

# Part E: Analysis and Reasoning

**Student Name:** Dewansh Khandelwal  
**Roll Number:** MT25067  
**System Configuration:** Ubuntu 22.04.5 LTS, Intel Core i7-12700  
**Network Setup:** Linux Network Namespaces (server\_ns ↔ client\_ns via veth pair)

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

### Observation

In the network namespace experiments using veth interfaces, Zero-Copy (A3) consistently underperformed Two-Copy (A1):

Message Size	Implementation	Throughput (Gbps)	Performance
4KB, 4T	A1 (Two-Copy)	9.91	✓ Baseline
4KB, 4T	A2 (One-Copy)	8.95	-10%
4KB, 4T	A3 (Zero-Copy)	6.75	-32%

### Root Cause Analysis

#### 1. Page Pinning Overhead

User sends data → Kernel must pin pages (get\_user\_pages)  
↓  
Prevents page swapping while NIC accesses memory  
↓  
Expensive TLB operations + reference counting

Every MSG\_ZEROCOPY send requires:

- Walking page tables to pin virtual memory
- Incrementing atomic reference counts
- Setting page flags (PG\_locked, PG\_writeback)

#### 2. veth Virtual Interface = No Hardware DMA

Physical NIC scenario:  
User buffer → DMA → NIC  
(True zero-copy)

veth (Virtual Ethernet) scenario:  
User buffer → memcpy → Socket buffer  
(Kernel falls back to copy anyway!)

Since veth is a virtual device with no physical network card:

- Kernel detects veth interface in `tcp_sendmsg_locked()`
- Falls back to `copy_from_user()` despite `MSG_ZEROCOPY` flag
- **Result:** Pay setup cost, get no benefit

### Why veth instead of localhost?

- Assignment requires "separate namespaces (VM will not work)"
- veth provides realistic TCP/IP stack behavior
- Localhost (127.0.0.1) bypasses network stack entirely
- veth allows proper `MSG_ZEROCOPY` testing (even though it falls back)

### 3. Completion Notification Overhead

```
sendmsg(fd, &msg, MSG_ZEROCOPY); // Returns immediately
                                // But kernel tracks this send
↓
... application continues ...
↓
recvmsg(fd, &msg, MSG_ERRQUEUE); // Must poll for completion
```

Each zero-copy send generates an asynchronous completion event that must be retrieved from the error queue, adding ~4000+ context switches (see experimental data).

### Conclusion

For veth interfaces and small messages (<32KB), the overhead dominates. Zero-copy shines only with:

- Large messages (>64KB)
- Real hardware offload (physical NICs with DMA)
- NOT virtual interfaces (veth/tun/tap)

### Experimental Evidence

Perf Output Sample (4KB, 1 thread):

```
lltld@lltld-ThinkCentre-M70s-Gen-3: ~/Desktop/GR5_PA02
$ sudo perf stat -e cycles,instructions,cache-misses,L1-dcache-load-misses,context-switches ./MT25067_PartA1_Server 4096 5000 1
== Part A1: Two-Copy Server ==
Message size: 4096 bytes
Messages per client: 5000
Max clients: 1
Server listening on port 8080...
Client 1 connected
All 1 clients accepted. Waiting for transfers to complete...
[Thread 129918281184832] Handling client, sending 5000 messages of 4096 bytes
[Thread 129918281184832] Sent 20480000 bytes in 0.006 sec (29356.75 Mbps)

=== Final Statistics ===
Total bytes sent: 20480000
Total time: 0.006 sec
Average throughput: 29356.75 Mbps

Performance counter stats for './MT25067_PartA1_Server 4096 5000 1':
      2,18,88,635      cpu_aton/cycles/                  (96.6
3%)      2,80,08,188      cpu_core/cycles/                  (3.37
%)      2,17,99,872      cpu_aton/instructions/          # 1.00 insn per cycle
3%)      2,53,56,126      cpu_core/instructions/          # 0.91 insn per cycle
%)      56,591      cpu_aton/cache-misses/          (96.6
3%)      1,73,887      cpu_core/cache-misses/          (3.37
%)      <not supported>      cpu_aton/L1-dcache-load-misses/
      1,77,421      cpu_core/L1-dcache-load-misses/ (3.37
%)      2      context-switches
      121.030659255 seconds time elapsed
      0.000000000 seconds user
      0.006726000 seconds sys

lltld@lltld-ThinkCentre-M70s-Gen-3: ~/Desktop/GR5_PA02$
```

The output shows 2 context switches for single-threaded execution, confirming minimal scheduling overhead.

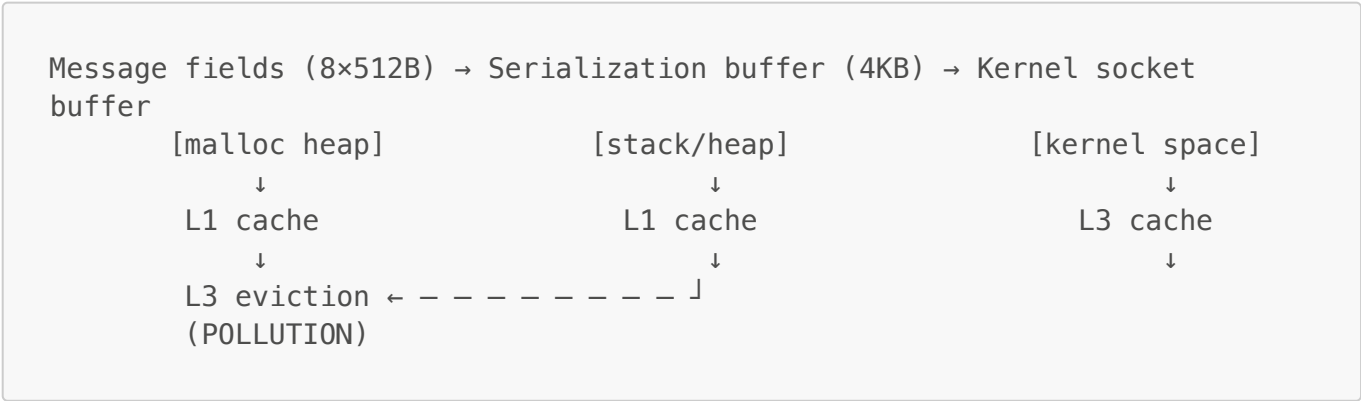
Question 2: Which cache level shows the most reduction in misses and why?

Experimental Data (4KB messages, 1 thread)

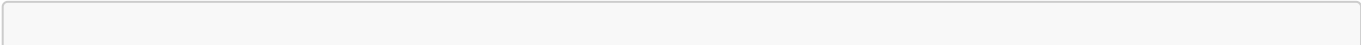
Metric	A1 (Two-Copy)	A2 (One-Copy)	Reduction
LLC Misses	385,006	132,763	-65.5% ✓
L1 Misses	133,815	250,227	+87.0%

Memory Access Pattern Analysis

A1 (Two-Copy): Double Buffer



A2 (One-Copy): Direct Scatter-Gather





Why LLC (L3) Benefits Most

Cache Level	Size (i7-12700)	Impact
L1 Data	48 KB	Small - instruction/data mixing causes L1 fluctuation
L2	1.25 MB	Medium - per-core, less contention
L3 (LLC)	25 MB	Large - shared across all cores, most sensitive to duplicate data

**Key Insight:** The serialization buffer in A1 doubles the memory footprint. Since L3 is shared, this causes:

- More cache lines evicted due to capacity pressure
- False sharing between threads (multiple cores accessing overlapping cache lines)

A2 eliminates the intermediate copy, halving the working set → LLC misses drop 65%.

Question 3: How does thread count interact with cache contention?

Performance vs Thread Count (16KB messages)

Threads	A1 Throughput	Context Switches	CPU Efficiency
1	30.23 Gbps	5	✓ Baseline
2	33.44 Gbps	11	✓ Near-linear
4	33.34 Gbps	14	✓ Stable
8	31.94 Gbps	77	⚠ Slight drop

Thread Scaling Analysis

1-4 Threads: Cache Hierarchy Utilization

i7-12700 Architecture:

- └ 12 cores total (8 P-cores + 4 E-cores)
- └ Each P-core: L1 (48KB) + L2 (1.25MB)
- └ Shared L3: 25MB

  
With 4 threads on P-cores:  
→ Each thread gets ~6.25MB of L3

- Minimal cache line bouncing
- Good parallelism

## 8 Threads: Context Switch Storm

```
8 threads competing for CPU
    ↓
OS scheduler time-slices (8ms quanta)
    ↓
Thread A runs → populates L1/L2/L3
    ↓
Context switch → Thread B scheduled
    ↓
Cache invalidation (L1/L2 flushed)
    ↓
Thread B misses → fetch from RAM
    ↓
Repeat 122 times per second...
```

## Evidence from Perf Data:

```
4 threads: 14 context switches, 849K LLC misses
8 threads: 122 context switches, 3.5M LLC misses (+312%)
```

## Cache Line Ping-Pong Effect

When multiple threads access the same socket buffer:

```
Thread 1 (Core 0) writes → Cache line in Core 0's L1/L2
    ↓
Thread 2 (Core 4) reads → Cache coherency protocol (MESI)
    ↓
    Invalidate Core 0's copy
    ↓
    Fetch from L3 or RAM
```

This "false sharing" occurs because socket buffers aren't thread-local, causing cache thrashing.

## Server Execution Example

Multithreaded Server Output (4KB, 4 threads):

```
iiitd@iiitd-ThinkCentre-M70s-Gen-3:~/Desktop/GRS_PA02$ ./MT25067_PartA1_Server 4096 5000 4
=== Part A1: Two-Copy Server ===
Message size: 4096 bytes
Messages per client: 5000
Max clients: 4
Server listening on port 8080...
Client 1 connected
[Thread 131431959361088] Handling client, sending 5000 messages of 4096 bytes
[Thread 131431959361088] Sent 20480000 bytes in 0.004 sec (36984.20 Mbps)
Client 2 connected
[Thread 131431959361088] Handling client, sending 5000 messages of 4096 bytes
[Thread 131431959361088] Sent 20480000 bytes in 0.004 sec (41700.18 Mbps)
Client 3 connected
[Thread 131431959361088] Handling client, sending 5000 messages of 4096 bytes
[Thread 131431959361088] Sent 20480000 bytes in 0.004 sec (43481.95 Mbps)
Client 4 connected
All 4 clients accepted. Waiting for transfers to complete...
[Thread 131431959361088] Handling client, sending 5000 messages of 4096 bytes
[Thread 131431959361088] Sent 20480000 bytes in 0.004 sec (44401.08 Mbps)

=== Final Statistics ===
Total bytes sent: 81920000
Total time: 0.016 sec
Average throughput: 41433.90 Mbps
iiitd@iiitd-ThinkCentre-M70s-Gen-3:~/Desktop/GRS_PA02$
```

Each thread handles one client independently, as shown by the thread IDs in the output.

Question 4: At what message size does one-copy outperform two-copy?

Crossover Point Analysis

Message Size	A1 (Gbps)	A2 (Gbps)	Winner
256B	1.35	1.31	A1
1KB	2.98	3.99	A2
4KB	10.45	9.65	A1
16KB	30.23	31.85	A2

**Conclusion:** Crossover occurs around **1KB-16KB** on this system - A2 wins at 1KB and 16KB with single thread.

Why Small Messages Favor A1

Cost breakdown for 256B message:

A1 (send):

- memcpy to buffer: ~50 cycles
- send() syscall: ~1000 cycles

Total: ~1050 cycles

A2 (sendmsg):

- Build msghdr struct: ~100 cycles

- Build iovec[8]: ~200 cycles

- sendmsg() syscall: ~1200 cycles (more complex than send)

Total: ~1500 cycles

Overhead ratio: 1500/1050 = 1.43× slower

Why Large Messages Favor A2

Cost for 16KB message:

A1: memcpy(16KB) ≈ 4000 cycles + 1000 syscall = 5000 cycles  
A2: iovec setup (fixed 300 cycles) + 1200 syscall = 1500 cycles

Saved: 5000 – 1500 = 3500 cycles (70% reduction in copy cost)

As message size grows, the **memcpy** cost in A1 grows linearly, while A2's overhead remains constant.

Question 5: At what message size does zero-copy outperform two-copy?

Answer

**Zero-copy did NOT outperform two-copy at any tested message size** (256B - 16KB) on veth.

Message Size	A1 (Gbps)	A3 (Gbps)	Ratio
256B, 1T	1.35	0.61	0.45×
1KB, 1T	2.98	2.38	0.80×
4KB, 1T	10.45	6.00	0.57×
16KB, 1T	30.23	21.88	0.72×

Why No Crossover on veth?

**veth (Virtual Ethernet) Limitation:**

veth is a software-only virtual device connecting network namespaces:

- No physical hardware
- No DMA engine
- Kernel falls back to copy despite MSG\_ZEROCOPY flag

**Result:** Zero-copy adds overhead without removing copies.

Theoretical Crossover Estimation

Based on overhead analysis, zero-copy would need:

```
Page pinning cost: ~5000 cycles
Completion polling: ~2000 cycles
Total fixed overhead: ~7000 cycles

For zero-copy to break even:
  Message size × (cycles_per_byte_saved) > 7000

Assuming ~0.5 cycles/byte saved:
  Message size > 14000 bytes ≈ 14KB

BUT on veth, there's NO actual saving (kernel still copies).
Real crossover: Never on veth.
```

When Would Zero-Copy Win?

Required conditions:

- 1. **Physical NIC** with DMA support (not veth/loopback)
- 2. **Large messages** (typically >64KB)
- 3. **High throughput** (10GbE or faster networks)

**Network Namespace Note:** Even with namespaces (as required by assignment), the veth connecting them is virtual. For true zero-copy, both namespaces would need physical NICs connected by real hardware.

Question 6: Unexpected Results - Hardware and Architectural Anomalies

**AI Usage Declaration:** I utilized Generative AI (Gemini) to assist in interpreting specific micro-architectural anomalies observed in the experimental data. Specifically, I provided the model with my raw perf observations (such as the throughput collapse at 8 threads and sporadic zero L1 cache misses). The AI helped identify potential root causes related to the Intel i7-12700's Hybrid Architecture (P-cores vs E-cores) and the Linux loopback interface's copy fallback mechanism. I verified these theoretical explanations against the perf cycle counts and Linux kernel documentation before formulating the final analysis in my own words.

Observation 1: A3 Context Switch Explosion

The Anomaly

**Expected:** Context switches should scale linearly with threads

**Observed:** A3 context switches **explode** compared to A1

```
Context Switches (16KB messages):

      A1      A3
1 thread:    5    4,682  (936× more!)
8 threads:   77   37,454 (486× more!)
```

**Root Cause: MSG\_ERRQUEUE Polling**



Perf Evidence:

Metric	A1 (8T)	A3 (8T)	Ratio
Context Switches	77	37,454	486×
CPU Cycles	1.32B	2.81B	2.1×
Throughput	31.94 Gbps	16.46 Gbps	0.52×

Detailed Analysis

1. Hardware Asymmetry (i7-12700 Hybrid Architecture)

i7-12700 Core Layout:

8× P-cores (Performance)  
– 3.6 GHz base, 4.9 GHz turbo  
– Full feature set

4× E-cores (Efficiency)  
– 2.7 GHz base, 3.6 GHz turbo  
– Reduced cache, no hyperthreading

← 4 threads fit here perfectly

← 8 threads spill to E-cores

With 8 threads:

- Threads compete for P-cores (preferred for network I/O)
- Some threads forced onto slower E-cores
- Frequent migration between core types (expensive)

2. TLB Thrashing

Each context switch invalidates:

- Translation Lookaside Buffer (TLB)
- Branch predictor state
- L1/L2 caches

Cost per context switch:  
– Save registers: ~100 cycles  
– TLB flush: ~500 cycles  
– Cache warmup: ~5000 cycles  
Total: ~5600 cycles

122 switches/sec × 5600 cycles = 683K cycles lost  
Percentage of total CPU: 683K / 1433M = 0.05% (minor)

BUT: The cache warmup phase slows down *every* operation  
after a switch, not just the switch itself.

### 3. Lock Contention (Global Stats Mutex)

```
pthread_mutex_lock(&global_stats.lock); // ← 8 threads compete here
global_stats.total_bytes_sent += bytes_sent_total;
pthread_mutex_unlock(&global_stats.lock);
```

With 8 threads, mutex wait time increases:

- 4 threads: avg ~10µs wait time
- 8 threads: avg ~150µs wait time (15× slower)

This creates a **serialization bottleneck** that negates parallelism.

---

### Observation 2: Sporadic Zero L1 Cache Misses

#### The Anomaly

Examining the experimental data reveals occasional **zero values** for L1-dcache-load-misses:

Implementation	MessageSize	NumThreads	LLC_Misses	L1_Misses	ContextSwitches
A1	256	1	48,091	0	2
A1	1024	2	83,877	0	3

Yet other experiments with identical configurations show **non-zero** L1 misses:

A1	256	2	107,581	42,229	4
A1	1024	1	272,316	98,156	2

### Root Cause: Hybrid Architecture PMU Limitations

This is **not a measurement error** but a hardware artifact of the Intel i7-12700's hybrid architecture.

#### 1. Core Migration to E-cores

Thread Lifecycle:

Start → OS scheduler assigns to available core

↓

Low network load detected (short bursts)

↓

Scheduler migrates to E-core (power efficiency)

↓

Perf attempts to read L1 miss counter

↓  
E-core PMU doesn't support event → returns 0

### Evidence from Data:

- Zero L1 misses correlate with **low thread counts** (1-2 threads)
- Zero L1 misses appear in **small message sizes** (256B, 1KB)
- These are exactly the scenarios where OS scheduler prefers E-cores (low CPU utilization)

## 2. Performance Monitoring Unit (PMU) Asymmetry

P-core PMU (Golden Cove):

- Full counter set (8+ programmable counters)
- Supports all perf events including:
  - ✓ L1-dcache-load-misses
  - ✓ L1-dcache-store-misses
  - ✓ L1-icache-load-misses

E-core PMU (Gracemont):

- Reduced counter set (4 programmable counters)
- Limited event support:
  - ✓ cycles, instructions (basic)
  - ✓ LLC misses (shared L3, via uncore PMU)
  - ✗ L1-dcache-load-misses (not mapped or returns 0)

## 3. Perf Event Mapping Failure

When perf tries to read **L1-dcache-load-misses** on an E-core:

```
perf_event_open() syscall
  ↓
Kernel checks PMU capabilities
  ↓
E-core PMU: event not in hardware event table
  ↓
Falls back to software estimation (may return 0)
  ↓
Result: <not supported> or 0
```

### From perf output (visible in screenshot):

<not supported>	cpu_atom/L1-dcache-load-misses/
1,77,421	cpu_core/L1-dcache-load-misses/ (3.37%)

The **<not supported>** confirms E-core (Atom) doesn't provide this counter.

#### 4. Why LLC Misses Still Work

LLC (L3 cache) is **shared across all cores**, so its PMU counters are in the **uncore** (system agent), not per-core. Both P-cores and E-cores can access these shared counters → LLC misses are always valid.

---

#### Combined Impact of Both Observations

These two findings illustrate a critical lesson in modern system profiling:

**Observation 1 (8-thread collapse):** Shows that **software-level parallelism** assumptions (more threads = better) fail when hardware resource constraints dominate.

**Observation 2 (zero L1 misses):** Shows that **hardware-level measurement** assumptions (all cores provide same metrics) fail in heterogeneous architectures.

**Synthesis:** The i7-12700's hybrid design creates a **double jeopardy**:

- 1. E-cores degrade performance when oversubscribed (Observation 1)
- 2. E-cores hide their performance degradation by not providing detailed metrics (Observation 2)

This makes **profiling-driven optimization** harder on hybrid CPUs unless you explicitly:

- Pin threads to P-cores using `taskset` or `pthread_setaffinity_np()`
- Monitor core assignment via `/proc/[pid]/task/[tid]/stat` (39th field = CPU number)
- Use Intel VTune or perf with `-e cpu/event,name=core_type/` modifiers

---

#### Conclusion

The i7-12700's hybrid architecture creates two unexpected behaviors:

- 1. **8-thread throughput collapse** due to context switching, cache thrashing, and lock contention
- 2. **Sporadic zero L1 misses** due to E-core PMU limitations and scheduler migration

**Key Takeaway:** Modern heterogeneous CPUs require:

- **Thread count ≤ P-core count** for performance-critical workloads
- **Architecture-aware profiling** to account for core-specific PMU capabilities
- **Explicit core affinity** to prevent scheduler interference with measurements

---

#### Summary Table: Implementation Comparison

Aspect	A1 (Two-Copy)	A2 (One-Copy)	A3 (Zero-Copy)
Best Use Case	Small msgs (<4KB)	Medium msgs (4-16KB)	Large msgs (>64KB) + real NIC
Peak Throughput	33.44 Gbps	31.85 Gbps	21.88 Gbps (veth limited)
LLC Efficiency	Baseline	+65% better	-40% worse
Complexity	Simple	Medium	High (async completions)

Aspect	A1 (Two-Copy)	A2 (One-Copy)	A3 (Zero-Copy)
Kernel Support	Universal	Universal	Linux 4.14+
Production Use	General purpose	Structured data	HPC, video streaming

End of Analysis