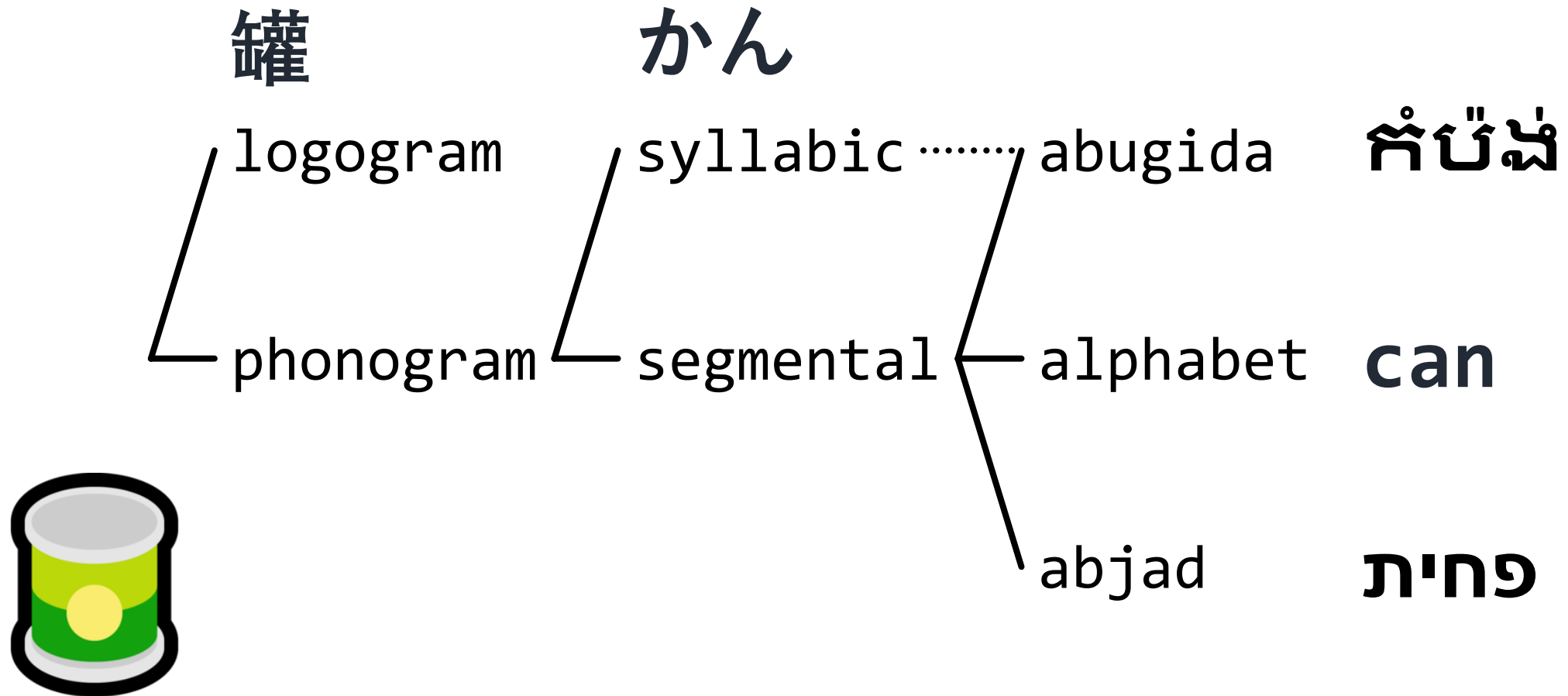


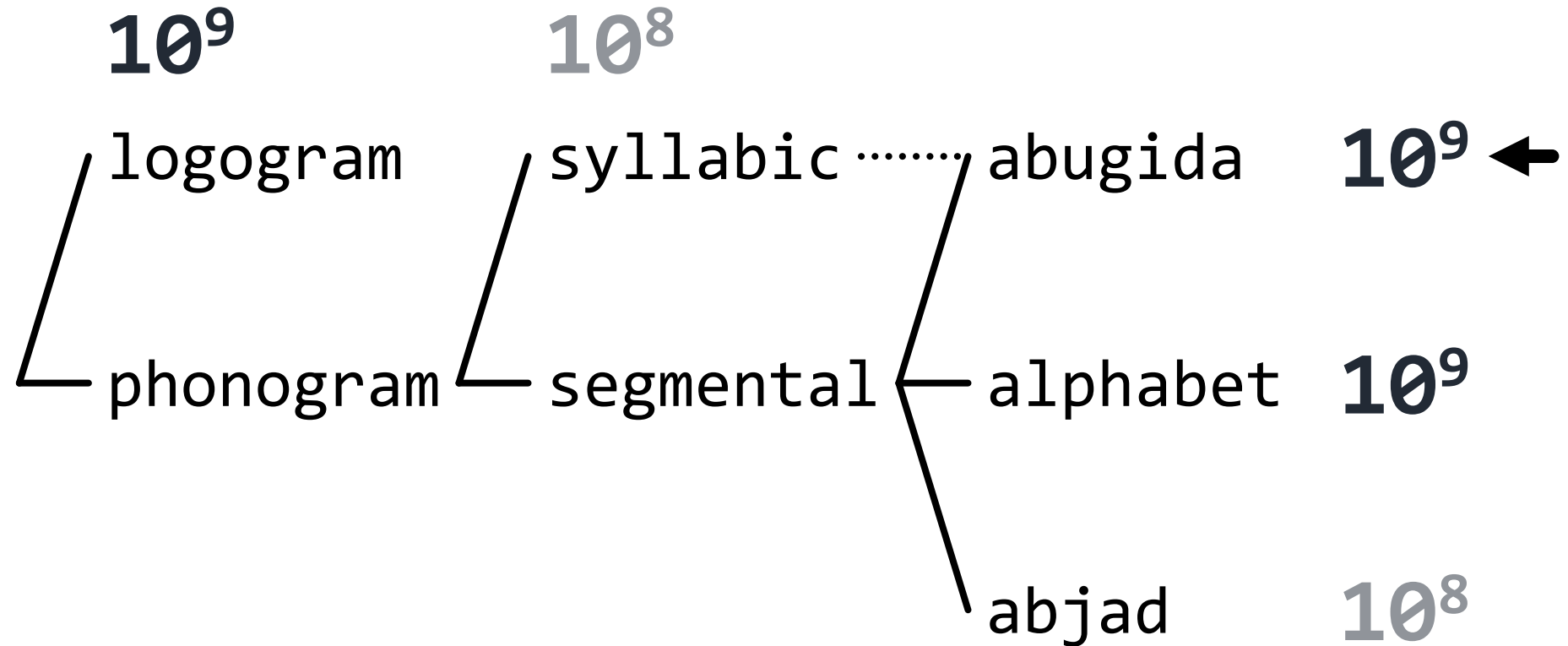
# Simplified Abugidas

Chenchen Ding, Masao Utiyama, and Eiichiro Sumita  
ASTREC, NICT, Japan

# Writing Systems : Hierarchy



# Writing Systems : Population



# Writing Systems : How to Input

Logogram > Syllabic > Abugida > Alphabet & Abjad

**Keyboard**



**Input method**



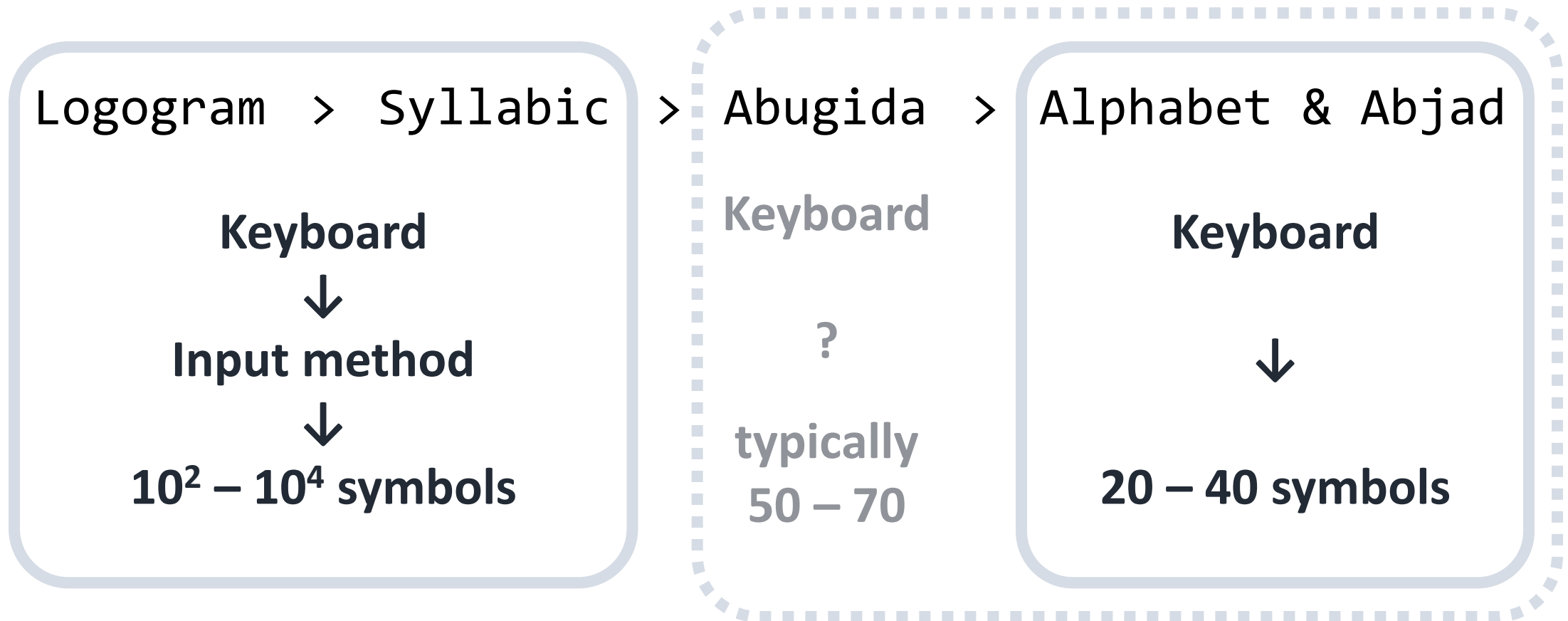
**$10^2 - 10^4$  symbols**

**Keyboard**



**20 – 40 symbols**

# Writing Systems : How to Input

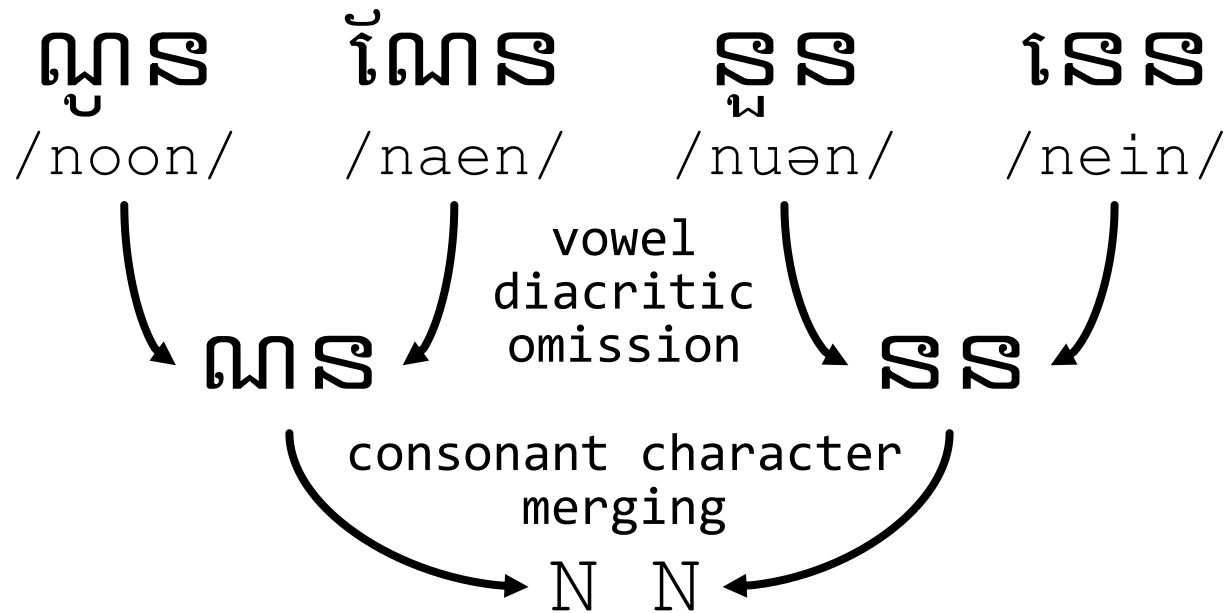


# Motivation of This Study

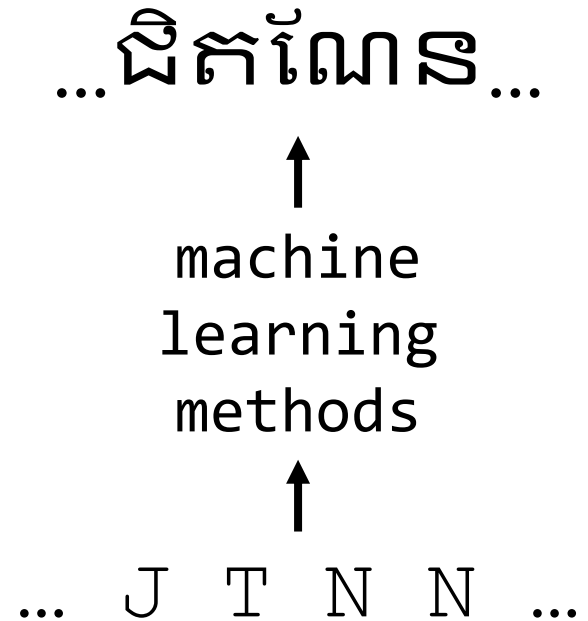
- Can abugidas be inputted more efficiently?
  - To insert a light layer of input method
  - To type less and to recover automatically
- Related work
  - Various approaches for Chinese and Japanese
  - To take advantage of redundancy in a writing system

# Outline of This Study

- Khmer script as an example



(a) ABUGIDA SIMPLIFICATION



(b) RECOVERY





# Abugida Simplification

- Khmer script as an example

ជិតណេន = ជ + ិ + ត + ណ + ៃ + ន → J T N N

$$\text{Leng.} = \frac{\text{len}(\text{"J T N N"})}{\text{len}(\text{"ជ ិ ត ណ ៃ ន"})} = \frac{4}{6} = 66.7\%$$

	Thai	Burmese	Khmer	Lao
Leng.	76.0%	74.0%	77.6%	72.8%

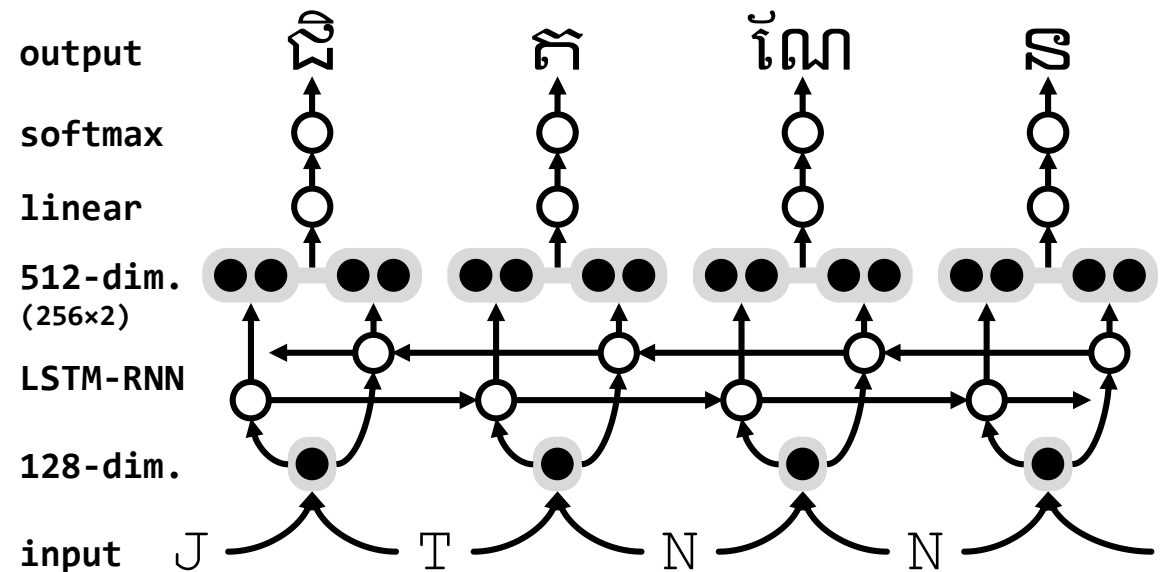
- Around **one quarter** characters (**diacritics**) saved

# Recovery Methods

- To formulate as a sequential labeling task
  - However, list-wise search as in **conditional random fields** is costing
- To solve by point-wise classification
  - **Support vector machine (SVM)** as a baseline
  - **Recurrent neural network (RNN)** as a state-of-the-art method
- Setting for the SVM baseline
  - Linear kernel with N-gram features
  - Using **LIBLINEAR** library
    - Wrapped by the **KyTea** toolkit

# RNN Structure and Settings

- Bi-gram of graphemes as input
  - **Embedding** → **Bi-directional LSTM** → Linear transform → Softmax
  - Original writing units as output
- 
- Implemented by **DyNet**
  - Trained by **Adam**
    - Initial learning rate 0.001
    - Controlled by a validation set
  - Multi-model ensemble



# Experimental Results

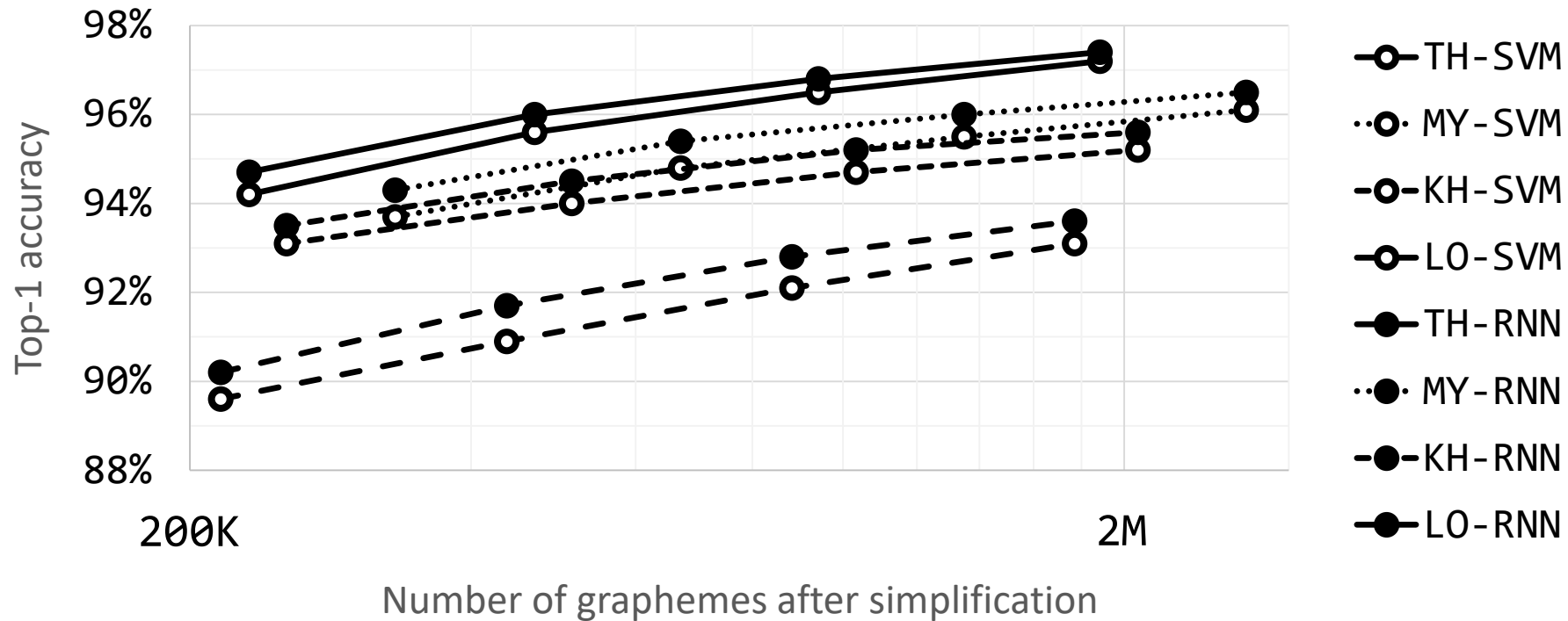
- Asian Lang. Treebank data
  - 20,000 sent., newswire
- **SVM**
  - Up to 5-gram for TH, KM, LO
  - 7-gram for Burmese
- **RNN**
  - $\oplus M$  : M-ensemble  
→ 16-ensemble is adequate
  - @N : Top-N results  
→ Top-4 is satisfactory

	Thai	Burmese	Khmer	Lao
SVM	97.2%	96.1%	95.2%	93.1%
RNN $_{\oplus 4}^{\textcircled{1}}$	97.2%	96.3% <sup>‡</sup>	95.5% <sup>‡</sup>	93.3% <sup>‡</sup>
RNN $_{\oplus 8}^{\textcircled{1}}$	97.3% <sup>†</sup>	96.4% <sup>‡</sup>	<b>95.6%<sup>‡</sup></b>	<b>93.6%<sup>‡</sup></b>
RNN $_{\oplus 16}^{\textcircled{1}}$	<b>97.4%<sup>‡</sup></b>	<b>96.5%<sup>‡</sup></b>	<b>95.6%<sup>‡</sup></b>	<b>93.6%<sup>‡</sup></b>
RNN $_{\oplus 16}^{\textcircled{2}}$	98.8%	98.4%	97.5%	96.6%
RNN $_{\oplus 16}^{\textcircled{4}}$	99.2%	98.8%	98.1%	97.7%
RNN $_{\oplus 16}^{\textcircled{8}}$	99.2%	98.9%	98.4%	97.9%

*† : p < 0.01, ‡ : p < 0.001*

→ **Embedding + bi-LSTM > N-gram features**

# Experimental Results : Training Data Size



→ RNN outperforms SVM, regardless of the training data size

# Manual Evaluation

- On Burmese and Khmer best results by RNN
- Conducted by native-speakers
- To classify errors into four-level
  - 0. acceptable, i.e., alternative spelling
  - 1. clear and easy to identify the correct result
  - 2. confusing but possible to identify the correct result
  - 3. incomprehensible

Level	0	1	2	3
Burmese	4.5%	51.0%	42.2%	2.2%
Khmer	22.5%	28.5%	16.3%	32.8%

# Conclusion and Future Work

- Abugidas can be simplified largely and recovered with high accuracy
  - Four Brahmic abugidas are investigated
  - Simplified into a compact symbol set (around 20 graphemes)
  - Recovered satisfactorily by standard machine learning method
  - Experimentally show the feasibility to encode abugidas in a lossy manner
- Future work
  - Language specific investigation
  - To integrate dictionary
  - To develop practical input method for abugidas

Thanks for your kind attention