

Clustering for Processing Rate Optimization*

Chuan Lin, Jia Wang, and Hai Zhou
Electrical and Computer Engineering
Northwestern University
Evanston, IL 60208

Abstract

Clustering (or partitioning) is a crucial step between logic synthesis and physical design in the layout of a large scale design. A design verified at the logic synthesis level may have timing closure problems at post-layout stages due to the emergence of multiple-clock-period interconnects. Consequently, a trade-off between clock frequency and throughput may be needed to meet the design requirements. In this paper, we find that the processing rate, defined as the product of frequency and throughput, of a sequential system is upper bounded by the reciprocal of its maximum cycle ratio, which is only dependent on the clustering. We formulate the problem of processing rate optimization as seeking an optimal clustering with the minimal maximum-cycle-ratio in a general graph, and present an iterative algorithm to solve it. Since our algorithm avoids binary search and is essentially incremental, it has the potential of being combined with other optimization techniques. Experimental results validate the efficiency of our algorithm.

1 Introduction

Circuit clustering (or partitioning) is often employed between logic synthesis and physical design to decompose a large circuit into parts. Each part will be implemented as a separate cluster that satisfies certain design constraints, such as the size of a cluster. Clustering helps to provide the first order information about interconnect delays as it classifies interconnects into two categories: intra-cluster ones are local interconnects due to their spatial proximity while inter-cluster ones may become global interconnects after floor-plan/placement and routing (also known as circuit layout).

Due to aggressive technology scaling and increasing operating frequencies, interconnect delay has become the main performance limiting factor in large scale designs. Industry data shows that even with interconnect optimization techniques such as buffer insertion, the delay of a global interconnect may still be longer than one clock period, and multiple clock periods are generally required to communicate such a global signal. Since global interconnects are not visible at logic synthesis when the functionality of the implementation is the major concern, a design that is correct at the logic synthesis level may have timing closure problems after layout due to the emergence of multiple-clock-period interconnects.

This gap has motivated recent research to tackle the problem from different aspects of view. Some of them resort to retiming [14], which is a traditional sequential optimization technique that moves flip-flops within a circuit without destroying its functionality. It was used in [20, 4, 16, 19, 15] to pipeline global interconnects so as to reduce the clock period. Although retiming helps to relieve the criticality of global interconnects, there is a lower bound of the clock periods that can be achieved because retiming cannot change the latency of either a (topological) cycle or an input-to-output path in the circuit. In case that the lower bound does not meet the frequency requirement, redesign and re-synthesis may have to be carried out.

One way to avoid redesign is to insert extra wire-pipelining units like flip-flops to pipeline long interconnects, as done within Intel [5] and IBM [13]. It can be shown that if the period lower bound is determined by an input-to-output path, pipelining can reduce the lower bound without affecting the functionality. However, if the period lower bound is given by a cycle, inserting extra flip-flops in it will change its functionality.

C -slow transformation [14] is a technique that slows down the input issue rate¹ of the circuit to accommodate higher frequencies. It was thus used in [17] to retain the functionality when extra flip-flops were inserted in cycles. In other words, throughput was sacrificed (became $1/C$) to meet the frequency requirement.

Instead of slowing down the throughput uniformly over the whole circuit, Latency Insensitive Design (LID) [2, 1], on the other hand, employs a protocol that slows down the throughput of a part of the circuit only when it is needed. As a result, LID can guarantee minimal throughput reduction while satisfying the frequency requirement.

We show in Section 2 that the aforementioned three approaches (retiming, pipelining with C -slow, and pipelining with LID) can be unified under the same objective function of maximizing the *processing rate*, defined as the product of frequency and throughput, as illustrated in Figure 1. In addition, the processing rate of a sequential system is upper bounded by the reciprocal of the maximum cycle ratio of the system, which is only dependent on the clustering. Therefore, we propose an optimal algorithm that finds a clustering with the minimal maximum-cycle-ratio.

The rest of this paper is organized as follows. Section 2 presents the problem formulation. Two previous works are

*This work was supported by the NSF under CCR-0238484.

¹The issue rate is defined as the number of clock periods between successive input changes. An issue rate of 1 indicates that the inputs can change every clock period.

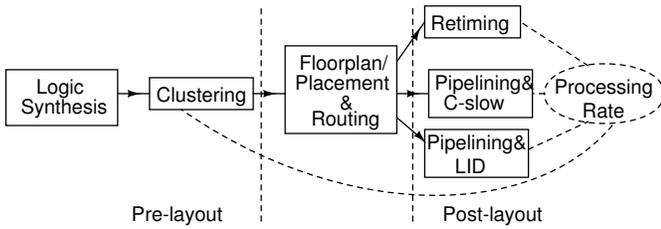


Figure 1: A logical and physical design flow.

reviewed in Section 3. Section 4 defines the notations and constraints used in this paper. Our algorithm is elaborated in Section 5, followed by the implementation details in Section 6. We present some experimental results in Section 7. Conclusions are given in Section 8.

2 Problem formulation

We consider clustering subject to a size limit for clusters. More specifically, each gate has a specified size, as well as each interconnect. We require that the size of each cluster, defined as the sum of the sizes of the gates and the interconnects in the cluster, should be no larger than a given constant A . Replication of gates is allowed, i.e., a gate may be assigned to more than one cluster in the layout. When a gate is replicated, its incident interconnects are also replicated so that the clustered circuit is logically equivalent to the original circuit. Figure 2 (taken from [18]) shows an example of gate replication in a clustering.

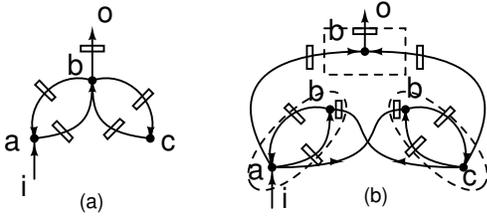


Figure 2: (a) An example circuit; (b) A clustering with 3 replicas of gate b .

Given a particular clustering c , we treat the replicas of gates and the original ones distinctly and denote them all as V_c . We use E_c to denote the set of interconnects among V_c . The clustered circuit is represented as $G_c = (V_c, E_c)$. In order for the circuit to operate at a specified clock period λ , additional wire-pipelining flip-flops are inserted. For all cycle o_c in G_c , let $d(o_c)$ denote the cycle delay, $w(o_c)$ and $w_\lambda(o_c)$ denote the number of flip-flops in o_c before and after additional pipelining flip-flops are inserted, respectively. Assuming $w(o_c) > 0$, the cycle ratio of o_c is defined as $\phi(o_c) = d(o_c)/w(o_c)$. Note that $\phi(o_c)$ is defined using $w(o_c)$, not $w_\lambda(o_c)$. The *maximum cycle ratio* over all the cycles in G_c is denoted as $\phi_c = \max_{o_c \in G_c} \phi(o_c)$.

We define *processing rate* as follows.

Definition 1 For a sequential system, processing rate is defined as the length of processed input sequence per unit time. In particular, it is the product of frequency and throughput in a synchronous system.

The larger the processing rate, the better the sequential system. Given the above definition, the approach of retim-

ing actually maximizes the processing rate by minimizing the period while keeping the throughput. It is interesting to notice that the approach of pipelining with C -slow transformation also maximizes the processing rate for a specified period by computing the least slowdown of the issue rate, which is transformed into throughput reduction. As an alternative, Latency Insensitive Design (LID) helps the clustered circuit reach the maximum throughput for a specified period. Therefore, all the three approaches can be unified under the same objective function of maximizing the processing rate.

It was shown in [3] that the maximum throughput ρ_λ of an LID for a specified period λ can be computed as

$$\rho_\lambda = \min_{\text{cycle } o_c \in G_c} \frac{w(o_c)}{w_\lambda(o_c)}.$$

On the other hand, the fact that the circuit can operate at the specified period λ after the insertion of additional flip-flops implies that $w_\lambda(o_c)\lambda \geq d(o_c)$, i.e., $\frac{1}{w_\lambda(o_c)} \leq \frac{1}{d(o_c)/\lambda}$, $\forall \text{cycle } o_c \in G_c$. Substitute this into the formula of ρ_λ to get

$$\rho_\lambda \leq \min_{o_c \in G_c} \frac{w(o_c)}{d(o_c)/\lambda} = \min_{o_c \in G_c} \frac{\lambda}{\phi(o_c)} = \frac{\lambda}{\phi_c}.$$

It follows that the maximum processing rate of an LID is upper bounded by $\frac{1}{\phi_c}$ since

$$\max \text{ processing rate} = \frac{1}{\lambda} \cdot \rho_\lambda \leq \frac{1}{\lambda} \cdot \frac{\lambda}{\phi_c} = \frac{1}{\phi_c}.$$

It is also an upper bound of the maximum processing rate obtained by the approach of retiming, as shown in [20]. In other words, all the three approaches share the same upper bound of their common objective.

To maximize the processing rate, one can either maximize the upper bound or try to achieve the upper bound. They are equally important. However, since achieving the upper bound requires further knowledge on physical design, such as buffer and flip-flop allowable regions [8, 20] while the upper bound itself is only dependent on the maximum cycle ratio of the clustered circuit, we will consider how to optimally cluster the circuit such that the upper bound is maximized, or equivalently, the maximum cycle ratio is minimized.

In order to compute the maximum cycle ratio, we need to know how to compute the delay of a cycle during clustering. Although local interconnect delays can be obtained using some delay models at synthesis, the delays of global interconnects are not available until layout. Therefore during clustering, we assume that each global interconnect induces an extra constant delay D , e.g., if interconnect (u, v) with delay $d(u, v)$ is assigned to be inter-cluster, then its delay becomes $d(u, v) + D$.

Since we want to minimize the maximum cycle ratio, the path delays from primary inputs (PIs) to primary outputs (POs) can be ignored since they can be mitigated by pipelining. This motivates us to formulate the problem in a strongly connected graph.

Problem 1 (Optimal Clustering Problem)

Given a directed, strongly connected graph $G = (V, E)$, where each vertex $v \in V$ has a delay $d(v)$ and a specified size, and each edge $(u, v) \in E$ has a delay $d(u, v)$, a specified size and a weight $w(u, v)$ (representing the number of flip-flops on it), find a clustering of vertices with possible vertex replication such that: 1. the size of each cluster is no larger than

a given constant A ; 2. each global interconnect induces an extra constant delay D ; 3. the maximum cycle ratio of the clustered circuit is minimized.

We assume that all delays are integral² and thus all cycle ratios are rational. In addition, we assume that each gate has unit size and the size of each interconnect is zero. Our proposed algorithm can be easily extended to handle various size scenarios.

3 Previous work

Pan *et al.* [18] proposed to optimally cluster a sequential circuit such that the lower bound of the period of the clustered circuit was minimized with retiming. However, the period lower bound may not come from a cycle ratio. In addition, their algorithm needs to start from PIs, thus cannot be used to solve our problem in a strongly connected graph. In this sense, they solved a different problem, even though it looks similar to ours.

Their problem was solved by binary search, using a test for feasibility as a subroutine. For each target period, they used a procedure called *labeling computation* to check the feasibility. The procedure starts with label assignment 0 for PIs and $-\infty$ for the other vertices, and repeatedly increases the label values until they all converge or the label value of some PO exceeds the target period, for which the target period is considered infeasible. For each vertex, the amount of increase in its label is computed using another binary search that basically selects the minimum from a candidate set. Because of the nested binary searches, their algorithm is relatively slow. In addition, the algorithm requires $O(|V|^2)$ space to store a pre-computed all-pair longest-path matrix, which is impractical for large designs. Cong *et al.* [9] improved the algorithm by tightening the candidate set to speed up the labeling computation, and by reducing the space complexity to linear dependency. But the improved algorithm still needs the nested binary searches.

Besides the difference in problem formulation, our algorithm differs from theirs in two algorithmic aspects. Firstly, our algorithm focuses on cycles, thus can work on any general graph. Secondly, no binary search is employed in our algorithm. As a result, our algorithm is efficient and essentially incremental. Like [9], our algorithm does not need pre-computed information on paths either.

Except for these differences, [18] revealed some important results on clustering, which we review here to simplify our notations.

- Each cluster has only one output, which is called the *root* of the cluster. If there is a cluster with more than one output, we can replicate the cluster enough times so that each copy of the cluster has only one output.
- For each vertex in V , there is at most one cluster rooted at it and its arrival time (defined in Section 4) is no larger than the arrival times of its replicas.
- If $u \in V$ is an input of the cluster rooted at $v \in V$, then the cluster rooted at v must not contain a replica of u .

²This assumption is not really restrictive in practice because computer works with rational numbers which we can convert to integers by multiplying by a suitably large number.

4 Notations and constraints

For a particular clustering c and a path $p_c = u \rightsquigarrow v$ in G_c , we use $w(p_c)$ to represent the number of flip-flops on p_c , which is the sum of the weights of p_c 's constituent edges. Similarly, $d(p_c)$ represents the delay along p_c , which is the sum of the delays of p_c 's constituent edges and vertices, except for $d(u)$. Note that the delay of an inter-cluster edge $(u, v) \in G_c$ is $d(u, v) + D$. When a path actually forms a cycle o_c , $w(o_c)$ includes the weight of each edge in the cycle only once. Similarly, $d(o_c)$ includes the delay of each edge and vertex in the cycle only once. We assume in this paper that $w(o) > 0$ for all cycle $o \in G$, thus $w(o_c) > 0$ for all cycle $o_c \in G_c$. We use $\phi(o_c)$ to denote the cycle ratio of o_c , and ϕ_c to denote the maximum cycle ratio of G_c .

Since we only need to consider clusters rooted at the vertices in V , at most one for each vertex, we use c_v to refer to the set of vertices that are included in the cluster rooted at $v \in V$. Let $i_v \subset V$ be the set of inputs of c_v and $r_v \subset V_c - V$ be the set of replicas of $v \in V$. In the remainder of this paper, when we say $u \in c_v$ ($u \neq v$), we mean that the cluster rooted at $v \in V$ contains a replica of $u \in V$. For example, Figure 3(a) shows a circuit before clustering. There are five vertices (a-e) and seven edges. Figure 3(b) illustrates a clustering of the circuit with size limit $A = 3$, where dashed circles represent clusters. For each cluster, the vertex whose index is outside the cluster indicates the root. For example, c_a contains replicas of vertices c and e with the input set $i_a = \{d\}$.

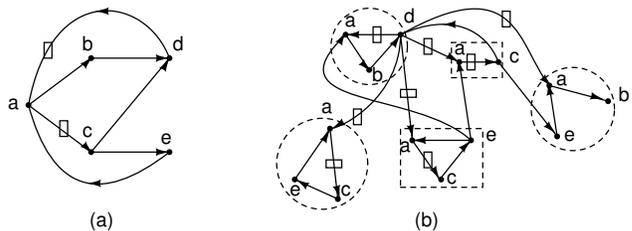


Figure 3: An example of clustering representation

We use a label $t : V \rightarrow \mathfrak{R}$ to denote the arrival time of the vertex. To ease the presentation, we will extend the domain of t to V_c to represent the arrival times of the replicas of the vertices. Based on this, a clustering that satisfies the cluster size requirement and has a maximum cycle ratio no larger than a given rational value ϕ can be characterized as follows.

$$t(v) \geq 0, \quad \forall v \in V \quad (1)$$

$$t(v) \leq t(v'), \quad \forall v' \in r_v, v \in V \quad (2)$$

$$t(v) \geq t(u) + d(u, v) + d(v) - w(u, v)\phi, \quad \forall (u, v) \in E_c, u, v \in c_x, x \in V \quad (3)$$

$$t(v) \geq t(u) + d(u, v) + D + d(v) - w(u, v)\phi, \quad \forall (u, v) \in E_c, u \in i_x, v \in c_x, x \in V \quad (4)$$

$$|c_v| \leq A, \quad \forall v \in V \quad (5)$$

where (1)-(4) guarantee that the arrival times are all achievable, and (5) is the cluster size requirement. In particular, (2) ensures that the arrival time of $v \in V$ is no larger than the arrival times of its replicas.

Following the convention, $(u, v) \in E_c$ is a *critical edge* under ϕ iff it is intra-cluster with $t(v) = t(u) + d(u, v) +$

$d(v) - w(u, v)\phi$, or it is inter-cluster with $t(v) = t(u) + d(u, v) + D + d(v) - w(u, v)\phi$. A *critical path* under ϕ refers to a path whose constituent edges are all critical under ϕ . Vertex u is a *critical input* of c_v under ϕ iff $u \in i_v$ and v can be reached by u through a critical path $p = u \rightarrow x \rightsquigarrow v$ under ϕ where the sub-path $x \rightsquigarrow v$ is in c_v . When a critical path actually forms a cycle, it is then called a *critical cycle*. Cycle o_c is critical under ϕ iff $d(o_c) = w(o_c)\phi$.

A *legal clustering* must satisfy (5). When the arrival times of a legal clustering satisfy (1)-(4) under ϕ , it is called a *feasible clustering* under ϕ . When a critical cycle is present in a feasible clustering under ϕ , it is called a *critical clustering* under ϕ . A given ϕ is *feasible* iff there exists a feasible clustering under ϕ . We must note that for a legal clustering, its maximum cycle ratio is feasible. In fact, any value larger than the maximum cycle ratio of a legal clustering is also feasible.

Consider a feasible clustering c under ϕ . For all $(u, v) \in E$, it is either in E_c with $u \in i_v$, or there is an edge (u', v) such that $u' \in r_u$. In either case, the following inequality is true by (2)-(4).

$$t(v) \geq t(u) + d(u, v) + d(v) - w(u, v)\phi, \forall (u, v) \in E \quad (6)$$

The following lemma provides a lower bound for ϕ .

Lemma 1 *A feasible ϕ is no smaller than the maximum cycle ratio of G , denoted as ϕ_{1b} .*

Proof: Since ϕ is feasible, then, by definition, there exists a clustering c satisfying (1)-(5) under ϕ . Since (6) is implied by (2)-(4), we have $d(o) \leq w(o)\phi$, for all cycle $o \in G$. Therefore, $\phi \geq \phi_{1b}$. \square

Define

$$\Delta(u, v, \phi) \triangleq \max_{p \in u \rightsquigarrow v \text{ in } G} (d(p) - w(p)\phi), \forall u, v \in V$$

Lemma 1 ensures that $\Delta(u, v, \phi)$ is well-defined on feasible ϕ 's.

5 Algorithm

5.1 Overview

The optimal clustering problem asks for a legal clustering with the minimal maximum-cycle-ratio. Since $A > 0$, the clustering with each vertex being a cluster is certainly legal. Starting from it, we will iteratively improve the clustering by reducing its maximum cycle ratio until the optimality is certified.

First of all, the maximum cycle ratio of a legal clustering is feasible and can be efficiently computed using Howard's algorithm [6, 11]. Given a feasible ϕ , we show that, unless ϕ is already the optimal solution, a particular legal clustering can be constructed whose maximum cycle ratio is smaller than ϕ . The smaller ϕ can be obtained by applying Howard's algorithm on the constructed clustering. Therefore, we alternate between applying Howard's algorithm and constructing a better clustering until the minimal ϕ is reached.

5.2 Clustering under a given $\phi > \phi_{1b}$

Given $\phi > \phi_{1b}$, if ϕ is feasible, we show in this section how to construct a feasible clustering under ϕ , i.e., a clustering satisfying (1)-(6) under ϕ , whose maximum cycle ratio is no larger than ϕ .

We choose to first satisfy (1) and (6) because they are independent on clustering, and iteratively update $t(v)$ and c_v to satisfy (2)-(5) while keeping (1) and (6).

Let T denote the arrival time vector, i.e.,

$$T = (t(1), t(2), \dots, t(|V|)).$$

A partial order (\leq) can be defined between two arrival time vectors T and T' as follows.

$$T \leq T' \triangleq t(v) \leq t'(v), \forall v \in V.$$

According to the lattice theory [12], if we treat assignment $t(v) = 0, \forall v \in V$ as the bottom element (\perp) and assignment $t(v) = \infty, \forall v \in V$ as the top element (\top), then the arrival time vector space $\mathbb{R}^{|V|}$ becomes a *complete partially ordered set*, that is, for all $T \in \mathbb{R}^{|V|}$, $\perp \leq T \leq \top$.

To satisfy (1), we set $t(v) = 0, \forall v \in V$. Then we apply Bellman-Ford's algorithm [10], denoted as BF , on E to satisfy (6) under ϕ . Bellman-Ford's algorithm is guaranteed to work as long as $\phi \geq \phi_{1b}$.

The resulting arrival time vector is denoted as

$$T_0 = BF(\perp, \phi).$$

In fact, T_0 is the *least* vector satisfying (1) and (6), as stated in the following lemma.

Lemma 2 $T_0 \leq T$, for all T satisfying (1) and (6).

Proof: Suppose we have a T satisfying (1) and (6) with $t(v) < t_0(v)$ for some $v \in V$. It follows that $t_0(v) > 0$ since $t(v) \geq 0$ by (1). Bellman-Ford's algorithm [10] guarantees that there exists a path $p = u \rightsquigarrow v$ in G such that $t_0(u) = 0$ and $t_0(v) = t_0(u) + d(p) - w(p)\phi$. Since T satisfies (6), we have $t(v) \geq t(u) + d(p) - w(p)\phi = t(u) + t_0(v) \geq t_0(v)$, which contradicts $t(v) < t_0(v)$. Therefore, such a T does not exist and the lemma is true. \square

In order to satisfy (2)-(5) while keeping (1) and (6), we define transformation $\mathcal{L} : (\mathbb{R}^{|V|}, \mathbb{R}) \rightarrow \mathbb{R}^{|V|}$ as follows.

For all $v \in V$, we will construct a new cluster c'_v rooted at v . The procedure starts with $c'_v = \{v\}$ and grows c'_v progressively by including one critical input at a time. Note that when a vertex is put in c'_v , its preceding vertex that is outside c'_v becomes an input of c'_v . Let $t(v)$ and $t'(v)$ denote the arrival time of v before and after c_v is replaced by c'_v , respectively. The procedure will stop only when either $|c'_v| = A$ or $t'(v) \leq t(v)$. If $t'(v) < t(v)$, we keep $t(v)$ and c_v unchanged; otherwise we update $t(v)$ and c_v with $t'(v)$ and c'_v respectively. The resulting arrival time of v is denoted as $\mathcal{L}_v(T, \phi)$. The next lemma helps to identify the critical input to be included at each time.

Lemma 3 *For all $x \notin c'_v$, $t'(v) \geq t(x) + D + \Delta(x, v, \phi)$. In particular, if u is a critical input of c'_v , then $t'(v) = t(u) + D + \Delta(u, v, \phi)$.*

Proof: For the sake of contradiction, we assume that $t'(v) < t(x) + D + \Delta(x, v, \phi)$ for some $x \notin c'_v$. Let p_x be the path where $\Delta(x, v, \phi) = d(p_x) - w(p_x)\phi$. Since $x \notin c'_v$, there exists a vertex $y \in i'_v$ on p_x and the sub-path from y to v has the form of $y \rightarrow z \rightsquigarrow v$, where $z \rightsquigarrow v$ is in c'_v . By (4), we have $t'(v) \geq t(y) + D + d(p_y) - w(p_y)\phi$. In addition, $t(y) \geq t(x) + d(p_x - p_y) - w(p_x - p_y)\phi$ by (6). Thus, $t'(v) \geq t(x) + D + d(p_x) - w(p_x)\phi = t(x) + D + \Delta(x, v, \phi)$, which is a contradiction. Therefore, $t'(v) \geq t(x) + D + \Delta(x, v, \phi)$ for all $x \notin c'_v$.

If c'_v has a critical input u , then, by definition, there exists a path $p = u \rightsquigarrow v$ such that $t(u) + D + d(p) - w(p)\phi = t'(v)$. Since $u \notin c'_v$, we have $t'(v) \geq t(u) + D + \Delta(u, v, \phi)$, thus $d(p) - w(p)\phi \geq \Delta(u, v, \phi)$. On the other hand, $d(p) - w(p)\phi \leq \Delta(u, v, \phi)$ since $\Delta(u, v, \phi)$ is the largest among all paths from u to v . Therefore, $d(p) - w(p)\phi = \Delta(u, v, \phi)$, which concludes our proof. \square

Define $\mathcal{L}(T, \phi)$ as the arrival time vector when all the $\mathcal{L}_v(T, \phi)$'s, $\forall v \in V$, are applied once, followed by Bellman-Ford's algorithm to ensure (6), expressed as

$$\mathcal{L}(T, \phi) \triangleq \text{BF}\left(\left(\mathcal{L}_1(T, \phi), \mathcal{L}_2(T, \phi), \dots, \mathcal{L}_{|V|}(T, \phi)\right), \phi\right).$$

The following lemma shows that \mathcal{L} is an order-preserving transformation.

Lemma 4 *For any T and \bar{T} satisfying (1) and (6), if $T \leq \bar{T}$, then $\mathcal{L}(T, \phi) \leq \mathcal{L}(\bar{T}, \phi)$.*

Proof: We first show that $\mathcal{L}_v(T, \phi) \leq \mathcal{L}_v(\bar{T}, \phi)$, $\forall v \in V$.

For the sake of contradiction, we assume that $\mathcal{L}_v(T, \phi) > \mathcal{L}_v(\bar{T}, \phi)$ for some $v \in V$. The procedure of \mathcal{L}_v guarantees that $\mathcal{L}_v(T, \phi) = \max(t'(v), t(v))$ and $\mathcal{L}_v(\bar{T}, \phi) = \max(\bar{t}'(v), \bar{t}(v))$. Since $t(v) \leq \bar{t}(v)$ by $T \leq \bar{T}$, we have $t'(v) > \bar{t}(v)$, otherwise $\mathcal{L}_v(T, \phi) = t(v) \leq \bar{t}(v) \leq \mathcal{L}_v(\bar{T}, \phi)$, which contradicts the assumption that $\mathcal{L}_v(T, \phi) > \mathcal{L}_v(\bar{T}, \phi)$. In addition, since T satisfies (1) and (6), Lemma 2 ensures that $t(v) \geq t_0(v)$, where $t_0(v)$ is the arrival time of v in vector $T_0 = \text{BF}(\perp, \phi)$. Thus, $t'(v) > t_0(v)$, which implies that cluster c'_v has a critical input u , otherwise the vertices that have critical paths to v are all inside c'_v and we have $t'(v) = t_0(v)$, which is a contradiction. The existence of a critical input u implies that $|c'_v| = A$, otherwise u should have been put in c'_v since $t'(v) > t(v)$. Let p be the path where $\Delta(u, v, \phi) = d(p) - w(p)\phi$. Figure 4(a) shows an example of c'_v .

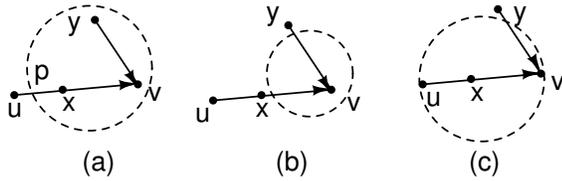


Figure 4: (a) Cluster c'_v ; (b) and (c) two cases of cluster \bar{c}'_v .

Now consider cluster \bar{c}'_v , there are two cases. Firstly, $u \notin \bar{c}'_v$. Thus there exists a vertex $x \in p$ such that $x \in \bar{c}'_v$, as illustrated in Figure 4(b). Since $x \in \bar{c}'_v$, we have $\bar{t}'(v) \geq \bar{t}(x) + D + \Delta(x, v, \phi)$ by Lemma 3. Considering the moment when x was about to be put in \bar{c}'_v , x was the critical input and $u \notin \bar{c}'_v$, thus $t(x) + D + \Delta(x, v, \phi) \geq t(u) + D + \Delta(u, v, \phi)$ by Lemma 3. Together with $t(x) \leq \bar{t}(x)$ (since $T \leq \bar{T}$), we have $\bar{t}'(v) \geq t(u) + D + \Delta(u, v, \phi) = t'(v)$.

Secondly, $u \in \bar{c}'_v$. Given that $|c'_v| = A$ and u is not in c'_v , we know that there exists a vertex $y \in c'_v$ such that $y \in \bar{c}'_v$, as shown in Figure 4(c). By the same argument, we can show that $\bar{t}'(v) \geq \bar{t}(y) + D + \Delta(y, v, \phi) \geq t(y) + D + \Delta(y, v, \phi) \geq t(u) + D + \Delta(u, v, \phi) = t'(v)$.

In either case, we have $\max(t(v), t'(v)) \leq \max(\bar{t}(v), \bar{t}'(v))$, i.e., $\mathcal{L}_v(T, \phi) \leq \mathcal{L}_v(\bar{T}, \phi)$, which is a contradiction. Therefore, the assumption is wrong and $\mathcal{L}_v(T, \phi) \leq \mathcal{L}_v(\bar{T}, \phi)$ is true, $\forall v \in V$. It is easy to verify that $\mathcal{L}(T, \phi) \leq \mathcal{L}(\bar{T}, \phi)$ after applying Bellman-Ford's algorithm. \square

We say that T is a *fixpoint* of \mathcal{L} under ϕ if and only if $T = \mathcal{L}(T, \phi)$. The following theorem bridges the existence of a fixpoint and the feasibility of ϕ .

Theorem 1 *ϕ is feasible if and only if \mathcal{L} has a fixpoint under ϕ .*

Proof: (\rightarrow): If ϕ is feasible, then, by definition, there exists a legal clustering c whose arrival time vector T satisfies (1)-(4) and (6) under ϕ . We claim that $\mathcal{L}_v(T, \phi) \leq t(v)$, $\forall v \in V$. Otherwise, $t'(v) > t(v) \geq t_0(v)$ for some $v \in V$, and c'_v has a critical input u , as shown in Figure 4(a). We can conduct a similar case study as Figure 4(b) and 4(c) to show that $t(v) \geq t'(v)$, which is a contradiction. On the other hand, $t(v) \leq \mathcal{L}_v(T, \phi)$ by the procedure of $\mathcal{L}_v(T, \phi)$. Therefore, $\mathcal{L}_v(T, \phi) = t(v)$, $\forall v \in V$. Given that T satisfies (6), applying Bellman-Ford's algorithm gives $T = \mathcal{L}(T, \phi)$, i.e., T is a fixpoint of \mathcal{L} under ϕ .

(\leftarrow): If \mathcal{L} has a fixpoint under ϕ , then, by the definition of \mathcal{L} , the constructed clustering is legal and the arrival time vector satisfies (1)-(4) and (6) under ϕ . Therefore, ϕ is feasible. \square

In fact, according to the lattice theory [12], if \mathcal{L} , defined on a complete partially ordered set, has a fixpoint under ϕ , then it has a *least fixpoint* T^ϕ , defined as

$$(T^\phi = \mathcal{L}(T^\phi, \phi)) \wedge (\forall T : T = \mathcal{L}(T, \phi) : T^\phi \leq T).$$

We use c^ϕ to denote the clustering constructed by $\mathcal{L}(T^\phi, \phi)$. In fact, if $t^\phi(v) > t_0(v)$, then there is a critical path from $u \in V$ to v with $t^\phi(u) = t_0(u)$. This is made precise in the following lemma.

Lemma 5 *If $t^\phi(v) > t_0(v)$, then there exists a sequence of vertices $x_i \in V$, $i = 0, 1, \dots, k-1$ such that $t_0(x_0) = t^\phi(x_0)$, $x_k = v$, and x_i is a critical input of cluster $c^\phi_{x_{i+1}}$.*

Proof: Since $t^\phi(v) > t_0(v)$, we know that cluster c^ϕ_v has critical inputs, otherwise the vertices that have critical paths to v are all inside c^ϕ_v and we have $t_0(v) = t^\phi(v)$, which is a contradiction.

Suppose otherwise that such a sequence does not exist, namely, all the critical paths terminating at v are actually critical cycles, where each constituent vertex $u \in V$ has $t^\phi(u) > t_0(u) \geq 0$. Choose the one that contains all of them, denoted as O . For the example in Figure 5, we will choose O to be $a \rightsquigarrow e \rightsquigarrow d \rightsquigarrow c \rightsquigarrow g \rightsquigarrow f \rightsquigarrow d \rightsquigarrow c \rightsquigarrow b \rightsquigarrow a$. Now

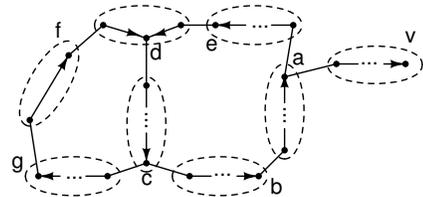


Figure 5: Vertices that have critical paths to v .

consider any incoming edge of O (from a vertex outside of O to a vertex in O), it must be non-critical, otherwise we can trace back from this edge and find another critical cycle that is not in O , which is a contradiction. Since the arrival times of the vertices in O are all greater than zero, we can decrease them simultaneously while keeping the arrival times of other vertices unchanged until some incoming edge of O becomes

critical or the arrival time of some vertex in O is reduced to zero. For either case, we obtain a fixpoint less than T^ϕ , which is a contradiction. Therefore, the lemma is true. \square

To reach a fixpoint, iterative method can be used on \mathcal{L} . It starts with T_0 as the initial vector, iteratively computes new vectors from previous ones $T_1 = \mathcal{L}(T_0, \phi)$, $T_2 = \mathcal{L}(T_1, \phi)$, \dots until it finds a T_n such that $T_n = T_{n-1}$. The following lemma states that applying iterative method on \mathcal{L} will converge to its least fixpoint in a finite number of iterations.

Lemma 6 *If ϕ is feasible, applying iterative method on \mathcal{L} will converge to T^ϕ in a finite number of iterations.*

Proof: Since we start with $T = T_0 \leq T^\phi$, Lemma 4 ensures that $T \leq T^\phi$ at each iteration. Therefore, if \mathcal{L} converges, the fixpoint has to be the least fixpoint. What remains is to show that \mathcal{L} is finitely convergent.

By Lemma 3 and 5, if $t^\phi(v) > t_0(v)$, then $t^\phi(v)$ can be written as

$$t^\phi(v) = t_0(x_0) + \sum_{0 \leq i \leq k-1} (D + \Delta(x_i, x_{i+1}, \phi)),$$

where $x_i \in V$ and $x_k = v$. Given that each vertex in V has at most one cluster rooted at it, we know that $k \leq |V|$, thus $t^\phi(v) \leq U$, where $U = (|V| - 1)D + |V| \max_{u,v \in V} \Delta(u, v, \phi)$.

On the other hand, since ϕ is a rational number, it can be expressed as p/q , where p and q are integers and $q \neq 0$. If $t(v)$ is increased during the iteration, the amount of increase will be at least $1/q$. Therefore, if \mathcal{L} does not converge after $|V|Uq$ iterations, then there exists a vertex $v \in V$ whose $t(v) > U \geq t^\phi(v)$, which contradicts $T \leq T^\phi$, which concludes our proof. \square

The next result is a corollary of Lemma 4-6.

Lemma 7 *($\forall v \in V : t(v) > t_0(v)$) implies that ϕ is infeasible.*

Proof: Suppose otherwise that ϕ is feasible. Then, by Lemma 4 and 6, when T^ϕ is reached, we have $(\forall v \in V : t^\phi(v) > t_0(v))$, which contradicts Lemma 5. Therefore, ϕ is infeasible. \square

5.3 Optimality checking

Given a legal clustering, its maximum cycle ratio ϕ is feasible. If ϕ is not optimal, then we can find a feasible $\phi' < \phi$, which is specified in the following lemma.

Lemma 8 *Given that ϕ is the maximum cycle ratio of a legal clustering, if ϕ is not optimal, then $\phi - 1/(|V|N_{\text{ff}})^2$ is also feasible, where N_{ff} is the maximum number of flip-flops on any acyclic path in G .*

Proof: Let o denote the cycle with the maximum cycle ratio, that is, $\phi = d(o)/w(o)$. If ϕ is not optimal, it means that there exists another legal clustering whose maximum cycle ratio ϕ' is smaller than ϕ . Let o' be the cycle with $\phi' = d(o')/w(o')$. The difference between ϕ' and ϕ can be written as

$$\phi - \phi' = \frac{d(o)}{w(o)} - \frac{d(o')}{w(o')} = \frac{d(o)w(o') - d(o')w(o)}{w(o)w(o')} \geq \frac{1}{w(o)w(o')},$$

since all delays are integers. In addition, since each vertex in V has at most one cluster rooted at it, we know that

both $w(o)$ and $w(o')$ are no larger than $|V|N_{\text{ff}}$, where N_{ff} is the maximum number of flip-flops on any acyclic path in G . Therefore, $\phi - \phi' \geq 1/(|V|N_{\text{ff}})^2$. In other words, $\phi - 1/(|V|N_{\text{ff}})^2$ is also feasible. \square

It implies that we can certify the optimality of ϕ by checking the feasibility of $\phi - 1/(|V|N_{\text{ff}})^2$. The algorithm for finding the optimal ϕ is presented in Figure 6. It first computes a feasible ϕ by treating each vertex as a cluster, and computes a lower bound of ϕ_{lb} by Lemma 1. After that, it checks the feasibility of $\phi - 1/(|V|N_{\text{ff}})^2$ by iterative method on \mathcal{L} . If \mathcal{L} converges, it means that we find a better clustering whose maximum cycle ratio is at most $\phi - 1/(|V|N_{\text{ff}})^2$ and can be computed by Howard's algorithm. The evidence of $(\forall v \in V : t(v) > t_0(v))$ or the fact that ϕ is reduced below $\phi_{\text{lb}} + 1/(|V|N_{\text{ff}})^2$ immediately certifies the optimality of the current feasible ϕ .

Algorithm Optimal clustering
Input: A directed graph $G = (V, E)$, A , D .
Output: A clustering c^{opt} with ϕ^{opt} .

```

 $c_v^{\text{opt}} \leftarrow c_v \leftarrow \{v\}, \forall v \in V;$ 
 $\phi^{\text{opt}} \leftarrow \phi \leftarrow$  maximum cycle ratio of  $G_c;$ 
 $\phi_{\text{lb}} \leftarrow$  maximum cycle ratio of  $G;$ 
While ( $\phi \geq \phi_{\text{lb}} + 1/(|V|N_{\text{ff}})^2$ ) do
   $\phi \leftarrow \phi - 1/(|V|N_{\text{ff}})^2;$ 
   $T \leftarrow T_0 \leftarrow BF(\perp, \phi);$ 
  While ( $(T \neq \mathcal{L}(T, \phi)) \wedge (\exists v \in V : t(v) = t_0(v))$ ) do
     $T \leftarrow \mathcal{L}(T, \phi);$ 
  If ( $\forall v \in V : t(v) > t_0(v)$ ) then
    break;
   $\phi \leftarrow$  maximum cycle ratio of  $G_{c\phi};$ 
   $c^{\text{opt}} \leftarrow c^\phi; \phi^{\text{opt}} \leftarrow \phi;$ 
Return  $c^{\text{opt}}$  and  $\phi^{\text{opt}};$ 

```

Figure 6: Pseudocode of optimal clustering algorithm

The correctness of the algorithm is stated in the following theorem.

Theorem 2 *The algorithm in Figure 6 will terminate with the optimal ϕ .*

Proof: First of all, ϕ is reduced during the execution of the outer while-loop in Figure 6. Since the amount of decrease in ϕ is at least $1/(|V|N_{\text{ff}})^2$ after each loop, the algorithm will terminate in $(\phi_{\text{ub}} - \phi_{\text{lb}})|V|^2N_{\text{ff}}^2$ loops, or $|V|^3N_{\text{ff}}^2D$ loops as $\phi_{\text{ub}} \leq \phi_{\text{lb}} + |V|D$.

When it terminates, we have either $\phi < \phi_{\text{lb}} + 1/(|V|N_{\text{ff}})^2$, or $(\forall v \in V : t(v) > t_0(v))$ under $\phi - 1/(|V|N_{\text{ff}})^2$. For the first case, Lemma 8 ensures that ϕ is optimal, otherwise $\phi - 1/(|V|N_{\text{ff}})^2$ is feasible, which contradicts Lemma 1. For the second case, $\phi - 1/(|V|N_{\text{ff}})^2$ is infeasible by Lemma 7, which, by Lemma 8, implies that ϕ is optimal. The optimal ϕ and the corresponding clustering are recorded in ϕ^{opt} and c^{opt} . \square

6 Implementation details

6.1 Implementation of \mathcal{L}_v

Our implementation of \mathcal{L}_v is similar to Cong *et al.* [8]. To characterize critical inputs, we introduce another label

$\delta(u)$, $\forall u \in V$. Before the construction of c'_v , we assign $\delta(u)$ with $-\infty$ for all $u \neq v$ in V while $\delta(v)$ with 0. At each time, the vertex $u \in i'_v$ with the largest $t(u) + \delta(u)$ is identified. If $t(u) + D + \delta(u) \leq t(v)$, the construction is completed. Otherwise, we put it in c'_v and update $\delta(x)$ with $\max(\delta(x), \delta(u) + d(u) + d(x, u) - w(x, u)\phi)$, for all $(x, u) \in E$. This procedure will iterate until either $|c'_v| = A$ or the last vertex u identified has $t(u) + D + \delta(u) \leq t(v)$.

To validate the above procedure, we need to show that it is equivalent to $\mathcal{L}_v(\mathbb{T}, \phi)$, or equivalently, to show that it can always identify the critical input of c'_v . This is fulfilled by the next lemma and corollary.

Lemma 9 For all $u \in c'_v$, $\delta(u) = \Delta(u, v, \phi)$.

Proof: We prove it by induction on the size of c'_v . At the beginning, $c'_v = \{v\}$ and $\delta(v) = 0 = \Delta(v, v, \phi)$. Suppose that the lemma is true for $|c'_v| \leq k < A$, we need to show that the lemma is also true for $|c'_v| = k + 1$. Therefore we start with $|c'_v| = k$ and let $u \in i'_v$ be the vertex with the largest $t(u) + D + \delta(u) > t(v)$.

We use p to denote the path where $d(p) - w(p)\phi = \Delta(u, v, \phi)$. Since $u \notin c'_v$, there exists a vertex $x \in i'_v$ on p such that the sub-path from x to v has the form of $x \rightarrow y \rightsquigarrow v$, where $y \rightsquigarrow v$ is in c'_v . Let p_1 and p_2 be the sub-path from u to x and from y to v , respectively. Since $d(p) - w(p)\phi = \Delta(u, v, \phi)$, we have $d(p_1) - w(p_1)\phi = \Delta(u, x, \phi)$ and $d(p_2) - w(p_2)\phi = \Delta(y, v, \phi)$.

Given that $u \in i'_v$ is the vertex with the largest $t(u) + \delta(u)$, we know that $t(u) + \delta(u) \geq t(x) + \delta(x)$. (6) ensures that $t(x) \geq t(u) + d(p_1) - w(p_1)\phi$. On the other hand, we have $\delta(x) \geq \delta(y) + d(x, y) - w(x, y)\phi$ since we updated $\delta(x)$ to be no less than $\delta(y) + d(x, y) - w(x, y)\phi$ when y was put in c'_v . Further, $\delta(y) = \Delta(y, v, \phi)$ by the inductive hypothesis. Consequently,

$$\begin{aligned}
& t(u) + \delta(u) \\
& \geq t(x) + \delta(x) \\
& \geq (t(u) + d(p_1) - w(p_1)\phi) + \delta(x) \\
& \geq (t(u) + d(p_1) - w(p_1)\phi) + (\delta(y) + d(x, y) - w(x, y)\phi) \\
& = (t(u) + d(p_1) - w(p_1)\phi) \\
& \quad + (\Delta(y, v, \phi) + d(x, y) - w(x, y)\phi) \\
& = (t(u) + d(p_1) - w(p_1)\phi) + \\
& \quad ((d(p_2) - w(p_2)\phi) + d(x, y) - w(x, y)\phi) \\
& = t(u) + d(p) - w(p)\phi \\
& = t(u) + \Delta(u, v, \phi),
\end{aligned}$$

or $\delta(u) \geq \Delta(u, v, \phi)$. However, $\delta(u) \leq \Delta(u, v, \phi)$ since $\Delta(u, v, \phi)$ is the largest among all paths from u to v . Therefore, $\delta(u) = \Delta(u, v, \phi)$. \square

Corollary 9.1 The vertex $u \in i'_v$ with the largest $t(u) + D + \delta(u) \geq t(v)$ is the critical input of c'_v .

Proof: Suppose u is not a critical input. Let p denote the path where $d(p) - w(p)\phi = \delta(u)$. Since p is not critical, we have $t'(v) > t(u) + D + d(p) - w(p)\phi = t(u) + D + \delta(u) \geq t(v)$. It implies that c'_v has a critical input x and a critical path $p' : x \rightarrow y \rightsquigarrow v$ such that $p'_1 : y \rightsquigarrow v$ is in c'_v and $t'(v) = t(x) + D + d(p') - w(p')\phi$. In addition, $t'(v) = t(x) + D + \Delta(x, v, \phi)$ by Lemma 3. Thus, $\Delta(x, v, \phi) = d(p') - w(p')\phi$, which implies that $\Delta(y, v, \phi) = d(p'_1) - w(p'_1)\phi$.

On the other hand, since $y \in c'_v$, we have $\delta(y) = \Delta(y, v, \phi)$ by Lemma 9. Hence $\delta(x) \geq \delta(y) + d(x, y) - w(x, y)\phi = d(p') - w(p')\phi = \Delta(x, v, \phi)$. Since $\Delta(x, v, \phi)$ is the largest among all paths from x to v , we have $\delta(x) = \Delta(x, v, \phi)$. It follows that $t(x) + D + \delta(x) = t'(v) > t(u) + D + \delta(u)$, which contradicts that u is the input with the largest $t(u) + \delta(u)$. Therefore, the lemma is true. \square

The pseudocode for computing $\mathcal{L}_v(\mathbb{T}, \phi)$ is given in Figure 7. It employs a heap Q for bookkeeping the vertices $u \in i'_v$ whose $t(u) + D + \delta(u) > t(v)$. At each iteration, it puts in c'_v the input $u \in Q$ with the largest $t(u) + D + \delta(u)$ and updates $\delta(x)$ for each fanin of u that becomes an input in i'_v . In our implementation, we choose Fibonacci heap [10] for Q .

Input: $G = (V, E)$, A , D , \mathbb{T} , ϕ , $v \in V$.
Output: $t'(v)$, c'_v .

```

 $\delta(u) \leftarrow -\infty, \forall u \in V; \delta(v) \leftarrow 0;$ 
 $Q \leftarrow \{v\}; c'_v \leftarrow \emptyset; t'(v) \leftarrow t(v);$ 
While  $((Q \neq \emptyset) \wedge (|c'_v| < A))$  do
   $u \leftarrow \text{extract from } Q \text{ with max } t(u) + \delta(u);$ 
   $c'_v \leftarrow c'_v \cup \{u\};$ 
  For  $e = (x, u) \in E$  with  $x \notin c'_v$  do
     $\delta(x) \leftarrow \max(\delta(x), \delta(u) + d(u) + d(x, u) - w(x, u)\phi);$ 
    If  $((x \notin Q) \wedge (t(x) + D + \delta(x) > t(v)))$  then
       $Q \leftarrow Q \cup \{x\};$ 
   $t'(v) \leftarrow \max t(u) + D + \delta(u)$  in  $Q$ , if  $Q \neq \emptyset;$ 
Return  $t'(v)$  and  $c'_v;$ 

```

Figure 7: Pseudocode of $\mathcal{L}_v(\mathbb{T}, \phi)$

6.2 Variations of \mathcal{L}

In Section 5.2, $\mathcal{L}(\mathbb{T}, \phi)$ is obtained by applying all the \mathcal{L}_v 's, $\forall v \in V$, once followed by Bellman-Ford's algorithm. In our implementation, all the $\mathcal{L}_v(\mathbb{T}, \phi)$'s are not computed at the same time. Intuitively, if previously computed \mathcal{L}_v 's can be taken into account in later computations of others, the convergence rate may be accelerated.

This motivates our study on a variation of \mathcal{L} , in which later computations of \mathcal{L}_v 's are based on previously computed ones, and each computation of \mathcal{L}_v is followed by Bellman-Ford's algorithm. Let $\mathcal{J}_v(\mathbb{T}, \phi)$ denote the vector after $t(v)$ is updated with $\mathcal{L}_v(\mathbb{T}, \phi)$, that is,

$$\mathcal{J}_v(\mathbb{T}, \phi) = (t(1), \dots, t(v-1), \mathcal{L}_v(\mathbb{T}, \phi), t(v+1), \dots, t(|V|)).$$

Define

$$\bar{\mathcal{L}}(\mathbb{T}, \phi) \triangleq BF\left(\mathcal{J}_{i_{|V|}}(\dots BF(\mathcal{J}_{i_1}(\mathbb{T}, \phi), \phi), \dots, \phi), \phi\right),$$

where $i_1, \dots, i_{|V|} \in V$. It can be seen that different evaluation orders of V give different $\bar{\mathcal{L}}$'s. However, they all satisfy the following relation.

Lemma 10 For any $\mathbb{T} \leq \mathbb{T}^\phi$ satisfying (1) and (6) under a feasible ϕ and any evaluation order of V , $\mathcal{L}(\mathbb{T}, \phi) \leq \bar{\mathcal{L}}(\mathbb{T}, \phi) \leq \bar{\mathcal{L}}(\mathbb{T}^\phi, \phi) = \mathbb{T}^\phi$.

Proof: Let $\bar{\mathcal{L}}_v(\mathbb{T}, \phi)$ denote the arrival time of $v \in V$ in $\bar{\mathcal{L}}(\mathbb{T}, \phi)$. The definition of $\bar{\mathcal{L}}(\mathbb{T}, \phi)$ implies that $\mathcal{L}_v(\mathbb{T}, \phi) \leq$

$\bar{\mathcal{L}}_v(\mathbb{T}, \phi)$, $\forall v \in V$, independent of the evaluation order. It follows that $\mathcal{L}(\mathbb{T}, \phi) \leq \bar{\mathcal{L}}(\mathbb{T}, \phi)$. On the other hand, since $\mathcal{L}_v(\mathbb{T}^\phi, \phi) = t^\phi(v)$, we have $\mathcal{J}_v(\mathbb{T}, \phi) = \mathbb{T}^\phi$, hence $\bar{\mathcal{L}}(\mathbb{T}^\phi, \phi) = \mathbb{T}^\phi$. What remains to show is $\bar{\mathcal{L}}(\mathbb{T}, \phi) \leq \mathbb{T}^\phi$. To this aim, we observe that $BF(\mathcal{J}_v(\mathbb{T}, \phi)) \leq \mathbb{T}^\phi$, $\forall v \in V$, provided that $\mathbb{T} \leq \mathbb{T}^\phi$. Based on this, we can show by induction that $\bar{\mathcal{L}}(\mathbb{T}, \phi) \leq \mathbb{T}^\phi$. Therefore, the lemma is true. \square

As a corollary, the next result ensures that we can apply iterative method on $\bar{\mathcal{L}}$ to reach \mathbb{T}^ϕ .

Corollary 10.1 *If ϕ is feasible, applying iterative method on $\bar{\mathcal{L}}$ will converge to \mathbb{T}^ϕ in a finite number of iterations, independent of the evaluation order of V .*

Proof: Since ϕ is feasible, \mathcal{L} is finitely convergent by Lemma 6. Let K be the number of iterations such that $\mathcal{L}^K(\mathbb{T}_0, \phi) = \mathbb{T}^\phi$. The corollary can be proved if we can show that $\bar{\mathcal{L}}^K(\mathbb{T}_0, \phi) = \mathbb{T}^\phi$, or equivalently, $\mathcal{L}^K(\mathbb{T}_0, \phi) \leq \bar{\mathcal{L}}^K(\mathbb{T}_0, \phi) \leq \mathcal{L}^K(\mathbb{T}^\phi, \phi)$.

The former part can be derived from Lemma 10 because $\bar{\mathcal{L}}^K(\mathbb{T}_0, \phi) \geq \mathcal{L}(\bar{\mathcal{L}}^{K-1}(\mathbb{T}_0, \phi)) \geq \dots \geq \mathcal{L}^K(\mathbb{T}_0, \phi)$. The latter part is also a consequence of Lemma 10 since $\mathbb{T}_0 \leq \mathbb{T}^\phi$. \square

6.3 Reduced clustering representation

It was shown in [11] that Howard's algorithm was by far the fastest algorithm for maximum cycle ratio computation. Given a clustered circuit $G_c = (V_c, E_c)$ with edge delays and weights specified, Howard's algorithm finds the maximum cycle ratio in $O(N_c|E_c|)$ time, where N_c is the product of the out-degrees of all the vertices in V_c . Since vertex replication is allowed, N_c and $|E_c|$ could be $|V|N$ and $|V||E|$ respectively, where N is the product of the out-degrees of the vertices in V .

To reduce the complexity, we propose a reduced clustering representation. For each cluster, we use edges from its inputs to its output (root) to represent the paths between them such that the delay and weight of an edge correspond to the delay and weight of an acyclic input-to-output path. Figure 8 shows the reduced representation of the clustered circuit in Figure 3(b).

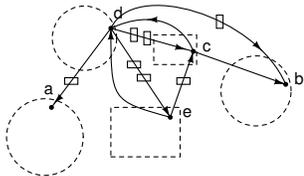


Figure 8: The reduced clustering representation of Figure 3(b).

Let ϕ_c^r denote the maximum cycle ratio of the reduced representation for clustering c . The following lemma formulates the relation among ϕ_c , ϕ_c^r and the lower bound ϕ_{1b} defined in Lemma 1.

Lemma 11 *For any clustering c , $\phi_c = \max(\phi_c^r, \phi_{1b})$.*

Proof: All the cycles in G_c can be classified into two groups according to whether they contain an inter-cluster edge or not. If a cycle contains only intra-cluster edges, its maximum cycle ratio is upper bounded by ϕ_{1b} . If a cycle contains inter-cluster edges, it is present in the reduced representation and thus is upper bounded by ϕ_c^r . \square

One benefit of the reduced clustering representation is that we can now represent the clustered circuit without explicit vertex replication, that is, using V instead of V_c . Let $G_c^r = (V, E_c^r)$ denote the reduced representation for clustering c . We call an edge in G_c^r *redundant* if its removal will not affect the maximum cycle ratio of G_c^r . The following lemma provides a criterion to prune the redundant edges so that Howard's algorithm can find the maximum cycle ratio of G_c^r more efficiently.

Lemma 12 *Let c denote a feasible clustering under ϕ , E_c^r denote its reduced representation, e_1 and e_2 denote two edges from $u \in V$ to $v \in V$ in E_c^r , $d(e_1)$ and $d(e_2)$ denote their delays respectively, and $w(e_1)$ and $w(e_2)$ denote their weights respectively. If $w(e_1) \geq w(e_2)$ and $d(e_1) - w(e_1)\phi \geq d(e_2) - w(e_2)\phi$, then e_2 can be pruned.*

Proof: Since c is feasible under ϕ , we know that $\phi \geq \phi_c \geq \phi_c^r$. If e_2 is not involved in any cycle in G_c^r , then e_2 can be safely pruned as it will not affect the computation of the maximum cycle ratio. Otherwise, let o_2 be a cycle in G_c^r involving e_2 , and o_1 be the cycle in G_c^r such that $o_1 = \{e_1\} \cup o_2 - \{e_2\}$. What remains to show is that if o_2 is a critical cycle under ϕ_c^r , so is o_1 .

By definition, if o_2 is a critical cycle under ϕ_c^r , then $d(o_2) - w(o_2)\phi_c^r = 0$. Since o_1 differs o_2 in e_1 only, we have

$$\begin{aligned} & d(o_1) - w(o_1)\phi_c^r \\ &= d(o_2) - d(e_2) + d(e_1) - (w(o_2) - w(e_2) + w(e_1))\phi_c^r \\ &= d(e_1) - w(e_1)\phi_c^r - (d(e_2) - w(e_2)\phi_c^r) \\ &= d(e_1) - w(e_1)\phi - (d(e_2) - w(e_2)\phi) \\ &\quad + (w(e_1) - w(e_2))(\phi - \phi_c^r) \\ &\geq 0, \end{aligned}$$

provided that $d(e_1) - w(e_1)\phi \geq d(e_2) - w(e_2)\phi$ and $w(e_1) \geq w(e_2)$. On the other hand, $d(o_1) \leq w(o_1)\phi_c^r$ since ϕ_c^r is the maximum cycle ratio. Therefore, $d(o_1) = w(o_1)\phi_c^r$, i.e., o_1 is also a critical cycle under ϕ_c^r . This implies that we can safely remove e_2 and the resulting representation has the same maximum cycle ratio as G_c^r . \square

In our implementation, we employ another two parameters $\text{pd} : V \rightarrow \{R\}$ and $\text{pw} : V \rightarrow \{Z\}$ to record the path delays and weights from the inputs of a cluster to its output, respectively. More specifically, we set $\text{pw}(u) = \text{pd}(u) = \emptyset$, $\forall u \in V$, before $\mathcal{L}_v(\mathbb{T}, \phi)$ is about to be carried out for some $v \in V$. After that, whenever a vertex u is put in \mathcal{C}_v , we compute the pd and pw values of its preceding vertices based on $\text{pd}(u)$ and $\text{pw}(u)$, followed by pruning.

7 Experimental results

We implemented the algorithm in a PC with a 2.4 GHz Xeon CPU, 512 KB 2nd level cache memory and 1GB RAM. To compare with the algorithm in [18], we used the same test files, which were generated from the ISCAS-89 benchmark suite. For each test case, we introduced a flip-flop with directed edges from each PO to it and from it to each PI so that every PI-to-PO path became a cycle. As in [18], the size and delay of each gate was set to 1, intra-cluster delays were 0, and inter-cluster interconnects had delays $D = 2$. The circuits used are summarized in Table 1. We also list the maximum cycle ratio of the circuit before clustering in

Table 2: Optimal Maximum-cycle-ratio

Circuit	$A = 5\% V $			$A = 10\% V $			$A = 20\% V $		
	ϕ^{opt}	#step	time(s)	ϕ^{opt}	#step	time(s)	ϕ^{opt}	#step	time(s)
s208	13.00	4	0.20	11.00	4	0.09	10.00	3	0.04
s349	18.00	14	1.73	16.00	21	0.47	14.67	14	1.39
s420	14.00	3	0.33	13.00	3	0.34	12.00	2	0.00
s635	75.00	4	5.71	70.00	4	8.02	68.00	4	6.43
s838	17.00	3	0.69	16.00	2	0.01	16.00	2	0.01
s1196	26.00	3	2.86	25.00	3	3.40	24.00	2	0.02
s1423	55.00	3	2.77	53.00	2	0.07	53.00	2	0.06
s1512	23.78	15	166.36	22.50	8	0.79	22.50	8	0.68
s3330	14.33	11	11.34	14.00	10	1.82	14.00	10	1.82
s4863	30.25	8	133.54	30.00	5	3.85	30.00	5	3.70
s5378	21.00	3	0.94	21.00	3	0.73	21.00	3	0.73
s9234	38.00	3	3.78	38.00	3	2.98	38.00	3	2.98
s35932	27.00	4	48.13	27.00	4	47.79	27.00	4	46.59
s38584	48.00	2	21.92	48.00	2	21.90	48.00	2	21.46

Table 1: Sequential Circuits from ISCAS-89

Circuit	$ V $	$ E $	N_{ff}	ϕ_{lb}
s208	104	183	2	10.00
s349	161	285	8	14.00
s420	218	385	2	12.00
s635	286	478	2	66.00
s838	446	789	2	16.00
s1196	529	1024	3	24.00
s1423	657	1170	6	53.00
s1512	780	1286	3	22.50
s3330	1789	2890	5	14.00
s4863	2342	4093	9	30.00
s5378	2779	4262	13	21.00
s9234	5597	6932	9	38.00
s35932	16065	28590	48	27.00
s38584	19253	33061	30	48.00

column ϕ_{lb} , which provides a lower bound of the solution by Lemma 1.

Although theoretically the algorithm in Fig. 6 will reach the exact solution without being provided a precision, we have to consider the impact of floating point error introduced by practical finite precision arithmetic, due to the divisions involved in the maximum-cycle-ratio computation. In our experiments, we set the error to be 0.001. Since $1/(|V|N_{\text{ff}})$ is generally smaller than 0.001, we set the precision of ϕ^{opt} to be 0.01.

For each circuit, we tested three size bounds: A is 5%, 10% and 20% of the number of gates. The results are shown in Table 2. Since we approach the minimal maximum-cycle-ratio by gradual reduction in the algorithm, we also report the number of reductions for each scenario of A in column “#step”. Column “time(s)” lists the running time in seconds.

To compare the running time in [18], where the optimal clock period is integral, we set the precision of ϕ^{opt} to be 1 and run the algorithm again for $A = 5\%|V|, 10\%|V|, 20\%|V|$ respectively. The obtained ϕ^{opt} matches the result in [18] for all the scenarios of A . The only running time information given in [18] is the largest running time per step among the three scenarios, which we list in Figure 3 under column “[18]”. We then compute ours in column “presented”. Row “arith” (“geo”) gives the arithmetic (geometric) mean of the first 11 circuits. We must note that the running times listed

are already scaled to take into account the difference in CPU frequencies between their machine and ours.

Table 3: Running Time Comparison with [18]

Circuit	time/step (s)	
	[18]	presented
s208	0.01	0.00
s349	0.19	0.02
s420	0.03	0.00
s635	0.10	0.03
s838	0.13	0.01
s1196	0.08	0.02
s1423	0.65	0.04
s1512	1.22	0.06
s3330	1.07	0.17
s4863	82.30	1.04
s5378	7.58	0.31
s9234	n/a	1.20
s35932	n/a	11.46
s38584	n/a	11.07
arith	19.39X	1
geo	12.78X	1

It can be seen that our algorithm takes much less time per step than the algorithm in [18]. The improvements are greater for larger circuits. The average speed-up is more than one order of magnitude. In addition, for most of the circuits, our algorithm finds the optimal solution in just a few steps, which is generally less than the number of iterations conducted in a binary search, which are not given in [18].

8 Conclusion

Processing rate, defined as the product of frequency and throughput, is identified as an important metric for sequential circuits. We show that the processing rate of a sequential circuit is upper bounded by the reciprocal of its maximum cycle ratio, which is only dependent on the clustering of the circuit. The problem of processing rate optimization is formulated as seeking an optimal clustering with minimal maximum-cycle-ratio in a general graph. An iterative algorithm is proposed that finds the minimal maximum-cycle-ratio. Since our algorithm avoids binary search and is essentially incremental, it has the potential to be combined with

other optimization techniques, such as gate sizing, budgeting, etc., thus can be used in incremental design methodologies [7]. In addition, since maximum cycle ratio is a fundamental metric, the proposed algorithm can be adapted to suit other traditional designs.

References

- [1] L. P. Carloni, K. L. McMillan, A. Saldanha, and A. L. Sangiovanni-Vincentelli. A methodology for correct-by-construction latency insensitive design. In *ICCAD*, 1999.
- [2] L. P. Carloni, K. L. McMillan, and A. L. Sangiovanni-Vincentelli. Latency insensitive protocols. In *CAV*, 1999.
- [3] M. R. Casu and L. Macchiarulo. A new approach to latency insensitive design. In *DAC*, 2004.
- [4] C. Chu, E. F. Y. Young, D. K. Y. Tong, and S. Dechu. Retiming with interconnect and gate delay. In *ICCAD*, pages 221–226, 2003.
- [5] P. Cocchini. Concurrent flip-flop and repeater insertion for high performance integrated circuits. In *ICCAD*, pages 268–273, 2002.
- [6] J. Cochet-Terrasson, G. Cohen, S. Gaubert, M. McGettrick, and J.-P. Quadrat. Numerical computation of spectral elements in max-plus algebra. In *Proc. IFAC Conf. on Syst. Structure and Control*, Nantes, France, 1998.
- [7] J. Cong, O. Coudert, and M. Sarrafzadeh. Incremental CAD. In *ICCAD*, 2000.
- [8] J. Cong, T. Kong, and D. Z. Pan. Buffer block planning for interconnect-driven floorplanning. In *ICCAD*, pages 358–363, 1999.
- [9] J. Cong, H. Li, and C. Wu. Simultaneous circuit partitioning/clustering with retiming for performance optimization. In *DAC*, pages 460 – 465, 1999.
- [10] T. H. Cormen, C. E. Leiserson, and R. H. Rivest. *Introduction to Algorithms*. MIT Press, 1989.
- [11] A. Dasdan, S. S. Irani, and R. K. Gupta. Efficient algorithms for optimum cycle mean and optimum cost to time ratio. In *DAC*, pages 37–42, 1999.
- [12] B. A. Davey and H. A. Priestley. *Introduction to Lattices and Order*. Cambridge, 1990.
- [13] S. Hassoun and C. J. Alpert. Optimal path routing in single and multiple clock domain systems. In *ICCAD*, pages 247–253, 2002.
- [14] C. E. Leiserson, F. M. Rose, and J. B. Saxe. Optimizing Synchronous Circuitry by Retiming. In *Advanced Research in VLSI: Proc. of the Third Caltech Conf.*, pages 86–116. Computer Science Press, 1983.
- [15] C. Lin and H. Zhou. Optimal wire retiming without binary search. In *ICCAD*, pages 452–458, 2004.
- [16] C. Lin and H. Zhou. Wire retiming for system-on-chip by fixpoint computation. In *DATE*, pages 1092–1097, 2004.
- [17] V. Nookala and S. S. Sapatnekar. A method for correcting the functionality of a wire-pipelined circuit. In *DAC*, pages 570–575, 2004.
- [18] P. Pan, A. K. Karandikar, and C. L. Liu. Optimal clock period clustering for sequential circuits with retiming. *IEEE TCAD*, 17(6):489–498, June 1998.
- [19] D. K. Y. Tong and E. F. Y. Young. Performance-driven register insertion in placement. In *ISPD*, pages 53–60, 2004.
- [20] H. Zhou and C. Lin. Retiming for wire pipelining in system-on-chip. *IEEE TCAD*, 23(9):1338–1345, September 2004.