

A Composable Framework for Secure Multi-Modal Access to Internet Services from Post-PC Devices

Steven J. Ross, Jason L. Hill, Michael Y. Chen, Anthony D. Joseph, David E. Culler, and Eric A. Brewer

Computer Science Department, University of California - Berkeley

E-mail: {stevross, jhill, mikechen, adj, culler, brewer}@cs.berkeley.edu

Abstract

The Post-PC revolution is bringing information access to a wide-range of devices beyond the desktop, such as public kiosks, and mobile devices like cellular telephones, PDAs, and voice based vehicle telematics. However, existing deployed Internet services are geared toward the secure rich interface of private desktop computers. We propose the use of an infrastructure-based secure proxy architecture to bridge the gap between the capabilities of Post-PC devices and the requirements of Internet services. By combining generic content and security transformation functions with service-specific rules, the architecture decouples device capabilities from service requirements and simplifies the addition of new devices and services. Security and protocol specifics are abstracted into re-usable components. Additionally, the architecture offers the novel ability to deal with untrusted public Internet access points by providing fine-grain control over the content and functionality exposed to the end device, as well as support for using trusted and untrusted devices in tandem. Adding support for a deployed Internet service requires a few hundred lines of scraping scripts. Similarly, adding support for a new device requires a few hundred lines of stylesheets for the device format. The average latency added by proxy transformations is around three seconds in our unoptimized Java implementation.

1 Introduction

Internet services for electronic communication and commerce have permeated our society. People routinely access and exchange private and sensitive information when they use Internet services, such as electronic mail, the World-Wide Web, electronic banking, on-line shopping, E-commerce, and stock trading. At the same time, we are moving into a “Post-PC” era, where people use many different types of devices to access these Internet services, including pay-per-use public Internet kiosks, laptops, PDAs,

cellular telephones, and two-way pagers, each with different characteristics and capabilities. Three critical requirements for users of Internet services in these Post-PC environments are *secure access to information*, *dynamic content adaptation*, and the *fusion of multiple devices*. The *fusion of multiple devices* represents an extension of the split trust model [5] where multiple devices work in tandem. Meeting these requirements will preserve the privacy and integrity of users’ transactions and enable the use of new and interesting services and access devices.

Secure, wide-area access to information is of paramount importance on the Internet, but existing solutions do not address Post-PC requirements (e.g., traffic can be intercepted and user identities can be forged). Protection against such attacks is provided by strong authentication and end-to-end encryption. A user’s *trusted* application establishes a secure connection to the service of interest. For example, when accessing a secure web-based banking service, the user’s web browser establishes a secure connection to the web server by using HTTP over Secure Socket Layer (SSL) [24].

The existing security model makes two basic assumptions: first, that both the user’s access device and the software running on it can be trusted not to intercept or send private information elsewhere; second, the access device has the computational resources to secure the connection (e.g., to perform public key SSL operations) and to display the service’s content (e.g., to render complex color graphics). While these assumptions are reasonable for users of private desktop computers, they do not hold for Post-PC devices, where users access services from both untrusted devices (such as public kiosks), and devices with limited resources (such as PDAs).

Untrusted kiosks are problematic for secure Internet services, as all service data is available in unencrypted form to the kiosk. Possible attacks may be as straightforward as recording all keystrokes (e.g., Personal Identification Numbers, passwords, and login information), or as subtle as recording users’ personal information for later fraudulent use (e.g., account numbers). Additionally, a kiosk can “hijack” a secure connection and perform active attacks (e.g.,

make bank transfers or send forged e-mails). Untrusted endpoints should not be allowed to see a user's personal information; *instead, the content value of such data must be reduced*. Likewise, access to sensitive service functionality from these endpoints must also be guarded.

PDA's are also problematic because they are generally low-power, computationally-limited devices with limited memory and networking capabilities. To perform the industry-standard SSL handshake phase on one such device, a 3Com Palm Pilot V [33] requires several seconds. This latency imposes an intolerable delay for connection setup, which is particularly undesirable if network connectivity is intermittent. An SSL implementation that uses elliptic curve cryptography [9] is feasible on a Palm Pilot V, but few Internet services support that option. Moreover, even if a PDA is capable of performing SSL, one may still opt for a protocol that is faster and more power-efficient. Finally, entering data to fill out forms using a pen-based interface is tedious at best and even more cumbersome using number pads on devices such as cellular telephones.

Demand for continuous access to Internet content is increasing and Internet services are becoming available in new environments such as automobiles [26] and kiosks in airplane seats [25]. These environments demand multi-modal access to content either through traditional HTML, PDA thin client browsers, WML enabled phones, and voice.

We propose the use of a trusted infrastructure-based proxy service to provide secure multi-modal access to Internet service content from any device. An advantage of this proxy-based approach is the layer of indirection it provides which allows transparent support for widely deployed systems owned and operated by others. For untrusted terminals, sensitive user information is removed from the information stream going to the terminal. For example, real names and account numbers can be screened out. The proxy also mediates access to the service by filtering control data in order to limit allowable actions (e.g., disallowing all stock trades on untrusted kiosks). For access from computationally-limited devices, the proxy service transcodes the security protocol into one which is more efficiently executed on the device. Additionally, the proxy can distill the service content into a format more suitable for the device [16]. For example, a HTML service can be rendered as WML for a WAP enable phone, or even as voice.

A drawback of the proxy approach is that it requires users to place a significant amount of trust in the proxy infrastructure. The components in the infrastructure require access to all data flowing through them to perform filtering and adaptation transformations. Additionally, users are required to store private information and preferences in the infrastructure to enable these operations. This trust can be mediated by colocating the proxy with a particular Internet service in which the user already has established trust

of sensitive information. For example, an Internet e-mail or stock-trading service can provide its own proxy service for access from mobile devices and public terminals. Alternately, the proxies can be treated as orthogonal to the services which they enable access; a user need only establish an account with a secure multi-modal proxy provider, much as they do with an Internet service provider today.

In addition to enabling basic access, our security proxy can be used to combine the capabilities of both untrusted public terminals and mobile personal devices. For example, the limited GUI of a PDA can be supplemented by the richer, larger display of a public kiosk, while the PDA is used only for sensitive operations. The security proxy splits trust between the PDA and the public terminal by fusing the devices together to provide one logical channel with secure access to the end service. Consider the case of users accessing their stock trading accounts from public access terminals. Instead of relying on the terminal to protect their secure information, the users could direct private or sensitive information (e.g., portfolio value and account information) to their portable device, while using the GUI capabilities of the public terminal to initiate requests and display generic stock information (e.g., stock price fluctuations and historical graphs). The users initiate trading operations through the untrusted public terminals and then confirm them using their trusted portable devices. The connections to the mobile devices can be provided by the environment (e.g., kiosks with infrared network connections for PDA's) or by the devices (e.g., receiving a call on a cell-phone, or a message on a two-way pager). The proxy allows data to be encrypted with a protocol more appropriate for the capabilities of computationally limited and power constrained small devices such as elliptical curve or shared key techniques. This capability bridges the gap to existing services that do not currently provide this option for smaller devices.

We present a secure proxy architecture that:

- Protects users from untrusted access points.
- Enables users to use untrusted access points in tandem with trusted mobile devices for security-enhanced service interactions.
- Allows users to use the access device of *their* choice regardless of the device's computational abilities.
- Simplifies the tailoring a service to multiple device formats to simple content authoring of style sheets and scraping scripts.

Our secure proxy incorporates a component-oriented, rule-based architecture which minimizes the amount of service-specific code that must be written. In particular, service developers (and users) can specify security rules that make the appropriate modifications to the data and control

actions flowing in both directions between the user and the service. Adding support for new services only requires authoring a XML representation of the semantic content and a WebL script to extract the content from the service. Adding support for new devices is done simply through XSL style-sheets that render the content to the devices' format.

2 Architecture

In this section, we present an overview of our architecture for an infrastructure-based proxy for secure multi-modal access to Internet service content, followed by a detailed design discussion of each component.

2.1 Design overview

Our secure proxy for multi-modal access to Internet services consists of a few building block components and the canonical path between them. The proxy provides a secure level of indirection that can be used to modify content flowing between clients and services. It is a component-based system for creating customizable building blocks that control and modify users' connections to secure services. These components can easily be composed together to bridge the gap between traditional services and Post-PC devices that have different trust levels and to simplify the process of adding support for new services and devices. This model provides a partitioning of computation that allows the components to be placed in a scalable, fault-tolerant, and highly-available execution environment such as [20] that meets the critical requirements of Internet-scale services.

Our architecture allows the user to take on different *trust profiles* based upon the location, device, capabilities, or trust in an access device. For example, these profiles include connecting from home devices (fully trusted), work computers (partially trusted), public terminals (not trusted), PDAs (trusted, but not powerful), or with an untrusted device in tandem with a trusted one (split-trust).

By dividing the functionality of our architecture into structural and semantic transformations, it is relatively easy to add new services and devices. *Format transformations* adapt content to device capabilities (e.g., changing the resolution, color depth, or encoding format of an image), while *semantic transformations* protect users' data and prevent various actions from being performed from untrusted end devices. Additionally, the use of rule-based components to add service-specific semantics to generic transformations simplifies the incorporation of new services.

The architecture consists of the following components:

- **Security Adaptors (SA):** These components bridge the gap between the secure access protocol required by a service (server-side) and access devices (client-side).

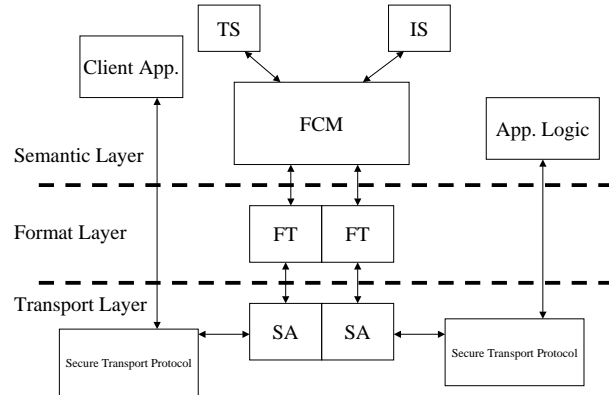


Figure 1. Secure proxy architecture

- **Filter and Control Modifier (FCM):** A rule-driven engine that performs semantic transformations to reduce the content value of data and restrict control actions.
- **Format Transcoders (FT):** These components perform structural transformations to convert data flowing to and from services or devices into a representation format that the FCM component can operate upon.
- **Identity Service (IS):** The IS a secure central repository that stores persistent state on behalf of users. The state includes identities for accessing services and rules for displaying and modifying service content based upon the users' role and mode of access. Additionally, the IS stores identifying information for Internet services. This information consist of scripts for scraping the content from the service, and style sheets for rendering the content to device formats.
- **Transient Store (TS):** The TS stores temporary state, such as current session identifiers, mappings for obfuscated or hidden data, Uniform Resource Locators (URLs), and client request information.

Figure 1 illustrates the architecture of our proxy. Each component unifies a collection of devices or protocols into one format, isolating service-specific details from the logical operations required to transform data. This abstraction capability becomes increasingly important as the number of devices and services grows. A new device or service can be enabled by simply creating a few device-specific components, easing composability versus modern vertical web-based applications that require service authors to individually tailor services to each end device. Our proxy simplifies the delivery of content to multiple devices by: reducing it to the task of authoring style sheets to render content, and delivering the content through building block components that are implemented once to support a specified security

or Internet protocol. Without these abstractions, authoring is an $N \times M$ task, where the introduction of a new device type requires changes to every service to add support for the device and vice versa (e.g., web-clipping [1]).

2.2 Security Adaptors (SA)

Security Adaptors allow devices capable of performing one security protocol to access services that require a different protocol. Consider a mobile user with a low-power device capable of shared-key authentication and Blowfish [34] encryption, but not capable of performing the complex computations required by SSL [24]. However, a service may only support authenticated access over SSL. In this case, the client-side Security Adaptor is a shared-key adaptor that obtains a user's key from the Identity Service and uses it to perform challenge-response authentication and to encrypt all communication with the client. The server side Security Adaptor communicates with the end service over SSL. Security Adaptors must be trusted as they have access to unencrypted user and service data and authenticate users, services, and devices to the infrastructure.

When connecting from an untrusted public key encryption-capable device, users authenticate themselves to the infrastructure using a pseudonym identity that identifies the trust profile and access device, and a one-time use password mechanism (e.g., S/KEY or SecurID [29, 19, 13]). This information does not need to be hidden from the untrusted endpoint as it can only be used once.

Additionally, the client-side Security Adaptor can provide persistent connections, even if the service does not provide them, enabling a client to establish a secure connection to the infrastructure once, while the infrastructure (potentially) establishes multiple connections to the service. If the connection setup time or round trip latency is high, this technique will yield a substantial performance improvement. Also, the infrastructure can establish multiple connections to different services to aggregate information (e.g., to perform a meta-search of multiple web sites).

2.3 Format Transcoders (FT)

Format Transcoders (FTs) transform data between external service and device formats and the FCM's XML-based language, abstracting service- and device-specific information from the FCM. There are two types of Format Transcoders: server-side FTs and client-side FTs.

2.3.1 Server-Side Format Transcoders

The server-side Format Transcoder receives a service's content and transcodes it into an XML [8] representation of the semantic data. We choose XML as it is the industry standard technology for semantic mark up. Additionally, the

XSL [15] presentation standard, which separates content from presentation, allows us to render the XML representation to any device format. Our architecture allows content authors to choose the XML representation of their choice, enabling authors to provide multi-modal access to services without waiting for service providers to agree on a standard XML DTD for their content.

The proxy extracts a service's semantic representation using "screen scraping"¹. Scraping is a service-specific operation, so each site requires its own content extraction scripts. The server-side Format Transcoder is a generic execution environment for these scripts.

By abstracting the logic for scraping service content into scripts, we create a single place where all service-specific code is located and enable the FT component to be generically reused across multiple services. By semantically tagging data the FCM can apply rules, based on the user's current trust profile, to these tags and appropriately modify the data.

The author of the server-side Format Transcoder scraping script must identify the semantic content of displayed data and the actions that can be performed by a service. For web-based services, this is a substantial task that requires manual analysis. Tools such as WebL [12] help to automate this task. Products such as Epicentric portal server [14] and Cohera [10] include web-scraping building blocks.

When operating on data flowing towards a service, the format transformation merely transforms the XML representations of control actions (e.g., HTTP POST or GET requests for web pages, or RMI for Java based services) into the services' target format. The code to perform these operations is reusable across multiple services.

2.3.2 Client-Side Format Transcoders

Client-side Format Transcoders perform transformations on data flowing to and from the client into and out of the FCM's XML representation. As data travels to a service, HTTP POST and GET requests, RMI calls, etc. are transformed by the client-side Format Transcoder into an XML representation that can be operated on by the FCM to apply control filtering rules.

Consider, for example, the buy request that a user would execute in the stock trading scenario. The original HTTP POST request is translated into a representation where each form item is given a semantic XML tag. Additionally, the intermediate XML format allows HTTP based clients to access RMI based services and vice versa.

For data traveling to the client, the client-side Format Transcoder uses XSL style sheets to render the the FCM-modified XML for the particular end device. Using XSL

¹Extracting information from a web page's HTML code.

gives content authors a standard tool to tailor the presentation to the device format. For the stock trading example, the format for the kiosk is HTML, while the format for the PDA could be WML, a thin client browser format, or another custom format.

2.4 Filter and Control Modifier (FCM)

The FCM implements application-level logic for adapting data flowing between devices and services to the desired security parameters by applying the rules associated with users' current trust profiles to perform content and control filtering on the XML representation of the semantic content. The FCM also applies rules for split-trust control functionality, where a trusted portable device is used to authorize and authenticate actions from an untrusted device.

2.4.1 Using content filters to obfuscate data

Infrastructure content filters alter or remove sensitive service content before it reaches the untrusted device, effectively decreasing the data's privacy or security level.

The FCM applies the user's rules by matching each XML tag to the appropriate rule, and performing the specified action. In order for a data item on the page to be displayed, the security level at which it can be rendered must be explicitly specified by the user's trust profile. This requirement follows the principle of least privilege.

There are six types of rules for data value reduction: *allow-content*, *hide*, *obfuscate-well-known*, *obfuscate-mapping*, *obfuscate-form*, and *obfuscate-cookie*. The *allow-content* rule explicitly states that a content item is allowed to be rendered on the page in its current format. This rule prevents service format changes from causing secure information to be accidentally released. The *hide* rule replaces sensitive information in the content with an uncorrelated number of “[” characters, within a reasonable limit of the original size of the data. The *obfuscate-well-known* rule replaces well-known sensitive information with its description rather than its value. For example, rather than displaying the user's home address “777 Main Street”, the text “HOME ADDRESS” is displayed. This information is meaningful to the end user, but leaks no useful information to the kiosk.

The *obfuscate-mapping* rule produces a random table of mappings, stores the table in the Transient Store, and replaces sensitive data with an obscured name from the table. If a trusted device is being used in tandem with an untrusted one, the FCM sends the mapping table to the user's trusted device, so the user can make sense of the obfuscated data. In the stock trading example, to prevent the kiosk from obtaining the list of stocks contained in a portfolio, the FCM applies the *obfuscate-mapping* rule to the portfolio data. For example, the result could be a table that maps stock

names to the names of states (e.g., “IBM” maps to “Kentucky”, “AOL” to “Michigan”, etc.).

A rule that is similar to the *obfuscate-mapping* rule is the *verify* rule, which leverages the trusted mobile device to verify content displayed on the untrusted endpoint. The *verify* rule can be used to make sure that the untrusted device does not modify useful public information in a malicious manner. For example, the price of an item on an on-line shopping web site is of dubious value if it can be modified by an untrusted endpoint. The *verify* rule uses the portable trusted device to display small amounts of data, thereby allowing the user to determine whether the untrusted kiosk has tampered with the data. The *verify* rule is particularly attractive for mobile devices with cumbersome interfaces as it allows the majority of the interaction to occur on the rich interface of the untrusted public device.

The *obfuscate-form* rule removes sensitive information encoded in form fields and hyper-links by generating temporary obfuscated names for form action names, form input fields, and hidden fields for state. Mappings between obfuscated values and actual data are held in the Transient Store. The rule also modifies hyper-links and Java Script in HTML pages containing sensitive session or user data. These transformations are necessary to prevent the leakage of sensitive data to kiosks, as well as to prevent kiosks from being able to post form requests directly to the end service. The *obfuscate-cookie* rule performs a similar transformation for any cookies that may be sent to the untrusted endpoint.

2.4.2 Using content rewriting to reduce the entry of sensitive information

Just as we do not trust endpoints to receive sensitive data, we do not want users to enter sensitive information into untrusted kiosks. Thus, the FCM uses replacement rules to exchange placeholder information with sensitive information that is stored only in the trusted infrastructure. There are four rules for content rewriting: *replace-well-known*, *replace-mapping*, *replace-form*, and *replace-cookie*.

The *replace-well-known* rule replaces well-known identifiers with the actual data, such as “CREDIT CARD”, with the user's credit card number, and “HOME ADDRESS” with the user's real home address. This type of form replacement is useful for supporting e-commerce purchase requests that need account numbers and shipping addresses that can't be entered on untrusted clients. This functionality is also useful for trusted clients as it simplifies the process of filling out forms.

The *replace-mapping* rules replace obfuscated values with mapped values using information stored in the Transient Store by *obfuscate-mapping* rules (e.g., a user selling a stock enters a U.S. state abbreviation rather than the stock symbol, and the FCM performs the reverse mapping).

The *replace-form* rule is used to replace obfuscated form action names, field names, and hidden text in HTTP GET/POST requests using mappings previously stored in the Transient Store. The *replace-cookies* rule works similarly for cookies.

2.4.3 Using control filtering to protect service functionality

Many services provide information with varying degrees of sensitivity (e.g., looking up a stock quote is much less sensitive than trading a stock), however, they provide all or nothing access to service functionality, as they assume that the endpoint is trusted. The proxy infrastructure supports access from untrusted endpoints by providing fine-grain access control over service functionality. The FCM provides this capability using control filtering rules that control data and commands flowing from the client to the server.

The FCM uses two rules for control filtering: *allow-control* and *authorize*. The first rule states whether or not a request to submit a form or view a page can be transmitted to the end service. In the stock trading example, the control functions on the service's forms and pages (even in the split-trust case) must be explicitly allowed by the profile.

The *authorize* rule provides fine-grain control over sensitive actions, such as buying stock from an untrusted endpoint by using an a trusted device to explicitly authorize the request. In the stock trading example, a request to trade stocks is intercepted and held by the FCM until authorization is received from the user's trusted companion device (e.g., a PDA, cellular telephone, or two-way pager). If a user wants actions to be authorized through a device, they need to first register the authorization device with the infrastructure. Then, when an action is performed that requires confirmation, the FCM will intercept and hold it, and contact the user via the specified device. The trade action is then sent to the device and confirmation is requested. Once it is authorized, the FCM releases the request.

2.5 Identity Service (IS)

The Identity Service is the secure repository of persistent data associated with users' *trust profiles* and their security preferences for access to services. The Identity Service stores the following types of information:

- Identities and credentials (e.g., usernames, passwords, encryption keys, etc.) for accessing secure services.
- Rules and preferences for content and control filtering.
- Personal information (e.g., addresses and account numbers) for automated form filling from untrusted devices.

- Pseudonym identities and one-time passwords.
- Service information such as content scraping scripts and XSL style-sheets.

The primary motivation for the Identity Service derives from the need to hide private data from untrusted endpoints, where there is no way for a user to securely transmit private information. The Identity Service provides an alternative where users' passwords and other private information are stored in the infrastructure and merged with any transmission after the data has left the untrusted device and entered the trusted infrastructure. One could imagine a world in which identity services become a common piece of the public Internet infrastructure. However, until existing Internet services migrate to adopt this functionality our proxy can transparently provide it.

The Identity Service provides additional benefits for users accessing secure services from trusted endpoints. The replication of security credentials on several devices, some of which may easily be stolen (e.g., small PDAs), increases the probability that credentials may be compromised. It is also unlikely that most users will be savvy enough to understand the security capabilities of several devices, operating systems, and file systems. Similarly, changes to a credential require that it be updated on all devices. The IS solves these problems by providing a single source for such data.

Obviously, the security model of the Identity Service must be carefully considered and it must be well-protected against attacks. We argue that storing the data in one professionally-administered, physically-secured Identity Service is more secure than having it managed across several physically insecure devices by naive end users. Users can control access to their private information in the Identity Service as each of the data entries has an associated Access Control List (ACL).

2.6 Transient Store (TS)

The Transient Store is used to store soft state associated with the current session of the user, such as mappings of obfuscated to original URLs, pending requests for trusted device authorization, or session information. Data can be stored using index keys that are generated by applying a secure hash function (e.g., SHA1 [35]). Secure access can be provided by using either a large key space such that a key can never be guessed or by performing access control.

3 Implementation

The design we propose in this paper is a generalization of the ideas we gained from an initial proof of concept implementation of the system. The proof of concept implementation provided access to the Datek Online Stock Trading

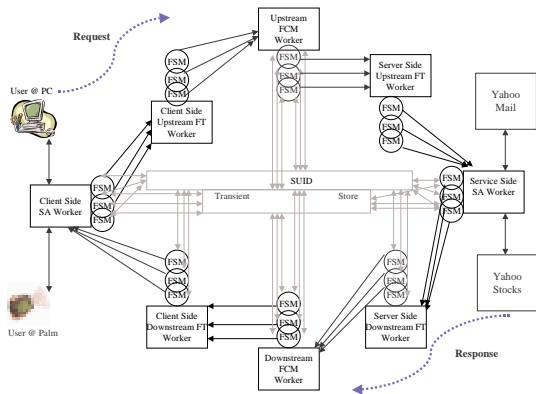


Figure 2. Transformation pipeline

service for users of untrusted kiosks and computationally limited PDAs. The lessons learned in the proof of concept implementation were valuable and led to the design decisions for the current implementation, based upon the Ninja vSpace Platform.

The vSpace platform is version two of the Ninja infrastructure for building scalable, highly available, fault tolerant, cluster based applications. Among the most notable features of vSpace are asynchronous RPC, distributed data structures, and an event driven programming model. Additionally, the vSpace execution environment provides support for load-balancing, and fail-over.

Our proxy service is implemented as a pipeline of ten vSpace Workers composed through the Task/Completion asynchronous Ninja RPC mechanism as depicted in figure 2. The vSpace platform is essential for high concurrency performance of the composable pipeline. In thread based RMI environments, blocking interfaces break down as threads are tied up on the callers. For example, when object A remotely calls object B that remotely calls object C the thread on A has to wait until the call to C finishes. In the asynchronous task based vSpace environment, the thread becomes freed immediately after object A remotely calls object B. The vSpace platform supports high concurrency with fewer threads and hence less computational overhead due to thread overhead and context switching.

3.1 Current worker implementations

In the current implementation, Client Side Security Adaptors (CSSA) listen on well known ports using asynchronous TCP server sockets² modified to support security protocols. Currently, SSL is supported, using Sun's JSSE library, as well as a shared key protocol based on Blowfish. SSSA session objects handle HTTP redirections directly to

²provided by the vSpace platform

avoid costly round trips through the potentially thin network to the client device. SSSA Session objects also support cookies, storing them in the Transient Store.

The Format Transcoder implementation consists of four vSpace workers supporting HTTP based services. The Client Side Upstream Format Transcoder (CSUFT) uses a simple parser to transform HTTP requests into an XML representation. The Service Side Upstream Format Transcoder (SSUFT) uses Sun's SAXP XML parser to transform the XML representation into the appropriate HTTP request for accessing the end service. The Service Side Downstream Format Transcoder (SSDFT) uses WebL [12] to scrape the content from HTML pages and return it in an semantic XML representation. A drawback in this implementation is that a new WebL process is exec'd for each page. The Client Side Downstream Format Transcoder (CSDFT) uses Apache's Xalan [3] XSL renderer.

The Filter and Control Modifier workers use Sun's SAXP XML parser and apply rules in one pass. Currently, the following subset of rules is implemented: hide, allow-content, allow-control, and replace-well-known. Rules are specified per profile and can be set by either a text file or through the JSP based UI.

The Identity Service (SUID) is implemented as an API layer on top of a distributed hash-table [21] and works across a cluster. The SUID Worker provides a Ninja RPC interface for the specific tasks required by the proxy service to manage user profile information, service transcoding rules, and UI functionality. Since the state transitions in the SUID are relatively simple, rather than creating a new FSM object for each request, it manages requests internally. The Transient Store is implemented in a similar fashion.

3.2 Discussion

Our initial design changed little during the implementation phase. The most notable change is using the WebL scripting language to create a generic Server Side Downstream FT component, a capability that requires the storage of service-specific information, such as WebL scraping scripts and XSL style sheets, in the Identity Service. A future change would be to separate out the XSL rendering and web scraping functionality into individual workers, rather than having the functionality embedded in the FT worker state machines. The further decomposition will create useful building blocks for other Ninja services. It will also make it easy to replace the workers with alternative web scraping and XSL tools.

Another interesting lesson learned in the current implementation is the need to provide pluggable protocol parsing in the security adaptors. We found the CSSAWapSession to be a useful place to do pre-processing of parsing of headers for protocols and gathering data packets.

Symbol	Last Trade	Chg	Volume	Shrs	Value	Value Change	Paid	Gain	More Info
SSCASH			75,916.25		\$75,916.25	-	-	-	N/A
MSFT	Sep 22 63 1/4	15 1/16	42,238,300	100	\$6,325.00	-\$93.75 -1.46%	62.3125	\$83.75 +1.34%	Chart, News, Mags, Profile, Research, Insider, Options
AOI	Sep 22 55 1/4	+1 3/16	10,005,800	100	\$5,525.00	\$175.00 +3.27%	54.3125	\$83.75 +1.54%	Chart, News, Mags, Profile, Research, Insider, Options
MSFT	Sep 22 63 1/4	15 1/16	42,238,300	100	\$6,325.00	-\$93.75 -1.46%	62.3125	\$83.75 +1.34%	Chart, News, Mags, Profile, Research, Insider, Options
INKT	Sep 22 125	15 1/16	1,292,200	50	\$6,250.00	-\$46.88 -0.74%	123	\$90.00 +1.46%	Chart, News, Mags, Profile, Research, Insider, Options
5 symbols					Totals(USD):	\$100,341.25	-\$59.38	-0.06%	\$341.25 +1.42%

Figure 3. Yahoo Contest holdings

symbol	last trade	change	volume	shares	value	value change	paid	gain	more info
SSCASH									N/A
MSFT	Sep 22 63 1/4	15 1/16	42,238,300	100	\$6,325.00	-\$93.75 -1.46%	62.3125	\$83.75 +1.34%	Chart, News, Mags, Profile, Research, Insider, Options
AOI	Sep 22 55 1/4	+1 3/16	10,005,800	100	\$5,525.00	+\$175.00 +3.27%	54.3125	\$83.75 +1.54%	Chart, News, Mags, Profile, Research, Insider, Options
MSFT	Sep 22 63 1/4	15 1/16	42,238,300	100	\$6,325.00	-\$93.75 -1.46%	62.3125	\$83.75 +1.34%	Chart, News, Mags, Profile, Research, Insider, Options
INKT	Sep 22 125	15 1/16	1,292,200	50	\$6,250.00	-\$46.88 -0.74%	123	\$90.00 +1.46%	Chart, News, Mags, Profile, Research, Insider, Options
5 symbols					Totals(USD):	\$100,341.25	-\$59.38	-0.06%	\$341.25 +1.42%

Figure 4. Yahoo Contest holdings filtered

Overall, the current implementation consists of approximately 26,000 lines of commented Java source. The management user interface is 2,100 lines of JSP code. Yahoo Contest and Yahoo Mail combined are 1,400 and 1,600 lines of XSL and WebL code respectively.

4 Analysis

In this section, we evaluate the effectiveness of our proxy-based approach to secure multi-modal access to Internet content from post-pc devices. The evaluation criteria are: ease of adding support for new services, ease of supporting new client formats, and overall performance of the system. As test cases, we use a service that provides secure access to Yahoo Contest [22], a stock trading simulator.

4.1 Adding new services

Adding support for a new service merely requires writing a site map, WebL scripts for each of the services pages, and a default rule set. Our implementation of proxied multi-modal access to Yahoo Contest service consists of five content pages: holdings, orderConfirmation, orderForm, orderVerification, and quotes. We also use a shared page that we include in each of these pages. The site map consists of the service name and unique page identifiers for each page.

A service author creates WebL scripts for each page and stores them in the SUID. The holdings page requires the most scraping and consists of a 250 line WebL script. Figure 3 displays the original Yahoo Contest holdings page. Writing the Yahoo Contest scraping scripts took three weeks of part time effort in a pass-fail course by an undergraduate unfamiliar with WebL. After gaining familiarity with WebL, authoring service scraping scripts became a simple task that was done in a few days.

The following example illustrates the default rule set for the tags found on the holdings page of Yahoo Contest. Only

the rules for hiding data are shown as the rest of the rules for this page are allow-content.

```
hide <Username>
hide <Shares>
hide <Value>
hide <Paid>
hide <TotalValue>
hide <TotalPaid>
```

Additional rules are needed for the Upstream FCM to determine which control actions to allow by default in a given profile. In this case, the actions map to URLs in the service.

```
allow-control <action> "http://finance.yahoo.com/p"
allow-control <action> "http://contest.finance.../t2"
allow-control <action> "http://contest.finance.../t3"
allow-control <action> "http://finance.yahoo.com/q"
```

4.2 Adding support for new device formats

Adding support for a new client device requires writing an XSL style sheet to render the content for that device. Security adaptors need to be implemented, if they do not already exist for the desired security protocol. However, since security adaptation and content adaptation functionality is separated in our design, security adaptors can be re-used and only need to be written once per new client format.

Figure 4 shows the modified version of the holdings page after the security transformation rules have been applied. Figure 5 shows the corresponding WML version of the holdings page. The figures illustrate multiple device fusion for service access, as sensitive information hidden from the kiosk is displayed on the trusted PDA.

Two XSL style sheets are used to render the holdings page. The first generates an HTML representation for use by either a secure home machine, or a web-based kiosk. The second generates a WML representation that can be viewed on mobile devices with WML browsers. The HTML style sheet was easy to author based on the site's original HTML presentation. The WML style sheet is a simple adaptation of the HTML style sheet. For the holdings page, the HTML and WML style sheets are 370 and 100 lines respectively, plus a shared file of 100 lines.

As illustrated in the examples, adding support for an additional device format only requires writing a short XSL style sheet to output the content in a suitable format for that device. This powerful capability allows content authors using the proxy to tailor the content for device screen layout and other factors. Additionally, the semantically tagged representation enables powerful transformations such as meaningful text-to-speech synthesis of the XML content.

However, these benefits come at a cost, with a trade-off versus the automated approach taken by other proxies such as ProxiNet for delivering content to PDAs. The ProxiNet approach does not require any style sheets and performs generic transformations of HTML into its custom thin client

browser format. Thus, instant access to any web-site is provided to PDAs that run its browser. However, generic approaches do not offer the ability for content authors to customize pages for small devices, or do transformations such as voice rendering. To enable access to more services, we could provide the functionality of the ProxiNet approach by bypassing the semantic transformation layer. However, this would disable the fine-grain access control, content-filtering, and multiple device fusion capabilities.

4.3 Performance

Performance was measured in two configurations. The first consisted of all the workers running on a single node. The second consisted of the workers distributed across the cluster, each node running one worker.

Both of these experiments were performed on the UC-Berkeley Millennium cluster [6]. Each cluster node has two 500 MHz Intel Pentium III processors with 512KB cache, 512 MB of RAM, two 9 GB disks, and a 100 Mb/s Ethernet connection. The nodes run Red Hat Linux version 6.0 release 2.2.5. The Java environment is Sun's JDK 1.2.2 production release for Linux running with green threads and the Inprise just in time compiler³.

For the experiments, a client used the proxy to request the holdings and quotes page 20 times each. The cold start run exhibited extremely high latencies (more than an order of magnitude higher than the average times after warmup). These results can be attributed to Java class loading and initial JIT compilation. As a result, the average times only include warm runs. During the experiments, the vSpace workers were run with minimal logging.

The single node configuration returned pages in an average of 3.75 seconds. The holdings and quotes pages round trip times (including authentication and redirections to the Yahoo service) averaged 4.28 s and 3.23 s respectively.

The long latency is caused by inefficiencies in the current implementation. The workers exec new processes for WebL (331.3 ms) and Xalan (99.94 ms). Inefficiencies in the size and serialization of the format transcoder collection also contributes to this time due to the slow read from disk and repeated serializations. The disk writes also affected the storage components for temporary state as the DDS uses two phase commit. A large portion, 34.82%, of the round trip time is not accounted for by these measurements. Factors that contribute to this time are task dispatch Java serialization time, queue wait time, and debug output. An optimized implementation should be able to make a substantial improvement in the latency of the proxy.

The distributed configuration performed in a similar manner. Pages were returned in 3.87 s on average. The round trip times for the holdings and quotes pages were 4.18

³The current Java environment only supports a single processor.

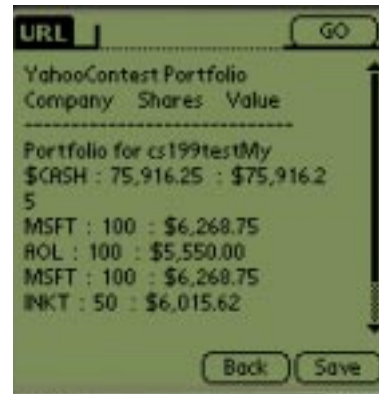


Figure 5. Yahoo Contest holdings WML

s and 3.55 s respectively. We expected the distributed version to outperform the single node as we observed the Java based proxy to be memory and CPU intensive during development on smaller workstations. However, the resources for the single millennium node proved adequate for this benchmark. The proxy was clearly I/O bound in this environment.

The distributed configuration displayed considerable tolerance to faults. Workers sometimes caused the JVM running the vSpace to crash. When the JVM was restarted, the system recovered gracefully with a higher latency for the initial class loading and compiling of the restarted worker. Future versions of vSpace, which will include automatic restart of workers, can take advantage of the properties of this configuration.

At the time of this writing, the current implementation is not stable enough to perform scalability experiments. Scalability results will be obtained after further optimization and cleanup of our proxy service, and the vSpace platform. We hypothesize the distributed configuration will have a higher throughput than the single node configuration.

4.4 Discussion

Overall, our proxy approach allows considerable customization capabilities and easy scaling to multiple devices and services. The development effort required is minimal: the simple authoring of WebL scripts for content scraping and XSL style sheets for device specific rendering. This development work can be done by an independent service provider running the proxy, and specializing in this work, rather than waiting for end services to provide support for multiple devices. Alternatively, content providers can provide content in an XML representation to bypass web-scraping and enable the content-filtering and digital wallet capabilities provided by the proxy.

In our experience, the time to author components for new protocols and style sheets for new devices is minimal. The

CSSAWapSession and WML pages for YahooContest were written in a few days on top of the existing implementation. This included the time to learn WML and discover the differences between the WML browser protocol and HTTP. Additionally, scraping the web site and generating a XML representation is also minimal. Support for the Yahoo Mail service was added in a few days. The user interface for customization of rules for user profiles should make it easy for users to choose a content and control filtering level they feel comfortable with from their mode of access.

The current implementation was developed on a prototype version of the Ninja vSpace platform. The majority of the development effort to get the current version running consisted of understanding the vSpace programming model and flushing out bugs in the prototype platform. Without time for optimization, the initial performance numbers of a 3 second latency are promising. Our experience with fault recovery in the distributed configuration also shows promise. Now that the platform is stabilizing, we should be able to significantly improve our performance, stability, and scalability using the understanding we gained in building the current implementation.

An open issue is the psychological acceptability of our approach. Mobile commerce is just starting to hit the public. Will mobile users be satisfied with the cumbersome interfaces of small devices or trusting of public terminals in the environment? One may argue that people aren't as concerned enough with privacy (or aware enough) as they are willing to use shopper affinity cards, revealing detailed purchase histories to supermarkets in exchange for saving a few dollars. However, we feel these same users will be sensitive and concerned when access to credit accounts and other personal information required to complete a commerce transaction are at stake. We also feel that these sensitive users will find the act of connecting to the proxy service and authenticating with one-time passwords an acceptable part of their mobile commerce routine.

5 Related work

Several related projects have focused on data filtering at different levels of granularity for increased security. Wiederhold and Billelo [38] examine the problem of protecting the release of data from databases for collaboration purposes, where the internal databases are not organized according to external access criteria. The authors note that, in these environments, traditional authentication and authorization tools and secure transmission fail to protect the release of inappropriate data. They propose a solution based upon a rule-driven security-mediator that performs result-checking to filter both queries and, more importantly, their results. Their rules for filtering query results contain traditional role-based access control, site blocking, and more

advanced features, such as the replacement of query result components with non-identifying (obfuscated) terms.

Our architecture's approach to content and control filtering is based upon having specific semantic knowledge of the service data and control actions that are being filtered or modified. However, such an approach does not work for protecting unstructured information, such as the content of e-mail messages. Instead, we believe that we could apply techniques from medical research that focus on generalizing, substituting, and removing information from medical data to protect privacy without obscuring "important" details of the data [36].

Related to our idea of multiple device fusion is the notion of *splitting trust* between PCs and mobile devices to increase security [5]. In the split trust paradigm, an application's security crucial parts execute on a small trusted device, while other parts execute on a more powerful, but untrusted device. The authors explore the split trust paradigm in the context of e-mail using a Palm Pilot as a PKCS11 smart card for Netscape Communicator. Our architecture employs a similar technique in the network using out-of-band channels between the trusted and untrusted devices.

Another related project is the work by the World Wide Web Consortium on the Platform for Privacy Preferences (P3P) [18] and the associated preference exchange language – APPEL [27]. P3P provides a standard means for service providers to disclose their practices regarding the collection and dissemination of personal information and allows users to make informed decisions regarding the use of this information. This research is primarily focused on protecting the information sent from users to services, varying from click-stream data to users' personal information, such as names and addresses. We explore the inverse of this problem: allowing users to express their preferences about the type of data to be displayed and the functionality to be allowed as a function of the users' access device. Also related to this research area are trust-management systems, such as KeyNote [7], that serve as engines for applications to determine whether or not a potentially dangerous operation conforms to a user's security policy.

Our infrastructure-based proxy model is also related to research on proxies to more efficiently use network resources, reduce cost, and increase security focusing on filtering to optimize performance along the last link in mobile environments [2, 40, 41]. Similar to these works, we stress the ease of deploying an intermediary instead of changing the protocols at end points and place considerable trust in the proxy.

Related systems have examined the problem of managing different types of identity data; for example, Apple's KeyChain [4] stores all username-password pairs locally on the client PC. Lucent's web proxy, ProxyMate [28], takes this capability a step further by generating stronger user

ID's and passwords transparently to the user, and using them to login to the actual site. Novell's DigitalMe [31] and Yodlee [39] store passwords on a central server⁴. Microsoft's Passport [30] service as well as other e-wallet services hold personally identifying information to simplify form filling in e-commerce transactions. This functionality is similar to our architecture; however, we make such information always available by keeping it in the infrastructure, rather than on a user's end device. CommerceNet's IdentitySafe [11] provides users with pseudo identities for purchasing from online merchants and only releases personally identifiable information to sites on a need to know basis.

Our application of security protocol adaptation is derived from the ideas presented in [23], which suggested that security transformations for small devices would be desirable and presented proxy-based access to Kerberized services. Their design places the bulk of Kerberos authentication work on an infrastructure-based proxy, providing indirect authentication of clients and simplifying the requirements of small devices.

Finally, to format content for a range of devices with limited display capabilities, we could leverage related research, such as TACC [16, 17]. By abstracting device functionality into the client-side Format Transcoder, we simplify the process of adding new types of devices. In addition, the Wireless Application Protocol(WAP) [37] presents a standard format for small devices. Other related work in this area includes Oracle's Portal-To-Go which promises any service to any device access [32].

6 Conclusion

In this paper, we present a novel infrastructure-based architecture for providing secure multi-modal access to wide-area services. The architecture uses rule-driven content and security transformation functions to enable access from a wide variety of end devices by decoupling device capabilities from service requirements. This approach greatly simplifies and reduces the amount of work required to support a new device or service.

Providing access to a new service merely requires writing a WebL script to scrape the content into an XML representation. Our example stock service consists of 6 scripts totaling 1000 lines. After gaining familiarity with WebL, support for the Yahoo Mail service was added in a few days. Yahoo Mail support consists of 5 scripts totaling 610 lines.

Support for new device formats is easily provided by writing XSL style sheets. For HTML, the most complex

⁴These systems provide convenience, but do not enhance security from untrusted endpoints, as the master password still needs to be entered and thus, can be captured or hijacked by the untrusted endpoint. In contrast, our architecture's Identity Service provides one-time passwords and fine-grain control over security.

format, there are 1100 lines of code total in all of the style sheets for Yahoo Contest. Yahoo Mail requires 360 lines respectively. Adding support for the WML browser and style sheets for WML pages only required a couple of days effort. The custom WML version of the holdings page was adapted from the HTML version and is only 100 lines compared to the 370 used for HTML tables. Although some development is required, this approach allows considerable customization capabilities. Furthermore, it provides the ability to support interesting formats such as custom WML and voice access which have yet to be seen from generic HTML transformation proxies.

The current implementation returns pages in an average time of approximately 3-4 seconds. The implementation has not yet been optimized and there are many opportunities for improvement. Implementing the system as a vSpace service required a steep learning curve initially. However, the cluster wide management of persistent state provided by the DDS, and auto-discovery of workers are very valuable building blocks. Experimenting with vSpace v1 workers in the distributed configuration showed promising fault tolerance results. We are looking forward to the added features provided by vSpace v2.

We provide a generic "any device to any service" model that is in stark contrast to the traditional approaches of vertical service and device integration and implementation. Our hypothesis is that such a generic model provides rapid support for new devices and services, a critical requirement for the Post-PC era.

Furthermore, the architecture also supports a new model of secure interaction from untrusted public Internet access points found in mobile environments. To our knowledge, our architecture is the first to address this interaction model. By providing a generic control and content rewriting capability, the architecture provides users with precise control over the exposure of information, both as a function of their access devices and as a function of the users' adopted roles. Additionally, the architecture supports the fusion of multiple devices, using trusted portable devices for secure authorization of sensitive requests. Further non-security related applications of multiple device fusion could allow voice access to a service via the car, displaying graphical information on a passenger's PDA.

References

- [1] 3Com. Web Clipping Applications Tutorials. <http://www.palm.com/devzone/webclipping>, January 2000.
- [2] E. Amir, S. McCanne, and H. Zhang. An Application Level Video Gateway. In *Proc. ACM Multimedia 95*, San Francisco, CA, 1995.

- [3] Apache. Apache XML Project. <http://www.apache.org>, March 2000.
- [4] Apple. Keychain. <http://www.apple.com/macos/feature4.htm>, January 2000.
- [5] D. Balfanz and E. Felten. Hand-held computers can be better smart cards. In *Proc. of the Eighth USENIX Security Symposium*, Berkeley, CA, August 1999.
- [6] UC Berkeley. Millennium. <http://www.millennium.berkeley.edu/>, 2000.
- [7] M. Blaze, J. Feigenbaum, J. Ioannidis, and A. Keromytis. RFC 2704: The KeyNote Trust-Management System Version 2. <http://www.cryptocom.com/papers/rfc2704.txt>, September 1999.
- [8] T. Bray, J. Paoli, and C. M. Sperberg-McQueen. Extensible Markup Language (XML) 1.0. <http://www.w3.org/TR/REC-xml>, 2000.
- [9] Certicom. Elliptic Curve Cryptography for Palm VII. <http://www.certicom.com/press/98/dec0298.htm>, December 1998.
- [10] Cohera. Cohera. <http://www.cohera.com>, 2000.
- [11] CommerceNet. IdentitySafe. <http://www.commerce.net/project/hotsheet.html>, January 2000.
- [12] Compaq. WebL User Manual. <http://www.compaq.com/WebL>, March 2000.
- [13] Security Dynamics. SecurID. <http://www.rsasecurity.com>, January 2000.
- [14] Epicentric. Epicentric Portal Server 3.0 Datasheet. <http://www.epicentric.com>, 2000.
- [15] A. Adler et. al. Extensible Stylesheet Language (XSL). <http://www.w3.org/TR/xsl/>, 2000.
- [16] A. Fox et. al. Adapting to client variability via on-demand dynamic distillation. In *Proc. of the 7th ACM Inter. Conference on Architectural support for Programming Languages and Operating Systems*, Cambridge, Massachusetts, October 1996.
- [17] A. Fox et. al. Scalable network services. In *Proc. of the 16th ACM Symp. on Operating Systems Principals (SOSP-16)*, St. Malo, France, October 1997.
- [18] L. Cranor et. al. Platform for Privacy Preferences (P3P1.0) Specification. <http://www.w3.org/TR/P3P>, 2000.
- [19] N. Maller et. al. A One-Time Password System. <http://www.ietf.org/rfc/rfc2289.txt>, February 1998.
- [20] S. Gribble et. al. The MultiSpace: an Evolutionary Platform for Infrastructural Services. In *Proc. of the 1999 Usenix Technical Conference*, 1999.
- [21] S. Gribble et. al. Scalable, distributed data structures for internet service construction. In *Proc. of the Fourth Symp. on Operating Systems Design and Implementation, OSDI*, October 2000.
- [22] Yahoo Finance. Yahoo Finance Investment Challenge. <http://contest.finance.yahoo.com/t1?u/>, 2000.
- [23] A. Fox and S. Gribble. Security on the Move: Indirect Authentication using Kerberos. In *Proc. of the 2nd ACM Inter. Conference on Mobile Computing and Networking*, Rye, New York, November 1996.
- [24] A. Frier, P. Karlton, and P. Kocher. The SSL 3.0 Protocol. <http://www.netscape.com/eng/ssl3/ssl-toc.html>, March 1996.
- [25] InfoWorld. Boeing to put Net in the air. <http://www.infoworld.com/articles/hn/xml/00/04/27/000427enboeing.xml>, April 2000.
- [26] InfoWorld. E-cars take to the streets; wireless connections link road warriors to the Net. <http://www.infoworld.com/articles/hn/xml/00/03/13/000313hnauto.xml>, March 2000.
- [27] M. Langheinrich. A P3P Preference Exchange Language (APPEL) Working Draft. <http://www.w3.org/TR/WD-P3P-preferences>, 1998.
- [28] Lucent. Proxymate. <http://www.proxymate.com/>, January 2000.
- [29] N. Maller. The S/KEY One-Time Password System. <http://www.ietf.org/rfc/rfc1760.txt>, February 1995.
- [30] Microsoft. Passport. <http://www.passport.com>, January 2000.
- [31] Novell. Digitalme. <http://www.digitalme.com/>, January 2000.
- [32] Oracle. Oracle Portal-to-Go Any Service to Any Device. <http://www.oracle.com/mobile/panbwp.pdf>, October 1999.
- [33] Palm.com. Palm V PDA. <http://www.palm.com/products/vseries.html>, 1999.
- [34] B. Schneier. Description of a New Variable-Length Key, 64-Bit Block Cipher (Blowfish). In *Fast Software Encryption, Cambridge Security Workshop Proceedings*, December 1993.
- [35] B. Schneier. *Applied Cryptography*. John Wiley and Sons, Inc., 1996.
- [36] L. Sweeney. Guaranteeing anonymity when sharing medical data, the Datafly system. In *Proc. of the American Medical Informatics Assoc. Symposium.*, Washington, DC, August 1998.
- [37] WAP. WAP Forum Specifications. <http://www.wapforum.org/what/technical.htm>, January 2000.
- [38] G. Wiederhold and M. Bilello. Protecting Inappropriate Release of Data from Realistic Databases. In *Proc. of Data and Expert Systems (DEXA) Security Workshop*, August 1998.
- [39] Yodlee. Yodlee. <http://www.yodlee.com/>, January 2000.
- [40] B. Zenel and D. Duchamp. General purpose proxies: solved and unsolved problems. In *Proc. of the Sixth Workshop on Hot Topics in Operating Systems*, Los Alamitos, CA, USA: IEEE Comput. Soc. Press, 1997.
- [41] B. Zenel and D. Duchamp. A general purpose proxy filtering mechanism applied to the mobile environment. In *Proc. of the 3rd ACM/IEEE Conference on Mobile Computing and Networking*, New York, NY USA, 1997.