

Node Replicated Kernel?

Node Replicated Kernel (NRK) is a research prototype OS Research, but is now being developed collaboratively by [a](#) and [academia](#). It is intended as a basis to explore ideas and systems for hardware of the future. NRK is written from scratch (assembly), and it runs on x86 platforms.

How is NRK different from Linux, Windows, etc?

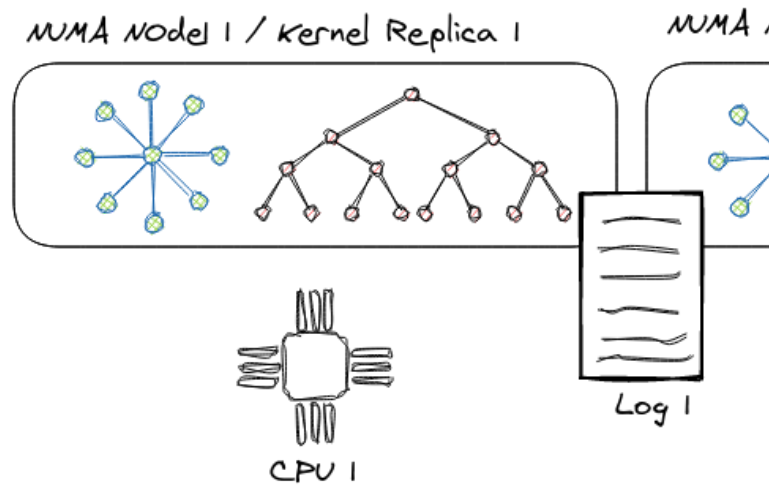
At present, NRK lacks many of the features of an operating system anyone other than systems researchers taking measurements. For example, there is currently no GUI or Shell, and only very basic networking. It's probably easier to compare NRK in its current form to a fully-featured OS.

From an architectural point of view, the NRK kernel is a new [multi-kernel OS](#):

A multikernel operating system treats a multi-core machine as if it were a distributed system. It does not assume a shared memory, and implements inter-process communications as message passing.

Unfortunately, such a model also brings unnecessary complexity. [multi-kernel \(Barrelfish\)](#) relied on per-core communication (2PC, 1PC) to achieve agreement, replication and sharing.

[We overcome this complexity in NRK by using logs](#): The kernel uses replicated data structures which are automatically replicated. Our logs make sure the replicas are always synchronized. Our approach bears resemblance to state machine replication in distributed systems. We transform replicated data-structures into linearizable, concurrent structures. We provide details on scenarios where this approach can heavily outperform lock-based or lock-free data-structures.



Schematic view of the NRK kernel. Core kernel data-system with a log.

In user-space, NRK runs applications in [ring 3](#): Each application is an isolated process/container/lightweight VM with little operating system processes. To run existing, well-known applications like [rumpkernel](#), which instantiates a user-space libraries (libc, libpthread etc.) within the process to provide decent support for POSIX.

Finally, NRK is written from scratch in Rust: We take advantage of a rich type-system for OS implementation, to gain better guarantees at compile time, while not impacting performance.

Kernel Architecture

The NRK kernel is a small, light-weight (multi-)kernel that [virtual memory](#), a [coarse-grained scheduler](#), as well as an

One key feature of the kernel is how it scales to many cor data-structure replication with operation logging. We exp for this in the [Node Replication](#) and [Concurrent Node Rep](#)

Node Replication (NR)

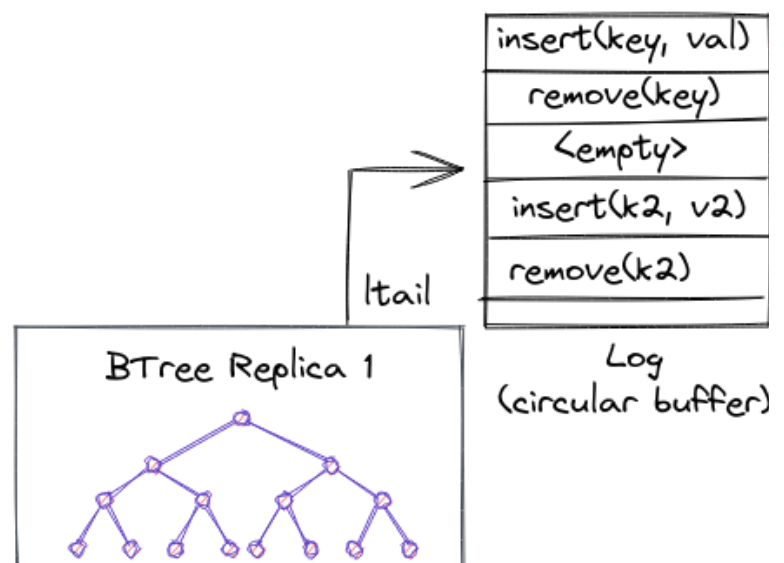
NR creates linearizable NUMA-aware concurrent data structures. NR replicates the sequential data structure operation log to maintain consistency between the replicas; concurrency using a readers-writer lock and from write conflicts using *flat combining*. In a nutshell, flat combining batches operations executed by a single thread (*the combiner*). This thread also updates the log; other replicas read the log to update their internal data structures.

Next, we explain in more detail the three main techniques: flat combining, logs, scalable reader-writer locks and flat combining.

The operation log

The log uses a circular buffer to represent the abstract state of the data structure. Each entry in the log represents a mutating operation, an index to the last operation added to the log is the *log tail*.

A replica can lazily consume the log by maintaining a pointer to how far from the log the replica has been updated. The log is so that entries can be garbage collected and reused. NR can be applied on each replica. This means at least one replica must be executing operations on the data structure, otherwise other replicas would not be able to make any more progress.

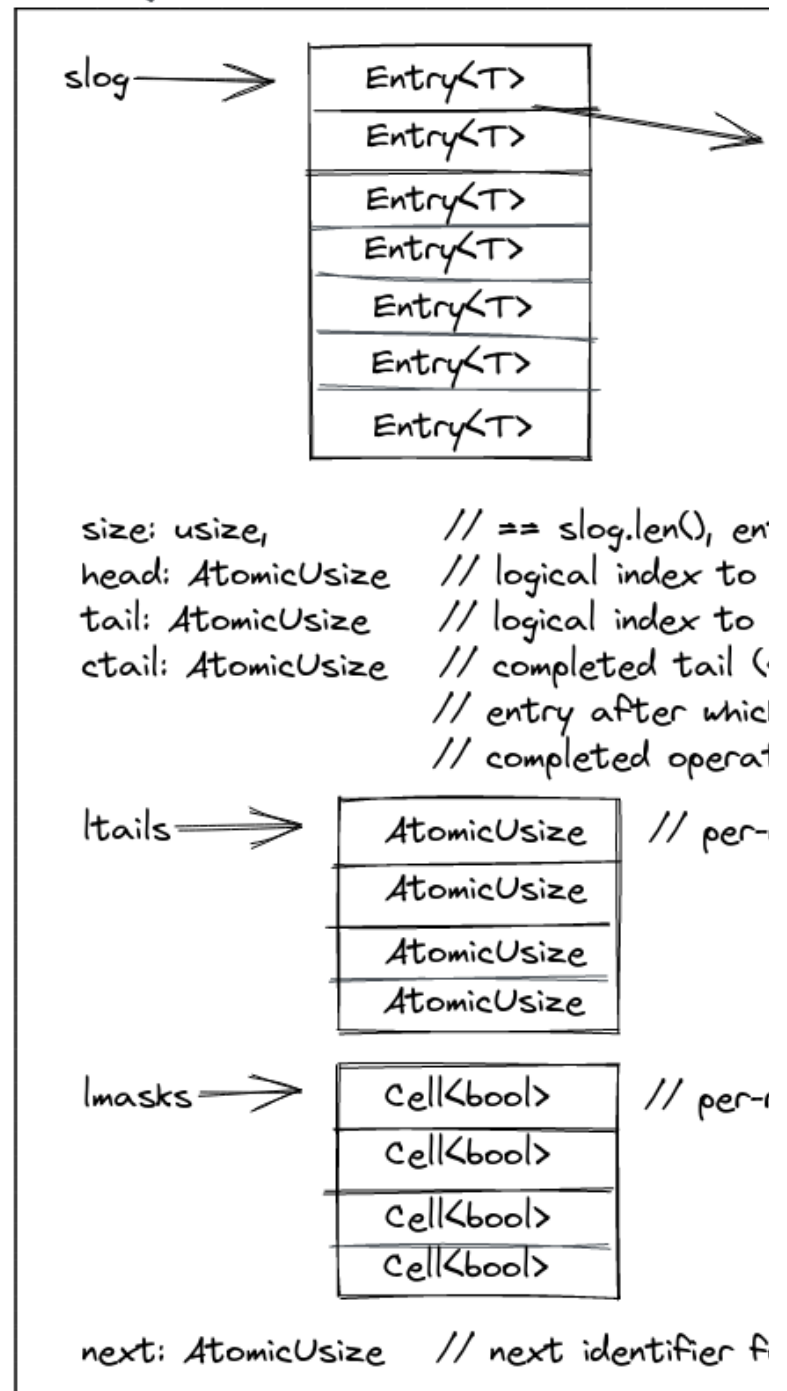


Shows a log containing mutable operations for a 2x contains a local tail (ltail1 and ltail2) pointer into the oldest entry that is still outstanding (i.e., needs to be replica before it can be overwritten).

Next, we're going to explain the log in more detail for the directly working on this, likely these subsections are prett

Log memory layout

NR Log<T>



The log structure and its members as they appear in

The `log` contains an array or slice of entries, `slog`, where the replicas. It also maintains a global `head` and `tail` which indicate the current operations in the log that still need to be processed or have not been garbage collected by the log. The `head` still needs to be processed, and the `tail` to the newest, which tracks per-replica progress. At any point in time, an subregion given by `head..tail`. `lmask` are generation bit replica wraps around in the circular buffer.

`next` is used internally by the library to hand out a register with the log (this is done only in the beginning, during set consuming from/appending to the log).

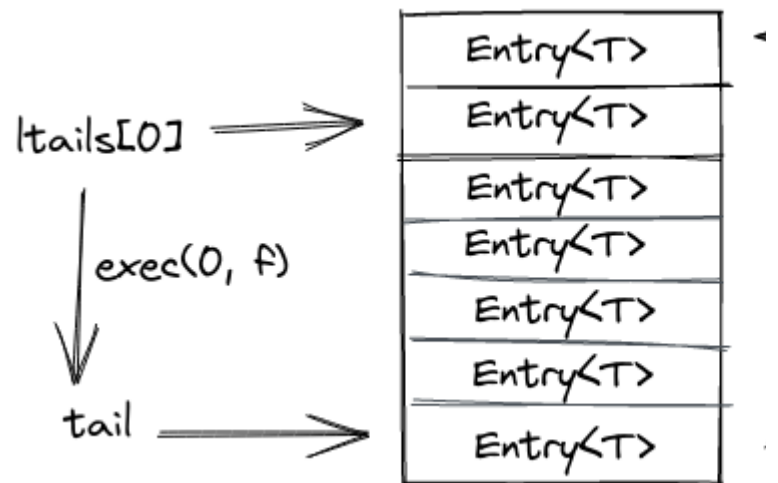
The `ctail` is an optimization for reads and tracks the max log entry (`< tail`) after which there are no completed op

Consuming from the log

The `exec method` is used by replicas to execute any outstanding entries. It takes a closure that is executed for each outstanding entry that calls `exec`.

Approximately, the method will do two things:

1. Invoke the closure for every entry in the log which is afterwards sets `ltail = tail`.
2. Finally, it may update the `ctail` if `ltail > ctail`



Indicates what happens to `ltail[0]` and `ctail` when replica 0 is updated to be equal to `tail` and the `ctail` (currently is `ltail[1]`) will point to a new max which is the same as `tail`. It also applies all the entries between `ltail[0]` and `tail` to replica 0.

Appending to the log

The `append operation` is invoked by replicas to insert multiple new operations coming from the replica and then apply some outstanding operations. This is necessary as t

insert more operations. Then, we'd first have to update our replicas to make progress).

Approximately, the method will end up doing these steps

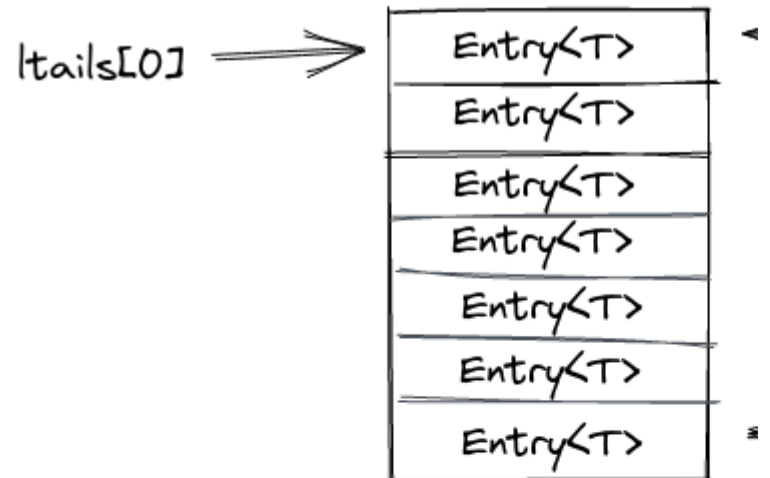
1. If there isn't enough space to insert our new batch of entries from the log until we have enough space available
2. Do a compare exchange on the log `tail` to reserve space
3. Insert entries into the log
4. See if we can collect some garbage aka old log entries

Garbage collecting the log

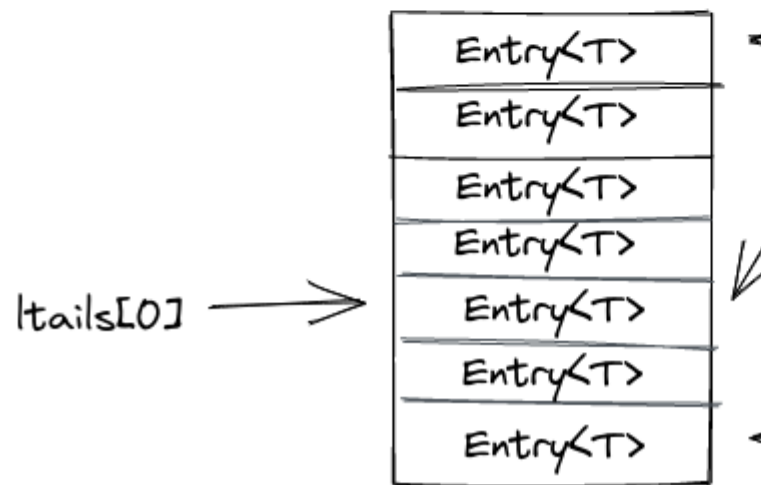
`advance_head` has the unfortunate job to collect the log garbage little as periodically advancing the head pointer.

For that it computes the minimum of all `ltails`. If this minimum head, it waits (and calls `exec` to try and advance its local head). If the head, it will update the head to that new minimum. It

advance_head can't advance:



advance_head can advance:



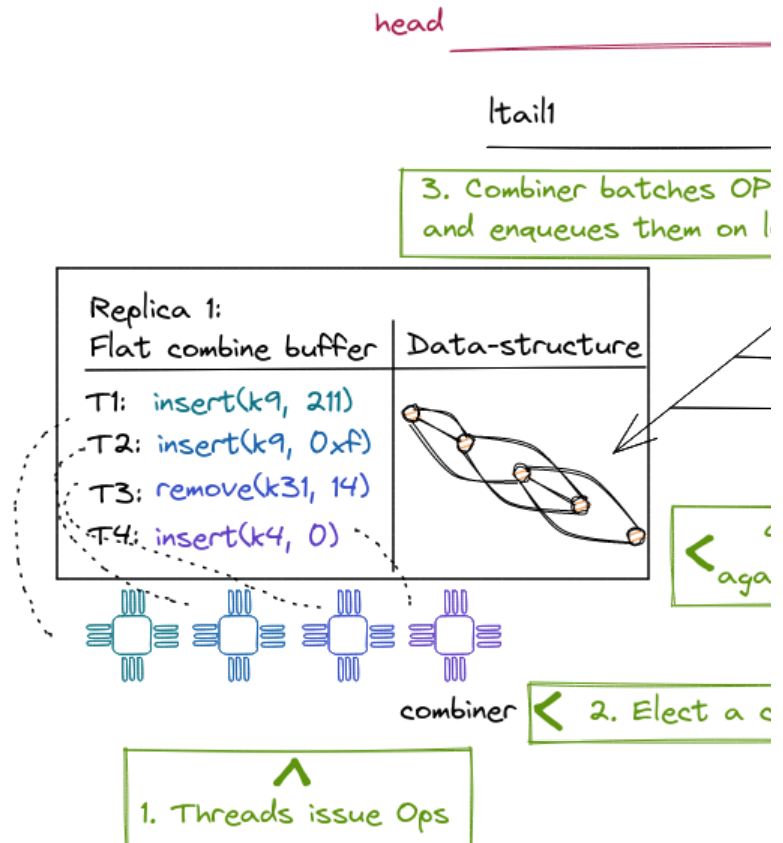
Shows the two possible scenarios in advance head: replica 0 points to the current head, so we can't advance to wait. In the bottom diagram, both ltail[0] and ltail head, so we don't need the entries between head..lt advanced.

Flat combining

NR uses [flat combining](#) to allow threads running on the server resulting in better cache locality both from flat combining local to the node's last-level cache.

1. The combiner can batch and combine multiple operations at a lower cost than executing each operation individually. For example, a priority queue can be done with a single atomic instruction for each operation.
2. A single thread accesses the data structure for multiple operations, ensuring cache locality for the whole batch of operations.

NR also benefits from combining operations placed on the single atomic operation for the whole batch instead of on



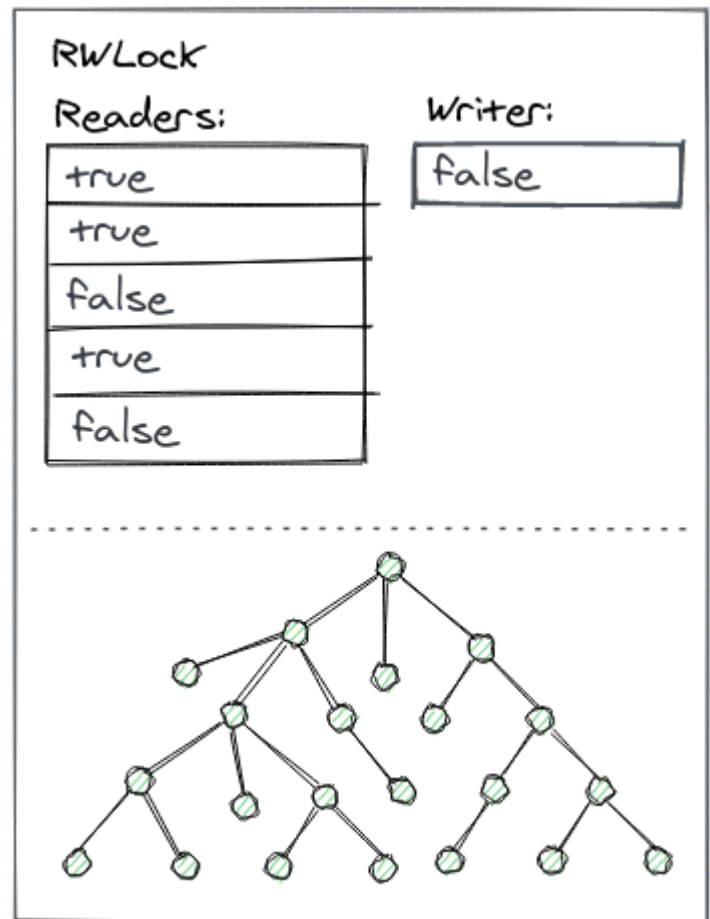
Flat combining multiple operations from threads on a thread is elected who will aggregate outstanding operations, put them on the log in one batch (3) and then execute (4), after executing any previously outstanding operations.

The optimized readers-writer lock

NR uses a writer-preference variant of the distributed RW synchronization of the combiner and reader threads when

Each reader acquires a local per-thread lock; the writer has to acquire the writer lock. A reader first acquires its local lock if there is no writer; then it checks the writer lock (and releases the reader lock if it notices the writer is active (writer preference). The combiner also acquires the writer lock. Thus, we give highest priority to the writer, because we don't want to let readers read stale data. This lock allows readers to read a local replica while the combiner performs operations to the log, increasing parallelism.

If there is no combiner, a reader might have to acquire the writer lock from the log to avoid a stale read, but this situation is rare.



The RWLock used to protect the data-structure in every cache-line sized atomic booleans for reader threads and one for the writer thread. It allows concurrent readers but only a single writer. Readers also have to make sure to advance the replication log.

A concurrent non-mutating operation

A non-mutation operation (read) can execute on its thread entry, but it needs to ensure that the replica is not stale. It reads the current tail when the operation begins, and waits until a combine operation updates the observed tail, or acquires the lock and updates the replica.

If there is no combiner, the read thread will acquire the write lock.

A concurrent mutating operation (update)

A thread *T* executing a mutating operation (update) needs to acquire the lock. If thread *T* executing this operation fails to acquire it, another thread may update the replica. Thread *T* then spin-waits to receive its operation's result from the existing thread.

However, *T* must still periodically attempt to become the combiner to complete its operation.

If *T* acquires the lock, it becomes *the combiner* and executes its operation. It then allows concurrent threads on the same replica. To do so, *T* appends its operation with a single atomic operation and ensures that the replica is up-to-date with all operations and returning their results to the waiting threads.

Source and code example

[Node-replication \(NR\)](#) is released as a stand-alone library.

To give an idea, here is an example that transforms an `HashMap` (from the Rust standard library) into a concurrent, replicated hash-table.

```

//! A minimal example that implements a replicated
use std::collections::HashMap;
use node_replication::Dispatch;
use node_replication::Log;
use node_replication::Replica;

/// The node-replicated hashmap uses a std hashmap
#[derive(Default)]
struct NrHashMap {
    storage: HashMap<u64, u64>,
}

/// We support mutable put operation on the hashmap
#[derive(Clone, Debug, PartialEq)]
enum Modify {
    Put(u64, u64),
}

/// We support an immutable read operation to look
#[derive(Clone, Debug, PartialEq)]
enum Access {
    Get(u64),
}

/// The Dispatch trait executes `ReadOperation` (
/// and `WriteOperation` (our `Modify` enum) again
/// data-structure.
impl Dispatch for NrHashMap {
    type ReadOperation = Access;
    type WriteOperation = Modify;
    type Response = Option<u64>;

    /// The `dispatch` function applies the immutable
    fn dispatch(&self, op: Self::ReadOperation) ->
        match op {
            Access::Get(key) => self.storage.get(&key)
        }
    }

    /// The `dispatch_mut` function applies the mutable
    fn dispatch_mut(&mut self, op: Self::WriteOperation) ->
        match op {
            Modify::Put(key, value) => self.storage.put(key, value)
        }
    }
}

```

As we can see we need to define two operation enums (`Modify` and `Access`) to end up on the log, and `Access` for immutable/read operations. We then implement the `Dispatch` trait from NR for our newly defined `NrHashMap` responsible to route the operations defined in `Access` and `Modify` to the appropriate

structure. [The full example](#), including how to create replic repository.

Node Replication (NR)

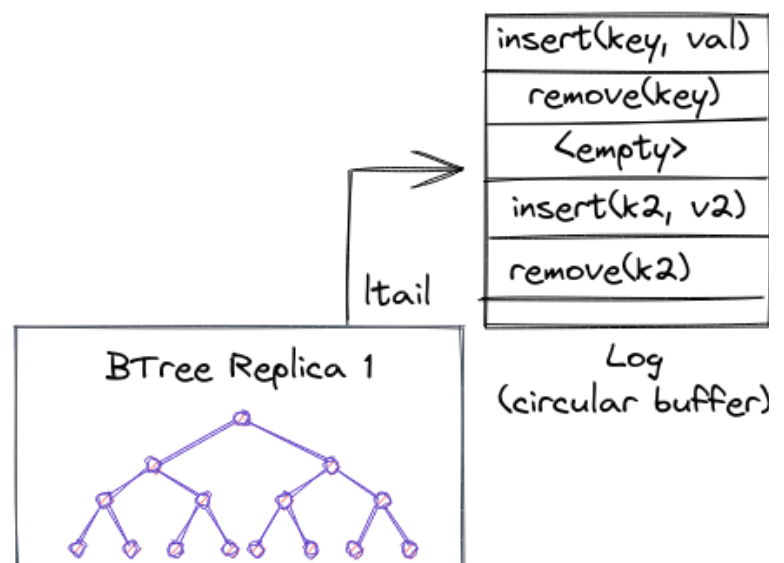
NR creates linearizable NUMA-aware concurrent data structures. NR replicates the sequential data structure operation log to maintain consistency between the replicas; concurrency using a readers-writer lock and from write conflicts using *flat combining*. In a nutshell, flat combining batches operations executed by a single thread (*the combiner*). This thread also updates the log; other replicas read the log to update their internal data structures.

Next, we explain in more detail the three main techniques: logs, scalable reader-writer locks and flat combining.

The operation log

The log uses a circular buffer to represent the abstract state of the data structure. Each entry in the log represents a mutating operation, an *log tail* gives the index to the last operation added to the log.

A replica can lazily consume the log by maintaining a pointer to how far from the log the replica has been updated. The log is so that entries can be garbage collected and reused. NR can have not been applied on each replica. This means at least one replica must be executing operations on the data structure, otherwise other replicas would not be able to make any more progress.

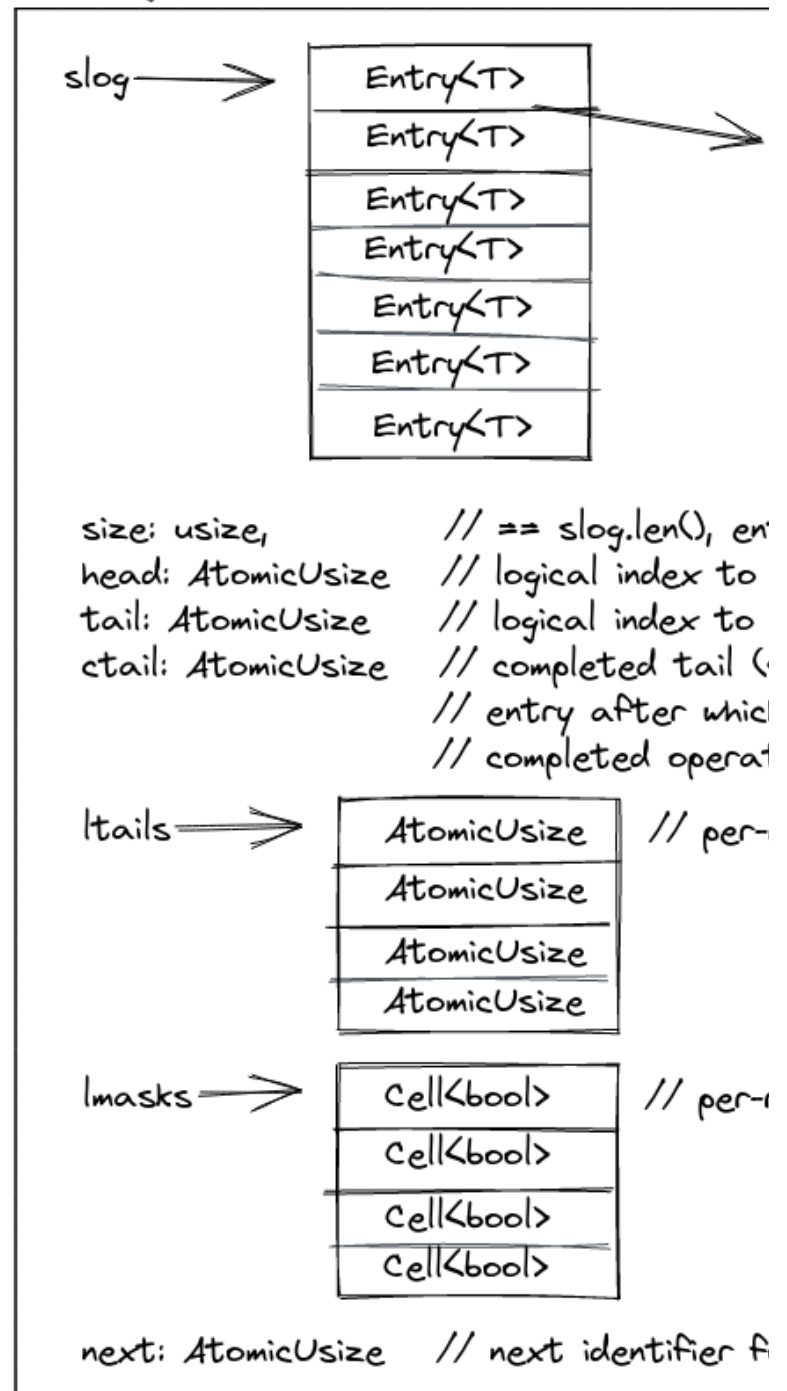


Shows a log containing mutable operations for a 2x contains a local tail (ltail1 and ltail2) pointer into the oldest entry that is still outstanding (i.e., needs to be replica before it can be overwritten).

Next, we're going to explain the log in more detail for the directly working on this, likely these subsections are prett

Log memory layout

NR Log<T>



The log structure and it's members as they appear in

The `log` contains an array or slice of entries, `slog`, where the replicas. It also maintains a global `head` and `tail` which indicate the current operations in the log that still need to be processed or have not been garbage collected by the log. The `head` tracks the oldest entry still needs to be processed, and the `tail` to the newest. `ltails` which tracks per-replica progress. At any point in time, an active subregion given by `head..tail`. `lmask` are generation bitmasks which replica wraps around in the circular buffer.

`next` is used internally by the library to hand out a register with the log (this is done only in the beginning, during set consuming from/appending to the log).

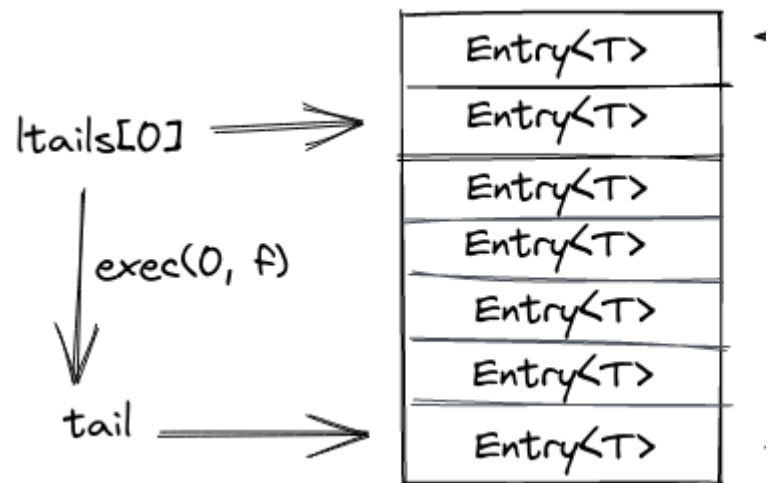
The `ctail` is an optimization for reads and tracks the max log entry (`< tail`) after which there are no completed op

Consuming from the log

The `exec method` is used by replicas to execute any outstanding entries. It takes a closure that is executed for each outstanding entry that calls `exec`.

Approximately, the method will do two things:

1. Invoke the closure for every entry in the log which is afterwards sets `ltail = tail`.
2. Finally, it may update the `ctail` if `ltail > ctail`



Indicates what happens to `ltail[0]` and `ctail` when replica 0 is updated to be equal to `tail` and the `ctail` (currently is `ltail[1]`) will point to a new max which is the same as `tail`. It also applies all the entries between `ltail[0]` and `tail` to replica 0.

Appending to the log

The `append operation` is invoked by replicas to insert multiple new operations coming from the replica and then apply some outstanding operations. This is necessary as t

insert more operations. Then, we'd first have to update our replicas to make progress).

Approximately, the method will end up doing these steps

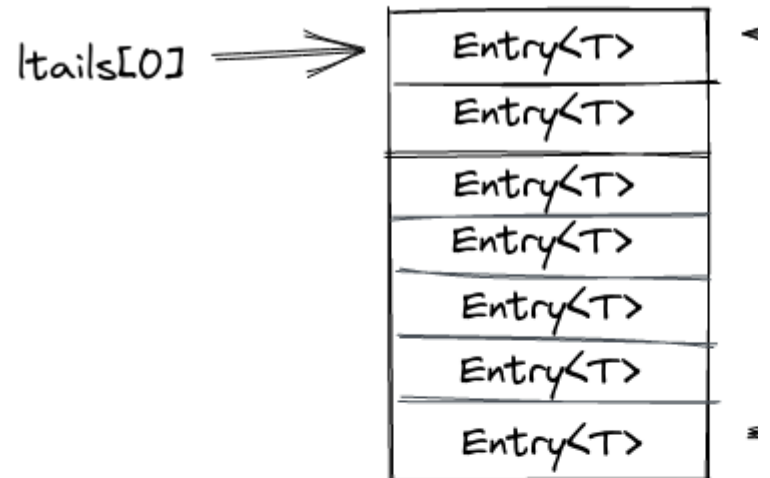
1. If there isn't enough space to insert our new batch of entries from the log until we have enough space available
2. Do a compare exchange on the log `tail` to reserve space
3. Insert entries into the log
4. See if we can collect some garbage aka old log entries

Garbage collecting the log

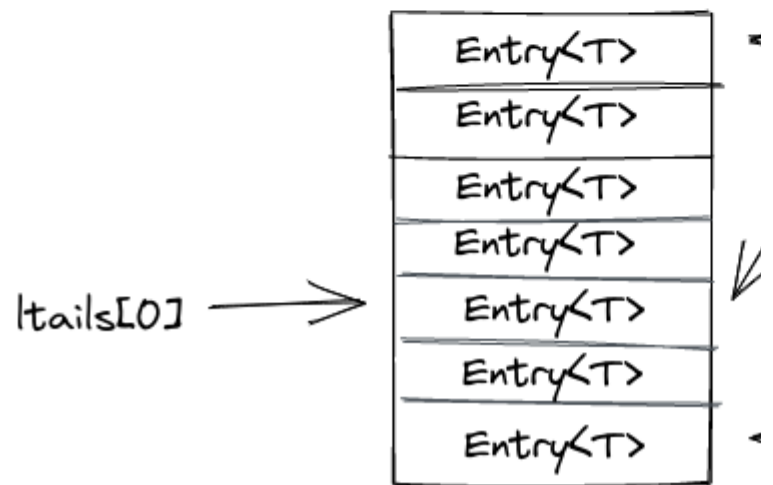
`advance_head` has the unfortunate job to collect the log entries as little as periodically advancing the head pointer.

For that it computes the minimum of all `l_tails`. If this minimum is ahead of the head, it waits (and calls `exec` to try and advance its local head). If the head, it will update the head to that new minimum. It

advance_head can't advance:



advance_head can advance:



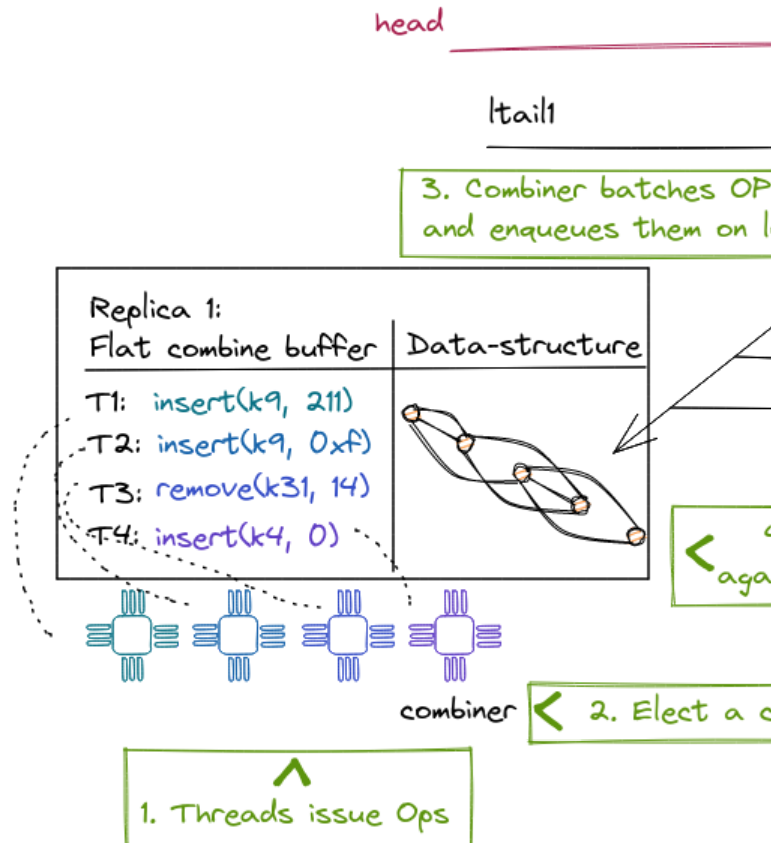
Shows the two possible scenarios in advance head: replica 0 points to the current head, so we can't advance to wait. In the bottom diagram, both ltail[0] and ltail head, so we don't need the entries between head..lt advanced.

Flat combining

NR uses [flat combining](#) to allow threads running on the same node to combine their results, resulting in better cache locality both from flat combining local to the node's last-level cache.

1. The combiner can batch and combine multiple operations at a lower cost than executing each operation individually. For example, a priority queue can be done with a single atomic instruction for each operation.
2. A single thread accesses the data structure for multiple operations, ensuring cache locality for the whole batch of operations.

NR also benefits from combining operations placed on the single atomic operation for the whole batch instead of on



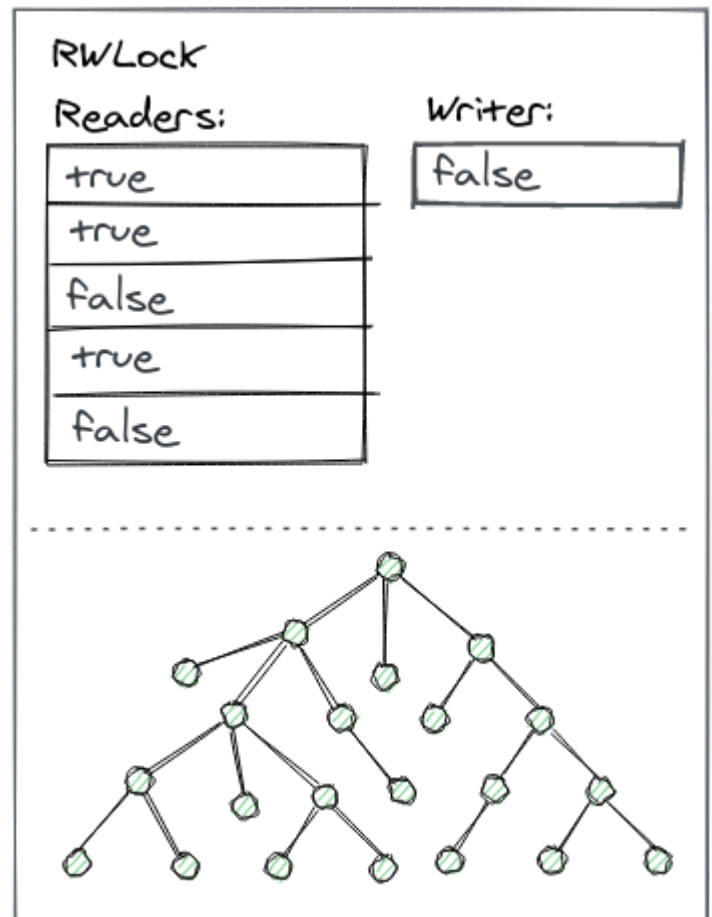
Flat combining multiple operations from threads on a thread is elected who will aggregate outstanding operations, put them on the log in one batch (3) and then execute (4), after executing any previously outstanding operations.

The optimized readers-writer lock

NR uses a writer-preference variant of the distributed RW synchronization of the combiner and reader threads when

Each reader acquires a local per-thread lock; the writer has readers its intent to acquire the writer lock. A reader first acquires its local lock if there is no writer; then it checks the lock (and releases the reader lock if it notices the writer desired (writer preference). The combiner also acquires the writer lock. Thus, we give highest priority to the writer, because we do not want to let readers starve the writer. This lock allows readers to read a local replica while the combiner performs operations to the log, increasing parallelism.

If there is no combiner, a reader might have to acquire the lock from the log to avoid a stale read, but this situation is rare.



The RWLock used to protect the data-structure in every cache-line sized atomic booleans for reader threads thread. It allows concurrent readers but only a single thread. Readers also have to make sure to advance the replication the log.

A concurrent non-mutating operation

A non-mutation operation (read) can execute on its thread entry, but it needs to ensure that the replica is not stale. It reads the current tail when the operation begins, and waits until a combine operation is observed, or acquires the lock and updates the replica.

If there is no combiner, the read thread will acquire the lock and update the tail.

A concurrent mutating operation (update)

A thread *T* executing a mutating operation (update) needs to acquire the lock. If thread *T* executing this operation fails to acquire it, another thread may acquire it. Thread *T* then spin-waits to receive its operation's result from the existing thread.

However, *T* must still periodically attempt to become the combiner to complete its operation.

If *T* acquires the lock, it becomes *the combiner* and executes its operation. It then allows concurrent threads on the same replica. To do so, *T* appends its operation with a single atomic operation and ensures that the replica is up-to-date with all operations and returning their results to the waiting threads.

Source and code example

[Node-replication \(NR\)](#) is released as a stand-alone library.

To give an idea, here is an example that transforms an `HashMap` (from the Rust standard library) into a concurrent, replicated hash-table.

```

//! A minimal example that implements a replicated
use std::collections::HashMap;
use node_replication::Dispatch;
use node_replication::Log;
use node_replication::Replica;

/// The node-replicated hashmap uses a std hashmap
#[derive(Default)]
struct NrHashMap {
    storage: HashMap<u64, u64>,
}

/// We support mutable put operation on the hashmap
#[derive(Clone, Debug, PartialEq)]
enum Modify {
    Put(u64, u64),
}

/// We support an immutable read operation to look
#[derive(Clone, Debug, PartialEq)]
enum Access {
    Get(u64),
}

/// The Dispatch trait executes `ReadOperation` (
/// and `WriteOperation` (our `Modify` enum) again
/// data-structure.
impl Dispatch for NrHashMap {
    type ReadOperation = Access;
    type WriteOperation = Modify;
    type Response = Option<u64>;

    /// The `dispatch` function applies the immutable
    fn dispatch(&self, op: Self::ReadOperation) ->
        match op {
            Access::Get(key) => self.storage.get(&key)
        }
    }

    /// The `dispatch_mut` function applies the mutable
    fn dispatch_mut(&mut self, op: Self::WriteOperation) ->
        match op {
            Modify::Put(key, value) => self.storage.put(key, value)
        }
    }
}

```

As we can see we need to define two operation enums (`Modify` and `Access`) to end up on the log, and `Access` for immutable/read operations. We then implement the `Dispatch` trait from NR for our newly defined `NrHashMap` responsible to route the operations defined in `Access` and `Modify` to the appropriate

structure. [The full example](#), including how to create replic repository.

Concurrent Node Replicati

Some OS subsystems can become limited by [node-replica](#) frequently mutating but would otherwise naturally comm

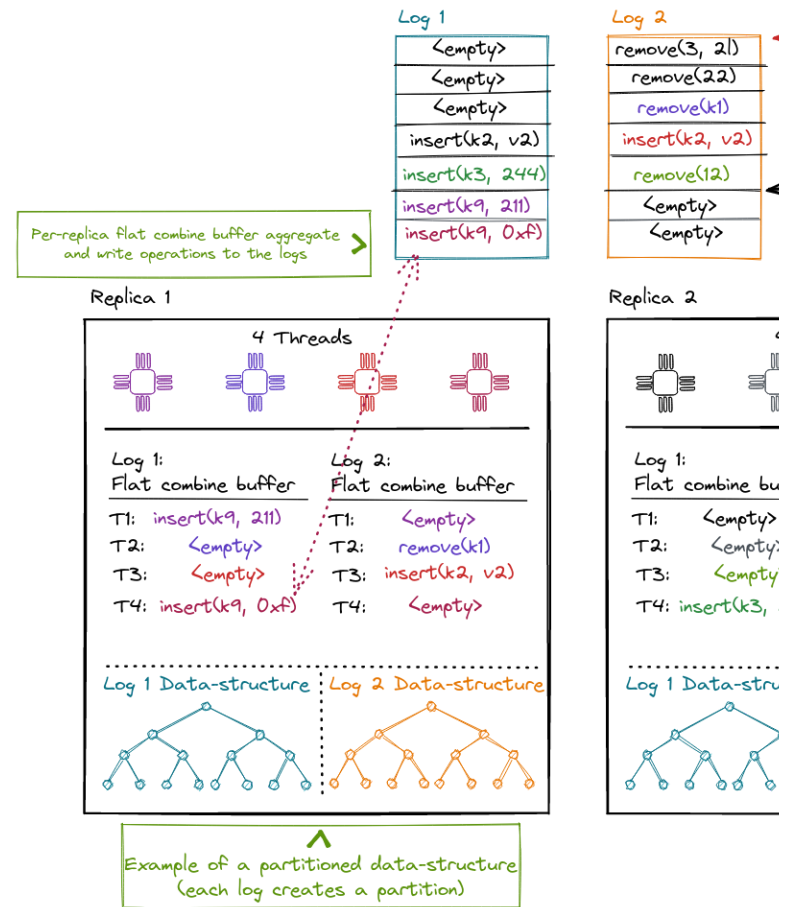
NR allows multiple combiners from different replicas to n
scalability is limited because

1. all combiner threads are operating on a single, shar
2. each replica is a sequential data structure, which rec
writer lock.

To address these problems, we can use CNR, a technique
approach by leveraging operation commutativity present
operations are said to be *commutative* if executing them i
structure in the same abstract state. Otherwise, we say op

As NR, CNR replicates a data structure across NUMA node
the replicas. However, CNR uses commutativity to scale th
logs, by assigning commutative operations to different log
operations always use the same log and thus have a total
concurrent or partitioned data structures for replicas, wh
combiners on each replica -- one for each shared log. This
writer lock and scales access to the data structure.

CNR lifts *an already concurrent data structure* to a NUMA-a
original data structure can be a concurrent (or partitioned
small number of threads (4-8 threads) within a single NUM
lock-free or lock-based and may result in poor performan
concurrent data structure that works well for a large num
across NUMA nodes, and that is resilient to contention.



In CNR a replica can distribute commuting operation replica maintains one flat-combining buffer per log aggregated. One elected combiner thread commits a batch to the log and then applies the ops against the

Compared to NR, the replicated data-structure is no reader-writer lock by default. Instead, the data-structure as in this diagram, use a lock-free approach, or rely on

CNR Operations and Linearizability

CNR is useful for *search data structures*, with operations *insert(x)* and *range-scan(x, y)*. These operations often benefit from both on the abstract operation type and on its input arguments. **boosting**, CNR considers the abstract data type for establishing data structure implementation.

Consider, for example, the *insert(x)* operation. Two operations operate on distinct arguments: e.g., *insert(x)* and *insert(y)*. A concrete implementation of the data structure could be a

that `insert(x)` and `insert(x+1)` are not commutative to memory locations. However, the original data structure at shared memory locations, due to its concurrent nature. H CNR and can be safely executed concurrently.

Interface

CNR's interface is similar to NR, but it adds *operation class* mutable and immutable operations. CNR relies on the previous conflicting operations by assigning them to the same operation information to allocate conflicting operations to the same the same NUMA node, to the same combiner too. In contrast, executed by different combiners and can use different shared memory executed concurrently.

As with NR, CNR executes different steps for mutating (update) operations. Each of these operations uses only one of the according to its log's order. Different logs are not totally conflicting, different logs are commutative.

In addition, CNR special-cases another type of operation, operation class. These are operations that conflict with mutable operation has to read the entire data structure to determine operation to a single class, all other operations need to be commutativity benefit.

Scan-type operations span multiple logs and need a *consistent* logs involved in the operation, obtained *during the lifetime* consistent state, the thread performing the scan collects a snapshot inserting the operation in these logs. This atomic snapshot point.

We show the CNR API, with the additional traits implemented using our earlier example in the NR section:

```

use chashmap::CHashMap;
use cnr::{Dispatch, LogMapper};

/// The replicated hashmap uses a concurrent hashmap
pub struct CNRHashMap {
    storage: CHashMap<usize, usize>,
}

/// We support a mutable put operation on the hashmap
#[derive(Debug, PartialEq, Clone)]
pub enum Modify {
    Put(usize, usize),
}

/// This `LogMapper` implementation distributes the work
/// in a round-robin fashion. One can change the
/// data locality based on the data structure layout.
impl LogMapper for Modify {
    fn hash(&self, nlogs: usize, logs: &mut Vec<usize>) {
        debug_assert!(logs.capacity() >= nlogs, "guarantee or
        debug_assert!(logs.is_empty(), "guarantee or

        match self {
            Modify::Put(key, _val) => logs.push(*key
        }
    }
}

/// We support an immutable read operation to look up a value
#[derive(Debug, PartialEq, Clone)]
pub enum Access {
    Get(usize),
}

/// `Access` follows the same operation to log map
/// ensures that the read and write operations for
/// the same log.
impl LogMapper for Access {
    fn hash(&self, nlogs: usize, logs: &mut Vec<usize>) {
        debug_assert!(logs.capacity() >= nlogs, "guarantee or
        debug_assert!(logs.is_empty(), "guarantee or

        match self {
            Access::Get(key) => logs.push(*key % nlogs
        }
    }
}

/// The Dispatch traits executes `ReadOperation` and
/// and `WriteOperation` (our Modify enum) against the
/// data-structure.
impl Dispatch for CNRHashMap {
    type ReadOperation = Access;

```

```

type WriteOperation = Modify;
type Response = Option<usize>;

/// The `dispatch` function applies the immutal
fn dispatch(&self, op: Self::ReadOperation) ->
    match op {
        Access::Get(key) => self.storage.get(&
    }
}

/// The `dispatch_mut` function applies the mut
fn dispatch_mut(&self, op: Self::WriteOperation) ->
    match op {
        Modify::Put(key, value) => self.storage
    }
}
}

```

CNR is available as a stand-alone rust library together wit

Comparison to NR and Notation

CNR benefits from the NR techniques, such as flat combinational operation description, which reduce contention and inter

As NR, CNR has two drawbacks: increased memory footprint from re-executing each operation on each replica. CNR has a single shared log can be split into multiple, smaller shared logs to increase parallelism within each NUMA node by using a concurrent log. CNR increases parallelism across NUMA nodes by using multiple

Memory

Kernel address space

The kernel address space layout follows a simple scheme mapped with a constant offset `KERNEL_BASE` in the kernel physical address can be accessed in the kernel by adding

All physical memory is always accessible in the kernel and mapped/unmapped at runtime. The kernel binary is linked loaded into physical memory and then relocated to run at (`KERNEL_BASE` + physical address).



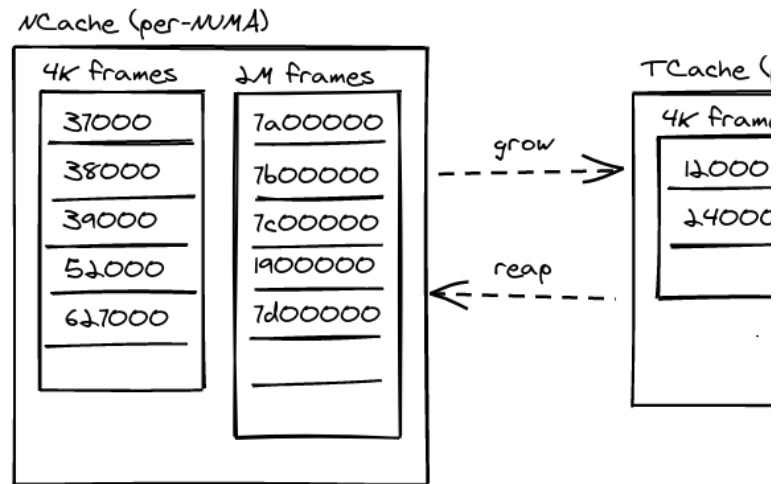
A view of different address spaces in nrk (physical, k

Physical memory

Physical memory allocation and dynamic memory allocation basic subsystems that do not rely on NR. Replicated subsystems but that allocation operation itself should not be replicated

mapping in a page table, each page table entry should refer to all replicas (though, each replica should have its own page table replicated, each allocation operation would be repeated on all replicas).

At boot time, the affinity for memory regions is identified, and memory is placed in NUMA node caches (FrameCacheLarge). The FrameCacheLarge is further divided into two classes of 4 KiB and 2 MiB frames. Every node has 4 KiB and 2 MiB frames for fast, no-contention allocation with small size. If it is empty, it refills from its local FrameCacheLarge. FrameCacheSmall and FrameCacheLarge implement a cache hierarchy that controls the flow between TCaches and NCaches.



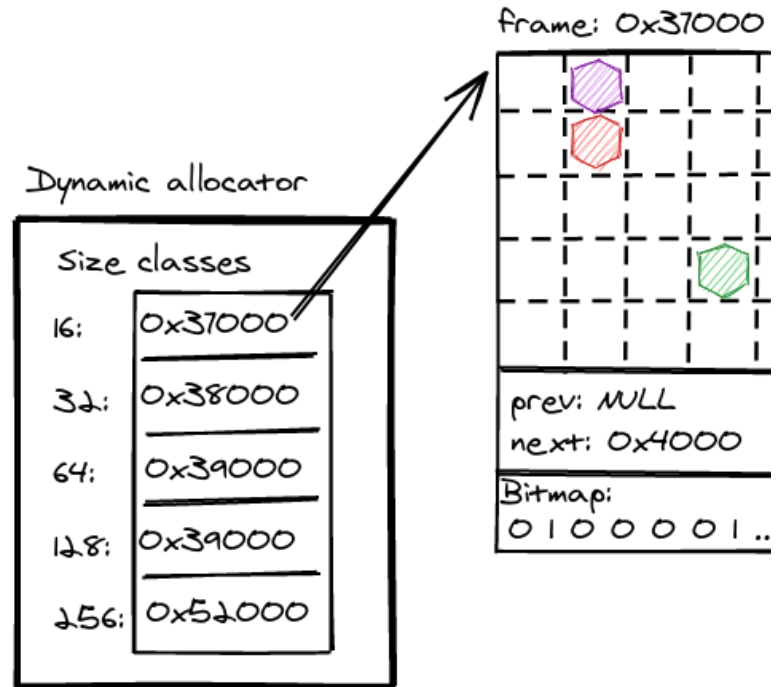
Shows a global per-NUMA FrameCacheLarge and a per-core FrameCacheSmall. Cores allocate 4K or 2M pages directly from the FrameCacheSmall. If the FrameCacheSmall is empty (grow), TCaches allocate frames in stacks to allow for quick allocation and deallocation.

Dynamic memory

Since NRK is implemented in Rust, memory management is handled by the Rust compiler to track the lifetime of allocated objects. This eliminates the need for manual memory management (e.g., after-free, uninitialized memory *etc.*), but the kernel still has to manage memory. NRK uses fallible allocations and intrusive data structures to handle errors gracefully.

The dynamic memory allocator in nrk provides an implementation of the [memory allocator interface](#). It uses size classes and different allocators per size class (e.g., 4 KiB or 2 MiB frames are used which are suitable for different size classes).

given class. A bitfield at the end of every frame tracks the frame (e.g., to determine if its allocated or not).



The dynamic memory allocator for kernel objects in containing two frames for less than 16 byte objects. allocated slots along with per-frame meta-data (prev indicate allocated blocks. Typically, one dynamic me instantiated.

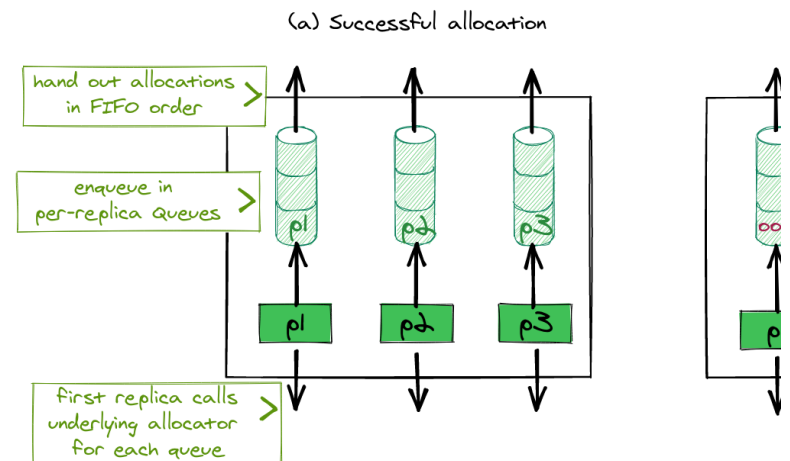
Deterministic memory

The kernel has to explicitly handle running out of memory handle out-of-memory errors gracefully by returning an e almost all cases in the kernel (with some exceptions durir 3rd party dependencies (e.g., to parse ELF binaries) are n allocations yet.

Another issue is that handling out-of-memory errors in pr becomes a little more challenging: Allocations which happ deterministic (e.g. they should either succeed on all replic would end up in an inconsistent state if after executing ar successful and some had unsuccessful allocations. Making equal amounts of memory available is infeasible because times and meanwhile allocations can happen on other co this problem in nrk by requiring that all memory allocatio

through a deterministic allocator. In the deterministic allocation request allocates memory on behalf of all other allocators remembers the results temporarily, until they are running behind. If an allocation for any of the replicas enqueue the error for all replicas, to ensure that all replicas Allocators in nrk are chainable and it is sufficient for the coordinator in the chain so it doesn't necessarily have to be invoked for request. Our implementation leverages custom allocators heap allocator is used for individual data-structures.

Deterministic Allocator



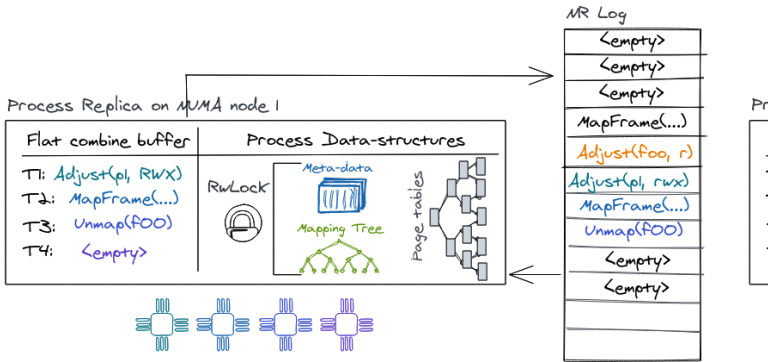
The deterministic memory allocator generally forwards requests to an underlying memory allocator. However, the first request is processed on behalf of all other replicas and store the results temporarily (one per replica). The deterministic allocator ensures that all replicas have a successful allocation for a given request or return a result back.

Process structure

Virtual memory

NRK relies on the MMU for isolation. Like most conventional OSes, NRK uses a per-process mapping database (as a B-Tree) to construct the process's hardware page tables. NRK currently uses a B-Tree for the mapping database. Both the B-Tree and the hardware page tables are wrapped behind the NR interface for concurrency.

Therefore, the mapping database and page tables are probably the biggest part of the process abstraction. A process exports several operations to modify its address space: `MapFrame` (to insert new mappings); and `Adjust` (to change permissions of a mapping). `Adjust` also supports a non-mutating `Resolve` operation (that advances the address space state).



Schematic diagram of a node-replicated process in NRK

Virtual memory and NR

There are several aspects of the virtual-memory design that are specific to NRK:

For example, the virtual-memory has to consider out-of-bounds address walks. Normally a read operation would go through all outstanding operations from the log first. However, a process must have this capability. A race arises if a process maps a page in replica B accesses that mapping in userspace before replica A has finished. This can be handled since it generates a page fault. In order to advance the replica by issuing a `Resolve` operation.

corresponding mapping of the virtual address generating process can be resumed since the `Resolve`-operation will mapping is found, the access was an invalid read or write

`Unmap` or `Adjust` (e.g., removing or modifying page-table entries on cores where the process is active to ensure TLB software, by the OS, and commonly referred to as perform core will start by enqueueing the operation for the local replica the unmap (or adjust) operation has been performed at local replica and is enqueue as a future operation on the log of processor interrupts (IPIs) to trigger TLB flushes on all cores process. As part of the IPI handler the cores will first acknowledge they must make sure to advance their local replica with operation forces the unmap/adjust if not already applied), then poll information about the regions that need to be flushed, and Meanwhile the initiator will invalidate its own TLB entries acknowledgments from other cores before it can return to

Scheduler

In NRK, the kernel-level scheduler is a coarse-grained scheduler for processes. Processes make system calls to request for more resources. The kernel notifies processes core allocations and deallocations. A process allocates executor objects (*i.e.*, the equivalent of a thread) to dispatch a given process on a CPU. An executor mainly consists of an upcall handler and one for the initial stack) and a set of metadata. Executors are allocated lazily but a process keeps them over time.

In the process, a userspace scheduler reacts to upcalls from the kernel, and it makes fine-grained scheduling decisions by dispatching processes to cores. This design means that the kernel is only responsible for coarse-grained scheduling and implements a global policy of core allocation to processes.

The scheduler uses a sequential hash table wrapped with process structure and to map process executors to cores. For each process, a process; to allocate and deallocate executors for a process on a given core.

File System

The NrFS is a simple, in-memory file system in nrk that supports standard file operations (`open`, `pread`, `pwrite`, `close`, *etc.*).

NrFS tracks files and directories by mapping each path to each inode number to an in-memory inode. Each inode has a list of file pages. The entire data structure is wrapped in a mutex for replication.

User Space

The main user-space components are explained in more

KPI: Kernel Public Interface

The Kernel Public Interface (KPI) is the lowest level user-space interface. As the name suggests it is the common interface between user-space programs and the kernel. It is special because it is the only library that is part of the kernel code.

The KPI contains the syscall interface and various structures that define the boundary between the kernel and user-space. If in the future, we cannot try to keep the syscall ABI compatible but would rather define a new boundary.

Typically, the KPI functionality will rarely be accessed directly. Instead, parts of it are re-exported or wrapped by the [vibrio](#) library, `lib/kpi`.

Lineup

Lineup is a user-space, cooperative thread scheduler that threads). It supports many synchronization primitives (mutexes, barriers etc.), thread-local storage, and has some basic support for compiler-assisted context-switching. The scheduler

Upcalls

The kernel can notify the scheduler about events through upcalls. The scheduler can then use these upcalls to notify about more available cores (or removal of cores) or to handle page-faults from kernel to user-space.

The mechanism for this is inspired by scheduler activation. When a process arrives at the kernel, it will save the current CPU context (registers, stack, etc.) into a common save area and resume the process with a new (mostly empty) context. The upcall handler gets triggered by the interrupt through function arguments so it can decide how to react. After the event is handled, the upcall handler restores the context (from before the interruption) from the common save area and resumes the computation left off before the upcall (or decide not to continue).

Vibrio

Vibrio is the user-space library that provides most of the functions for applications in user-space. It is found in `lib/vibrio`.

Memory

The user-space memory manager provides a `malloc` and `GlobalAlloc` implementation for rust programs. We rely on it for small--medium sized blocks (between 0 and 2 MiB) by allocating memory with the `map` syscall.

RumpRT

A rumpkernel is a componentized NetBSD kernel that can run in various environments. It contains file systems, a POSIX system call interface, a SCSI protocol stack, virtio, a TCP/IP stack, `libc` and `libpthread`.

Vibrio has a `rump` module which provides the necessary interface to run a rumpkernel inside a user-space process (e.g., the `rumpus` module). The advantage that it's possible to run many POSIX compatible applications is by building a fully-fledged POSIX compatibility layer into `NrO`.

- Bare-metal and Xen implementations for [rump](#)
- [Some supported applications](#)
- [PhD thesis about rumpkernels](#)

Vibrio dependency graph

Vibrio uses the following crates / dependencies:

vibrio

- arrayvec
- bitflags
- crossbeam-utils
- cstr_core
- hashbrown
- kpi
- lazy_static
- lineup
- log
- rawtime
- rumpkernel
- serde_cbor
- slabmalloc
- spin
- x86

rkapps

The rkapps directory (in `usr/rkapps`) is an empty project build file contains the steps to clone and build a few different (memcached, LevelDB, Redis etc.) that we use to test user

The checked-out program sources and binaries are placed in your build directory:

```
target/x86_64-nrk-none/<debug | release>/build/rkapps
```

The build for these programs can be a bit hard to understand

1. Clone the packages repo which has build instructions for running on rumpkernels.
2. For each application that we want to build (enabled/disabled) in the respective directory. This will compile the application using the rumpkernel toolchain. The toolchain can be found in a subdirectory `target/x86_64-nrk-none/<debug | release>/build/rkapps`.
3. Linking binaries with vibrio which provides the low-level interface.

For more information on how to run rkapps applications, see the [application](#) section.

Development

This chapter should teach you how to [Build](#), [Run](#), [Debug](#),

Configuration

Some tips and pointers for setting up and configuring the

VSCode

VSCode generally works well for developing nrk. The `rust` `rls` which often has build issues due to the project not h

Git

For first time git users or new accounts, you'll have to con

```
git config --global user.name "Gerd Zellweger"  
git config --global user.email "mail@gerdzellweger"
```

To have better usability when working with submodules, `git pull` etc.

```
git config --global submodule.recurse true
```

Fetch multiple submodules in parallel:

```
git config --global submodule.fetchJobs 20
```

We don't allow merge requests on master, to always keep can be helpful:

```
[alias]  
  purr = pull --rebase
```

Adding a new submodule to the repository

1. `cd lib`
2. `git submodule add <path-to-repo> <foldername>`

Removing a submodule in the repository

1. Delete the relevant section from the `.gitmodules` file
2. Stage the `.gitmodules` changes: `git add .gitmodules`
3. Delete the relevant section from `.git/config`.
4. Run `git rm --cached path_to_submodule` (no trailing slash)
5. Run `rm -rf .git/modules/path_to_submodule` (no trailing slash)
6. Commit changes

Styleguide

Code format

We rely on [rustfmt](#) to automatically format our code.

Code organization

We organize/separate imports into three blocks (all separ

- 1st block for core language things: `core`, `alloc`, `st`
- 2nd block for libraries: `vibrio`, `x86`, `lazy_static` (
- 3rd block for internal imports: `crate::*`, `super::*`
- 4th block for re-exports: `pub(crate) use::*` etc.
- 5th block for modules: `mod foo`; etc.

Afterwards a `.rs` file should (roughly) have the following

- 1st `type` declarations
- 2nd `const` declarations
- 3rd `static` declarations
- 4th `struct`, `fn`, `impl` etc. declarations

Visibility

Avoid the use of `pub` in the kernel. Use `pub(crate)`, `pub` code elimination.

Assembly

We use AT&T syntax for assembly code (`options(att_syr`

Cargo features

Libraries and binaries only have non-additive / non-conflicting compilation problems quickly (e.g. with `cargo build --a`

Errors

The `KError` type is used to represent errors in the kernel should only be used once/in a single location (to be easy to find a descriptive name).

Formatting Commit Messages

We follow the conventions on [How to Write a Git Commit](#)

Be sure to include any related GitHub issue references in the commit message for referencing issues and commits.

Github pull requests & history

Since github doesn't do fast-forward merges through the command line to keep the same commit hashes of the

```
git checkout master
git merge --ff-only feature-branch-name
```

Building

There are two sets of dependencies required for the development of nrk. We typically build, develop and test using nrk in QEMU. Other Linux systems will probably work but not all dependencies. Other operating systems likely won't work without adjustments for code and the build-process.

Check-out the source tree

Check out the nrk sources first:

```
git clone <repo-url>
cd nrk
```

The repository is structured using [git submodules](#). You'll have to checkout submodules separately:

In case you don't have the SSH key of your machine registered, you will need to convert all submodule URLs to use the https protocol. Run this sed script before proceeding:

```
sed -i'' -e 's/git@github.com:/https:\\\\github.com/'
```

```
git submodule update --init
```

Dependencies

If you want to build without [Docker](#), you can install both kernel and user-space dependencies by executing `setup.sh` in the root of the repository directly (tested on *latest Ubuntu LTS*). The script will install all required OS packages and some additional rust programs and dependencies. To run the tests, you also have to install the [DCM-based scheduler dependencies](#).

The build dependencies can be divided into these categories:

- Rust (nightly) and the `rust-src` component for cor
- `python3` (and some python libraries) to execute the
- Test dependencies (qemu, corealloc, dhcpd, redis-b
- Rumpkernel dependencies (gcc, zlib1g etc.)
- Build for documentation ([mdbook](#))

See `scripts/generic-setup.sh` function `install_build`.

Use Docker

We provide scripts to create a docker image which contain

To use Docker, it needs to be installed in your system.
steps:

```
sudo apt install docker.io
sudo service docker restart
sudo addgroup $USER docker
newgrp docker
```

To create the image execute the following command in th

```
bash ./docker-run.sh
```

This will create the docker image and start the container.
running inside the Docker container. You can build the OS
dependencies natively.

The script will create a user inside the docker container
the host system (same username and user ID).

You can rebuild the image with:

```
bash ./docker-run.sh force-build
```

To exit the container, just type `exit` to terminate the she

Build without running

To just build the OS invoke the `run.py` script (in the `kernel` directory) with the `-n` flag (no-run flag).

```
python3 kernel/run.py -n
```

If you want to run the build in a docker container, run `build_docker.sh` beforehand. The source directory tree will be mounted in

Using run.py

The `kernel/run.py` script provides a simple way to build, settings and configuration. For a complete set of parameters `run.py --help` instructions.

As an example, the following invocation

```
python3 run.py --kfeatures test-userspace --cmd='1
rkapps init --ufeatures rkapps:redis --machine qemu
qemu-cores 2
```

will

- compile the kernel with Cargo feature `test-userspace`
- pass the kernel the command-line arguments `log=info` (sets logging to info and starts redis.bin for testing)
- Compile two user-space modules `rkapps` (with cargo features)
- Deploy and run the compiled system on `qemu` with 2 cores allocated to the VM

If Docker is used as build environment, it is necessary to find required features inside the Docker container:

```
python3 run.py --kfeatures test-userspace --mods rkapps:redis -n
```

Afterwards, the aforementioned command can be used to run the container with the given configuration. The `run.py` script has already been build and will directly start `qemu`.

Sometimes it's helpful to know what commands are actually used to figure out what the exact qemu command line invocation should be supplied.

Depending on the underlying system configuration NRK may not be able to establish a local network. In this case, the following issue:

1. Disable AppArmor. Detailed instructions can be found in the [AppArmor](#) section.
2. Manually start the DHCP server immediately after NRK is started.

```
sudo dhcpcd -f -d tap0 --no-pid -cf ./kernel/t
```

Baremetal execution

The `kernel/run.py` script supports execution on baremetal with the `--baremetal` argument:

```
python3 run.py --machine b1542 --verbose --cmd "l
```

This invocation will try to run nrk on the machine described in the TOML file:

A TOML file for a machine has the following format:

```
[server]
# A name for the server we're trying to boot
name = "b1542"
# The hostname, where to reach the server
hostname = "b1542.test.com"
# The type of the machine
type = "skylake2x"
# An arbitrary command to set-up the PXE boot environment
# This often involves creating a hardlink of a file
# of the machine and pointing it to some pxe boot
pre-boot-cmd = "./pxeboot-configure.sh -m E4-43-4f

# run.py support only booting machines that have a BIOS
[idrac]
# How to reach the ilo/iDRAC interface of the machine
hostname = "b1542-ilo.test.com"
# Login information for iDRAC
username = "user"
password = "pass"
# Serial console which we'll read from
console = "com2"
# Which iDRAC version we're dealing with (current is 3)
idrac-version = "3"
# Typical time until machine is booted
boot-timeout = 320

[deploy]
# Server where binaries are deployed for booting via PXE
hostname = "ipxe-server.test.com"
username = "user"
ssh-pubkey = "~/.ssh/id_rsa"
# Where to deploy kernel and user binaries
ipxe-deploy = "/home/gz/public_html/"
```

An iPXE environment that the machine will boot from needs to be compiled with UEFI and ELF support for running on bare-metal machines.

Note that the current support for bare-metal execution on machines with an iDRAC management console (needs redfish or SNMP support) will be added in the future.

Compiling the iPXE bootloader

TBD.

Debugging

Currently the debugging facilities are not as good as on a However, there are some options available: `gdb`, `printf` - at code. We will discuss the options in this chapter.

GDB support in the kernel

tldr: To use `gdb`, add `--kgdb` to `run.py`.

NRK provides an implementation for the `gdb` remote protocol communication. This means you can use `gdb` to connect to

To use it, start `run.py` with the `--kgdb` argument. Once

```
Waiting for a GDB connection on I/O port 0x2f8...
Use `target remote localhost:1234` in gdb session
```

Next, connect with GDB to the kernel, using:

```
$ cd kernel
$ gdb ../target/x86_64-uefi/<debug | release>/esp,
[...]
(gdb) target remote localhost:1234
Remote debugging using localhost:1234
[...]
```

If you execute `gdb` in the kernel directory, the `.gdbinit` `target remote localhost:1234` for you. But you have "trusted" path by adding this line to `$HOME/.gdbinit`:

```
add-auto-load-safe-path <REPO-BASE>/kernel/.gdb
```

The [GDB dashboard](#) works as well, just insert `target remote` `.gdbinit` file.

Breakpoints

tldr: use `break` or `hbreak` in gdb.

Currently the maximum limit of supported breakpoints (a

Why? Because we use the x86-64 debug registers for break registers. Our gdb stub implements both software and hardware registers.

An alternative technique would be to either insert `int3` instructions and let gdb do it automatically if software interrupts are marked. However, this is a bit more complicated because we need to handle the interrupt handler (e.g., the debug interrupt handler should ideally not hit a debug register, this is fairly easy to achieve, as we just have to enable them when we resume, whereas the `int3` approach requires a bunch of `.text` offsets. On the plus side it would enable this ever becomes necessary.

Watchpoints

Again the maximum limit is four watchpoints (and breakpoints).

Use `watch -l <variable>` to set a watchpoint. The `-l` option specifies the memory location of the variable/expression rather than the variable name. This is not supported as gdb may try to overwrite `.text` location (which would not execute) in the kernel.

printf debugging with the log crate

Here are a few tips:

- Change the log-level of the kernel to info, debug, or trace:
`cmd='log=info'`
- Logging can also be enabled per-module basis. For example, to enable logging for just the `gdbstub` library and the `gdb` module in the `cmd` invocation for `run.py` would look like: `--cmd "log='gdbstub=trace,nrk::arch::gdb=trace'"`

- Change the log-level of the user-space libOS in vibrio
- Make sure the [Tests](#) run (to see if something broke).

Figuring out why things failed

Maybe you'll encounter failures, for example like this one

```
[IRQ] GENERAL PROTECTION FAULT: From Any memory read
checks.
```

```
No error!
```

```
Instruction Pointer: 0x534a39
```

```
ExceptionArguments { vec = 0xd exception = 0x0 rip =
0x13206 rsp = 0x5210400928 ss = 0x1b }
```

```
Register State:
```

```
Some(SaveArea
```

```

rax =                0x0 rdx =                0x0
0x0
rsi =                0x0 rdi =                0x5210400a50
0x5210400928
r8  =                0x2 r9  =                0x5202044c00
0x28927a
r12 =                0x520e266810 r13 =                0x7d8ac0
0x686680
rip =                0x534a39 rflags = FLAGS_RF | FLAGS_
FLAGS_IOPL3 | FLAGS_IF | FLAGS_PF | FLAGS_A1)
stack[0] = 0x5210400958
stack[1] = 0x53c7fd
stack[2] = 0x0
stack[3] = 0x0
stack[4] = 0x0
stack[5] = 0x0
stack[6] = 0x52104009b8
stack[7] = 0x534829
stack[8] = 0x5210400a50
stack[9] = 0x5210400a50
stack[10] = 0x0
stack[11] = 0x268
```

The typical workflow to figure out what went wrong:

1. Generally, look for the instruction pointer (`rip`) which
2. If the instruction pointer (and `rsp` and `rbp`) is below user-space when the failure happened (you can also it's easier to tell from the other registers).
3. Determine exactly where the error happened. To do which was running. Those are usually located in `tar`

```
uefi/<release|debug>/esp/<binary> .
```

4. Use `addr2line -e <path to binary> <rip>` to see
5. If the failure was in kernel space, make sure you adjust the PIE offset where the kernel binary was executing in the following line `INFO: Kernel loaded at address: 0` by the bootloader early during the boot process. Subtract the offset in the ELF file.
6. Sometimes `addr2line` doesn't find anything, it's good to give more context: `objdump -S --disassemble --c uEFI/<release|debug>/esp/<binary> | less`
7. The function that gets reported might not be useful in some case, look for addresses that could be return addresses too (e.g., `0x534829` looks suspiciously like a return address).
8. If all this fails, something went wrong in a bad way, try debugging.

Always find the first occurrence of a failure in the serial log. The log is not very robust, it still quite often triggers cascading failures that are not relevant.

Debugging rumpkernel/NetBSD console

The `nrk` user-space links with a rather large (NetBSD) code-base. Since it's in there, it's sometimes helpful to temporarily change or fix the C code.

You can edit that code-base directly since it gets checked out. For example, to edit the `rump_init` function, open the file `target/x86_64-nrk-none/release/obj/netbsd/sys/rump/librump/rumpkern/rump.c`

Make sure to identify the correct `$HASH` that is used for the multiple `rumpkernel-*` directories in the build dir, otherwise

After you're done with edits, you can manually invoke the

As a simple example you can search for `rump_init(void)` in `none/release/build/rumpkernel-$HASH/out` and add a `printf`. The following steps should ensure the print also appears on the

```
cd target/x86_64-nrk-none/release/build/rumpkernel
./build-rr.sh -j24 nrk -- -F "CFLAGS=-w"
# Delete ../target/x86_64-nrk-none/debug/build/in
# Invoke run.py again...
```

If you change the compiler/rustc version, do a clean build as your changes might be overridden as the sources exist in `target`. It's a good idea to save changes somewhere important.

Debugging in QEMU/KVM

If the system ends up in a dead-lock, you might be able to debug it by asking qemu. Deadlocks with our kernel design (locking APIs) it can definitely happen.

The following steps should help:

1. Add `--qemu-monitor` to the `run.py` invocation to start a monitor.
2. Connect to the monitor in a new terminal with `telnet`.
3. You can use `info registers -a` to get a dump of the registers or any other command to query the hypervisor state.
4. If you're stuck in some loop, getting a couple registers by invoking `info registers` just once.

When developing drivers that are emulated in qemu, it can be useful to use the interface in QEMU to see what state the device is in. For `vmxnet3` in the sources, you can change the `#undef` in `hw/net/vmxnet_debug.h` to `#define` and recompile the code (the look similar to this snippet below):

```
#define VMXNET_DEBUG_CB
#define VMXNET_DEBUG_INTERRUPTS
#define VMXNET_DEBUG_CONFIG
#define VMXNET_DEBUG_RINGS
#define VMXNET_DEBUG_PACKETS
#define VMXNET_DEBUG_SHMEM_ACCESS
```

Testing

If you've found and fixed a bug, we better write a test for it and methodologies to ensure everything works as expected.

- Regular unit tests: Those can be executed running `cargo test`. Sometimes adding `RUST_TEST_THREADS=1` is necessary depending on the runner/frameworks used. This should be indicated in the test's documentation.
- A slightly more exhaustive variant of unit tests is provided by `cargo test --test s00_core` to make sure that the implementation of kernel subsystems is correct.
- Integration tests are found in the kernel, they typically [reexpect](#) to interact with the guest.
- Fuzz testing: TBD.

Running tests

To run the unit tests of the kernel:

1. `cd kernel`
2. `RUST_BACKTRACE=1 RUST_TEST_THREADS=1 cargo test`

To run the integration tests of the kernel:

1. `cd kernel`
2. `RUST_TEST_THREADS=1 cargo test --test '*'`

If you would like to run a specific integration test you can use:

1. `RUST_TEST_THREADS=1 cargo test --test '*' --test s00_core`

If you would like to run a specific set of integration tests, you can use:

1. `RUST_TEST_THREADS=1 cargo test --test s00_core`

In case an integration test fails, adding `--nocapture` at the end of the command will make sure that the underlying `run.py` invocations are printed to the console, which is helpful to figure out the exact `run.py` invocation that a test is failing. You can also try yourself manually for debugging.

Parallel testing for the kernel is not possible at the moment for testing.

The `commitable.sh` script automatically runs the unit and

```
cd kernel
bash commitable.sh
```

Writing a unit-test for the kernel

Typically these can just be declared in the code using `#[test]` to run under the `unix` platform. A small hack is necessary to also run under unix too: When run on a x86-64 unix platform the kernel in `arch/x86_64/` will be included as a module name `kernel` would be `arch`. This is a double-edged sword: we can now run metal code (great), but we can also easily crash the test program MSR for example (e.g, things that would require ring 0 privileges).

Writing an integration test for the kernel

Integration tests typically spawn a QEMU instance and boot it in a virtual space with a custom set of Cargo feature flags. Then it passes the expected output. Part of those custom compile flags is a `test` function that is different than the one you're seeing (which will go off to the kernel programs for example).

There are two parts to the integration test.

- The host side (that will go off and spawn a qemu instance and run the tests. It is found in `kernel/tests`).
- The corresponding main functions in the kernel that are tested. For example are located at `kernel/src/integration_main.rs`.

To add a new integration test the following steps may be required:

1. Modify `kernel/Cargo.toml` to add a feature (under

2. Optional: Add a new `xmain` function and test impler `kernel/src/integration_main.rs` with the used fe, also be possible to re-use an existing `xmain` function name used to include it.
3. Add a runner function to one of the files in `kernel/` cargo feature runs it and checks the output.

Integration tests are divided into categories and named a tests run in a sensible order):

- `s00_*` : Core kernel functionality like boot-up and fa
- `s01_*` : Low level kernel services: SSE, memory alloc
- `s02_*` : High level kernel services: ACPI, core booting
- `s03_*` : High level kernel functionality: Spawn cores,
- `s04_*` : User-space runtimes
- `s05_*` : User-space applications
- `s06_*` : Rackscale (distributed) tests

Benchmarks are named as such:

- `s10_*` : User-space applications benchmarks
- `s11_*` : Rackscale (distributed) benchmarks

The `s11_*` benchmarks may be configured with two feati

- `baseline` : Runs NrOS configured similarly to racksc
- `affinity-shmem` : Runs the `ivshmem-server` using `s` option requires preconfiguring `hugetlbfs` with `sudo` having a kernel with 2MB huge pages enabled, and t node, with a command like: `echo <page-num> | suc /proc/sys/vm/nr_hugepages_mempolicy` The numbe verified with `numastat -m`.

Network

nrk has support for three network interfaces at the mom and e1000 are available by using the respective `rumpkerr` `vmxnet3` is a standalone implementation that uses `smolt` capable of running in ring 0.

Network Setup

The integration tests that run multiple instances of nrk re those integration tests, the test framework calls run.py wi destroy existing conflicting tap interfaces and create new the number of hosts in the test. Then, to run the nrk insta `no-network-setup` flag.

To setup the network for a single client and server (`--wor` following command:

```
python3 run.py --kfeatures integration-test --cmd
workers 2 --network-only
```

Ping

A simple check is to use ping (on the host) to test the netw Adaptive `ping -A`, flooding `ping -f` are good modes to stack work and can handle an "infinite" amount of packet

Some expected output if it's working:

```
$ ping 172.31.0.10
64 bytes from 172.31.0.10: icmp_seq=1 ttl=64 time=
64 bytes from 172.31.0.10: icmp_seq=2 ttl=64 time=
64 bytes from 172.31.0.10: icmp_seq=3 ttl=64 time=
64 bytes from 172.31.0.10: icmp_seq=4 ttl=64 time=
```

For network tests, it's easiest to start a DHCP server for th IP by communicating with the server:

```
# Stop apparmor from blocking a custom dhcp instar
service apparmor stop
# Terminate any (old) existing dhcp instance
sudo killall dhcpd
# Spawn a dhcp server, in the kernel/ directory de
sudo dhcpd -f -d tap0 --no-pid -cf ./tests/dhcpd.c
```

A fully automated CI test that checks the network using pi invoked with the following command:

```
RUST_TEST_THREADS=1 cargo test --test '*' -- s04_u
```


socat and netcat

`socat` is a helpful utility on the host to interface with the port and print on incoming packets on the command line

```
socat UDP-LISTEN:8889,fork stdout
```

Similarly we can use `netcat` to connect to a port and send

```
nc 172.31.0.10 6337
```

The integration tests `s05_redis_smoke` and `s04_userspa` tool to verify that networking is working as expected.

tcpdump

`tcpdump` is another handy tool to see all packets that are For debugging nrk network issues, this command is useful

```
tcpdump -i tap0 -vvv -XX
```

Tracing

Use Intel PT (processor-trace). TBD.

Benchmarking

This chapter provides notes and pointers on how to set-up benchmarking and run various OS micro-benchmarks.

Microbenchmarks

File-system

The code contains an implementation of the [fxmark](#) benchmark located at `usr/init/src/fxmark`.

To run the fxmark benchmarks invoke the following command:

```
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_
```

fxmark supports several different file benchmarks:

- *drbh*: Read a shared block in a shared file
- *drbl*: Read a block in a private file.
- *dwol*: Overwrite a block in a private file.
- *dwom*: Overwrite a private block in a shared file.
- *mwrl*: Rename a private file in a private directory.
- *mwrn*: Move a private file to a shared directory.
- *mix*: Access/overwrite a random block (with fixed pattern).

By default the integration test might not run all benchmarks. To change what benchmarks are run or study it to determine the correct arguments to `run.py`.

Address-space

The following integration tests benchmark the address-space:

- `s10_vmops_benchmark`: This benchmark repeatedly increases the number of pages in the process' address space, while varying the number of core workers in its own partition of the address space to measure throughput (operations per second).
- `s10_vmops_latency_benchmark`: Same as `s10_vmops_benchmark` instead of throughput.

- `s10_vmops_unmaplat_latency_benchmark` : The benchmark maps a memory space, then spawns a series of threads on other cores, then unmaps the frame and measures the latency of the threads. The benchmark is dominated by completing the TLB shutdown protocol.
- `s10_shutdown_simple` : The benchmark measures the time taken for programming the APIC and sending IPIs to initiate a shutdown.

The benchmark code is located at `usr/init/src/vmops/`.

```
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_
nocapture
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_
--nocapture
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_
```

Network

TBD.

Benchmarking Redis

Redis is a simple, single-threaded key-value store written in C. It is designed for high performance and single-threaded performance of the system.

Automated integration tests

The easiest way to run redis on nrk, is to invoke the redis

- `s05_redis_smoke` will spawn nrk with a redis instance and run a few commands to test basic functionality.
- `s10_redis_benchmark_virtio` and `s10_redis_benchmark_e1000` will spawn a redis instance and launch the `redis-benchmark` CLI. The results obtained by `redis-benchmark` are parsed and saved to `redis_benchmark.csv`. The `virtio` and `e1000` suffixes are used.

```
cd kernel
# Runs both _virtio and _e1000 redis benchmark tests
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_
```

Launch redis manually

You can also do the steps that the integration test does manually. The `apparmor teardown` is necessary if you don't have a `dhcpd.conf` at this location.

```
cd kernel
sudo service apparmor teardown
sudo dhcpd -f -d tap0 --no-pid -cf ./tests/dhcpd.conf
```

Next run the redis server in nrk (exchange the nic parameter with the NIC):

```
python3 run.py \
  --kfeatures test-userspace \
  --nic e1000 \
  --cmd "log=info init=redis.bin" \
  --mods rkapps \
  --ufeatures "rkapps:redis" \
  --qemu-settings="-m 1024M"
```

Finally, execute the redis-benchmark on the host.

```
redis-benchmark -h 172.31.0.10 -n 10000000 -p 6379
```

You should see an output similar to this:

```
===== SET =====
10000000 requests completed in 10.29 seconds
50 parallel clients
3 bytes payload
keep alive: 1

0.31% <= 1 milliseconds
98.53% <= 2 milliseconds
99.89% <= 3 milliseconds
100.00% <= 4 milliseconds
100.00% <= 4 milliseconds
972100.75 requests per second

===== GET =====
10000000 requests completed in 19.97 seconds
50 parallel clients
3 bytes payload
keep alive: 1

0.14% <= 1 milliseconds
6.35% <= 2 milliseconds
77.66% <= 3 milliseconds
94.62% <= 4 milliseconds
97.04% <= 5 milliseconds
99.35% <= 6 milliseconds
99.76% <= 7 milliseconds
99.94% <= 8 milliseconds
99.99% <= 9 milliseconds
99.99% <= 10 milliseconds
100.00% <= 11 milliseconds
100.00% <= 11 milliseconds
500726.03 requests per second
```

redis-benchmark

`redis-benchmark` is a closed-loop benchmarking tool that can be installed on Ubuntu by installing the `redis-tools` package:

```
sudo apt-get install redis-tools
```

Example invocation:

```
redis-benchmark -h 172.31.0.10 -t set,ping
```

For maximal throughput, use pipelining (`-P`), and the virtio

```
redis-benchmark -h 172.31.0.10 -n 10000000 -p 6379
```

Redis on Linux

You'll need a Linux VM image, see the [Compare against](#) to create one.

Before starting the VM we can re-use the DHCP server configuration of the DHCP server on the host that configures the network of the host.

```
# Launch a DHCP server (can reuse nrk config)
cd nrk/kernel
sudo dhcpcd -f -d tap0 --no-pid -cf ./tests/dhcpd.conf
```

Next, start the VM. To have the same benchmarks conditions, launch the VM like this (select either `e1000` or `virtio`, generate emulated `e1000`):

e1000 NIC:

```
qemu-system-x86_64 \
  --enable-kvm -m 2048 -k en-us --smp 2 \
  -boot d ubuntu-testing.img -nographic \
  -net nic,model=e1000,netdev=n0 \
  -netdev tap,id=n0,script=no,ifname=tap0
```


virtio NIC:

```
qemu-system-x86_64 \
  --enable-kvm -m 2048 -k en-us --smp 2 \
  -boot d ubuntu-testing.img -nographic \
  -net nic,model=virtio,netdev=n0 \
  -netdev tap,id=n0,script=no,ifname=tap0
```

Inside the Linux VM use the following steps to install redis:

```
sudo apt install vim git build-essential libjemalloc
git clone https://github.com/antirez/redis.git
cd redis
git checkout 3.0.6
make
```

Finally, start redis:

```
cd redis/src
rm dump.rdb && ./redis-server
```

Some approximate numbers to expect on a Linux VM and

e1000, no pipeline:

- SET 50k req/s
- GET 50k req/s

virtio, -P 29 :

- SET 1.8M req/s
- GET 2.0M req/s

Redis on the rumprun unikernel

Install the toolchain:

```
git clone https://github.com/rumpkernel/rumprun.git
cd rumprun
# Rumprun install
git submodule update --init
./build-rr.sh hw -- -F CFLAGS='-w'
. "/root/rumprun/./obj-amd64-hw/config-PATH.sh"
```

Build the redis unikernel:

```
# Packages install
git clone https://github.com/gz/rumprun-packages.git
cd rumprun-packages

cp config.mk.dist config.mk
vim config.mk

cd redis
make -j8
rumprun-bake hw_generic redis.bin bin/redis-server
```

Run the unikernel

```
# Run using virtio
rumprun kvm -i -M 256 -I if,vioif,'-net tap,script=
-W if,inet,dhcp -b images/data.iso,/data -- redis
# Run using e1000
rumprun kvm -i -M 256 -I if,wm,'-net tap,ifname=tap0
nic,model=e1000' -W if,inet,dhcp -b images/data.iso
```

Run the benchmark

```
redis-benchmark -h 172.31.0.10
```

Approximate numbers to expect:

- virtio: PING ~100k req/s
- e1000 PING ~30k req/s

Benchmarking Memcached

Yet another key--value store written in C, but compared to

Automated integration test

The easiest way to run memcached on nrk, is to invoke the

```
cd kernel
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_
```

This test will spawn memcached on one, two and four threads and measure latency with [memaslap](#).

Launch memcached manually

Start the server binary on the VM instance:

```
cd kernel
python3 run.py \
  --kfeatures test-userspace-smp \
  --cmd 'log=info init=memcached.bin' \
  --nic virtio \
  --mods rkapps \
  --qemu-settings='-m 1024M' \
  --ufeatures 'rkapps:memcached' \
  --release \
  --qemu-cores 4 \
  --verbose
```

As usual, make sure `dhcpd` is running on the host:

```
cd kernel
sudo service apparmor teardown
sudo dhcpd -f -d tap0 --no-pid -cf ./tests/dhcpd.conf
```

Start the load-generator on the host:

```
memaslap -s 172.31.0.10 -t 10s -S 10s
```

memaslap: Load generator

memaslap measures throughput and latency of a memca this:

```
memaslap -s 172.31.0.10:11211 -B -S 1s
```

Explanation of arguments:

- `-B` : Use the binary protocol (faster than the ASCII version)
- `-S 1s` : Dump statistics every X seconds

The other default arguments the tool assumes are:

- 8 client threads with concurrency of 128 sockets
- 1000000 requests
- SET proportion: 10%
- GET proportion: 90%

Unfortunately, the memaslap binary does not come with the source code. Follow the [steps in the CI guide](#) to install it from source

LevelDB

And, yet another key--value store written in C, but this one (unlike and [memcached](#) or [Redis](#)).

Automated integration test

The easiest way to run LevelDB on nrk, is to invoke the int

```
cd kernel
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_
```

This test will run the db-bench binary for LevelDB which runs the LevelDB process. Our test is configured to create a database size of 64 KiB, and then perform 100k random lookups. This is increasing the amount of cores/threads.

Launch dbbench manually

An example invocation to launch the db-bench binary directly

```
python3 run.py --kfeatures test-userspace-smp \
  --cmd "log=info init=dbbench.bin initargs=28 :
benchmarks=fillseq,readrandom --reads=100000 --num
--nic virtio --mods rkapps --ufeatures rkapps:
--release --qemu-cores 28 --qemu-nodes 2 --qem
--qemu-affinity --qemu-prealloc
```

dbbench example output

Running dbbench should ideally print an output similar to

```

LevelDB:      version 1.18
Keys:         16 bytes each
Values:       100 bytes each (50 bytes after compress)
Entries:      1000000
RawSize:      110.6 MB (estimated)
FileSize:     62.9 MB (estimated)
WARNING: Snappy compression is not enabled
-----

```

```

fillseq      :      1.810 micros/op;   61.1 MB/s
fillsync     :      0.000 micros/op;    inf MB/s
fillrandom   :      1.350 micros/op;   81.9 MB/s
overwrite    :      2.620 micros/op;   42.2 MB/s
readrandom   :      0.000 micros/op; (1000000 of
readrandom   :     10.440 micros/op; (1000000 of
readseq      :      0.206 micros/op;  537.4 MB/s
readreverse  :      0.364 micros/op;  303.7 MB/s
compact      :    60000.000 micros/op;
readrandom   :      2.100 micros/op; (1000000 of
readseq      :      0.190 micros/op;  582.1 MB/s
readreverse  :      0.301 micros/op;  367.7 MB/s
fill100K     :     390.000 micros/op;  244.6 MB/s
crc32c       :      5.234 micros/op;   746.3 MB/s
snappycomp   :      0.000 micros/op; (snappy fail
snappyuncomp :      0.000 micros/op; (snappy fail
acquireload  :      0.000 micros/op; (each op is

```

Build steps

Some special handling is currently encoded in the build-pdbench is a C++ program and C++ uses libunwind. However Rust also uses libunwind and this leads to duplicate symbols. The C++ toolchain provides it (the non-hacky solution would provide unwind symbols).

Implications:

- We have a `-L${RUMPRUN_SYSROOT}/../../obj-amd64` LevelDB Makefile (`$CXX` variable)
- We pass `-Wl,-allow-multiple-definition` to `rump` now defined twice (vibrio and NetBSD unwind lib)

See code in `usr/rkapps/build.rs` which adds flag for this

If you'll ever find yourself in a situation where you need to not necessary except when debugging build), you can use

```
cd nrk/target/x86_64-nrk-none/<release | debug>/bu
export PATH=`realpath ../../../../rumpkernel-$HASH/ou

RUMPRUN_TOOLCHAIN_TUPLE=x86_64-rumprun-netbsd make
RUMPRUN_TOOLCHAIN_TUPLE=x86_64-rumprun-netbsd make
RUMPSHAKE_ENV="-wL,-allow-multiple-definition" RUMPRUN
rumprun-netbsd rumprun-bake nrk_generic ../../../../.
```

You might also want to delete the `rm -rf build` and target of the Makefile if you want to call `clean` to recom

Run LevelDB on the rumprun unikernel

Build unikernel:

```
git clone https://github.com/rumpkernel/rumprun.git
cd rumprun
./build-rr.sh hw -- -F CFLAGS='-w'
. "/PATH/TO/config-PATH.sh"
```

Build LevelDB:

```
# Packages install
git clone https://github.com/gz/librettos-packages
cd librettos-packages/leveldb

RUMPRUN_TOOLCHAIN_TUPLE=x86_64-rumprun-netbsd make
RUMPRUN_TOOLCHAIN_TUPLE=x86_64-rumprun-netbsd make
RUMPRUN_TOOLCHAIN_TUPLE=x86_64-rumprun-netbsd rumprun
bin/db_bench
```

Run it in a VM:

```
rm data.img
mkfs.ext2 data.img 512M
rumprun kvm -i -M 1024 -g '-nographic -display cui
TEST_TMPDIR=/data dbbench.bin
```

Artifact Evaluation

Thank you for your time and picking our paper for the art

This file contains the steps to run experiments used in our Replication and Sharing in an Operating System.

All the experiments run smoothly on the c6420 CloudLab issues if one tries to run the experiments on some other configuration.

Reserve a cloudlab machine

Please follow the given steps to reserve a machine to run

1. Setup an account on [CloudLab](#), if not already present
2. Log in to [CloudLab](#) and setup a [password-less](#) ssh key
3. Start an experiment by clicking on `Experiments` on
4. Use the node type [c6420](#) (by entering `Optional phy` the node with the Ubuntu 20.04 disk image.

Download the code and setup the environment

Download and checkout the sources:

```
cd $HOME
git clone https://github.com/vmware-labs/node-replication
cd nrk
git checkout osdi21-ae-v2
bash setup.sh
```

Configure the lab machine

Add password-less sudo capability for your user (scripts r


```
sudo visudo
# Add the following line at the bottom of the file
$YOUR_USERNAME_HERE ALL=(ALL) NOPASSWD: ALL
```

Add yourself to the KVM group:

```
sudo adduser $USER kvm
```

Disable apparmor, an annoying security feature that blocks during testing. You can also set-up a rule to allow this on the test machine:

Most likely apparmor is not installed if you're using clo... will fail and you can ignore that.

```
sudo systemctl stop apparmor
sudo systemctl disable apparmor
sudo apt remove --assume-yes --purge apparmor
```

Unfortunately, for apparmor and kvm group changes to take

```
sudo reboot
```

Do a test run

Note: Most of our benchmarks takes a while to finish, so in a tmux session, or increase the session timeout to avoid

To check if the environment is setup properly, run

```
source $HOME/.cargo/env
cd $HOME/nrk/kernel
python3 ./run.py --release
```

The script downloads needed crates, compiles the OS and can take a few minutes).

If everything worked, you should see the following last line

```
[...]
[DEBUG] - init: Initialized logging
[DEBUG] - init: Done with init tests, if we came h
[SUCCESS]
```

Figure 3: NR-FS vs. tmpfs

Please follow the given steps to reproduce Figure 3 in the

NrFS results

To execute the benchmark, run:

```
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_
```

The command runs all NR-FS microbenchmarks and store `fxmark_benchmark.csv`. This step can take a while (~30-6

If everything worked, you should see an output like this o

```
[...]
Invoke QEMU: "python3" "run.py" "--kfeatures" "tes
"log=info initargs=32X8XmixX100" "--nic" "e1000" '
"fxmark" "--release" "--qemu-cores" "32" "--qemu-r
"49152" "--qemu-affinity"
Invoke QEMU: "python3" "run.py" "--kfeatures" "tes
"log=info initargs=32X12XmixX100" "--nic" "e1000"
"fxmark" "--release" "--qemu-cores" "32" "--qemu-r
"49152" "--qemu-affinity"
Invoke QEMU: "python3" "run.py" "--kfeatures" "tes
"log=info initargs=32X16XmixX100" "--nic" "e1000"
"fxmark" "--release" "--qemu-cores" "32" "--qemu-r
"49152" "--qemu-affinity"
ok

test result: ok. 1 passed; 0 failed; 0 ignored; 0
finished in 2769.78s
```

Linux tmpfs results

You can also generate the `tmpfs` result on Linux:

```
cd $HOME
git clone https://github.com/gz/vmopsbench
cd vmopsbench
git checkout c011854
bash scripts/run.sh
```

The above command runs the benchmark and writes the `fsops_benchmark.csv`.

Plot Figure 3

All the plot scripts are in a github repository, execute the

```
cd $HOME
git clone https://github.com/ankit-iitb/plot-scripts
```

To install the required dependencies, run:

```
cd $HOME/plot-scripts
sudo apt install python3-pip
pip3 install -r requirements.txt
```

Plot the Figure 3 by running:

```
# python3 fsops_plot.py <Linux fsops csv> <NrOS fsops csv>
python3 fsops_plot.py $HOME/vmopsbench/fsops_benchmark.csv
$HOME/nrk/kernel/fxmark_benchmark.csv
```

Arguments given in the plot scripts assume that the repository is cloned to the current directory. Please use the argument order given in the command for some reason.

Figure 4: LevelDB

Figure 4 in the paper compares LevelDB workload performance on Linux and NrOS.

LevelDB on NrOS

To run the LevelDB benchmark on NrOS execute:

```
cd $HOME/nrk/kernel
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_
```

This step will take ~15-20min. If everything worked, you s
the end:

```
[...]
Invoke QEMU: "python3" "run.py" "--kfeatures" "tes
"log=info init=dbbench.bin initargs=32 appcmd='\--
benchmarks=fillseq,readrandom --reads=100000 --nur
nic" "virtio" "--mods" "rkapps" "--ufeatures" "rk
"--qemu-cores" "32" "--qemu-nodes" "2" "--qemu-mer
"--qemu-prealloc"
readrandom    : done: 3200000,  949492.348 ops/sec;
ok

test result: ok. 1 passed; 0 failed; 0 ignored; 0
finished in 738.67s
```

The command runs benchmarks and stores the results in

LevelDB on Linux

To run the LevelDB benchmark on Linux follow the steps
repository in a different path than NrOS.

```
cd $HOME
git clone https://github.com/amytai/leveldb.git
cd leveldb
git checkout 8af5ca6
bash run.sh
```

The above commands run the benchmarks and writes the
linux_leveldb.csv.

Plot the LevelDB figure

Make sure that steps to download the plot scripts and ins
already been performed as explained in [Plot Figure 3](#) bef

Run the following commands to plot the Figure 4.

```
cd $HOME/plot-scripts
# python3 leveldb_plot.py <Linux leveldb csv> <NrOS leveldb csv>
python3 leveldb_plot.py $HOME/levelldb/linux_levelldb.csv
$HOME/nrk/kernel/levelldb_benchmark.csv
```

Figure 5 / 6a / 6c

Figure 5 in the paper compares address-space insertion throughput on NrOS with Linux.

NR-VMem

To run the throughput benchmark (Figure 5) on NrOS execute the following:

```
cd $HOME/nrk/kernel
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_throughput
```

This step will take ~3min. If everything worked, you should see the following output:

```
Invoke QEMU: "python3" "run.py" "--kfeatures" "test"
"log=info initargs=32" "--nic" "e1000" "--mods" "mmio"
"--release" "--qemu-cores" "32" "--qemu-nodes" "2"
qemu-affinity"
ok
```

```
test result: ok. 1 passed; 0 failed; 0 ignored; 0
finished in 118.94s
```

The results will be stored in `vmops_benchmark.csv`.

To run the latency benchmark (Figure 6a) on NrOS execute the following:

```
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_latency
nocapture
```

This step will take ~2min. If everything worked, you should see the following output:

```
Invoke QEMU: "python3" "run.py" "--kfeatures" "tes
"log=info initargs=32" "--nic" "e1000" "--mods" "-
vmops,latency" "--release" "--qemu-cores" "32" "--
"32768" "--qemu-affinity"
ok
```

```
test result: ok. 1 passed; 0 failed; 0 ignored; 0
finished in 106.67s
```

The results will be stored in `vmops_benchmark_latency.c`

To run the unmap latency benchmark (Figure 6c) on NrOS

```
cd $HOME/nrk/kernel
RUST_TEST_THREADS=1 cargo test --test s10* -- s10_
--nocapture
```

This step will take ~2min. If everything worked, you should end:

Be aware unmap latency numbers might be impacted

```
Invoke QEMU: "python3" "run.py" "--kfeatures" "tes
"log=info initargs=32" "--nic" "e1000" "--mods" "-
unmaplat,latency" "--release" "--qemu-cores" "32"
memory" "32768" "--qemu-affinity"
ok
```

```
test result: ok. 1 passed; 0 failed; 0 ignored; 0
finished in 97.38s
```

The results will be stored in `vmops_unmaplat_benchmark_`

Linux VMem

To run the benchmark on Linux follow the steps below.

```
cd $HOME/vmopsbench
git checkout master
bash scripts/linux.bash throughput
bash scripts/linux.bash latency
bash scripts/linux-tlb.bash latency
```

The results for Figure 5, 6a, and 6c will be store in:

- Figure 5 in `vmops_linux_maponly-isolated-shared_threads_all_throughput_results.csv`
- Figure 6a in `Linux-Map_latency_percentiles.csv`
- Figure 6c in `Linux-Unmap_latency_percentiles.csv`

Plot Figure 5 and 6a and 6c

Go to the plot-scripts repository:

```
cd $HOME/plot-scripts
```

Plot Figure 5:

```
# python3 vmops_throughput_plot.py <linux vmops <
python3 vmops_throughput_plot.py $HOME/vmopsbench
shared_threads_all_throughput_results.csv $HOME/nr
```

Plot Figure 6a:

```
# python3 map_latency_plot.py <linux map-latency <
python3 map_latency_plot.py $HOME/vmopsbench/Linux-Map-
latency_percentiles.csv $HOME/nrk/kernel/vmops_benchmark_latency.csv
```

Plot Figure 6c:

```
# python3 mapunmap_latency_plot.py <linux unmap-latency <
csv>
python3 mapunmap_latency_plot.py $HOME/vmopsbench/Linux-Unmap-
latency_percentiles.csv $HOME/nrk/kernel/vmops_unmaplat_benchmark_latency.csv
```

Baseline Operating System

Contains steps to get other operating systems compiled a

Compare against Linux

To get an idea if nrk is competitive with Linux performance, we can create an image. The following steps create an `ubuntu-minimal` installer:

```
wget http://archive.ubuntu.com/ubuntu/dists/bionic/
amd64/current/images/netboot/mini.iso
qemu-img create -f vmdk -o size=20G ubuntu-testing
kvm -m 2048 -k en-us --smp 2 --cpu host -cdrom mir
# Follow installer instructions
```

Afterwards the image can be booted using `kvm`:

```
kvm -m 2048 -k en-us --smp 2 -boot d ubuntu-testir
```

Switch to serial output

One step that makes life easier is to enable to serial input graphical QEMU interface. To enable serial, edit the grub file that follows in the VM:

```
GRUB_CMDLINE_LINUX_DEFAULT=""
GRUB_TERMINAL='serial console'
GRUB_CMDLINE_LINUX="console=tty0 console=ttyS0,115200"
GRUB_SERIAL_COMMAND="serial --speed=115200 --unit=0"
```

Then you must run `update-grub` to update the menu entries in the VM using (not the `-nographic` option):

```
qemu-system-x86_64 --enable-kvm -m 2048 -k en-us -
testing.img -nographic
```


Compare against Barrelfish

TBD.

Compare against sv6

To clone & build the code (needs an older compiler versio

```
git clone https://github.com/aclements/sv6.git
sudo apt-get install gcc-4.8 g++-4.8
CXX=g++-4.8 CC=gcc-4.8 make
```

Update `param.h`:

```
QEMU      ?= qemu-system-x86_64 -enable-kvm
QEMUSMP   ?= 56
QEMUMEM   ?= 24000
```

Run:

```
CXX=g++-4.8 CC=gcc-4.8 make qemu`
```

Rackscale

One of the baselines for rackscale is NrOS. To run the rac
corresponding NrOS baselines, run them with `--feature |`

Environment

This chapter contains various notes on configuration and the hypervisor (QEMU) to use either pass-through or emu nrkernel and develop for it.

Install QEMU from sources

Make sure the QEMU version for the account is is ≥ 6 . If build it from scratch, if it the Ubuntu release has a lesser

First, make sure to uncomment all `#deb-src` lines in `/etc/apt/sources.list` uncommented. Then, run the following commands:

For any build:

```
sudo apt update
sudo apt install build-essential libpmem-dev libd
apt source qemu
sudo apt build-dep qemu
```

For non-rackscale:

```
wget https://download.qemu.org/qemu-6.0.0.tar.xz
tar xvJf qemu-6.0.0.tar.xz
cd qemu-6.0.0
```

For non-rackscale OR rackscale:

```
git clone https://github.com/hunhoffe/qemu.git qer
cd qemu
git checkout --track origin/dev/ivshmem-numa
```

For any build:

```
./configure --enable-rdma --enable-libpmem
make -j 28
sudo make -j28 install
sudo make rdmactm-mux
```

```
# Check version (should be  $\geq 6.0.0$ )
qemu-system-x86_64 --version
```

You can also add `--enable-debug` to the configure script (useful for source information when stepping through qe

Note that sometimes `make install` doesn't actually repl
`ivshmem-server` to find the current location and then ove
`qemu/build/contrib/ivshmem-server/ivshmem-server` .

Use NVDIMM in QEMU

tldr: The `--qemu-pmem` option in `run.py` will add persistence. If you want to customize further, read on.

Qemu has support for NVDIMM that is provided by a memory backend-ram. A simple way to create a vNVDIMM device is using the following command-line options:

```
-machine pc,nvdimmm
-m $RAM_SIZE,slots=$N,maxmem=$MAX_SIZE
-object memory-backend-file,id=mem1,share=on,memdev=mem1
-device nvdimmm,id=nvdimmm1,memdev=mem1
```

Where,

- the `nvdimmm` machine option enables vNVDIMM feature.
- `slots=$N` should be equal to or larger than the total number of vNVDIMM devices, e.g. `$N` should be `>= 2` here.
- `maxmem=$MAX_SIZE` should be equal to or larger than the total memory of vNVDIMM devices.
- `object memory-backend-file,id=mem1,share=on,memdev=mem1` creates a backend storage of size `$NVDIMM_SIZE`.
- `share=on/off` controls the visibility of guest writes. If `share=on`, multiple guests will be visible to each other.
- `device nvdimmm,id=nvdimmm1,memdev=mem1` creates a vNVDIMM device whose storage is provided by above memory backend.

Guest Data Persistence

Though QEMU supports multiple types of vNVDIMM backends, the one that can guarantee the guest write persistence is:

- DAX device (e.g., `/dev/dax0.0`, ...) or
- DAX file (mounted with `dax` option)

When using DAX file (A file supporting direct mapping of p write persistence is guaranteed if the host kernel has sup mmap system call and additionally, both 'pmem' and 'sha backend.

NVDIMM Persistence

Users can provide a persistence value to a guest via the o command line option:

```
-machine pc,accel=kvm,nvdimmm,nvdimmm-persistence=cp
```

There are currently two valid values for this option:

`mem-ctrl` - The platform supports flushing dirty data from NVDIMMs in the event of power loss.

`cpu` - The platform supports flushing dirty data from the event of power loss.

Emulate PMEM using DRAM

Linux systems allow emulating DRAM as PMEM. These de Memory Region by the OS. Usually, these devices are fast not provide any persistence. So, such devices are used or

On Linux, to find the DRAM region that can be used as PM

```
dmesg | grep BIOS-e820
```

The viable region will have "usable" word at the end.

```
[    0.000000] BIOS-e820: [mem 0x00000000100000000-
```

This means that the memory region between 4 GiB (0x000 (0x0000000053ffffff) is usable. Say we want to reserve a 16 need to add this information to the grub configuration file

```
sudo vi /etc/default/grub
GRUB_CMDLINE_LINUX="memmap=16G!4G"
sudo update-grub2
```

After rebooting with our new kernel parameter, the `dmesg` persistent memory region like the following:

```
[    0.000000] user: [mem 0x000000001000000000-0x0000000012)
```

We will see this reserved memory range as `/dev/pmem0` . It is ready to use. Mount it with the `dax` option.

```
sudo mkdir /mnt/pmem0
sudo mkfs.ext4 /dev/pmem0
sudo mount -o dax /dev/pmem0 /mnt/pmem0
```

Use it as a `mem-path=/mnt/pmem0` as explained [earlier](#).

Configure and Provision NVDIMMs

The NVDIMMs need to be configured and provisioned before use. The Intel `ipmctl` tool can be used to discover and provision them.

To show all the NVDIMMs attached to the machine, run:

```
sudo ipmctl show -dimm
```

To show all the NVDIMMs attached on a socket, run:

```
sudo ipmctl show -dimm -socket SocketID
```

Provisioning

NVDIMMs can be configured both in volatile (MemoryMode) and non-volatile (PersistentMode) modes or a mix of two using `ipmctl` tool on Linux.

We are only interested in using the NVDIMMs in AppDirect mode. NVDIMMs can be configured in two ways; AppDirect and PersistentMode.

mode, the data is interleaved across multiple DIMMs, and AppDirectNotInterleaved is used. To configure multiple D run:

```
sudo ipmctl create -goal PersistentMemoryType=AppD
```

Reboot the machine to reflect the changes made using th creates a region on each socket on the machine.

```
ndctl show --regions
```

To show the DIMMs included in each region, run:

```
ndctl show --regions --dimms
```

Each region can be divided in one or more namespaces in the operating system. To create the namespace(s), run:

```
sudo ndctl create-namespace --mode=[raw/sector/fsc
```

The namespace can be created in different modes like raw default mode is fsdax.

Reboot the machine after creating the namespaces, and t /dev/depending on the mode. For example, if the mode is fsc /dev/pmem.

Mount these devices:

```
sudo mkdir /mnt/pmem0
sudo mkfs.ext4 /dev/pmem0
sudo mount -o dax /dev/pmem0 /mnt/pmem0
```

These mount points can be used directly in the userspace machine as explained [earlier](#).

Use RDMA support in QEMU

tldr: The `-pvrDMA` option in `run.py` will enable RDMA support. You will manually have to run `rdmacm-mux` and unload the Mellanox drivers.

QEMU has support for `pvrDMA` (a para-virtual RDMA drive cards (like Mellanox). In order to use it (aside from the `--rdmacm-mux` during building), the following steps are needed:

Install Mellanox drivers (or any other native drivers for your host).

```
wget https://content.mellanox.com/ofed/MLNX_OFED-5.4-1.0.3.0-ubuntu20.04-x86_64.tgz
tar zxvf MLNX_OFED_LINUX-5.4-1.0.3.0-ubuntu20.04-x86_64.tgz
cd MLNX_OFED_LINUX-5.4-1.0.3.0-ubuntu20.04-x86_64
./mlnxofedinstall --all
```

Before running the `rdmacm-mux` make sure that both `ib_core` and `rdma_cm` aren't loaded, otherwise the `rdmacm-mux` service will fail.

```
sudo rmmod ib_ipoib
sudo rmmod rdma_cm
sudo rmmod ib_cm
```

Start the QEMU `rdmacm-mux` utility (before launching a qemu instance):

```
./rdmacm-mux -d mlx5_0 -p 0
```


Inter-VM Communication using shared memory

tldr: Use the `--qemu-ivshmem` and `--qemu-shmem-path` to enable VM shared-memory support in QEMU.

This section describes how to use shared memory to communicate between VMs. First, create a shared memory file (with hugepages):

```
echo 1024 > /sys/kernel/mm/hugepages/hugepages-2048kB
sudo mkdir -p /mnt/hugepages
sudo mount -t hugetlbfs pagesize=2GB /mnt/hugepages
sudo chmod 777 /mnt/hugepages
```

Now, use a file on this mount point to create a shared memory file. Qemu allows two types of configuration:

- Just the shared memory file: `ivshmem-plain`.
- Shared memory plus interrupts: `ivshmem-doorbell`.

We use the plain shared memory configuration as the goal for this guide. Add the following parameters to the Qemu command line:

```
-object memory-backend-file,size=2G,mem-path=/mnt/hugepages/ivshmem-plain,share=on,id=HMB \
-device ivshmem-plain,memdev=HMB
```

Discover the shared memory file inside the VM

Qemu exposes the shared memory file to the kernel by creating a PCI device. Run the following command to discover to check if the PC

```
lspci | grep "shared memory"
```

Running `lspci` should show something like:

```
00:04.0 RAM memory: Red Hat, Inc. Inter-VM shared memory
```

Use the `lspci` command to know more about the PCI de

```
lspci -s 00:04.0 -nvv
```

This should print the BAR registers related information. There are three BARs (depending on shared memory or interrupt de

- BAR0 holds device registers (256 Byte MMIO)
- BAR1 holds MSI-X table and PBA (only ivshmem-doo
- BAR2 maps the shared memory object

Since we are using the plain shared memory configuration, we use the BAR0 and BAR2 as Region 0 and Region 1.

```
00:04.0 0500: 1af4:1110 (rev 01)
    Subsystem: 1af4:1100
    Physical Slot: 4
    Control: I/O+ Mem+ BusMaster- SpecCycle- M
Stepping- SERR+ FastB2B- DisINTx-
    Status: Cap- 66MHz- UDF- FastB2B- ParErr-
<MAbort- >SERR- <PERR- INTx-
    Region 0: Memory at febf1000 (32-bit, non-
    Region 2: Memory at 2800000000 (64-bit, pre
```

Use the shared memory file inside C

If you only need the shared memory part, BAR2 suffices. You can access the shared memory in the guest and can use it as you see fit. The shared memory is at `2800000000` with the size of `2G`.

Here is a sample C program that writes to the shared me

```

#include<stdio.h>
#include<stdint.h>
#include<unistd.h>
#include<fcntl.h>
#include<sys/mman.h>

int main() {
    void *baseaddr = (void *) 0x280000000; // BAR1
    uint64_t size = 2147483648; // BAR2 size

    int fd = open("/sys/bus/pci/devices/0000:00:04:00", O_RDWR);
    void *retaddr = mmap(baseaddr, size, PROT_READ|PROT_WRITE,
0);
    if (addr == MAP_FAILED) {
        printf("mmap failed");
        return 0;
    }

    uint8_t *addr = (uint8_t *)retaddr;
    addr[0] = 0xa;

    munmap(retaddr, size);
    close(fd);
}

```

Compile and run the program (use `sudo` to run).

Perform the similar steps to read the shared memory file

Discover CXL devices in Lin

This document aims to list out steps to discover CXL type
 Since there is no hardware available, the only way to achi
 Unfortunately, even the Qemu mainstream branch does r
 tutorial uses a custom version of Qemu that supports CXL

Build custom Qemu version

First, download and build the custom Qemu version on yo

```
sudo apt install build-essential libpmem-dev libd
cd ~/cxl
git clone https://gitlab.com/bwidawsk/qemu.git
cd qemu
git checkout cxl-2.0v4
./configure --enable-libpmem
make -j 16
```

Check the version:

```
./build/qemu-system-x86_64 --version
```

```
QEMU emulator version 6.0.50 (v6.0.0-930-g1839565)
Copyright (c) 2003-2021 Fabrice Bellard and the QEMU
```

Build custom Linux Kernel

Next, download the latest kernel version and build an ima

```
cd ~/cxl
git clone https://git.kernel.org/pub/scm/linux/kernel
cd linux
make defconfig
```

defconfig generates default configuration values and st
 kernel requires some special configuration changes to ha
 these configuration flags are present in the .config file,
 these flags.

```
CONFIG_ACPI_HMAT=y
CONFIG_ACPI_APEI_PCIEAER=y
CONFIG_ACPI_HOTPLUG_MEMORY=y
CONFIG_MEMORY_HOTPLUG=y
CONFIG_MEMORY_HOTPLUG_DEFAULT_ONLINE=y
CONFIG_MEMORY_HOTREMOVE=y
CONFIG_CXL_BUS=y
CONFIG_CXL_MEM=y
```

Once the configuration related changes are done, compile the image file.

```
make -j 16
cd ~/cxl
```

The image file should be in `linux/arch/x86_64/boot/bzImage`.

Run Qemu with CXL related parameters

Qemu provides `-kernel` parameter to use the kernel image.

```
qemu/build/qemu-system-x86_64 -kernel linux/arch/x86_64/boot/bzImage
-nographic -append "console=ttyS0" -m 1024 --enable
```

The kernel runs until it tries to find the root fs; then it crashes to create a ramdisk.

```
mkinitramfs -o ramdisk.img
```

Press `Ctrl+a` and `c` to exit kernel and then press `q` to quit Qemu.

Now, run qemu with the newly created ramdisk:

```
qemu/build/qemu-system-x86_64 -kernel linux/arch/x86_64/boot/bzImage
-append "console=ttyS0" -m 1024 -initrd ramdisk.img
```

The kernel runs properly this time. Now, it is time to add CXL support to running Qemu:

```

qemu/build/qemu-system-x86_64 -kernel linux/arch/x86/boot/bzImage \
    -append "console=ttyS0" -initrd ramdisk.img -m 1024 \
    -m 1024,slots=12,maxmem=16G -M q35,accel=kvm,cpu=host \
    -object memory-backend-file,id=cxl-mem1,share=on \
    -object memory-backend-file,id=cxl-label1,share=on \
    -object memory-backend-file,id=cxl-label2,share=on \
    -device pxb-cxl,id=cxl.0,bus=pcie.0,bus_nr=52,uid=0 \
    base[0]=0x4c0000000000,memdev[0]=cxl-mem1 \
    -device cxl-rp,id=rp0,bus=cxl.0,addr=0.0,chassis=0 \
    -device cxl-rp,id=rp1,bus=cxl.0,addr=1.0,chassis=0 \
    -device cxl-type3,bus=rp0,memdev=cxl-mem1,id=cxl-0 \
    -device cxl-type3,bus=rp1,memdev=cxl-mem1,id=cxl-1

```

Qemu exposes the CXL devices to the kernel and the kernel exposes the devices by running:

```
ls /sys/bus/cxl/devices/
```

or

```
dmesg | grep '3[45]:00'
```

References

- [Booting a custom linux kernel in Qemu](#)
- [CXL 2.0 support in Linux](#)
- [CXL 2.0 + Linux + Qemu](#)

Continuous integration (CI)

We run tests using the github-actions infrastructure. The setup involves provisioning up a new runner machine (and connect a github repo to it).

Steps to add a new CI machine:

1. [Install github-runner software on a new test machine](#)
2. [Give access to the benchmark repository](#)
3. [Configure software for the github-runner account](#)
4. [Disable AppArmor](#)
5. [Install a recent QEMU](#)
6. [Do a test-run](#)
7. [Start the new runner](#)

Install github-runner software on a

Create a github-runner user first:

```
sudo useradd github-runner -m -s /bin/zsh
```

Add sudo capability for github-runner:

```
sudo visudo
# github-runner  ALL=(ALL) NOPASSWD: ALL
```

For better security with self-hosted code execution, make sure to enable `Require approval for all actions on your repository` in the `Settings -> Runners -> Require approval for all actions on your repository` github repo settings!

Other than that, follow the steps listed under `Settings -> Runners` in the github runner :

```
sudo su github-runner
cd $HOME
<< steps from Web-UI >>
```

When asked for labels, make sure to give it a machine specific label. Use the following labels `skylake2x`, `skylake4x`, `cascade`

machine type and the number of sockets/NUMA nodes. Machines should have the same tag to allow parallel test execution.

If you add a new machine label, make sure to also add `_scripts` folder.

Don't launch the runner yet with `run.sh` (this happens fu

Give access to the benchmark repos

Benchmark results are uploaded automatically to git.

Generate a key for accessing the repository or use an existing account. Also add the user to the KVM group. Adding you logout/reboot which we do in later steps.

```
sudo adduser github-runner kvm
ssh-keygen
```

Then, add the pub key (`.ssh/id_rsa.pub`) to the github C

Configure software for the github ru

Install necessary software for use by the runner:

```
git clone git@github.com:vmware-labs/node-replicat
cd nrk/
bash setup.sh
source $HOME/.cargo/env
```

Install a recent qemu

[Follow the steps in the Environment chapter.](#)

Install memaslap

The memcached benchmark uses the `memaslap` binary that is not included in the Ubuntu `libmemcached-tools` deb package, so you have to get it from the sources:

```
cd $HOME
sudo apt-get build-dep libmemcached-tools
wget https://launchpad.net/libmemcached/1.0/1.0.18/+download/libmemcached-1.0.18.tar.gz
tar zxvf libmemcached-1.0.18.tar.gz

cd libmemcached-1.0.18/
LDFLAGS='-lpthread' CPPFLAGS='-fcommon -fpermissive' ./configure --enable-memaslap
CPPFLAGS='-fcommon' make -j12
sudo make install
sudo ldconfig

which memaslap
```

Disable AppArmor

An annoying security feature that blocks our DHCP server. You can add a rule for allowing this but it's easiest to just get rid of it.

```
sudo systemctl stop apparmor
sudo systemctl disable apparmor
sudo apt remove --assume-yes --purge apparmor
# Unfortunately for apparmor and kvm group changes you need to
reboot:
sudo reboot
```

Do a test-run

After the reboot, verify that the nrk tests pass (this will take a while to succeed too):

```
# Init submodules if not done so already:
cd nrk
git submodule update --init
source $HOME/.cargo/env

cd kernel
RUST_TEST_THREADS=1 cargo test --features smoke --
```

Start the runner

Finally, launch the runner:

```
cd $HOME/actions-runner
source $HOME/.cargo/env
./run.sh
```

Start runner as systemd service

```
cd $HOME/actions-runner
sudo ./svc.sh install
sudo ./svc.sh start
```

Check the runner status with:

```
sudo ./svc.sh status
```

Stop the runner with:

```
sudo ./svc.sh stop
```

Uninstall the service with:

```
sudo ./svc.sh uninstall
```

Repository settings

If the repo is migrated to a new location, the following set

1. Under Settings -> Secrets: Add secret `WEBSITE_DEPLOY_KEY` and push the generated documentation to the correct repository
2. Under Settings -> Options: Disable "Allow merge commits"
3. Under Settings -> Branches: Add branch protection for `main` with the following settings:
 - Require pull request reviews before merging
 - Dismiss stale pull request approvals when new approvals are added to commit
 - Require status checks to pass before merging
 - Require branches to be up to date before merging
 - Require linear history
4. Under Settings -> Actions -> Runners:
 - "Require approval for all outside collaborators"

Related Work

NRK takes inspiration from decades of academic research exhaustive:

Operating Systems

- [Barrelfish](#)
- [Tornado](#)
- [K42](#)
- [sv6](#)
- [Rumpkernel](#)
- [LibrettOS](#)
- [seL4](#)
- [Disco](#)
- [Mitosis](#)

Scalable Data structures

- [Flat combining](#)
- [Read-Log-Update and RLU with multi-versioning](#)
- [Predictive Log Synchronization](#)
- [OpLog](#)

Log based designs

- [ScaleFS](#)
- [Corfu](#)
- [Raft](#)

Contributors

Here is a list of the contributors who helped to build NRK

- [Amy Tai](#)
- [Ankit Bhardwaj](#)
- [Chinmay Kulkarni](#)
- [Christian Menges](#)
- [Erika Hunhoff](#)
- [Gerd Zellweger](#)
- [Irina Calciu](#)
- [Reto Achermann](#)
- [Ryan Stutsman](#)
- [Sanidhya Kashyap](#)
- [Stanko Novakovic](#)
- [Zack McKevitt](#)

If you feel you're missing from this list, feel free to add yo