



Fast Random Access Log (FRAL)

[MS Thesis Defense]

David Maaghul

Outline

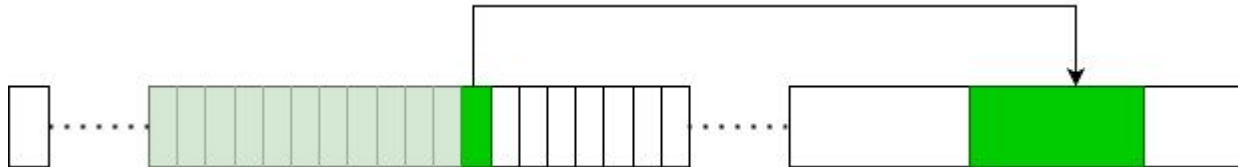
- Overview
- Core Architecture
- Performance Testing: C++ and Python
- Simple Sample Usage
- 3 Advanced Use Cases
- Demo
- Future Work

Overview

- Want to create data structure that allows multiple processes to read from and write to a shared array-indexed log as fast as possible
- Design Considerations:
 - Log is allocated over shared memory
 - Our current implementation uses memory-mapped files
 - Non-blocking writes
 - Reads are random access
- Implementation details:
 - Core engine implemented in C++
 - Python bindings provided and tested for convenient use (with less performance constraints)

Core Architecture

- The interface is very small and simple
- Three main (non-blocking) functions:
 - **allocate:** allocates n contiguous bytes of the memory space (a blob) and returns a pointer to the beginning of the allocation. This is similar to `malloc` in C.
 - **append:** appends the address of a memory allocation (blob) to an indexed log, updates the size of the log and returns the index of the appended blob.
 - **load:** provided a log-entry index, returns the address of the corresponding memory allocation
- There is no **store** function; appends are permanent but blobs can change
- Virtual Memory = different addressing for different processes
 - Store offset from beginning of shared-memory allocation to array
 - Offset corresponds to array index = random access



Non-Blocking Allocate

- Our first atomic variable:
 - h : the location of the next free memory address
- Algorithm uses fetch-add primitive
 - Offset of next free memory address is an atomic variable
- Caveat:
 - Process failure (e.g., the OS kills the process) after fetch-add step (line 5) results in wasted space.
 - Still structurally consistent

Algorithm 1 Non-Blocking Allocate

```
1:  $s$ : allocation size (bytes)
2:  $S$ : size of memory space (bytes)
3:  $h$ : offset of next free memory address
4: function ALLOCATE( $s$ )
5:    $c = \text{FetchAdd}(h, s)$ 
6:   if  $c + s > S$  then
7:     return null
8:   end if
9:   return  $c$ 
10: end function
```

Non-Blocking Append

- Motivation for approach:
 - How we expect log tiling to be implemented is critical in our design of our append algorithm

Algorithm 2 Tail Log

```
1:  $S$ : max size of log
2:  $fral$ : the random access log
3: function READ( $fral$ )
4:    $i \leftarrow 0$ 
5:   for  $i < S$  do
6:      $blob = \text{load}(fral, i)$ 
7:     if  $blob$  is not null then
8:       ProcessBlob( $blob$ )
9:        $i \leftarrow i + 1$ 
10:    end if
11:  end for
12: end function
```

Non-Blocking Append (Continued)

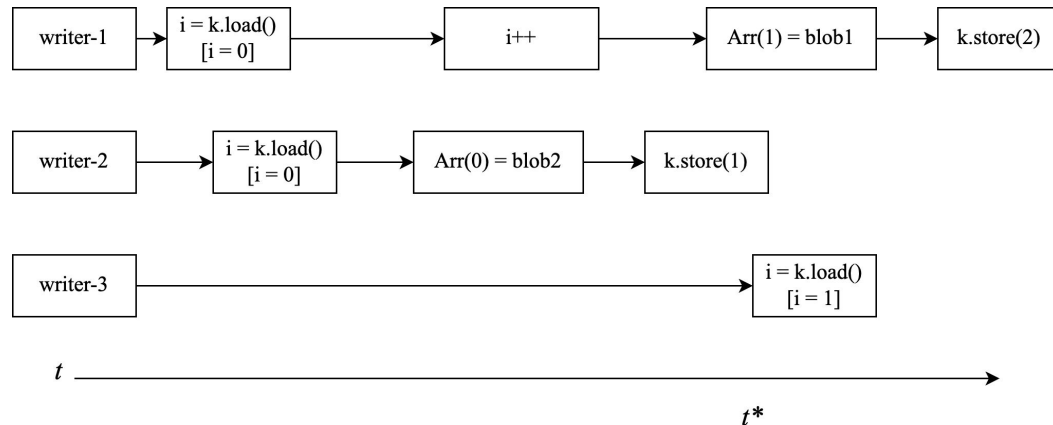
- Two new atomic variables:
 - $Arr[i]$: The i^{th} entry of the array of offsets
 - k : The next available array index
- Use compare-and-swap primitive:
 - Only store offset value if corresponding entry is available
- A process can fail at any point during the append
 - Array entries are never wasted
 - Iterating over log will never yield infinite loop (previous slide)
- In the case of n writers in parallel:
 - A given writer has a worst case operation count of $O(n)$ and the total operation count between the n writers is $O(n^2)$
 - We could achieve this with $O(1)$ and $O(n)$, respectively, but would then lose the promise of not wasting array entries
 - Asymptotic gain didn't materialize in performance testing

Algorithm 3 Non-Blocking Append

```
1: blob: allocation to be appended
2: S: max size of log
3: k: the next available index
4: a: address of start of shared memory space
5: Arr: the log of offsets
6: function APPEND(blob)
7:    $z \leftarrow BlobOffset(a, blob)$ 
8:    $i \leftarrow AtomicLoad(k)$ 
9:   for  $i < S$  do
10:    if CompareAndSwap( $Arr[i], 0, z$ ) then
11:      AtomicStore( $k, i + 1$ )
12:      return  $i$ 
13:    end if
14:     $i \leftarrow i + 1$ 
15:  end for
16:  return -1
17: end function
```

Non-Blocking Append (Continued)

- Problem: can't atomically update two variables at the same time
- The next available index is estimated (relaxed):
 - Example with 3 writers:
 - writer-3 loads $k = 1$ at time t^* , but entry 1 has already been appended to by writer-1
 - writer-3 needs two compare-and-swap operations to realize and write to the next available index ($k = 2$)



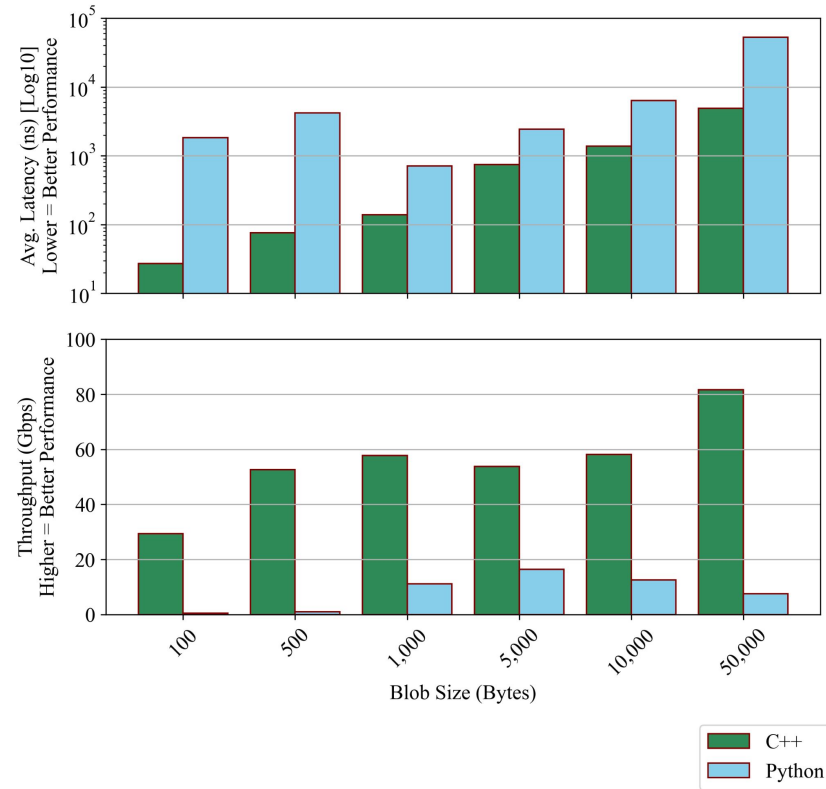
Performance Testing Results

- *Two tests of core architecture performance:*

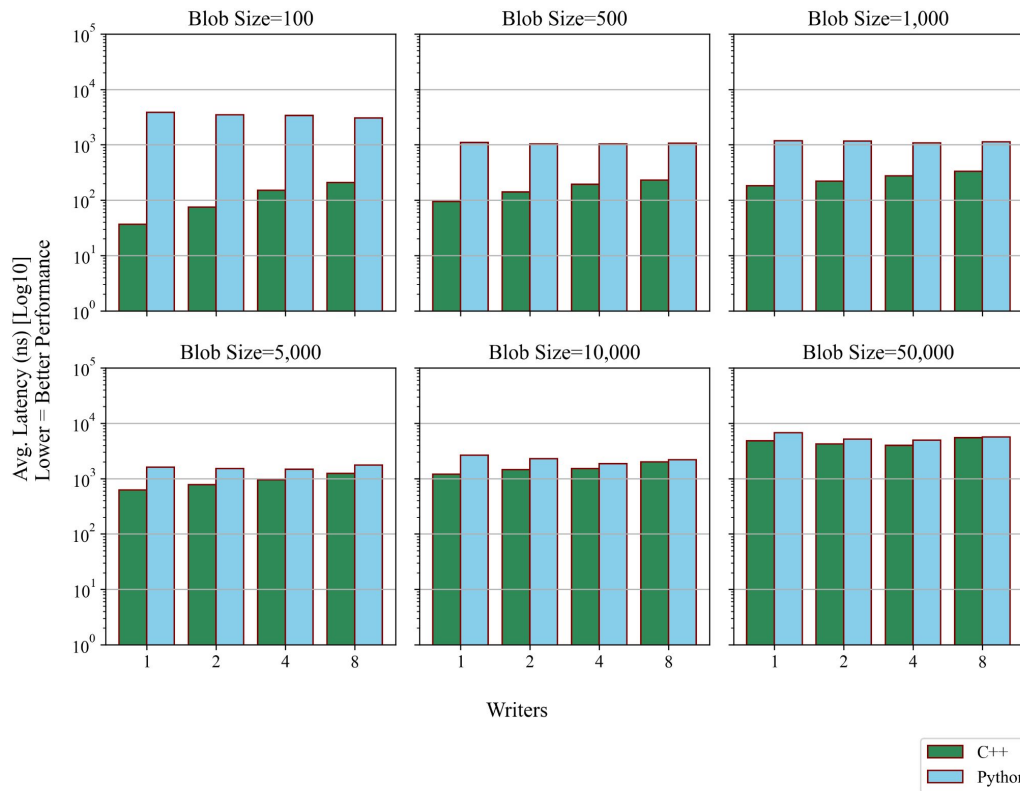
- Look to measure throughput and latency in various ways
- *Write Test:*
 - Measures the speed and efficiency of writing N GB of data to the log in n -byte blobs
 - Assesses performance of high volume of write operations by one writer
 - The timer starts before the first write and ends after the last; latency is an average
- *Producer-Consumer Test:*
 - Measures the speed and efficiency of simultaneously writing and reading N GB of data to the log in n -bytes blobs with k writers
 - Assesses performance of concurrent read and write operations from multiple sources
 - The timer starts before the first write and ends after the last read; latency, again, is an average

*Note that all testing was done with an 8-core Apple M2 machine with 8GB of unified RAM.

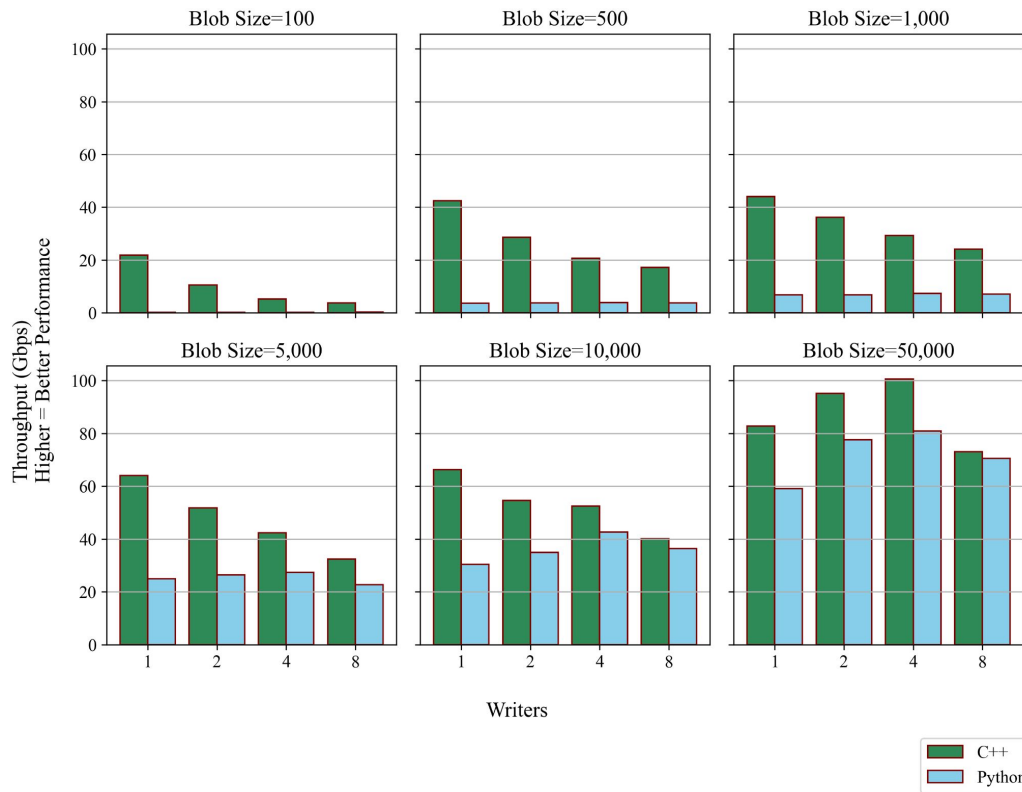
1GB Write Test Results (Latency + Throughput)



1GB Producer-Consumer Test Results (Latency)



1GB Producer-Consumer Test Results (Throughput)



Simple Usage (C++)

```
1 //Process 1
2 auto ralA = fral::FRAL("test.bin");
3 const char *TEST_STR = "TEST";
4
5 auto blob = (char *) ralA.allocate(strlen(TEST_STR));
6 strcpy(blob, TEST_STR);
7 printf("%s\n", blob);
8
9 ralA.append(blob);
```

```
1 //Process 2
2 auto ralB = fral::FRAL("test.bin");
3
4 while(true){
5     auto blob2 = (char *) ralB[0]; //load operation
6     if(blob2){
7         break;
8     }
9 }
10
11 printf("%s\n", blob2);
12 blob2[0] = 'D';
13
14 auto blob3 = (char *) ralB[0]; //load operation
15 printf("%s\n", blob3);
```

```
TEST
TEST
DEST
```

Simple Usage (Python)

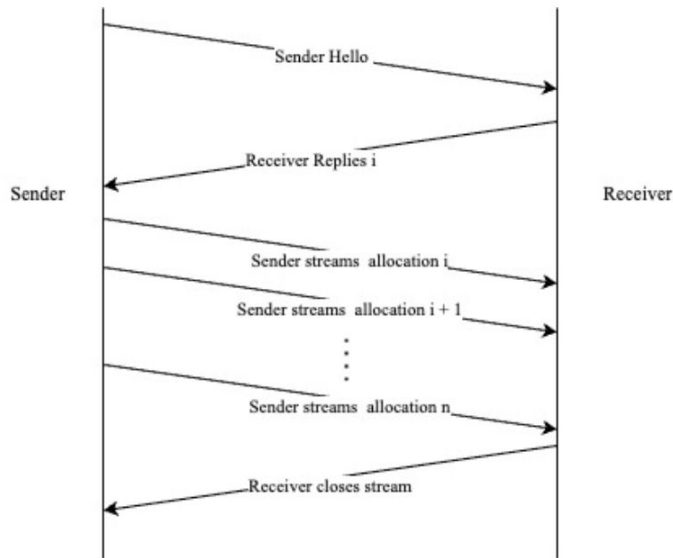
```
1 # Process 1
2 ral_A = FRAL("test.bin")
3 test_bytes = "TEST".encode()
4
5 test_blob = ral_A.allocate(len(test_bytes))
6 test_blob[:len(test_bytes)] = test_bytes
7
8 print(bytes(test_blob).decode())
9 ral_A.append(test_blob)
```

```
1 # Process 2
2 ral_B = FRAL("test.bin")
3
4 while True:
5     test_blob2 = ral_B[0] # load operation
6     if test_blob2:
7         break
8
9 print(bytes(test_blob2).decode())
10 test_blob2[0:1] = 'D'.encode()
11
12 test_blob3 = ral_B[0] # load operation
13 print(bytes(test_blob3).decode())
```

```
TEST
TEST
DEST
```

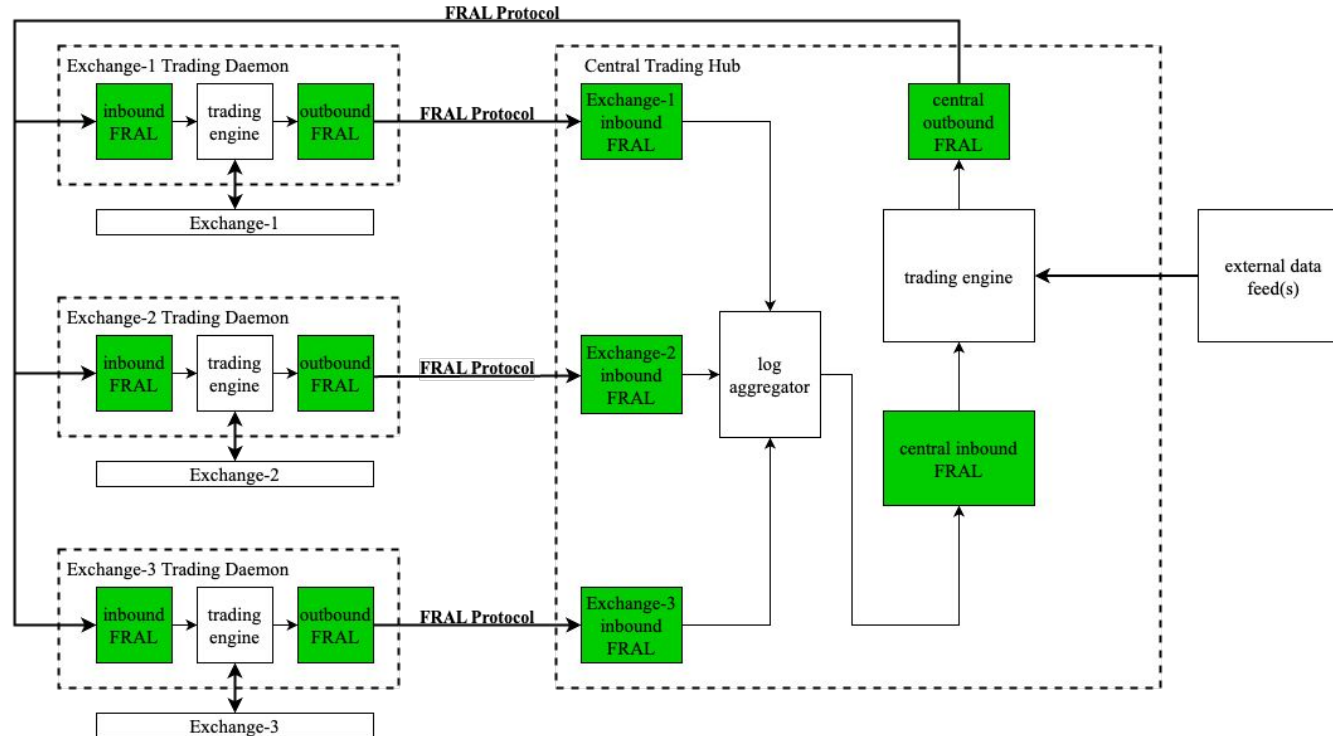
Syncing Over The Network

- Design Considerations:
 - Non-intrusive, i.e., don't need to alter the underlying log structure or add overhead
 - Order of entries is maintained over network
 - Robust to failure for both client and server; mirroring log is identical even with intermittent failures
- Other details:
 - Initial design Implemented over gRPC and Protobuf
 - Future work: non-intrusive network sync tooling for many-to-one log mappings
- Network Test:
 - Similar to write test, but over the network:
 - Performance degrades (likely to client side buffering by gRPC)
 - Write and Read N entries of various sizes

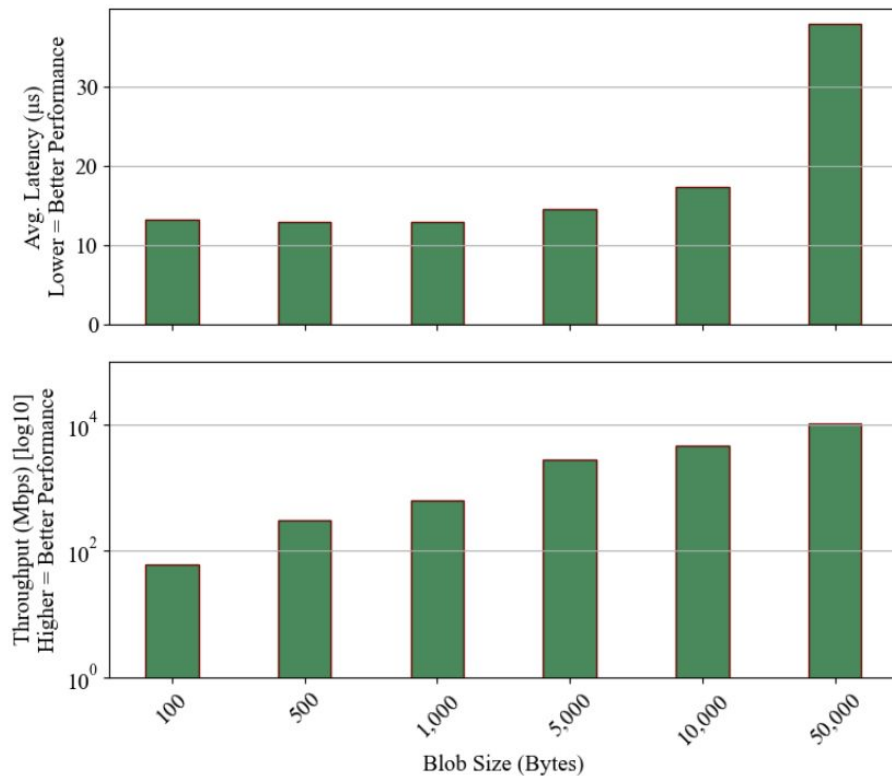


Syncing Over The Network (Use Case)

- High Frequency Trading:* An entity trades on multiple exchanges in different geographic locations, while needing to keep track of net holdings in a central place



Syncing Over The Network (10k Entry Network Test)



Language-Agnostic Demo

- Language-Agnostic Framework
 - Shared memory opens door for language-agnostic data structures
 - Good for systems with microservices that vary in performance constraints
 - Use C++ for computationally-intensive tasks
 - Capitalize on rapid prototyping in development with Python
 - Example: HFT back-office tasks in Python, trading in C++
 - Demo (video):
 - Read integer from one language, increment by one and write
 - Python initiates process by writing and sending (appending 1)
 - Standard output for writing 5 entries:

```
Python sending 1 to C++!  
C++ received 1 from Python, sending 2!  
Python received 2 from C++, sending 3!  
C++ received 3 from Python, sending 4!  
Python received 4 from C++, sending 5!  
C++ received 5 from Python, done!
```

Future Work

- Processor affinity functionality
 - Assigning one processor to handle page faults (prefetching)
- Networking performance improvements
- Reliably syncing multiple logs to one over the network
 - Non-intrusive, i.e., don't add overhead to allocations or underlying log structure as a whole (for accounting)
- High-performance local log aggregator
- Other general performance improvements
 - How does moving away from memory-mapped files affect performance?
- Automatic extension
 - How can we efficiently extend the size of the structure when we run out of space?