

Improving the accessibility of scientific documents

Current state, user needs, and a system solution to enhance scientific PDF accessibility for blind and low vision users

LUCY LU WANG*, Allen Institute for AI

ISABEL CACHOLA*†, The Johns Hopkins University

JONATHAN BRAGG, Allen Institute for AI

EVIE YU-YEN CHENG, Allen Institute for AI

CHELSEA HAUPT, Allen Institute for AI

MATT LATZKE, Allen Institute for AI

BAILEY KUEHL, Allen Institute for AI

MADELEINE VAN ZUYLEN, Allen Institute for AI

LINDA WAGNER, Allen Institute for AI

DANIEL S. WELD, Allen Institute for AI and University of Washington

The majority of scientific papers are distributed in PDF, which pose challenges for accessibility, especially for blind and low vision (BLV) readers. We characterize the scope of this problem by assessing the accessibility of 11,397 PDFs published 2010–2019 sampled across various fields of study, finding that only 2.4% of these PDFs satisfy all of our defined accessibility criteria. We introduce the SciA11y system to offset some of the issues around inaccessibility. SciA11y incorporates several machine learning models to extract the content of scientific PDFs and render this content as accessible HTML, with added novel navigational features to support screen reader users. An intrinsic evaluation of extraction quality indicates that the majority of HTML renders (87%) produced by our system have no or only some readability issues. We perform a qualitative user study to understand the needs of BLV researchers when reading papers, and to assess whether the SciA11y system could address these needs. We summarize our user study findings into a set of five design recommendations for accessible scientific reader systems. User response to SciA11y was positive, with all users saying they would be likely to use the system in the future, and some stating that the system, if available, would become their primary workflow. We successfully produce HTML renders for over 12M papers, of which an open access subset of 1.5M are available for browsing at scia11y.org.

CCS Concepts: • **Human-centered computing** → **Empirical studies in accessibility**; **Accessibility systems and tools**; *HCI design and evaluation methods*; *Accessibility design and evaluation methods*.

Additional Key Words and Phrases: accessibility, accessible reader, scientific documents, blind and low vision readers, science of science, user study

*Denotes equal contribution

†Work done while at the Allen Institute for AI

Authors' addresses: Lucy Lu Wang, lucyw@allenai.org, Allen Institute for AI, Seattle, WA, 98103; Isabel Cachola, icachola@cs.jhu.edu, The Johns Hopkins University, Baltimore, MD, 21218; Jonathan Bragg, jbragg@allenai.org, Allen Institute for AI, Seattle, WA, 98103; Evie Yu-Yen Cheng, eviec@allenai.org, Allen Institute for AI, Seattle, WA, 98103; Chelsea Haupt, chealseah@allenai.org, Allen Institute for AI, Seattle, WA, 98103; Matt Latzke, mattl@allenai.org, Allen Institute for AI, Seattle, WA, 98103; Bailey Kuehl, baileyk@allenai.org, Allen Institute for AI, Seattle, WA, 98103; Madeleine van Zuylen, madeleinev@allenai.org, Allen Institute for AI, Seattle, WA, 98103; Linda Wagner, lindaw@allenai.org, Allen Institute for AI, Seattle, WA, 98103; Daniel S. Weld, danw@allenai.org, Allen Institute for AI, Seattle, WA, 98103, University of Washington, Seattle, WA, 98103.

1 INTRODUCTION

Scientific literature is most commonly available in the form of PDFs, which pose challenges for accessibility [6, 34]. When researchers, students, and other individuals who are blind or low vision (BLV) interact with scientific PDFs through screen readers, the availability of document structure tags, labeled reading order, labeled headers, and image alt-text are necessary to facilitate these interactions. However, these features must be painstakingly added by authors using proprietary software tools, and as a result, are often missing from papers. Low vision or dyslexic readers who interact with PDFs through screen magnification or text-to-speech may also find the complexity of certain academic paper PDF formats challenging, e.g., non-linear layout can interrupt the flow of text in a magnifying tool. Inaccessible paper PDFs can lead to high cognitive overload, frustration, and abandonment of reading for BLV readers.

Unfortunately, we find that the majority of scientific PDFs lack basic accessibility features. We estimate based on a sample of 11,397 PDFs from multiple fields of study that only around 2.4% of paper PDFs released in the last decade satisfy all of the aforementioned accessibility requirements. Accessibility challenges for academic PDFs are largely due to three factors: (1) the complexity of the PDF file format, which make it less amenable to certain accessibility features, (2) the dearth of tools, especially non-proprietary tools, for creating accessible PDFs, and (3) the dependency on volunteerism from the community with minimal support or enforcement [6]. The intent of the PDF file format is to support faithful visual representation of a document for printing, a goal that is inherently divergent from that of document representation for the purposes of accessibility. Though some professional organizations like the Association for Computing Machinery (ACM) have encouraged PDF accessibility through standards and writing guidelines,¹ uptake among academic publishers and disciplines more broadly has been limited.

While policy changes help, the fact remains that most academic PDFs produced today, and historically, are inaccessible, yet remain as the dominant way to read those papers. A long-range solution will necessitate buy-in from multiple stakeholders—publishers, authors, readers, technologists, granting agencies, and the like. But in the interim, there are technological solutions that can be offered as a sort of “band-aid” to the problem. We use this paper to offer an in-depth qualitative and quantitative description of the problem as it stands, and to introduce one such technological solution: the SciA11y system that automatically extracts semantic information from paper PDFs and re-renders this content in the form of an accessible HTML document. Though the process is imperfect and can introduce errors, we demonstrate the ability of the rendered HTMLs to reduce cognitive load and facilitate in-paper navigation and interactions for BLV users.

The goals and contributions of this paper are three-fold:

- (1) We characterize the state of academic-paper PDF accessibility by estimating the degree of adherence to accessibility criteria for papers published in the last decade (2010–2019), and describe correlations between year, field of study, PDF typesetting software, and PDF accessibility.
- (2) We propose an automated approach for extracting the content of academic PDFs and displaying this content in a more accessible HTML document format. We build a prototype that re-renders 12 million PDFs in HTML, and describe the design decisions, features, and quality of the renders (assessed as faithfulness to the source PDF). We perform expert grading of the rendered HTML and report an error analysis. A demo of our system is available at scia11y.org, which makes available 1.5M HTML renders of open access PDFs.
- (3) We conduct an exploratory user study with six BLV scholars to better understand the challenges they experience when reading academic papers and how our proposed tool might augment their current workflow. During

¹<https://www.acm.org/publications/authors/submissions>

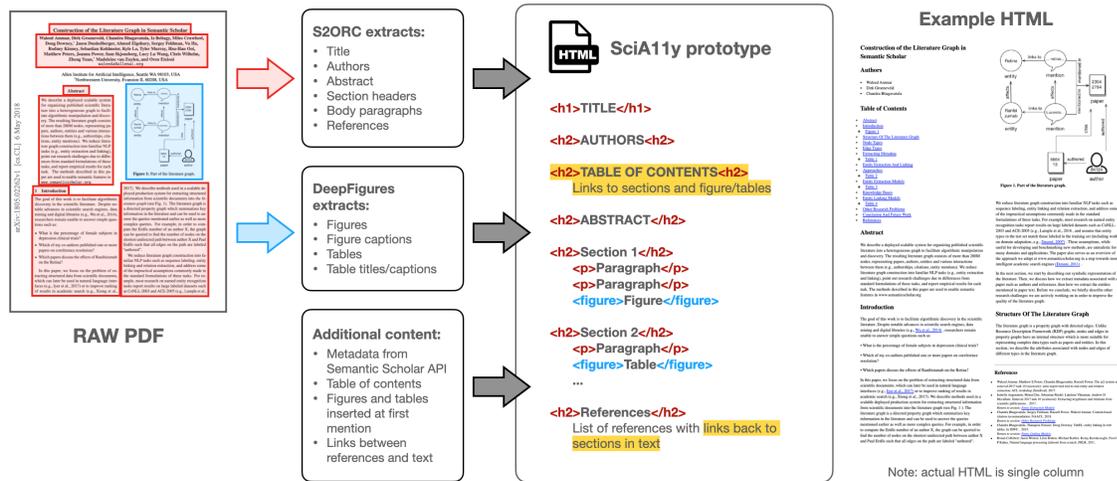


Fig. 1. A schematic for creating the SciA11y HTML render from a paper PDF. Starting with the raw two-column PDF on the left, S2ORC [24] is used to extract title, authors, abstract, section headers, body text, and references. S2ORC also identifies links between inline citations and references to figures and table objects. DeepFigures [43] is used to extract figures and tables, along with their captions. The output of these two models are merged with metadata from the Semantic Scholar API. Heuristics are used to construct a table of contents, to insert figures and tables in the appropriate places in the text, and to repair broken URLs. We add HTML headers as illustrated (header tags for sections, paragraph tags for body text, and figure tags for figures and tables); highlighted components (table of contents and links in references) are not in the PDF and novel navigational features that we introduce to the HTML render. An example HTML render of parts of a paper document is show to the right (actual render is single column, which is split here for presentation).

the study, we ask users to interact with the prototype and offer feedback for its improvement. We perform open coding of interviews to identify existing reading challenges, coping mechanisms, as well as positive and negative responses to prototype features. We summarize the findings of this user study into a set of design recommendations.

Our analysis reveals that PDF accessibility adherence is low across all fields of study. Of the five accessibility criteria we assess, only 2.4% of the PDFs we assess demonstrate full compliance. Though compliance for several criteria seems to be increasing over time, author awareness and contribution to accessibility remains low, as Alt-text has the lowest compliance of the five criteria at between 5–10% (Alt-text is the only criterion of the five that *requires* author intervention in all cases using current tools). We also find that typesetting software is strongly associated with accessibility compliance, with LaTeX and publishing software like Arbortext APP producing low compliance PDFs, while Microsoft Word is generally associated with higher compliance.

To offset the reading challenges of inaccessible papers for BLV researchers, we propose and test the SciA11y system for rendering academic PDFs into accessible HTML documents. As shown in Figure 1, our prototype integrates several machine learning text and vision models to extract the structure and semantic content of papers. The content is represented as an HTML document with headings and links for navigation, figures and tables, as well as other novel features to assist in document structure understanding. Our evaluation of the SciA11y system identifies common classes of extraction problems, and finds that though many papers exhibit some extraction errors, the majority (55%) have no major problems that impact readability, and another 32% have only some problems that impact readability.

Through our user study, we identify numerous challenges faced by BLV users when reading paper PDFs, including some that affect the whole document or limit navigation, and many that affect the ability of the reader to understand text or various elements of a paper like math content or tables. Responses to SciA11y were positive; participants especially liked navigation features such as headings, the table of contents, and bidirectional links between inline citations and references. Of the extraction errors in SciA11y, missed or incorrectly extracted headings were the most problematic, as these impact the user’s ability to navigate between sections and fully trust the system. All users reported being likely to use the system in the future. When asked how the system might be integrated into their workflow, one participant replied “I think it would become the workflow.” Another participant said, “for inaccessible PDFs, this is life-changing.” We condense these findings into a set of recommendations for designing and engineering accessible reading systems (Section 6.3). Most importantly, documents should be structured to match a reader’s mental model, objects should be properly tagged, and care should be taken to reduce the reader’s cognitive load and increase trust in the system. Features that emulate the external memory that visual layout provides to sighted users can be especially beneficial.

This paper is organized as follows. Following a description of related work in Section 2, we first provide a meta-scientific analysis of the current state of academic PDF accessibility in Section 3. In Section 4, we document our pipeline for converting PDF to HTML and describe the SciA11y prototype for rendering papers. An evaluation of HTML render quality and faithfulness is provided in Section 5. Section 6 describes our user study and findings. We recognize that no PDF extraction system is perfect, and many open research challenges remain in improving these systems. However, based on our findings, we believe SciA11y can dramatically improve screen reader navigation of most papers compared to PDFs, and is well-positioned to assist BLV researchers with many of their most common reading use cases. Our hope is that a system such as SciA11y can improve BLV researcher access to the content of academic papers, and that these design recommendations can be leveraged by others to create better, more faithful, and ultimately more usable tools and systems for scholars in the BLV community.

2 RELATED WORK

Accessibility is an essential component of computing, which aims to make technology broadly accessible to as many users as possible, including those with differing sets of abilities. Improvements in usability and accessibility falls to the community, to better understand the needs of users with differing abilities, and to design technologies that play to this spectrum of abilities [48]. In computing, significant strides have been made to increase the accessibility of web content. For example, various versions of the Web Content Accessibility Guidelines (WCAG) [8, 10] and the in-progress working draft for WCAG 3.0,² or standards such as ARIA from the W3C’s Web Accessibility Initiative (WAI)³ have been released and used to guide web accessibility design and implementation. Similarly, positive steps have been made to improve the accessibility of user interfaces and user experience [5, 35, 36, 46], as well as various types of media content [19, 29, 32].

We take inspiration from accessibility design principles in our effort to make research publications more accessible to users who are blind and low vision. Blindness and low vision are some of the most common forms of disability, affecting an estimated 3–10% of Americans depending on how visual impairment is defined [18]. BLV researchers also make up a representative sample of researchers in the United States and worldwide. A recent Nature editorial pushes the scientific community to better support researchers with visual impairments [41], since existing tools and resources can be limited. There are many inherent accessibility challenges to performing research. In this paper, we engage with one of these challenges that affects all domains of study, accessing and reading the content of academic publications.

²<https://www.w3.org/TR/wcag-3.0/>

³<https://www.w3.org/WAI/standards-guidelines/aria/>

BLV users interact with papers using screen readers, braille displays, text-to-speech, and other assistive tools. A WebAIM survey of screen reader users found that the vast majority (75.1%) of respondents indicate that PDF documents are very or somewhat likely to pose significant accessibility issues.⁴ Most papers are published in PDF, which is inherently inaccessible, due in large part to its conflation of visual layout information with semantic content [6, 34]. Bigam et al. [6] describe the historical reasons we use PDF as the standard document format for scientific publications, as well as the barriers the format itself presents to accessibility. Prior work on scientific accessibility have made recommendations for how to make PDFs more accessible [11, 38], including greater awareness for what constitutes an accessible PDF and better tooling for generating accessible PDFs. Some work has focused on addressing components of paper accessibility, such as the correct way for screen readers to interpret and read mathematical equations [1, 4, 16, 17, 26, 44, 45], describe charts and figures [12–14], automatically generate figure captions [9, 37], or automatically classify the content of figures [21]. Other work applicable to all types of PDF documents aims to improve automatic text and layout detection of scanned documents [31] and extract table content [15, 39]. In this work, we focus on the issue of representing overall document structure, and navigation within that structure. Being able to quickly navigate the contents of a paper through skimming and scanning is an essential reading technique [28], which is currently under-supported by PDF documents and PDF readers when reading these documents by screen reader.

There also exists a variety of automatic and manual tools that assess and fix accessibility compliance issues in PDFs, including the Adobe Acrobat Pro Accessibility Checker⁵, Common Look⁶, ABBYY FineReader⁷, PAVE⁸, and PDFa Inspector⁹. To our knowledge, PAVE and PDFa Inspector are the only non-proprietary, open-source tools for this purpose. Based on our experiences, however, all of these tools require some degree of human intervention to properly tag a scientific document, and tagging and fixing must be performed for each new version of a PDF, regardless of how minor the change may be.

Guidelines and policy changes have been introduced in the past decade to ameliorate some of the issues around scientific PDF accessibility. Some conferences, such as The ACM CHI Virtual Conference on Human Factors in Computing Systems (CHI) and The ACM SIGACCESS Conference on Computers and Accessibility (ASSETS), have released guidelines for creating accessible submissions.¹⁰ The ACM Digital Library¹¹ provides some publications in HTML format, which is easier to make accessible than PDF [20]. Ribera et al. [40] conducted a case study on DSAI 2016 (Software Development and Technologies for Enhancing Accessibility and Fighting Infoexclusion). The authors of DSAI were responsible for creating accessible proceedings and identified barriers to creating accessible proceedings, including lack of sufficient tooling and lack of awareness of accessibility. The authors recommended creating a new role in the organizing committee dedicated to accessible publishing. These policy changes have led to improvements in localized communities, but have not been widely adopted by all academic publishers and conference organizers.

Table 1 lists prior studies that have analyzed PDF accessibility of academic papers, and shows how our study compares. Prior work has primarily focused on papers published in Human-Computer Interaction and related fields, specific to certain publication venues, while our analysis tries to quantify paper accessibility more broadly. Brady et al. [7] quantified the accessibility of 1,811 papers from CHI 2010-2016, ASSETS 2014, and W4A, assessing the presence of

⁴<https://webaim.org/projects/screenreadersurvey8/>

⁵<https://www.adobe.com/accessibility/products/acrobat/using-acrobat-pro-accessibility-checker.html>

⁶<https://monsido.com/monsido-commonlook-partnership>

⁷<https://pdf.abbyy.com/>

⁸<https://pave-pdf.org/faq.html>

⁹<https://github.com/pdfae/PDFaInspector>

¹⁰See <http://chi2019.acm.org/authors/papers/guide-to-an-accessible-submission/> and https://assets19.sigaccess.org/creating_accessible_pdfs.html

¹¹<https://dl.acm.org/>

Prior work	PDFs analyzed	Venues	Year	Accessibility checker
Brady et al. [7]	1811	CHI, ASSETS and W4A	2011–2014	PDFa Inspector
Lazar et al. [23]	465 + 32	CHI and ASSETS	2014–2015	Adobe Acrobat Action Wizard
Ribera et al. [40]	59	DSAI	2016	Adobe PDF Accessibility Checker 2.0
Nganji [33]	200	<i>Disability & Society, Journal of Developmental and Physical Disabilities, Journal of Learning Disabilities, and Research in Developmental Disabilities</i>	2009–2013	Adobe PDF Accessibility Checker 1.3
Our analysis	11,397	Venues across various fields of study	2010–2019	Adobe Acrobat Accessibility Plug-in Version 21.001.20145

Table 1. Prior work has investigated PDF accessibility for papers published in specific venues such as CHI, ASSETS, W4A, DSAI, or various disability journals. Several of these works were conducted manually, and were limited to a small number of papers, while the more thorough analysis was conducted for CHI and ASSETS, two conference venues focused on accessibility and HCI. Our study expands on this prior work to investigate accessibility over 11,397 PDFs sampled from across different fields of study.

document tags, headers, and language. They found that compliance improved over time as a response to conference organizers offering to make papers accessible as a service to any author upon request. Lazar et al. [23] conducted a study quantifying accessibility compliance at CHI from 2010 to 2016 as well as ASSETS 2015, confirming the results of Brady et al. [7]. They found that across 5 accessibility criteria, the rate of compliance was less than 30% for CHI papers in each of the 7 years that were studied. The study also analyzed papers from ASSETS 2015, an ACM conference explicitly focused on accessibility, and found that those papers had significantly higher rates of compliance, with over 90% of the papers being tagged for correct reading order and no criteria having less than 50% compliance. This finding indicates that community buy-in is an important contributor to paper accessibility. Nganji [33] conducted a study of 200 PDFs of papers published in four disability studies journals, finding that accessibility compliance was between 15-30% for the four journals analyzed, with some publishers having higher adherence than others. To date, no large scale analysis of scientific PDF accessibility has been conducted outside of disability studies and HCI, due in part to the challenge of scaling such an analysis. We believe such an analysis is useful for establishing a baseline and characterizing routes for future improvement. Consequently, as part of this work, we conduct an analysis of scientific PDF accessibility across various fields of study, and report our findings relative to prior work.

3 ANALYSIS OF ACADEMIC PDF ACCESSIBILITY

To capture and better characterize the scope and depth of the problems around academic PDF accessibility, we perform a broad meta-scientific analysis. We aim to measure the extent of the problem (e.g., what proportion of papers have accessible PDFs?), whether the state of PDF accessibility is improving over time (e.g., are papers published in 2019 more likely to be accessible than those published in 2010?), and whether the typesetting software used to create a paper is associated with the accessibility of its PDF (e.g., are papers created using Microsoft Word more or less accessible than papers created with other software?).

Prior studies on PDF accessibility have been limited to papers from specific publication venues such as CHI, ASSETS, W4A, DSAI, and journals in disability research. Notably, these venues are closer to the field of accessible computing, and

are consequently more invested in accessibility.¹² We expand upon this work by investigating accessibility trends across various fields of study and publication venues. Our goal is to characterize the overall state of paper PDF accessibility and identify ongoing challenges to accessibility going forward.

3.1 Data & methods

We sample PDFs from the Semantic Scholar literature corpus [3] for analysis. We construct a dataset of papers by sampling PDFs published in the years of 2010–2019 stratified across the 19 top level fields of study defined by Microsoft Academic Graph [42, 47]. Examples of fields include Biology, Computer Science, Physics, Sociology, and others. This dataset allows us to investigate the overall state of PDF accessibility for academic papers, and to study the relationship between field of study and PDF accessibility.

For each field of study, we sample papers from the top venues by total citation count, along with some documents without venue information, which include things like books and book chapters. The resulting papers come from 1058 unique publication venues; for each field of study, between 29 and 110 publication venues are represented, with Art on the minimum end, and Economics and Computer Science on the maximum end. Each field is represented by an average of 65 different publication venues. The vast majority of documents sampled into our dataset are published papers, rather than preprints or other non-peer-reviewed manuscripts. Publication venues represented in our sample are generally highly reputable journals, for example, *The Lancet* or *Neurology* for Medicine, *The Astrophysical Journal* and *Physical Review Letters* for Physics, or various IEEE publications for Computer Science and Engineering. In some cases, the mapping between publication venue and field of study can be unclear; for example, the publication venue *Mathematical Problems in Engineering* is associated with Mathematics in our sample rather than Engineering. From an examination of the data, classifications seem reasonable and could be justified. We estimate that around 2.2% of the sample are conference papers, 6.1% are book chapters, reports, or lecture notes, less than 0.5% are preprints, and the remaining majority are journal publications. We believe this is a reasonably representative sample of paper-like documents available to scholars and researchers.

We analyze the PDFs in our dataset using the Adobe Acrobat Pro DC PDF accessibility checker.¹³ Though this checker is proprietary and requires a paid license, it is the most comprehensive accessibility checker available and has been used in prior work on accessibility [23, 33, 40]. Alternatively, non-proprietary PDF parsers such as PDFBox¹⁴ do not consistently extract accessibility criteria from sample PDFs, even when the criteria are met. We also prefer Adobe’s checker to PDFa Inspector, used by Brady et al. [7], because PDFa Inspector only analyzes three criteria, whereas we are interested in other accessibility attributes as well, like the presence of alt-text.

For each PDF, the Adobe accessibility checker generates a report that includes whether or not the PDF passes or fails tests for certain accessibility features, such as the inclusion of figure alt-text or properly tagged headings for navigation. Because there is no API or standalone application for the Adobe accessibility checker, it can only be accessed through the user interface of a licensed version of Adobe Acrobat Pro. We develop an AppleScript program that enables us to automatically process papers through the Adobe checker. Our program requires a dedicated computer running MacOS and a licensed version of Adobe Acrobat Pro. It takes 10 seconds on average to download and process each PDF, which

¹²See submission and accessibility guidelines for ASSETS (https://assets19.sigaccess.org/creating_accessible_pdfs.html), CHI (<https://chi2021.acm.org/for-authors/presenting/papers/guide-to-an-accessible-submission>), W4A (<http://www.w4a.info/2021/submissions/technical-papers/>) and DSAI (<http://dsai.ws/2020/submissions/>).

¹³<https://www.adobe.com/accessibility/products/acrobat/using-acrobat-pro-accessibility-checker.html>

¹⁴<https://github.com/apache/pdfbox>

Criterion	CHI 2010[23]	Ours-CHI 2010	Ours-All (11,397)
Alt-text	3.6%	4.0%	7.5%
Table headers	0.7%	1.0%	13.3%
Tagged PDF	6.3%	7.4%	13.4%
Default language	2.3%	3.0%	17.2%
Tab order	0.3%	1.0%	9.3%
Adobe-5 Compliance	-	-	2.4%

Table 2. We reproduce the analysis conducted by Lazar et al. [23] on PDFs of papers published in CHI, showing the percentage of papers that satisfy each of the five accessibility criteria. We find similar compliance rates, indicating that our automated accessibility checker pipeline is comparable to previous analysis methods. We also show the percentage of papers in our full dataset of 11,397 PDFs that satisfy each criterion, along with the percent that satisfy Adobe-5 Compliance.

enables us to scale up our analysis to tens of thousands of papers. Accessibility reports from the checker are saved in HTML format for subsequent analysis.

Each report contains a total of 32 accessibility criteria, marked as “Passed,” “Failed,” or “Needs manual check.”¹⁵ Following Lazar et al. [23], we analyze the following five criteria¹⁶:

- Alt-text: Figures have alternate text.
- Table headers: Tables have headers.
- Tagged PDF: The document is tagged to specify the correct reading order.
- Default language: The document has a specified reading language.
- Tab order: The document is tagged with correct reading order, used for navigation with the tab key.

For our analysis, we also report *Total Compliance*, which refers to the sum number of accessibility criteria met (e.g. if a paper has met 3 out of the 5 criteria we specify, then Total Compliance is 3). In some cases, we report the *Normalized Total Compliance*, which is computed as the Total Compliance divided by 5, and can be interpreted as the proportion of the 5 criteria which are satisfied. We also report *Adobe-5 Compliance*, a binary value of whether a paper has met all 5 criteria we specify (1 if all 5 criteria are met, 0 if any are not met), and the rate of Adobe-5 Compliance for papers in our dataset.

In addition to running the accessibility checker, we also extract metadata for each PDF, focusing on metadata related to the PDF creation process. PDF metadata are generated by the software used to create each file, and we analyze the associations between different PDF creation software and the accessibility of the resulting PDF document. Our hypothesis is that some classes of software (such as Microsoft Word) produce more accessible PDFs.

3.2 Accuracy of our automated accessibility checker

Previous work employed different versions of the Adobe Accessibility Checker to generate paper accessibility reports. To confirm the accuracy of our checker, as well as the automated script we create to perform the analysis, we run our checker on CHI 2010 papers to reproduce the results of Lazar et al. [23]. We identify CHI papers using DOIs reported by the ACM, and resolve these to PDFs in the Semantic Scholar corpus [3]. We identify 3,248 CHI papers in the corpus, and generate accessibility reports for these using our automated checker.

¹⁵Please see <https://helpx.adobe.com/acrobat/using/create-verify-pdf-accessibility.html> for a description of the accessibility report.

¹⁶For papers containing no tables and/or no figures, the Adobe checker can still return both pass or fail for the Table header and Alt-text criteria respectively. When objects in the PDF are *not* tagged, the checker will fail these criteria even when the paper has no tables and/or no figures. When objects in the PDF *are* tagged and the PDF is accessible, the checker will pass these criteria even when the paper has no tables or no figures.

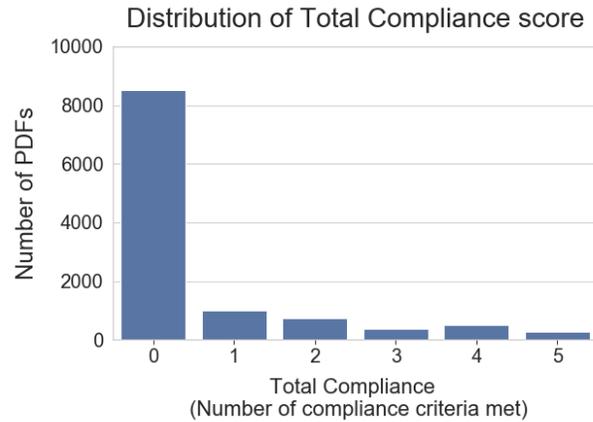


Fig. 2. The distribution of numbers of PDFs in our dataset that meet our defined accessibility compliance criteria. A large majority (8519) of PDFs in our sample meet 0 out of 5 accessibility criteria. Of those meeting 1 criterion (Total Compliance = 1), the most commonly met criterion is Default Language (793 of 1010, 78.5%). Of those meeting 4 criteria (Total Compliance = 4), the most common missing criterion is Alt-text (396 of 494, 80.2%).

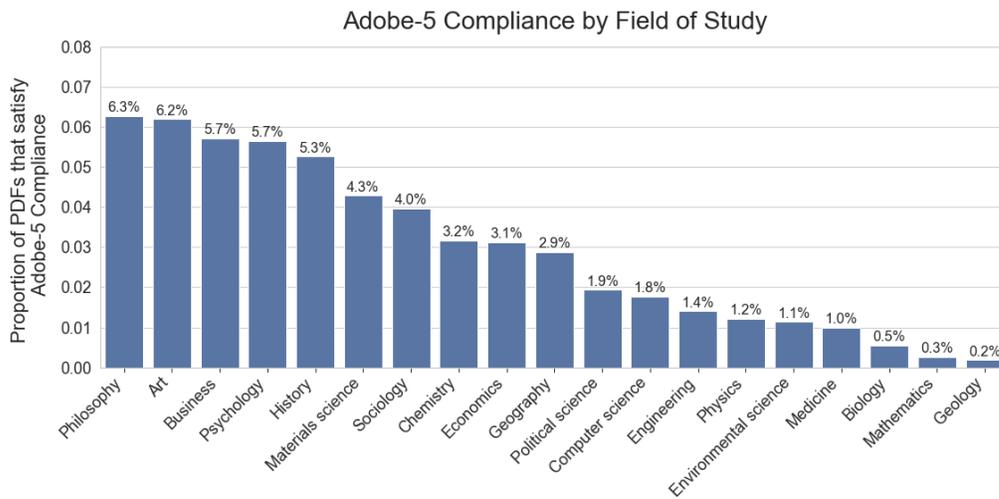


Fig. 3. Percent of papers per field of study that meet all 5 accessibility criteria defined in Adobe-5 Compliance. Philosophy, Art, and Psychology have the highest rates of Adobe-5 Compliance satisfaction while Biology, Mathematics, and Geology have the lowest rates. None of the fields had more than 6.5% of PDFs satisfying Adobe-5 Compliance.

Our results shows similar rates of compliance compared to what was measured by Lazar et al. [23] (see Table 2 for results). This indicates that our automated accessibility checker produces comparable results to previous studies.

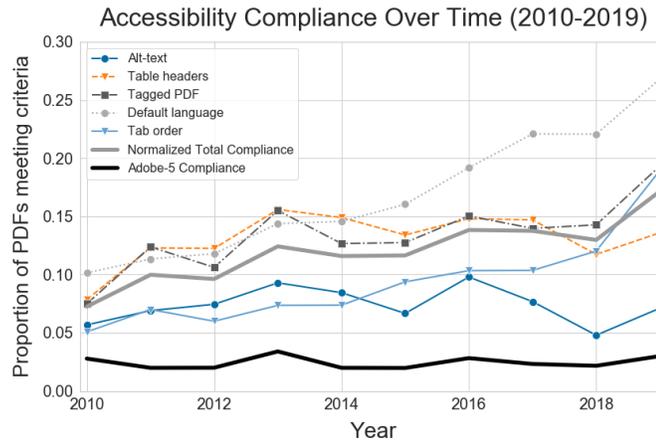


Fig. 4. Accessibility compliance over time (2010-2019). The rate of Adobe-5 Compliance has remained relatively stable over the last decade, at around 2–3%. Compliance along several criteria have improved over time, though the largest improvements have been in Default Language, the simplistic criteria to meet. Modest improvements are seen for Table headers, Tagged PDFs, and Tab order. The presence of alt-text has remained stable and lower, around 5–10%.

3.3 Proportion of papers with accessible PDFs

Around 1.6% of PDFs we attempted to process failed in the Adobe checker (i.e., we could not generate an accessibility report). The accessibility checker most commonly fails because the PDF file is password protected, or the PDF file is corrupt. In both of these cases, the PDF is inaccessible to the user. We exclude these PDFs from subsequent analysis.

Accessibility compliance over all papers is low. Table 2 shows the percent of papers meeting each of the five criteria, as well as the Adobe-5 Compliance rate associated with this sample of papers. Figure 2 shows that the vast majority of papers do not meet any of the five accessibility criteria (8519 papers, 74.7% do not meet any criteria) and very few (275 papers, 2.4%) meet all five. Of those PDFs meeting 1 criterion, the most commonly met criterion is Default Language (793 of 1010, 78.5%). Of those PDFs meeting 4 criteria, the most common *missing* criterion is Alt-text (396 of 494, 80.2%). In fact, only 854 PDFs (7.5%) in the whole dataset have alt-text for figures. This is intuitive as Alt-text is the only criterion that *always* requires author input to achieve, while the other four criteria can be derived from the document or automatically inferred, depending on the software used to generate the PDF.

As shown in Figure 3, all fields have an Adobe-5 Compliance of less than 7%. The fields with the highest rates of compliance are Philosophy (6.3%), Art (6.2%), Business (5.7%), Psychology (5.7%), and History (5.3%) while the fields with the lowest rates of compliance are Geology (0.2%), Mathematics (0.3%), and Biology (0.6%). Fields associated with higher compliance tend to be closer to the humanities, and those with lower levels of compliance tend to be science and engineering fields. The prevalence of different document editing and typesetting software by field of study may explain some of these differences, and we explore these associations in Section 3.5.

3.4 Trends in paper accessibility over time

We show changes in compliance for all fields of study over time in Figure 4. With the exception of Default Language, all accessibility criteria demonstrate slowly increasing or stable compliance rates over the past decade, with increases seen in Tagged PDFs and Tab order over time. Default language compliance is increasing most rapidly, from around

Typesetting Software	Count (%)
Adobe InDesign	1591 (14.0%)
LaTeX	1431 (12.6%)
Arbortext APP	1374 (12.1%)
Microsoft Word	1318 (11.6%)
Printer	1021 (9.0%)
Other	4662 (40.9%)

Table 3. Count of papers per Typesetting Software. “Other” includes PDFs created with an additional 24 unique software programs, each with counts of less than 350, as well as those created with an unknown typesetting software.

10% compliance in 2010 to more than 25% in 2019. This may be due to changes in PDF generation defaults in various typesetting software. Though this improvement is good, Default Language is the easiest of the five criteria to bring into compliance, and arguably the least valuable in terms of improving the accessible reading experience. The criterion with the lowest rate of compliance is Alt-text, which has remained stable between 5–10% and has been lower in recent years. Since Alt-text is the only criterion of the five which always necessitates author intervention, we believe this is a sign that authors have not become more attuned to accessibility needs, and that at least some of the improvements we see over time can be attributed to typesetting software or publisher-level changes.

3.5 Association between typesetting software and paper accessibility

Typesetting software is extracted from PDF metadata and manually canonicalized. We extract values for three metadata fields: `xmp:CreatorTool`, `pdf:docinfo:creator_tool`, and `producer`. All unique PDF creation tools associated with more than 20 PDFs in our dataset are reviewed and mapped to a canonical typesetting software. For example, the values (`latex`, `pdftex`, `tex live`, `tex`, `vtex pdf`, `xetex`) are mapped to the LaTeX cluster, while the values (`microsoft`, `for word`, `word`) and other variants are mapped to the Microsoft Word cluster. We realize that not all Microsoft Word versions, LaTeX distributions, or other versions of typesetting software within a cluster are equal, but this normalization allows us to generalize over these software clusters. For analysis, we compare the five most commonly observed typesetting software clusters in our dataset, grouping all others into a cluster called Other.

We report the distribution of typesetting software in Table 3. The most popular PDF creators are Adobe InDesign, LaTeX, Arbortext APP, Microsoft Word, and Printer. “Printer” refers to PDFs generated by a printer driver (by selecting “Print” → “Save as PDF” in most operating systems); unfortunately, creating a PDF through printing provides no indicator of what software was used to typeset the document, and is generally associated with very low accessibility compliance. The “Other” category aggregates papers created by all other clusters of typesetting software; each of these clusters is associated with less than 350 PDFs, i.e., the falloff is steep after the Printer cluster. For the following analysis, we present a comparison between the five most common PDF creator clusters.

Figure 5 shows histograms of the Total Compliance score for PDFs in the five most common typesetting software clusters. While the vast majority of papers do not meet any accessibility criteria, it is clear that Microsoft Word produces the most accessible PDFs, followed by Adobe InDesign. To determine the significance of this difference, we compute the ANOVA and Kruskal-Wallis [22] statistics with the PDF typesetting software clusters as the sample groups and the Total Compliance as the measurements for the groups. We compute an ANOVA statistic of 2587.1 ($p < 0.001$) and

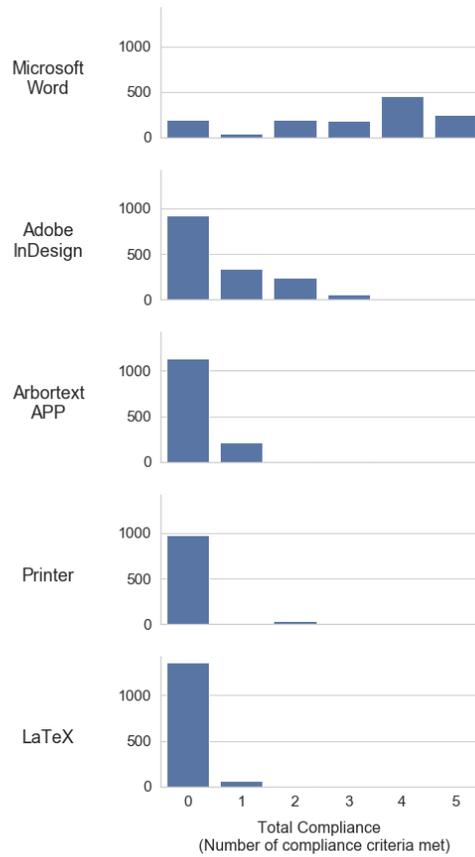


Fig. 5. Histograms showing the distribution of Total Compliance scores for each of the top 5 typesetting software, ordered by decreasing mean Total Compliance. Microsoft Word stands out as producing PDFs with significantly higher Total Compliance than other typesetting software. Three of the top five PDF typesetting software clusters, Arbortext APP, Printer, and LaTeX, produce PDFs with low Total Compliance, with the majority of PDFs at 0 compliance.

a Kruskal-Wallis H statistic of 4422.0 ($p < 0.001$). This indicates a significant difference in the distribution of Total Compliance scores between the five most common PDF typesetting software.

In Figure 6, we observe again that usage of Microsoft Word is highly correlated with accessibility compliance. Here, we plot the proportion of Microsoft Word usage per field of study and the corresponding mean normalized Total Compliance rates for those fields. Higher rates of Microsoft Word usage are statistically correlated with higher mean normalized Total Compliance ($r = 0.89$, $p < 0.001$).

In Figure 7, we show the proportion of usage of each of the five typesetting software over time. In recent years, Adobe InDesign, LaTeX, and Microsoft Word usage are proportionally increasing, while the proportion of Printer-created PDFs is declining. The increase in Adobe InDesign and Microsoft Word have likely driven the increase in rates of Total Compliance over time, since these typesetting software are the most associated with higher accessibility compliance.

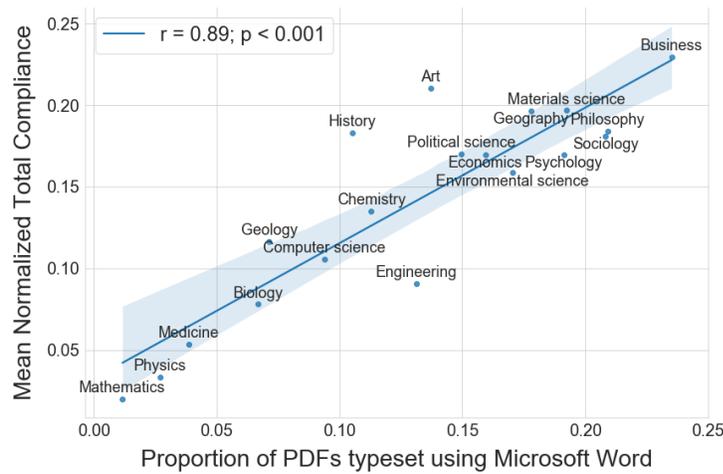


Fig. 6. There is a strong correlation ($r = 0.89$, $p < 0.001$, 95% CI shown) between the proportion of PDFs typeset using Microsoft Word and the mean normalized Total Compliance of papers by field of study. Fields such as Business, Philosophy, Sociology, Materials science, and Psychology use Microsoft Word around or over 20% of the time, and have correspondingly higher mean accessibility compliance. On the other end of the spectrum are fields like Mathematics, Physics, and Medicine, where Microsoft Word is rarely used, and which have very low levels of mean compliance.

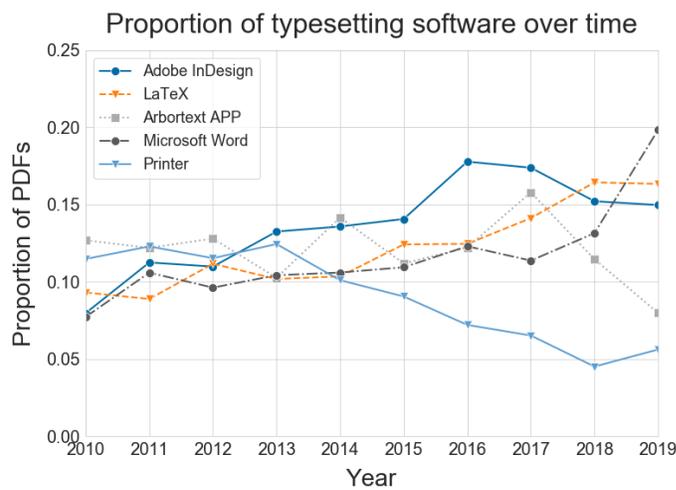


Fig. 7. The proportion of PDFs typeset by the five most common typesetting software over time. Software such as Adobe InDesign, LaTeX, and Microsoft Word are increasing in popularity over time.

3.6 Summary of analyses

Overall, accessibility compliance over the past decade and across all fields of study have slowly improved. Full compliance based on Adobe-5 Compliance, however, has remained around 2.4% on average and does not show trends towards improving. Improvements in several compliance criteria are observed, with Default Language being the most improved,

neering 30% coverage in 2019. However, Default Language is the easiest criteria to meet, and arguably produces the least amount of accessibility improvement in user experience. Criteria such as Tagged PDFs, Tab order, and Headers show modest improvements over time, though only between 10–15% of papers in our sample meet any one of these individual criteria. Alt-text compliance is the lowest of our measured criteria, and as the only criterion of the five requiring author intervention in all cases, the lack of alt-text may be indicative of the general lack of author awareness and contribution to accessibility efforts for scientific papers.

Based on our analysis, typesetting software plays a large role in document accessibility. Of the most common PDF creator software, Microsoft Word appears to produce the most accessibility-compliant PDFs, while LaTeX produces PDFs with the lowest compliance. Microsoft has recently made investments in the accessibility of their Office 365 Suite.¹⁷ It is clear that software can help increase accessibility compliance by prioritizing accessibility concerns during document creation, and we encourage other developers of typesetting and publishing software to prioritize accessibility concerns in their development process.

Improvements in accessibility compliance have stalled over the past decade, likely because accessibility concerns are considered marginal, and are outside of the awareness of most publishing authors and researchers. Significant changes in the authorial and publication processes are needed to change this status quo, and to increase the accessibility of scientific papers for BLV users going forward. Though we believe and encourage change in the academic paper authorial and publication process in relation to accessibility, the likelihood of rapid improvement is low and these changes will not impact the many millions of academic PDFs that have already been published. Therefore, we introduce a technological solution that may mitigate some of the accessibility challenges of existing paper PDFs, and aim to understand how this solution and others like it could serve the immediate needs of the BLV research community.

4 CONVERTING PDF TO HTML: THE SCIA11Y PIPELINE

To address the broad accessibility challenges described in Section 3, we propose and prototype a system for extracting semantic content from paper PDFs and re-rendering this content as accessible HTML. HTML is widely accepted as a more accessible document format than PDFs. In the 2019 Access SIGCHI Report, the authors discuss the reasoning behind switching CHI publications to a new HTML5 proceedings format to improve accessibility [27]. By rendering the content of paper PDFs as HTML, and introducing proper reading order and accessibility features such as section headings, links, and figure tags, we can offset many of the issues of reading from an inaccessible PDF. Our PDF to HTML rendering system is named SciA11y after the community-adopted numeronym for digital accessibility.¹⁸

Figure 1 provides a schematic for the approach. SciA11y leverages the two open source PDF processing projects S2ORC [24] and DeepFigures [43], the Semantic Scholar API,¹⁹ and a custom Flask application for rendering the extracted content of the PDF as HTML. The S2ORC project [24] integrates the Grobid machine learning library [25] and a custom XML to JSON parser²⁰ to produce a structured representation of paper text. We use a version of the S2ORC pipeline that is based on Grobid v0.6.0. The resulting JSON representation includes metadata fields like title, authors, and affiliations, and paper content fields such as abstract, section headers, body text organized into paragraphs, bibliography entries, and figure and table objects (though not the figure images themselves). The output also contains links between inline citations and figure/table references respectively to bibliography entries and figure/table objects.

¹⁷<https://www.microsoft.com/en-us/accessibility/microsoft-365>

¹⁸<https://www.a11yproject.com/>

¹⁹<https://api.semanticscholar.org/>

²⁰Available at <https://github.com/allenai/s2orc-doc2json>

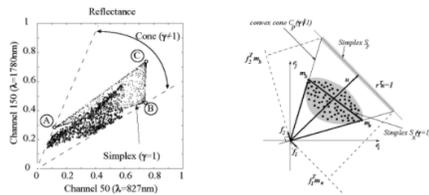


Figure 1. 2-D scatter-plot of mixtures of the three endmembers shown in Fig. 2; Circles denote pure materials. Right: illustration of the VCA algorithm.



Figure 2. Not extracted; please refer to original document.

EQUATION (1): Not extracted; please refer to original document.

Fig. 8. A successfully extracted figure from Nascimento and Bioucas-Dias [30] is shown with its corresponding figure caption (*top left*). When figures are not extracted and inferred to exist (handle mentioned in text or number between two extracted figures), a placeholder image is shown along with a message referencing the failed extraction (*top right*). Similarly, when an equation is detected to be present in the PDF and not extracted, we insert text signaling the failed extraction and refer the user to the source document (*bottom*).

DeepFigures [43], on the other hand, leverages a computer vision model to extract images of figures and tables as well as their corresponding captions from the source PDF.

The outputs of S2ORC and DeepFigures are stitched together to form the HTML render as in Figure 1. We place header tags (`<h1> . . . </h1>`, `<h2> . . . </h2>`) around the title, authors, abstract, section headings, and reference heading. Paragraphs of body text are enclosed in `<p> . . . </p>` tags in order within their appropriate sections. Bibliography entries are provided in an unordered list under the reference heading. Figures and tables are enclosed in `<figure> . . . </figure>` tags and placement is inferred based on mentions in the text. A figure or table is placed immediately after the paragraph in which its handle is first mentioned (e.g. “In Fig. 1, we show...” is the first mention of Figure 1 and the figure is placed directly after the paragraph with this mention). Figure and table captions are attached to their corresponding image objects, so that correspondences between the caption text and image are made explicit (in PDFs, this is usually not the case). Any figures or tables which are not mentioned in order in the text are placed in order nonetheless; in other words, if paragraph 1 mentions Figure 1 and paragraph 2 mentions Figure 3, both Figure 1 and 2 will be placed directly following paragraph 1 and Figure 3 following paragraph 2. This ensures that the layout for the HTML render closely approximates the intended reading order. We justify this decision based on user feedback from our pilot study, which is discussed in Section 6.

In some cases, we are able to successfully process a PDF through S2ORC to extract textual content but DeepFigures either fails to process the PDF or fails to extract some or all figures from the PDF. To mitigate the cognitive dissonance around figure or table mentions without corresponding figure or table objects, we insert placeholder objects into the HTML render as in Figure 8. For example, if “Figure 2” is mentioned in the text but is not successfully extracted by DeepFigures, we would insert a placeholder image for the figure based on the logic described in the previous paragraph along with the text “Figure 2. Not extracted; please refer to original document.” Similarly, mathematical equations that we cannot currently extract are acknowledged with the same placeholder text.

We add links between inline citations and the corresponding reference entry where possible. We insert links at each inline citation in the body text that link to the corresponding bibliography entry. Following each bibliography entry, we provide links back to the first mention of that entry in each section of the paper in which it was mentioned. For example, if bibliography entry [1] is cited in the “II. Related Works” section and the “III. Methods” section, we provide two links following the entry in the bibliography to the corresponding citation locations in sections II and III, as in:

[1] Last name et al. Paper title. Venue. DOI.
[Link to return to Section II](#), [Link to return to Section III](#)

This allows users to navigate back to their reading location in the document after clicking through to a bibliography entry. A user may otherwise hesitate to resolve a link, because it may result in losing their place and train of thought. Finally, we introduce a table of contents near the beginning of the HTML render to facilitate better understanding of overall document structure. The table of contents includes all section titles, linked to the corresponding sections, as well as figures and tables nested under their respective section headers. The table of contents provides a rapid overview of the structure of the document, and facilitates rapid navigation to the reader’s desired sections.

In the current iteration of the HTML render, we do not display author affiliations, footnotes, or mathematical equations due to the difficulty of extracting these pieces of information from the PDF. Though some of the elements are extracted in S2ORC, the overall quality of the extractions for these elements is lower, and is currently insufficient for surfacing in the prototype (see Section 5 for details). Future work includes investigating the possibility of extracting and exposing these elements, either by improving current models or training new models targeted towards the extraction of specific paper elements.

We leverage the feedback we received during our pilot studies (see Section 6) to make improvements prior to the main user study. We denote the versions of the prototype as v0.1 (initial version; version seen by P1), v0.2 (version seen by P2), and v0.3 (version seen by all other participants in the main user study). Features implemented in v0.1 include the primary components of the HTML render such as title, authors, abstract, body text with section headers, figures and tables, references, and links between inline citations and references. In v0.1, figures and tables were placed in a separate section following the main body of the paper. Following P1, for version v0.2, we implemented the table of contents, inserted placeholders for objects that we could not extract, and began inserting figures and tables into the body text adjacent to their first mentions. This last change was made in response to P1’s feedback that navigating away to figures caused him to lose his reading location. Following P2, for version v0.3, we implemented only minor changes. P2 signaled during his session that URLs in the bibliography were not being correctly extracted, so we patched the data to correctly extract and display URLs in bibliography entries.

Based on our evaluation of the quality of these HTML renders (Section 5) and user feedback and response (Section 6), we believe our approach can dramatically increase the screen reader navigability and accessibility of scientific papers across all disciplines by providing an alternate and more accessible HTML version of these papers. Properly tagged section headings allow for quick navigation and skimming of a paper, links between inline citations and bibliography entries allow users to browse to cited papers without losing their place, and figure tags for figure and table objects allow for direct navigation to these in-paper objects. We now discuss the quality of our PDF extractions (Section 5) and user response to the prototype (Section 6) in detail.

5 HTML RENDER QUALITY EVALUATION

Extracting semantic content from PDF is an imperfect process. Though re-rendering a PDF as HTML can increase a document’s accessibility, the process relies on machine learning models that can make mistakes when extracting information. As we glean from user studies, BLV users may have some tolerance for error, but there is an inherent trade-off between errors and perceived trust in the system. We conduct a study to estimate the (1) faithfulness of the HTML renders to the source PDFs, and (2) overall readability of the resulting HTML renders. We define *faithfulness* as how accurately the HTML render represents different facets of the PDF document, such as displaying the correct title, section headers, and figure captions. These facets are measured as the number of errors that are made in rendering, e.g., mistakenly parsing one figure caption into the body text is counted as one error towards that facet. *Readability*, on the other hand, is an ordinal variable meant to capture the overall usability of the parse. Documents are given one of three grades, those with no major problems, some problems, and many problems impacting readability.

To evaluate readability and faithfulness, we first perform open coding on a small sample of document PDFs and corresponding SciA11y HTML renders. The purpose of this exercise is to identify facets of extraction that impact the ability to read a paper. A rubric is then designed based on these identified facets. The process taken to design the evaluation rubric, the rubric’s content, and annotation instructions are detailed in Section 5.2. We then annotate a sample of 385 papers across different fields of study using this rubric. For each category of errors identified during open coding, we compute the overall error rates seen in our sample. We also present the overall assessed readability, reported in aggregate over our sample and by fields of study. Results of this evaluation are presented in Section 5.3.

5.1 Open coding of document facets

One author performed open coding on a sample of papers, comparing the PDF and SciA11y HTML renders to identify inconsistencies and facets that impact the faithfulness of document representation. Papers are sampled from the Semantic Scholar API²¹ using various search terms, and selecting the top 3 results for each search term for which a PDF and S2ORC parse are available. Search terms were selected to achieve coverage over different domains, and the top papers are sampled to select for relevant publications. The author stopped sampling papers upon reaching saturation, resulting in 8 search terms and 24 papers. The search terms used were: human computer interaction, epilepsy, quasars, language model, influenza epidemiology, anabolic steroids, social networks, and arctic snow cover.

For each paper, the author evaluated the PDF and HTML render side-by-side, scanning through the document to identify points of difference between the two document representations. Specifically, the author looked for any text in the PDF that is not shown in the HTML, any text from the PDF that is mixed into the main text of the HTML (e.g. figure captions, headers, or footnotes that should be separate from the main text but are mixed in, interrupting the reading flow), and other parsing mistakes (e.g. errors with math, missing lists and tables etc). These observations are detailed qualitatively, and each facet is assessed for its faithfulness to the original PDF document as well as its overall impact on readability.

5.2 Evaluation rubric

Observations from open coding are coalesced into an evaluation rubric and form for grading the quality and faithfulness of the HTML render. The evaluation form attempts to capture errors in PDF extraction that affect each of the primary

²¹<https://api.semanticscholar.org/>

Category	Description	Common errors
TITLE	The title and subtitle of the paper	Missing words Extra words
AUTHORS	A list of authors who wrote the paper; this includes affiliation, though we do not explicitly evaluate affiliation in this study	Missing authors Extra authors Misspellings
ABSTRACT	The abstract of the paper	Some text not extracted Other text incorrectly extracted as abstract
SECTION HEADINGS	The text of section headings	Some headings not extracted (part of body text) Other text incorrectly extracted as headings
BODY TEXT	The main text of the paper, organized by paragraph under each section heading	Some paragraphs not extracted (missing) Some text not extracted Other text incorrectly extracted as body text
FIGURES	Images, captions, and alt-text of each figure	Figure not extracted Caption text not extracted (part of body text) Other text incorrectly extracted as caption text
TABLES	Caption/title and content of each table	Table not extracted (not part of body text) Table not extracted (part of body text) Caption text not extracted (part of body text) Other text incorrectly extracted as caption text
EQUATIONS	Mathematical formulas, represented in TeX or Math ML; note: our current pipeline does not extract math	Some equations not extracted Some equations incorrectly extracted
BIBLIOGRAPHY	Bibliography entries in the reference section	Some bibliography entries not extracted Some bibliography entries incorrectly extracted Other text incorrectly extracted as bibliography
INLINE CITATIONS	Inline citations from the body text to papers in the bibliography section	Some inline citations not detected Some inline citations incorrectly linked
HEADERS, FOOTERS & FOOTNOTES	Page headers and footers, footnotes, endnotes, and other text that is not a part of the main body of the document	Some headers and footers incorrectly extracted into body text

Table 4. Categories of paper objects identified for evaluation along with the common errors seen for each category.

semantic categories identified for proper reading. These semantic categories and common extraction errors are given in Table 4.

Questions in the form are designed to capture each type of faithfulness error, while allowing annotators to qualify their responses. We also include a question to capture the overall readability of the HTML render. Instructions for completing the annotation form are provided in Appendix A.1; the final version of the form is replicated in Appendix A.1; and the rubric for overall readability evaluation is given in Appendix A.3.

Three authors iterated twice on the content of the evaluation form, until they came to a consensus that all evaluation categories were adequately addressed using a minimum set of questions. Two authors then participated in pilot annotations, where each person independently annotated the same set of five papers sampled from the set labeled by the third author during open coding. Answers to all numeric questions were within ± 1 for these five papers when comparing the two authors' annotations. All three authors discussed discrepancies in overall readability score, iterating on the rubric defined in Appendix A.3 and coming to a consensus. The finalized form and rubric are used for evaluation.

Of the categories and errors described in Table 4, our current pipeline does not extract table content and equations. Tables are extracted as images by DeepFigures [43], which do not contain table semantic information. Regarding equations, we distinguish between inline equations (math written in the body text) and display equations (independent line items that can usually be referenced by number); for this work, we evaluated a small sample of papers for successful extraction of display equations. Though some display equations are recognized, the quality of equation extraction is low, usually resulting in missing tokens or improper math formatting. Therefore, we decided to replace display equations in the prototype with the equation placeholder shown in Figure 8. Since problems with mathematical formulae are among those most mentioned by users in our study, equation extraction is among our most urgent future goals, and we discuss some options in Section 7.1.

5.3 Evaluation results

We start with the dataset of 11,397 papers we analyze in Section 3, and subsample 535 documents stratified by field of study. Two expert annotators with undergraduate science training code papers from this sample, with an aim of annotating around 20 papers per field of study. Though we achieve the target number for most fields, we missed this target for some fields closer to the humanities because more of these documents are difficult to manually annotate within our time and resource constraints. For example, documents are deemed unsuitable for annotation if they are not papers (i.e., they are books, posters, abstracts, etc), if they are too long, or if they are not in English. In these cases, the annotators can skip the document. Detailed guidance on suitability is provided in the annotation instructions (see Appendix A.1).

The two annotators annotated 385 unique papers and skipped 137. The resulting annotated sample consists of papers from 195 unique publication venues. Each paper takes 5–10 minutes to grade. Documents are skipped primarily due to language (paper not in English), length, or the document is not a paper. Inter-annotator agreement is computed over a sample of 20 papers over each of the evaluated facets. We report Cohen’s Kappa for categorical questions such as those on the extraction of title, authors, abstract, and bibliography. For numerical questions such as counting the occurrence of extraction errors related to figures, tables, section headings, and body paragraphs etc, we report the intraclass correlation coefficient (ICC) as well as the average difference of values between the two annotators. See Table 5 for these results. Agreement was high for most element-level annotator questions. Annotators had the highest levels of disagreement on the evaluation of header/footer/footnote errors, section heading errors, and body paragraph errors, likely due to these being text-based and the most numerous; though the average differences reported between annotators on these questions are only between 1-2. Likewise, agreement on overall readability score is modest, at 0.55; we note, however, that neither annotator labeled any paper as having no major readability problems when the other annotator labeled it as having lots of readability problems.

All results and statistics are reported on the set of 385 annotated papers. Figure 9 shows the breakdown of each type of error and the frequency at which it occurs. Metadata elements like title, authors, and abstract are successfully extracted the majority of the time. For figure and table elements, approximately 25% of papers in our evaluation sample do not include figures, and around 45% do not have tables. Of those that have figures, the majority (201, 69.1% of 291) do not have extraction or parsing errors; around half of documents with errors have errors that only relate to one figure. Similarly, the majority of tables and table captions are correctly identified as tables and table captions, and are not incorrectly mixed into the body text. We note that the lack of an error here does not indicate that the table is extracted correctly in an accessible manner, just that it is not incorrectly parsed as body text.

Evaluation criteria	Number of classes	Agreement	Cohen's Kappa	ICC	Mean Difference (\pm SD)
Title	3	0.87	0.33	-	-
Authors	3	1.00	1.00	-	-
Abstract	3	0.95	0.64	-	-
Number of figures	-	1.00	-	1.00	0.00 \pm 0.00
Figure extraction errors	-	0.89	-	1.00	0.11 \pm 0.31
Figure caption errors	-	0.89	-	1.00	0.11 \pm 0.31
Number of tables	-	0.92	-	0.98	0.12 \pm 0.43
Table extraction errors	-	0.89	-	0.98	0.17 \pm 0.50
Table caption errors	-	0.78	-	0.94	0.33 \pm 0.67
Header/footer/footnote errors	-	0.40	-	0.60	1.88 \pm 2.12
Section heading errors	-	0.71	-	0.79	0.71 \pm 1.70
Body paragraph errors	-	0.46	-	0.66	1.50 \pm 2.22
Bibliography extraction	4	0.94	0.82	-	-
Inline citation linking	4	0.80	0.11	-	-
Overall score	3	0.55	0.07	-	-

Table 5. Inter-rater agreement for evaluation. For categorical questions, such as title, author, abstract, bibliography, inline citation, and overall score, we report the number of classes available for annotation, along with annotator agreement and Cohen's Kappa. For numerical questions, such as the number of each type of extraction error, we report agreement, the intraclass correlation coefficient (ICC), and the average difference and standard deviation of the values between the two annotators.

Unsurprisingly, errors in text element parsing are the most prevalent, especially for headers/footers/footnotes and section headings. The most common type of header/footer/footnote error observed are when these texts are mixed into the body text around page breaks, interrupting reading flow. These types of errors are also observed frequently during screen reader use when reading directly from an untagged PDF. For section headings, in particular, the majority of papers have errors; around 67% of papers have between 1–5 errors (either missed headings or extraneous headings), and 9% have more than 5 errors. Due to the large number of section headings in papers, parsing errors are more frequent, and unfortunately, these errors impact the ability to properly navigate the HTML parse. Errors in body text extraction also negatively impact readability, in this case, select text in the document is being missed completely in the HTML render. We see that though the majority of parses have no body text errors, around 33% of papers have between 1-5 missing paragraphs.

Figure 9(d) shows grading results for bibliography elements. Our pipeline is quite good at extracting bibliography entries, extracting all or most entries in the vast majority of cases, and successfully linking inline citations to these bibliography entries also in a large majority of cases. When bibliography extraction fails, it tends to fail catastrophically, resulting in no or few extractions.

The overall readability score is provided in Figure 9(e). A majority of papers (54.5%, 210 papers) have no major problems impacting readability. Another 31.7% (122) of papers have some problems impacting readability, and 13.8% (53) of papers have lots of readability problems. We are encouraged that a majority of HTML renders have no major problems, though our results necessitate further understanding of the papers with which our extraction pipeline has difficulty. If papers with lots of problems can be identified *a priori*, we can prevent surfacing these low quality parses to

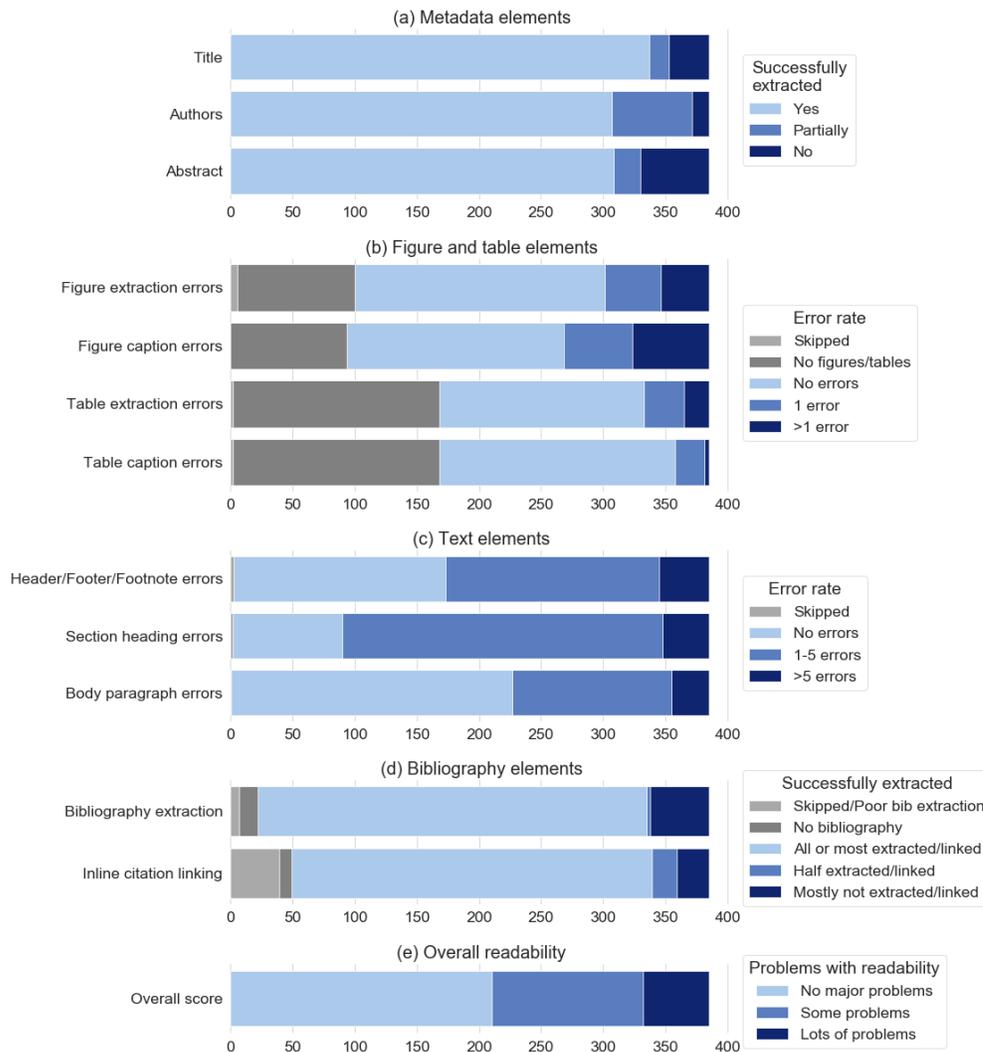


Fig. 9. Evaluation results for various document components. Corresponding numbers are provided in Table 13 in Appendix B.

the user. We perform some preliminary experiments to identify paper features that are more correlated with readability problems, though no features stood out as being predictive; we present those results in Appendix C.

In Figure 10, we show the breakdown of overall readability by field of study, plotting the proportion of papers per field that are classified as having no major problems, some problems, and lots of problems impacting readability. Many fields have similar distributions compared to the overall evaluation set. However, we note that some fields such as Art, Business, Economics, and Environmental science to some degree, have significantly lower quality extraction results. We posit that this may be due to biases in our PDF extraction pipeline. Some of the machine learning modules we use are primarily trained on paper data from the biomedical and Computer Science domains, where large scale labeled

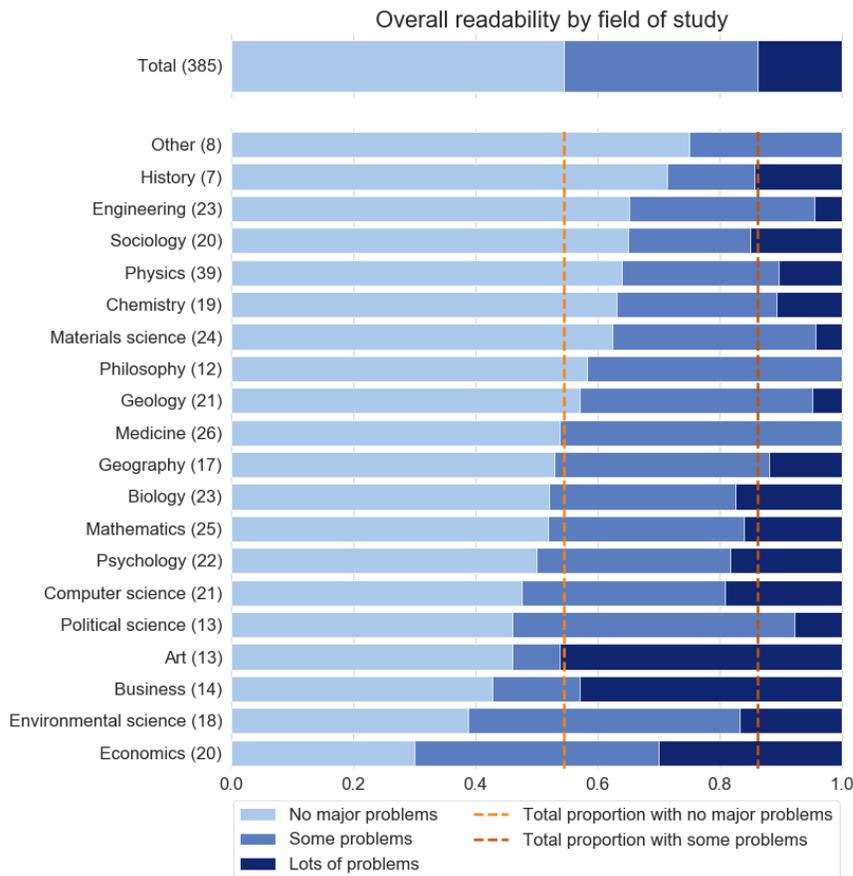


Fig. 10. Overall readability results as proportion of total split by field of study, sorted by the percentage of papers with no major problems. The number of documents analyzed in each field is given, ranging from $N=7$ (History) to $N=39$ (Physics). The fields of study with the worst parse quality (Economics, Environmental science, Business, Art, and Political science) tend to be closer to the humanities, and may be due to the under-representation of papers from these fields in the data used to train the PDF parsers we use in our extraction pipeline. Corresponding numbers are provided in Table 14 in Appendix B.

PDF extraction datasets can be found. Humanities-adjacent fields like Art and Business have very different publication norms, and the different layouts and content of papers and documents in these fields may provide additional challenges to our system, resulting in lower quality extraction and rendering.

6 USER STUDY

We conduct an exploratory user study to better understand the needs of BLV scientists when reading papers, and to assess whether our prototype supports these needs. The study consists of a preliminary questionnaire and semi-structured video interview. Interviews are conducted remotely on Zoom.²² All recruitment materials, questionnaires, and the interview plan are reviewed and approved by the internal review board at anonymized. We recruit and interview

²²<https://zoom.us/>

six users, with a pilot involving two users, and a main study involving four users. Modifications to the prototype between pilots and the main study can be found in Section 4. We report results from all six participants in any analysis that does not involve the prototype, and for analysis that directly involves the prototype, we denote all cases where prototype modifications between the pilot and main study may impact our results.

The inclusion criteria for participants are:

- The participant is over 18 years of age;
- The participant identifies as blind or low vision;
- The participant reads scientific papers regularly (more than 5 per year);
- The participant must have used a screen reader to read a paper in the last year; and
- The participant must complete the pre-interview questionnaire.

Participants were recruited through mailing lists, word-of-mouth, and snowball sampling. Prior to each interview, the participant was asked to provide several keywords corresponding to their subject areas of interest, and between 3–5 papers where they experienced difficulty reading the PDF. Among the 3–5 papers, we selected one paper to use for the study, based on the availability of an HTML render, and maximizing the features that would be seen during the user study (e.g., given a choice between a paper with figures and a paper without figures but where both otherwise demonstrate the same paper components, we would select the paper with figures). Each study session was 75 minutes, consisting of three phases:

Phase I: Capturing challenges with current work flow

The primary research questions we investigate in this phase are:

- What methods and/or tools do BLV researchers use to assist in reading the literature?
- What main accessibility challenges do BLV researchers face?
- How do BLV researchers cope with these challenges?

We first asked the participant to describe their current workflow and the challenges they face when reading papers, clarifying how the user copes with challenges when their workflow does not adequately address the problem. We then asked the participant to demonstrate how they currently read a paper, by opening a paper PDF and walking us through the usage of their tools (PDF viewer, screen reader, magnifier, speech-to-text, etc). Participants kept their computer audio on so we could hear the output of their reader tools. The participant was asked to think aloud and describe their actions when reading the paper. We asked the participant to demonstrate any reading challenges they described in their pre-interview questionnaire. At the end of this phase, we asked the participant to assess how easy or difficult it was to read the paper with their current reading pipeline.

Phase II: Interaction with prototype

The primary research questions we investigate in this phase are:

- What features of the HTML render resonated positively with the participant?
- What problems can be identified in the HTML render?

The goal of this phase was to understand how helpful or not helpful the HTML render is to the participant. The participant was asked to interact with an HTML render of the same paper they read in Phase I in the SciA11y prototype. We first provided an introduction to the prototype, then allowed the participant to proceed uninterrupted for several minutes interacting with the render. The participant was asked to think aloud during their interactions. Towards the end of this phase, we prompted the participant to interact with any features in

the HTML render they may have skipped over. At the end of this phase, we asked the participant to assess how easy or difficult it was to read the paper with the HTML render.

Phase III: Q&A and discussion

The primary objectives of this phase are to answer the questions:

- How likely is the participant to use the HTML render in the future?
- How can the HTML render be improved to best meet the participant’s needs moving forward?

The participant was given further opportunities to ask questions or discuss the prototype. The participant was asked to describe their perceived pros and cons of the prototype, and to provide suggestions of missing features, ordered by priority. We asked the participant whether they would use this prototype if it were available, and if not, what features would need to be implemented to change that decision.

The interviews were conducted by one author, with two other authors observing and participating during Phase III. All interviews were recorded for followup analysis, and participants were compensated with a \$150 USD gift card for their time. The questions used to guide the semi-structured interview are provided in Appendix D.3.

We follow a grounded theory approach to identify themes and concepts from the participant interviews. We first perform open coding to identify relevant concepts, then axial coding to group these concepts under broad themes. These themes are 1) the technologies employed by users, 2) challenges in their current reading pipeline, and 3) mitigation or coping strategies, and in relation to the SciA11y prototype: 4) positive features, 5) negative features or issues with the prototype, and 6) suggestions for improvement. Interviews are selectively coded a second time to identify all concepts falling under each theme. We also employ the same method to code issues raised by participants in the pre-interview questionnaire.

Themes and concepts are arrived upon by two authors following detailed reading of the interviews. In several cases, we further define attributes associated with some concepts, such as defining whether the technologies used were in relation to opening PDFs, screen reading, or other tasks; or whether the challenges identified affect the whole document, navigation, text, or a particular in-paper element. These delineations are described further in their respective results sections.

6.1 Study participants

Participants are graduate students, PhD students, and faculty members from predominantly English-speaking countries, whose primary research areas are in computer science, though also spanning neuroscience and mathematics. We interviewed two participants during the pilot phase and four participants during the main phase of our study. We report findings from all six participants for all themes captured in Phase I of the study. Since only minor changes were made to the prototype between the pilot and main study, we report findings from all participants for Phase II and III as well, making note of features that changed following the pilots. Three of six participants study human-computer interaction and accessibility, which may be due in part to our sampling methodology, but may also reflect the relevance of accessibility research to BLV researchers. Other study participants conduct research in the areas of machine learning, neuroscience, software engineering, and blockchain. All but one participant reported having more than one year of experience using screen readers. The tools employed by participants are summarized in Table 6 along with the version of the SciA11y prototype with which they interacted.

ID	Study	Prototype Version	Current Tools
P1	Pilot	v0.1	NVDA Screen Reader, Adobe Acrobat Reader
P2*	Pilot	v0.2	Mac Text-to-speech, Mac Magnifying Glass (sighted navigation), Mac Preview
P3	Main	v0.3	Braille display, Mac VoiceOver, JAWS/NVDA on Windows, Mac Preview, Adobe Acrobat Reader
P4	Main	v0.3	Mac VoiceOver, Mac Preview or Adobe Acrobat Reader
P5	Main	v0.3	Microsoft Narrator, Adobe Acrobat Reader
P6	Main	v0.3	Braille display, InftyReader, Mac VoiceOver, Mac Preview

Table 6. User study participants, the prototype versions they interacted with, and the tools they currently use for reading papers. *P2 is low vision and uses sighted navigation tools in conjunction with a screen reader.

6.2 Study findings

Summary of current experience. Of the six participants, three users have experience with screen readers on the Windows OS, such as NVDA, JAWS, and Microsoft Narrator, and three users use VoiceOver on MacOS. Two users use braille display in conjunction with their screen reader. One participant (P2) is low vision and uses a combination of text-to-speech and a magnifying glass to perform sighted navigation; P2’s primary reading interaction involves selecting blocks of text in the PDF and using text-to-speech. Adobe Acrobat Reader is the most common software for opening PDFs; though several participants use Preview in MacOS, with one participant (P4) explicitly stating a preference for Preview over Acrobat. One participant uses a proprietary tool called InftyReader, which converts PDFs into ASCII text and math formulas into MathML, which is accessible.

Challenges of current PDF reading pipeline. Table 7 lists the challenges recognized by all participants in their current PDF reading pipeline. Some of these challenges affect the entire document, e.g., when a document lacks heading markup, it affects the ability to navigate the whole document. Others pertain to specific elements in PDFs, like inaccessible math formulas or lack of figure alt-text. All six users discussed the inaccessibility of math formulas. Unfortunately, document elements like math, figures, tables, and algorithm blocks are used to convey a significant amount of the information content of a paper, and the inability to access their content can produce negative impacts on the reader’s ability to understand the paper.

Coping mechanisms. The coping mechanisms employed by BLV researchers to read inaccessible PDFs are wide-ranging, often involving trying tools outside of their primary workflow, soliciting help from others, or in the worst case, giving up and moving on. We describe these in Table 8. Several users reported trying certain tools like alternate PDF readers, browsers, or optical character recognition (OCR), even though the tools usually do not result in a significant improvement over their standard pipeline; when asked why, several participants reported feeling “hopeful” that a tool might work (P1) or hoping to get lucky (P3).

Several of these coping mechanisms involved other people. For example, three participants reported needing to ask sighted colleagues or family members to copy text, or to explain select paper content, especially figures and equations. Asking for PDF remediation was also a possibility for several participants; in this process, workers at the researcher’s host institution convert a PDF into an accessible format, manually correcting equation representation and writing descriptions for figures. The output of the remediation process is seen as “ideal” (P4), but the process takes significant

Issue description	Affects	Raised by user
Scanned PDFs cannot be read without remediation	Document	P3, P4, P5*
No headings/sub-headings for navigation	Navigation	P1, P3, P5
Figures are not annotated as figures	Navigation	P1, P5
Losing cursor focus when switching away from the PDF	Navigation	P1
Headings are not hierarchical (no sub-headings)	Navigation	P5
Text is read as single string (no spaces or punctuation)	Text	P1, P4, P5
Headers/footers/footnotes mixed into text	Text	P1, P4, P5
Words with ligatures are mispronounced	Text	P1, P3
Words split at line breaks are mispronounced	Text	P2, P3
Reading order is incorrect	Text	P3, P5
Text before and after figures sometimes skipped	Text	P4
Text on some pages not recognized at all	Text	P4
Math content is inaccessible	Element	P1, P2, P3, P4, P5, P6
Tables are inaccessible	Element	P1, P2*, P3, P5, P6
Figures lack alt-text	Element	P1, P3, P5, P6
Figure captions are not associated with figures	Element	P1, P5
Characters or words in figures are read and do not make sense	Element	P4, P5
Figure alt-text (when provided) is not descriptive	Element	P5
Code blocks are inaccessible	Element	P2, P4

Table 7. Challenges to PDF reading identified by participants during interviews. *Only identified as an issue during pre-interview questionnaire.

time (several weeks for any PDF) and may not fit into a researcher’s schedule and timeline. Additionally, this process may only be available to researchers affiliated with a significantly large and resourced institution, and as P6 discusses, may no longer be a viable option for those who work outside of academia. In some cases, BLV researchers may also message authors directly to gain access to the source documents (P3 and P4). Both LaTeX source and Word documents are more accessible than PDFs, and access to these source documents can greatly improve the ability to read these papers.

Perhaps most disheartening is how often BLV researchers may simply give up in the face of an inaccessible paper. P1 says that by the time he has spent several hours making a paper readable, he may have already lost interest and motivation to read it. When asked how often papers are abandoned, P3 responds 60–70% of the time. Though P4 does not discuss abandonment directly, P4 shares the following relevant sentiment: “reading papers is the hardest part of research” for a BLV researcher, and if papers were more accessible, there would be more blind researchers.

Response to HTML render. All user interviews were analyzed to extract positive and negative responses to various features or flaws of the prototype. We summarize these features and flaws in Table 9. Among the participants’ favorite features are links between inline citations and references (all 6 participants), section headings for navigation (5 participants), the table of contents (4 participants), and figures tagged as figures with associated figure captions (3 participants). Regarding links between inline citations and references, several participants were especially supportive of the return links that allow the reader to return back to their reading context after following a citation link. P3 said that the links acted as external memory, allowing BLV users to essentially “glance” at the bibliography and back, like a sighted user might. Similar sentiments were shared by P5 and P6, although P5 also proposed the possibility of preserving

Coping mechanism	Raised by user	What users said
Give up, abandon the paper	P1, P3, P5	P3: when asked how often they abandon papers, answers “60–70% of the time” P5: sometimes the only option is to “sit down and start crying” (jokingly, though the sentiment is true)
Try other conversion tools	P1, P3, P6	
Download LaTeX source or Word document if available	P3, P4, P6	
Ask sighted colleagues or family members to read	P3, P5, P6	
Ask for remediation / convert to braille	P4, P5, P6	P4: 10 day turnaround is on the quick side, which is not good enough for research P5: process takes a long time, around 1-2 weeks
Try other PDF readers or browsers	P1, P6	P1: may try Microsoft Edge browser even though it usually does not help, but he feels “hopeful”
Message authors to get source document	P3, P4	P4: sometimes the author manuscript is accessible but the camera-ready version is not; fault of the conferences and publishers, not the authors

Table 8. Coping mechanisms discussed by users for dealing with challenging papers.

the context even further by providing bibliography information inline rather than navigating back and forth between the main text and references section.

Among the negative features observed by participants, most have to do with imperfect extraction, for example, incorrectly extracted headings (3 participants), missed headings (2 participants), and various extraction issues with code blocks, tables, equations, and more. Many of these issues are known and quantified in Section 5. Of these issues, problems with heading extraction were most notable, likely because the heading structure is the first element of the document with which the participants interact, and it provides a mental model of the overall document structure. Mistakes in heading extraction are obvious and erode trust in our overall system. As P5 says, “it’s really important that I trust it,” and errors of this nature, both false positive and false negative extractions, can reduce trust. Similarly, though we describe in our introductory material that our system currently does not extract equations, P6 points out that it is unclear whether the system extracts equations because occasionally math can be found in the body text. This type of conflict between what is described and what is seen can also reduce trust. However, one may be able to build trust even in the face of extraction errors by indicating to the user when content is not extracted; as P4 says regarding the placeholders for not extracted items, “at least I know there was an equation here.”

Difficulty scale. The responses of the users to the difficulty of their current pipeline versus the HTML render are shown in Table 10. We ask the following question: *On a scale of 1 to 5, how easy or difficult was it to read this paper with the HTML render, and why? (Answers: 1 = Very easy; 2 = Easy; 3 = Neutral; 4 = Difficult; 5 = Very difficult)*

All participants in the main study reported that the HTML render is easier for reading than their current pipeline. Reductions in difficulty rating ranged from 0.5 to 3.0. Most of our participants rated their current pipeline as difficult (4 participants) or neutral (1 participant), with one participant who is low vision (P2) reporting that their current pipeline is easy. During our pilot sessions, users reported that the HTML render was difficult to use. For the main study, users reported the HTML render as neutral or easy to use.

Feature	Raised by user	What users said
POSITIVE		
Bidirectional links between inline citations and references	P1, P2, P3, P4, P5, P6	P3: "very few research teams actually get this and get this right, so well done"; "crucial piece of the puzzle" P4: "Headings are the best thing ever"; makes it very clear what section you are in
Headings for easy navigation	P1, P2, P3, P4, P6	
Table of contents*	P2, P3, P5, P6	P4: "at least I know there was an equation here"
Figures are tagged as figures, and captions are associated	P4, P5, P6	
Can use browser and OS features like find/copy/paste	P1, P4	
Simple typography for reading	P2	
Can interact with headings word-by-word or letter-by-letter	P4	
Not extracted items are noted as missing	P4	
NEGATIVE		
Some headings extracted incorrectly	P1, P3, P5	P5: "it's really important that i trust it"; "there [should be] *no* false negatives"
Some headings missed in extraction	P3, P5	
Code block not extracted	P2, P4	P6: Not sure if this system extracts equations because sometimes there is some math in the body text
Tables are extracted as figures	P2, P6	
Equations not extracted	P4, P6	
Figures placed away from text*	P1	
No alt-text extracted	P1	
URLs missing from bibliography entries**	P2	
Some information not surfaced (keywords, footnotes)	P3	
Some headers/footers/footnotes mixed in text	P4	
Headings are not hierarchical	P5	

Table 9. Positive and negative features identified in the prototype. *The feature was implemented or the issue addressed in v0.2 following P1 pilot. **The issue was addressed in v0.3 following P2 pilot.

P2 is the only participant to report the HTML render as being more difficult to use than their current pipeline; we note that P2 is sighted and did not engage with most of the navigation features we designed and implemented for screen reader-based navigation. Because P2 primarily interacted with papers through sighted navigation, text highlighting, and text-to-speech, they were able to interact with section headers, figures, tables, and equations in the original PDF using the magnifier tool, and found any missing content in the HTML render to be significantly detrimental to their reading experience.

The overall median difference in difficulty scores between the PDF and HTML render is modest, at 0.75. This modest change may be due to the conflation of interface design and system errors when asking participants to rate the difficulty of use. In general, all users responded very positively to the interface design, especially around the navigational features we introduce. Issues were raised around extraction accuracy and the propagation of these errors to the interface. We may be able to offset some of the latter issues by detecting and removing papers that suffer from more extraction errors, though we leave this to future work.

Future usage. At the end of each session, we ask users whether they would be likely to use the prototype in the future if it were made publicly available on a range of papers. We ask specifically: *On a scale of 1 to 5, how likely are you to use the HTML render, if it is available to you in the future? (Answers: 1 = Very unlikely, 2 = Unlikely, 3 = Neutral, 4 = Likely, 5 = Very likely)* If the answer is unlikely or neutral, we ask what changes would need to be made to the tool such that they would use it.

All users reported that they would use the prototype in the future. Five users responded 5, that they would be very likely to use it; one user (P5) responded 3 to the prototype as it currently is, and 5 if some of the issues for heading extraction were addressed. P1, who participated in an early pilot with fewer implemented features, said that this would

ID	Study	Current pipeline	HTML render	Difference	Would use in future
P1	Pilot	4.0	4.0	0.0	Yes
P2	Pilot	2.0	4.0	-2.0	Yes
P3	Main	3.0	2.0	1.0	Yes
P4	Main	4.0	1.0	3.0	Yes
P5	Main	4.0	3.0	1.0	Yes*
P6	Main	4.0	3.5	0.5	Yes

Table 10. Participant ratings on the difficulty scale (1 = very easy, 2 = easy, 3 = neutral, 4 = difficult, 5 = very difficult) and whether they would use the tool in the future. All participants reported a change from more difficult to more easy when moving from their current pipeline to the HTML render except P2, who uses sighted navigation. The median reduction in difficulty score for all participants is 0.75. All participants reported that they would be very likely to use the system in the future were it to be available; P5’s response is contingent on improvements in section heading extraction.

become a tool in the toolbox, but he would not be able to rely solely on it due to incomplete extractions. P5 expressed a similar sentiment, that in its current state, he may try the prototype system when his current workflow fails, but if issues around heading extraction were addressed, he would be very likely to use it. P3 replies when asked how the system might be integrated into their workflow, “I think it would become the workflow.” P4 says “for unaccessible PDFs, this is life-changing.”

6.3 Design recommendations

We distill our learnings into a set of five design recommendations for BLV user-friendly paper reading systems. Figure 11 summarizes the following recommendations:

1. **Document structure should match the mental model of the user.** Structure is necessary for providing an overview of a document and is essential to navigation. Headings in a paper should be tagged as such and the hierarchy of the headings should match the mental model of the user, i.e., top level headings should be tagged `<h1>` or `<h2>`, and lower level headings `<h3>` through `<h6>` accordingly. Reading order should be specified, as to not interject non-body text objects into the body text, e.g., headers, footers, and footnotes often disrupt the main flow of text because they visually break paragraphs. Similarly, a user expects a natural flow to a paper, beginning with the title, authors, abstract, introduction etc, and ending with conclusions and references. Papers with various elements interspersed are disruptive of this mental model and can interfere with the reader’s understanding of the document.
2. **Objects in the paper should be tagged appropriately.** Self-explanatory. Headings should be tagged as headings, figures as figures, tables as tables, lists as lists and so on. Appropriate tagging allows a user to take advantage of the screen reader’s capabilities for navigating to specific types of objects, e.g., most screen readers have shortcuts for navigating headings, and to figures or lists. Proper tagging emulates a sighted user’s ability to detect visually distinct objects such as headings, figures, and tables. When objects are not appropriately tagged, a screen reader user must scroll through the whole document each time to identify the desired sections.
3. **The system should act as external memory for the user.** Visual layout can act as a source of external memory for sighted users, who can quickly derive reading context and object types from visual cues. For BLV users, strategies for emulating such external memory can be beneficial. For example, bi-directional navigation

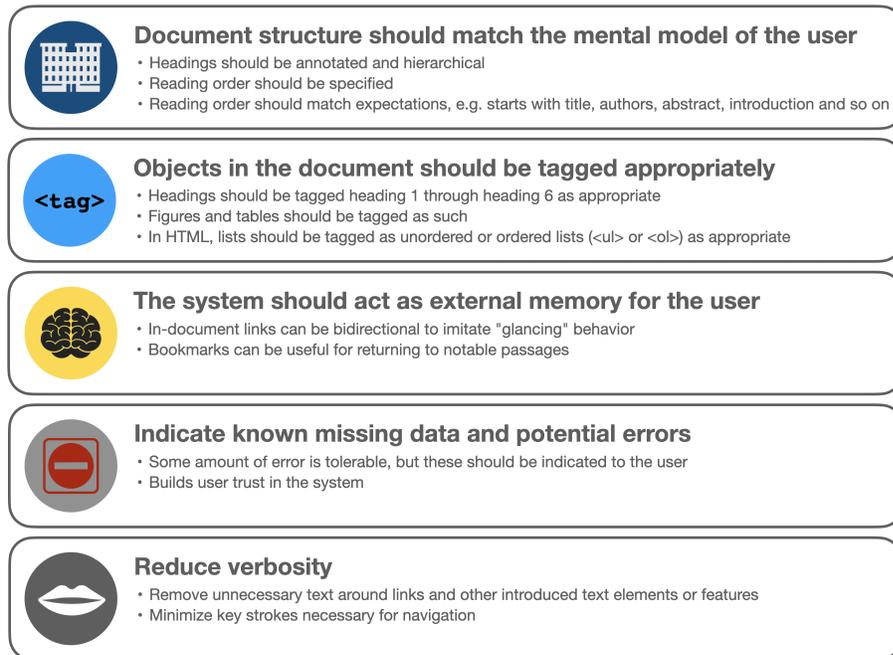


Fig. 11. Design recommendations for screen reader friendly paper reading systems. A system should aim to provide the document structure in a way that matches the mental model of the user, and to tag all elements appropriately. These aspects are achievable through proper tagging of a paper, including in PDF format. Additionally, a system should aim to act as external memory for the user, minimizing the amount of cognitive load needed to return to their reading context. To improve trust, a system should indicate when there is known missing data in the extraction or a possibility of missing or incorrect data. Finally, a system should reduce verbosity, ensuring that as few keystrokes as possible are necessary for the user to perform their desired task.

for all in-document links are a type of “glancing” feature. With this feature, a user no longer needs to commit text to memory in order to rediscover their previous reading context after navigating away. P3, in particular, emphasizes that these features are a “crucial piece of the puzzle.” Other memory features like bookmarking or note-taking may also be helpful for returning the user to their reading context.

4. **Indicate known missing data and potential errors.** To facilitate trust in the system, the system should indicate the presence of missing and erroneous data to users. Some degree of fault tolerance is permitted, as long as the overall benefit to the user is greater. However, as these systems rely on statistical methods, extraction quality is rarely perfect. Most users indicate a preference for knowing when the system fails, rather than dealing with the uncertainty of figuring out whether the issue is with the underlying paper, or with the extraction and reading interface.
5. **Reduce verbosity.** Any minimization of unnecessary text and spaces between links can simplify navigation for BLV users. Though these extra commas and spaces may seem innocuous for sighted users, they require extra keystrokes for screen readers. Reduction of unnecessary verbosity around links and introduced features can save time for screen reader users.

The overarching themes of these recommendations are to reduce user cognitive load and improve trust in the system. Regarding cognitive load, interruptions to reading flow for BLV users are especially disruptive, since there are no visual markers to help identify reading context. Paper reading systems for BLV users should therefore attempt to mitigate cognitive load caused by loss of context, by allowing users to quickly navigate back to their reading context when following any links, and by avoiding any disruption of reading flow. Regarding this latter point, properly labeled reading order, headings for navigation, and appropriately tagged objects all contribute to mitigating disruptions. Further, it is also important to remove interjections from headers, footers, footnotes, figure and table captions, and other text, all of which interrupt the natural flow of reading.

Regarding user trust in the system: this should be a priority of any system builder. Because PDF extraction and document rendering are imperfect processes, some degree of error is expected. Though all participants in our user study expressed that some degree of error is tolerable, one can mitigate the conversion of errors to distrust by clearly indicating known errors and missing content in the system. For example, in some cases our system is unable to extract a figure caption; if the caption for Figure 3 is not extracted, rather than skipping from Figure 2 to Figure 4 and causing confusion for the reader, it is better to indicate that Figure 3 is missing in the extraction.

A system that responds quickly to user requests is obviously more desirable. However, several participants indicated that some wait time is acceptable, especially if a longer wait time corresponds to a higher quality reading experience. Though we report this finding, we ask readers to take it with a grain of salt. This point may not hold for all or even a majority of users, since several users also remark on the PDF remediation process (which usually takes 1–2 weeks) as being too long to adequately support their research workflow.

Though we derive these design recommendations in the scope of paper reading, they are generalizable to other classes of documents. In fact, several of these design principles echo available guidelines for human-AI interaction [2], especially in indicating the capabilities and limitations of the system (recommendation 4). A number of our recommendations are simply good practice, such as exposing the structure of a document and tagging document objects appropriately, and are covered by current guidelines for creating accessible documents. Other recommendations focus on emulating the types of advantages that sighted users derive from layout and visual information, but to implement them in such a way that BLV users can benefit, e.g. using the system as a source of external memory.

7 DISCUSSION

In this work, we present the results of several studies that aim to characterize the current state of accessibility for academic paper PDFs, to learn the challenges faced by BLV researchers when reading papers, and to demonstrate how our SciA11y system that renders PDFs into accessible HTML can be used to mitigate many of these challenges.

Based on our analysis, the current state of paper accessibility is grim, with an average of 2.4% of papers across all fields of study satisfying our five assessed accessibility criteria. Though there is some improvement seen over time, we are not optimistic that these improvements are due to authors prioritizing accessibility when writing papers, since the presence of figure alt-text (the only of the five criteria that requires author intervention) remains low. Rather, the commitment to accessibility made by certain typesetting software providers such as Microsoft Word may be responsible for a portion of these improvements. Given the strong correlation between PDF creation software and accessibility compliance, we encourage conferences, publishers, and authors to consider the tools they are using to generate PDFs, and to integrate accessibility requirements during the publication process.

Given the scope and magnitude of the problem, and how PDF is still the dominant file type used for distributing scientific papers, there are clear needs for immediate technological solutions. We propose the SciA11y system, which

integrates several text and vision machine learning models to extract the content from paper PDFs and render this content as HTML. The system adds tags and infers reading order, thereby improving the navigational capabilities of BLV users. Of course, no extractive pipeline is perfect, and we quantify and qualify extraction quality through an evaluation study and user study. Our intrinsic evaluation of extraction quality indicates that most extractions have no major problems affecting readability (86.2% have no or only some problems). The most common extraction problems are incorrectly extracted or missed section headings, as well as headers, footers, and footnotes being improperly mixed into the body text, which can interrupt reading flow. Participants in our user study responded positively to SciA11y, preferring its navigational features and tagging to working with PDFs. Though the various types of extraction mistakes made by our system are noted by participants, most participants reported an improvement from their current reading pipeline, and all participants expressed an interest in using the system in the future.

We present the challenges, coping mechanisms, and positive and negative features identified by participants. We also summarize the collective themes into a set of five design recommendations for other researchers and practitioners looking to design and build systems for accessible reading. The recommendations include (1) matching the document structure to the mental model of the user, (2) tagging all objects within the document appropriately, (3) acting as external memory for the user, (4) indicating known missing data or extraction errors, and (5) reducing verbosity. The first two of these recommendations are related to proper and correct representation of the document structure and in-paper objects. Both are necessary components of an accessible document. The third recommendation is to provide additional navigation features that are otherwise encoded in the visual layout of the document and inaccessible to BLV users. The fourth recommendation is related to error tolerance and user trust. For any machine learning-based document parsing system, errors are inevitable; managing user expectations for these systems is crucial. This recommendation echoes previously published guidelines for human-AI interaction, which suggest communicating to the user the capabilities and limitations of the AI system [2]. Setting expectations correctly and referring the user back to the original source document when the extractive procedure fails can help mitigate inappropriate reliance on the system. The final recommendation aims to reduce verbosity and the number of keystrokes needed for performing any task, which can speed up the use of such a system.

We hope these design recommendations will facilitate further conversations around the needs of BLV users, and that they may result in systems that ease the reading burden for these users. As one participant puts it, “reading papers is the hardest part of research” for researchers who are blind or low vision, and if papers were more accessible, “there would be more blind researchers.” It is a duty of the entire community to facilitate this, and to design, prototype, and build systems to support the needs of the BLV research community.

7.1 Limitations & Future Work

This work focuses on rendering PDF papers in HTML to improve document navigation and provide a more intuitive reading order. There are many other aspects of accessibility with which we do not contend, such as providing figure alt-text, accessible math, or tagging tables. Future work involves investigating various ways to improve or provide these features automatically, or by harnessing the power of the community to provide some of these features for papers as they are requested. For example, we may integrate element-specific reading features for mathematical equations [4, 17, 26, 45] or graphs and charts [12–14], or create a crowd-sourcing pipeline to solicit alt-text annotations for figures that lack descriptions.

PDF parsing remains an open research problem with many challenges. Our reliance on these technologies necessarily introduce error into our pipeline and system. We attempt to describe and quantify these errors in Appendix C, but

found no strong correlation between any particular type of error and the overall quality assessment. Unfortunately, this means that there is no obvious mitigation strategy for identifying low-quality extractions before they are shown to users. Further work remains to automatically or semi-automatically identify low-quality parses prior to surfacing them. For example, we could investigate other paper features as predictors of parse quality. With more labeled data, we could also train a neural classifier to identify low-quality parses.

In this work, we focus on processing PDFs and making them accessible. Some papers are available in XML, HTML, or other structured markup languages; and LaTeX or Word document source can be found for others. Our system could take advantage of these alternatives to PDFs when they are publicly available, for example, by rendering the semantic content of the paper as extracted from these other document representations, as in arXiv Vanity²³ for arXiv LaTeX source or Pubmed Central's PubReader,²⁴ which renders JATS XML. Though S2ORC [24] contains LaTeX parses derived from arXiv for over 1 million papers, further study is necessary to determine whether these parses are suitable for HTML rendering in our system.

Though we conduct a user study to better understand the challenges of BLV users and their responses to our prototype, the number of participants involved is small. Consequently, we focus on identifying qualitative learnings from these user studies. These learning, when combined with our evaluation and analysis of the current state of scholarly PDF accessibility, provide a more complete portrait of the challenges and issues BLV scholars face when reading papers. To more fully assess the benefits and flaws of our system, a broader user study and testing period is needed. We hope to achieve this in future work.

Lastly, PDFs have been repeatedly called out as being inaccessible, not only for screen readers, but broadly for reading, especially on mobile and other devices with small screen sizes [34]. Dissociating publishing from PDFs continues to be a good goal for the future. In recent years, alternative publication formats have risen in popularity, such as eLife's dual publication in PDF and HTML,²⁵ the interactive HTML papers at distill.pub,²⁶ or the ACM Digital Library's very own dual publication (PDF and HTML) process,²⁷ which is now available for many of the ACM's computing conferences and journals. We have no doubt that viable alternatives to PDF have and will arise, and encourage the community to explore these options when making publication decisions.

8 CONCLUSION

Based on our findings, most academic papers are inaccessible and significant challenges remain for BLV researchers when interacting with and reading these papers. Though some improvements in accessibility have been seen over time, these changes may not be reflective of author actions directly. In the meantime, we offer a potential solution for the millions of PDFs that have already been published and which still remain the dominant form of distribution for academic papers. We introduce the SciA11y system for rendering PDFs as accessible HTML documents. The system extracts the content of PDFs, tagging headings and objects and inferring reading order, which results in a more navigable and accessible document. Though the extraction pipeline is imperfect and can result in errors, our evaluation suggests that for the majority of papers, the resulting HTML render has no major problems that impact readability. We confirm these findings in our user study, where all users responded positively to the prototype system, claiming that they would be likely or very likely to use the system were it to be available in the future. Participants described the system as likely to

²³<http://www.arxiv-vanity.com/>

²⁴<https://www.ncbi.nlm.nih.gov/pmc/about/pubreader/>

²⁵<https://reviewer.lifesciences.org/author-guide/post>

²⁶<https://distill.pub/>

²⁷<https://www.acm.org/publications/authors/submissions>

“become the workflow” or “life-changing,” indicating both a strong favorable response and particular need for these types of solutions.

We do not claim that SciA11y solves all (or even close to all) accessibility problems for BLV researchers, but it is a step in the right direction. SciA11y is a technological solution that can mitigate many of the challenges experienced by BLV researchers at this moment. Though a longer term solution would surely require more dialogue between all stakeholders and a potential revolution in the way in which scholars publish and distribute their research findings, we encourage researchers to prioritize and address these challenges with whatever tools they have in their toolbox right now. We especially encourage others to take into account our findings on the needs and challenges of BLV researchers when designing and engineering new systems and tools for reading the scholarly literature.

ACKNOWLEDGMENTS

This work was supported in part by ONR grant N00014-18-1-2193, NSF RAPID grant 2040196, and the University of Washington WRF/Cable Professorship. We thank Jeff Bigham, Leah Findlater, Jon Froehlich, and Venkatesh Potluri for their valuable feedback on study design and recruitment. We thank Oren Etzioni and Doug Raymond for valuable feedback on the project. We thank Bryan Newbold for providing feedback on earlier drafts of the manuscript. We thank Sam Skjonsberg for help with the demo, and Michal Guerquin and Michael Schmitz for feedback on demo deployment. We thank the Semantic Scholar team for assisting with data access and system infrastructure. Finally, we thank the users who participated in our study, who offered invaluable feedback and suggestions.

REFERENCES

- [1] D. Ahmetovic, T. Armano, C. Bernareggi, M. Berra, A. Capietto, S. Coriasco, N. Murru, Alice Ruighi, and E. Taranto. 2018. Aaccessibility: a LaTeX Package for Mathematical Formulae Accessibility in PDF Documents. *Proceedings of the 20th International ACM SIGACCESS Conference on Computers and Accessibility* (2018).
- [2] Saleema Amershi, Daniel S. Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi T. Iqbal, P. Bennett, Kori Inkpen Quinn, J. Teevan, Ruth Kikin-Gil, and E. Horvitz. 2019. Guidelines for Human-AI Interaction. *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (2019).
- [3] Waleed Ammar, Dirk Groeneveld, Chandra Bhagavatula, Iz Beltagy, Miles Crawford, Doug Downey, Jason Dunkelberger, Ahmed Elgohary, Sergey Feldman, Vu A. Ha, Rodney Michael Kinney, Sebastian Kohlmeier, Kyle Lo, Tyler C. Murray, Hsu-Han Ooi, Matthew E. Peters, Joanna L. Power, Sam Skjonsberg, Lucy Lu Wang, Christopher Wilhelm, Zheng Yuan, Madeleine van Zuylen, and Oren Etzioni. 2018. Construction of the Literature Graph in Semantic Scholar. In *NAACL-HLT*.
- [4] E. Bates and D. Fitzpatrick. 2010. Spoken Mathematics Using Prosody, Earcons and Spearcons. In *ICCHP*.
- [5] Jeffrey P. Bigham. 2014. Making the web easier to see with opportunistic accessibility improvement. *Proceedings of the 27th annual ACM symposium on User interface software and technology* (2014).
- [6] Jeffrey P. Bigham, E. Brady, Cole Gleason, Anhong Guo, and D. Shamma. 2016. An Uninteresting Tour Through Why Our Research Papers Aren't Accessible. *Proceedings of the 2016 CHI Conference Extended Abstracts on Human Factors in Computing Systems* (2016).
- [7] E. Brady, Y. Zhong, and Jeffrey P. Bigham. 2015. Creating accessible PDFs for conference proceedings. *Proceedings of the 12th Web for All Conference* (2015).
- [8] B. Caldwell, M. Cooper, Loretta Guarino Reid, and G. Vanderheiden. 2008. Web Content Accessibility Guidelines (WCAG) 2.0.
- [9] Chen Chen, Ruiyi Zhang, Sungchul Kim, S. Cohen, T. Yu, R. Rossi, and Razvan C. Bunescu. 2019. Neural caption generation over figures. *Adjunct Proceedings of the 2019 ACM International Joint Conference on Pervasive and Ubiquitous Computing and Proceedings of the 2019 ACM International Symposium on Wearable Computers* (2019).
- [10] W. Chisholm, G. Vanderheiden, and Ian Jacobs. 2001. Web content accessibility guidelines 1.0. *Interactions* 8 (2001), 35–54.
- [11] Alireza Darvishy. 2018. PDF Accessibility: Tools and Challenges. In *ICCHP*.
- [12] Stephanie Elzer, E. J. Schwartz, S. Carberry, D. Chester, Seniz Demir, and Peng Wu. 2008. Accessible bar charts for visually impaired users.
- [13] Christin Engel, David Gollasch, Meinhardt Branig, and G. Weber. 2017. Towards Accessible Charts for Blind and Partially Sighted People. In *Mensch & Computer*.
- [14] Christin Engel, E. Müller, and G. Weber. 2019. SVGPlot: an accessible tool to generate highly adaptable, accessible audio-tactile charts for and from blind and visually impaired people. *Proceedings of the 12th ACM International Conference on Pervasive Technologies Related to Assistive Environments*

- (2019).
- [15] Miao Fan and Doo Soon Kim. 2015. Table Region Detection on Large-scale PDF Files without Labeled Data. *ArXiv abs/1506.08891* (2015).
 - [16] H. Ferreira and D. Freitas. 2004. Enhancing the Accessibility of Mathematics for Blind People: The AudioMath Project. In *ICCHP*.
 - [17] S. Flores, M. Andrade-Ar echiga, Alfonso Flores-Barriga, and Juan Lazaro-Flores. 2010. MathML to ASCII-Braille and Hierarchical Tree Converter. In *ICCHP*.
 - [18] Center for Disease Control and Prevention. [n.d.]. The Burden of Vision Loss. <https://www.cdc.gov/visionhealth/risk/burden.htm>. Accessed: 2021-01-31.
 - [19] Cole Gleason, A. Pavel, Emma McCamey, Christina Low, Patrick Carrington, Kris M. Kitani, and Jeffrey P. Bigham. 2020. Twitter A11y: A Browser Extension to Make Twitter Images Accessible. *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (2020).
 - [20] Miquel T ermens i Graells, M. B. Cerrej on, M. D. Boladeras, D. Murillo, P. Asensio, and Mireia Ribera Turr o. 2007. Estudio de la accesibilidad de los documentos cientıficos en soporte digital. *Revista Espanola De Documentacion Cientifica* 31 (2007), 552–572.
 - [21] E. Kim and Kathleen F. McCoy. 2018. Multimodal Deep Learning using Images and Text for Information Graphic Classification. *Proceedings of the 20th International ACM SIGACCESS Conference on Computers and Accessibility* (2018).
 - [22] W. Kruskal and W. A. Wallis. 1952. Use of Ranks in One-Criterion Variance Analysis. *J. Amer. Statist. Assoc.* 47 (1952), 583–621.
 - [23] J. Lazar, E. Churchill, T. Grossman, G. V. D. Veer, Philippe A. Palanque, J. Morris, and Jennifer Mankoff. 2017. Making the field of computing more inclusive. *Commun. ACM* 60 (2017), 50 – 59.
 - [24] Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. 2020. S2ORC: The Semantic Scholar Open Research Corpus. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 4969–4983. <https://doi.org/10.18653/v1/2020.acl-main.447>
 - [25] P. Lopez and Laurent Romary. 2015. GROBID - Information Extraction from Scientific Publications. *ERCIM News* 2015 (2015).
 - [26] M. Mackowski, P. Brzoza, M. Zabka, and D. Spi nczyk. 2017. Multimedia platform for mathematics’ interactive learning accessible to blind people. *Multimedia Tools and Applications* 77 (2017), 6191–6208.
 - [27] Jennifer Mankoff, Anne Spencer Ross, Cynthia Bennett, Katta Spiel, Megan Hofmann, and Jennifer Rode. 2020. 2019 Access SIGCHI Report. *SIGACCESS Access. Comput.* 126, Article 7 (March 2020), 1 pages. <https://doi.org/10.1145/3386280.3386287>
 - [28] M. Maxwell. 1972. Skimming and Scanning Improvement: The Needs, Assumptions and Knowledge Base. *Journal of Literacy Research* 5 (1972), 47 – 59.
 - [29] S. Mirri, S. Peroni, P. Salomoni, F. Vitali, and Vincenzo Rubano. 2017. Towards accessible graphs in HTML-based scientific articles. *2017 14th IEEE Annual Consumer Communications & Networking Conference (CCNC)* (2017), 1067–1072.
 - [30] Jos e M. P. Nascimento and J. Bioucas-Dias. 2005. Vertex component analysis: a fast algorithm to unmix hyperspectral data. *IEEE Transactions on Geoscience and Remote Sensing* 43 (2005), 898–910.
 - [31] A. Nazemi, Iain Murray, and David A. McMeekin. 2014. Practical Segmentation Methods for Logical and Geometric Layout Analysis to Improve Scanned PDF Accessibility to Vision Impaired.
 - [32] A. Nengroo and K. Kuppusamy. 2017. Accessible images (AIMS): a model to build self-describing images for assisting screen reader users. *Universal Access in the Information Society* 17 (2017), 607–619.
 - [33] J. Nganji. 2015. The Portable Document Format (PDF) accessibility practice of four journal publishers. *Library & Information Science Research* 37 (2015), 254–262.
 - [34] Jakob Nielsen and Anna Kaley. 2020. PDF: Still Unfit for Human Consumption, 20 Years Later. <https://www.nngroup.com/articles/pdf-unfit-for-human-consumption/>. Accessed: 2021-01-31.
 - [35] M. Peissner and Rob Edlin-White. 2013. User Control in Adaptive User Interfaces for Accessibility. In *INTERACT*.
 - [36] M. Peissner, Dagmar H abe, Doris Janssen, and T. Sellner. 2012. MyUI: generating accessible user interfaces from multimodal design patterns. In *EICS ’12*.
 - [37] Xin Qian, E. Koh, F. Du, Sungchul Kim, and J. Chan. 2020. A Formative Study on Designing Accurate and Natural Figure Captioning Systems. *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems* (2020).
 - [38] A. J. Rajkumar, J. Lazar, J. B. Jordan, Alireza Darvishy, and H. Hutter. 2020. PDF Accessibility of Research Papers: What Tools are Needed for Assessment and Remediation?. In *HICSS*.
 - [39] Roya Rastan, H. Paik, and John Shepherd. 2019. TEXUS: A unified framework for extracting and understanding tables in PDF documents. *Inf. Process. Manag.* 56 (2019), 895–918.
 - [40] M. Ribera, R. Pozzobon, and S. Sayago. 2019. Publishing accessible proceedings: the DSAI 2016 case study. *Universal Access in the Information Society* (2019), 1–13.
 - [41] Naheda Sahtout. 2020. How science should support researchers with visual impairments. <https://www.cdc.gov/visionhealth/risk/burden.htm>. Accessed: 2021-01-31.
 - [42] Zhihong Shen, Hao Ma, and Kuansan Wang. 2018. A Web-scale system for scientific knowledge exploration. In *Proceedings of ACL 2018, System Demonstrations*. Association for Computational Linguistics, Melbourne, Australia, 87–92. <https://doi.org/10.18653/v1/P18-4015>
 - [43] N. Siegel, Nicholas Lourie, R. Power, and Waleed Ammar. 2018. Extracting Scientific Figures with Distantly Supervised Neural Networks. *Proceedings of the 18th ACM/IEEE on Joint Conference on Digital Libraries* (2018).
 - [44] P. Sojka, M. Ruzicka, Maro  Kucbel, and Martin Jarmar. 2013. Accessibility Issues in Digital Mathematical Libraries.

- [45] V. Sorge, C. Chen, T. Raman, and David Tseng. 2014. Towards making mathematics a first class citizen in general screen readers. In *W4A*.
- [46] Terrill Thompson. 2014. Improving the user interface for people with disabilities. *CHI '14 Extended Abstracts on Human Factors in Computing Systems* (2014).
- [47] Kuansan Wang, Zhihong Shen, Chiyuan Huang, Chieh-Han Wu, Darrin Eide, Yuxiao Dong, Junjie Qian, Anshul Kanakia, Alvin Chen, and Richard Rogahn. 2019. A Review of Microsoft Academic Services for Science of Science Studies. *Frontiers in Big Data* 2 (2019).
- [48] J. Wobbrock, S. Kane, Krzysztof Z Gajos, S. Harada, and Jon Froehlich. 2011. Ability-Based Design: Concept, Principles and Examples. *ACM Trans. Access. Comput.* 3 (2011), 9:1–9:27.

A EVALUATION FORMS

This section contains forms and documents used to evaluate the quality of HTML renders produced by our system.

A.1 Evaluation instructions

Instructions for annotators are reproduced verbatim below.

Goal: Identify and quantify the prevalence of different parse issues in S2ORC parses to assess their suitability for accessibility purposes. This will help us decide whether S2ORC parses can help meet screen reader accessibility needs.

Number of papers: 500 papers sampled across different domains of science

You will be presented with a spreadsheet of scientific papers, each with a pair of links. One link goes to a PDF of the paper. One link goes to an HTML representation of the same article. For each pair of links, we would like to know how faithfully the HTML representation captures the information on the PDF.

INSTRUCTIONS:

1. Open the two links side by side.
2. If the two links do not seem to correspond to the same paper, STOP. Make a note in the spreadsheet and SKIP.
3. If the PDF shows a paper that is not suitable, STOP. Make a note in the spreadsheet and SKIP. Non-suitable may include:
 - It is not a scientific paper.
 - It is spam or a fake paper.
 - It is slides, a poster, or other such non-paper document.
 - It is just an abstract.
 - It is a series of articles (e.g. conference proceedings, journal issue etc).
 - It is a book.
 - It is supplementary material;
note: some supplementary material is solely made up of figure or images.
 - Something else that makes you pause. If you're not sure, SKIP it.
4. Copy the paper identifier corresponding to this paper into the first question on this form. Please make sure the identifier matches the paper you are evaluating.
5. Answer each of the questions in this form as best as you can, treating the PDF as gold. There is no need to review every word or line of text. We are just trying to get an overall assessment of parse quality. For any question that asks for a number, enter '0' if there are no obvious problems with those extractions.
6. Submit the form and mark the row in the spreadsheet as complete.

Note: Display equations (those that are in their own paragraph) are currently not preserved in S2ORC, so we ask the annotator to ignore issues around missing display equations. Inline equations (those that are inline within a paragraph) are converted to token streams, which may not be faithful to the original PDF (e.g. fractions may not be preserved). The annotator can provide a description of issues around equation parsing when there is a notable issue.

A.2 Evaluation questions

Questions asked in the evaluation form are reproduced in Table 11.

A.3 Quality rubric

The quality rubric for the final question in the evaluation form is given in Table 12. This rating attempts to capture the overall readability and usability of the HTML render. Three authors discussed and converged upon this rubric following initial pilot annotations.

B EVALUATION RESULTS

Raw counts for each type of error detected during the evaluation of HTML renders are provided in Table 13. The overall quality score split by field of study is shown in Table 14.

C ASSOCIATION BETWEEN PAPER FEATURES AND OVERALL READABILITY

To investigate the possibility of identifying paper extractions with major problems, we fit a Logistic Regression classifier using element specific evaluation results as input features, and whether or not a paper has major problems as the target class for classification. Element specific questions are converted into 43 binary input variables; for example, the title element is mapped to three binary variables, whether the title is extracted correctly (`title_yes`), extracted partially (`title_partially`), or extracted incorrectly (`title_no`). We collapse the targets into two binary classes, 1 if the paper has major problems, and 0 if it has no major problems or some problems. The classifier is trained using 5-fold cross validation, with balanced class weights, and achieves a mean accuracy of 0.69, and area under the ROC of 0.65. This performance is not particularly notable or good; the labeled training sample is small, and due to the complexity of what makes a document problematic to read, we did not expect there to be a clear way to predict extractions with major problems based on a small number of element-level features. Something we aim to explore more in the future is whether the raw tokens on the PDF or publisher metadata can be leveraged to better predict when our extractive parse has failed.

The top 10 predictive features and their logistic regression weights are:

Abstract extracted incorrectly:	0.42
One table extraction error:	0.18
No table caption errors:	0.15
One figure extraction error:	0.15
One figure caption extraction error:	0.15
Bibliography extraction is very bad:	0.13
More than one figure extraction error:	0.13
Authors extracted incorrectly:	0.12
More than one table extraction error:	0.11
Authors extracted correctly:	0.09

The most predictive feature is when abstracts are extracted incorrectly. Given the prevalence of abstracts in various literature databases, abstract quality could be easily assessed through external verification. In other words, if the abstract we extract is different from the abstract found for the same paper in other databases or in the publisher metadata, perhaps we can avoid surfacing this paper. However, the distribution of weights among various other element-level features suggests that this feature alone would be insufficient, and that the contributions of these various features are complex, denying us an easy way of identifying paper parses with major problems.

D USER STUDY MATERIALS

Documents used for the user study are provided in this Appendix.

D.1 Recruitment email

The following email was sent and forwarded to several mailing lists to recruit participants.

The anon_corpus Research Team at the anonymized is conducting an experiment to evaluate the screen reader accessibility of scientific papers.

We are looking for participants who are age 18 or older, who identify as blind or low vision, and who have experience using screen readers to interact with scientific papers. If you are interested in participating, please complete the following form to determine eligibility: [link](#)

Participation in this study is entirely voluntary. If you do decide to participate, your individual data will be kept strictly confidential and will be stored without personal identifiers.

The study involves an informational interview to better understand screen reader needs around scientific papers. Each participant will also be asked to interact with papers on a web interface developed by the team. The study will take approximately 75 minutes, and participants will receive a \$150 Amazon gift card for their time.

Location: Online (Zoom)

Please contact the authors if you have any questions or concerns about this study. Thank you in advance for your time! Please help us spread the word by forwarding as appropriate.

D.2 Pre-interview questionnaire

Prior to each user study interview, the participant was asked to complete the following form:

Share 3 to 5 scientific papers that are difficult to read due to accessibility issues

Thank you for volunteering to take part in this study! Please take a few minutes to supply us with some subject keywords you are interested in, and a list of 3 to 5 scientific papers you have found difficult to read due to accessibility issues. This would help us better plan the study based on your experience.

1. Your name (First name, last initial)

2. Please give a few examples of subject keywords you care about.

For example, computing hardware, analog computer, etc.

3. Share one paper you have had difficulty reading due to accessibility issues by answering the following questions.

- Paper title & link

For example, "What every Researcher should know about Searching - Clarified Concepts, Search Advice, and an Agenda to improve Finding in Academia" (<https://pubmed.ncbi.nlm.nih.gov/33031639/>)

- On a scale of 1 to 5, how easy or difficult was it for you to read this paper?

(1 = very easy, 5 = very difficult)

- Briefly describe why you chose the rating

4-7. *Repeat 3.*

D.3 Interview questions

The following discussion guide is used to provide structure for user interviews.

Phase I – Warmup:

- Can you tell us a little bit about yourself? (Background, what kind of research do you do)
- Tell us about how you normally read papers - What is your workflow like? What tools do you use?
 - Do you usually read PDFs directly or do you read papers in other ways?
- If you need to read a paper and it is not accessible, what do you do now?
 - How long does the process take?
 - How often is it successful?
- Can you give a few examples of the main challenges you face when reading papers? (For example, are there certain features or attributes of papers that make them particularly difficult to read?)
- In your opinion, are there any resources that provide papers that are easier to read by screen readers? (For example, any journals, conferences, or search engines?)
- Overall, how do you feel about your current experience of reading papers?

Phase I – Current workflow:

- Based on the list of papers you provided, walk us through how you would read the paper [abc]. Use the screen reader of your choice, and any additional tools or extensions that are part of your usual process.
- Instructions: Please share your whole screen, think aloud and walk me through your thinking process
- What kind of information were you looking for, and how did you explore the page to find the information?
- On a scale of 1 to 5, how easy or difficult was it to read this paper with the tools, and why? (1 = very easy, 5 = very difficult)
- If you could change anything, how could this best meet your needs?

Phase II:

- We are currently working on an experimental prototype to make papers more easily read by screen readers. Please take a minute to read the about page first: [link](#)

- Based on the list of papers you provided, walk us through how you would read the paper [abc] using this HTML render. You can also use the screen reader of your choice, and any additional tools or extensions that are part of your usual process.
- Instructions:
 - We are working with prototypes so not everything works
 - Please think aloud and walk me through your thinking process
 - Feel free to provide as many feedback as you can, good or bad
- Please take a few more minutes to explore the other parts of this prototype. (e.g. References)
- On a scale of 1 to 5, how easy or difficult was it to read this paper with the tools, and why? (1 = very easy, 5 = very difficult)

Phase III:

- On a scale of 1 to 5, how likely are you to use the HTML render, if it is available to you in the future? (1 = very unlikely, 5 = very likely)
- Which features do you consider to be most helpful?
- Is there anything it would need to have, or change to convince you to use it?
- How do you envision yourself using this tool? How might it fit into your workflow? (For example, would it be an additional extension that is part of your usual process?)
- If you could search for papers and view them in this format, what do you think?
- If you could upload any PDF and create an HTML page like this, what do you think? (Would that be helpful for you, or something you might use, why or why not?)
- Are you aware of any other tools that display papers in any way besides PDF?
- Do you have any additional feedback about the HTML render or anything else that you would like to share?
- Thank you

Answer	Questions
y/p/n text	Is the TITLE correctly extracted? <i>Comment (clarify if answer is “partially” or “no”)</i>
y/p/n text	Are the AUTHOR(S) correctly extracted? <i>Comment (clarify if answer is “partially” or “no”)</i>
y/p/n text	Is the ABSTRACT correctly extracted? <i>Comment (clarify if answer is “partially” or “no”)</i>
y/n	Does this paper contain a substantial number of math EQUATIONS (more than 5 display equations)?
number	How many FIGURES are in the PDF? (Enter ‘0’ if none)
number	How many FIGURES are correctly extracted? (Enter ‘0’ if no figures in paper)
number	How many FIGURE CAPTIONS are correctly extracted? (Enter ‘0’ if no figures in paper)
number	Approximately how many FIGURE captions are **INCORRECTLY** parsed into the body text (should be a figure caption but is mixed in with the body text)? (Enter ‘0’ if they are all correct or if no figures)
text	<i>Comment (Optional – Note anything here about FIGURES or FIGURE CAPTIONS, e.g. which figures are not extracted, which figure captions are not extracted, which figure captions are incorrectly extracted into the body text etc.)</i>
number	How many TABLES are in the PDF? (Enter ‘0’ if none)
number	How many TABLES are correctly extracted? (Enter ‘0’ if no tables in paper)
number	How many TABLE TITLES / CAPTIONS are correctly extracted? (Enter ‘0’ if no tables in paper or if those tables do not have titles / captions)
number	Approximately how many TABLE titles are **INCORRECTLY** parsed into the body text (should be a table title / caption but is mixed in with the body text)? (Enter ‘0’ if they are all correct or if no tables or table titles / captions)
number	Approximately how many TABLES have content that is **INCORRECTLY** parsed into the body text (content of table is mixed with the body text)? (Enter ‘0’ if they are all correct or if no tables)
text	<i>Comment (Optional – Note anything here about TABLES or TABLE TITLES / CAPTIONS, e.g. which tables are not extracted, which table captions are not extracted, which table title / captions / content are incorrectly extracted into the body text etc.)</i>
number	Approximately how many times are page HEADERS or FOOTERS **INCORRECTLY** mixed into the body text? This also includes margin content such as arXiv watermarks. (Enter ‘0’ if all okay or no headers or footers)
text	<i>Comment (Optional – Note anything interesting here about incorrectly parsed headers or footers; no need to provide page numbers)</i>
number	Approximately how many SECTION HEADINGS are **INCORRECTLY** extracted? (Enter ‘0’ if they are all correct or no section headings)
text	<i>Comment (Optional – Note anything interesting about the section heading extractions, no need to list exhaustively)</i>
number	Approximately how many BODY TEXT PARAGRAPHS are **MISSING** from the extraction? (Enter ‘0’ if they are all there or there is no body text)
text	<i>Comment (Optional – Note anything interesting about the body text extractions)</i>
choice	Are BIBLIOGRAPHY entries extracted correctly? (options: all correct, mostly correct, half correct, mostly incorrect, incorrect, no bibliography)
text	<i>Comment (Optional – Note anything interesting about the bibliography extractions; no need to list exhaustively)</i>
choice	Are INLINE CITATIONS linked to bibliography entries? (Please answer this questions considering only the bibliography entries that were extracted) (options: all linked, majority linked, half linked, most unlinked, none linked, no bibliography)
text	<i>Comment (Optional – Note anything interesting about the inline citation linking; no need to list exhaustively)</i>
text	<i>Are there any other problems with the HTML parse that are not covered by one of the above questions? Please describe. (Optional)</i>
choice	Please rate the overall full text quality in the HTML render (options: no major problems, some problems, lots of problems – see rubric in Section A.3)

Table 11. Evaluation questions. Optional questions are in *italics*.

Rating	Criteria
No major problems that impact readability	<ul style="list-style-type: none"> • No errors or relatively few errors • No missing paragraphs, but a few insertions into paragraphs or incorrect headers okay • Any errors impact only a couple of paragraphs
Some problems that impact readability	<ul style="list-style-type: none"> • Few missing paragraphs (<1 per 5 pages) OR Several figure/table insertions into paragraphs or incorrect headers • Errors can impact multiple paragraphs
Lots of problems that impact readability	<ul style="list-style-type: none"> • Difficult to read • Multiple missing paragraphs OR multiple figure/table insertions that make some paragraphs unreadable • Errors impact majority of paragraphs

Table 12. Rubric for HTML parse quality assessment (final question in evaluation questionnaire).

Metadata Element	Yes	Partially	No		
Title	337	16	32		
Authors	307	64	14		
Abstract	308	22	55		
Figure/Table Element	Skipped	No figures/tables	No errors	1 error	>1 error
Figure extraction errors	6	94	201	45	39
Figure caption errors	0	94	174	55	62
Table extraction errors	2	166	165	32	20
Table caption errors	2	166	190	23	4
Text Element	Skipped	No errors	1-5 errors	>5 errors	
Header/Footer/Footnote errors	3	170	172	40	
Section heading errors	2	88	258	37	
Body paragraph errors	1	226	128	30	
Bibliography Element	Skipped/poor bib extraction	No bibliography	All or most correct	Half correct	Mostly incorrect
Bibliography extraction	7	15	313	3	47
Inline citation linking	39	10	290	20	26
Overall Readability	Good	Okay	Bad		
Overall score	210	122	53		

Table 13. Assessment count for all evaluation paper elements. Corresponds to distributions shown in Figure 9.

Overall Readability	Number of papers	Good	Okay	Bad
All papers	385	210	122	53
Art	13	6	1	6
Biology	23	12	7	4
Business	14	6	2	6
Chemistry	19	12	5	2
Computer science	21	10	7	4
Economics	20	6	8	6
Engineering	23	15	7	1
Environmental science	18	7	8	3
Geography	17	9	6	2
Geology	21	12	8	1
History	7	5	1	1
Materials science	24	15	8	1
Mathematics	25	13	8	4
Medicine	26	14	12	0
Other	8	6	2	0
Philosophy	12	7	5	0
Physics	39	25	10	4
Political science	13	6	6	1
Psychology	22	11	7	4
Sociology	20	13	4	3

Table 14. Distribution of overall quality scores for readability, split by field of study. Corresponds to distributions shown in Figure 10.

Class	Precision	Recall	F1-score	Support
No major problems / Some problems	0.91	0.71	0.80	332
Major problems	0.23	0.55	0.32	53

Table 15. Precision, recall, and F1-scores for classification. The classifier does not perform well at identifying papers with major problems from element-based features (F1 = 0.32).