# Sparse & irregular processing on GPUs using functional languages
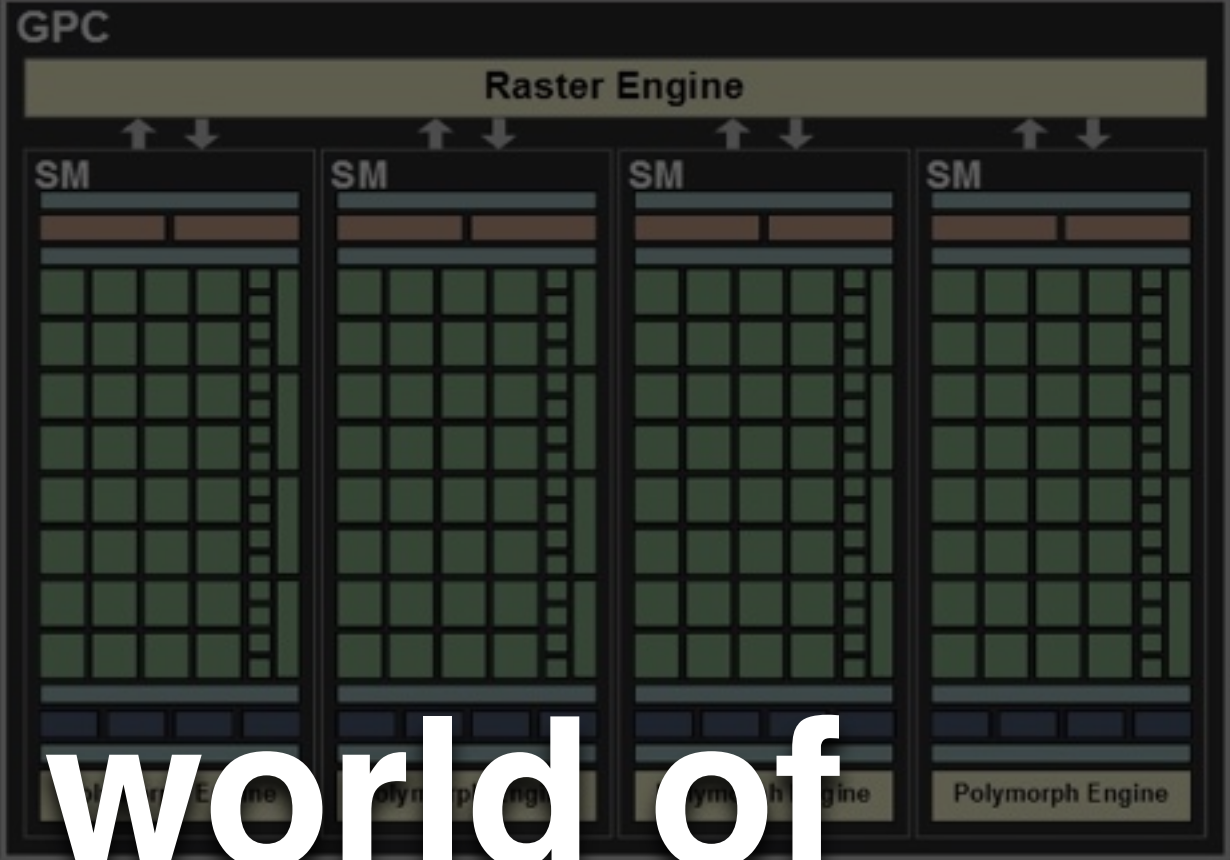
**Adam Harries**
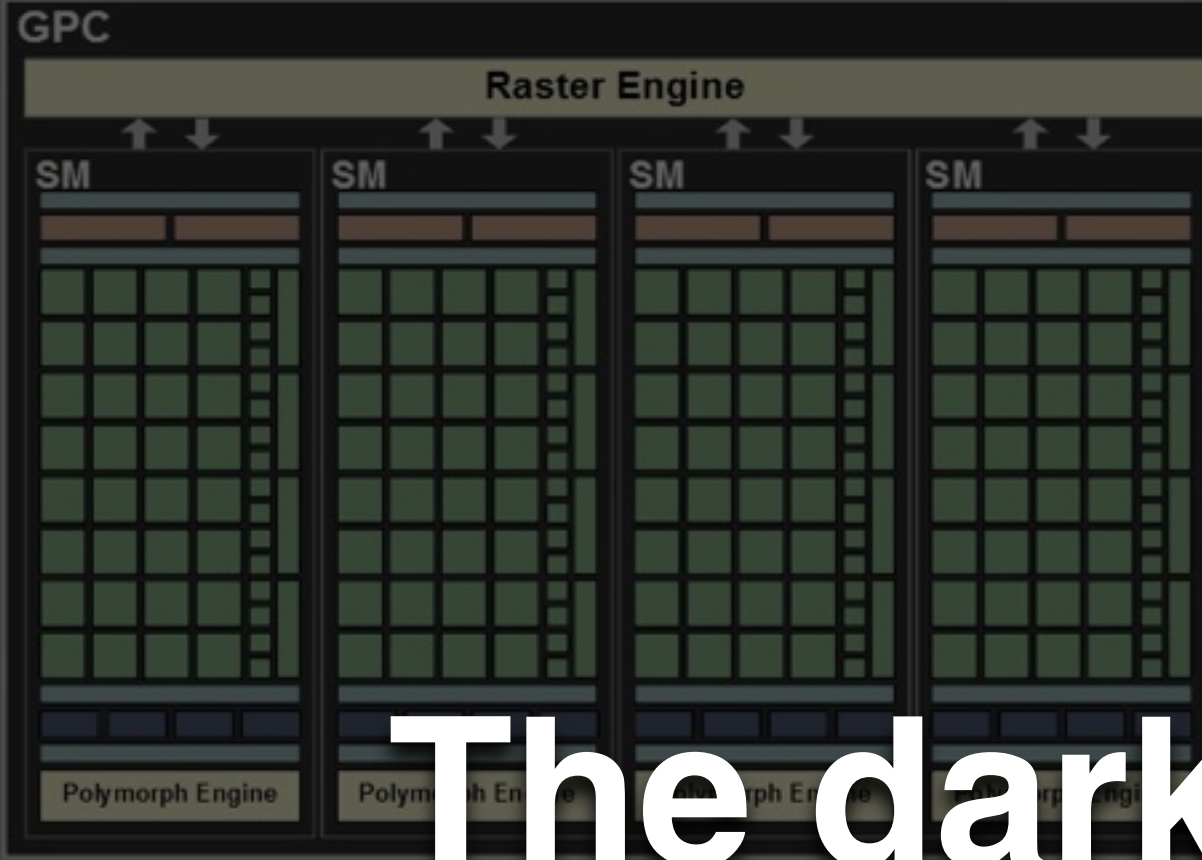
School of Informatics, University of Edinburgh
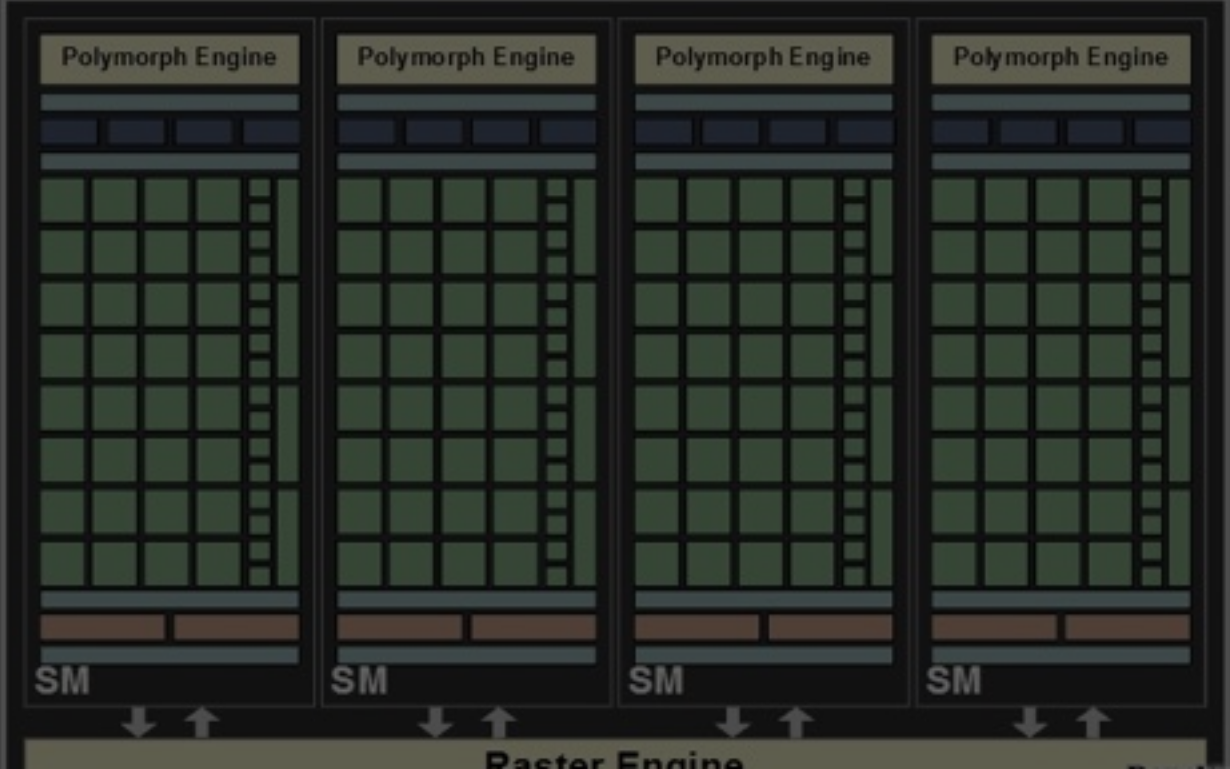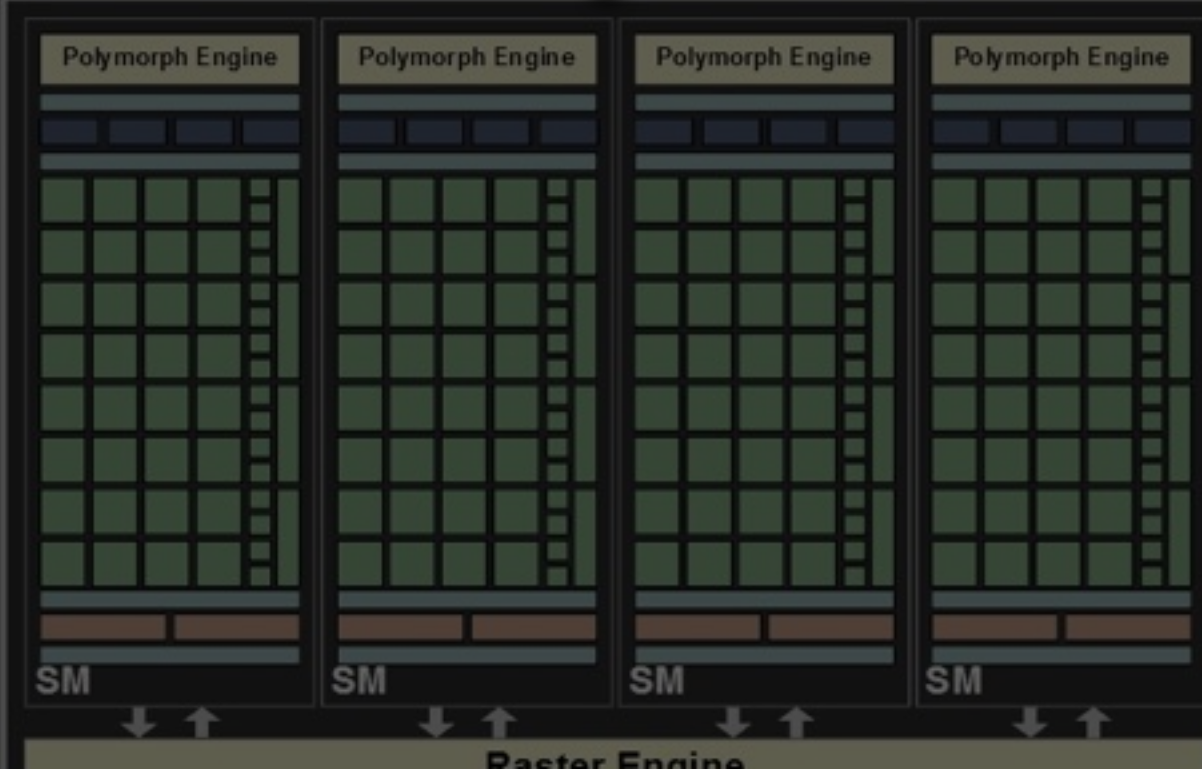
THE UNIVERSITY of EDINBURGH
**informatics**

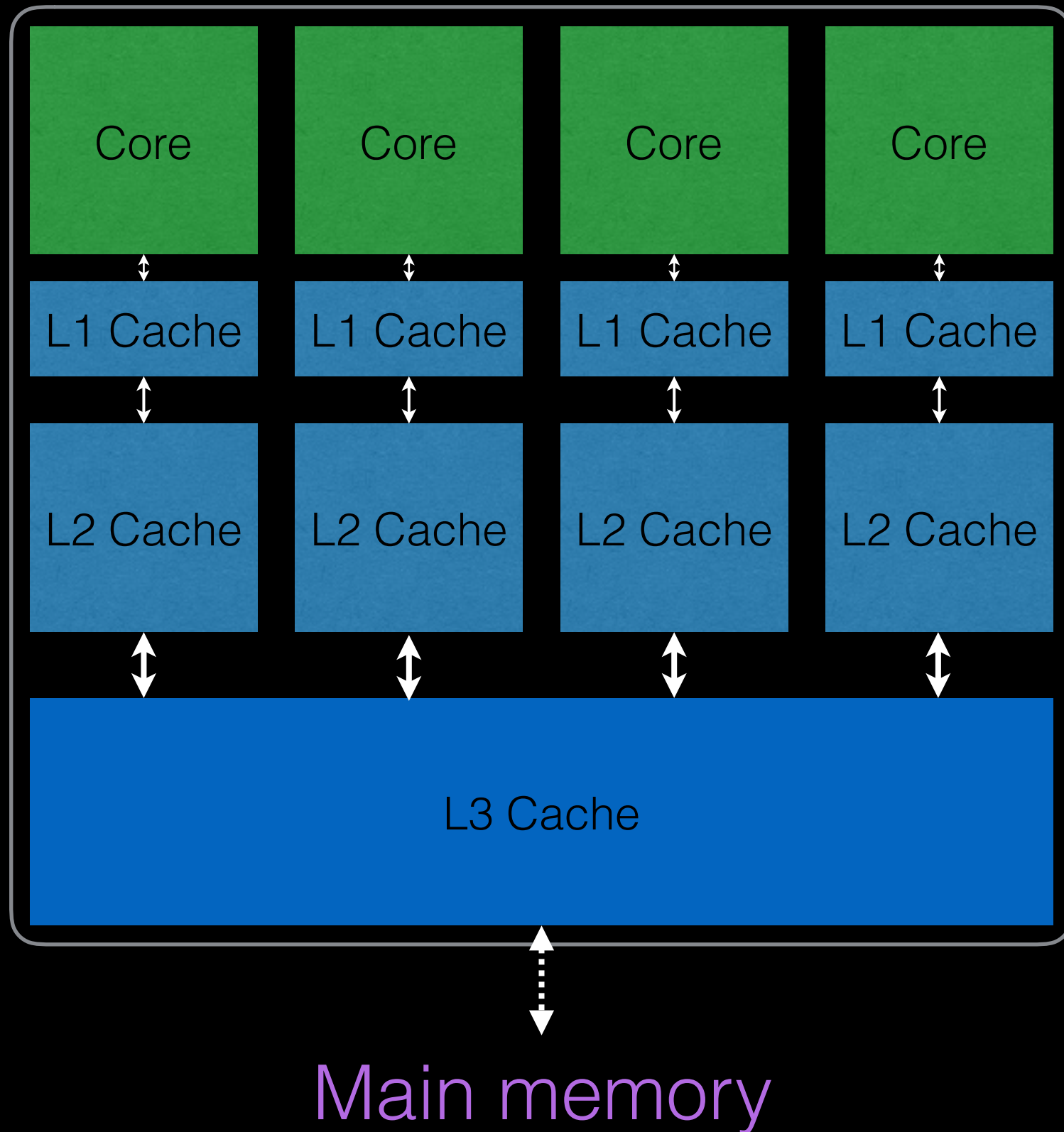EPSRC **Centre for Doctoral Training in**
**Pervasive Parallelism**

EPSRC
Engineering and Physical Sciences
Research Council

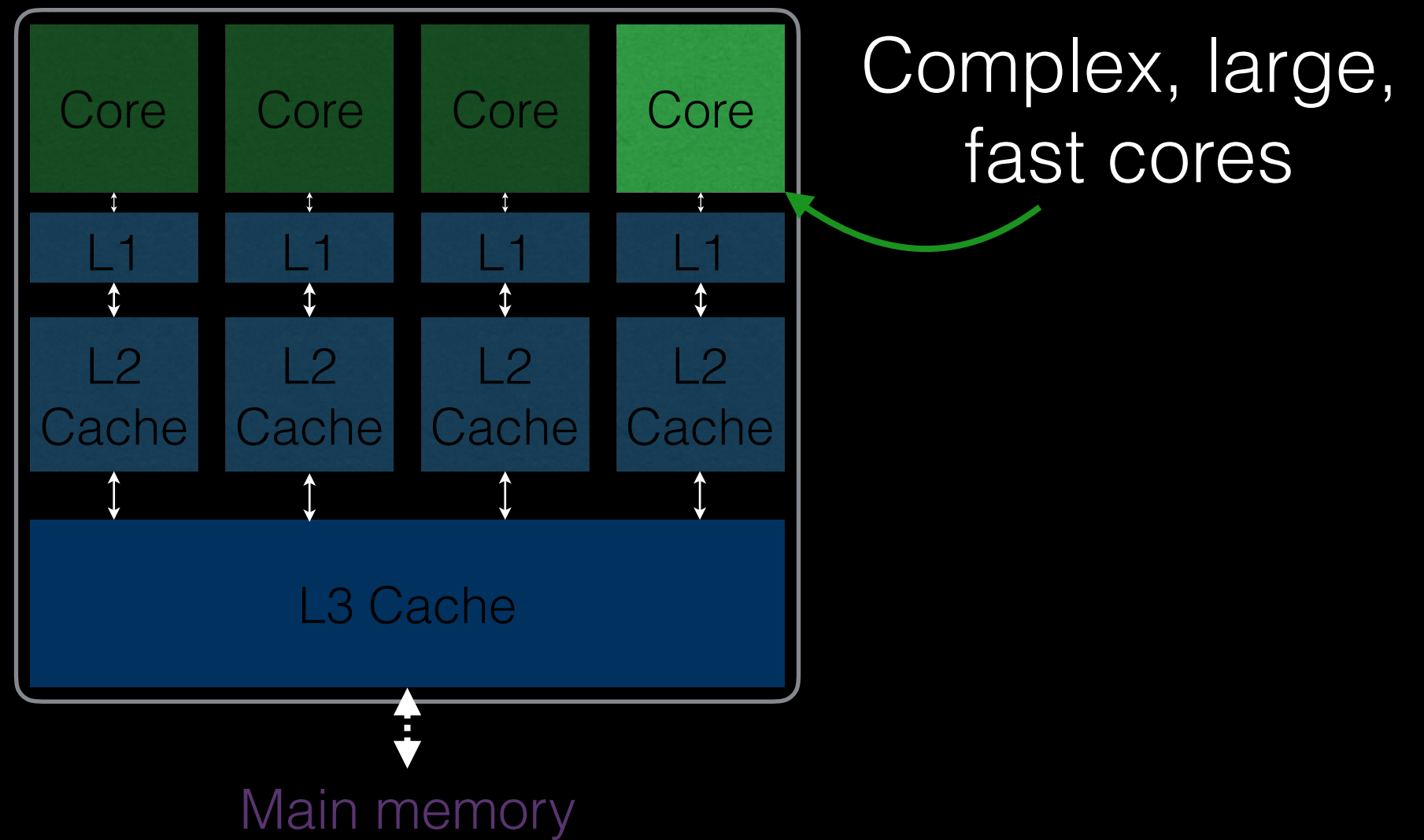# The dark world of computer architecture

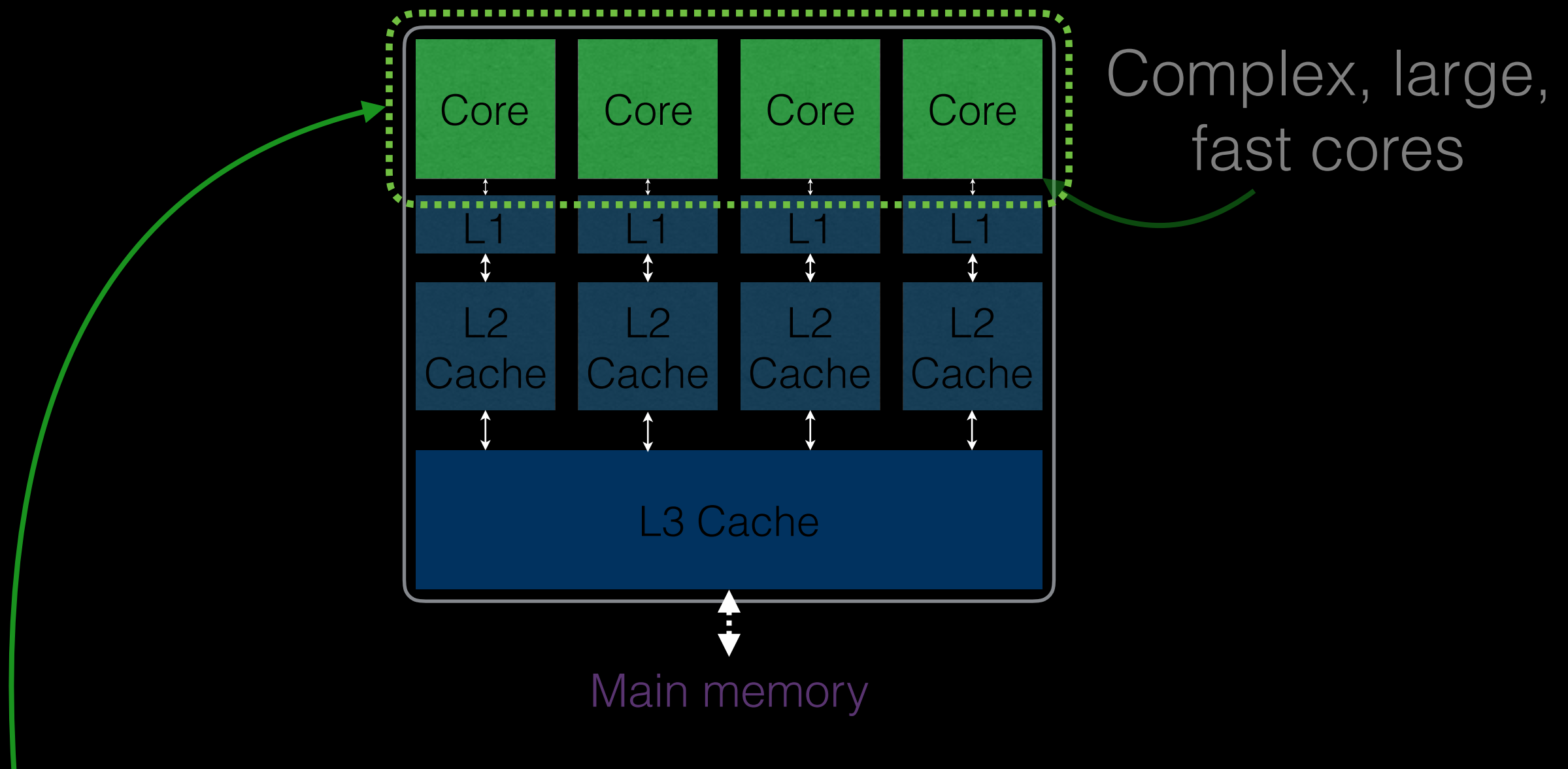# Current CPUs

# Current CPUs

| | | | |
|---|---|---|---|
| Core | Core | Core | Core |
| L1 | L1 | L1 | L1 |
| L2 Cache | L2 Cache | L2 Cache | L2 Cache |
| L3 Cache | | | |

Complex, large, fast cores

Main memory

# Current CPUs

| Core | Core | Core | Core |

L1   L1   L1   L1

| L2 Cache | L2 Cache | L2 Cache | L2 Cache |

L3 Cache

Main memory

Complex, large, fast cores
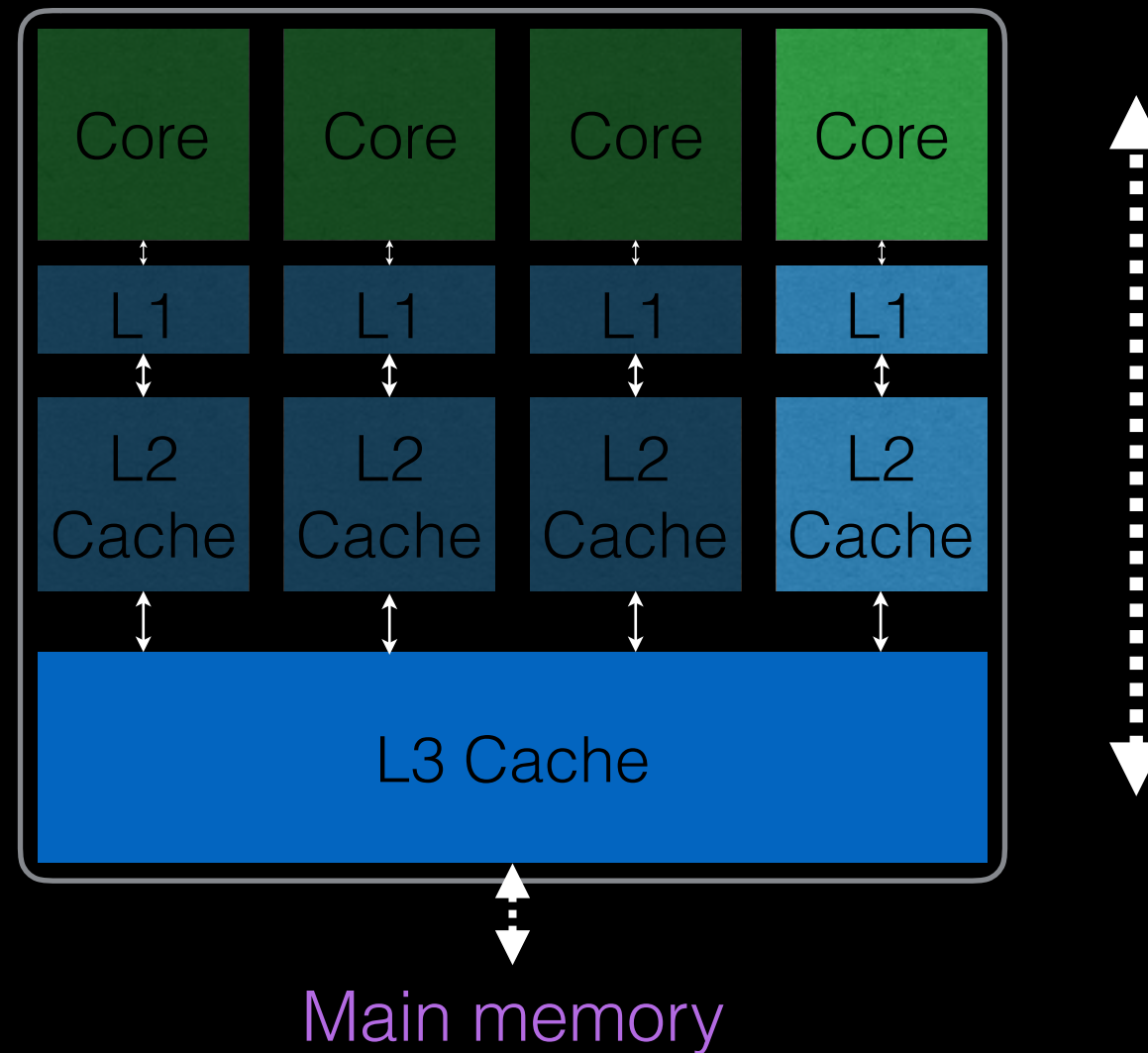
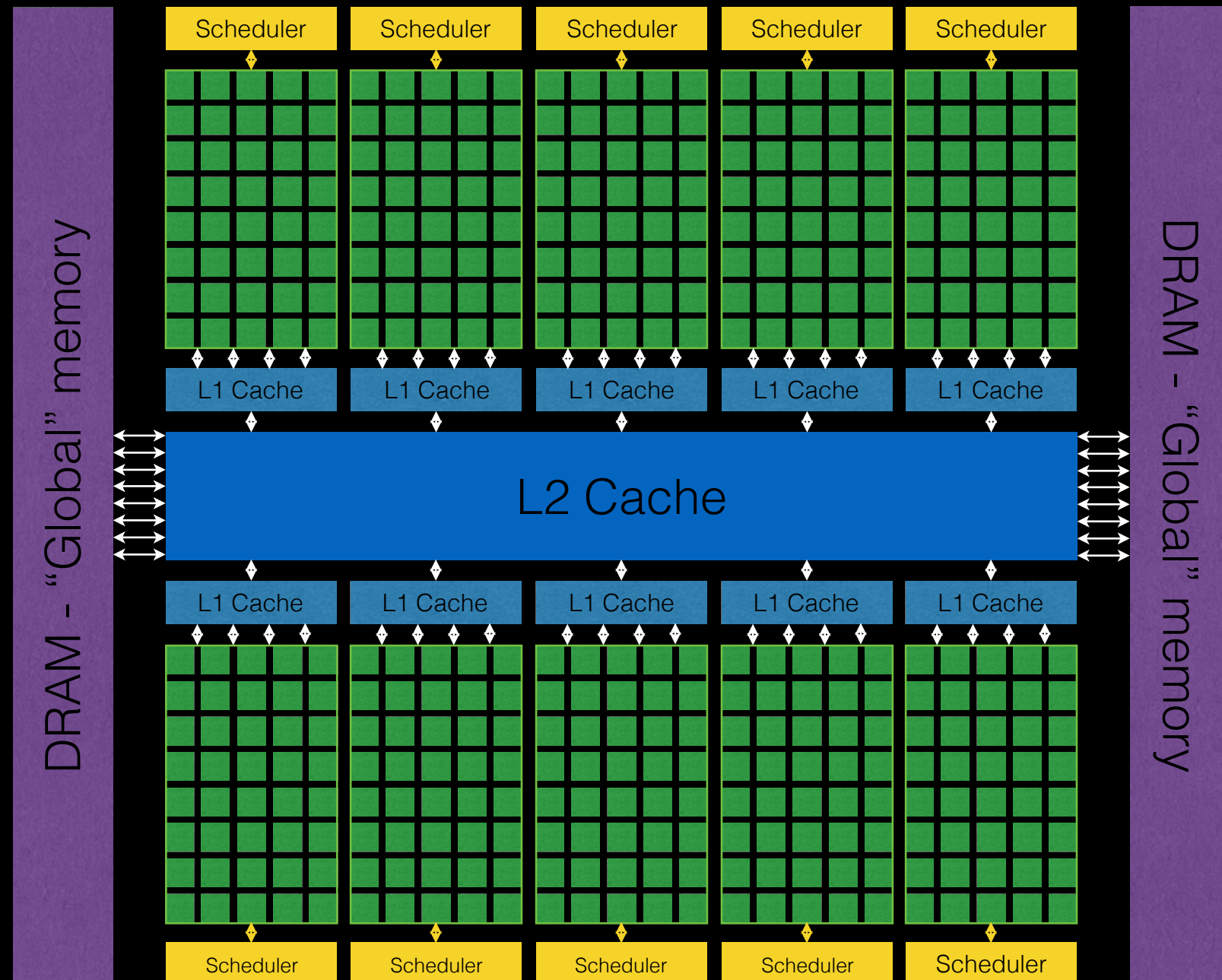Coarse parallelism: **low thread count** and **heavyweight threads**

# Current CPUs



Context switching between threads is **expensive**:
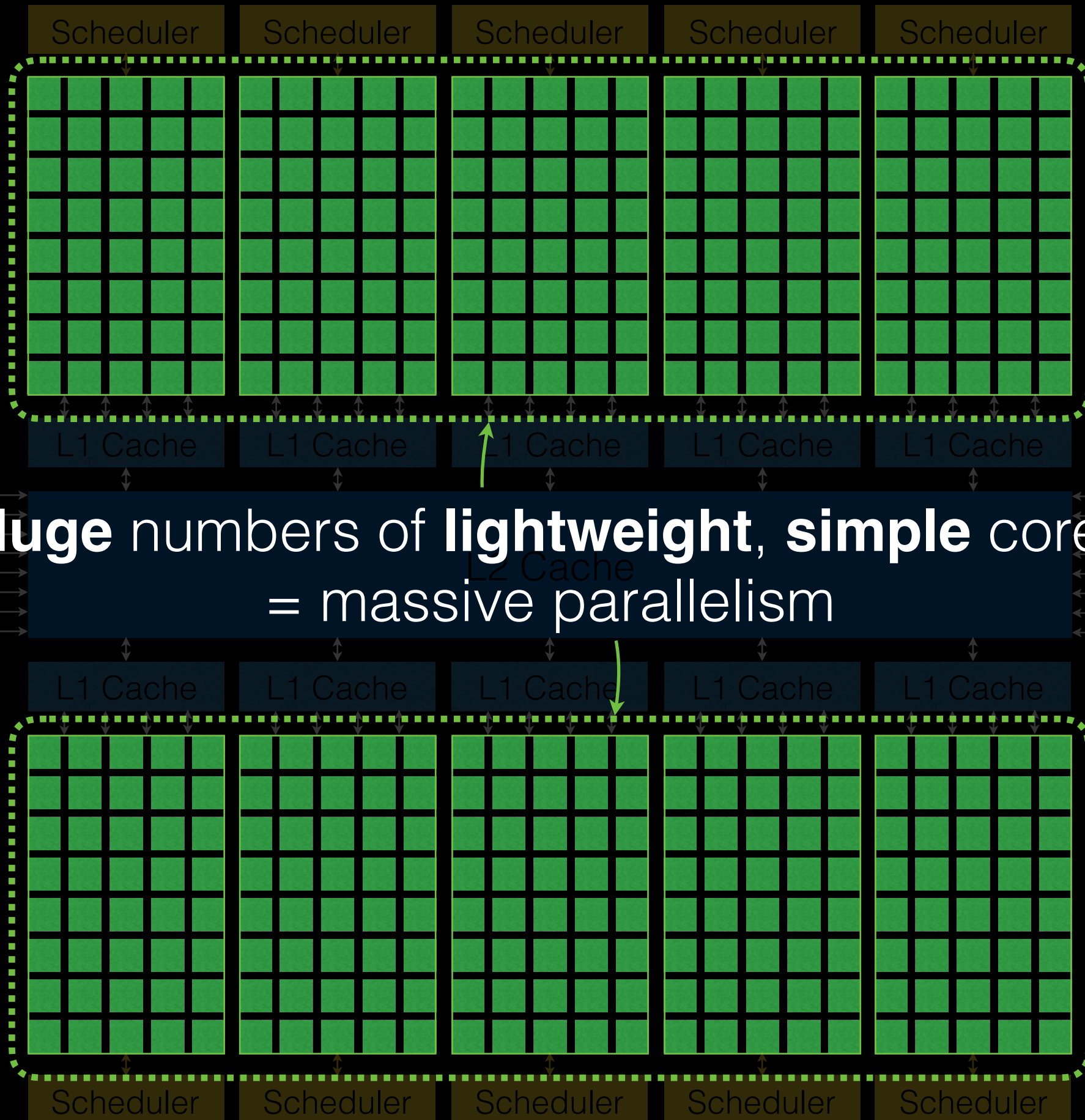architecture optimised for single threaded performance

# Current GPUs

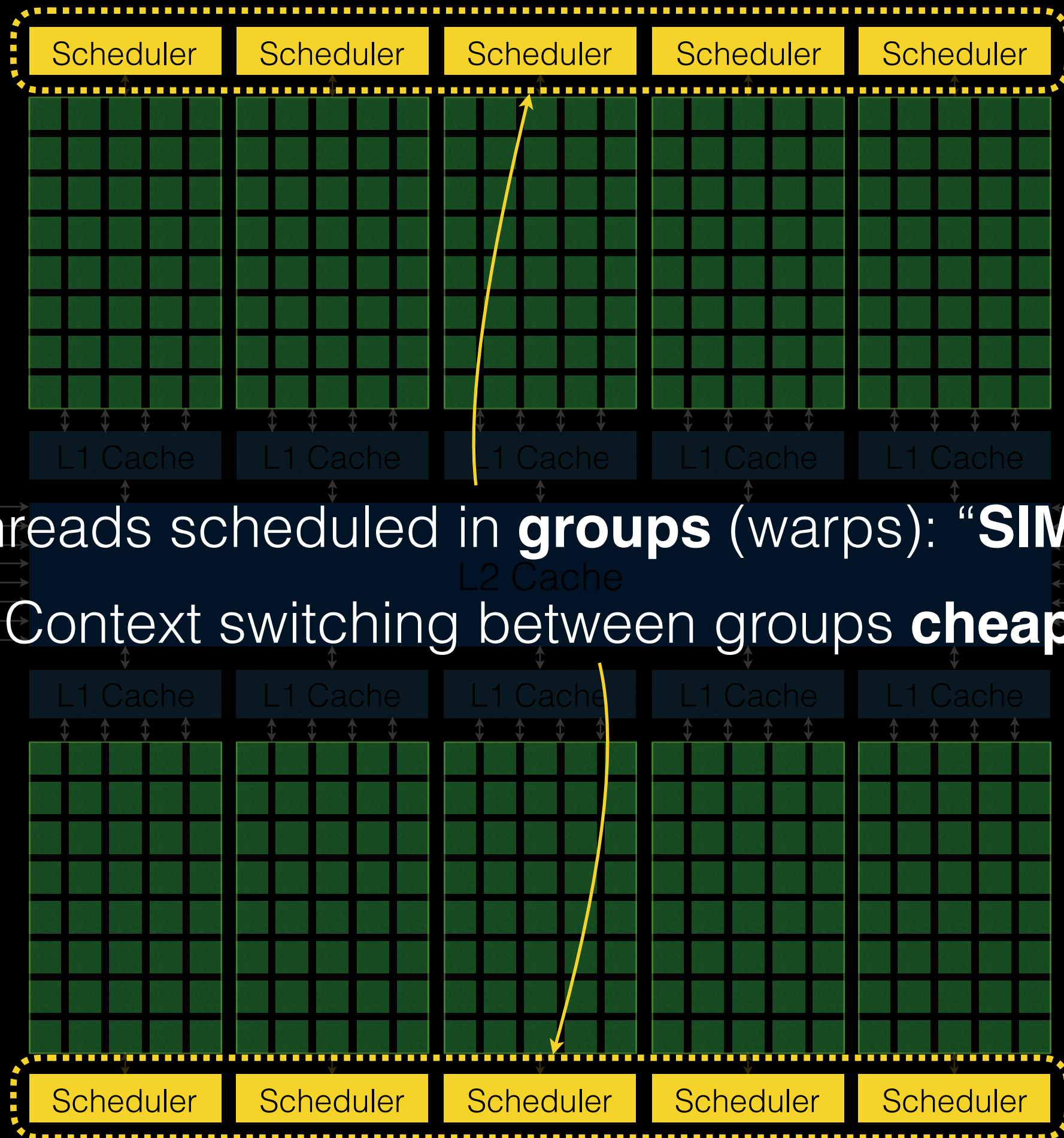**Conclusion:** GPUs give you *greater*, more *flexible* parallelism

**Conclusion:** GPUs give you *greater*, more *flexible* parallelism…right?

**Conclusion:** GPUs give you *greater*, more *flexible* parallelism…right?

**Why don't we use them for everything?**

**No caches:** Programmers must be careful with intelligently fetching data

**Grouped Parallelism:** Programmers should avoid thread divergence within groups

**Simple cores:** Programmers must parallelise with fine grained *and statically* distributed workloads

**No caches:** Programmers must be careful with intelligently fetching data

**Grouped Parallelism:** Programmers should avoid thread divergence within groups

**Simple cores:** Programmers must parallelise with fine grained *and statically* distributed workloads

GPUs encourage **static** *array based* parallelism

**No caches:** Programmers must be careful with intelligently fetching data

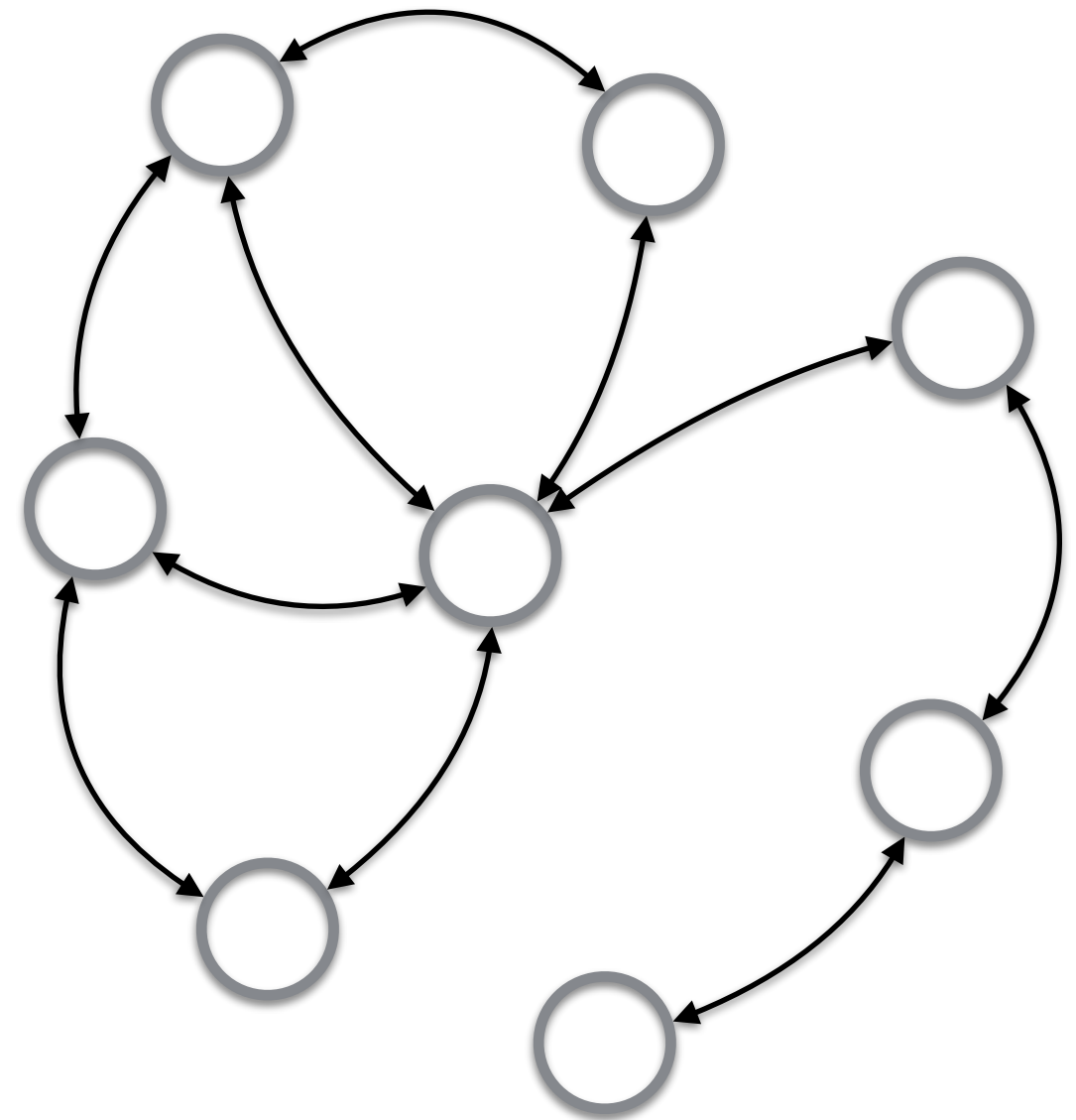**Grouped Parallelism:** Programmers should avoid thread divergence within groups

**Simple cores:** Programmers must parallelise with fine grained *and statically* distributed workloads

GPUs encourage **static** *array based* parallelism
How can we **efficiently** parallelise
*other domains*?

Domain of interest:
**Graph algorithms**
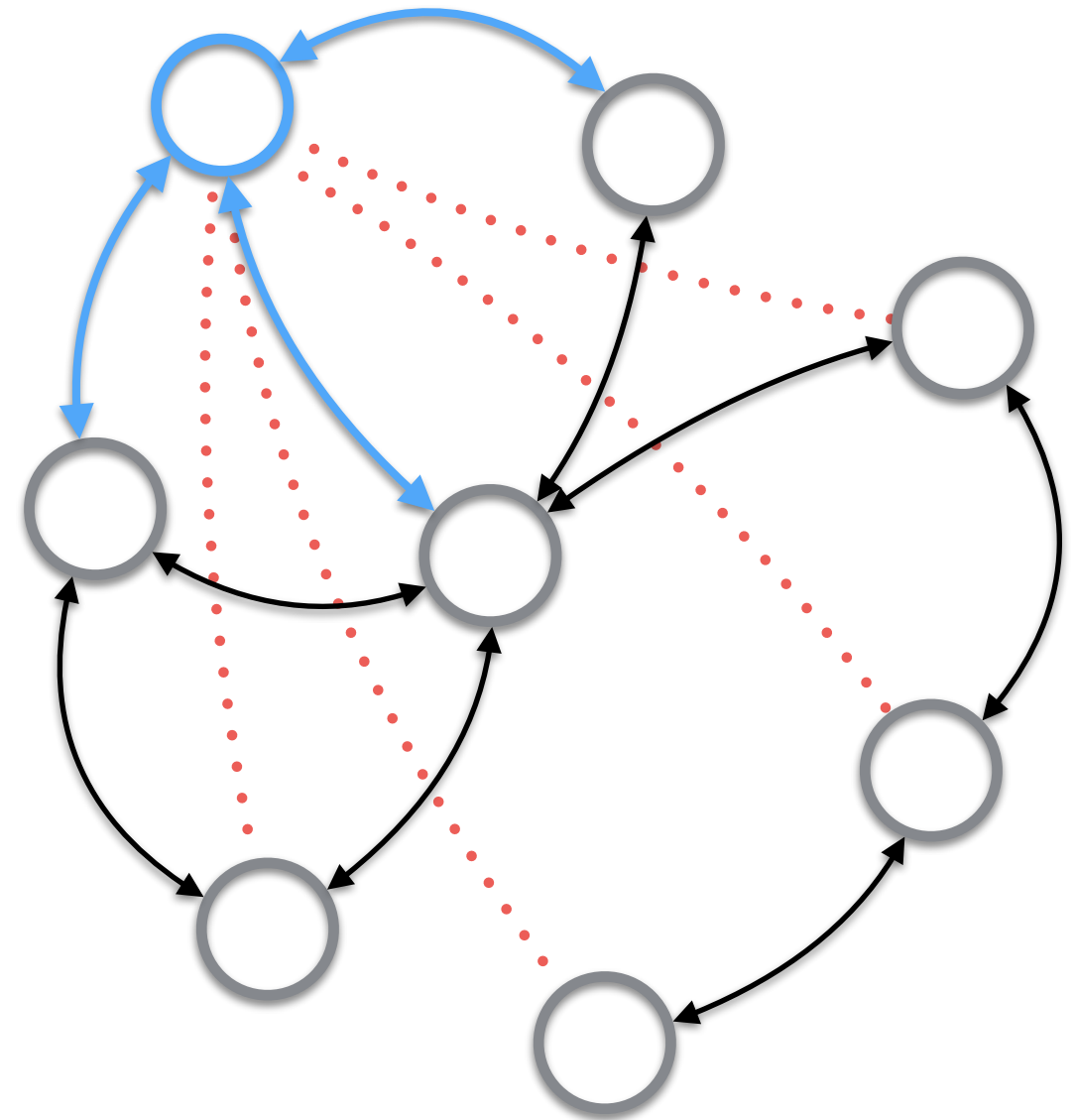
Aim: accelerate
processing using
**GPU parallelism**

What are the
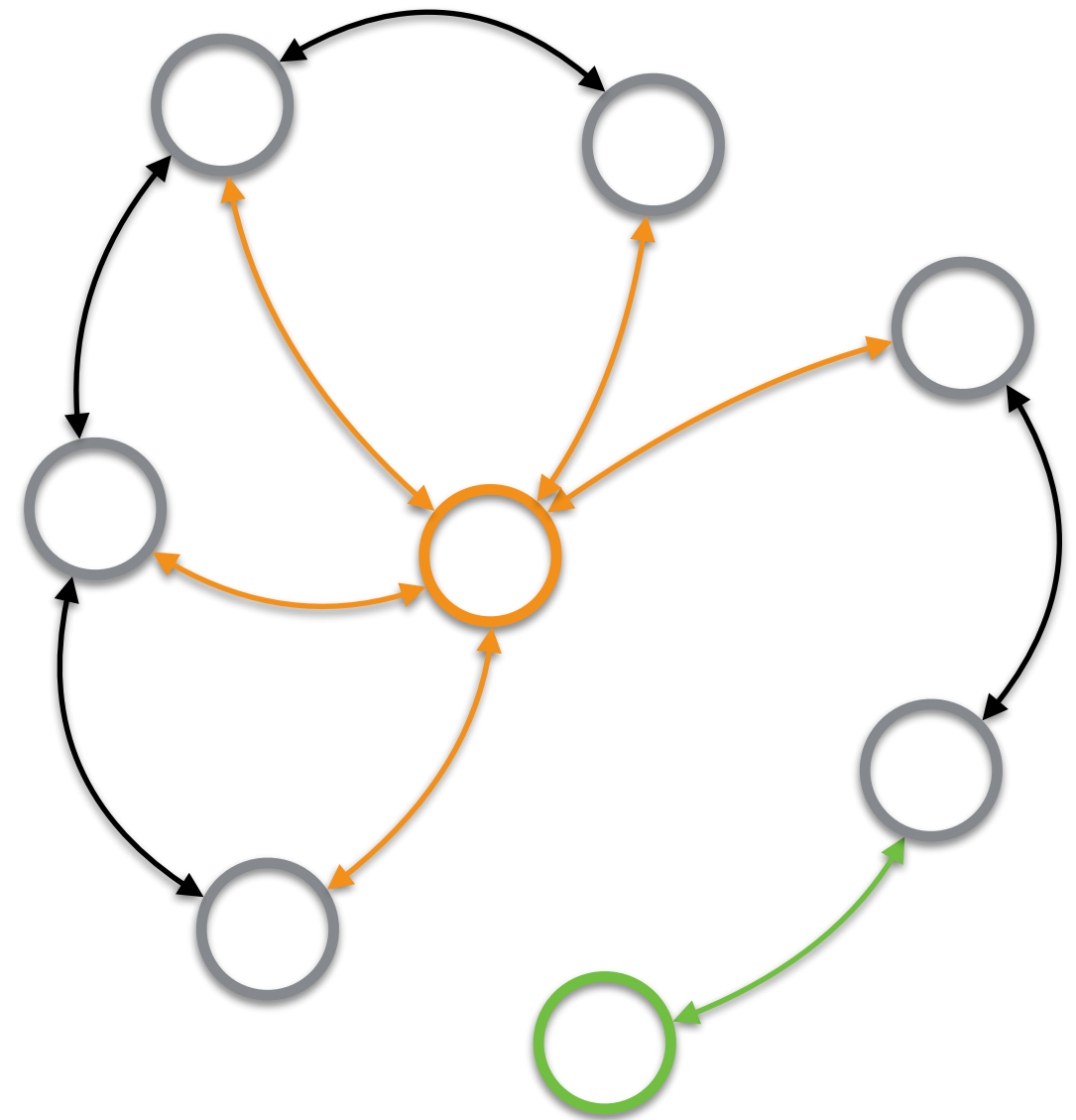**problems**?

# What are the **problems**?

**Sparsity:**
- Complex data structures
- Random data access patterns
- Low compute to data ratio

# What are the **problems**?

**Irregularity:**
- Unbalanced workload within data sets
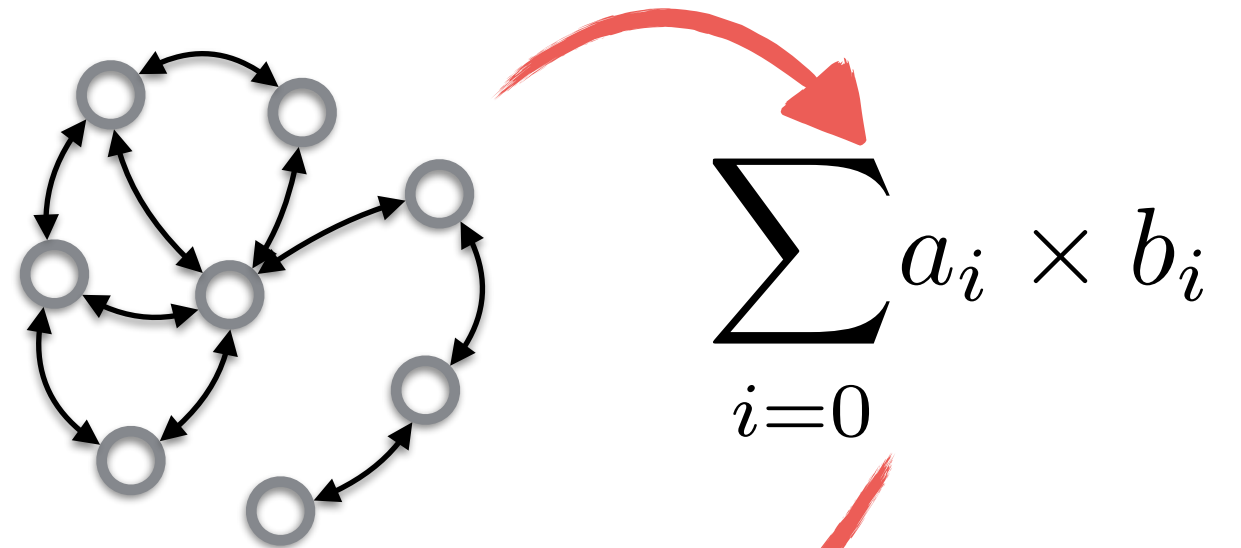- Control flow divergence
- Coarse granularity of parallelism

# Our proposal:

Translate problem to simpler domain: (sparse) **Linear Algebra**
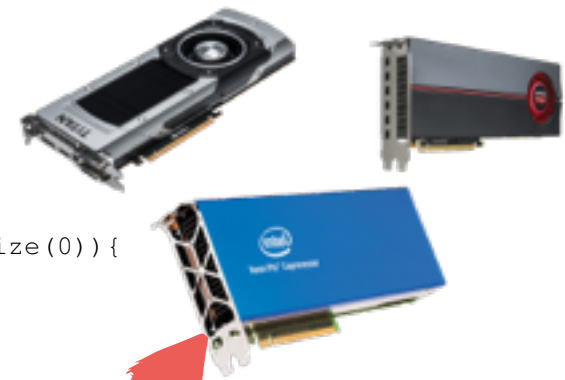
$$\sum_{i=0} a_i \times b_i$$

Express algorithm with **high level parallelism** (aka, functional languages)

Reduce(add, 0) ○ Map(mult)

**Explore** space of implementations to find **high performance** solution
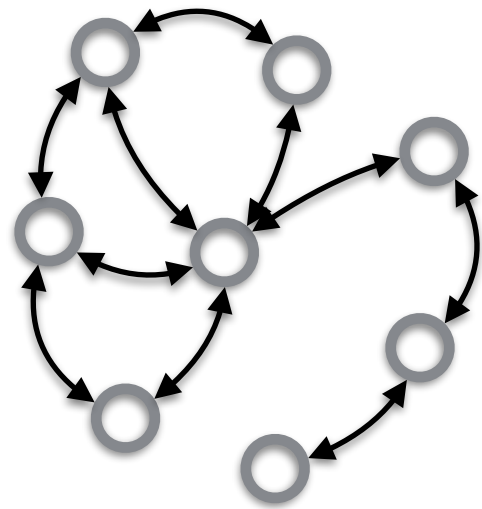
```
__kernel void dotProduct(
    __global const elem_t* a,
    __global const elem_t* b,
    __global const elem_t* res,
    __global const int len){
  int id = get_global_id(0);
  elem_t tmp;
  for(int i = id;i<len;i+=get_global_size(0)){
    tmp = a[i]*b[i];
    atomic_add(res, tmp);
  }
}
```
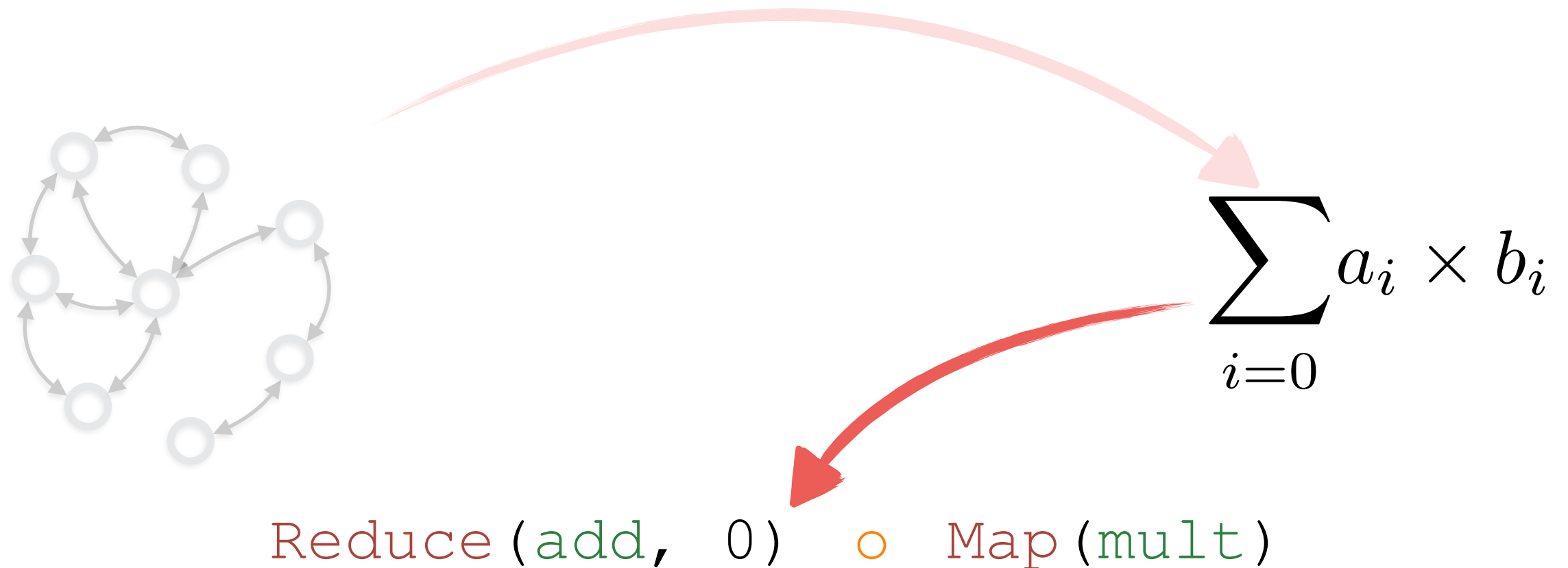
# Our proposal:



$$\sum_{i=0} a_i \times b_i$$

Translate problem to simpler domain:
(sparse) **Linear Algebra**

(intuition: operations on adjacency matrix of graph)

# Our proposal:

$$\sum_{i=0} a_i \times b_i$$

`Reduce(add, 0) o Map(mult)`

Implement linear algebra in parallel using a high level **functional** language
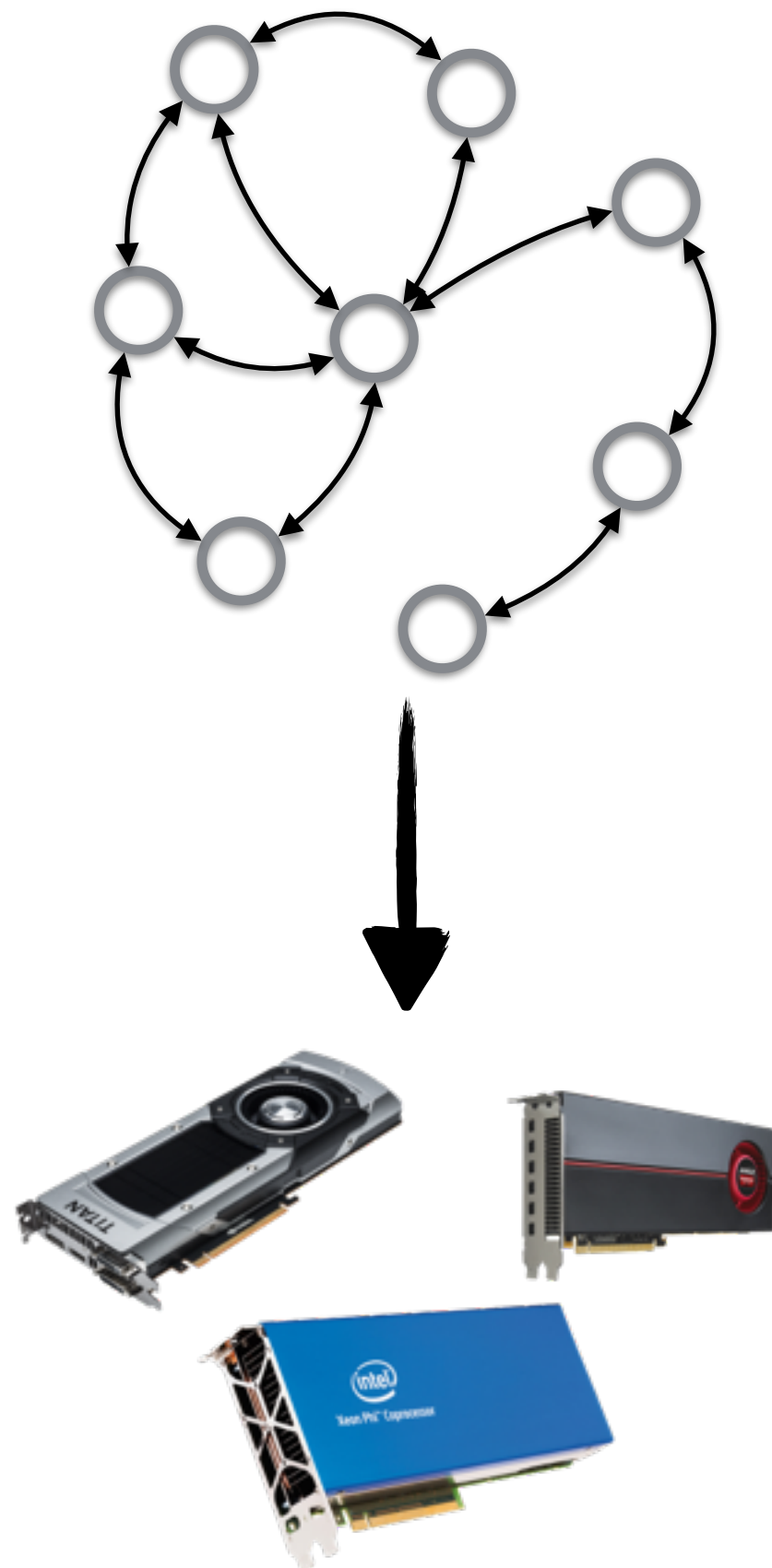
# Our proposal:

Use rewrite rules to automatically generate **high performance** OpenCL code

$$\sum_{i=0} a_i \times b_i$$

Reduce(add, 0)  ∘  Map(mult)

```
__kernel void dotProduct(
    __global const elem_t* a,
    __global const elem_t* b,
    __global const elem_t* res,
    __global const int len){
    int id = get_global_id(0);
    elem_t tmp;
    for(int i = id;i<len;i+=get_global_size(0)){
        tmp = a[i]*b[i];
        atomic_add(res, tmp);
    }
}
```

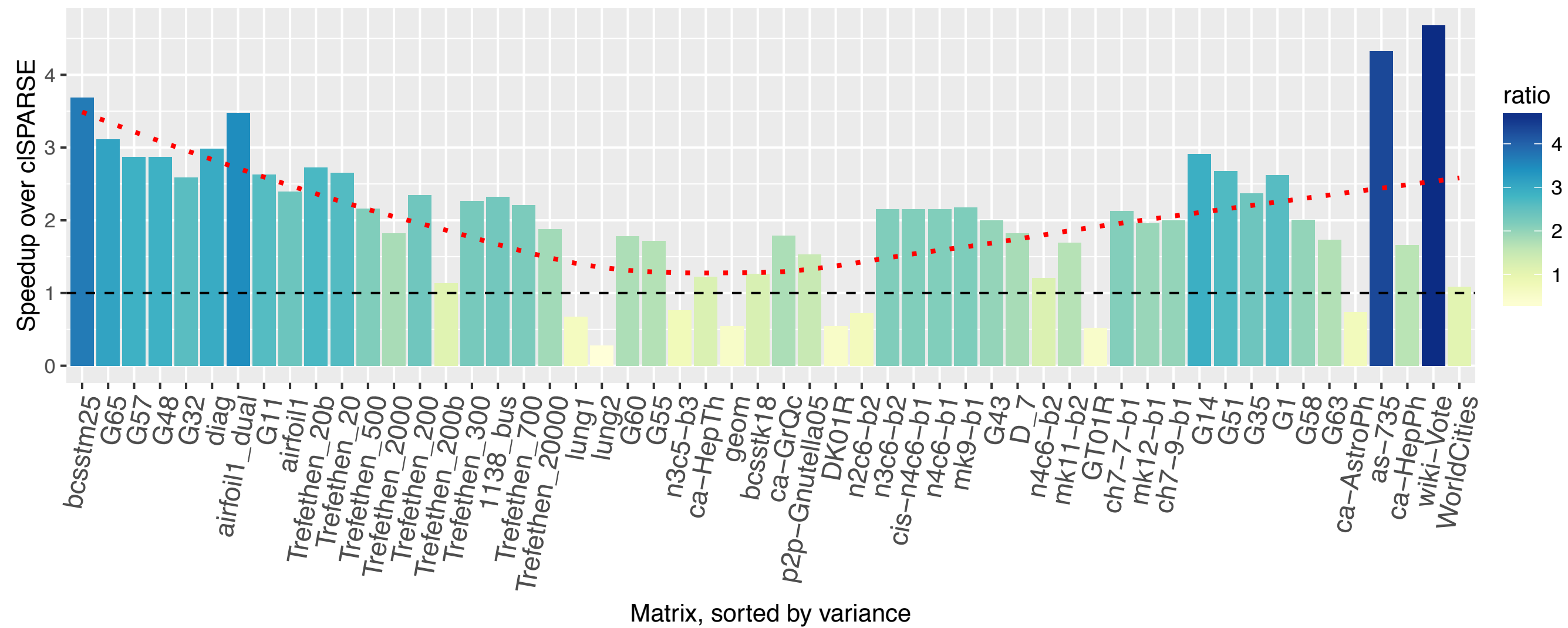**Sparsity:** Translate complex sparse domain to simpler dense domain

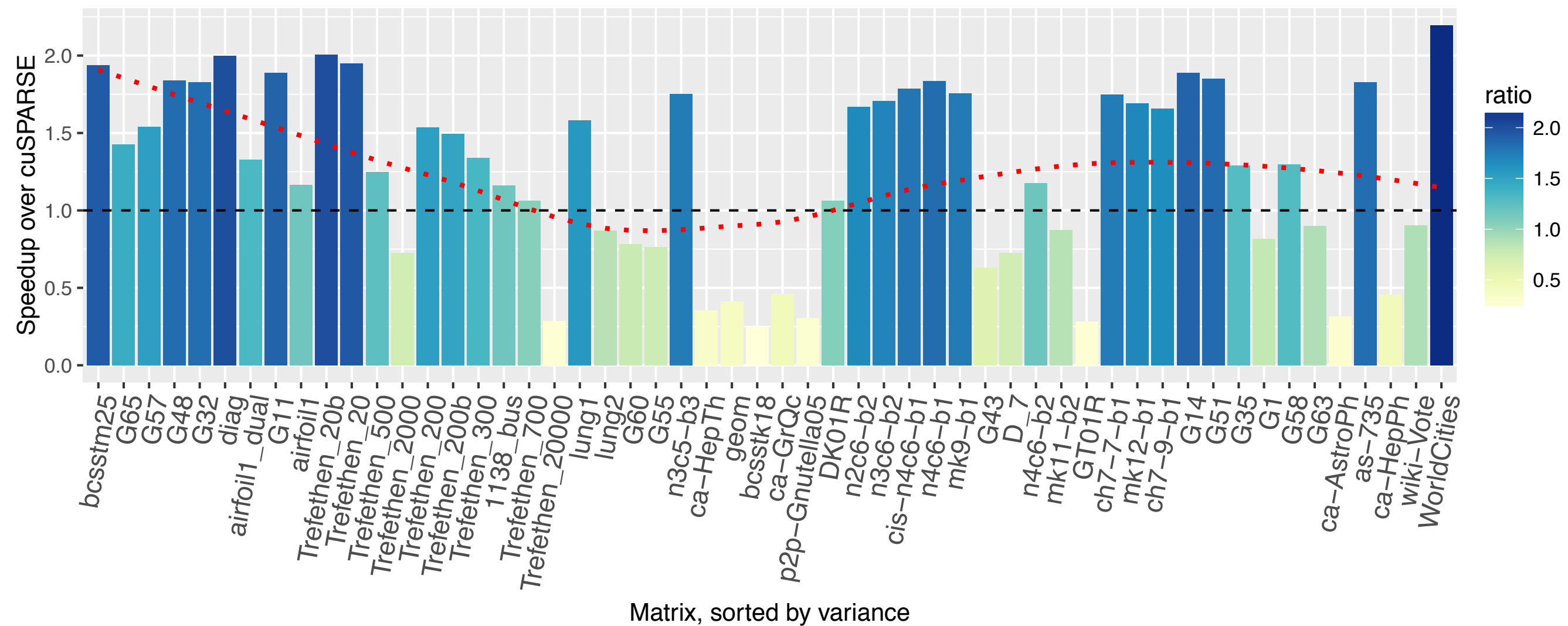**Irregularity:** Regularise code using intelligent compilation techniques

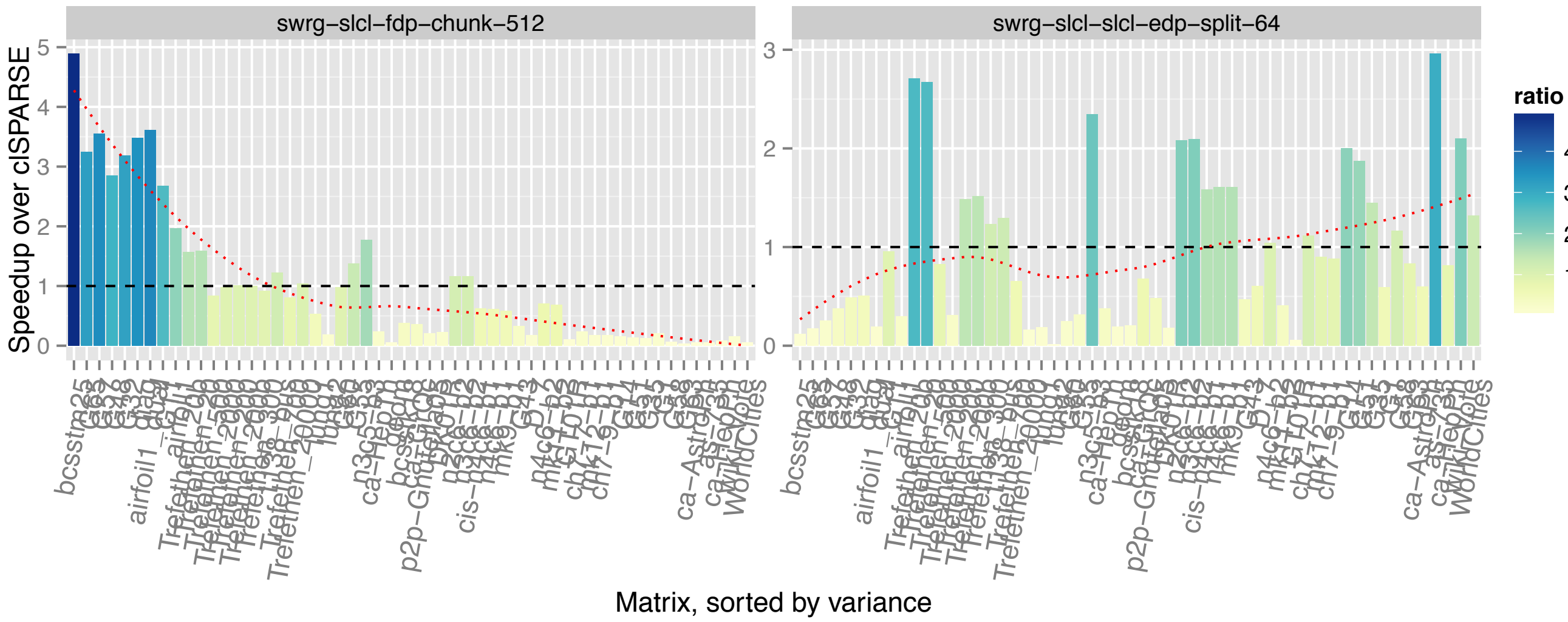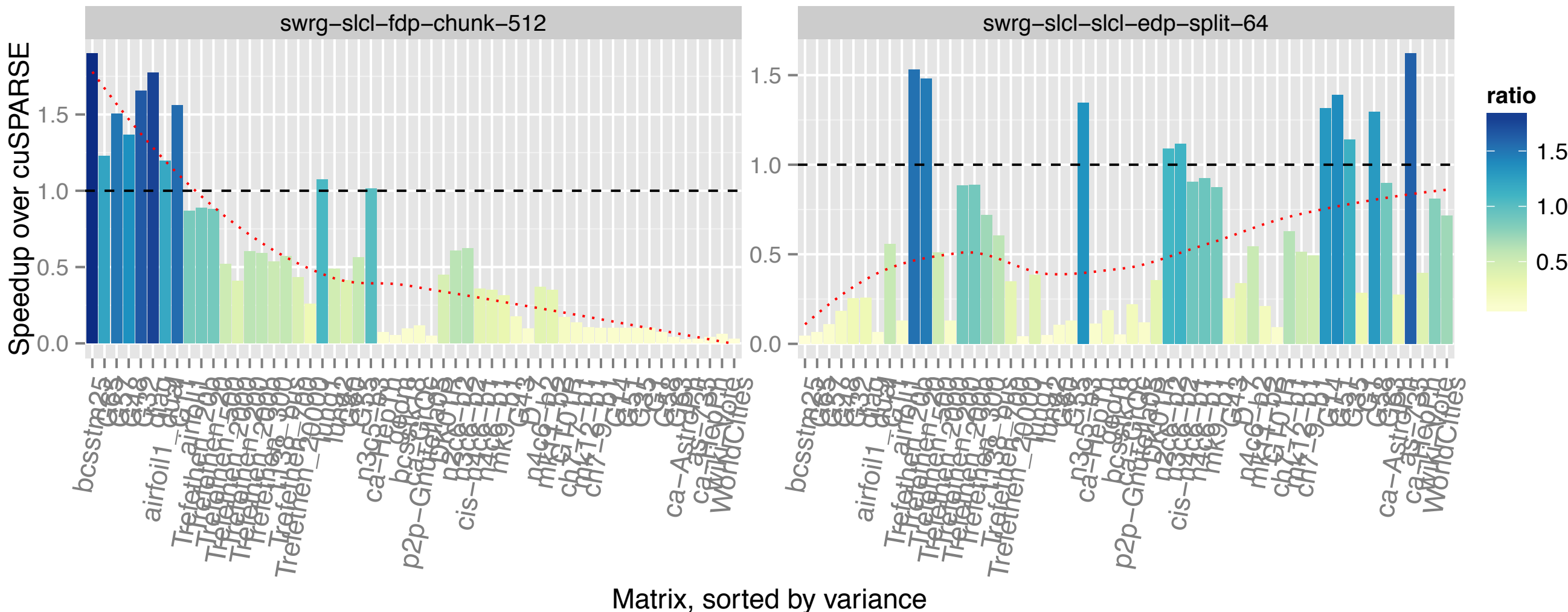**Current work:** extending to and evaluating on other domains.

AMD

NVIDIA

AMD