

# FlowPools: Lock-Free Deterministic Concurrent Data-Flow Queues

Authors

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

**Abstract.** Implementing correct and deterministic parallel programs is hard (additional motivation sentence after this). We present the design and implementation of a fundamental data structure for deterministic parallel data-flow computation. Additionally, we provide a proof of correctness, showing that the implementation is linearizable, lock-free, and deterministic. Finally, we provide microbenchmarks which compare our *flow pools* against corresponding operations on other popular concurrent data structures, in addition to performance benchmarks on a real *XYZ* application using real data. (Keep abstract between 70 and 150 words).

**Keywords:** data-flow, concurrent data-structure, determinism

## 1 Introduction

### 1.1 Motivation

- we want a deterministic model - we do not block in the programming model (i.e. there are no operations which cause blocking until a value becomes available) - we want a non-blocking data-structure (i.e. the operations on the data-structures should themselves be non-blocking) - programs run indefinitely => we need to GC parts of the data structure we no longer need - we want to reduce heap allocation and inline the datastructure as much as possible => lower memory consumption, better cache behaviour and fewer GC cycles

Obligatory multicore motivation paragraph.

Lock-free is better, and why.

Introduction and motivation for data-flow programming model.

## 2 Model of Computation

Producer-consumer parallelism. Description and image of queue/stream of values, producer, and multiple consumers.

### 3 Programming Model

The FlowPool suite supports the following operations:

- `Builder.<<(x: T): Builder`  
Inserts an element into the underlying FlowPool.
- `Builder.seal(n: Int): Unit`  
Seals the underlying FlowPool at `n` elements. The need for the size argument is explained below. A sealed FlowPool may only contain `n` elements. This allows for callback cleanup and termination.
- `FlowPool.doForAll(f: T => Unit): Future[Int]`  
Instructs the FlowPool to execute the closure `f` exactly once for each element inserted into the FlowPool (asynchronously). The returned future contains the number of elements in the FlowPool and completes once `f` has been executed for all elements.
- `FlowPool.mappedFold[U, V <: U](acc: V)(cmb: (U,V) => V)(map: T => U): Future[(Int, V)]` Reduces the FlowPool to a single value of type `V`, by first mapping each element to an internal representation using `map` and then aggregating using `cmb`. No guarantee is given about synchronization or order. Returns a future containing a tuple with number of elements in the FlowPool and the aggregated value.
- `FlowPool.builder`  
Returns a builder for this FlowPool.
- `Future.map[U](f: T => U)`  
Maps this future to another future executing the function `f` exactly once when the first future completes.
- `future[T](f: () => T): Future[T]`  
Asynchronously dispatch execution of `f` and return a future with its result.

*Determinism of seal* We will show that the final size of the FlowPool is required as an argument to the `seal` method in order to satisfy the determinism property of the FlowPool. Look at the following program:

```
val p = new FlowPool[Int]()
val b = p.builder

future {
  for (i <- 1 to 10) { b << i }
  b.seal
}

future { for (i <- 1 to 10) { b << i } }
```

Depending on which for-loop completes first, this program completes successfully or yields an error. A similar program with `b.seal(20)` will always succeed.

*Generators* In the following we'll present a couple of generators for FlowPools based on common generators in the Scala standard library.

```
def iterate[T](start: T, len: Int)(f: (T) => T) = {
  val p = new FlowPool[T]
  val b = p.builder
  future {
    var e = start
    for (i <- 1 to len) { b << e; e = f(e) }
    b.seal(len)
  }; p
}
```

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

*Monadic Operations* In the following we'll present some monadic operations on top of the basic FlowPool operations. This will also show some use-cases of the futures as result type of `doForAll`.

```
def map[S](f: T => S) = {
  val fp = new FlowPool[S]
  val b = fp.builder
  doForAll { x =>
    b << f(x)
  } map { b.seal _ }
  fp
}
```

```

def filter(f: T => Boolean) = {
  val fp = new FlowPool[T]
  val b = fp.builder

  mappedFold(0)(_ + _) { x =>
    if (f(x)) { b << x; 1 } else 0
  } map { case (c,fc) => b.seal(fc) }

  fp
}

def flatMap[S](f: T => FlowPool[S]) = {
  val fp = new FlowPool[S]
  val b = fp.builder

  mappedFold(future(0))(_ <+> _) { x =>
    f(x).doForAll(b << _)
  } map { case (c,cfut) => cfut.map(b.seal _) }

  fp
}

```

where `<+>` is the future of the sum of two `Future[Int]`.

## 4 Implementation

## 5 Proofs

### 5.1 Abstract Pool Semantics

Alex

### 5.2 Linearizability

Alex

### 5.3 Lock-Freedom

### 5.4 Determinism

## 6 Experimental Results

## 7 Related Work

Things to probably cite: Oz, gpars, Java CLQ, our futures writeup.

Things we should probably have a look at: Microsoft TPL, Dataflow Java, FlumeJava...

Forcing a bib, [?], [?], [?], [?], [?], [?]

## 8 Conclusion

```

1  def append(elem: Elem)
2    b = READ(current)
3    idx = READ(b.index)
4    nextobj = READ(b.array(idx + 1))
5    curobj = READ(b.array(idx))
6    if (check(b, idx, curobj, nextobj))
7      if (CAS(b.array(idx + 1), nextobj, curobj))
8        if (CAS(b.array(idx), curobj, elem))
9          WRITE(b.index, idx + 1)
10         invokeCallbacks(elem, curobj)
11       else append(elem)
12     else append(elem)
13   else
14     advance()
15     append(elem)
16
17
18 def check(b: Block, idx: Int, curobj: Object, nextobj: Object)
19   // The check on the index is done implicitly in the real code
20   if (idx > LASTELEMPOS) return false
21   else curobj match
22     elem: Elem =>
23       return false
24     term: Terminal =>
25       if (term.sealed == NOTSEALED) return true
26     else
27       if (totalElems(b, idx) < term.sealed) return true
28       else error("sealed")
29
30   null =>
31     error("unreachable")
32
33
34 def advance()
35   b = READ(current)
36   idx = READ(b.index)
37   if (idx > LASTELEMPOS) expand(b)
38   else
39     obj = READ(b.array(idx))
40     if (obj is Elem) WRITE(b.index, idx + 1)
41
42
43 def expand(b: Block)
44   nb = READ(b.next)
45   if (nb is null)
46     nb = createBlock(b.blockindex + 1)
47     if (CAS(b.next, null, nb))
48       expand(b)
49   else
50     CAS(current, b, nb)
51
52
53 def totalElems(b: Block, idx: Int)
54   return b.blockindex * (BLOCKSIZE - 1) + idx
55
56 def invokeCallbacks(elem: Elem, term: Terminal)
57   for (f <- term.callbacks) future
58     f(elem)
59
60
61 def seal(size: Int)
62   b = READ(current)
63   idx = READ(b.index)
64   if (idx <= LASTELEMPOS)
65     curobj = READ(b.array(idx))

```

```

66     curobj match
67       term: Terminal =>
68         tryWriteSeal(term, b, idx, size)
69     elem: Elem =>
70       WRITE(b.index, idx + 1)
71       seal(size)
72     null =>
73       error("unreachable")
74
75   else
76     expand(b)
77     seal(size)
78
79
80   def tryWriteSeal(term: Terminal, b: Block, idx: Int, size: Int)
81     val total = totalElems(b, idx)
82     if (total > size) error("too many elements")
83     if (term.sealed == NOTSEALED)
84       nterm = new Terminal
85       sealed = size
86       callbacks = term.callbacks
87
88       CAS(b.array(idx), term, nterm)
89     else if (term.sealed != size)
90       error("already sealed with different size")
91
92
93   def doForAll(f: Elem => Unit)
94     future
95       asyncDoForAll(f, start, 0)
96
97
98   def asyncDoForAll(f: Elem => Unit, b: Block, idx: Int)
99     if (idx <= LASTELEMPOS)
100       obj = READ(b.array(idx))
101       obj match
102         term: Terminal =>
103           nterm = new Terminal
104           sealed = term.sealed
105           callbacks = f :: term.callbacks
106
107           if (!CAS(b.array(idx), term, nterm)) asyncDoForAll(f, b, idx)
108       elem: Elem =>
109         f(elem)
110         asyncDoForAll(f, b, idx + 1)
111       null =>
112         error("unreachable")
113
114   else
115     // In the real code we take a shortcut when preparing the new block
116     expand(b)
117     asyncDoForAll(f, b.next, 0)
118
119

```

**Definition 1 (FlowPool).** *The **FlowPool** is defined as the references start and current. The **FlowPool state** is defined as the configuration of objects transitively reachable from the reference start.*

*We define the following relations:*

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

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

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

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

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

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  $\text{willBeCalled}(\text{start}, e, f)$  holds.

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

*TODO formal definitions of these relations (involving the flowpool datatypes)*

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

**Definition 2 (Invariants).** We define the following invariants for the **Flow-Pool**:

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

**INV2**  $\forall b \in \text{reachable}(\text{start}), b \notin \text{following}(b)$

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

**INV4**  $\forall b = \text{last}(\text{start}), b.\text{array} = p \cdot t \cdot n$ , where:

$$p = X^P, t = t_1 \cdot t_2, n = \text{null}^N$$

$$x \in \text{Elem}, t_1 \in \text{Terminal}, t_2 \in \{\text{null}\} \cup \text{Terminal}$$

$$P + N + 2 = \text{BLOCKSIZE}$$

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

**Definition 3 (Validity).** A FlowPool state  $\mathbb{S}$  is **valid** if and only if the invariants [INV1-5] hold for that state.

**Definition 4 (Abstract pool).** An **abstract pool**  $\mathbb{P}$  is a function from time  $t$  to a tuple of sets  $(\text{elems}, \text{callbacks}, \text{seal})$  such that:

$$\text{seal} \in \{\emptyset, \{w\}\}, w \in \mathbb{N}$$

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

$$\text{called} \subseteq \text{elems} \subseteq \text{Elem}$$

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

**Definition 5 (Abstract pool operations).** We say that an abstract pool operation  $op$  applied to an abstract pool  $P_0 = (\text{elems}_0, \text{callbacks}_0, \text{seal}_0)$  at some time  $t$  **changes** the state of the abstract pool to  $P = (\text{elems}, \text{callbacks}, \text{seal})$  if  $\exists t_0, \forall \tau, t_0 < \tau < t, \mathbb{P}(\tau) = P_0$  and  $\mathbb{P}(t) = P$ .

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

$\exists t_1 > t_0, \forall t_2 > t_1, \mathbb{P}(t_2) = (elems_2, (f, called_2) \cup callbacks_2, seal_2)$ , where:

$called_2 \subseteq elems_2 \subseteq elems$

Append...

Seal...

**Definition 6 (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:

$\forall e \in Elem, hasElem(start, e) \Leftrightarrow e \in elems$

$\forall f \in Elem \Rightarrow Unit, \forall e \in Elem, willBeCalled(start, e, f) \Leftrightarrow \exists t_1 \geq t_0, \mathbb{P}(t_1) = (elems_1, (f, called_1) \cup callbacks_1, seal_1), elems \subseteq called_1$

$\forall s \in \mathbb{N}, sealedAt(start, s) \Leftrightarrow s \in seal$

A *FlowPool* operation *op* completing at some time  $t_0$  is consistent with an abstract pool operation *op'* if and only if *op* changes the state of the *FlowPool* from  $\mathbb{S}_1$  to  $\mathbb{S}_2$ , where  $\mathbb{S}_1$  and  $\mathbb{S}_2$  are consistent with the abstract pool states  $\mathbb{A}_1$  and  $\mathbb{A}_2$ , respectively, and *op'* changes the state of the abstract pool from  $\mathbb{A}_1$  to  $\mathbb{A}_2$ .

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

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

**Lemma 1 (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 8 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 8 writes a value of type *Elem*, no other write can change it.

**Lemma 2 (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 1 it follows that  $\forall t \geq t_1, b.array(idx - 1) \in Elem$ .

**Corollary 1 (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 lemmas 1 and 2, and the fact that the CAS in line 8 only writes to array entries  $idx$  for which it previously read the value from  $b.index$ .



**Definition 7 (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:*

$$f : v_0 \rightarrow v_1$$

**Lemma 3 (Freshness).** *For all blocks  $b \in \text{reachable}(\text{start})$ , and for all  $x \in b.\text{array}$ , function  $x$  has unique elements in its string of transitions.*

*Proof.* CAS instruction in line 8 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 88 and 107, shows that their expected values are of type *Terminal* or are *null*. Their new values are always freshly allocated.

The more difficult part is to prove that CAS instruction in line 7 respects the statement of the lemma.

Since the CAS instructions in lines 88 and 107 are preceded by a read of  $\text{idx} = b.\text{index}$ , from 2 it follows that  $b.\text{array}(\text{idx} - 1)$  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.

A successful CAS in line 7 overwrites a value  $cb_0$  at  $\text{idx} + 1$  read in line 4 at  $t_0$  with a new value  $cb_2$  at time  $t_2$ . Value  $cb_2$  was read in line 5 at  $t_1$  from the entry  $\text{idx}$ . The string of transitions of values at  $\text{idx}$  is composed of unique values at least since  $t_1$ , since there is a value of type *Elem* at the index  $\text{idx} - 1$ .

Now assume that a thread T1 successfully CAS in line 7 at  $\text{idx} + 1$  at time  $t_2$  overwrites  $cb_0$  with a value  $cb_2$  read from  $\text{idx}$  at  $t_1$ , and that another thread T2 subsequently successfully completes the CAS in line 7 at  $\text{idx} + 1$  at time  $t_2$  with a value  $cb_{\text{prev}}$  which was at  $\text{idx} + 1$  at some time  $t < t_0$ . That would mean that the thread T2 read the value  $cb_{\text{prev}}$  in line 5 at some time  $t_{\text{prev}1} < t_1$  and successfully completed the CAS at time  $t_{\text{prev}2} > t_2$ . If the CAS was successful, then the read in line 4 by T2 occurred at  $t_{\text{prev}0} < t_{\text{prev}1} < t_1$ . TODO finish this

**Lemma 4 (Lifecycle).** *For all blocks  $b \in \text{reachable}(\text{start})$ , and for all  $x \in b.\text{array}$ , function  $x$  goes through and only through the prefix of the following transitions:*

$$\text{null} \rightarrow cb_1 \rightarrow \dots \rightarrow cb_n \rightarrow \text{elem}, \text{ where:}$$

$$cb_i \in \text{Terminal}, i \neq j \Rightarrow cb_i \neq cb_j, \text{elem} \in \text{Elem}$$

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

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

The CAS in line 8 has a value  $\text{curobj} \in \text{Terminal}$  as an expected value and writes an  $\text{elem} \in \text{Elem}$ . This means that the only transition that this CAS can cause is of type  $cb_i \in \text{Terminal} \rightarrow \text{elem} \in \text{Elem}$ .

We will now prove that the CAS in line 7 at time  $t_2$  is successful if and only if the entry at  $idx + 1$  is *null* or  $nextobj \in Terminal$ . 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 4, because of lemma 3 and the fact that CAS in line 7 is assumed to be successful. We know that during the read in line 5 at time  $t_1$ , such that  $t_0 < t_1 < t_2$ , the entry at  $idx$  was  $curobj \in Terminal$ , by trivial analysis of the **check** procedure. It follows from corollary 1 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 7 always overwrites a one value of type *Terminal* (or *null*) with another value of type *Terminal*. We still have to show that it never overwrites the value  $cb_0$  with a value  $cb_2$  that was at  $b.array(idx)$  at an earlier time. This follows from the 3 directly.

Finally, note that the statement for CAS instructions in lines 88 and 107 follows directly from the proof for lemma 3.

**Lemma 5 (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  $\mathbb{S}$ . We analyze each write.

**Lemma 6 (Housekeeping).** *Given a FlowPool in state  $\mathbb{S}$  consistent with some abstract pool state  $\mathbb{A}$ , CAS instructions in lines 7, 47 and 50 do not change the state of the abstract pool.*

*Proof.*

**Lemma 7 (Append correctness).** *Given a FlowPool in state  $\mathbb{S}$  consistent with some abstract pool state  $\mathbb{A}$ , a successful CAS in line 8 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:*

$$\mathbb{A} = (elems, callbacks, seal)$$

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

*Furthermore, given a fair scheduler, there exists a time  $t_1 > t_0$  such that all the callback functions are called for elem.*

*Proof.*

**Definition 8 (Obstruction-freedom).** *Given a FlowPool in a valid state, any operation*

**Lemma 8 (Obstruction-free operations).** *The FlowPool operations are obstruction-free.*

*Proof.* By trivial sequential code analysis.