

Prepare Phase

Total Time: 9~11 minutes | **Storage:** 3 GB

Three-Stage Process

Stage 1: Data Ingestion (30s)

- Load 245M rows from CSV into DuckDB

Stage 2: Type Partitioning (3~4mins)

- Split data into 4 columnar Parquet files by event type (impression, click, purchase, serve)
- Enables 75% I/O reduction for type-filtered queries
- Parquet compression achieves 5-10x size reduction

Stage 3: Pre-Aggregation (5~7mins)

- Build 51 specialized aggregation tables covering common query patterns
- **Temporal:** Daily (366 rows), weekly (53), hourly (8,760), minute (525k)
- **Dimensional:** Country (12), publisher (1,114), advertiser (1,654)
- **Multi-dimensional:** Publisher×country×daily (4.87M rows), advertiser×hourly (14.43M rows), etc.

The core insight: Ad analytics queries follow predictable patterns (80% time-based, 60% type-filtered). By pre-computing these combinations, we trade 9 minutes of preparation for sub-100ms query execution.

A Test Config

Dataset: 5 given sample queries

Run: 10 runs (drop min/max)

Cache: Cleared between cold runs, persistent for warm runs

Prepare Phase:

	All Values (10 runs)
Baseline (ms)	68854, 70560, 66619, 71870, 66899, 71184, 70254, 74182, 83441, 81800
Cold (ms)	88.268, 98.286, 89.613, 100.308, 89.282, 89.469, 89.051, 89.090, 89.602, 96.720
Warm (ms)	47.423, 46.622, 47.133, 47.275, 47.507, 47.339, 47.370, 47.193, 45.611, 48.956

Metric	Mean	Std Dev	Speedup
Baseline	70480 ms	± 2340	/
Cold start	91.389 ms	± 3.802	771x
Warm start	47.233 ms	± 0.275	1492x

Correctness Check:

- **All queries results match with baseline results**

(Note: Minor column name formatting differences (count_star()) vs COUNT(*)) do not affect numerical accuracy)

B Test Config

Dataset: 10 generate similar queries

Run: 5 runs (drop min/max)

Cache: Cleared between cold runs, persistent for warm runs

	All Values (5 runs)
Baseline (ms)	145298, 154602, 174047, 152595, 175037
Cold (ms)	119.188, 117.016, 121.583, 120.873, 127.269
Warm (ms)	73.888, 70.962, 72.391, 74.287, 76.852

Metric	Mean	Std Dev	Speedup
Baseline	160415 ms	± 11849	/
Cold start	120.548 ms	± 1.230	1,330x
Warm start	73.522 ms	± 1.000	2,182x

Correctness Check:

- **All queries results match with baseline results**

(Note: Minor column name formatting differences (count_star()) vs COUNT(*)) do not affect numerical accuracy)

C Test Config

Dataset: 70 total queries (65 generate + 5 given)

Run: 3 runs

Cache: Cleared between cold runs, persistent for warm runs

	All Values (3 runs)
Baseline (ms)	1387852, 1198912, 1408873
Cold (ms)	299.675, 306.157, 304.423
Warm (ms)	212.688, 214.721, 218.002

Metric	Mean	Std Dev	Speedup
Baseline	1331879 ms	± 115631	/
Cold start	303.418 ms	± 3.356	4,389x
Warm start	215.137 ms	± 2.681	6,190x

Note: Crazy speed up due to some duplication in the 70 queries so cache hits!!

Correctness Check:

- **All queries results match with baseline results**
(Note: Minor column name formatting differences (count_star()) vs COUNT(*)) do not affect numerical accuracy)

Key Findings

Consistent extreme speedup: hugely faster across all test configurations

Sub-100ms performance: Even cold start averages <140ms for 10 complex queries

Excellent stability: low std demonstrates production-ready consistency

Scalability validated: Performance holds on 14x larger query set (Test C)

Warm start helps a lot: Performance by warm start almost twice as cold start