

# GunNLP at the NTCIR-16 Real-MedNLP Task: Collaborative filtering-based similar case identification method via structured data “case matrix”

Rei Noguchi, Gunma University Hospital

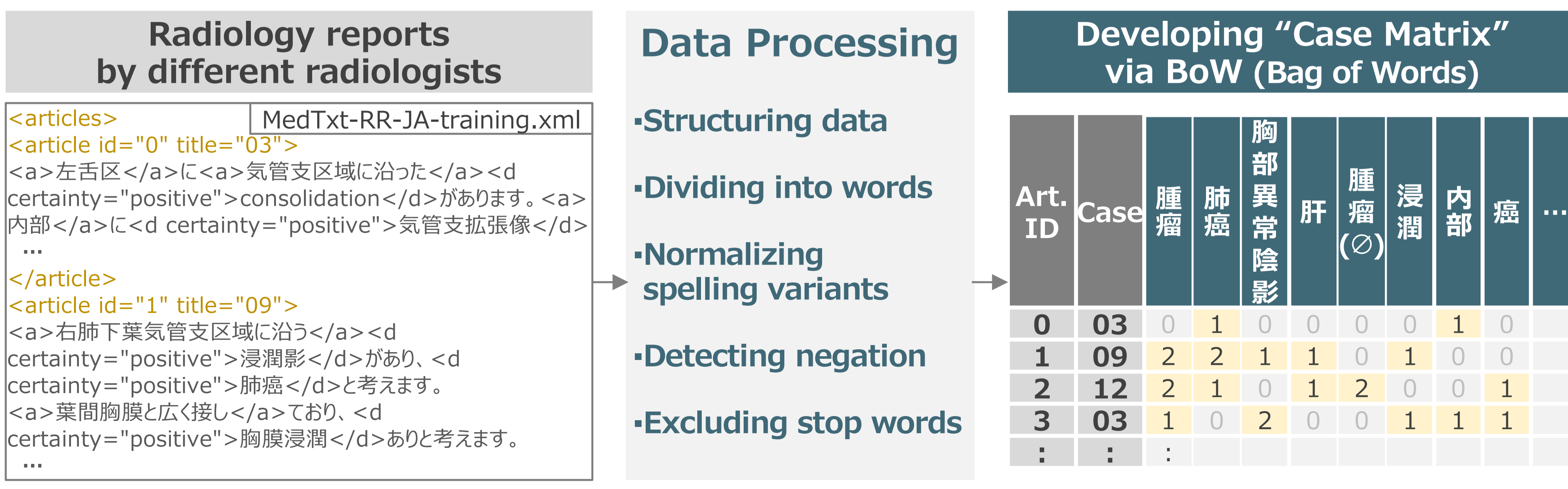
## Introduction

- The clinical text data such as radiology reports are highly expected to be utilized directly for diagnosis support or similar case search using natural language processing and machine learning model.
- Especially in medical fields, the explainability of a machine learning model is extremely important.
- In this study, I propose an explainable case identification model for radiology reports by structuring the reports into a “case matrix” and by applying a collaborative filtering algorithm.

## Method

- I tried to apply **user-based collaborative filtering**, which is mainly used for personalized recommendations, to case identification.
- For the application, given training data of free text description were processed into **structured data “Case Matrix”**, in which each record was unique for each case, and the presence or absence of each symptom was stored in separate columns.
- Normalization of spelling variants** using Manbyo-Dictionary, **detection of negation** and its simple dependency relationship, and **exclusion of stop words** were also conducted.

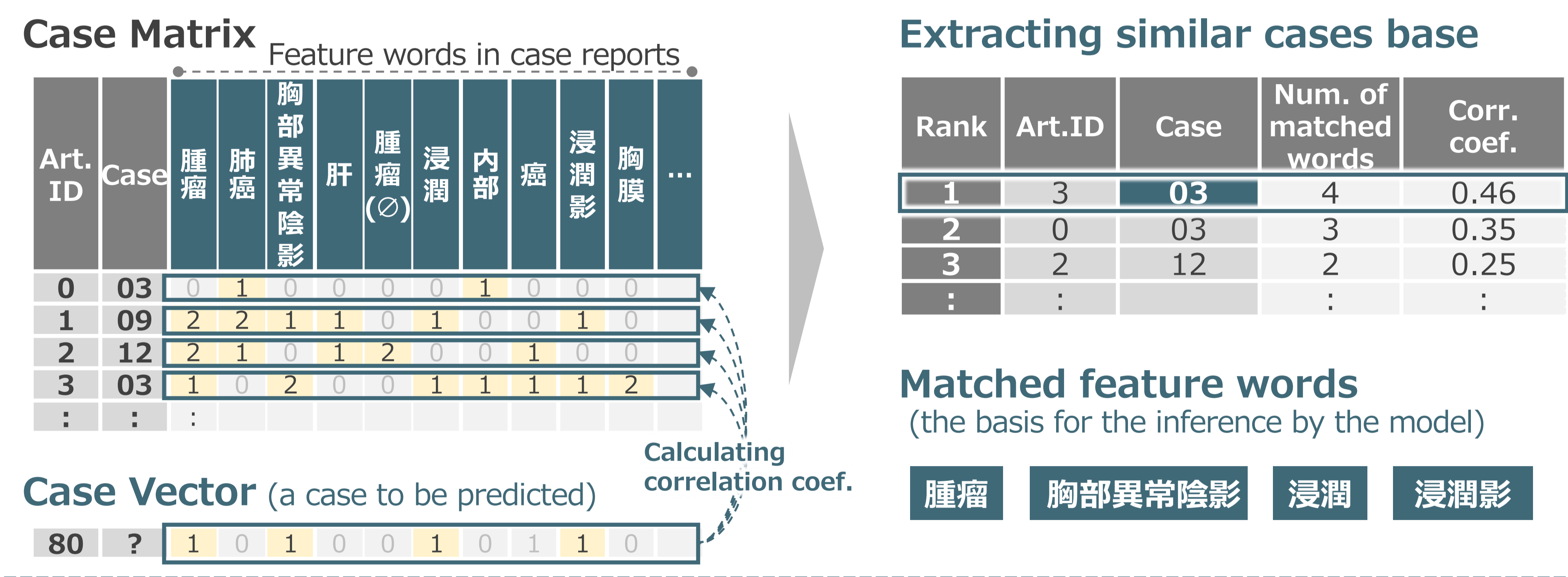
### Scheme of developing “Case Matrix”



## Results

- For training data, a matching rate between the inferred similar case numbers and true ones was **60% (43 of 72)**.
- Dimensionality reduction by sparse modeling decreased the variables from 782 to 55 and improved the matching rate to **71% (51 of 72)**.
- Moreover, this method was able to **provide the basis for the inference by the model** based on the matched words.
- In contrast, for test data, the evaluation result notified by the organizer was 0.3569 in terms of normalized mutual information.

### Scheme of collaborative filtering-based similar case identification



## Conclusion

- The case matrix and collaborative filtering-based method worked for case identification **with explainability** in radiology reports.
- However, the accuracy for test data was much lower than that for training data. It might be because of the “cold start problem”, the high dimensionality, and the sparseness of the case matrix.
- I will consider the combination with model-based methods such as non-negative matrix factorization to improve accuracy and enhance functions such as the prediction of symptoms and complications.