

FILE ALLOCATION ON LOCAL COMPUTER SYSTEMS WITH TWO-LEVEL MULTI-ACCESS NETWORKS

ABSTRACT

The file allocation problem (FAP) on local computer systems with two-level multi-access networks is studied in this paper. FAP is less complicated, although it is still NP-hard, on such networks than that on general networks because the communication overhead for remote file access is reduced from a random matrix to a three-value variable. If the communication overhead dominates other overheads, the simple file allocation problem becomes polynomially solvable. For optimal solution, the problem is transformed into a pure zero-one integer programming problem and solved using a variation of Balas's additive search algorithm. A simple heuristic algorithm based on the 0-1 knapsack problem is proposed with analytical and empirical evaluations. The performance degradation of the proposed heuristic algorithm is shown to be small in both evaluations.

Keywords: Distributed computing, distributed database, file allocation, local computer system, two-level local multi-access network.

FILE ALLOCATION ON LOCAL COMPUTER SYSTEMS WITH TWO-LEVEL MULTI-ACCESS NETWORKS

Yao-Nan Lien
AT&T Bell Laboratories
101 Crawfords Corner Rd.
Holmdel, NJ 07733
E-mail: yaonan.lien@att.com

Yih-Long Chang*
School of Management
Georgia Institute of Technology
Atlanta, GA 30332-0520
E-mail: YihLong.Chang@som.gatech.edu

Benjamin W. Wah
Department of Electrical and Computer Engineering
and the Coordinated Science Laboratory
University of Illinois at Urbana-Champaign
Urbana, Illinois 61801
E-mail: wah@aquinas.csl.uiuc.edu

1. INTRODUCTION

This paper studies the file allocation problems on local computer systems with two-level multi-access networks. As computer and local-area networks (LANs) technologies advance, users are moving toward "downsizing" their mainframe systems to distributed computing systems (DCS) which interconnects a number of smaller computers with a local area network. In such environment, both users and resources are physically dispersed in a local area, it requires special system control and resource management strategies to deliver the effective aggregated computing capability to the users. One particular example of the DCS is Merrill Lynch case reported in

Berman and Nigam [1992]. To avoid a premature performance saturation, it is important to efficiently manage the following scarce resources with a balanced utilization: processing capacity (including processor and memory), disk, and network. Among many other resource management strategies, a good file allocation strategy is essential to reduce remote data-access overhead and to increase data availability. Compared to wide-area-networks (WAN), the performance of DCS on LAN environment is much more sensitive to the quality of file allocation strategy due to its lower space locality. (This is because (1) users are able to access the system from many different nodes and (2) the locations of processes and data are often controlled by the system for reasons such as load balancing and location transparency.)

The placement of files has been extensively studied as the *file allocation problem* (FAP) in distributed databases (Aperis, 1988; Chu, 1969; Dowdy and Foster, 1982). The FAP entails the distribution of replicated files to a set of sites on a computer network to minimize the overall overhead or cost. To reduce remote data-access overhead and to increase data availability, multiple copies of files are allocated to sites where demands are strong, however, at the expense of extra overhead in consistency maintenance and data storage. Originated by Chu (1969), FAP was formulated as an integer programming problem in which the objective is to minimize the overall operating overhead with the constraints of response time and storage. Subsequently, much work has been done in this area (Lien and Wah, 1987; Lien et al., 1988; Wah, 1984; Wah and Lien, 1984). A special case of FAP is the *simple file allocation problem* (SFAP) (Ramamooorthy and Wah, 1983), in which multiple copies of a single file are to be allocated and the effects of queries, updates, and data storage are represented as costs. Another variation of FAP is to dynamically migrate files to requesting sites in order to preserve or improve system performance. The file migration problems have also been widely studied (Smith, 1981a,b,c; Procar, 1982; Liu Sheng, 1986, 1992a, 1992b). This paper focuses on the static file allocation problems on local computer systems with two-level multi-access networks.

Both FAP and SFAP on general point-to-point networks are known to be NP-hard. These problems in local area environments were studied by Wah and Lien (1985). They found that, along with some other interesting results, the broadcast characteristics of this kind of networks significantly reduces the complexities of FAP and SFAP, resulting more efficient solutions. The scalability of a single LAN system may be limited by the LAN bandwidth. It is a natural step to scale up the system capacity by expanding a LAN system into a multi-LAN system where a higher bandwidth backbone network is used to interconnect a number of regular LANs. (Because of the enhanced network capacity and user population, such a network environment may spread around a metropolitan area referred to as a MAN. This paper makes no difference between the two as far as the physical and logical characteristics are concerned.) This paper studies the file allocation problem on this kind of network environments. (Note that the traditional diskless workstation and single file server environments do not need to solve FAP problem because all files are located in the file server.)

The model of the problem is proposed in Section 2. In Section 3, the general FAP problem is transformed into a linear pure zero-one integer programming problem and is solved by a search

2. MODELS

We assume that a data site takes the following steps to process a remote read-only query:

1. Interpret the query (processing overhead),
2. search data from disk (disk I/O overhead), and
3. transmit the result back to the originating site (communication overhead).

Assume that an update is processed in the following simplified steps:

1. A data site interpret the query (processing overhead), and search data from disk (disk I/O overhead in one site);
2. All data sites proceed with a concurrency and transaction control (communication overhead); and
3. All data sites update the data without a search (nominal I/O overhead on each data site).

The transaction control in a real DBMS is much more complicated than the above procedure. However, since FAP only concerns with the statistical characteristics of all overhead components, there is no need to elaborate the implementation details.

2.1 Throughput vs. Response Time

Two major optimization objectives of a DCS system design are maximizing throughput and minimizing response time. This paper chooses to minimize the response time for the following reasons:

1. The two objectives are not completely independent to each other. If a good scheduling strategy is in place, optimizing one may not necessarily compromise the other significantly.
2. The throughput is easier to scale up in a DCS environment by adding new nodes into the system. On the other hand, adding new nodes may not necessarily reduce response time unless the system is always saturated.

2.2 Cost Analysis

2.2.1 Disk I/O Overhead vs. Communication Overhead

Compared to WANs, the communication overheads are less dominant in LAN environments. (Disk access time is in the range of 10 ms and the transfer rate is in the range of 5 Mbytes; and transmitting a short message from one process to a remote process on the same LAN is in the range of 5-20 ms.) It is a natural guess that disk overhead should be a major cost component in FAP. A careful investigation which will be shown in the following paragraph tells us that this is not quite true. Disk I/O overhead actually can be ignored from the FAP cost model

method based on Balas's additive algorithm [Balas, 1965]. The algorithm is suitable for evaluation of small to moderate sized problems. In Section 4, the problem of a communication-dominant system is presented. The practical considerations of the file allocation problems in real systems are discussed in Section 5.

simply because it is not sensitive to the file location. Except some extreme cases, the disk I/O overhead is actually independent of where the data is stored. In other words, FAP cannot help to reduce disk I/O overhead even though it is very important. (There are many ways to reduce I/O overhead such as indexing, using good query processing algorithms, adding new disks, using faster disks, etc.)

The following discussion will explain why disk I/O overhead is not sensitive to file allocation. First, the access time are fairly uniform across all types of disks (assuming they all use up-to-date technology.) Secondly, for read-only queries, approximately the same amount of disk access will be generated whether it is stored locally or remotely when the same query algorithm is used in each of the data sites. For update requests, this assumption may not be severely violated even when the system is update-intensive and data is replicated because (1) it is the concurrency control and transaction control that make an update request require longer response time, and (2) only one data site needs to process the update request to identify the data items (tuples) that are to be updated.

2.2.2 Processing Overhead

Processing overhead referred to the CPU (or processors) overhead. It was ignored in all other FAP study, either in WAN or in LAN environments, based on the following assumptions:

1. In a WAN environment, the processing overhead is so insignificant as compared to other overheads.
2. In a LAN environment, even if it is comparable with other overheads, it will become very insignificant in the near future as the processor speed almost doubled every two or three years.

Unfortunately, after a careful study, we found that the second assumption may not be valid. For the following reasons, the processing overhead may not always be ignored in LAN environments.

The recent memory technology advances are all on the storage density enhancement, not on the access speed enhancement. Typical memory access time is about 50-80 ns at chip level and 100-500 ns at bus level. In the context of FAP, the processing overhead will impact the communication cost which will be explained in the following section. For the same reason as I/O overhead, the impact to the query processing itself is insignificant.

2.2.3 Communication Overhead and High Bandwidth Networks

It is another surprise that higher bandwidth networks such as FDDI (Fiber Distributed Data Interface) or ATM (Asynchronous Transfer Mode) does not significantly reduce communication overhead (in terms of response time).

In general, the communication cost consists of four components: propagation delay, transmission delay, queuing delay, and protocol delay. Propagation delay is negligible on LAN environment; transmission delay is a function of bandwidth; queuing delay depends on the network workload; while protocol delay depends on the complexity of protocols and processing speed. Higher bandwidth network can only reduce transmission delay and queuing delay. Standard communication protocols such as TCP/IP require several layers of program interface before a message can be transferred from the program context to the communication media and vice versa. This procedure may induce a significant overhead due to the memory copying procedure and program execution in each layer. As mentioned above, memory access overhead can not be ignored and, therefore, communication overhead cannot be ignored even if a higher bandwidth LAN is used.

2.3. One-Level Local Multi-access Networks

On a one-level local multi-access network, any message transmitted to the network can be received by every site on the network directly with approximately the same communication overhead. A message can be received by all sites (broadcast) or some selected sites (multicast) simultaneously. The transmission of messages is controlled by an access control protocol. Detailed physical characteristics can be easily found elsewhere. Their logical properties with respect to the data intensive distributed computing applications are as follows.

Because the communication overhead may depend on processing speed, we separate the protocol processing overhead from communication overhead in our model to obtain a more accurate model. However, in real systems, there is no practical way to measure them separately. It would be more practical to measure the communication delay regardless of how it is constituted. A simplified model that considers only communication overhead can be used, which will be described in Section 4, although it may be less accurate.

2.4. Homogeneous Two-Level Local Multi-access Networks

The network concerned in this paper consists of three components: a backbone network, a set of local subnets, and a set of homogeneous local systems (sites). The per-unit disk-access costs and the processing costs are identical for all sites. The backbone network and all local subnets are one-level local multi-access networks as described in the last section, except that the backbone network may have a higher data rate. The networks are organized into a two-level hierarchy: each local system is connected to a local subnet; and all local subnets are connected to

1. The per-unit cost of inter-site communication is independent of source and destination sites due to the multi-access/broadcast capability.
2. The communication overhead is assumed proportional to the volume of data transmitted. (Shorter controlling messages can be ignored.)
3. The communication cost of updating a file is independent of the number of times that file is replicated (see section 2.2.1).
4. The per-unit disk-access costs and the processing (CPU) costs are identical for all sites (homogeneous environments).
5. Synchronization and consistency maintenance, such as concurrency control, is simplified because the message order is preserved in all sites. (Although the FAP does not depend on this property, it is worth to mention that this property is very useful in a distributed environment. It makes the synchronization, concurrency control, and dynamic query processing much easier than that in point-to-point networks (Gouda and Dayal, 1981; Hevner et al., 1985; Kerschberg et al., 1982; Lien, 1986; Lien and Wah, 1987; Masuyama et al., 1983; Sacco, 1984; Wah and Lien, 1985).)
6. The overhead to update a file is independent of system load. (On point-to-point networks, the update overhead may be very sensitive to the system load due to the extra synchronization needed for increased resource contention. However, on local multi-access networks, simplified concurrency control protocols that are insensitive to the system load by taking the advantage of broadcast/multicast capability and property (5) has been found (Wah and Lien, 1985).)

the backbone network through gateways (or routers).^{*} Properties (1) and (5) of one-level networks are not preserved in two-level multi-access networks, but many problems involving distributed control can still be simplified. This will be discussed later in the paper. An example of a two-level local multi-access network is shown in Figure 1.

Since the transmissions on different subnets are independent, the total bandwidth is higher than the bandwidth of a single local subnet. A message may have to be stored in the buffer storage at a gateway before it can be forwarded from one network to another. The local subnet to which a particular site is connected is called the *home subnet* of the site. With respect to this site, other local subnets are *remote subnets*. A query (resp., update) initiated from a site is called a *local query* (resp., *update*) if it is directed to local files, a *neighboring query* (resp., *update*) if it is directed to files in another site in the home subnet, and an *internet query* (resp., *update*) if it is directed to files in a remote subnet. Neighboring and internet queries (resp., update) are *remote queries* (resp., *updates*). With respect to a file f , the sites (resp. local subnets) that are allocated with the file are referred to as the *f-filed sites* (resp., *f-filed subnets*) or *filed sites* (resp., *filed subnets*) in short.

Only one transmission is needed to transmit a message from one site to another site on the same local subnet. Processors in both sites and the home subnet are involved in this transmission. Three transmissions are needed to transmit a message from a given site to a remote subnet: one on the home subnet, one on the backbone network, and one on the target subnet. The message is first transmitted to the home subnet in which the sender resides. It is then forwarded to the backbone network, through one gateway, and, finally, to the remote subnet through another gateway. Processors in both sites, the two local subnets, the backbone network, and the gateways for the two local subnets are involved in this transmission.

2.5. A Model for File Allocation Problems

Assuming homogeneity, the disk overhead to access or modify a file is independent of the location of the file. Although the disk overhead depends on the number of file copies in an update, it is usually insignificant when compared to the overhead in identifying the target data and in executing concurrency control. The problem can therefore be simplified by ignoring the disk overhead (see Section 2.2.1). Thus, local file access has a zero cost since it does not involve any communication overhead.

The cost parameters in the proposed model consist of overheads in processors, local subnets, backbone networks, and gateways. We do not see the cost as a real dollar value but the overhead or time spent. The following notations are defined for the modeling purpose.

- N - set of local subnets
- S - set of sites in the system
- S_n - set of sites in the local subnet n
- F - set of files in the database, $|F| = m$
- (n, s) - site s of subnet $n, s \in S_n$
- λ_{fns} - query load originating at site (n, s) for file f per unit time
- ϕ_{fns} - update load originating at site (n, s) for file f per unit time
- d_α - per unit communication cost on the backbone network for a remote query

^{*} The interface between two communication networks is usually referred to as a *router* if two networks use the same communication protocol; otherwise, it is called a *gateway*.

- d_β - per unit communication cost on a local subnet for a remote query
- d_p - per-unit processing cost for intersite communications in a remote query
- d_g - per-unit gateway processing cost for intersite communications in a remote query
- α - per-unit communication cost on the backbone network in a remote update
- β - per-unit communication cost on a local subnet in a remote update
- d_p - per-unit processing cost for intersite communications in a remote update
- d_g - per-unit gateway processing cost for intersite communications in a remote update
- c_{fns} - storage cost per unit time of file f at site (n, s)
- l_f - length of file f
- C_{ns} - storage capacity at site (n, s)
- I_{ns} - 1 if file f is allocated to site (n, s) ; 0, otherwise
- Y_{fn} - 1 if $\sum_{s \in S_n} I_{ns} \geq 1$; 0, otherwise

Since the protocols for handling queries and updates may be different, their corresponding overheads may also be different. Local queries do not involve networking overhead. A cost of $d_\beta + 2d_p$ is incurred for each data unit in a neighboring query and $(d_\alpha + 2(d_\beta + d_p + d_g))$ for an internet query. The cost of an update is different from that of a query in two aspects. First, the per-unit cost of updates may be higher than that of queries. Secondly, the number of network transmissions is dependent on the number of file copies. Considering the update cost for each data unit, there is no cost if a unique copy is allocated at the site where the update is initiated. When the home subnet is the unique f -filed subnet, the overhead is d_β and a cost of d_p for each f -filed site except the initiation site. The cost for an internet update is more complicated: one transmission on the home subnet, one transmission on the backbone network, and one transmission on each of the remote filed subnets.

Obviously, at least one copy of every file must be allocated in the system. Further, the total size of the allocated files in a site must not exceed the storage capacity of the site. The problem is formulated as follows.

Problem P

$$(1) \text{ Minimize } C(I) = \sum_{n \in N} \sum_{s \in S_n} \lambda_{fns} c_{fns}(I) + \phi_{fns} d_{fns}(I) + \alpha_{fns} I_{fns} \quad (1)$$

$$\text{Subject to } \sum_{n \in N} \sum_{s \in S_n} I_{ns} \geq 1 \quad \forall f \in F \quad (2)$$

$$\sum_{f \in F} I_{fns} \leq C_{ns} \quad \forall n \in N, s \in S_n \quad (3)$$

where

$$d_{fns}(I) = 0 \quad \text{if } I_{fns} = 1$$

$$= d_\beta + 2d_p \quad \text{if } I_{fns} = 0 \text{ and } Y_{fn} = 1$$

$$= d_\beta + 2(d_\beta + d_p + d_g) \quad \text{if } Y_{fn} = 0$$

$$d_{fns}(I) = 1 \text{ and } \sum_{n' \in N} \sum_{s' \in S_{n'}} I_{fn's'} = 1$$

4. FILE ALLOCATION ON COMMUNICATION DOMINANT SYSTEMS

In communication dominant systems, the processing cost is ignored, and $d_{fns}(i)$ and $d_{fns}(i)$ read as follows:

$$d_{fns}(i) = \begin{cases} 0 & \text{if } I_{fns} = 1; \\ d_{\beta} & \text{if } I_{fns} = 0 \text{ and } Y_{fm} = 1 \\ d\alpha + 2d_{\beta} & \text{if } Y_{fm} = 0 \\ 0 & \text{if } I_{fns} = 1; \\ 0 & \text{if } I_{fns} = 0 \text{ and } \sum_{s' \in S_n} I_{fns'} = 1; \\ d_{\beta} & \text{if } \sum_{s' \neq s} I_{fns'} \geq 1 \text{ and } \sum_{n' \in N} Y_{fn'} = 1; \\ d\alpha + (1 + \sum_{n' \neq n} Y_{fn'})d_{\beta} & \text{otherwise} \end{cases}$$

As we will show in Section 4.1, the SFAP in this case is polynomially solvable. The general problem with storage constraints is still NP-hard. A slight modification on the linearization shown in Section 3 can be used here. In Section 4.2, an efficient heuristic solution with guaranteed performance is developed.

4.1. Simple File Allocation

In the following discussion, the index f is dropped because the problem is defined with respect to a single file.

$$C(i) = \sum_{n \in N} \sum_{s \in S_n} (\lambda_{ns} d_{ns}(i) + \phi_{ns} d'_{ns}(i) + \sigma_{ns} I_{ns})$$

In a communication dominant system, the update cost is dependent on the number of filed subnets, but not on the number of filed sites. Only one transmission on each subnet is needed to query or update a file. As shown below, the problem becomes polynomially solvable with respect to the number of sites and the number of local subnets.

To solve the problem, three mutual exclusive allocation alternatives are considered:

- Single-copy allocation:** One copy of the file is allocated to a single site after the optimal placement is determined. In this case, no remote update is needed for changes initiated by the filed site.
- Single-subnet allocation:** Multiple copies of the file are allocated to a single local subnet after the best subnet is determined. No internet update is needed for the changes initiated by the filed site.
- Multiple-subnet allocation:** The file is allocated to more than one local subnet. Update cost depends on the number of filed subnets.

The optimal solution for a given problem can be discovered by comparing all alternatives.

Single-copy allocation. In the following equation for single-copy allocation, τ denotes the total communication cost when all queries and updates are internet. τ is defined as:

$$\tau = \sum_{n \in N} \sum_{s \in S_n} (\lambda_{ns} (d\alpha + 2d_{\beta}) + \phi_{ns} (d\alpha + 2d_{\beta}))$$

Δ_{ns} denotes the reduction of cost if the file is only allocated to site (n, s) . The equation is

$$\Delta_{ns} = \sum_{s' \in S_n} (\lambda_{ns'} (d\alpha + 2d_{\beta}) + \phi_{ns'}) - (d\alpha + 2d_{\beta}) + \lambda_{ns} d_{\beta} + \phi_{ns} d_{\beta} - \sigma_{ns} \quad (12)$$

The first term in Eq. (12) is the reduction of the cost for neighboring queries and updates. The second term is the additional reduction for the requests initiated from site (n, s) . The third term is the storage cost. Obviously, the file should be allocated to the site with the maximum Δ_{ns} . The cost of single-copy allocation becomes:

$$C(i) = \tau - \text{Max}_{n \in N} \sum_{s \in S_n} (\Delta_{ns})$$

Single-subnet allocation. In the single-subnet allocation, where τ remains the same. The cost reduction is

$$\Delta_n = \sum_{s \in S_n} (\lambda_{ns} (d\alpha + d_{\beta}) + \phi_{ns} (d\alpha + d_{\beta})) + \sum_{s \in S_n} (\lambda_{ns} d_{\beta} - \sigma_{ns}) I_{ns} \quad (13)$$

The first term in Eq. (13) is the reduction of the cost for neighboring queries and updates. The second term is the additional cost reduction for local queries. There is no local update since a neighboring update is inevitable for every update made in the home subnet. To maximize Eq. (13) for subnet n , every site in subnet n having positive values of $(\lambda_{ns} d_{\beta} - \sigma_{ns})$ should be allocated with a copy. In the subnets that have zero or one filed site, Eq. (13) is not applicable and such subnets should not be considered in the single subnet allocation since such an allocation will not be better than the single-copy solution. The cost reduction for the filed subnet n is:

$$\Delta_n = \sum_{s \in S_n} (\lambda_{ns} (d\alpha + d_{\beta}) + \phi_{ns} (d\alpha + d_{\beta})) + \sum_{s \in S_n} \text{Max}(0, (\lambda_{ns} d_{\beta} - \sigma_{ns}))$$

The subnet with the maximum Δ_n should be selected as the filed subnet. Every site in that subnet with positive $(\lambda_{ns} d_{\beta} - \sigma_{ns})$ should be allocated a copy. Unless all local subnets are discarded, the cost of a single-subnet multiple-copy allocation is then:

$$C(i) = \tau - \text{Max}_{n \in N} (\Delta_n)$$

Multiple-subnet allocation. In multiple-subnet allocation, the cost of update depends on the number of filed subnets. The equation becomes:

$$d'_{fns}(i) = \begin{cases} 0 & \text{if } I_{ns} = 1 \text{ and } \sum_{n' \in N} \sum_{s' \in S_{n'}} I_{ns'} \neq s \text{ or } I_{ns'} = 0; \\ d\alpha + q d_{\beta} & \text{if } \sum_{s' \neq s} I_{ns'} \geq 1; \\ d\alpha + (1+q)d_{\beta} & \text{otherwise} \end{cases}$$

generally difficult to evaluate. In this section we will discuss several approaches for heuristic algorithm design and propose one with guaranteed performance.

The following approaches all belong to the category of "divide-and-conquer", a methodology frequently used to obtain approximate solutions for complicated combinatorial problems.

Approach 1. Allocate each file independently using simple file allocation algorithms until all files are allocated. It may be necessary to reallocate some files if the storage capacity constraints are violated.

Approach 2. Allocate the first copy of each file to the system to generate an initial solution, then improve the solution by adding more copies. Allocation of extra copies to each site can be solved as a standard 0-1 knapsack problem using either optimal algorithms or heuristic algorithms.

Approach 3. Allocate files to each site as a standard 0-1 knapsack problem (optimally or heuristically), then, identify the missing files and modify the solution by adding missing files and removing redundant files.

Approach 4. A combination of the above approaches.

Although the knapsack problem itself is an NP-hard problem, there does exist pseudo-polynomial algorithms (Gary and Johnson, 1979) and good polynomial time suboptimal algorithms. Since the FAP is a design problem, the time to solve FAP with an pseudo-polynomial optimal algorithms may be tolerable. It is difficult to compare these approaches. In the rest of this section, a heuristic algorithm based on Approach 2 will be shown in the rest of this section. It is simply because we can estimate how good a solution will be.

Heuristic Algorithm FAP-HEUR

(a) An arbitrary site is selected to allocate the first copy of each file such that the total length of all files in a given site is at most $\lfloor \frac{L}{|S|} \rfloor L$, where L is the maximum length of the files. (The single-copy allocation algorithm developed for SFAP in Section 3.1 can be used.)

(b) Additional copies are allocated to each site as a standard 0-1 knapsack problem to fill the remaining storage capacity. In solving the 0-1 knapsack problem, the files are considered the objects; the storage capacity of a site is the capacity of the corresponding knapsack; the length of a file is the weight of the corresponding object; and the profit of allocating a file to site (n,s) is $(\lambda_{ns} d_j - \sigma_{ns})$, which is the cost reduction when file f is allocated to more than one subnet.

For simplicity, we assume that the storage capacity of each site is at least $\lfloor \frac{L}{|S|} \rfloor L$. Otherwise, Step (a) must be modified. Under this assumption, Step (a) is always feasible since the number of files each site can hold is at least $\lfloor \frac{L}{|S|} \rfloor$. In step (b), either optimal or heuristic 0-1 knapsack algorithms can be used as we have mentioned.

where q is the number of filed subnets ($q = \sum_{n \in N} Y_n^m$). The case in which there are exactly q filed subnets, called the *q-subnet-allocation*, can be solved easily. For a particular q , the multiple-subnet allocation can be solved by selecting the best q subnets. (The allocation of each subnet can be solved in a way similar to the single-subnet allocation.) The original multiple-subnet allocation can be solved by comparing all n cases.

Let r_q be the total communication cost when all queries are internet queries, and the copies in all q filed subnets are updated by internet communications in a q -subnet-allocation. Thus,

$$r_q = \sum_{n \in N} \sum_{s \in S_n} (\alpha + 2d_j) + \phi_{ns} (\alpha + (q+1)d_j) \quad (14)$$

If subnet n is a filed subnet, the reduction of the cost for remote queries and updates is:

$$\sum_{s \in S_n} (\alpha + d_j) + \phi_{ns} d_j \quad (15)$$

The first term in the summation is the reduction of the cost for remote queries. The second term is the reduction of the cost for neighboring updates since there are only $(q-1)$ filed subnets. For each copy allocated to site (n,s) , a cost of $(\lambda_{ns} d_j - \sigma_{ns})$ can be saved as well. Obviously, the optimum allocation for each subnet is to allocate a copy to each site that has a positive $(\lambda_{ns} d_j - \sigma_{ns})$. If there is no such site, it is best to allocate a copy to the site with the maximum of $(\lambda_{ns} d_j - \sigma_{ns})$. If $\sum_{s \in S_n} (\alpha + d_j) + \phi_{ns} d_j + \text{Max}_s (\lambda_{ns} d_j - \sigma_{ns}) > 0$, or not to allocate the file to the subnet. If fewer than q subnets can be allocated, the q -subnet-allocation case should be discarded since it is dominated by either the (q) -subnet case, where $q > q$, or by the single-subnet allocation. The optimal solution for the multiple-subnet allocation can be obtained by comparing the costs of all q -subnet allocations. Similar to the single-subnet allocation, the multiple-subnet allocation should be discarded if there are one or fewer filed subnets. The cost of an optimum allocation is then:

$$C(1) = r - \sum_{n \in N} \sum_{s \in S_n} (\alpha + d_j) + \phi_{ns} d_j - \sum_{s \in S_n} (\lambda_{ns} d_j - \sigma_{ns}) Y_n \quad (16)$$

Global optimum solution. A global optimum solution can be obtained by comparing the results obtained from the three different allocations. The time complexity for single-copy allocation and single-subnet allocation is $O(|S|)$, where $|S|$ is the total number of sites. The complexity for multiple-subnet allocation is $O(|S|^{m+1})$, where $|S|$ is the number of local subnets. Therefore, the complexity for the algorithm is $O(|S|^{m+1})$.

The solution algorithm is summarized in algorithm SFAP-OPT in Appendix A. A demonstration of the algorithm is shown in Appendix B.

4.2. Heuristic Solutions for File Allocation Problems

Fast heuristic algorithms are needed when the number of variables is more than what a system can handle in a reasonable amount of time and when the locality changes rapidly. It is not difficult to obtain good heuristics for file allocation problems. Their performances, however, are

available constraints. The relaxed-FAP can be solved optimally by allocating files to each site as a standard 0-1 knapsack problem. The profit of allocating file f to site (n,s) is P_{fns}^* .

We denote P^* , p_h and P^{**} as the total profit generated by an optimal solution of FAP, by algorithm FAP-HEUR, and by an optimal solution of relaxed-FAP, respectively. We also denote P_{fns}^{**} as the profit of site (n,s) as a knapsack problem instance and P_{ph}^{ns} as the profit produced by site (n,s) in algorithm FAP-HEUR. In other words, $P^{**} = \sum_{s \in N} \sum_{f \in F} P_{fns}^{**}$ and $P_h = \sum_{s \in N} P_{ph}^{ns}$. It is easy to prove that $P^{**} - P_h \leq \sum_{s \in N} \sum_{f \in F} (P_{fns}^{**} - P_{ph}^{ns})$.

Denote P_{fns} as the profit produced by allocating files to site (n,s) when the storage capacity of that site is $(1 + \frac{|S|}{|F|})L$ and the profit of allocating a file f to site (n,s) is $max_{-}P_{fns}$ where L is the maximum length of all files.

Theorem 1. The performance degradation of algorithm FAP-HEUR is at most

$$(a) \sum_{n \in N} \sum_{s \in S_n} (P_{fns}^{**} + \sum_{f \in F} \Delta P_{fns}), \text{ or}$$

$$(b) P^* \sum_{n \in N} \sum_{s \in S_n} (\epsilon_{ns} P_{fns}^{**} / P^*) \text{ if } \sum_{f \in F} \Delta P_{fns} > P_{fns}^*, \forall f \in F$$

$$\text{where } \Delta P_{fns} = (max_{-}P_{fns} - min_{-}P_{fns}), \epsilon_{ns} = (q+1)L/C_{ns}, q = \left\lceil \frac{|F|}{|S|} \right\rceil, \text{ and } L = \max\{P_f\}.$$

Proof. The proof for this theorem is shown in Appendix C.

In practice, part (a) is difficult to use since several instances of the knapsack problem must be evaluated. In contrast, part (b) is much easier to use in some cases. To use part (b), the summation of ΔP_{fns} is required to be small when compared to P_{fns}^* . This seems to be too strong

in practice. However, after examining the derivation of Theorem 1, we found that ΔP_{fns} really presents an upper bound of the profit estimation error in allocating file f to site (n,s) . Since we expect that P_{fns}^{**} occurs more frequently than P_{fns}^* , P_{fns}^* , and P_{fns}^* , the real value of $\sum_{f \in F} \Delta P_{fns}$ is probably negligible as compared to P_{fns}^* . If the discrepancy between P^* and $\sum_{f \in F} \Delta P_{fns}$ is small, $\epsilon_{ns} P_{fns}^{**} / P^*$ can be replaced by $\epsilon_{ns} P_{fns}^{**} / P^*$, making the solution even easier. The relationship between the estimated upper bound of performance degradation and the various parameters is demonstrated in Figure 2 on the following conditions:

- (a) P_{fns}^{**} is similar in all sites; (b) $\sum_{f \in F} \Delta P_{fns} < < P_{fns}^*$; and (c) $P_{fns}^{**} \equiv P^*$.

4.4 SIMULATION EXPERIMENTS

To compare the heuristic method and zero-one formulation (Problem Q), we randomly generated multiple instances of various sizes. The details for these data sets are as follows.

4.3. How Good is Algorithm FAP-HEUR?

Although FAP-HEUR is not the best of its type, its performance can be estimated based on the properties of knapsack problems (Lien, 1987). The evaluation is explained as follows. To evaluate the algorithm, the FAP is first transformed into a *multiple-choice-knapsack problem* (MCKKP), in which some classes of objects are to be allocated to a set of knapsacks with a goal of maximizing the total profit. Each class has an unlimited number of identical objects. At most one object from each class can be allocated to a knapsack; and at least one object from each class must be allocated to one of the knapsacks. When allocating an object to a particular knapsack, one of the following four profits is obtained:

$$(a) P_{fns}^*(1) = \sum_{s' \in S_n} (\lambda_{fns}^{s'} d_{\beta} + \phi_{fns}^{s'} d_{\beta}) + \phi_{fns}^{s'} d_{\beta} - \sigma_{fns}$$

$$(b) P_{fns}^*(2) = \sum_{s' \in S_n} (\lambda_{fns}^{s'} (d_{\alpha} + d_{\beta}) + \phi_{fns}^{s'} d_{\beta}) + \phi_{fns}^{s'} d_{\beta} - \sigma_{fns}$$

$$(c) P_{fns}^*(3) = \sum_{s' \in S_n} (\lambda_{fns}^{s'} (d_{\alpha} + d_{\beta})) + (\lambda_{fns}^{s'} d_{\beta} - \sigma_{fns})$$

$$(d) P_{fns}^*(4) = (\lambda_{fns}^{s'} d_{\beta} - \sigma_{fns})$$

We denote, among the four different profits, the maximum, the minimum, and the one with respect to a particular optimal allocation as $max_{-}P_{fns}$, $min_{-}P_{fns}$, and P_{fns}^* , respectively. Each of these profits is nothing more than the cost reduction from the following fixed cost when one object (file) is allocated to a knapsack.

$$\tau = \sum_{n \in N} \sum_{s \in S_n} \lambda_{fns} (d_{\alpha} + 2d_{\beta}) + \phi_{fns} (d_{\alpha} + d_{\beta})$$

Maximizing the profit of an MCKKP instance minimizes the cost of the corresponding instance in FAP. The profit $P_{fns}^*(1)$ is the cost reduction of allocating file f to site (n,s) when it is the only copy allocated. The third term in $P_{fns}^*(1)$ is the communication cost on subnet n for the updates generated by remote subnets, which is not counted in τ . Profit $P_{fns}^*(2)$ is the cost reduction when file f is allocated to site (n,s) where subnet n is the unique f -filed subnet and site (n,s) is the first f -filed site in the subnet. Profit $P_{fns}^*(3)$ is the cost reduction of allocating file f to site (n,s) when there are two or more f -filed subnets and site (n,s) is the first f -filed site in subnet n . Profit $P_{fns}^*(4)$ is the cost reduction of allocating file f to site (n,s) when the conditions in $P_{fns}^*(1)$, $P_{fns}^*(2)$, and $P_{fns}^*(3)$ are not applicable.

After the transformation, the problem becomes the maximization of total profit since (total cost) = τ - (total profit). In the rest of this section, the minimization of overall cost is referred to as the maximization of total profit. The evaluation of FAP-HEUR involves finding the maximum discrepancy between the profit generated by FAP-HEUR and the optimal solution, referred to as the *performance degradation* of FAP-HEUR. The performance degradation is bounded by the discrepancy between the FAP-HEUR and the relaxed-FAP, which is the original FAP without

(1) For $|N| = 2$ and $|S| = 2$, 10 instances for each of the parameter $|F| = 2, 3, 4, 6, 8, 10, 15, 20, 30,$ and 40 were generated. That is, we randomly generated 100 instances in total for different number of files.

(2) For each problem, the coefficients λ_{fns} , ϕ_{fns} , σ_{fns} , ρ_f , and C_{ns} were uniformly determined as follows:

(a) $\lambda_{fns} = U(10, 100)$ (i.e. uniform distribution from 10 to 100.)

(b) $\phi_{fns} = U(1, 50)$

(c) $\rho_f = U(1, 100)$

(d) $\sigma_{fns} = \rho_f \cdot U(0.1, 1) + 1$

(e) $C_{ns} = 100|F|/3 + 100|F|U(0, 1)$

The result of the experiment is shown in Table 1. Since the zero-one formulation is NP-hard, it is not possible to obtain optimal solutions in practicable time when $|F|$ is greater than 10. Nonetheless, the experiment does show that the heuristic performs better when the problem size increases and the heuristic solution can be found in a reasonable time although it is not a polynomial algorithm. It also shows that the computation time of the heuristic is very stable.

Table 1. Experiment results for $|N| = 2$ and $|S| = 2$.

F	Problem Q		HC2BFAP	
	Objective function ¹	CPU Time ³	Objective function ²	CPU Time ⁴
2	100	22 sec.	150	38 sec.
3	100	55 sec.	140	43 sec.
4	100	4 min.	149	45 sec.
6	100	10 min.	122	42 sec.
8	100	35 min.	121	42 sec.
10	100(*)	40 min.	114	43 sec.
15	100(*)	50 min.	109	54 sec.
20	100(*)	60 min.	106	55 sec.
30	100(*)	90 min.	100	68 sec.
40	100(*)	120 min.	96	93 sec.

1. 10 problems average, normalized to 100%
 2. 10 problems average, compared to Problem Q solutions
 3. Run on a VAX 8600 with a VMS operating system
 4. Run on a Paramid 98x with a UNIX 4.2BSD operating system
 *. Not solved to optimal

5. IMPLEMENTATION CONSIDERATIONS

5.1. Applicability of The FAP Model

Not every distributed computing system can benefit from FAP. The model is best suited to the environment in which most of the files are shared by a number of users such as database systems. In the environment in which most of the files are private files, a simple strategy such as cache-on-demand may be better than the complicated integer programming model. A careful evaluation must be made before the model is actually applied to a real system.

5.2. Dynamic File Allocation

The instantiation of the model seems to be a serious problem in applying the FAP model to a real system. A real system is very unlikely to have stable access rates for the whole period between two consecutive file reorganizations. Furthermore, the accuracy of access rate estimation presents another problem in applying the model, especially, in the local area environments where the user/process/data mobility is higher than wide-area-networks as mentioned in Section 1. However, the dynamic file allocation strategy, known as the *file migration problem* (FMP), to relocate the files periodically in response to the changes in the file access pattern is very useful to alleviate the problem. The basic model of FMP is to reallocate the files as a file allocation problem at the end of each time "window", which is a fixed length time interval, based on the observed references during the last window (Gavish and Liu Sheng, 1990a, 1990b; Levin, 1982; Levin and Morgan, 1978; Liu Sheng, 1986, 1992a, 1992b; Porcar, 1982; Segall, 1976; Segall and Sandell, 1979).

5.3. File Preallocation and Reduction of Problem Size

In practice, the large number of files may make FAP impossible to solve optimally in a reasonable amount of time. Optimal algorithms, however, are still useful for the following reasons:

- (a) Optimal algorithms provide a boundary within which one can arrive at approximate solutions. For example, the optimal algorithm we have proposed earlier may be modified to use a heuristic knapsack algorithm instead of an optimal algorithm.
- (b) For those systems operating under steady state, it is not necessary to solve the file allocation problem in real time. The allocation is only needed when the file system or the database is reorganized. Therefore, a long execution time to solve the FAP is tolerable as long as the ratio of the execution time to the mean time between reorganization is reasonable.

Problem size can be reduced by allocating the following files separately: (1) rarely used files; (2) files with strong locality in a few sites; and (3) frequently-used files with weak locality. For rarely used files, a single copy should be allocated. Files in category (2) should have a copy allocated to each site. For the files in category (3), replicating files in each system can be beneficial. Simple allocation heuristics such as a best-first search are useful in these preallocations.

6. CONCLUSION

The file allocation problem on homogeneous local computer systems with two-level local multi-access networks is studied in this paper. The problem is less complicated on two-level networks than on general point-to-point networks because the communication cost of remote file access is reduced from a random matrix to a three-value variable. The general file allocation with storage-limit constraints is still an NP-hard problem. It is transformed into a linear zero-one integer programming problem which is solved by a search method based on Balas's additive algorithm. In the case in which the system is dominated by the communication cost, the simple file allocation problem is solved in $O(|S||N|)$ time complexity, where $|F|$ is the total number of sites in the network and $|N|$ is the total number of local subsets in the system. A simple heuristic algorithm based on the 0-1 knapsack problem is proposed with analytical and empirical evaluations. The performance degradation of the proposed heuristic algorithm is shown to be small in both evaluations.

References

1. Aperi, Peter, "Data Allocation in Distributed Database Systems," *ACM Trans. on Database Systems*, vol. 13, No. 3, pp. 263-304, Sept. 1988.
2. Balas, Egon, "An Additive Algorithm for Solving Linear Programs with Zero-one Variables," *Operating Research*, vol. 13, pp. 517-549, July-August 1965.
3. Berman, L., and Nigam, R., "Optimal Partitioning of Data Bases across Multiple Servers in a LAN," *Interface*, vol. 22, no. 2, pp. 18-27, March-April, 1992.
4. Chu, W. W., "Multiple File Allocation in a Multiple Computer System," *IEEE Trans. on Comp.*, vol. C-18, no. 10, pp. 885-889, Oct. 1969.
5. Dowdy, L. W. and D. V. Foster, "Comparative Models of the File Assignment Problem," *ACM Computing Surveys*, vol. 14, no. 2, pp. 287-313, June 1982.
6. Fortin, Michael, S. J. Kao, and Douglas Kerr, "Benchmarking Workstations from the User's Perspective: Taking into Account the Environment," Working Paper, Dept. of Comp. and Info. Sci., The Ohio State Univ., Columbus, Ohio, Sept. 8, 1988.
7. Garey, Michael R. and David S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, Freeman, Reading, England, 1979.
8. Gavish, B., and Liu Sheng, O. R., "Adaptive Heuristics for Dynamic File Migration Policies in Distributed Computer Systems," Working Paper, University of Arizona, 1990a.
9. Gavish, B., and Liu Sheng, O. R., "Dynamic File Migration Policies in Distributed Computer Systems," *Comm. ACM*, vol. 33, no. 2, pp. 177-189, 1990b.
10. Gouda, M. G. and U. Dayal, "Optimal Semi-join Schedules for Query Processing in Local Distributed Database Systems," *Proc. ACM SIGMOD Conference*, pp. 164-175, May 1981.
11. Hevner, A. R., O. Q. Wu, and S. B. Yao, "Query Optimization on Local Area Networks," *ACM Trans. on Office Information Systems*, vol. 3, no. 1, pp. 35-62, Jan. 1985.
12. Kerschberg, Larry, Peter D. Ting, and S. Bing Yao, "Query Optimization in Star Computer Networks," *ACM Transactions on Database Systems*, vol. 7, no. 4, pp. 678-711, Dec. 1982.
13. Levin, K. D. and H. L. Morgan, "A Dynamic Optimization Model For Distributed Databases," *Oper. Res.*, vol. 26, no. 5, pp. 824-835, Sept.-Oct. 1978.
14. Levin, K. D., "Adaptive Structuring of Distributed Database," *Proc. National Computer Conf.*, pp. 691-696, 1982.
15. Lien, Yao-Nan, *Distributed Databases on Local Multi-access Computer Systems*, School of Electrical Engineering, Purdue University, Aug. 1986.
16. Lien, Yao-Nan, "Some Properties of 0-1 Knapsack Problems," *Conference on Combinatorics and Complexity*, p. 105, Chicago, IL, May 1987.
17. Lien, Yao-Nan and Benjamin Wah, "Design and Performance Study of DDRLM: A Distributed Database on a Local Computer System," *Proceedings of HICSS-20*, Volume 2, pp. 407-418, Kona, Hawaii, Jan. 1987.
18. Lien, Y. N., Y. L. Chang, and B. W. Wah, "File Allocation Problems on Homogeneous Two-level Local Broadcast Networks," *Proc. of Fourth Int'l Conf. on Data Engineering*, pp. 92-99, Los Angeles, CA, Feb. 1988.
19. Liu Sheng, O. R., *Models for Dynamic File Migration in Distributed Computer Systems*, Ph.D. Dissertation, W. E. School of Business Administration, Univ. of Rochester, New York, 1986.
20. Liu Sheng, O. R., "Analysis of Optimal File Migration Policies in Distributed Computer Systems," *Management Science*, Vol. 38, No. 4, April, 1992a, pp. 459-482.
21. Liu Sheng, O. R., "Optimization of File Migration in Distributed Computer Systems," *Computer and Operation Research*, (forthcoming), 1992b.
22. Masuyama, S., S. Muro, T. Mizutani, T. Ibaraki, and T. Hasegawa, "Shortest Semijoin Schedules for Local Area Distributed Database Systems," *Proc. 16th Annual Hawaiian Int'l Conf. on System Sciences*, pp. 284-293, IBBE, 1983.
23. Morris, James H., Mahadev Satyanarayanan, Michael H. Connor, John H. Howard, David S. H. Rosenthal and F. Donelson Smith, "ANDREW: A Distributed Personal Computing Environment," *Commun. of the ACM*, vol. 29, no. 3, pp. 184-201, March 1986.
24. Nguyen, Gia Toan, "Distributed Query Management for a Local Network Database System," *Proceedings of the Second International Conference on Distributed Computing Systems Paris*, April 1981.
25. Notkin, David, Norman Hutchinson, Jan Saniolo, and Michael Schwartz, "Heterogeneous Computing Environments: Report on The ACM SIGOPS Workshop on Accommodating Heterogeneity," *Commun. ACM*, vol. 30, no. 2, pp. 132-140, Feb. 1987.

APPENDIX A

```

/* Algorithm SFAP-OPT: to solve SFAP
on two-level local multi-access networks.
-----*/
max(Δ = -∞
FORALL n IN N DO
FORALL s IN Sn DO
IF (Δns > maxΔ) /* Eq. (12) */
THEN { maxΔ = Δns; sc.allocation = (n,s) }
}
sc_cost = t - maxΔ
/* (Single-copy Allocation) */
/* (Single-network Allocation) */
maxΔ = -∞
FORALL n IN N DO
{
number_of_copy = 0
FORALL s IN Sn DO
{
Δn = Σs (Δns (dα + dβ) + φns (dα + dβ)) /* Eq. (13) */
FORALL s IN Sn DO {
saving = (Δns dβ - σns)
IF (saving > 0)
THEN {
Δn = Δn + saving
allocation[n] = allocation[n] ∪ (n,s)
number_of_copy = number_of_copy + 1
}
}
}
}
IF (number_of_copy <= 1)
THEN { Δn = -∞; allocation[n] = ∅; }
IF (Δn > maxΔ)
THEN { maxΔ = -∞; sn_allocation = allocation[n] }
}
}
sn_cost = t - maxΔ
/* (Multiple-network Allocation) */
mn_cost = ∞
FOR q FROM 2 TO N DO
{
number_of_subnet = 0
FORALL n IN N DO
{
number_of_copy = 0
}
}
}

```

26. Porcar, J. M., *File Migration in Distributed Computer Systems*, Ph.D. Dissertation, Lawrence Berkeley Lab, Univ. of California, Berkeley, July 1982.
27. Ramamoorthy, C. V. and B. W. Wah, "The Isomorphism of Simple File Allocation," *IEEE Trans. on Computers*, vol. C-32, no. 3, pp. 221-232, March 1983.
28. Sacco, G. M., "Distributed Query Evaluation in Local Area Networks," *Proc. of Int'l Conf on Data Engineering*, pp. 510-516, Los Angeles, CA, April 1984.
29. Segall, A., "Dynamic File Assignment in a Computer Network," *IEEE Trans Auto, Control*, vol. AC-21, no. 2, April 1976.
30. Segall, A. and N.R. Sandell Jr., "Dynamic File Assignment in a Computer Network - Part II: Decentralized Control," *IEEE Trans Auto, Control*, vol. AC-24, no. 5, Oct. 1979.
31. Smith, A. J., "Analysis of Long Term File Reference Patterns for Application to File Migration Algorithms," *IEEE Trans. Soft Eng.*, vol. 7, no. 4, pp. 403-417, July 1981.
32. Smith, A. J., "Long Term File Migration: Development and Evaluation of Algorithms," *Comm. ACM*, 24 (1981b), pp. 521-532.
33. Smith, A. J., "Optimization of I/O Systems by Cache Disks and File Migration: A Summary," *Performance Evaluation*, 1 (1981c), pp. 249-262.
34. Wah, B. W., "File Placement on Distributed Computer Systems," *IEEE Computer*, vol. 17, no. 1, pp. 23-32, Jan. 1984.
35. Wah, B. W. and Y. N. Lien, "The File-Assignment and Query-Processing Problems in Local Multiaccess Networks," *Proc. of Int'l Conf. on Data Engineering*, pp. 228-235, Los Angeles, CA, April 1984.
36. Wah, B. W. and Y. N. Lien, "Design of Distributed Databases on Local Computer Systems with A Multiaccess Network," *IEEE Trans. on Software Eng.*, vol. SE-11, no. 7, pp. 606-619, July 1985.

APPENDIX B

An example showing the simple file allocation algorithm

$d\beta$	$d\beta'$	$d\alpha$	$d\alpha'$
1	2	3	4

	$n=1$	$n=2$
	$s=1$	$s=3$
	$s=2$	$s=4$
C	50	40
	50	50
$ F =1$	λ	λ
	6	7
	7	4
	4	5
	3	6
	5	2
$ F =2$	ϕ	ϕ
	7	3
	1	4
	2	5
	σ	σ
	2	1
$L=20$	σ	σ
	2	4

File 1 is to be allocated in the simple file allocation.

Single-copy allocation:

$\tau = 230$
 $\Delta n_s = [115, 120, 78, 81]$
 The file is allocated to site 2 with $C(I) = 110$.

Single-network allocation:

$\tau = 230$
 $\Delta n = [111, 75]$
 The file is allocated to site 1 and 2 in network 1 with $C(I)=119$.

Multiple-network allocation:

$\tau_q = 260$

sites 1,2 in network 1 and sites 3,4 in network 2 are allocated with $C(I)=134$.
 Finally, three cases are compared and the single copy is allocated to site 2 with a cost 110.

$$\Delta q_n = \sum_s \in S_n (\lambda_{ns} (d\alpha + d\beta) + \phi_{ns} d\beta) \quad /* \text{ Eq. (15) } *$$

```

maxΔ = -∞
FORALL s IN Sn DO {
    saving = (λns dβ - σns)
    IF (saving > 0)
        THEN {
            Δqn = Δqn + saving
            net_allocation[n] = net_allocation[n] ∪ {n,s}
            number_of_copy = number_of_copy + 1
        }
    IF (maxΔ < saving < 0 AND number_of_copy = 0)
        THEN { maxΔ = saving; Smax = {n,s} }
    /* s */
    IF (number_of_copy = 0)
        THEN {
            Δqn = Δqn + maxΔ
            IF (Δqn > 0)
                THEN {
                    net_allocation[n] = net_allocation[n] ∪ Smax
                    number_of_subnet = number_of_subnet + 1
                }
            ELSE number_of_subnet = number_of_subnet + 1
            /* n */
        }
    IF (number_of_subnet ≥ q)
        THEN {
            picked_subnet = first q subnet in descending order of Δqn
            cost[q] = ∑s ∈ Ns (λns (dα + 2dβ) + φns (dα + (q+1)dβ)) /* Eq. (14) */
            FORALL n IN picked_subnet DO {
                q_allocation[q] = q_allocation[q] ∪ net_allocation[n]
                cost[q] = cost[q] - Δqn
            }
            ELSE cost[q] = ∞
            THEN { mn_cost = cost[q]; mn_allocation = q_allocation[q] }
            /* q */
        }
    /* Compare all cases */
}
Select the minimum among sc_cost, sn_cost, and mn_cost

```

APPENDIX C

Proof of Theorem 1

The properties found in Lien (1987) are very useful to prove Theorem 1 shown in Section 4.

C.1. Some Properties of 0-1 Knapsack Problems.

The standard 0-1 knapsack problem is as follows: given a set of objects $N = \{1, \dots, n\}$, the associated positive weights and profits, $\{w_i | i \in N\}$, $\{p_i | i \in N\}$, and a positive capacity c , find x_1, x_2, \dots, x_n so as to maximize the total profit as follows:

$$Z = \max \left\{ \sum_{i=1}^n p_i x_i \mid \sum_{i \in N} w_i x_i \leq c; x_i \in \{0, 1\}; i \in N \right\}.$$

For easy to refer, the instance of allocating objects N to a knapsack with a capacity c is denoted as (c, N) ; an optimal allocation of the instance and the associated total profit are denoted as A_c^N and Z_c^N respectively.

Theorem C1 described below is useful in answering the following question: given a knapsack instance (c, N) , what will be the maximum profit degradation if the objects in \mathcal{Q} , a subset of N , is preallocated in the knapsack? The following notations are defined for Theorem C1.

$$w_{\mathcal{Q}} = \sum_{i \in \mathcal{Q}} w_i.$$

$$Z_{\mathcal{Q}}^N$$

$$Z_{\mathcal{Q}} = [Z_{c-w_{\mathcal{Q}}}^{c-w_{\mathcal{Q}}} + \sum_{i \in \mathcal{Q}} p_i].$$

where $Z_{c-w_{\mathcal{Q}}}^{c-w_{\mathcal{Q}}}$ is the maximum profit of allocating the remaining objects $(N-\mathcal{Q})$ to the knapsack with the remaining capacity $(c-w_{\mathcal{Q}})$ when all the objects in \mathcal{Q} have been allocated in the knapsack already. Actually, $Z_{\mathcal{Q}}$ can be restated as follows.

$$Z_{\mathcal{Q}} = \max \left\{ \sum_{i=1}^n p_i x_i \mid \sum_{i \in N-\mathcal{Q}} w_i x_i \leq c; x_i \in \{0, 1\}; i \in N-\mathcal{Q}; x_i = 1, \forall i \in \mathcal{Q} \right\}.$$

Theorem C1:

For any knapsack instance (c, N) , and \mathcal{Q} , a subset of N ,

$$(a) \quad Z_{\mathcal{Q}+Z}^{c-w_{\mathcal{Q}+Z}} \geq Z_1,$$

$$(b) \quad Z_{\mathcal{Q}} \geq \left[1 - \frac{c}{w_{\mathcal{Q}+k}} \right] Z_1, \text{ and}$$

$$(c) \quad \text{either } \frac{Z_{\mathcal{Q}}}{c-w_{\mathcal{Q}}} > \frac{c}{c-w_{\mathcal{Q}}} \text{ or } \frac{Z_{\mathcal{Q}}}{c-w_{\mathcal{Q}}} \equiv \frac{c}{c-w_{\mathcal{Q}}}, \text{ if } c-w_{\mathcal{Q}} \gg k,$$

$$\text{where } k = \max \{w_i\} \leq c, \text{ where } i \in N$$

Theorem C2 described below is useful in answering questions such as the following. Imagine that everyone in a supermarket contest is to fill a 10 lb. bag with food. The price tags were taken away and everyone has to make his best guess. The one who packs the highest dollar value of food wins. Assume one solves the problem using an optimal knapsack algorithm, then what is the discrepancy between the value he actually earns and the value he could earn if he knew the prices in advance? This scenario can be thought as a 0-1 multiple-choice knapsack problem solved as a standard 0-1 knapsack problem. Therefore, the properties in Theorem C2 are useful to estimate this discrepancy.

A formal description of the problem is as follows. Given two problem instances (c, N_1) and (c, N_2) , assume that the objects in N_1 and N_2 are the same with different profits, (the objects in N_1 and N_2 are considered the same if they have the same index), what would be the possible profit degradation as compared to the maximum profit of (c, N_1) if the knapsack is filled with N_1 in such a way: an object of N_1 is allocated to the knapsack if the corresponding object in N_2 is in $A_c^{N_2}$, which is an optimal allocation of (c, N_2) .

More notations for Theorem C2 are defined as follows.

$$N_t \quad \text{the } t\text{-th set of objects, } N_t = \{i_t | i_t = 1, 2, \dots, n\}, t \in \{1, 2\};$$

$$w_{i_t} \quad \text{the weight of object } i_t;$$

$$p_{i_t} \quad \text{the profit of object } i_t;$$

$$A_{2-1} \quad \{i_1 | i_1 \in N_1; i_2 \in A_c^{N_2}\}, \text{ which is the set of objects of } N_1 \text{ corresponding to the}$$

$$Z_{2-1} \quad \sum_{i_1 \in A_{2-1}} p_{i_1}.$$

Theorem C2:

Given two knapsack problem instances (c, N₁) and (c, N₂), |N₁| = |N₂| = n, if w_{i1} = w_{i2}, i = 1, 2, ..., n, then

$$(a) \quad Z_{2-1} + \sum_{i=1}^n \max\{0, (p_{i2} - p_{i1})\} \geq Z_{2-1}^c$$

$$(b) \quad Z_{2-1} + \sum_{i=1}^n |p_{i1} - p_{i2}| \geq Z_{2-1}^c$$

C.2. Proof of Theorem 1

Define the basic-relaxed-FAP as the relaxed FAP except that the profit of allocating file *f* to site (*n*, *s*) is (A_{ns}δ_f - σ_{ns}). Denote the maximum profit of a basic-relaxed-FAP and the maximum profit associated with each site as P^{***}_{ns} and P^{***}_{ns} respectively. Therefore,

$$P^{**} = \sum_{n \in N} \sum_{s \in S_n} P^{***}_{ns} \geq P^{***} = \sum_{n \in N} \sum_{s \in S_n} P^{***}_{ns} \geq P^{ph} = \sum_{n \in N} \sum_{s \in S_n} (P^{***}_{ns} - p^{ph}_{ns}).$$

(a.1) From Theorem 2C, we know P^{***}_{ns} - P^{***}_{ns} ≤ ΔP_{ns}.

(a.2) From Theorem 1C, we get P^{***}_{ns} - P^{***}_{ns} ≤ P^{***}_{ns}.

(a.3) From (a.1) and (a.2), we get P^{***}_{ns} - P^{***}_{ns} ≤ P^{***}_{ns} + P^{***}_{ns}.

(b) Since the first step of algorithm FAP-HEUR is to allocate files to each site with a size at most $\left\lfloor \frac{|L|}{|S|} \right\rfloor L$ so that W_Q is the total length of the files allocated to a site in the first step. From Theorem C1 and the assumption that ΔP_{ns} << P^{***}_{ns}, we know that P^{***}_{ns} - P^{***}_{ns} = P^{***}_{ns} - P^{***}_{ns} ≤ P^{***}_{ns} $\frac{C_{ns}}{W_Q + k}$, where P^{***}_{ns} is the part of P^{***} that is associated with site (*n*, *s*). (That is, P^{***} = $\sum_{n \in N} \sum_{s \in S_n} P^{***}_{ns}$.) It is not difficult to prove that:

$$P^{**} - P^{ph} \leq P^{**} (\sum_{n \in N} \sum_{s \in S_n} (\epsilon_{ns} \frac{P^{***}_{ns}}{P^{**}})),$$

$$\text{where } \Delta P_{ns} = (max-p_{ns} - min-p_{ns}), \epsilon_{ns} = \frac{C_{ns}}{(q+1)L}, \text{ and } q = \left\lfloor \frac{|L|}{|S|} \right\rfloor.$$

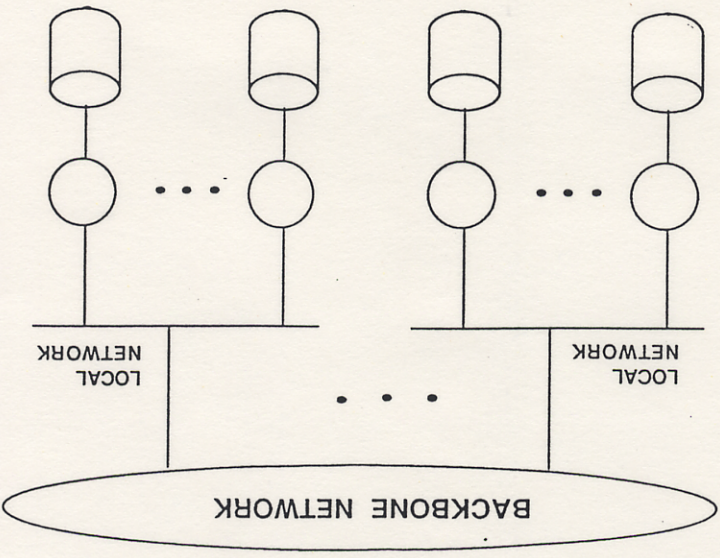


Figure 1 An example of Two-level Local Broadcast Networks.

HCBFAP ERROR ESTIMATION

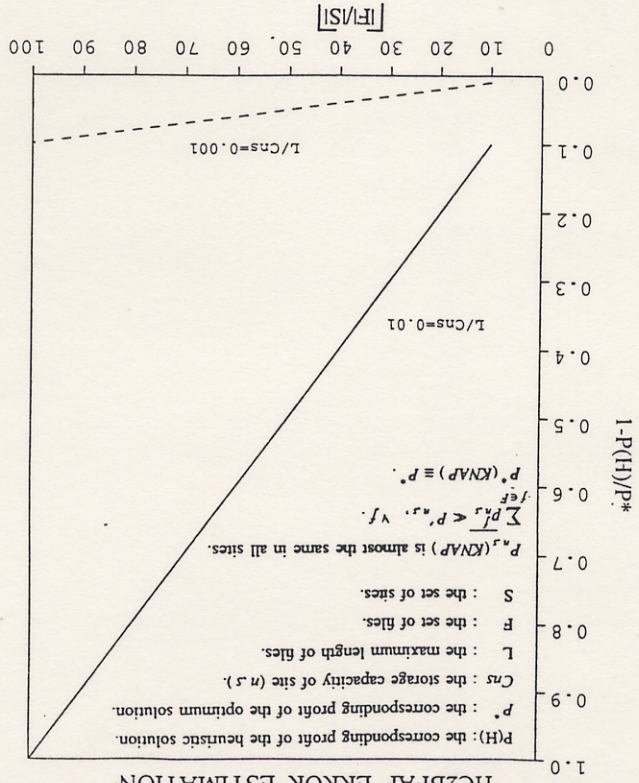


Figure 2 Error Estimation of Algorithm HCBFAP.