

# PCFGs Can Do Better: Inducing Probabilistic Context-Free Grammars with Many Symbols

Songlin Yang<sup>1</sup>, Yanpeng Zhao<sup>2</sup>, Kewei Tu<sup>1</sup>

<sup>1</sup>School of Information Science and Technology, ShanghaiTech University

<sup>2</sup>ILCC, University of Edinburgh



## Introduction

Neural PCFGs have been proved effective in unsupervised constituency parsing. However, due to the high complexity of the inside algorithm, they cannot scale to a relatively large number of grammar symbols. We use CP decomposition to decompose the binary rule probability tensor to decrease the complexity from cubic to at most quadratic in the number of grammar symbols. We enforce row/column normalization on the decomposed matrices to ensure the validity of the binary rule probability.

## CP decomposition on PCFGs

Start rule:

$$S \rightarrow A \quad A \in \mathcal{N}$$

$$r_{h_A} = \pi_{S \rightarrow A}, \quad r \in \mathbb{R}^n$$

Binary rule:

$$A \rightarrow BC, \quad A \in \mathcal{N}, \quad B, C \in \mathcal{N} \cup \mathcal{P}$$

$$T_{h_A, h_B, h_C} = \pi_{A \rightarrow BC}, \quad T \in \mathbb{R}^{n \times m \times m}$$

Unary rule:

$$T \rightarrow w, \quad T \in \mathcal{P}, \quad w \in \Sigma$$

$$Q_{h_T, h_w} = \pi_{T \rightarrow w}, \quad Q \in \mathbb{R}^{p \times q}$$

Recursive form of the inside algorithm

$$s_{i,j}^A = \sum_{k=i}^{j-1} \sum_{B,C} \pi_{A \rightarrow BC} \cdot s_{i,k}^B \cdot s_{k+1,j}^C \quad (1)$$

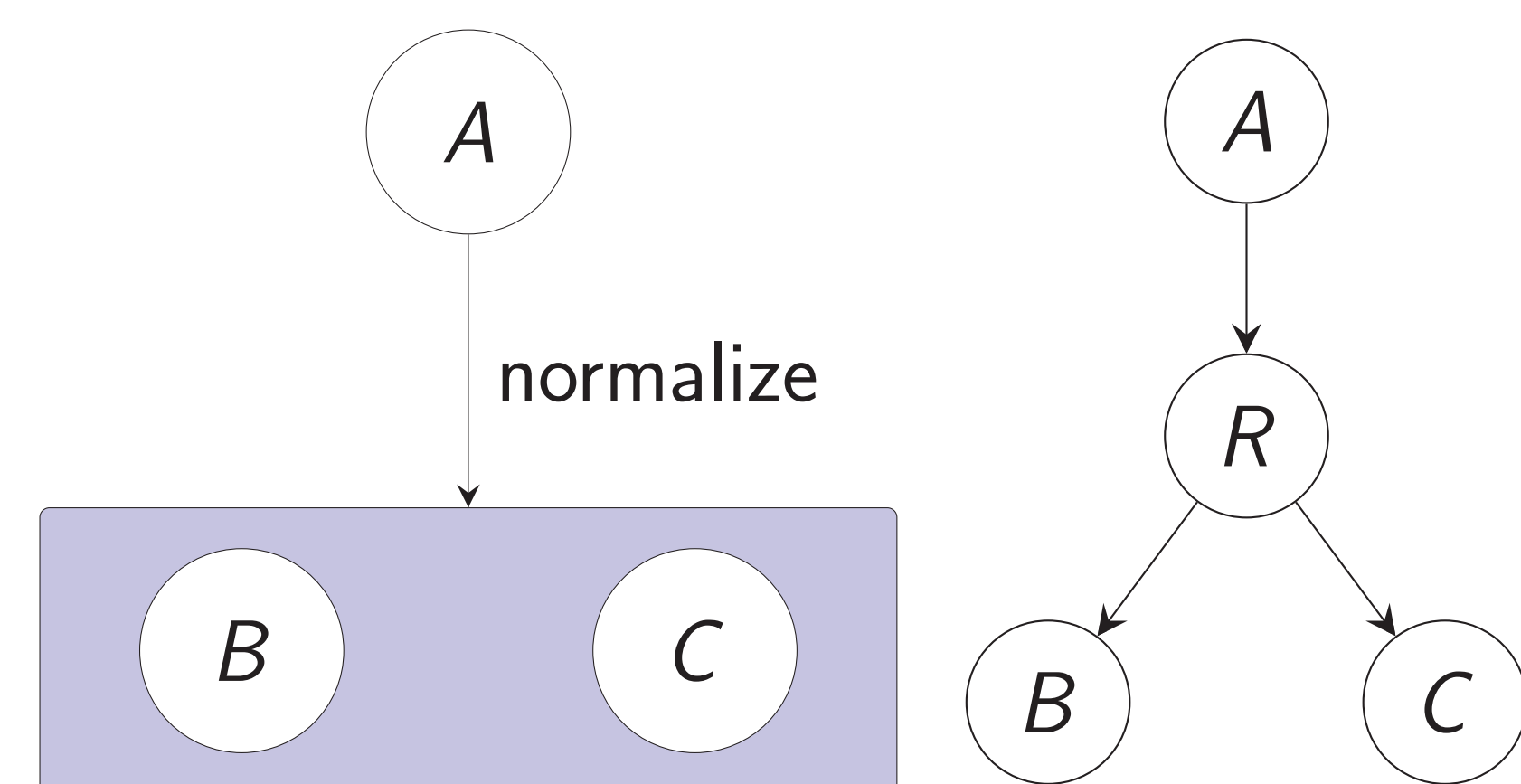
Kruskal form of the binary rule probability tensor:

$$T = \sum_{l=1}^r T^{(l)}, \quad T^{(l)} = u^{(l)} \otimes v^{(l)} \otimes w^{(l)},$$

Eq. 1 becomes the following with complexity  $O(l^3 m + l^2 m r)$ :

$$s_{i,j} = U \cdot \sum_{k=i}^{j-1} ((V^T s_{i,k}) \odot (W^T s_{k+1,j}))$$

## CP decomposition in probabilistic perspective



$$p(A \rightarrow B, C) = \sum_R p(A \rightarrow R) p(R \rightarrow B) p(R \rightarrow C)$$

$$p(A \rightarrow B, C) = \sum_{l=1}^r u^{(l)} \otimes v^{(l)} \otimes w^{(l)}$$

The newly introduced dimension in CP decomposition can be regarded as a latent variable  $R$ .  $\|R\|$  is tensor rank.

Left: normalize  $A \rightarrow B, C$ , total complexity  $O(m^3)$ . Right: normalize  $A \rightarrow R, R \rightarrow B, R \rightarrow C$ , total complexity:  $O(mr)$ .

Row normalization of  $U$  is equivalent to normalize  $p(A \rightarrow R)$ ,

Column normalization of  $V, W$  is equivalent to normalize  $p(R \rightarrow B), p(R \rightarrow C)$  respectively.

## Neural Parameterization

Use distributed representation for each symbol. Use neural networks to calculate grammar rule probabilities.

Use **Softmax** function to ensure the row/column normalization of  $U, V, W$ .

**Without** neural parameterization, performance drops significantly (even underperforms right-branching baseline in WSJ)

Activation functions other than **ReLU** perform much worse.

## Experiments

Comparison on WSJ test set.

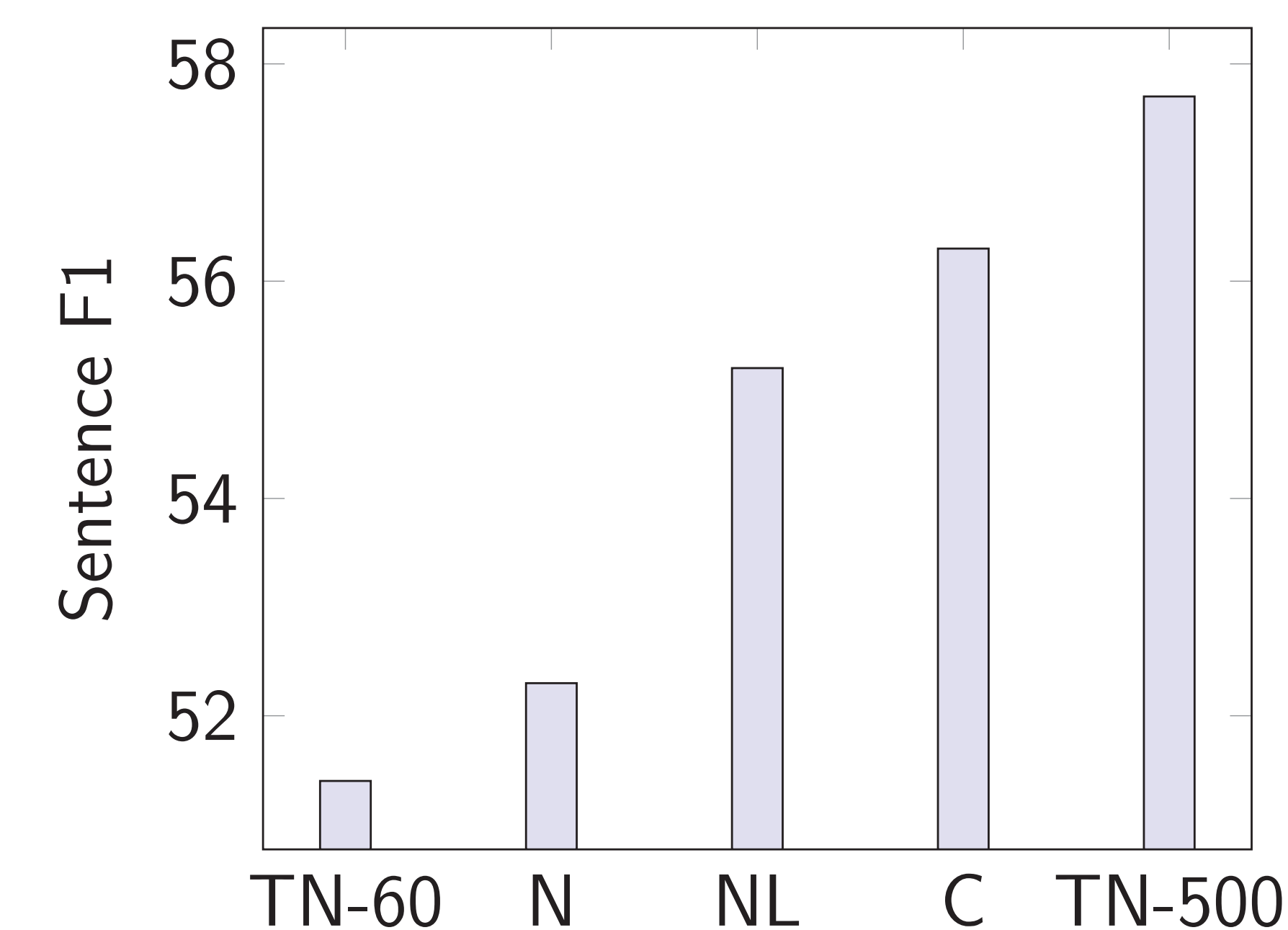


Figure 3. We omit the PCFG suffix. TN: our proposed method. -X: X perterminals. N: Neural Lexicalized. C: Compound

Influence of symbol number

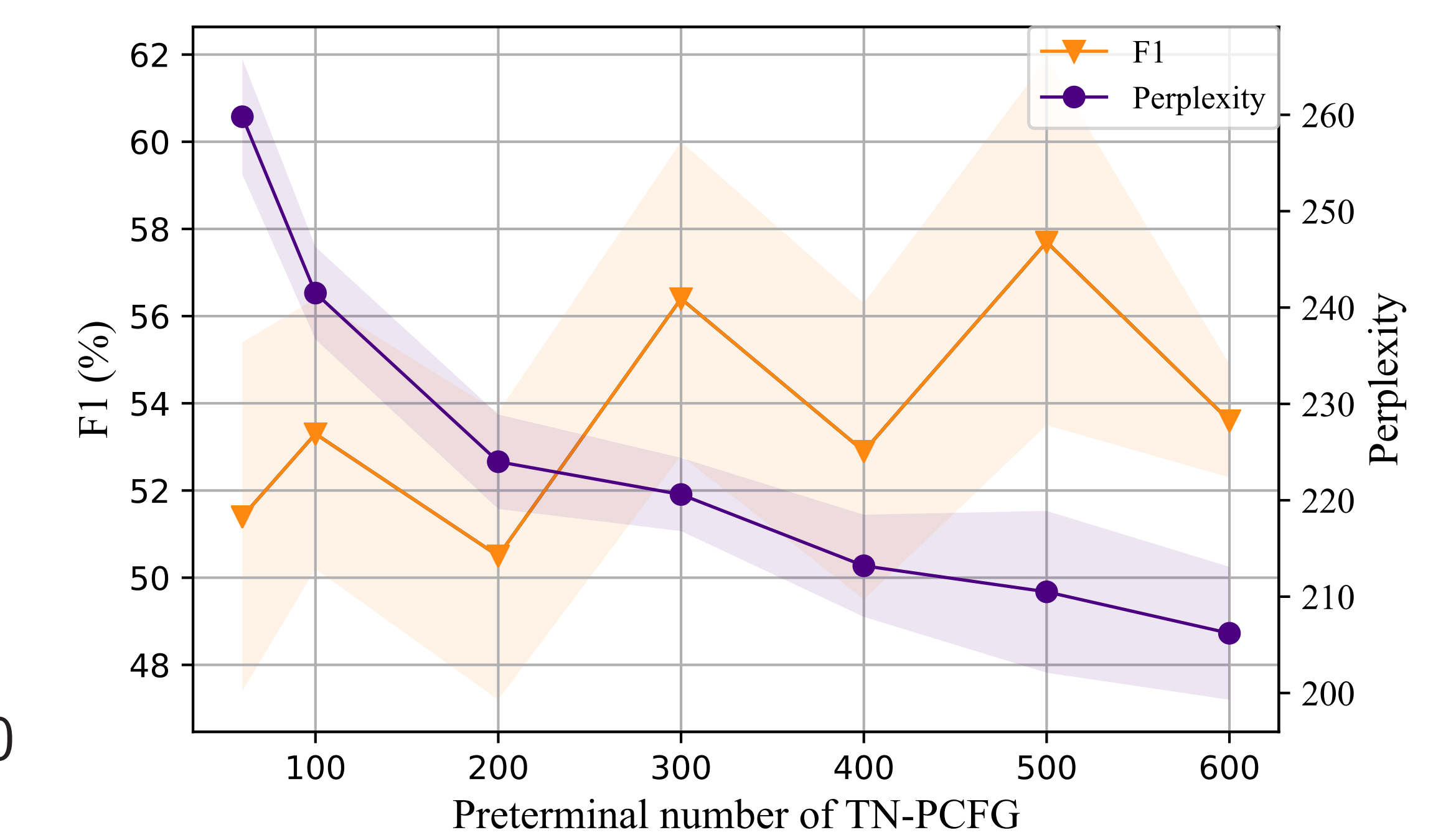


Figure 4. We set the ratio of nonterminals: perterminals as 1:2.

Multilingual experiments

Model	Chinese	Basque	German	French	Hebrew	Hungarian	Korean	Polish	Swedish	Mean
Left Branching <sup>†</sup>	7.2	17.9	10.0	5.7	8.5	13.3	18.5	10.9	8.4	11.2
Right Branching <sup>†</sup>	25.5	15.4	14.7	26.4	30.0	12.7	19.2	34.2	30.4	23.2
Random Trees <sup>†</sup>	15.2	19.5	13.9	16.2	19.7	14.1	22.2	21.4	16.4	17.6
N-PCFG w/ MBR	26.3 $\pm$ 2.5	35.1 $\pm$ 2.0	42.3 $\pm$ 1.6	<b>45.0</b> $\pm$ 2.0	<b>45.7</b> $\pm$ 2.2	43.5 $\pm$ 1.2	28.4 $\pm$ 6.5	43.2 $\pm$ 0.8	17.0 $\pm$ 9.9	36.3
C-PCFG w/ MBR	38.7 $\pm$ 6.6	<b>36.0</b> $\pm$ 1.2	43.5 $\pm$ 1.2	<b>45.0</b> $\pm$ 1.1	45.2 $\pm$ 0.5	<b>44.9</b> $\pm$ 1.5	30.5 $\pm$ 4.2	43.8 $\pm$ 1.3	33.0 $\pm$ 15.4	40.1
TN-PCFG $p = 500$	<b>39.2</b> $\pm$ 5.0	<b>36.0</b> $\pm$ 3.0	<b>47.1</b> $\pm$ 1.7	39.1 $\pm$ 4.1	39.2 $\pm$ 10.7	43.1 $\pm$ 1.1	<b>35.4</b> $\pm$ 2.8	<b>48.6</b> $\pm$ 3.1	<b>40.0</b> $\pm$ 4.8	<b>40.9</b>