Industrial and Manufacturing Systems Engineering Conference Proceedings and Posters

Industrial and Manufacturing Systems Engineering

12-15-2017

# Rapid Tagging and Reporting for Functional Language Extraction in Scientific Articles

Mahdi Ramezani
*Iowa State University*, ramezani@iastate.edu

Vijay K. Kalivarapu
*Iowa State University*, vkk2@iastate.edu

Stephen B. Gilbert
*Iowa State University*, gilbert@iastate.edu

Sarah R. Huffman
*Iowa State University*, shuffman@iastate.edu

Elena Cotos
*Iowa State University*, ecotos@iastate.edu

*See next page for additional authors*

Follow this and additional works at: https://lib.dr.iastate.edu/imse_conf

Part of the Communication Technology and New Media Commons, Ergonomics Commons, Graphics and Human Computer Interfaces Commons, and the Publishing Commons

# Rapid Tagging and Reporting for Functional Language Extraction in Scientific Articles

**Abstract**

This paper describes the development of a web-based application for tagging scientific articles, in part to create machine learning training datasets for automated functional language identification and extraction (AFLEX). The initial intent for this work was to provide a new member of the ecosystem of tools that facilitate the structured automation of systematic reviews, an area of work that typically requires critical analysis of multiple research studies and provides an exhaustive summary of literature related to a research question. However, the tool's modular interface allows use across disciplines. A user may upload PDF or text documents and quickly tag selected parts of the document with a customizable set of discipline-specific tags, and export results to CSV or JSON formats. An integrated back-end database stores tagging data for comparison between taggers or visual display of results on the web browser. While other discipline-specific text tagging tools exist, the authors have not encountered a cloud-based customizable tool for PDF and text annotation as flexible as the AFLEX Tag Tool developed by the authors.

**Keywords**

Text tagging, PDF annotation, user-centered design, systematic review, tagging tools

**Disciplines**

Communication Technology and New Media | Computer Sciences | Ergonomics | Graphics and Human Computer Interfaces | Publishing

**Authors**

Mahdi Ramezani, Vijay K. Kalivarapu, Stephen B. Gilbert, Sarah R. Huffman, Elena Cotos, and Annette M. O'Connor

# Rapid Tagging and Reporting for Functional Language Extraction in Scientific Articles

### M. Ramezani
Iowa State University
1620 Howe Hall 537 Bissel Road
Ames, IA 50011-2274
ramezani@iastate.edu

### V. Kalivarapu
Iowa State University
1620 Howe Hall 537 Bissel Road
Ames, IA 50011-2274
vkk2@iastate.edu

### S. B. Gilbert
Iowa State University
1620 Howe Hall 537 Bissel Road
Ames, IA 50011-2274
gilbert@iastate.edu

### S. Huffman
Iowa State University
1137 Pearson 505 Morrill Rd
Ames, IA 50011
shuffman@iastate.edu

### E. Cotos
Iowa State University
1137 Pearson 505 Morrill Rd
Ames, IA 50011
ecotos@iastate.edu

### A. O'Conner
Iowa State University
2424 Vet Med 1800 Christensen Dr
Ames, IA 50011-1134
oconnor@iastate.edu

## ABSTRACT
This paper describes the development of a web-based application for tagging scientific articles, in part to create machine learning training datasets for automated functional language identification and extraction (AFLEX). The initial intent for this work was to provide a new member of the ecosystem of tools that facilitate the structured automation of systematic reviews, an area of work that typically requires critical analysis of multiple research studies and provides an exhaustive summary of literature related to a research question. However, the tool's modular interface allows use across disciplines. A user may upload PDF or text documents and quickly tag selected parts of the document with a customizable set of discipline-specific tags, and export results to CSV or JSON formats. An integrated back-end database stores tagging data for comparison between taggers or visual display of results on the web browser. While other discipline-specific text tagging tools exist, the authors have not encountered a cloud-based customizable tool for PDF and text annotation as flexible as the AFLEX Tag Tool developed by the authors.

## CCS CONCEPTS
• **Human-centered computing~Web-based interaction** • *Human-centered computing~Graphical user interfaces* • **Applied computing~Annotation** • *Applied computing~Document analysis*

## KEYWORDS

Text tagging, PDF annotation, user-centered design, systematic review, tagging tools

## 1 INTRODUCTION
Tagging is a very useful method to annotate unstructured text and documents in preparation for data analysis and retrieval, pattern recognition, grouping similarities (i.e., clustering), etc. Historically, the practice of annotation, or affiliating unstructured text such as commentary with a chosen passage of original text, began with medieval authors of manuscripts. These scholars used notes in the white spaces in documents to conduct debate and commentary [1].

Tagging can be viewed as a subset of annotation, in that tagging typically involves the affiliation of specific structured machine-readable labels (tags) with a chosen passage of original text. These tags are typically used to aid in indexing, search, and categorization of documents.

Although tagging can be done manually for quick processing and retrieval of smaller datasets, it presents a far superior use when such tags can be used as training datasets for machine learning classification, so that once trained, an intelligent system can conduct textual analysis automatically. This approach can efficiently classify large volumes of text in a relatively short period. This approach has been taken by the Automatic Functional Language Extraction (AFLEX) group at the Iowa State University. The goal was to assist in efficiently conducting systematic reviews (SRs), an increasingly popular extensive literature review process in which researchers summarize everything that is known about a

particular topic from thousands of research publications, particularly in the field of evidence-based medicine [2]. One of the required steps in conducting SRs entails classifying existing publications (scientific and technical reports) into pre-set categories, such as randomized controlled trials vs. small case studies. This has predominantly been performed through manual efforts and is very time consuming. This step is of critical interest because machine learning provides an opportunity to automate the SR process by requiring manual classification of a much smaller training set of studies [3].

While the AFLEX Tag Tool was designed to assist the SR process, the cloud-based architecture is flexible and modular to accommodate not just annotation and tagging of medical research articles but any PDF documents and text data alike, across disciplines. Relevant literature, user-centered methods to build the interface, and cross-discipline use cases are discussed below.

## 2 BACKGROUND

### 2.1 Tagging Tools

Annotations and tagging is not a new concept. A researcher may annotate hundreds of articles by manually identifying textual excerpts. Since there are many different tools which vary according to different purposes, the authors introduce just a few examples of tools to illustrate the current procedure of tagging documents.

BRAT [4] is a good example of a modern tagging tool. The tool is designed for tagging plain text and exporting the tagged text as figures or PDF files. It supports linking tags together, and side-by-side comparison of tagged versions of the same text by different taggers, and offers a clean, WYSWYG interface. Also, the BRAT tool uses a standoff data structure, which means that its tags/annotations are stored separately from the original text, but are linked to their corresponding source text spans by character offsets. This approach works well when it is important to avoid editing the source text. However, the BRAT tagging tool does not work with PDF documents because PDFs do not maintain a standard approach to sequencing text blocks, which prevents the use of the character offset approach.

A more traditional approach to annotation of PDFs can be seen in Adobe Acrobat's built-in commenting feature. This procedure is time intensive and vulnerable to human error, since the user is required to type in annotation data via keyboard rather than assigning preset tags. Since the annotations are stored inside the PDF file, they cannot be easily or automatically exported to different formats to be used by other systems. Moreover, extracting any tagging work performed by a user requires manually opening the file in an Acrobat viewer, and its licensing prevents applying third-party scripting to access annotation data. Finally, there is no easy method of comparison between taggers with this type of annotation.

Another more powerful example of a tagging tool is Callisto Annotation Workbench [5], a Java-based open source tagging system designed for linguistic annotation. Callisto's information architecture allows significant data manipulation and comparison, as well as the design of plug-ins for other language based processing. However, Callisto is limited to importing plain text, and its usability for taggers presents a significant challenge, requiring multiple user clicks from dropdown menus to tag a single passage.

### 2.2 Automated Tagging Tools

As noted above, an additional goal of the AFLEX Tag Tool is to facilitate training machine learning classifiers that work with text. Part of that training consists of tagging documents to create a training set, but another key component of the training process includes human-in-the-loop feedback on a classifier's performance. The Tag Tool's user interface, if well designed to support human tagging and review of other users' tags, could presumably be adjusted in small ways to display the results of machine learning automated tagging and allow a user to approve or correct the classifier's judgments. Thus, it is worth exploring other tools that automatically tag documents.

RobotReviewer [6] is an open-source web-based tool that is capable of automatically generating annotations and extracting data from clinical trial reports uploaded in PDF formats, and was developed to aid systematic reviews. The tool however does not provide a mechanism to manually tag documents or to correct the classifier's decisions. As such, it is a view-only tool.

One the leading plagiarism detection tools, iThenticate, www.ithenticate.com (last visited Apr. 29, 2017), is worth noting because of its approach to automatically highlighting plagiarized text in a document. The instances of plagiarism are analogous to tags, in that the tool tags each highlighted passage with its original source and citation. The tool also offers numerical statistics, e.g., the percent plagiarized, and several settings to adjust the automated tagging, e.g., "Don't include quotations."

In both tools, their overall look and feel, with the PDF at left with tagged text highlighted and tag information at right that can be browsed and adjusted, is well designed and offered some inspiration to the design of the AFLEX tag tool's look and feel.

Although these tagging tools and many others have been described in literature [7-9] the following desired features were either non-existent or unavailable in a single application: a) An intuitive interface to quickly tag both PDF and text documents, b) easy exporting of the tagged dataset for use with other systems (e.g., for machine learning), c) a user profile system to crowd source tagging information from multiple subject domain experts and compare their work, and d) a user-centered design that anticipates user input errors and allows easy correction. These features hence formed the baseline requirements for the development of the work presented in this paper.

## 3 AFLEX TAG TOOL USER NEED DISCOVERY

### 3.1 Use Cases

Based on informal interviews with three colleagues in the medical pathology field and two in the linguistics field, we established the following use cases and designed the user interface to support them. Note that sometimes, users wanted to tag text within a PDF. Other times users wanted to tag text within multiple plain text passages.

Rapid Tagging and Reporting for Functional Language Extraction in Scientific Articles

WOSP 2017, June 19, 2017, Toronto, ON, Canada

Thus, two similar but different tagging interfaces were created: Few to Many (Figure 1) and Many to One (Figure 2).

**Table 1: Use Case 1: Tag a few text strings in a PDF with many tags.**

| # | User Functional Requirements | System Support & Response |
|---|---|---|
| 1 | Read PDF and visually locate text. | PDF rendered at readable size in browser with zooming features |
| 2 | Select text within PDF | Selected text highlighted in PDF; plain text extracted to Annotation Box sidebar. |
| 2a | Option: select additional non-contiguous text and they will be grouped together | Additional text highlighted in PDF and plain text extracted and added to Annotation Box. |
| 2b | Option: Delete text selection just made | Text passages in Annotation Box can be deleted individually |
| 3 | Select one or more tags | Tags highlight when selected |
| 3a | Option: deselect one or more tags | Selected tags de-highlight when clicked again |
| 3b | Option: add free response comment text | Textbox accepts plain text. |
| 4 | Click OK to complete tagging | Text passages and tags stored in database and added to Work History in sidebar. Highlights clear from PDF. Tags reset to unselected. Annotation Box cleared. |
| 5 | Note the time spent on tagging | The system records the time taken between OK button clicks to measure the time spent to tag each text excerpt. It is also possible to aggregate the times for a specific document or part of the corpus. |

*3.1.1 Tag Text: Few Text Strings to Many Tags.* The user would like to select one string of text or several non-contiguous strings of text and assign multiple tags to them, e.g., a medical pathology user selects two passages within a medical research PDF and chooses the tags "blinded" and "one-arm parallel design" to indicate that those passages provide evidence that the article is discussing a one-arm parallel experimental design with blinded assignment. Assumption 1: The same text passages or excerpts of them can be tagged multiple times with different tags. Assumption 2: "Few" is less than 5. While technologically there is not a limit to the number of text strings chosen, the user interface is designed to easily display approximately 5 strings, depending on their length. This use case is described in Table 1.

The authors learned from financial colleagues that this use case could also apply for financial analysts, who frequently tag financial documents of publically traded companies using XBRL [7], a specific financial tag set. Portions of a company's income statement might be tagged with XBRL tags such as CashCashEquivalentsEndingBalance, NetCashFlowsUsedOperatingActivities, or IncreaseDecreaseTradeOtherReceivables, for example.

*3.1.2 View Tags by Tagged Text.* The user would like to review what tags have been affiliated with a previously tagged text passage. This is essentially a database query, and the use case is described in Table 2.

**Table 2: Use Case 2: View Tags by Tagged Text.**

| # | User Functional Requirements | System Support & Response |
|---|---|---|
| 1 | Select a passage of text within Work History. | Work History item highlights. Text passages are highlighted in PDF. PDF is scrolled so first text passage is in view. Tags associated with the text are highlighted in sidebar. Tool is now in Edit mode. |
| 1a | Revise the tags associated with an excerpt of the text. | The system allows tags for that excerpt to be toggled on and off and for the comment to be edited. |

*3.1.3 View Tagged Text by Tag.* The user would like to review what text passages have been tagged with a specific tag. This is essentially a database query, and is described in Table 3.

**Table 3: Use Case 3: View Tagged Text by Tag.**

| # | User Functional Requirements | System Support & Response |
|---|---|---|
| 1 | Select a tag to query. | Tag highlights. Text passages associated with that tag are highlighted in PDF. PDF is scrolled so first text passage is in view. Badge number appears on tag indicating number of query results in PDF. |

*3.2.4 Tag Text: Many Text Strings to One Tag.* The user would like to select multiple strings of text in a plain text passage and tag them with a single tag, e.g., a linguistics user selects all the prepositional phrases in a sentence and tags them "prepositional phrase." Assumption: the multiple strings of text affiliated with the single tag will not overlap. This use case is described in Table 4 below.

*3.2.5 View Tagged Text by Tag.* The user would like to review what text strings have been tagged with a specific tag. This is essentially a database query, and is described in Table 5 below**.**

## 4    AFLEX USER INTERFACE

Version 1.0 of this tool follows techniques derived from Agile Unified Process (AUP) [10], and fulfills most of the use cases described above. A handful of freely available server and client side software libraries were used in development including: a) Hypertext Preprocessor (PHP) for server side scripting, b) MariaDB database server, an open source variant of MySQL (mariadb.com), c) Twitter's bootstrap for client side CSS styling (getbootstrap.com), d) JQuery by Google for JavaScript, and e) PDF.js for rendering PDF documents on the browser (github.com/mozilla/pdf.js).

**Table 4: Use Case 4: Tag Many Strings to One Tag.**

| # | User Functional Requirements | System Support & Response |
|---|---|---|
| 1 | Read plain text passage. | Plain text rendered at readable size in browser. |
| 2 | Select tag to work with. | Tag highlighted in sidebar. |
| 3 | Select string of text. | Text highlights when selected |
| 3a | Option: Select an additional string of text that is not selected. | Selected tags de-highlight when clicked again |
| 3b | Option: Delete highlight from a string. | A selected highlight shows a button for deletion. |
| 4 | Click OK to complete tagging | Text passages and tag stored in database. Highlights clear from text passage. Tag changes to "tag used" color and gains badge number with count of strings tagged. |

**Table 5: Use Case 5: View Tagged Text by Tag.**

| # | User Functional Requirements | System Support & Response |
|---|---|---|
| 1 | Select a tag to query. | Tag highlights. Text strings in the text passage associated with that tag are highlighted. |

Although OAuth user authentication from Google is currently implemented, the code is designed for use with other authentication schemes (e.g., single sign-on shibboleth authentication from an educational institution). Extracting a clean body of text from a PDF is not always trivial, as demonstrated by the large literature from the field of digital libraries on how best to do so. Our method, when extraction of the entire PDF text is required, is documented in [11], and is based on font size analysis. For example, the text characters in a PDF document whose font size occurs the most can be construed as the main body of the text. Such font based heuristics can be applied to extract different sections of a PDF. The following

sub-sections detail the implementation of the AFLEX Tag Tool so that it meets the requirements described in Section 2.

Figure 1 shows a screenshot of the AFLEX Tag Tool, with various regions on the screen identified by their functionality. PDF.js library renders a PDF document on the browser as illustrated by boxed region (a) in the figure. Features built into PDF.js libraries such as text selection and highlighting, and search for text can be accessed via PHP scripts and are implemented within the AFLEX Tag Tool, as can be seen in Box (a). Box (b) of the screenshot features an annotation box with text passages that have been highlighted within the PDF. Each selected text passage is marked as one annotated statement and multiple text passages can be selected and added to the list in the annotations box. Subject-specific tags pre-populated from a MySQL database are shown in box (c). Multiple sentences can simultaneously be assigned a certain tag set. For example, the arrows in box (c) indicate that three tags are assigned to the two text excerpts listed in the annotation box (b). The tag names themselves can be edited, new ones added or existing ones deleted on the fly by the user, which also updates the titles in the server database. A user can also add specific comments for a tagging session within the interface, as seen in Figure 3. Box (d) displays the Work History for a certain user, where a list of all sentences that were tagged can be either re-worked or deleted if deemed unsuitable. Also, a user can access any document that he/she has tagged in the past from his/her user profile and re-work or modify tagging as suited. In addition to PDF documents, the AFLEX interface was designed to work with text input as well, meaning that a user can either upload text files or copy paste text into the browser and tag the text.

The interface also keeps track of the time spent by a user in tagging activities. A JavaScript based clock ticks as a user begins a mouse click activity and the elapsed time is saved in the database along with other tagging data. This feature is added to perform qualitative analysis of time spent by multiple subject matter experts on a single research article, intended for a later use.

The AFLEX Tag Tool interface also supports exporting tagged data into CSV and JSON formats. These files can serve as an input to machine learning computational algorithms for classification. With the assumption that tagging is performed by subject domain experts, the architecture was designed so that a user has access to tagging data from all users that tagged a specific PDF document. This architecture will enable users to compare, contrast and agree on tagging a certain document so that a machine learning algorithm receives the most accurate numerical information for classification.

While no formal evaluation of the AFLEX Tag Tool has been conducted, two systematic reviewers who have used it have commented that they appreciated being able to tag sentences with multiple design elements, as well as having pre-populated tags so that there would be no spelling mistakes while tagging. They also appreciated the ability to add their own comments to any annotated text in addition to tags.
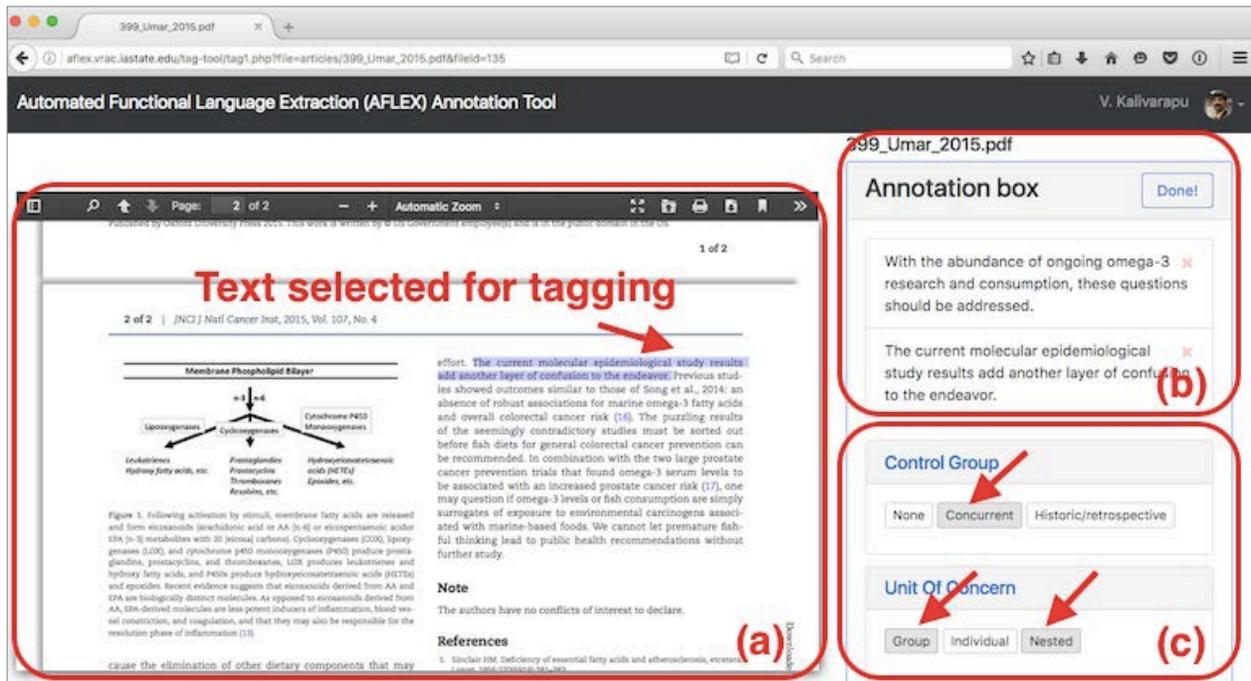
**Figure 1: AFLEX Tag Tool for Few Strings to Many Tags: two sentences in [12] (b) are assigned three tags.**
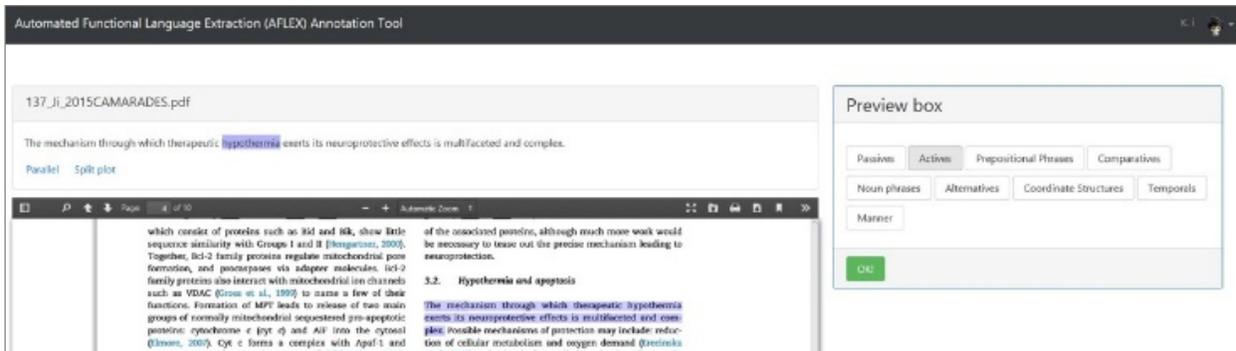


**Figure 2: Many to One tagging interface: A plain text sentence is extracted from a PDF file [13], and the user highlights multiple non-overlapping excerpts of the sentence with a single ta**g such as "Prepositional Phrase."

## 5 CONCLUSION

A web-based interactive tagging/annotating tool for functional language extraction in scientific articles was developed and described in this paper. While tagging tools exist, a simple, easy-to-use cloud-based system for tagging both PDFs and text did not appear to exist. The AFLEX Tag Tool meets these needs with an interface designed to be efficient to use.

With this work, the authors do not intend to replace other established annotation or systematic review automation tools, but rather augment them with this tool, which will have web-based APIs in the future so that third-party tools can send data to it and receive from it. The Tag Tool provides a rapid and intuitive means to tag PDFs or text segments across disciplines. Its output can aid/assist systematic review experts to collaboratively tag articles or for generating training datasets for machine learning and other systems. Because of its modularity and flexibility for adding tag sets, it should be useful for tagging activities in multiple disciplines, including, but not limited to, analysis of data tables, analysis of a range of communication genres (e.g., grant proposals, cover letters), or provision of commentary on texts for critique purposes (e.g., peer review, literary criticism).

## 5.1 Future work

The next steps for the AFLEX Tag Tool include several features. First, a side-by-side conflict visualization tool with agreement metrics will be helpful, to identify interrater reliability and differences in tagging between multiple human taggers or a human and an automated agent. Next, integration with machine learning agents will be helpful so that those systems can effectively tag the text and allow feedback from users refining the classification. Users can provide the feedback to the algorithm by indicating the correct and incorrect tagged sentences, which will help improving the AI tagging process. An interface that promotes user trust in the machine learning classifier is needed for this purpose.
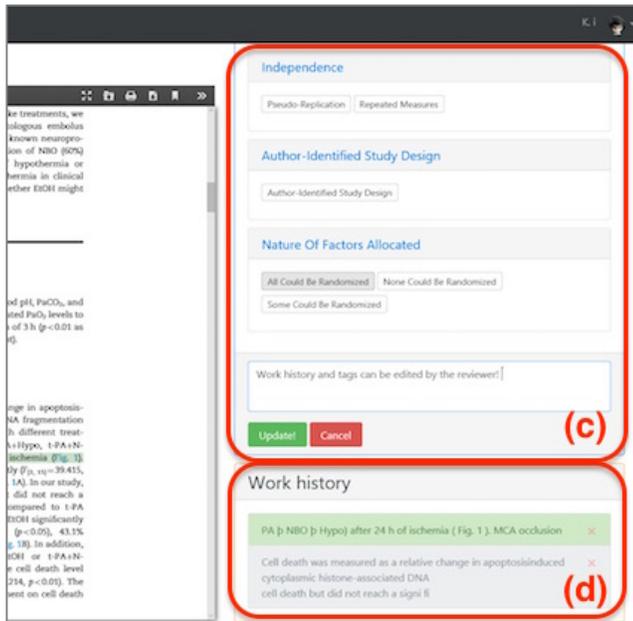


**Figure 3: Tag and Work History boxes.**

Also, future work contains an evaluation effort, in which users who annotate PDFs or text will be observed and have their workflow timed with and without the Tag Tool to establish its gain in efficiency. The authors believe that by taking a user-centered design approach, the Tag Tool, especially when powered by a machine learning engine, can dramatically decrease the time to annotate documents.

## REFERENCES

[1] Joanna Wolfe. 2002. Annotation technologies: A software and research review. Computers and Composition, 19, 4, 471-497.

[2] J.P.T. Higgins and S. Green (eds.). 2011. Cochrane Handbook for Systematic Reviews of Interventions Version 5.1.0 [updated March 2011]. The Cochrane Collaboration, 2011.

[3] G. Tsafnat, P. Glasziou, M. K. Choong, A. Dunn, F. Galgani and E. Coiera. 2014. Systematic review automation technologies. Syst Rev, 3, 74.

[4] Pontus Stenetorp, Sampo Pyysalo, Goran Topić, Tomoko Ohta, Sophia Ananiadou and Jun'ichi Tsujii. 2012. BRAT: a web-based tool for NLP-assisted text annotation. In Proceedings of the Demonstrations at the 13th Conference of the European Chapter of the Association for Computational Linguistics, Association for Computational Linguistics, 102-107.

[5] David S Day, Chad McHenry, Robyn Kozierok and Laurel Riek. 2004. Callisto: A Configurable Annotation Workbench. In LREC.

[6] Iain J Marshall, Joël Kuiper and Byron C Wallace. 2015. RobotReviewer: evaluation of a system for automatically assessing bias in clinical trials. Journal of the American Medical Informatics Association, ocv044.

[7] Seán O'Riain, Edward Curry and Andreas Harth. 2012. XBRL and open data for global financial ecosystems: A linked data approach. International Journal of Accounting Information Systems, 13, 2, 141-162.

[8] Francesco Osborne, Angelo Salatino, Aliaksandr Birukou and Enrico Motta. 2016. Automatic classification of springer nature proceedings with smart topic miner. In International Semantic Web Conference, Springer, 383-399.

[9] Francesco Ronzano and Horacio Saggion. 2015. Dr. inventor framework: Extracting structured information from scientific publications. In International Conference on Discovery Science, Springer, 209-220.

[10] Kent Beck, Mike Beedle, Arie Van Bennekum, Alistair Cockburn, Ward Cunningham, Martin Fowler, James Grenning, Jim Highsmith, Andrew Hunt and Ron Jeffries. 2001. Manifesto for agile software development.

[11] Stephen Gilbert, Nirav Kamdar, Vijay Kalivarapu, Mostafa Amin-Naseri, Elena Cotos and Annette O'Connor. 2016. Extraction of Relevant Text from PDF Research Articles Using Font Analysis. In Workshop on Mining Scientific Publications at the Joint Conference on Digital Libraries.

[12] Asad Umar, Ellen Richmond and Barnett S Kramer. 2015. Colorectal cancer prevention and fishful thinking, Oxford University Press US.

[13] Zhili Ji, Kayin Liu, Lipeng Cai, Changya Peng, Ruiqiang Xin, Zhi Gao, Ethan Zhao, Radhika Rastogi, Wei Han and Jose A Rafols. 2015. Therapeutic effect of tPA in ischemic stroke is enhanced by its combination with normobaric oxygen and hypothermia or ethanol. Brain research, 1627, 31-40.