

On Broadcasting in Heterogenous Networks ^{*}

Samir Khuller [†]

Yoo-Ah Kim [‡]

Abstract

In this paper we study a well known broadcasting heuristic for heterogenous networks of workstations, called *fastest node first*. We show that this heuristic produces an optimal solution for minimizing the sum of all completion times, and in addition produces a 1.5 approximation for the problem of minimizing the maximum completion time. We also develop a polynomial time approximation scheme for this problem. We extend these results to show that the same bounds can be obtained for the multicast operation on such heterogenous networks. In addition we show that the problem of minimizing the maximum completion time is *NP*-hard, which settles the complexity of this open problem.

1 Introduction

Networks of Workstations (NOWs) are an extremely popular alternative to massively parallel machines and are widely used (for example the Condor project at Wisconsin [20] and the Berkeley NOW project [19]). By simply using off-the-shelf PC's, a very powerful workstation cluster can be created, and this can provide a high amount of parallelism at relatively low cost. Since NOWs are put together over time, the machines tend to have different capabilities and this leads to a *heterogenous* collection of machines, rather than a *homogenous* collection, where all the machines have identical capabilities.

One fundamental operation that is used in such clusters, is that of *broadcast* (this is a primitive in many message passing systems such as MPI [1, 6, 8]). In addition it is used as a primitive in many parallel algorithms. The main objective of a broadcast operation is to quickly distribute the input data to the entire network for processing. Another situation is when the system is performing a highly parallel search, then the successful processor needs to inform all other processors that the search has concluded successfully. Various models for heterogenous environments have been proposed in the literature. One general model is the one proposed by Bar-Noy et al [3] where the communication costs between links are not uniform. In addition, the sender may engage in another communication before the current one is complete. An approximation factor with a guarantee of $O(\log k)$ is given for the operation of performing a multicast. Other popular models in the theory literature generally assume an underlying communication graph, with the property that only nodes adjacent in this graph may communicate.

Broadcasting efficiently is an essential operation and many works are devoted to this (see [21, 10, 11, 4, 5] and references therein). In addition, for emergency notification an understanding of how to perform broadcast quickly is essential.

Most of the algorithmic (theoretical) work done on such problems is of little interest to practitioners, because the approximation algorithms tend to be fairly complex and slow, thus defeating the purpose of performing broadcast quickly. On the other hand, a simple model and algorithm was proposed by

^{*}Research supported by NSF Awards CCR-9820965 and CCR-0113192. Preliminary versions of this work were published in [13] and [14].

[†]Department of Computer Science, University of Maryland, College Park, MD 20742. E-mail : samir@cs.umd.edu.

[‡]Department of Computer Science, University of Maryland, College Park, MD 20742. E-mail : ykim@cs.umd.edu.

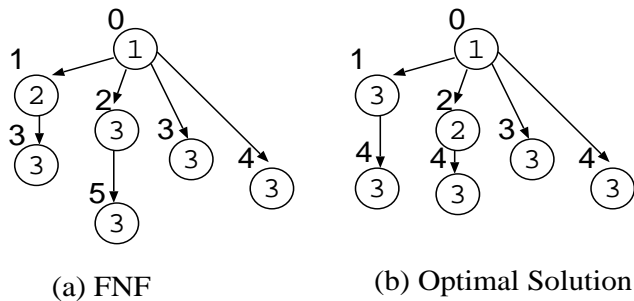


Figure 1: An example that FNF does not produce an optimal solution. Transmission times of processors are inside the circles. Times at which nodes receive a message are also shown.

Banikazemi et al [2]. In this model, heterogeneity among processors is modeled by a non-uniform speed of the sending processor. A heterogenous cluster is defined as a collection of processors p_1, p_2, \dots, p_n in which each processor is capable of communicating with any other processor. Each processor has a transmission time which is the time required to send a message to any other processor in the cluster. Thus the time required for the communication is a function of only the sender. Each processor may send messages to other processors in order, and each processor may be receiving only one message at a time.

Thus a broadcast operation is implemented as a broadcast tree. Each node in the tree represents a processor of the cluster. The root of the tree is the source of the original message. The children of a node p_i are the processors that receive the message from p_i . The completion time of a node is the time at which it completes receiving the message from its parent. The completion time of the children of p_i is $c_i + j \cdot t_i$, where c_i is the completion time of p_i , t_i is the transmission time of p_i and j is the child number. In other words, the first child of p_i has a completion time of $c_i + t_i$ ($j = 1$), the second child has a completion time of $c_i + 2t_i$ ($j = 2$) etc. See Figure 1 for an example.

A commonly used method to find a broadcast tree is referred to as the “Fastest Node First” (FNF) technique [2]. This works as follows: In each iteration, the algorithm chooses a sender from the set of processors that have received the message (set S) and a receiver from the set of processors that have not yet received the message (set R). The algorithm then picks the sender from $s \in S$ so that s can finish the transmission as early as possible, and chooses the receiver $r \in R$ as the processor with the minimum transmission time in R . Then r is moved from R to S and the algorithm continues. The intuition is that sending the message to fast processors first is a more effective way to propagate the message quickly. This technique is very effective and easy to implement. In practice it works extremely well (using simulations) and in fact frequently finds optimal solutions as well [2]. However, there are situations when this method also fails to find an optimal solution. A simple example is shown in Figure 1.

Despite several non-trivial advances in an understanding of the *fastest node first* method by Liu [16] (see also work by Liu and Sheng [18] in SPAA 2000), it was still open as to what the complexity of this problem is. For example, is there a polynomial time algorithm to compute a solution that minimizes the broadcast time? Can we show that in all instances the FNF heuristic will find solutions close to optimal? The second question is of interest, regardless of the complexity of the problem; for example, if there was a complex polynomial time algorithm to find the optimal solution, it is unlikely that it would be used in practice and an understanding of the FNF heuristic is still very useful.

Liu [16] (see also [18]) shows that if there are only two classes of processors, then FNF produces an optimal solution. In addition, if the transmission time of every slower processor is a multiple of the transmission time of every faster processor, then again the FNF heuristic produces an optimal solution.

So for example, if the transmission time of the fastest processor is 1 and the transmission time of all other processors are powers of 2, then the algorithm produces an optimal solution. It immediately follows that by rounding all transmission times to powers of 2 we can obtain a solution using FNF whose cost is at most twice the cost of an optimal solution. However, this still does not explain the fact that this heuristic does much much better in practice. Is there a better worst case guarantee on this simple method?

In this paper we show seven results.

- We show that the FNF heuristic actually produces an optimal solution for the problem of minimizing the sum of completion times (Section 3).
- We use this to show that the FNF method has a performance ratio of at most 1.5 when compared to the optimal solution for minimizing broadcast time (Section 4).
- We show that the performance ratio of FNF is at least $\frac{25}{22}$ ¹ (Section 5).
- As a corollary of the above approximation result, we are able to show that if the transmission times of the fastest $\frac{n}{2}$ processors are in the range $[1 \dots C]$ then FNF produces a solution with makespan at most $T_{OPT} + C$ (Section 4).
- We prove that there is a polynomial time approximation scheme (PTAS) for minimizing the broadcast time (Section 6).
- We extend the above results to the problem of multicasting (Section 7).
- We also establish that the problem is *NP*-hard, so unless $P = NP$ there is no polynomial time algorithm for finding an optimal solution (Section 8).

The problem to minimize broadcast time was shown to be NP-hard when each processor has an arbitrary release time by Hall et al. [9]. We prove that it is NP-hard even when all processors are available at the beginning. Hall et al. also proved that FNF minimizes the sum of completion times, and we provide a short proof for the problem. It was shown by Liu et al. [17] and Libeskind-Hadas et al. [15] independently that FNF gives a 2-approximation for minimizing broadcast time. We prove that FNF gives a 1.5-approximation. Libeskind-Hadas et al. also showed that a schedule with the minimum broadcast time can be obtained in $O(n^{2k})$ using dynamic programming, where k is the number of different sending speeds of processors.

We show that the FNF heuristic actually produces an optimal solution for the problem of minimizing the sum of completion times. (This itself is a little surprising, since for most problems, minimizing the sum of completion times is usually harder than minimizing makespan.) We use this property of the algorithm to show that it actually maximizes the number of “fractional blocks²” in the schedule. This property leads to a lower bound on the optimal schedule based on the FNF schedule. Using this lower bound we can prove a guarantee of 1.5. We also prove that there is a polynomial time approximation scheme (PTAS) for this problem. However, this algorithm is not practical due to its high running time, albeit polynomial.

¹Observe that in the example shown in Figure 1 the performance ratio of FNF is $\frac{5}{4}$, but this example does not appear to generalize easily to arbitrarily large instances.

²This concept will be explained later.

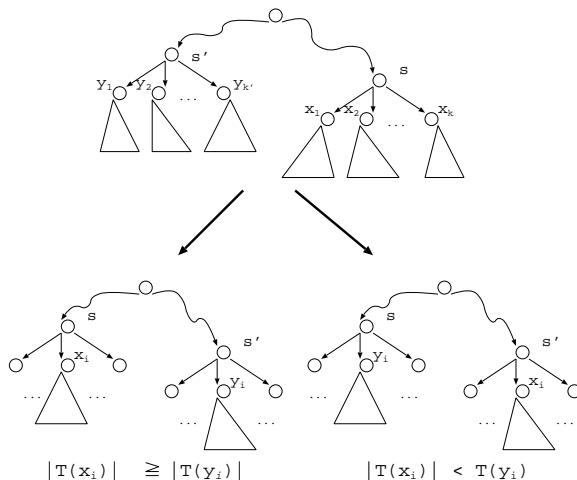


Figure 2: Figure shows how to modify the schedule of $T(s)$ and $T(s')$.

2 Problem Definition

We are given a set of processor p_0, p_1, \dots, p_n . There is one message to be broadcast from p_0 to all other processors p_1, p_2, \dots, p_n . Each processor p_i can send a message to another processor with transmission time t_i once it has received the message. Each processor can be either sending a message or receiving a message at any point of time. Without loss of generality, we assume that $t_1 \leq t_2 \leq \dots \leq t_n$ and $t_0 = 1$.

We define the completion time of processor p_i to be the time when p_i has received the message. Our objective is to find a schedule that minimizes $C_{\max} = \max_i c_i$ where c_i is the completion time of processor p_i . In other words, we want to find a schedule that minimizes the time required to send the message to all the processors.

Our proof makes use of the following results from [2].

THEOREM 2.1. [2] *There exists an optimal broadcast tree in which all processors send messages without delay.*

THEOREM 2.2. [2] *There exists an optimal broadcast tree in which every processor has transmission time no less than its parent (except the source node).*

3 Minimizing the Sum of Completion Times

In this section, we show that the FNF scheme finds an optimal schedule to minimize the sum of completion times of all the nodes, i.e., it minimizes $C_{\text{sum}} = \sum_i c_i$.

Suppose that we are given an (optimal) schedule that minimizes the sum of completion times and the sum of completion times is C . If there is a processor $p_i (i \geq 1)$ whose completion time is later than processor $p_j (i < j)$, i.e., $c_i > c_j$ in the schedule, then we can find a modified schedule such that processor $p_i (i \geq 1)$ has completion time no later than processor p_j and the sum of completion times is no more than C .

Let us define a permutation $\pi : \{1, 2, \dots, n\} \rightarrow \{1, 2, \dots, n\}$. Then any schedule can be represented as an ordered list of processors $(p_{\pi(1)}, p_{\pi(2)}, \dots, p_{\pi(n)})$ by sorting in non-decreasing order of completion times (ties broken in accordance with their indices).

We first prove a lemma similar to Theorem 2.2 for sum of completion times. The proof is the same as Theorem 2.2 [2]. We include the proof for the completeness.

LEMMA 3.1. *There exists a broadcast tree such that it minimizes the sum completion times and for any processor p_i other than p_0 , its children have the transmission times no less than p_i .*

Proof. Suppose that in an optimal solution for sum of completion times, there are processors p_i and p_j ($i, j \neq 0$) with $t_i < t_j$ and p_i is a child of p_j . Then we can modify the solution and show that the sum of completions times for the modified solution is smaller than the given optimal solution, which is a contradiction.

Let c_i and c'_i be the completion time of p_i in the given optimal solution and the modified solution, respectively. Let $P(c_i)$ denote a set of processors p_j such that $c_j \leq c_i$. We first modify the schedule *only for processors in $P(c_i)$* . The schedule is the same as the original schedule but p_i and p_j are exchanged. In other words, p_i receives the message in place of p_j and sends the message to children of p_j in $P(c_i)$. Clearly $c'_j < c_i$ and $c'_i = c_j$. For the rest of processors p_k in $P(c_i)$, c'_k is no later than c_k since p_i is faster than p_j .

For the processors not in $P(c_i)$, we can broadcast the message using the original schedule. Since $c'_k \leq c_k$ for all $p_k \in P(c_i) \setminus \{p_i, p_j\}$ and both p_i and p_j are available to send a message at time c_i , the modified completion times are no later than the original schedule. Therefore, we have $\sum_k c'_k < \sum_k c_k$, which is a contradiction. ■

LEMMA 3.2. *Let $(p_{\pi(1)}, p_{\pi(2)}, \dots, p_{\pi(n)})$ be an optimal schedule for the problem with sum of completion times C . Let s be the smallest index such that $\pi(s) \neq s$. Then we can find a schedule with sum of completion times no more than C and $i = \pi(i)$ for all $1 \leq i \leq s$.*

Proof. In the given optimal broadcast tree, let us call the node corresponding to processor p_i as node i . We denote the subtree rooted at node i as $T(i)$. Consider subtrees $T(s)$ and $T(s')$ where $s' = \pi(s)$. We know that s cannot be an ancestor of s' as $c_{s'} \leq c_s$. Also s' cannot be an ancestor of s as $t_s < t_{s'}$ by Lemma 3.1. Therefore, $T(s)$ and $T(s')$ are disjoint.

Let node s have children x_1, x_2, \dots, x_k and node s' have $y_1, y_2, \dots, y_{k'}$ as shown in Figure 2. We change the schedule as follows. First we exchange s and s' . In other words, the modified completion time of p_s becomes $c_{s'}$ and the completion time of $p_{s'}$ becomes c_s . Clearly, this does not increase the sum of completion times. For all i ($1 \leq i \leq \max(k, k')$), we compare the size of subtree $T(x_i)$ and $T(y_i)$ and attach the bigger one to s and the smaller one to s' as i -th child. (if there does not exist a child, simply consider the size of the subtree as zero)

We can prove that this modification does not increase the sum of completion times. The difference of the sum of completion times for subtree $T(x_i)$ and $T(y_i)$ depends on which parent they are attached to. In case that $|T(x_i)| \geq |T(y_i)|$, the completion times of processors in $T(x_i)$ are decreased by $c_s - c_{s'}$ since the completion time of p_{x_i} is changed from $c_s + i \cdot t_s$ to $c_{s'} + i \cdot t_s$. The completion times of processors in $T(y_i)$ are increased by $c_s - c_{s'}$ since the completion time of p_{y_i} is changed from $c_{s'} + i \cdot t_{s'}$ to $c_s + i \cdot t_{s'}$. Therefore the difference is

$$(c_s - c_{s'}) (|T(y_i)| - |T(x_i)|) \leq 0$$

In case that $|T(x_i)| < |T(y_i)|$, the completion times of processors in $T(x_i)$ are increased by $i \cdot (t_{s'} - t_s)$ because the completion of x_i is modified from $c_s + i \cdot t_s$ to $c_s + i \cdot t_{s'}$. and the completion times of processors in $T(y_i)$ are decreased by $i \cdot (t_{s'} - t_s)$ because the completion of x_i is modified from $c_{s'} + i \cdot t_{s'}$ to $c_{s'} + i \cdot t_s$. Therefore the difference is

$$i \cdot (t_{s'} - t_s) (|T(x_i)| - |T(y_i)|) \leq 0$$

Therefore we have the lemma. ■

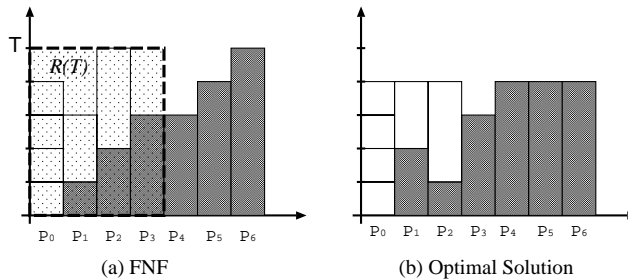


Figure 3: An example of Bar charts corresponding to the schedules created by the instance specified in Figure 1.

By repeating the procedure in Lemma 3.2, we can find a schedule such that processors receive the message in non-decreasing order of their transmission times and the sum of completion times is no more than the optimal. We thus conclude:

THEOREM 3.1. *Algorithm FNF minimizes the sum of completion times.*

In fact, we can prove even stronger result by applying the same procedure only to a subset of processors p_i ($1 \leq i \leq k$) for any $k \leq n$.

COROLLARY 3.1. *Algorithm FNF minimizes the sum of completion times over all processors p_i ($1 \leq i \leq k$) for any $k \leq n$.*

Proof. We do the same procedure as in Lemma 3.2 except that we only count processors p_i ($1 \leq i \leq k$) when we compute the sizes of subtrees. ■

We will use this corollary in the next section to prove that the FNF scheme gives a 1.5-approximation for minimizing broadcast time.

4 1.5-approximation for Minimizing Broadcast Time

In this section, we prove that FNF scheme gives 1.5-approximation.

Let us consider the bar chart as shown in Figure 3 (a), where processors are listed in non-decreasing order of transmission time in the horizontal line and each processor has a block whose height corresponds to its completion time in a schedule (These two charts correspond to the schedule created on the instance specified in Figure 1). We call these bars *grey blocks*. Each processor p_i can start sending messages as soon as it receives the message, and send to one processor in each t_i time unit. A *white block* corresponds to each message sent by p_i . Therefore, the height of a white block is the same as the transmission times of the corresponding processors.

DEFINITION 4.1. *We define the number of fractional blocks as the number of (fractional) messages a processor can send by the given time. In other words, given a time T , the number of fractional blocks of processor p_i is $(T - c_i)/t_i$.*

For example, the number of fractional blocks of p_1 by time T is 2 in the FNF schedule (see Figure 3(a)).

DEFINITION 4.2. *We define $R(T)$ as the rectangular region bounded by time 0 and T and including processors from p_0 to $p_{\lfloor n/2 \rfloor}$ in the bar chart.*

An example is shown in Figure 3(a).

We only include processors $p_0, p_1, \dots, p_{\lfloor n/2 \rfloor}$ in $R(T)$ because of the following lemma.

LEMMA 4.1. *There is an optimal schedule in which only the $\lfloor n/2 \rfloor$ fastest processors and the source processor send messages, that is, processor p_i ($\lfloor n/2 \rfloor + 1 \leq i \leq n$) need not send any messages.*

Proof. We prove this by showing that there is an optimal broadcast tree in which every internal node (except the source) has at least one child that is a leaf. Suppose an internal node s ($\neq p_0$) does not have a leaf child. Then we move the processor which receives the message last in subtree $T(s)$ to a child of s . It is easy to see that this does not increase the makespan of the schedule by Theorem 2.2. By repeatedly applying this modification to the given broadcast tree, we can find an optimal broadcast tree satisfying the property. ■

LEMMA 4.2. *Algorithm FNF maximizes the number of fractional blocks in $R(T)$ for any T .*

To prove this lemma, we first prove the following proposition.

PROPOSITION 4.1. $\sum_{i=1}^m a_i/b_i \geq \sum_{i=1}^m a'_i/b_i$ if $\sum_{i=1}^l a_i \geq \sum_{i=1}^l a'_i$ for all $1 \leq l \leq m$ and $0 < b_i \leq b_{i+1}$ for all $1 \leq i \leq m-1$.

Proof. We will show that for all $1 \leq l \leq m$

$$\sum_{i=1}^m \frac{a_i}{b_i} \geq \sum_{i=1}^l \frac{a'_i}{b_i} + \sum_{i=l+1}^m \frac{a_i}{b_i} + \sum_{i=1}^l \frac{a_i - a'_i}{b_l}.$$

Then if we set l as m , we have the proposition as $\sum_{i=1}^m a_i \geq \sum_{i=1}^m a'_i$.

We prove this by induction. For $l = 1$, it is clearly true since

$$\sum_{i=1}^m \frac{a_i}{b_i} = \frac{a_1}{b_1} + \sum_{i=2}^m \frac{a_i}{b_i} + \frac{a_1 - a'_1}{b_1}.$$

Suppose that it is true when $l = k$. Then

$$\begin{aligned} \sum_{i=1}^m \frac{a_i}{b_i} &\geq \sum_{i=1}^k \frac{a'_i}{b_i} + \sum_{i=k+1}^m \frac{a_i}{b_i} + \sum_{i=1}^k \frac{a_i - a'_i}{b_k} \\ &\geq \sum_{i=1}^k \frac{a'_i}{b_i} + \sum_{i=k+1}^m \frac{a_i}{b_i} + \sum_{i=1}^k \frac{a_i - a'_i}{b_{k+1}} \\ &= \sum_{i=1}^{k+1} \frac{a'_i}{b_i} + \sum_{i=k+2}^m \frac{a_i}{b_i} + \sum_{i=1}^{k+1} \frac{a_i - a'_i}{b_{k+1}}. \end{aligned}$$

Proof of Lemma 4.2. Let us denote the completion time of p_i in FNF and any other given schedule as c_i and c'_i , respectively. We need to show that $\sum_{i=0}^{\lfloor n/2 \rfloor} (T - c_i)/t_i \geq \sum_{i=0}^{\lfloor n/2 \rfloor} (T - c'_i)/t_i$. In fact, since $c_0 = c'_0 = 0$, it is enough to show that $\sum_{i=1}^{\lfloor n/2 \rfloor} (T - c_i)/t_i \geq \sum_{i=1}^{\lfloor n/2 \rfloor} (T - c'_i)/t_i$. Since we have Proposition 4.1 and $t_i \leq t_{i+1}$, it is enough to show that $\sum_{i=1}^l (T - c_i) \geq \sum_{i=1}^l (T - c'_i)$ for all $1 \leq l \leq \lfloor n/2 \rfloor$. This is true because $\sum_{i=1}^l c_i \leq \sum_{i=1}^l c'_i$ for all $1 \leq l \leq \lfloor n/2 \rfloor$ by Corollary 3.1. ■

Let the makespan of an optimal schedule and FNF be T_{OPT} and T_{FNF} , respectively. Then we have the following lemma.

LEMMA 4.3. *The number of fractional blocks by FNF in $R(T_{FNF})$ is at most $3n/2$.*

Proof. If a processor receives the message from processor p_i , then it can be mapped to a white block of p_i . To finish the schedule, we should send the message to n processors and therefore, there are n white blocks in $R(T_{FNF})$. In addition, each processor $p_i (0 \leq i \leq \lfloor n/2 \rfloor)$ may have a fraction of block which is not finished by time T_{FNF} . But at least one processor should have no incomplete block since the makespan of FNF is T_{FNF} . Thus the number of fractional blocks is at most $n + \lfloor n/2 \rfloor \leq 3n/2$. ■

THEOREM 4.1. *There are at most n fractional blocks in $R(\frac{2}{3}T_{FNF})$ in any schedule.*

Proof. It is enough to show that FNF can have at most n fractional blocks in $R(\frac{2}{3}T_{FNF})$ since FNF maximizes the number of blocks by Lemma 4.2. By time $\frac{2}{3}T_{FNF}$, each processor can have only $\frac{2}{3}$ of the number of blocks it has in $R(T_{FNF})$. Let f_i be the number of fractional blocks of processor p_i in $R(T_{FNF})$ and f'_i be the number of fractional blocks in $R(\frac{2}{3}T_{FNF})$. Since $f_i = \frac{T_{FNF}-c_i}{t_i}$ and $f'_i = \frac{\frac{2}{3}T_{FNF}-c_i}{t_i}$, we have $f'_i \leq \frac{2}{3}(\frac{T_{FNF}-\frac{3}{2}c_i}{t_i}) \leq \frac{2}{3}f_i$. Therefore, we have at most $\frac{2}{3} \cdot 3n/2 = n$ fractional blocks in $R(\frac{2}{3}T_{FNF})$. ■

COROLLARY 4.1. *Algorithm FNF gives a 1.5-approximation.*

Proof. Since we need to send the message to n processors in the optimal schedule, we should have n blocks in $R(T_{OPT})$. It implies that $T_{OPT} \geq \frac{2}{3}T_{FNF}$. ■

When transmission times are in a small range, the FNF heuristic has a better bound.

THEOREM 4.2. *Suppose $C' = n/(2 \sum_{i=1}^{\lfloor n/2 \rfloor} 1/t_i)$ then the FNF heuristic finds a solution of cost at most $T_{OPT} + C'$.*

Proof. In the bar chart of an optimal solution, the number of fractional blocks we can have between height T_{FNF} and T_{OPT} is $t/t_0 + t/t_1 + \dots + t/t_{\lfloor n/2 \rfloor}$ where $t = T_{FNF} - T_{OPT}$. Since we have at least n fractional blocks in $R(T_{OPT})$ and at most $3n/2$ fractional blocks in $R(T_{FNF})$, $\sum_{i=1}^{\lfloor n/2 \rfloor} t/t_i \leq n/2$. Therefore $T_{FNF} - T_{OPT} \leq n/(2 \sum_{i=1}^{\lfloor n/2 \rfloor} 1/t_i)$. ■

THEOREM 4.3. *If the transmission times of the fastest $\lfloor n/2 \rfloor$ processors (except p_0) is at most C , then the FNF heuristic finds a solution of cost at most $T_{OPT} + C$.*

Proof. In the bar chart of an optimal schedule, let f_i be the number of fractional blocks of processor p_i by time T_{OPT} and g_i be the number of complete blocks of processor p_i by T_{OPT} . In the bar chart of FNF schedule, let f'_i be the number of fractional blocks of processor p_i by time T_{OPT} .

Define g'_i to be the number of complete blocks by time $T_{OPT} + C$ and we need to prove that $\sum_{i=0}^{\lfloor n/2 \rfloor} g'_i \geq n$. For $1 \leq i \leq \lfloor n/2 \rfloor$, since $t_i \leq C$, we have $g'_i \geq \lceil f'_i \rceil$. Therefore, the number of complete blocks in $R(T_{OPT} + C)$ is

$$\begin{aligned} \sum_{i=0}^{\lfloor n/2 \rfloor} g'_i &\geq g'_0 + \sum_{i=1}^{\lfloor n/2 \rfloor} \lceil f'_i \rceil \geq g_0 + \sum_{i=1}^{\lfloor n/2 \rfloor} f'_i \\ &\geq g_0 + \sum_{i=1}^{\lfloor n/2 \rfloor} f_i \geq g_0 + \sum_{i=1}^{\lfloor n/2 \rfloor} g_i. \end{aligned}$$

The third inequality comes from the fact that FNF maximizes the number of fractional blocks in $R(T)$ for any T and $f_0 = f'_0$.

We also have $\sum_{i=0}^{\lfloor n/2 \rfloor} g_i \geq n$ by definition of T_{OPT} . Therefore, $T_{FNF} \leq T_{OPT} + C$. ■

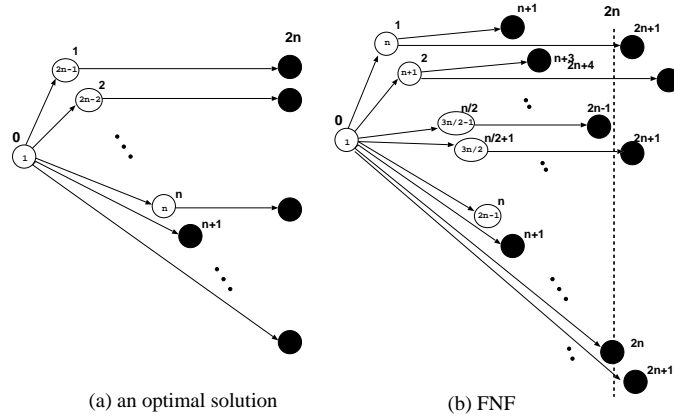


Figure 4: A bad example of FNF. The number inside a node is the transmission time of the processor and the number next to a node is the time it received the message. Black nodes are very slow processors. At time $2n$ (the dotted line), the optimal solution finishes broadcasting but in the FNF schedule, $\frac{n}{2}$ processors have not received the messages as yet.

5 Bad Example

Bhat et al. [5] gave an example and proved that the broadcast time by FNF can be $\frac{17}{16}$ times the optimal in the example. In this section, we show that in fact, FNF gives the broadcast time of $\frac{25}{22}$ times the optimal on the same example.

Consider the example shown in Figure 4. We have the source with transmission time 1 and $2n$ processors with a very large transmission time. Also there are n processors with transmission times $n, n+1, \dots, 2n-1$. In the optimal schedule, the source should send messages to processors with transmission time $2n-1, 2n-2, \dots, n$, respectively. In other words, at time i the node with transmission time $2n-i$ receives the message from the source. Immediately after receiving the message, each of these processors send a message to one of the slow processors. The schedule completes at time $2n$.

In the FNF schedule, the source sends messages to processors with transmission time $n, n+1, \dots, 2n-1$, respectively and again immediately after receiving the message, each of these processors send a message to one of the slow processors. At time $2n$, $\frac{n}{2}$ of the slow processors have not yet received the message. After time $2n$, processor with transmission time $n+i-1$ will send another message to a slow processor at time $2n+3i-2$. At the same time, processor with transmission time $\frac{3n}{2}+i-1$ send a message at time $2n+2i-1$. The source sends a message every time unit. Therefore, if t is the time needed to send messages to the remaining $\frac{n}{2}$ processors, then we have

$$\lfloor \frac{t+2}{3} \rfloor + \lfloor \frac{t+1}{2} \rfloor + t \geq \frac{n}{2}.$$

Therefore, we need at least $\frac{3n}{11} - \frac{7}{11}$ additional time. That means that the broadcast time of FNF can be $\frac{25}{22}$ times the optimal for large values of n .

6 Polynomial Time Approximation Scheme

We now describe a polynomial time approximation scheme for the problem of performing broadcast in the minimum possible time. Unfortunately, the algorithm has a very high running time when compared to the *fastest node first* heuristic.

We will assume that we know the broadcast time T of the optimal solution. Since $t_0 = 1$, we know that the minimum broadcast time T is between 1 and n , and we can try all possible values of the form $(1 + \epsilon)^j$ for some fixed $\epsilon > 0$ and $j = 1 \dots \lceil \frac{\log n}{\log(1+\epsilon)} \rceil$. In this guessing process we lose a factor of $(1 + \epsilon)$.

Let $\epsilon' > 0$ be a fixed constant. We define a set of fast processors F as all processors whose transmission time is at most $\epsilon'T$. Formally, $F = \{p_j | t_j \leq \epsilon'T, 1 \leq j \leq n\}$. Let S be the set of remaining (slow) processors. We partition S into collections of processors of similar transmissions speeds. For $i = 1 \dots k$, define $S_i = \{p_j | \epsilon'T(1 + \epsilon')^{i-1} < t_j \leq \epsilon'T(1 + \epsilon')^i, 1 \leq j \leq n\}$ where k is $\lceil \frac{\log(1/\epsilon')}{\log(1+\epsilon')} \rceil$. Since $t_1 \leq t_2 \leq \dots \leq t_n$, $F = \{p_1, \dots, p_{|F|}\}$ and $S = \{p_{|F|+1}, \dots, p_n\}$.

We first send messages to processors in F using FNF. We prove that there is a schedule with broadcast time at most $(1 + O(\epsilon))T$ such that all processors in F receive the message first. We then find a schedule for slow processors based on a dynamic programming approach.

Schedule for F : We use the FNF heuristic to send the message to processors in set F . Assume that the schedule for F has a broadcast time of T_{FNF} . In this schedule every processor $p_j \in F$ becomes idle at some time between $T_{FNF} - t_j$ and T_{FNF} .

We will prove that there is a schedule with broadcast time at most $(1 + O(\epsilon))T$ such that all processors in F receive the message first, and then send it to the slow processors. In an optimal schedule, let P_1 be the first $|F|$ processors (except the source) which finish receiving the message and T_1 be the time when all processors in P_1 finish receiving the message. The following lemma relates T_{FNF} with T_1 .

LEMMA 6.1. $T_1 \geq T_{FNF} - \epsilon'T$.

Proof. We prove this by contradiction. Consider a FNF schedule with set P_1 . If an optimal schedule is able to have all processors in P_1 finish receiving the message before $T_{FNF} - \epsilon'T$, then FNF can finish sending the messages to P_1 before time T_{FNF} (by Corollary 4.3). Since set F includes the fastest $|F|$ processors, this means that FNF can finish broadcasting for F before time T_{FNF} ; it is a contradiction since T_{FNF} is the earliest time that processors in F receive the message in FNF. ■

At time T_1 there can be some set of processors P_2 which have received the message *partially*. Note that $|P_2| \leq |P_1|$ since every processor in P_2 should receive the message from a processor in P_1 or from the source and at least one processor P_1 completes receiving the message at time T_1 .

LEMMA 6.2. *There is a schedule in which all processors in F receive the message no later than any processor in S and the makespan of the schedule is at most $(1 + 3\epsilon')T$.*

Proof. The main idea behind the proof is to show that an optimal schedule can be modified to have a certain form. Notice that by time T_1 the optimal schedule can finish broadcasting for P_1 and partially send the message to P_2 , and in additional time $T - T_1$ the optimal schedule can finish broadcasting for the remaining processors.

In FNF schedule, all processors in F have the message at time T_{FNF} . Since $|P_1 \cup P_2| \leq 2|F|$ and any processor in F has speed at most $\epsilon'T$, in additional $2\epsilon'T$ time we can finish broadcasting the message to $P_1 \cup P_2$. Once we broadcast the message to $P_1 \cup P_2$, we send the message to the remaining processors as in the optimal solution.

The broadcast time of this schedule is at most $T_{FNF} + 2\epsilon'T + T - T_1 \leq T_{FNF} + 2\epsilon'T + T - (T_{FNF} - \epsilon'T) = (1 + 3\epsilon')T$. ■

Create all possible trees of S : For the processors in S , we will produce a set \mathcal{S} of labeled trees \mathcal{T} .

A tree \mathcal{T} is any possible tree with broadcast time at most T consisting of a subset of processors in S . Then we label a node in the tree as i if the corresponding processor belongs to S_i ($i = 1 \dots k$). We prove that the number of different trees is constant.

LEMMA 6.3. *The size of \mathcal{S} is constant for fixed $\epsilon' > 0$.*

Proof. First consider the size of a tree \mathcal{T} (that is, the number of processors in the tree). Let us denote it as $|\mathcal{T}|$. Since the transmission time of processors in S is greater than $\epsilon'T$, we need at least $\epsilon'T$ time units to double the number of processors that received the message. It means that given a processor as a root of the tree, within time T we can have at most $2^{1/\epsilon'}$ processors receive the message. Therefore, $|\mathcal{T}| \leq 2^{1/\epsilon'}$. Now each node in the tree can have different label $i = 1 \dots k$. To obtain an upperbound of the number of different trees, given a tree \mathcal{T} we transform it to a complete binomial tree of size $2^{1/\epsilon'}$ by adding nodes labeled as 0. Then the number of different trees is at most $(k + 1)^{2^{1/\epsilon'}}$. ■

Attach \mathcal{T} to F : Let the completion time of every processor $p_j \in F$ be c_j . Each processor p_j in F sends a message to a processor in S every t_j time unit. Therefore, a fast processor p_j can send messages to at most $X_j = \lfloor \frac{T-c_j}{t_j} \rfloor$ other processors. Let $X = \sum_{p_j \in F} X_j$. Let us consider the time x_i of each sending point in X . We sort those x_i in nondecreasing order and attach a tree from \mathcal{S} to each point (See Figure 5). Note that we can attach at most $|X|$ trees of slow processors. Clearly $|X| \leq n$.

We check if an attachment is feasible, using dynamic programming. Recall that we partition slow processors into a collection of processors S_1, S_2, \dots, S_k ($k = \frac{\log(1/\epsilon')}{\log(1+\epsilon)}$). Let s_i denote the number of processors in set S_i . We define a state $s[j, n_1, n_2, \dots, n_k]$ ($0 \leq j \leq |X|$, $0 \leq n_i \leq s_i$) to be true if there is a set of j trees in \mathcal{S} that we can attach to first j sending points and the corresponding schedule satisfies the following two conditions: i) the schedule completes by time T and ii) exactly n_i processors in S_i appear in j trees in total. Our goal is to find out whether $s[j, s_1, s_2, \dots, s_k]$ is true for some j , which means that there is a feasible schedule with makespan at most T . The number of states is at most $O(n^{k+1})$ since we need at most n trees ($|X| \leq n$) and $s_i \leq n$.

Now we prove that each state can be computed in constant time. Given $s[j - 1, \dots]$, we compute $s[j, n_1, n_2, \dots, n_k]$ as follows. We try to attach all possible trees in \mathcal{S} to x_j . Then $s[j, n_1, n_2, \dots, n_k]$ is true if there exists a tree \mathcal{T}' such that the makespan of \mathcal{T}' is at most $T - x_j$ and $s[j - 1, n_1 - m_1, n_2 - m_2, \dots, n_k - m_k]$ is true where \mathcal{T}' has m_i slow processors belonging to set S_i . It can be checked in constant time since the size of \mathcal{S} is constant (Lemma 6.3).

THEOREM 6.1. *Given a value T , if a broadcast tree with broadcast time T exists, then the above algorithm will find a broadcast tree with broadcast time at most $(1 + \epsilon')(1 + 3\epsilon')T$.*

Proof. Consider the best schedule among all schedules in which processors in F receive the message first. By Lemma 6.2, the broadcast time of this schedule is at most $(1 + 3\epsilon')T$. We round up the transmission time of p_j in S_i to $\epsilon'T(1 + \epsilon')^i$ where i is the smallest integer such that $t_j \leq \epsilon'(1 + \epsilon')^i T$. By this rounding, we increase the broadcast time by factor of at most $1 + \epsilon'$. Therefore, the broadcast time of our schedule is at most $(1 + \epsilon')(1 + 3\epsilon')T$. ■

THEOREM 6.2. *The algorithm takes as input the transmission times of the n processors, and constants $\epsilon, \epsilon' > 0$. The algorithm finds a broadcast tree with broadcast time at most $(1 + \epsilon)(1 + \epsilon')(1 + 3\epsilon')T$ in polynomial time.*

Proof. We try the above algorithm for all possible value of the form $T = (1 + \epsilon)^j$ for $j = 1 \dots \lceil \frac{\log n}{\log(1+\epsilon)} \rceil$. This will increase the broadcast time by factor of at most $1 + \epsilon$. Therefore the broadcast time of our schedule is at most $(1 + \epsilon)(1 + \epsilon')(1 + 3\epsilon')T$.

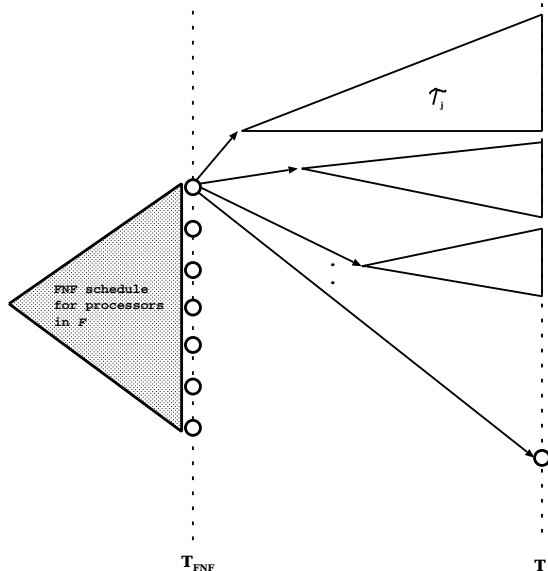


Figure 5: Attach trees for slow processors to fast processors

For each given value $(1 + \epsilon)^j$, we find FNF schedule for processors in F (it takes at most $O(n \log n)$) and attach trees of slow processors to processors in F , using dynamic programming. As we discussed earlier, the number of states is $O(n^{k+1})$ and each state can be checked if it is feasible in $O((k + 1)^{2^{1/\epsilon'}})$ time, which is constant. Thus the running time of our algorithm is $O(\lceil \frac{\log n}{\log(1+\epsilon)} \rceil (n \log n + (k + 1)^{2^{1/\epsilon'} + 1} \cdot n^{k+1}))$ where k is $\lceil \frac{\log(1/\epsilon')}{\log(1+\epsilon)} \rceil$. ■

7 Multicast

A multicast operation involves only a subset of processors. By utilizing fast processors which are not in the multicast group, we can reduce the multicasting time significantly. For example, suppose that we have m processors with transmission time t_1 and m more processors with transmission time t_2 where $t_1 < t_2$. Let we want to multicast a message to all processors with transmission time t_2 . If we only use processors in the multicast group, it will take $t_2 \cdot \log m$ time. But if we utilize processors with transmission time t_1 , we can finish the multicast in $t_1 \cdot (\log m + 1)$. Therefore, when $t_1 \ll t_2$, the speed-up is significant.

THEOREM 7.1. *Suppose that we have a ρ -approximation algorithm for broadcasting. Then we can find a ρ -approximation algorithm for multicasting.*

Proof. Note that if an optimal solution utilizes k processors not in the multicast group, then those processors are the k fastest ones. Therefore, if we know how many processors participate in multicasting, we can use our ρ -approximation algorithm for broadcasting. By trying all possible k and taking the best one, we have ρ -approximation for multicasting. ■

THEOREM 7.2. *We have a polynomial time approximation scheme for multicasting.*

Proof. The proof is similar to Theorem 7.1. We find approximation schemes with k fastest processors not in the multicast group for all possible k , and take the best one. ■

8 NP-Completeness

We first show that the following problem is *NP*-complete. Then we show how to use this problem to derive a proof for the *NP*-hardness of minimizing the broadcast time.

The regular 3-PARTITION problem can be stated as follows: Given a finite set A of $3m$ elements, an integer B and a size $s(a)$ for each element in A , such that $\frac{B}{4} < s(a) < \frac{B}{2}$. Can A be partitioned into m disjoint sets S_1, \dots, S_m such that for $1 \leq i \leq m$, $\sum_{a \in S_i} s(a) = B$?

Notice that the constraints imply that each set S_i contains exactly three elements.

In Appendix A we prove that the following special version of 3-PARTITION is *NP*-complete. We call it the SEEDED 3-PARTITION problem. Given a finite set A of $3m$ elements, where m is a power of 2, an integer B and a size $s(a)$ for each element in A , such that $\frac{B}{4} < s(a) < \frac{B}{2}$. Also specified is a (seed) subset A' of A , such that each element $a' \in A'$ has $s(a') < B(\frac{1}{4} + \epsilon)$ for a small positive constant ϵ , and each element $a \in A \setminus A'$ has $B(\frac{3}{8} - \epsilon) < s(a) < B(\frac{3}{8} + \epsilon)$. Can A be partitioned into m disjoint sets S_1, \dots, S_m such that for $1 \leq i \leq m$, $\sum_{a \in S_i} s(a) = B$ and each S_i contains exactly one element from A' ?

We show how to reduce SEEDED 3-PARTITION to minimizing the broadcast time in heterogenous networks. Formally, given n processors, each with message transmission time t_i , is there a broadcast schedule to broadcast a message from one source processor to all other processors in T units of time?

The reduction works as follows: We create a set S of m processors each with transmission time 1. In addition we have a set S' containing $2m$ processors, m of them have transmission time $1.5 - 8\epsilon$, and m of them have transmission time $1.5 + 8\epsilon$. We also have a set S'' of m processors each with transmission time 2. Corresponding to each element in $a \in A$ we create a processor $p_a \in A$ with transmission time $t_a = \frac{s(a)}{B/4}$. For each element $a \in A'$ we create an additional set of processors $\bar{p}_a \in \bar{A}$ with transmission time $\bar{t}_a = \frac{B-2s(a)}{B/4}$. We also create a set M of $9m$ processors with extremely high transmission times. The processors in M cannot be used to send any other processors a message if the broadcast has to complete in time T . The source of the message is a processor in S . We set $T = \log m + 5$. (Since m is a power of 2, T is an integer.)

First observe that all the processors have transmission times between 1 and 2. For the processors in set A since $\frac{B}{2} > s(a) > \frac{B}{4}$, this follows. In fact, for the processors in A' the transmission times are strictly between 1 and $1 + 4\epsilon$. For the processors in $A \setminus A'$ the transmission times are strictly between $1.5 - 4\epsilon$ and $1.5 + 4\epsilon$. For the processors in \bar{A} , notice that the transmission times are almost 2.

We first show that if there is a solution to the SEEDED 3-PARTITION problem, then there is a solution to the broadcasting problem with $T = (\log m + 1) + 4$. The source message first broadcasts the message to processors in S . This takes $\log m$ transmission steps (each of which takes 1 unit of time). Suppose the solution of SEEDED 3-PARTITION is represented by sets S_i (see Figure 6). Let x_i^1, x_i^2, x_i^3 be the 3 elements in S_i , with $s(x_i^1) \leq s(x_i^2) \leq s(x_i^3)$. We can assume that $x_i^1 \in A'$ since this was a property of SEEDED 3-PARTITION. In the next transmission step, all the processors of S can simultaneously send messages to the processors corresponding to x_i^1 for $i = 1 \dots m$. This takes another 1 unit of time. We have used up $\log m + 1$ time units. At this stage there are 4 time units remaining. We have processor $p_{x_i^1}$ send its first message to $p_{x_i^2}$, its second message sent to processor $\bar{p}_{x_i^1}$, which in turn sends a message to a processor in M . The third message sent by $p_{x_i^1}$ goes to a processor in M . Meanwhile $p_{x_i^2}$ sends its first message to $p_{x_i^3}$ and the second message to a processor in M . Also $p_{x_i^3}$ sends the message to a processor in M .

Since $s(x_i^1) + s(x_i^2) + s(x_i^3) = B$, this transmission chain takes exactly 4 units of time. Also note that $s(x_i^1) + 2s(x_i^2) \leq B$ and thus this chain has length at most 4 as well. The path from $p_{x_i^1}$ to $\bar{p}_{x_i^1}$ has length $2\frac{s(x_i^1)}{B/4} + \frac{B-2s(x_i^1)}{B/4} = 4$ as well.

Since there are 4 time units remaining at time T' , each processor in S can send out 4 more messages.

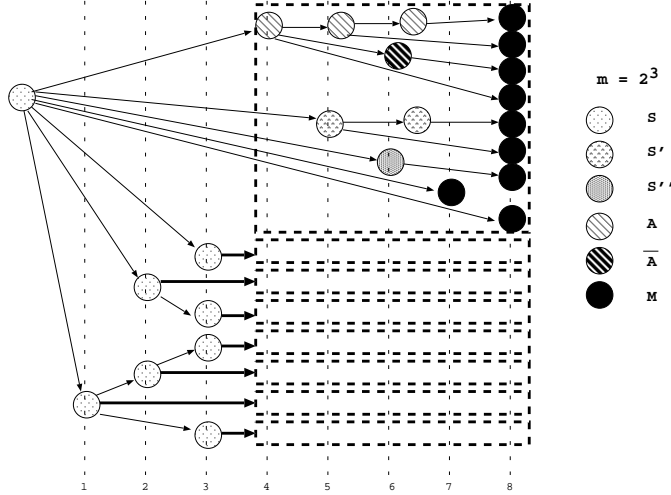


Figure 6: Reduction from SEEDED 3-PARTITION to minimizing the broadcast time in heterogeneous networks. Set S, S', S'' include processors with transmission time 1, $\{1.5 - 8\epsilon, 1.5 + 8\epsilon\}, 2$, respectively. Set A has processors with transmission time $t_{x_i^1}, t_{x_i^2}, t_{x_i^3}$ and processors in \bar{A} have transmission time $\bar{t}_{x_i^1}$. M includes very slow processors.

The first message goes to the faster processor in S' . Thus half the processors in S' have completion time $(1 + \log m) + 1$. These in turn then send the message to the other half (slower) processors in S' . Thus half the processors in S' have completion time $(1 + \log m) + 1 + 1.5 - 8\epsilon$. At this time the nodes in S' can send messages to $2m$ nodes from M in the remaining time.

The second message sent by processors in S goes to processors in S'' . These have completion time $(1 + \log m) + 2$. In the remaining time, the processors in S'' can each send a message to a node in M .

The third and fourth messages from processors in S go to processors in M . Thus $2m$ processors in M get their messages directly from processors in S .

We now prove that if there is a solution to the broadcast problem in T time units, then there is a solution to the SEEDED 3-PARTITION problem. Let $T' = \log m + 1$. Consider an optimal broadcasting schedule that takes $T' + 4$ time units. We first recall the following theorem shown by Liu [16].

THEOREM 8.1. [16] *There exists an optimal broadcast sequence (ordered by completion times) in which all the fastest processors appear before all of the other processors.*

We break our schedule into two groups - Group I consists of the processors that receive the message strictly before time T' . Group II consists of the processors that receive the message at or after T' (and clearly before T). Using Theorem 8.1 it is easy to see that Group I consists of the processors in S in an optimal solution, since all these have transmission time 1. Group II consists of the processors in S', S'', A, \bar{A} and M . For the Group II processors, we wish to argue that they have a particular form.

We analyze the subtree $T(i)$ rooted at a single processor $i \in S$ at time $T' - 1$. This subtree includes the set of processors that receive a message, either directly or indirectly through this processor in S . Since the transmission time is 1, at each time $t = T', \dots, T' + 4$ some processor receives a message. We refer to the processor that receives a message at time $T' + j, j = 0 \dots 4$ as G_j (see Figure 7).

We now examine the maximum sizes each subtree rooted at G_j can have. Recall that the total number of processors in $S' \cup S'' \cup A \cup \bar{A} \cup M$ is $16m$. Thus the average number of processors in the subtree $T(i)$ is 16. We will argue that the maximum number of processors is exactly 16.

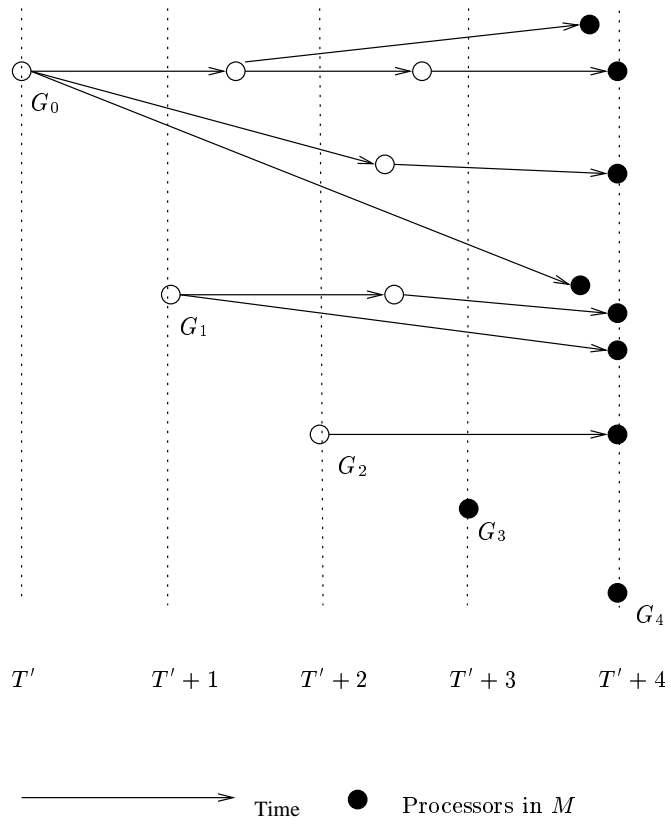


Figure 7: Figure to show last 4 steps of the broadcast schedule.

The subtree rooted at G_0 has size at most 8. This is because since all the processors in this subtree have transmission time > 1 , we can only have 3 rounds of transmissions, and not 4.

Similarly, the subtree rooted at G_1 has size at most 4, since there are only 2 rounds of transmission, and not 3. The subtree rooted at G_2 has size at most 2, and the subtrees rooted at G_3 and G_4 each have size at most 1. This implies an upper bound of 16 on the number of processors in $T(i)$. Thus each subtree $T(i)$ must have exactly 16 processors.

This subtree has 9 leaves, and these must all be chosen from M as there are a total of $9m$ nodes in M , and 9 are chosen by each $T(i)$ subtree. (Recall that these processors cannot do any transmissions as they are very slow.)

Observe that a processor with speed 2 (from S''), must be chosen as G_2 , otherwise we cannot complete the transmission to 16 processors in $T(i)$. In fact each $T(i)$ group can have only one processor from S'' for this reason. Processors G_3 and G_4 are both chosen from M , as are all the leaves in the subtree $T(i)$.

Note that only processors from A' can be G_0 since the transmission time of G_0 cannot be greater than $4/3$. Let us denote it as $p_{x_i^1}$. Then w.l.o.g, we can assume that $\bar{p}_{x_i^1}$ is the second child of $p_{x_i^1}$ (if not, we can exchange the second child with $\bar{p}_{x_i^1}$ and keep the broadcast time within $\log m + 5$).

In the subtree rooted at G_0 , the first child and first grandchild of G_0 are not fixed and the broadcast time of this chain (starting from G_0 to M) should be at most 4. In the subtree rooted at G_1 , G_1 and its first child is available and the broadcast time of this chain should be at most 3. In fact since we have to place all processors from A and S' in these two chains and the sum of transmission times of those processors is $7m$. We can show the chain starting from G_0 should be exactly 4 and the chain starting from G_1 exactly 3 (otherwise, we cannot include all processors).

Now we prove that processor from S' cannot be in the subtree rooted at G_0 . This is because any combination of processors from A and S' cannot sum to exactly 4. Once we establish this, it is easy to see that all three processors corresponding to the 3 elements of the 3-set are all in one chain in the G_0 subtree. This implies that $\frac{s(x_i^1)+s(x_i^2)+s(x_i^3)}{B/4} = 4$. Thus the sum of the sizes is exactly B .

9 Acknowledgement:

We thank Maxim Sviridenko for useful discussions. We are also grateful to Gerhard Woeginger for collaborating with us on developing a polynomial time approximation scheme.

References

- [1] *Message Passing Interface Form*, <http://www.mpi-forum.org>, March 1994.
- [2] M. Banikazemi, V. Moorthy and D. K. Panda. Efficient Collective Communication on Heterogeneous Networks of Workstations. *International Conference on Parallel Processing*, pp. 460–467, 1998.
- [3] A. Bar-Noy, S. Guha, J. Naor and B. Schieber. Multicasting in Heterogeneous Networks. *Proceedings of the 13th Annual ACM Symposium on Theory of Computing*, pp. 448–453, 1998.
- [4] A. Bar-Noy and S. Kipnis. Designing broadcast algorithms in the Postal Model for Message-passing Systems. *Mathematical Systems Theory*, 27(5), pp. 431–452, 1994.
- [5] P. Bhat, C. Raghavendra and V. Prasanna. Efficient Collective Communication in Distributed Heterogeneous Systems. *Proceedings of the International Conference on Distributed Computing Systems*, pp. 15–24, 1999.
- [6] J. Bruck, D. Dolev, C. Ho, M. Rosu and R. Strong, Efficient Message Passing Interface(MPI) for Parallel Computing on Clusters of Workstations. *J. Parallel Distributed Computing*, 40 (1997), pp. 19–34.
- [7] M.R.Garey and D.S.Johnson *Computer and Intractability: A Guide to the Theory of NP-completeness*, Freeman, New York, 1979.
- [8] W. Gropp, E. Lusk, N. Doss and A. Skjellum. A High-performance, portable Implementation of the MPI: a Message Passing Interface Standard. *Parallel Computing* 22(1996), pp. 789–828.

- [9] N. Hall, W.-P. Liu, and J. Sidney. Scheduling in broadcast networks, *Networks* 32, pp. 233-253, 1998.
- [10] S.M.Hedetniemi, S.T. Hedetniemi and A.L.Liestman. A Survey of Gossiping and Broadcasting in Communication Networks, *Networks* 18(1991), pp.129-134.
- [11] R. Karp, A. Sahay, E. Santos and K.E.Schauser. Optimal Broadcast and Summation in the Logp Model. *Proceedings of 5th Annual Symposium on Parallel Algorithms and Architectures*, 1993 pp. 142-153.
- [12] R. Kesavan, K. Bondalapati and D. Panda. Multicast on Irregular Switch-based Networks with Wormhole Routing. *Proceedings of the International Symposium on High Performance Computer Architectures*, 1997, pp. 48-57.
- [13] S. Khuller and Y-A. Kim. On Broadcasting in Heterogenous Networks. *Proceedings of the 15th Annual ACM/SIAM Symposium on Discrete Algorithms*, 2004, pp. 1004–1013.
- [14] S. Khuller, Y-A. Kim and G. Woeginger. Approximation Schemes for Broadcasting in Heterogenous Networks. *Proceedings of the Workshop on Approximation Algorithms*, LNCS 3122, 2004, pp. 163–170.
- [15] R. Libeskind-Hadas, J. R. K. Hartline, P. Boothe, G. Rae, and J. Swisher. On Multicast Algorithms for Heterogeneous Networks of Workstations. *Journal of Parallel and Distributed Computing* 61(11), pp. 1665 - 1679, 2001.
- [16] P. Liu. Broadcasting Scheduling Optimization for Heterogeneous Cluster Systems. *Journal of Algorithms* 42, pp.135-152, 2002.
- [17] P. Liu and D. Wang. Reduction Optimization in Heterogeneous Cluster Environments. *Proceedings of the International Parallel and Distributed Processing Symposium*. 2000, pp. 477-482.
- [18] P. Liu and T-H. Sheng. Broadcasting Scheduling Optimization for Heterogeneous Cluster Systems. *ACM Symp. on Parallel Algorithms and Architectures (SPAA)*, 2000, pp 129–136.
- [19] D. A. Patterson, D. E. Culler and T. E. Anderson. A case for NOWs (Networks of Workstations). *IEEE Micro*, 15(1) (1995), pp.54–64.
- [20] J. Pruyne and M. Livny. Interfacing Condor and PVM to Harness the Cycles of Workstation Clusters. *Journal on Future Generations of Computer Systems*, 12:53–65, 1996.
- [21] D. Richards and A. L. Liestman. Generalization of broadcasting and Gossiping. *Networks* 18(1988), pp. 125-138.

A NP-Completeness Proofs

The goal of this section is to prove that SEEDED 3-PARTITION is *NP*-complete.

First some definitions.

3-DIMENSIONAL MATCHING: Given a set $M \subseteq W \times X \times Y$, where W, X, Y are all disjoint sets each having q elements. Does M contain a matching, i.e., a subset $M' \subset M$ such that $|M'| = q$ and no two elements of M' agree in any co-ordinate?

We first show that the following problem is *NP*-complete by a reduction from 3-DIMENSIONAL MATCHING. This problem is a variation of the standard “4 PARTITION” problem.

SMALL-LARGE 4 PARTITION: Given finite sets A_1 and A_2 each having $2m$ elements, and a bound B and a size $s(a)$ for each element, such that $\frac{B}{5} < s(a) < \frac{B}{3}$. Moreover, $s(a) > \frac{B}{4}$ for $a \in A_1$ and $s(a) < \frac{B}{4}$ for $a \in A_2$. Can we find m disjoint sets S_1, \dots, S_m such that $\sum_{a \in S_i} s(a) = B$ and each set S_i includes *exactly* two elements from each A_j ?

The proof for the *NP*-completeness of the SMALL-LARGE 4 PARTITION problem is almost the same as the proof of the standard 4 PARTITION problem [7]. We have to change some essential parameters in the reduction.

Proof. SMALL-LARGE 4-PARTITION clearly belongs in *NP* since we can verify in polynomial time that the given partition has the stated properties.

In fact, we will show that we can find a reduction where $\max_a s(a) \leq poly(|A_1|^4)$. Let $W = \{w_1, \dots, w_q\}$, $X = \{x_1, \dots, x_q\}$, $Y = \{y_1, \dots, y_q\}$, and $M \subseteq W \times X \times Y$ denote an arbitrary instance of 3DM. We assume w.l.o.g that $|M| \geq q$. Our instance of SMALL-LARGE 4 PARTITION has $|A_j| = 2|M|$ elements. Corresponding to each triple in M there are two elements in each A_j set. The elements corresponding to each $z \in W \cup X \cup Y$ will be denoted by $z[1], z[2], \dots, z[N(z)]$, where $N(z)$ denotes the

number of triples from M in which z occurs. We regard $z[1]$ as the actual element corresponding to z and $z[2]$ through $z[N(z)]$ as the dummy elements. Let $r = 128q$.

We generate set A_1 as follows:

$$\begin{aligned} s(w_i[1]) &= 33r^4 + ir + 1 \text{ for } 1 \leq i \leq q. \\ s(w_i[l]) &= 34r^4 + ir + 1 \text{ for } 1 \leq i \leq q, 2 \leq l \leq N(w_i). \\ s(x_j[1]) &= 33r^4 + jr^2 + 2 \text{ for } 1 \leq j \leq q. \\ s(x_j[l]) &= 34r^4 + jr^2 + 2 \text{ for } 1 \leq j \leq q, 2 \leq l \leq N(x_j). \end{aligned}$$

We generate A_2 as follows:

$$\begin{aligned} s(y_k[1]) &= 26r^4 + kr^3 + 4 \text{ for } 1 \leq k \leq q. \\ s(y_k[l]) &= 24r^4 + kr^3 + 4 \text{ for } 1 \leq k \leq q, 2 \leq l \leq N(y_k). \end{aligned}$$

For each triple $m_p = (w_i, x_j, y_k) \in M$ we create an element u_p such that $s(u_p) = 28r^4 - kr^3 - jr^2 - ir + 8$.

We set the value of $B = 120r^4 + 15$. Note that all the sizes are greater than $\frac{B}{5} = 24r^4 + 3$. All the sizes are smaller than $\frac{B}{3} = 40r^4 + 5$. The sizes of elements in A_1 are $> \frac{B}{4} = 30r^4 + 15/4$ and the sizes of the elements in A_2 are $< \frac{B}{4}$.

The size of each element is at most $35r^4 \leq 35(2^7q)^4$. The proof of the correctness of this reduction is almost the same as in Garey and Johnson [7]. \blacksquare

We then use the SMALL-LARGE 4 PARTITION problem to show that the SEEDED 3-PARTITION problem is NP -complete. This proof is a modification of the proof in Garey and Johnson [7].

SEEDED 3-PARTITION: Given a finite set A of $3m$ elements, where m is a power of 2, an integer B' and a size $s'(a)$ for each element in A , such that $\frac{B'}{4} < s'(a) < \frac{B'}{2}$. Also specified is a (seed) subset A' of A , such that each element $a' \in A'$ has $s'(a') < B'(\frac{1}{4} + \epsilon)$ for some small positive constant ϵ . For each element $a \in A \setminus A'$ we have $B'(\frac{3}{8} - \epsilon) < s'(a) < B'(\frac{3}{8} + \epsilon)$. Can A be partitioned into m disjoint sets S_1, \dots, S_m such that for $1 \leq i \leq m$, $\sum_{a \in S_i} s'(a) = B'$ and each S_i contains exactly one element from A' ?

Proof. The SEEDED 3-PARTITION problem is clearly in NP . We prove that it is NP -complete by reducing from the SMALL-LARGE 4 PARTITION problem. Here we do not show that m is a power of 2. This can be shown using standard padding techniques of adding dummy elements.

We assume that we have an instance of SMALL-LARGE 4 PARTITION specified by A_1, A_2 and the value B .

We show how to generate a set A and a set A' as follows. The parameter used is B' .

Assume that each set A_i has $2n$ elements. We create a set A with $3m$ elements where $m = 8n^2 - n$. We also specify a subset A' with m elements. We will create one element corresponding to each element in $\cup_i A_i$, two elements for each pair of elements in $\cup_i A_i$ and $8n^2 - 3n$ "filler" elements. The subset of elements in A' correspond to all the filler elements F , and the small elements corresponding to the ones from A_2 . Note that $|A'| = |A_2| + |F|$ has $2n + 8n^2 - 3n = m$ elements. $|A \setminus A'| = |A_1| + 2p$, where p is the number of pairs. Thus $|A \setminus A'| = 2n + 4n(4n - 1) = 2(n + 8n^2 - 2n) = 2m$. The goal is to find a solution to SEEDED 3-PARTITION where each 3-set chooses one element from A' and two elements from $A \setminus A'$.

Corresponding to each $a \in A_1$ we create an element w with size $s'(w) = 1499B + 4s(a) + 1$.

Corresponding to each $a \in A_2$ we create an element w with size $s'(w) = 1000B + 4s(a) + 1$.

Corresponding to each pair of elements $a_p, a_q \in \cup_i A_i$ we create two pairing elements $u[p, q]$ and $\bar{u}[p, q]$ with $s'(u[p, q]) = 1501B - 4s(a_p) - 4s(a_q) + 2$ and $s'(\bar{u}[p, q]) = 1497B + 4s(a_p) + 4s(a_q) + 2$.

The filler elements u_k^* all have size $s'(u_k^*) = 1002B$. The bound B' for SEEDED 3-PARTITION is chosen as $4000B + 4$.

First note that all the elements have sizes between $\frac{B'}{4} = 1000B + 1$ and $\frac{B'}{2} = 2000B + 2$. The elements of A' consist of the filler elements that have size $1002B < B'(\frac{1}{4} + \frac{1}{2000})$. The elements in A_2 are strictly smaller than $\frac{B'}{4}$ and thus the corresponding elements are clearly at most $1000B + B + 1$ and are at most the size of the filler elements. Thus the property that all elements in A' are at most $B'(\frac{1}{4} + \epsilon)$ is true.

Note that the sizes of the elements in $A \setminus A'$ are in the range $1498B$ and $1502B$. Since $B'(\frac{3}{8} - \frac{1}{2000}) < s'(a) < B'(\frac{3}{8} + \frac{1}{2000})$. Thus the bound follows with $\epsilon = \frac{1}{2000}$. (Note that by changing the parameters, we can make ϵ as small as we wish.)

Suppose there is a solution to SMALL-LARGE 4 PARTITION, then assume that a set in the partition has the form $\{a_i, a_j, a_k, a_l\}$ with the elements having the property that $a_i, a_j \in A_1, a_k, a_l \in A_2$. Our SEEDED 3-PARTITION would contain the sets $\{a_k, a_i, u[i, k]\}$ and $\{a_l, a_j, \bar{u}[i, k]\}$. The first element in each set is from A' . Notice that the sum of each set is $4000B + 4$. Doing this for each of the n 4-sets, we get $2n$ 3-sets using no filler elements. For the remaining pairing elements we can pair them together with one filler element each. Thus we would get a set $\{u_k^*, u[p, q], \bar{u}[p, q]\}$. Each has size $4000B + 4$; in addition each has one element from A' .

To prove the converse, assume that there is a solution to the SEEDED 3-PARTITION problem. We have a collection of m 3-sets. Each 3-set has one element from A' and two elements from $A \setminus A'$. Each filler element from A' pairs up with two pairing elements from $A \setminus A'$ (since they each are congruent to 2 mod 4). Each element corresponding to A_2 pairs up with two elements from $A \setminus A'$ as well. The element from A_2 is congruent to 1 mod 4, it must pair up with one pairing element and one non-pairing element from A_1 . By merging two 3-sets that use $u[i, j]$ and $\bar{u}[i, j]$ we can obtain a 4-set that has two elements from A_1 and two elements from A_2 . It is an easy exercise to check that in each 3-set the smallest element is from A' . The rest of the proof is the same as in Garey and Johnson [7]. ■