

대규모 병렬 컴퓨팅

Massively Parallel Computing with CUDA

biztripcru@gmail.com

© 2021. biztripcru@gmail.com. All rights reserved.

CUDA 소개

Introduction to CUDA

본 동영상과, 본 동영상 촬영에 사용된 발표 자료는 저작권법의 보호를 받습니다.

본 동영상과 발표 자료는 공개/공유/복제/편집/상업적 이용 등, **개인 수강 이외의 다른 목적 사용을 금지합니다.**

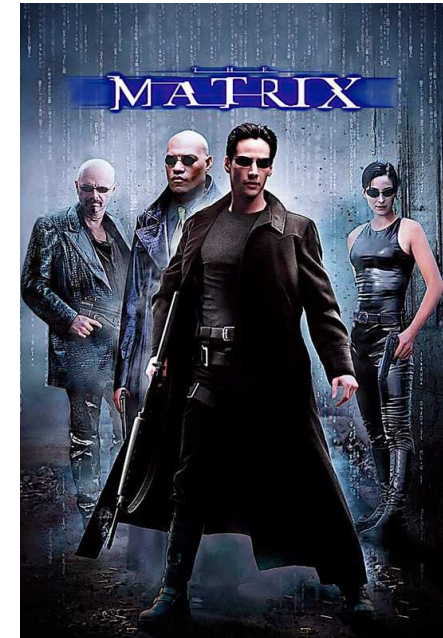
© 2021. biztripcru@gmail.com. All rights reserved.

내용 contents

- CUDA의 개발 배경 – 컴퓨터 그래픽스
- CUDA의 구성
- CUDA의 성능 – 슈퍼 컴퓨터 동향
- CUDA ^{쿠다} = Compute Unified Device Architecture, from NVIDIA 엔비디아

컴퓨터 그래픽스 computer graphics

- 현실과 똑같은 이미지 image
 - 복잡한 모델 model
 - 물리학, 광학 법칙 적용
 - 부드러운 동작 motion 생성
- 게다가, (매우) 빠르게
 - 실시간 처리 목표 → 영화 "매트릭스 Matrix" 시리즈
- 결국, 괴물 하드웨어 monster hardware 사용
 - 괴물 하드웨어를 계산용으로도 사용하자!
 - 대규모 병렬 컴퓨팅 MPC 으로 발전

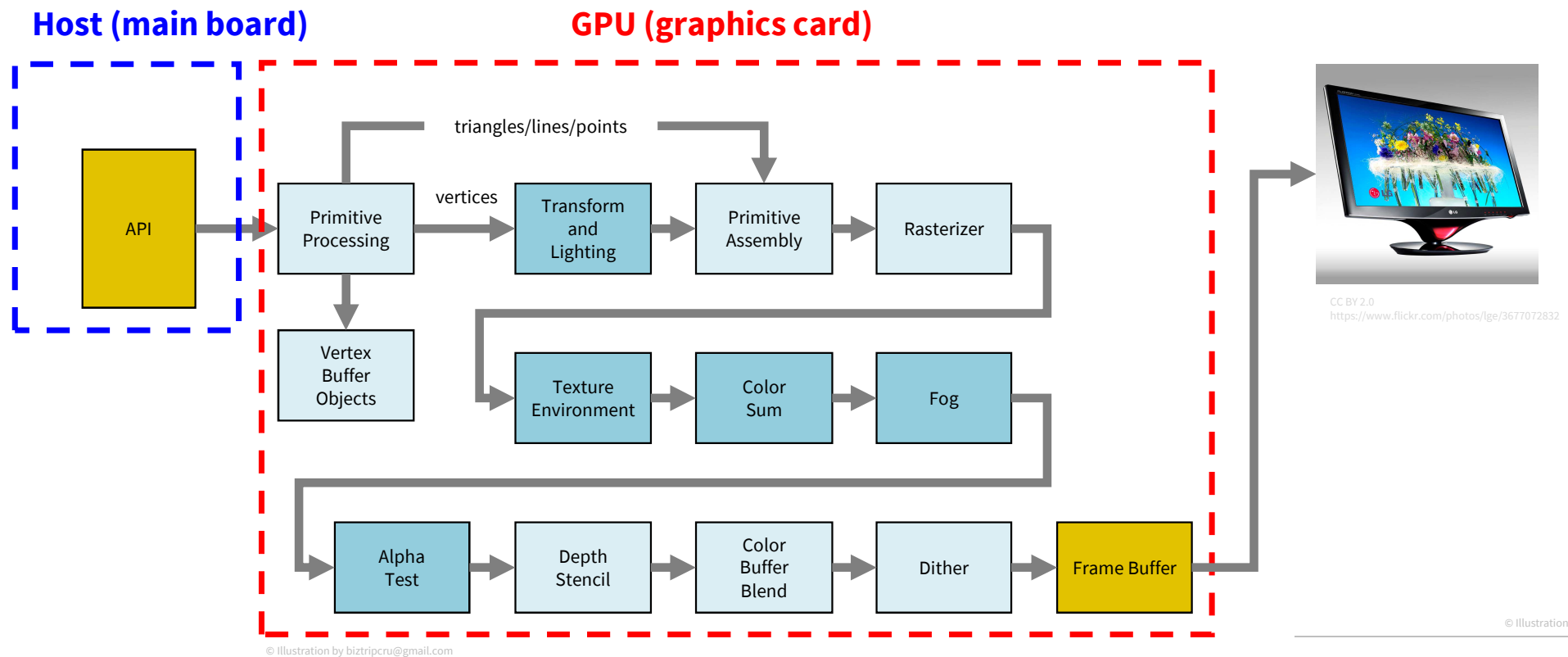


public domain
<https://www.flickr.com/photos/stevetroughton/17072638696>

영화 “매트릭스” 포스터
가상현실의 완벽한 구현

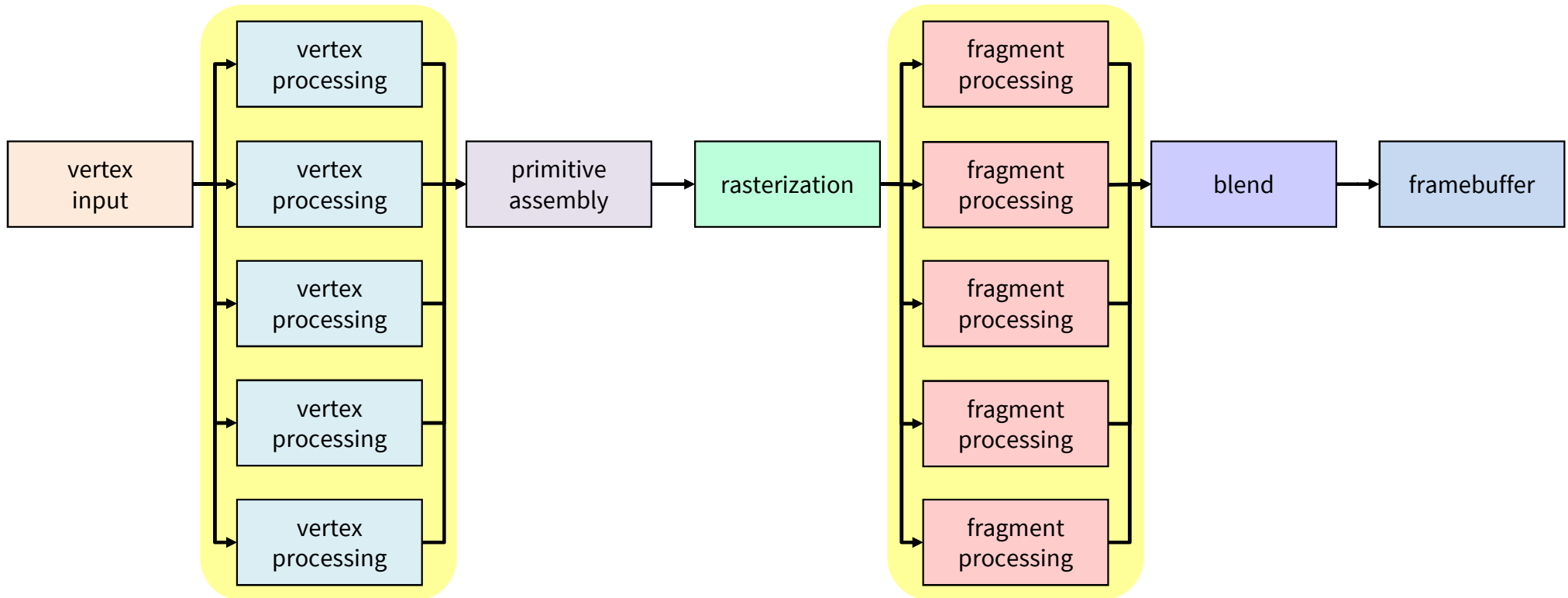
3D 그래픽스 파이프라인 graphics pipeline

- 매우 많은 데이터를 단계적으로 처리



병렬 처리를 도입

- 중요 단계마다 병렬 처리 parallel processing 로 가속 acceleration

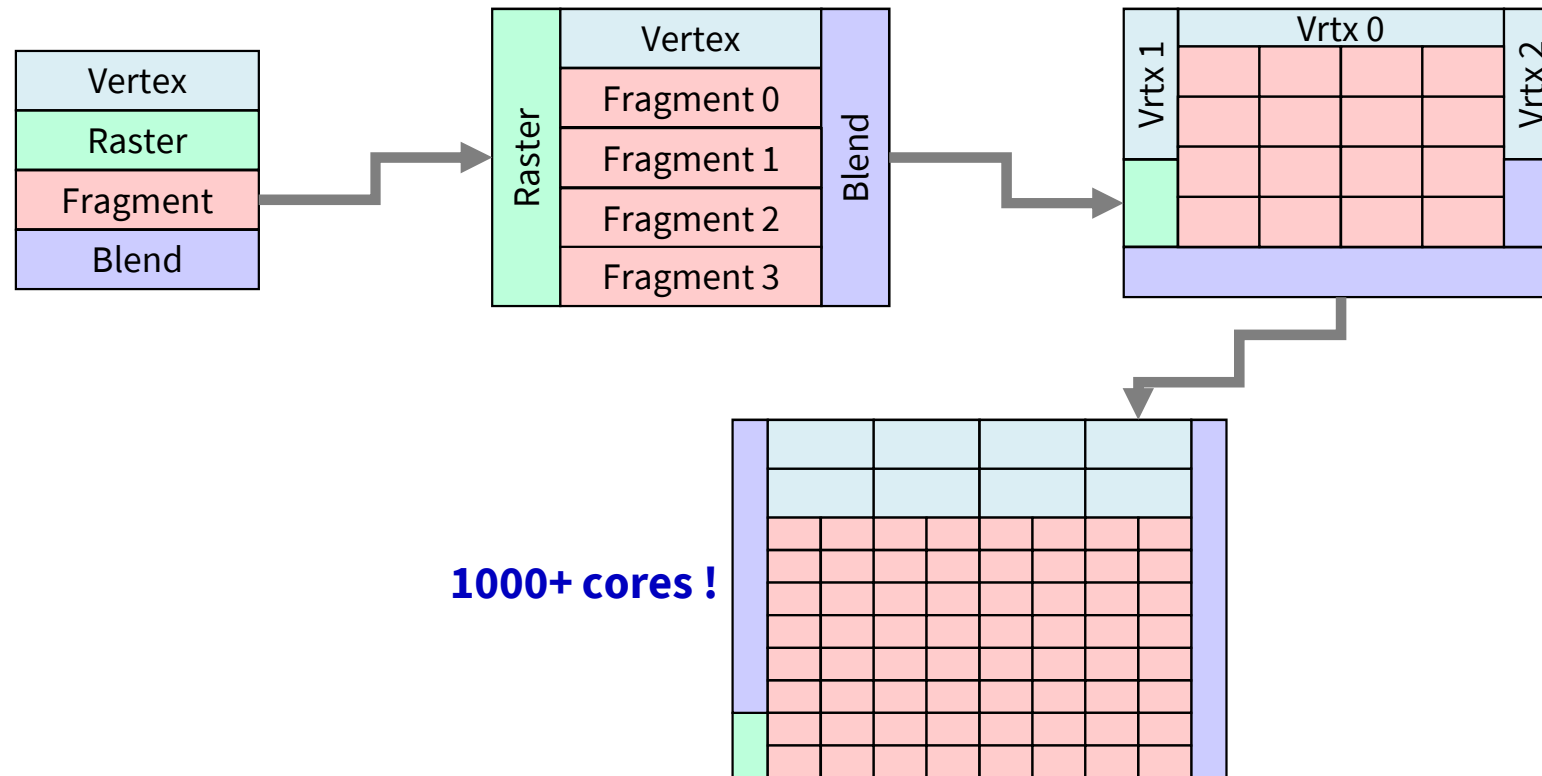


© Illustration by biztripcru@gmail.com

© Illustration by biztripcru@gmail.com

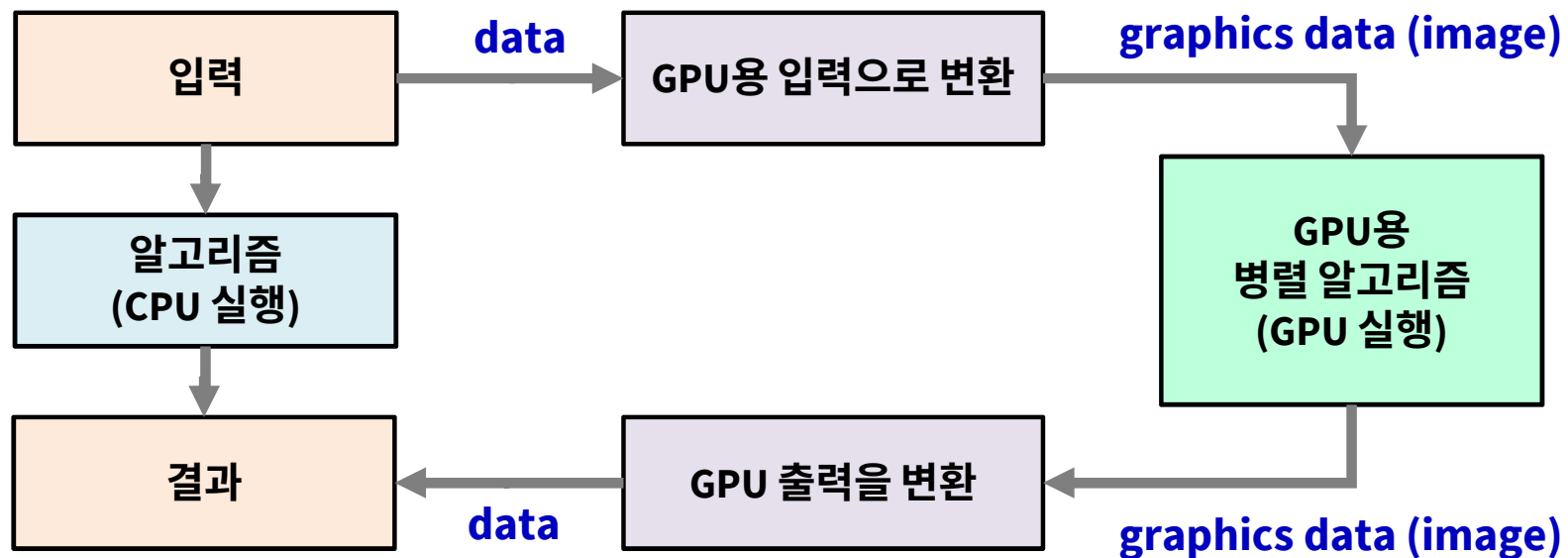
반도체 기술의 발전

- 무어의 법칙에 따라 더 많은 트랜지스터 → 더 많은 병렬처리



GPGPU의 도입

- **GPGPU: general purpose graphics processing unit**
 - GPU의 괴물 하드웨어를 계산용으로 사용하는 테크닉



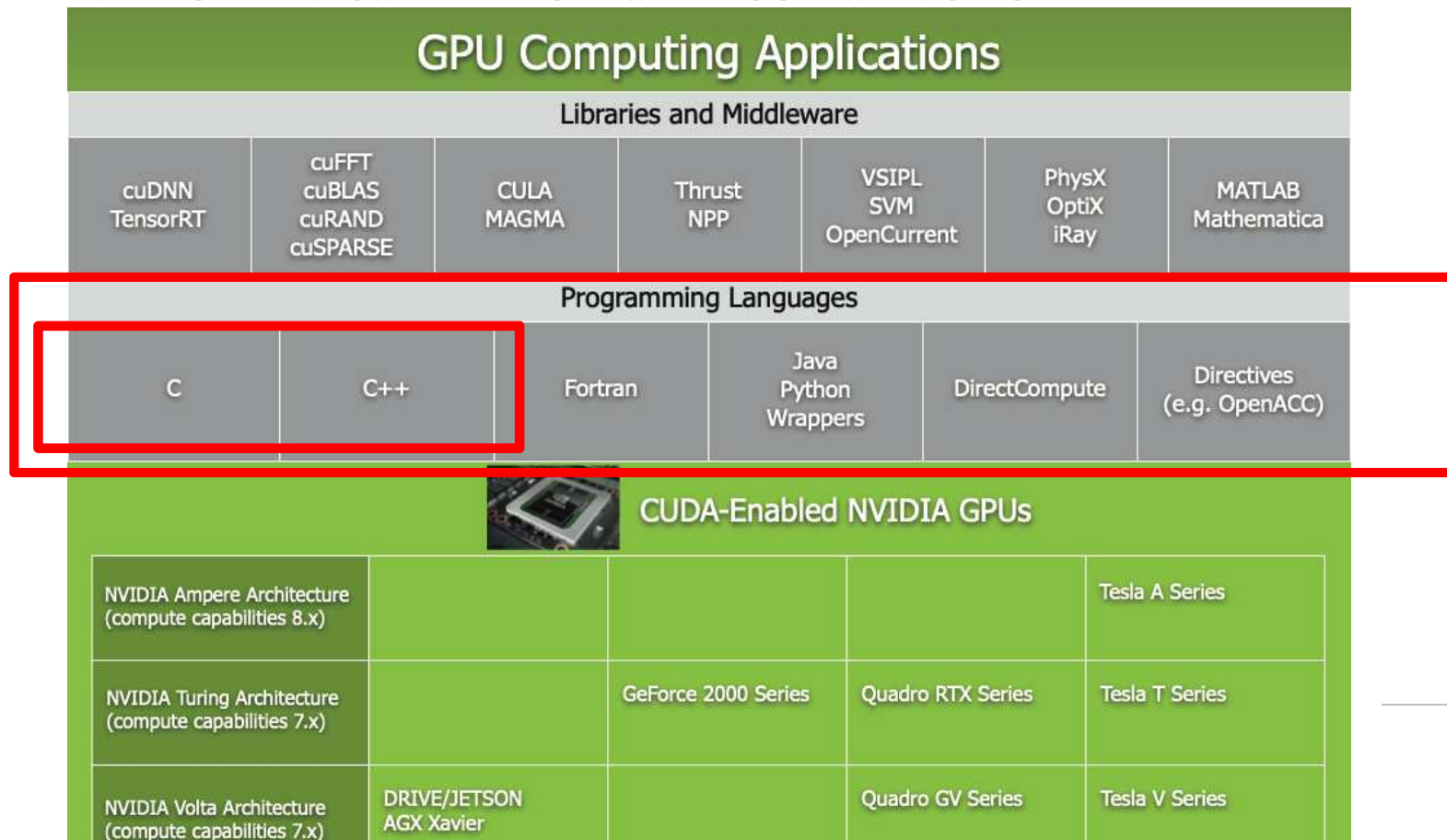
CUDA 쿠다

- **Compute Unified Device Architecture**
 - 2006년, NVIDIA 전용으로 출시
- **범용 general purpose 프로그래밍 모델**
 - GPU에서 대규모 쓰레드 thread 를 실행
 - GPU = 대규모 병렬처리 코프로세서 **massively data parallel co-processor**
 - **모델 = 디바이스 / 컴퓨터 구조 + 프로그래밍 언어 + 컴파일러 + much more !**
- **GPU 를 범용으로 사용하는 툴킷 toolkit 으로 구성**
 - CUDA 드라이버 → GPU 구동
 - CUDA 라이브러리 → API 함수들
 - GPU 기능을 직접 제어 가능 → 최고 효율 획득

CUDA의 구성

- 다양한 언어와 라이브러리를 제공 → 가장 기본은 **C/C++ API**

Copied from "Figure 2. GPU Computing Applications. CUDA is designed to support various languages and application programming interfaces." in CUDA Toolkit Documentation, v11.4, NVIDIA



슈퍼 컴퓨터 super computer

- 가장 빠른 컴퓨터들

- 초당 1조 (10^{12}) 번 이상 계산
- 초당 1,000조 (10^{15}) 번 → **Peta 급**



public domain
https://ko.wikipedia.org/wiki/%ED%94%8C%EB%A0%88%EC%9D%B4%EC%95%84%EB%8D%B0%EC%8A%A4_%EC%84%B1%EB%8B%A8#/media/%ED%8C%EC%9D%BC:Pleiades_large.jpg

플레이아데스 성단 open star cluster



public domain
[https://en.wikipedia.org/wiki/Pleiades_\(supercomputer\)#/media/File:Pleiades_supercomputer.jpg](https://en.wikipedia.org/wiki/Pleiades_(supercomputer)#/media/File:Pleiades_supercomputer.jpg)

플레이아데스 Pleiades
- NASA가 보유한 Peta급 슈퍼컴퓨터

슈퍼 컴퓨터 super computer 계속

- 가장 빠른 컴퓨터들
 - 초당 1조 (10^{12}) 번 이상 계산
 - 초당 1,000조 (10^{15}) 번 → **Peta 급**
- 특징: 계산 속도를 높이기 위해서
 - CPU/GPU 1,000개 이상을 동시 사용
 - 액화 질소 냉각, 건물 1개 층 이상
- 용도: 엄청난 계산이 꼭 필요한 분야
 - 기상대 (일기 예보)
 - 과학 계산 (물리학, 분자생물학)
 - 시뮬레이션 (핵폭발, 태풍) 등등



public domain
[https://en.wikipedia.org/wiki/Pleiades_\(supercomputer\)#/media/File:Pleiades_supercomputer.jpg](https://en.wikipedia.org/wiki/Pleiades_(supercomputer)#/media/File:Pleiades_supercomputer.jpg)

플레이아데스 Pleiades
- NASA가 보유한 Peta급 슈퍼컴퓨터

FLOPS

- **FLOPS : floating-point operations per second**
 - 컴퓨터의 성능을 나타내는 지표 중 하나
 - 초당 floating-point operation 횟수
- **Intel Skylake-X architecture (September 2017)**
 - Intel Core i9 7900X (10 cores @ 3.30GHz) : 638.9 GFlops
 - Intel Core i9 7980XE (18 cores @ 2.60GHz) : 977.0 GFlops
 - benchmark results from “Intel Core-i9 7900X and 7980XE”
<https://www.pugetsystems.com/labs/hpc/Intel-Core-i9-7900X-and-7980XE-Skylake-X-Linux-Linpack-Performance-1059/>
- **PC CPU의 최대 성능은 약 1 TFLOPS**

슈퍼 컴퓨터 동향

- from <http://www.top500.org/lists/>

- #1 : 442 PFlops = 442,010 TFlops

captured from www.top500.org, with an e-mail permission.

Rank	System	Cores	Rmax (TFlop/s)	Rpeak (TFlop/s)	Power (kW)
1	Supercomputer Fugaku - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu RIKEN Center for Computational Science Japan	7,630,848	442,010.0	537,212.0	29,899
2	Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, IBM DOE/SC/Oak Ridge National Laboratory United States	2,414,592	148,600.0	200,794.9	10,096
3	Sierra - IBM Power System AC922, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, IBM / NVIDIA / Mellanox DOE/NNSA/LLNL United States	1,572,480	94,640.0	125,712.0	7,438
4	Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway, NRCPC National Supercomputing Center in Wuxi China	10,649,600	93,014.6	125,435.9	15,371
5	Perlmutter - HPE Cray EX235n, AMD EPYC 7763 64C 2.45GHz, NVIDIA A100 SXM4 40 GB, Slingshot-10, HPE	706,304	64,590.0	89,794.5	2,528



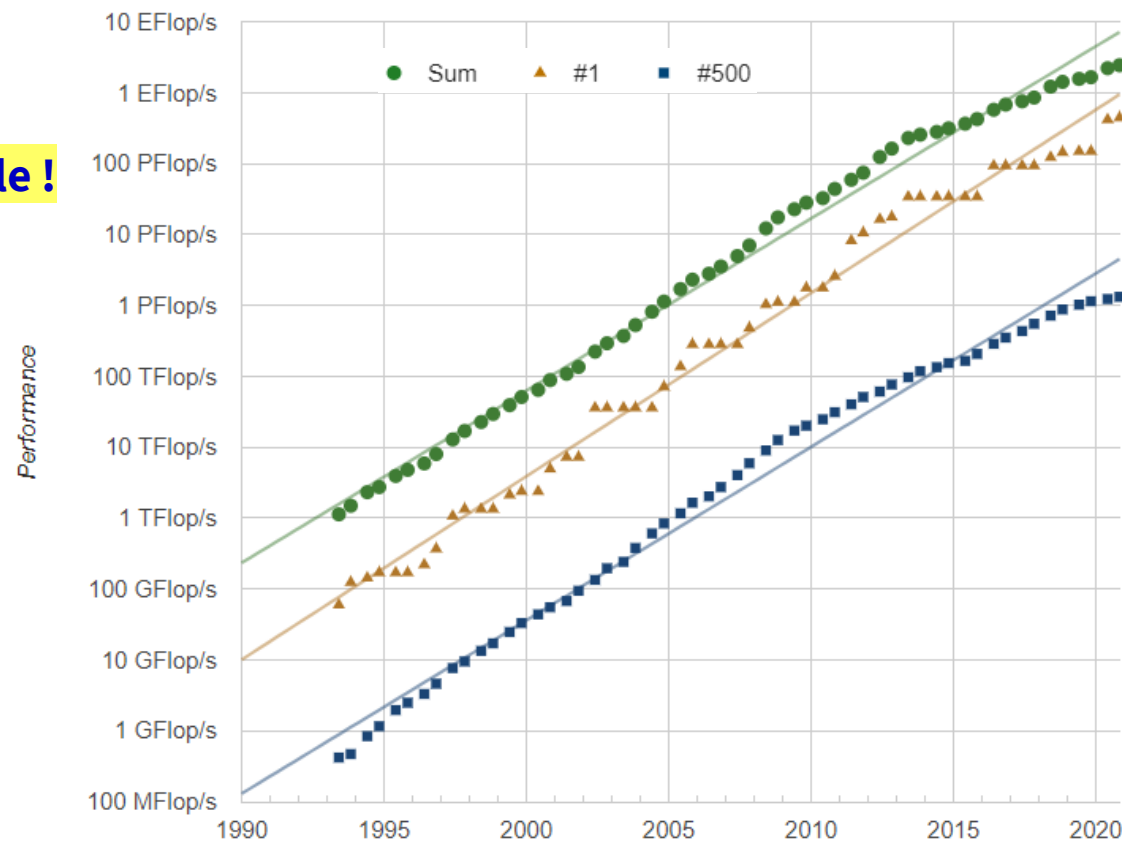
captured from www.top500.org, with an e-mail permission.

슈퍼 컴퓨터 동향 계속

- from <https://www.top500.org/statistics/perfdevel/>

captured from www.top500.org, with an e-mail permission.

Projected Performance Development



captured from www.top500.org, with an e-mail permission.

PC에서의 슈퍼 컴퓨팅

- **Tesla V100-SXM2-16GB**

- 2017년 6월 출시
- 최대 28.26 TFLOPS
- 딥러닝 특화 시, 최대 112 ~ 125 TFLOPS



CC BY-SA 4.0
<https://commons.wikimedia.org/wiki/File:Data-center-tesla-v100-pcie-625-ud@2x.jpg>

- **GeForce RTX 3090**

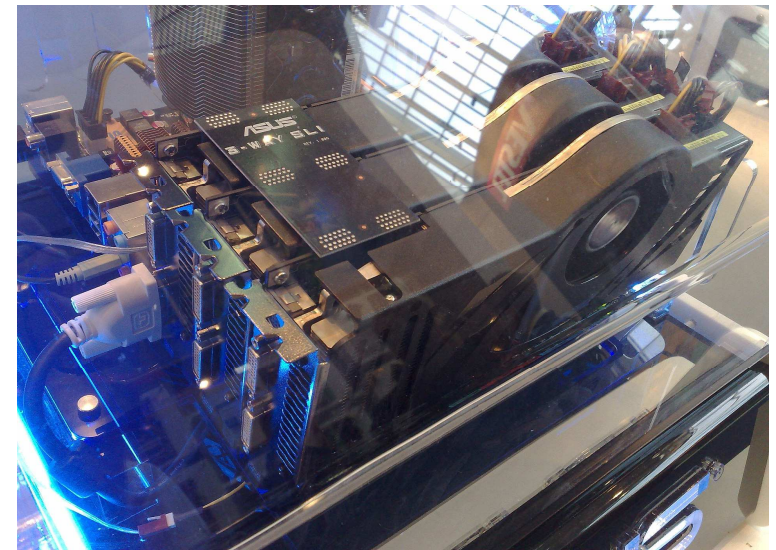
- 2020년 9월 출시
- 최대 36 TFLOPS



CC BY 4.0
https://commons.wikimedia.org/wiki/File:Gigabyte_GeForce_RTX_3090_Eagle_OC_24G_24576_MiB_GDDR6X_legend_Backplate_20201114_DSC5945_Neu.jpg

NVLink, SLI, Crossfire 기술

- NVLink : NVIDIA link
- SLI : Scalable Link Interface, from NVIDIA
- Crossfire Technology, from AMD (former ATI)
 - GPU 2 ~ 4개를 병렬 연결
 - 1개의 GPU 처럼 작동 가능



CC BY-SA 4.0
https://en.wikipedia.org/wiki/Scalable_Link_Interface#/media/File:3-way-SLI.jpg

PC에서의 슈퍼 컴퓨팅

- **Tesla V100-SXM2-16GB**

- 2017년 6월 출시
- 최대 28.26 TFLOPS
- 딥러닝 특화 시, 최대 112 ~ 125 TFLOPS
→ **4개 연결 시, 최대 ~500 TFLOPS**



CC BY-SA 4.0
<https://commons.wikimedia.org/wiki/File:Data-center-tesla-v100-pcie-625-ud@2x.jpg>

- **GeForce RTX 3090**

- 2020년 9월 출시
- 최대 36 TFLOPS
→ **4개 연결 시, 최대 ~144 TFLOPS**



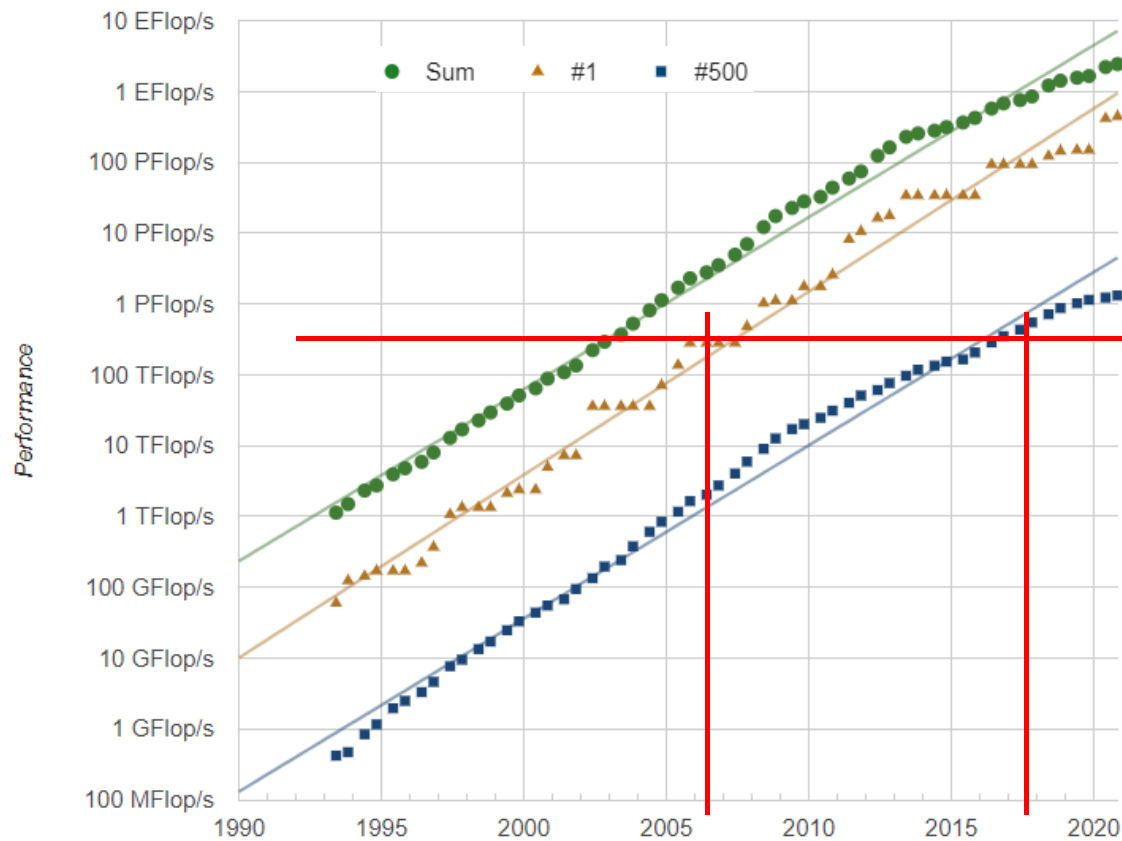
CC BY 4.0
https://commons.wikimedia.org/wiki/File:Gigabyte_GeForce_RTX_3090_Eagle_OC_24G_24576_MiB_GDDR6X_legend_Backplate_20201114_DSC5945_Neu.jpg

슈퍼 컴퓨터 동향 계속

- from <https://www.top500.org/statistics/perfdevel/>

captured from www.top500.org, with an e-mail permission.

Projected Performance Development



captured from www.top500.org, with an e-mail permission.

또다른 시도

- 암호 화폐 채굴기 **crypto mining machine**
 - PCI Express 버스로 GPU 연결 가능



CC BY 2.0
<https://www.flickr.com/photos/bitcoin-crypto/42160336234>

내용 contents

- CUDA의 개발 배경 – 컴퓨터 그래픽스
- CUDA의 구성
- CUDA의 성능 – 슈퍼 컴퓨터 동향
- CUDA 쿼다 = Compute Unified Device Architecture, from NVIDIA 엔비디아

CUDA 소개

폰트 정상 출력 → 큰 교자 타고 혼례 치른 날
정참판 양반댁 규수 큰 교자 타고 혼례 치른 날
정참판 양반댁 규수 큰 교자 타고 혼례 치른 날
본고딕 Noto Sans KR

The quick brown fox jumps over the lazy dog
The quick brown fox jumps over the lazy dog
The quick brown fox jumps over the lazy dog
Source Sans Pro