

Learning Word Embeddings

Navid Rekabsaz



navid.rekabsaz@idiap.ch

Agenda

- Neural Networks - Recap
- Language Modeling
- Neural Word Embedding - Revisited
 - word2vec Skip-Gram with Negative Sampling
 - Bias in Word Embeddings
- Contextualized Word Embeddings
 - ELMo
 - BERT

Agenda

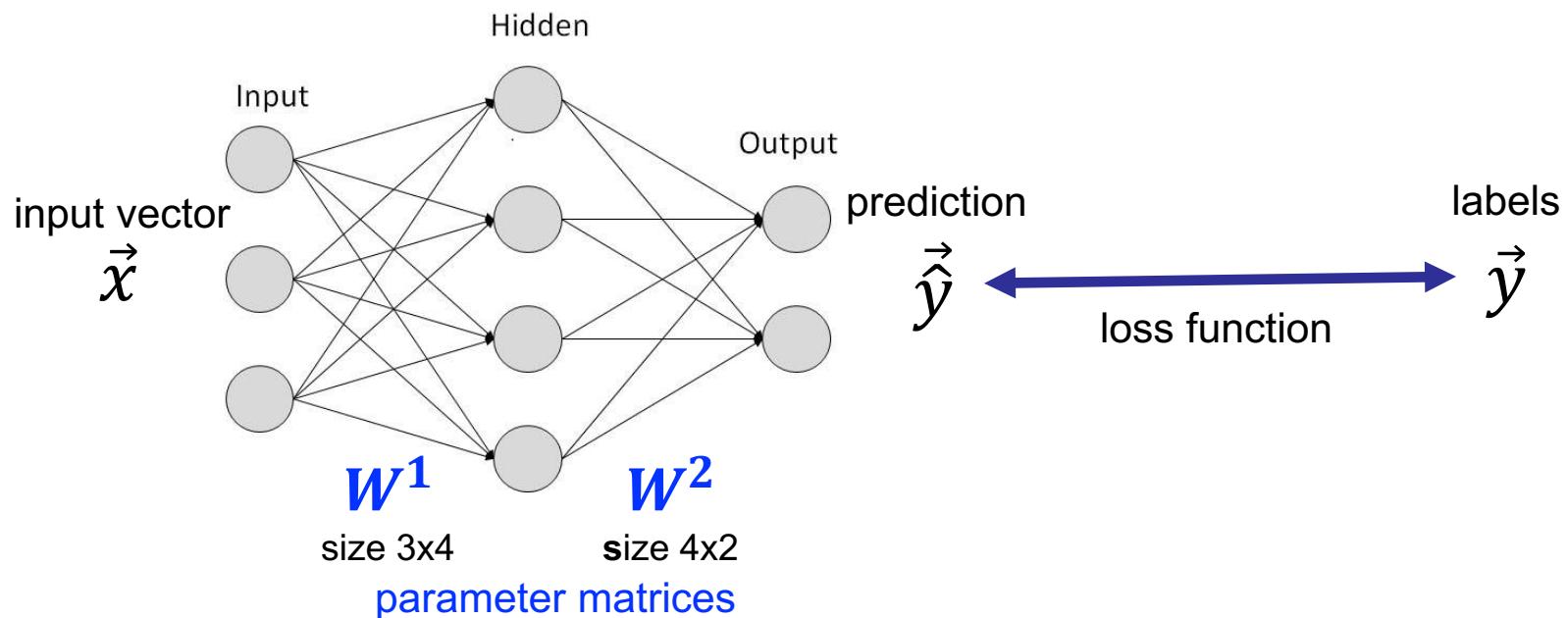
- Neural Networks - Recap
- Language Modeling
- Neural Word Embedding - Revisited
 - word2vec Skip-Gram with Negative Sampling
 - Bias in Word Embeddings
- Contextualized Word Embeddings
 - ELMo
 - BERT

Linear Algebra

- Scalar a
- Vector \vec{b}
- Matrix W
- Tensor: generalization to higher dimensions
- Dot product
 - $\vec{a} \cdot \vec{b}^T =$
dimensions: $1 \times d \cdot d \times 1 =$
 - $\vec{a} \cdot W =$
dimensions: $1 \times d \cdot d \times e =$
 - $A \cdot B =$
dimensions: $l \times m \cdot m \times n =$
- Element-wise Multiplication
 - $\vec{a} \odot \vec{b} =$

Neural Networks

- Neural Networks are **non-linear functions** with many parameters
$$\vec{\hat{y}} = f(\vec{x})$$
- They composed of several simple **non-linear operations**
- Normally, the objective is to **maximize likelihood**, namely
$$p(y|x, \theta)$$
- Optimized using gradient-based algorithms like Stochastic Gradient Descent

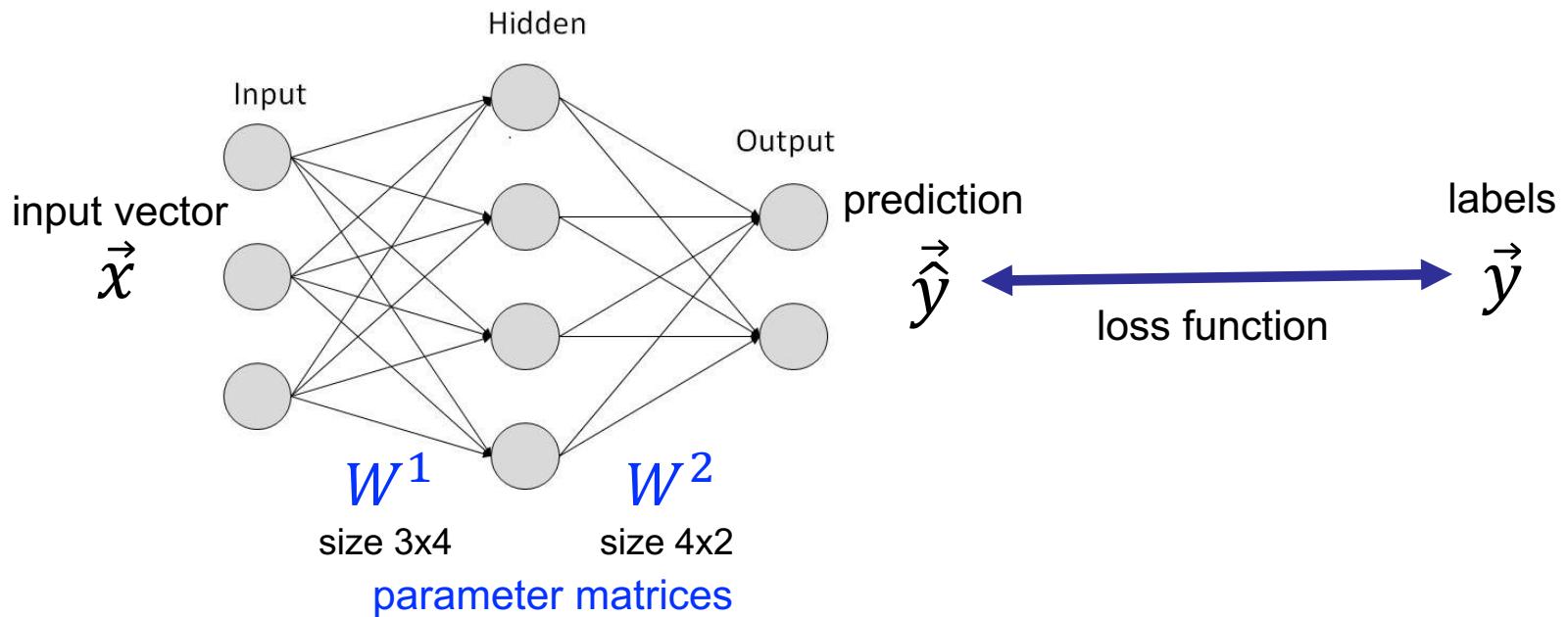


Training a Neural Network

Initialize parameters

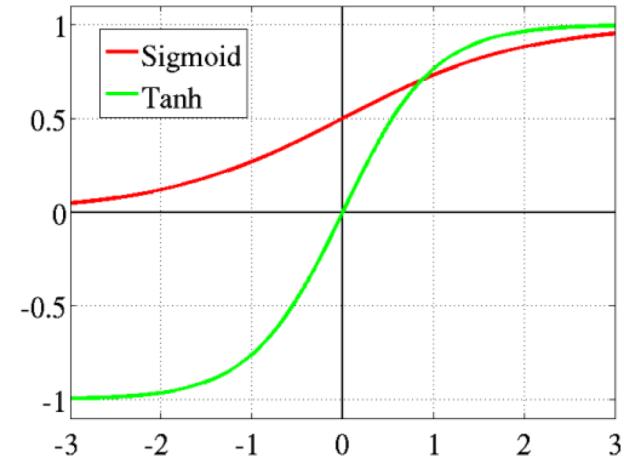
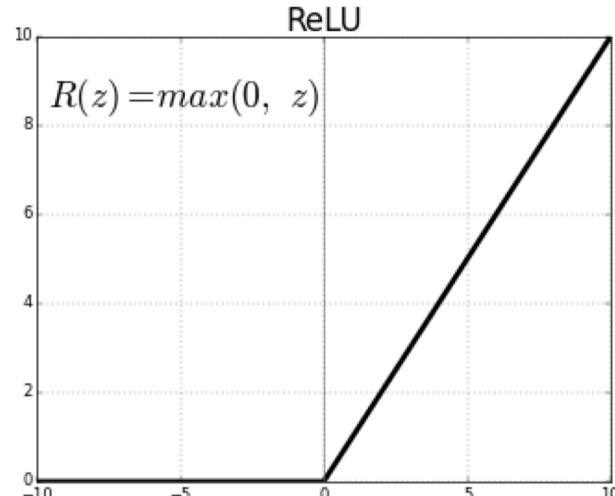
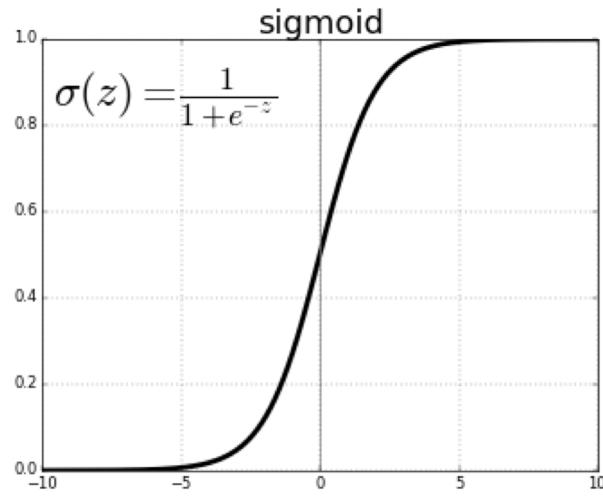
Loop over training data (or minibatches)

1. Do **forward pass**: given input \vec{x} predict output \hat{y}
2. Calculate **loss** by comparing \hat{y} with labels y
3. Do **backpropagation**: calculate the gradient of each parameter in regard to the loss function
4. Update parameters in the direction of gradient
5. Exit if some stopping criteria are met



Non-linearities

- Sigmoid
 - Projects input to value between 0 to 1 → becomes like a probability value
- ReLU (Rectified Linear Units)
 - Suggested for deep architectures to prevent vanishing gradient
- Tanh



Softmax

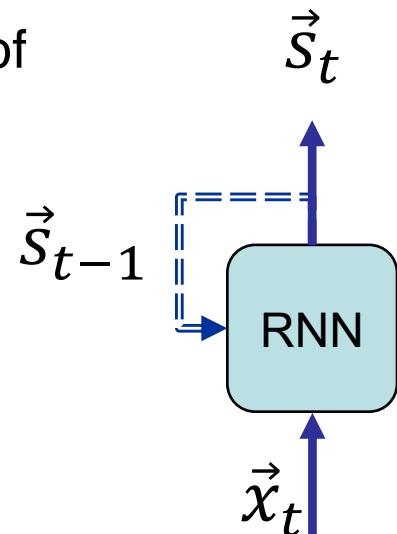
- softmax turns a vector to a probability distribution
 - The vector values become in the range of 0 to 1 and sum of all the values is equal 1

$$\text{softmax}(\vec{v})_i = \frac{e^{v_i}}{\sum_{k=1}^d e^{v_k}}$$

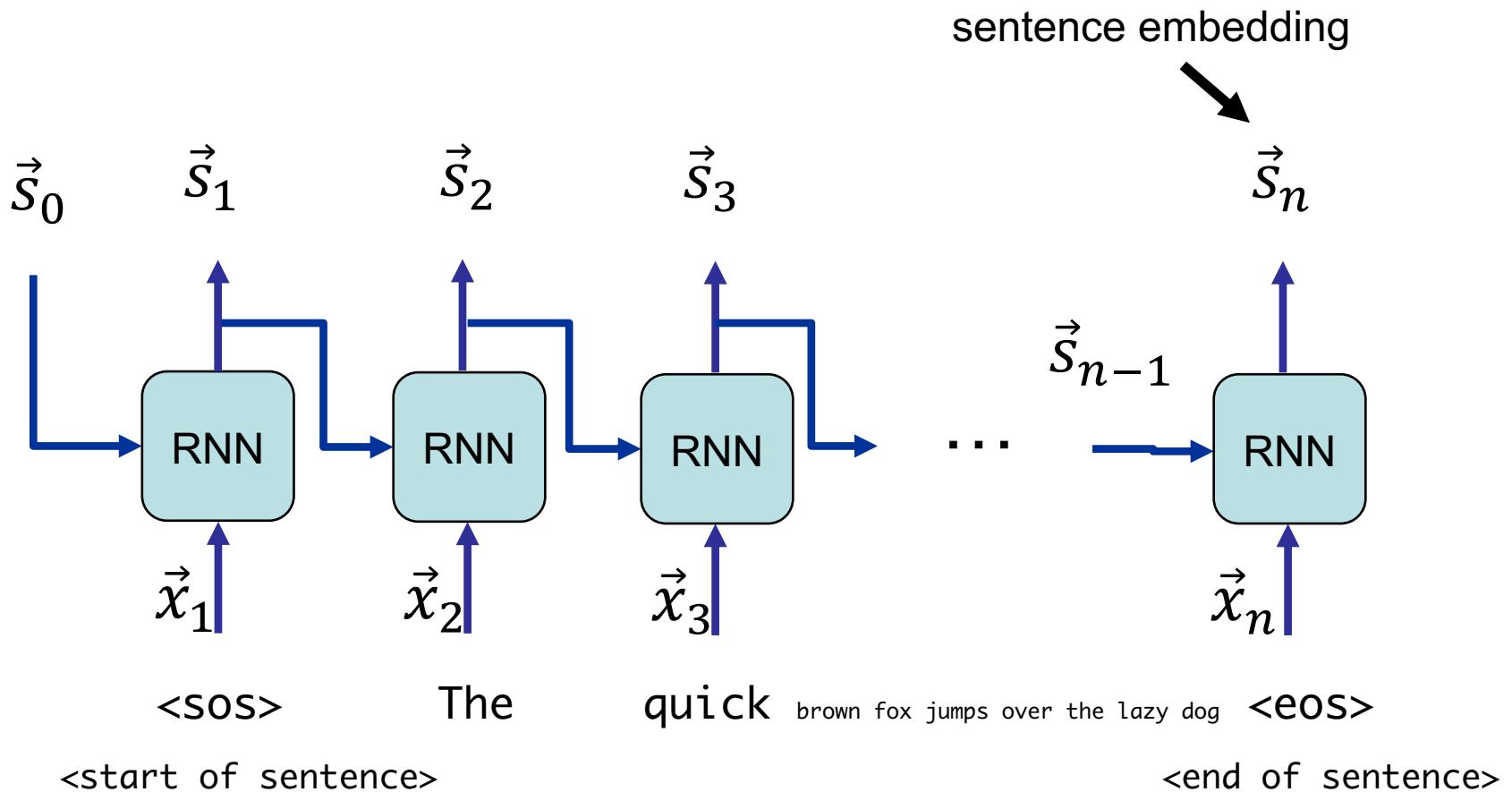
- Normally applied to the output layer and provide a probability distribution over output classes
- For example, given four classes:
 $\vec{y} = [2, 3, 5, 6]$ $\text{softmax}(\hat{y}) = [0.01, 0.03, 0.26, 0.70]$

Recurrent Neural Networks

- Encodes/Embeds a **sequence of entities** (vectors) such as ...
 - Sequence of word vectors
 - Time series
- ... to a **final composed vector** as well as **intermediary** vectors on each time step
- The output \vec{s}_t is a function of input and the output of the previous time step
- Output is also called **hidden state**
- Hidden state is in fact a “memory” of the **previous entities**



RNN - unfolded

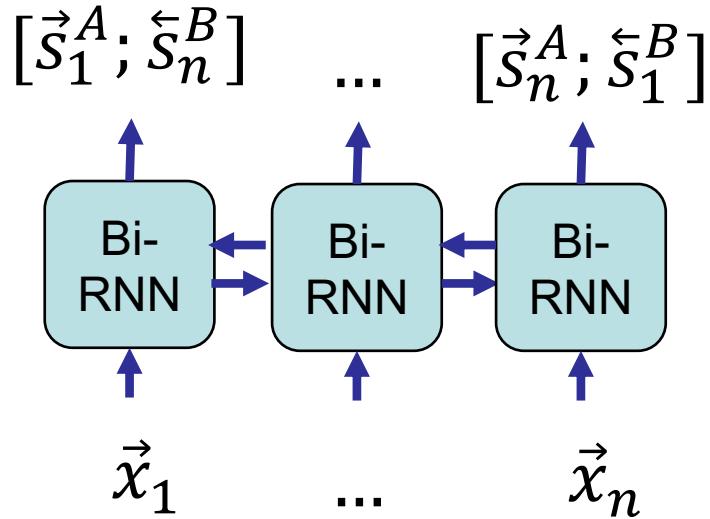


Types and Bidirectional

- Common types of RNN
 - Standard (Elman) RNN
 - Gated Recurrent Unit (GRU)
 - Long Short-Term Memory (LSTM)

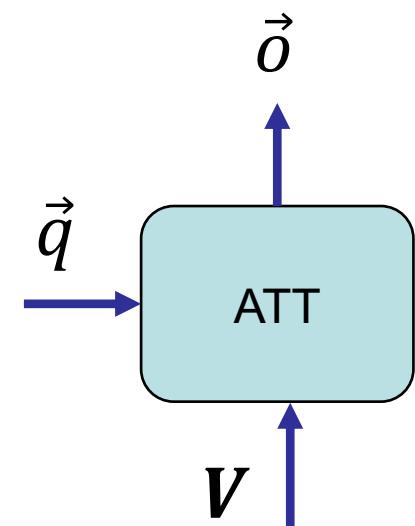
Bidirectional RNN

- Two RNNs together: A and B
- Reading the sequence from
 - Beginning to end with RNN A
 - End to beginning with RNN B
- Output at time step t is the concatenation of two hidden states \vec{s}_t^A and \vec{s}_{n-t+1}^B



Attention Networks

- Encodes/Embeds a set of vectors to a composed vector
- Given a query vector \vec{q} and a matrix of values V , an attention network assigns weights (attentions) to the vectors of V , and produces output vector \vec{o}



- General form of an attention network

$$\vec{o} = ATT(\vec{q}, V)$$

➤ In general, query is also a matrix. Here it is assumed as a vector for simplicity

Attention Networks - details

- Given the query, an attention network learns to assign some amount of **attention** α_i on each value vector \vec{v}_i using the attention function f

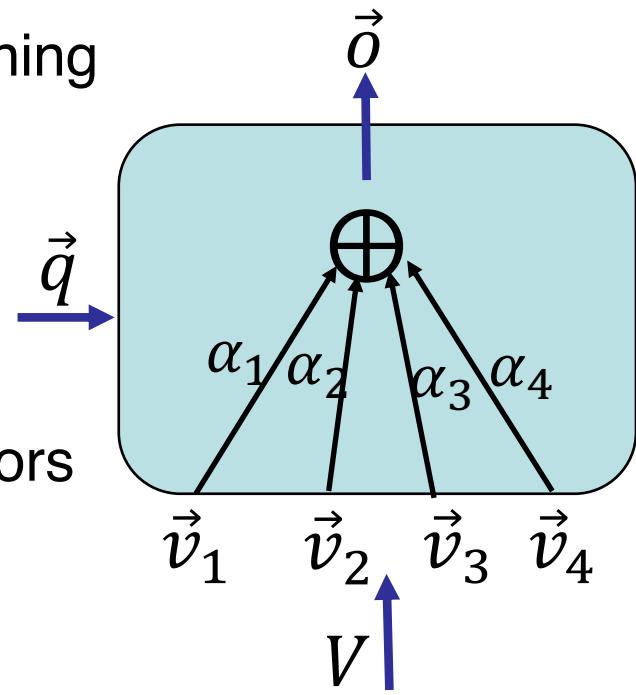
$$\alpha_i = f(\vec{q}, \vec{v}_i)$$

- where α is a **probability distribution**, meaning

$$\sum_{i=1}^n \alpha_i = 1$$

- Output is the **weighted sum** of value vectors

$$\vec{o} = \sum_{i=1}^n \alpha_i \cdot \vec{v}_i$$



Multi-Layer Perceptron Attention

- First non-normalized attention \tilde{a}_i is calculated by a neural network

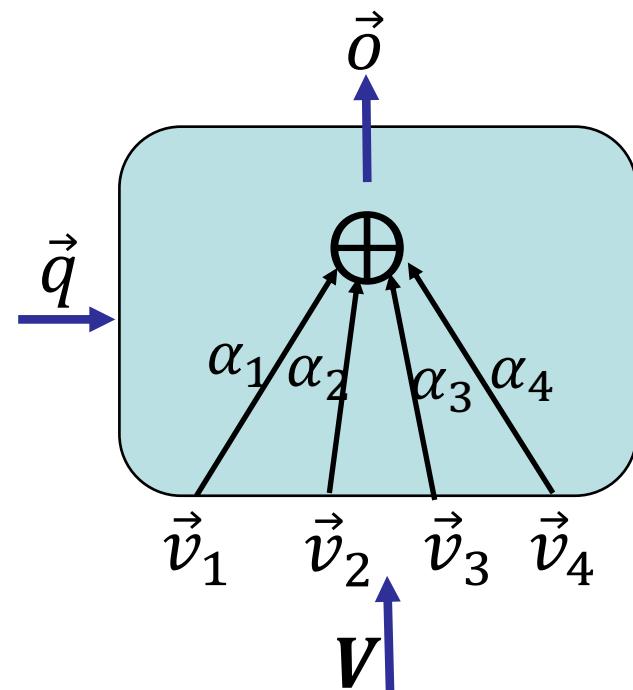
$$\tilde{a}_i = \vec{u} \cdot \tanh(\vec{q} \cdot W_1 + \vec{v}_i \cdot W_2)$$

- As before, attentions are normalized with softmax

$$\alpha_i = \text{softmax}(\tilde{a})_i$$

- and output is ...

$$\vec{o} = \sum_{i=1}^n \alpha_i \cdot \vec{v}_i$$



➤ Model parameters are shown in red

Agenda

- Neural Networks - Recap
- Language Modeling
- Neural Word Embedding - Revisited
 - word2vec Skip-Gram with Negative Sampling
 - Bias in word embeddings
- Contextualized Word Embeddings
 - ELMo
 - BERT

Language Model

- Ideally, a Language Model aims to model the language!
- In practice, Language Models use the statistics of the collection to calculate the probability:

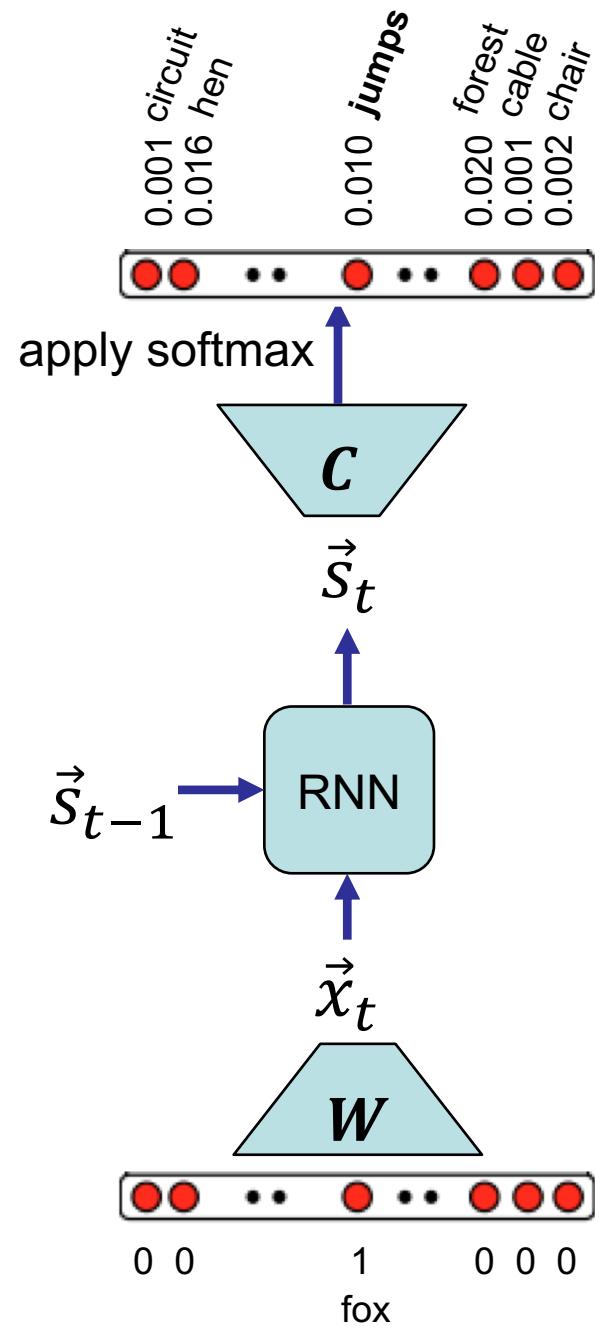
$$P(w|\text{context})$$

- A Language Model can answer the questions like
 - Given “The quick brown fox jumps over the lazy ...”, what is the probability of the next word:
 $P(w|\text{The quick brown fox jumps over the lazy})$
- Neural Networks is one way of calculating this probability

Training Language Model with RNN

- Input:
 - a text corpus with n words and $|\mathbb{V}|$ unique words: $[w_1, w_2, \dots, w_n]$
 - a RNN model
 - encoder embedding \mathbf{W} with dimensions: $|\mathbb{V}| \times d$
 - decode embedding \mathbf{C} with dimensions: $|\mathbb{V}| \times d$
- Loop over the text corpus; at time step t :
 - a. Create the **one-hot vector** of the word w_t
 - b. Calculate dot product of the one-hot vector to \mathbf{W} , resulting to \vec{w}_t
 - c. Input \vec{w}_t and \vec{s}_{t-1} to RNN and calculate \vec{s}_t
 - d. Calculate dot product of \vec{s}_t to \mathbf{C}^T , followed by **softmax**, achieving
$$P(w_{t+1} | w_t, \vec{s}_{t-1})$$
 - e. Calculate **Negative Log Likelihood** loss $\mathcal{L} = -\log p(w_{t+1} | w_t, \vec{s}_{t-1})$
 - f. Update parameters using the gradient of the loss

- Calculate loss
- Update the parameters



Transfer Learning with Embeddings

- Language modeling is a common aspect in various language processing tasks
- We train neural embedding models on large-scale text corpora using the **unsupervised objective of language modeling**
 - takes relatively longer time, but is done once
- The pre-trained parameters of the models are then used in various **supervised tasks**
 - Almost all language processing tasks such as Information Retrieval, Summarization, Question Answering, etc.
 - Usually followed by a fine-tuning of the pre-trained parameters

Agenda

- Neural Networks - Recap
- Language Modeling
- Neural Word Embedding - Revisited
 - word2vec Skip-Gram with Negative Sampling
 - Bias in Word Embeddings
- Contextualized Word Embeddings
 - ELMo
 - BERT

Recap: Skip-Grams

- Word pairs found in the text within a context window of size 2
- Language Model objective → predict context word from the center word: $P(c|w)$

Source Text	Training Samples
The quick brown fox jumps over the lazy dog. →	(the, quick) (the, brown) (w, c)
The quick brown fox jumps over the lazy dog. →	(quick, the) (quick, brown) (quick, fox)
The quick brown fox jumps over the lazy dog. →	(brown, the) (brown, quick) (brown, fox) (brown, jumps)
The quick brown fox jumps over the lazy dog. →	(fox, quick) (fox, brown) (fox, jumps) (fox, over)

Tesgüino!

bottle

likes

drink

everybody

corn

Tesgüino

makes

you

we

make

Example from J. R. Firth

bottle

fermentation

grain

brew

Ale

medieval

bar

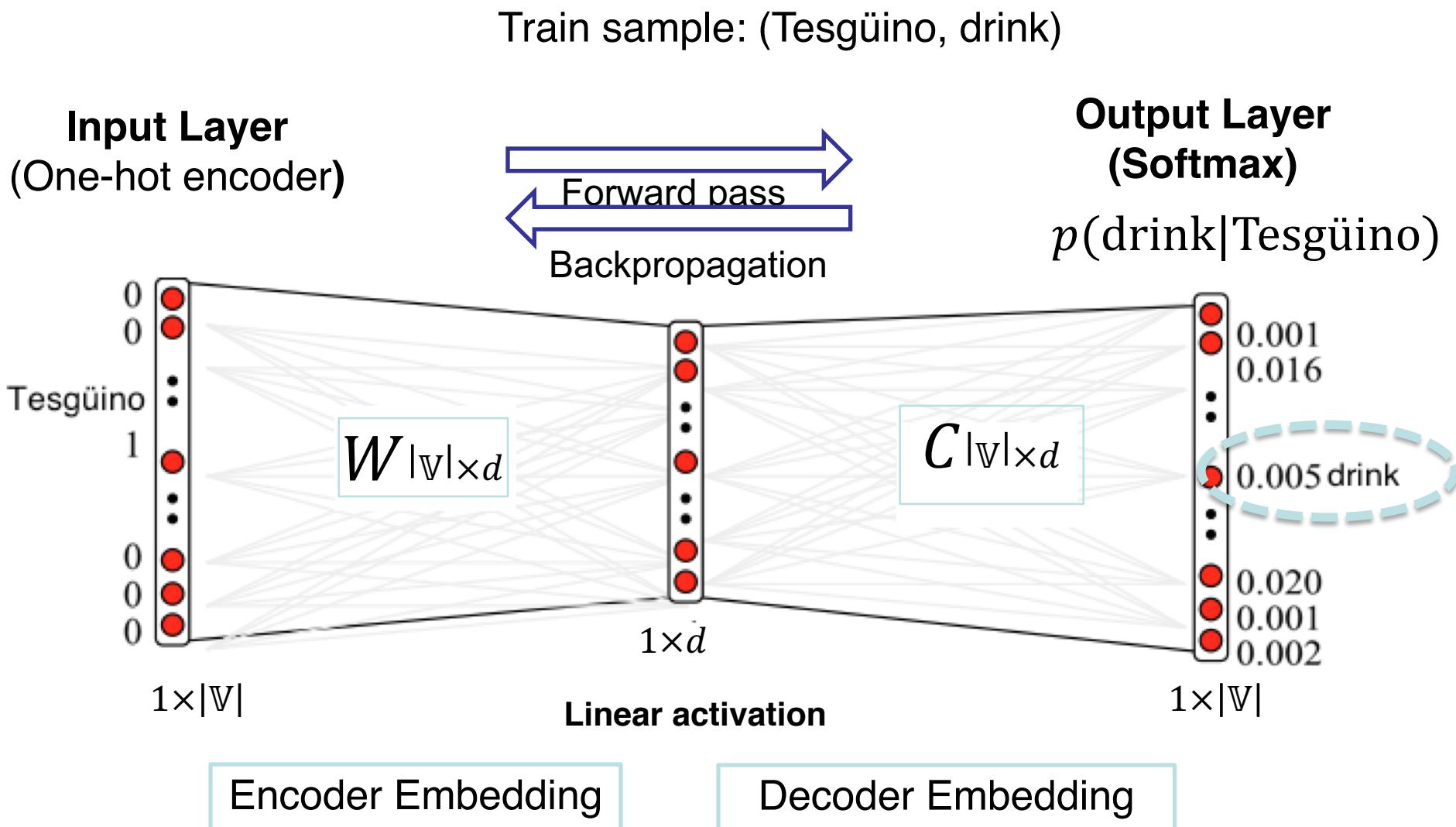
pale

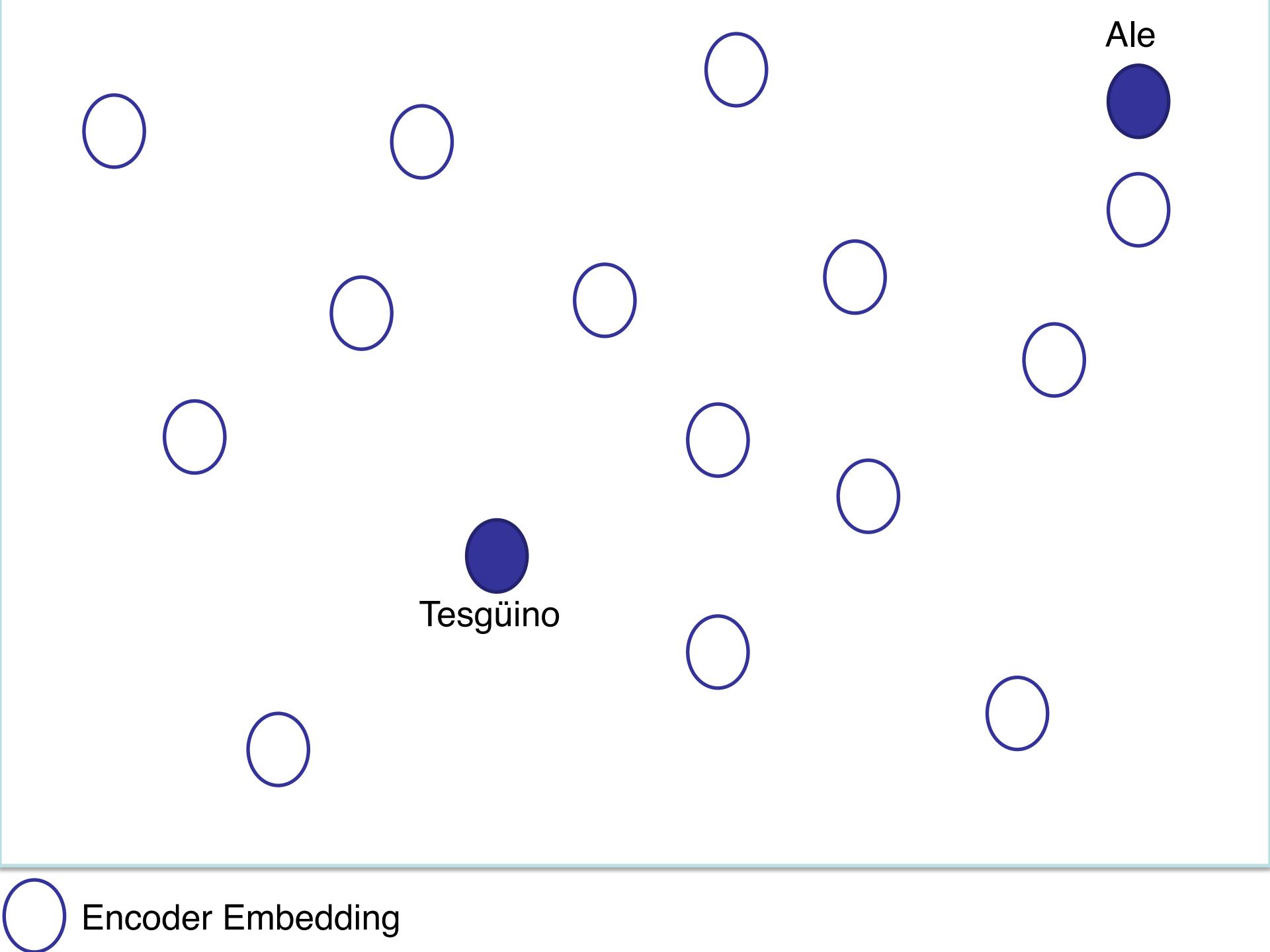
drink

Ale



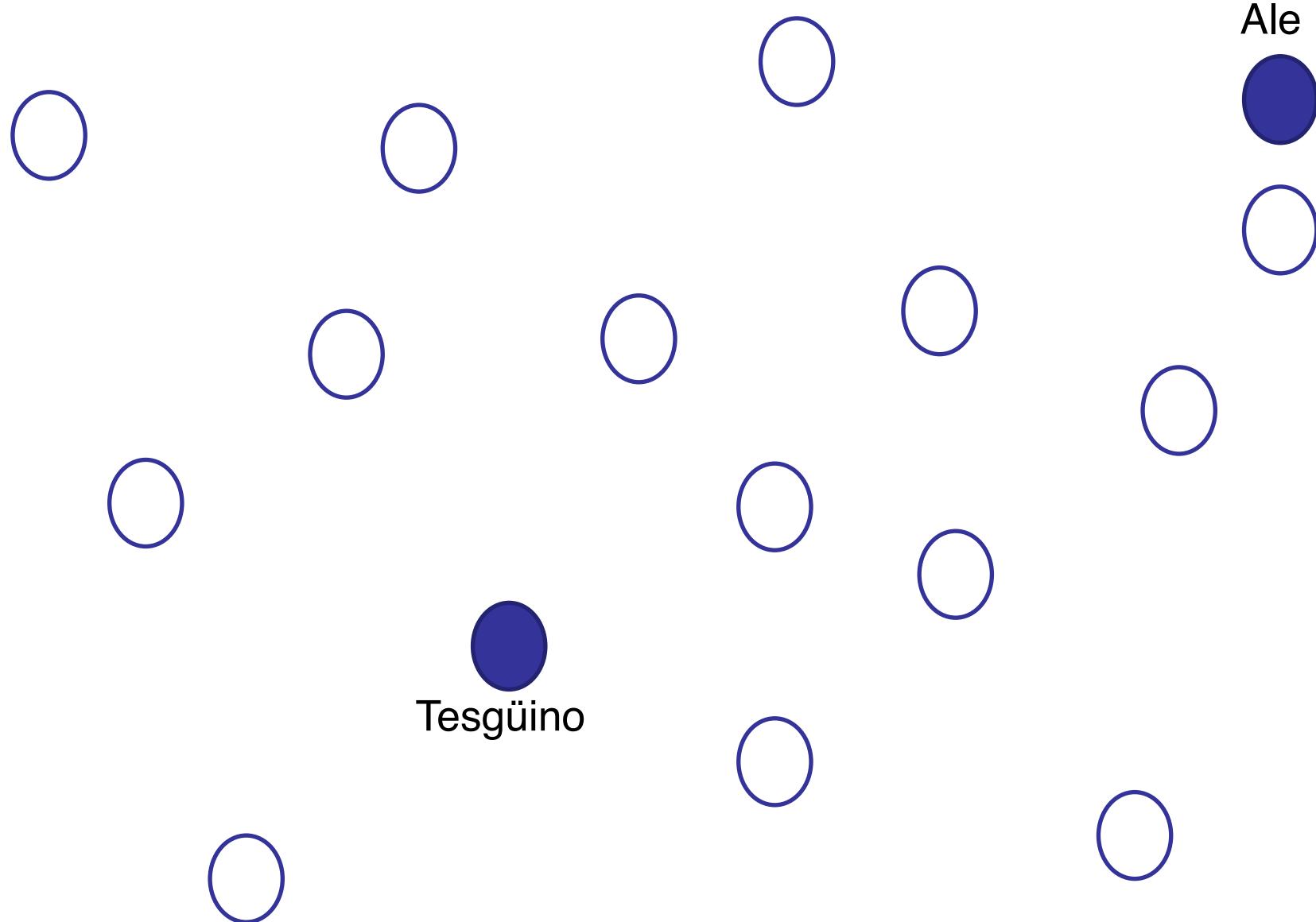
Recap: Neural Word Embedding Architecture





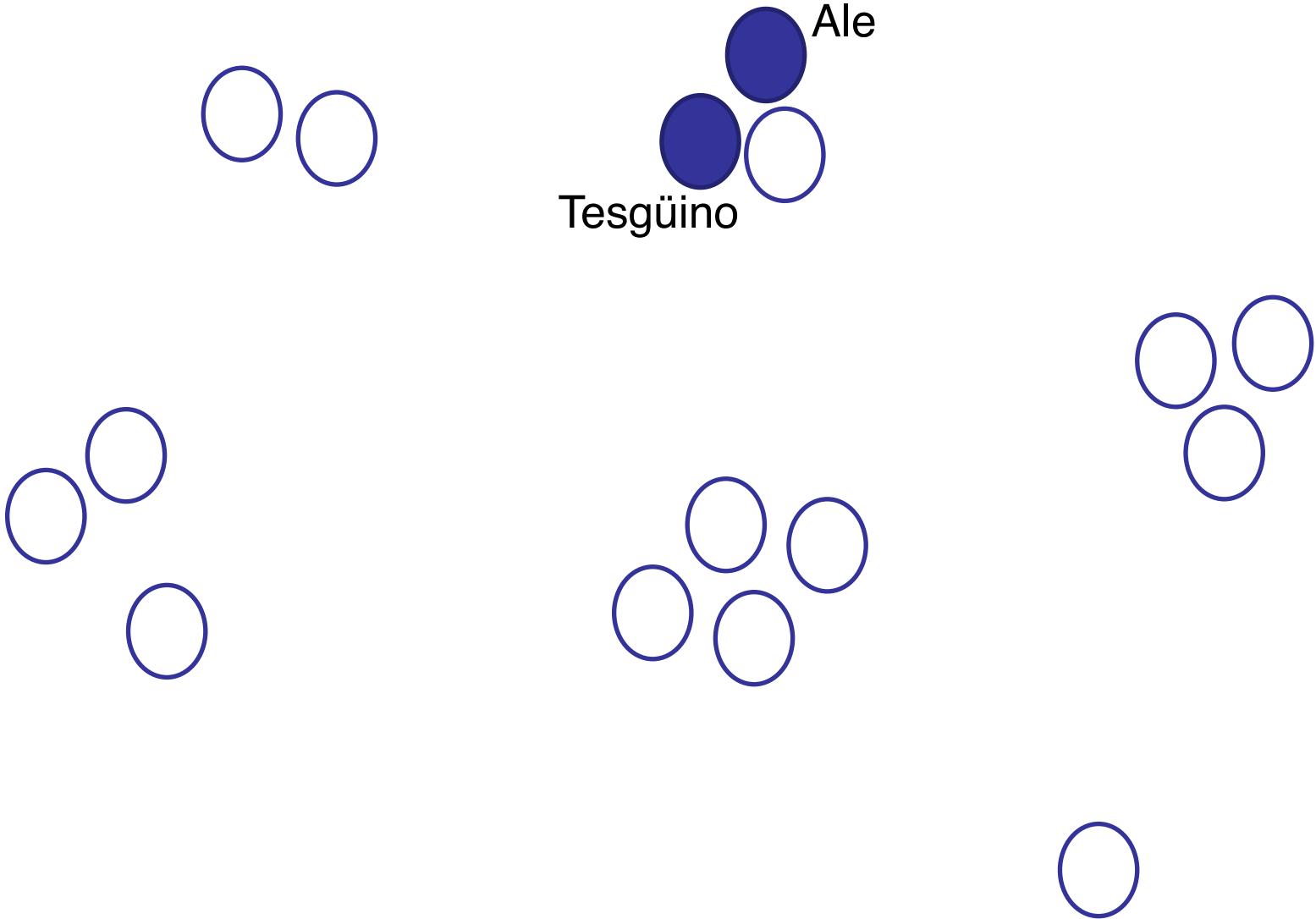
Ale

Tesgüino





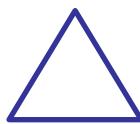
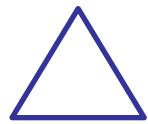
Encoder Embedding



Ale



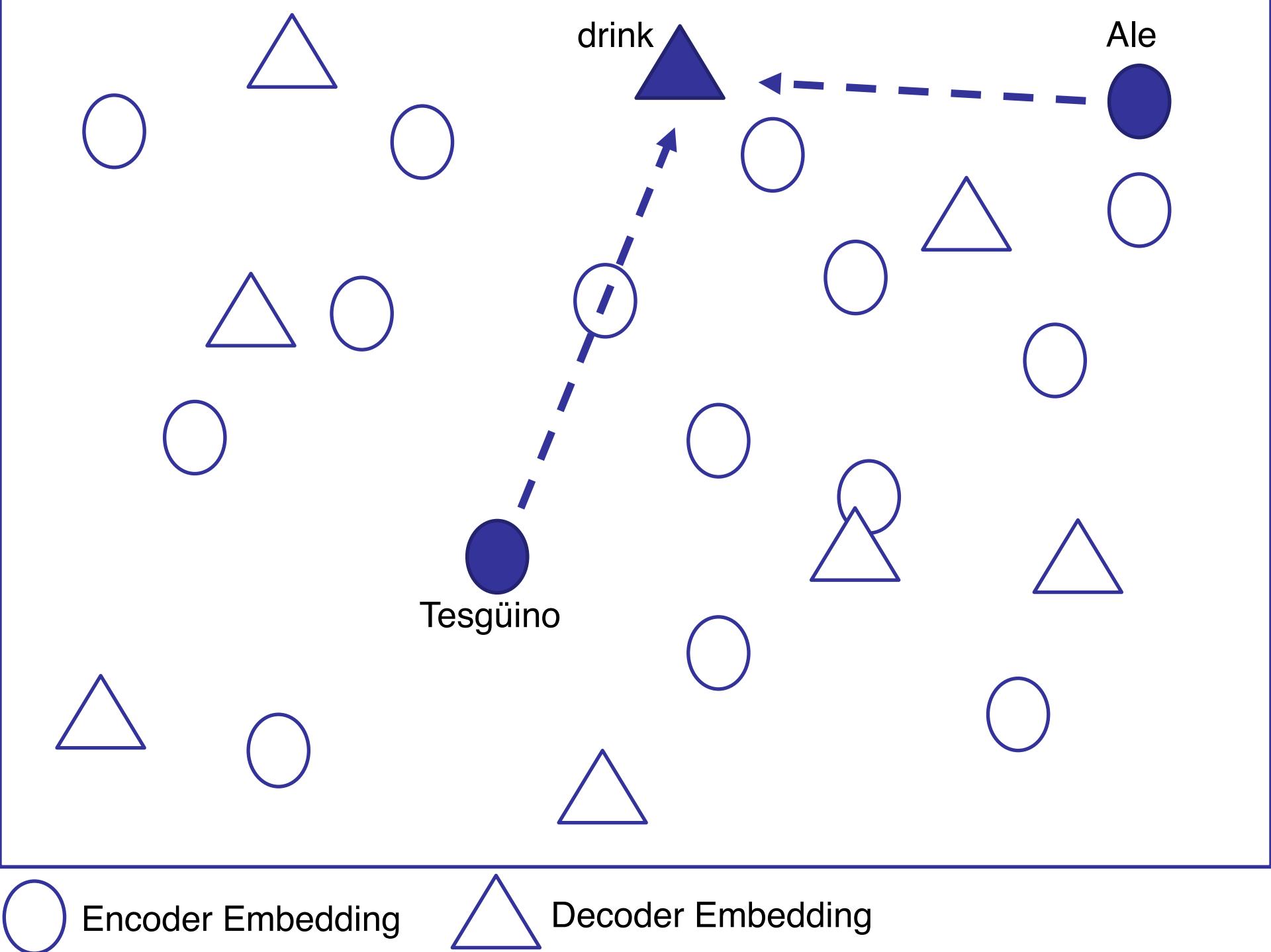
Tesgüino



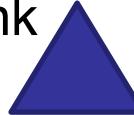
Encoder Embedding



Decoder Embedding

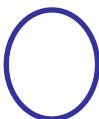
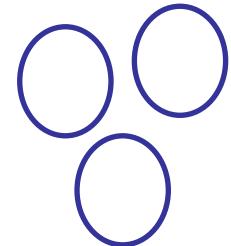
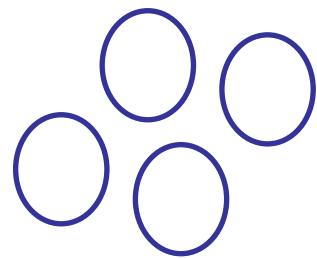
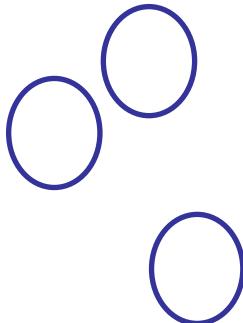


drink

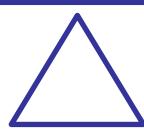


Ale

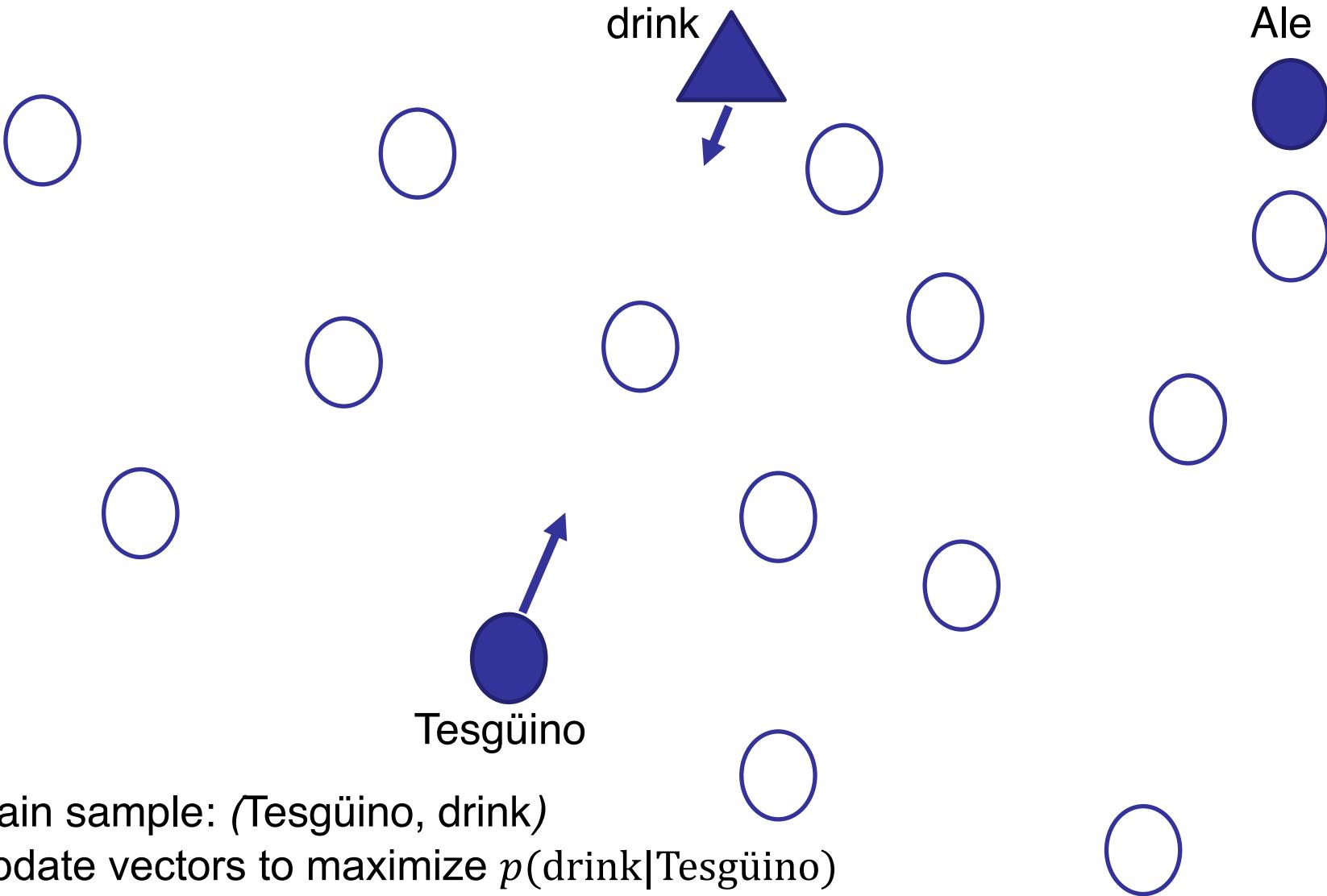
Tesgüino



Encoder Embedding



Decoder Embedding



- Train sample: (Tesgüino, drink)
- Update vectors to maximize $p(\text{drink}|\text{Tesgüino})$

○ Word vector ▲ Context vector

Neural Word Embedding - Summary

- Output value is equal to: $\vec{w}_{\text{Tesgüino}} \cdot \vec{c}_{\text{drink}}$

- Output layer is normalized with **Softmax**

$$p(\text{drink}|\text{Tesgüino}) = \frac{\exp(\vec{w}_{\text{Tesgüino}} \cdot \vec{c}_{\text{drink}})}{\sum_{v \in \mathbb{V}} \exp(\vec{w}_{\text{Tesgüino}} \cdot \vec{c}_v)}$$

\mathbb{V} is the set of vocabularies



Sorry! Denominator is too expensive!

- Loss function is the **Negative Log Likelihood** over all training samples T

$$\mathcal{L} = -\frac{1}{T} \sum_1^T \log p(c|w)$$

word2vec Skip-Gram with Negative Sampling

- word2vec is an **efficient** and **effective** algorithm
- Instead of $p(c|w)$, word2vec measures $p(y = 1|w, c)$: the probability of **genuine co-occurrence** of (w, c)

$$p(y = 1|w, c) = \sigma(\vec{w}_w \cdot \vec{c}_c)$$

↓
sigmoid

- When two words (w, c) appear in the training data, it is a **positive sample**
- word2vec algorithm tries to distinguish between the co-occurrence probability of a **positive sample** from any **negative sample**
- To do it, word2vec draws k **negative samples** by randomly sampling from the words distribution → why randomly?

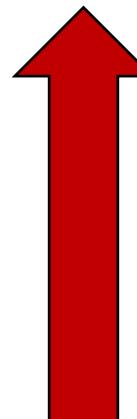
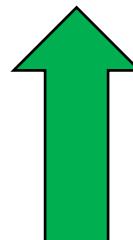
word2vec Skip-Gram with Negative Sampling

- The objective function
 - increases the probability for the **positive sample** (w, c)
 - decreases the probability for the k **negative samples** (w, \check{c})
- Loss function:

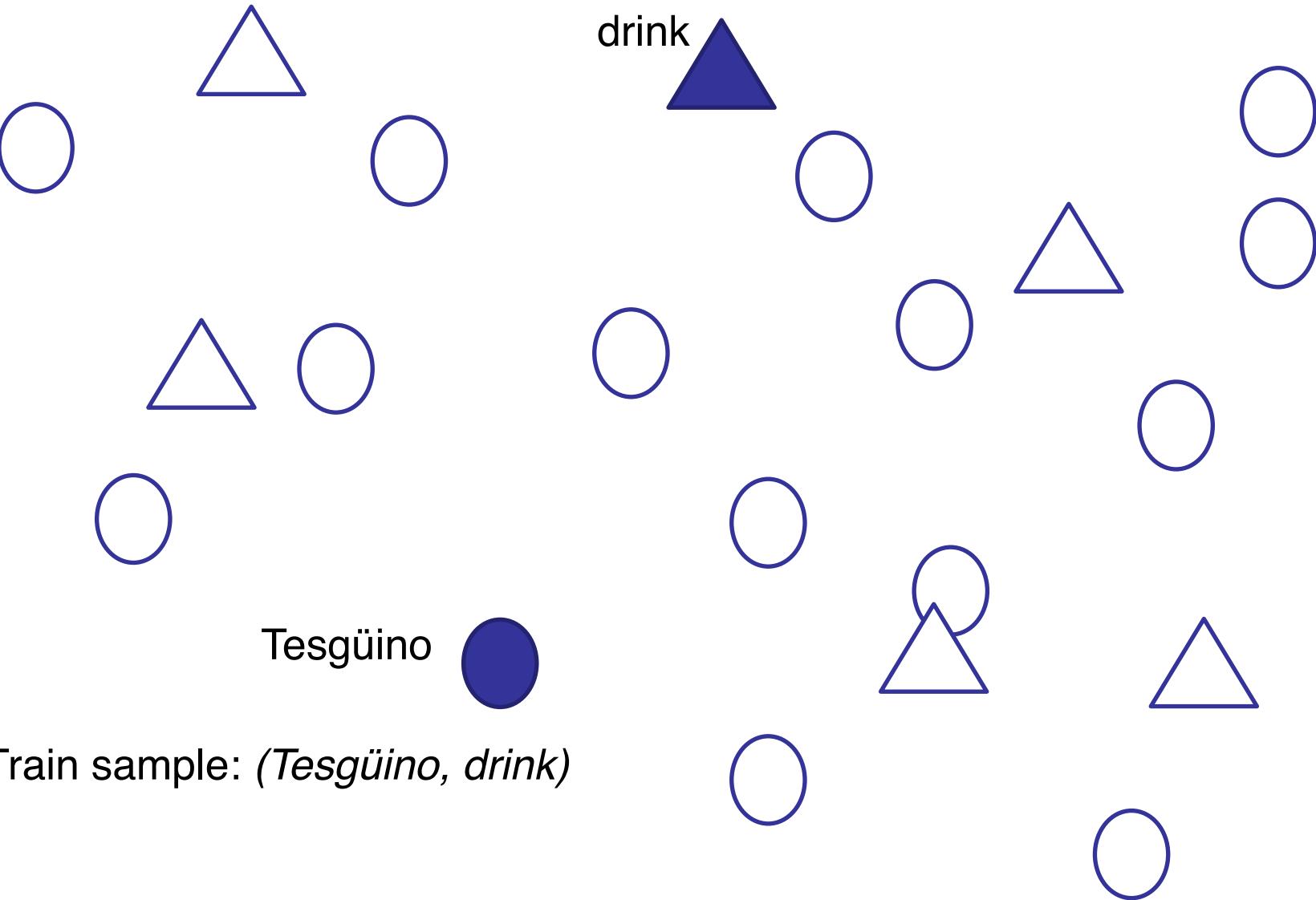
$$L = -\frac{1}{T} \sum_1^T \left[\log p(y = 1|w, c) - \sum_{i=1}^k \log p(y = 1|w, \check{c}) \right]$$

$k \sim 2-10$

Training Samples



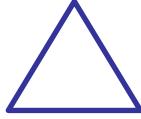
Negative Samples



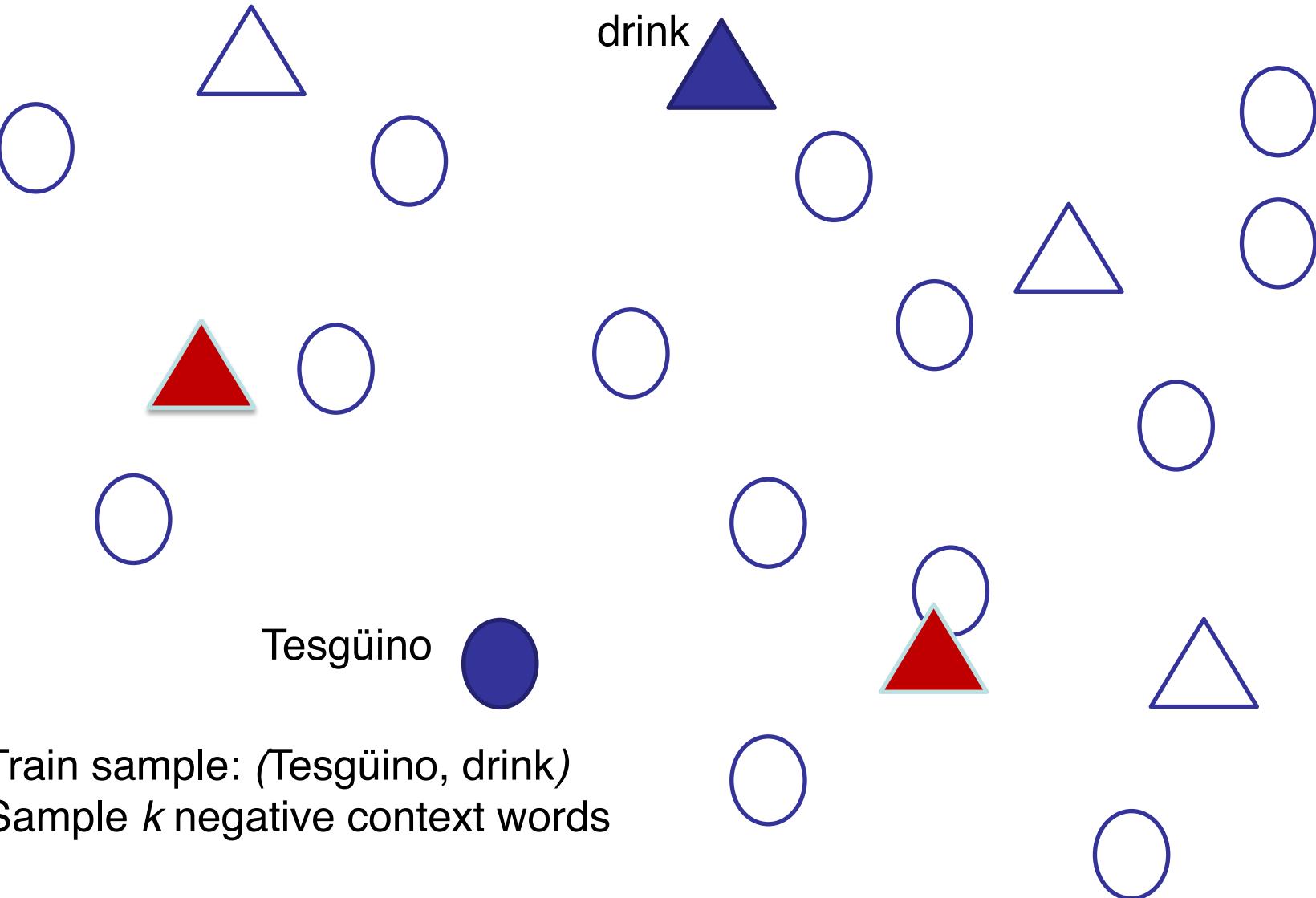
- Train sample: (*Tesgüino*, *drink*)



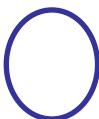
Encoder Embedding



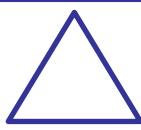
Decoder Embedding



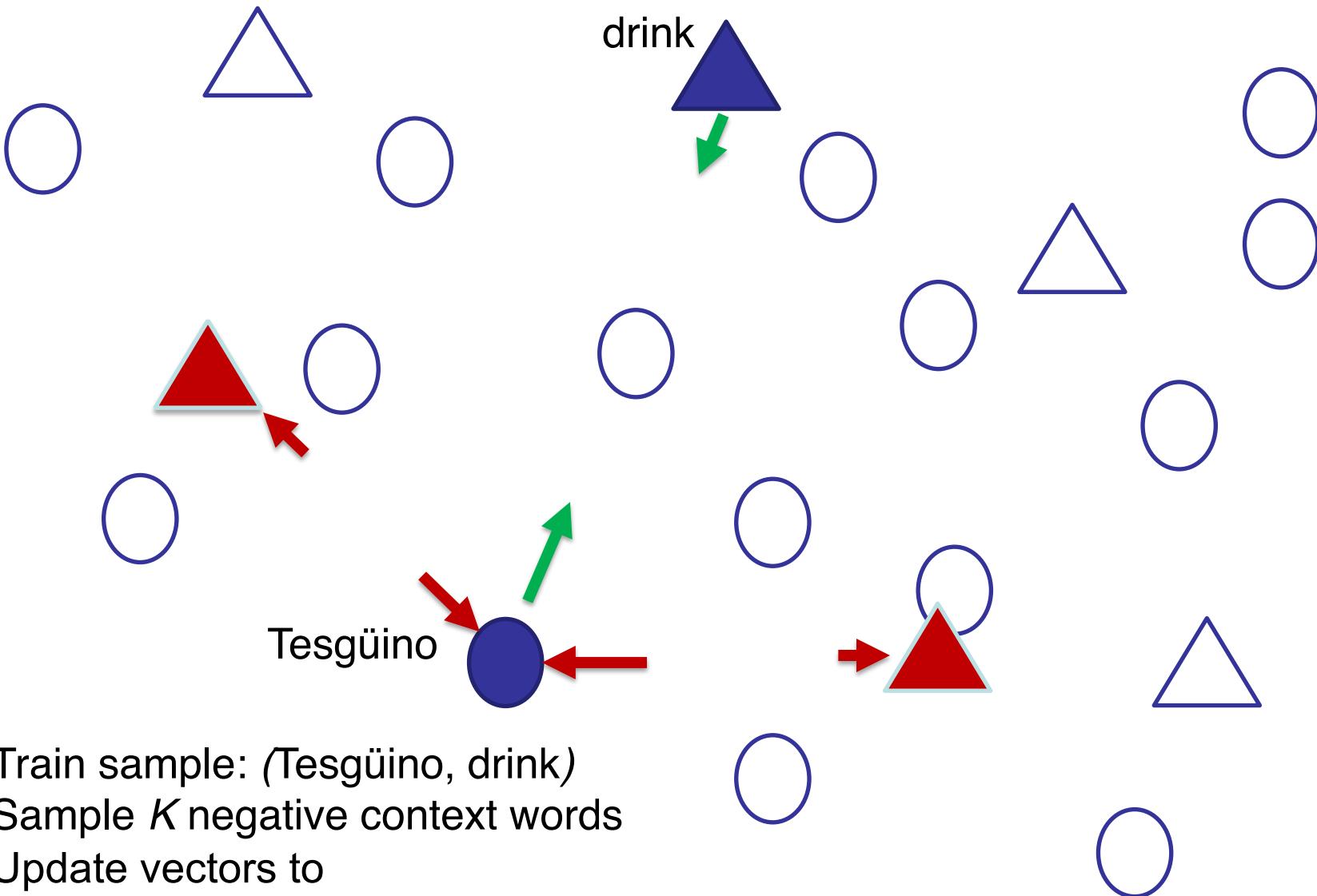
- Train sample: (Tesgüino, drink)
- Sample k negative context words



Encoder Embedding



Decoder Embedding



Encoder Embedding



Decoder Embedding

Discussion about Bias in Data

- A word embedding model captures intrinsic patterns of the given text corpus
- If the data is biased towards a specific demographic, the embedding model also encodes it in the vectors
- This bias in word embeddings can be propagated to down-stream tasks



"I think your test grading is biased in favor of students who answer the test questions correctly."

Bias in Machine Translation



Elaheh Raisi @elaheh_raisi · Oct 3

Bias in google translate from Persian to English 😢 (Persian uses the gender-neutral pronoun)

PERSIAN - DETECTED

ENGLISH



GERMAN

ENGLISH

FRENCH



او مدیر است
او خدمتکار است



26/5000



He is the manager
She is a maid



same gender-neutral pronoun

Bias in Information Retrieval

Search: **nurse**



All

Images

News

Videos

Maps

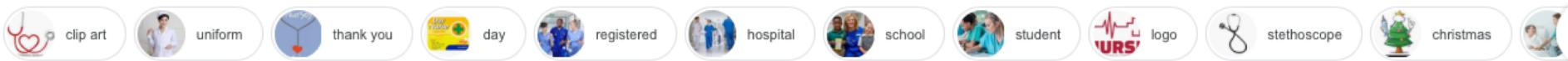
More

Settings

Tools

Collections

SafeSearch ▾



Foreign-Educated Nurse
nurse.org



documenting properly on patient charts
nurse.com



One Nurse at a Time Helps ...
online.nursing.georgetown.edu



Mountain Valleys Health Centers
mtnvalleyhc.org



Nurse Staffing
nursingworld.org



Magic Back to Nursing in 2019
nurse.org



Changes to Nurse Licensure Compact ...
taledmed.com



Nursing Jobs
nursingworld.org



Practical Nurse Program - A...
abccott.edu



Auto Insurance for Nurses from ...
mycalcas.com



An Irish nurse: 'I began to hate my job ...
thejournal.ie



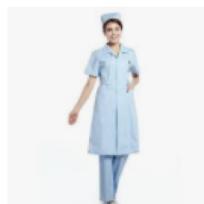
Nurse interview questions and answers ...
snagajob.com



an Emergency Room Nurse
americanr.com



Men in Nursing Archives - Minority Nurse
minoritynurse.com



Hospital Nurse Uniform at Rs...
indiamart.com



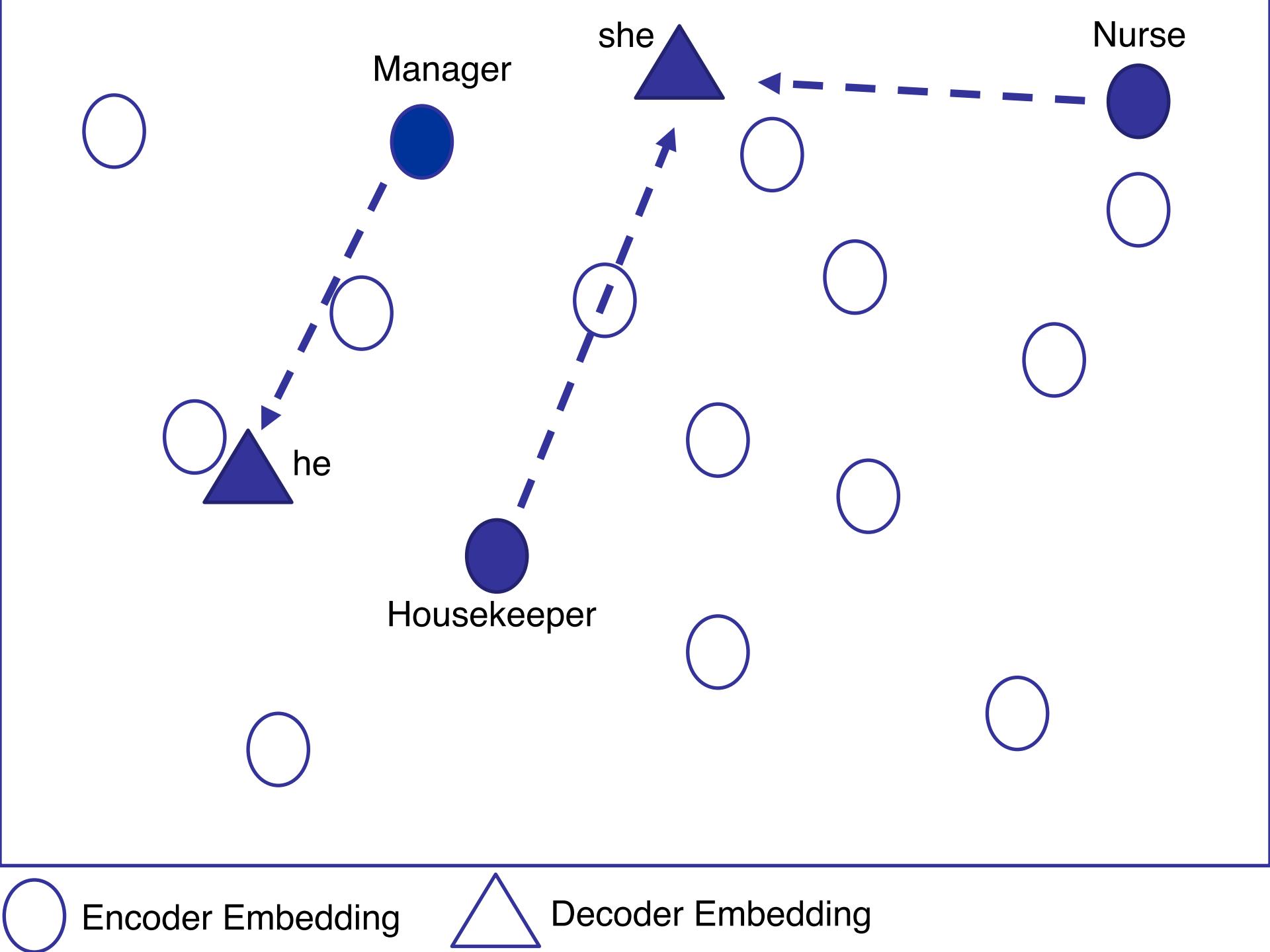
Johnson Looking for Nurse Innovators ...
elitemca.com

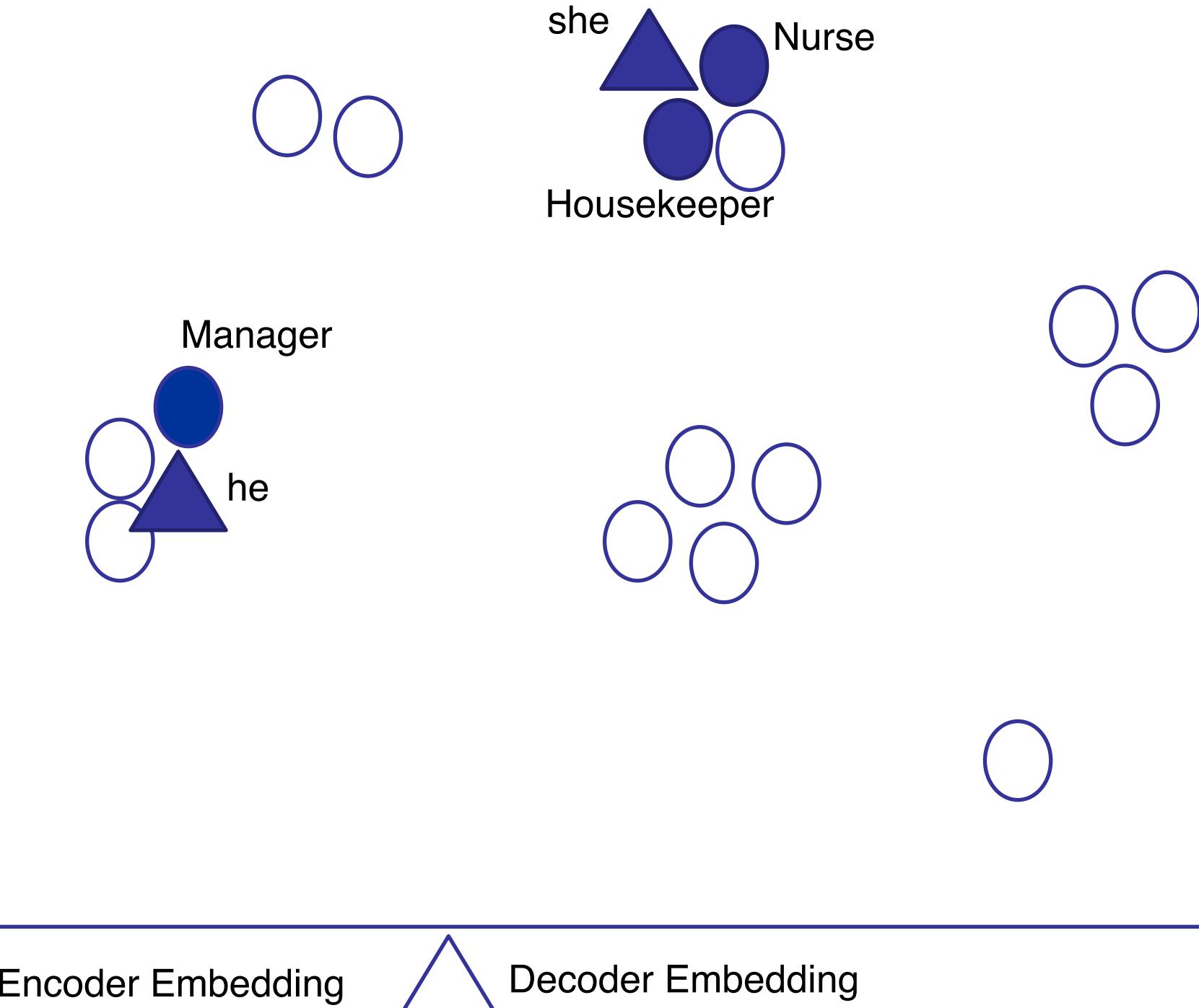


IoT vendors. When a paediatric nurse ...
theranister.co.uk



know if Nursing is Your Dream Job ...
advanced-care.us





Gender Bias in Wikipedia

- The bias of 500 occupations to female/male in the word2vec Skip-Gram model, created using an English Wikipedia corpus



Agenda

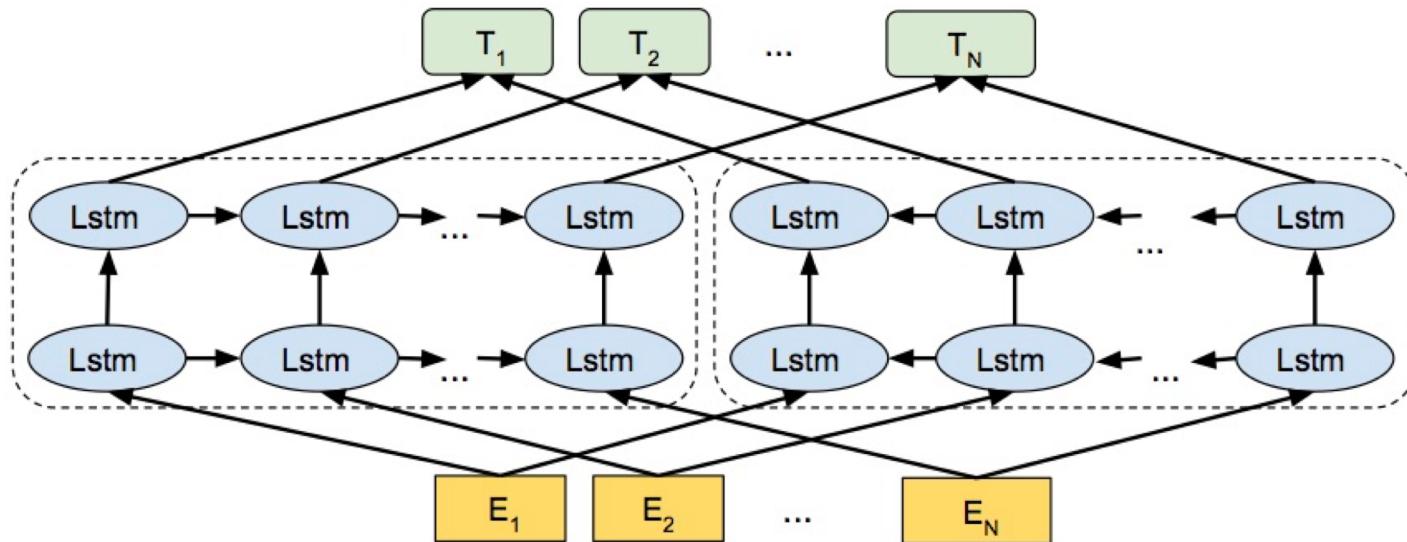
- Neural Networks - Recap
- Language Modeling
- Neural Word Embedding - Revisited
 - word2vec Skip-Gram with Negative Sampling
 - Bias in word embeddings
- Contextualized Word Embeddings
 - ELMo
 - BERT

Contextualized Word Embedding

- What is *Apple*?!
 - a fruit or a company?
- The meanings of words become more clear, when we read them in their contexts
 - *eating an apple* vs. *share of the Apple company*
- In the ‘standard’ word embeddings, various meanings (senses) of a word are all encoded in one vector
- In **contextualized word embedding**, the representation of each word is defined in its context
- The input is a **set of words** (like a sentence) and the output is a **set of vectors**, each vector representing a word

ELMo (Embeddings from Language Models)

- ELMo is a Multi-layer Bidirectional LSTM



- Training is similar to training a Language Model
 - The given sentence is read one time from left to right and one time from right to left (bidirectional)
 - In each step, the model inputs a word and predicts its next word using the hidden state of the previous step

ELMo

- In a supervised task (such as document classification), the final representations are the weighted sum of the hidden states of the intermediary RNN layers
- For an ELMo with L layers, the representation of the word w is:

$$\text{ELMo}_w = \gamma \sum_{j=1}^L \theta_j \vec{s}_w^j$$

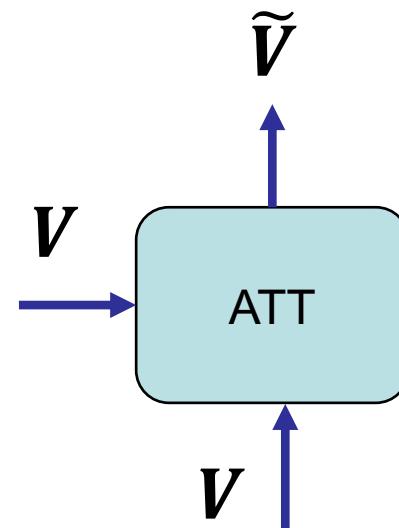
where γ and θ are parameters, learned during the task

Self-Attention Networks

- The query matrix is the same as the value matrix
- Self-attention outputs a new matrix \tilde{V} , such that each vector of V attends to all other vectors V
- In contrast to RNNs, all the input vectors are fed to the model at the same time and not one by one

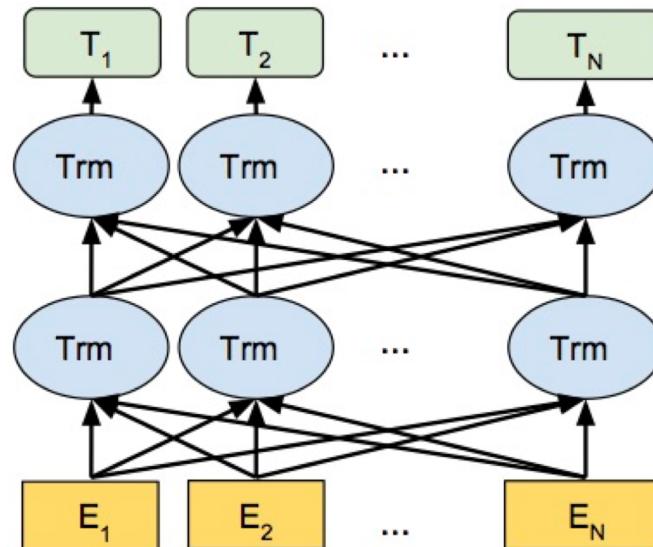
- Transformers are a special architecture of (self-)attention networks
 - Reference for more details:

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. "**Attention is all you need.**" In *Advances in neural information processing systems*,



BERT (Bidirectional Encoder Representation)

- Multi-layer Transformer network



- For training, for each sentence, the model masks some of the input words and tries to predict the masked ones

Input: The man went to the [MASK]₁ . He bought a [MASK]₂ of milk .

Labels: [MASK]₁ = store; [MASK]₂ = gallon

BERT

- BERT is quite big!
 - **BERT-Base:** 12-layer, 110M parameters
 - **BERT-Large:** 24-layer, 340M parameters
- Though the results of fine-tuning in down-stream tasks such as Information Retrieval, Natural Language Inference, Question/Answering, Document classification, etc. are impressive
- Topics for further learning:
 - Learning sentence relations
 - Positional embeddings
 - Fine-tuning procedure