

이기종 시스템에서의 캐시 최적화 연구 동향

양서현¹, 김나현¹, 김예진¹, 김지훈¹, 김현지¹, 오현영^{2*}

¹가천대학교 AI 소프트웨어학부 학부생

²가천대학교 AI 소프트웨어학부 교수

qk325672@gachon.ac.kr, nh5102@gachon.ac.kr, vvvsfhk@naver.com, jihungim500@gmail.com,
hjkim1770@gmail.com, hyoh@gachon.ac.kr

Cache Optimization Research Trends in Heterogeneous Systems

Seo-hyeon Yang, Na-hyeon Kim, Ye-jin Kim, Ji-hun Kim, Hyeon-ji Kim, Hyunyoung Oh
Dept. of AI · Software, Gachon University

요 약

최근 이기종 메모리 및 캐시 구조가 컴퓨터 시스템에서 자주 활용되고 있다. 그러나 이기종 환경에서 기존의 균일한 캐시 관리 방식을 사용하여 성능을 향상하는 데에는 한계가 있으며, 이기종 시스템의 캐시 구조를 고려한 최적화 방법이 필요하다. 본 논문에서는 이기종 시스템에서의 캐시 최적화 방법 동향에 대해 살펴본다. 특히, 데이터 배치 최적화, 동적 캐시 할당, 그리고 에너지 효율적인 메모리 아키텍처에 초점을 맞춘 최신 연구들을 분석한다.

1. 서론

최근의 컴퓨팅 시스템에서는 이기종(heterogeneous) 메모리 및 캐시 구조의 활용이 급격히 증가하고 있다. 예를 들어, 인텔의 Optane DC 영구 메모리는 DRAM과 SSD 사이의 성능 특성을 제공하며, 기존 메모리 계층 구조에 새로운 층을 추가한다. 이러한 이기종 메모리 및 캐시 구조는 다양한 지연 시간과 대역폭 특성을 가진 메모리 구성 요소를 통합하여, 기존의 균일한 캐시 관리 방식으로는 성능 향상에 한계가 있음을 시사한다. 본 논문에서는 이기종 메모리 및 캐시 구조를 고려한 성능 최적화 연구의 두 가지 주요 동향을 살펴보려 한다: 이기종 캐시 및 메모리 환경에서 데이터 배치를 최적화하는 접근 방식[1,2]과 이기종 캐시 구조를 적극적으로 활용하여 성능과 에너지 효율성을 향상시키는 연구[3,4]이다.

2. 이기종 시스템(Heterogeneous System)

이기종 시스템(Heterogeneous System)은 서로 다른 특성을 가진 다양한 컴퓨팅 요소들을 단일 시스템 내에 통합한 컴퓨팅 아키텍처를 의미한다. 이러한 시스템은 각 컴포넌트의 강점을 활용하여 전체 시스템의 성능, 에너지 효율성, 그리고 특정 작업에 대한 최적화를 달성하는 것을 목표로 한다. 주요 구성 요소로는 범용 프로세서(CPU), 그래픽 처리 장치(GPU), 특수 목적 가속기(예: TPU, DSP), FPGA, 그리고 다양한

유형의 메모리(SRAM, DRAM, HBM, NVM 등)가 있다. 실제 예시로는 모바일 SoC, 데이터 센터 서버, 자율주행 차량 컴퓨팅 시스템 등이 있다. 이기종 시스템의 주요 장점은 작업별 최적화된 성능, 에너지 효율성, 그리고 유연성이지만, 복잡한 프로그래밍 모델, 리소스 관리의 복잡성, 메모리 일관성 유지 등의 도전 과제도 존재한다.

3. 최근 이기종 환경이 고려된 캐시 최적화 연구

3.1 CachedArrays: 이기종 메모리 시스템에서의 데이터 이동 최적화

Mark Hildebrand et al.이 제안한 CachedArrays[1]은 이기종 메모리 시스템에서 데이터 이동을 최소화하고 대규모 워크로드에서 하드웨어 캐시의 한계를 극복하기 위한 기법이다. CachedArrays의 핵심 아이디어는 애플리케이션의 데이터 접근, 데이터 이동 메커니즘, 그리고 데이터 이동 정책을 분리하는 것이다. 이를 통해 애플리케이션은 데이터 이용 메커니즘과 직접 상호작용하지 않고도 의미적 정보를 활용할 수 있게 된다. CachedArrays는 데이터 관리자를 통해 메모리 영역의 할당 및 해제, 영역 간 고성능 메모리 복사, 그리고 객체의 기본 영역 재할당을 지원한다. 특히, 애플리케이션이 "이 객체를 읽을 예정이다"와 같은 힌트를 제공하면, 데이터 관리 API를 통해 객체와 영역의 상태를 효율적으로 조작한다. CachedArrays는 캐시 라인 단위가 아닌 객체 단위로 데이터를 이동함으

* 교신저자

로써 더 효율적인 메모리 사용을 가능하게 한다. DRAM 과 NVRAM 이 혼합된 실제 하드웨어 시스템에서 CNN 학습 작업을 수행한 결과, CachedArrays 는 기존 방식 대비 1.4~2 배의 성능 향상을 달성했다.

3.2 Chopt(Choice-aware OPT): 오프라인 데이터 배치 최적화

Lei Zhang et al.이 제안한 Chopt[2]는 이기종 메모리 환경에서의 오프라인 데이터 배치 최적화 알고리즘이다. Chopt 는 두 가지 주요 관찰에 기반한다: 첫째, 요청된 데이터를 항상 더 빠른 메모리에 캐싱할 필요는 없으며 때로는 느린 메모리에서 직접 접근하는 것이 효율적일 수 있다는 점, 둘째, 단순히 캐시 미스 비율을 최소화하는 기존 방식은 최적의 지연 시간 관리를 하지 못한다는 점이다. Chopt 는 메모리 계층 구조에서의 데이터 배치 결정을 네트워크 흐름 문제로 모델링하고, 이를 MCMF(Minimum-Cost Maximum-Flow) 문제로 변환하여 해결한다. 알고리즘은 각 데이터 요청에 대한 지연 시간 절감을 계산하여 최적의 데이터 배치를 결정하며, L1 캐시부터 시작하여 각 메모리 계층의 사용률을 최대화하는 것을 목표로 한다. 실제 작업 환경에서 Chopt 를 테스트한 결과, 선행 연구들과 비교했을 때 평균 요청 지연 시간을 8.2%에서 44.8%까지 개선했다.

3.3 McCore: 하이엔드 이기종 멀티코어 시스템을 위한 동적 캐시 관리

Jaewon Kwon et al.이 제안한 McCore[3]는 하이엔드 이기종 멀티코어 시스템에서 하드웨어 기반 강화 학습(RL) 스케줄러를 통해 메모리 계층 구조를 동적으로 재구성하고 최적화하는 기법이다. McCore 는 하드웨어 성능 카운터와 소프트웨어 프로파일링을 사용하여 캐시 액세스를 실시간으로 모니터링한다. 이렇게 수집된 데이터를 분석하여 다양한 워크로드의 캐시 활용 패턴을 식별하고, 이에 기반하여 각 코어에 할당된 캐시의 크기를 동적으로 조정한다. 이를 통해 고성능 코어는 필요 시 더 많은 캐시를 확보하고, 저전력 코어는 캐시 할당을 줄여 전체적인 성능과 에너지 효율을 최적화한다. McCore 는 Shadow Indexing(실제 캐시의 동작을 모방하는 병렬 그림자 구조), Sub-ranked DRAM(단일 버스트에 두 개의 서로 다른 캐시 라인을 채움), Adaptive granularity memory systems(워크로드에 따라 메모리 블록의 크기와 구조를 동적으로 조정) 등의 기술을 활용하여 메모리 액세스를 세밀하게 제어한다. 이러한 기법들을 통해 McCore 시스템은 기존 기법 대비 평균 25%의 성능 향상을 달성했다.

3.4 COMPAD(Compaction + Scratchpad): 임베디드 시스템을 위한 에너지 효율적 메모리 아키텍처

Tommaso Marinelli et al.이 제안한 COMPAD[4]는 임베디드 시스템에서의 에너지 효율성 개선을 위한 스크래치패드 메모리(SPM) 기반 아키텍처이다. 전통적인 캐시 시스템이 데이터 교체와 불필요한 에너지 소비로 인해 임베디드 시스템에서 비효율적이라는 점에

착안하여 개발되었다. COMPAD 는 자주 접근되는 데이터를 SPM 에 저장하여 에너지 소비를 줄이는 것을 목표로 한다. 연구진은 메모리 접근 패턴을 분석하여 특정 데이터가 반복적으로 접근되는 과정에서 불필요한 데이터가 포함되어 있음을 발견하고, 필요한 데이터만을 저장하는 방식으로 메모리 레이아웃을 압축하는 기법을 제안했다. 또한, 데이터의 재사용 가능성과 접근의 연속성을 평가하여 SPM 으로 매핑할 데이터를 선정하는 알고리즘을 개발했다. 루프 기반 벤치마크를 사용한 실험 결과, COMPAD 아키텍처는 전통적인 캐시 전용 솔루션과 비교하여 동적 에너지 소비 측면에서 평균 43%의 개선을, 전체적으로는 13%의 성능 향상을 달성했다.

4. 결론 및 향후 계획

본 논문에서 분석한 연구들은 이기종 시스템에서의 캐시 최적화가 단순한 캐시 미스 비율 최소화를 넘어 시스템의 전체 성능과 효율성을 극대화하는 방향으로 발전하고 있음을 보여준다. 그러나 향후 연구에서는 몇 가지 중요한 과제들이 남아있다. 첫째, 다양한 워크로드와 하드웨어 구성에서의 포괄적인 성능 평가가 필요하다. 이는 제안된 기법들의 일반화 가능성을 검증하는 데 중요하다. 둘째, 전력 소모와 열 관리 측면에서의 최적화 연구가 요구된다. 특히 에너지 효율성이 중요한 환경에서 이는 핵심적인 고려사항이 될 것이다. 셋째, 온라인과 오프라인 환경에서 모두 효과적인 하이브리드 접근 방식의 연구가 필요하다. 마지막으로, 복잡한 이기종 시스템에서의 동적 최적화를 위한 기계 학습 기반 접근 방식의 적용이 요구된다. 이러한 과제들을 해결함으로써, 미래의 이기종 컴퓨팅 시스템에서 더욱 효율적이고 강력한 캐시 최적화 기법들이 개발될 수 있을 것으로 기대된다.

사사문구

이 논문은 2024 년도 정부(산업통상자원부)의 재원으로 한국 산업기술기획평가원의 지원(No. RS-2024-00406121, 자동차보안취약점기반위험분석시스템개발(R&D))과 정부(과학기술정보통신부)의 재원으로 한국연구재단의 지원(No. RS-2022-00166529)을 받고 과기정통부 정보통신기획평가원의 정보보호핵심원천기술개발사업(No. RS-2024-00337414)으로 수행한 결과임.

참고문헌

- [1] Mark Hildebrand et al., "CachedArrays: Optimizing Data Movement for Heterogeneous Memory Systems," IPDPS, 2024
- [2] Lei Zhang et al., "Optimal Data Placement for Heterogeneous Cache, Memory, and Storage Systems," Proceedings of the ACM on Measurement and Analysis of Computing Systems, 2020
- [3] Jaewon Kwon et al., "McCore: A Holistic Management of High-Performance Heterogeneous Multicores," MICRO, 2023
- [4] Tommaso Marinelli et al., "COMPAD: A heterogeneous cache-scratchpad CPU architecture with data layout compaction for embedded loop-dominated applications," Journal of Systems Architecture, 2023