

BPEmb: Tokenization-free Pre-trained Subword Embeddings in 275 Languages

Benjamin Heinzerling^{1,2} Michael Strube²

¹AIPHES ²Heidelberg Institute for Theoretical Studies



Heidelberg Institute for Theoretical Studies



Myxomatosis

-osis

“state, abnormal condition, or action”

-oma / -omato

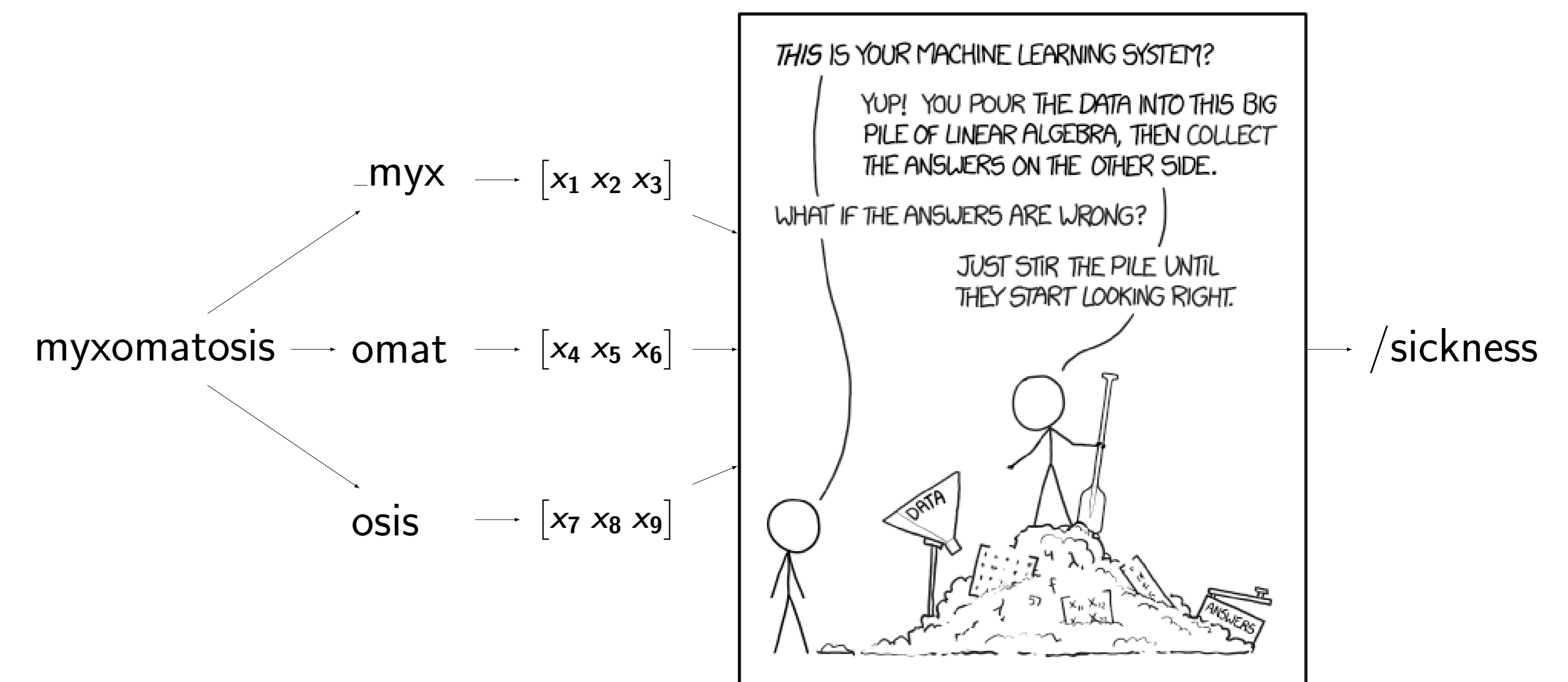
Forming nouns indicating tumors or masses

myxo

From Ancient Greek *múxa*, “mucus”



Subword-based Entity Typing



Computational Approximations to Morphological Analysis

1. Split into subwords

characters: *m, y, x, o, m, a, t, o, s, i, s*

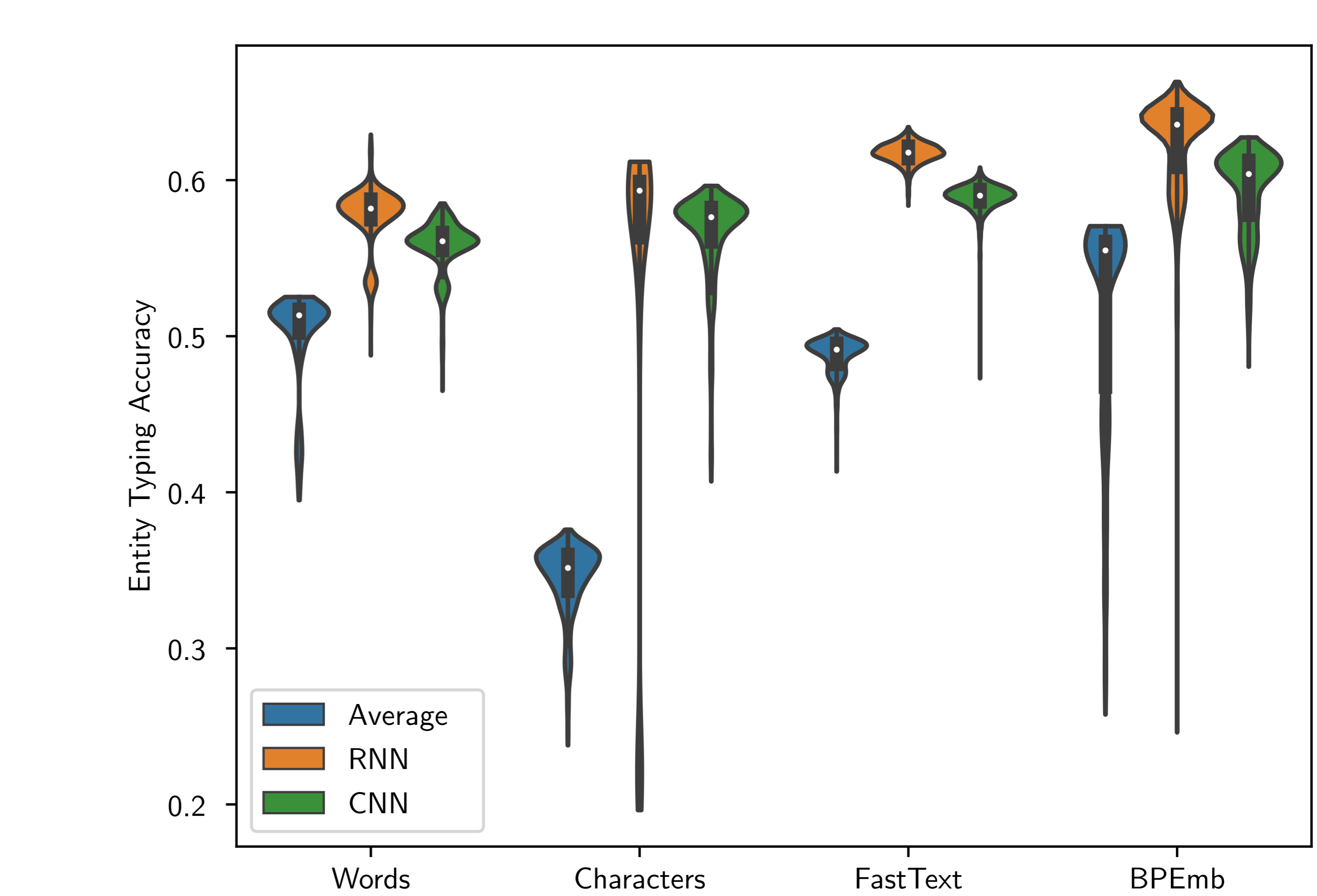
ngrams, e.g.: *myx, yxo, xom, oma, mat, ato, tos, osi, sis*

FastText: *myx + yxo + xom + oma + mat + ato + tos + osi + sis + myxo + yxom + xoma + omat + mato + atos + tosi + osis + myxom + ... + xomato + omatos + matosi + atosis*

byte pairs

2. Learn function that infers word meaning from subwords

Results: English Entity Typing



Byte-Pair Encoding (BPE) (Gage, 1994)

A B A B C A B C D

Most frequent pair: A B

Merge pair A B to X:

X X C X C D

Most frequent pair: X C

Merge pair X C to Y:

X Y Y D

Symbol table:

X: A B

Y: X C

BPE for Text (Sennrich et al., 2016)

A B A B C A B C D

Most frequent pair: A B

Merge symbol pair A B into AB:

AB AB C AB C D

Most frequent pair: AB C

Merge pair AB C into ABC:

AB ABC ABC D

Symbol table:

A B: AB

AB C: ABC

BPE Applied to English Wikipedia

_ t → _t

_ a → _a

h e → he

i n → in

_t he → _the

er

_s

on

_c

re

_o

_w

is

an

_in

...

ough

series

int

ai

stit

ery

ister

...

igo

osis

_jose

...

omato

...

Unsupervised Segmentation with BPE

Merge ops	Byte-pair encoded text
1000	to y od a _station .is a _r ail way station _on the _ch ū ō _main _l ine
3000	to y od a _station .is a _railway _station _on the _ch ū ō _main _line
10000	toy oda _station .is a _railway _station _on the _ch ū ō _main _line
50000	toy oda _station .is a _railway _station _on the _ch ū ō _main _line
100000	toy oda _station .is a _railway _station _on the _ch ū ō _main _line
Tokenized	toyoda station is a railway station on the chūō main line
10000	豊田 站 是 東 日 本 旅 客 鐵 道 (JR 東 日 本) 中 央 本 線 的 鐵 路 車 站
25000	豊田 站 是 東 日 本 旅 客 鐵 道 (JR 東 日 本) 中 央 本 線 的 鐵 路 車 站
50000	豊田 站 是 東 日 本 旅 客 鐵 道 (JR 東 日 本) 中 央 本 線 的 鐵 路 車 站
Tokenized	豊田站 是 東日本 旅客 鐵道 (JR 東日本) 中央本線 的 鐵路車站
5000	豊 田 駅 (と よ だ え き) は 、 東 京 都 日 野 市 豊 田 四 丁 目 に 在 る
10000	豊 田 駅 (と よ だ え き) は 、 東 京 都 日 野 市 豊 田 四 丁 目 に 在 る
25000	豊 田 駅 (と よ だ え き) は 、 東 京 都 日 野 市 豊 田 四 丁 目 に 在 る
50000	豊 田 駅 (と よ だ え き) は 、 東 京 都 日 野 市 豊 田 四 丁 目 に 在 る
Tokenized	豊田 駅 (とよだえき) は、東京都日野市豊田四丁目にある

Download Embeddings and BPE Models in 275 Languages



<https://github.com/bheinzerling/bpemb>