

# One-shot Meta-learning with Matching Networks



Adrian Gonzalez-Martin<sup>1</sup>, Matthew Lee<sup>1</sup>, Adrian Daniel Szwarc<sup>1</sup>, Talip Uçar<sup>1</sup>

Computer Science Department, University College London, Gower Street, London, WC1E 6BT

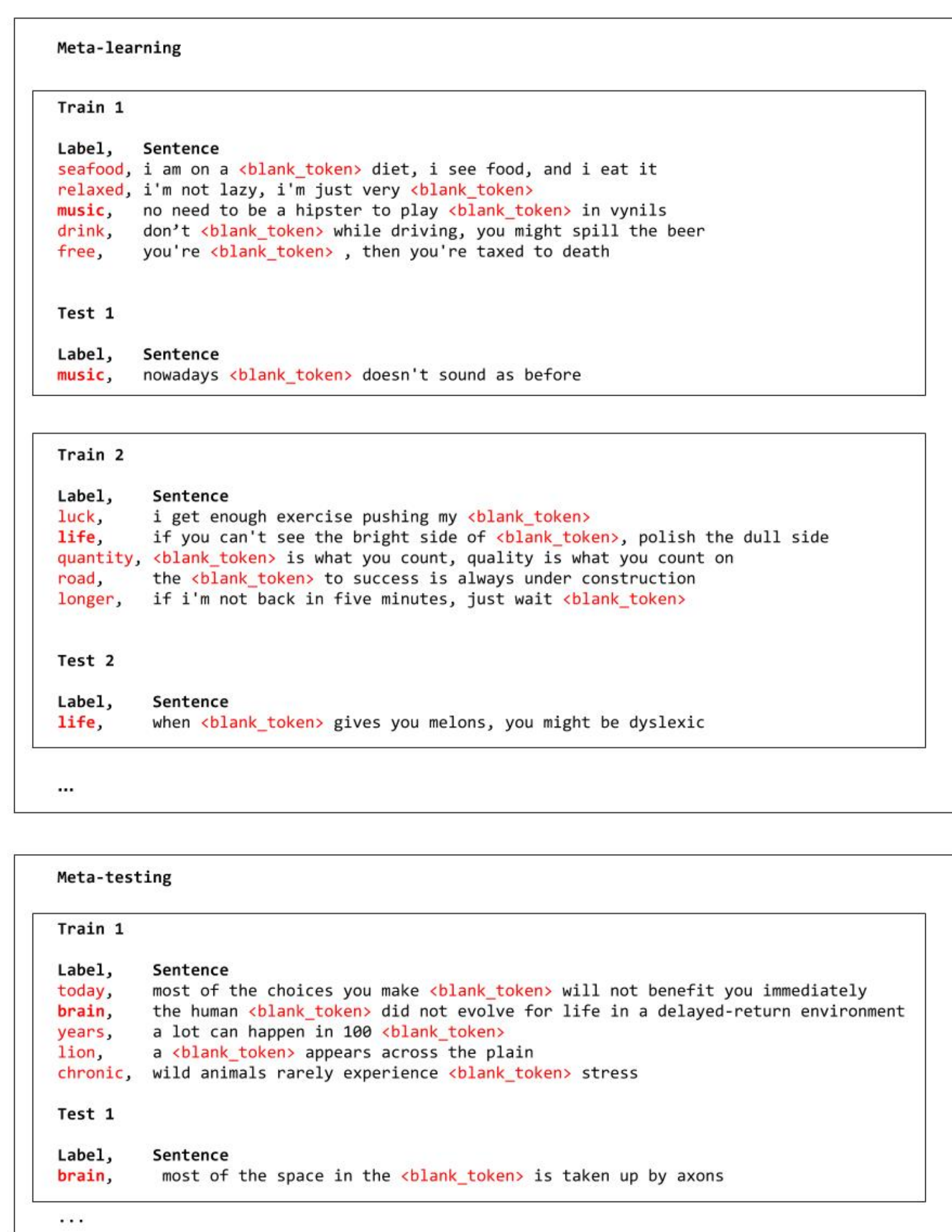
talip.ucar.16@ucl.ac.uk

## 1. Introduction

One-shot learning is learning from one example. It is an important property of natural intelligence. Humans can often learn new concepts using one, or few examples. With advances in areas of machine learning concerning memory, attention and various optimization techniques, there has been a renewed interest in one-shot learning approaches in recent years. But, it has been mostly investigated in the context of vision tasks, rather than language ones. Thus, our work focuses on one-shot learning in a language task, particularly predicting missing words in sentences. We aim to establish a benchmark by using the WikiText-2 [1] dataset, and comparing different one-shot approaches. We use Matching Networks [2] to find our baseline, and then improve on it, using new embeddings and distance metrics such as cosine, Euclidean, and Poincaré.

## 2. Meta-learning Framework

The success of recent one-shot learning approaches has not been purely down to advances in neural network architectures but also clever meta-learning frameworks.



## 3. Preprocessing

"The fox is sick of jumping over dogs."

-> Lowercase, replace numbers with "N", remove punctuation  
"the fox is sick of jumping over dogs"

-> Choose missing words and split into labels and sentences  
"sick, the fox is <blank\_token> of jumping over dogs"

-> Tokenization  
["sick", "the", "fox", "is", "<blank\_token>", "jumping", "over", "dogs"]

-> (vocab\*)

{"the" -> 0,	-> (numericalise)	-> (basic
"fox"-> 1,	[0, 1, 2, ...]	embedding)
"is" -> 2,		[
...}		[1, 0, 0, ...],
		[0, 1, 0, ...],
		[0, 0, 1, ...]
		]

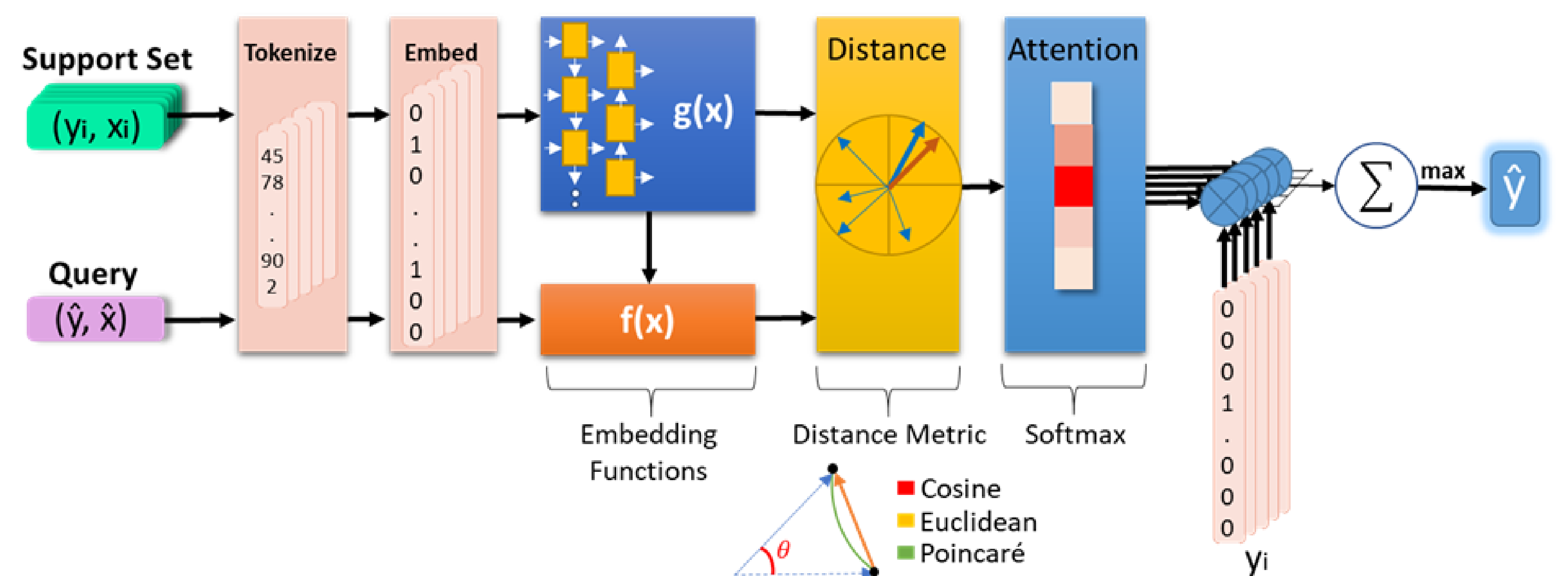
\* can have separate vocabs for labels and sentences

## 6. References

- [1] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *CoRR*, abs/1609.07843, 2016.
- [2] Oriol Vinyals, Charles Blundell, Timothy Lillicrap, koray kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 3630–3638. Curran Associates, Inc., 2016.

## 4. Matching Network Architecture

Embeddings + Similarity Measure + Attention



## 5. Embedding, Attention and Distance Metrics

### Embedding: Full Context Embeddings

- $g(x)$  employs a bidirectional LSTM, with the support set  $S$  treated as a sequence. The  $i$ th stage of the LSTM has as input  $g'(x_i)$  - the initial embedding of the previous step. The LSTM is also a function of the last hidden state  $\vec{h}_{i-1}$  and last output vector  $\vec{c}_{i-1}$  as usual.

$$\vec{h}_i, \vec{c}_i = \text{LSTM}(g'(x_i), \vec{h}_{i-1}, \vec{c}_{i-1})$$

$$\overleftarrow{h}_i, \overleftarrow{c}_i = \text{LSTM}(g'(x_i), \overleftarrow{h}_{i-1}, \overleftarrow{c}_{i-1})$$

- $f(\hat{x})$  is a one way LSTM with an attention mechanism that unfolds over  $K$  steps:

$$f(\hat{x}) = \text{attLSTM}(f'(\hat{x}), g(S), K)$$

where like  $g'$ ,  $f'$  is our initial embedding but for the query;  $g(S)$  represents  $g(x_i)$  already applied to all  $x_i$  and  $K$  is the number of steps as noted.

### Attention:

- The attention kernel applies the cosine distance,  $c(\cdot)$ , between the embedded vectors of two data points before the softmax function normalizes the resulting weights:

$$a(\hat{x}, x_i) = \frac{e^{c(f(\hat{x}), g(x_i))}}{\sum_{j=1}^k e^{c(f(\hat{x}), g(x_j))}}$$

### Distance Metrics:

- Investigating three distance metrics: **Cosine**, **Euclidean**, and **Poincaré**.

