# Self-supervised speech models contain linear morpho-phonological structures
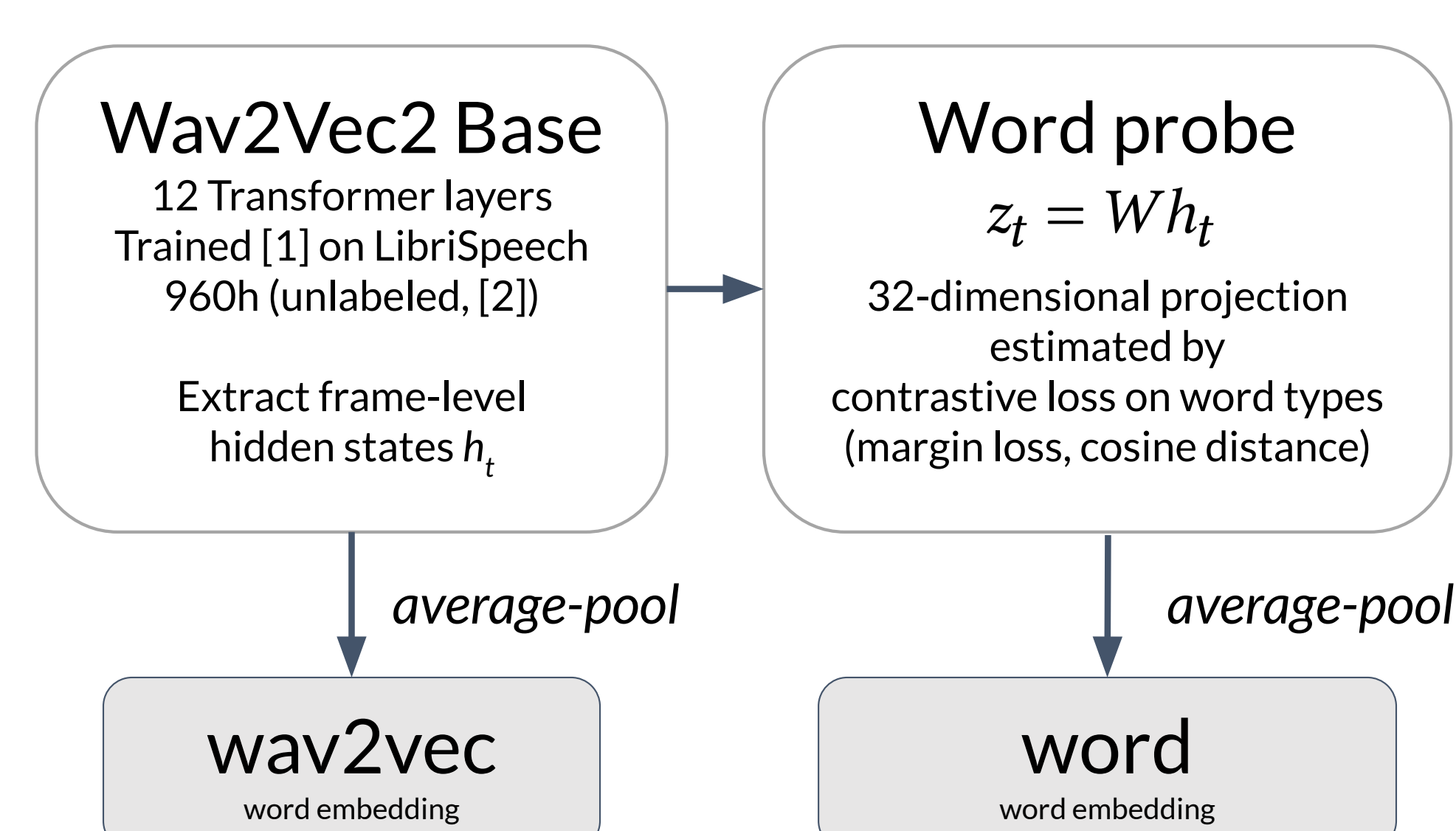
What linguistic representations are necessary to **recognize spoken words**? We study acoustic word embeddings of two models:

Wav2Vec2 Base
12 Transformer layers
Trained [1] on LibriSpeech
960h (unlabeled, [2])
Extract frame-level
hidden states $h_t$

*average-pool*

wav2vec
word embedding

Word probe
$z_t = W h_t$
32-dimensional projection
estimated by
contrastive loss on word types
(margin loss, cosine distance)

*average-pool*

word
word embedding

## APPROACH

Study model encodings of word-final /z/, /s/, and /ɪz/, which are governed by overlapping morphological and phonological processes:

| | Allomorph | Base | Inflected |
|---|---|---|---|
| **Noun plural** | /z/ | daughter | daughter<u>s</u> |
| | /s/ | lip | lip<u>s</u> |
| | /ɪz/ | age | ag<u>es</u> |
| **Verb 3rd-person singular** | /z/ | bring | bring<u>s</u> |
| | /s/ | speak | speak<u>s</u> |
| | /ɪz/ | please | pleas<u>es</u> |

Does a regular geometry link these word pairs in acoustic word embedding space? Is this geometry sensitive to **phonological** or **morphological** distinctions?
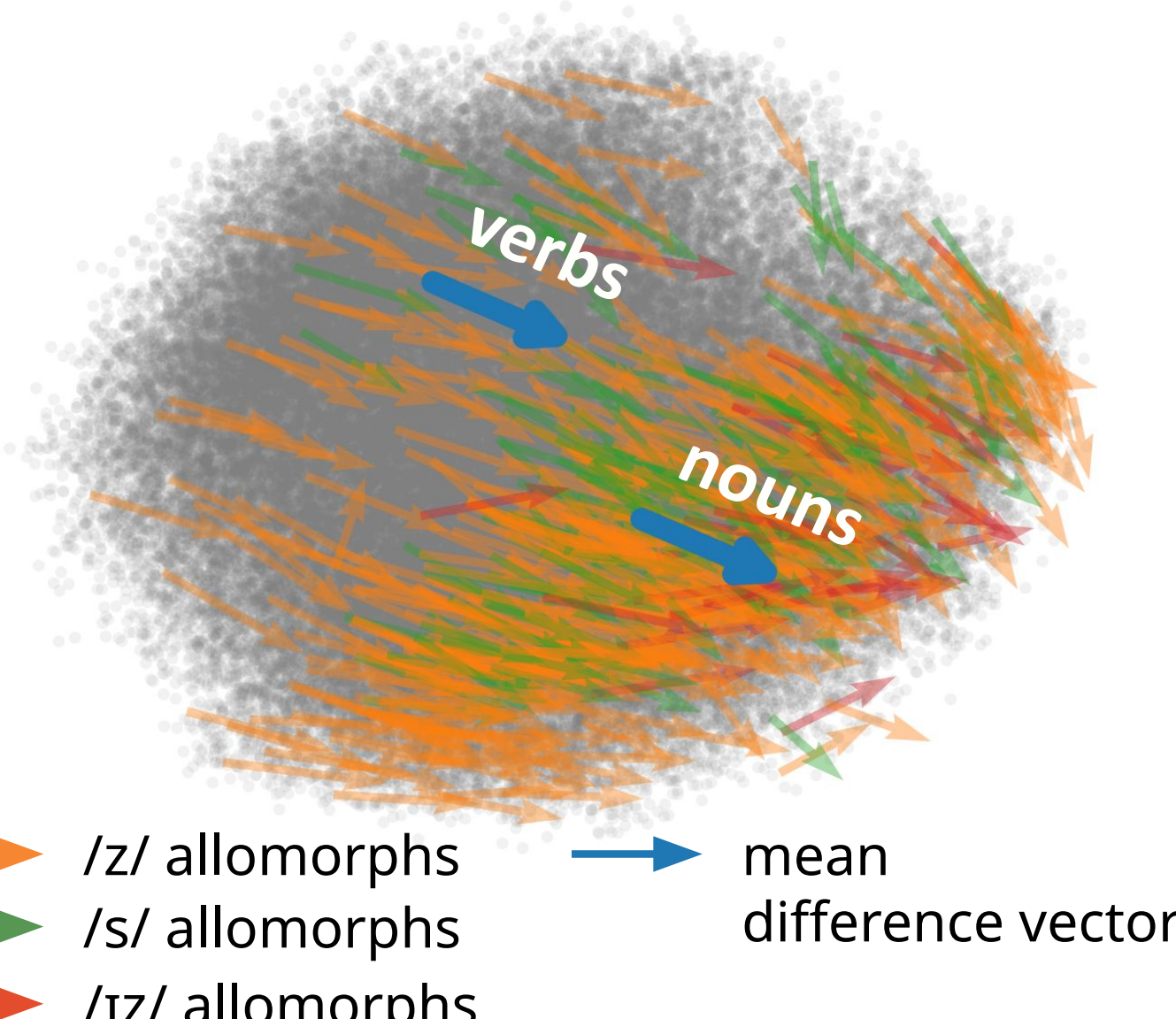
daughters : daughter :: speaks : speak

$$f(\text{daughters}) - f(\text{daughter}) \approx f(\text{speaks}) - f(\text{speak})$$

Predict analogy solutions from token-level
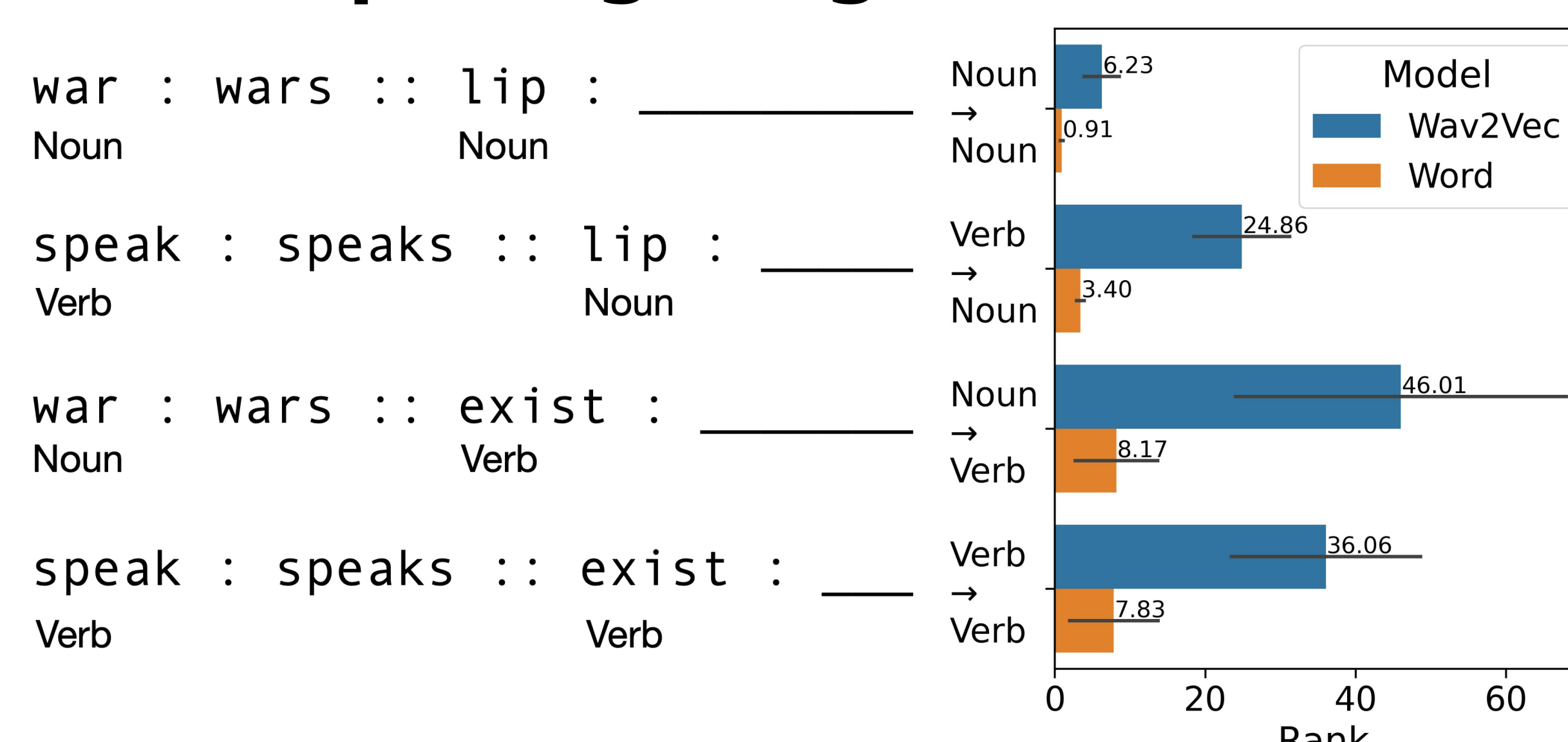em| $\hat{y} = f(\text{daughters}) - f(\text{daughter}) + f(\text{speak})$

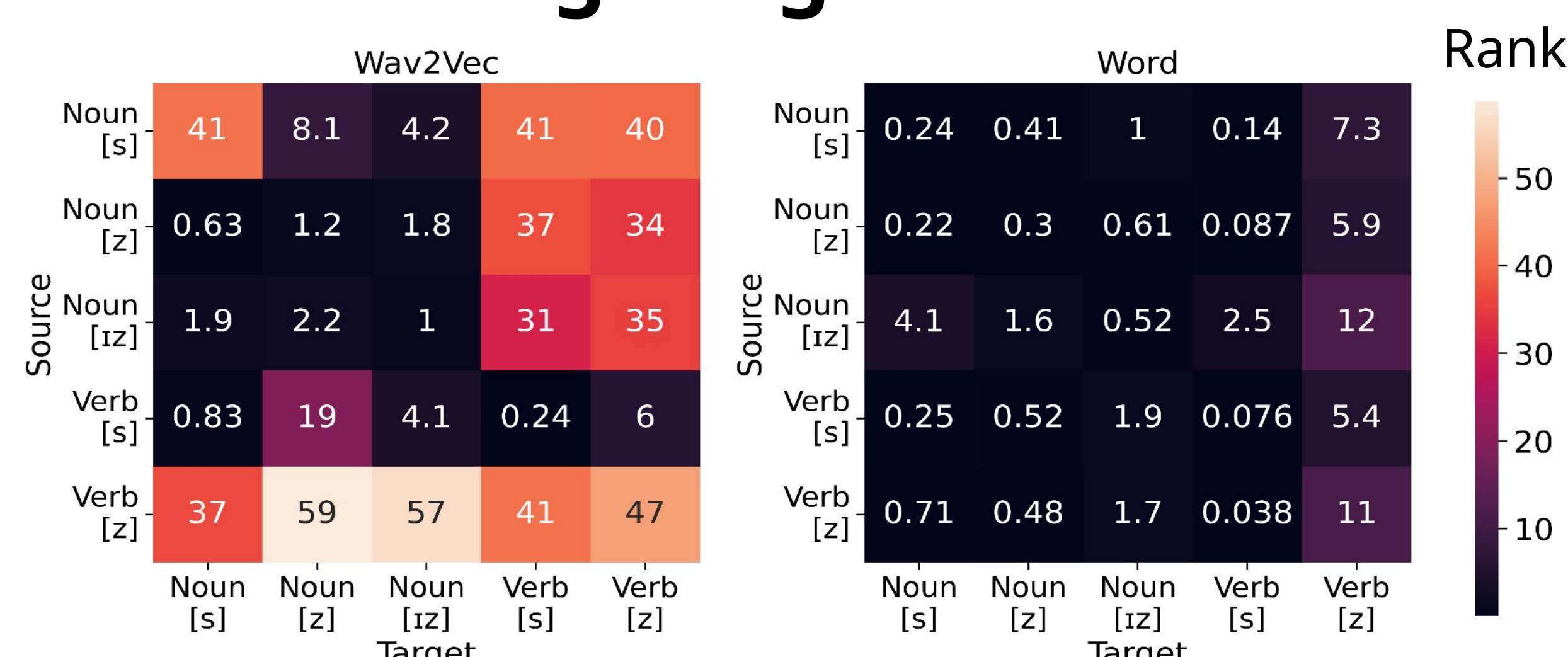Evaluate predictions by **rank** of the nearest valid

## RESULTS

- Both embedding spaces show a regular direction linking base and inflected forms
- Relations in the fine-tuned word embedding space generalize across both **morphological** and **phonological** distinctions; this is not true of the pretrained Wav2Vec embeddings
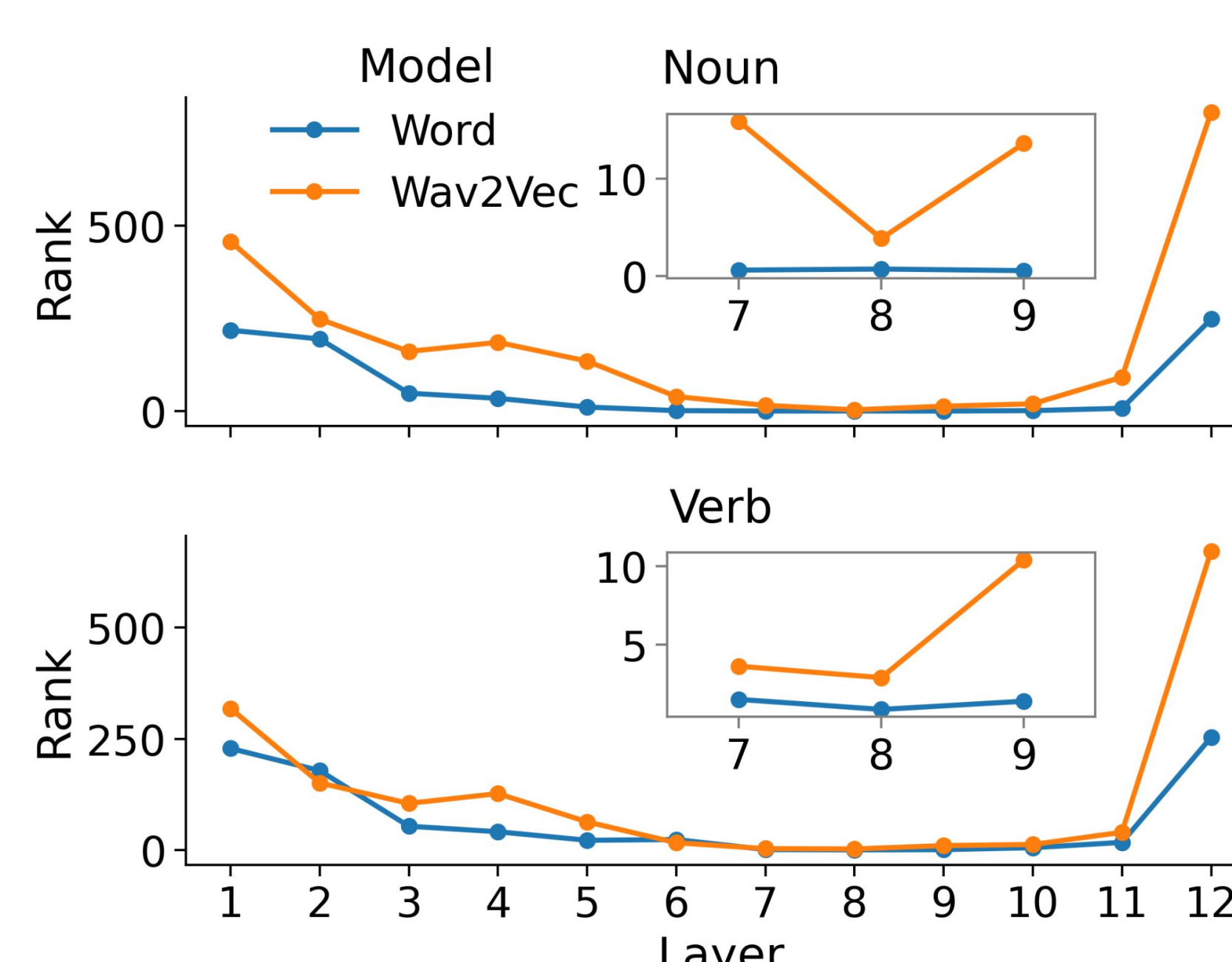


verbs
nouns

→ /z/ allomorphs
→ /s/ allomorphs
→ /ɪz/ allomorphs
→ mean difference vector

### Morphological generalization



war : wars :: lip : _____
Noun        Noun

speak : speaks :: lip : _____
Verb            Noun

war : wars :: exist : _____
Noun        Noun

speak : speaks :: exist : ____
Verb            Verb

Noun → Noun: 6.23 / 0.91
Verb → Noun: 24.86 / 3.40
Noun → Verb: 46.01 / 8.17
Verb → Verb: 36.06 / 7.83

Model: Wav2Vec / Word
Rank

### Phonological generalization



| | Noun [z] | Noun [z] | Noun [ɪz] | Verb [s] | Verb [z] |
|---|---|---|---|---|---|

Wav2Vec

| Source | | | | | |
|---|---|---|---|---|---|
| Noun [s] | 41 | 8.1 | 4.2 | 41 | 40 |
| Noun [z] | 0.63 | 1.2 | 1.8 | 37 | 34 |
| Noun [ɪz] | 1.9 | 2.2 | 1 | 31 | 35 |
| Verb [s] | 0.83 | 19 | 4.1 | 0.24 | 6 |
| Verb [z] | 37 | 59 | 57 | 41 | 47 |

Word

| Source | | | | | |
|---|---|---|---|---|---|
| Noun [s] | 0.24 | 0.41 | 1 | 0.14 | 7.3 |
| Noun [z] | 0.22 | 0.3 | 0.61 | 0.087 | 5.9 |
| Noun [ɪz] | 4.1 | 1.6 | 0.52 | 2.5 | 12 |
| Verb [s] | 0.25 | 0.52 | 1.9 | 0.076 | 5.4 |
| Verb [z] | 0.71 | 0.48 | 1.7 | 0.038 | 11 |

Target

Rank

### Per layer results

Results above use the highest performing layer for both models (layer 8).



Model: Word, Wav2Vec
Noun
Verb
Rank
Layer

## FORCED-CHOICE ANALYSIS

This movement in embedding space effectively adds a sound which is **consistent with the phonological rules** exhibited in noun and verb inflections.

[z]
own : owns :: bay : { bays **(consistent)** / base **(inconsistent)** }



Proportion of phonologically consistent choices

Source allomorph: S, Z, IZ

## TAKEAWAYS

- Acoustic word embeddings contain a linear subspace tracking an abstract morpho-phonological relationship between words
- This subspace is discoverable by optimizing for **word contrast**
- These results are not sensitive to embedding method (see paper)
- This suggests possible representations supporting human spoken word recognition

### Citations

**[1]** Baevski et al. NeurIPS 2020.
**[2]** Panayotov et al. ICASSP 2015.

**Jon Gauthier, Canaan Breiss,** Matthew Leonard, Edward Chang

UCSF University of California San Francisco

USC University of Southern California

Poster, slides, and code available here ➤