

"We overcome the multicore/manycore programming wall by means of compiler, architecture, and OS techniques"



Multicore Computing Research Laboratory

멀티코어 컴퓨팅 연구실

서울대학교 301동 515호
Phone: 02-880-1837
Email: chief@aces.snu.ac.kr
Web: http://aces.snu.ac.kr/

고성능 컴퓨터 시스템



슈퍼컴퓨터 '천동'

고성능 컴퓨터 시스템의 구축은 병렬 프로그래밍 모델 및 소프트웨어의 성능을 평가하는데 있어 매우 중요합니다. 멀티코어 컴퓨팅 연구실에서는 병렬 프로그래밍 모델 SnuCL의 성능평가를 위해 2012년에 이종 슈퍼컴퓨터 '천동'을 설계하고 구축하였습니다. 천동은 자체 수냉 시스템을 탑재하였으며 낮은 가격 및 저전력으로 높은 성능을 달성하는데 초점이 맞춰져 있습니다. 현재 천동 1.5가 대중에게 공개되어 있고, 변화하는 하드웨어 요구사항에 맞추어 끊임없이 개선되고 있습니다.

기록

TOP500 277위 등재 (2012년 11월)

Green500 32위 등재 (2012년 11월)

단 7억원의 시스템 구축 비용 (1US달러 당 159 MFLOPS)

TOP500에서 7번째로 전력 효율이 높음 (2012년 11월)

TOP500의 412개 클러스터 중 노드당 성능이 가장 높음 (2012년 11월)



하드웨어 사양

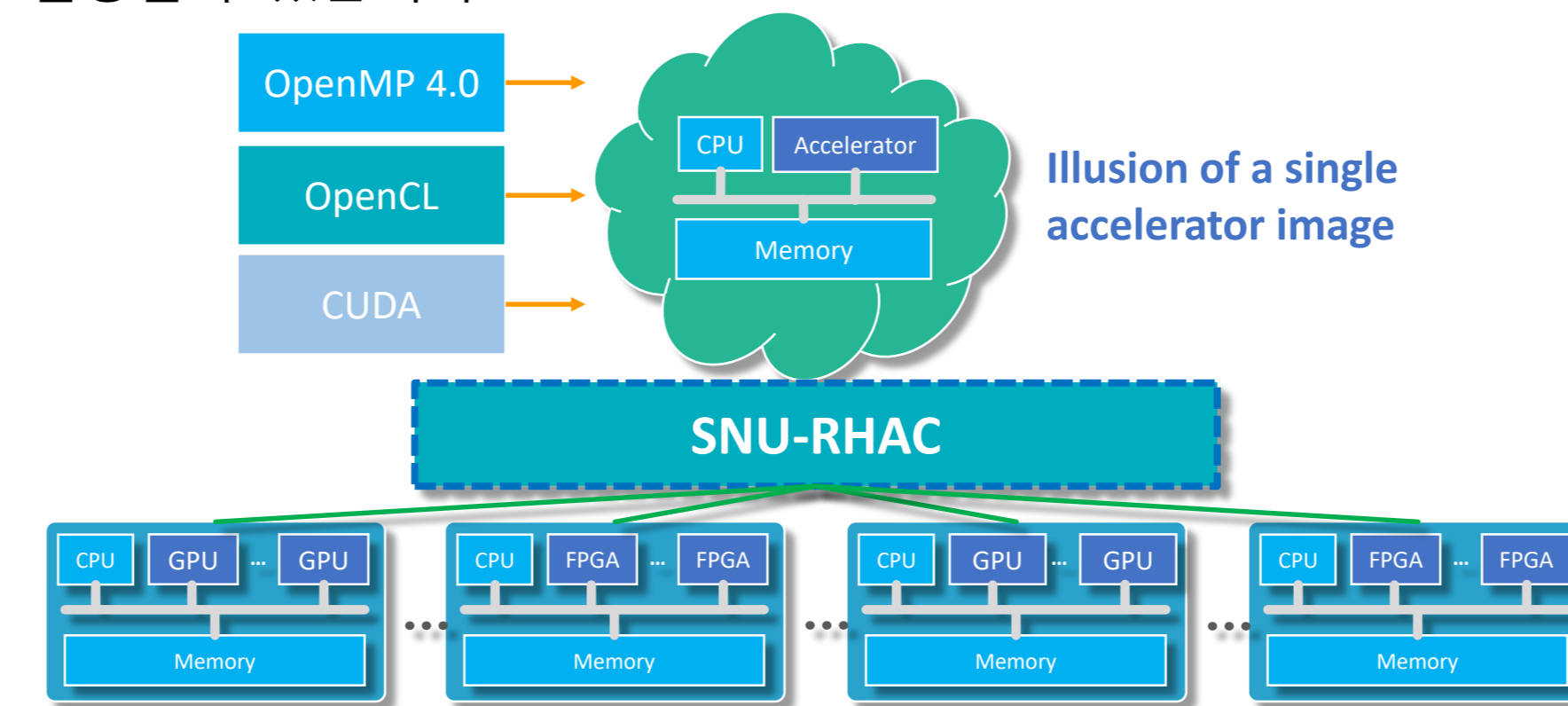
52 nodes ×
 2× Intel Xeon E5-2650 CPUs (8-core)
 4× AMD Radeon HD 7970 GPUs
 or AMD Radeon R9 Nano GPUs
 or AMD Radeon R9 290X GPUs
 or Nvidia GeForce GTX 1080 GPUs
 128 GB main memory
 +
 InfiniBand 4x QDR
 interconnection network
 144 TB storage space
 Water cooling system

천동은 2013년 4월부터 국내 사용자에게 연구 및 교육 목적으로 사용할 수 있도록 서비스를 제공하고 있습니다. <http://chundoong.snu.ac.kr/>. 지금까지 330명 이상의 사용자가 천동을 이용하였습니다.

범용 프로그래밍 모델

SNU Runtime for Heterogenous Accelerators in a Cluster (SNU-RHAC)

SNU-RHAC는 연구실에서 개발중인 이종 클러스터를 위한 범용 프로그래밍 모델입니다. SNU-RHAC를 통해 사용자는 여러 노드에 장착된 다양한 종류의 가속기를 마치 하나의 가속기인 것처럼 사용할 수 있습니다. 실제로 어떤 가속기가 사용되는지에 관계 없이 사용자는 OpenMP, OpenCL, CUDA 등의 원하는 프레임워크로 개발할 수 있습니다. 또한, MPI-OpenCL처럼 여러 프로그래밍 모델을 복잡하게 혼합하여 사용하는 불편함 없이 클러스터의 모든 자원을 활용할 수 있습니다.



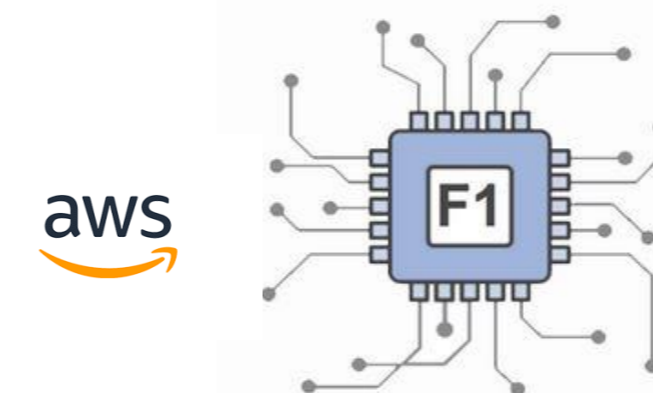
SNU-RHAC의 전신인 SnuCL은 오픈소스 소프트웨어입니다. <http://snucl.snu.ac.kr/>. SnuCL은 top-tier 컴퓨터 학회에서 10개의 논문과 9번의 튜토리얼을 통해 소개되었습니다.

FPGA를 위한 OpenCL/CUDA 프레임워크

FPGA는 뛰어난 에너지 효율성으로 인해 고성능 컴퓨팅을 위한 차세대 가속기로서 주목받고 있습니다. 하지만 FPGA 프로그래밍의 어려움이 FPGA가 범용 계산에 사용되는 데에 큰 걸림돌이 되고 있습니다. 멀티코어 컴퓨팅 연구실에서는 FPGA를 위한 OpenCL/CUDA 프레임워크를 개발하고 있습니다. 이 프레임워크의 목표는 컴파일러 최적화 기술과 유연한 메모리 서브시스템을 활용하여, Verilog 또는 VHDL로 쓰여진 프로그램과 비슷한 수준의 성능을 달성하는 것입니다.



Microsoft는 검색 엔진, 딥 러닝 등 다양한 워크로드를 가속하기 위하여 자사 데이터센터에 FPGA를 탑재하였습니다.

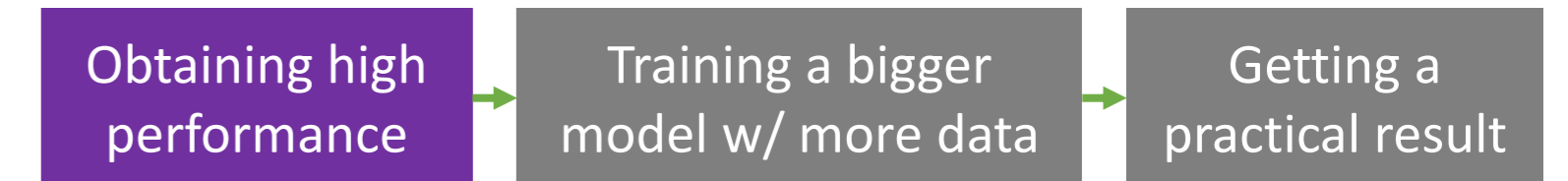


Amazon은 FPGA가 탑재된 EC2 F1 인스턴스를 서비스하고 있습니다.

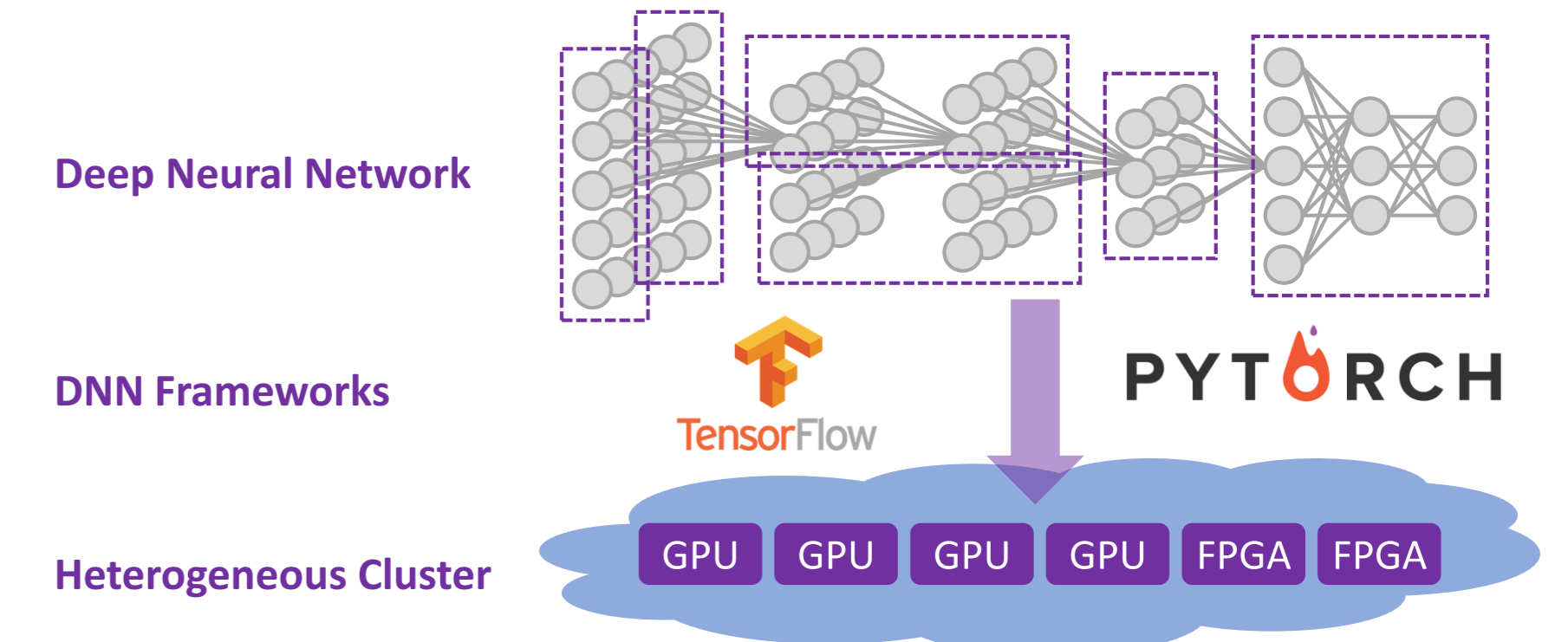
딥 러닝 프로그래밍 모델

클러스터에서의 딥 러닝

딥 러닝의 성공은 대량의 데이터와 컴퓨팅 파워의 향상을 통해 이루어졌습니다. 특히, GPU의 대규모 병렬 처리 능력은 딥 뉴럴 네트워크 모델이 합리적인 시간 안에 훈련되는 것을 가능하게 만들었습니다.



그러나 널리 쓰이는 딥 러닝 프레임워크들은 여전히 멀티 GPU 또는 멀티 노드 환경에서 성능 확장성을 가지지 못합니다. 멀티코어 컴퓨팅 연구실에서는 딥 러닝 프로그램의 계산 및 메모리 접근 패턴을 파악하여 대규모 클러스터에서 확장성 있는 훈련 성능을 달성하기 위한 연구를 진행하고 있습니다. 또한, 해당 기법을 널리 쓰이는 딥 러닝 프레임워크에 적용하면 기존에 사용하던 딥 러닝 프로그램의 수정 없이 고성능을 달성할 수 있습니다.



FPGA를 이용한 딥 러닝

최근 높은 에너지 효율을 달성하기 위해 FPGA를 이용한 딥 러닝이 주목받고 있습니다. 멀티코어 컴퓨팅 연구실에서는 사용자가 딥 뉴럴 네트워크의 명세를 입력하면 FPGA에서 실행할 수 있도록 자동으로 코드를 생성하는 딥 러닝 컴파일러를 개발하고 있습니다.

또한, 많은 수의 FPGA 딥 러닝 가속기가 고정소수점 데이터의 추론만 지원하는데 반해, 멀티코어 컴퓨팅 연구실에서는 부동소수점 및 추론/학습을 모두 지원하는 가속기를 구현한 바 있습니다. 이 가속기는 강화학습의 한 종류인 Asynchronous Advantage Actor-Critic(A3C)에서 GPU보다도 뛰어난 성능과 에너지 효율을 보여주었습니다.