# Tutorial 4: Compilers, linkers and performance
## Informatik elective: GPU Computing

Pratik Nayak

# NVIDIA GPU architectures

| Arch → | Tesla | Fermi | Kepler | Maxwell | Pascal | Volta | Turing | Ampere | Ada Lovelace | Hopper | Blackwell |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Compute capability → | 1.0 | 2.0 | 3.0 | 5.x | 6.x | 7.0 | 7.5 | 8.0 | 8.9 | 9.0 | Blackwell |
| GPUs → | | | | | P100 | V100 | | A100 | | H100 | |

New architectures → new features, more performance.

Summit supercomputer (ORNL, USA)

JUPITER (Jülich, DE)
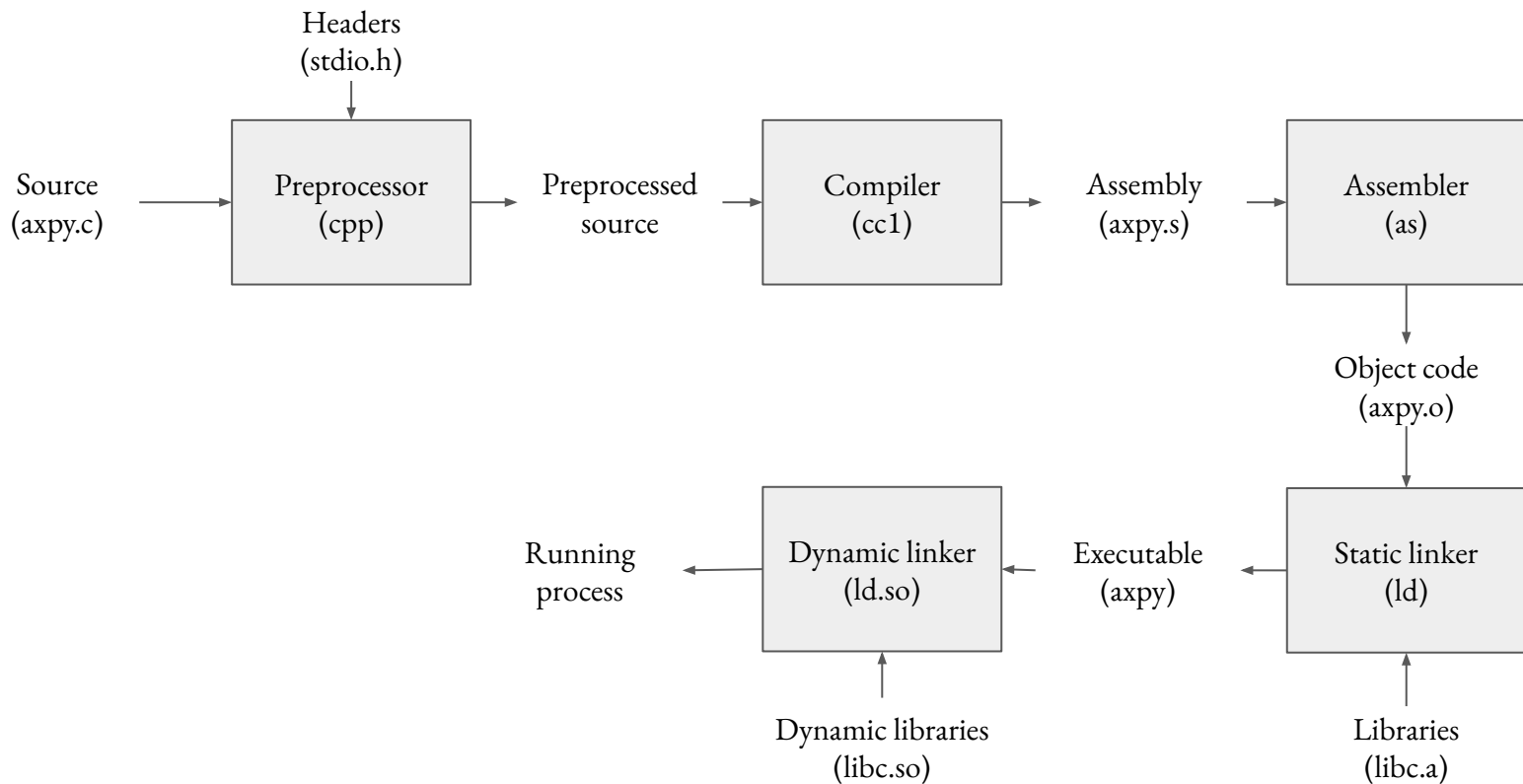
Pratik Nayak - GPU Computing

Computational Mathematics Group (CIT)

# NVIDIA GPU architectures

| Feature support (unlisted features are supported for all compute capabilities) | Compute capability (version) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | 1.0, 1.1 | 1.2, 1.3 | 2.x | 3.0 | 3.2 | 3.5, 3.7, 5.x, 6.x, 7.0, 7.2 | 7.5 | 8.x | 9.0 |
| Warp vote functions (__all(), __any()) | No | Yes | | | | | | | |
| Warp vote functions (__ballot()) | No | | Yes | | | | | | |
| Memory fence functions (__threadfence_system()) | | | | | | | | | |
| Synchronization functions (__syncthreads_count(), __syncthreads_and(), __syncthreads_or()) | | | | | | | | | |
| Surface functions | | | | | | | | | |
| 3D grid of thread blocks | | | | | | | | | |
| Warp shuffle functions | | No | | Yes | | | | | |
| Unified memory programming | | | | | | | | | |
| Funnel shift | | No | | | Yes | | | | |
| Dynamic parallelism | | | No | | | Yes | | | |
| Uniform Datapath[57] | | | | No | | | Yes | | |
| Hardware-accelerated async-copy | | | | | No | | | Yes | |
| Hardware-accelerated *split arrive/wait barrier* | | | | | | | | | |
| Warp-level support for reduction ops | | | | | | | | | |
| L2 cache residency management | | | | | | | | | |
| DPX instructions for accelerated dynamic programming | | | | | No | | | | Yes |
| Distributed shared memory | | | | | | | | | |
| Thread block cluster | | | | | | | | | |
| Tensor memory accelerator (TMA) unit | | | | | | | | | |
| Feature support (unlisted features are supported for all compute capabilities) | 1.0,1.1 | 1.2,1.3 | 2.x | 3.0 | 3.2 | 3.5, 3.7, 5.x, 6.x, 7.0, 7.2 | 7.5 | 8.x | 9.0 |
| | Compute capability (version) | | | | | | | | |

[Wikimedia: wikipedia.org/wiki/CUDA]

# C build process

Headers
(stdio.h)

Source
(axpy.c) → Preprocessor
(cpp) → Preprocessed
source → Compiler
(cc1) → Assembly
(axpy.s) → Assembler
(as)

Object code
(axpy.o)

Running
process ← Dynamic linker
(ld.so) ← Executable
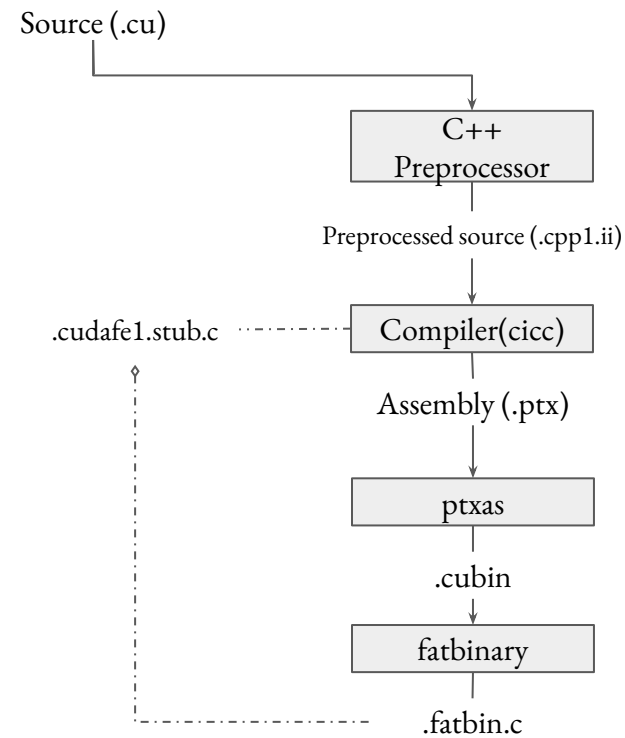(axpy) ← Static linker
(ld)

Dynamic libraries
(libc.so)

Libraries
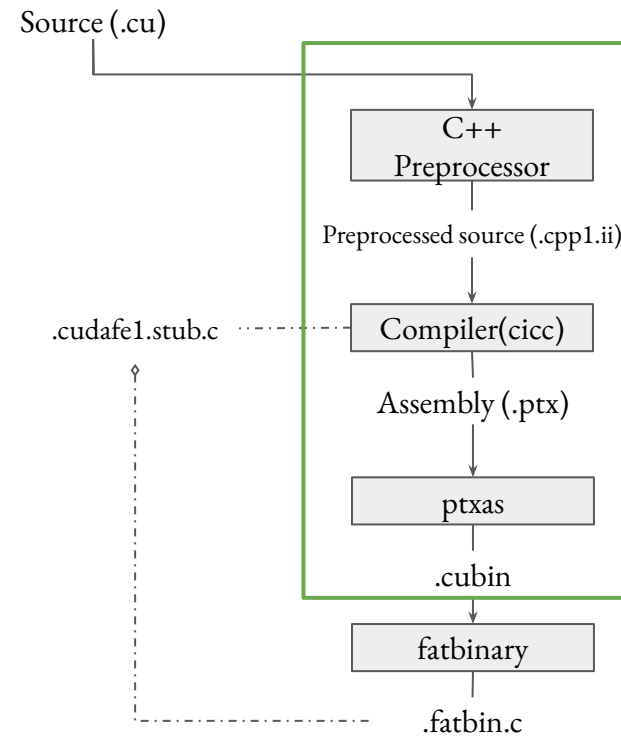(libc.a)

# NVIDIA CUDA compilation process

- NVIDIA CUDA compiler driver: `nvcc`

  - Performs the necessary steps to build a CUDA executable/library.

  - Includes compilation of both host code and device code.

- See https://docs.nvidia.com/cuda/cuda-compiler-driver-nvcc/index.html for documentation and more

  details.

Pratik Nayak - GPU Computing Computational Mathematics Group (CIT)

# CUDA build process

Source (.cu)

C++
Preprocessor

Preprocessed source (.cpp1.ii)

.cudafe1.stub.c ········· Compiler(cicc)

Assembly (.ptx)

ptxas

.cubin

fatbinary

.fatbin.c

# CUDA build process

Source (.cu)

C++
Preprocessor

Preprocessed source (.cpp1.ii)

.cudafe1.stub.c ----- Compiler(cicc)

Assembly (.ptx)

ptxas

.cubin

fatbinary

.fatbin.c

# CUDA build process

Source (.cu)

| | |
|---|---|
| C++ Preprocessor | C++ Preprocessor |

Preprocessed source (.cpp4.ii)          Preprocessed source (.cpp1.ii)

Compiler (cudafe++)          .cudafe1.stub.c          Compiler(cicc)

Assembly (.ptx)

.cudafe1.cpp

ptxas

.cubin

C++ compiler

fatbinary

.o/.obj          .fatbin.c

# CUDA build process

Source (.cu)

C++ Preprocessor

Preprocessed source (.cpp4.ii)

Compiler (cudafe++)

.cudafe1.cpp

C++ compiler

.o/.obj

.cudafe1.stub.c

C++ Preprocessor

Preprocessed source (.cpp1.ii)

Compiler(cicc)

Assembly (.ptx)

ptxas

.cubin

fatbinary

.fatbin.c

# CUDA build process

- Repeat ☐ for each architecture
- Repeat ⬚ for each .cu file



executable

Host Linker

a_dlink.o/a_dlink.obj

C++ compiler

link.stub

a_dlink.fatbin.c

fatbinary

a_dlink.cubin

a_dlink.reg.c

nvlink

Source (.cu)

C++ Preprocessor

Preprocessed source (.cpp4.ii)

Compiler (cudafe++)

.cudafe1.cpp

.cudafe1.stub.c

C++ compiler

.o/.obj

.fatbin.c

C++ Preprocessor

Preprocessed source (.cpp1.ii)

Compiler(cicc)

Assembly (.ptx)

ptxas

.cubin

fatbinary

Pratik Nayak - GPU Computing

Computational Mathematics Group (CIT)

# Helpful tools (Linux)

- Dump assembly from object file:

```
objdump -d file.o                                    # See documentation: man objdump
```

- Look at available symbols in the object file

```
nm file.o                                            # See documentation: man nm
```

- Look at linked libraries (dependencies on other shared objects)

```
ldd executable                                       # See documentation: man ldd
```

Pratik Nayak - GPU Computing

Computational Mathematics Group (CIT)

# Measuring performance (best practices)

- Compilation flags are appropriate: For example with all the required optimization flags and no debug flags.

- Eliminate noise by repeating measurement a few times and averaging the overall time.

- Ensure operation is complete (synchronize explicitly if necessary).

- Ensure clock used is steady and cannot be arbitrarily updated (system_clock v/s steady_clock).

- Ensure operation is actually being performed (not skipped due to predication, or due to some clause).

- Try to eliminate common sources of interference: other users, some unrelated I/O, system effects.

Highly recommended read: https://blogs.fau.de/hager/archives/category/fooling-the-masses

Pratik Nayak - GPU Computing

Computational Mathematics Group (CIT)