

# SOCRATex - a Staged Optimization system for Curation, Regularization, and Annotation of clinical Text

Jimyung Park<sup>1</sup>, Seng Chan You, M.D., M.S.<sup>2</sup>, Jin Roh, M.D.<sup>3</sup>, Ph.D, Dongsu Park<sup>2</sup>, Kwang Soo Jeong<sup>2</sup>, Rae Woong Park, M.D., Ph.D<sup>1,2</sup>

<sup>1</sup>Dept. of Biomedical Sciences, Ajou University Graduate School of Medicine, Yeongtong-gu, Suwon, South Korea; <sup>2</sup>Dept. of Biomedical Informatics, Ajou University School of Medicine, Yeongtong-gu, Suwon, South Korea; <sup>3</sup>Dept. of Pathology, Ajou University Hospital, Yeongtong-gu, Suwon, South Korea

Is this the first time you have submitted your work to be displayed at any OHDSI Symposium?

Yes  No

## Abstract

*Unstructured clinical narrative reports may contain invaluable information. To extract precious information from clinical reports, free text must be converted into machine-readable data via Natural Language Processing (NLP). However, Korean clinical reports can be difficult to process for reasons including multilingual contents and the use of different formats. There is still a lack of NLP tools for Korean clinical reports. Hence, a user-friendly NLP tool which can be accommodate the diverse formats of Korean clinical reports is needed. In this research, we aimed to develop a system called SOCRATex, which is optimizing the stages of NLP, including preprocessing, exploration, and annotation of clinical narrative text. SOCRATex is developed based on OMOP-CDM (Observational Medical Outcome Partnership – Common Data Model), so it can be used in other hospitals and for different types of clinical reports. To assess the feasibility of SOCRATex, we applied it to the analysis of pathology reports on colorectal cancer at Ajou University Medical Center.*

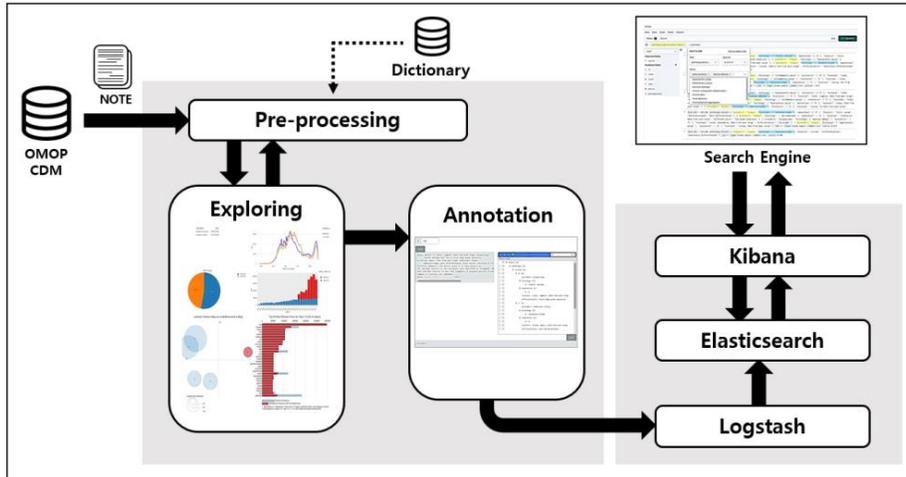
## Introduction

Since unstructured clinical reports are in the form of natural language, they need to be processed and annotated using Natural Language Processing (NLP) for computational accessibility and analysis<sup>1,2</sup>. However, this annotation requires medical experts to review all of the reports, an extremely time-consuming process. The analysis and storage of annotated data is not always possible, due to a lack of suitable tools and databases. Our aim was to establish a usable system by optimizing preprocessing, exploration, annotation, and search engine access based on OMOP-CDM, to enable multi-center distributable research.

## Methods

The pathology reports for colorectal cancer patients from 2014 to 2017 were obtained from Ajou University Medical Center. Colorectal cancer was defined using the SNOMED-CT codes, which were converted from the C18-20 codes of ICD-10 (International Classification of Disease, 10<sup>th</sup> revision). The pathology reports of other anatomical sites were excluded.

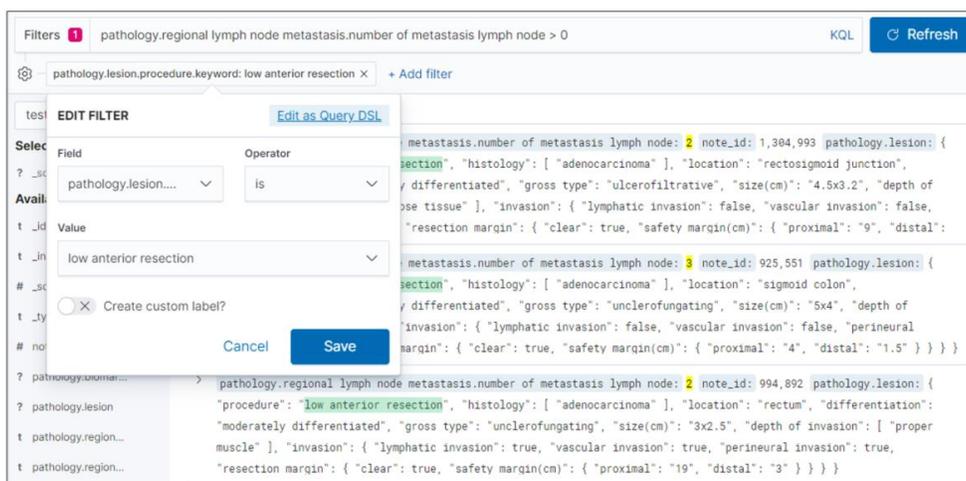
The pipeline of SOCRATex consists of four main processes: preprocessing; exploration; annotation; and searching engine (Figure 1). The algorithms used for preprocessing were aimed at the removal of punctuation, the establishment of a user-defined dictionary, and similar data cleansing. During exploration, Latent Dirichlet Allocation (LDA) can be applied to the preprocessed data for topic clustering<sup>3</sup>.



**Figure 1.** The architecture of the SOCRATex, consisting of preprocessing, exploration, annotation, and ELK stack (Elasticsearch, Logstash, Kibana)

After the identification of topic clusters showing sample documents, researchers can define JSON (JavaScript Object Notation) objects and templates (JSON schema). JSON data can be validated using JSON schema, which forces the use of JSON data types and lists according to the user-definition. Finally, a clinical text search engine can be constructed by sending the JSON data to Elasticsearch. Elasticsearch was developed using the open source software Apache Lucene, and offers real-time searching on the inverted indices of JSON data.

The user-friendly system was developed for researchers who are unfamiliar with NLP techniques or advanced programming skills, using *R version 3.4.4* and the interactive web application development framework *Rshiny*.



**Figure 2.** Sample view of Kibana, searching for pathology reports of low anterior resection with metastasis lymph nodes.

## Results

Overall, 1,985 pathology reports from 1,929 patients with colorectal cancer from 2014 to 2017 were extracted. Only the pathology reports were extracted, by user selection. Abbreviations and stop words such as *a-colon*, *ti*, *dsc*, *bx*, were deleted or replaced using the dictionary and other preprocessing algorithms. LDA, a generative probabilistic model, was used to cluster text into five topics. The results of the topic clustering are listed in (Table 1). The JSON schema was defined mainly using a pathology result. The pathology result is composed of regional lymph nodes, biomarkers, and lesions consist of annotations such as histology, location, differentiation, gross type, invasion, and resection margins. The JSON data is sent to Elasticsearch to construct the clinical text search engine (Figure 2).

## Conclusion

SOCRATex was used to demonstrate the process of clinical document annotation and construction of the search engine using validation with pathology reports about colorectal cancer. SOCRATex allows users to extract, process, and annotate clinical text in OMOP-CDM. To validate the usability and feasibility of the system, 1,985 pathology reports from Ajou University Hospital were analyzed. The reports were transformed into JSON format and sent to Elasticsearch. Consequently, a scalable clinical text search engine based on OMOP-CDM was created and validated using 1,985 pathology reports.

**Table 1.** Top most 30 terms of topics from Latent Dirichlet Allocation

Topics	Terms
Topic1	<i>biopsy</i> , <i>all</i> , <i>consists</i> , <i>xxcm</i> , <i>embedded</i> , <i>mucosal</i> , <i>received</i> , <i>measuring</i> , <i>diagnosis</i> , <i>sections</i> , <i>tissue</i> , <i>pieces</i> , <i>labelled</i> , <i>gross</i> , <i>biopsied</i> , <i>adenocarcinoma</i> , <i>cancer</i> , <i>differentiated</i> , <i>colon</i> , <i>moderately</i> , <i>rectal</i> , <i>verge</i> , <i>rectum</i> , <i>anal</i> , <i>four</i> , <i>sigmoid</i> , <i>endoscopic</i> , <i>largest</i> , <i>five</i> , <i>one</i>
Topic2	<i>anal</i> , <i>verge</i> , <i>colon</i> , <i>one</i> , <i>tubular</i> , <i>adenoma</i> , <i>low</i> , <i>grade</i> , <i>dysplasia</i> , <i>biopsy</i> , <i>transverse</i> , <i>polypectomy</i> , <i>containers</i> , <i>each</i> , <i>ascending</i> , <i>identified</i> , <i>consists</i> , <i>two</i> , <i>polyp</i> , <i>largest</i> , <i>sigmoid</i> , <i>descending</i> , <i>polypoid</i> , <i>hyperplastic</i> , <i>mucosal</i> , <i>proximal</i> , <i>endoscopic</i> , <i>polyps</i> , <i>xxcm</i> , <i>three</i>
Topic3	<i>margin</i> , <i>resection</i> , <i>mass</i> , <i>lymph</i> , <i>invasion</i> , <i>node</i> , <i>regional</i> , <i>xcm</i> , <i>metastasis</i> , <i>distal</i> , <i>apart</i> , <i>len</i> , <i>pericollic</i> , <i>identified</i> , <i>proximal</i> , <i>carcinoma</i> , <i>circumference</i> , <i>nodes</i> , <i>fresh</i> , <i>some</i> , <i>free</i> , <i>illdefined</i> , <i>state</i> , <i>cut</i> , <i>instability</i> , <i>test</i> , <i>msimicrosatellite</i> , <i>bat</i> , <i>invades</i> , <i>iple</i>
Topic4	<i>invasion</i> , <i>adenoma</i> , <i>resection</i> , <i>margin</i> , <i>submitted</i> , <i>consu</i> , <i>ation</i> , <i>hampe</i> , <i>grade</i> , <i>histopathologic</i> , <i>stained</i> , <i>size</i> , <i>adenocarcinoma</i> , <i>dysplasia</i> , <i>high</i> , <i>tumor</i> , <i>tubovillous</i> , <i>type</i> , <i>depth</i> , <i>low</i> , <i>biopsy</i> , <i>gross</i> , <i>well</i> , <i>tubular</i> , <i>labelled</i> , <i>differentiated</i> , <i>polypectomy</i> , <i>colon</i> , <i>endoscopic</i> , <i>whitish</i>
Topic5	<i>kras</i> , <i>mutation</i> , <i>analysis</i> , <i>dna</i> , <i>realtime</i> , <i>clamping</i> , <i>pcr</i> , <i>codon</i> , <i>comments</i> , <i>antiegfr</i> , <i>rapy</i> , <i>msi</i> , <i>using</i> , <i>genomic</i> , <i>isolated</i> , <i>mediated</i> , <i>paraffinembedded</i> , <i>target</i> , <i>cetuximab</i> , <i>panitumumab</i> , <i>marker</i> , <i>pnamediated</i> , <i>materials</i> , <i>erlotinib</i> , <i>gefitinib</i> , <i>kinase</i> , <i>tyrosine</i> , <i>inhibitor</i> , <i>pna</i> , <i>additional</i>

## Acknowledgement

This work was supported by a grant of the Korea Health Technology R&D Project through the Korea Health Industry Development Institute (KHIDI), funded by the Ministry of Health & Welfare, Republic of Korea [grant number : HI16C0992] and the Bio Industrial Strategic Technology Development Program (20005021) funded By the Ministry of Trade, Industry & Energy (MOTIE, Korea).

## References

1. Wang Y, Wang L, Rastegar-Mojarad M, Moon S, Shen F, Afzal N, et al. Clinical information extraction applications: A literature review. *J Biomed Inform.* 2018;77:34-49.
2. Hassanpour S, Langlotz CP. Information extraction from multi-institutional radiology reports. *Artif Intell Med.* 2016;66:29-39.
3. Sievert C, Shirley KE. LDAvis: A method for visualizing and interpreting topics. *Proceedings of the Workshop on Interactive Language Learning, Visualization, and Interfaces.* 2014;pp 63–70.