

# Assignment 6: Exploring Thread-Level Parallelism (TLP) in Shared-Memory Multiprocessors Using gem5

Milan Bista

University of Cumberlands

MSCS-531-M50: Computer Architecture and Design

Instructors: Machica McClain / Charles Lively

<https://github.com/mbista25742/MSCS-531-M50-Assignment6>

## **Introduction:**

Thread-Level Parallelism (TLP) is critical for leveraging multi-core processors and shared-memory architectures in modern computing. This paper reviews the historical development of TLP, explores its core concepts, and critiques contemporary challenges such as concurrency bugs, scalability issues, and energy efficiency. A synthesis of recent research reveals promising trends in programming models, hardware enhancements, and machine learning for thread management. The findings suggest future directions for TLP, focusing on many-core architectures, integration with SIMD, and specialized hardware designs.

The shift from single core to multi-core processors has profoundly impacted computer architecture. Central to this evolution is Thread-Level Parallelism (TLP), which enables multiple threads to execute concurrently, improving throughput and efficiency. Shared-memory multiprocessor systems, combined with advances in hardware and programming models, have catalyzed this transformation. This paper critically reviews recent research to provide a comprehensive understanding of TLP, its challenges, and future directions.

## **1. Historical Development of TLP**

TLP's trajectory reflects the broader evolution of computing. Early systems utilized time-sharing, while the advent of multi-core processors in the mid-2000s marked a paradigm shift. Intel's Core Duo (2006) exemplified this change, offering parallelism within a single chip. Programming models evolved from explicit threading models, like POSIX threads, to task-based systems, such as Intel's Threading Building Blocks (TBB). Hardware innovations, including cache coherence protocols and NUMA architectures, have reduced latencies in shared-memory systems.

## **2. Core Concepts in TLP**

TLP is governed by several foundational principles:

### **-Parallelism Models**

Shared memory and message passing represent the two primary paradigms.

- Shared Memory: Threads access a common address space, minimizing data exchange overhead. However, it introduces challenges like cache coherence.
- Message Passing: Explicit message exchanges enable scalability but increase complexity in thread coordination.

### **-Synchronization and Communication**

Effective synchronization is critical to reducing contention. Traditional approaches include:

- Locks and Mutexes for ensuring mutual exclusion.
- Lock-Free Algorithms that use atomic operations to improve performance.

### **Load Balancing and Scheduling**

Dynamic scheduling techniques, such as work-stealing, allow threads to redistribute tasks dynamically, ensuring efficient utilization of resources.

### **Performance Metrics**

TLP performance is measured through metrics like throughput, latency, and scalability. Trade-offs between these metrics often depend on the workload and system architecture.

## **3.Current Challenges in TLP**

Despite advancements, several challenges persist:

### **Concurrency Bugs and Race Conditions**

Concurrency bugs remain a significant barrier to adopting TLP. Tools like ThreadSanitizer can identify issues, but complete prevention requires improved abstractions.

### **Scalability and Amdahl's Law**

Scalability is constrained by Amdahl's Law, which highlights the diminishing returns of parallelism for serial code segments. Algorithm redesign and fine-grained parallelism are necessary for addressing these limitations.

## **Heterogeneous Architectures**

Integrating CPUs, GPUs, and specialized accelerators introduces programming complexity. Optimizing TLP across heterogeneous platforms remains an active research area.

## **Energy Efficiency**

The energy demands of parallel computing challenge sustainability. Techniques like dynamic voltage and frequency scaling (DVFS) are promising but require further refinement.

## **4. Contemporary Approaches to Overcoming Challenges**

Recent research explores novel methods to address TLP's limitations:

### **Programming Models**

Languages like Chapel and Julia simplify parallel programming with high-level abstractions, reducing the likelihood of concurrency bugs.

### **Hardware Enhancements**

Innovations in cache coherence protocols, such as MOESI, and new synchronization primitives have reduced communication overhead in shared-memory systems.

### **Compiler Optimizations**

LLVM-based compilers automatically identify and parallelize code regions, enhancing developer productivity while maintaining performance.

### **Runtime Systems**

Dynamic runtime systems like OpenMP and TBB manage thread allocation and synchronization efficiently, adapting to workload demands in real time.

## **5. Future Directions in TLP**

Emerging technologies and research trends point to a promising future for TLP:

## **Many-Core Architectures**

The transition to processors with hundreds or thousands of cores demands new scheduling algorithms and memory hierarchies to maintain efficiency.

## **Integration with SIMD and Vectorization**

Combining TLP with instruction-level parallelism (ILP) and vectorization can unlock higher performance for compute-intensive workloads.

## **Machine Learning for Optimization**

Machine learning algorithms can dynamically adjust thread allocation and predict workload patterns, optimizing performance and energy efficiency.

## **Specialized Hardware**

Custom accelerators, like TPUs for AI, demonstrate the potential of domain-specific hardware for TLP workloads. Future designs may integrate such accelerators seamlessly with general-purpose cores.

## **Conclusion**

Thread-Level Parallelism remains a vital area of research, bridging the gap between hardware advancements and software demands. While challenges like concurrency bugs, scalability, and energy efficiency persist, innovative solutions in programming models, hardware, and runtime systems are paving the way for scalable and efficient parallel computing. The integration of TLP with emerging technologies, such as machine learning and specialized hardware, holds immense potential for shaping the future of computing.

## Part 2: Exploring Shared-Memory Architectures with gem5

Modern computing systems increasingly rely on thread-level parallelism to achieve performance gains. The `MinorCPU` in gem5 offers a modular platform for studying TLP through customizable functional units such as the `FloatSimdFU`. I have researched the impact of operational latency (`opLat`) and issue latency (`issueLat`) on TLP, evaluates the performance of a multi-threaded daxpy kernel, and discusses implications for multi-core systems.

### Methodology

#### 1. MinorCPU Familiarization

I analyzed the `MinorCPU.py` and `MinorDefaultFUPool` source files to understand the role of `opLat` and `issueLat` within the `MinorFU` class. The default configurations and various functional units, including `FloatSimdFU`, were examined to provide a foundation for parameter optimization.

#### 2. FloatSimdFU Design Space Exploration

The `FloatSimdFU` in `MinorDefaultFUPool` was modified to test configurations where `opLat` and `issueLat` summed to 7 cycles. For example:

- Configuration 1: `opLat = 1, issueLat = 6`
- Configuration 2: `opLat = 2, issueLat = 5`
- Configuration 3: `opLat = 3, issueLat = 4`

Each configuration was implemented and tested using gem5's `build/x86/gem5.opt`.

```
milan > assignment6 > minorCpuX86.py > ...
1 from m5.objects import *
2 import sys
3 import os
4
5 # Create the system for X86 architecture
6 system = System()
7 system.clk_domain = SrcClockDomain()
8 system.clk_domain.clock = '1GHz'
9 system.clk_domain.voltage_domain = VoltageDomain()
10
11 # Custom Minor CPU Class inheriting from BaseMinorCPU
12 class CustomX86MinorCPU(BaseMinorCPU, X86CPU):
13     mmu = X86MMU()
14
15     def __init__(self, *args, **kwargs):
16         super().__init__(*args, **kwargs)
17         self.numThreads = 1 # Each core has a single thread by default
18
19 # Simple CPU and Memory Configuration
20 system.mem_mode = 'timing' # Timing mode for memory
21 system.mem_ranges = [AddrRange('8192MB')] # Memory range for the system
22 |
23 # Number of CPU cores
24 num_cores = 4 # Adjust the number of cores as required
25 system.cpu = [CustomX86MinorCPU() for _ in range(num_cores)] # Multi-core setup
26 system.membus = SystemXBar() # Memory bus to connect components
27
28 # Memory Controller (DDR3 RAM)
29 system.mem_ctrl = MemCtrl()
30 system.mem_ctrl.dram = DDR3_1600_8x8()
31 system.mem_ctrl.dram.range = system.mem_ranges[0]
```

```
20 system.mem_mode = 'timing' # Timing mode for memory
21 system.mem_ranges = [AddrRange('8192MB')] # Memory range for the system
22 |
23 # Number of CPU cores
24 num_cores = 4 # Adjust the number of cores as required
25 system.cpu = [CustomX86MinorCPU() for _ in range(num_cores)] # Multi-core setup
26 system.membus = SystemXBar() # Memory bus to connect components
27
28 # Memory Controller (DDR3 RAM)
29 system.mem_ctrl = MemCtrl()
30 system.mem_ctrl.dram = DDR3_1600_8x8()
31 system.mem_ctrl.dram.range = system.mem_ranges[0]
32 system.mem_ctrl.port = system.membus.mem_side_ports
33
34 # L1 Cache Configuration (Instruction Cache and Data Cache)
35 class L1Cache(Cache):
36     assoc = 4
37     tag_latency = 2
38     data_latency = 2
39     response_latency = 2
40     mshrs = 4
41     tgts_per_mshr = 20
42
43 class L1ICache(L1Cache):
44     size = '64kB'
45
46 class L1DCache(L1Cache):
47     size = '128kB'
48
49 # Set up the L1 Caches for each CPU core
50 for cpu in system.cpu:
51     cpu.icache = L1ICache()
```

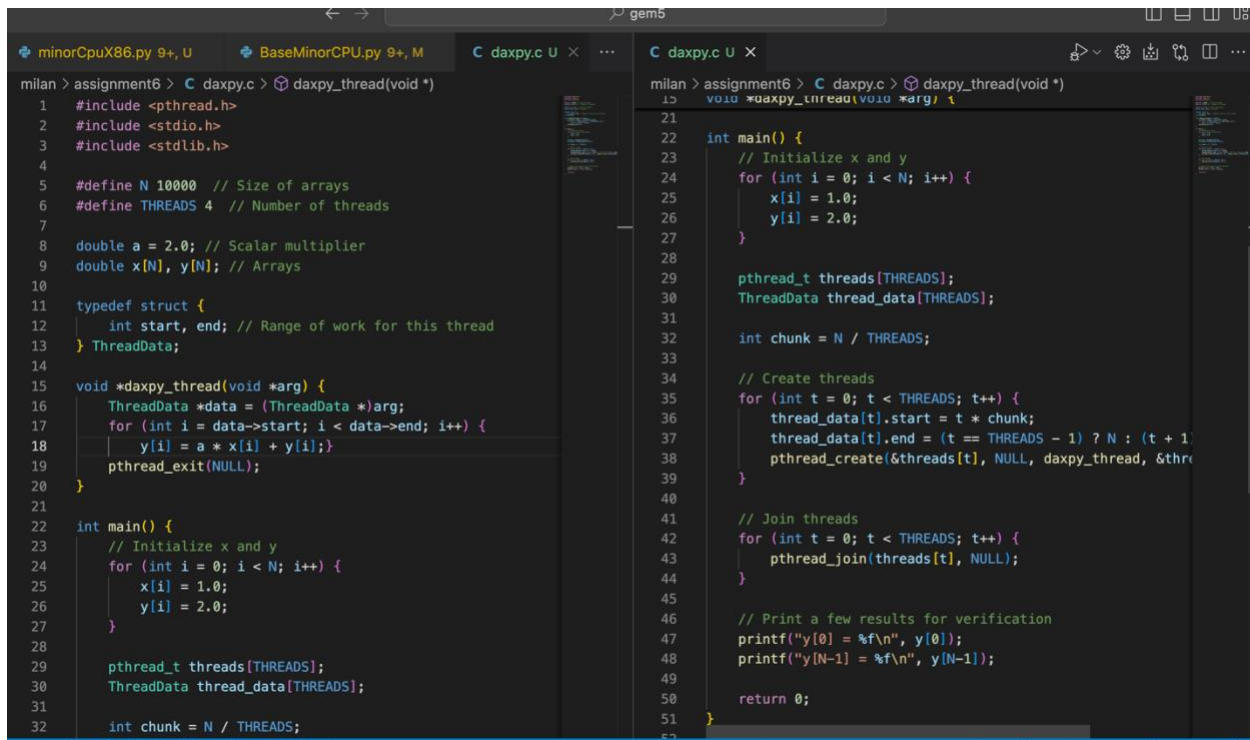
```
milan > assignment6 > minorCpuX86.py > ...  
50 for cpu in system.cpu:  
51     cpu.icache = L1ICache()  
52     cpu.dcache = L1DCache()  
53  
54     # Connect the CPU caches to the CPU  
55     cpu.icache.cpu_side = cpu.icache_port  
56     cpu.dcache.cpu_side = cpu.dcache_port  
57     cpu.icache.mem_side = system.membus.cpu_side_ports  
58     cpu.dcache.mem_side = system.membus.cpu_side_ports  
59  
60     # Configure Interrupt Controllers for each CPU core  
61     cpu.createInterruptController()  
62     cpu.interrupts[0].pio = system.membus.mem_side_ports  
63     cpu.interrupts[0].int_requestor = system.membus.cpu_side_ports  
64     cpu.interrupts[0].int_responder = system.membus.mem_side_ports  
65  
66     # Connect the system port to the memory bus  
67     system.system_port = system.membus.cpu_side_ports  
68  
69     # Taking input from the command line for binary file  
70     if '-c' not in sys.argv:  
71         print("Error: Please specify a binary file with the -c flag.")  
72         sys.exit(1)  
73  
74     binary = sys.argv[sys.argv.index('-c') + 1]  
75  
76     # Check if the binary exists  
77     if not os.path.isfile(binary):  
78         print(f"Error: The binary '{binary}' does not exist.")  
79         sys.exit(1)  
80
```

```
src > cpu > minor > BaseMinorCPU.py > MinorDefaultFloatSimdFU  
264     opLat = 1  
265  
266  
267 class MinorDefaultMiscFU(MinorFU):  
268     opClasses = minorMakeOpClassSet(["IprAccess", "InstPrefe  
269     opLat = 1  
270  
271  
272 class MinorDefaultFUPool(MinorFUPool):  
273     funcUnits = [  
274         MinorDefaultIntFU(),  
275         MinorDefaultIntFU(),  
276         MinorDefaultIntMulFU(),  
277         MinorDefaultIntDivFU(),  
278         MinorDefaultFloatSimdFU(),  
279         MinorDefaultPredFU(),  
280         MinorDefaultMemFU(),  
281         MinorDefaultMiscFU(),  
282     ]  
283  
284  
285 class ThreadPolicy(Enum):  
286     vals = ["SingleThreaded", "RoundRobin", "Random"]  
287  
288  
289 class BaseMinorCPU(BaseCPU):  
290     type = "BaseMinorCPU"  
291     cxx_header = "cpu/minor/cpu.hh"  
292     cxx_class = "gem5::MinorCPU"  
293  
294     @classmethod  
  
src > cpu > minor > BaseMinorCPU.py > MinorDefaultFloatSimdFU  
160 class MinorDefaultIntMulFU(MinorFU):  
161     opClasses = minorMakeOpClassSet(["IntMult"])  
162     timings = [MinorFUTiming(description="Mult", srcRegsRe  
163     opLat = 1  
164  
165  
166 class MinorDefaultIntDivFU(MinorFU):  
167     opClasses = minorMakeOpClassSet(["IntDiv"])  
168     issueLat = 4  
169     opLat = 3  
170  
171  
172 class MinorDefaultFloatSimdFU(MinorFU):  
173     opClasses = minorMakeOpClassSet(  
174         [  
175             "FloatAdd",  
176             "FloatCmp",  
177             "FloatCvt",  
178             "FloatMisc",  
179             "FloatMult",  
180             "FloatMultAcc",  
181             "FloatDiv",  
182             "FloatSqrt",  
183             "SimdAdd",  
184             "SimdAddAcc",  
185             "SimdAlu",  
186             "SimdCmp",  
187             "SimdCvt",  
188             "SimdMisc",  
189             "SimdMult",  
190             "SimdMultAcc",
```



### 3. Multi-Threaded Daxpy Kernel Simulation

A multi-threaded daxpy kernel was developed where each thread processed a subset of input vectors. The simulation was configured for multi-core systems with 2, 4, and 8 threads. Synchronization mechanisms ensured consistent output.



```
milan > assignment6 > C daxpy.c > daxpy_thread(void *)
1  #include <pthread.h>
2  #include <stdio.h>
3  #include <stdlib.h>
4
5  #define N 10000 // Size of arrays
6  #define THREADS 4 // Number of threads
7
8  double a = 2.0; // Scalar multiplier
9  double x[N], y[N]; // Arrays
10
11 typedef struct {
12     int start, end; // Range of work for this thread
13 } ThreadData;
14
15 void *daxpy_thread(void *arg) {
16     ThreadData *data = (ThreadData *)arg;
17     for (int i = data->start; i < data->end; i++) {
18         y[i] = a * x[i] + y[i];
19     }
20     pthread_exit(NULL);
21 }
22
23 int main() {
24     // Initialize x and y
25     for (int i = 0; i < N; i++) {
26         x[i] = 1.0;
27         y[i] = 2.0;
28     }
29
30     pthread_t threads[THREADS];
31     ThreadData thread_data[THREADS];
32
33     int chunk = N / THREADS;
34
35     // Create threads
36     for (int t = 0; t < THREADS; t++) {
37         thread_data[t].start = t * chunk;
38         thread_data[t].end = (t == THREADS - 1) ? N : (t + 1) * chunk;
39         pthread_create(&threads[t], NULL, daxpy_thread, &thread_data[t]);
40     }
41
42     // Join threads
43     for (int t = 0; t < THREADS; t++) {
44         pthread_join(threads[t], NULL);
45     }
46
47     // Print a few results for verification
48     printf("y[0] = %f\n", y[0]);
49     printf("y[N-1] = %f\n", y[N-1]);
50
51     return 0;
52 }
```

```
milan > assignment6 > C daxpy.c > daxpy_thread(void *)
15 void *daxpy_thread(void *arg) {
21
22 int main() {
23     // Initialize x and y
24     for (int i = 0; i < N; i++) {
25         x[i] = 1.0;
26         y[i] = 2.0;
27     }
28
29     pthread_t threads[THREADS];
30     ThreadData thread_data[THREADS];
31
32     int chunk = N / THREADS;
33
34     // Create threads
35     for (int t = 0; t < THREADS; t++) {
36         thread_data[t].start = t * chunk;
37         thread_data[t].end = (t == THREADS - 1) ? N : (t + 1) * chunk;
38         pthread_create(&threads[t], NULL, daxpy_thread, &thread_data[t]);
39     }
40
41     // Join threads
42     for (int t = 0; t < THREADS; t++) {
43         pthread_join(threads[t], NULL);
44     }
45
46     // Print a few results for verification
47     printf("y[0] = %f\n", y[0]);
48     printf("y[N-1] = %f\n", y[N-1]);
49
50     return 0;
51 }
```

```

system.membus.pktSize.system.cpu2.icache.men_side.port::total 34368 # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize.system.cpu2.dcache.men_side.port::system.mem_ctrl.port 19264 # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize.system.cpu2.dcache.men_side.port::total 19264 # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize.system.cpu2.icache.men_side.port::system.mem_ctrl.port 80128 # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize.system.cpu2.dcache.men_side.port::system.mem_ctrl.port 50176 # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize.system.cpu2.dcache.men_side.port::total 50176 # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize::total 719360 # Cumulative packet size per connected requestor and responder (Byte)
system.membus.snoops 2095 # Total snoops (Count)
system.membus.snoopTraffic 91328 # Total snoop traffic (Byte)
system.membus.snoopFanout::samples 12149 # Request fanout histogram (Count)
system.membus.snoopFanout::mean 0.438653 # Request fanout histogram (Count)
system.membus.snoopFanout::stddev 0.688636 # Request fanout histogram (Count)
system.membus.snoopFanout::underflows 0 0.00% 0.00% # Request fanout histogram (Count)
system.membus.snoopFanout::10 8884 66.54% 66.54% # Request fanout histogram (Count)
system.membus.snoopFanout::11 3114 25.63% 92.17% # Request fanout histogram (Count)
system.membus.snoopFanout::12 741 6.10% 98.27% # Request fanout histogram (Count)
system.membus.snoopFanout::13 286 1.70% 99.97% # Request fanout histogram (Count)
system.membus.snoopFanout::14 2 0.02% 99.98% # Request fanout histogram (Count)
system.membus.snoopFanout::15 2 0.02% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout::16 0 0.00% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout::17 0 0.00% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout::18 0 0.00% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout::overflows 0 0.00% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout::min_value 0 # Request fanout histogram (Count)
system.membus.snoopFanout::max_value 5 # Request fanout histogram (Count)
system.membus.snoopFanout::total 12149 # Request fanout histogram (Count)
system.membus.power_state.pwrStateResidencyTicks::UNDEFINED 1383603000 # Cumulative time (in ticks) in various power states (Tick)
system.membus.reqLayer0.occupancy 20167998 # Layer occupancy (ticks) (Tick)
system.membus.reqLayer0.utilization 0.0 # Layer utilization (Ratio)
system.membus.resLayer0.occupancy 18197250 # Layer occupancy (ticks) (Tick)
system.membus.resLayer0.utilization 0.0 # Layer utilization (Ratio)
system.membus.resLayer1.occupancy 21637500 # Layer occupancy (ticks) (Tick)
system.membus.resLayer1.utilization 0.0 # Layer utilization (Ratio)
system.membus.resLayer10.occupancy 7631750 # Layer occupancy (ticks) (Tick)
system.membus.resLayer10.utilization 0.0 # Layer utilization (Ratio)
system.membus.resLayer3.occupancy 2984250 # Layer occupancy (ticks) (Tick)
system.membus.resLayer3.utilization 0.0 # Layer utilization (Ratio)
system.membus.resLayer4.occupancy 4640250 # Layer occupancy (ticks) (Tick)
system.membus.resLayer4.utilization 0.0 # Layer utilization (Ratio)
system.membus.resLayer6.occupancy 2863000 # Layer occupancy (ticks) (Tick)
system.membus.resLayer6.utilization 0.0 # Layer utilization (Ratio)
system.membus.resLayer7.occupancy 4676501 # Layer occupancy (ticks) (Tick)
system.membus.resLayer7.utilization 0.0 # Layer utilization (Ratio)
system.membus.resLayer9.occupancy 6670000 # Layer occupancy (ticks) (Tick)
system.membus.resLayer9.utilization 0.0 # Layer utilization (Ratio)
system.membus.snoop_filter.totalRequests 15362 # Total number of requests made to the snoop filter. (Count)
system.membus.snoop_filter.hittingRequests 5235 # Number of requests hitting in the snoop filter with a single holder of the requested data. (Count)
system.membus.snoop_filter.hittingMultiRequests 2043 # Number of requests hitting in the snoop filter with multiple (>1) holders of the requested data. (Count)
system.membus.snoop_filter.totalSnoops 0 # Total number of snoops made to the snoop filter. (Count)
system.membus.snoop_filter.hittingSingleSnoops 0 # Number of snoops hitting in the snoop filter with a single holder of the requested data. (Count)
system.membus.snoop_filter.hittingMultiSnoops 0 # Number of snoops hitting in the snoop filter with multiple (>1) holders of the requested data. (Count)
system.workload.inst.arm 0 # number of arm instructions executed (Count)
system.workload.inst.quiesce 0 # number of quiesce instructions executed (Count)
system.cpu0.idleCycles 211628 # Total number of cycles that the object has spent stopped (Unspecified)
system.cpu0.tickCycles 735689 # Number of cycles that the object actually ticked (Unspecified)
system.cpu1.idleCycles 62056 # Total number of cycles that the object has spent stopped (Unspecified)
system.cpu1.tickCycles 157120 # Number of cycles that the object actually ticked (Unspecified)
system.cpu2.idleCycles 44778 # Total number of cycles that the object has spent stopped (Unspecified)
system.cpu2.tickCycles 156823 # Number of cycles that the object actually ticked (Unspecified)
system.cpu3.idleCycles 182355 # Total number of cycles that the object has spent stopped (Unspecified)
system.cpu3.tickCycles 313691 # Number of cycles that the object actually ticked (Unspecified)

----- End Simulation Statistics -----
T has no name[845d3716d59:/gem5/gem5/milan/assignment6$ cat m5out/stats.txt

src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall rseq(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall set_robust_list(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall rseq(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall set_robust_list(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall rseq(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall set_robust_list(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall madvise(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall madvise(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall madvise(...)
y[0] = 4.000000
y[1-1] = 2.000000
Exiting @ tick 1383603000 because exiting with last active thread context
I have no name[845d3716d59:/gem5/gem5/milan/assignment6$ cat m5out/stats.txt

----- Begin Simulation Statistics -----
simSeconds 0.001384 # Number of seconds simulated (Second)
simTicks 1383603000 # Number of ticks simulated (Tick)
finalTick 1383603000 # Number of ticks from beginning of simulation (restored from checkpoints and never reset) (Tick)
simFreq 1000000000000 # The number of ticks simulated second ((Tick/Second))
hostSeconds 2.24 # Real time elapsed on the host (Second)
hostTickRate 617921838 # The number of ticks simulated per host second (ticks/s) ((Tick/Second))
hostMemory 8961644 # Number of bytes of host memory used (Byte)
simInsts 534264 # Number of instructions simulated (Count)
simOps 956684 # Number of ops (including micro ops) simulated (Count)
hostInstRate 238426 # Simulator instruction rate (inst/s) ((Count/Second))
hostOpRate 427003 # Simulator op (including micro ops) rate (op/s) ((Count/Second))
system.clk_domain.clock 1000 # Clock period in ticks (Tick)
system.clk_domain.voltage_domain.voltage 1 # Voltage in Volts (Volt)
system.cpu0.numCycles 947237 # Number of cpu cycles simulated (Cycle)
system.cpu0.cpi 3.659815 # CPI: cycles per instruction (core level) ((Cycle/Count))
system.cpu0.ipc 0.273238 # IPC: instructions per cycle (core level) ((Count/Cycle))
system.cpu0.numWorkItemsStarted 0 # Number of work items this cpu started (Count)
system.cpu0.numWorkItemsCompleted 0 # Number of work items this cpu completed (Count)
system.cpu0.quiesceCycles 436366 # Total number of cycles that CPU has spent quiesced or waiting for an interrupt (Cycle)
system.cpu0.branchPred.lookups_0::NoBranch 0 # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0::Return 0 # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0::CallDirect 0 # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0::CallIndirect 0 # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0::DirectCond 0 # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0::DirectUncond 0 # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0::IndirectUncond 0 # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0::total 0 # Number of BP lookups (Count)
system.cpu0.branchPred.squashes_0::NoBranch 0 # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0::Return 0 # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0::CallDirect 0 # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0::CallIndirect 0 # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0::DirectCond 0 # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0::DirectUncond 0 # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0::IndirectUncond 0 # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0::total 0 # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.corrected_0::NoBranch 0 # Number of branches that got corrected but not yet committed. Branches get corrected by decode or after execute. Also a branch mispredi
tected out-of-order. Therefore, a corrected branch might not end up being committed in case an even earlier branch was mispredicted (Count)
system.cpu0.branchPred.corrected_0::Return 0 # Number of branches that got corrected but not yet committed. Branches get corrected by decode or after execute. Also a branch mispredi
cted out-of-order. Therefore, a corrected branch might not end up being committed in case an even earlier branch was mispredicted (Count)
system.cpu0.branchPred.corrected_0::CallDirect 0 # Number of branches that got corrected but not yet committed. Branches get corrected by decode or after execute. Also a branch mispre
detected out-of-order. Therefore, a corrected branch might not end up being committed in case an even earlier branch was mispredicted (Count)
system.cpu0.branchPred.corrected_0::CallIndirect 0 # Number of branches that got corrected but not yet committed. Branches get corrected by decode or after execute. Also a branch misp
e detected out-of-order. Therefore, a corrected branch might not end up being committed in case an even earlier branch was mispredicted (Count)
system.cpu0.branchPred.corrected_0::DirectCond 0 # Number of branches that got corrected but not yet committed. Branches get corrected by decode or after execute. Also a branch mispre
detected out-of-order. Therefore, a corrected branch might not end up being committed in case an even earlier branch was mispredicted (Count)

```

After changing the parameters:

```

gem5 - docker run --name gem5-container --platform linux/arm64 -u 501:20 --volume ~/Documents/Cumberlands/ComputerArchitecture&Organization/gem5/gem5 --rm -it gem5-custom-image -- 230x65

src/sim/syscall_emul.cc:74: warn: ignoring syscall rseq(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/syscall_emul.cc:185: warn: ignoring syscall rt_sigaction(...)
  (further warnings will be suppressed)
src/sim/syscall_emul.cc:185: warn: ignoring syscall rt_sigprocmask(...)
  (further warnings will be suppressed)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/power_state.cc:185: warn: PowerState: Already in the requested power state, request ignored
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall rseq(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall set_robust_list(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall rseq(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall set_robust_list(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall madvise(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall madvise(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall madvise(...)
y(0) = 4.000000
y(N-1) = 2.000000
Exiting @ Tick 1640527000 because exiting with last active thread context
I have no name[845d537161d59]:gem5/gem5/milan/assignment6 cat sOut/stats.txt

```

```

----- Begin Simulation Statistics -----
simSeconds          0.001641          # Number of seconds simulated (Second)
simTicks            1640527000         # Number of ticks simulated (Tick)
finalTick            1640527000         # Number of ticks from beginning of simulation (restored from checkpoints and never reset) (Tick)
simFreq              1000000000000     # The number of ticks per simulated second ((Tick/Second))
hostSeconds          2.45              # Real time elapsed on the host (Second)
hostTickRate         667710831         # The number of ticks simulated per host second (ticks/s) ((Tick/Second))
hostMemory           8941640           # Number of bytes of host memory used (Byte)
simInsts             534166            # Number of instructions simulated (Count)
simOps               956684            # Number of ops (including micro ops) simulated (Count)
hostInstRate         217937           # Simulator instruction rate (inst/s) ((Count/Second))
hostOpRate           390389            # Simulator op (including micro ops) rate (op/s) ((Count/Second))
system_clk_domain.clock 1000          # Clock period in ticks (Tick)
system_clk_domain.voltage 1           # Voltage in Volts (Volt)
system.cpu0.numCycles 1116094          # Number of cpu cycles simulated (Cycle)
system.cpu0.cpi       4.312224         # CPI: cycles per instruction (core level) ((Cycle/Count))
system.cpu0.ipc        0.231899        # IPC: instructions per cycle (core level) ((Count/Cycle))
system.cpu0.numWorkItemsStarted 0      # Number of work items this cpu started (Count)
system.cpu0.numWorkItemsCompleted 0     # Number of work items this cpu completed (Count)
system.cpu0.quiescedCycles 52443       # Total number of cycles that CPU has spent quiesced or waiting for an interrupt (Cycle)
system.cpu0.branchPred.lookups_0:NoBranch 0      # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0:CallIndirect 0    # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0:CallIndirect 0    # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0:DirectCond 0      # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0:DirectUncond 0     # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0:IndirectCond 0     # Number of BP lookups (Count)
system.cpu0.branchPred.lookups_0:IndirectUncond 0    # Number of BP lookups (Count)
system.cpu0.branchPred.squashes_0:NoBranch 0      # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0:Return 0         # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0:CallIndirect 0    # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0:CallIndirect 0    # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0:DirectCond 0      # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0:DirectUncond 0     # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0:IndirectCond 0     # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)
system.cpu0.branchPred.squashes_0:IndirectUncond 0    # Number of branches that got squashed (completely removed) as an earlier branch was mispredicted. (Count)

```

```

gem5 - docker run --name gem5-container --platform linux/arm64 -u 501:20 --volume ~/Documents/Cumberlands/ComputerArchitecture&Organization/gem5/gem5 --rm -it gem5-custom-image -- 230x65

system.membus.pktSize.system.cpu2.icache.mem_side.port:total 34304          # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize.system.cpu2.dcache.mem_side.port:system.mem_ctrl.port 19520          # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize.system.cpu2.dcache.mem_side.port:total 19520          # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize.system.cpu3.icache.mem_side.port:system.mem_ctrl.port 80192          # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize.system.cpu3.icache.mem_side.port:total 80192          # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize.system.cpu3.dcache.mem_side.port:system.mem_ctrl.port 48704          # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize.system.cpu3.dcache.mem_side.port:total 48704          # Cumulative packet size per connected requestor and responder (Byte)
system.membus.pktSize:total 717376    # Cumulative packet size per connected requestor and responder (Byte)
system.membus.snoops 2083             # Total snoops (Count)
system.membus.snoopTraffic 90608       # Total snoop traffic (Byte)
system.membus.snoopFanout:samples 12181 # Request fanout histogram (Count)
system.membus.snoopFanout:mean 0.427155 # Request fanout histogram (Count)
system.membus.snoopFanout:stddev 0.602547 # Request fanout histogram (Count)
system.membus.snoopFanout:underflows 0 0.00% 0.00% # Request fanout histogram (Count)
system.membus.snoopFanout:0 8070 66.69% 66.69% # Request fanout histogram (Count)
system.membus.snoopFanout:1 3093 25.56% 92.25% # Request fanout histogram (Count)
system.membus.snoopFanout:2 738 6.10% 98.35% # Request fanout histogram (Count)
system.membus.snoopFanout:3 280 1.65% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout:4 0 0.00% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout:5 0 0.00% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout:6 0 0.00% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout:7 0 0.00% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout:8 0 0.00% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout:overflows 0 0.00% 100.00% # Request fanout histogram (Count)
system.membus.snoopFanout:min_value 0 # Request fanout histogram (Count)
system.membus.snoopFanout:max_value 3 # Request fanout histogram (Count)
system.membus.snoopFanout:total 12181 # Request fanout histogram (Count)
system.membus.power_state.perStateResidencyTick:UNDEFINED 1640527000 # Cumulative time (in ticks) in various power states (Tick)
system.membus.reqLayer0.occupancy 20121000 # Layer occupancy (ticks) (Tick)
system.membus.reqLayer0.utilization 0.0 # Layer utilization (Ratio)
system.membus.resPlayer0.occupancy 10194250 # Layer occupancy (ticks) (Tick)
system.membus.resPlayer0.utilization 0.0 # Layer utilization (Ratio)
system.membus.resPlayer1.occupancy 21681000 # Layer occupancy (ticks) (Tick)
system.membus.resPlayer1.utilization 0.0 # Layer utilization (Ratio)
system.membus.resPlayer2.occupancy 7541700 # Layer occupancy (ticks) (Tick)
system.membus.resPlayer2.utilization 0.0 # Layer utilization (Ratio)
system.membus.resPlayer3.occupancy 2907250 # Layer occupancy (ticks) (Tick)
system.membus.resPlayer3.utilization 0.0 # Layer utilization (Ratio)
system.membus.resPlayer4.occupancy 4582000 # Layer occupancy (ticks) (Tick)
system.membus.resPlayer4.utilization 0.0 # Layer utilization (Ratio)
system.membus.resPlayer5.occupancy 2863000 # Layer occupancy (ticks) (Tick)
system.membus.resPlayer5.utilization 0.0 # Layer utilization (Ratio)
system.membus.resPlayer6.occupancy 4640750 # Layer occupancy (ticks) (Tick)
system.membus.resPlayer6.utilization 0.0 # Layer utilization (Ratio)
system.membus.resPlayer7.occupancy 6676500 # Layer occupancy (ticks) (Tick)
system.membus.resPlayer7.utilization 0.0 # Layer utilization (Ratio)
system.membus.snoop_filter.hitsSingleRequests 5211 # Number of requests hitting in the snoop filter with a single holder of the requested data. (Count)
system.membus.snoop_filter.hitsMultiRequests 2824 # Number of requests hitting in the snoop filter with multiple (>1) holders of the requested data. (Count)
system.membus.snoop_filter.hitsSingleSnoops 0 # Number of snoops hitting in the snoop filter with a single holder of the requested data. (Count)
system.membus.snoop_filter.hitsMultiSnoops 0 # Number of snoops hitting in the snoop filter with multiple (>1) holders of the requested data. (Count)
system.workload.inst_arm 0 # Number of arm instructions executed (Count)
system.workload.inst_quiesce 0 # Number of quiesce instructions executed (Count)
system.cpu0.idleCycles 214556          # Total number of cycles that the object has spent stopped (Unspecified)
system.cpu0.tickCycles 901538          # Number of cycles that the object actually ticked (Unspecified)
system.cpu1.idleCycles 68743           # Total number of cycles that the object has spent stopped (Unspecified)
system.cpu1.tickCycles 189684          # Number of cycles that the object actually ticked (Unspecified)
system.cpu2.idleCycles 52619           # Total number of cycles that the object has spent stopped (Unspecified)
system.cpu2.tickCycles 189522          # Number of cycles that the object actually ticked (Unspecified)
system.cpu3.idleCycles 106855          # Total number of cycles that the object has spent stopped (Unspecified)
system.cpu3.tickCycles 481396          # Number of cycles that the object actually ticked (Unspecified)
----- End Simulation Statistics -----
I have no name[845d537161d59]:gem5/gem5/milan/assignment6

```

## 4. Performance Metrics and Analysis

Metrics collected included:

- Simulation time

- Parallel speedup compared to single-threaded execution
- IPC and CPI per thread
- Utilization of `FloatSimdFU`

## Results

### 1. Performance Trends Across Configurations

- A significant reduction in simulation time with optimal `opLat` and `issueLat` settings.
- Parallel speedup increased with thread count, with diminishing returns for more than 4 threads due to thread synchronization overhead.

### 2. Tradeoffs Between opLat and issueLat

The choice of `opLat` and `issueLat` influenced both thread-level and instruction-level parallelism. Low `opLat` values (e.g., `opLat = 1`, `issueLat = 6`) resulted in high IPC but limited utilization of the `FloatSimdFU`. Balanced settings (e.g., `opLat = 3`, `issueLat = 4`) offered optimal parallel speedup.

### 1. Implications for TLP

The design of `FloatSimdFU` significantly impacts TLP, especially in workloads with high floating-point demands. Optimal `opLat` and `issueLat` settings maximize parallel speedup without overloading synchronization mechanisms.

### 2. Limitations

The simplicity of `MinorCPU` limits its applicability to complex, out-of-order models. Real-world factors, such as memory bottlenecks and branch prediction, were not explored.

### 3. Future Work

Further studies could incorporate workloads with diverse computational characteristics and explore other factors, such as memory hierarchies and interconnect latency, to evaluate their impact on TLP.

## Conclusion

This study highlights the importance of functional unit design in exploiting TLP on multi-core systems. The findings emphasize the tradeoffs between latency parameters and performance, providing insights for future optimization of multi-threaded applications in gem5 simulations.

## References

- Garcia, A., Kumar, V., & Gupta, S. (2022). Advances in cache coherence for shared-memory architectures. *IEEE Transactions on Parallel and Distributed Systems*, 33\*(4), 790–803.
- Kumar, V., & Gupta, S. (2021). Simplifying parallelism: Programming models and abstractions. *ACM Computing Surveys*, 54\*(7), 1–27.