

DEEP LEARNING ALAPÚ SZERZŐAZONOSÍTÁS

FÜLEKI FÁBIÁN, JANI BALÁZS GÁBOR, TORNER MÁRTON

CSAPAT: LOREMIPSUM

SZERZŐAZONOSÍTÁS

- EGYEDI STÍLUS FELISMERÉSE
 - KÉZÍRÁS
 - SZAVAK ÉS SZÓFORDULATOK HASZNÁLATA
 - MONDAT SZERKEZET
- STILOMETRIA
- PLÁGIUMDETEKCIÓ

KORÁBBI MUNKÁK

- SENTENCE LEVEL VS. ARTICLE LEVEL MEGKÖZELÍTÉS (46% vs 69.1%)
- GRU, LSTM (69.1% vs 62.7%)
- (SZIÁMI HÁLÓZAT)
 - PLÁGIUMDETEKCIÓ 99.8%!

DATASET - REUTER_50_50 (C50)

- NYERS SZÖVEG
- 50 SZERZŐ 50+50 CIKKE (5000DB)
- ÖSSZEHASONLÍTHATÓSÁG MÁS CIKKEKKEL
- [HTTPS://ARCHIVE.ICS.UCI.EDU/ML/DATASETS/REUTER_50_50](https://archive.ics.uci.edu/ml/datasets/REUTER_50_50)

ADATOK ELŐKÉSZÍTÉSE

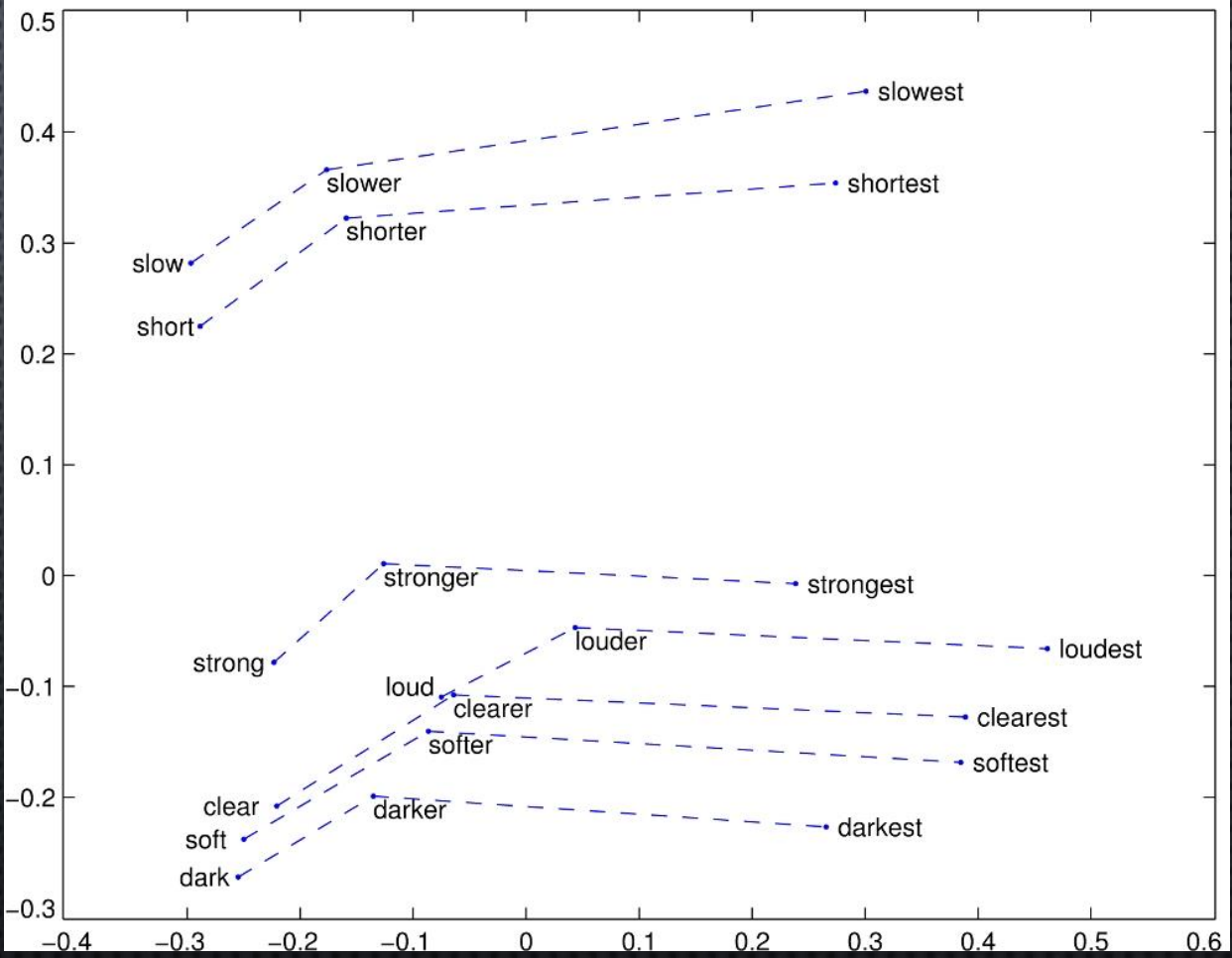
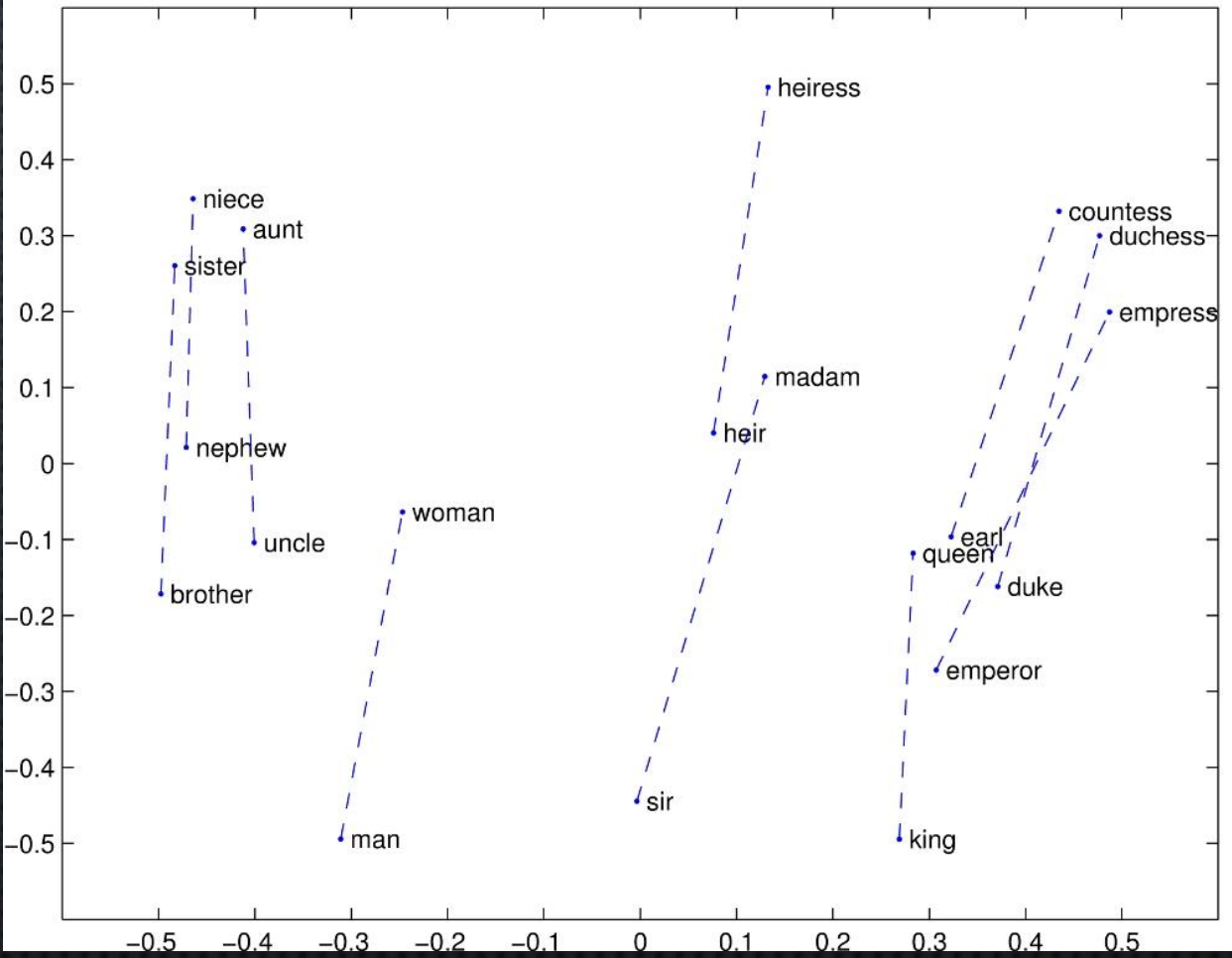
- SZAVAK REPRESENTÁLÁSA SZÁMOKKÉNT
 - SZÓVEKTOROK
 - PART OF SPEECH TAG-EK
- AZONOS HOSSZÚSÁGÚ MONDATOK
- TOVÁBBI TRÜKKÖK

SPACY

- GLOVE
- PART OF SPEECH TAGGING
- 301 DIMENZIÓS VEKTOR
 - SZÓ REPRESENTÁCIÓ + POS TAG
- KISZŰRT RÉSZEK
 - STOP WORDS, PUNCTUATION
 - RITKA SZAVAK (NINCS MEGFELELŐ REPRESENTÁCIÓ AZ ADATBÁZISBAN)

GLOBAL VECTOR FOR WORDS

- „GLOVE IS AN UNSUPERVISED LEARNING ALGORITHM FOR OBTAINING VECTOR REPRESENTATIONS FOR WORDS.”
- A SZAVAK REPREZENTÁLÁSA VEKTOROKKÉNT
- ÁLTALUNK HASZNÁLT BÁZIS: 685 000 EGYEDI VEKTOR (SZÓ)
 - 300 DIMENZIÓS VEKTOROK
- [HTTPS://NLP.STANFORD.EDU/PUBS/GLOVE.PDF](https://nlp.stanford.edu/pubs/glove.pdf)
- [HTTPS://NLP.STANFORD.EDU/PROJECTS/GLOVE/](https://nlp.stanford.edu/projects/glove/)



	0	1	2	3	4	5	6	7	8	...
text	The	Internet	may	be	overflowing	with	new	technology	but	...
vector	[0.27204, -0.06203, -0.1884, 0.023225, -0.0181...	[-0.50955, 0.088231, -0.32273, -0.40398, 0.003...	[-0.042501, 0.090773, -0.11918, 0.12372, -0.19...	[-0.059177, 0.10653, -0.21613, -0.086178, 0.00...	[0.074908, -0.036973, 0.082992, -0.31622, 0.22...	[-0.099534, 0.028202, -0.23189, 0.094477, 0.12...	[0.34046, 0.13752, -0.20643, -0.4555, 0.19251, ..	[-0.32298, 0.38883, 0.4586, -0.5227, -0.064451...	[-0.01689, 0.17402, -0.30247, -0.30063, 0.2141...	...
pos_str	DET	NOUN	VERB	VERB	VERB	ADP	ADJ	NOUN	CCONJ	...

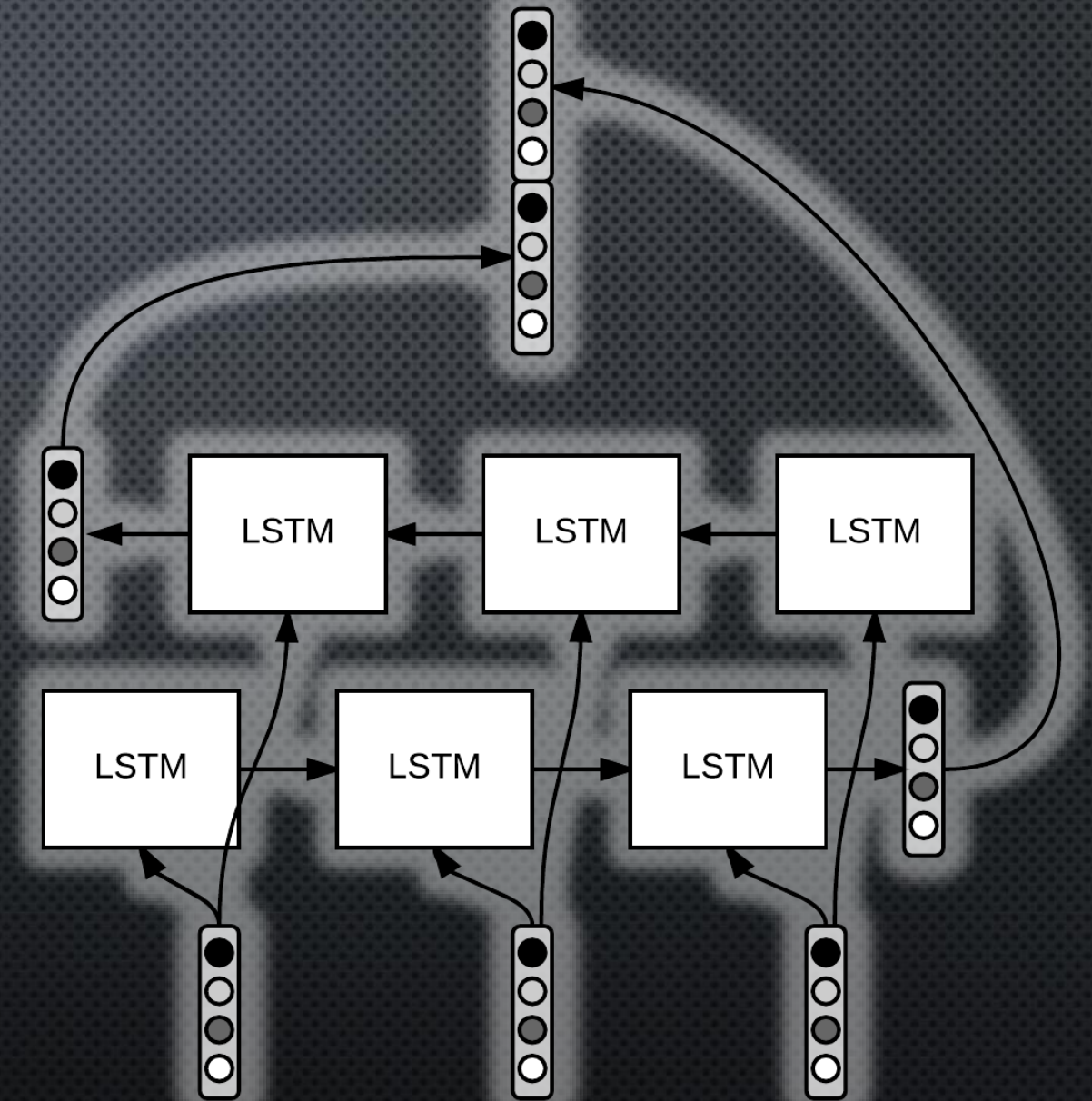
MONDATOK ÖSSZEFOGÁSA

- AZONOS HOSSZÚSÁG
 - RÖVIDRE VÁGÁS
 - BŐVÍTÉS/KIEGÉSZÍTÉS
- 3 MONDAT EGYÜTTES HASZNÁLATA
 - KIFEJEZÉSEK ÁTFEDÉSE

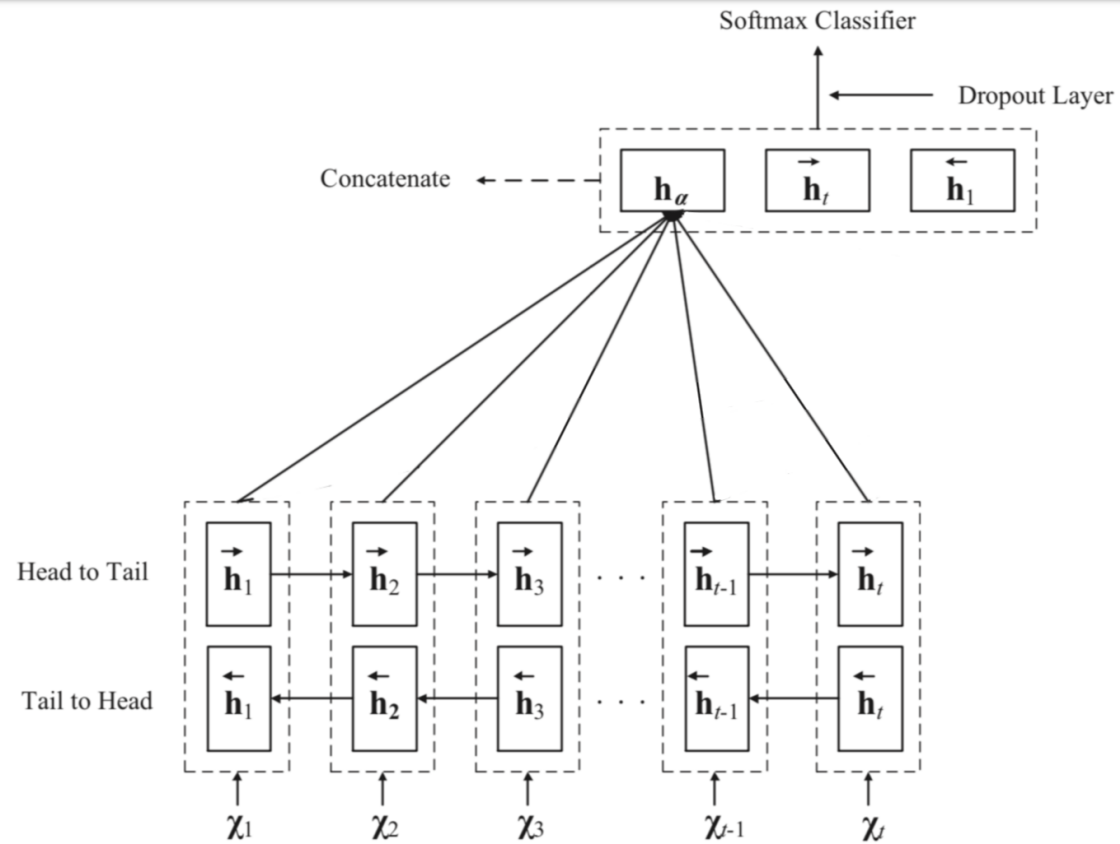


ARCHITEKTÚRÁK

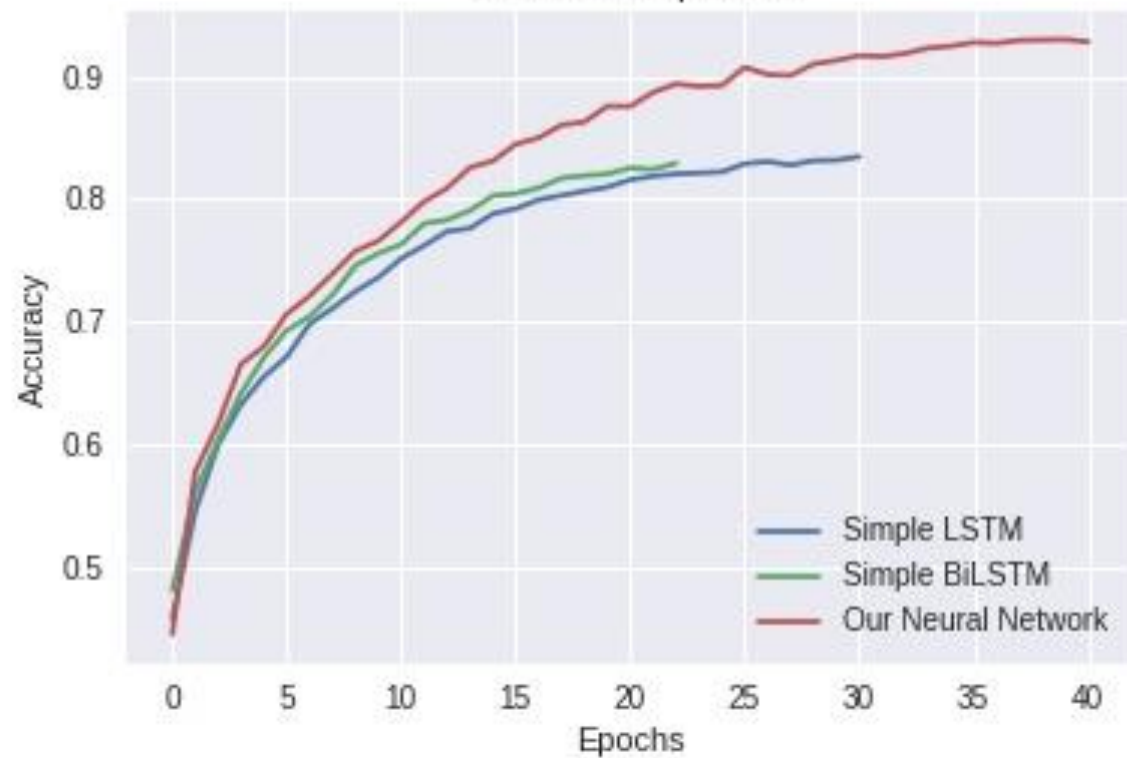
- LSTM / BiLSTM
 - TIME DISTRIBUTED
 - AVERAGE POOLING
 - SOFTMAX CLASSIFIER
- BiLSTM
 - KATAFORÁK ÉS ANAFORÁK



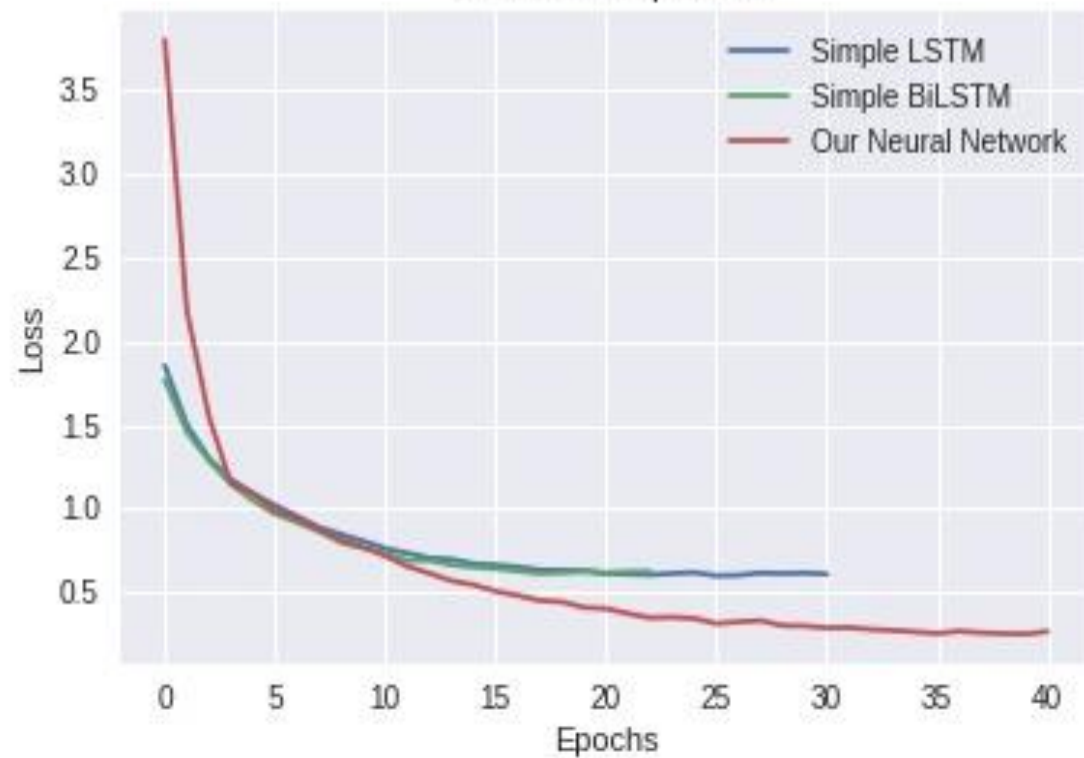
SAJÁT MEGKÖZELÍTÉS



Network comparison



Network comparison



HIPEROPTIMALIZÁLÁS PARAMÉTEREI

- BATCH-EK MÉRETE: 256, 512, 1024, 2048
- NEURONOK SZÁMA: 128, 256, 512, 1024
- DROPOUT, REKURRENS DROPOUT: $0 \rightarrow 1$
- OPTIMIZER: SGD, RMSPROP, ADAGRAD, ADAM

HIPEROPTIMALIZÁLÁS KONKLÚZIÓI

- LEGJOBB OPTIMALIZÁCIÓS MÓDSZEREK: ADAM, RMSPROP
- LEGJOBB FUTTATÁS PARAMÉTEREI: 512-ES BATCH MÉRET, 512 NEURON, (DROPOUT: 0,5 KÖRÜL), ADAM
- TÚL SOK NEURON / TÚL KICSI BATCH MÉRET: SOKKAL LASSABBAN TANUL CSAK EGY KICSIT JOBBAN

ÖSSZEHASONLÍTÁS

- QIAN, CHEN, TIANCHANG HE, AND RAO ZHANG.
"DEEP LEARNING BASED AUTHORSHIP IDENTIFICATION."
 - ARTICLE LEVEL LSTM – 69.1% ACCURACY (BEST)
 - [HTTP://WEB.STANFORD.EDU/CLASS/CS224N/REPORTS/2760185.PDF](http://web.stanford.edu/class/cs224n/reports/2760185.pdf)
- SAJÁT: 94%
 - 3 SENTENCE BUNDLES
 - 512 BATCH SIZE, 512 NEURON, 0.53 DROPOUT, 0.5 RECURRENT DROPOUT

MERRE TOVÁBB?

- ISMERETLEN CSOPORT
- ATTENTION
- MÁΣ ELŐFELDOLGOZÁS (NAGYOBB BUNDLE-ÖK, ARTICLE LEVEL, STB.)
 - PLUSZ SZEMANTIKAI INFORMÁCIÓK

KÖSZÖNJÜK A FIGYELMET!