

# Response Time Properties of Some Asynchronous Circuits

Jo Ebergen\* and Robert Berks  
Department of Computer Science  
University of Waterloo  
Waterloo, Ontario, Canada N2L 3G1  
{jeborgen,rtberks}@maveric.uwaterloo.ca

6 Jan, 1997

## Abstract

We discuss response time properties of linear arrays of cells with various handshake communication behaviours. The response times of a linear array are the delays between requests and succeeding acknowledgments for the first cell. We derive simple formulas for the worst-case response time and amortized response time of linear arrays using a general variable-delay model, where delays may vary between a lower and upper bound. The properties are independent of any particular implementation of the cells of the network.

## 1 Introduction

We present some properties on the response time of some regular asynchronous networks. The basic problems that we address can be illustrated with the following example.

Consider a linear array of  $L + 1$  cells as indicated in Figure 1 (where  $L = 3$ ). We are only

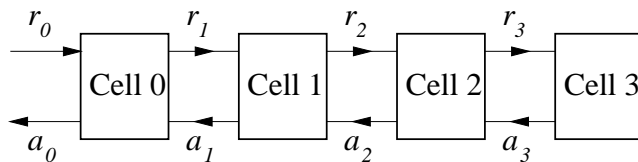


Figure 1: A linear array of cells

interested in the communication behaviour at the interface of each cell as an ordering of events. Each interface between two cells has the same communication behaviour, viz., a repetition of a request  $r$  followed by an acknowledgment  $a$ . A request followed by an acknowledgment is called a handshake. The end cell handshakes at its only interface, the other cells perform handshakes at

---

\*Currently on leave at Sun Microsystems Laboratories, Mountain View, USA.

both interfaces. Each non-end cell in this array has a communication behaviour that synchronizes the handshakes at both sides of the cell in some fashion. One of the simplest synchronizations is that for each handshake on the left there is exactly one handshake on the right.

In these communication behaviours, the requests and acknowledgments are considered as instantaneous events, and delays can be incurred in the cells or the environment of the pipeline. For each non-end cell we assume that the delays between the receipt of the *last* input and the production of any of the two outputs is always bounded from below by  $\delta$  and from above by  $\Delta$ . For the last cell we assume that its delays are bounded from below by  $\delta_L$  and from above by  $\Delta_L$ . The reason for distinguishing the last cell from the other cells is explained later. Finally, we assume that the environment can provide a next request  $r_0$  after an acknowledgment  $a_0$  after a minimum delay of  $\delta$ . There is no upper bound for the environment delays.

The delays of the cells do not have to be fixed, but may vary between a fixed upper and lower bound. Delays may depend on the data received, or they may vary over different instances of the same cell, or they may vary over time. Furthermore, the delay model does not depend on particular implementations of the basic cells. The only requirement is that the different implementations satisfy the upper bounds and lower bounds for the delays of the basic cells. If we take the upper and lower bound to be equal, as often happens in implementations using data bundling, we get the results for a fixed delay network.

Under these delay assumptions we would like to know what you can say about the worst-case response time of the linear network, that is, what is the maximum delay between any request  $r_0$  and succeeding acknowledgment  $a_0$  over all possible delay distributions? The next question then is whether the worst-case response time depends on the length of the array  $L$  or any of the bounds  $\delta$ ,  $\Delta$ ,  $\delta_L$ , and  $\Delta_L$ ? For example, we show that the worst-case response time for the linear network above is bounded from above by  $\Delta_L + (L + 1)(\Delta - \delta)$  and that this bound is tight. We also show under which delay distributions this bound can be achieved.

The second problem we consider is calculating an upper bound for the amortized response time of a network. The amortized response time gives an upper bound for the average response time of a network. The amortized response time can differ significantly from the worst-case response time. For example, we show that for the linear network above the amortized response time is bounded from above by  $\Delta_L + \Delta - \delta$ , which is independent of the length of the array  $L$ . This bound is also tight.

We can consider more general behaviours than the one considered above. For example, we can consider cells where for every  $n, n > 1$ , handshakes with the left neighbour there is at most one handshake with the right neighbour. What can we say about the worst-case response time and amortized response time in these cases? This is the third problem that we address. We show that the following properties hold for the worst-case response time  $WR$  and amortized response time  $AR$ .

- $WR \leq \Delta + (L + 1)(\Delta - \delta)$  provided that  $\Delta_L \leq 2\delta(n^{L-1} - 1) + \Delta$ .

Thus, even if the delay through the last cell increases exponentially with  $L$  (as in  $2\delta n^L$ ), the worst-case response time of the network increases at most linearly with  $L$  (as in  $L(\Delta - \delta)$ ).

- $AR \leq (3\Delta + \delta)/2$  provided that

$$\Delta_L \leq 2\delta(n^{L-1} - 1) + \Delta.$$

- $WR \leq 2\Delta$  provided that  $\Delta_L \leq \Delta n^{L-1} - (L - 1)(\Delta - \delta)$  and  $\Delta \leq 2\delta$ .

In the networks above we have parameterized the maximum and minimum delay for the last cell by  $\Delta_L$  and  $\delta_L$  respectively. Thus we can find out to what extent these bounds on the delay of the last cell can influence the response times of the network. Knowing the extent of this influence may give a designer more (or less) freedom in choosing an implementation for the end cell. For example, for linear arrays where  $n > 1$  applies, a bounded worst-case and amortized response time is maintained even if the maximum delay  $\Delta_L$  of the end cell increases exponentially with the length of the array. In such cases, a designer has the freedom to choose a few fast non-end cells and a very slow end cell so that we still achieve a bounded worst-case and bounded amortized response time. For other linear arrays where  $n = 1$ , the worst-case and amortized response time linearly depend on the maximum delay of the end cell. In these cases it may be important to choose a fast end cell.

The end cell is also a special cell for other reasons. It could be a part of an environment in which the non-end cells are placed. For example, for a simple FIFO the end cell could be seen as the “get” environment. In such cases, the formulas indicate to what extent the delays at the “get” environment can influence the response time at the “put” environment.

Symbol	Usage
$n$	Handshake multiplication factor. ( $n = 1$ is typical of micropipelines)
$\delta$	Lower bound on non-end cell delays and environment delays
$\Delta$	Upper bound on non-end cell delays
$\delta_L$	Lower bound on the end cell delay
$\Delta_L$	Upper bound on the end cell delay
$L$	Number of non-end cells and index of end cell
$WR$	worst-case response time
$AR$	amortized response time

Table 1: Summary of symbols used.

## 2 Related Work

The delay assumptions and dependency graphs (also called process graphs) in this paper are inspired by the model used in [5, 6]. In [5, 6] algorithms are described to obtain exact bounds for the delay between two given occurrences of events in a given dependency graph under a variable delay model. A much larger class of networks and communication behaviours is allowed than we use here. In contrast, we have focused on giving formula expressions for a small set of regular process graphs.

Such formulas have a few parameters, like the depth of the network and the upper and lower bounds for certain delays. These formulas allow for quick back-of-the-envelope calculations.

More detailed performance analysis and optimization techniques, all the way down to transistor sizing, are given in [1]. Techniques for analyzing the throughput and latency of micropipelines and rings are proposed in [14], where formulas are derived for the throughput as a function of the forward and backward latencies of the stages and the number of tokens in the ring. However, a fixed-delay model is used. See also [11], where these techniques have been applied. A fixed delay model is also used in [9], where an algorithm is presented to compute the average cycle time of arbitrary Signal Graphs. In [4] formulas for upper bounds of the utilization of pipelines are presented under various assumptions for the delay distributions for the latencies of the cells. In [10] yet another model is used for analyzing the response time of implementations. The model is based on the notion of sequence functions and is more restrictive than a variable delay model. Kearney and Bergmann [8] perform an analysis of a pipeline with variable delays and a multiplication factor of  $n = 1$ . Their results are largely based on simulation, and they do not give closed-form formulas.

### 3 Communication Behaviours

Consider the cell in Figure 2 from a linear array. At interface  $i$ , the communication behaviour for

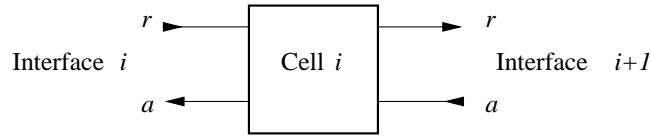


Figure 2: One cell

the handshakes can be given by  $*[r_i?; a_i!]$ . That is, we take an input request,  $r_i?$ , and then we output an acknowledgement  $a_i!$ . (? denotes input, while ! denotes output.) This behaviour then repeats (indicated by  $*[..]$ ). At interface  $i + 1$  the communication behaviour is given by  $*[r_{i+1}!; a_{i+1}?)$ .

Cell  $i$  synchronizes the handshakes at interface  $i$  and  $i + 1$ . In general, the communication behaviour for a cell can take many different forms. We restrict ourselves to behaviours for cells where for every  $n$  handshakes on interface  $i$  there is one handshake on interface  $i + 1$ . More precisely, the communication behaviours we look at are given by

$$(r_i?; *[ (a_i!; r_i?)^{n-1}; ((a_i!; r_i?) || (r_{i+1}!; a_{i+1}?) ) ])$$

In other words, only every  $n$ -th occurrence of the segment  $a_i!; r_i?$  takes place in parallel with a segment  $r_{i+1}!; a_{i+1}?$ .

For  $n = 1$  we have a communication behaviour that can be represented by

$$(r_i?; *[(a_i!; r_i?) || (r_{i+1}!; a_{i+1}?) ])$$

An unfolded process graph for  $n = 1$  for a cell is given in Figure 3.

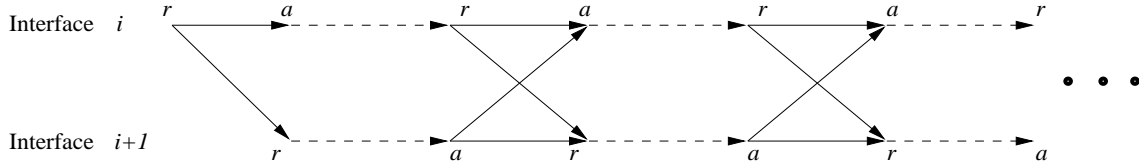


Figure 3: Unfolded communication behaviour for a cell

Let us denote the  $k$ -th occurrence of event  $a_i$  by  $(a_i, k)$ , where  $k \geq 0$ . A process graph gives a precedence relation for occurrences of events. A precedence between occurrences of events  $e_0$  and  $e_1$  is denoted by  $e_0 \rightarrow e_1$ . For example, in Figure 3 we have the precedences

$$(r_i, j) \rightarrow (a_i, j) \quad \text{and} \quad (a_{i+1}, j-1) \rightarrow (a_i, j)$$

These precedences indicate that occurrence  $j$  of event  $a_i$  can only happen after occurrence  $j$  of event  $r_i$  and occurrence  $j-1$  of event  $a_{i+1}$  has happened.

Solid arrows indicate precedence relations brought about by the cell. Dashed arrows indicate precedence relations brought about by the environment of the cell.

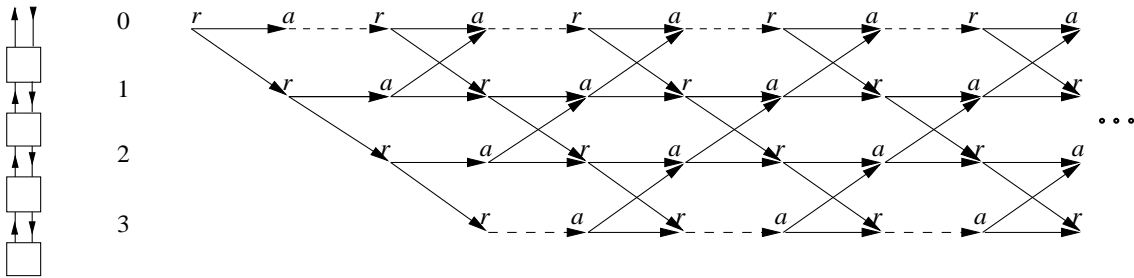


Figure 4: Unfolded communication behaviour for a linear array of four cells and  $n = 1$

If we consider the complete network of 3 non-end cells, one end cell, and the environment of the pipeline, then we can depict its behaviour in a process graph (Figure 4), where all edges represent direct precedences brought about by some component in the network. The dashed arrows represent direct precedences caused by the environment of the pipeline (top) or by the end cell (bottom).

For  $n = 2$  we have a communication behaviour that can be represented by

$$(r_i?; *[a_i!; r_i?; ((a_i!; r_i?) || (r_{i+1}!; a_{i+1}?!))])$$

An unfolded process graph for  $n = 2$  for a cell is given in Figure 5 and the corresponding graph for a linear array of four cells is given in Figure 6. For later use we have labelled the nodes in this graph with occurrence indexes.

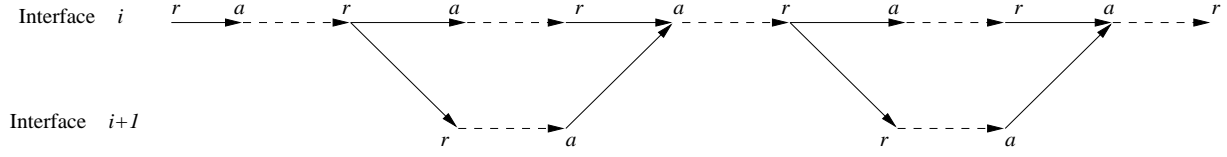


Figure 5: Unfolded communication behaviour for a cell with  $n = 2$

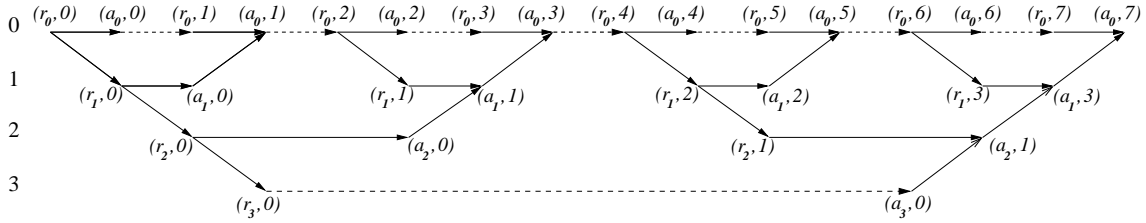


Figure 6: Part of an unfolded process graph for a linear array of four cells with  $n = 2$

## 4 The Delay Model

Once we have a process graph we can assign a time stamp to each node in the graph, where nodes are occurrences of events. This timing assignment,  $T$ , to nodes in the graph has to satisfy certain restrictions given by the delay model. The delay restrictions translate into lower and upper bounds for delays that can be incurred in each arc of the process graph and into the times that can be assigned to a node.

Let  $c$  be a node in a (unfolded) process graph with  $c$  an output of a non-end cell. Our delay model prescribes that any timing assignment  $T$  of the process graph must satisfy the following constraints.

$$T_{pred}(c) + \delta \leq T(c) \leq T_{pred}(c) + \Delta$$

where  $T(c)$  is the time assigned to node  $c$  and

$$T_{pred}(c) = \max\{T(b) \mid b \text{ is a direct predecessor of } c\}$$

This condition stipulates that the delay is always measured from the time  $T_{pred}(c)$ , when the *last* directly preceding event occurs, to the actual occurrence of the event  $c$ . The constraints prescribe that all delays through non-end cells have upper bound  $\Delta$  and lower bound  $\delta$ , that is

$$\delta \leq \text{Non - end cell delays} \leq \Delta$$

The delays incurred in the end cell can be expressed as  $T(a_L, j) - T(r_L, j)$  for occurrence  $j, j \geq 0$ , of output  $a_L$  and input  $r_L$ . We assume that the delays of the end cell have a lower bound  $\delta_L$  and upper bound  $\Delta_L$ , that is

$$\delta_L \leq \text{End cell delays} \leq \Delta_L$$

Furthermore we assume that  $\delta \leq \delta_L$  and  $\Delta \leq \Delta_L$ .

Finally, we assume that the delay through the environment of the pipeline is at least  $\delta$ , that is,  $T(r_0, j+1) - T(a_0, j) \geq \delta$  for all  $j \geq 0$ .

When a dependency graph and a timing assignment is given, we can define the useful notion of *critical paths*. We define an edge  $a \rightarrow b$  to be a *critical edge* if  $\delta \leq T(b) - T(a) \leq \Delta$  for  $b$  an output of a non-end cell, or if  $\delta_L \leq T(b) - T(a) \leq \Delta_L$  for  $b$  an output of an end cell. Note that the lower bound  $\delta$  is automatic by our statement that edge delays must be at least  $\delta$ . The upper bound may not always be true, because other dependencies may make event  $b$  occur later. We call a directed path a *critical path*, if all its edges are critical edges.

We can view the critical edges as pieces of elastic which can “stretch” between lengths  $\delta$  and  $\Delta$ , or  $\delta_L$  and  $\Delta_L$ . Note that for each node at least one incoming edge must be a critical edge, because of the delay constraints. If we stretch an edge beyond  $\Delta$  it breaks. On a critical path, we know that no bit of elastic on that path is broken, since, by the choice of the edges on the path, none of the bits of elastic is stretched beyond its upper bound  $\Delta$  or  $\Delta_L$ .

## 5 A Timing Property

Before we formulate some properties about response times, we define the concept of *response depth* of an occurrence of an event. Let the following path exist in a process graph for a behaviour  $B$

$$(a_i, j_i) \rightarrow (a_{i-1}, j_{i-1}) \rightarrow \dots (a_{k+1}, j_{k+1}) \rightarrow (a_k, j_k)$$

with  $i \geq k$  and occurrences  $j_i < j_{i-1} < j_{i-2} < \dots < j_{k+1} < j_k$ . That is, the acknowledgment  $a_k$  depends on  $a_i$  via a sequence of acknowledgments, one for each interface. Let furthermore  $T$  be a timing assignment for behaviour  $B$ . We say that  $(a_k, j_k)$  has *response depth*  $i - k$  for behaviour  $B$  and timing assignment  $T$  iff the path  $(a_i, j_i) \rightarrow (a_{i-1}, j_{i-1}) \rightarrow \dots (a_{k+1}, j_{k+1}) \rightarrow (a_k, j_k)$  is a critical path, and the edge  $(a_{i+1}, j_{i+1}) \rightarrow (a_i, j_i)$  does not exist, or, if it does exist, it is not critical. We call this critical path the “response path” for acknowledgment  $(a_k, j_k)$ . For example, if we have a response path with  $k = 0$ , and  $i = L$ , then we have a critical path of acknowledgments from the end-cell to the environment, with response depth  $L$ .

We now state the main theorem.

**Theorem 1** *If the response depth of  $(a_i, j_i)$  is  $k$ ,  $1 \leq k \leq L - i$ , then*

$$T(a_i, j_i) - T(a_i, j_i - n^{k-1}) \leq C_{i+k} + \Delta + k(\Delta - \delta)$$

where

$$C_k = \begin{cases} \Delta & \text{if } 0 \leq k < L \\ \Delta_L & \text{if } k = L \end{cases}$$

□

This theorem says that if acknowledgment  $(a_i, j_i)$  has a response depth of  $k$ ,  $1 \leq k \leq L - i$ , then the duration of the  $n^{k-1}$  handshake cycles at interface  $i$  that end in  $(a_i, j_i)$  is at most  $C_{i+k} + \Delta + k(\Delta - \delta)$ . The factor  $n^{k-1}$  in the above equation will prove useful, because we give a linear bound to the duration between two events that have an exponential number of events between them (if  $n > 1$ ).

Here is a brief example of the main theorem. Suppose  $n = 2$ ,  $k = 2$ ,  $L = 3$ , and that we are looking at event  $a_0$  ( $i = 0$ ) as illustrated in Figure 6. We will choose  $j_0 = 7$ . Therefore we wish to show that

$$\begin{aligned} T(a_0, 7) - T(a_0, 7 - 2^{(2-1)}) &\leq C_2 + \Delta + 2(\Delta - \delta) \\ T(a_0, 7) - T(a_0, 7 - 2) &\leq \Delta + \Delta + 2(\Delta - \delta) \\ T(a_0, 7) - T(a_0, 5) &\leq 4\Delta - 2\delta \end{aligned}$$

Since  $k = 2$ , we know that  $(a_0, 7)$  is dependent on an acknowledgment 2 levels deeper in the pipeline, viz.,  $(a_2, 1)$ , and not on acknowledgments any deeper. Therefore we can ignore the dependency from  $(a_3, 0)$ . So, tracing two steps further back,  $(a_0, 7)$  must be dependent on  $(r_1, 2)$ . The path from  $(r_1, 2)$  to  $(a_0, 7)$  via  $(a_2, 1)$  cannot be longer than  $4\Delta$ . The path from  $(r_1, 2)$  to  $(a_0, 5)$  cannot be shorter than  $2\delta$ . So we get

$$\begin{aligned} &T(a_0, 7) - T(a_0, 7 - 2^1) \\ &\leq (T(r_1, 2) + 4\Delta) - (T(r_1, 2) + 2\delta) \end{aligned}$$

which is equivalent to

$$T(a_0, 7) - T(a_0, 5) \leq 4\Delta - 2\delta$$

Before we prove the main theorem, we give some lemmas.

**Lemma 2** *If we have the direct dependency  $b \rightarrow c$ , then*

$$T(c) - T(b) \geq \delta \quad \text{and, hence } T(b) - T(c) \leq -\delta$$

*Proof.* The result follows directly from the definition of the timing assignment.  $\square$

Lemma 2 says that if we know of the existence of a dependency, then the separation between the source and target event must be at least  $\delta$ . The next lemma says that if an acknowledgment to a request to cell  $k$  depends only on cell  $k$  itself and not on cells  $k + 1, \dots, L$ , then the only delay incurred is the delay of cell  $k$ .

**Lemma 3** *For any  $k, 0 \leq k \leq L$ , if the response depth of  $(a_k, j_k)$  is 0, then*

$$T(a_k, j_k) - T(r_k, j_k) \leq C_k$$

*Proof.* In any behaviour  $B$  of a pipeline with the given restrictions, there are at most two direct dependencies that end in  $(a_k, j_k)$ , viz.

$$(r_k, j_k) \rightarrow (a_k, j_k) \quad \text{and possibly } (a_{k+1}, j_{k+1}) \rightarrow (a_k, j_k)$$

If the response depth of  $(a_k, j_k)$  is 0, then either the dependency  $(a_{k+1}, j_{k+1}) \rightarrow (a_k, j_k)$  does not exist or, if it does exist,

$$T(a_k, j_k) - T(a_{k+1}, j_{k+1}) > \Delta$$



This is equivalent to saying that the edge  $(a_{k+1}, j_{k+1}) \rightarrow (a_k, j_k)$ , if it exists, is not a critical edge. Since each node must have at least one incoming critical edge, the edge  $(r_k, j_k) \rightarrow (a_k, j_k)$  is a critical edge. Furthermore, we observe that the delay through any edge of the form  $(r_k, j_k) \rightarrow (a_k, j_k)$  is not an environment delay. Hence, by our definition of the timing model, the delay through edge  $(r_k, j_k) \rightarrow (a_k, j_k)$  must be bounded above by  $\Delta$  for a non-end cell and by  $D$  for an end cell.

□

With the preceding lemmas we can now prove the main theorem.

*Proof of Theorem 1.* By induction on  $k$ .

*Basis.* Assume that  $(a_i, j_i)$  has a response depth of  $k = 1$ . Let  $(a_{i+1}, j_{i+1}) \rightarrow (a_i, j_i)$ , then  $(a_{i+1}, j_{i+1})$  has response depth 0. From the restrictions on communication behaviours as explained in Section 3, it follows that we have a dependency graph for levels  $i$  and  $i+1$  as in Figure 7. From the

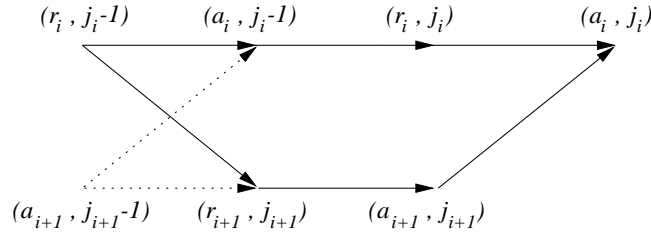


Figure 7: A part of the behaviour graph. Dotted edges are not present in every instance.

definition of  $T_{pred}(c)$  and the dependency graph in Figure 7, we have  $T_{pred}(r_{i+1}, j_{i+1}) = T_{pred}(a_i, j_i - 1)$ . In the proof below we take  $T_{pred} = T_{pred}(a_i, j_i - 1)$ . We wish to derive an upper bound for  $T(a_i, j_i) - T(a_i, j_i - 1)$ . We observe

$$\begin{aligned}
& T(a_i, j_i) - T(a_i, j_i - 1) \\
= & \quad \{ \text{Break path into 4 edges.} \} \\
& (T(a_i, j_i) - T(a_{i+1}, j_{i+1})) + (T(a_{i+1}, j_{i+1}) - T(r_{i+1}, j_{i+1})) + (T(r_{i+1}, j_{i+1}) - T_{pred}) + (T_{pred} - T(a_i, j_i - 1)) \\
\leq & \quad \{ (a_i, j_i) \text{ has positive response depth so timing dependency exists. Also, ack/ack pair cannot be environment or end-cell delays} \} \\
& \Delta + (T(a_{i+1}, j_{i+1}) - T(r_{i+1}, j_{i+1})) + (T(r_{i+1}, j_{i+1}) - T_{pred}) + (T_{pred} - T(a_i, j_i - 1)) \\
\leq & \quad \{ (a_{i+1}, j_{i+1}) \text{ has resp. depth 0, Lemma 3} \} \\
& \Delta + C_{i+1} + (T(r_{i+1}, j_{i+1}) - T_{pred}) + (T_{pred} - T(a_i, j_i - 1)) \\
= & \quad \{ T_{pred} = T_{pred}(a_i, j_i - 1) = T_{pred}(r_{i+1}, j_{i+1}) \} \\
& \Delta + C_{i+1} + (T(r_{i+1}, j_{i+1}) - T_{pred}(r_{i+1}, j_{i+1})) + (T_{pred} - T(a_i, j_i - 1)) \\
\leq & \quad \{ \delta \leq \text{cell delays} \leq \Delta, \text{ Lemma 2} \} \\
& \Delta + C_{i+1} + \Delta + (-\delta)
\end{aligned}$$

$$= \{ \text{calc.} \}$$

$$C_{i+1} + \Delta + (\Delta - \delta)$$

*Step.* Assume the theorem holds for  $k > 0$  and that  $(a_i, j_i)$  has response depth  $k + 1$ . Then there is a node  $(a_{i+1}, j_{i+1})$  such that dependency  $(a_{i+1}, j_{i+1}) \rightarrow (a_i, j_i)$  exists and  $(a_{i+1}, j_{i+1})$  has response depth  $k$ . From the induction hypothesis it then follows that node  $(a_{i+1}, j_{i+1} - n^{k-1})$  exists. Given that node  $(a_{i+1}, j_{i+1} - n^{k-1})$  and dependency  $(a_{i+1}, j_{i+1}) \rightarrow (a_i, j_i)$  exist, we show that node  $(a_i, j_i - n^k)$  and dependency  $(a_{i+1}, j_{i+1} - n^{k-1}) \rightarrow (a_i, j_i - n^k)$  also exist and that we can depict the dependencies for levels  $i$  and  $i + 1$  only as in Figure 8.

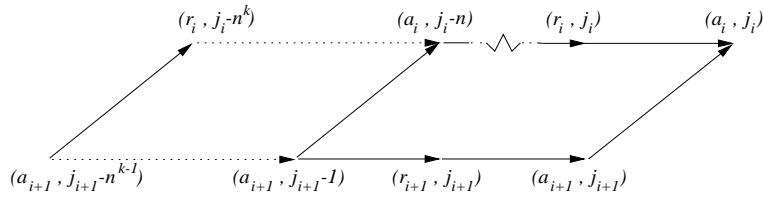


Figure 8: Dependency graph. The jagged line consists of  $2n - 1$  edges. Dotted lines represent paths.

Let us look at the behaviour of a cell at level  $i$ .

$$(r_i?; *[ (a_i!; r_i?)^{n-1}; ((a_i!; r_i?) || (r_{i+1}!; a_{i+1}?) ) ])$$

From this behaviour it follows that between any two successive edges  $a_{i+1} \rightarrow a_i$  there are  $n$  edges  $a_i \rightarrow r_i$ . It follows that dependency  $(a_{i+1}, j_{i+1} - 1) \rightarrow (a_i, j_i - n)$  exists. Repeating this step  $n^{k-1}$  times yields the dependency graph above with direct dependency  $(a_{i+1}, j_{i+1} - n^{k-1}) \rightarrow (a_i, j_i - n^k)$ .

Finally we observe

$$T(a_i, j_i) - T(a_i, j_i - n^k)$$

$$= \{ \text{Break path into three: edge, path, edge.} \}$$

$$((T(a_i, j_i)) - T(a_{i+1}, j_{i+1})) + (T(a_{i+1}, j_{i+1}) - T(a_{i+1}, j_{i+1} - n^{k-1})) + (T(a_{i+1}, j_{i+1} - n^{k-1}) - T(a_i, j_i - n^k))$$

$$\leq \{ (a_{i+1}, j_{i+1}) \text{ has response depth } k, \text{ induction hypothesis} \}$$

$$((T(a_i, j_i)) - T(a_{i+1}, j_{i+1})) + C_{i+k+1} + \Delta + k(\Delta - \delta) + (T(a_{i+1}, j_{i+1} - n^{k-1}) - T(a_i, j_i - n^k))$$

$$\leq \{ \text{from Lemma 2} \}$$

$$((T(a_i, j_i)) - T(a_{i+1}, j_{i+1})) + C_{i+k+1} + \Delta + k(\Delta - \delta) + (-\delta)$$

$$\leq \{ (a_i, j_i) \text{ has positive response depth} \}$$

$$\Delta + C_{i+k+1} + \Delta + k(\Delta - \delta) + (-\delta)$$

$$\leq \{ \text{calc.} \}$$

$$C_{i+k+1} + \Delta + (k + 1)(\Delta - \delta)$$

□

## 6 Worst-Case Response Time

From the last theorem we can infer a number of interesting properties. We can look at worst-case cycle time and worst-case response time. Worst-case cycle time is the largest value that  $T(a_0, j + 1) - T(a_0, j)$  can take over all possible values of  $j$  and all valid delay distributions. The worst-case cycle time is not that interesting in this delay model, since the environment of the pipeline can stretch the cycle time as much as possible by waiting to send a request. (The environment delays do not have an upper bound.) The worst-case response time  $WR$  is more interesting. We have the following theorem.

**Theorem 4** *The worst-case response time,  $WR$ , for any pipeline (as defined in section 3 operating under a delay model of section 4) is bounded from above as*

$$WR \leq \max(\{C_k + \Delta + k(\Delta - \delta) - 2\delta n^{k-1} + \delta \mid 1 \leq k \leq L\} \cup \{C_0\})$$

*Proof.* We have to find an upper bound for  $T(a_0, j) - T(r_0, j)$  over all  $j \geq 0$  and all valid delay distributions. Each acknowledgment has a response depth  $k$ , where  $0 \leq k \leq L$ . For a response depth  $k = 0$ , we have from Lemma 3,  $T(a_0, j) - T(r_0, j) \leq C_0$ . For a response depth of  $k > 0$ , we use the result of the previous section. We observe

$$\begin{aligned} & T(a_0, j) - T(r_0, j) \\ = & \quad \{ \text{calc.} \} \\ & (T(a_0, j) - T(a_0, j - n^{k-1})) + (T(a_0, j - n^{k-1}) - T(r_0, j)) \\ \leq & \quad \{ (a_0, j) \text{ has response depth } k, \text{ Theorem 1} \} \\ & C_k + \Delta + k(\Delta - \delta) + (T(a_0, j - n^{k-1}) - T(r_0, j)) \end{aligned}$$

Notice that there are  $n^{k-1}$  handshakes at level 0 between occurrences  $(a_0, j)$  and  $(a_0, j - n^{k-1})$ . The handshakes at level 0 experience delays through the environment and through cell 0. Assuming that the environment's and the cell's minimum response time is  $\delta$ , the minimum duration of one cycle through environment and cell is  $2\delta$ . Therefore the minimum duration from  $T(a_0, j - n^{k-1})$  to  $T(r_0, j)$  is  $2\delta n^{k-1} - \delta$ . This then leads to the inequality

$$T(a_0, j) - T(r_0, j) \leq C_k + \Delta + k(\Delta - \delta) - 2\delta n^{k-1} + \delta$$

Maximizing these upper bounds over all  $k, 0 \leq k \leq L$  gives the desired result.

□

The above theorem holds for all pipelines with a multiplication factor  $n > 0$ . The case when  $L = 0$  is not especially interesting, as this is simply an environment and an end cell. In this case the formula resolves to  $WR \leq \Delta_L$ . For  $L > 0$ , we consider two special cases:  $n = 1$  and  $n > 1$ . Instantiating the above theorem for  $n = 1$ , we derive

**Theorem 5** *The worst-case response time  $WR$  for a linear array with multiplication factor  $n = 1$  satisfies*

$$WR \leq \Delta_L + (L + 1)(\Delta - \delta)$$

*Proof.*

$$\begin{aligned}
& WR \\
& \leq \{ \text{Theorem 1 above, } n = 1 \} \\
& \quad \max(\{C_k + (k + 1)(\Delta - \delta) \mid 1 \leq k \leq L\} \cup \{C_0\}) \\
& = \{ \text{calc., } \Delta \geq \delta, \Delta_L \geq \Delta \} \\
& \quad \Delta_L + (L + 1)(\Delta - \delta)
\end{aligned}$$

□

In other words, for a pipeline with multiplication factor 1, the worst-case response time depends linearly on  $\Delta_L$ ,  $L$ , and  $\Delta - \delta$ . With the proper distribution of delays, this upper bound can indeed be attained. See Figure 9. All solid thin lines have a delay of  $\delta$ , and all solid thick lines have a delay of  $\Delta$ . The dashed, thick line has a delay of  $\Delta_L$ . The dashed, thin line has a delay  $\delta_L$ . The delay between  $(r_0, 4)$  and  $(a_0, 4)$  is  $\Delta_L + 4(\Delta - \delta)$ . Obviously, this example can be generalized for a pipeline of any length.

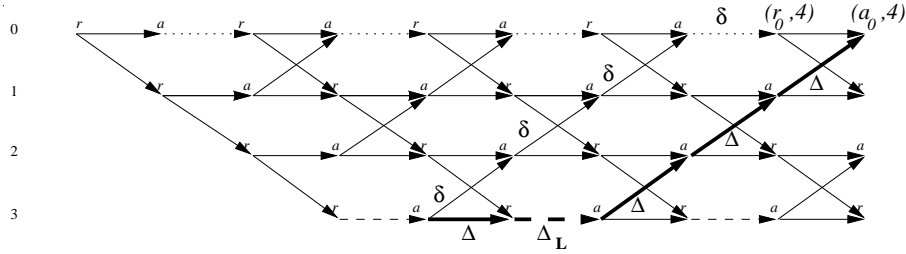


Figure 9: Example with worst-case response

For the case where the multiplication factor  $n > 1$ , we take a function  $D_L$  to serve as an upper bound for the delay through the end cell. Let

$$D_L = 2\delta(n^{L-1} - 1) + \Delta$$

for  $1 \leq L$ . If we assume that  $\Delta_L \leq D_L$ , we derive with Theorem 4

**Theorem 6** *The worst-case response time  $WR$  for a linear array with multiplication factor  $n \geq 1$  satisfies*

$$WR \leq \Delta + (L + 1)(\Delta - \delta)$$

*provided the delay through the end cell satisfies*

$$\Delta_L \leq 2\delta(n^{L-1} - 1) + \Delta$$

*Proof.* Let  $D_k = 2\delta(n^{k-1} - 1) + \Delta$  with  $1 \leq k \leq L$ . Then  $C_k \leq D_k, 1 \leq k \leq L$ . This also implies that  $D_L \leq 2\delta(n^{L-1} - 1) + \Delta$ .

$$\begin{aligned}
WR &\leq \max(\{C_k + \Delta + k(\Delta - \delta) - 2\delta n^{k-1} + \delta \mid 1 \leq k \leq L\} \cup \{C_0\}) \\
&= \max\{C_k + \Delta + k(\Delta - \delta) - 2\delta n^{k-1} + \delta \mid 1 \leq k \leq L\} \quad \{\text{Since } L > 0\} \\
&\leq \max\{D_k + \Delta + k(\Delta - \delta) - 2\delta n^{k-1} + \delta \mid 1 \leq k \leq L\} \quad \{\text{Since } C_k \leq D_k\} \\
&= \max\{2\delta(n^{k-1} - 1) + \Delta + \Delta + k(\Delta - \delta) - 2\delta n^{k-1} + \delta \mid 1 \leq k \leq L\} \\
&= \max\{-2\delta + 2\Delta + k(\Delta - \delta) + \delta \mid 1 \leq k \leq L\} \\
&= \max\{2\Delta - \delta + k(\Delta - \delta) \mid 1 \leq k \leq L\} \\
&= \max\{\Delta + (k+1)(\Delta - \delta) \mid 1 \leq k \leq L\} \\
&= \Delta + (L+1)(\Delta - \delta)
\end{aligned}$$

□

In words, if the delay in the last cell increases at most exponentially with  $L$  (as in  $2\delta n^L$ ), the worst-case response time of the network increases at most linearly with  $L$  (as in  $k(\Delta - \delta)$ ).

A special case occurs when the delays in the cells are fixed ( $\Delta = \delta$ ), which often happens in bundled data implementations where matching delays are fixed, then the upper bound for the worst-case response time is bounded by  $\Delta$ .

If the upper and lower bounds for the delays differ, we can still maintain a bounded response time for  $n > 1$  if we are willing to let the delay in the last cell increase a little bit less than  $D_L$ . How much less increase is needed depends on the ratio of  $\Delta$  and  $\delta$ . For example, we have

**Theorem 7** *If  $\Delta \leq 2\delta$  then the worst-case response time  $WR$  for a linear array with multiplication factor  $n > 1$  satisfies*

$$WR \leq 2\Delta$$

*provided the delay through the end cell satisfies*

$$\Delta_L \leq \Delta n^{L-1} - (L-1)(\Delta - \delta)$$

*Proof.* Let  $D_k = \Delta n^{k-1} - (k-1)(\Delta - \delta)$  with  $1 \leq k \leq L$ . Then  $C_k \leq D_k, 1 \leq k \leq L$ . This also implies that  $\Delta_L \leq \Delta n^{L-1} - (L-1)(\Delta - \delta)$ .

$$\begin{aligned}
R &\leq \max(\{C_k + \Delta + k(\Delta - \delta) - 2\delta n^{k-1} + \delta \mid 1 \leq k \leq L\} \cup \{C_0\}) \\
&= \max(\{C_k + \Delta + k(\Delta - \delta) - 2\delta n^{k-1} + \delta \mid 1 \leq k \leq L\}) \quad \{\text{Since } L > 0\} \\
&\leq \max(\{D_k + \Delta + k(\Delta - \delta) - 2\delta n^{k-1} + \delta \mid 1 \leq k \leq L\}) \quad \{\text{Since } C_k \leq D_k\} \\
&= \max(\{(\Delta \cdot n^{k-1} - (k-1) \cdot (\Delta - \delta)) + \Delta + k(\Delta - \delta) - 2\delta n^{k-1} + \delta \mid 1 \leq k \leq L\}) \\
&= \max(\{(\Delta \cdot n^{k-1} - k(\Delta - \delta) + (\Delta - \delta)) + \Delta + k(\Delta - \delta) - 2\delta n^{k-1} + \delta \mid 1 \leq k \leq L\}) \\
&= \max(\{(\Delta - 2\delta) \cdot n^{k-1} + (\Delta - \delta) + \Delta + \delta \mid 1 \leq k \leq L\}) \\
&= \max(\{(\Delta - 2\delta) \cdot n^{k-1} + 2\Delta \mid 1 \leq k \leq L\})
\end{aligned}$$

If  $\Delta \leq 2\delta$ , the expression  $(\Delta - 2\delta) \leq 0$ , and so the maximum occurs when  $n^{k-1}$  is minimized. For  $k \geq 1$ , the maximum occurs when  $k = 1$  with value

$$(\Delta - 2\delta) + 2\Delta$$

$$\begin{aligned}
&\leq (2\delta - 2\delta) + 2\Delta \quad \{\text{Since } \Delta \leq 2\delta \} \\
&= 2\Delta
\end{aligned}$$

□

## 7 Amortized Response Time

The amortized response time is the average response time under the worst-case delay distribution. In other words, the average response time for any delay distribution is at most the amortized response time. We show that the amortized response time is bounded by constants independent of  $L$  for pipelines with multiplication factor  $n > 1$ , as we would expect, but also for pipelines with  $n = 1$ . In other words, although the worst-case response time for a pipeline with  $n = 1$  depends linearly on  $L$ , the average response time for such a pipeline is bounded from above by a constant independent of  $L$ .

We consider two cases  $n = 1$  and  $n > 1$  and use the notion of a critical response path to prove an upper bound for the amortized response time for these pipelines. Here is the theorem for the case  $n = 1$ .

**Theorem 8** *The amortized response time  $AR$  of a pipeline with multiplication factor  $n = 1$  is bounded from above as*

$$AR \leq \Delta_L + \Delta - \delta$$

*Proof.* Consider a pipeline with multiplication factor  $n = 1$  and a timing assignment for its dependency graph as illustrated in Figure 4. Take an occurrence of an acknowledgment, say  $(a_0, j_0)$ . From acknowledgment  $(a_0, j_0)$  there is a critical path from some request  $(r_0, i_0)$  with  $i_0 \leq j_0$  and on this critical path there is no other event at level 0. Let the length of this path be  $h$ . Notice that if  $h = 1$ , then  $j_0 = i_0$  and we have  $T(a_0, j_0) - T(r_0, i_0) \leq \Delta$ . For a path of length  $h > 1$ , a delay of  $\Delta_L$  can be obtained in at most every other edge of the path and only in edges at level  $L$ . In the other edges a delay of at most  $\Delta$  can be obtained. Let  $h = 2k + 1$  for some  $k \geq 0$ . For  $k < L$  we then get an upper bound of

$$T(a_0, j_0) - T(r_0, i_0) \leq (2k + 1)\Delta$$

Furthermore, between  $(r_0, i_0)$  and  $(a_0, j_0)$  at level 0 there are  $k = (h - 1)/2$  handshake cycles plus one cell delay. (A handshake cycle consists of one environment and one cell delay.) Environment delays are at least  $\delta$ . In  $k$  handshake cycles, the total time for the environment delays is at least  $k\delta$ . Hence, the total time for the pipeline responses between  $T(r_0, i_0)$  and  $T(a_0, j_0)$  is at most  $(2k + 1)\Delta - k\delta$ . So the average response time  $AR_k$  for the pipeline responses between  $T(r_0, i_0)$  and  $T(a_0, j_0)$  satisfies

$$AR_k \leq ((2k + 1)\Delta - k\delta)/(k + 1) = 2\Delta - \delta - (\Delta - \delta)/(k + 1)$$

for  $0 \leq k < L$ . Maximizing this upper bound over all  $k$ ,  $0 \leq k < L$ , yields an upper bound for the amortized response time,  $AR_0$ .

$$AR_0 \leq 2\Delta - \delta - (\Delta - \delta)/L$$

For  $k \geq L$  we get a critical path of length  $2k + 1$  with maximum delay if there are  $2L$  edges with delay  $\Delta$ ,  $k - L$  edges with delay  $\Delta$ , and  $k + 1 - L$  edges with delay  $\Delta_L$ .

$$T(a_0, j_0) - T(r_0, i_0) \leq 2L\Delta + \Delta_L + (k - L)(\Delta_L + \Delta)$$

Using a similar reasoning as above, we find for the average response time,  $AR_k$ , for the pipeline responses between  $T(r_0, i_0)$  and  $T(a_0, j_0)$ .

$$\begin{aligned} & AR_k \\ \leq & \quad \{ \text{From above} \} \\ & (2L\Delta + \Delta_L + (k - L)(\Delta_L + \Delta) - k\delta)/(k + 1) \\ = & \quad \{ \text{calc.} \} \\ & \Delta_L + \Delta - \delta - (\Delta - \delta + L(\Delta_L - \Delta))/(k + 1) \end{aligned}$$

Maximizing this upper bound over all  $k \geq L$ , and taking into account that  $\Delta \geq \delta$  and  $\Delta_L \geq \Delta$ , yields an upper bound for the amortized response time,  $AR1$ .

$$AR1 \leq \Delta_L + \Delta - \delta$$

The maximum of  $AR0$  and  $AR1$  gives as upper bound for the amortized response time  $AR$ .

$$AR \leq \Delta_L + \Delta - \delta$$

□

This bound is tight, since we can get arbitrarily close to it with a proper delay distribution and a sufficiently long sequence of handshakes. For example, if we take a delay distribution, where each delay through a non-end cell is  $\Delta$ , each delay through an end cell is  $\Delta_L$ , and each delay through the environment is  $\delta$ , and we take an infinite number of handshakes at cell 0, then this bound is achieved.

For a multiplication factor  $n > 1$  we can use the following theorem. It states that if the delay through the last cell increases at most exponentially with  $L$  (as in  $2\delta n^L$ ), then the amortized response time for a pipeline with multiplication factor  $n > 1$  is bounded by  $(3\Delta + \delta)/2$ .

**Theorem 9** *The amortized response time  $AR$  of a pipeline with multiplication factor  $n > 1$  is bounded from above as*

$$AR \leq (3\Delta + \delta)/2$$

*provided the delay in the last cell satisfies*

$$\Delta_L \leq 2\delta(n^{L-1} - 1) + \Delta$$

*Proof.* A derivation for an upper bound of the amortized response time for  $n > 1$  is similar to the derivation for  $n = 1$ . Part of the dependency graph for when  $n = 2$  is shown in Figure 6. Let there be a critical path of length  $h$  from  $(a_0, j_0)$  back to  $(r_0, i_0)$  with  $i_0 \leq j_0$  and no other occurrence of an event on the path is at level 0. As before we have  $h = 2k + 1$ . Moreover, for a dependency

graph with  $n > 1$ , we infer that there are  $k$  “up-edges,”  $k$  “down-edges,” and one “level-edge” in the path. (See Figure 6 for  $n = 2$ , for example.) The total delay therefore for this path is at most

$$2k\Delta + \Delta_L$$

when  $k = L$ . For  $k < L$  the upper bound is

$$2k\Delta + \Delta$$

By definition of  $C_k$ , for  $0 \leq k \leq L$ , the upper bound can be expressed as

$$2k\Delta + C_k$$

The number of handshake cycles between  $(r_0, i_0)$  and  $(a_0, j_0)$  is

$$N_k = \left( \sum_{i=0}^{k-1} n^i \right) = (n^k - 1)/(n - 1)$$

where we define  $N_0 = 0$ . Let  $D_k = 2\delta(n^{k-1} - 1) + \Delta$  with  $1 \leq k \leq L$ . Then  $C_k \leq D_k$ ,  $1 \leq k \leq L$ . This also implies that  $\Delta_L \leq 2\delta(n^{L-1} - 1) + \Delta$ . Similar to the proof of the previous theorem we get the following derivation for the amortized response time  $AR$ .

$$\begin{aligned}
& AR \\
& \leq \{ \text{calc.} \} \\
& \quad \max\{(2k\Delta + C_k - N_k\delta)/(N_k + 1) \mid 0 \leq k \leq L\} \\
& = \{ L > 0, N_0 = 0, \text{ and } C_0 \leq \Delta + (C_1 - \delta)/2 \} \\
& \quad \max\{(2k\Delta + C_k - N_k\delta)/(N_k + 1) \mid 1 \leq k \leq L\} \\
& = \{ \text{calc.} \} \\
& \quad \max\{(2k\Delta + C_k + \delta)/(N_k + 1) - \delta \mid 1 \leq k \leq L\} \\
& \leq \{ C_k \leq D_k, 1 \leq k \leq L \} \\
& \quad \max\{(2k\Delta + D_k + \delta)/(N_k + 1) - \delta \mid 1 \leq k \leq L\} \\
& = \{ \text{Def. of } D_k \} \\
& \quad \max\{(2k\Delta + 2\delta(n^{k-1} - 1) + \Delta + \delta)/(N_k + 1) - \delta \mid 1 \leq k \leq L\} \\
& = \{ \text{calc.} \} \\
& \quad \max\{((2k + 1)\Delta - \delta) + 2\delta n^{k-1}/(N_k + 1) - \delta \mid 1 \leq k \leq L\} \\
& \leq \{ \text{Exponential term reaches maximum of } 2\delta \text{ when } k \rightarrow \infty \} \\
& \quad \max\{((2k + 1)\Delta - \delta)/(N_k + 1) + 2\delta - \delta \mid 1 \leq k \leq L\} \\
& \leq \{ \text{Remaining term is maximized when } k = 1 \} \\
& \quad (3\Delta - \delta)/2 + \delta \\
& = \{ \text{calc.} \}
\end{aligned}$$



$$(3\Delta + \delta)/2$$

This bound is not quite tight, because we maximize at  $k = 1$  for one term and at  $k \rightarrow \infty$  for another term.

□

## 8 Concluding Remarks

Given a linear pipeline, various communication behaviours, and bounds on the delays of the cells, we have given formulas for the amortized and worst-case response times of the pipeline. Notably, we have shown that we can achieve bounded amortized response times independent of the length of the pipeline. Furthermore, we have shown under what conditions a bounded worst-case response time can be achieved.

Not only do these formulas allow for quick back-of-the-envelope calculations, but they also allow for a better understanding of which parameters are influential and in what sense they influence the response time of the design. For example, we have shown how slow the end cell can be without affecting the response time. Such knowledge can be exploited conveniently in the design of the last cell by trading speed for better energy efficiency, for example.

There are many pipeline designs for various computations all with multiplication factor  $n = 1$  and where data is flowing in one or both directions, see [2, 3, 7, 10, 12, 13] for example. If we know what the values are for each of the parameters for certain implementations, we can quickly determine what the upper bounds for the worst-case and amortized response times are. If we know that delays are fixed, that is,  $\Delta = \delta$ , then we get a worst-case and amortized response time that are bounded by constants.

The results for the linear arrays and the detailed proofs form a convenient starting point for various generalizations. We mention just a few of the generalizations that are under investigation. For example, how do the results generalize if one permits more parallelism between handshakes on the left-hand side and on the right-hand side? How do the results generalize if one allows behaviours where handshakes on the right-hand side are conditional on communicated data values? And finally, how do the results generalize if we consider tree-like networks.

## 9 Acknowledgments

The second author would like to thank Ed Carr for his helpful comments involving the proofs.

## References

- [1] S. M. Burns and A. J. Martin. Performance analysis and optimization of asynchronous circuits. In C. H. Séquin, editor, *Advanced Research in VLSI: Proceedings of the 1991 UC Santa Cruz Conference*, pages 71–86. MIT Press, 1991.
- [2] J. Ebergen and R. Hoogerwoord. A derivation of a serial-parallel multiplier. *Science of Computer Programming*, 15:201–215, 1990.

- [3] J. C. Ebergen, J. Segers, and I. Benko. Parallel program and asynchronous circuit design. In G. Birtwistle and A. Davis, editors, *Proceedings Banff VIII Workshop: Asynchronous Digital Circuit Design*, Workshops in Computing. Springer-Verlag, 1995.
- [4] M. R. Greenstreet and K. Steiglitz. Bubbles can make self-timed pipelines fast. *Journal of VLSI Signal Processing*, 2(3):139–148, Nov. 1990.
- [5] H. Hulgaard, S. Burns, T. Amon, and G. Borriello. An algorithm for exact bounds on the time separation of events in concurrent systems. Technical Report TR-94-02-02, University of Washington, Department of Computer Science and Engineering, Feb. 1994.
- [6] H. Hulgaard and S. M. Burns. Bounded delay timing analysis of a class of CSP programs with choice. In *Proc. International Symposium on Advanced Research in Asynchronous Circuits and Systems*, pages 2–11, Nov. 1994.
- [7] M. B. Josephs, P. G. Lucassen, J. T. Udding, and T. Verhoeff. Formal design of an asynchronous DSP counterflow pipeline: A case study in handshake algebra. In *Proc. International Symposium on Advanced Research in Asynchronous Circuits and Systems*, pages 206–215, Nov. 1994.
- [8] D. Kearney and N. W. Bergman. Performance evaluation of asynchronous logic pipelines with data dependant processing delays. In *Second Working Conference on Asynchronous Design Methodologies*, pages 4–13. IEEE Computer Society Press, May 1995.
- [9] C. D. Nielsen and M. Kishinevski. Performance analysis based on timing simulation. In *Proc. ACM/IEEE Design Automation Conference*, pages 70–76, June 1994.
- [10] M. Rem. Trace theory and systolic computations. In J. W. de Bakker, A. J. Nijman, and P. C. Treleaven, editors, *PARLE: Parallel Architectures and Languages Europe, Vol. I*, volume 258 of *Lecture Notes in Computer Science*, pages 14–33. Springer-Verlag, 1987.
- [11] J. Sparsø and J. Staunstrup. Design and performance analysis of delay insensitive multi-ring structures. In *Proc. Hawaii International Conf. System Sciences*, volume I, pages 349–358. IEEE Computer Society Press, Jan. 1993.
- [12] I. E. Sutherland. Micropipelines. *Communications of the ACM*, 32(6):720–738, June 1989.
- [13] K. van Berkel and M. Rem. VLSI programming of asynchronous circuits for low power. In G. Birtwistle and A. Davis, editors, *Proceedings Banff VIII Workshop: Asynchronous Digital Circuit Design*, Workshops in Computing. Springer-Verlag, 1995.
- [14] T. E. Williams. Analyzing and improving the latency and throughput performance of self-timed pipelines and rings. In *Proc. International Symposium on Circuits and Systems*, May 1992.