

다중처리 마이크로프로세서를 이용한 객체 인식

정용화¹, 박경, 한우중
하드웨어구조연구팀, 한국전자통신연구원

Object Recognition using an On-Chip Multiprocessing Microprocessor

Yongwha Chung, Kyoung Park, and Woo-Jong Hahn
Hardware Architecture Team, ETRI

요 약

객체 인식은 고성능 컴퓨팅을 필요로 하는 흥미있는 응용 분야이다. 현재 대부분의 고성능 컴퓨터는 슈퍼스칼라 구조의 범용 마이크로프로세서를 채택하고 있으나, 반도체 집적도가 증가함에 따라 슈퍼스칼라 구조를 대신할 새로운 마이크로프로세서 구조가 제안되고 있다. 본 논문에서는 최근 새로운 마이크로프로세서 구조로 급부상하고 있는 다중처리 마이크로프로세서 구조가 객체 인식 응용에 적합한지를 분석한다. 성능 특성을 확인하기 위하여 먼저 프로그램 구동방식의 마이크로프로세서 시뮬레이터와 프로그래밍 환경을 개발하였다. 이를 기반으로 시뮬레이션을 수행한 결과, 다중처리 마이크로프로세서가 작은 오버헤드로 쓰레드 수준의 병렬성을 적절히 활용하고 있어 객체 인식 응용에 적합한 구조임을 확인하였다.

1. 서론

이미지로부터 관심 있는 객체를 인식하는 프로세스로 정의되는 객체 인식(object recognition)은 지난 30년간 컴퓨터 비전 분야의 꾸준한 연구로 공업용 검사(industrial inspection), 자동 운행(vehicle navigation), 항공사진 해석(aerial photo interpretation), 자동 목표물 인식(automatic target recognition) 등의 여러 응용에서 좋은 결과를 보여주고 있다. 그런데 이러한 객체 인식 태스크는 일반적으로 매우 많은 양의 계산을 필요로 한다. 예를 들어, 512×512 크기의 픽셀을 갖는 일련의 컬러(3 바이트/픽셀) 이미지가 표준 프레임율(30 프레임/초)로 입력되는 경우를 가정해보자. 이 경우 초당 24Mbyte의 데이터를 처리해야 하는데, 간단한 특징점 추출 알고리즘이 픽셀 당 수천번의 오퍼레이션을 필요로 하며 일반적인 비전 시스템의 경우는 더욱 복잡한 계산을 필요로 한다. 따라서 이러한 문제를 해결하기 위해서는 고성능 처리(high performance computing)가 필수적임을 알 수 있다.

이러한 이유로 미국의 HPCC(High Performance Computing and Communication) 프로젝트에서는 컴퓨터 비전을 Grand Challenge 문제 중의 하나로 정의하였는데, 이러한 Grand Challenge 문제를 해결하기 위해서는 100~1,000 billion operations/second 정도의 계산 능력을 필요로 한다. 그러나 컴퓨터 비전의 계산 특징은 다른 과학계산용 응용들과 매우 다르다. 예를 들어, 3,000×3,000×11 마일의 영역에 대한 일기예보의 경우에 100G 개의 데이터, 1000T 번의 오퍼레이션을 필요로 하며, 이를 1GFLOPS급 Cray에서 수행할 경우 약 280시간이 소요된다. 그러나 일반적인 비전 시스템의 경우, 1K×1K 크기의 픽셀 이미지에 대해서 10K~1M 개의 데이터를 처리해야 하는데, 이는 SUN 워크스테이션에서 약 한시간 정도의 계산을 필요로 한다. 즉, 일기예보의 경우 몇 일 걸리는 계산을 병렬처리에 의해 몇 시간으로(batch) 단축시키는 것인데 반하여, 비전 시스템의 경우 몇 시간 또는 몇 십분 걸리는 계산을 병렬처리로 몇 초 이내로(interactive 또는 realtime) 단축시키는 것을 그 목적으로 한다. 따라서 컴퓨터 비전에서의 성능 요구사항을 만족시키기 위해서는 다른 과학계산용 응용을 병렬화할 때

무시될 수 있었던 요소들이 컴퓨터 비전을 병렬화할 때는 매우 조심스럽게 고려되어야만 한다[1].

현재 대부분의 고성능 컴퓨터는 *Intel Pentium, Compaq Alpha21164, IBM PowerPC620, Sun UltraSparc, HP PA8000, MIPS R10000* 등 슈퍼스칼라 (superscalar) 구조[2]의 범용 마이크로프로세서를 채택하고 있다. 이러한 슈퍼스칼라 구조는 명령어 수준의 병렬성 (Instruction-Level Parallelism, ILP)을 활용하여 하나의 싸이클에 여러 개의 명령어를 수행할 수 있으나, 일반적으로 순차 프로그램내의 ILP가 크지 않은 제약에 대한 성능상의 한계를 갖고 있다[3]. 특히, 이렇게 복잡한 마이크로프로세서를 설계하는데 엄청난 비용이 든다는 사실은 슈퍼스칼라 구조를 대신할 새로운 마이크로프로세서 구조에 대한 연구를 가속화 시키고 있다[4-9]. 이중 다중처리 마이크로프로세서 (on-chip multiprocessing microprocessor)[6-8]는 ILP의 한계를 극복함과 동시에 구현의 용이성을 제공한다는 점에서 차세대 마이크로프로세서 구조로 급부상하고 있다. 즉, ILP 외에 스레드 수준의 병렬성 (Thread-Level Parallelism, TLP)[4-9]을 추가로 제공하고, 반도체 집적도가 증가함에 따라 간단한 구조의 프로세서 유닛을 여러 개 장착함으로써 슈퍼스칼라 구조의 단점을 보완할 수 있다. 그러나 이러한 다중처리 마이크로프로세서가 과학계산 응용[10]에서는 우수한 성능을 나타내는 것으로 발표되고 있지만 [6-8], 아직까지 컴퓨터 비전 응용에 적용시킨 결과는 발표되지 않고 있다. 따라서 본 논문에서는 새로운 마이크로프로세서 구조로 급부상하고 있는 다중처리 마이크로프로세서 구조가 객체 인식 응용에 적합한지를 살펴본다.

다중처리 마이크로프로세서를 이용한 객체 인식의 성능 특성을 분석하기 위하여 먼저 프로그램 구동 방식의 전용 아키텍처 시뮬레이터로서 RapSim을 개발하였다. RapSim은 명령어 시뮬레이터인 전처리기 (Pre-Processing Unit)와 각 프로세서 모델에 대한 성능 시뮬레이터인 후처리기 (Post-Processing Unit)로 구성된다. 또한, TLP를 위한 Simultaneous MultiThreading (SMT)를 지원하기 위하여 프로그래밍 환경을 개발하였다. 이를 기반으로 시뮬레이션을 수행한 결과, 다중처리 마이크로프로세서가 작은 오버헤드로 TLP를 적절히 활용하고 있어 객체 인식 응용에 우수한 성능을 제공할 수 있다는 것을 확인하였다.

본 논문의 구성은 다음과 같다. 먼저 2장에서 대상 응용인 객체 인식에 대해서 서술하고, 3장에서는 성능 분석 대상 시스템인 다중처리 마이크로프로세서를 설명한다. 그리고 성능 분석을 위해 개발된 프로그램 구동 방식의 마이크로프로세서 시뮬레이터와 프로그래밍 환경을 4장에서 기술하고, 이를 기반으로 수행한 시뮬레이션 결과를 5장에서 설명하며, 마지막으로 6장에서 결론을 맺는다.

2. 객체 인식

일반적으로 비전 시스템을 구성하는 태스크는 계산의 특성상 하위레벨 (low-level), 중간 레벨 (intermediate-level), 상위 레벨 (high-level)의 세가지 레벨로 구분된다. 하위 레벨에서 특징점 (feature) 추출은 이미징으로부터 에지 (edge)나 코너 (corner) 등과 같은 기본적인 특징점을 추출해 낸다. 그러나 센서 데이터는 일반적으로 잡음이나 계량화에 따른 에러를 수반하기 때문에, 추출된 특징점이 완벽하다고 간주할 수 없다. 이러한 어려움을 해결하기 위해 대부분의 비전 시스템은 "hypothesize & verify" 방법을 취한다. 즉, 여러 개의 코너로부터 직사각형의 가정 (hypothesis)을 생성해 내고, 상위 레벨 분석에 의하여 그러한 가정들을 선택하고 검증하여 최종적으로 객체를 인식하게 된다.

본 논문에서는 이러한 비전 시스템의 예로 DARPA Image Understanding 벤치마크[11]에서 정의한 객체 인식 시스템을 설정하였다. DARPA Image Understanding 벤치마크는 컴퓨터 비전 응용을 수행할 때 병렬 시스템의 성능을 비교하는데 많이 이용되는 벤치마크로서, 강도 (intensity) 및 범위 (range) 센서로부터 생성된 이미지가 주어졌을 때 2 1/2 차원의 "모빌 (mobile)" 구조를 인식하는 태스크들로 구성되어 있다. 여기서 두개의 센서는 상호 보완적인 역할을 하며, 하나의 센서 이미지만으로는 주어진 모빌 구조를 정확히 인식해 내지 못하게 되어 있다.

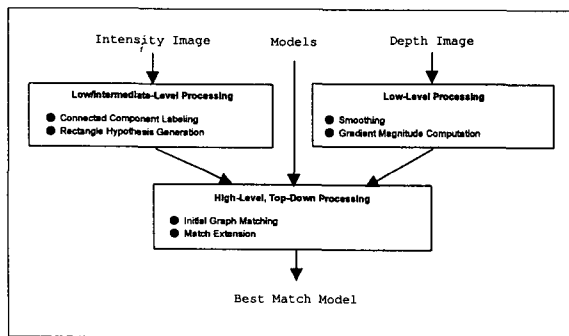
인식해야 할 모빌 구조는 여러 가지 크기, 밝기, 방위, 깊이를 가지는 2차원 직사각형의 집합으로 구성되며, 각각의 직사각형은 Z축 (시선 축)에 직각을 이루면서 상수 깊이의 표면을 갖는다. 따라서 각각의 직사각형은 따로 깊이 정보를 갖지 않으며, 깊이는 직사각형간의 공간적 관계에 기인한다 (이러한 이유로 모빌 구조가 2 1/2 차원을 갖는다고 함). 또한 장면 (scene) 내의 간섭체 (clutter)는 또 다른 직사각형으로 구성되는데, 모빌 구조와 유사한 크기, 밝기, 방위, 깊이를 갖는다. 직사각형은 다른 직사각형에 의하여 부분적 또는 전체적으로 가리워 (occlude)질 수 있는데, 하나의 직사각형은 바로 뒤에서 동일한 밝기 또는 약간 큰 깊이를 갖는 또 다른 직사각형이 존재할 때 완전히 사라질 수 있다.

모빌 인식은 인식해야 할 모빌과 유사한 구조들을 나타내는 모델들이 입력으로 주어지면 이중 주어진 장면과 가장 일치하는 모델을 찾아내는 것으로, 이때 모델은 직사각형의 크기, 방위, 깊이, 공간적 위치 등의 차이를 허용하는 대략의 표현이다. 또한 이러한 모델은 트리 구조를 갖는데, 여기서 링크는 모빌 구조의 링크를 의미하고, 노드는 한 직사각형의 크기, 방위, 깊이, 강도에 대한 정보를 가지고 있다. 따라서 한 노드의 자손을 연결하는 링크는 그 노드와 바로 밑에 있는 노드의 공간적 위치를 묘사한다.

DARPA Image Understanding 벤치마크를 구성하

면서 상정한 시나리오는 위에서 보아 보이지 않는 링크를 갖는 모빌에 다른 모빌의 일부가 날려 들어오는 것이다. 초기 상태는 전처리 과정에서 가능한 매치의 범위를 몇 개의 유사한 구조로 이미 좁혀 놓았으며, 기존의 이미지와 매치시키기 위하여 방위를 조정해 놓은 상태이다. 그러나 객체(object)가 이미 움직였기 때문에, 이 마지막 단계 전에 새로운 이미지를 취한다.

계산은 강도와 깊이 이미지에 대한 하위 레벨 비전 처리로부터 시작하여, 후보 직사각형을 추출하기 위해 강도 데이터에 대한 분류(grouping) 처리를 수행한다. 그리고 후보 직사각형들은 저장된 모델과의 부분 매치를 생성하는데 사용된다. 이때 각각의 모델에 대하여 여러 개의 가상 포즈(pose)가 생성될 수 있으며, 각각의 모델 포즈에 대하여 강도와 깊이 이미지를 탐침(probe)하기 위하여 저장된 정보를 이용한다. 여기서 각각의 탐침은 이미지의 주어진 위치에 직사각형이 존재하는지를 시험하기 위함이다. 직사각형이 실제로 존재하지 않는다는 확실한 증거가 있을 때 그 직사각형의 가정을 기각하는데, 이는 해당 모델 포즈를 삭제하는 결과를 낳는다. 반면 가정의 확인은 그 직사각형에 대한 매치의 세기를 계산하고, 새로운 크기, 방위, 위치 정보로 그 모델 포즈의 표현을 갱신한다. 여기서 "매치의 세기"는 하나의 직사각형이 다른 것에 의하여 완전히 가리워 졌을 때와 같이, 매치에 대한 어떠한 증거나 그 직사각형이 존재하지 않는다는 증거가 모두 없을 때 0이 된다. 모델 포즈 리스트 내의 매치되지 않은 모든 직사각형들에 대한 탐침이 수행되면, 제거되지 않은 각각의 포즈에 대하여 매치의 정도에 대한 평균이 계산된다. 여기서 가장 높은 평균을 갖는 모델 포즈가 최적의 매치로 간주된다. 이러한 계산 절차들을 <그림 1>에 나타낸다.



<그림 1> DARPA Image Understanding 벤치마크에서 사용한 객체 인식 시스템의 흐름도

3. 다중처리 마이크로프로세서

명령어 수준의 병렬성(Instruction-Level

Parallelism, ILP) 지원을 위한 마이크로프로세서 구조로 슈퍼스칼라(superscalar)와 Very Long Instruction Word(VLIW)가 활발히 연구되고 있으며, 현재 발표되고 있는 상용 마이크로프로세서들은 대부분 슈퍼스칼라 방식을 채택하고 있다. 슈퍼스칼라[2] 방식의 가장 큰 특징은 동일한 파이프라인 구조를 갖는 연산장치를 여러 개 병렬로 사용하여 여러 개의 명령어를 동시에 수행시키는 기법이다. 현재 사용되고 있는 상용 마이크로프로세서 중 Pentium, PA7100, PowerPC603, Alpha는 동시에 2개의 명령어를 수행하며(2-way), 4개 이상의 명령어를 동시에 수행시키는 Power2, UltraSparc, PowerPC620, PA8000 등도 최근에 개발되었다. 특히, 슈퍼스칼라 방식은 선형으로 구성된 명령어 흐름에서 하드웨어가 상호의존성이 없는 명령어를 골라내어 동시에 수행시키는 방식이기 때문에, 명령어간 상호의존성 조사를 위한 하드웨어가 필요하다. 또한, 슈퍼스칼라 방식의 성능을 극대화하기 위하여 분기 예측, 레지스터 재명명, 비순차 실행 등의 기술들이 요구된다. 따라서 슈퍼스칼라 방식은 동시에 처리하는 명령어의 수가 증가할수록 하드웨어 복잡도가 기하급수적으로 증가하여 실행 사이클 시간을 증가시킨다. 또한, 명령어 구성 비율에 따라 자원활용도의 편중이 발생하여, 일반적으로 8-way 이상에 대해서는 회의적인 전망이 지배적이다.

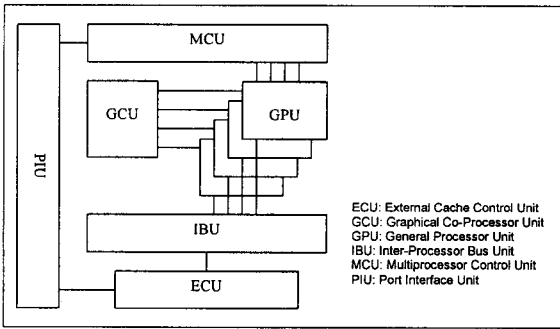
슈퍼스칼라 방식이 ILP를 위하여 하드웨어에 의한 동적 스케줄링 방식을 사용하는 반면, VLIW[12] 방식은 컴파일러에 의한 정적 스케줄링 방식을 사용한다. 즉, VLIW 방식에서는 컴파일러가 선형 명령어 흐름에서 병렬 수행 가능한 명령어들을 찾아내고, 이 명령어들을 하나의 긴 명령어로 만든다. 이때 하드웨어는 복수 개의 단위 명령어로 구성된 하나의 긴 명령어를 받아 단위 명령어들을 동시에 수행시킨다. 이러한 VLIW 방식은 슈퍼스칼라 방식에 비하여 하드웨어 복잡도가 낮은 장점이 있지만, 컴파일러 개발에 어려움이 많고 마이크로프로세서 상호간 호환성을 보장할 수 없는 단점이 있다.

따라서 이러한 ILP 기반의 마이크로프로세서 구조를 대신할 새로운 구조에 대한 연구가 활발히 진행중에 있으며[4-9], 이중 다중처리 마이크로프로세서(on-chip multiprocessing microprocessor)[6-8]는 ILP의 한계를 극복함과 동시에 구현의 용이성을 제공한다는 점에서 차세대 마이크로프로세서 구조로 급부상하고 있다. 즉, 간단한 구조의 프로세서 유닛트를 이용하여 한 쓰레드 내의 크기 않은 ILP를 지원하면서, 여러 개의 프로세서 유닛트를 이용하여 여러 개의 쓰레드를 동시에 실행시킴으로써 쓰레드 수준의 병렬성(Thread-Level Parallelism, TLP)[4-9]을 추가로 제공한다.

본 논문에서는 다중처리 마이크로프로세서로 Raptor[13]를 사용한다. Raptor는 한국전자통신연구원에서 차세대 마이크로프로세서로 개발한 단일 칩 다중처리 마이크로프로세서로, General Processor Unit(GPU)라는 4개의 독립적인 프로세서 유닛트와 Graphic Co-processor Unit(GCU)라는 하나의 그래픽 연산 유닛

트로 구성된다. Raptor의 주요 특징은 다음과 같다:

- 각각의 프로세서 유닛은 SPARC V9 명령어 세트 구조를 가지며, 64비트 정수 및 부동 소수 점 연산을 한다.
- 4개의 독립된 프로세서 유닛이 그래픽 연산 유닛을 공유하여 그래픽 연산을 처리한다.
- 크기가 비교적 작은 내장 1차 캐쉬와 대용량의 외부 2차 캐쉬로 구성된 다중 캐쉬 구조를 갖는다.
- 2차 캐쉬 제어기를 내장하며 여러 개의 Raptor 프로세서를 이용한 시스템 구현을 지원한다.



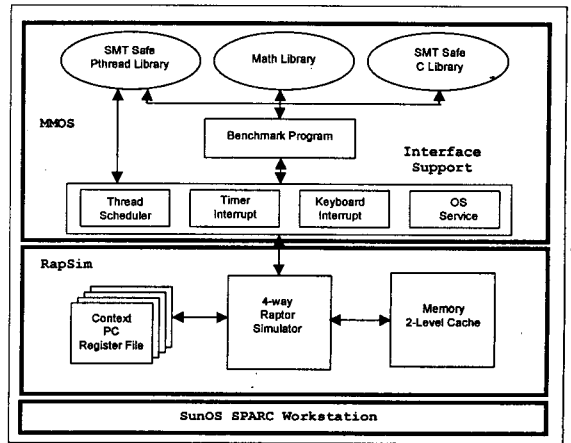
<그림 2> Raptor 마이크로프로세서의 블록 구성도

이러한 특징을 갖는 Raptor의 블록 구성도를 <그림 2>에 나타내었다. 여기서 Inter-processor Bus Unit(IBU)은 GPU와 External cache Control Unit(ECU)을 연결하는 공유 버스이다. 그리고 Multiprocessor Control Unit(MCU)은 인터럽트를 GPU로 분배하고 GPU간의 동기를 지원한다. 반면, Port Interface Unit(PIU)는 칩 내부에서 발생하는 메모리 접근 요구, 입출력 장치 접근 요구, 인터럽트 및 기타 유틸리티 신호선 사용 요구를 처리한다. 특히, 부가적인 외부 로직없이 다중처리 시스템을 구성할 수 있는 multiprocessor-ready 형태의 인터페이스를 제공한다. 그리고 4개의 GPU는 각각의 레지스터 파일과 프로그램 카운터로 그래픽 명령어를 제외한 모든 명령어를 실행하지만, IBU를 통하여 ECU를 공유한다. 마지막으로 GCP는 4개의 GPU에 의해 Single Instruction stream Multiple Data stream(SIMD) 방식으로 그래픽 명령어를 처리한다. Raptor에 대한 보다 자세한 내용은 [13]에 설명되어 있다.

4. 시뮬레이션 환경

다중처리 마이크로프로세서인 Raptor를 정성적으로 분석하기 위하여 전용 시뮬레이터인 RapSim을 개발하였

다. RapSim은 4개의 GPU와 이들에 의해 공유되는 메모리 계층을 모델링하는 프로그램 구동 방식의 마이크로아키텍처 시뮬레이터이다. 여기서 각각의 GPU 모델은 명령어 세트 시뮬레이터인 전처리기(Pre-Processing Unit)와 성능 시뮬레이터인 후처리기(Post-Processing Unit)로 구성된다. 또한 GPU간의 SMT를 지원하기 위하여 MMOS(Multithreaded Mini-OS)라고 하는 프로그래밍 환경을 개발하였다. RapSim과 MMOS의 전체적인 관계를 <그림 3>에 나타내었다.



<그림 3> 시뮬레이션 환경

4.1. RapSim

RapSim의 주요 특징은 다음과 같다:

- SPARC V9 명령어와 그래픽 연산 유닛용 명령어의 수행
- 시간 정보를 갖는 프로그램 구동 방식의 시뮬레이터
- 4개의 프로세서 유닛 및 1개의 그래픽 연산 유닛으로 구성된 다중처리 모델
- 프로세서 유닛의 비순차(out-of-order) 실행 지원
- SMT 프로그래밍 모델 지원
- 성능 측정을 위한 정보 수집

먼저 전처리기는 명령어를 수행하는 프로세서 모델, 레지스터 파일을 위한 자료구조, 시스템호출 처리를 위한 프록시(proxy) 모델, 1차 캐쉬 모델 등을 갖는 명령어 세트 시뮬레이터이다. 즉, 전처리기는 2차 캐쉬 모델을 포함하는 공유메모리 계층으로부터 명령어와 데이터를 접근하고, 그 명령어를 수행하며, 후처리기가 수행하는 트레이스를 동적으로 생성해낸다.

전처리기는 컴파일되어 정적으로 MMOS 라이브러리에 링크된 벤치마크 실행파일을 로딩함으로써 시뮬레

이션을 개시한다. 이러한 로딩 동작 중에 시작 프로그램 카운터가 적절한 값으로 세트되고, 트랩 테이블 및 트랩 핸들러가 초기화되며, 스택이 메모리 모델에 만들어진다. 그리고 프로세서 모델은 실행기, 레지스터 파일, 1차 캐쉬 등의 내부 자원을 이용하여 실행을 시작한다.

전처리기는 명령어 스트림을 실행하는 동안, 수행된 명령어의 순서로 정의된 트레이스를 동적으로 생성해 낸다. 여기서 트레이스 각각의 엔트리는 충분한 정보를 포함하여, 후처리가 그 트레이스를 입력으로 받아 시뮬레이션을 수행할 수 있어야 한다. 즉, 후처리는 전처리가 생성한 명령어 트레이스를 이용하여 시뮬레이션을 수행하는 RISC 파이프라인 모델이다. 이때 후처리는 비순차 실행을 지원하기 위해 Reservation Station(RS)과 Re-Order Buffer(ROB)를 포함하는 2-way 슈퍼스칼라로 모델링된다.

하나의 사이클에 트레이스 버퍼에 있는 2개의 명령어가 접근되어 디코드를 위한 전처리 절차를 수행한다. 명령어 버퍼 내에서 이렇게 전처리된 명령어는 디코드되어 적절한 RS에 할당되고, ROB는 동시에 갱신된다. 이때 각각의 실행기는 적절한 RS로부터 데이터 종속문제를 해결한 명령어를 수행한다. 수행 결과는 ROB에 반영되고, ROB에서 종료된 엔트리는 레지스터 파일로 갱신된다. 본 논문에서 사용한 구체적인 아키텍처 변수의 기본값을 <표 1>에 나타내었다.

<표 1> 아키텍처 변수

항목	기본값
GPU 수	1, 2, 4
GPU 이슈 크기	2
1차 캐쉬 크기	16KB I-cache, 16KB D-cache
2차 캐쉬 크기	4MB
1차 캐쉬 일관성 유지	Write Through
2차 캐쉬 일관성 유지	Write Back

4.2. MMOS

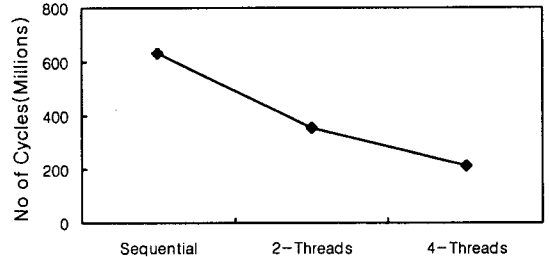
MMOS는 4개의 GPU를 효과적으로 사용하기 위하여 RapSim에 SMT 프로그래밍 환경을 제공한다. 즉, MMOS는 SMT-safe Pthread 라이브러리, SMT-safe C 라이브러리, RapSim 인터페이스를 제공한다.

여기서 SMT-safe C 라이브러리는 여러 개의 쓰레드가 공유 C 라이브러리를 동기문제 없이 액세스할 수 있게 해주고, SMT-safe Pthread 라이브러리는 여러 개 쓰레드간의 동기 및 스케줄링 기능을 제공해 준다. 그리고 RapSim 인터페이스는 MMOS를 RapSim에 연결하여 4개의 쓰레드로 4개의 프로세서를 시뮬레이션할 수 있게 하고, 쓰레드를 스케줄하여 RapSim의 프로세서 모델로 할당한다.

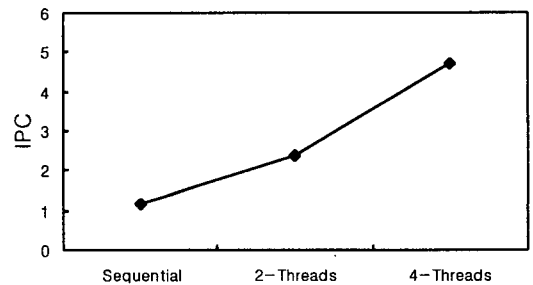
5. 시뮬레이션 결과

본 논문에서는 1개의 GPU상에서 수행한 순차 실행(Sequential), 2개의 GPU상에서 2개의 쓰레드를 이용한 실행(2-Threads), 4개의 GPU상에서 4개의 쓰레드를 이용한 실행(4-Threads)에 대한 시뮬레이션을 수행하였다. 그리고 다중처리 마이크로프로세서의 성능지수로는 수행 사이클 수, Instruction Per Cycle(IPC), 병렬 오버헤드를 설정하였다.

먼저 쓰레드 수를 증가시킬 때 수행 사이클 수와 IPC의 변화를 살펴보자. <그림 4>, <그림 5>에서 알 수 있듯이, 쓰레드 수를 증가시키면 수행 사이클 수는 선형적으로 감소하고 반대로 IPC는 선형적으로 증가한다. 즉, 작은 오버헤드로 쓰레드 수준의 병렬성을 적절히 활용하여 선형적인 성능 향상을 기대할 수 있다. 여기서 병렬 쓰레드 이용시 사이클 수는 parent process를 수행하는 GPU의 사이클 수를 의미한다.



<그림 4> 순차 및 쓰레드를 이용한 실행시 수행 사이클 수 비교



<그림 5> 순차 및 쓰레드를 이용한 실행시 IPC 비교

다음에 쓰레드 수준의 병렬성을 활용하는데 따른 오버헤드를 분석하기 위하여 다음과 같은 오버헤드 함수를 정의한다. 즉,

$$\text{Overhead}(1, 2) = (W2 - W1) / W1 \times 100\%$$

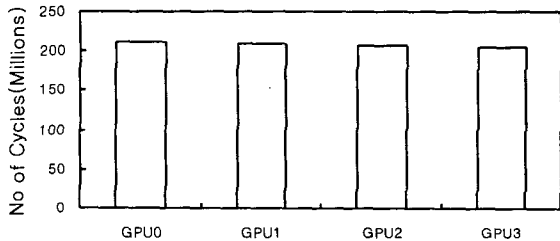
여기서 W_i 는 실행환경 i 에서의 작업부하라고 하며, 각 GPU에서 실행한 수행 사이클의 합으로 정의된다. 이러한 오버헤드 함수를 요약하면 <표 2>와 같다. 즉,

쓰레드 수를 증가 시킬수록 쓰레드 동기 등의 오버헤드에 따른 오버헤드 함수값이 증가하지만 그 값이 크지 않음을 알 수 있다.

<표 2> 병렬 오버헤드

구분	값
Overhead(Sequential, 2-Threads)	11.4%
Overhead(2-Threads, 4-Threads)	18.0%

그리고 여러 개의 쓰레드를 이용할 경우 GPU의 활용도(utilization)를 확인하기 위하여 4개 쓰레드 이용한 실행시 수행 사이클 수를 <그림 6>에 나타내었다. GPU0가 순차실행 부분 때문에 사이클 수가 약간 많지만, 대체로 GPU간 균형이 이루어짐을 알 수 있다.



<그림 6> 4개 쓰레드를 이용한 실행시 GPU간 수행 사이클 수의 분포

6. 결론

다중처리 마이크로프로세서는 현존하는 슈퍼스칼라 구조가 갖는 명령어 수준 병렬성의 한계를 극복할 수 있으며, 수십억개의 트랜지스터가 집적된 차세대 반도체에 적합한 구조로 최근 급부상하고 있다. 즉, 하나의 마이크로프로세서 내에 간단한 프로세서 유닛을 여러 개 장착하여, 명령어 및 쓰레드 수준의 병렬성을 모두 제공할 수 있는 구조를 제공한다.

본 논문에서는 이러한 다중처리 마이크로프로세서를 이용하여 객체 인식 응용을 수행할 때 어떠한 성능 특성을 나타내는지 분석하였다. 성능 분석을 위하여 먼저 마이크로프로세서 시뮬레이터와 프로그래밍 환경을 개발하였다. 이를 기반으로 프로그램 구동 시뮬레이션을 수행하여 수행 사이클 수, IPC, 병렬 오버헤드를 분석하였다. 시뮬레이션 수행 결과, GPU 수가 증가함에 따라 IPC도 거의 선형적으로 증가함을 확인하였다. 즉, 4개의 GPU를 사용할 경우 4.7 이상의 IPC를 얻을 수 있었다. 반면 GPU 수를 증가시킬수록 쓰레드 동기 등의 병렬 오버헤드가 증가함을 확인하였다.

끝으로 본 논문에서는 하나의 다중처리 마이크로프로세서에 객체 인식 작업부하를 부과한 시뮬레이션 결과

를 나타내었는데, 여러 개의 다중처리 마이크로프로세서가 대규모 병렬머신으로 구현된 경우의 성능도 향후 흥미로운 연구분야가 될 것으로 기대된다.

참고문헌

- [1] C. Weems, "Architectural Requirements of Image Understanding with respect to Parallel Processing", *IEEE Proceedings*, Vol. 79, No. 4, pp. 537-547, 1991.
- [2] M. Slater, "The Microprocessor Today," *IEEE Micro*, Vol. 16, No. 6, pp. 32-45, 1996.
- [3] D.Wall, "Limits of Instruction Level Parallelism", *WRL Research Report*, Digital Western Research Laboratory, 1993.
- [4] J. Wilson, "Challenges and Trends in Processor Design," *IEEE Computer*, Vol. 31, No. 1, pp. 39-50, 1998.
- [5] S. Egger, et al., "Simultaneous Multithreading: A Platform for Next-Generation Processor," *IEEE Micro*, Vol. 17, No. 5, pp. 12-19, 1997.
- [6] B. Nayfe, L. Hammond, and K. Olukotun, "Evaluation of Design Alternatives for Multiprocessor Microprocessor," *Proc. of Int'l Symp. on Computer Architecture*, pp. 66-77, 1996.
- [7] L. Hammond, et al., "A Single-Chip Multiprocessor," *IEEE Computer*, Vol. 30, No. 9, pp. 79-85, 1997.
- [8] S. Amarashinhe, et al., "Multiprocessors From a Software Perspective," *IEEE Micro*, Vol. 16, No. 3, pp. 52-61, 1996.
- [9] A. Agarwal, et al., "Performance Tradeoff in Multithreading Processors," *IEEE Tr. on Parallel and Distributed Systems*, Vol. 3, No. 5, pp. 525-539, 1992.
- [10] J. Singh, W. Weber, and A. Gupta, "SPLASH: Stanford Parallel Applications for Shared Memory," *Computer Architecture News*, Vol. 20, No. 1, pp. 5-44, 1992.
- [11] C. Weems, et al., "The DARPA Image Understanding Benchmark for Parallel Computers", *J. of Parallel and Distributed Computing*, Vol. 11, No. 1, pp. 1-24, 1991.
- [12] J. Fisher, "Very Long Instruction Word Architecture and the ELI-512", *Proc. of Int'l Symp. on Computer Architecture*, pp. 140-150, 1983.
- [13] 박경 외 3인, "다중처리 마이크로프로세서 설계," *대한전자공학회 추계학술대회 논문집*, pp. 717-720, 1996.