

# Assignment 4: Exploring Instruction-Level Parallelism (ILP) in Modern Processors

**Milan Bista**

University of Cumberlands

MSCS-531-M50: Computer Architecture and Design

Instructors: Machica McClain / Charles Lively

Github: <https://github.com/mbista25742/MSCS-531-M50-Assignment4>

## **Abstract**

Instruction-Level Parallelism (ILP) is a crucial measure of a processor's capability to execute multiple instructions simultaneously by overlapping their execution. This paper explores the evolution of ILP, its implementation techniques, and the challenges faced in modern processor designs. It also discusses performance metrics related to ILP, current challenges, and future directions for research, highlighting the importance of ILP in enhancing computational throughput without increasing clock frequency.

## **Introduction**

Instruction-Level Parallelism (ILP) is defined as the ability of a processor to execute multiple instructions concurrently, capitalizing on the inherent parallelism found in instruction sequences. By exploiting the independence of instructions, ILP maximizes computational throughput within a single processor core. Increasing the number of instructions processed per clock cycle enhances performance without necessitating an increase in clock frequency, which is vital due to physical and thermal constraints on clock speeds.

## **Techniques for Implementing ILP**

ILP is realized through various techniques, including:

1. **Pipelining:** This technique allows different stages of instruction execution to occur simultaneously, thereby increasing efficiency.
2. **Superscalar Execution:** In this approach, multiple instructions are dispatched to parallel pipelines, enabling multiple instructions to be processed per clock cycle.
3. **Out-of-Order Execution:** This allows instructions to be executed as soon as their required data becomes available, rather than strictly following program order, thereby improving throughput.
4. **Speculative Execution:** Processors preemptively execute instructions based on predicted paths of conditional branches, maintaining a steady flow of operations.

Each of these techniques addresses independent instructions that can be processed concurrently, thus reducing idle cycles and enhancing overall efficiency.

## **Dependencies in Instruction Streams**

ILP is fundamentally reliant on the identification and management of dependencies within instruction streams, which include:

- **Data Dependencies:** Occur when one instruction requires the output of a previous instruction (e.g., read-after-write).
- **Control Dependencies:** Arise from branches or conditional instructions affecting the flow of execution.
- **Resource Conflicts:** Occur when hardware resources are insufficient for parallel execution.

By addressing these dependencies, ILP has become essential in modern processor design, enhancing performance across various applications.

### **Historical Context of ILP**

The evolution of ILP has significantly advanced modern processor design, beginning with basic pipelining techniques in the 1960s. Early pipeline architectures, such as the IBM Stretch and CDC 6600, laid the groundwork by allowing multiple stages of instruction execution to overlap.

In the 1980s, the introduction of superscalar architectures, exemplified by the IBM System/360 Model 91, allowed processors to execute multiple instructions per clock cycle. This marked a significant shift in how processors utilized parallelism within a single thread of instructions.

The 1990s saw the introduction of out-of-order execution in high-performance processors like the Intel Pentium Pro, enabling instructions to be processed as their data became available. This adaptive management of instruction sequences further increased ILP.

### **Techniques for Exploiting ILP**

Modern processors leverage various techniques to maximize ILP:

- **Dynamic Scheduling and Out-of-Order Execution:** Allow instructions to be executed as soon as dependencies are resolved.
- **Register Renaming:** Helps resolve false dependencies by assigning unique hardware registers to instructions.
- **Branch Prediction and Speculative Execution:** Predicts the outcomes of branches and executes instructions along predicted paths.
- **Superscalar Execution:** Dispatches multiple instructions per clock cycle across functional units.
- **Vectorization:** Supports parallel execution within instructions, ideal for data-heavy operations.

### **Challenges in ILP**

Despite its benefits, ILP faces several challenges:

1. **Increasing Complexity:** Modern processors' architectures have become complex, making effective ILP implementation challenging.
2. **Diminishing Returns:** Traditional ILP techniques have shown diminishing returns as transistor scaling slows down.
3. **Power Constraints:** High ILP designs can lead to increased power usage and thermal issues.
4. **Memory Bottlenecks:** High demand for memory bandwidth can hinder the full exploitation of ILP.

### Addressing ILP Challenges

To address these challenges, researchers are exploring various techniques:

- **Heterogeneous Architectures:** Integrating specialized cores alongside traditional CPU cores for optimized execution.
- **Machine Learning-Based Optimizations:** Utilizing machine learning to enhance scheduling, resource allocation, and branch prediction.
- **Dynamic Voltage and Frequency Scaling (DVFS):** Adjusting power consumption dynamically based on workload requirements.
- **Novel Instruction Set Architectures (ISAs):** Developing ISAs that support efficient parallel execution.

### Future Directions for ILP Research

Emerging trends suggest promising directions for future ILP research, including:

1. **Specialized Accelerators:** The rise of domain-specific accelerators for workloads like AI is likely to redefine ILP considerations.
2. **Quantum Computing:** The implications of quantum algorithms on ILP could yield significant insights.
3. **Enhanced Memory Architectures:** Innovations in memory technologies could alleviate bottlenecks and enable better ILP exploitation.
4. **Software-Driven Approaches:** Sophisticated software tools and compilers can optimize instruction scheduling and resource management.
5. **Energy-Efficient Designs:** Research into energy-efficient designs that maintain high performance will be crucial as power constraints become increasingly pressing.

## **Conclusion**

The evolution of Instruction-Level Parallelism (ILP) reflects a continual effort to maximize computational efficiency in processors, adapting to the evolving demands of hardware and software. This journey has driven innovations that have not only improved processor performance but also laid the groundwork for the complex architectures seen in today's computing landscape. As challenges persist, ongoing research will play a vital role in harnessing ILP to meet future computing needs.

## **Part 2: Practical Exploration of ILP Techniques**

In this part, I am going to simulate Instruction Level Parallelism in gem5

### **gem5 Configuration:**

I have created a simple out of order cpu model that will help us simulate and visualize fetch decode and execute instructions

```
EXPLORER
...
out_of_order_cpu.py 9+, U X
C hello.c .../assignment3 U
C hello.c .../assignment4 U

DEM5
gem5
> .devcontainer
> .github
> .vscode
> build
> build_opts
> build_tools
> configs
> ext
> include
> m5out
> milan
  > assignment3
  > assignment4
    > m5out
    ≡ copyrighted software U
    ≡ hello U
    ≡ C hello.c U
    ≡ out_of_order_cpu.py 9+, U
    ≡ pipeview.out U
    ≡ sing trace... done! U
  > customUtils
> milan-local-resources
> milan-test
> residency
> site_scons
> src

out_of_order_cpu.py
18 class L1DCache(L1Cache):
19     size = '128kB'
20
21 #Create the system for our architecture simulation
22 system = System()
23 system.clk_domain = SrcClockDomain()
24 system.clk_domain.clock = '1GHz'
25 system.clk_domain.voltage_domain = VoltageDomain()
26
27 # Set the block size (cache line size) for the entire system
28 system.cache_line_size = 128
29
30 #Simple CPU and Memory Config
31 system.mem_mode = 'timing'
32 system.mem_ranges = [AddrRange('8192MB')]
33 system.cpu = Deriv03CPU() # Out-of-order CPU model
34 system.cpu.numThreads = 1
35 # Configure ILP parameters for Deriv03CPU
36 system.cpu.fetchWidth = 4 # Number of instructions fetched per cycle
37 system.cpu.decodeWidth = 4 # Number of instructions decoded per cycle
38 system.cpu.issueWidth = 4 # Number of instructions issued per cycle
39 system.cpu.commitWidth = 4 # Number of instructions committed per cycle
40 system.cpu.dispatchWidth = 4 # Number of instructions dispatched per cycle
41
42
43 system.membus = SystemXBar()
44
45 #Memory Control
46 system.mem_ctrl = MemCtrl()
47 system.mem_ctrl.dram = DDR3_1600_8x8()
48 system.mem_ctrl.dram.range = system.mem_ranges[0]
49 system.mem_ctrl.port = system.membus.mem_side_ports
50
51 #Creating a Simple L1 Cache System
```

**Run Simulation:** I have created a simple binary hello program and ran the script as follows

```
../build/x86/gem5.opt --debug-flags=O3PipeView --debug-file=trace.out out_of_order_cpu.py -  
c hello
```

```
gem5 — docker run --name gem5-container --platform linux/arm64 -u 501:20 --volume ~/Documents/Cumberlands/ComputerArchitecture&Organization/gem5:/gem5 --f...
I have no name!@950eff622438:/gem5/gem5/milan/assignment4$ ../../build/x86/gem5.opt --debug-flags=O3PipeView --debug-file=trace.out out_of_order_cpu.py -c hello
gem5 Simulator System: https://www.gem5.org
gem5 is copyrighted software; use the --copyright option for details.

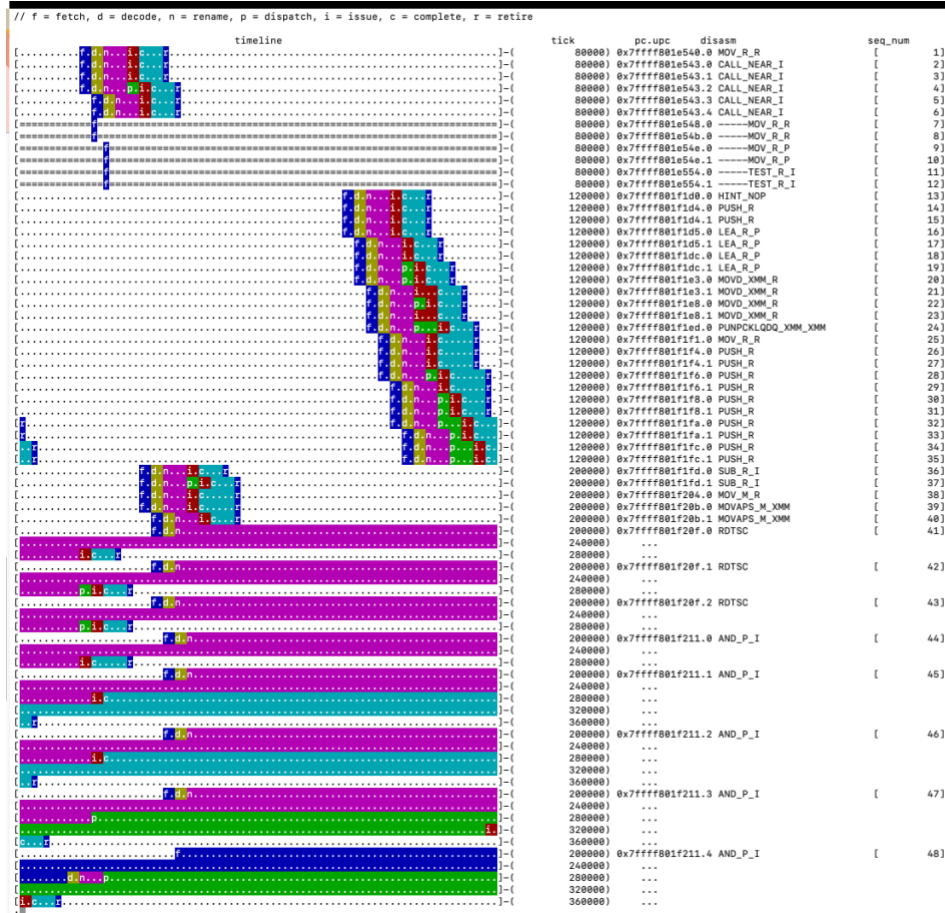
gem5 version 24.0.0.1
gem5 compiled Sep  6 2024 17:16:33
gem5 started Oct 26 2024 19:52:08
gem5 executing on 950eff622438, pid 194
command line: ../../build/x86/gem5.opt --debug-flags=O3PipeView --debug-file=trace.out out_of_order_cpu.py -c hello

Global frequency set at 1000000000000 ticks per second
src/base/statistics.hh:279: warn: One of the stats is a legacy stat. Legacy stat is a stat that does not belong to any statistics::Group. Legacy stat is deprecated.
system.remote_gdb: Listening for connections on port 7000
Beginning of the Simulation!!!!
src/sim/simulate.cc:199: info: Entering event queue @ 0. Starting simulation...
src/sim/mem_state.cc:448: info: Increasing stack size by one page.
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 mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
Hello MilanExiting @ tick 159645000 because exiting with last active thread context
I have no name!@950eff622438:/gem5/gem5/milan/assignment4$
```

**Visualize:**

I have used gem5 graphical pipeline viewer tool from utils class to visualize the output as follows

```
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(...)
Hello MilanExiting @ tick 159645000 because exiting with last active thread context
I have no name!@950eff622438:/gem5/gem5/milan/assignment4$ ../../util/o3-pipeview.py -c 500 -o pipeview.out --color m5out/trace.out
Processing trace... done!
I have no name!@950eff622438:/gem5/gem5/milan/assignment4$
```



## Cpu Instructions Per cycle can be visualized as

Hello MilanExiting @ tick 159645000 because exiting with last active thread context  
I have no name@958eff622438:/gem5/gem5/milan/assignment4\$ cat m5out/stats.txt

Begin Simulation Statistics			
simSeconds	0.000160	# Number of seconds simulated (Second)	
simTicks	159645000	# Number of ticks simulated (Tick)	
finalTick	159645000	# 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	98.50	# Real time elapsed on the host (Second)	
hostTickRate	1620786	# The number of ticks simulated per host second (ticks/s) ((Tick/Second))	
hostMemory	8877668	# Number of bytes of host memory used (Byte)	
simInsts	96840	# Number of instructions simulated (Count)	
simOps	186471	# Number of ops (including micro ops) simulated (Count)	
hostInstRate	983	# Simulator instruction rate (inst/s) ((Count/Second))	
hostOpRate	1893	# 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.cpu.numCycles	159646	# Number of cpu cycles simulated (Cycle)	
system.cpu.cpi	1.648554	# CPI: cycles per instruction (core level) ((Cycle/Count))	
system.cpu.ipc	0.606592	# IPC: instructions per cycle (core level) ((Count/Cycle))	
system.cpu.numWorkItemsStarted	0	# Number of work items this cpu started (Count)	
system.cpu.numWorkItemsCompleted	0	# Number of work items this cpu completed (Count)	
system.cpu.instsAdded	245212	# Number of instructions added to the IQ (excludes non-spec) (Count)	
system.cpu.nonSpecInstsAdded	68	# Number of non-speculative instructions added to the IQ (Count)	
system.cpu.instsIssued	229073	# Number of instructions issued (Count)	
system.cpu.squashedInstsIssued	932	# Number of squashed instructions issued (Count)	
system.cpu.squashedInstsExamined	58863	# Number of squashed instructions iterated over during squash; mainly for profiling (Count)	
system.cpu.squashedOperandsExamined	92622	# Number of squashed operands that are examined and possibly removed from graph (Count)	
system.cpu.squashedNonSpecRemoved	32	# Number of squashed non-spec instructions that were removed (Count)	
system.cpu.numIssueDist::samples	135035	# Number of insts issued each cycle (Count)	
system.cpu.numIssueDist::mean	1.694397	# Number of insts issued each cycle (Count)	
system.cpu.numIssueDist::stdev	1.614741	# Number of insts issued each cycle (Count)	
system.cpu.numIssueDist::underflows	0	0.00%	0.00% # Number of insts issued each cycle (Count)

I have created a custom script to read stats.txt to analyze the metrics as



```
t_of_order_cpu.py 9+, U  analyze_stats.py U x  C hello.c .../assignment3 U  C hello.c .../assignment4 >
gem5 > milan > assignment4 > analyze_stats.py > ...
3
4 # File path to stats.txt
5 stats_file = "m5out/stats.txt"
6
7 ipc = 0
8 total_cycles = 0
9 total_insts = 0
10
11 with open(stats_file, "r") as file:
12     for line in file:
13         # Find the IPC
14         if "system.cpu.ipc" in line:
15             ipc = float(line.split()[1])
16             print(f"Instruction Per Cycle (IPC): {ipc}")
17
18         # Find the total cycles
19         if "system.cpu.numCycles" in line:
20             total_cycles = int(line.split()[1])
21             print(f"Number of Cycles: {total_cycles}")
22
23         # Find the total instructions
24         if "system.cpu.instsIssued" in line:
25             total_insts = int(line.split()[1])
26             print(f"Total Number of Instructions {total_insts}")
27
28 # Calculate latency in cycles
29 if total_insts > 0:
30     latency = total_cycles / total_insts
31 else:
32     latency = None
33
34 # Output results
35 print(f"Instruction Throughput (IPC): {ipc}")
36 print(f"Instruction Latency (Cycles per Instruction): {latency}")
```

```
gem5 — docker run --name gem5-container --platform linux/arm64 -u 501:20 --volume ~/Documents
Q system.cpu.ipc
[I have no name!@950eff622438:/gem5/gem5/milan/assignment4$ python3 analyze_stats.py
Number of Cycles: 159646
Instruction Per Cycle (IPC): 0.606592
Total Number of Instructions 229073
Instruction Throughput (IPC): 0.606592
Instruction Latency (Cycles per Instruction): 0.6969219419137131
I have no name!@950eff622438:/gem5/gem5/milan/assignment4$
```

## Impact of Branch Prediction

Now I am going to add a branch predictor to the above configuration

```
27
28 # Set the block size (cache line size) for the entire system
29 system.cache_line_size = 128
30
31 # Simple CPU and Memory Config
32 system.mem_mode = 'timing'
33 system.mem_ranges = [AddrRange('8192MB')]
34 system.cpu = Deriv03CPU() # Out-of-order CPU model
35 system.cpu.numThreads = 1
36
37 # Configure ILP parameters for Deriv03CPU
38 system.cpu.fetchWidth = 4 # Number of instructions fetched per cycle
39 system.cpu.decodeWidth = 4 # Number of instructions decoded per cycle
40 system.cpu.issueWidth = 4 # Number of instructions issued per cycle
41 system.cpu.commitWidth = 4 # Number of instructions committed per cycle
42 system.cpu.dispatchWidth = 4 # Number of instructions dispatched per cycle
43
44 # Add a simple branch predictor
45 system.cpu.branchPred.loopPredictor = LoopPredictor()
46 system.cpu.branchPred.globalPredictorSize = 16384
47
48 system.membus = SystemXBar()
49
50 # Memory Control
51 system.mem_ctrl = MemCtrl()
52 system.mem_ctrl.dram = DDR3_1600_8x8()
53 system.mem_ctrl.dram.range = system.mem_ranges[0]
54 system.mem_ctrl.port = system.membus.mem_side_ports
55
56 # Creating a Simple L1 Cache System
57 system.cpu.icache = L1ICache() # L1 Instruction Cache
58 system.cpu.dcache = L1DCache() # L1 Data Cache
59
```

```
gem5 — docker run --name gem5-container --platform linux/arm64 -u 501:20 --volume ~/Documents/Cumberla
```

```
Q predic
```

```
I have no name!@950eff622438:/gem5/gem5/milan/assignment4$ ../../build/x86/gem5.opt --debug-flags=O3PipeView --d
gem5 Simulator System.  https://www.gem5.org
gem5 is copyrighted software; use the --copyright option for details.

gem5 version 24.0.0.1
gem5 compiled Sep  6 2024 17:16:33
gem5 started Oct 27 2024 19:09:42
gem5 executing on 950eff622438, pid 354
command line: ../../build/x86/gem5.opt --debug-flags=O3PipeView --debug-file=trace.out out_of_order_cpu_bp.py -c

Global frequency set at 1000000000000 ticks per second
src/base/statistics.hh:279: warn: One of the stats is a legacy stat. Legacy stat is a stat that does not belong t
system.remote_gdb: Listening for connections on port 7000
Beginning of the Simulation!!!!
src/sim/simulate.cc:199: info: Entering event queue @ 0. Starting simulation...
src/sim/mem_state.cc:448: info: Increasing stack size by one page.
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 mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
src/sim/syscall_emul.cc:74: warn: ignoring syscall mprotect(...)
Printng 10 times 1
Printng 10 times 2
Printng 10 times 3
Printng 10 times 4
Printng 10 times 5
Printng 10 times 6
Printng 10 times 7
Printng 10 times 8
Printng 10 times 9
Printng 10 times 10
Exiting @ tick 169490000 because exiting with last active thread context
I have no name!@950eff622438:/gem5/gem5/milan/assignment4$ python3 analyze_stats.py
Number of Cycles: 169491
Instruction Per Cycle (IPC): 0.624428
Total Number of Instructions 247940
Instruction Throughput (IPC): 0.624428
Instruction Latency (Cycles per Instruction): 0.6835968379446641
I have no name!@950eff622438:/gem5/gem5/milan/assignment4$ █
```

We implemented the DerivO3CPU model in our gem5 simulation, which provided notable performance improvements. By leveraging its out-of-order execution capabilities, we observed enhanced instruction throughput and better utilization of resources, allowing for higher instruction-level parallelism (ILP) compared to simpler CPU models. The advanced features of DerivO3CPU, such as dynamic scheduling and branch prediction, contributed significantly to reducing execution time and improving overall simulation efficiency.

Performance can be improved by using Simultaneous Multithreading (SMT), which enables multiple threads to share the same processor resources. With SMT, the CPU can execute instructions from different threads concurrently, maximizing the utilization of functional units, registers, and pipeline stages. This approach reduces idle cycles and increases instruction throughput, allowing the processor to complete more instructions per cycle. As a result, SMT enhances overall system efficiency, particularly in workloads with parallel threads, by minimizing resource bottlenecks and improving execution speed.

```
system.cpu.fetchWidth = 4      # Number of instructions fetched per cycle
system.cpu.decodeWidth = 4     # Number of instructions decoded per cycle
system.cpu.issueWidth = 4      # Number of instructions issued per cycle
system.cpu.commitWidth = 4     # Number of instructions committed per cycle
system.cpu.dispatchWidth = 4   # Number of instructions dispatched per cycle

# Add a simple branch predictor
system.cpu.branchPred.loopPredictor = LoopPredictor()
system.cpu.branchPred.globalPredictorSize = 16384

# system.cpu.smtNumFetchingThreads=
system.cpu.smtLSQThreshold=200
system.cpu.smtIQThreshold = 200
system.cpu.smtROBThreshold = 200
system.cpu.numThreads = 4

system.membus = SystemXBar()

# Memory Control
system.mem_ctrl = MemCtrl()
system.mem_ctrl.dram = DDR3_1600_8x8()
system.mem_ctrl.dram.range = system.mem_ranges[0]
system.mem_ctrl.port = system.membus.mem_side_ports

# Creating a Simple L1 Cache System
system.cpu.icache = L1ICache() # L1 Instruction Cache
```

Throughout these experiments, it's essential to consider how different Instruction-Level Parallelism (ILP) techniques, such as pipelining, out-of-order execution, and branch prediction, interact to improve performance. These techniques often complement each other by allowing

multiple instructions to be executed simultaneously and addressing data dependencies more efficiently. However, limitations arise due to factors like instruction dependencies, pipeline hazards, and resource contention, which can restrict the degree of achievable parallelism. Balancing the complexity of ILP techniques with performance benefits requires evaluating each technique's impact on execution time and hardware resources. Achieving optimal ILP often means selectively implementing techniques that maximize throughput without excessively increasing power consumption or design complexity.

## References

- Hennessy, J. L., & Patterson, D. A. (2017). *Computer architecture: A quantitative approach* (6th ed.). Morgan Kaufmann.
- Stallings, W. (2021). *Computer organization and architecture* (11th ed.). Pearson.
- Rabaey, J. M., Chandrakasan, A. P., & Nikolic, B. (2022). *Digital integrated circuits: A design perspective* (3rd ed.). Pearson.