# Arhitecturi Paralele
## Introducere GPU
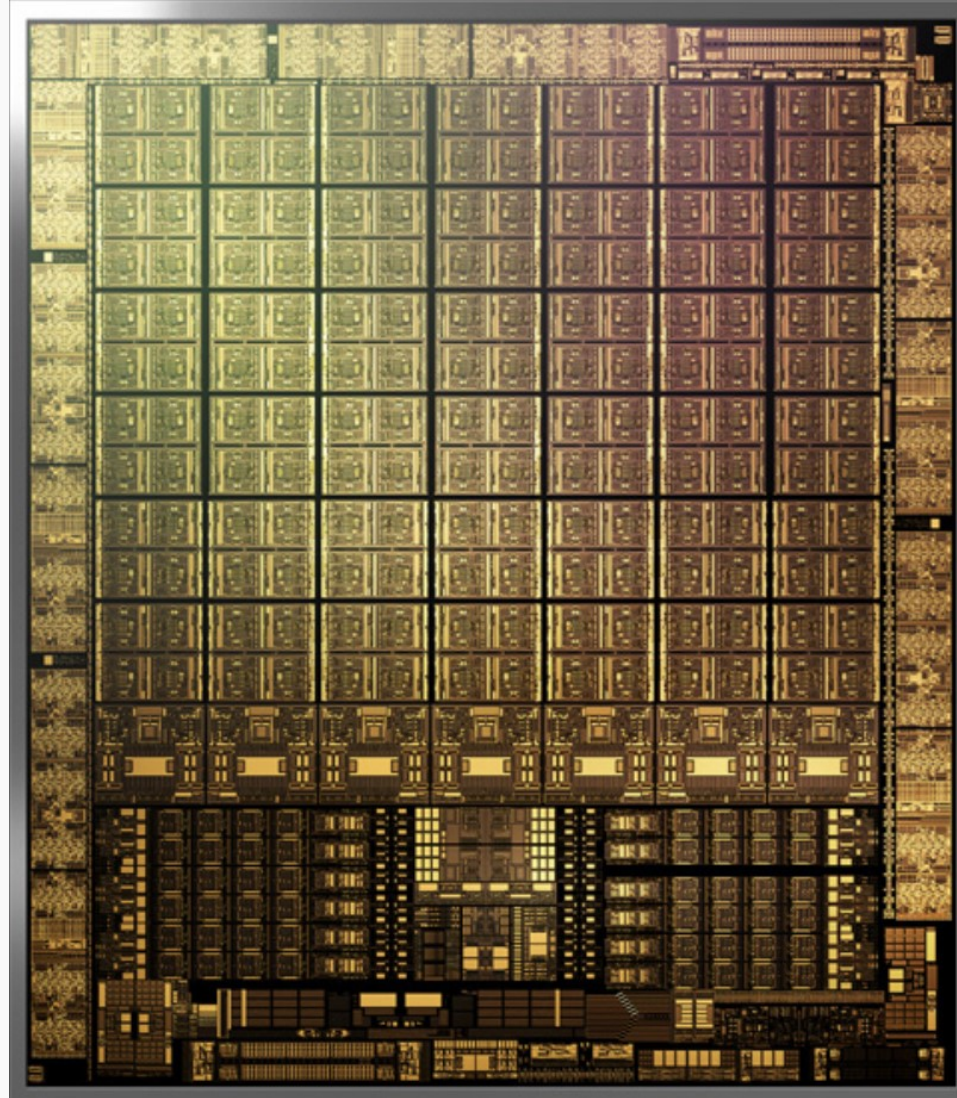
Lect. Dr. Ing. Cristian Chilipirea – cristian.chilipirea@mta.ro
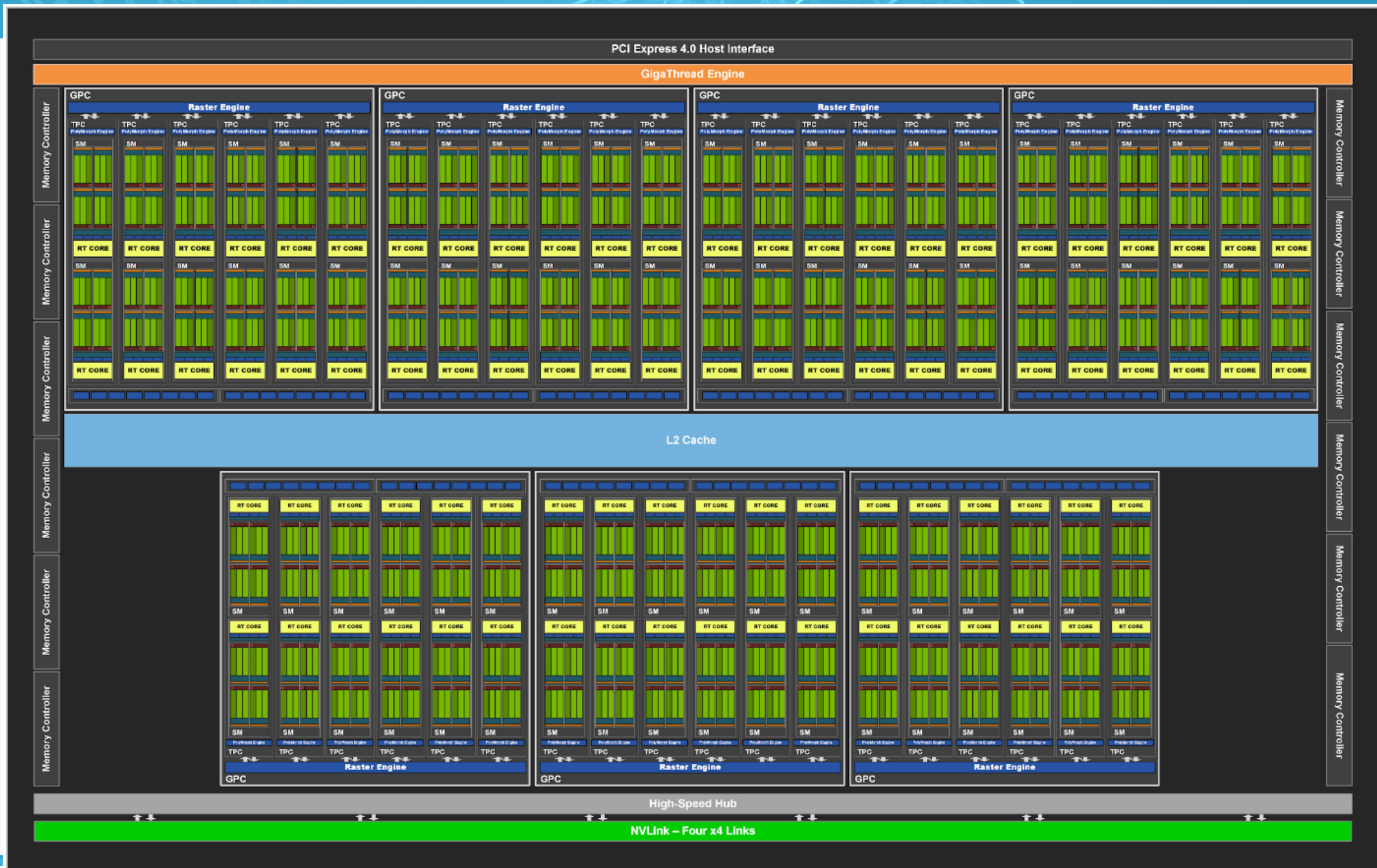
# NVIDIA Ampere

# NVIDIA Ampere v2

| | | GEFORCE RTX 3080 Ti | GEFORCE RTX 3080 |
|---|---|---|---|
| **Specificații GPU:** | Nuclee NVIDIA CUDA® | 10240 | 8704 |
| | Frecvență Boost (GHz) | 1.67 | 1.71 |
| | Frecvență de bază (GHz) | 1.37 | 1.44 |
| **Specificații memorie:** | Configurație memorie standard | 12 GB GDDR6X | 10 GB GDDR6X |
| | Lățime interfață memorie | 384 biți | 320 biți |
| **Tehnologii integrate:** | Nuclee cu ray-tracing | Cea de-a doua generație | Cea de-a doua generație |
| | Nuclee Tensor | Cea de-a treia generație | Cea de-a treia generație |
| | Arhitectură NVIDIA | Ampere | Ampere |

# RTX 3080

# NVIDIA A40

## SPECIFICATIONS

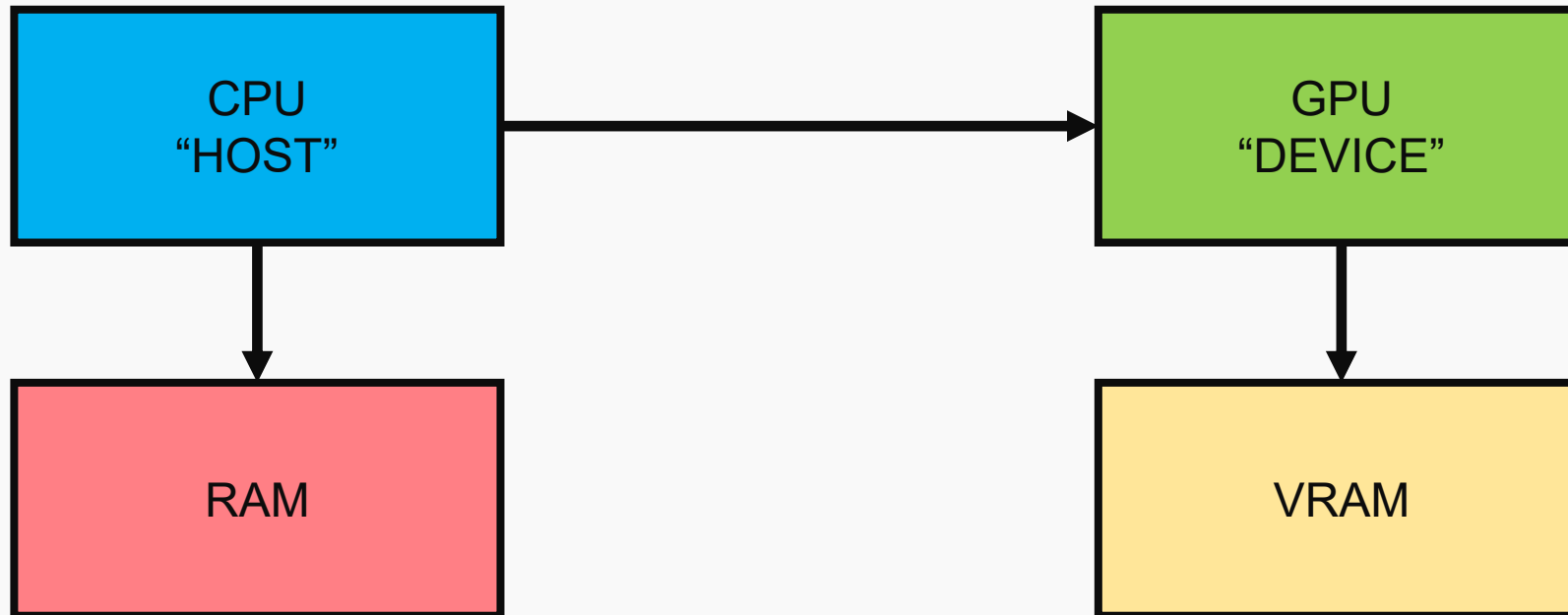| | |
|---|---|
| GPU architecture | NVIDIA Ampere architecture |
| GPU memory | 48 GB GDDR6 with ECC |
| Memory bandwidth | 696 GB/s |
| Interconnect interface | NVIDIA® NVLink® 112.5 GB/s (bidirectional)³ PCIe Gen4 31.5 GB/s (bidirectional) |
| NVIDIA Ampere architecture-based CUDA Cores | 10,752 |
| NVIDIA second-generation RT Cores | 84 |
| NVIDIA third-generation Tensor Cores | 336 |

# CPU vs GPU?

- Cores
- Frequency
- Core complexity

# Arhitectura system heterogen

# Typical Program

- CPU alocă memorie pe GPU (în VRAM)
- CPU copiază date din RAM în VRAM
- CPU pornește **kernelul** pe GPU
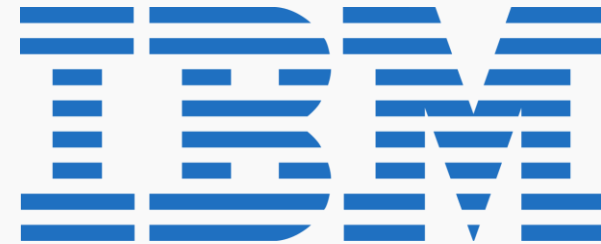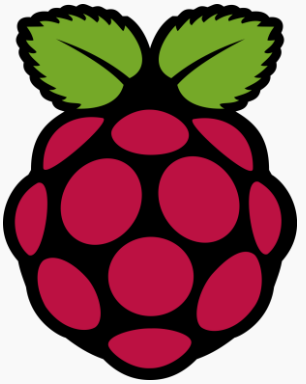- CPU copiază date din VRAM pe RAM

# Open Standard for Parallel Programming of Heterogeneous Systems

Cristian Chilipirea – Arhitecturi Paralele
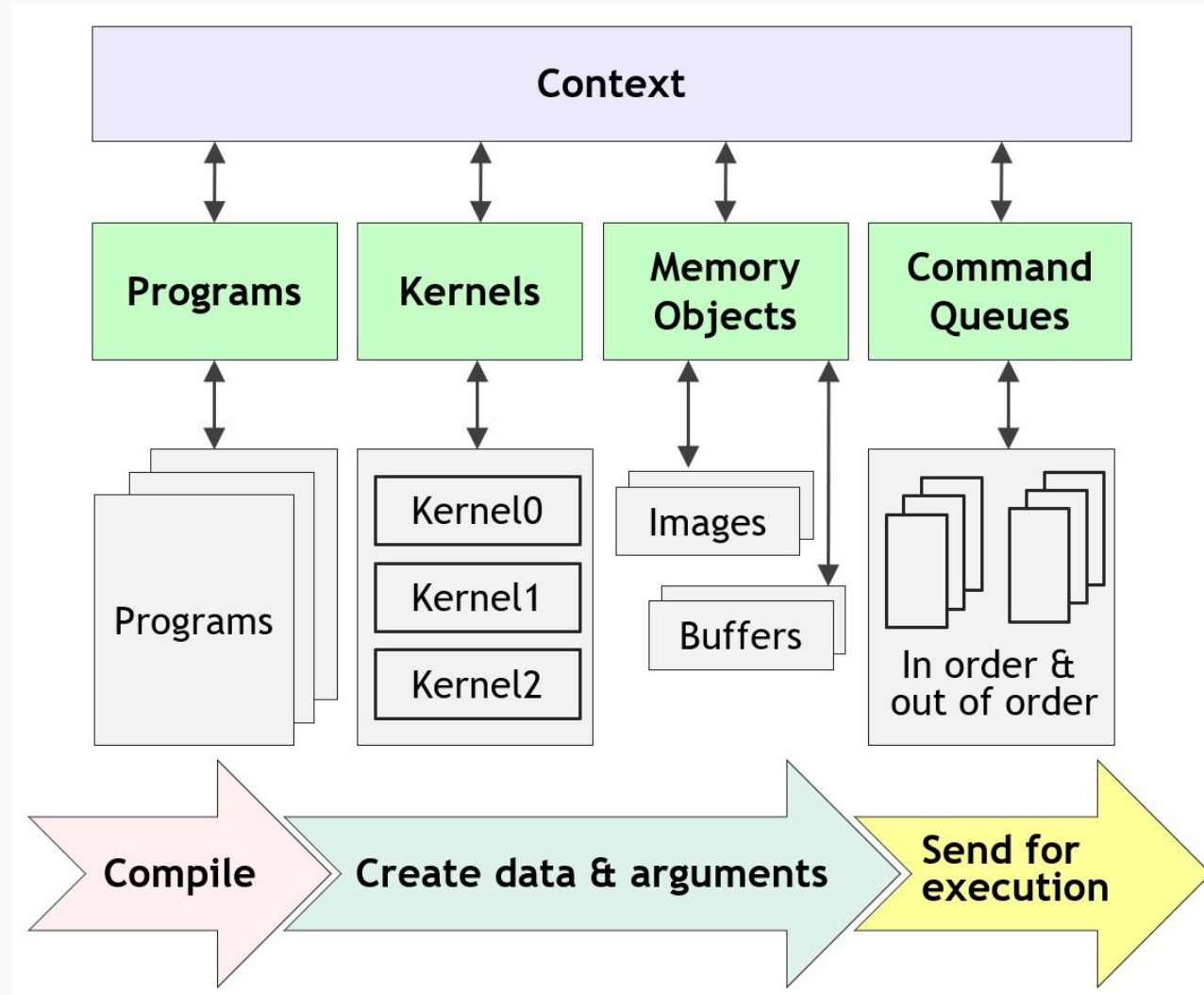
# Sequence for Executing OpenCL Kernels

# A complete sequence for executing an OpenCL program

- Query for available OpenCL platforms and devices
- Create a context for one or more OpenCL devices in a platform
- Create and build programs for OpenCL devices in the context
- Select kernels to execute from the programs
- Create memory objects for kernels to operate on
- Create command queues to execute commands on an OpenCL device
- *Enqueue* data transfer commands into the memory objects, if needed
- *Enqueue* kernels into the command queue for execution
- *Enqueue* commands to transfer data back to the host, if needed

# Traditional Versus OpenCL Programming Using OpenCL C Kernels

# OpenCL Language Ecosystem Enabled With SPIR-V
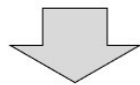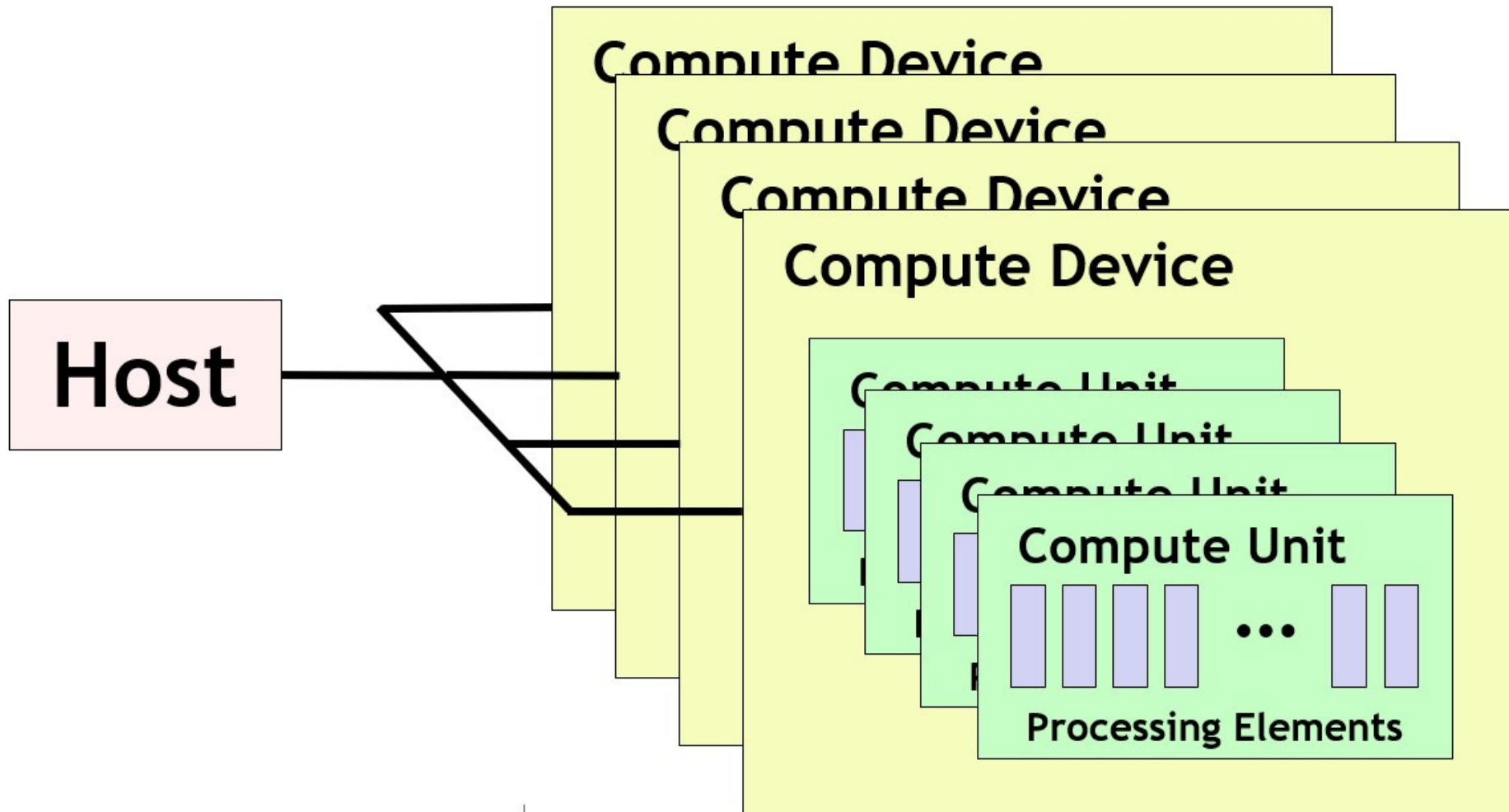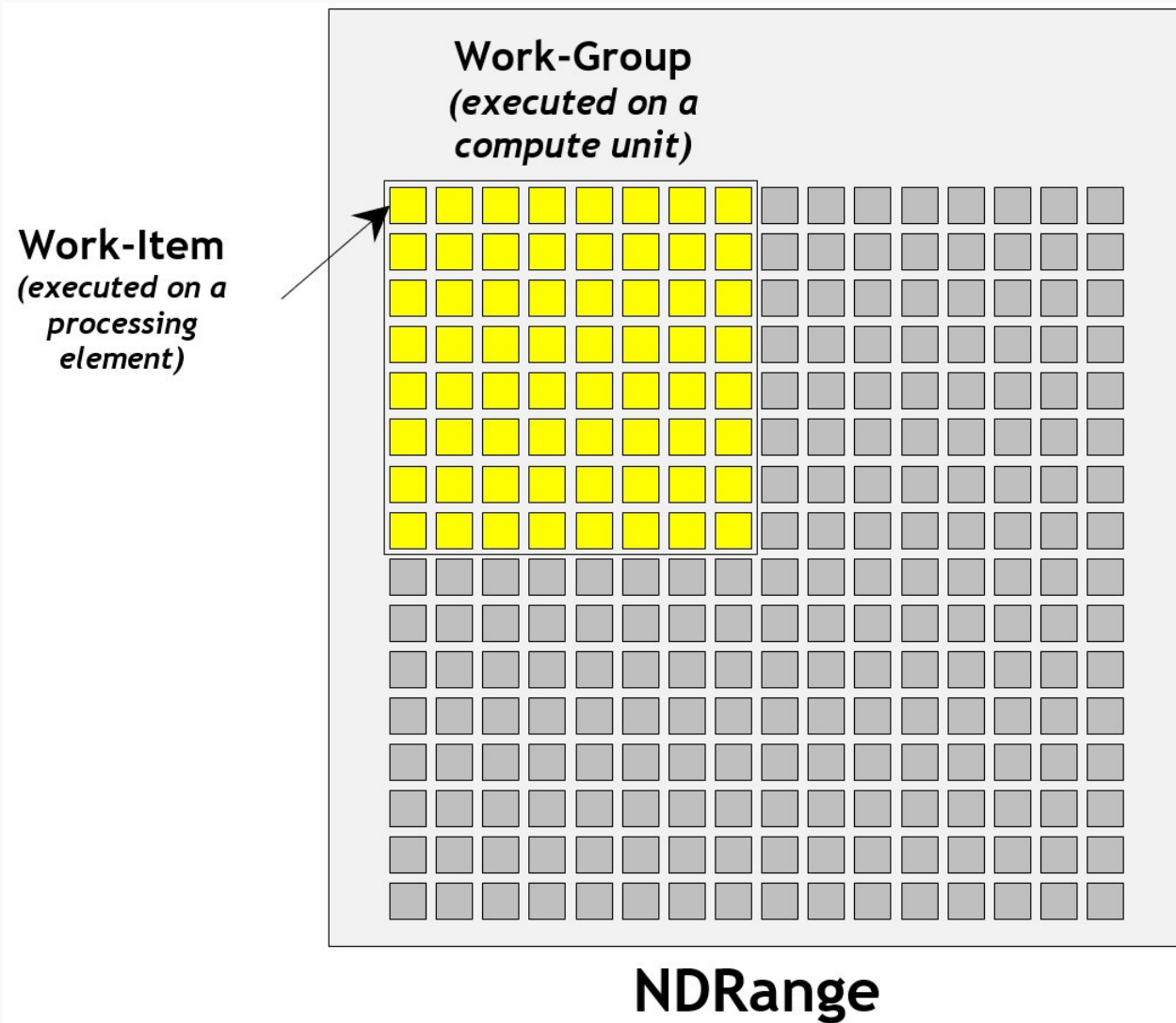
# OpenCL Platform Model

# A 2D Image as an Example NDRange

# OpenCL Memory Model