

Overview

- **Objective:** Create a Japanese morphological analyzer (word segmentation + POS tagging) that is robust and adaptable to new domains
- **Approach:** Use pointwise prediction, which estimates all tags independently of other tags
- **Pointwise prediction:**
 - *Robust:* does not rely on dictionaries as much as previous methods
 - *Adaptable:* it can be learned from single annotated words, not full sentences
 - *Works with active learning:* Single words to annotate can be chosen effectively
- **Evaluation** on Japanese morphological analysis shows improvement over traditional methods

1

Features for Pointwise MA

- Specify features using character n-grams, character type n-grams, length-annotated dictionary presence

Boundary Point

	<i>hon</i>	<i>zai</i>	<i>wo</i>	<i>tou</i>	<i>yo</i>	<i>su</i>	<i>ru</i>
	本	剤	を	投	与	す	る
WS	Char 1-gram			X0 投	X1 与		
	Char 2-gram	X-1 を投	X0 投与	X1 与す			
	Char 3-gram	X-1 を投与	X0 投与す				
	Type 1-gram		T0K	T1K			
	Type 2-gram		T-1HK	T0KK	T1KH		
	Type 3-gram	T-2KHK	T-1HKK	T0KKH	T1KHH		
	Dictionary	D0L1(投)	D0R1(与)	D0I2(投与)			

POS Char n-gram + Type n-gram
Word Identity W 投与 + Dictionary DN DV

- **Key point:** None of the features require word boundaries or surrounding tags

3

Experiments

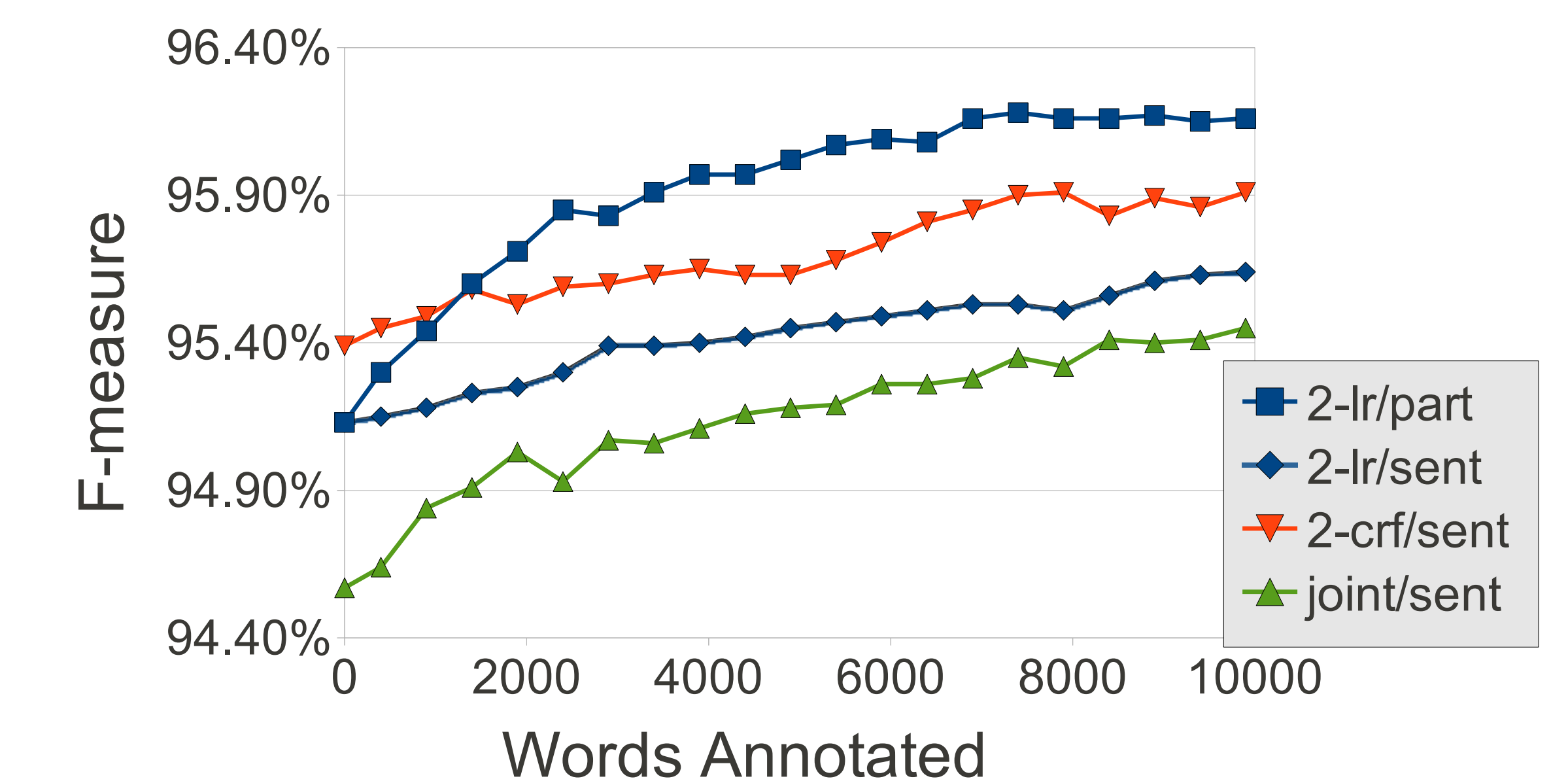
- Experiments performed on the Balanced Corpus of Contemporary Written Japanese (BCCWJ)
- **General domain:** News, white papers, books
- **Target domain:** Web text

	Train	Test
General	782k	87.5k
Target	153k	17.3k

- Tested three systems
- **Joint:** Kudo et al.'s CRF-based method, as implemented by the MeCab toolkit
- **2-CRF:** The 2 step method using CRFs as a solver
- **2-LR:** 2 step pointwise method using LR

	Train	Test	Joint	2-CRF	2-LR
General	General	General	97.31%	98.08%	98.03%
General	General	Target	94.57%	95.39%	95.13%
Gen+Tar	General	Target	96.45%	96.91%	96.82%

- 2-LR slightly worse than 2-CRF, better than Joint
- Tested 2 annotation strategies for domain adaptation
- **Sentence:** annotate the sentence with lowest overall posterior probability
- **Partial:** annotate the word with lowest prob. margin

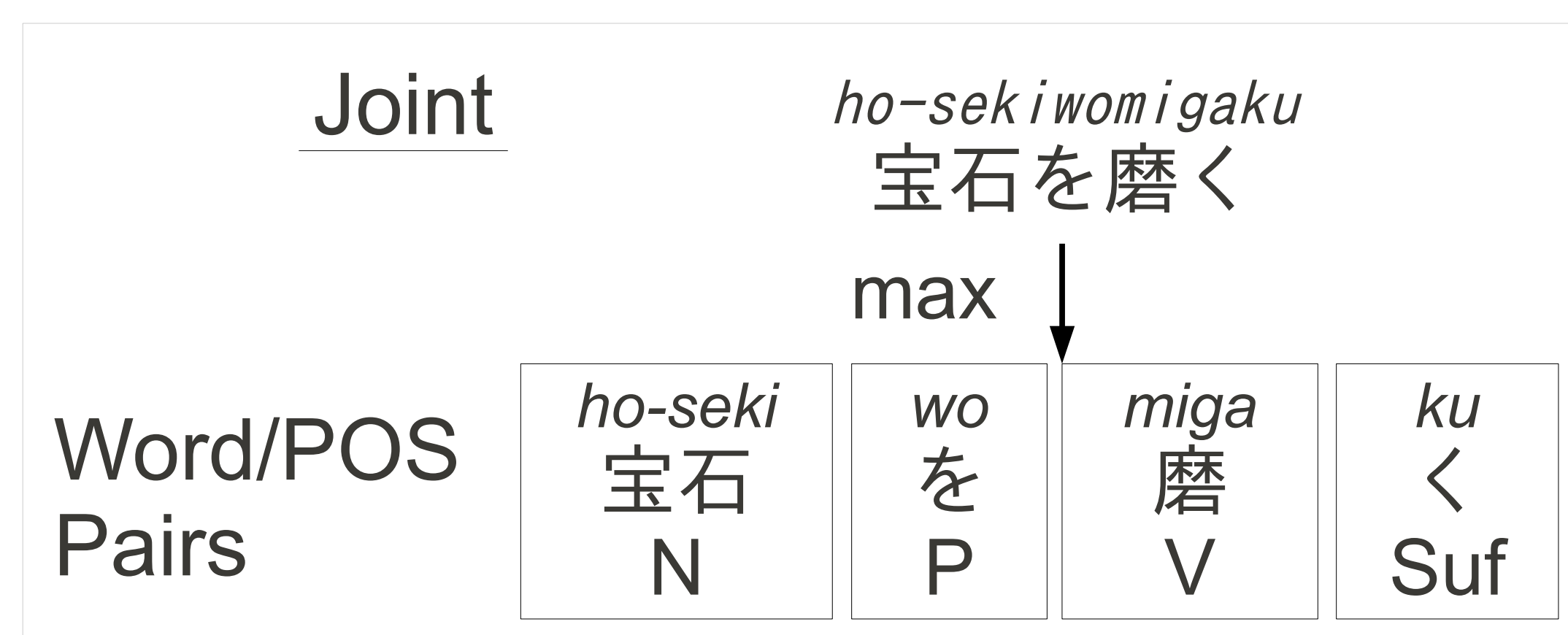


- Partial annotation much more effective

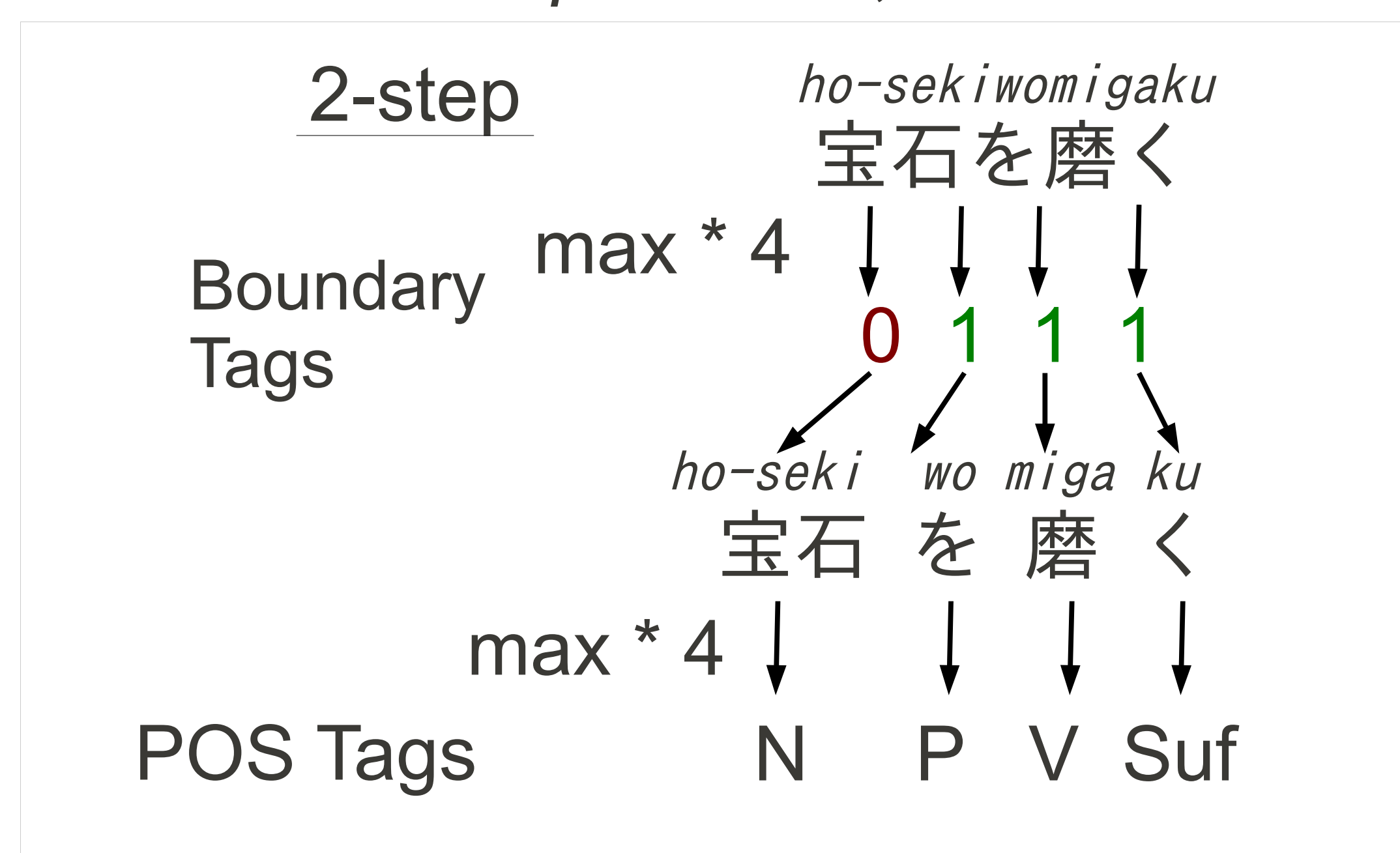
5

Morphological Analysis Methods

- **Joint:** Predict word boundaries+tags simultaneously
- Use HMMs, CRFs, or language models



- **2-Step:** First predict word boundaries, then POSs
- Can use Logistic Regression, SVM, CRF
- LR and SVM are *pointwise*, CRF not



2

Annotation Methods

- Morphological analysis underperforms on out-of-domain text → we would like to adapt
- We have an in-domain unannotated text, and some annotator time
- Goal is to maximize the effect for annotator time
- Use active learning to choose data to annotate

Reference	本	剤	を	投	与	す	る	た	め
Automatic Result	本	Pre	剤	を	投	与	す	る	た
	0.8	0.91	1.0	0.98	0.94	0.998	0.997		

- **Full annotation:** Choose sentences with low prob.
- Can train any model on this annotated data

本 剤 /N を /P 投 与 /N す る /V た め /N

Annotated (5)

- **Partial annotation:** Choose words with low prob.
- Only pointwise prediction can be used

本 剤 /N を 投 与 /N す る た め

Annotated (2) Unannotated

4

Available Open Source!

<http://www.phontron.com/kytea/>

Chinese models, Japanese pronunciation estimation also available