

# Compositional Morphology for Word Representations and Language Modelling

Jan Botha and Phil Blunsom  
Department of Computer Science, University of Oxford



## CLBL++ Model Definition

The log bilinear model (LBL) assigns  $n$ -gram probabilities using distributed feature vectors for words and smooth scoring functions.

**We extend the standard LBL in two ways:**

- compose word vectors from morpheme vectors (LBL++; see (5))
- partition vocabulary into word classes for fast normalisation (CLBL)

Predict next representation  $\mathbf{p}$  given preceding word vectors  $\tilde{\mathbf{q}}_j$ :

$$\mathbf{p} = \sum_{j=1}^{n-1} \tilde{\mathbf{q}}_j C_j \quad (1)$$

Score next word  $w$  and its class  $c$  (word vector  $\tilde{\mathbf{r}}_w$ , class vector  $\mathbf{s}_c$ ):

$$\nu(w) = \mathbf{p} \cdot \tilde{\mathbf{r}}_w + b_w \quad (2)$$

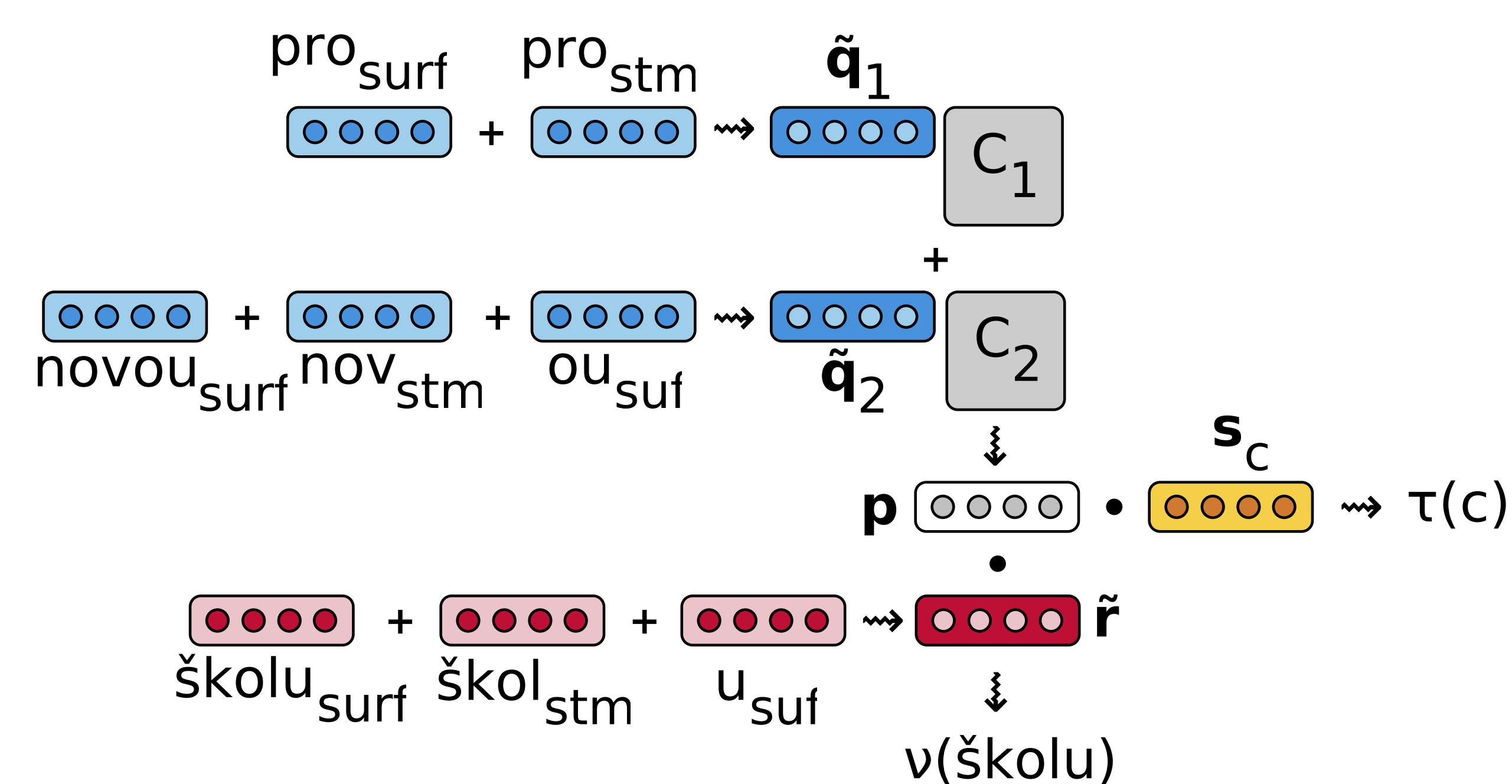
$$\tau(c) = \mathbf{p} \cdot \mathbf{s}_c + t_c \quad (3)$$

Compute probability of word  $w$  under model:

$$P(w \mid w_{i-n+1}^{i-1}) = \underbrace{P(\text{class}(w) \mid w_{i-n+1}^{i-1})}_{\sum_{c' \in \{\text{classes}\}} \frac{\exp(\tau(c'))}{\sum_{c' \in \{\text{classes}\}} \exp(\tau(c'))}} \underbrace{P(w \mid w_{i-n+1}^{i-1}, \text{class}(w))}_{\frac{\exp(\nu(w))}{\sum_{v' \in \{\text{words in class } c\}} \exp(\nu(v'))}} \quad (4)$$

We train the model against an  $L_2$ -regularised maximum likelihood objective function using adaptive gradient descent.

## Model Diagram



CLBL++ model illustrated for the Czech trigram 'pro novou školu'.

## Addition as Composition

- Map each word type  $v$  to a sequence of surface-level factors  $\mu(v)$ .
- Define word type vector  $\tilde{\mathbf{r}}_v$  as sum of its factor vectors:

$$\tilde{\mathbf{r}}_v \equiv \sum_{f \in \mu(v)} \mathbf{r}_f \quad (\text{notation: } \overrightarrow{\text{word}} \equiv \overrightarrow{\text{factor}_1} + \overrightarrow{\text{factor}_2} + \dots) \quad (5)$$

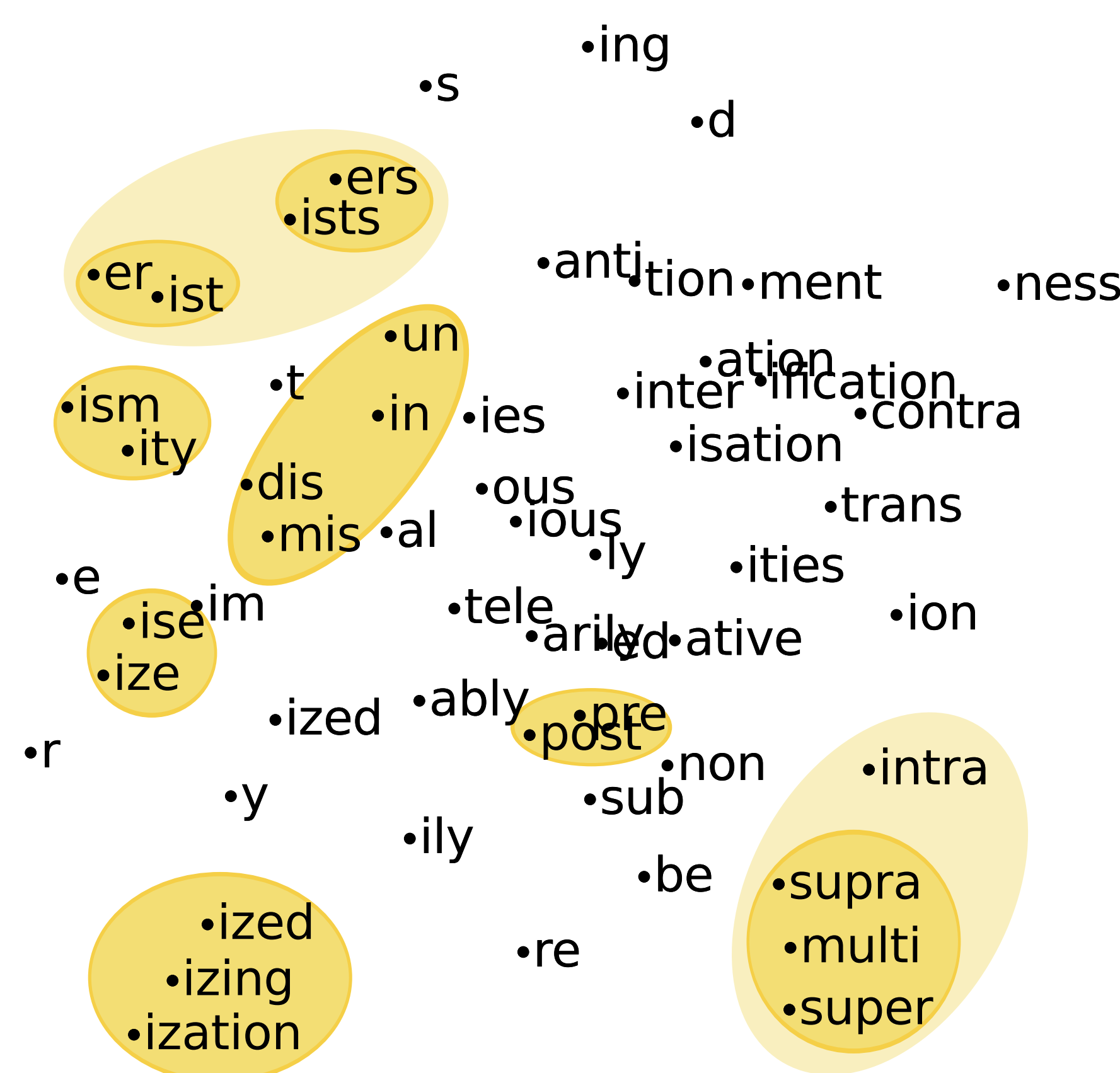
**Effect is to tie words with shared morphemes, e.g.**

$$\begin{aligned} \overrightarrow{\text{unexpected}} &= \overrightarrow{\text{unexpected}_{surf}} + \overrightarrow{\text{un}_{pre}} + \overrightarrow{\text{expect}_{stm}} + \overrightarrow{\text{ed}_{suf}} \\ \overrightarrow{\text{expectations}} &= \overrightarrow{\text{expectations}_{surf}} + \overrightarrow{\text{expect}_{stm}} + \overrightarrow{\text{ation}_{suf}} + \overrightarrow{s}_{suf} \end{aligned}$$

**Including surface form as factor means our approach**

- appreciates order:  $\overrightarrow{\text{hangover}} \neq \overrightarrow{\text{overhang}}$
- handles non-compositionality:  $\overrightarrow{\text{greenhouse}} \neq \overrightarrow{\text{green}} + \overrightarrow{\text{house}}$
- overcomes noise from unsupervised segmentor (we used Morfessor).

## Morpheme Embeddings Acquired



*t-SNE projections of English affix vectors learnt by CLBL++, with shading added manually for emphasis.*

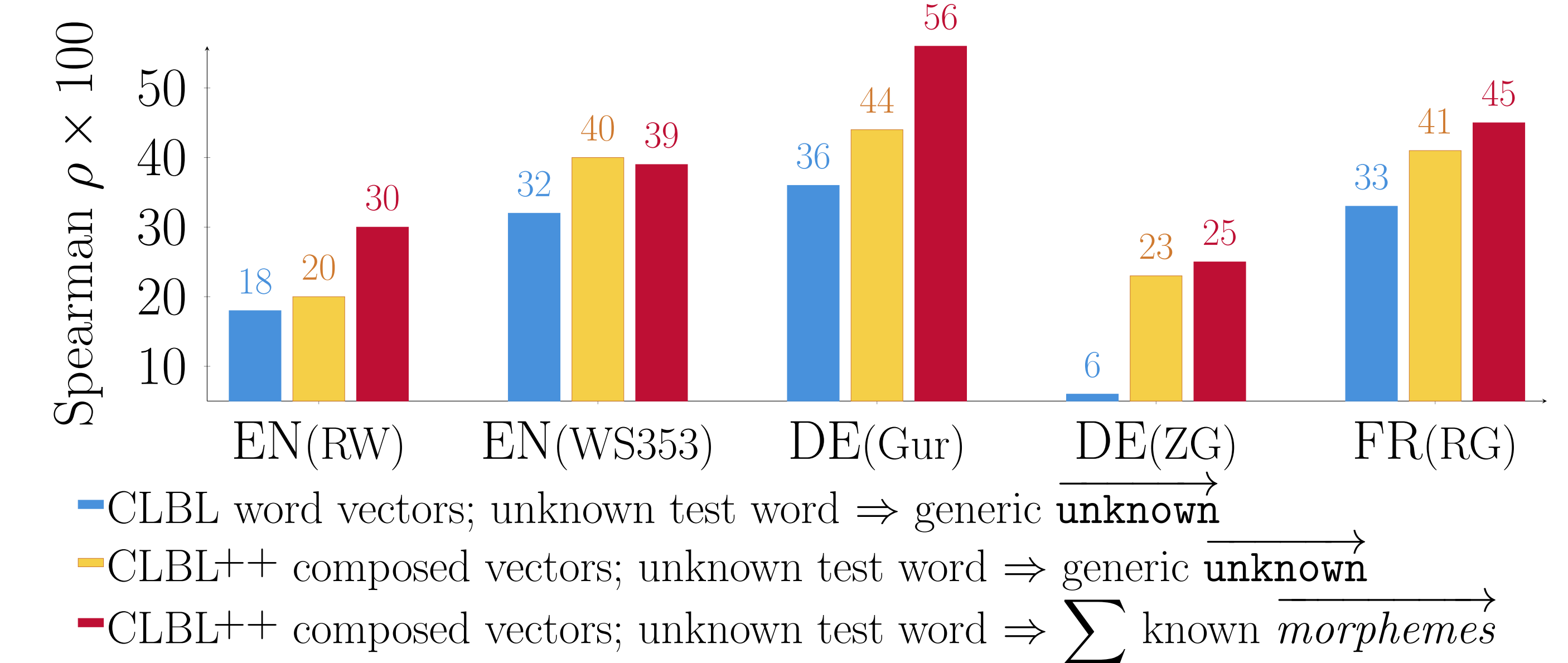
## Summary

**Integrate morphology into vector-based language models**

- Simple, scalable, unsupervised method
- Learns morpheme vectors as part of model
- Word classes enable integration into MT decoder
- Improvements in three evaluation settings and multiple languages

## Word-pair Similarity Rating Task

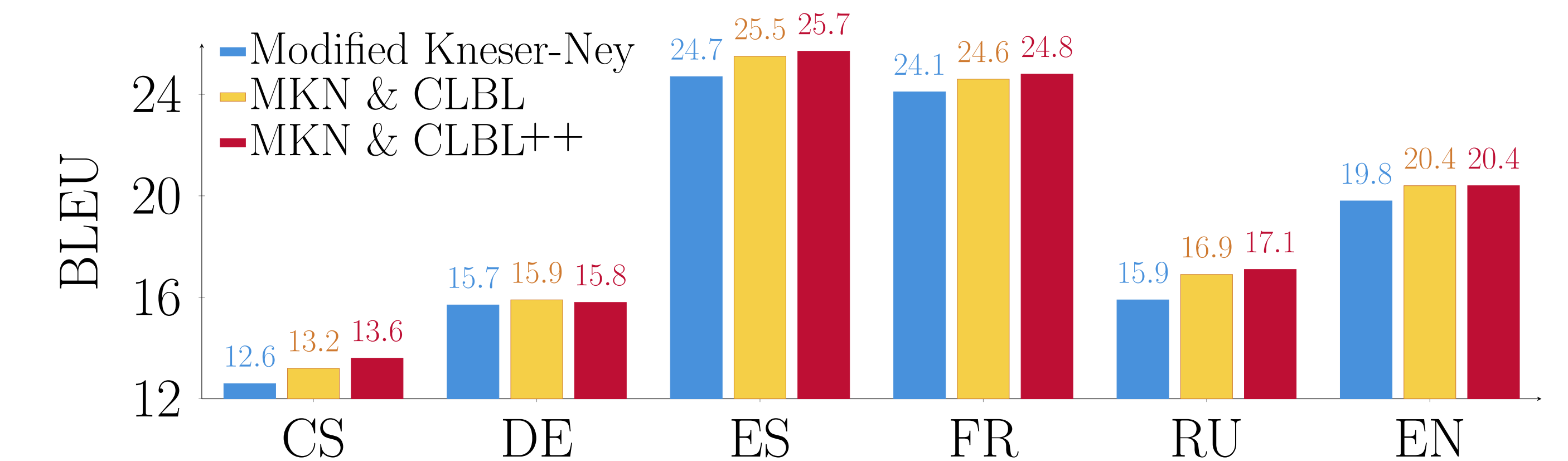
The CLBL++ model learns vectors that obtain stronger correlation with human judgements of word-pair similarity. Morpheme vectors allow more subtle handling of unknown words.



*Results for measuring similarity of word-pairs using learnt vectors.*

## Machine Translation Task

Class-based partitioning of vocabulary speeds up computation of normalised language model probabilities. This is crucial for the large vocabularies of morphologically rich languages, and enabled integration of the CLBL/CLBL++ inside an MT decoder.



*Results for translation into different languages, varying the LMs.*