De-obfuscation and reusability of scientific publications

An Optical Character Recognition and heuristic approach

Luke Darlow

How and why things have changed

With research came... more questions!

- Found that bioinformatics supplementary data wasn't stored in an easily reusable way
 - PDFs Extracting data is a nightmare
 - Reusability and repetition are core to the scientific process

Goals

Build a proof of concept system for supplementary data extraction

- Finding the supplements: web scraping
- Extracting the data (Largest chunk of research) assuming tables and that a PDF page only contains this
 - Excel (easy) and PDFs (not so easy)
- Providing reusability
- Allow for user intervention
- Explore different techniques (OCR) and test viability
- Learn where things can change and improve

What is being done?

Current default techniques fail unless carefully infection AIDS and HIV infection chr2 >5902390 customized: AIDS and HIV Ъ, chr2 infection 56066880 Τ AIDS and HIV Nobody uses OCR chr2 infection Ψ AIDS and HIV chr2 infection or image processing AIDS and HIV



What I decided to try

- Used Scrapy to show it is possible to find certain document links
- Used xlrd to extract from excel spreadsheets
- Approached PDFs differently
 - Turned a page into an image
 - Used image processing and heuristics to find table dimensions
 - Used Tesseract OCR with approximate string matching to extract cell contents
- Built a simple user interface











Algorithmic nuances

- Row fixing algorithm
- Dark pixel counts
- OCR tweaks single characters
- Fuzzy string matching

Findings

- Using OCR isn't always accurate enough
 - The text exists in a readable form
 - Need to develop better technique
- Cell dimension finding needs more robustness – smoothing pixel counts could help
- Accurate automated information extraction is made difficult by the popular PDF
- Dynamic resolution of links is a challenge when scraping

What's next?

- Improving the table dimension finding
 - Possible use of AI algorithms
- Implementing a coordinate to element extraction instead of OCR
- Building a robust user interface
- Moving from proof of concept to development

