# KerBS: Kernelized Bayesian Softmax for Text Generation

Ning Miao,   Hao Zhou,   Chengqi Zhao,   Wenxian Shi,   Lei Li

## Motivation

- Softmax based generation models assume context embeddings of a word should concentrate around its word embedding as in Figure 1.a.
- However, embeddings of some words actually exhibit properties such as **multi-sense** (Figure 1.b) or **varying variance** (Figure 1.c).
- We need to explicitly model the properties for more accurate word prediction!
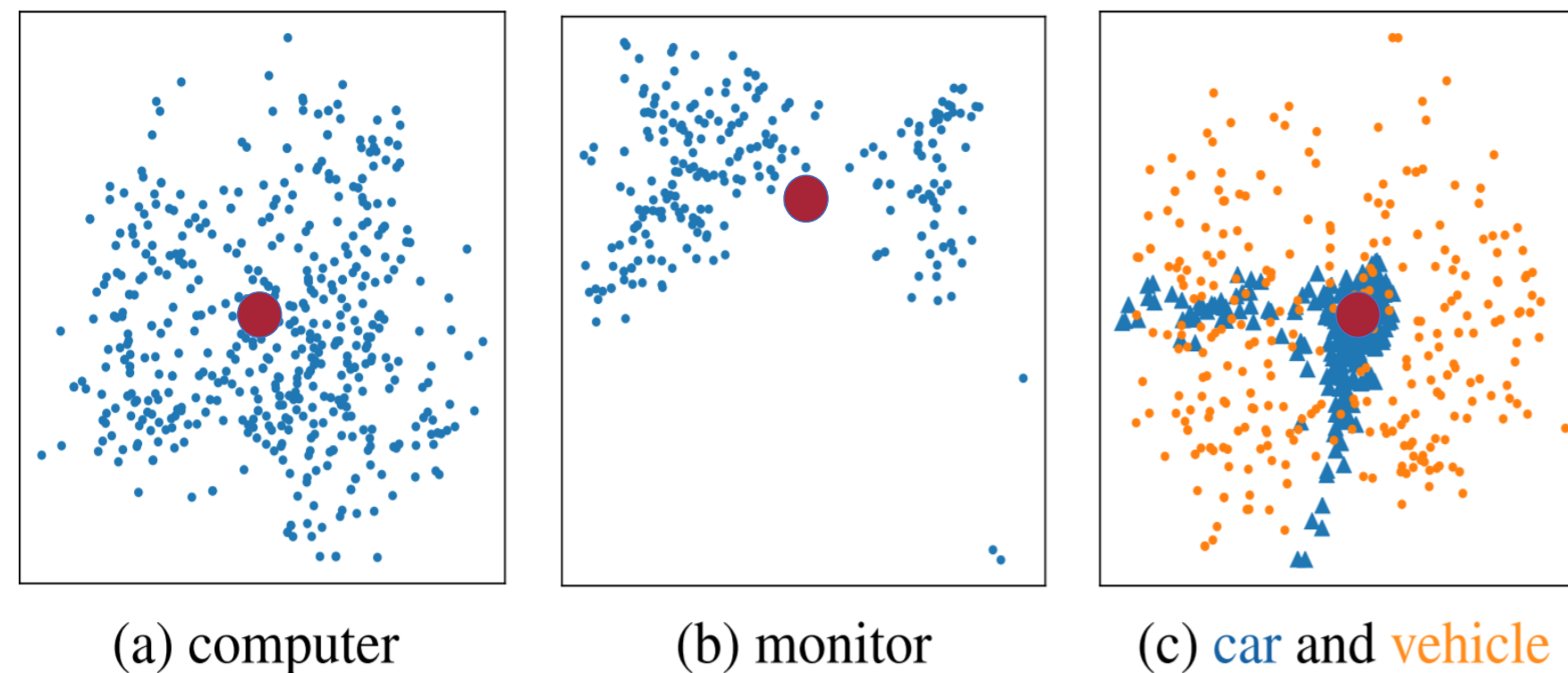


(a) computer     (b) monitor     (c) car and vehicle

Figure 1: Context embeddings and word embeddings (●) of different word.



## Method

- ### Multi-sense
  - We replace word embeddings with **sense embeddings**, and sums all sense probabilities of a word to get the probability of the word.
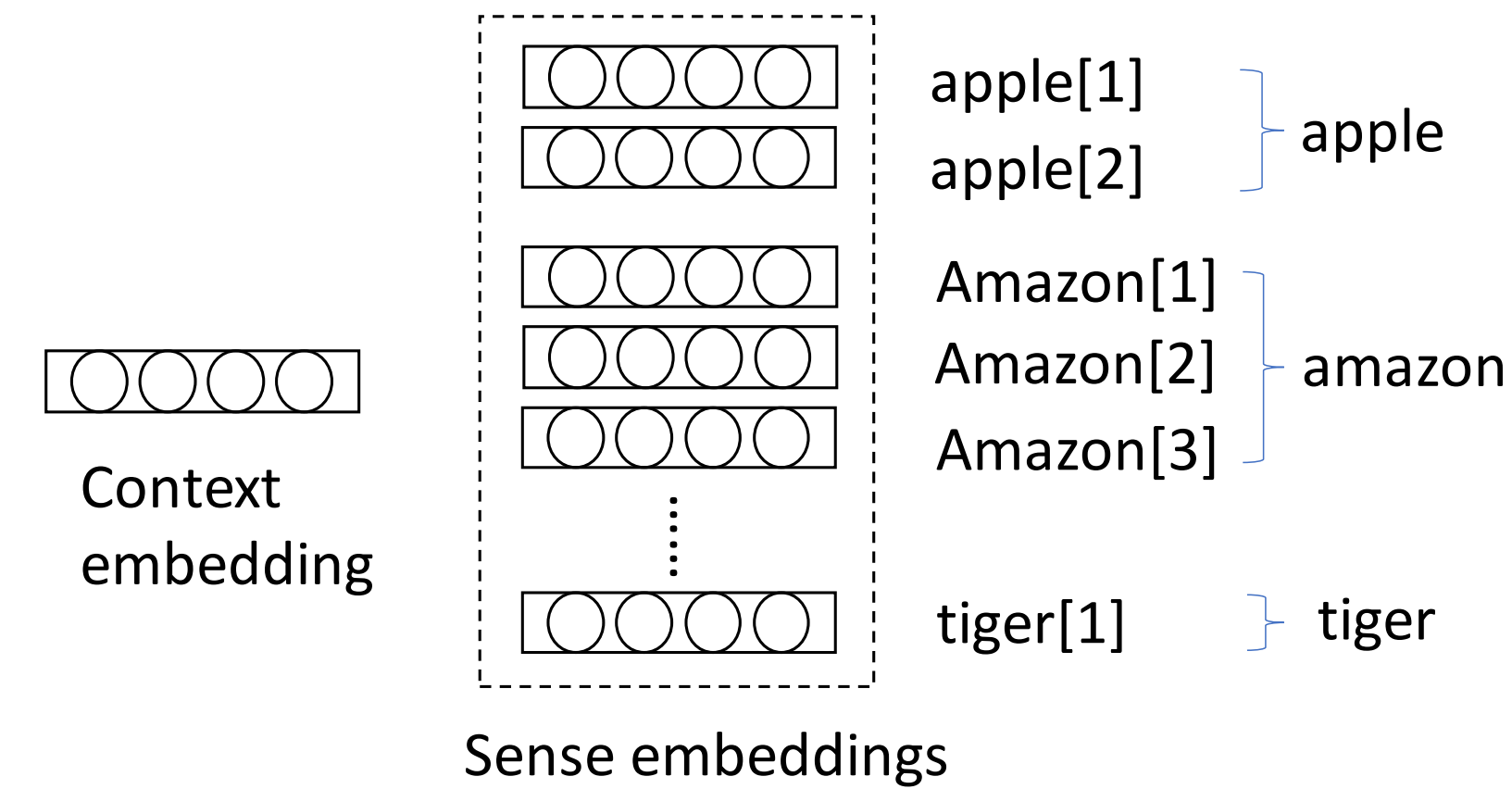


apple[1]
apple[2]  } apple

Amazon[1]
Amazon[2]  } amazon
Amazon[3]

tiger[1]  } tiger

Context embedding

Sense embeddings

Figure 2: Illustration of sense embedding.

- ### Kernelized Variances
  - Since Gaussian distribution is instable in high-dimensional space, we designed a **kernel function** to model **different variances** of each word's context embeddings.

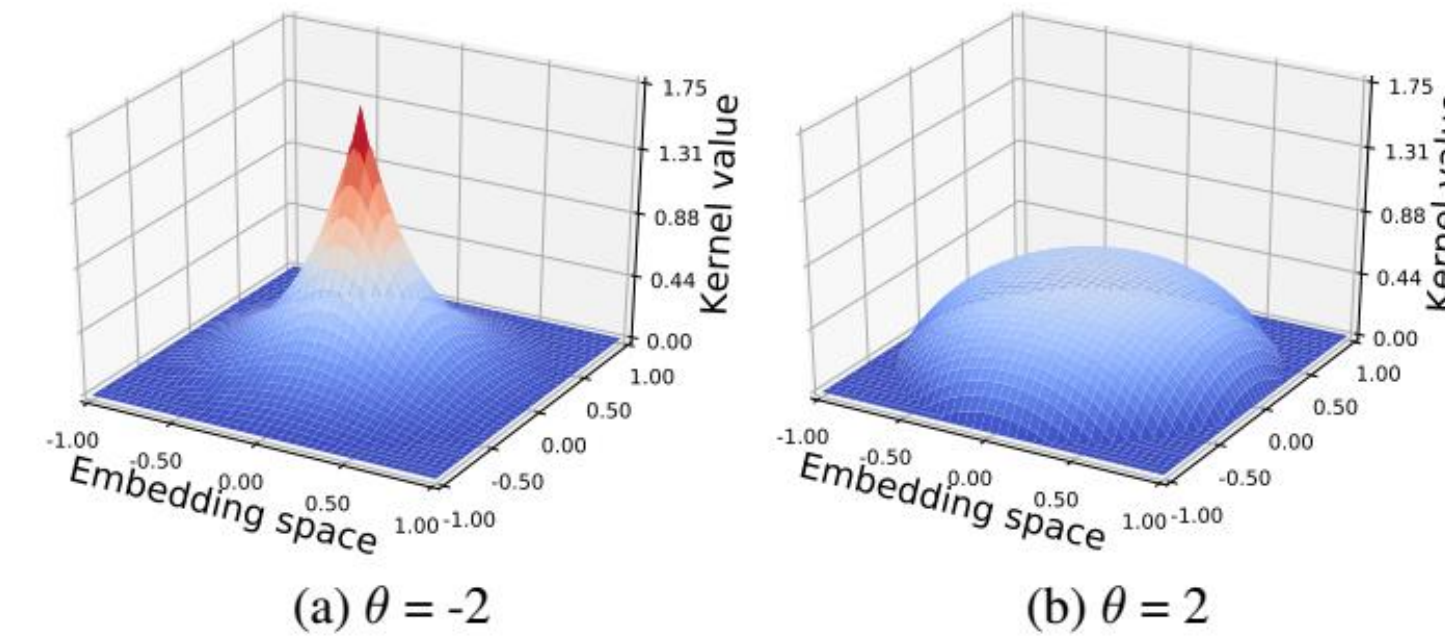$$\mathcal{K}_\theta(h, e) = |h||e|(a\exp(-\theta\cos(h, e)) - a)$$



(a) $\theta = -2$     (b) $\theta = 2$

Figure 2: Kernels with different $\theta$.

- ### Dynamic Sense Allocation
  - During training, we record the **usage** of each sense and **log prediction accuracy** of each word.
  - After every M steps, we **reallocate** the least used senses to the poorly predicted words.



$Word_1$ (-1.3)   A
$Word_2$ (-3.2)
$Word_3$ (-1.1)   B
$Word_4$ (-1.6)
$Word_5$ (-2.9)

→

$Word_1$
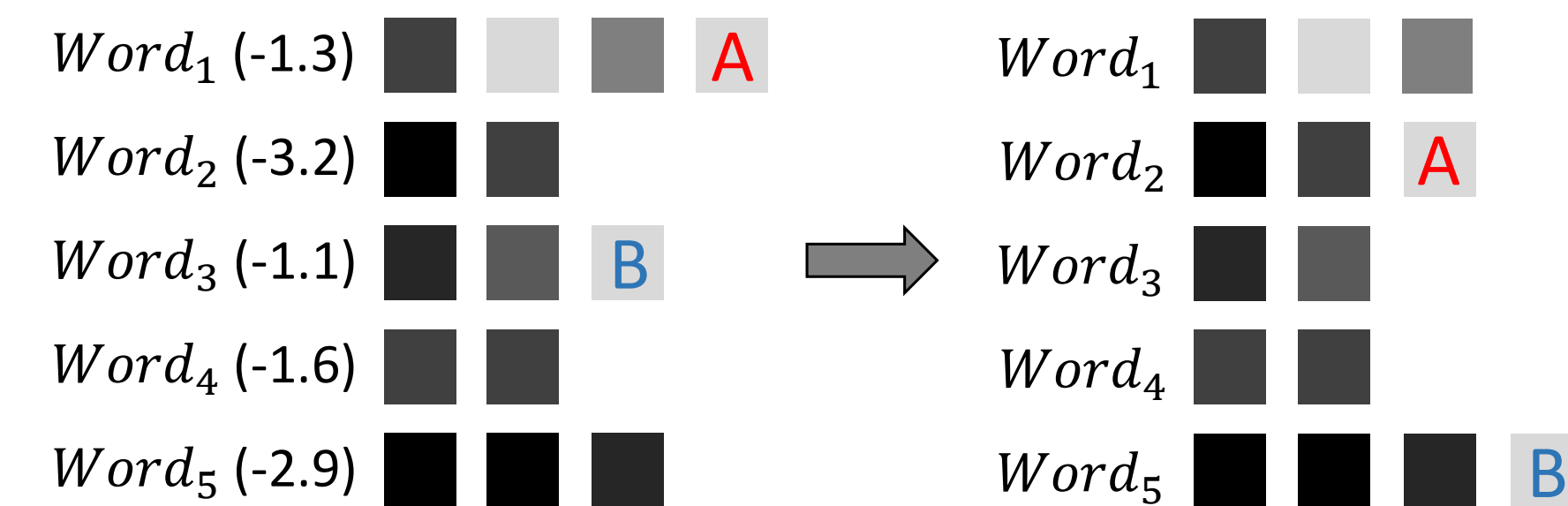$Word_2$   A
$Word_3$
$Word_4$
$Word_5$   B

Figure 3: Illustration of sense reallocation. Numbers in the brackets are log prediction accuracies of each word. We use squares to represent senses and darker color means higher usage.

## Experiments and Analyses

- Better performance than baselines.

| Tasks | Base models | Metrics | Base | Base+Mos | Bas+KerBS |
|-------|-------------|---------|------|----------|-----------|
| MT | Transformer | BLEU-4 | 29.61 | 28.54 | **30,90** |
| | Seq2Seq | BLEU-4 | 25.91 | 26.45 | **27.28** |
| LM | GRU | PPL | 103.12 | 102.72 | **102.17** |
| Dialog | Transformer | BLEU-1 | 10.61 | 9.81 | **10.90** |
| | Seq2Seq | BLEU-1 | 16.56 | 13.73 | **17.85** |
| | | Human | 1.24 | 1.04 | **1.40** |

- More sense are allocated to words with complex meanings.

| Sense | 1 | 2 | 3 | 4 |
|-------|---|---|---|---|
| word | Redwood heal structural theoretical rotate | particular figure during known size | open order amazing sound base | they work body power change |

- $\theta$ reflects the sematic scopes of words.



Beijing    China    earth

monkey  cat    animal

Jeep  Ford    car

-0.5   -0.1   0   0.1   $\theta$