

01 Chapter 9 02 Remote Usability Evaluation: Discussion 03 of a General Framework and Experiences 04 from Research with a Specific Tool 05 06

07
08 Fabio Paternò and Carmen Santoro
09
10
11
12

13 **Abstract** The goal of this chapter is to present a design space for tools and meth-
14 ods supporting remote usability evaluation of interactive applications. This type of
15 approach is acquiring increasing importance because it allows usability evaluation
16 even when users are in their daily environments. Several techniques have been devel-
17 oped in this area for addressing various types of applications that can be used in
18 different contexts. We discuss them within a unifying framework that can be used to
19 compare the weaknesses and strengths of the various approaches and identify areas
20 that require further research work to exploit all the possibilities opened up by remote
21 evaluation.
22

23 24 9.1 Introduction 25

26 In this chapter, we present and discuss a design space for remote usability evalua-
27 tion. This type of approach to usability evaluation is characterized by the fact that
28 users and evaluators are separated in time and/or space (Hartson, et al. 1996). Thus,
29 it still requires the involvement of these two actors (user and evaluator), but it relaxes
30 the constraint that they need to be present at the same time in the same place. The
31 motivations for remote evaluation are various:
32

- 33 • Usability laboratories can be expensive to set up because they require dedicated
34 sites with specific equipment
- 35 • Moving users to the usability laboratory can be difficult and expensive as well, in
36 particular for expert users, whose time is costly. Indeed, it can be difficult to find
37 an adequate number of users willing to move to a usability lab for a test
- 38 • Remote evaluation can be useful to analyze user behavior in their daily environ-
39 ment (e.g., workplace, home, and so on), thus in more realistic settings
- 40 • It facilitates the possibility of a continuous evaluation, even after the first release
41 of the application
42

43 Some studies have investigated to what extent remote evaluations yield results
44 similar to lab testing. For example, Tullis (Tullis, et al. 2002) found that remote
45 evaluation in the field yielded results that were largely similar to studies in the lab.
46 There are several methods that support some kind of remote usability evaluation.

01 They differ in the type of information that is made available to the evaluator and
02 how it is provided to them. Ivory and Hearst (2001) wrote an interesting review of
03 the state-of-the-art on automating usability evaluation of user interfaces, in which
04 some methods and tools in the area of remote evaluation were considered as well. In
05 this chapter, we provide a more updated and focused discussion of the state-of-the-
06 art in remote evaluation through a more refined framework for this area that high-
07 lights the important aspects to consider when analyzing approaches within it. More
08 specifically, the framework proposed in this chapter is defined by analyzing various
09 dimensions. The first one considers the type of interaction that occurs between the
10 user and the evaluator, and is strongly connected with the possibility of having a
11 co-presence (in terms of time) between the user and the evaluator. Another dimen-
12 sion involves the techniques that can be used to gather information on user sessions
13 (server/proxy/client logs, webcams, eye-trackers, and other sensing technologies)
14 and the information provided, which are useful for the evaluation. Another interest-
15 ing dimension we consider is the type of platform used for interaction. In this regard,
16 we plan to distinguish access through desktop or mobile devices, for example, and
17 discuss how the choice of platform affects the aspects to consider in the evaluation.
18 The last dimension regards the type of application considered (for instance: Java-
19 based, Web-based, etc.). A discussion about the potential correspondences between
20 such dimensions should shed some light on which techniques/technologies evalu-
21 ators should direct their attention to obtain the desired information, therefore pro-
22 viding them with a better understanding of techniques for remote evaluation of user
23 sessions and how to use them to identify problematic parts of interactive applications
24 and make improvements accordingly (when necessary).

25 To summarize, the relevant dimensions we have identified for analyzing the dif-
26 ferent methods for assessing remote usability evaluation are:

- 27
- 28 ● The type of interaction between the user and the evaluator
- 29 ● The platform used for the interaction (desktop, mobile, vocal, etc.)
- 30 ● The techniques used for collecting information about the users and their behavior
- 31 (graphical logs, voice and/or Webcam recordings, eye-tracking, etc.)
- 32 ● The type of application considered in terms of implementation environment
- 33 (Web, java-based, .NET, etc.)
- 34 ● The type of the evaluation results (task performance, emotional state) provided
- 35

36

37 In the next sections, we use the framework composed of such dimensions to
38 discuss a number of techniques that can be used to perform remote evaluation of
39 user sessions (logging technology, interaction platform, semantic analysis), along
40 with a review of the most relevant works in the area together with a discussion
41 about issues that have been resolved from the perspective of usability evaluation
42 and problems that are still open.

43 To make the discussion more concrete, we also discuss our own experiences in
44 this area, including our method (and the related tool) for remote usability evaluation
45 of websites that considers information from user tasks, log files, videos recorded
46 during user tests, and data collected by an eye-tracker (Paternò, et al. 2006).

9.2 The Type of Interaction between the User and the Evaluator

There are different methods and techniques that can be applied to perform a remote evaluation. One important dimension that can be used to classify them is how users and evaluators actually interact between themselves.

Bearing in mind that remote evaluation assumes that users and evaluators are separated in *space*, the type of interaction occurring between users and evaluators strongly depends on the type of synchronization occurring on *time*. Indeed, while the *asynchronous* evaluation method assumes that evaluators might not be necessarily available at the time when the user session is taking place (and therefore, there is no possibility for the evaluator/moderator to deliver immediate input for impromptu changes), it is the opposite with *synchronous* evaluations. As an exemplary case of synchronous evaluation, we mention collaborative usability evaluation methods via the network, in which evaluators in usability labs are connected to remote users via commercially available teleconferencing software (e.g., Microsoft Netmeeting), supporting real-time application sharing, audio links, shared drawing tools, and/or file transfer capabilities. Below, we describe the various possibilities for both synchronous (first bullet) and asynchronous evaluation (second-fourth bullets):

- *Remote Observation*—this implies that users and evaluators are separated in space but are active at the same time and connected through some tool (for example, video conference tools) that allows the evaluator to observe the actual user behavior in real time
- *Remote Questionnaires*—this is a technique that allows users to provide their feedback through a series of questions made available electronically
- *Critical Incidents Reported by the User*—in this case, the user directly reports to the evaluator when an incident occurs
- *Automatic Data Collection*—this is the method that has stimulated the most interest because there are many ways to collect data regarding user behavior and then analyze it. The potential information ranges from browser logs, to videos taken by Web-cams, to eye-tracking data. This case also includes our approach (Paternò, et al. 2006), which will be described in Section 9-7.

To assess the pros and cons of such options, we can notice that on the one hand, remote observation provides the evaluator with more capabilities for observing the session and also for intervening *during* the session. Furthermore, the simultaneous presence of evaluator and user brings the additional advantage of not requiring a particularly strong effort for an *a posteriori* analysis of the collected data, because most of this work should be already carried out by evaluators during the session. On the other hand, remote observation strongly limits the number of users that can be evaluated at a time and, additionally, it might also happen that the behavior of the users might be affected to some extent by their awareness of being currently observed by the evaluator.

Remote questionnaires and critical incidents are useful because they report aspects that the users themselves noticed and judged relevant from the point of

01 view of usability. However, the result of such techniques might be compromised
02 by the fact that the reporting time is generally postponed with respect to when
03 the problem appeared. Therefore, retrospective reporting and questionnaires might
04 be subjected to loss of detail, which can hinder the reconstruction of the original
05 problem.

06 The last technique (automatic data collection) on the one hand guarantees gath-
07 ering a vast amount of detailed data, which, on the other hand, generally claim a
08 non-irrelevant effort and time for being correctly interpreted by humans, in absence
09 of appropriate automatic data analysis techniques.

10
11

12 **9.3 The Platform Used for the Application Interaction**

13
14
15 One of the main characteristics of the rapid evolution of information and communi-
16 cation technologies is the wide availability of various types of interaction platforms.
17 The desktop is no longer the only device that even nonprofessionals use for access-
18 ing their applications. There is a wide variety of interaction platforms on the market,
19 which can largely differentiate in terms of interaction resources (such as screen size,
20 etc.) and supported modalities.

21 Heterogeneous platforms raise specific issues that should not be neglected for
22 the purposes of remote evaluation. For instance, mobile systems are typically used
23 in highly dynamic contexts, and remotely evaluating mobile users requires the use
24 of specific techniques able to capture and identify usability problems that might
25 be experienced in mobile use. One exemplary issue in remote usability evalua-
26 tion involving mobile users is that they are *physically* moving, and such changes
27 in the context might imply a number of known and unknown variables poten-
28 tially affecting the set up (for instance, when increasing the amount of physical
29 activity, a significantly increased subjective workload might be experienced by the
30 users). In addition, the use of a particular platform should also be considered, with
31 the objective of identifying the appropriate means for collecting user data in the
32 remote site. For instance, eye-tracking systems are clearly useless for recording
33 user interactions with only-voice applications. Therefore, a current issue for this
34 dimension is represented by the capability of the different techniques for remote
35 evaluation to dynamically vary the information that should be collected about the
36 users, so as to cope with the potential issues that the specific platform in use can
37 introduce. The expected objective is providing the evaluator with the most com-
38 prehensive picture of all the aspects that might have affected the interaction, to
39 always be in a position to correctly derive the potential causes of a usability problem
40 that occurred on the client side. As previously mentioned, gathering information
41 about the current environment is extremely important for a mobile user, because
42 the environment can often change, ~~and~~ it becomes less important for a station-
43 ary user interacting with a desktop application because the environment is almost
44 fixed.

45 In this section, we analyze how the remote evaluation methods address the issues
46 raised by the specific platforms.

9.3.1 Desktop Applications

01 **9.3.1 Desktop Applications**
02
03 ~~Because~~ most applications have been developed for the desktop, the majority of
04 remote evaluation methods have addressed this type of platform. In Hartson, et al.
05 (1996), one of the first examples of remote-control evaluation is described. The
06 remote-control method checks a local computer from another computer at a remote
07 site. The user is separated from the evaluator in space and possibly in time. The
08 two computers can be connected through the Internet, or through a direct dial-up
09 telephone line with commercially available software (e.g., Timbuktu TM, PC Any-
10 where). Using this method, the evaluator's computer is located in the usability lab
11 where a video camera or scan converter captures the users' actions. The remote users
12 remain in their work environment and audio capture is performed via the computer
13 or telephone. If the audio capture is via telephone, the evaluator and user remain
14 connected at the same time. Alternatively, equipment in the usability lab could be
15 configured to automatically activate data-capture tools based on the use of a partic-
16 ular application. This is an example of quite a flexible technique for asynchronous
17 remote evaluation, which is restricted to use on desktop systems due to some soft-
18 ware limitations on the underlying hardware (e.g., PC Anywhere only operates on
19 PC platforms).
20
21
22

9.3.2 Vocal Applications

23 **9.3.2 Vocal Applications**
24
25
26 As for the vocal platform, interest is arising in studying this modality because the
27 associated technology is becoming more reliable and robust in different systems
28 in everyday use. The most exemplary case is that of the automatic response sys-
29 tems commonly used for completing user tasks such as banking, paying bills, and
30 receiving train/flight information. Such systems can accept both speech and touch
31 tone inputs, and in response provide relevant information through voice, email,
32 text messaging, and fax. Most companies are seeing such systems as an immedi-
33 ate cost saver because call centers are becoming too expensive to be operated by
34 humans.

35 As speech applications advance, so does the need for a means to evaluate vocal
36 user interfaces (VUIs), to be able to assess how a user interacts with a vocal appli-
37 cation. In this respect, we have to say that sometimes methods that are commonly
38 applied in GUI evaluation are also applied to the evaluation of VUIs, although this
39 translation is not a perfect fit. Indeed, the sequential nature of speech means that
40 VUIs are inherently more restrictive than GUIs, and therefore fewer choices can be
41 explored with a VUI in a certain time interval, with respect to what can be done
42 with a corresponding GUI. One of the consequences from the point of view of
43 usability evaluation is that, the number of tasks carried out in a certain interval of
44 time by interacting with a VUI will deliver a lower value when compared with a
45 corresponding GUI, without being necessarily a sign of a bad vocal user interface
46 usability.

01 Furthermore, it should be kept in mind that only using speech as an interaction
02 medium might represent a burden on users' memory, meaning that not only VUI
03 users should be focused on a smaller set of choices, and in a narrower context, but
04 also that without visual cues and a well-established mental model, they are even
05 unlikely to understand what choices are available to them. The consequence is that,
06 without careful design, these limitations can severely diminish the general usability
07 of the vocal application.

08 Despite the limitations noted above, well-designed voice applications have
09 proven to be both engaging and effective. Some novel evaluation methods for
10 these interfaces are under development, and several experiments in the lab have
11 already been done—although, with the proliferation of cellular and wireless phones,
12 evaluations of VUIs in lab environments suffer from unrealistic settings because
13 they are very different from the real contexts of use.

14 Several techniques can be envisaged for evaluating vocal user interfaces and,
15 more specifically, automatic response systems. Among them, we cite surveys, call
16 recordings, and call logs. Surveys are issued after a call is completed, but some-
17 times callers do not complete the call, hence never reaching the survey. As for call
18 recordings, most VUI systems record caller interactions, because call recordings
19 tell the VUI designer exactly what happened during each call. They have several
20 shortcomings—they cannot tell the designer what the caller was trying to do, how
21 the caller felt, or why the caller did what s/he did. Another shortcoming is the mas-
22 sive amount of effort required to analyze the calls. Lastly, in call logs, every inter-
23 active voice recognition (IVR) platform comes with extensive call-logging capabil-
24 ities. While surveys and call recordings typically result in qualitative data, call-log
25 data is typically quantitative (e.g., average call length, time on hold, abandon rate,
26 etc.). Due to the enormous amount of data that can be collected, data-mining tech-
27 niques are suitable for processing such data. Call logs identify where in the VUI
28 callers have difficulty, but this is only part of the picture. Call logs do not provide
29 a lot of context for helping to interpret the results. Therefore, because surveys, call
30 recordings, and call logs provide different information for in-use situation analysis,
31 it seems that careful consideration of their combined use as compensation for their
32 various advantages and drawbacks may be a viable solution for the purposes of
33 evaluating VUIs.

34 An example of tools for VUIs is ClickFox (2002), which aims to answer ques-
35 tions like: “What is the main cause of customer hang-ups? What are callers doing
36 most frequently at critical decision points? Are callers using the system in the
37 way that you expected?” Another example is provided by IQ Services (2006),
38 which is able to log and record each call, allowing IQ Services' analysts to dupli-
39 cate and experience system errors. After the test is completed, designers receive
40 online test results, step-by-step logs, and online playback of each digital test call
41 recording.

42 To conclude, while there are not many works on remote usability evaluation for
43 vocal applications, the naturalness of this kind of interaction, and its quick diffusion
44 in a number of applications covering different devices, make us expect that further
45 research will be done on this subject.
46

9.3.3 Mobile Applications

In mobile applications, it is important to understand the influence of the context of use, which is composed of three main parts: the user, the device, and the environment. Thus, one issue is to understand how usability is affected by dynamic changes of any of these components. Regarding evaluating interaction with mobile devices, the work of Denis and Karsenty (2003) focuses on the usability of a multidevice system, and introduces the concept of interusability to designate the ease with which users can reuse their knowledge and skills for a given functionality when switching to other devices. In their paper, a framework for achieving interusability between devices is proposed. It is based on two components: 1) a theoretical analysis of the cognitive processes underlying device transitions, and 2) an exploratory empirical study of the problems in using functionalities across multiple devices. Another work in this area is the paper by Waterson, et al. (Waterson, et al.), where the authors discuss a pilot usability study using wireless, Internet-enabled personal digital assistants (PDAs), in which they compare usability data gathered in traditional lab studies with a proxy-based clickstream logging and analysis tool. They found that this remote testing technique can more easily gather many of the content-related usability issues, whereas device-related issues are more difficult to capture. Lastly, the work of Stoica, et al. (2005) is worth mentioning, in which the authors describe a usability evaluation study of a system which permits collaboration of small groups of museum visitors through mobile handheld devices (PDAs). As the authors point out, techniques to measure usability-related factors generally include 1) inspection methods, 2) testing methods, and 3) inquiry methods. For systems including mobile devices, a combination of these techniques is sometimes used. As usability evaluation methodology, they propose a combination of a logging mechanism and an analysis tool—the ColAT environment (Avouris, et al. 2004), which permits mixing of multiple sources of observational data, a necessary requirement in evaluation studies involving mobile technology when users move about in physical space and are difficult to track. The museum system evaluated is based on a client-server architecture and an important characteristic of the application is that the server produces a centralized XML log file of the actions that take place during the visit. This log file can be combined with a video recording of the visit allowing evaluation of activity during the visit. In the experiment shown in the paper, different teams gathered the clues and then each group had to discuss and discover collaboratively what the combined clues were to solve the problem. The experiment was recorded by three video and two audio recorders for further analysis, using the ColAT analysis tool that interrelates activity and logs video and observers notes in the same environment. So, through ColAT, the actions that the users performed during the use of the PDAs, which were logged by the server, were synchronized with the videos. The methodology was able to deliver data useful for deriving quantitative information (e.g., total and average times for solving the puzzles, etc.), aspects related to group activities (number of exchanges between the group, and strategies used for solving the puzzles), and behavioral patterns of participants.

01 Indeed, the importance of performing a comprehensive evaluation that can take
02 into account data derived from multiple sources to adequately gain insight into large
03 bodies of multisource data, especially when mobile applications are considered, is
04 quite clear. An example of this trend can also be found in the work of Tennent, et al.
05 (2006), in which the authors present Replayer, which consists of a number of tools
06 (two video players, an audio player, an aggregate log visualization, a text search
07 tool, and a playback control tool), and a collaborative tool for analysis of recorded
08 data of mobile applications. The tool was designed ~~in the effort~~ to provide analysts
09 from a variety of disciplines (each using distinct sets of skills to focus on specific
10 aspects of the problem) with the ability to work cooperatively.

11 One of the emerging needs in this area is for tools that are better able to support
12 analysis of how task performance varies depending on the context change.

15 **9.4 The Techniques for Collecting Information about** 16 **the User Behavior**

17
18
19 In this section, we discuss the various techniques available for collecting informa-
20 tion regarding the user behavior (task performance, use of mouse and keyboard,
21 facial expressions, verbal comments, gestures, gazes, etc.). In this category, we
22 include several techniques for logging low-level user actions, other techniques for
23 gathering users' physiological information, and others capable of recording verbal
24 (and nonverbal) cues coming from the user's side (collected through a webcam
25 and/or a microphone).

26 It is worth pointing out that, while there are techniques that rely on commonly
27 available support, and can be used almost without any regard to the particular plat-
28 form considered (see, for instance, server-side logging techniques), other techniques
29 (e.g., eye-tracking) require specific hardware, whose use cannot neglect the partic-
30 ular platform in use.

33 **9.4.1 Logging (Server Side)**

34
35
36 This technique refers to Web-based applications and allows for collecting data at
37 the server side. Its effectiveness is strongly limited by the impossibility to capture
38 local user interaction with the user interface techniques (menus, buttons, fill-in text,
39 use of anchor links within the same page or Back button, . . .) and by the validity of
40 the server logs that cannot capture the accesses to the pages stored into the proxy
41 servers and the browser cache. For instance, if the requested page is in the browser
42 cache, then the request will never reach the server and is thus not logged. Moreover,
43 multiple people can also share the same IP address, making it difficult to distinguish
44 who is actually requesting what pages. Dynamically assigned IP addresses, where a
45 computer's IP address changes every time it connects to the Internet, can also make
46 it quite difficult to determine what an individual user is doing because IP addresses

01 are often used as identifiers. Thus, interpreting the actions of an individual user is
02 extremely difficult, because methods for capturing and generating Web usage logs
03 are not designed for gathering useful usability data, as pointed out by some works
04 (Etgen and Cantor 1999; Davison 1999; Pitkow and Pirolli 1999; Choo, et al. 1998;
05 Tauscher 1999).

06 Another method is to ask surfers to register online at the first visit and log
07 on with every subsequent visit. In this setting, the Web server can construct an
08 individual profile for each visitor, and track all user behaviors without ambiguity.
09 The Web server stores users' log-on names and their personal information, such
10 as age, gender, and occupation, and the visited pages. Such datasets are very rich,
11 and statistics on types of Web surfers, their interests and their browsing habits can
12 be generated with the Web mining process. This technique is widely adopted by
13 firms selling digital information products (e.g., online newspapers), which request
14 the users to log on before enabling file downloads. However, there are two main
15 limitations. First, Web visitors' choices are greatly reduced if they are required to
16 log on every time they visit the site. It becomes a serious issue for online firms and,
17 even for websites providing free registration, online users may reregister or provide
18 fake details. The statistics will become blurred, and this will result in invalid and
19 confusing conclusions. Second, the online firms cannot keep track of the visitors
20 once they leave to go to other websites. All generated knowledge is limited to only
21 a single website.

22 23 24 **9.4.2 Proxy-Based Logging** 25

26
27 This solution still supports Web-based applications through an intermediate server
28 between the client and the content server. Proxy servers are even less intrusive and
29 do not require any modification in the Web application to evaluate, but they limit
30 their analysis to the accessed page and are not able to capture the local user interac-
31 tions. The proxy approach has three key advantages over the server-side approach.
32 First, the proxy represents a separation of concerns. Any modifications needed for
33 tracking purposes can be done on the proxy, leaving the application server to deal
34 with just serving content, which makes it easier to deploy because the application
35 server and its content do not have to be modified. Second, the proxy allows anyone
36 to run usability tests on any website, even if they do not own that website. Lastly,
37 having testers go through a proxy allows Web designers to *tag* and uniquely identi-
38 fy each tester. Furthermore, a proxy logger also has advantages over client-side
39 logging. For example, it does not require any special software on the client side
40 beyond a Web browser, making it faster and much simpler to deploy. Therefore, the
41 proxy makes it easier to test a site with different test participants, operating systems,
42 and Web browsers than a client-side logger does, so allowing testing with a more
43 realistic sample.

44 An example of this kind of solution can be found in WebQuilt (Hong & Landay
45 2001), which uses a proxy logger to capture user accesses on the Web. As a proxy, it
46 lies between clients and content servers, with the assumption that clients will make

01 all requests through the proxy. Traditionally, proxies are used for things like caching
02 and firewalls. In WebQuilt, the Web proxy is used for usability purposes, with spe-
03 cial features to make the logging more useful for usability analysis. Although the
04 proxy-based technique seems quite appealing, there are still limitations on what the
05 WebQuilt proxy logger can capture. The most pressing of these cases is links or
06 redirects created dynamically by JavaScript and other browser-scripting languages.
07 As a consequence, the JavaScript-generated pop-up windows and DHTML menus
08 popular on many websites are not captured by the proxy. Another situation that
09 WebQuilt cannot handle is server-side image maps. Other elusive cases include
10 embedded page components such as Java applets and Flash animations. As tech-
11 nologies change and develop, the proxy will need to be updated to handle these new
12 cases.

15 **9.4.3 Logging (Client Side)**

18 In this category, various techniques are considered. Before analyzing them, it is
19 important to remember that client-logging is a technique that can be applied not
20 only to Web applications but also to Java and Microsoft applications with similar
21 results, because many tools have been developed for this purpose, as well.

22 In addition, it has been pointed out that through logging user interactions with a
23 given application, we can infer patterns of user behavior that indicate usability prob-
24 lems or other design deficiencies. This possibility has obvious attractions for Web
25 designers, but in HCI usability research some issues have been raised regarding the
26 possibility of identifying usability problems without access to the use context—to
27 the user's tasks and goals and to the user's own reports of what counts as a problem
28 for them. Thus, logging techniques alone are unlikely to provide useful results to
29 the evaluators.

30 *Cookies.* One method is to install cookies at Web client computers. A cookie
31 is a small text file that the Web server embeds in the browser for identifying the
32 user. If the user provides his name when he comes to a new site supporting cookies,
33 his name is stored in a plain text file at the client computer. No data is stored at
34 the server side, but every time the same browser asks for the page or the same
35 website, HTTP sends the cookie to the Web server, which uses it to identify the
36 user and display personalized information, such as name-calling greetings. One of
37 the advantages of using cookies is the ease of implementation. However, there are
38 two drawbacks. First, the amount of information stored in cookies is limited (the
39 average size is about 4K) and therefore, strictly speaking, no Web-mining process
40 can be performed based on such limited information. Second, because the cookies
41 are saved as plain text, they can be easily retrieved at the client computer. Hence,
42 security and privacy can be at risk.

43 *Client-side Logs.* They capture more accurate, comprehensive usage data than
44 server-side logs because they allow all browser events to be recorded, and it might
45 provide useful insight for usability evaluation. One alternative to gathering data on
46 the server is to collect it on the client side. Clients are instrumented with special

01 software so that all usage transactions will be captured. More specifically, clients
02 can be modified either by running software that transparently records user actions
03 whenever the Web browser is being used (as in Choo, et al. 1998), by modify-
04 ing an existing Web browser (as in Tauscher 1999), or by creating a custom Web
05 browser specifically for capturing usage information (as with Vividence 2000). The
06 advantage of client-side logging is that literally everything can be recorded, from
07 low-level events such as keystrokes and mouse clicks, to higher-level events such
08 as page requests. All of this is valuable usability information. However, there are
09 some potential drawbacks to client-side logging. First, special software must be
10 installed on the client, which end-users may be unwilling or unable to do. This can
11 severely limit the usability test participants to experienced users, which may not be
12 representative of the target audience. Second, there needs to be some mechanism
13 for sending the logged data back to the team that wants to collect the logs. Third,
14 the software, in some cases, is platform-dependent, meaning that the software only
15 works for a specific operating system or a specific browser.

16 Paganelli and Paternò (2003) developed a tool for performing client-logging of
17 Web applications: the main advantages are that it does not require expensive equip-
18 ment, and facilitates the problem of modifying the evaluated pages because it auto-
19 matically includes JavaScript code in all the pages that have to be evaluated. Such
20 Javascript snippets are able to adapt to the various features of different browsers.
21 Using a browser's log-based analysis, the evaluator can accurately measure time
22 spent on tasks or particular pages, as well as study the use of the Back button
23 and user clickstreams. It is also possible to precisely identify the downloading time
24 and the time when the page is visible to the users. In addition, their tool is able to
25 automatically analyze the information contained in Web browser logs and compare
26 it with task models specifying the designer model of the possible users' behav-
27 iors when interacting with the application to identify whether and where users'
28 interactions deviate from those envisioned by the system design and represented
29 in the model. Within this client-side technique, we also cite the work (Ho 2005)
30 developed in the e-commerce domain area, which is about the use of a *user remote*
31 *tracker* to examine Web users' characteristics, trying to draw a linkage between
32 Web customers' characteristics and their browsing behaviors. The authors propose
33 a user-remote tracking framework based on Web services and XML to track every
34 HTTP request from client computers to understand surfers' characteristics. The
35 user-remote tracker is a piece of software installed in the users' browser to keep
36 track of every keyboard input and mouse click from the users. No matter what the
37 users input, all HTTP requests and responses are tracked by the software program,
38 including interactions with Java Applet programs. This program will automatically
39 send the activity log file, together with the user identity, to a central machine for
40 Web-mining (instead of sending such information directly to the Web server). It is
41 that central machine that analyzes clickstreams and generates navigation rules for
42 these users through algorithms. There are several advantages with this user-remote
43 tracker. First, it can follow users everywhere. Second, while server-logging cannot
44 track the interaction between a user and an applet program, the tracker can solve
45 this problem. Third, in the traditional data collection method, it is possible to get
46 little information once the users enter the secure websites (i.e., websites started with

01 https://). Here, because the user-remote tracker uses low-level programs to track
02 every user input signal, the activities can be tracked even in this case.

03
04
05
06
07

05 **9.4.4 Eye-Trackers**

08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30

Eye trackers are a technique for measuring users' eye movements so that it is possible to know both where a person is looking at any given time, and the sequence in which their eyes are shifting from one location to another, on the screen. Tracking people's eye movements can help evaluators understand visual information processing and the factors that may impact the usability of system interfaces, thereby providing an objective source of data that can inform the design of improved interfaces. Evaluators using eye-tracking, however, should take into account the limits of such technology and how such limits impact the data collected. For example, an appropriate minimum threshold time for a fixation should be carefully identified, because interpretations can vary a lot according to the time set to detect a fixation in the eye-tracking system. Moreover, eye trackers might have difficulty tracking participants who have lenses. Furthermore, visual distractions (e.g., colorful or moving objects around the screen or in the testing environment) should also be eliminated, as these will inevitably contaminate the eye-movement data. Also, eye-tracking generates huge amounts of data, so it is essential to automatically perform filtering and analysis. Eye-tracking technology, however, has evolved in recent years and there are now more systems that can be used for remote evaluation (see the Tobii system, <http://www.tobii.com/>) because they can be transported in suitcases and do not require that users wear intrusive equipment. It is only necessary to hold an initial standard training exercise. Nevertheless, one of the most relevant problems with the eye-tracking technique remains the fact that it is possible to know what users see but not what users think about what they see. In other words, how data is actually being processed by the person.

31
32

33
34
35

33 **9.4.5 Webcam/Audio Recorders or Microphones**

36
37
38
39
40
41
42
43
44
45
46

The use of webcams and audio recorders allows for acquiring more contextual information about the data collected. Indeed, as it has been previously mentioned, through logging keystrokes and webpages on a given site, we could infer patterns of user behavior that indicate usability problems or other design deficiencies. In HCI usability research, however, it has been argued that it is not possible to identify usability problems without access to the context of use, to the user's tasks and goals, and to the user's own reports of what counts as a problem for them. Webcam-based videos are very valuable when further analysis is necessary when an error is found, because the evaluator can analyze the video clip and convert it into a usability problem description, or use it in any case to understand the reason of a usability problem. For instance, videos can be valuable in capturing facial movements/expressions,

01 verbal/vocal signals and expressions, non-verbal communication, body language,
02 and posture. Moreover, facial expressions may provide indications of the immediate
03 appreciation of the system by showing the instantaneous reactions to the system, and
04 also might reflect the subject's considerations about the system. Furthermore, the
05 use of more than one camera is valuable for capturing some environmental condi-
06 tions occurring in the testing environment. Work by Lister (2003) has been oriented
07 to using audio and video capture for qualitative analysis performed by evaluators on
08 the results of usability testing.

09 Also in the work of Paternò et al. (2006), webcams are used to record the users
10 (not the users' screens) to provide valuable information for interpreting problematic
11 parts of the user interaction. For instance, in this work, videos are also used to check
12 user behavior whenever some measurements (e.g., time needed for completing a
13 task) captured by another software component provide unexpected values.

16 9.4.6 Sensors

17
18
19 ~~By the term sensors, we mean more sophisticated~~ research solutions for data acqui-
20 sition and analysis of some physiological data. Recently, a number of sensors are
21 being used more and more for the evaluation of user interfaces, trying to take into
22 account the emotional dimension of computer-human interaction (e.g., affective user
23 interfaces). Among such measures, we cite physiological signals like ECG, respira-
24 tion, galvanic skin response, heart rate, and skin temperature. Most of them, such
25 as galvanic skin response (GSR), heart rate (HR) and blood volume pulse (BVP)
26 are generally chosen as good, physically non-invasive indicators of stress (under
27 stress, GSR and HR increase, while BVP decreases), and are also easy to be mea-
28 sured with specialized equipment. In this respect, we mention the ProComp system
29 manufactured by Thought Technology, Ltd. (<http://www.thoughttechnology.com>), or
30 the BIOPAC system (<http://www.biopac.com/>), which allows for recording different
31 kinds of data—physiological signals, vocal/verbal signals, and non-verbal signals
32 (posture, gaze direction, facial movements). Unfortunately, the use of sensors in
33 remote usability evaluation is currently suffering the limitation of the highly spe-
34 cialized equipment necessary, which cannot be assumed available in users' daily
35 environments (although it is slowly appearing and used more and more in telemedicine
36 applications). However, more research effort is envisaged in the next years on this sub-
37 ject for the useful information that it can provide to analyze the user's emotional state.

38 To summarize, almost all the results obtained with each technique indicated in
39 this section requires additional knowledge about the user from the evaluator to be
40 actually useful for the purposes of the evaluation. Therefore, the big issue is that
41 such data is not informative *per se* about possible usability problems, but requires
42 further comparison with supplementary information. One of the few exceptions can
43 be identified in, for instance, recording users positively (or negatively) commenting
44 on the session while interacting with the application in a remote think-aloud session
45 (which should theoretically provide the evaluator with immediate feedback about
46 the user's satisfaction). In almost all the other cases, a further contextualization

01 (and integration) of the data collected is needed to correctly evaluate the session
02 state (think about, for example, the uselessness of logging mouse and keyboard
03 actions without contextualizing such actions within the current user intention). One
04 of the current issues is identifying techniques enabling an easy synchronisation and
05 aggregation of all such different sources of information in some semantic context,
06 to facilitate the evaluator's work.

09 9.5 The Type of Application Considered

12 In this section, we analyze another dimension of the proposed framework—the type
13 of application, considered in terms of the underlying software environment. As with
14 already analyzed dimensions, in this case the consideration of this dimension is also
15 not completely independent from the other ones. Indeed, the type of application
16 considered may prevent (or strongly promote) applying specific techniques men-
17 tioned in the previous section, as well as the use of particular interaction platforms.
18 For instance, while in the case of Web-based applications we have seen that there
19 are several options about where the logging tool should work (e.g., server, client,
20 proxy, ..) regardless of the particular platform at hand, the consideration of .NET-
21 based applications for remote evaluation has only occurred in recent years. It is
22 almost always connected with stationary platforms because only recently prototyp-
23 ical tools that support the evaluation of .NET applications for mobile devices have
24 appeared, and are still limited in terms of the information they are actually able to
25 provide.

26 Indeed, the first applications that were evaluated with some type of remote
27 evaluation were graphical applications, often implemented in languages such as
28 Java (e.g., Paternò & Ballardin 2000). Then, with the advent of the Web and the
29 related ease of performing a remote evaluation when the Web is considered (due to
30 the related simplicity in involving a high number of testers with little effort), the
31 majority of methods have considered websites as their primary evaluation targets.
32 Java-based applications indeed have been taken into account, sometimes as a sort of
33 side-effect of the desire to improve the flexibility of techniques considered for Web
34 applications whenever applets are also included. An example of this can be found
35 in the already mentioned work of Ho (2005), developed in the e-commerce domain
36 area. It is about the use of a user-remote tracker to examine Web users' character-
37 istics, trying to draw a linkage between Web customers' characteristics and their
38 browsing, and with the capability of tracking client-side logs, including interactions
39 with Java Applet programs. Microsoft .NET applications have been considered as
40 well (i.e., PDA devices), for which they often provide more robust and supported
41 solutions with respect to Java. An example of logging tools for Microsoft envi-
42 ronments is the VibeLog logging tool (<http://research.microsoft.com/vibe/>), which
43 has been developed at Microsoft Research to evaluate the ways that work practice
44 might change as users move between various-sized displays during their work day.
45 The logging tool is married with ethnographic research data, which should provide
46 good indications of what parts of the Windows and Office designs do not scale well
across different display sizes. This analysis is used to understand where they should

01 orient their research efforts in novel visualization and interaction development, with
02 an eye toward designing more elegant UIs.

03
04

05 **9.6 The Type of Evaluation Results**

06

07
08 Before analyzing the last dimension of the proposed framework in depth—the type
09 of results an evaluation can deliver (i.e., qualitative vs. quantitative data)—we judge
10 it useful to mention the work by Petrie et al. (2006) about remote evaluation. In
11 this work, the authors highlight how both formative and summative evaluations can
12 be supported by remote techniques. Indeed, in *summative* evaluations, one of the
13 main goals is to understand whether the users can install and run a system on their
14 own and on their own machines, and how they rate the key functions of the system.
15 Therefore, the disparate environments and configurations that can be reached
16 with remote evaluation can provide highly reliable data in this respect. In *formative*
17 evaluations, the objective is to collect information about design flaws and inform
18 redesign. Therefore, it is particularly important that participants feel free to criticize
19 a system and avoid evaluator bias, and this may be easier if they are in the privacy
20 of their own environment, rather than the potentially more threatening situation of
21 the usability laboratory.

22 However, the authors take it a step further, claiming that, in particular, remote
23 asynchronous techniques in which the evaluator cannot intervene during the user
24 sessions are especially useful for *summative* evaluation. To support this idea, the
25 authors report on two evaluations conducted with disabled users. In both cases, they
26 performed both a local and a remote evaluation. The technique used for remote evaluation
27 was, in one case, making notes on problems encountered and then sending
28 them to the evaluator, and in the other case, recording problems encountered and
29 then sending them along with ratings of the accessed websites. Both remote and
30 local evaluations provided considerable quantities of qualitative data, but the local
31 evaluations provided far richer data because the researchers were able to record
32 problems that the participant may not have been aware of, and are in a position to
33 prompt the participant to explore these problems, comment on them, and analyze
34 what had caused them.

35 On one hand, achieving the rich interaction between participants, researchers and
36 developers, as requested by *formative* evaluation, is very difficult in remote evaluation
37 situations. With high quality video conferencing, broadband connections, and
38 remote recording systems, however, it might be possible to conduct remote evaluations
39 that capture a rich set of data. On the other hand, if the evaluation is *summative*,
40 a remote evaluation may be quite appropriate because it adequately shows real user
41 behavior.

42 We agree with this position to some extent because, in our opinion, remote evaluation
43 can provide different types of results, which can, in turn, be used for different
44 purposes, both summative and formative. In this section, we are going to analyze
45 the type of results an evaluation can deliver. In particular, a discussion about the
46 type of information that can be useful to obtain in order to analyze the multimodal
data regarding user sessions is provided. This information can be quantitatively

01 determined by specific software and highlighted during the evaluation (tasks not
02 completed, errors occurring during the performance of tasks, time for completing
03 a task, etc.) together with other information that deals with intrinsic qualities of
04 the user interface (e.g., time needed for the performance of a task). It is not sur-
05 prising that some relations exist between the evaluation techniques mentioned in
06 Section 9.4 and the evaluation results analyzed in this section (for instance, sensing
07 technologies deliver quantitative data about the user's emotional and physical state),
08 while in other case such correspondence is not so straightforward.

11 **9.6.1 Task-Related Information**

13 Many applications are task-oriented and therefore some important aspects to con-
14 sider are whether the users are able to accomplish the desired tasks, and information
15 regarding task performance (task duration, number of actions, ...). One issue is
16 how to know what the desired tasks are. One possible solution is to ask users to
17 explicitly indicate the tasks at the beginning of the session. The issues associated
18 with user errors are related to the issues associated with task performance, because
19 user errors are actions not useful for the current task. The errors are good indicators
20 of bad usability and difficulties in task accomplishment. Tasks can be considered at
21 various granularities. In some cases, it can be interesting to analyze the performance
22 of short basic tasks, and in other cases it is important to focus on the performance
23 of high-level complex activities.

24 During an analysis of task performance, it can be useful to analyze when it devi-
25 ates from the ideal expected behavior, and to what extent. The evaluator then has to
26 understand the reasons for such mismatch and needs to go back and analyze what
27 happened for each action in the user session and what factor triggered the deviation.

31 **9.6.2 Qualitative Information**

33 Under this heading, we mean all the techniques that allow evaluators to collect
34 qualitative data from the users. As we already mentioned, qualitative data are quite
35 relevant, especially if the kind of evaluation is *formative*, therefore the richness of
36 the qualitative data is very important in understanding how to improve the sys-
37 tem. For example, gathering informal and spontaneous comments in natural lan-
38 guage from the users undoubtedly offers valuable information to the evaluators
39 for improving the resulting design. Also, because this information can provide a
40 rich contextual knowledge about the situation currently occurring during the user's
41 interaction, it may also be used for cross-checking other collected data that may be
42 too ambiguous to be interpreted—generally quantitative data. An example in which
43 this strategy might disambiguate other data, is the case of a user spending a long
44 time visiting a page. Considering only this quantitative information (registered by
45 the browser-logging tool) would not allow the evaluator to assess whether the users
46 found the information very important or just had problems in finding the concerned

01 information—in this case, the webcam can help in correctly interpret the feeling of
02 the user (engaging or not the visit).

03
04

05 **9.6.3 Presentation-Related Data**

06

07 In this section, we analyze the results that the evaluation should deliver about the
08 usability of the user-interface presentation (e.g., for GUIs—layout, choice of wid-
09 gets, colors, labels, etc.). There are tools that link the task performance with the
10 user-interface elements supporting such performance. In other cases, the tools are
11 able to provide reports that highlight the user-interface elements that might be prob-
12 lematic from the usability point of view. For example, WebQuilt (Waterson & others
13 2002) provides representations consisting of nodes representing visited webpages,
14 and arrows representing the traffic between the pages. Entry pages are green and
15 exit pages are cyan. Thicker arrows represent heavier traffic. Arrow color is used
16 to indicate time spent on a page before transitioning, where the closer the arrow is
17 to red, the longer the user spent in transition. The designer's path is highlighted in
18 blue. There is a slider along the left-hand side that allows the designer to zoom into
19 the graph, viewing actual images of the pages users saw and where they clicked.

20
21

22 **9.6.4 Quantitative Cognitive-Physiological Information**

23
24

25 Quantitative psycho-physiological measurements can provide useful information
26 about more general, qualitative information on a human's feeling in a specific situa-
27 tion. For instance, with a growing population of elderly persons today, this result
28 is expected to be more and more applied in the field of elderly care/assistance,
29 where there has been an increasing interest in investigating algorithms to enable
30 the possibility of assessing elderly mood in a non-intrusive manner. To make state-
31 of-mind information available, sensor technology can be employed. Various psycho-
32 physiological signals are known in literature that can convey the presence of strong
33 emotions or stress (Cacioppo, et al. 2000)—skin conductance, muscle tension, heart
34 rate, and heart rate variability. Such signals can now be measured in an unobtrusive
35 manner. The measured signals have to be analyzed to reliably convey short-term
36 mood changes (that might be relevant for the relatives and form a basis for an
37 enhanced feeling of connectedness), as well as long-term trends. When the shape
38 of people is mostly visible, computer vision tools can be used to classify their post-
39 ure and gait, and posture changes over time. This information can be exploited,
40 for example, to predict (by gait analysis) and detect (by analyzing posture changes)
41 falls. Computer vision techniques can be used to detect the head position in real-
42 time, and classify the facial orientation (frontal, profile) to provide the process of
43 facial expression analysis with suitable data. In addition, faces are processed for
44 expression/recognition/authentication. If a person is not visible or the user does
45 not like a camera to be used (e.g., in the bathroom), speech/audio tracking is an
46 alternative.

01 Eye-tracking systems can provide many interesting pieces of information derived
02 from fixations and saccades. Long fixations can indicate that users spend too much
03 effort to interpret or to process what they are looking at. The number of fixa-
04 tions is often related to the user efforts to process the content of the screen area
05 being analyzed. The duration of the scanpath is a productivity measure and can
06 be compared with a theoretical optimal duration. Even the ratio between saccades
07 and fixations can be a useful index for comparing the percentage of time spent in
08 looking for information (saccades) and that during which information is acquired
09 (fixations).

10
11

12 **9.7 An Example Tool for Remote Evaluation:** 13 **MultiModalWebRemUsine**

14
15

16 In this section, we discuss an example of a tool for remote evaluation according to
17 the framework presented in the chapter. The basic idea of this tool is to analyze user
18 logs through the semantic information contained in task models. Thus, on the one
19 hand, we have a task model that describes how designers expect users to perform
20 their activities, and on the other hand, there are logs indicating the actions performed
21 by the users while interacting with the application. Each user session can be defined
22 through the sequence of the corresponding user actions, which can be associated
23 with a corresponding sequence of basic task performance to achieve the user's goal.
24 If the performed task sequence diverts from those enabled by the task model, there
25 is clearly a mismatch that needs to be analyzed by the evaluators. Either the task
26 model is too rigid or there is something unclear in the user interface, which prevents
27 the user from performing the expected sequences of tasks.

28 Various versions of the tool have been developed, which vary for the type of
29 application addressed and the type of results provided. The first version, USINE
30 (Lecerof & Paternò 1998), mainly addressed the issue of using task models for
31 analyzing user logs without considering its use as remote evaluation. The next
32 version, RemUSINE (Paternò & Ballardini 2000), was developed for remotely evalu-
33 ating desktop Java applications and was tested in industrial sites, providing useful
34 information regarding its possibilities even in comparison with other methods. It
35 was, for example, compared with a video-based evaluation. It turned out that for
36 evaluating a small number of sessions, the video-based evaluation was more efficient
37 because RemUSINE required some time to enable the automatic evaluation given
38 that the evaluator first has to provide the task model of the designed application and
39 create mappings between basic tasks and log events. On the positive side, it was
40 noted that, in some cases, video analysis is not able to detect quick user actions
41 (such as some user clicks) and is not usable for evaluations when users are located
42 far from the evaluator.

43 Given the explosion of the Web, which has become the most common user
44 interface, we thought it useful to develop a new version (WebRemUSINE), aimed
45 at evaluating this type of application (Paganelli & Paternò 2003). We also had to
46 decide how to log user interactions. For this purpose, we implemented an efficient,

01 interoperable, client-side logging system. In addition to information regarding task
02 performance, the Web-oriented version provides a lot of information regarding the
03 Web pages analyzed—visited pages, never visited pages, extent of scrolling and
04 resizing, page patterns, and download and visit time. Some information is provided,
05 along with summary data regarding the content of the page. Thus, the visit time
06 is provided and also indicates the number of forms, links, and words in the page
07 so the evaluator can compare the visit time with the quantity of information avail-
08 able in the page. The latest version of the tool (MultimodalWebRemUSINE) aims
09 to exploit the possibilities opened up by recent technologies to gather a richer set
10 of information regarding user behavior. Thus, the traditional graphical logs can be
11 analyzed together with the logs from webcams and portable eye-trackers, which do
12 not require the use of intrusive equipment.

13 In summary, the changes in the tool mainly aimed at fulfilling the evolving needs
14 of usability evaluators. Indeed, the tool started from the original vision of providing
15 cost-effective techniques for usability evaluation for analyzing data about product
16 usage in a real-world environment. To this end, remote evaluation is valuable when
17 trying to keep budgets down while staying competitive in the marketplace (which is
18 especially relevant for companies). The necessity to reach larger, more diverse and
19 dispersed pools of participants also stimulated the attention to Web applications. In
20 these times of global customers and development organizations, there is a clear cor-
21 relation between the globalization of the product market and the potential (and chal-
22 lenges) of remote evaluation. Next, the tool kept evolving in these directions with
23 an eye toward the available improvements (in terms of robustness and affordabil-
24 ity) of technology and broadband infrastructure, which were efficiently exploited
25 for enriching the tool with multimodal information on the user's behavior. The
26 objective was to compensate the recognized evaluator's decreased ability—typically
27 connected with remote usability evaluation techniques—to interpret the motivations
28 underlying a certain user behavior, due to separation in space (and sometimes also
29 in time) between the user and the evaluator.

30 In general, MultimodalWebRemUsine is based on a comparison of planned user
31 behavior and actual user behavior. Information about the planned logical behavior
32 of the user is contained in a (previously developed) task model, while data about
33 actual user behavior is provided by the other modules supposedly available within
34 the client environment (the logging tool, the webcam and the eye-tracker).

35 Before starting the test, users have to explicitly indicate the target task. After that,
36 all the user actions will be automatically recorded. The evaluation then analyzes the
37 user's sequences of actions to determine whether the user has correctly performed
38 the tasks in accordance with the temporal relationships defined in the task model,
39 or if some errors occurred. In addition, the tool evaluates whether the user is able
40 to reach the goals and if the actions performed are actually useful to reach the pre-
41 defined goals, by means of an internal task model simulator. For each action in the
42 log, the corresponding basic task is first identified and then there is a check to see
43 whether that task was logically enabled. If no error occurs, the list of the basic tasks
44 that have been enabled after its performance is provided, together with the updated
45 list of high-level tasks already accomplished, to allow the evaluator to check ~~to see~~
46 if the target task has been completed. Otherwise, some error will be notified in the

01 report analyzing the user session. An example of error is a *precondition error*, which
02 means that the actual user's task performance did not respect the relations defined
03 in the system design model. For example, if people want to access a remote service
04 (such as Web access to emails), they usually have to provide username and password
05 and then activate the request through a button. If the user interfaces elements are not
06 located in such a way that the user can easily realize that both fields have to be filled
07 in before connecting to the mailbox, then some precondition errors can occur (for
08 example, the user sends the request without first proving the password). Such types
09 of errors can be detected through this type of approach.

10 From the log analysis, the tool can generate various indications:

- 11 • *Success*—the user has been able to perform a set of basic tasks required to accom-
12 plish the target task and thus achieve the goal
- 13 • *Failure*—the user starts the performance of the target task but is not able to
14 complete it
- 15 • *Useless uncritical task*—the user performs a task that is not strictly useful to
16 accomplish the target task but does not prevent its completion
- 17 • *Deviation from the target task*—in a situation where the target task is enabled
18 and the user performs a basic task whose effect is to disable it. This shows a
19 problematic situation since the user is getting farther away from the main goal in
20 addition to performing useless actions
- 21 • *Inaccessible task*—when the user is never able to enable a certain target task
22

23 Recently, we have paid attention to how to represent user sessions and related
24 data in such a way that eases their analysis. Figure 9.1 shows the type of represen-
25 tations designed. It is possible to show data related to several sessions in different
26 ways at the same time. In Figure 9.1, we analyze the parts of the sessions about
27 users who want to become member of an association. As you can see from the
28 selected radio buttons, the deviation graph is shown for the first two users, while
29 the state graph is visualized for the other ones. In both types of graphs, the white
30 circles are associated with the basic tasks performed, and their positions indicate
31 when they have been accomplished. In the first diagram (deviation diagram), there
32 are three lines—one for the basic tasks correctly performed, one for those uselessly
33 performed, and one for the tasks that have diverted the user from achieving the
34 current goal. In the state diagram, the color of the line underlying the white circles
35 is used to indicate whether the user is correctly or incorrectly performing the task.

36 A further type of information considered during the evaluation regards the
37 task-execution time. In case of tasks correctly performed, the tool calculates the
38 global time of performance. This information is calculated by examining the tem-
39 poral information associated with each event and stored in the logs. The duration
40 is calculated for both high-level and basic tasks. The set of results regarding the
41 execution time can provide information useful to understanding what the most
42 complicated tasks are or what tasks require a longer time to be performed. In
43 Figure 9.2, a screenshot of the tool is presented. As you can see, whenever an
44 inexplicably lengthy time period for carrying out a certain task is registered by the
45 tool, the evaluator can activate the related video recorded through a webcam to
46 gather further information.

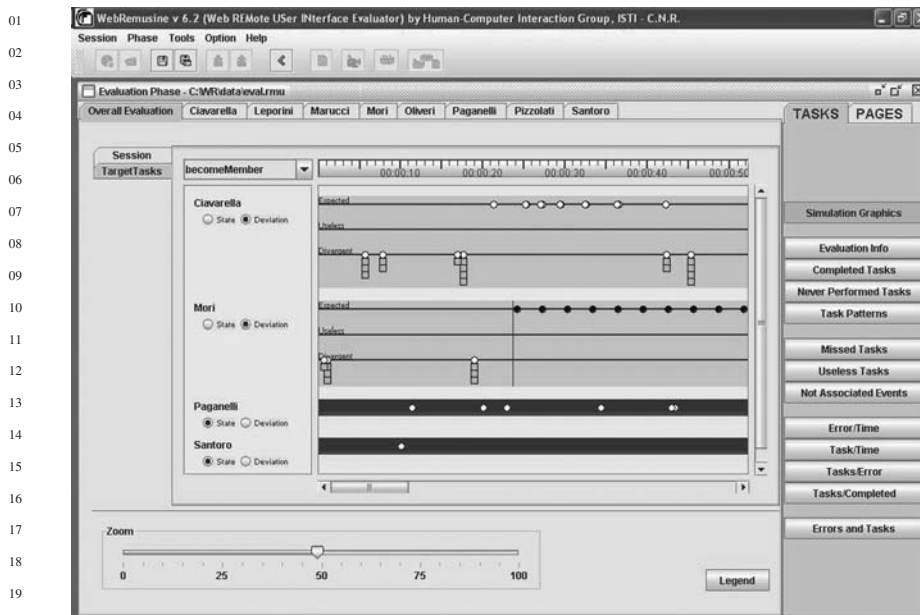


Fig. 9.1 Representation of user sessions in MMWebRemUsine

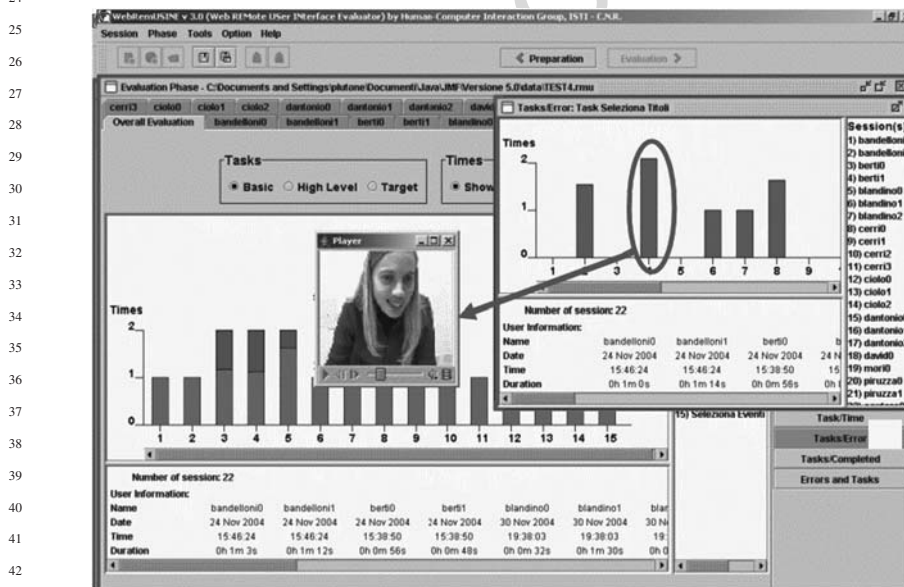


Fig. 9.2 Highlighting a video within the MultiModalWebRemusine environment

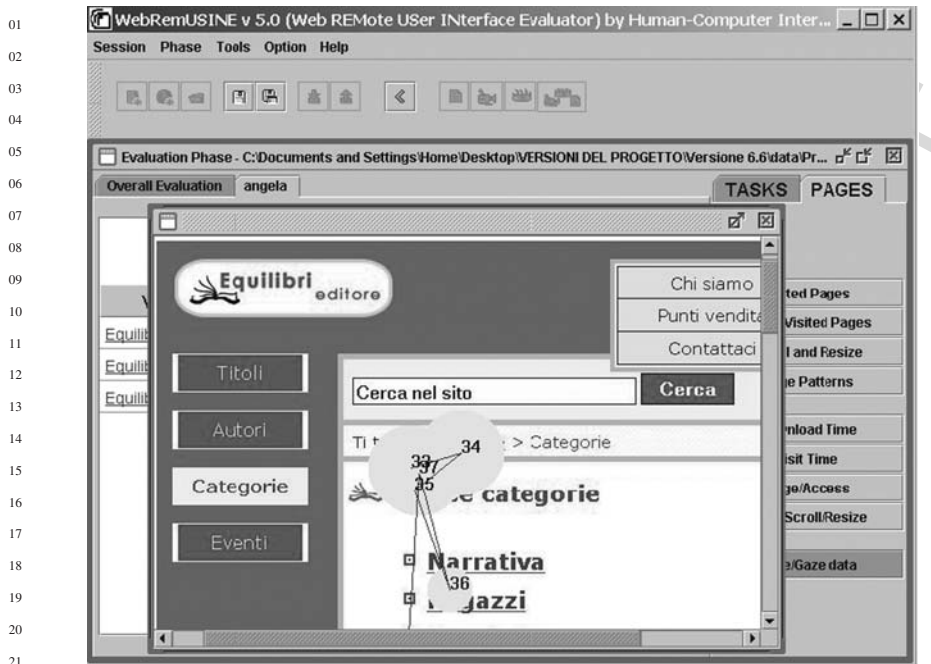


Fig. 9.3 An example of visualization of fixations (the yellow areas) and scanpaths (the paths in-between) registered by the eye-tracker

The approach supported by MultiModalWebremusine is also able to provide usability data associated with every single presentation. Moreover, it is worth noting that because this approach can correlate task-based measures with presentation-related data, it can also analyze the usability of the website from both viewpoints. For example, the tool can compare the time to perform a task with that for loading the page(s) involved in such a performance.

Eye-tracking is a technique that has been used within MultiModalWebRemUsine. In Figure 9.3, a screenshot has been taken that shows how the registrations of the eye-tracker are visualized within the MultiModalWebRemUsine environment. As you can see, each fixation is represented by an area whose size is proportional to its recorded duration, while the lines connecting such areas (scanpath) highlight the path the user followed while visiting the page.

9.8 Discussion and Interrelationships between the Framework Dimensions

We can discuss possible interrelationships among the different dimensions that we have identified while analyzing several contributions in the area of remote usability evaluation. First of all, the choice of the platform type might substantially limit—in

01 terms of hardware and software—the type of technology and/or techniques that can
02 be used for remote evaluation, as well as put less/more emphasis on the relevance
03 of particular information for the purposes of the evaluation. Indeed, apart from the
04 well-known differences between the type of applications that might be supported by
05 a cellphone and by a desktop system due to the diversities in hardware and software
06 capabilities, the type of platform used may also affect the relevance of a piece of
07 information with respect to another one. This is the case, for instance, with the
08 user context, which is very important for mobile applications and less important for
09 desktop systems.

10 Another observation is the fact that the same type of information useful for
11 the evaluation (e.g., user workload) can be gained with different techniques. For
12 instance, a possible indication for user workload might be a user blinking very
13 frequently (information that can be gained from an eye-tracker), but some physi-
14 ological data can—more reliably—signal this workload. Another example regards
15 user feedback, which can be gained with different means. A user nodding (captured
16 by the webcam) might be a sign of a good feedback, as well as a user using some
17 positive vocal expressions. In other cases, there is information useful for evaluation
18 (e.g., task-based data) that cannot be derived without explicitly including additional
19 information (task specification).

20 Moreover, the type of interaction between the evaluator and the user might also
21 affect the use of the particular technique(s) adopted, as well as the quantity (and
22 also the quality) of information collected for the evaluation. For instance, while the
23 use of remote questionnaires indicates a specific type of technique for collecting
24 information about the user, in automatic collection of data the range of techniques
25 can greatly vary, as does the range of the evaluation results that can be derived from
26 interpreting the collected data.

27 Lastly, as we already noticed, the type of application considered may prevent
28 (or strongly encourage) the application of certain techniques, as well as the use
29 of specific interaction platforms. For instance, in the case of Web-based applica-
30 tions, the use of server-side and client-side logging techniques is a well-known and
31 established approach, while the consideration of user interactions with additional
32 software components like .NET for remote evaluation only belongs to recent years,
33 and is often restricted only to specific types of platforms.

34
35

36 9.9 Conclusions and Future Challenges

37
38

39 In this chapter, we have described a framework composed of different dimensions
40 that we have identified as relevant in the area of remote usability evaluation. This
41 type of evaluation is becoming more and more important in a time of globalization
42 of companies and their customers. We have used such a framework to review a large
43 spectrum of methods that have been proposed in this area. These methods have been
44 receiving more and more interest due to the improvements in techniques that are able
45 to capture information regarding user behavior and the validity of the data that are
46 collected *in the field*. Therefore, the chapter tries to shed some light on the different

01 methods through a common framework in which it is possible to compare and con-
 02 trast current works in the area of remote evaluation, as well as delineate possible
 03 future trends in the research agenda of remote usability evaluation. In addition, this
 04 chapter is useful in identifying which are the current strategies for compensating
 05 some traditional weaknesses in this type of evaluation. For instance, future work
 06 should be dedicated to extending the gathered data regarding user behavior and
 07 state—including emotional state—so a more complete analysis of what happens
 08 during user sessions can be done and potential usability issues can be better iden-
 09 tified. Another novel emerging application area is that of mobile applications, in
 10 which is important ~~in understanding~~ how task performance varies depending on the
 11 changes in the context of use. To make the discussion more concrete, we have also
 12 reported our experience with our tool in the area of remote evaluation and analyzed
 13 it according to the dimensions of the logical framework proposed.

17 References

- 18
- 19
- 20 Avouris N., Komis V., Margaritis M., & Fiotakis G., (2004), An environment for studying col-
 21 laborative learning activities, *Journal of Educational Technology & Society, Special Issue on*
 22 *Technology – Enhanced Learning*, 7 (2), 34–41.
- 23 Cacioppo, J.T., Berntson, G.G., Larsen, J.T., Poehlmann, K.M., & Ito, T.A. (2000). The psy-
 24 chophysiology of emotion. In M. Lewis & R.J.M. Haviland-Jones (Eds.), *The Handbook of*
 25 *Emotions* (2nd Ed.) (pp. 173–191). New York: Guilford Press.
- 26 Card, S., Pirolli, P., Van der Wege, M., Morrison, J., Reeder, R., Schraedley, P., & Boshart, J.
 27 (2001). Information scent as a driver of Web behavior graphs: Results of a protocol analysis
 28 method for Web usability, *Proceedings ACM CHI 2001* (pp. 498–504).
- 29 Choo, C. W., Detlor, B., & Turnbull, D. (1998). A behavioral model of information seeking on
 30 the Web—preliminary results of a study of how managers and IT specialists use the Web. In
 31 *Proceedings of the 61st ASIS Annual Meeting*, 35, (pp. 290–302).
- 32 ClickFox (2002), ClickFox Inc., <http://www.clickfox.com>.
- 33 Davison, B. (1999). Web traffic logs: An imperfect resource for evaluation. In: *Proceedings of*
 34 *Ninth Annual Conference of the Internet Society (INET'99)*, San Jose, CA, June 1999.
- 35 Denis, C., & Karsenty, L. (2003). Inter-usability of multi-device systems: A conceptual framework.
 36 In A. Seffah & H. Javahery (Eds.), *Multiple User Interfaces: Engineering and Application*
 37 *Framework* (pp. 375–385). New Jersey: John Wiley and Sons.
- 38 Etgen, M., & Cantor, J. (1999). What does getting WET (Web Event-Logging Tool) mean for Web
 39 usability? In *Proceedings of the Fifth Conference on Human Factors and the Web*. Gaithersburg,
 40 MD, June.
- 41 Hartson, R. H., Castillo, J. C., Kelso, J. T., & Neale W. C. (1996). Remote evaluation: The network
 42 as an extension of the usability laboratory. In: *Proceedings of CHI 1996* (pp. 228–235).
- 43 Ho, S.Y. (2005). An exploratory study of using a user remote tracker to examine Web users' per-
 44 sonality traits. In *Proceedings of the 7th International Conference on Electronic commerce,*
 45 *ICEC'05* (pp. 659–665), August 15–17, 2005, Xi'an, China. ACM Press.
- 46 Hong, J.I., & Landay, J.A. (2001). WebQuilt: a framework for capturing and visualizing the Web
 experience. In *Proceedings of WWW 2001 Conference* (pp. 717–724).
- IQ Services (2006). *IQ Services: Interactive Quality Services, Inc.* Access at: <http://www.iq-services.com/>
- Ivory M. Y., & Hearst M. A., (2001) The state of the art in automating usability evaluation of user
 interfaces. *ACM Computing Surveys*, 33(4), 470–516.

- 01 Lecerof, A., & Paternò, F. (1998). Automatic support for usability evaluation. *IEEE Transactions*
02 *on Software Engineering*, 24 (10), 863–888.
- 03 Lister M. (2003). Streaming format software for usability testing. In *Proceedings ACM CHI 2003,*
04 *Extended Abstracts* (pp. 632–633).
- 05 Paganelli, L., & Paternò F. (2003). Tools for remote usability evaluation of Web applications
06 through browser logs and task models. *Behavior Research Methods, Instruments, and Com-*
07 *puters*, 35 (3), 369–378.
- 08 Paternò, F., & Ballardini, G. (2000). RemUSINE: a bridge between empirical and model-based
09 evaluation when evaluators and users are distant. *Interacting with Computers*, 13(2), 229–251.
- 10 Paternò, F., Piruzza, A., & Santoro, C. (2006). Remote Web usability evaluation exploiting multi-
11 modal information on user behavior. In *Proceedings CADUI 2006*, Bucharest. Springer-Verlag.
- 12 Petrie, H., Hamilton, F., King, N., & Pavan, P. (2006). Remote usability evaluations with disabled
13 people. In *Proceedings of CHI 2006* (pp. 1131–1141), Montréal, Québec, Canada, April 22–27,
14 2006.
- 15 Scholtz, J., Laskowski, S., & Downey L. (1998) Developing usability tools and techniques for
16 designing and testing Websites. In *Proceedings HFWeb'98* (Basking Ridge, NJ, June 1998).
17 Access at: [http://www.research.att.com/conf/hfWeb/](http://www.research.att.com/conf/hfWeb/proceedings/scholtz/index.html) proceedings/scholtz/index.html
- 18 Stoica, A., Fiotakis, G., Simarro-Cabrera, J., Frutos, H.M., Avouris, N., & Dimitriadis, Y. (2005).
19 Usability evaluation of handheld devices: A case study for a museum application. In *Proceed-*
20 *ings of PCI 2005*, Volos, November 2005.
- 21 Tauscher, L.M. (1999). *Evaluating history mechanisms: An empirical study of reuse patterns*
22 *in WWW navigation*. MS Thesis, Department of Computer Science, University of Calgary,
23 Alberta, Canada.
- 24 Tennent, P., Chalmers, M., & Morrison, A. (2006). Replayer: Collaborative evaluation of mobile
25 applications. Presented in *CHI'06 Workshop on Information Visualization and Interaction Tech-*
26 *niques for Collaboration across Multiple Displays*, Montreal, Canada.
- 27 Tullis, T., Fleischman, S., McNulty, M, Cianchette, C., & Bergel, M. (2002). An empirical compar-
28 ison of lab and remote usability testing of websites. In: *Proceedings of Usability Professionals*
29 *Conference*, Pennsylvania, 2002.
- 30 Waterson, S., Landay, J.A., & Matthews, T. (2002). In the lab and out in the wild: Remote Web
31 usability testing for mobile devices. In: *Proceedings of CHI 2002* (pp. 796–797), April 20–25,
32 Minneapolis, USA.
- 33 West, R., & Lehman, K.R. (2006). Automated summative usability studies: An empirical eval-
34 uation. In *Proceedings of CHI 2006* (pp. 631–639), April 22–27, 2006, Montréal, Québec,
35 Canada, ACM Press.
- 36
37
38
39
40
41
42
43
44
45
46

01 **Chapter-9**

02

03 Query No.	Page No.	Line No.	Query
05 AQ1	203	35	In who's paper-Stoica?

06

07

08

09

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

UNCORRECTED PROOF