

NLU lecture 6:

Compositional character language modeling (& word2vec)

representations

Probability simply requires us to obey the following rules
(remember: V is finite):

Adam Lopez
alopez@inf.ed.ac.uk

Credits: Clara Vania
2 Feb 2018

$$P : V \rightarrow \mathcal{R}_+$$

$$\sum_{w \in V} P(w | w_{i-n+1}, \dots, w_{i-1}) = 1$$

When does this assumption make sense for language modeling?

But words are not a finite set!

What if we could scale softmax to the training data vocabulary? Would that help?

- Bengio et al.: “Rare words with frequency ≤ 3 were merged into a single symbol, reducing the vocabulary size to $|V| = 16,383$.”
- Bahdanau et al.: “we use a shortlist of 30,000 most frequent words in each language to train our models. Any word not included in the shortlist is mapped to a special token ([UNK]).”

SOFTMAX ALL THE WORDS



Src | 日本 の 主要 作物 は 米 で ある 。
Ref | the main crop of japan is rice .
Hyp | the _UNK is popular of _UNK . _EOS

Idea: scale by partitioning

- Partition the vocabulary into smaller pieces.

$$p(w_i|h_i) = p(c_i|h_i)p(w_i|c_i, h_i)$$

Class-based LM



Idea: scale by partitioning

- Differentiated softmax: assign more parameters to more frequent words, fewer to less frequent words.

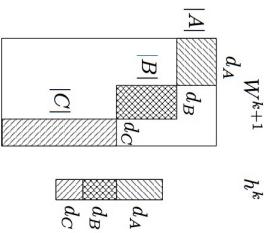


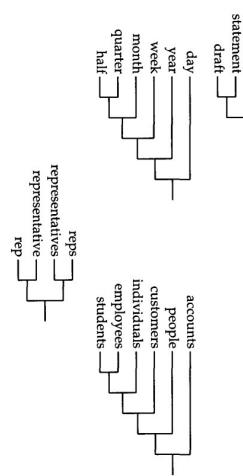
Figure 1: Final weight matrix W^{k+1} and hidden layer h^k for differentiated softmax for partitions A, B, C of the output vocabulary with embedding dimensions d_A, d_B, d_C ; non-shaded areas are zero.

Source: Strategies for training large vocabulary language models.
Chen, Auli, and Grangier, 2015

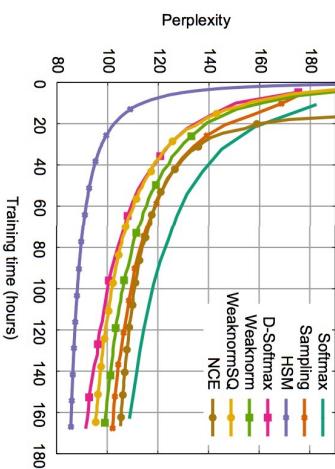
Idea: scale by partitioning

- Partition the vocabulary into smaller pieces hierarchically (*hierarchical softmax*).

Brown clustering:
hard clustering
based on mutual
information



Partitioning helps



| Dataset | Train | Test | Vocab | OOV |
|-----------|--------|-------|-------|------|
| PTB | 1M | 0.08M | 10k | 5.8% |
| gitaword | 4.631M | 27.9M | 100k | 5.6% |
| billlionW | 7.99M | 8.1M | 793k | 0.3% |

Table 1: Dataset statistics. Number of tokens for train and test set, vocabulary size and ratio of out-of-vocabulary words in the test set.

Source: Strategies for training large vocabulary language models.
Chen, Auli, and Grangier, 2015

Partitioning helps... but

could be better

Partitioning helps... but

| KN | PTB | gigaword | billionW |
|---------------|-------|----------|----------|
| KN | 141.2 | 57.1 | 70.2 |
| Softmax | 123.8 | 56.5 | 108.3 |
| D-Softmax | 121.1 | 52.0 | 91.2 |
| Sampling | 124.2 | 57.6 | 101.0 |
| HSM | 138.2 | 57.1 | 85.2 |
| NCE | 143.1 | 78.4 | 104.7 |
| Weaknorm | 124.4 | 56.9 | 98.7 |
| WeaknormSQ | 122.1 | 56.1 | 94.9 |
| KN+Softmax | 108.5 | 43.6 | 59.4 |
| KN+D-Softmax | 107.0 | 42.0 | 56.3 |
| KN+Sampling | 109.4 | 43.8 | 58.1 |
| KN+HSM | 115.0 | 43.9 | 55.6 |
| KN+NCE | 114.6 | 49.0 | 58.8 |
| KN+Weaknorm | 109.2 | 43.8 | 58.1 |
| KN+WeaknormSQ | 108.8 | 43.8 | 57.7 |

Table 2: Test perplexity of individual models and interpolation with Kneser-Ney.

Partitioning helps... but
could be better

Partitioning helps... but
could be better

| KN | PTB | gigaword | billionW |
|---------------|-------|----------|----------|
| KN | 141.2 | 57.1 | 70.2 |
| Softmax | 123.8 | 56.5 | 108.3 |
| D-Softmax | 121.1 | 52.0 | 91.2 |
| Sampling | 124.2 | 57.6 | 101.0 |
| HSM | 138.2 | 57.1 | 85.2 |
| NCE | 143.1 | 78.4 | 104.7 |
| Weaknorm | 124.4 | 56.9 | 98.7 |
| WeaknormSQ | 122.1 | 56.1 | 94.9 |
| KN+Softmax | 108.5 | 43.6 | 59.4 |
| KN+D-Softmax | 107.0 | 42.0 | 56.3 |
| KN+Sampling | 109.4 | 43.8 | 58.1 |
| KN+HSM | 115.0 | 43.9 | 55.6 |
| KN+NCE | 114.6 | 49.0 | 58.8 |
| KN+Weaknorm | 109.2 | 43.8 | 58.1 |
| KN+WeaknormSQ | 108.8 | 43.8 | 57.7 |

Table 2: Test perplexity of individual models and interpolation with Kneser-Ney.

Partitioning helps... but
could be better

Partitioning helps... but
could be better

| KN | PTB | gigaword | billionW |
|---------------|-------|----------|----------|
| KN | 141.2 | 57.1 | 70.2 |
| Softmax | 123.8 | 56.5 | 108.3 |
| D-Softmax | 121.1 | 52.0 | 91.2 |
| Sampling | 124.2 | 57.6 | 101.0 |
| HSM | 138.2 | 57.1 | 85.2 |
| NCE | 143.1 | 78.4 | 104.7 |
| Weaknorm | 124.4 | 56.9 | 98.7 |
| WeaknormSQ | 122.1 | 56.1 | 94.9 |
| KN+Softmax | 108.5 | 43.6 | 59.4 |
| KN+D-Softmax | 107.0 | 42.0 | 56.3 |
| KN+Sampling | 109.4 | 43.8 | 58.1 |
| KN+HSM | 115.0 | 43.9 | 55.6 |
| KN+NCE | 114.6 | 49.0 | 58.8 |
| KN+Weaknorm | 109.2 | 43.8 | 58.1 |
| KN+WeaknormSQ | 108.8 | 43.8 | 57.7 |

Table 2: Test perplexity of individual models and interpolation with Kneser-Ney.

Table 2: Test perplexity of individual models and interpolation with Kneser-Ney.

V is not finite

What set is finite?

- Practical problem: softmax computation is linear in vocabulary size.

- **Theorem.** The vocabulary of word types is infinite.

Proof 1. productive morphology, loanwords, “fleek”

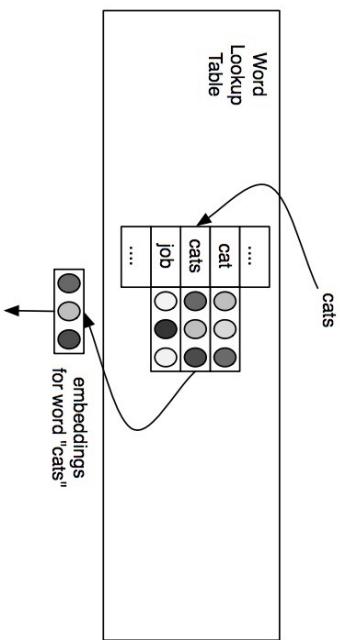
Proof 2. 1, 2, 3, 4, ...

Characters.
More precisely, unicode code points.

💡 Are you sure? 🤔

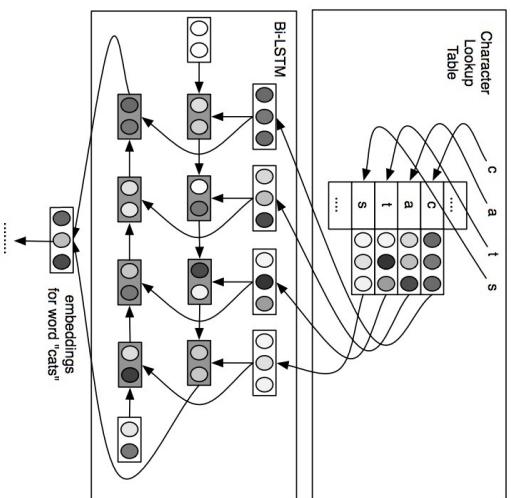
Not all characters are the same, because not all languages have alphabets. Some have syllabaries (e.g. Japanese kana) and/ or logographies (Chinese hàanzi).

Rather than look up word representations...



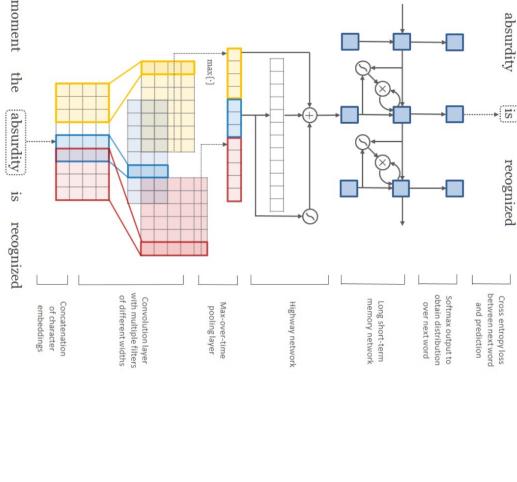
Source: Finding function in form: compositional character models for open vocabulary word representation, Ling et al. 2015

Compose character representations into word representations with LSTMs



Source: Finding function in form: compositional character models for open vocabulary word representation, Ling et al. 2015

Compose character representations into word representations with CNNs



Source: Character-aware neural language models, Kim et al. 2015

Character models actually work. Train them long enough, they generate words

antrest
artifactive
capacited
capitaling
compensive
dermitories
despator
dividement
extremilated
faxemary
follect

hamburgo
identimity
ipoteca
nightmale
orience
patholism
pinguenas
sammitment
tasteman
understrumental
wisholver

Wow, the disconversated vocabulations of their system are fantastics!
—Sharon Goldwater

Character models actually work. Train them long enough, they generate words

Table 2: Most-similar in-vocabulary words under the C2W model; the two query words on the left are in the training vocabulary, those on the right are nonce (invented) words.

| <i>increased</i> | <i>John</i> | <i>Noahshire</i> | <i>phding</i> |
|------------------|-------------|------------------|---------------|
| reduced | Richard | Nottinghamshire | mixing |
| improved | George | Bucharest | modelling |
| expected | James | Saxony | styling |
| decreased | Robert | Johannesburg | blaming |
| targeted | Edward | Gloucestershire | christening |

How good are character-level NLP models?

tation algorithm. Ling et al. (2015b) indeed states that “[In]uch of the prior information regarding morphology, cognates and rare word translation among others, should be incorporated”.

It however becomes unnecessary to consider these prior information, if we use a neural network, be it recurrent, convolution or their combination, directly on the unsegmented character sequence. The possibility of using a sequence of un-

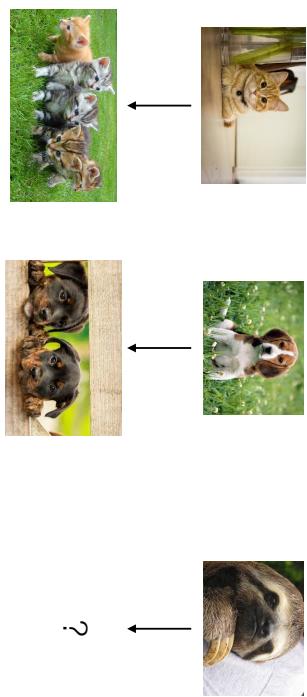


Implied(?): character-level neural models learn everything they need to know about language.

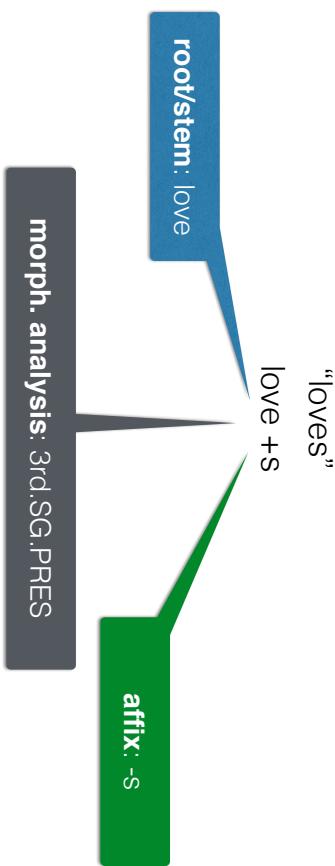
Source: Finding function in form: compositional character models for open vocabulary word representation, Ling et al. 2015

Word embeddings have obvious limitations

- Closed vocabulary assumption



Morpheme: the smallest meaningful unit of language



The ratio of morphemes to words varies by language

Analytic languages

one morpheme per word

Hai *đều* *bo?* *nhanh* *là* *tại* *giêng*-đinh *thông* *chồng*.
two individual leave each-other be because-of family guy husband.
'They divorced because of his family'

Vietnamese

Synthetic languages

many morphemes per word

English

And **we** know a lot about linguistic structure

Posit-angl-i-unmarr-para
understand-not-completely-1SG.SBJ;3SG.OBL;IND
ilta-juma-suuit
come-want-2SG.PTCP
'I didn't understand at all that you wanted to come along'

West Greenlandic

Derivational morphology

love (VB), loves (VB), loving (VB), loved (VB)

Turkish

lover (NN), lovely (ADJ), lovable (ADJ)

Morphemes can represent one or more *features*

one or more *features*

Agglutinative languages

one feature per morpheme

(Turkish)
oku-**r-sa-m**

read-AOR.COND.1SG
'if I read ...'

Fusional languages

many features per morpheme

(English)
read-**s**

read-3SG.SG
'reads'

rescue + LNK + helicopter + emergency + landing + place
'Rescue helicopter emergency landing pad'

Inflection is not limited to affixation

There are many different ways to compute word representations from *subwords*

Words can have more than one stem

Affixation

one stem per word

studying
study + ing

(English)
read-AOR.COND.1SG
'if I read ...'

Compounding

many stems per word

(German)
Rettungshubschraubernotlandeplatz

Rettung + s + **hubschrauber** + **not** + **lande**+ **platz**

rescue + LNK + helicopter + emergency + landing + place
'Rescue helicopter emergency landing pad'

Base Modification

drink, **drank**, **drunk**

(English)

Root & Pattern

k(a)t(a)b(a)
write-PST.3SG.M
'he wrote'

(Arabic)

Reduplication

kemerah~merahan

(Indonesian)

Basic Units of Representation

Compositional Function
basic unit(s)

- Characters (Ling et al., 2015, Kim et al., 2016, Lee et al., 2016)

- Character n-grams (Sperr et al., 2013, Wieting et al., 2016, Bojanowski et al., 2016)

- Morphemes (Luong et al., 2013, Botha & Blunsom, 2014, Sennrich et al., 2016)

Bidirectional LSTMs,
Convolutional NN,
...

- Morphological analysis (Cotterel & Schütze, 2015, Kahn & Schütze, 2016)



We've revised morphology, so we have some questions about character models

- How do representations based on morphemes compared with those based on characters?
- What is the best way to compose subword representations?
- Do character-level models the same **predictive utility** as models with knowledge of morphology?
- How do different representations interact with languages of different morphological typologies?

Open vocabulary history $\sigma(\text{the})$ $\sigma(\text{dogs})$ $\sigma(\text{are})$ subword units $f(\cdot)$ $f(\cdot)$ $f(\cdot)$ composition function word representation

Closed vocabulary prediction $P(\text{playing, are, dogs, the, \dots})$

Open vocabulary prediction is interesting, but our goal is to understand representations, not build a better neural LM.

Variable: Subword Unit

Variable: Composition Function

- Vector addition (except for characters)

- Bidirectional LSTMs
- Convolutional NN

| Unit | Examples |
|--------------|--|
| Morfessor | <code>\want, \$s\$</code> |
| BPE | <code>\w, antis\$</code> |
| char-trigram | <code>\wa, wan, ant, nts, ts\$</code> |
| character | <code>\, w, a, n, t, s, \$</code> |
| analysis | <code>want+VB, +3rd, +SG, +Pres</code> |

The last row is part of an oracle experiment: suppose you had an oracle that could tell you the true morphology. In this case, the oracle is a human annotator.

Prediction problem: neural language modeling

Variable: Language Typology

Summary of perplexity: use bi-LSTMs over character trigrams

Fusional (English)

Agglutinative (Turkish)

| Language | word | character | char-trigrams | BPE | Morfessor | bi-LSTM | %imp |
|-------------|-------------------|-----------|---------------|---------|-----------|---------|-------|
| | bi-LSTM | CNN | add | bi-LSTM | add | bi-LSTM | |
| read-s | oku-r-sa-m | 33.59 | 49.96 | 33.74 | 47.74 | 36.87 | 18.98 |
| read-3SG,SG | read-AOR.COND.1SG | 42.97 | 47.51 | 43.3 | 49.72 | 49.72 | 7.39 |
| 'reads' | 'if I read ...' | 27.72 | 40.1 | 28.52 | 39.6 | 31.31 | 20.64 |

Root&Pattern (Arabic)

Reduplication (Indonesian)

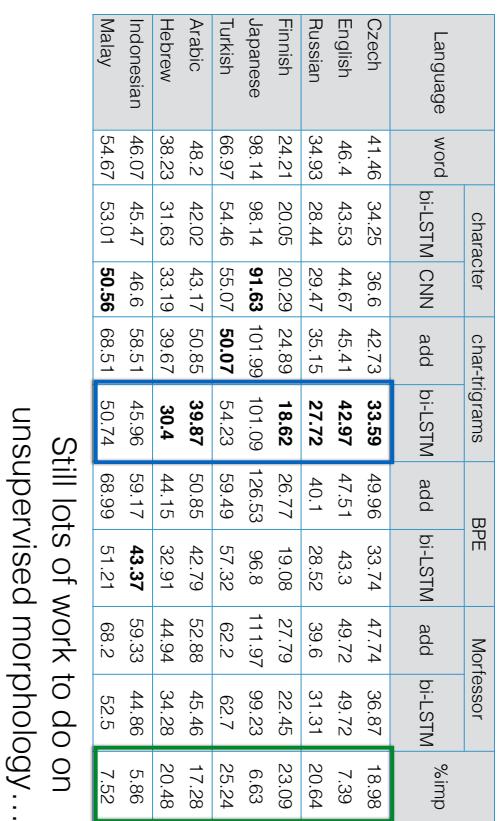
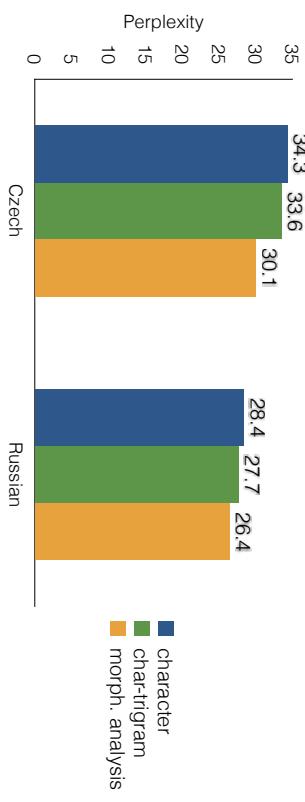
| Language | word | character | char-trigrams | BPE | Morfessor | bi-LSTM | %imp |
|-----------------|------------|-----------|---------------|---------|-----------|---------|-------|
| | bi-LSTM | CNN | add | bi-LSTM | add | bi-LSTM | |
| k(a)t(a)b(a) | anak~anak | 91.63 | 101.99 | 126.53 | 96.8 | 111.97 | 99.23 |
| write-PST.3SG.M | child-PL | 50.07 | 55.07 | 54.23 | 59.49 | 57.32 | 62.2 |
| 'he wrote' | 'children' | 39.87 | 43.17 | 50.85 | 42.79 | 52.88 | 45.46 |

Do character-level models have the predictive utility of models with access to actual morphology? **NO**

no morphology

actual morphology

(^, r, e, a, d, s, \$) —→ (read, VB, 3rd, SG, Present)



Still lots of work to do on unsupervised morphology...

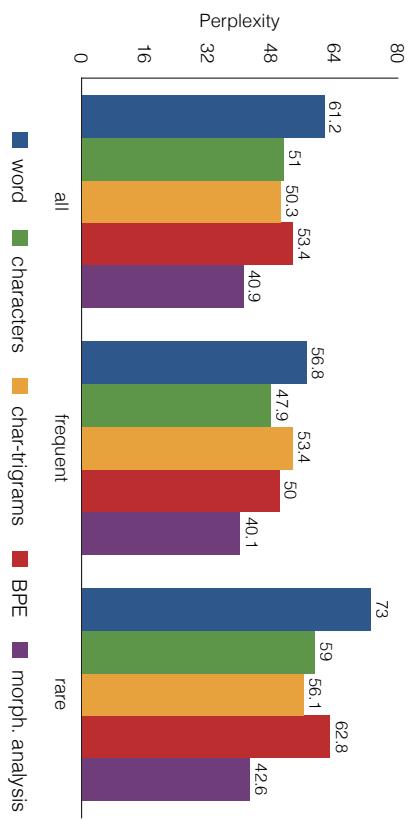
How do we know that the morphological annotations that make the difference?

- Measure **Targeted perplexity**: perplexity on specific subset of words in the test data

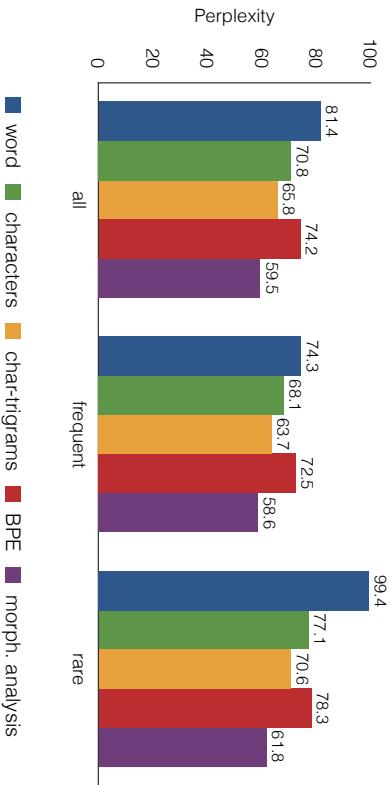
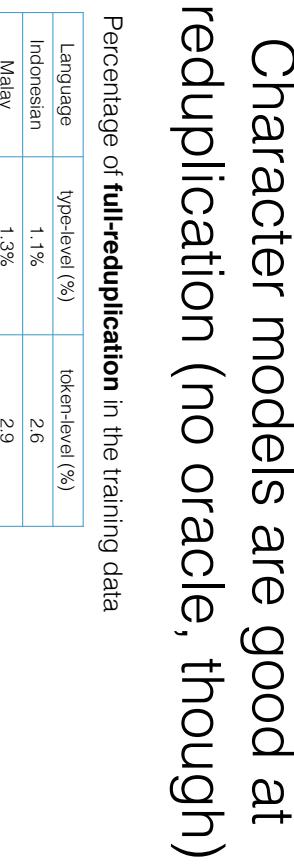
- Analyze perplexities when the **inflected words** of interest are in the most recent history: **nouns** and **verbs**

Green **tea** or white **tea**?

*The sushi **is** great, and they **have** a great **selection**.*



Targeted perplexity of Czech verbs is lower when we use morphology

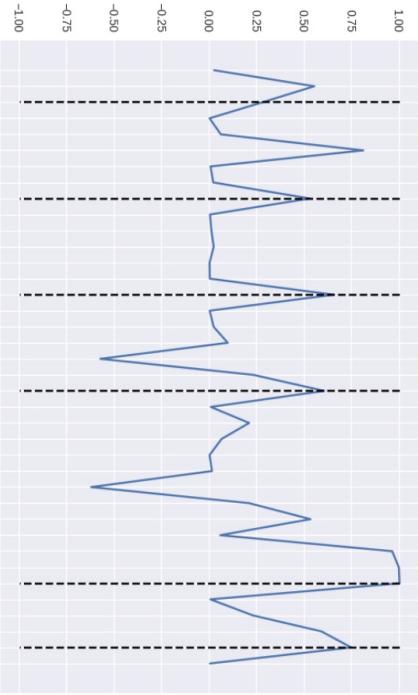
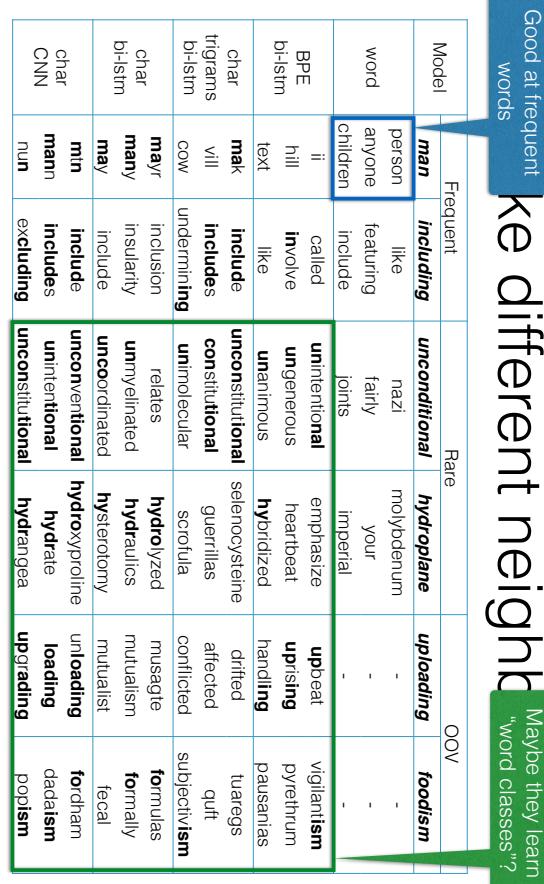


Different representations make different neighbors

| Model | Frequent | | Rare | | Unknown | |
|----------|-------------|--------------------|-------------------------|-------------------|------------------|---------------------|
| | <i>man</i> | <i>including</i> | <i>unconditional</i> | <i>hydroplane</i> | <i>uploading</i> | <i>foodism</i> |
| word | person | like | nazi | molybdenum | - | - |
| | anyone | featuring | fairly | your | - | - |
| | children | include | joints | imperial | - | - |
| BPE | ii | called | unintentional | emphasize | upbeat | vigilantism |
| bi-lstm | hill | involve | ungenerous | heartbeat | uprising | pyrethrum |
| | text | like | unanimous | hybridized | handling | pausanas |
| char | mak | include | unconstitutional | selenocysteine | drifted | tuaregs |
| trigrams | vill | includes | constitutional | guerrillas | affected | quft |
| bi-lstm | cow | undermining | unimolecular | scrofula | conflicted | subjectivism |
| char | mayr | inclusion | relates | hydrolyzed | musagie | formulas |
| bi-lstm | many | insularity | unmyelinated | hydraulics | mutualism | formally |
| | may | include | uncoordinated | hysterotomy | mutualist | fecal |
| char | min | include | unconventional | hydroxyproline | fordham | fordham |
| CNN | mann | includes | unintentional | hydrate | loading | dacidism |
| | num | excluding | unconstitutional | hydrangea | upgrading | popism |

Character NLMs learn word boundaries.

...and memorize POS tags

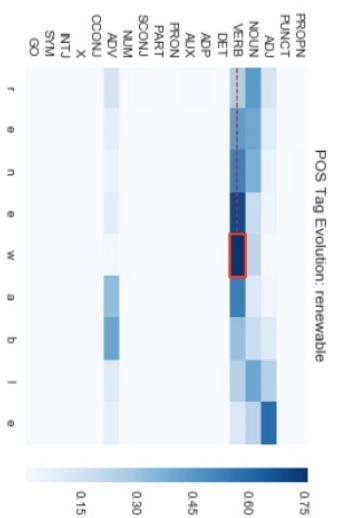


Source: Yova Kementchedzhieva, Morpho-syntactic awareness in a character-level language model
2017 Informatics M.Sc. thesis

Source: Yova Kementchedzhieva, Morpho-syntactic awareness in a character-level language model
2017 Informatics M.Sc. thesis

...and memorize POS tags

What do NLMs learn about morphology?



Source: Yova Kementchedzhieva, Morpho-syntactic awareness in a character-level language model
2017 Informatics M.Sc. thesis

What do we know about what NNs know about language?

- Still very little.
- Evidence suggests: nothing surprising. Lots of memorization, local generalization.
- NNs are great for simplicity of specification and end-to-end learning.
- But these things are not magic! We still don't have enough data, and these models could be better if they knew about morphology.
- But how do we do that?

- Character-level NLMs are great! Across typologies, but especially for agglutinative morphology.
- However, they **do not match** predictive accuracy of model with explicit knowledge of morphology (or POS).
- Qualitative analyses suggests that they learn **orthographic similarity** of affixes, and **forget meaning of root morphemes**.
- More generally, they appear to **memorize frequent subpatterns**.