

AI 모델 학습 시 GPU 소비 전력과 최적화 방안 분석

조민성 권오현 신동건(멘토)

대전대신고등학교(Daejeon Daeshin High.) A.C.T.(KE)

Analysis of GPU power consumption and optimization strategies during AI model training

ABSTRACT: 본 연구는 국가 디지털 트윈 플랫폼(V-World)에서 제공하는 고해상도 DEM 을 기반으로 오픈소스 GIS 소프트웨어(QGIS)와 공간분석 도구를 활용하여, 학교(대전대신고등학교) 인접 사면의 산사태 및 토석류 위험도를 정량적으로 평가하고 우선순위에 따라 실무적 공학적 완화 대책을 제안한다. 분석 인자로는 경사도(slope), 유량 집적(flow accumulation), 지형 습윤 지수(TWI)를 사용하였으며, 각 인자는 국내 설계·안전 기준과 통계적 분포를 참고해 1(매우 안전)~5(매우 위험)로 점수화하였다. 중첩(가중합)으로 최종 위험지수를 산출하고, 현장 적용 가능성을 고려한 공학적 대책(사방댐/체크댐, 옹벽/유도 배수로, 표면배수 및 식생복구 등)을 설계 기준과 함께 제시한다. 결론적으로 V-World 기반 디지털 트윈 자료는 예방적 학교 안전관리와 행정적 의사결정에 실무적 근거를 제공할 수 있음을 보였다.

1. 서론

- 문제 및 연구 배경

답러닝 모델의 규모와 활용이 폭발적으로 증가하면서, 이러한 모델들을 학습시키는 데 필요한에너지 소비와 GPU 전력 소모 문제가 중요한 화두로부상하고 있다. 특히 GPT-3 와 같은 초거대 모델학습에는 막대한 전기가 소모되었으며, GPT-3 를학습시킨 GPU들은 약 1,300 MWh의 전력을 사용한것으로 추산되는데 이는 미국 가정 1,450 가구의 한달치 전력 소비량에 맞먹는 수준이다. 데이터센터분야에서도 AI 모델 학습으로 인한 전력 수요가급증하여, 2030 년에는 데이터센터가 세계 전력의 21%까지 소비할 것이라는 전망도 나온다. 이러한배경에서, "AI 모델 학습 시 GPU 소비 전력과 최적화방안 분석"에 대한 심층 연구는 지속가능한 AI를구현하기 위한 필수 과제가 되었다.

- 연구 목적

본 탐구에서는 최신 논문과 신뢰할 수 있는 기술 보고서를 바탕으로, 딥러닝 모델 학습 시 GPU 의 전력 소비 특성과 이를 최적화하는 방법들을 체계적으로 분석한다. 주요 GPU 아키텍처(A100, V100, H100, RTX 3090 등)의 소비 전력을 비교하고, GPU 전력소모의 원인 요인을 심충적으로 살펴볼 것이다. 또한모델 규모, 학습시간, 배치 크기 등이 전력 소비에미치는 영향을 고찰하며, 하드웨어/소프트웨어측면에서의 전력 최적화 전략을 정리한다. 마지막으로 DeepMind, OpenAI, Meta 등의 실제사례를 통해 이론적 최적화 기법들이 산업 현장에서어떻게 적용되고 있는지 알아보고자 한다. 이를 통해고등학생 수준의 보고서이지만 내용의 전문성은대학수준으로 유지하면서, 독자들이 GPU 전력 효율향상의 중요성과 구체적 방법을 이해할 수 있도록하는 것이 본 연구의 목표이다.

- 연구 범위 및 방법

본 연구는 문헌 조사를 기반으로 한 정성적 메타분석 형태로 수행되었다. 최신 AI 하드웨어의 스펙 비교부터 딥러닝 학습의 전력 최적화 기법까지 폭넓은 주제를 다루기 위해 체계적인 분석 절차를 적용하였다.

먼저 자료 수집 단계에서는 2023 년 이후 발표된학술 논문(IEEE, USENIX, arXiv), 기업 기술보고서(NVIDIA, Meta), 그리고 기관 연구자료(MIT Lincoln Laboratory 보고서 등)를 폭넓게 수집하였다. 특히 GPU 전력 측정 데이터와 전력 최적화 효과를정량적으로 다룬 자료를 우선적으로 선별하였다.

다음으로 내용 분류 단계에서는 수집된 자료를 본연구의 다섯 가지 세부 주제에 따라 체계적으로 분류하였다. (1) GPU 모델별 소비 전력 비교, (2) GPU 전력 소모의 주요 원인, (3) 모델 규모 및 학습 조건과 소비 전력 간의 상관관계, (4) 하드웨어 및 소프트웨어 수준의 전력 최적화 전략, (5) 실제 산업 적용 사례가이에 해당한다.

심층 분석 단계에서는 각 분야별 핵심 데이터를 발췌하고 상호 보완적으로 검토하였다. 예를 들어 GPU 별 열설계전력(TDP)과 MLPerf 벤치마크를 활용하여 성능 대비 전력 효율(performance-per-



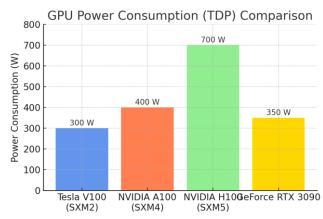
watt)을 계산하였으며, 각 논문에서 제시된 최적화 기법의 에너지 절감률을 비교 분석하였다.

마지막으로 검증 및 종합 단계에서는 인용된 자료의 신뢰성을 검토하고, 동일 주제에 대해 여러 출처의 결과가 일관되는지를 교차 확인하였다. 이러한 검증 과정을 통해 도출된 결론을 바탕으로 각 주제별 주요 논점을 종합 정리하고, 이해를 돕기 위한 시각 자료를 함께 구성하였다.

2. 본론

- GPU 모델별 소비 전력 비교 분석

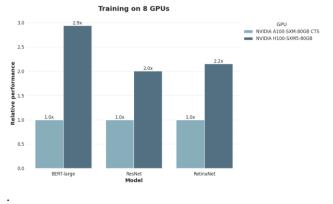
답러닝 학습에 널리 사용되는 주요 GPU 들인 NVIDIA Tesla V100 (Volta), A100 (Ampere), H100 (Hopper) 및 소비자용 RTX 3090 (Ampere)의 전력 소모 특성을 비교하면 다음과 같다. 각 GPU 의 공식적인 TDP(Thermal Design Power, 최대소비전력)를 비교한 결과, H100(SXM5)의 TDP 는 무려 700W 에 달하여 이전 세대인 A100(SXM4)의 400W 보다 훨씬 높고, V100(SXM2)은 300W, RTX 3090 은 350W 수준임을 확인할 수 있다. 아래 그림은 이러한 주요 GPU 들의 전력 한도를 시각화한 것이다.

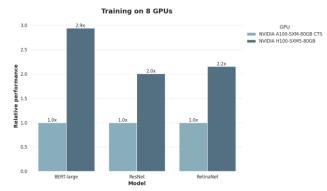


주요 GPU 모델의 공칭 최대전력 비교 (V100, A100, H100 은 서버용 SXM 모델 기준, RTX 3090 은 소비자용 제품 기준). H100 의 최대전력은 A100 대비약 1.75 배에 이르지만, 이를 통해 얻는 성능 향상도 크다.

서버급 GPU 인 V100/A100/H100 의 경우 PCIe 카드형과 SXM 모듈형에 따라 전력 제한이 다소 다른데, 예를 들어 A100 PCIe 버전은 250W, SXM 버전은 400W 로 동작한다. 마찬가지로 H100 도 PCIe 는 약 350W, SXM 은 700W 까지 소비한다. 한편 V100은 PCIe 형 250W, SXM 형 300W로 설계되었으며, 메모리 용량이 두 배로 늘어난 32GB SXM2 모델의 경우 350W 로 TDP 가 상승했다. RTX 3090 소비자 GPU 는 공칭 350W 로, 고성능이지만 발열과 전력소모가 매우 큰 편이다.

그러나 단순 전력 수치만으로 GPU 효율을 평가할 수는 없다. 최신 GPU 일수록 더 많은 전력을 쓰지만 성능(연산 처리량)도 대폭 향상되어. 와트당 연산능력(performance-per-watt) 지표에서는 오히려 효율이 좋아지는 경향을 보인다. 예를 들어 NVIDIA Hopper(H100) 아키텍처는 Ampere(A100) 대비 와트당 성능이 최대 3 배에 이른다고 알려져 있다. 실제 MLPerf 벤치마크에서 H100 은 A100 대비 약 60% 높은 에너지 효율을 보이며, FP16 기준 A100 이 와트당 최대 10TFOPS 를 처리할 때 H100 은 와트당 20TFOPS 까지 달성한다. 즉, H100 은 전력을 75% 더 소모하지만 연산 성능은 2~3 배에 달해 종합적인 전력 효율이 향상되는 것이다. 예를 들어 ResNet-50 학습 작업에서 H100 은 A100 보다 약 1.6 배 높은 효율(동일 전력당 처리 이미지 수 기준)을 보인다. BERT 대형 모델 학습에서는 H100 이 A100 대비 2.9 배 빠른 학습 속도를 내는 사례도 보고되었다.





MLPerf 기준 AI 모델 학습 성능 비교: 동일한 8-GPU 시스템에서 A100(밝은색)을 1.0 배로 볼 때 H100(짙은색)은 BERT-Large 학습에서 약 2.9×,



ResNet 에서는 2.0×, RetinaNet 에서는 2.2×의 성능을 달성했습니다. 향상된 성능 대비 소비 전력을 고려하면 H100 의 와트당 성능은 A100 대비 크게 향상되었음을 알 수 있습니다.

요약하면, H100 은 전력 소모가 크지만 그 이상의성능 향상을 통해 전체 효율을 높였고, A100 도전세대 V100 대비 효율이 대폭 개선되었습니다. 반면 RTX 3090 등 소비자용 GPU 는 절대전력은 높지만데이터센터 GPU 대비 최적화가 덜되어 와트당성능은 낮은 편입니다. 이러한 비교는 GPU 선택 시 "절대 성능"뿐 아니라 "에너지 효율"도 고려해야함을 시사합니다.

- GPU 전력 소모의 주요 워인

GPU 의 전력 소모는 여러 요인이 복합적으로 작용한 결과이다. 대규모 병렬 연산을 수행하는 GPU 의 구조적 특성상, 연산 장치의 동시 활성화와 빈번한 메모리 접근이 전력 소비의 주요 원인으로 작용한다.

먼저 병렬 처리 연산 부하는 GPU 전력 소모의 핵심 요소이다. 딥러닝 학습 과정에서는 수만 개의 CUDA 코어와 Tensor 코어가 동시에 매트릭스 곱셈 등대규모 연산을 수행한다. 특히 Transformer 계열모델은 거대한 행렬 연산을 반복적으로 수행하기때문에 모든 코어가 동시에 동작하며 큰 전류가흐른다. 이에 따라 GPU 는 시스템 전체 전력의 약70%를 소비하는 것으로 보고되며, 이는 GPU의 연산밀도가 다른 부품보다 훨씬 높기 때문이다.

메모리 접근과 데이터 이동 또한 전력 소모의 주요 요인이다. GPU 내부의 메모리 전송뿐 아니라 GPU 와 외부 메모리 간의 I/O 과정에서도 많은 에너지가 사용된다. 거대 Transformer 모델의 경우, 대규모 파라미터를 지속적으로 불러오고 키-값 캐시(KV cache)를 유지하는 데 상당한 전력이 필요하다. 2025 년 발표된 연구에서는 GPT-3 와 같은 초거대 언어모델이 전체 에너지의 대부분을 메모리 접근에 소비한다고 보고했으며, 반면 ResNet 과 같은 CNN 모델은 연산량이 많을수록 연산(MAC) 에너지 비중이 증가한다고 밝혔다. 이는 Transformer 모델이 메모리 바운드(memory-bound), CNN 모델이 연산 바운드(compute-bound)로 구분된다는 점을 보여준다. 연산 정밀도 또한 GPU 전력 효율에 큰 영향을 미친다. FP32(32 비트 부동소수점) 연산에 비해 FP16 또는 BF16 과 같은 저정밀 연산은 한 번에 더 많은 데이터를 병렬 처리할 수 있고, 전용 Tensor 코어를 활용하여 연산 속도를 높이면서 전력 소모를 줄일 수 있다. 실제로 45nm 공정에서 FP16 연산 한 번은 약 1.5pJ의 에너지를 소비하지만 INT8 연산은 0.23pJ로 6 배 이상 효율적이다. 정밀도가 높을수록 더 많은 트랜지스터 전환과 데이터 이동이 필요하므로 전력 소모가 커지며, 혼합 정밀도(Mixed Precision) 학습을 적용하면 주요 연산을 FP16 으로 수행하고 일부민감한 부분만 FP32 로 유지하여 전체 연산 시간을 줄이고 에너지 효율을 향상시킬 수 있다.

GPU 의 아키텍처와 클럭 속도도 전력 소비를 결정짓는 요소이다. 클럭 주파수와 전압이 높을수록 연산 속도는 증가하지만 전력 소모와 발열은 비선형적으로 커진다. 예를 들어 NVIDIA V100 은 "Maximum Efficiency Mode"를 통해 전압과 클럭을 낮춰 절반의 전력으로 약 80%의 성능을 낼 수 있었다. 이는 최고 성능이 필요하지 않은 상황에서는 낮은 전력으로도 충분한 효율을 달성할 수 있음을 보여준다. 현대 GPU 들은 P-State 기반의 전력 관리기능을 통해 부하에 따라 전력을 조절하지만, 딥러닝학습 시에는 대부분 최대 성능으로 작동하여 높은 전력 상태를 유지하는 경우가 많다.

종합하면 GPU 전력 소모는 연산량, 메모리 접근 빈도, 연산 정밀도, 클릭 설정의 상호작용에 따라 결정된다. 계산량이 많고 데이터 이동이 빈번하며, 높은 정밀도와 최대 클럭으로 동작할수록 전력 소모는 급격히 증가한다. 또한 GPU 활용률(Utilization)이 낮을수록 전력 대비 성능 효율이 떨어지기 때문에, 하드웨어 자원을 얼마나 효율적으로 활용하느냐가 에너지 최적화의 핵심 요소라고 할 수 있다.



- 모델 규모 및 학습조건과 전력 소모의 상관관계

모델의 파라미터 수, 연산 복잡도(FLOPs), 학습시간(에폭수), 배치 크기(batch size) 등 학습 조건은 GPU 의 총 에너지 소비량을 결정짓는 핵심 요인이다. 일반적으로 모델이 클수록, 학습 데이터가 많을수록, 배치가 클수록 에너지 소모가 증가하지만, 그 관계는 단순한 선형이 아니라 비선형적으로 변화한다.

먼저 모델 파라미터 수와 FLOPs 는 연산량과 메모리 접근량을 직접 결정한다. 모델 크기가 커지면 연산량이 기하급수적으로 증가하며, 매미니배치마다 더 많은 가중치 데이터를 불러와야한다. 예를 들어 GPT-3 175B와 같은 초거대 모델은 작은 모델보다 수백 배 이상의 연산과 에너지를 필요로 한다. 실제로 175 억 파라미터 모델과 7 억 파라미터 모델을 비교했을 때 파라미터 수는 약 25 배차이지만, 연산량은 약 128 배에 달한다. 이는 모델이 커질수록 연산 효율이 떨어지고, 분산 학습 시 통신 오버헤드까지 발생해 에너지 소비가 비선형적으로 증가함을 의미한다. 따라서 거대 모델일수록 학습효율 개선 전략을 병행하지 않으면 전력 부담이급격히 커질 수 있다.

학습 시간 역시 GPU 에너지 소비에 직접적인 영향을 미친다. 학습 반복 횟수(에폭 수)가 늘어날수록 총 에너지는 증가하지만, 학습률 조정, 모델 구조, 배치 크기 등의 설정에 따라 목표 성능에 도달하는 시간은 크게 달라질 수 있다. 잘 설계된 러닝레이트 스케줄이나 적절한 하이퍼파라미터 설정은 빠른 수렴을 유도하여 불필요한 연산을 줄이고 에너지를 절약할 수 있다. 반대로 비효율적인 설정은 학습 시간을 불필요하게 늘려 전력을 낭비하게 된다. Early Stopping(조기 종료)은 검증 오차가 증가하기 시작하는 시점에서 학습을 중단함으로써 과적합을 방지하고 연산 낭비를 줄이는 효과적인 방법이다. MIT Lincoln Laboratory 의 한 연구에서는 GPU 전력을 150W 로 제한하고 BERT 모델을 학습했을 때, 학습 시간이 단 2 시간 늘어나는 대신 전체 에너지 소비가 한 가정의 일주일치 전력 사용량에 해당하는 수준으로 감소했다고 보고하였다. 이처럼 약간의 학습 시간 증가로 큰 에너지 절감을 얻는 것은 효율적 선택이라 할 수 있다.

배치 크기 또한 GPU 자원 활용도와 전력 효율을 결정하는 주요 변수이다. 작은 배치를 사용하면 각스텝의 연산량은 줄지만 전체 스텝 수가 늘어나고, 큰배치를 사용하면 스텝 수는 줄지만 한 번에 처리할 연산이 많아진다. 일정 수준까지 배치를 키우면 GPU의 연산 장치를 더 효율적으로 활용할 수 있어

에너지 효율이 향상되지만, 너무 큰 배치는 학습 안정성을 떨어뜨리고 메모리 및 통신 병목을 유발해 효율이 다시 저하된다. NSDI 2023 의 Zeus 연구에서는 배치 크기를 최적화하는 것만으로도 기본 설정 대비 3.4%에서 최대 65%까지 에너지 절감을 달성할 수 있다고 보고하였다. 즉, 단순히 배치 크기조정만으로도 GPU 학습 전체 에너지의 절반 이상을 절약할 잠재력이 있다는 것이다.

결국 모델의 규모, 연산 복잡도, 학습 시간, 배치 크기 등은 모두 GPU의 전력 소비를 결정짓는 중요한 요인으로, 모델이 커지고 학습이 길어질수록 에너지 소모는 급격히 증가한다. 그러나 효율적인 하이퍼파라미터 설정, 조기 종료, 적절한 배치 크기조정 등의 전략을 적용하면 이 증가율을 완화할 수 있다. 에너지 효율적인 딥러닝을 위해서는 모델 설계 단계부터 연산 효율과 자원 활용을 고려해, 최소한의 전력으로 최대 성능을 달성하는 방향으로 접근하는 것이 필요하다.

- 전력 최적화 전략

GPU 를 활용한 딥러닝 학습의 전력 효율을 개선하기 위해 하드웨어적인 접근과 소프트웨어적인 접근 모두 중요합니다. 하드웨어 측면에서는 더효율적인 GPU 아키텍처와 시스템 환경(냉각 등)의선택이 핵심이며, 소프트웨어 측면에서는 연산최적화 기법과 알고리즘적 개선으로 같은 작업을 덜쓰는 에너지로 해내는 것이 목표입니다. 아래에서는 하드웨어, 소프트웨어, 그리고 학습 과정 자체의최적화 기법으로 나누어 전략들을 정리합니다.

하드웨어 측면 최적화 전략:

고효율 GPU 선택은 딥러닝 학습에서 에너지 효율을 높이기위한 가장 직접적인 방법 중 하나이다. 최신 GPU 일수록 성능 대비 전력 효율이 향상되는 경향이 뚜렷하게 나타난다. 예를 들어 NVIDIA H100은 이전 세대인 A100 대비 와트당 성능이 2배이상 향상되어, 초기 비용은 높더라도 동일한 작업을 더 적은 전력으로 수행할 수 있다. 따라서 예산이 허용된다면 전력 효율이 높은 GPU 로업그레이드하는 것이 장기적인 관점에서 유리하다. 실제로 NVIDIA 는 Ampere(A100)에서 Hopper(H100)로 세대를 교체하면서 성능뿐 아니라에너지 효율 향상을 핵심 개선점으로 강조하였다. 또한 구형 GPU 여러 대를 동시에 사용하는 것보다



최신 GPU 소수로 동일한 작업을 처리하는 편이 전력 사용량과 유지 관리 비용 측면에서 더 효율적이다.

전력 제한(Power Capping)의 활용도 GPU 전력 최적화에서 중요한 전략이다. NVIDIA 의 nvidia-smi 명령어를 이용하면 GPU 가 사용할 수 있는 최대 전력을 제한할 수 있는데, 이를 통해 약간의 성능 저하를 감수하면서도 전체 에너지 소비를 줄일 수 있다. MIT Lincoln Laboratory 의 연구에 따르면 GPU 의 전력 상한을 150W 로 제한할 경우 약12~15%의 에너지를 절감하면서도 작업 시간 증가는 3%에 불과하였다. 이러한 트레이드오프는 약간의속도 손실을 허용하는 대신 총 전력 소비를 줄이는 효과를 가져온다. 데이터센터 환경에서는 워크로드특성에 맞추어 전력 캡을 자동 조정함으로써 피크 전력 부하와 냉각 부담을 완화하고, Slurm 과 같은스케줄러와 연계하여 효율적인 에너지 관리가가능하다.

효율적인 냉각 및 전력 인프라 구축 또한 간접적인 전력 최적화에 기여한다. GPU 는 높은 발열로 인해 냉각에 많은 전력을 소모하므로, 수냉식 쿨링이나 AI 기반 냉각 제어를 도입하면 냉각 전력 사용량을 크게 줄일 수 있다. 실제로 DeepMind 는 구글 데이터센터의 냉각 시스템에 AI 제어 기술을 적용하여 냉각 에너지를 40% 절감하는 성과를 거두었다. 또한 서버실 온도와 공조 환경을 최적화하거나, 냉각 효율이 높은 시간대(야간이나 겨울철)에 집중적으로 학습 작업을 스케줄링하는 방식도 전력 절감에 효과적이다.

Hopper 와 Ampere 아키텍처의 비교에서도 최신세대의 전력 효율성이 두드러진다. NVIDIA Hopper(H100)는 Ampere(A100) 대비 Transformer Engine(FP8 정밀도 지원), 더 큰 L2 캐시, HBM3 메모리 등을 통해 동일한 연산을 더 적은 데이터 이동으로 수행하도록 설계되었다. FP8 정밀도는 Transformer 모델 학습 시 일부 연산을 낮은 정밀도로 수행하여 연산량과 메모리 대역폭 요구를 줄여 에너지를 절감한다. 또한 NVLink 4.0 과 NVSwitch 개선으로 GPU 간 통신 효율이 향상되었고, 2세대 MIG(Multi-Instance GPU) 기능을 통해 하나의 GPU 를 여러 작업으로 나누어 활용할 수 있어 유휴 자원을 최소화할 수 있다. 이러한 최신 아키텍처의 기능을 적극적으로 활용하면 하드웨어 차원에서의 에너지 효율을 극대화할 수 있다.

소프트웨어 측면 최적화 전략:

혼합 정밀도(Mixed Precision) 학습은 현재 딥러닝 훈련의 표준으로 자리 잡은 핵심 기술이다. FP16(BF16)과 FP32 연산을 혼용함으로써 메모리 사용량을 줄이고 Tensor 코어를 활용할 수 있어 학습 단축되고 에너지 효율이 향상된다. NVIDIA 의 기술 보고서에 따르면 혼합 정밀도 사용 시 최대 2 배의 학습 속도 향상이 가능하며, 이는 동일한 시간 동안 더 적은 전력으로 작업을 완료할 수 있음을 의미한다. 또한 2022 년 USENIX 학회에서 발표된 Campo 연구에서는 연산 캐스팅 비용까지 최적화한 혼합 정밀도 학습을 통해 에너지 효율이 개선되었다고 보고하였다. 요약하면 21.4% FP32 만으로 학습하는 것은 전력 낭비에 가깝고, 가능한 한 FP16 또는 BF16 을 사용하되 필요한 연산에서만 FP32 를 활용하는 것이 가장 효율적인 전략이다.

그래디언트 체크포인팅(Gradient Checkpointing)은 일부 중간 활성값을 저장하지 않고 필요 시재계산함으로써 메모리 사용량을 줄이는 기법이다. 이를 통해 GPU 메모리 여유가 확보되어 더 큰 배치나모델을 한 번에 처리할 수 있으며, 그 결과 같은 작업을 더 적은 반복(iteration)으로 수행할 수 있다. 이방식은 GPU 자원 활용도를 높이고 전체 전력 소비를 줄이는 데 기여하지만, 재계산으로 인한 연산 증가로약간의 속도 저하가 발생할 수 있다. 따라서 메모리가부족한 대규모 모델 학습(GPT-3 등)에 특히 유용하며, GPU 수나 배치를 줄임으로써 결과적으로 총 에너지절감을 달성할 수 있다.

희소성 및 프루닝(Sparsity & Pruning)은 불필요한 연산을 제거하여 계산량과 메모리 접근을 줄이는 대표적인 최적화 기법이다. NVIDIA Ampere 아키텍처는 2:4 구조적 희소성을 지원하여 0 값이 많은 행렬을 효율적으로 처리하며, 이로 인해 동일한 연산을 절반의 시간과 전력으로 수행할 수 있다. 이는에너지 절감에 직접적인 효과를 가지며, 하드웨어가 구조적 희소성을 지원할수록 효과가 극대화된다. 예를 들어 Ampere 이상의 GPU에서는 sparse Tensor Core 를 활용한 희소 연산 최적화가 가능하다. 또한 대규모 언어모델에서는 MoE(Mixture of Experts) 구조를 통해 입력마다 일부 전문가 네트워크만 활성화함으로써 연산량과 전력 소모를 줄이는 연구도 활발히 진행되고 있다.

컴파일러 및 커널 최적화는 소프트웨어 측면에서 전력 효율을 높이는 또 다른 핵심 전략이다. NVIDIA 의 cuDNN, TensorRT, PyTorch XLA 등과 같은 저수준 최적화 라이브러리를 활용하면 메모리 접근



패턴과 연산 순서를 개선하여 불필요한 데이터 이동을 줄이고 GPU 유휴 시간을 최소화할 수 있다. 연산 fusion, 메모리 preload, 동적 연산 최적화 등의 기술은 전체적인 전력 소모를 낮추는 데 기여한다. 최근에는 AutoML 과 결합된 컴파일러 최적화 연구가 활발히 진행되어, 주어진 하드웨어에서 에너지 효율이 최대화되는 커널 실행 방식을 자동으로 탐색하는 기술도 개발되고 있다.

이와 같이 혼합 정밀도, 그래디언트 체크포인팅, 희소성, 그리고 커널 최적화는 각각 독립적으로도 에너지 효율 향상에 기여하지만, 함께 적용할 경우 상승 효과를 통해 GPU 학습 과정의 전력 소비를 크게 줄일 수 있다.

- 학습 과정 자체의 최적화 기법

하이퍼파라미터 최적화와 스케줄러 설정은 학습 효율과 전력 소비에 직접적인 영향을 미치는 핵심 요소이다. 학습률, 옵티마이저, 그리고 스케줄링 전략을 적절히 조율하면 필요한 에폭 수를 줄여 학습 시간을 단축시키고, 그만큼 에너지를 절약할 수 있다. 예를 들어 사이클릭 러닝레이트(one-cvcle)나 warmup 후 감속(cosine decay) 스케줄은 빠른 수렴을 유도하여 목표 정확도에 더 짧은 시간 내 도달하게 한다. 또한 과적합을 방지하기 위한 정규화(regularization)나 드롭아웃(dropout) 기법 역시 불필요한 연산을 줄여 전력 효율을 높이는 방법이다. 가접적으로 하이퍼파라미터 탐색 과정에서는 모든 조합을 전부 실험하는 대신, 피델리티(fidelity) 전략을 적용해 유망한 설정만 정밀 검증함으로써 실험 횟수와 에너지를 줄일 수 있다. MIT LLSC 의 연구에 따르면 학습 초반의 성능 곡선을 분석해 비효율적인 실험을 조기에 중단한 결과, 최대 80%의 에너지를 절감할 수 있었다고 보고하였다. 이는 하이퍼파라미터 최적화나 AutoML 에서 에너지 인지형 전략의 중요성을 보여준다.

조기 종료(Early Stopping)는 학습 중 검증 오차가 증가하기 시작하는 시점에서 학습을 멈추는 기법으로, 불필요한 학습을 방지해 시간과 전력을 절약하는 가장 직접적인 방법 중 하나이다. 이 방법을 적용하면 전체 학습 연산량의 5~30% 이상을 줄일 수 있으며, 특히 대형 모델의 경우 에폭 하나당 에너지소비량이 매우 크기 때문에 몇 번의 에폭만 줄여도 절대적인 절감 효과가 크다.

분산 학습과 병렬화 최적화 또한 전력 효율에 큰 영향을 미친다. 여러 GPU를 사용할 때 통신 병목을 줄이고 부하를 균등하게 분산하면 자원을 효율적으로 활용할 수 있다. 예를 들어 Gradient Accumulation 을 통해 통신 빈도를 줄이거나, 데이터 병렬과모델 병렬을 혼합해 각 GPU 의 메모리 및 연산 자원을 최대한 활용하는 것이 좋다. 최근에는 유휴 GPU 를 자동으로 껐다가 필요할 때 다시 가동하는 탄력적 분산 학습 방식도 연구되고 있으며, 이는 클러스터 단위의 에너지 절약에 효과적이다. 또한 NVLink 를 통한 동일 섀시 내 고속 통신이나 그래디언트 압축(gradient compression) 기법을 적용하면 노드 간 통신으로 인한 추가 전력 소모를 줄일 수 있다.

결국 하드웨어, 소프트웨어, 학습기법을 종합적으로 최적화하는 것이 중요하다. "더 적은 자원으로 같은 모델을 학습시키는 것"이 목표이며, 개별기법보다 여러 전략을 병행할 때 가장 큰 효과를 얻을 수 있다. 예를 들어 H100 GPU 에서 혼합 정밀도학습, 전력 제한(Power Capping), 대형 배치, 조기종료를 동시에 적용한다면, 각 요소의 절감 효과가누적되어 상당한 수준의 에너지 효율 향상을 달성할수 있을 것이다.

- 실제 기업/연구기관의 전력 최적화 사례

이론적으로 제시된 전력 최적화 기법들은 실제 산업계와 연구 현장에서도 활발히 적용되고 있다. 대형 기술 기업들과 주요 연구기관들은 초거대 AI 모델을 운용하면서 전력 효율 문제에 직면하였고, 이를 해결하기 위해 다양한 혁신적 방안을 도입하고 있다.

먼저 DeepMind(구글)는 데이터센터 냉각 효율을 높이기위해자사의 강화학습기술을 HVAC 시스템에 적용하여 서버실 냉각 에너지를 약 30~40% 절감하였다. 이는 AI 를 데이터센터 운영에 역으로 활용한 대표적인 사례로, GPU 자체의 전력 소모를 줄인 것은 아니지만 부수적인 전력 소비를 대폭줄이는 효과를 거두었다. 또한 풍력 발전량 예측 등 전력 그리드 운영 측면에서도 AI를 활용하여 에너지효율화를 실현하고 있으며, 이러한 시도는 구글의 탄소중립 전략과도 맞물려 진행되고 있다.

OpenAI 는 GPT-3 와 GPT-4 의 학습을 위해 마이크로소프트와 협력하여 수만 개의 GPU 로 구성된 초대형 HPC(고성능 컴퓨팅) 인프라를 구축하였다. 이 시스템에는 최적화된 전력 분배 구조, 수냉식 냉각, 고효율 전원공급장치가 적용되어 있다.



또한 OpenAI 는 모델 효율화를 위해 희소 Transformer(Sparse Transformer), 모델 압축, 지식 증류(knowledge distillation) 등을 연구하고 있으며, 학습 과정에서는 FSDP(Fully Sharded Data Parallel) 기술을 통해 메모리 사용을 줄이고, 정밀도 자동조정으로 연산 효율을 향상시켰다. 비록 공식적인에너지 절감 수치를 발표하지는 않았으나, OpenAI는 GPT-3 학습 이후 탄소 배출 상쇄를 위한 조치를 취했으며, 다양한 크기의 모델을 제공함으로써성능과 효율 간의 균형을 고려하는 운영 방식을 채택하고 있다.

Meta AI 는 2023 년부터 AI 연구 과정에서의 에너지와 탄소 배출량을 투명하게 공개하기 시작하였다. 특히 LLaMA 모델 발표 당시 학습에 사용된 전력량과 탄소 배출량을 함께 보고하여 연구자들의 인식을 제고하였다. Meta 는 또한 "시스템적 에너지·탄소 발자국 보고 프레임워크"를 개발하여 머신러닝 실험의 에너지 소비를 실시간으로 추적하고 표준화된 형태로 보고할 수 있는 도구를 제시하였다. 더불어 Meta AI Lab 에서는 하드웨어 자원에 따라 작업 부하를 자동으로 최적화하는 알고리즘을 연구하여, GPU 와 CPU 의 분리함으로써 탄소 효율을 10~20% 향상시키는 성과를 거두었으며, 해당 기술은 2023년 HPDC 학회에서 최고 논문상을 수상하였다.

이외에도 Microsoft 는 전력이 저렴하거나 잉여 전력이 발생하는 시간대에 AI 학습을 집중시키는 '지능형 부하 이동(intelligent workload shifting)' 전략을 적용하여 효율적인 자원 운용을 실현하였다. NVIDIA 는 MLPerf 벤치마크를 통해 자사 하드웨어와 소프트웨어 스택을 최적화하여 와트당 최고 성능을 기록하였으며, Tesla 등 자율주행 기업은 훈련용 슈퍼컴퓨터(DOJO)와 ASIC 기반 하드웨어를 자체 개발하여 전력 효율을 높이고 있다. 학계 역시 "그린 AI(Green AI)" 흐름 속에서 모델의 정확도뿐 아니라 에너지 효율과 탄소 배출 저감 기여도를 함께 평가하는 추세로 전환되고 있다.

이처럼 산업계와 학계 전반에서 진행되는 이러한 노력들은 AI의 지속 가능성을 확보하고 운영 비용을 절감하는 동시에, 환경적 책임을 실천하기 위한 방향으로 발전하고 있다.

3. 결론

- 연구 결과 요약

본 탐구는 "AI 모델 학습 시 GPU 소비 전력과 최적화 방안"을 주제로 GPU 전력 소모의 현황과 개선 전략을 심층적으로 분석한 연구이다. 연구 결과, GPU 전력 소모는 AI 시대의 핵심 과제로, GPT-3 와 같은 초거대 모델의 학습 과정에서 막대한 에너지와 탄소 배출이 발생함이 확인되었다. 하드웨어 측면에서는 최신 GPU 로의 업그레이드, 전력 관리, 냉각 효율 향상이 에너지 절감의 핵심 요소로 나타났으며, 특히 H100 과 같은 신형 GPU는 이전 세대보다 에너지당 성능이 크게 향상된 것으로 분석되었다. 소프트웨어 및 알고리즘 측면에서는 혼합 정밀도, 희소화, 그래디언트 체크포인팅, 조기 종료 등의 기법을 적용함으로써 연산량과 메모리 사용을 줄여 전력 효율을 높일 수 있었다. 또한 산업계에서는 이미 이러한 최적화 노력을 실무에 반영하고 있으며, 에너지 효율이 AI 연구와 개발의 새로운 평가 기준으로 부상하고 있다. 결론적으로, 성능 향상과 에너지 효율은 상호 보완적으로 추구되어야 하며, "친환경 AI"는 선택이 아닌 필수적인 발전 방향임이 확인되었다. 본 연구는 이러한 변화 속에서 지속 가능한 AI 기술 발전을 위한 기초 자료로서 의의가 있다.

- 한계 및 향후 연구

본 조사 연구는 최신 자료들을 망라하여 전반적인 흐름과 전략을 제시하였으나, 정량적 실험 검증이 미흡한 한계가 있다. 예를 들어 각 최적화 기법의 상호 작용 효과(시너지)를 실측으로 확인하지는 못하였다. 향후에는 혼합정밀도, 전력 제한(Power Capping), 대형 배치를 동시에 적용했을 때의 에너지 효율 변화를 정량적으로 측정하는 연구가 이루어질 필요가 있다. 또한 TPU, IPU, ASIC 등 GPU 대안 하드웨어와의 비교나, 클라우드 AI 서비스의 멀티테넌시 환경에서의 에너지 최적화 방안 역시 후속 연구로서 가치가 크다. 그럼에도 본 연구는 현재 시점에서 알려진 다양한 방법들을 통합적으로 고찰하였다는 점에서 의의가 있으며, 에너지 효율적 AI 개발에 관심을 가진 연구자와 엔지니어들에게 유용한 가이드가 될 것이다.

4. 참고문허

1. Kylie Foy, "AI models are devouring energy. Tools to reduce consumption are here, if data centers will adopt.", MIT Lincoln Laboratory News, Sep 2023

2.



3. Siddharth Samsi et al., "Understanding and Optimizing GPU Energy Consumption of DNN Training.", USENIX NSDI 2023 (Project Zeus)

4.

5. Ilpyung Yoon, Jihwan Mun, Kyeong-Sik Min, "Comparative Study on Energy Consumption of Neural Networks by Scaling of Weight-Memory Energy Versus Computing Energy.", Electronics, vol.14, no.13, 2718, 2025

6.

7. Dawson Lear, "NVIDIA H100 vs NVIDIA A100", AMAX Engineering Whitepaper, Jan 2024

8.

9. NVIDIA, *"Tesla V100 GPU Accelerator Datasheet"*, Mar 2018

10.

- 11. Xin He et al., "Campo: Cost-Aware Performance Optimization for Mixed-Precision Training", USENIX ATC 2022
- 12. Joh min soeng(professor of), Research of Babo munghungi; damn AI