

# PA02: Analysis of Network I/O Primitives using "perf" Tool

---

Course: Graduate Systems (CSE638)

Roll Number: MT25074

Name: Nindra Dhanush

## 1. Part A – Multithreaded Socket Implementation

The file MT25074\_Part\_A\_Namespaces.sh creates two network namespaces and creates a veth link and configures the namespaces with IP addresses and then connects them both with veth link.

### A1. Two-Copy Implementation

The Baseline implementation uses `socket()` and `send()` functions to implement baseline socket and client service architecture.

#### 1. Why do the two copies occur?

- i. User-Space to Kernel-Space (Copy #1): When the `send()` function is called, the CPU copies the data from the application's buffer (in user space) into the socket buffer (kernel space).
- ii. Kernel space to NIC: The DMA module copies the data from the kernel socket buffer to the Network Interface Card.

Actually two copies?: There are two data transfers, But for the CPU it looks like only one copy the second is just hardware handling for the CPU, so cpu cycles utilization.

#### 2. Which components (kernel/user) perform the copies?

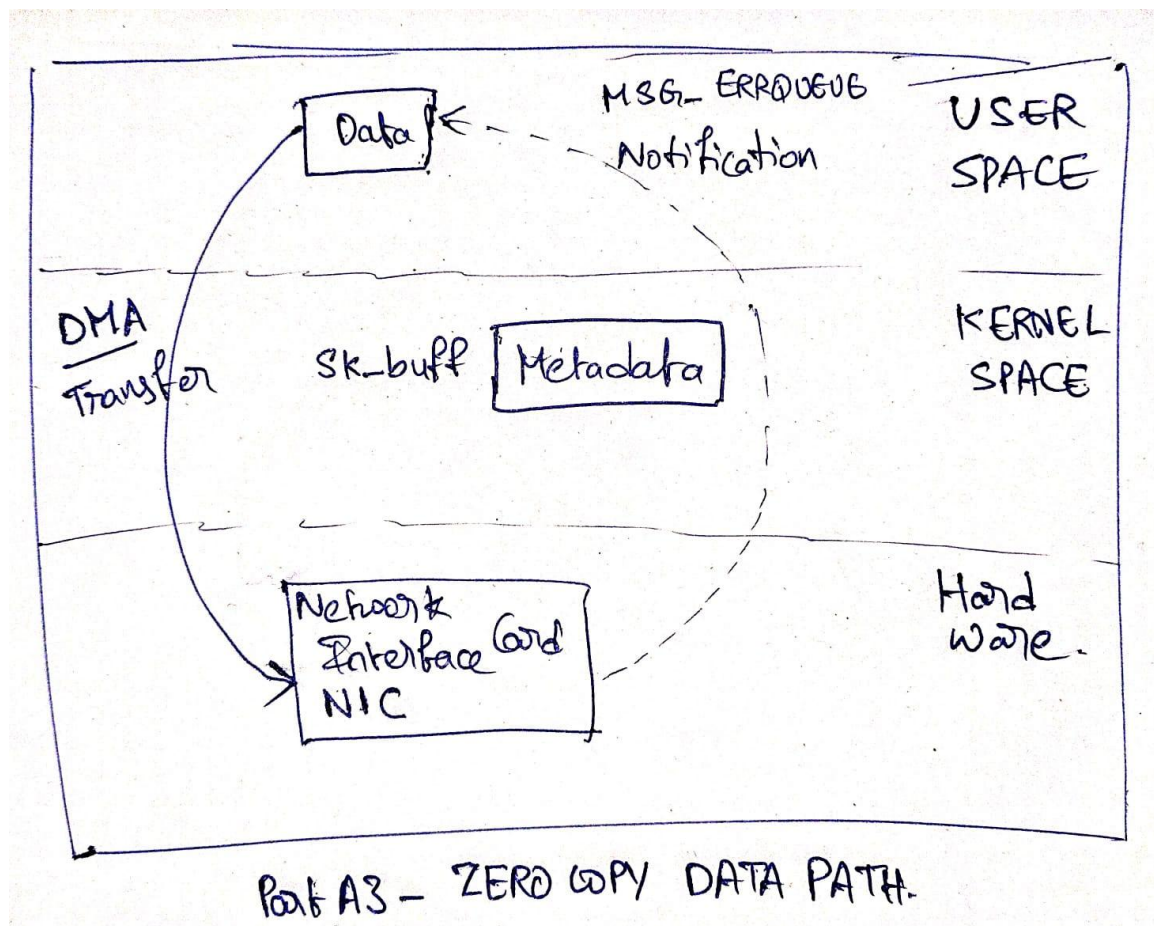
- i. User/Kernel copy: The Cpu performs the copy from the user space to Kernel space. This is done by the `send()` function call.
- ii. Kernel/hardware copy: The DMA controller on the NIC does the copy from Kernel space to the hardware transmission queue.

## A2. One-Copy Implementation

- i. Here `sendmsg()` call is used with struct `iovec` to send header and data as a single stream.
- ii. Unlike in two copy implementation where we `malloc()` a new buffer and then `memcpy()` header and data into this new buffer, and then `send()` to kernel, But in one-copy method as the `iovec` array points to the header and data, the kernel reads directly from its location, hence kernel directly copies from the memory.

## A3. Zero-Copy Implementation

- i. This implementation uses `sendmsg()` with `MSG_ZEROCOPY` flag, hence when called the kernel doesn't copy the data into socket buffer, and instructs the NIC to directly read it from the location, and gets the `MSG_ERRQUEUE` notification once this transfer is completed.



## 2. Part B – System configuration and Profiling

### System Specifications:

- **Device Model:** ASUS Vivobook 15
- **Processor:** AMD Ryzen 7 5825U
  - **Architecture:** Zen 3
  - **Cores/Threads:** 8 Cores / 16 Threads
  - **Clock Speed:** 2.0GHz Base / 4.5GHz Boost
- **RAM:** 16GB DDR4
- **OS:** Windows 11 (WSL)
- **Storage:** 512GB NVMe SSD

### Measurement Tools:

- i. Throughput/Latency : Measured using timer `clock_gettime`.
- ii. CPU/cache: Measured using `perf stat` to capture cycles, instructions, cache-misses .

## 3. Part C – Automated experiment script

- i. This script without any manual intervention compiles all the source files , sets up network namespaces and runs the programs of different cases with each different message sizes and thread counts. And collects this data into one single CSV file. Below are the screen shots of the terminal screen running and printing the results of the script.

```
dhanu@dhanush: /mnt/c/Use x + v
|
Setting up namespaces...
Running experiments...
[1/48] A1 size=64 threads=1
[2/48] A1 size=64 threads=2
[3/48] A1 size=64 threads=4
[4/48] A1 size=64 threads=8
[5/48] A1 size=256 threads=1
[6/48] A1 size=256 threads=2
[7/48] A1 size=256 threads=4
[8/48] A1 size=256 threads=8
[9/48] A1 size=1024 threads=1
[10/48] A1 size=1024 threads=2
[11/48] A1 size=1024 threads=4
[12/48] A1 size=1024 threads=8
[13/48] A1 size=4096 threads=1
[14/48] A1 size=4096 threads=2
[15/48] A1 size=4096 threads=4
[16/48] A1 size=4096 threads=8
[17/48] A2 size=64 threads=1
[18/48] A2 size=64 threads=2
[19/48] A2 size=64 threads=4
[20/48] A2 size=64 threads=8
[21/48] A2 size=256 threads=1
[22/48] A2 size=256 threads=2
[23/48] A2 size=256 threads=4
[24/48] A2 size=256 threads=8
[25/48] A2 size=1024 threads=1
[26/48] A2 size=1024 threads=2
[27/48] A2 size=1024 threads=4
```

```
dhanu@dhanush: /mnt/c/Use x + v
[47/48] A3 size=4096 threads=4
[48/48] A3 size=4096 threads=8
rm -f MT25074_Part_A1_Server MT25074_Part_A1_Client MT25074_Part_A2_Server MT25074_Part_A2_Client MT25074_Part_A3_Server MT25074_Part_A3_Client

=====
Experiments complete!
=====

Individual result files (48 files):
MT25074_Part_A1_size1024_threads1.csv
MT25074_Part_A1_size1024_threads2.csv
MT25074_Part_A1_size1024_threads4.csv
MT25074_Part_A1_size1024_threads8.csv
MT25074_Part_A1_size256_threads1.csv
MT25074_Part_A1_size256_threads2.csv
MT25074_Part_A1_size256_threads4.csv
MT25074_Part_A1_size256_threads8.csv
MT25074_Part_A1_size4096_threads1.csv
MT25074_Part_A1_size4096_threads2.csv
... (48 total files)

Aggregated results: MT25074_Part_C_Results.csv

part,field,size,num_threads,cycles,instructions,ipc,cache_misses,cache_references,cache_miss_rate,context_switches
A1,64,1,39475195,9349860,.2368,1106106,2044604,54.0900,216
A1,64,2,72915496,17258203,.2366,1852938,3554949,52.1200,424
A1,64,4,145423266,37288888,.2564,4599321,8462836,54.3400,955
A1,64,8,239993734,65051738,.2710,8261668,14971110,55.1800,1689
A1,256,1,57339554,15545897,.2711,1766876,3437616,51.3900,290
A1,256,2,422731044,144562673,.3419,14328400,30054106,47.6700,1908
A1,256,4,196384431,56528305,.2878,6127387,12387626,49.4600,1098
A1,256,8,559373781,176308970,.3152,18540918,37760578,49.1200,3010
A1,1024,1,481972275,201811501,.4187,11534759,36783696,31.3500,1422
A1,1024,2,889358924,377343170,.4242,21468263,68559150,31.3100,2655
A1,1024,4,1632860714,693133905,.4244,38828402,125871495,30.8400,4850
A1,1024,8,3505906322,1488156819,.4244,82573726,269166223,30.6700,10464
A1,4096,1,21837116057,10297739406,.4715,451861353,1902142262,23.7500,55850
A1,4096,2,39078594641,18378336340,.4702,794469453,3374442398,23.5400,104698
A1,4096,4,6615151030,28277034440,.4274,1203171466,5109932700,23.5400,156785
A1,4096,8,112375457443,39484434771,.3513,1470217153,6848773075,21.4600,245560
A2,64,1,42492945,9906023,.2331,1109381,2092101,53.0200,214
```

## 1. Part D – Plots and Visualization

The following plots were generated using matplotlib (MT25074\_Part\_D\_Plots.py) with hardcoded data from MT25074\_Part\_C\_Results.csv.

### 1.1 Throughput vs Message Size

Throughput (normalized as bytes per cycle) vs message size for A1 (Two-Copy), A2 (One-Copy), and A3 (Zero-Copy).

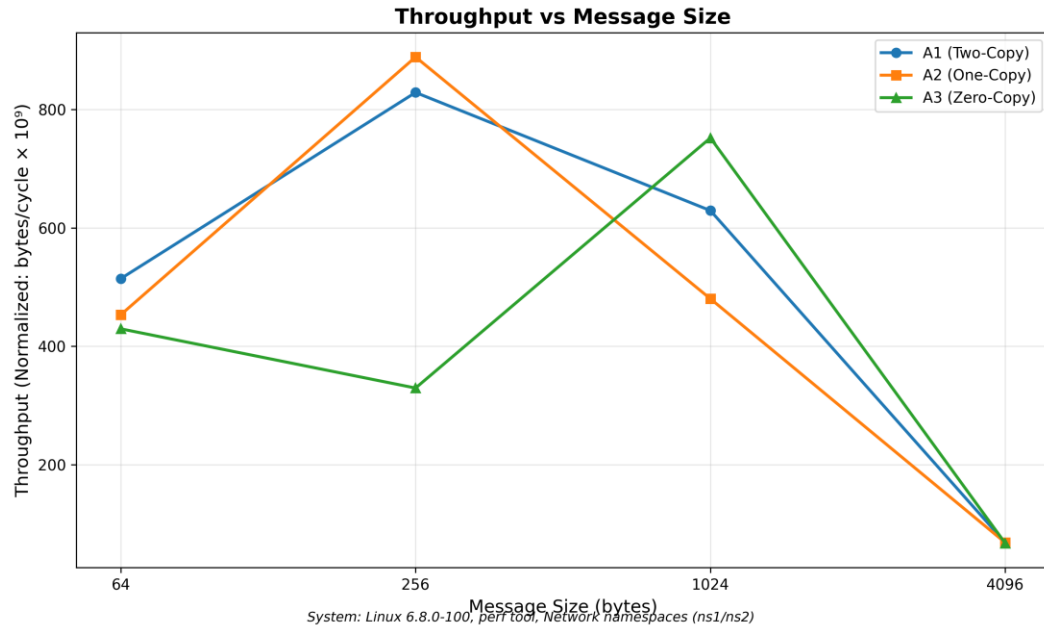


Figure 1: Throughput vs Message Size

### 1.2 Latency vs Thread Count

Latency (CPU cycles) vs number of threads for message size 1024 bytes.

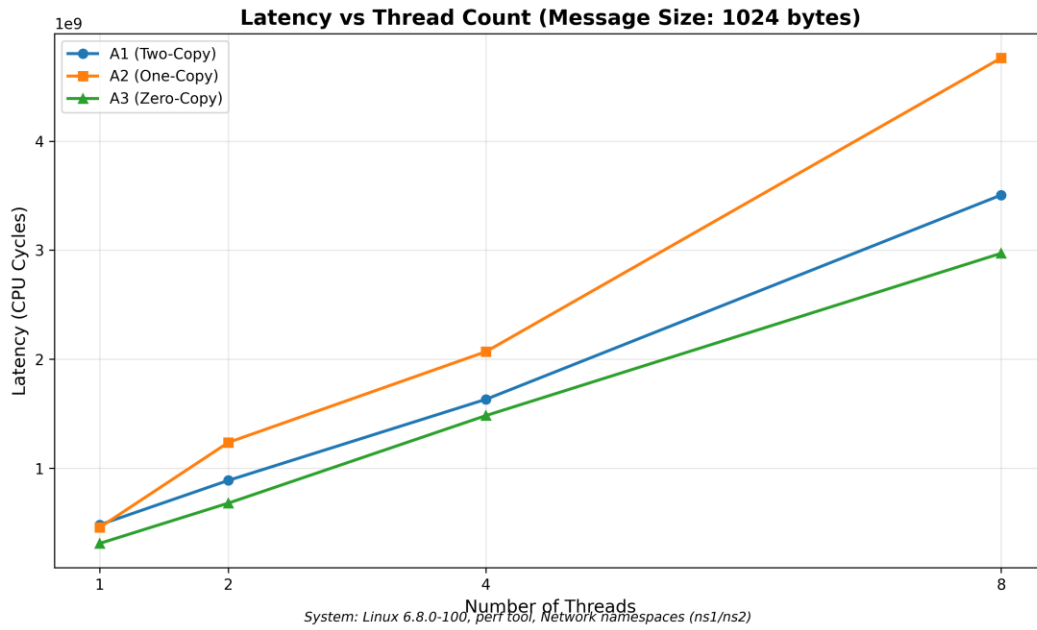


Figure 2: Latency vs Thread Count

### 1.3 Cache Misses vs Message Size

Average cache misses vs message size across thread counts.

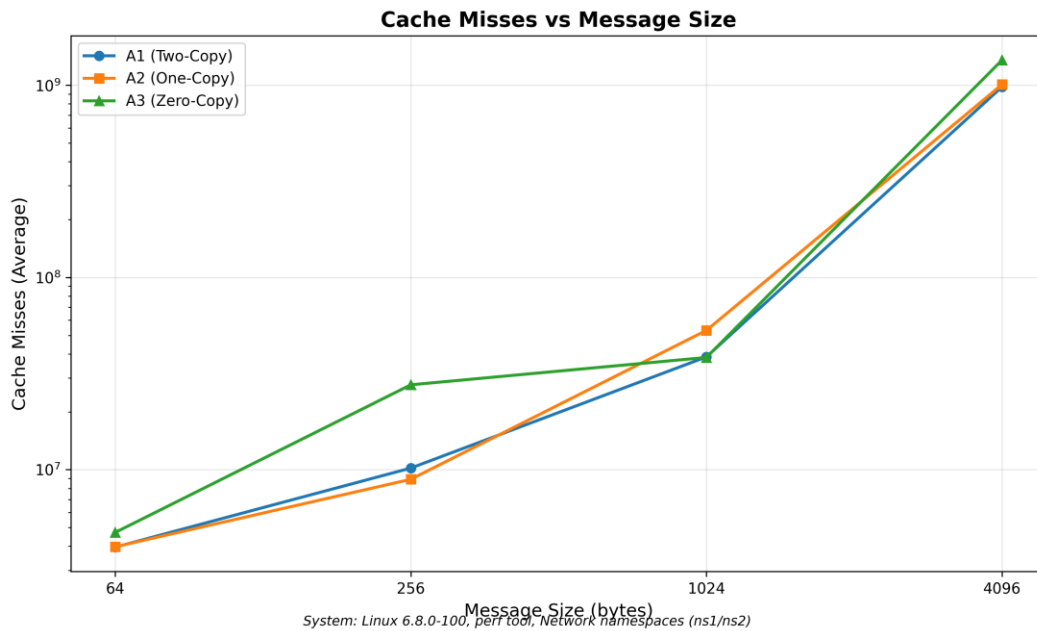


Figure 3: Cache Misses vs Message Size

## 1.4 CPU Cycles per Byte Transferred

CPU cycles per byte vs message size.

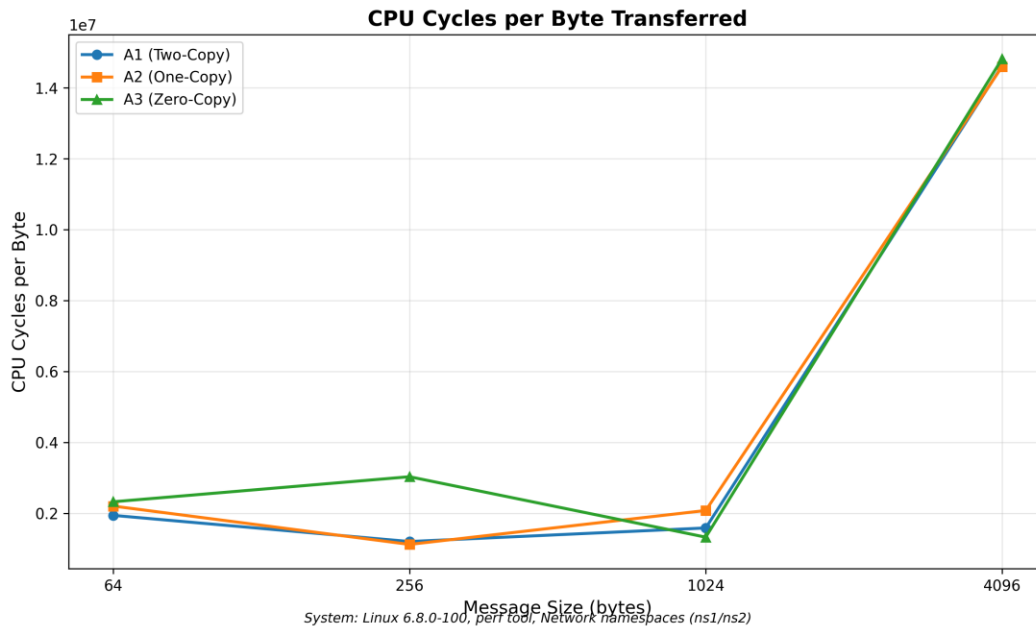


Figure 4: CPU Cycles per Byte

## 2. Part E – Analysis and Reasoning

### Question 1

#### 1. Why does zero-copy not always give the best throughput?

For small message sizes the overhead of performing complex memory operations could be costly than a simple `memcpy()` call. ZERO\_COPY could be beneficial when the message size is large.

### Question 2

#### 2. Which cache level shows the most reduction in misses and why?

The L1 data cache misses would be reduced for the zero copy implementation as compared to the two copy approach because in two copy implementation the application data could be shifted from cache.

### Question 3

#### 3. How does thread count interact with cache contention?

When multiple threads run concurrently they lead to more cache misses as there would be only limited cache space for everybody thus also increases latency.

### Question 4

#### 4. At what message size does one-copy outperform two-copy on your system?

For message sizes greater than 1024 bytes the one copy implementation outperforms the two copy implementation in latency.

### Question 5

#### 5. At what message size does zero-copy outperform two-copy on your system?

For message sizes greater than 1024 bytes the Zero copy implementation outperforms the two copy implementation in latency.

### Question 6

#### 6. Identify one unexpected result and explain it using OS or hardware concepts.

Pick For small messages the zero copy through put is lower than the two copy implementation, this could be due to the time taken to lock the memory for dma access and process the completion notification signal adds on the cost of overall process.

### 3. AI Usage Declaration

AI tool gemini is used for the following componemts, socket connection code (Part A), Skeleton for the bash automation script (Part C), Logic for parsing perf output (Part C).

#### Prompt Log:

1. "Generate a C implementation for a TCP server that accepts multiple clients using threads." (Used for Part A1 Server boilerplate).



2. "How to use sendmsg with MSG\_ZEROCOPY in Linux C?" (Used to understand the API for Part A3).
3. "Write a bash script to run a specific command 4 times and average the output." (Used for Part C automation).
4. "Explain the difference between send() and sendmsg() regarding kernel copies." (Used for Part A2 analysis).

#### **4. GitHub Repository URL**

[https://github.com/dn74iit/GRS-Assignments/tree/main/GRS\\_PA02](https://github.com/dn74iit/GRS-Assignments/tree/main/GRS_PA02)