

NEETS Domain Modeling v 2

Previously we attempted this using an IKVM version of TIKA, but were unable to extract images properly. This new version uses Beaker 1.7.1 which is able to load TIKA classes properly, which Beaker 1.6 was unable to do. However 1.7.1 has problems with Groovy and Scala coexisting (maybe Java too). So while the Java below for using TIKA must run in 1.7.1 for class loading purposes, it may be necessary to run the rest in 1.6 until the Beaker bug is fixed.

- Using TIKA, convert PDF to HTML and extract images

Creates a folder for each pdf with extracted images and html. The images are a mixture of pngs and tifs. Tifs will not render in html, so we must follow this with a tif to png conversion

Use TIKA to convert PDF to HTML and extract images

Jv Java

```
1 /**
2 * Licensed under the Apache License, Version 2.0 (the "License");
3 * you may not use this file except in compliance with the License.
4 * You may obtain a copy of the License at
5 *
6 * http://www.apache.org/licenses/LICENSE-2.0
7 *
8 * Unless required by applicable law or agreed to in writing, software
9 * distributed under the License is distributed on an "AS IS" BASIS,
10 * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
11 * See the License for the specific language governing permissions and
12 * limitations under the License.
13 */
14
15 package org.apache.tika.example;
16
17
18 import java.io.IOException;
```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```
22 import java.io.FileWriter;
23 import java.io.PrintWriter;
24 import java.io.File;
25 import java.io.FileInputStream;
26 import java.io.FileOutputStream;
27 import java.io.OutputStreamWriter;
28 import java.io.BufferedWriter;
29 import java.nio.charset.Charset;
30
31 import org.apache.tika.config.TikaConfig;
32 import org.apache.tika.detect.Detector;
33 import org.apache.tika.exception.TikaException;
34 import org.apache.tika.extractor.EmbeddedDocumentExtractor;
35 import org.apache.tika.extractor.ParsingEmbeddedDocumentExtractor;
36 import org.apache.tika.io.FilenameUtils;
37 import org.apache.tika.metadata.Metadata;
38 import org.apache.tika.mime.MediaType;
39 import org.apache.tika.mime.MimeTypeException;
40 import org.apache.tika.parser.AutoDetectParser;
41 import org.apache.tika.parser.ParseContext;
42 import org.apache.tika.parser.Parser;
43 import org.apache.tika.sax.BodyContentHandler;
44 import org.xml.sax.ContentHandler;
45 import org.xml.sax.SAXException;
46
47 import org.apache.tika.sax.ToXMLContentHandler;
48 import org.apache.tika.parser.pdf.PDFParserConfig;
49
50 public class ExtractEmbeddedFiles {
51
52     private Parser parser = new AutoDetectParser();
53     private Detector detector = ((AutoDetectParser)parser).getDetector();
54     private TikaConfig config = TikaConfig.getDefaultConfig();
55
56     //public void extract(InputStream is, Path outputDir) throws
57     //SAXException, TikaException, IOException {
58     public void extract(String inputPath) throws SAXException,
59     TikaException, IOException {
60         /*
61         Metadata m = new Metadata();
62         ParseContext c = new ParseContext();
63         ContentHandler h = new BodyContentHandler(-1);
64
65         c.set(Parser.class, parser);
66         EmbeddedDocumentExtractor ex = new
67         MyEmbeddedDocumentExtractor(outputDir, c);
68         c.set(EmbeddedDocumentExtractor.class, ex);
69         */
70         File inputFile = new File(inputPath);
```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```
12         InputStream inputStream = new FileInputStream(inputFile);
73         File outputDirectory = new File(parentDirectory, inputFile.getName() +
    "-extracted");
74         System.out.println(outputDirectory.getPath()); //junk
75
76         Parser parser = new AutoDetectParser();
77         ToXMLContentHandler handler = new
    org.apache.tika.sax.ToXMLContentHandler();
78
79         PDFParserConfig pdfConfig = new PDFParserConfig();
80         pdfConfig.setExtractInlineImages(true);
81
82         ParseContext parseContext = new ParseContext();
83
84         parseContext.set(PDFParserConfig.class, pdfConfig);
85         parseContext.set(Parser.class, parser);
86
87         EmbeddedDocumentExtractor ex = new
    MyEmbeddedDocumentExtractor(outputDirectory.toPath(), parseContext);
88         parseContext.set(EmbeddedDocumentExtractor.class, ex);
89
90         Metadata metadata = new Metadata();
91
92         parser.parse(inputStream, handler, metadata, parseContext);
93
94         String text = handler.toString().trim();
95
96         File outputFile = new File(outputDirectory, inputFile.getName() +
    ".xhtml");
97         System.out.println(outputFile.getPath()); //junk
98         //PrintWriter printer = new PrintWriter( new FileWriter(
    outputFile.getPath() ) );
99
100        PrintWriter printer = new PrintWriter( new BufferedWriter (new
    OutputStreamWriter(
101            new FileOutputStream( outputFile.getPath() ), 
102            Charset.forName("UTF-8").newEncoder()
103          ));
104        printer.print( text );
105        printer.close();
106
107 }
108
109 private class MyEmbeddedDocumentExtractor extends
    ParsingEmbeddedDocumentExtractor {
110        private final Path outputDir;
111        private int fileCount = 0;
112
113        private MyEmbeddedDocumentExtractor(Path outputDir, ParseContext
    context) {
114            super(context);
```

```
118     @Override
119     public boolean shouldParseEmbedded(Metadata metadata) {
120         return true;
121     }
122
123     @Override
124     public void parseEmbedded(InputStream stream, ContentHandler handler,
125                               Metadata metadata, boolean outputHtml)
126                               throws SAXException, IOException {
127
128         //try to get the name of the embedded file from the metadata
129         String name = metadata.get(Metadata.RESOURCE_NAME_KEY);
130
131         if (name == null) {
132             name = "file_" + fileCount++;
133         } else {
134             //make sure to select only the file name (not any directory paths
135             //that might be included in the name) and make sure
136             //to normalize the name
137             name = FilenameUtils.normalize(FilenameUtils.getName(name));
138         }
139
140         //now try to figure out the right extension for the embedded file
141         MediaType contentType = detector.detect(stream, metadata);
142
143         if (name.indexOf('.')==-1 && contentType!=null) {
144             try {
145                 name += config.getMimeRepository().forName(
146                     contentType.toString()).getExtension();
147             } catch (MimeTypeException e) {
148                 e.printStackTrace();
149             }
150             //should add check to make sure that you aren't overwriting a file
151             Path outputFile = outputDir.resolve(name);
152
153             //get parent, convert to file, make directory
154             //outputFile.getParent().toFile().mkdir();
155             //do a better job than this of checking
156             Files.createDirectories(outputFile.getParent());
157             Files.copy(stream, outputFile);
158
159
160     }
161 }
162 }
```

Run

Jv Java ►

```
1 import org.apache.tika.example.ExtractEmbeddedFiles;
2 import java.io.File;
3 import java.nio.file.Files;
4 import java.nio.file.Path;
5 import java.io.InputStream;
6 import java.io.FileInputStream;
7
8 File pdfDirectory = new
9     File("/z/aolney/research_projects/braintrust/materials/NEETS/pdf/");
10 File[] pdfFiles = pdfDirectory.listFiles((d, name) ->
11     name.endsWith(".pdf"));
12 for ( File pdfFile : pdfFiles ) {
13     ExtractEmbeddedFiles ex = new ExtractEmbeddedFiles();
14     ex.extract( pdfFile.getPath() );
15 }
```

Error

2 lines of stderr, BeakerDisplay

Use ImageMagick to convert TIF to PNG ►

Jv Java ►

```
1 import java.io.File;
2 import java.nio.file.Files;
3 import java.nio.file.Path;
4 import java.io.InputStream;
5 import java.io.FileInputStream;
6 import java.lang.ProcessBuilder;
7 import java.io.IOException;
8 import java.io.PrintWriter;
9 import java.io.OutputStreamWriter;
10 import java.io.Writer;
11 import java.io.FileOutputStream;
12
```

Beaker

*/z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

```
File View Notebook Help edited
10     File tempScript = createTempScript(commands);
11
12
13
14
15
16
17
18     try {
19         ProcessBuilder pb = new ProcessBuilder("bash",
19 tempScript.toString());
20         pb.inheritIO();
21         Process process = pb.start();
22         process.waitFor();
23     } catch (java.lang.InterruptedException e ) {
24         System.err.println( e );
25     } finally {
26         tempScript.delete();
27     }
28 }
29
30 private static File createTempScript(String commands ) throws IOException {
31     File tempScript = File.createTempFile("script", null);
32
33     Writer streamWriter = new OutputStreamWriter(new FileOutputStream(
34             tempScript));
35     PrintWriter printWriter = new PrintWriter(streamWriter);
36
37     printWriter.println("#!/bin/bash");
38     printWriter.println(commands);
39
40     printWriter.close();
41
42     return tempScript;
43 }
44 }
45
```

Run

com.twosigma.beaker.javash.bkrfb7a9607.Bash

Jv Java

```
1 import java.io.File;
2 import java.nio.file.Files;
3 import java.nio.file.Path;
4 import java.io.InputStream;
5 import java.io.FileInputStream;
6
```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```

9 for ( File directory : htmlDirectories ) {
10    //do tif to png conversion here, assumes imagemagick is installed and
11    //we have bash shell
12    Bash.executeCommands("cd '" + directory.getPath() + "' ;for f in
13    *.tif; do echo \"Converting $f\"; convert \"$f\" \"$(basename \"$f\""
14    ".tif).png\"; done");
15 }
```

- Apply domain model, regenerate HTML ai ▶ JSON

Domain model ▶

This was developed in monodevelop and pasted below. Code below runs fine, but the complexity of the code is such that it is advisable develop in the orginal solution and treat this as archival.

Two modifications must be made

- Comment out EntryPoint
- Add main [] as last line so main is called

CSS for the output HTML was hand authored.

F# FSharp ▶

```

1 type PageType =
2   | TOCPage
3   | TitlePage
4   | PrefacePage
5   | MainPage
6   | AppendixPage
7   | IndexPage
8   | AssignmentPage
9   with override this.ToString() = sprintf "%A" this
10
11 type Header =
12   | Chapter of string * int
13   | LearningObjectives of string
14   | Section of string
```

Beaker

*/z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

```
File View Notebook Help edited
18    with override this.ToString() = sprintf "%A" this
19
20 type PageElement =
21   | ImageCaption of string
22   | ImageURL of string
23   | Header of Header
24   | Question of string
25   | Answer of string
26   | LearningObjective of string
27   | Paragraph of string
28   with override this.ToString() = sprintf "%A" this
29
30 //Helpers to classify page type
31 type Page =
32   {
33     Number : int option; //preface pages are negative
34     NumberString : string;
35     Type : PageType
36     Elements : ResizeArray<PageElement>
37   }
38   with override this.ToString() =
39     let pageNumber =
40       match this.Number with
41       | Some( x ) -> x.ToString()
42       | None -> "no-number"
43     pageNumber + " " + this.Type.ToString() + " " + System.String.Join(
44       " ", this.Elements)
```

Run



F# FSharp



```
1 #r "/z/aolney/repos/HtmlAgilityPack.1.4.9.5/lib/Net40/HtmlAgilityPack.dll"
2 #r "/z/aolney/repos/Newtonsoft.Json.9.0.1/lib/net40/Newtonsoft.Json.dll"
3
4 open HtmlAgilityPack
5
6 let (|IsTOC|_|) (elements:ResizeArray<PageElement>) =
7   elements |> Seq.tryFind(
8     function
9       | Header (Section hs) when hs.Contains("TABLE OF CONTENTS") ->
10         true
```

```
13
14 let (|IsAppendix|_|) (elements:ResizeArray<PageElement>) =
15     elements |> Seq.tryFind(
16         function
17             | Header (Section hs) when hs.Contains("APPENDIX") -> true
18             | _ -> false
19         )
20
21 let (|IsIndex|_|) (elements:ResizeArray<PageElement>) =
22     elements |> Seq.tryFind(
23         function
24             | Header (Section hs) when hs.Contains("INDEX") -> true
25             | _ -> false
26         )
27
28 let (|IsAssignment|_|) (elements:ResizeArray<PageElement>) =
29     elements |> Seq.tryFind(
30         function
31             | Header (Section hs) when hs.StartsWith("ASSIGNMENT") ||
hs.StartsWith("Assignment") -> true
32             | _ -> false
33         )
34
35 //helpers to classify page elements
36 let (|IsImageCaption|_|) (text:string,html:string, _) =
37     if text.StartsWith("Figure") && text.Contains("-") then
38         Some(text)
39     else
40         None
41
42 let (|IsImageUrl|_|) (text:string,html:string,_) =
43     if html.Contains("img src=") then
44         Some( html )
45     else
46         None
47
48 //It appears that headers never have clause final punctuation
49 let (|IsHeader|_|) (rawText:string,html:string,_) =
50     let text = rawText.Trim()
51     if text = "" || text.Contains(".") || text.Contains("?") ||
text.Contains("!") || text.EndsWith(":") then
52         None
53     else
54         Some(text)
55
56 let questionRegex = new System.Text.RegularExpressions.Regex("^\\s*Q\\d+")
57 //ignores that multiple questions can be in one paragraph
58 let (|IsQuestion|_|) (text:string,html:string,_) =
59     if questionRegex.IsMatch(text) then
60         Some(text)
61     else
```

File View Notebook Help

edited

```
55 let (|IsAnswer|_|) (text:string, num:string, _) =
56     if answerRegex.IsMatch( text ) then
57         Some( text )
58     else
59         None
60
61 let learningObjectiveRegex = new
62     System.Text.RegularExpressions.Regex("^\\s*\\d+")
63 let (|IsLearningObjective|_|) (text:string, html:string, lastHeader:Header)
64 =
65     let trimmed = text.Trim()
66     let patternMatch = learningObjectiveRegex.IsMatch( trimmed ) &&
67     trimmed.EndsWith(".")
68     match lastHeader, patternMatch with
69     | Header.LearningObjectives(text), true -> Some(text)
70     | _ -> None
71
72
73 let capitalLetterRegex = new System.Text.RegularExpressions.Regex("^[A-
74 Z\\s-]+$")
75 let pageNumberRegex = new System.Text.RegularExpressions.Regex("^\s*
76 [ivxIVXANDEX0-9-]+\s*$") //may need adjusting for appendix
77
78 //assumes we have already tested this is in fact a header
79 let HeaderClassifier (text:string) (html:string) =
80     if text.StartsWith("CHAPTER") then
81         let s = text.Split(' ')
82         match System.Int32.TryParse( s.[1] ) with
83         | true, number -> Header.Chapter (text, number)
84         | false, _ -> Section text
85     else if text.StartsWith("LEARNING OBJECTIVES") then
86         Header.LearningObjectives text
87     else if text.StartsWith("SUMMARY") then
88         Header.Summary text
89     else if pageNumberRegex.IsMatch(text) then
90         PageNumber text
91     else if capitalLetterRegex.IsMatch(text) then
92         Section text
93     else
94         Subsection text
95
96 let GetElementByChapterList( pageList : ResizeArray<Page> ) (
97     elementMatcher : PageElement -> string option ) =
98     //for each chapter, build a list of following sections
99     let mutable sectionList = new ResizeArray<string*int>()
100    let mutable chapterList = new
101        ResizeArray<string*int*ResizeArray<string*int>_()
102        let chapterHash = new System.Collections.Generic.HashSet<string>()
103        //specifically to solve chapters appearing in appendix
104        let mutable currentChapter = (0, "")
```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```

111         //then see if our page number is still increasing
112         match pageList.[index].Number,pageList.[index+1].Number with
113             | Some(current),Some(next) when next > current -> true
114             | _ -> false
115         else
116             false
117
118     let mutable i = 0
119     while( i < pageList.Count - 1 ) do //ignore last page
120         let chapterText = pageList.[i].Elements |> Seq.tryPick( function |
121             Header ( Chapter (text,_) ) -> Some( text ) | _ -> None)
122
123             if chapterText.IsSome && chapterHash.Contains(
124                 chapterText.Value.Trim() ) |> not then //b/c closures and mutation
125                 currentChapter <- (i,chapterText.Value)
126
127             while ShouldContinue i pageList do
128                 //this only works because on a page, sections never come
129                 before chapters
130
131                 let sections = pageList.[i].Elements |> Seq.choose(
132                     elementMatcher )
133                     for sectionText in sections do sectionList.Add(
134                         sectionText, i )
135                         i <- i + 1
136
137                 chapterList.Add( currentChapter |> snd, currentChapter |> fst,
138                 sectionList )
139
140                 chapterHash.Add( (currentChapter |> snd).Trim() ) |> ignore
141
142                 sectionList <- new ResizeArray<string*int>()
143                 i <- i + 1
144
145             else
146                 i <- i + 1
147
148         //return
149         chapterList
150
151     let whitespaceRegex = new System.Text.RegularExpressions.Regex("\s+")
152
153 //render an entire file to html
154 let RenderToHTML (fileName : string) (pageList : ResizeArray<Page> ) =
155
156     //helpers for internal link ids
157     let SectionLinkId (text:string) (page:Page) =
158         whitespaceRegex.Replace("section-" + text + " " +
159         page.NumberString, "-")
160
161     let PageLinkId (page:Page)=
162         whitespaceRegex.Replace("page-" + page.Type.ToString() + " " +
163         page.NumberString,"-")
164
165     //helpers for images
166     let imageNamePrefix = fileName.Replace(".xhtml","")

```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

<p>File View Notebook Help</p>	edited
---	--------

```

155     | image0.png | image1.png -> image-logo
156     | _ -> "image-regular"
157
158     //html for an element
159     let GetElementHTML (page:Page) (element:PageElement) =
160         match element with
161         | ImageCaption text -> "<p class=\"imagecaption\">" + text + "
162             </p>" +
163             | ImageURL html ->
164                 let regexMatch = imageUrlRegex.Match( html )
165                 //we've already converted tifs to pngs at this stage
166                 let imageFile = regexMatch.Groups.
167                     ["imageFile"].Value.Replace(".tif", ".png")
168                     "<img class=' " + (ImageClass imageFile) + "' src=\"" + imageFile
169                     + "\" alt=\"" + imageFile + "\" />" +
170                     | Header ( Chapter (text,_) ) -> "<h1 class=\"chapter\" id=\""
171                         (SectionLinkId text page) + "\">" + text + "</h1>" +
172                     | Header ( LearningObjectives text ) -> "<h2
173                         class=\"learningobjectivesection\">" + text + "</h2>" +
174                     | Header ( Section text ) -> "<h2 class=\"section\" id=\""
175                         (SectionLinkId text page) + "\">" + text + "</h2>" +
176                     | Header ( Subsection text ) -> "<h3 class=\"subsection\">" + text
177                     + "</h3>" +
178                     | Header ( PageNumber text ) -> "<p class=\"pagenumber\">" + text
179                     + "</p>" +
180                     | Header ( Summary text ) -> "<h1 class=\"summary\">" + text + "
181                         </h1>" +
182                         | Question text -> "<p class=\"question\">" + text + "</p>" +
183                         | Answer text -> "<p class=\"answer\">" + text + "</p>" +
184                         | LearningObjective text -> "<p class=\"learningobjective\">" +
185                         text + "</p>" +
186                         | Paragraph text -> "<p class=\"paragraph\">" + text + "</p>" +
187
188                         //page helper function
189                         let RenderPageElementsToHTML fileName ( page : Page ) =
190
191                             let sb = System.Text.StringBuilder(page.Elements.Count)
192
193                             //split out image elements from rest
194                             let imageList,otherList = page.Elements |> Seq.toList |>
195                             List.partition( function | ImageCaption _ -> true | ImageURL _ -> true | _
196                               -> false )
197                             let urlList,captionList = imageList |> List.partition(function |
198                               ImageURL _ -> true | _ -> false )
199                             let urlElements = new ResizeArray<PageElement>( urlList )
200                             let captionElements = new ResizeArray<PageElement>( captionList )
201
202                             //create image div for all image elements on this page, it goes at
203                             top
204                             sb.Append("<div class=\"imageblock\">"") |> ignore
205                             //for each image/caption, create a div

```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```
195      //add a caption if it exists
196      if i < captionElements.Count then
197          sb.Append( GetElementHTML page captionElements.[i] ) |>
198      ignore
199          sb.Append("</div>") |> ignore
200          sb.Append("</div>") |> ignore
201
202          //add remaining page elements
203          //for element in page.Elements do
204          for element in otherList do
205              sb.Append( GetElementHTML page element ) |> ignore
206          //return everything in the string builder
207          sb.ToString()
208
209          //body helper function
210          let RenderHTMLBody fileName (pageList : ResizeArray<Page>) =
211              let GetPageNav ( index : int ) (pageList : ResizeArray<Page>) =
212                  let previousPageId = if index > 0 then (index-1).ToString()
213                  else "" //PageLinkId pageList.[index-1] else ""
214                  let nextPageId = if index < pageList.Count - 1 then
215                      (index+1).ToString() else "" //PageLinkId pageList.[index+1] else ""
216                  """<div class="centered">
217                  <div class="pagination">
218                      <a href='####' + previousPageId + "##'></a>
219                      <a href='####' + nextPageId + "##'></a>
220                  </div></div>"""
221
222          let sb = System.Text.StringBuilder(pageList.Count)
223          //the body is just a series of pages
224          for i = 0 to pageList.Count - 1 do
225              let page = pageList.[i]
226              //create an internal link for navigation purposes
227              //each page is a div with classes for page and pageType along
228              with a pagination based internal link
229              sb.Append("<div class=\"page " +
230              page.Type.ToString().ToLower() + "\" id=\"" + i.ToString() + "\">")
231              ignore
232              sb.Append( GetPageNav i pageList ) |> ignore
233              sb.Append( RenderPageElementsToHTML fileName page ) |> ignore
234              sb.Append("</div>\n") |> ignore
235              //return everything in the string builder
236              sb.ToString()
237
238          //get a list of chapters and sections within them
239          let chapterList = GetElementByChapterList pageList ( function | Header
240          ( Section (text) ) -> Some( text ) | _ -> None)
241
242          //return main html with body expanded
243          let sb = System.Text.StringBuilder()
244          sb.Append(
```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```
242 <!DOCTYPE XHTML= http://www.w3.org/1999/xhtml>
243 <head>
244     <title>"" + fileName.Replace(".pdf.xhtml", "") + """</title>
245     <link rel="stylesheet" type="text/css" href="../neetstyle.css"/>
246 </head>
247 <body>
248     <!-- Site navigation menu -->
249     <header>
250         <div class="nav">
251             <ul>""") |> ignore
252
253     //add a heading for each chapter
254     for text,index,sections in chapterList do
255         sb.Append("<li class='main-menu'><a href='#" + index.ToString() +
256         "'>" + text + "</a>") |> ignore
257         //add a sublist to heading for each section
258         sb.Append("<ul>") |> ignore
259         for sectionText,sectionIndex in sections do
260             sb.Append("<li><a href='#" + sectionIndex.ToString() + "'>" +
261             sectionText + "</a></li>") |> ignore
262             sb.Append("</ul>") |> ignore //close sublist
263             sb.Append("</li>") |> ignore //close heading
264         sb.Append("</ul>") |> ignore //close all headings
265         sb.Append("</div>") |> ignore //close nav
266         sb.Append(
267             (RenderHTMLBody fileName pageTitle) +
268             """</body>
269             </html>""") |> ignore
270         sb.ToString()
271
272 // Learn more about F# at http://fsharp.net
273 // See the 'F# Tutorial' project for more help.
274 // [<EntryPoint>]
275 let main argv =
276
277     let directoryInfo = new
278         System.IO.DirectoryInfo(@"/z/aolney/research_projects/braintrust/materials
279         /NEETS/xhtml")
280
281     let files = directoryInfo.GetDirectories() |> Array.collect( fun x ->
282         x.GetFiles("*.xhtml"))
283
284     //delete existing pretty files
285     let prettyFiles = directoryInfo.GetDirectories() |> Array.collect( fun
286         x -> x.GetFiles("*pretty*"))
287
288     for file in prettyFiles do
289         file.Delete()
290
291     for file in files do
292         let filePath = file.FullName
293
```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```
290         doc.LoadFromFile
291         doc.DocumentNode.Element("html").Element("body")
292
293
294         //each div is a page
295         //our goal is to transform this div list into a sequence of
296         //structured page objects
297         let divList = new ResizeArray<HtmlNode>( seq { for node in
298             body.Elements("div") do yield node } )
299         let pageList = new ResizeArray<Page>()
300
301         //we need some state because a section may only be marked on the
302         //first page
303         let mutable firstPage = true
304         let mutable prefacePage = 0
305         let mutable lastPageType = TitlePage
306
307         //loop over all divs to construct pages
308         for div in divList do
309
310             //prepare a list of pageElements corresponding to nodes on
311             //page
312             let pageElementList = new ResizeArray<PageElement>()
313
314             //get a list of nodes we can traverse sanely, ignoring #text
315             //nodes
316             let nodeList = ResizeArray<HtmlNode>( seq { for node in
317                 div.ChildNodes do if node.Name <> "#text" then yield node } )
318
319             //construct a page element for each node
320             let mutable lastHeader = Summary "" //we need some state for
321             //learning objectives, seed with junk
322
323             for node in nodeList do
324                 //construct a PageElement
325                 let html = node.OuterHtml
326                 let text = node.InnerText
327                 let pageElement =
328                     match (text,html,lastHeader) with
329                     | IsImageCaption t -> ImageCaption( text )
330                     | IsImageUrl t -> ImageURL(html)
331                     | IsHeader t ->
332                         let header = HeaderClassifier text html
333                         lastHeader <- header
334                         Header( header )
335                     | IsQuestion t -> Question(text)
336                     | IsAnswer t -> Answer(text)
337                     | IsLearningObjective t -> LearningObjective( text )
338                     | _ -> Paragraph(text)
339                 pageElementList.Add(pageElement)
```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```
337        let pageNumber,pageNumberString =
338            let pageNumberElement =
339                pageElementList |> Seq.tryFind(
340                    function
341                      | Header (PageNumber pn) -> true
342                      | _ -> false
343                      )
344
345            match pageNumberElement with
346            | Some( Header (PageNumber pn) ) ->
347                let s = pn.Split('-')
348                if s.Length > 1 then
349                    match System.Int32.TryParse( s.[1] ) with
350                      | true,page -> Some page,pn
351                      | false,_ -> prefacePage <- prefacePage - 1; Some
352                      prefacePage,pn
353
354            else if s.Length = 1 then
355                match System.Int32.TryParse( s.[0] ) with
356                    | true,page -> Some page,pn
357                    | false,_ -> prefacePage <- prefacePage - 1; Some
358                      prefacePage,pn
359
360
361
362            let pageType =
363                //have we hit a page transition
364            let transitionPage =
365                match pageElementList with
366                | IsTOC e -> TOCPage
367                | IsAppendix e -> AppendixPage
368                | IsIndex e -> IndexPage
369                | IsAssignment e -> AssignmentPage
370                | _ -> MainPage //default
371
372                //special cases
373            if firstPage then
374                firstPage <- false
375                TitlePage
376            else if pageNumber.IsSome && pageNumber.Value < 0 &&
377            transitionPage <> TOCPage then
378                PrefacePage
379            else if transitionPage = AppendixPage then
380                AppendixPage
381            else if transitionPage = IndexPage then
382                IndexPage
383            else if transitionPage = AssignmentPage then
384                AssignmentPage
```

```
File      View     Notebook    Help          edited
380                         AppendixPage
387             else if transitionPage = MainPage && lastPageType =
IndexPage then
388                 IndexPage
389             else if transitionPage = MainPage && lastPageType =
AssignmentPage then
390                 AssignmentPage
391             else if transitionPage = TOCPage then
392                 TOCPage
393             else
394                 MainPage
395
396             pageList.Add(
{Number=pageNumber;NumberString=pageNumberString;Type=pageType;Elements=pa
geElementList})
397
398             lastPageType <- pageType
399
400             //Serialize to json
401             let json = Newtonsoft.Json.JsonConvert.SerializeObject(pageList,
Newtonsoft.Json.Formatting.Indented )
402             System.IO.File.WriteAllText( filePath + ".json", json,
System.Text.Encoding.UTF8 )
403
404             //generate a "nice" html page
405             let html = RenderToHTML file.Name pageList
406             System.IO.File.WriteAllText( filePath + "-pretty.html", html,
System.Text.Encoding.UTF8)
407
408             printfn "%A" argv
409             0 // return an integer exit code
```



Most reasonable approach is to section the text by chapters and parse those (but not ► headers). Parsing at the section level is possible but seems a bit too severe. The exact grain size/unit for parsing at this point is debateable. It is only relevant for coreference and discourse parsing; everything else is sentence level.

NOTE THE TEXT HAS WHITESPACE STRIPPED AND IS UNIDECODED, SO IT NO LONGER IS BINARY MATCHED TO THE ORIGINAL

F# FSharp ►

```
1 #r "/z/aolney/repos/Newtonsoft.Json.9.0.1/lib/net40/Newtonsoft.Json.dll"
2 #r "/z/aolney/repos/UnidecodeSharp.1.0.0.0/lib/net35/UnidecodeSharp.dll"
3 open BinaryAnalysis.UnidecodeSharp
4
5 ///tuples so we can align parsed text to pages
6 let cleanRegex = new System.Text.RegularExpressions.Regex("[\s\r\n]+")
7 let GetTextPageTuples ( sections : ResizeArray<string*int> ) =
8   sections
9   |> Seq.map( fun (text,page) -> cleanRegex.Replace( text, " "
10   ".Trim().Unidecode(), page.ToString() )
11   |> Seq.filter( fun (text,page) -> text.Length > 0)
12 let directoryInfo = new
13   System.IO.DirectoryInfo(@"/z/aolney/research_projects/braintrust/materials/
14   NEETS/xhtml")
15 let files = directoryInfo.GetDirectories() |> Array.collect( fun x ->
16   x.GetFiles("*.json") )
17
18 //go from file to json to object
19 let json = System.IO.File.ReadAllText(filePath)
20 let pageList =
21   Newtonsoft.Json.JsonConvert.DeserializeObject<ResizeArray<Page>>(json)
22
23 //get only paragraphs of text we want, write to file
24 let chapterList = GetElementByChapterList pageList ( function | Paragraph
25   (text) -> Some( text ) | _ -> None)
26
27 //for each chapter in chapter list, create a parseable file
28 for text,index,sections in chapterList do
29   let texts,pages = sections |> GetTextPageTuples |> Seq.toList |>
30   List.unzip
31   System.IO.File.WriteAllLines( filePath + "-" + text.Replace("\n","") +
32   ".toparse", texts)
```

Use CLU on input text (WARNING: THIS TAKES ABOUT 1 HR) ▶

ONLY PERFORM IF YOU NEED TO REGENERATE SER

Using .toparse files prepared by fsharp

- load each file
- parse using clu
- serialize as .ser

Sc **Scala**

```
1 import java.io.File
2 def recursiveListFiles(f: File): Array[File] = {
3   val these = f.listFiles
4   these ++ these.filter(_.isDirectory).flatMap(recursiveListFiles)
5 }
6
```



```
import java.io.File
recursiveListFiles: (f: java.io.File)Array[java.io.File]
```

Sc **Scala**

```
1 import edu.arizona.sista.processors.corenlp.CoreNLPProcessor
2 import java.io._
3 import edu.arizona.sista.utils.Files
4 import scala.collection.mutable.ListBuffer
5 import edu.arizona.sista.processors.Processor
6 import edu.arizona.sista.utils.Serializer
7 import java.io.File
8 import scala.io.Source
```

```

11
12 val filesToParse = recursiveListFiles( new File( rootDirectory ) ).filter(
13   f=> """.*\.\toparse$""".r.findFirstIn(f.getName).isDefined)
14
15 //300 word limit is empirically determined based on some sample sentences
16 val proc:Processor = new CoreNLPProcessor(withDiscourse =
17   true,maxSentenceLength = 300)
18
19 //ASSUMES EACH LINE IS A SENTENCE -- CLASS5 ASSUMPTION
20 /*
21 def fileToSentences(file:File):List[String] = {
22   val sents = new ListBuffer[String]
23   io.Source.fromFile(file).getLines().foreach(sents += _)
24   sents.toList
25 }
26 for( file <- filesToParse ) {
27   println(s"Parsing and serializing $file...")
28   //val doc = proc.annotateFromSentences( fileToSentences( file ) )
29   val text = Source.fromFile(file, "utf-8").getLines().mkString(" ")
30   val doc = proc.annotate( text )
31   Serializer.save(doc, file + ".ser")
32 }
33
34

```

9 lines of stdout, 10 lines of stderr

– Map Parse to JSON and JSON to F#

TODO add a bit here to add SRL annotation (or do the SRL parse). See
[/z/aolney/research_projects/csal/analysis/cloze-memory/cloze-generation-0416.bkr](#)

File View Notebook Help edited

```
1 import java.io.PrintWriter
2 import java.io.File
3 import org.json4s._
4 import org.json4s.JsonDSL._
5 import org.json4s.native.JsonMethods._
6 import edu.arizona.sista.discourse.rstparser.DiscourseTree
7 import edu.arizona.sista.discourse.rstparser.RelationDirection
8 import edu.arizona.sista.discourse.rstparser.TokenOffset
9 import scala.collection.mutable.ListBuffer
10 import edu.arizona.sista.utils.Serializer
11 import edu.arizona.sista.struct.DirectedGraphEdgeIterator
12 import edu.arizona.sista.utils.Files
13
14
15 def dTreeJSON (dTree:DiscourseTree) : JValue = {
16   val json = dTreeJSONRec( dTree, JObject(), true, true, 0 )
17   return json
18   //compact( render( json ) )
19 }
20
21 def dTreeJSONRec ( dTree:DiscourseTree, argJson:JValue,
22   printChildren:Boolean, printKind:Boolean, depth:BigInt): JValue = {
23   var json = argJson
24
25   if (printKind) {
26     json = json merge JObject(JField("kind",
27       JString(dTree.kind.toString())))
28   }
29
30   if (dTree.relationLabel.length > 0) {
31     json = json merge JObject(JField("relLabel",
32       JString(dTree.relationLabel)))
33     if (dTree.relationDirection != RelationDirection.None) {
34       json = json merge JObject(JField("relDir",
35         JString(dTree.relationDirection.toString())))
36     }
37
38   json = json merge JObject(JField("depth", JInt( depth )))
39   json = json merge JObject(JField("sentence",
40     JInt(dTree.firstToken.sentence)))
41   json = json merge JObject(JField("firstToken",
42     JInt(dTree.firstToken.token)))
43   json = json merge JObject(JField("lastToken",
44     JInt(dTree.lastToken.token)))
45
46   if (dTree.rawText != null) {
47     json = json merge JObject(JField("text", JString(dTree.rawText)))
48   }
49
50   if (printChildren) {
```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```
4/         if (kids.length > 0)
48            json = json merge JObject(JField("kids", JArray(kids.toList)))
49         }
50     }
51     return json
52 }
53
54 def safeGet(x: Option[Array[String]], i: Int) = x match {
55     case Some(s) => s(i)
56     case None => ""
57 }
58
59 //-----
60 //loop over docs in working directory
61
62 val parsedFiles = Files.findFiles(beaker.workingDirectory.toString(),
63     "ser")
63 val rootDirectory =
64     "/z/aolney/research_projects/braintrust/materials/NEETS/xhtml"
64 val parsedFiles = recursiveListFiles( new File( rootDirectory ) ).filter(
65     f=> """.*\.ser$""".r.findFirstIn(f.getName).isDefined)
66
67 for( file <- parsedFiles ) {
68     //deserialize parse
69     val doc =
70         Serializer.load[edu.arizona.sista.processors.corenlp.CoreNLPDocument](
71         file.getPath() )
72     //we will write the parse to JSON file
72     var jsonOutput = JObject()
73
74     //output discourse parse as json; others
75     https://github.com/clulab/processors/blob/37392ced3a0ebdaf0a7481dccf0ce26991820dba/src/main/scala/edu/arizona/sista/processors/visualizer/DiscourseParserRunner.scala
76     val jsonDiscourseTree = doc.discourseTree map { dTree => dTreeJSON(
77         dTree )} // { dTree => dTree.visualizerJSON() }
76     jsonOutput = jsonOutput merge JObject(JField("discourse",
77         jsonDiscourseTree ))
78
78     //coref as json
79     var chainCorefJson = ListBuffer[JObject]()
80     var chainId = 0
81     doc.coreferenceChains.foreach(chains => {
82         for (chain <- chains.getChains) {
83             //random color for each chain
84             //val color = new Color( (Math.random() *
85                 0x1000000).asInstanceOf[Int])
85             //val colorString = "rgb(" + color.getRed() + "," + color.getGreen()
86                 + "," + color.getBlue() + ")"
```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```

89         //nack: nope that text is distinctive; store associated color; we
90         don't have indices in dTree json
91         val text =
92             doc.sentences(mention.sentenceIndex).words.slice(mention.startOffset,
93                 mention.endOffset).mkString( " " )
94             //json = json merge JObject(JField("color", JString(colorString)))
95             json = json merge JObject(JField("sentence",
96                 JInt(mention.sentenceIndex)))
97             json = json merge JObject(JField("head",
98                 JInt(mention.headIndex)))
99             json = json merge JObject(JField("start",
100                JInt(mention.startOffset)))
101             json = json merge JObject(JField("end",
102                JInt(mention.endOffset)))
103             json = json merge JObject(JField("chainLength",
104                JInt(length)))
105             json = json merge JObject(JField("chainId",
106                JInt(chainId)))
107             json = json merge JObject(JField("text",
108                JString(text)))
109             chainCorefJson += json
110         }
111         chainId += 1
112     }
113 }
114
115 jsonOutput = jsonOutput merge JObject("coreference",
116     chainCorefJson ))
117
118 /*
119  //TODO add SRL with MatePlus
120  //output srl as json
121 var srlJson = new ListBuffer[ListBuffer[JObject]]()
122 var sentenceCount = 0
123 for (sentence <- srlDoc.sentences) {
124     srlJson += new ListBuffer[JObject]()
125     sentence.semanticRoles.foreach(dependencies => {
126         val iterator = new DirectedGraphEdgeIterator[String](dependencies)
127         while(iterator.hasNext) {
128             val dep = iterator.next
129             var json = JObject()
130             json = json merge JObject(JField("sentence",
131                 JInt(sentenceCount)))
132             json = json merge JObject(JField("head",
133                 JInt(dep._1)))
134             json = json merge JObject(JField("token",
135                 JInt(dep._2)))
136             json = json merge JObject(JField("label",
137                 JString(dep._3)))
138             srlJson(sentenceCount) += json
139         }
140     })
141     sentenceCount += 1
142 }
143 */
144
145 //word level as json
146 var wordJson = new ListBuffer[ListBuffer[JObject]]()
147 var sentenceCount = 0

```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```
137     var json = JObject()
138     json = json merge JObject(JField("token",
139       JString(sentence.words(i))))
140     json = json merge JObject(JField("lemma",
141       JString(safeGet(sentence.lemmas,i))))
142     json = json merge JObject(JField("tag", JString(
143       safeGet(sentence.tags,i))))
144     json = json merge JObject(JField("entity", JString(
145       safeGet(sentence.entities,i))))
146     wordJson( sentenceCount ) += json
147   }
148   sentenceCount += 1
149 }
150 jsonOutput = jsonOutput merge JObject(JField("word", wordJson ))
151 //syntactic dep as json
152 var synJson = new ListBuffer[ListBuffer[JObject]]()
153 sentenceCount = 0
154 for (sentence <- doc.sentences) {
155   synJson += new ListBuffer[JObject]()
156   sentence.dependencies.foreach(dependencies => {
157     val iterator = new DirectedGraphEdgeIterator[String](dependencies)
158     while(iterator.hasNext) {
159       val dep = iterator.next
160       var json = JObject()
161       json = json merge JObject(JField("sentence", JInt(sentenceCount)))
162       json = json merge JObject(JField("head", JInt(dep._1)))
163       json = json merge JObject(JField("dependent", JInt(dep._2)))
164       json = json merge JObject(JField("label", JString(dep._3)))
165       synJson(sentenceCount) += json
166     }
167   })
168   sentenceCount += 1
169 }
170 jsonOutput = jsonOutput merge JObject(JField("dependencies", synJson ))
171 //write the combined object of all features for this session Id to file
172 //NOTE: this will require some post processing to make features
173 val pw = new PrintWriter(new File( file.getPath() + ".json" ))
174 pw.write( pretty(render(jsonOutput)) )
175 pw.close
177 }
```

Run

F# FSharp ▶

```
1 type DiscourseTree = {
2     kind : string
3     relLabel : string
4     relDir : string
5     sentence : int
6     firstToken : int
7     lastToken : int
8     text : string
9     kids : DiscourseTree array
10 }
11
12 type Coreference = {
13     sentence : int
14     head : int
15     start : int
16     ``end`` : int
17     chainLength : int
18     chainId : int
19     text : string
20 }
21
22 type Word = {
23     token : string
24     lemma : string
25     tag : string
26     entity : string
27 }
28
29 type Dependency = {
30     sentence : int
31     head : int
32     dependent : int
33     label : string
34 }
35
36 type CluResult = {
37     discourse : DiscourseTree
38     coreference : Coreference array
39     word : Word array array
40     dependencies: Dependency array array
41 }
42
43 type PageParagraphSentenceIndex = {
44     page : int array
45     paragraph : int array
46     sentence : int array
```

Fsharp clean up Scala JSON and output final objects we will use

F# FSharp



```
1 ///Trivially we can traverse the tree and map to sentence indices. However,
2 since kind is broken, we must generate a new tree during traversal with
3 kind corrected
4 let MapDiscourseTreeToSentences( clu : CluResult ) =
5   let sentenceToTreeMap =
6     System.Collections.Generic.Dictionary<int,ResizeArray<DiscourseTree>>()
7
8   let rec discourseDFS( tree:DiscourseTree ) (isNucleus:bool) =
9     let kind =
10       if isNucleus then
11         "nucleus"
12       else
13         "satellite"
14
15     //for now, only map leaves. This excludes the structural information,
16     e.g. relLabel
17
18     //if we wanted we could track the rellabel path to the leaf; see class5
19     code for example
20
21     if tree.text <> null then
22       if sentenceToTreeMap.ContainsKey( tree.sentence ) |> not then
23         sentenceToTreeMap.Add( tree.sentence, new ResizeArray<DiscourseTree>() )
24       sentenceToTreeMap.[tree.sentence].Add( { tree with kind = kind } )
25
26     if tree.kids <> null then
27       discourseDFS tree.kids.[0] (tree.relDir = "LeftToRight")
28       discourseDFS tree.kids.[1] (tree.relDir = "RightToLeft")
29
30   //kick off traversal of tree
31   discourseDFS clu.discourse true
32
33   //return the map
34   sentenceToTreeMap
35
36 //Clu returns a mess of coref. Index by sentence and chain
37 let MapCorefToSentencesAndChains( clu : CluResult ) =
```

Beaker */z/aolney/research_projects/braintrust/analysis/neets-021717.bkr

File View Notebook Help edited

```
33      |> Map.ofSeq
34
35      let chainToCorefMap =
36          clu.coreference
37          |> Seq.groupBy( fun x -> x.chainId )
38          |> Map.ofSeq
39
40      // sentenceToCorefMap,chainToCorefMap
41
42
43
44 let SerializeToJson (filePath:string) ( o ) =
45     let json = Newtonsoft.Json.JsonConvert.SerializeObject(o,
46         Newtonsoft.Json.Formatting.Indented )
47     System.IO.File.WriteAllText( filePath, json, System.Text.Encoding.UTF8
48 )
49 //-----
50 let xhtmlDirectory =
51     @"/z/aolney/research_projects/braintrust/materials/NEETS/xhtml"
52 //builds pretty html and json in preparation of parsing
53 //TraverseHTML (xhtmlDirectory)
54
55 //after parsing, read in JSON from parser and build f# structures
56 let outputList = new ResizeArray<string>()
57
58 let jsonFiles = System.IO.Directory.GetFiles(xhtmlDirectory, "*.ser.json",
59                                              System.IO.SearchOption.AllDirectories)
60
61 for jsonFile in jsonFiles do
62     //load up parse
63     let clu = Newtonsoft.Json.JsonConvert.DeserializeObject<CluResult>(
64         System.IO.File.ReadAllText( jsonFile ) )
65
66     //linearize the discourse tree, align with sentences
67     let discourseBySentences = MapDiscourseTreeToSentences clu
68     SerializeToJson (jsonFile.Replace(".toparse.ser.json",".discourse-
69     sentence")) discourseBySentences
70
71     //align coref with sentences; separately align with chains
72     let sentenceCorefMap,chainCorefMap = MapCorefToSentencesAndChains clu
73     SerializeToJson (jsonFile.Replace(".toparse.ser.json",".coref-sentence"))
74     sentenceCorefMap
75     SerializeToJson (jsonFile.Replace(".toparse.ser.json",".coref-chain"))
76     chainCorefMap
77
78     //now map the parse output to pages (i.e. map sentences to pages
79     let mutable paragraphIndex = 0
80     let mutable lengthAccumulator = 0
```

File View Notebook Help

edited

```
78    for sentenceIndex = 0 to clu.word.Length - 1 do
79        let fauxSentence = clu.word.[sentenceIndex] |> Seq.map( fun x ->
80            x.token ) |> String.concat "" //remove spaces b/c of tokenization
81        lengthAccumulator <- lengthAccumulator + fauxSentence.Length
82        sentenceToParagraphArray.[sentenceIndex] <- paragraphIndex
83        if lengthAccumulator >= toParseLines.[paragraphIndex].Length then
84            lengthAccumulator <- 0
85            paragraphIndex <- paragraphIndex + 1
86        //write page, paragraph, sentence index to file
87        let paragraphPageMap = jsonFile.Replace(".toparse.ser.json", ".page") |>
88            System.IO.File.ReadAllLines |> Array.map System.Int32.Parse
89        let pageArray,paragraphArray,sentenceArray =
90            sentenceToParagraphArray
91            |> Seq.mapi( fun sentence paragraph -> paragraphPageMap.[paragraph],
92                paragraph, sentence )
93            |> Seq.toArray
94            |> Array.unzip3
95        let pageParagraphSentenceIndex =
96            {page=pageArray;paragraph=paragraphArray;sentence=sentenceArray}
97        SerializeToJson (jsonFile.Replace(".toparse.ser.json", ".page-paragraph-
98            sentence"))  pageParagraphSentenceIndex
99
```

Run

Insert FSharp Cell

code ▾

text

section ▾