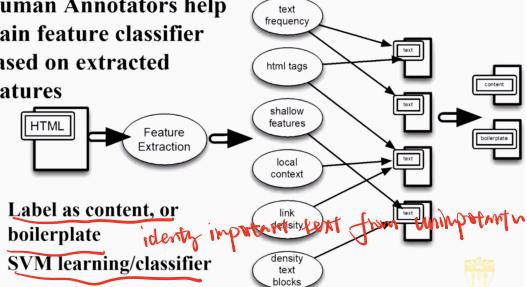


BoilerPipe Algorithm

- Human Annotators help train feature classifier based on extracted features



Readability – Browser Plugin

get, read, analyze HTML

- <https://web.archive.org/web/20110427053555/http://code.google.com/p/arc90labs-readability/>
- Examines HTML and relies on article text, title, images, and well formed HTML to extract article text, even if segmented across multiple pages
- Ports to Python and Ruby, though originally developed in Javascript
- Arc90 experiment – NYC based design and technology shop started in 2009
- Authors: Philip Forget and Chris Schomaker
- <http://www.readability.com/>

Other Key Insight: Applications

- You need this:



to read this:

choose parser.

Some Design Points

Quantitative

- Project Activity (num active committers; num commits)
- Computational Performance *3x*
- Memory Performance
- Reliability

Qualitative

- Support for corner cases
- Support for advanced text extraction
- Metadata Support
- License for the library
 - Permissive versus Reciprocal

Tika General Approach

Curated mapping of IANA MIME types to parsers

- Parser libraries that support MIME types are added to Tika, either in pure library form, or by writing small parser code

Parser selection strategy

- First parser selected that declares support for MIME type *for composite filetypes*
- FallBack Parser, CompositeParser, ExternalParser
- Other methodologies are being explored

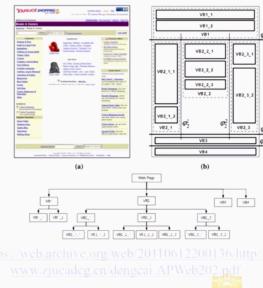
polite

Visual Page Segmentation

- Cai et al, APWeb 2003

Approach

- Extract DOM blocks and assign them a value based on visual coherence
- Build hierarchical model of the site based on this
- Vertical/Horizontal separators
- Identify content blocks



<https://web.archive.org/web/20110612200136/http://www.zjiaide.en.dengtan/APWeb2003.pdf>

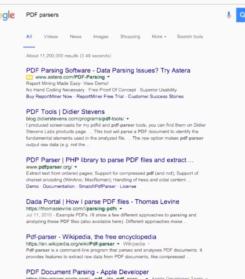
Some Key Insights and Observations

HTML is great way to structure documents. It's about each part mean.

- Several tried and true approaches work when you can deal with HTML oriented documents
 - Is there some way to transform documents into HTML, especially well formed HTML?
- Lots of exploitation of text oriented features, and text summarization techniques
 - Metadata features not really used, since focus in on content – we'll get to metadata later
- Many libraries exist to perform extraction
 - Some are more powerful than others, more reliable, etc.

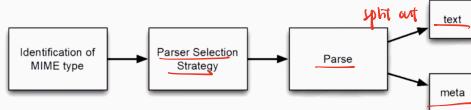
Finding Parsing Libraries is Difficult

- Google search for PDF parsers
- 50+ parsers on the first few pages
- Which one is the right one?
- Note: NONE of these are the ones that Tika uses (PDFBox)



Apache Tika – Content Extraction Process

- Looks something like this



- Incorporates Parser Selection based on MIME hierarchy

- Requires mapping of Parsers to MIME types that they support

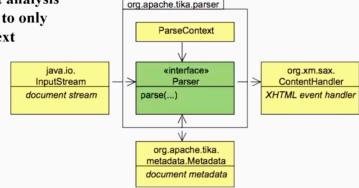
Design Goal: Convention over Configuration

AutoDetectParser algorithm

- Load all parsers
 - For any parser that fails to initialize or where JAR file and/or class isn't available, ignore or warn and move on
- Load all mappings to MIME types
- Incoming file, stream, URL
 - Detect MIME Type with Accuracy
 - Find first parser that satisfies MIME type, and invoke it
 - Parser could be Façade to
 - Fallback Parser (try one; fall back to other); Composite Parser (call many parsers in some sequence); External Parser (external program)

XHTML as an intermediate output format

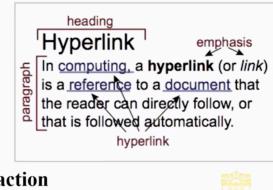
- Wealth of research on how to process text from HTML documents
- XHTML => structured HTML, well formed unlike most of Internet architecture
- Can leverage streaming SAX (Simple API for XHTML Processing)
- Enables wide text analysis and isolates Tika to only extracting base text



Another benefit

- XHTML allows annotations to be captured surrounding text

- Annotations are key input e.g., to Tag Frequency algorithms, to Feature Detection Algorithms and to Text Analytics
- Also extremely useful from a metadata perspective as annotations can enhance metadata extraction



Some Current Thrusts and Research Areas

Optical Character Recognition

- Tesseract
 - <https://code.google.com/p/tesseract-ocr/>
 - Great and Accurate Toolkit, Apache License, version 2 ("ALv2")
 - Many recent improvements by Google and Support for Multiple Languages
- Integrated Tika and Tesseract
 - <http://issues.apache.org/jira/browse/TIKA-93>
 - Thank you to Grant Ingersoll (original patch) and Tyler Palsulich for taking the work the rest of the way to get it contributed

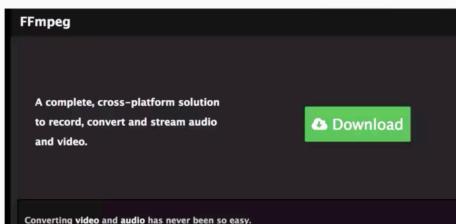


Scientific Literature

- GeneRation Of Bibliographic Data (GROBID)
 - P. Lopez et al, Research and Advanced Technology for Digital Libraries, 2009
- Uses Conditional Random Fields and trained on 1000 examples for identifying header information (authors, institutions, publication year, affiliations) and bibliographic information (citations) and measurements from scientific literature
 - GROBID code on Github:
 - <http://grobid.readthedocs.org/en/latest/Introduction/>

Image / Video / Multimedia

- Extraction of content using traditional image and video analysis tools
 - FFMPEG, EXIF Tool



Tika + Tesseract

- <https://cwiki.apache.org/confluence/display/tika/TikaOCR>
- brew install tesseract --all-languages
- tika -t /path/to/tiff/file.tif
 - Tika will automatically discern whether you have Tesseract installed or not
- Try it from the Tika REST server!
 - In another window, start Tika server
 - java -jar /path/to/tika-server-1.11.jar
 - In another window, issue a cURL request
 - curl -T /path/to/tiff/image.tif http://localhost:9998/tika --header "Content-type: image/tiff"

GROBID is cool

- Header extraction and parsing from article in PDF format. The extraction here covers the usual bibliographical information (e.g. title, abstract, authors, affiliations, keywords, etc.).
- References extraction and parsing from articles in PDF format, around .87 F1-score against on an independent PubMed Central set of 1943 PDF containing 90,125 references, and around .89 on a similar bioRxiv set. All the usual publication metadata are covered (including DOI, PMID, etc.).
- Citation contexts recognition and resolution of the full bibliographical references of the article. The accuracy of citation contexts resolution is above .78 f-score (which corresponds to both the correct identification of the citation callout and its correct association with a full bibliographical reference).
- Parsing of references in isolation (above .90 F1-score at instance-level, .95 F1-score at field level).
- Parsing of names (e.g. person title, forenames, middlename, etc.), in particular author names in header, and author names in references (two distinct models).
- Parsing of affiliation and address blocks.
- Parsing of dates, ISO normalized day, month, year.
- Full text extraction and structuring from PDF articles, including a model for the overall document segmentation and models for the structuring of the text body (paragraph, section titles, reference callout, figure, table, etc.).
- Consolidation/resolution of the extracted bibliographical references using the biblio-glutton service or the CrossRef REST API. In both cases, DOI resolution performance is higher than 0.95 F1-score from PDF extraction.
- Extraction and parsing of patent and non-patent references in patent publications.
- PDF coordinates for extracted information, allowing to create "augmented" interactive PDF.

Scientific Analysis

- GeoTopicParsing
 - Text Analysis with Apache OpenNLP to identify Location entities in text
 - Lucene Geo Gazetteer
 - Entities resolved to Geonames and to Latitudes and Longitudes
 - <https://github.com/christmann/lucene-geo-gazetteer>
 - Result is: “reading” the text to spot locations and to return their coordinates
 - Started as a CSCI 572 Directed Research Project!
- GDAL (Geospatial Data Abstraction Library)

Tika Research in Advanced Extraction

- <https://cwiki.apache.org/confluence/display/tika/Home-AdvancedContentExtractionwithTika-Integration>

Advanced Content Extraction with Tika - Integration

- Getting Tika up and Running with Pooled Time Series - How to use Tika with the Pooled Time Series video descriptor similarity code.
- Getting Tika up and Running with Apache cTAKES - How to use Tika with Apache cTAKES the clinical text biomedical knowledge extraction framework.
- Getting Tika up and Running with Apache NLP API - How to use Tika with EXTRAFONT.
- Getting Tika up and Running with FFmpeg - How to use Tika with FFmpeg.
- Getting Tika up and Running with the GROBID PDF Journal parser - How to use Tika with the GROBID PDF journal parser.
- Getting Tika up and Running with the GeoTopicParser based on Geonames.org, Lucene, and OpenNLP.
- Getting Tika up and Running with OCR - How to use Tika with OCR from Tesseract.
- Getting Tika up and Running with the Geospatial Data Abstraction Library (GDAL) - How to use Tika with GDAL to parse/extract geospatial data files.
- Getting Tika up and running with Stanford Core NLP and with OpenNLP - How to use Tika with Stanford NERNLP and with Apache Open NLP.

Summary

- How to extract text from any document?
- Basic Approach
 - Identification of Text
 - Analysis of Text
 - Featurization
- How is it done in Tika?
 - Architecture
 - Implementation and Tradeoffs
- Current directions in extraction of text

API Approaches for Text Extraction

Alchemy

- <http://www.alchemyapi.com/api/text-extraction>
 - HTML to XML, JSON and RDF

Textise

- Any page or URL to text
- <http://www.textise.net/>

Givemetext.org Open Knowledge Foundation

- <http://givemetext.okfnlabs.org/>
- Based on Apache Tika