# Who am I, and why I talking about this?

- I've been cleaning up some IT shit since the previous millennium. ;)

- Last 5 years in 24x7x365 mostly ClickHouse

**Maintainer/Contributor**

- clickhouse-backup

- clickhouse-grafana

- clickhouse-operator

- altinity-mcp

2

https://github.com/Slach

# ClickHouse Quiz: Raise your hand!

1. Who knows what ClickHouse is?

2. Who works with ClickHouse right now?

3. Who has compressed data more than 10TB right now?

# Common usage of resources: IO, CPU, RAM, NET.

ClickHouse is a data excavator: for each query, it tries to allocate the maximum allowable resources to process your query for maximum efficiency.

Spike workloads. Many concurrent queries will struggle for resources.

Background merges, message broker streaming, and dictionaries loading also allocate resources.

Don't mix clickhouse with other software components in production (including keeper) on the same hardware.
Recommendations of proportions: 10 TB of compressed data per server with 128 GB RAM and 32 CPU cores.
10 Gbit network is always better than a 1 Gbit network. ;)

# ClickHouse store a lot of internal telemetry

```
# metrics related              # runtime                    # replication and clusters
metrics                        processes                    distributed_ddl_queue
events                         merges                       distribution_queue
metric_log                     moves                        part_moves_between_shards
asynchronous_metrics           mutations                    replicas
asynchronous_metric_log                                     replicated_fetches
                               user_processes               replication_queue
dimensional_metrics            view_refreshes
histogram_metrics              dictionaries                 # disk space usage
                                                            disks
latency_buckets                asynchronous_inserts         parts
latency_log                    asynchronous_loader          tables
                                                            parts_columns
# stack related                query_cache                  projection_parts
trace_log                      query_condition_cache        projection_parts_columns
stack_trace                                                 projections
                                                            data_skipping_indices
# logs                         # data streaming
text_log                       azure_queue                  backups
                               s3queue                      backup_log
errors                         kafka_consumers
warnings                       iceberg_history
```

# Internal basic monitoring - part 1/4

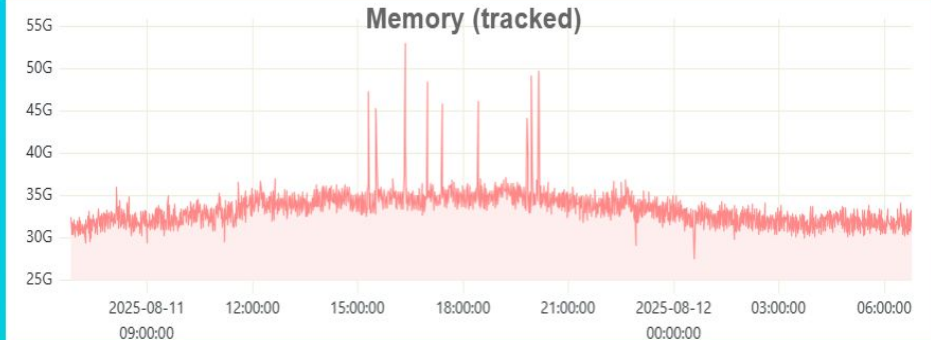http://localhost:8123/dashboards
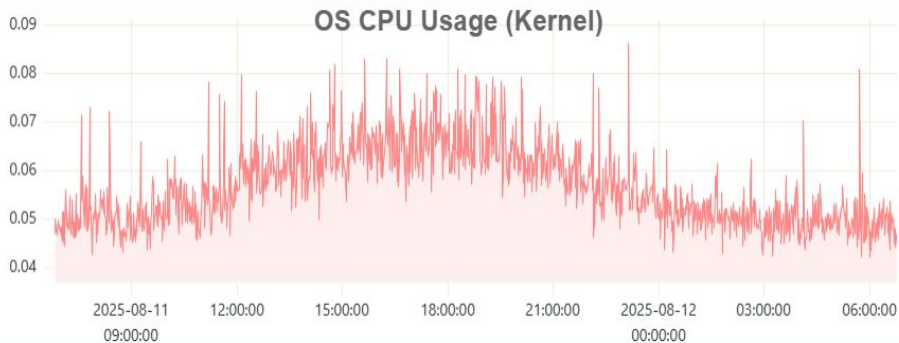
# Internal basic monitoring - part 2/4

http://localhost:8123/dashboards

# Internal basic monitoring - part 3/4

http://localhost:8123/dashboards

# Internal basic monitoring - part 4/4

# Heatmap queries by frequency

https://github.com/zix99/rare + https://github.com/Slach/clickhouse-timeline

# Heatmap error:query_hash by frequency

# Heatmap queries by CPU Usage

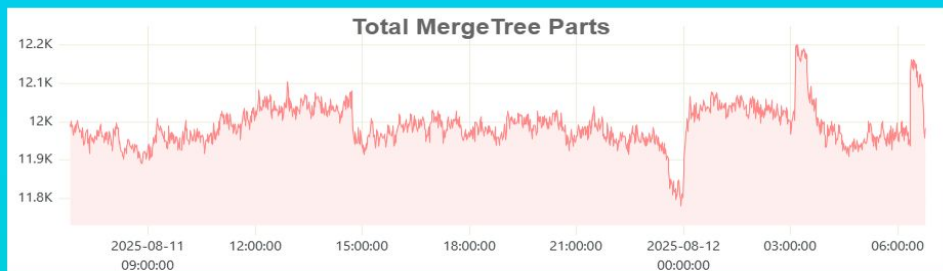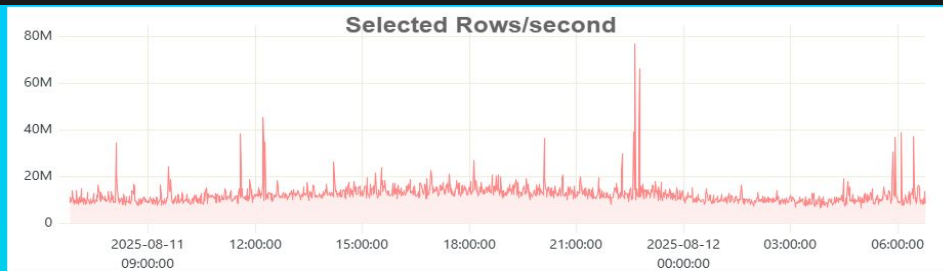# Heatmap queries by Memory Usage

# Heatmap queries by read_bytes



bytes read per minutes

# Heatmap queries by written_bytes



15

# Found bad `normalized_query_hash` and what?

Make flamegraph https://github.com/Slach/clickhouse-flamegraph

```
clickhouse-client -q "SELECT 'clickhouse-flamegraph
--query-ids=' ||
arrayStringConcat(groupArray(10)(query_id),',') ||
'\n' FROM system.query_log WHERE
normalized_query_hash=? AND event_date=? AND
event_time BETWEEN ? AND ? ORDER BY
query_duration_ms DESC LIMIT 10
FORMAT TSVRaw" | bash
```

# Flamegraph: what is it?

https://github.com/YS-L/flamelens

# Found bad `normalized_query_hash` - what next?

Look percentiles for ProfileEvents

```
┌Profile Events: 9689294398787416682 (2025-08-13 17:10:00 +04:00 to 2025-08-13 17:20:00 +04:00)┐
            Event              Count       p50            p90            p99
CompressedReadBufferBlocks      744        4.00           7.00          41.57
CompressedReadBufferBytes       744   80.46 thousand  200.05 thousand   2.54 million
ConcurrencyControlSlotsAcquired 1118      3.00           3.00           3.00
ConcurrencyControlSlotsAcquiredNonCompeting 1118  1.00   1.00           1.00
ConcurrencyControlSlotsGranted  1118      1.00           1.00           1.00
ContextLock                     1118     29.00          29.00          29.00
ContextLockWaitMicroseconds       1       1.00           1.00           1.00
FunctionExecute                 1118      2.00           2.00           2.00
GlobalThreadPoolJobs            1118      5.00           5.00           5.00
GlobalThreadPoolLockWaitMicroseconds 201  6.00          7.00          11.00
IOBufferAllocBytes              744    1.09 million   1.10 million    1.11 million
IOBufferAllocs                  744       3.00           3.00           3.00
InitialQuery                    1118      1.00           1.00           1.00
InsertQuery                     1118      1.00           1.00           1.00
InsertedBytes                   1118  80.00 thousand  200.00 thousand   2.63 million
InsertedRows                    1118  10.00 thousand   10.00 thousand  10.00 thousand
InterfaceNativeReceiveBytes     1118  38.43 thousand  120.06 thousand  210.06 thousand
InterfaceNativeSendBytes        1118   5.56 thousand    5.80 thousand   6.05 thousand
LocalThreadPoolExpansions       1118      4.00           4.00           4.00
LocalThreadPoolJobs             1118      4.00           4.00           4.00
LocalThreadPoolLockWaitMicroseconds 25    5.00          6.60           7.00
LocalThreadPoolShrinks          1118      4.00           4.00           4.00
LocalThreadPoolThreadCreationMicroseconds 1118  22.00   28.00         32.00
NetworkReceiveBytes             1118  38.43 thousand  120.06 thousand  210.06 thousand
NetworkReceiveElapsedMicroseconds 1118 47.64 thousand 195.05 thousand 298.03 thousand
NetworkSendBytes                1118   5.56 thousand    5.80 thousand   6.05 thousand
NetworkSendElapsedMicroseconds  1118     58.00          65.00          74.83
OSCPUVirtualTimeMicroseconds    1118    920.50       1.09 thousand    3.00 thousand
OSCPUWaitMicroseconds             37     13.00          25.20          62.04
OSReadChars                     1118  41.86 thousand  123.80 thousand  213.96 thousand
OSWriteChars                    1118   5.66 thousand    6.12 thousand   6.44 thousand
Query                           1118      1.00           1.00           1.00
```

┌Normalized Query──────
`INSERT INTO`
`temporary_roundtrip_table ( ``col``,``id`` ) VALUES`

when p99 >= 4*p50

it means query data is unstable

# MEMORY_LIMIT_EXCEED - How Memory Tracker works

When memory allocates, clickhouse looks to the memory tracker counter.
If value is more than the soft limit value, then wait
memory_usage_overcommit_max_wait_microseconds. If memory is not freed, then any next
allocation will throw an exception, and any other query could be canceled.

Memory limit exceeded (for query)
Tweak max_bytes_before_external_sort, max_bytes_before_external_group_by,
max_threads, max_insert_threads, max_memory_usage
Memory limit exceeded (for user) - 90% use cases another queries.
max_memory_usage_for_user
Memory limit exceeded (total)
Try to decrease max_server_memory_usage_ratio (to avoid OOM)
https://kb.altinity.com/altinity-kb-setup-and-maintenance/altinity-kb-who-ate-my-memory/

# Who ate my memory



| Group | Name | chi-github-github-0-0-0 |
|---|---|---|
| OS | OSMemoryTotal | 123.55 GB |
| OS | OSMemoryAvailable | 116.75 GB |
| OS | OSMemoryFreePlusCached | 115.42 GB |
| OS | OSMemoryCached | 105.33 GB |
| OS | OSMemoryFreeWithoutCached | 10.09 GB |
| OS | OSMemoryBuffers | 1.19 GB |
| Process | MemoryVirtual | 294.69 GB |
| Process | MemoryDataAndStack | 293.71 GB |
| Process | MemoryResidentMax | 61.11 GB |
| Process | MemoryResidentWithoutPageCache | 4.06 GB |
| Process | MemoryResident | 4.06 GB |
| Process | MemoryCode | 279.38 MB |
| Process | MemoryShared | 199.44 MB |
| Caches | FilesystemCacheBytes | 200.00 GB |
| MMaps | MMappedFileBytes | 535.01 MB |
| StorageBuffer | StorageBufferBytes | 0 B |
| MemoryTables | Memory | 0 B |
| Dictionaries | Flat | 517.62 MB |
| Dictionaries | ComplexKeyHashed | 24.13 MB |
| Dictionaries | ComplexHashedArray | 1.29 MB |
| Dictionaries | Hashed | 256.45 KB |
| Dictionaries | | 0 B |
| PrimaryKeys | db:default | 99.18 MB |
| PrimaryKeys | db:eth | 89.17 MB |
| PrimaryKeys | db:tiered | 288.72 KB |
| PrimaryKeys | db:system | 261.01 KB |
| PrimaryKeys | db:fts | 0 B |
| PrimaryKeys | db:git | 0 B |
| PrimaryKeys | db:demo | 0 B |
| PrimaryKeys | db:staging | 0 B |
| PrimaryKeys | db:mindsdb | 0 B |
| PrimaryKeys | db:eth_final | 0 B |
| PrimaryKeys | db:clickbench | 0 B |
| PrimaryKeys | db:import | 0 B |
| PrimaryKeys | db:eth_old | 0 B |
| PrimaryKeys | db:training | 0 B |

## Memory usage

| Group | Name | chi-github-github-0-0-0 |
|---|---|---|
| Queries | SELECT | 4.00 MB |
| UserMemoryTracking | slach | 4.00 MB |
| UserMemoryTracking | | 27.24 KB |
| UserMemoryTracking | default | 12.59 KB |
| UserMemoryTracking | updater | 4.80 KB |
| UserMemoryTracking | clickhouse_operator | 0 B |
| UserMemoryTracking | dmitrii | 0 B |
| UserMemoryTracking | demo | 0 B |
| UserMemoryTracking | blockchain_explorer | 0 B |
| UserMemoryTracking | altinity | 0 B |
| UserMemoryTracking | rill | 0 B |
| QueryCacheBytes | | 427.83 KB |
| FileBuffersVirtual | OpenFileForRead | 54 MB |
| FileBuffersVirtual | OpenFileForWrite | 6 MB |
| ThreadStacksVirtual | GlobalThread | 9.05 GB |
| MemoryTracking | total | 2.33 GB |

20

# Explain bad normalized_query_hash

**Query Text**

```sql
SELECT
`event_type`,COUNT(*) as `__Count`
FROM
github_events

WHERE
1=1
GROUP BY
`event_type`
```

**EXPLAIN PLAN indexes=1, projections=1**

```
Expression ((Project names + Projection))
  Aggregating
    Expression ((Before GROUP BY + Change column names to column identifiers))
      ReadFromMergeTree (default.github_events)
      Indexes:
        MinMax
          Condition: true
          Parts: 45/45
          Granules: 1258981/1258981
        Partition
          Condition: true
          Parts: 45/45
          Granules: 1258981/1258981
        PrimaryKey
          Condition: true
```

**EXPLAIN PIPELINE**

```
(Expression)
ExpressionTransform × 32
  (Aggregating)
  Resize 32 → 32
    AggregatingTransform × 32
      (Expression)
      ExpressionTransform × 32
        (ReadFromMergeTree)
        MergeTreeSelect(pool: ReadPool, algorithm: Thread) × 32 0 → 1
```

**EXPLAIN ESTIMATE**

| database.table | parts | rows | marks |
|---|---|---|---|
| default.github_events | 45 | 10G | 1M |

# Thanks for your attention :)
## Please ask questions!



@BLOODJAZMAN



Altinity