# Nvidia GPU microarchitecures

# Nvidia microarchitectures

- Tesla
- **Fermi**
- **Kepler**
- Maxwell
- **Pascal**
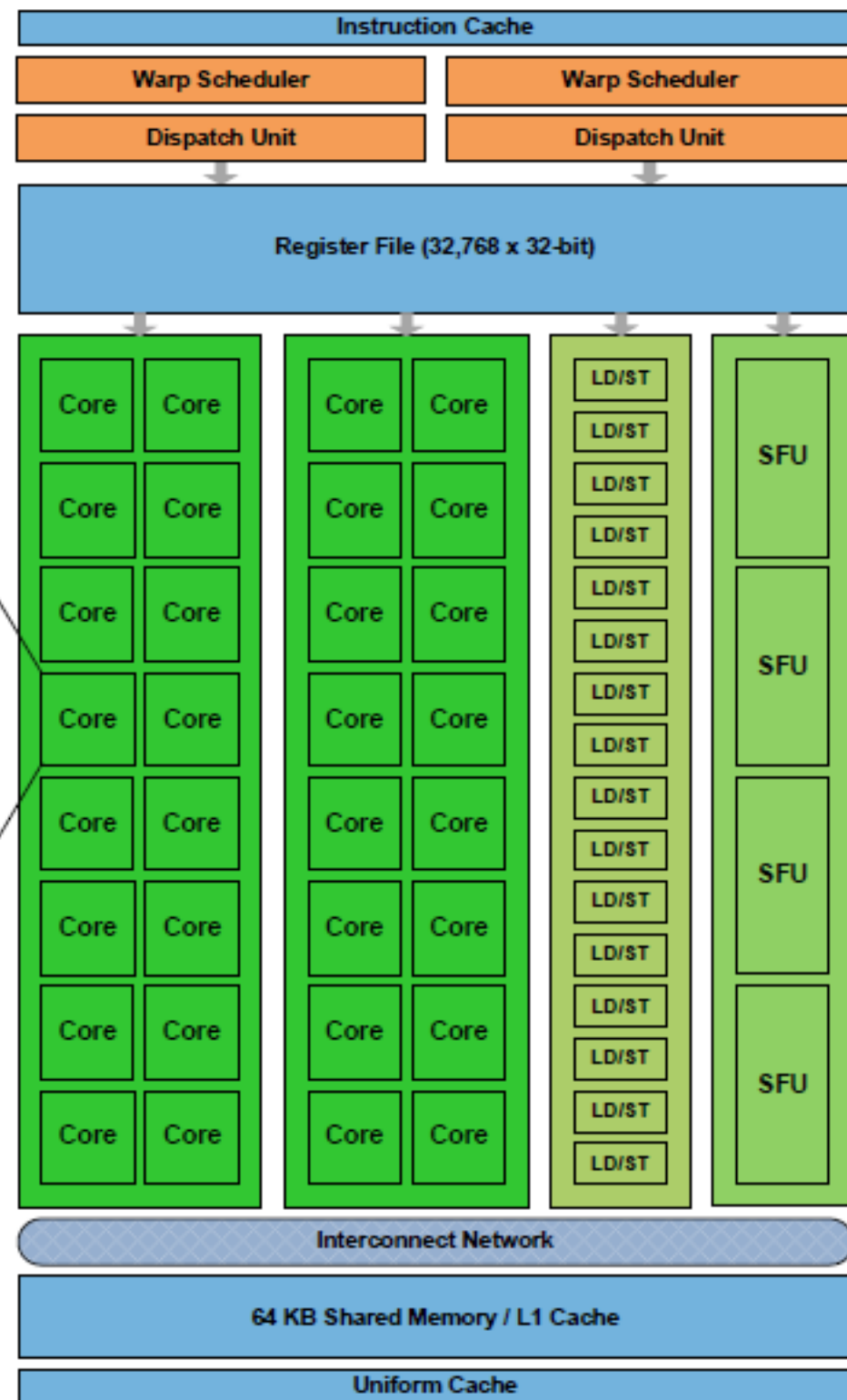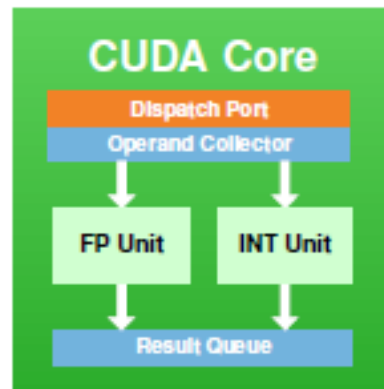- **Volta**
- Turing
- **Ampere**

# Kepler GPU

# Streaming Multiprocessor(SM)

- Basic building block

- Contains
  - Multiple CUDA cores
  - Special Function Units (SFUs)
  - Register file
  - L1 cache / Shared memory
  - ...

# Fermi

# Kepler

- 192 single-precision CUDA cores
- 64 double-precision units
- 32 special function units
- 32 load/store units
- 4 warp schedulers
- 8 instruction dispatch units
- 48 KB Read-Only Data Cache
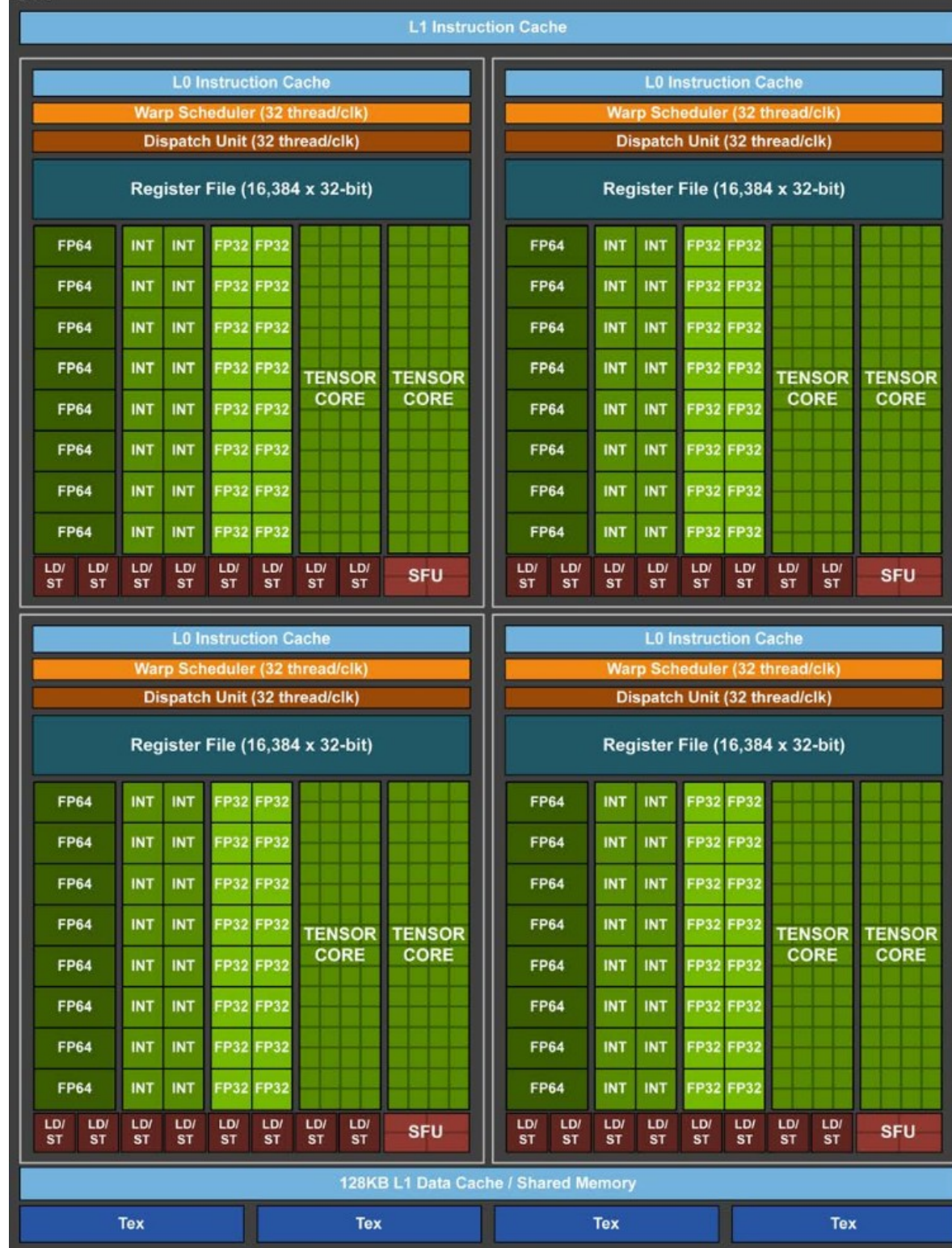- 64/128 KB L1 cache / shared memory

# Pascal

# Pascal (contd)

- two processing blocks, each with
  - 32 single-precision CUDA Cores
  - 16 double precision (FP64) CUDA Cores
  - one instruction buffer
  - one warp scheduler
  - two dispatch units

# Volta

Four processing blocks

# Volta(contd)

each processing block has
- 16 FP32 Cores
- 8 FP64 Cores
- 16 INT32 Cores
- 2 new mixed-precision Tensor Cores
- one L0 instruction cache
- one warp scheduler
- one dispatch unit
- 64 KB Register File

Ampere