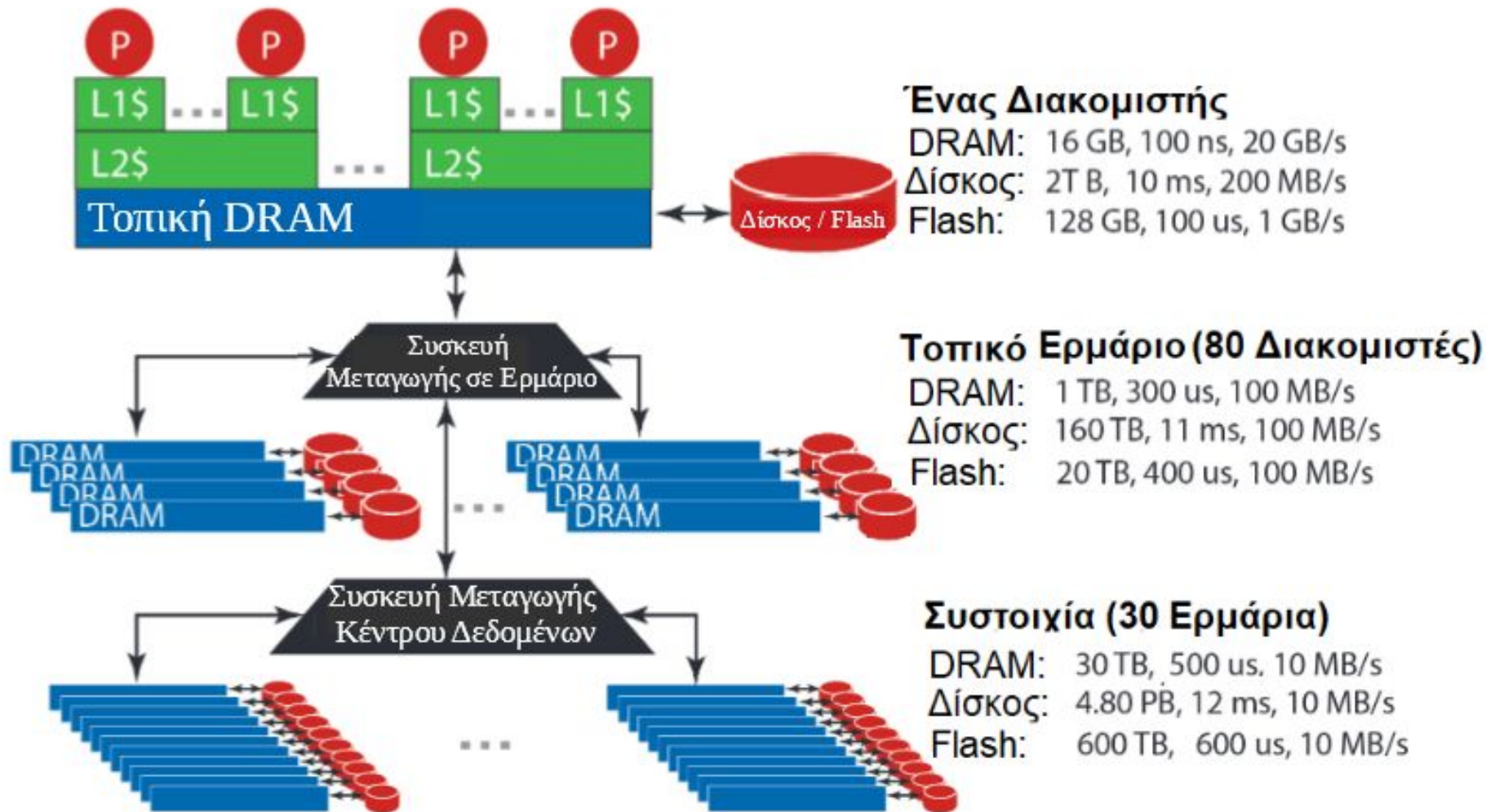
# Η ΔΟΜΗΣΗ ΤΗΣ ΔΙΚΤΥΩΣΗΣ

---

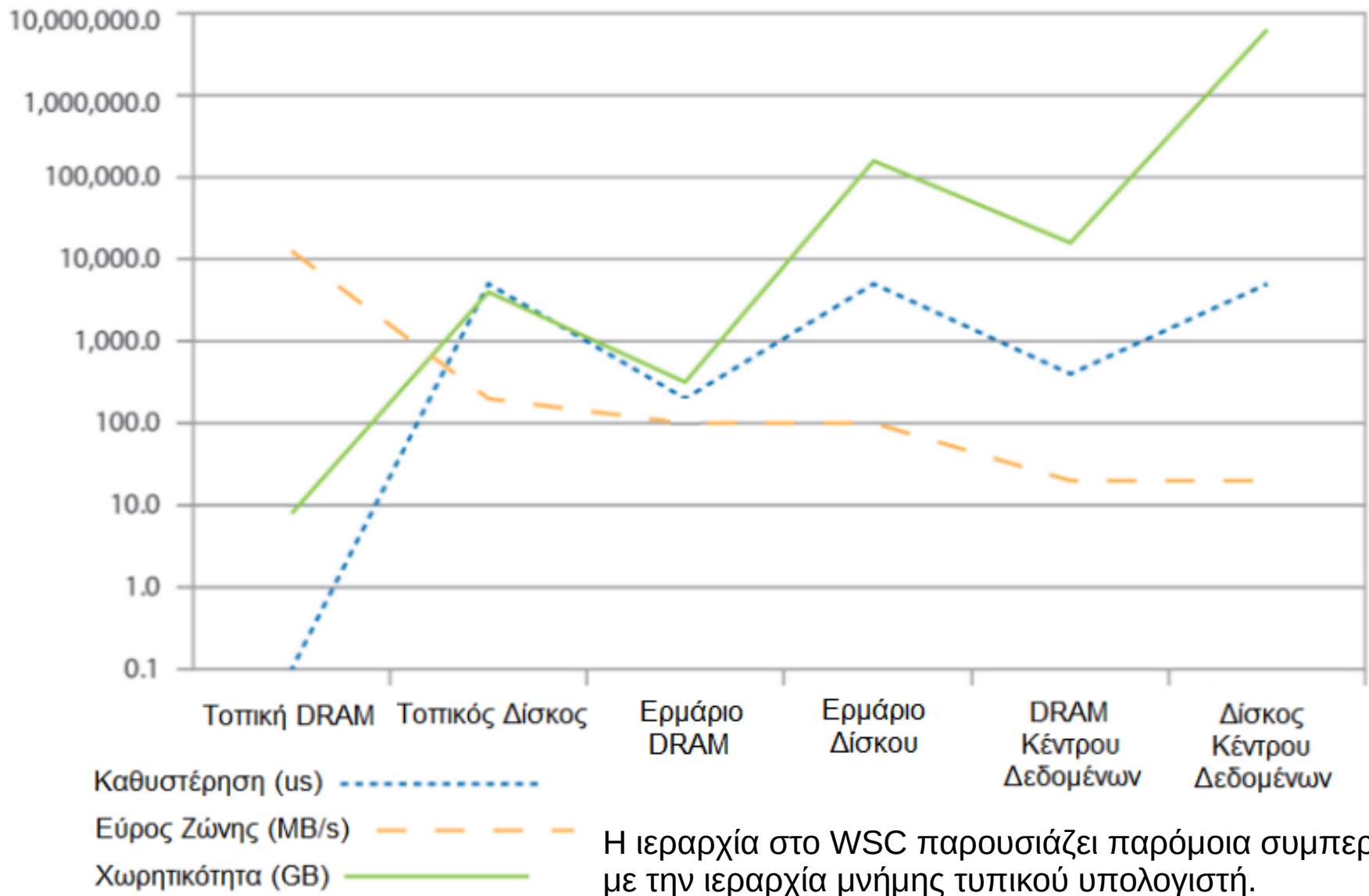
- Η επιλογή δόμησης δικτύωσης για WSCs συνεπάγεται μια εναλλαγή μεταξύ της ταχύτητας, της κλίμακας και του κόστους.
- Το πόσα χρήματα δαπανώνται για τη δικτύωση έναντι της δαπάνης του αντίστοιχου ποσού για την αγορά περισσότερων εξυπηρετητών ή χώρου αποθήκευσης είναι μια ερώτηση που αφορά συγκεκριμένες εφαρμογές και δεν έχει καμία σωστή απάντηση.



# ΙΕΡΑΡΧΙΑ ΑΠΟΘΗΚΕΥΣΗΣ

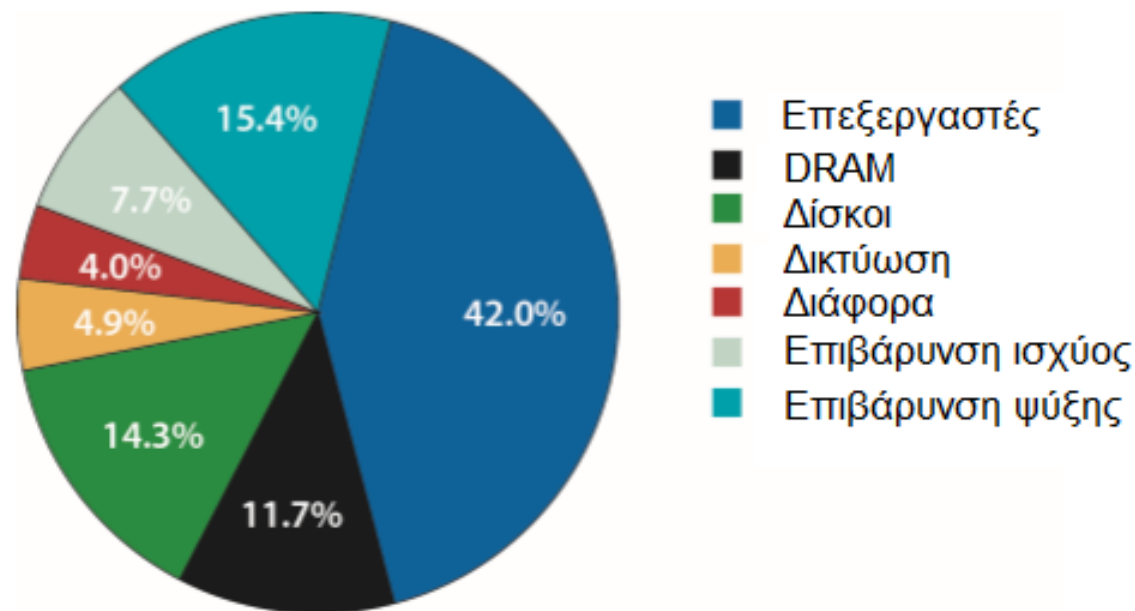


# ΚΑΘΥΣΤΕΡΗΣΗ ΜΕΤΑΦΟΡΑΣ, ΜΕΤΑΦΟΡΑ ΔΕΔΟΜΕΝΩΝ ΚΑΙ ΧΩΡΗΤΙΚΟΤΗΤΑ



# ΧΡΗΣΗ ΕΝΕΡΓΕΙΑΣ ΣΤΑ WSCs

Η κατανάλωση ενέργειας και η ισχύς είναι επίσης σημαντικές ανησυχίες για το σχεδιασμό των WSCs, επειδή το κόστος που σχετίζεται με την ενέργεια έχει καταστεί σημαντικό στοιχείο του συνολικού κόστους αυτής της κατηγορίας συστημάτων.



Η τεράστια κλίμακα των WSCs απαιτεί ότι το λογισμικό υπηρεσιών Internet ανέχεται σχετικά υψηλά ποσοστά σφαλμάτων. Οι μονάδες δίσκων, για παράδειγμα, μπορούν να εμφανίζουν ετήσια ποσοστά αποτυχίας υψηλότερα από 4%.



# ΚΕΝΤΡΟ ΔΕΔΟΜΕΝΩΝ vs ΕΠΙΤΡΑΠΕΖΙΟΣ ΥΠΟΛΟΓΙΣΤΗΣ (1/6)

---

- Άφθονος παραλληλισμός: Οι τυπικές υπηρεσίες Διαδικτύου παρουσιάζουν ένα μεγάλο μέρος παραλληλισμού που απορρέει τόσο από τα δεδομένα όσο και από το επίπεδο των απαιτήσεων.
  - Ο παραλληλισμός των δεδομένων προκύπτει από τα μεγάλα σύνολα δεδομένων από σχετικά ανεξάρτητα αρχεία που χρειάζονται επεξεργασία, όπως συλλογές δισεκατομμυρίων ιστοσελίδων ή δισεκατομμύρια γραμμών κειμένων.



# ΚΕΝΤΡΟ ΔΕΔΟΜΕΝΩΝ vs ΕΠΙΤΡΑΠΕΖΙΟΣ ΥΠΟΛΟΓΙΣΤΗΣ (2/6)

- Ο παραλληλισμός επιπέδου αιτήματος προέρχεται από τις εκατοντάδες ή χιλιάδες αιτήσεις ανά δευτερόλεπτο που λαμβάνουν οι δημοφιλείς υπηρεσίες Internet.
- Οι συναλλαγές μέσω ηλεκτρονικού ταχυδρομείου Web τροποποιούν τα δεδομένα χρήστη. Τα αιτήματα διαφόρων χρηστών είναι ουσιαστικά ανεξάρτητα το ένα από το άλλο, δημιουργώντας φυσικές μονάδες καταμερισμού δεδομένων και ταυτόχρονης ταυτότητας. Εφόσον το ποσοστό ενημέρωσης είναι χαμηλό, ακόμη και τα συστήματα με υψηλά διασυνδεδεμένα δεδομένα (όπως τα backends κοινωνικής δικτύωσης) μπορούν να επωφεληθούν από τον υψηλό παραλληλισμό των αιτημάτων.



# ΚΕΝΤΡΟ ΔΕΔΟΜΕΝΩΝ vs ΕΠΙΤΡΑΠΕΖΙΟΣ ΥΠΟΛΟΓΙΣΤΗΣ (3/6)

- Εναλλαγή όγκου δουλειάς: Οι χρήστες των υπηρεσιών Διαδικτύου απομονώνονται από τις λεπτομέρειες εφαρμογής της υπηρεσίας με σχετικά καλά καθορισμένα και σταθερά API υψηλού επιπέδου (π.χ. απλές διευθύνσεις URL (Uniform Resource Locator) ), καθιστώντας πολύ πιο εύκολη την ταχεία ανάπτυξη του νέου λογισμικού.
  - Το περιβάλλον δημιουργεί σημαντικά κίνητρα για την ταχεία καινοτομία των προϊόντων, αλλά καθιστά δύσκολο για έναν σχεδιαστή συστήματος να εξαγάγει χρήσιμα σημεία αναφοράς ακόμη και από τις καθιερωμένες εφαρμογές.





# ΚΕΝΤΡΟ ΔΕΔΟΜΕΝΩΝ vs ΕΠΙΤΡΑΠΕΖΙΟΣ ΥΠΟΛΟΓΙΣΤΗΣ (4/6)

- Οι υπηρεσίες Διαδικτύου εξακολουθούν να είναι σχετικά νέο πεδίο, εμφανίζονται συχνά νέα προϊόντα και υπηρεσίες και η επιτυχία τους με τους χρήστες επηρεάζει άμεσα τον προκύπτοντα συνδυασμό φόρτου εργασίας στο κέντρο δεδομένων.
- Ομοιογένεια πλατφόρμας: Το κέντρο δεδομένων είναι γενικά ένα πιο ομοιογενές περιβάλλον από την επιφάνεια εργασίας ως πλατφόρμα - στόχο για την ανάπτυξη λογισμικού.
  - Σημαντική ετερογένεια προκύπτει πρωτίστως από τα κίνητρα για την ανάπτυξη αποδοτικότερων από πλευράς κόστους στοιχείων που καθίστανται διαθέσιμα με την πάροδο του χρόνου. Η ομοιογένεια στο πλαίσιο μιας γενιάς πλατφόρμας απλοποιεί τον προγραμματισμό και την εξισορρόπηση φορτίου σε επίπεδο ομάδας και μειώνει το φορτίο συντήρησης λογισμικού πλατφόρμας (πυρήνες, προγράμματα οδήγησης κλπ.).



# ΚΕΝΤΡΟ ΔΕΔΟΜΕΝΩΝ vs ΕΠΙΤΡΑΠΕΖΙΟΣ ΥΠΟΛΟΓΙΣΤΗΣ (5/6)

- Ομοίως, η ομοιογένεια μπορεί να επιτρέψει πιο αποτελεσματικές αλυσίδες εφοδιασμού και αποτελεσματικότερες διαδικασίες επισκευής, διότι οι αυτόματες και μη αυτόματες επισκευές επωφελούνται από την απόκτηση μεγαλύτερης εμπειρίας με λιγότερους τύπους συστημάτων.
  - Το λογισμικό για επιτραπέζια συστήματα μπορεί να κάνει λίγες υποθέσεις σχετικά με την πλατφόρμα υλικού ή λογισμικού στην οποία αναπτύσσονται και η πολυπλοκότητα και τα χαρακτηριστικά απόδοσης τους ενδέχεται να υποφέρουν από την ανάγκη υποστήριξης χιλιάδων ή ακόμα και εκατομμυρίων διαμορφώσεων λογισμικού υλικού και συστήματος.



# ΚΕΝΤΡΟ ΔΕΔΟΜΕΝΩΝ vs ΕΠΙΤΡΑΠΕΖΙΟΣ ΥΠΟΛΟΓΙΣΤΗΣ (6/6)

---

- Λειτουργία χωρίς σφάλματα: Επειδή οι εφαρμογές υπηρεσιών Internet λειτουργούν σε συστοιχίες χιλιάδων μηχανών - καθένα από αυτά δεν είναι δραματικά πιο αξιόπιστο από το υλικό της κλάσης PC.
  - Οι υπηρεσίες Διαδικτύου πρέπει να λειτουργούν σε περιβάλλον όπου τα σφάλματα αποτελούν μέρος της καθημερινής ζωής.



# ΕΡΓΑΛΕΙΑ ΑΠΟΔΟΣΗΣ - ΔΙΑΘΕΣΙΜΟΤΗΤΑΣ

ΤΕΧΝΙΚΗ	ΒΑΣΙΚΟ ΠΛΕΟΝΕΚΤΗΜΑ
Αντικατάσταση	Απόδοση και διαθεσιμότητα
Κωδικοί Reed-Solomon	Διαθεσιμότητα και εξοικονόμηση χώρου
Διαχωρισμός (partitioning)	Απόδοση και διαθεσιμότητα
Ισορροπία φόρτωσης	Απόδοση
Έλεγχος για προβλήματα και για χρονικά περιθώρια	Διαθεσιμότητα
Έλεγχοι ακεραιότητας	Διαθεσιμότητα
Εφαρμογή και συγκεκριμένη συμπίεση	Απόδοση
Ενδεχόμενη συνέπεια	Απόδοση και διαθεσιμότητα
Κεντρικός έλεγχος	Απόδοση
Canaries (Λογισμικό ανίχνευσης επίθεσης)	Διαθεσιμότητα
Εφεδρική εκτέλεση και αντοχή στην ουρά	Απόδοση



- Η βασική εικόνα του συστήματος λογισμικού που εκτελείται σε διακομιστή ενός WSC δεν είναι πολύ διαφορετική από αυτή που θα περίμενε κανείς σε μια κανονική πλατφόρμα διακομιστή.
- Δεδομένου του υψηλού βαθμού ομοιογένειας στις διαμορφώσεις υλικού των διακομιστών WSC, μπορούμε να βελτιώσουμε την ανάπτυξη και τον έλεγχο του firmware και του προγράμματος οδήγησης συσκευών, καθώς θα υπάρχουν λιγότεροι συνδυασμοί συσκευών.



- Η πλειοψηφία των συνδέσεων δικτύωσης από ένα διακομιστή WSC θα είναι σε άλλες μηχανές εντός του ίδιου κτιρίου και θα έχει χαμηλότερες απώλειες πακέτων σε σχέση με τις συνδέσεις Internet μεγάλων αποστάσεων.
- Ένα εικονικό μηχάνημα VM (Virtual Machine) παρέχει μια συνοπτική και φορητή επαφή για τη διαχείριση τόσο της ασφάλειας και απόδοσης της εφαρμογής ενός πελάτη όσο και για την ταυτόχρονη συνύπαρξη πολλών λειτουργικών συστημάτων με περιορισμένη επιπλέον πολυπλοκότητα.
- Η απλότητα της εγκατάστασης VM καθιστά επίσης ευκολότερη την υλοποίηση της ζωντανής «μετακίνησης».



- Ελέγχει την τοποθέτηση των εργασιών του χρήστη σε πόρους - πηγές, επιβάλλει προτεραιότητες και παρέχει βασικές υπηρεσίες διαχείρισης εργασιών.
- Καταγραφή Υλικού και άλλες βασικές υπηρεσίες
  - Παραδείγματα αποτελούν η αξιόπιστη κατανεμημένη αποθήκευση, η μετάδοση μηνυμάτων και ο συγχρονισμός σε επίπεδο συστοιχίας.
- Ανάπτυξη και συντήρηση
- Όρια προγραμματισμού



- Η αναζήτηση στον παγκόσμιο ιστό ήταν μια από τις πρώτες υπηρεσίες διαδικτύου μεγάλης κλίμακας και η οργάνωση αυτού του τεράστιου όγκου πληροφοριών αρχίζει σιγά σιγά να δημιουργεί προβλήματα.
- Ωστόσο, καθώς η συνδεσιμότητα δικτύων με τα σπίτια και τις επιχειρήσεις συνεχίζει να βελτιώνεται, γίνεται πιο ελκυστική η προσφορά νέων υπηρεσιών μέσω του Διαδικτύου.
- Οι χάρτες και οι υπηρεσίες ηλεκτρονικού ταχυδρομείου είναι στοιχεία της τάσης που επικρατεί.





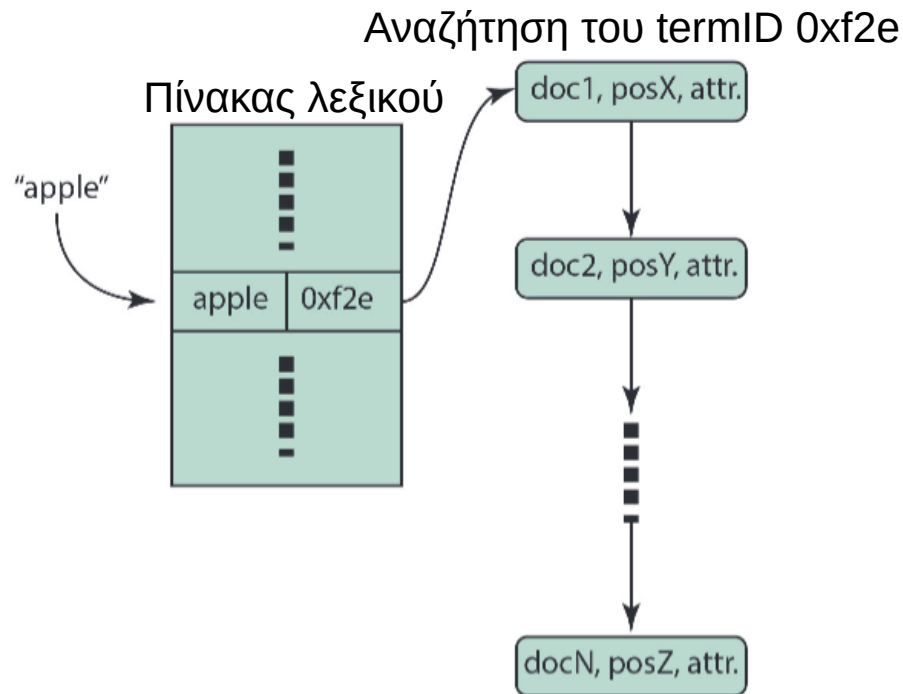
- Μόλις ληφθούν υπόψη οι ποικίλες απαιτήσεις πολλαπλών υπηρεσιών, καθίσταται σαφές ότι το κέντρο δεδομένων πρέπει να είναι ένα σύστημα πληροφορικής γενικής χρήσης.
- Παρόλο που οι εξειδικευμένες λύσεις υλικού μπορεί να είναι κατάλληλες για μεμονωμένα τμήματα υπηρεσιών, το εύρος των απαιτήσεων καθιστά λιγότερο πιθανό το εξειδικευμένο υλικό να έχει μεγάλη συνολική επίδραση στη λειτουργία.
- Ένας παράγοντας εναντίον της εξειδίκευσης του υλικού είναι η ταχύτητα φόρτωσης του φορτίου.



# ΠΑΡΑΔΕΙΓΜΑ ΟΓΚΟΥ ΕΡΓΑΣΙΑΣ (1/2)

ONLINE: Αναζήτηση στο διαδίκτυο

Παρακάτω φαίνεται το λογικό διάγραμμα ενός ευρετηρίου στο διαδίκτυο.



## ΠΑΡΑΔΕΙΓΜΑ ΟΓΚΟΥ ΕΡΓΑΣΙΑΣ (2/2)

- ONLINE: Αναζήτηση στο διαδίκτυο
  - Μια δομή αρχείων συνδέει μια ID σε κάθε όρο στον αποθηκευτικό χώρο.
  - Ο αλγόριθμος αναζήτησης πρέπει να διασχίσει τις δημοσιευμένες λίστες για κάθε όρο μέχρι να βρει όλα τα έγγραφα που περιέχονται σε όλες τις λίστες καταχώρισης.
  - Δεδομένου του τεράστιου μεγέθους του δείκτη, αυτός ο αλγόριθμος αναζήτησης μπορεί να τρέξει σε μερικές χιλιάδες μηχανές.
  - Η υψηλή απόδοση είναι επίσης μια βασική μέτρηση απόδοσης, επειδή μια δημοφιλής υπηρεσία μπορεί να χρειαστεί να υποστηρίξει πολλές χιλιάδες ερωτήματα ανά δευτερόλεπτο.
  - Τέλος, επειδή η αναζήτηση στο Web είναι μια ηλεκτρονική υπηρεσία, υποφέρει από κανονικές διακυμάνσεις της επισκεψιμότητας.



- Οι διαχειριστές συστημάτων πρέπει να παρακολουθούν πόσο καλά μια υπηρεσία Διαδικτύου πληροί το επίπεδο εξυπηρέτησής της.
- Οι υπηρεσίες μεγάλης κλίμακας συχνά χρειάζονται πιο εξελιγμένη και κλιμακούμενη υποστήριξη παρακολούθησης, καθώς ο αριθμός των πρωτοκόλλων μπορεί να είναι αρκετά μεγάλος και χρειάζονται περισσότερα σήματα για να χαρακτηρίσουν την «υγεία» της υπηρεσίας.



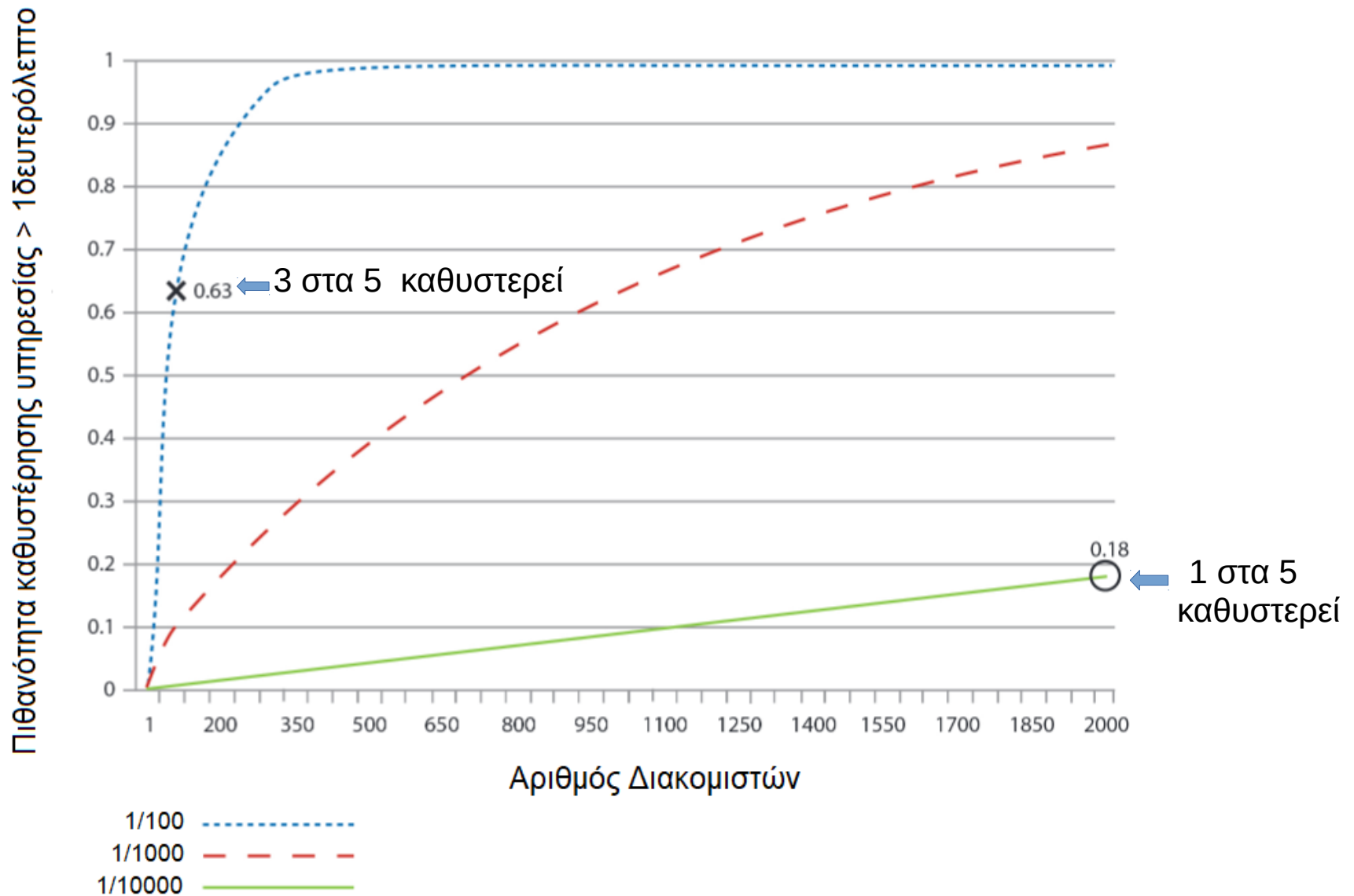
- Αυτά τα εργαλεία επιχειρούν να προσδιορίσουν όλη την εργασία που γίνεται σε ένα κατανεμημένο σύστημα για λογαριασμό ενός συγκεκριμένου εκκινήτη (όπως ένα αίτημα χρήστη) και να αναλύσουν τις αιτίες ή τις προσωρινές σχέσεις μεταξύ των διαφόρων εμπλεκόμενων στοιχείων.
- Αυτά τα εργαλεία χωρίζονται σε δύο μεγάλες κατηγορίες: τα συστήματα παρακολούθησης black box και τα συστήματα εφαρμογών.
- Οι CPUs (Central Processing Unit) που βασίζονται σε δείγματα των μετρητών απόδοσης Υλικού έχουν αποδειχθεί εξαιρετικά επιτυχείς για να βοηθήσουν τους προγραμματιστές να κατανοήσουν τα φαινόμενα απόδοσης των μικροκατασκευών.



- Τα καταναμεημένα εργαλεία εντοπισμού του συστήματος και οι πίνακες ελέγχου του επιπέδου εφαρμογών μετράνε την ομαλή λειτουργία και την απόδοση των εφαρμογών. Αυτά τα εργαλεία μπορούν να συμπεράνουν ότι μια συνιστώσα υλικού μπορεί να είναι παρεκτρέπουσα, αλλά και αυτή είναι ακόμα μια εκτίμηση.
- Όλο το σύνολο της στοίβας του λογισμικού βρίσκεται υπό τον έλεγχο του προγραμματιστή.

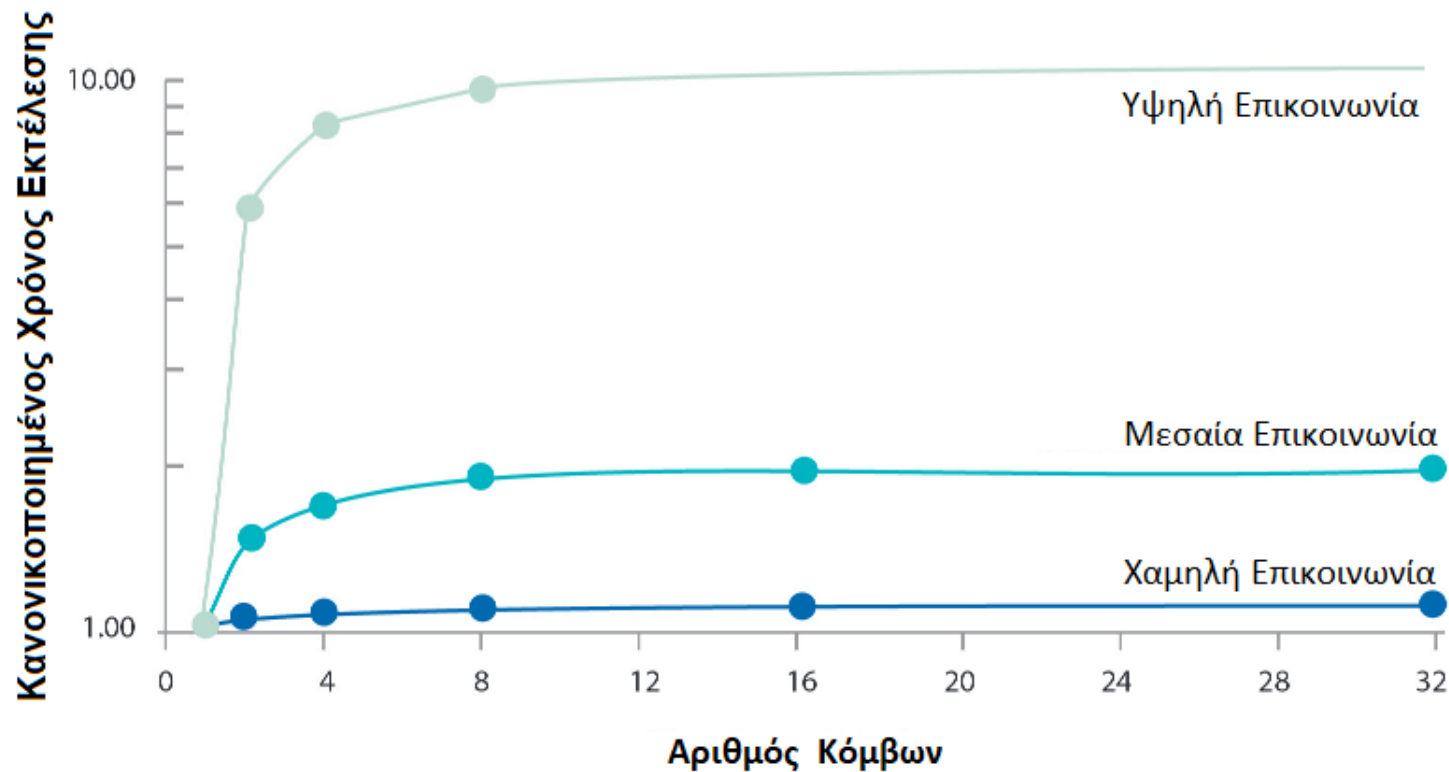


# ΕΠΙΡΡΟΗ ΤΗΣ ΚΑΘΥΣΤΕΡΗΣΗΣ ΥΠΗΡΕΣΙΑΣ ΓΙΑ ΔΙΑΦΟΡΑ ΣΕΝΑΡΙΑ



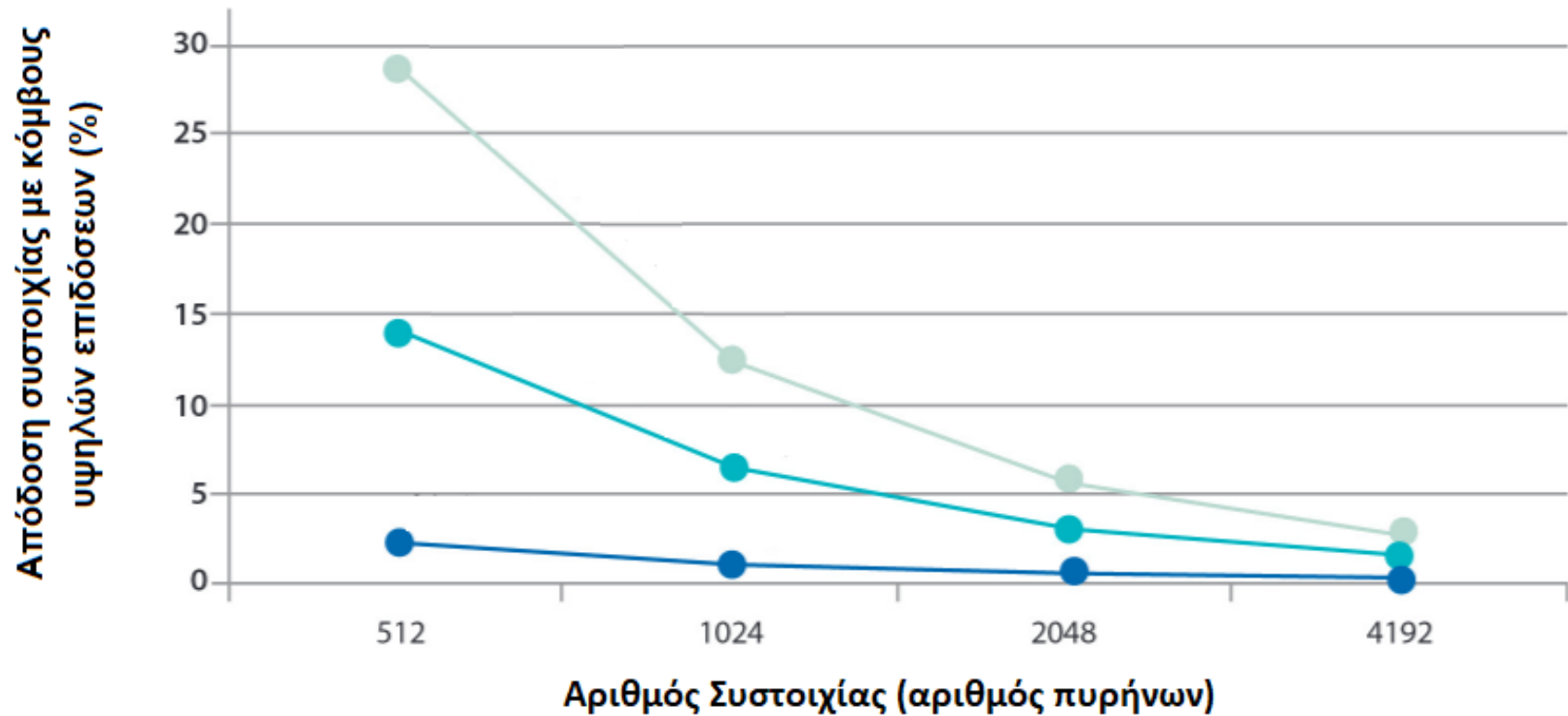
# ΑΠΟΤΕΛΕΣΜΑΤΙΚΟΤΗΤΑ ΤΟΥ ΥΛΙΚΟΥ ΤΩΝ ΔΙΑΚΟΜΙΣΤΩΝ

- Οι επιπτώσεις στο χρόνο εκτέλεσης από ένα μεγάλο συμμετρικό σύστημα πολλαπλών επεξεργαστών (SMP) στην απόδοση της επικοινωνίας:
  - Χρόνος εκτέλεσης =  $1 \text{ ms} + f * [100 \text{ ns} / \# \text{ κόμβοι} + 100 \text{ μs} * (1 - 1 / \# \text{ κόμβοι})]$





# ΑΠΟΔΟΣΗ ΜΕ ΒΑΣΗ ΤΟΝ ΑΡΙΘΜΟ ΠΥΡΗΝΩΝ



Υψηλή Επικοινωνία  
Μεσαία Επικοινωνία  
Χαμηλή Επικοινωνία



- Τα πλεονεκτήματα της χρήσης μικρότερων και βραδύτερων επεξεργαστών είναι παρόμοια με τα επιχειρήματα για τη χρήση διακομιστών βασικών προϊόντων μεσαίας εμβέλειας αντί για υψηλών επιδόσεων SMPs.
  - Οι μεγάλες CPUs σε διακομιστές μεσαίας εμβέλειας φέρουν συνήθως κάποιο κόστος/απόδοση διαφορετικό από επεξεργαστές χαμηλότερου επιπέδου.
  - Πολλές εφαρμογές είναι περιορισμένες σε μνήμη, με αποτέλεσμα πιο γρήγορες CPUs να μην αποδίδουν για μεγάλες εφαρμογές.
  - Οι πιο αργές CPUs είναι πιο αποδοτικές από πλευράς ισχύος.



- Οι αρχιτέκτονες υπολογιστών εκπαιδεύονται για να λύσουν το πρόβλημα της εξεύρεσης του σωστού συνδυασμού απόδοσης και χωρητικότητας από τα διάφορα μέρη που συνθέτουν ένα WSC.
  - Οι έξυπνοι προγραμματιστές μπορεί να είναι σε θέση να αναδιαρθρώσουν τους αλγόριθμους τους για να ταιριάζουν καλύτερα με μια πιο οικονομική εναλλακτική λύση.



- Η πιο οικονομικά αποδοτική και ισορροπημένη διαμόρφωση για το Υλικό μπορεί να ταιριάζει με τις συνδυασμένες απαιτήσεις πολλαπλού όγκου εργασίας και όχι απαραίτητα από μια μόνο συγκεκριμένη εργασία.
- Οι μόνιμες πηγές τείνουν να χρησιμοποιούνται περισσότερο.
- Το σωστό σημείο σχεδιασμού εξαρτάται περισσότερο από τη δομή υψηλού επιπέδου του ίδιου του όγκου εργασίας, διότι το μέγεθος των δεδομένων και η δημοτικότητα των υπηρεσιών διαδραματίζουν σημαντικό ρόλο.



# ΑΠΟΘΗΚΕΥΣΗ ΤΩΝ WSCs (1/3)

- Ανοργάνωτη αποθήκευση
  - Το GFS της Google είναι ένα παράδειγμα ενός συστήματος αποθήκευσης με απλή αφαίρεση αρχείων (το σύστημα Colossus της Google έχει αντικαταστήσει το GFS, αλλά ακολουθεί μια παρόμοια αρχιτεκτονική φιλοσοφία και έτσι επιλέγουμε να περιγράψουμε το πιο γνωστό GFS).
  - Το GFS σχεδιάστηκε για να υποστηρίζει το σύστημα ευρετηρίου αναζήτησης ιστού και ως εκ τούτου επικεντρώνεται στην υψηλή απόδοση για χιλιάδες ταυτόχρονους αναγνώστες/συγγραφείς και στην ισχυρή απόδοση σε σχέση με τα υψηλά ποσοστά αποτυχίας του Υλικού. Οι χρήστες του GFS συνήθως χειρίζονται μεγάλες ποσότητες δεδομένων και έτσι το GFS βελτιστοποιείται περαιτέρω ώστε να βοηθούν ακόμα και για μεγάλες επιχειρήσεις.
  - Παρόλο που η αρχική έκδοση του GFS υποστηρίζει μόνο την απλή αντιγραφή, η σημερινή έκδοση (Colossus) έχει προσθέσει υποστήριξη για κώδικες Reed-Solomon (πιο αποδοτικούς από άποψη χώρου).



- Δομημένη – Οργανωμένη Αποθήκευση
  - Η απλή αφαίρεση αρχείων των GFS και Colossus μπορεί να επαρκεί για συστήματα που χειρίζονται μεγάλες ποσότητες δεδομένων, αλλά οι προγραμματιστές εφαρμογών χρειάζονται επίσης το WSC ως ισοδύναμο των λειτουργιών μιας βάσης δεδομένων, όπου τα σύνολα δεδομένων μπορούν να δομηθούν και να χρησιμοποιηθούν για εύκολες μικρές ενημερώσεις ή σύνθετα ερωτήματα.
  - Σε σύγκριση με τα παραδοσιακά συστήματα βάσεων δεδομένων, τα BigTable και Dynamo θυσιάζουν ορισμένα χαρακτηριστικά, όπως τη δημιουργία σχήματος και την ισχυρή συνοχή, υπέρ της υψηλότερης απόδοσης και της διαθεσιμότητας σε τεράστιες κλίμακες.



## ΑΠΟΘΗΚΕΥΣΗ ΤΩΝ WSCs (3/3)

- Η επιλογή της σταθερότητας στο BigTable και το Dynamo μετατοπίζει το βάρος της επίλυσης προσωρινών σφαλμάτων στις εφαρμογές που χρησιμοποιούν αυτά τα συστήματα.
  - Megastore και Spanner: Είναι συστήματα υψηλών απαιτήσεων και υψηλών αποδόσεων.
- Αλληλεπίδραση του χώρου αποθήκευσης και της τεχνολογίας δικτύου
    - Το χάσμα μεταξύ του δικτύου και της απόδοσης του δίσκου έχει διευρυνθεί στο σημείο στο οποίο η τοποθεσία του δίσκου δεν είναι πλέον σχετική σε υπολογισμούς εντός του κέντρου δεδομένων.



## ΔΙΚΤΥΟ WSC (1/5)

- Οι διακομιστές πρέπει να είναι συνδεδεμένοι και καθώς οι επιδόσεις των διακομιστών αυξάνονται με την πάροδο του χρόνου, η ζήτηση για εύρος ζώνης μεταξύ διακομιστών αυξάνεται επίσης.
- Όμως, ενώ μπορούμε να διπλασιάσουμε τη συνολική υπολογιστική χωρητικότητα ή να διπλασιάσουμε τη συσσωρευμένη αποθήκευση απλά διπλασιάζοντας τον αριθμό των υπολογιστικών στοιχείων ή στοιχείων αποθήκευσης, η δικτύωση δεν έχει απλή λύση οριζόντιας κλίμακας.



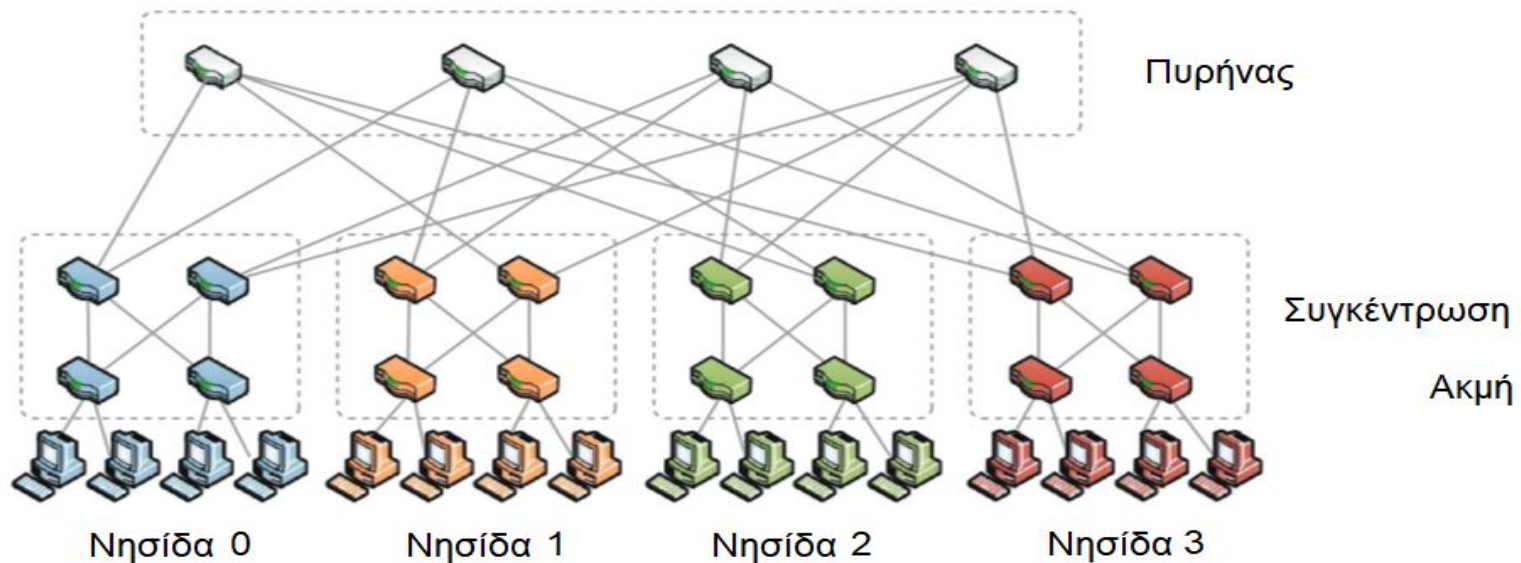


- Ο διπλασιασμός του εύρους ζώνης φύλλων είναι εύκολος - με δύο φορές περισσότερους διακομιστές θα έχουμε διπλάσιες θύρες δικτύου και συνεπώς διπλάσιο εύρος ζώνης.
- Αλλά αν υποθέσουμε ότι κάθε διακομιστής θέλει να μιλήσει σε κάθε άλλο διακομιστή, πρέπει να διπλασιάσουμε όχι μόνο το εύρος ζώνης του διαμερίσματος, αλλά το εύρος ζώνης διχοτόμησης.



## ΔΙΚΤΥΟ WSC (3/5)

- Δυστυχώς, ο διπλασιασμός του εύρους ζώνης διχοτόμησης είναι δύσκολος επειδή δεν μπορούμε απλώς να αγοράσουμε (ή να δημιουργήσουμε) ένα μεγαλύτερο μεταγωγέα λειτουργίας.
- Σε αυτό το δέντρο, μπορεί να διανεμηθεί η ίδια ποσότητα και από έναν μεμονωμένο μεταγωγέα.



## ΔΙΚΤΥΟ WSC (4/5)

---

- Για να αποφεύγεται η πληρωμή χιλιάδων δολαρίων σε κόστος δικτύωσης ανά μηχανή, οι υπεύθυνοι υλοποίησης WSC συχνά μειώνουν το κόστος με την υπερκάλυψη του δικτύου με ένα μεταγωγέα top-of-rack.
- Ένας άλλος τρόπος αντιμετώπισης της επεκτασιμότητας του δικτύου είναι η εκφόρτωση κάποιας κυκλοφορίας σε δίκτυο ειδικού σκοπού.
- Τα WSCs που χρησιμοποιούνται, δημιουργούν νέες προκλήσεις στα δίκτυα, αφού τα τελικά σημεία της σύνδεσης (δηλαδή οι διευθύνσεις IP (Internet Protocol) /συνδυασμοί θυρών) μπορούν να μετακινηθούν από μία φυσική μηχανή σε άλλη χωρίς προβλήματα.



## ΔΙΚΤΥΟ WSC (5/5)

---

- Τα WSCs που χρησιμοποιούν VM δημιουργούν νέες προκλήσεις στα δίκτυα, αφού τα τελικά σημεία της σύνδεσης (δηλαδή οι διευθύνσεις IP/συνδυασμοί θυρών) μπορούν να μετακινηθούν από μία φυσική μηχανή σε άλλη.
- Επιπλέον, οι διακομιστές είναι πιο εύκολο να προγραμματιστούν, προσφέροντας πλουσιότερα περιβάλλοντα προγραμματισμού και πολύ πιο ισχυρό εξοπλισμό - από το 2013, όπου ένας τυπικός επεξεργαστής ελέγχου δρομολογητή αποτελείται από επεξεργαστή μονού πυρήνα 2 GHz με 4 GB μνήμης RAM (Random Access Memory).



- Βαθμίδες και προδιαγραφές ενός κέντρου δεδομένων
  - Ο σχεδιασμός ενός κέντρου δεδομένων συχνά ταξινομείται ως μέλος της κατηγορίας “Επίπεδο I-IV”.
  - Επίπεδο I: Μια ενιαία διαδρομή για κατανομή ισχύος, UPS (Uninterruptible Power Supply) και διανομή ψύξης, χωρίς περιττές συνιστώσες.
  - Επίπεδο II: Προσθέτει πλεονάζοντα εξαρτήματα σε αυτό το σχέδιο (N+1), βελτιώνοντας τη διαθεσιμότητα.
  - Επίπεδο III: Διαθέτουν μία ενεργή και μία εναλλακτική διαδρομή διανομής σε βοηθητικά προγράμματα.
  - Επίπεδο IV: Διαθέτουν δύο ταυτόχρονα ενεργές διαδρομές διανομής ενέργειας και ψύξης.



# ΤΑ ΚΥΡΙΑ ΣΥΣΤΑΤΙΚΑ ΕΝΟΣ ΤΥΠΙΚΟΥ ΚΕΝΤΡΟΥ ΔΕΔΟΜΕΝΩΝ



- Περιέχει έναν μεταγωγέα μεταφοράς που επιλέγει την ενεργή πηγή εισόδου (είτε την πηγή που παρέχεται από το δίκτυο είτε τη γεννήτρια).
- Περιέχει κάποια μορφή αποθήκευσης ενέργειας (ηλεκτρικής ή μηχανικής) για να γεφυρώσει τη βλάβη του δικτύου.
- Ρυθμίζει την εισερχόμενη τροφοδοσία και τις παραμορφώσεις στην τροφοδοσία εναλλασσόμενου ρεύματος.
- Είναι δυνατή η χρήση συστημάτων UPS όχι μόνο σε διακοπές ρεύματος, αλλά και ως συμπληρωματική ρύθμιση για την διαχείριση ενέργειας.



# Ο ΡΟΛΟΣ ΤΩΝ PDUs

---

- Η έξοδος UPS δρομολογείται σε PDUs (Power Distribution Units) στην αρχή, στο πάτωμα του κέντρου δεδομένων.
- Τα PDU μοιάζουν με ηλεκτρολογικούς πίνακες σε κατοικίες, αλλά μπορούν επίσης να ενσωματώνουν μετασχηματιστές για τελικές ρυθμίσεις τάσης.
- Παίρνουν μια μεγαλύτερη τροφοδοσία εισόδου και την διασπούν σε πολλά μικρότερα κυκλώματα που διανέμουν ισχύ στους πραγματικούς διακομιστές στο «πάτωμα» του κέντρου δεδομένων.





## Ο ΡΟΛΟΣ ΤΩΝ PDUs

- Στη Βόρεια Αμερική, η είσοδος στην PDU είναι τυπικά τροφοδοσία ισχύος τριών φάσεων τάσης 480 V. Αυτό απαιτεί από το PDU να εκτελέσει μετασχηματισμό για να παραδώσει την επιθυμητή έξοδο 110 V για τους διακομιστές.
- Στην ΕΕ, η είσοδος στην PDU είναι τυπικά τροφοδοτική ισχύς 400 V.
- Με κατάλληλες μετατροπές δίνεται το επιθυμητό 230 V.
- Χρησιμοποιώντας τις ίδιες τεχνικές στη Βόρειο Αμερική απαιτείται ο εξοπλισμός υπολογιστών να δέχεται 277 V, η οποία δυστυχώς υπερβαίνει την ανώτερη γκάμα τυποποιημένων τροφοδοτικών.



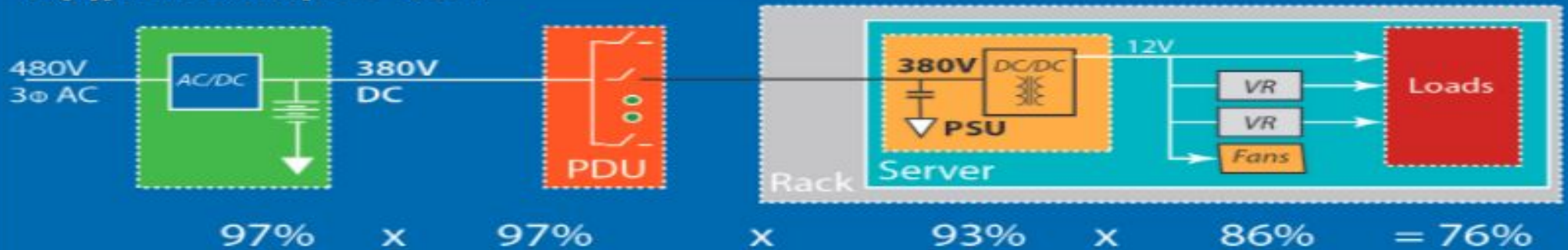
# ΜΕΤΑΦΟΡΑ ΕΝΕΡΓΕΙΑΣ ΤΟΥ ΚΕΝΤΡΟΥ ΔΕΔΟΜΕΝΩΝ (1/3)

AC Vs DC ΑΠΟΔΟΤΙΚΟΤΗΤΑ ΜΕΤΑΦΟΡΑΣ: 76% απόδοση μετατροπής DC 380 V προς 68% απόδοση μετατροπής συμβατικής αρχιτεκτονικής AC

Βελτίωση της απόδοσης της μεταφοράς ενέργειας  
Συμβατική αρχιτεκτονική AC



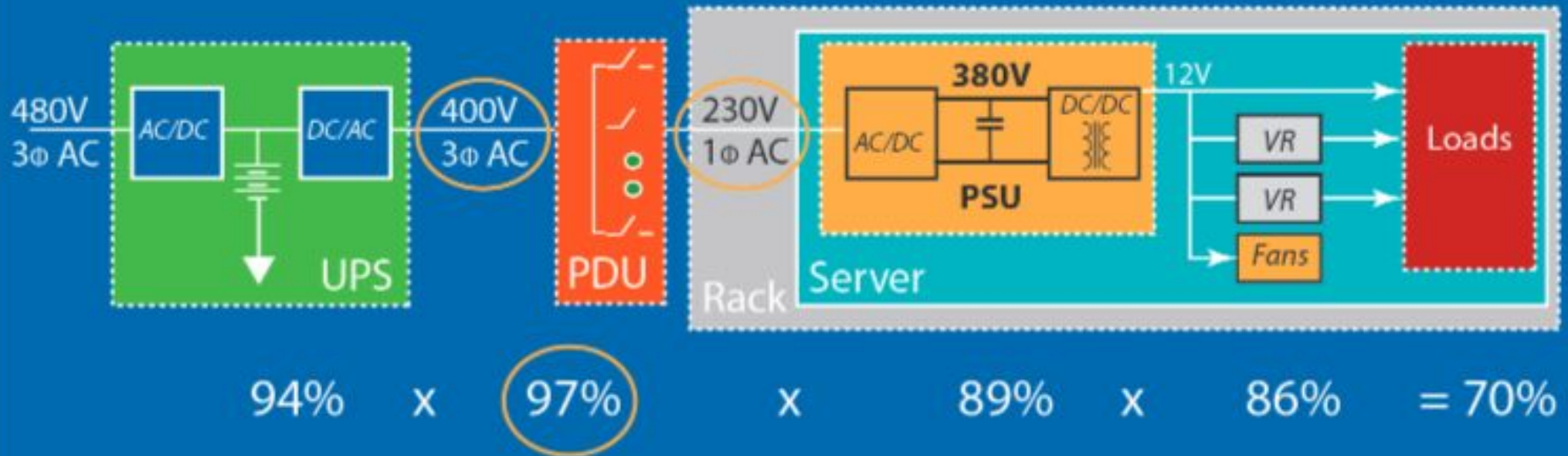
Σύγκριση της αρχιτεκτονικής μεταφοράς ενέργειας  
Αρχιτεκτονική DC 380V



# ΜΕΤΑΦΟΡΑ ΕΝΕΡΓΕΙΑΣ ΤΟΥ ΚΕΝΤΡΟΥ ΔΕΔΟΜΕΝΩΝ (2/3)

Αποτελεσματικότητα μεταφοράς 400V AC: **70%** απόδοση μετατροπής DC 400 V

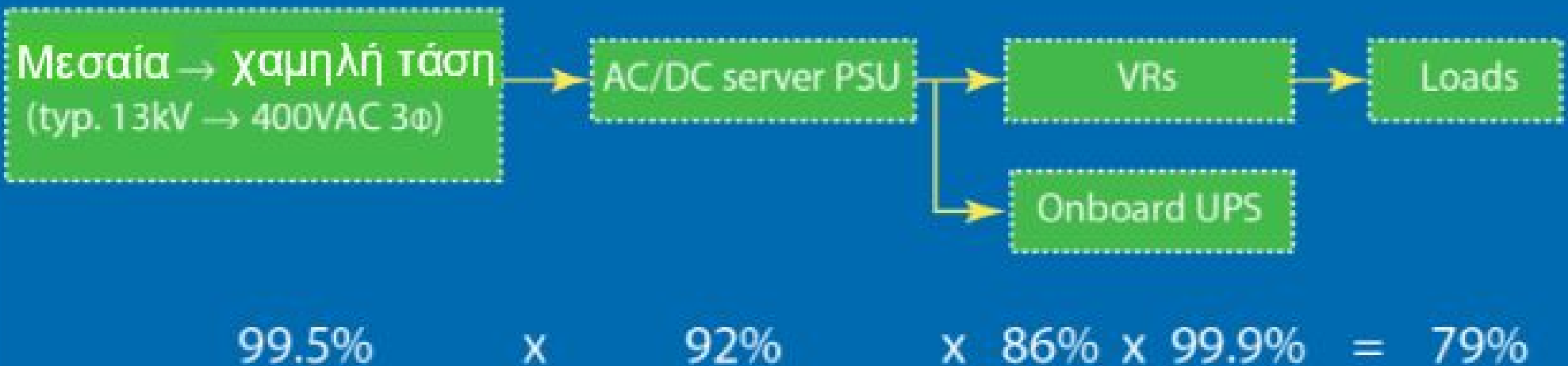
## Σύγκριση της αρχιτεκτονικής μεταφοράς ενέργειας Αρχιτεκτονική DC 400V



# ΜΕΤΑΦΟΡΑ ΕΝΕΡΓΕΙΑΣ ΤΟΥ ΚΕΝΤΡΟΥ ΔΕΔΟΜΕΝΩΝ (3/3)

Διάγραμμα απόδοσης της διανομής AC της Google για κέντρα δεδομένων: **Η καλύτερη απόδοση αρχιτεκτονικής είναι η 79% απόδοση της αρχιτεκτονικής Google AC**

## Αρχιτεκτονική GOOGLE AC



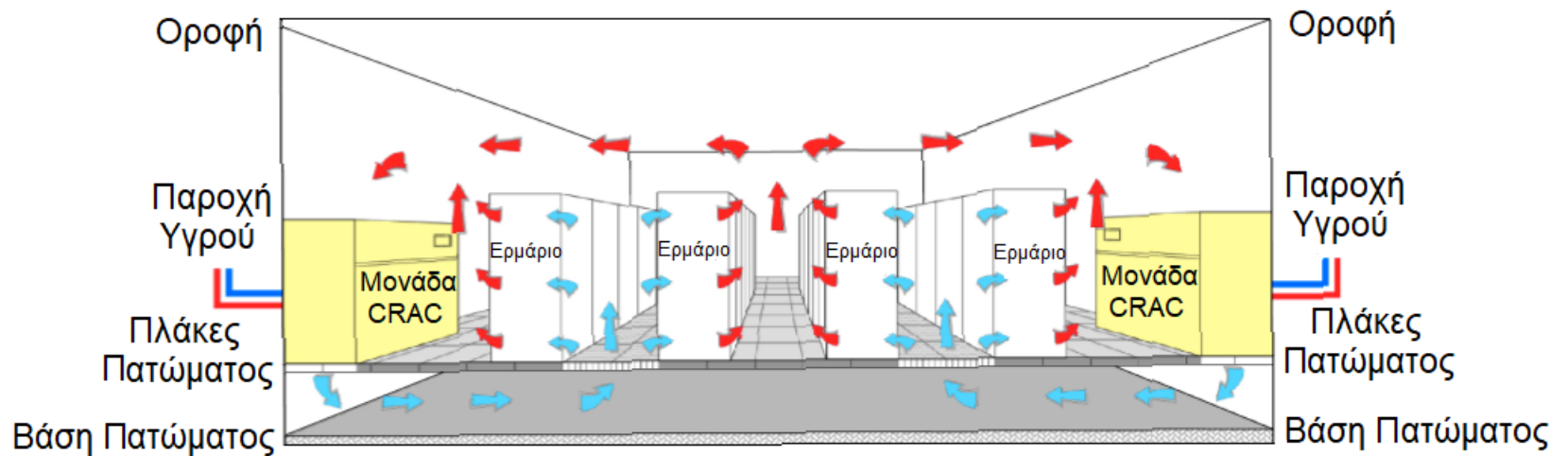
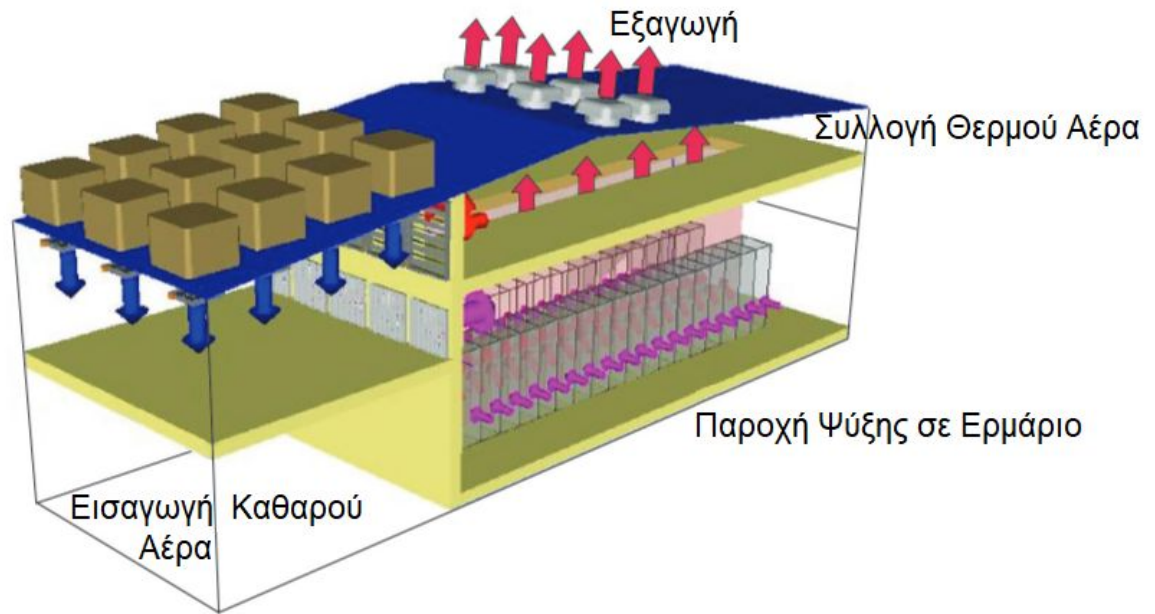
- Τα συστήματα ψύξης κέντρων δεδομένων αφαιρούν τη θερμότητα που παράγεται από όλο τον εξοπλισμό.
- Για να απομακρυνθεί η θερμότητα, ένα σύστημα ψύξης πρέπει να χρησιμοποιήσει κάποια ιεραρχία συστημάτων βρόχου, καθένα από τα οποία φέρνει ένα κρύο μέσο που θερμαίνεται μέσω κάποιας μορφής ανταλλαγής θερμότητας.
- Ένα σύστημα κλειστού βρόχου επανακυκλοφορεί το ίδιο μέσο ξανά και ξανά, μεταφέροντας τη θερμότητα σε γειτονικό άνω βρόχο, σε εναλλάκτη θερμότητας και τελικά στο περιβάλλον.
- Όλα τα συστήματα πρέπει να περιέχουν βρόχο στο εξωτερικό περιβάλλον για απόλυτη απόρριψη της θερμότητας.



# ΣΥΣΤΗΜΑΤΑ ΨΥΞΗΣ ΚΕΝΤΡΩΝ ΔΕΔΟΜΕΝΩΝ (2/3)

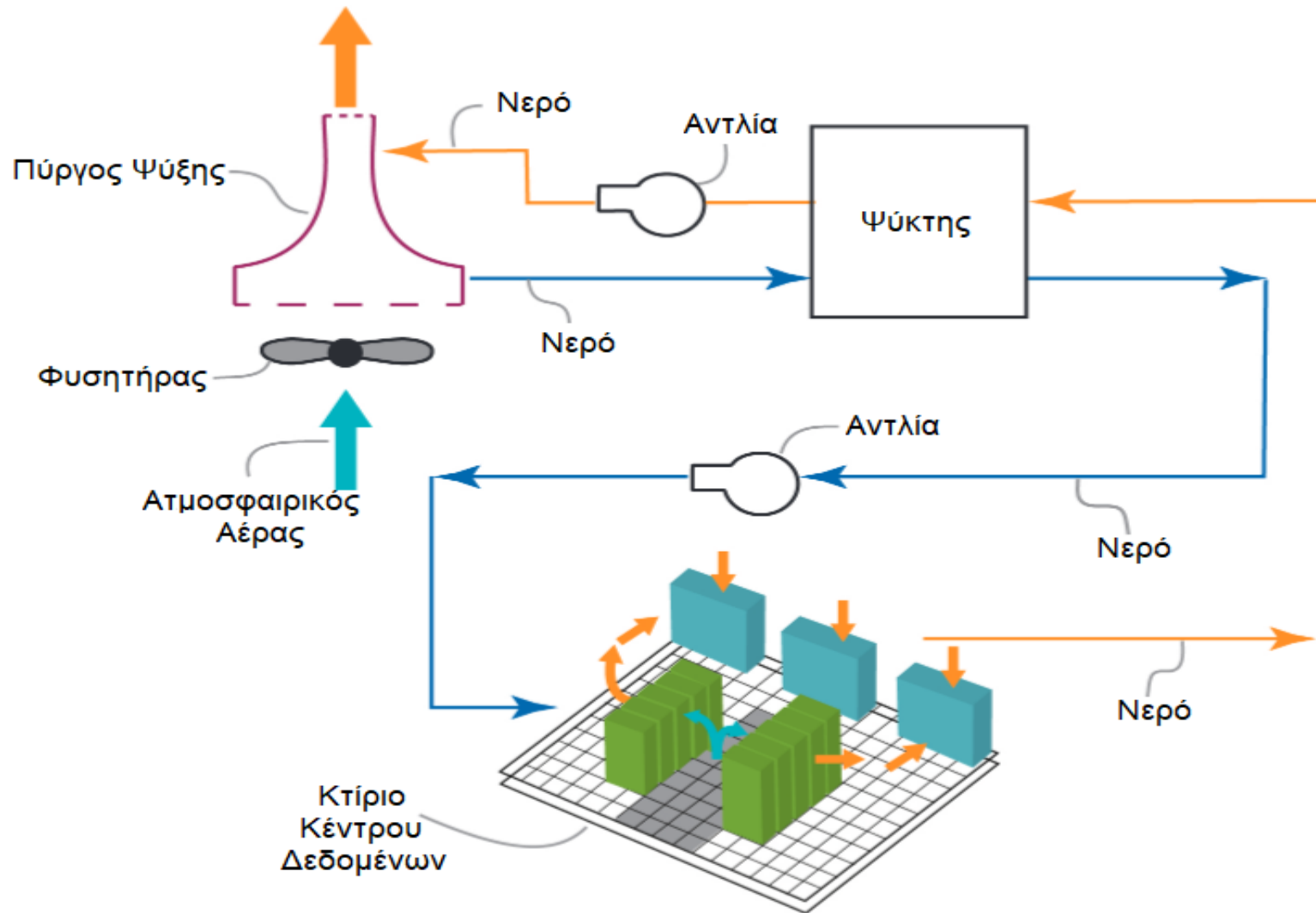
- Σχέδιο ροής αέρα ενός κέντρου δεδομένων

- Το υψωμένο πάτωμα ενός κέντρου δεδομένων που έγινε με τη βοήθεια κρύου διαδρόμου



# ΣΥΣΤΗΜΑΤΑ ΨΥΞΗΣ ΚΕΝΤΡΩΝ ΔΕΔΟΜΕΝΩΝ (3/3)

## Σύστημα ψύξης κέντρων δεδομένων τριών βρόχων



# Ο ΡΟΛΟΣ ΤΩΝ CRACS

---

- Όλα τα CRACS περιέχουν εναλλάκτη θερμότητας, έλεγχο του αέρα και χειριστήρια. Διαφέρουν ανάλογα με τον τύπο ψύξης που χρησιμοποιούν:
  - Άμεση επέκταση DX (Direct Expansion): Μια μονάδα DX είναι ένα split κλιματιστικό με πηνία ψύξης (εξατμιστής) μέσα στο CRAC και πηνία εξαγωγής θερμότητας (συμπυκνωτής) εκτός του κέντρου δεδομένων.
  - Υγρό διάλυμα
  - Ψύξη με νερό





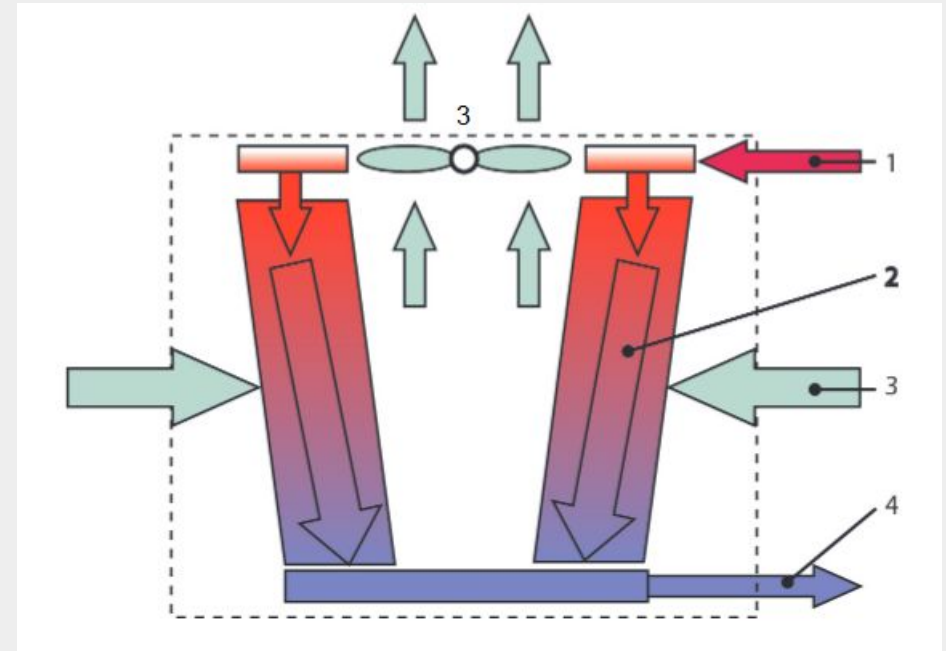
# ΨΥΚΤΕΣ

Κλιματιστικό με νερό: Βυθίζει τους εξατμιστήρες και τα πηνία συμπυκνωτή σε νερό σε δύο μεγάλα χωριστά διαμερίσματα, τα οποία συνδέονται μέσω του συστήματος ψύξης που είναι τοποθετημένο στο επάνω μέρος και αποτελείται από ένα συμπιεστή, μια βαλβίδα εκτόνωσης και σωληνώσεις. Καθώς το νερό ρέει πάνω από τα βυθισμένα πηνία, ψύχεται ή θερμαίνεται ανάλογα με την πλευρά του ψυκτικού συγκροτήματος.



## ΠΥΡΓΟΙ ΨΥΞΗΣ (1/2)

- Πώς λειτουργεί ένας πύργος ψύξης:

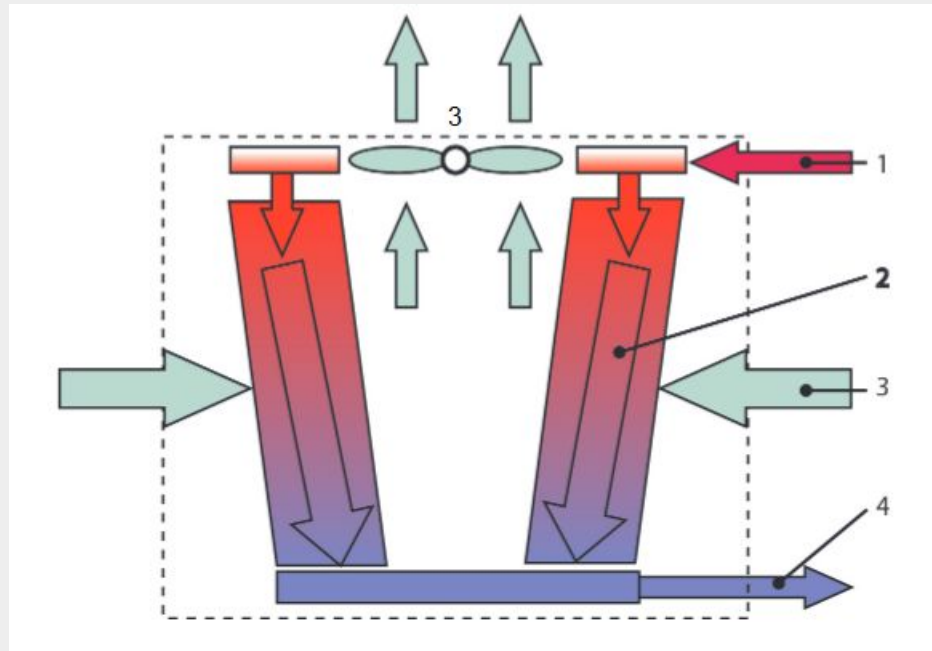


- 1. Το ζεστό νερό από το κέντρο δεδομένων ρέει από την κορυφή του πύργου ψύξης σε υλικό «πλήρωσης» μέσα στον πύργο. Το γέμισμα δημιουργεί πρόσθετη επιφάνεια για να βελτιώσει την απόδοση εξατμίσησης.
- 2. Καθώς το νερό ρέει κάτω από τον πύργο, κάποιο μέρος εξατμίζεται αντλώντας ενέργεια από το υπόλοιπο νερό.



## ΠΥΡΓΟΙ ΨΥΞΗΣ (2/2)

- Πώς λειτουργεί ένας πύργος ψύξης:



- 3. Ένας ανεμιστήρας στην κορυφή αντλεί αέρα μέσω του πύργου για να βοηθήσει στην εξάτμιση. Ο ξηρός αέρας εισέρχεται στις πλευρές και ο υγρός αέρας εξέρχεται από την κορυφή.
- 4. Το κρύο νερό συλλέγεται στη βάση του πύργου και επιστρέφει στο κέντρο δεδομένων.



# ΕΛΕΥΘΕΡΗ ΨΥΞΗ

---

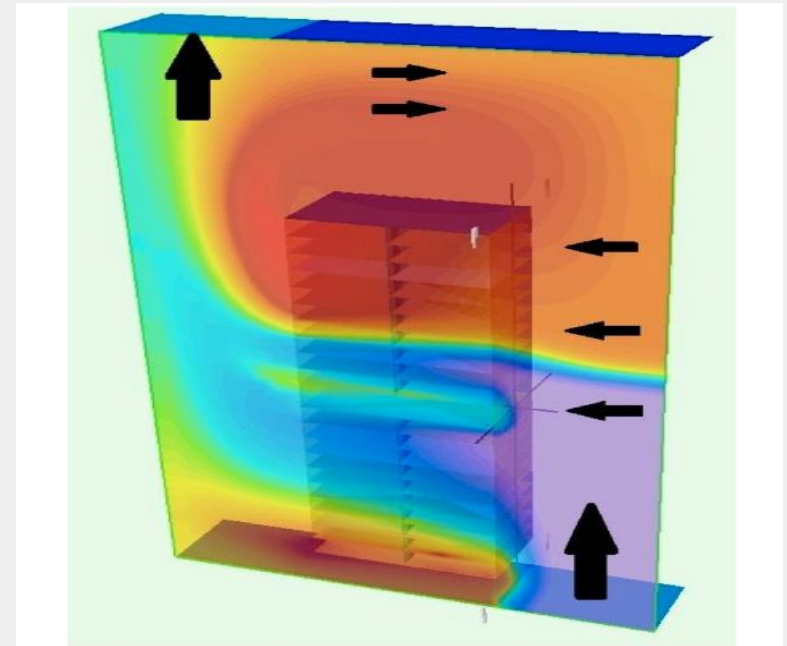
- Η ψύξη αυτή χρησιμοποιεί χαμηλές εξωτερικές θερμοκρασίες για να παράγει κρύο νερό ή χρησιμοποιεί εξωτερικό αέρα για να δροσίζει διακομιστές.
- Η ελεύθερη ψύξη δεν είναι πλήρως ελεύθερη, αλλά είναι πολύ αποτελεσματική σε σύγκριση με τη χρήση ψυκτικού συγκροτήματος.



# ΡΟΗ ΤΟΥ ΑΕΡΑ

- Τα περισσότερα κέντρα δεδομένων χρησιμοποιούν τη ρύθμιση του υπερυψωμένου δαπέδου. Για να γίνει αλλαγή της ποσότητας ψύξης σε ένα συγκεκριμένο ερμάριο ή σειρά, προσαρμόζουμε τον αριθμό των διάτρητων πλακιδίων αντικαθιστώντας τα πλακάκια με διάτρητα ή αντίστροφα. Για να λειτουργήσει καλά η ψύξη, η κρύα ροή αέρα που περνάει από την οροφή πρέπει να ταιριάζει με την οριζόντια ροή αέρα μέσω των διακομιστών στο ερμάριο.

- Μοντέλο CFD (Computational Fluid Dynamics) που δείχνει διαδρομές ανακύκλωσης και διαστρωμάτωση θερμοκρασίας για ερμάριο με ανεπαρκή ροή αέρα.



## IN-RACK ΚΑΙ IN-ROW ΨΥΞΗ (1/2)

---

- Η ψύξη στο ερμάριο (in-rack cooling) μπορεί να αυξήσει την πυκνότητα ισχύος. Μπορεί να αφαιρέσει μόνο μέρος της θερμότητας ή όλη τη θερμότητα, αντικαθιστώντας αποτελεσματικά τα CRACs.
- Η ψύξη στη σειρά (in-row cooling) λειτουργεί όπως η ψύξη στο ερμάριο, εκτός από τα πηνία ψύξης που δεν βρίσκονται στο ερμάριο, αλλά δίπλα σε αυτό. Ένας συλλέκτης κατευθύνει τον ζεστό αέρα στα πηνία και αποτρέπει τη διαρροή στο ψυχρό διάδρομο.

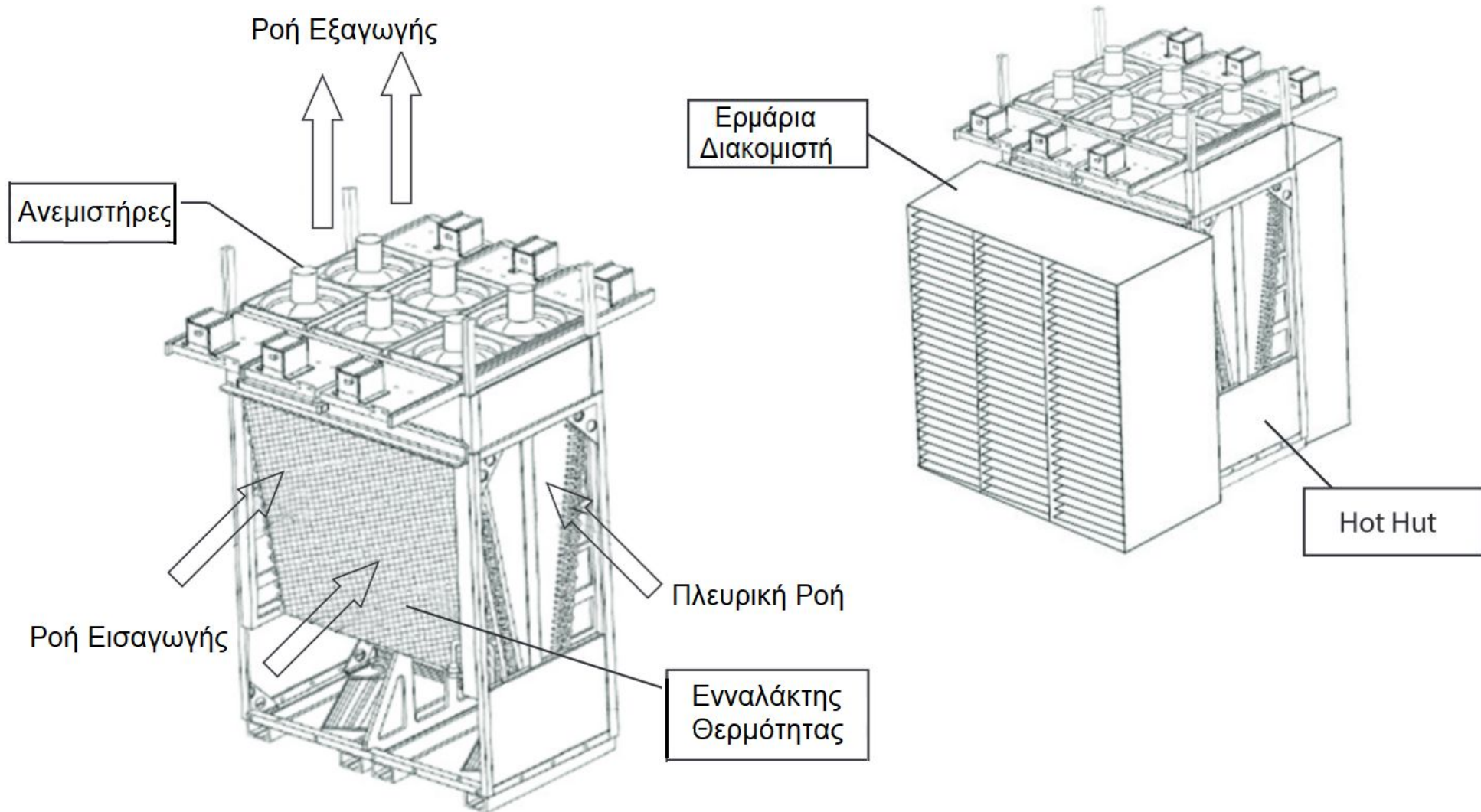


- Συνήθως, δεν είναι πρακτικό να ψύχονται όλα τα εξαρτήματα με ψυχρές πλάκες, έτσι αντί αυτού στοχεύουμε τα εξαρτήματα υψηλότερης ισχύος, όπως οι CPUs, και χρησιμοποιούμε αέρα για να αφαιρέσουμε το υπόλοιπο της θερμότητας.
- Αν και οι ψυχρές πλάκες μπορούν να απομακρύνουν τα υψηλά τοπικά φορτία θερμότητας, είναι σχετικά σπάνια λόγω του κόστους και της πολυπλοκότητας του σχεδιασμού των πλακών και των συνδέσμων που απαιτούνται για τη σύνδεση και αποσύνδεση του βρόχου.



# ΤΟ ΣΥΣΤΗΜΑ IN-ROW ΤΗΣ GOOGLE (1/3)

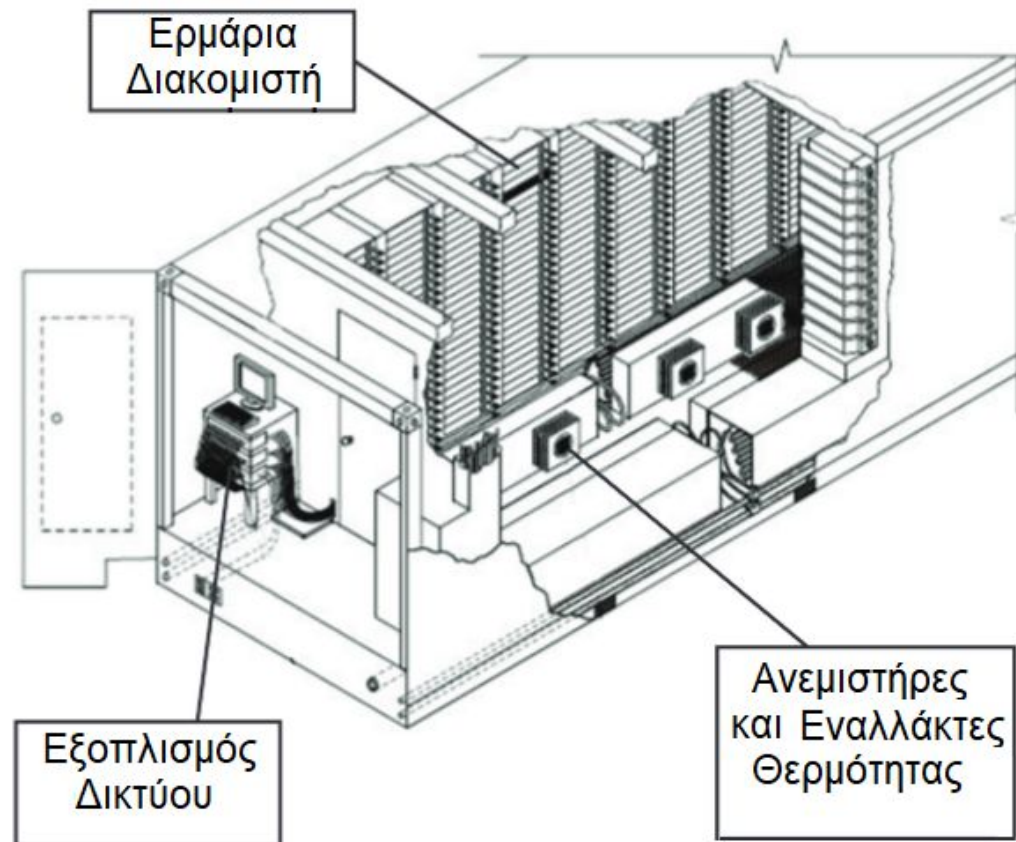
Google's "Hot Hut", μια εν σειρά, υδρόψυκτη μονάδα ανεμιστήρα





## ΤΟ ΣΥΣΤΗΜΑ IN-ROW ΤΗΣ GOOGLE (2/3)

Ο σχεδιασμός των containers της Google περιλαμβάνει όλη την υποδομή του δαπέδου του κέντρου δεδομένων.



## ΤΟ ΣΥΣΤΗΜΑ IN-ROW ΤΗΣ GOOGLE (3/3)

Το κέντρο δεδομένων της Google με βάση τα containers



## Δυνατότητα Ενεργειακής Απόδοσης Κέντρου Δεδομένων

$$\text{Αποδοτικότητα} = \frac{\text{Υπολογιστική}}{\text{Συνολική Ενέργεια}} = \left( \frac{1}{\text{PUE}} \right) \times \left( \frac{1}{\text{SPUE}} \right) \times \left( \frac{\text{Υπολογιστική}}{\text{Συνολική Ενέργεια στα Ηλεκτρονικά Εξαρτήματα}} \right)$$



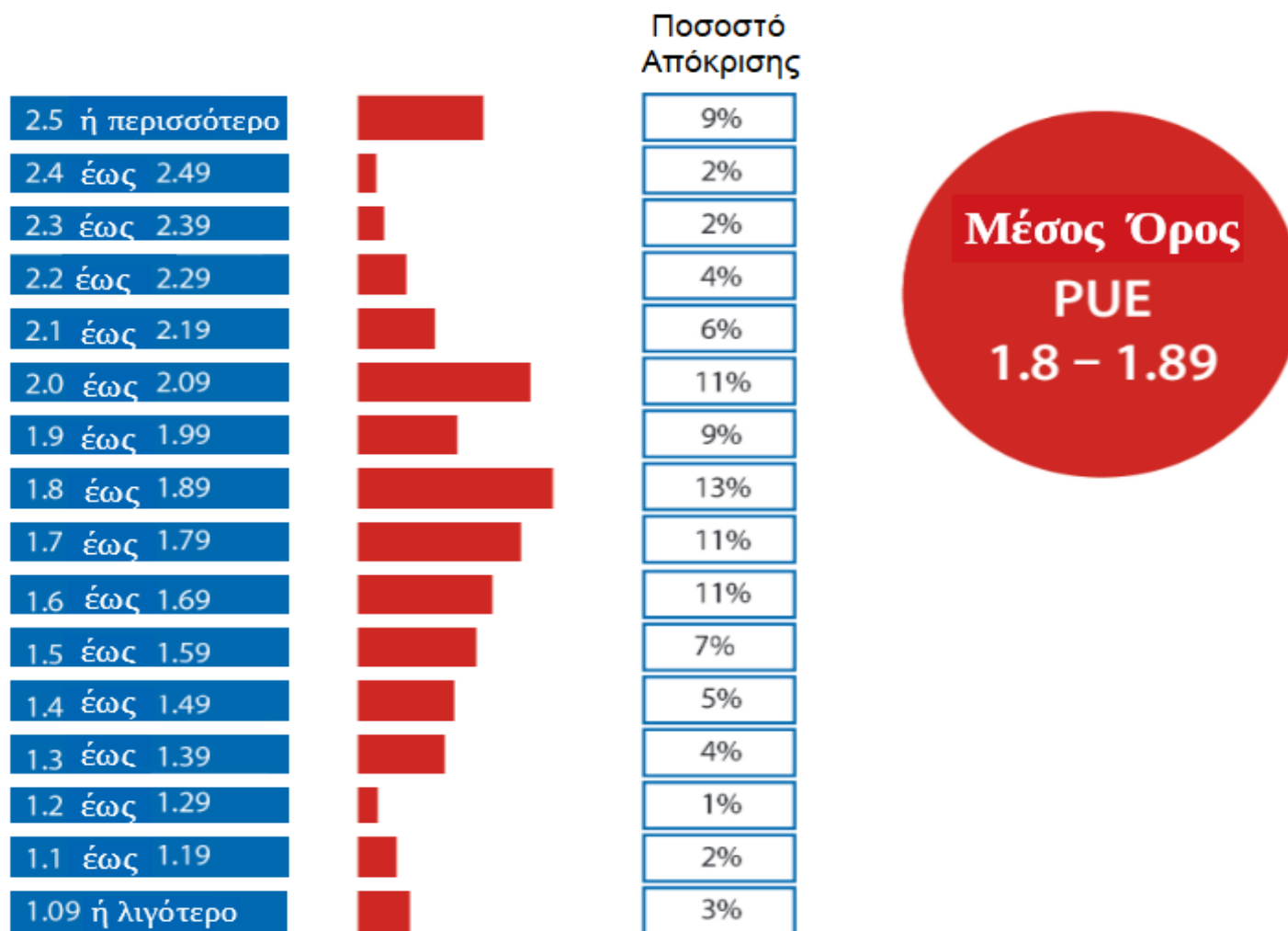
## ΑΠΟΔΟΣΗ ΕΝΕΡΓΕΙΑΣ (1/4)

---

- Αντανακλά την ποιότητα της ίδιας της υποδομής του κέντρου δεδομένων και καταγράφει την αναλογία της συνολικής ισχύος κτιρίου με την ισχύ της τεχνολογίας πληροφοριών.
- Η PUE (Power Usage Effectiveness) έχει γίνει πολύ δημοφιλής ως μέτρηση αποτελεσματικότητας του κέντρου δεδομένων δεδομένου ότι η ευρεία αναφορά άρχισε γύρω στο 2009. Μπορούμε εύκολα να μετρήσουμε την PUE προσθέτοντας ηλεκτρικούς μετρητές στις γραμμές που τροφοδοτούν τα διάφορα τμήματα ενός κέντρου δεδομένων, προσδιορίζοντας έτσι πόση ισχύς χρησιμοποιείται από ψυκτικά συγκροτήματα ή UPS.



## ΑΠΟΔΟΣΗ ΕΝΕΡΓΕΙΑΣ (2/4)



Αποτελέσματα έρευνας του 2012 για 1100 κέντρα δεδομένων σχετικά με το μέσο όρο της απόδοσης PUE.



## ΑΠΟΔΟΣΗ ΕΝΕΡΓΕΙΑΣ (3/4)

---

- Σημαντικοί παράγοντες που μπορούν να εξουδετερώσουν τις τιμές της PUE είναι οι ακόλουθοι:
  - Όλες οι μετρήσεις PUE δεν περιλαμβάνουν τα ίδια γενικά έξοδα.
  - Οι στιγμιαίες PUE διαφέρουν από τις μέσες PUE. Κατά τη διάρκεια μιας ημέρας ή ενός έτους, η PUE μιας εγκατάστασης μπορεί να ποικίλει σημαντικά.



# ΒΕΛΤΙΩΣΗ ΤΗΣ ΕΝΕΡΓΕΙΑΚΗΣ ΑΠΟΔΟΣΗΣ ΤΩΝ ΚΕΝΤΡΩΝ ΔΕΔΟΜΕΝΩΝ (1/2)

- Προσοχή στον χειρισμό της ροής του αέρα: Πρέπει να γίνεται διαχωρισμός του θερμού αέρα που εξαντλείται από τους διακομιστές από τον κρύο αέρα και να διατηρείται η διαδρομή προς το ψυκτικό πηνίο έτσι ώστε να καταναλώνεται ελάχιστη ενέργεια που κινεί το κρύο ή το ζεστό αέρα σε μεγάλες αποστάσεις.
- Αυξημένες θερμοκρασίες: Πρέπει ο κρύος διάδρομος να παραμένει στους 25-30°C παρά στους 18-20°C. Οι υψηλότερες θερμοκρασίες καθιστούν πολύ πιο εύκολη την αποτελεσματική δρομολόγηση των κέντρων δεδομένων. Ουσιαστικά κανένας εξοπλισμός διακομιστή ή δικτύου δεν χρειάζεται θερμοκρασίες εισόδου 20°C και δεν υπάρχουν στοιχεία ότι οι υψηλότερες θερμοκρασίες προκαλούν περισσότερες αστοχίες των εξαρτημάτων.



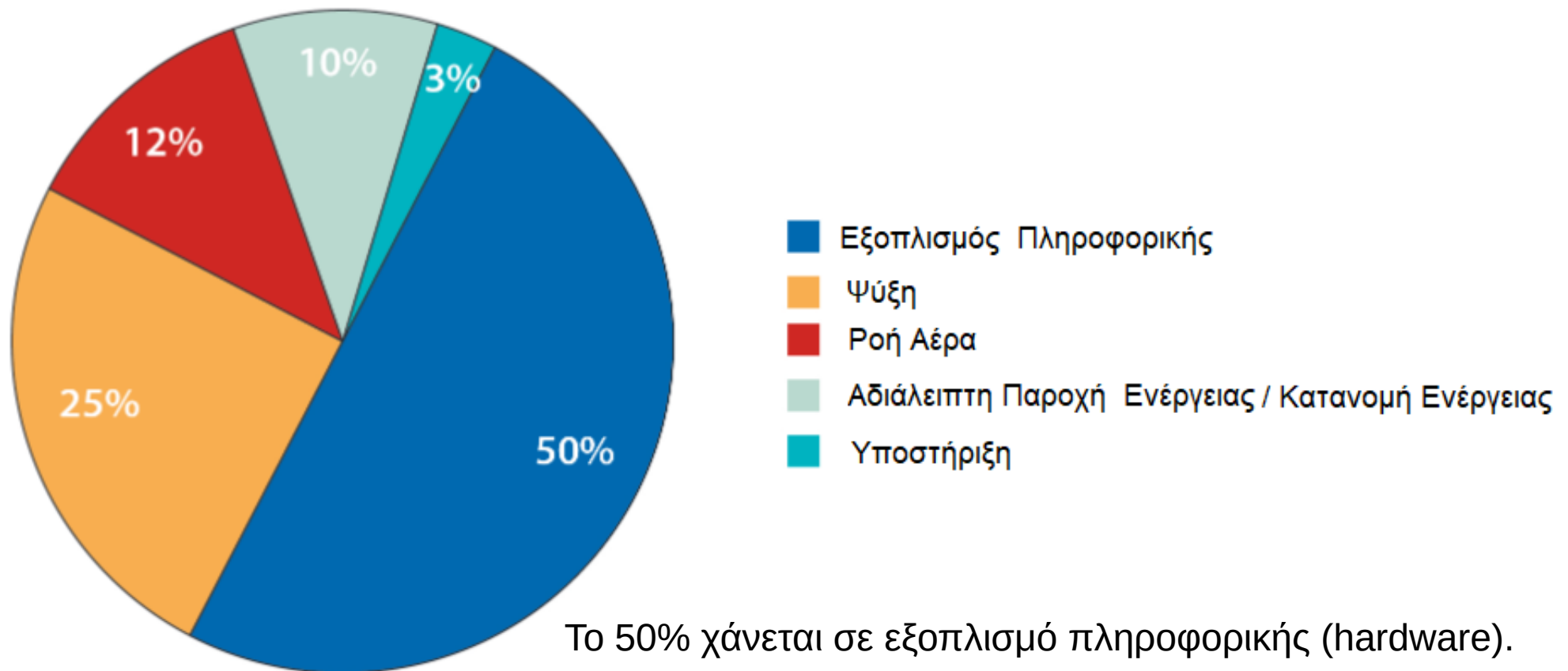
- Ορισμένες PUEs δεν αντικατοπτρίζουν την πραγματικότητα. Συχνά οι πωλητές δημοσιεύουν PUE «σχεδιασμού» που υπολογίζονται με βάση τις βέλτιστες συνθήκες λειτουργίας και τις ονομαστικές τιμές απόδοσης ή δημοσιεύουν μια τιμή που μετράται κατά τη διάρκεια μιας δοκιμής μικρού φορτίου υπό βέλτιστες συνθήκες.
- Συνήθως, οι τιμές PUE που παρέχονται χωρίς λεπτομέρειες εμπίπτουν σε αυτήν την κατηγορία.





# ΠΗΓΕΣ ΑΠΩΛΕΙΑΣ ΑΠΟΤΕΛΕΣΜΑΤΙΚΟΤΗΤΑΣ ΣΕ ΚΕΝΤΡΑ ΔΕΔΟΜΕΝΩΝ

Απώλεια ισχύος σε ένα κλασικό κέντρο δεδομένων



# ΒΕΛΤΙΩΣΗ ΤΗΣ ΕΝΕΡΓΕΙΑΚΗΣ ΑΠΟΔΟΣΗΣ ΤΩΝ ΚΕΝΤΡΩΝ ΔΕΔΟΜΕΝΩΝ (2/2)

- Ελεύθερη ψύξη: Στα πιο μέτρια κλίματα, η ελεύθερη ψύξη μπορεί να εξαλείψει το μεγαλύτερο μέρος του χρόνου λειτουργίας του ψυκτικού συγκροτήματος ή να εξαλείψει συνολικά τους ψυκτικούς θαλάμους.
- Καλύτερη αρχιτεκτονική συστήματος ισχύος: Οι απώλειες UPS και κατανομής ισχύος μπορούν συχνά να μειωθούν σημαντικά με την επιλογή εργαλείων υψηλότερης απόδοσης.



# ΥΠΟΛΟΓΙΣΜΟΣ ΤΗΣ ΑΠΟΔΟΣΗΣ

$$\text{Αποδοτικότητα} = \frac{\text{Υπολογιστική}}{\text{Συνολική Ενέργεια}} = \left( \frac{1}{\text{PUE}} \right) \times \left( \frac{1}{\text{SPUE}} \right) \times \left( \frac{\text{Υπολογιστική}}{\text{Συνολική Ενέργεια στα Ηλεκτρονικά Εξαρτήματα}} \right)$$

(α)                      (β)                      (γ)

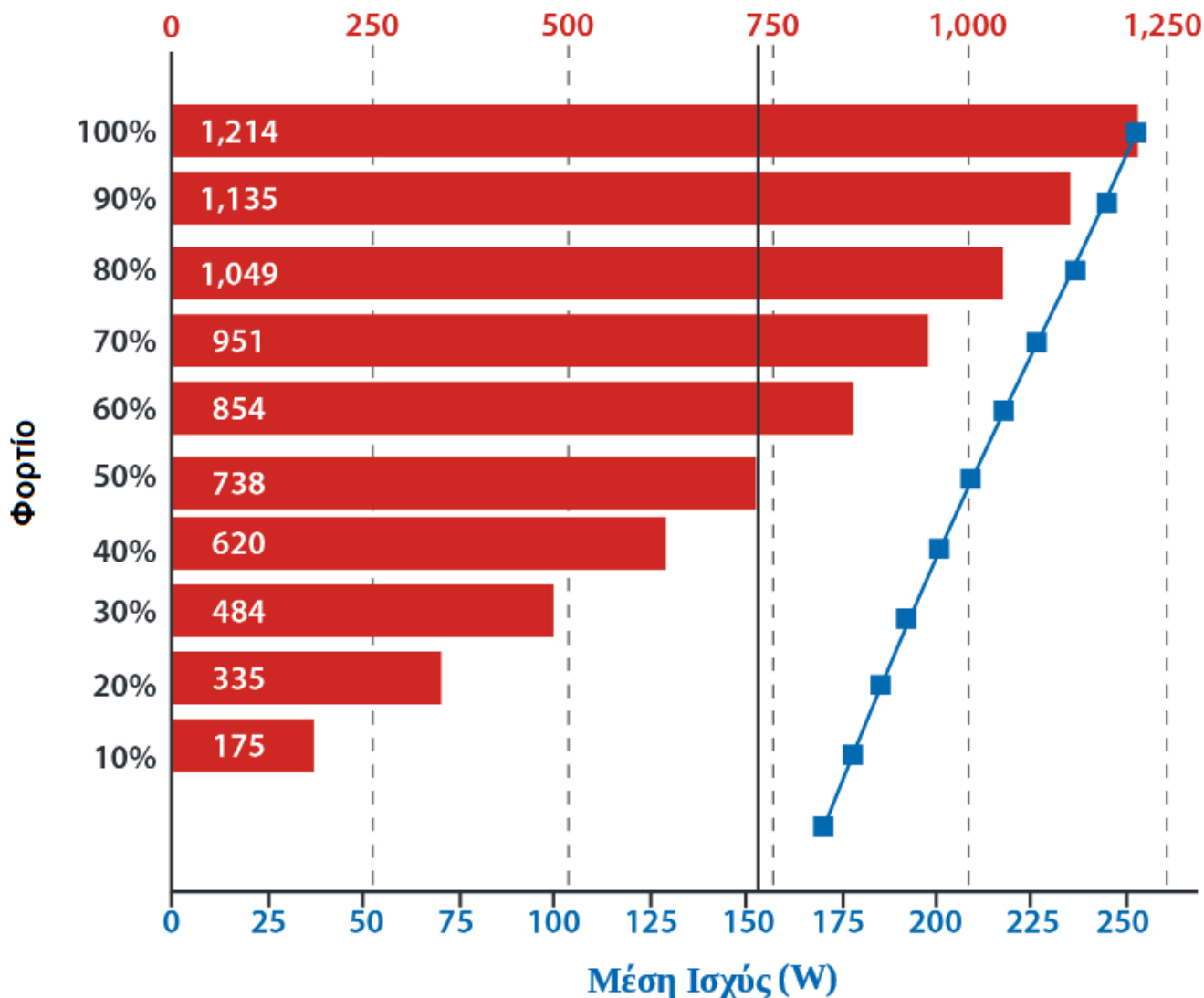
- Ο δεύτερος όρος (β) καλύπτει τα γενικά έξοδα εντός διακομιστών ή άλλου εξοπλισμού πληροφορικής που χρησιμοποιεί μια μέτρηση ανάλογη με την PUE, την server PUE (SPUE).
- Η SPUE αποτελείται από το λόγο της συνολικής ισχύος εισόδου του διακομιστή στη ωφέλιμη του ισχύ, όπου η ωφέλιμη ισχύς περιλαμβάνει μόνο την ισχύ που καταναλώνουν τα ηλεκτρονικά εξαρτήματα που εμπλέκονται άμεσα στον υπολογισμό: μητρική πλακέτα, δίσκοι, CPU, DRAM (Dynamic Random Access Memory), κάρτες εισόδου/εξόδου κ.ο.κ. Μπορεί να χαθούν σημαντικές ποσότητες ενέργειας στην τροφοδοσία ρεύματος του διακομιστή, στις μονάδες ρυθμιστή τάσης (VRM), Voltage Regulator Modules, και στους ανεμιστήρες ψύξης.



# ΕΝΕΡΓΕΙΑΚΗ ΑΠΟΔΟΣΗ ΤΟΥ SERVER

Η σχέση απόδοσης ισχύος είναι η μέγιστη σε φορτίο 100%

Σχέση Απόδοσης Ισχύος



Παράδειγμα συγκριτικής αξιολόγησης για το SPECpower\_ssj2008.

Η ενεργειακή απόδοση υποδεικνύεται με μπάρες, ενώ η κατανάλωση ενέργειας υποδεικνύεται από τη γραμμή.

Το σύστημα περιλαμβάνει ένα single-chip 2.83 GHz τεσσάρων πυρήνων Intel Xeon, 4 GB DRAM και μία μονάδα δίσκου SATA 3,5" ΙΝΤΣΩΝ.

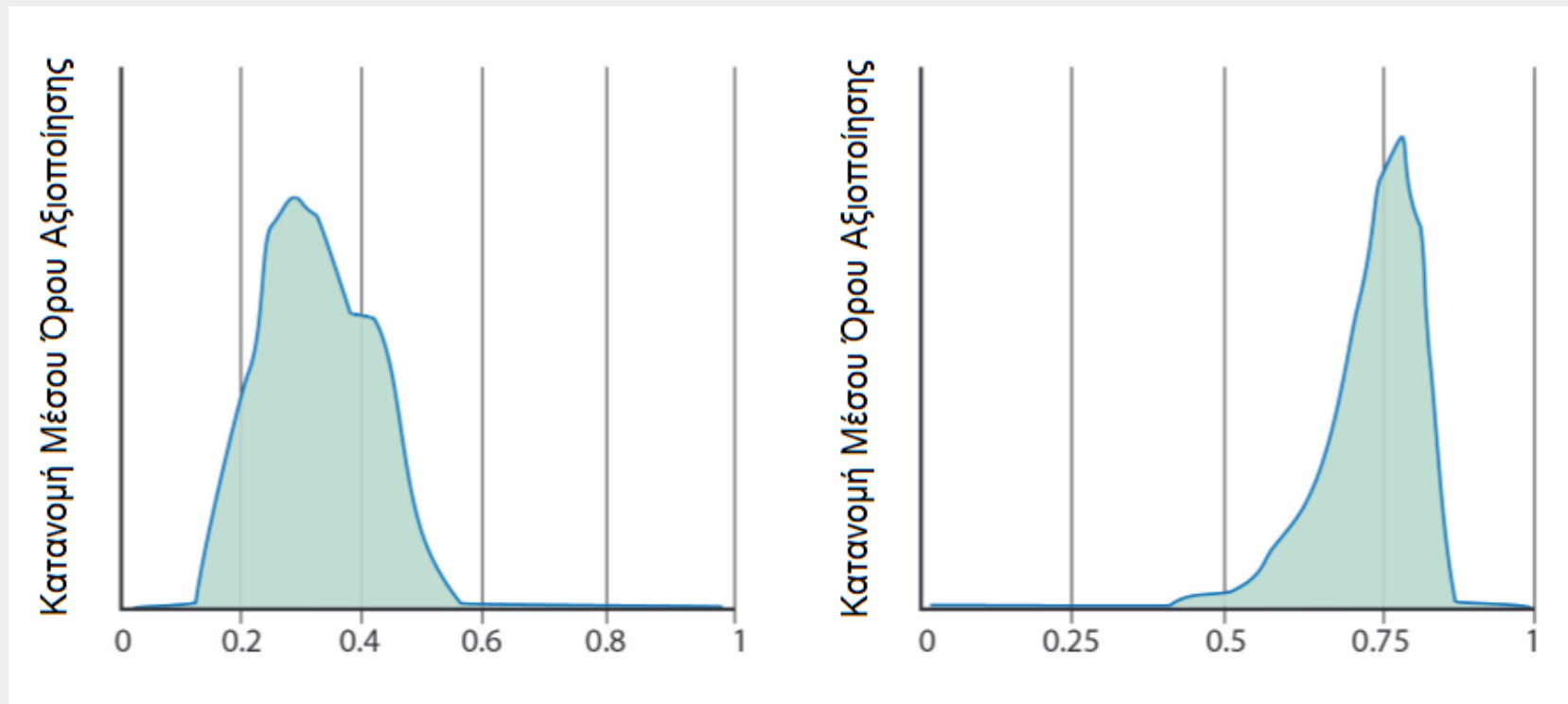


- Συνήθως κάποιος θέλει να μετρήσει την αξία που λαμβάνεται από την ενέργεια που καταναλώνεται στον υπολογισμό.
- Για παράδειγμα, για να συγκρίνει τη σχετική απόδοση δύο WSCs ή να καθοδηγήσει τις επιλογές σχεδίασης για τα νέα συστήματα.



## ΠΡΟΦΙΛ ΧΡΗΣΗΣ ΤΩΝ WSC

Μέση κατανομή δραστηριότητας ενός δείγματος από δύο ομάδες Google, το καθένα από τα οποία περιλαμβάνει πάνω από 20.000 διακομιστές, για μια περίοδο 3 μηνών (Ιανουάριος-Μάρτιος 2013).



- Προτείνεται να προστεθεί η ενεργειακή αναλογικότητα ως στόχος σχεδιασμού για τον υπολογισμό των στοιχείων.
- Στην ιδανική περίπτωση, τα ενεργειακά αναλογικά συστήματα καταναλώνουν σχεδόν καθόλου ενέργεια όταν είναι αδρανή (ιδιαίτερα στις ενεργές αδρανείς καταστάσεις όπου είναι ακόμα διαθέσιμα για να δουλέψουν) και καταναλώνουν σταδιακά περισσότερη ισχύ καθώς αυξάνεται το επίπεδο δραστηριότητας. Ένας απλός τρόπος να αιτιολογηθεί αυτή η ιδανική καμπύλη είναι να υποθέσουμε γραμμικότητα.



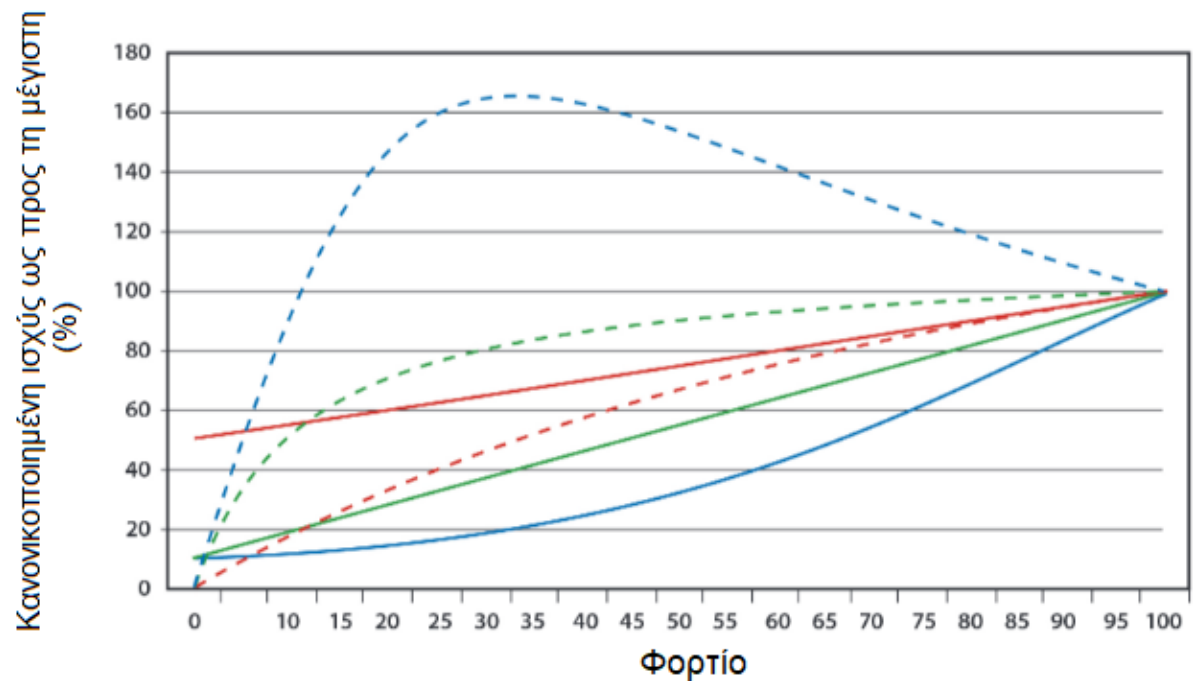
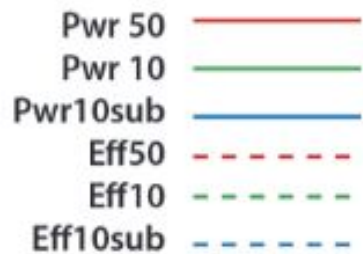
## ΕΝΕΡΓΕΙΑ – ΑΝΑΛΟΓΙΚΟΣ ΥΠΟΛΟΓΙΣΜΟΣ (2/2)

Η ισχύς και η αντίστοιχη απόδοση ισχύος τριών υποθετικών συστημάτων.

Οι συμπαγείς γραμμές αντιπροσωπεύουν ισχύ % (κανονικοποιημένη σε μέγιστη ισχύ). Οι διακεκομμένες γραμμές αντιπροσωπεύουν την απόδοση ως ποσοστό της ενεργειακής απόδοσης στην κορυφή.

Το Pwr αναφέρεται στη ισχύ 50% και 10% αντίστοιχα.

Το Eff αναφέρεται στην απόδοση 50% και 10% αντίστοιχα.





- Αν και οι CPUs έχουν ιστορικά κακή φήμη όσον αφορά τη χρήση ενέργειας, δεν είναι αναγκαστικά ο κύριος ένοχος για την κακή ενεργειακή αναλογικότητα.
- Τα τελευταία χρόνια, οι σχεδιαστές CPU έδωσαν μεγαλύτερη προσοχή στην ενεργειακή απόδοση.
- Ταυτόχρονα, γίνεται μετάβαση σε αρχιτεκτονικές πολλαπλών χρήσεων αντί να συνεχίζεται η αυξανόμενη τάση για μέγιστες ταχύτητες.

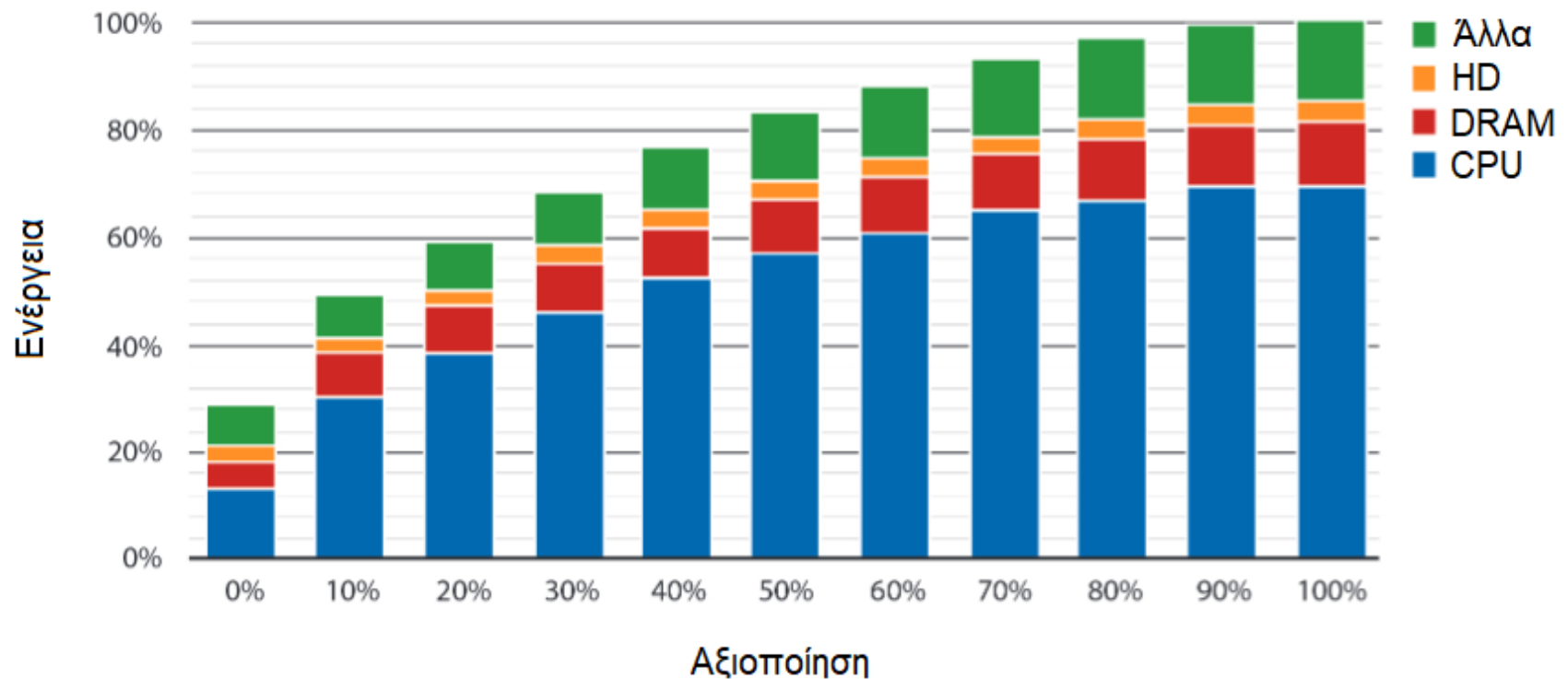


- Η σχετική συμβολή του συστήματος μνήμης στη συνολική κατανάλωση ενέργειας μειώθηκε τα τελευταία 5 χρόνια όσον αφορά τη χρήση ενέργειας από την CPU, αντιστρέφοντας την τάση υψηλότερου ενεργειακού προφίλ DRAM σε όλη τη διάρκεια της προηγούμενης δεκαετίας. Η μείωση του κλάσματος της ενέργειας που χρησιμοποιείται στα συστήματα μνήμης οφείλεται σε ένα συνδυασμό παραγόντων:
  - η νεότερη τεχνολογία DDR3 είναι ουσιαστικά πιο αποδοτική από την προηγούμενη τεχνολογία (FBDIMMs)
  - τα επίπεδα τάσης τσιπ DRAM έχουν πέσει από 1,8 V σε λιγότερο από 1,5 V
  - νέα τσιπ CPU χρησιμοποιούν περισσότερη ενέργεια



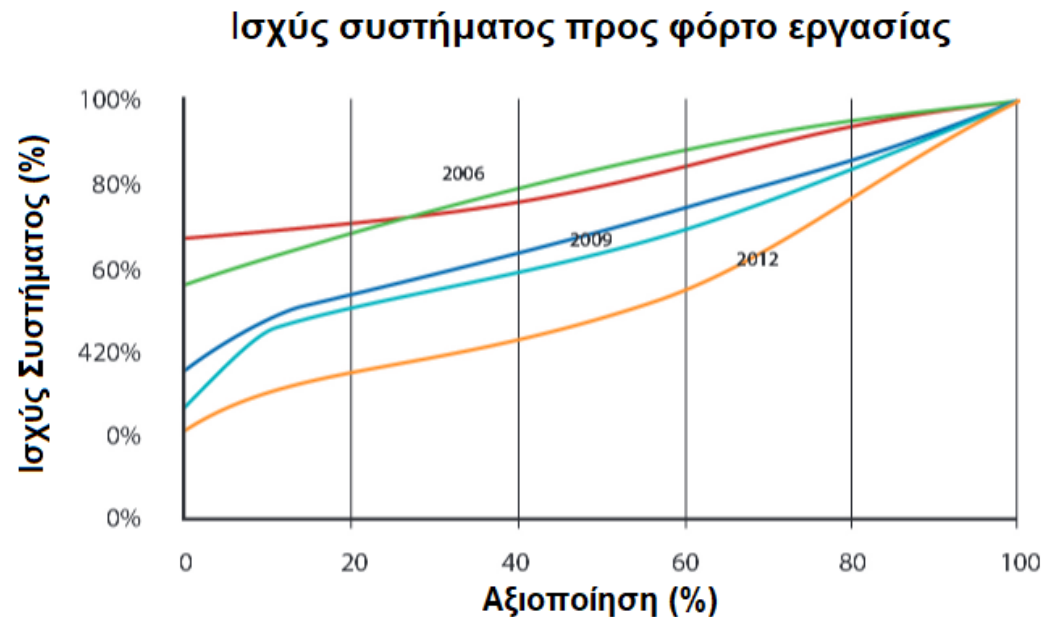
## ΑΙΤΙΕΣ ΦΤΩΧΗΣ ΕΝΕΡΓΕΙΑΚΗΣ ΑΝΑΛΟΓΙΚΟΤΗΤΑΣ (3/3)

Η χρήση ενέργειας υποσυστήματος σε ένα διακομιστή x86, καθώς το φορτίο υπολογίζεται από την αδράνεια στην πλήρη χρήση.



# ΒΕΛΤΙΩΣΗ ΤΗΣ ΑΝΑΛΟΓΙΚΟΤΗΤΑΣ ΤΗΣ ΕΝΕΡΓΕΙΑΣ

- Κανονικοποιημένη ισχύς συστήματος έναντι χρήσης σε διακομιστές της Intel από το 2006 και το 2012 (ευγενική παραχώρηση του Winston Sounders, Intel).
- Οι πράσινες και κόκκινες γραμμές αντιπροσωπεύουν διακομιστές 2006, οι δύο μπλε είναι οι διακομιστές 2009 και η πορτοκαλί γραμμή είναι πιο πρόσφατη. Το διάγραμμα δείχνει ότι οι διακομιστές της Intel έχουν γίνει πιο ενεργειακά ανάλογες σε αυτή την επταετή περίοδο.



Πληροφορίες από SPEC.org



- Ενώ η αναλογικότητα της ενέργειας του επεξεργαστή έχει βελτιωθεί, απαιτείται ακόμη μεγαλύτερη προσπάθεια για DRAM, αποθήκευση και δικτύωση.
  - Carrera et al: εξέταση της ενεργειακής επίδρασης των μονάδων πολλαπλών ταχυτήτων και των συνδυασμών μονάδων διακομιστών και φορητών υπολογιστών για την επίτευξη αναλογικής ενεργειακής συμπεριφοράς.
  - Sankar et al: εξέταση διαφορετικών αρχιτεκτονικών για δίσκους, παρατηρώντας ότι επειδή οι κινήσεις της κεφαλής είναι σχετικά ενεργειακά ανάλογες, ένας δίσκος με χαμηλότερη ταχύτητα περιστροφής και πολλαπλές κεφαλές μπορεί να επιτύχει παρόμοια απόδοση και χαμηλότερη ισχύ σε σύγκριση με ένα μόνιμα υψηλό RPM.



- Είναι πιθανό να προβλέψουμε ότι ο εξοπλισμός δικτύωσης είναι υπεύθυνος για το 10-20% των χρημάτων για την ενέργεια της εγκατάστασης. Σε εκείνο το σημείο, η έλλειψη αναλογικότητας θα είναι σοβαρό πρόβλημα. Για να επεξηγήσουμε αυτό το σημείο, ας υποθέσουμε ένα σύστημα που παρουσιάζει ένα γραμμικό προφίλ χρήσης ενέργειας ως συνάρτηση της χρήσης ( $u$ ):
  - $P(u) = P_i + u(1-P_i)$
  - Στην παραπάνω εξίσωση, το  $P_i$  αντιπροσωπεύει την αδρανή ισχύ του συστήματος και η μέγιστη ισχύς κανονικοποιείται στο 1,0. Σε ένα τέτοιο σύστημα, η ενεργειακή απόδοση γίνεται  $u/P(u)$ , η οποία μειώνεται βάσει του νόμου Amdahl.

$$E(u) = \frac{1}{1 - P_i + P_i/u}$$



# ΣΧΕΤΙΚΗ ΑΠΟΤΕΛΕΣΜΑΤΙΚΟΤΗΤΑ ΤΩΝ ΛΕΙΤΟΥΡΓΙΩΝ ΧΑΜΗΛΗΣ ΤΡΟΦΟΔΟΣΙΑΣ (1/2)

---

- Η ύπαρξη διαστημάτων μεγάλης αδράνειας θα επέτρεπε την επίτευξη υψηλότερης ενεργειακής αναλογικότητας χρησιμοποιώντας διάφορα είδη sleep modes.
- Οι αδρανείς λειτουργίες χαμηλής κατανάλωσης αναπτύχθηκαν αρχικά για κινητές και ενσωματωμένες συσκευές και είναι πολύ επιτυχείς σε αυτόν τον τομέα.
- Οι περισσότερες από αυτές τις τεχνικές είναι ανεπαρκείς για συστήματα WSCs.



# ΣΧΕΤΙΚΗ ΑΠΟΤΕΛΕΣΜΑΤΙΚΟΤΗΤΑ ΤΩΝ ΛΕΙΤΟΥΡΓΙΩΝ ΧΑΜΗΛΗΣ ΤΡΟΦΟΔΟΣΙΑΣ (2/2)

- Μεγάλη εξοικονόμηση ενέργειας είναι διαθέσιμη από αδρανείς λειτουργίες χαμηλής κατανάλωσης ενέργειας, όπως οι μονάδες δίσκων με περιστροφή.
- Τα προτεινόμενα συστήματα, PowerNap και IdleCap, υποθέτουν ότι οι συνιστώσες δεν έχουν χρήσιμες λειτουργίες χαμηλής ισχύος εκτός από την πλήρη αδράνεια και διαμορφώνουν μεταβάσεις σε αδράνεια σε όλες τις επιμέρους συνιστώσες. Αυτό γίνεται προκειμένου να μειωθεί η ισχύς σε χαμηλότερες τιμές, περιορίζοντας ταυτόχρονα την απόδοση.





- Γενικά τα εξαρτήματα υλικού πρέπει να υποστούν σημαντικές βελτιώσεις στην ενεργειακή αναλογικότητα για να καταστήσουν πιο ενεργειακά αποδοτικά συστήματα WSCs.
- Ωστόσο, η πιο έξυπνη διαχείριση της ενέργειας και η προγραμματισμένη υποδομή λογισμικού παίζει σημαντικό ρόλο. Για ορισμένους τύπους στοιχείων, η επίτευξη τέλει ενεργειακής αναλογικής συμπεριφοράς μπορεί να μην είναι εφικτός στόχος.
- Οι σχεδιαστές θα πρέπει να εφαρμόσουν στρατηγικές λογισμικού για έξυπνη χρήση των λειτουργιών διαχείρισης ισχύος σε υπάρχον υλικό, χρησιμοποιώντας χαμηλής λειτουργίας εναέρια ή ενεργά προγράμματα χαμηλής κατανάλωσης ισχύος, καθώς και εφαρμόζοντας προγραμματισμό των εργασιών για την ενίσχυση της ενεργειακής αναλογικότητας των συστημάτων υλικού.



- Το λογισμικό πρέπει να ξεπεράσει δύο βασικές προκλήσεις: την ενσωμάτωση και την απόδοση.
- Οι μηχανισμοί που ασχολούνται με την ενέργεια πρέπει να ενσωματώνονται σε υπομονάδες χαμηλότερου επιπέδου, ώστε να ελαχιστοποιείται η έκθεση των προγραμματιστών εφαρμογών σε πρόσθετη πολυπλοκότητα υποδομής.
- Αυτή η εντατική διαχείριση ισχύος θέτει προβλήματα μη ελέγχου.



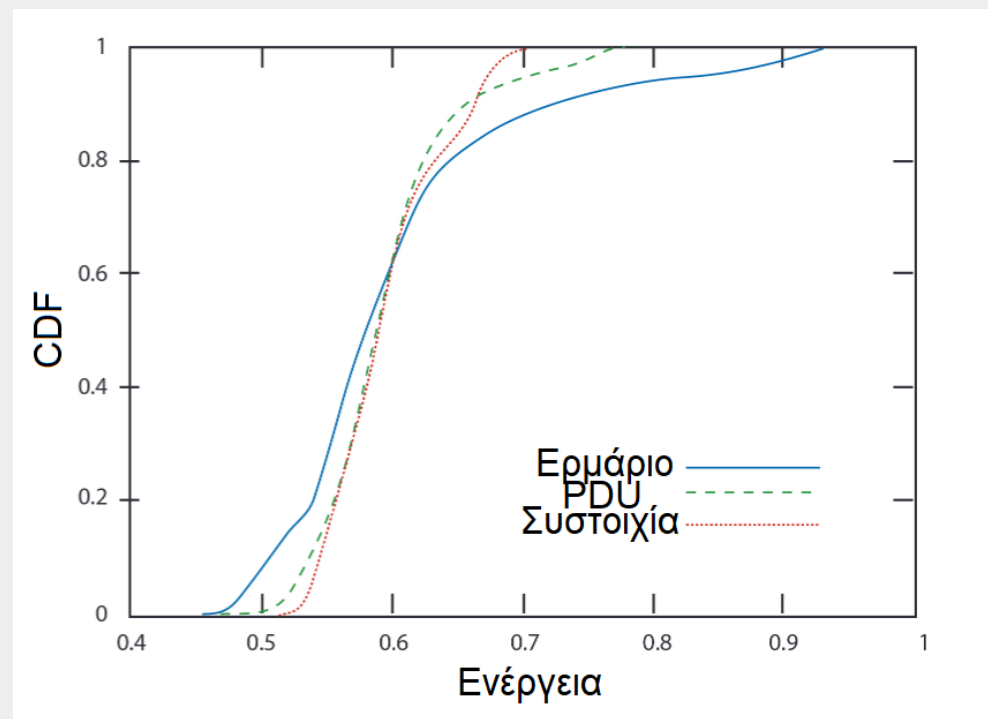
- Πρώτον, οι προδιαγραφές διακομιστή παρέχουν συνήθως πολύ συντηρητικές τιμές για τη μέγιστη κατανάλωση ενέργειας.
- Δεύτερον, η πραγματική κατανάλωση ενέργειας ποικίλλει σημαντικά με το φορτίο (χάρη στην ενεργειακή αναλογικότητα) και μπορεί να είναι δύσκολο να προβλεφθεί η μέγιστη κατανάλωση ισχύος μιας ομάδας διακομιστών.



# ΥΠΕΡΚΑΛΥΨΗ ΤΗΣ ΙΣΧΥΟΣ ΕΓΚΑΤΑΣΤΑΣΗΣ (1/3)

- Μόλις χρησιμοποιήσουμε οτιδήποτε άλλο εκτός από την πιο συντηρητική εκτίμηση της κατανάλωσης ενέργειας του εξοπλισμού για την ανάπτυξη συμπλεγμάτων, υπάρχει κάποιος κίνδυνος να υπερβούμε τη διαθέσιμη ποσότητα ισχύος, δηλαδή να υπερκαλύψουμε την ισχύ της εγκατάστασης. Η επιτυχής εφαρμογή της υπερκάλυψης ισχύος αυξάνει τη συνολική αξιοποίηση του προϋπολογισμού ισχύος του κέντρου δεδομένων, ενώ ελαχιστοποιεί τον κίνδυνο καταστάσεων υπερφόρτωσης.

- Το σχήμα δείχνει την αθροιστική κατανομή της χρήσης ενέργειας με την πάροδο του χρόνου για ομάδες 80 διακομιστών (ερμάριο), 800 διακομιστών (PDU) και 5.000 διακομιστών (συστοιχία).



## ΥΠΕΡΚΑΛΥΨΗ ΤΗΣ ΙΣΧΥΟΣ ΕΓΚΑΤΑΣΤΑΣΗΣ (2/3)

- Η μελέτη αυτή αξιολογεί επίσης το δυναμικό των μηχανών με μεγαλύτερη ενεργειακή αναλογία για τη μείωση της μέγιστης κατανάλωσης ισχύος στο επίπεδο της εγκατάστασης.
- Προτείνει ότι η μείωση της αδρανούς ισχύος από 50% σε 10% της κορυφής μπορεί να μειώσει περαιτέρω τη χρήση μέγιστης ισχύος συστοιχίας κατά περισσότερο από 30%. Αυτό θα ισοδυναμούσε με πρόσθετη αύξηση της δυναμικότητας φιλοξενίας εγκαταστάσεων κατά 40%.

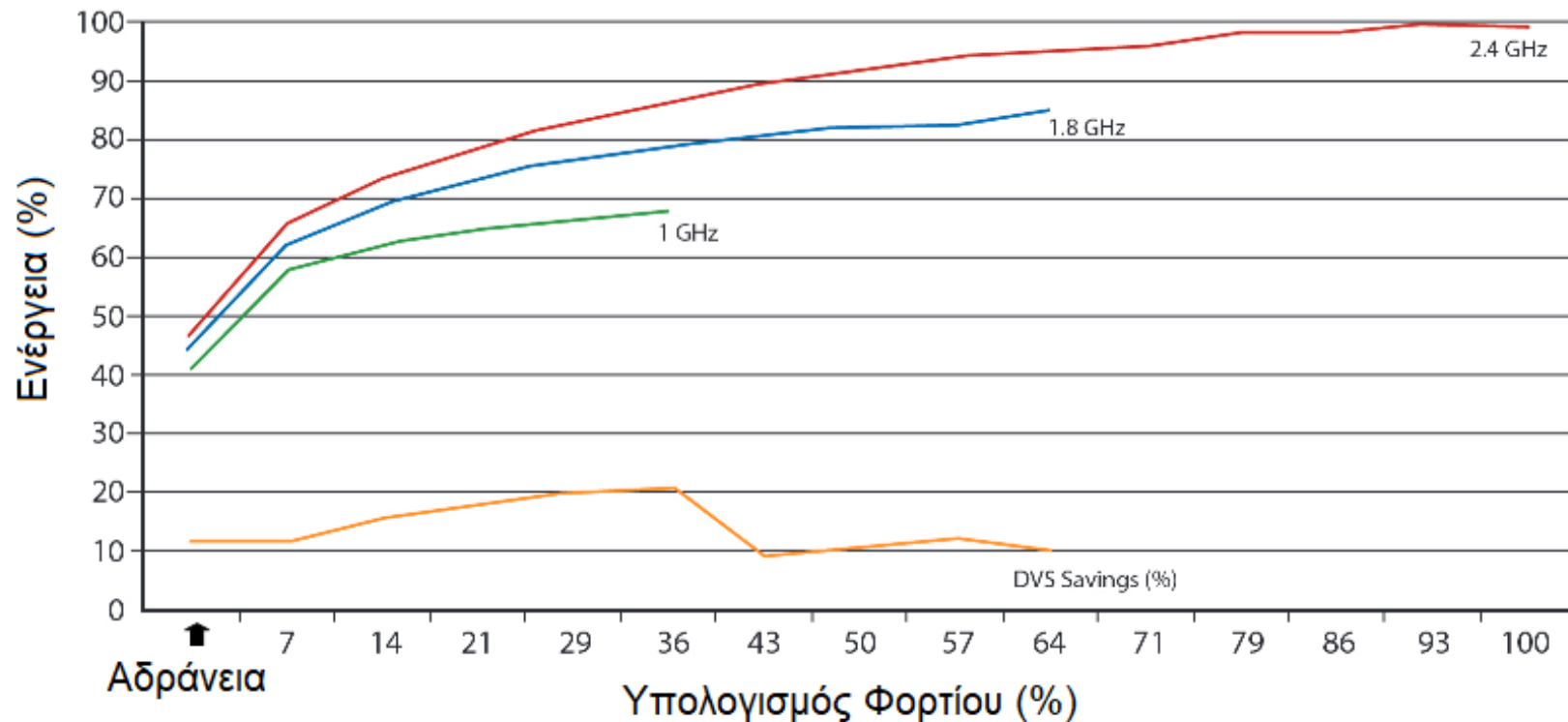


- Η μελέτη διαπίστωσε επίσης ότι η ανάμειξη διαφορετικών φορτίων στο εσωτερικό ενός συμπλέγματος αύξησε τις δυνατότητες υπερκάλυψης της ισχύος, επειδή αυτό μειώνει την πιθανότητα συγχρονισμένων σημείων ισχύος σε όλες τις μηχανές.
- Σε μια πραγματική εξάπλωση, είναι εύκολο να καταλήξουμε σε μια ανεπαρκώς αξιοποιημένη εγκατάσταση ακόμα και όταν δίνουμε προσοχή στις σωστές τιμές ισχύος.



# Η ΤΑΣΗ ΣΤΗ ΧΡΗΣΗ ΕΝΕΡΓΕΙΑΣ ΤΟΥ SERVER

Ενώ στο παρελθόν η δυναμική κλιμάκωση τάσης και συχνότητας (DVFS), Dynamic Voltage and Frequency Scaling, ήταν ο κυρίαρχος μηχανισμός για τη διαχείριση της κατανάλωσης ενέργειας σε διακομιστές, σήμερα αντιμετωπίζουμε ένα διαφορετικό και πιο περίπλοκο σενάριο.



# ΧΡΗΣΗ ΤΗΣ ΑΠΟΘΗΚΕΥΣΗΣ ΕΝΕΡΓΕΙΑΣ ΓΙΑ ΔΙΑΧΕΙΡΙΣΗ ΙΣΧΥΟΣ (1/2)

---

- Πολλές πρόσφατες μελέτες έχουν προτείνει τη χρήση ενέργειας που αποθηκεύεται στα εφεδρικά συστήματα της εγκατάστασης (π.χ. μπαταρίες UPS) για τη βελτιστοποίηση της απόδοσης της εγκατάστασης ή τη μείωση του ενεργειακού κόστους.
- Κατά τη γνώμη μας, η πιο ελπιδοφόρα χρήση της αποθήκευσης ενέργειας στη διαχείριση ενέργειας συνίσταται στη διαχείριση των βραχέων αιχμών της ζήτησης.





# ΧΡΗΣΗ ΤΗΣ ΑΠΟΘΗΚΕΥΣΗΣ ΕΝΕΡΓΕΙΑΣ ΓΙΑ ΔΙΑΧΕΙΡΙΣΗ ΙΣΧΥΟΣ (2/2)

---

- Η ανάπτυξη ενός τέτοιου συστήματος είναι δύσκολη και ενδεχομένως δαπανηρή. Εκτός από την πολυπλοκότητα του ελέγχου, το πρόσθετο κόστος των συσσωρευτών μπορεί να είναι σημαντικό, αφού δεν μπορούμε απλώς να ξαναχρησιμοποιούμε την υπάρχουσα ικανότητα UPS για διαχείριση ενέργειας, καθώς κάτι τέτοιο θα καθιστούσε την εγκατάσταση πιο ευάλωτη σε μια διακοπή λειτουργίας.
- Κανένα τέτοιο σύστημα διαχείρισης ισχύος δεν έχει χρησιμοποιηθεί ακόμα στα συστήματα παραγωγής.



- Πρώτον, η ισχύς και η ενέργεια πρέπει να διαχειρίζονται καλύτερα ώστε να ελαχιστοποιούν το λειτουργικό κόστος.
- Δεύτερον, το σημερινό υλικό δεν προσαρμόζει με ευκολία τη χρήση ισχύος του σε μεταβαλλόμενες συνθήκες φόρτωσης και ως εκ τούτου η αποδοτικότητα ενός διακομιστή υποβαθμίζεται σοβαρά υπό ελαφρύ φορτίο.
- Τρίτον, η βελτιστοποίηση της ενέργειας είναι ένα περίπλοκο πρόβλημα από άκρο σε άκρο, απαιτώντας έναν περίπλοκο συντονισμό μεταξύ του υλικού, των λειτουργικών συστημάτων, των εικονικών μηχανών, των ενδιάμεσων λογισμικών, των εφαρμογών και των επιχειρήσεων.



- Το κόστος κατασκευής κέντρου δεδομένων ποικίλλει σημαντικά ανάλογα με το σχεδιασμό, το μέγεθος, την τοποθεσία και την επιθυμητή ταχύτητα κατασκευής.
- Ο χαρακτηρισμός του κόστους σε δολάρια ανά watt είναι κάτι λογικό για τα μεγαλύτερα κέντρα δεδομένων (όπου τα σταθερά κόστη ανεξάρτητα από το μέγεθος είναι ένα σχετικά μικρό τμήμα του συνολικού κόστους), επειδή όλα τα κύρια στοιχεία του κέντρου δεδομένων - ισχύς, ψύξη και χώρος - είναι σχεδόν γραμμικά με τα watt.



- Το κόστος εκφράζεται σε δολάρια ανά ρεύμα, δηλαδή ανά watt που μπορεί πραγματικά να χρησιμοποιηθεί από τον εξοπλισμό.
- Το μηνιαίο κόστος απόσβεσης (ή το κόστος απόσβεσης) που προκύπτει από το αρχικό κόστος κατασκευής εξαρτάται από τη διάρκεια απομείωσης της επένδυσης (η οποία σχετίζεται με την αναμενόμενη διάρκεια ζωής) και το επιτόκιο.
- Τα έξοδα διακομιστή υπολογίζονται παρομοίως, εκτός από το ότι οι διακομιστές έχουν μικρότερη διάρκεια ζωής και κατά συνέπεια αποσβένονται συνήθως σε διάστημα 3-4 ετών.



## ΛΕΙΤΟΥΡΓΙΚΑ ΚΟΣΤΗ (1/2)

- Το κέντρο δεδομένων Orex είναι πιο δύσκολο να χαρακτηριστεί γιατί εξαρτάται σε μεγάλο βαθμό από τα επιχειρησιακά πρότυπα καθώς και από το μέγεθος του κέντρου δεδομένων (τα μεγαλύτερα κέντρα δεδομένων είναι φθηνότερα, διότι οι δαπάνες αποσβένονται καλύτερα). Το κόστος μπορεί επίσης να ποικίλει ανάλογα με τη γεωγραφική θέση (κλίμα, φόροι, επίπεδα μισθών κλπ.) και το σχεδιασμό και την ηλικία του κέντρου δεδομένων.
- Οι διακομιστές έχουν λειτουργικό κόστος. Επειδή εστιάζουμε μόνο στο κόστος λειτουργίας της ίδιας της υποδομής, θα επικεντρωθούμε μόνο στη συντήρηση και στις επισκευές υλικού, καθώς και στο κόστος ηλεκτρικής ενέργειας.



## ΛΕΙΤΟΥΡΓΙΚΑ ΚΟΣΤΗ (2/2)

---

- Επίσης, στα παραδοσιακά περιβάλλοντα πληροφορικής, ο κύριος όγκος του κόστους λειτουργίας έγκειται στις εφαρμογές, δηλαδή στις άδειες λογισμικού και στο κόστος των διαχειριστών συστημάτων, των διαχειριστών βάσεων δεδομένων, των μηχανικών δικτύων κλπ.
- Εξαιρούμε αυτές τις δαπάνες εδώ, επειδή εστιάζουμε στο κόστος λειτουργίας της φυσικής υποδομής, αλλά και επειδή το κόστος εφαρμογής ποικίλλει σημαντικά ανάλογα με την κατάσταση.



## ΠΕΡΙΠΤΩΣΕΙΣ ΧΡΗΣΗΣ ΚΕΝΤΡΟΥ ΔΕΔΟΜΕΝΩΝ (1/3)

- Δεδομένου του μεγάλου αριθμού μεταβλητών που εμπλέκονται, είναι καλύτερο να επεξηγηθεί το φάσμα των συντελεστών κόστους εξετάζοντας ένα μικρό αριθμό μελέτης περιπτώσεων που αντιπροσωπεύουν διαφορετικά είδη ανάπτυξης.
- Θεωρούμε ένα τυπικό νέο κέντρο δεδομένων πολλών μεγαβάτ στις Ηνωμένες Πολιτείες (κάτι πιο κοντά στην ταξινόμηση της βαθμίδας 3 του Uptime Institute).
- Για αυτό το παράδειγμα επιλέξαμε ένα Dell PowerEdge R520 με 2 επεξεργαστές, 48 GB μνήμης RAM και τέσσερις δίσκους. Αυτός ο διακομιστής αντλεί 340 W και κοστίζει περίπου 7.700 δολάρια από το 2012. Οι υπόλοιπες παράμετροι της βασικής περίπτωσης επιλέχθηκαν ως εξής:



## ΠΕΡΙΠΤΩΣΕΙΣ ΧΡΗΣΗΣ ΚΕΝΤΡΟΥ ΔΕΔΟΜΕΝΩΝ (2/3)

---

- Το κόστος της ηλεκτρικής ενέργειας είναι το μέσο βιομηχανικό ποσοστό των ΗΠΑ 6,7 cents/kWh για το 2012.
- Το επιτόκιο που μια επιχείρηση πρέπει να πληρώσει για τα δάνεια είναι 8%, και χρηματοδοτούμε τους διακομιστές με τριετή δάνειο μόνο για τόκους.
- Το κόστος κατασκευής του κέντρου δεδομένων είναι 10\$/W που αποσβένεται σε 12 χρόνια.
- Το κέντρο δεδομένων OpeX είναι 0,04\$/W ανά μήνα.





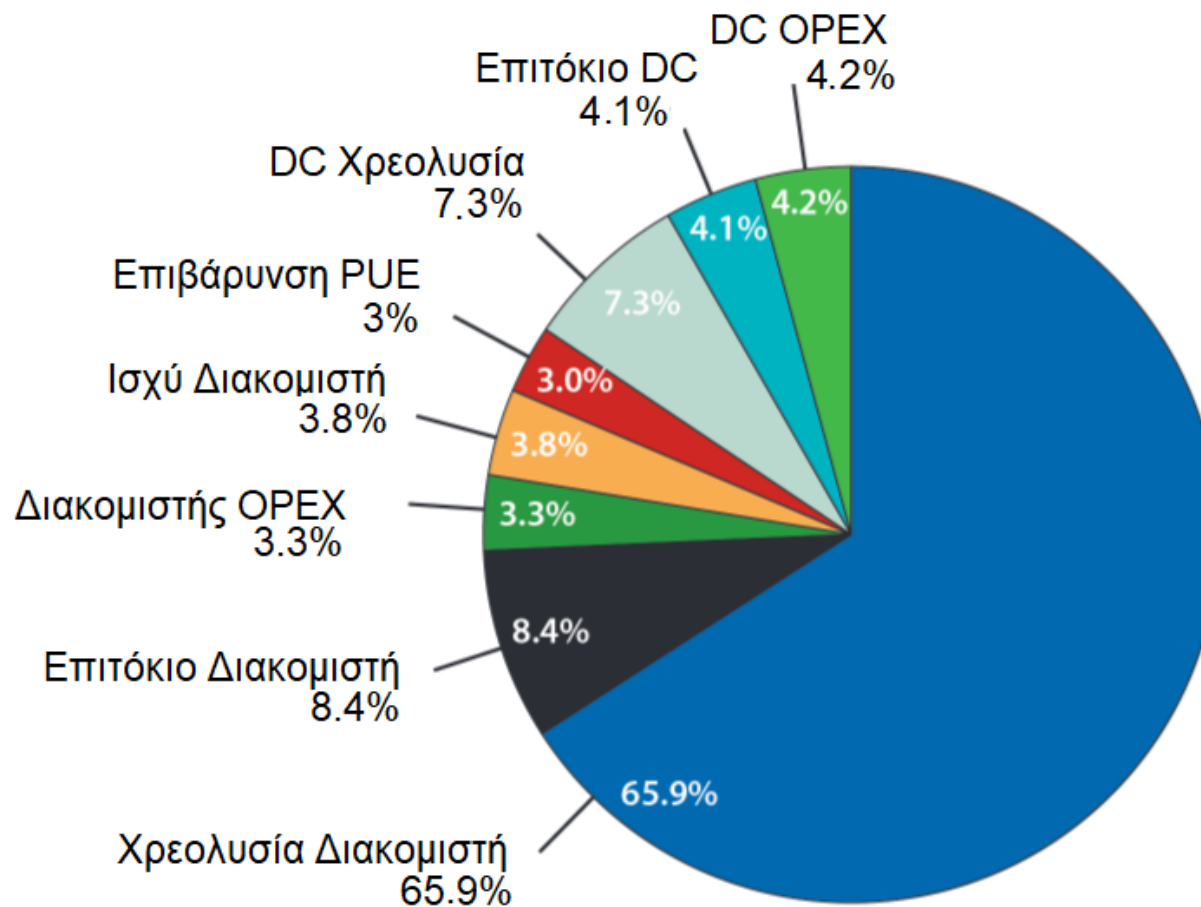
## ΠΕΡΙΠΤΩΣΕΙΣ ΧΡΗΣΗΣ ΚΕΝΤΡΟΥ ΔΕΔΟΜΕΝΩΝ (3/3)

---

- Το κέντρο δεδομένων έχει αποτελεσματικότητα χρήσης ενέργειας (PUE) 1,8, τον τρέχοντα μέσο όρο της βιομηχανίας. (Χρησιμοποιήστε το υπολογιστικό φύλλο στη διεύθυνση <http://goo.gl/eb6Ui> για να επαναπροσδιορίσετε τα παραδείγματα με ένα PUE της κατηγορίας 1.1 που είναι πιο τυπικό για τα WSC της Google.)
- Η διάρκεια ζωής του διακομιστή είναι 3 έτη και η επισκευή και συντήρηση του είναι 5% του Capex ανά έτος.
- Η μέση ισχύς του διακομιστή είναι 75% της μέγιστης ισχύος.



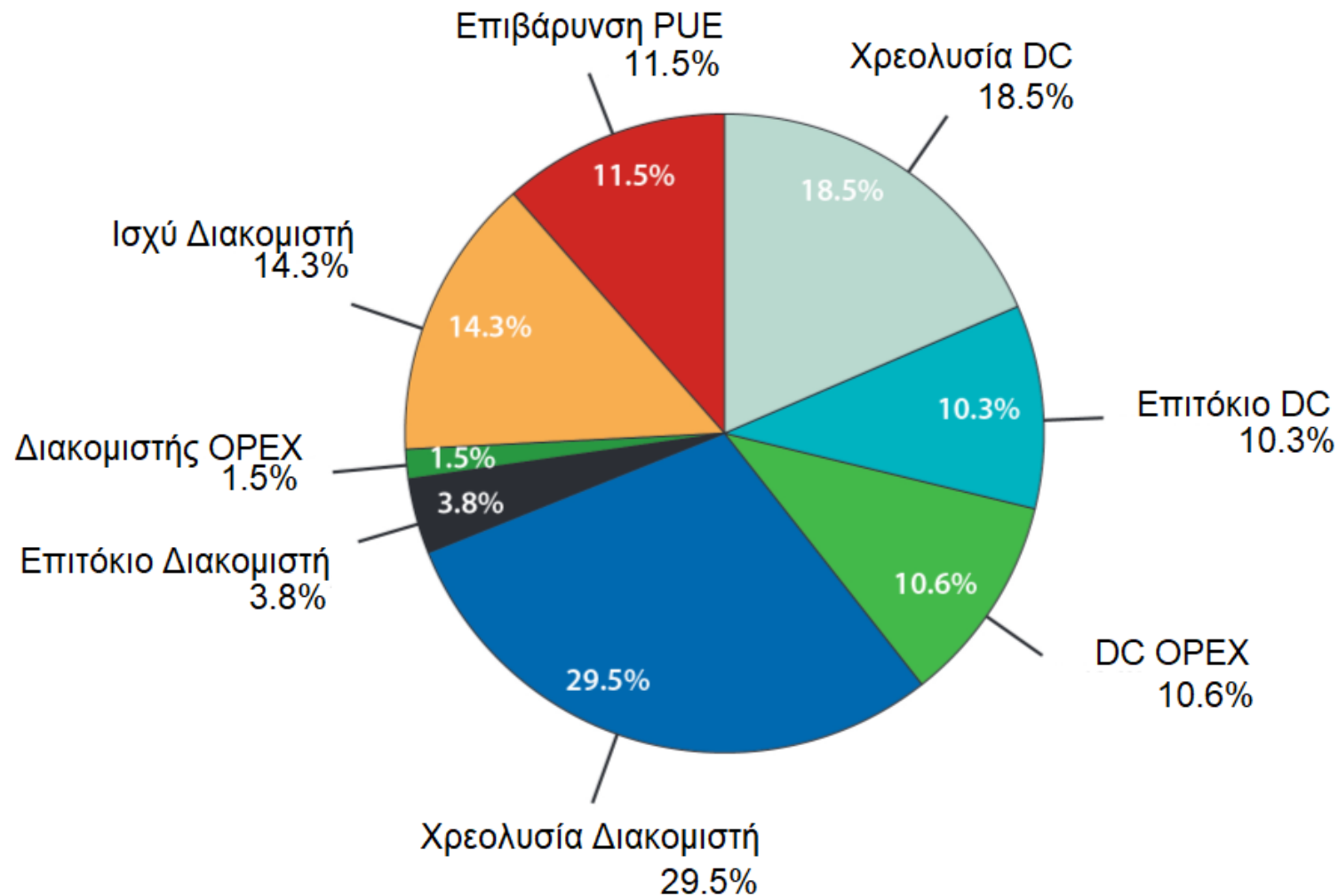
# ΚΑΤΑΝΟΜΗ ΚΟΣΤΟΥΣ ΠΕΡΙΠΤΩΣΗΣ Α



Ανάλυση του συνολικού κόστους ιδιοκτησίας (TCO, Total Cost of Ownership) για την περίπτωση ενός κλασικού κέντρου δεδομένων.



# ΚΑΤΑΝΟΜΗ ΚΟΣΤΟΥΣ ΠΕΡΙΠΤΩΣΗΣ Β



Ανάλυση του συνολικού κόστους ιδιοκτησίας (TCO, Total Cost of Ownership) για την περίπτωση ενός κέντρου δεδομένων με server υψηλής ισχύος.



# ΤΟ ΚΟΣΤΟΣ ΚΕΝΤΡΟΥ ΔΕΔΟΜΕΝΩΝ ΣΤΗΝ ΠΡΑΓΜΑΤΙΚΟΤΗΤΑ (1/2)

---

- Στην πραγματικότητα, το κόστος κέντρων δεδομένων είναι ακόμη υψηλότερο.
- Όλα τα μοντέλα που παρουσιάζονται μέχρι στιγμής υποθέτουν ότι το κέντρο δεδομένων είναι 100% πλήρες και ότι οι διακομιστές είναι αρκετά απασχολημένοι (το 75% της μέγιστης ισχύος αντιστοιχεί σε χρήση CPU περίπου 50%). Στην πραγματικότητα, αυτό συχνά δεν συμβαίνει.

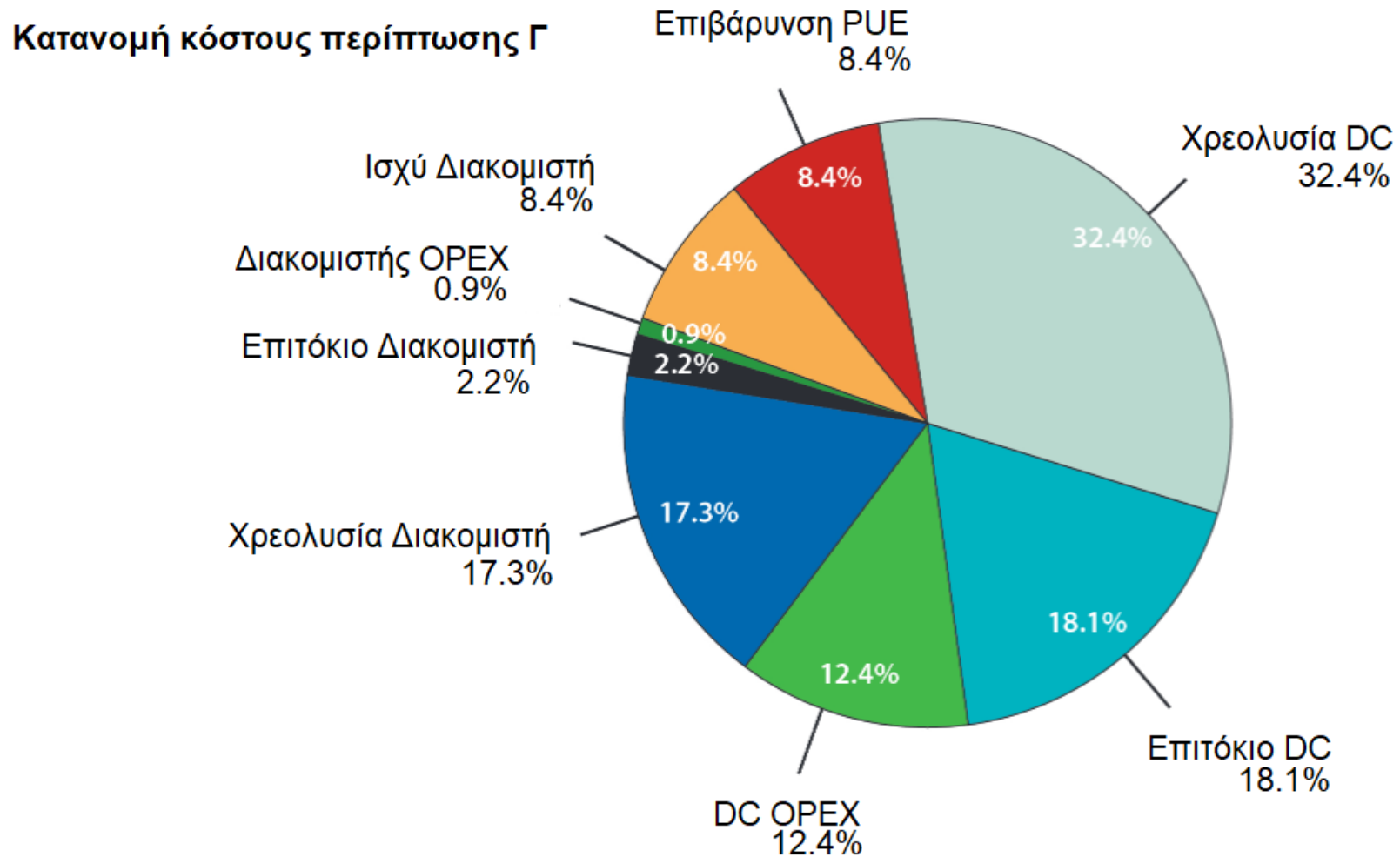


## ΤΟ ΚΟΣΤΟΣ ΚΕΝΤΡΟΥ ΔΕΔΟΜΕΝΩΝ ΣΤΗΝ ΠΡΑΓΜΑΤΙΚΟΤΗΤΑ (2/2)

- Η επίτευξη υψηλής αξιοποίησης της ισχύος του κέντρου δεδομένων δεν είναι τόσο απλή όσο φαίνεται. Ακόμη και αν ο προμηθευτής παρέχει μια αριθμομηχανή ισχύος για να υπολογίσει την πραγματική μέγιστη ισχύ για μια συγκεκριμένη διαμόρφωση, αυτή η τιμή θα υποθέσει τη χρησιμοποίηση 100% της CPU. Εάν εγκαταστήσουμε διακομιστές που βασίζονται σε αυτή την τιμή και λειτουργούν με μέσο όρο CPU μόνο 30% (καταναλώνουν 200 W αντί 300 W), αφήσαμε μόλις 30% της χωρητικότητας του κέντρου δεδομένων.
- Από την άλλη πλευρά, αν εγκαταστήσουμε με βάση τη μέση τιμή των 200 W και στο τέλος του μήνα, οι διακομιστές τρέχουν πραγματικά για σχεδόν πλήρη χωρητικότητα για κάποιο χρονικό διάστημα, το κέντρο δεδομένων μας θα υπερθερμανθεί ή θα βλάψει έναν μεταγωγέα.



# ΜΕΘΟΔΟΣ ΜΕΡΙΚΗΣ ΠΛΗΡΟΤΗΤΑΣ ΚΕΝΤΡΟΥ ΔΕΔΟΜΕΝΩΝ (1/2)



Ανάλυση του συνολικού κόστους ιδιοκτησίας (TCO, Total Cost of Ownership) για την περίπτωση ενός εν μέρη πλήρους κέντρου δεδομένων (50% εκμετάλλευση).



# ΜΕΘΟΔΟΣ ΜΕΡΙΚΗΣ ΠΛΗΡΟΤΗΤΑΣ ΚΕΝΤΡΟΥ ΔΕΔΟΜΕΝΩΝ (2/2)

- Οι μερικώς χρησιμοποιούμενοι διακομιστές επηρεάζουν επίσης το λειτουργικό κόστος με θετικό τρόπο, επειδή ο διακομιστής χρησιμοποιεί λιγότερη ενέργεια.
- Φυσικά, αυτές οι εξοικονομήσεις είναι αμφισβητήσιμες επειδή οι εφαρμογές που εκτελούνται σε αυτούς τους διακομιστές είναι πιθανό να παράγουν λιγότερη αξία.
- Το μοντέλο TCO (Total Cost of Ownership) δεν μπορεί να καταγράψει αυτό το αποτέλεσμα επειδή βασίζεται μόνο στο κόστος της φυσικής υποδομής και αποκλείει την εφαρμογή που εκτελείται σε αυτό το υλικό. Για να μετρήσουμε αυτήν την απόδοση από άκρο σε άκρο, μπορούμε να μετρήσουμε ένα διακομιστή μεσολάβησης για την τιμή εφαρμογής (π.χ. τον αριθμό των ολοκληρωμένων τραπεζικών συναλλαγών ή τον αριθμό των αναζητήσεων στο Web).



## ΤΟ ΚΟΣΤΟΣ ΤΩΝ ΔΗΜΟΣΙΩΝ ΣΥΝΝΕΦΩΝ (1/2)

- Αντί να δημιουργήσετε το δικό σας κέντρο δεδομένων και διακομιστή, μπορείτε να νοικιάσετε μια εικονική μηχανή από έναν δημόσιο προμηθευτή σύννεφων, όπως το Compute Engine της Google ή το EC2 του Amazon.
- Η τοπική τιμολόγηση είναι “pay as you go” τιμολόγηση - μπορείτε να ξεκινήσετε και να σταματήσετε ένα VM ανά πάσα στιγμή, οπότε αν χρειαστείτε ένα μόνο για λίγες μέρες το χρόνο, η τιμολόγηση κατά παραγγελία θα είναι πολύ φθηνότερη από οποιαδήποτε άλλη εναλλακτική λύση.
- Εάν χρειάζεστε ένα διακομιστή για μεγάλο χρονικό διάστημα, οι δημόσιοι προμηθευτές σύννεφων θα μειώσουν την ωριαία τιμή σε αντάλλαγμα για μια μακροπρόθεσμη δέσμευση καθώς και μια προκαταβολή.





- Πολλά από τα λειτουργικά έξοδα είναι σχετικά ανεξάρτητα από το μέγεθος του κέντρου δεδομένων: εάν θέλετε έναν υπεύθυνο ασφαλείας ή έναν τεχνικό εγκαταστάσεων 24x7, το κόστος θα είναι το ίδιο είτε το σύστημα είναι 1 MW είτε 5 MW.
- Επιπλέον, τα κεφαλαιουχικά έξοδα ενός παρόχου σύννεφων για διακομιστές και κτίρια είναι πιθανόν χαμηλότερα από δικά σας, αφού αγοράζουν (και χτίζουν) σε μεγάλη κλίμακα.
- Η Google, για παράδειγμα, σχεδιάζει τους δικούς της διακομιστές και κέντρα δεδομένων για να μειώσει το κόστος.



- Επιπτώσεις από το λογισμικό:
  - Όσο το δυνατόν περισσότερο, θα πρέπει να προσπαθήσουμε να εφαρμόσουμε ένα στρώμα υποδομής λογισμικού ανεκτικό σε σφάλματα, το οποίο μπορεί να κρύψει ένα μεγάλο μέρος αυτής της πολυπλοκότητας αποτυχίας από το λογισμικό σε επίπεδο εφαρμογής.
  - Μόλις τα ελαττώματα υλικού είναι ανεκτά χωρίς αδικαιολόγητη διακοπή μιας υπηρεσίας, οι αρχιτέκτονες υπολογιστών έχουν κάποια ελευθερία να επιλέξουν το επίπεδο αξιοπιστίας του υλικού που μεγιστοποιεί τη συνολική αποδοτικότητα του συστήματος. Αυτό το περιθώριο επιτρέπει να εξεταστεί, για παράδειγμα, η χρήση μη δαπανηρού υλικού PC-class για μια πλατφόρμα διακομιστή αντί για mainframe-class υπολογιστές.



- Η βασική ιδιότητα που εκμεταλλεύεται εδώ είναι ότι σε αντίθεση με τις παραδοσιακές ρυθμίσεις διακομιστή, δεν είναι πλέον απαραίτητο να διατηρείται ένας εξυπηρετητής σε λειτουργία με κάθε κόστος. Αυτή η απλή μετατόπιση απαιτήσεων επηρεάζει σχεδόν κάθε πτυχή της ανάπτυξης, από τη σχεδίαση του μηχανήματος/του κέντρου δεδομένων έως τις λειτουργίες, επιτρέποντας συχνά ευκαιρίες βελτιστοποίησης που διαφορετικά δεν θα υπήρχαν στο τραπέζι.
- Ένα άλλο χρήσιμο παράδειγμα αφορά τα σχεδιαστικά αντισταθμίσματα για ένα αξιόπιστο σύστημα αποθήκευσης. Μία εναλλακτική λύση είναι η δημιουργία πολύ αξιόπιστων κόμβων αποθήκευσης μέσω της χρήσης μονάδων δίσκου πολλαπλών δίσκων σε μια κατοπτρική ή RAIDed διαμόρφωση, έτσι ώστε να μπορούν να διορθωθούν μερικά λάθη δίσκου.



## ΑΝΤΙΜΕΤΩΠΙΣΗ ΑΠΟΤΥΧΙΩΝ ΚΑΙ ΕΠΙΣΚΕΥΩΝ (3/3)

- Σε ένα σύστημα που μπορεί να ανεχτεί πολλές αποτυχίες σε επίπεδο λογισμικού, η ελάχιστη απαίτηση που έχει επιβληθεί στο επίπεδο υλικού είναι ότι τα ελαττώματά του εντοπίζονται πάντοτε και αναφέρονται στο λογισμικό εγκαίρως, ώστε να επιτρέπουν στην υποδομή λογισμικού να το συγκρατεί και να λαμβάνει κατάλληλες ενέργειες ανάκτησης.
- Ωστόσο, η χαλάρωση της απαίτησης ανίχνευσης σφαλμάτων υλικού θα ήταν πολύ πιο δύσκολη, διότι σημαίνει ότι κάθε στοιχείο λογισμικού θα επιβαρυνόταν από την ανάγκη να ελέγξει τη σωστή εκτέλεσή του. Σε ένα πρώιμο σημείο της ιστορίας της, η Google είχε να αντιμετωπίσει διακομιστές που είχαν DRAM που δεν είχαν ακόμη έλεγχο της ιστοιμίας. Η παραγωγή ενός ευρετηρίου αναζήτησης ιστού αποτελείται ουσιαστικά από μια πολύ μεγάλη λειτουργία ταξινόμησης ανακατεύθυνσης/συγχώνευσης, χρησιμοποιώντας πολλά μηχανήματα για μεγάλο χρονικό διάστημα.



## ΒΑΘΜΟΣ ΑΥΣΤΗΡΟΤΗΤΑΣ ΑΠΟΤΥΧΙΩΝ (1/2)

- Κατατάσσουμε ευρέως τις αποτυχίες σε επίπεδο υπηρεσιών στις ακόλουθες κατηγορίες, οι οποίες αναφέρονται με μειωμένο βαθμό αυστηρότητας:
  - Αλλοιωμένη: δεσμευμένα δεδομένα που είναι αδύνατον να αναγεννηθούν, χαθούν ή αλλοιωθούν.
  - Ανέφικτη: η υπηρεσία είναι εκτός λειτουργίας ή αλλιώς δεν είναι δυνατή από τους χρήστες.
  - Υποβαθμισμένη: η υπηρεσία είναι διαθέσιμη αλλά σε κάποια υποβαθμισμένη λειτουργία.
  - Αποτυχία με χρήση μασκών: εμφανίζονται σφάλματα, αλλά αποκρύπτονται εντελώς από τους χρήστες από μηχανισμούς λογισμικού/υλικού ανθεκτικότητας σε σφάλματα.



## ΒΑΘΜΟΣ ΑΥΣΤΗΡΟΤΗΤΑΣ ΑΠΟΤΥΧΙΩΝ (2/2)

- Η διαθεσιμότητα των υπηρεσιών είναι πολύ σημαντική, κυρίως επειδή τα έσοδα από τις υπηρεσίες Internet συχνά σχετίζονται κατά κάποιο τρόπο με τον όγκο της κυκλοφορίας.
- Η μέτρηση της διαθεσιμότητας της υπηρεσίας σε απόλυτο χρόνο είναι λιγότερο χρήσιμη για τις υπηρεσίες Διαδικτύου που συνήθως εμφανίζουν μεγάλες διακυμάνσεις της ημερήσιας, εβδομαδιαίας και εποχικής κυκλοφορίας. Μια πιο κατάλληλη μέτρηση διαθεσιμότητας είναι το κλάσμα των αιτημάτων που ικανοποιείται από την υπηρεσία.
- Τέλος, μια ιδιαίτερα ζημιογόνος κλάση αποτυχιών είναι η απώλεια ή η αλλοίωση των δεσμευμένων ενημερώσεων σε κρίσιμα δεδομένα, ιδιαίτερα τα δεδομένα χρηστών, τα κρίσιμα επιχειρησιακά ημερολόγια ή τα σχετικά δεδομένα που είναι δύσκολο ή αδύνατο να αναγεννηθούν.



## ΑΙΤΙΕΣ ΣΦΑΛΜΑΤΩΝ ΥΠΗΡΕΣΙΩΝ (1/2)

---

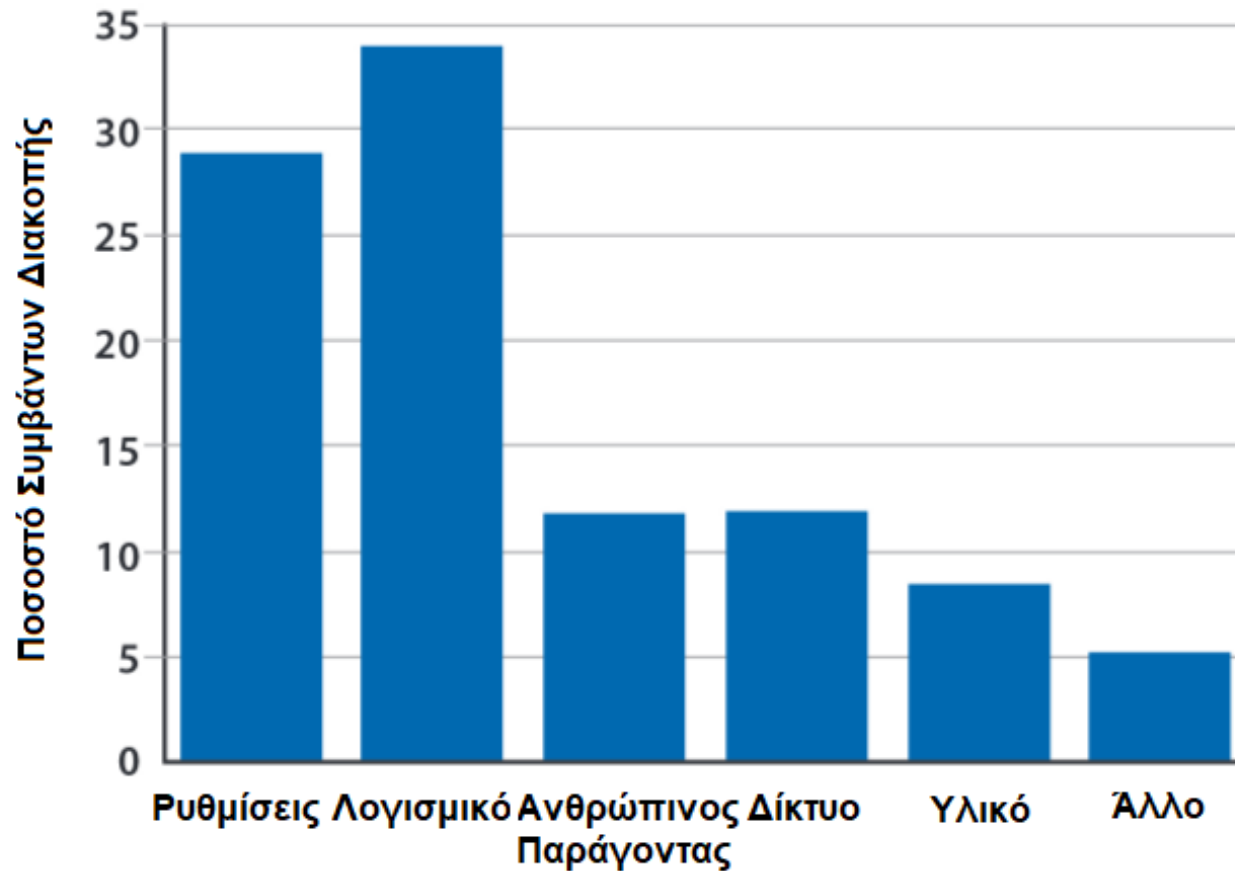
- Τα δεδομένα του Orpenheimer είναι κάπως συνεπή με το σημαντικό έργο του Gray, το οποίο δεν εξετάζει τις υπηρεσίες Διαδικτύου, αλλά εξετάζει τα δεδομένα πεδίου από τους εξαιρετικά ανεκτικούς σε σφάλματα διακομιστές Tandem μεταξύ 1985 και 1990.
- Είναι κάπως απροσδόκητο, αρχικά, να παρατηρηθούν σφάλματα υλικού που συμβάλλουν σε τόσο λίγα γεγονότα διακοπών σε αυτά τα δύο πολύ διαφορετικά συστήματα.



## ΑΙΤΙΕΣ ΣΦΑΛΜΑΤΩΝ ΥΠΗΡΕΣΙΩΝ (2/2)

Κατανομή των συμβάντων διακοπής της υπηρεσίας από την πιο πιθανή αιτία σε μία από τις κύριες υπηρεσίες της Google, που συλλέχθηκαν για περίοδο έξι εβδομάδων από τον Robert Stroud της Google.

Η πιθανότερη αιτία είναι η δυσλειτουργία του λογισμικού.





## ΑΤΥΧΗΜΑΤΑ ΕΠΙΠΕΔΟΥ ΜΗΧΑΝΗΜΑΤΩΝ (1/3)

- Ένας σημαντικός παράγοντας για το σχεδιασμό κατανομημένων συστημάτων ανεκτικών σε σφάλματα είναι η κατανόηση της διαθεσιμότητας σε επίπεδο διακομιστή. Εδώ θεωρούμε ότι οι αποτυχίες σε επίπεδο μηχανής είναι όλες οι καταστάσεις που οδηγούν σε μείωση του διακομιστή, όποια και αν είναι η αιτία.
- Οι Schroeder και Gibson μελέτησαν στατιστικά στοιχεία αποτυχίας από συστήματα υπολογιστών υψηλής απόδοσης στο Εθνικό Εργαστήριο του Los Alamos. Αν και αυτές δεν είναι η τάξη των υπολογιστών που μας ενδιαφέρει εδώ, αποτελούνται από κόμβους που μοιάζουν με μεμονωμένους διακομιστές σε WSCs. Έτσι τα δεδομένα τους είναι συναφή με την κατανόηση των αποτυχιών σε επίπεδο μηχανής στο πλαίσιο μας.



## ΑΤΥΧΗΜΑΤΑ ΕΠΙΠΕΔΟΥ ΜΗΧΑΝΗΜΑΤΩΝ (2/3)

- Παρακάτω θα παρουσιαστούν τα στατιστικά στοιχεία αποτυχίας και διακοπής σε επίπεδο μηχανής της Google.
- Τα δεδομένα βασίζονται σε μια παρατήρηση 6 μηνών για όλα τα συμβάντα επανεκκίνησης του μηχανήματος και τον αντίστοιχο χρόνο διακοπής τους, όπου ο χρόνος διακοπής αντιστοιχεί σε ολόκληρο το χρονικό διάστημα όπου ένα μηχάνημα δεν είναι διαθέσιμο για υπηρεσία, ανεξάρτητα από την αιτία.
- Αυτές οι στατιστικές είναι σε όλες τις μηχανές της Google. Για παράδειγμα, περιλαμβάνουν μηχανήματα που βρίσκονται στον αγωγό επισκευών, προγραμματισμένο χρόνο αναμονής για αναβαθμίσεις, καθώς και όλα τα είδη συντριβών μηχανών.

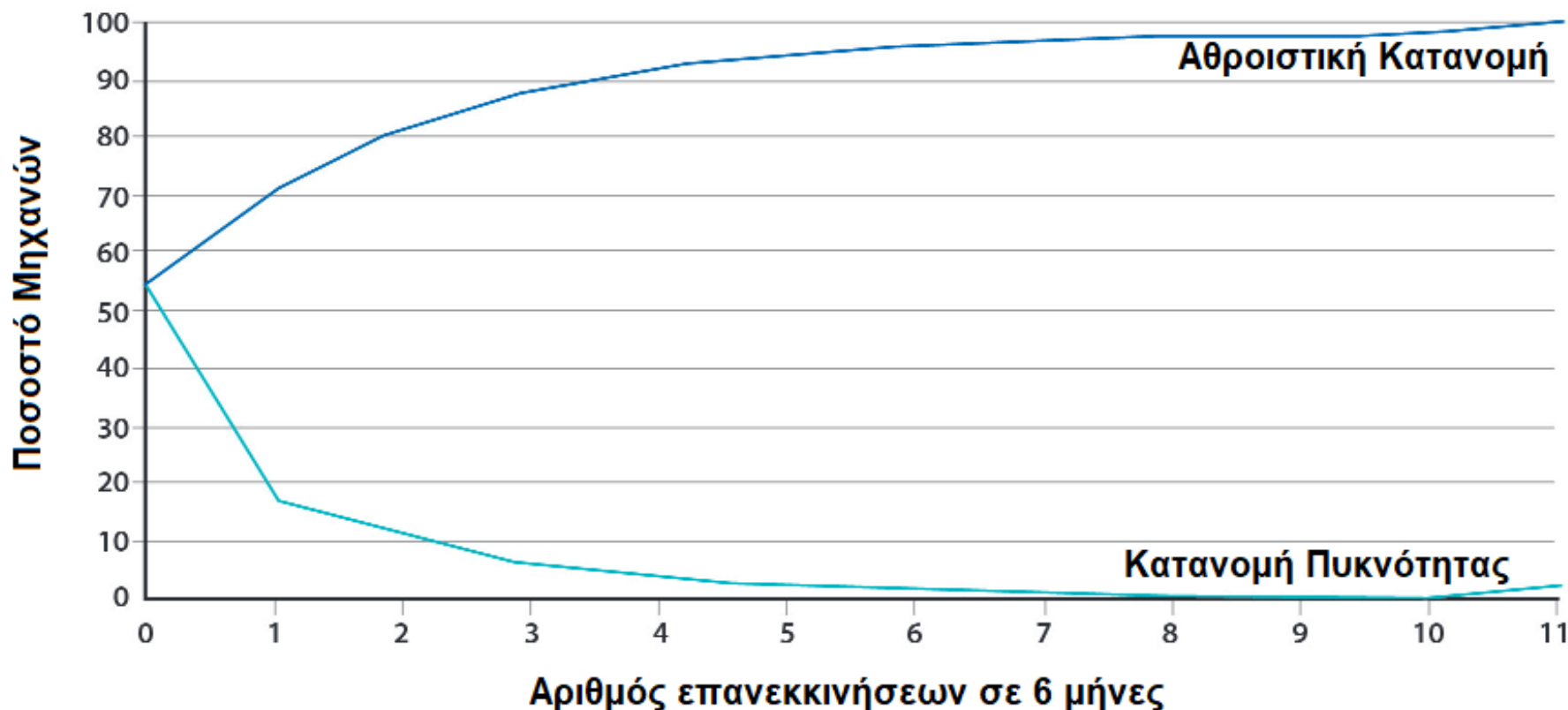


- Η εμπειρία μας στο Google είναι γενικά σύμφωνη με την ταξινόμηση του Orpenheimer, ακόμα κι αν οι ορισμοί των κατηγοριών δεν είναι απολύτως συνεπείς. Θα παρουσιαστεί παρακάτω μια ακατέργαστη ταξινόμηση όλων των συμβάντων που αντιστοιχούσαν σε αισθητές διαταραχές σε επίπεδο υπηρεσιών σε μία από τις μεγάλες ηλεκτρονικές υπηρεσίες της Google.
- Τέλος, μια ιδιαίτερα ζημιογόνος κλάση αποτυχιών είναι η απώλεια ή η αλλοίωση των δεσμευμένων ενημερώσεων σε κρίσιμα δεδομένα, ιδιαίτερα τα δεδομένα χρηστών, τα κρίσιμα επιχειρησιακά ημερολόγια ή τα σχετικά δεδομένα που είναι δύσκολο ή αδύνατο να αναγεννηθούν.



# ΕΠΑΝΕΚΚΙΝΗΣΗ ΜΗΧΑΝΩΝ ΤΗΣ GOOGLE

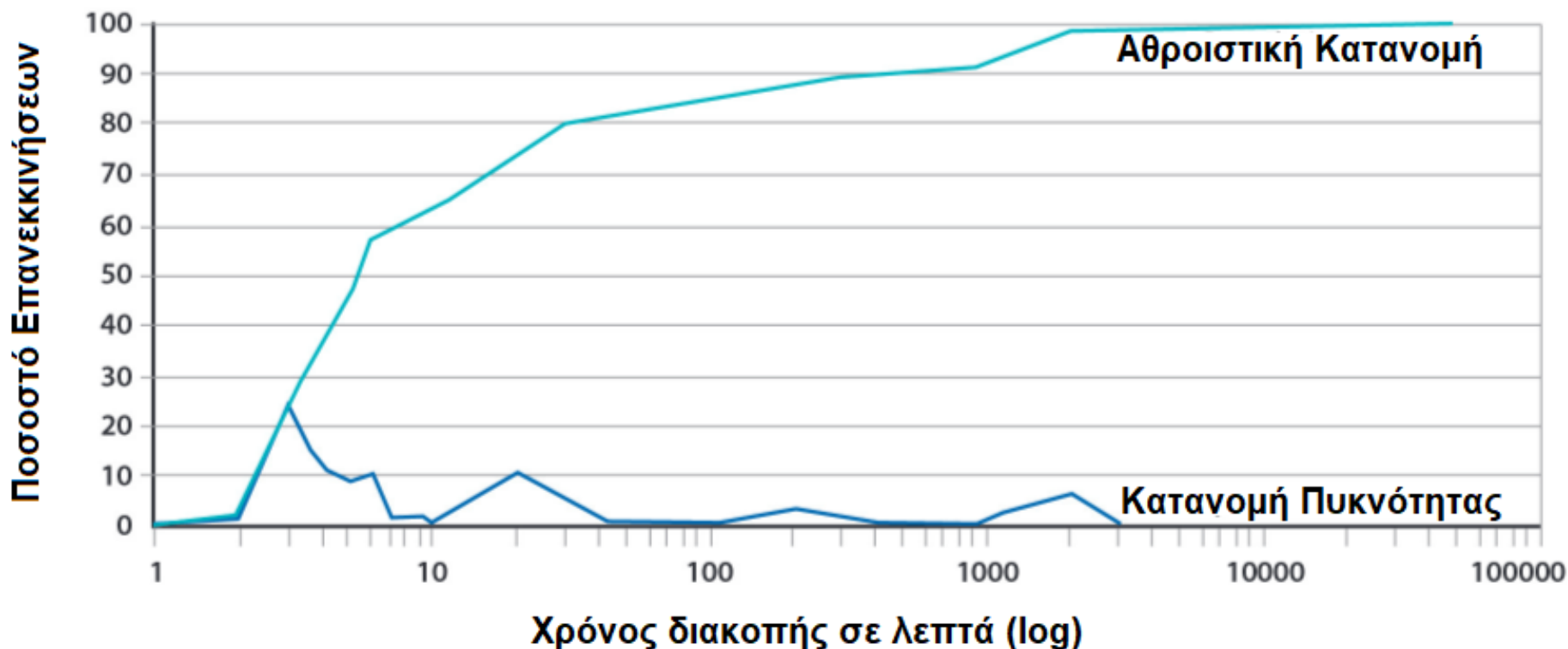
Η κατανομή των επανεκκινήσεων των μηχανών της Google για 6 μήνες



Οι περισσότεροι από τους μισούς διακομιστές είναι σε λειτουργία όλο το διάστημα παρατήρησης, και περισσότερο από το 95% των μηχανών επανεκκινούνται λιγότερο από μία φορά τον μήνα.



# ΚΑΤΑΝΟΜΗ ΧΡΟΝΟΥ ΔΙΑΚΟΠΗΣ ΜΗΧΑΝΗΜΑΤΩΝ



Κατανομή του χρόνου διακοπής της μηχανής, που παρατηρήθηκε στο Google για 6 μήνες. Ο μέσος ετήσιος ρυθμός επανεκκίνησης σε όλες τις μηχανές είναι 4,2, που αντιστοιχεί σε ένα μέσο χρόνο μεταξύ των επανεκκινήσεων μόλις λιγότερο από 3 μήνες.



- **Απλά σφάλματα DRAM.** Παρόλο που υπάρχουν λίγα διαθέσιμα δεδομένα πεδίου σχετικά με αυτό το θέμα, πιστεύεται γενικά ότι τα ποσοστά απλών σφαλμάτων DRAM είναι εξαιρετικά χαμηλά όταν χρησιμοποιούνται σύγχρονα ECC (Error Correction Coding).
- **Σφάλματα δίσκου.** Μελέτες που βασίζονται σε δεδομένα από συσκευές δικτύου/ Πανεπιστήμιο του Wisconsin, Carnegie Mellon και Google πρόσφατα έριξαν φως στα χαρακτηριστικά αποτυχίας των σύγχρονων μονάδων δίσκου.



## ΤΙ ΠΡΟΚΑΛΕΙ ΤΗΝ ΚΑΤΑΣΤΡΟΦΗ ΤΗΣ ΜΗΧΑΝΗΣ; (2/2)

---

- Οι αριθμοί υποδεικνύουν ότι ο μέσος όρος των μηχανών που συντρίβονται ετησίως λόγω βλαβών του υποσυστήματος δίσκου ή μνήμης θα πρέπει να είναι μικρότερο από το 10% όλων των μηχανών. Αντίθετα, παρατηρούμε ότι οι συγκρούσεις είναι πιο συχνές και ευρύτερα κατανεμημένες σε ολόκληρο τον πληθυσμό της μηχανής.
- Είναι σημαντικό να αναφέρουμε ότι ένα βασικό χαρακτηριστικό ενός καλά σχεδιασμένου λογισμικού ανθεκτικότητας σε σφάλματα είναι η ικανότητά του να επιβιώνει με μεμονωμένα σφάλματα είτε προκαλούνται από σφάλματα υλικού ή λογισμικού.



- Η ικανότητα πρόβλεψης μελλοντικών βλαβών του μηχανήματος ή των εξαρτημάτων εκτιμάται ιδιαίτερα, διότι θα μπορούσε να αποφευχθεί η πιθανή διακοπή των απρογραμμάτιστων διακοπών.
- Είναι προφανές ότι τα μοντέλα που μπορούν να προβλέψουν τα περισσότερα παραδείγματα μιας συγκεκριμένης κατηγορίας σφαλμάτων με πολύ χαμηλά ποσοστά αποτυχίας μπορεί να είναι πολύ χρήσιμα, ειδικά όταν αυτές οι προβλέψεις περιλαμβάνουν το άμεσο μέλλον.





## ΠΡΟΒΛΕΨΕΙΣ ΜΕΛΛΟΝΤΙΚΩΝ ΒΛΑΒΩΝ (2/2)

---

- Το Pinheiro et al περιγράφει μία από τις προσπάθειες της Google να δημιουργήσει μοντέλα πρόβλεψης για αποτυχίες μονάδας δίσκου με βάση τις παραμέτρους υγείας του δίσκου που είναι διαθέσιμες μέσω του προτύπου Τεχνολογικής Ανάλυσης Αυτοματισμού και Αναφοράς.
- Καταλήγουν στο συμπέρασμα ότι τέτοια μοντέλα είναι απίθανο να προβλέψουν τις περισσότερες αποτυχίες.
- Η γενική μας εμπειρία είναι ότι μόνο ένα μικρό υποσύνολο τάξεων αποτυχίας μπορεί να προβλεφθεί με αρκετά υψηλή ακρίβεια για να παραχθούν χρήσιμα επιχειρησιακά μοντέλα για WSCs.



## ΕΠΙΔΙΟΡΘΩΣΕΙΣ ΒΛΑΒΩΝ ΤΩΝ WSCs (1/2)

- Η αποτελεσματική διαδικασία επισκευής είναι κρίσιμη για τη συνολική αποδοτικότητα κόστους των WSCs.
- Όταν ένα μηχάνημα επισκευάζεται είναι ουσιαστικά εκτός λειτουργίας, οπότε όσο περισσότερο η μηχανή βρίσκεται σε επισκευές τόσο χαμηλότερη είναι η συνολική διαθεσιμότητα του συνόλου.
- Επίσης, οι δράσεις επιδιόρθωσης είναι δαπανηρές όσον αφορά τόσο τα ανταλλακτικά όσο και το εξειδικευμένο εργατικό δυναμικό.
- Τέλος, η ποιότητα επισκευής - πόσο πιθανό είναι ότι μια ενέργεια επιδιόρθωσης θα διορθώσει πραγματικά ένα πρόβλημα ενώ προσδιορίζει με ακρίβεια ποια (αν υπάρχει) συνιστώσα είναι σε σφάλμα - επηρεάζει τόσο τα έξοδα συνιστωσών όσο και τη μέση αξιοπιστία της μηχανής.

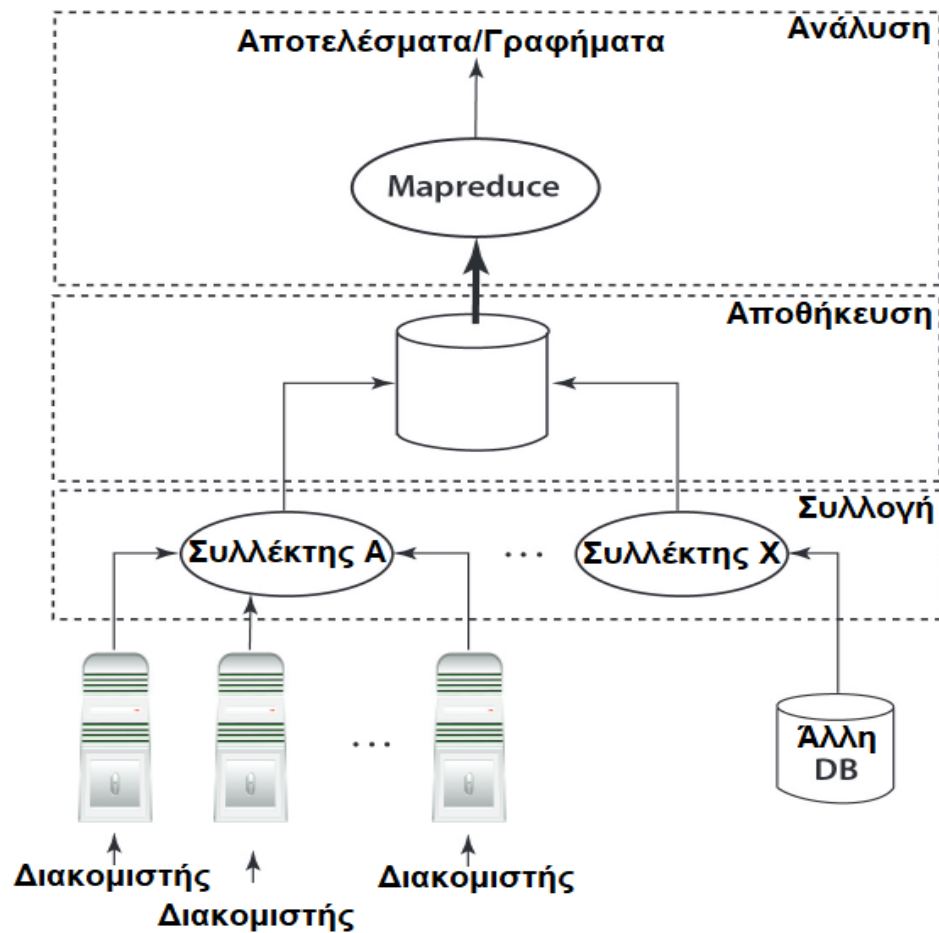


## ΕΠΙΔΙΟΡΘΩΣΕΙΣ ΒΛΑΒΩΝ ΤΩΝ WSCs (2/2)

- Υπάρχουν δύο χαρακτηριστικά των WSCs που επηρεάζουν άμεσα την αποδοτικότητα των επισκευών.
- Πρώτον, λόγω του μεγάλου αριθμού σχετικών διακομιστών χαμηλού επιπέδου και της ύπαρξης στρώματος αντοχής σφαλμάτων λογισμικού, δεν είναι τόσο κρίσιμο να υπάρχει άμεση ανταπόκριση σε μεμονωμένες περιπτώσεις επισκευής, επειδή είναι απίθανο να επηρεάσουν τη συνολική υγεία των υπηρεσιών.
- Αντί αυτού, ένα κέντρο δεδομένων μπορεί να εφαρμόσει ένα πρόγραμμα που κάνει την πιο αποτελεσματική χρήση του χρόνου ενός τεχνικού κάνοντας μια καθημερινή σάρωση όλων των μηχανών που χρειάζονται προσοχή. Η φιλοσοφία είναι να αυξηθεί ο ρυθμός επισκευών, διατηρώντας ταυτόχρονα την καθυστέρηση των επισκευών σε αποδεκτά επίπεδα.



# ΥΠΟΔΟΜΗ ΠΑΡΑΚΟΛΟΥΘΗΣΗΣ ΚΑΙ ΑΝΑΛΥΣΗΣ ΛΕΙΤΟΥΡΓΙΑΣ ΤΗΣ GOOGLE



Η υποδομή παρακολούθησης και ανάλυσης λειτουργίας της Google είναι ένα παράδειγμα ενός συστήματος παρακολούθησης που εκμεταλλεύεται αυτή την τεράστια πηγή δεδομένων.



# ΑΝΟΧΗ ΣΦΑΛΜΑΤΩΝ

- Η ικανότητα ενός καλά σχεδιασμένου λογισμικού αντοχής σε σφάλματα για την κάλυψη μεγάλου αριθμού αποτυχιών με σχετικά μικρή επίπτωση στις μετρήσεις σε επίπεδο υπηρεσίας θα μπορούσε να έχει απροσδόκητα επικίνδυνες παρενέργειες.
- Στις περιγραφές συστημάτων, χρησιμοποιούμε συχνά  $N$  για να δηλώσουμε τον αριθμό των εξυπηρετητών που απαιτούνται για την παροχή μιας υπηρεσίας με πλήρες φορτίο, οπότε το  $N+1$  περιγράφει ένα σύστημα με ένα επιπλέον αντίγραφο για πλεονασμό.



- Υλικό των WSCs
- Λογισμικό των WSCs
- Οικονομικά των WSCs



## ΥΛΙΚΟ ΤΩΝ WSCs (1/3)

---

- Τα δομικά στοιχεία επιλογής για WSCs είναι μηχανές κατηγορίας server-class, μονάδες δίσκων για καταναλωτές ή επιχειρήσεις και δίκτυα δικτύωσης βασισμένα σε Ethernet.
- Με γνώμονα τον όγκο αγορών εκατοντάδων εκατομμυρίων καταναλωτών και μικρών επιχειρήσεων, τα συστατικά των βασικών προϊόντων επωφελούνται από τις πιο μικρές σε κλίμακα κατασκευές και συνεπώς παρουσιάζουν σημαντικά καλύτερες αναλογίες τιμών/απόδοσης από ότι οι αντίστοιχες εταιρείες υψηλού επιπέδου.



- Η υψηλότερη αξιοπιστία του εξοπλισμού υψηλής τεχνολογίας είναι λιγότερο σημαντική σε αυτόν τον τομέα, επειδή απαιτείται στρώμα λογισμικού ανεκτικής βλάβης για την παροχή αξιόπιστης υπηρεσίας Διαδικτύου ανεξάρτητα από την ποιότητα του υλικού - σε συστοιχίες με δεκάδες χιλιάδες συστήματα, ακόμη και συστοιχίες με εξαιρετικά αξιόπιστους διακομιστές να αντιμετωπίζουν βλάβες υπερβολικά συχνά, ώστε το λογισμικό να αναλαμβάνει τη λειτουργία χωρίς προβλήματα.
- Επιπλέον, οι μεγάλες και σύνθετες υπηρεσίες Διαδικτύου αποτελούνται συχνά από πολλαπλές μονάδες ή στρώματα λογισμικού που δεν είναι bug-free και μπορούν να αποτύχουν σε ακόμη υψηλότερα ποσοστά από τα εξαρτήματα υλικού.





## ΥΛΙΚΟ ΤΩΝ WSCs (3/3)

---

- Η απόδοση της δικτύωσης και το υποσύστημα αποθήκευσης μπορεί να είναι πιο σχετικά με τους προγραμματιστές WSC από τα υποσυστήματα CPU και DRAM, αντίθετα με τα τυπικά συστήματα μικρότερης κλίμακας.
- Το σχετικά υψηλό κόστος (ανά gigabyte) μνήμης DRAM ή FLASH το καθιστά απαγορευτικά ακριβό για μεγάλα σύνολα δεδομένων ή σπάνια προσπελάσιμα δεδομένα. Ως εκ τούτου, οι δίσκοι εξακολουθούν να χρησιμοποιούνται περισσότερο.



## ΛΟΓΙΣΜΙΚΟ ΤΩΝ WSCs

- Τα WSCs είναι πιο πολύπλοκοι στόχοι προγραμματισμού από τα παραδοσιακά συστήματα υπολογιστών εξαιτίας της τεράστιας κλίμακας τους, της πολυπλοκότητας της αρχιτεκτονικής τους (όπως φαίνεται από τον προγραμματιστή) και της ανάγκης να ανεχθούν συχνές αποτυχίες.
- Η ανάπτυξη λογισμικού για υπηρεσίες Διαδικτύου διαφέρει επίσης από το παραδοσιακό μοντέλο πελάτη/διακομιστή με διάφορους τρόπους.
  - Ample parallelism
  - Workload churn
  - Platform homogeneity
  - Fault-free operation



## ΟΙΚΟΝΟΜΙΚΑ ΤΩΝ WSCs (1/3)

- Η απαράδεκτη ζήτηση για μεγαλύτερη αποδοτικότητα κόστους στον υπολογισμό έναντι υψηλότερων υπολογιστικών επιδόσεων κόστισε την πρωταρχική μετρική στο σχεδιασμό των συστημάτων WSCs.
- Επιπλέον, η αποδοτικότητα του κόστους πρέπει να οριστεί ευρέως για να ληφθούν υπόψη όλες οι σημαντικές συνιστώσες του κόστους, συμπεριλαμβανομένων των κεφαλαίων εγκατάστασης υποδομών και των λειτουργικών εξόδων (που περιλαμβάνουν την παροχή ηλεκτρικού ρεύματος και το ενεργειακό κόστος), το υλικό, το λογισμικό, το προσωπικό διαχείρισης και τις επισκευές.
- Τα κόστη ισχύς και ενέργειας είναι ιδιαίτερα σημαντικά για τα WSCs λόγω του μεγέθους τους.



- Τα χαρακτηριστικά χρήσης των WSCs, τα οποία δαπανούν λίγο χρόνο σε πλήρη αδράνεια ή σε πολύ υψηλά επίπεδα φορτίου, απαιτούν τα συστήματα και τα στοιχεία να είναι ενεργειακά αποδοτικά σε ένα ευρύ φάσμα φορτίου και ιδιαίτερα σε χαμηλά επίπεδα χρήσης. Η ενεργειακή απόδοση των εξυπηρετητών και των WSCs συχνά υπερεκτιμάται με τη χρήση σημείων αναφοράς τα οποία υποθέτουν επίπεδα απόδοσης κορυφαίας λειτουργίας.
- Επιπλέον, τα ίδια τα κέντρα δεδομένων δεν είναι ιδιαίτερα αποτελεσματικά. Η αποτελεσματικότητα της χρήσης του κτιρίου (PUE) είναι ο λόγος της συνολικής κατανάλωσης ενέργειας που διαιρείται με χρήσιμη ισχύ (διακομιστή). Για παράδειγμα, ένα κέντρο δεδομένων με PUE 2.0 χρησιμοποιεί 1 W ισχύ για κάθε ρεύμα ισχύος διακομιστή.



- Δυστυχώς, πολλές υπάρχουσες εγκαταστάσεις λειτουργούν σε PUEs 2 ή μεγαλύτερες, ενώ οι PUEs του 1,5 είναι σπάνιες. Είναι σαφές ότι υπάρχουν σημαντικές ευκαιρίες βελτίωσης της απόδοσης όχι μόνο σε επίπεδο διακομιστών αλλά και σε επίπεδο κτιρίων, όπως αποδείχθηκε από το ετήσιο 1,13 PUE της Google σε όλες τις προσαρμοσμένες εγκαταστάσεις της στα τέλη του 2012.
- Ωστόσο, το κόστος παροχής ενέργειας, δηλαδή το κόστος κατασκευής μιας εγκατάστασης ικανής να παρέχει και να ψύχει αποδοτικά, μπορεί να είναι ακόμη σημαντικότερο από το ίδιο το κόστος ηλεκτρικής ενέργειας.



# ΣΥΜΠΕΡΑΣΜΑΤΑ

---

- Ο υπολογισμός μετακινείται στο σύννεφο και επομένως σε WSC.
- Οι αρχιτέκτονες λογισμικού και υλικού πρέπει να γνωρίζουν τα συστήματα από άκρο σε άκρο για να σχεδιάζουν καλές λύσεις.
- Δεν σχεδιάζουμε πλέον ξεχωριστά εφαρμογές ενός διακομιστή και δεν μπορούμε πλέον να αγνοούμε τους φυσικούς και οικονομικούς μηχανισμούς που σχετίζονται με μια αποθήκη γεμάτη υπολογιστές.



---

Σας ευχαριστώ για την προσοχή σας

