



# SOCRATex

## Staged Optimization of 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. Ph.D<sup>4</sup>, Rae Woong Park M.D. Ph.D<sup>1,2,3</sup>

1 Dept. of Biomedical Sciences, Ajou University Graduate School of Medicine, Yeongtong-gu, Suwon, 16499

2 Dept. of Biomedical Informatics, Ajou University School of Medicine, Yeongtong-gu, Suwon, 16499

3 FEEDER-NET(Federated E-Health Big Data for Evidence Renovation Network)

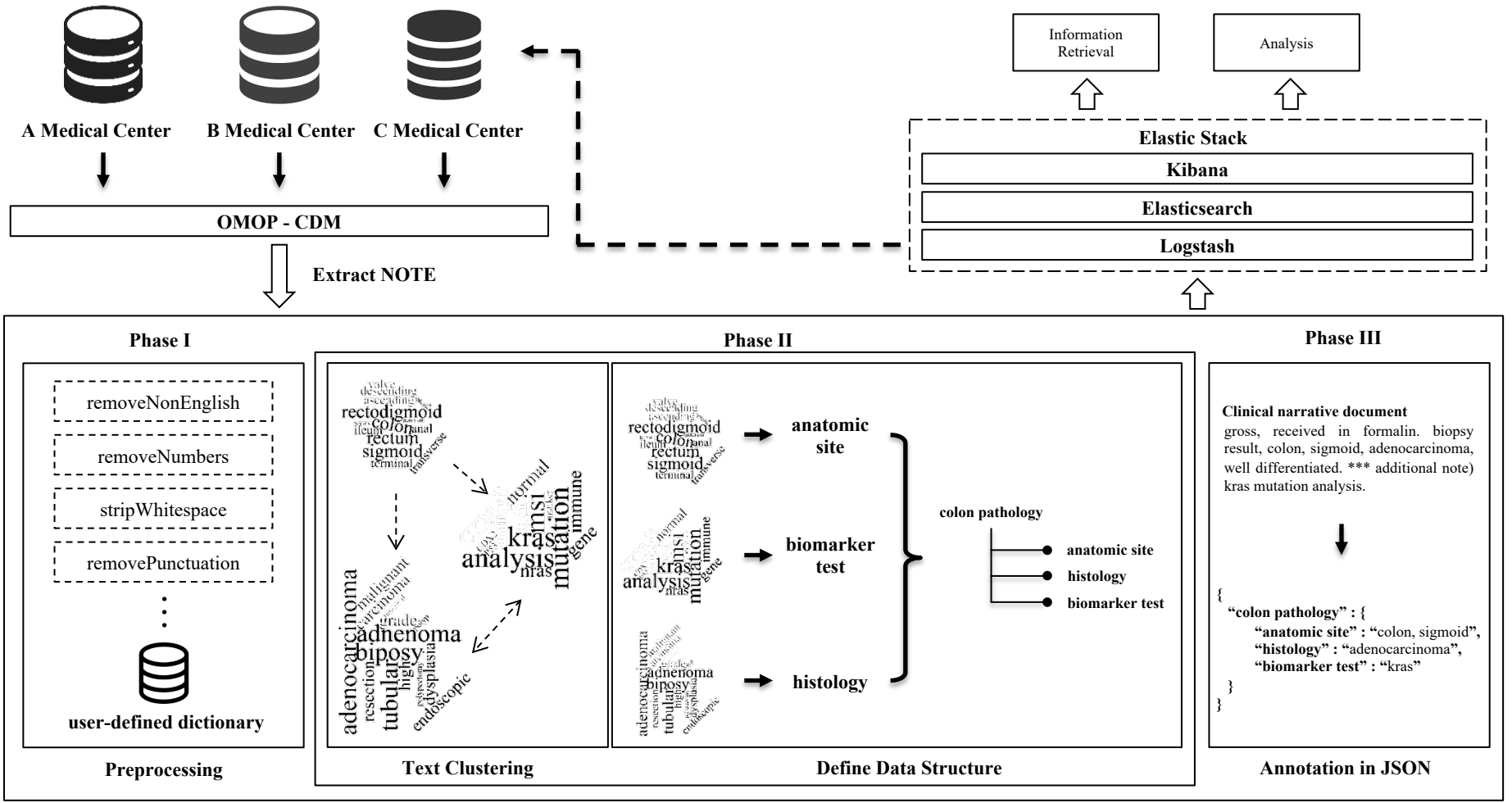
4 Dept. of Pathology, Ajou University Hospital, Yeongtong-gu, Suwon, 16499



# Background

- State-of-the-art (SOTA) methods made great stories on Natural Language Processing (NLP) tasks
- Yet, the SOTA methods usually require **massive amounts of labeled-data** to learn
- SOCRATex is a Natural Language Processing (NLP) system which helps users to **understand**, **annotate**, and **retrieve** their text documents

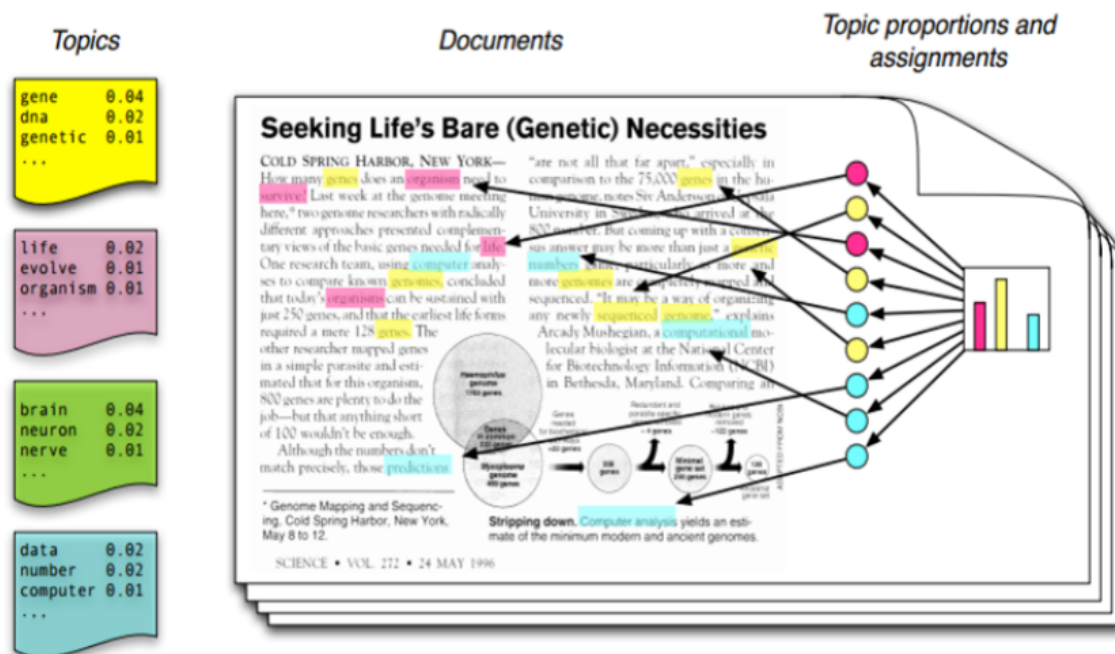
# System Architecture of SOCRATex





# Latent Dirichlet Allocation

- Latent Dirichlet Allocation (LDA) is a statistical topic model method which captures **latent topics** in documents



- I. LDA is an unsupervised method which can reduce **time** and **cost** of reviewing the reports
- II. LDA results can give an insight **how to annotate** and **re-construct** their text data



# Elastic Stack

- **Elasticsearch** is a open-source indexing framework for searching, which can accept tree-based documents such as JSON documents.
- **Kibana** helps to visualize indexed data stored in Elasticsearch

ELASTICSEARCH + KIBANA

## Combining powerful products and features

Built on an open source foundation, Elasticsearch and Kibana pave the way for diverse use cases that start with logging and span as far as your imagination takes you. Elastic features like [machine learning](#), [security](#), and [reporting](#) compound that value — and since they're made for Elastic, you'll only find them from us. [See a full list of Elastic Stack features.](#)



### Elasticsearch

Elasticsearch is a distributed, JSON-based search and analytics engine.

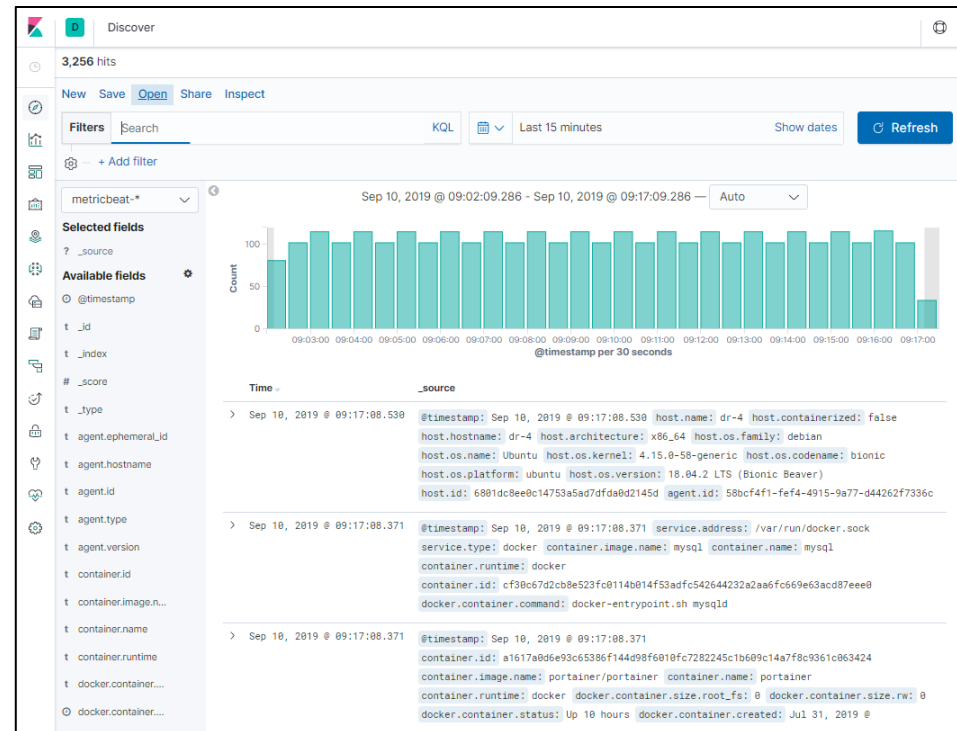
[Learn more](#)



### Kibana

Kibana gives shape to your data and is the extensible user interface.

[Learn more](#)

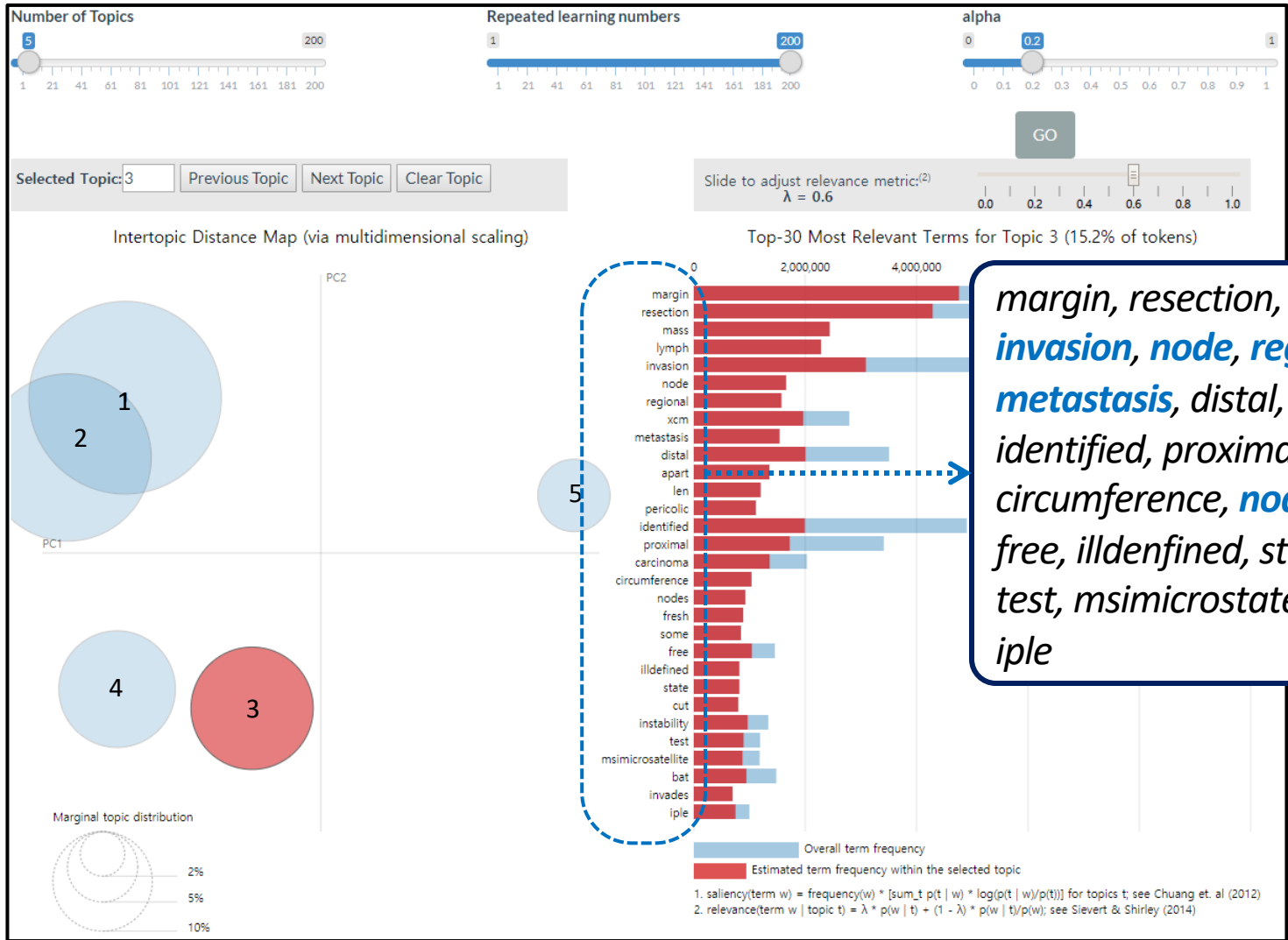




# Data Source

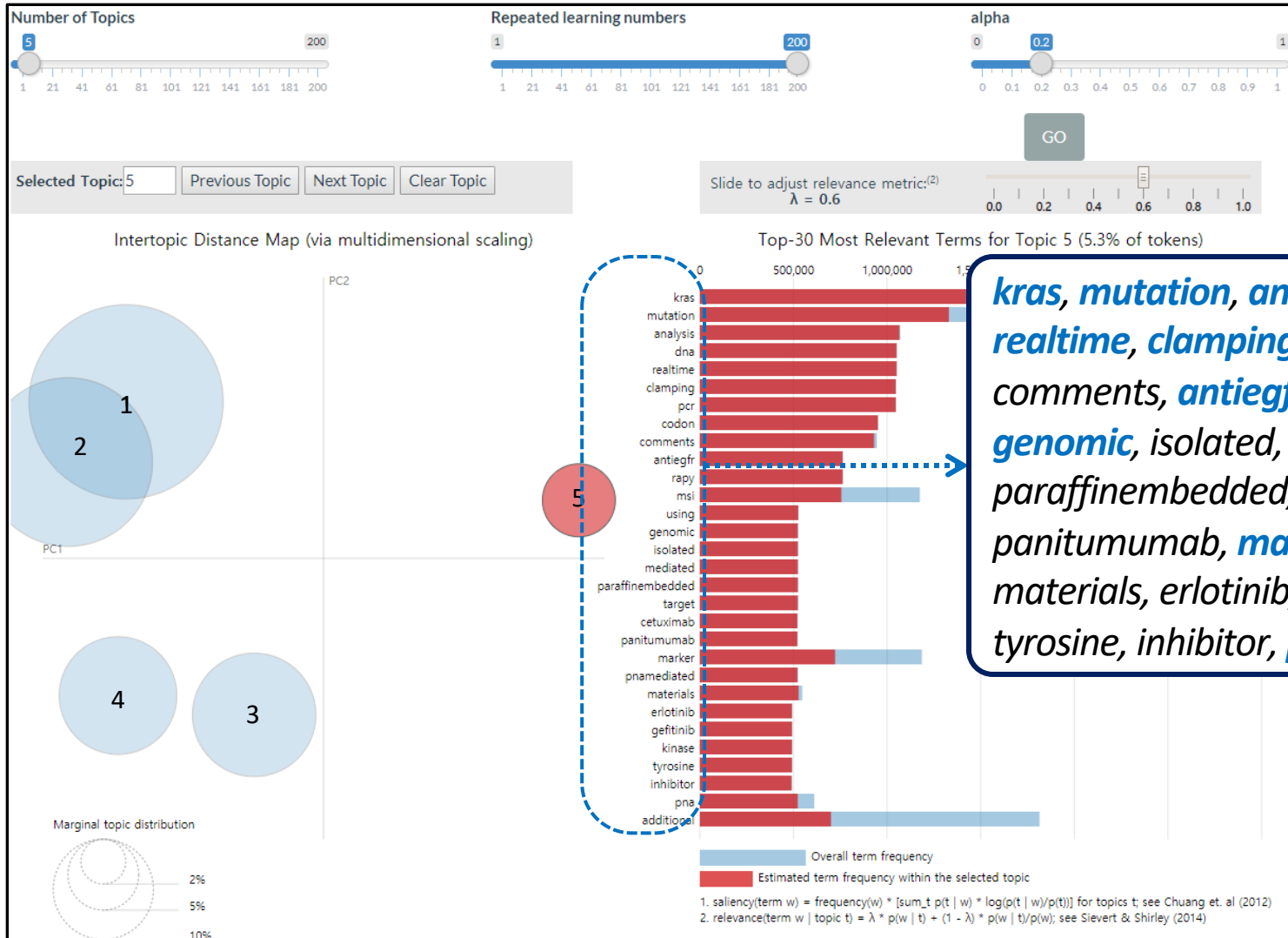
- Ajou University Medical Center
- ICD-10<sup>th</sup> C18-20 diagnosed patients from 2014-2017 were included
- 1,989 pathology reports on colorectal cancer of 1,929 patients were included

# Results of LDA



*margin, resection, mass, lymph, invasion, node, regional, xcm, metastasis, distal, apart, len, pericolc, identified, proximal, carcinoma, circumference, nodes, fresh, some, free, illdenfined, state, cut, instability, test, msimicrostatellite, bat, invades, iple*

# Results of LDA



*kras, mutation, analysis, dna, realtime, clamping, pcr, codon, comments, antiegfr, rapy, msi, using, genomic, isolated, mediated, paraffinembedded, target, cetuximab, panitumumab, marker, pnamediated, materials, erlotinib, gefitinib, kinase, tyrosine, inhibitor, pna, additional*



# Results of LDA

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



# Defining JSON structure based on LDA

## Results

### Topic1

Rectum, endoscopic biopsy: tubulovillous adenoma with high grade dysplasia and adenocarcinoma



- Endoscopic biopsy
- Tubulovillous adenoma
- High grade dysplasia
- Adenocarcinoma



- Procedure
- Histology
- Annotation
- Location
- Differentiation
- Gross type
- Size(cm)
- Depth of invasion
- Underlying pathology

### Topic2

Colon, proximal transverse, polypectomy: Hyperplastic polyp. Tubular adenoma with low grade dysplasia and clear resection margin.



- Polypectomy
- Hyperplastic polyp
- Clear resection



- Procedure
- Histology
- Annotation
- Location
- Differentiation
- Gross type
- Size(cm)
- Depth of invasion
- Underlying pathology

### Topic3

Colon, radical sigmoid colectomy: Adenocarcinoma  
...  
Regional lymph node metastasis: no metastasis in all 17 regional lymph node nodes (pN0) (pericolic 0/17) Lymphatic invasion: not identified



- Colectomy
- Regional lymph node metastasis
- Lymphatic invasion



- Number of metastasis lymph node
- Number of whole lymph node
- Invasion
  - Lymphatic invasion
  - Vascular invasion
  - Perineural invasion
  - Resection margin
    - Clear
    - Proximal
    - Distal
    - Radial

### Topic4

Colon and rectum, Hartman's operation: Adenocarcinoma, moderately differentiated. 1. Location: rectum 2. Gross type: ulcerofungating 3. Size: 7.5x6cm 4. invaded perirectal adipose tissue



- Hartman's operation
- Adenocarcinoma



- Number of metastasis lymph node
- Number of whole lymph node
- Invasion
  - Lymphatic invasion
  - Vascular invasion
  - Perineural invasion
  - Resection margin
    - Clear
    - Proximal
    - Distal
    - Radial

### Topic5

Additional report. Kras Mutation Analysis Report; Kras mutation is not detected by PNA mediated real-time PCR clamping method. Materials: Genomic DNA isolated from paraffin-embedded tissue



- Kras mutation analysis
- Not detected
- PNA mediated real-time PCR clamping method

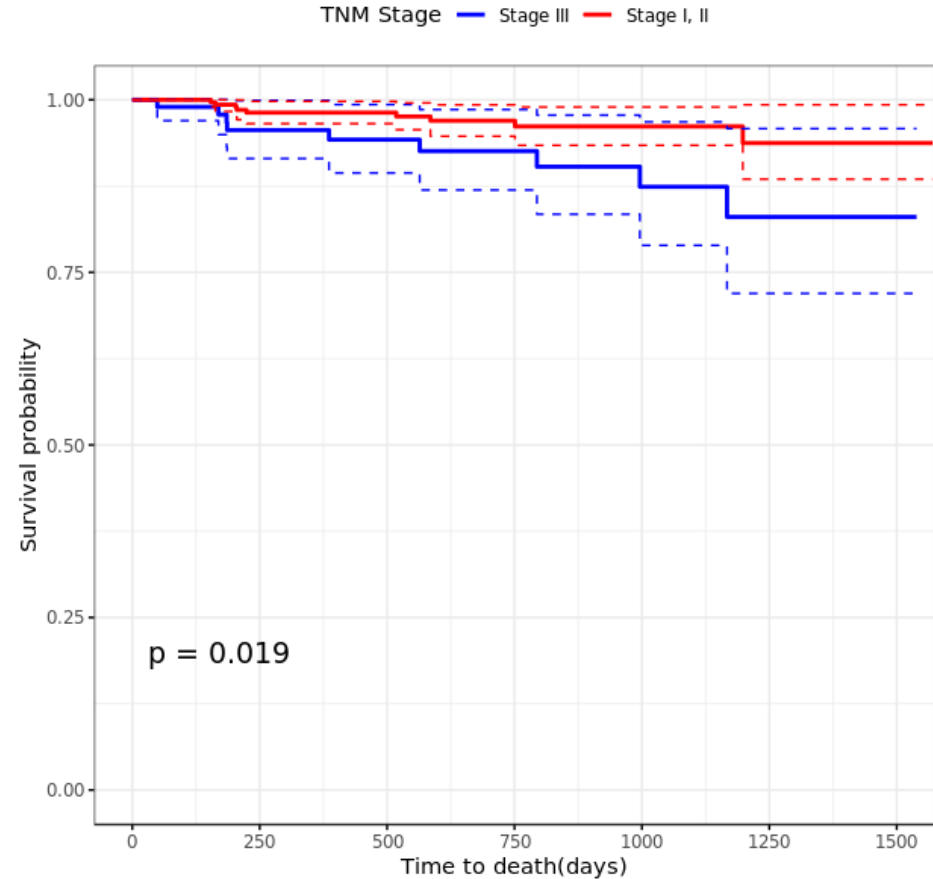


- Biomarker test
- Biomarker method
- Biomarker result



# Clinical Analysis using the Annotated JSON Documents

- Using two keys of JSON structure, **TNM stage** was easily extracted, and **5-year survival analysis** was conducted
  - depth of invasion
  - metastasis lymph node
- This research shows **SOCRATex** can actually be combined with **OMOP CDM** and help to conduct **clinical analysis**
- **SOCRATex** is available at <https://github.com/ABMI/SOCRATex>



## Numbers at risk

Stage I, II	350	242	180	118	73	34	2
Stage III	107	84	59	44	30	15	1



**Thank You**