SOCIOCOM at the NTCIR-18 RadNLP Main Task: **Zero-Shot LLM Approaches for Lung Cancer Staging** Yuki Tashiro¹Yuta Nakamura² Eiji Aramaki¹

1 Nara Institute of Science and Technology, Japan 2 The University of Tokyo, Japan

Introduction [Nakamura+2025]

- Writing radiology reports is one of the most critical workflows for radiologists
- Although radiologists communicate with physicians by providing the results of imaging studies, radiology reports do not always specify the stage explicitly
- This imposes extra workload on human experts for careful manual information extraction, which can be aided by automation
- Task Objective: To classify the stage of lung cancer from radiology reports

Scope of RadNLP 2024

A tumor with a diameter of 12cm is observed spreading across the upper and lower lobes of the left lung, suggesting known lung cancer. It is in extensive contact with the pleura and is accompanied by the destruction of the left 3rd rib. Rib and parietal pleural infiltration is suspected. There are small nodules in the left upper lobe, suspecting secondary tumor nodules. The left mediastinal and bilateral hilar lymph nodes are enlarged, suspecting metastasis. No pleural effusion is observed. No obvious abnormalities are observed in the upper abdominal organs within the imaging range.







Lung cancer **CT** image

Dataset

- **Training Set**: Comprises 108 documents for 12 cases
- Validation Set: Includes 54 documents for 6 cases
- **Test Set**: Contains 216 documents for 24 cases



Our Method: Zero-shot LLM Approaches

- This document was inferenced by GPT-40 in each category, such as T, N, M
- As post-processing, if a predicted label is not in the allowed list, we replace it with the most frequent label for the category



Normal appearance of the visualized upper abdomen.

definition of label

- "N0": "No regional lymph node metastasis", "N1": "Metastasis in ipsilateral peribronchial …",
- "N2": "Metastasis in ipsilateral mediastinal and/or subcarinal lymph node(s)", "N3": "Metastasis in contralateral mediastinal, contralateral hilar ···"

Baseline

- Fine-tuned bi-encoder models: BERT, JMEDRoBERTa
- Data augmentation: Shuffled and rearranged sentences on a per-sentence basis

Results

- GPT-40 achieved the highest scores on the joint task
- GPT-40 is not necessary to post-processing
- Demonstrated performance that significantly improved bi-encoder models

Our model shows less overfitting than other teams' systems

Validation vs Test Joint Accuracy (JPN)

	Joint		Т		Ν		М	
Model Name	fine	coarse	fine	coarse	fine	coarse	fine	coarse
Baseline (most frequency label)	0.0	0.0	0.3333	0.3333	0.3704	0.3704	0.5	0.5
BERT	0.0741	0.2407	0.2593	0.4259	0.7963	0.7963	0.7037	0.7778
BERT (DA)	0.1852	0.3889	0.4074	0.5185	0.7222	0.7222	0.7407	0.8519
JMEDRoBERTa (DA, WordPiece)	0.0926	0.2037	0.4074	0.5741	0.5000	0.5000	0.5926	0.8519
IMEDRoBERTa (DA SentencePiece)	0 2037	0 4630	0 4074	0.6111	0 7222	0 7222	0.6111	0 8889



Conclusions

- Our results demonstrated that GPT-40 achieved higher accuracy compared to traditional bi-encoder models ullet
- Shared a simple, yet robust, GPT-4-mini baseline method with the community \bullet
- These findings highlight the potential of LLMs in medical natural language processing tasks

Nakamura et al. NTCIR-18 RadNLP 2024 Overview: Dataset and solutions for automated lung cancer staging. In Proceedings of the NTCIR-18 Conference, 2025