

Let's Make PDFs Accessible for All: A Survey on PDF-Accessibility in Switzerland



Prof. Dr.
Alireza Darvishy



Prof. Dr.
Rolf Sethe



Ines
Engler



Felix
Schmitt-Koopmann



Oriane
Pierrès

Agenda

Part I: Introduction

- What makes a PDF document accessible?
- The problem with accessible PDFs
- PDF Tags
- Creating a PDF UA
- PDF remediation with AI

Part II: The state of accessible scientific PDFs in Switzerland

Before we begin...

What makes a PDF document accessible?

Introduction

Example: Reading Order

Accessible PDFs: Applying Artificial Intelligence for Automated Remediation of STEM PDFs

Felix M. Schmitt-Koopmann
University of Zurich, Department of Informatics and Zurich University of Applied Science, Institute of Applied Information Technology
fschmitt@ifi.uzh.ch

Prof. Dr. Elaine M. Huang
University of Zurich, Department of Informatics
huang@ifi.uzh.ch

Prof. Dr. Alireza Darvishy
Zurich University of Applied Science, Institute of Applied Information Technology
alireza.darvishy@zhaw.ch

ABSTRACT
People with visual impairments use assistive technology, e.g., screen readers, to navigate and read PDFs. However, such screen readers need extra information about the logical structure of the PDF, such as the reading order, header levels, and mathematical formulas, described in readable form to navigate the document in a meaningful way. This logical structure can be added to a PDF with tags. Creating tags for a PDF is time-consuming, and requires awareness and expert knowledge. Hence, most PDFs are left untagged, and as a result, they are poorly readable or unreadable for people who rely on screen readers. STEM documents are particularly problematic with their complex document structure and complicated mathematical formulae. These inaccessible PDFs present a major barrier for people with visual impairments wishing to pursue studies or careers in STEM fields, who cannot easily read studies and publications from their field. The goal of this Ph.D. is to apply artificial intelligence for document analysis to reasonably automate the remediation process of PDFs and present a solution for large mathematical formulae accessibility in PDFs. With these new methods, the Ph.D. research aims to lower barriers to creating accessible scientific PDFs, by reducing the time, effort, and expertise necessary to do so, ultimately facilitating greater access to scientific documents for people with visual impairments.

CCS CONCEPTS
• **Human-centered computing** → Accessibility; Accessibility systems and tools; Accessibility; Accessibility technologies; • **Applied computing** → Document management and text processing; Document capture; Document analysis.

KEYWORDS
Accessibility, PDF Accessibility, Tagged PDF, PDF/UA, Math Viewer, Document Analysis, Formula Recognition, Page Object Detection, Reading Order

ACM Reference Format:
Felix M. Schmitt-Koopmann, Prof. Dr. Elaine M. Huang, and Prof. Dr. Alireza Darvishy. 2022. Accessible PDFs: Applying Artificial Intelligence

This work is licensed under a Creative Commons Attribution International 4.0 License.

ASSETS '22, October 23–26, 2022, Athens, Greece
© 2022 Copyright held by the owner/authors(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9228-9/22/10
<https://doi.org/10.1145/3517428.3550407>

for Automated Remediation of STEM PDFs. In *The 24th International ACM SIGACCESS Conference on Computers and Accessibility (ASSETS '22)*, October 23–26, 2022, Athens, Greece. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/3517428.3550407>

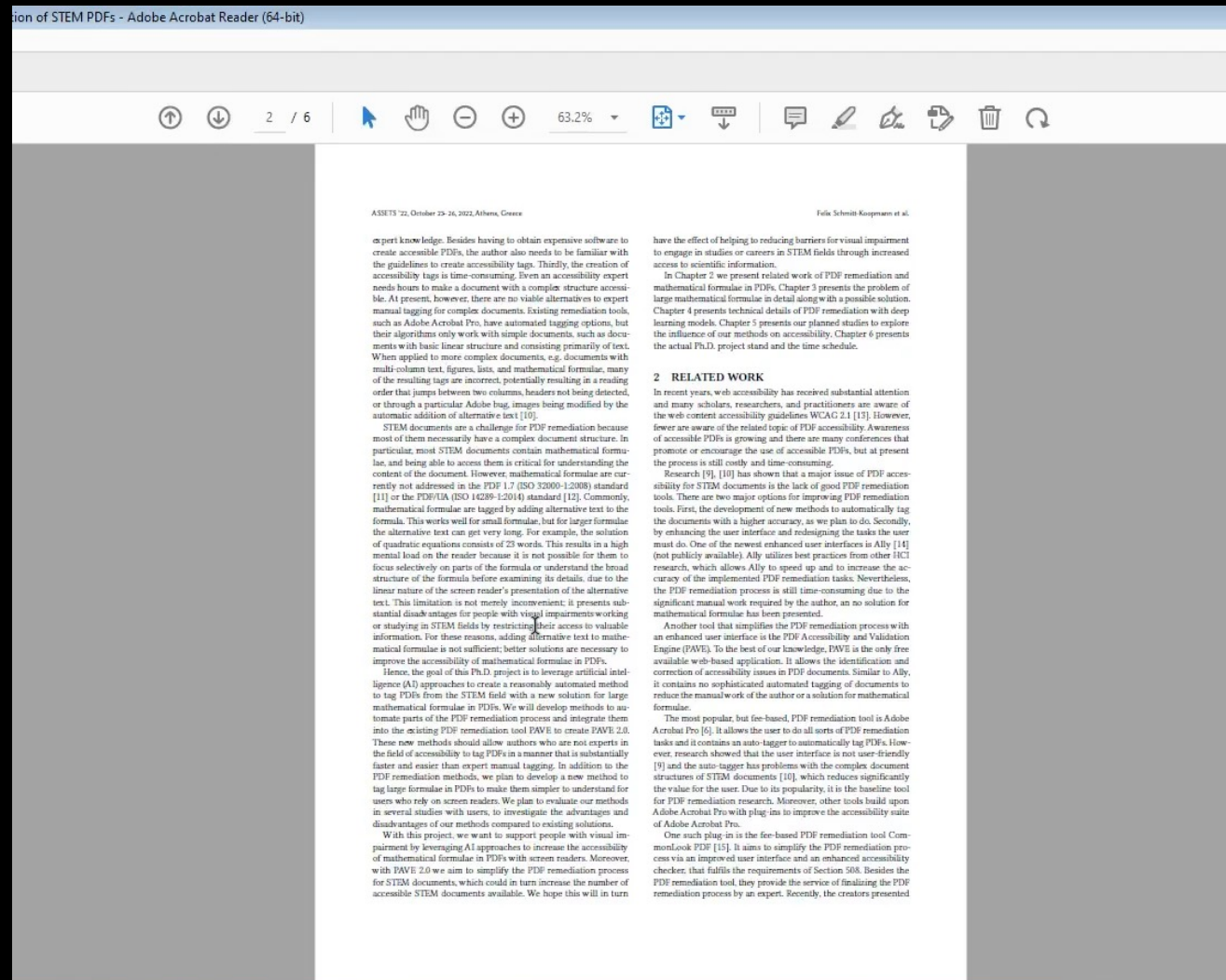
1 INTRODUCTION
Since 1998, the US Rehabilitation Act section 508 [1] requires US Federal departments and agencies to make electronic and information technology accessible to people with disabilities. Additionally, the 2008 United Nations Convention on the Rights of Persons with Disabilities [2], and the 2019 European Accessibility Act [3], require that critical products and services be usable by people with disabilities. The members of the European Union must implement these requirements by 2025. One element of these acts is document accessibility.

The Portable Document Format (PDF) is the most popular document format, especially for scientific papers. Adobe created it in 1993, and since 2008 it is an open format managed by the PDF Association [4]. The PDF format was developed to display documents independent of the software and hardware used, which is one of the reasons for the format's popularity. In 2012, the PDF Association introduced the ISO 14289 standard, which is better known as the PDF/Universal Accessibility (UA) standard. It specifies that a PDF must be tagged to be accessible with assistive tools, such as screen readers. The tags contain information about the logical structure of the PDF, e.g. what is the header and which header level is it. This logical structure allows screen readers to correctly process content objects, such as headers, tables, and lists, and read the objects in the correct reading order. However, most PDFs do not meet the UA standard, and therefore are not easily readable for people with visual impairments who rely on screen readers.

Different tools exist to create accessible PDFs. These tools can be separated into two groups. The first group allows the tagging of existing PDFs, a process known as PDF remediation. The second group supports the generation of tags during the creation of the PDF e.g., with special add-ins for Microsoft Word or Microsoft PowerPoint [5]. In this Ph.D. project, we investigate PDF remediation, because PDF remediation can be applied to all PDFs and it is not software specific. For PDF remediation, the creator of the PDF can use programs such as Adobe Acrobat Pro [6], FAVE [7], [8], and others to tag their PDFs.

Nevertheless, most authors do not use these tools to create accessible PDFs. Research has shown there are three main reasons why many PDFs contain no tags [9]. The first reason is that authors lack awareness about accessible PDFs and do not know this problem exists. A second related reason is that PDF remediation requires

Example: Headings

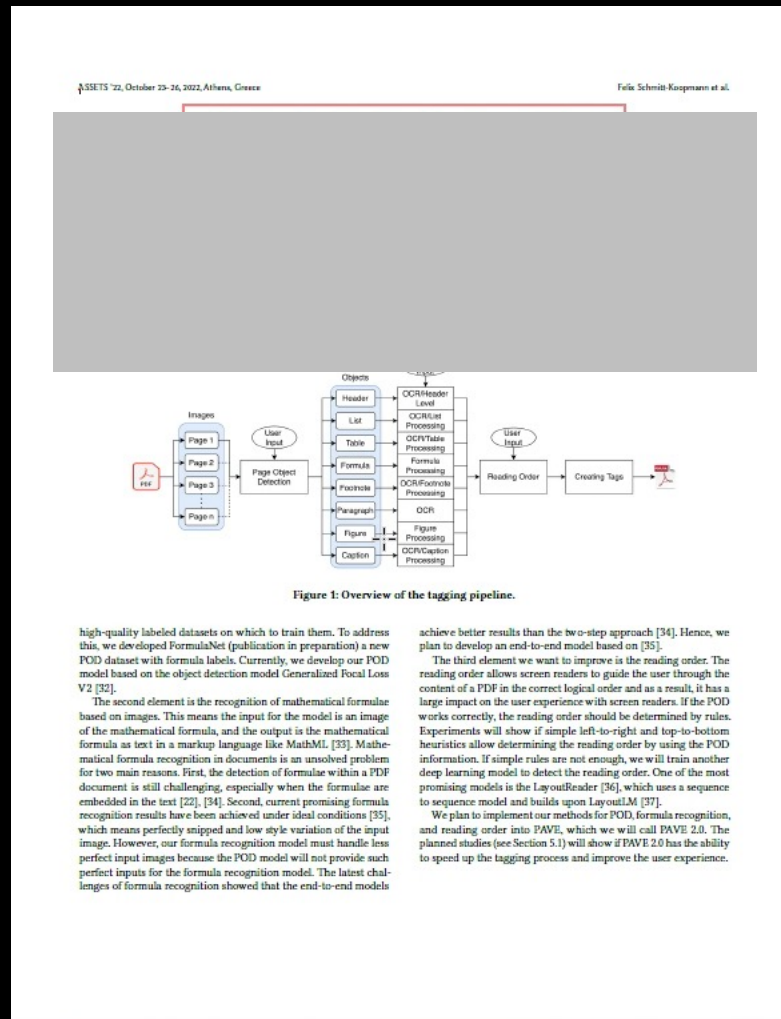


Example: Visuals



P8 – Make visuals accessible

Example: Tables



Example: Mathematical Formulas



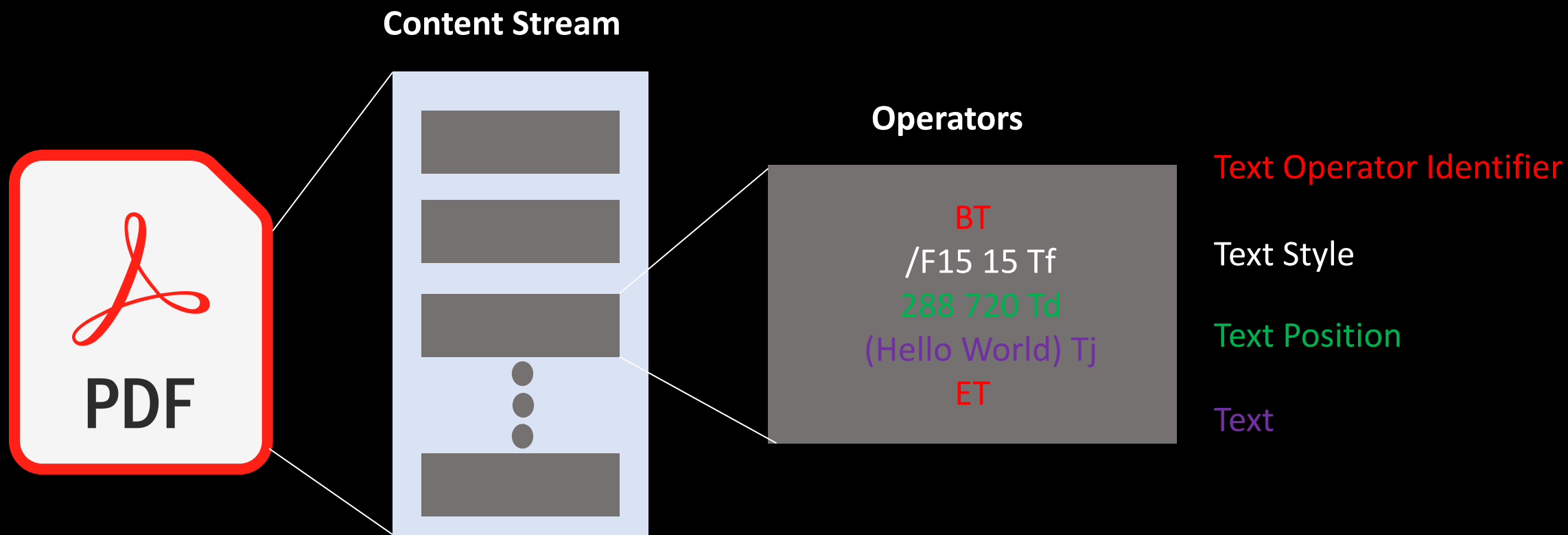
one can result in substantial cognitive load on the person with a visual impairment.

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} \quad (1)$$

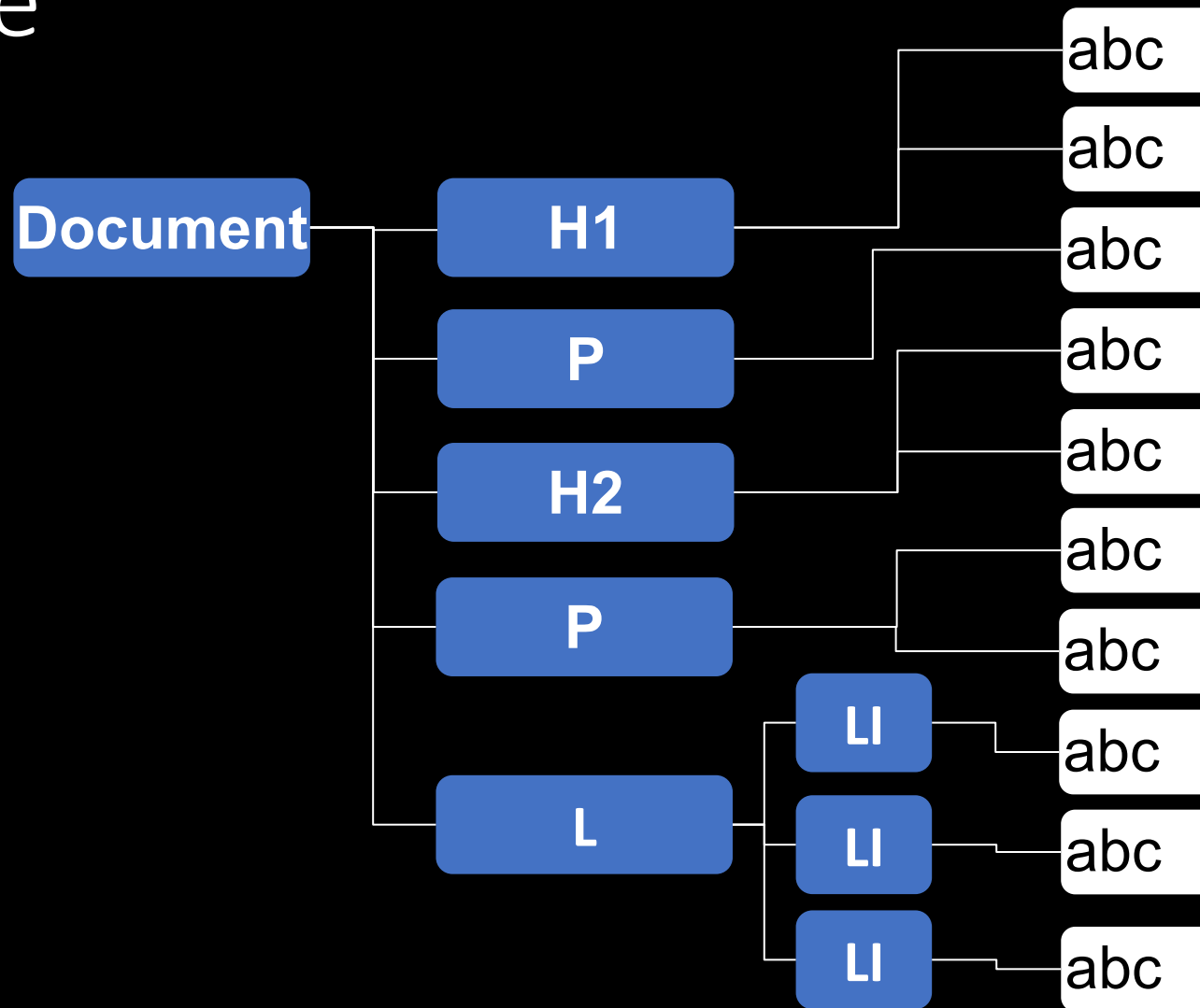
For other formats such as websites, there are solutions with “math viewers,” such as the math viewer provided with JAWS [30]

What have you noticed about the slides?

Why is the PDF format problematic?



Struct Tree



Heading – Hn tags

Heading 1 (H1)

Heading 1.1 (H2)

Heading 1.1.1 (H3)

Heading 1.2 (H2)

Heading 2 (H1)

H1

H2

H3

H4

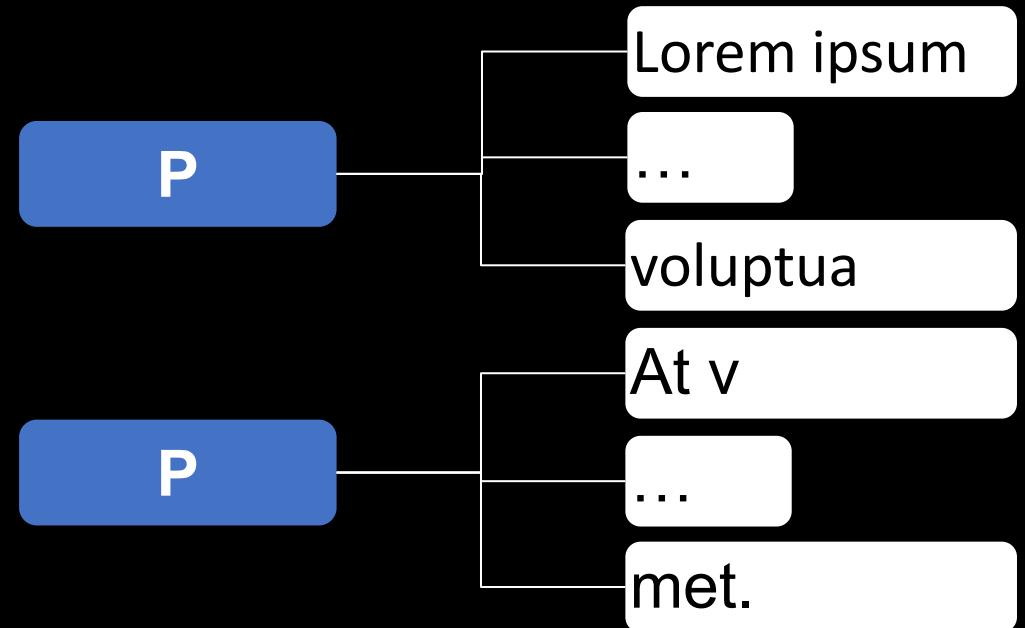
H5

H6

Paragraph – P tags

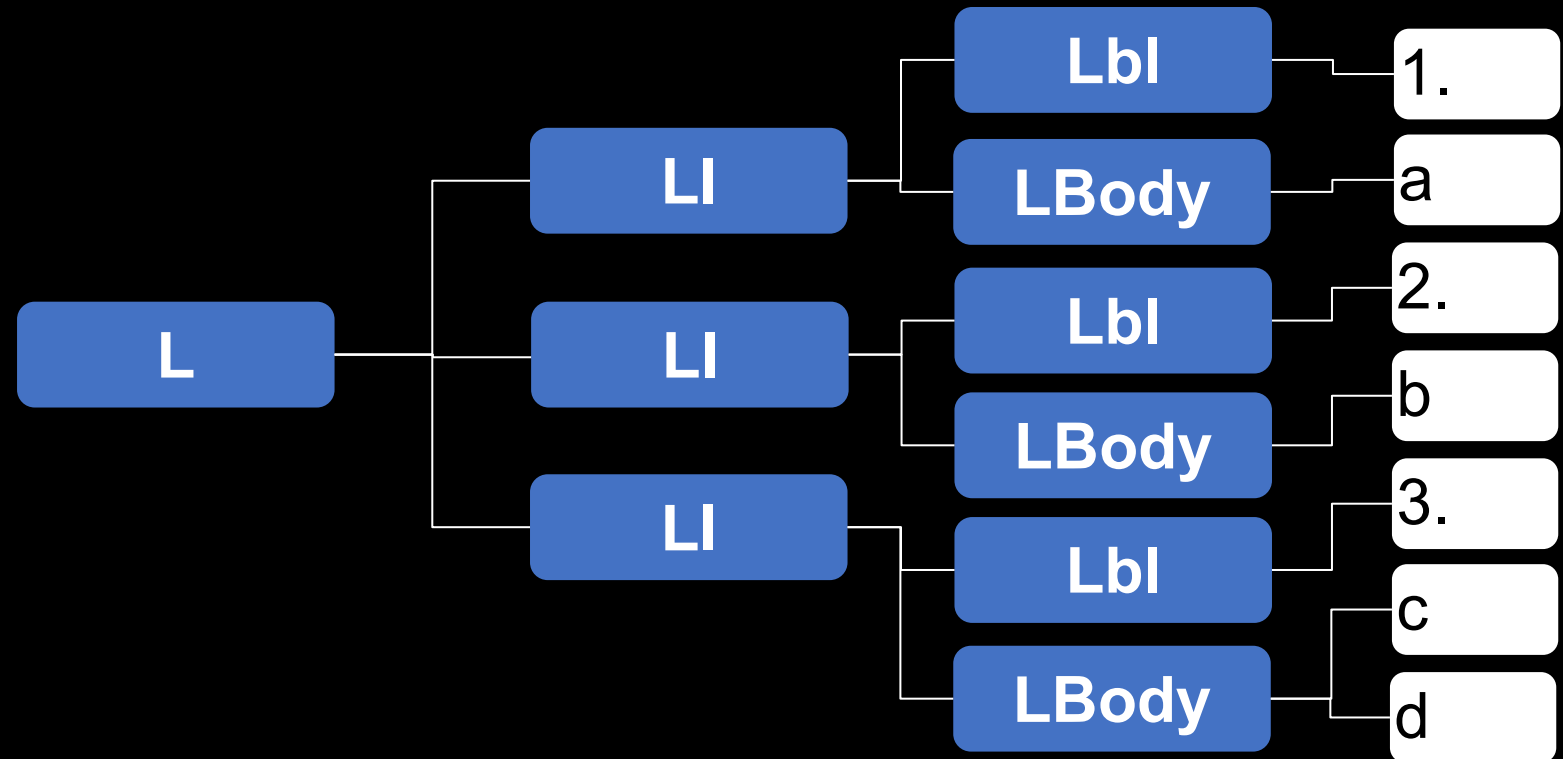
(P) Lorem ipsum dolor sit amet, consetetur sadipscing elitr, sed diam nonumy eirmod tempor invidunt ut labore et dolore magna aliquyam erat, sed diam voluptua.

(P) At vero eos et accusam et justo duo dolores et ea rebum. Stet clita kasd gubergren, no sea takimata sanctus est Lorem ipsum dolor sit amet.



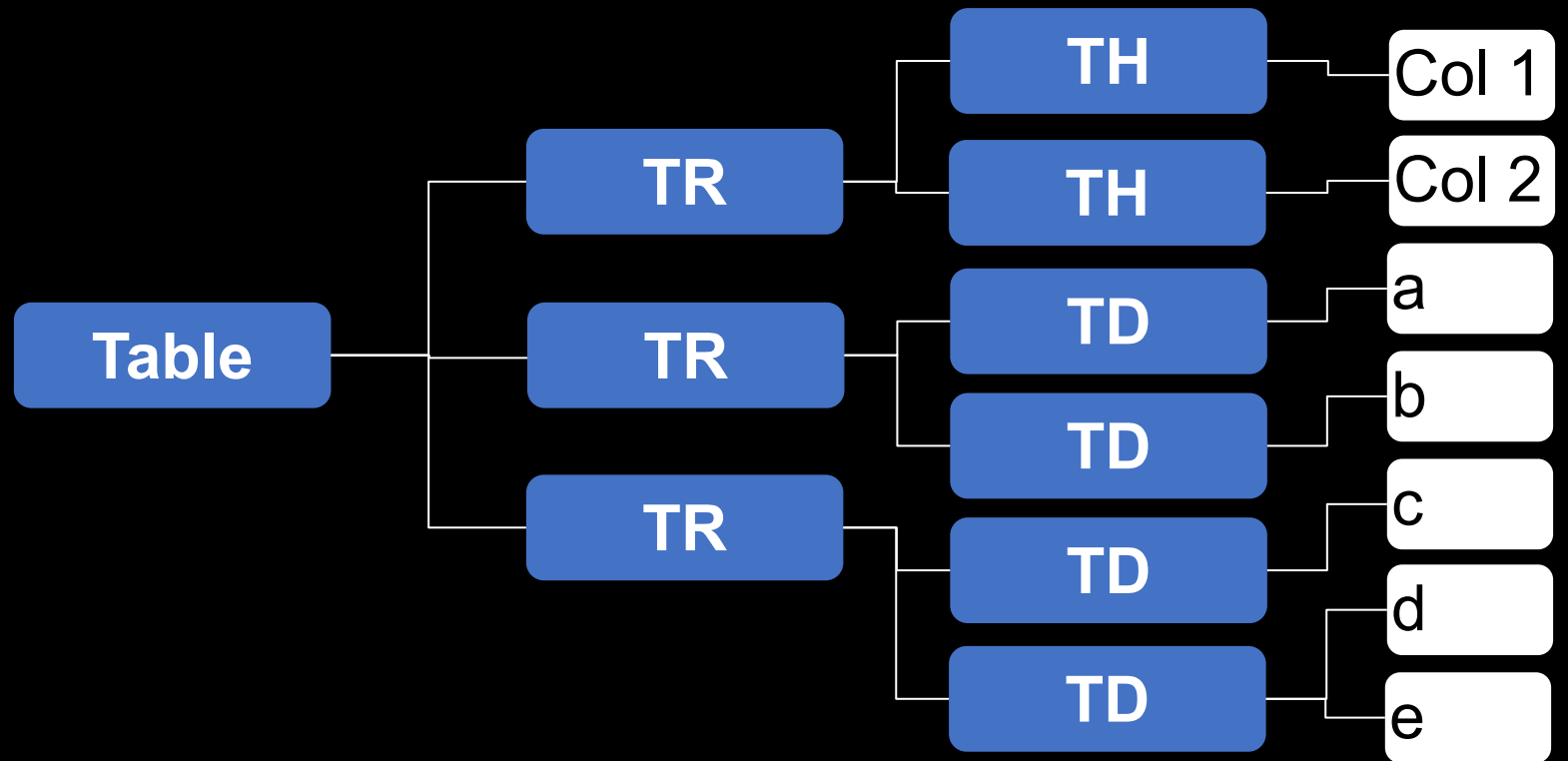
Lists – L tags

- 1. a
- 2. b
- 3. cd

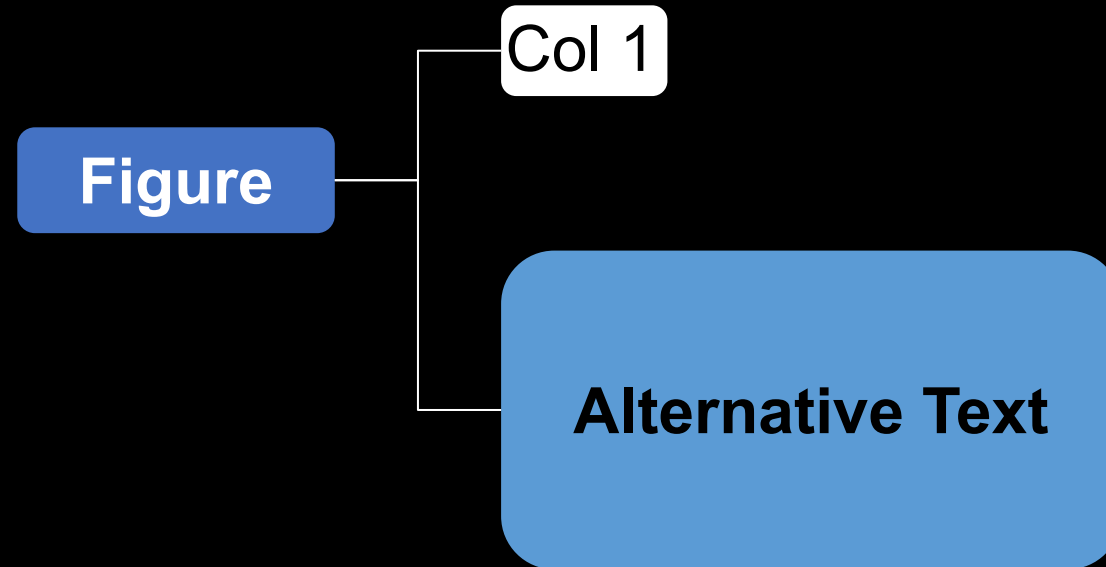


Tables – table tags

Col 1	Col 2
a	b
c	de



Figures – figure tags



Other tags

The screenshot shows the 'Accessible PDF' website. The header includes the site logo, the title 'Accessible PDF', the subtitle 'How to create PDF/UA compliant documents', and navigation links for 'Basics', 'Troubleshooting', 'Glossary', 'EN', and 'DE'. The main content area is divided into a left sidebar and a main article. The sidebar has a 'Basics articles' section with 'General tutorials' and 'InDesign specific tutorials'. Under 'General tutorials', there are three items: 'Structure with the help of multi-level headings', 'Check semantics and logical reading order', and 'Overview of the PDF tags' (which is highlighted with a red border). The main article is titled 'Overview of the PDF tags', written by Stefan Brechbühl, and last updated on 12/7/2020. The article text states: 'This overview shows the most important tags from the PDF 1.7 standard. The reference helps you to choose the correct and semantic tags. The tags listed below correspond to the [ISO standard PDF 1.7](#). In 2018 the newer standard PDF 2.0 has been published. In this standard, some of the tags described here have been removed and new ones have been added. Since the standard is not yet widely used and a revision of the PDF/UA standard is still open, this overview still corresponds to PDF 1.7.'

Overview of the PDF Tags
→ <https://accessible-pdf.info/>

Reading Order

1 **My Best-Selling Novel**

2 By Iaam Meeh

3 **Chapter 1**

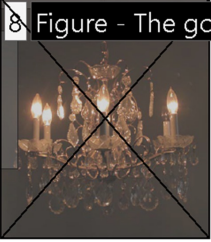
4 It was a dark and stormy night, when the chimes of the ancient grandfather clock that stood tall and elegant in the foyer, sounded for the last time. Bong. Bong. Bong. Bong. Bong. Bong. Bong. Bong. Bong. Bong. Bong. Bong. Bong. Bong.

5 **The stage is set**

6 It was midnight. A full moon poured its crystal photons through the grand palladium window at the top of the landing, bathing the gold-plated spiral staircase banister with its steely glow while the Austrian crystal chandelier dangling above on its gossamer-like brass chain tinkled gently in the soft breeze from the French doors that stood ajar below.

7 Suddenly, a shot rang out. A hound wailed with fear in the distant village. And a tiny crystal from the chandelier plummeted to the marble mosaic floor 20 feet below where it shattered into a million shards ... and drowned in a pool of fresh blood.

8 **Figure - The gold and crystal chandelier hanging in the foyer has a crystal missing from the right-most lamp. It's shattered on the floor below and covered in the hapless victim's blood. Who done it? And why? Mystery!**



9

10 **plot thickens**

11 Our mysterious journey continues as our hero enters the ghastly scene beneath the crystal chandelier. To be continued ...

Source: https://www.pubcom.com/blog/2020_08-18_ReadingOrder/ordertree_01.png

Creating an accessible PDF

If the PDF exists:

- Adobe Acrobat Pro
- www.pave-pdf.org

If the source file exists:

- Export the PDF with tags

PDF remediation with PAVE

The screenshot displays the PAVE website interface. At the top, the PAVE logo is shown, followed by navigation links for Introduction, FAQ, Background Information, and Contact. Language options for DE and EN are also present. The main heading is "Four steps to an accessible PDF". Below this, a diagram illustrates the process: 1. Upload a PDF document (represented by a document icon), 2. Automatic corrections (represented by gears), 3. Manual corrections (represented by a pencil), and 4. Download the accessible PDF (represented by a document icon with an ear). The text explains that PAVE makes PDFs accessible and interpretable by conventional readers, developed by the ICT Accessibility Lab at ZHAW. A "Start PAVE" button is provided, along with a link to a YouTube introduction video. A disclaimer states that ZHAW cannot guarantee the accuracy of the generated PDFs. A section titled "Why are accessible PDF documents important?" explains the benefits for visually impaired users. The footer includes copyright information for 2014-2023 ZHAW - ICT-Accessibility Lab and a link to legal information.

PAVE

Introduction FAQ Background Information Contact DE EN

Four steps to an accessible PDF

You can use PAVE to make your PDF documents accessible and to interpret conventional reader programs correctly. It does not change the visual layout of your PDF. The [ICT Accessibility Lab](#) of the ZHAW School of Engineering which developed PAVE, is making it available free of charge for personal use. Give it a try! If you want to make a large volume of PDFs accessible, please [contact us](#).

1. Upload your PDF document to PAVE.
Please note: The maximal allowed file size is **5 megabytes**.
2. PAVE will make the automatic corrections.
3. Simply make the remaining corrections yourself in PAVE.
4. Now you can download the accessible PDF document. The PDF document will remain on the PAVE server for a maximum of three weeks, unless you delete it manually beforehand.

[Start PAVE](#)

[Watch an introduction video on YouTube](#)

Disclaimer

While we make every effort to avoid changes to content or presentation of the document when edited in PAVE, we cannot provide a guarantee. The user is responsible for reviewing the PDF documents generated. ZHAW assumes no liability.

Why are accessible PDF documents important?

Persons with visual impairments can have electronic documents read aloud to them with a special software. However, this works well with PDF documents only if they are properly tagged – i.e. have the required metadata embedded. However, this is not frequently the case. There is now a simple solution to this widespread problem: Make your PDF documents accessible directly in PAVE.

Although the [ICT Accessibility Lab](#) has previously developed plug-ins to create accessible PDF documents from MS Word and MS PowerPoint, you cannot use these to make existing PDFs accessible. The gap is now closed thanks to PAVE.

© 2014 - 2023 ZHAW - ICT-Accessibility Lab | [Legal Information](#) [Top of Page](#)

zhaw School of Engineering

SBV Schweizerischer Blinden- und Sehbehindertenverband

PAVE has been developed by the [ICT-Accessibility Lab](#) of the [Zurich University of Applied Sciences \(ZHAW\)](#) in collaboration with the [Swiss Federation of the Blind and Visually Impaired](#).

PAVE Wins First Prize in International Competition

The [International Conference on Computers Helping People with Special Needs](#) honors software projects that help persons with disabilities. At this year's final in Paris, the [ICT Accessibility Lab](#) team impressed the international jury with PAVE, its groundbreaking open source software, and took first place.

[More about the award ceremony](#)

[SS12 Competition 2014](#)

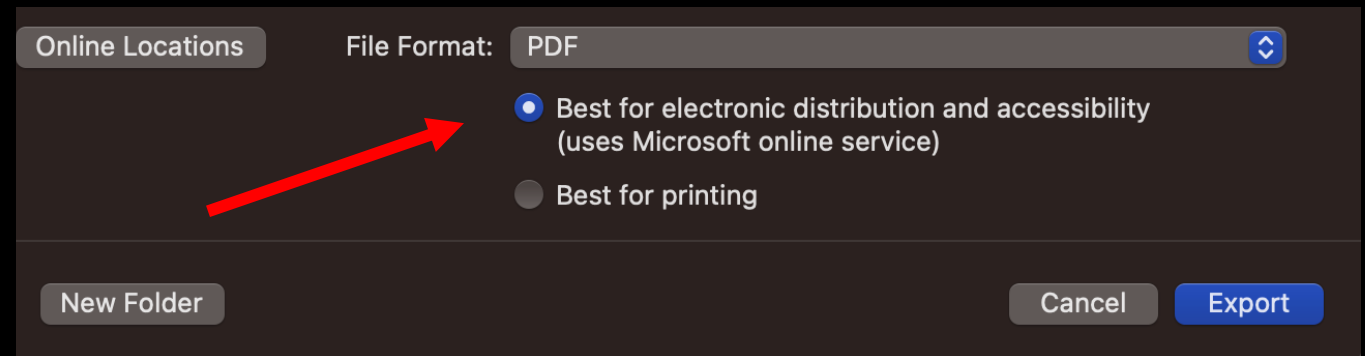
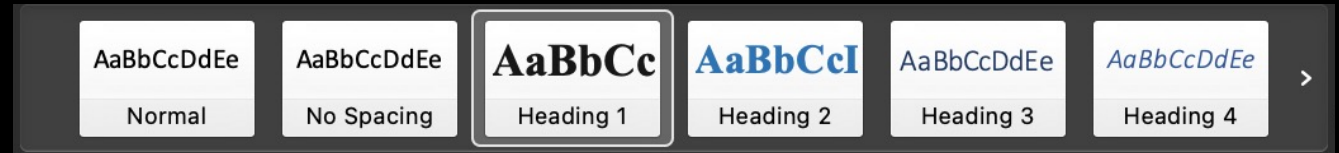
Contact

✉ Email: alireza.darvishy@zhaw.ch

➔ [Contact Information](#)

Export PDF with tags from Word

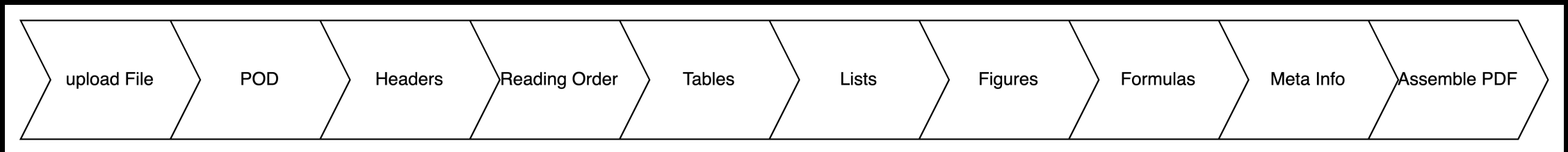
- Use the style templates
- Set alternative text
- Export the PDF with tags



How AI can simplify PDF remediation

Project: Accessible Scientific PDFs for All

- AI allows automating remediation steps for complex documents



Alternative text for mathematical formulae

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

Alternative Text

x equals start fraction minus b plus-or-minus square root of b squared minus 4 a c end root over 2 a end fraction

Agenda

Part II: The state of scientific PDF accessibility in Switzerland

- Survey of accessibility in Swiss repositories
- What does the future look like?

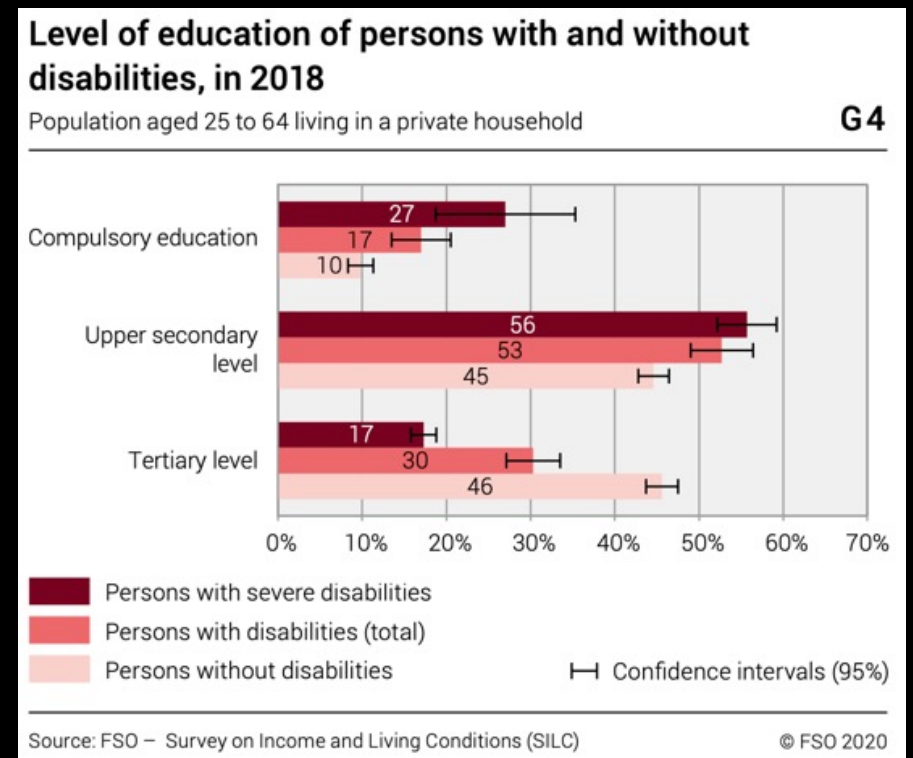
Motivation

People with disabilities in Switzerland

2018	As %	Number of persons
Persons with disabilities, severely limited	5.0%	347 000
Persons with disabilities, limited but not severely	17.2%	1 204 000
Total of persons with disabilities	22.2%	1 551 000

Source: FSO – Statistics on Income and Living Conditions (SILC) © FSO 2020

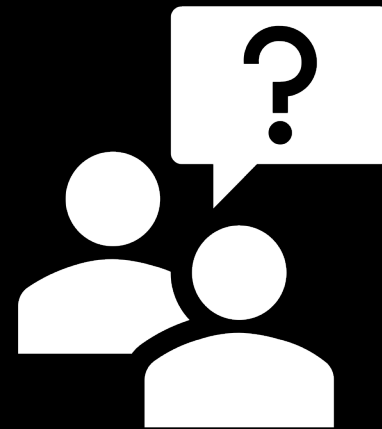
- Less people with disabilities pursue higher education than people without disabilities
- Around 30% of people with disabilities have a university level education



A mixed-method study on the state of scientific PDF accessibility in Swiss repositories



Quantitative analysis



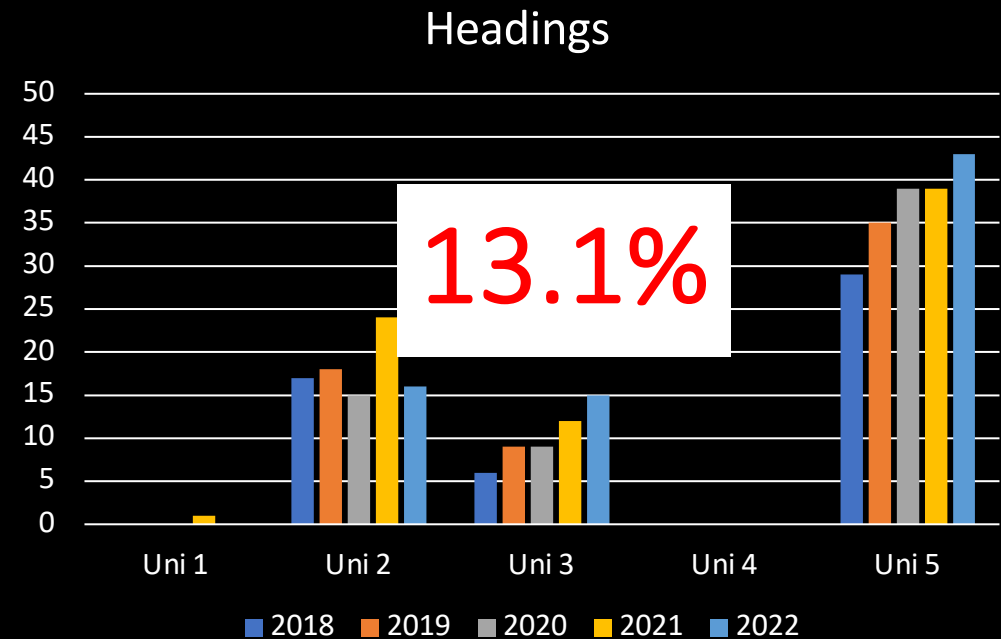
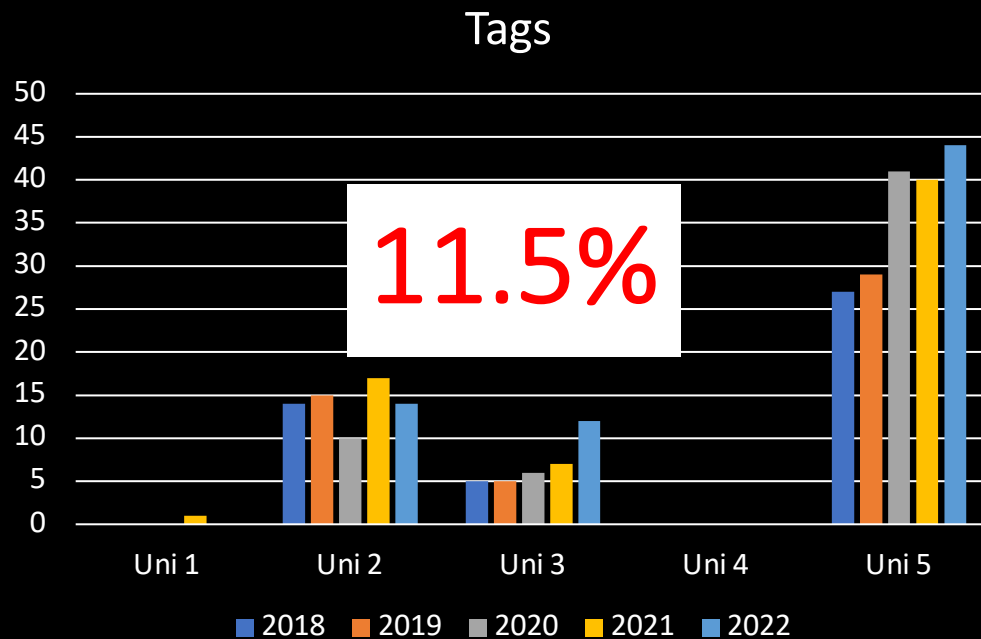
Qualitative analysis:
semi-structured interviews

Quantitative Analysis: Methods



- 5 repositories of German-speaking Swiss universities
- Papers from 2018 – 2022
- All Open Access papers were downloaded
 - Random sample of 100 papers per year per university
 - = in total 2500 papers were analyzed
- 2 minimal accessibility features → tags & headings
- Automatically checked & counted

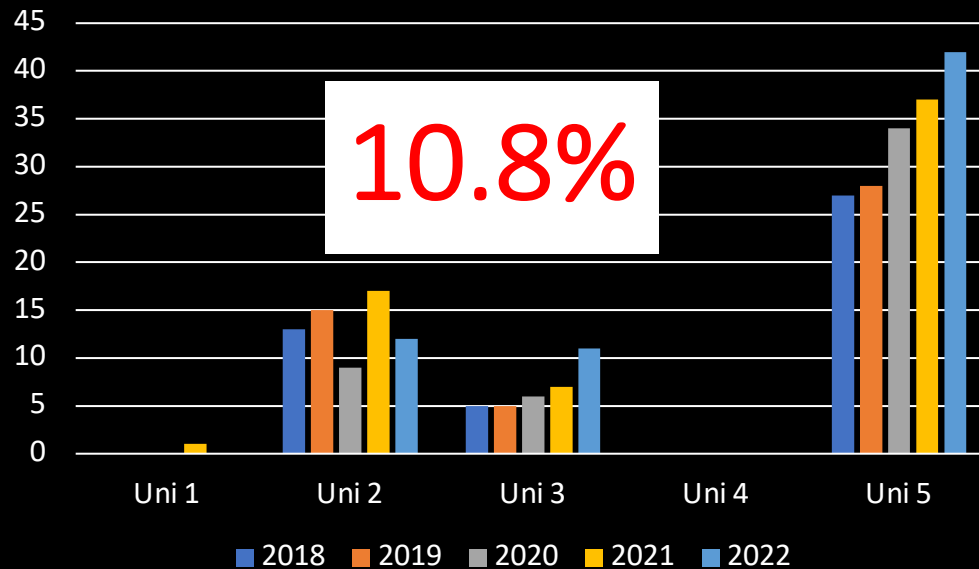
Quantitative Analysis: Results



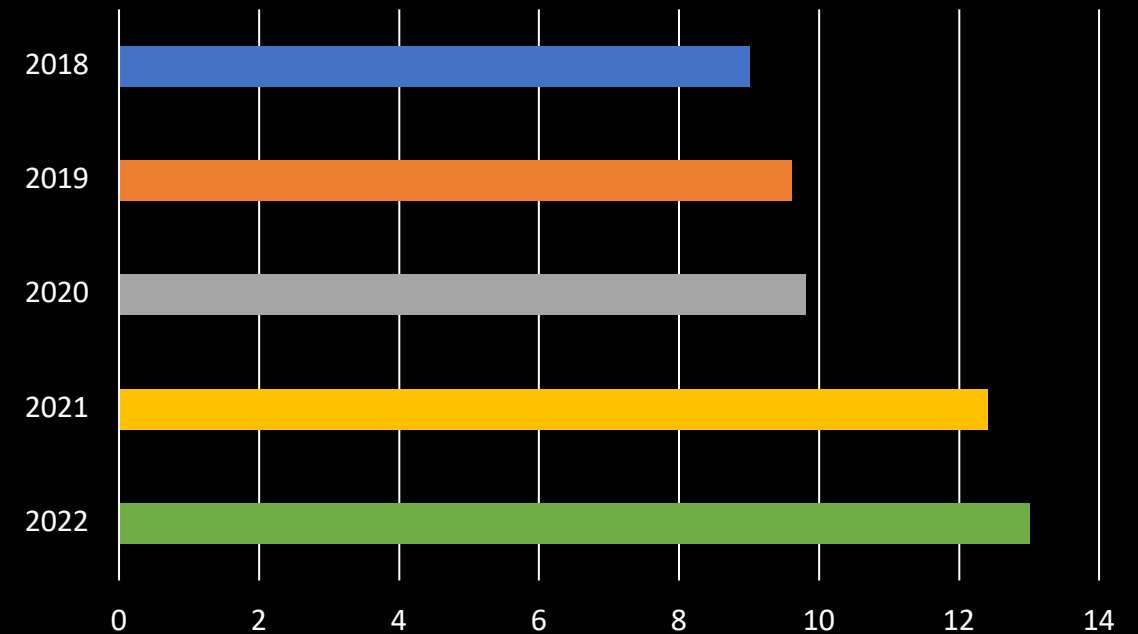
Quantitative Analysis: Results



Tags & Headings



Tags & Headings; total per year in %



Quantitative Analysis: Discussion



- No university reaches a satisfactory % of accessible papers
- Big differences in % between universities
- Somewhat of an upward trend



- What could the repositories do additionally?
- Is the problem that there is a lack of knowledge or know-how?
- With what priority should this issue be addressed?

Qualitative Analysis: Methods



- 9 Swiss universities
- Semi-structured interviews with the heads of the libraries / repositories, managers of accessibility projects, etc.
- Anonymous participation
- Set of 12 questions in 4 topics
 - Knowledge & Opinions
 - Priority & Sensitization
 - Accessibility Measures
 - Future Plans

Qualitative Analysis: Methods



1) Could you please briefly introduce yourself and your role in the repository?

Knowledge and Opinions

- 2) What do you know about the topic of accessible PDF documents?
- 3) In your opinion, what is an accessible PDF document?
- 4) How would you rate the accessibility of the publications in your repository?
 - a. Do you keep internal statistic on the accessibility of accessible ~~publications~~?

Priority and Awareness

- 5) With what priority do you address the issue of accessibility in your repository?
- 6) How are the employees at the repository sensitized on the topic of “accessibility”?
 - a. Are there trainings or conferences on the topic of “accessibility” for the employees?
 - b. Is attendance in such trainings or conferences mandatory?
- 7) Is there any exchange between the repositories of universities in Switzerland – through associations, conferences, etc.)?
If yes...
 - a. To what extent is “accessibility” addressed in the discussions?

Accessibility Measures

- 8) What are specific measures underway in your repository to keep it as accessible as possible (in terms of accessibility of publications)?
 - a. Do you have internal guidelines to ensure accessibility of published publications (mostly in PDF format)?
- 9) Do you check the accessibility of new incoming documents before publishing them on the repository?
If yes...
 - a. How do you check the accessibility of the documents (manually, etc.)?
 - b. Since when do you check the accessibility of the documents?
- 10) How do you handle it when, if for example, a visually impaired person requests an accessible version of a document?
 - a. Do you offer a service of providing alternative document formats?

Wishes and Conclusion

- 11) In your opinion, what fundamental change is needed at universities and research institutions for science in general to become more accessible?
- 12) What do you intend to do in the future to improve accessibility in your repository?

Qualitative Analysis: Results (1/2)



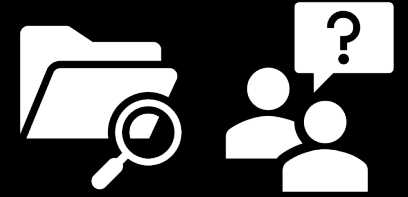
- Definition of "accessible PDF document": mainly focused on Open Access & structure
- Estimated amount of accessible documents = low
 - No statistic on the number of accessible documents
- Almost no repository provides trainings on accessibility for their employees
- Discussions within Open Access Working Group (Arbeitsgruppe)
 - Document accessibility is not a topic

Qualitative Analysis: Results (2/2)



- Accessibility has a low priority
 - Accessibility not seen as a responsibility
 - "Resources" as a big hindrance
- No measures to improve document accessibility
 - Accessibility measures usually concern the platform and UI
 - Almost no repository has internal guidelines
 - One exception with dissertations
- No check of the accessibility of incoming documents
- Reactive approach → ready to remediate documents on a case-by-case basis, but no established services
- Future plans don't include increasing accessibility of PDFs

Qualitative and quantitative Analysis: Combined Results



- Repositories of universities with accessibility guidelines had more PDFs with tags & headings
 - CAVE → concrete (!) guidelines necessary, i.e. all documents from within the university have to be accessible
- The more knowledge, the likelier it was for the repository to be willing to have more accessible PDFs
 - BUT willingness to be more accessible didn't translate into factual accessibility
- Certain common practices of some universities hinder creating accessible PDFs
 - e.g. using LaTeX

Recommendations



- Seek out knowledge & raise awareness on accessibility of PDF documents
 - Include people with disabilities in decisions and proposals
 - Participate in trainings and conferences
 - Collaborate with other universities to establish accessibility goals



- Make accessibility a requirement
 - Internal change: make accessibility a priority and a requirement (for authors, publishers, and repository managers)
 - Keep statistics on accessibility features of documents of own platform
 - Check accessibility of incoming documents before publication



- Use and establish guidelines
 - Existing ones: [ZHAW guidelines](#), [Swissuniversities](#), [Webaim](#), etc.
 - Build new guidelines within own institution



Don't only rely on laws to push you to make a change! → Willingness and open mindset are important

Questions? Comments? Ideas?