

Abstract

Today's Internet clients vary widely with respect to both hardware and software properties: screen size, color depth, effective bandwidth, processing power, and the ability to handle different data formats. The order-of-magnitude span of this variation is too large to hide at the network level, making application-level techniques necessary. We show that on-the-fly adaptation by transformational proxies is a widely applicable, cost-effective, and flexible technique for addressing all these types of variations. To support this claim, we describe our experience with datatype-specific distillation (lossy compression) in a variety of applications. We also argue that placing adaptation machinery in the network infrastructure, rather than inserting it into end servers, enables incremental deployment and amortization of operating costs. To this end, we describe a programming model for large-scale interactive Internet services and a scalable cluster-based framework that has been in production use at UC Berkeley since April 1997. We present a detailed examination of TranSend, a scalable transformational Web proxy deployed on our cluster framework, and give descriptions of several handheld-device applications that demonstrate the wide applicability of the proxy-adaptation philosophy.

Adapting to Network and Client Variation Using Infrastructural Proxies: Lessons and Perspectives

ARMANDO FOX, STEVEN D. GRIBBLE, YATIN CHAWATHE,
AND ERIC A. BREWER
UNIVERSITY OF CALIFORNIA, BERKELEY

The current Internet infrastructure includes an extensive range and number of clients and servers. Clients vary along many axes, including screen size, color depth, effective bandwidth, processing power, and ability to handle specific data encodings (e.g., GIF, PostScript, or MPEG). As shown in Tables 1 and 2, each type of variation often spans orders of magnitude. High-volume devices such as smart phones [1] and smart two-way pagers will soon constitute an increasing fraction of Internet clients, making the variation even more pronounced.

These conditions make it difficult for servers to provide a level of service that is appropriate for every client. Application-level adaptation is required to provide a *meaningful* Internet experience across the range of client capabilities. Despite continuing improvements in client computing power and connectivity, we expect the high end to advance roughly in parallel with the low end, effectively maintaining a gap between the two; hence, the need for application-level adaptation.

The Approach: Infrastructural Proxy Services

We argue for a *proxy-based approach* to adaptation, in which proxy agents placed between clients and servers perform aggressive computation and storage on behalf of clients. The proxy approach stands in contrast to the *client-based approach*, which attempts to bring all clients up to a least-common-denominator level of functionality (e.g., text-only, HTML-subset compatibility for thin-client Web browsers), and the *server-based approach*, which attempts to insert adaptation machinery at each end server. We believe the proxy approach directly confers three advantages over the client and server approaches:

- *Leveraging the installed infrastructure through incremental deployment.* The enormous installed infrastructure, and its attendant base of existing content, is too valuable to waste; yet some clients cannot handle certain data types effective-

ly. A compelling solution to the problem of client and network heterogeneity should allow interoperability with existing servers, thus enabling incremental deployment while evolving content formats and protocols are tuned and standardized for different target platforms. A proxy-based approach lends itself naturally to transparent incremental deployment, since an application-level proxy appears as a server to existing clients and as a client to existing servers.

- *Rapid prototyping during turbulent standardization cycles.* Software development on "Internet time" does not allow

Platform	SPEC92/ memory	Screen size	Bits/pixel
High-end PC	200/64 M	1280 x 1024	24
Midrange PC	160/32 M	1024 x 768	16
Typical laptop	110/16 M	800 x 600	8
Typical PDA	Low/2 M	320 x 200	2

■ Table 1. Physical variation among clients.

Network	Bandwidth (b/s)	Round-trip time
Local Ethernet	10–100 M	0.5–2.0 ms
ISDN	128 K	10–20 ms
Wireline modem	14.4–56 K	350 ms
Cellular/CDPD	9.6–19.2 K	0.1–0.5 s

■ Table 2. Typical network variation.

for long deployment cycles. Proxy-based adaptation provides a smooth path for rapid prototyping of new services, formats, and protocols, which can be deployed to servers (or clients) later if the prototypes succeed.

- *Economy of scale.* Basic queuing theory shows that a large central (virtual) server is more efficient in both cost and utilization (though less predictable in per-request performance) than a collection of smaller servers; standalone desktop systems represent the degenerate case of one “server” per user. This supports the argument for network computers [2] and suggests that collocating proxy services with infrastructural elements such as Internet points of presence (POPs) is one way to achieve effective economies of scale.

Large-scale network services remain difficult to deploy because of three fundamental challenges: scalability, availability, and cost effectiveness. By scalability, we mean that when the load offered to the service increases, an incremental and linear increase in hardware can maintain the same per-user level of service. By availability, we mean that the service as a whole must be available 24/7, despite transient partial hardware or software failures. By cost effectiveness, we mean that the service must be economical to administer and expand, even though it potentially comprises many workstation nodes operating as a centralized cluster or “server farm.” In the third section we describe how we have addressed these challenges in our cluster-based proxy application server architecture.

Contributions and a Map of the Article

In the following section we describe our measurements and experience with datatype-specific distillation and refinement, a mechanism that has been central to our proxy-based approach to network and client adaptation. We then introduce a generalized “building block” programming model for designing and implementing adaptive applications, describe our implemented cluster-based application server that instantiates the model, and present detailed measurements of a particular production application: TranSend, a transformational Web proxy service. We present case studies of other services we have built using our programming model, some of which are in daily use by thousands of users, including the Top Gun Wingman graphical Web browser for the 3Com PalmPilot handheld device. We discuss related work, and then attempt to draw some lessons from our experience and guidelines for future research in the last section.

Adaptation via Datatype-Specific Distillation

We propose three design principles that we believe are fundamental for addressing client variation most effectively.

Adapt to Client Variation via Datatype-Specific Lossy Compression – Datatype-specific lossy compression mechanisms can achieve much better compression than “generic” compressors, because they can make intelligent decisions about what information to throw away based on the semantic type of the data. For example, lossy compression of an image requires discarding color information, high-frequency components, or pixel resolution. Lossy compression of video can additionally include frame rate reduction. Less obviously, lossy compression of formatted text requires discarding some formatting information but preserving the actual prose. In all cases, the goal is to preserve the information that has the highest semantic value. We refer to this process generically as *distilla-*

Semantic type	Specific encodings	Distillation axes
Image	GIF, JPEG, PPM, Postscript	Resolution, color depth, color palette
Text	Plain, HTML, Postscript, PDF	Richness (heavily formatted vs. simple markup vs. plain text)
Video	NV, H.261, VQ, MPEG	Resolution, frame rate, color depth, progression limit (for progressive encodings)

■ **Table 3.** Three important types and the distillation axes corresponding to each.

tion. A distilled object allows the user to decide whether it is worth asking for a *refinement*: for instance, zooming in on a section of a graphic or video frame, or rendering a particular page containing PostScript text and figures without having to render the preceding pages.

Perform Adaptation on the Fly – To reap the maximum benefit from distillation and refinement, a distilled representation must target specific attributes of the client. The measurements that we later report show that for typical images and rich text, distillation time is small in practice, and end-to-end latency is reduced because of the much smaller number of bytes transmitted over low-bandwidth links. *On-demand distillation* provides an easy path for incorporating support for new clients, and also allows distillation aggressiveness to track, for example, such significant changes in network bandwidth as might occur in vertical handoffs between different wireless networks [3]. We have successfully implemented useful distillation “workers” that serve clients spanning an order of magnitude in each area of variation, and we have generalized our approach into a common framework, which we discuss later.

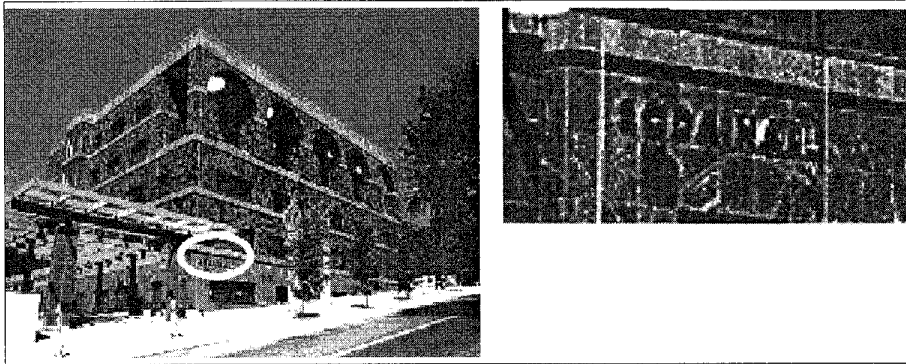
Move Complexity Away from Both Clients and Servers – Application partitioning arguments have long been used to keep clients simple [4]. However, adaptation through a shared infrastructural proxy enables incremental deployment and legacy client support, as we argued earlier. Therefore, on-demand distillation and refinement should be done at an intermediate proxy that has access to substantial computing resources and is well-connected to the rest of the Internet.

Table 3 lists the “axes” of compression corresponding to three important datatypes: formatted text, images, and video streams. We have found that order-of-magnitude size reductions are often possible without destroying the semantic content of an object (e.g., without rendering an image unrecognizable to the user).

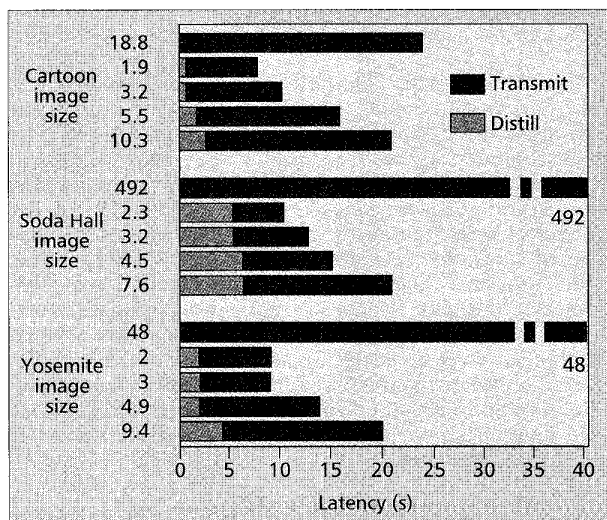
The Performance of Distillation and Refinement On Demand

We now describe and evaluate datatype-specific distillers for images and rich text.¹ The goal of this section is to support our claim that in the majority of cases, *end-to-end latency is reduced by distillation*, that is, the time to produce a useful distilled object on today’s workstation hardware is small enough to be more than compensated for by the savings in transmission time for the distilled object relative to the original.

¹ A distiller for real-time network video streams is described separately in [7].



■ **Figure 1.** Distillation example: left) a distilled image of Soda Hall; right) illustrates refinement. The left occupies 17 kbytes at 320 x 200 pixels in 16 grays, compared with the 492-kbyte 880 x 600 pixel, 249-color original (not shown). The refinement (right) occupies 12 kbytes. Distillation took 6 s on a SPARCstation 20/71, and refinement took less than 1 s.



■ **Figure 2.** End-to-end latency for images with and without distillation.

Images — We have implemented an image distiller called *gifmunch*, which implements distillation and refinement for GIF [5] images, and consists largely of source code from the NetPBM Toolkit [6]. Figure 1 shows the result of running *gifmunch* on a large color GIF image of the Berkeley Computer Science Division's home building, Soda Hall. The image of Fig. 1 (left) measures 320 x 200 pixels — about 1/8 the total area of the original 880 x 610 — and uses 16 grays, making it suitable for display on a typical handheld device.

Due to the degradation of quality, the writing on the building is unreadable, but the user can request a refinement of the subregion containing the writing, which can then be viewed at full resolution (Fig. 1, right).

Image distillation can be used to address all three areas of client variation:

- **Network variation:** The graphs in Fig. 2 depict end-to-end client latency for retrieving the original and each of four distilled versions of a selection of GIF images: the top set of bars is for a cartoon found on a popular Web page, the middle set corresponds to a large photographic image, and the bottom to a computer-rendered image. Each group of bars represents one image with five levels of distillation; the top bar represents no distillation at all. The y-axis number is the distilled size in kilobytes (so the top bar gives the

original size). Note that two of the undistilled images are off the scale; the Soda Hall image is off by an order of magnitude. The images were fetched using a 14.4 kb/s modem with standard compression (V.42bis and MNP-5) through the UC Berkeley PPP gateway, via a process that runs each image through *gifmunch*.² Each bar is segmented to show the distillation latency and transmission latency separately. Clearly, even though distillation adds latency at the proxy, it can result in

greatly reduced end-to-end latency. This shows that on-the-fly distillation is not prohibitively expensive.

- **Hardware variation:** A “map to 16 grays” operation would be appropriate for PDA-class clients with shallow grayscale displays. We can identify this operation as an effective lossy compression technique precisely because we know we are operating on an image, regardless of the particular encoding, and the compression achieved is significantly better than the 2–4 times compression typically achieved by “generic” lossless compression.
- **Software variation:** Handheld devices such as the 3Com PalmPilot frequently have built-in support for proprietary image encodings only. The ability to convert to this format saves code space and decoding latency on the client.

Rich Text — We have also implemented a rich-text distiller that performs lossy compression of PostScript-encoded text using a third-party PostScript-to-text converter [8]. The distiller replaces PostScript formatting information with HTML markup tags or a custom rich-text format that preserves the position information of the words. PostScript is an excellent target for a distiller because of its complexity and verbosity: both transmission over the network and rendering on the client are resource-intensive. Table 4 compares the features available in each format. Figure 3 shows the advantage of rich text over PostScript for screen viewing. As with image distillation, PostScript distillation yields advantages in all three categories of client variation:

- **Network variation:** Again, distillation reduces the required bandwidth and thus the end-to-end latency. We achieved an average size reduction of a factor of 5 when going from compressed PostScript to gzipped HTML. Second, the pages of a PostScript document are pipelined through the distiller, so the second page is distilled while the user views the first page. In practice, users only experience the latency of the first page, so the difference in perceived latency is about a factor of 8 for a 28.8 K modem. Distillation typically took about 5 s for the first page and about 2 s for subsequent pages.
- **Hardware variation:** Distillation reduces decoding time by delivering data in an easy-to-parse format, and results in better-looking documents on clients with lower-quality displays.

² The network and distillation latencies reflect significant overhead due to the naive implementation of *gifmunch* and the latency and slow-start effects of the PPP gateway, respectively. Later we discuss how to overcome some of these problems, but it is worth noting that end-to-end latency is still substantially reduced, even in this naive prototype implementation.

- *Software variation:* PostScript distillation allows clients that do not directly support PostScript, such as handhelds, to view these documents in HTML or our rich-text format. The rich-text viewer could be an external viewer similar to *ghostscript*, an applet for a Java-capable browser, or a browser plug-in rendering module.

Overall, rich-text distillation reduces end-to-end latency, results in more readable presentation, and adds new abilities to low-end clients, such as PostScript viewing. The latency for the appearance of the first page was reduced on an average by a factor of 8 using the proxy and PostScript distiller. Both HTML and our rich-text format are significantly easier to read on screen than rendered PostScript, although they sacrifice some layout and graphics accuracy compared to the original PostScript.

Summary

High client variability is an area of increasing concern that existing servers do not handle well. We have proposed three design principles we believe are fundamental to addressing variation:

- Datatype-specific distillation and refinement achieve better compression than does lossless compression, while retaining useful semantic content and allowing network resources to be managed at the application level.
- When the proxy-to-client bandwidth is substantially smaller than the proxy-to-server bandwidth (as is the case, e.g., in wireless networks or with consumer wire-line modems), on-demand distillation and refinement reduce end-to-end latency perceived by the client (sometimes by almost an order of magnitude), are more flexible than reliance on precomputed static representations, and give low-end clients new abilities such as PostScript viewing.
- Performing distillation and refinement in the network infrastructure rather than at the endpoints separates technical as well as economic concerns of clients and servers.

1.2 The Remote Queue Model

We introduce *Remote Queues (RQ)*, which provides a general abstraction for low-level systems of three basic elements. First, one or more receiving nodes. Second, an *enqueue* operation `enqueue(n, q, arg0, ..., argn, stuff)` that causes `arg0` through `argn`, followed

1.2 The Remote Queue Model

We introduce *Remote Queues (RQ)*, which provides a general abstraction for low-level systems of three basic elements. First, one or more receiving nodes. Second, an *enqueue* operation `enqueue(n, q, arg0, ..., argn, stuff)` that causes `arg0` through `argn`, followed

■ **Figure 3.** Screen snapshots of our rich text (top) versus ghostview (bottom). The rich text is easier to read because it uses screen fonts.

Feature	HTML	Rich text	PostScript
Different fonts	Y	Y	Y
Bold and italics	Y	Y	Y
Preserves font size	Headings	Y	Y
Preserves paragraphs	Y	Y	Y
Preserves layout	N	Y	Y
Handles equations	N	Some	Y
Preserves tables	Y	Y	Y
Preserves graphs	N	N	Y

■ **Table 4.** Features for PostScript distillation.

Scalable Internet Application Servers

In order to accommodate compute-intensive adaptation techniques by putting resources in the network infrastructure, we must address two important challenges:

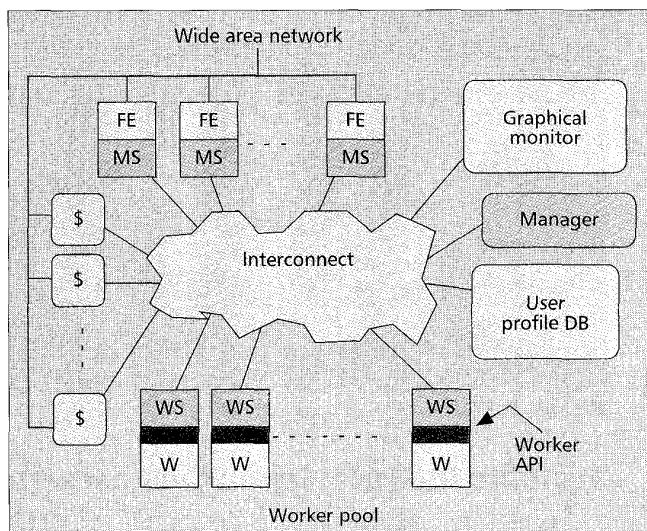
- Infrastructural resources are typically shared, and the sizes of user communities sharing resources such as Internet POPs is increasing exponentially. A shared infrastructural service must therefore scale gracefully to serve large numbers of users.
- Infrastructural resources such as the IP routing infrastructure are expected to be reliable, with availability approaching 24/7 operation. If we place application-level computing resources such as distillation engines into this infrastructure, we should be prepared to meet comparable expectations.

In this section, we focus on the problem of deploying adaptation-based proxy services to large communities (tens of thousands of users, representative of the subscription size of a medium-sized Internet service provider). In particular, we discuss a cluster-friendly programming model for building interactive and adaptive Internet services, and measurements of our implemented prototype of a scalable, cluster-based server that instantiates the model. Our framework reflects the implementation of three real services in use today: TranSend, a scalable transformation proxy for the 25,000 UC Berkeley dialup users (connecting through a bank of 600 modems); Top Gun Wingman, the only graphical Web browser for the 3Com PalmPilot handheld device (commercialized by ProxiNet); and the Inktomi search engine (commercialized as HotBot), which performs over 10 million queries per day against a database of over 100 million Web pages. Although HotBot does not demonstrate client adaptation, we use it to validate particular design decisions in the implementation of our server platform, since it pioneered many of the cluster-based scalability techniques generalized in our scalable server prototype. We focus our detailed discussion and measurements on TranSend, a transformational proxy service that performs on-the-fly lossy image compression. TranSend applies the ideas explored in the preceding section to the World Wide Web.

TACC: A Programming Model for Internet Services

We focus on a particular subset of Internet services, based on *transformation* (distillation, filtering, format conversion, etc.), *aggregation* (collecting and collating data from various sources, as search engines do), *caching* (both original and transformed content), and *customization* (maintenance of a per-user preferences database that allows workers to tailor their output to the user's needs or device characteristics).

We refer to this model as *TACC*, from the initials of the four elements above. In the TACC model, applications are



■ **Figure 4.** Architecture of a cluster-based TACC server; components include frontends (FEs), a pool of TACC workers (Ws), some of which may be caches (\$), a user profile database, a graphical monitor, and a fault-tolerant load manager, whose functionality logically extends into the manager stubs (MSs) and worker stubs (WSs).

built from building blocks interconnected with simple application programming interfaces (APIs). Each building block, or *worker*, specializes in a particular task, for example, scaling/dithering of images in a particular format, conversion between specific data formats, or extracting “landmark” information from specific Web pages. Complete applications are built by *composing* workers; roughly speaking, one worker can *chain* to another (similar to processes in a UNIX pipeline), or a worker can *call* another as a subroutine or coroutine. This model of composition results in a very general programming model that subsumes transformation proxies [9], proxy filters [10], customized information aggregators, and search engines.

A *TACC server* is a platform that instantiates TACC workers, provides dispatch rules for routing network data traffic to and from them, and provides support for the inter-worker calling and chaining APIs. Similar to a UNIX shell, a TACC server provides the mechanisms that insulate workers from having to deal directly with low-level concerns such as data routing and exception handling, and gives workers a clean set of APIs for communicating with each other, the caches, and the customization database (described below). We describe our prototype implementation of a scalable, commodity-PC cluster-based TACC server later in this section.

Cluster-Based TACC Server Architecture

We observe that clusters of workstations have some fundamental properties that can be exploited to meet the requirements of large-scale network services (scalability, high availability, and cost effectiveness). Using commodity PCs as the unit of scaling allows the service to ride the leading edge of the cost/performance curve; the inherent redundancy of clusters can be used to mask transient failures; and “embarrassingly parallel” network service workloads map well onto networks of commodity workstations.

However, developing cluster software and administering a running cluster remain complex. A primary contribution of our work is the design, analysis, and implementation of a layered framework for building adaptive network services that addresses this complexity while realizing the sought-after economies of scale. New services can use this framework as an

off-the-shelf solution to scalability, availability, and several other problems, and focus instead on the content of the service being developed.

We now describe our proposed system architecture and service-programming model for building scalable TACC servers using clusters of PCs. The architecture attempts to address the challenges of cluster computing (unwieldy administration, managing partial failures, and the lack of shared state across components) while exploiting the strengths of cluster computing (support for incremental scalability, high availability through redundancy, and the ability to use commodity building blocks). A more detailed discussion of the architecture can be found in [11].

The goal of our architecture is to separate the *content* of network services (i.e., what the services do) from their implementation by encapsulating the “scalable network service” (SNS) requirements of high availability, scalability, and fault tolerance in a reusable layer with narrow interfaces. Application writers program to the TACC APIs alluded to in the previous section, without regard to the underlying TACC server implementation; the resulting TACC applications automatically receive the benefits of linear scaling, high availability, and failure management when run on our cluster-based TACC server.

The software-component block diagram of a scalable TACC server is shown in Fig. 4. Each physical workstation in a network of workstations (NOW) [12] supports one or more software components in the figure, but each component in the diagram is confined to one node. In general, the components whose tasks are naturally parallelizable are replicated for scalability, fault tolerance, or both.

Frontends provide the interface to the TACC server as seen by the outside world (e.g., HTTP server). They “shepherd” incoming requests by matching them up with the appropriate user profile from the customization database, and queuing them for service by one or more workers. Frontends maximize system throughput by maintaining state for many simultaneous outstanding requests, and can be replicated for both scalability and availability.

The *worker pool* consists of caches (currently Harvest [13]) and service-specific modules that implement the actual service (data transformation/filtering, content aggregation, etc.). Each type of module may be instantiated zero or more times, depending on offered load. The TACC API allows all cache workers to be managed as a single virtual cache by providing URL hashing, automatic failover, and dynamic growth of the cache pool.

The *customization database* stores user profiles that allow mass customization of request processing. The manager balances load across workers and spawns additional workers as offered load fluctuates or faults occur. When necessary, it may assign work to machines in the overflow pool, a set of backup machines (perhaps on desktops) that can be harnessed to handle load bursts and provide a smooth transition during incremental growth.

The *load balancing/fault tolerance manager* keeps track of what workers are running where, autostarts new workers as needed, and balances load across workers. Its detailed operation is described next, in the context of the TranSend implementation. Although it is a centralized agent, [11] describes the various mechanisms, including multicast heartbeat and process-peer fault tolerance, that keep this and other system components running and allow the system to survive transient component failures.

The *graphical monitor* for system management supports asynchronous error notification via e-mail or pager, temporary

disabling of system components for hot upgrades, and visualization of the system's behavior using Tcl/Tk [14]. The benefits of visualization are not just cosmetic: we can immediately detect by looking at the visualization panel what state the system as a whole is in, whether any component is currently causing a bottleneck (such as cache-miss time, distillation queuing delay, interconnect), what resources the system is using, and similar behaviors of interest.

The *interconnect* provides a low-latency, high-bandwidth, scalable interconnect, such as switched 100-Mb/s Ethernet or Myrinet [15]. Its main goal is to prevent the interconnect from becoming the bottleneck as the system scales.

Components in our TACC server architecture may be replicated for fault tolerance or high availability, but we also use replication to achieve scalability. When the offered load to the system saturates the capacity of some component class, more instances of that component can be launched on incrementally added nodes. The duties of our replicated components are largely independent of each other (because of the nature of the Internet services' workload), which means the amount of additional resources required is a linear function of the increase in offered load.

Analysis of the TranSend Implementation

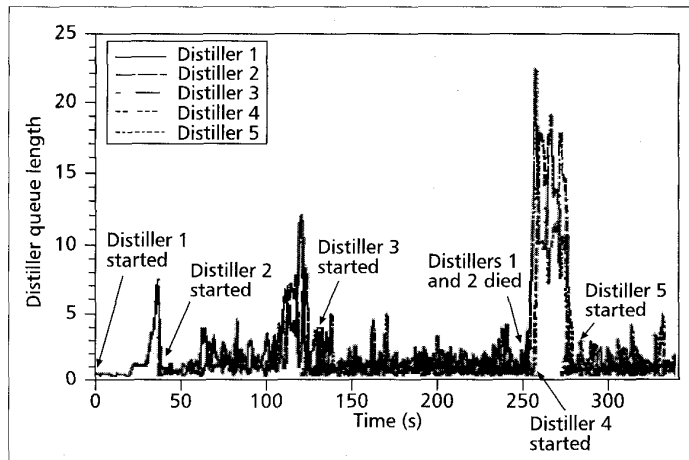
TranSend [9], a TACC reimplementation of our earlier prototype called Pythia [15], performs lossy Web image compression on the fly. Each TranSend worker handles compression or markup for a specific MIME type; objects of unsupported types are passed through to the user unaltered.³

We took measurements of TranSend using a cluster of 15 Sun SPARC Ultra-1 workstations connected by 100 Mb/s switched Ethernet and isolated from external load or network traffic. For measurements requiring Internet access, the access was via a 10 Mb/s switched Ethernet network connecting our workstation to the outside world. Many of the performance tests are based on HTTP trace data from the 25,000 UC Berkeley dialup IP users [17], played back using a high-performance playback engine of our own design that can either generate requests at a constant rate or faithfully play back a trace according to the timestamps in the trace file. In the following subsections we report on experiments that stress TranSend's fault tolerance, responsiveness, and scalability.

Self-Tuning and Load Balancing – As mentioned previously, the load balancing and fault tolerance manager is

Requests/s	# FEs	# workers	Element that saturated
0–24	1	1	Workers
25–47	1	2	Workers
48–72	1	3	Workers
73–87	1	4	FE Ethernet
88–91	2	4	Workers
92–112	2	5	Workers
113–135	2	6	Workers + FE Ethernet
136–159	3	7	Workers

■ **Table 5.** Results of the scalability experiment. FE: frontend.



■ **Figure 5.** Worker queue lengths observed over time as the load presented to the system fluctuates, and workers are manually brought down.

charged with spawning and reaping workers and distributing internal load across them. The mechanisms by which this is accomplished, which include monitoring worker queue lengths and applying some simple hysteresis, are described in [11].

Figure 5 shows the variation in worker queue lengths over time. The system was bootstrapped with one frontend and the manager, and a single demand-spawned worker. Continuously increasing the load caused the manager to spawn a second and later a third worker. We then manually killed the first two workers; the sudden load increase on the remaining worker caused the manager to spawn one and later another new worker, to stabilize the queue lengths.

Scalability – To demonstrate the scalability of the system, we performed the following experiment:

- We began with a minimal instance of the system: one frontend, one worker, the manager, and a fixed number of cache partitions. (Since for these experiments we repeatedly requested the same subset of images, the cache was effectively not tested.)
- We increased the offered load until some system component saturated (e.g., worker queues grew too long, frontends no longer accepted additional connections).
- We then added more resources to the system to eliminate this saturation (in many cases the system does this automatically, as when it recruits overflow nodes to run more workers), and we recorded the amount of resources added as a function of the increase in offered load, measured in requests per second.
- We continued until the saturated resource could not be replenished (i.e., we ran out of hardware), or until adding more of the saturated resource no longer resulted in a linear or close-to-linear improvement in performance.

Table 5 presents the results of this experiment. At 24 requests/s, as the offered load exceeded the capacity of the single available worker, the manager automatically spawned one additional worker, and then subsequent workers as necessary. (In addition to using faster hardware, the performance engineering of the cluster-based server has caused a large reduction in the amortized cost of distillation for a typical image, compared to the values suggested by Fig. 2.) At 87 requests/s, the Ethernet segment leading into the frontend

³ The PostScript-to-rich-text worker described earlier has not yet been added to TranSend.

saturated, requiring a new frontend to be spawned. We were unable to test the system at rates higher than 159 requests/s, because all of our cluster's machines were hosting workers, frontends, or playback engines. We did observe nearly perfectly linear growth of the system over the scaled range: a worker can handle approximately 23 requests/s, and a 100 Mb/s Ethernet segment into a frontend can handle approximately 70 requests/s. We were unable to saturate the frontend or the cache partitions, or fully saturate the interior interconnect during this experiment. We draw two conclusions from this result:

- Even with a commodity 100 Mb/s interconnect, linear scaling is limited primarily by bandwidth into the system rather than bandwidth inside the system.
- Although we originally deployed TranSend on four SPARC 10s, a single Ultra-1-class machine would suffice to serve the entire dialup IP population of UC Berkeley (25,000 users officially, over 8000 of whom surfed during the trace).

Other TACC Applications

We now discuss several examples of new services in various stages of deployment, showing how each exploits the TACC model and discussing some of our experiences with the applications. Rather than providing detailed measurements as we did for TranSend in the previous section, the present goal is to demonstrate the flexibility of the TACC framework in accommodating an interesting range of applications, while providing consistent guidelines for approaching application partitioning decisions.

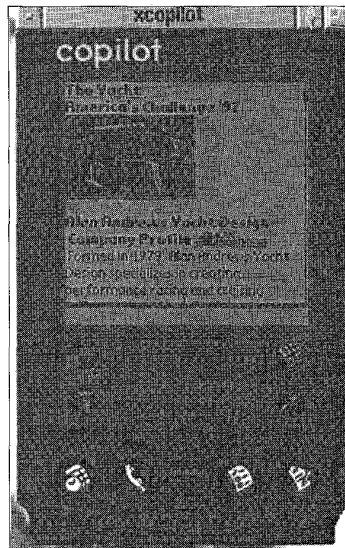
We restrict our discussion here to services that can be implemented using the proxy model (i.e., transparent interposition of computation between clients and servers). (Some of our services do not communicate via HTTP but are conceptually similar.) Also, although we have developed a wider range of applications using the TACC model as part of a graduate seminar [18], we concentrate on those applications that enable adaptation to network and client variation. These services share the following common characteristics, which make them amenable to implementation on our cluster-based framework:

- Compute-intensive transformation or aggregation
- Computation that is parallelizable with granularity of a few CPU seconds
- Substantial value added by mass customization

TranSend as a TACC Application

TranSend is one of the simplest TACC applications we have produced. The dispatch rules simply match the MIME type of the object returned from the origin server to the list of known workers, which (as in all TACC applications) can be updated dynamically. In particular, TranSend does not exploit TACC's ability to compose workers by chaining them into a "pipeline" or having one worker call others as coroutines.

Transformed objects are stored in the cache with "fat URLs" that encode a subset of the transformation parameters, saving the work of retransforming an original should another user ask for the same degraded version later. Each user can select a desired level of aggressiveness for the lossy compression and choose between HTML and Java-based interfaces for modifying their preferences.



■ **Figure 6.** Screenshot of the *Top Gun Wingman* browser. This screenshot is taken from the "xcopilot" hardware-level Pilot simulator [21].

The main difference between TranSend and commercial products based on its ideas (such as Intel's recently announced QuickWeb [19]) is extensibility: adding support for new datatypes to TranSend is as simple as adding a new worker, and composing workers is as simple as modifying the dispatch rules (or modifying existing workers to hint to the TACC server that it should fall through to new workers).

In fact, we have generalized TranSend into a "lazy fixations" system [20] in which users could select from among a variety of available formats for viewing an object; this was implemented by a "graph search" worker that treated all the transformation workers as edges in a directed graph and performed a shortest-paths search to determine which sequence of workers should be run to satisfy a particular request.

One of the goals of TACC is to exploit modularity and composition to make new services easy to prototype by reusing existing building blocks. TranSend's HTML and

JPEG workers consist almost entirely of off-the-shelf code, and each took an afternoon to write. A pair of anecdotes illustrates the flexibility of the TACC APIs in constructing responsive services. Our original HTML parser was a fast C language implementation from the W3C. Debugging the pathological cases for this parser was spread out over a period of days—since our prototype TACC server masks transient faults by bypassing original content "around" the faulting worker, we could only deduce the existence of bugs by noticing (on the graphical monitor display) that the HTML worker had been restarted several times over a period of hours, although the service as a whole was continuously available.

We later wrote a much slower but more robust parser in Perl to handle proprietary HTML extensions such as inline JavaScript. All HTML pages are initially passed to the slower Perl parser, but if it believes (based on page length and tag density) that processing the page will introduce a delay longer than one or two seconds, it immediately throws an exception and indicates that the C parser should take over. Because the majority of long pages tend to be papers published in HTML rather than complex pages with weird tags, this scheme exploits TACC composition and dispatch to handle common cases well while keeping HTML processing latency barely noticeable.

Top Gun Wingman

Top Gun Wingman is the only graphical Web browser available for the 3Com PalmPilot, a typical "thin client" device. Based on file downloads, we estimate that 8000 to 10,000 users are using the client software and UC Berkeley's experimental cluster; ProxiNet, Inc. has since commercialized the program and deployed a production cluster to serve additional users. Figure 6 shows a screenshot of the browser.

Previous attempts to provide graphical Web browsing on such small devices have foundered on the severe limitations imposed by small screens, limited computing capability, and austere programming environments, and virtually all have fallen back to simple text-only browsing. Our adaptation approach, combined with the composable-workers model provided by TACC, allows us to approach this problem from a different perspective. The core of Top Gun Wingman consists of three TACC workers: HTML layout, image conversion, and

conversion of intermediate-form layout to device-specific data format. These three workers address the three areas of variation introduced earlier:

- **Hardware and software adaptation** — We have built TACC workers that output simplified binary markup and scaled-down images ready to be “spoon fed” to a thin-client device, given knowledge of the client’s screen dimensions, image format, and font metrics. This greatly simplifies client-side code since no HTML parsing, layout, or image processing is necessary, and as a side benefit, the smaller and more efficient data representation reduces transmission time to the client. The image worker delivers 2-b/pixel images, since that is what the PalmPilot hardware supports, and the HTML parsing and layout worker ensures that no page description larger than about 32 kbytes is delivered to the client, since that is the approximate heap space limit imposed by the PalmPilot’s programming environment. We have also added three “software upgrades” at the proxy since Wingman was first deployed: a worker that delivers data in AportisDoc [22] format (a popular PalmPilot e-book format), a worker that extracts and displays the contents of software archives for download directly to the PalmPilot, and an improved image-processing module contributed by a senior graphics hacker. In terms of code footprint, Wingman weighs in at 40 kbytes of code (compared with 74 kbytes and 109 kbytes for HandWeb and Palm-escape 5.0, respectively, neither of which currently support image viewing).
- **Network protocol adaptation:** In addition to delivering data in a more compact format and exploiting datatype-specific distillation, we have replaced HTTP with a simpler, datagram-oriented protocol based on application-level framing [23]. The combined effect of these optimizations is that Wingman is two to four times faster than a desktop browser loading the same Web pages over the same bandwidth, and Wingman’s performance on text-only pages often exceeds that of HTML/HTTP-compliant browsers on the same platform, especially on slow (< 56 kb/s) links.

Top Gun Mediaboard

Top Gun Mediaboard is an electronic shared whiteboard application for the PalmPilot. This is a derivative of the desktop *mediaboard* application, which uses Scalable Reliable Multicast (SRM) as the underlying communication protocol. A reliable multicast proxy (RMX) TACC worker participates in the SRM session on behalf of the PDA clients, performing four main types of client adaptation:

- **Transport protocol conversion** — The PalmPilot’s network stack does not support IP multicast. The RMX converts the SRM data into a unicast TCP stream that the client can handle.
- **Application protocol adaptation** — To keep the client implementation simple, all the complexities of the mediaboard command protocol are handled by the RMX. The protocol adapter transforms the entire sequence of mediaboard commands into a “pseudo-canvas” by executing each command and storing its result in the canvas, transmitting only a sequence of simple draw-ops to the client. The protocol and data format for transmitting the draw-ops is a direct extension of the Top Gun Wingman datagram protocol.
- **On-demand distillation** — The RMX converts specific data objects according to the client’s needs. For example, it transforms the GIF and JPEG images that may be placed on the mediaboard into simpler image representations that the PalmPilot can understand, using the same worker that is part of Wingman. The client application can refine (zoom in on) specific portions of the canvas.

- **Intelligent rate limiting** — Since the proxy has complete knowledge of the client’s state, the RMX can perform intelligent forwarding of data from the mediaboard session to the client. By eliminating redundant draw-ops (e.g., *create* followed by *delete* on the same object) before sending data to the client, the RMX reduces the number of bytes that must be sent over the low-bandwidth link. Moreover, although a whiteboard session can consist of a number of distinct pages, the RMX forwards only the data associated with the page currently being viewed on the client.

Top Gun Mediaboard is in pre-alpha use at UC Berkeley, and performs satisfactorily even over slow links such as the Metricom Ricochet wireless packet radio modem [24].

Charon: Indirect Authentication for Thin Clients

Although not yet rewritten as a TACC application, Charon [25] illustrates a similar use of adaptation by proxy, for performing indirect authentication. In particular, Charon mediates between thin clients and a Kerberos [26] infrastructure. Charon is necessary because, as we describe in [25], the computing resources required for a direct port of Kerberos to thin clients are forbidding. With Charon, Kerberos can be used to authenticate both clients to the proxy service, and the proxied clients to Kerberized servers. Charon relieves the client of a significant amount of Kerberos protocol processing, while limiting the amount of trust that must be placed in the proxy; in particular, if the proxy is compromised, existing user sessions may be hijacked but no new sessions can be initiated, since new sessions require cooperation between the client and proxy. Our Charon prototype client for the Sony MagicLink [27], a once-popular PDA, had a client footprint of only 45 kbytes, including stack and heap usage.

Related Work

At the network level, various approaches have been used to shield clients from the effects of poor (especially wireless) networks [28, 29]. At the application level, data transformation by proxy interposition has become particularly popular for HTTP, whose proxy mechanism was originally intended for users behind security firewalls. The mechanism has been harnessed for anonymization [30], Kanji transcoding [31, 32], application-specific stream transformation [33], and personalized “associates” for Web browsing [34, 35]. Some projects provide an integrated solution with both network-level and application-level mechanisms [36, 37, 10], although none propose a uniform application-development model analogous to TACC.

Rover [38], Coda [39], and Wit [4] differ in their respective approaches to partitioning applications between a thin or poorly connected client and a more powerful server. In particular, Rover and Coda provide explicit support for disconnected operation, unlike our TACC work. We find that Rover’s application-specific, toolkit-based approach is a particularly good complement to our own; although the TACC model provides a reasonable set of guidelines for thinking about partitioning (leave the client to do what it does well, and move as much as possible of the remaining functionality to the back end), we are working on integrating Rover into TACC to provide a rich abstraction for dealing with disconnection in TACC applications.

SmartClients [40] and SWEB++ [41] have exploited the extensibility of client browsers via Java and JavaScript to enhance scalability of network-based services by dividing labor between the client and server. We note that our system does not preclude and, in fact, benefits from exploiting intelligence

and computational resources at the client; we discuss various approaches we have tried in [42].

Lessons and Conclusions

We proposed three design principles for adapting to network and client variation and delivering a meaningful Internet experience to impoverished clients: datatype-specific distillation and refinement, adaptation on demand, and moving complexity into the infrastructure. We also offered a high-level description of the TACC programming model (transformation, aggregation, caching, customization) that we have evolved for building adaptive applications, and presented measurements of our scalable, highly available, cluster-based TACC server architecture, focusing on the TranSend Web accelerator application. Finally, we described other applications we have built that are in daily use, including some that push the limits of client adaptation (such as Top Gun Wingman and Top Gun Mediaboard). In this section we try to draw some lessons from what we have learned from building these and similar applications and experimenting with our framework.

Aggressively pushing the adaptation-by-proxy model to its limits, as we have tried to do with Top Gun Wingman and Top Gun Mediaboard, has helped us validate the proxy-interposition approach for serving thin clients. Our variation on the theme of application partitioning has been to split the application between the client and the proxy, rather than between the client and the server. This has allowed our clients to access existing content with no server modifications. Our guideline for partitioning applications has been to allow the client to perform those tasks it does well in native code, and move as much as possible of the remaining work to the proxy. For example, since most thin clients support some form of toolkit for building graphical interfaces, sending HTML markup is too cumbersome for the client, but sending screen-sized bitmaps is unnecessarily cumbersome for the proxy.

A frequent objection raised against our partitioning approach is that it requires that the proxy service be available at all times, which is more difficult than simply maintaining the reliability of a bank of modems and routers. This observation motivated our work on the cluster-based scalable and highly available server platform described earlier, and in fact the TranSend and Wingman proxy services have been running for several months at UC Berkeley with high stability, except for a two-week period in February 1998 when the cluster was affected by an OS upgrade. Other than one part-time undergraduate assistant, the cluster manages itself, yet thousands of users have come to rely on its stability for using Top Gun Wingman, validating the efficacy of our cluster platform. This observation, combined with the current trends toward massive cluster-based applications such as HotBot [43], suggests to us that the adaptive proxy style of adaptation will be of major importance in serving convergent "smart-phone"-like devices.

Acknowledgments

This project has benefited from the detailed and perceptive comments of countless anonymous reviewers, users, and collaborators. Ken Lutz and Eric Fraser configured and administered the test network on which the TranSend scaling experiments were performed. Cliff Frost of the UC Berkeley Data Communications and Networks Services group allowed us to collect traces on the Berkeley dialup IP network and has worked with us to deploy and promote TranSend within UC Berkeley. Undergraduate researchers Anthony Polito, Benjamin Ling, Andrew Huang, David Lee,

and Tim Kimball helped implement various parts of TranSend and Top Gun Wingman. Ian Goldberg and David Wagner helped us debug TranSend, especially through their implementation of the Anonymous Rewebber [44]. Ian implemented major parts of the client side of Top Gun Wingman, especially the 2-bit-per-pixel hacks. Paul Haeberli of Silicon Graphics contributed image processing code for Top Gun Wingman. Murray Mazer at the Open Group Research Institute has provided much useful insight on the structure of Internet applications and future extensions of this work. We also thank the patient students of UCB Computer Science 294-6, *Internet Services*, Fall 1997, for being the first real outside developers on our TACC platform and greatly improving the quality of the software and documentation. We have received much valuable feedback from our UC Berkeley colleagues, especially David Culler, Eric Anderson, Trevor Pering, Hari Balakrishnan, Mark Stemm, and Randy Katz. This research is supported by DARPA contracts #DAAB07-95-CD154 and #J-FBI-93-153, the California MICRO program, the UC Berkeley Chancellor's Opportunity Fellowship, the NSERC PGS-A fellowship, Hughes Aircraft Corp., and Metricom Inc.

References

- [1] Nokia Corp. and Geoworks Inc., Nokia 9000 Communicator, <http://www.geoworks.com/devices/9000>.
- [2] Tom R. Halfhill, "Inside the web pc," *Byte Mag.*, Mar. 1996, pp. 44-56.
- [3] M. Stemm and R. H. Katz, "Vertical handoffs in wireless overlay network," *ACM Mobile Networking (MONET)*, Special Issue on Mobile Networking in the Internet, Fall 1997.
- [4] T. Watson, "Application Design for Wireless Computing," *Mobile Comp. Sys. and Appl. Wksp.*, Aug. 1994.
- [5] Graphics interchange format version 89a (GIF), CompuServe Inc., Columbus, OH, July 1990.
- [6] J. Poskanzer, Netpbm release 7, <ftp://wuarchive.wustl.edu/graphics/graphics/packages/NetPBM>, 1993.
- [7] E. Amir, S. McCanne, and H. Zhang, "An application level video gateway," *Proc. ACM Multimedia 1995*, 1995.
- [8] P. MacJones, DEC SRC pers. commun., PostScript-to-text converter.
- [9] A. Fox et al., "TranSend web accelerator proxy," free service deployed by UC Berkeley; <http://transend.cs.berkeley.edu>, 1997.
- [10] B. Zenel, "A Proxy Based Filtering Mechanism for the Mobile Environment," Thesis proposal, Mar. 1996.
- [11] A. Fox et al., "Cluster-Based Scalable Network Services," *Proc. 16th ACM Symp. Op. Sys. Principles*, St.-Malo, France, Oct. 1997.
- [12] T. E. Anderson, D. E. Culler, and D. Patterson, "The case for NOW (networks of workstations)," *IEEE Micro*, vol. 12, no. 1, Feb. 1995, pp. 54-64.
- [13] A. Chankunthod et al., "A hierarchical internet object cache," *Proc. 1996 Usenix Annual Tech. Conf.*, Jan. 1996, pp. 153-63.
- [14] J. K. Ousterhout, *Tcl and the Tk Toolkit*, Addison-Wesley, 1994.
- [15] Myricom, "Myrinet: A Gigabit Per Second Local Area Network," *IEEE Micro*, Feb. 1995.
- [16] A. Fox and E. A. Brewer, "Reducing WWW Latency and Bandwidth Requirements via Real-Time Distillation," *Proc. 5th Int'l World Wide Web Conf.*, Paris, France, May 1996.
- [17] S. D. Gribble and E. A. Brewer, "System Design Issues for Internet Middleware Services: Deductions from a Large Client Trace," *Proc. 1997 USENIX Symp. Internet Tech. and Sys.*, Monterey, CA, Dec. 1997.
- [18] A. Fox and E. A. Brewer, "CS 294-6: Internet services, class proceedings," Fall 1997, <http://www.cs.berkeley.edu/~fox/cs294>.
- [19] Intel Corp. QuickWeb Web Accelerator.
- [20] A. Fox and S. D. Gribble, "DOLF: Digital objects with lazy fixations," unpublished manuscript, CS 294-5 Digital Libraries Seminar, Spring 1996.
- [21] I. Curtis, xcopilot Pilot simulator, 1998.
- [22] Aportis Inc., "AportisDoc Overview," 1998, <http://www.aportis.com/products/AportisDoc/benefits.html>.
- [23] D. D. Clark and D. L. Tennenhouse, "Architectural Considerations for a New Generation of Protocols," *Comp. Commun. Rev.*, vol. 20, no. 4, Sept. 1990, pp. 200-8.
- [24] Metricom Corp., "Ricochet wireless modem," 1998, <http://www.ricochet.net>.
- [25] A. Fox and S. D. Gribble, "Security On the Move: Indirect Authentication Using Kerberos," *Proc. 2nd Int'l. Conf. Wireless Networking and Mobile Comp. (MobiCom '96)*, Rye, NY, Nov. 1996.
- [26] J. G. Steiner, C. Neuman, and J. I. Schiller, "Kerberos: An authentication service for open network systems," *Proc. USENIX Winter Conf.*, 1988, Dallas, TX, Feb. 1988, pp. 191-202.
- [27] Sony Corp., "The Sony MagicLink PDA."

- [28] H. Balakrishnan et al., "A comparison of mechanisms for improving tcp performance over wireless links," *Proc. ACM SIGCOMM '96*, Stanford, CA, Aug. 1996.
- [29] H. Balakrishnan et al., "Improving tcp/ip performance over wireless networks," *Proc. 1st ACM Conf. Mobile Comput. and Networking*, Berkeley, CA, Nov. 1995.
- [30] C2net, Web anonymizer.
- [31] Y. Sato, DeleGate Server, Mar. 1994; [http:// wall.etl.go.jp/delegate/](http://wall.etl.go.jp/delegate/).
- [32] K. P. Yee, Shoduoka Mediator Service, 1995, <http://www.shoduoka.com>.
- [33] C. Brooks, M. S. Mazer, S. Meeks, and J. Miller, "Application-Specific Proxy Servers as HTTP Stream Transducers," *Proc. 4th Int'l. World Wide Web Conf.*, Dec. 1995.
- [34] R. Barrett, P. P. Maglio, and D. C. Kellem, "How To Personalize the Web," Conf. Human Factors in Comp. Sys. (CHI '95), Denver, CO, May 1995. WBI, developed at IBM Almaden; <http://www.raleigh.ibm.com/wbi/wbisoft.htm>.
- [35] M. A. Schickler, M. S. Mazer, and C. Brooks, "Pan-browser support for annotations and other metainformation on the world wide web," *Proc. 5th Int'l. World Wide Web Conf. (WWW-5)*, May 1996.
- [36] M. Liljeberg et al., "Enhanced services for world wide web in mobile WAN environment," Tech. rep. C-1996-28, Univ. of Helsinki CS Dept., Apr. 1996.
- [37] WAP Forum, Wireless application protocol (WAP) forum, <http://www.wapforum.org>.
- [38] A. D. Joseph et al., "Rover: A toolkit for mobile information access," *Proc. 15th ACM Symp. Op. Sys. Principles*, Copper Mountain Resort, CO, Dec. 1995.
- [39] J. J. Kistler and M. Satyanarayanan, "Disconnected Operation in the Coda File System," *ACM Trans. Comp. Sys.*, vol. 10, Feb. 1992, pp. 3-25.
- [40] C. Yoshikawa et al., "Using smart clients to build scalable services," *Proc. Winter 1997 USENIX Tech. Conf.*, Jan. 1997.
- [41] D. Andresen et al., "Scalability issues for high performance digital libraries on the world wide web," *Proc. IEEE ADL '96*, Forum on Research and Technology Advances in Digital Libraries, Washington, DC, May 1996.
- [42] A. Fox et al., "Orthogonal Extensions to the WWW User Interface Using Client-Side Technologies," *User Interface Software and Technology (UIST) '97*, Banff, Canada, Oct. 1997.
- [43] Inktomi Corp., The HotBot search engine.

- [44] I. Goldberg and D. Wagner, "TAZ servers and the rewebber network: Enabling anonymous publishing on the world wide web," Unpublished manuscript available at <http://www.cs.berkeley.edu/~daw/cs268/>, May 1997.

Biographies

ARMANDO FOX (fox@cs.berkeley.edu, fox@alum.mit.edu) received a B.S.E.E. from MIT and an M.S.E.E. from the University of Illinois, has worked as a CPU architect at Intel Corp., and is currently completing his Ph.D. at the University of California, Berkeley, as a researcher in the Daedalus/BARWAN and InfoPad projects. He will be an assistant professor at Stanford University starting in January 1999. His primary interests are application-level support for adaptive mobile computing, multimedia networks for mobile computing, and user interface issues related to mobile computing.

STEVEN GRIBBLE (gribble@cs.berkeley.edu) received a combined computer science and physics B.S. degree from the University of British Columbia, Vancouver, Canada, in 1995, and a computer science M.S. degree from the University of California, Berkeley in 1997. He is currently pursuing his Ph.D. at Berkeley, and expects to graduate in May 2000. His interests include application-level support for adaptive mobile computing, and system and compiler support for scalable, highly available infrastructure services.

YATIN CHAWATHE (yatin@cs.berkeley.edu) is a doctoral student at the University of California, Berkeley. He received a Bachelor of Engineering (computer engineering) degree from the University of Bombay, India, and an M.S. (computer science) from the University of California, Berkeley. His primary interests are internet services, application support for reliable multicast in heterogeneous environments, and multimedia networking.

ERIC A. BREWER (brewer@cs.berkeley.edu) is an assistant professor of computer science at U.C. Berkeley, and received his Ph.D. in computer science from MIT in 1994. His research interests include mobile and wireless computing (the InfoPad and Daedalus projects), scalable servers (the NOW, Inktomi, and TACC projects), and application- and system-level security (the ISAAC project). His previous work includes multiprocessor-network software and topologies, and high-performance multiprocessor simulation.