

Parallel Processing-Parallel Memory Approach for Super Fast Design of Future Microprocessor

Yaakob Karomy. Hanna
Al-Mustansiriyah University, College of Engineering
Department of Computer & Software Engineering
Baghdad, Iraq

Abstract

The early design of the microprocessor (μP) used a single ALU with a single unit of memory. The development of the microprocessor design generates a multi-ALUs microprocessor that is a parallel processing with multi-units of memory.

The parallel processing approach will increase the speed of the processing but this speed up is non-linear with increasing the number of processors that are used in the system (μP). However, the efficiency of the parallel processing is non-linear and depends on some factors such as the parallel processing type, the overall design, the programming approach and the applications, yet in general the parallel processing efficiency will decrease by the increase of the number of processors in the system.

The history of the Intel μP 's will be used as an example to trace and analyze the growth of the μP . This tracing will disclose that the main future problem in the μP is the storage not the processing of the data. This problem is generated because the shared memory in the parallel processing will capture the processors in the system. The processor in this parallel processing system is not free to use the memory but it shares a single memory with other processors in this system.

This paper propose a novel approach designing a parallel memory that gives the processors in the parallel processing system a higher freedom to use the memory and eventually increases the efficiency of each processor, that end result will increase the total speed of the parallel processing system because it will become a parallel processing parallel memory (PPPM) system.

This approach will apply to the Intel processor P^{ξ} , which will show that it is able to increase the speed of the P^{ξ} processor for more than four times. These results are used to propose a future design strategy as a first step to implement a super fast processor and then a super fast PC.

The proposed processor is PPPM system with 256 ALUs, this processor is expected to enhance the strategy of the management and control units to become a successful super fast processor with speed up to 7.1 times over the Intel P^{ξ} .

Keyword: Intel μP , parallel processing, parallel memory, cache memory, DMA, memory cycle, P ϵ , prediction, fetching, Net Burst, shared memory, delay, latency.

الخلاصة

التصاميم الاولية للمعالجات الدقيقة استخدمت وحدة واحدة فقط من وحدات الحساب والمنطق (ALU) مع ذاكرة واحدة. التطوير في تصميم المعالج الدقيق خلق معالجات دقيقة بعدة وحدات من الـ (ALU) تعمل بصورة متوازية مع بعضها البعض مع استخدام عدة وحدات من الذاكرة (multi-units of memory).

طريقة المعالجة المتوازية تزيد من سرعة المعالجة ولكن هذه الزيادة لا تكون زيادة خطية مع الزيادة في عدد المعالجات المستخدمة في النظام. ان كفاءة نظام المعالجة المتوازية ليست خطية بسبب اعتمادها على عوامل اخرى مثل، نوع المعالجة المتوازية، التصميم النهائي للنظام، نوع البرمجيات مع التطبيقات، لكن بصورة عامة ان كفاءة نظام المعالجة المتوازية تقل كلما ازداد عدد المعالجات في النظام.

ان تاريخ تطور المعالجات من نوع (Intel) ستستخدم كمثال لبحث وتحليل نمو وتطور المعالج الدقيق. هذا البحث سيكشف ان المشكلة المستقبلية الرئيسية في المعالج الدقيق ستكون في خزن المعلومات وليس في معالجتها. هذه المشكلة تظهر بسبب ان الذاكرة المشتركة في نظام المعالجة المتوازية ستحدد عمل المعالجات في النظام حيث ان المعالج في هذا النظام سيكون ليس حرا في استخدام الذاكرة ولكنه سيتقاسم تلك الذاكرة مع المعالجات الاخرى في النظام. البحث يقترح نهجا جديدا لتصميم الذاكرة المتوازية التي تعطي المعالجات في نظام المعالجة المتوازية حرية اكبر في استخدام الذاكرة مما يزيد من كفاءة كل معالج، ونتيجة لذلك ستزداد السرعة الكلية للنظام بسبب نظام المعالجة المتوازية والذاكرة المتوازية (PPPM).

هذا النهج سيطبق على المعالج الدقيق نوع (Intel P ϵ)، والذي سوف يزيد سرعة المعالج P ϵ لأكثر من أربع مرات. هذه النتائج تستخدم لاقتراح استراتيجية تصميم مستقبلية كخطوة أولى لتنفيذ معالج بسرعة فائقة ومن ثم حاسوب شخصي (PC) يمثل هذه السرعة.

المعالج المقترح يستخدم نظام معالجة متوازية مع ذاكرة متوازية (PPPM) بـ ٢٥٦ ALUs وحدة حساب ومنطق، ومن المتوقع ان هذا المعالج سيعزز استراتيجية وحدات الادارة والسيطرة لتصبح بنجاح معالج فائق السرعة تفوق سرعة معالج (Intel P ϵ) بـ (٦٠) مرة.

١. Introduction

What are the requirements to implement a fast PC? The basic answer is the implementation of a fast PC requires a fast μP ; but this answer is not complete because the fast μP requires to a suitable components to help it. The fast PC required at least to:

١. Fast μP .
٢. Large size of memory with suitable heretically implementation.
٣. Flexible I/O controller.

This answer arise a new question, what is the fast μP ? It's a μP that has at least:

١. Fast clock.
٢. Wide data bus (ALU).
٣. Wide address bus.

- ξ. Suitable control signals.
- ο. Large set of instructions.
- ϒ. Good control circuits.

The application of all the above points in the design of the μP is not enough to implement a fast μP ; the key of the next step to increase the μP speed is the parallel processing with extra help units such as fetching unit, MMX, SSE, Net Burst,..... [1]

This paper will trace the steps of the processor growing from the history of the Intel processors Section ϒ.

Section ϓ shows the types of the parallel processing while the efficiency measurement of the shared parallel processing system will be summarized in section ξ.

Section ο shows the design of the P^ξ processor as suitable example for Intel μP that uses the parallel processing approach.

Section ϒ will give the key of the novel method of the parallel memory, while the effect of the PPPM approach will present in examples in section ϗ that will estimate the time delays using a simple basic equations then simulate the simple parallel processing system for much better results.

Section ^ gives the main lines to design a super processor in the future, the conclusions and a simple future work of this paper are in sections ϑ and ϑ.

2. Growing History of Intel microprocessors

The history of the Intel microprocessors will be summarized in tables (1, 2 and 3); these tables show how the growing of the μP is joined with the other factors in the PC. The light points for this history are noted in bold lines in the tables like the use of DMA in 8008 μP .

The first row in table (1) shows the ξοοξ μP in 1971 which is the first general IC (the early form of μP) used to implement micro-controllers such as the simple calculator. this IC manage 1ξο Byte of memory and uses the ξοοϑ IC as I/O manager. The expansion of the memory size to 1ϒ KB in 8008 will require a special memory management and I/O terminals. This unit is called Direct Memory Access (DMA). The 8008 IC has the same architecture of the ξοοξ IC with twice number of bytes for data bus and address bus. The 1ϒ bit address bus was presented in 8008 that expands the memory to 1ξ KB. The 8008 in 1971 is a glaring point in the Intel processor history; this IC is a real μP that has a complex set of instructions 1ϒ bit address bus with new input (maskable interrupt) controller, so it has three internally data-buses with all necessary registers set. The weakness in this μP is the limited memory (1ξ KB) that is not suitable to save the real programs and the data for the real problems.

Expansion of the address bus to ϓο bit expands the memory size to 1 MB. The new size for memory requires a large size memory word and registers to capture the addresses of the instructions and the data in the memory, this problem is solved by using the *segmentation* approach. The electronic growth helps to implement fast and large ICs, the designer use this

advantage to increase the data bus to 16 bit, and hence to implement the first element for the successful family of processor that called 8086 in 1978. The 8086 μ P has speed up to 8 MHz, maskable interrupt and management 1 MB memory will be used to realize the dream to implement the early PC that called IBM PS/2.

Table (1). The early generations Intel μ Ps.

Year	μ P	Data bus	Address bus	Clock (Hz)	MIPS	Memory (Byte)	I/O	Used
1971	8008	8	8	0.06 MHz	0.06	640	8009	Calculator
1971	8080	8	8	0.06 MHz	0.06	1024	8207	Calculator
1972	8085	8	8	0.06 MHz	0.06	16K	DMA	-----
1974	8088	8	16	2M	0.64	64K	DMA	Traffic light
1976	8088	8	16	2M	0.37	64K	maskable interrupt	
1978	8086	16	20	8M	0.66	1M segmentation	maskable interrupt	IBM PS/2

The implementation of μ P which is suitable to implement the practical PCs starts with the 8086 μ P in 1978 as in table (1). The real challenge in that time was the software manager not the hardware, this problem remained about 3 years then the first real software manager for the PC that called Microsoft Software Disk Operating System (MSDOS) has been born, this software with the DMA and maskable interrupt input built a successful strategy for PC, that will be growing with the μ P later.

The 80286 μ P that has 24 bit in address bus and 16 MB memory with MSDOS is able to implement a useful PC that can be used for long time. The new problem in that time was the limited memory size in comparison with the requirements for solved problems in PC. The solution comes with 386DX in 1985 that had 32 bit for data bus and 32 bits for address bus to expand the memory to 4 GB with additional feature called *virtual memory* that expanded the external memory of the PC (hard disk) to 16 TB.

The large memory size in 386DX required a long read and write times. this problem was solved by adding a small (8 KB) intermediate memory called *cache* memory in 486DX.

The Pentium technology opens a new gate to implement a very large and very high speed ICs (3,1 million transistors with 323 pin), the data bus became 32 bit with 16 KB cache in the first generation of these processors which called P5 or P6. This processor with all the next generations uses two new technologies, the first is the *prediction* unit which represents the early class of the parallel processing, while the second is the imbedded processor to increase the speed independent of the memory and I/O terminals manager. The new technology (Pentium) with the

two new features in the processor gives the PC a high speed that makes the PCs able to run the advances operating systems (AOS) programs such as windows.

The 2 ALUs with two levels cache memory in Pentium Processor represent the early attempt to use the *parallel* processing approach that increases the power of processor. The additional important unit emerged in 1997 with PII used MMX internally unit that made the processor able to run a large set of instructions that called MMX set instruction, therefore sometimes this processor is called Pentium MMX. The PIII expands to run a new set of instructions called SSE also it used 4 ALUs.

Table (2). The practical generations Intel μPs for PC.

Year	μP	Data bus	Address bus	Clock (Hz)	Cache (Byte)	Memory (Byte)	I/O	Other factors
1982	80186	16	20	5M	no	1 M seg.	DMA Interrupt	MSDOS
1982	80286	16	24	10M	no	16 M seg.	DMA Interrupt	MSDOS
1980	80386DX	32	32	16M	no	4 G virtual 64 T paged	DMA Interrupt	MSDOS
1989	80486DX	32	32	20M	L1 8K	4 G 64 T	DMA Interrupt	MSDOS
1993	P0 (PI)	64	32	60M	L1 16K	4 G 64 T	Imbedded processor	Predict +AOS
1995	Pentium Processor	64	32	100M	L1 16K L2 256K	4 G 64 T	Imbedded processor	2 ALUs
1997	PII	64	32	300M	L1 32K L2 512K	4 G 64 T	Imbedded processor	MMX
1999	PIII	64	32	500M	L1 256K L2 2M	4 G 64 T	Imbedded processor	4 ALUs SSE

The ALUs in Pentium Processor., PII and PIII are half parallel processing because they are connected as a parallel input serial output, the real parallel processing starts with *Net Burst* in P4 as in table (3) it has 4 ALUs work as parallel input and parallel output units.

Table (3) shows how the micro-architecture of the classic μP reaches saturation and the improve start to use the other factors as in the last column in table (3). P4 satisfies the parallel processing using *Net Burst* and enhances the set of instructions using *SSE/SSE2* sets.

Itanium 2 (Dual core) grouped two P4 μPs internally and expands its instruction set to new set called SSE2. The growing generates a new technology called Intel 64 that is used to implement the Intel Core 2 type Kentfield which has two dual cores internally and expands it to the Intel Core 2 type Yorkfield that has 2 cores internally with additional instructions set called SSSE2.

Core i3 technology (Nehalem type Lynnfield as example) used 3- physical cores with 3-channel DDR3 memory also its designer separates the cache memory L3 to two parts one for data

and the other for instructions, expand L₂ to 4 sections each one is 256 KB adding the third common cache level L₃ that has 8 MB capacity.

The Dual core, Two Dual core, Core 2 Quad, 4-physical cores with 2-channel DDR2 and in the last the 6-physical cores with 4-channel DDR2 and the *Turbo Boost Technology* in Pentium Core™ i7 (Westmere type Gulftown) μP strategies in March 2010 enhanced the parallel processing but all the new technologies will not add high features for the total power of the μP because the limitations of the parallel processing that will be discussed later. For more detail about the history of microprocessor see Ref [1-5].

Table (2). The modern generations of Intel μPs.

Year	μP	Data bus	Address bus	Clock (Hz)	Cache (Byte)	Memory (Byte)	Other factors
2002	P4	64	32	2, 4G	L1 16 K L2 256 K	4 G 64 T	Net Burst SSE/SSE2
2002	Itanium 2	64	32	2, 4G	L1 16 K L2 4 M	4 G 64 T	Dual core SSE2
2006	Intel 74, Intel Core 2, Kentfield	64	32	3G	L1 16 K L2 4x4 M	4 G 64 T	Two Dual Core SSSE2
2008	Intel 74, Intel Core 2, Yorkfield	64	32	3G	L1 32K L2 4x4M	4 G 64 T	Core 2 Quad SSSE2
2009	Nehalem, Core i3, Lynnfield	64	32	3, 4G	L1 4x32+32 K L2 4x256 K L3 8 M	4 G 64 T	4- physical cores 2-channel DDR2
2010	Westmere, Core™ i7, Gulftown	64	32	3, 4G	L1 4x32+32 K L2 4x256 K L3 8 M	4 G 64 T	6- physical cores 4 × DDR2 Turbo Boost Technology

3. Parallel Processing

There are two main types of the parallel processing, each one has advantages and disadvantages over the other type. These two types are the Distributed Memory and the Shared Memory. Also these types are improved to become hybrid to generate a large amount of parallel processing systems. For more detail about the parallel processing see Ref [4, 6].

3.1. Distributed Memory

Distributed memory systems vary widely but share a common characteristic. It requires a communication network to connect inter-processor memory as in Fig (1).

The processors have their own local memory. Memory addresses in one processor do not map to another processor, so there is no concept of global address space across all processors. Because each processor has its own local memory, it operates independently. Changes that made to any local memory have no effect on the memory of other processors. Hence, the concept of cache coherency does not apply. [4]

When a processor needs access to a data in another processor, it is usually the task of the programmer to explicitly define how and when data is communicated. Synchronization between tasks is likewise the programmer's responsibility. The network "fabric" used for data transfer varies widely, though it can be as simple as Ethernet.

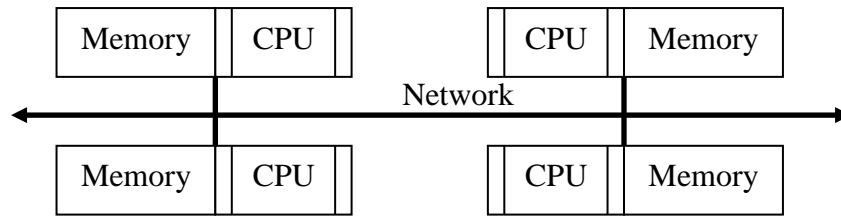


Fig (1): Distributed memory systems parallel processing.

The advantage of this method is each processor can rapidly access its own memory without interference and without the overhead incurred with trying to maintain cache coherency. While the disadvantage of this method, it is used for networks and for independent processors. [4]

3.2. Shared Memory

The shared memory parallel computers Fig (2) vary widely, but generally have in common the ability for all processors to access all memory as global address space. Multiple processors can operate independently but share the same memory resources. Changes in a memory location effected by one processor are visible to all other processors.

The advantages of this method are Global address space provides a user-friendly programming perspective to memory and the data Sharing between tasks is both fast and uniform due to the proximity of memory to CPUs.

The disadvantages of this method are the lack of scalability between memory and CPUs. Adding more CPUs can geometrically increase traffic on the shared memory-CPU path and for cache coherent systems and geometrically increase traffic associated with cache/memory management, programmer responsibility for synchronization constructs that insure "correct" access of global memory. [4]

This paper will use the shared memory technique for all the next sections and call for simplicity "parallel processing".

3.3. Multi Levels Parallel Processing

The shared memory method is active for a small number of processors (max. 8 or 16 processors). The system with a large number of processors will use a multi levels parallel processing as shown in Fig (3). The total result is like a parallel inside parallel processing. [5]

The number of processors in each group is m ; this number in default is equal for all groups because that gives optimal results and simple control. These (n) groups will be connected to a common memory in form of shared memory, this system is indicated to as m/n . Theoretically there is no condition that limits the number of levels or the number of groups in each level or the number of processors in the higher level, but practically the symmetry saves the best case.

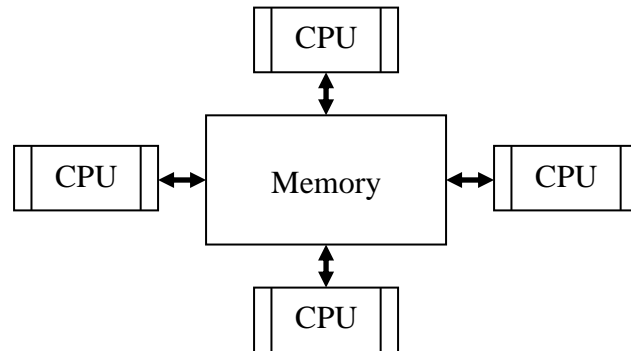


Fig (2): Shared memory systems parallel processing.

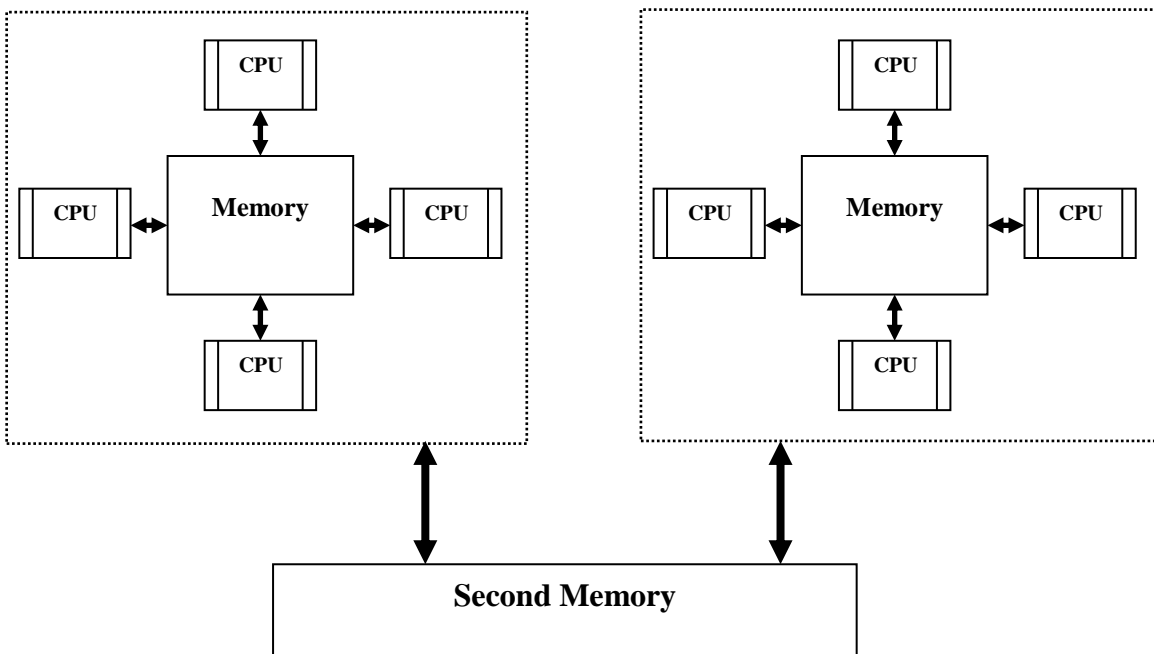


Fig (3): Two levels (ξ/η) shared memory parallel processing system.

3, ξ . Virtual Processor

Some types of processors such as P^ξ use multiple levels of memories (L^1 and L^η) to interchange the data. Hence, the relation between them is not memory to memory but it is like the

relations between processor and the memory, therefore L^y will regard the L^x as a processor; hence, the memory unit L^x in this case will be called a *virtual processor*.

4. Efficiency Measurement of the Parallel Processing System

There are many factors that affect the efficiency of the parallel processing system like the number of processors in the system, the distribution form of these processors, the design strategy, type of connection, the processor characteristics, the ratio of used memory, memory speed, the size of the memory and other different factors that depend on the system and problems. [7, 8]

However, for simplicity let:

1. All the processors are of the same type.
2. All processors and memory have same speed (equal clocks time).
3. All processors have same memory used ratio (P_{mi}).
4. The value of P_{mi} is less than $1/n$ where n is the number of processors in the system.
5. Neglect all other factors.

Under these conditions the total speed of the system can be written as in the following equation (1). [7, 8]

$$S_s \leq S_p \times \frac{n+1}{2} \quad \text{if } \sum P_{mi} \leq 1 \quad (1)$$

Where S_s is the speed of the system, S_p is the speed of a single processor.

Define a new factor R_s as the maximum rational system speed that is given in equation (1) as:

$$R_s = \frac{n+1}{2} \quad \text{if } \sum P_{mi} \leq 1 \quad (2)$$

The solutions for equation (1) show that the use of two parallel processors will give a system of a single processor with virtual speed less than 1.5 times of the original speed of a single processor (S_p). Table (4) shows some results of a simple simulation (enhanced with some practical results) for low rate data systems ($P_{mi} < 0.5$).

From table (4) we can note that, the efficient case is the system of two parallel processors with low traffic rate in memory. Another important note is the increase of the number of parallel processor will decrease the efficiency of the system for all cases and forms. The maximum practical number of processor is 16 processor and maximum acceptable levels is 3 levels. So any specifications over this range will reduce the efficiency of the processor to not useful range.

The practical results are lower than these values except for some special problems that accept the parallel computation. For example the efficiency of a 3 parallel processors full in rang from 40% to 80% and the efficiency of a 4 parallel processors full in rang from 45% to 60%. [9]

5. P4 Intel Processor

The designs of Intel processors are very different as we noted in section 4, which means that there is no common design for these processors. Hence, we will select the design of the P4 processor as a good example to show the simple design of a modern parallel processor.

The architecture of P4 and its working strategy is very complex but we will simplify the main important points of the architecture and focus on the parallel processing features. Then we discuss its main algorithm of the memory cycle that will help to estimate the speed and efficiency of the processor.

Table (4). Rational system speeds for parallel processing (shared memory).

Number of processor (shared memory)	Form	Rational speed	Efficiency Rs/n*100
2	Single level	1,80	93%
4	Single level	2,70	79%
6	Single level	3,0	68%
8	Single level	4	60%
8	4/2 or 2/4	4,6	68%
16	Single level	4,9	31%
16	8/2	6,3	40%
16	2/8	6,3	40%
16	4/2/2	7,6	48%
16	4/4	7,1	44%
64	8/8	14,4	23%
64	4/4/4	16,9	26%

5.1. The Architecture of P4 Processor

The architecture of the Intel P4 processor [4] is very complex as in Fig (4) but the main important advantages are:

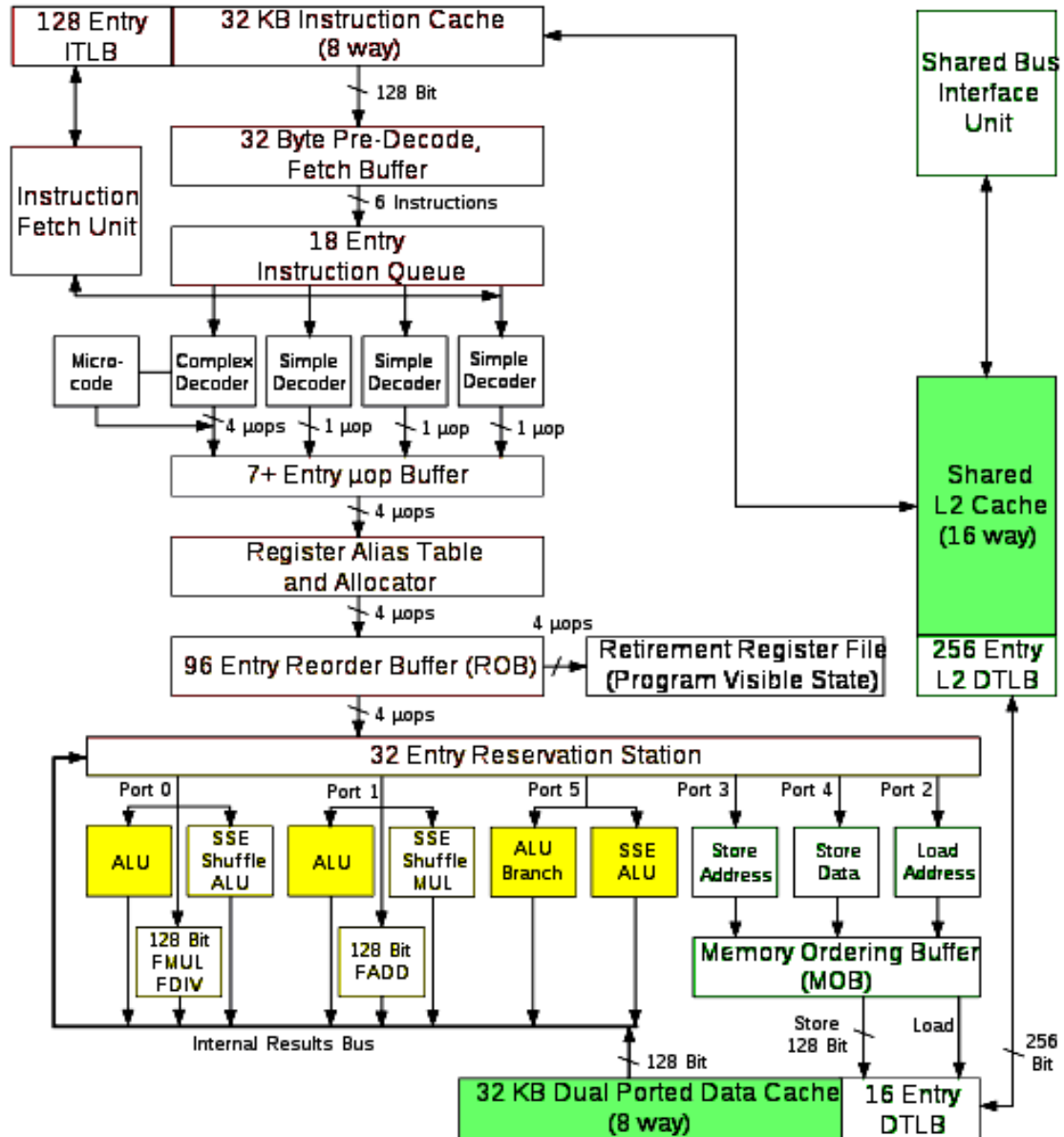
1. It has 4 ALU's, 4 ways allocation registers and Net Burst unit to work as 4 parallel processors.
2. The four ALU's are not symmetric but each one has special properties as indicated in Fig (4).
3. The allocation and the 22 entry station unit garnet the parallel in the input side of the 4 ALU's.
4. The Net Burst (impeded circuit in the control unit) garnet the parallel in the output side of the 4 ALU's.
5. There are three auxiliary ALU's used for auxiliary missions such as memory load and store, these ALU's do not appear in the calculation of the efficiency.
6. It has two levels of internally caches, the first L1 has 8-ways and the second L2 has 16-ways. The size of these caches is different from type to another, so the sizes of their blocks are different in these types. But the default values for these factors are that used in table (5) and in the example later.
7. These four processors are connected as a shared memory with L1. The L1 memory connected with L2 (as virtual processor) the final form is like *two levels (4/1) shared memory systems parallel processing*.

^ The processor will be connected with a single main memory in PC. The result will become like three levels (ε/1/1) shared memory systems parallel processing.

◦, 2. The Memory Cycle of Pε Processor

The traditional memory cycle [9] is as in Fig (◦) summarized in four main steps:

- 1. Start with the data request from instructions fetch unit to L1 the hit in this point will return the successful order to read/write order in step 2, while the miss will freeze the processing operations and go to step 3.



Intel Core 2 Architecture

Fig (ε): micro architecture of Intel Pε processor.

- ϕ. Check the write instruction algorithm to move the new data in L^1 to L^2 and later to main memory in the ideal memory time.
- ϕ. Start the search in L^2 cache memory, if L^2 hit move the wanted data block from L^2 to L^1 and go to step ϑ, while the miss sends processor to step ζ.
- ζ. Start the search in the main memory, if the main hit moves the data from the main to L^2 to L^1 and go to step ϑ, while the miss sends processor to step ϖ.
- ϖ. Defreeze the processor to start processing with the next instruction and send order to DMA unit to start the search in the hard disk of PC.
- ϗ. When DMA find the wanted data start to move the data from the hard disk to the main memory and indicate the processor.
- Ϙ. Refreeze the processor and move the data from main to L^2 to L^1 and go to step ϑ.

Note that the write instruction algorithm to move the new data from L^1 to L^2 and to main memory will add a latency time because it works in the ideal memory time.

The miss in L^1 and L^2 causes a delay in the processor, because of the freeze in processing in this time. The freeze order is compulsory because the two caches are busy to move the wanted data to L^1 .

The miss in the main memory is latency because the defreeze order, but this time will become delay if the next instruction depends on the present instruction or if it requires data at the same block as we will discuss in section ϙ.

ϙ. Proposed Parallel Memory Technique

The deep analysis for the P^z processor shows that, the weakness point in this system is the *crowding* of the virtual processors (ALU's and L^2) on the L^1 cache. This problem is solved partially in P^z by using \wedge way cache L^1 .

The \wedge way cache L^1 [\wedge] means the L^1 is divided into \wedge independent parts but this solution will reduce the problem by factor less than 90% because the stream of instructions is stored in the same block so their data are also stored in the same block generally.

At this time we need to *crumble* the single block to independent bytes form, this solution (no blocking) is very complex and impractical for large memories such as in P^z .

This paper proposes a simple solution with a small additional cost and a small additional complexity for the control circuit in the μP but eliminates the problem of crowding in L^1 in about 90%.

The basic idea depends on the use of *multi independent internally channels*, each channel has *data bus and address bus* inside the memory as in Fig (ϙ a) [ϙ ⋅] that shows the architecture for the control of the single cell in the present memory while Fig (ϙ b) shows the architecture for the control of the single cell in the proposed memory with two channels. The channel that uses

the memory doesn't care to the other channels. Hence, each channel is free to use the memory, so that each processor will become *free* to use the memory.

The collision is found when two channels need the same byte, this problem has very low probability so it will be solved using simple controller inside the memory giving a priority for channels. The controller will delay the lower priority in one cycle (very short time). The low probability and very short time delay means the effect of this case can be neglected in the calculations of the time processing.

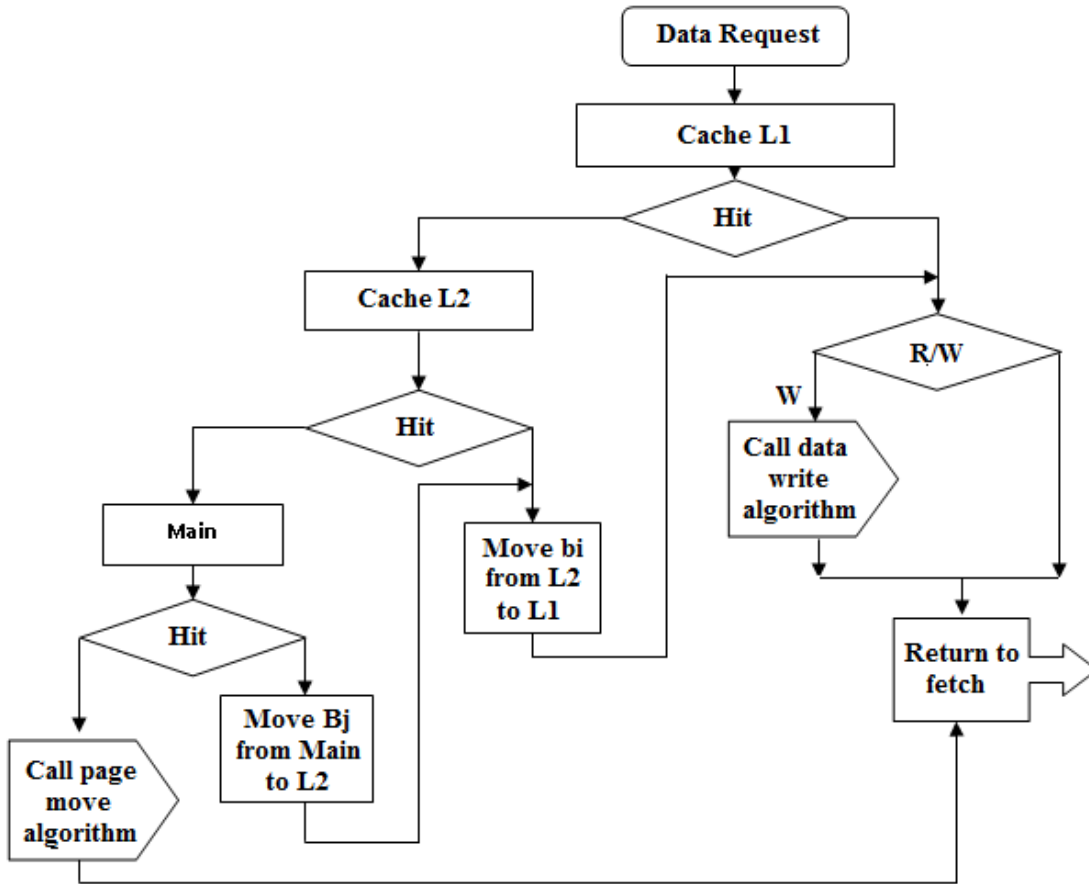


Fig (°). The memory cycle in two levels caches Intel processor.

Hence, this approach will be used in all the memories of the μ P and PC. The effect of this idea will show in the processor example in the next section.

4. Time Delay Calculations and Simulation

4.1. Delay and Latency

If we have a queue of instructions under execution in a μ P, so any frozen for an instruction in this queue because, generally; it requires data as the miss in L₁ case in memory

cycle. This freeze will be solved by stopping all the next instructions. This stopping time is called *Delay* and all the next instructions results will be delayed. Sometime this delay is large as in case of miss in the main memory, the solution to illuminate this delay is by putting aside this instruction and executing the next instructions until the required data is prepared which is necessary to run this instruction, this time delay is called *Latency*. The latency delays the result of

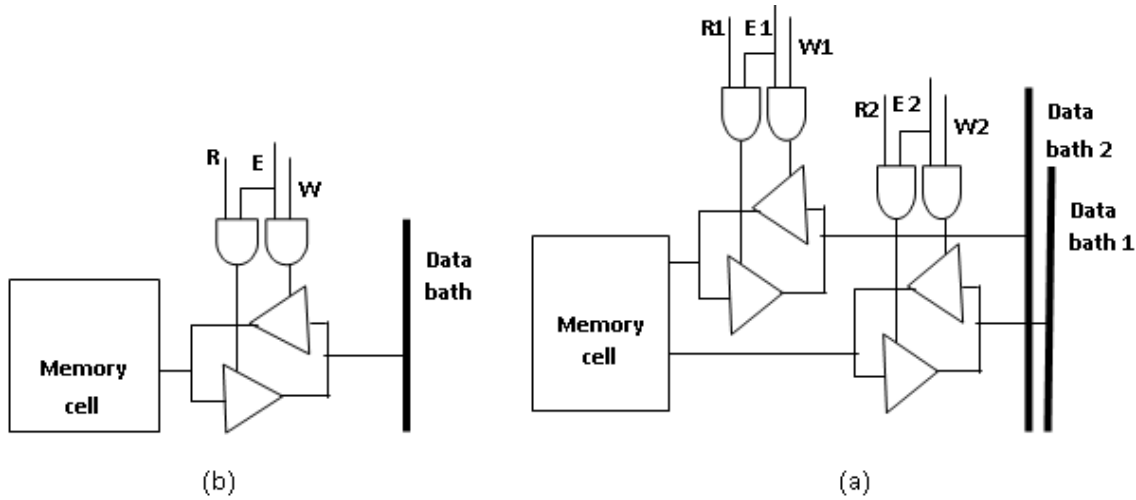


Fig (1). Memory cell architecture, (a) single channel (b) proposed two channels.

the present instruction while it passes the results of the next instructions. Some times one (or more) instruction from the next instructions queue is depends on the aside instruction or sometimes it share the same block with it. In theses two cases all the latency of these instruction will become a delay.

4.2- Processor Example

The practical speed calculations of the processor with and without the parallel memory require practical values. These proposed values are the default values of P^z processor [1, 11] in documents and statistics. They will be classified for simplicity into three groups as follows:

Group A: memory specifications

1. The block size in L¹ (bi) is 16 B.
2. The block size in L² (Bj) is 1 kB.
3. The page size in main memory (Pk) is 16 kB.

Group B: processor times

1. t is the clock time = the time of executing one instruction in fast ALU in the processor and let t=1.
2. The time to read or write any data from L¹ is 1 (1t).
3. The time to read or write any data from L² is 10 (10t).
4. The time to read or write any data from main memory is 160 (160t).
5. The time to transfer one block from (to) L¹ to (from) L² is Tb=10 (10t).

- ٦. The time to transfer one block from (to) L^x to (from) main memory is $TB=1 \dots (1 \dots t)$.
- ٧. The time to transfer one page from (to) hard disk to (from) main memory is $Tp=1 \dots (1 \dots t)$.

Group C: probabilities of the work

- ١. The number of run instruction in the CPU is $1 \dots$.
- ٢. See table (٥) for the used probabilities. These values are the most used in P^z processor.

Table (٥). Definition of the probabilities set and their default values in the example.

Symbo l	Value	No. of instructions	Definition	Note
Pm	٠,٢	٢٠٠٠٠٠	Processor need to external data	Pm + Pp = ١
Pp	٠,٨	٨٠٠٠٠٠	Processor not need to data	
Ph	٠,٩	١٨٠,٠٠٠	Hit in L^1	Ph + Pmh + Pmmh + Pmmm = ١
Pmh	٠,٠٩	١٨٠,٠٠٠	Miss in L^1 Hit in L^x	
Pmmh	٠,٠٠٩	١٠٨,٠٠٠	Miss in L^1 Miss in L^x Hit in Main	
Pmmm	٠,٠٠١	٢٠٠	Miss in L^1 Miss in L^x Miss in Main	
Pr	٠,٥	١٠٠٠٠٠	Read data	Pr + Pw = ١
Pw	٠,٥	١٠٠٠٠٠	Write data	
Pd	٠,٠٠١	٢٠٠	Dependency*	
Pwo ^١	٠,٠٢	--	Over flow in L^1	
Pwo ^x	٠,٠٠١	--	Over flow in L^x	
Prw	[٠,٠١]		Run in the waiting time	Practical value

* Dependency probability means the run of the next instruction depends on the results of the present instruction or the data of the two instructions are inserted in the same block (bi).

The application of the values in table (٥) over the memory cycle in Fig (٥) will result to calculate the average required time to run a single instruction in P^z processor Tav . These calculations of the average time bellow will shows that the approximated time without parallel-memory approach required is about ٢.٨٧ cycles while the required time using parallel-memory approach is about ٠.٦٣ cycles. In other words the parallel-memory approach will increase the speed of the processor in more than ٤ times.

a. The calculations of the average required time in the example without parallel-memory approach.

$$Tav = Pp \times 0.25 \times (1 - Prw) + Pm \times \{ [Ph \times 2 + Pmh \times (2 + Tb) + Pmmh \times (2 + Tb + TB) + Pmmm \times (2 + Pmmm \times Tp)] + Pw \times [Pwo1 \times Tb + Pwo2 \times TB] \}$$

$$T_{av} = 0.8 \times 0.25 \times (1 - Pr_w) + 0.2 \times \{ [0.9 \times 2 + 0.09 \times 22 + 0.009 \times 1022 + 0.001 \times (2 + 0.001 \times 150,000)] + 0.5 \times [0.02 \times 20 + 0.001 \times 1000] \}$$

$$= 0.2 \times (1 - Pr_w) + 2.626 + 0.14$$

$$T_{av} = 2.866t$$

$$= 0.2 * (1 - 0.5) + 2.626 + 0.14$$

b. The calculations of the average required time in the example using parallel-memory approach.

$$T_{av} = P_p \times 0.25 \times (1 - Pr_w) + P_m \times \{ P_h \times 2 + P_{mh} \times (2 + P_{mh} \times T_b) + P_{mmh} \times (2 + P_{mh} \times T_b + P_{mmh} \times T_B) + P_{mmm} \times (2 + P_{mmm} \times T_p) \}$$

$$T_{av} = 0.8 \times 0.25 \times (1 - Pr_w) + 0.2 \times \{ 0.9 \times 2 + 0.09 \times (2 + 0.09 \times 20) + 0.009 \times (2 + 0.09 \times 20 + 0.009 \times 1000) + 0.001 \times (2 + 0.001 \times 150,000) \}$$

$$= 0.2 \times (1 - Pr_w) + 0.12184 = 0.2 \times (1 - 0.1) + 0.449764$$

$$T_{av} = 0.629764t$$

4.3. The Results of the Simulation

This sub-section will summarize the results of a simple simulation for the memory cycle of the Intel processors. This simulation will use the default values in the previous sub-section and the information of processors in table (3). The simulation will run over a proposal case (multi-channels parallel memory) to show the effects of the parallel-memory approach. Hence, the number of channels will be indicated in two arcs {} in table (4).

The times in table (4) measured for a universal clock for all processors that write in CPO (clock per operation) while the inverse term OPC (operation per clock) will be use to measure the speeds, the results in the last column is the percentage ratio for the ideal case.

The first row in table (4) shows the properties of the present P^ε processor that describe as ε{ }/^{ }/^{ } that's mean it has ε ALUs connected to L¹ which is connected to L² and with L² connected to main memory. Hence, all these memories (L¹, L² and main) have a single internally channel.

The row 1 (row ε) satisfies the ideal case for the processor that has an ideal time 0,20 (0,0,ε2) for the P^ε (Gulftown) processor but the ALUs are not symmetrical and not all of them are general but each one is specialized for a type of instructions as in Fig (ε) and these instructions cannot be balance for all applications, therefore some ALUs in processor will be idle for some time. The idle time in the case of no miss means this case is not ideal case.

4. Design of Super Future Processor

The design of a super future processor requires a deep study of the results in table (4), and then tries some designs under simulation to select the most practical one.

row 8 in the table shows the form of the imminent future processor that has 16 ALUs and speed 16 time over P ξ (for the same clock) ideally but the simulation shows its practical speed is less than 7 times over the practical speed of P ξ , with very low efficiency which is about 7%. The enhancement of the imminent future processor using parallel-memory approach (row 9) will jump up to 24% efficiency and the speed to 10 OPC.

The try and error for simulation shows the maximum acceptable number for ALUs in processor is 206. The row 10 shows a processor 100 times faster than P ξ while the row 11 shows a processor 62 times faster than P ξ . The two proposed super processors have 206 ALUs but the second design has lesser memory layers with more complexity control circuits.

Table (1): The results of a simple simulation for the memory cycle of Intel processors.

No.	Type of μP	Parallel form	Time range CPO	Tav CPO	Speed OPC	Efficiency
1	P ξ without miss.	$\xi\{1\}/\{1\}/\{1\}$	0,32*	0,32	3,120	79%
2	P ξ with single channel	$\xi\{1\}/\{1\}/\{1\}$	0,76 - 1,8	1,11	0,9	23%
3	P ξ with two channels	$\xi\{2\}/\{2\}/\{2\}$	0,34 - 0,79	0,39	2,6	60%
4	Gulftown without miss.	$\xi\{1\}/\{1\}/\{1\}/\{1\}$	0,091*	0,091	11	46%
5	Gulftown with single channel	$\xi\{1\}/\{1\}/\{1\}/\{1\}$	0,19 - 0,71	0,37	2,7	11%
6	Gulftown with two channels	$\xi\{2\}/\{2\}/\{2\}/\{2\}$	0,12 - 0,27	0,19	0,3	22%
7	Gulftown with two channels and ξ channels in L γ	$\xi\{2\}/\{2\}/\{2\}/\{2\}$	0,096 - 0,23	0,160	6	20%
8	Future processor	$\lambda\{1\}/\{1\}/\lambda\{1\}/\{1\}$	0,10 - 0,37	0,21	4,8	7%
9	Proposed processor 1	$\lambda\{\xi\}/\{2\}/\lambda\{6\}/\{2\}$	0,32 - 0,12	0,66	10	24%
10	Proposed processor 2	$16\{16\}/\xi\{8\}/\xi\{8\}/\{2\}$	0,14 - 0,05	0,2	00	19%
11	Proposed processor 3	$32\{64\}/\lambda\{8\}/\{2\}$	0,11 - 0,06	0,18	06	22%

The large number of processor (206 ALUs) in PPPM approach cannot be successful without the good additional units. This process requires at least to:

1. An enhanced prediction unit that predicates more than one instruction as example set prediction unit.
2. Extra flexible allocate unit.
3. Enhanced Net Burst circuit.
4. Suitable classification and selection for ALUs with suitable distribution.

9. Conclusions

Having done this work we can conclude some points that help the designer to indicate the optimal way to improve the microprocessor design in the future.

1. Increase the number of parallel processing will be useful for limited number of processors; maximum number is N processors for single level shared memory systems parallel processing.
2. Increase the number of levels also be useful for a limited number of levels; maximum number is K .
- 3- The crowded ALUs over L in the μP is the main problem in the design of present processors.
- 4- The PPPM approach is a very useful technique that will increase the speed of processor to about K times.
- 5- The design of the super processors in the future will not be satisfied by using more ALUs and more cores but it requires radical new ideas.
- 6- The success of the PPPM approach required to enhance and expand the other units in the μP such as prediction unit, allocate unit and Net Burst circuit.

10. Future Work

The key implementation of the super fast processor that covers the future requirements and applications can be summarized in three main layers:

1. Use the PPPM proposed approach in this paper with enhancement and expand the other units in the processors such as the prediction unit, the allocation unit and the Net Burst unit.
2. Use the assistant special processors; this idea is like the use of the math co. processor with the μP in some PCs. The proposed assistant processors are coding/decoding processor, 3D image processor and complex function processor.
3. Discuss and deduce new ideas for the data processing, which are suitable with burst data processing. These ideas must be left the classic idea of the ALU that depends on the using a single point (ALU) to execute a complete single instruction and use algorithms able to execute a large number of instructions partially at the same time. Hint, the burst data processing approach may burn a complex and slow processors but it is able to grow faster than the classical ideas.

References

- [1] List of Intel microprocessors, Wikipedia, the free encyclopedia, "http://en.wikipedia.org/wiki/List_of_Intel_microprocessors".
- [2] Intel Consumer Desktop Pc Microprocessor History Timeline, <http://www.intel.com/intel/intelis/museum/exhibit/>.
- [3] Intel® μP and IA-32 Architectures Optimization Reference Manual, Order Number: 248966-020, November 2009.
- [4] Blaise Barney, "Introduction to Parallel Computing", Lawrence Livermore National Laboratory, "<http://www /Introduction to Parallel Computing.htm>".
- [5] Trupti Patil, "Evaluation Of Multi-Core Architectures For Image Processing Algorithms", Thesis Presented to the Graduate School of Clemson University August 2009.

- [6] B. L. Buzbee, "The Efficiency of Parallel Processing", Frontiers of Supercomputing, Los Alamos Science Fall 1983, pp 51.
- [7] Quentin Smith, "Efficiency of Parallel Processing with JCilk", under the direction of Professor Charles E. Leiserson Massachusetts Institute of Technology Research Science Institute, August 2, 2000.
- [8] Glenn Hinton, Dave Sager, Mike Upton, Darrell Boggs, Doug Carmean, Alan Kyker, Patrice Roussel, "The Microarchitecture of the Pentium 4 Processor", Intel Technology Journal Q1, 2001, pp 1.
- [9] Mostafa Abd-El-Barr, Hesham El-Rewini, "Fundamentals Of Computer organization And Architecture", Copyright # 2000 by John Wiley & Sons, Inc. All rights reserved.
- [10] Samuel Naffziger, Blaine Stackhouse, Tom Grutkowski, Doug Josephson, Jayen Desai, Elad Alon, and Mark Horowitz, "The Implementation of a 2-Core, Multi-Threaded Itanium Family Processor" IEEE journal of solid-state circuits, vol. 41, no. 1, January 2006, pp 197.
- [11] Trent Rolf, "Cache Organization and Memory Management of the Intel Nehalem Computer Architecture", University of Utah Computer Engineering, CS 6810 Final Project, December 2009.