

# Extending the RISC-V ISA for exploring advanced reconfigurable SIMD instructions

Philippou Papaphilippou , Paul H. J. Kelly, Wayne Luk  
Department of Computing, Imperial College London, UK  
{pp616,p.kelly,w.luk}@imperial.ac.uk

## ABSTRACT

This paper presents a novel, non-standard set of vector instruction types for exploring custom SIMD instructions in a softcore. The new types allow simultaneous access to a relatively high number of operands, reducing the instruction count where applicable. Additionally, a high-performance open-source RISC-V (RV32 IM) softcore is introduced, optimised for exploring custom SIMD instructions and streaming performance. By providing instruction templates for instruction development in HDL/Verilog, efficient FPGA-based instructions can be developed with few low-level lines of code. In order to improve custom SIMD instruction performance, the softcore's cache hierarchy is optimised for bandwidth, such as with very wide blocks for the last-level cache. The approach is demonstrated on example memory-intensive applications on an FPGA. Although the exploration is based on the softcore, the goal is to provide a means to experiment with advanced SIMD instructions which could be loaded in future CPUs that feature reconfigurable regions as custom instructions. Finally, we provide some insights on the challenges and effectiveness of such future micro-architectures.

## KEYWORDS

FPGAs, RISC-V, softcore, SIMD, cache hierarchy, reconfigurable, custom instructions, big data, streaming, sorting, prefix scan

## 1 INTRODUCTION

Modern general purpose processors (CPUs) support a wide range of single-instruction-multiple data (SIMD) instructions [15] as a way to accelerate applications that exhibit data-level parallelism. While there can be notable performance gains in certain applications [5], the instruction set extensions become bloated with overspecialised instructions [9], and sometimes it is difficult to express efficiently a parallel task using a fixed set of SIMD extensions [8].

As an alternative means to acceleration, FPGAs achieve unparalleled processing capabilities in specialised tasks [28, 41]. They have been getting attention for datacenter use, with numerous academic and industrial solutions focusing on processing big data [20, 32, 35]. However, the combination of specialisation and their placement in highly heterogeneous systems has high development and deployment costs.

FPGAs are often left behind in terms of main memory bandwidth, leading to a bandwidth bottleneck for big data accelerators [6, 11, 32]. Also, high-end FPGAs are usually based on non-uniform memory access (NUMA) systems, and the communication techniques are mostly inconvenient and expensive: First, PCIe, the most widely-used interconnection protocol has a high access latency [40], and is not appropriate for low-latency applications. Second, vendors promote high-level synthesis (HLS) tools in an effort to

abstract the complexity of communication, at the expense of enforcing generalised programming models. Last, unlike CPUs, any cache memory hierarchy, such as for improving random memory accesses, is usually implemented from scratch on FPGAs, and this complexity is often reflected in designs [11].

In combination with the openness of the RISC-V instruction set, and its support for custom instructions [43], this is a great time to start considering custom SIMD instructions on general purpose CPUs. Small FPGAs can be integrated, to implement custom instructions [9] and are demonstrated to improve the performance over existing extensions significantly [30, 31]. In the literature, the exploration of custom instructions on CPUs, and specifically SIMD, is rather limited, even though small reconfigurable regions working as instructions is a promising idea and possibly the future of acceleration.

In this paper, we present novel instruction types, Verilog templates and an open-source framework that are optimised for exploring custom SIMD instructions. In order to achieve high throughput on the provided softcore for streaming applications, the focus was given on the cache hierarchy and communication. The framework allows easy integration of custom vector instructions and evaluation in simulation and hardware. Finally, we evaluate examples of custom instructions and provide insights on introducing small FPGAs as execution units in CPUs. Our contributions are as follows:

- A set of experimental non-standard instruction types to enable optional access to a high number of registers for expressing complex tasks with fewer instructions (section 2.1).
- An open-source<sup>1</sup> softcore framework to evaluate novel SIMD custom instructions, and design choices to maximise streaming performance (sections 3,4.1).
- Defining a clean Verilog template for custom acceleration as an SIMD ISA extension in a lightweight softcore (section 2.2).
- A demonstration of the approach with novel SIMD instructions for sorting and prefix sum (section 4.3).

## 2 CUSTOM SIMD INSTRUCTIONS

In order to support adding and using SIMD instructions on the proposed softcore, we introduce (A) new instruction types that refer to the vector registers and (B) HDL (Verilog) code templates for implementing the instructions in hardware.

### 2.1 Vector instruction types

In addition to the RV32IM standards we propose two additional non-standard instruction types for supporting the custom SIMD instructions. Originally there are 4 main instruction types: R, I, S/B, U/J in RV32I, that define the format of the immediate field and the operands. The proposed instruction types I' and S' are variations of the I and S types respectively.

The official draft RISC-V "V" vector extension has not been followed in this work, as it seems to target high-end/ hardened processors. For example, it requires a fixed number of 32 vector registers and features hundreds of instructions [2], also reducing the number of registers an instruction can access. For our use case, this would be contrary to the idea of having small reconfigurable regions as instructions, rather than supporting hundreds of intrinsics.

The modification repurposes the space used for the immediate field for the vector register names. **Up to 6 registers** (vector and non-vector) can be accessed by the same instruction, reducing the instruction count for complex or stateful (through registers) instructions. The use of the immediate field for the vector registers was also convenient for minimal interference and modification to the RISC-V GNU toolchain in GCC's binutils for supporting inline assembly of custom instructions. There is currently no official tool support for vector registers due to the draft status of the vector extension. We opted to use the opcodes dedicated to custom applications for all custom vector instructions, for inline assembly.

Type I' provides access to the register operands of the I-type, that is one source and one destination 32-bit register. It also provides access to two source (*vrs1* and *vrs2*) and two destination (*vr1* and *vr2*) vector registers. Type S' exchanges the space used by *vrs2* and *vr2* for access to an additional 32-bit source register *rs2*. The latter would be useful, for example, for load and store instructions, for breaking loop indexes into two registers and potentially reducing the instruction count in some loops, as in other ISAs. The proposed variations are summarised in Figure 1.

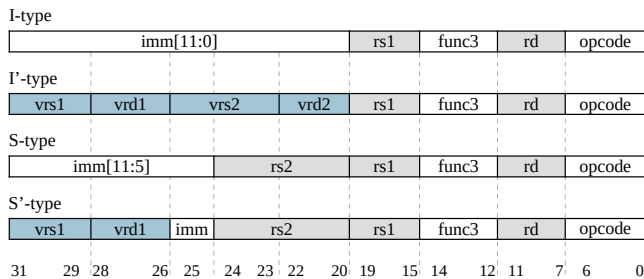


Figure 1: Two variations of the I and S instruction types

As shown in Figure 1, the fields for each vector register are three bits wide, which sets the maximum number of vector registers to 8. Vector register 0 is conveniently assigned to a **constant value of 0**, similarly to the 32-bit base registers. It is useful for the proposed many-register instructions types, because not all register operands may need to be accessed at once, allowing different combinations of operands using the same type. In software, this can be achieved with aliasing the unused operands with register 0, and was **not** supported in "V", as vector 0 originally represents a register.

An interesting feature in the "V" specification is the ability to chain vector instructions, hence the need for high number of registers. Since our solution enables custom instruction **pipelines of arbitrary length**, a lower number of registers was considered satisfactory, as the **need for chaining is minimised**.

## 2.2 Instruction templates

The custom instruction templates are placeholder modules inside the softcore codebase, for adding user code for specialised SIMD instruction implementations. Algorithm 1 is written in Verilog and shows part of the template for I'-type instructions, plus an example user code, marked in yellow. On each cycle, the instruction module also accepts the destination register names (*rd*, *vr1* and *vr2*) to provide them later, when the data result is ready, after the specified pipeline length (*c1\_cycles*). In this way, with a pipelined implementation the module can process **multiple calls** one after another. Blocking instructions are also supported with minor modification.

Algorithm 1 Verilog template for I'-type, with user-provided code highlighted in yellow

```

1  `define XLEN 32 // 32-bit base registers
2  `define VLEN 128 // 128-bit vector registers
3  `define c1_cycles 3 // pipeline length of c1 custom instruction
4
5  module c1 [...]
6      input clk, reset, in_valid; // valid bit for input
7
8      // Destination register names (1 base and 2 vector)
9      input [4:0] rd; input [2:0] vr1, vr2;
10
11     input [`XLEN-1:0] in_data; // 32-bit input
12     input [`VLEN-1:0] in_vdata1, in_vdata2; // 128-bit input
13
14     // (Delayed) output valid bit and output register names
15     output out_v;
16     output [4:0] out_rd; output [2:0] out_vr1, out_vr2;
17
18     output [`XLEN-1:0] out_data; // 32-bit output
19     output [`VLEN-1:0] out_vdata1, out_vdata2; // 128-bit output
20
21     // Shift register logic to delay rd, vr1, vr2 and in_valid
22     // by c1_cycles, to out_rd, out_vr1, out_vr2 and out_v resp.
23     [...]
24     /// User code ///
25     wire [31:0] net [15:0]; // 32-bit wires inside the sorting network
26     // On pos. edge: a, b -> min(a,b), max(a,b)
27     CAS cas0(clk, net[0], net[1], net[4], net[5]);
28     CAS cas1(clk, net[2], net[3], net[6], net[7]);
29
30     CAS cas2(clk, net[4], net[7], net[8], net[11]);
31     CAS cas3(clk, net[5], net[6], net[9], net[10]);
32
33     CAS cas4(clk, net[8], net[9], net[12], net[13]);
34     CAS cas5(clk, net[10], net[11], net[14], net[15]);
35
36     // Assigning input and output to wires in the sorting network
37     // (only using 1 input and 1 output reg. for this instruction)
38     for (i=0; i<4; i=i+1) assign net[i]=in_vdata1[32*(i+1)-1.:32];
39     assign out_vdata1={net[12], net[13], net[14], net[15]};
40
41 endmodule

```

The example instruction implementation in Algorithm 1 is a bitonic sorter of 4 inputs. Such sorting networks are parallel and pipelinable algorithms for sorting a list of *N* values. In each parallel step there is a number of compare-and-swap (CAS) units, that collectively sort the entire input list, as the input moves along the network. The odd-even merge sorter and the bitonic sorter [4] are two similar sorting network topologies, both consisting of  $\Theta(\log^2(N))$  parallel steps. For a vector register width of 128-bit, this bitonic sorter sorts four 32-bit values in 3 cycles.

The template of the S'-type instructions is similar to Algorithm 1, but with an interface that reflects the correct operands (2 base and 1 vector registers for input and 1 base and 1 vector for output). One S' type instruction for loading and storing VLEN-sized vectors (*c0\_lv* and *c0\_sv* respectively) is provided by default.

### 3 SOFTCORE MICROARCHITECTURE

The proposed softcore supports the RISC-V 32-bit base integer instruction set (RV32I v. 2.1), plus the “M” extension for integer multiplication and division [43]. The novel features of our approach include a series of design choices to: (1) enable high-performance for custom vector instructions and streaming applications; (2) allow efficient implementation on recent FPGAs by enhancing the block RAM (BRAM) organisation and the behaviour of the inter-chip communication.

#### 3.1 Cache hierarchy optimisations

On the first level, there is an instruction cache (IL1) and a data cache (DL1), and on the second level there is a unified last-level cache (LLC). LLC responds to requests from both IL1 and DL1. It communicates in bursts to DRAM through an interconnect such as AXI. It resembles a modified Harvard architecture, as the address space is common between data and instructions. Figure 2 provides a high-level example for the data communication throughout the cache hierarchy.

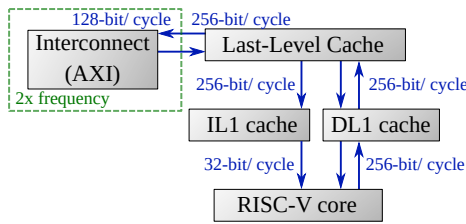


Figure 2: Data movement in an example cache configuration

All caches use the writeback policy, with the exception of the IL1 cache, where writing is not needed. This is achieved by storing a dirty bit alongside each stored block, to acknowledge modification. While DL1 and LLC are set-associative caches, the IL1 is direct-mapped for fast lookup of the next instruction, to avoid a stall on instruction hits.

The LLC is implemented in block memory (BRAM), to accommodate the high capacity of the last level. The IL1 is implemented in registers for a reduced latency, in order to provide the successor instruction immediately on the next cycle (and avoiding being the critical path) on hits. The DL1 is implemented in BRAM by default, although changing the directive to registers yields similar performance and utilisation results, due to its relatively low size.

At the set-associative caches (DL1 and LLC), each block can be allocated into multiple possible ways, represented by different parallel block RAM sections. The block replacement policy for these caches is not-recently-used (NRU). It uses one bit of meta-information per block [42], but closely resembles the Least-Recently-Used (LRU) eviction policy. The choice of a replacement policy can be crucial to the performance of streaming applications, due to the wide data blocks and the reduced cache space on the FPGA. For instance, a random policy would stagnate the bandwidth for memory copying (`memcpy()`), when the source and destination are aligned.

A series of design choices are presented for optimising the performance and applicability for our purposes.

**3.1.1 Level-1 block size.** One optimisation is to set the block size of the DL1 to be equal to the **vector register width**, such as 256 bits. This is because a wider block size would require an additional read on each write, from the cache of the higher level, so that the entire block becomes valid. When the data are from vector registers and are properly aligned, there is no need to wait for fetching that block on a write miss, because the whole block will contain new information.

The IL1 uses the same block size for easier arbitration between DL1 blocks, at the cache of the higher-level (LLC). Additionally, since IL1 is direct-mapped, using a wider-block than 32-bits is also beneficial to performance, as it can also be seen as a natural way of prefetching.

**3.1.2 LLC block size.** An important feature for increased streaming performance is **very wide blocks** for LLC, such as 8192-bit wide. This is in contrary to today’s CPUs with a 512-bit (64-byte) block size. The idea is that on write-back to/ fetch from main memory, it achieves a higher speed because of longer bursts. Longer bursts are shown to have significant impact on the overall throughput, such as in heterogeneous systems with AXI [22], and this is especially useful for streaming. Associating entire LLC blocks with bursts was a convenient and practical organisation choice, because of interconnect protocol limitations, such as for not crossing the 4KB address boundary in AXI [14].

**3.1.3 LLC strobe functionality.** A naive implementation of the LLC in BRAM, would be to read the (wide) blocks in their entirety in a single cycle, as in the DL1 case. However, BRAM is organised in chunks of certain width and length, such as 36-bit wide. With a LLC of just a single wide-enough block, the BRAM capacity of the FPGA can be exceeded, or stagnate timing performance. For this reason, the proposed LLC stores blocks in consecutive BRAM **locations of narrower size**. There is an internal notion of sets that corresponds to the address of the block memory, where each requested data can reside. The tag array only stores the tags of entire blocks.

There is no overhead in access latency by using sub-blocks, as it still takes a single cycle to read an I/DL1-sized block from LLC. Another advantage of this technique is that, on fetch, the requested I/DL1 block can be provided before a read burst from DRAM finishes, since the LLC blocks are stored progressively.

**3.1.4 Doubling the frequency of the interconnect.** In contrast to the timing characteristics of this softcore, as well with other well-known softcores [12], the operating frequency of the interconnect on FPGAs can be relatively much higher [17]. Given that the port data widths are rather narrow (e.g. 128 bits/ cycle), this directly impacts the throughput for streaming applications. This optional optimisation involves setting double rate for the interconnect, to **emulate double data width** by fetching or writing twice per cycle, and saturate [22] the bandwidth more easily (see Figure 2).

#### 3.2 Main core

The core has a single pipeline stage, even though more advanced instructions such as pipelined vector instructions have their own pipeline. Almost all instructions in RV32I consume 1 cycle and the result is available on the immediately next cycle. In practice, this has a similar effect to operand forwarding in pipelined processors,

as consecutive dependent instructions are executed sequentially without stalls.

The load and store instructions are handled by the cache system independently. On a data cache hit there is a latency of 3 cycles until the dependent command gets executed. The 3 cycles can be seen as a small pipeline with one cycle for memory access, one for fetching the data and one for updating the registers. This effectively yields a latency of 2 cycles for cache hits, when the next instruction is data-dependent on the load, as the execution is in-order.

Having a **single pipeline stage**, so that most instructions complete in a single cycle, is useful for simplifying the dependency checks. The output of simple instructions such as add, addi, etc. is not tracked for dependencies. Of course, there are alternative approaches, but the current implementation mapped well in our evaluation platform and facilitated the SIMD functionality rather efficiently.

The implementation of the SIMD instructions follow the templates of section 2), that allow a variable pipeline length, abstracted through a *ready* signal for when the result is available. Apart from the 32 base 32-bit registers (as per RV32I), there are up to 8 VLEN-wide registers, such as 256-bit-wide for the SIMD instructions. In both sets of registers, the register 0 is driven by the constant 0.

## 4 EVALUATION

The exploration is divided in three parts according to the outcomes of each set of experiments: (4.1) justifies important design choices related to streaming performance, (4.2) shows that the performance is still acceptable when no SIMD instruction is used and (4.3) explores the behaviour and efficiency of example novel custom SIMD instructions.

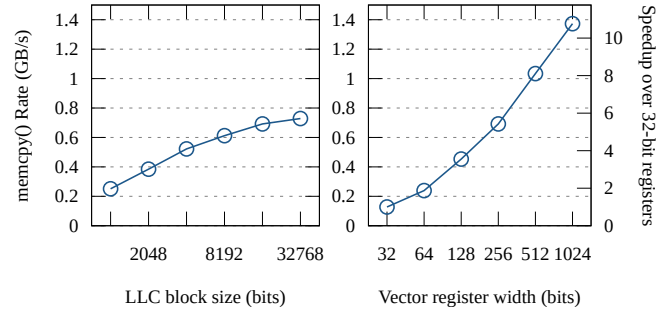
The evaluation platform is Ultra96, which features the Xilinx UltraScale+ ZU3EG device. The FPGA on the device shares the same 2GB DDR4 main memory with the 4 ARM cores. The ARM cores run Linux, but the kernel address space is manually configured to end at the 1GB mark, so that the other 1GB is dedicated to the FPGA, that includes the softcore.

### 4.1 Design Space Exploration

The target application is memory copying (*memcpy()*), as its performance is (indirectly) detrimental to big data processing and related evaluations have a long history in HPC applications [26]. *memcpy()* here is manually implemented with the custom instructions for load vector and store vector, instead of a library implementation using base registers. The data length is 256 MiB, in order to surpass the cache sizes.

Figure 3 (right) illustrates the impact of the vector register size on *memcpy()*. The 1024-bit softcore achieved a *memcpy()* rate of 1.37 GB/s. Though, we opt for 256-bit (VLEN) registers, with a rate of 0.69 GB/s, as 512-bit and beyond seemed more challenging to route efficiently when incorporating more complex custom instructions. These designs used a 16384-bit-wide LLC block.

One other experiment (Figure 3 left) measures the impact of the block size in last-level cache (LLC). Wider LLC blocks seem to be a considerable contributor to memory performance, as they relate to the burst size. The improvement starts to plateau after longer bursts at around 8192 bits. All implementations reached timing closure



**Figure 3: Memcpy() read and write (bidirectional) throughput for different last-level cache block sizes (left) and vector register widths (right)**

for a frequency 150 MHz, except the 1024-bit configuration that was clocked at 125 MHz. Table 1 summarises the selected baseline configuration for the remainder of the evaluation.

**Table 1: Selected configuration**

IL1		DL1		LLC		VLEN	$f_{max}$
sets	block (bits)	sets	ways block (bits)	sets	ways block (bits)	(bits)	(MHz)
64	256 (=2KiB)	32	4 (=4KiB)	32	4 16384 (=256KiB)	256	150

### 4.2 Performance as a RV32IM core

In order to show that there is no significant bottleneck when compared with other non-SIMD cores, we overview some other works with a similar specification. Table 2 presents common benchmark metrics alongside previously reported numbers using FPGAs. Note that this is not for direct comparison, as each work used a different FPGA family, cache configuration and compilation environment.

**Table 2: Indicative comparison ignoring SIMD**

	DMIPS/MHz	Coremark/MHz	$f_{max}$	FPGA architecture
RVCOREP/radix-4[18]	1.25	1.69	169	Xilinx Artix-7
RVCOREP/DSP[18]	1.4	2.33	169	Xilinx Artix-7
PicoRV32[44]	0.52	N/A	N/A	(simulation)
RSD/hdiv[23]	2.04	N/A	95	Zynq
BOOM/hdiv [3, 23]	1.06	N/A	76	Zynq
Taiga[12, 25]	>1	2.53	~200	Xilinx Virtex-7
<i>This work</i>	1.47	2.26	150	Zynq UltraScale+

Additionally, the performance of our proposal is measured for memory-intensive situations, without the use of SIMD. STREAM [26] is an established benchmark suite measuring the memory performance, especially in HPC. Figure 4 shows the obtained throughput in MB/s for each of the 4 kernels.

On the same FPGA, we place PicoRV32 [44], as a drop-in replacement that supports AXI (Lite). Although it was not designed for performance, it achieves high operating frequencies (300 MHz in our platform), partly mitigating for its low IPC [12]. It does not have a cache, although this does not directly impact memory bandwidth, as the data reuse is practically zero. (The steps in Figure 4 are from

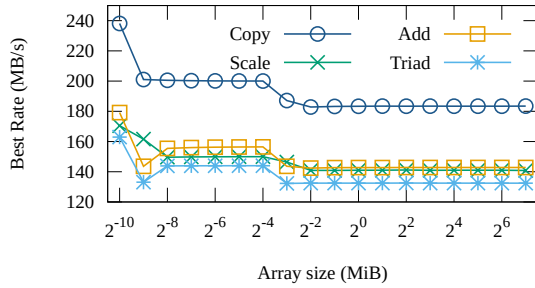


Figure 4: Adapted STREAM benchmark results, no SIMD

reusing data from the initialisation). The results of PicoRV32 were 4.8, 3.6, 4.4 and 4 MB/s for Copy, Scale, Add and Triad consistently across the array size range. This makes our approach **38x** faster for Copy at 183.4 MB/s, or **144x** faster if we consider the 256-bit *memcpy()* performance. This also highlights the importance of optimising communication for streaming applications.

### 4.3 Custom SIMD instruction use cases

**4.3.1 Sorting (32-bit integers).** Sorting is a widely applicable big data application. Existing SIMD intrinsic solutions are based on algorithms such as sorting networks [4], radix sort [16], mergesort [8], quicksort [5] and combinations.

The algorithm of our solution is merge sort, with the help of sorting networks for introducing parallelism. Sorting networks were adapted for both software [5, 8] and hardware [10, 21, 34, 39] solutions for sorting arbitrarily long input.

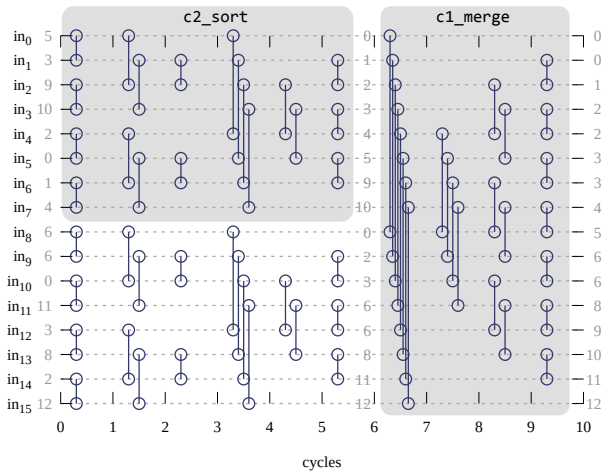


Figure 5: Two new custom instructions based on odd-even mergesort

In order to accelerate sorting in the softcore for arbitrary-sized input, a sorting network is used first to sort the entire list first in small chunks, as in [5]. Then, a traditional recursive merge sort approach is performed, but instead of merging each two sublists by comparing one element by one, it uses a parallel merge block. The merge block (the last  $\log_2(N)$  layers of odd-even mergesort) is to merge two already-sorted lists together, as demonstrated in

a numerical example in Figure 5. In our implementation, we add one more stage in the beginning to enable merging arbitrarily long lists progressively, and the algorithm is inline with the intrinsics merge algorithm [8].

For brevity, we only elaborate on the sort-in-chunks loop. Figure 6 illustrates the instruction start and end times for this loop during a simulated run. From this figure we can observe the pipelining effect, as two instances of *c2\_sort* take place simultaneously, to sort two octuples. The second sort is shifted by two cycles, as it still waits its operand *v2* from the second load. Then, the merge instruction (*c1\_merge*) merges the registers *v1* and *v2* and stores the upper and lower half back to *v1* and *v2* respectively, for sorted chunks of size 16.

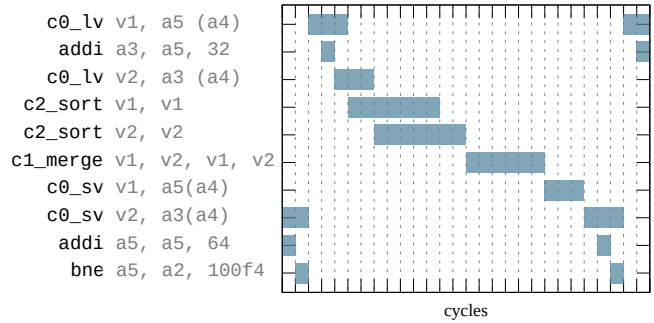


Figure 6: Iteration cycles from the sorting-in-chunks loop

The performance of the resulting mergesort function is compared against non-vectorised code on the softcore, running at 150 MHz, as well as on the ARM A53 core, running at 1.2 GHz. The baseline is *qsort()* from C’s standard library. The obtained speedup is **12.1x** and 1.8 times over the *qsort()* on the softcore and ARM respectively, for 64 MiB random input. Comparison with more optimised code such as multi-threaded NEON-based for ARM, as well as other SIMD algorithms [5, 33] would also be possible, but are out of the scope of this work.

**4.3.2 Prefix sum.** Another fundamental operator is prefix sum, and has numerous applications in databases, including in radix hash joins and parallel filtering [48]. The prefix sum for a series of values is the cumulative sum up to each value inclusive, (i.e.  $out_k = \sum_{i=0}^k in_i$  for  $k \in \{0, 1, \dots, N - 1\}$ ), where  $N$  in the number of inputs. The serial implementation of prefix sum is trivial and easy for compiling efficient code. Each element is read one by one and is added to a counter initialised with 0. On every read, the value of the counter is written back as output for the corresponding position.

Figure 7 presents our custom instruction for the task. A widely-used algorithm for parallelising prefix sum is from Hillis and Steele [13], and is used in recent SIMD-based software [48]. The first  $\log N$  steps contain a pipelined version of this algorithm, plus one additional stage that adds the cumulative sum of the previous batch, that also happens to be the cumulative sum of the entire input up to that batch. In this way, it can calculate the prefix sum of an arbitrarily long input in a pipelined and non-blocking way. For 64 MiB input, vectorising prefix sum yielded a speedup of **4.1x** over the serial version, though it had 0.4x the speed of ARM A53.

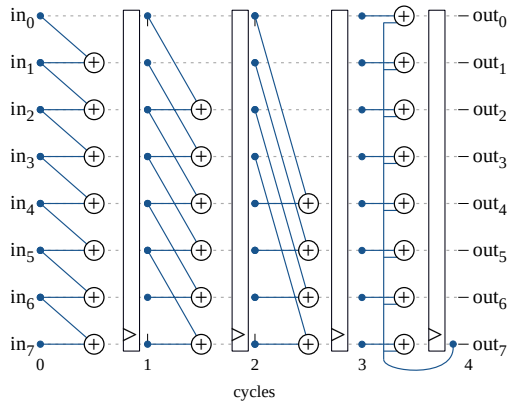


Figure 7: New SIMD instruction for prefix sum

## 5 RELATED WORK

There are various pieces of research trying to introduce FPGAs to vector processors and overlays, but are not suitable for exploring modular SIMD instructions for CPUs. For example, Cheung et al. introduced a simple RISC-V extension for coupling FPGA-based accelerators, but the new instructions are only for passing and retrieving exclusive control to a dedicated accelerator logic [29]. Other works elaborated on a fixed vector instruction set [7, 19, 37, 46, 47] but not for adding custom instructions, and of a custom pipeline complexity. There are also frameworks that eliminate the base instruction set [38] in favour of only executing a specialised task. Additionally, there are works on accelerator generators that are based on SIMD instructions, but there is a complete absence of instructions in the implementation itself [36].

With respect to the softcore implementation, the closest research are from Dao et al. who presented FlexBex [9], an open-source framework for embedding small FPGAs in a modified Ibex RISC-V core [45]. Although it provides a form of SIMD functionality, it is for embedded solutions and without a cache, and the operation is done on multiple 32-bit registers. This makes it limiting for benefiting memory-intensive applications, which is our target. On a similar note, Ordaz et al. [31] developed a closed-source 128-bit wide SIMD engine shared between RISC-V cores [25]. Unlike our solution, the memory interface was much narrower than the SIMD engine, and both of these works mostly focus on the fabrication aspect rather than instruction development and streaming performance.

A narrow datapath and related restrictions are also the case with multiple other FPGA-based softcores [12, 18, 24, 27], where a 32-bit memory interface limitation is often hard-coded. In such cases it could require considerable effort altering the framework to support our methodology for exploring reconfigurable SIMD instructions.

## 6 DISCUSSION

One of the most useful insights from such exploration is about the reduction in the number of instructions and cycles required for a task. For instance, if we look at the `c2_sort` instruction, it is able to sort a list of 8 32-bit elements in 6 cycles. In contrast, a sorting network implementation of only 4 32-bit inputs in older Intel processors required 13 SIMD instructions and 26 cycles [8]. This **13x** and **4.3x** reduction of instructions and cycles respectively,

while solving a bigger problem, is due to the unavailability of such specialised intrinsics. Even with AVX-512 [15], for each layer of compare-and-swap (CAS) units, a pair of separate instructions `min` and `max` are required, as well as a few calls of `shuffle` that permute the inputs for correct alignment [8]. A similar discussion can be made for prefix sum [48], though parallel prefix sum uses more comparisons than the serial case, hence the less notable speedups.

This approach aims mostly at exploration rather than acceleration. Ideally, a higher-end CPU should be hardened, and provide reconfigurability only for instructions. The same is true for the ISA extension, where a 64-bit variant could further reduce the trade-off between the number of registers and the number of operands an instruction can access, if a higher number of registers is preferred (see section 2 for the related design choices). A reduced version of the V extension sharing some of the described features could also be appropriate for reconfigurable use in lighter-weight targets.

Sometimes it cannot compete with dedicated FPGA accelerators, as with a sorter that achieves up to 49x speedup on the same platform [34]. This gap is expected and relates to the presence of instructions in general, such as by trying to avoid internal states. As sorting is not a purely streaming task, many accelerators try to “internalise” processing to reduce the number of times the data is read, but not on CPUs [8].

As a future work it would be appropriate to explore ways to include (or avoid) internal states in such custom instructions, to maximise the benefits of using FPGAs for general purpose computing. In this work, holding a state is allowed in template 1 because this softcore is simple enough to not be of concern, such as with no wrong execution paths. However, in higher-end systems, such as with multi-processors that support context-switching, it would be more challenging to support arbitrary instructions holding states.

CPUs operate at higher frequencies than FPGAs, and this is the reason why NEON-based `mempcpy()` implementations can achieve high bandwidth on ARM [1]. Given that isolated or out-of-context FPGA designs can run much faster than when integrated in bigger systems [18], hardening all of communication could further close the gap between hardened logic and FPGA-based instructions.

## 7 CONCLUSION

The provided instruction types and templates provide the ability to develop advanced SIMD instructions with a few lines of HDL code, and minimise the instruction count for increased performance. This softcore, in combination with the proposed optimisations can be used to explore novel high-performance SIMD instructions. It is demonstrated that custom SIMD instructions can provide an order of magnitude of speedup over serial implementations for memory-intensive applications. The availability of small reconfigurable regions as instructions in future generations of CPUs could be a more efficient use of silicon and processing cycles, and also simplify designs and solve the main memory bottleneck found in today’s FPGA-based datacenter accelerators.

## ACKNOWLEDGMENT

This research was sponsored by dunnhumby. The support of Microsoft and the United Kingdom EPSRC (grant number EP/L016796/1, EP/I012036/1, EP/L00058X/1, EP/N031768/1 and EP/K034448/1), European Union Horizon 2020 Research and Innovation Programme (grant number 671653) is gratefully acknowledged.

## REFERENCES

- [1] 2013. A.3.2. Cortex-A8 Data memory access. In *NEON Programmer's Guide*. arm, 5–7. <https://developer.arm.com/documentation/den0018/a>
- [2] 2020. RISC-V "V" Vector Extension, Version 0.9. (2020).
- [3] Krste Asanovic, David A Patterson, and Christopher Celio. 2015. *The berkeley out-of-order machine (boom): An industry-competitive, synthesizable, parameterized risc-v processor*. Technical Report. University of California at Berkeley Berkeley United States.
- [4] Kenneth E Batchner. 1968. Sorting networks and their applications. In *Proceedings of the April 30–May 2, 1968, spring joint computer conference*. ACM, 307–314.
- [5] Berenger Bramas. 2017. A novel hybrid quicksort algorithm vectorized using AVX-512 on Intel Skylake. *International Journal of Advanced Computer Science and Applications (IJACSA)* 8, 10 (2017), 337–344.
- [6] Jared Casper and Kunle Olukotun. 2014. Hardware acceleration of database operations. In *Proceedings of the 2014 ACM/SIGDA international symposium on Field-programmable gate arrays*. 151–160.
- [7] Chen Chen, Xiaoyan Xiang, Chang Liu, Yunhai Shang, Ren Guo, Dongqi Liu, Yimin Lu, Ziyi Hao, Jiahui Luo, Zhijian Chen, et al. 2020. Xuantie-910: A commercial multi-core 12-stage pipeline out-of-order 64-bit high performance risc-v processor with vector extension: Industrial product. In *2020 ACM/IEEE 47th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 52–64.
- [8] Jatin Chhugani, Anthony D Nguyen, Victor W Lee, William Macy, Mostafa Hagog, Yen-Kuang Chen, Akram Baransi, Sanjeev Kumar, and Pradeep Dubey. 2008. Efficient implementation of sorting on multi-core SIMD CPU architecture. *Proceedings of the VLDB Endowment* 1, 2 (2008), 1313–1324.
- [9] Nguyen Dao, Andrew Attwood, Bea Healy, and Dirk Koch. 2020. FlexBex: A RISC-V with a Reconfigurable Instruction Extension. In *2020 International Conference on Field-Programmable Technology (ICFPT)*. IEEE, 190–195.
- [10] Elsayed A Elsayed and Kenji Kise. 2020. High-Performance and Hardware-Efficient Odd-Even Based Merge Sorter. *IEICE Transactions on Information and Systems* 103, 12 (2020), 2504–2517.
- [11] Robert J Halstead, Ildar Absalyamov, Walid A Najjar, and Vassilis J Tsotras. 2015. FPGA-based Multithreading for In-Memory Hash Joins. In *CIDR*.
- [12] Carsten Heinz, Yannick Lavan, Jaco Hofmann, and Andreas Koch. 2019. A Catalog and In-Hardware Evaluation of Open-Source Drop-In Compatible RISC-V Software Processors. In *2019 International Conference on ReConfigurable Computing and FPGAs (ReConFig)*. IEEE, 1–8.
- [13] W Daniel Hillis and Guy L Steele Jr. 1986. Data parallel algorithms. *Commun. ACM* 29, 12 (1986), 1170–1183.
- [14] A Holdings. [n.d.]. AMBA AXI and ACE protocol specification. *Tech. rep.* 2011 ([n. d.]).
- [15] Intel (R). [n.d.]. Intel Intrinsics Guide. <https://software.intel.com/sites/landingpage/IntrinsicsGuide/>
- [16] Intel (R). [n.d.]. Intel(R) Integrated Performance Primitives: Developer Reference, Volume 1: Signal Processing. (Accessed on 11/01/2021). <https://software.intel.com/content/www/us/en/develop/documentation/ipp-dev-reference/top/volume-1-signal-and-data-processing.html>
- [17] AXI Interconnect. 2017. v2. 1 LogiCORE IP Product Guide. *PG059, Xilinx, December 20* (2017).
- [18] Md Ashrafal Islam, Hiromu Miyazaki, and Kenji Kise. 2020. RVCOREP-32IM: An effective architecture to implement mul/div instructions for five stage RISC-V soft processors. *arXiv preprint arXiv:2010.16171* (2020).
- [19] Matthew Johns and Tom J Kazmierski. 2020. A Minimal RISC-V Vector Processor for Embedded Systems. In *2020 Forum for Specification and Design Languages (FDL)*. IEEE, 1–4.
- [20] Kaan Kara, Zeke Wang, Ce Zhang, and Alonso Gustavo. 2019. doppiODB 2.0: hardware techniques for improved integration of machine learning into databases. *Proceedings of the VLDB Endowment* 12, 12 (2019), 1818–1821.
- [21] Ryohei Kobayashi and Kenji Kise. 2015. Face: Fast and customizable sorting accelerator for heterogeneous many-core systems. In *2015 IEEE 9th International Symposium on Embedded Multicore/Many-core Systems-on-Chip*. IEEE, 49–56.
- [22] Kristiyan Manev, Anuj Vaishnav, and Dirk Koch. 2019. Unexpected Diversity: Quantitative Memory Analysis for Zynq UltraScale+ Systems. In *International Conference on Field-Programmable Technology (ICFPT)*. IEEE, 179–187.
- [23] Susumu Mashimo, Akifumi Fujita, Reoma Matsuo, Seiya Akaki, Akifumi Fukuda, Toru Koizumi, Junichiro Kadamoto, Hidetsugu Irie, Masahiro Goshima, Koji Inoue, et al. 2019. An Open Source FPGA-Optimized Out-of-Order RISC-V Soft Processor. In *2019 International Conference on Field-Programmable Technology (ICFPT)*. IEEE, 63–71.
- [24] Eric Matthews, Alec Lu, Zhenman Fang, and Lesley Shannon. 2019. Rethinking integer divider design for fpga-based soft-processors. In *2019 IEEE 27th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 289–297.
- [25] Eric Matthews and Lesley Shannon. 2017. Taiga: A new risc-v soft-processor framework enabling high performance cpu architectural features. In *2017 27th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 1–4.
- [26] John D McCalpin. 1995. Stream benchmark. *Link: www.cs.virginia.edu/stream/ref.html#what* 22 (1995).
- [27] Junya Miura, Hiromu Miyazaki, and Kenji Kise. 2020. A portable and Linux capable RISC-V computer system in Verilog HDL. *arXiv preprint arXiv:2002.03576* (2020).
- [28] Hiroki Nakahara, Zhiqiang Que, and Wayne Luk. 2020. High-Throughput Convolutional Neural Network on an FPGA by Customized JPEG Compression. In *28th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 1–9.
- [29] Ho-Cheung Ng, Cheng Liu, and Hayden Kwok-Hay So. 2016. A Soft Processor Overlay with Tightly-coupled FPGA Accelerator. *arXiv:1606.06483 [cs.AR]*
- [30] Jose Raul Garcia Ordaz and Dirk Koch. 2016. soft-NEON: A study on replacing the NEON engine of an ARM SoC with a reconfigurable fabric. In *27th International Conference on Application-specific Systems, Architectures and Processors (ASAP)*. IEEE, 229–230.
- [31] Ordaz, Jose Raul Garcia and Koch, Dirk. 2018. A Soft Dual-Processor System with a Partially Run-Time Reconfigurable Shared 128-Bit SIMD Engine. In *29th International Conference on Application-specific Systems, Architectures and Processors (ASAP)*. IEEE, 1–8.
- [32] Muhsen Owaida, David Sidler, Kaan Kara, and Gustavo Alonso. 2017. Centaur: A framework for hybrid CPU-FPGA databases. In *2017 IEEE 25th Annual International Symposium on Field-Programmable Custom Computing Machines (FCCM)*. IEEE, 211–218.
- [33] Philippos Papaphilippou, Chris Brooks, and Wayne Luk. 2018. FLiMS: Fast Lightweight Merge Sorter. In *2018 International Conference on Field-Programmable Technology (FPT)*. IEEE, 78–85.
- [34] Philippos Papaphilippou, Chris Brooks, and Wayne Luk. 2020. An Adaptable High-Throughput FPGA Merge Sorter for Accelerating Database Analytics. In *2020 30th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 65–72.
- [35] Philippos Papaphilippou and Wayne Luk. 2018. Accelerating database systems using FPGAs: A survey. In *2018 28th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 125–130.
- [36] Mehmet Alp Sarkisla and Arda Yurdakul. 2021. SIMDify: Framework for SIMD-Processing with RISC-V Scalar Instruction Set. In *2021 Australasian Computer Science Week Multiconference*. 1–10.
- [37] Aaron Severance and Guy Lemieux. 2012. VENICE: A compact vector processor for FPGA applications. In *2012 International Conference on Field-Programmable Technology*. IEEE, 261–268.
- [38] Aaron Severance and Guy GF Lemieux. 2013. Embedded supercomputing in FPGAs with the VectorBlox MXP matrix processor. In *2013 International Conference on Hardware/Software Codesign and System Synthesis (CODES+ ISSS)*. IEEE, 1–10.
- [39] Wei Song, Dirk Koch, Mikel Lujan, and Jim Garside. 2016. Parallel hardware merge sorter. In *Field-Programmable Custom Computing Machines (FCCM), 24th Annual International Symposium on*. IEEE, 95–102.
- [40] Jeffrey Stuecheli, William J Starke, John D Irish, L Baba Arimilli, D Dreps, Bart Blaner, Curt Wollbrink, and Brian Allison. 2018. IBM POWER9 opens up a new era of acceleration enablement: OpenCAPI. *IBM Journal of Research and Development* 62, 4/5 (2018), 8–1.
- [41] Kosuke Tatsumura, Alexander Dixon, and Hayato Goto. 2019. FPGA-Based Simulated Bifurcation Machine. In *29th International Conference on Field Programmable Logic and Applications (FPL)*. IEEE, 59–66.
- [42] T2 UltraSPARC. 2006. Supplement to the UltraSPARC architecture 2007.
- [43] Andrew Waterman and Krste Asanovic. 2020. The RISC-V instruction set manual, volume I: Unprivileged ISA document, version 20191214-draft. *RISC-V Foundation, Tech. Rep* (2020).
- [44] C. Wolf. 2019. Picorv32-a size-optimized risc-v cpu.
- [45] www.lowrisc.org. [n.d.]. Ibex RISC-V. ([n. d.]). <https://github.com/lowRISC/ibex>
- [46] Peter Yiannacouras, J Gregory Steffan, and Jonathan Rose. 2008. VESPA: portable, scalable, and flexible FPGA-based vector processors. In *Proceedings of the 2008 international conference on Compilers, architectures and synthesis for embedded systems*. 61–70.
- [47] Jason Yu, Guy Lemieux, and Christopher Eagleston. 2008. Vector processing as a soft-core CPU accelerator. In *Proceedings of the 16th international ACM/SIGDA symposium on Field programmable gate arrays*. 222–232.
- [48] Wangda Zhang, Yanbin Wang, and Kenneth A Ross. [n.d.]. Parallel Prefix Sum with SIMD. *Algorithms* 5 ([n. d.]), 31.