

Optimality for Staggered Scheduling

Kwangil Ko*,
Thomas Robertazzi*
and Larry Wittie**

*Dept. of Electrical and Computer
Engineering,

**Dept. of Computer Science,
Stony Brook University,
Stony Brook, NY 11794

e-mail: tom@ece.sunysb.edu

Abstract —

Staggered scheduling is compared to “arbitrary” scheduling using matrix based techniques. It is shown that staggered scheduling always has a smaller solution time compared to any other possible arbitrary schedule for the same sequence of load distribution and solution reporting. This work is significant in using the arbitrary schedule concept in much the same way as the general distribution is used in queueing theory.

I. INTRODUCTION

Over the past 15 years there has been a good amount of research [1],[2] on divisible load scheduling. A divisible load is a computation and communication load that can be arbitrarily partitioned among processors and links. Such loads arise in data intensive computing, grid computing and metacomputing. A generally linear theory has been developed that provides tractable analytic solutions for optimal load allocation for a given scheduling strategy and interconnection network.

In this paper staggered scheduling is compared to “arbitrary” scheduling using matrix based techniques. Here staggered scheduling means that a computational load is well scheduled so as to minimize communication delay (this is elaborated on in the paper). It is proved that staggered scheduling always has a smaller solution time compared to any other possible arbitrary schedule for the same sequence of load distribution and solution reporting. This work is novel and significant in using the arbitrary schedule concept for scheduling, in much the same way the general distribution is used in queueing theory. It is also novel in its use of a matrix based analytical methodology. A final novel feature of this work is the optimization of solution reporting sequence back to the load originating processor.

An earlier work on divisible load scheduling optimality is [3].

This work is done in the context of divisible load modeling in a single level tree network with sequential load distribution and sequential solution reporting.

II. ARBITRARY LOAD SCHEDULING

In this section load scheduling which is not necessarily optimal, called arbitrary load scheduling, is discussed. The computation of finish time is obtained when arbitrary load is distributed to each processor. In arbitrary load scheduling, the load is not scheduled to each processor but allocated arbitrarily; that is, we allow any possible load fraction allocation to

processors as long as a specified sequence of load distribution and solution reporting is kept. We use the word “arbitrarily” in the same sense as a M/G/1 queue has a “general” solution. That is, any scheduling policy may be substituted for an arbitrary scheduling and the results still hold. The following notation will be used throughout this paper:

- p_i : The i^{th} processor.
- α_i : The fraction of the entire processing load assigned to the i^{th} processor.
- w_i : A constant inversely proportional to the computation speed of the i^{th} processor.
- z_i : A constant inversely proportional to the channel speed of the i^{th} link
- T_{cp} : Computing intensity constant. The entire load is processed in $w_i T_{cp}$ seconds by the i^{th} processor.
- T_{cm} : Communication intensity constant. The entire load can be transmitted in $z_i T_{cm}$ seconds over the i^{th} link.
- T_{cm}^{sol} : Solution reporting communication intensity constant. The entire solution report can be transmitted in $z_i T_{cm}^{sol}$ seconds over the i^{th} link.

Assume that there are $(M+1)$ processors named as follows:

$$\mathbf{P} = \{p_0, p_1, p_2, \dots, p_i, \dots, p_M\}$$

An arbitrary load distribution is defined as an $(M+1)$ tuple, $\tilde{\alpha}$, given by the fractions of the total load allocated to each processor, p_i for $i = 0, 1, 2, \dots, M$ and selected arbitrarily.

$$\tilde{\alpha} = (\tilde{\alpha}_0, \tilde{\alpha}_1, \tilde{\alpha}_2, \dots, \tilde{\alpha}_i, \dots, \tilde{\alpha}_M) \quad (1)$$

Here, the tilde is used to denote an arbitrary schedule. Further, the load is normalized.

$$\sum_{i=0}^M \tilde{\alpha}_i = 1 \quad (2)$$

Fig. 1 shows a timing diagram for arbitrary load distribution. The computation takes place over a single level tree network where each processor is equipped with a front-end processor for communication off-loading. The use of front end processors for communication off-loading means that communication and computation can proceed simultaneously for each processor. In the diagram computation time appears above each axis and communication time appears below each axis.

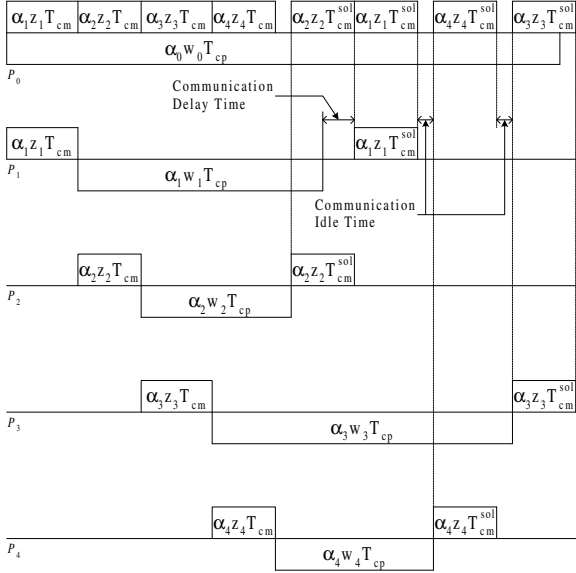


Figure 1: Timing Diagram for Arbitrary Scheduling of Homogeneous Single Level Tree Network.

Without loss of generality, it is assumed that load is distributed from originator (root) processor, p_0 . As soon as each child processor completely receives its assigned segment of the load it begins processing that segment. After each processor finishes this computation it requires the use of its link for a period of time to broadcast the result of its computation back to the originator. It is assumed that the scheduling must be done in such a way that only one processor transmits at one time and processors have sufficient time to perform their computations before transmitting their results. Thus, in the diagram processor 1's solution reporting is delayed until processor 2 reports its solution. Also note the presence of idle time after processor 1 reports its solution as processors 3 and 4 are still computing.

The originator processor, p_0 distributes fractions to M processors from p_1 to p_M in sequence. For simplicity we do not consider multi-installment load distribution strategies as in [2]. Processors start to process their fractions as soon as processors receive their fraction. For example, the i^{th} processor waits until the previous processors receive their fractions, and then starts to receive its fraction, $\tilde{\alpha}_i$. After this reception, the i^{th} processor can process its fraction.

First, it is necessary to calculate the time that it takes for the i^{th} processor to receive the fraction of the arbitrarily scheduled load, $\tilde{\alpha}_i$, and to process this load. That time is defined as \tilde{T}_{rp,p_i} (the "rp" stands for receive and process).

$$\tilde{T}_{rp,p_i} = \sum_{j=1}^i \tilde{\alpha}_j z_j T_{cm} + \tilde{\alpha}_i w_i T_{cp} \quad (3)$$

Note that the above expressions do not include solution reporting time. The first term and the second term in equation (3) are the reception time and computation time for the i^{th} processor, respectively.

However, processors finish processing in some arbitrary order as shown in Fig. 1. The reporting order is not predetermined but depends on \tilde{T}_{rp,p_i} for $i = 1, 2, \dots, M$. The solution for the i^{th} processor can be sent to the originator proces-

sor after \tilde{T}_{rp,p_i} . The processors are assumed to send their solutions back to the root in the order of finishing their work under arbitrary scheduling. Thus, sorting $\{\tilde{T}_{rp,p_i} : i = 1, 2, \dots, M\}$ is needed to determine the solution reporting order in the case of arbitrary load scheduling. The result of ascending order sorting $\{\tilde{T}_{rp,p_i} : i = 1, 2, \dots, M\}$ can be represented by:

$$\tilde{\mathbf{T}}_{rp}^{(sort)} = \{\tilde{T}_{rp}^{(1)}, \dots, \tilde{T}_{rp}^{(i)}, \dots, \tilde{T}_{rp}^{(j)}, \dots, \tilde{T}_{rp}^{(M)}\} \quad (4)$$

Here for $i < j$:

$$\tilde{T}_{rp}^{(i)} \leq \tilde{T}_{rp}^{(j)} \quad (5)$$

Further, $\tilde{\alpha}^{(sort)}$ is the version of the sorting $\{\tilde{\alpha}_i : i = 1, 2, \dots, N\}$ with respect to $\tilde{\mathbf{T}}_{rp}^{(sort)}$.

$$\tilde{\alpha}^{(sort)} = \{\tilde{\alpha}^{(1)}, \dots, \tilde{\alpha}^{(i)}, \dots, \tilde{\alpha}^{(j)}, \dots, \tilde{\alpha}^{(M)}\} \quad (6)$$

Consequently, $\tilde{\alpha}^{(i)}$ is the load fraction of the entire processing load assigned to $p^{(i)}$. Here $p^{(i)}$ is defined as the i^{th} reporting processor. Also, $z^{(i)}$ is the inverse speed of the link connected to $p^{(i)}$. Then the following property holds:

$$\sum_{i=1}^M \tilde{\alpha}^{(i)} z^{(i)} = \sum_{i=1}^M \tilde{\alpha}_i z_i \quad (7)$$

The reporting order is that $p^{(i)}$ is followed by $p^{(i+1)}$ for $i = 1, 2, \dots, M-1$. In Fig. 1, the reporting order of processors is p_2, p_1, p_4 , and then p_3 .

"Finish time" is the time for the system to complete processing its entire load. Naturally, in general finish time is equal to the time when the originator processor finishes processing or the last reporting processor completes processing and reporting.

$$\tilde{T}_f(M) = \max [\tilde{\alpha}_0 w_0 T_{cp} \quad \tilde{T}_s^{(M)}] \quad (8)$$

The term, $\tilde{\alpha}_0 w_0 T_{cp}$ is the processing time of the originator processor. Here, $\tilde{\alpha}_0$ is arbitrary generated. The term, $\tilde{T}_s^{(M)}$, is the time taken by the last reporting processor to complete reporting.

In this section, we develop an expression for the finish time of arbitrary scheduling. Now, $\tilde{T}_s^{(0)}$ is the time taken by the originator processor to distribute load fractions to all of the processors.

$$\tilde{T}_s^{(0)} = \sum_{j=1}^M \tilde{\alpha}_j z_j T_{cm} \quad (9)$$

Also, $\tilde{T}_s^{(i)}$ is the time taken by the i^{th} reporting order processor to complete sending its solution according to this reporting order. Here $p^{(i)}$, the i^{th} reporting order processor, tries to report its solution after $\tilde{T}_{rp}^{(i)}$. The i^{th} reporting order processor can immediately send its solution if $p^{(i-1)}$ already finished reporting. However, if $p^{(i-1)}$ didn't finish reporting its solution yet, $p^{(i)}$ should wait until $p^{(i-1)}$ finishes. Consequently, $p^{(i)}$ can report its solution after a time equal to the maximum of $\tilde{T}_s^{(i-1)}$ and $\tilde{T}_{rp}^{(i)}$. For $i = 1, 2, \dots, M$:

$$\tilde{T}_s^{(i)} = \max [\tilde{T}_s^{(i-1)}, \tilde{T}_{rp}^{(i)}] + \tilde{\alpha}^{(i)} z^{(i)} T_{cm}^{sol} \quad (10)$$

The last term in the above equation is the time taken by $p^{(i)}$ to report. Note that the first reporting processor can't send

its solution while the originator processor distributes fractions to children processors.

Let G_1 be the time difference between the instant of the start of the transmission of the solution of the first reporting processor, $\tilde{T}_{rp}^{(1)}$ and the instant that the originator processor finishes distributing load fractions, $\tilde{T}_s^{(0)}$. Also let G_i be the time difference between $\tilde{T}_{rp}^{(i)}$, the instant of the start of transmission of the solution of the i^{th} reporting processor, and $\tilde{T}_s^{(i-1)}$, the instant when the transmission of the solution of the $(i-1)^{th}$ reporting processor is complete. For $i = 2, 3, \dots, M$:

$$G_i = \tilde{T}_{rp}^{(i)} - \tilde{T}_s^{(i-1)} \quad (11)$$

It is assumed that the root can transmit or receive with only one child processor at a time. In the case that $p^{(i)}$ were to begin to report its solution while $p^{(i-1)}$ is transmitting its solution, $(\tilde{T}_{rp}^{(i)} < \tilde{T}_s^{(i-1)})$, a conflict would occur. Also all links are idle when $p^{(i)}$ cannot immediately send its solution, after $p^{(i-1)}$ has already transmitted $\tilde{\alpha}^{(i-1)}$, $(\tilde{T}_{rp}^{(i)} > \tilde{T}_s^{(i-1)})$. A negative G_i indicates a conflict and a positive G_i indicates a time when all of the links are idle. In the case of a conflict $(G_i < 0 \text{ or } \tilde{T}_{rp}^{(i)} < \tilde{T}_s^{(i-1)})$, $p^{(i)}$ can immediately report its solution as soon as $p^{(i-1)}$ finishes transmitting its solution. Thus equation (10) is represented as follows:

$$\tilde{T}_s^{(i)} = \tilde{T}_s^{(i-1)} + \tilde{\alpha}^{(i)} z^{(i)} T_{cm}^{sol} \quad (12)$$

In the case of idleness $(G_i > 0 \text{ or } \tilde{T}_{rp}^{(i)} > \tilde{T}_s^{(i-1)})$, equation (10) is represented as follow:

$$\tilde{T}_s^{(i)} = \tilde{T}_{rp}^{(i)} + \tilde{\alpha}^{(i)} z^{(i)} T_{cm}^{sol} \quad (13)$$

Using equation (11) the last two equations can be combined as follows:

$$\tilde{T}_s^{(1)} = \tilde{T}_s^{(0)} + \tilde{\alpha}^{(1)} z^{(1)} T_{cm}^{sol} + G_1 \text{sign}(G_1) \quad (14)$$

⋮

$$\tilde{T}_s^{(i)} = \tilde{T}_s^{(i-1)} + \tilde{\alpha}^{(i)} z^{(i)} T_{cm}^{sol} + G_i \text{sign}(G_i) \quad (15)$$

$$\tilde{T}_s^{(i+1)} = \tilde{T}_s^{(i)} + \tilde{\alpha}^{(i+1)} z^{(i+1)} T_{cm}^{sol} + G_{i+1} \text{sign}(G_{i+1})$$

⋮

$$\tilde{T}_s^{(M)} = \tilde{T}_s^{(M-1)} + \tilde{\alpha}^{(M)} z^{(M)} T_{cm}^{sol} + G_M \text{sign}(G_M) \quad (17)$$

Here $\text{sign}(x)$ is defined as follows:

$$\text{sign}(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases} \quad (18)$$

In equation (15), the time for the i^{th} reporting processor to finish reporting its solution includes three terms. One is the time for the $(i-1)^{th}$ reporting processor to finish reporting. The second term is the time for the i^{th} reporting processor to report the solution. The last term is the potential idle time of reporting between $p^{(i)}$'s and $p^{(i-1)}$'s reports.

The solution of the recursion of equation (14) to equation (17), $\tilde{T}_s^{(M)}$, can be written as follows:

$$\tilde{T}_s^{(M)} = \tilde{T}_s^{(0)} + \sum_{i=1}^M \tilde{\alpha}^{(i)} z^{(i)} T_{cm}^{sol} + \sum_{i=1}^M G_i \text{sign}(G_i) \quad (19)$$

Substituting equation (7) and (9) into (19), $\tilde{T}_f^{(M)}$ can be rewritten as follows:

$$\tilde{T}_s^{(M)} = \sum_{i=1}^M \tilde{\alpha}^{(i)} z^{(i)} T_{cm} + \sum_{i=1}^M \tilde{\alpha}^{(i)} z^{(i)} T_{cm}^{sol} + \sum_{i=1}^M G_i \text{sign}(G_i) \quad (20)$$

The above equation expresses the time at which all processors finish reporting their solution to the originator processor. Note that $\tilde{T}_f(M)$, defined as the finish time for arbitrary load scheduling, is the maximum of the processing time of the originator processor and the completion time of reporting all solutions. The processing time of the originator processor is determined by $\tilde{\alpha}_0$ generated arbitrary and the completion time of reporting all solutions is obtained from equation (20). Thus, $\tilde{T}_f(M)$ from equation (8) can be calculated.

Next, our attention turns to developing the concept of staggered scheduling.

III. STAGGERED SCHEDULING

The basic idea of staggered scheduling is that $p^{(i)}$, the i^{th} reporting processor, begins to transmit its solution immediately after the transmission of the solution of $p^{(i-1)}$, the $(i-1)^{th}$ reporting processor for $i = 2, 3, \dots, M$ while the originator (root) processor finishes its processing as soon as receiving the solution of the last reporting order processor, $p^{(M)}$.

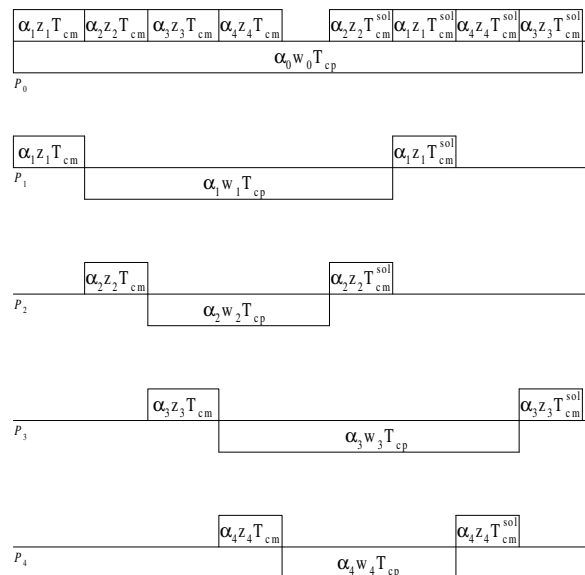


Figure 2: Timing Diagram for Staggered Scheduling of Homogeneous Single Level Tree Network.

Fig. 2 is a timing diagram of staggered scheduling. From equation (19) the first term and $G_1 \text{sign}(G_1)$ can be represented as follows:

$$\begin{aligned} T_s^{(0)} + G_1 \text{sign}(G_1) &= \tilde{T}_s^{(0)} + (T_{rp}^{(1)} - \tilde{T}_s^{(0)}) \cdot \text{sign}(T_{rp}^{(1)} - \tilde{T}_s^{(0)}) \quad (21) \end{aligned}$$

Again, $T_s^{(0)}$ is the time taken by the originator processor to distribute load fractions to all of the processors so as to create a well posed problem. The originator processor can distribute fractions before the first solution is received $(T_{rp}^{(1)} \geq \tilde{T}_s^{(0)})$.

The above equation can then be reduced to $T_{rp}^{(1)}$. This can be done since $T_{rp}^{(1)}$ is greater than $\tilde{T}_s^{(0)}$ for staggered scheduling. Thus equation (19) can be rewritten as follows using equation (7):

$$\tilde{T}_s^{(M)} = \tilde{T}_{rp}^{(1)} + \sum_{j=1}^M \tilde{\alpha}_j z_j T_{cm}^{sol} + \sum_{j=2}^M G_j \text{sign}(G_j) \quad (22)$$

If reporting solutions back to the originator processor is staggered so that there is no idle time, then the last term of the above equation is zero.

$$T_s^{(M)} = T_{rp}^{(1)} + \sum_{j=1}^M \alpha_j z_j T_{cm}^{sol} \quad (23)$$

Note we use α_i (no tilde) to denote a staggered scheduled load fraction.

In addition, the originator processor finishes its processing as soon as it receives the solution of $p^{(M)}$, the last reporting processor.

$$\begin{aligned} T_s^{(M)} &= T_f(M) \\ &= \alpha_0 w_0 T_{cp} \end{aligned} \quad (24)$$

Here, $\tilde{\alpha}_i$ for $i = 0, 1, \dots, M$ should be adjusted for staggering.

$$\alpha_i = \tilde{\alpha}_i - \alpha_i' \quad (25)$$

Again, $\tilde{\alpha}_i$ is the arbitrary scheduling load fraction for the i^{th} processor and now α_i is the staggered scheduling load fraction for the i^{th} processor. The normalization equation of load fractions for staggered scheduling is applied.

$$\sum_{j=0}^M \alpha_j = 1 \quad (26)$$

From the two normalization equation, equation (2) and equation (26), the following equation is obtained.

$$\sum_{j=0}^M \alpha_j' = 0 \quad (27)$$

The problem is to find α_i for $i = 0, 1, \dots, M$ so that G_i for $i = 2, 3, 4, \dots, M$ equals zero, which means that all processors are staggered. This problem is addressed in the next section.

Here, $M!$ staggered scheduling cases exist because of the reporting sequence. In this section, the staggered finish time for all cases will be obtained. The smallest one of the $M!$ staggered finish times is the optimal schedule. Each child processor has four time phases; receiving, processing, reporting the solution, and resting after reporting solution until the originator processor receives all solutions. The processor which sends its solution last has no rest time. Thus, the sum of the four time phases is equal to the finish time for any $i = 1, 2, \dots, M$ and is represented as follows:

$$\begin{aligned} T_f(M) &= \alpha_0 w_0 T_{cp} \\ &= T_s^{(M)} \end{aligned} \quad (28)$$

For any $i = 0, 1, \dots, M$:

$$\begin{aligned} T_s^{(M)} &= \sum_{j=1}^i \alpha_j z_j T_{cm} + \alpha_i w_i T_{cp} \\ &+ r_{i,i} \alpha_i z_i T_{cm}^{sol} + \sum_{j=1, j \neq i}^M r_{i,j} \alpha_j z_j T_{cm}^{sol} \end{aligned} \quad (29)$$

Here, we suppress that the dependence of $T_s^{(M)}$ on i as $T_s^{(M)}$ is the same for all i . The first two terms of equation (29) are the time taken by the i^{th} processor to receive its fraction and then to process its fraction, respectively. The processor reports its solution back to the originator (root) immediately after processing (third term). The last term is the resting time until all processors complete processing and reporting their solutions to the root. The $r_{i,j}$ is defined as a resting coefficient. Here, $r_{i,j}$ is one if the i^{th} processor is resting when the j^{th} processor reports its solution. Further $r_{i,j}$ is zero if the i^{th} processor is processing. The coefficient $r_{i,j}$ can be determined from the solution reporting order of processors and, again, could be zero or one. If the i^{th} processor reports its solution before the j^{th} processor, $r_{i,j}$ is one and $r_{j,i}$ is zero. If the i^{th} processor reports its solution after the j^{th} processor, $r_{i,j}$ is zero and $r_{j,i}$ is one. Thus, one of the $r_{i,j}$ and $r_{j,i}$ is one, and the other zero. Also the i^{th} processor should wait until it sends its own solution so $r_{i,i} = 1$ in the third term of equation (29). If the i^{th} processor is the k^{th} solution reporting processor, then the following equation holds:

$$\sum_{j=1}^M r_{j,i} = k \quad (30)$$

$$= S_i \quad (31)$$

Here, S_i is the reporting order of the i^{th} processor. The properties of the resting coefficient are summarized as follows:

1. $r_{i,j} = 1$ or 0
2. $r_{i,j} + r_{j,i} = 1$
3. $r_{i,j} \cdot r_{j,i} = 0$
4. $(r_{i,j})^2 = r_{i,j}$
5. $r_{i,j} = 1$ if $j = i$

Equation (28) can be rewritten as follows for $i = 1, 2, \dots, M$

$$\begin{aligned} & -\alpha_0 w_0 T_{cp} + \sum_{j=1}^i \alpha_j z_j T_{cm} + \alpha_i w_i T_{cp} \\ & + \alpha_i z_i T_{cm}^{sol} + \sum_{j=1, j \neq i}^M r_{i,j} \alpha_j z_j T_{cm}^{sol} \\ & = 0 \end{aligned} \quad (32)$$

Equation (32) for $i = 1, 2, \dots, M$ and (26) form $(M+1)$ linear equations.

$$[\mathbf{D} + \mathbf{W} + \mathbf{R}^{(k)}] \boldsymbol{\alpha} = \mathbf{E} \quad (33)$$

Here,

$$\boldsymbol{\alpha} = [\alpha_0 \quad \alpha_1 \quad \dots \quad \alpha_{M-1} \quad \alpha_M] \quad (34)$$

$$\mathbf{E} = [T_{cp} \quad 0 \quad 0 \quad \dots \quad 0] \quad (35)$$

$$\mathbf{D} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & z_1 & \dots & 0 \\ 0 & \vdots & \ddots & \vdots \\ 0 & z_1 & \dots & z_M \end{bmatrix} T_{cm} \quad (36)$$

$$\mathbf{W} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ -w_0 & w_1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ -w_0 & 0 & \dots & w_M \end{bmatrix} T_{cp} \quad (37)$$

$$\mathbf{R}^{(k)} = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & z_1 & \cdots & r_{1,M}z_M \\ 0 & \vdots & \ddots & \vdots \\ 0 & r_{M,1}z_1 & \cdots & z_M \end{bmatrix} T_{cm}^{sol} \quad (38)$$

Here, $r_{i,j}$ is for $i, j = 0, 1, 2, \dots, M$, and the reporting sequence matrix, $\mathbf{R}^{(k)}$ for $k = 1, 2, \dots, M!$ reporting sequences are specified. Now α_i for $i = 0, 1, 2, \dots, M$ can be obtained using Cramer's rule.

$$\alpha_i = \frac{\det(\mathbf{E}_{i+1})}{\det[\mathbf{D} + \mathbf{W} + \mathbf{R}^{(k)}]} \quad (39)$$

Here \mathbf{E}_i is obtained from $[\mathbf{D} + \mathbf{W} + \mathbf{R}]$ by replacing the i^{th} column of $[\mathbf{A} + \mathbf{B} + \mathbf{C}]$ by \mathbf{E} for $i = 0, 1, \dots, M$. Then α_0 is:

$$\alpha_0 = \frac{\det(\mathbf{E}_1)}{\det[\mathbf{D} + \mathbf{W} + \mathbf{R}^{(k)}]} \quad (40)$$

Now the finish time for computing a unit load with M children processor is obtained as follows for any $k = 1, 2, \dots, M!$:

$$T_f(M) = \frac{\det(\mathbf{E}_1)}{\det[\mathbf{D} + \mathbf{W} + \mathbf{R}^{(k)}]} w_0 T_{cp} \quad (41)$$

There are $M!$ possible cases depending on the solution reporting sequence. In this section the finish time of each possible order has been obtained. It should be noted that for a specific reporting sequence recursive solutions for the staggered schedule load fractions can be developed as in [2].

IV. OPTIMALITY PROOF

In this section, both arbitrary scheduling and staggered scheduling are assumed to follow the same reporting sequence. We will show that if all of the processors are staggered for a *given* reporting sequence, the finish time of staggered scheduling is always less than the finish time of arbitrary scheduling with the same reporting sequence.

Let the hyperplane H_0 in $(M+1)$ -dimensions be a set of the form $\{\mathbf{x} : \mathbf{n} \cdot \mathbf{x} = 1\}$ where \mathbf{n} is a 1 by $(M+1)$ ones vector.

Definition 1 Any point on the hyperplane H_0 can be defined as $\tilde{\alpha}$ which is an arbitrary load distribution.

$$\{\tilde{\alpha} : \mathbf{n} \cdot \tilde{\alpha} = 1\} \quad (42)$$

Here:

$$\tilde{\alpha} = (\tilde{\alpha}_0, \tilde{\alpha}_1, \tilde{\alpha}_2, \dots, \tilde{\alpha}_i, \dots, \tilde{\alpha}_M) \quad (43)$$

The set of $\{\tilde{\alpha} : \mathbf{n} \cdot \tilde{\alpha} = 1\}$ is a normalization equation.

Consider staggered scheduling. There exist $M!$ cases depending on reporting sequences. Let $\mathbf{Q}^{(k)}$ for $k = 1, 2, \dots, M!$ be $[\mathbf{D} + \mathbf{W} + \mathbf{R}^{(k)}]$ in the linear equation (33). There exists a unique solution (through normalization) on the hyperplane H_0 given reporting sequence, $\mathbf{R}^{(k)}$, for any $k = 1, 2, \dots, M!$ satisfying equation (33).

Definition 2 In the $(M+1)$ -dimensional space, the unique solution of the linear equation (33) is defined as α^k for any $k = 1, 2, \dots, M!$.

$$\mathbf{Q}^{(k)} \cdot \alpha^k = \mathbf{E} \quad (44)$$

Here:

$$\alpha^k = (\alpha_0^k, \alpha_1^k, \dots, \alpha_M^k) \quad (45)$$

Specifically, α^k is staggered scheduling for a given reporting sequence, $\mathbf{R}^{(k)}$. Thus α^k is the fraction of load for the staggered scheduling of the $\mathbf{Q}^{(k)}$ matrix.

Definition 3 Let $\bar{\alpha}$ be any point in a set of the form:

$$\{\bar{\alpha} : \bar{\alpha}_0 \neq \alpha_0^k \text{ and } \bar{\alpha}_i^k = m\alpha_i^k \text{ for all } i = 1, 2, \dots, M\} \quad (46)$$

Thus $\bar{\alpha}$ can be defined as follows:

$$\bar{\alpha} = (\bar{\alpha}_0, m\alpha_1^k, \dots, m\alpha_i^k, \dots, m\alpha_M^k) \quad (47)$$

Here, m is any arbitrary real number and $0 < m < \frac{1}{1-\alpha_0^k}$.

The upper limit of m is chosen so that the $\bar{\alpha}$ is on a hyperplane. Each load fraction of the children processors in the load distribution $\bar{\alpha}$, is proportional to the load fraction of the children processors in the load distribution α^k . Thus, the reporting time of all of the children processors with load distribution $\bar{\alpha}$, is also staggered. In the load distribution, $\bar{\alpha}$, it is not necessary that the computation of the originator processor and reporting of the last reporting processor finish at the same time.

Lemma 1 For any reporting sequence $\mathbf{R}^{(k)}$ for $k = 1, 2, \dots, M!$, staggered scheduling has a smaller finish time than scheduling with load distribution, $\bar{\alpha}$.

Lemma 1 can be written as follow equation: The finish time of scheduling with $\bar{\alpha}$ can be written as follows:

$$\begin{aligned} \bar{T}_f(M) &\geq mT_s^{(M)} \\ &> T_s^{(M)} = T_f(M) \end{aligned} \quad (48)$$

Here, $\bar{T}_f(M)$ is the finish time of scheduling with $\bar{\alpha}$, and $T_f(M)$ is the finish time of staggered scheduling.

Any point on the hyperplane H_0 can be denoted as $\tilde{\alpha}$ which is an arbitrary load distribution.

$$\tilde{\alpha} = (\tilde{\alpha}_0, \tilde{\alpha}_1, \tilde{\alpha}_2, \dots, \tilde{\alpha}_i, \dots, \tilde{\alpha}_M) \quad (49)$$

Let $\tilde{\mathbf{E}}$ be the result from multiplying $\mathbf{Q}^{(k)}$ by $\tilde{\alpha}$

$$\mathbf{Q}^{(k)} \cdot \tilde{\alpha} = \tilde{\mathbf{E}} \quad (50)$$

Here $\tilde{\mathbf{E}}$ can be expressed as follows:

$$\tilde{\mathbf{E}} = [T_{cp} \quad \tilde{e}_1 \quad \tilde{e}_2 \quad \cdots \quad \tilde{e}_M]' \quad (51)$$

Subtracting both side of equation (44) from (50)

$$\mathbf{Q}^{(k)} \cdot (\tilde{\alpha} - \alpha) = \tilde{\mathbf{E}} - \mathbf{E} \quad (52)$$

or

$$\mathbf{Q}^{(k)} \cdot \hat{\alpha} = \hat{\mathbf{E}} \quad (53)$$

Here:

$$\hat{\alpha} = \tilde{\alpha} - \alpha \quad (54)$$

and

$$\hat{\mathbf{E}} = \tilde{\mathbf{E}} - \mathbf{E} \quad (55)$$

$$= [0 \quad \hat{e}_1 \quad \hat{e}_2 \quad \cdots \quad \hat{e}_M]' \quad (56)$$

Lemma 2 If two distributions, $\tilde{\alpha}$ and $\bar{\alpha}$ are on the hyperplane then at least, one of \hat{e}_i for $i = 1, 2, \dots, M$ is greater than or equal to zero.

Lemma 3 In a given sequence, if $\bar{\alpha}_0$ equals to $\tilde{\alpha}_0$ then the load distribution with $\bar{\alpha}$, has a smaller finish time than arbitrary load distribution, $\tilde{\alpha}$.

Proof.

From Lemma 2, assume that \hat{e}_i for any $i = 1, 2, \dots, M$ is greater than or equal to zero. Also from equation (53), \hat{e}_i can be expressed as multiplying $(\tilde{\alpha} - \bar{\alpha})$ by the $(i+1)^{th}$ row of $\mathbf{Q}^{(k)}$.

$$\begin{aligned} \hat{e}_i &= \mathbf{q}_{i+1}^{(k)} \cdot \tilde{\alpha} - \mathbf{q}_{i+1}^{(k)} \cdot \bar{\alpha} \\ &\geq 0 \end{aligned} \quad (57)$$

Here, $\mathbf{q}_{i+1}^{(k)}$ is defined as the $(i+1)^{th}$ row of $\mathbf{Q}^{(k)}$. The second term on the right side of the above equation can be expressed as equation (33).

$$\begin{aligned} \mathbf{q}_{i+1}^{(k)} \cdot \bar{\alpha} &= -\bar{\alpha}_0 w_0 T_{cp} + m \sum_{j=1}^i \alpha_j^k z_j T_{cm} \\ &\quad + m \alpha_i^k w_i T_{cp} + m \alpha_i^k z_i T_{cm}^{sol} \\ &\quad + m \sum_{j=1, j \neq i}^M r_{i,j} \alpha_j^k z_j T_{cm}^{sol} \end{aligned} \quad (58)$$

Similarly, $\mathbf{q}_i^{(k)} \cdot \tilde{\alpha}$ can be expressed as follows:

$$\begin{aligned} \mathbf{q}_i^{(k)} \cdot \tilde{\alpha} &= -\tilde{\alpha}_0 w_0 T_{cp} + \sum_{j=1}^i \tilde{\alpha}_j z_j T_{cm} \\ &\quad + \tilde{\alpha}_i w_i T_{cp} + \tilde{\alpha}_i z_i T_{cm}^{sol} \\ &\quad + \sum_{j=1, j \neq i}^M r_{i,j} \tilde{\alpha}_j z_j T_{cm}^{sol} \end{aligned} \quad (59)$$

Thus:

$$\begin{aligned} \hat{e}_i &= -\tilde{\alpha}_0 w_0 T_{cp} + \bar{\alpha}_0 w_0 T_{cp} \\ &\quad + \left[\sum_{j=1}^i \tilde{\alpha}_j z_j T_{cm} + \tilde{\alpha}_i w_i T_{cp} + \tilde{\alpha}_i z_i T_{cm}^{sol} \right. \\ &\quad \left. + \sum_{j=1, j \neq i}^M r_{i,j} \tilde{\alpha}_j z_j T_{cm}^{sol} \right] \\ &\quad - m \cdot \left[\sum_{j=1}^i \alpha_j^k z_j T_{cm} + \alpha_i^k w_i T_{cp} + \alpha_i^k z_i T_{cm}^{sol} \right. \\ &\quad \left. + \sum_{j=1, j \neq i}^M r_{i,j} \alpha_j^k z_j T_{cm}^{sol} \right] \end{aligned} \quad (60)$$

Note that \hat{e}_i is greater than or equal to zero and $\bar{\alpha}_0$ equals to $\tilde{\alpha}_0$ by assumption. Thus:

$$\begin{aligned} &\left[\sum_{j=1}^i \tilde{\alpha}_j z_j T_{cm} + \tilde{\alpha}_i w_i T_{cp} + \tilde{\alpha}_i z_i T_{cm}^{sol} \right. \\ &\quad \left. + \sum_{j=1, j \neq i}^M r_{i,j} \tilde{\alpha}_j z_j T_{cm}^{sol} \right] \end{aligned} \quad (61)$$

$$\begin{aligned} &\geq m \cdot \left[\sum_{j=1}^i \alpha_j^k z_j T_{cm} + \alpha_i^k w_i T_{cp} + \alpha_i^k z_i T_{cm}^{sol} \right. \\ &\quad \left. + \sum_{j=1, j \neq i}^M r_{i,j} \alpha_j^k z_j T_{cm}^{sol} \right] \end{aligned} \quad (62)$$

The above inequality may be summarized as follows:

$$\tilde{T}_s^{(M)} > \bar{T}_s^{(M)} \quad (63)$$

Recall that the finish time of arbitrary load scheduling:

$$\tilde{T}_f(M) = \max \left[\tilde{\alpha}_0 w_0 T_{cp}, \tilde{T}_s^{(M)} \right] \quad (64)$$

If the finish time with distribution $\bar{\alpha}$ is determined by $\bar{T}_s^{(M)}$:

$$\bar{T}_f(M) = \bar{T}_s^{(M)} \quad (65)$$

Then, the following result can be obtained from equation (63) and (64).

$$\tilde{T}_f(M) \geq \tilde{T}_s^{(M)} \geq \bar{T}_s^{(M)} = \bar{T}_f(M) \quad (66)$$

If the finish time with distribution $\bar{\alpha}$ is determined by $\bar{\alpha}_0 w_0 T_{cp}$:

$$\bar{T}_f(M) = \bar{\alpha}_0 w_0 T_{cp} \quad (67)$$

Then, from equation (64),

$$\tilde{T}_f(M) \geq \tilde{\alpha}_0 w_0 T_{cp} = \bar{\alpha}_0 w_0 T_{cp} = \bar{T}_f(M) \quad (68)$$

Therefore, the finish time, $\bar{T}_f(M)$, is less than that of arbitrary scheduling.

$$\tilde{T}_f(M) \geq \bar{T}_f(M) \quad (69)$$

Theorem 4 *Staggered load scheduling has a smaller finish time than arbitrary load scheduling.*

Proof. From inequality (48) and (69) the following result is obtained:

$$T_f(M) \leq \bar{T}_f(M) \leq \tilde{T}_f(M) \quad (70)$$

V. CONCLUSION

This work is significant in using the arbitrary schedule concept for scheduling in much the same way that the general distribution is used in queueing theory. The proof methodology used here can be applied to other divisible load optimality problems.

REFERENCES

- [1] V. Bharawaj, D. Ghose and T.G. Robertazzi, "Divisible Load Theory: A New Paradigm for Load Scheduling in Distributed Systems," in special issue of *Cluster Computing on Divisible Load Scheduling* (D. Ghose and T. Robertazzi, editors), vol. 6, no. 1, Jan. 2003, pp. 7-17.
- [2] V. Bharadwaj, D. Ghose, V. Mani and T. G. Robertazzi, *Scheduling Divisible Loads in Parallel and Distributed Systems*, IEEE Computer Society Press (now distributed by Wiley), 1996.
- [3] J. Sohn and T.G. Robertazzi, "Optimal Load Sharing for a Divisible Job on a Bus Network," *IEEE Transactions on Aerospace and Electronic System*, vol. 32, no. 1, Jan. 1996, pp. 34-40.