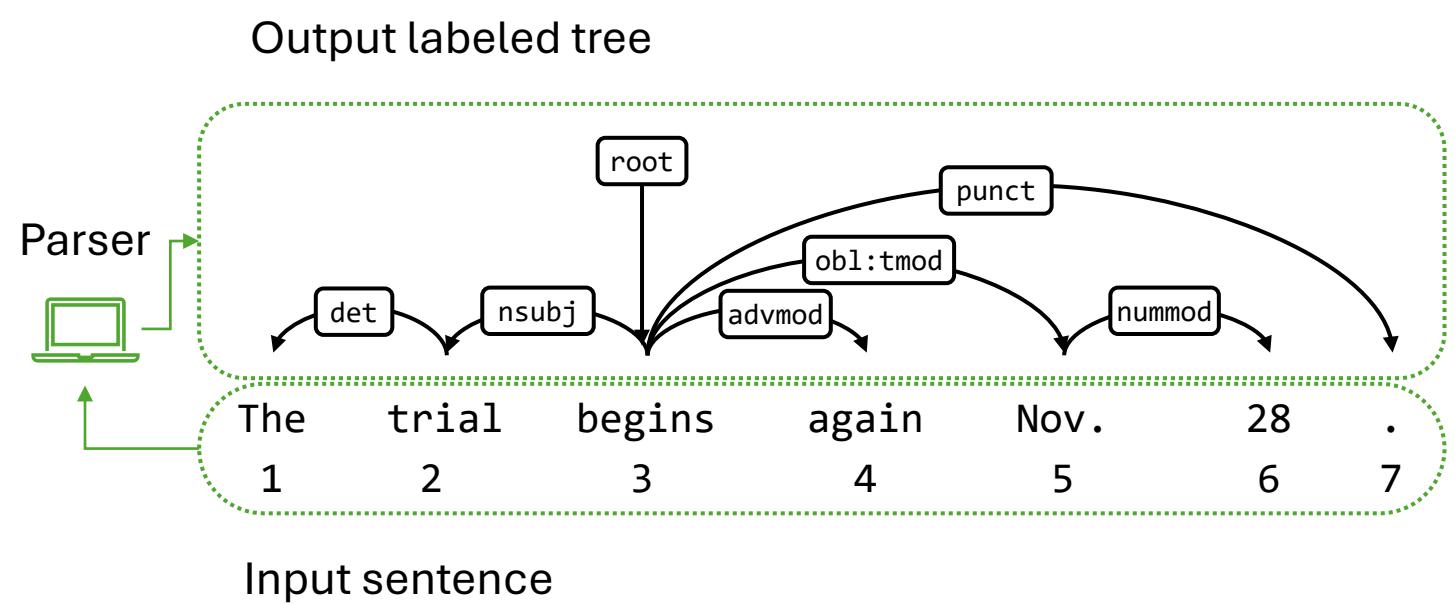


# Better Benchmarking LLMs for Zero-Shot Dependency Parsing

Ana Ezquerro, Carlos Gómez-Rodríguez and David Vilares

## ZERO-SHOT DEPENDENCY PARSING

- Evaluation of LLMs in a zero-shot setting (prompt).
- Difficult structured-task.



- Prompt with the CoNLL format:

In dependency parsing the CoNLL format for the sentence <The trial begins again Nov. 28 .> is:									
ID	FORM	LEMMA	UPOS	XPOS	FEATS	HEAD	DEPREL	DEPS	MISC
1	The	-	-	-	-	2	det	-	-
2	trial	-	-	-	-	3	nsubj	-	-
3	begins	-	-	-	-	0	root	-	-
4	again	-	-	-	-	3	advmod	-	-
5	Nov.	-	-	-	-	3	obl:tm	-	-
6	28	-	-	-	-	5	nummod	-	-
7	.	-	-	-	-	3	Punct	-	-

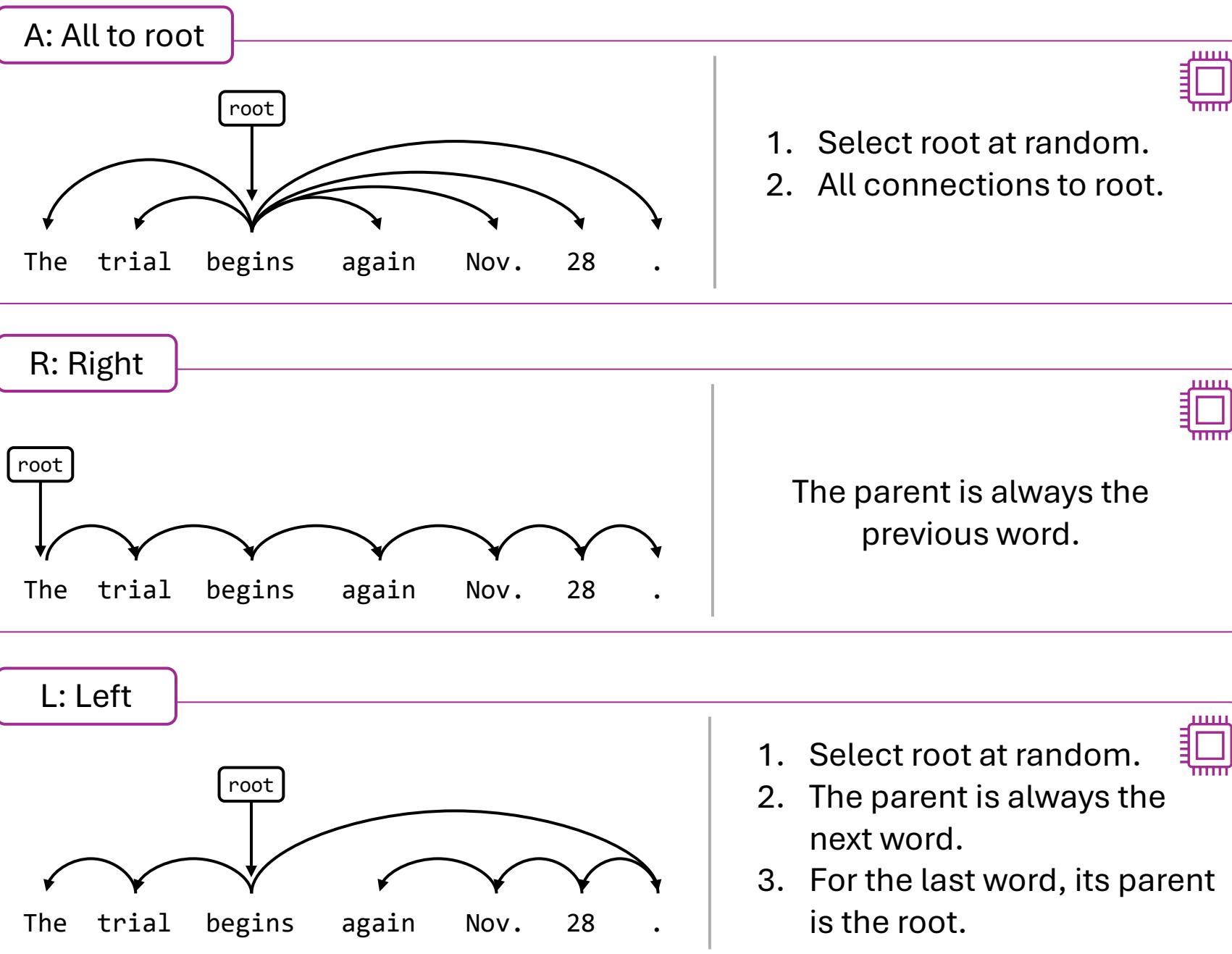
Now return the CoNLL format for the sentence: <What if Google Morphed into GoogleOS?>

Sure! This is the CoNLL format for the sentence <What if Google Morphed into GoogleOS?>:  
1 What \_ \_ \_ 0 nsubj \_ \_  
2 if \_ \_ \_ 4 mark \_ \_ \_  
3 Google \_ \_ \_ 4 nsubj  
4 Morphed \_ \_ \_ 0 root \_ \_  
5 into \_ \_ \_ 6 case \_ \_  
6 GoogleOS \_ \_ \_ 8 nmod \_ \_  
7 ? \_ \_ \_ 4 punct \_ \_  
Let me know if I can help you again!

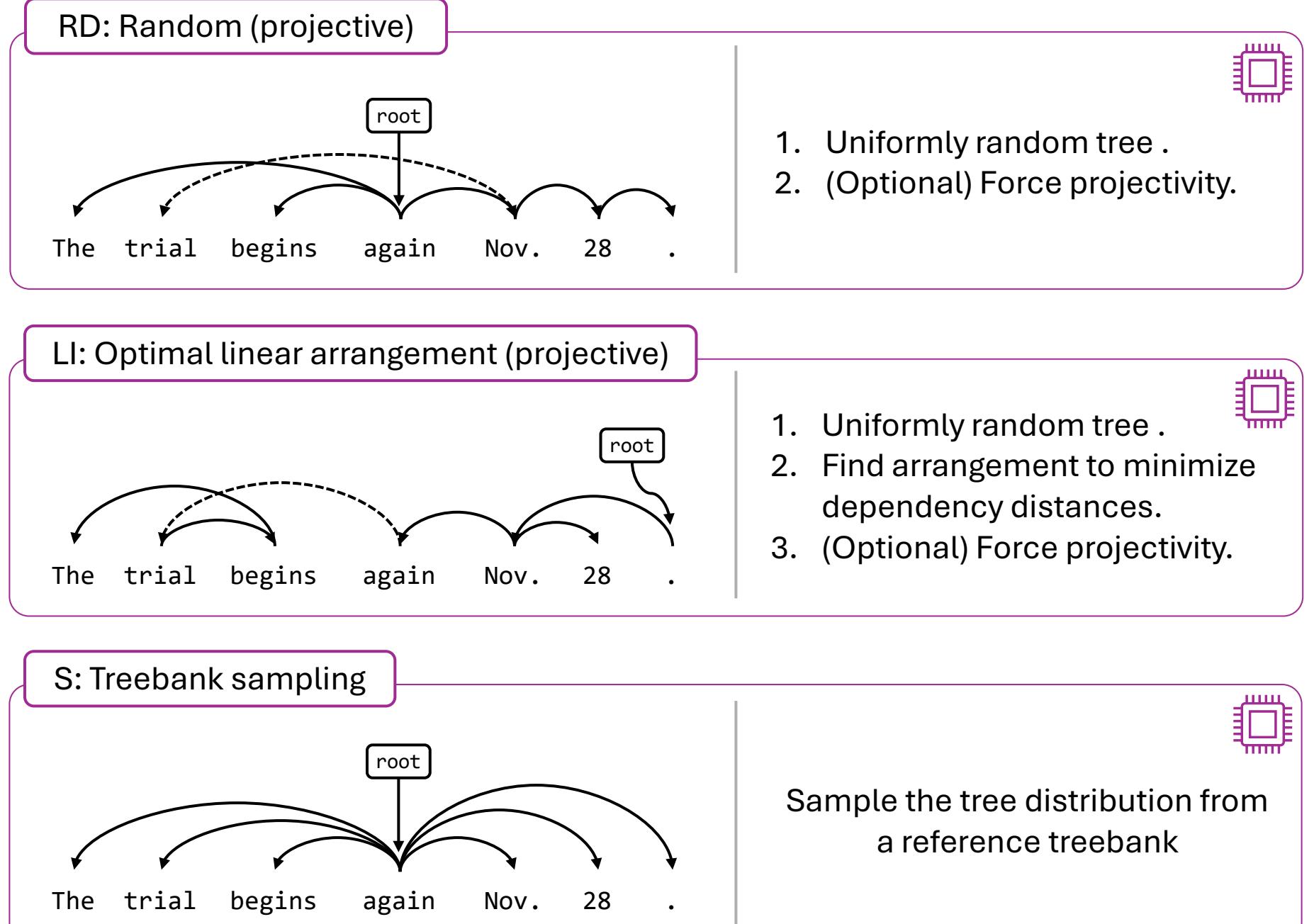
<https://github.com/anaezquerro/naipar>

## BETTER BENCHMARKING WITH NAÏVE BASELINES

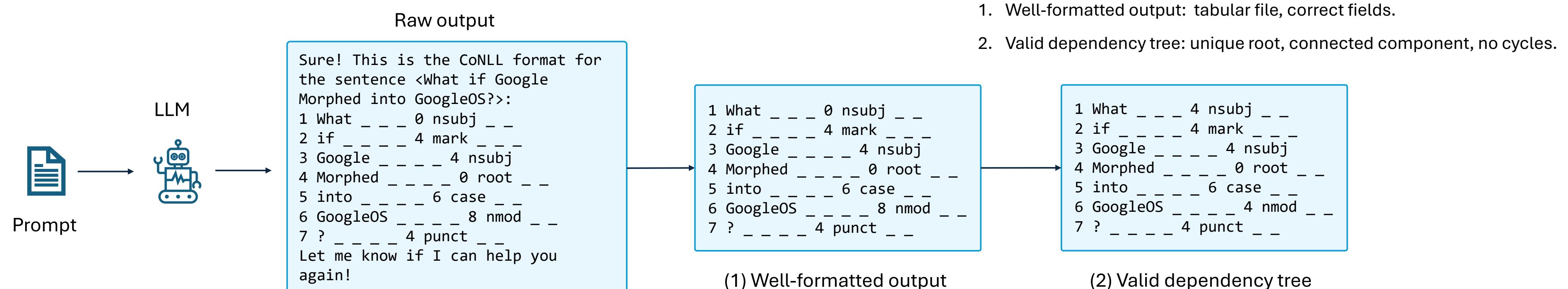
**Key idea:** Do LLMs perform better than uninformed baselines?



- No information about the input sentence (fully unsupervised).
- Only treebank sampling considers the length of the sentence and a reference treebank.



## POST-PROCESSING PIPELINE



## EXPERIMENTS

- Multilingual benchmark : English, French, German and Hindi from Universal Dependencies.
- Different instruction-based LLMs: Gemma (G), LLaMA (O) and Mistral (D).
- Half precision (minimal impact).

- Only the largest models (LLaMA and Mistral) outperform baselines.
- Considerable improvement with the latest versions of LLaMA (v3.2 vs v2).
- Mistral (22B) offers a trade-off between size and zero-shot performance.

