

# VeloxDB - A Locally Grown Database

Ivan Savu

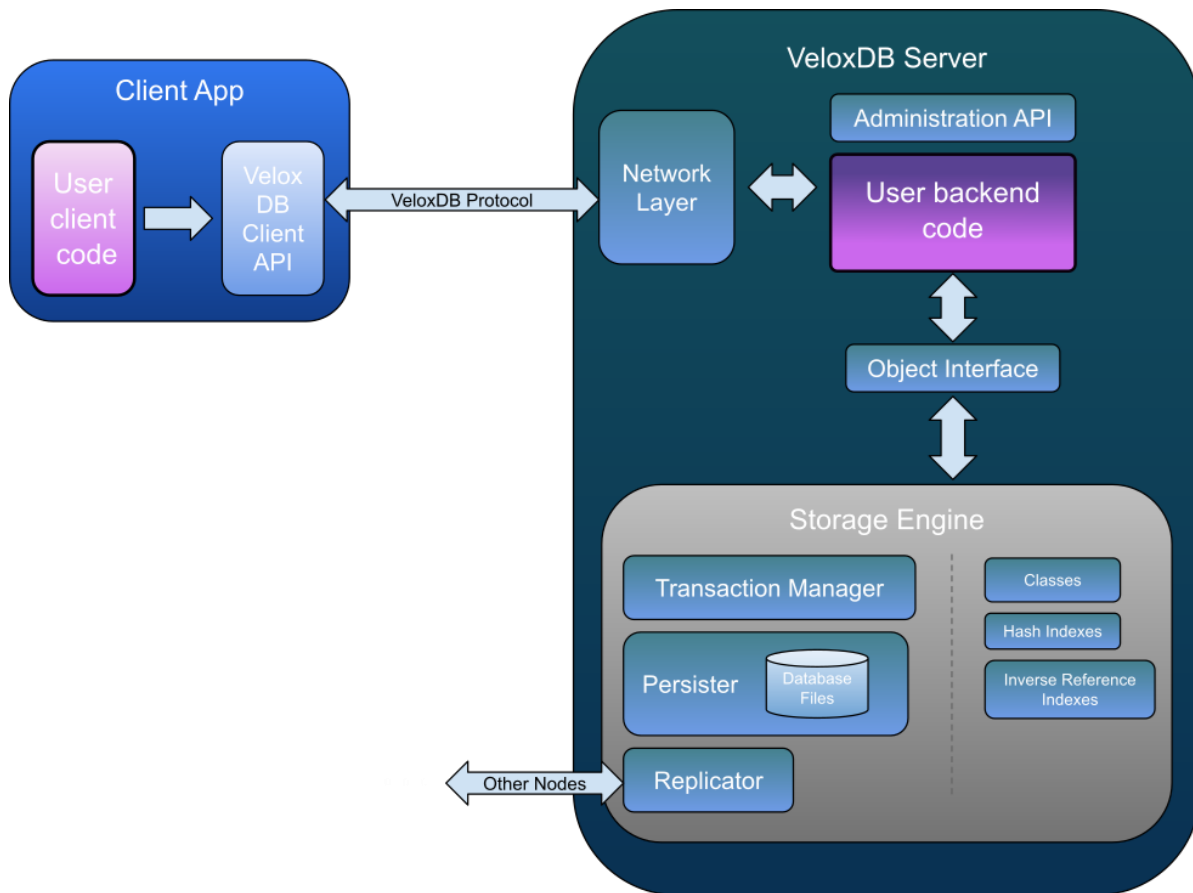


## ●Members:

- Nikola Morača
- Denis Balog
- Ivan Savu



- Fast (2.5 million trans/second on AWS c6a.8xlarge)
- ACID
- Object Oriented (ORM as first class citizen)
- .NET based
- Cross platform
- Open Source (MIT license)
- High Availability
- Read scale out
- In - Memory



- User provides model as C# abstract classes
- VeloxDB implements abstract classes with code generation at runtime
- Each class has a pointer to unmanaged data.
- Getters are implemented by directly reading data
- Setters accumulate changes for writing to disk and replication

## Example model:

```
[DatabaseClass]
public abstract class Post : DatabaseObject
{
    [DatabaseProperty]
    public abstract string Title { get; set; }

    [DatabaseProperty]
    public abstract string Content { get; set; }
}
```

## Example database operation:

```
[DbAPIOperation]  
public long CreatePost(ObjectModel om, PostDTO  
post)  
{  
    Post newPost = om.CreateObject<Post>();  
    newPost.Title = post.Title;
```

- Complexity
- Performance problems
- N+1 query problem
- Incomplete abstraction
- Difficult to debug



- Uses Multi Version Concurrency Control (MVCC) for transaction management.
- Handles updates by treating them as inserts with new versions.
- Stores multiple versions of the same object (Multiversioning).

T1: AddBook("The Great Gatsby", "F. Scott Fitzegarld")

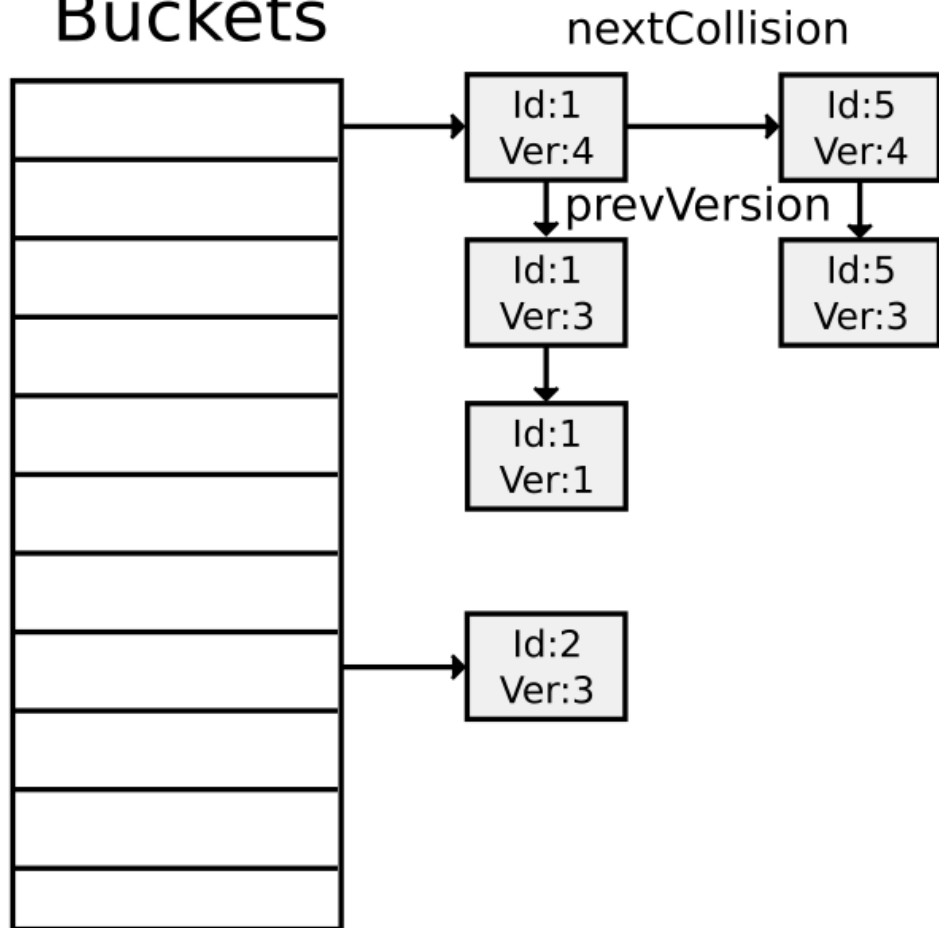
T2: AddBook("To Kill a Mockingbird", "Harper Lee")

T3: AddBook("One Hundred Years of Solitude", "Gabriel Garcia Marquez")

T4: UpdateBookAuthor(1, "F. Scott Fitzegarld")

<b>Id</b>	<b>Title</b>	<b>Author</b>	<b>Version</b>
1	The Great Gatsby	F. Scott Fitzegarld	1
1	The Great Gatsby	F. Scott Fitzgerald	4
2	To Kill a Mockingbird	Harper Lee	2
3	One Hundred Years of Solitude	Gabriel Garcia Marquez	3

# Buckets

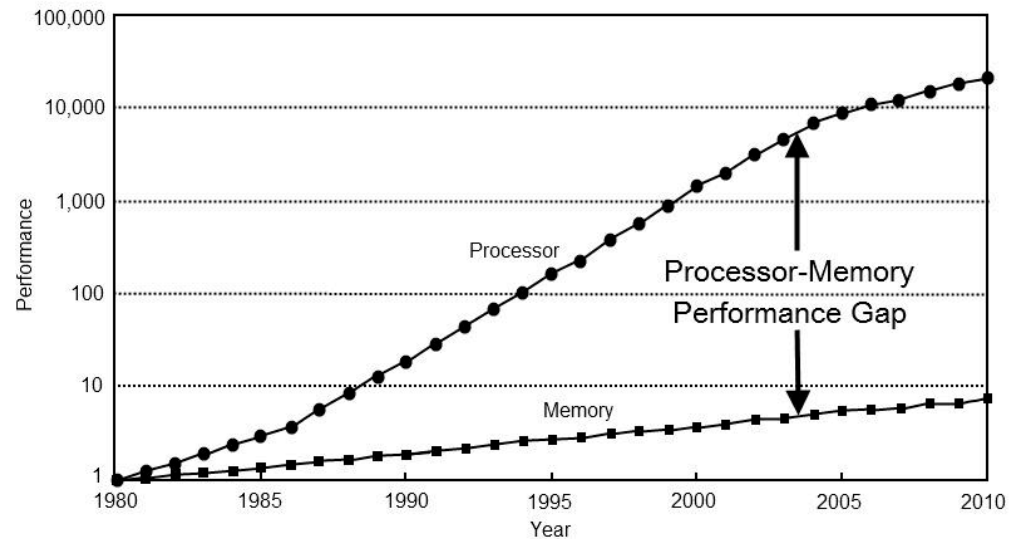


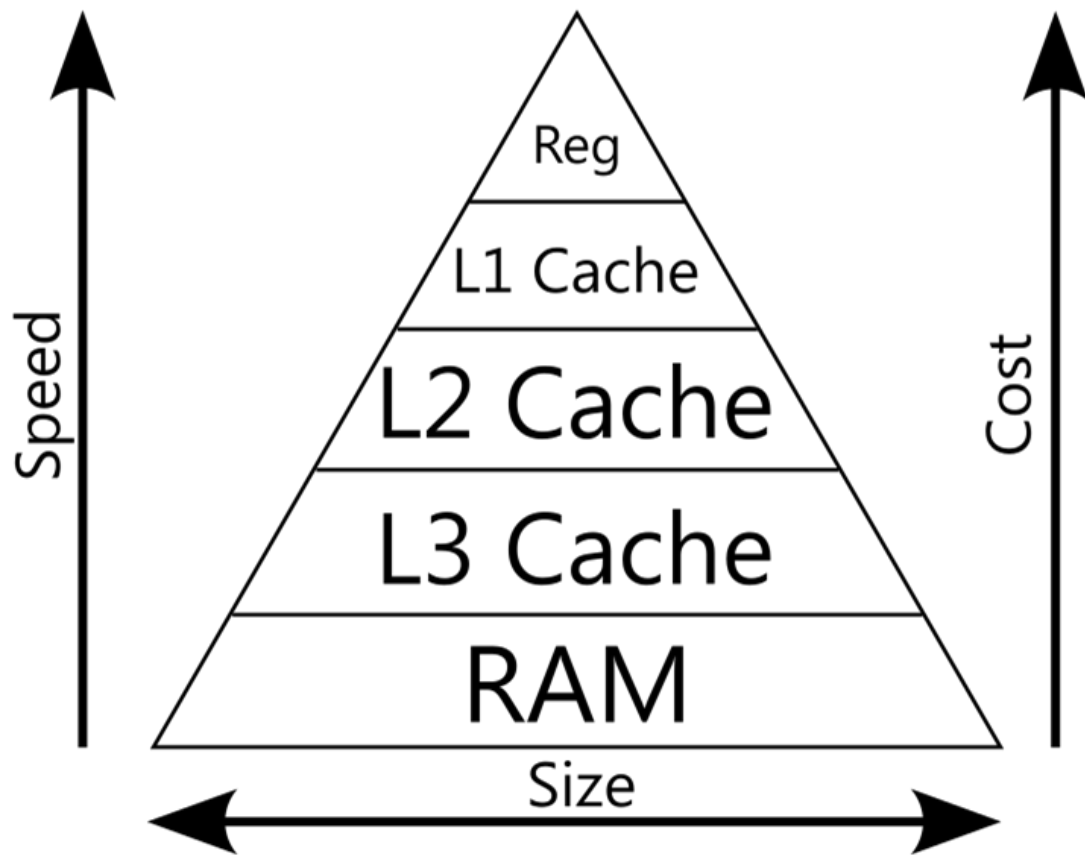
- We use very fine grained locks - employing a significant number of locks
- .NET locks are objects, each object has 24-byte overhead
- Spinlock
  - 1 bit needed
  - Suitable for rare contentions
  - Pure userspace

```
void EnterLock()
{
    while(true)
        observedState = state;

    if (GetBit(observedState, LOCK_BIT) != LOCKED)
    {
        // Calculate lock state with lock bit set
        lockedState = SetBit(observedState, LOCK_BIT, LOCKED)
        if (CompareAndExchange(ref state, lockedState, observedState) == observedState):
            return // Lock acquired
    }
}
```

- CPUs are improving at faster rate than memory
- Bottlenecks:
  - Bandwidth
  - Latency







# Processor 0

Core 0

L1 Cache

L2 Cache

Core 1

L1 Cache

L2 Cache

Core 2

L1 Cache

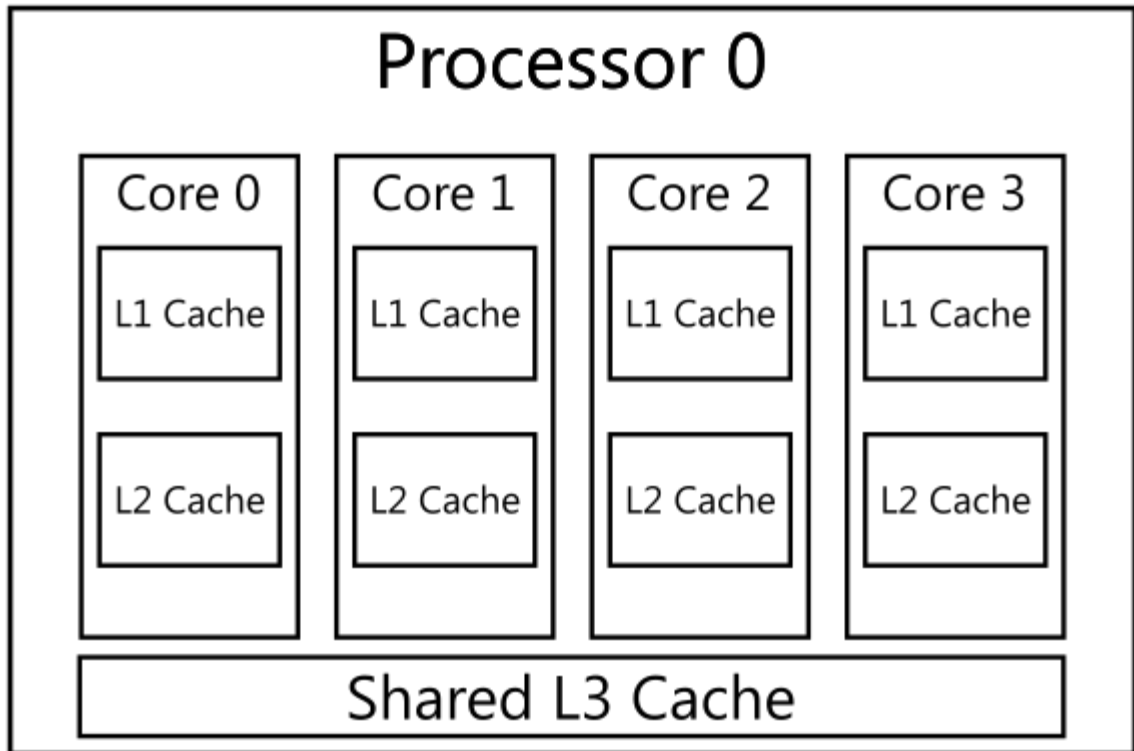
L2 Cache

Core 3

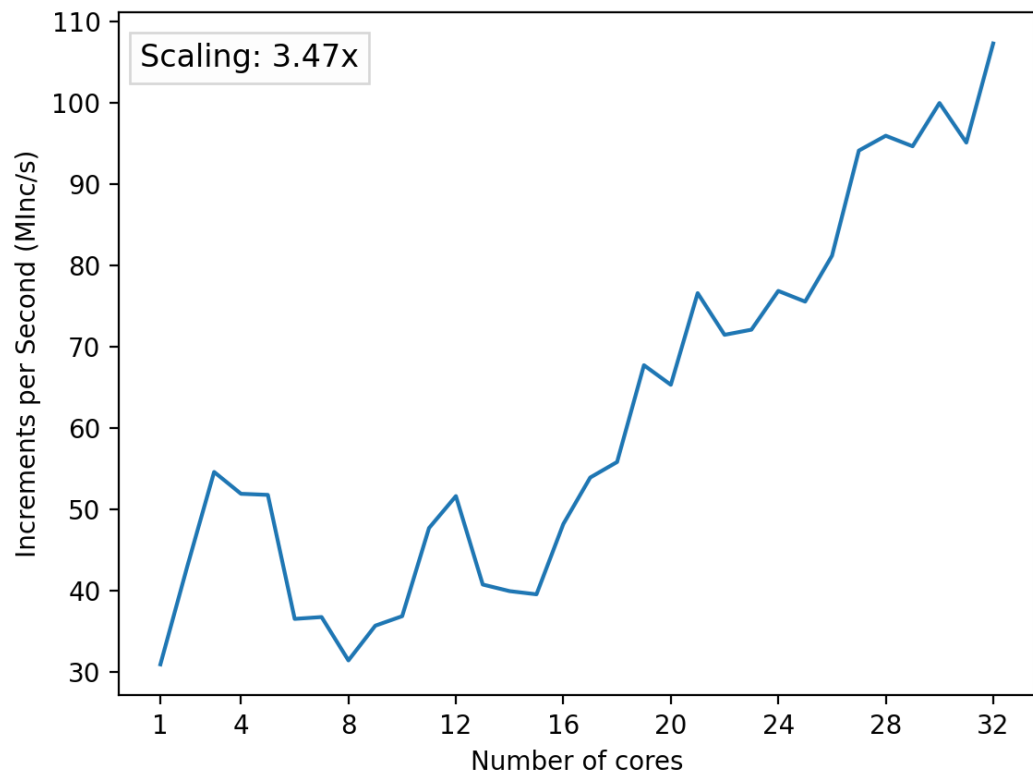
L1 Cache

L2 Cache

Shared L3 Cache

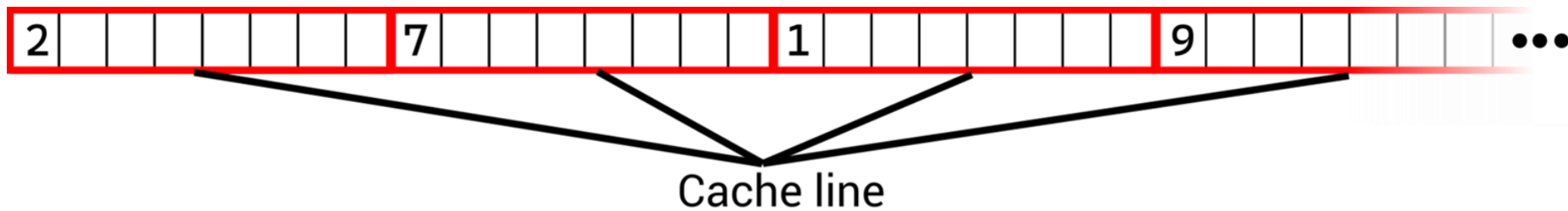


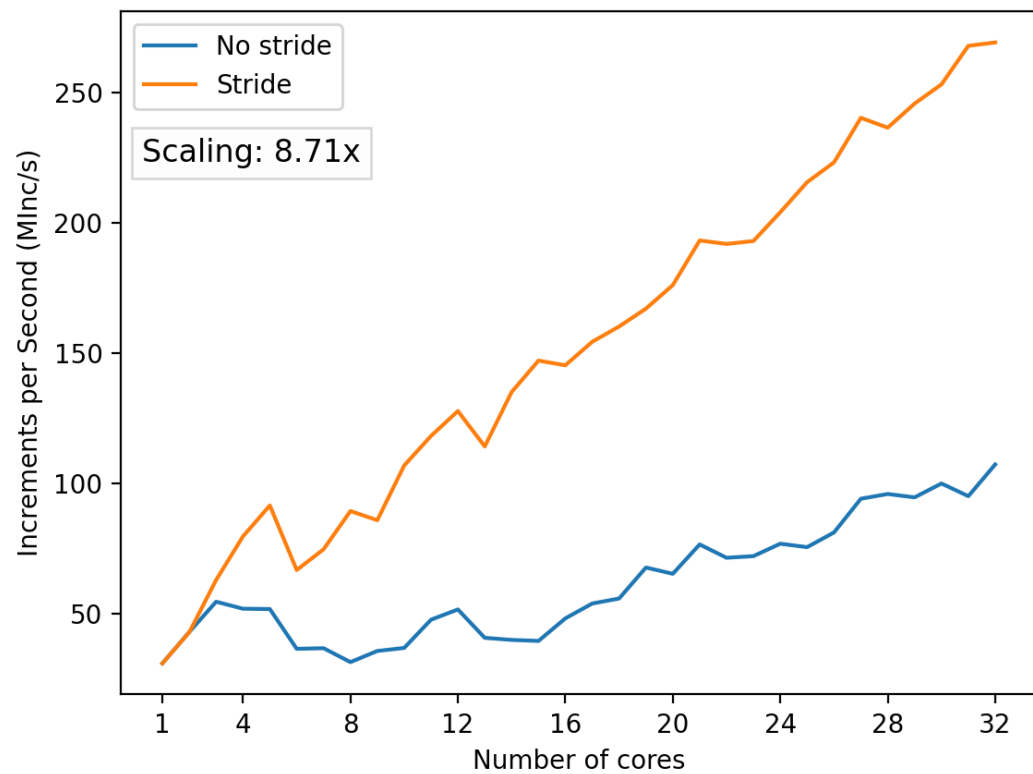


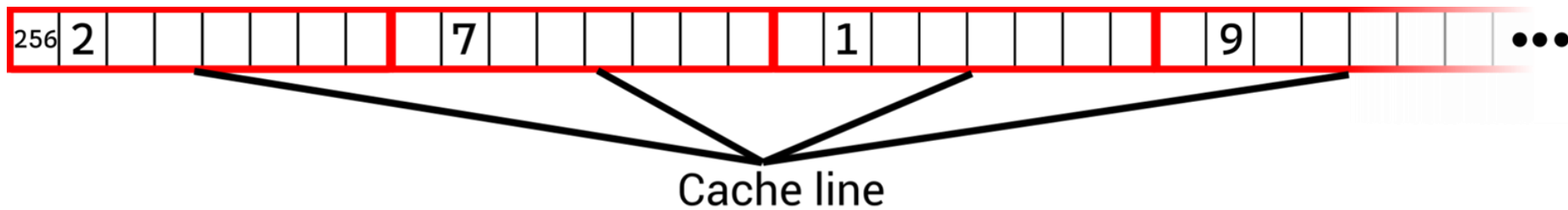


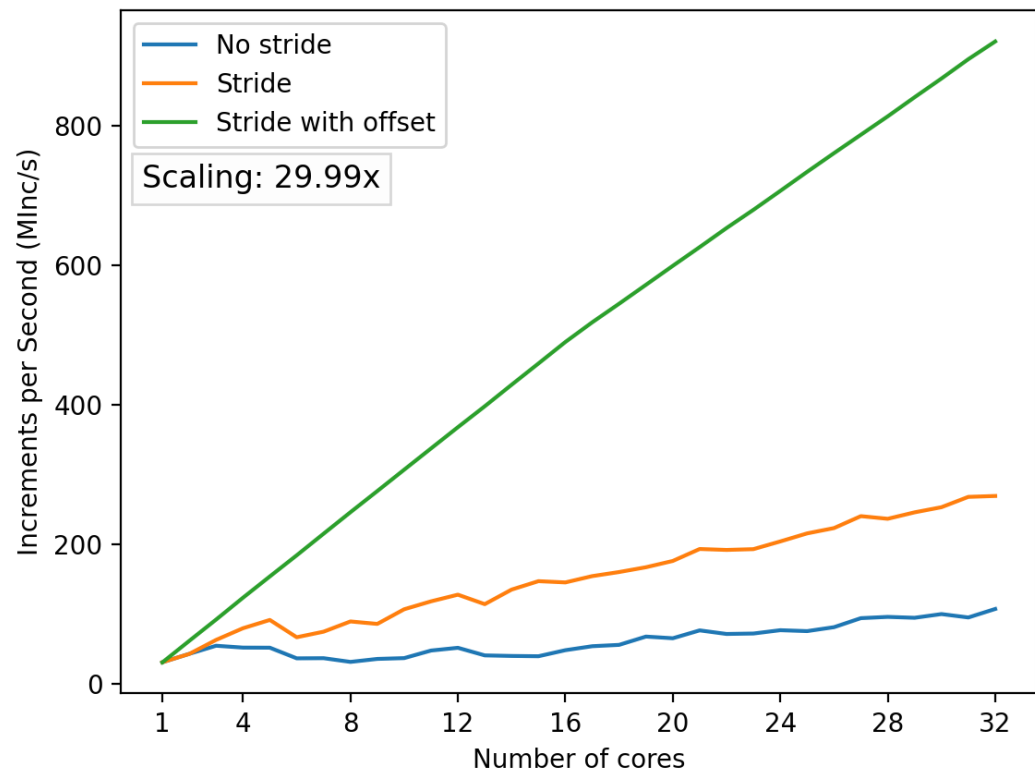
2	7	1	9	4	3	5	2	9	1	4	4	8	5	7	9	9	9	1	2	1	3	4	4	5	4	3	8	7	9	6	2
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Cache line





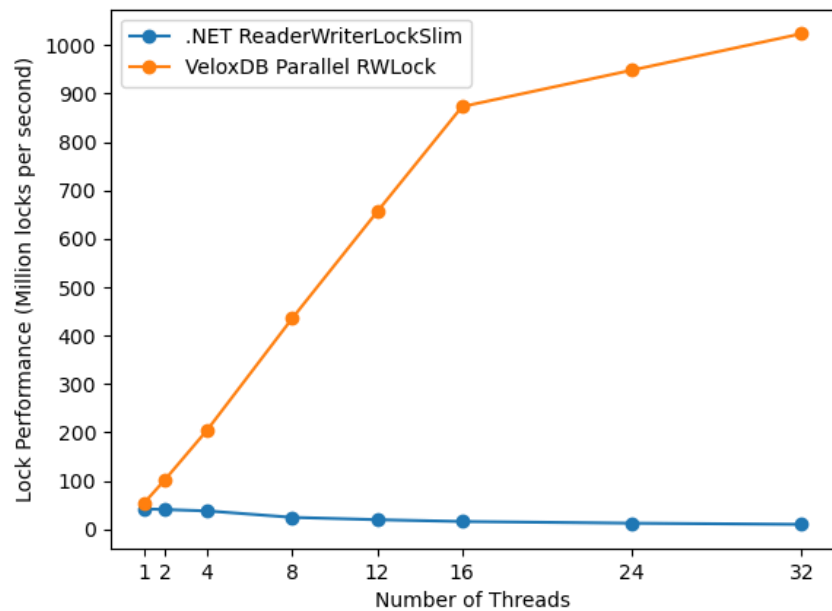




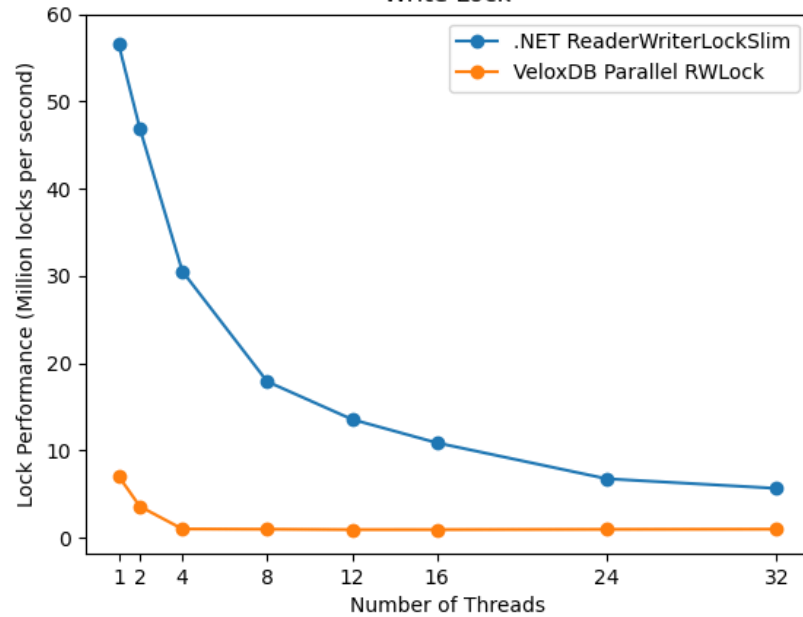


- Very fast read locking
- No contentions when reading
- Slower write
- Each CPU gets its own lock, on separate cache-line
- Read lock acquires only per-cpu lock
- Writes acquire all locks

Read Lock



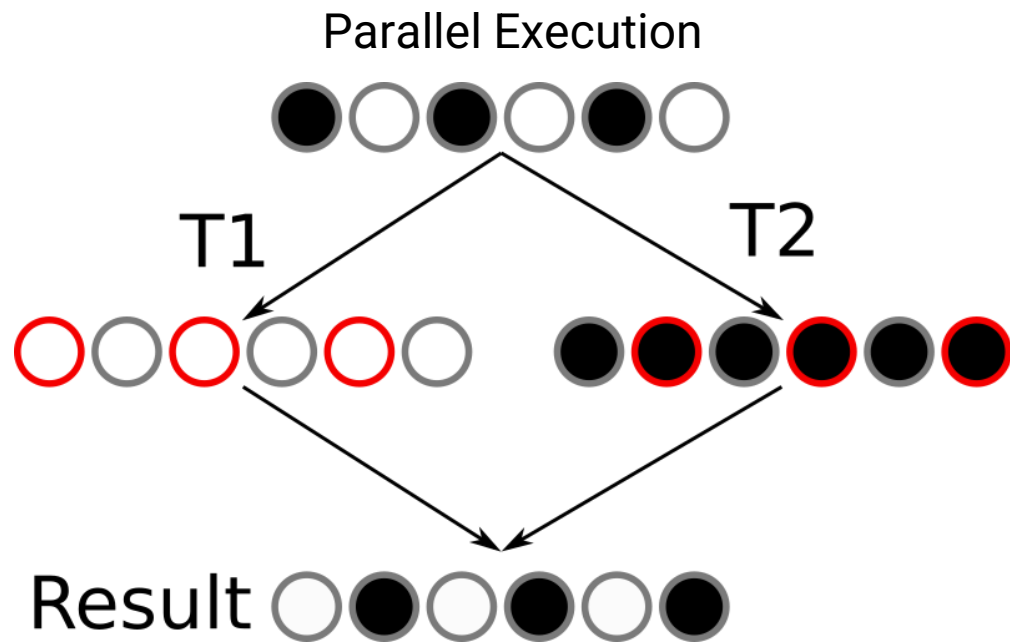
Write Lock



- Strict Serializability - Transactions view the database as if they are the only ones executing
- Transactions are executed in parallel
- Optimistic Concurrency Control
- Conflict - When transactions cannot be executed concurrently.
  - Write-Write conflict - Two transactions trying to modify the same object
  - Read-Write conflict - A transaction trying to write data that has been read by another transaction

- Table with black/white balls
- T1 - Scans the table and flips black to white
- T2 - Scans the table and flips white to black





- Multiversion Concurrency Control (MVCC)
- Objects must keep track of every transaction that reads them
- Lists are expensive
- Small 8-byte structure that can hold 3 transaction IDs and a count
- If more than 3 transactions read the same object, a dedicated list is created

- Better understanding of data lifecycle
  - Data stored in database is going to be long lived
  - Fewer distinct types
  - Easier to detect what data can be discarded
- Reduces pressure on the .NET GC



- Each table has its own memory manager
- Each memory manager is composed of smaller per-cpu memory managers
- Fixed width allocators
- Freed objects are linked in free object list for reuse

- It's important to verify that database behaves correctly
- Large body of tests (unit, component, end to end)
- Most useful - “stress” tests
  - Concurrent
  - Semantic Fuzzing
  - Verified using Reference database

- Challenging bugs (concurrent and memory corruption)
- Hard to reproduce
- We built our own time travel solution

1 reference

```
public void CreateObjectDiff(ClassObject* obj, ulong commonVersion, ChangesetWriter writer,
    IdSet partnerIds, List<long> otherIds, GenerateAlignDelegate alignDelegate)
{
    TTTrace.Write(TraceId, ClassDesc.Id, obj->id, obj->version, obj->IsDeleted, commonVersion);

    if (obj->IsDeleted)
        return;

    bool sbyMissing = partnerIds != null && !partnerIds.Contains(obj->id);
    if (obj->version > commonVersion || sbyMissing)
    {
        // If the standby deleted the object (due to split brain scenario) we need to pack an entire object which
        // will be achieved with the commonVersion being set to zero (this forces strings and blobs to be packed as well).
        TTTrace.Write(TraceId, ClassDesc.Id, obj->id, obj->version, commonVersion, sbyMissing);
        ulong version = sbyMissing ? 0 : commonVersion;
        alignDelegate(obj, stringStorage, blobStorage, writer, version);
        writer.LastValueWritten();
    }
    else
    {
        You, 14 months ago • Initial commit
        TTTrace.Write(TraceId, ClassDesc.Id, obj->id, obj->version, commonVersion);
        otherIds ??= new List<long>();
        otherIds.Add(obj->id);
    }
}
```

Visual Studio Code interface showing a C# file named `MemoryManager.cs` in the `VeloxDB` project. The code is in the `Free` method, which is annotated with `[MethodImpl(MethodImplOptions.AggressiveInlining)]`. The code includes logic for handling memory buffers and freeing memory.

```
190 uint originalSize = *((uint*)buffer);
191 if (originalSize > bufferSize[sizeIndex] - 8)
192     throw new CriticalDatabaseException();
193
194 uint originalSize2 = *((uint*)(buffer + 4 + originalSize));
195 if (originalSize != originalSize2)
196     throw new CriticalDatabaseException();
197 #endif
198
199 #if TEST_BUILD
200 byte* limit2 = (byte*)buffer + bufferSize[sizeIndex];
201 ITrace.Write((ulong)buffer, (ulong)limit2);
202 Utils.FillMemory(buffer, bufferSize[sizeIndex], 0xeb);
203 #endif
204
205 int procNum = ProcessorNumber.GetCore();
206 perCPUData[procNum]-->Free(buffer, sizeIndex, freeLists[sizeIndex]);
207 }
208
209 [MethodImpl(MethodImplOptions.AggressiveInlining)]
210 public void SafeFree(Action freeAction)
211 {
212     lock (disposeSync)
213     {
214         if (disposed)
215             return;
216
217         freeAction();
218     }
219 }
220
221 private static int FindFixedAllocator(int size)
```

The right sidebar shows the **Time Travel Control Panel** with the following sections:

- File:** testhost\_0.trd - 100.0%  
Clear
- Navigation:** Navigation icons (back, forward, etc.)  
Current Position: 166778  
Jump To
- Values:**

Threadid	14	
(ulong)buffer	2554986824192	
(ulong)limit2	2554986824576	

  
Copy To Clipboard
- Breakpoints:** Clear

At the bottom, the **Developer PowerShell** tab is active, showing "No issues found".

# Questions