

SOME DEPENDENCY PARSING WORK

MARK ANDERSON

OVERVIEW

Developing

- Chunk-and-Pass
- Distillation

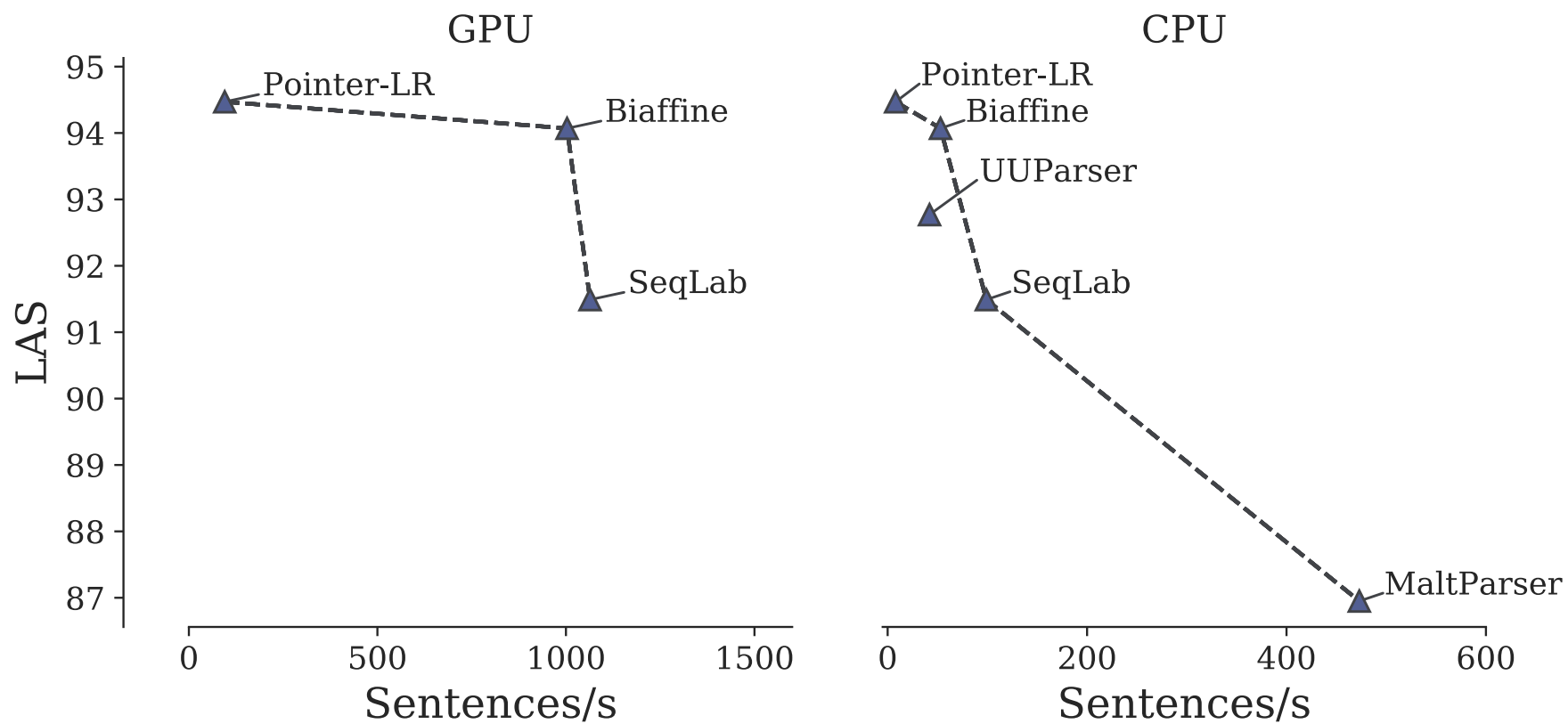
Evaluating

- Edge displacement
- POS tags

PART I

DEVELOPING PARSERS

PARETO FRONT FOR PTB (WSJ)



CHUNK AND PASS

CHUNK AND PASS

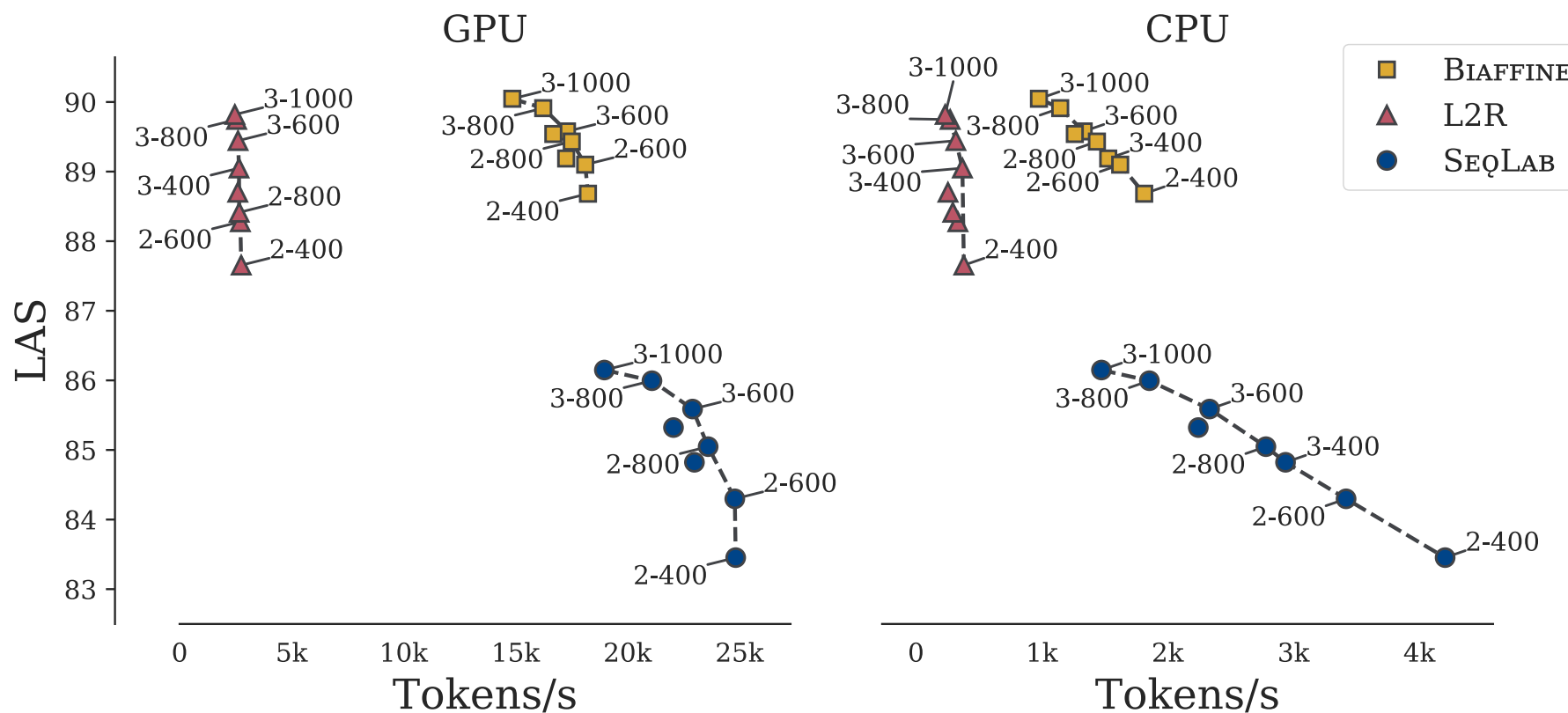
"Now-or-never" bottleneck¹

3 steps

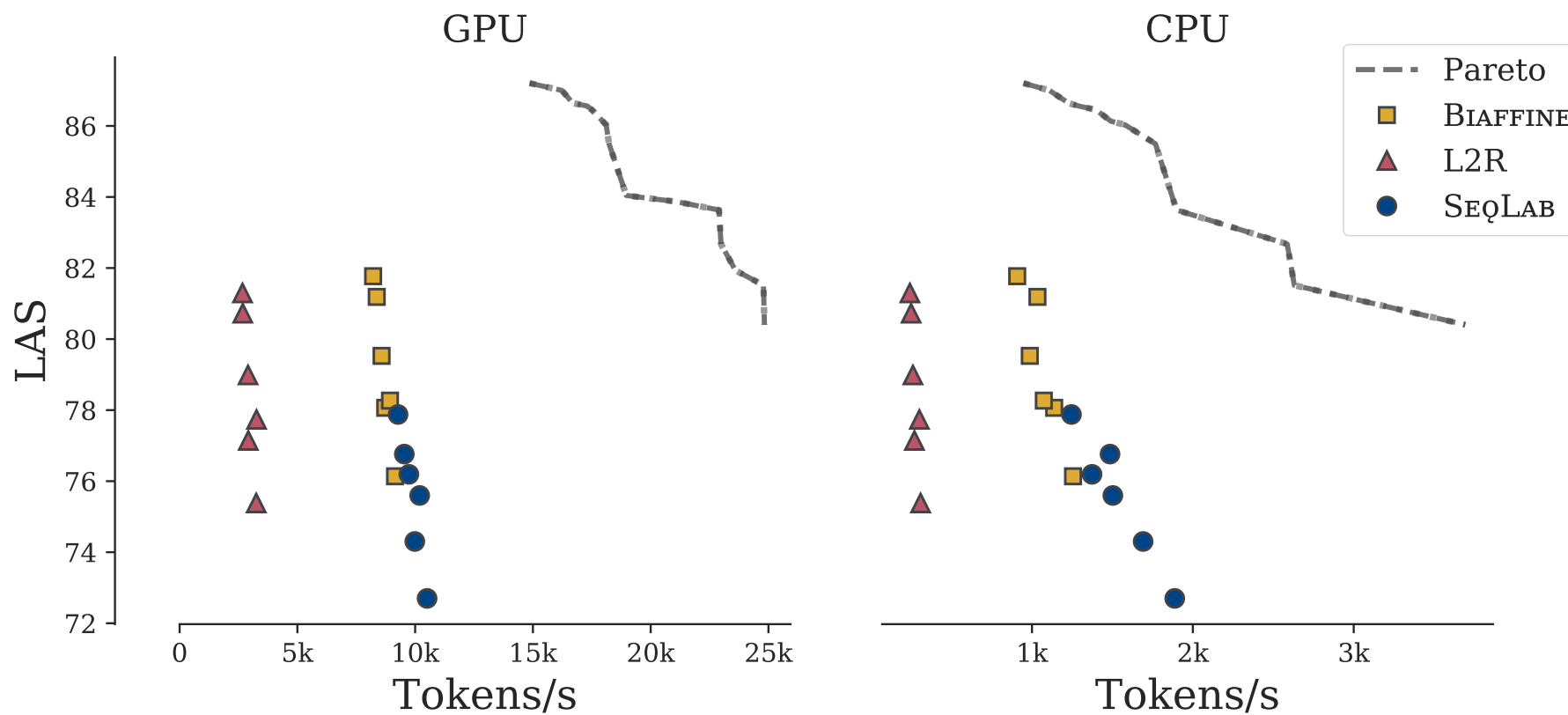
1. Shallow parse (chunk) --> Compressed data and shallow syntax
2. Parse chunks --> Higher-order syntactic relations
3. Collate --> Combined 1 and 2 for full parse

¹ Christiansen, M.H. and Chater N., *The Now-or-Never bottleneck: A fundamental constraint on language*, 2015

PARETO FRONT FOR UD (ZH, HI, KO, PL)



PARETO FRONT FOR UD (ZH, HI, KO, PL)



CHUNKING — Each word is labelled, **B**, **I**, or **O**.

B

A token that begins a chunk.

Suffixed with chunk phrase type.

E.g. B-NP for a noun phrase.

I

A token inside of a chunk.

Also suffixed with chunk phrase type.

E.g. I-VP for a verb phrase.

O

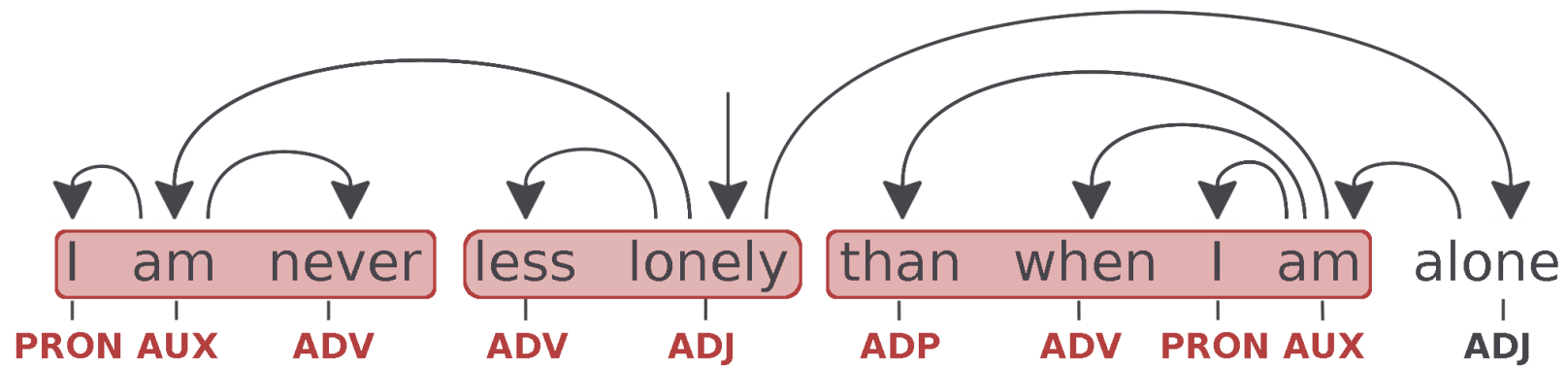
Anything outside of a chunk.

EXTRACTING CHUNKS

- 1. Evolutionary technique (slow...)**
- 2. Information theory**

CHUNK CANDIDATE CRITERIA

1. The components are syntactically linked
2. There is only one level of dependency (one head and its dependents)
3. The components are continuous.
4. No dependents within a chunk has a dependent outwith the chunk.



(DET ADJ NOUN)

(PRON AUX ADV)

(PART VERB)

(ADP ADV PRON AUX)

(SCONJ ADV VERB)

(AUX AUX VERB)

(PRON PROPN VERB)

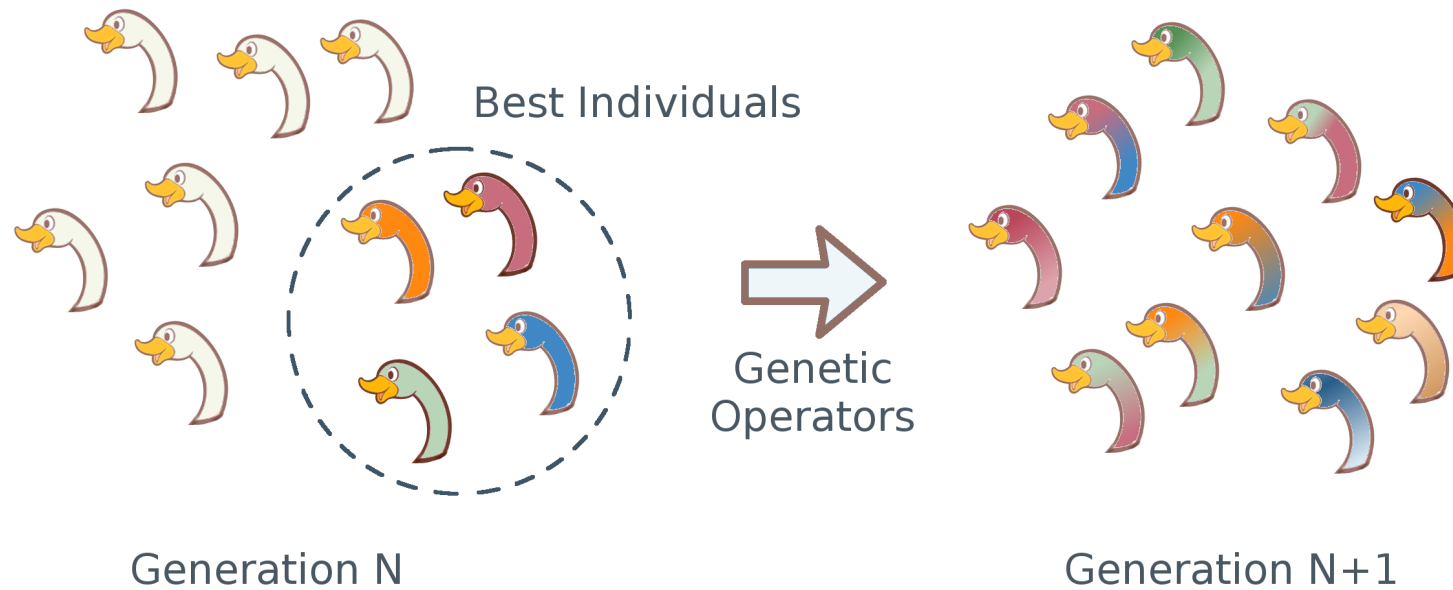
(CCONJ PRON AUX DET ADJ NOUN)

Extract 2615 unique rules from UD English EWT treebank v2.3

**512 occur more than 5
times.**

**1.34×10^{154} different
rule sets.**

EVOLUTIONARY SEARCH



EVOLUTIONARY SEARCH

Individual = [0,1,0,0,1 0,1]



(DET ADJ NOUN) (PRON AUX ADV) (PART VERB)

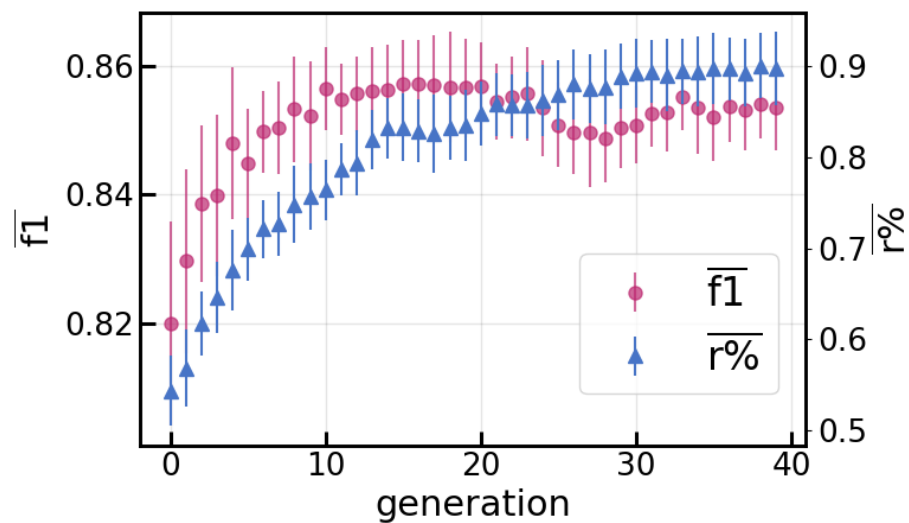
Fitness = Chunking F1-score
+ 0.5 x proportion of
max compression

PROPORTION OF MAX COMPRESSION

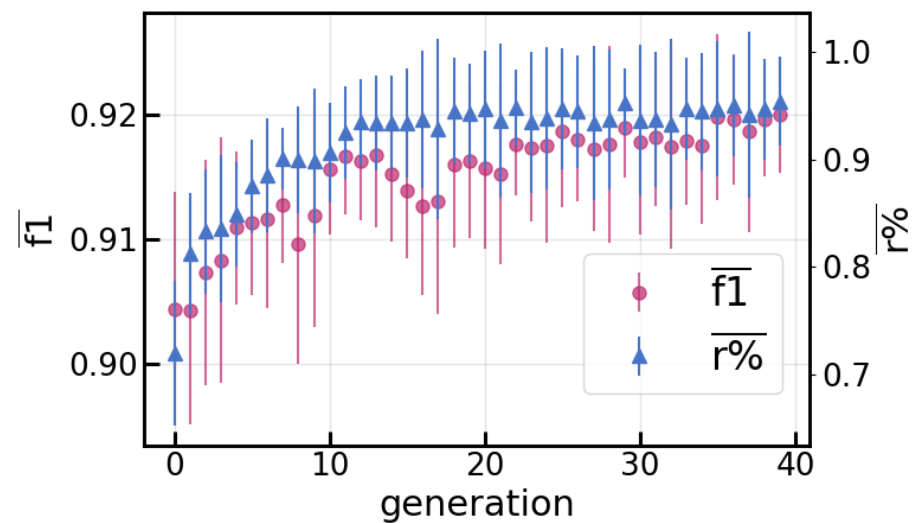
$$r = \frac{N_{tokens}}{N_{tokens} + N_{chunks}}$$

$$r_{prop} = \frac{r_{subset} - 1}{r_{all} - 1}$$

English-EWT



Japanese-GSD

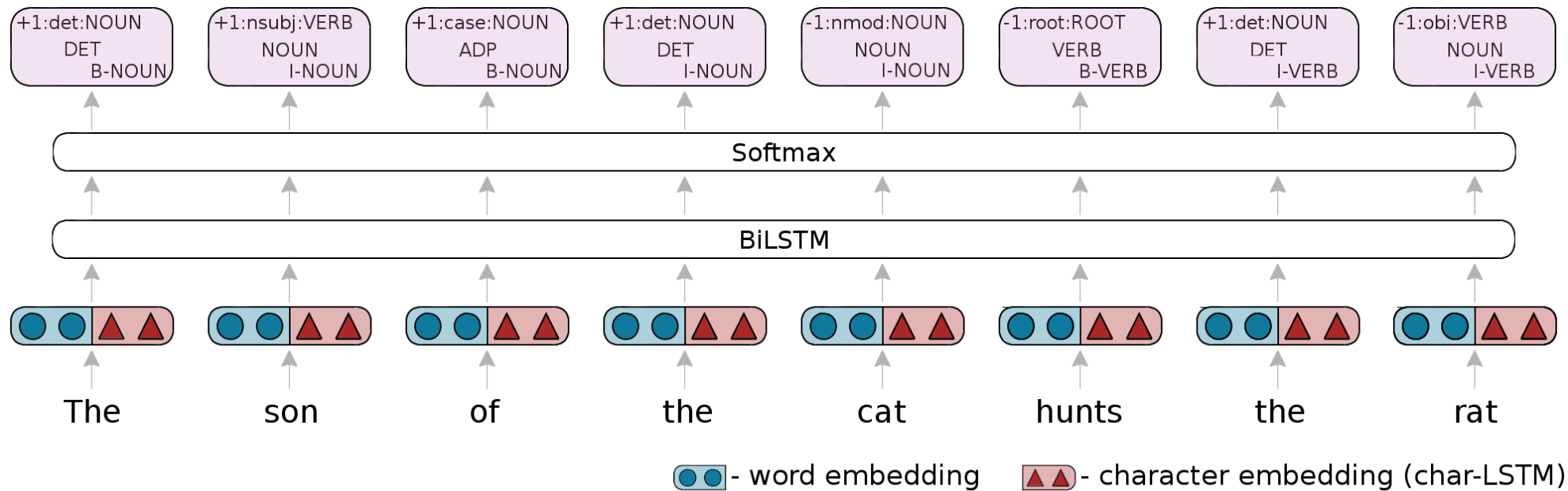


SYSTEM DETAILS.

We use a neural sequence labelling toolkit, NCRF++.¹ And a relative POS tag position encoding for sequence labelling parsing.²

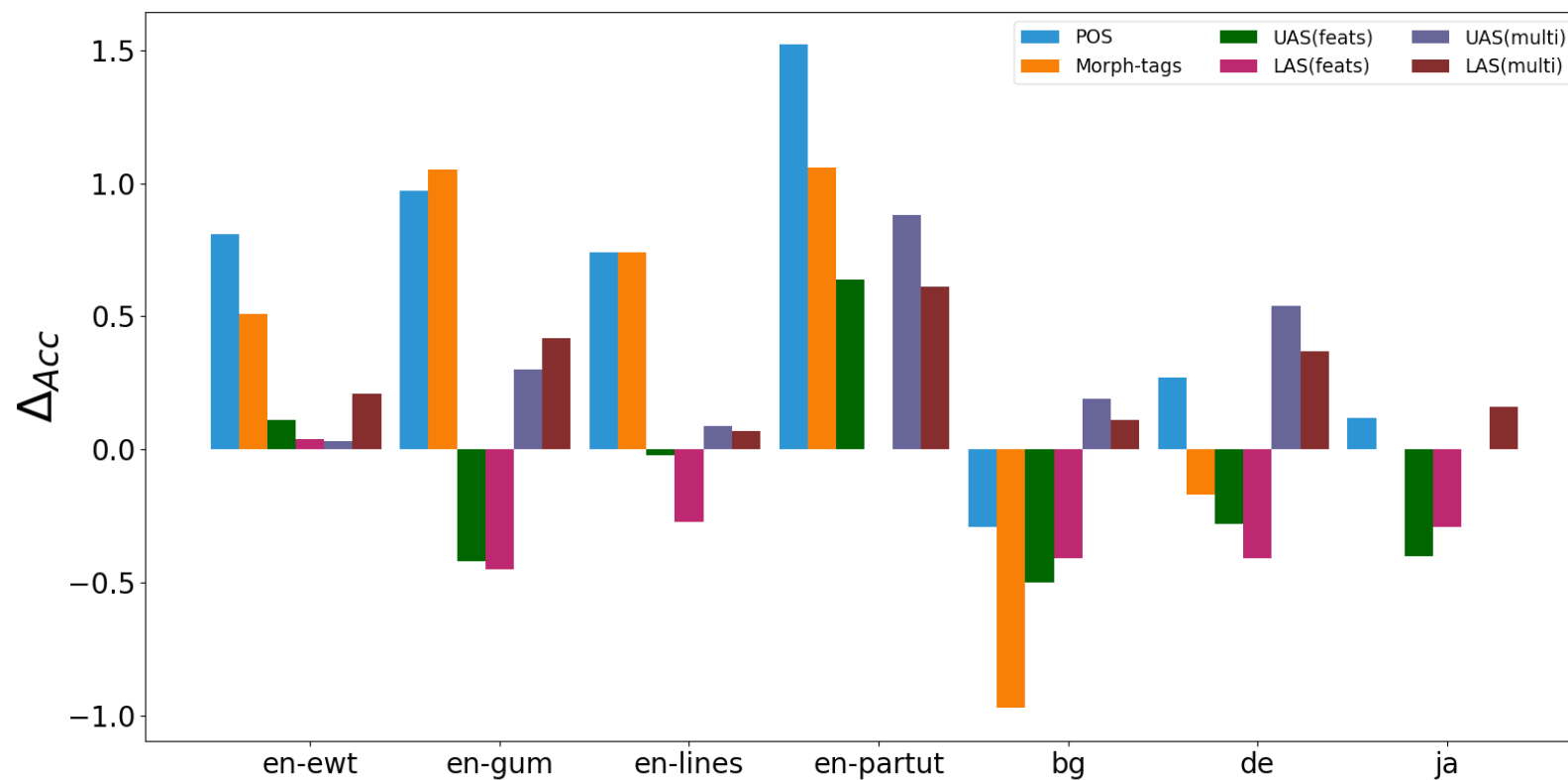
¹Yang, J. and Zhang Y., *NCRF++: An Open-source Neural Sequence Labeling Toolkit*, 2018

²Spoustová, D.J. and Spousta M. *Dependency Parsing as a Sequence Labeling Task*, 2010



BROAD RESULTS

Difference between **with** and **without** chunks.

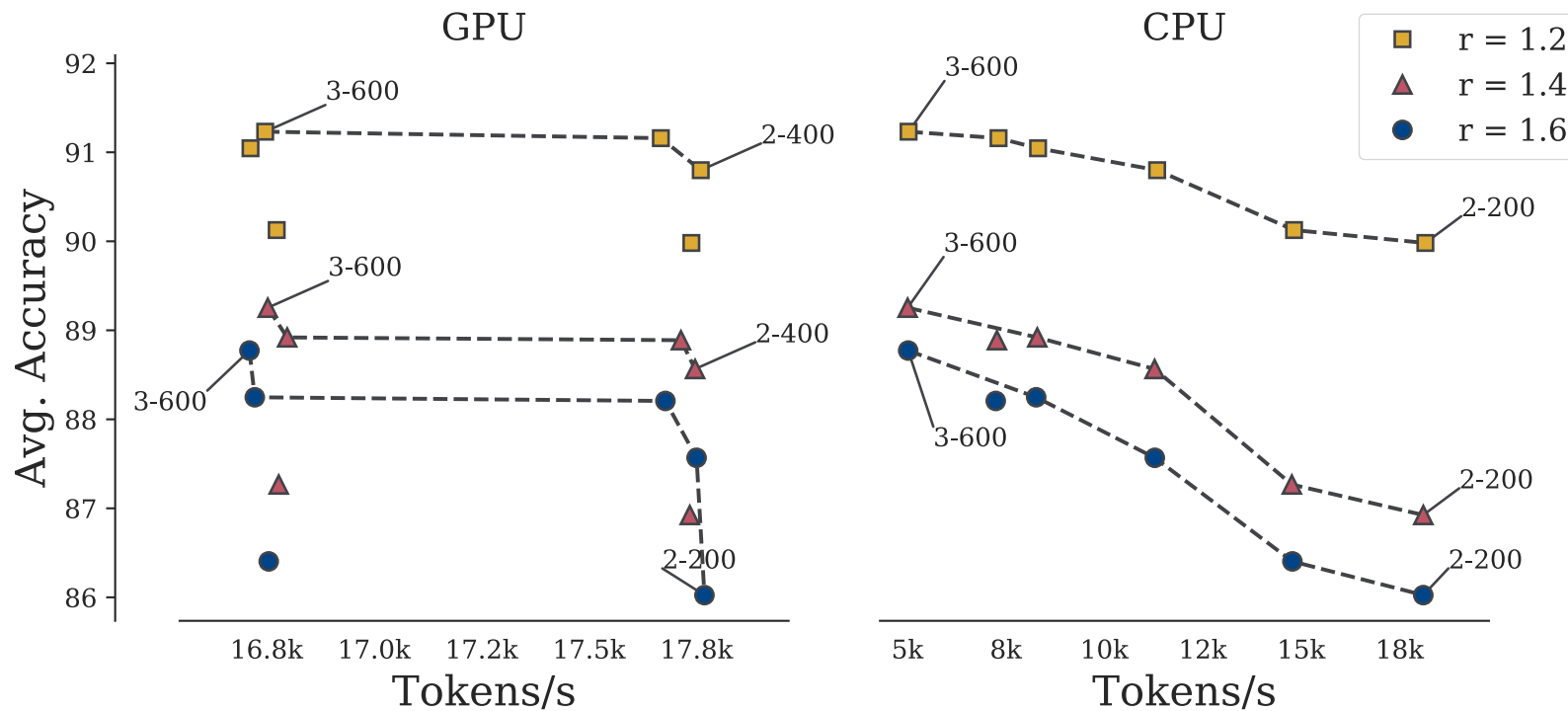


CHUNK AND PASS PARSING

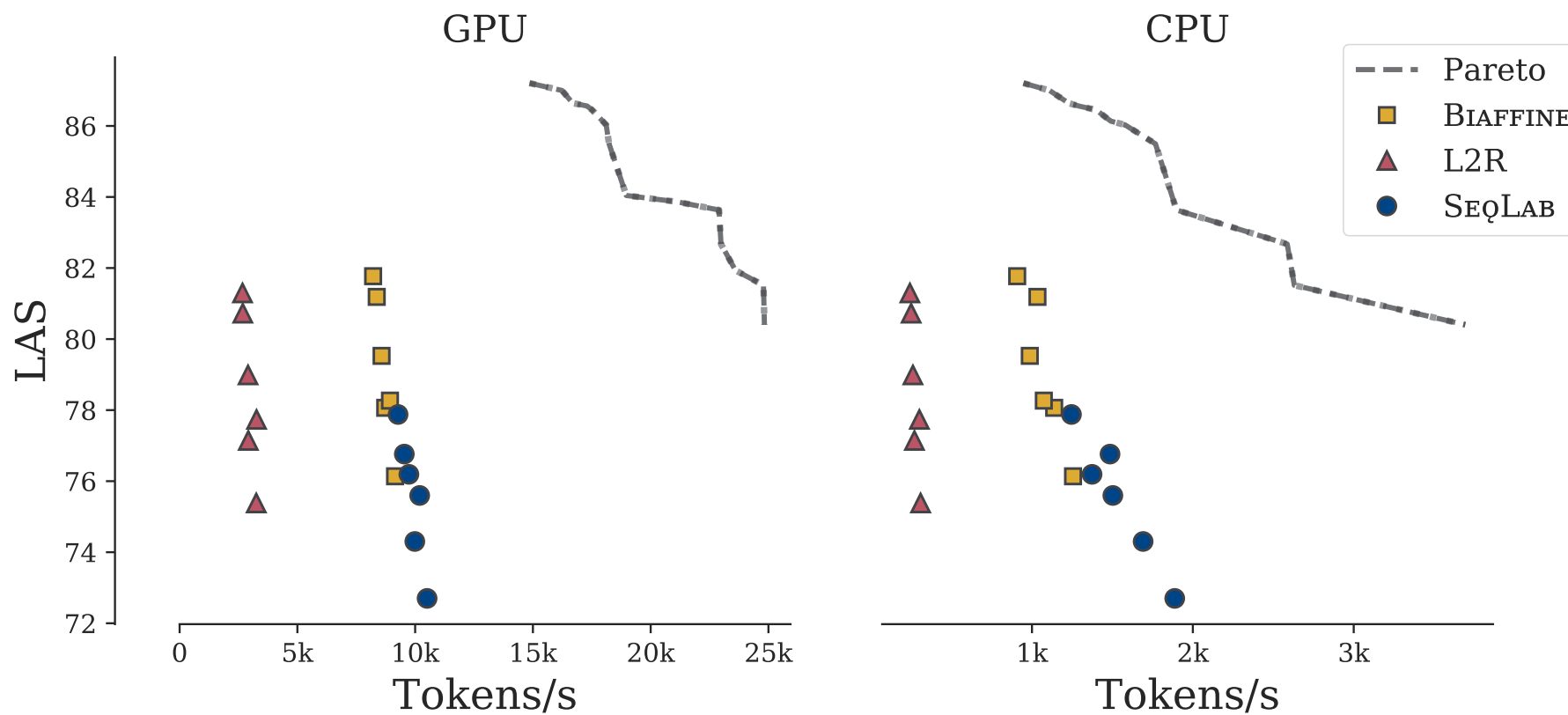
- Normalised PMI (used various thresholds to obtain different amount of compression)
- Compared to/with leading systems (biaffine, l2r, and sequence labelling).
- All same network type (BiLSTM).

Chunker performance (BiLSTM w/ fastText and char embeddings)

Varying compression and network size. Reasonable performance.
Extends to labelled edge predictions.
(Chinese, Hindi, Korean, and Polish)

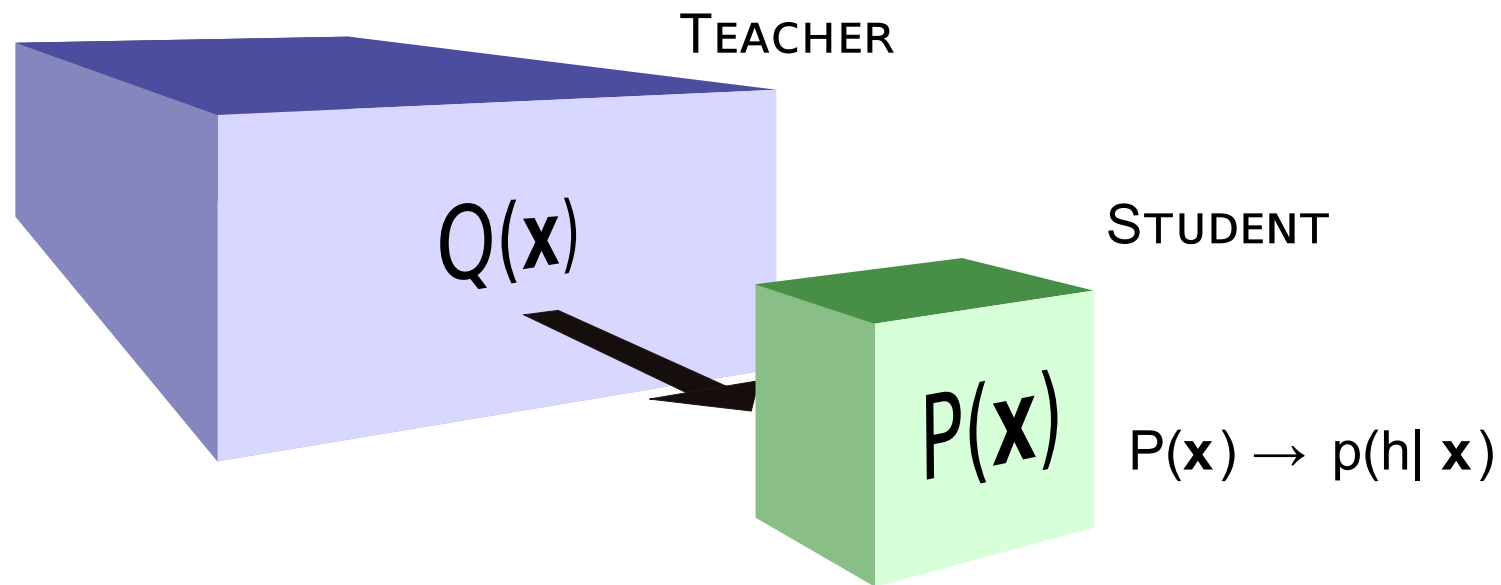


PARETO FRONT FOR UD (ZH, HI, KO, PL)



DISTILLATION

TEACHER-STUDENT DISTILLATION



$$\mathcal{L}_{\text{KL}}(Q(\mathbf{x}), P(\mathbf{x})) + \mathcal{L}_{\text{CE}}(p(h|\mathbf{x}))$$

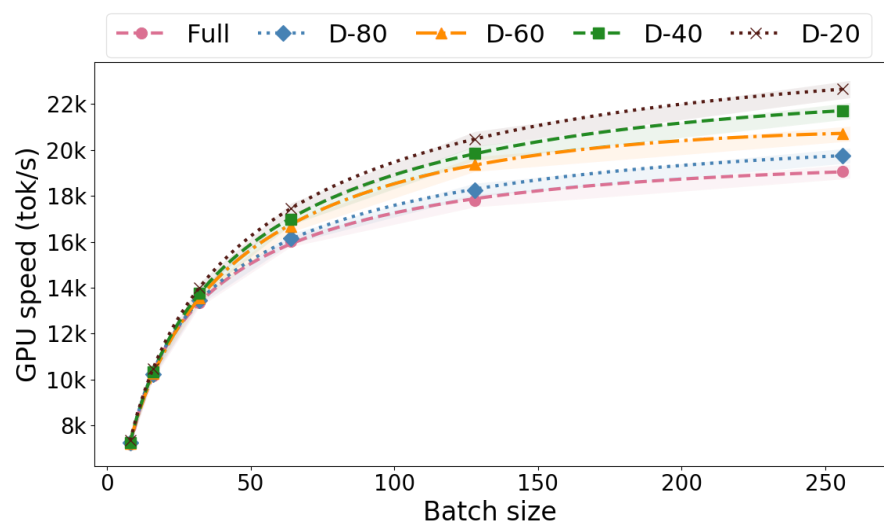
EXPERIMENTS

- Used Biaffine (fast and accurate).
- Compared against teacher and small models.
- Used UDv2.4 (Ancient Greek, Chinese, English, Finnish, Hebrew, Russian, Tamil, Uyghur, and Wolof). Treebanks based on subset used in de Lhoneux et al. (2017).¹
- Compress to 20%, 40%, 60%, and 80% of original model.
- fastText and gold POS tags. (°_°)

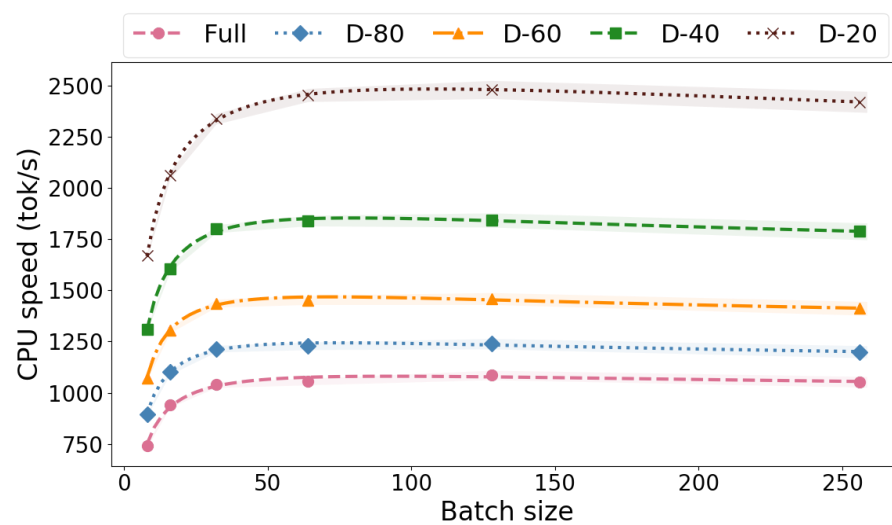
¹de Lhoneux, M., Stymne, S. and Nivre, J. *Old school vs. new school: Comparing transition-based parsers with and without neural network enhancement*, 2017

SPEED

GPU

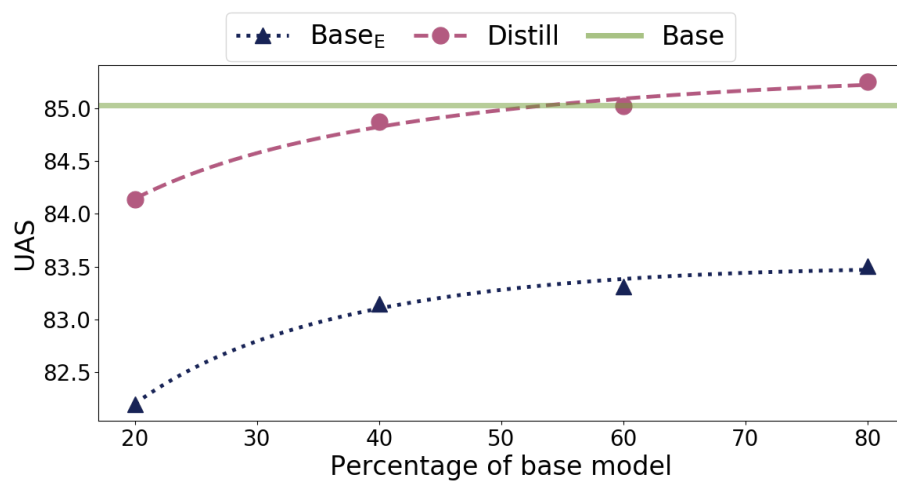


CPU

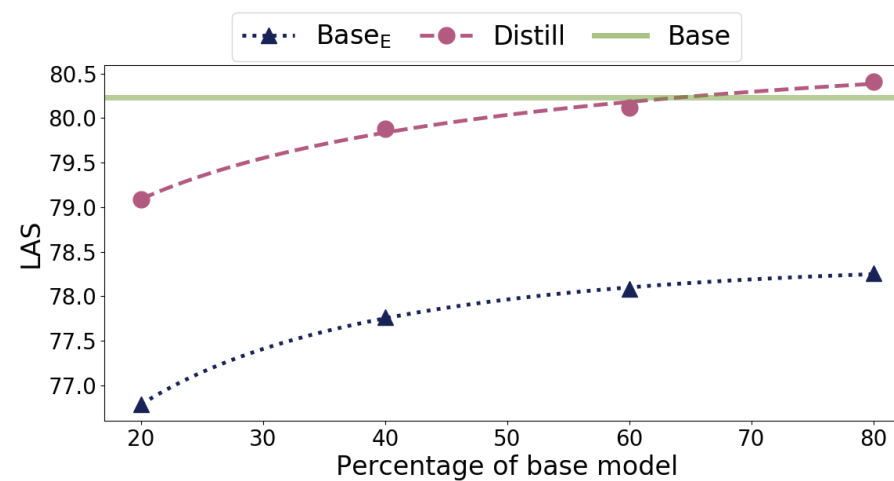


ACCURACY

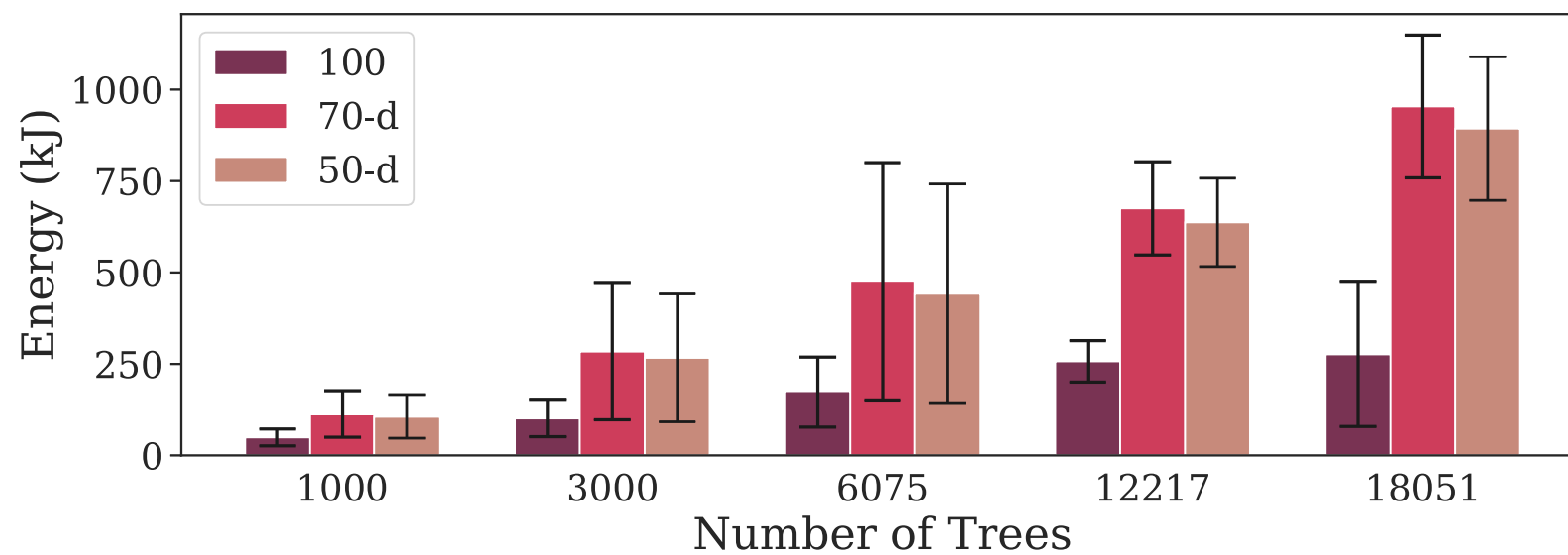
UAS



LAS



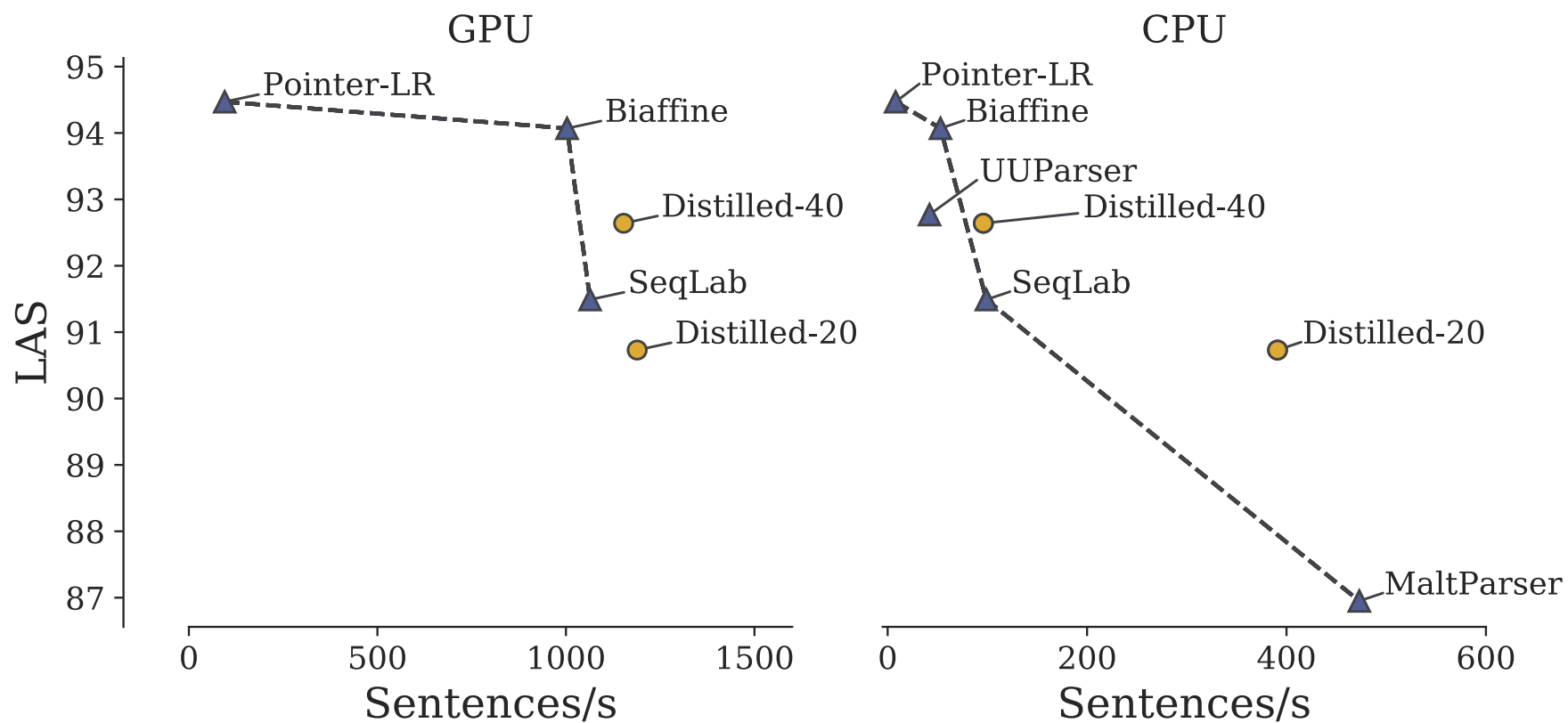
ENERGY COST



INCONSISTENT

- Subsequent work using only randomly initialised word and character.
- Didn't outperform small baselines.

PARETO FRONT FOR PTB (WSJ)



PART II

EVALUATING PARSERS

DEPENDENCY DISPLACEMENT

TRANSITION-BASED ALGORITHMS

Different transition-based algorithms perform differently on different treebanks

Perhaps because certain algorithms are inherently biased to creating edges that match given languages/treebanks?

- Used a non-NN parser, MaltParser.¹
- Contains multiple algorithms.
- Differences between algorithms observed still seen in NN implementations.²
- 76 treebanks from UD v2.2.

¹Nivre, J. et al., *MaltParser: A language-independent system for data-driven dependency parsing*, 2007.

²de Lhoneux, M., Stymne, S. and Nivre, J. *Old school vs. new school: Comparing transition-based parsers with and without neural network enhancement*, 2017

Transition-based algorithms

STACK	BUFFER
	Estoy muy ca..
SHIFT b0 to top of STACK	
Estoy	muy cansado
SHIFT b0 to top of STACK	
Estoy muy	cansado por...
REDUCE(right-advmod) remove s0 from STACK	
Estoy	cansado por...

advmod
muy cansado

ALGORITHMS

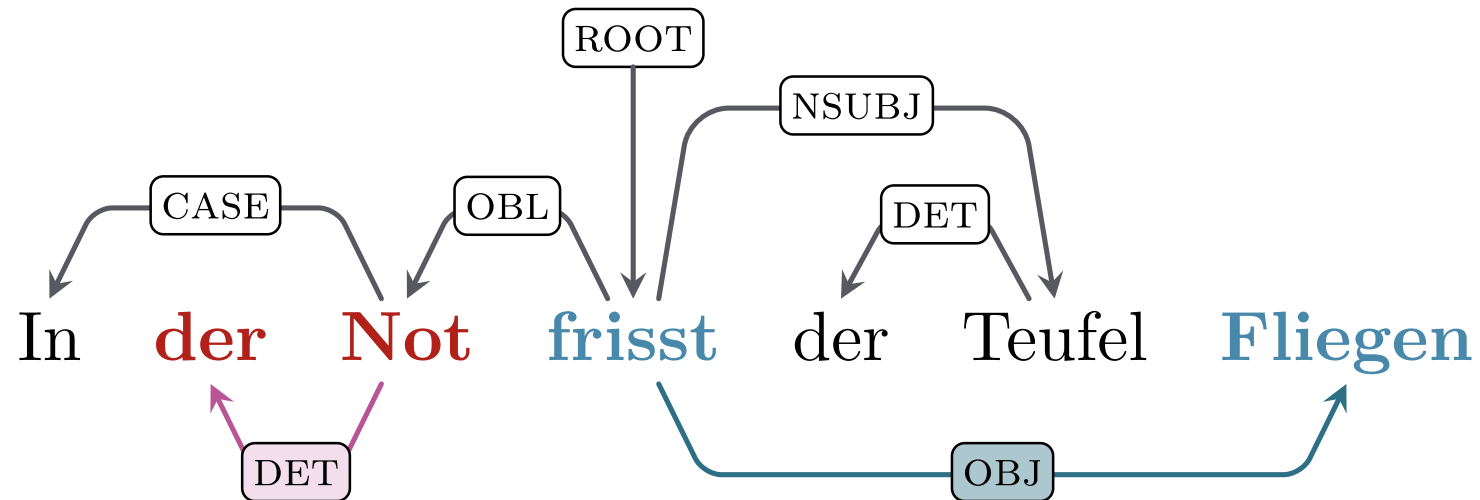
Projective

- Arc Standard
- Arc Eager
- Covington Projective

Non-projective

- Arc Swap
- Covington Non-projective

DEPENDENCY DISPLACEMENT

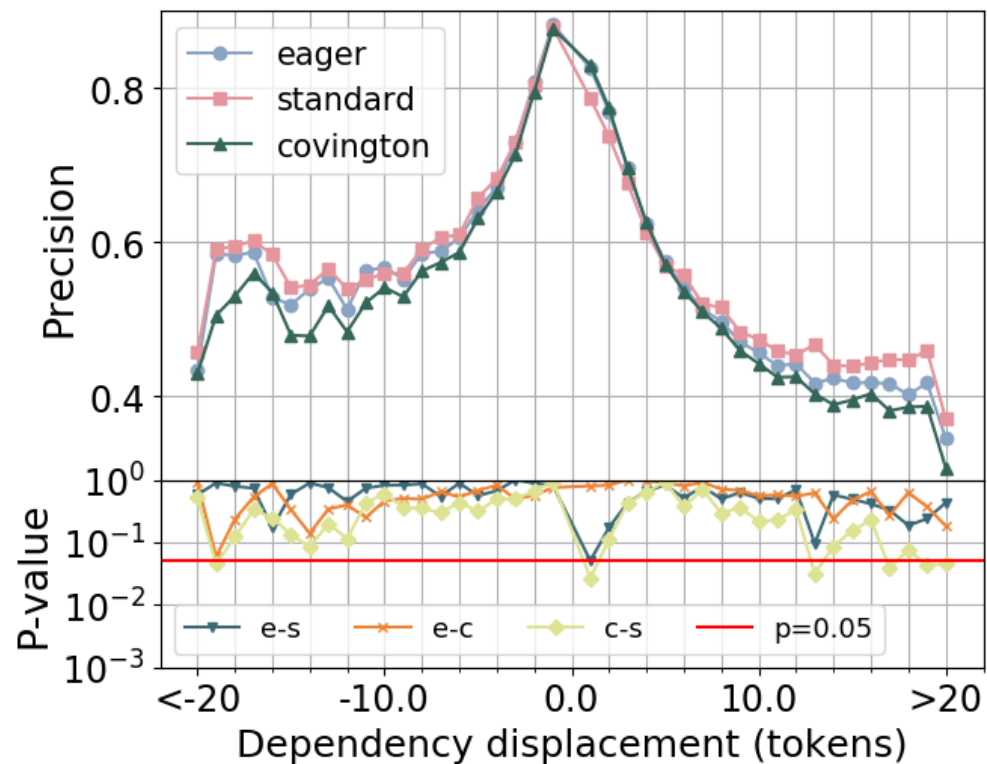


DET edge between **der** and **Not**: -1 (2-3)

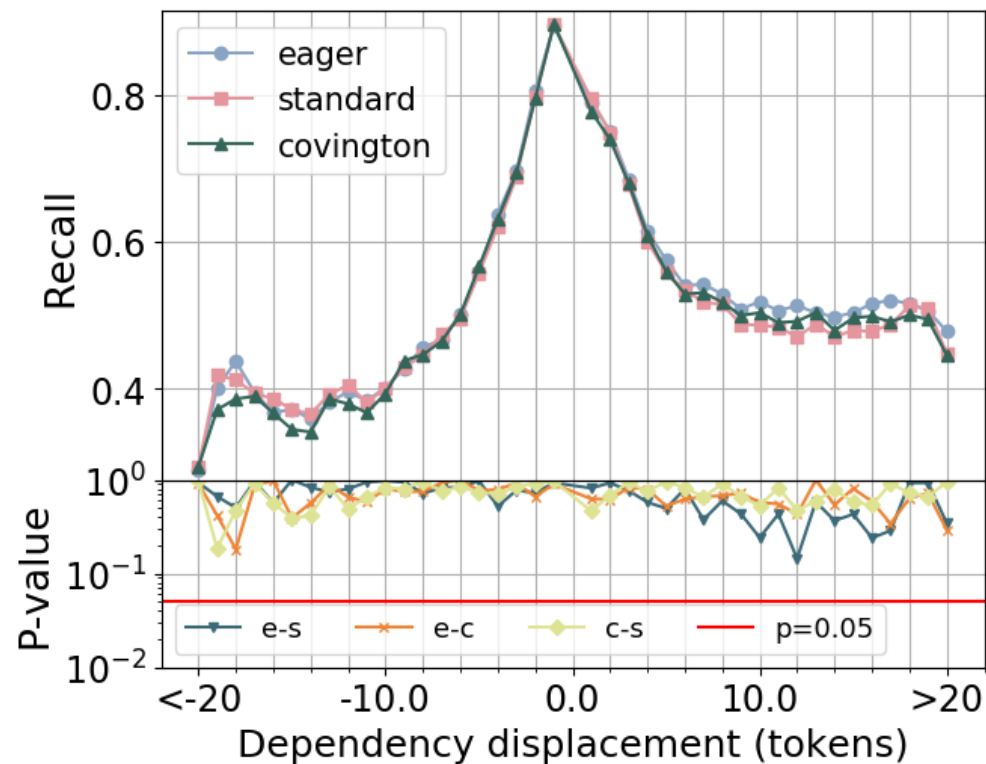
OBJ edge between **frisst** and **Fliegen**: 3 (7-4)

PROJECTIVE ALGORITHMS

Precision

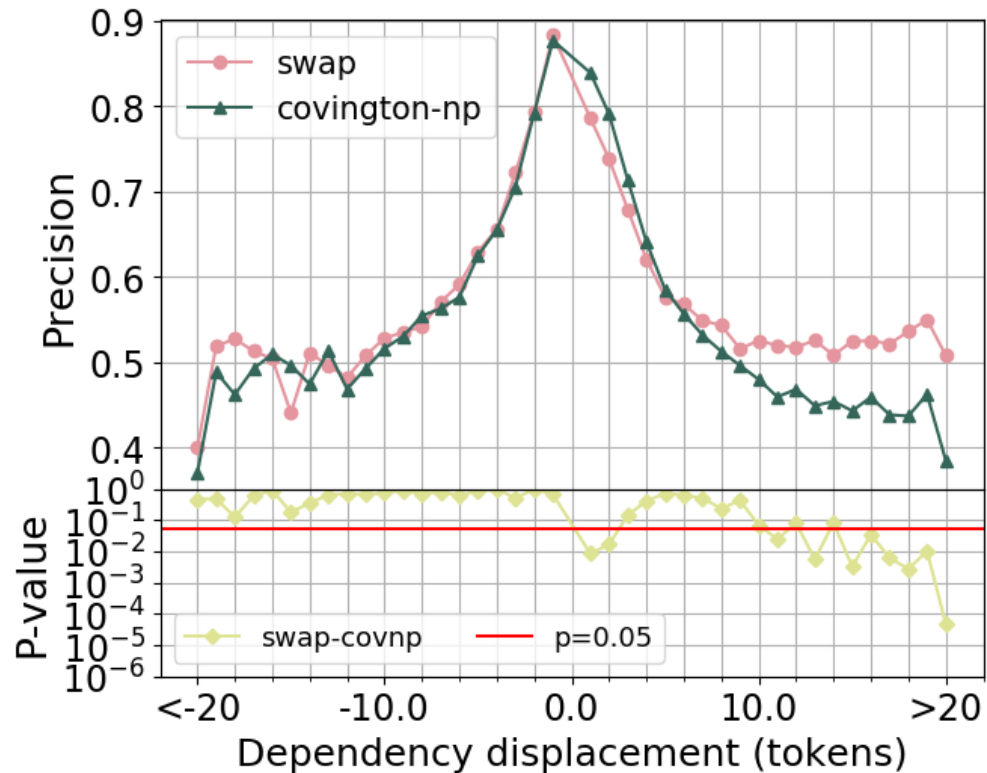


Recall

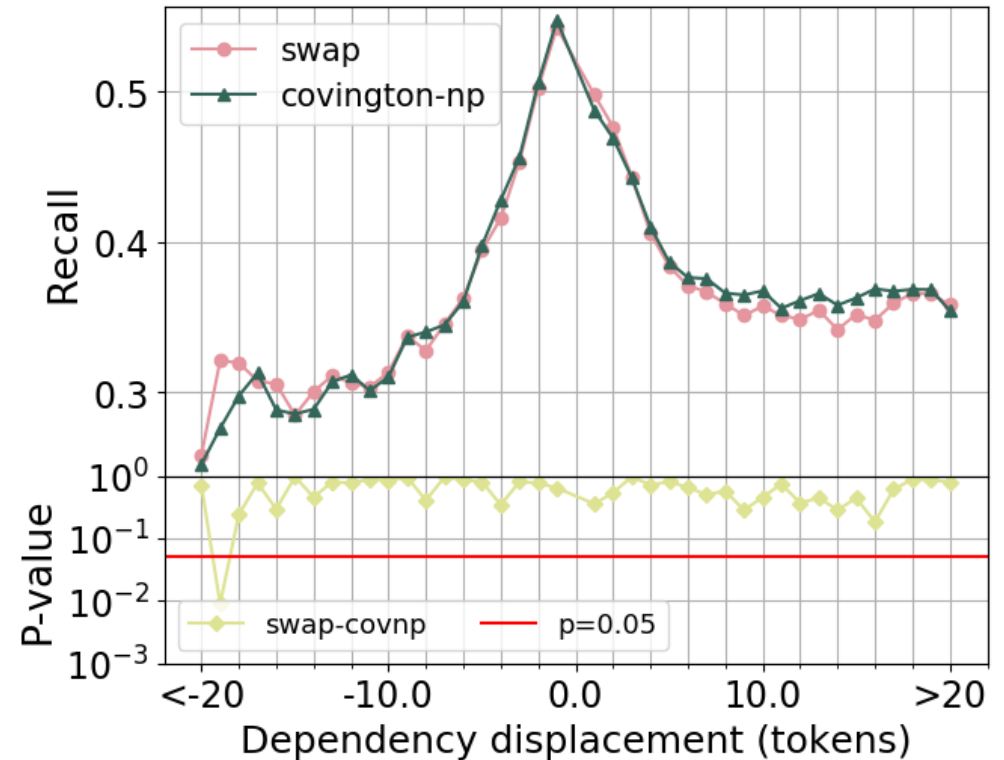


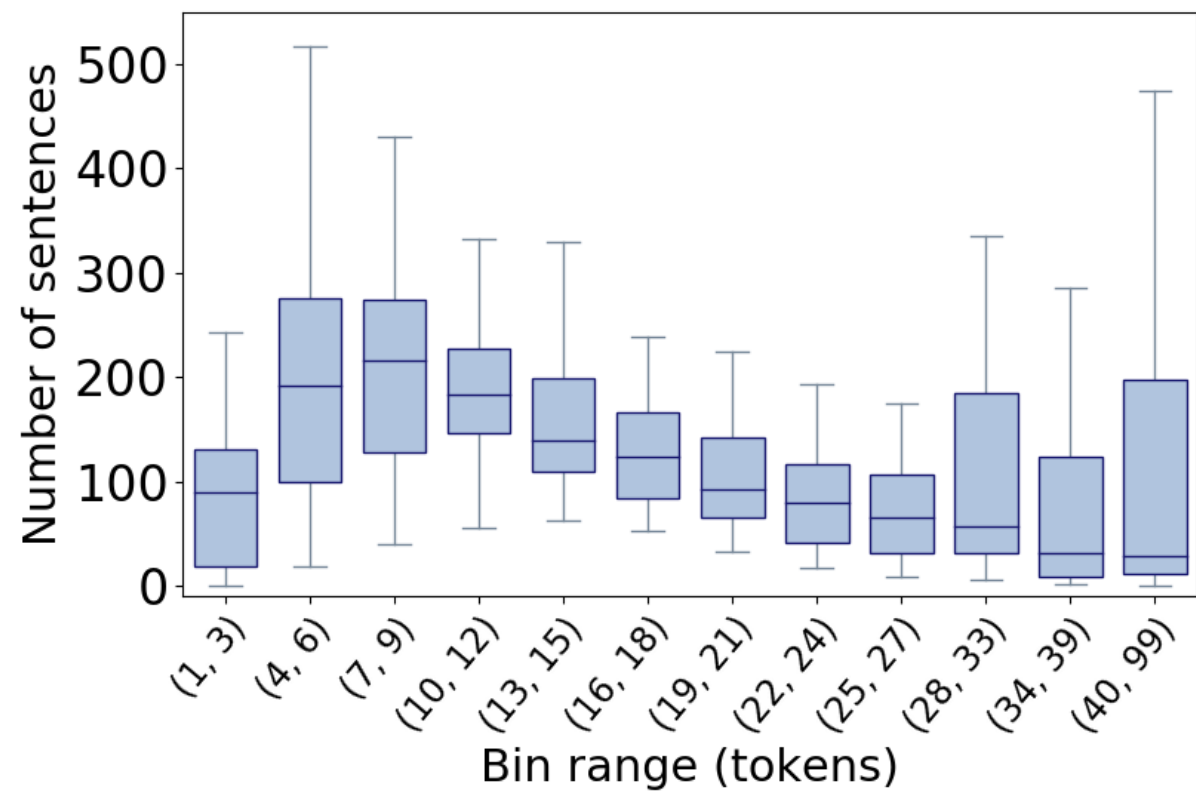
NON-PROJECTIVE ALGORITHMS

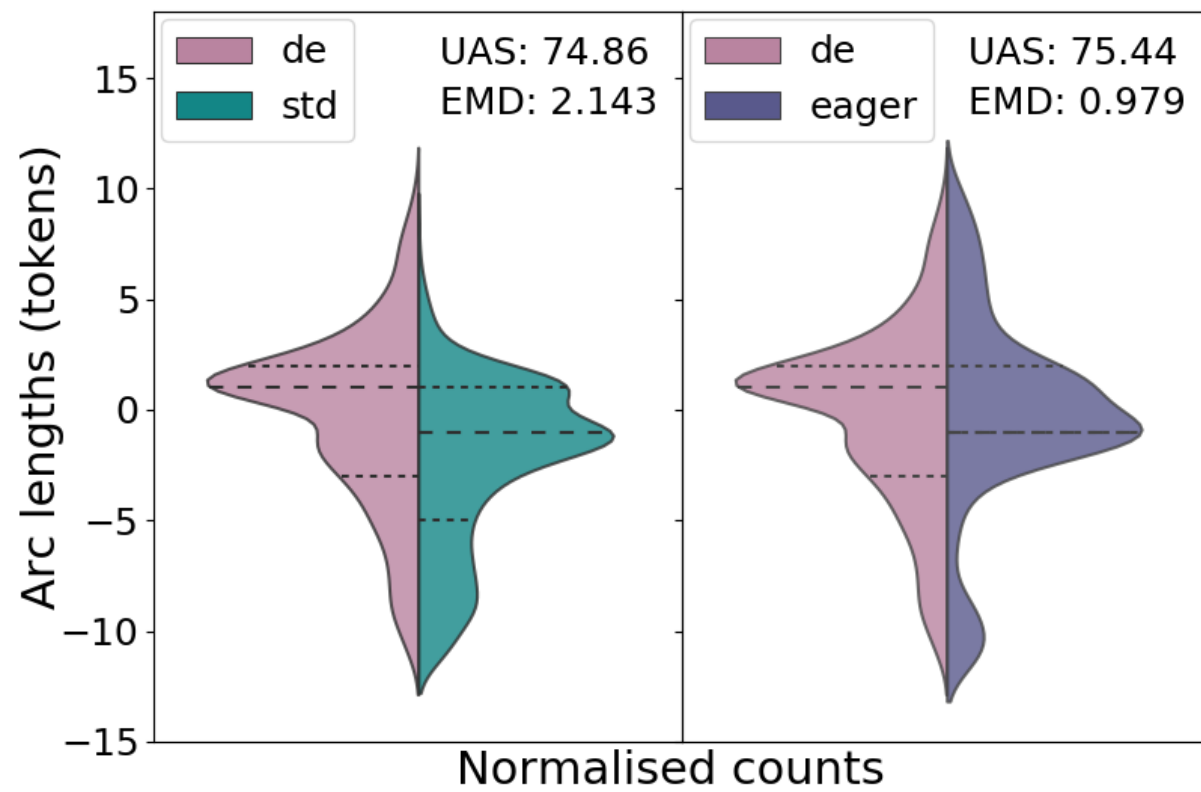
Precision

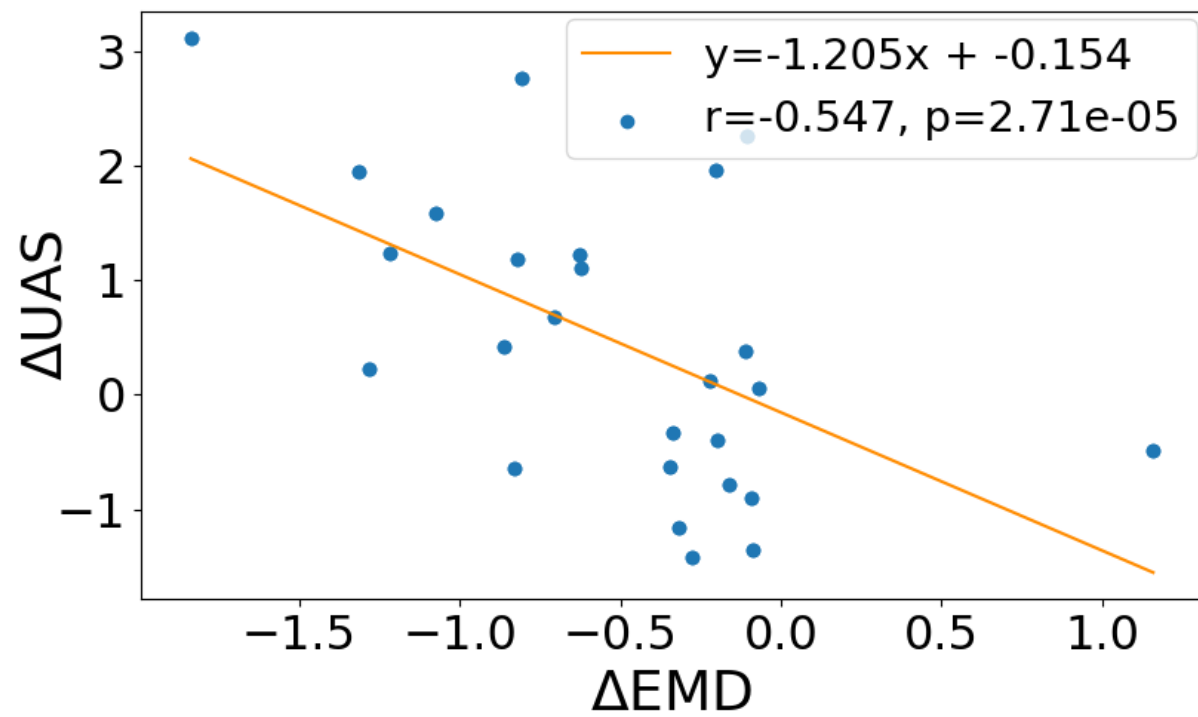


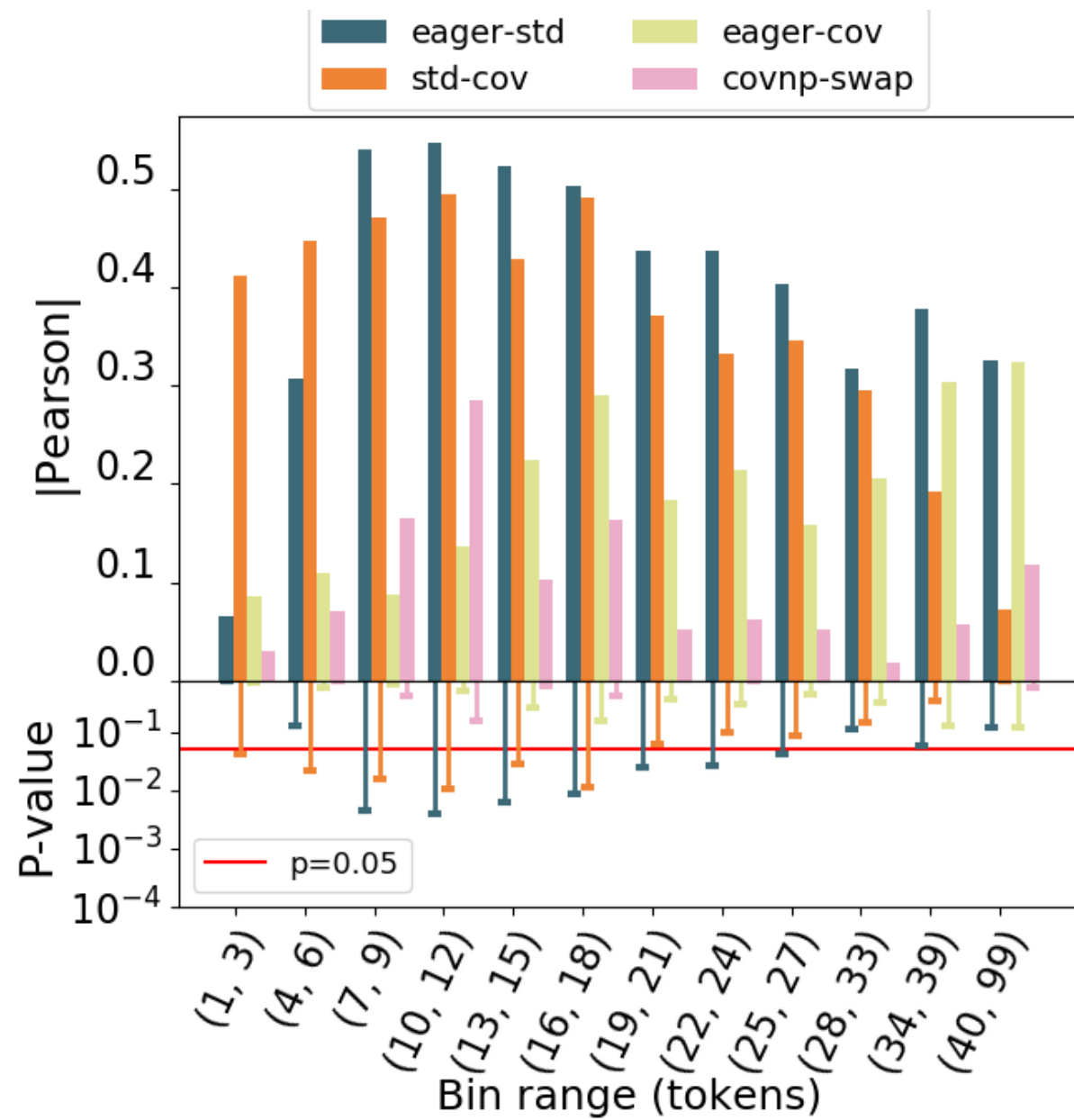
Recall



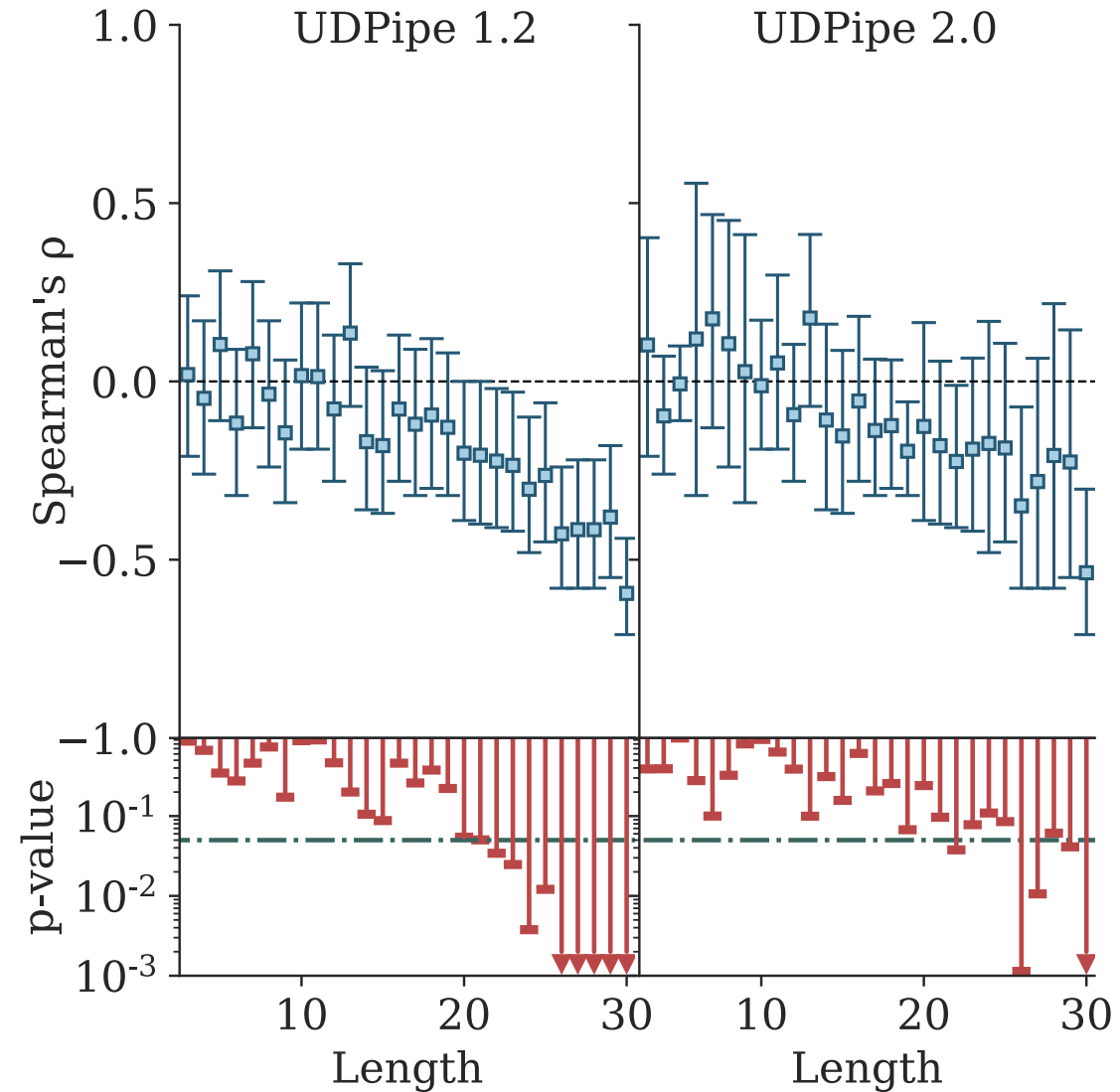




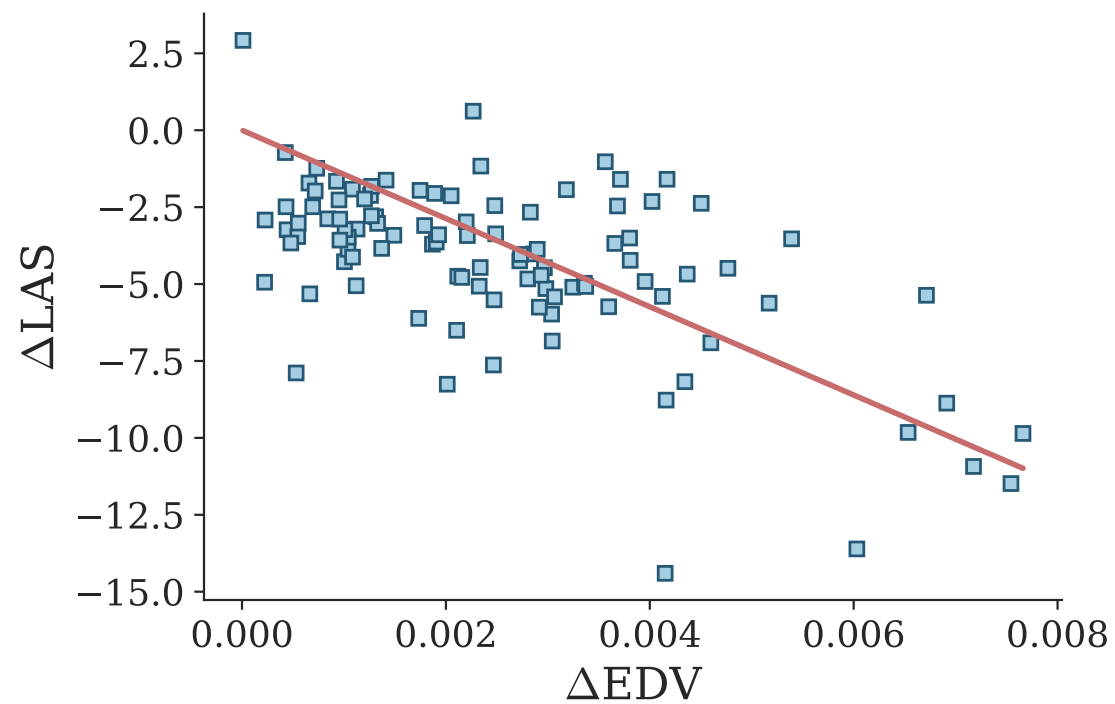




TRAINING VS TEST DATA



ADVERSARIAL SPLITS



UNIVERSAL POS TAGS

PARSERS

UUParser

Transition-based parser
out of Uppsala developed
from TB BIST Parser.^{1,2}

Biaffine

Graph-based parser
developed from GB BIST
Parser.³

word⊕char⊕upos

External pre-trained word embeddings, mainly fastText.
Same treebanks from distillation work.

¹Kiperwasser, E. and Goldberg, Y. *Simple and accurate dependency parsing using bidirectional LSTM feature representations*, 2016

²Smith, A., de Lhoneux, M., Stymne, S. and Nivre, J. *An investigation of the interactions between pre-trained word embeddings, character models and POS tags in dependency parsing*, 2018

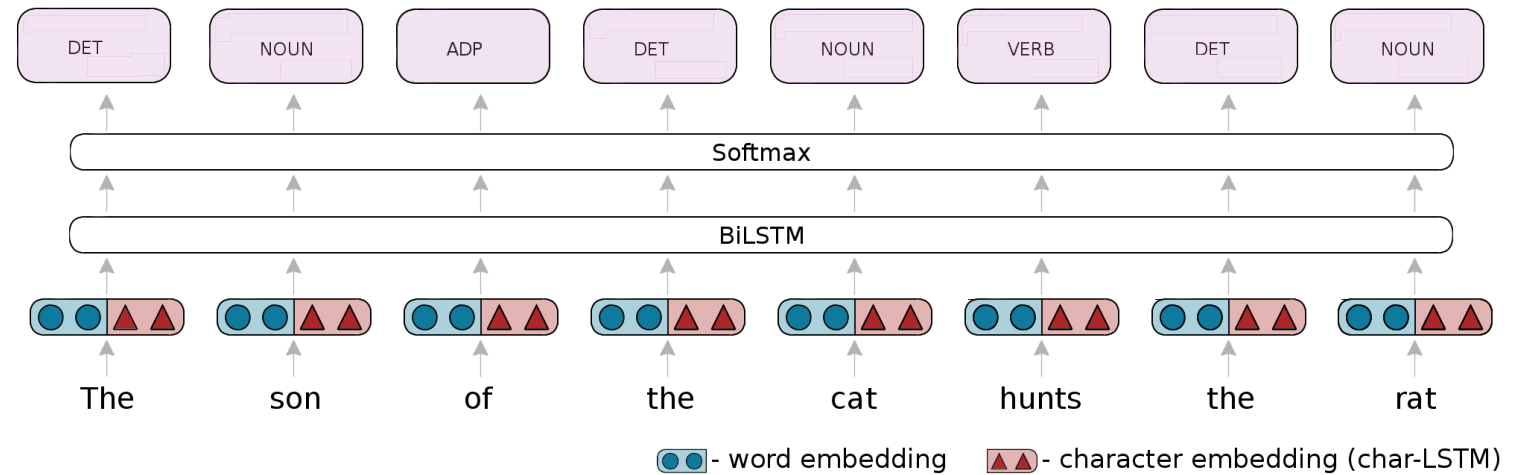
³Dozat, T. Manning, C.D., *Deep biaffine attention for neural dependency parsing*, 2017

CONTROLLING UPOS ACCURACY

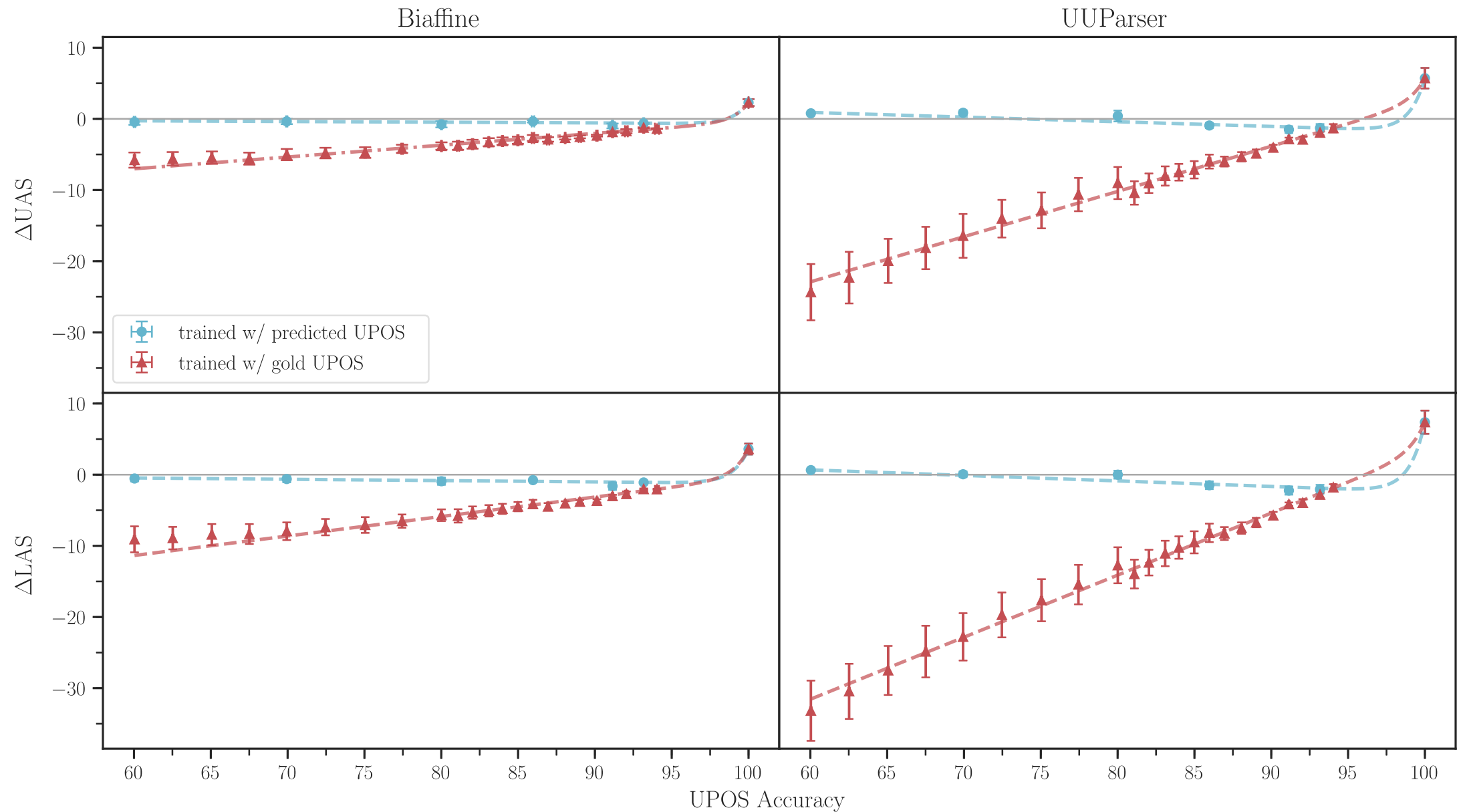
NCRF++ BiLSTM SL framework

BINS:: 2.5 ± 0.3 from 60 to 80 and 1 ± 0.3 from 80 onwards.

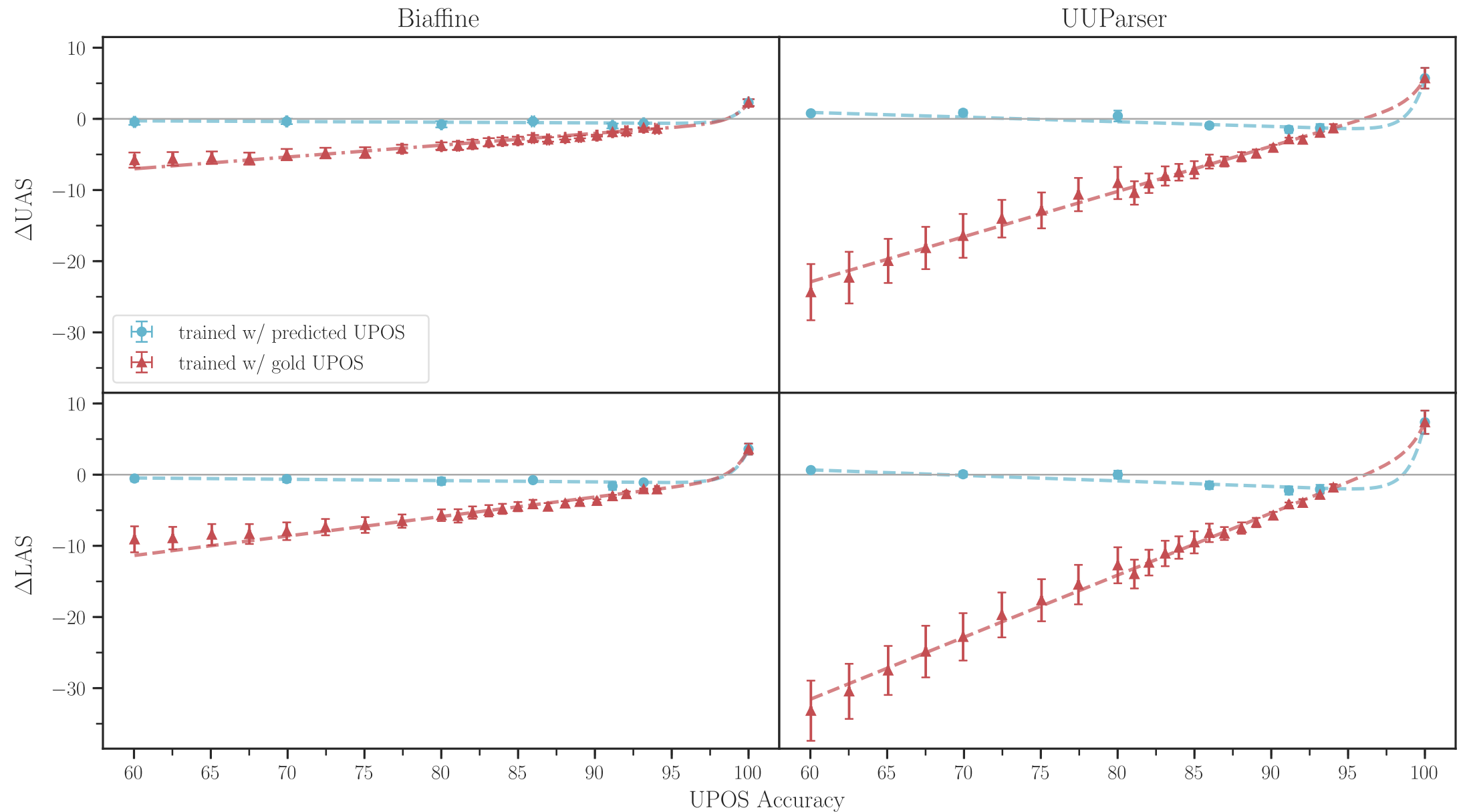
PARSERS: For training used gold tags and a subset of accuracy bins (60, 70, 80, 86, 91, and 93).



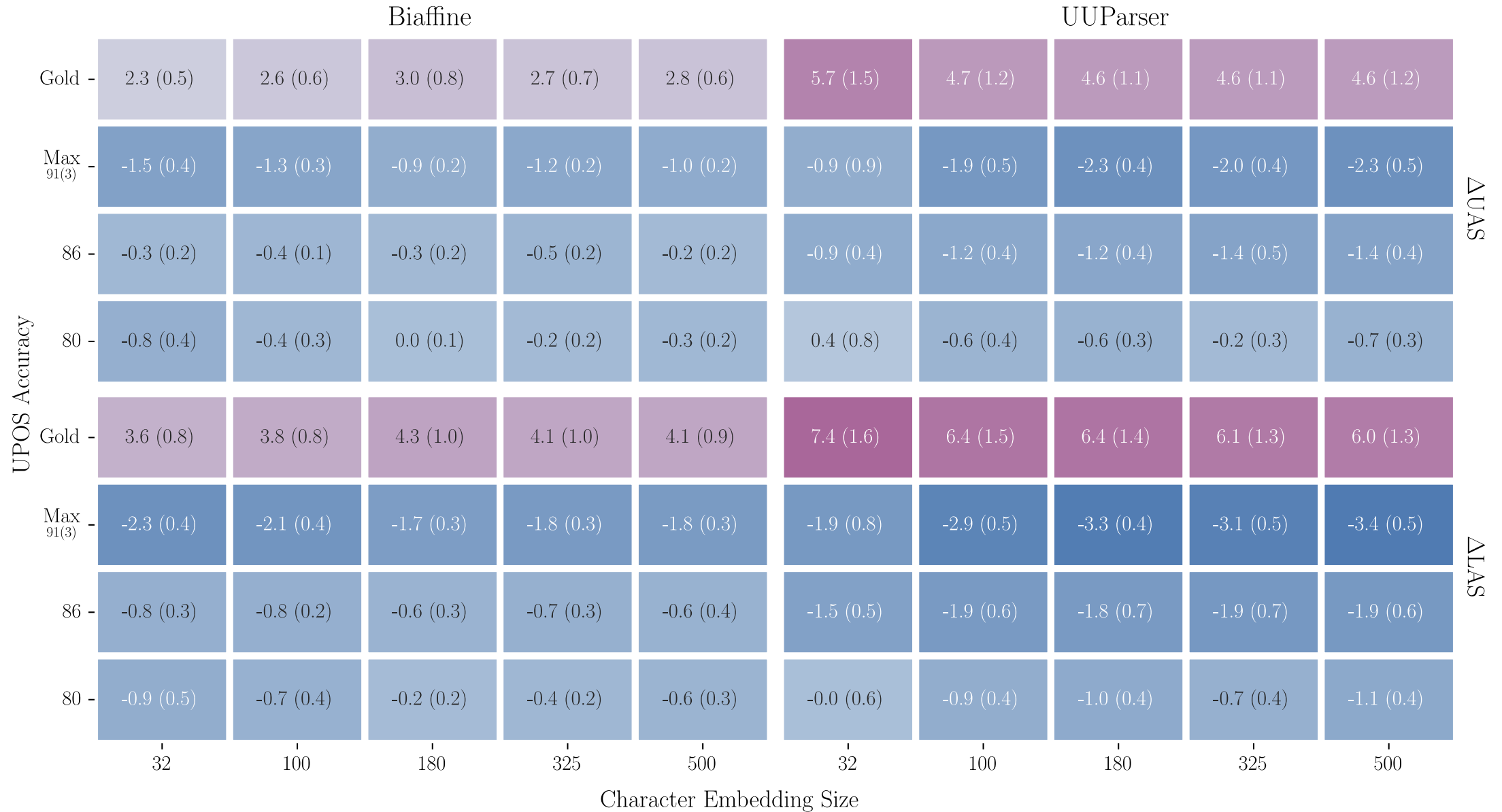
EXPERIMENT 1



EXPERIMENT 1

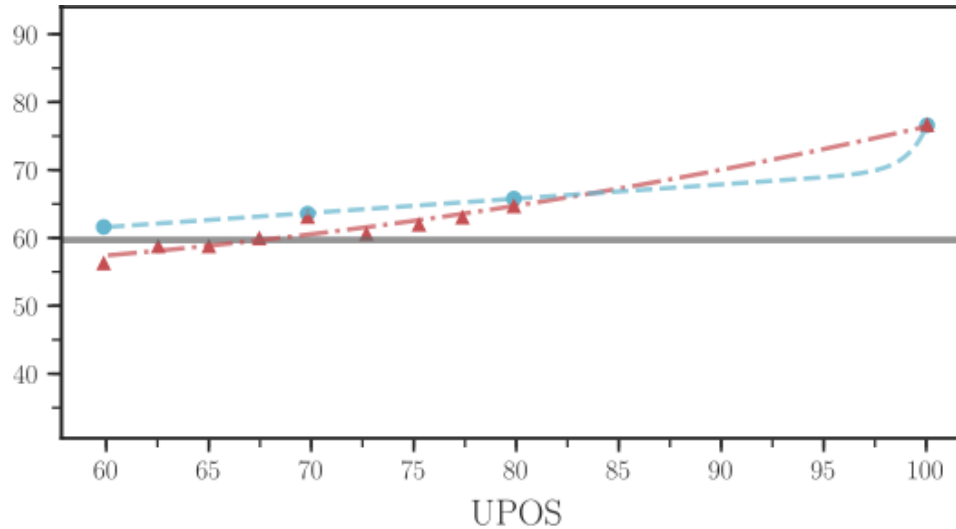


EXPERIMENT 2

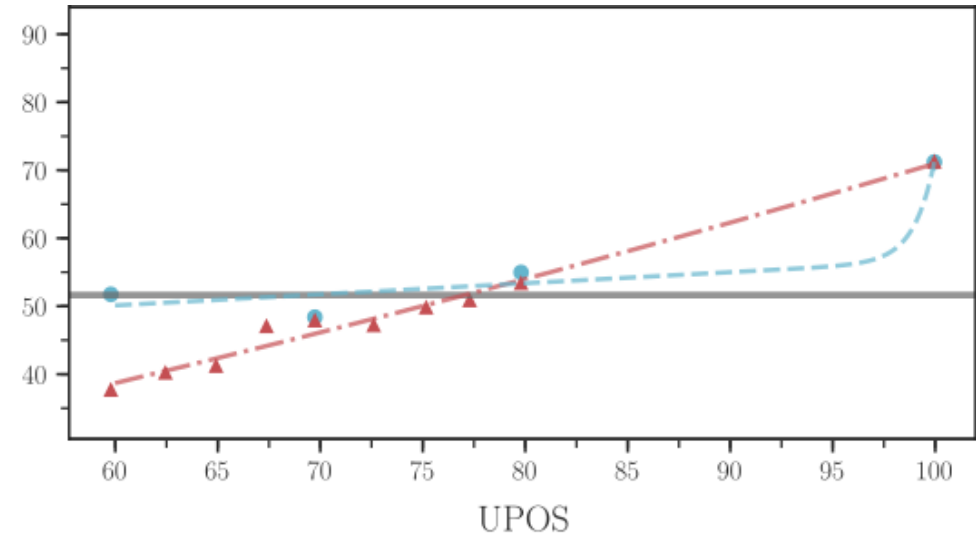


TAMIL RESULTS (~400 SENTENCES)

UAS



LAS



Still some improvement with low-accuracy taggers.

VERY LOW-RESOURCE

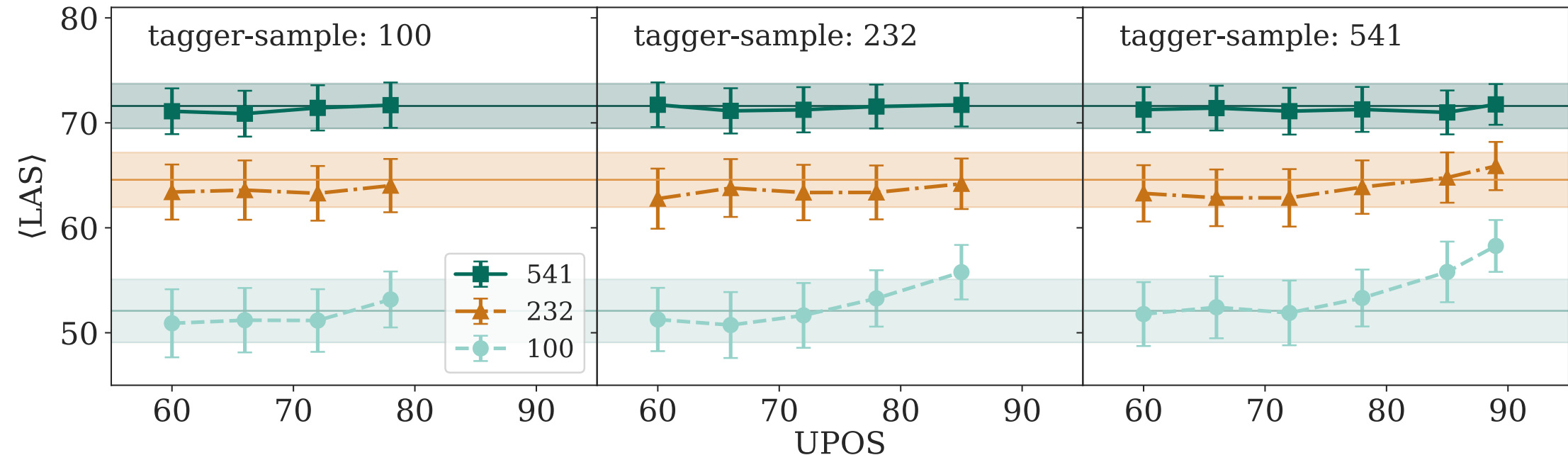
	UPOS			LAS		
	Single	Multi	None	Pred	Gold	Multi
bxr	48.72	48.34	10.45	12.36	20.31	14.41
kk	53.37	52.14	22.48	21.63	36.66	23.50
kmr	50.56	53.73	19.16	18.31	35.54	21.58
olo	37.84	37.37	9.74	10.89	17.54	7.59
hsb	53.44	47.28	18.36	20.03	41.88	14.66
avg	48.79	47.77	16.04	16.64	30.39	16.25

FAIRLY LOW-RESOURCED

	UPOS			LAS		
	Single	Multi	None	Pred	Gold	Multi
be	92.82	87.29	61.82	64.91	68.87	62.28
gl	93.54	88.56	70.60	72.73	79.06	70.54
lt	79.25	71.51	37.17	35.94	48.30	38.96
mr	80.58	76.46	57.04	58.74	64.32	56.31
orv	87.77	81.60	49.53	51.34	60.24	50.33
ta	86.88	79.23	63.85	62.75	74.31	63.15
cy	91.77	86.41	72.10	72.93	80.71	73.00
avg	85.89	77.77	55.24	56.52	64.13	55.10

ARTIFICIAL LOW-RESOURCE

Indonesian GSD, Irish IDT, Japanese GSD, and Wolof WTB.

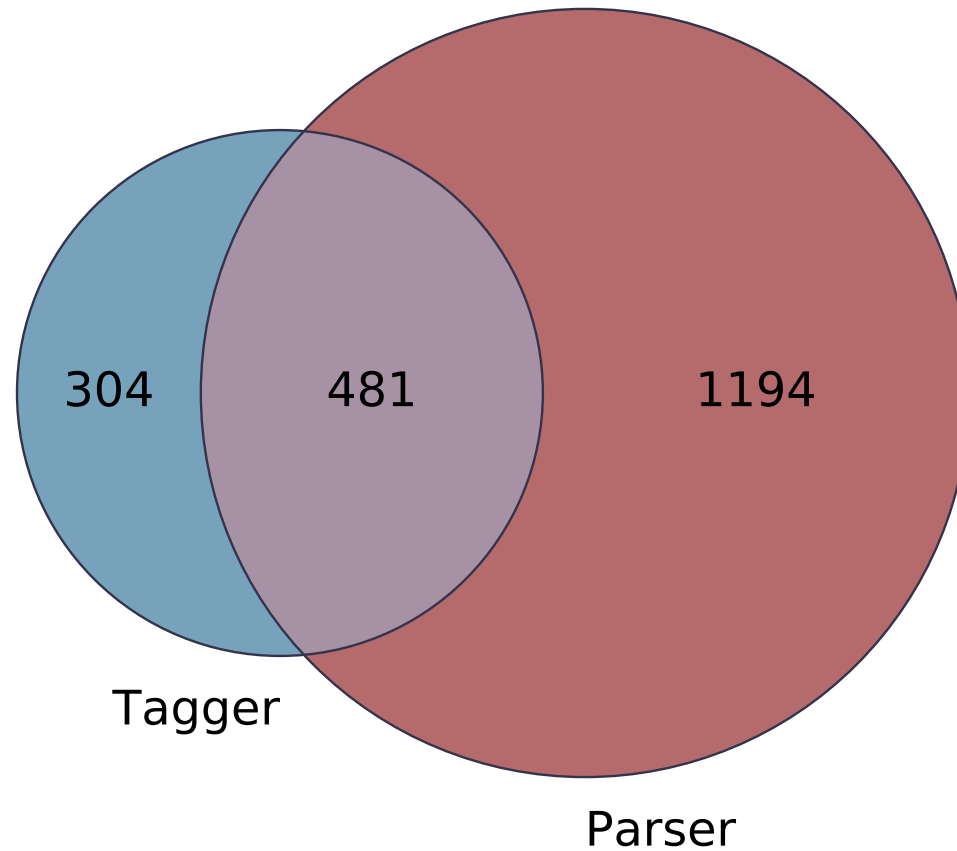


PROBING EXPERIMENT

Arabic, Basque, Finnish, Indonesian, Irish, Japanese, Korean, Tamil,
Turkish, Vietnamese, and Wolof

- Trained parser (Biaffine)
- Finetuned parser to predict POS tags
- Trained taggers with same system

PROBING EXPERIMENT



Union of errors (average over treebanks).

	None	Pred.	$M\text{-}E_T$	$M\text{-}E_P$	$M\forall E_T$	Gold
ar	83.29	82.87	84.17	84.06	84.45	84.73
eu	81.12	81.14	82.33	82.62	83.13	84.45
fi	85.96	86.04	86.88	87.09	87.61	88.80
id	79.04	78.95	82.20	82.69	81.08	82.95
ga	76.13	76.57	76.62	76.65	77.46	77.90
ja	93.15	92.72	94.41	94.38	94.39	95.30
ko	85.40	85.86	87.53	87.82	87.44	88.52
ta	65.61	64.50	70.24	66.67	66.01	71.95
tr	66.67	67.68	67.62	67.66	67.84	68.86
vi	58.43	60.09	65.42	66.75	65.18	70.87
wo	77.87	78.49	82.03	81.39	81.11	85.41
avg	77.52	77.72	79.95	79.80	79.61	81.79

Masking Experiment

- None — no tags.
- Pred. — predicted tags.
- $M\text{-}E_T$ — gold tags except tagger errors.
- $M\text{-}E_P$ — gold tags except parser errors.
- $M\forall E_T$ — gold tags only tagger errors.
- Gold — all gold tags.

END

CONCLUSION

Developing

- Chunk-and-Pass
- Distillation

Evaluating

- Edge displacement
- POS tags

COLLABORATORS (WORK PRESENTED)

Carlos Gómez Rodríguez

Mathieu Dehouck

David Vilares