

# JGLUE: Japanese General Language Understanding Evaluation

Kentaro Kurihara<sup>1</sup>, Daisuke Kawahara<sup>1</sup>, Tomohide Shibata<sup>2</sup>

Waseda University<sup>1</sup>, Yahoo Japan Corporation<sup>2</sup>

## Overview

- Construct the first NLU benchmark in Japanese.
- Tasks are chosen to cover GLUE [Wang+ 18] and SuperGLUE [Wang+ 19] tasks.

Task	Dataset	Train	Valid	Test
Text Classification	MARC – ja	187,528	5,654	5,639
	JCoLA [Someya+ 22]	---	---	---
Sentence Pair Classification	JSTS	12,451	1,457	1,589
	JNLI	20,073	2,434	2,508
QA	JSQuAD	63,870	4,475	4,470
	JCommonsenseQA	9,012	1,119	1,118

## Introduction

- To develop high-performance NLU models, a benchmark is necessary.
- In the case of English, the **GLUE Benchmark** is publicly available.
- While benchmarks for languages other than English have been constructed, **there is no benchmark for Japanese**.
  - Chinese: CLUE [Xu+ 19], French: FLUE [Le+ 20], etc.

## Construction of Japanese Benchmark

- Due to the linguistic characteristics of Japanese, findings in other languages cannot necessarily be applied.
  - The Japanese alphabet includes *hiragana*, *katakana*, Chinese characters, and the Latin alphabet.
  - There are no spaces between words.

私は今朝パンを2枚食べました。  
(I had two pieces of bread this morning.)

- Problems with existing Japanese datasets

- Translation (e.g., JSNLI [Yoshikoshi+ 20])
  - Unnaturalness of Japanese in machine / manual translation
- Specific domains (e.g., JRTE [Hayashibe+ 20])
  - Not suitable for evaluating NLU ability in general domain.

From Scratch

In General Domain

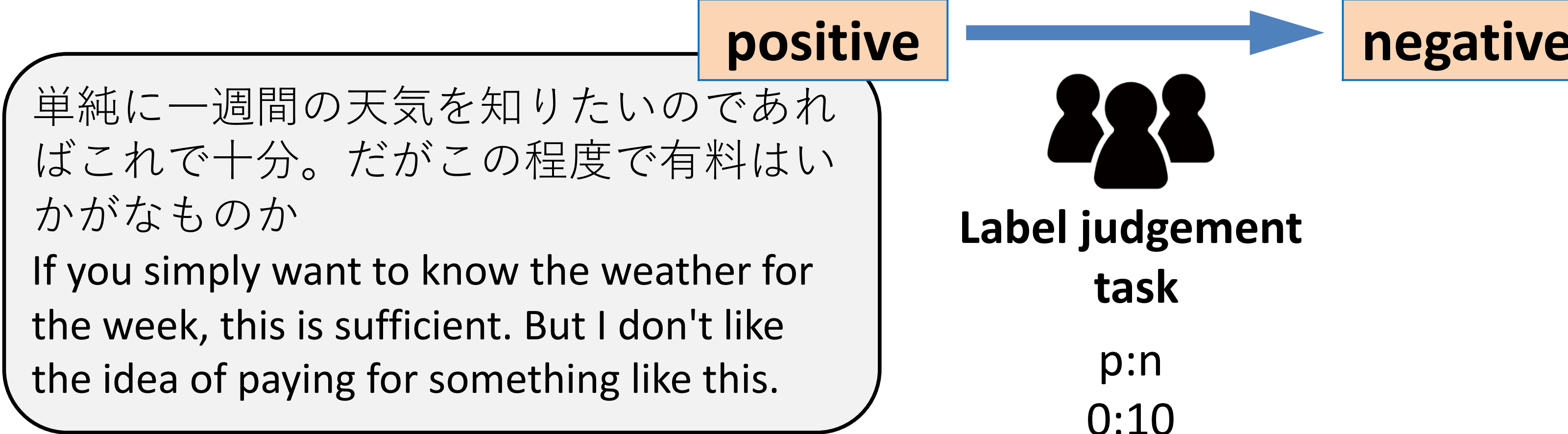
**Build a Japanese Language Understanding Benchmark and facilitate NLU research for Japanese.**

## Evaluation using JGLUE

- Evaluate the performance of pre-trained models using JGLUE.
- Findings
  - Overall, XLM-RoBERTa<sub>LARGE</sub> performed the best.
    - 1. LARGE model size, 2. CommonCrawl > Wikipedia
  - Tokenization: subword base > character base
  - Models pre-trained using CommonCrawl performed well.

## MARC – ja

- A Sentiment classification task (positive or negative) for product reviews.
- Build on the Japanese part of MARC (Multilingual Amazon Reviews Corpus) [Keung+ 20].
- Label modification through a crowdsourced label judgement task.



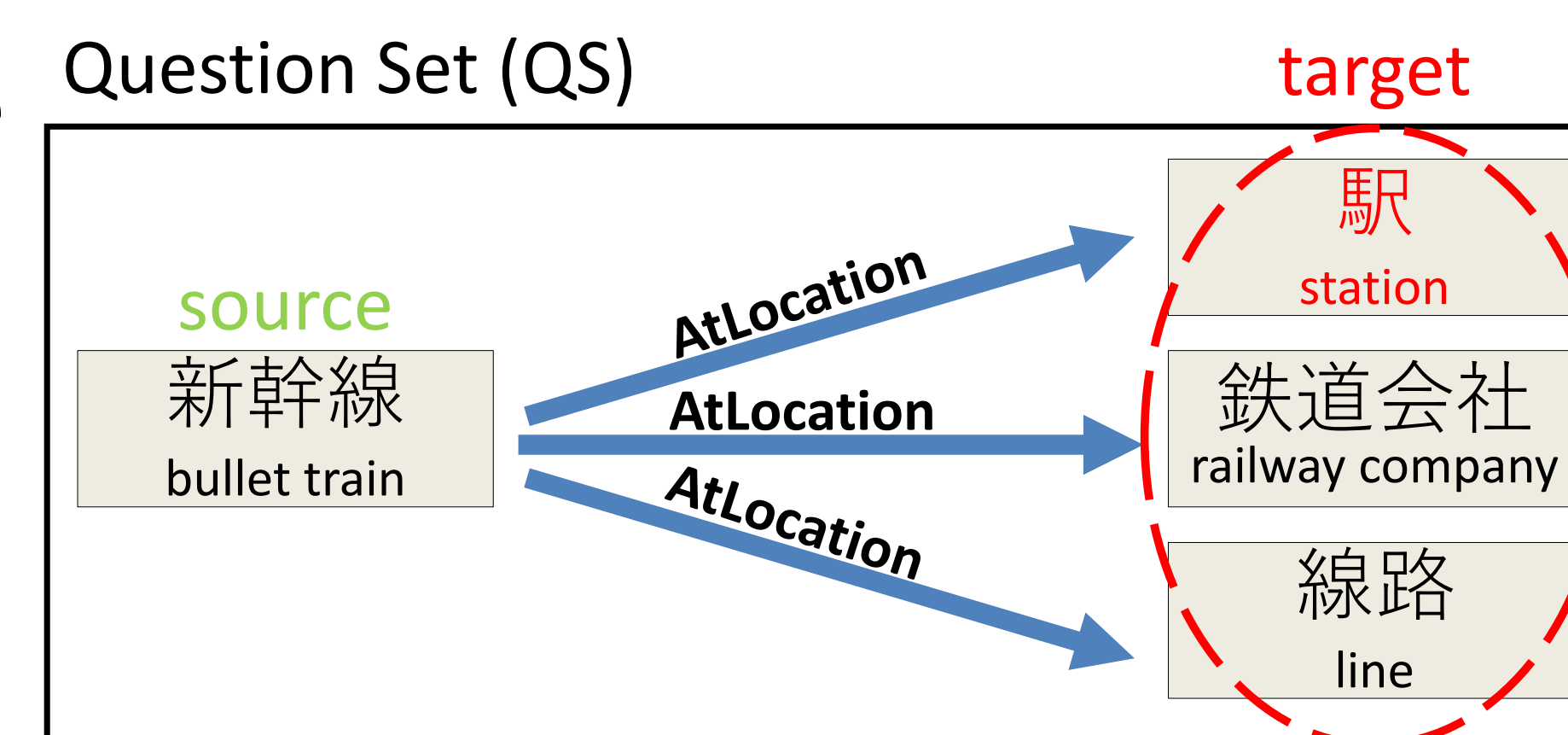
## JSQuAD

- A Japanese version of SQuAD [Rajpurkar+ 16], a general domain QA dataset constructed using Wikipedia.
- Our construction method is based on SQuAD.

## JCommonsenseQA

- Five choice QA to evaluate commonsense reasoning ability.
- Our construction method basically follows CommonsenseQA [Talmor+ 19].
- Remove Qs that contain synonyms in targets to make the answer unique.

庭 (garden) ≡ 庭園 (garden)



Question  
電車に人が乗り降りする場所を何という？  
What do you call a place where people get on and off the train?

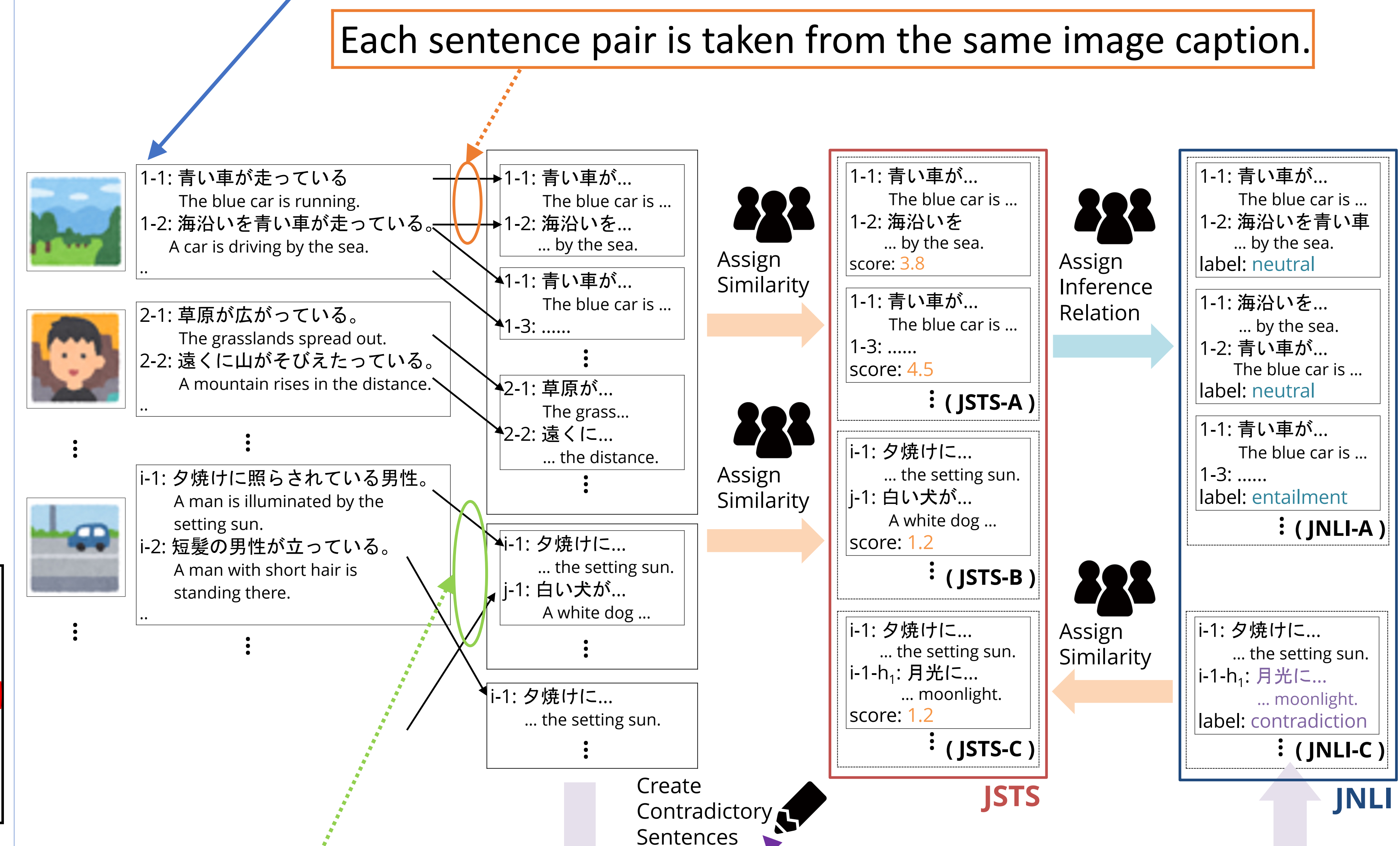
Choices  
駅, 鉄道会社, 線路, 空港, 港  
station, railway company, line, airport, port  
Add distractors

## Datasets in JGLUE

### JSTS / JNLI

- JSTS**
  - A task of predicting the semantic similarity between two sentences with a value between 0 (completely different in meaning) and 5 (equivalent in meaning).
- JNLI**
  - A task of recognizing the inference relation that a premise sentence has to a hypothesis sentence based on 3 labels (“entailment”, “neutral”, and “contradiction”).

We basically extract sentence pairs in JSTS and JNLI from YJ Captions Dataset [Miyazaki+ 16].



There are few non-similar sentence pairs collected from captions of the same image. → Collect sentence pairs from different image captions.

There are few contradictory sentence pairs from captions of the same image. → A worker creates contradictory sentences given a caption.

## Conclusion / Future work

- Constructed JGLUE, the first NLU benchmark in Japanese.
- JGLUE is available at <https://randd.yahoo.co.jp/en/softwaredata#jglue>
- Plan to build Japanese benchmarks for generation tasks such as GLGE [Liu+ 21].



	MARC-ja	JSTS	JNLI	JSQuAD	JComQA
Models	acc	Pearson	acc	F1	acc
Human	0.990	0.909	0.917	0.946	0.988
Tohoku BERT <sub>BASE</sub>	0.957	0.901	0.876	0.946	0.782
Tohoku BERT <sub>BASE</sub> (char)	0.957	0.889	0.861	0.937	0.728
Tohoku BERT <sub>LARGE</sub>	0.961	0.907	0.878	0.950	0.822
NICT BERT <sub>BASE</sub>	0.960	0.909	0.881	<b>0.952</b>	0.807
Waseda RoBERTa <sub>BASE</sub>	0.962	0.901	0.876	0.926	<b>0.849</b>
XLM-RoBERTa <sub>LARGE</sub>	<b>0.965</b>	<b>0.916</b>	<b>0.902</b>	—	0.842