

Usability and Transparency in Automatic Support for Web Accessibility Validation

Iannuzzi Nicola, Marco Manca, Fabio Paternò, Carmen Santoro

CNR-ISTI, HIIS Laboratory, Pisa Italy

{iannuzzi.nicola, marco.manca, Fabio.paterno, carmen.santoro}@isti.cnr.it

The importance of guaranteeing accessible Web applications is becoming widely recognised, and supported by national and international legislation. This implies an increasing need for large scale validations, which can be obtained only through automatic support. Thus, there is a need for a new generation of accessibility validation tools able to validate against the continuously evolving guidelines and report the associated issues, considering that their results will be used by many people with varied backgrounds and goals. Such tools should be able to support monitoring of Web sites' accessibility, provide user-friendly information suitable for various purposes, and be transparent in terms of what they are actually able to evaluate. They should also consider how the technologies for implementing Web sites have evolved in recent years. In this paper, we provide a description of the design dimensions that characterise such tools in order to address the emerging needs, indicate how we have addressed and implemented them in a specific tool, and report on a first usability test providing encouraging feedback.

Keywords: Accessibility, automatic validation tools, accessibility monitoring

1 INTRODUCTION

The importance of providing accessible Web applications for all, including people with cognitive or physical disabilities, has become increasingly recognised. This is confirmed by the indications provided by national and international legislations to support it ([Gulliksen et al., 2010], [Lazar & Olalere, 2021], [Paternò and Schiavone, 2015]). A major initiative to address such aspects is the EU Directive on the “Accessibility of the Websites and Mobile Applications of Public Sector Bodies” that came into force on 26/10/2016, also known as the Web Accessibility Directive (WAD) [EU Commission, 2016], which establishes accessibility requirements for the websites and mobile applications of public sector bodies. One aspect that is particularly stimulated by this directive is monitoring, which should be done more systematically in terms of the number of Web pages involved, and with some level of frequency.

In parallel, the guidelines for accessibility of Web Sites, which are developed by W3C in the Web Accessibility Initiative, are continuously evolving considering the need for better addressing the various possible disabilities and the evolution of Web Technologies. The current version (2.1) is structured into principles, guidelines, success criteria, and techniques, which have increased with respect to the previous versions. Thus, the WCAG 2.1 added 17 new success criteria to address mobile access and some disabilities better, so that now there are overall 82 success criteria.

All such aspects imply that thorough accessibility validation requires considerable effort for the number of elements and aspects that have to be checked. Thus, the interest in automatic support for this activity is continuously increasing, and stimulates further research and development in this area because of its potential to support the collection and analyses of data on the effective application of the accessibility guidelines, detect non-compliance in a consistent manner, and provide relevant information on how to address possible problems. At the same time, it is important to be aware that not all

accessibility issues can be automatically detected, some of them require manual checking from accessibility experts, and subjective feedback is still important to consider [Power et al., 2012].

In order to provide automatic support for accessibility, many proposals have been put forward. Bobby was probably the first accessibility tool widely used. Developed and released in 1996 by the Center for Applied Special Technology (CAST), Bobby was a free online accessibility tool that was used to evaluate against WCAG 1.0, WAI and Section 508. In about 2006 it was acquired by IBM and then removed from service. Over time several researches and development efforts have been carried out in this area (e.g. [Beirekdar et al., 2002], [Ivory et al, 2003], [Beirekdar et al., 2005], [Fuertes et al., 2009], [Fernandes et al., 2014], [Nietzio et al., 2011], [Schiavone and Paternò, 2015]). As of May 2022, the W3C Web Accessibility Evaluation Tools list¹ contains 161 elements, and further tools, which are not included in it, have been put forward. However, many of them have limited impact for several reasons: some have not been able to evolve in order to address the most recent WCAG version (e.g. [Gay and Li, 2010], [Mirri et al, 2011]); some address only specific aspects, such as colour contrast or readability [Miniukovich et al., 2019] or lexical simplification [Moreno et al., 2019]; some provide only information in national languages and there is not an English version.

Abascal et al. [2019] provide a set of useful criteria to analyse the support provided by accessibility tools: the type of license (free versus commercial); the platform where they can be executed; the evaluation scope (ranging from single pages to entire websites); the support provided for repairing identified issues; how the evaluation results, guidelines supported, and detected issues are rendered and exported. However, there are at least three emerging important aspects missing in such classification: *accessibility monitoring*, *tool transparency*, and *support for dynamic Web sites*. The first means the ability to indicate a set of Web pages and periodically check the level of accessibility in order to inform relevant stakeholders about how it evolves. The second aims to address one important problem that users of automatic validators often encounter when using multiple tools: they may provide different results. For example, a study on automatic Web accessibility evaluation [Abduganiev, 2017], which only considered support for the previous WCAG 2.0 guidelines, analysed eight popular and free online automated Web accessibility evaluation tools finding significant differences in terms of various aspects (coverage, completeness, correctness, validity, efficiency and capacity). Users of such tools are often disoriented by such differences, and find them somewhat unclear. Thus, it becomes important that the tools be transparent and indicate in detail what they are actually able to validate [Parvin et al., 2021]. The last aspect aims to address the increasing use of development frameworks that implement dynamic Web sites, such as Angular or Vue.js. In these cases, the Web pages that are created and arrive to the browsers contain a few elements, and depend on JavaScripts that are executed at loading time and deeply modify the actual Web page content. Thus, a traditional validator that only analyses the initial version of the page would be able to actually consider very few elements. Equipping the validator with functionalities able to perform server-side rendering, and then analyse the corresponding results would provide more meaningful results. Overall, we think it is time for a new generation of automatic tools for accessibility validation characterised by the ability to address such aspects.

In this paper, we present the design aspects that have to be considered in order to address such issues, and how we have addressed them to be included in a previous validation tool [Broccia et al., 2020], which has a modular approach to managing guidelines (which has been adopted also in other tools [Pelzetter, 2021]), it is publicly available and has a large community of users. In particular, after discussing related work, we indicate the requirements that have driven the novel work, present the corresponding design and implementation in the tool, and report on a first user test, which provided encouraging feedback. To our knowledge, such aspects have not been reported in the literature. There are some commercial

¹ <https://www.w3.org/WAI/ER/tools/>

tools that provide some support in this direction but their authors have not described how they obtain it. Thus, this paper can be useful for tool and application developers in order to understand how the emerging requirements can be addressed, and for all the community interested in Web accessibility validation to better understand the actual possibilities of automatic support.

2 RELATED WORK

While there are several tools that offer some support to locate and visualise errors in some way, the issue of supporting monitoring functionalities has been addressed in a more limited manner till now. Indeed, interest in the possibility of monitoring Web sites on a periodical basis has increased recently, also thanks to the EU Directive that enforced recurrent monitoring of websites of public bodies and organisations (which typically include many pages). Prior work already discussed and compared different automatic accessibility evaluation tools (see e.g. [Abascal et al., 2019], [Vigo et al. 2013]), but only limited attention has been put up to now to evaluate how such approaches support monitoring the accessibility of Web sites over time even though there is increasing need of having effective means for such aspects. A first exploration of how to support monitoring of Web sites accessibility at a geopolitical level was discussed in [Mirri et al., 2011], but that tool provided only some limited representations in a tabular format of the accessibility levels detected for older WCAG versions. More recently [Burkard et al., 2021] have presented a comparative analysis of four commercial accessibility validation tools. The evaluated tools were SiteImprove, PopeTech, aXe Monitoring and ARC Monitoring, which were evaluated by the authors by using trial versions sent to them by their respective vendor companies. In particular, in that work, the analysis was done in analytical and empirical manner. In analytical manner, by using a set of evaluation criteria (i.e. coverage of web pages, coverage of success criteria, correctness, support for localisation of errors, degree of implementing gamification patterns) with specific weights assigned to such criteria. In empirical manner by involving 15 users in a study in which each participant had to freely explore every tool on their own for some time (i.e. without concrete tasks to carry out). The evaluation criteria were partly based on those used in other studies ([Abduganiev and Gaibullojonovich, 2017], [Padure and Pribeanu, 2019], [Vigo et al., 2013]) for comparing automatic Web accessibility evaluation tools. In addition, they also considered aspects aiming to understand how user-friendly and motivating they are to use. However, their analysis did not focus on the aspects that are addressed by our proposal (monitoring, transparency, and support for dynamic sites).

Some tools, such as Adaplugin, Accessibility Scanning & Monitoring by UserWay, and AccessiBe, offer some monitoring support and support recurrent scan. However, they are commercial ones and therefore not directly and easily exploitable by all interested/relevant stakeholders. Generally, they provide free demos or e-tours, or in the best case, a free trial for a limited period of time (i.e. 10-14 days) with several limitations (i.e. scan just a single page, one time per month, such as Userway), as the only alternative to purchasing one of the available subscription packages. We aim at providing better support in a tool that can be freely and easily accessed and used by different relevant stakeholders interested in monitoring the accessibility of their websites.

Generally, it is easy to see that when applying different validation tools to the same Web content, they provide different results, and users have difficulties understanding the reasons for such variability, and to what extent the results are meaningful. Thus, also for improving their usability, there is a need for more transparency to help users better interpret their results [Parvin et al., 2021]. In another work, [Vigo et al., 2013] analysed the effectiveness of six frequently used accessibility evaluation tools in terms of coverage, completeness, and correctness concerning the WCAG 2.0 guidelines. They found that coverage was narrow as, at most, 50% of the success criteria were covered, and similarly, completeness ranged between 14% and 38%. In addition, some of the tools that exhibit higher completeness scores produced lower

correctness scores (66-71%) because catching as many violations as possible can lead to an increase in false positives. Lastly, they indicated that the effectiveness in terms of coverage and completeness could be boosted if the right combination of tools is employed for each success criterion. A further study on automatic Web accessibility evaluation [Abduganiev, 2017], which only considered support for the previous WCAG 2.0 guidelines, has analysed eight popular and free online automated Web accessibility evaluation tools finding significant differences in terms of various aspects (coverage, completeness, correctness, validity, efficiency and capacity). More recently, [Padure et al., 2019] compared five automatic tools for assessing accessibility. The result of the study indicates that the combined use of two of the considered tools would increase the completeness and reliability of the assessment. Frazao and Duarte [2020] focused their analysis of accessibility on validation plugins extensions for the Chrome Web browser. They found that individual tools still provide limited and varied coverage of the success criteria. After analysing their results, they recommend using more than one tool and complementing automated evaluation with manual checking. Solutions based only on plugins are able to address issues related to validation of dynamic pages since they access directly the DOM in the browser but are not suitable for monitoring Web sites. In general, none of such studies focused on the transparency aspects and how to help users understand how the accessibility evaluation tools work by providing clear information about their coverage and working. Indeed, in a survey [Ysilada, 2015], respondents strongly agreed that accessibility must be grounded on user-centred practices and that accessibility evaluation is more than just inspecting source code.

3 THE DESIGN OF THE TOOL

3.1 Requirements

In our work, we started by identifying a set of requirements that should characterise a new generation of tools for supporting accessibility validation. For this purpose, we analysed the state of art in the scientific literature in this area, in which some studies discussed the usability problems of accessibility evaluation tools (e.g. [Molinero et al., 2006], [Petrie et al., 2007], [Brajnik et al., 2012], [Salehnamadi et al., 2021]). In addition, in previous work [Paternò et al., 2020] some requirements elicitation activities (online questionnaires, interviews and workshops) for new tools for accessibility validation were carried out in the context of the WADCHER European project. We also considered our direct experience with tools in research and international projects, collaboration with the national agency for accessibility, teaching accessibility validation in HCI courses, and, more generally, with the analysis of current accessibility validation practices, and observations of feedback gathered from interaction with accessibility experts and users of the previous version of our tool. The resulting requirements are listed below.

R1. *The tool should be able to support different types of stakeholders.* Managing a public Web site typically involves several people with different backgrounds and goals. In general, we can identify *Web commissioners*, who are responsible for the site and indicate its purpose, but often have little knowledge about the design and implementation techniques for accessibility. *Web developers* are those who actually implement the Web site, and often need support to understand how to address the accessibility guidelines. Sometimes there is also the involvement of *User Interface Designers*, who have knowledge on user experience and how to support it, but often do not know accessibility guidelines. In some cases, there are *Accessibility Experts*, who can manually check whether the Web elements satisfy the guidelines, but so far, there is still a rather limited number of this type of experts, and their work is problematic since modern Web sites contain many elements to check, and this is quite difficult to do manually.

R2. *The tool should be able to validate single and groups of Web pages.* The support provided by the tool should be flexible in terms of the granularity of the pages that are validated, also because at different times there may be different

types of requests. They can range from on-the-fly validation of a single page, which for some reason is of particular interest because someone may have complained about its accessibility, or it is particularly important and requires a more complete verification, to groups of pages, up to entire Web sites. In this perspective, the EU directive 2016/212 has indicated a methodology (EU 2018/1524) whereby the pages should be validated according to two modalities. The *in-depth* one is intended to thoroughly verify whether a website satisfies all the requirements identified in the standards and technical specifications referred to in the EU Directive. Therefore it shall evaluate the interaction with forms, interface controls and dialogue boxes, the confirmations for data entry, the error messages and other feedback resulting from user interaction when possible, as well as the behaviour of the website or mobile application when applying different settings or preferences. Then, there is a *simplified* modality, to detect instances of non-compliance with a sub-set of the requirements in the standards and technical specifications referred to in the EU Directive, and related to the requirements of perceivability, operability, understandability and robustness.

R3. *The tool should be able to monitor over time the level of accessibility.* Often Web sites are modified because the content and the functionalities need to be updated, or they are modified because some usability or accessibility problems have been signaled by end users, thus there is a need for periodically checking the accessibility level to see whether it has improved. In addition, the EU directive requires that all the European countries provide monitoring of their levels of accessibility, also indicating the number of websites to monitor depending on the number of inhabitants.

R4. *The tool should be able to connect the errors identified to both the page code and the page user interface.* Once errors are identified, on the one hand web developers need immediate support to localise them in the code to understand how to fix them. However, especially for stakeholders different from the technical ones (e.g. web commissioners), it is also useful to indicate the user interface element affected by that error (locate the error within the user interface of the page) to better understand its actual impact on the user, and how its modification can be addressed.

R5. *The tool should be able to provide quantitative summary information on the actual level of accessibility identified.* One characteristic of the accessibility validators is that often they generate long lists of issues detected; in some cases they are errors, in others they are warnings, several of them are minor issues, thus their impact can significantly vary and in the end people have difficulties estimating the overall accessibility level. Thus, some summary quantitative estimations can be useful. Some efforts in the area of accessibility validation have been put forward; however, some of the metrics proposed are complicated and require deep technical and accessibility knowledge that often people do not have, thus resulting rather obscure.

R6. *The tool should be transparent indicating what it is actually able to check.* One common issue that people who use accessibility validators encounter is that for a given Web site, they often provide different results. One of the main reasons for such differences is that they vary in terms of the techniques that they are able to validate, or differently implement the accessibility check. However, this is not immediately understandable because usually such tools claim to support a given set of guidelines (e.g. WCAG 2.1) without indicating to what extent they actually support it. Indeed, the validation of WCAG 2.1 is implemented by analysing many techniques, and some of them cannot even be automatically checked. If we focus on those that can be checked we still have a large number of techniques, and it is rather rare that the validators provide an indication of those they are able to address, and more generally how much the validation performed is complete.

R7. *The tool should provide practical indications about how to solve the identified problems.* Once the issues are identified, their correction requires some technical knowledge, and some immediate access to relevant resources would certainly be appreciated, and make more efficient the correction process. However, in this case the challenge is that sometimes there is no general one-fits-all solution for solving a specific issue, as a meaningful modification can additionally require understanding the context of the error (i.e. the content and state of the relevant part of the web page).

R8 *The tool should be able to validate dynamic Web pages.* Often modern Web sites are developed with JavaScript frameworks that deeply modify the content of the pages through scripts that are executed at browser's loading time. Thus, a complete validation should be able to address also the dynamically generated content, even though this increases the complexity of the validation since the generated content is usually much larger and more varied than the initial static version.

3.2 The Tool Functionalities and How they are Presented to the Users

This section reports the tool characteristics linking them to the requirements identified in the previous section. By using the tool, it is possible to evaluate a single page or create a project, which is a set of pages to validate using the same validation settings consistently within an "audit" (R2). An audit is an accessibility inspection that can be done either once or can be scheduled for being periodically and automatically activated at a defined interval of times decided by the user. The introduction of "projects" allows users to request multiple evaluations of the same group of web pages (belonging to a project), to be able to compare the obtained results and monitor the evolution of accessibility over time (R3). Users can configure a project according to a set of parameters based on the type of evaluation to carry out, we defined three different project types:

- *Single Page Project.* This is the simplest case of project, and it should be selected when the user wants to monitor the evolution of accessibility over time of just one web page. For creating a Single Page Project, it is necessary to specify just its URL, the level of conformance requested, the type of device (desktop, smartphone, tablet) and the user agent (operating system and browser) to simulate when accessing the Web page. It is also possible to carry out the validation of a single page *without creating an associated project*: in this case it is possible to choose whether the validation will consider only static content or the dynamic one (by using the server-side rendering validation). Instead, when validating a *project*, the tool considers the dynamic content obtained through server-side rendering (R8), even if this takes longer time, given that the results are provided asynchronously.
- *Simplified Project.* It evaluates a subset of pages belonging to a website (see Figure 1) starting from a base URL, and selected according to two parameters: the maximum number of pages to consider, and the maximum depth to reach when selecting the pages starting from the home page in the target directory.
- *In-depth Project.* It corresponds to the most thorough validation: in this case the tool evaluates a list of representative pages such as Home, Login, Sitemap, Accessibility Statement and verifies whether the website satisfies all the requirements identified in the standards and technical specifications referred to in Article 6 of EU Directive 2016/2102.

For both types of projects (simplified and in-depth), it is also possible to schedule a specific frequency by selecting the days to automatically perform the evaluation (e.g. Sundays, Mondays), and whether it should be performed every week, every two weeks, or four weeks. It is possible to select the particular set of guidelines to use (i.e. WCAG 2.1 or WCAG 2.0), the conformance level, and also the device and user agent to consider for the validation.

In addition, on a page called "Info", the tool reports the list of Success Criteria and the associated Techniques that it is able to evaluate for providing transparent information regarding what it is actually able to validate (R6). It also provides access to the XML-based specification of how each technique is interpreted by the tool for validation purposes. Thus, users can better understand how and when the tool determines erroneous cases.

Figure 1: Creating a Project

After having created a project, it will be added to the user’s workspace together with the other projects available for that user. The tool contains a specific section showing the information associated with each specific project. Figure 2 shows the top-most part, namely the one dedicated to providing a summary information of the selected project, from left to right: the details of the project (creation date and main configuration parameters), the list of audit(s) associated with it, and the URLs of the web pages considered in this project.

When the user selects a specific *audit* of a project (through the link “View audit” in the central panel of the window shown in Figure 2), it is possible to get the information associated with it. The details are shown divided into the following sections:

- **Results:** a summary of the results produced by the considered audit, in terms of the average number of techniques (calculated over the pages belonging to that audit) with error/warning/success/non-applicable results;
- **Page Results:** a graph visualising, for each page, the number of errors (or warnings) found: the X axis corresponds to the pages belonging to the audit, the Y axis to the number of occurrences of the different types of issues found (errors/warnings). By hovering the mouse over each dot shown in the graph it is possible to display the corresponding information (the URL of the page considered, number of errors or warnings);
- **List of Evaluated pages:** the list of pages assessed in the audit is provided in the last section. Further information on the evaluation of each page can be seen by selecting the corresponding link.

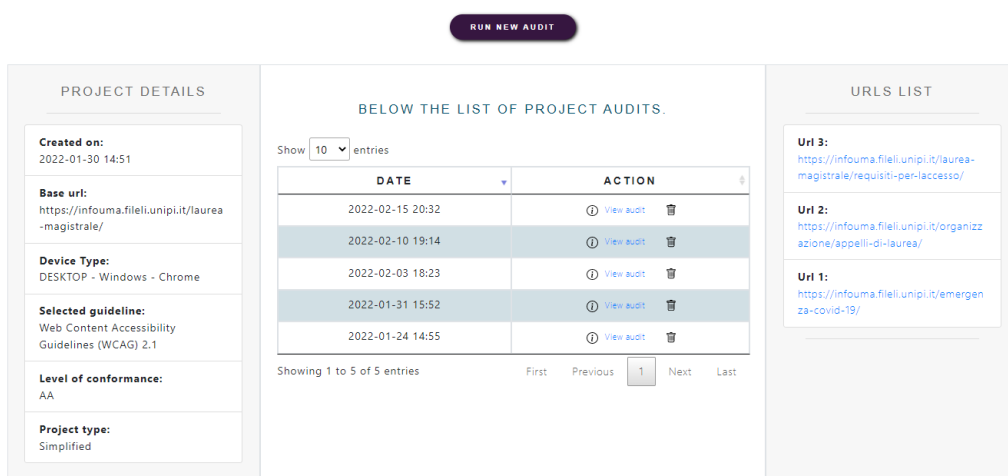


Figure 2: The top part of the page dedicated to showing the information about a project (summary information)

When the user selects a specific page evaluated by the tool (through one of the links provided in the section “List of Evaluated Pages”) different pieces of information are provided via multiple tabs (see Figure 3):

- **Evaluation Summary (R5)**, which provides the parameters specified by the user at validation creation time, the summary of the corresponding results in terms of number of errors, warning, success and not applicable techniques, the possibility to download the evaluation report in PDF or in EARL² format [Abou-Zahra, 2017] (the W3C standard to represent test results), and two metrics that provide overall information about the accessibility level of the evaluated pages;
- **End User View (R1)**, which provides the results grouped according to various parameters, i.e. principles (e.g. perceivable, operable), categories (e.g. ARIA, content) and code type (e.g. HTML, CSS), which can be understood even by people with limited technical knowledge;
- **Live Preview (R1, R4)**, which allows users to visually locate the identified issues directly within the web page user interface, and connect them to the corresponding code;
- **Web Developer View (R1)**, which highlights the errors/warnings within the code, particularly useful for Web developers.

It is worth noting that the goal of the last three views (End-User, Live Preview and Web Developer views) is to support different stakeholders that need to access accessibility information for different purposes (R1).

² <https://www.w3.org/WAI/standards-guidelines/earl/>

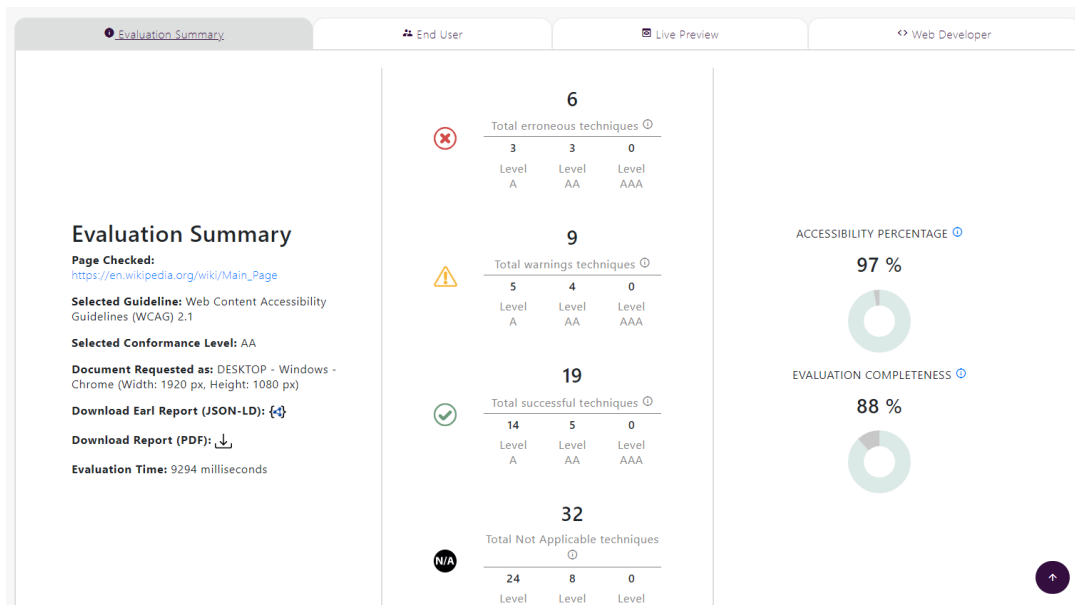


Figure 3: Evaluation Summary

The **Evaluation Summary** provides an overview of the issues found in the validations of the various techniques. In particular, it groups them into four categories, counting the number of i) techniques that in at least some cases resulted as violated (*errors*); ii) techniques that resulted in successful application (*success*); iii) techniques that were applied but their evaluation did not give a definitive answer, therefore a manual check is needed (*warnings*); iv) techniques that was not possible to apply because not relevant for the considered web page (*non-applicable*). It also provides such numbers according to the conformance levels.

The two metrics are visualised using a doughnut chart representation. They are calculated based on the following values: S=number of elements successfully evaluated; E= Number of elements unsuccessfully evaluated; W= Number of elements that require a manual evaluation. They are both percentages, and are:

- **Accessibility Percentage** = $S/(S+E)$: Number of distinct techniques successfully evaluated out of the total number of techniques for which the tool was able to make a successful or unsuccessful evaluation; The sum of elements unsuccessfully evaluated is weighted by taking into account the conformance level (<https://www.w3.org/TR/WCAG21/#cc1>) of each evaluated technique: a level A technique is considered with weight 1, level AA with weight 0.6 and level AAA with level 0.2. This has been done because errors generated by techniques of level AAA are considered less 'important' than level A techniques and consequently their impact on the accessibility level calculation is lower.
- **Evaluation Completeness** = $(S+E)/(S+E+W)$: Number of distinct techniques for which the tool was able to carry out a successful or unsuccessful evaluation compared to the total number of techniques evaluated. We introduced this metric to make it more transparent to the users that the tool is not able to decide on the accessibility of all the web elements analysed. Thus, even if the accessibility percentage score is 100%, they still have to check the evaluation completeness to understand whether the automatic validation has been able to decide on the accessibility of all the analysed elements.

Such metrics are helpful to provide a compact indication of the current level of accessibility (accessibility percentage), but also remind their users that the automatic evaluation may not be complete, and provide an indication of the level of completeness (evaluation completeness).

The goal of the “**End User**” view (see Figure 4) is to provide information useful for Web commissioners, or in general people with limited technological knowledge but still interested in understanding the aspects more problematic of their website from an accessibility perspective. Here the results are provided according to the involved WCAG 2.1 accessibility principles (namely: perceivable, understandable, robust, operable): for each type of issue found (among errors and warning) the technique that was involved is reported as well as the corresponding success criterion, and this information is reported under the specific WCAG 2.1 principle that is involved. It is also possible to select the “Categories” tab to see the type of elements involved in the errors, and the “Code Type” tab to see which ones are HTML or CSS errors.

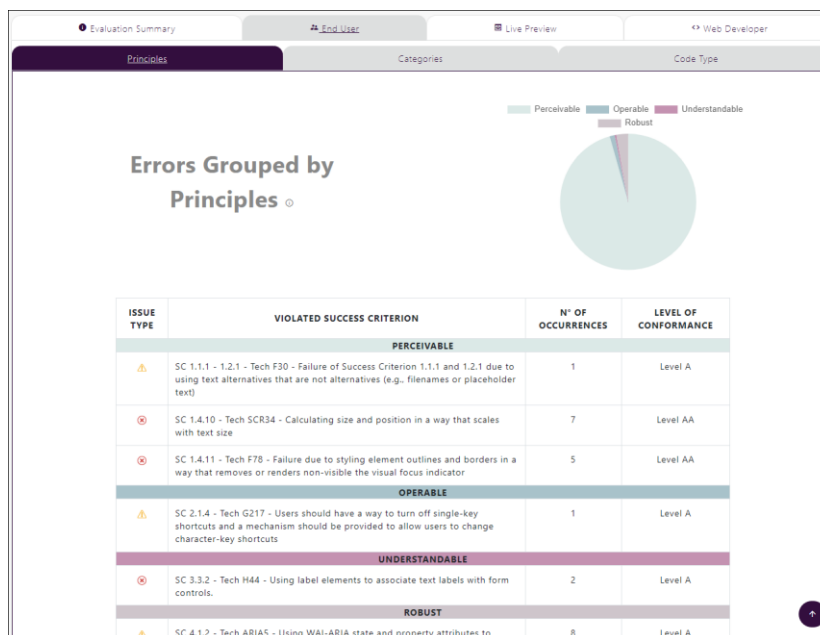


Figure 4: End User View with errors grouped by principles

The “**Live Preview**” aims to provide end users with the possibility to easily localise a specific error directly within the user interface of the considered page, to have more information about the error, and better understand its potential impact on the user (Figure 5). There is a left-hand panel with two parts. One is dedicated to filtering the elements to analyse: errors or warnings, HTML or CSS elements, errors related to specific WCAG principles. The other part shows the filtered elements: for each type of error the number of occurrences and the list of such occurrences are indicated. The user can select a specific occurrence of an error (by clicking on the associated eye-shaped icon) and then automatically the associated part of the web page is highlighted in red within the page. If the user hovers the mouse over that highlighted element, guidance information about how to solve the associated issue is displayed (R7 - see the “How to solve?” tooltip in Figure 5).

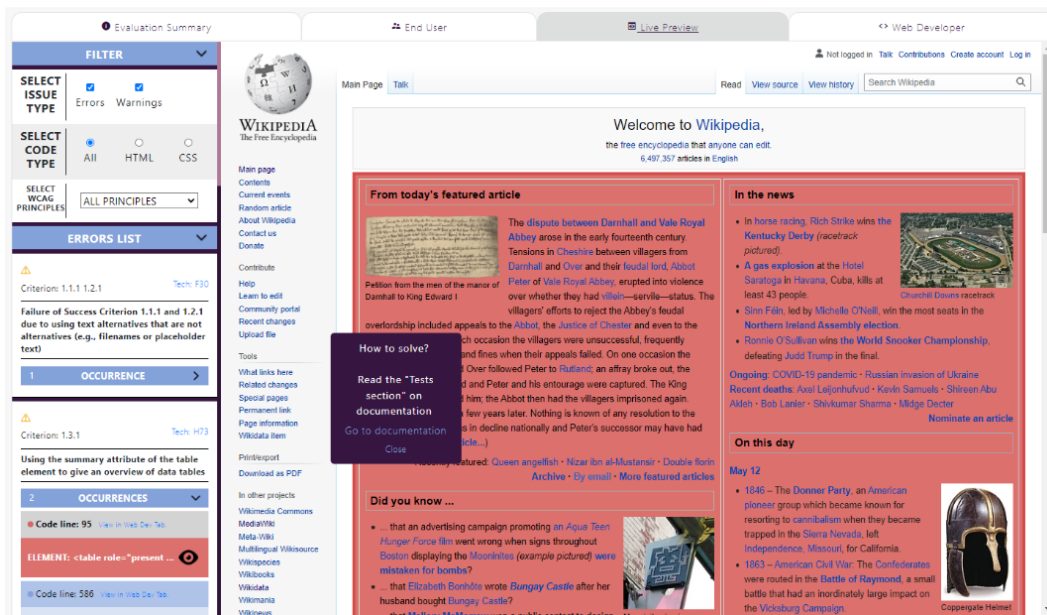


Figure 5: The Live Preview showing (left part) the list of accessibility issues found by the tool and (main panel) a visualisation of the involved page

The **Web Developer View** (Figure 6) presents the HTML code listing, highlighting in different colours the rows that have been affected by accessibility issues: red for errors and yellow for warnings. The rows that are affected by accessibility issues are interactive: by selecting one of them, the user is redirected to the corresponding W3C page referring to the involved WCAG technique. In the right-bottom part of the page there is a small interactive arrow linking directly to the top of the page (for easier access to the various panels).

The screenshot shows a web developer tool interface with a 'Source code' tab selected. The code is HTML for a Wikipedia page. Several accessibility error messages are highlighted in yellow:

- Line 47: **SC 4.1.2 - Tech ARIA16 [WCAG 2.1 (A)]** Using aria-labelledby to provide a name for user interface control
- Line 48: **SC 4.1.2 - Tech ARIAS [WCAG 2.1 (A)]** Using WAI-ARIA state and property attributes to expose the state of a user interface component
- Line 56: **SC 1.4.1 - Tech F73 [WCAG 2.1 (A)]** Failure of Success Criterion 1.4.1 due to creating links that are not visually evident without color vision
- Line 62: **SC 1.4.1 - Tech F73 [WCAG 2.1 (A)]** Failure of Success Criterion 1.4.1 due to creating links that are not visually evident without color vision
- Line 63: **SC 1.4.1 - Tech F73 [WCAG 2.1 (A)]** Failure of Success Criterion 1.4.1 due to creating links that are not visually evident without color vision
- Line 64: **SC 1.4.1 - Tech F73 [WCAG 2.1 (A)]** Failure of Success Criterion 1.4.1 due to creating links that are not visually evident without color vision

Figure 6: The Web Developer View

Figure 2 shows the top-most part of the page dedicated to the information about each project, namely the one with summary information on that project. However, in that page also additional pieces of more detailed information are provided, namely:

- **Trend of results (R3):** This appears through a line chart (see Figure 7) showing three lines indicating the number of elements (Y axis) evaluated respectively as errors, warnings and successes for each of the audits (X axis). It is also possible to interactively exclude one or more lines by acting on the respective legends, and also to zoom in/out the graph;
- **Accessibility percentage and evaluation completeness:** it is a bar chart in which two bars are shown for each audit, one for each metric. In a way similar to the previous graph, the user can exclude some data from the graph, and also zoom in/out.

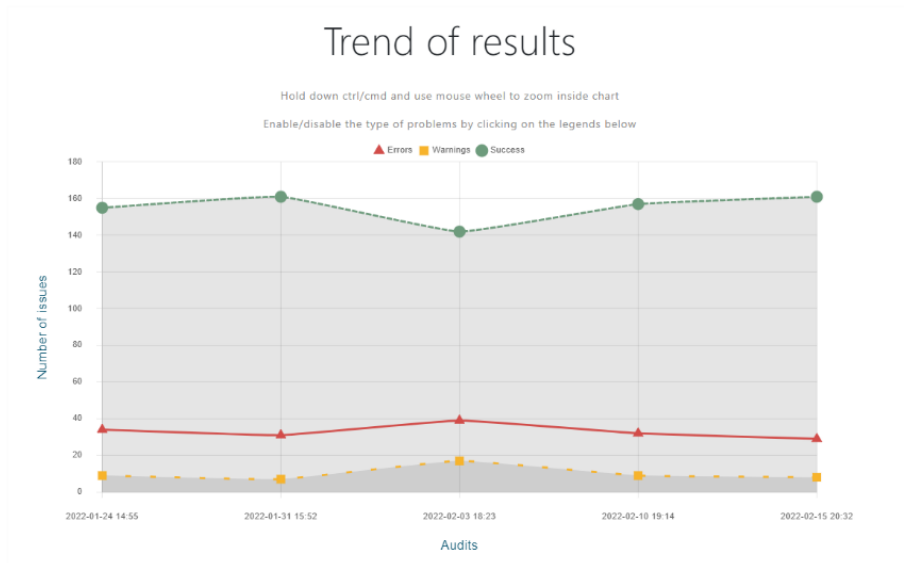


Figure 7: Trend of results for a project

4 THE TOOL ARCHITECTURE

We aim to a solution that is open, supports standard semantic interpretations of the accessibility rules, extensible with limited effort to new guidelines, flexible in defining what to validate (single pages, groups of pages, entire web sites) and in reporting results tailored for different relevant roles, and able to support validation according to the hierarchy of accessibility requirements (principles, guidelines, success criteria, techniques) and for different types of devices. In addition, beyond the possibility to carry out a validation just once, the tool offers the possibility to monitor the accessibility of a web site *over time*, which means automatically scheduling and running accessibility evaluations of a Web site (i.e. audits) at specific times, so that users can follow its evolution over time. To support this, we introduced the possibility to create “projects” in the tool: a project groups together the various audits conducted at different times on a specific group of pages, where each audit corresponds to the evaluation of the web pages belonging to that project at a specific time.

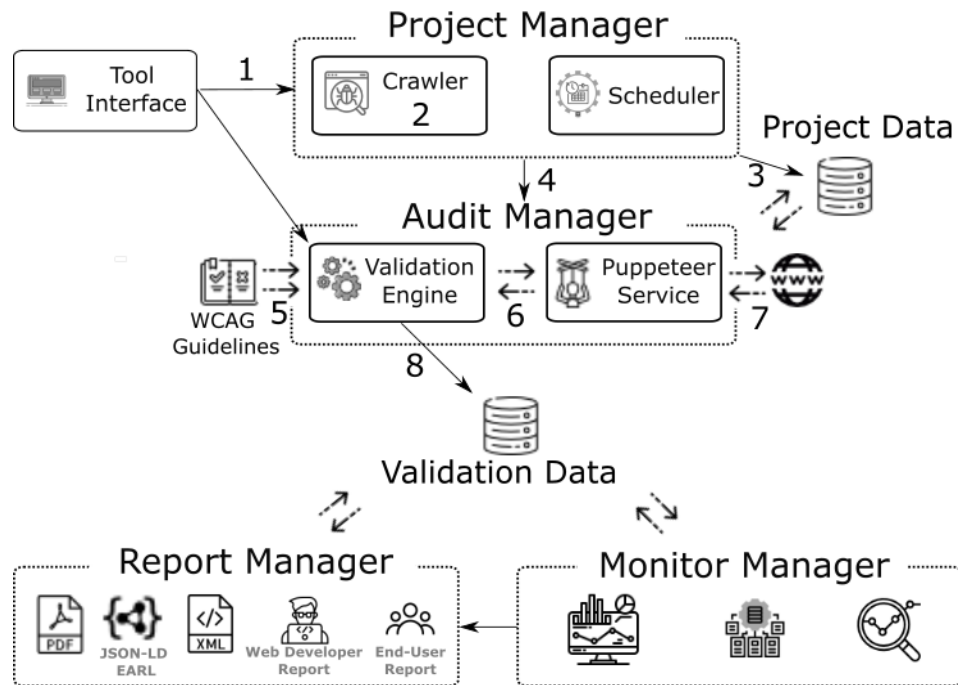


Figure 8: The Tool Architecture

Figure 8 describes the architecture of the new tool to support accessibility evaluation monitoring and server-side rendering. It is mainly composed of:

- a *Project Manager*, which discovers the actual pages of interest for the validation (using a Crawler) according to the parameters specified by the user, and, through the Scheduler, allows for planning future audits in such a way that they will automatically run according to the information specified by the user at project’s creation time;
- an *Audit Manager*, which carries out the actual validation of the selected pages according to a specific set of guidelines (e.g. WCAG 2.1), also managing the validation when such pages have dynamic content;
- a *Monitor Manager*: this module takes the results associated with the audits belonging to a project, aggregates them and provides interactive representations of such information within the tool, to allow the user to analyse such results and have an overview of the evolution of accessibility of the considered web pages.
- a *Report Manager*, which generates validation results in various formats (PDF, Web, EARL), aimed to various types of users.

By using the tool's Web interface, users can specify some parameters (1) that will be used to create different types of projects (i.e. simplified or in-depth). In the case of a “simplified” project (which is the case actually explained in Figure 1), the pages included in the project are those belonging to a specific Web site and identified by a crawler (2) starting from a base URL specified by the user at the project’s creation time. It is worth noting that in the case of an in-depth project the user will specifically indicate the pages that should be involved in the validation. The Project Data such as the crawler configuration and the evaluation settings (WCAG version, target device and user agent, etc.) are stored in a database (3), and will be used by the Project Manager and the Audit Manager modules.

To implement the crawler, we exploited crawler4j³, an open-source multi-thread Web crawler for Java that provides a simple interface for crawling the Web. The crawler needs a base URL that is used as a seed: this is the first page fetched, and from such URL, the crawler will start following the links discovered within the page, and iteratively analyse the newly discovered pages. The crawler takes as input two main configuration parameters provided by users during the project creation: *Crawler depth* and *Maximum Number of Pages* that should be fetched. The depth of crawling describes the extent to which this module discovers the target pages; for example, let us start from the seed page "A", which links to page "B", which links to page "C". Thus, we have the following link structure: A -> B -> C ->[...]; since "A" is a seed page, it will have a depth of 0, "B" will have a depth of 1 and so on.

The crawler implements a "Should Visit" policy to specify whether a given URL should be crawled or not; we configured this policy to ignore URLs associated with CSS, JavaScript code, images, videos, PDF and all the content types that are not immediately relevant for the validation. Moreover, we configured the crawler in such a way that it should not follow URLs belonging to domains different from the one specified in the seed URL.

The URLs of the pages discovered/identified by the Crawler (according to the *Crawler depth* and *Maximum Number of Pages* parameters) at project's creation time are then saved in a database (3), and afterwards they are provided to the Audit Manager, which will actually evaluate them (4). As we will see later on in this section, since the actual pages associated with these URLs might depend on some JavaScripts, another module (Puppeteer) actually retrieves the HTML code currently associated with these URLs in order to perform a specific audit.

In order to support the functionality of monitoring the accessibility over time, we developed a module called Scheduler. The Scheduler is responsible for scheduling a new evaluation of the discovered pages in the time interval specified at project creation time (e.g. each Monday every week). The Validation Engine is part of the Audit Manager: it takes as input (5) the WCAG Guidelines specification defined in an XML-based language, and it exploits an HTML&CSS parser⁴ to select all the page elements considered in the guidelines, and finally it verifies whether such elements respect the accessibility checks. In this regard, the tool has been designed in such a way to facilitate its update in case of changes in the relevant guidelines: indeed, the checking of the guidelines is not hardcoded in the tool but it is performed by the Validation engine, which is able to check any guidelines written in an XML-based language (developed previously). The HTML&CSS parser exploits the parameter (provided by the user) specifying the target device, which defines different viewport dimensions associated with each device type. Starting from such viewport dimensions, the parser then applies the eligible media queries to the DOM elements of the Web page, to identify what content is relevant for the current validation. Thus, the accessibility evaluation for a desktop version can provide a different result from a mobile version evaluation of the same Web site.

The Validation Engine evaluates the dynamic content of Web pages by exploiting the functionalities provided by the Puppeteer Service (6), which implements the Server-Side Rendering -based validation. The tool allows the user to select between two types of validation: static validation, and server-side rendering validation. In the Static Web Page validation, the tool downloads the HTML and the CSS code of the page, and then parses and validates the corresponding DOM. Instead, using the Server-Side Rendering validation, the engine does not parse just the static Web page code: indeed, it exploits the Puppeteer library (<https://developers.google.com/web/tools/puppeteer>) to load the HTML (and CSS) code (7) within a headless version of the Chrome browser, to simulate the page being opened within the user's browser. This provides the validation tool with a more complete page, as its content is also populated with the result of the execution of

³ <https://github.com/yasserg/crawler4j>

⁴ <https://github.com/radkovo/jStyleParser>

the JavaScript code included in the page. The Server-Side Rendering Validation is always used by the tool in case of validations of projects. Instead, when users select a single page validation, they can decide the type of validation to use (either static or server side rendering one). While the static one is faster, the server-side rendering validation can be more complete but it might require more time for its execution. We further extended the Validation Engine to be fully compliant with the EARL outcome specification⁵ by adding the ability to calculate the WCAG techniques that have not been applied because there are no elements that are relevant for their application within the considered page.

The Monitor Manager provides users with interactive representations of the validation results (i.e. trend of results, metrics, most frequent errors). Finally, there is a Report Manager which gives the accessibility results in form of reports provided through different formats (8): through some web interfaces (e.g. the Live Preview for the end-user, the “Web developer” view), but also by providing some downloadable PDF, EARL and XML-based reports.

5 THE USABILITY STUDY

5.1 Organisation of the Test

In order to gather user feedback on the tool and especially about the recently introduced features associated with the monitoring support, a usability test was carried out in which the users had to perform some tasks using the tool, and also answer related questions. We preferred to administer a specific ad-hoc questionnaire rather than a standard usability evaluation scale, to have more focused feedback from users on the usability of the tool.

We recruited 15 users (8 females; average age: 41.6; standard deviation: 12.8), by sending an email in the area in which our research institute belongs, inviting people with specific competencies/skills to participate in the test (i.e. being or having been a web developer or a web commissioner). To describe their profile, users were allowed to specify one or more than one option. In the end, eleven users selected “web developer” as one of their main profiles, five users selected “web commissioner”, two selected “accessibility expert”, three users selected other profiles such as UI designer or researcher. They were asked whether they have used some validation tools before the test: 8 users replied “Yes”, mentioning Siteimprove (3 users), Wave (3 users), Mauve (2 users), W3C tools (2 users), Google Lighthouse (1 user); the other 7 answered “No”.

The test was carried out remotely with the support of a videoconference system, and each test session was video-recorded with users' permission collected in an informed consent that users had to fill in and sign before the test. Some days before the test, users also received some documentation introducing the main aspects of the WCAG guidelines, and of the tool.

Just before starting the actual test, the moderator briefly introduced its goal. Then, the user received the link of the tool and the credentials to use for the test (all the users used the same account), and the link to a Web page which contained the tasks, alternating them with related usability questions. To better follow the test, the participants were asked to share with the moderator the browser's window where the tool was visualised. Since the users had to fill in the questionnaire while running the test, to prevent any influence on their responses, at the beginning the participants were asked to open the questionnaire page on a browser window different from the one that was shared with the moderator. The duration of each test session ranged between half an hour and 1 hour.

⁵ <https://www.w3.org/TR/EARL10-Schema/#OutcomeValue>

5.2 Tasks and Questionnaire

The users had to perform the following tasks, which were intertwined with filling in a related questionnaire section. The tasks were identified in such a way to cover the identified requirements: thus, below, for each task we indicate the associated requirement (R1-R8). However, for requirements R6 and R8, while the tool supports them (i.e. there is a specific "Info" page in which the techniques supported are listed, and the tool is able to support the accessibility validation of dynamic web pages), we have not gathered specific information through the test (by introducing specific tasks or questions) in order not to prolong it too much.

Task1. Single Page Evaluation – Analyse of the various views provided by the tool (R2, R7)

The first task required carrying out a "Single web page evaluation" of a Wikipedia page (https://it.wikipedia.org/wiki/Pagina_principale), using WCAG 2.1 guidelines, Level of Conformance = AAA, server-side rendering validation, and selecting "Desktop - Windows – Edge" as the device/user agent considered for the evaluation. Then, they had to start the validation using the tool, get an overview of the content of the four views provided by it ("Evaluation Summary", "End User", "Live preview" and "Web developer"), and answer the following related questions:

- Q1: By analysing the results, how many distinct types of problems categorised as "error" has the tool identified?
- Q2: Within the "error" category problems, which type of problem occurs most frequently? How many times does it occur?
- Q3: Did you encounter any difficulty in understanding one or more pieces of information included in the views provided by the tool ("Evaluation Summary", "End User", "Live preview" and "Web developer")? If so, what did you find difficult to understand or unclear, and why?

Task2. Analysis of the Live Preview (R1, R4, R7)

After completing the previous task, in the "Live preview" panel, users were asked to filter the results only to the CSS errors affecting the "perceivable" principle (thereby the parameters to set were: issue type=" errors", code type="CSS", principle="Perceivable"). Next, in the "Error List" section, they had to identify the errors associated with the violation of Criterion 1.4.10, SCR34 technique, analyse its first 5 occurrences, also paying attention to how the tool presents the results when the element affected by the error *is* or *is not* visible in the page. Then, they were asked to rate their agreement with the following statements using a 1-5 Likert scale (1= I strongly disagree, 5= I strongly agree):

- S1: When it was possible to view the error within the Web page, the view offered to the user is useful for locating the identified problem
- S2: When the tool was able to highlight the error within the Web page, the tool also provides some information to solve the error. The view offered to the user is useful to solve the identified problem.
- S3: When the occurrence of the error is not visible within the page, the tool provides a link to the "Web developer" view for the analysis of the source code. It is useful to be able to analyse the error directly in the source code.

Then, they had to ask the question:

- Q4: Would you have any suggestions on how to improve the Live Preview, or what to change to improve it?

Task3. Creation of an in-depth Project (R2)

The users had to create a new project of type "In depth", by configuring it with the following URLs and parameters:

URLs:

- Homepage: <https://www.unipi.it/>

- Login: <https://unipi.idp.cineca.it/idp/profile/SAML2/Redirect/SSO?execution=e2s1>
- Forms: <https://unimap.unipi.it/cercapersone/cercapersone.php>
- Accessibility Statement: <https://www.unipi.it/index.php/documenti-ateneo/item/14764>

Parameters:

The guidelines to select are WCAG 2.1, the Level of Conformance is AA, the audit should be repeated just once, and it should be planned for the current day, also selecting “Desktop - Windows – Chrome” for the user agent/device.

After having created the project, an audit automatically started in the tool (as associated with this newly created project): this was because it was planned for the current day (as per the specified parameters). The users had then to answer this question:

- Q5: Have you experienced particular difficulties in creating a project or in understanding the parameters to specify? If so, please indicate where.

Task4. Analysis of a Single Audit (R2)

Users were asked to select the newly created project and, within it, open the audit just started and analyse the information provided by it (within the "Results" and "Page Results" sections): in particular, they could analyse the details of the evaluation results associated with the various pages associated with the project using the links in the "List of Evaluated pages" section. Then, they had to rate their agreement with the following statements (using a 1-5 Likert scale, 1= I strongly disagree, 5= I strongly agree):

- S4: The views presented by the tool and associated with an audit are clear and understandable

Then, they had to answer the following open-ended question:

- Q6: Would you have suggestions on what to change/how to improve the information associated with each audit?

Task5. Analysis of information associated with a project (R3, R5)

Differently from the previous task in which users had to assess the information provided by a *single evaluation* (audit) over a specific set of pages, in this task users had to assess the information that the tool provides after evaluating a set of pages *over time* (i.e. by carrying out audits in different times). To allow users to perform this task during the test, in the tool’s workspace of the participant’s account, we made available a project containing data that simulated various audits repeated in different times of the year on a specific group of web pages. Users had to open and analyse this project, explore the various sections of the page ("Trend of results" / “Accessibility percentage and evaluation completeness" / “Most frequent errors”) and they have to answer the following questions.

- Q7: In which audit was the Evaluation Completeness highest?
- Q8: Which audit had the least number of successes?
- Q9: What is the most frequent type of error found in the last audit?
- Q10: If you have experienced any difficulty in understanding one or more result included in the various views associated with the project, or in interacting with one of the graphs displayed, could you specify in what situation?
- Q11: Do you have any suggestions on what you would change to improve the info associated with each project?

They had also to rate their agreement level to this statement:

- S5: The views presented in the tool and associated with a project are clear and understandable

After this task, the questionnaire presented a number of final, open-ended questions for evaluating the overall experience of using the tool. The questions were:

- Q12: What are the three aspects of the tool you liked the most?
- Q13: What are the three aspects of the tool you liked the least?
- Q14: In general, do you think that that the tool clearly gives all the information necessary to get an overview of the accessibility problems of a site, to be able to monitor its accessibility over time?
- Q15: Is there any information that you think would be useful but that you have not found in the tool?
- Q16: Do you have any further suggestion for improving the tool?

5.3 Results

In this section we report the users' responses to the questions indicated above.

For question Q1, users had to *indicate how many types of problems the tool identified within the "error" category*: the vast majority of them (13 out of 15) correctly replied to this question. For question Q2 users had to *indicate, within the "error" category, which type of problem occurred most frequently*: 11 (out of 15) correctly answered. Some of those who provided an incorrect answer just misunderstood the question statement, which included the "category" word: as a consequence, they looked in the tool where the errors are grouped "by category". Still in Q2, users had to indicate *how many times the most recurrent problem occurred*. Nine users correctly replied, three did not provide an answer, the remaining three provided a wrong one (which was connected to the incorrect answer provided just before).

For Q3, ten users explicitly reported that the information included in the various views was clear, and they had no particular difficulty in understanding it. As a suggestion, one of them recommended giving more visibility and structure to the left-hand panel of the "Live Preview", also expressing appreciation for the fact that in the "Web developer" view the identification and categorisation of the issues (i.e. in warning and errors) was easy to follow, thanks to the different colours used. The remaining five reported the following comments (some provided more than one suggestion). Two reported difficulties in understanding the connections between the information provided in the "Evaluation Summary" (which shows the issues categorised in terms of techniques that resulted erroneous, successful, warnings, and not applicable) and the information available in the "End user" tab, which details the number of occurrences (categorised according to principles) for each technique of type "warning" or "error". In particular, one of them suggested adding further explanation about how to read these connected views. Two users suggested providing the possibility of ordering the column showing the number of occurrences within the "End User" view, to have more readily available the information about the most recurring error/warning. One participant suggested changing the current representation (which is a pie chart) used to show the occurrences of errors/warnings grouped by principles, because in some cases the slices may be difficult to perceive. A user reported some difficulties with the (standard W3C) acronyms (i.e. "SCR34", "G212") used in the tool to identify the various techniques/success criteria.

Next (Q4), the participants had to *provide suggestions on how to improve the Live Preview, or about what to change to improve it*. Six users judged the view as being already very comprehensible and useful. As for the others: i) three commented about the part supporting filtering: one suggested distinguishing more clearly the filtering part from the one presenting the results, also using different colours. Two suggested adding a button to explicitly support applying the filtering criteria (currently they are automatically applied whenever the user changes them); ii) five users gave suggestions

about the type of visualisation provided by the tool: one suggested changing the eye-shaped icon of the button supporting the localisation of the selected error within the Web page when the issue *cannot* be visualised in the page, by using a more intuitive one (e.g. through a strikethrough eye-based icon, to highlight that the error cannot be seen); two users said that, when an issue is not visible, the link to the “Web developer” view was judged as not particularly evident; another user did not find it intuitive to use mouse-over to trigger the appearance of the tip. Another user did not realise that the eye-shaped icon is interactive.

For the next question (Q5), ten users declared that there was no particular difficulty in creating a project. One user just suggested better explaining the various sections. Five users provided additional remarks to better highlight the various parts in the project creation form, also indicating the mandatory fields within it.

For Q6, eight users declared not having any further suggestions to improve the information associated with each audit. One user suggested adding another graph visualising the number of erroneous techniques, possibly superimposed on the graph visualising the number of errors in the Page Results section (which visualises the number of error/warnings occurrences found in the page using a bar chart). The same user suggested better explaining some expressions (e.g. the “Average number of techniques” in the “Results” section did not clearly indicate which elements the average was calculated on). One user suggested making the graphs even more interactive (e.g. further details could be shown after clicking on a bar in the “Page Results”).

Users also had to indicate some results provided by the tool on the evaluation of multiple pages (corresponding to a “project”) over multiple audits. This was done to understand whether such information was actually comprehensible. In particular, 14 users replied correctly to Q7 (*in which audit was the Evaluation Completeness highest?*), 1 incorrectly; all the users correctly answered Q8 (the audit with the least number of successes); 14 users replied correctly to Q9 (most frequent type of error in the last audit), 1 user incorrectly. So, we can infer that the key information associated with a project is generally understandable by users.

Furthermore, five users explicitly stated not having had any *difficulty with the info associated with each project* (Q10), while the others took this opportunity to suggest further improvements. Three users expressed concerns about the “Most frequent errors” section. In particular, they judged that the way such information was described could be improved. Indeed, to indicate respectively the most frequent/the second-most frequent/the third-most frequent type of error, we used the (probably too compact) expressions “first/second/third”, which resulted a bit ambiguous for some of them (e.g. one user thought it referred to the audits). Two users reported some slight difficulties in identifying the lowest number of successes. As for the line chart “Trend of results”, one user suggested deleting the greyish area under the line and keeping just the lines; four users suggested slightly increasing the size of the circles visualised in that graph; another user suggested replacing this graph with a histogram, to further enhance its clarity.

As for Q11, nine users explicitly said that they had no particular *suggestion to improve the information associated with each project*. One suggested deleting the (decorative) icons currently shown under the title “List of Evaluated pages” as they do not completely reflect the actual information that is shown in that part of the page, and to use consistently the colours that are shared by other graphs (e.g. the orange colour is generally used in the tool to indicate a warning). In the section “Accessibility percentage and evaluation completeness” one user suggested having just two series of bars (one for

each of the metrics considered: accessibility percentage and evaluation completeness) instead of two bars for each audit (one bar for the first metric and one for the second). Even though currently it is possible to hide a bar series associated with a specific metric, the user interestingly pointed out that a comparison between the bars referring to each metric could be useful to have a more immediate manner in some cases. This would suggest that, in order to allow an effective use in a diverse range of situations, accessibility validation tools should support as much flexibility as possible.

Finally, in the last part of the questionnaire, we asked users to provide some more general feedback about the overall experience. Among *the three things they liked the most about the tool* (Q12), four said that the information provided is well/clearly structured. The possibility to perform multiple audits and see the trends over time was mentioned also four times. The clearness/interactiveness of graphs was mentioned four times. Three users liked the possibility to see the issues in the rendered page and in the code, while three users liked the graphical style of the tool. Three users mentioned the clarity of the tool, two its intuitiveness, two appreciated that the tool provides objective results and in a detailed manner; two the fact that the tool provides results promptly; two liked the easy connection with related W3C documentation to help users solve an issue. Additional aspects were: the completeness of the information, the interactivity of the tool, the fact that it allows for analysing different aspects of accessibility, the possibility for even a non-expert user to carry out an analysis. They also mentioned the easy to understand results, the Web developer view, the Live Preview, the possibility to assess a single page or go deeper to assess an entire Web site, the versatility/flexibility of the tool, and the fact that it is responsive. One user also highlighted that the tool provides the user with a clear workflow to create a new project.

Below we report a table in which we summarise the main descriptive statistics metrics for the Likert scale ratings associated with statements S1-S5.

	Min	Max	Median	Mean	St. Dev.
S1	4	5	4	4.5	0.5
S2	3	5	4	4.4	0.6
S3	4	5	5	4.5	0.5
S4	2	5	4	4.1	0.7
S5	3	5	4	4.1	0.6

Table 1: Summary of *min, max, range, median, mean and st. dev. of Likert scale ratings associated with statements S1-S5*

While this table shows generally positive user appreciation of the tool, the minimum value (2) corresponds to S4 (“The views presented by the tool and associated with an audit are clear and understandable”), which was provided by an accessibility expert, who suggested using different types of visualisations (e.g. histograms) for highlighting the errors and warnings detected in each page whereas currently the tool shows a line chart (Trend of Results diagram).

We also asked about *the three things they liked least about the tool* (Q13). Four users declared not being able to identify any aspect in this regard. Three users mentioned some issues related to the legends associated with some of the graphs: while they appreciated the possibility of customising a part of a visualisation (e.g. hiding a portion) to better focus on specific aspects of interest. They reported some usability issues with the way in which it was currently supported. For instance, one user complained about having to select a legend to hide the related part in a graph, another user found the legends too small (compared to the size of the graph), another suggested enhancing the contrast of labels used besides

some legends (to be more visible), also noting that some labels are too spaced apart. Two users said that some functionalities are not easily visible, mentioning the link supporting the visualisation of errors within the “Web Developer” view; the same applies to the “View project” / “View audit” links. Two users mentioned adding further information to the various sections in the tool, as it is not immediately clear which aspects they refer to.

All the users replied affirmatively to question Q14 (overview of the accessibility problems of a site, to be able to monitor its accessibility over time). While six users did not provide any further detail, five users further highlighted some aspects they found important: the possibility to create projects is best suited to the kind of analysis supported, the fact that it is a tool usable by both web developers and non- web developers, it provides further explanations when needed, it offers a list of errors ordered by some priority. One user highlighted that, while it is certainly a useful tool, sometimes it uses language that can result too specialised.

Finally (Q15-Q16), we asked whether, in user’s opinion, there was any useful information that they have not found in the tool or whether they had any suggestions to improve the tool. All the users declared not be able to identify any further information that the tool currently does not provide. One user mentioned that it could add some information about the Web page loading time. Another mentioned that the tool could provide some information about non-working links.

5.4 Discussion

The results that we gathered through this test were overall encouraging. First, since the participants of the test had different kinds of profiles (both technical and non-technical) and they were generally able to complete the various tasks in a satisfactorily manner, this result would confirm that the tool can easily support various types of stakeholders, which was one of our main requirements (R1). Also, by analysing the results gathered in association with Task1 (which dealt with a single page evaluation), Task 3 (creation of a project), and Task 4 (analysis of a single audit), overall it seems that the users satisfactorily used the tool for conducting both the evaluation of a single page and that of multiple pages (also one of the requirements, see R2). Thus, the tool seems to be able to support both the needs of those who have a specific focus on a particular web page for their own goals (e.g. a personal home page), but also the needs of organisations that require a more comprehensive, site-wide evaluation of the accessibility of their web applications. The answers provided to Q7-Q11 and the level of user’s agreement to statement S5 also show that, while the tool can be improved in describing the information it provides, it is overall able to offer to its users key information to monitor how the accessibility of a site/group of pages can vary over time both in a detailed manner (i.e. in terms of number of errors/warning detected in different audits) and in terms of more general metrics that summarise the level of accessibility of a site over time (see requirements R3 and R5), as well as addressing the needs of different types of stakeholders (e.g. a summative metric can likely be of more interest for non-technical users). Mainly based on the data gathered in the ratings to statements S1, S2, S3, and the answers to question Q4, while also considering more general comments, the participants seem to appreciate the level of support that the tool offers for both localising (requirement R4) and solving (requirement R7) the errors identified during an accessibility evaluation, especially through the provided validation results views (e.g. Live Preview, Web Developer View) able to target different types of stakeholders, as well as the provision of direct link/reference to the relevant W3C documentation. However, further work should be done to provide more specific indications about how to solve some of the issues that the tool detects (even by supporting user’s manual intervention).

To sum up, the participants seemed to be quite satisfied with the tool, which is particularly encouraging especially considering that the vast majority of them used it for the first time, none had used the new functionalities (e.g. the monitoring support) previously, and no familiarisation phase with the tool was carried out before the test. On the one hand, they especially appreciated the fact that it allows monitoring the evolution of different audits over time and found clear the information provided, including the graphs; they also liked that the tool provides objective results, and in a prompt and detailed manner. On the other hand, some of them suggested improving aspects connected with the clarity of the tool (e.g. improve some legends using a less specialised language), and also make more visible some of the provided features.

6 CONCLUSIONS AND FUTURE WORK

In this paper, we have presented a set of requirements that should characterise automatic validation tools able to address the emerging needs derived from the increasing demand for accessible Web sites and the evolution of the Web technologies. We have shown how such design dimensions can be supported and implemented in a tool, which is thus able to provide users with monitoring functionalities, clear indications of its capabilities, and validation of dynamic content. We have also reported on a first user test that provided positive feedback and suggestions for small adjustments, such as in the choice of the graphs for visualizing the validation results.

The tool is available for open access, and future work will be dedicated to modifying it by taking into considerations the suggestions provided by users, adding functionalities that will allow accessibility experts to integrate their manual validations with those automatically generated (to decide the correctness of the elements that cannot be validated automatically), and produce accessibility reports focused on specific disabilities or application areas.

REFERENCES

- Abascal, J., Argue, M., & Valencia, X. (2019). Tools for Web accessibility evaluation. *Web Accessibility* (pp. 479-503). London: Springer.
- Siddikjon Gaibullojonovich Abduganiev. 2017. Towards automated Web accessibility evaluation: a comparative study. *Int. J. Inf. Technol. Comput. Sci.(IJITCS)* 9, 9 (2017), 18–44.
- Shadi Abou-Zahra. 2017. Evaluation and Report Language (EARL). Retrieved February 2, 2017 from <https://www.w3.org/TR/EARL10-Schema/#OutcomeValue>
- Beirekdar, A., Vanderdonck, J., & Noirhomme-Fraiture, M. (2002). Kwaresmi—Knowledge-based Web Automated Evaluation with REconfigurable guidelineS optimisation. (Springer, Ed.) *DSV-IS*, 2545, 362-376.
- Beirekdar A., Keita M., Noirhomme M., Randolet F., Vanderdonck J., Mariage C. (2005) Flexible Reporting for Automated Usability and Accessibility Evaluation of Web Sites. In: Costabile M.F., Paternò F. (eds) *Human-Computer Interaction - INTERACT 2005*. INTERACT 2005. Lecture Notes in Computer Science, vol 3585. Springer, Berlin, Heidelberg
- Brajnik, G., Yesilada, Y., & Harper, S. (2012). Is accessibility conformance an elusive property? A study of validity and reliability of WCAG 2.0. *ACM Transactions on Accessible Computing (TACCESS)*, 4(2), 1-28.
- Broccia, B., Manca, M., Paternò, F., Pulina, F. Flexible Automatic Support for Web Accessibility Validation. *Proc. ACM Hum. Comput. Interact.* 4(EICS): 83:1-83:24 (2020)
- Burkard, A., Zimmermann, G., and Schwarzer, B. 2021. Monitoring Systems for Checking Websites on Accessibility. *Frontiers in Computer Science* 3 (2021), 2.
- EU Commission. (2016, October 26). Directive (EU) 2016/2102 of the European Parliament and of the Council. Retrieved from <https://eur-lex.europa.eu/https://eur-lex.europa.eu/eli/dir/2016/2102/oj>
- Fernandes, N., Kaklanis, N., Votis, K., Tzovaras, D., & Carriço, L. (2014). An analysis of spersonalised Web accessibility. *Proceedings of the 11th Web for All Conference* (p. 19). ACM.
- Frazão, T., and Duarte, C. 2020. Comparing accessibility evaluation plug-ins. In *Proceedings of the 17th International Web for All Conference (W4A '20)*. Association for Computing Machinery, New York, NY, USA, Article 20, 1–11. DOI:<https://doi.org/10.1145/3371300.3383346>
- Fuertes, J. L., González, R., Gutiérrez, E., & Martínez, L. (2009). Hera-FFX: a Firefox add-on for semi-automatic Web accessibility evaluation. *Proceedings of the 2009 International Cross-Disciplinary Conference on Web Accessibility (W4A)* (pp. 26-35). ACM.
- Gay, G., and Qi Li, C. 2010. AChecker: open, interactive, customisable, Web accessibility checking. In *Proceedings of the 2010 International Cross-Disciplinary Conference on Web Accessibility (W4A)*. 1–2.

- Gulliksen, J., Von Axelson, H., Persson, H., & Göransson, H. (2010). Accessibility and public policy in Sweden. *Interactions*, 17(3), 26-29.
- Ivory, MY, Mankoff, J., and Le, A. 2003. Using automated tools to improve Web site usage by users with diverse abilities. *Human-Computer Interaction Institute* (2003), 117.
- Lazar, J., & Olalere, A. (2011). Investigation of best practices for maintaining section 508 Compliance in US federal Web sites. *International Conference on Universal Access in Human-Computer Interaction* (pp. 498-506). Berlin: Springer.
- Miniukovich, A., Scaltritti, M., Sulpizio, S., and De Angeli, A. 2019. Guideline-Based Evaluation of Web Readability. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems (CHI '19)*. Association for Computing Machinery, New York, NY, USA, Paper 508, 1–12. DOI:<https://doi.org/10.1145/3290605.3300738>
- Mirri, S., Muratori, L. A., & Salomoni, P. (2011). Monitoring accessibility: large scale evaluations at a geo political level. *The proceedings of the 13th international ACM SIGACCESS conference on Computers and accessibility* (pp. 163-170). New York: ACM.
- Moreno, L., Alarcon, R., Segura-Bedmar, I., and Martínez, P. 2019. Lexical simplification approach to support the accessibility guidelines. In *Proceedings of the XX International Conference on Human Computer Interaction (Interaccion '19)*. Association for Computing Machinery, New York, NY, USA, Article 14, 1–4. DOI:<https://doi.org/10.1145/3335595.3335651>
- Molineró, AM., Kohun, FG, and Morris, R. 2006. Reliability in Automated Evaluation Tools for Web Accessibility Standards Compliance. *issues in Information Systems* 7, 2 (2006), 218–222.
- Mucha, J., Snaprud, M., & Nietzio, A. (2016). Web page clustering for more efficient website accessibility evaluations. *International Conference on Computers Helping People with Special Needs* (pp. 259-266). Springer.
- Nietzio, A., Eibegger, M., Goodwin, M., & Snaprud, M. (2011). Towards a score function for WCAG 2.0 benchmarking. *Proceedings of W3C Online Symposium on Website Accessibility Metrics*. Retrieved from <https://www.w3.org/WAI/RD/2011/metrics/paper11>
- Pädure, M., and Pribeau, C. 2019. Exploring the differences between five accessibility evaluation tools. (2019).
- Parvin, P., Palumbo, V., Manca, M., Paternò, F. 2021. The Transparency of Automatic Accessibility Evaluation Tools. In *Proceedings of the 18th International Web for All Conference (W4A '21)*, April 19–20, 2021, Ljubljana, Slovenia. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3430263.3452436>
- Paternò F., Pulina F., Santoro C., Gappa H., Mohamad Y. (2020) Requirements for Large Scale Web Accessibility Evaluation. In: Miesenberger K., Manduchi R., Covarrubias Rodriguez M., Peñáz P. (eds) *Computers Helping People with Special Needs. ICCHP 2020. Lecture Notes in Computer Science*, vol 12376. Springer, Cham. https://doi.org/10.1007/978-3-030-58796-3_33
- Paternò F., Schiavone A., The role of tool support in public policies and accessibility. *ACM Interactions* 22(3): 60-63 (2015)
- Pelzetter, J. A Declarative Model for Web Accessibility Requirements and its Implementation. *Frontiers Comput. Sci.* 3: 605772 (2021)
- Petrie, H., King, N., Velasco, C., Gappa, H., Nordbrock, G.: The usability of accessibility evaluation tools. In: Stephanidis, C. (ed.) *UAHCI 2007. LNCS*, vol. 4556, pp. 124–132. Springer, Heidelberg (2007). https://doi.org/10.1007/978-3-540-73283-9_15
- Power, C., Freire, A., Petrie, H., & Swallow, D. (2012). Guidelines are only half of the story: accessibility problems encountered by blind users on the Web. *Proceedings of the SIGCHI conference on human factors in computing systems* (pp. 433-442). ACM.
- Schiavone A., Paternò F., An extensible environment for guideline-based accessibility evaluation of dynamic Web applications, *Universal Access in the Information Society*, Springer Verlag, 14(1): 111-132, 2015.
- Salehnamadi, N., Alshayban, A., Lin, JW, Ahmed, I., Branham, S., and Malek, S. 2021. Latte: Use-Case and Assistive-Service Driven Automated Accessibility Testing Framework for Android. In *CHI Conference on Human Factors in Computing Systems (CHI '21)*, May 8– 13, 2021, Yokohama, Japan. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3411764.3445455>
- Vigo, M., Brown, J., and Conway, V. 2013. Benchmarking Web accessibility evaluation tools: measuring the harm of sole reliance on automated tests. In *Proceedings of the 10th International Cross-Disciplinary Conference on Web Accessibility*. 1–10.
- Yesilada, Y., Brajnik, G., Vigo, M., Harper, S.: Exploring perceptions of Web accessibility: a survey approach. *Behav. Inf. Technol.* 34(2), 119–134 (2015)
- W3C WAETL, Web Accessibility Evaluation Tools List, <https://www.w3.org/WAI/ER/tools/> (last accessed 20 January 2022).