

The Weakest Failure Detectors to Boost Obstruction-Freedom*

Rachid Guerraoui^{1,2} Michał Kapalka² Petr Kouznetsov³

¹ Computer Science and Artificial Intelligence Laboratory, MIT

² School of Computer and Communication Sciences, EPFL

³ Max Planck Institute for Software Systems

13th July 2006

Abstract

This paper determines necessary and sufficient conditions to implement *wait-free* and *non-blocking* contention managers in a shared memory system. The necessary conditions hold even when universal objects (like compare-and-swap) or random oracles are available, whereas the sufficient ones assume only registers.

We show that failure detector $\diamond\mathcal{P}$ is the weakest to convert any obstruction-free algorithm into a wait-free one, and Ω^* , a new failure detector which we introduce in this paper, and which is strictly weaker than $\diamond\mathcal{P}$ but strictly stronger than Ω , is the weakest to convert any obstruction-free algorithm into a non-blocking one.

1 Introduction

Multiprocessor systems are becoming more and more common nowadays. Multithreading thus becomes the norm and studying scalable and efficient synchronization methods is essential, for traditional locking-based techniques do not scale and may induce priority inversion, deadlock and fault-tolerance issues when a large number of threads is involved.

Wait-free synchronization algorithms [13] circumvent the issues of locking and guarantee individual progress even in presence of high contention. *Wait-freedom* is a liveness property which stipulates that every process completes every op-

eration in a finite number of its own steps, regardless of the status of other processes, i.e., contending or even crashed. Ideal synchronization algorithms would ensure *linearizability* [16, 2], a safety property which provides the illusion of instantaneous operation executions, together with *wait-freedom*.

Alternatively, a liveness property called *non-blockingness*¹ may be considered instead of *wait-freedom*. *Non-blockingness* guarantees global progress, i.e., that some process will complete an operation in a finite number of steps, regardless of the behavior of other processes. *Non-blockingness* is weaker than *wait-freedom* as it does not prevent some processes from starvation.

Wait-free and *non-blocking* algorithms are, however, notoriously difficult to design [18, 3], especially with the practical goal to be fast in low contention scenarios, which are usually considered the most common in practice. An appealing principle to reduce this difficulty consists in separating two concerns of a synchronization algorithm: (1) ensuring linearizability with a minimal conditional progress guarantee, and (2) boosting progress. More specifically, the idea is to focus on algorithms that ensure linearizability together with a weak liveness property called *obstruction-freedom* [15], and then combine these algorithms with separate generic oracles that boost progress, called *contention managers* [14, 20, 21, 9]. This separation lies at the heart of modern (obstruction-

*EPFL Technical Report LPD-REPORT-2006-007. Elements of this work are to appear in a paper with the same title in the Proceedings of the 20th International Symposium on Distributed Computing (DISC'06).

¹The term *non-blocking* is defined here in the traditional way [13]: "some process will complete its operation in a finite number of steps, regardless of the relative execution speeds of the processes." This term is sometimes confused with the term *lock-free*. Note that *non-blocking* implementations provide a weaker liveness guarantee than *wait-free* implementations.

free) software transactional memory (STM) frameworks [14].

With obstruction-free (or OF, for short) algorithms, progress is ensured only for every process that executes in isolation for sufficiently long time. In presence of high contention, however, OF algorithms can livelock, preventing any process from terminating. Contention managers are used precisely to cope with such scenarios. When queried by a process executing an OF algorithm, a contention manager can delay the process for some time in order to boost the progress of other processes. The contention manager can neither share objects with the OF algorithm, nor return results on its behalf. If it did, the contention manager could peril the safety of the OF algorithm, hampering the overall separation of concerns principle.

In short, the goal of a contention manager is to provide processes with enough time without contention so that they can complete their operations. In its simplest form, a contention manager can be a randomized back-off protocol. More sophisticated contention management strategies have been experimented in practice [20, 21, 10]. Precisely because they are entirely devoted to progress, they can be combined or changed on the fly [9]. Most previous strategies were *pragmatic*, with no aim to provide *worst case guarantees*. In this paper we focus on contention managers that provide such guarantees. More specifically, we study contention managers that convert any OF algorithm into a non-blocking or wait-free one, and which we call, respectively, *non-blocking* or *wait-free* contention managers.

Two wait-free contention managers have recently been proposed [6, 8]. Both rely on timing assumptions to detect processes that fail in the middle of their operations. This suggests that *some* information about failures might inherently be needed by any wait-free contention manager. But this is not entirely clear because, in principle, a contention manager could also use randomization to schedule processes, or even powerful synchronization primitives like compare-and-swap, which is known to be *universal*, i.e., able to wait-free implement any other object [13]. In the parlance of [5], we would like to determine whether a *failure detector* is actually needed to implement a contention manager with worst case guarantees, and if it is, what is the *weakest* one [4]. Besides the theoretical interest, determining the minimal conditions under which a contention manager can ensure certain guarantees is, we believe, of practical

relevance, for this might help portability and optimization.

We show that the eventually perfect failure detector $\diamond\mathcal{P}$ [5] is the weakest to implement a wait-free contention manager.² We also introduce a failure detector Ω^* , which we show is the weakest to implement a non-blocking contention manager. Failure detector Ω^* is strictly weaker than $\diamond\mathcal{P}$, and strictly stronger than failure detector Ω [4], known to be the weakest to wait-free implement the (universal) consensus object [13].³

It might be surprising that Ω is not sufficient to implement a wait-free or even a non-blocking contention manager. For example, the seminal Paxos algorithm [19] uses Ω to transform an OF implementation of consensus into a wait-free one. Each process that is eventually elected a leader by Ω is given enough time to run alone, reach a decision and communicate it to the others. This approach does not help, however, if we want to make sure that processes make progress regardless of the actual (possibly long-lived) object and its OF implementation. Intuitively, the leader elected by Ω may have no operation to perform while other processes may livelock forever. Because a contention manager cannot make processes help each other, the output of Ω is not sufficient: this is so even if randomized oracles or universal objects are available. Intuitively, wait-free contention managers need a failure detector that would take care of *every* non-crashed process with a pending operation so that the process can run alone for sufficiently long time. As for non-blocking contention managers, at least *one* process that never crashes, among the ones with pending operations, should be given enough time to run alone.

The paper is organized as follows. Section 2 presents our system model and formally defines wait-free and non-blocking contention managers. These definitions are, we believe, contributions in their own rights, for they capture precisely the interaction between a contention manager and an obstruction-free algorithm. In Sect. 3 and 4, we prove our weakest failure detector results. In each case, we first present (necessary part) a *reduction* algorithm [4] that *extracts* the output of failure detector Ω^* (respectively $\diamond\mathcal{P}$) using a non-blocking (respectively wait-free) contention manager im-

² $\diamond\mathcal{P}$ ensures that eventually: (1) every failure is detected by every correct (i.e., non-faulty) process and (2) there is no false detection.

³ Ω ensures that eventually all correct (i.e., non-faulty) processes elect the same correct process as their leader.

plementation. When devising our reduction algorithms, we do not restrict what objects (or random oracles) can be used by the contention manager or the OF algorithm. Then (sufficient part), we present algorithms that implement the contention managers using the failure detectors and registers. These algorithms are devised with the sole purpose of proving our sufficiency claims. We do not seek to minimize the overhead of the interaction between the OF algorithm and the contention manager, nor do we discuss how the failure detector can itself be implemented with little synchrony assumptions and minimal overhead, unlike the transformations presented in [6]. However, as we show in [11], our algorithms can be easily extended to meet these challenges.

2 Preliminaries

Processes and Failure Detectors. We consider a set of n processes $\Pi = \{p_1, \dots, p_n\}$ in a shared memory system [13, 17]. A process executes the (possibly randomized) algorithm assigned to it, until the process *crashes (fails)* and stops executing any action. We assume the existence of a global discrete clock that is, however, inaccessible to the processes. We say that a process is *correct* if it never crashes. We say that process p_i is *alive* at time t if p_i has not crashed by time t .

A *failure detector* [5, 4] is a distributed oracle that provides every process with some information about failures. The output of a failure detector depends only on which and when processes fail, and not on computations being performed by the processes. A process p_i queries a failure detector \mathcal{D} by accessing local variable $\mathcal{D}\text{-output}_i$ —the output of the module of \mathcal{D} at process p_i . Failure detectors can be partially ordered according to the amount of information about failures they provide. A failure detector \mathcal{D} is *weaker than a failure detector \mathcal{D}'* , and we write $\mathcal{D} \preceq \mathcal{D}'$, if there exists an algorithm (called a *reduction* algorithm) that transforms \mathcal{D}' into \mathcal{D} . If $\mathcal{D} \preceq \mathcal{D}'$ but $\mathcal{D}' \not\preceq \mathcal{D}$, we say that \mathcal{D} is *strictly weaker than \mathcal{D}'* , and we write $\mathcal{D} \prec \mathcal{D}'$.

Base and High-Level Objects. Processes communicate by invoking primitive operations (which we will call *instructions*) on *base* shared objects and seek to implement the *operations* of a *high-level* shared object O . Object O is in turn used by an application, as a high-level inter-process communication mechanism. We call invocation and re-

sponse events of a high-level operation op on the implemented object O *application events* and denote them by, respectively, $inv(op)$ and $ret(op)$ (or $inv_i(op)$ and $ret_i(op)$ at a process p_i).

An *implementation* of O is a distributed algorithm that specifies, for every process p_i and every operation op of O , the sequences of *steps* that p_i should take in order to complete op . Process p_i *completes* operation op when p_i returns from op . Every process p_i may complete any number of operations but, at any point in time, at most one operation op can be *pending* (started and not yet completed) at p_i .

We consider implementations of O that combine a sub-protocol that ensures a minimal liveness property, called *obstruction-freedom*, with a sub-protocol that boosts this liveness guarantee. The former is called an *obstruction-free (OF)* algorithm A and the latter a *contention manager* CM . We focus on *linearizable* [16, 2] implementations of O : every operation appears to the application as if it took effect instantaneously between its invocation and its return. An implementation of O involves two categories of steps executed by any process p_i : those (executed on behalf) of CM and those (executed on behalf) of A . In each step, a process p_i either executes an instruction on a base shared object or (in case p_i executes a step on behalf of CM) queries a failure detector.

Obstruction-freedom [15, 14] stipulates that if a process that invokes an operation op on object O and from some point in time executes steps of A alone⁴, then it eventually completes op . *Non-blockingness* stipulates that if some correct process never completes an invoked operation, then some other process completes infinitely many operations. *Wait-freedom* [13] ensures that every correct process that invokes an operation eventually returns from the operation.

Interaction Between Modules. OF algorithm A , executed by any process p_i , communicates with contention manager CM via *calls* try_i and *resign* _{i} implemented by CM (see Fig. 1). Process p_i invokes try_i just after p_i starts an operation, and also later (even several times before p_i completes the operation) to signal possible contention. Process p_i invokes $resign_i$ just before returning from an operation, and always eventually returns from this call (or crashes). Both calls, try_i and $resign_i$, return *ok*.

⁴I.e., without encountering *step contention* [1].

An example OF algorithm that uses this model of interaction with a contention manager is presented in Algorithm 1. The algorithm implements a timestamping mechanism and is based on the implementation of a splitter. It is not meant to be practical or efficient—it just shows how calls *try* and *resign* should be used.

A discussion about overhead of wait-free/non-blocking contention managers that explains when calls to *try/resign* can be omitted for efficiency reasons can be found in [11].

Algorithm 1: An example OF algorithm implementing a timestamping mechanism

uses: $A[1, \dots]$ —unbounded array of registers, $B[1, \dots]$ —unbounded array of single-bit registers, L —a register

initially: $A[1, \dots] \leftarrow \perp$, $B[1, \dots] \leftarrow false$, $L \leftarrow 1$

```

1.1 upon of-getTimestamp do
1.2   CM.tryi
1.3   j ← L
1.4   while true do
1.5     A[j] ← i
1.6     if B[j] = false then
1.7       B[j] ← true
1.8       if A[j] = i then
1.9         L ← j
1.10        CM.resigni
1.11        return j
1.12   CM.tryi
1.13   j ← j + 1

```

We denote by $B(A)$ and $B(CM)$ the sets of base shared objects, always *disjoint*, that can be possibly accessed by steps of, respectively, A and CM , in every execution, by every process. Calls *try* and *resign* are thus the only means by which A and CM interact. The events corresponding to invocations of, and responses from, *try* and *resign* are called *cm-events*. We denote by try_i^{inv} and $resign_i^{inv}$ an invocation of call try_i and $resign_i$, respectively (at process p_i), and by try_i^{ret} and $resign_i^{ret}$ —the corresponding responses.

Executions and Histories. An *execution* of an OF algorithm A combined with a contention manager CM is a sequence of *events* that include steps of A , steps of CM , cm-events and application events. Every event in an execution is associated with a unique time at which the event took place. Every execution e induces a *history* $H(e)$ that in-

cludes only application events (invocations and responses of high-level operations). The corresponding *CM-history* $H_{CM}(e)$ is the subsequence of e containing only application events and cm-events of the execution, and the corresponding *OF-history* $H_{OF}(e)$ is the subsequence of e containing only application events, cm-events, and steps of A . For a sequence s of events, $s|i$ denotes the subsequence of s containing only events at process p_i .

We say that a process p_i is *blocked* at time t in an execution e if (1) p_i is alive at time t , and (2) the latest event in $H_{CM}(e)|i$ that occurred before t is try_i^{inv} or $resign_i^{inv}$. A process p_i is *busy* at time t in e if (1) p_i is alive at time t , and (2) the latest event in $H_{CM}(e)|i$ that occurred before t is try_i^{ret} . We say that a process p_i is *active* at t in e if p_i is either busy or blocked at time t in e . We say that a process p_i is *idle* at time t in e if p_i is not active at t in e .⁵ A process *resigns* when it invokes *resign* on a contention manager.

We say that p_i is *obstruction-free* in an interval $[t, t']$ in an execution e , if p_i is the only process that takes steps of A in $[t, t']$ in e and p_i is not blocked infinitely long in $[t, t']$ (if $t' = \infty$). We say that process p_i is *eventually obstruction-free* at time t in e if p_i is active at t or later and p_i either resigns after t or is obstruction-free in the interval $[t', \infty)$ for some $t' > t$. Note that, since algorithm A is obstruction-free, if an active process p_i is eventually obstruction-free, then p_i eventually resigns and completes its operation.

Well-Formed Executions. We impose certain restrictions on the way an OF algorithm A and a contention manager CM interact. In particular, we assume that no process takes steps of A while being blocked by CM or idle, and no process takes infinitely many steps of A without calling CM infinitely many times. Further, a process must inform CM that an operation is completed by calling *resign* before returning the response to the application.

Formally, we assume that every execution e is *well-formed*, i.e., $H(e)$ is linearizable [16, 2], and, for every process p_i , (1) $H_{CM}(e)|i$ is a prefix of a sequence $[op_1][op_2], \dots$, where each $[op_k]$ has the form $inv_i(op_k), try_i^{inv}, try_i^{ret}, \dots, try_i^{inv}, try_i^{ret}, resign_i^{inv}, resign_i^{ret}, ret_i(op_k)$; (2) in $H_{OF}(e)|i$, no step of A is executed when p_i is blocked or idle, (3) in $H_{OF}(e)|i$, inv_i can only be followed by try_i^{inv} , and

⁵Note that every process that has crashed is permanently idle.

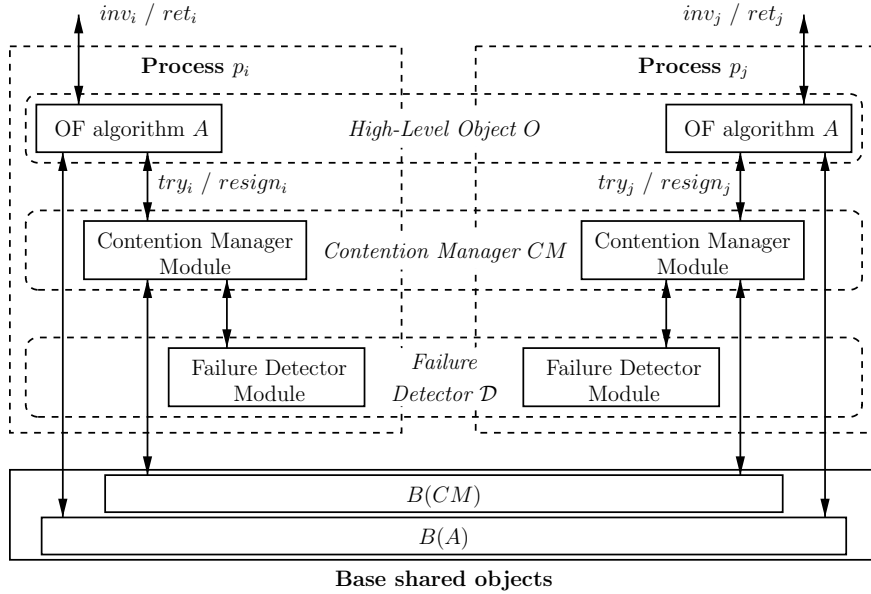


Figure 1: The OF algorithm/contention manager interface

ret_i can only be preceded by $resign_i^{ret}$; (4) if p_i is busy at time t in e , then at some $t' > t$, process p_i is idle or blocked. The last condition implies that every busy process p_i eventually invokes try_i (and becomes blocked), resigns or crashes. Clearly, in a well-formed execution, every process goes through the following cyclical order of modes: *idle, active, idle, ...*, where each *active* period consists itself of a sequence *blocked, busy, blocked, ...*

Non-blocking Contention Manager. We say that a contention manager CM *guarantees non-blockingness for an OF algorithm A* if in each execution e of A combined with CM the following property is satisfied: if some correct process is active at a time t , then at some time $t' > t$ some process resigns.

A *non-blocking contention manager* guarantees non-blockingness for every OF algorithm. Intuitively, this will happen if the contention manager allows at least one active process to be obstruction-free (and busy) for sufficiently long time, so that the process can complete its operation. More precisely, we say that a contention manager CM is *non-blocking* if, for every OF algorithm A , in every execution of A combined with CM the following property is ensured at every time t :

Global Progress. If some correct process is active at t , then some correct process is eventually obstruction-free at t .

Theorem 1 *A contention manager CM guarantees non-blockingness for every OF algorithm if and only if CM is non-blocking.*

Proof. (\Rightarrow) Consider a contention manager CM that guarantees non-blockingness for every OF algorithm. Let A be any OF algorithm and e be any execution of A combined with CM . Let some correct process be active at time t in e . Since CM guarantees non-blockingness, some active process resigns at some future time, and the Global Progress property is trivially ensured.

(\Leftarrow) By contradiction, assume that there exists a non-blocking contention manager CM such that, for some OF algorithm A , there is an execution e of A combined with CM , such that some correct process is active at t , and no active process resigns after t . By Global Progress, some correct active process p_i is eventually obstruction-free at t . Since A is obstruction-free and p_i takes infinitely many steps of A in isolation, p_i must complete its operation and resign—a contradiction. \square

Wait-Free Contention Manager. We say that a contention manager CM *guarantees wait-freedom for an OF algorithm A* if in every execution e of A combined with CM , the following property is satisfied: if a process p_i is active at a time t , then at some time $t' > t$, p_i becomes idle. In other words, every operation executed by a correct process eventually returns.

A *wait-free contention manager* guarantees wait-

freedom for every OF algorithm. Intuitively, this will happen if the contention manager makes sure that every correct active process is given “enough” time to complete its operation, regardless of how other processes behave. More precisely, a contention manager CM is wait-free if, for every OF algorithm A , in every execution of A combined with CM , the following property is ensured at every time t :⁶

Fairness. If a correct process p_i is active at t , then p_i is eventually obstruction-free at t .

Theorem 2 *A contention manager CM guarantees wait-freedom for every OF algorithm if and only if CM is wait-free.*

Proof. (\Rightarrow) Consider a contention manager CM that guarantees wait-freedom for every OF algorithm. Let A be any OF algorithm and e be any execution of A combined with CM . Since in e every active process is eventually idle, every correct active process eventually resigns in e , and so the Fairness property is trivially satisfied.

(\Leftarrow) Let CM be a wait-free contention manager, and A be any OF algorithm. Consider any execution e of A combined with CM .

Suppose, by contradiction, that some correct process p_i is active at time t and never completes its operation thereafter. But then, by Fairness, p_i is eventually obstruction-free at t and so p_i is obstruction-free in period $[t', \infty)$ for some $t' > t$. Therefore, since A is obstruction-free and p_i takes infinitely many steps of A in isolation, p_i must eventually resign and complete its operation—a contradiction. \square

In the following, we seek to determine the *weakest* [4] failure detector \mathcal{D} to implement a non-blocking (resp. wait-free) contention manager CM . This means that (1) \mathcal{D} implements such a contention manager, i.e., there is an algorithm that implements CM using \mathcal{D} , and (2) \mathcal{D} is *necessary* to implement such a contention manager, i.e., if a failure detector \mathcal{D}' implements CM , then $\mathcal{D} \preceq \mathcal{D}'$. In our context, a reduction algorithm that transforms \mathcal{D}' into \mathcal{D} uses the \mathcal{D}' -based implementation of the corresponding contention manager as a “black box” and read-write registers.

⁶This property is ensured by wait-free contention managers from the literature [6, 8].

3 Non-blocking Contention Managers

Let $S \subseteq \Pi$ be a non-empty set of processes. Failure detector Ω_S outputs, at every process, an identifier of a process (called a *leader*), such that all correct processes in S eventually agree on the identifier of the same *correct* process in S .⁷

Failure detector Ω^* is the composition $\{\Omega_S\}_{S \subseteq \Pi, S \neq \emptyset}$: at every process p_i , Ω^* -output $_i$ is a tuple consisting of the outputs of failure detectors Ω_S . We position Ω^* in the hierarchy of failure detectors of [5] by proving the following theorem:

Theorem 3 $\Omega \prec \Omega^* \prec \diamond\mathcal{P}$.

Proof. It is immediate that Ω is weaker than Ω^* : Ω_Π is equivalent to Ω . In a system of three or more processes, Ω is strictly weaker than Ω^* . Indeed, consider a system of three processes, p_1 , p_2 , and p_3 , and assume, by contradiction, that Ω^* is weaker than Ω , i.e., that there exists a reduction algorithm $T_{\Omega \rightarrow \Omega^*}$ which extracts the output of Ω^* using Ω . Take an execution e of $T_{\Omega \rightarrow \Omega^*}$ in which p_3 is correct, p_2 is faulty, Ω always outputs p_3 at every process and consider the emulated output of $\Omega_{\{p_1, p_2\}}$. Since p_1 is the only correct process in $\{p_1, p_2\}$, there is a finite prefix e' of e in which $\Omega_{\{p_1, p_2\}}$ outputs p_1 at p_1 . But this finite execution is indistinguishable from a finite execution e'' in which p_2 is correct but slow. Now consider a finite extension of e'' in which p_1 fails, and thus eventually $\Omega_{\{p_1, p_2\}}$ outputs p_2 at p_2 . But this finite execution is indistinguishable from a finite execution in which p_1 is correct but slow. By repeating this argument, we obtain an infinite execution of $T_{\Omega \rightarrow \Omega^*}$ in which both p_1 and p_2 are correct, and the output $\Omega_{\{p_1, p_2\}}$ never stabilizes at a single correct process—a contradiction.

It is immediate that Ω^* is weaker than $\diamond\mathcal{P}$: eventually each correct process p_i has complete and accurate information about failures of all other processes, so p_i can perform an eventually perfect leader election in each subset of processes p_i belongs to.

To show that Ω^* is *strictly* weaker than $\diamond\mathcal{P}$, consider a system of two processes, p_1 and p_2 , and assume, by contradiction, that $\diamond\mathcal{P}$ is weaker than

⁷ Ω_S can be seen as a restriction of the eventual leader election failure detector Ω [4] to processes in S . The definition of Ω_S resembles the notion of Γ -accurate failure detectors introduced in [12]. Clearly, Ω_Π is Ω .

Ω^* , i.e., that there exists a reduction algorithm $T_{\Omega^* \rightarrow \diamond \mathcal{P}}$ which extracts the output of $\diamond \mathcal{P}$ using Ω^* .

Using $T_{\Omega^* \rightarrow \diamond \mathcal{P}}$, we implement $\diamond \mathcal{P}$ in the *asynchronous system*, establishing a contradiction with [7, 5]. In the implementation, the processes run two parallel algorithms, T_1 and T_2 . The algorithm T_i ($i = 1, 2$) is identical to $T_{\Omega^* \rightarrow \diamond \mathcal{P}}$, except that, instead of querying Ω^* , it assumes that Ω^* always outputs p_i at every process. Note that every finite execution of T_i is also a finite execution of $T_{\Omega^* \rightarrow \diamond \mathcal{P}}$. If p_i is correct, then every (even infinite) execution of T_i is also an execution of $T_{\Omega^* \rightarrow \diamond \mathcal{P}}$. Thus, in both cases, p_i obtains a valid output of $\diamond \mathcal{P}$.

Hence, we obtain an implementation of $\diamond \mathcal{P}$, in an asynchronous system of two processes, contradicting [7, 5]. \square

To show that Ω^* is necessary to implement a non-blocking contention manager, it suffices to prove that, for every non-empty $S \subseteq \Pi$, Ω_S is necessary to implement a non-blocking contention manager. Let CM be a non-blocking contention manager using failure detector \mathcal{D} . We show that $\Omega^* \preceq \mathcal{D}$ by presenting an algorithm $T_{\mathcal{D} \rightarrow \Omega_S}$ (Algorithm 2) that, using CM and \mathcal{D} , emulates the output of Ω_S .

The algorithm works as follows. Every process $p_i \in S$ runs two parallel tasks T_i and F_i . In task T_i , process p_i periodically (1) gets blocked by CM after invoking try_i (line 2.5), and (2) once p_i gets busy again, announces itself a leader for set S by writing its id in L (line 2.6). In task F_i , process p_i periodically determines its leader by reading register L (line 2.2).⁸

Thus, no process ever resigns and every correct process in S is permanently active from some point in time. Intuitively, this signals a possible livelock to CM which has to eventually block all active processes except for one that should run obstruction-free for sufficiently long time. By Global Progress, CM cannot block *all* active processes forever and so if the elected process crashes (and so becomes idle), CM lets another active process run obstruction-free. Eventually, all correct processes in S agree on the same process in S . Processes outside S are permanently idle and permanently output their own ids: they do not access CM .

This approach contains a subtlety. To make sure that there is a time after which the same correct leader in S is permanently elected by the correct

processes in S , we do not allow the elected leader to resign (the output of Ω_S has to be eventually stable). This violates the assumption that processes using CM run an obstruction-free algorithm, and, thus, a priori, CM is not obliged to preserve Global Progress. However, as we show below, since CM does not “know” how much time a process executing an OF algorithm requires to complete its operation, CM has to provide some correct process with *unbounded* time to run in isolation.

Theorem 4 *Every non-blocking contention manager can be used to implement failure detector Ω^* .*

Proof. Let $S \subseteq \Pi$, $S \neq \emptyset$ and consider any execution of Algorithm 2. If S contains no correct process, then Ω_S -output $_i$ (for every process $p_i \in S$) trivially satisfies the property of Ω_S . Now assume that there is a correct process in S . We claim that CM eventually lets exactly one correct process in S run obstruction-free while blocking forever all the other processes in S .

Suppose not. We obtain an execution in which every correct process in S is allowed to be obstruction-free only for bounded periods of time. But the CM -history of this execution corresponds to an execution of some OF algorithm A combined with CM in which no active process ever completes its operation because no active process ever obtains enough time to run in isolation. Thus, no active process is eventually obstruction-free in that execution. This contradicts the assumption that CM is non-blocking.

Therefore, there is a time after which exactly one correct process $p_j \in S$ is periodically busy (others are blocked or idle forever) and, respectively, register L permanently stores the identifier of p_j . Thus, eventually, every correct process in S outputs p_j : the output of Ω_S is extracted. \square

We describe an implementation of a non-blocking contention manager using Ω^* and registers in Algorithm 3. The algorithm works as follows. All active processes, upon calling try , participate in the leader election mechanism using Ω^* in lines 3.3–3.5. The active process p_i that is elected a leader returns from try and is (eventually) allowed to run obstruction-free until p_i resigns. Once p_i resigns, the processes elect another leader. Failure detector Ω^* guarantees that if an active process is elected and crashes before resigning, another active process is eventually elected.

Lemma 5 *Contention manager shown in Algorithm 3 guarantees non-blockingness for every OF algorithm.*

⁸If a process is blocked in one task, it continues executing steps in parallel tasks.

Algorithm 2: Extracting Ω_S from a non-blocking contention manager (code for processes from set S ; others are permanently idle)

uses: L —register

initially: $\Omega_S\text{-output}_i \leftarrow p_i, L \leftarrow \text{some process in } S$

Launch two parallel tasks: T_i and F_i

2.1 **parallel task** F_i

2.2 $\lfloor \Omega_S\text{-output}_i \leftarrow L$

2.3 **parallel task** T_i

2.4 \lfloor **while** *true* **do**

2.5 $\lfloor \lfloor$ issue *try* _{i} and wait until busy (i.e., until call *try* _{i} returns)

2.6 $\lfloor \lfloor L \leftarrow p_i$ // announce yourself a leader

Algorithm 3: A non-blocking contention manager using $\Omega^* = \{\Omega_S\}_{S \subseteq \Pi, S \neq \emptyset}$

uses: $T[1, \dots, n]$ —array of single-bit registers

initially: $T[1, \dots, n] \leftarrow \text{false}$

3.1 **upon** *try* _{i} **do**

3.2 $\lfloor T[i] \leftarrow \text{true}$

3.3 **repeat**

3.4 $\lfloor \lfloor S \leftarrow \{p_j \in \Pi \mid T[j] = \text{true}\}$

3.5 $\lfloor \lfloor$ **until** $\Omega_S\text{-output}_i = p_i$

3.6 **upon** *resign* _{i} **do**

3.7 $\lfloor T[i] \leftarrow \text{false}$

Proof. Assume, by contradiction that there exists an OF algorithm A for which contention manager CM implemented by Algorithm 3 does not guarantee non-blockingness, i.e., there exists an execution e of A combined with CM in which there are a correct process p_i and a time t , such that p_i is active at t but for all $t' > t$, no active process resigns at t' .

Take any time $t' > t$. Let us denote by $S(t')$ the set of all processes p_j such that $T[j] = \text{true}$ at time t' in e . Since no active process resigns after t , there is a time $t^* \geq t$ and a set S , such that for all $t' > t^*$, $S(t') = S$. By the algorithm, p_i eventually sets $T[j]$ to *true*. Thus, p_i is in S , i.e., S includes at least one correct process. At every correct process in S , Ω_S eventually outputs the same correct process p_j in set S (a leader).

Since every active process eventually invokes *try*, resigns or crashes (by the properties of OF algorithms), and no process resigns after t^* , there is a time $t' > t^*$ after which every correct process except for p_j gets permanently blocked in lines 3.3–3.5. That is because p_j does not resign after t and so p_j does not reset $T[j]$ to *false* thereafter and remains

the leader for set S forever. Thus, p_j is eventually obstruction-free at t . Since p_j runs an obstruction-free algorithm A , it eventually resigns and completes its operation—a contradiction. \square

From Theorem 1 and Lemma 5 we immediately obtain a proof of the following theorem:

Theorem 6 *Algorithm 3 implements a non-blocking contention manager.*

4 Wait-Free Contention Managers

We prove here that the weakest failure detector to implement a wait-free contention manager is $\diamond\mathcal{P}$. Failure detector $\diamond\mathcal{P}$ [5] outputs, at each time and every process, a set of *suspected* processes. There is a time after which (1) every crashed process is permanently suspected by every correct process and (2) no correct process is ever suspected by any correct process.

We first consider a wait-free contention manager CM using a failure detector \mathcal{D} , and we exhibit a reduction algorithm $T_{\mathcal{D} \rightarrow \diamond\mathcal{P}}$ (Algorithm 4) that, using CM and \mathcal{D} , emulates the output of $\diamond\mathcal{P}$.

We run several instances of CM . These instances use disjoint sets of base shared objects and do not directly interact. Basically, in each instance, only two processes are active and all other processes are idle. One of the two processes, say p_j , gets active and never resigns thereafter, while the other, say p_i , permanently alternates between being active and idle. To CM it looks like p_j is always obstructed by p_i . Thus, to guarantee wait-freedom, the instance of CM has to eventually block p_i and let p_j run obstruction-free until p_j resigns or crashes. Therefore, when p_i is blocked, p_i

Algorithm 4: Extracting $\diamond\mathcal{P}$ from a wait-free contention manager

uses: $R[1, \dots, n]$ —array of registers

initially: $\diamond\mathcal{P}\text{-output}_i \leftarrow \Pi - \{p_i\}$, $k \leftarrow 0$, $R[i] \leftarrow 0$

Launch $n(n-1)$ parallel instances of CM: C_{jk} , $j, k \in \{1, \dots, n\}$, $j \neq k$

Launch $2n-1$ parallel tasks: T_{ij} , T_{ji} , $j \in \{1, \dots, n\}$, $i \neq j$, and F_i

4.1 **parallel task** F_i

4.2 \lfloor **while true do** $R[i] \leftarrow R[i] + 1$ // ‘‘heartbeat’’ signal

4.3 **parallel task** T_{ij} , $j = 1, \dots, i-1, i+1, \dots, n$

4.4 \lfloor **while true do**

4.5 \lfloor $x_j \leftarrow R[j]$

4.6 \lfloor $\diamond\mathcal{P}\text{-output}_i \leftarrow \diamond\mathcal{P}\text{-output}_i - \{p_j\}$ // stop suspecting p_j

4.7 \lfloor issue try_i^{ij} (in C_{ij}) and wait until busy

4.8 \lfloor issue resign_i^{ij} (in C_{ij}) and wait until idle

4.9 \lfloor $\diamond\mathcal{P}\text{-output}_i \leftarrow \diamond\mathcal{P}\text{-output}_i \cup \{p_j\}$ // start suspecting p_j

4.10 \lfloor wait until $R[j] > x_j$ // wait until p_j takes a new step

4.11 **parallel task** T_{ji} , $j = 1, \dots, i-1, i+1, \dots, n$

4.12 \lfloor **while true do** issue try_i^{ji} (in C_{ji}) and wait until busy

can assume that p_j is alive and when p_i is busy, p_i can suspect p_j of having crashed, until p_i eventually observes p_j 's ‘‘heartbeat’’ signal, which p_j periodically broadcasts using a register. This ensures the properties of $\diamond\mathcal{P}$ at process p_i , provided that p_j never resigns.

As in Sect. 3, we face the following issue. If p_j is correct, p_i will be eventually blocked forever and p_j will thus be eventually obstruction-free. Hence, in the corresponding execution, obstruction-freedom is violated, i.e., the execution cannot be produced by any OF algorithm combined with CM. One might argue then that CM is not obliged to preserve Fairness with respect to p_j . However, we show that, since CM does not ‘‘know’’ how much time a process executing an OF algorithm requires to complete its operation, CM has to provide p_j with *unbounded* time to run in isolation.

More precisely, the processes in Algorithm 4 run $n(n-1)$ parallel instances of CM, denoted each CM_{jk} , where $j, k \in \{1, \dots, n\}$, $j \neq k$. We denote the events that process p_i issues in instance CM_{jk} by try_i^{jk} and resign_i^{jk} . Besides, every process p_i runs $2n-1$ parallel tasks: T_{ij} , T_{ji} , where $j \in \{1, \dots, n\}$, $i \neq j$, and F_i . Every task T_{ij} executed by p_i is responsible for detecting failures of process p_j . Every task T_{ji} executed by p_i is responsible for preventing p_j from falsely suspecting p_i . In task F_i , p_i periodically writes ever-increasing ‘‘heartbeat’’

values in a shared register $R[i]$.

In every instance CM_{ij} , there can be only two active processes: p_i and p_j . Process p_i cyclically gets active (line 4.7) and resigns (line 4.8), and process p_j gets active once and keeps getting blocked (line 4.12). Each time before p_i gets active, p_i removes p_j from the list of suspected processes (line 4.6). Each time p_i stops being blocked, p_i starts suspecting p_j (line 4.9) and waits until p_i observes a ‘‘new’’ step of p_j (line 4.10). Once such a step of p_j is observed, p_i stops suspecting p_j and gets active again.

Theorem 7 *Every wait-free contention manager can be used to implement failure detector $\diamond\mathcal{P}$.*

Proof. Consider any execution e of $T_{\mathcal{D} \rightarrow \diamond\mathcal{P}}$, and let p_i be any correct process. We show that, in e , $\diamond\mathcal{P}\text{-output}_i$ satisfies the properties of $\diamond\mathcal{P}$, i.e., p_i eventually permanently suspects every non-correct process and stops suspecting every correct process. (Note that if a process p_i is not correct, then $\diamond\mathcal{P}\text{-output}_i$ trivially satisfies the properties of $\diamond\mathcal{P}$.)

Let p_j be any process distinct from p_i . Assume p_j is not correct. Thus p_i is the only correct active process in instance CM_{ij} . By the Fairness property of CM, p_i is eventually obstruction-free every time p_i becomes active, and so p_i cannot be blocked infinitely long in line 4.7. Since there is a time after which p_j stops taking steps, eventually p_i starts

suspecting p_j (line 4.9) and suspends in line 4.10, waiting until p_j takes a new step. Thus, p_i eventually suspects p_j forever.

Assume now that p_j is correct. We claim that p_i must eventually get permanently blocked so that p_j would run obstruction-free from some point in time forever. Suppose not. But then we obtain an execution in which p_i alternates between active and idle modes infinitely many times, and p_j stays active and runs obstruction-free only for bounded periods of time. But the CM-history of this execution could be produced by an execution e' of some OF algorithm combined with CM in which p_j never completes its operation because p_j never runs long enough in isolation. Thus, Fairness is violated in execution e' and this contradicts the assumption that CM is wait-free. Hence, eventually p_i gets permanently blocked in line 4.7. Since each time p_i is about to get blocked, p_i stops suspecting p_j in line 4.6, there is a time after which p_i never suspects p_j .

Thus, there is a time after which, if p_j is correct, then p_j stops being suspected by every correct process, and if p_j is non-correct, then every correct process permanently suspects p_j . \square

We describe an implementation of a wait-free contention manager using $\diamond\mathcal{P}$ and registers in Algorithm 5. The algorithm relies on a (wait-free) primitive $GetTimestamp()$ that generates unique, locally increasing timestamps and makes sure that if a process gets a timestamp ts , then no process can get timestamps lower than ts infinitely many times (this primitive can be implemented in an asynchronous system using read-write registers). The idea of the algorithm is the following. Every process p_i that gets active receives a timestamp in line 5.2 and announces the timestamp in register $T[i]$. Every active process that invokes try repeatedly runs a leader election mechanism (lines 5.3–5.6): the non-suspected (by $\diamond\mathcal{P}$) process that announced the lowest (non- \perp) timestamp is elected a leader. If a process p_i is elected, p_i returns from try_i and becomes busy. $\diamond\mathcal{P}$ guarantees that eventually the same correct active process is elected by all active processes. All other active processes stay blocked until the process resigns and resets its timestamp in line 5.8. The leader executes steps obstruction-free then. Since the leader runs an OF algorithm, the leader eventually resigns and resets its timestamp in line 5.8 so that another active process, which now has the lowest timestamp in T , can become a leader.

Lemma 8 *Contention manager implemented by Algorithm 5 guarantees wait-freedom for all OF algorithms.*

Proof. Consider an execution e of any OF algorithm A combined with contention manager CM implemented by Algorithm 5. By contradiction, assume that in e , some correct process is active at some time t , and never resigns after t . Let V denote the non-empty set of correct processes that are active at t but never resign (in line 5.8) and complete their operations thereafter, i.e., that remain active after t forever. Recall that every process in V either invokes try infinitely many times or invokes try and stays blocked forever (by the properties of OF algorithms). Let $t^* > t$ be time at which every process in V invoked try and reached line 5.3 at least once. For every $p_j \in V$, let ts_j^* denote the value of $T[j]$ at time t^* . Note that since every $ts_j^* \neq \perp$ and no process in V resigns after time t^* , $T[j] = ts_j^*$ at all times $t' \geq t^*$.

Let p_i be the process in V having the lowest timestamp in $\{ts_k^* \mid p_k \in V\}$ (there is exactly one such process since timestamps are unique). We establish a contradiction by showing that p_i has to eventually resign.

Let us consider time $t' > t^*$ after which:

- at every correct process, failure detector $\diamond\mathcal{P}$ outputs the list of all non-correct processes (by the properties of $\diamond\mathcal{P}$, this eventually happens),
- all non-correct processes have crashed,
- for every correct process $p_j \neq p_i$, if $T[j] \neq \perp$, then $T[j] > ts_i^*$.

The last condition eventually holds, because timestamps are unique, no process can receive a timestamp lower than ts_i^* infinitely many times and p_i has the lowest timestamp among processes in V (that retain their timestamps infinitely long).

Thus, after t' , p_i is always elected a leader, and every correct process p_j other than p_i that gets blocked after time t' will remain blocked in lines 3.3–3.5, as long as p_i does not resign.

Hence, eventually p_i will be the only active process that is not blocked and thus p_i will be given unbounded time to perform steps of A in isolation. Since A is obstruction-free, p_i eventually resigns and completes its operation—a contradiction. \square

From Theorem 2 and Lemma 8 we immediately obtain a proof of the following theorem:

Theorem 9 *Algorithm 5 implements a wait-free contention manager.*

Algorithm 5: A wait-free contention manager using $\diamond\mathcal{P}$

uses: $T[1, \dots, N]$ —array of registers (other variables are local)

initially: $T[1, \dots, N] \leftarrow \perp$

```
5.1 upon  $try_i$  do
5.2   if  $T[i] = \perp$  then  $T[i] \leftarrow GetTimestamp()$ 
5.3   repeat
5.4      $sact_i \leftarrow \{j \mid T[j] \neq \perp \wedge p_j \notin \diamond\mathcal{P}\text{-output}_i\}$ 
5.5      $leader_i \leftarrow \operatorname{argmin}_{j \in sact_i} T[j]$ 
5.6   until  $leader_i = i$ 

5.7 upon  $resign_i$  do
5.8    $T[i] \leftarrow \perp$ 
```

Acknowledgements.

We are very grateful to Hagit Attiya, Maurice Herlihy, Bastian Pochon, Faith Fich, Victor Luchangco, Mark Moir and Nir Shavit for interesting discussions on the topic of this paper. We would also like to thank the anonymous reviewers of DISC'06 for helpful comments.

References

- [1] H. Attiya, R. Guerraoui, and P. Kouznetsov. Computing with reads and writes in the absence of step contention. In *Proceedings of the 19th International Symposium on Distributed Computing (DISC'05)*, 2005.
- [2] H. Attiya and J. L. Welch. *Distributed Computing: Fundamentals, Simulations and Advanced Topics (2nd edition)*. Wiley, 2004.
- [3] B. N. Bershad. Practical considerations for non-blocking concurrent objects. In *Proceedings of the 14th IEEE International Conference on Distributed Computing Systems (ICDCS'93)*, pages 264–273, 1993.
- [4] T. D. Chandra, V. Hadzilacos, and S. Toueg. The weakest failure detector for solving consensus. *Journal of the ACM*, 43(4):685–722, July 1996.
- [5] T. D. Chandra and S. Toueg. Unreliable failure detectors for reliable distributed systems. *Journal of the ACM*, 43(2):225–267, March 1996.
- [6] F. Fich, V. Luchangco, M. Moir, and N. Shavit. Obstruction-free algorithms can be practically wait-free. In *Proceedings of the 19th International Symposium on Distributed Computing (DISC'05)*, 2005.
- [7] M. J. Fischer, N. A. Lynch, and M. S. Paterson. Impossibility of distributed consensus with one faulty process. *Journal of the ACM*, 32(3):374–382, April 1985.
- [8] R. Guerraoui, M. Herlihy, M. Kapalka, and B. Pochon. Robust contention management in software transactional memory. In *Proceedings of the Workshop on Synchronization and Concurrency in Object-Oriented Languages (SCOOL); in conjunction with the ACM Conference on Object-Oriented Programming, Systems, Languages and Applications (OOPSLA'05)*, October 2005.
- [9] R. Guerraoui, M. Herlihy, and B. Pochon. Polymorphic contention management. In *Proceedings of the 19th International Symposium on Distributed Computing (DISC'05)*, pages 303–323. LNCS, Springer, September 2005.
- [10] R. Guerraoui, M. Herlihy, and B. Pochon. Toward a theory of transactional contention managers. In *Proceedings of the 24th Annual ACM Symposium on Principles of Distributed Computing (PODC'05)*, 2005.
- [11] R. Guerraoui, M. Kapalka, and P. Kouznetsov. Boosting obstruction-freedom with low overhead. Technical report, EPFL, 2006. Submitted for publication.
- [12] R. Guerraoui and A. Schiper. “ Γ -accurate” failure detectors. In *Proceedings of the 10th International Workshop on Distributed Algorithms (WDAG'96)*. Springer-Verlag, 1996.

- [13] M. Herlihy. Wait-free synchronization. *ACM Transactions on Programming Languages and Systems*, 13(1):124–149, January 1991.
- [14] M. Herlihy, V. Luchangco, M. Moir, and W. N. Scherer III. Software transactional memory for dynamic-sized data structures. In *Proceedings of the 22nd Annual ACM Symposium on Principles of Distributed Computing (PODC'03)*, pages 92–101, 2003.
- [15] M. Herlihy, V. Luchangco, and M. Moir. Obstruction-free synchronization: Double-ended queues as an example. In *Proceedings of the 23rd IEEE International Conference on Distributed Computing Systems (ICDCS'93)*, pages 522–529, 2003.
- [16] M. Herlihy and J. M. Wing. Linearizability: a correctness condition for concurrent objects. *ACM Transactions on Programming Languages and Systems*, 12(3):463–492, June 1990.
- [17] P. Jayanti. Robust wait-free hierarchies. *Journal of the ACM*, 44(4):592–614, 1997.
- [18] A. LaMarca. A performance evaluation of lock-free synchronization protocols. In *Proceedings of the 13th Annual ACM Symposium on Principles of Distributed Computing (PODC'94)*, pages 130–140, 1994.
- [19] L. Lamport. The part-time parliament. *ACM Transactions on Computer Systems*, 16(2):133–169, May 1998.
- [20] W. N. Scherer III and M. L. Scott. Contention management in dynamic software transactional memory. In *PODC Workshop on Concurrency and Synchronization in Java Programs*, July 2004.
- [21] W. N. Scherer III and M. L. Scott. Advanced contention management for dynamic software transactional memory. In *Proceedings of the 24th Annual ACM Symposium on Principles of Distributed Computing (PODC'05)*, 2005.