# Analyzing the Surprising Variability in Word Embedding Stability Across Languages
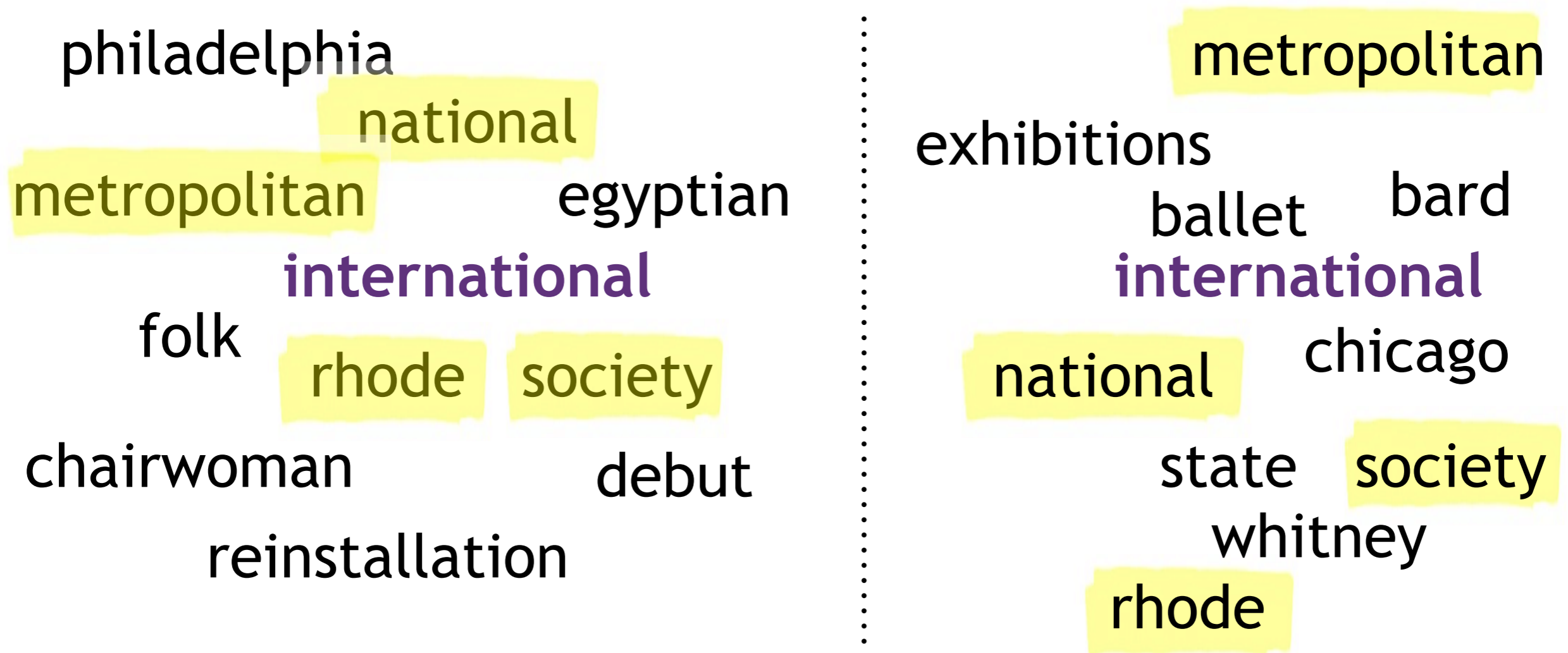
**Laura Burdick**, Jonathan K. Kummerfeld, Rada Mihalcea
*University of Michigan*

# What is Stability?

**Stability** = *percent overlap between ten nearest neighbors in an embedding space*

philadelphia
national
metropolitan    egyptian
**international**
folk
rhode    society
chairwoman    debut
reinstallation

metropolitan
exhibitions
ballet    bard
**international**
national    chicago
state    society
whitney
rhode

**Stability = 40%**

# This Work

*Does stability vary for different languages?*

*Is stability associated with linguistic properties?*

▸ **Data**
  ▸ **Wikipedia** (40 languages)

  ▸ **Bible** (97 languages)

  ▸ **World Atlas of Language Structures (WALS),** phonological, lexical, and grammatical properties (>2,000 languages)
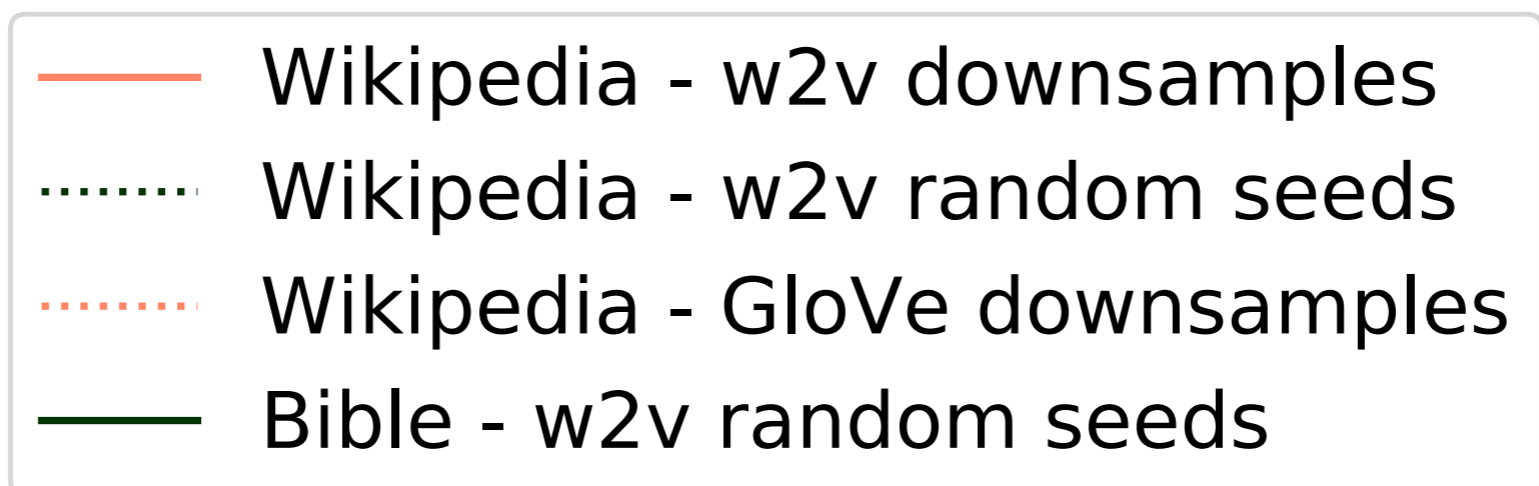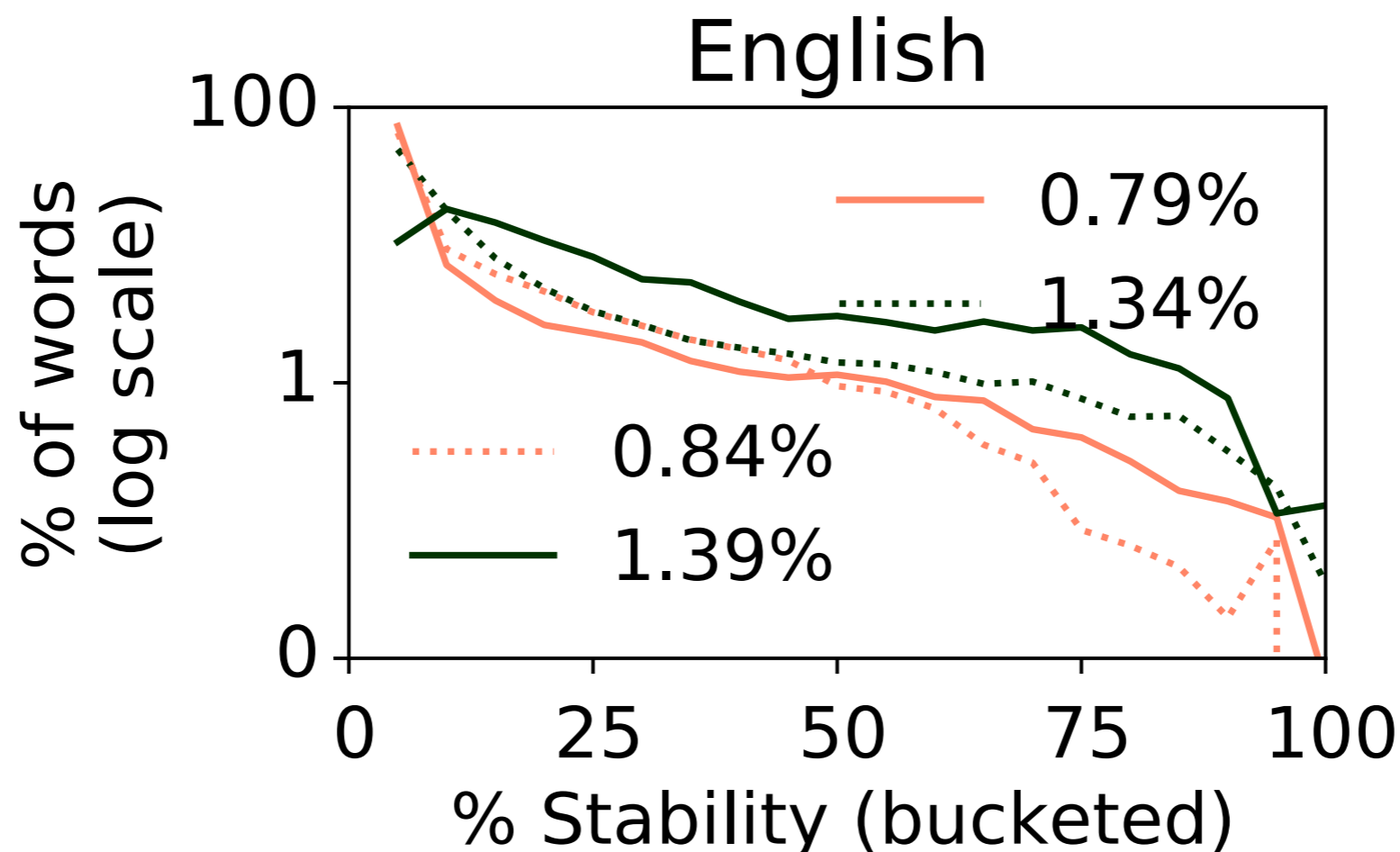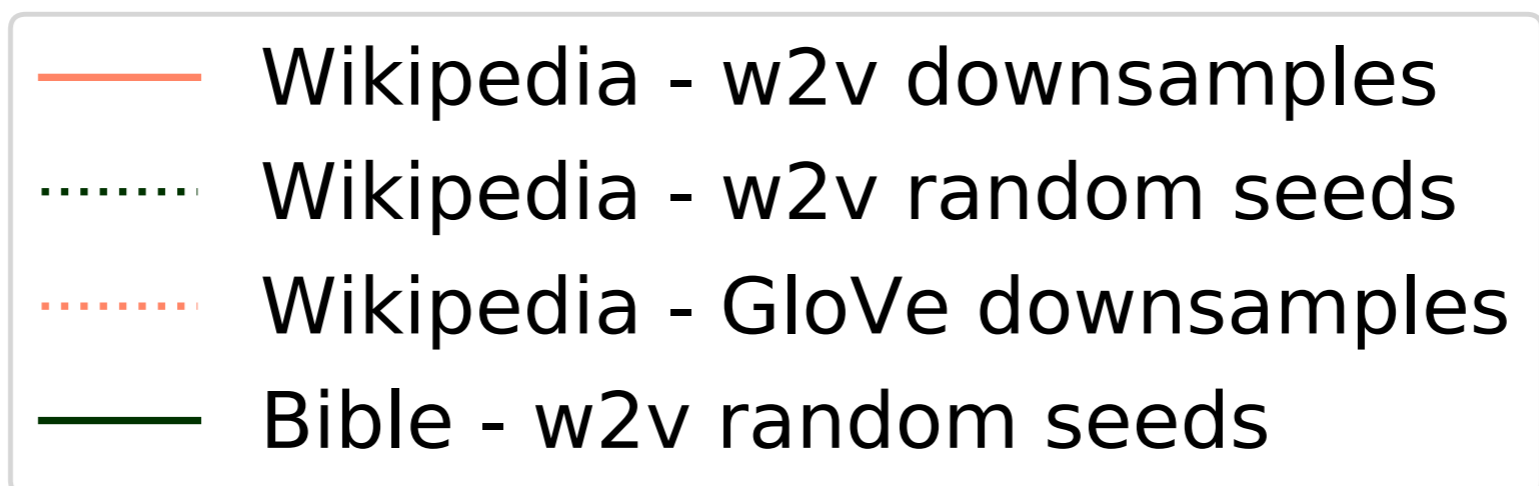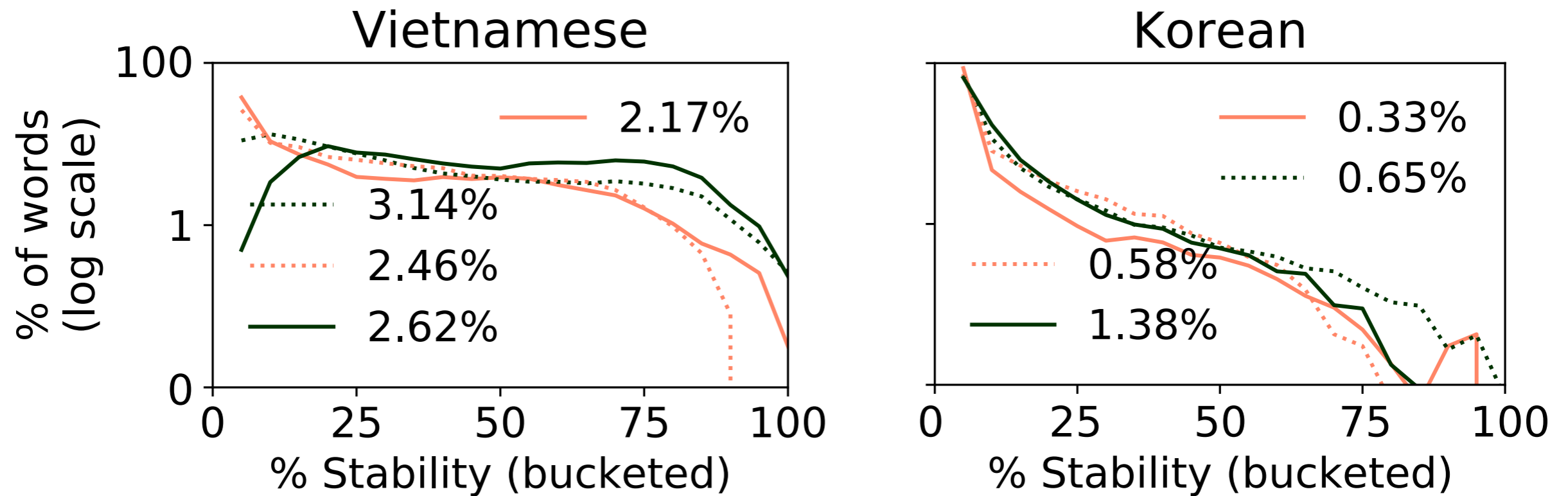
# Stability for Wikipedia and the Bible

*We compare the stability of embeddings for 26 languages.*

- **Wikipedia (3 settings):** Stability of...
  - GloVe embeddings across 5 downsampled corpora
  - word2vec (w2v) embeddings across 5 downsampled corpora
  - w2v using 5 random seeds on 1 downsampled corpus

- **One setting for the Bible:** Stability of w2v embeddings using 5 random seeds on 1 downsampled corpus

- Each downsampled corpora 100,000 sentences

# Stability for Wikipedia and the Bible



English

% of words (log scale)

100

1

0

0    25    50    75    100

% Stability (bucketed)

0.79%

1.34%

0.84%

1.39%

—— Wikipedia - w2v downsamples
······ Wikipedia - w2v random seeds
······ Wikipedia - GloVe downsamples
—— Bible - w2v random seeds

# Stability for Wikipedia and the Bible



**Vietnamese**

% of words (log scale)

100
1
0

2.17%
3.14%
2.46%
2.62%

% Stability (bucketed)

**Korean**

0.33%
0.65%
0.58%
1.38%

% Stability (bucketed)

— Wikipedia - w2v downsamples
····· Wikipedia - w2v random seeds
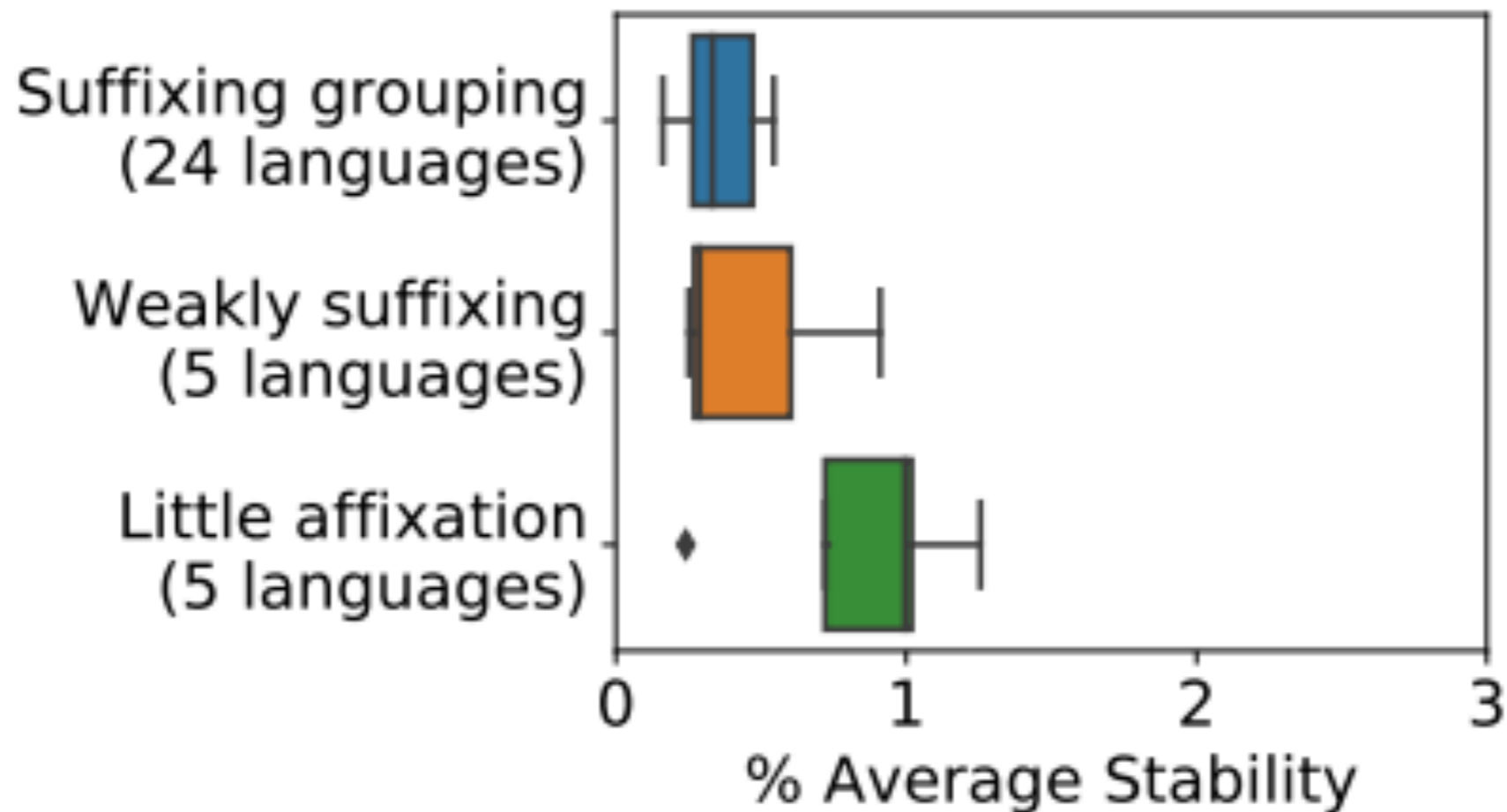····· Wikipedia - GloVe downsamples
— Bible - w2v random seeds

# Regression Modeling

*We use a regression model to predict stability in a language using linguistic properties.*

- Ridge regression

- 37 languages

- Input: 97 WALS properties
- Output: Average stability of all the words in a language

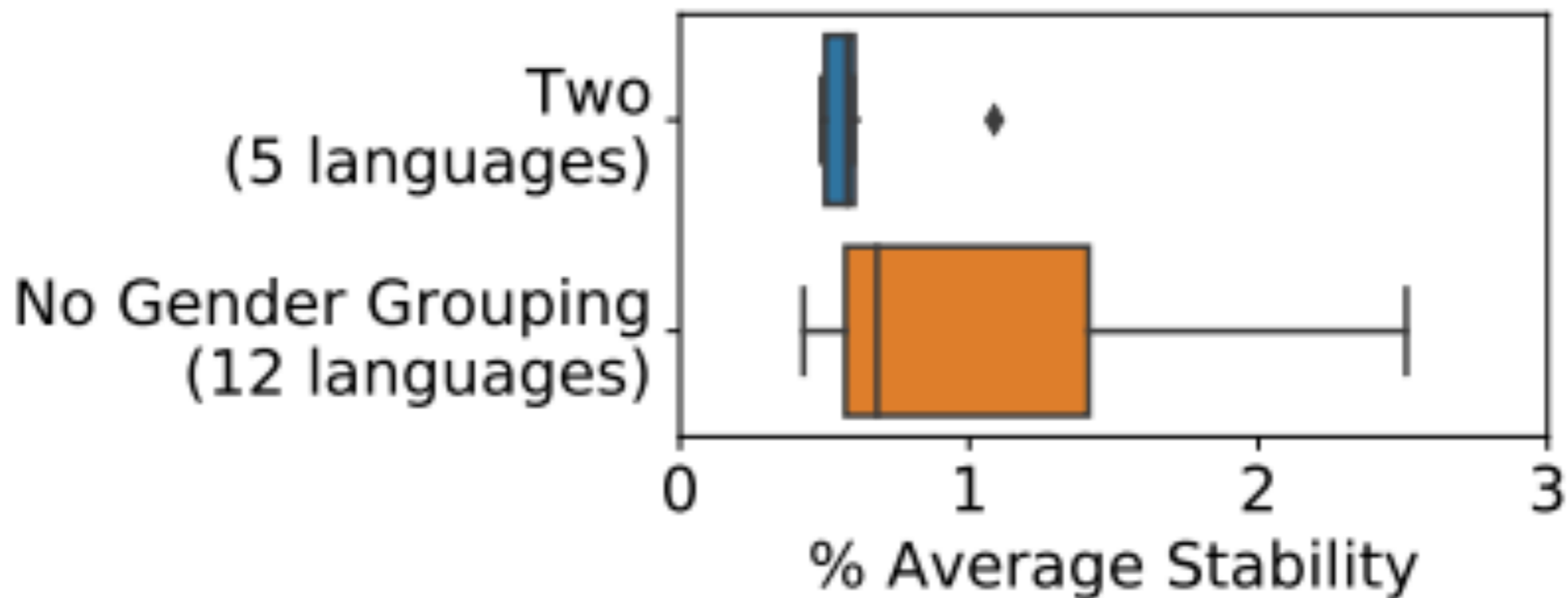- High $R^2$ score of 0.96 ± 0.00

# Regression Modeling

- More affixing (suffixing and prefixing) associated with lower stability
  - Affixes cause increased word variation



*Prefixing v. Suffixing in Inflectional Languages*

# Regression Modeling

▸ Languages with no gender system associated with higher stability

  ▸ Languages with gender systems have more word forms



*Number of Genders*

# Final Thoughts

- *Languages with more affixing tend to have less stable embeddings*

- *Languages with no gender systems tend to have more stable embeddings*

- *Future embedding design needs to take into account gendered words and morphologically rich words with affixes*

*Download our code:*

`http://lit.eecs.umich.edu/downloads.html`