
Disclaimer: "The contents of this document are scribe notes for The University of Texas at Austin EE382V Spring 2007, Computer Architecture: User System Interplay. The notes capture the class discussion and may contain erroneous and unverified information and comments.*

Piranha: A Scalable Architecture Based on Single-Chip Multiprocessing

Lecture #18: Wednesday, 28 March 2007
Lecturer: Dr. Mattan Erez
Scribe: Bob Ascott
Reviewer: Dr. Mattan Erez

The discussion of the Piranha Single Chip Processor was limited to approximately 30 minutes as the class adjourned to attend the lecture from AMD on their new Barcelona Quad-Core microprocessor.

This paper represents a CMP (Chip Multi Processor) approach to parallelism which is very different from papers we have studied previously.

1 What is the Problem Being Solved?

- The design of Piranha was aimed at Commercial server workloads which can be characterized by a large data and instruction footprint which can lead to many memory stalls. Also data dependencies in these applications do not work well in multi-issue designs, and floating point units are not usually required. Finally very little instruction level parallelism is found in these applications.
- High TLP (Thread level parallelism) is the primary focus of the parallelism provided by the Piranha processing system.
- A design goal was established to minimize the cost of the design through the use of ASIC's, Libraries, and Modular-reuse.
- Full backward compatibility was targeted to maximize the number of systems and applications which could benefit from this design.

*Copyright 2007 Bob Ascott and Mattan Erez, all rights reserved. This work may be reproduced and redistributed, in whole or in part, without prior written permission, provided all copies cite the original source of the document including the names of the copyright holders and "The University of Texas at Austin EE382V Spring 2007, Computer Architecture: User System Interplay".

2 Who are the intended users?

- Users of this system are Commercial/Server processing applications. This can include system programmers, application programmers, data center personnel, and ultimately the recipients of data from these commercial applications.
- (The discussion of this topic was truncated due to time constraints.)

3 What is unique about the suggested solution?

- Non-inclusive L2 cache means that data moved to the (large) L1 caches is not replicated in the L2 cache. This results in the L2 cache being treated as a victim cache.
- L1 Tags in L2 assist the coherency logic to determine what data is in the L1 caches when lookups from other processors are performed.
- A sophisticated Cross-Bar switch for on-chip network is employed. This switch has 27 sources and consumers of the data packets.
- No Negative acknowledgements (NAKS) are employed in the cross bar network. Further characteristics of the network include hardware for Priority Service, Hot Potato Routing, and Age based priorities.
- The class engaged in a side bar discussion about Live-Lock and Dead-Lock. Live-Lock occurs when messages are moving through the network but not arriving at their intended destination. Dead-Lock occurs when nothing is moving.
- The memory system used the RAMBUS, which allowed narrower memory bus to handle the many busses off the chip. This was key as a wider bus would have limited the number of memory subsystems that could have been supported.
- The processing subsystem provided very high bandwidth both on the chip and off the chip. It was noted that this bandwidth appears to be significantly higher than the processing rate of the CPU.
- Because ECC is provided on the main storage subsystem, ECC bits in the CACHE were reused for directory data which supported the coherency system.
- The coherence engines were programmable allowing changes in the coherency protocols.

4 How is the idea evaluated? (from scribes's notes - not discussed in class)

- Both Simulation and implementation of design were performed.
- For completeness, the design was projected into more aggressive technology to move the comparison towards the future.
- Comparison were made against an aggressive OOO (out of order) multi-issue processor.
- Benchmarks OLTP - TPC-B and DSS (Decision Support) were used to evaluate the design.

5 Was the evaluation in line with the stated problem or solution?

- Yes - with qualifications
- The evaluation did not compare against the primary competitive vehicle - SMT (Simultaneous Multi-Threading) processors.
- No measurements of the inter-connection system were conducted perhaps due to imbalance in network speed vs processing performance.

6 Questions 6-9

- The questions covering technology, software, other requirements, and context of class were not discussed due to the short class period. This processor represents another attempt to exploit parallelism to achieve performance improvements. The focus on commercial applications with multiple threads defined as part of the application provided the focus for both the problem and solution.