

# **B3CC: Concurrency**

## *16: Conclusion*

---

Ivo Gabe de Wolff

# Final exam!

- Final exam

- Tuesday 30 January @ 13:30
- Olympos Hal 2
- Mix of multiple choice and open questions
- Covers material from second half of the course:
  - From lecture 9 (Parallelism) through lecture 15 (Work & Span)
- You don't need to write Accelerate code
  - But you may be asked to design a parallel algorithm in terms of the parallel patterns
- Remindo has a calculator, no physical calculators allowed



# **Brief course summary**

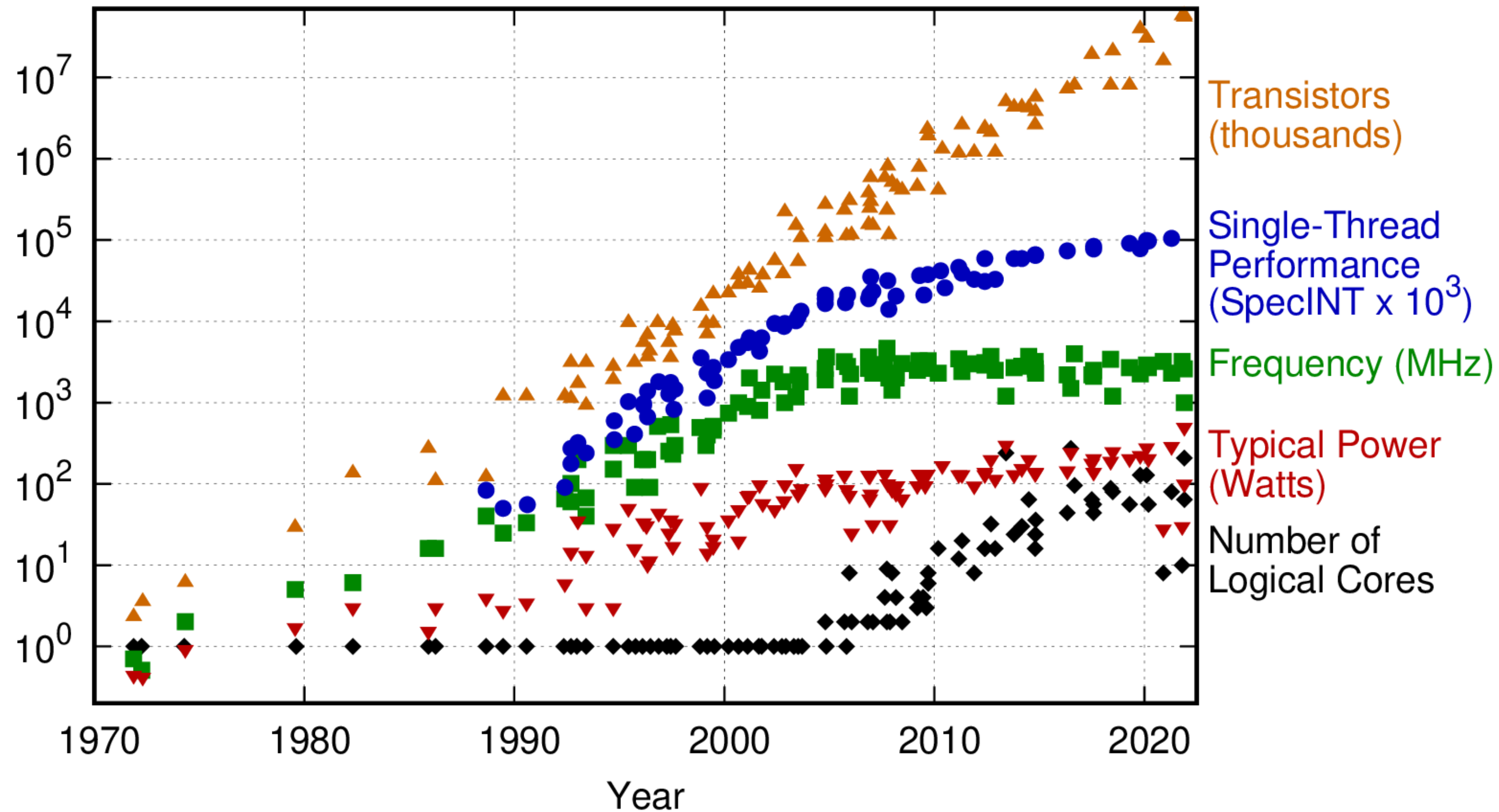
---

**What?**

# **Parallelism & Concurrency**

# Why?

## 50 Years of Microprocessor Trend Data



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten  
New plot and data collected for 2010-2021 by K. Rupp

# Where?

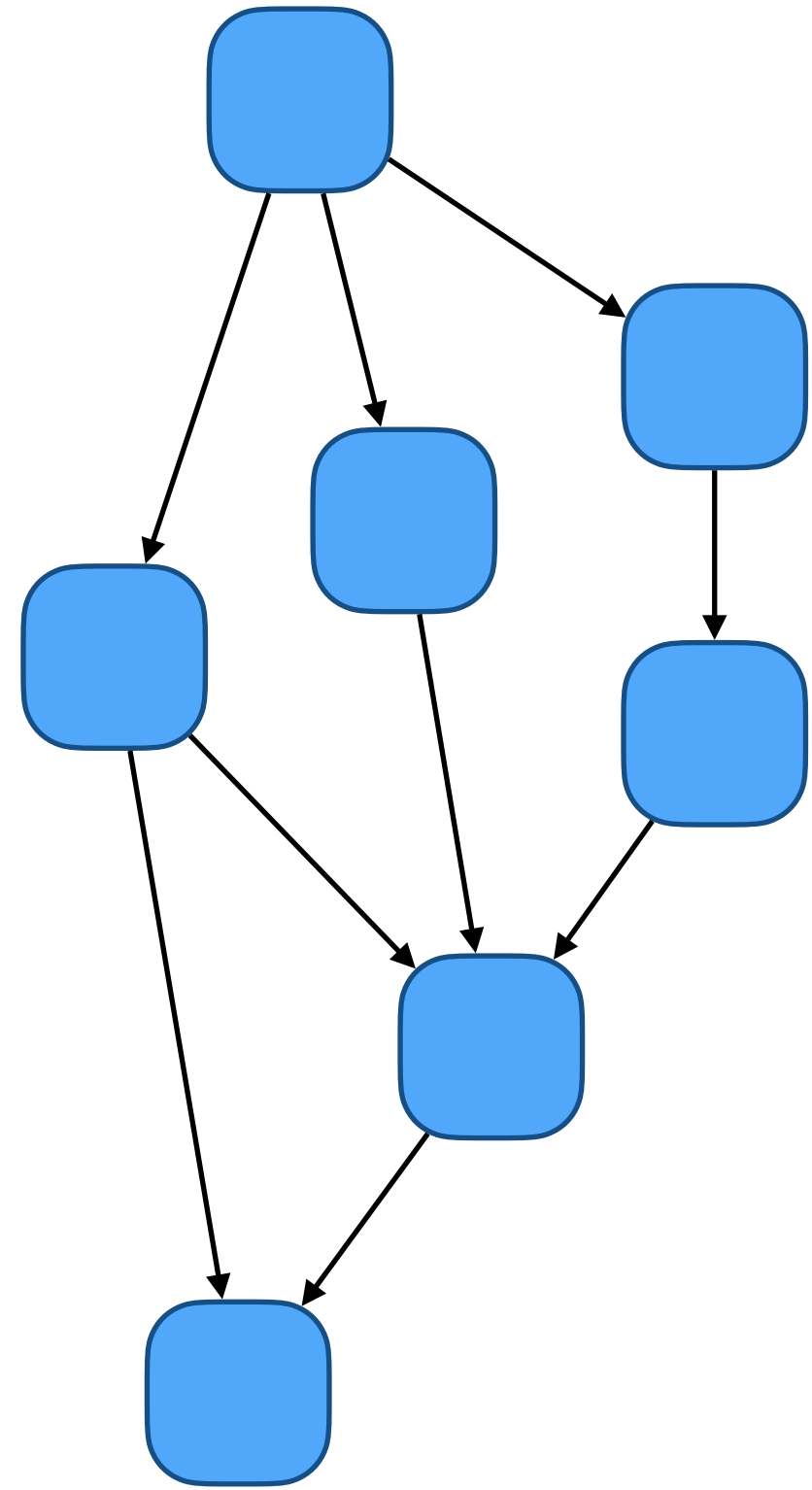
- Three kinds of code:
  - Gameplay simulation
    - Models the state of the game world as interacting entities
  - Numeric computation
    - Physics, collision detection, path finding, scene graph traversal, etc.
  - Shading
    - Pixel & vertex attributes; runs on the GPU



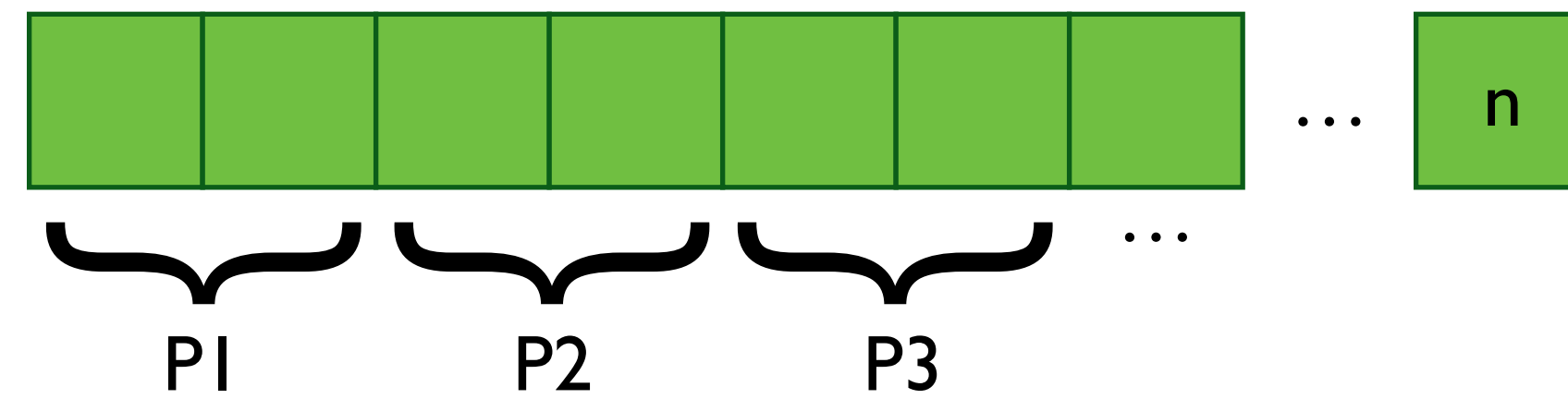
# How?

	<b>Game Simulation</b>	<b>Numeric Computation</b>	<b>Shading</b>
<b>Languages</b>	C++, scripting	C++	GC, HLSL
<b>CPU Budget</b>	10%	90%	n/a
<b>Lines of Code</b>	250.000	250.000	10.000
<b>FPU Usage</b>	0.5 GFLOPS	5 GFLOPS	500 GFLOPS
<b>Concurrency/ Parallelism</b>	<b>STM</b>	<b>SIMD</b>	<b>GPU</b>

# Kinds of parallelism



**Task parallelism**

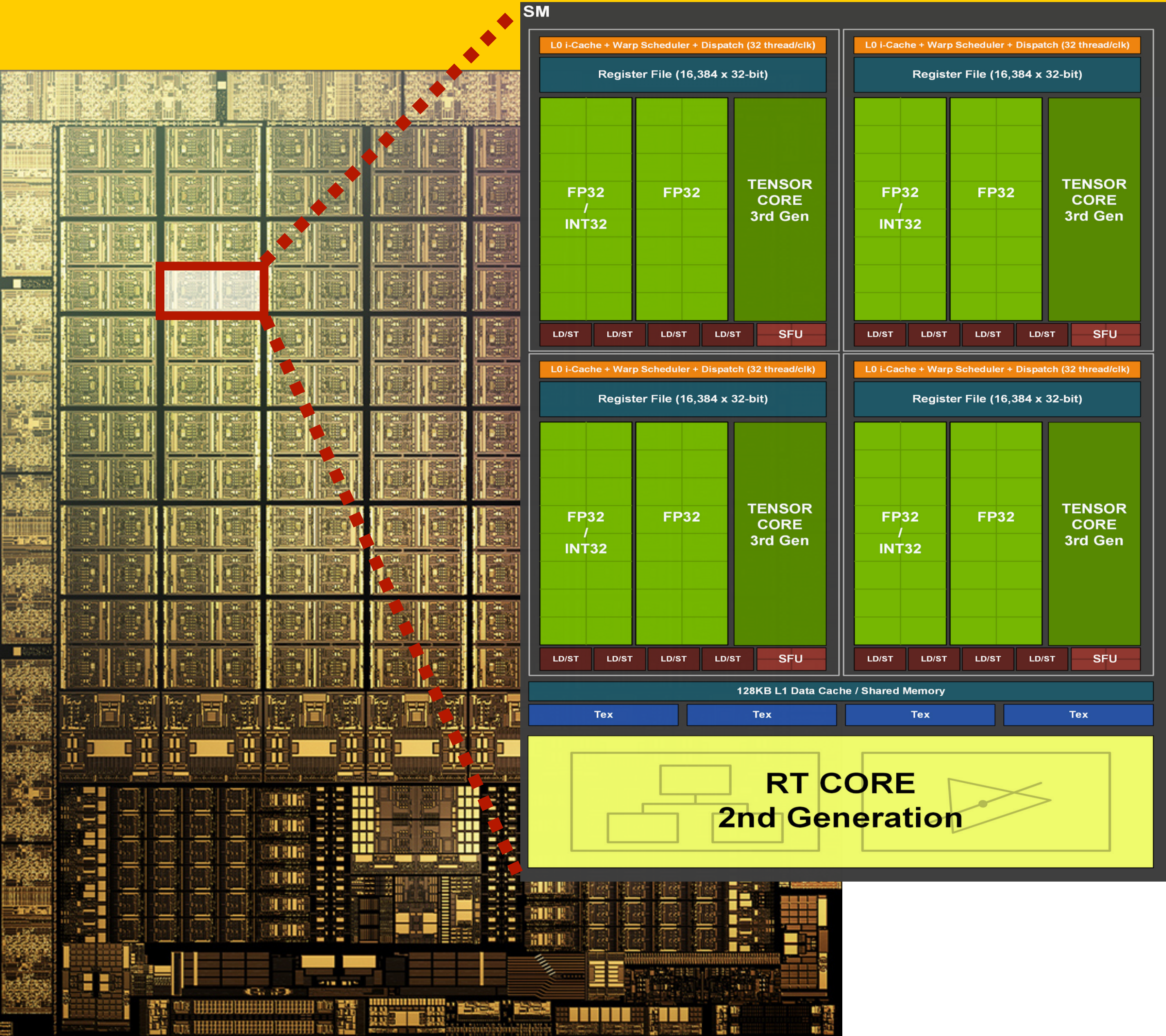


**Data parallelism**

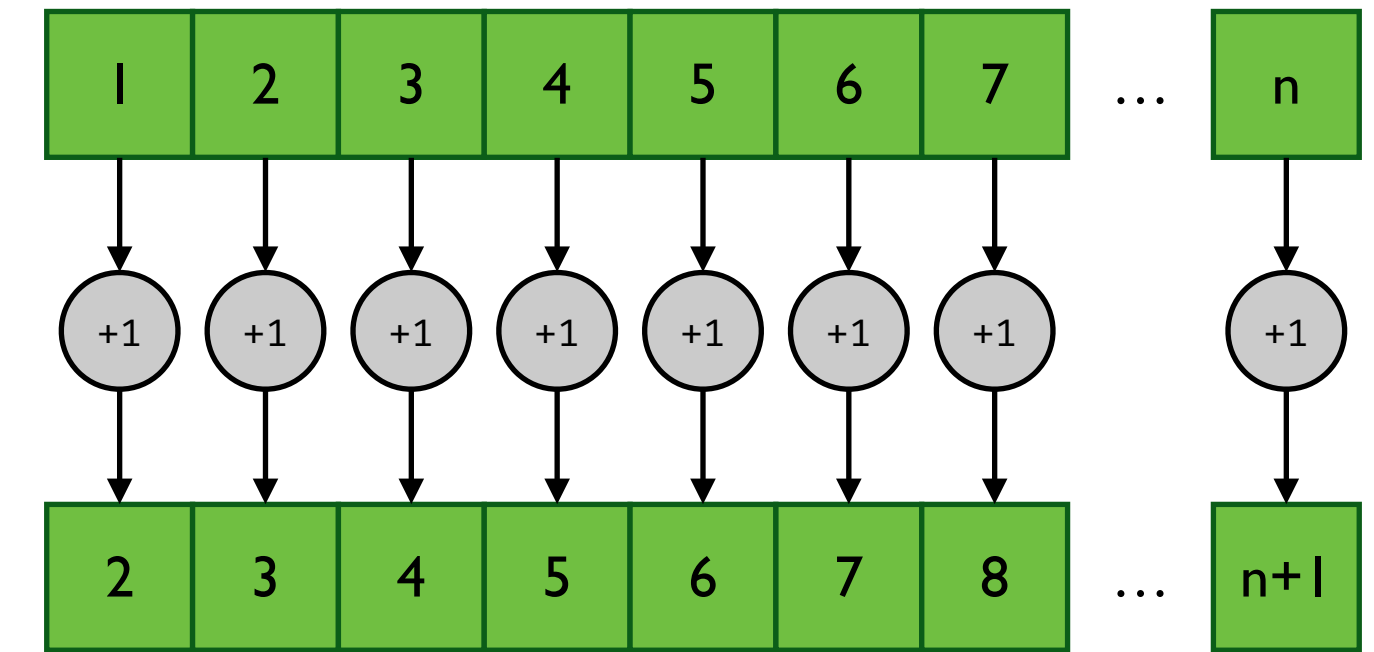


# GPGPU

- How the parallel patterns we have talked about map to GPU code
- Difference between CPU and GPU
  - What each is designed for; strengths and weaknesses
  - What the GPU programming model (CUDA) is designed for

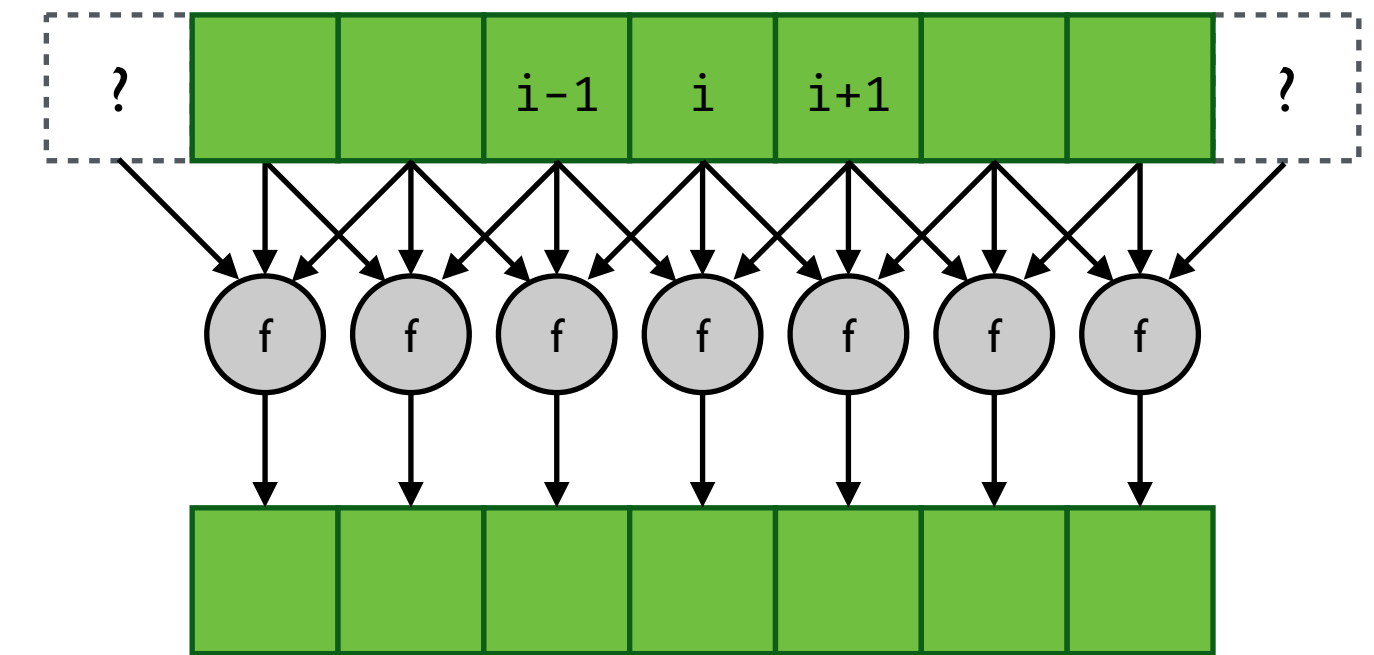


# Patterns: map



- Apply a function to every element of an array, independently
- This one is (hopefully) straightforward...

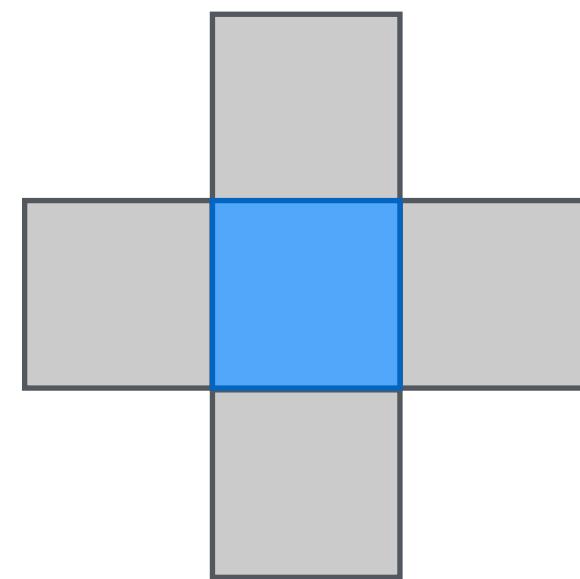
# Patterns: stencil



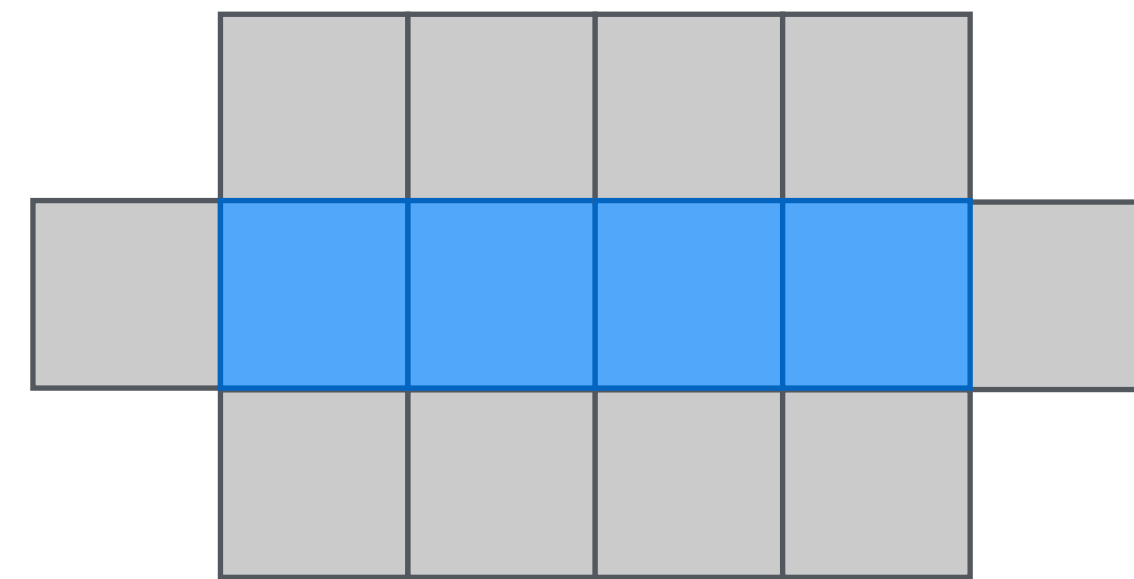
- A map with access to the surrounding neighbourhood
- What are the difficulties/limitations?
  - The ghost region (halo) and how/why to use it
  - Optimisations (tiling, strip mining, etc.)

# Stencil optimisations

- Use a different kernel for the interior and border regions
  - In the gaussian blur example of a 512x512 pixel image, 98% of the pixels do not require in-bounds checks
- Optimise data locality & reuse through tiling
  - *Strip mining* is an optimisation that groups elements in a way that avoids redundant memory access and aligns accesses with cache lines



4 x (5 reads + 1 write)

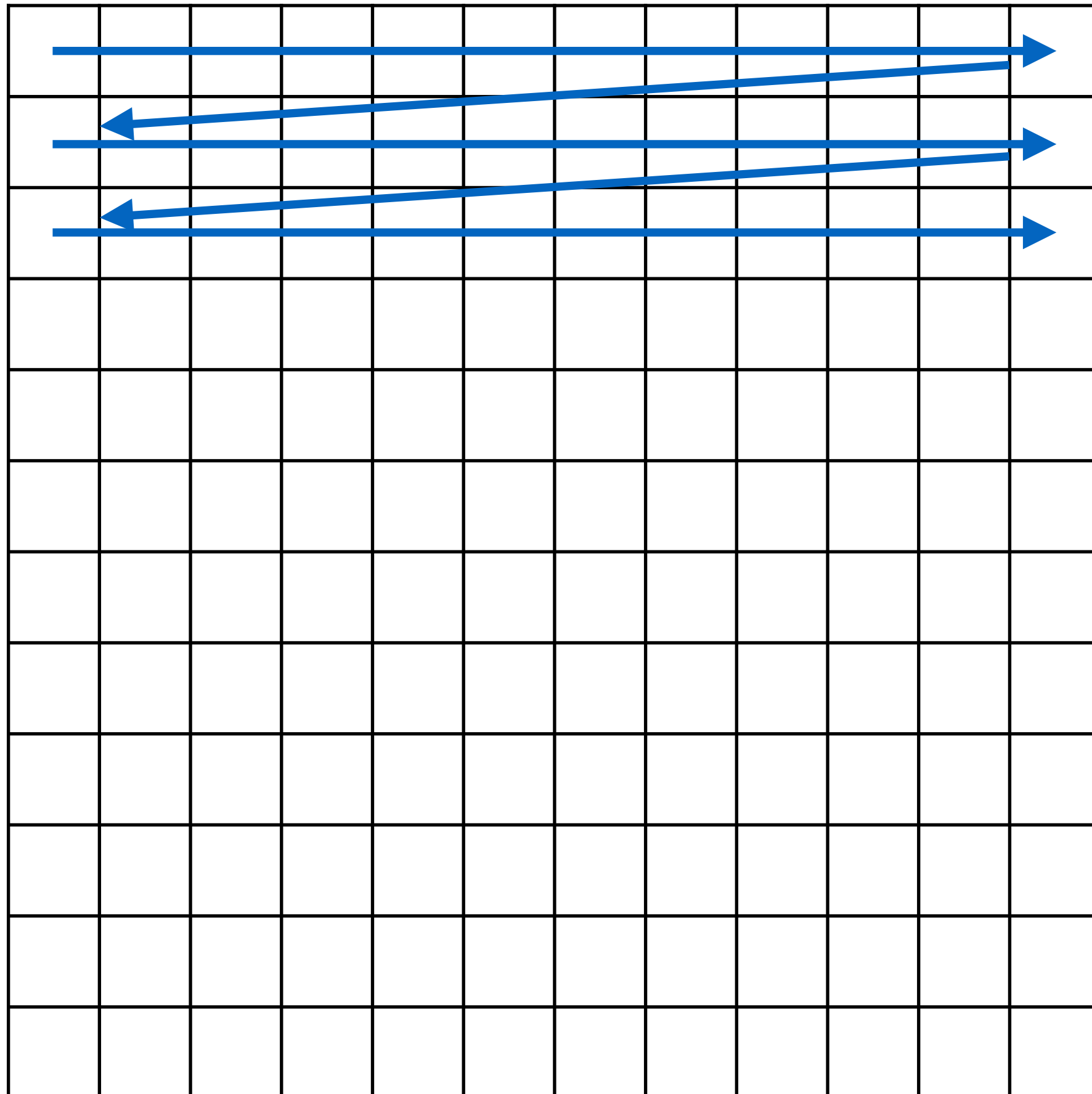


14 reads + 4 writes

# Stencil optimisations

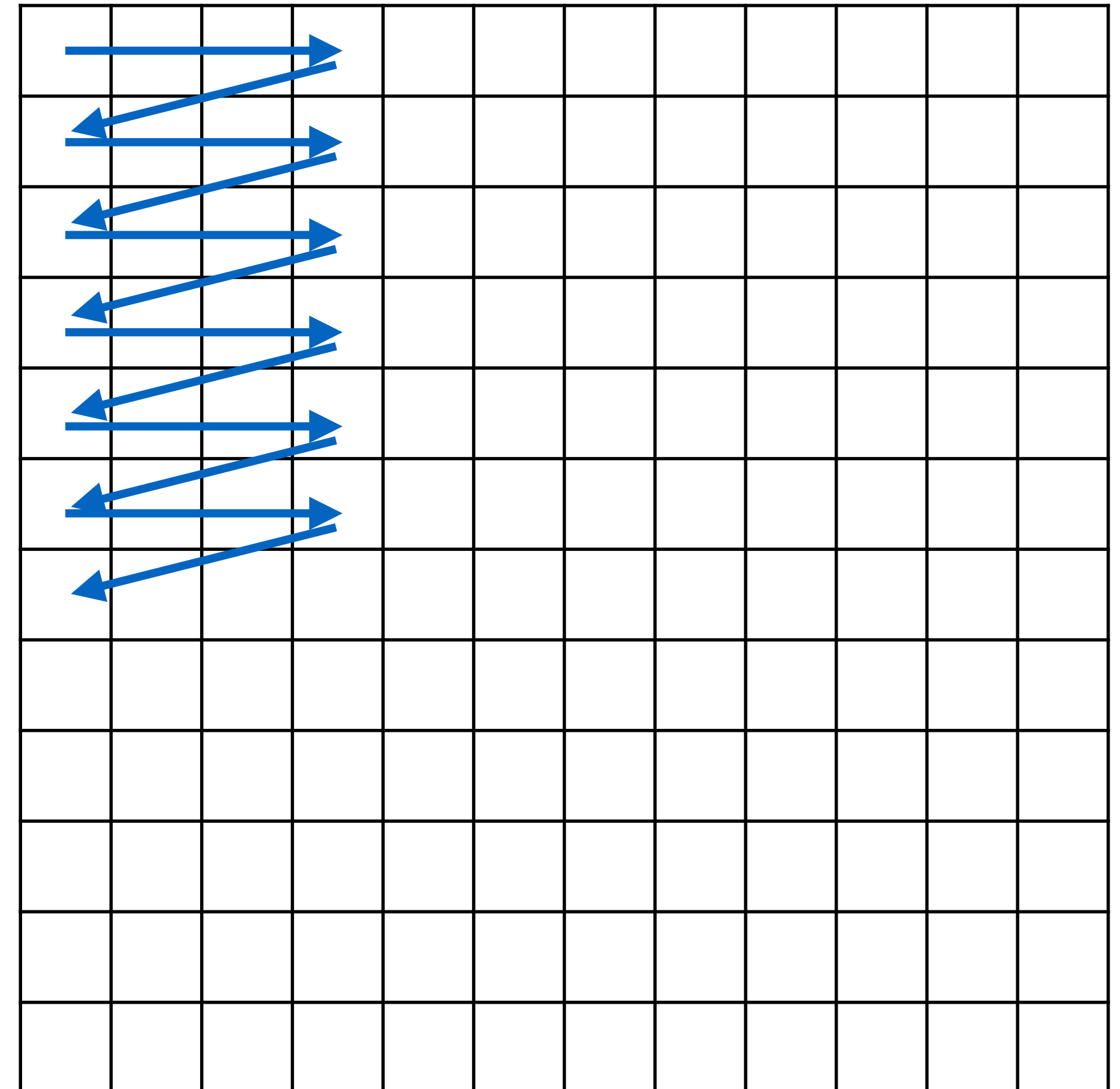
## Without tiling

- When handling row 0, row 1 is loaded in cache.
- First values of row 1 may already be out of cache, when handling row 1



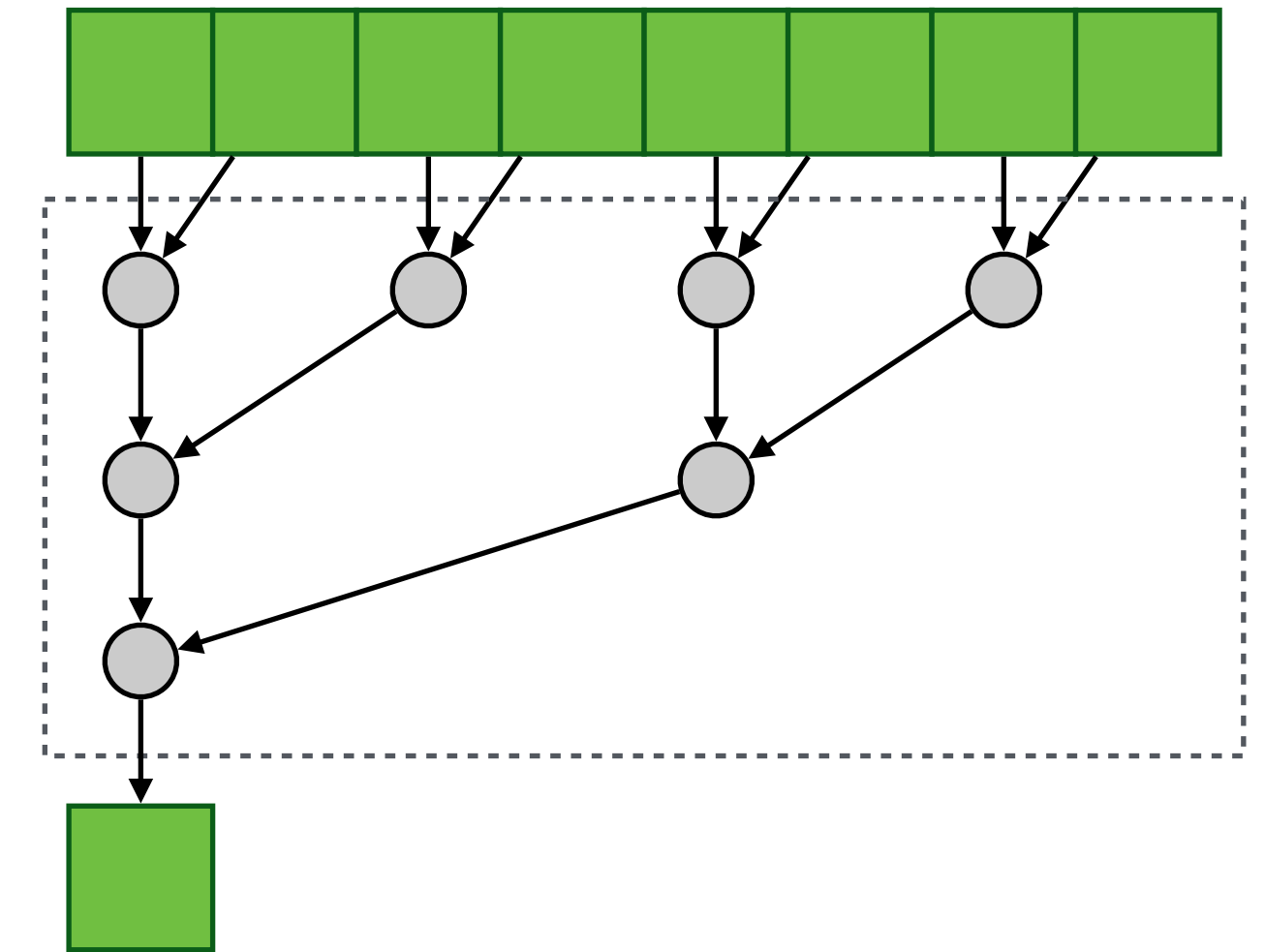
## With tiling

- Previously loaded row is still in cache
- Tile width is usually a power of 2, on GPUs often the warp size (32)



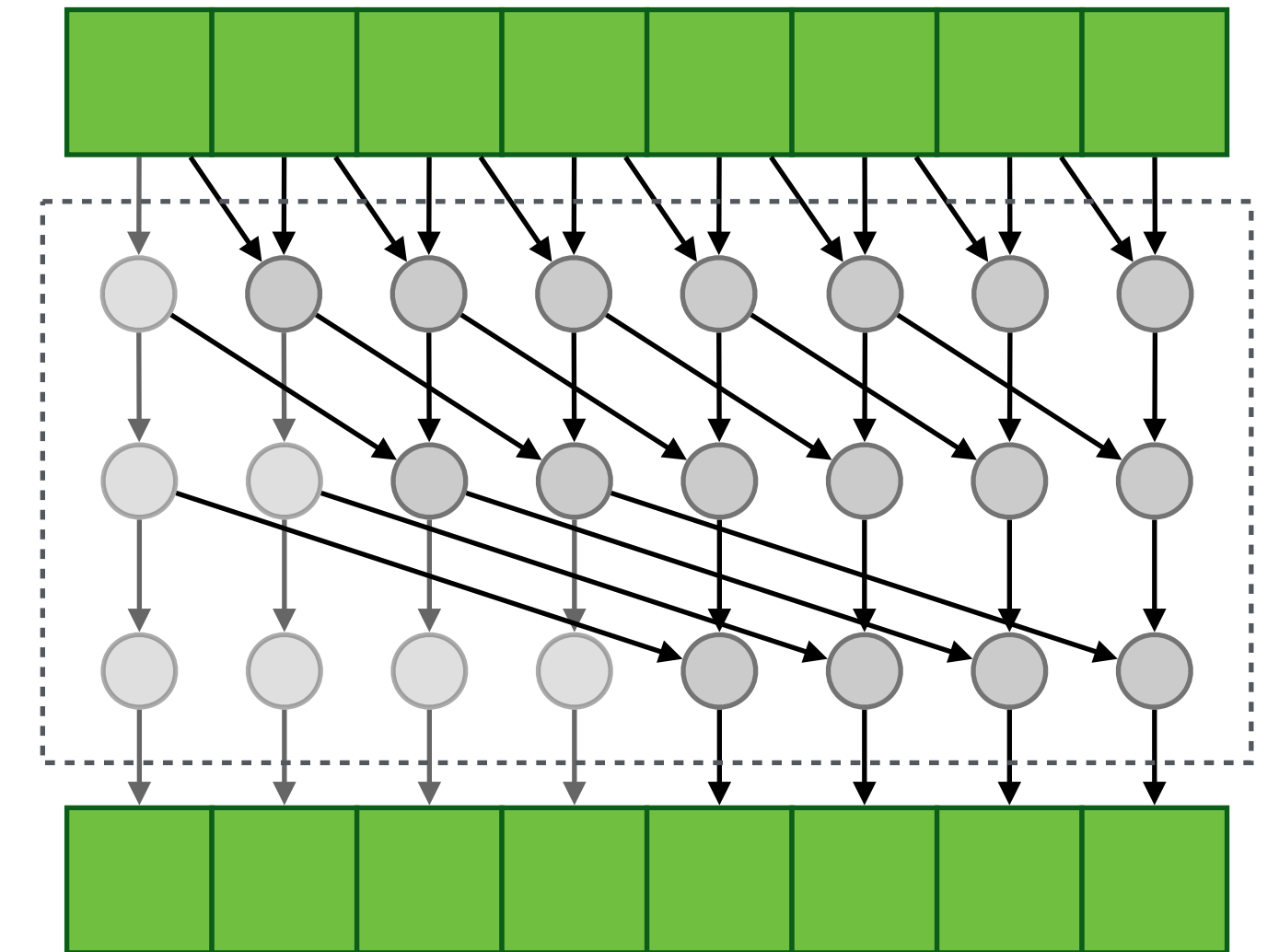
# Patterns: fold

- Reduce an array to a summary value
- How to implement this in parallel
  - What kinds of restrictions are necessary?
  - What additional restrictions can be leveraged to improve it further?



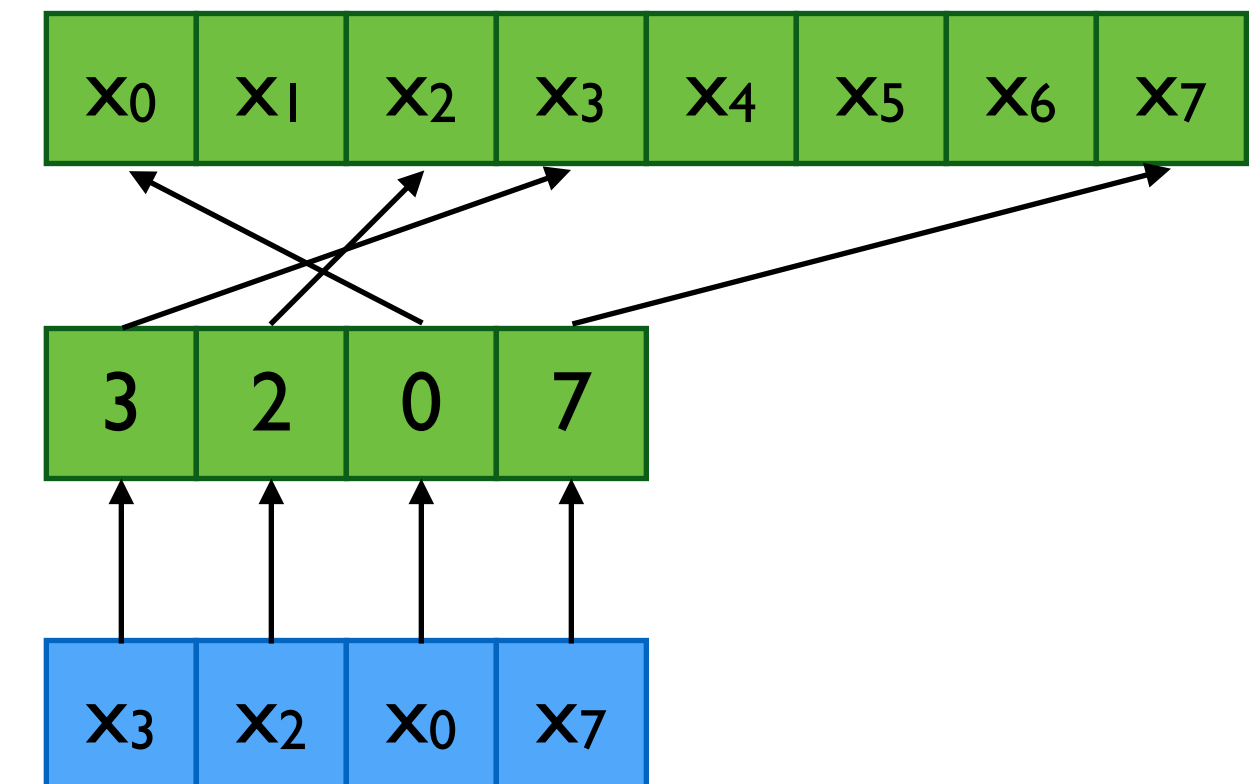
# Patterns: scan

- All partial reductions of an array
- Varieties
  - Inclusive vs. exclusive etc.
- Parallel implementation
  - Restrictions, etc.



# Patterns: gather

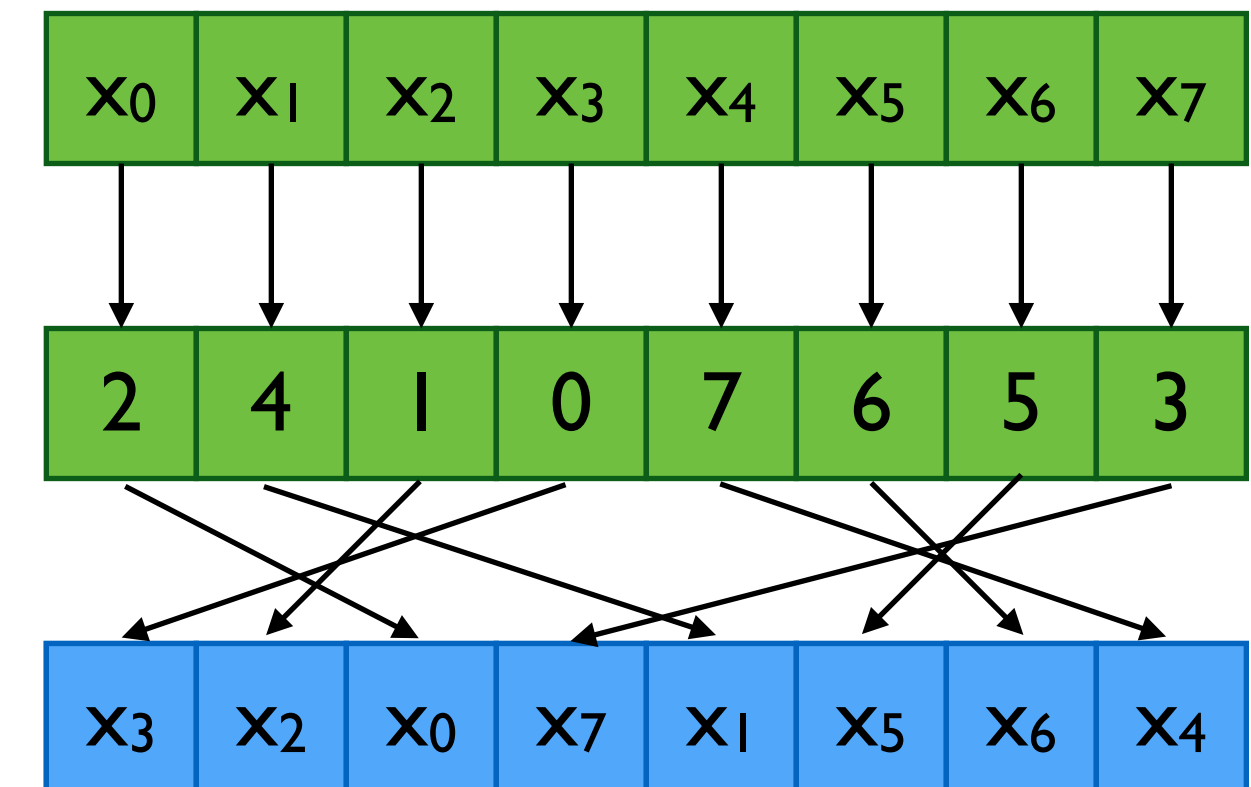
- Parallel random read
- Implications for memory access patterns
  - Optimisations for special cases (e.g. transpose), like tiling
  - Implications for the GPU, caches, etc.





# Patterns: scatter

- Parallel random write
- How to handle collisions in the index permutation function
  - Performance implications of collisions, false sharing, etc.
  - Scatter vs. gather



# Patterns: gather vs scatter

- Random reads (gather) are slower than structured reads
- Random writes (scatter) are slower than structured writes
- This problem is larger for scatter,  
as the processor needs to perform more synchronization between cores
- In general, use gather instead of scatter if both are possible

# Patterns

- You should be able to:
  - Give examples for each pattern
  - Recognise these patterns and where they can be used
    - e.g. given a problem description, give an implementation in terms of these patterns
    - Use Accelerate code, pseudocode or an explanation in text
      - Especially for the latter, make sure your explanation is concrete

# Work & Span

- We analysed the performance of algorithms using the *work* and *span*:
  - **Work** =  $T_1$   
How long to execute on a single processor
  - **Span** =  $T_\infty$   
How long to execute on an infinite number of processors
    - The longest dependence chain / critical path length / computational depth
    - Example:  $O(\log n)$  for summing an array

# Efficient & optimal

- The parallelisation *overhead* of an algorithm is its work divided by the cost of the best sequential algorithm
  - For this parallel scan we have to put  $O(n \log n)$  work into something which can be done sequentially in linear  $O(n)$  time: the overhead is logarithmic
  - A parallel algorithm is:
    - *Efficient* when the span is poly-logarithmic and the overhead is also poly-logarithmic
    - *Optimal* when the span is poly-logarithmic and the overhead is constant

# Master Theorem

- The master theorem provides a solution to recurrence relations of the form
  - For constants  $a \geq 1$  and  $b > 1$  and  $f$  asymptotically positive

$$T(n) = aT\left(\frac{n}{b}\right) + f(n)$$

- The master theorem has three cases:

Recursion dominates

If  $f(n) = O(n^{\log_b a - \epsilon})$   
for some  $\epsilon > 0$ ,  
then  $T(n) = \Theta(n^{\log_b a})$

Both contribute

If  $f(n) = \Theta(n^{\log_b a})$ , then  
 $T(n) = \Theta(n^{\log_b a} \log n)$

$f$  dominates

If  $f(n) = \Omega(n^{\log_b a + \epsilon})$   
for some  $\epsilon > 0$ , and  
 $af(n/b) \leq cf(n)$   
for some  $c < 1$   
for all  $n$  sufficiently large,  
then  $T(n) = \Theta(f(n))$

# Analysis of parallel algorithms

- You should be able to:
  - Compute the work and span given a problem description/code
  - Compare parallel algorithms
    - Efficient & optimal
    - Parallel speedup (Amdhal vs. Gustafson-Baris)

# Questions?

---



# Finally...

- Please fill out the Thermometer survey!
  - All constructive feedback is welcome
  - <https://caracal.uu.nl/35916/Respond>

# success!

exam

you

me

claudiopiccoli.com