

Evaluation Dataset and System for Japanese Lexical Simplification

kajiwara@jnlp.org

http://www.jnlp.org/SNOW

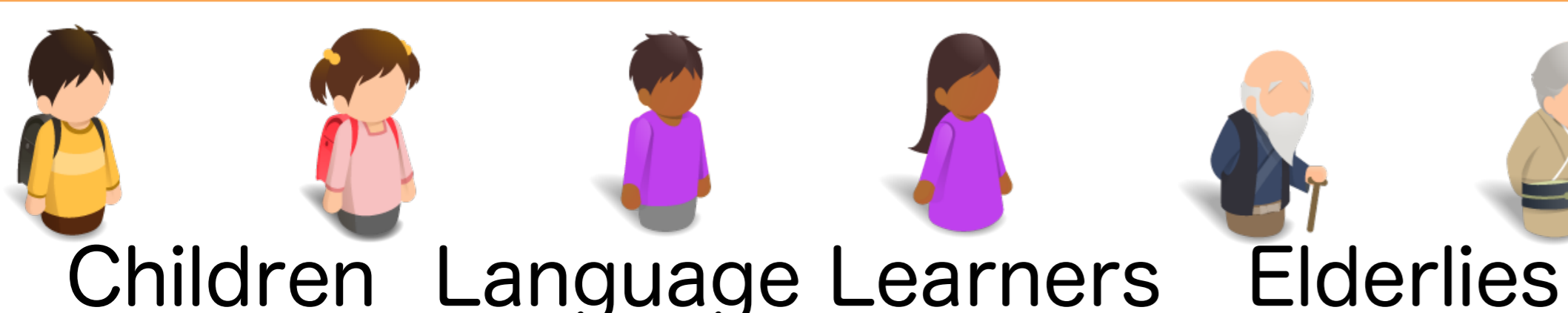
Tomoyuki Kajiwara and Kazuhide Yamamoto

Nagaoka University of Technology, Japan

Motivation

Extensive / various forms of texts

Hitler committed terrible **atrocities** during the second World War.
Hitler committed terrible **cruelties** during the second World War.



Task: Lexical Simplification

- Substitutes a complex word or phrase in a sentence with a simpler synonym
- Supports the reading comprehension of a wide range of readers

Problems in Japanese

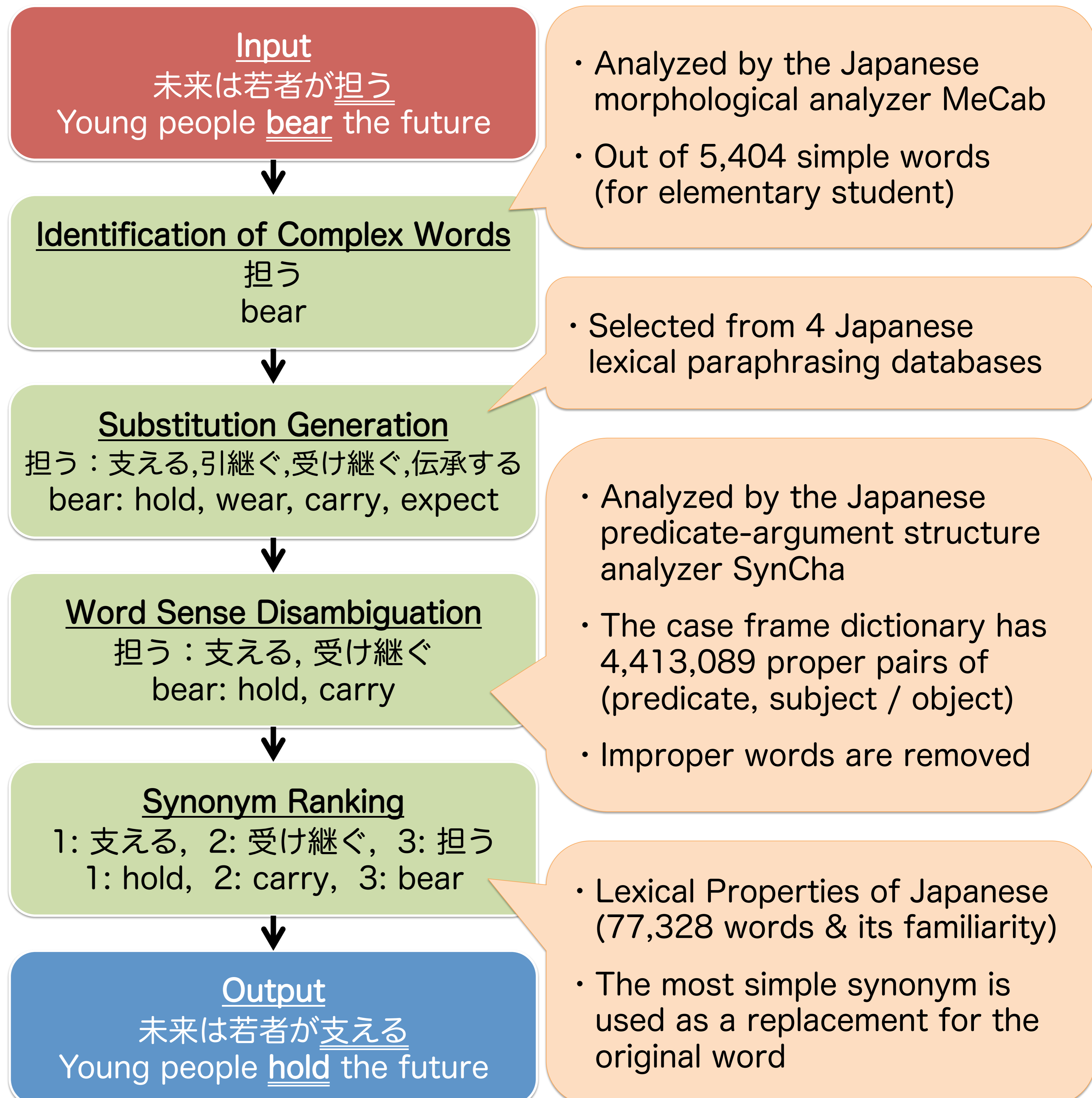
➤ Unpublished system

It is difficult for people who need reading assistance to obtain simple Japanese sentences

➤ Unpublished dataset

It is difficult for researchers and developers to evaluate the performance of different systems

Lexical Simplification System



Evaluation Dataset

1. Constructing Japanese Lexical **Substitution** Dataset

- Collecting Substitutions (crowdsourcing: 5 workers, 17.8%)
- Evaluating Substitutions (crowdsourcing: 5 workers, 66.4%)

2. Transforming it into Lexical **Simplification** Dataset

- Ranking Substitutions (crowdsourcing: 5 workers, 33.2%)
- Merging All Rankings

- This dataset consists of 2,330 sentences, 233 target words each with 10 sentences as contexts
- These contexts were randomly selected from Japanese newspaper articles

Sample: Young people bear the future.

- **Lexical Substitutions:** carry, hold
- **Rank of Simple Level:** 1.hold, 2.carry, 3.bear

The gold-standard annotations were generated by averaging the annotations from all annotators

Dataset	Sentence	Noun	Verb	Adjective	Adverb
SemEval 2012 Task1	2,010	580 (28.9%)	520 (25.9%)	560 (27.9%)	350 (17.4%)
Ours	2,330	630 (27.0%)	720 (30.9%)	500 (21.5%)	480 (20.6%)

System	Precision	Recall	F-measure
Our Original	0.89	0.08	0.15
w/o WSD	0.84	0.71	0.77

Properties of the dataset	
The average number of substitutions	4.50
The average number of levels of difficulty	4.94
There are synonyms that are more simpler	69.4 %
There are synonyms that are more complex	83.5 %

Context dependency ratio		
①: context pairs	10,485	100 %
②: ① with same list	1,593	15 %
③: ② with different rankings	948	60 %
④: ③ with different top word	463	49 %