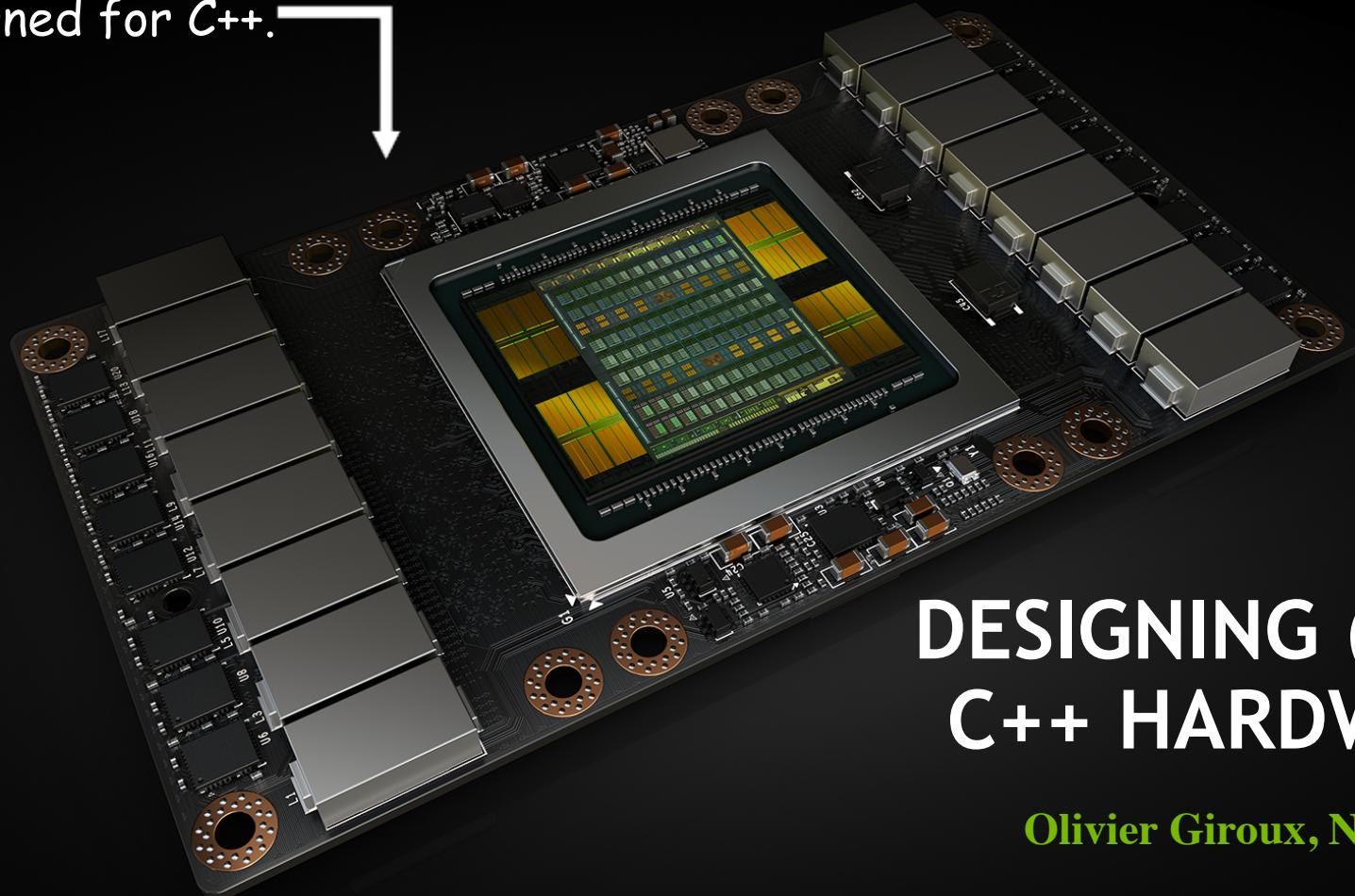


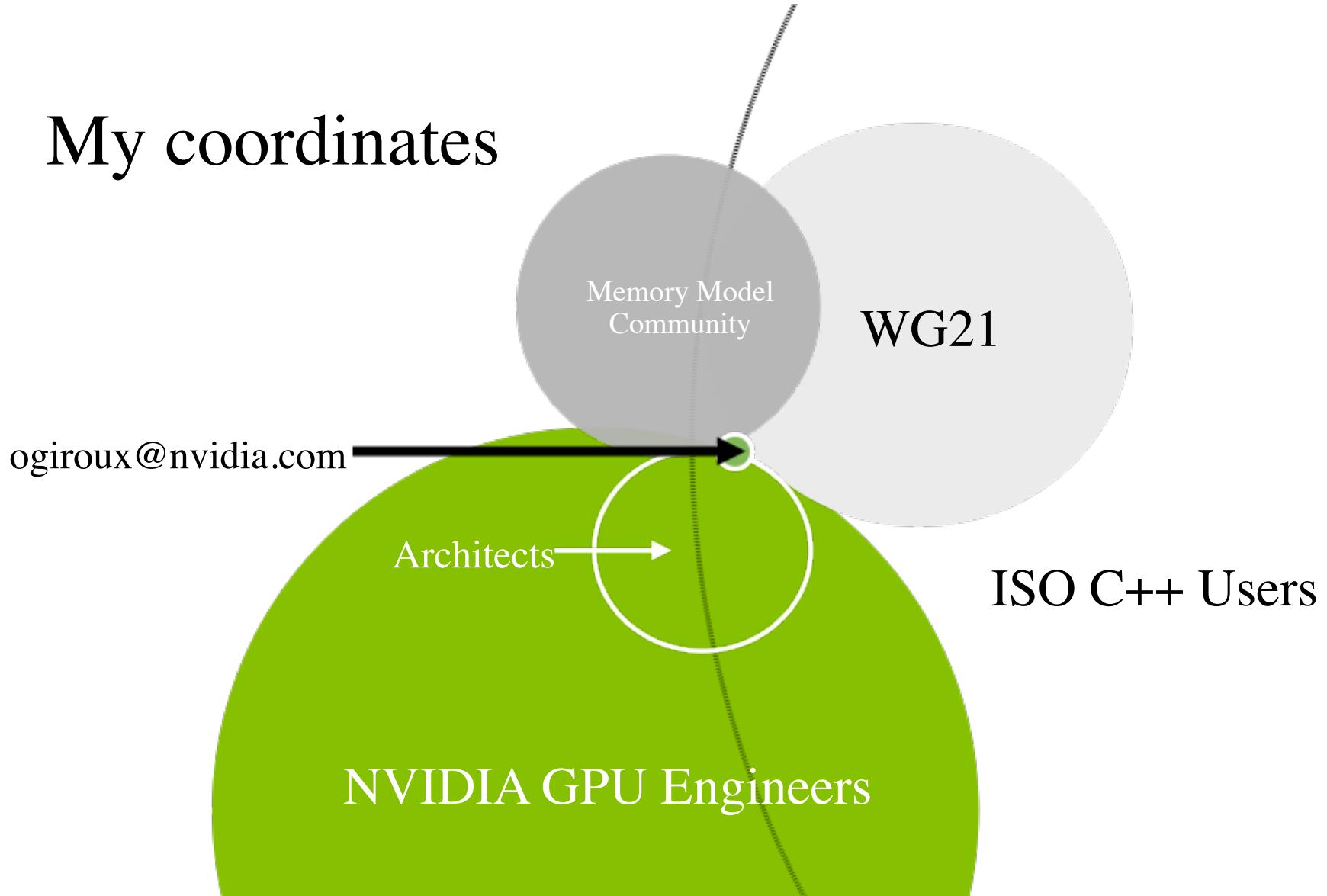
Designed for C++.



DESIGNING (NEW) C++ HARDWARE

Olivier Giroux, NVIDIA

My coordinates



AGENDA

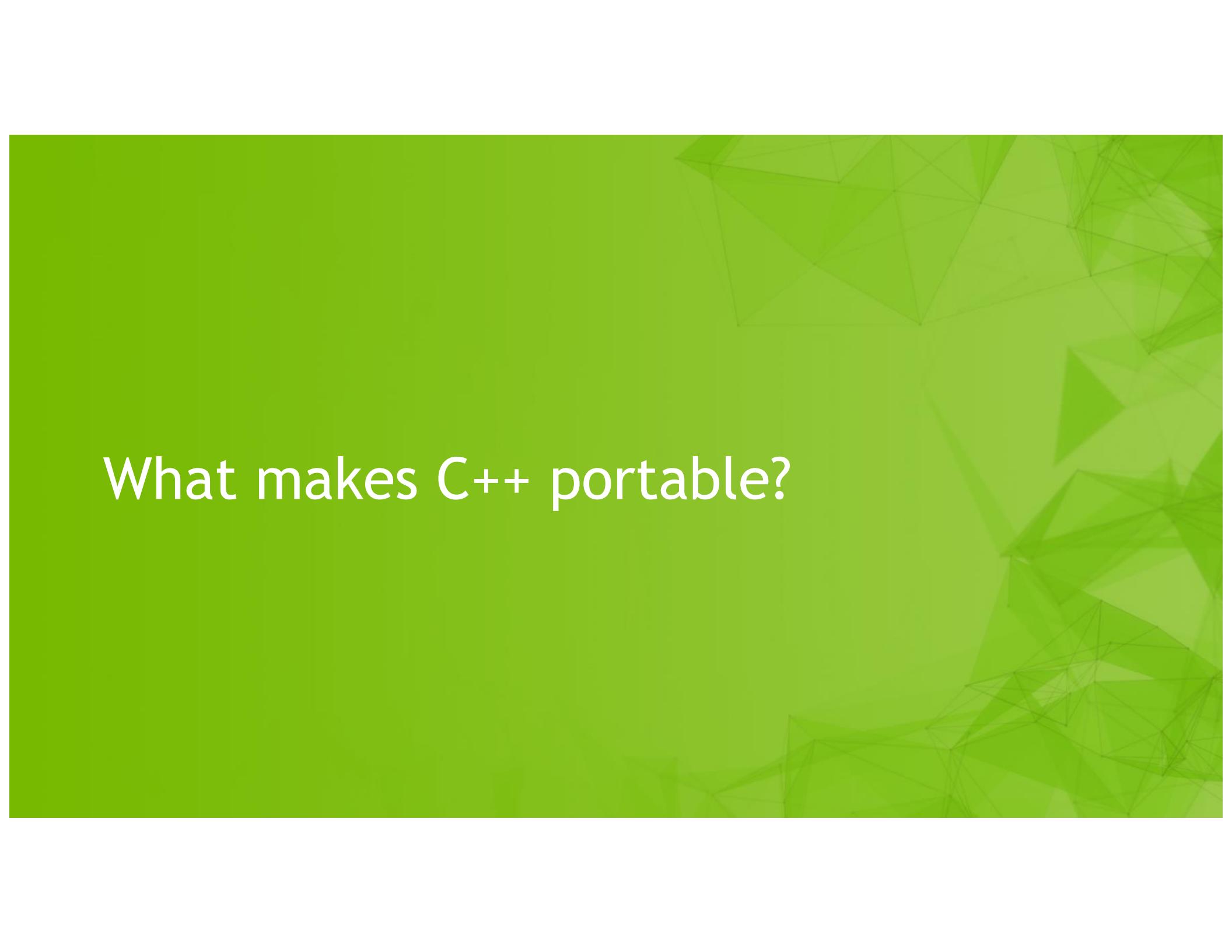
Not code

Not code

Not code

Bad code

Not code

The background of the slide features a subtle, abstract geometric pattern composed of numerous thin, light-green lines forming a triangular mesh. This pattern is more dense and prominent in the lower right quadrant, while the upper left quadrant is a solid, darker shade of green.

What makes C++ portable?

WHAT MAKES C++ PORTABLE?

SUPPORT IN C++

Non-8-bit char

Noncommittal sizeof

Non-2's comp. int

Non-IEEE float

Non-endian addressing

Aligned addressing

Segmented memory

WHAT MAKES C++ PORTABLE?

SUPPORT IN C++	STORAGE
Non-8-bit char	✓
Noncommittal sizeof	✓
Non-2's comp. int	✓
Non-IEEE float	✓
Non-endian addressing	✓
Aligned addressing	✓
Segmented memory	✓

WHAT MAKES C++ PORTABLE?

SUPPORT IN C++	STORAGE	ARITHMETIC (UB)
Non-8-bit char	✓	✓
Noncommittal sizeof	✓	✓
Non-2's comp. int	✓	✓
Non-IEEE float	✓	✓
Non-endian addressing	✓	✓
Aligned addressing	✓	✓
Segmented memory	✓	✓

WHAT MAKES C++ PORTABLE?

SUPPORT IN C++	STORAGE	ARITHMETIC (UB)	20TH CENTURY TRIVIAL PURSUIT
Non-8-bit char	✓	✓	✓
Noncommittal sizeof	✓	✓	✓
Non-2's comp. int	✓	✓	✓
Non-IEEE float	✓	✓	✓
Non-endian addressing	✓	✓	✓
Aligned addressing	✓	✓	Useful.
Segmented memory	✓	✓	Useful.

WHY IS MOST OF THIS NOT HELPFUL?

FALSE CHOICES

Most options are dictated.

New CPU? Match AARCH64.

New GPU? Match the host.

GPUs match *all the hosts*.

BAD CHOICES

Most alternatives are bad.

Negligible area savings.

Negligible power savings.

Programmer surprise.

IT'S TOO RISKY IN 2017

NVIDIA TAPE-OUTS *versus* C++ REVISIONS

C++ implementations go all the way down to atoms.

C++98

C++03

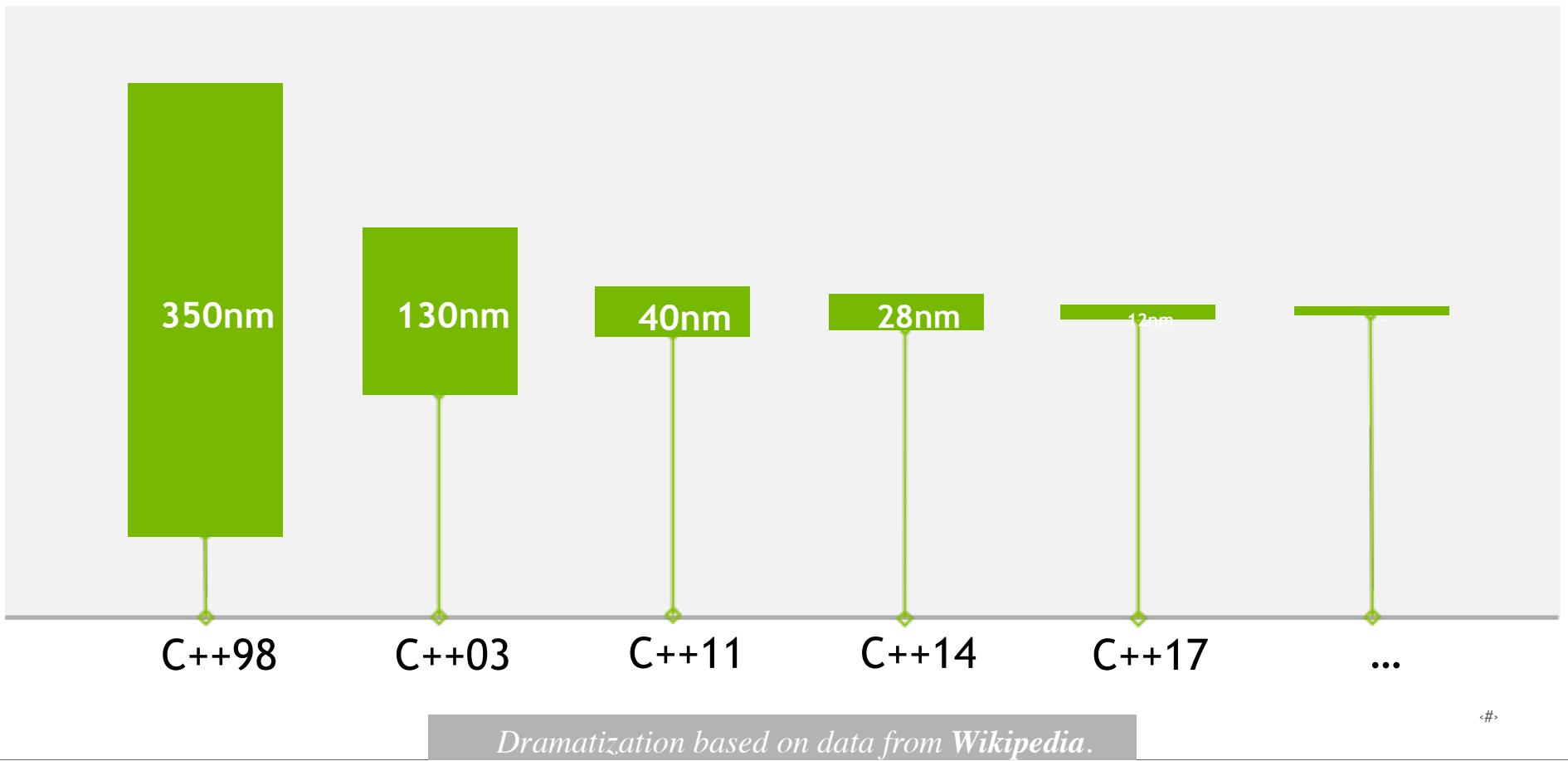
C++11

C++14

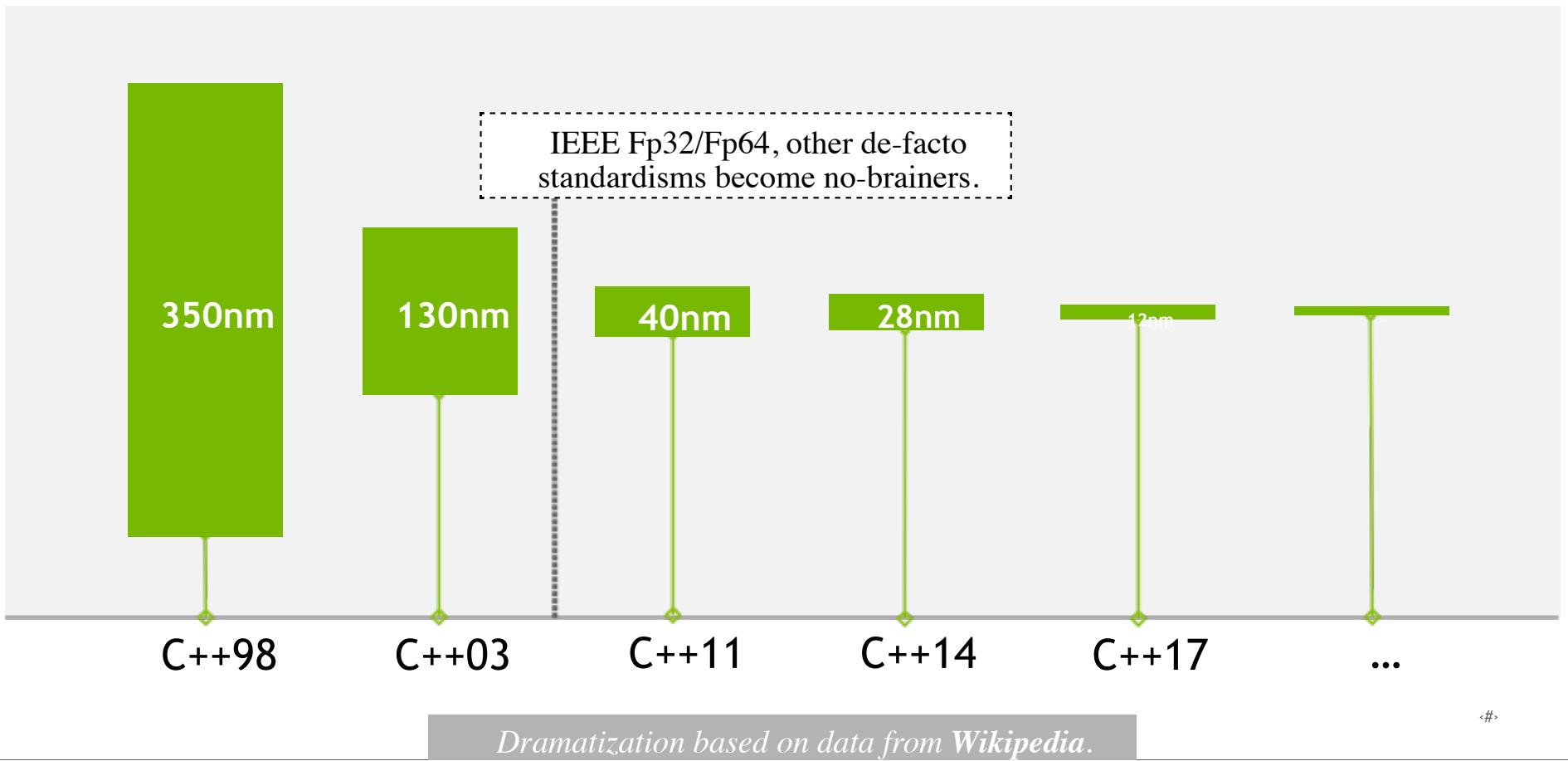
C++17

...

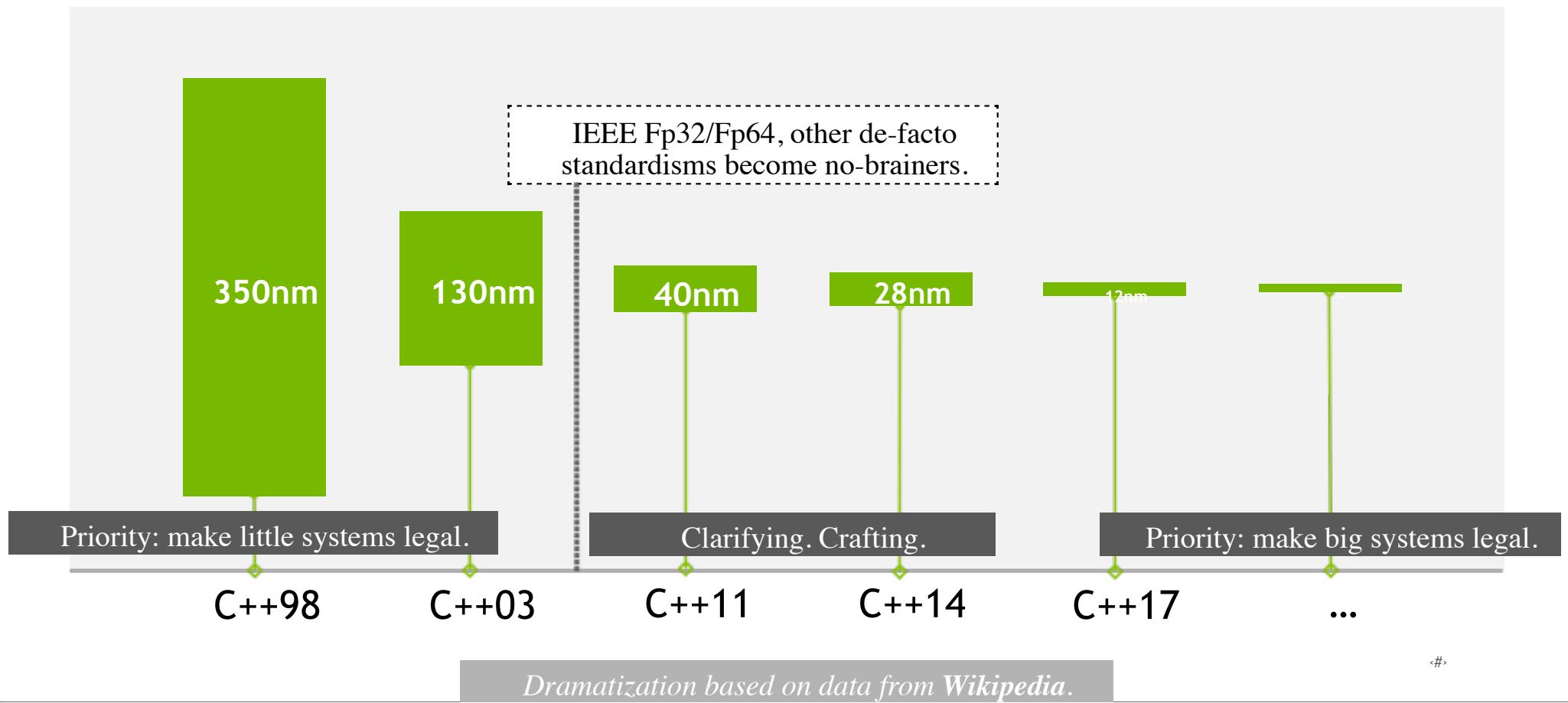
NVIDIA TAPE-OUTS *versus* C++ REVISIONS



NVIDIA TAPE-OUTS *versus* C++ REVISIONS



NVIDIA TAPE-OUTS *versus* C++ REVISIONS

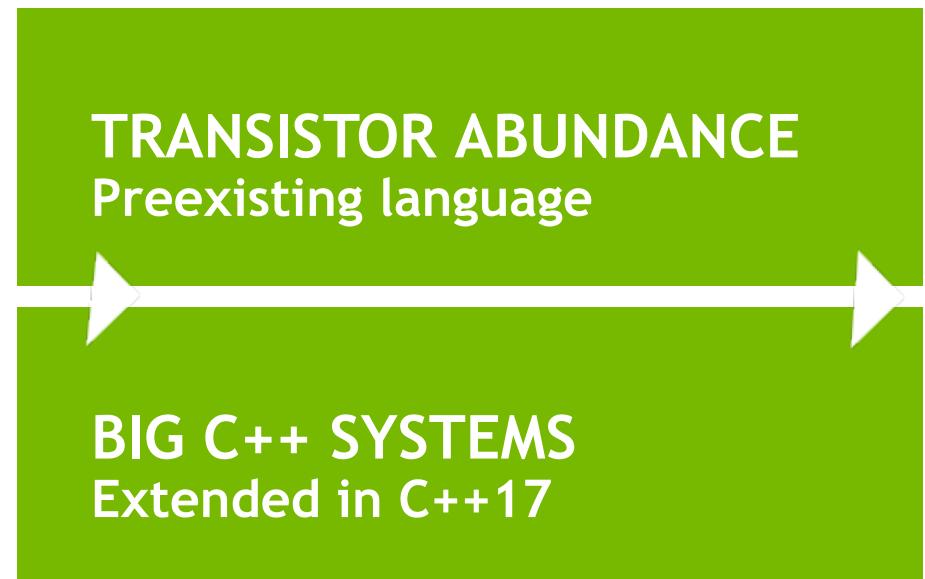


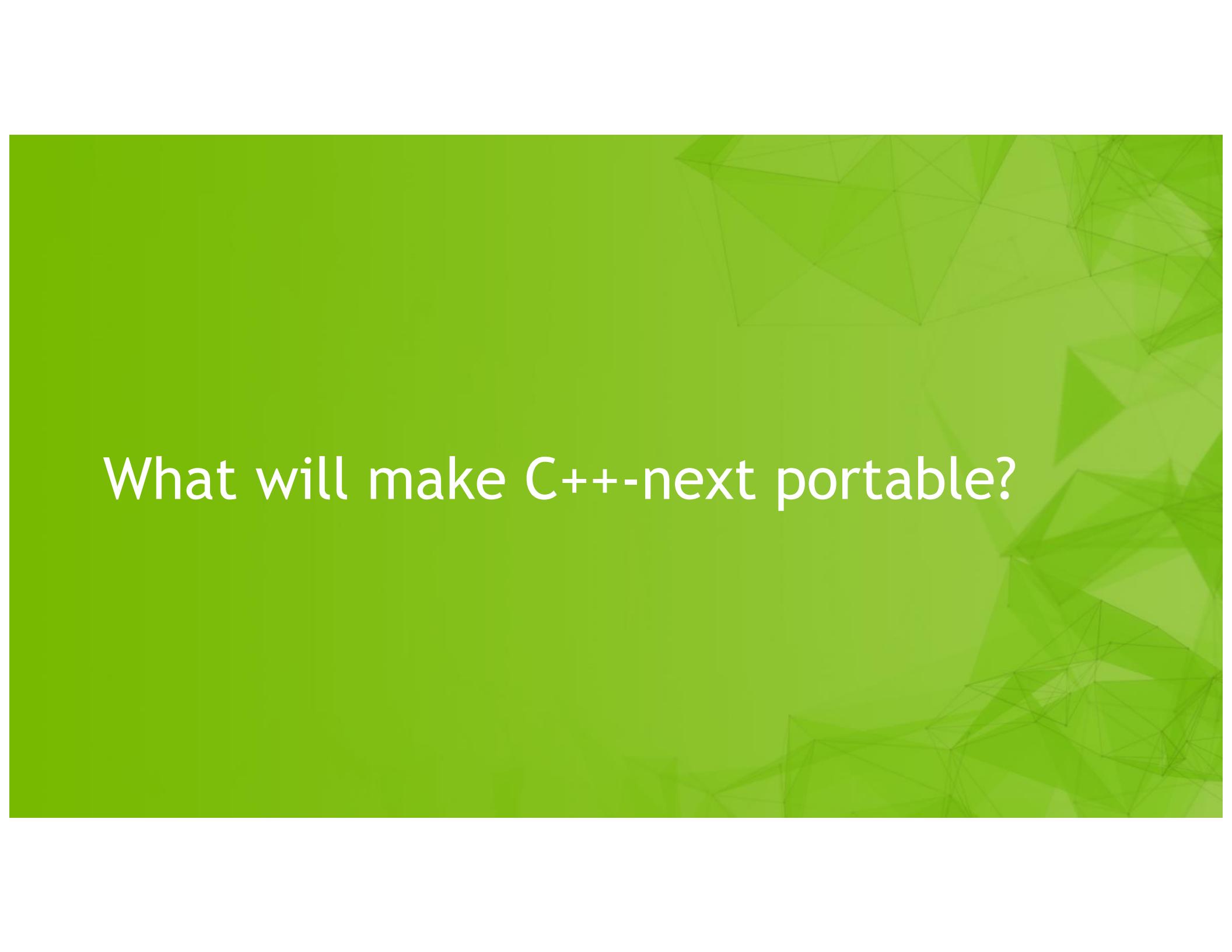
TOWARDS NON-'CLASSIC' PORTABILITY

Then:



Now:



The background of the slide features a subtle, abstract geometric pattern composed of numerous thin, light-green lines forming a triangular mesh. This pattern is more dense in the lower right quadrant and becomes lighter towards the top left.

What will make C++-next portable?

TOPOLOGY

CPU

DRAM

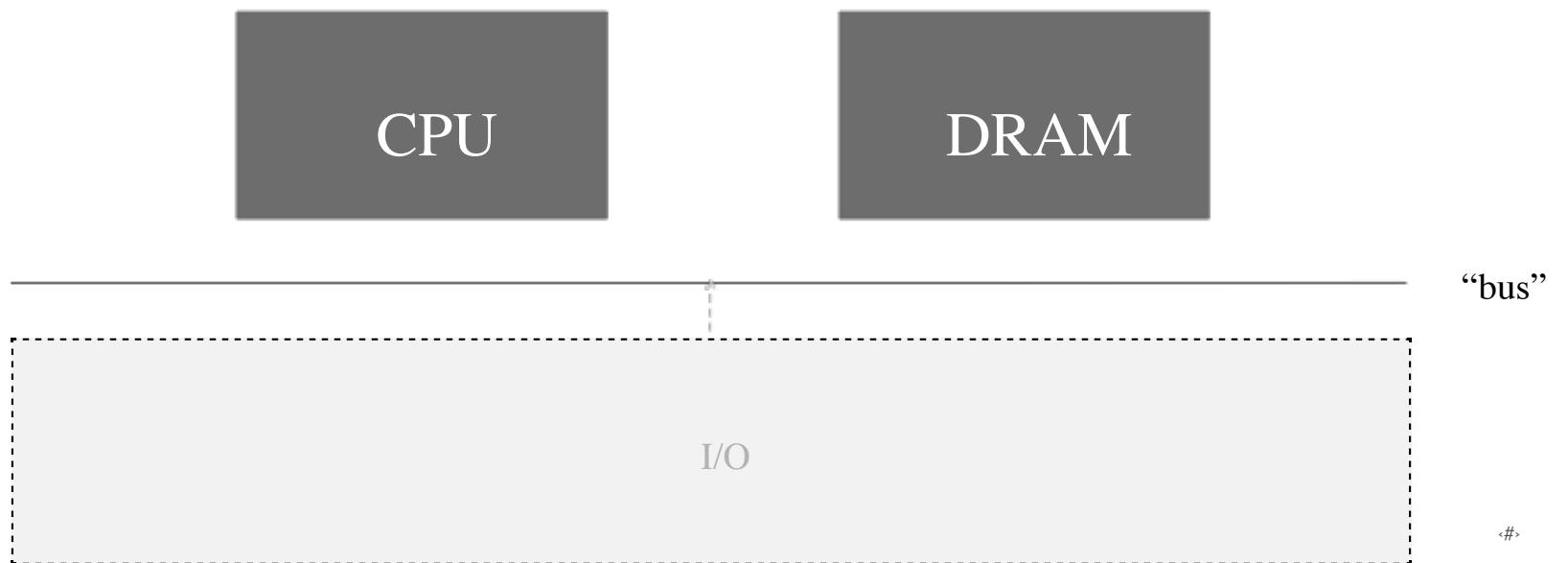
“bus”

1980'S TOPOLOGY

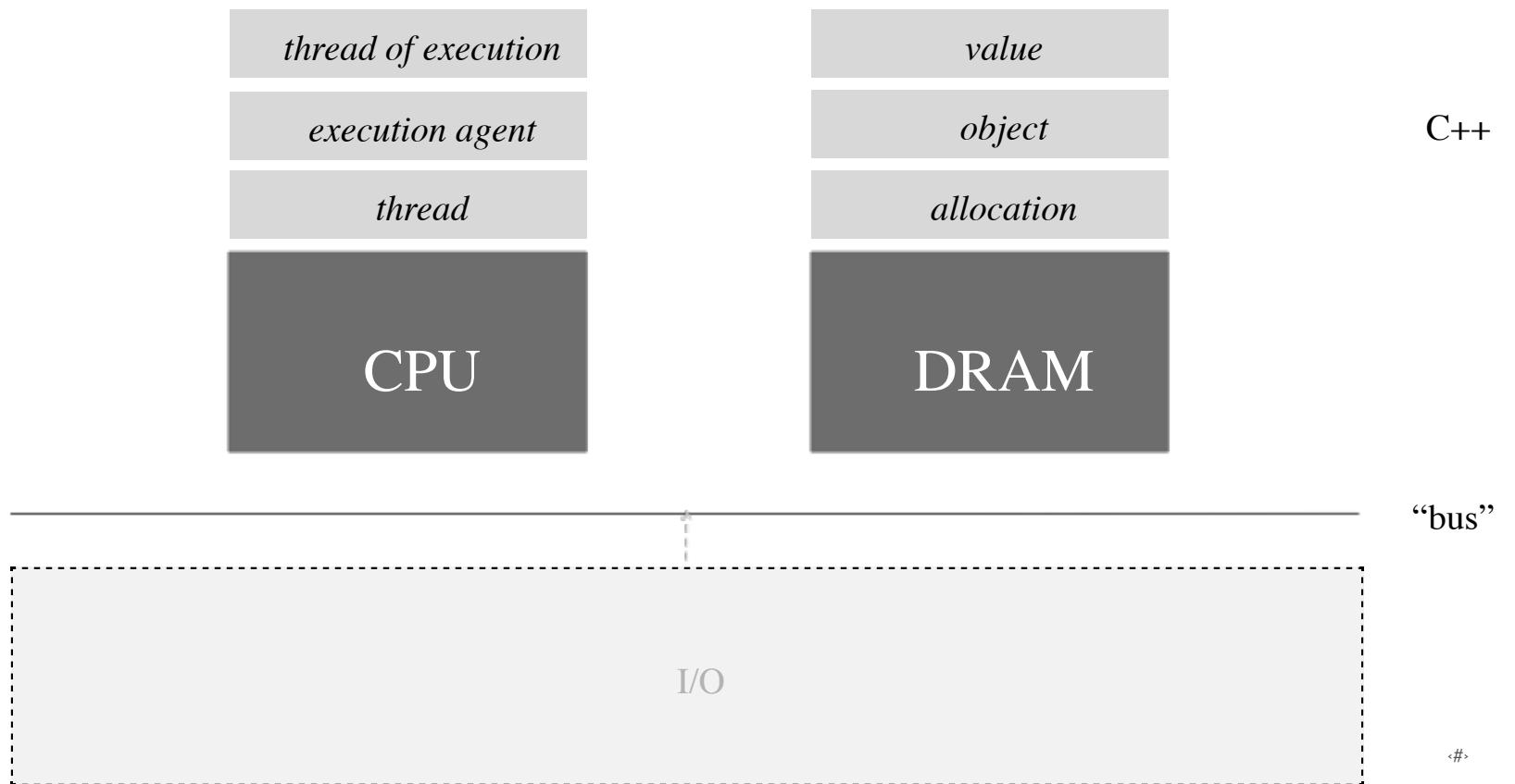


“bus”

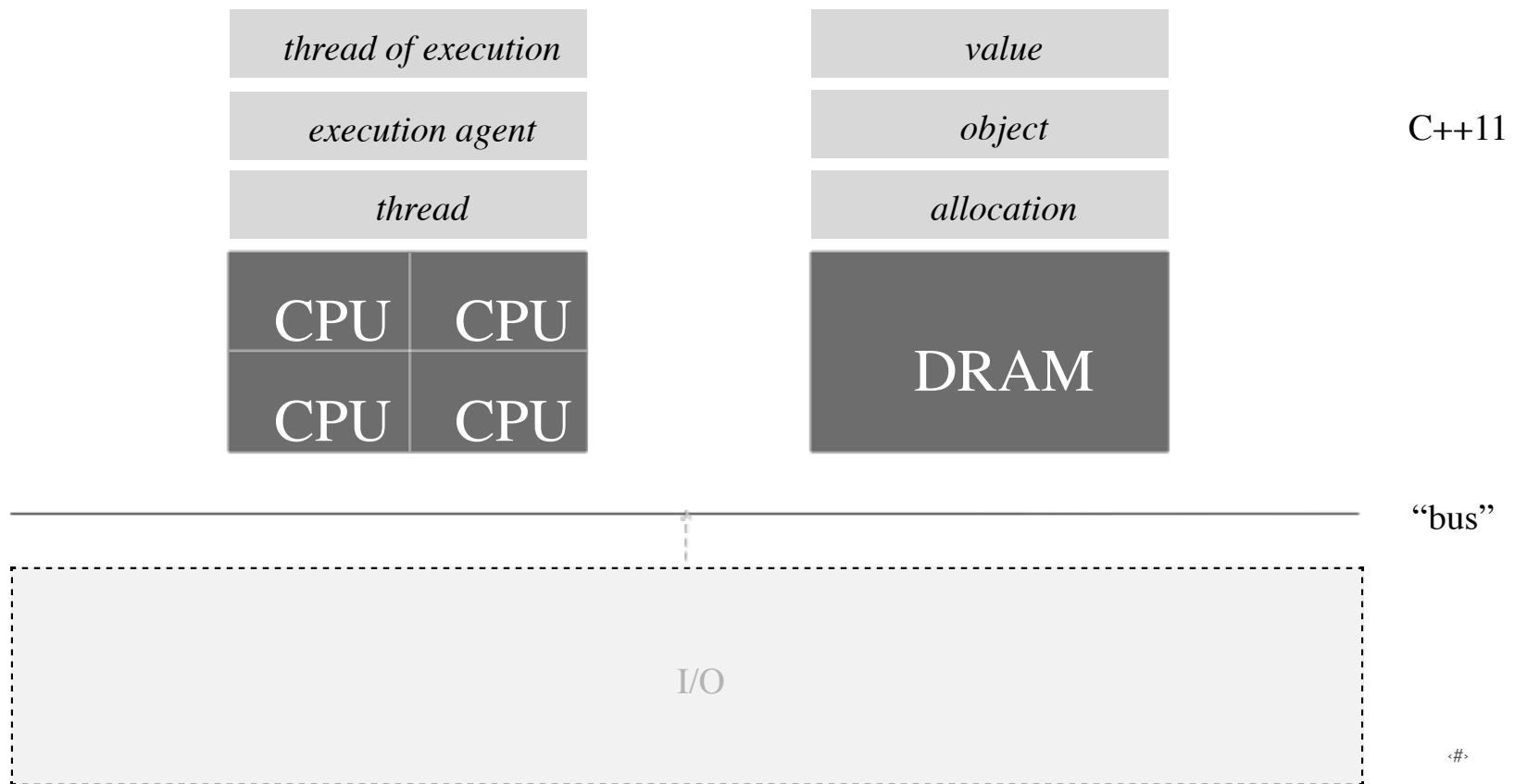
1980'S TOPOLOGY



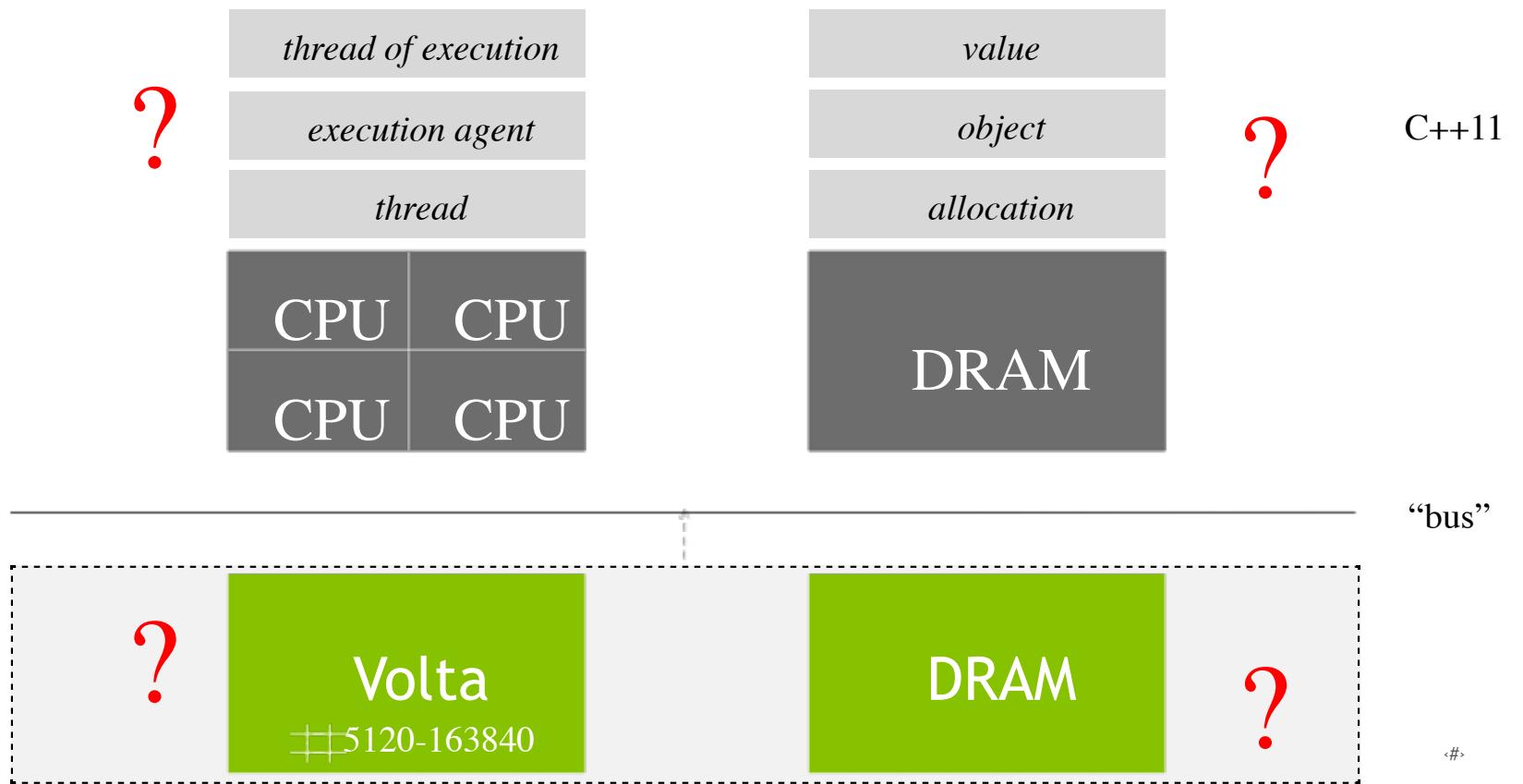
1980'S TOPOLOGY



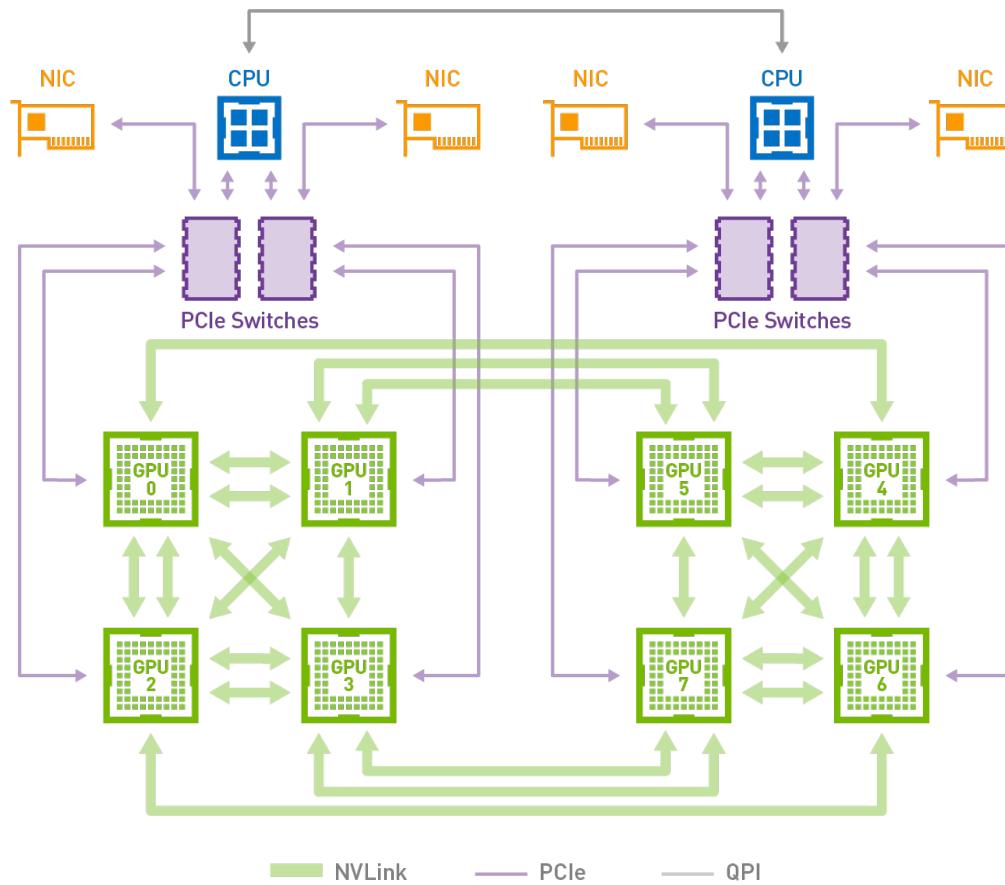
1990-2000'S TOPOLOGY



2010'S TOPOLOGY



SHOW ❤ FOR TOPOLOGY

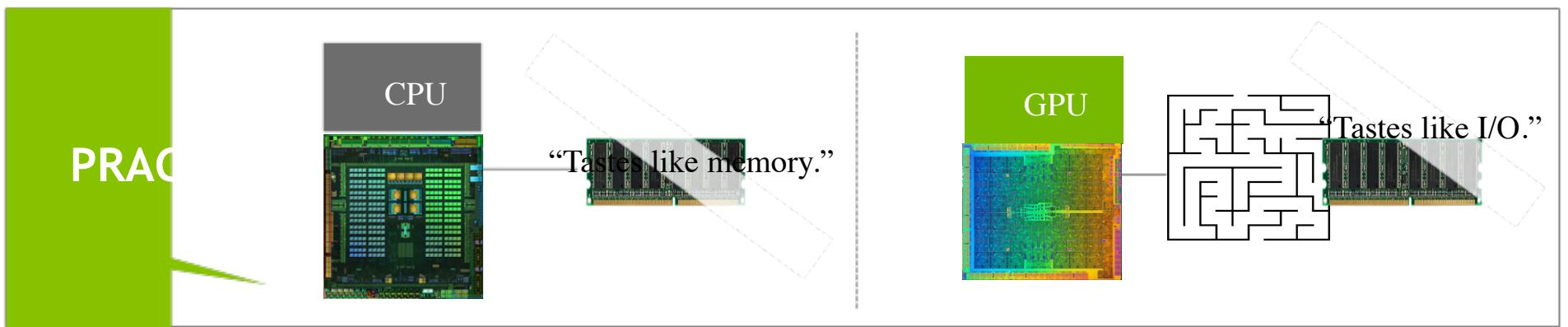


Designing Volta for C++ memory

COHERENCY



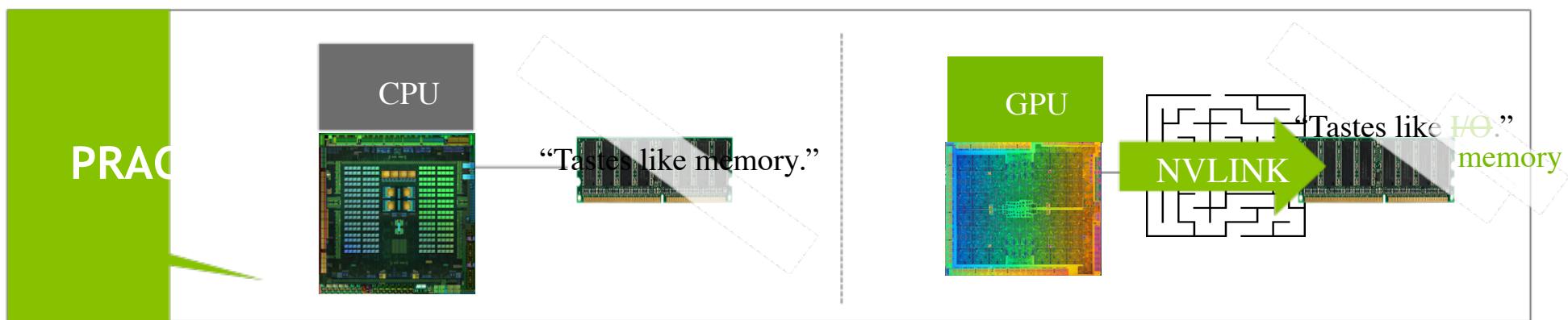
COHERENCY



COHERENCY

ARCHITECTURE	CUDA	X86	ARM & POWER
Tesla & Fermi	1+	cudaMalloc & cudaMemcpy	
Kepler & Maxwell	6+	cudaMallocManaged + symmetric heap	
Pascal & Volta “Tastes like memory.”	8+	cudaMallocManaged + paging	cudaMallocManaged (NVLINK)
	Linux HMM	malloc + paging	malloc (NVLINK)

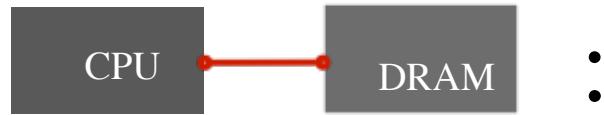
COHERENCY



COHERENCY

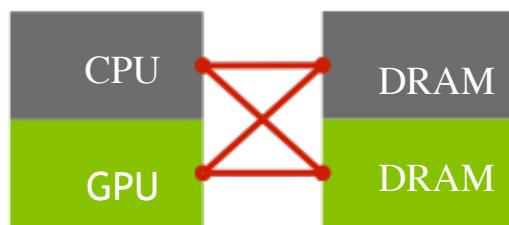
ARCHITECTURE	CUDA	X86	ARM & POWER
Tesla & Fermi	1+	cudaMalloc & cudaMemcpy	
Kepler & Maxwell	6+	cudaMallocManaged + symmetric heap	
Pascal & Volta “Tastes like memory.”	8+	cudaMallocManaged + paging	cudaMallocManaged + NVLINK
	Linux HMM	malloc + paging	malloc + NVLINK

CONSISTENCY FOR CPUS



C++11	load	store	exchange	fence
(not atomic)	✓	✓	-	-
relaxed	✓	✓	✓	-
consume	▼	-	▼	-
acquire	✓	-	✓	✓
release	-	✓	✓	✓
acq_rel	-	-	✓	✓
seq_cst	✓	✓	✓	✓

CONSISTENCY FOR VOLTA = SAME



•
•

C++11	load	store	exchange	fence
(not atomic)	✓	✓	-	-
relaxed	✓	✓	✓	-
consume	✗	-	✗	-
acquire	✓	-	✓	✓
release	-	✓	✓	✓
acq_rel	-	-	✓	✓
seq_cst	✓	✓	✓	✓

COMPLETELY NEW memory model for Volta. Outline similar to POWER.

See the PTX 6.0 ISA programming guide, chapter 8.

WE ❤ THE MEMORY MODEL IN C++11

Good example of “what makes C++ portable”.

Not: `volatile`, was tyranny for us architects.

Open issues for C++ memory

FOOTPRINT

16
CPU

128GB
DRA

163840
Volta

16GB
DRA

FOOTPRINT

16 CPU : 128GB DRA = 1 : 8GB

163840 Volta : 16GB DRA = 1 : 102KB

FOOTPRINT

$$163840 \text{ Volta} : 16\text{GB DRA} = 1 : 102\text{KB}$$


Shall we spend all this on `thread_local`? (y/N)

Whoever invented TLS didn't think this far ahead.

FOOTPRINT



MEMORY PROGRESS REPORT

PROBLEM	GPU DISPOSITION	C++ DISPOSITION
Coherency	Pascal-Volta.	
Consistency	Volta, for assembly. (CUDA C++ TBD.)	P0668.
Footprint	-	P0072.

Designing Volta for C++ execution

PROGRESS

thread of execution

= A chain of evaluations in your code.

execution agent

= A thing that runs your code.

thread

= A particularly onerous example of that thing.

CPU	CPU
CPU	CPU

GPU

PROGRESS

thread of execution

= A chain of evaluations in your code.

execution agent

= A thing that runs your code.

thread

= A particularly onerous example of that thing.

CPU	CPU
CPU	CPU

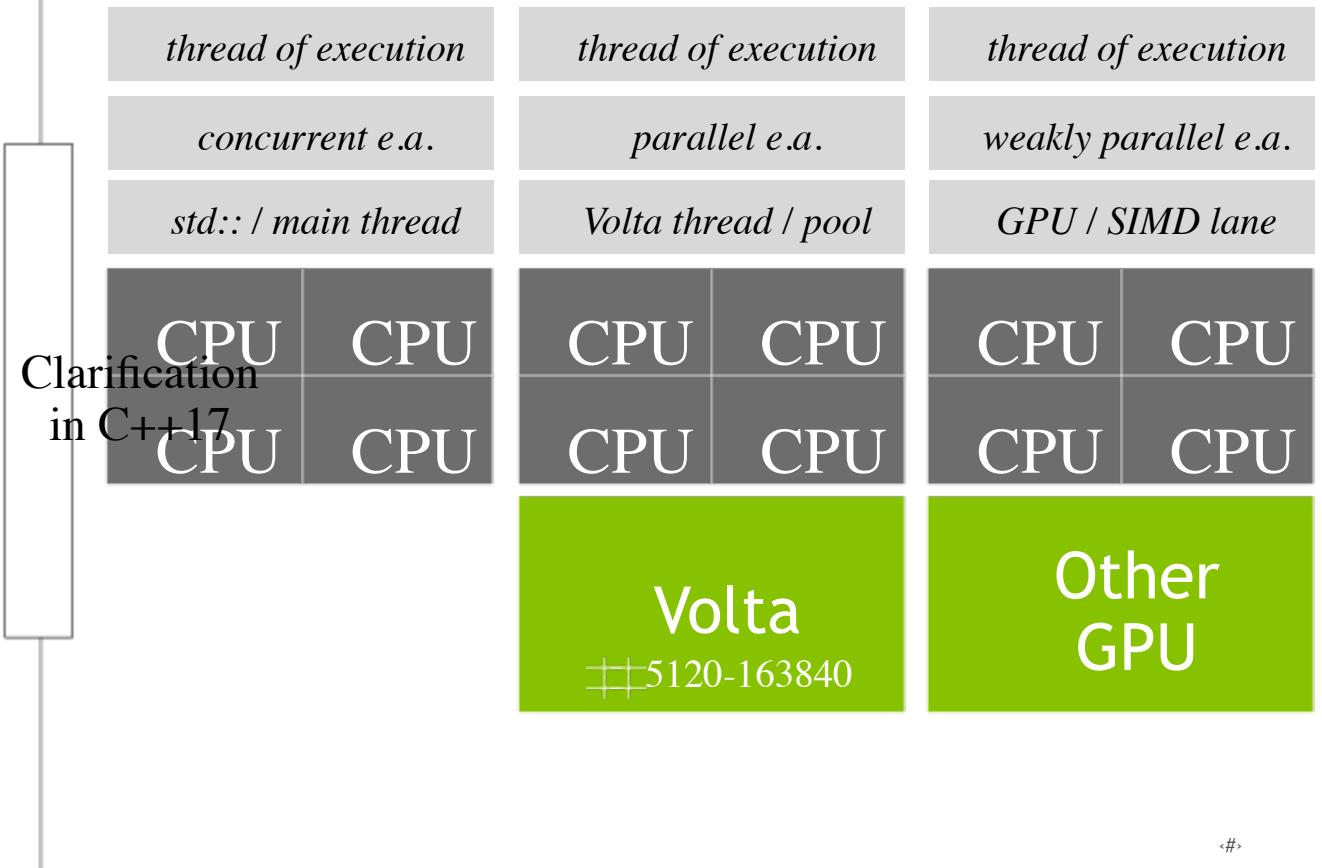
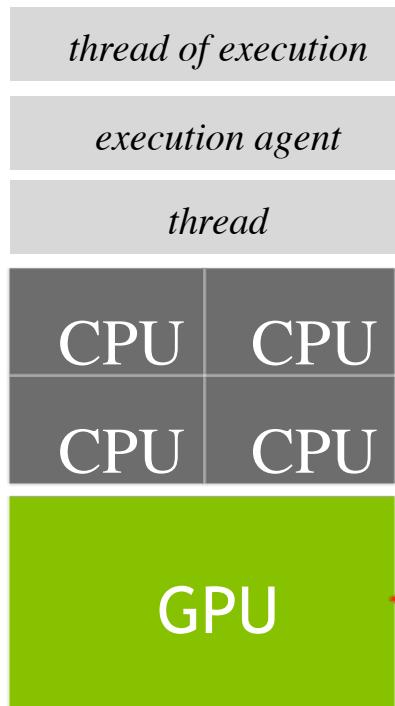
= Runs anything.

GPU

= Runs things that aren't onerous.



PROGRESS

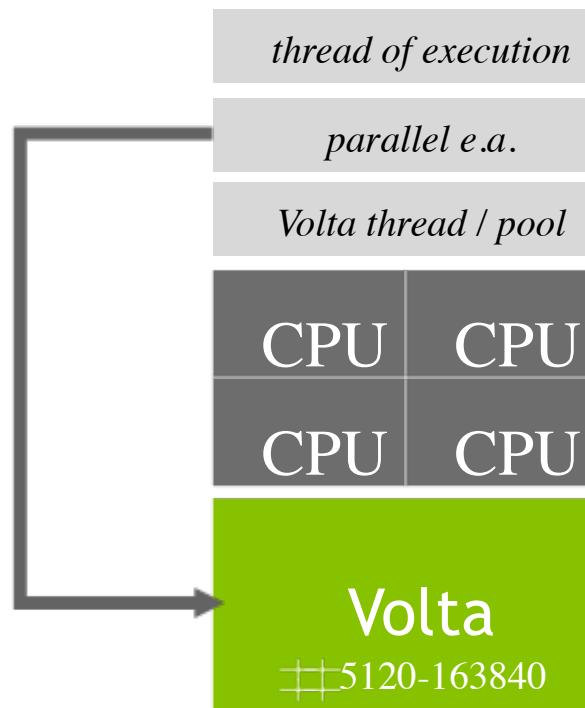


PROGRESS

Not “business as usual”.

A concerted effort by
dedicated engineers.

Volta is alone of its kind.



CLARIFYING SIMT C++

Sorry, CUDA `atomic<T>` isn't ready yet.

CLARIFYING SIMT C++

Sorry, CUDA `atomic<T>` isn't ready yet.

```
enum { unlocked = 0, locked = 1 };

//atomic<int> mutex = ATOMIC_VAR_INIT(unlocked);
volatile int mutex = unlocked;

void demo() {

    //int old;
    //while(!mutex.compare_exchange_strong(old=unlocked, locked))
    while(atomicCAS(&mutex,unlocked,locked)==locked)
        /* some kind of backoff here */
    ;

    /* Among throughput processors, only Volta is
       guaranteed to run this critical section. */
}
```

Open issues for C++ execution

STATE OF HETEROGENEITY

host device

Short version: nothing to report today.

This (and lack of compiler integration) is behind most conformance issues.

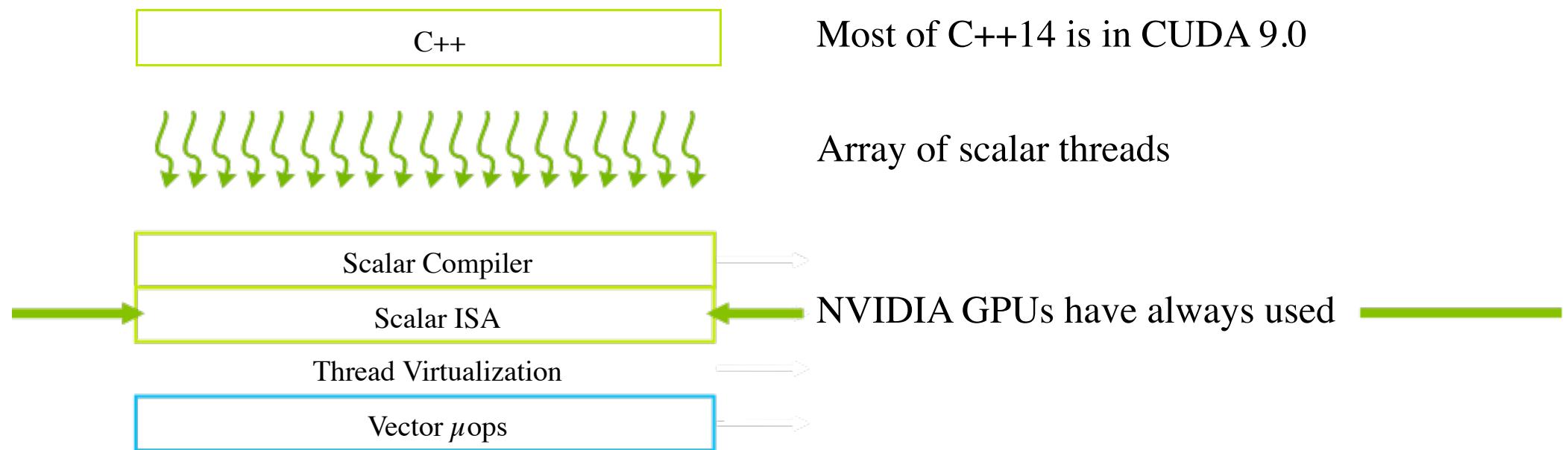
e.g.: no `std::vector::operator[]` or anything from `std::` because of this.

EXECUTION PROGRESS REPORT

PROBLEM	GPU DISPOSITION	C++ DISPOSITION
Progress	Clarified in Volta.	Clarified in C++17. 
SIMT C++	Clarified in Volta.	-
Heterogeneity		TBD

VOLTA C++ ARCHITECTURE WRAP-UP

VOLTA: MANY SCALAR C++ THREADS





QUESTIONS.