

# Nested Named Entity Recognition

2020/04/21

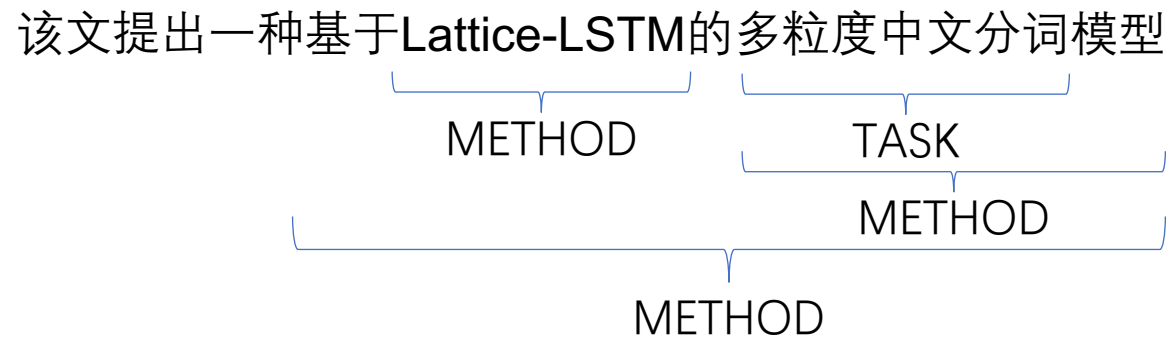
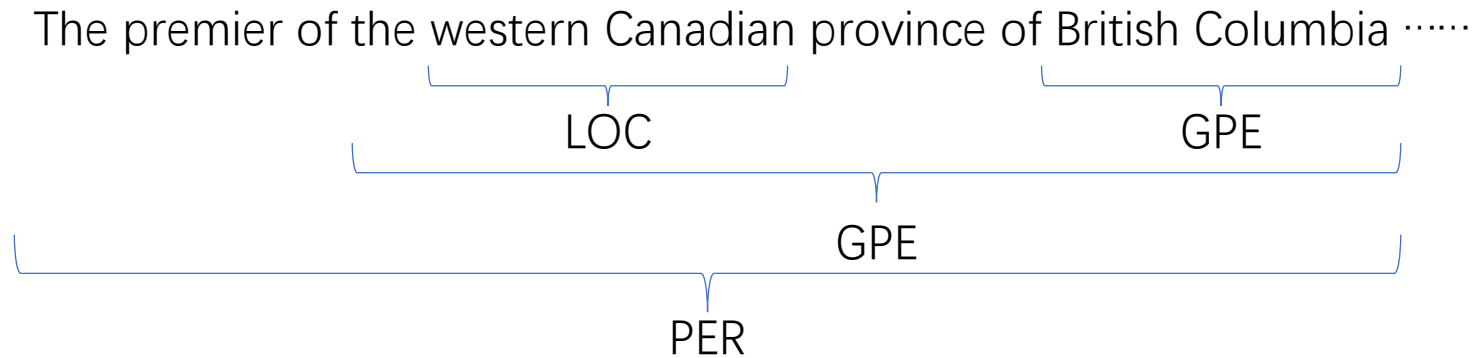
CQU 1701 Online Meeting

# Contents

- Problem Definition
- Datasets
- Sequence Labeling Methods
- Sub-Sequence Classifier Methods
- Summary

# Problem Definition

Embedded names which are included in other entities



# Problem Definition

- Give a sequence  $X = \{x_1, x_2, \dots, x_n\}$
- Predict a sequence of label  $Y = \{y_1, y_2, \dots, y_n\}$
- $y_n$  contains a list of labels not a single label:  
$$y_n = \{y_n^1, y_n^2, \dots, y_n^m\}, \quad m = \textit{nested\_layer}$$

A multi-label classification problem

# Nested Type

- A nested entity contains more than 1 entity:
  - the western Canadian province of British Columbia
- A nested entity extended from 1 flat entity + key words:
  - 多粒度中文分词模型

# Datasets

- GENIA
- ACE2005
- NNE

|         | Text Type       | Doc Number | Mentions | Entity Types | Nested Level |
|---------|-----------------|------------|----------|--------------|--------------|
| GENIA   | Biomedical Text | 2,000      | 92,681   | 36           | 4            |
| ACE2005 | News            | 464        | 30,966   | 7            | 6            |
| NNE     | News            | 2,312      | 279,795  | 114          | 6            |

# Raw & Processed format

Mice<sup>T5</sup> transgenic for the human T cell leukemia virus<sup>T7</sup> (HTLV-I<sup>T8</sup>) Tax<sup>T9</sup> gene<sup>T6</sup> develop fibroblastic tumors<sup>T10</sup> that express NF-kappa B-inducible early genes<sup>T11</sup>.

<sentence id="S2"><term id="T5" sem="Multicellular\_organism">Mice</term> transgenic for the <term id="T6" sem="DNA\_domain\_or\_region"><term id="T7" sem="Virus">human T cell leukemia virus</term> (<term id="T8" sem="Virus">HTLV-I</term>) <term id="T9" sem="Protein\_molecule">Tax</term> gene</term> develop <term id="T10" sem="Tissue">fibroblastic tumors</term> that express <term id="T11" sem="DNA\_family\_or\_group">NF-kappa B-inducible early genes</term>.</sentence>

```
These 0 0 0 0
data 0 0 0 0
indicate 0 0 0 0
that 0 0 0 0
IL B-protein 0 0 0
- I-protein 0 0 0
4 I-protein 0 0 0
suppresses 0 0 0 0
the 0 0 0 0
induction 0 0 0 0
of 0 0 0 0
transcription B-protein 0 0 0
factors I-protein 0 0 0
in 0 0 0 0
human 0 B-cell_type 0 0
activated 0 I-cell_type 0 0
monocytes B-cell_type I-cell_type 0 0
. 0 0 0 0
```

# Sequence Labeling Methods

- Combined Label
- Neural Layered Model



# BIOES Encoding(Combined Label)

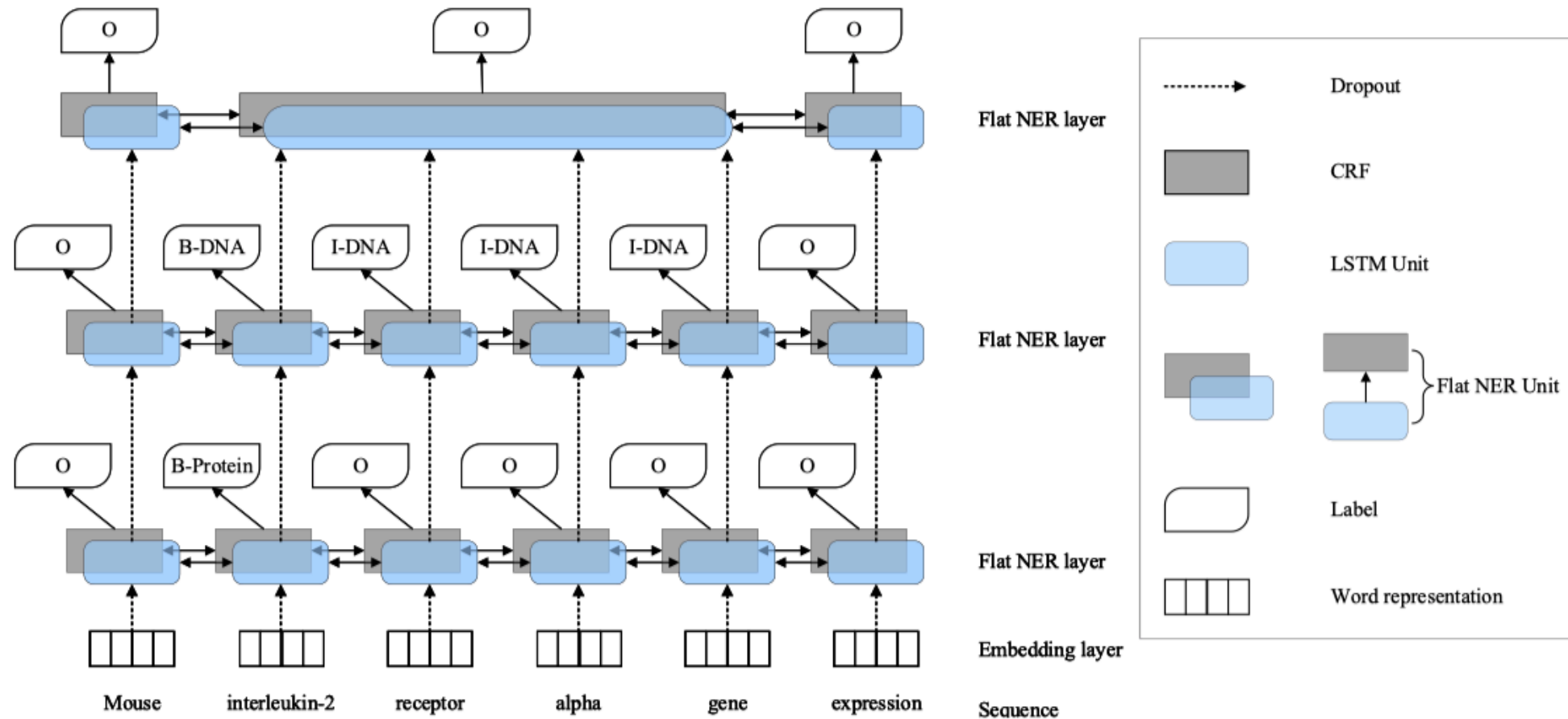
- Flat NER:

- Begin, Inside, Outside
  - Type1, Type2, Type3
- } B-Type1, I-Type1,.....

- Nested NER:

|          |               |
|----------|---------------|
| in       | O             |
| the      | B-ORG         |
| US       | I-ORG   U-GPE |
| Federal  | I-ORG         |
| District | I-ORG   U-GPE |
| Court    | I-ORG         |
| of       | I-ORG         |
| New      | I-ORG   B-GPE |
| Mexico   | L-ORG   L-GPE |
| .        | O             |

# Neural Layered Model



A Neural Layered Model for Nested Named Entity Recognition, NAACL-HLT 2018

# Problems

- Combined Label:
  - The number of Labels grows exponentially
  - Label distribution is too sparse
- Layered Model:
  - Error propagation
  - Can not train model parallelly

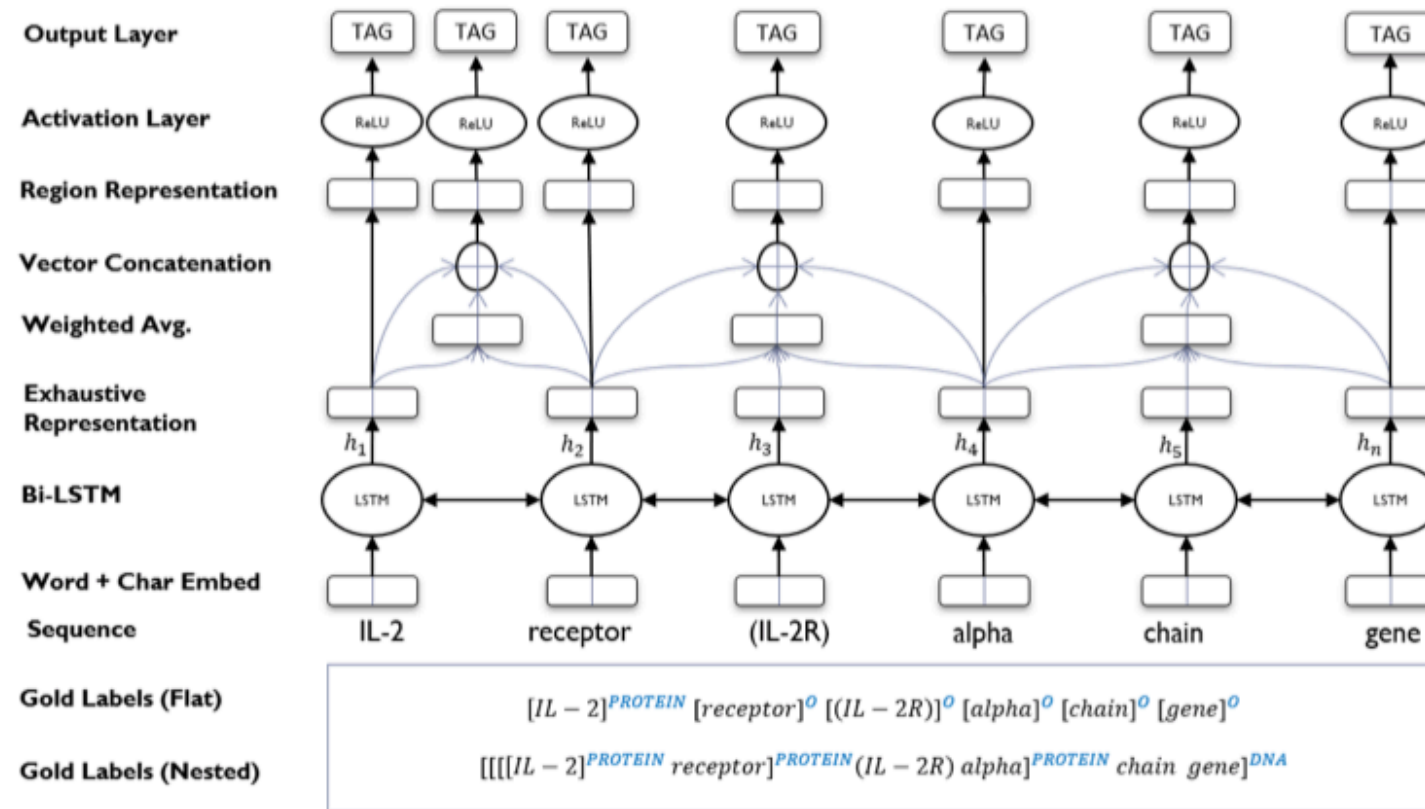
# Sub-Sequence Classifier Methods

- Deep Exhaustive Model
- Boundary-aware Neural Model
- Connection-aware Model

# Simple Idea

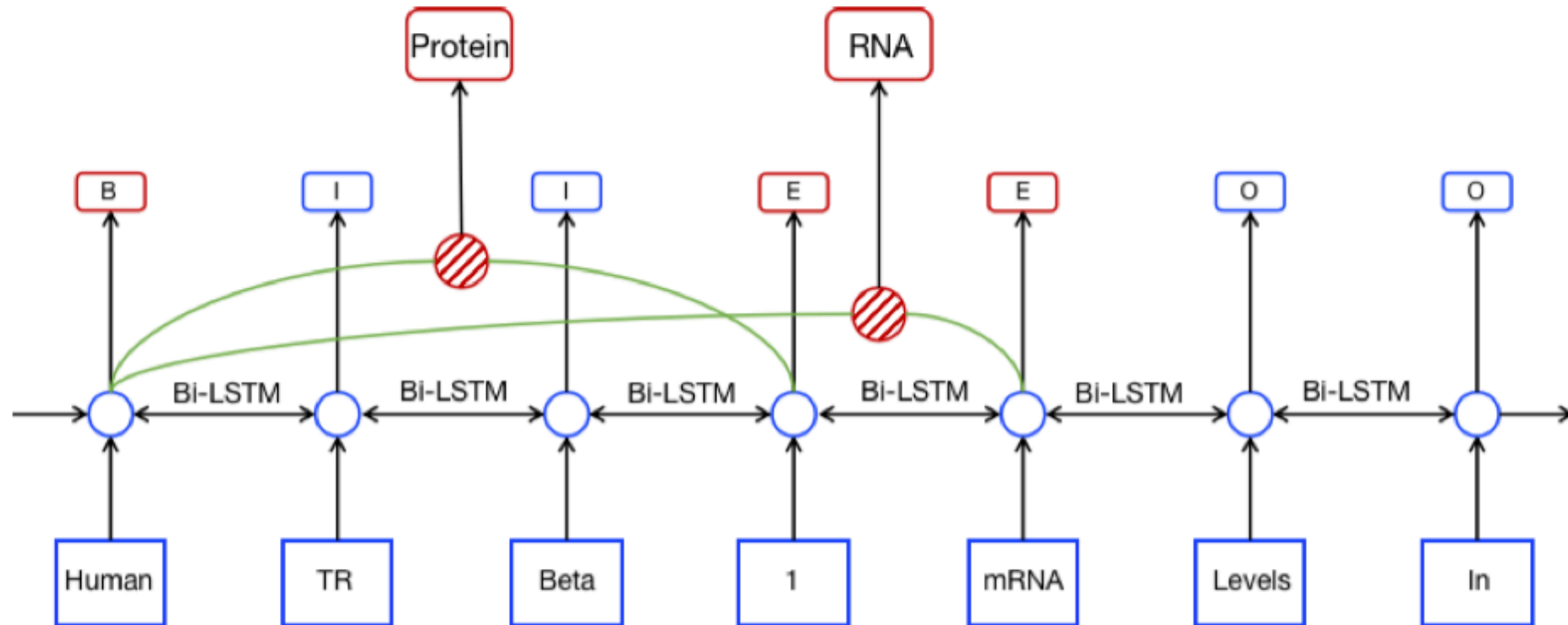
- Give a sequence  $X = \{x_1, x_2, \dots, x_n\}$
- Enumerate all sub-sequences of  $X$ ,  
$$S = \{s_1, s_1s_2, s_1s_2s_3, \dots, s_2, s_2s_3, \dots, s_n\}$$
- Train a Classifier  $C$ , predict the label of each sub-sequence:  
$$Y = \{y_1, y_2, \dots, y_m\}, \quad m = \text{len}(S)$$

# Deep Exhaustive Model



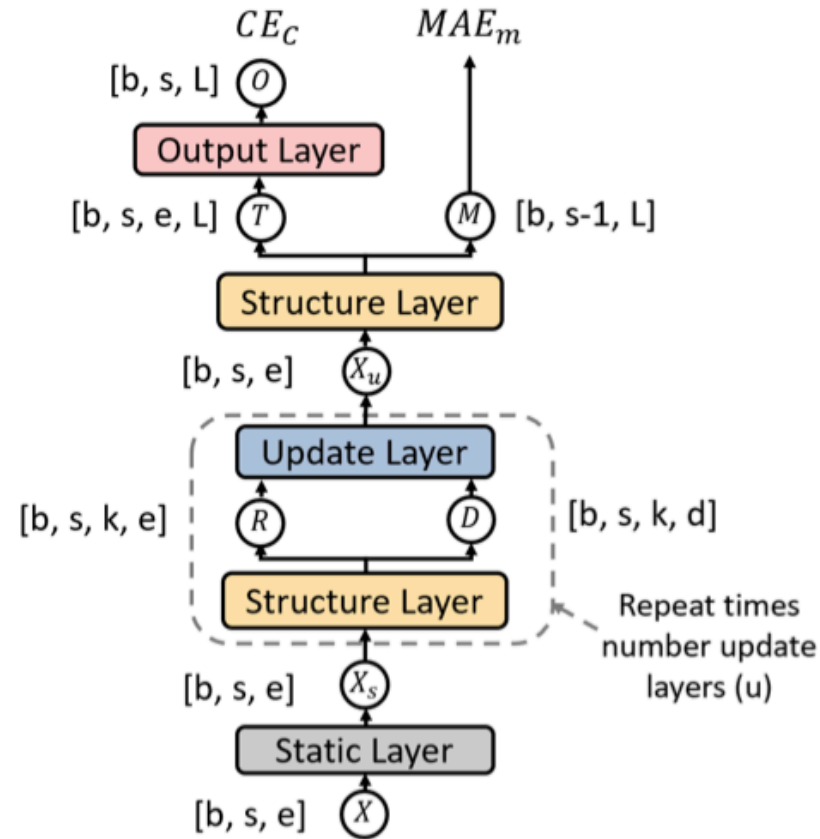
Deep Exhaustive Model for Nested Named Entity Recognition, EMNLP 2018

# Boundary-aware Model



A Boundary-aware Neural Model for Nested Named Entity Recognition, EMNLP-IJCNLP 2019

# Connection-aware Model





# Problems

- Negative Samples
- Length of sub-sequence
- High complexity

# Results on GENIA

|                      | P    | R    | F     |
|----------------------|------|------|-------|
| Neural Layered       | 78.5 | 71.3 | 74.7  |
| Deep Exhaustive      | 73.3 | 68.3 | 70.7  |
| Linearization(Flair) | /    | /    | 78.31 |
| Boundary-aware       | 75.9 | 73.6 | 74.7  |

# Results on ACE2005

|                      | P    | R    | F     |
|----------------------|------|------|-------|
| Neural Layered       | 74.2 | 70.3 | 72.2  |
| Linearization(Flair) | /    | /    | 84.33 |
| Merge and Label      | 82.7 | 82.1 | 82.4  |

# Summary

- The trend is use sub-sequence classifier
- Most work aims to reduce the number of negative samples and complexity
- Layered model still work
- Combined label performs well due to BERT/ELMo.