# Spatial computation on a homogeneous, many-core architecture
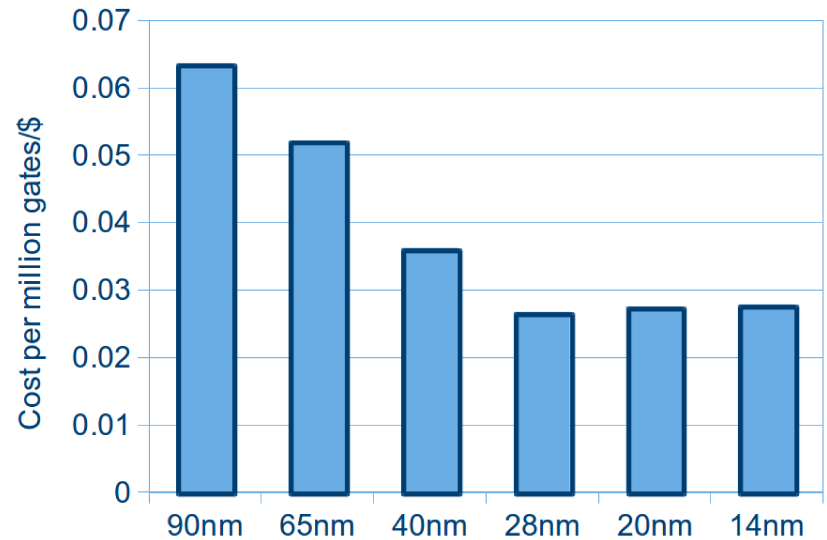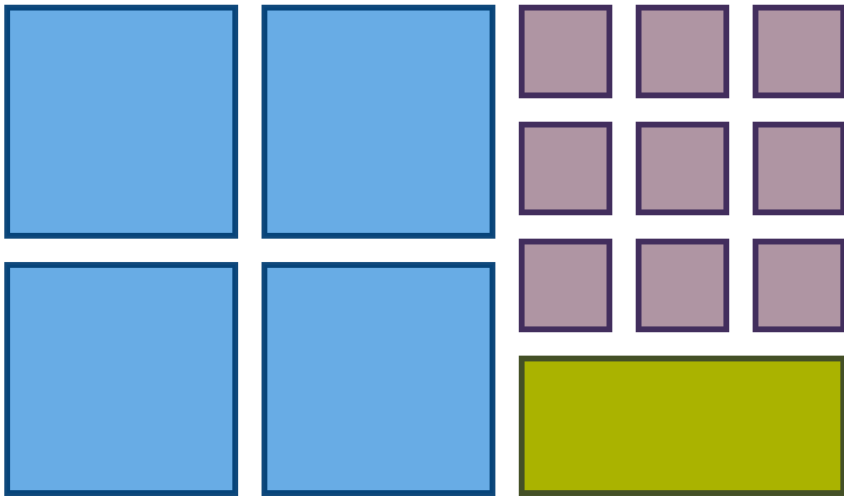
**Daniel Bates, Alex Bradbury, Andreas Koltes and Robert Mullins**

UNIVERSITY OF CAMBRIDGE

**PRISM-2**

# Motivation: ever-changing tradeoffs

Heterogeneity won't be optimal forever
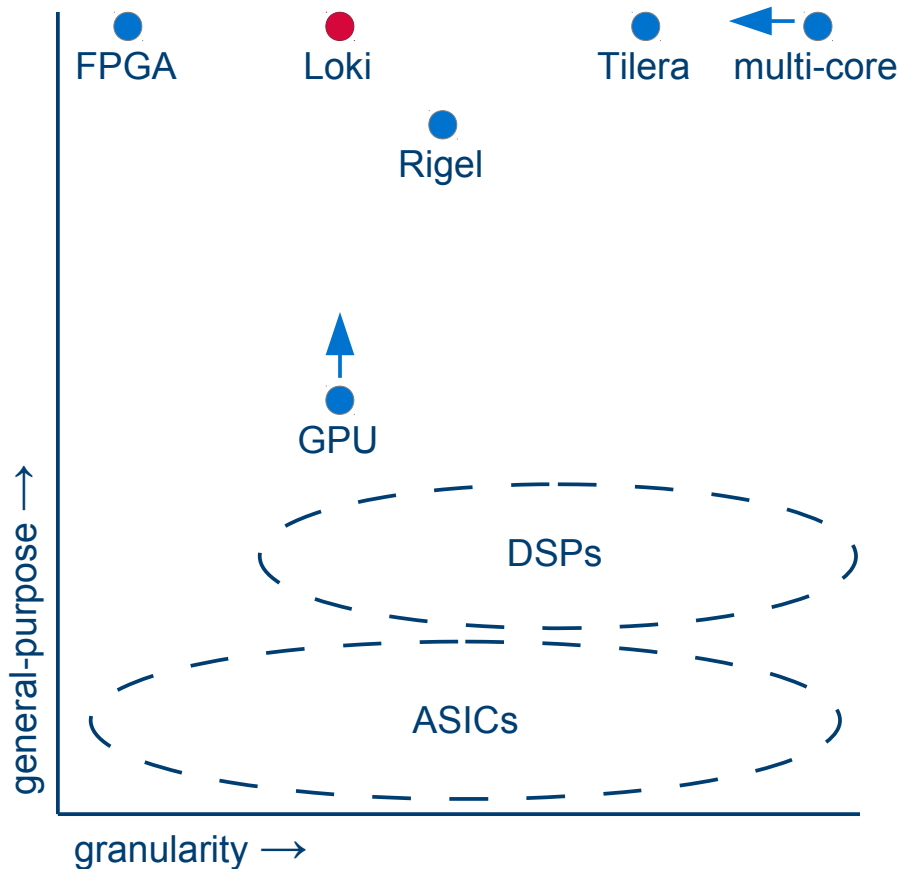
- Rising complexity

- Rising transistor costs

Source: *Feature dimension reduction slowdown*, by Handel Jones
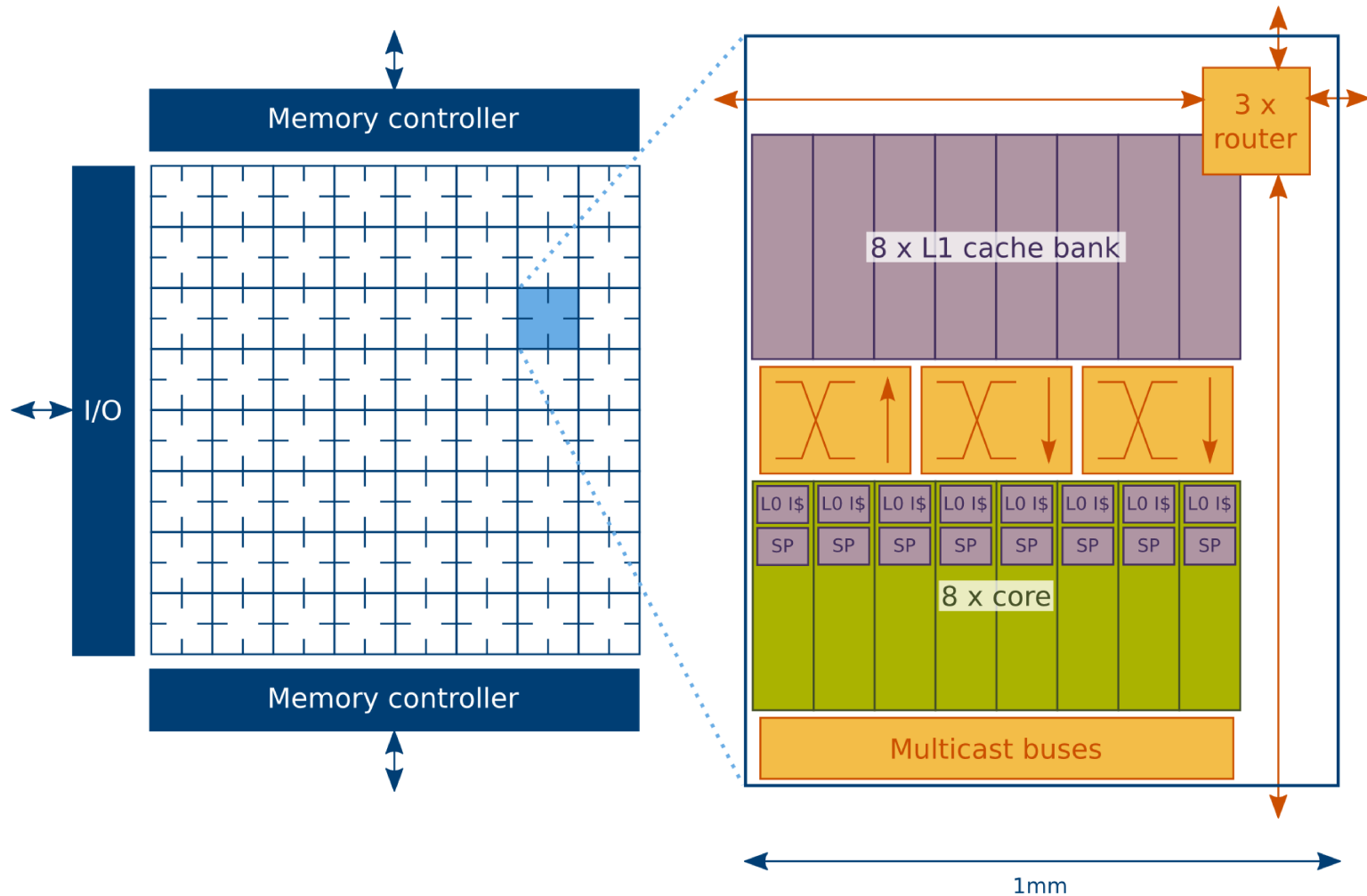
UNIVERSITY OF CAMBRIDGE

# Related work

Many projects in this field:

- Elm

- Raw/Tilera

- Rigel

- TRIPS

- …

None have the same focus on flexibility and cooperation between cores.
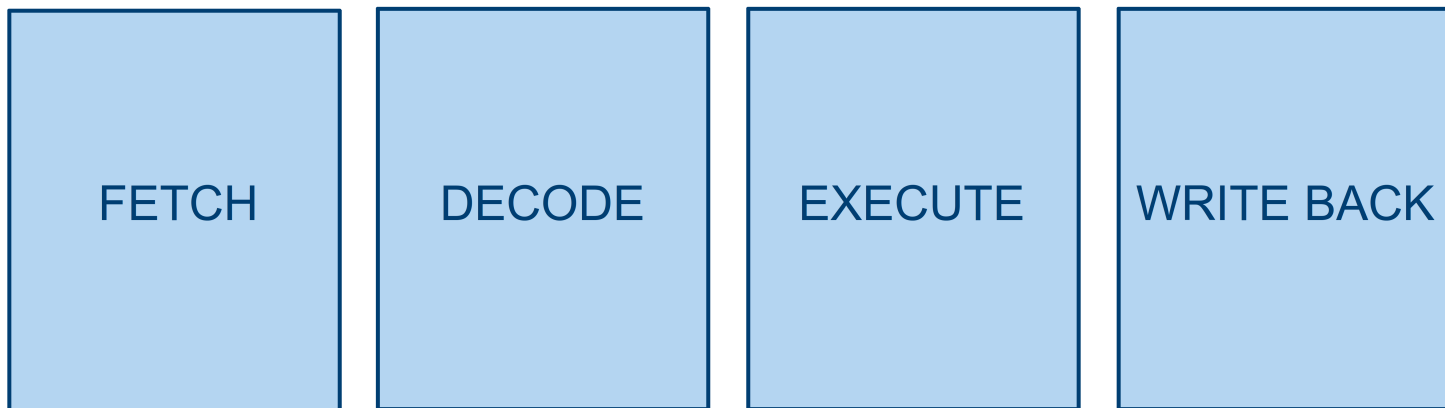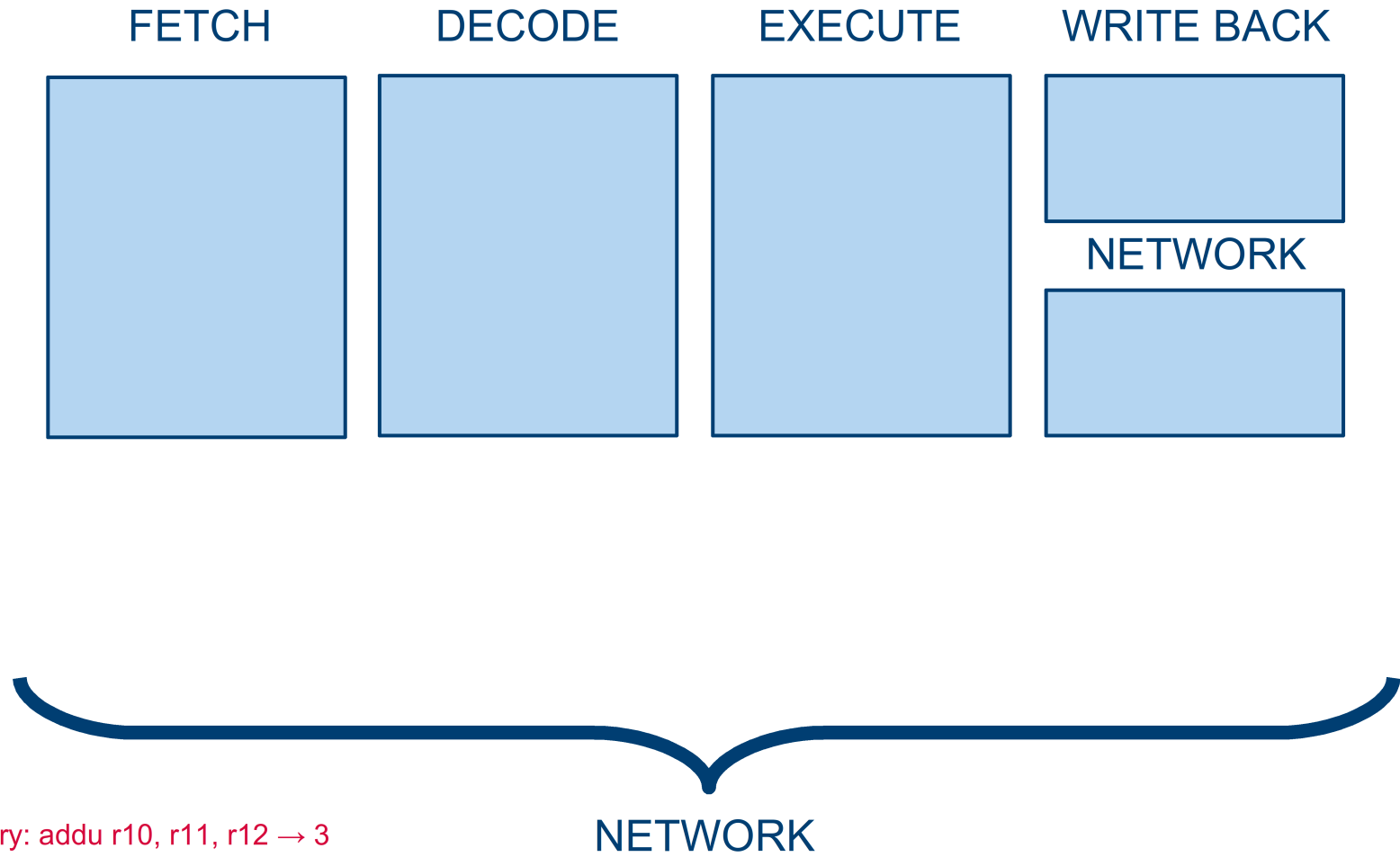
# Loki architecture: chip and tile
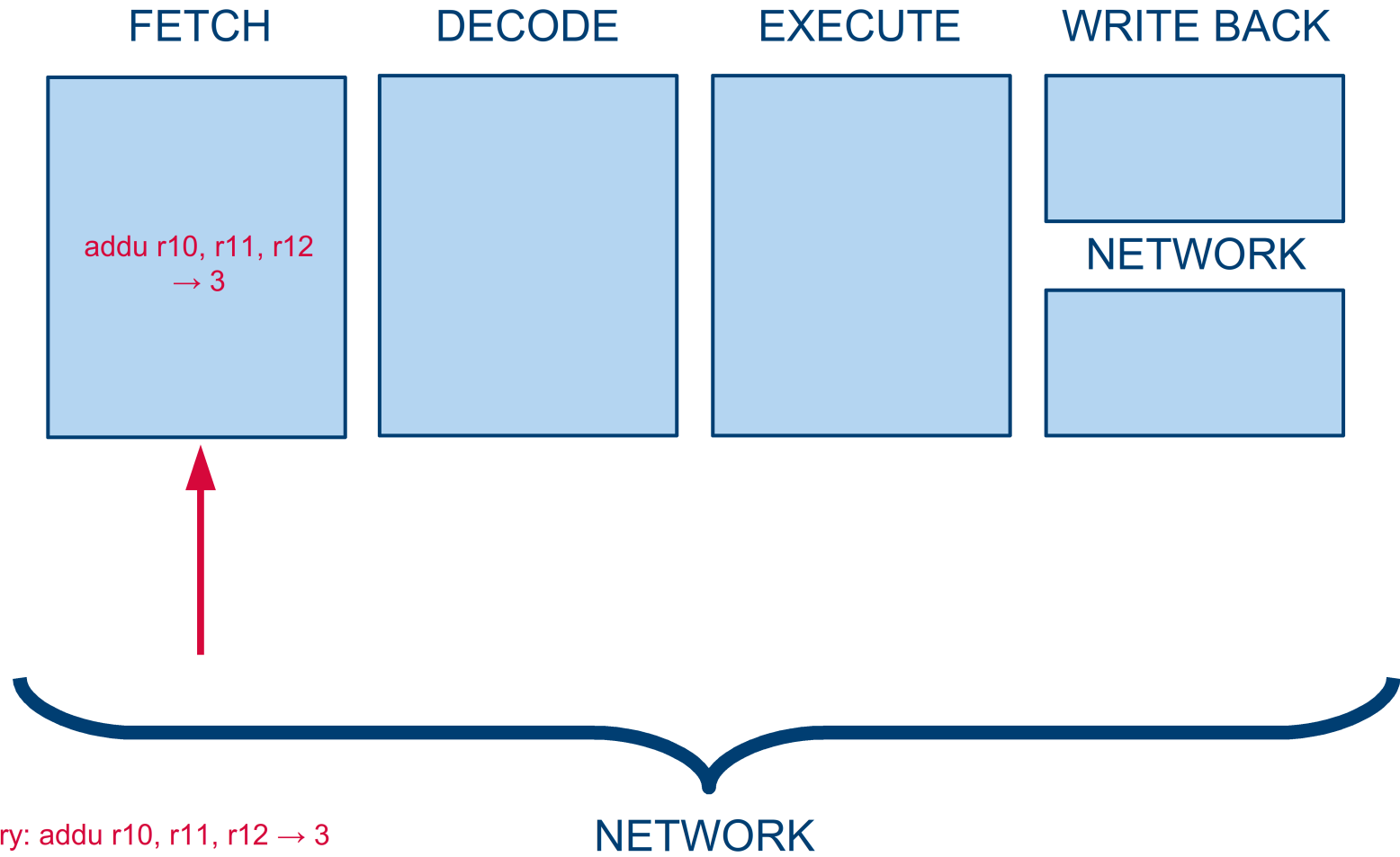
# Loki architecture: key features

- 4-stage, in-order, single-issue pipeline

- Deep network integration
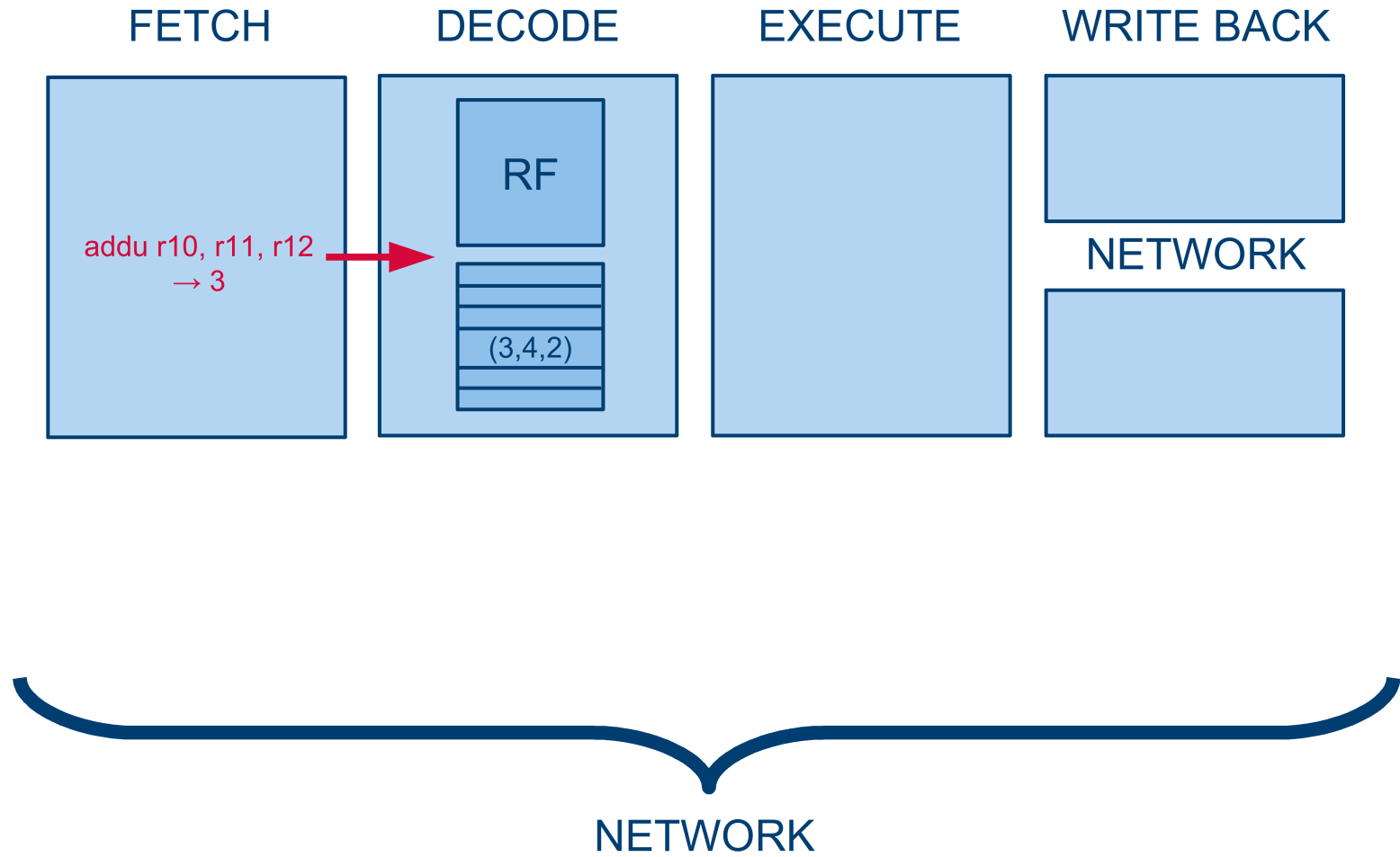
  - No memory access stage

- Indirect network access

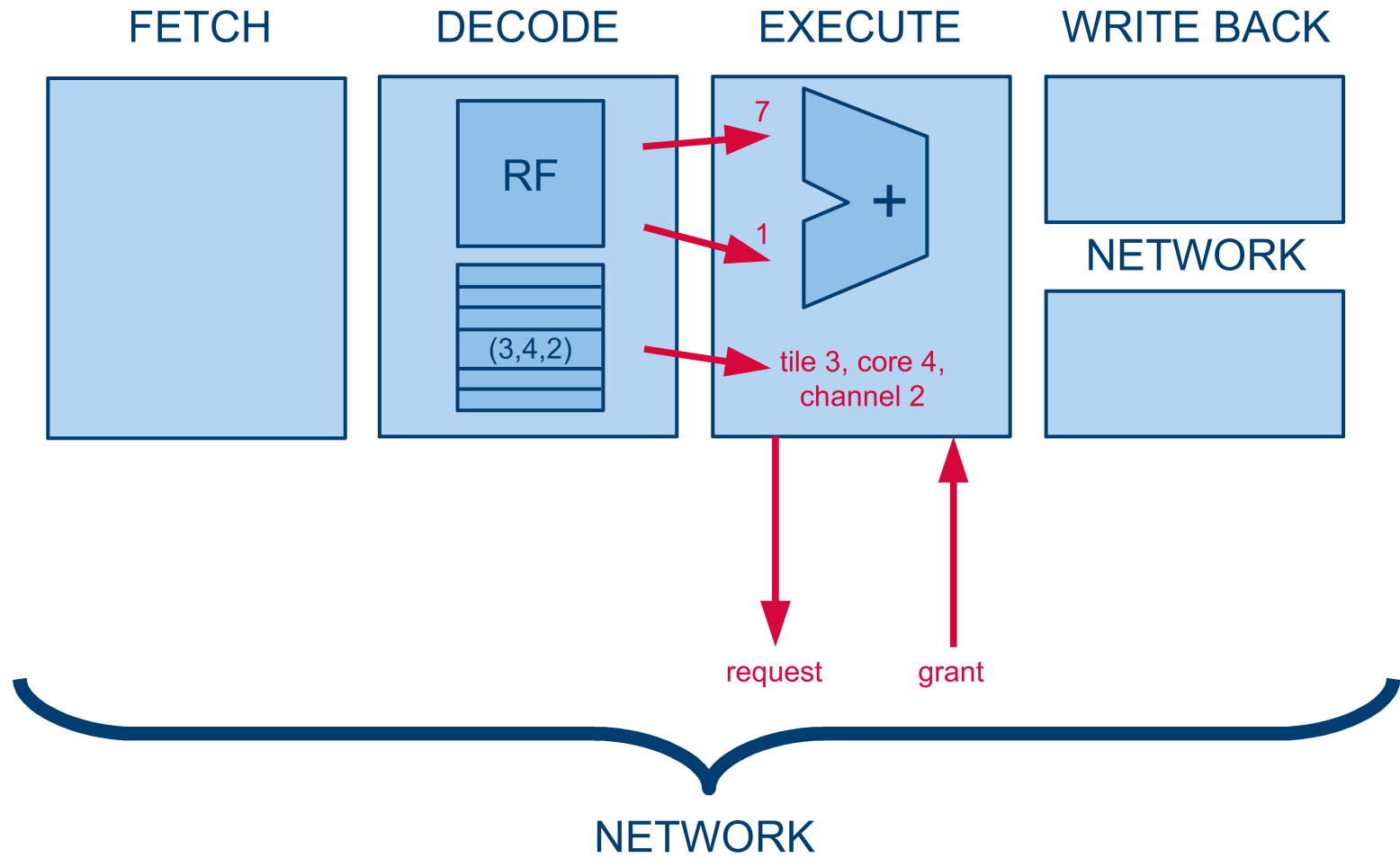| FETCH | DECODE | EXECUTE | WRITE BACK |

# Loki architecture: key features

FETCH          DECODE          EXECUTE          WRITE BACK

NETWORK

Memory: addu r10, r11, r12 → 3

NETWORK

UNIVERSITY OF CAMBRIDGE

# Loki architecture: key features

# Loki architecture: key features

FETCH            DECODE            EXECUTE            WRITE BACK

RF

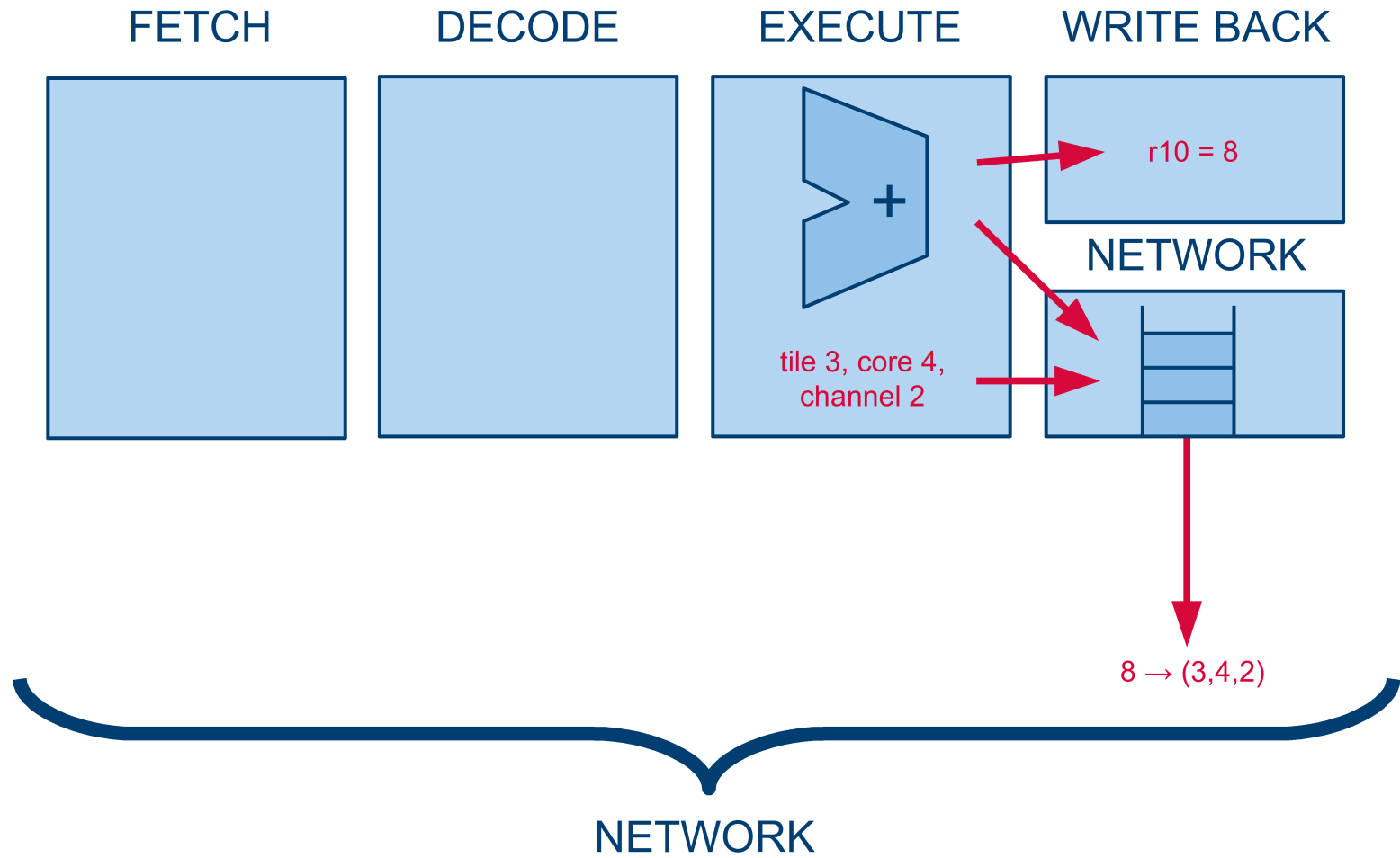addu r10, r11, r12
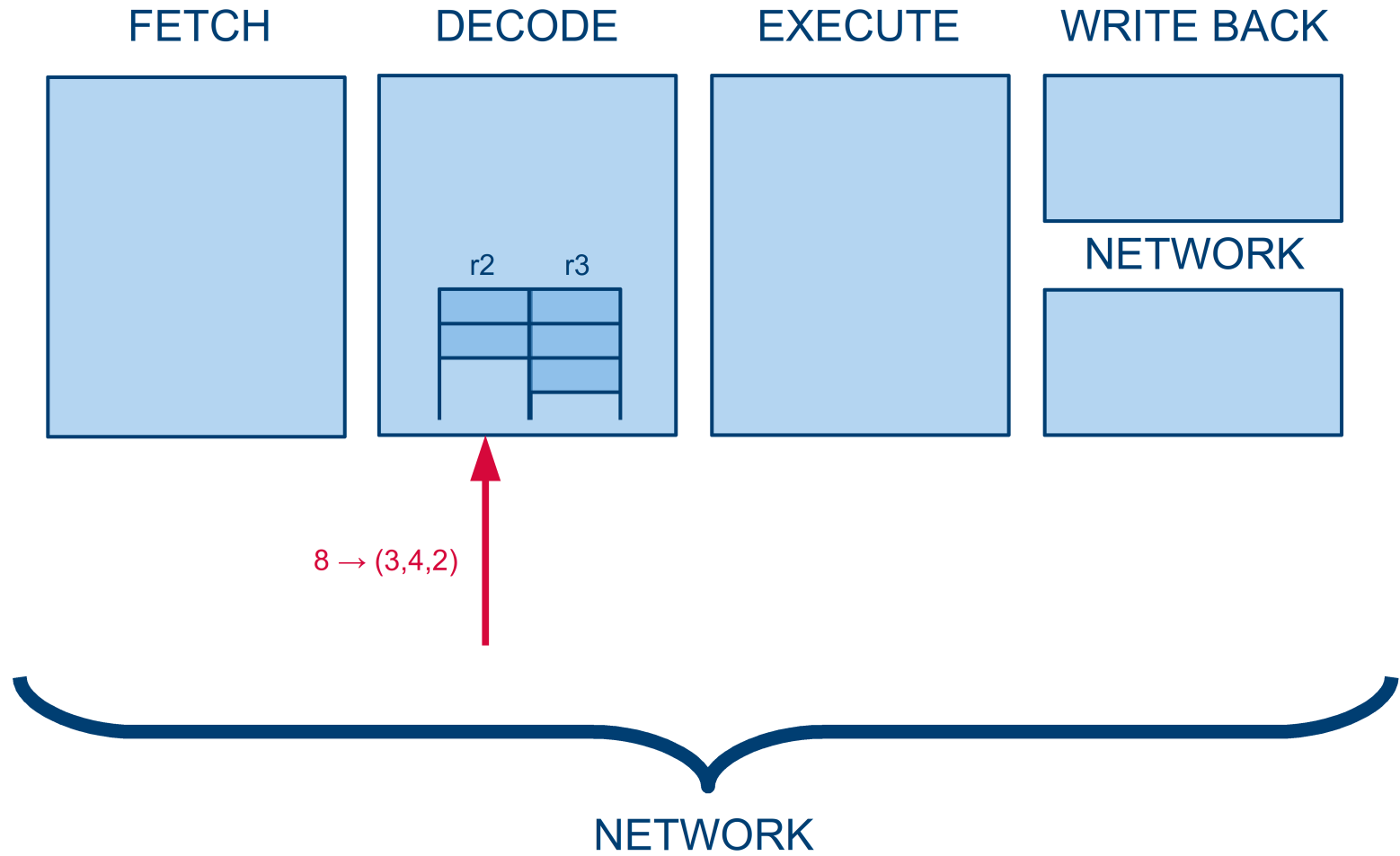→ 3

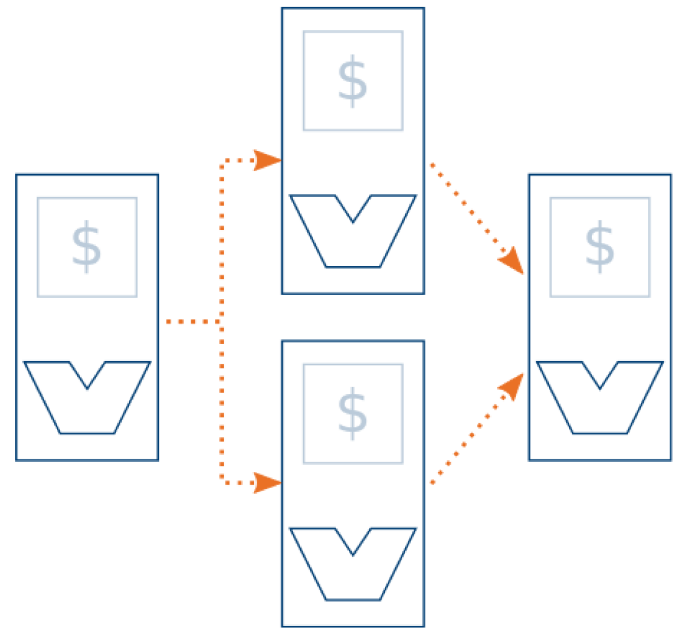(3,4,2)

NETWORK

NETWORK

UNIVERSITY OF CAMBRIDGE
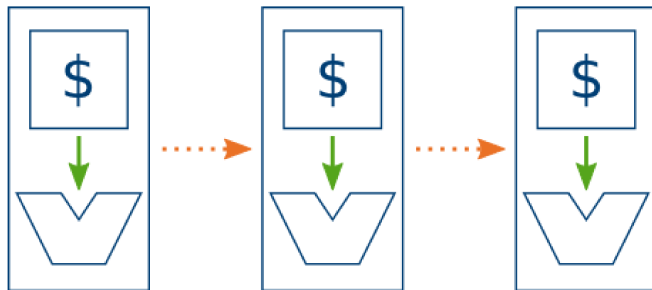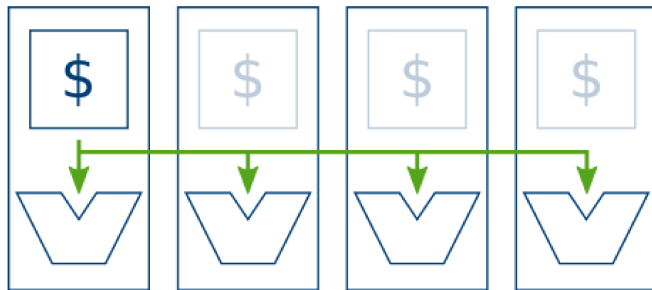
# Loki architecture: key features

# Loki architecture: key features

FETCH

DECODE

EXECUTE

WRITE BACK

r10 = 8

NETWORK

tile 3, core 4, channel 2

8 → (3,4,2)

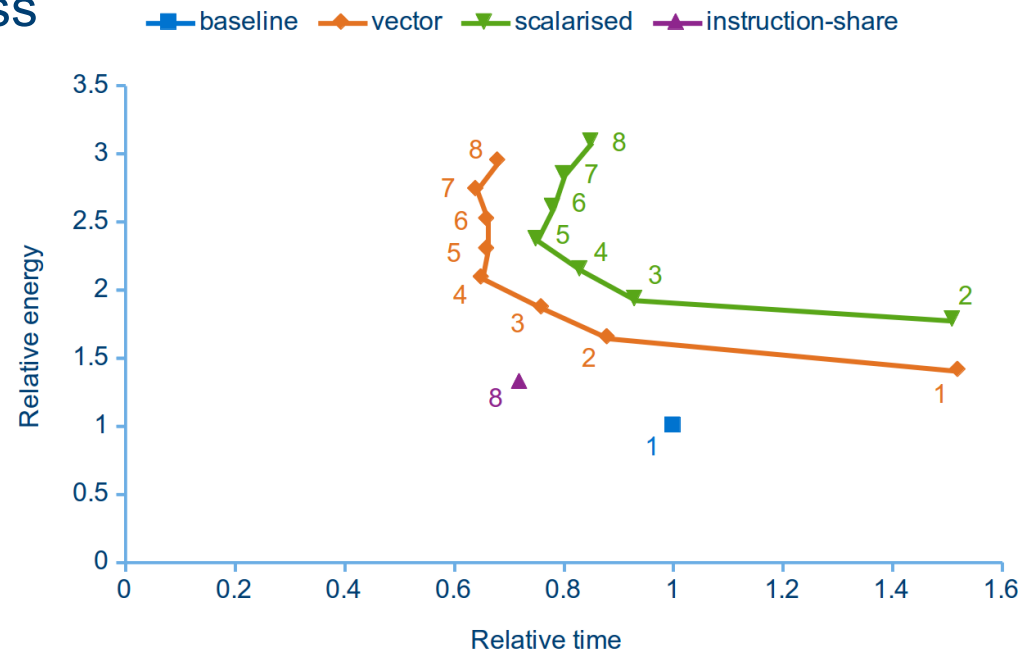NETWORK

UNIVERSITY OF CAMBRIDGE

# Loki architecture: key features
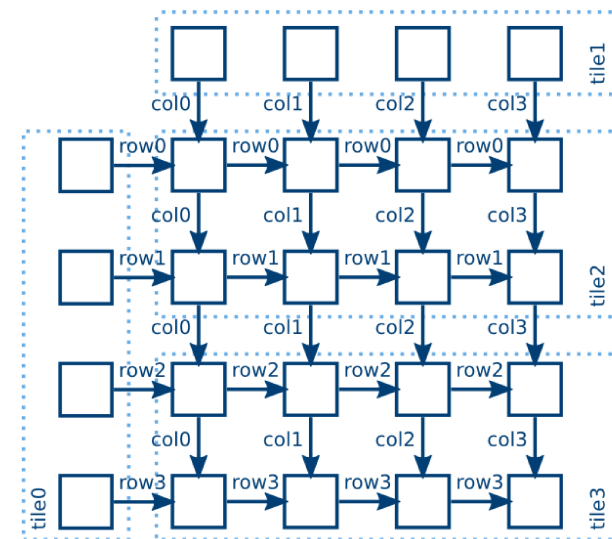
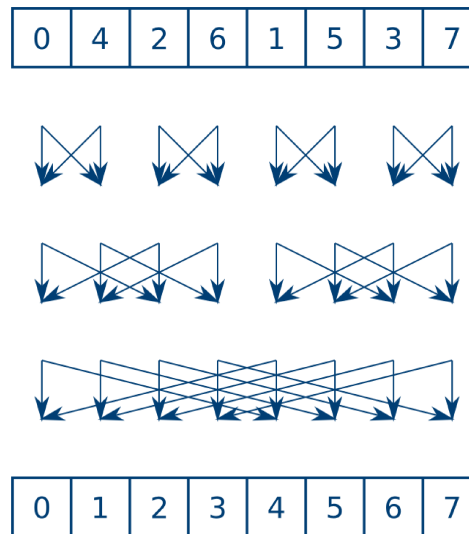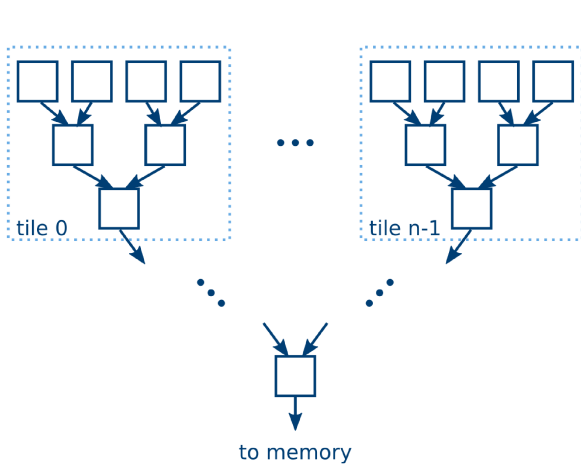# Previous work: execution patterns

# Previous work: flexibility

- Mappings within one tile (8 cores)

- Complementary to this work

- Now exploring mappings across multiple tiles

# This work: case studies

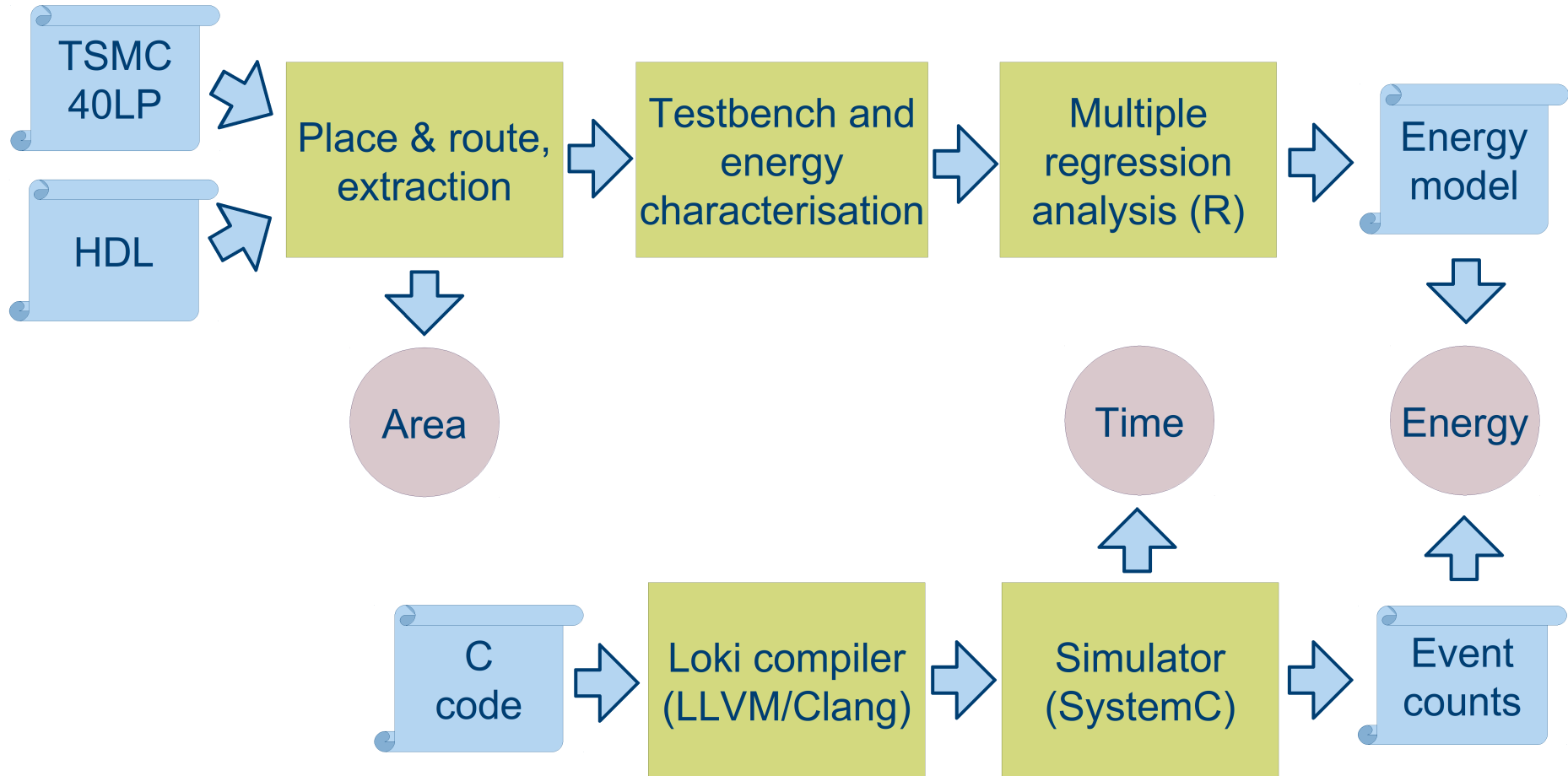- Three kernels: sorting, FFT, matrix multiplication

- Explore different mappings

- Choose energy-performance-area tradeoff

# Methodology

# Comparison

| | Loki | ARM1176 |
|---|---|---|
| Technology | 40nm LP | 40nm LP |
| Frequency/MHz | 435 | 700 |
| Area/mm² (core + L1) | ~0.125 | ~1 |
| Energy per instruction/pJ | ~5-20 | ~140 |
| Floating point support | Software only | HW double-precision |

# Sorting



tile 0

tile n-1

to memory

Critical path: 9 instructions

loki    arm1176

32  16  8  4  2  1

Energy/mJ

3
2.5
2
1.5
1
0.5
0

0.001    0.01    0.1    1

Time/s

Input: $2^{16}$ integers

UNIVERSITY OF CAMBRIDGE

# Sorting: reducing the bottleneck

Triggered instructions (Parashar et al, ISCA 2013)

```
compare:
    when (!p0 and !in1.finished and !in2.finished)
        p1 = in1.read() < in2.read()
        p0 = 1
send1:
    when (p0 and p1)
        out.write(in1.read())
        p0 = 0
send2:
    when (p0 and !p1)
        out.write(in2.read())
        p0 = 0
```

# ISA improvements

Features communication-centric architectures can borrow:

- Blocking communication

- Peek operation

- Data streams

- Triggered branch

Critical path: 4 instructions

```
while true
    if in1.hasData and in2.hasData and out.hasSpace
        a = in1.read()
        b = in2.read()

    if numRead1 < listSize and numRead2 < listSize
        if a < b
            out.write(a)
            numRead1++
        else
            out.write(b)
            numRead2++

    . . .
```
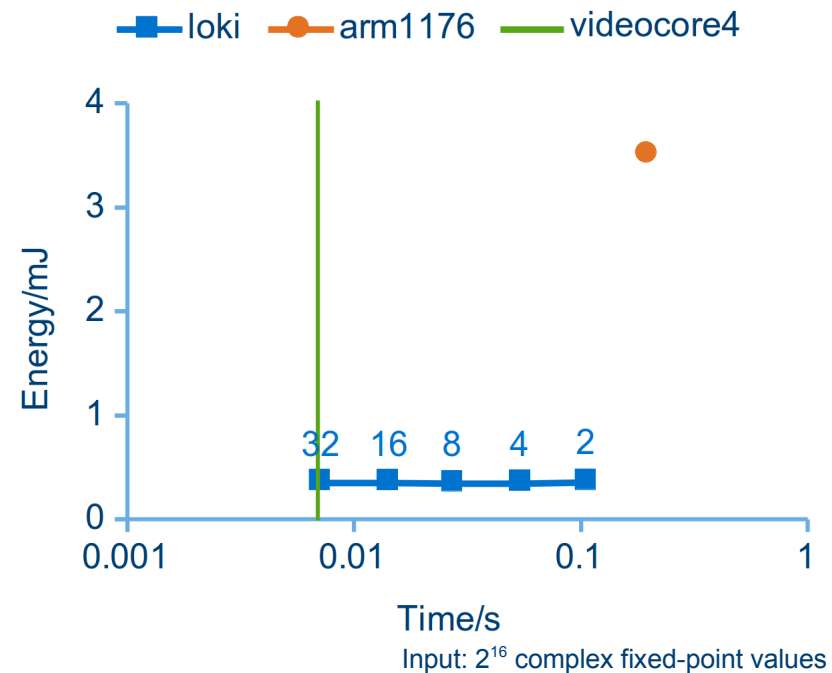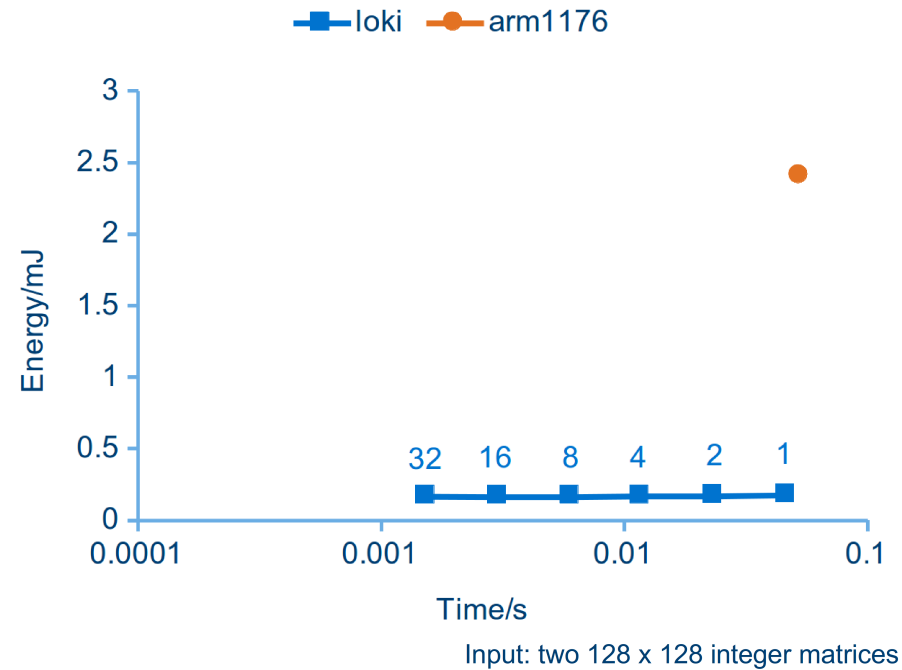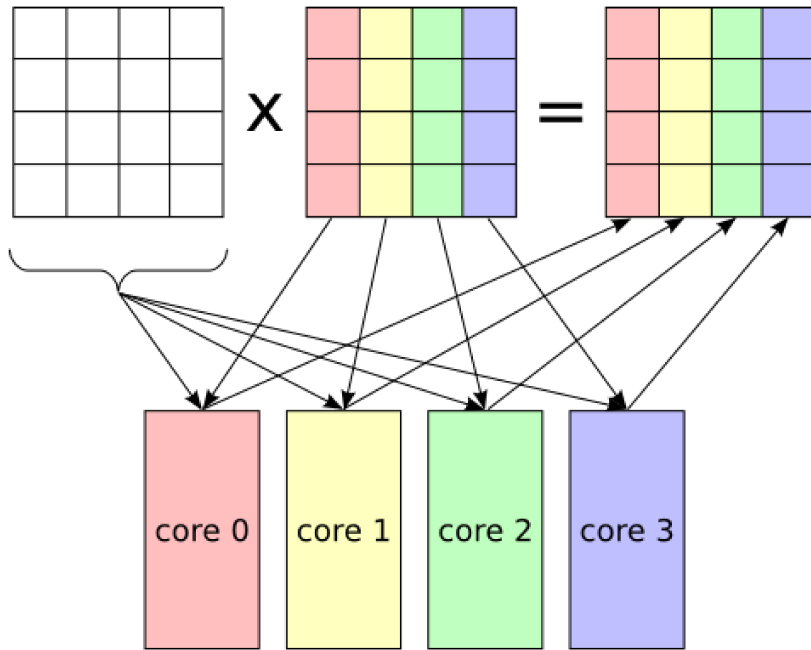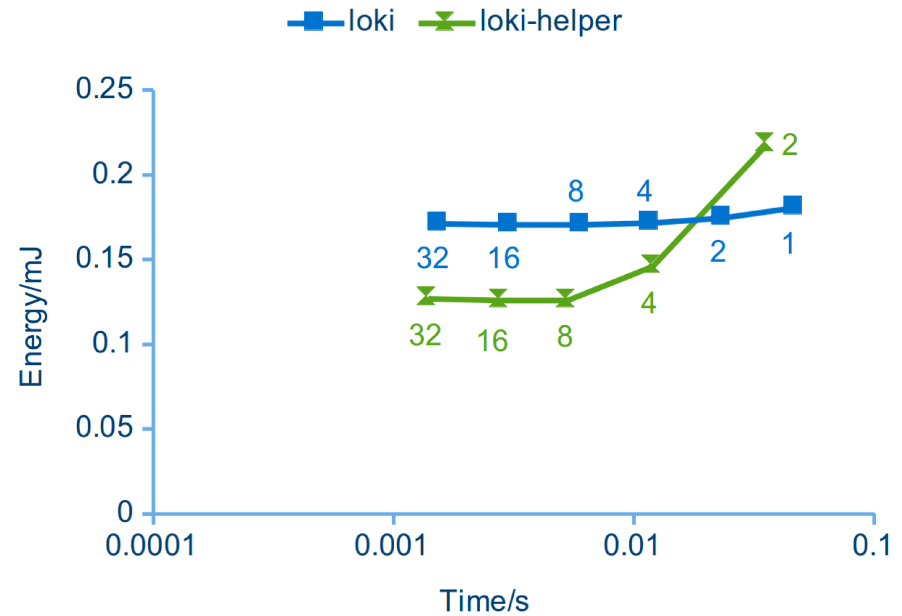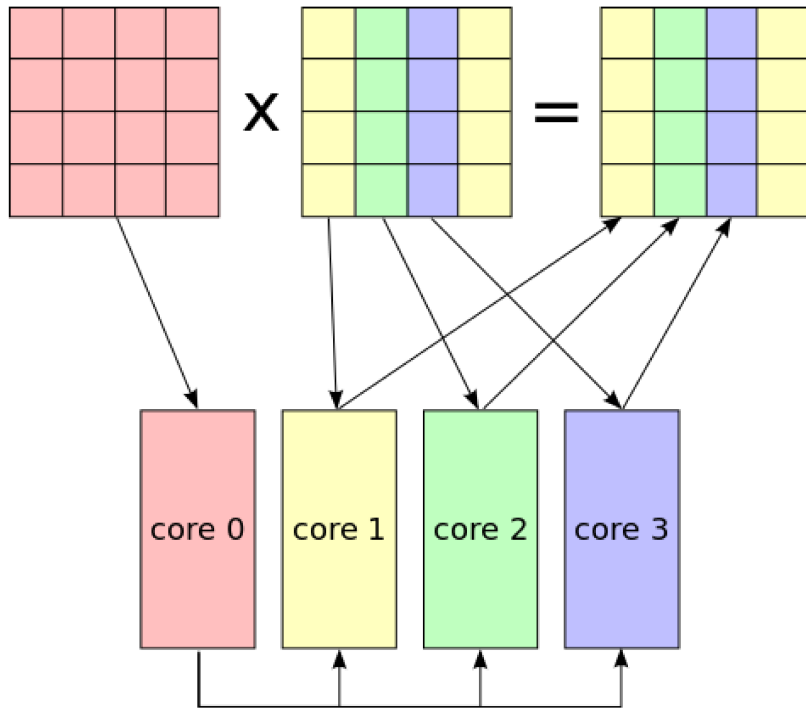
# FFT

# Software coherence

- Take the SWEL approach (Pugsley et al, PACT 2010)

    - Cache data in the lowest-level shared cache

    - Requires L1 bypass
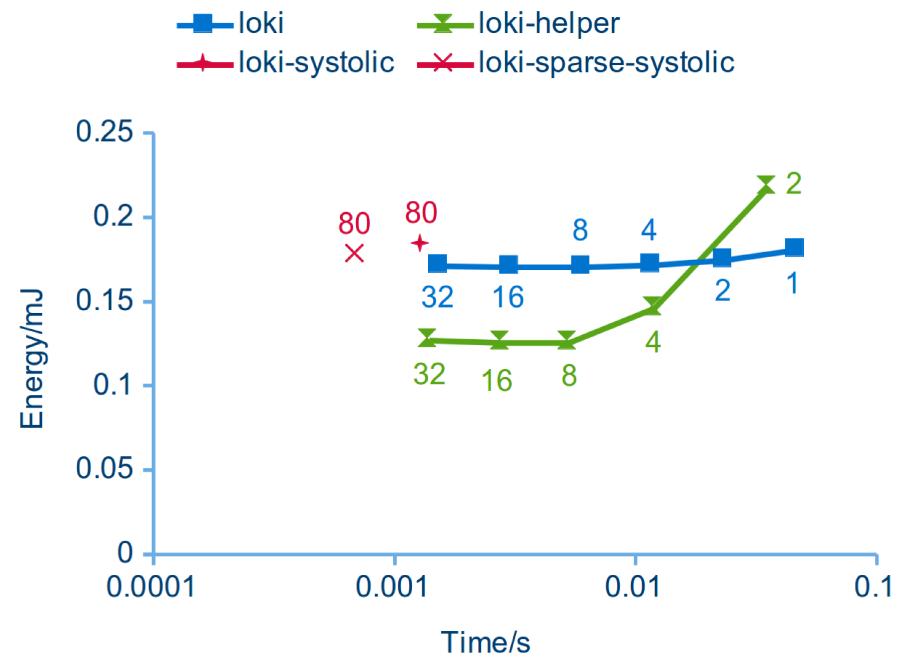
- Also looking into configurable cache hierarchies



Input: $2^{16}$ complex fixed-point values

# Matrix multiplication - vector



core 0    core 1    core 2    core 3
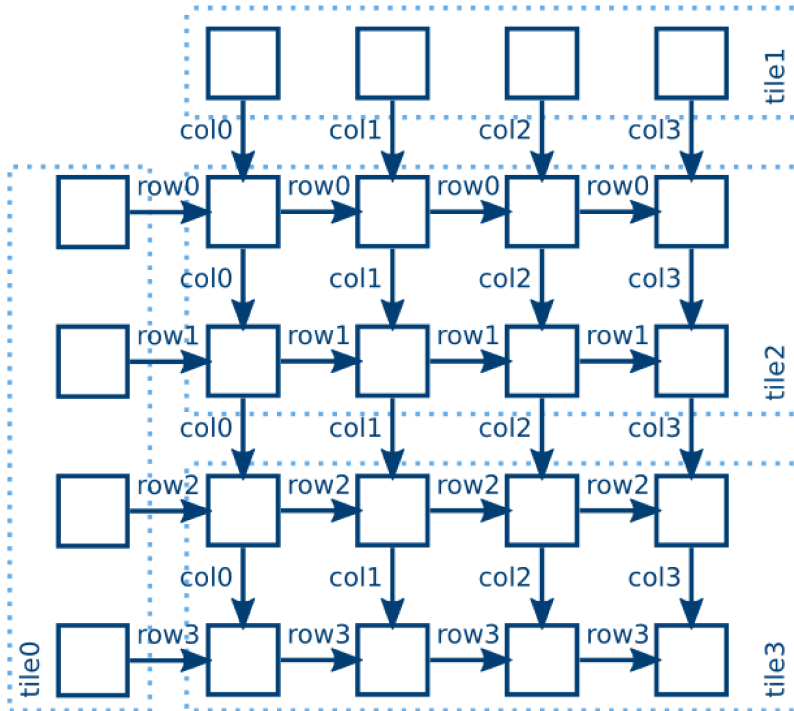
loki    arm1176

Energy/mJ

Time/s

Input: two 128 x 128 integer matrices

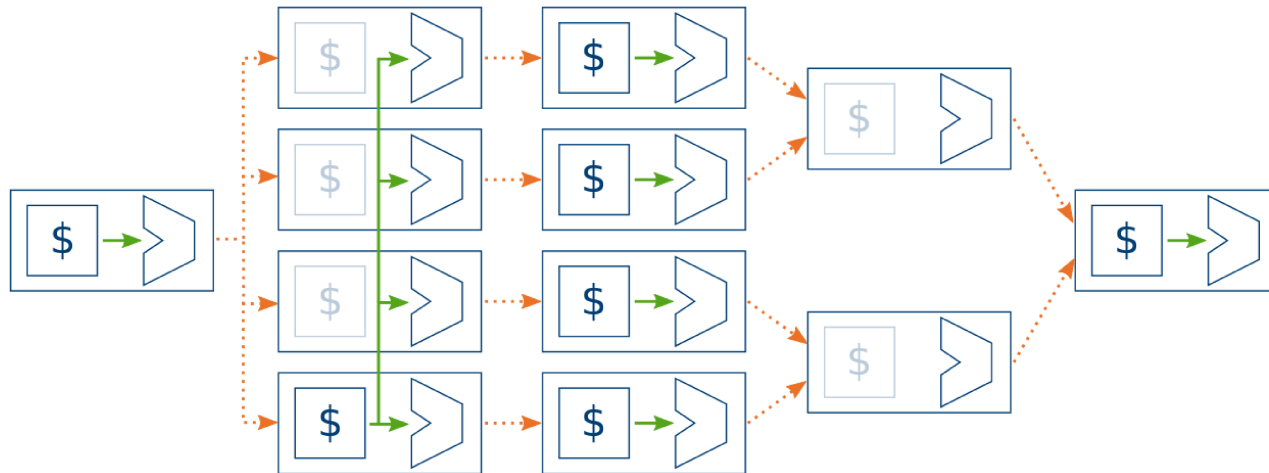# Matrix multiplication - scalarised

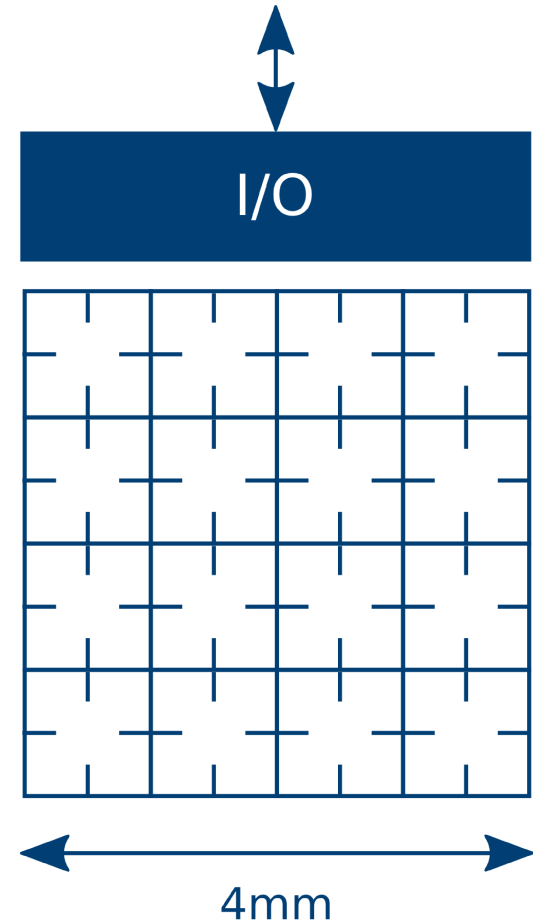# Matrix multiplication – systolic array

# Summary

- We can use large numbers of cores effectively

- Promising performance and energy figures

- A step towards a fully-homogeneous system

- Flexible communication is key

# What's next?

- Memory system

- Network protocols

- Compiler optimisations

- Programmable accelerators

- Test chip (2015)

  - 4 x 4 tiles x 8 cores/tile = 128 cores

  - Hoping to share development boards

- Operating system



UNIVERSITY OF
CAMBRIDGE

# Spatial computation on a homogeneous, many-core architecture

**Daniel Bates, Alex Bradbury, Andreas Koltes and Robert Mullins**

**Daniel.Bates@cl.cam.ac.uk**

**www.cl.cam.ac.uk/~rdm34/loki/**

**PRISM-2**