

# The Impact of Timing on Linearizability in Counting Networks

Marios Mavronicolas\*   Marina Papatriantafidou†   Philippas Tsigas‡

## Abstract

*Counting networks* form a new class of distributed, low-contention data structures, made up of interconnected *balancers* and are suitable for solving a variety of multiprocessor synchronization problems that can be expressed as counting problems. A *linearizable* counting network guarantees that the order of the values it returns respects the real-time order they were requested. Linearizability significantly raises the capabilities of the network, but at a possible price in network size or synchronization support. In this work we further pursue the systematic study of the impact of timing on linearizability for counting networks, along the line of research recently initiated by Lynch *et al.* in [18]. We consider two basic timing models, the instantaneous balancer model, in which the transition of a token from an input to an output port of a balancer is modeled as an instantaneous event, and the periodic balancer model, where balancers send out tokens at a fixed rate. We also consider lower and upper bounds on the delays incurred by wires connecting the balancers. We present necessary and sufficient conditions for linearizability in the form of precise inequalities that not only involve timing parameters, but also identify structural parameters of the counting network, which may be of more general interest. Our results extend and strengthen previous impossibility and possibility results on linearizability in counting networks [13, 18].

---

\*Department of Informatics, University of Cyprus, Nicosia, CY-1678, Cyprus. (mavronic@turing.cs.ucy.ac.cy)  
Partially supported by funds for the promotion of research at University of Cyprus (research project “Parallel, Concurrent and Distributed Computations”). Part of the work of this author was performed while visiting Max-Planck Institut für Informatik, Saarbrücken, Germany.

†Max-Planck-Institut für Informatik, Im Stadtwald, 66123 Saarbrücken, Germany (ptrianta@mpi-sb.mpg.de)

‡**Corresponding author:** Max-Planck-Institut für Informatik, Im Stadtwald, 66123 Saarbrücken, Germany. (tsigas@mpi-sb.mpg.de) *Tel.* (+49)-681-9325122, *Fax* (+49)-681-9325199

# 1 Introduction

In the *counting problem*, a number of concurrent processors repeatedly assign themselves successive values from a given range, such as memory addresses or destinations on an interconnection network. A solution is said to be *linearizable* [14] if the order of the assigned values reflects the real-time order in which they were requested. Linearizable counting provides the ground for a number of concurrent solutions to significant multiprocessor synchronization problems, such as time-stamp generation, multi-version database handling, scheduling of multi-threaded computations, implementation of data structures, dynamic load balancing, and buffer management (see, e.g., [7, 9, 11, 23]).

A *counting network* [3] is a highly concurrent data structure used to implement a counter. Roughly speaking, a counting network is a directed graph whose nodes are simple computing elements called *balancers*, and whose edges are called *wires*. A request for a counter value is represented by a *token*, which enters on one of the network’s input wires, propagates through the network asynchronously by traversing a sequence of balancers, and leaves on an output wire. Counting networks are among the very few counting techniques that are known to be scalable, since they minimize contention (“hot-spots”) as concurrency increases by distributing memory accesses, thus increasing parallelism and throughput (see, e.g., [3, 6, 12, 22, 21]).

In order to enhance the design of concurrent counting techniques, so that they both are scalable and support effective specification and analysis of MIMD shared memory algorithms—that rely on linearizability for correctness—it would be desirable to construct *linearizable counting networks*. Alternatively, it would be useful to study the possibility of using additional software constructs in order to extend a given counting network to become linearizable [13], or the conditions under which implemented counting networks always exhibit a linearizable behavior [18].

In this work we pursue the latter approach, following a direction pointed by a recent watershed paper [18], which studied linearizability properties of *uniform* counting networks relatively to timing assumptions on traffic speed. We further continue the systematic study of the impact of *timing* assumptions on linearizability for counting networks. More specifically, we study the boundaries between linearizable and non-linearizable behaviors of *any* counting network with respect to speed variations of its tokens and balancers, in a hope to provide practitioners with additional formal tools to support decision making in the phase of design. We provide necessary and sufficient conditions for a counting network to be linearizable, in the form of precise inequalities expressed in terms of specific graph-theoretic parameters and their relation to variations in traffic speed.

In more detail, we consider two basic timing models for balancer implementations in either shared memory or message passing. In both models, we follow Lynch *et al.* [18] and consider (non-zero) lower and upper bounds  $c_{min}$  and  $c_{max}$ , respectively, on the time it takes a token to traverse a wire from balancer to balancer. In the *instantaneous balancer* model, introduced and studied by Lynch *et al.* [18], the transition of a token from an input to an output port of a balancer is modeled as an instantaneous event. As pointed out also in [18], this is equivalent to the  $c_{min}$  and  $c_{max}$  bounds including the traversal of a node, but the output being determined at the instant of the token arrival. However, in some implementations of

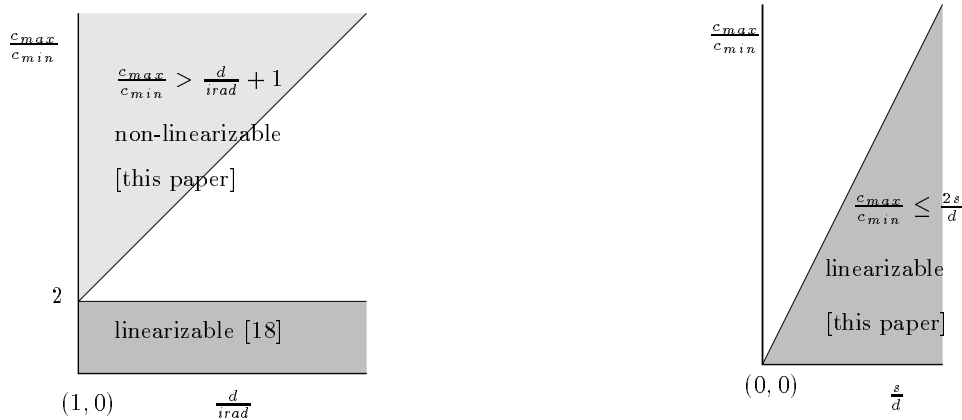


Figure 1: Pictorial summary of the results for the instantaneous balancer model

balancers there may be restrictions due to bandwidth or clock rates. In shared memory implementations of balancers memory accesses to variables implementing the balancers may require a constant number of steps to be completed due to restrictions in bandwidth, while in message passing implementations, processors that use messages to “simulate” the balancers may have access to clocks running at a bounded or fixed rate (see, e.g., [3, 12, 21, 22]). We model these “non-instantaneous” implementations by introducing a new timing model, called the *periodic balancer* model, assuming a *constant* period at which a balancer forwards tokens to its outputs. This assumption is motivated by periodic constraints commonly used in many real-time problems (especially in scheduling real-time tasks on multiprocessors [15, 17]), and resembles a timing model for periodic processes studied by Rhee and Welch [20]. The periodic balancer model is more realistic in that it models balancer delay to be proportional to the number of tokens concurrently traversing a balancer; this modeling is aligned with the concept of *stalls* introduced and used by Dwork *et al.* [6] in their elegant framework for analyzing contention in counting networks. We use  $r_{min}$  and  $r_{max}$  to denote the minimum and maximum balancer’s periods, respectively, over all balancers in the network.

We study, in particular, uniform and non-uniform counting networks; in *uniform* networks [18], each node lies on some path from inputs to outputs, and all paths from inputs to outputs have constant lengths. Our study introduces and identifies two crucial graph-theoretic parameters of a counting network, the first one being called its *influence radius* and denoted  $irad$ ; roughly speaking, the influence radius is defined to be the length of the maximum disjoint part between two maximal paths from an inner node of the network to any two output nodes, and captures the maximum degree of influence two output nodes can receive in common. The second parameter that our study identifies and uses is the *shallowness* of the network, which is the length of the shortest directed path from any of its inputs to any of its outputs (i.e. the opposite of the *depth* of the network). It turns out that the influence radius or the shallowness of the network together with its depth determine in a precise, quantitative way whether a counting network is linearizable under various timings. Our specific results and their relation to previous work are as follows.

**Necessary conditions** For the instantaneous balancer model, we show that any uniform

counting network of depth  $d$  is *not* linearizable if  $c_{max}/c_{min} > 1 + d/irad$ . Besides, by applying in any case of uniform counting networks, this result, for the case of diffracting trees [22], where  $irad = d = \log n$ , constitutes an alternative proof of an impossibility result of Lynch *et al.* [18, Theorem 4.1].

For the limit case where  $c_{max}$  tends to infinity, corresponding to completely asynchronous traffic, our impossibility result shows that infinite depth (hence, infinite size) is a necessary condition for linearizability in uniform counting networks – this constitutes an alternative proof of an impossibility result of Herlihy *et al.* [13, Theorem 5.1].

Moreover, our previous necessary condition justifies a method proposed in [18, Corollary 3.12] for turning *any* counting network of depth  $d$  such that  $c_{max}/c_{min} = k$ , for any constant  $k > 2$ , into a linearizable network of depth  $O(d)$ ; the method prepends each of the network’s inputs with a simple path consisting of  $dk$  single-input single-output balancers. This can be seen as a mechanism for safeguarding the appropriate ratio  $d/irad$ .

We next turn to the periodic balancer model, for which we show the equivalent necessary condition for linearizability, which depends on the influence radius of the network and the product of fan-outs of balancers along a crucial path from inputs to outputs. The proof of this condition follows the same structure as the one for the instantaneous balancer model, but, because of the more delicate timing assumptions in the periodic balancer model, it requires substantially more careful timing arguments.

**Sufficient condition** We show that *any* counting network (i.e. uniform or not) of depth  $d$  and shallowness  $s$  is linearizable if  $c_{max}/c_{min} \leq 2s/d$ ; this essentially says that the less “equilateral” the network is, the smaller are the variations in token speeds under which it can retain linearizability. This result is an important generalization and extension of the result by Lynch *et al.* which was tight for uniform networks, to the non-uniform cases; i.e. specialized in the case of uniform counting networks, where  $s = d$ , it yields that  $c_{max}/c_{min} \leq 2$  is sufficient for linearizability — the condition shown in [18].

Our results, which for the instantaneous balancer model are depicted in Figure 1, agree well with and provide a complement to known results [13, 18] on linearizability properties of counting networks, while they extend them significantly, in the ways explained in the previous paragraphs. Essentially, they imply that given a (uniform) counting network, we can determine traffic classes for which the network is linearizable or not by simply computing its depth and its shallowness or its influence radius, respectively. More important, our necessary and sufficient conditions together imply that, in general, linearizability may not be dictated by local conditions [18, Sections 3 and 4], but, rather, by conditions which need to involve graph-theoretic parameters describing the structure of the network. We remark that our impossibility results are shown using very simple, lock step and round-robin executions, which are expected to be common in practice. Besides, our proof techniques may be of independent interest.

We proceed by formally defining counting networks and the model of computation in Section 2. Sections 3 and 4 contain our necessary and sufficient conditions, respectively, for linearizable counting networks. We conclude in Section 4, with a discussion of our results and some open problems.

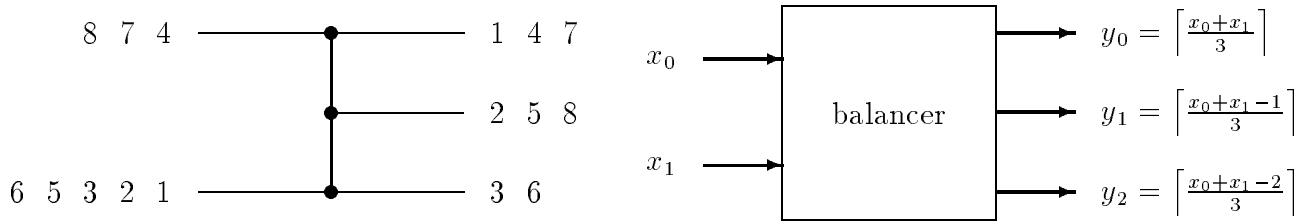


Figure 2: A (2,3)-balancer: symbolic and node representations

## 2 Preliminaries

Our definitions for balancing networks are standard and follow those in [1, 2, 3, 8, 10, 13, 18]. For self-containment, we define in this section, all terms used throughout the paper.

A  $(f_{in}, f_{out})$ -balancer, or *balancer* for short, is a computing element receiving tokens on  $f_{in}$  input wires, and sending out tokens to  $f_{out}$  output wires (see Figure 2);  $f_{in}$  and  $f_{out}$  are called the *fan-in* and *fan-out* of the balancer, respectively. Tokens arrive on the balancer’s input wires at arbitrary times and are output on its output wires. Intuitively, a balancer resembles a toggle mechanism which, given a stream of input tokens, alternately forwards them to its output wires, from top to bottom; thus, a balancer effectively balances the number of tokens that have been output on its output wires. We denote by  $x_i$ ,  $i \in [f_{in}]$ , the number of tokens ever received on the  $i$ -th input wire of a balancer, and by  $y_j$ ,  $j \in [f_{out}]$ , the number of tokens ever sent on its  $j$ -th output wire. (Throughout the paper, we will often abuse notation and use  $x_i$  (resp.,  $y_j$ ) as both the name of the  $i$ -th input wire (resp., output wire) and the number of tokens received (resp., sent) on the wire.) The *state* of a balancer at a given time is the collection of tokens on its input and output wires partitioned per input or output wire. In the initial state all wires contain no tokens. (For clarity we assume that all tokens are distinct.) A state of a  $(f_{in}, f_{out})$ -balancer is *quiescent* if  $\sum_{i=0}^{f_{in}-1} x_i = \sum_{j=0}^{f_{out}-1} y_j$ ; that is, the number of tokens that entered the balancer is equal to the number of tokens that left it. The following formal properties are required for a  $(f_{in}, f_{out})$ -balancer.

1. *Safety property*: In any state,  $\sum_{i=0}^{f_{in}-1} x_i \geq \sum_{j=0}^{f_{out}-1} y_j$  (a balancer never creates tokens).
2. *Liveness property*: Given any finite number of tokens  $m = \sum_{i=0}^{f_{in}-1} x_i$  that have been input to a  $(f_{in}, f_{out})$ -balancer, the balancer reaches within a finite amount of time a *quiescent state*, i.e.  $\sum_{i=0}^{f_{in}-1} x_i = \sum_{j=0}^{f_{out}-1} y_j = m$  (a balancer never “swallows” tokens).
3. *Step property*: In any quiescent state,  $0 \leq y_i - y_j \leq 1$  for any pair  $i$  and  $j$  such that  $0 \leq i < j \leq f_{out} - 1$ ; that is, the output has the step property.

A  $(w_{in}, w_{out})$ -balancing network is a collection of interconnected balancers; such a network is associated with a directed graph  $G$ , with three kinds of nodes:  $w_{in}$  source nodes,  $x_0, \dots, x_{w_{in}-1}$ , and  $w_{out}$  sink nodes,  $y_0, \dots, y_{w_{out}-1}$  which represent the input and output wires of the network, and a finite number of *inner* nodes, which represent the balancers of the network. The edges of the graph are the wires of the network’s balancers, the  $f_{in}$  incoming and  $f_{out}$  outgoing edges of a

node being the corresponding  $(f_{in}, f_{out})$ -balancer's input and output wires, respectively; the sink and source nodes have degree 1. We denote by  $\tilde{G}$  the non-directed version of  $G$ . Throughout the paper we consider acyclic networks. The *size* of a balancing network is the total number of its balancers. For any wire  $z$  in a balancing network, its *depth*, denoted  $depth(z)$ , is defined to be zero if  $z$  is an input wire of the network and  $\max_{i \in [f_{in}]} depth(z_i) + 1$  if  $z$  is the output wire of a balancer with input wires  $z_0, z_1, \dots, z_{f_{in}-1}$ . For any balancer  $b$  in a balancing network, its *depth*, denoted  $depth(b)$ , is the maximal wire depth over all of its input wires. The *depth*  $d$  of a balancing network is the maximal depth over all of its balancers. Each maximal set of balancers having the same depth  $l$  is called the *level*  $l$  of the network.

A balancing network is *uniform* [18] if each node of the network lies on some path from inputs to outputs, and all paths from inputs to outputs have the same length. Note that for any balancer  $b$  in a uniform balancing network,  $depth(b) = dist(x, b)$  for any source node  $x$  in  $G$ . The *configuration* of a balancing network at a given time is defined as the tuple of states of its component balancers at that time. For a configuration  $\sigma$ , denote by  $state_\sigma(b)$  the state of balancer  $b$  in  $\sigma$ . A configuration is *initial* if all component states are initial. A configuration of a balancing network is *quiescent* if  $\sum_{i=0}^{w_{in}-1} x_i = \sum_{j=0}^{w_{out}-1} y_j$ ; that is, the number of tokens that entered the network is equal to the number of tokens that left it. The safety and liveness properties of a balancing network follow naturally from its definition and the safety and liveness properties of a balancer.

A  $(w_{in}, w_{out})$ -*counting network* is a  $(w_{in}, w_{out})$ -balancing network for which, in any quiescent state,  $0 \leq y_j - y_k \leq 1$ , for any pair of indices  $j$  and  $k$  such that  $0 \leq j < k \leq w - 1$ ; that is, the output of a counting network has the *step property*. Each one of the  $w_{out}$  outputs of a balancing network is connected to an atomic counter (sink node), identified with the name of the respective wire. The tokens exiting the network through wire  $y_j$ , are consecutively assigned the integers  $j, j + w, j + 3w, \dots$ . The integer assigned to a token  $T$  by a counter is called the *returned value*, or *value* for short, of the token, and is denoted by  $val(T)$ .

We briefly describe below our model of multiprocessor computation, following [4, 14, 19]. We model computations of the system as sequences of (atomic) *events*, or *events* for short. Each event is either a *balancer transition event*, denoted by  $trans\langle T, b \rangle$ , representing transition of token  $T$  through balancer  $b$ , or a *wire transition event*, denoted by  $trans\langle T, b_1, b_2 \rangle$  representing transition of token  $T$  through a wire connecting an output of  $b_1$  to an input of  $b_2$ . An *execution*  $\mathcal{E}$  of a balancing network is a (possibly infinite) sequence of alternating configurations and events  $\sigma_0, e_1, \sigma_1, e_2, \sigma_2, \dots$ , where  $\sigma_0$  is the initial configuration.

Let  $c_{min}$  and  $c_{max}$ , such that  $0 < c_{min} \leq c_{max} < \infty$ , be the minimum and maximum time, respectively, that it takes for a token to traverse a link of the network. In the *asynchronous* model, the ratio  $c_{max}/c_{min}$  is unbounded, while in the completely synchronous it equals 1. A balancer  $b$  passes input tokens to its outputs at a constant rate, denoted by  $r(b)$ , ( $0 \leq r(b) < \infty$ ). Denote  $r_{min} = \min r(b)$  and  $r_{max} = \max r(b)$ , where the minimum and maximum are taken over all balancers. In the *instantaneous balancer* model [18] all  $r(b) = 0$ , and the upper and lower bounds of traversing a node may be included in the bounds  $c_{min}, c_{max}$ , the output of the node being determined at the instant of the token arrival. In the *periodic balancer* model all  $r(b) > 0$ , i.e. there are no bounds on the time that it takes for a token to traverse a node.

A *timed event* is a pair  $(t, e)$ , where  $t$ , the “time”, is a nonnegative real number, and  $e$  is an event. A *timed sequence* is an infinite sequence of alternating configurations and timed events  $\sigma_0, (t_1, e_1), \sigma_1, (t_2, e_2), \sigma_2, \dots, (t_j, e_j), \sigma_j, \dots$ , where the times are nondecreasing and unbounded. A timed sequence  $\mathcal{E}$  is a *timed execution* provided that  $\sigma_0, e_1, \sigma_1, \dots, e_j, \sigma_j, \dots$  is an execution and the following all hold:

1. (*Balancer transition time*)
  - (a) if the  $j$ th event is  $(t_j, \text{trans}\langle T, b_1, b_2 \rangle)$  and there are no input tokens in  $\text{state}_{\sigma_{j-1}}(b_2)$ , then there exists a  $k > j$  with  $t_k = t_j + r(b_2)$  such that the  $k$ th event is  $(t_k, \text{trans}\langle T, b_2 \rangle)$ ;
  - (b) if the  $j$ th event is  $(t_j, \text{trans}\langle T, b \rangle)$  and there exists an input token  $T'$  in  $\text{state}_{\sigma_{j-1}}(b)$ , then there exists a  $k > j$  with  $t_k = t_j + r(b)$  such that the  $k$ th event is  $(t_k, \text{trans}\langle T', b \rangle)$ ;
2. (*Wire transition time*)
  - (a) (Lower bound) if the  $j$ th event is  $(t_j, \text{trans}\langle T, b_1 \rangle)$ , then there is no  $k > j$  with  $t_k < t_j + c_{\min}$  such that the  $k$ th event is  $(t_k, \text{trans}\langle T, b_1, b_2 \rangle)$  for any balancer  $b_2$ ;
  - (b) (Upper bound) if the  $j$ th event is  $(t_j, \text{trans}\langle T, b_1 \rangle)$ , then there exists a  $k > j$  with  $t_k \leq t_j + c_{\max}$  such that the  $k$ th event is  $(t_k, \text{trans}\langle T, b_1, b_2 \rangle)$  for some balancer  $b_2$ .

Finally, the original definition of *linearizability* proposed by Herlihy and Wing [14] is adopted to counting networks in the natural way (cf. [13]). Given a timed execution  $\mathcal{E}$ , for any token  $T$ , define  $t_{\text{in}}(T, \mathcal{E})$  to be the least  $t$  such that  $(t, \text{trans}\langle T, b_1, b_2 \rangle)$  is an event in  $\mathcal{E}$ , and  $t_{\text{out}}(T, \mathcal{E})$  to be the greatest  $t$  such that  $(t, \text{trans}\langle T, b \rangle)$  is an event in  $\mathcal{E}$ . We say that token  $T_1$  *precedes* token  $T_2$  in  $\mathcal{E}$  if  $t_{\text{in}}(T_1, \mathcal{E}) < t_{\text{out}}(T_2, \mathcal{E})$ ; we write  $T_1 \xrightarrow{\mathcal{E}} T_2$  to denote this precedence. A timed execution  $\mathcal{E}$  of a balancing network is *linearizable* if for every pair of tokens  $T_1$  and  $T_2$  such that  $T_1 \xrightarrow{\mathcal{E}} T_2$ , it holds that  $\text{val}(T_1, \mathcal{E}) < \text{val}(T_2, \mathcal{E})$ . A balancing network is *linearizable* if each of its timed executions is linearizable.

### 3 Necessary Conditions

Consider an arbitrary uniform counting network  $G$  and its undirected version  $\tilde{G}$ . For any pair of sink nodes  $y_j$  and  $y_k$ , where  $0 \leq j, k \leq w_{\text{out}} - 1$ , let  $\text{dist}(y_j, y_k)$  denote the length (measured as the number of edges) of a shortest path in  $\tilde{G}$  connecting  $y_j$  and  $y_k$  (there is at least one simple path connecting them, since  $\tilde{G}$  is connected); each such shortest path is called a *geodesic* between  $y_j$  and  $y_k$  and denoted by  $\gamma(y_j, y_k)$ . Notice that for uniform networks, the length of any geodesic is even; hence, there exists on each geodesic  $\gamma(y_j, y_k)$  a node  $v_\gamma$  such that  $\text{dist}(v_\gamma, y_j) = \text{dist}(v_\gamma, y_k) = \text{dist}(y_j, y_k)/2$ ; call  $v_\gamma$  a *closest common ancestor* of  $y_j$  and  $y_k$ . The *influence radius* between  $y_j$  and  $y_k$ , denoted  $\text{irad}(y_j, y_k)$ , is the distance between  $y_j$  (or  $y_k$ ) and a closest common ancestor of  $y_j$  and  $y_k$ ; notice that in uniform networks  $\text{irad}(y_j, y_k) = \text{dist}(y_j, y_k)/2$ . Notice also that if two nodes  $u$  and  $v$  are both closest common ancestors for the same pair of sink nodes  $\text{depth}(u) = \text{depth}(v)$ . The *influence radius* of the network, denoted  $\text{irad}$ , is the maximum influence radius among any pair of its sink nodes.

Sections 3.1 and 3.2 present necessary conditions for linearizability for uniform networks, in the instantaneous and periodic balancer models, respectively. in terms of timing and the

depth and influence radius of the network. In all proofs of necessary conditions, we use a timed execution in which tokens propagate through the network in lock step, and each token traverses any link after delay exactly  $c_{min}$  and each balancer propagates tokens at a rate  $1/r_{min}$ . It is called a *fast, synchronous* timed execution.

### 3.1 Instantaneous Balancer Model

We prove the following theorem:

**Theorem 3.1** *In the instantaneous balancer model, a linearizable, uniform counting network does not exist if*

$$\frac{c_{max}}{c_{min}} > \frac{d}{irad} + 1.$$

**Proof.** Assume, by way of contradiction, that a linearizable, uniform  $(w_{in}, w_{out})$ -counting network  $G$  exists while  $c_{max}/c_{min} > d/irad + 1$ .

The following is an informal outline of the proof. We start with a fast, synchronous timed execution of  $G$  in which two distinguished tokens exit  $G$  through two antipodal sink nodes. By “retiming”, we slow down the token receiving the least value, while maintaining the propagation of the other token through  $G$ , thus, the latter token receives the same value after “retiming”. The “retimed” timed execution is further “augmented” to include a sufficient number of tokens fed in the network after the “fast” tokens exit it, and performing a fast traversal. The assumption on the timing parameters implies that at least one of the additional tokens will bypass the “slow” token of the previous timed execution and attain the value it received before retiming. This value is smaller than that of one of the “fast” tokens, contradicting linearizability. We now present the details of the formal proof.

Consider an arbitrary pair of antipodal sink nodes  $y_j$  and  $y_k$ , where  $0 \leq j < k \leq w_{out}$ ; thus,  $dist(y_j, y_k) = 2irad$ . Let  $\mathcal{E}$  be a timed execution in which we feed the network with  $k + 1$  tokens  $T_0, \dots, T_k$ , each  $T_i$ ,  $0 \leq i \leq k$ , entering the network through source node  $x_i$  and executing in lockstep in the maximum speed  $(1/c_{min})$ . Consider the network after it reaches a quiescent state; the output will have the step property and henceforth  $y_i$  equals 1 for  $i = 0, \dots, k$  and 0 for all other  $i < w_{out}$ . Let  $T_\kappa$  and  $T_j$  be the tokens that exit the network from  $y_k$  and  $y_j$ , respectively. Tokens  $T_j$  and  $T_\kappa$  in  $\mathcal{E}$  traverse the network through the paths  $\pi(T_j, \mathcal{E}) = b_{j,0} \rightsquigarrow b_{j,1} \rightsquigarrow \dots \rightsquigarrow b_{j,d-1}$  and  $\pi(T_\kappa, \mathcal{E}) = b_{\kappa,0} \rightsquigarrow b_{\kappa,1} \rightsquigarrow \dots \rightsquigarrow b_{\kappa,d-1}$ , respectively.

Now, “perturb”  $\mathcal{E}$  to create another timed execution  $\mathcal{E}'$  of  $G$  in the following way: let the same  $k + 1$  tokens enter the network from the same input nodes as in  $\mathcal{E}$  and keep the same timing as in  $\mathcal{E}$  until  $T_j$  goes through balancer  $b_{j,d-irad}$ ; then slowdown only  $T_j$  to the minimum speed  $1/c_{max}$  until it exits the network. Now  $T_j$  and  $T_\kappa$  traverse the network through the paths  $\pi(T_j, \mathcal{E}') = b'_{j,0} \rightsquigarrow b'_{j,1} \rightsquigarrow \dots \rightsquigarrow b'_{j,d-1}$  and  $\pi(T_\kappa, \mathcal{E}') = b'_{\kappa,0} \rightsquigarrow b'_{\kappa,1} \rightsquigarrow \dots \rightsquigarrow b'_{\kappa,d-1}$ , respectively. Notice that by the construction of  $\mathcal{E}'$ ,  $b'_{j,l} = b_{j,l}$  for  $0 \leq l \leq d - irad$ .

**Lemma 3.2** *There is no directed path in  $G$  from  $b'_{j,d-irad+1}$  to any node of  $\pi(T_\kappa, \mathcal{E})$ .*

**Proof.** Since the traversal of  $b_{j,d-irad}$  by  $T_j$  is not retimed in  $\mathcal{E}'$ , it follows that  $b'_{j,d-irad+1} = b_{j,d-irad+1}$ . Clearly, there is no directed path from  $b_{j,d-irad+1}$  to any  $b_{\kappa,i}$  with  $0 \leq i \leq d - irad + 1$ .



The existence of a path from  $b_{j,d-irad+1}$  to some  $b_{\kappa,i}$  with  $i > d - irad + 1$  would imply that  $dist(y_k, y_j) < 2irad$ , a contradiction, since  $y_k, y_j$  are antipodal.  $\square$

**Lemma 3.3** *In  $\mathcal{E}'$  each balancer is visited by the same number of tokens as in  $\mathcal{E}$ .*

**Proof.** Since the total number of tokens entering the network in  $\mathcal{E}'$  is the same as in  $\mathcal{E}$ , the lemma holds for each balancer at level 0. By a simple inductive argument using the liveness property of the balancers in the network, it holds for all levels.  $\square$

**Lemma 3.4**  $\pi(T_\kappa, \mathcal{E}') = \pi(T_\kappa, \mathcal{E})$ .

**Proof.** We use induction on the number of levels. It holds that  $b_{\kappa,0} = b'_{\kappa,0}$ . Assume that for all  $l$ ,  $0 \leq l < d$ , it is  $b_{\kappa,l} = b'_{\kappa,l}$ ; from lemma 3.3 it is known that  $b'_{\kappa,l}$  in  $\mathcal{E}'$  is traversed by the same number of tokens as in  $\mathcal{E}$ ; none of these tokens is  $T_j$  (this follows from lemma 3.2 and the fact that until level  $d - irad$  the timing is the same in both executions). Moreover, from the construction of  $\mathcal{E}'$  all the tokens but  $T_j$  are not retimed in  $\mathcal{E}'$ . Hence,  $T_\kappa$  follows the same output port after traversing  $b_{\kappa,l}$  as in  $\mathcal{E}$ . This shows that  $T_\kappa$  will go to the same balancer of level  $l + 1$  as in  $\mathcal{E}$ , i.e.  $b'_{\kappa,l+1} = b_{\kappa,l+1}$ , q.e.d.  $\square$

**Lemma 3.5**  $r = val(T_j, \mathcal{E}') < k$ .

**Proof.** Since the network is counting, after a quiescence state is reached, the output must have the step property. Since only  $k + 1$  ( $k < w_{out}$ ) tokens enter the network and since, by lemma 3.4,  $T_\kappa$  exits from the  $k$ -th output node,  $T_j$  must exit from a node with a lower index.  $\square$

We now perturb  $\mathcal{E}'$  to obtain a timed execution  $\mathcal{E}''$  of  $G$  which is not linearizable. Let

$$F_{out} = \prod_{i=0}^{d-1} f_{out}(b'_{j,i})$$

i.e.  $F_{out}$  is the product of the fan-outs of all the balancers on  $\pi(T_j, \mathcal{E}')$ . Let a set of tokens  $\Upsilon$ , where  $|\Upsilon| = F_{out}$ , enter the network from source node  $x_{j \bmod w_{in}}$ , in time  $dc_{min} + \epsilon$  (recall that  $dc_{min}$  is the time that  $T_\kappa$  exits the network;  $\epsilon$  is an arbitrarily small constant) and propagate in the network in lockstep and at maximum speed ( $1/c_{min}$ ). We will prove that there is at least one token  $T_v \in \Upsilon$  that will follow the same path  $\pi(T_j, \mathcal{E}')$ , which  $T_j$  followed in  $\mathcal{E}'$ , and that will bypass  $T_j$  before  $T_j$  exits the network.

**Lemma 3.6** *For each  $l$ ,  $0 \leq l \leq d - 1$ , at least  $\prod_{i=l+1}^{d-1} f_{out}(b'_{j,i})$  tokens have exited balancer  $b'_{j,l}$  and have been forwarded to  $b'_{j,l+1}$  (or the equivalent sink node if  $l = d - 1$ ) by time  $dc_{min} + \epsilon + (l + 1)c_{min}$  in  $\mathcal{E}''$ .*

**Proof.** We show this by induction on the number of levels. All  $F_{out}$  tokens of  $\Upsilon$  go through  $b'_{j,0}$ . Since they execute in lockstep, by time  $dc_{min} + \epsilon + c_{min}$ , a fraction  $1/f_{out}(b'_{j,0})$  of them

(i.e.  $\prod_{i=1}^{d-1} f_{out}(b'_{j,i})$  tokens) have exited balancer  $b'_{j,0}$  and have been forwarded to  $b'_{j,1}$ , which proves the base case. Assume that the lemma holds for all  $l$ ,  $0 \leq l < d - 1$ . Then at least  $\prod_{i=l+1}^{d-1} f_{out}(b'_{j,i})$  tokens are through balancer  $b'_{j,l}$  by time  $d c_{min} + \epsilon + (l + 1) c_{min}$ , ready to enter  $b'_{j,l+1}$ . A fraction  $1/f_{out}(b'_{j,l})$  of these tokens are output on the same port of  $b'_{j,l+1}$  that  $T_j$  is output in  $\mathcal{E}'$  and are, hence, forwarded to  $b'_{j,l+2}$ . Since all delays in link traversals involving these tokens are equal to  $c_{min}$ , it follows that by time  $d c_{min} + \epsilon + (l + 2) c_{min}$  they will have been forwarded to  $b'_{j,l+2}$ , thus showing the lemma for  $l + 1$ , as well.  $\square$

Lemma 3.6 shows that by time  $2 d c_{min} + \epsilon$  at least one token in  $\Upsilon$  will exit the network from the sink node  $y_r$ , where  $T_j$  exits from in  $\mathcal{E}'$ . The tokens in  $\Upsilon$  enter the network immediately after  $T_\kappa$  exits and, hence  $T_\kappa$  precedes them in  $\mathcal{E}''$ . At the time that  $T_\kappa$  exits,  $T_j$  still has  $irad(1 - c_{min}/c_{max})$  links to traverse before exiting the network, for which it will need  $irad(c_{max} - c_{min})$  time units. The tokens in  $\Upsilon$  need  $d c_{min}$  time units to traverse the network. If  $d c_{min} < irad(c_{max} - c_{min})$  (which is equivalent to the condition assumed, that  $c_{max}/c_{min} > d/irad + 1$ ) then  $T_j$  will be bypassed by at least one  $T_v \in \Upsilon$ , which will be the first to exit from  $y_r$  and will thus get  $val(T_v, \mathcal{E}'') = r$ ; from lemma 3.5 it is known that  $r < k$ . Since  $T_\kappa$  has exited the network before the tokens in  $\Upsilon$  entered, it follows that it will again, as in  $\mathcal{E}'$ , exit from wire  $y_k$  and get  $val(T_\kappa, \mathcal{E}'') = k$ . But now  $T_\kappa \xrightarrow{\mathcal{E}''} T_v$ , hence the fact that  $r < k$  implies that the linearizability condition is violated.  $\square$

## 3.2 Periodic Balancer Model

We prove the following theorem:

**Theorem 3.7** *In the periodic balancer model, a linearizable, uniform counting network does not exist if there exists a path to an output node with fan-out product  $F_{out}$  such that*

$$irad r_{max} + irad c_{max} > [irad + d(2w_{out} + F_{out} - 1)] r_{min} + (irad + d) c_{min}$$

**Proof.** Assume, by way of contradiction, that a linearizable, uniform counting network  $G$  exists while the condition of the theorem holds. We construct a timed execution of  $G$  which is not linearizable. The structure of the proof is similar to the one of theorem 3.1 for the instantaneous balancer model, but, because of the more delicate timing assumptions in the periodic balancer model, it requires substantially more careful timing arguments.

Fix again any pair of antipodal sink nodes  $y_j$  and  $y_k$  (i.e.  $dist(y_j, y_k) = 2irad$ ), where  $0 \leq j < k \leq w_{out} - 1$ . Let  $\mathcal{E}$  be a fast, synchronous timed execution with  $k + 1$  tokens, where each  $T_l$ ,  $0 \leq l \leq k$ , enters the network through source node  $x_l$ , each balancer outputs one token per  $r_{min}$  time, and all links incur a delay of  $c_{min}$ . By the step property for counting networks, when  $G$  reaches a quiescent configuration,  $y_l = 1$  for  $0 \leq l \leq k$  and 0 otherwise. Let  $T_j$  and  $T_\kappa$  be the tokens that exit  $\mathcal{B}$  through the sink nodes  $y_j$  and  $y_k$ , respectively, so that  $val(T_j, \mathcal{E}) = j$  and  $val(T_\kappa, \mathcal{E}) = k$ . Assume that in  $\mathcal{E}$ ,  $T_j$  and  $T_\kappa$  traverse the network through the paths  $\pi(T_j, \mathcal{E}) = b_{j,0} \rightsquigarrow b_{j,1} \rightsquigarrow \dots \rightsquigarrow b_{j,d-1}$  and  $\pi(T_\kappa, \mathcal{E}) = b_{\kappa,0} \rightsquigarrow b_{\kappa,1} \rightsquigarrow \dots \rightsquigarrow b_{\kappa,d-1}$ , respectively.

Notice that  $T_\kappa$  encounters  $d$  balancers while traversing the network. Since  $\mathcal{E}$  involves  $k+1 \leq w_{out}$  tokens,  $T_\kappa$  will wait at most  $w_{out}r_{min}$  time to traverse each of the  $d$  balancers, and  $c_{min}$  time to traverse each of the  $d$  links to exit completely the network, so that

**Fact 3.1**  $t_{out}(T_\kappa, \mathcal{E}) \leq d w_{out} r_{min} + d c_{min}$

We “perturb”  $\mathcal{E}$  to obtain another timed execution  $\mathcal{E}'$  of  $G$  so that each event occurring in  $\mathcal{E}$  no later than event  $trans\langle T_j, b_{j,d-irad} \rangle$  is not retimed, while later events are retimed so that each balancer and link encountered by  $T_j$  in  $\mathcal{E}'$  outputs one token per  $r_{max}$  time, and incurs a  $c_{max}$  delay, respectively. Assume that in  $\mathcal{E}'$ ,  $T_j$  and  $T_\kappa$  traverse the network through the paths  $\pi(T_j, \mathcal{E}') = b'_{j,0} \rightsquigarrow b'_{j,1} \rightsquigarrow \dots \rightsquigarrow b'_{j,d-1}$  and  $\pi(T_\kappa, \mathcal{E}') = b'_{\kappa,0} \rightsquigarrow b'_{\kappa,1} \rightsquigarrow \dots \rightsquigarrow b'_{\kappa,d-1}$ , respectively.

Notice that, by the construction of  $\mathcal{E}'$ ,  $b'_{j,l} = b_{j,l}$  for  $0 \leq l \leq d-irad$ . Also, since traversal of  $b_{j,d-irad}$  by  $j$  in  $\mathcal{E}$  is not retimed in  $\mathcal{E}'$ ,  $j$  follows the same output wire after traversing  $b_{j,d-irad}$  in either  $\mathcal{E}$  or  $\mathcal{E}'$ ; hence,  $b'_{j,d-irad+1} = b_{j,d-irad+1}$ . Hence, since  $y_j$  and  $y_k$  are antipodal, it holds that  $dist(b'_{j,irad+1}, y_k) > irad - 1$ . Hence, it follows that:

**Claim 3.8** *There is no directed path in  $G$  from  $b'_{j,d-irad+1}$  to any node of  $\pi(T_\kappa, \mathcal{E})$ .*

We continue to show:

**Lemma 3.9**  $\pi(T_\kappa, \mathcal{E}') = \pi(T_\kappa, \mathcal{E})$

**Proof.** Since tokens propagate through  $G$  in lock step, and no event occurring earlier than the traversal of the level  $d-irad$  is retimed in  $\mathcal{E}'$ , it follows that for each  $l$ ,  $0 \leq l \leq d-irad$ ,  $b'_{\kappa,l} = b_{\kappa,l}$ . We proceed to show by induction on  $l$ , where  $d-irad+1 \leq l \leq d-1$ , that  $b_{\kappa,l} = b'_{\kappa,l}$ . For the base case where  $l = d-irad+1$ , notice that since traversal of  $b_{\kappa,d-irad}$  by  $T_\kappa$  in  $\mathcal{E}$  is not retimed in  $\mathcal{E}'$ ,  $T_\kappa$  follows the same output wire after traversing  $b_{\kappa,d-irad}$  in either  $\mathcal{E}$  or  $\mathcal{E}'$ ; hence,  $b'_{\kappa,d-irad+1} = b_{\kappa,d-irad+1}$ . Assume inductively that  $b'_{\kappa,l} = b_{\kappa,l}$ , where  $d-irad+1 \leq l < d-1$ , and consider the balancer  $b'_{\kappa,l+1}$ . Since the execution is lock step, all not retimed in  $\mathcal{E}'$  tokens will reach  $b'_{\kappa,l}$  as before; also, by Claim 3.8, no retimed token will reach  $b_{\kappa,l}$ . It follows that each token visiting  $b_{\kappa,l}$  in  $\mathcal{E}'$  will follow the same link out of  $b_{\kappa,l}$  as in  $\mathcal{E}$ ; in particular, this implies that  $b'_{\kappa,l+1} = b_{\kappa,l+1}$ , as needed.  $\square$

Since only  $k+1 \leq w_{out}$  tokens are involved in  $\mathcal{E}'$  and  $G$  is a counting network, the step property immediately implies:

**Claim 3.10**  $r = val(T_j, \mathcal{E}') < k$

We now “perturb”  $\mathcal{E}'$  to obtain a third timed execution  $\mathcal{E}''$  of  $G$  which is not linearizable. Let

$$F_{out} = \prod_{r=0}^{d-1} f_{out}(b'_{j,r})$$

that is,  $F_{out}$  is the product of fan-outs of balancers in  $\pi(T_j, \mathcal{E}')$ . In  $\mathcal{E}''$ , new tokens  $\tilde{T}_1, \tilde{T}_2, \dots, \tilde{T}_{F_{out}}$  enter the network  $\mathcal{B}$  through source node  $x_{j \bmod w_{in}}$  strictly after  $T_\kappa$  exits it in execution  $\mathcal{E}'$ .

Since  $t_{out}(T_\kappa, \mathcal{E}') \leq d w_{out} r_{min} + d c_{min}$ , which holds due to lemma 3.9 and the fact that  $T_\kappa$  is again fast in  $\mathcal{E}'$ , for any arbitrarily small constant  $\epsilon > 0$ , for each  $l$ ,  $1 \leq l \leq F_{out}$ ,

$$t_{in}(\tilde{T}_l, \mathcal{E}'') = t_{out}(T_\kappa, \mathcal{E}') + \epsilon \leq w_{out} r_{min} + d c_{min} + \epsilon$$

All additional tokens traverse the network in FCFS order, each balancer they encounter outputs one token per  $r_{min}$  time, and all link traversals involving them are equal to  $c_{min}$ . We prove:

**Lemma 3.11** *For each  $l$ ,  $0 \leq l \leq d - 1$ , at least  $\prod_{r=l+1}^{d-1} f_{out}(b'_{j,r})$  tokens have exited balancer  $b'_{j,l}$  in  $\mathcal{E}''$  and have been forwarded to balancer  $b'_{j,l+1}$  (or an atomic counter if  $l = d - 1$ ) by time*

$$d w_{out} r_{min} + d c_{min} + \epsilon + (l + 1) [(F_{out} + w_{out}) r_{min} + c_{min}]$$

**Proof.** By induction on  $l$ . For the base case, where  $l = 0$ , notice that by construction of  $\mathcal{E}''$   $\prod_{r=0}^{d-1} f_{out}(b'_{j,r}) (= F_{out})$  tokens enter balancer  $b'_{j,0}$  by time  $d w_{out} r_{min} + d c_{min} + \epsilon$ . A fraction  $1/f_{out}(b'_{j,0})$  of them (i.e.  $\prod_{r=1}^{d-1} f_{out}(b'_{j,r})$  tokens) will exit balancer  $b'_{j,0}$  and be forwarded to balancer  $b'_{j,1}$ . Since  $\mathcal{E}''$  involves  $F_{out} + k + 1 \leq F_{out} + w_{out}$  tokens in total, and  $b'_{j,0}$  outputs one token per time  $r_{min}$ , while each wire incurs  $c_{min}$  delay to each token in this execution, it will take at most  $(F_{out} + w_{out}) r_{min} + c_{min}$  time for the last of these tokens to exit  $b'_{j,0}$  and be forwarded to  $b'_{j,1}$ , which proves the base case.

Assume inductively that the claim holds for any integer  $l$ , where  $0 \leq l < d - 1$ , that is at least  $\prod_{r=l+1}^{d-1} f_{out}(b'_{j,r})$  tokens have exited balancer  $b'_{j,l}$  in  $\mathcal{E}''$  and have been forwarded to balancer  $b'_{j,l+1}$  by time  $d w_{out} r_{min} + d c_{min} + \epsilon + (l + 1) [(F_{out} + w_{out}) r_{min} + c_{min}]$ . A fraction  $1/f_{out}(b'_{j,l+1})$  of them (i.e.  $\prod_{r=l+2}^{d-1} f_{out}(b'_{j,r})$  tokens) will be forwarded from  $b'_{j,l+1}$  to  $b'_{j,l+2}$ . Since  $\mathcal{E}''$  involves  $F_{out} + k + 1 \leq F_{out} + w_{out}$  tokens in total, and  $b'_{j,l+1}$  outputs one token per time  $r_{min}$ , while each wire incurs  $c_{min}$  delay to each token, it will take at most  $(F_{out} + w_{out}) r_{min} + c_{min}$  time for the last of these tokens to exit  $b'_{j,l+1}$  and be forwarded to  $b'_{j,l+2}$ , which proves that the lemma holds for  $l + 1$ , as well.  $\square$

By Lemma 3.11, at least one token (say  $\tilde{T}$ ) has been forwarded to exit  $y_j$  by time

$$t_{out}(\tilde{T}, \mathcal{E}'') \leq 2 d c_{min} + (2 d w_{out} + F_{out}) r_{min} + \epsilon.$$

By construction of  $\mathcal{E}''$ ,  $T_j$  is “fast” in  $\mathcal{E}''$  till it passes  $b'_{j,d-irad}$ , and “slow” afterwards. Hence, even if  $T_j$  never waits at a balancer *en route* due to other tokens concurrently traversing the same balancer, it holds that

$$t_{out}(T_j, \mathcal{E}'') \geq (d - irad) (r_{min} + c_{min}) + irad (r_{max} + c_{max})$$

Hence, the difference  $t_{out}(T_j, \mathcal{E}'') - t_{out}(\tilde{T}, \mathcal{E}'')$  is greater than

$$irad r_{max} + irad c_{max} - [d (2 w_{out} + F_{out} - 1) + irad] r_{min} - (d + irad) c_{min}$$

which is positive as we assumed. Hence,  $\tilde{T}$  exits from  $y_j$  before  $T_j$  in  $\mathcal{E}''$  and receives  $val(\tilde{T}, \mathcal{E}'') = r$ . Since  $T_\kappa$  has exited the network before the additional  $F_{out}$  tokens entered, it follows that it will again, as in  $\mathcal{E}'$ , exit from wire  $y_k$  and get  $val(T_\kappa, \mathcal{E}'') = k$ . Since by claim 3.10  $k > r$ , linearizability implies that  $\tilde{T} \xrightarrow{\mathcal{E}''} T_\kappa$ , a contradiction.  $\square$

## 4 Sufficient Conditions

In this section, we present our sufficient condition for linearizability in counting networks, which is an important generalization and extension of the respective result in [18] — which was tight for uniform networks — for the case of non-uniform networks, too. Consider any arbitrary  $(w_{in}, w_{out})$ -counting network  $G$  (uniform or not), and let its *shallowness*,  $s = \min_{i,j} \text{dist}(x_i, y_j)$  (i.e. the “opposite” of its depth); that is,  $s$  is the length of the shortest directed path in  $G$ . We prove:

**Theorem 4.1** *In the instantaneous balancer model, a counting network is linearizable if*

$$\frac{c_{max}}{c_{min}} \leq \frac{2s}{d}.$$

**Proof.** Assume, by way of contradiction, that there exists a counting network  $G$  of depth  $d$  and shallowness  $s$  which is not linearizable for  $c_{max}/c_{min} \leq 2s/d$ . By definition of linearizability, there exists a timed execution  $\mathcal{E}$  of  $G$  such that for a pair of tokens  $T_\kappa$  and  $T_\nu$ ,  $T_\kappa \xrightarrow{\mathcal{E}} T_\nu$ , while  $val(T_\kappa, \mathcal{E}) > val(T_\nu, \mathcal{E})$ .

We start with some auxiliary definitions. Following [18], we associate with each balancer  $b$  and token  $T$  in a timed execution, auxiliary *history variables*  $\mathcal{H}_b(t)$  and  $\mathcal{H}_T(t)$ , respectively, which formalize the “knowledge” that  $b$  and  $T$  have, respectively, at time  $t$  about other tokens in the network. Formally, at time 0,  $\mathcal{H}_b(0) = \emptyset$  and  $\mathcal{H}_T(0) = T$ . Each time a token traverses a balancer, the knowledge of the two is combined; formally, if a token  $T$  traverses a balancer  $b$  at time  $t$ , then  $\mathcal{H}_b(t) = \mathcal{H}_T(t) = \mathcal{H}_b(t^-) \cup \mathcal{H}_T(t^-)$ , where  $t^-$  is the time of occurrence of the timed event immediately preceding this in the timed execution. For a token  $T$  traversing  $G$  through the path  $b_0 \rightsquigarrow b_2 \rightsquigarrow \dots \rightsquigarrow b_{d-1}$ , define the *history sequence numbers*  $S_T^0, \dots, S_T^{d-1}$ , where  $S_T^i$  is the number of tokens that have been through  $b_i$  by the time instant  $t_i$  in which it crosses  $b_i$  in the timed execution.

Let  $y_k$  be the output node through which  $T_\kappa$  exits  $G$  in  $\mathcal{E}$ . The proof of [18, Lemma 3.1] does not rely on uniformity; hence, it applies to non-uniform networks as well, to yield:

**Claim 4.2** *If  $T_\kappa$  is the  $a^{\text{th}}$  token to exit  $G$  through  $y_k$ , then,  $|\mathcal{H}_{T_\kappa}(t_{out}(T_\kappa, \mathcal{E}))| \geq w(a-1) + k + 1$ .*

We now “perturb”  $\mathcal{E}$  to obtain a timed execution  $\mathcal{E}'$  which contains only the tokens in  $\mathcal{H}_{T_\kappa}(t_{out}(T_\kappa, \mathcal{E}))$  following the same timing in traversing  $G$ . Since no events were retimed and we only removed tokens about which  $T_\kappa$  did not “know” in  $\mathcal{E}$ , token  $T_\kappa$  still exits  $G$  through  $y_k$  at time  $t_{out}(T_\kappa, \mathcal{E}') = t_{out}(T_\kappa, \mathcal{E})$  in  $\mathcal{E}'$ . Note also that since in  $\mathcal{E}$   $T_\nu$  enters  $G$  after  $T_\kappa$  has exited,  $T_\nu$  is not in  $\mathcal{H}_{T_\kappa}(t_{out}(T_\kappa, \mathcal{E}'))$ ; hence,  $T_\nu$  is not participating in  $\mathcal{E}'$ . Now consider the token  $T_\mu$  for which  $val(T_\mu, \mathcal{E}') = val(T_\nu, \mathcal{E})$ . By construction,  $t_{in}(T_\mu, \mathcal{E}') = t_{in}(T_\mu, \mathcal{E})$ . Since  $T_\kappa$  “knows” about  $T_\mu$  at the time it is exiting  $G$ , and this information must have been propagated through  $G$  from some input wire to  $y_k$ , it must have traversed at least  $s$  wires. Hence,  $T_\mu$  must have entered  $G$  by time  $s c_{min}$  before in the latest. Hence, it follows:

**Claim 4.3**  $t_{in}(T_\mu, \mathcal{E}) = t_{in}(T_\mu, \mathcal{E}') \leq t_{out}(T_\kappa, \mathcal{E}) - s c_{min} = t_{out}(T_\kappa, \mathcal{E}') - s c_{min}$

The rest of the proof shows that  $T_\mu$  must have been bypassed in  $\mathcal{E}$  by some faster token, which, in turn, has similarly been bypassed by some other faster token; repeating this argument yields that  $T_\mu$  was bypassed by  $T_v$ ; note that  $T_v$  returns the same value that  $T_\mu$  returns in  $\mathcal{E}'$ . Next, we use the assumption  $c_{max}/c_{min} \leq 2s/d$  to show that for this to be possible  $T_v$  must have entered  $G$  before  $T_\kappa$  exited it in  $\mathcal{E}$ ; this contradicts the assumption that  $T_\kappa \xrightarrow{\mathcal{E}} T_v$  and completes the proof. These arguments are more formally presented in the following lemmas.

Let  $\mathcal{V}$  be the set of balancers visited by tokens during  $\mathcal{E}'$ . Then,

**Lemma 4.4** *In  $\mathcal{E}$  each balancer  $b \in \mathcal{V}$  processes at least the same number of tokens that it processes in  $\mathcal{E}'$ .*

**Proof.** Since the total number of tokens entering the network in  $\mathcal{E}'$  is less than in  $\mathcal{E}$ , the lemma holds for each balancer at level 0. By a simple inductive argument using the safety property which holds for all the balancers in the network, it holds for all levels.  $\square$

Token  $T_\mu$  during  $\mathcal{E}'$  traverses the network through a sequence of balancers  $b_0, \dots, b_{l-1}$  along path  $\pi(T_\mu, \mathcal{E}')$ , by going through  $b_i$  at the time instant denoted by  $t_i^\mu$ . We say that a token  $T_{\mu_j}$  in  $\mathcal{E}$  *simulates steps* of  $T_\mu$  of  $\mathcal{E}'$  on balancers  $b_{j_1}, \dots, b_{j_x}$  of  $\pi(T_\mu, \mathcal{E}')$ , if  $T_{\mu_j}$  goes through these balancers in  $\mathcal{E}$  and its history sequence numbers corresponding to them equal the respective history sequence numbers of  $T_\mu$  in  $\mathcal{E}'$  (naturally,  $T_\mu$  in  $\mathcal{E}$  may simulate itself). From the previous lemma we conclude that there exist tokens  $T_{\mu_1}, \dots, T_{\mu_\varepsilon}$ , which in  $\mathcal{E}$  simulate consecutive steps of  $T_\mu$  of  $\mathcal{E}'$ . Note that  $T_{\mu_\varepsilon} = T_v$ . The following lemma is essential for the completeness of the proof of our theorem.

**Lemma 4.5** *For any  $T_{\mu_j}$ , which simulates in  $\mathcal{E}$  steps of  $T_\mu$  of  $\mathcal{E}'$  on balancers  $b_{i_j}, \dots, b_{i_x}$  of  $\pi(T_\mu, \mathcal{E}')$ , the time  $t_i^{\mu_j}$  that it goes through balancer  $b_i \in \{b_{i_j}, \dots, b_{i_x}\}$  is such that  $t_i^{\mu_j} \leq t_{in}(T_\mu, \mathcal{E}) + \text{depth}(b_{i_j}) c_{max}$ .*

**Proof.** By induction on the length of the path  $\pi(T_\mu, \mathcal{E}')$ . At the first wire  $T_\mu$  simulates itself. Tokens can not traverse links slower than  $1/c_{max}$  and a token that is simulating  $T_\mu$  must have bypassed a token that was already simulating  $T_\mu$ , which, must have bypassed before another token that was simulating  $T_\mu$ , and so on. This “chain” of bypasses ends with a token that is simulating  $T_\mu$  because it bypassed  $T_\mu$  itself when crossing some balancer on  $\pi(T_\mu, \mathcal{E}')$ .  $\square$

From the last lemma we have that:

$$\begin{aligned} t_{in}(T_v, \mathcal{E}) + s c_{min} &\leq t_{out}(T_v, \mathcal{E}) \leq t_{in}(T_\mu, \mathcal{E}) + \text{length}(\pi(T_\mu, \mathcal{E}')) \\ &\leq t_{out}(T_\mu, \mathcal{E}) \leq t_{in}(T_\mu, \mathcal{E}) + d c_{max} \end{aligned}$$

Combining the above with the inequality of claim 4.3 we have that:

$$t_{in}(T_v, \mathcal{E}) \leq t_{out}(T_\kappa, \mathcal{E}) - 2s c_{min} + d c_{max}.$$

For  $2s c_{min} - d c_{max} > 0$ , it follows that  $t_{in}(T_v, \mathcal{E}) < t_{out}(T_\kappa, \mathcal{E})$ , contradicting the assumption that  $T_\kappa \xrightarrow{\mathcal{E}} T_v$ .  $\square$

Theorem 4.1 essentially says that the less “equilateral” the network is, the smaller are the variations in token speeds under which it can retain linearizability; specialized in the case of uniform counting networks, where  $s = d$ , it yields that  $c_{max}/c_{min} \leq 2$  is sufficient for linearizability – a result shown in [18].

## Discussion

We presented necessary and sufficient conditions for linearizability in counting networks, under different timing assumptions on balancers and wires. Although we do not yet have a complete characterization of linearizability for the specific timing models we consider, our results demonstrate how the possibility of achieving linearizability depends on both timing parameters of the model and structural parameters of the network.

We remark that the proofs of our necessary conditions can be extended to apply to other classes of balancing networks too, suggesting that it is not the requirement for the step property that has been the main obstacle to implementing linearizable counting networks, but, rather, the requirement for linearizability.

Our work leaves open several interesting problems. Can the necessary conditions be extended to non-uniform counting networks? An obvious open problem is to prove a sufficient condition for linearizability in the periodic balancer model. It would also be interesting to understand how much non-linearizable a counting network may be in case linearizability is impossible; Lynch *et al.* [18, Theorem 4.4] take the first step in this direction by providing a lower bound on the *non-linearizability fraction* for the special case of the bitonic counting network. Our necessary conditions should yield similar results for *any* uniform network in the models we studied. Does our sufficient condition for linearizability hold also for counting networks required to handle both tokens and *antitokens* [21]? A wide avenue for further research includes formalizing other possible systems and timing variations, and studying the possibility of linearizability in those models.

## References

- [1] E. Aharonson and H. Attiya, “Counting Networks with Arbitrary Fan-Out,” *Distributed Computing*, Vol. 8, pp. 163–169, 1995. Preliminary version: *Proceedings of the 3rd Annual ACM–SIAM Symposium on Discrete Algorithms*, pp. 104–113, January 1992.
- [2] W. Aiello, R. Venkatesan and M. Yung, “Coins, Weights and Contention in Balancing Networks,” *Proceedings of the 13th Annual ACM Symposium on Principles of Distributed Computing*, pp. 193–205, August 1994.
- [3] J. Aspnes, M. Herlihy and N. Shavit, “Counting Networks,” *Journal of the ACM*, Vol. 41, No. 5, pp. 1020–1048, September 1994. Preliminary version: “Counting Networks and Multi-Processor Coordination,” *Proceedings of the 23rd Annual ACM Symposium on Theory of Computing*, pp. 348–358, 1991.
- [4] H. Attiya and M. Mavronicolas, “Efficiency of Semi-Synchronous versus Asynchronous Networks,” *Mathematical Systems Theory*, Vol. 27, No. 6, pp. 547–571, November/December 1994.
- [5] C. Busch and M. Mavronicolas, “A Combinatorial Treatment of Balancing Networks,” *Proceedings of the 13th Annual ACM Symposium on Principles of Distributed Computing*, pp. 206–215, August 1994.
- [6] C. Dwork, M. Herlihy and O. Waarts, “Contention in Shared Memory Algorithms,” *Proceedings of the 25th Annual ACM Symposium on Theory of Computing*, pp. 174–183, 1993.
- [7] C. S. Ellis and T. J. Olson, “Algorithms for Parallel Memory Allocation,” *Journal of Parallel Programming*, Vol. 17, No. 4, pp. 303–345, August 1988.
- [8] E. W. Felten, A. LaMarca and R. Ladner, “Building Counting Networks from Larger Balancers,” Technical Report 93-04-09, Department of Computer Science and Engineering, University of Washington, April 1993.
- [9] A. Gottlieb, B. D. Lubachevsky and L. Rudolph, “Basic Techniques for the Efficient Coordination of Very Large Numbers of Cooperating Sequential Processors,” *ACM Transactions on Programming Languages and Systems*, Vol. 5, No. 2, pp. 164–189, April 1983.
- [10] N. Hardavellas, D. Karakos and M. Mavronicolas, “Notes on Sorting and Counting Networks,” *Proceedings of the 7th International Workshop on Distributed Algorithms (WDAG-93)*, Lecture Notes in Computer Science, Vol. # 725 (A. Schiper, ed.), Springer-Verlag, pp. 234–248, Lausanne, Switzerland, September 1993.
- [11] M. Herlihy, “A Methodology for Implementing Highly Concurrent Data Structures,” *Proceedings of the 2nd Annual ACM Symposium on Principles and Practice of Parallel Programming*, pp. 197–206, March 1990.
- [12] M. Herlihy, B.-C. Lim and N. Shavit, “Low Contention Load Balancing on Large-Scale Multiprocessors,” *Proceedings of the 4th Annual ACM Symposium on Parallel Algorithms and Architectures*, pp. 219–227, July 1992.



- [13] M. Herlihy, N. Shavit and O. Waarts, “Linearizable Counting Networks,” *Distributed Computing*, to appear, 1996. Preliminary version: “Low Contention Linearizable Counting Networks,” *Proceedings of the 32nd Annual IEEE Symposium on Foundations of Computer Science*, pp. 526–535, October 1991.
- [14] M. Herlihy and J. Wing, “Linearizability: A Correctness Condition for Concurrent Objects,” *ACM Transactions on Programming Languages and Systems*, Vol. 12, No. 3, pp. 463–492, July 1990.
- [15] K. Jeffay, D. F. Stanat and C. U. Martel, “On Optimal, Non-Preemptive Scheduling of Periodic and Sporadic Tasks,” *Proceedings of the 12th IEEE Real-Time Systems Symposium*, pp. 129–139, December 1991.
- [16] M. Klugerman and C. Plaxton, “Small-Depth Counting Networks,” *Proceedings of the 24th Annual ACM Symposium on Theory of Computing*, pp. 417–428, May 1992.
- [17] C. L. Liu and J. W. Layland, “Scheduling Algorithms for Multiprogramming in a Hard Real-Time Environment,” *Journal of the ACM*, Vol. 20, No. 1, pp. 46–61, January 1973.
- [18] N. Lynch, N. Shavit, A. Shvartsman and D. Touitou, “Counting Networks are Practically Linearizable,” *Proceedings of the 15th Annual ACM Symposium on Principles of Distributed Computing*, pp. 280–289, May 1996.
- [19] N. Lynch and M. Tuttle, “An Introduction to Input/Output Automata,” *CWI Quarterly*, Vol. 2, No. 3, pp. 219–246, September 1989.
- [20] I. Rhee and J. L. Welch, “The Impact of Time on the Session Problem,” *Proceedings of the 11th Annual ACM Symposium on Principles of Distributed Computing*, pp. 191–202, August 1992.
- [21] N. Shavit and D. Touitou, “Elimination Trees and the Construction of Pools and Stacks,” Preliminary version: *Proceedings of the 7th Annual ACM Symposium on Parallel Algorithms and Architectures*, pp. 54–63, July 1995.
- [22] N. Shavit and A. Zemach, “Diffracting Trees,” *Proceedings of the 6th Annual ACM Symposium on Parallel Algorithms and Architectures*, pp. 167–176, June 1994.
- [23] H. S. Stone, “Database Applications of the Fetch-and-Add Instruction,” *IEEE Transactions on Computers*, Vol. C-33, No. 7, pp. 604–612, July 1984.