

# Text2Story Lusa

## A Dataset for Narrative Analysis in European Portuguese News Articles

Sérgio Nunes, Alípio Jorge, Evelin Amorim, Hugo Sousa, António Leal, Purificação Silvano, Inês Cantante, Ricardo Campos

Narratives have been the subject of extensive research across various scientific fields such as linguistics and computer science. However, the scarcity of freely available datasets remains a significant obstacle.

To address this gap, we developed the Text2Story Lusa datasets, which consist of two independent datasets:

**Text2Story Lusa** and the **Text2Story Lusa Annotated Corpus**.

### Construction

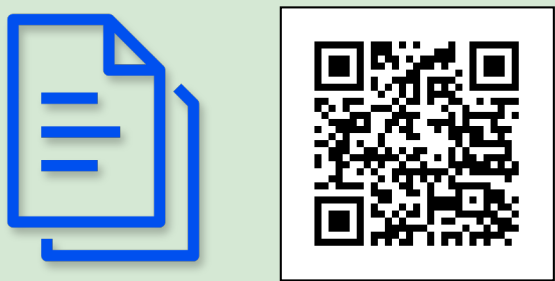
**1st selection criteria:**

- size (50 to 200 words)
- language (PT only)
- no repetitions

**2nd selection criteria:**

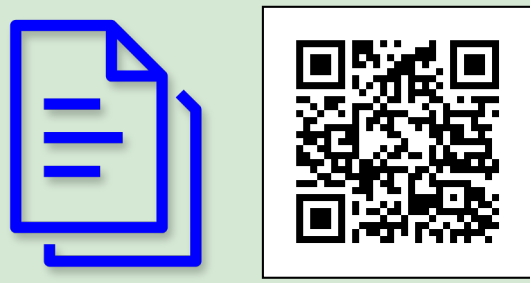
- + set of keywords

**Text2Story Lusa**  
357 news articles



### Annotation

**Text2Story Lusa Subset**  
117 news articles



*Annotation Layers*

Time and Events

Semantic Roles

Spatial

Referential

**Linguists**

**Annotation process**

- individual annotation
- synchronization meetings

**Linguists**

**Incremental development**

- group hands-on sessions
- hand validation

**Annotation Guidelines**

### Annotation Scheme

Time and Events

EVENT

TEMPORAL\_LINK

ASPECTUAL\_LINK

TIME

SUBORDINATION\_LINK

Semantic Roles

SEMANTIC\_ROLE\_LINK

Spatial

SPATIAL\_RELATION

SPATIAL\_LINK

MEASURE

Referential

PARTICIPANT

OBJECTAL\_LINK

### Examples

```
1 {
2   "articles": [
3     {
4       "id": "102",
5       "location": "Lisboa",
6       "publication_time": "2020-12-12",
7       "headline": "Homem socorrido no Rio Tejo está livre de perigo",
8       "content": "Um homem de 68 anos caiu hoje no Rio Tejo, junto ao Cais das Colunas, [...] onde ainda se encontra internado, mas livre de perigo."
9     }
10  ]
11 }
```

Listing 1: Text2Story Lusa news article example.

```
1 T64 Event 930 937 Segundo
2 A91 Class T64 Reporting
3 T53 Participant 938 943 a PSP
4 A286 Lexical_Head T53 Noun
5 A287 Individuation_Domain T53 Individual
6 A288 Participant_Type_Domain T53 Org
7 A289 Involvement T53 I
8 T54 Participant 945 953 o detido
9 A290 Lexical_Head T54 Noun
10 A291 Individuation_Domain T54 Individual
11 A292 Participant_Type_Domain T54 Per
12 A293 Involvement T54 I
13 T21 Time 956 960 hoje
14 A145 Time_Type T21 Date
15 A146 TemporalFunction T21 Publication_Time
16 T22 Event 961 996 presente às Autoridades Judiciárias
17 A147 Class T22 Occurrence
18 A148 Event_Type T22 Transition
19 A149 Pos T22 Verb
20 A150 Tense T22 Present
21 A154 Polarity T22 Pos
22 R86 SRLINK_agent Arg1:T53 Arg2:T64
23 R67 TLINK_after Arg1:T22 Arg2:T64
24 R78 SRLINK_patient Arg1:T22 Arg2:T54
25 R19 TLINK_isIncluded Arg1:T22 Arg2:T21
```

Listing 2: An annotation excerpt, in the BRAT standoff format, from the Text2Story Lusa Dataset. The complete annotated sentence is "Segundo a PSP, o detido é hoje presente às Autoridades Judiciárias." (Translated: "According to the PSP, the detainee is present today at the Judicial Authorities").

### Characterization

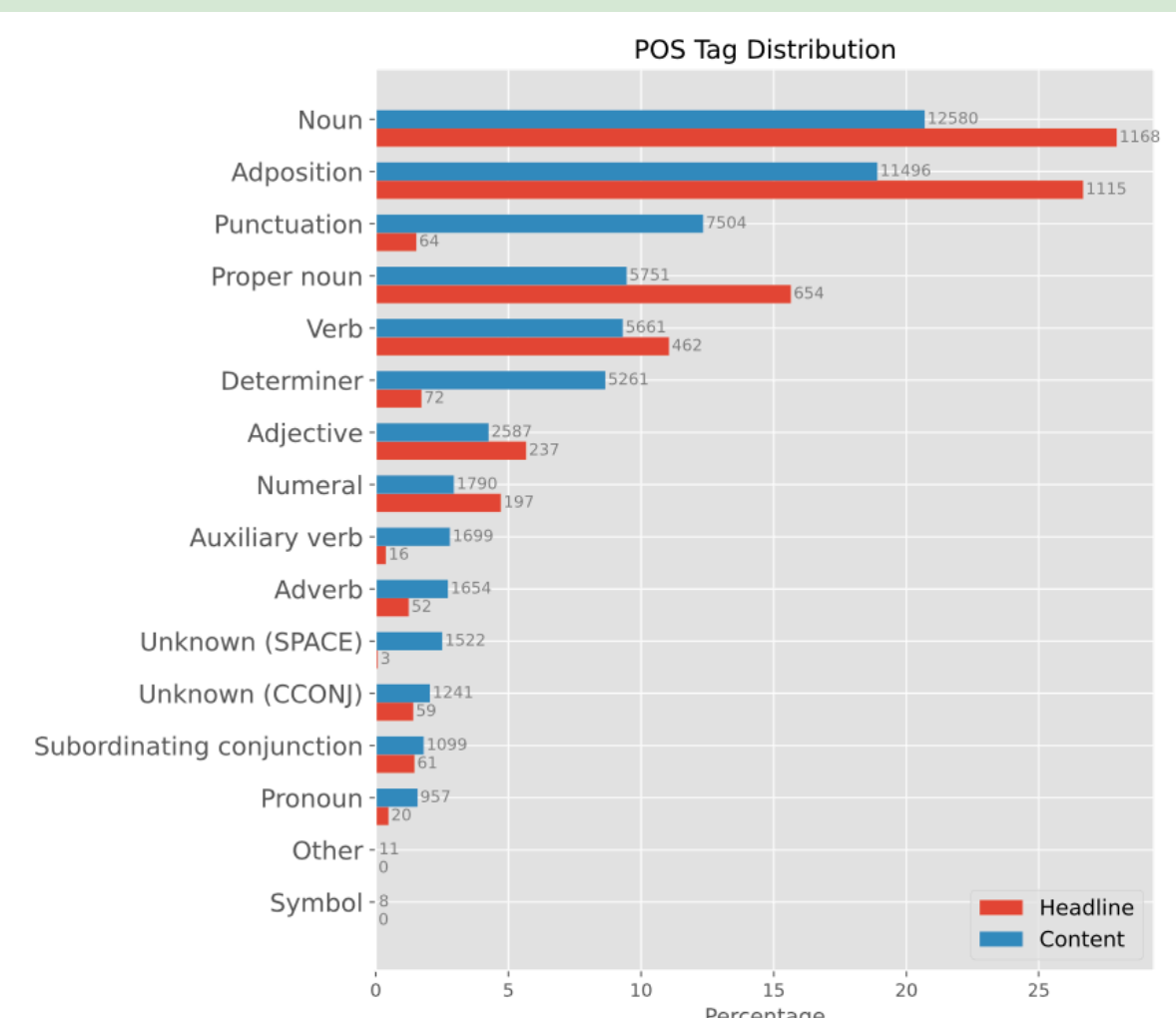


Figure 3: Distribution of POS types in the 'headline' and the 'content' fields.

Narrative Component	Frequency
Participants	3,530
Events	3,027
Spacial Relations	512
Times	438
Measure	9
Semantic Role Links	4,764
Temporal Links	3,522
Objectal Links	2,224
Qualitative Spatial Link	927
Subordination Link	429
Movement Link	125
Aspectual Link	18

Table 1: Frequency of each narrative component in the Text2Story annotated corpus.

### Conclusions

Text2Story Lusa is a linguistic resource comprised of two datasets developed to support research in various fields.

**Text2Story Lusa** includes a manually curated collection of news articles written in European Portuguese, obtained from the Lusa news agency.

**Text2Story Lusa Annotated Corpus** includes manual annotations for a subset of the original articles. The annotation adopts a multilayer semantic scheme comprising entity structures and link structures.