

AN INVESTIGATION INTO THE DESIGN OF MANYCORE ARCHITECTURES ON A MODERN FPGA

By: Thomas Bain (trb1g19@soton.ac.uk) Supervisors: David Thomas, Graeme Bragg



Introduction

ManyCore systems on FPGA have potential for highly parallel computational throughput with good energy efficiency. However, traditional processors now have core counts exceeding that of previous ManyCore work. To remain relevant new ManyCore architectures must allow scaling into the thousands of cores. Modern FPGAs provide new opportunities to support designs of these larger scales but also bring new challenges:

- Stacked Silicon Interconnect FPGAs have higher resource counts, but lower interconnect availability when crossing between chiplets.
- HBM can handle increased bandwidth demands of larger core counts, but the many channels requires different design approaches to DDR.

The table below shows a selection of previous FPGA ManyCore work and a modern x86 processor for comparison. Designs using the HBM have custom processors and interconnects to best exploit the bandwidth. Those also using SSI FPGAs have found their designs are constrained by the limited chiplet crossings.

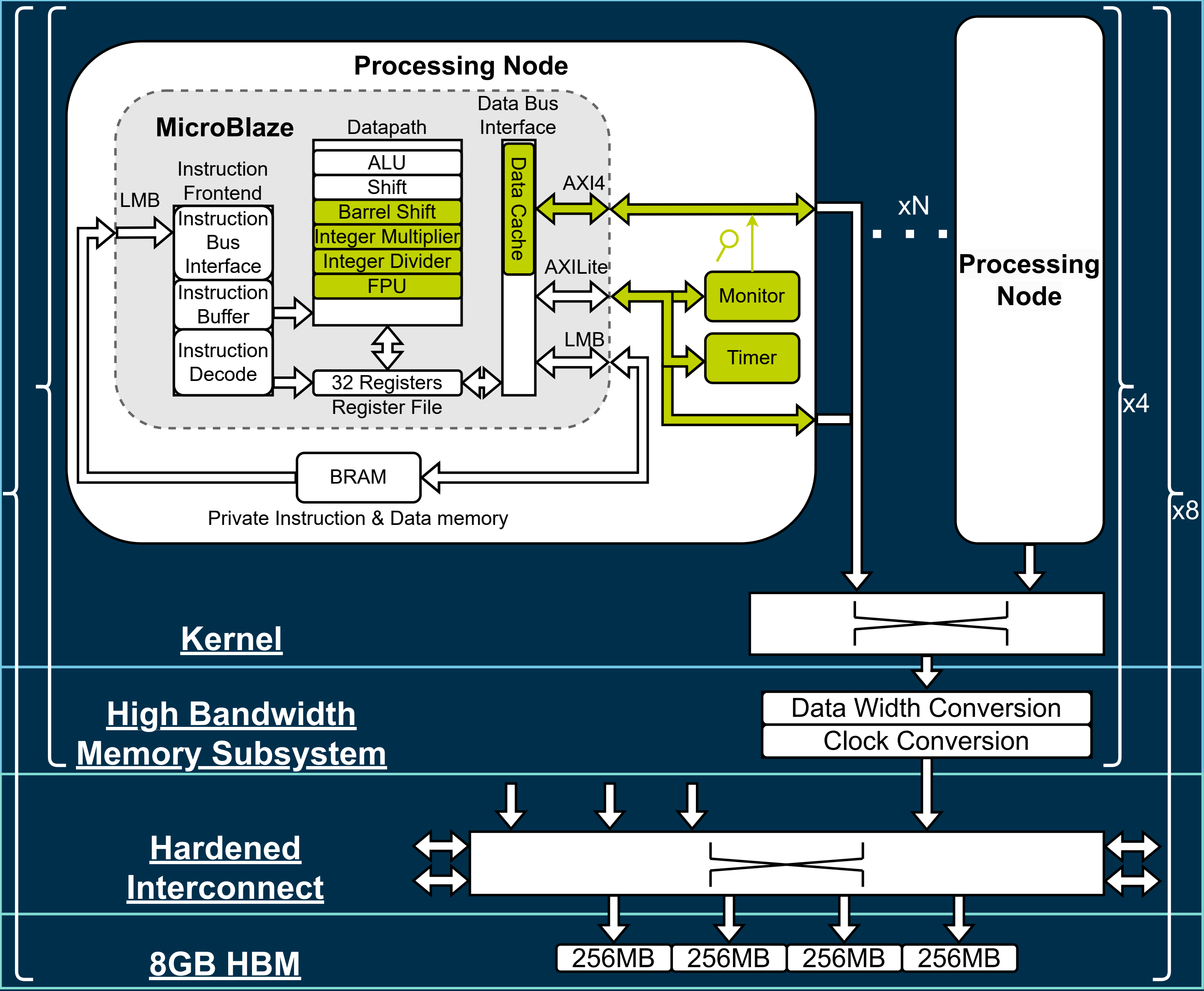
Design	ISA	Processor	Cores	Freq.(MHz)	DRAM	SSI	Year
DRAGON[1]	Custom	Custom	144	130	HBM	Yes	2021
Prakash et al.[2]	Custom	Custom	24	256	HBM	Yes	2022
AGILER[3]	RISCV	OS	32	120	DDR	Yes	2022
AMD EPYC 9754	x86_64	Custom	128	2250	DDR		2023

Initial Investigation

We performed an initial investigation to explore the strong and weak scaling performance of an architecture utilising the HBM when general purpose commercially available IP is used. The design has the following goals:

- Configurable number of processors to investigate scaling.
- Utilises all HBM channels to evaluate achievable bandwidth and latency.
- Spans all chiplets to explore the interconnect availability.
- Use COTS IP to reduce design effort and allow faster exploration.

The architecture is shown below, with a configurable number of processing nodes instantiated in the kernel providing 32 AXI ports for data access. Shaded features within the Processing Node are optional and exposed as parameters to the build script.

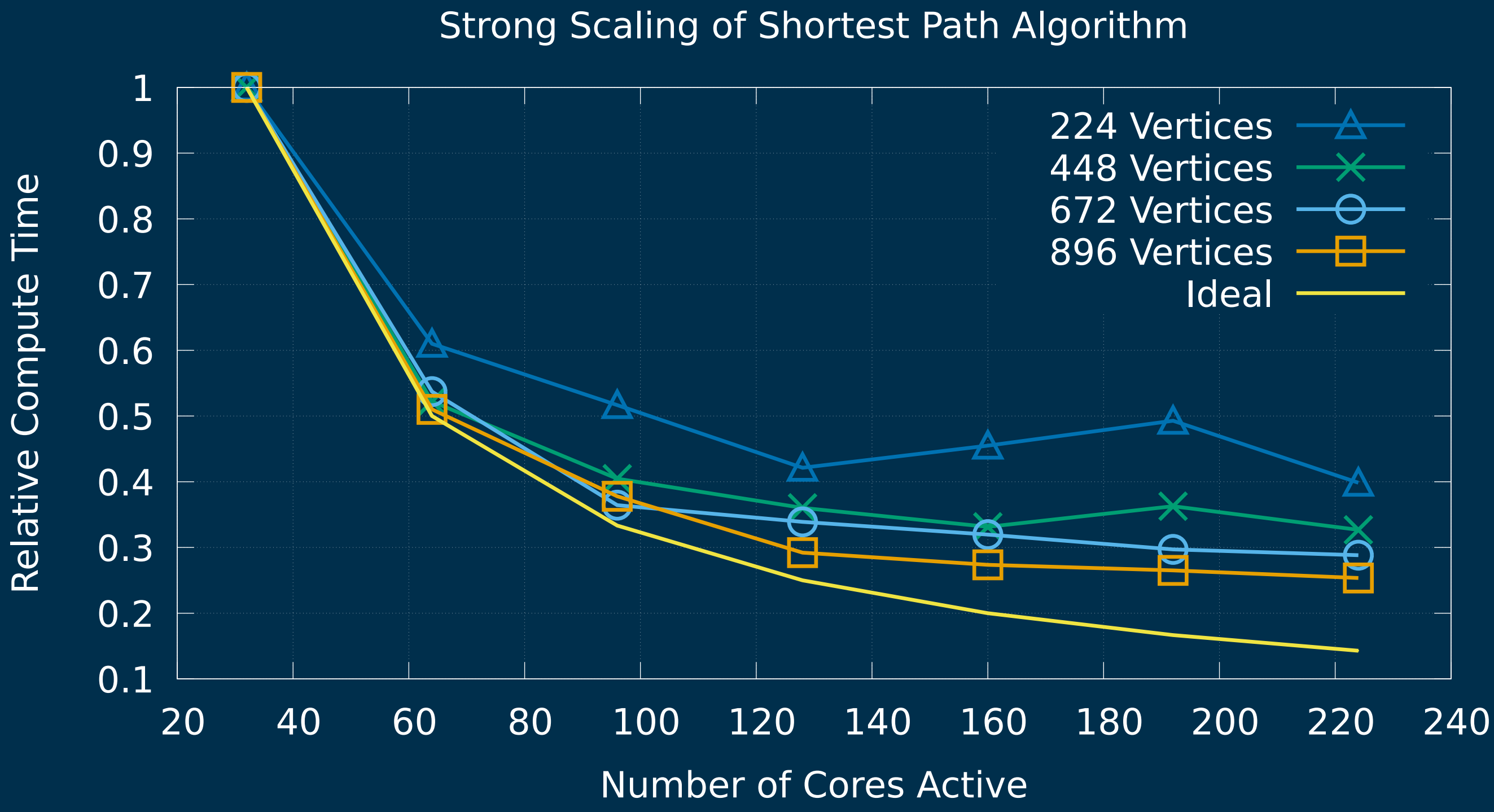
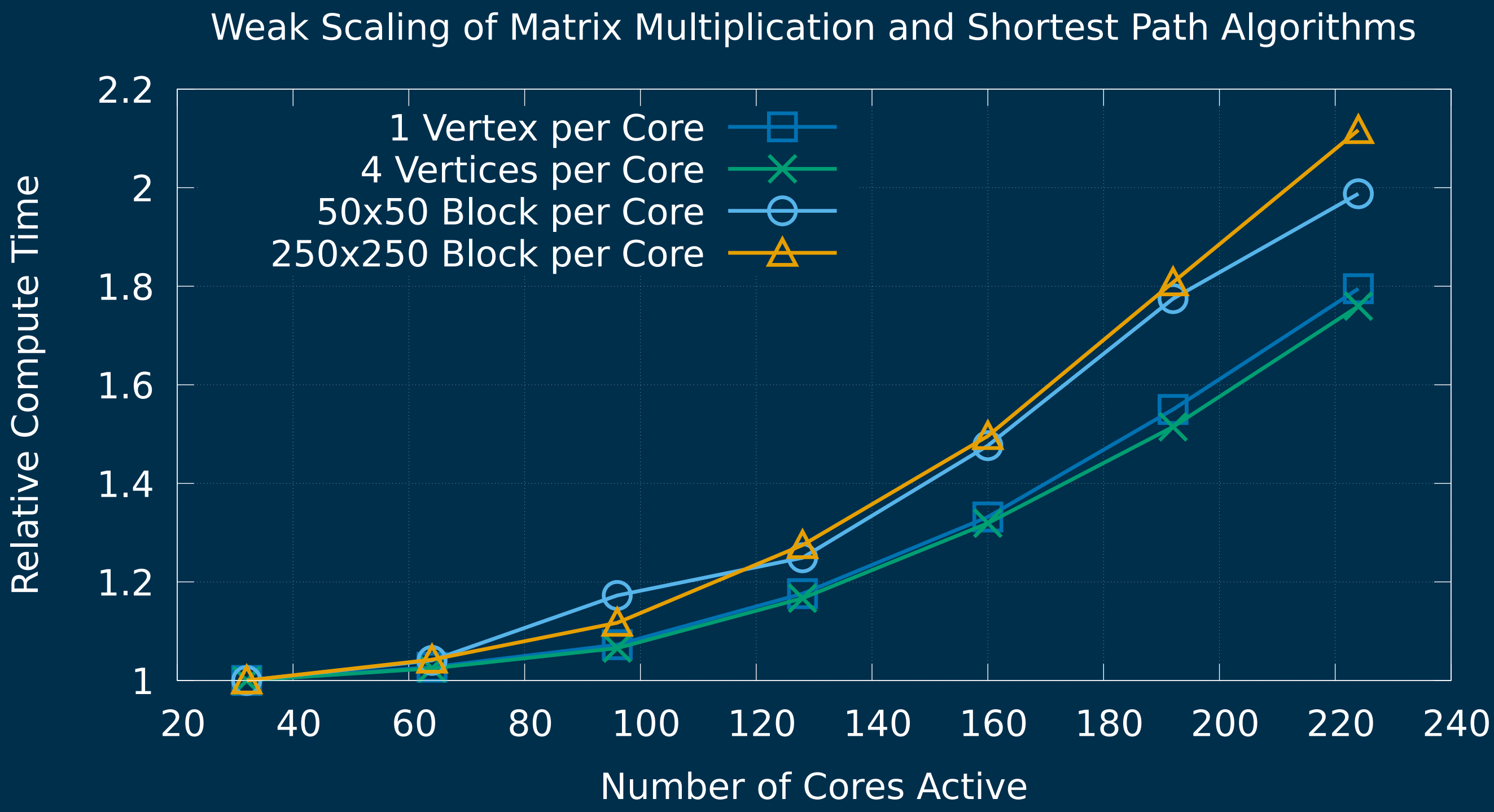


Initial Results

A 224 core configuration was implemented on a Xilinx U280 FPGA at 200MHz. Synthetic HBM benchmarks produced a total 21GB/s read BW from local channels, 360ns latency to local channel and 413ns latency from port 0 to channel 31.

The graphs below show the scaling performance of two multithreaded example applications:

- Blocked matrix multiplication
- Shortest path computation of a graph in adjacency list representation.



Future Work

Future work will continue exploring the optimisation of HBM bandwidth and access latency in a ManyCore system comprised of general purpose processing cores. Explorations in hardware take significant time and so an approach utilising simulation and mathematical models is being targeted. This model will be used to explore changes such as:

- Multi-Threaded processors to hide HBM latency.
- Higher level shared caches.
- Increasing data bus widths.
- Using DMA prefetchers.

We will aim to create a generalisable model that also allows exploration of how future FPGA developments could impact the problem with possible changes being:

- Increasing numbers of chiplets, chiplet crossing lines or HBM channels.
- 2D grid arrangements of chiplets.

These explorations aim to push the limits of what is possible with current technology and hope to find design methods that maximise the usage of available technology while managing the limitations. The model's generality will allow other researchers in this space to perform early analysis of expected performance before undertaking low level design and implementation.