# Analyzing the Performance of Lock-Free Data Structures: A Conflict-based Model

Aras Atalar, Paul Renaud-Goud and Philippas Tsigas

Chalmers University of Technology

{aaras|goud|tsigas}@chalmers.se

### Abstract

This paper considers the modeling and the analysis of the performance of lock-free concurrent data structures. Lock-free designs employ an optimistic conflict control mechanism, allowing several processes to access the shared data object at the same time. They guarantee that at least one concurrent operation finishes in a finite number of its own steps regardless of the state of the operations. Our analysis considers such lock-free data structures that can be represented as linear combinations of fixed size retry loops.

Our main contribution is a new way of modeling and analyzing a general class of lock-free algorithms, achieving predictions of throughput that are close to what we observe in practice. We emphasize two kinds of conflicts that shape the performance: (i) hardware conflicts, due to concurrent calls to atomic primitives; (ii) logical conflicts, caused by simultaneous operations on the shared data structure.

We show how to deal with these hardware and logical conflicts separately, and how to combine them, so as to calculate the throughput of lock-free algorithms. We propose also a common framework that enables a fair comparison between lock-free implementations by covering the whole contention domain, together with a better understanding of the performance impacting factors. This part of our analysis comes with a method for calculating a good back-off strategy to finely tune the performance of a lock-free algorithm. Our experimental results, based on a set of widely used concurrent data structures and on abstract lock-free designs, show that our analysis follows closely the actual code behavior.

## CONTENTS

# I. INTRODUCTION

Lock-free programming provides highly concurrent access to data and has been increasing its footprint in industrial settings. Providing a modeling and an analysis framework capable of describing the practical performance of lock-free algorithms is an essential, missing resource necessary to the parallel programming and algorithmic research communities in their effort to build on previous intellectual efforts. The definition of lock-freedom mainly guarantees that at least one concurrent operation on the data structure finishes in a finite number of its own steps, regardless of the state of the operations. On the individual operation level, lock-freedom cannot guarantee that an operation will not starve.

The goal of this paper is to provide a way to model and analyze the practically observed performance of lock-free data structures. In the literature, the common performance measure of a lock-free data structure is the throughput, *i.e.* the number of successful operations per unit of time. It is obtained while threads are accessing the data structure according to an access pattern that interleaves local work between calls to consecutive operations on the data structure. Although this access pattern to the data structure is significant, there is no consensus in the literature on what access to be used when comparing two data structures. So, the amount of local work (that we will refer as parallel work for the rest of the paper) could be constant ([MS96], [SL00]), uniformly distributed ([HSY10], [DLM13]), exponentially distributed ([Val94], [DB08]), null ([KH14], [LJ13]), *etc.*, and more questionably, the average amount is rarely scanned, which leads to a partial covering of the contention domain.

We propose here a common framework enabling a fair comparison between lock-free data structures, while exhibiting the main phenomena that drive performance, and particularly the contention, which leads to different kinds of conflicts. As this is the first step in this direction, we want to deeply analyze the core of the problem, without impacting factors being diluted within a probabilistic smoothing. Therefore, we choose a constant local work, hence constant access rate to the data structures. In addition to the prediction of the data structure performance, our model provides a good back-off strategy, that achieves the peak performance of a lock-free algorithm.

Two kinds of conflict appear during the execution of a lock-free algorithm, both of them leading to additional work. Hardware conflicts occur when concurrent operations call atomic primitives on the same data: these calls collide and conduct to stall time, that we name here *expansion*. Logical conflicts take place if concurrent operations overlap: because of the lock-free nature of the algorithm, several concurrent operations can run simultaneously, but only one retry can logically succeed. We show that the additional work produced by the failures is not necessarily harmful for the system-wise performance.

We then show how throughput can be computed by connecting these two key factors in an iterative way. We start by estimating the expansion probabilistically, and emulate the effect of stall time introduced by the hardware conflicts as extra work added to each thread. Then we estimate the number of failed operations, that in turn lead to additional extra work, by computing again the expansion on a system setting where those two new amounts of work have been incorporated, and reiterate the process; the convergence is ensured by a fixed-point search.

We consider the class of lock-free algorithms that can be modeled as a linear composition of fixed size retry loops. This class covers numerous extensively used lock-free designs such as stacks [Tre86] (Pop, Push), queues [MS96] (Enqueue, Dequeue), counters [DLM13] (Increment, Decrement) and priority queues [LJ13] (DeleteMin).

To evaluate the accuracy of our model and analysis framework, we performed experiments both on synthetic tests, that capture a wide range of possible abstract algorithmic designs, and on several reference implementations of extensively studied lock-free data structures. Our evaluation results reveal that our model is able to capture the behavior of all the synthetic and real designs for

all different numbers of threads and sizes of parallel work (consequently also contention). We also evaluate the use of our analysis as a tool for tuning the performance of lock-free code by selecting the appropriate back-off strategy that will maximize throughput by comparing our method with against widely known back-off policies, namely linear and exponential.

The rest of the paper is organized as follows. We discuss related work in Section II, then the problem is formally described in Section III. We consider the logical conflicts in the absence of hardware conflicts in Section IV, while in Section V, we firstly show how to compute the expansion, then combine hardware and logical conflicts to obtain the final throughput estimate. We describe the experimental results in Section VI.

## II. Related Work

Anderson *et al.* [ARJ97] evaluated the performance of lock-free objects in a single processor real-time system by emphasizing the impact of retry loop interference. Tasks can be preempted during the retry loop execution, which can lead to interference, and consequently to an inflation in retry loop execution due to retries. They obtained upper bounds for the number of interferences under various scheduling schemes for periodic real-time tasks.

Intel [Int13] conducted an empirical study to illustrate performance and scalability of locks. They showed that the critical section size, the time interval between releasing and re-acquiring the lock (that is similar to our parallel section size) and number of threads contending the lock are vital parameters.

Failed retries do not only lead to useless effort but also degrade the performance of successful ones by contending the shared resources. Alemany *et al.* [AF92] have pointed out this fact, that is in accordance with our two key factors, and, without trying to model it, have mitigated those effects by designing non-blocking algorithms with operating system support.

Alistarh *et al.* [ACS14] have studied the same class of lock-free structures that we consider in this paper. The analysis is done in terms of scheduler steps, in a system where only one thread can be scheduled (and can then run) at each step. If compared with execution time, this is particularly appropriate to a system with a single processor and several threads, or to a system where the instructions of the threads cannot be done in parallel (*e.g.* multi-threaded program on a multi-core processor with only read and write on the same cache line of the shared memory). In our paper, the execution is evaluated in terms of processor cycles, strongly related to the execution time. In addition, the "parallel work" and the "critical work" can be done in parallel, and we only consider retry-loops with one Read and one CAS, which are serialized. In addition, they bound the asymptotic expected system latency (with a big O, when the number of threads tends to infinity), while in our paper we estimate the throughput (close to the inverse of system latency) for any number of threads.

## III. Problem Statement

### A. Running Program and Targeted Platform

In this paper, we aim at evaluating the throughput of a multi-threaded algorithm that is based on the utilization of a shared lock-free data structure. Such a program can be abstracted by the Procedure AbstractAlgorithm (see Figure 1) that represents the skeleton of the function which is called by each spawned thread. It is decomposed in two main phases: the *parallel section*, represented on line 3, and the *retry loop*, from line 4 to line 7. A *retry* starts at line 5 and ends at line 7.

As for line 1, the function Initialization shall be seen as an abstraction of the delay between the spawns of the threads, that is expected not to be null, even when a barrier is used. We then consider that the threads begin at the exact same time, but have different initialization times.
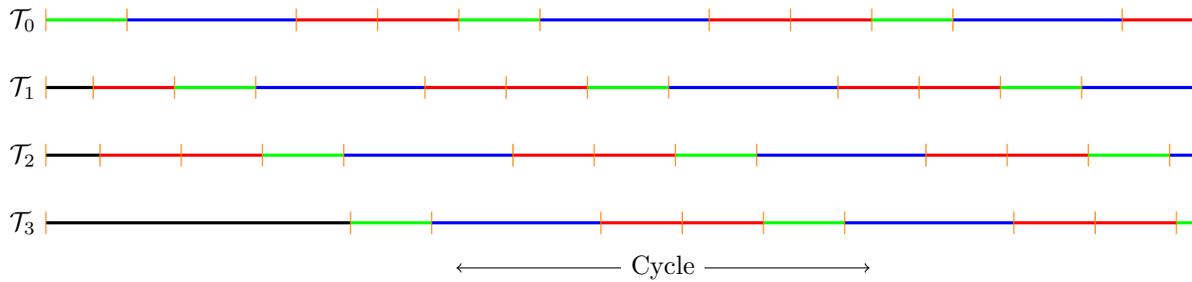
The parallel section is the part of the code where the thread does not access the shared data structure; the work that is performed inside this parallel section can possibly depend on the value

**Procedure** AbstractAlgorithm

---

**1** Initialization();
**2 while** *!* done **do**
**3** | Parallel_Work();
**4** | **while** *!* success **do**
**5** | | current ← Read(AP);
**6** | | new ← Critical_Work(current);
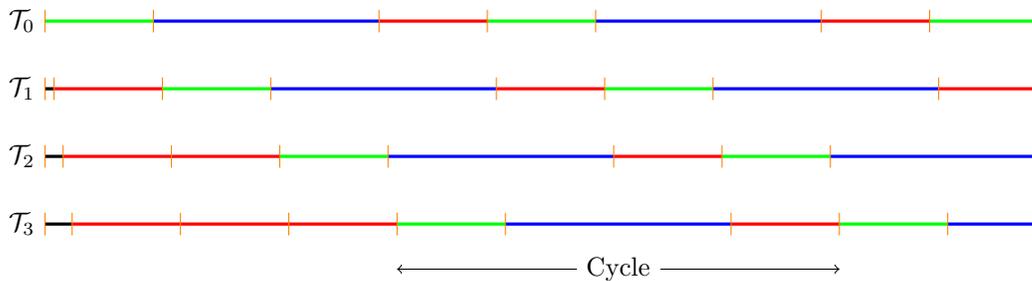**7** | | success ← CAS(AP, current, new);

---

**Figure 1:** Thread procedure



**Figure 2:** Execution with one wasted retry, and one inevitable failure

that has been read from the data structure, *e.g.* in the case of processing an element that has been dequeued from a FIFO (First-In-First-Out) queue.

In each retry, a thread tries to modify the data structure, and does not exit the retry loop until it has successfully modified the data structure. It does that by firstly reading the access point AP of the data structure, then according to the value that has been read, and possibly to other previous computations that occurred in the past, the thread prepares the new desired value as an access point of the data structure. Finally, it atomically tries to perform the change through a call to the *Compare-And-Swap* (*CAS*) primitive. If it succeeds, *i.e.* if the access point has not been changed by another thread between the first *Read* and the *CAS*, then it goes to the next parallel section, otherwise it repeats the process. The retry loop is composed of at least one retry, and we number the retries starting from 0, since the first iteration of the retry loop is actually not a retry, but a try.

We analyze the behavior of AbstractAlgorithm from a throughput perspective, which is defined



**Figure 3:** Execution with minimum number of failures

as the number of successful data structure operations per unit of time. In the context of Procedure AbstractAlgorithm, it is equivalent to the number of successful *CAS*s.

The throughput of the lock-free algorithm, that we denote by $T$, is impacted by several parameters.
- *Algorithm parameters*: the amount of work inside a call to Parallel_Work (resp. Critical_Work) denoted by $pw$ (resp. $cw$).
- *Platform parameters*: *Read* and *CAS* latencies ($rc$ and $cc$ respectively), and the number $P$ of processing units (cores). We assume homogeneity for the latencies, *i.e.* every thread experiences the same latency when accessing an uncontended shared data, which is achieved in practice by pinning threads to the same socket.

### B. Examples and Issues

We first present two straightforward upper bounds on the throughput, and describe the two kinds of conflict that keep the actual throughput away from those upper bounds.

*1) Immediate Upper Bounds:* Trivially, the minimum amount of work $rlw^{(\text{-})}$ in a given retry is $rlw^{(\text{-})} = rc + cw + cc$, as we should pay at least the memory accesses and the critical work $cw$ in between.

**Thread-wise:** A given thread can at most perform one successful retry every $pw + rlw^{(\text{-})}$ units of time. In the best case, $P$ threads can then lead to a throughput of $P/(pw + rlw^{(\text{-})})$.

**System-wise:** By definition, two successful retries cannot overlap, hence we have at most 1 successful retry every $rlw^{(\text{-})}$ units of time.

Altogether, the throughput $T$ is bounded by

$$T \leq \min\left(\frac{1}{rc + cw + cc}, \frac{P}{pw + rc + cw + cc}\right), i.e.$$

$$T \leq \begin{cases} \frac{1}{rc+cw+cc} & \text{if } pw \leq (P-1)(rc + cw + cc) \\ \frac{P}{pw+rc+cw+cc} & \text{otherwise.} \end{cases} \tag{1}$$

*2) Conflicts:*

*Logical conflicts:* Equation 1 expresses the fact that when $pw$ is small enough, *i.e.* when $pw \leq (P-1)rlw^{(\text{-})}$, we cannot expect that every thread performs a successful retry every $pw + rlw^{(\text{-})}$ units of time, since it is more than what the retry loop can afford. As a result, some logical conflicts, hence unsuccessful retries, will be inevitable, while the others, if any, are called *wasted*.

However, different executions can lead to different numbers of failures, which end up with different throughput values. Figures 2 and 3 depict two executions, where the black parts are the calls to Initialization, the blue parts are the parallel sections, and the retries can be either unsuccessful — in red — or successful — in green. We experiment different initialization times, and observe different synchronizations, hence different numbers of wasted retries. After the initial transient state, the execution depicted in Figure 3 comprises only the inevitable unsuccessful retries, while the execution of Figure 2 contains one wasted retry.

We can see on those two examples that a cyclic execution is reached after the transient behavior; actually, we show in Section IV that, in the absence of hardware conflicts, every execution will become periodic, if the initialization times are spaced enough. In addition, we prove that the shortest period is such that, during this period, every thread succeeds exactly once. This finally leads us to define the additional failures as wasted, since we can directly link the throughput with this number of wasted retries: a higher number of wasted retries implying a lower throughput.
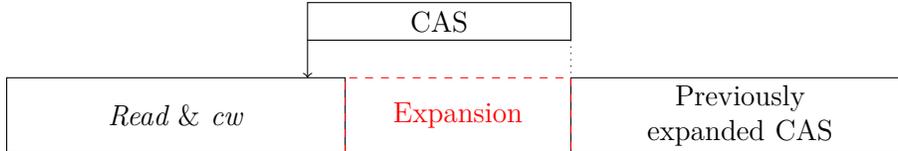
**Figure 4:** Expansion

*Hardware conflicts:* The requirement of atomicity compels the ownership of the data in an exclusive manner by the executing core. This fact prohibits concurrent execution of atomic instructions if they are operating on the same data. Therefore, overlapping parts of atomic instructions are serialized by the hardware, leading to stalls in subsequently issued ones. For our target lock-free algorithm, these stalls that we refer to as expansion become an important slowdown factor in case threads interfere in the retry loop. As illustrated in Figure 4, the latency for *CAS* can expand and cause remarkable decreases in throughput since the *CAS* of a successful thread is then expanded by others; for this reason, the amount of work inside a retry is not constant, but is, generally speaking, a function depending on the number of threads that are inside the retry loop.

*3) Process:* We deal with the two kinds of conflicts separately and connect them together through the fixed-point iterative convergence.

In Section V-A, we compute the expansion in execution time of a retry, noted $e$, by following a probabilistic approach. The estimation takes as input the expected number of threads inside the retry loop at any time, and returns the expected increase in the execution time of a retry due to the serialization of atomic primitives.

In Section IV, we are given program without hardware conflict described by the size of the parallel section $pw^{(+)}$ and the size of a retry $rlw^{(+)}$. We compute upper and lower bounds on the throughput $T$, the number of wasted retries $w$, and the average number of threads inside the retry loop $P_{rl}$. Without loss of generality, we can normalize those execution times by the execution time of a retry, and define the parallel section size as $pw^{(+)} = q + r$, where $q$ is a non-negative integer and $r$ is such that $0 \leq r < 1$. This pair (together with the number of threads $P$) constitutes the actual input of the estimation.

Finally, we combine those two outcomes in Section V-B by emulating expansion through work not prone to hardware conflicts and obtain the full estimation of the throughput.

## IV. Execution without hardware conflict

We show in this section that, in the absence of hardware conflicts, the execution becomes periodic, which eases the calculation of the throughput. We start by defining some useful concepts: $(f, P)$-cyclic executions are special kind of periodic executions such that within the shortest period, each thread performs exactly $f$ unsuccessful retries and 1 successful retry. The *well-formed seed* is a set of events that allows us to detect an $(f, P)$-*cyclic execution* early, and the *gaps* are a measure of the quality of the synchronization between threads. The idea is to iteratively add threads into the game and show that the periodicity is maintained. Theorem 1 establishes a fundamental relation between gaps and well-formed seeds, while Theorem 2 proves the periodicity, relying on the disjoint cases of Lemma 1, 3, and 4. Finally, we exhibit upper and lower bounds on throughput and number of failures, along with the average number of threads inside the retry loop.

*A. Setting*

*1) Initial Restrictions:*

**Remark 1.** Concerning correctness, we assume that the reference point of the *Read* and the *CAS* occurs when the thread enters and exits any retry, respectively.

**Remark 2.** We do not consider simultaneous events, so all inequalities that refer to time comparison are strict, and can be viewed as follows: time instants are real numbers, and can be equal, but every event is associated with a thread; also, in order to obtain a strict order relation, we break ties according to the thread numbers (for instance with the relation $<$).

*2) Notations and Definitions:* We recall that $P$ threads are executing the pseudo-code described in Procedure AbstractAlgorithm, one retry is of unit-size, and the parallel section is of size $pw^{(+)} = q + r$, where $q$ is a non-negative integer and $r$ is such that $0 \leq r < 1$. Considering a thread $\mathcal{T}_n$ which succeeds at time $S_n$; this thread completes a whole retry in 1 unit of time, then executes the parallel section of size $pw^{(+)}$, and attempts to perform again the operation every unit of time, until one of the attempt is successful.

**Definition 1.** *An execution with $P$ threads is called $(C, P)$-cyclic execution if and only if (i) the execution is periodic, i.e. at every time, every thread is in the same state as one period before, (ii) the shortest period contains exactly one successful attempt per thread, (iii) the shortest period is $1 + q + r + C$.*

**Definition 2.** *Let $\mathcal{S} = (\mathcal{T}_i, S_i)_{i \in [\![0, P-1]\!]}$, where $\mathcal{T}_i$ are threads and $S_i$ ordered times, i.e. such that $S_0 < S_1 < \cdots < S_{P-1}$. $\mathcal{S}$ is a seed if and only if for all $i \in [\![0, P-1]\!]$, $\mathcal{T}_i$ does not succeed between $S_0$ and $S_i$, and starts a retry at $S_i$.*

*We define $f(\mathcal{S})$ as the smallest non-negative integer such that $S_0 + 1 + q + r + f(\mathcal{S}) > S_{P-1} + 1$, i.e. $f(\mathcal{S}) = \max(0, \lceil S_{P-1} - S_0 - q - r \rceil)$. When $\mathcal{S}$ is clear from the context, we denote $f(\mathcal{S})$ by $f$.*

**Definition 3.** *$\mathcal{S}$ is a well-formed seed if and only if for each $i \in [\![0, P-1]\!]$, the execution of thread $\mathcal{T}_i$ contains the following sequence: firstly a success beginning at $S_i$, the parallel section, $f$ unsuccessful retries, and finally a successful retry.*

Those definitions are coupled through the two natural following properties:

**Property 1.** *Given a $(C, P)$-cyclic execution, any seed $\mathcal{S}$ including $P$ consecutive successes is a well-formed seed, with $f(\mathcal{S}) = C$.*

*Proof:* Choosing any set of $P$ consecutive successes, we are ensured, by the definition of a $(f, P)$-cyclic execution, that for each thread, after the first success, the next success will be obtained after $f$ failures. The order will be preserved, and this shows that a seed including our set of successes is actually a well-formed seed. ∎

**Property 2.** *If there exists a well-formed seed in an execution, then after each thread succeeded once, the execution coincides with an $(f, P)$-cyclic execution.*

*Proof:* By the definition of a well-formed seed, we know that the threads will first succeed in order, fails $f$ times, and succeed again in the same order. Considering the second set of successes in a new well-formed seed, we observe that the threads will succeed a third time in the same order, after failing $f$ times. By induction, the execution coincides with an $(f, P)$-cyclic execution. ∎

Together with the seed concept, we define the notion of *gap* that we will use extensively in the next subsection. The general idea of those gaps is that within an $(f, P)$-cyclic execution, the period is higher than $P \times 1$, which is the total execution time of all the successful retries within the period. The difference between the period (that lasts $1 + q + r + f$) and $P$, reduced by $r$ (so that we obtain an integer), is referred as *lagging time* in the following. If the threads are numbered according to their order of success (modulo $P$), as the time elapsed between the successes of two given consecutive
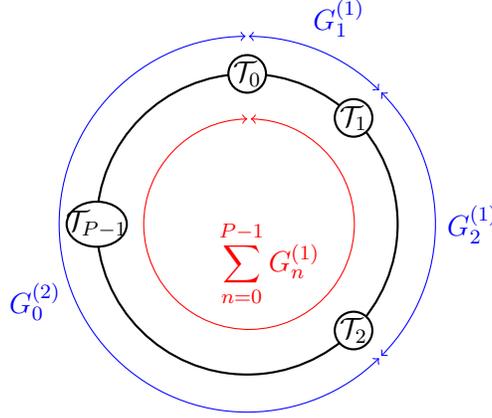
**Figure 5:** Gaps

threads is constant (during the next period, this time will remain the same), this lagging time can be seen in a circular manner (see Figure 5): the threads are represented on a circle whose length is the lagging time increased by $r$, and the length between two consecutive threads is the time between the end of the successful retry of the first thread and the begin of the successful retry of the second one. More formally, for all $(n,k) \in [\![0, P-1]\!]^2$, we define the gap $G_n^{(k)}$ between $\mathcal{T}_n$ and its $k^{\text{th}}$ predecessor based on the gap with the first predecessor:

$$\left\{ \begin{array}{l} \forall n \in [\![1, P-1]\!] \quad ; \quad G_n^{(1)} = S_n - S_{n-1} - 1 \\ G_0^{(1)} = S_0 + q + r + f - S_{P-1} \end{array} \right. ,$$

which leads to the definition of higher order gaps:

$$\forall n \in [\![0, P-1]\!] \quad ; \quad \forall k > 0 \quad ; \quad G_n^{(k)} = \sum_{j=n-k+1}^{n} G_{j \bmod P}^{(1)}.$$

For consistency, for all $n \in [\![0, P-1]\!]$, $G_n^{(0)} = 0$.

Equally, the gaps can be obtained directly from the successes: for all $k \in [\![1, P-1]\!]$,

$$G_n^{(k)} = \left\{ \begin{array}{ll} S_n - S_{n-k} - k & \text{if } n > k \\ S_n - S_{P+n-k} + 1 + q + r + f - k & \text{otherwise} \end{array} \right. \tag{2}$$

Note that, in an $(f, P)$-cyclic execution, the lagging time is the sum of all first order gaps, reduced by $r$.

Now we extend the concept of well-formed seed to weakly-formed seed.

**Definition 4.** *Let $\mathcal{S} = (\mathcal{T}_i, S_i)_{i \in [\![0, P-1]\!]}$ be a seed.*

*$\mathcal{S}$ is a weakly-formed seed for $P$ threads if and only if: $(\mathcal{T}_i, S_i)_{i \in [\![0, P-2]\!]}$ is a well-formed seed for $P-1$ threads, and the first thread succeeding after $\mathcal{T}_{P-2}$ is $\mathcal{T}_{P-1}$.*

**Property 3.** *Let $\mathcal{S} = (\mathcal{T}_i, S_i)_{i \in [\![0, P-1]\!]}$ be a weakly-formed seed.*

*Denoting $f = f\left((\mathcal{T}_i, S_i)_{i \in [\![0, P-2]\!]}\right)$, for each $n \in [\![0, P-1]\!]$, $G_n^{(f)} < 1$.*

*Proof:* We have $S_{P-2} + 1 < S_{P-1} < R_0^f$, and if we note indeed $\widetilde{G}_n^{(k)}$ the gaps within $(\mathcal{T}_i, S_i)_{i \in [\![0, P-2]\!]}$, the previous well-formed seed with $P-1$ threads, we know that for all $n \in [\![1, P-2]\!]$, $\widetilde{G}_n^{(1)} = G_n^{(1)}$, and $G_{P-1}^{(1)} + G_0^{(1)} = \widetilde{G}_0^{(1)}$, which leads to $G_n^{(k)} \leq \widetilde{G}_n^{(k)}$, for all $n \in [\![0, P-1]\!]$ and $k$; hence the weaker property. ∎
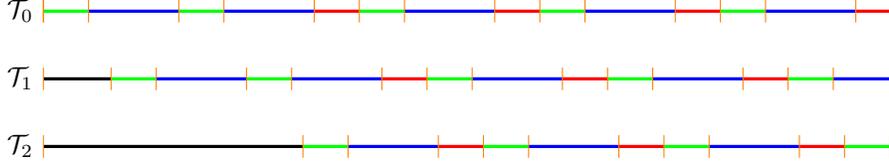
**Figure 6:** Lemma 1 configuration

**Lemma 1.** *Let $\mathcal{S}$ be a weakly-formed seed, and $f = f\left((\mathcal{T}_i, S_i)_{i \in [\![0, P-2]\!]}\right)$. If, for all $n \in [\![0, P-1]\!]$, $G_n^{(f+1)} < 1$, then there exists later in the execution a well-formed seed $\mathcal{S}'$ for $P$ threads such that $f(\mathcal{S}') = f + 1$.*

*Proof:* The proof is straightforward; $\mathcal{S}$ is actually a well-formed seed such that $f(\mathcal{S}) = f + 1$. Since $R_0^f - S_{P-1} < G_0^{(1)} < 1$, the first success of $\mathcal{T}_0$ after the success of $\mathcal{T}_{P-1}$ is its $f + 1^{\text{th}}$ retry. ∎

*B. Cyclic Executions*

**Theorem 1.** *Given a seed $\mathcal{S} = (\mathcal{T}_i, S_i)_{i \in [\![0, P-1]\!]}$, $\mathcal{S}$ is a well-formed seed if and only if for all $n \in [\![0, P-1]\!]$, $0 \le G_n^{(f)} < 1$.*

*Proof:*

Let $\mathcal{S} = (\mathcal{T}_i, S_i)_{i \in [\![0, P-1]\!]}$ be a seed.

($\Leftarrow$) We assume that for all $n \in [\![0, P-1]\!]$, $0 < G_n^{(f)} < 1$, and we first show that the first successes occur in the following order: $\mathcal{T}_0$ at $S_0$, $\mathcal{T}_1$ at $S_1$, ..., $\mathcal{T}_{P-1}$ at $S_{P-1}$, $\mathcal{T}_0$ again at $R_0^f$. The first threads that are successful executes their parallel section after their success, then enters their second retry loop: from this moment, they can make the first attempt of the threads, that has not been successful yet, fail. Therefore, we will look at which retry of which already successful threads could have an impact on which other threads.

We can notice that for all $n \in [\![0, P-1]\!]$, if the first success of $\mathcal{T}_n$ occurs at $S_n$, then its next attempts will potentially occur at $R_n^k = S_n + 1 + q + r + k$, where $k \ge 0$. More specifically, thanks to Equation 2, for all $n \le f$, $R_n^k = S_{P+n-f} + G_n^{(f)} + k$. Also, for all $k \le f - n$,

$$
\begin{aligned}
R_n^k - S_{P+n-f+k} &= -(S_{P+n-f+k} - S_{P+n-f} - k) + G_n^{(f)} \\
&= G_n^{(f)} - G_{P+n-f+k}^{(k)} \\
R_n^k - S_{P+n-f+k} &= G_n^{(f-k)},
\end{aligned}
\tag{3}
$$

and this implies that if $k > 0$,

$$
S_{P+n-f+k} - R_n^{k-1} = 1 - G_n^{(f-k)}.
\tag{4}
$$

We know, by hypothesis, that $0 < G_n^{(f-k)} < 1$, equivalently $0 < 1 - G_n^{(f-k)} < 1$. Therefore Equation 3 states that if a thread $\mathcal{T}_{n'}$ starts a successful attempt at $S_{P+n-f+k}$, then this thread will make the $k^{\text{th}}$ retry of $\mathcal{T}_n$ fail, since $\mathcal{T}_n$ enters a retry while $\mathcal{T}_{n'}$ is in a successful retry. And Equation 4 shows that, given a thread $\mathcal{T}_{n'}$ starting a new retry at $S_{P+n-f+k}$, the only retry of $\mathcal{T}_n$ that can make $\mathcal{T}_{n'}$ fail on its attempt is the $(k-1)^{\text{th}}$ one. There is indeed only one retry of $\mathcal{T}_n$ that can enter a retry before the entrance of $\mathcal{T}_{n'}$, and exit the retry after it.

$\mathcal{T}_0$ is the first thread to succeed at $S_0$, because no other thread is in the retry loop at this time. Its next attempt will occur at $R_0^0$, and all thread attempts that start before $S_{P-f}$ (included) cannot fail because of $\mathcal{T}_0$, since it runs then the parallel section. Also, since all gaps are positive, the threads $\mathcal{T}_1$ to $\mathcal{T}_{P-f}$ will succeed in this order, respectively starting at times $S_1$ to $S_{P-f}$.

Then, using induction, we can show that $\mathcal{T}_{P-f+1}$, ..., $\mathcal{T}_{P-1}$ succeed in this order, respectively starting at times $S_{P-f+1}$, ..., $S_{P-1}$. For $j \in [\![0, f-1]\!]$, let $(\mathcal{P}_j)$ be the following property: for all $n \in [\![0, P-f+j]\!]$, $\mathcal{T}_n$ starts a successful retry at $S_n$. We assume that for a given $j$, $(\mathcal{P}_j)$ is true, and we show that it implies that $\mathcal{T}_{P-f+j+1}$ will succeed at $S_{P-f+j+1}$. The successful attempt of $\mathcal{T}_{P-f+j}$ at $S_{P-f+j}$ leads, for all $j' \in [\![0, j]\!]$, to the failure of the $j'^{\text{th}}$ retry of $\mathcal{T}_{j-j'}$ (explanation of Equation 3). But for each $\mathcal{T}_{j'}$, this attempt was precisely the one that could have made $\mathcal{T}_{P-f+j+1}$ fail on its attempt at $S_{P-f+j+1}$ (explanation of Equation 3). Given that all threads $\mathcal{T}_n$, where $n > P-f+j+1$, do not start any retry loop before $S_{P-f+j+1}$, $\mathcal{T}_{P-f+j+1}$ will succeed at $S_{P-f+j+1}$. By induction, $(\mathcal{P}_j)$ is true for all $j \in [\![0, f-1]\!]$.

Finally, when $\mathcal{T}_{P-1}$ succeeds, it makes the $(f-1-n)^{\text{th}}$ retry of $\mathcal{T}_n$ fail, for all $n \in [\![0, f-1]\!]$; also the next potentially successful attempt for $\mathcal{T}_n$ is at $R_n^{f-n}$. (Naturally, for all $n \in [\![f, P-1]\!]$, the next potentially successful attempt for $\mathcal{T}_n$ is at $R_n^0$.)

We can observe that for all $n < P$, $j \in [\![0, P-1-n]\!]$, and all $k \geq j$,

$$R_{n+j}^{k-j} - R_n^k = S_{n+j} + k - j - (S_n + k)$$
$$R_{n+j}^{k-j} - R_n^k = G_{n+j}^{(j)}, \tag{5}$$

hence for all $n \in [\![1, f]\!]$, $R_n^{f-n} - R_0^f = G_n^{(n)} > 0$.

$$R_n^{f-n} - R_0^f = G_n^{(n)} > 0.$$

As we have as well, for all $n \in [\![f+1, P-1]\!]$, $R_n^0 > R_f^0$, we obtain that among all the threads, the earliest possibly successful attempt is $R_0^f$. Following $\mathcal{T}_{P-1}$, $\mathcal{T}_0$ is consequently the next successful thread in its $f^{\text{th}}$ retry.

To conclude this part, we can renumber the threads ($\mathcal{T}_{n+1}$ becoming now $\mathcal{T}_n$ if $n > 0$, and $\mathcal{T}_0$ becoming $\mathcal{T}_{P-1}$), and follow the same line of reasoning. The only difference is the fact that $\mathcal{T}_{P-1}$ (according to the new numbering) enters the retry loop $f$ units of time before $S_{P-1}$, but it does not interfere with the other threads, since we know that those attempts will fail.

There remains the case where there exists $n \in [\![0, P-1]\!]$ such that $G_n^{(f)} = 0$. This implies that $f = 0$, thus we have a well-formed seed.

($\Rightarrow$) We prove now the implication by contraposition; we assume that there exists $n \in [\![0, P-1]\!]$ such that $G_n^{(f)} > 1$ or $G_n^{(f)} < 0$, and show that $\mathcal{S}$ is not a well-formed seed.

We assume first that an $f^{\text{th}}$ order gap is negative. As it is a sum of $1^{\text{st}}$ order gaps, then there exists $n'$ such that $G_{n'}^{(1)}$ is negative; let $n''$ be the highest one.

If $n'' > 0$, then either the threads $\mathcal{T}_0, \ldots, \mathcal{T}_{n''-1}$ succeeded in order at their $0^{\text{th}}$ retry, and then $\mathcal{T}_{n''-1}$ makes $\mathcal{T}_{n''}$ fail at its $0^{\text{th}}$ retry (we have a seed, hence by definition, $S_{n''-1} < S_{n''}$, and $G_{n''}^{(1)} < 0$, thus $S_{n''-1} < S_{n''} < S_{n''-1} + 1$ ), or they did not succeed in order at their first try. In both cases, $\mathcal{S}$ is not a well-formed seed.

If $n'' = 0$, let us assume that $\mathcal{S}$ is a well-formed seed. Let also a new seed be $\mathcal{S}' = (\mathcal{T}_i, S_i')_{i \in [\![0, P-1]\!]}$, where for all $n \in [\![0, P-2]\!]$, $S_{n+1}' = S_n$, and $S_0' = S_{P-1} - (q+1+f+r)$. Like $\mathcal{S}$, $\mathcal{S}'$ is a well-formed seed; however, $G_1^{(1)}$ is negative, and we fall back into the previous case, which shows that $\mathcal{S}'$ is not a well-formed seed. This is absurd, hence $\mathcal{S}$ is not a well-formed seed.

We assume now that every gap is positive and choose $n_0$ defined by: $n_0 = \min\{n \; ; \; \exists k \in [\![0, P-1]\!] / G_{n+k}^{(k)} > 1\}$, and $f_0 = \min\{k \; ; \; G_{n_0+k}^{(k)} > 1\}$: among the gaps that exceed 1, we pick those that concern the earliest thread, and among them the one with the lowest order.

Let us assume that threads $\mathcal{T}_0, \ldots, \mathcal{T}_{P-1}$ succeed at their $0^{\text{th}}$ retry in this order, then $\mathcal{T}_0, \ldots, \mathcal{T}_{n_0}$ complete their second successful retry loop at their $f^{\text{th}}$ retry, in this order. If this is not the case, then $\mathcal{S}$ is not a well-formed seed, and the proof is completed. According to Equation 5, we have, on the one

hand, $R_{n_0+1}^{f_0-1} - R_{n_0}^{f_0} = G_{n_0+1}^{(1)}$, which implies $R_{n_0+1}^{f_0} - 1 - R_{n_0}^{f_0} = G_{n_0+1}^{(1)}$, thus $R_{n_0+1}^{f} - (R_{n_0}^{f}+1) = G_{n_0+1}^{(1)}$; and on the other hand, $R_{n_0+f_0}^{0} - R_{n_0}^{f_0} = G_{n_0+f_0}^{(f_0)}$ implying $R_{n_0+f_0}^{f-f_0} - (R_{n_0}^{f}+1) = G_{n_0+f_0}^{(f_0)} - 1$. As we know that $G_{n_0+f_0}^{(f_0)} - G_{n_0+1}^{(1)} = G_{n_0+f_0}^{(f_0-1)} < 1$ by definition of $f_0$ (and $n_0$), we can derive that $R_{n_0+1}^{f} - (R_{n_0}^{f}+1) > R_{n_0+f_0}^{f-f_0} - (R_{n_0}^{f}+1)$. We have assumed that $\mathcal{T}_{n_0}$ succeeds at its $f^{\text{th}}$ retry, which will end at $R_{n_0}^{f} + 1$. The previous inequality states then that $\mathcal{T}_{n_0+1}$ cannot be successful at its $f^{\text{th}}$ retry, since either a thread succeeds before $\mathcal{T}_{n_0+f_0}$ and makes both $\mathcal{T}_{n_0+f_0}$ and $\mathcal{T}_{n_0+1}$ fail, or $\mathcal{T}_{n_0+f_0}$ succeeds and makes $\mathcal{T}_{n_0+1}$ fail. We have shown that $\mathcal{S}$ is not a well-formed seed. ∎

**Lemma 2.** *Assuming $r \neq 0$, if a new thread is added to an $(f, P)$-cyclic execution, it will eventually succeed.*

*Proof:*

Let $R_P^0$ be the time of the $0^{\text{th}}$ retry of the new thread, that we number $\mathcal{T}_P$. If this retry is successful, we are done; let us assume now that this retry is a failure, and let us shift the thread numbers (for the threads $\mathcal{T}_0, \ldots, \mathcal{T}_{P-1}$) so that $\mathcal{T}_0$ makes $\mathcal{T}_P$ fail on its first attempt. We distinguish two cases, depending on whether $G_0^{(P)} > R_P^0 - S_0$ or not.

We assume that $G_0^{(P)} > R_P^0 - S_0$. We know that $n \mapsto G_n^{(n)}$ is increasing on $[\![0, P-1]\!]$ and that $G_0^{(0)} = 0$, hence let $n_0 = \min\{n \in [\![0, P-1]\!] \; ; \; G_n^{(n)} < R_P^0 - S_0\}$. For all $k \in [\![0, n_0]\!]$, we have $R_P^k - S_k = k + R_P^0 - (G_k^{(k)} + S_0 + k) = R_P^0 - S_0 - G_k^{(k)}$ hence $R_P^k - S_k > 0$ and $R_P^k - S_k < R_P^0 - S_0 < 1$. This shows that $\mathcal{T}_0, \ldots, \mathcal{T}_{n_0}$, because of their successes at $S_0, \ldots, S_{n_0}$, successively make $0^{\text{th}}, \ldots, n_0^{\text{th}}$ retries (respectively) of $\mathcal{T}_P$ fail. The next attempt for $\mathcal{T}_P$ is at $R_P^{n_0+1}$, which fulfills the following inequality: $R_P^{n_0+1} - (S_{n_0} + 1) < S_{n_0+1} - (S_{n_0} + 1)$ since

$$R_P^{n_0+1} - S_{n_0+1} = (n_0 + 1 + R_P^0) - (G_{n_0+1}^{(n_0+1)} + S_0 + n_0 + 1)$$
$$R_P^{n_0+1} - S_{n_0+1} > 0.$$

$\mathcal{T}_{n_0+1}$ should have been the successful thread, but $\mathcal{T}_P$ starts a retry before $S_{n_0+1}$, and is therefore succeeding.

We consider now the reverse case by assuming that $G_0^{(P)} < R_P^0 - S_0$. With the previous line of reasoning, we can show that $\mathcal{T}_0, \ldots, \mathcal{T}_{P-1}$, because of their successes at $S_0, \ldots, S_{P-1}$, successively make $0^{\text{th}}, \ldots, (P-1)^{\text{th}}$ retries (respectively) of $\mathcal{T}_P$ fail. Then we are back in the same situation when $\mathcal{T}_0$ made $\mathcal{T}_P$ fail for the first time ($\mathcal{T}_0$ makes $\mathcal{T}_P$ fail), except that the success of $\mathcal{T}_0$ starts at $S_0' = S_0 + G_0^{(P)}$. As $G_0^{(P)} = q + r + f - P > 0$ and $q$, f and $P$ are integers, we have that $G_0^{(P)} \geq r$. By the way, if we had $G_0^{(P)} > r$, we would have $G_0^{(P)} \geq 1 + r > R_P^0 - S_0$, which is absurd. $S_0$ makes indeed $R_P^0$ fail, therefore $G_0^{(P)}$ should be less than 1. Consequently, we are ensured that $G_0^{(P)} = r$. We define

$$k_0 = \left\lfloor \frac{R_P^0 - S_0}{r} \right\rfloor ;$$

also, for every $k \in [\![1, k_0]\!]$, $r < R_P^0 - (S_0 + k \times r)$ and $r > R_P^0 - (S_0 + (k_0 + 1) \times r)$: the cycle of successes of $\mathcal{T}_0, \ldots, \mathcal{T}_{P-1}$ is executed $k_0$ times. Then the situation is similar to the first case, and $\mathcal{T}_P$ will succeed.

∎

**Lemma 3.** *Let $\mathcal{S}$ be a weakly-formed seed, and $f = f\left((\mathcal{T}_i, S_i)_{i \in [\![0, P-2]\!]}\right)$. If $G_f^{(f+1)} > 1$, and if the second success of $\mathcal{T}_{P-1}$ does not occur before the second success of $\mathcal{T}_{f-1}$, then we can find in the execution a well-formed seed $\mathcal{S}'$ for $P$ threads such that $f(\mathcal{S}') = f$.*
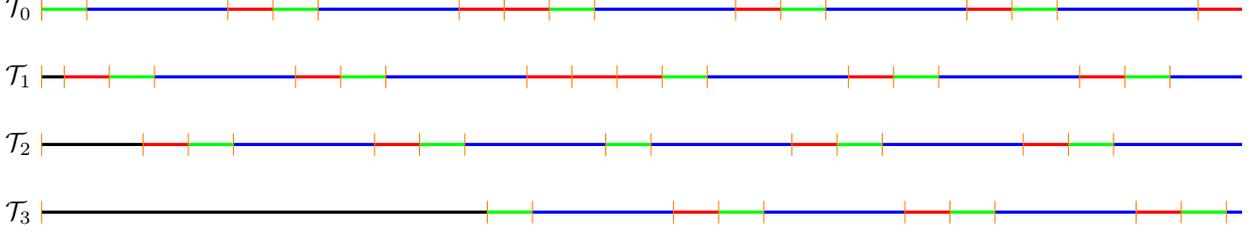
*Proof:*

**Figure 7:** Lemma 3 configuration

Let us first remark that, by the definition of a weakly-formed seed, all threads will succeed once, in order. Then two ordered groups of threads will compete for each of the next successes, until $\mathcal{T}_{f-1}$ succeeds for the second time.

Let $e$ be the smallest integer of $[\![f, P-1]\!]$ such that the second success of $\mathcal{T}_e$ occurs after the second success of $\mathcal{T}_{f-1}$. Let then $\mathcal{S}_1$ and $\mathcal{S}_2$ be the two groups of threads that are in competition, defined by

$$\mathcal{S}_1 = \{\mathcal{T}_n \; ; \; n \in [\![0, f-1]\!]\}$$
$$\mathcal{S}_2 = \{\mathcal{T}_n \; ; \; n \in [\![f, e-1]\!]\}$$

For all $n \in [\![0, e-1]\!]$, we note

$$rank(n) = \begin{cases} G_n^{(n+1)} & \text{if } \mathcal{T}_n \in \mathcal{S}_1 \\ G_n^{(n+1)} - 1 & \text{if } \mathcal{T}_n \in \mathcal{S}_2 \end{cases}.$$

We define $\sigma$, a permutation of $[\![0, e-1]\!]$ that describes the reordering of the threads during the round of the second successes, such that, for all $(i, j) \in [\![0, e-1]\!]^2$, $\sigma(i) < \sigma(j)$ if and only if $rank(i) < rank(j)$.

We also define a function that will help in expressing the $\sigma^{-1}(k)$'s:

$$m_2 : \quad [\![0, e-1]\!] \quad \longrightarrow \quad [\![f, e-1]\!]$$
$$k \quad \longmapsto \quad \max\{\ell \in [\![f, e-1]\!] \; ; \; \mathcal{T}_\ell \in \mathcal{S}_2 \; ; \; \sigma(\ell) \leq k\}.$$

We note that $rank|_{[\![0, f-1]\!]}$ is increasing, as well as $rank|_{[\![f, e-1]\!]}$. This shows that $\#\{\mathcal{T}_\ell \in \mathcal{S}_2 \; ; \; \sigma(\ell) \leq k\} = m_2(k) - (f-1)$. Consequently, if $\mathcal{T}_{\sigma^{-1}(k)} \in \mathcal{S}_2$, then

$$m_2(k) = \#\{\mathcal{T}_\ell \in \mathcal{S}_2 \; ; \; \sigma(\ell) \leq k\} + f - 1$$
$$= \#\{\mathcal{T}_\ell \in \mathcal{S}_2 \; ; \; \ell \leq \sigma^{-1}(k)\} + f - 1$$
$$= \sigma^{-1}(k) - f + 1 + f - 1$$
$$m_2(k) = \sigma^{-1}(k).$$

Conversely, if $\mathcal{T}_{\sigma^{-1}(k)} \in \mathcal{S}_1$, among $\{\mathcal{T}_{\sigma(n)} \; ; \; n \in [\![0, k]\!]\}$, there are exactly $m_2(k) - f + 1$ threads in $\mathcal{S}_2$, hence

$$\sigma^{-1}(k) = k + 1 - (m_2(k) - f + 1) - 1 = f + k - m_2(k) - 1.$$

In both cases, among $\{\mathcal{T}_{\sigma(n)} \; ; \; n \in [\![0, k]\!]\}$, there are exactly $m_2(k) - f + 1$ threads in $\mathcal{S}_2$, and $m_1(k) = k - (m_2(k) - f)$ threads in $\mathcal{S}_1$.

We prove by induction that after this first round, the next successes will be, respectively, achieved by $\mathcal{T}_{\sigma^{-1}(0)}, \mathcal{T}_{\sigma^{-1}(1)}, \ldots, \mathcal{T}_{\sigma^{-1}(e-1)}$. In the following, by "$k^{\text{th}}$ success", we mean $k^{\text{th}}$ success after the first success of $\mathcal{T}_{P-1}$, starting from 0, and the $R_i^j$'s denote the attempts of the second round.

Let $(\mathcal{P}_K)$ be the following property: for all $k \leq K$, the $k^{\text{th}}$ success is achieved by $\mathcal{T}_{\sigma^{-1}(k)}$ at $R_{\sigma^{-1}(k)}^{f+k-\sigma^{-1}(k)}$. We assume $(\mathcal{P}_K)$ true, and we show that the $(K+1)^{\text{th}}$ success is achieved by $\mathcal{T}_{\sigma^{-1}(K+1)}$ at $R_{\sigma^{-1}(K+1)}^{f+K+1-\sigma^{-1}(K+1)}$.

We first show that if $\mathcal{T}_{\sigma^{-1}(K)} \in \mathcal{S}_1$, then

$$R_{m_2(K)+1}^{m_1(K)-1} > R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} > R_{m_2(K)}^{m_1(K)}. \tag{6}$$

On the one hand,

$$
\begin{aligned}
R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} &= K - \sigma^{-1}(K) + R_{\sigma^{-1}(K)}^{f} \\
&= K - \sigma^{-1}(K) + R_0^{f} + \sigma^{-1}(K) + G_{\sigma^{-1}(K)}^{(\sigma^{-1}(K))} \\
&= K + S_{P-1} + 1 + G_0^{(1)} + G_{\sigma^{-1}(K)}^{(\sigma^{-1}(K))} \\
R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} &= K + S_{P-1} + 1 + G_{\sigma^{-1}(K)}^{(\sigma^{-1}(K)+1)}.
\end{aligned}
$$

On the other hand,

$$
\begin{aligned}
R_{m_2(K)}^{f+K-m_2(K)} &= (m_2(K) - f) + R_f^{K-(m_2(K)-f)} + G_{m_2(K)}^{(m_2(K)-f)} \\
&= (m_2(K) - f) + K - (m_2(K) - f) + R_f^{0} + G_{m_2(K)}^{(m_2(K)-f)} \\
&= (m_2(K) - f) + K - (m_2(K) - f) + S_{P-1} + 1 + (G_f^{(f+1)} - 1) + G_{m_2(K)}^{(m_2(K)-f)} \\
R_{m_2(K)}^{f+K-m_2(K)} &= K + S_{P-1} + 1 + G_{m_2(K)}^{(m_2(K)+1)} - 1.
\end{aligned}
$$

Therefore,

$$
\begin{aligned}
R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} - R_{m_2(K)}^{m_1(K)} &= R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} - R_{m_2(K)}^{f+K-m_2(K)} \\
&= G_{\sigma^{-1}(K)}^{(\sigma^{-1}(K)+1)} - \left( G_{m_2(K)}^{(m_2(K)+1)} - 1 \right) \\
R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} - R_{m_2(K)}^{m_1(K)} &= rank\left( \sigma^{-1}(K) \right) - rank\left( m_2(K) \right).
\end{aligned}
$$

In a similar way, we can obtain that if $\mathcal{T}_{\sigma^{-1}(K)} \in \mathcal{S}_2$, then

$$R_{m_1(K)}^{m_2(K)} > R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} > R_{m_1(K)-1}^{m_2(K)+1}. \tag{7}$$

In addition, we recall that if $\mathcal{T}_{\sigma^{-1}(K)} \in \mathcal{S}_2$, $\sigma^{-1}(K) = m_2(K)$, thus the second inequality of Equation 6 becomes an equality, and if $\mathcal{T}_{\sigma^{-1}(K)} \in \mathcal{S}_1$, $\sigma^{-1}(K) = f + K - m_2(K) - 1$, hence the second inequality of Equation 7 becomes an equality.

Now let us look at which attempt of other threads $\mathcal{T}_{\sigma^{-1}(K)}$ made fail. From now on, and until explicitly said otherwise, we assume that $\mathcal{T}_{\sigma^{-1}(K)} \in \mathcal{S}_1$. According to Equation 6, we have

$$
\begin{array}{ccc}
R_{m_2(K)+1}^{m_1(K)-1} > & R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} > & R_{m_2(K)}^{m_1(K)} \\
R_{m_2(K)+j}^{m_1(K)-j} - R_{m_2(K)+1}^{m_1(K)-1} < & R_{m_2(K)+j}^{m_1(K)-j} - R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} < & R_{m_2(K)+j}^{m_1(K)-j} - R_{m_2(K)}^{m_1(K)} \\
G_{m_2(K)+j}^{(j-1)} < & R_{m_2(K)+j}^{m_1(K)-j} - R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} < & G_{m_2(K)+j}^{(j)}
\end{array}
$$

This holds for every $j \in [\![1, m_1(K)]\!]$, implying $j \leq f$, since there could not be more than $f$ threads in $\mathcal{S}_1$. Therefore, as by assumptions gaps of at most $f^{\text{th}}$ order are between 0 and 1,

$$0 < R_{m_2(K)+j}^{m_1(K)-j} - R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} < 1;$$

showing that the success of $\mathcal{T}_{\sigma^{-1}(K)}$ makes thread $\mathcal{T}_{m_2(K)+j}$ fail on its attempt at $R_{m_2(K)+j}^{m_1(K)-j}$, for all $j \in [\![1, m_1(K)]\!]$.

Since $\mathcal{T}_{\sigma^{-1}(K)} \in \mathcal{S}_1$, $\sigma^{-1}(K) = m_1(K) - 1$. Also, for all $j \in [\![0, f-1-m_1(K)]\!]$,

$$
\begin{aligned}
R_{m_1(K)+j}^{m_2(K)-j} - R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} &= R_{m_1(K)+j}^{m_2(K)-j} - R_{m_1(K)-1}^{m_2(K)+1} \\
&= \left( R_{m_1(K)-1}^{m_2(K)-j} + (j+1) + G_{m_1(K)+j}^{(j+1)} \right) - \left( R_{m_1(K)-1}^{m_2(K)-j} + (j+1) \right)
\end{aligned}
$$
$$
R_{m_1(K)+j}^{m_2(K)-j} - R_{\sigma^{-1}(K)}^{f+K-\sigma^{-1}(K)} = G_{m_1(K)+j}^{(j+1)}
$$

As a result, $\mathcal{T}_{\sigma^{-1}(K)}$ makes $\mathcal{T}_{m_1(K)+j}$ fail on its attempt at $R_{m_1(K)+j}^{m_2(K)-j}$, for all $j \in [\![0, f-1-m_1(K)]\!]$, and the next attempt will occur at $R_{m_1(K)+j}^{m_2(K)-j+1}$.

Altogether, the next attempt after the end of the success of $\mathcal{T}_{\sigma^{-1}(K)}$ for $\mathcal{T}_{m_1(K)+j}$ is $R_{m_1(K)+j}^{m_2(K)-j+1}$, for $j \in [\![0, f-1-m_1(K)]\!]$, and for $\mathcal{T}_{m_2(K)+j}$ is $R_{m_2(K)+j}^{m_1(K)-j+1}$, for all $j \in [\![1, m_1(K)]\!]$.

Additionally, a thread will begin a new retry loop, the $0^{\text{th}}$ retry being at $R_{m_2(K)+m_1(K)+1}^0 = R_{f+K+1}^0$. We note that $f + K + 1$ could be higher than $P-1$, referring to a thread whose number is more than $P-1$. Actually, if $n > P-1$, $R_n^j$ refers to the $j^{\text{th}}$ retry of $\mathcal{T}_{rank(n-P+1)}$, after its first two successes.

The two heads, *i.e.* the two smallest indices, of $\mathcal{S}_1 \cap \sigma^{-1}([\![K+1, e-1]\!])$ and $\mathcal{S}_2 \cap \sigma^{-1}([\![K+1, e-1]\!])$ will then compete for being successful. Indeed, within $\mathcal{S}_1$, for $j \in [\![0, f-1-m_1(K)]\!]$,

$$
R_{m_1(K)+j}^{m_2(K)-j+1} - R_{m_1(K)}^{m_2(K)+1} = G_{m_1(K)+j}^{(j)} > 0,
$$

thus if someone succeeds in $\mathcal{S}_1$, it will be $\mathcal{T}_{m_1(K)}$. In the same way, for all $j \in [\![1, m_1(K)+1]\!]$,

$$
R_{m_2(K)+j}^{m_1(K)-j+1} - R_{m_2(K)+1}^{m_1(K)} = G_{m_2(K)+j}^{(j-1)} > 0,
$$

meaning that if someone succeeds in $\mathcal{S}_2$, it will be $\mathcal{T}_{m_2(K)+1}$.

Let us compare now those two candidates:

$$
\begin{aligned}
R_{m_1(K)}^{m_2(K)+1} - R_{m_2(K)+1}^{m_1(K)} &= m_2(K) + 1 - f + S_{P-1} + m_1(K) + G_{m_1(K)}^{(m_1(K)+1)} \\
&\quad - \left( m_1(K) + R_f^0 + m_2(K) + 1 - f + G_{m_2(K)+1}^{(m_2(K)+1-f)} \right) \\
&= S_{P-1} - 1 + G_{m_1(K)}^{(m_1(K)+1)} \\
&\quad - \left( S_{P-1} + G_f^{(f+1)} - 1 + G_{m_2(K)+1}^{(m_2(K)+1-f)} \right) \\
&= G_{m_1(K)}^{(m_1(K)+1)} - \left( G_{m_2(K)+1}^{(m_2(K)+2)} - 1 \right) \\
R_{m_1(K)}^{m_2(K)+1} - R_{m_2(K)+1}^{m_1(K)} &= rank(m_1(K)) - rank(m_2(K)+1).
\end{aligned}
$$

By definition, $\sigma^{-1}(K+1)$ is either $m_1(K)$ or $m_2(K) + 1$ and corresponds to the next successful thread. We can follow the same line of reasoning in the case where $\mathcal{T}_{\sigma^{-1}(K)} \in \mathcal{S}_2$ and prove in this way that $(\mathcal{P}_{K+1})$ is true.

$(\mathcal{P}_0)$ is true, and the property spreads until $(\mathcal{P}_{e-1})$, where all threads of $\mathcal{S}_1$ and $\mathcal{S}_2$ have been successful, in the order ruled by $\sigma^{-1}$, *i.e.* $\mathcal{T}_{\sigma^{-1}(0)}$, ..., $\mathcal{T}_{\sigma^{-1}(e-1)}$. And before those successes the threads $\mathcal{T}_{e-1} = \mathcal{T}_{\sigma^{-1}(e-1)}$, ..., $\mathcal{T}_{P-1}$ have been successful as well. The seed composed of those successes is a well-formed seed. Given a thread, the gap between this thread and the next one in the new order could indeed not be higher than the gap in the previous order with its next thread. Also the $f^{\text{th}}$ order gaps remain smaller than 1. And as $\mathcal{T}_{e-1}$ succeeds the second time after $f$ failures, it means that the new seed $\mathcal{S}''$ is such that $f(\mathcal{S}'') = f$.
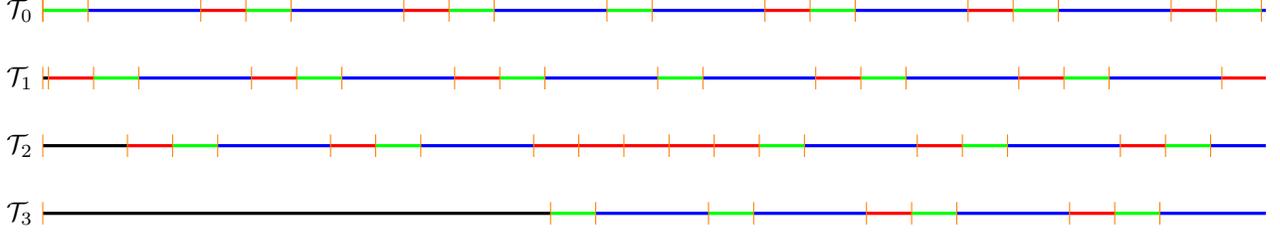
**Figure 8:** Lemma 4 configuration

∎

**Lemma 4.** *Let $\mathcal{S}$ be a weakly-formed seed, and $f = f\left((\mathcal{T}_i, S_i)_{i\in[\![0,P-2]\!]}\right)$. If $G_f^{(f+1)} > 1$ and if the second success of $\mathcal{T}_{P-1}$ occurs before the second success of $\mathcal{T}_{f-1}$, then we can find in the execution a well-formed seed $\mathcal{S}'$ for $P$ threads such that $f(\mathcal{S}') = f$.*

*Proof:* Until the second success of $\mathcal{T}_{P-1}$, the execution follows the same pattern as in Lemma 3. Actually, the case invoked in the current lemma could have been handled in the previous lemma, but it would have implied tricky notations, when we referred to $\mathcal{T}_{rank(n-P+1)}$. Let us deal with this case independently then, and come back to the instant where $\mathcal{T}_{P-1}$ succeeds for the second time.

We had $0 < R_{f-1}^0 - S_{P-1} = G_{f-1}^{(f)} < 1$. For the thread $\mathcal{T}_{\sigma(j)}$ to succeed at its $k^{\text{th}}$ retry after the first success of $\mathcal{T}_{P-1}$ and before $\mathcal{T}_{f-1}$, it should necessary fill the following condition: $j + 1 < R_{\sigma(j)}^k - S_{P-1} < j + 1 + G_{f-1}^{(f)}$. This holds also for the second success of $\mathcal{T}_{P-1}$, which implies that $P' < S_{P-1} + 1 + q + r + h - S_{P-1} < P' + G_{f-1}^{(f)}$, where $h$ is the number of failures of $\mathcal{T}_{P-1}$ before its second success and $P'$ is the number of successes between the two successes of $\mathcal{T}_{P-1}$. As $G_{f-1}^{(f)} < 1$, and $q$, $P'$ and $h$ are non-negative integers, we have $r < G_{f-1}^{(f)}$ and $h = P' - 1 - q$.

To conclude, as any gap at any order is less than the gap between the two successes of $\mathcal{T}_{P-1}$, which is $r < 1$, we found a well-formed seed for $P'$ threads.

Finally any other thread will eventually succeed (see Lemma 2). We can renumber the threads such that $\mathcal{T}_{P'}$ is the first thread that is not in the well-formed seed to succeed, and the threads of the well-formed seed succeeded previously as $\mathcal{T}_0, \ldots, \mathcal{T}_{P'-1}$. As explained before, for all $(k, n) \in [\![0, P'-1]\!]^2$, $G_n^{(k)} < G_n^{(n)} = r$. With the new thread, the first order gaps are changed by decomposing $G_0^{(1)}$ into $G_{P'}^{(1)}$ and the new $G_0^{(1)}$. All gaps can only be decreased, hence we have a new well-formed seed for $P' + 1$ threads. We repeat the process until all threads have been encountered, and obtain in the end $\mathcal{S}'$, a well-formed seed with $P$ threads such that $f(\mathcal{S}') = P - 1 - q$, which is an optimal cyclic execution.

Still, as $\mathcal{T}_f$ succeeds between two successes of $\mathcal{T}_{P-1}$ that are separated by $r$, we had, in the initial configuration: $G_{P-1}^{(P-1-f)} < r$. As, in addition, we have both $G_{f-1}^{(f)} < 1$ and $G_f^{(1)} < 1$, we conclude that the lagging time was initially less than $2 + r$. By hypothesis, we know that $G_f^{(f+1)} > 1$, which implies that, before the entry of the new thread, the lagging time was $1 + r$. In the final execution with one more thread, the lagging time is $r$ and we have one more success in the cycle, thus $f(\mathcal{S}') = f$. ∎

**Theorem 2.** *Assuming $r \neq 0$, if a new thread is added to an $(f, P-1)$-cyclic execution, then all the threads will eventually form either an $(f, P)$-cyclic execution, or an $(f+1, P)$-cyclic execution.*

*Proof:* According to Lemma 2, the new thread will eventually succeed. In addition, we recall that Properties 1 and 2 ensure that before the first success of the new thread, any set of $P - 1$

consecutive successes is a well-formed seed with $P-1$ threads. We then consider a seed (we number the threads accordingly, and number the new thread as $\mathcal{T}_{P-1}$) such that the success of the new thread occurs between the success of $\mathcal{T}_{P-2}$ and $\mathcal{T}_0$; we obtain in this way a weakly-formed seed $\mathcal{S} = (\mathcal{T}_n, S_n)_{n \in [\![0, P-1]\!]\&}$. We differentiate between two cases.

Firstly, if for all $n \in [\![0, P-1]\!]$, $G_n^{(f+1)} < 1$, according to Lemma 1, we can find later in the execution a well-formed seed $\mathcal{S}'$ for $P$ threads such that $f(\mathcal{S}') = f+1$, hence we reach eventually an $(f+1, P)$-cyclic execution.

Let us assume now that this condition is not fulfilled. There exists $n_0 \in [\![0, P-1]\!]$ such that $G_{n_0}^{(f+1)} > 1$. We shift the thread numbers, such that $n_0$ is now $f$, and we have then $G_f^{(f+1)} > 1$. Then two cases are feasible. If the second success of $\mathcal{T}_{P-1}$ occurs before the second success of $\mathcal{T}_{f-1}$, then Lemma 3 shows that we will reach an $(f, P)$-cyclic execution. Otherwise, from Lemma 3, we conclude that an $(f, P)$-cyclic execution will still occur. ∎

### C. Throughput Bounds

Firstly we calculate the expression of throughput and the expected number of threads inside the retry loop (that is needed when we gather expansion and wasted retries). Then we exhibit upper and lower bounds on both throughput and the number of failures, and show that those bounds are reached. Finally, we give the worst case on the number of wasted retries.

**Lemma 5.** *In an $(f, P)$-cyclic execution, the throughput is*

$$T = \frac{P}{q + r + 1 + f}. \tag{8}$$

*Proof:* By definition, the execution is periodic, and the period lasts $q + r + 1 + f$ units of time. As $P$ successes occur during this period, we end up with the claimed expression. ∎

**Lemma 6.** *In an $(f, P)$-cyclic execution, the average number of threads $P_{rl}$ in the retry loop is given by*

$$P_{rl} = P \times \frac{f+1}{q + r + f + 1}.$$

*Proof:* Within a period, each thread spends $f+1$ units of time in the retry loop, among the $q + r + f + 1$ units of time of the period, hence the Lemma. ∎

**Lemma 7.** *The number of failures is not less than $f^{(\text{-})}$, where*

$$f^{(\text{-})} = \begin{cases} P - q - 1 & \text{if } q \leq P - 1 \\ 0 & \text{otherwise} \end{cases}, \text{ and accordingly,} \qquad T \leq \begin{cases} \frac{P}{P+r} & \text{if } q \leq P - 1 \\ \frac{P}{q+r+1} & \text{otherwise.} \end{cases} \tag{9}$$

*Proof:* According to Equation 8, the throughput is maximized when the number of failures is minimized. In addition, we have two lower bounds on the number of failures: (i) $f \geq 0$, and (ii) $P$ successes should fit within a period, hence $q + 1 + f \geq P$. Therefore, if $P - 1 - q < 0$, $T \leq P/(q + r + 1 + 0)$, otherwise,

$$T \leq \frac{P}{q + r + 1 + P - 1 - q} = \frac{P}{P + r}.$$

∎

**Remark 3.** We notice that if $q > P - 1$, the upper bound in Equation 9 is actually the same as the immediate upper bound described in Section III-B1. However, if $q \leq P - 1$, Equation 9 refines the immediate upper bound.

**Lemma 8.** *The number of failures is bounded by*

$$f \leq f^{(+)} = \left\lfloor \frac{1}{2} \left( (P - 1 - q - r) + \sqrt{(P - 1 - q - r)^2 + 4P} \right) \right\rfloor, \text{ and accordingly,}$$

*the throughput is bounded by*

$$T \geq \frac{P}{q + r + 1 + f^{(+)}}.$$

*Proof:* We show that a necessary condition so that an $(f, P)$-cyclic execution, whose lagging time is $\ell$, exists, is $f \times (\ell + r) < P$. According to Property 1, any set of $P$ consecutive successes is a well-formed seed with $P$ threads. Let $\mathcal{S}$ be any of them. As we have $f$ failures before success, Theorem 1 ensures that for all $n \in [\![0, P-1]\!]$, $G_n^{(f)} < 1$. We recall that for all $n \in [\![0, P-1]\!]$, we also have $G_n^{(P)} = \ell + r$.

On the one hand, we have

$$\sum_{n=0}^{P-1} G_n^{(f)} = \sum_{n=0}^{P-1} \sum_{j=n-f+1}^{n} G_{j \bmod P}^{(1)}$$

$$= f \times \sum_{n=0}^{P-1} G_n^{(1)}$$

$$\sum_{n=0}^{P-1} G_n^{(f)} = f \times (\ell + r).$$

On the other hand, $\sum_{n=0}^{P-1} G_n^{(f)} < \sum_{n=0}^{P-1} 1 = P$.

Altogether, the necessary condition states that $f \times (\ell + r) < P$, which can be rewritten as $f \times (q + 1 + f - P + r) < P$. The proof is complete since minimizing the throughput is equivalent to maximizing the number of failures. ∎

**Lemma 9.** *For each of the bounds defined in Lemmas 7 and 8, there exists an $(f, P)$-cyclic execution that reaches the bound.*

*Proof:* According to Lemmas 7 and 8, if an $(f, P)$-cyclic execution exists, then the number of failures is such that $f^{(-)} \leq f \leq f^{(+)}$.
We show now that this double necessary condition is also sufficient. We consider $f$ such that $f^{(-)} \leq f \leq f^{(+)}$, and build a well-formed seed $\mathcal{S} = (\mathcal{T}_i, S_i)_{i \in [\![0, P-1]\!]}$.
For all $n \in [\![0, P-1]\!]$, we define $S_i$ as

$$S_n = n \times \left( \frac{q + 1 + f - P + r}{P} + 1 \right).$$

We first show that $f(\mathcal{S}) = f$. By definition, $f(\mathcal{S}) = \max(0, \lceil S_{P-1} - S_0 - q - r \rceil)$; we have then

$$f(\mathcal{S}) = \max\left( 0, \left\lceil (P-1) \times \left( \frac{q + 1 + f - P + r}{P} + 1 \right) - q - r \right\rceil \right)$$

$$= \max\left( 0, \left\lceil (P - 1 - q - r) + (q + 1 + f - P + r) - \frac{q + 1 + f - P + r}{P} \right\rceil \right)$$

$$f(\mathcal{S}) = \max\left( 0, \left\lceil f - \frac{q + 1 + f - P + r}{P} \right\rceil \right).$$

Firstly, we know that $q + 1 + f - P \geq 0$, thus if $f = 0$, then the second term of the maximum is not positive, and $f(\mathcal{S}) = 0 = f$. Secondly, if $f > 0$, then according to Lemma 7, $(q + 1 + f - P + r)/P < 1/f \leq 1$. As we also have $(q + 1 + f - P + r)/P \geq 0$, we conclude that $f(\mathcal{S}) = \left\lceil f - \frac{q+1+f-P+r}{P} \right\rceil = f$.

Additionally, for all $n \in [\![0, P - 1]\!]$,

$$G_n^{(f)} = \begin{cases} S_n - S_{n-f} - f & \text{if } n > f \\ S_n - S_{P+n-f} + 1 + q + r & \text{otherwise} \end{cases}$$

$$= \begin{cases} n \times \left(\frac{q+1+f-P+r}{P} + 1\right) - (n - f) \times \left(\frac{q+1+f-P+r}{P} + 1\right) - f \\ n \times \left(\frac{q+1+f-P+r}{P} + 1\right) - (P + n - f) \times \left(\frac{q+1+f-P+r}{P} + 1\right) + 1 + q + r \end{cases}$$

$$= \begin{cases} f \times \frac{q+1+f-P+r}{P} \\ -(P - f) - (q + 1 + f - P + r) + f \times \frac{q+1+f-P+r}{P} + 1 + q + r \end{cases}$$

$$G_n^{(f)} = f \times \frac{w + r}{P}$$

As $w \leq 0$ and $f \leq 0$, $G_n^{(f)} > 0$. Since $f \leq f^{(+)}$, $G_n^{(f)} < 1$. Theorem 1 implies that $\mathcal{S}$ is a well-formed seed that leads to an $(f, P)$-cyclic execution.

We have shown that for all $f$ such that $f^{(-)} \leq f \leq f^{(+)}$ there exists an $(f, P)$-cyclic execution; in particular there exist an $(f^{(+)}, P)$-cyclic execution and an $(f^{(-)}, P)$-cyclic execution. ∎

**Corollary 1.** *The highest possible number of wasted repetitions is $\left\lceil \sqrt{P} - 1 \right\rceil$ and is achieved when $P = q + 1$.*

*Proof:*

The highest possible number of wasted repetitions $\tilde{w}(P)$ with $P$ threads is given by

$$\tilde{w}(P) = f^{(+)} - f^{(-)} = \left\lfloor \frac{1}{2} \left(-a(P) + \sqrt{a(P)^2 + 4P}\right) - f^{(-)} \right\rfloor.$$

Let $a$ and $h$ be the functions respectively defined as $a(P) = q + 1 - P + r$, which implies $a'(P) = -1$, and $h(P) = (-a(P) + \sqrt{a(P)^2 + 4P})/2 - f^{(-)}$, so that $\tilde{w}(P) = \lfloor h(P) \rfloor$.

Let us first assume that $a(P) > 0$. In this case, $q \leq P - 1$, hence $f^{(-)} = 0$. We have

$$2h'(P) = 1 + \frac{-2a(P) + 4}{2\sqrt{a(P)^2 + 4P}}$$

$$2h'(P) = 2 \times \frac{2 - a(P) + \sqrt{a(P)^2 + 4P}}{2\sqrt{a(P)^2 + 4P}}$$

Therefore, $h'(P)$ is negative if and only if $\sqrt{a(P)^2 + 4P} < a(P) - 2$. It cannot be true if $a(P) < 2$. If $a(P) \geq 2$, then the previous inequality is equivalent to $a(P)^2 + 4P < a(P)^2 - 4a(P) + 4$, which can be rewritten in $q + 1 + r < 1$, which is absurd. We have shown that $h$ is increasing in $]0, q + 1]$.

Let us now assume that $a(P) \leq 0$. In this case, $q > P - 1$, hence $f^{(-)} = P - q - 1$, and $h'(P) = \left(a(P) + \sqrt{a(P)^2 + 4P}\right)/2 - r$. Assuming $h'(P)$ to be positive leads to the same absurd inequality $q + 1 + r < 1$, which proves that $h$ is decreasing on $[q + 2, +\infty[$.

Also, the maximum number of wasted repetitions is achieved as $P = q + 1$ or $P = q + 2$. Since

$$h(q + 1) = \frac{1}{2} \left(-r + \sqrt{r^2 + 4P}\right) > \frac{1}{2} \left(-(r + 1) + \sqrt{r^2 + 4P}\right) = h(q + 2),$$

the maximum number of wasted repetitions is $\tilde{w}(q + 1)$. In addition,

$$\begin{aligned} \frac{1}{2}\left(-r + \sqrt{4P}\right) &< h(q + 1) &< \frac{1}{2}\left(-r + \sqrt{r^2} + \sqrt{4P}\right) \\ \sqrt{P} - \frac{r}{2} &< h(q + 1) &< \sqrt{P} \\ \sqrt{P} - 1 &\leq h(q + 1) &< \sqrt{P} \end{aligned}$$

We conclude that the maximum number of wasted repetitions is $\left\lceil \sqrt{P} - 1 \right\rceil$. ∎

## V. Expansion and Complete Throughput Estimation

### A. Expansion

Interference of threads does not only lead to logical conflicts but also to hardware conflicts which impact the performance significantly. We model the behavior of the cache coherency protocols which determine the interaction of overlapping *Read*s and *CAS*s. By taking MESIF [GH09] as basis, we come up with the following assumptions. When executing an atomic *CAS*, the core gets the cache line in exclusive state and does not forward it to any other requesting core until the instruction is retired. Therefore, requests stall for the release of the cache line which implies serialization. On the other hand, ongoing *Read*s can overlap with other operations. As a result, a *CAS* introduces expansion only to overlapping *Read* and *CAS* operations that start after it, as illustrated in Figure 4. As a remark, we ignore memory bandwidth issues which are negligible for our study.

Furthermore, we assume that *Read*s that are executed just after a *CAS* do not experience expansion (as the thread already owns of the data), which takes effect at the beginning of a retry following a failing attempt. Thus, read expansions need only to be considered before the $0^{\text{th}}$ retry. In this sense, read expansion can be moved to parallel section and calculated in the same way as *CAS* expansion is calculated.

To estimate expansion, we consider the delay that a thread can introduce, provided that there is already a given number of threads in the retry loop. The starting point of each *CAS* is a random variable which is distributed uniformly within an expanded retry. The cost function $d$ provides the amount of delay that the additional thread introduces, depending on the point where the starting point of its *CAS* hits. By using this cost function we can formulate the expansion increase that each new thread introduces and derive the differential equation below to calculate the expansion of a *CAS*.

**Lemma 10.** *The expansion of a CAS operation is the solution of the following system of equations:*

$$
\begin{cases}
e'\left(P_{rl}\right) &= cc \times \dfrac{\frac{cc}{2} + e\left(P_{rl}\right)}{rc + cw + cc + e\left(P_{rl}\right)} \\
e\left(P_{rl}^{(0)}\right) &= 0
\end{cases}
, \qquad
\begin{array}{l}
\textit{where } P_{rl}^{(0)} \textit{ is the point where} \\
\textit{expansion begins.}
\end{array}
$$

*Proof:*

We compute $e\left(P_{rl} + h\right)$, where $h \leq 1$, by assuming that there are already $P_{rl}$ threads in the retry loop, and that a new thread attempts to *CAS* during the retry, within a probability $h$.

$$e\left(P_{rl}+h\right) = e\left(P_{rl}\right) + h\times \int_{0}^{rlw^{(+)}} \frac{d\left(t\right)}{rlw^{(+)}}\, dt$$

$$= e\left(P_{rl}\right) + h\times \Big( \int_{0}^{rc+cw-cc} \frac{d\left(t\right)}{rlw^{(+)}}\, dt$$

$$+ \int_{rc+cw-cc}^{rc+cw} \frac{d\left(t\right)}{rlw^{(+)}}\, dt$$

$$+ \int_{rc+cw}^{rc+cw+e(P_{rl})} \frac{d\left(t\right)}{rlw^{(+)}}\, dt$$

$$+ \int_{rc+cw+e(P_{rl})}^{rlw^{(+)}} \frac{d\left(t\right)}{rlw^{(+)}}\, dt \Big)$$

$$= e\left(P_{rl}\right) + h\times \Big( \int_{rc+cw-cc}^{rc+cw} \frac{t}{rlw^{(+)}}\, dt$$

$$+ \int_{rc+cw}^{rc+cw+e(P_{rl})} \frac{cc}{rlw^{(+)}}\, dt \Big)$$

$$e\left(P_{rl}+h\right) = e\left(P_{rl}\right) + h\times \frac{\frac{cc^2}{2} + e\left(P_{rl}\right)\times cc}{rlw^{(+)}}$$

This leads to $\dfrac{e\left(P_{rl}+h\right) - e\left(P_{rl}\right)}{h} = \dfrac{\frac{cc^2}{2} + e\left(P_{rl}\right)\times cc}{rlw^{(+)}}$. When making $h$ tend to 0, we finally obtain

$$e'\left(P_{rl}\right) = cc \times \frac{\frac{cc}{2} + e\left(P_{rl}\right)}{rc + cw + cc + e\left(P_{rl}\right)}.$$

∎

### B. Throughput Estimate

There remains to combine hardware and logical conflicts in order to obtain the final upper and lower bounds on throughput. We are given as an input an expected number of threads $P_{rl}$ inside the retry loop. We firstly compute the expansion accordingly, by solving numerically the differential equation of Lemma 10. As explained in the previous subsection, we have $pw^{(+)} = pw + e$, and $rlw^{(+)} = rc + cw + e + cc$. We can then compute $q$ and $r$, that are the inputs (together with the total number of threads $P$) of the method described in Section IV. Assuming that the initialization times of the threads are spaced enough, the execution will superimpose an $(f, P)$-cyclic execution. Thanks to Lemma 6, we can compute the average number of threads inside the retry loop, that we note by $h_f(P_{rl})$. A posteriori, the solution is consistent if this average number of threads inside the retry loop $h_f(P_{rl})$ is equal to the expected number of threads $P_{rl}$ that has been given as an input.

Several $(f, P)$-cyclic executions belong to the domain of the possible outcomes, but we are interested in upper and lower bounds on the number of failures $f$. We can compute them through Lemmas 7 and 8, along with their corresponding throughput and average number of threads inside the retry loop. We note by $h^{(+)}(P_{rl})$ and $h^{(-)}(P_{rl})$ the average number of threads for the lowest number of failures and highest one, respectively. Our aim is finally to find $P_{rl}^{(-)}$ and $P_{rl}^{(+)}$, such that $h^{(+)}(P_{rl}^{(+)}) = P_{rl}^{(+)}$ and $h^{(-)}(P_{rl}^{(-)}) = P_{rl}^{(-)}$. If several solutions exist, then we want to keep the smallest, since the retry loop stops to expand when a stable state is reached.

Note that we also need to provide the point where the expansion begins. It begins when we start to have failures, while reducing the parallel section. Thus this point is $(2P-1)rlw^{(-)}$ (resp. $(P-1)rlw^{(-)}$) for the lower (resp. upper) bound on the throughput.

**Theorem 3.** *Let $(x_n)$ be the sequence defined recursively by $x_0 = 0$ and $x_{n+1} = h^{(+)}(x_n)$. If $pw \geq rc + cw + cc$, then*

$$P_{rl}^{(+)} = \lim_{n \to +\infty} x_n.$$

*Proof:* First of all, the average number of threads belongs to $]0, P[$, thus for all $x \in [0, P]$, $0 < h^{(+)}(x) < P$. In particular, we have $h^{(+)}(0) > 0$, and $h^{(+)}(P) < P$, which proves that there exist one fixed point for $h^{(+)}$.

In addition, we show that $h^{(+)}$ is a non-decreasing function. According to Lemma 6,

$$h^{(+)}(P_{rl}) = P \times \frac{1 + f^{(-)}}{q + r + f^{(-)} + 1},$$

where all variables except $P$ depend actually on $P_{rl}$. We have

$$q = \left\lfloor \frac{pw + e}{rlw^{(-)} + e} \right\rfloor \text{ and } r = \frac{pw + e}{rlw^{(-)} + e} - q,$$

hence, if $pw \geq rlw^{(-)}$, $q$ and $r$ are non-increasing as $e$ is non-decreasing, which is non-decreasing with $P_{rl}$. Since $f^{(-)}$ is non-decreasing as a function of $q$, we have shown that if $pw \geq rlw^{(-)}$, $h^{(+)}$ is a non-decreasing function.

Finally, the proof is completed by the theorem of Knaster-Tarski. ∎

The same line of reasoning holds for $h^{(-)}$ as well. As a remark, w point out that when $pw < rlw^{(-)}$, we scan the interval of solution, and have no guarantees about the fact that the solution is the smallest one; still it corresponds to very extreme cases.

### C. Several Retry Loops

We consider here a lock-free algorithm that, instead of being a loop over one parallel section and one retry loop, is composed of a loop over a sequence of alternating parallel sections and retry loops. We show that this algorithm is equivalent to an algorithm with only one parallel section and one retry loop, by proving the intuition that the longest retry loop is the only one that fails and hence expands.

*1) Problem Formulation:* In this subsection, we consider an execution such that each spawned thread runs Procedure Combined in Figure 9. Each thread executes a linear combination of $S$ independent retry loops, *i.e.* operating on separate variables, interleaved with parallel sections. We note now as $rlw_i^{(+)}$ and $pw_i^{(+)}$ the size of a retry of the $i^{th}$ retry loop and the size of the $i^{th}$ parallel section, respectively, for each $i \in [\![1, S]\!]$. As previously, $q_i$ and $r_i$ are defined such that $pw_i^{(+)} = (q_i + r_i) \times rlw_i^{(+)}$, where $q_i$ is a non-negative integer and $r_i$ is smaller than 1.

---

**Procedure** Combined

**1** Initialization();
**2** **while** *!* done **do**
**3**  **for** $i \leftarrow 1$ **to** S **do**
**4**    Parallel_Work($i$);
**5**    **while** *!* success **do**
**6**      current $\leftarrow$ Read(AP[$i$]);
**7**      new $\leftarrow$ Critical_Work($i$,current);
**8**      success $\leftarrow$ CAS(AP, current, new);

---

**Figure 9:** Thread procedure with several retry loops

The Procedure Combined executes the retry loops and parallel sections in a cyclic fashion, so we can normalize the writing of this procedure by assuming that a retry of the $1^{\text{st}}$ retry loop is the longest one. More precisely, we consider the initial algorithm, and we define $i_0$ as

$$i_0 = \min \operatorname{argmax}_{i \in [\![1, S]\!]} rlw_i^{(+)}.$$

We then renumber the retry loops such that the new ordering is $i_0, \ldots, S, 1, \ldots, i_0 - 1$, and we add in Initialization the first parallel sections and retry loops on access points from $1$ to $i_0$ — according to the initial ordering.

One success at the system level is defined as one success of the last *CAS*, and the throughput is defined accordingly. We note that in steady-state, all retry loops have the same throughput, so the throughput can be computed from the throughput of the $1^{\text{st}}$ retry loop instead.

*2) Wasted Retries:*

**Lemma 11.** *Unsuccessful retry loops can only occur in the $1^{\text{st}}$ retry loop.*

    *Proof:*

We note $(t_n)_{n \in [1,+\infty[}$ the sequence of the thread numbers that succeeds in the $1^{\text{st}}$ retry loop, and $(s_n)_{n \in [1,+\infty[}$ the sequence of the corresponding time where they exit the retry loop. We notice that by construction, for all $n \in [1, +\infty[$, $s_n < s_{n+1}$. Let, for $i \in [\![2, S]\!]$ and $n \in [1, +\infty[$, $(\mathcal{P}_{i,n})$ be the following property: for all $i' \in [\![2, i]\!]$, and for all $n' \in [\![1, n]\!]$, the thread $\mathcal{T}_{t_{n'}}$ succeeds in the $i^{\text{th}}$ retry loop at its first attempt.

We assume that for a given $(i, n)$, $(\mathcal{P}_{i+1,n})$ and $(\mathcal{P}_{i,n+1})$ is true, and show that $(\mathcal{P}_{i+1,n+1})$ is true. As the threads $\mathcal{T}_{t_n}$ and $\mathcal{T}_{t_{n+1}}$ do not have any failure in the first $i$ retry loops, their entrance time in the $i + 1^{\text{th}}$ retry loop is given by

$$s_n + \sum_{i'=1}^{i} (rlw_{i'}^{(+)} + pw_{i'}^{(+)}) + pw_{i+1}^{(+)} = X_1 \text{ and } s_{n+1} + \sum_{i'=1}^{i} (rlw_{i'}^{(+)} + pw_{i'}^{(+)}) + pw_{i+1}^{(+)} = X_2,$$

respectively. Thread $\mathcal{T}_{t_n}$ does not fail in the $i + 1^{\text{th}}$ retry loop, hence exits at

$$X_1 + rlw_{i+1}^{(+)} < X_1 + rlw_1^{(+)} = s_n + X_2 - s_{n+1} < X_2.$$

As the previous threads $\mathcal{T}_{n-1}, \ldots, \mathcal{T}_1$ exits the $i^{\text{th}}$ retry loop before $\mathcal{T}_n$, and next threads $\mathcal{T}_{n'}$, where $n' > n+1$, enters this retry loop after $\mathcal{T}_{n+1}$, this implies that the thread $\mathcal{T}_{t_{n+1}}$ succeeds in the $i + 1^{\text{th}}$ retry loop at its first attempt, and $(\mathcal{P}_{i+1,n+1})$ is true.

Regarding the first thread that succeeds in the first retry loop, we know that he successes in any retry loop since there is no other thread to compete with. Therefore, for all $i \in [\![2, S]\!]$, $(\mathcal{P}_{i,1})$ is true. Then we show by induction that all $(\mathcal{P}_{2,n})$ is true, then all $(\mathcal{P}_{3,n})$, *etc.*, until all $(\mathcal{P}_{S,n})$, which concludes the proof. ∎

**Theorem 4.** *The multi-retry loop Procedure Combined is equivalent to the Procedure Abstract-Algorithm, where*

$$pw^{(+)} = pw_1^{(+)} + \sum_{i=2}^{S} \left( pw_i^{(+)} + rlw_i^{(+)} \right) \quad and \quad rlw^{(+)} = rlw_1^{(+)}.$$

*Proof:* According to Lemma 11 there is no failure in other retry loop than the first one; therefore, all retry loops have a constant duration, and can thus be considered as parallel sections. ∎

*3) Expansion:* The expansion in the retry loop starts as threads fail inside this retry loop. When threads are launched, there is no expansion, and Lemma 11 implies that if threads fail, it should be inside the first retry loop, because it is the longest one. As a result, there will be some stall time in the memory accesses of this first retry loop, *i.e.* expansion, and it will get even longer. Failures will thus still occur in the first retry loop: there is a positive feedback on the expansion of the first retry loop that keeps this first retry loop as the longest one among all retry loops. Therefore, in accordance to Theorem 4, we can compute the expansion by considering the equivalent single-retry loop procedure described in the theorem.

## VI. Experimental Evaluation

We validate our model and analysis framework through successive steps, from synthetic tests, capturing a wide range of possible abstract algorithmic designs, to several reference implementations of extensively studied lock-free data structure designs that include cases with non-constant parallel section and retry loop.

### A. Setting

We have conducted experiments on an Intel ccNUMA workstation system. The system is composed of two sockets, that is equipped with Intel Xeon E5-2687W v2 CPUs with frequency band 1.2-3.4. GHz The physical cores have private L1, L2 caches and they share an L3 cache, which is 25 MB. In a socket, the ring interconnect provides L3 cache accesses and core-to-core communication. Due to the bi-directionality of the ring interconnect, uncontended latencies for intra-socket communication between cores do not show significant variability.

Our model assumes uniformity in the *CAS* and *Read* latencies on the shared cache line. Thus, threads are pinned to a single socket to minimize non-uniformity in *Read* and *CAS* latencies. In the experiments, we vary the number of threads between 4 and 8 since the maximum number of threads that can be used in the experiments are bounded by the number of physical cores that reside in one socket.

In all figures, y-axis provides the throughput, which is the number of successful operations completed per millisecond. Parallel work is represented in x-axis in cycles. The graphs contain the high and low estimates (see Section IV), corresponding to the lower and upper bound on the wasted retries, respectively, and an additional curve that shows the average of them.

As mentioned before, the latencies of *CAS* and *Read* are parameters of our model. We used the methodology described in [DGT13] to measure latencies of these operations in a benchmark program by using two threads that are pinned to the same socket. The aim is to bring the cache line into the state used in our model. Our assumption is that the *Read* is conducted on an invalid line. For *CAS*, the state of the cache line could be exclusive, forward, shared or invalid. Regardless of the state of the cache line, *CAS* requests it for ownership, that compels invalidation in other cores, which in turn incurs a two-way communication and a memory fence afterwards to assure atomicity. Thus, the latency of *CAS* does not show negligible variability with respect to the state of the cache line, as also revealed in our latency benchmarks.

As for the computation cost, the work inside the parallel section is implemented by a dummy for-loop of *Pause* instructions.

### B. Synthetic Tests

*1) Single retry loop:* For the evaluation of our model, we first create synthetic tests that emulate different design patterns of lock-free data structures (value of *cw*) and different application contexts (value of *pw*). As described in the previous subsection, in the Procedure AbstractAlgorithm, the
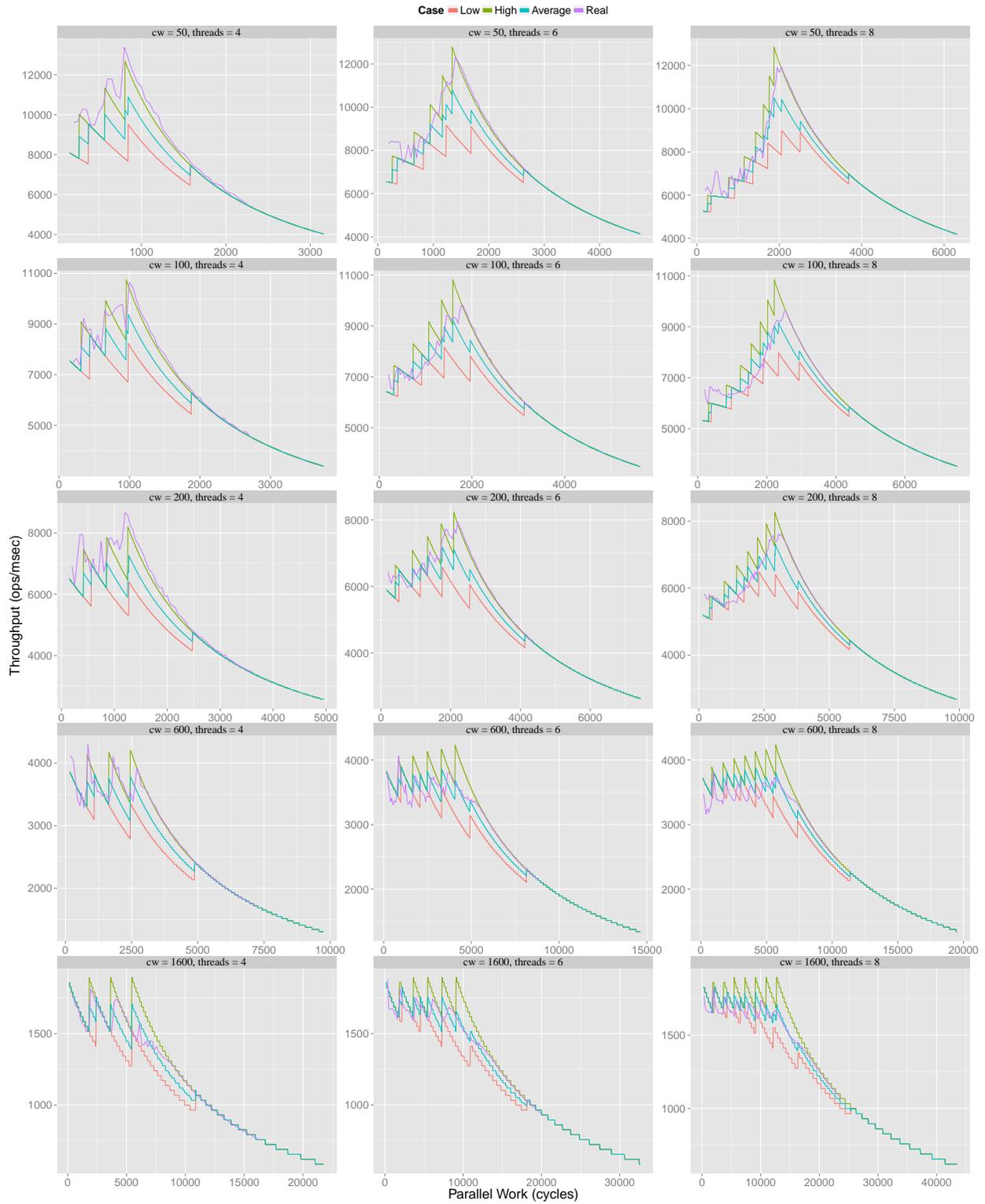
**Figure 10:** Synthetic program

amount of work in both the parallel section and the retry loop are implemented as dummy loops, whose costs are adjusted through the number of iterations in the loop.

Generally speaking, in Figure 10, we observe two main behaviors: when *pw* is high, the data structure is not contended, and threads can operate without failure. When *pw* is low, the data structure is contended, and depending on the size of *cw* (that drives the expansion) a steep decrease in throughput or just a roughly constant bound on the performance is observed.

The position of the experimental curve between the high and low estimates, depends on *cw*. It can be observed that the experimental curve mostly tends upwards as *cw* gets smaller, possibly because the serialization of the *CAS*s helps the synchronization of the threads.

Another interesting fact is the waves appearing on the experimental curve, especially when the number of threads is low or the critical work big. This behavior is originating because of the variation of $r$ with the change of parallel work, a fact that is captured by our analysis.
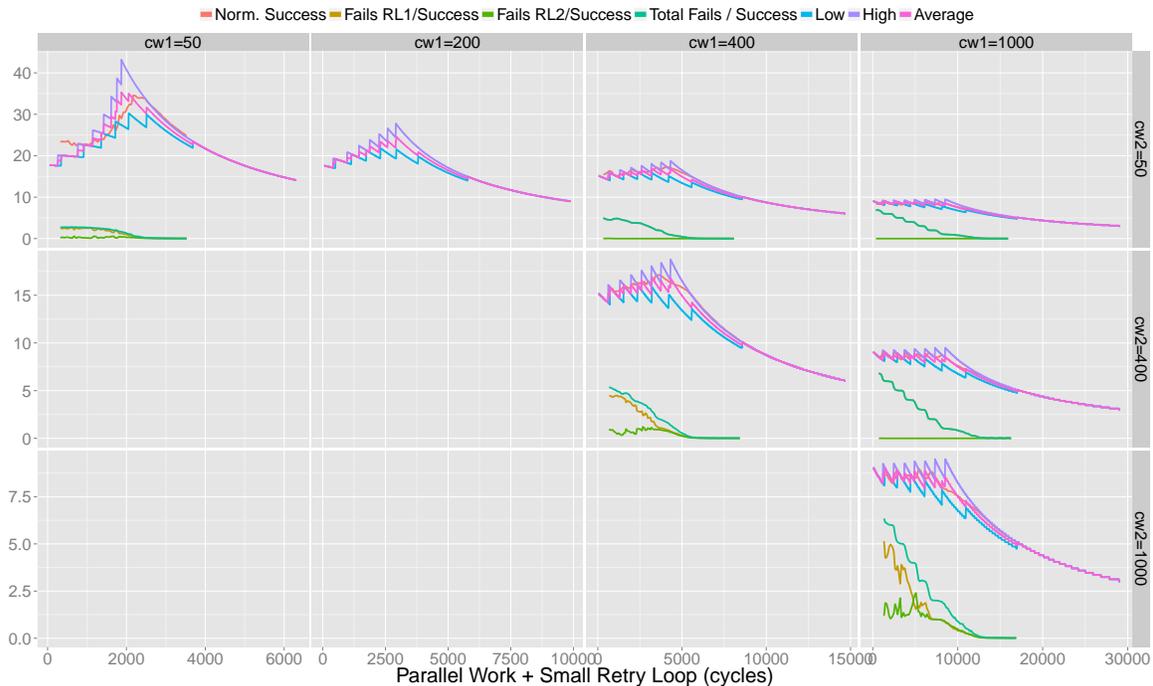


**Figure 11:** Multiple retry loops with 8 threads

*2) Several retry loops:* We have created experiments by combining several retry loops, each operating on an independent variable which is aligned to a cache line. In Figure 11, results are compared with the model for single retry loop case where the single retry loop is equal to the longest retry loop, while the other retry loops are part of the parallel section. The distribution of fails in the retry loops are illustrated and all throughput curves are normalized with a factor of 175 (to be easily seen in the same graph). Fails per success values are not normalized and a success is obtained after completing all retry loops.

### C. Treiber's Stack

The lock-free stack by Treiber [Tre86] is one of the most studied efficient data structures. Pop and Push both contain a retry loop, such that each retry starts with a *Read* and ends with *CAS* on the shared top pointer. In order to validate our model, we start by using Pops. From a stack which is
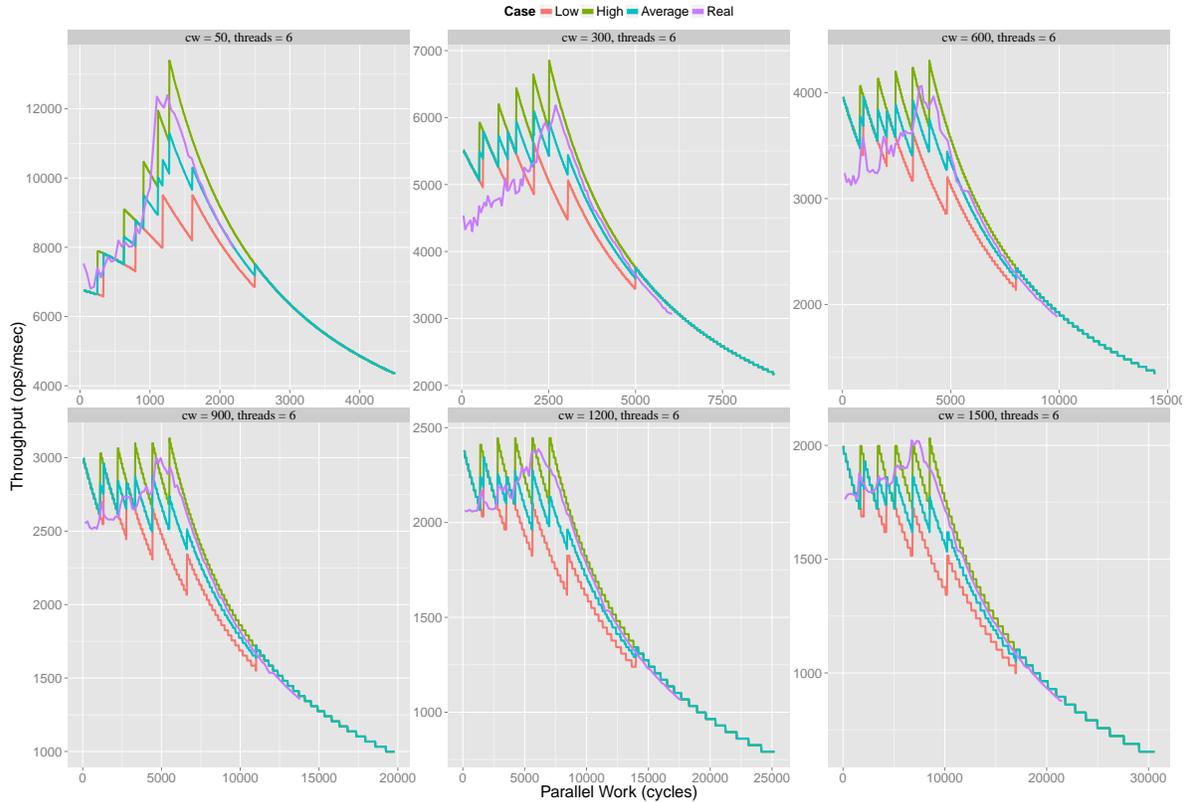
**Figure 12:** Pop on Treiber's stack

initiated with 50 million elements, threads continuously pop elements for a given amount of time. We count the total number of pop operations per millisecond. Each Pop first reads the top pointer and gets the next pointer of the element to obtain the address of the second element in the stack, before attempting to *CAS* with the address of the second element. The access to the next pointer of the first element occurs in between the *Read* and the *CAS*. Thus, it represents the work in *cw*. This memory access can possibly introduce a costly cache miss depending on the locality of the popped element.

To validate our model with different *cw* values, we make use of this costly cache miss possibility. We allocate a contiguous chunk of memory and align each element to a cache line. Then, we initialize the stack by pushing elements from contiguous memory either with a single or large stride to disable the prefetcher. When we measure the latency of *cw* in Pop for single and large stride cases, we obtain the values that are approximately 50 and 300 cycles, respectively. As a remark, 300 cycles is the cost of an L3 miss in our system when it is serviced from the local main memory module. To create more test cases with larger *cw*, we extended the stack implementation to pop multiple elements with a single operation. Thus, each access to the next element could introduce an additional L3 cache miss while popping multiple elements. By doing so, we created cases in which each thread pops 2, 3, *etc.* elements, and *cw* goes to 600, 900, *etc.* cycles, respectively. In Figure 12, comparison of the experimental results from Treiber's stack and our model is provided.

As a remark, we did not implemented memory reclamation for our experiments but one can implement a stack that allows pop and push of multiple elements with small modifications using hazard pointers [Mic04]. Pushing can be implemented in the same way as single element case. A

---

**Algorithm 1:** Multiple Pop

---

**1** Pop (multiple)
**2** **while** *true* **do**
**3**   t = Read(top);
**4**   **for** *multiple* **do**
**5**     **if** *t = NULL* **then**
**6**     │ return EMPTY;
**7**     hp* = t;
**8**     **if** *top != t* **then**
**9**     │ break;
**10**    hp++;
**11**    next = t.next;
**12**   **if** *CAS(&top, t, next)* **then**
**13**   │ break;
**14** RetireNodes (t, multiple);

---

**Pop** requires some modifications for memory reclamation. It can be implemented by making use of hazard pointers just by adding the address of the next element to the hazard list before jumping to it. Also, the validity of top pointer should be checked after adding the pointer to the hazard list to make sure that other threads are aware of the newly added hazard pointer. By repeating this process, a thread can jump through multiple elements and pop all of them with a *CAS* at the end.

### D. *Shared Counter*

In [DLM13], the authors have implemented a "scalable statistics counters" relying on the following idea: when contention is low, the implementation is a regular concurrent counter with a *CAS*; when the counter starts to be contended, it switches to a statistical implementation, where the counter is actually incremented less frequently, but by a higher value. One key point of this algorithm is the switch point, which is decided thanks to the number of failed increments; our model can be used by providing the peak point of performance of the regular counter implementation as the switch point. We then have implemented a shared counter which is basically a *Fetch-and-Increment* using a *CAS*, and compared it with our analysis. The result is illustrated in Figure 13, and shows that the parallel section size corresponding to the peak point is correctly estimated using our analysis.

### E. DeleteMin *in Priority List*

We have applied our model to DeleteMin of the skiplist based priority queue designed in [LJ13]. DeleteMin traverses the list from the beginning of the lowest level, finds the first node that is not logically deleted, and tries to delete it by marking. If the operation does not succeed, it continues with the next node. Physical removal is done in batches when reaching a threshold on the number of deleted prefixes, and is followed by a restructuring of the list by updating the higher level pointers, which is conducted by the thread that is successful in redirecting the head to the node deleted by itself.

We consider the last link traversal before the logical deletion as critical work, as it continues with the next node in case of failure. The rest of the traversal is attributed to the parallel section as the threads can proceed concurrently without interference. We measured the average cost of a traversal under low contention for each number of threads, since traversal becomes expensive with
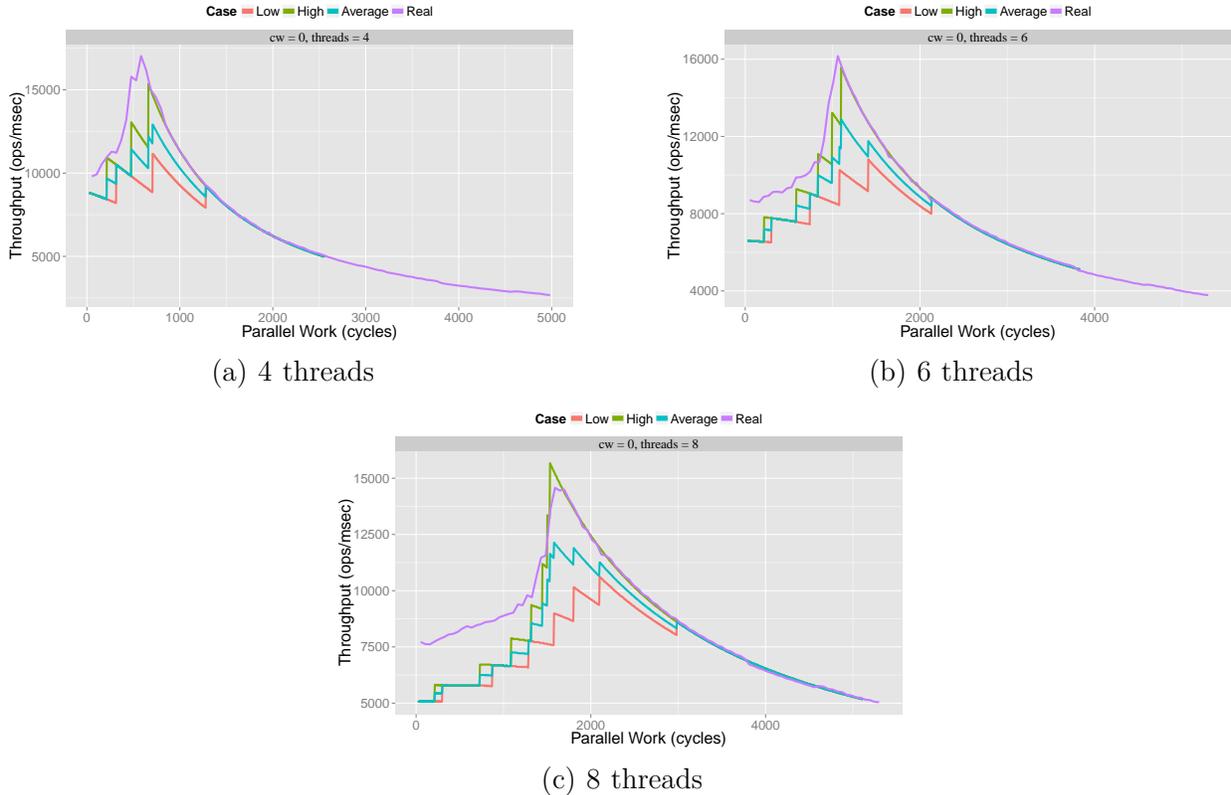
(a) 4 threads



(b) 6 threads



(c) 8 threads

**Figure 13:** Increment on a shared counter

more threads. In addition, average cost of restructuring is also included in the parallel section since it is executed infrequently by a single thread.

We initialize the priority queue with a large set of elements. As illustrated in Figure 14, the smallest $pw$ value is not zero as the average cost of traversal and restructuring is intrinsically included. The peak point is in the estimated place but the curve does not go down sharply under high contention. This presumably occurs as the traversal might require more than one steps (link access) after a failed attempt, which creates a back-off effect.

*F. Enqueue-Dequeue on a Queue*

In order to demonstrate the validity of the model with several retry loops, and that the results covers a wider spectrum of application and designs from the ones we focused in our model, we studied the following setting: the threads share a queue, and each thread enqueues an element, executes the parallel section, dequeues an element, and reiterates. We consider the queue implementation by Michael and Scott [MS96], that is usually viewed as the reference queue while looking at lock-free queue implementations.

Dequeue operations fit immediately into our model but Enqueue operations need an adjustment due to the helping mechanism. Note that without this helping mechanism, a simple queue implementation would fit directly, but we also want to show that the model is malleable, *i.e.* the fundamental behavior remains unchanged even if we divert slightly from the initial assumptions. We consider an equivalent execution that catches up with the model, and use it to approximate the performance of the actual execution of Enqueue.
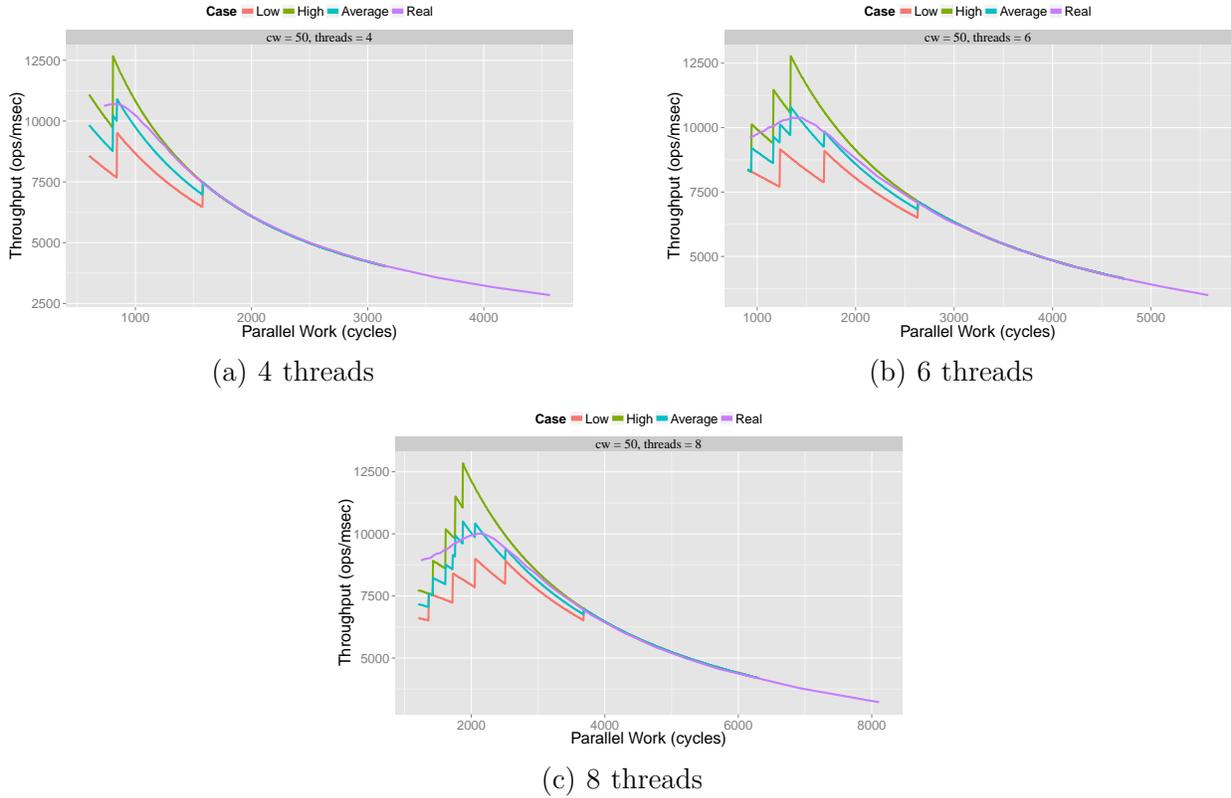
(a) 4 threads



(b) 6 threads



(c) 8 threads

**Figure 14:** DeleteMin on a priority list

Enqueue is composed of two steps. Firstly, the new node is attached to the last node of the queue via a *CAS*, that we denote by $CAS_A$, leading to a transient state. Secondly, the tail is redirected to point to the new node via another *CAS*, that we denote by $CAS_B$, which brings back the queue into a steady state.

A new Enqueue can not proceed before the two steps of previous success are completed. The first step is the linearization point of operation and the second step could be conducted by a different thread through the helping mechanism. In order to start a new Enqueue, concurrent Enqueues help the completion of the second step of the last success if they find the queue in the transient state. Alternatively, they try to attach their node to the queue if the queue is in the steady state at the instant of check. This process continues until they manage to attach their node to the queue via a retry loop in which state is checked and corresponding CAS is executed.

The flow of an Enqueue is determined by this state checks. Thus, an Enqueue could execute multiple $CAS_B$ (successful or failing) and multiple $CAS_A$ (failing) in an interleaved manner, before succeeding in $CAS_A$ at the end of the last retry. If we assume that both states are equally probable for a check instant which will then end up with a retry, the number of CAS s that ends up with a retry are expected to be distributed equally among $CAS_A$ and $CAS_B$ for each thread. In addition, each thread has a successful $CAS_A$ (which linearizes the Enqueue) and a $CAS_B$ at the end of the operation which could either be successful or failed by a concurrent helper thread.

We imitate such an execution with an equivalent execution in which threads keep the same relative ordering of the invocation, return from Enqueue together with same result. In equivalent execution, threads alternate between $CAS_A$ and $CAS_B$ in their retries, and both steps of successful operation is
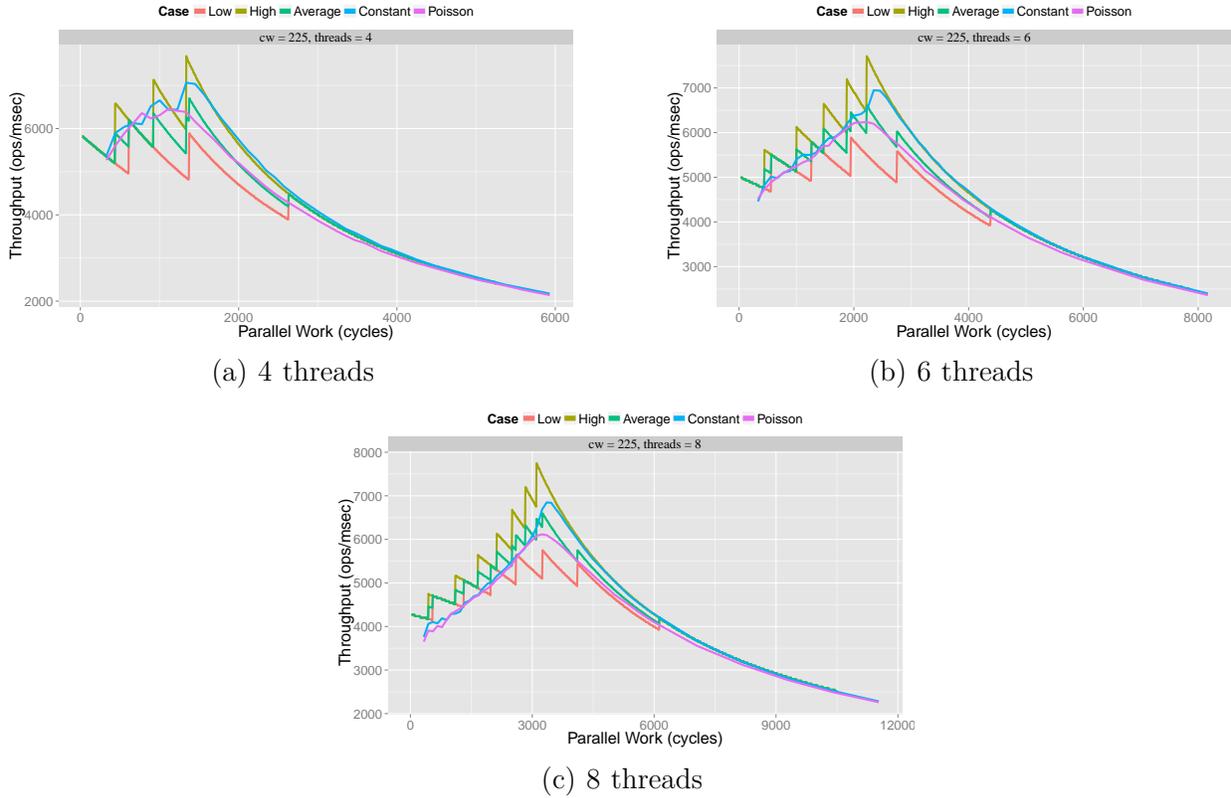
**Figure 15:** Enqueue-Dequeue on Michael and Scott queues

conducted by the same thread. The equivalent execution can be obtained by thread-wise reordering of CAS s that leads to a retry and exchanging successful $CAS_B$ s with the failed counterparts at the end of an Enqueue, as the latter ones indeed fail because of this success of helper threads. The model can be applied to this equivalent execution by attributing each $CAS_A$-$CAS_B$ couple to a single iteration and represent it as a larger retry loop since the successful couple can not overlap with another successful one and all overlapping ones fail. With a straightforward extension of the expansion formula, we accomodate the $CAS_A$ in the critical work which can also expand, and use $CAS_B$ as the *CAS* of our model.

In addition, we take one step further outside the analysis by including a new case, where the parallel section follows a Poisson distribution, instead of being constant. *pw* is chosen as the mean to generate Poisson distribution instead of taking it constant. The results are illustrated in Figure 15. Our model provides good estimates for the constant *pw* and also reasonable results for the Poisson distribution case, although this case deviates from (/extends) our model assumptions. The advantage of regularity, which brings synchronization to threads, can be observed when the constant and Poisson distributions are compared. In the Poisson distribution, the threads start to fail with larger *pw*, which smoothes the curve around the peak of the throughput curve.

### G. Discussion

In this subsection we discuss the adequacy of our model, specifically the cyclic argument, to capture the behavior that we observe in practice. Figure 16 illustrates the frequency of occurrence of
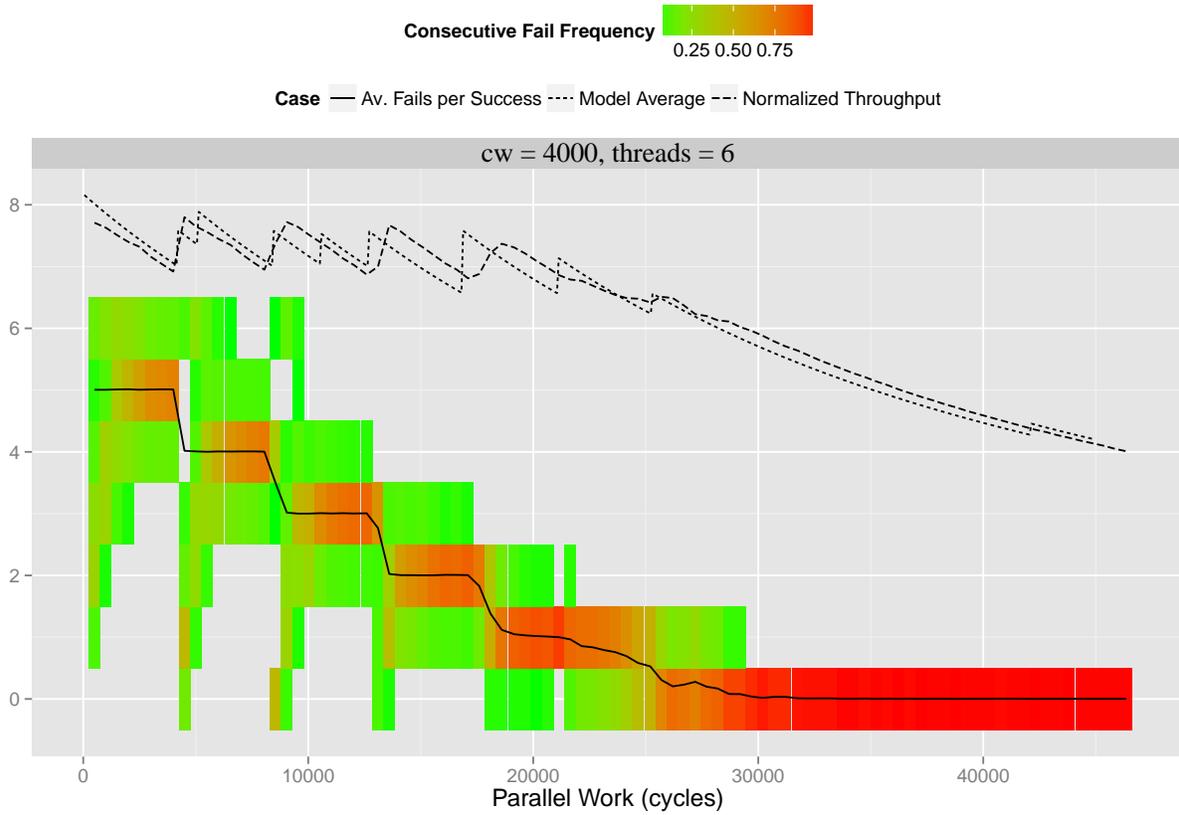
**Figure 16:** Consecutive Fails Frequency

a given number of consecutive fails, together with average fails per success values and the throughput values, normalized by a constant factor so that they can be seen on the graph. In the background, the frequency of occurrence of a given number of consecutive fails before success is presented. As a remark, the frequency of 6+ fails is gathered with 6. We expect to see a frequency distribution concentrated around the average fails per success value, within the bounds computed by our model.

While comparing the distribution of failures with the throughput, we could conjecture that the bumps come from the fact that the failures spread out. However, our model captures correctly the throughput variations and thus strips down the right impacting factor. The spread of the distribution of failures indicates the violation of a stable cyclic execution (that takes place in our model), but in these regions, $r$ actually gets close to 0, as well as the minimum of all gaps. The scattering in failures shows that, during the execution, a thread is overtaken by another one. Still, as gaps are close to 0, the imaginary execution, in which we switch the two thread IDs, would create almost the same performance effect. This reasoning is strengthened by the fact that the actual average number of failures follows the step behavior, predicted by our model. This shows that even when the real execution is not cyclic and the distribution of failures is not concentrated, our model that results in a cyclic execution remains a close approximation of the actual execution.

## H. Back-Off Tuning

Together with the analysis comes a natural back-off strategy: we estimate the $pw$ corresponding to the peak point of the average curve, and when the parallel section is smaller than the corresponding
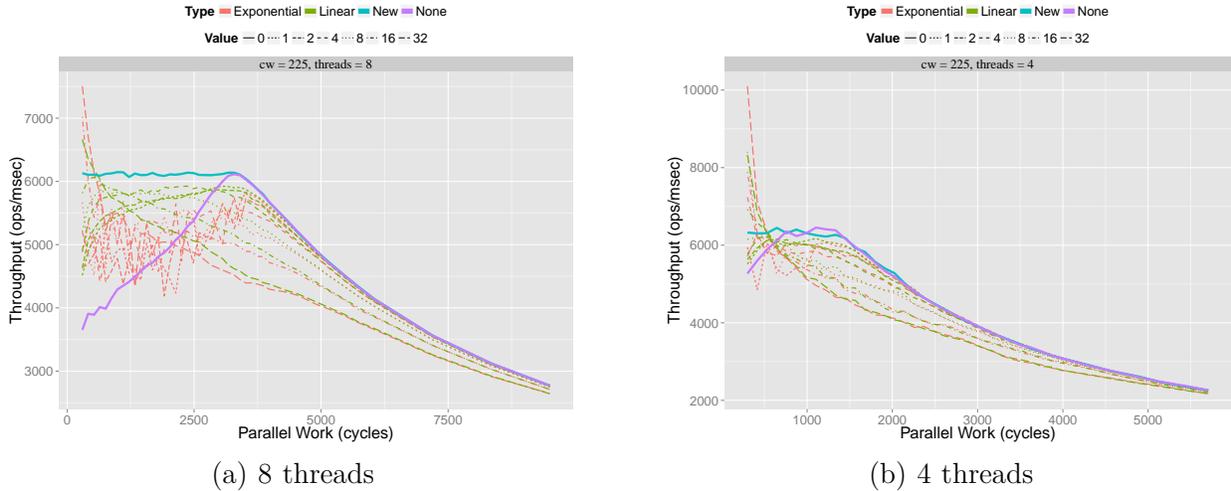
(a) 8 threads

(b) 4 threads

**Figure 17:** Comparison of back-off schemes for Poisson Distribution

$pw$, we add a back-off in the parallel section, so that the new parallel section is at the peak point.

We have applied exponential, linear and our back-off strategy to the Enqueue/Dequeue experiment specified above. Our back-off estimate provides good results for both types of distribution. In Figure 17 (where the values of back-off are steps of 115 cycles), the comparison is plotted for the Poisson distribution, which is likely to be the worst for our back-off. Our back-off strategy is better than the other, except for very small parallel sections, but other back-off strategies should be tuned for each value of $pw$.

We obtained the same shapes while removing the distribution law and considering constant values. The results are illustrated in Figure 18.
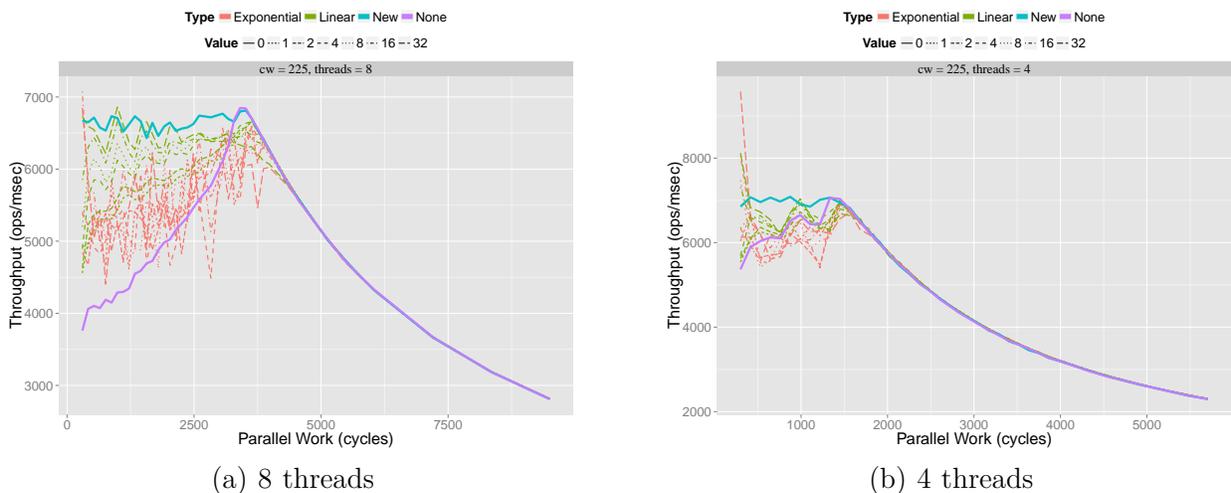


(a) 8 threads

(b) 4 threads

**Figure 18:** Comparison of back-off schemes for constant $pw$

## VII. Conclusion

In this paper, we have modeled and analyzed the performance of a general class of lock-free algorithms, and have so been able to predict the throughput of such algorithms, on actual executions. The analysis rely on the estimation of two impacting factors that lower the throughput: on the one hand, the expansion, due to the serialization of the atomic primitives that take place in the retry loops; on the other hand, the wasted retries, due to a non-optimal synchronization between the running threads. We have derived methods to calculate those parameters, along with the final throughput estimate, that is calculated from a combination of these two previous parameters. As a side result of our work, this accurate prediction enables the design of a back-off technique that performs better than other well-known techniques, namely linear and exponential back-offs.

As a future work, we envision to enlarge the domain of validity of the model, in order to cope with data structures whose operations do not have constant retry loop, as well as the framework, so that it includes more various access patterns. The fact that our results extend outside the model allows us to be optimistic on the identification of the right impacting factors. Finally, we also foresee studying back-off techniques that would combine a back-off in the parallel section (for lower contention) and in the retry loops (for higher robustness).

## References

[ACS14]   Dan Alistarh, Keren Censor-Hillel, and Nir Shavit. Are lock-free concurrent algorithms practically wait-free? In David B. Shmoys, editor, *Symposium on Theory of Computing (STOC)*, pages 714–723. ACM, June 2014.

[AF92]    Juan Alemany and Edward W. Felten. Performance issues in non-blocking synchronization on shared-memory multiprocessors. In Norman C. Hutchinson, editor, *Proceedings of the ACM Symposium on Principles of Distributed Computing (PoDC)*, pages 125–134. ACM, 1992.

[ARJ97]   James H. Anderson, Srikanth Ramamurthy, and Kevin Jeffay. Real-time computing with lock-free shared objects. *ACM Transactions on Computer Systems (TOCS)*, 15(2):134–165, 1997.

[DB08]    Kristijan Dragicevic and Daniel Bauer. A survey of concurrent priority queue algorithms. In *Proceedings of the International Parallel and Distributed Processing Symposium (IPDPS)*, pages 1–6, April 2008.

[DGT13]   Tudor David, Rachid Guerraoui, and Vasileios Trigonakis. Everything you always wanted to know about synchronization but were afraid to ask. In Michael Kaminsky and Mike Dahlin, editors, *Proceedings of the ACM Symposium on Operating Systems Principles (SOSP)*, pages 33–48. ACM, November 2013.

[DLM13]   Dave Dice, Yossi Lev, and Mark Moir. Scalable statistics counters. In Guy E. Blelloch and Berthold Vöcking, editors, *Proceedings of the ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*, pages 43–52. ACM, July 2013.

[GH09]    James R. Goodman and Herbert Hing Jing Hum. Mesif: A two-hop cache coherency protocol for point-to-point interconnects. Technical report, University of Auckland, November 2009.

[HSY10]   Danny Hendler, Nir Shavit, and Lena Yerushalmi. A scalable lock-free stack algorithm. *Journal of Parallel and Distributed Computing (JPDC)*, 70(1):1–12, 2010.

[Int13]   Intel. Lock scaling analysis on Intel® Xeon® processors. Technical Report 328878-001, Intel, April 2013.

[KH14]    Alex Kogan and Maurice Herlihy. The future(s) of shared data structures. In Magnús M. Halldórsson and Shlomi Dolev, editors, *Proceedings of the ACM Symposium on Principles of Distributed Computing (PoDC)*, pages 30–39. ACM, July 2014.

[LJ13]    Jonatan Lindén and Bengt Jonsson. A skiplist-based concurrent priority queue with minimal memory contention. In Roberto Baldoni, Nicolas Nisse, and Maarten van Steen, editors, *Proceedings of the International Conference on Principle of Distributed Systems (OPODIS)*, volume 8304 of *Lecture Notes in Computer Science*, pages 206–220. Springer, December 2013.

[Mic04]   Maged M. Michael. Hazard pointers: Safe memory reclamation for lock-free objects. *IEEE Transactions on Parallel and Distributed Systems (TPDS)*, 15(6):491–504, 2004.

[MS96]    Maged M. Michael and Michael L. Scott. Simple, fast, and practical non-blocking and blocking concurrent queue algorithms. In James E. Burns and Yoram Moses, editors, *Proceedings of the ACM Symposium on Principles of Distributed Computing (PoDC)*, pages 267–275. ACM, May 1996.

[SL00]    Nir Shavit and Itay Lotan. Skiplist-based concurrent priority queues. In *Proceedings of the International Parallel and Distributed Processing Symposium (IPDPS)*, pages 263–268, May 2000.

[Tre86]   R. Kent Treiber. *Systems programming: Coping with parallelism.* International Business Machines Incorporated, Thomas J. Watson Research Center, 1986.

[Val94]    J. D. Valois. Implementing Lock-Free Queues. In *Proceedings of International Conference on Parallel and Distributed Systems (ICPADS)*, pages 64–69, December 1994.