

DESIGN AND ANALYSIS OF CACHE ARCHITECTURE

COMPUTER ARCHITECTURE

AUTHOR:

Yusuf Taha ÖNCÜ

TRACE ONE: (16, 24, 128, 16, 24, 128, 16, 24, 128, 16, 24, 128, 16, 24, 128, 16, 24, 128)

Part I: Workload Characterization & Locality Hypothesis (Trace One)

1. Manual Analysis: The trace shows a repeating pattern of three memory addresses: 16, 24, and 128. The sequence is periodic with a period of 3.

2. Hypothesize Locality: This trace exhibits excellent **temporal locality**. Since the same three addresses are accessed repeatedly in a short loop, once these blocks are brought into the cache, they are likely to be referenced again very soon. **Spatial Locality:** There is **some spatial locality** between addresses 16 and 24. The difference is only 8 bytes. If the block size is large enough and the blocks are aligned, these two addresses might reside in the same cache block. Address 128 is far from the others, showing no spatial locality with the first two.

3. Hypothesize Workload Type: This pattern likely represents a tight **loop execution** or an iterative process accessing a small set of variables (e.g., `for(i=0; i<N; i++) { sum += A[i]; }` where instructions and data variables are constantly reused). It suggests a program working on a very small "working set" of data.

Part II: Experimental Simulation & Data Collection (Trace One)

***Methodology Note:** To conduct a controlled experiment, I adopted a "change one variable at a time" approach. I initially fixed the Replacement Policy to LRU (the standard benchmark for temporal locality) while varying hardware parameters. This ensures that any observed changes in hit rates are strictly due to the hardware architecture changes, not the replacement algorithm.*

STEP A: TESTING SPATIAL LOCALITY

- The size of cache and associativity will remain same. Block size will be tested if 16 and 24 tries to enter the same block.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 1	128 Bytes	Direct Mapped (1-way)	2 Bytes	LRU
Test 2	128 Bytes	Direct Mapped (1-way)	4 Bytes	LRU
Test 3	128 Bytes	Direct Mapped (1-way)	8 Bytes	LRU

Analysis Of Varying Block Size

Test ID	Cache Size	Associativity	Block Size	Policy	Total Accesses	Hits	Misses	Hit Rate
Test 1	128 B	Direct Mapped	2 B	LRU	18	15	3	83.3%
Test 2	128 B	Direct Mapped	4 B	LRU	18	15	3	83.3%
Test 3	128 B	Direct Mapped	8 B	LRU	18	15	3	83.3%

Observation: Across Tests 1, 2, and 3, varying the Block Size () from 2 bytes to 8 bytes resulted in identical hit/miss patterns. In all cases, the first iteration of the loop ($16 \rightarrow 24 \rightarrow 128$) generated misses, while all subsequent iterations resulted in hits.

Detailed Interpretation:

- **Compulsory Misses:** The three misses observed at the beginning of each simulation are **compulsory (cold) misses**. Since the cache is initially empty, the first access to addresses $0x10$ (16), $0x18$ (24), and $0x80$ (128) must fetch data from the main memory, regardless of the block size or associativity.
- **Temporal Locality Success:** The 100% hit rate in subsequent iterations proves that the workload exhibits strong **temporal locality**. The total working set of this trace consists of only 3 unique memory blocks. Since the cache capacity (128 Bytes) is significantly larger than the working set size (even with the largest block size of 8 bytes, the working set is only 24 bytes), the cache can easily hold all active blocks without eviction.
- **The "Spatial Locality" Limit:** A key objective of increasing Block Size is to exploit spatial locality loading neighbor addresses (like 16 and 24) in a single fetch. However, the data shows that **increasing the block size to 8 bytes failed to capture addresses 16 and 24 in the same block.**
 - **Mathematical Explanation:**
 - Address 16 ($0x10$): (Start of a block).
 - Address 24 ($0x18$): (Start of the *next* block). Because the distance between the two addresses is exactly 8 bytes, and they are aligned to 8-byte boundaries, they map to **consecutive but distinct cache blocks**. To exploit spatial locality and eliminate the miss for address 24, a block size of at least **16 bytes** would be required. Since the simulation was limited to a maximum , we observed no reduction in the miss rate.

STEP B: TESTING CONFLICT

- The size of cache and block will remain same. Associativity will be tested if there is any conflict between traces.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 4	128 Bytes	2-way	8 Bytes	LRU
Test 5	128 Bytes	4-way	8 Bytes	LRU

Analysis Of Varying Associativity

Test ID	Cache Size	Associativity	Block Size	Policy	Hits	Misses	Hit Rate
Test 4	128 B	2-Way	8 B	LRU	15	3	83.3%
Test 5	128 B	4-Way	8 B	LRU	15	3	83.3%

Observation: Increasing the associativity from Direct-Mapped (Test 3) to 2-Way (Test 4) and 4-Way (Test 5) yielded **no improvement** in the hit rate. The system consistently incurred 3 compulsory misses followed by a stream of hits.

Detailed Interpretation:

- **Absence of Conflict Misses:** The primary purpose of increasing associativity is to reduce **conflict misses**, which occur when multiple active blocks map to the same set index.
 - In this specific trace, the active blocks correspond to addresses 16 ($0x10$), 24 ($0x18$), and 128 ($0x80$).
 - With and , the set indices are calculated as follows (for Direct Mapped):
 - Address 16 Set 2
 - Address 24 Set 3
 - Address 128 Set 0
 - Since each memory block naturally maps to a **unique set**, there were zero conflicts even in the simplest Direct-Mapped configuration.
- **Diminishing Returns & Cost Analysis:** Moving to 2-Way or 4-Way associativity provided no performance benefit because the underlying problem was not conflict-related. However, increasing associativity introduces **hardware overhead**.
 - **Complexity:** Requires more comparators and complex multiplexers.
 - **Energy/Time:** Associative searches consume more power and can slightly increase hit time compared to Direct-Mapped caches.
 - **Conclusion for Optimization:** For this specific workload, a **Direct-Mapped cache is superior** to 2-Way or 4-Way designs because it delivers identical performance at a lower hardware cost.

STEP C: TESTING CAPACITY SIZE OF CACHE

- The trace we have is small compared to real life situations. So we may not see any difference at the efficiency of hit and miss rate.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 6	64 Bytes	2-way	8 Bytes	LRU
Test 7	256 Bytes	2-way	8 Bytes	LRU

Analysis Of Varying Size Of Cache

Test ID	Cache Size	Associativity	Block Size	Policy	Hits	Misses	Hit Rate
Test 6	64 B	2-Way	8 B	LRU	15	3	83.3%
Test 7	256 B	2-Way	8 B	LRU	15	3	83.3%

Observation: Varying the cache size from 64 Bytes to 256 Bytes resulted in **identical performance metrics**. The hit/miss counts remained constant across all tested capacities.

Detailed Interpretation:

- **Working Set Analysis:** The results confirm that the "working set" of this specific workload is extremely small. The trace repeatedly accesses only 3 distinct memory

blocks (addresses 16, 24, and 128). With a block size of 8 bytes, the total memory footprint required to store the active data is only:

- **Capacity Threshold:**
 - **Test 6 (64 Bytes):** Even the smallest cache tested (64 Bytes) is nearly **3x larger** than the required footprint (24 Bytes). Therefore, the cache never suffers from "Capacity Misses."
 - **Test 7 (256 Bytes):** Increasing the size to 256 Bytes provides **zero marginal benefit**. It simply results in unused cache lines ("dead space").
- **Conclusion:** For this workload, a 64-Byte cache is already sufficient. Using a 256-Byte cache would be an inefficient design choice, increasing area and leakage power costs without improving performance.

STEP D: EFFECT OF REPLACEMENT POLICY

- The replacement policy might be crucial for our efficiency. Hardware parameters will remain same while varying replacement policy.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 8	128 Bytes	2-way	8 Bytes	FIFO
Test 9	128 Bytes	2-way	8 Bytes	Random

Analysis Of Varying Replacement Method

Test ID	Cache Size	Assoc	Block Size	Policy	Hits	Misses	Hit Rate
Test 4 (Ref)	128 B	2-Way	8 B	LRU	15	3	83.3%
Test 8	128 B	2-Way	8 B	FIFO	15	3	83.3%
Test 9	128 B	2-Way	8 B	Random	15	3	83.3%

Observation: Changing the replacement policy from LRU to FIFO or Random had **zero impact** on the hit rate. All policies produced identical results.

Detailed Interpretation:

- **The "No Eviction" Scenario:** Replacement policies dictate which block to evict when a set is full and a new block needs to be brought in. However, in this experiment, the cache sets never reached their full capacity.
- **Set Occupancy:**
 - With Associativity = 2 and Block Size = 8, the cache has multiple sets.
 - Addresses 16, 24, and 128 map to different sets or do not conflict enough to fill any single set beyond its capacity of 2 ways.
- **Conclusion:** Since no blocks ever needed to be evicted, the replacement algorithm was never triggered. The simulator simply filled the empty slots (cold misses) and then hit on subsequent accesses.

Note on Trace Length and Hit Rate: The observed hit rate of 83.3% is heavily influenced by the short length of the trace). Since the misses are purely "compulsory," extending the simulation duration would cause the hit rate to asymptotically approach 100%. The 16.7% miss rate represents the "warm-up" cost, not a flaw in the cache design.

Part III: Optimal Design Recommendation & Justification (Trace One)

Recommendation: Based on the experimental data collected in Part II, the **single best cache configuration** for this specific workload is:

- **Cache Size:** 64 Bytes
- **Associativity:** Direct-Mapped (1-Way)
- **Block Size:** 8 Bytes
- **Replacement Policy:** LRU

Justification: The selection of this configuration is driven by a **cost-benefit analysis**. Since the simulation results showed identical hit rates (83.3%) across all tested configurations (varying associativity, cache size, and block size), the optimal design is the one that minimizes hardware complexity, power consumption, and cost without sacrificing performance.

- **Why this Configuration is Optimal:** The workload exhibits a very small working set (only 3 unique memory blocks) and a strictly repeating pattern. A **64-Byte** cache is sufficient to hold the entire working set, making larger caches (e.g., 256 Bytes) wasteful in terms of silicon area and static power leakage. Furthermore, since there are no conflict misses, increasing associativity provides no benefit. A **Direct-Mapped** cache is chosen because it has the fastest hit time (no multiplexer delay) and the lowest logic complexity compared to set-associative designs.
- **Block Size and Spatial Locality:** Although the trace exhibits spatial locality between addresses 16 and 24, the experimental data showed that a block size of **8 Bytes** was insufficient to capture both in a single block due to alignment boundaries (Address 16 and 24 mapped to consecutive but distinct blocks). However, I recommend maintaining the **8-Byte** block size over smaller sizes (2 or 4 Bytes). Larger blocks reduce the total number of cache lines for a given capacity, thereby significantly reducing the **tag overhead** (memory required to store address tags), making the design more area-efficient.
- **Associativity and Conflict Mitigation:** The data revealed **zero conflict misses**. The addresses 16, 24, and 128 naturally mapped to different sets in a Direct-Mapped configuration. Therefore, increasing associativity to 2-way or 4-way offered **diminishing returns**. It would introduce unnecessary hardware overhead (comparators and complex replacement logic) without improving the hit rate.
- **Hypothesis Confirmation & Surprises:** The data confirmed the initial hypothesis of **strong temporal locality**, as evidenced by the hits in repeated loop iterations. The surprise was the **ineffectiveness of spatial locality optimizations**; the specific memory alignment of the variables prevented the larger block size from reducing compulsory misses.
- **Trade-offs (Real-World Factors):**

- **Why not a larger cache?** A larger cache (e.g., 256B or 512B) would increase access latency and energy consumption per access. Since the hit rate is already maximized for the steady state, extra capacity is redundant resources ("dead space").
- **Why not higher associativity?** Higher associativity increases the **hit time** (critical path delay). Since we achieve the same hit rate with a Direct-Mapped cache, the simpler design is preferred for higher clock frequency potential.

Part IV: The AI Sanity Check (Trace One)

Prompt: Analyze this memory access trace (16, 24, 128, 16, 24, 128...) and recommend an optimal cache configuration for a system with a 1024-byte memory space.

AI Response: "Based on the trace 16, 24, 128 repeating, the workload shows strong **temporal locality** (the same addresses are accessed repeatedly) and potential **spatial locality** between addresses 16 and 24 (which are close to each other).

Recommendation:

- **Cache Size:** 128 Bytes. This is small enough to be fast but large enough to hold the working set of 3 variables comfortably.
- **Block Size:** 16 Bytes. Addresses 16 ($0x10$) and 24 ($0x18$) are 8 bytes apart. A larger block size like 16 bytes is recommended to try and capture both variables in a single cache line, reducing compulsory misses.
- **Associativity:** 2-Way Set Associative. This is a safe choice to prevent potential conflict misses between the addresses, especially since 128 is a power-of-two multiple which often causes mapping conflicts in direct-mapped caches.
- **Replacement Policy:** LRU (Least Recently Used) is the standard policy to handle loops efficiently."

Critical Analysis (Human vs. AI)

Critique: The AI provided a theoretically sound "textbook" answer, but my experimental simulation revealed specific engineering constraints that the AI missed.

Where the AI got it right: The AI correctly identified the **Workload Type** (a loop with strong temporal locality) and recognized the potential for spatial locality between addresses 16 and 24. Its recommendation for a small cache size was also accurate.

Where the AI fell short (The "Human" Insight):

1. **Block Size & Alignment Failure:** The AI recommended a 16-Byte block to capture spatial locality. However, it failed to calculate the specific **address alignment**. In my simulation (Part II), I discovered that even with an 8-Byte block, addresses 16 and 24 mapped to *different* blocks due to the boundary alignment (and). The AI assumed a

generic benefit of large blocks without verifying the specific mapping arithmetic. My analysis showed that smaller blocks (8 Bytes) are equally effective and more cost-efficient for this specific alignment.

2. **Unnecessary Associativity (Over-Engineering):** The AI recommended a **2-Way Set Associative** cache as a "safe choice" to avoid conflicts. My specific conflict analysis (Tests 4-5) proved this to be **over-engineering**. By tracing the set indices manually, I proved that addresses 16, 24, and 128 *never* conflict in a Direct-Mapped cache. The AI suggested a more expensive hardware solution (2-Way) for a problem that didn't exist. My recommendation of **Direct-Mapped** saves hardware cost and access time.
3. **Lack of Cost Awareness:** The AI focused on maximizing theoretical hit rates. My analysis focused on the **Cost/Performance trade-off**. I chose the configuration (Direct-Mapped, 64B Cache) that delivered the *same* hit rate as the AI's complex suggestion but with significantly fewer transistors and lower power consumption.

TRACE TWO: (128, 132, 136, 140, 144, 148, 152, 156, 160, 164, 168, 172, 176, 180, 184, 188)

Part I: Workload Characterization & Locality Hypothesis (Trace Two)

1. Manual Analysis: The trace consists of a strictly increasing sequence of memory addresses starting at 128 and incrementing by exactly 4 bytes at each step (.). The sequence is linear and non-repetitive; no address is accessed more than once within the observed window.

2. Hypothesize Locality:

- **Temporal Locality:** This trace exhibits **extremely poor (or non-existent) temporal locality**. Since the program accesses each address exactly once and immediately moves to the next, data brought into the cache is never referenced again. A standard caching mechanism relying solely on reuse will struggle here.
- **Spatial Locality:** The trace exhibits **excellent spatial locality**. The stride (distance between accesses) is consistently 4 bytes, which corresponds to the size of a standard 32-bit integer. This implies that the requested data resides in contiguous memory locations. If a cache block is large enough to hold multiple 4-byte words, a single fetch could successfully pre-load subsequent requests.

3. Hypothesize Workload Type: This pattern is characteristic of a **Linear Array Traversal** or a "Streaming" workload. It likely represents a loop iterating sequentially through an array, vector, or a contiguous data structure (e.g., `for(int i=0; i<N; i++) { sum += Array[i]; }`). The CPU is reading data item-by-item from a continuous block of memory.

Part II: Experimental Simulation & Data Collection (Trace Two)

STEP A: TESTING SPATIAL LOCALITY

- This trace does not loop so our expectations to place the elements side by side. We will try to achieve this expectation with varying block size

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 1	128 Bytes	Direct Mapped (1-way)	2 Bytes	LRU
Test 2	128 Bytes	Direct Mapped (1-way)	4 Bytes	LRU
Test 3	128 Bytes	Direct Mapped (1-way)	8 Bytes	LRU

Analysis Of Varying Block Size

Test ID	Cache Size	Associativity	Block Size	Policy	Total Accesses	Hits	Misses	Hit Rate
Test 1	128 B	Direct Mapped	2 B	LRU	16	0	16	0%
Test 2	128 B	Direct Mapped	4 B	LRU	16	0	16	0%
Test 3	128 B	Direct Mapped	8 B	LRU	16	8	8	50%

Observation: In Tests 1 and 2 (Block Sizes of 2 and 4 Bytes), the system failed to achieve a single cache hit, resulting in a **0% hit rate**. However, increasing the Block Size to 8 Bytes (Test 3) caused a dramatic shift, resulting in a **50% hit rate** with a perfectly alternating "Miss-Hit-Miss-Hit" pattern.

Detailed Interpretation:

- **Failure of Small Blocks (The "Stride" Problem):** The trace accesses memory with a stride of **4 Bytes** (e.g.,).
 - **Test 1 (B=2):** The block is smaller than the data object (assuming 4-byte integers) or the stride. Every access maps to a new block.
 - **Test 2 (B=4):** Each cache block holds exactly one 4-byte word.
 - Access to **0x80** (Miss) Loads Block holding **0x80-0x83**.
 - Next Access **0x84** This address is outside the previous block. It requires a new fetch (Miss). Because the stride equals the block size, there is no room to "prefetch" the neighbor.
- **Success of Larger Blocks (Exploiting Spatial Locality):** In Test 3, the Block Size is **8 Bytes**. This allows two consecutive 4-byte words to fit into a single cache line.
 - **Access 1 (0x80): MISS.** The cache fetches the block containing bytes **0x80** to **0x87**.
 - **Access 2 (0x84): HIT.** Since **0x84** falls within the same loaded block range (**0x80-0x87**), it is found in the cache.
 - **Access 3 (0x88): MISS.** Starts a new block (**0x88-0x8F**).
 - **Access 4 (0x8C): HIT.** Found in the current block. This pattern confirms that **Spatial Locality** is the primary performance driver for this sequential workload.

STEP B: TESTING CONFLICT

- We already tested block size, most efficient one is 8, so block size will remain 8. Only associativity will vary. Also we already tried direct associativity, two test will be enough.

Deney No	Cache Size	Associativity	Block Size	Replacement
Test 4	128 Bytes	2-way	8 Bytes	LRU
Test 5	128 Bytes	4-way	8 Bytes	LRU

Analysis Of Varying Associativity

Test ID	Cache Size	Associativity	Block Size	Policy	Hits	Misses	Hit Rate
Test 4	128 B	2-Way	8 B	LRU	8	8	50%
Test 5	128 B	4-Way	8 B	LRU	8	8	50%

Observation: Increasing the associativity from Direct-Mapped (Test 3) to 2-Way (Test 4) and 4-Way (Test 5) yielded **no change** in performance. The hit rate remained constant at 50%.

Detailed Interpretation:

- **Nature of Conflict Misses:** Conflict misses occur when multiple addresses compete for the same cache set and evict useful data that is needed later.
- **Why Associativity Failed:** This workload is strictly sequential (). It never looks back.
 - Even if an old block is evicted due to a set mapping conflict (in Direct-Mapped), that eviction is harmless because the program **never returns** to access that old address again.
 - Therefore, the "safety net" provided by higher associativity (keeping old blocks longer) provides no value for a streaming workload that has **zero temporal locality**.

STEP C: TESTING CAPACITY SIZE OF CACHE

- Even though we reduce the size of cache our data is streaming so our prediction is nothing will change the efficiency of the hit rate

Deney No	Cache Size	Associativity	Block Size	Replacement
Test 6	64 Bytes	2-way	8 Bytes	LRU
Test 7	256 Bytes	2-way	8 Bytes	LRU

Analysis Of Varying Capacity Size Of Cache

Test ID	Cache Size	Associativity	Block Size	Policy	Hits	Misses	Hit Rate
Test 6	64 B	2-Way	8 B	LRU	8	8	50%
Test 7	256 B	2-Way	8 B	LRU	8	8	50%

Observation: Varying the cache capacity had **zero impact**. A tiny 64-Byte cache performed exactly the same as a larger 256-Byte cache.

Detailed Interpretation:

- **Capacity Irrelevance:** Capacity misses occur when the working set is too large to fit in the cache. However, for a streaming workload without loops, the "active working set" is effectively just the **current cache block**.
- **One-Block Requirement:** As long as the cache is large enough to hold a single block (8 Bytes), it can capture the spatial locality (the neighbor element). Storing previous blocks in a large 256-byte cache is useless because the CPU never requests them again. This confirms the hypothesis of **poor temporal locality**.

STEP D: EFFECT OF REPLACEMENT POLICY

- The replacement policy might be crucial for our efficiency. Hardware parameters will remain the same while varying the replacement policy.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 8	128 Bytes	2-way	8 Bytes	FIFO
Test 9	128 Bytes	2-way	8 Bytes	Random

Analysis Of Varying Replacement Method

Test ID	Cache Size	Assoc	Block Size	Policy	Hits	Misses	Hit Rate
Test 4 (Ref)	128 B	2-Way	8 B	LRU	8	8	50%
Test 8	128 B	2-Way	8 B	FIFO	8	8	50%
Test 9	128 B	2-Way	8 B	Random	8	8	50%

Observation: Changing the replacement policy from LRU to FIFO or Random resulted in **identical hit rates (50%)**.

Detailed Interpretation:

- **Irrelevance of Eviction Strategy:** Replacement policies determine which data to discard to make room for new data. For policies like LRU or FIFO to be effective, there must be a possibility that "old" data will be reused.

- **Sequential Access Characteristic:** Since this trace is purely sequential () and never loops back to previous addresses, **evicted data is never requested again**.
- **Conclusion:** Whether the cache evicts the oldest block (FIFO), the least recently used block (LRU), or a random block is mathematically irrelevant. The penalty for eviction is zero because the evicted block has zero future utility.

Part III: Optimal Design Recommendation & Justification (Trace Two)

Recommendation: Based on the experimental data, the **single best cache configuration** for this sequential workload is:

- **Cache Size:** 64 Bytes (Smallest tested)
- **Associativity:** Direct-Mapped (1-Way)
- **Block Size:** 8 Bytes (Or larger if allowed)
- **Replacement Policy:** Random (Cheapest to implement)

Justification: The optimization strategy for this workload is purely "**Cost Minimization**", as performance is strictly limited by the sequential nature of the access pattern.

- **Block Size is the Only Driver:** The only parameter that improved performance from 0% to 50% was the Block Size. By using an **8-Byte block**, we exploited the spatial locality of the consecutive integers. A smaller block size (2 or 4 Bytes) resulted in total failure.
- **Minimal Hardware Complexity:** Since increasing Associativity (to 2-way or 4-way) and Cache Size (to 256 Bytes) provided **zero performance gain**, the simplest possible hardware is the optimal choice. A **Direct-Mapped, 64-Byte cache** consumes the least power and silicon area while delivering the maximum achievable hit rate (50%) for this specific block size.
- **Policy Selection:** Since the replacement policy has no impact on hit rate, **Random** replacement is recommended over LRU. Random logic requires fewer gates (less hardware overhead) than maintaining LRU bits/counters, further reducing the cost.

Part IV: The AI Sanity Check (Trace Two)

Prompt: "Analyze the memory trace 128, 132, 136, 140... and recommend a cache config."

AI Response: "This trace represents a sequential scan of an array. It exhibits strong **spatial locality** but no temporal locality."

Recommendation:

- **Block Size:** Maximize the block size (e.g., 32 or 64 bytes). This will load more upcoming words in a single fetch, drastically reducing compulsory misses.
- **Prefetching:** A hardware prefetcher is highly recommended for this pattern.
- **Cache Size:** Keep it small; large caches won't help since data isn't reused.

- **Associativity:** Direct-Mapped is sufficient as there are no conflict misses in a linear scan."

Critical Analysis (Human vs. AI)

Critique: The AI correctly identified the workload type, but its recommendations drifted into "theoretical" solutions that were outside the scope of my constrained simulation.

Where the AI got it right: The AI correctly identified that **Spatial Locality** is the key and that **Temporal Locality** is nonexistent. It also correctly advised a small cache size and Direct-Mapped associativity, aligning with my cost-benefit analysis.

Where the AI fell short (The "Constraint" Insight):

- **The "Prefetching" Cop-out:** The AI suggested "Hardware Prefetching," which is a valid real-world solution but wasn't a configurable parameter in our experiment. It used this as a "magic bullet" answer.
- **Block Size Generalization:** The AI suggested "Maximize Block Size." While theoretically true, my analysis in **Part II (Step A)** showed the specific sensitivity of this parameter. I proved *experimentally* that 4 Bytes was useless and 8 Bytes was the tipping point. The AI gave a generic rule of thumb; I found the specific engineering threshold.

TRACE THREE: (20, 480, 312, 800, 128, 940, 516, 624, 2, 712, 340, 998)

Part I: Workload Characterization & Locality Hypothesis (Trace Three)

1. Manual Analysis: The trace consists of a seemingly random sequence of memory addresses ranging from 2 to 998 (e.g.,). There is no discernable pattern, stride, or repetition. The addresses jump across large gaps in the memory space.

2. Hypothesize Locality:

- **Temporal Locality:** This trace exhibits **zero temporal locality**. No address is accessed more than once. Once a block is loaded, it is essentially "dead weight" as the CPU never requests it again.
- **Spatial Locality:** The trace exhibits **extremely poor spatial locality**. The distance between consecutive accesses is erratic and large (e.g., jumping from 20 to 480).
 - *Potential Exception:* There is a slight proximity between address 20 (`Trace[0]`) and address 2 (`Trace[8]`), and perhaps 312 and 340. However, they are accessed far apart in time, making it unlikely for a standard cache block to capture them together unless the block size is very large and the block is not evicted.

3. Hypothesize Workload Type: This pattern is characteristic of a **Random Access** workload, such as a Hash Table lookup, a Pointer Chasing algorithm (Traversing a Linked List where nodes

are scattered in heap memory), or a sparse matrix operation. It represents a "Cache Hostile" scenario.

Part II: Experimental Simulation & Data Collection (Trace Three)

STEP A: TESTING SPATIAL LOCALITY

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 1	128 Bytes	Direct Mapped (1-way)	2 Bytes	LRU
Test 2	128 Bytes	Direct Mapped (1-way)	4 Bytes	LRU
Test 3	128 Bytes	Direct Mapped (1-way)	8 Bytes	LRU

Analysis Of Varying Block Size

Test ID	Cache Size	Associativity	Block Size	Policy	Total Accesses	Hits	Misses	Hit Rate
Test 1	128 B	Direct Mapped	2 B	LRU	12	0	12	0%
Test 2	128 B	Direct Mapped	4 B	LRU	12	0	12	0%
Test 3	128 B	Direct Mapped	8 B	LRU	12	0	12	0%

Observation: The simulation resulted in a **0% hit rate** across all tested block sizes. Every single memory access resulted in a miss.

Detailed Interpretation:

- **Total Absence of Locality:** This trace represents a "Worst-Case Scenario" for caching.
 - **No Temporal Reuse:** The trace never accesses the same address twice. Therefore, once a block is loaded into the cache (compulsory miss), it occupies space but provides no future value.
 - **Spatial Discontinuity:** The memory accesses jump across large gaps (e.g., from 20 to 480, then to 312). Even with the largest tested **Block Size of 8 Bytes**, the "neighbors" of the requested data were never the targets of subsequent requests. For example, accessing address 20 loads bytes 16-23. The next request is 480, which is far outside this block. Address 2 (0x02) and 20 (0x14) are relatively close, but their distance (18 bytes) still exceeds the 8-byte block limit, causing them to map to different blocks.

STEP B: TESTING CONFLICT

- No reuse so there is no logical meaning to keep the data in the data 2-way or 4-way but still we give it a shot.

Deney No	Cache Size	Associativity	Block Size	Replacement
Test 4	128 Bytes	2-way	8 Bytes	LRU
Test 5	128 Bytes	4-way	8 Bytes	LRU

Analysis Of Varying Associativity

Test ID	Cache Size	Associativity	Block Size	Policy	Hits	Misses	Hit Rate
Test 4	128 B	2-Way	8 B	LRU	0	12	0%
Test 5	128 B	4-Way	8 B	LRU	0	12	0%

Observation: Increasing the associativity to 2-Way or 4-Way had absolutely **no effect**. The hit rate remained at 0%.

Detailed Interpretation:

- **No Conflicts to Solve:** Associativity is designed to reduce conflict misses (where two addresses map to the same index). However, for a conflict miss to occur, the CPU must eventually *return* to a previously evicted address.
- **Trace Characteristic:** Since this trace never accesses the same address twice (no temporal locality), every single miss is a "**Compulsory Miss**". No amount of associativity can prevent a compulsory miss; the data must be fetched from memory the first time it is seen. Therefore, complex associative hardware is wasted on this workload.

STEP C: TESTING CAPACITY SIZE OF CACHE

- We can build the cache even with 1TB does not matter, no reuse of the data reduce our ability to build efficient cache

Deney No	Cache Size	Associativity	Block Size	Replacement
Test 6	64 Bytes	2-way	8 Bytes	LRU
Test 7	256 Bytes	2-way	8 Bytes	LRU

Analysis Of Varying Size Of Cache

Test ID	Cache Size	Associativity	Block Size	Policy	Hits	Misses	Hit Rate
Test 6	64 B	2-Way	8 B	LRU	0	12	0%
Test 7	256 B	2-Way	8 B	LRU	0	12	0%

Observation: Varying the cache capacity yielded identical results (0% hit rate).

Detailed Interpretation:

- **Capacity is Not the Bottleneck:** A larger cache is useful only if the working set is reused. Here, the "active" data is never reused.
- **Dead Blocks:** Once a block is loaded into the cache, it sits there occupying space but is never referenced again. Whether the cache can hold 8 blocks (64B) or 32 blocks (256B) is irrelevant because the program simply marches forward to new, random addresses.

STEP D: EFFECT OF REPLACEMENT POLICY

- The replacement policy unlikely be crucial for our efficiency. Hardware parameters will remain the same while varying the replacement policy.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 8	128 Bytes	2-way	8 Bytes	FIFO
Test 9	128 Bytes	2-way	8 Bytes	Random

Analysis Of Varying Replacement Method

Test ID	Cache Size	Assoc	Block Size	Policy	Hits	Misses	Hit Rate
Test 4 (Ref)	128 B	2-Way	8 B	LRU	0	12	0%
Test 8	128 B	2-Way	8 B	FIFO	0	12	0%
Test 9	128 B	2-Way	8 B	Random	0	12	0%

Observation: Changing the replacement policy had absolutely no effect. The hit rate remained at 0% for all policies.

Detailed Interpretation:

- **No Re-Use, No Policy Impact:** Replacement policies decide what to keep in the cache for *future use*. Since this workload (Random Access) never requests the same data twice, it doesn't matter what we keep or what we throw away.
- **Cache Occupancy:**
 - Total unique blocks accessed: 12
 - Space required: .
 - Cache Capacity: 128 Bytes.
 - Since , the cache **never filled up**. Therefore, no block was ever evicted, and the replacement logic (LRU/FIFO/Random) was never even executed.

Part III: Optimal Design Recommendation & Justification (Trace Three)

Recommendation: Based on the experimental data, the **single best cache configuration** for this specific workload is:

- **Cache Size:** 64 Bytes (Smallest possible)
- **Associativity:** Direct-Mapped (1-Way)
- **Block Size:** 8 Bytes (Or smaller)
- **Replacement Policy:** Random

Justification: The optimization strategy for this workload is "**Damage Control**" (Minimizing Waste).

- **The "Zero Benefit" Reality:** The simulation proved that the cache provides **0% performance benefit** (0% Hit Rate) regardless of the configuration. The workload is purely random and sparse ("Cache Hostile").
- **Cost Minimization Strategy:** Since we cannot improve performance, we must minimize **cost**.
 - A **Direct-Mapped, 64-Byte** cache is the cheapest, smallest, and lowest-power option available.
 - Using a larger cache (e.g., 256 Bytes) or higher associativity (4-Way) would consume more silicon area and energy *for absolutely no return on investment*.
- **Theoretical Recommendation:** In a real-world scenario specifically for *only* this trace, the optimal design would be to **bypass the cache entirely** (Uncached Access) to avoid the latency overhead of checking tags for guaranteed misses. However, within the constraints of selecting a cache design, we choose the "cheapest placeholder."

Part IV: The AI Sanity Check (Trace Three)

Prompt: "Analyze the memory trace 20, 480, 312, 800... and recommend a cache config."

AI Response: "The trace appears to be a **sparse/random access** pattern with large jumps between addresses.

Recommendation:

- **Associativity:** Use **Fully Associative** or high associativity (e.g., 8-Way). Random accesses cause erratic conflict misses, and high associativity helps mitigate this.
- **Cache Size:** Increase the cache size significantly (e.g., 1KB or more) to try and capture the entire 'working set' since locality is poor.
- **Block Size:** Keep block size moderate (e.g., 32 bytes) to balance spatial locality potential without wasting bandwidth."

Critical Analysis (Human vs. AI)

Critique: The AI failed to recognize the futility of caching for this specific sparse dataset and suggested a "Brute Force" approach that leads to massive waste.

Where the AI got it right: It correctly identified the workload as "sparse/random."

Where the AI fell short (The "Cost" Insight):

- **The "Throw Money at the Problem" Fallacy:** The AI suggested increasing Cache Size and Associativity to fight the random misses. My simulation (Part II) proved this to be false. Even with 4-Way associativity, the hit rate was 0%. The AI's recommendation would result in a very expensive, high-power cache that still yields near-zero performance.
- **Misunderstanding Compulsory Misses:** The AI tried to solve "Conflict Misses" (via Associativity). However, my analysis showed that all misses were **Compulsory** (First-time accesses). No amount of associativity can fix compulsory misses. I correctly identified that the cheapest hardware is the best choice when performance is impossible to improve.

TRACE FOUR: (0,4,8,12,16,20, 256,260,264,268,272,276, 512,516,520,524,528,532, 768,772,776,780,784,788)

Part I: Workload Characterization & Locality Hypothesis (Trace Four)

1. Manual Analysis: The trace exhibits a "Chunked" or "Strided" access pattern. It accesses data in 4 distinct groups (chunks): 0–20, 256–276, 512–532, and 768–788.

- **Within each chunk:** The access is purely sequential with a stride of 4 bytes.
- **Between chunks:** There is a large, fixed stride of **256 bytes** between the starting addresses of each group () .

2. Hypothesize Locality:

- **Temporal Locality:** Poor. The trace touches a data range and moves on to the next "chunk" without returning. There is no reuse within the observed window.
- **Spatial Locality: Excellent (Locally).** Within each 6-element chunk, the addresses are consecutive integers. This suggests that loading a block will likely capture the immediate neighbors required for the next few cycles.

3. Hypothesize Workload Type: This pattern is highly characteristic of **Matrix Processing** or **Multi-Buffer Operations**.

- Example: Processing the first few columns of a matrix row-by-row, where the row width (stride) is 256 bytes.
- Imagine a loop: `for (r=0; r<4; r++) { process(Matrix[r][0...5]); }`.

Part II: Experimental Simulation & Data Collection (Trace Four)

Important Prediction: If the cache size is 128 or 256 bytes, then addresses 0, 256, 512, and 768 can all be mapped to the same set index (Set 0). But since the trace goes "forward" and doesn't

return, will these conflicts cause us problems? Probably not (the logic in Trace 2). Still, we'll test it and see.

STEP A: TESTING SPATIAL LOCALITY

- Our expectations are similar with Trace two. The overall look does not seem promising.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 1	128 Bytes	Direct Mapped (1-way)	2 Bytes	LRU
Test 2	128 Bytes	Direct Mapped (1-way)	4 Bytes	LRU
Test 3	128 Bytes	Direct Mapped (1-way)	8 Bytes	LRU

Analysis Of Varying Block Size

- Amazing. Block size 8 saved us again just like Trace two. The results are expected.

Test ID	Cache Size	Associativity	Block Size	Policy	Total Accesses	Hits	Misses	Hit Rate
Test 1	128 B	Direct Mapped	2 B	LRU	24	0	24	0%
Test 2	128 B	Direct Mapped	4 B	LRU	24	0	24	0%
Test 3	128 B	Direct Mapped	8 B	LRU	24	12	12	50%

Observation: Similar to Trace 2, small block sizes (2 and 4 Bytes) resulted in total failure (0% hit rate). However, increasing the Block Size to **8 Bytes** instantly improved the performance to **50%**, exhibiting a perfect "Miss-Hit" alternation.

Detailed Interpretation:

- **Spatial Locality Within Chunks:** The trace operates in localized "chunks" (e.g., 0, 4, 8 . . . then 256, 260 . . .). Within these chunks, the memory access stride is exactly **4 bytes**.
 - **In Test 2 (Block = 4B):** The cache loads exactly one integer. When the CPU asks for the *next* integer (4 bytes away), it falls into the *next* block, causing a Miss.
 - **In Test 3 (Block = 8B):** A single cache miss loads **8 bytes** of data.
 - Access 0x00: **Miss**. Loads 0x00 – 0x07.
 - Access 0x04: **Hit**. It is already in the loaded block.
 - Access 0x08: **Miss**. Loads 0x08 – 0x0F.
 - Access 0x0C: **Hit**.
- **Insensitivity to Large Gaps:** The trace jumps by 256 bytes between groups (e.g., from 20 to 256). While this is a large jump, it does not negatively impact the *Block Size* experiment because the cache simply treats the new address as a new compulsory miss. The benefit comes strictly from the *local* sequential access inside each group.

STEP B: TESTING CONFLICT

- There's a very interesting detail here: Address Conflict. However, since Trace "doesn't return" (i.e., it doesn't request 0 again after passing 256), this conflict doesn't cause a problem. We discard the old data (0) and write the new data (256), and life goes on.

Deney No	Cache Size	Associativity	Block Size	Replacement
Test 4	128 Bytes	2-way	8 Bytes	LRU
Test 5	128 Bytes	4-way	8 Bytes	LRU

Analysis Of Varying Associativity

Test ID	Cache Size	Associativity	Block Size	Policy	Hits	Misses	Hit Rate
Test 4	128 B	2-Way	8 B	LRU	12	12	50%
Test 5	128 B	4-Way	8 B	LRU	12	12	50%

Observation: Increasing the associativity yielded **no change** in performance. The hit rate remained constant at 50%.

Detailed Interpretation:

- Harmless Conflicts:** In a Direct-Mapped cache (Test 3), addresses separated by 128 bytes (or multiples) map to the same set. Here, the jump is 256 bytes (), meaning address 0 and address 256 map to the exact same cache index (Set 0).
- Why Associativity Didn't Help:** Technically, accessing address 256 *evicts* the block containing address 0. However, since the trace is "forward-moving" and never returns to read address 0 again, this eviction is harmless.
- Conclusion:** The "safety net" of 2-Way or 4-Way associativity is unnecessary because the evicted data is dead (no temporal locality). The misses are purely **compulsory** (first time access to a new block), not conflict-induced.

STEP C: TESTING CAPACITY SIZE OF CACHE

Deney No	Cache Size	Associativity	Block Size	Replacement
Test 6	64 Bytes	2-way	8 Bytes	LRU
Test 7	256 Bytes	2-way	8 Bytes	LRU

Analysis Of Varying Size Of Cache

Test ID	Cache Size	Associativity	Block Size	Policy	Hits	Misses	Hit Rate
Test 6	64 B	2-Way	8 B	LRU	12	12	50%
Test 7	256 B	2-Way	8 B	LRU	12	12	50%

Observation: Varying the cache capacity yielded identical results. A small 64-Byte cache performed equally as well as a larger 256-Byte cache.

Detailed Interpretation:

- **Local Working Set:** The trace processes data in small "chunks" of 24 bytes (addresses 0–20).
- **Dead Data:** Once the program finishes processing a chunk (e.g., 0–20) and jumps to the next chunk far away (256–276), the data from the previous chunk is never accessed again.
- **Capacity Irrelevance:** Storing the "old" chunk in a larger cache provides no benefit because the CPU never looks back. A 64-Byte cache is more than enough to hold the currently active chunk, making any extra capacity redundant.

STEP D: EFFECT OF REPLACEMENT POLICY

- We are expecting to remain the same hit rate.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 8	128 Bytes	2-way	8 Bytes	FIFO
Test 9	128 Bytes	2-way	8 Bytes	Random

Analysis Of Varying Replacement Method

Test ID	Cache Size	Assoc	Block Size	Policy	Hits	Misses	Hit Rate
Test 4 (Ref)	128 B	2-Way	8 B	LRU	12	12	50%
Test 8	128 B	2-Way	8 B	FIFO	12	12	50%
Test 9	128 B	2-Way	8 B	Random	12	12	50%

Observation: Changing the replacement policy had no effect on the hit rate (constant at 50%).

Detailed Interpretation:

- **Forward-Moving Pattern:** The trace moves strictly forward from one chunk to the next () .

- **Eviction Doesn't Matter:** When the cache fills up (or a set conflict occurs), the policy decides what to throw out. Since the "old" data (e.g., from address range 0-20) is dead and will never be requested again, it doesn't matter if we evict the oldest (FIFO), least recently used (LRU), or a random block. The penalty for eviction is zero.

Part III: Optimal Design Recommendation & Justification (Trace Four)

Recommendation: Based on the experimental data, the **single best cache configuration** for this workload is:

- **Cache Size:** 64 Bytes (Smallest tested)
- **Associativity:** Direct-Mapped (1-Way)
- **Block Size:** 8 Bytes (Crucial for Spatial Locality)
- **Replacement Policy:** Random

Justification: The optimization strategy relies on exploiting **Local Spatial Locality** while minimizing hardware cost for this specific "strided" pattern.

- **Block Size Success:** The only parameter that influenced performance was Block Size. Using **8 Bytes** allowed the cache to capture the sequential nature of the data *within* each chunk, achieving a 50% hit rate (versus 0% for smaller blocks).
- **Cost Efficiency:**
 - **Associativity:** Since conflicts between the large strides (0 vs 256) were harmless (no reuse), increasing associativity to 2-Way or 4-Way is a waste of resources. **Direct-Mapped** is chosen for speed and simplicity.
 - **Cache Size:** The active working set is very small (one chunk at a time). A **64-Byte** cache is sufficient. Larger caches just hold dead data.
- **Conclusion:** This design captures the maximum possible performance (limited by compulsory misses at the start of each block) with the minimum possible transistor count.

Part IV: The AI Sanity Check (Trace Four)

Prompt: "Analyze the memory trace 0, 4, 8, 12... 256, 260... and recommend a cache config."

AI Response: "This trace exhibits a **strided** or **chunked** access pattern, likely iterating through rows of a matrix or distinct buffers. **Recommendation:**

- **Block Size:** Use a large block size (e.g., 32 or 64 bytes). This is critical to maximize hits within each sequential chunk.
- **Prefetching:** Enable 'Stride Prefetching' to detect the jump between chunks (0 to 256).
- **Associativity:** Use 2-Way Set Associative. The stride of 256 bytes is a power of 2, which often causes 'Thrashing' in Direct-Mapped caches (mapping to the same set). Associativity will prevent this."

Critical Analysis (Human vs. AI)

Critique: The AI correctly identified the pattern but gave a generic advice regarding associativity that was technically correct but practically useless for *this specific* trace.

Where the AI got it right: It correctly identified the "Chunked" nature and the importance of **Block Size** to capture the sequential part of the access.

Where the AI fell short (The "Reuse" Insight):

- **The "Conflict" Fallacy:** The AI recommended **2-Way Associativity** to prevent thrashing between address 0 and 256. While it is true that they map to the same set, my simulation proved that this conflict is **irrelevant** because the program never returns to address 0. The AI tried to solve a problem that doesn't impact performance, leading to an over-engineered solution.
- **Cost vs. Benefit:** My analysis showed that a Direct-Mapped cache performs identically to the AI's suggested 2-Way cache but is cheaper to build.

TRACE FIVE: (0,4,8,12, 256,260,264,268, 0,4,8,12, 256,260,264,268, 0,4,8,12,256,260,264,268)

Part I: Workload Characterization (Trace Five)

1. Manual Analysis: The trace consists of an alternating pattern between two distinct memory regions: a lower chunk (0-12) and a higher chunk (256-268). This pattern repeats 3 times () .

2. Hypothesize Locality:

- **Temporal Locality: Excellent.** The data blocks are reused multiple times. After processing the high chunk (256), the program returns to the low chunk (0). Ideally, the cache should retain the low chunk for the next iteration.
- **Spatial Locality: Excellent.** Within each chunk, accesses are sequential with a 4-byte stride.
- **The Conflict Hazard:** Addresses 0 and 256 are exactly 256 bytes apart. In many Direct-Mapped cache configurations (where 256 is a multiple of the cache size), these addresses map to the **same set index**. This creates a high risk of "**Thrashing**," where the two chunks constantly evict each other.

3. Hypothesize Workload Type: This pattern represents a "**Ping-Pong**" workload or alternating buffer processing. It simulates a scenario where two data structures (e.g., source and destination arrays) are accessed in an interleaved manner within a loop.

Part II: Experimental Simulation Plan (Trace Five)

STEP A: TESTING SPATIAL LOCALITY

- We are expecting to see similarities with trace two and four. But again, our hit rate will stuck at 50% because of conflict does not allow us to use temporal locality.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 1	128 Bytes	Direct Mapped (1-way)	2 Bytes	LRU
Test 2	128 Bytes	Direct Mapped (1-way)	4 Bytes	LRU
Test 3	128 Bytes	Direct Mapped (1-way)	8 Bytes	LRU

Analysis Of Varying Block Size

Test ID	Cache Size	Associativity	Block Size	Policy	Total Accesses	Hits	Misses	Hit Rate
Test 1	128 B	Direct Mapped	2 B	LRU	24	0	24	0%
Test 2	128 B	Direct Mapped	4 B	LRU	24	0	24	0%
Test 3	128 B	Direct Mapped	8 B	LRU	24	12	12	50%

Observation: Test 3 (8 Byte Block) achieved a **50% hit rate**, while smaller blocks failed completely. However, despite the repetitive nature of the trace, the hit rate did not exceed 50% in the Direct-Mapped configuration.

Detailed Interpretation:

- **Spatial Locality Success:** The jump to 50% in Test 3 is due to Block Size. Loading address 0 (Miss) brings in 0-7. The next access to 4 is a Hit. This accounts for exactly half of the accesses being hits.
- **Temporal Locality Failure (The "Thrashing" Problem):** The trace alternates between address 0 (Group A) and 256 (Group B).
 - In a **Direct-Mapped** cache with 128B size and 8B blocks, there are 16 Sets.
 - Address 0 maps to **Set 0**.
 - Address 256 () also maps to **Set 0 ()**.
- **Conclusion:** Every time the program switches groups, the new data **evicts** the old data from Set 0. Even though the program returns to Group A immediately, it has already been kicked out. This phenomenon is known as **Cache Thrashing**.

STEP B: TESTING CONFLICT

- 0 and 256 can be next to each other if the associativity is more than 2. We expect hit rate to reach a pretty high score.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 4	128 Bytes	2-way	8 Bytes	LRU
Test 5	128 Bytes	4-way	8 Bytes	LRU

Analysis Of Varying Associativity

Test ID	Cache Size	Associativity	Block Size	Policy	Hits	Misses	Hit Rate
Test 4	128 B	2-Way	8 B	LRU	20	4	83.3%
Test 5	128 B	4-Way	8 B	LRU	20	4	83.3%

Observation: Increasing the associativity to 2-Way (Test 4) caused a massive performance leap from 50% to **83.3%**.

Detailed Interpretation:

- **Solving the Conflict:** By changing the cache to **2-Way Set Associative**, each Set now has 2 "ways" (slots) to store data.
 - **Group A (Address 0):** Loads into Set 0, Way 0.
 - **Group B (Address 256):** Loads into Set 0, Way 1.
- **Result:** Both working sets can coexist in the cache simultaneously without evicting each other.
- **The Hit Rate Logic:**
 - **Pass 1:** 4 Compulsory Misses (Loading the unique blocks for the first time).
 - **Pass 2 & 3:** 100% Hits (All data is now safely in cache).
 - Total: 20 Hits / 24 Accesses = **83.3%**.
- **Conclusion:** This experiment perfectly demonstrates the value of Associativity. It eliminated the conflict misses that crippled the Direct-Mapped design.

STEP C: ANALYSIS OF VARYING CACHE SIZE

- Even let the size of cache be 64, our data maximum is 32. So this combination handles the data well. We don't expect any change of hit rate.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 6	64 Bytes	2-way	8 Bytes	LRU
Test 7	256 Bytes	2-way	8 Bytes	LRU

Analysis Of Varying Size Of Cache

Test ID	Cache Size	Associativity	Block Size	Policy	Hits	Misses	Hit Rate
Test 6	64 B	2-Way	8 B	LRU	20	4	83.3%
Test 7	256 B	2-Way	8 B	LRU	20	4	83.3%

Observation: Varying the cache capacity had **zero impact** on performance. The hit rate remained constant at 83.3%.

Detailed Interpretation:

- **Working Set vs. Capacity:** The workload consists of two active groups: Group A (0-12) and Group B (256-268).
 - Each group occupies approx. 16 bytes.
 - Total Working Set 32 Bytes.
- **Why 64B is Enough:** Even the smallest tested cache (64 Bytes) is double the size of the working set. Therefore, the cache is never "full" in terms of total capacity.
- **The Real Bottleneck was Conflict, not Capacity:** The misses in previous tests (Direct-Mapped) were due to *mapping conflicts* (both groups fighting for Set 0), not a lack of total space. Since 2-Way associativity (Step B) solved the mapping conflict, adding more raw capacity (256B) provides no additional benefit.

STEP D: EFFECT OF REPLACEMENT POLICY

- The replacement policy unlikely to be crucial for our efficiency. Hardware parameters will remain the same while varying the replacement policy.

Exp. No	Cache Size	Associativity	Block Size	Replacement
Test 8	128 Bytes	2-way	8 Bytes	FIFO
Test 9	128 Bytes	2-way	8 Bytes	Random

Analysis Of Varying Replacement Method

Test ID	Cache Size	Assoc	Block Size	Policy	Hits	Misses	Hit Rate
Test 4 (Ref)	128 B	2-Way	8 B	LRU	20	4	83.3%
Test 8	128 B	2-Way	8 B	FIFO	20	4	83.3%
Test 9	128 B	2-Way	8 B	Random	20	4	83.3%

Observation: Changing the replacement policy resulted in identical hit rates (83.3%).

Detailed Interpretation:

- **Associativity Solved the Eviction Issue:** In the 2-Way configuration, Address 0 and Address 256 can sit side-by-side in the same set (Set 0).
- **No "Forced" Evictions:** Since the working set fits perfectly into the 2 ways available, the cache rarely needs to evict a *useful* block to make room for another useful block. The "ping-pong" effect is handled by the hardware slots, not the replacement logic.
- **Conclusion:** When the hardware resources (ways) are sufficient to hold the contending blocks, the intelligence of the replacement algorithm (LRU vs. Random) becomes irrelevant.

Part III: Optimal Design Recommendation & Justification (Trace 5)

Recommendation Based on the experimental data, the **single best cache configuration** for this workload is:

- **Cache Size:** 64 Bytes (Smallest tested)
- **Associativity:** 2-Way Set Associative (**Critical**)
- **Block Size:** 8 Bytes
- **Replacement Policy:** Random (or LRU)

Justification The optimization strategy here is "**Conflict Resolution.**" Unlike previous traces where Associativity was useless, here it is the single most important factor.

- **Why 2-Way is Essential:** The trace exhibits a "Ping-Pong" pattern between addresses 0 and 256, which map to the same set index. In a Direct-Mapped cache, this caused **Thrashing** (50% Hit Rate). Switching to **2-Way Associativity** eliminated this conflict, boosting performance to **83.3%**. This is the only configuration change that provided a tangible benefit.
- **Why 64 Bytes:** Once the conflict was resolved by associativity, the total data size (32 Bytes) fit comfortably within the smallest 64-Byte cache. Larger caches are unnecessary.
- **Cost/Performance Balance:** A **64-Byte, 2-Way** cache strikes the perfect balance. It provides the complex mapping logic (2-Way) needed to prevent thrashing, but keeps the total memory size (64B) minimal to save power and area.

Part IV: The AI Sanity Check (Trace 5)

Prompt: "Analyze the memory trace 0, 4, 8... 256, 260... 0, 4... and recommend a cache config."

AI Response: "This trace shows an alternating access pattern between two memory regions (0-12 and 256-268).**Recommendation:**

- **Associativity: 2-Way or 4-Way Set Associative** is strongly recommended. Addresses 0 and 256 are multiples of 256, which likely causes them to map to the same cache line. A Direct-Mapped cache would suffer from severe **thrashing**.
- **Cache Size:** A small cache is sufficient since the dataset is small.
- **Replacement Policy:** LRU is best to ensure the most recently used 'chunk' stays in memory, though with high associativity, this matters less."

Critical Analysis (Human vs. AI)

Critique: This is the one instance where the AI and my experimental analysis aligned perfectly. The AI correctly predicted the "Thrashing" phenomenon, which was the defining characteristic of this trace.

Where the AI got it right: The AI accurately identified the mathematical conflict between address 0 and 256 (). It correctly advised that **Associativity** was the key to solving this problem.

My experiment (Step B) confirmed this, showing a massive jump from 50% (Direct) to 83.3% (2-Way).

Where the Human Insight Added Value: While the AI gave a qualitative prediction ("Thrashing will occur"), my simulation quantified the *impact* of that thrashing. I proved that Thrashing capped the performance at exactly 50% (due to spatial locality saving the day), and that 2-Way associativity restored it to the maximum theoretical limit (minus compulsory misses). The AI provided the theory; I provided the engineering proof.

References:

1. D. A. Patterson and J. L. Hennessy, *Computer Organization and Design: The Hardware/Software Interface*, 5th ed. Morgan Kaufmann, 2013.
2. Istanbul Health and Technology University, "COE-301 Computer Architecture Course Lecture Notes & Project Guidelines," Fall 2025-2026.
3. Cache Simulator Tool. [Online].
Available: <https://courses.cs.washington.edu/courses/cse351/cachesim/>
4. OpenAI, "ChatGPT (GPT-4 Model)," Used for workload analysis comparison (Part IV), 2025. [Online]. Available: <https://chat.openai.com>
5. Google, "Gemini (Pro Model)," Used for experimental data verification and report structuring, 2025. [Online]. Available: <https://gemini.google.com>