

FlowPools: Lock-Free Deterministic Concurrent Dataflow Queues

Authors

EPFL, Switzerland
firstname.lastname@epfl.ch
<http://lamp.epfl.ch>

Abstract. Implementing correct and deterministic parallel programs is challenging. Even though concurrency constructs exist in popular programming languages to facilitate the task of deterministic parallel programming, they are often too low level, or do not compose well due to underlying blocking mechanisms. In this paper, we present the design and implementation of a fundamental data structure for composable deterministic parallel dataflow computation through the use of functional programming abstractions. Additionally, we provide a correctness proof, showing that the implementation is linearizable, lock-free, and deterministic. Finally, we show experimental results which compare our *FlowPool* against corresponding operations on other concurrent data structures, and show that in addition to offering new capabilities, FlowPools reduce insertion time by 49 – 54% on a 4-core i7 machine with respect to comparable concurrent queue data structures in the Java standard library.

Keywords: dataflow, concurrent data-structure, deterministic parallelism

1 Introduction

Multicore architectures have become ubiquitous— even most mobile devices now ship with multiple core processors. Yet parallel programming has yet to enter the daily workflow of the mainstream developer. One significant obstacle is an undesirable choice programmers must often face when solving a problem that could greatly benefit from leveraging available parallelism. Either choose a non-deterministic, but performant, data structure or programming model, or sacrifice performance for the sake of clarity and correctness.

Programming models based on *dataflow* [1, 2] have the potential to simplify parallel programming, since the resulting programs are deterministic. Moreover, dataflow programs can be expressed more declaratively than programs based on mainstream concurrency constructs, such as shared-memory threads and locks, as programmers are only required to specify data and control dependencies. This allows one to reason sequentially about the intended behavior of their program, meanwhile enabling the underlying framework to effectively extract parallelism.

In this paper, we present the design and implementation of FlowPools, a fundamental dataflow collections abstraction which can be used as a building

block for larger and more complex *deterministic* and parallel dataflow programs. Our FlowPool abstraction is backed by an efficient non-blocking data structure, which we go on to prove is lock-free.

As a result, our data structure benefits from the increased robustness provided by lock-freedom [9], since its operations are unaffected by thread delays. We provide a lock-freedom proof, which guarantees progress regardless of the behavior, including the failure, of concurrent threads.

In combining lock-freedom with a functional interface, we go on to show that FlowPools are *composable*. That is, using prototypical higher-order functions such as `foreach` and `aggregate`, one can concisely form dataflow graphs, in which associated functions are executed asynchronously in a completely non-blocking way, as elements of FlowPools in the dataflow graph become available.

Finally, we show that FlowPools are able to overcome practical issues, such as out-of-memory errors, thus enabling programs based upon FlowPools to run indefinitely. By using a *builder* abstraction, instead of something like iterators or streams (which can lead to non-determinism) we are able to garbage collect parts of the data structure we no longer need, thus reducing memory consumption.

Our contributions are the following:

1. The design and Scala implementation¹ of a parallel dataflow abstraction and underlying data structure that is deterministic, lock-free, & composable.
2. Proofs of lock-freedom, linearizability, and determinism.
3. Detailed benchmarks comparing the performance of our FlowPools against other popular concurrent data structures.

2 Model of Computation

FlowPools can be thought of as similar to a typical collections abstraction. That is, operations invoked on a FlowPool are executed on its individual elements. However, FlowPools do not only act as a data container full of elements—unlike a typical collection, FlowPools also act as nodes and edges of a directed acyclic computation graph (DAG), in which the operations to be performed are registered with the FlowPool.

Nodes in this directed acyclic graph are data containers which are first class values. This makes it possible to use FlowPools as function arguments or to receive them as return values. Edges, on the other hand, can be thought of as combinators or higher-order functions whose associated functions are the previously-mentioned operations that are registered with the FlowPool. In addition to providing composability, this means that the DAG does not have to be specified at compile time, but can be generated dynamically at run time instead.

This structure allows for complete asynchrony, allowing the runtime to extract parallelism as a result. That is, elements can be asynchronously inserted, all registered operations can be asynchronously executed, and new operations can be asynchronously registered.

¹ <http://ASSEMBLA-REPO-URL-HERE>

$t ::=$	terms	
	$p \ll v$	append
	create p	pool creation
	p foreach f	foreach
	p seal n	seal
	$t_1 ; t_2$	sequence
	$p \in \{(vs, \sigma, cbs) \mid vs \subseteq Elem, \sigma \in \{-1\} \cup \mathbb{N}, cbs \subset Elem \Rightarrow Unit\}$	
	$f \in Elem \Rightarrow Unit$	
		$v \in Elem$ $n \in \mathbb{N}$

Fig. 1. Syntax

Therefore, invoking several higher-order functions in succession on a given FlowPool doesn't add barriers between nodes in the DAG, it only extends the DAG. This means that individual elements can *flow* through different edges of the DAG independently of each other.

Properties of FlowPools. In our model, FlowPools have a few important properties which work together to ensure that resulting programs are deterministic.

1. Data elements in FlowPools are unordered.
2. Operations passed to FlowPools must be pure and associative.
3. All callbacks are eventually executed on all elements.

Determinism. FlowPools are deterministic in the sense that all execution schedules either lead to non-termination (*e.g.*, an exception), or no difference can be observed in the final state of the resulting data structures. For a more formal definition and proof, see section 5.

3 Programming Interface

A FlowPool can be viewed as a concurrent pool data structure, and as such has no guarantees on ordering. In this section, we describe the semantics of a handful of functional combinators and other basic operations defined on FlowPools.

Append (\ll). The most fundamental of all operations on FlowPools is the append operation. As its name suggests, it simply takes an argument of type `Elem` and appends it to a given FlowPool.

Foreach and Aggregate. A pool containing a set of elements is of little use if the elements cannot be manipulated in some manner. One of the basic data structure operations is element traversal, often provided by iterators or streams – stateful objects storing the current position in the data structure. Since they can be manipulated by several threads, they allow nondeterministic execution.

Another way to traverse the elements is to provide a higher-level `foreach` operation which takes a user-specified function as an argument and applies it to every element. To ensure determinism, it is called for every element that is eventually in the FlowPool, rather than only those present when `foreach` is called. To ensure non-blocking behaviour, it must not wait until additional elements

are added to the FlowPool. For these reasons, the `foreach` operation must execute asynchronously and is eventually applied to every element. Its signature is `def foreach[U](f: T => U): Future[Int]`. The return type `Future[Int]` is an integer value which becomes available once `foreach` traverses all the elements added to the pool. It denotes the number of times the `foreach` has been called. The user-specified function return value of type `U` is ignored.

The `aggregate` operation aggregates the elements of the pool and has the following signature: `def aggregate[S](zero: =>S)(cb: (S, S) => S)(op: (S, T) => S): Future[S]`, where `zero` is the initial aggregation, `cb` is an associative operator which combines several aggregations, `op` is an operator that adds an element to the aggregation, and `Future[S]` is the final aggregation of all the elements which becomes available once all the elements have been added. The `aggregate` operator divides elements into subsets and applies the aggregation operator `op` to aggregate elements in each subset starting from the `zero` aggregation, and then combines aggregations from different subsets with the `cb` operator. In essence, the first part of `aggregate` defines the commutative monoid and the functions involved must be non-side-effecting. In contrast, the operator `op` is guaranteed to be called only once for each element and it can have side effects.

While in an imperative programming model `foreach` and `aggregate` are equivalent in the sense that one can be implemented in terms of the other, in a single-assignment programming model `aggregate` is more expressive than `foreach`. The `foreach` operation can be implemented using `aggregate`, but not vice versa, due to the absence of side effects.

Builders. The FlowPool described so far must maintain a reference to all the elements at all times to implement the `foreach` operation correctly. Since elements are never removed, the pool may grow indefinitely and run out of memory.

At the same time, appending new elements does not require a reference to any of the existing elements. We use this observation to factor out the `<<` operation into a different abstraction called a `builder`. A typical application starts by registering all the `foreach` operations and then releases the references to FlowPools. In a managed environment the GC then implicitly takes care of discarding the no longer needed objects.

Seal. When clients are sure that no more elements will be added to the pool they can disallow further appends by calling `seal`. This has the advantage of discarding the registered `foreach` operations. More importantly, `aggregate` can only complete its resulting future once it is certain that no more elements will be added.

A `seal` which just closes the FlowPool at the moment when it is called, however, yields a nondeterministic programming model. Imagine a thread that attempts to seal the pool executing concurrently with a thread that appends an element. In one execution, the append precedes the seal, and in the other the append follows the seal causing an error. To avoid nondeterminism, there has to be an agreement on the state of the pool when it is sealed. A convenient and

sufficient way to make `seal` deterministic is to provide the expected pool size as an argument. The semantics of `seal` is then such that it fails if the pool is already sealed with a different size or the number of elements is greater than the desired size.

Higher-level operations. We now show how the basic abstractions above can be used to build higher-level abstractions. To begin with, in addition to a default `FlowPool` constructor, it is convenient to have generators which create certain types of pools. In a dataflow graph `FlowPools` created using generators are source nodes. As a simple example, `tabulate` (see below) creates a sequence of elements by applying a user-specified function `f` to each natural number. One can imagine more complicated generators, adding elements from a network socket or a file.

```
def tabulate[T]
  (n: Int, f: Int => T)
  val p = new FlowPool[T]
  val b = p.builder
  def recurse(i: Int) {
    b << f(i)
    if i < n recurse(i + 1)
  }
  future { recurse(0) }
  p

def map[S](f: T => S)
  val p = new FlowPool[S]
  val b = p.builder
  for (x <- this) {
    b << f(x)
  } map {
    sz => b.seal(sz)
  }
  p

def foreach[U](f: T => U)
  aggregate(0)(_ + _) {
    (acc, x) =>
      f(x)
      acc + 1
  }
```

The `tabulate` generator starts by creating a `FlowPool` of an arbitrary type `T` and creating its builder instance. It then starts an asynchronous computation using the `future` construct, which recursively applies `f` to each number and appends it to the builder. Finally, the reference to the pool is returned.

A typical higher-level collection operation `map` is used to map each element of a dataset to produce a new one. This corresponds to chaining the nodes in a dataflow graph. We implement it by traversing the elements of the receiver `FlowPool` (`this`) and appending each mapped element to the builder. The `for` loop is syntactic sugar for calling the `foreach` method on `this`. We assume that the `foreach` return type `Future[Int]` has `map` and `flatMap` operations, which are executed once the future value becomes available. The `Future.map` above ensures that if the current pool (`this`) is ever sealed, the mapped pool is sealed to the appropriate size.

As argued before, `foreach` can be expressed in terms of `aggregate` by accumulating the number of elements and invoking the callback `f` each time. We further show that some patterns cannot be expressed in terms of a mere `foreach`. The `filter` combinator filters out the elements for which a specified predicate does not hold. Appending the elements to a new pool can proceed as before, but the `seal` needs to know the exact number of elements added. The `aggregate` accumulator is thus used to track the number of added elements.

```

type Terminal {
  sealed: Int
  callbacks: List[Elem => Unit]
}
type Elem

type Block {
  array: Array[Elem]
  next: Block
  index: Int
  blockindex: Int
}

type FlowPool {
  start: Block
  current: Block
}
LASTELEMPOS = BLOCKSIZE - 2
NOSEAL = -1

```

Fig. 2. FlowPool data-types

```

def filter
(pred: T => Boolean)
val p = new FlowPool[T]
val b = p.builder
aggregate(0)(_ + _) {
  (acc, x) => if pred(x) {
    b << x
    1
  } else 0
} map { sz => b.seal(sz) }
p

def flatMap[S]
(f: T => FlowPool[S])
val p = new FlowPool[S]
val b = p.builder
aggregate(future(0))(add) {
  (af, x) =>
  val sf = for (y <- f(x))
    b << y
  add(af, sf)
} map { sz => b.seal(sz) }
p

def union[T]
(that: FlowPool[T])
val p = new FlowPool[T]
val b = p.builder
val f = for (x <- this) b << x
val g = for (y <- that) b << y
for (s1 <- f; s2 <- g)
  b.seal(s1 + s2)
p

def add(f: Future[Int], g: Future[Int]) =
  for (a <- f; b <- g) yield a + b

```

The `flatMap` operation retrieves a pool for each element of `this` pool and adds its elements to the resulting pool. Given two `FlowPools`, it can be used to generate the cartesian product of their elements. The implementation is similar to that of `filter`, but we reduce the size on the future values of the sizes, because each intermediate pool may not yet be sealed. The operation `q union r` produces a new pool which has elements of both pool `q` and pool `r`.

The last two operations correspond to joining nodes in the dataflow graph. Note that if we could somehow merge the two different `foreach` loops to implement the third join type `zip`, we would obtain a nondeterministic operation. The programming model does not allow us to do so, however. The `zip` function is more suited for data structures with deterministic ordering, such as Oz streams – which would in turn have a nondeterministic `union`.

4 Implementation

We now describe the `FlowPool` and its basic operations. In doing so, we omit the details not relevant to the algorithm² and focus on a high-level description of a non-blocking data structure. One straightforward way to implement a growing pool is to use a linked list of nodes that wrap elements. As we are concerned about the memory footprint and cache-locality, we store the elements into arrays instead, which we call blocks. Whenever a block becomes full, a new block is allocated and the previous block is made to point to the `next` block. This way, most writes amount to a simple array-write, while allocation occurs only

² Specifically the builder abstraction and the `aggregate` operation. The `aggregate` can be implemented using `foreach` with a side-effecting accumulator.

```

1 def create()
2   new FlowPool {
3     start = createBlock(0)
4     current = start
5   }
6
7 def createBlock(bidx: Int)
8   new Block {
9     array = new Array(BLOCKSIZE)
10    index = 0
11    blockindex = bidx
12    next = null
13  }
14
15 def append(elem: Elem)
16   b = READ(current)
17   idx = READ(b.index)
18   nexto = READ(b.array(idx + 1))
19   curo = READ(b.array(idx))
20   if check(b, idx, curo) {
21     if CAS(b.array(idx + 1), nexto, curo) {
22       if CAS(b.array(idx), curo, elem) {
23         WRITE(b.index, idx + 1)
24         invokeCallbacks(elem, curo)
25       } else append(elem)
26     } else append(elem)
27   } else {
28     advance()
29     append(elem)
30   }
31
32 def check(b: Block, idx: Int, curo: Object)
33   if idx > LASTELEMPPOS return false
34   else curo match {
35     elem: Elem =>
36       return false
37     term: Terminal =>
38       if term.sealed = NOSEAL return true
39       else {
40         if totalElems(b, idx) < term.sealed
41           return true
42         else error("sealed")
43       }
44     null =>
45       error("unreachable")
46   }
47
48 def advance()
49   b = READ(current)
50   idx = READ(b.index)
51   if idx > LASTELEMPPOS
52     expand(b, b.array(idx))
53   else {
54     obj = READ(b.array(idx))
55     if obj is Elem WRITE(b.index, idx + 1)
56   }
57
58 def expand(b: Block, t: Terminal)
59   nb = READ(b.next)
60   if nb is null {
61     nb = createBlock(b.blockindex + 1)
62     nb.array(0) = t
63     if CAS(b.next, null, nb)
64       expand(b, t)
65   } else {
66     CAS(current, b, nb)
67   }
68 def totalElems(b: Block, idx: Int)
69   return b.blockindex * (BLOCKSIZE - 1) + idx
70
71 def invokeCallbacks(e: Elem, term: Terminal)
72   for (f <- term.callbacks) future {
73     f(e)
74   }
75
76 def seal(size: Int)
77   b = READ(current)
78   idx = READ(b.index)
79   if idx <= LASTELEMPPOS {
80     curo = READ(b.array(idx))
81     curo match {
82       term: Terminal =>
83         if !tryWriteSeal(term, b, idx, size)
84           seal(size)
85       elem: Elem =>
86         WRITE(b.index, idx + 1)
87         seal(size)
88       null =>
89         error("unreachable")
90     }
91   } else {
92     expand(b, b.array(idx))
93     seal(size)
94   }
95
96 def tryWriteSeal(term: Terminal, b: Block,
97   idx: Int, size: Int)
98   val total = totalElems(b, idx)
99   if total > size error("too many elements")
100  if term.sealed = NOSEAL {
101    nterm = new Terminal {
102      sealed = size
103      callbacks = term.callbacks
104    }
105    return CAS(b.array(idx), term, nterm)
106  } else if term.sealed ≠ size {
107    error("already sealed with different size")
108  } else return true
109
110 def foreach(f: Elem => Unit)
111   future {
112     asyncFor(f, start, 0)
113   }
114
115 def asyncFor(f: Elem => Unit, b: Block, idx: Int)
116   if idx <= LASTELEMPPOS {
117     obj = READ(b.array(idx))
118     obj match {
119       term: Terminal =>
120         nterm = new Terminal {
121           sealed = term.sealed
122           callbacks = f ∪ term.callbacks
123         }
124         if !CAS(b.array(idx), term, nterm)
125           asyncFor(f, b, idx)
126     elem: Elem =>
127       f(elem)
128     asyncFor(f, b, idx + 1)
129     null =>
130       error("unreachable")
131   }
132   } else {
133     expand(b, b.array(idx))
134     asyncFor(f, b.next, 0)
135   }

```

Fig. 3. FlowPool operations pseudocode

occasionally. Each block contains a hint `index` to the first free entry in the array, i.e. one that does not contain an element. An `index` is a hint, since it may actually reference an earlier index. The FlowPool maintains a reference to the first block called `start`. It also maintains a hint to the last block in the chain of blocks, called `current`. This reference may not always be up to date, but it always points to some block in the chain.

Each FlowPool is associated with a list of callbacks which have to be called in the future as new elements are added. Each FlowPool can be in a sealed state, meaning there is a bound on the number of elements it stores. This information is stored as a `Terminal` value in the first free entry of the array. At all times we maintain the invariant that the array in each block starts with a sequence of elements, followed by a `Terminal` delimiter. From a higher-level perspective, appending an element starts by copying the `Terminal` value to the next entry and then overwriting the current entry with the element being appended.

The `append` operation starts by reading the `current` block and the `index` of the free position. It then reads the `nexto` after the first free entry, followed by a read of the `curo` at the free entry. The `check` procedure checks the bounds conditions, whether the FlowPool was already sealed or if the current array entry contains an element. In either of these events, the `current` and `index` values need to be set – this is done in the `advance` procedure. We call this the **slow path** of the `append` method. Notice that there are several causes that trigger the slow path. If some other thread completes the `append` method but is preempted before updating the value of the hint `index`, then the `curo` will have the type `Elem`. The same happens if a preempted thread updates the value of the hint `index` after additional elements have been added, via unconditional write in line 23. Finally, reaching an end of block triggers the slow path.

Otherwise, the operation executes the **fast path** and appends an element. It first copies the `Terminal` value to the next entry with a CAS instruction in line 21, with `nexto` being the expected value. If it fails (e.g. due to a concurrent CAS), the append operation is restarted. Otherwise, it proceeds by writing the element to the current entry with a CAS in line 22, the expected value being `curo`. On success it updates the `b.index` value and invokes all the callbacks (present when the element was added) with the `future` construct. In the implementation we do not schedule an asynchronous computation for each element. Instead, the callback invocations are batched to avoid the scheduling overhead – the array is scanned for new elements until there are no more left.

Interestingly, inverting the order of the reads in lines 18 and 19 would cause a race in which a thread could overwrite a `Terminal` value with some older `Terminal` value if some other thread appended an element in between.

The `seal` operation continuously increases the `index` in the block until it finds the first free entry. It then tries to replace the `Terminal` value there with a new `Terminal` value which has the seal size set. An error occurs if a different seal size is set already. The `foreach` operation works in a similar way, but is executed asynchronously. Unlike `seal`, it starts from the first element in the pool and calls the callback for each element until it finds the first free entry. It then replaces the

`Terminal` value with a new `Terminal` value with the additional callback. From that point on the `append` method is responsible for scheduling that callback for subsequently added elements. Note that all three operations call `expand` to add an additional block once the current block is empty, to ensure lock-freedom.

Multi-Lane FlowPools. Using a single block sequence (i.e. lane) to implement a FlowPool doesn't fully take advantage of the lack of ordering guarantees and – as elements are inserted in sequence – may cause slowdowns due to collisions when multiple concurrent writers are present. Multi-Lane FlowPools solve this limitation by having a lane for each CPU, where each lane is following the same implementation as for a normal FlowPool.

This has several implications:

- CAS failures during insertion are to a high extent avoided.
- Less memory contention due to “private” blocks for each processor.
- `seal` needs to be globally synchronized in a non-blocking fashion.
- Callbacks for `aggregate` have to be added to each lane individually and aggregated once all of them have completed.

In practice, a hash on the ID of the inserting thread has to be used instead of the CPU index, as such information is currently not available through the Java framework. Also note that the remaining slots resulting from a `seal` have to be split up amongst the lanes and a writer might switch lanes if his lane is full (and hence CAS failures may happen).

You find an evaluation of the performance of Multi-Lane FlowPools in section 6.

5 Correctness

We give an outline of the correctness proof here. More formal definitions, a complete set of lemmas and proofs can be found in the appendix.

We define the notion of an abstract pool $\mathbb{A} = (elems, callbacks, seal)$ of elements in the pool, callbacks and the seal size. Given an abstract pool, abstract pool operations produce a new abstract pool. The key to showing correctness is to show that an abstract pool operation corresponds to a FlowPool operation – that is, it produces a new abstract pool corresponding to the state of the FlowPool after the FlowPool operation has been completed.

Lemma 5.1 Given a FlowPool consistent with some abstract pool, CAS instructions in lines 21, 63 and 66 do not change the corresponding abstract pool.

Lemma 5.2 Given a FlowPool consistent with an abstract pool $(elems, cbs, seal)$, a successful CAS in line 22 changes it to the state consistent with an abstract pool $(\{elem\} \cup elems, cbs, seal)$. There exists a time $t_1 \geq t_0$ at which every callback $f \in cbs$ has been called on $elem$.

$t ::=$	terms	
	$p \ll v$	append
	create p	pool creation
	p foreach f	foreach
	p seal n	seal
	$t_1 ; t_2$	sequence
	$p \in \{(vs, \sigma, cbs) \mid vs \subseteq Elem, \sigma \in \{-1\} \cup \mathbb{N}, cbs \subset Elem \Rightarrow Unit\}$	
	$f \in Elem \Rightarrow Unit$	
		$v \in Elem$ $n \in \mathbb{N}$

Fig. 4. Syntax

Lemma 5.3 Given a FlowPool consistent with an abstract pool $(elems, cbs, seal)$, a successful CAS in line 125 changes it to the state consistent with an abstract pool $(elems, (f, \emptyset) \cup cbs, seal)$ There exists a time $t_1 \geq t_0$ at which f has been called for every element in $elems$.

Lemma 5.4 Given a FlowPool consistent with an abstract pool $(elems, cbs, seal)$, a successful CAS in line 125 changes it to the state consistent with an abstract pool $(elems, cbs, s)$, where either $seal = -1 \wedge s \in \mathbb{N}_0$ or $seal \in \mathbb{N}_0 \wedge s = seal$.

Theorem 5.5 [Safety] FlowPool operations **append**, **foreach** and **seal** are consistent with the abstract pool semantics.

Theorem 5.6 [Linearizable operations] FlowPool operations **append** and **seal** are linearizable.

Lemma 5.7 [Non-consistency changing instructions] After invoking a FlowPool operation **append**, **seal** or **foreach**, if a non-consistency changing CAS instruction in lines 21, 63, or 66 fails, they must have already been completed by another thread since the FlowPool operation began.

Lemma 5.8 [Consistency changing instructions] After invoking a FlowPool operation **append**, **seal** or **foreach**, if a consistency-changing CAS instruction in lines 22, 105, or 124 fails, then some thread has successfully completed a consistency changing CAS after some finite number of steps.

Lemma 5.9 [Consistency changing operations] After invoking a FlowPool operation **append**, **seal** or **foreach**, a consistency changing instruction will be completed after a finite number of steps.

Lemma 5.10 Assuming some concurrent FlowPool operation is started. If some thread completes a consistency changing CAS instruction, then some concurrent operation is guaranteed to be completed.

Theorem 5.11 [Lock-freedom] FlowPool operations **append**, **foreach** and **seal** are lock-free.

Determinism. We claim that the FlowPool abstraction is *deterministic* in the sense that a program computes the same result (which can also be an error) regardless of the interleaving of concurrent operations. We give an outline of the determinism proof – a complete formal proof can be found in the appendix.

The following definitions and the determinism theorem are based on the language shown in Figure 4. The semantics of our core language is defined using reduction rules which define transitions between *execution states*. An execution state is a pair $T \mid P$ where T is a set of concurrent threads and P is a set of FlowPools. Each thread executes a *term* of the core language (typically a sequence of terms). State of a thread is represented as the (rest of) the term that it still has to execute; this means there is a one-to-one mapping between threads and terms. For example, the semantics of `append` is defined by the following reduction rule (a complete summary of all the rules can be found in the appendix):

$$\frac{t = p \ll v ; t' \quad p = (vs, -1, cbs) \quad p' = (\{v\} \cup vs, -1, cbs)}{t, T \mid p, P \longrightarrow t', T \mid p', P} \text{ (APPEND1)}$$

Append simply adds the value v to the pool p , yielding a modified pool p' . Note that this rule can only be applied if the pool p is not sealed (the seal size is -1). The rule for `foreach` modifies the set of callback functions in the pool:

$$\frac{t = p \text{ foreach } f ; t' \quad p = (vs, n, cbs) \quad T' = \{g(v) \mid g \in \{f\} \cup cbs, v \in vs\} \quad p' = (vs, n, \{f\} \cup cbs)}{t, T \mid p, P \longrightarrow t', T, T' \mid p', P} \text{ (FOREACH2)}$$

This rule only applies if p is sealed at size n , meaning that no more elements will be appended later. Therefore, an invocation of the new callback f is scheduled for each element v in the pool. Each invocation creates a new thread in T' .

Programs are built by first creating one or more FlowPools using `create`. Concurrent threads can then be started by (a) appending an element to a FlowPool, (b) sealing the FlowPool and (c) registering callback functions (`foreach`).

Definition 5.12 [Termination] A term t terminates with result P if its reduction ends in execution state $\{t : t = \{\epsilon\}\} \mid P$.

Definition 5.13 [Interleaving] Consider the reduction of a term $t: T_1 \mid P_1 \longrightarrow T_2 \mid P_2 \longrightarrow \dots \longrightarrow \{t : t = \{\epsilon\}\} \mid P_n$. An *interleaving* is a reduction of t starting in $T_1 \mid P_1$ in which reduction rules are applied in a different order.

Definition 5.14 [Determinism] The reduction of a term t is *deterministic* iff either (a) t does not terminate for any interleaving, or (b) t always terminates with the same result for all interleavings.

Theorem 5.15 [FlowPool Determinism] Reduction of terms t is deterministic.

6 Evaluation

We evaluate our implementation (single-lane and multi-lane FlowPools) against the `LinkedTransferQueue` [11] for all benchmarks and the `ConcurrentLinkedQueue` [14] for the `insert` benchmark, both found in JDK 1.7, on three different architectures – a quad-core 3.4 GHz i7-2600, 4x octa-core 2.27 GHz Intel Xeon x7560 (both with hyperthreading) and an octa-core 1.2GHz UltraSPARC T2 with 64 hardware threads. The level of parallelism P we vary as follows: On the first architecture we have $P \in [1, 8]$, on the latter two $P \in [1, 64]$. Here, we discuss the scaling of the said data structures (see fig 5). For additional evaluations, please refer to section D in the appendix.

In the *Insert* benchmark, we evaluate the ability to write concurrently – often the most difficult task – by distributing the work of inserting N elements into the data structure concurrently across P threads. We can see that both single-lane FlowPools and concurrent queues do not scale well with the number of concurrent threads, particularly on the i7 architecture. They seem to slow down rapidly due to cache line collisions and CAS failures. On the other hand, multi-lane FlowPools scale well, as threads write to different lanes and hence different cache lines in most occasions, while avoiding CAS failures. This may reduce execution time for insertions up to 54% on 4-core i7, 63% on 32-core Xeon and 92% on UltraSPARC T2.

Usage of the inserted data is evaluated in the *Reduce*, *Map* and *Histogram* benchmark, where the latter – serving as a “real life” example here – uses the former two primitives, whose evaluation is discussed in the appendix.

In the *Histogram* benchmark, P threads produce a total of N elements and add it to the FlowPool. The `aggregate` operation is then used to produce 10 different histograms concurrently with a different number of bins. Each separate histogram is constructed by its own thread (or up to P , for multi-lane FlowPools). A crucial difference between queues and FlowPools in the latter benchmark is that with FlowPools, multiple histograms are produced by invoking several `aggregate` operations, while queues require writing each element to several queues – one for each histogram. Without additional synchronization, reading a single queue is not an option since elements have to be removed from the queue eventually and it is not clear to each reader when an element is no longer needed. With FlowPools, elements are implicitly garbage collected when no longer needed.

Finally, to validate the last claim of garbage being implicitly collected, in the *Comm* benchmark we create a pool in which a large number of elements N is added concurrently by P threads. Each element is then processed by one of P threads through the use of the `aggregate` operation. With linked transfer queues, P threads concurrently remove elements from the queue and process it.

7 Related Work

An introduction to linearizability, lock-freedom and basic concurrent programming techniques is given by Herlihy and Shavit [10]. A detailed overview of

Operations on FlowPools Across Architectures

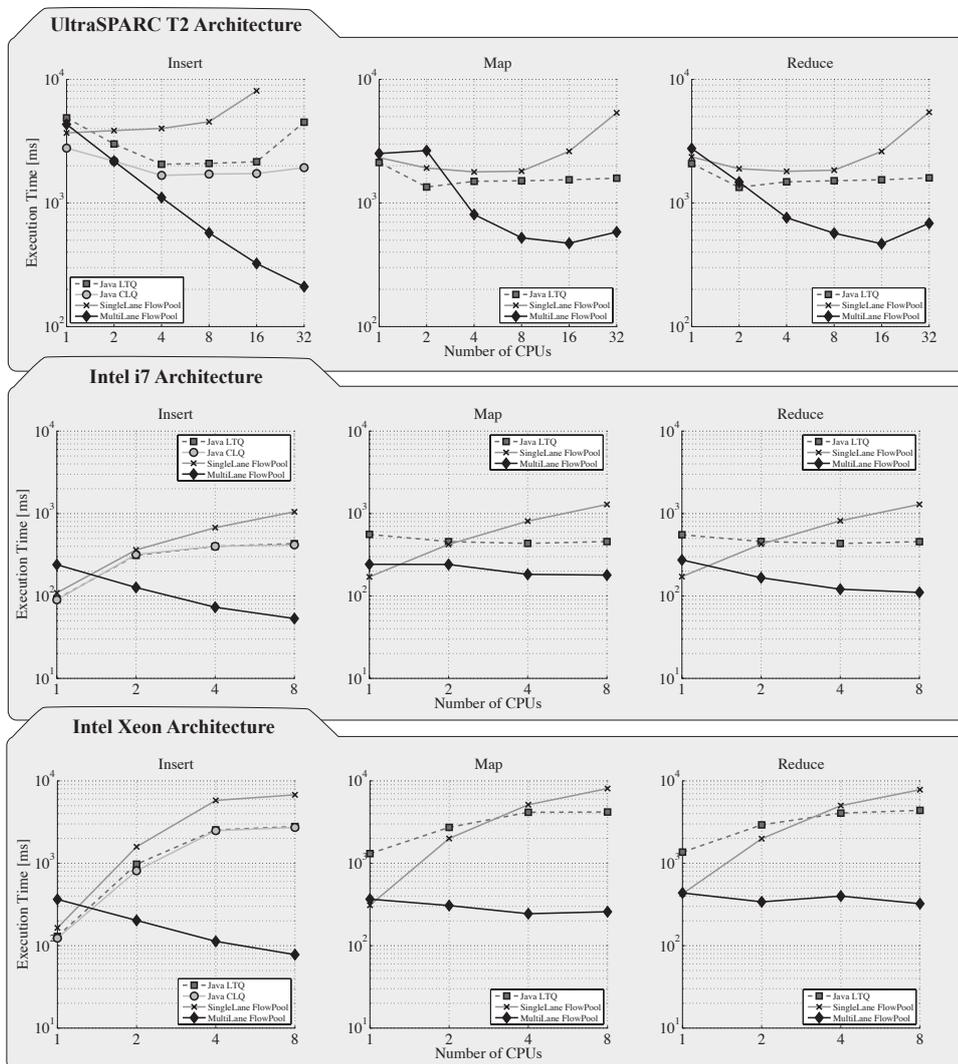


Fig. 5. Execution time vs parallelization across three different architectures on three important FlowPool operations; insert, map, reduce.

concurrent data structures is given by Moir and Shavit [15]. To date, concurrent data structures remain an active area of research – it is out of scope to list them all, so we restrict ourselves to those relevant to this work.

Concurrently accessible queues have been present for a while – an implementation is described by [13]. Non-blocking concurrent linked queues are described by Michael and Scott [14]. This CAS-based linked-list queue implementation is cited and used widely today, a variant of which is present in the Java standard library. More recently, Scherer, Lea and Scott [11] describe a synchronous

queues which internally hold both data and requests. Both approaches above entail blocking (or spinning) at least on the consumer's part when the queue is empty.

While the abstractions above fit well in the concurrent imperative model, they have the disadvantage that the programs written using them are inherently nondeterministic. Book by Roy and Haridi [17] describes the Oz programming language a subset of which yields programs deterministic by construction. The Oz dataflow streams are built on top of single-assignment variables and are deterministically ordered. They allow multiple consumers, but only one producer at a time. Oz has its own runtime which implements blocking using continuations. Blocking is less efficient on the JVM, CLR and many native environments.

The concept of single-assignment variables is also embodied in futures proposed by Baker and Hewitt [8], and promises first mentioned by Friedman and Wise [6]. Futures were first implemented in MultiLISP [7], and have been employed in many languages and frameworks since. Scala futures [?] and Finagle futures [?] are of interest, because they define monadic operators and a number of high-level combinators which produce new futures. This provides an API which avoids blocking. Futures have been generalized to data-driven futures, which can provide additional information to the scheduler [19]. Many frameworks have constructs which start an asynchronous computation and yield a future holding its result. For example, Habanero Java [3] offers an `async` construct and Scala futures offer a `future` construct.

A number of other models and frameworks recognized the need to embed the concept of futures into other data-structures. Single-assignment variables have been generalized to I-Structures proposed by Arvind [1] which are essentially single-assignment arrays. The CnC model [4] [2] is a parallel programming model influenced by dynamic dataflow, stream-processing and tuple spaces. In CnC the user provides high-level operations along with the ordering constraints that form a computation dependency graph. FlumeJava [5] is a distributed programming model which relies heavily on the concept of collections containing a set of futures. An issue that often arises with dataflow programming models are unbalanced loads. This is often solved through the use of bounded buffers which prevent the producer from overflowing the consumer. Analytical approaches to modeling pipelined applications have also been addressed [16].

Opposed to the correct-by-construction determinism described thus far, a type-systematic approach can also ensure that concurrent executions have deterministic results. Recently, work on Deterministic Parallel Java showed that a region-based type system can ensure determinism [12]. X10's constrained-based dependent types can similarly ensure determinism and deadlock-freedom [18].

8 Conclusion

We have shown that deterministic higher-level operations can be build on top of `<<` and `foreach`, which corresponds to sequential collections. The prerequisite is that `<<` can be invoked concurrently and that `foreach` executes asynchronously.

We emphasize a property of FlowPools which goes beyond a dataflow model such as Oz (or some other future/promise based model) in which more complex data structures such as streams are built on top of single-assignment variables. With single-assignment pools multiple threads can add elements without agreeing on where the added element should structurally be, as is the case with dataflow streams. A consequence of this is a more flexible programming model.

In conclusion, we postulate the existence of a range of other concurrent collection types with deterministic semantics fitting the presented correct-by-construction single-assignment model, such as bounded buffers, dataflow streams and maps, all of which have higher-level operations expressed in terms of the same basic primitives. On the implementation level, we anticipate the need of embedding the callbacks within the data-structure itself, as is the case with callback-based futures and FlowPools – this offers a particular benefit on platforms which do not by default support efficient continuations.

References

1. Arvind, R. S. Nikhil, and K. K. Pingali. I-structures: Data structures for parallel computing. *ACM Transactions on Programming Languages and Systems*, 11(4):598–632, Oct. 1989.
2. Z. Budimlic, M. G. Burke, V. Cavé, K. Knobe, G. Lowney, R. Newton, J. Palsberg, D. M. Peixotto, V. Sarkar, F. Schlimbach, and S. Tasirlar. Concurrent collections. *Scientific Programming*, 18(3-4):203–217, 2010.
3. Z. Budimlic, V. Cavé, R. Raman, J. Shirako, S. Tasirlar, J. Zhao, and V. Sarkar. The design and implementation of the habanero-java parallel programming language. In *OOPSLA Companion*, pages 185–186, 2011.
4. M. G. Burke, K. Knobe, R. Newton, and V. Sarkar. Concurrent collections programming model. In *Encyclopedia of Parallel Computing*, pages 364–371. 2011.
5. C. Chambers, A. Raniwala, F. Perry, S. Adams, R. R. Henry, R. Bradshaw, and N. Weizenbaum. FlumeJava: easy, efficient data-parallel pipelines. *ACM SIGPLAN Notices*, 45(6):363–375, June 2010.
6. D. Friedman and D. Wise. The impact of applicative programming on multiprocessing. In *International Conference on Parallel Processing*, 1976.
7. J. R. H. Halstead. MultiLISP: A language for concurrent symbolic computation. *ACM Transactions on Programming Languages and Systems*, 7(4):501–538, Oct. 1985.
8. J. Henry C. Baker and C. Hewitt. The incremental garbage collection of processes. In *Proceedings of the 1977 symposium on Artificial intelligence and programming languages*, 1977.
9. M. Herlihy. A methodology for implementing highly concurrent data structures. In *PPoPP*, pages 197–206, 1990.
10. M. Herlihy and N. Shavit. *The Art of Multiprocessor Programming*. Morgan Kaufman, Apr. 2008.
11. W. N. S. III, D. Lea, and M. L. Scott. Scalable synchronous queues. *Commun. ACM*, 52(5):100–111, 2009.
12. R. L. B. Jr., V. S. Adve, D. Dig, S. V. Adve, S. Heumann, R. Komuravelli, J. Overbey, P. Simmons, H. Sung, and M. Vakilian. A type and effect system for deterministic parallel java. In *OOPSLA*, pages 97–116, 2009.

13. J. M. Mellor-Crummey. Concurrent queues: Practical fetch-and- algorithms. 1987.
14. M. M. Michael and M. L. Scott. Simple, fast, and practical non-blocking and blocking concurrent queue algorithms. In *PODC*, pages 267–275, 1996.
15. Moir and Shavit. Concurrent data structures. In Mehta and Sahni, editors, *Handbook of Data Structures and Applications*, Chapman & Hall/CRC, 2005. 2005.
16. A. G. Navarro, R. Asenjo, S. Tabik, and C. Cascaval. Analytical modeling of pipeline parallelism. In *PACT*, pages 281–290, 2009.
17. P. V. Roy and S. Haridi. *Concepts, Techniques, and Models of Computer Programming*. MIT Press, 2004.
18. V. A. Saraswat, V. Sarkar, and C. von Praun. X10: concurrent programming for modern architectures. In *PPOPP*, page 271, 2007.
19. S. Tasirlar and V. Sarkar. Data-driven tasks and their implementation. In *ICPP*, pages 652–661, 2011.

A Proof of Correctness

Definition A.1 [Data types] A **Block** b is an object which contains an array $b.array$, which itself can contain elements, $e \in Elem$, where **Elem** represents the type of e and can be any countable set. A given block b additionally contains an index $b.index$ which represents an index location in $b.array$, a unique index identifying the array $b.blockIndex$, and $b.next$, a reference to a successor block c where $c.blockIndex = b.blockIndex + 1$. A **Terminal term** is a sentinel object, which contains an integer $term.sealed \in \{-1\} \cup \mathbb{N}_0$, and $term.callbacks$, a set of functions $f \in Elem \Rightarrow Unit$.

We define the following functions:

$$following(b : Block) = \begin{cases} \emptyset & \text{if } b.next = \text{null}, \\ b.next \cup following(b.next) & \text{otherwise} \end{cases}$$

$$reachable(b : Block) = \{b\} \cup following(b)$$

$$last(b : Block) = b' : b' \in reachable(b) \wedge b'.next = \text{null}$$

$$size(b : Block) = |\{x : x \in b.array \wedge x \in Elem\}|$$

Based on them we define the following relation:

$$reachable(b, c) \Leftrightarrow c \in reachable(b)$$

Definition A.2 [FlowPool] A **FlowPool** $pool$ is an object that has a reference $pool.start$, to the first block b_0 (with $b_0.blockIndex = 0$), as well as a reference $pool.current$. We sometimes refer to these just as $start$ and $current$, respectively.

A **scheduled callback invocation** is a pair (f, e) of a function $f \in Elem \Rightarrow Unit$ and an element $e \in Elem$. The programming construct that adds such a pair to the set of *futures* is `future { f(e) }`.

The **FlowPool state** is defined as a pair of the directed graph of objects transitively reachable from the reference $start$ and the set of scheduled callback invocations called *futures*.

A **state changing** or **destructive** instruction is any atomic write or CAS instruction that changes the FlowPool state.

We say that the FlowPool **has an element** e at some time t_0 if and only if the relation $hasElem(start, e)$ holds.

$$hasElem(start, e) \Leftrightarrow \exists b \in reachable(start), e \in b.array$$

We say that the FlowPool **has a callback** f at some time t_0 if and only if the relation $hasCallback(start, f)$ holds.

$$\begin{aligned} hasCallback(start, f) \Leftrightarrow \forall b = last(start), b.array = x^P \cdot t \cdot y^N, x \in Elem, \\ t = Terminal(seal, callbacks), f \in callbacks \end{aligned}$$

We say that a callback f in a FlowPool **will be called** for the element e at some time t_0 if and only if the relation $willBeCalled(start, e, f)$ holds.

$$willBeCalled(start, e, f) \Leftrightarrow \exists t_1, \forall t > t_1, (f, e) \in futures$$

We say that the FlowPool is **sealed** at the size s at some t_0 if and only if the relation $sealedAt(start, s)$ holds.

$$\begin{aligned} sealedAt(start, s) \Leftrightarrow s \neq -1 \wedge \forall b = last(start), b.array = x^P \cdot t \cdot y^N, \\ x \in Elem, t = Terminal(s, callbacks) \end{aligned}$$

FlowPool operations are **append**, **foreach** and **seal**, and are defined by pseudocodes in figures ...

Definition A.3 [Invariants] We define the following invariants for the **FlowPool**:

INV1 $start = b : Block, b \neq null, current \in reachable(start)$

INV2 $\forall b \in reachable(start), b \notin following(b)$

INV3 $\forall b \in reachable(start), b \neq last(start) \Rightarrow size(b) = LASTELEMPOS \wedge \\ b.array(BLOCKSIZE - 1) \in Terminal$

INV4 $\forall b = last(start), b.array = p \cdot c \cdot n$, where:

$$p = X^P, c = c_1 \cdot c_2, n = null^N$$

$$x \in Elem, c_1 \in Terminal, c_2 \in \{null\} \cup Terminal$$

$$P + N + 2 = BLOCKSIZE$$

INV5 $\forall b \in reachable(start), b.index > 0 \Rightarrow b.array(b.index - 1) \in Elem$

Definition A.4 [Validity] A FlowPool state \mathbb{S} is **valid** if and only if the invariants [INV1-5] hold for that state.

Definition A.5 [Abstract pool] An **abstract pool** \mathbb{P} is a function from time t to a tuple $(elems, callbacks, seal)$ such that:

$$seal \in \{-1\} \cup \mathbb{N}_0$$

$$callbacks \subset \{(f : Elem \Rightarrow Unit, called)\}$$

$$called \subseteq elems \subseteq Elem$$

We say that an abstract pool \mathbb{P} is **in state** $\mathbb{A} = (elems, callbacks, seal)$ at time t if and only if $\mathbb{P}(t) = (elems, callbacks, seal)$.

Definition A.6 [Abstract pool operations] We say that an **abstract pool operation** op that is applied to some abstract pool \mathbb{P} in abstract state $\mathbb{A}_0 = (elems_0, callbacks_0, seal_0)$ at some time t **changes** the abstract state of the abstract pool to $\mathbb{A} = (elems, callbacks, seal)$ if $\exists t_0, \forall \tau, t_0 < \tau < t, \mathbb{P}(\tau) = \mathbb{A}_0$ and $\mathbb{P}(t) = \mathbb{A}$. We denote this as $\mathbb{A} = op(\mathbb{A}_0)$.

Abstract pool operation $foreach(f)$ changes the abstract state at t_0 from $(elems, callbacks, seal)$ to $(elems, (f, \emptyset) \cup callbacks, seal)$. Furthermore:

$$\begin{aligned} \exists t_1 \geq t_0, \forall t_2 > t_1, \mathbb{P}(t_2) &= (elems_2, callbacks_2, seal_2) \\ \wedge \forall (f, called_2) \in callbacks_2, elems &\subseteq called_2 \subseteq elems_2 \end{aligned}$$

Abstract pool operation $append(e)$ changes the abstract state at t_0 from $(elems, callbacks, seal)$ to $(\{e\} \cup elems, callbacks, seal)$. Furthermore:

$$\begin{aligned} \exists t_1 \geq t_0, \forall t_2 > t_1, \mathbb{P}(t_2) &= (elems_2, callbacks_2, seal_2) \\ \wedge \forall (f, called_2) \in callbacks_2, (f, called) &\in callbacks \Rightarrow e \in called_2 \end{aligned}$$

Abstract pool operation $seal(s)$ changes the abstract state of the FlowPool at t_0 from $(elems, callbacks, seal)$ to $(elems, callbacks, s)$, assuming that $seal \in \{-1\} \cup \{s\}$ and $s \in \mathbb{N}_0$, and $|elems| \leq s$.

Definition A.7 [Consistency] A FlowPool state \mathbb{S} is **consistent** with an abstract pool $\mathbb{P} = (elems, callbacks, seal)$ at t_0 if and only if \mathbb{S} is a valid state and:

$$\begin{aligned} \forall e \in Elem, hasElem(start, e) &\Leftrightarrow e \in elems \\ \forall f \in Elem \Rightarrow Unit, hasCallback(start, f) &\Leftrightarrow f \in callbacks \\ \forall f \in Elem \Rightarrow Unit, \forall e \in Elem, willBeCalled(start, e, f) &\Leftrightarrow \exists t_1 \geq t_0, \mathbb{P}(t_1) = \\ &(elems_{s_1}, (f, called_1) \cup callbacks_{s_1}, seal_1), elems \subseteq called_1 \\ \forall s \in \mathbb{N}_0, sealedAt(start, s) &\Leftrightarrow s = seal \end{aligned}$$

A FlowPool operation op is **consistent** with the corresponding abstract state operation op' if and only if $\mathbb{S}' = op(\mathbb{S})$ is consistent with an abstract state $\mathbb{A}' = op'(\mathbb{A})$.

A **consistency change** is a change from state \mathbb{S} to state \mathbb{S}' such that \mathbb{S} is consistent with an abstract state \mathbb{A} and \mathbb{S}' is consistent with an abstract set \mathbb{A}' , where $\mathbb{A} \neq \mathbb{A}'$.

Proposition 1. *Every valid state is consistent with some abstract pool.*

Theorem A.8 [Safety] FlowPool operation **create** creates a new FlowPool consistent with the abstract pool $\mathbb{P} = (\emptyset, \emptyset, -1)$. FlowPool operations **foreach**, **append** and **seal** are consistent with the abstract pool semantics.

Lemma A.9 [End of life] For all blocks $b \in reachable(start)$, if value $v \in Elem$ is written to $b.array$ at some position idx at some time t_0 , then $\forall t > t_0, b.array(idx) = v$.

Proof. The CAS in line 22 is the only CAS which writes an element. No other CAS has a value of type *Elem* as the expected value. This means that once the CAS in line 22 writes a value of type *Elem*, no other write can change it.

Corollary 1. *The end of life lemma implies that if all the values in $b.array$ are of type *Elem* at t_0 , then $\forall t > t_0$ there is no write to $b.array$.*

Lemma A.10 [Valid hint] For all blocks $b \in reachable(start)$, if $b.index > 0$ at some time t_0 , then $b.array(b.index - 1) \in Elem$ at time t_0 .

Proof. Observe every write to $b.index$ – they are all unconditional. However, at every such write occurring at some time t_1 that writes some value idx we know that some previous value at $b.array$ entry $idx-1$ at some time $t_0 < t_1$ was of type *Elem*. Hence, from Lemma A.9 it follows that $\forall t \geq t_1, b.array(idx - 1) \in Elem$.

Corollary 2 (Compactness). *For all blocks $b \in reachable(start)$, if for some idx $b.array(idx) \in Elem$ at time t_0 then $b.array(idx - 1) \in Elem$ at time t_0 . This follows directly from the Lemma A.9 and Lemma A.10, and the fact that the CAS in line 22 only writes to array entries idx for which it previously read the value from $b.index$.*

Definition A.11 [Transition] If for a function $f(t)$ there exist times t_0 and t_1 such that $\forall t, t_0 < t < t_1, f(t) = v_0$ and $f(t_1) = v_1$, then we say that the function f goes through a **transition** at t_1 . We denote this as:

$$f : v_0 \xrightarrow{t_1} v_1$$

Or, if we don't care about the exact time t_1 , simply as:

$$f : v_0 \rightarrow v_1$$

Definition A.12 [Monotonicity] A function of time $f(t)$ is said to be **monotonic**, if every value in its string of transitions occurs only once.

Lemma A.13 [Freshness] For all blocks $b \in reachable(start)$, and for all $x \in b.array$, function x is monotonic.

Proof. CAS instruction in line 22 writes a value of type *Elem*. No CAS instruction has a value of type *Elem* as the expected value.

Trivial analysis of CAS instructions in lines 105 and 125, shows that their expected values are of type *Terminal*. Their new values are always freshly allocated.

The more difficult part is to show that CAS instruction in line 21 respects the statement of the lemma.

Since the CAS instructions in lines 105 and 125 are preceded by a read of $idx = b.index$, from Lemma A.10 it follows that $b.array(idx - 1)$ contains a value of type *Elem*. These are also the only CAS instructions which replace a *Terminal* value with another *Terminal* value. The new value is always unique, as shown above.

So the only potential CAS to write a non-fresh value to $idx + 1$ is the CAS in line 21.

A successful CAS in line 21 overwrites a value cb_0 at $idx + 1$ read in line 18 at t_0 with a new value cb_2 at time t_2 . Value cb_2 was read in line 19 at t_1 from the entry idx . The string of transitions of values at idx is composed of unique values at least since t_1 (by Lemma A.9), since there is a value of type *Elem* at the index $idx - 1$.

The conclusion above ensures that the values read in line 19 to be subsequently used as new values for the CAS in line 21 form a monotonic function $f(t) = b.array(idx)$ at t .

Now assume that a thread T1 successfully overwrites cb_0 via CAS in line 21 at $idx + 1$ at time t_2 to a value cb_2 read from idx at t_1 , and that another thread T2 is the **first** thread (since the FlowPool was created) to subsequently successfully complete the CAS in line 21 at $idx + 1$ at time $t_{prev2} > t_2$ with some value cb_{prev2} which was at $idx + 1$ at some time $t < t_0$.

That would mean that $b.array(idx + 1)$ does not change during $\langle t_0, t_2 \rangle$, since T2 was the first thread the write a non-fresh value to $idx + 1$, and any other write would cause the CAS in line 21 by T1 to fail.

Also, that would mean that the thread T2 read the value cb_{prev2} in line 19 at some time $t_{prev1} < t_1$ and successfully completed the CAS at time $t_{prev2} > t_2$. If the CAS was successful, then the read in line 18 by T2 occurred at $t_{prev0} < t_{prev1} < t_1$. Since we assumed that T2 is the first thread to write a value cb_{prev2} to $idx + 1$ at time t_{prev2} which was previously in $idx + 1$ at some time $t < t_0$, then the CAS in line 21 at time t_{prev2} could not have succeeded, since its expected value is cb_{prev0} read at some time t_{prev0} , and we know that the value at $idx + 1$ was changed at least once in $\langle t_{prev0}, t_{prev2} \rangle$ because of the write of a fresh value by thread T1 at $t_2 \in \langle t_{prev0}, t_{prev2} \rangle$. This value is known to be fresh because $b.array(idx)$ is a monotonic function at least since t_{prev1} , and the read of the new value written by T1 occurred at $t_1 > t_{prev1}$. We also know that there is no other thread T3 to write the value cb_{prev0} during $\langle t_{prev0}, t_{prev2} \rangle$ back to $idx + 1$, since we assumed that T2 is the first to write a non-fresh value at that position.

Hence, a contradiction shows that there is no thread T2 which is the **first** to write a non-fresh value via CAS in line 21 at $idx + 1$ for any idx , so there is no thread that writes a non-fresh value at all.

Lemma A.14 [Lifecycle] For all blocks $b \in reachable(start)$, and for all $x \in b.array$, function x goes through and only through the prefix of the following transitions:

$$\begin{aligned} & null \rightarrow cb_1 \rightarrow \dots \rightarrow cb_n \rightarrow elem, \text{ where:} \\ & cb_i \in Terminal, i \neq j \Rightarrow cb_i \neq cb_j, elem \in Elem \end{aligned}$$

Proof. First of all, it is obvious from the code that each block that becomes an element of $reachable(start)$ at some time t_0 has the value of all $x \in b.array$ set to *null*.

Next, we inspect all the CAS instructions that operate on entries of $b.array$.

The CAS in line 22 has a value $curo \in Terminal$ as an expected value and writes an $elem \in Elem$. This means that the only transition that this CAS can cause is of type $cb_i \in Terminal \rightarrow elem \in Elem$.

We will now prove that the CAS in line 21 at time t_2 is successful if and only if the entry at $idx + 1$ is *null* or *nexto* \in *Terminal*. We know that the entry at $idx + 1$ does not change $\forall t, t_0 < t < t_2$, where t_0 is the read in line 18, because of Lemma A.13 and the fact that CAS in line 21 is assumed to be successful. We know that during the read in line 19 at time t_1 , such that $t_0 < t_1 < t_2$, the entry at idx was *curo* \in *Terminal*, by trivial analysis of the **check** procedure. It follows from corollary 2 that the array entry $idx + 1$ is not of type *Elem* at time t_1 , otherwise array entry idx would have to be of type *Elem*. Finally, we know that the entry at $idx + 1$ has the same value during the interval $\langle t_1, t_2 \rangle$, so its value is not *Elem* at t_2 .

The above reasoning shows that the CAS in line 21 always overwrites a one value of type *Terminal* (or *null*) with another value of type *Terminal*. We have shown in Lemma A.13 that it never overwrites the value cb_0 with a value cb_2 that was at $b.array(idx)$ at an earlier time.

Finally, note that the statement for CAS instructions in lines 105 and 125 also follows directly from the proof for Lemma A.13.

Lemma A.15 [Subsequence] Assume that for some block $b \in reachable(start)$ the transitions of $b.array(idx)$ are:

$$b.array(idx) : null \rightarrow cb_1 \rightarrow \dots \rightarrow cb_n \xrightarrow{t_0} elem : Elem$$

Assume that the transitions of $b.array(idx + 1)$ up to time t_0 are:

$$b.array(idx + 1) : null \rightarrow cb'_1 \rightarrow \dots \rightarrow cb'_m$$

The string of transitions $null \rightarrow cb'_1 \rightarrow \dots \rightarrow cb'_m$ is a subsequence of $null \rightarrow cb_1 \rightarrow \dots \rightarrow cb_n$.

Proof. Note that all the values written to $idx + 1$ before t_0 by CAS in line 21 were previously read from idx in line 19. This means that the set of values occurring in $b.array(idx + 1)$ before t_0 is a subset of the set of values in $b.array(idx)$. We have to prove that it is actually a subsequence.

Assume that there exist two values cb_1 and cb_2 read by threads T1 and T2 in line 19 at times t_1 and $t_2 > t_1$, respectively. Assume that these values are written to $idx + 1$ by threads T1 and T2 in line 21 in the opposite order, that is at times t_{cas1} and $t_{cas2} < t_{cas1}$, respectively. That would mean that the CAS by thread T1 would have to fail, since its expected value cb_0 has changed between the time it was read in line 18 and the t_{cas1} at least once to a different value, and it could not have been changed back to cb_0 as we know from the Lemma A.13.

Notice that we have actually proved a stronger result above. We have also shown that the string of values written at $idx + 1$ by CAS in line 21 successfully is a subsequence of **all** the transitions of values at idx (not just until t_0).

Lemma A.16 [Valid writes] Given a FlowPool in a valid state, all writes in all operations produce a FlowPool in a valid state.

Proof. A new FlowPool is trivially in a valid state.

Otherwise, assume that the FlowPool is in a valid state \mathbb{S} . In the rest of the proof, whenever some invariant is trivially unaffected by a write, we omit mentioning it. We start by noting that we already proved the claim for atomic writes in lines 23, 55 and 86 (which only affect [INV5]) in Lemma A.10. We proceed by analyzing each atomic CAS instruction.

CAS in line 63 at time t_1 maintains the invariant [INV1]. This is because its expected value is always *null*, which ensures that the lifecycle of $b.next$ is $null \rightarrow b' : Block$, meaning that the function $reachable(start)$ returns a monotonically growing set. So if $current \in reachable(start)$ at t_0 , then this also holds at $t_1 > t_0$. It also maintains [INV2] because the new value nb is always fresh, so $\forall b, b \notin following(b)$. Finally, it maintains [INV3] because it is preceded with a bounds check and we know from corollary 2 and the Lemma A.9 that all the values in $b.array[idx], idx < LASTELEMPOS$ must be of type *Elem*.

CAS in line 66 at time t_1 maintains the invariant [INV1], since the new value for the $current \neq null$ was read from $b.next$ at $t_0 < t_1$ when the invariant was assumed to hold, and it is still there at t_1 , as shown before.

For CAS instructions in lines 22, 125 and 105 that write to index idx we know from Lemma A.10 that the value at $idx - 1$ is of type *Elem*. This immediately shows that CAS instructions in lines 125 and 105 maintain [INV3] and [INV4].

For CAS in line 22 we additionally know that it must have been preceded by a successful CAS in line 21 which previously wrote a *Terminal* value to $idx + 1$. From Lemma A.14 we know that $idx + 1$ is still *Terminal* when the CAS in line 22 occurs, hence [INV4] is kept.

Finally, CAS in line 21 succeeds only if the value at $idx + 1$ is of type *Terminal*, as shown before in Lemma A.14. By the same lemma, the value at idx is either *Terminal* or *Elem* at that point, since $idx - 1$ is known to be *Elem* by Lemma A.10. This means that [INV4] is kept.

Lemma A.17 [Housekeeping] Given a FlowPool in state \mathbb{S} consistent with some abstract pool state \mathbb{A} , CAS instructions in lines 21, 63 and 66 do not change the abstract pool state \mathbb{A} .

Proof. Since none of the relations $hasElem$, $hasCallback$, $willBeCalled$ and $sealedAt$ are defined by the value of $current$ CAS in line 66 does not change them, hence it does not change the abstract pool state.

No CAS changes the set of scheduled futures, nor is succeeded by a *future* construct so it does not affect the $willBeCalled$ relation.

It is easy to see that the CAS in line 63 does not remove any elements, nor make any additional elements reachable, since the new block nb which becomes reachable does not contain any elements at that time. Hence the $hasElem$ relation is not affected. It does change the value $last(start)$ to nb , but since $nb.array = t \cdot null^{BLOCKSIZE-1}$, where $t \in Terminal$ was previously the last non-null element in $b.array$, it does changes neither the $sealedAt$ nor the $hasCallback$ relation.

The CAS in line 21 does not make some new element reachable, hence the $hasElem$ relation is preserved.

Note now that this CAS does not change the relations *hasCallback* and *sealedAt* as long as there is a value of type *Terminal* at the preceding entry *idx*. We claim that if the CAS succeeds at t_2 , then either the value at *idx* is of type *Terminal* (trivially) or the CAS did not change the value at $idx + 1$. In other words, if the value at *idx* at time t_2 is of type *Elem*, then the write by CAS in line 21 does not change the value at $idx + 1$ at t_2 . This was, in fact, already shown in the proof of Lemma A.15.

The argument above proves directly that relations *hasCallback* and *sealedAt* are not changed by the CAS in line 21.

Lemma A.18 [Append correctness] Given a FlowPool in state \mathbb{S} consistent with some abstract pool state \mathbb{A} , a successful CAS in line 22 at some time t_0 changes the state of the FlowPool to \mathbb{S}_0 consistent with an abstract pool state \mathbb{A}_0 , such that:

$$\begin{aligned}\mathbb{A} &= (elems, callbacks, seal) \\ \mathbb{A}_0 &= (\{elem\} \cup elems, callbacks, seal)\end{aligned}$$

Furthermore, given a fair scheduler, there exists a time $t_1 > t_0$ at which the FlowPool is consistent with an abstract pool in state \mathbb{A}_1 , such that:

$$\begin{aligned}\mathbb{A}_1 &= (elems_1, callbacks_1, seal_1), \text{ where:} \\ \forall (f, called_1) \in callbacks_1, (f, called) \in callbacks &\Rightarrow elem \in called_1\end{aligned}$$

Proof. Assume that the CAS in line 22 succeeds at some time t_3 , the CAS in line 21 succeeds at some time $t_2 < t_3$, the read in line 19 occurs at some time $t_1 < t_2$ and the read in line 19 occurs at some time $t_0 < t_1$.

It is easy to see from the invariants, *check* procedure and the corollary 1 that the CAS in line 22 can only occur if $b = last(start)$.

We claim that for the block $b \in reachable(start)$ such that $b = last(b)$ the following holds at t_2 :

$$b.array = elem^N \cdot cb_1 \cdot cb_2 \cdot null^{BLOCKSIZE-N-2}$$

where $cb_1 = cb_2$, since there was no write to *idx* after cb_1 , otherwise the CAS in line 22 at t_3 would not have been successful (by lemma Lemma A.13).

Furthermore, $cb_1 = cb_2$ at t_3 , as shown in the Lemma A.15. Due to the same lemma, the entries of *b.array* stay the same until t_3 , otherwise the CAS in line 22 would not have been successful. After the successful CAS at t_3 , we have:

$$b.array = elem^N \cdot e \cdot cb_1 \cdot null^{BLOCKSIZE-N-2}$$

where $e : Elem$ is the newly appended element— at t_3 the relation *hasElem*(*start*, *e*) holds, and *sealedAt*(*start*, *s*) and *hasCallback*(*start*, *f*) did not change between t_2 and t_3 .

It remains to be shown that *willBeCalled*(*start*, *e*, *f*) holds at t_3 . Given a fair scheduler, within a finite number of steps the future store will contain a request for an asynchronous computation that invokes *f* on *e*. The fair scheduler ensures that the future is scheduled within a finite number of steps.

Lemma A.19 [Foreach correctness] Given a FlowPool in state \mathbb{S} consistent with some abstract pool state \mathbb{A} , a successful CAS in line 125 at some time t_0 changes the state of the FlowPool to \mathbb{S}_0 consistent with an abstract pool state \mathbb{A}_0 , such that:

$$\begin{aligned}\mathbb{A} &= (elems, callbacks, seal) \\ \mathbb{A}_0 &= (elems, (f, \emptyset) \cup callbacks, seal)\end{aligned}$$

Furthermore, given a fair scheduler, there exists a time $t_1 \geq t_0$ at which the FlowPool is consistent with an abstract pool in state \mathbb{A}_1 , such that:

$$\begin{aligned}\mathbb{A}_1 &= (elems_1, callbacks_1, seal_1), \text{ where:} \\ elems &\subseteq elems_1 \\ \forall(f, called_1) \in callbacks_1, elems &\subseteq called_1\end{aligned}$$

Proof. From Lemma A.13 and the assumption that the CAS is successful we know that the value at $b.array(idx)$ has not changed between the read in line 117 and the CAS in line 125. From Lemma A.10 we know that the value at $idx-1$ was of type *Elem* since $b.index$ was read. This means that neither $hasElem(start, e)$ nor $sealedAt$ have changed after the CAS. Since after the CAS there is a *Terminal* with an additional function f at idx , the $hasCallback(start, f)$ holds after the CAS. Finally, the $willBeCalled(start, e, f)$ holds for all elements e for which the $hasElem(e)$ holds, since the CAS has been preceded by a call $f(e)$ in line 127 for each element. The Lemma A.9 ensures that for each element f was called for stays in the pool indefinitely (i.e. is not removed).

Trivially, the time t_1 from the statement of the lemma is such that $t_1 = t_0$.

Lemma A.20 [Seal correctness] Given a FlowPool in state \mathbb{S} consistent with some abstract pool state \mathbb{A} , a successful CAS in line 105 at some time t_0 changes the state of the FlowPool to \mathbb{S}_0 consistent with an abstract pool state \mathbb{A}_0 , such that:

$$\begin{aligned}\mathbb{A} &= (elems, callbacks, seal), \text{ where } seal \in \{-1\} \cup \{s\} \\ \mathbb{A}_0 &= (elems, callbacks, \{s\})\end{aligned}$$

Proof. Similar to the proof of Lemma A.19.

Definition A.21 [Obstruction-freedom] Given a FlowPool in a valid state, an operation op is **obstruction-free** if and only if a thread T executing the operation op completes within a finite number of steps given that no other thread was executing the operation op since T started executing it.

We say that thread T executes the operation op **in isolation**.

Lemma A.22 [Obstruction-free operations] All operations on FlowPools are obstruction-free.

Proof. By trivial sequential code analysis supported by the fact that the invariants (especially [INV2]) hold in a valid state.

Proof (Safety). From Lemma A.17, Lemma A.31, Lemma A.19 and Lemma A.20 directly, along with the fact that all operations executing in isolation complete after a finite number of steps by Lemma A.22.

Definition A.23 [Linearizability] We say that an operation op is linearizable if every thread observes that it completes at some time t_0 after it was invoked and before it finished executing.

Theorem A.24 [Linearizable operations] FlowPool operations `append` and `seal` are linearizable.

Proof (Linearizable operations). This follows directly from statements about CAS instructions in Lemma A.17, Lemma A.31 and Lemma A.20, along with the fact that a CAS instruction itself is linearizable.

Note that `foreach` starts by executing an asynchronous computation and then returns the control to the caller. This means that the linearization point may happen outside the execution interval of that procedure – so, `foreach` is not linearizable.

Definition A.25 [Lock-freedom] In a scenario where some finite number of threads are executing a concurrent operation, that concurrent operation is *lock-free* if and only if that concurrent operation is completed after a finite number of steps by some thread.

Theorem A.26 [Lock-freedom] FlowPool operations `append`, `seal`, and `foreach` are lock-free.

We begin by first proving that there are a finite number of execution steps before a consistency change occurs.

By Lemma A.31, after invoking `append`, a consistency change occurs after a finite number of steps. Likewise, by Lemma A.34, after invoking `seal`, a consistency change occurs after a finite number of steps. And finally, by Lemma A.35, after invoking `foreach`, a consistency change likewise occurs after a finite number of steps.

By Lemma A.36, this means a concurrent operation `append`, `seal`, or `foreach` will successfully complete. Therefore, by Definition A.25, these operations are lock-free.

Note. For the sake of clarity in this section of the correctness proof, we assign the following aliases to the following CAS and WRITE instructions:

- $CAS_{append-out}$ corresponds to the outer CAS in `append`, on line 21.
- $CAS_{append-inn}$ corresponds to the inner CAS in `append`, on line 22.
- $CAS_{expand-nxt}$ corresponds to the CAS on `next` in `expand`, line 63.
- $CAS_{expand-curr}$ corresponds to the CAS on `current` in `expand`, line 66.
- CAS_{seal} corresponds to the CAS on the `Terminal` in `tryWriteSeal`, line 105.
- $CAS_{foreach}$ corresponds to the CAS on the `Terminal` in `asyncFor`, line 124.
- $WRITE_{app}$ corresponds to the WRITE on the new `index` in `append`, line 23.
- $WRITE_{adv}$ corresponds to the WRITE on the new `index` in `advance`, line 55.
- $WRITE_{seal}$ corresponds to the WRITE on the new `index` in `seal`, line 86.

Lemma A.27 After invoking an operation `op`, if non-consistency changing CAS operations $CAS_{append-out}$, $CAS_{expand-nxt}$, or $CAS_{expand-curr}$, in the pseudocode fail, they must have already been successfully completed by another thread since `op` began.

Proof. Trivial inspection of the pseudocode reveals that since $CAS_{append-out}$ makes up a check that precedes $CAS_{append-inn}$, and since $CAS_{append-inn}$ is the only operation besides $CAS_{append-out}$ which can change the expected value of $CAS_{append-out}$, in the case of a failure of $CAS_{append-out}$, $CAS_{append-inn}$ (and thus $CAS_{append-out}$) must have already successfully completed or $CAS_{append-out}$ must have already successfully completed by a different thread since op began executing.

Likewise, by trivial inspection $CAS_{expand-nxt}$ is the only CAS which can update the $b.next$ reference, therefore in the case of a failure, some other thread must have already successfully completed $CAS_{expand-nxt}$ since the beginning of op .

Like above, $CAS_{expand-curr}$ is the only CAS which can change the *current* reference, therefore in the case of a failure, some other thread must have already successfully completed $CAS_{expand-curr}$ since op began. \square

Lemma A.28 [Expand] Invoking the *expand* operation will execute a non-consistency changing instruction after a finite number of steps. Moreover, it is guaranteed that the *current* reference is updated to point to a subsequent block after a finite number of steps. Finally, *expand* will return after a finite number of steps

Proof. From inspection of the pseudocode, it is clear that the only point at which $expand(b)$ can be invoked is under the condition that for some block b , $b.index > LASTELEMPOS$, where $LASTELEMPOS$ is the maximum size set aside for elements of type *Elem* in any block. Given this, we will proceed by showing that a new block will be created with all related references $b.next$ and *current* correctly set.

There are two conditions under which a non-consistency changing CAS instruction will be carried out.

- **Case 1:** if $b.next = null$, a new block nb will be created and $CAS_{expand-nxt}$ will be executed. From Lemma A.27, we know that $CAS_{expand-nxt}$ must complete successfully on some thread. Afterwards recursively calling *expand* on the original block b .
- **Case 2:** if $b.next \neq null$, $CAS_{expand-curr}$ will be executed. Lemma A.27 guarantees that $CAS_{expand-curr}$ will update *current* to refer to $b.next$, which we will show can only be a new block. Likewise, Lemma A.27 has shown that $CAS_{expand-nxt}$ is the only state changing instruction that can initiate a state change at location $b.next$, therefore, since $CAS_{expand-nxt}$ takes place within Case 1, Case 2 can only be reachable after Case 1 has been executed successfully. Given that Case 1 always creates a new block, therefore, $b.next$ in this case, must always refer to a new block.

Therefore, since from Lemma A.27 we know that both $CAS_{expand-nxt}$ and $CAS_{expand-curr}$ can only fail if already completed guaranteeing their finite completion, and since $CAS_{expand-nxt}$ and $CAS_{expand-curr}$ are the only state changing operations invoked through *expand*, the *expand* operation must complete in a finite number of steps.

Finally, since we saw in Case 2 that a new block is always created and related references are always correctly set, that is both $b.next$ and $current$ are correctly updated to refer to the new block, it follows that $numBlocks$ strictly increases after some finite number of steps. \square

Lemma A.29 [$CAS_{append-inn}$] After invoking $append(elem)$, if $CAS_{append-inn}$ fails, then some thread has successfully completed $CAS_{append-inn}$ or CAS_{seal} (or likewise, $CAS_{foreach}$) after some finite number of steps.

Proof. First, we show that a thread attempting to complete $CAS_{append-inn}$ can't fail due to a different thread completing $CAS_{append-out}$ so long as $seal$ has not been invoked after completing the read of $currobj$. We address this exception later on.

Since after $check$, the only condition under which $CAS_{append-out}$, and by extension, $CAS_{append-inn}$ can be executed is the situation where the current object $currobj$ with index location idx is the *Terminal* object, it follows that $CAS_{append-out}$ can only ever serve to duplicate this *Terminal* object at location $idx+1$, leaving at most two *Terminals* in block referred to by $current$ momentarily until $CAS_{append-inn}$ can be executed. By Lemma A.27, since $CAS_{append-out}$ is a non-consistency changing instruction, it follows that any thread holding any element $elem'$ can execute this instruction without changing the expected value of $currobj$ in $CAS_{append-inn}$, as no new object is ever created and placed in location idx . Therefore, $CAS_{append-inn}$ cannot fail due to $CAS_{append-out}$, so long as $seal$ has not been invoked by some other thread after the read of $currobj$.

This leaves only two scenarios in which consistency changing $CAS_{append-inn}$ can fail:

- **Case 1:** Another thread has already completed $CAS_{append-inn}$ with a different element $elem'$.
- **Case 2:** Another thread completes an invocation to the $seal$ operation after the current thread completes the read of $currobj$. In this case, $CAS_{append-inn}$ can fail because CAS_{seal} (or, likewise $CAS_{foreach}$) might have completed before, in which case, it inserts a new *Terminal* object $term$ into location idx (in the case of a $seal$ invocation, $term.sealed \in \mathbb{N}_0$, or in the case of a $foreach$ invocation, $term.callbacks \in \{Elem \Rightarrow Unit\}$).

We omit the proof and detailed discussion of $CAS_{foreach}$ because it can be proven using the same steps as were taken for CAS_{seal} . \square

Lemma A.30 [Finite Steps Before State Change] All operations with the exception of $append$, $seal$, and $foreach$ execute only a finite number of steps between each state changing instruction.

Proof. The $advance$, $check$, $totalElems$, $invokeCallbacks$, and $tryWriteSeal$ operations have a finite number of execution steps, as they contain no recursive calls, loops, or other possibility to restart.

While the `expand` operation contains a recursive call following a CAS instruction, it was shown in Lemma A.28 that an invocation of `expand` is guaranteed to execute a state changing instruction after a finite number of steps.

□

Lemma A.31 [Append] After invoking `append(elem)`, a consistency changing instruction will be completed after a finite number of steps.

Proof. The `append` operation can be restarted in three cases. We show that in each case, it's guaranteed to either complete in a finite number of steps, or leads to a state changing instruction:

- **Case 1:** The call to `check`, a finite operation by Lemma A.30, returns *false*, causing a call to `advance`, also a finite operation by Lemma A.30, followed by a recursive call to `append` with the same element *elem* which in turn once again calls `check`.

We show that after a finite number of steps, the `check` will evaluate to *true*, or some other thread will have completed a consistency changing operation since the initial invocation of `append`. In the case where `check` evaluates to *true*, Lemma A.29 applies, as it guarantees that a consistency changing CAS is completed after a finite number of steps.

When the call to the finite operation `check` returns *false*, if the subsequent `advance` finds that a *Terminal* object is at the current block index *idx*, then the next invocation of `append` will evaluate `check` to *true*. Otherwise, it must be the case that another thread has moved the *Terminal* to a subsequent index since the initial invocation of `append`, which is only possible using a consistency changing instruction.

Finally, if `advance` finds that the element at *idx* is an *Elem*, *b.index* will be incremented after a finite number of steps. By *INV1*, this can only happen a finite number of times until a *Terminal* is found. In the case that `expand` is meanwhile invoked through `advance`, by Lemma A.28 it's guaranteed to complete state changing instructions $CAS_{expand-nxt}$ or $CAS_{expand-curr}$ in a finite number of steps. Otherwise, some other thread has moved the *Terminal* to a subsequent index. However, this latter case is only possible by successfully completing $CAS_{append-inn}$, a consistency changing instruction, after the initial invocation of `append`.

- **Case 2:** $CAS_{append-out}$ fails, which we know from Lemma A.27 means that it must've already been completed by another thread, guaranteeing that $CAS_{append-inn}$ will be attempted. If $CAS_{append-inn}$ fails, after a finite number of steps, a consistency changing instruction will be completed. If $CAS_{append-inn}$ succeeds, as a consistency changing instruction, consistency will have clearly been changed.
- **Case 3:** $CAS_{append-inn}$ fails, which, by Lemma A.29, indicates that either some other thread has already completed $CAS_{append-inn}$ with another element, or another consistency changing instruction, CAS_{seal} or $CAS_{foreach}$ has successfully completed.

Therefore, `append` itself as well as all other operations reachable via an invocation of `append` are guaranteed to have a finite number of steps between *consistency* changing instructions. \square

Lemma A.32 [CAS_{seal}] After invoking `seal(size)`, if CAS_{seal} fails, then some thread has successfully completed CAS_{seal} or $CAS_{append-inn}$ after some finite number of steps.

Proof. Since by Lemma A.29, we know that $CAS_{append-out}$ only duplicates an existing *Terminal*, it can not be the cause for a failing CAS_{seal} . This leaves only two cases in which CAS_{seal} can fail:

- **Case 1:** Another thread has already completed CAS_{seal} .
- **Case 2:** Another thread completes an invocation to the `append(elem)` operation after the current thread completes the read of *currobj*. In this case, CAS_{seal} can fail because $CAS_{append-inn}$ might have completed before, in which case, it inserts a new *Elem* object *elem* into location *idx*. \square

Lemma A.33 [$WRITE_{adv}$ and $WRITE_{seal}$] After updating *b.index* using $WRITE_{adv}$ or $WRITE_{seal}$, *b.index* is guaranteed to be incremented after a finite number of steps.

Proof. For some index, *idx*, both calls to $WRITE_{adv}$ and $WRITE_{seal}$ attempt to write *idx* + 1 to *b.index*. In both cases, it's possible that another thread could complete either $WRITE_{adv}$ or $WRITE_{seal}$, once again writing *idx* to *b.index* after the current thread has completed, in effect overwriting the current thread's write with *idx* + 1. By inspection of the pseudocode, both $WRITE_{adv}$ and $WRITE_{seal}$ will be repeated if *b.index* has not been incremented. However, since the number of threads operating on the FlowPool is finite, *p*, we are guaranteed that in the worst case, this scenario can repeat at most *p* times, before a write correctly updates *b.index* with *idx* + 1. \square

Lemma A.34 [Finite Steps Before Consistency Change] After invoking `seal(size)`, a consistency changing instruction will be completed after a finite number of steps, or the initial invocation of `seal(size)` completes.

Proof. The `seal` operation can be restarted in two scenarios.

- **Case 1:** The check $idx \leq LASTELEMPOS$ succeeds, indicating that we are at a valid location in the current block *b*, but the object at the current index location *idx* is of type *Elem*, not *Terminal*, causing a recursive call to `seal` with the same size *size*.
In this case, we begin by showing that the atomic write of *idx* + 1 to *b.index*, required to iterate through the block *b* for the recursive call to `seal`, will be correctly incremented after a finite number of steps.
Therefore, by both the guarantee that, in a finite number of steps, *b.index* will eventually be correctly incremented as we saw in Lemma A.33, as

well as by *INV1* we know that the original invocation of `seal` will correctly iterate through b until a *Terminal* is found. Thus, we know that the call to `tryWriteSeal` will be invoked, and by both Lemma A.30 and Lemma A.31, we know that either `tryWriteSeal`, will successfully complete in a finite number of steps, in turn successfully completing `seal(size)`, or $CAS_{append-inn}$, another consistency changing operation will successfully complete.

- **Case 2:** The check $idx \leq LASTELEMPOS$ fails, indicating that we must move on to the next block, causing first a call to `expand` followed by a recursive call to `seal` with the same size $size$.

We proceed by showing that after a finite number of steps, we must end up in Case 1, which we have just showed itself completes in a finite number of steps, or that a consistency change must've already occurred.

By Lemma A.28, we know that an invocation of `expand` returns after a finite number of steps, and $pool.current$ is updated to point to a subsequent block. If we are in the recursive call to `seal`, and the $idx \leq LASTELEMPOS$ condition is *false*, trivially, a consistency changing operation must have occurred, as, the only way for the condition to evaluate to *true* is through a consistency changing operation, in the case that a block has been created during an invocation to `append`, for example.

Otherwise, if we are in the recursive call to `seal`, and the $idx \leq LASTELEMPOS$ condition evaluates to *true*, we enter Case 1, which we just showed will successfully complete in a finite number of steps.

□

Lemma A.35 [Foreach] After invoking `foreach(fun)`, a consistency changing instruction will be completed after a finite number of steps.

We omit the proof for `foreach` since it proceeds in the exactly the same way as does the proof for `seal` in Lemma A.34.

Lemma A.36 Assume some concurrent operation is started. If some thread completes a consistency changing CAS instruction, then some concurrent operation is guaranteed to be completed.

Proof. By trival inspection of the pseudocode, if $CAS_{append-inn}$ successfully completes on some thread, then that thread is guaranteed to complete the corresponding invocation of `append` in a finite number of steps.

Likewise by trivial inspection, if CAS_{seal} successfully completes on some thread, then by Lemma A.30, `tryWriteSeal` is guaranteed to complete in a finite number of steps, and therefore, that thread is guaranteed to complete the corresponding invocation of `seal` in a finite number of steps.

The case for $CAS_{foreach}$ is omitted since it follows the same steps as for the case of CAS_{seal}

B Proof of Determinism

B.1 Reduction rules

$$\frac{t = \mathbf{create} \ p ; \ t' \ p = (\emptyset, -1, \emptyset)}{t, T \mid P \longrightarrow t', T \mid p, P} \quad (\text{CREATE})$$

$$\frac{t = p \ll v ; \ t' \ p = (vs, -1, cbs) \quad p' = (\{v\} \cup vs, -1, cbs)}{t, T \mid p, P \longrightarrow t', T \mid p', P} \quad (\text{APPEND1})$$

$$\frac{t = p \ll v ; \ t' \ p = (vs, n, cbs) \quad |\{v\} \cup vs| \leq n \\ p' = (\{v\} \cup vs, n, cbs) \quad T' = \{f(v) \mid f \in cbs\}}{t, T \mid p, P \longrightarrow t', T, T' \mid p', P} \quad (\text{APPEND2})$$

$$\frac{t = p \ \mathbf{foreach} \ f ; \ t' \ p = (vs, -1, cbs) \quad p' = (vs, -1, \{f\} \cup cbs)}{t, T \mid p, P \longrightarrow t', T \mid p', P} \quad (\text{FOREACH1})$$

$$\frac{t = p \ \mathbf{foreach} \ f ; \ t' \ p = (vs, n, cbs) \\ T' = \{g(v) \mid g \in \{f\} \cup cbs, v \in vs\} \quad p' = (vs, n, \{f\} \cup cbs)}{t, T \mid p, P \longrightarrow t', T, T' \mid p', P} \quad (\text{FOREACH2})$$

$$\frac{t = p \ \mathbf{seal} \ n ; \ t' \ p = (vs, -1, cbs) \\ T' = \{g(v) \mid g \in cbs, v \in vs\} \quad p' = (vs, n, cbs)}{t, T \mid p, P \longrightarrow t', T, T' \mid p', P} \quad (\text{SEAL1})$$

$$\frac{t = p \ \mathbf{seal} \ n ; \ t' \ p = (vs, n, cbs)}{t, T \mid p, P \longrightarrow t', T \mid p, P} \quad (\text{SEAL2})$$

B.2 Definitions

Definition B.1 [Termination] A term t terminates with result P if its reduction ends in execution state $\{t : t = \{\epsilon\}\} \mid P$.

Definition B.2 [Interleaving] Consider the reduction of a term $t: T_1 \mid P_1 \longrightarrow T_2 \mid P_2 \longrightarrow \dots \longrightarrow \{t : t = \{\epsilon\}\} \mid P_n$. An *interleaving* is a reduction of t starting in execution state $T_1 \mid P_1$ in which reduction rules are applied in a different order.

Definition B.3 [Valid Interleaving] An interleaving S', c, S'' of a reduction sequence S is *valid iff* for any reduction step c using rule (CREATE), if c creates pool p , p is not used in any reduction step in S' .

Definition B.4 [Determinism] The reduction of a term t is *deterministic iff* either (a) t does not terminate for any valid interleaving, or (b) t always terminates with the same result for all valid interleavings.

B.3 Theorems and Lemmas

Theorem B.5 [Determinism] The reduction of terms t is deterministic.

To prove this, we show that the reduction steps of every valid interleaving can be reordered in a way that does not change the termination, nor the result. The proof is based on the commutativity of transitions in the reduction.

Lemma B.6 [Commutativity] Consider the reduction R of a term t and the two subsequent reduction steps based on reduction rules R_1 and R_2 :

$$\dots \longrightarrow T_n \mid P_n \xrightarrow{R_1} T_{n+1} \mid P_{n+1} \xrightarrow{R_2} T_{n+2} \mid P_{n+2} \longrightarrow \dots$$

We list the reduction rules for which it is true that switching these two reduction steps yields the following reduction R' :

$$\dots \longrightarrow T_n \mid P_n \xrightarrow{R'_2} T'_{n+1} \mid P'_{n+1} \xrightarrow{R'_1} T_{n+2} \mid P_{n+2} \longrightarrow \dots$$

where the prefix and suffix of the reduction R' are the same as in R .

Any pair of rules R_1 and R_2 applies above if the reduction steps operate on different pools. If they operate on the same pool, then for any two (not necessarily different) threads executing these reduction steps, the following table of rules applies.

R_1	R_2	R'_1	R'_2
APPEND1	APPEND1	APPEND1	APPEND1
APPEND1	FOREACH1	FOREACH1	APPEND1
APPEND1	SEAL1	SEAL1	APPEND2
APPEND2	FOREACH2	FOREACH2	APPEND2
APPEND2	SEAL2	SEAL2	APPEND2
FOREACH1	APPEND1	APPEND1	FOREACH1
FOREACH1	FOREACH1	FOREACH1	FOREACH1
FOREACH1	SEAL1	SEAL1	FOREACH2
FOREACH2	APPEND2	APPEND2	FOREACH2
FOREACH2	SEAL2	SEAL2	FOREACH2
SEAL1	APPEND2	APPEND1	SEAL1
SEAL1	FOREACH2	FOREACH1	SEAL1
SEAL1	SEAL2	SEAL1	SEAL2
SEAL2	APPEND2	APPEND2	SEAL2
SEAL2	FOREACH2	FOREACH2	SEAL2

Proof. By straightforward inspection of each pair of rules.

As an example, we pick the third row from the table. Consider the following execution schedule fragment:

$$\begin{aligned}
& p \ll v ; t', p \text{ seal } n ; t'', T_k \mid p = (vs, -1, cbs), P_k \\
& \xrightarrow{\text{APPEND}^1} t', p \text{ seal } n ; t'', T_k \mid p = (\{v\} \cup vs, -1, cbs), P_k \\
& \xrightarrow{\text{SEAL}^1} \{cbs(w) : w \in \{v\} \cup vs\}, t', t'', T_k \mid p = (\{v\} \cup vs, n, cbs), P_k
\end{aligned}$$

If we reorder these two reduction steps, we get:

$$\begin{aligned}
& p \ll v ; t', p \text{ seal } n ; t'', T_k \mid p = (vs, -1, cbs), P_k \\
& \xrightarrow{\text{SEAL}^1} \{cbs(w) : w \in vs\}, p \ll v ; t', t'', T_k \mid p = (vs, n, cbs), P_k \\
& \xrightarrow{\text{APPEND}^2} cbs(v), \{cbs(w) : w \in vs\}, t', t'', T_k, \mid p = (\{v\} \cup vs, n, cbs), P_k
\end{aligned}$$

The final execution states in these two execution fragments are the same, hence the rest of the interleaving stays the same.

Definition B.7 [Canonical interleaving] Assume, without the loss of generality, that each CREATE rule assigns a unique identifier to each created pool. Assume any ordering on this set of identifiers.

We define the following ordering **comes-before** between threads as follows. There exists one **main thread** defined by the term t being reduced. The main thread comes before any other thread. For all other threads T_p and T_q let p and q be the terms with which the threads first appear in some execution state. Thread T_p comes before T_q iff p comes before q in the lexicographic ordering.

Note that in any interleaving we can identify the first reduction step in which an arbitrary thread T_p appears. We can recursively find the rest of its reduction steps. This means that we can assign the index p to every thread in the set of threads in every execution state later³, and thus uniquely identify every thread in any execution state.

We define the **canonical order of the reduction steps of a thread** as follows:

1. First apply all the CREATE steps in the order of the pool identifiers.
2. Then apply all the APPEND steps in the order of the pool identifiers and some ordering of the elements being appended⁴.
3. Then apply all the FOREACH steps in the order of the pool identifiers and the lexicographic ordering on the callback term.
4. Finally, apply all the SEAL steps in the order of the pool identifiers and the sizes of the seal.

³ We claim that in our programming model this choice is unique, but even if this was not the case, we could pick either thread in the case of an ambiguity and do so without the loss of generality.

⁴ Any countable set has some total ordering.

Or, more formally:

$$\begin{aligned}
 & T_0 \mid P_0 \xrightarrow{\text{CREATE}} T_0 \mid p_{i_1}, P_0 \xrightarrow{\text{CREATE}} \dots \xrightarrow{\text{CREATE}} T_0 \mid p_{i_n}, p_{i_{n-1}}, \dots, p_{i_1}, P_0 \\
 & \xrightarrow{\text{APPEND}} T_1 \mid (\{v_{j_1}\} \cup vs_1, n_1, cbs_1)_{j_1}, P_{j_1} \xrightarrow{\text{APPEND}} \dots \\
 & \xrightarrow{\text{APPEND}} T_m \mid (\{v_{j_m}\} \cup vs_m, n_m, cbs_m)_{j_m}, P_{j_m} \\
 & \xrightarrow{\text{FOREACH}} T_{m+1} \mid (vs_{m+1}, n_{m+1}, \{t_{k_1}\} \cup cbs_{m+1})_{k_1}, P_{k_1} \xrightarrow{\text{FOREACH}} \dots \\
 & \xrightarrow{\text{FOREACH}} T_{m+f} \mid (vs_{m+f}, n_{m+f}, \{t_{k_f}\} \cup cbs_{m+f})_{k_f}, P_{k_f} \\
 & \xrightarrow{\text{SEAL}} T_{m+f+1} \mid (vs_{m+f+1}, n_{m+f+1}, cbs_{m+f+1})_{l_1}, P_{l_1} \xrightarrow{\text{SEAL}} \dots \\
 & \xrightarrow{\text{SEAL}} T_{m+f+s} \mid (vs_{m+f+s}, n_{m+f+s}, cbs_{m+f+s})_{l_s}, P_{l_s}
 \end{aligned}$$

where (with $\overset{a}{<}$ being the lexicographic ordering and $\dot{<}$ being some ordering on the set of elements):

$$\begin{aligned}
 & i_1 < i_2 < \dots < i_n \\
 & x < y \Rightarrow v_{j_x} \dot{<} v_{j_y} \vee (v_{j_x} = v_{j_y} \wedge j_x < j_y) \\
 & x < y \Rightarrow t_{k_x} \overset{a}{<} t_{k_y} \vee (t_{k_x} = t_{k_y} \wedge k_x < k_y) \\
 & x < y \Rightarrow n_x < n_y \vee (n_x = n_y \wedge l_x < l_y)
 \end{aligned}$$

We define the **canonical order of the reduction steps** as follows:

1. First apply all the reduction steps of the main thread t in the canonical order.
2. Identify all the threads created in the previous step. Take each of these threads in the above defined comes-before ordering and apply all their reduction steps in the canonical order.
3. Repeat the last step until there are no more applicable reduction steps.

Or, more formally:

$$\begin{aligned}
 & \{t\} \mid \emptyset \xrightarrow{S_C(t)} \{\epsilon, T_1, T_2, \dots, T_n\} \mid P(t) \\
 & \xrightarrow{S_C(T_1)} \{\epsilon, \epsilon, T_2, \dots, T_n, T_{1,1}, \dots, T_{1,n_1}\} \mid P(T_1), P(t) \\
 & \xrightarrow{S_C(T_2)} \dots \xrightarrow{S_C(T_n)} \\
 & \xrightarrow{S_C(T_n)} \{\epsilon, \dots, \epsilon, T_{1,1}, \dots, T_{1,n_1}, T_{2,1}, \dots, T_{n-1,n_{n-1}}, T_{n,1}, \dots, T_{n,n_n}\} \\
 & \mid P(T_1), \dots, P(T_n), P(t) \\
 & \xrightarrow{S_C(T_{1,1})} \dots \xrightarrow{S_C(T_{1,n})} \dots \xrightarrow{S_C(T_{2,1})} \dots \xrightarrow{S_C(T_{n-1,n_{n-1}})} \dots \xrightarrow{S_C(T_{n,1})} \dots
 \end{aligned}$$

where $S_C(T)$ is the sequence of reduction steps of thread T in the canonical order and t is the term being reduced, i.e. the main thread. For each T_p and T_q if $p \stackrel{a}{<} q$, then $T_p \stackrel{cb}{<} T_q$, where $\stackrel{cb}{<}$ is the comes-before ordering.

A **canonical interleaving** S_c is an interleaving of the reduction of a term t , such that the reduction steps are applied in the canonical order.

Proposition 2. *Every canonical interleaving S_c is a valid interleaving.*

Lemma B.8 [Canonicity] Every valid interleaving of the reduction of a term t can have its execution steps reordered to the canonical interleaving S_c so that neither the termination nor the result is changed.

Proof. Consider an arbitrary valid interleaving S . We know that the starting execution state is $\{t\} \mid \emptyset$, where t is the term being reduced. Here, t is the main thread. As argued before, we can recursively identify all the reduction steps in the interleaving S in which the main thread executes.

We start by identifying all the CREATE steps of the main thread. Relying on the results from Lemma B.6 we create a new interleaving in which the CREATE steps appear at the beginning in the canonical order. We know that the state after the last CREATE step in the original interleaving does not change by reordering the steps. Therefore, neither the termination nor the result are changed in the new interleaving. The new interleaving is trivially valid.

We then identify all the APPEND steps of the main thread. We know from our programming model that each of these steps can only refer to FlowPools from the CREATE steps we have already moved to the beginning of the interleaving – we cannot append to a pool created by some other subsequently created thread, due to scoping rules. For this reason we can reorder the APPEND steps of the main thread so that they appear after the CREATE steps and maintain a valid interleaving. Again, from Lemma B.6 the new interleaving has the same termination and the result. We repeat the same for the FOREACH and SEAL steps.

At this point the new interleaving is the prefix of a canonical interleaving. From the commutativity rules we know that the main thread still creates the same set of threads as in the original interleaving. We order the newly created threads according to the comes-before ordering introduced earlier and identify each of these threads in the original interleaving. We recursively apply the same reordering of steps for the newly created threads, traversing them in the comes-before ordering.

We claim that every reduction in our programming model has a finite number of steps, so this procedure will be completed eventually, yielding a canonical interleaving.

More generally, if our programming model supported recursion, we could have infinite reductions. In that case, every time we produce a new interleaving in the process above, we extend the prefix of the interleaving which is the same as in the canonical interleaving. Furthermore, every time we produce a new interleaving in the process above, there is an execution state after the last

reordered reduction step which is the same in both the new and the old interleaving, meaning that all the subsequent execution states are the same in both interleavings (by Lemma B.6). This means that for any reduction of a term t we can produce a new interleaving with a prefix of an arbitrary length which corresponds to the prefix of the canonical interleaving, in the same time not changing the non-termination. Hence, every reduction of a term t does not terminate.

Proof (Determinism). Directly from the definition of determinism and Lemma B.8.

C Syntax and examples

This section the used syntax in more detail and presents a range of programming examples. The syntax is based on languages such as Groovy, Scala and Ruby. For conciseness and reasons of space, we omit the block braces and use indentation to denote block boundaries.

Methods. We use the `def` keyword to declare methods. After declaring the method name, the parameter list follows. Each parameter is given a name and a type behind a colon. Here is an example of declaring a `max` method which returns the greater of the two integers:

```
def max(a: Int, b: Int)
  if a > b a else b
```

Each method may either be standalone or defined within some object, in our case a `FlowPool`. If the method is defined within the object, it can refer to the object instance using the `this` keyword. It can additionally call the methods of the `this` object without prefixing their names with a `this` and a dot. Otherwise, invoking methods belonging to object instances must be prefixed with the object instance name and a dot, as in most object-oriented languages. Methods can be generic in their types and this is expressed with a list of type parameters in square brackets behind the method name. The following generic method just returns its parameter.

```
def id[T](x: T)
  x

id[Int](0)
id(0)
id(true)
id("Ok!")
```

Notice that above we did not have to put a type parameter value `T` when invoking the method – we assume that type parameters are inferred from the types of regular method parameters.

Methods can also nest, as in the following example of a method which computes the sum of first `n` numbers:

```
def sum(n: Int)
  def subsum(i: Int)
    if i == n n
    else i + sum(i + 1)
  subsum(0)
```

Finally, methods can have multiple parameter lists, which provides a nice syntax for methods such as `aggregate`.

Values. Values are declared using the `val` keyword. Once declared, the value does not change. The following method creates an empty `FlowPool`:

```
def empty[T]
  val fp = new FlowPool[T]
  fp.builder.seal(0)
  fp
```

First-class functions. First-class functions can be simulated with anonymous classes in most object-oriented languages, but we provide special syntax for reasons of conciseness. The type of the function which takes parameters of type T_1 to T_N and returns a value of type R is denoted as $(T_1, T_2, \dots, T_N) \Rightarrow R$. The function values are expressed by writing the name of the parameter, followed by a \Rightarrow keyword and the method body.

Here is an example of a generic method which takes a value of type T and applies a custom function on it twice:

```
def twice[T](x: T)(f: T => T)
  f(f(x))
```

Since `twice` has multiple parameter lists, we can invoke with either of the following two notations:

```
twice(0)(x => x + 1)
twice(0) {
  x => x + 1
}
```

Additionally, we can omit the $x \Rightarrow$ prefix when defining the function value to make notation even more concise:

```
twice(0) {
  _ + 1
}
```

We can do this only if the parameter occurs in the method body only once. Each subsequent occurrence of a `_` keyword denotes a different parameter.

For-comprehensions. We define syntactic sugar for element traversal, best described through a couple of examples. The following for-loop which is supposed to print numbers from 0 until 10:

```
for (i <- 0 until 10) {
  println(i)
}
```

is desugared into the following expression:

```
(0.until(10)).foreach {
  i => println(i)
}
```

Above, the expression following the `<-` keyword must be an object containing the `foreach` method. This method must take a single parameter function value – in this case the block that prints a given number. We assume that the object produced by the expression `0.until(10)` is predefined for integer values.

This mechanism allows us to use the same for-loop notation both for traversing numbers, installing callback handlers on future value or asynchronously traversing `FlowPool` elements.

In some cases, given a set of values being traversed, we want to produce a new set of values. The syntax for this involves the `yield` keyword, as in the following example:

```

val fp = new FlowPool[Int]()
fp.builder << 1
fp.builder << 2
fp.builder << 3

for (i <- fp) yield {
  i * i
}

```

The for-loop above is translated into a call to the `map` method on FlowPools:

```

fp.map {
  i => i * i
}

```

The `map` method must take a single argument function. The `map` method on FlowPools will return a new FlowPool with every element mapped.

A similar method called `flatMap` takes a single argument function which returns another traversable object. It can be used to compose traversals over several objects within one for-loop. The following for-loop which traverses two flow pools to produce a Cartesian product of their elements:

```

for (x <- fp1; y <- fp2) yield (x, y)

```

is translated into the following calls:

```

fp1.flatMap {
  x => fp2.map {
    y => (x, y)
  }
}

```

Futures. Futures are values which will become available at some point in the future. We distinguish between a value of type `Future[T]`, where `T` is the type of the value which becomes available (for example, an integer – `Int`) and an asynchronous computation which completes a future value. This asynchronous computation may be started using the `future` construct:

```

val p = new FlowPool[Int]
val f = future {
  p << 1
}
p.seal(1)

```

Above, the last line may be executed before or after appending the element 1 at runtime. Also, the `future` construct returns a future completed with the value its associated block computes.

Futures can be used in for-comprehensions, i.e. they define `foreach`, `map` and `flatMap` methods. A `foreach` loop on a future is an asynchronous computation which executes once the value of the future becomes available – it is a for-loop which traverses only a single value, and only once it becomes available. More complex patterns can also arise. Given two futures `f` and `g`, we can produce a third future which whose value becomes available only after values of `f` and `g` are available:

```

val f = future {
  sum(10)
}
val p = new FlowPool[Int]
val g = p.aggregate(0)(_ + _)(_ + _)
for (i <- 0 until 10) p.builder << i * i
p.seal(10)

val h: Future[Double] = for (x <- f; y <- g) yield {
  sqrt(y / 10 - x * x / 100)
}

```

Above, the real-valued future `h` contains the standard deviation of the first ten integers. The for-comprehension on futures is desugared into:

```

f.flatMap {
  x => g.map {
    y => sqrt(y / 10 - x * x / 100)
  }
}

```

Examples. We start off with a couple of methods that create new FlowPools. The `iterate` method iteratively applies a function to the starting value and puts it into a FlowPool. The `range` method creates a fixed size FlowPool from a range of numbers. The `fill` method creates a fixed size FlowPool filled with exactly the same element. Given the syntax primer above, it should be straightforward to grasp their contents.

```

def iterate[T]
  (s: T, f: T => T)
  val p = new FlowPool[T]
  val b = p.builder
  def recurse(x: T) {
    b << x
    recurse(f(x))
  }
  future { recurse(s) }
  p

def range
  (from: Int, end: Int)
  val p = new FlowPool[Int]
  val b = p.builder
  future {
    for (i <- start to end)
      b << i
  }
  b.seal(n)
  p

def fill[T]
  (n: Int, elem: =>T)
  val p = new FlowPool[T]
  val b = p.builder
  future {
    for (i <- 1 to n)
      b << elem
  }
  b.seal(n)
  p

```

We now show that `aggregate` can be used to implement a range of other methods based on reducing the values. Notice that the `fold` we define is less general than the `aggregate`, since it places a type bound on the type parameter, meaning it can only operate on supertypes of the values in the pool⁵.

```

def exists
  (pred: T => Boolean)
  aggregate(false)(_ ∨ _) {
    (acc, x) =>
    acc ∨ pred(x)
  }

def forall
  (pred: T => Boolean)
  aggregate(true)(_ ∧ _) {
    (acc, x) =>
    acc ∧ pred(x)
  }

def min()
  def min(a: Int, b: Int)
  if a < b a else b
  aggregate(0)(min) {
    min
  }

```

⁵ We have not precisely defined what type bounds and inheritance are, but it suffices to say that in a language with subtyping an unrestricted fold cannot be expressed in terms of aggregate we have defined.

```

def sum()
  aggregate(0)(_ + _) {
    - + -
  }

def product()
  aggregate(1)(_ * _) {
    - * -
  }

def count
  (pred: T => Boolean)
  aggregate(0)(_ + _) {
    (acc, x) =>
      if pred(x) 1 else 0
  }

def fold[U >: T]
  (zero: U, op: (U, U) => U)
  aggregate(zero)(op)(op)

def flatten[S]
  (q: FlowPool[FlowPool[S]])
  val p = new FlowPool[S]
  val b = p.builder
  aggregate(future(0))(add) {
    (af, r) =>
      val sf = for (x <- q) b << x
      add(af, sf)
  } map { sz => b.seal(sz) }
  p

```

The `flatten` method turns a `FlowPool` whose elements are `FlowPools` themselves into a `FlowPool` containing all the elements of the nested `FlowPools`. It could also have been expressed in terms of the `flatMap` method described earlier.

Sometimes there is a need to convert between futures and `FlowPools`. Lets assume we have a list of futures and we wish to have a `FlowPool` of their values instead. The method `toFlowPool` does this by calling a `foreach` on every future. We assume the existence of a traversable `List` datatype.

The converse is not so easy, since we do not know how many elements will there be available in the `FlowPool`. This means we cannot create a list of futures, because we do not know how many future there will be. We can instead create a future which holds the final list of values. This is what the method `toList` does. Notice that the `aggregate` invocation does not form a commutative monoid, so the results will be nondeterministic, which is reflected in the order of elements in the resulting list. We could still maintain determinism given that we use a `Set` datatype instead of a `List`.

```

def toFlowPool[T]
  (fs: List[Future[T]])
  val p = new FlowPool[T]
  val b = p.builder
  for (f <- fs; x <- f)
    b << x
  p

def toList[T]
  (p: FlowPool[T])
  p.aggregate(())(_ :: _) {
    (acc, x) => x :: acc
  }

def toList2[T]
  (p: FlowPool[T], n: Int)
  type LFT = List[Future[T]]
  val fs = List.fill(n)
    (new Future[T]())
  def complete(l: LFT, x: T)
    if !l.head.tryComplete(x)
      complete(l.tail, x)
    else l
  p.builder.seal(n)
  p.aggregate(fs) {
    (gs, hs) => gs
  } {
    (acc, x) =>
      complete(acc, x)
  }
  fs

def intersect
  (that: FlowPool[T])
  val p = new FlowPool[T]
  val b = p.builder
  for (x <- this) {
    that.aggregate(false)(_ ∨ _) {
      (acc, y) =>
        if (!acc) {
          b << x
          true
        } else false
    }
  }
  p

```

However, given that we assume the number of elements in the `FlowPool`, we can do better than that. The method `toList2` starts by sealing the pool to the expected number of elements. If successful, the predefined futures are completed as the elements arrive into the `FlowPool`. We have to use the `tryComplete` method on futures, which completes the future only given that it has not already been completed. This potentially yields nondeterministic computations, but we cannot avoid this without more expressive abstractions, such as pools for which we know that all of the elements they hold are the same and some sort of a

`forany` construct which invokes a callback on any element of the pool, roughly speaking.

The `intersect` method above produces a new `FlowPool` with elements that appear both in the current `FlowPool` `this` and another `FlowPool` `that`. It is not very efficient, however – a more efficient implementation would require a more expressive single-assignment abstraction such as a single-assignment map.

Architecture	Elements	FlowPool $t[ms]$	P	Queue $t[ms]$	P	decr.
4-core i7	2M	17	8	35	1	51%
4-core i7	5M	44	8	87	1	49%
4-core i7	15M	118	8	258	1	54%
UltraSPARC T2	1M	23	32	111	4	79%
UltraSPARC T2	2M	34	64	224	4	84%
UltraSPARC T2	5M	62	64	556	4	88%
UltraSPARC T2	15M	129	64	1661	4	92%
32-core Xeon	2M	30	8	50	1	40%
32-core Xeon	5M	46	64	120	1	61%
32-core Xeon	15M	126	64	347	1	63%

Table 1. Execution times for insert benchmark for multi-lane FlowPool and concurrent linked queues, including execution time decrease percentage.

D Additional Evaluation

In this section, some additional evaluation results are presented.

Map and Reduce The *Reduce* benchmark starts P threads which concurrently insert a total of N elements. The **aggregate** operation is used to reduce the set of values inserted into the pool. Note that in the FlowPool implementation there may be as many threads computing the aggregation as there are different lanes – elements from different lanes are batched together once the pool is sealed.

The *Map* benchmark is similar to the *Reduce* benchmark, but instead of reducing a value, each element is mapped into a new one and added to a second pool.

Scaling in Input Size In figure 6 we can see that the Input, Map, Reduce and Histogram benchmark all scale linearly in the input size with any parallelism level. The Comm benchmark has not been tested for different sizes.

Multi-Lane Scaling By default, the number of lanes is set to the parallelism level P , corresponding to the number of used CPUs. However, since the implementation has to use hashing on the thread IDs instead of the real CPU index, we tested whether varying the number of lanes to $1.5P$, $2P$, $3P$ and $4P$ results in performance gain due to fewer collisions. Benchmarks have shown (see fig 7) that this yields no observable gain – in fact, this sometimes even decreased performance slightly.

Performance Gain As stated in the abstract, FlowPools – or more precisely multi-lane FlowPools – may reduce execution time by 49 – 54% on 4-core i7. These figures have been obtained by comparing medians of execution times for insertions between multi-lane FlowPools and concurrent linked queues (which were always faster than linked transfer queues), where each structure was evaluated on its optimal parallelization level. The resulting data is shown in table 1.

Methodology All the presented configurations have been measured 20 times, where the 5 first values have been discarded to let the JIT stabilize. Aggregated values are always medians. The benchmarks have been written using

`scala.testing.Benchmark` and executed through SBT⁶ using the following flags for the JVM: `-Xmx2048m -Xms2048m -XX:+UseCondCardMark -verbose:gc -XX:+PrintGCDetails -server`.

⁶ Simple Build Tool

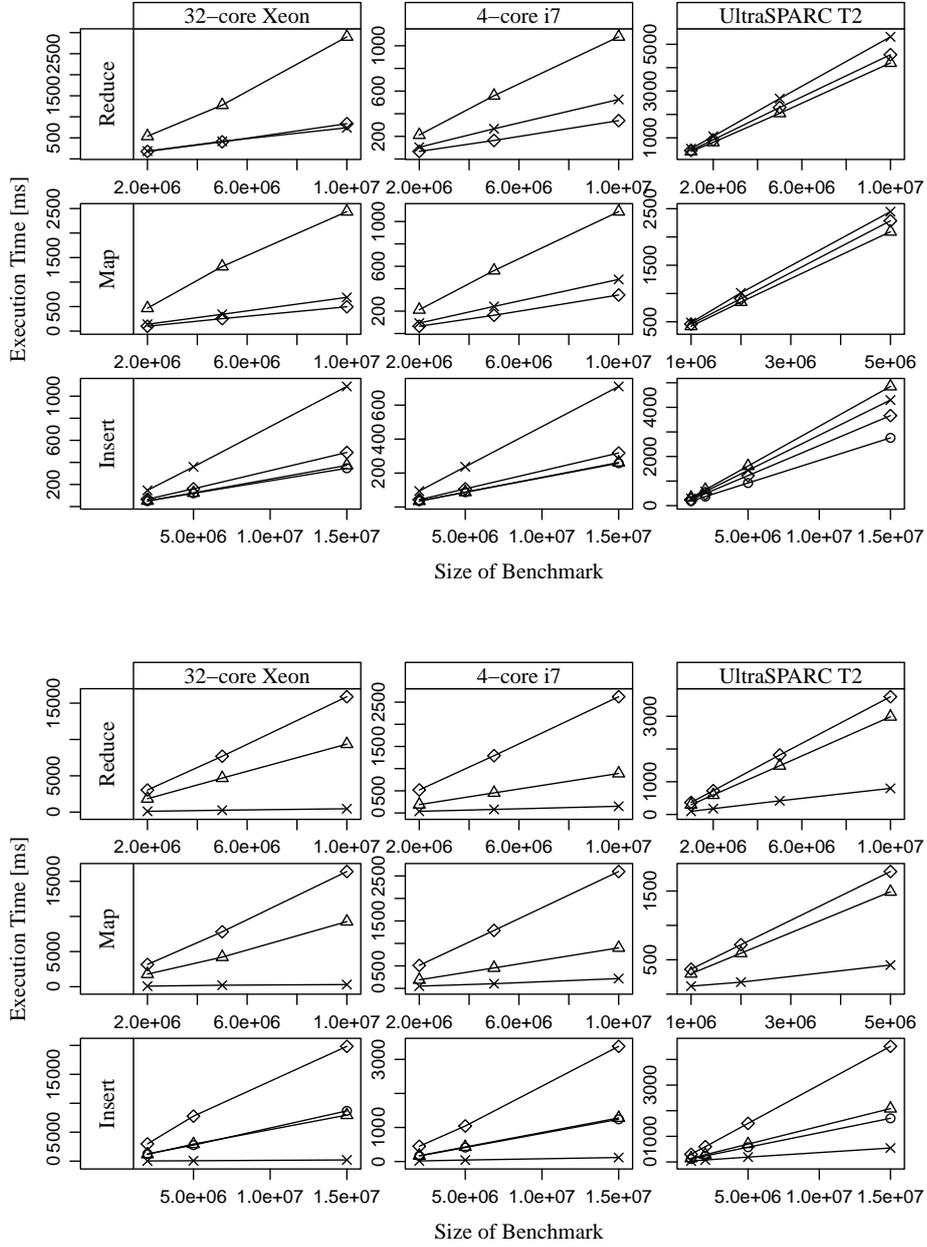


Fig. 6. Execution time vs benchmark size ($P = 1, 8$). \diamond single-lane FlowPool, \times multi-lane FlowPool, \triangle linked transfer queue, \circ concurrent linked queue

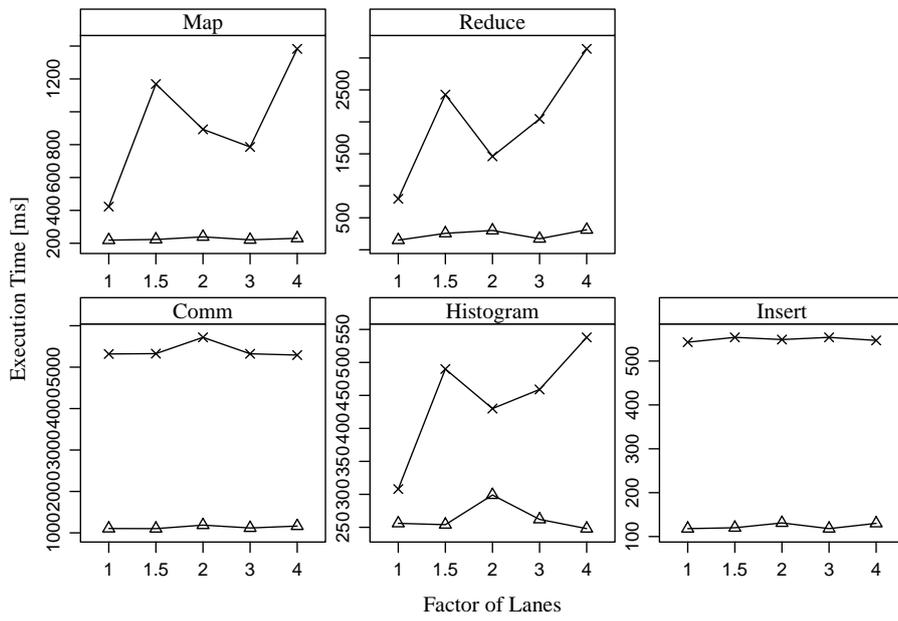


Fig. 7. Execution times vs lane-factor for multi-lane FlowPools. × UltraSPARC T2, △ 4-core i7