

## Task Overview

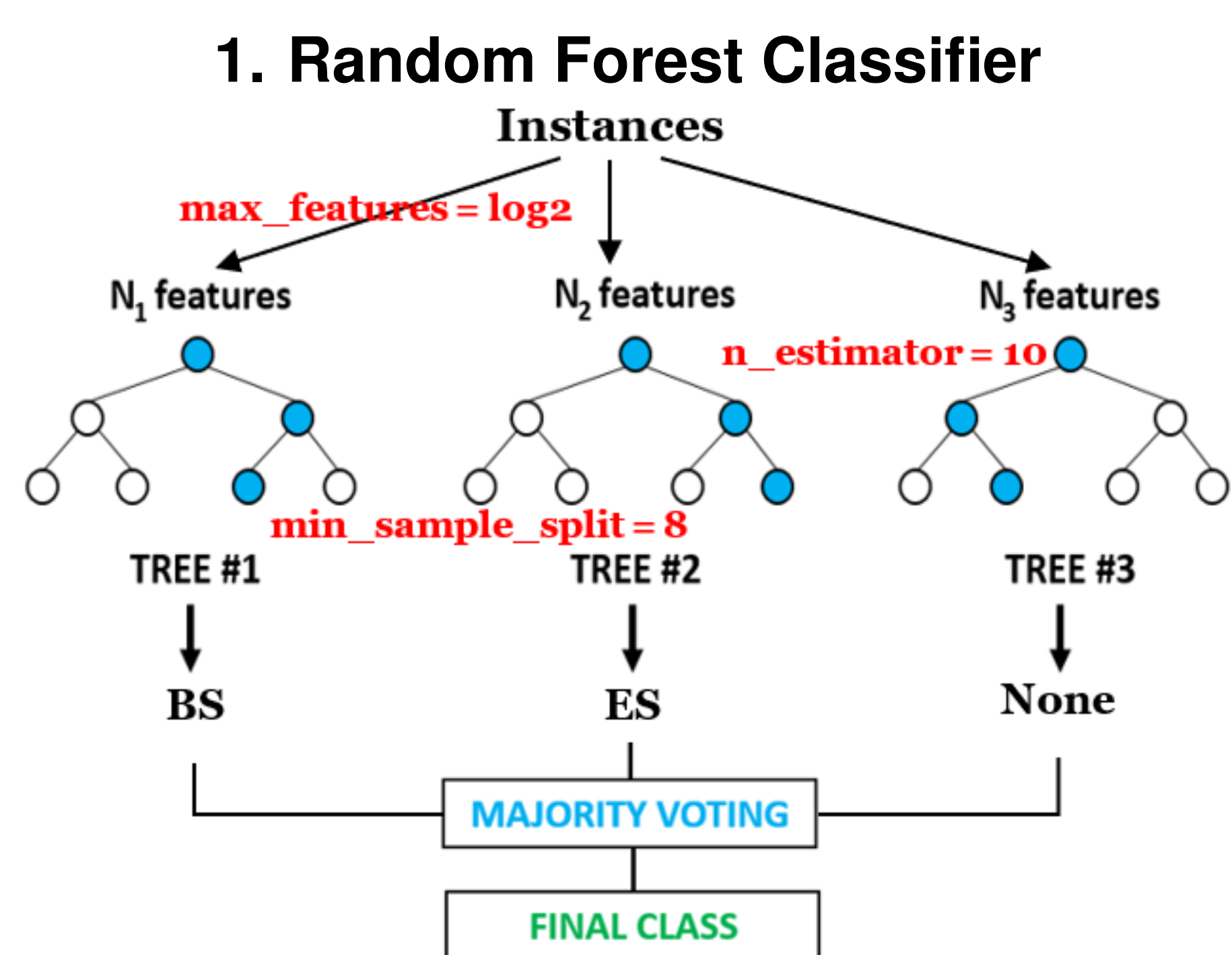
- **Purpose:**
  - to detect sentence boundaries (the begins and ends) from noisy texts of financial prospectuses.
- **Challenges:**
  - documents encoded in non-text formats, such as tables, images, etc.
  - conventional punctuation does not necessarily mark sentence boundary.
- **Our Approach:**
  - Machine Learning: a supervised task of three-class predictions + Feature Engineering: domain-specific and document-specific observation + Rule-based Validation: also known as domain knowledge enhancement.
- **Performance:**
  - English: 0.835 (avg. F-score of BS and ES)
  - French: 0.86 (avg. F-score of BS and ES)

## Feature Engineering

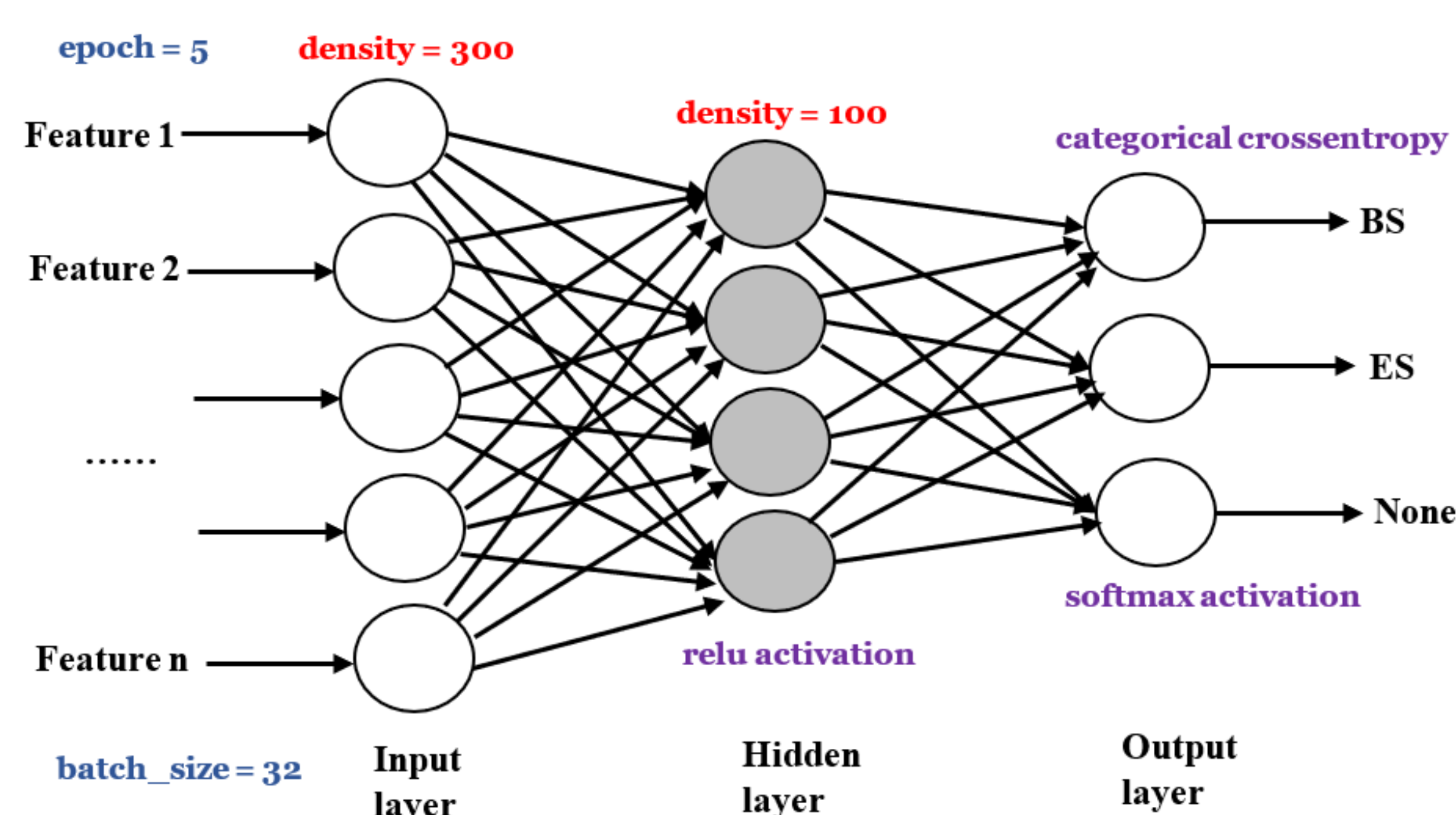
A fused feature set of 24 dimensions, including:

- **Two sets of punctuation:**
  - PUNC\_SET1 = [',']
  - PUNC\_SET2 = ['?', '!', ':', '%', '-', '/', '"', '\', ')', '(', '\*', '□', '<', '≥', '>', '≤', '•', '€', '\$', '£', '™', '©', '®']
- **Initially capitalized words**, e.g. “\_Enter\_END The\_BEGIN sales”;
- **Acronyms**, e.g. “UBS\_BEGIN” “co” or “kiid”;
- **Digital numbers**, ambiguity resolution as in “10.3”;
- **Letters or Roman numbers**, e.g. e.g. A-Z, a-z or I, II, ..., XII;
- **POS tags**, structural cues.

## Classifiers



## 2. Neural Network



## Domain Knowledge Enhancement

### Algorithm 1 Keyword Extraction

**Input:** dataset  
**Output:** keyword\_dict

```

1: for i in len(dataset) do
2:   curword = dataset[i].
3:   nxtword = dataset[i+1].
4:   if "END" in curword then
5:     if "BEGIN" in nxtword then
6:       continue
7:     else
8:       add keyword to keyword_dict.
9:       update frequency in keyword_dict.
10:    end if
11:  end if
12: end for
13: refine keyword_dict with length threshold.
14: refine keyword_dict with frequency threshold.
15: return keyword_dict

```

### Algorithm 2 Rule-based validation

**Input:** dataset, raw\_pred, keyword\_dict  
**Output:** updated\_prediction

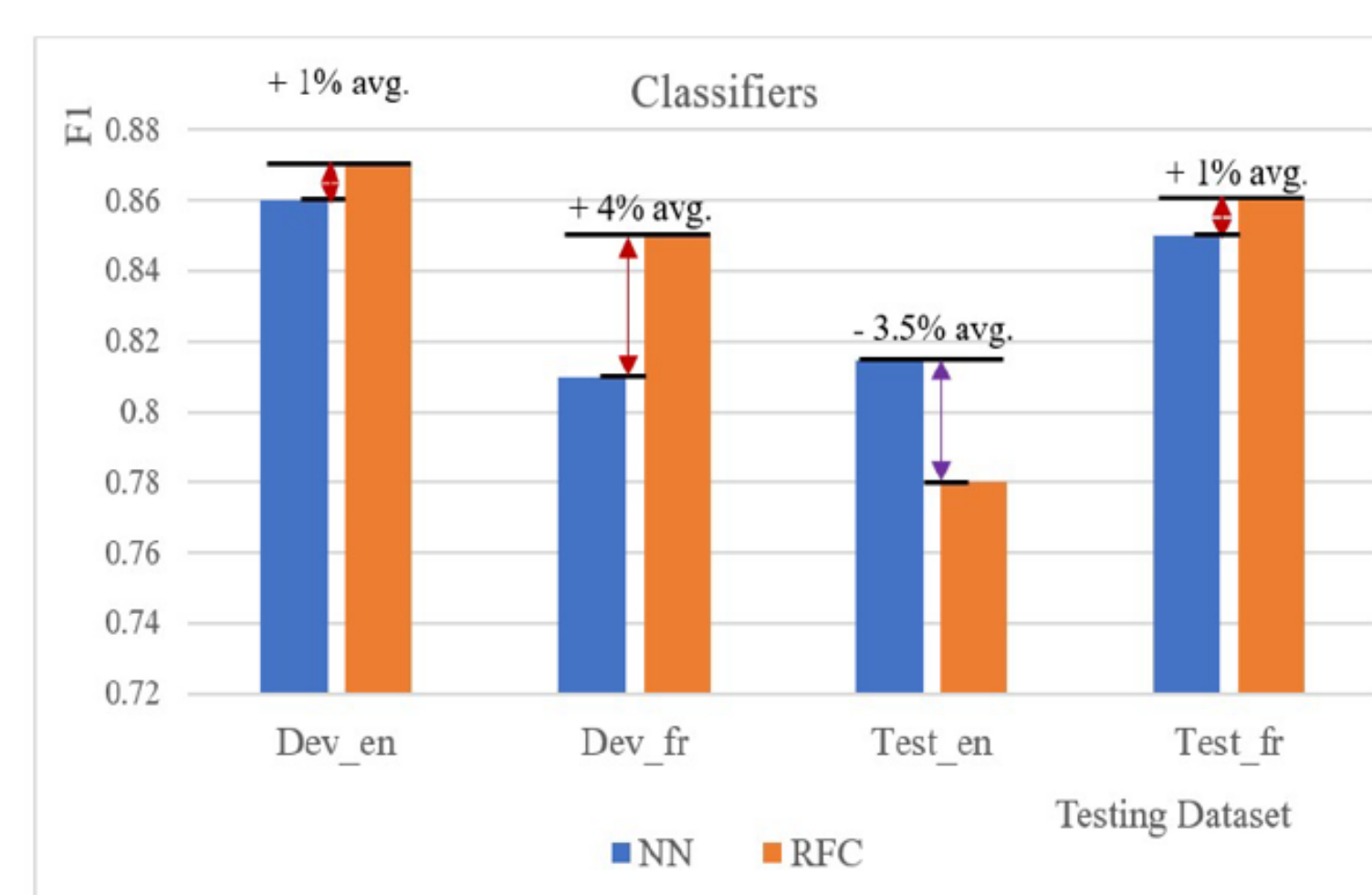
```

1: for keyword in keyword_dict.key do
2:   for i in len(dataset)-len(keyword) do
3:     if word sequence match keyword then
4:       update raw_pred with NO_BOUNDARY
5:     end if
6:   end for
7: end for
8: return raw_pred

```

- The two steps can correct false positive predictions caused by non-text sections, such as the titles, dates, tables, figures, etc., which fail to fall into the traditional category of sentence boundary.

## Results and Comparison



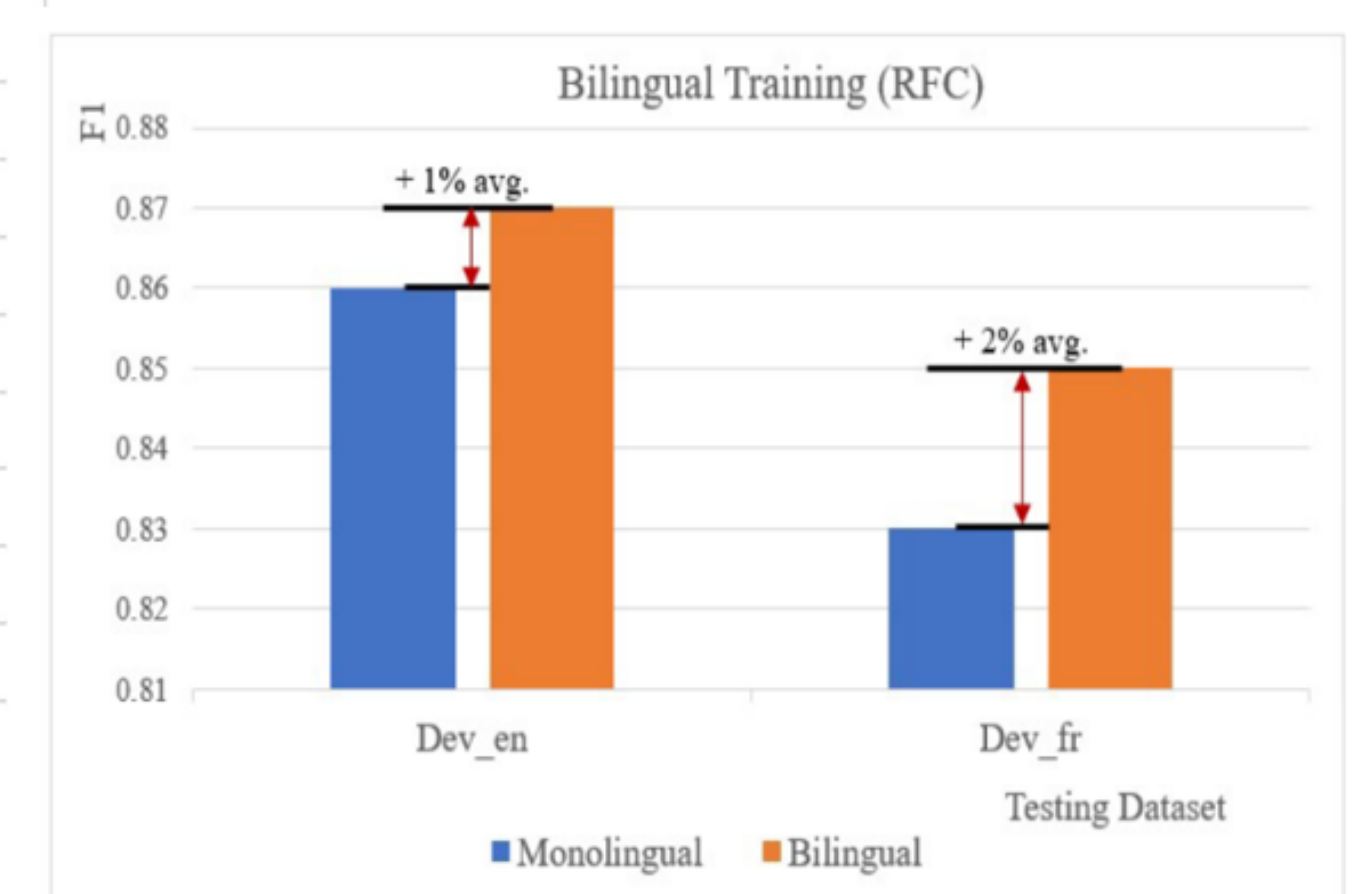
- RFC outperforms NN in 3/4 tasks
- NN outperforms RFC in the English Test set

Features\F1	BS <sup>a</sup>	ES <sup>b</sup>	Mean	$\delta(\%)^c$
Punc1	0.70	0.82	0.76	baseline
+Punc2	0.71	0.83	0.77	1 $\uparrow$
+Cap	0.72	0.86	0.79	2 $\uparrow$
+Acro	0.73	0.87	0.80	1 $\uparrow$
+Dig	0.73	0.89	0.81	1 $\uparrow$
+Lett	0.74	0.90	0.82	1 $\uparrow$
+POS	0.78	0.92	0.85	3 $\uparrow$
+Enter/All	0.80	0.92	0.86	1 $\uparrow$

<sup>a</sup> Beginning boundaries  
<sup>b</sup> Ending boundaries  
<sup>c</sup> F1 improvement in percentage

Table 1: Performance of feature mining in the English Dev set with RFC

- POS vector is the most useful feature.
- Knowledge enhancement is significantly helpful.



- Bilingual training consistently outperforms monolingual training with 1-2 % F gain.

	F1	NO <sup>a</sup>	YES <sup>b</sup>	$\delta(\%)^c$
Dev BS	0.83	0.83	0	
Dev ES	0.86	0.87	1 $\uparrow$	
Dev Mean	0.845	0.85	0.5 $\uparrow$	
Test BS	0.81	0.84	3 $\uparrow$	
Test ES	0.88	0.88	0	
Test Mean	0.845	0.86	1.5 $\uparrow$	

<sup>a</sup> Without keyword validation  
<sup>b</sup> With keyword validation

Table 2: Performance of RFC in the French Dev and Test sets in terms of keyword validation

## Adaptation to Test Set

Test Set	F1-Mean	Rank
Dev_en_rfc1	0.875	—
Test_en_rfc1	0.78*	16 $\diamond$
Test_en_rfc1_adapted	0.835*	10*
Dev_fr_rfc1	0.85	—
Test_fr_rfc1	0.86*	8 $\diamond$
Test_fr_rfc1_adapted	0.86*	8*

\* Result without adaptation to the test set  
 $\diamond$  Rank without adaptation to the test set  
\* Result with adaptation to the test set  
\* Rank with adaptation to the test set

Table 3: Performance of RFC w.r.t. knowledge adaptation

## Conclusions

- NN does not show advantage over statistical classifiers, but it outperforms RFC in terms of unknown features;
- Syntactic information may be helpful for sentence detection;
- Our system is highly effective with minimal training data;
- Future studies will further verify the effectiveness of this feature design and knowledge enrichment.