

An Analytical Approach to File Prefetching

Technical Report CUCS-031-96

Hui Lei and Dan Duchamp

Computer Science Department
Columbia University
New York, NY 10027
{lei, djd}@cs.columbia.edu

Abstract

File prefetching is an effective technique for improving file access performance. In this paper, we present a file prefetching mechanism that is based on on-line analytic modeling of interesting system events and is transparent to higher levels. The mechanism, incorporated into a client's file cache manager, seeks to build semantic structures, called access trees, that capture the correlations between file accesses. It then heuristically uses these structures to represent distinct file usage patterns and exploits them to prefetch files from a file server. We show results of a simulation study and of a working implementation. Measurements suggest that our method can predict future file accesses with an accuracy around 90%, that it can reduce cache miss rate by up to 47% and application latency by up to 40%. Our method imposes little overhead, even under antagonistic circumstances.

1 Introduction

This paper reports the effectiveness of a file cache management technique that prefetches from a file server automatically and without any sort of user-supplied information or intervention. The technique, which is incorporated into the client's cache manager, makes extra requests to the server, hopefully in advance of the actual need for the data. Prefetched data is then placed in the client's cache.

We hypothesize that there is pronounced regularity in file access patterns and that relatively simple algorithms can identify a file access pattern and quickly spot it when it re-emerges later during another run of the application. In particular, we seek to uncover the semantic structures underlying file accesses. These semantic structures, called *access trees*, capture potentially useful information concerning interrelationships and dependencies between files. An access tree for a program records all the files accessed during one execution of the program. By intercepting relevant system calls, we are able to build an access tree for each program execution. A number of access trees are maintained in virtual memory for each application representing distinct file usage pat-

terns. When an application is re-executed, we compare the access tree being formed by current activity against saved access trees to determine which usage pattern, if any, is recurring; we then prefetch the files remaining in the saved access tree.

File prefetching brings two major advantages. First, applications run faster because they hit more in the file cache. Second, there is less "burst" load placed on the network because prefetching is done only when there is network bandwidth available rather than on demand. On the other hand, there are two main costs of prefetching. The first is the CPU cycles expended by the client in determining when and what to prefetch. Cycles are spent both on overhead in gathering the information necessary to make prefetch decisions, and on actually carrying out the prefetch. The second cost is the network bandwidth wasted when prefetch decisions inevitably prove less than perfect.

Higher cache hit rates are generally good for application performance, but especially so if the cache is small and/or the network is slow. Therefore, we expect our technique to be particularly useful for mobile computers operating over low-bandwidth links. By prefetching soon-to-be-needed files during slack periods when the link is not in use, we can create the illusion that the link has lower latency than is actually the case.

Section 2 details the prefetching mechanism. Section 3 reports results from a trace-driven simulation. Section 4 describes the implementation and presents some initial performance data. Section 5 discusses related work.

2 Mechanism

2.1 Data Abstraction: Access Tree

Every program can invoke other programs. In UNIX-style operating systems, this is usually materialized by the executing program forking child processes, which in turn execute other programs. The programs may also open some data files. All the file accesses in an application can be formulated into a tree data structure, dubbed an access tree. Files (program files and data files) are the nodes, and an edge is drawn from parent A

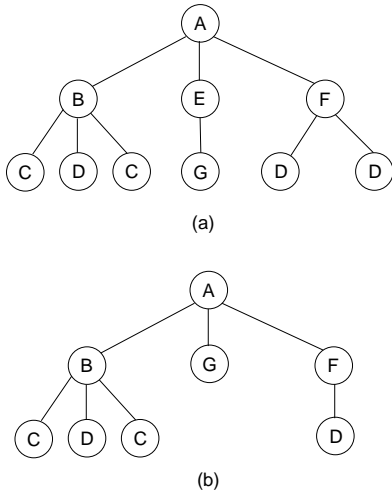


Figure 1: **Example Access Trees.** These graphs show the access trees generated from the events described in section 2. Graph (a) shows the version before compression. Graph (b) that after both vertical and horizontal compression, assuming E is a shell program.

to child B if either (1) program A invokes program B or (2) program A opens data file B. The order of siblings reflects the chronology of file accesses.

Figure 1(a) depicts the access tree for an application A that includes the following activities:

1. Program A invokes program B
2. B opens data files C and D, in that order
3. B opens C again
4. A invokes programs E and F
5. F opens D
6. E invokes program G
7. F opens D again

Two kinds of compression are applicable to an access tree. Vertical compression draws edges through UNIX shells, effectively cutting them out of the access tree. This is necessary because of the role shells play as command interpreters. Shells are invoked in a variety of circumstances and can generate a large number of file usage patterns. Since we perform file prefetching for each program file in the access tree and it is infeasible to prefetch accurately for shells, we choose to ignore them. Horizontal compression removes consecutive accesses of the same data file. The detail of consecutive accesses offers no help in prefetching and can be safely omitted. However, non-consecutive accesses of the same data file are preserved for use in a later phase. Assuming program E is a shell in the previous example, Figure 1(b) shows the access tree after compressions.

As an application proceeds, we construct an access tree for it by intercepting `fork`, `execve`, `open`, `chdir` and `exit` system calls. `Fork` and `execve` substantiate access of an executable file, while `open` that of a data file. Information on `chdir` calls enables us to resolve relative path names of files. The access tree is completed upon `exit` of the program execution. An access tree that is being constructed is called a *working tree*; a finished access tree that is saved to exemplify a file usage pattern is called a *pattern tree*.

2.2 The Prefetch Algorithm

For now we assume the existence of a heuristic function `compatibility(w_tree, p_tree)` which returns an indication between 0 and 1 of the degree to which a working tree `w_tree` and a pattern tree `p_tree` describe the same file access pattern. When the function's value is above constant `MATCH_THRESHOLD`, we say that the two trees *match*.

As mentioned earlier, a number of pattern trees are saved in virtual memory for each program; a working tree is constructed in the course of every program execution. Whenever a program references a file, a new child node is added in the working tree for the program, and some analysis is performed to find out whether any saved pattern trees can be prefetched. Our analysis follows a simple guideline: if there is no previously prefetched pattern tree or the current working tree no longer matches the prefetched pattern tree, compute the compatibility of the working tree and each of the pattern trees for the program. If any match is found, then prefetch the pattern tree with the highest compatibility. If more than one pattern tree bears the highest compatibility, then prefetch the one most recently saved.

The above analysis is carried out for each executable file inside the working tree whenever the executable initiates a file access. At this point, a pattern tree may have already been prefetched for the executable, either as the result of prefetch analysis incurred by earlier file accesses or in the form of a subtree of the larger tree we prefetched at a higher level. The executable can always prefetch another pattern tree, based on the analysis result. This effectively allows minor prefetch corrections to be made, reducing the cost of a bad guess.

Two complications may arise when the pattern tree selected for prefetching is large: prefetched files are evicted from the cache before they are actually referenced, and the cache is considerably destroyed when the prefetching guess is bad. To diminish the extent of these problems, we place an upper limit (`PREFETCH_CAPACITY`) on the number of files from a pattern tree we shall prefetch at one time. When the pattern tree selected is too large, we only prefetch those files in the initial portion of the tree so that the prefetch limit is not exceeded. We also record the immediate child node of the pattern tree we prefetch last. When later the working tree extends to this child, we will prefetch the remaining portion of the pattern tree, if the latest working tree still matches it.

On program exit, we compare the newly completed working tree with the saved pattern trees. If it doesn't match any of the pattern trees, the working tree is saved as new information. Otherwise, it is substituted for the pattern tree that it matches best.

We set the two algorithmic parameters to the following values:

- `MATCH_THRESHOLD`: 0.4
- `PREFETCH_CAPACITY`: 15

We have found that the behavior of our mechanism is insensitive to the value of these parameters. We shall illustrate this later.

2.3 Compatibility Computation

We now take a closer look at the `compatibility` function. This function abstracts out the complexities of prefetch analysis and pattern tree maintenance. We illustrate the definition with the example in Figure 2. A pattern tree is shown in 2(a).

Let us first consider the case that the working tree is a finished one. We try to pair up the immediate child nodes in the working tree with identical child nodes in the pattern tree, preserving the order of the nodes. Recall that a child node can represent either an executable file, which may root another access tree, or a data file. For the two trees to be considered to describe the same file access pattern, we require that there be a one-to-one correspondence between all the executable files, but not all the data files. However, the same executable file may root a different access subtree in the working tree than in the pattern tree. We define C_d as the percentage of data files that can be paired up, C_e as the percentage of pairs of executables that root the same access subtree. Intuitively, C_d suggests how “compatible” the data files are, C_e suggests how “compatible” the executable files are. We choose to use the average of the two values as the compatibility of the two trees in question.

For the finished working tree in Figure 2(b), E is the only data file appearing in both the working tree and the pattern tree, making C_d 0.4. Three out of four executable pairs (B, F and G) root identical access subtrees, so C_e is 0.75. The compatibility is thus 0.575.

The case of an unfinished working tree is similar except that one immediate child in the pattern tree is first determined to be the *pivot* node. The pivot corresponds to the most recently added child node in the working tree. Only the child nodes in the pattern tree that appear before the pivot are involved in compatibility computation. If the pattern tree is selected, we prefetch those files that follow the pivot in the sequence of pattern tree preordering, since those before the pivot probably have already been accessed. Given the unfinished working tree in Figure 2(c), node E is the pivot in the pattern tree. C_d and C_e are both 0.5, giving rise to a 0.5 compatibility. If we decide to prefetch this pattern tree, only the files in subtrees T_{F3} and T_{G1} will be prefetched.

Since the compatibility function is invoked often, it is important that it not be expensive. That is why we

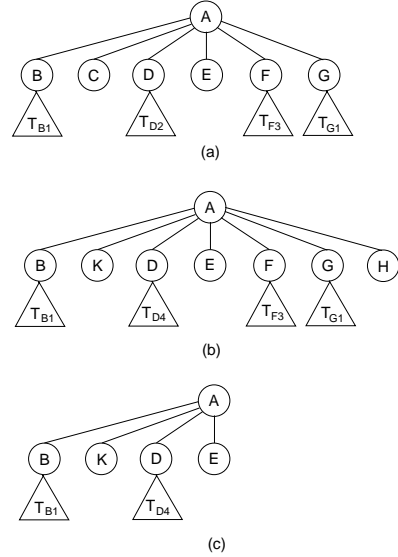


Figure 2: **Computation of Compatibility.** Graph (a) shows a pattern tree. Graph (b) shows a finished working tree, bearing a 0.575 compatibility with the given pattern tree. Graph (c) shows an unfinished working tree, so far bearing a 0.5 compatibility with the pattern tree; E is the pivot in the pattern tree, since it corresponds to the latest child node in the working tree being formed.

examine only immediate child nodes. The time complexity of the computation is proportional to the number of child nodes.

3 Simulation

Our initial assessments of the mechanism were obtained by trace-driven simulation.

3.1 Method

We gathered several file traces [23] on SunOS 4.0.3c. This version of SunOS offers a “C2 secure computing facility” that includes the ability to produce a system call audit trail. Using this feature, we gathered three traces of a volunteer user performing his normal work activity over a period of two weeks. The first trace contains 25,441 invocations of previously enumerated system calls captured over 72 hours. The second trace contains 23,858 invocations captured over 52 hours, while the numbers for the third trace are 85,962 and 86, respectively. During these hours activity varied widely and included compilations, document production, data analysis and display, large file searches, news reading, printing, and other operations.

Our simulated cache manager used an LRU replacement policy. It stepped through the traces, maintaining the cache in accordance with the user file accesses. Since our trace data lacks file sizes, we defined cache

Trace	Cache Size	No PF	PF	Ratio
1	30	34.13%	17.84%	1.91
	50	24.56%	13.58%	1.81
	100	21.37%	11.54%	1.85
2	30	42.10%	26.03%	1.62
	50	34.89%	22.90%	1.52
	100	27.32%	19.91%	1.37
3	30	12.77%	6.79%	1.88
	50	11.00%	5.88%	1.87
	100	8.81%	4.95%	1.78

Table 1: **Cache Miss Rate.** This table compares the accumulative cache miss rate with prefetching against that without prefetching. The ratio of two miss rates is given in the last column.

size by number of files. We varied the cache size and for each cache size, we compared the results of prefetching against those without prefetching.

3.2 Results

Our first metric was the cache miss rate. A summary of the results for the three traces appears in Table 1. In the table, the “No PF” column and “PF” column report the miss rates without prefetching and with prefetching, respectively. The “Ratio” column shows the ratio of miss rates. Prefetching delivers a substantially better miss rate with all cache sizes. The results are worse in every measure for the second trace because it includes a recursive directory traversal (i.e., the UNIX `find` program) over a hierarchy that includes thousands of files that are never accessed again.

In addition, we monitored the cache behavior at a finer grain: how well our mechanism worked over each run of 500 accesses. At the end of each run, we checked to see whether prefetching had beaten no-prefetching by measuring the number of misses during that run. As shown in Table 2, prefetching won this comparison easily and consistently. This shows that our mechanism’s superiority is steady and stable, and not simply a result of a few exceptionally fruitful prefetch sequences. The few losses due to bad prefetching guesses are more than offset by the many wins.

Both Tables 1 and 2 indicate that the increased intelligence of our mechanism is more effective in smaller caches. This was expected. Consider that, in the extreme case of an infinitely large cache, any file appearing in any access tree is still in the cache. There is no room for improvement from any prefetching algorithm that is based solely on information about the past.

We measured the accuracy of our prefetch decisions, defined as the percentage of file access predictions (*predicts*) that were actually used (*uses*). We also calculated one overhead of the mechanism, i.e., the percentage of the file fetches (*fetches*) that were initiated due to bad prefetch decisions (*wastes*). This reflects the network

Trace	Cache Size	Wins	Ties	Losses
1	30	31	6	0
	50	29	6	2
	100	24	11	2
2	30	31	4	0
	50	30	4	1
	100	23	9	3
3	30	92	47	2
	50	86	53	2
	100	70	67	4

Table 2: **A Finer-grained Comparison.** At the end of each run of 50 file accesses, we determined whether prefetching had beaten no-prefetching by measuring the number of cache misses during that run.

Trace	Loads	Improvers	Percent
1	1314	82	6.24%
2	1426	198	13.88%
3	4201	181	4.31%

Table 4: **Effectiveness of the compatibility Function.** This table shows, for each trace, the number of times we prefetched pattern trees (*loads*), and the number of times the prefetched pattern trees turned out to be inaccurate (*improvers*). The last column lists the percentage of loads that are improper.

bandwidth that was wasted. Table 3 shows these results. (We shall address another overhead, the CPU cycles expended by the prefetcher, in section 4.) In accordance with the LRU policy, bigger caches give better accuracy results because prefetched entries have a better chance to survive cache entry replacements. Similarly, bigger caches also give better overhead results because bad predictions have a better chance to have already existed in the cache; in such a case, no fetch is performed.

The **compatibility** function plays a critical role in our mechanism. The assumption underlying this function is that if the initial portions of two access trees bear a high compatibility, so will the two access trees in their entirety. In order to test this assumption, we counted the number of times we prefetched pattern trees (*loads*) and the number of times the prefetched pattern trees turned out to be inaccurate (*improvers*). An improper load was detected when a different pattern tree was selected by the prefetcher to take the place of the current one, or when the finished working tree did not match the pattern tree loaded. Our results in Table 4 suggest that the compatibility function is effective.

Finally, we illustrate the stability of our mechanism with regard to the algorithmic parameters. Table 5 gives three measurements of the algorithm under five different

Trace	Cache Size	Predicts	Uses	Accuracy	Fetches	Wastes	Overhead
1	30	7258	6795	93.62%	7066	386	5.46%
	50	7208	6842	94.92%	4942	234	4.73%
	100	7190	6875	95.62%	4227	176	4.16%
2	30	7402	6320	85.38%	8631	924	10.71%
	50	7111	6410	90.14%	6963	498	7.15%
	100	7052	6526	92.54%	5238	289	5.52%
3	30	16537	15421	93.25%	10194	924	9.06%
	50	16409	15615	95.16%	8486	565	6.66%
	100	16365	15730	96.12%	6745	388	5.75%

Table 3: **Prefetching Accuracy and Overhead.** The *accuracy* is defined as the percentage of file access predictions (*predicts*) that were actually used (*uses*). The *overhead* is defined as the percentage of the total file fetches (*fetches*) that were initiated due to bad prefetch decisions (*wastes*).

parameter settings. The measurements are the ratio of miss rates (*Ratio*), prediction accuracy (*accuracy*) and prefetch overhead (*overhead*). They were taken for all three traces with a cache size of 50. The two values in the “Setting” column are for `MATCH_THRESHOLD` and `PREFETCH_CAPACITY` respectively. The simulation results we have shown so far correspond to the first setting.

The simulation demonstrates that our algorithm can accurately predict future file accesses based on past file usage. The major limitation of the simulation study is that it does not account for the relative timing of events due to the absence of such information in the traces used. As a result of this, prefetched files are assumed to appear in the cache instantaneously. It remains to be determined whether a real system will have the resources to exploit the information on future file accesses. The limitations of the simulation motivated us to conduct a full implementation and further evaluations.

4 Implementation

We have implemented our mechanism in UX42 [8], a BSD UNIX server running on Mach 3.0 [1]. UX42 resides in user space and is organized as a collection of C threads [6]. Most threads handle BSD system calls. Among the others are NFS *async daemons*, which handle asynchronous NFS block I/O requests. Since we expect that network file accesses would be the performance bottleneck in a client-server architecture, particularly when the network link is wireless, we prefetch only NFS files opened for read.

4.1 Structure

Figure 3 shows the basic structure of the implementation, where each box stands for a C thread and the shaded area constitutes the prefetcher. The prefetcher consists of two pieces of code. The first is the system independent prefetch engine, which handles prefetch analysis, working tree construction and pattern tree maintenance. The same code was used in the simulation. In the implementation, the BSD service threads were modified

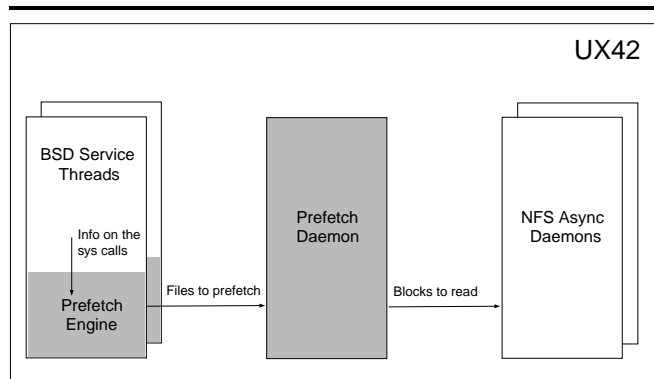


Figure 3: **Structure of Implementation.** A prefetch daemon is added to the collection of C threads in UX42. Each BSD service thread is extended to call the prefetch engine when a relevant system call is serviced.

to provide the prefetch engine information on each `fork`, `execve`, `open`, `chdir` and `exit` system call. The prefetch engine processes this information, makes prefetch decisions and enters the files to be prefetched in a queue. The second piece of prefetcher code is an added thread called the prefetch daemon. The prefetch daemon consumes the file queue and produces block read requests in another queue. These requests are then satisfied by the async daemons. Unlike a user-initiated read operation, a prefetch ends when the block is placed in the system buffer cache. No copy to user space is necessary.

The prefetch daemon takes advantage of the NFS read-ahead logic by prefetching only the first block of a file. When a requested block number is one more than the number of the block last read (`r_lastr`), NFS read-ahead logic speculates that the file is being accessed sequentially, and initiates an asynchronous read on the next block together with the requested block. At least two requested blocks are needed to establish the sequen-

Setting	Ratio			Accuracy			Overhead		
	Tr 1	Tr 2	Tr 3	Tr 1	Tr 2	Tr 3	Tr 1	Tr 2	Tr 3
0.4 / 15	1.81	1.52	1.87	94.92%	90.14%	95.16%	4.73%	7.15%	6.66%
0.3 / 10	1.80	1.50	1.88	95.46%	91.17%	94.55%	4.30%	6.79%	7.65%
0.5 / 10	1.79	1.49	1.76	95.42%	91.30%	94.01%	4.16%	6.70%	8.67%
0.5 / 20	1.80	1.51	1.77	94.73%	90.17%	92.04%	4.92%	7.65%	11.36%
0.3 / 20	1.80	1.52	1.89	94.89%	90.03%	92.94%	4.91%	7.72%	10.00%

Table 5: **Mechanism Stability.** This table presents selected measurements of the prefetching mechanism under five different settings of the algorithmic parameters (`MATCH_THRESHOLD` and `PREFETCH_CAPACITY`). The measurements were taken for all three traces with a cache size of 50.

tial access pattern before NFS starts read-ahead. Accordingly, when it removes an entry from the file queue, the prefetch daemon generates a read request for only the first file block (block 0). It also sets `r_lastr` to -1. Later when block 0 is actually accessed, the request will hit in the cache, and, since `r_lastr` is -1, a read-ahead on block 1 will be issued. This scheme moves on as the file is read block by block, which is the norm. Since the prefetch daemon issues only one block read for each file to be prefetched, the cost of prefetching is minimal.

We have ensured that prefetch I/O yields to regular NFS asynchronous I/O. All the NFS async daemons can be used towards regular I/O, but only up to a certain number of them towards prefetch I/O. Prefetch I/O will be started only if this limit has not been reached and there is no pending regular I/O. Thus, regular asynchronous I/O’s can be serviced promptly and when there are too many of these, the prefetcher will refrain from issuing any prefetch I/O’s. This ensures that prefetching halts when the system capacity limit is reached, and does not add extra load to an overload.

The implementation consists of approximately 3380 lines of C code. Of these, 90 lines have been added to the existent UX42 source files, 2280 lines are in separate “.c” files and 1010 lines in “.h” files.

4.2 Controlled Experiments

We started the evaluation of the implementation with two simple, controlled experiments. The following questions motivated our experiments:

- How much are the potential benefits of prefetching?
- Under antagonistic circumstances, what is the bearing of prefetching on performance?
- What is the CPU overhead due to prefetching?

The first experiment is a shell script that consists of tens of filter programs, each of which reads one parametric input file, performs some transformation, and writes one output file. Many well known UNIX programs are filters; included in our script are: `awk`, `compress`, `sed`, `sort`, `strings`, `uniq`, `uuencode`, and members of the `grep` family. All input files are the same size and reside remotely.

The second experiment is composed of several program builds. Although builds may seem ideally suited to prefetching, they are antagonistic, for two reasons. First, on our hardware platform a build is CPU intensive, leaving the prefetch mechanism relatively little excess CPU cycles with which to make prefetch decisions. Second, compilations often access header files rapid-fire, meaning that the time between a prefetch decision and the actual need for the data may not be sufficient to complete the prefetch I/O.

The experiments were run standalone. Both the client and the server are 486 processors with 16MB of memory. The client dedicates 1.57MB to its UNIX buffer cache. Since we are interested to find out how our mechanism behaves in relation to different network bandwidth, we ran the tests using both hardwired and wireless links directly connecting the client and server. The hardwired connection is an Ethernet (10 Mb/sec), while the wireless link is an NCR “WaveLAN” [25] radio link with a maximum 2 Mb/sec data rate. For each combination of experiment and network link, we ran the tests both with and without prefetching. Each number reported below is the mean of three trials.

Table 6(a) shows the results when the filters experiment is run with a wired link. We varied the work load by using different input file sizes, as given in the “File Size” column. Size is measured in multiples of the server’s preferred NFS block size, which is 4KB. For each workload, we list the UNIX buffer cache miss rate as well as the directory name lookup cache miss rate. Since our mechanism deals exclusively with network files, the cache miss rates are those of remote entries. The reduction of the miss rates due to prefetching is substantial. We also present the measurements of application latency, or total elapsed time, which we couldn’t do in the simulation. From a user’s standpoint, latency is the most important performance metric. The time measurements are in seconds, with standard deviations included in parentheses. The “Speedup” column is the ratio of no-prefetching latency to prefetching latency. The speedups are significant, particularly when the input files are relatively small, since the delay caused by a sequential file access mainly lies in the first one or two block reads.

File Size (blocks)	Cache Miss Rate				Latency		
	Buffer Cache		Name Cache		No PF	PF	Speedup
	No PF	PF	No PF	PF			
1	41.67%	0.00%	68.47%	5.87%	20.34 (0.29)	12.21 (0.52)	1.67
2	52.63%	0.00%	68.47%	4.26%	25.92 (0.17)	17.43 (0.53)	1.49
4	30.30%	0.00%	68.47%	4.64%	34.24 (1.18)	24.69 (0.56)	1.39
8	16.39%	0.10%	68.47%	4.45%	46.39 (0.27)	37.47 (0.82)	1.24

(a)

File Size (blocks)	Cache Miss Rate				Latency		
	Buffer Cache		Name Cache		No PF	PF	Speedup
	No PF	PF	No PF	PF			
1	41.67%	1.04%	68.47%	8.14%	23.94 (0.54)	14.33 (0.89)	1.67
2	52.63%	0.88%	68.47%	6.72%	30.07 (1.24)	20.34 (0.66)	1.48
4	30.30%	1.07%	68.47%	6.63%	38.76 (1.10)	28.38 (1.24)	1.37
8	16.39%	0.07%	68.47%	4.92%	54.04 (1.28)	42.94 (1.10)	1.26

(b)

Table 6: **Filters Experiment.** This table summarizes the performance results of the filters experiment. Part (a) presents the results with the wired link; part (b) with the wireless link.

The remaining blocks, if any, will be brought into cache by NFS read-ahead logic before they are needed.

Table 6(b) shows the results when the same experiment is run over a wireless link. The application latency is larger because of the slower link, but the speedups are comparable to those in a wired setup. The cache miss rates with prefetching are slightly higher over a wireless link than over a wired link since a small amount of prefetch operations cannot be completed in time for demanded use. Nevertheless, the bandwidth of the wireless link is adequate to perform NFS read-ahead on a timely basis, so it still suffices for the prefetcher to read only the first file block for each prediction it makes.

Figure 4 graphically illustrates the application latency of the filters experiment.

Our second experiment consists of builds of several UNIX utility programs. Table 7 contains the results of these tests. The application latency is also illustrated in Figure 5. When a large number of header files are opened in quick succession — common during compilation — the prefetcher often does not have enough time or available CPU cycles to make a fruitful prefetch even though it can speculate accurately which files will soon be referenced. The relatively small name cache miss rate improvements confirm this. However, prefetching still manages to reduce the buffer cache miss rate by 65% to 85%. There is no significant latency enhancement, suggesting that the builds are CPU-bound. We are pleased to see that no negative effects are observed in an adverse case such as this.

We have used simulation traces to demonstrate the prefetch overhead in terms of wasted network bandwidth. The implementation enables us to collect data on another type of cost, extra CPU consumption. Table

8 presents the results for the two experiments. The CPU time includes the system time and user time. Again the time measurements are in seconds and the standard deviations are given in parentheses. The CPU time overhead is negligible in all cases. We also show the ratio of CPU time to application latency ($CPU/latency$). It comes as no surprise that prefetching increases this ratio, substantially in the first experiment. A prefetcher reduces the total elapsed time by increasing the parallelism between CPU processing and I/O.

4.3 Discussion

There are three necessary conditions for prefetching to be useful. First, there must be spare capacity in the whole “data pipe” that extends from server’s disk to client’s cache. This pipe consists of the client’s network I/O interface, the network, the server’s CPU, and the server’s disk and network I/O paths. (Up to this point, we have been using the term “network bandwidth” loosely to refer to the data pipe.) Second, there must be system resources on the client side for the prefetcher’s use. Of special importance are CPU cycles. The prefetcher cannot fulfill its duty if it doesn’t acquire CPU cycles promptly, even though the amount of CPU time needed is low. Third, the workload must allow some interval between file accesses so that the prefetched I/O can be started and completed ahead of the demand.

Prefetching proves feasible on both wired and wireless network connections. Although our experiment results suggest that the speedups for these two situations are comparable, the speedup may be even better appreciated by the client at the end of a wireless link. With the slower link, a job lasts longer and more real time can be saved via prefetching. As the ratio of CPU speed

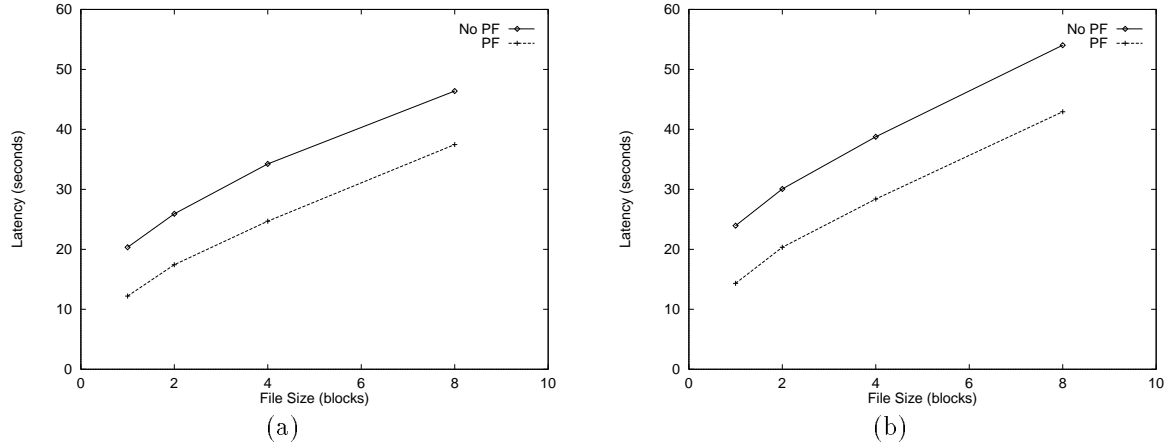


Figure 4: **Filters Experiment: Comparison of Application Latency.** These graphs illustrate the latency data in Table 6. Part (a) compares latency in the wired setting; part (b) in the wireless setting. Prefetching results in significant speedup in this experiment.

to network speed increases, prefetching should provide more benefit. However, there is some subtlety with our implementation scheme. When the ratio is lowered to a certain point, we will no longer be able to depend on the NFS read-ahead function to bring in subsequent file blocks quickly enough. This difficulty can be tackled by a simple generalization of our initial approach. The prefetcher can always prefetch the first N file blocks (N growing with the CPU/network ratio) and the NFS read-ahead code can be modified to read the next N -th file block, instead of the immediately next block. Currently, N is simply set to 1.

5 Related Work

Prefetching is an old idea. It has been studied extensively in various areas, including prepagging, prefetching of files and prefetching of database objects. Prepagging has not had a major impact in computer architecture because of the tight time and complexity constraints on paging hardware and software. However, prefetching of files and database objects is a more promising endeavor for two reasons. First, since file and database accesses are less frequent than page accesses, the speed with which the decision to prefetch must be made is not so much of the essence. Secondly, the resource most needed to arrange intelligent prefetching, namely client CPU cycles, is the resource most in excess in distributed systems now and in the likely future.

Some researchers have looked into prefetching blocks *within* files. Korner describes a method for detecting and exploiting block access patterns of individual files [13]. His approach depends on off-line processing by expert systems to discover patterns. The work of Kotz and Ellis (see, for example, [14]) focuses on the uses of prefetching to increase I/O bandwidth in MIMD shared-memory multiprocessors. They exploit the fact that many scientific and database applications running on multiproces-

sors exhibit simple and well-known patterns of sequential access coupled with read-only or write-only behavior; e.g., read every N th block of the entire file. Their prefetching methods are geared to these access patterns. In contrast to these methods, our work looks for access patterns *across* files.

Some file prefetching methods require that each application inform the operating system of its future demands. This includes the TIP project by Patterson *et al.* [20, 21] and the work of Cao, Felton, Karlin and Li [4, 5]. These researchers' outlook is more broad than ours. They consider prefetching from devices as well as from file servers. They also consider the interaction between prefetching and caching. TIP uses application-disclosed access patterns to dynamically allocate file buffers between the competing demands of prefetching and caching, based on a cost-benefit model. Cao *et al.* allow applications to pass down both prefetching hints and caching hints. They then employ an integrated algorithm for prefetching and caching, which is shown to be theoretically near-optimal. These "informed" approaches possess an advantage over ours in that prefetching is driven not by deductions made after snooping, but rather by certain knowledge provided in advance by higher levels. There is no danger that disastrously incorrect speculative prefetching might trash the cache. On the other hand, these approaches require recoding applications. Also, the prefetching mechanism must act within the interval between when the higher level learns of the need to do I/O and when it actually initiates I/O; this interval may not always be sufficient to perform the prefetch I/O.

Like ours, a number of other prefetching methods are completely transparent to clients. They use past accesses to predict future accesses. While we seek to build semantic structures, *i.e.*, access trees, that are endowed with application-level meaning, most of the other approaches use a probabilistic method to model the user

Utility	Cache Miss Rate				Latency		
	Buffer Cache		Name Cache		Latency		Speedup
	No PF	PF	No PF	PF	No PF	PF	
hash_info	32.09%	5.30%	8.39%	6.92%	21.59 (0.32)	19.53 (0.53)	1.11
snames	17.15%	2.85%	9.03%	7.84%	30.93 (0.44)	29.81 (0.49)	1.04
machid	10.57%	3.17%	10.35%	9.83%	64.11 (1.05)	63.63 (0.68)	1.01
machipc	23.42%	3.36%	8.48%	6.90%	25.48 (0.78)	24.70 (2.51)	1.03

(a)

Utility	Cache Miss Rate				Latency		
	Buffer Cache		Name Cache		Latency		Speedup
	No PF	PF	No PF	PF	No PF	PF	
hash_info	32.09%	5.61%	8.46%	6.99%	26.22 (0.66)	24.79 (0.47)	1.06
snames	17.20%	2.96%	9.09%	7.91%	38.17 (1.17)	33.89 (1.64)	1.13
machid	9.15%	3.21%	10.39%	9.86%	77.31 (3.25)	74.42 (2.13)	1.04
machipc	23.27%	4.03%	8.55%	7.02%	29.10 (0.36)	27.17 (0.30)	1.07

(b)

Table 7: **Builds Experiment.** This table summarizes the performance results of the builds experiment. Part (a) presents the results with the wired link; part (b) with the wireless link.

behavior. They examine strings of file accesses for patterns of the form “when file A is accessed, there is a 40% chance that file B will be accessed soon.” Unlike our tree-based approach, there is no attempt to build in an “understanding” of why files A and B are likely to be accessed together.

Some recent examples of work based on probabilistic methods are a paper by Curewitz *et al.* [7] on prefetching objects in an object-oriented database, work by Griffioen and Appleton [9, 10] and by Kroeger and Long [15] that prefetch whole files. Both [7] and [15] adapt context modeling techniques used in data compression to predict the next access. Their work is inspired by the theory that prediction is synonymous with data compression. Intuitively, in order to compress data well, one has to be able to predict future data well, and hence a good data compressor should be able to predict well for the purpose of prefetching. Griffioen and Appleton’s work, in comparison, employs a “probability graph” that for each file accumulates frequency counts of all files that are accessed shortly after the first file.

In the initial stages of our work, we considered probabilistic modeling. We implemented (in the simulator only) two extremely simple probabilistic methods, which we called “stupid pairs” and “smart pairs:”

- **Stupid pairs:** When file F is accessed, prefetch the file that was accessed immediately following F at the last time F was accessed.

Sample series of accesses: F, G (remember F-G), F (prefetch G, remember G-F), H (remember F-H), F (prefetch H, remember H-F), H (cache hit).

- **Smart pairs:** Keep track of all files that are accessed immediately after F, and when F is accessed

the next time, choose one according to the frequency distribution.

Sample series of accesses (only F’s pairs are shown): F, G (F-G1), F (prefetch G), H (F-[G1, H1]), F (prefetch G or H with 1:1 weighting), H (F-[G1, H2]), F (prefetch G or H with 1:2 weighting).

The stupid pair scheme worked better than the smart pair approach when applied to our traces. This unexpected result is explained by considering the locality of many accesses: often, a user works on one file or one group of files for some time, then moves on to similar operations with different files. Stupid pairs are well-equipped to handle this usage pattern, since they invariably prefetch the most recently used successor file. The seemingly superior intelligence of smart pairs actually becomes a liability when locality is strong; files no longer in active use may be given undue weight if they were heavily accessed in the past.

The phenomenon, that less information is better provided that it is more recent, identifies one common weakness of probabilistic methods: how to determine an appropriate window in the recent past from which to infer future accesses? None of the above methods addresses this issue. Instead they make predictions based on a global history of accesses. An additional weakness of methods like [7] and [15] is that they have limited lookahead: they predict the *single* next event. They are less likely to be widely applicable than methods with longer lookahead.

Similar to our work, Palmer and Zdonik’s work on Fido [19] also explicitly recognizes and maintains access patterns. Several important aspects make their work different. First, their work was conducted in the context of object-oriented database systems, whereas our context

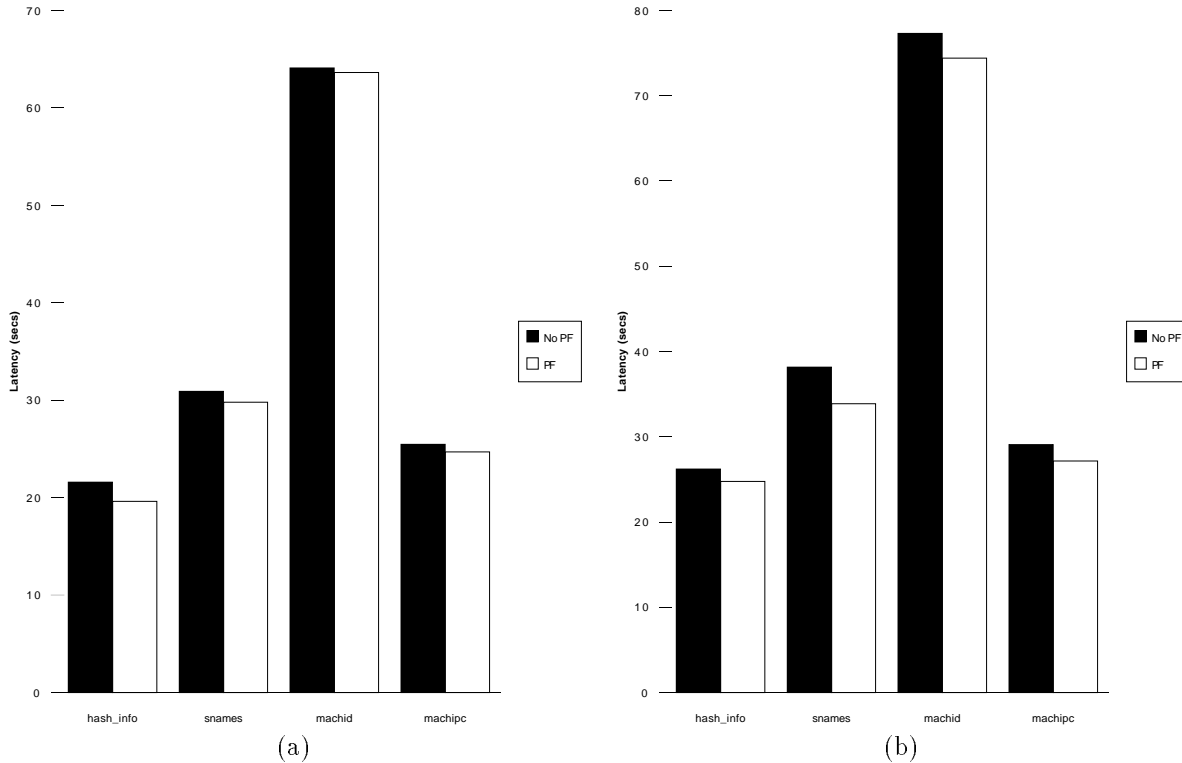


Figure 5: **Builds Experiment: Comparison of Application Latency.** These graphs illustrate the latency data in Table 7. Part (a) compares latency in the wired setting; part (b) in the wireless setting. There are no observable negative effects in this antagonistic experiment.

is file systems. Second, they represent access patterns with strings of object identifiers, with no semantics involved. We represent patterns with access trees. Third, they employ specialized pattern memory, while we store pattern trees in virtual memory. Finally and most notably, Fido requires a separate training phase after each user session, while our mechanism is more on-line: there is no off-line computation or periodic analysis needed.

An issue that is related to prefetching is hoarding [22, 11, 16, 24]. Both prefetching and hoarding involve anticipatory file fetches: bringing files from remote servers into a local cache before they are needed. These are not exactly the same techniques, however. Hoarding is a scheme designed to increase the likelihood that a mobile client will be able to continue working during periods of total disconnection from file servers. Since hoarding is a relatively infrequent operation performed only at a client’s request prior to disconnection, timing is not critical. On the other hand, prefetching is mainly concerned with improving performance and timing is important. With prefetching, the file server is assumed to be still accessible, although the network connectivity may be weak. A cache miss is much more catastrophic in disconnected operations, hence hoarding is typically willing to overfetch substantially in order to enhance the availability of files. Despite the differences, our idea of uncovering and exploiting the semantic structure under-

lying file accesses also applies to the hoarding problem, as shown in [24].

6 Conclusion

We have presented a technique for transparent on-line file prefetching. The technique analytically models interesting system calls and builds semantic structures that capture the interrelationships between file accesses. It makes accurate predictions of future file accesses, imposes little CPU overhead, doesn’t interfere with demand I/O, and delivers substantially lower client cache miss rates and elapsed time for I/O-intensive applications.

One central trait of the algorithm is that it spends client CPU cycles in return for more effective use of client cache space and fewer on-demand network operations. Another distinguishing aspect is that the algorithm’s lookahead ability is potentially much greater than that of previous work. Both of these traits help to couple application I/O performance more closely to CPU speed than to I/O device speed, thereby addressing a fundamental and longstanding problem in operating systems [18].

Our initial performance evaluation has been encouraging. We intend to conduct more sophisticated experiments so that we can study the prefetcher behavior

File Size	Wired Link				Wireless Link			
	CPU Time		CPU/latency		CPU Time		CPU/latency	
	No PF	PF	No PF	PF	No PF	PF	No PF	PF
1	3.82 (0.24)	4.50 (0.24)	18.80%	36.87%	3.69 (0.07)	4.37 (0.07)	15.41%	30.47%
2	5.39 (0.71)	5.80 (0.48)	20.81%	33.29%	4.77 (0.09)	5.83 (0.14)	15.86%	28.65%
4	7.88 (0.49)	8.78 (0.79)	23.00%	35.58%	7.39 (0.22)	8.60 (0.35)	19.07%	30.29%
8	13.24 (1.03)	14.45 (0.53)	28.54%	38.58%	13.14 (0.08)	14.58 (0.08)	24.32%	33.95%

(a)

Utility	Wired Link				Wireless Link			
	CPU Time		CPU/latency		CPU Time		CPU/latency	
	No PF	PF	No PF	PF	No PF	PF	No PF	PF
hash_info	5.19 (0.21)	5.84 (0.37)	24.03%	29.91%	5.17 (0.11)	5.79 (0.21)	19.72%	23.36%
snames	10.36 (0.67)	11.25 (0.80)	33.48%	37.73%	10.69 (0.10)	11.16 (0.40)	28.00%	32.94%
machid	30.29 (0.92)	32.08 (1.28)	47.24%	50.41%	30.96 (1.82)	32.58 (1.70)	40.04%	43.78%
machipc	6.74 (0.22)	7.63 (0.45)	26.45%	30.90%	6.91 (0.27)	7.80 (0.34)	23.73%	28.72%

(b)

Table 8: **CPU Time Consumption.** This table presents the CPU time with and without prefetching, over two different network links. Also given is the ratio of CPU time to application latency. Part (a) give the results for the filters experiment; part (b) for the builds experiment.

(gains and overheads) over a much wider spectrum of conditions (network capacity, client CPU capacity, work load etc.). In addition, we intend to explore applying the access tree notion to caching and investigate effective combination of three related techniques: prefetching, hoarding and caching. We speculate that an integrated file cache management strategy taking advantage of the semantics revealed in past file usage would work well.

7 Acknowledgements

This work was supported in part by the Advanced Research Projects Agency, ARPA order number B094, under contract N00014-94-1-0719, monitored by the Office of Naval Research; and in part by the Center for Telecommunications Research, an NSF Engineering Research Center supported by grant number ECD-88-11111.

The authors wish to thank Carl Tait for many constructive conversations and insightful comments.

References

- [1] M. Accetta, R. Baron, W. Bolosky, D. Golub, R. Rashid, A. Tevanian, and M. Young. Mach: A New Kernel Foundation for UNIX Development. In *Proc. 1986 USENIX Summer Conf.*, pages 93–112, June 1986.
- [2] R. Alonso, D. Barbara, and L. L. Cova. FACE: Enhancing Distributed File Systems for Autonomous Computing Environments. Technical Report CS-TR-214-89, Princeton University, March 1989.
- [3] R. Alonso, D. Barbara, and L. L. Cova. Augmenting Availability in Distributed File Systems. Technical Report CS-TR-234-89, Princeton University, October 1989.
- [4] P. Cao and E. W. Felten and A. Karlin and K. Li. A Study of Integrated Prefetching and Caching Strategies. In *Proc. 1995 ACM SIGMETRICS*, pages 171–182, June 1995.
- [5] P. Cao and E. W. Felten and A. Karlin and K. Li. Implementation and Performance of Integrated Application-Controlled Caching, Prefetching and Disk Scheduling. In *Proc. First USENIX Symposium on Operating Systems Design and Implementation*, pages 165–178, November 1994.
- [6] E. C. Cooper and R. P. Draves. C Threads. Technical Report CMU-CS-88-154, Carnegie Mellon University, June 1988.
- [7] K. M. Curewitz, P. Krishnan, and J. S. Vitter. Practical Prefetching via Data Compression. In *Proc. 1993 ACM SIGMOD*, pages 257–266, May 1993.
- [8] D. Golub, R. Dean, A. Forin, and R. Rashid. Unix as an Application Program. In *Proc. Summer 1990 USENIX Conf.*, USENIX, pages 87–95, June 1990.
- [9] J. Griffioen and R. Appleton. Reducing File System Latency Using a Predictive Approach. In *Proc. 1994 USENIX Summer Conf.*, pages 197–207, June 1994.
- [10] J. Griffioen and R. Appleton. Performance Measurements of Automatic Prefetching. In *Parallel and Distributed Computing Systems*, pages 165–170, IEEE, September 1995.

- [11] L. B. Huston and P. Honeyman. Disconnected Operation for AFS. In *Proc. First USENIX Symp. on Mobile and Location-Independent Computing*, pages 1–10, August 1993.
- [12] J. J. Kistler and M. Satyanarayanan. Disconnected Operation in the Coda File System. In *Proc. Thirteenth Symp. on Operating System Principles*, pages 213–225. ACM, October 1991.
- [13] K. Korner. Intelligent Caching for Remote File Service. In *Proc. Tenth Intl. Conf. on Distributed Computing Systems*, pages 220–226. IEEE, May 1990.
- [14] D. Kotz and C. S. Ellis. Practical Prefetching Techniques for Parallel File Systems. In *Proc. First Intl. Conf. on Parallel and Distributed Information Systems*, pages 182–189. ACM, December 1991.
- [15] T. M. Kroeger and D. D. E. Long. Predicting Future File-System Actions from Prior Events. In *Proc. 1996 USENIX Annual Technical Conf.*, pages 319–328, January 1996.
- [16] G. H. Kuenning. The Design of the Seer Predictive Caching System. in *Proc. Workshop on Mobile Computing Systems and Applications*, ACM/IEEE, pages 37–43, December 1994.
- [17] G. H. Kuenning, G. J. Popek, and P. L. Reiher. An Analysis of Trace Data for Predictive File Caching in Mobile Computing. In *Proc. 1994 USENIX Summer Conf.*, pages 291–303, June 1994.
- [18] J. K. Ousterhout. Why Aren't Operating Systems Getting Faster As Fast As Hardware? In *Proc. 1990 USENIX Summer Conf.*, pages 247–256, June 1990.
- [19] M. Palmer and S. Zdonik. Fido: A Cache That Learns to Fetch. In *Proc. 17th Intl. Conf. on Very Large Data Bases*, pages 255–264, September 1991.
- [20] R. H. Patterson, G. A. Gibson, and M. Satyanarayanan. A Status Report on Research in Transparent Informed Prefetching. *Operating Systems Review*, 27(2):21–34, April 1993.
- [21] R. H. Patterson, G. A. Gibson, E. Ginting, D. Stodolsky and J. Zelenka. Informed Prefetching and Caching. In *Proc. Fifteenth Symp. on Operating System Principles*, pages 79–95. ACM, December 1995.
- [22] M. Satyanarayanan, J. J. Kistler, P. Kumar, M. E. Okasaki, E. H. Siegel, and D. C. Steere. Coda: A Highly Available File System for a Distributed Workstation Environment. *IEEE Trans. on Computers*, 39(4):447–459, April 1990.
- [23] C. D. Tait and D. Duchamp. Detection and Exploitation of File Working Sets. In *Proc. Eleventh Intl. Conf. on Distributed Computing Systems*, pages 2–9. IEEE, May 1991.
- [24] C. D. Tait, H. Lei, S. Acharya and H. Chang. Intelligent File Hoarding for Mobile Computers. In *Proc. First Intl. Conf. on Mobile Computing and Networking*, pages 119–125, ACM, November 1995.
- [25] B. Tuch. An Engineer's Story of WaveLAN. In *Proc. First Virginia Tech Symp. on Wireless Personal Communications*, June 1991.