



Universidade Federal de Alagoas
Instituto de Computação

Doc2Vec

Prof. Thales Vieira

Word Vectors framework

Parecido com o CBOW, aqui tentamos prever a próxima palavra da frase.

$$p(w_t | w_{t-k}, \dots, w_{t+k}) = \frac{e^{y_{w_t}}}{\sum_i e^{y_i}}$$

$$y = b + Uh(w_{t-k}, \dots, w_{t+k}; W) \quad (1)$$

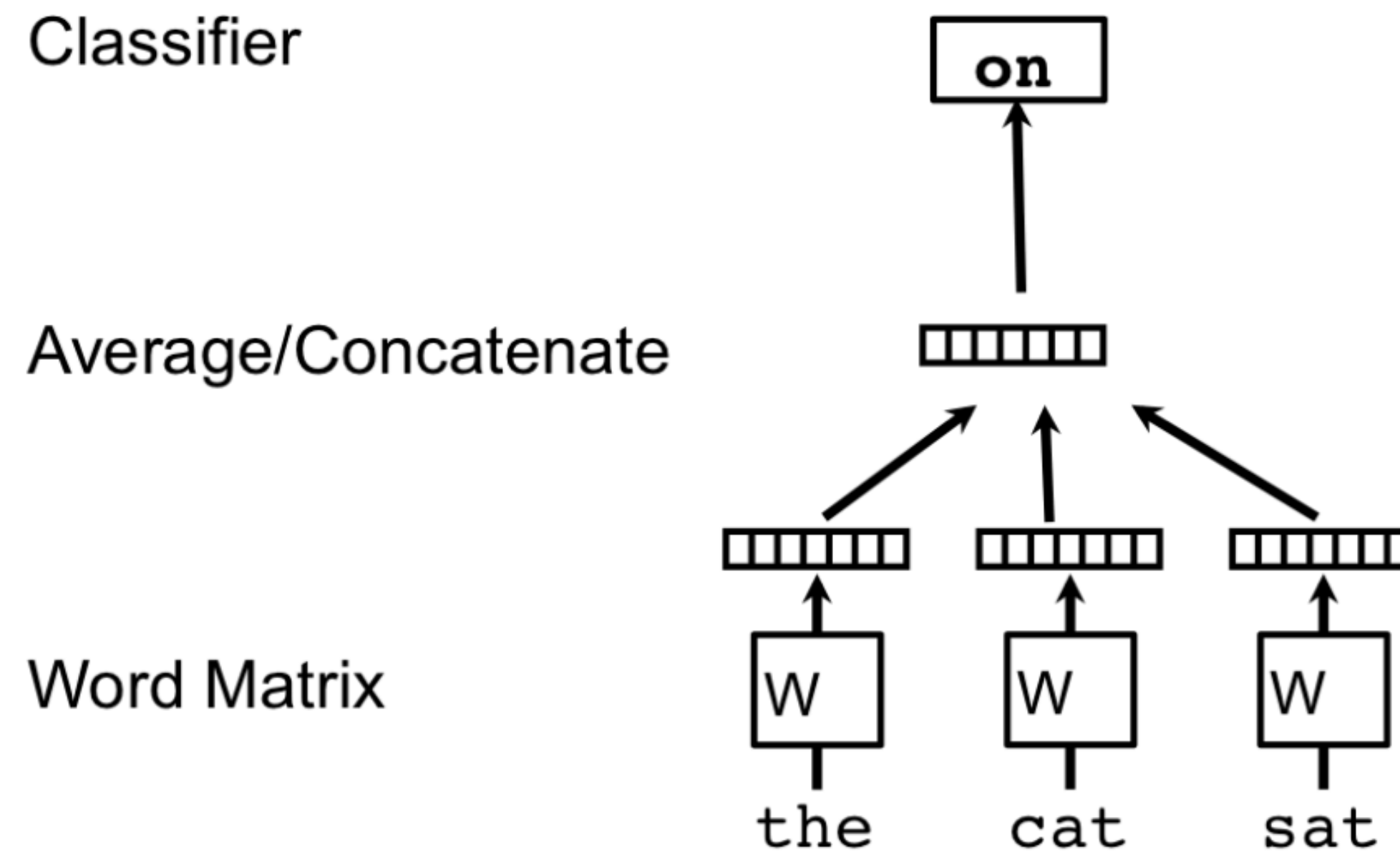
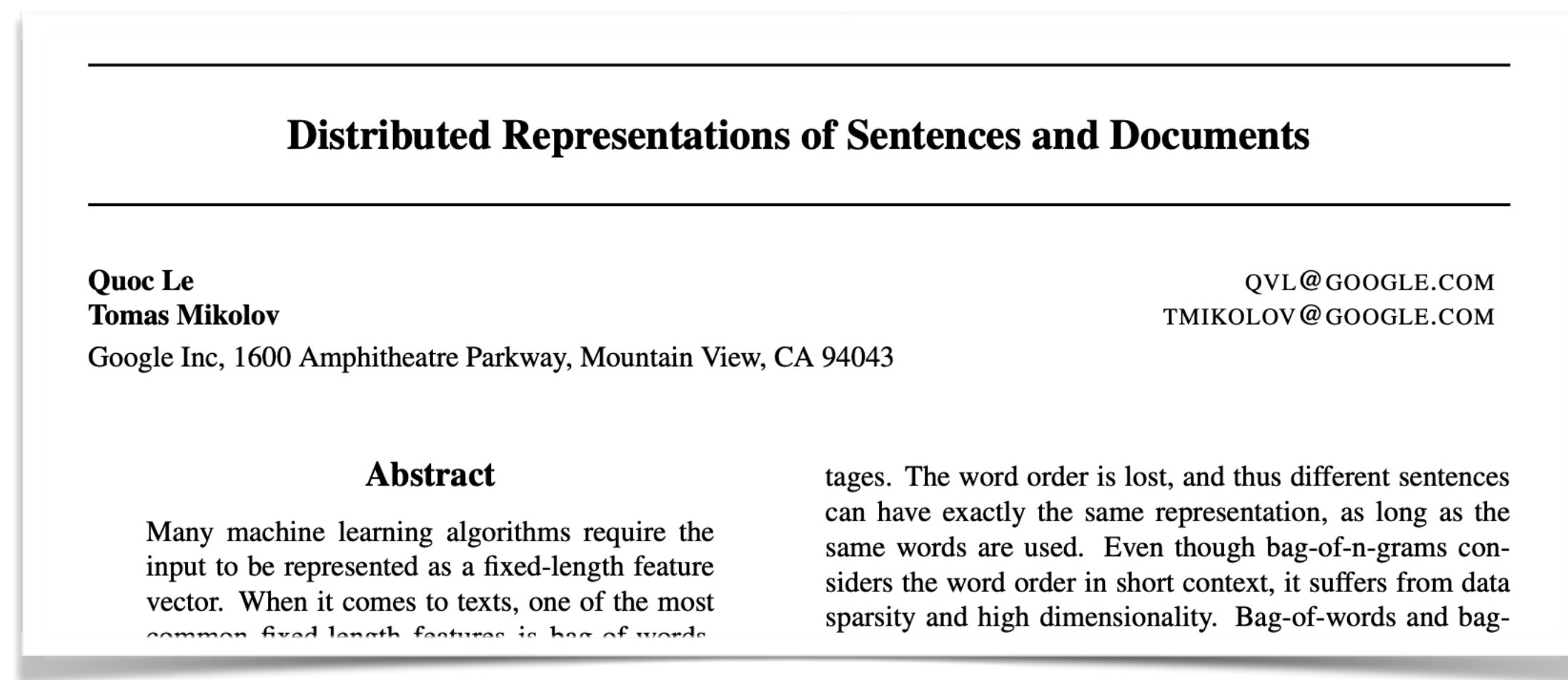


Figure 1. A framework for learning word vectors. Context of three words (“the,” “cat,” and “sat”) is used to predict the fourth word (“on”). The input words are mapped to columns of the matrix W to predict the output word.

Doc2Vec framework



"In our Paragraph Vector framework (see Figure 2), every paragraph is mapped to a unique vector, represented by a column in matrix **D** and every word is also mapped to a unique vector, represented by a column in matrix **W**. The paragraph vector and word vectors are averaged or concatenated to predict the next word in a context. In the experiments, we use concatenation as the method to combine the vectors."

"More formally, the only change in this model compared to the word vector framework is in equation 1, where h is constructed from W and D ."

$$y = b + Uh(w_{t-k}, \dots, w_{t+k}; W) \quad (1)$$

Classifier

Average/Concatenate

Paragraph Matrix----->

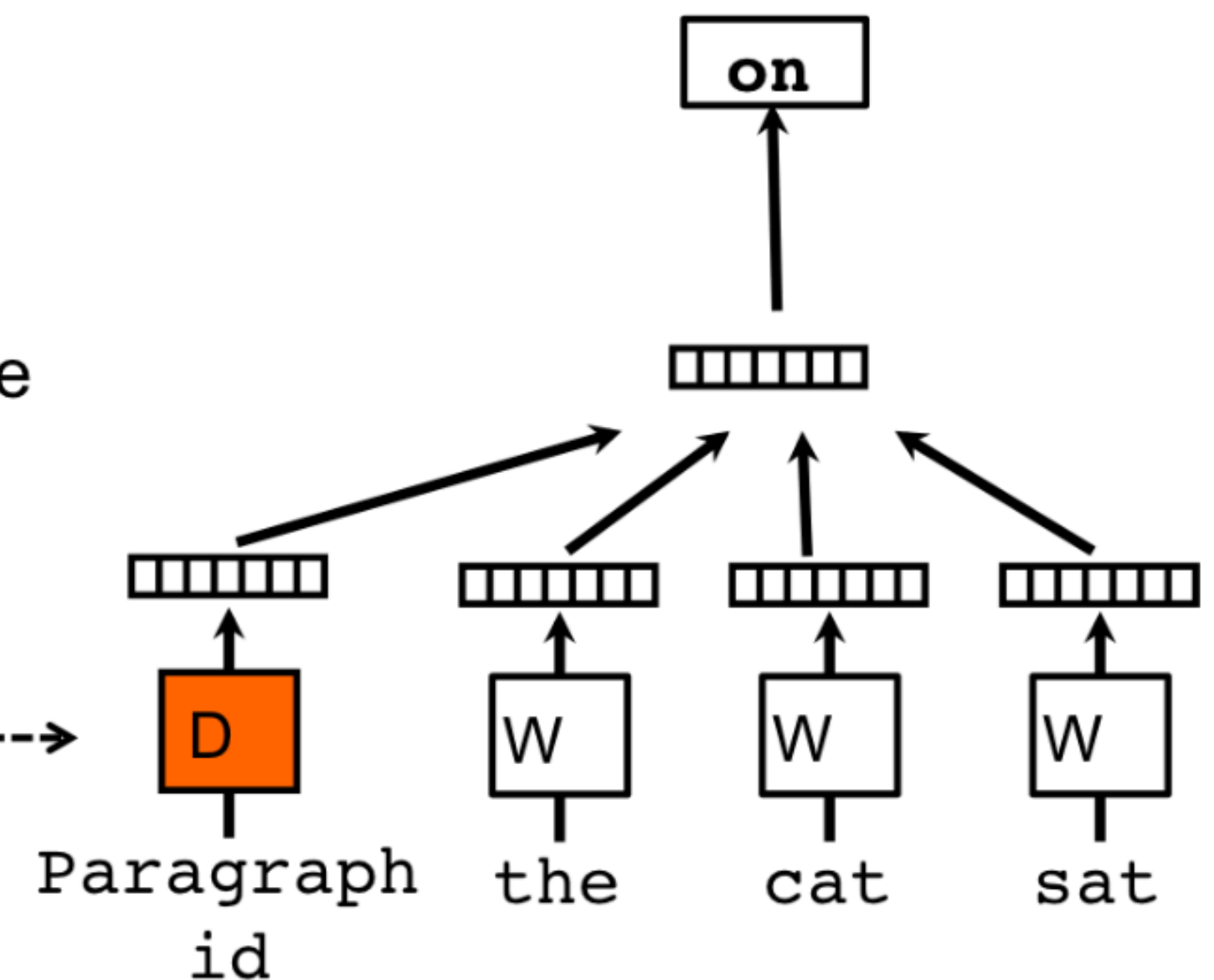


Figure 2. A framework for learning paragraph vector. This framework is similar to the framework presented in Figure 1; the only change is the additional paragraph token that is mapped to a vector via matrix D . In this model, the concatenation or average of this vector with a context of three words is used to predict the fourth word. The paragraph vector represents the missing information from the current context and can act as a memory of the topic of the paragraph.

Doc2Vec framework

Distributed Representations of Sentences and Documents

Quoc Le
Tomas Mikolov

Google Inc, 1600 Amphitheatre Parkway, Mountain View, CA 94043

QVL@GOOGLE.COM
TMIKOLOV@GOOGLE.COM

Abstract

Many machine learning algorithms require the input to be represented as a fixed-length feature vector. When it comes to texts, one of the most common fixed-length features is bag-of-words

tages. The word order is lost, and thus different sentences can have exactly the same representation, as long as the same words are used. Even though bag-of-n-grams considers the word order in short context, it suffers from data sparsity and high dimensionality. Bag-of-words and bag-

"The contexts are fixed-length and sampled from a sliding window over the paragraph. The paragraph vector is shared across all contexts generated from the same paragraph but not across paragraphs. The word vector matrix W , however, is shared across paragraphs. I.e., the vector for “powerful” is the same for all paragraphs.”

"The paragraph vectors and word vectors are trained using stochastic gradient descent and the gradient is obtained via backpropagation. At every step of stochastic gradient descent, one can sample a fixed-length context from a random paragraph, compute the error gradient from the network in Figure 2 and use the gradient to update the parameters in our model.”

Classifier

Average/Concatenate

Paragraph Matrix----->

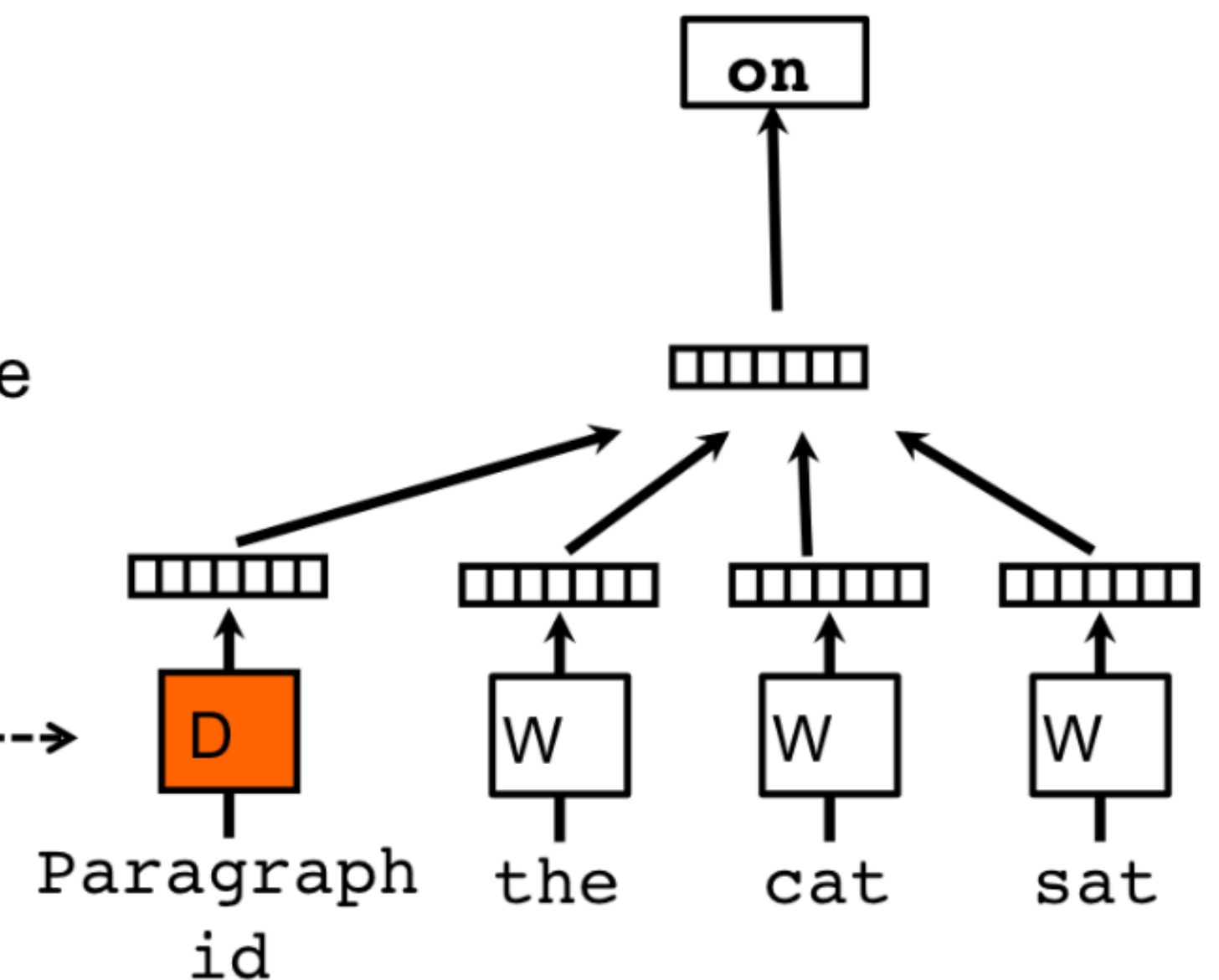
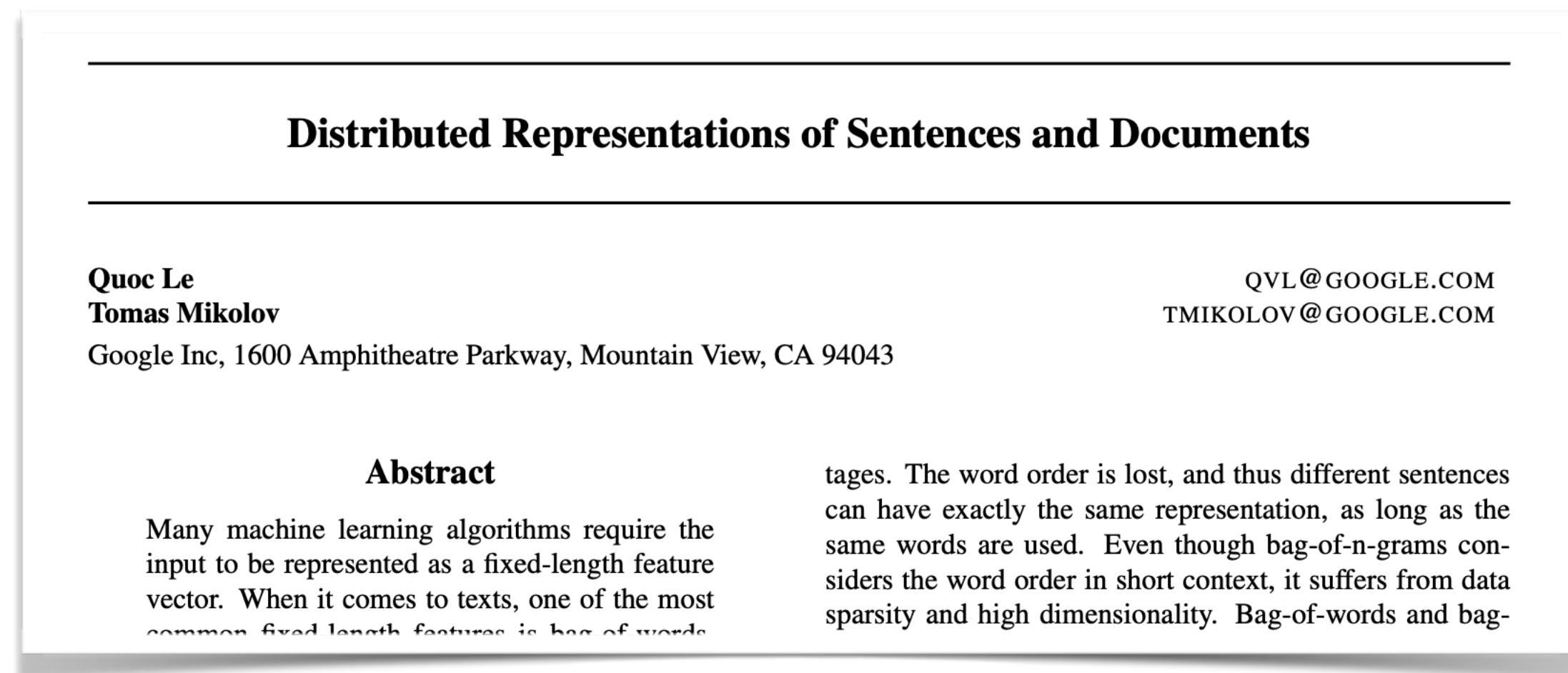


Figure 2. A framework for learning paragraph vector. This framework is similar to the framework presented in Figure 1; the only change is the additional paragraph token that is mapped to a vector via matrix D . In this model, the concatenation or average of this vector with a context of three words is used to predict the fourth word. The paragraph vector represents the missing information from the current context and can act as a memory of the topic of the paragraph.

Doc2Vec framework



"At prediction time, one needs to perform an inference step to compute the paragraph vector for a new paragraph. This is also obtained by gradient descent. In this step, the parameters for the rest of the model, the word vectors W and the softmax weights, are fixed."

"In summary, the algorithm itself has two key stages: 1) training to get word vectors W , softmax weights U , b and paragraph vectors D on already seen paragraphs; and 2) "the inference stage" to get paragraph vectors D for new paragraphs (never seen before) by adding more columns in D and gradient descending on D while holding W , U , b fixed. We use D to make a prediction about some particular labels using a standard classifier, e.g., logistic regression or SVM."

Classifier

Average/Concatenate

Paragraph Matrix----->

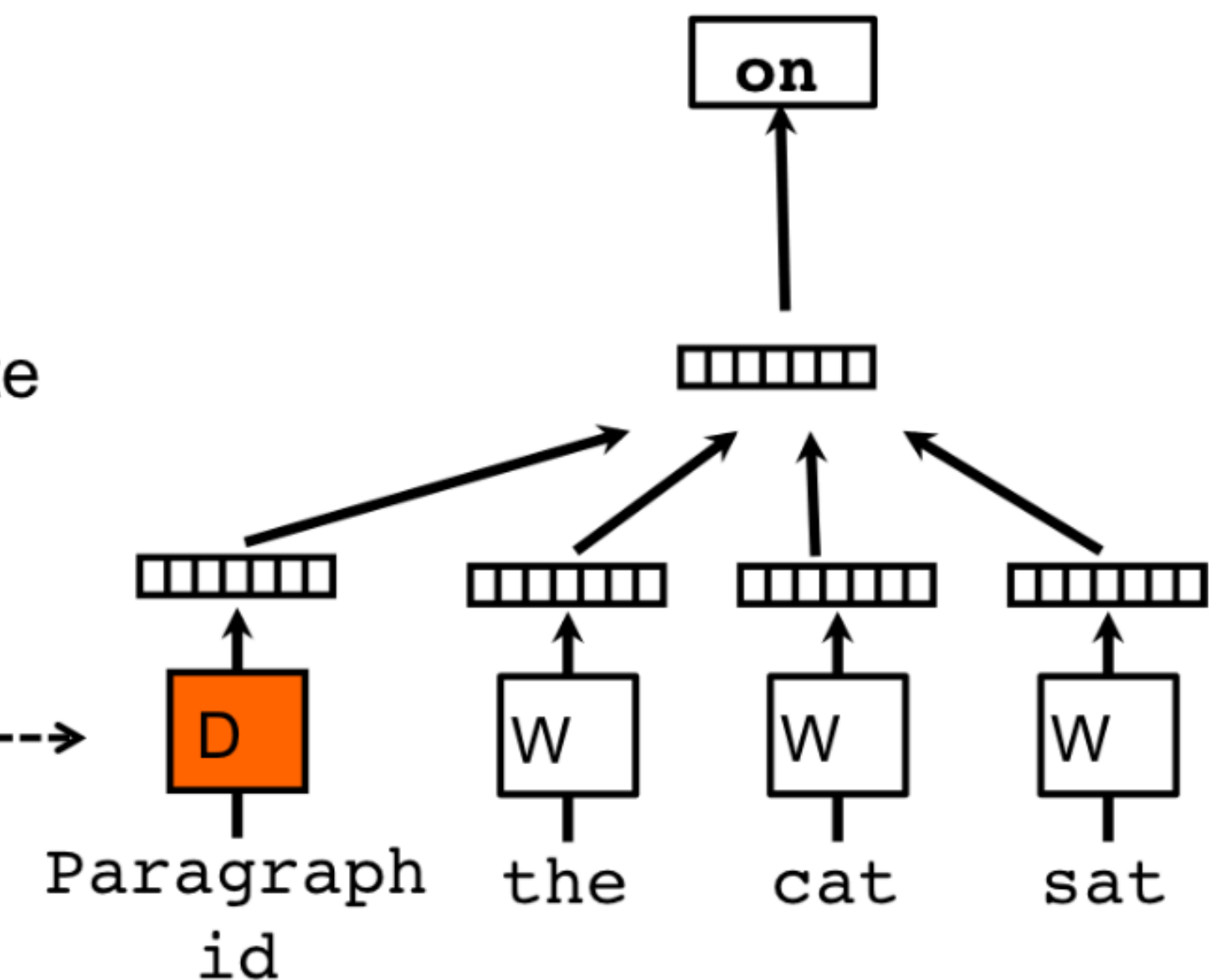


Figure 2. A framework for learning paragraph vector. This framework is similar to the framework presented in Figure 1; the only change is the additional paragraph token that is mapped to a vector via matrix D . In this model, the concatenation or average of this vector with a context of three words is used to predict the fourth word. The paragraph vector represents the missing information from the current context and can act as a memory of the topic of the paragraph.