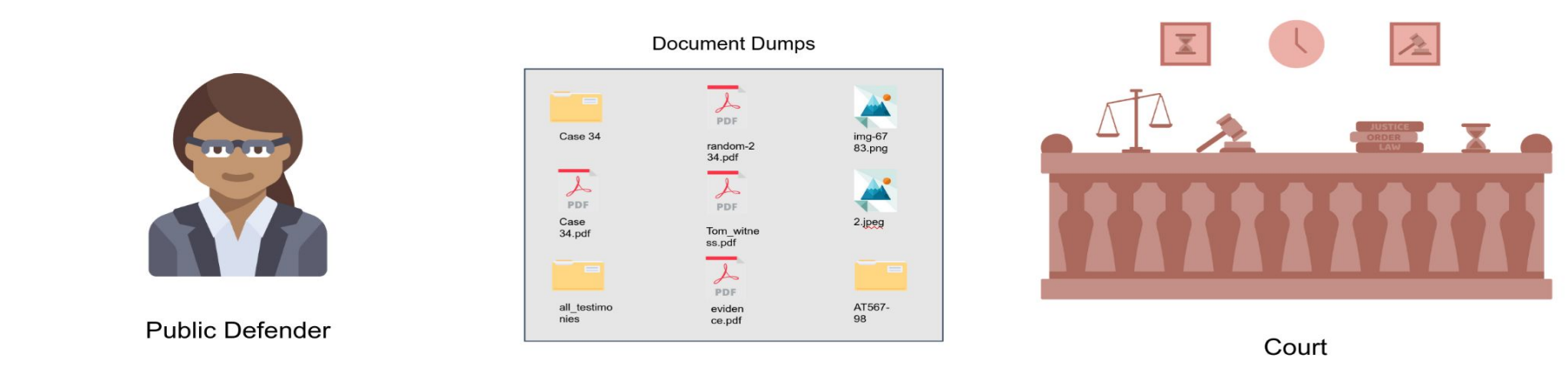


PRESENTER
Hellina Hailu Nigatu

Background

Domain experts like journalists and public defenders have to do a lot of manual data cleaning and processing when working with large document dumps. We collaborated with domain experts and built a document organization tool. We tested three programming paradigms to identify what works for our set of users.



Methods

Data Needs

Data Cleaning

- 1. Duplicates.
- 2. Quality of pages.

Data Extraction

- 1. Page Type
- 2. Agency Format

Data Organization

- 1. One case; several files.
- 2. Multiple cases might be found in a single PDF.

Programming Paradigm Study Results

- Text-based gives users low level control and exploration outside of designer provided abstraction.
- Visual gives designers the opportunity to provide information that users could not uncover on their own.
- PBE puts the focus on the data rather than the program structure.



Through collaborative design, we can make accessible programming tools for under-resourced domain experts.

*** Work to appear in FAccT23***

Design Goal	Constraints	Design Implications
<i>Human Control and Intervention</i>	<ul style="list-style-type: none">• Risks of mistakes in document classification are too high for this domain.• We found corrective actions to be more time and energy consuming than active, incremental decisions.	Design should prioritize supporting users exclusively in organizing the data rather than automating the whole process.
<i>Non-Interference with Existing Practices</i>	<ul style="list-style-type: none">• Cross-team differences in data management and handling practices.• Conflicting needs: Privacy and security concerns with using online platforms for one team conflicting with lack of local storage space for large data size in another team.• Pre-existing workflows for post-data organization tasks and file sharing.	Design should account for and be adoptable to pre-existing workflows and practices.
<i>Robustness to Data Variants</i>	<ul style="list-style-type: none">• Different teams with similar but not identical data.• Changes in data structure due to differences between LEAs themselves.	Potential solution would need to be resilient to changing formats and representations and inter-operable with similar but not identical datasets from others.
<i>High-level Abstractions</i>	<ul style="list-style-type: none">• Plain programming languages like Python or R require too much detailed technical knowledge to execute the required tasks.• Pre-built software solutions give limited flexibility to our users.	Tools should prioritize meeting users where they are with technical skills; requiring minimum training while allowing flexibility.
<i>Cost-Sensitive Solutions</i>	<ul style="list-style-type: none">• Resource-constrained teams lack the monetary resource to employ commercial software for their tasks.• Open-source software products do not produce same level of quality results.	When relying on open source software, tools should identify trade-offs with quality and ensure quality control with other schemes.

Lessons Learned: Co-designing

Cross-Team Transfer

Benefits of Co-Design

Long Term Engagement

Beyond Computational Tools

Programming Paradigms

