

ECE 5745 Complex Digital ASIC Design, Spring 2022

Lab 2: ASIC Sorting Accelerator

School of Electrical and Computer Engineering
Cornell University

revision: 2022-03-05-18-09

In this lab, you will explore a medium-grain hardware accelerator for sorting an array of integer values of unknown length. The baseline design is a pure-software sorting microbenchmark running on a pipelined processor with its own instruction and data cache. The alternative design augments the baseline design with a hardware accelerator and includes the necessary software to configure and potentially assist the accelerator. You will use a standard-cell ASIC toolflow to quantitatively analyze the area, energy, and performance of both the baseline and alternative designs. You are required to implement the alternative design, verify the design using an effective testing strategy, push all designs through the ASIC toolflow, and perform an evaluation comparing the various implementations. This lab is designed to give you experience with: (1) application-specific accelerator design; (2) software/hardware co-design; (3) state-of-the-art standard-cell ASIC toolflow for quantitatively analyzing the area, energy, and timing of a design.

This handout assumes that you have read and understand the course tutorials and the lab assignment logistics document. To get started, login to an `ecelinux` machine, source the setup script, and clone your lab group's remote repository from GitHub:

```
% source setup-ece5745.sh
% mkdir -p ${HOME}/ece5745
% cd ${HOME}/ece5745
% git clone git@github.com:cornell-ece5745/lab2-groupXX.git
```

where `XX` is your group number. **You should never fork your lab group's remote repository! If you need to work in isolation then use a branch within your lab group's remote repository.** If you have already cloned your lab group's remote repository, then use `git pull` to ensure you have any recent updates before running all of the tests. You can run all of the tests in the lab like this:

```
% cd ${HOME}/ece5745/lab2-groupXX
% git pull --rebase
% mkdir -p sim/build
% cd sim/build
% pytest ../lab2_xcel
```

All of the tests should pass except for the tests related to the sorting accelerator you will be implementing in this lab. For this lab, you will be making use of the following subprojects:

- `lab1_imul` – Staff solution to lab 1, used in `proc`
- `lab2_xcel` – Your solution to lab 2
- `sram` – SRAMs used in cache and potentially your accelerator
- `proc` – Pipelined, five-stage, TinyRV2 processor
- `cache` – Blocking, two-way, set-associative cache
- `pmx` – Processor, cache, accelerator composition

You will be mostly be working in the `lab2_xcel` subproject which includes the following files:

- `SortXcelFL.py` – FL sorting accelerator
- `SortXcelPRTL.py` – PyMTL RTL model of sorting accelerator
- `SortXcelVRTL.v` – Verilog RTL model of sorting accelerator
- `SortXcelRTL.py` – Wrapper to choose which RTL language
- `sort-xcel-sim` – Accelerator simulator for isolated eval
- `__init__.py` – Package setup
- `test/SortXcelFL_test.py` – FL sorting accelerator unit tests
- `test/SortXcelRTL_test.py` – RTL sorting accelerator unit tests
- `test/__init__.py` – Test package setup

1. Introduction

In ECE 4750, you gained experience designing, implementing, testing, and evaluating general-purpose processors, memories, and networks. In this lab, you will gain experience with a similar process for an application-specific medium-grain accelerator. Fine-grain accelerators are tightly integrated within the processor pipeline (e.g., a specialized functional unit for bit-reversed addressing useful in implementing an FFT), while coarse-grain accelerators are loosely integrated with a processor through the memory hierarchy (e.g., a graphics rendering accelerator sharing the last-level cache with a general-purpose processor). Medium-grain accelerators are often integrated as co-processors: the processor can directly send/receive messages to/from the accelerator with special instructions, but the co-processor is relatively decoupled from the main processor pipeline and can also independently interact with memory.

We will be accelerating a simple sorting application. The application is given a source array, a destination array, and a size. The application should sort the input array of unsigned integers in increasing numerical order and place the result in the given destination array. The application is allowed to modify the source array. The application must be able to handle both very small arrays (e.g., four elements) and large arrays (e.g., thousands of elements). You can assume the data in the source array will be uniformly distributed 32-bit random integers.

2. Baseline Design

The baseline design for this lab assignment is a pure-software sorting microbenchmark running on a pipelined processor with its own instruction and data cache. Figure 1 illustrates the overall system we will be using in this lab assignment. The processor includes eight latency insensitive `val/rdy` interfaces. The `mnggr2proc/proc2mnggr` interfaces are used for the test harness to send data to the processor and for the processor to send data back to the test harness. The `imemreq/imemresp` interfaces are used for instruction fetch, and the `dmemreq/dmemresp` interfaces are used for implementing load/store instructions. The system includes both instruction and data caches. The `xcelreq/xcelresp` interfaces are used for the processor to send messages to the accelerator. The `mnggr2proc/proc2mnggr` and `memreq/memresp` interfaces were all introduced in ECE 4750. For the baseline design, we provide a simple null accelerator which you can largely ignore.

We provide two implementations of the TinyRV2 processor. The FL model in `sim/proc/ProcFL.py` is essentially an instruction-set-architecture (ISA) simulator; it simulates only the instruction semantics and makes no attempt to model any timing behavior. The RTL model in `sim/proc/ProcPRTL.py` is similar to the alternative design for lab 2 in ECE 4750. It is a five-stage pipelined processor that implements the TinyRV2 instruction set and includes full bypassing/forwarding to resolve data hazards.

There are two important differences from the alternative design for lab 2 of ECE 4750. First, the new processor design uses a single-cycle integer multiplier. We can push the design through the flow and verify that the single-cycle integer multiplier does not adversely impact the overall processor cycle time. Second, the new processor design includes the ability to handle new CSRs for interacting with medium-grain accelerators. The datapath diagram for the processor is shown in Figure 2.

We provide an RTL model in `sim/cache/BlockingCachePRTL.py` which is very similar to the alternative design for lab 3 of ECE 4750. It is a two-way set-associative cache with 16B cache lines and a write-back/write-allocate write policy and LRU replacement policy. There are three important differences from the alternative design for lab 3 of ECE 4750. First, the new cache design is larger with a total capacity of 8 KB. Second, the new cache design carefully merges states to enable a single-cycle hit latency for both reads and writes. Note that writes have a two cycle occupancy (i.e., back-to-back writes will only be able to be serviced at half throughput). Third, the previous cache design used combinational-read SRAMs, while the new cache design uses synchronous-read SRAMs. Combinational-read SRAMs mean the read data is valid on the same cycle we set the read address. Synchronous-read SRAMs mean the read data is valid on the cycle *after* we set the read address. Combinational SRAMs simplify the design, but are not realistic. Almost all real SRAM memory generators used in ASIC toolflows produce synchronous-read SRAMs, and indeed the CACTI memory compiler discussed in the previous tutorial also produces synchronous-read SRAMs. Using synchronous-read SRAMs requires non-trivial changes to both the datapath and the control logic. The cache FSM must make sure all control signals for the SRAM are ready the cycle before we need the read data. The datapath and FSM diagrams for the new cache are shown in Figures 3 and 4. Notice in the FSM how we are able to stay in the `TAG_CHECK_READ_DATA` state if another request is ready; this is what produces the single-cycle hit latency.

Finally, we provide a reasonably optimized pure-software sorting microbenchmark which uses `quicksort` in `app/ubmark/ubmark-sort.c`. Note that we also provide more optimized versions of `quicksort` in (i.e., `-v2`, `-v3`). However, `-v1` is the one currently used in `ubmark-sort.c`, and this is the one you should use as your baseline design.

3. Alternative Design

The alternative design is to implement a medium-grain sorting accelerator and to develop a corresponding software microbenchmark that takes advantage of this new accelerator. This lab is completely open-ended. There are many ways to design such an accelerator. You can design a simple FSM-based accelerator that iteratively loads one or two values from memory does a comparison and writes one or two values back to memory algorithm to implement bubble sort, selection sort, merge sort, quicksort, or some other comparison sort. You might consider a non-comparison sort such as radix sort. You can also implement a more sophisticated accelerator that streams a block of elements into a temporary buffer in the accelerator, sort these values inside the accelerator, and then streams this block of elements back out to the memory system. You would need to merge the sorted blocks in either software or hardware. You can implement multiple “sub-accelerators” which map to different accelerator registers. You are free to use SRAMs.

We provide you an FL model of a sorting accelerator which uses the following accelerator registers:

- `xr0`: go/done
- `xr1`: base address of array
- `xr2`: number of elements in array

Using the accelerator protocol involves the following steps:

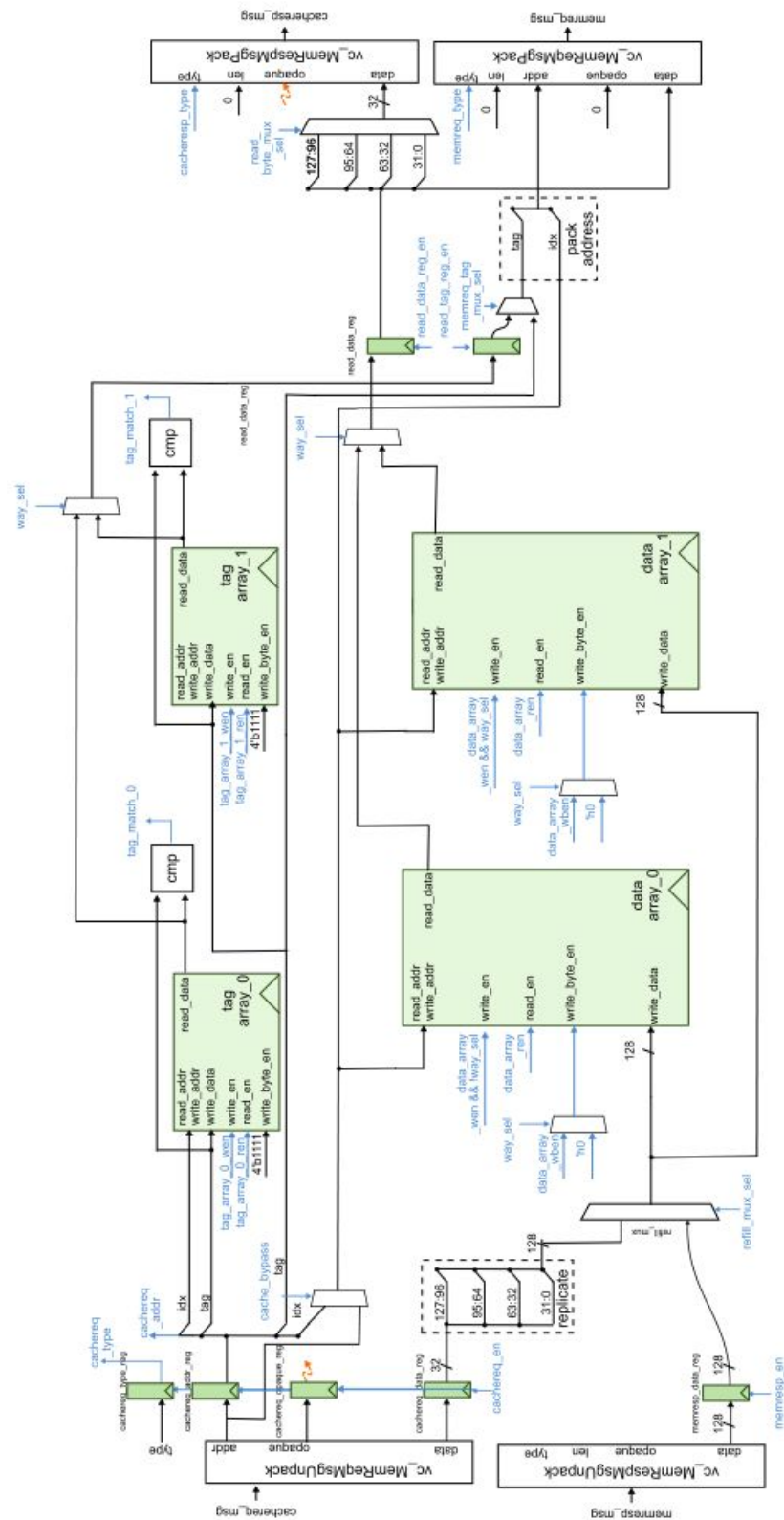


Figure 3: Baseline Cache Datapath

cycle and you must use the provided accelerator message interface; you can change the accelerator protocol but you must use the basic xcelreq/xcelresp interface. We also ask you not to flatten your design, nor make significant modifications to the scripts used in the automated ASIC flow. We ask you to keep the clock constraint at a specific value which will be provided by the instructor. If your accelerator significantly impacts the cycle time compared to the baseline design, then you should optimize your accelerator so that it is no longer on the critical path.

Recall that the sorting application must be able to handle both very small arrays (e.g., four elements) and large arrays (e.g., thousands of elements). It is fine for your accelerator to only handle a fixed array size, but you will need to ensure that your software can use this accelerator to sort an arbitrary array size.

You should not feel limited to implementing only one alternative design. Feel free to implement multiple alternative designs, or to experiment with different parameters in order to lay the foundation for a rich and compelling design-space exploration in your lab report.

4. Testing Strategy

We have provided you with some basic unit tests that test the sorting accelerator in isolation. These tests are defined in `sim/lab2_xcel/test/SortXcelFL_test.py` and they are reused in the test scripts for the RTL models. You will need to modify these tests if you change the sorting accelerator protocol. For example, if your sorting accelerator can only sort a fixed array size, you will need to modify the unit tests appropriately. You will almost certainly want to add more tests for specific corner cases related to your sorting accelerator implementation. The following commands illustrate how to run all of the tests for the entire project, how to run just the tests for this lab, and how to run just the tests for the FL and RTL models.

```
% cd ${HOME}/ece5745/lab2-groupXX/sim/build
% pytest ..
% pytest ../lab2_xcel
% pytest ../lab2_xcel/test/SortXcelFL_test.py --verbose
% pytest ../lab2_xcel/test/SortXcelRTL_test.py --verbose
```

Once your alternative design passes all of the provided tests, then you should verify that your RTL model can be successfully translated using the following command:

```
% cd ${HOME}/ece5745/lab2-groupXX/sim/build
% pytest ../lab2_xcel/test/SortXcelRTL_test.py --test-verilog
```

5. Evaluation

Once you have verified the functionality of your alternative design, you should then use the provided simulator to evaluate the cycle-level performance of both the baseline and alternative designs. You can build the microbenchmarks for the baseline and alternative designs like this:

```
% mkdir -p ${HOME}/ece5745/lab2-groupXX/app/build
% cd ${HOME}/ece5745/lab2-groupXX/app/build
% ../configure --host=riscv32-unknown-elf
% make ubmark-sort
% make ubmark-sort-xcel
% make ubmark-sort-dummy
```

We will use the dummy microbenchmark to estimate the energy of the initialization and verification overhead. You can run the simulator for the FL and RTL models of the baseline design like this:

```
% cd ${HOME}/ece5745/lab2-groupXX/sim/build
% ../pmx/pmx-sim ../../app/build/ubmark-sort
% ../pmx/pmx-sim --proc-impl rtl --cache-impl rtl \
  --stats ../../app/build/ubmark-sort
```

The next step is to test your sorting microbenchmark on an FL model of the processor and accelerator. This step verifies that your microbenchmark correctly uses the sorting accelerator protocol.

```
% cd ${HOME}/ece5745/lab2-groupXX/sim/build
% ../pmx/pmx-sim --xcel-impl sort-fl ../../app/build/ubmark-sort-xcel
```

Only after this works, should you try running your sorting microbenchmark on the RTL models. You might want to try running your sorting microbenchmark with and without the cache.

```
% cd ${HOME}/ece5745/lab2-groupXX/sim/build
% ../pmx/pmx-sim --proc-impl rtl --xcel-impl sort-rtl \
  ../../app/build/ubmark-sort-xcel
% ../pmx/pmx-sim --proc-impl rtl --cache-impl rtl --xcel-impl sort-rtl \
  ../../app/build/ubmark-sort-xcel
```

Since we do not have any “real” integration tests of the processor, cache, and accelerator composition, you might want to consider running your sorting microbenchmark on the FL and RTL models as part of your testing strategy. Once the microbenchmark passes, we can use the `--stats` option to determine the total number of cycles to execute the specified input dataset.

```
% cd ${HOME}/ece5745/lab2-groupXX/sim/build
% ../pmx/pmx-sim --proc-impl rtl --cache-impl rtl --xcel-impl sort-rtl \
  --stats ../../app/build/ubmark-sort-xcel
```

You should study the line traces (with the `--trace` option) and possibly the waveforms (with the `--dump-vcd` option) to understand the reason why each design performs as it does on the various patterns.

Once you have explored the cycle-level performance of the baseline and alternative designs, you should push the RTL models through the standard-cell ASIC toolflow to quantify the area, energy, and timing of each design. You should start by pushing just the accelerator through the flow. You can generate the Verilog RTL and test benches for the sorting accelerator in isolation like this:

```
% mkdir -p $TOPDIR/sim/build
% cd $TOPDIR/sim/build
% pytest ../lab2_xcel --test-verilog --dump-vtb
% ../lab2_xcel/sort-xcel-sim --impl rtl --input random --stats --translate --dump-vtb
```

You can then push the sorting accelerator through the flow using `mflowgen` like this:

```
% mkdir -p $TOPDIR/asic/build-lab2-sort-xcel
% cd $TOPDIR/asic/build-lab2-sort-xcel
% mflowgen run --design ../lab2-sort-xcel
% make brg-rtl-4-state-vcssim
```



```
% make brg-synopsys-dc-synthesis
% make post-synth-gate-level-simulation
% make post-synth-power-analysis
% make brg-cadence-innovus-signoff
% make brg-flow-summary
```

Make sure every step is successful before moving on to the next step. Once you know your accelerator can be successfully pushed through the flow and also meets timing, you can try pushing the baseline and alternative design through the flow. Note our flow scripts are currently only setup to push the processor and accelerator through the flow without the cache. Start by generating the Verilog RTL and test benches for both the baseline and alternative design.

```
% cd $TOPDIR/sim/build
% ../pmx/pmx-sim --proc-impl rtl --cache-impl null \
  --stats --translate --dump-vtb ../../app/build/ubmark-sort
% ../pmx/pmx-sim --proc-impl rtl --cache-impl null \
  --stats --translate --dump-vtb ../../app/build/ubmark-sort-dummy
% ../pmx/pmx-sim --proc-impl rtl --cache-impl null --xcel-impl sort-rtl \
  --stats --translate --dump-vtb ../../app/build/ubmark-sort-xcel
```

Make sure to use the execution time from the `pmx-sim` (which is only when stats are enabled) not the execution time that you will see in the summary reports from the ASIC flow (which is for the entire execution of the application)! Now push the baseline design through the flow using `mflowgen` like this:

```
% mkdir -p $TOPDIR/asic/build-lab2-pmx-null
% cd $TOPDIR/asic/build-lab2-pmx-null
% mflowgen run --design ../lab2-pmx-null
% make brg-rtl-4-state-vcssim
% make brg-synopsys-dc-synthesis
% make post-synth-gate-level-simulation
% make post-synth-power-analysis
% make brg-cadence-innovus-signoff
% make brg-flow-summary
```

Make sure every step is successful before moving on to the next step. Then you can push the alternative design through the flow using `mflowgen` like this:

```
% mkdir -p $TOPDIR/asic/build-lab2-pmx-sort
% cd $TOPDIR/asic/build-lab2-pmx-sort
% mflowgen run --design ../lab2-pmx-sort
% make brg-rtl-4-state-vcssim
% make brg-synopsys-dc-synthesis
% make post-synth-gate-level-simulation
% make post-synth-power-analysis
% make brg-cadence-innovus-signoff
% make brg-flow-summary
```

Make sure every step is successful before moving on to the next step. Given all of these results, here is how you should analyze the execution time, cycle time, area, and energy.

- **Execution Time:** Use the number of cycles reported from `pmx-sim`, this only counts cycles when stats are enabled. Use the line traces to understand performance overheads.
- **Cycle Time:** Ensure your accelerator in isolation along with the processor/accelerator composition all meet timing at 3.0ns (333MHz). Use the timing reports from the accelerator in isolation to determine the critical path of your accelerator.
- **Area:** Use the area of your accelerator in isolation and compare it to the area of the baseline processor design. Use the hierarchical area reports from the accelerator in isolation to determine how much area each module in your design requires.
- **Energy:** Subtract the energy of the dummy microbenchmark from the energy of the baseline design to estimate the energy of the software-only baseline. Subtract the energy of the dummy microbenchmark from the energy of the alternative design to estimate the energy of the alternative design. Use the hierarchical energy reports from the accelerator in isolation to determine how much energy each module in your design requires.

6. Pareto-Optimal Frontier Competition

We hope students will continue to modify their software and/or hardware to improve the performance of their sorting accelerator. To encourage such optimization, a small bonus will be given to those designs which lay on the pareto-optimal frontier in the area vs. performance space. We will run your sorting microbenchmark on your processor, cache, and accelerator composition using one or more of our own private datasets. These datasets will not be the same size as the dataset we give you. They may be smaller or they may be larger. Note that students should only attempt improving the software and/or hardware once everything is completely working and they have finished a draft of the lab report. We allow students to submit up to two different accelerator designs for the purposes of evaluating the pareto-optimal frontier.

Here are the steps we will use to determine the pareto-optimal frontier. We first clone your lab repo. Note that all of these steps must work from a clean clone, otherwise your design will not be considered for the pareto-optimal frontier.

```
% source setup-ece5745.sh
% mkdir -p $HOME/ece5745
% cd $HOME/ece5745
% git clone git@github.com:cornell-ece5745/lab2-groupXX
% cd lab2-groupXX
% TOPDIR=$PWD
```

We will then evaluate the area and cycle time of the accelerator in isolation.

```
% mkdir -p $TOPDIR/sim/build
% cd $TOPDIR/sim/build
% ../lab2_xcel/sort-xcel-sim --impl rtl --input random --translate --dump-vtb

% mkdir -p $TOPDIR/asic/build-lab2-sort-xcel
% cd $TOPDIR/asic/build-lab2-sort-xcel
% mkdir -p $TOPDIR/asic/build-lab2-sort-xcel
% cd $TOPDIR/asic/build-lab2-sort-xcel
% mflowgen run --design ../lab2-sort-xcel
% make
```

We will use the area of your accelerator in isolation as one metric for determining the pareto-optimal frontier. The reason we use the area of your accelerator in isolation as opposed to the area of the processor and accelerator composition, is because we want to avoid any variation in the way the tools synthesize and place-and-route the processor from impacting our evaluation of the area of the accelerator. The cycle time of your accelerator in isolation must be less than the cycle time of the baseline design. Please consult the instructors if the cycle time of your accelerator is a little greater than the cycle time of the baseline design; this may still be acceptable.

We will then evaluate the cycle time and the execution time of the processor and accelerator composition using your sorting microbenchmark.

```
% mkdir -p $TOPDIR/app/build
% cd $TOPDIR/app/build
% ../configure --host=riscv32-unknown-elf
% make ubmark-sort-xcel

% cd $TOPDIR/sim/build
% ../pmx/pmx-sim --proc-impl rtl --cache-impl null --xcel-impl sort-rtl \
  --stats --translate --dump-vtb ../../app/build/ubmark-sort-xcel

% mkdir -p $TOPDIR/asic/build-lab2-pmx-sort
% cd $TOPDIR/asic/build-lab2-pmx-sort
% mflowgen run --design ../lab2-pmx-sort
% make
```

If the cycle time is significantly larger than the baseline design, you probably need to revise your accelerator, but you should also discuss these results with the instructors.

For those groups that want to submit two different accelerator designs, we will also do the following steps. Note that you are responsible for: modifying `sort-xcel-sim` and `pmx-sim` so they can correctly instantiate both sorting accelerators; adding new 'flow.py' scripts to the `asic` directory; and adding a new microbenchmark which uses your second sorting accelerator.

```
% mkdir -p $TOPDIR/sim/build
% cd $TOPDIR/sim/build
% ../lab2_xcel/sort-xcel-sim --impl rtl2 --input random --translate --dump-vtb

% mkdir -p $TOPDIR/asic/build-lab2-sort2-xcel
% cd $TOPDIR/asic/build-lab2-sort2-xcel
% mflowgen run --design ../lab2-sort2-xcel
% make

% mkdir -p $TOPDIR/app/build
% cd $TOPDIR/app/build
% ../configure --host=riscv32-unknown-elf
% make ubmark-sort2-xcel

% cd $TOPDIR/sim/build
% ../pmx/pmx-sim --proc-impl rtl --cache-impl rtl --xcel-impl sort2-rtl \
  --stats --translate --dump-vtb ../../app/build/ubmark-sort2-xcel
```

```
% mkdir -p $TOPDIR/asic/build-lab2-pmx-sort2
% cd $TOPDIR/asic/build-lab2-pmx-sort2
% mflowgen run --design ../lab2-pmx-sort2
% make
```

7. Report

The *Lab Assignment Logistics* document provides general details about the requirements for the report. You must actually read this document to ensure you know what goes in your report and how to format it.

Alternative Design Section – In the alternative design section of your report you should discuss your sorting accelerator design in detail. Be sure to include one or more diagrams illustrating your design. If you also implemented more than one sorting accelerator, then you must describe all of these accelerators in the alternative design section.

Evaluation Section – In addition to tables and bar plots, you should include several energy vs. performance plots similar to what we use in lecture to compare the various designs. We recommend having two energy vs. performance plots; one plot for sorting a smaller array and one plot for sorting a larger array. Each plot should have two points: one for the baseline design and one for the alternative design. You could also have additional energy vs. performance plots for different datasets with the same input array size. These energy vs. performance plots should enable you to provide deep insight into the various trade-offs in the evaluation section of your report. If you also implemented more than one accelerator, then you must evaluate it in the alternative design section. **You must include at least two amoeba plots in your report: one of your accelerator in isolation and one of the processor and accelerator composition!**

Acknowledgments

This lab was created by Christopher Batten, Khalid Al-Hawaj, Jason Setter, Shunning Jiang, Moyang Wang, and Jack Brzozowski as part of the course ECE 5745 Complex Digital ASIC Design at Cornell University.