



ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ
ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ ΚΑΙ ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ
ΤΟΜΕΑΣ ΤΕΧΝΟΛΟΓΙΑΣ ΠΛΗΡΟΦΟΡΙΚΗΣ ΚΑΙ ΥΠΟΛΟΓΙΣΤΩΝ
ΕΡΓΑΣΤΗΡΙΟ ΥΠΟΛΟΓΙΣΤΙΚΩΝ ΣΥΣΤΗΜΑΤΩΝ
www.cslab.ece.ntua.gr

Διπλωματική Εργασία

*Μελέτη και αξιολόγηση διαφορετικών σχημάτων διαμοιρασμού
ιεραρχίας μνήμης για αρχιτεκτονικές πολλαπλών πυρήνων*

Καθηγητής: Νεκτάριος Κοζύρης (nkoziris@cslab.ece.ntua.gr)
Επικοινωνία: Αναστόπουλος Νίκος (anastop@cslab.ece.ntua.gr)

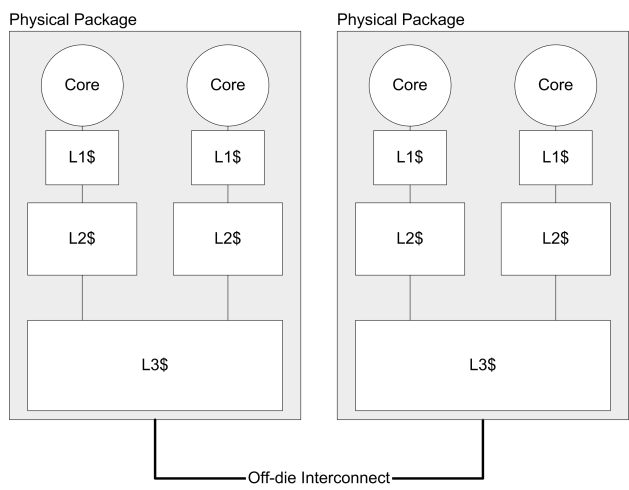
1. Περιγραφή

Η ενσωμάτωση πολλαπλών επεξεργαστικών πυρήνων σε ένα chip (multicores ή CMPs – Chip Multi-Processors) κερδίζει συνεχώς έδαφος στην κατεύθυνση της μεγιστοποίησης της απόδοσης των σύγχρονων επεξεργαστών. Σε συνδυασμό με τον εγγενή παραλληλισμό που μπορεί να εξαχθεί από την εκάστοτε εφαρμογή, εκτιμάται ότι μπορεί συνδράμει σημαντικά στη μείωση του χρόνου εκτέλεσης των εφαρμογών, περισσότερο από ό,τι άλλες παράμετροι που δεν επιδέχονται σε μεγάλο βαθμό περαιτέρω βελτιστοποιήσεις (π.χ. συχνότητα λειτουργίας ή μικροαρχιτεκτονικά χαρακτηριστικά όπως η υπερβαθμωτή και η out-of-order εκτέλεση εντολών ή η πρόβλεψη διακλαδώσεων). Αναμένεται ότι στα επόμενα χρόνια, ο αριθμός των πυρήνων που θα συμπεριλαμβάνονται σε ένα chip θα είναι της τάξης των δεκάδων ή ακόμα και των εκατοντάδων, οδηγώντας έτσι σε πολυ-πύρηνες αρχιτεκτονικές μεγάλης κλίμακας.

Ωστόσο, δεν είναι ακόμα σαφές ποιος είναι ο πλέον ενδεδειγμένος τρόπος διαμοιρασμού της ιεραρχίας μνήμης, σε ένα περιβάλλον όπου μεγάλος αριθμός νημάτων μίας παράλληλης εφαρμογής εκτελούνται υπό έναν κοινό χώρο διευθύνσεων, και μπορεί να μοιράζονται ή και να ανταγωνίζονται για τα ίδια δεδομένα.

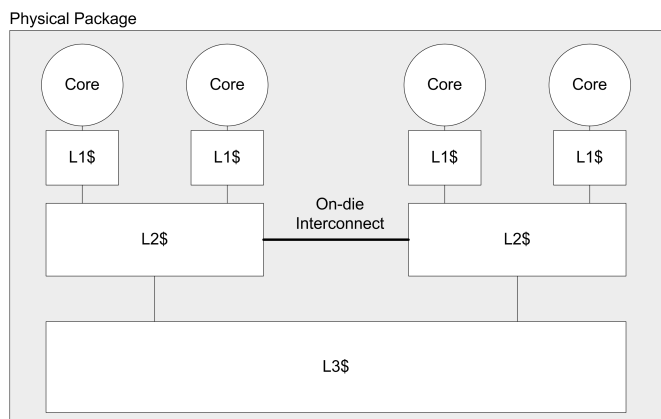
Από τη μια πλευρά, έχουμε διαμορφώσεις όπου υπάρχει χαμηλός (ή καθόλου) διαμοιρασμός, και κάθε πυρήνας έχει δικιά του τα υψηλότερα επίπεδα κρυφής μνήμης (Σχήμα 1). Αυτή είναι και η πολιτική που ακολούθηθηκε στις πρώτες υλοποιήσεις πολυπύρηνων επεξεργαστών που έκαναν την εμφάνισή τους τα τελευταία χρόνια. Το πλεονέκτημα τέτοιων σχημάτων είναι ότι κάθε πυρήνας έχει για αποκλειστική χρήση πολλαπλά επίπεδα κρυφής μνήμης, χωρίς να υπάρχει ο κίνδυνος τα δεδομένα πάνω στα οποία εργάζεται το αντίστοιχο νήμα να εκτοπιστούν λόγω συγκρούσεων με δεδομένα νημάτων από άλλους πυρήνες. Το μειονέκτημα είναι ότι η ενδο-επικοινωνία μεταξύ των νημάτων, για τους σκοπούς της ανταλλαγής δεδομένων ή του συγχρονισμού, πρέπει να περάσει από τα χαμηλότερα επίπεδα στην ιεραρχία (L3 cache ή κύρια μνήμη), με αποτέλεσμα να γίνεται με μεγάλη σχετικά

καθυστερήση (latency) και μικρό ρυθμό (bandwidth). Σε αυτά, προστίθεται και η επιπλέον επιβάρυνση εξαιτίας της λειτουργίας του πρωτοκόλλου συνοχής δεδομένων (coherency protocol), από τη στιγμή που υπάρχουν πολλαπλά αντίγραφα των ίδιων δεδομένων σε διαφορετικές κρυφές μνήμες.



Σχήμα 1: σενάριο χαμηλού διαμοιρασμού ιεραρχίας μνήμης

Στην άλλη πλευρά, έχουμε διαμορφώσεις με υψηλό (ή πλήρη) διαμοιρασμό, με τους πυρήνες να μοιράζονται τα περισσότερα επίπεδα κρυφής μνήμης (Σχήμα 2). Σε αυτή την περίπτωση, η ενδο-επικοινωνία είναι αποτελεσματικότερη, εφόσον η διασύνδεση γίνεται στο εσωτερικό του chip (σε υψηλές συχνότητες, μικρές καθυστερήσεις και μεγάλο εύρος δεδομένων, δηλαδή), ενώ δεν υπάρχει επιβάρυνση λόγω κάποιου πρωτοκόλλου συνοχής. Από τη στιγμή όμως που υπάρχει διαμοιρασμός, σε κάθε πυρήνα αναλογεί και ουσία μικρότερο μέρος της κρυφής μνήμης σε σχέση με την πρώτη περίπτωση, ενώ υπάρχει και ο κίνδυνος συγκρούσεων μεταξύ των δεδομένων των νημάτων.



Σχήμα 2: σενάριο υψηλού διαμοιρασμού ιεραρχίας μνήμης

2. Ζητούμενο

Στην διπλωματική εργασία αυτή, ο φοιτητής καλείται κατά κάποιο τρόπο να απαντήσει στο εξής ερώτημα: δεδομένου συγκεκριμένου αριθμού πυρήνων και συγκεκριμένης χωρητικότητας για κάθε επίπεδο κρυφής μνήμης, ποιος είναι ο πιο αποδοτικός τρόπος να διαμοιραστεί και να

απεικονιστεί το κάθε επίπεδο σε αυτούς;

Θα διερευνηθούν διάφορα πιθανά σενάρια διάταξης και διαμοιρασμού της κρυφής μνήμης, στη λογική των όσων αναφέρθηκαν στα παραπάνω. Για κάθε σχήμα, θα μελετηθεί ο τρόπος που επηρεάζουν την απόδοση μια σειρά από βασικές παράμετροι λειτουργίας της κρυφής μνήμης, όπως η χωρητικότητα, η συσχετιστικότητα, ο χρόνος επικοινωνίας με τον επεξεργαστή ή με χαμηλότερα επίπεδα στην ιεραρχία, κ.λπ.

3. Λεπτομέρειες Υλοποίησης

Για τους σκοπούς της αξιολόγησης, θα χρησιμοποιηθεί ο προσομοιωτής M5. Ο M5 έχει τη δυνατότητα να κάνει πλήρη προσομοίωση ενός υπολογιστικού συστήματος (full system simulation), από τον επεξεργαστή και τις μνήμες, μέχρι τις διάφορες περιφερειακές συσκευές (αποθήκευση, δίκτυο, κ.λπ.). Είναι αρκετά ευέλικτος, από την άποψη ότι ο χρήστης έχει τη δυνατότητα να ορίσει τις βασικές παραμέτρους λειτουργίας του εκάστοτε συστατικού στοιχείου που προσομοιώνεται (π.χ. L1 cache), καθώς και τον τρόπο που διασυνδέονται τα επιμέρους συστατικά στοιχεία.

Σαν μετροπρογράμματα για την αξιολόγηση της απόδοσης των διαφόρων σχημάτων ιεραρχίας μνήμης, θα χρησιμοποιηθούν φορτία εργασίας από επιστημονικούς κώδικες που ακολουθούν το μοντέλο της παράλληλης εκτέλεσης σε μοιραζόμενη μνήμη. Για την ανάπτυξη τέτοιων εφαρμογών, θα χρησιμοποιηθούν “cross” GNU εργαλεία (compiler, assembler, linker, κ.λπ.) τα οποία θα παράγουν κώδικα και εκτελέσιμα για την αρχιτεκτονική του συστήματος που προσομοιώνεται (Alpha ISA).

4. Προαπαιτούμενα

- καλή γνώση αρχιτεκτονικής υπολογιστών
- καλή γνώση στις αρχές παράλληλου προγραμματισμού και ειδικότερα στον πολυνηματικό προγραμματισμό (σε γλώσσα C)
- βασικές γνώσεις σε Linux και στα GNU εργαλεία ανάπτυξης λογισμικού

5. Σχετική Βιβλιογραφία - Links

- η σελίδα του M5 στο πανεπιστήμιο του Michigan:
http://m5.eecs.umich.edu/wiki/index.php/Main_Page
- πρόσφατο άρθρο σχετικά με τον M5 στο IEEE Micro (July/August):
The M5 Simulator: Modeling Networked Systems. N. Binkert, R. Dreslinski, L. Hsu, K. Lim, A. Saidi and Steven Reinhardt.
- άρθρα πάνω σε επεξεργαστές που ενσωματώνουν πολλαπλές ροές εκτέλεσης και διαφορετικά σχήματα διαμοιρασμού της κρυφής μνήμης:
 - *IBM Power5 Chip: A Dual-Core Multithreaded Processor*. R. Kalla, B. Sinharoy, J. Tendler – IEEE Micro, vol. 24, no. 2, 2004
 - *Hyper-Threading Technology Architecture and Microarchitecture*. D. Marr, F. Binns, D. Hill, G. Hinton, D. Koufaty, J. Miller, M. Upton – Intel Technology Journal,

vol.3, issue 1, 2002

- *Niagara: A 32-Way Multithreading Sparc Processor.* P. Kongetira, K. Aingaran, K. Olukotun – IEEE Micro, vol. 25, no. 2, 2005
- *CMP Implementation in Systems Based on the Intel Core Duo Processor.* A. Mendelson, J. Mandelblat, S. Gochman, A. Shemer, R. Chabukswar, E. Niemeyer, A. Kumar – Intel Technology Journal, vol. 10, issue 2, 2006