



**LREC-COLING 2024**

Lingotto Conference Centre - Torino (Italia)

20-25 May, 2024

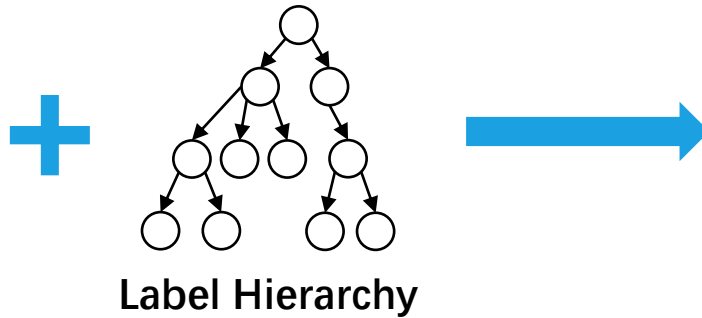
# **NER-guided Comprehensive Hierarchy-aware Prompt Tuning for Hierarchical Text Classification**

Fuhan Cai, Duo Liu, Zhongqiang Zhang, Ge Liu, Xiaozhe Yang, Xiangzhong Fang  
Shanghai Jiao Tong University, East China Normal University

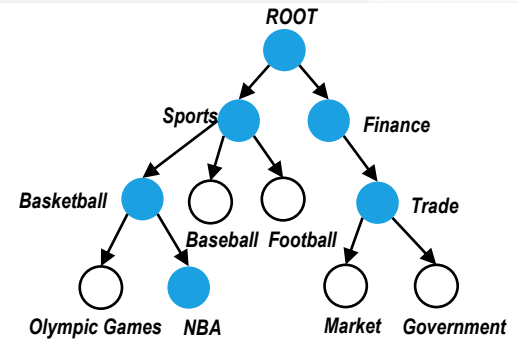
## Hierarchical Text Classification (HTC)

### Input&Output:

**Text:** James Harden traded to Clippers.



Label Hierarchy



**Labels:** Sports, Basketball, NBA, Finance, Trade

### Dataset:

- Web of Science (WOS): The input text is the abstracts of articles in the Web of Science citation database, fixed 2-layer, single-path;
- New York Times (NYT): The input text is a news article with up to 8 layers, a multi-path dataset;
- RCV1-V2: The input is Reuters news with up to 4 layers, a multi-path dataset.

Dataset	$ Y $	Depth	$Avg( y_i )$	Train	Dev	Test
WOS	141	2	2.0	30,070	7,518	9,397
NYT	166	8	7.6	23,345	5,834	7,292
RCV1-V2	103	4	3.24	20,833	2,316	781,265

Table 1: Data Statistics.  $|Y|$  is the number of classes.  $Avg(|y_i|)$  is the average number of classes per sample. Depth is the maximum level of label hierarchy.

## Prompt Tuning Framework for HTC

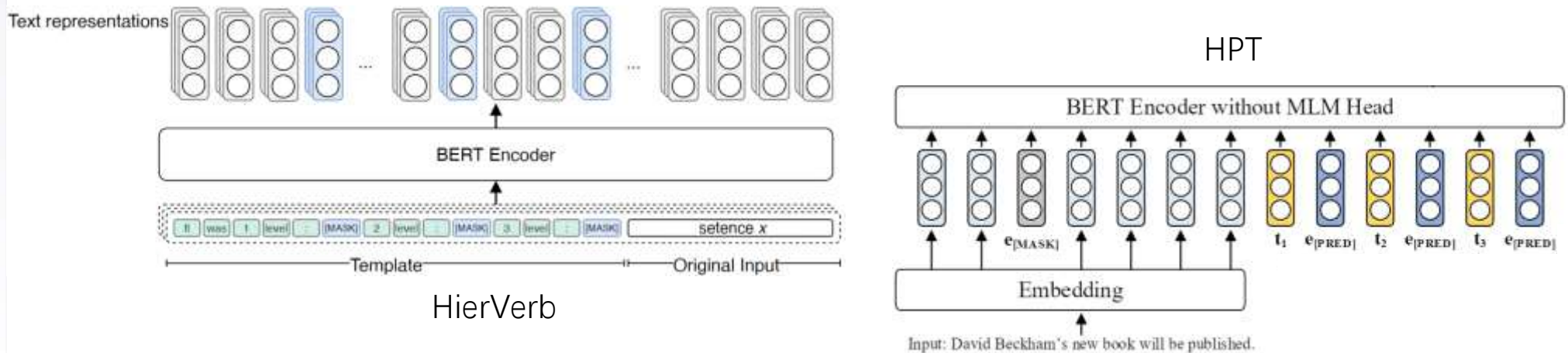
### Hard Prompt:

[CLS] sentence x [SEP] The text is about [MASK] [SEP]

[CLS] It was 1 level: [MASK] 2 level: [MASK] [SEP] sentence x [SEP]

### Soft Prompt:

[CLS] sentence x [SEP] [V1] [MASK] [V2] [MASK] ... [V8] [MASK] [SEP]



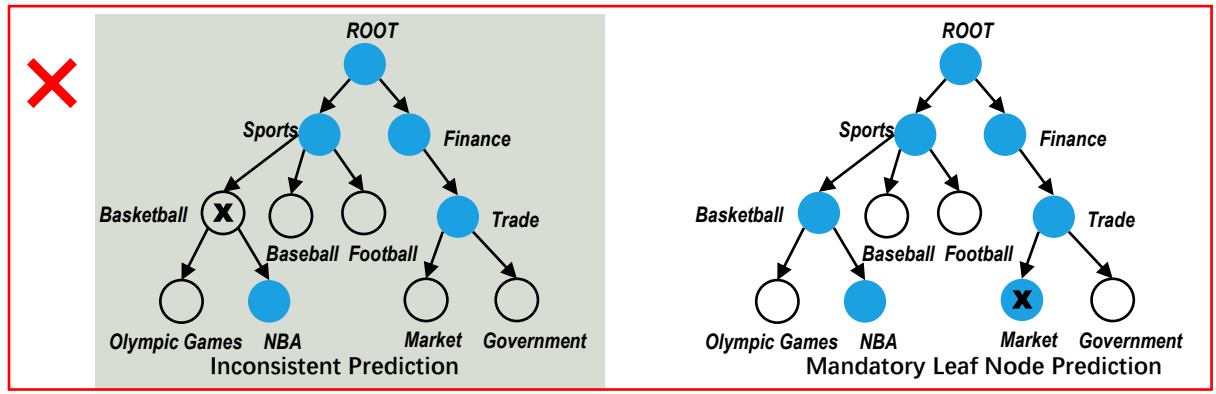
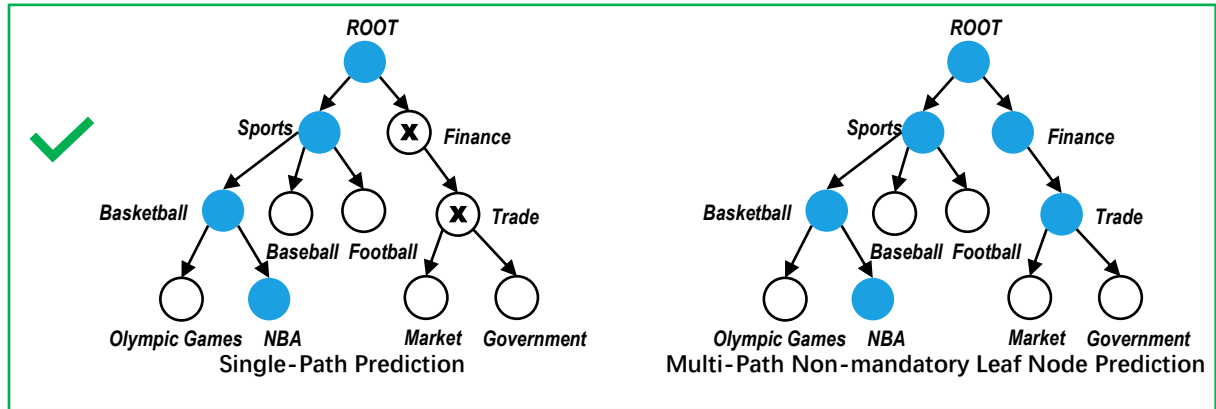
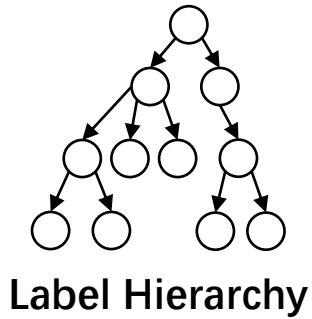


# Challenges

How to solve the hierarchical inconsistency problem in a simple and effective way



James Harden traded to Clippers.

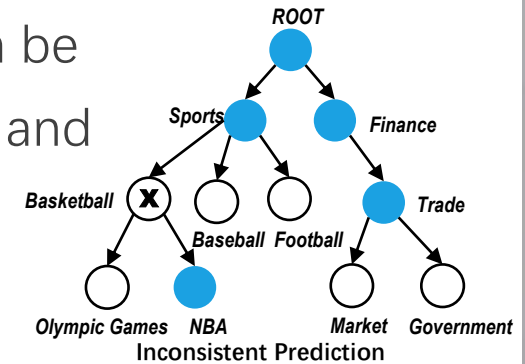


## Motivation

1. The solution to the problem of hierarchical consistency can be approached in terms of local relationships between parent and child nodes ;
2. The NER task can be analogized to the HTC task ;
3. NER tasks are more dependent on the relationship between neighboring tokens;



Therefore, it is possible to migrate the methods of NER tasks to HTC tasks to strengthen the dependency between parent and child nodes.

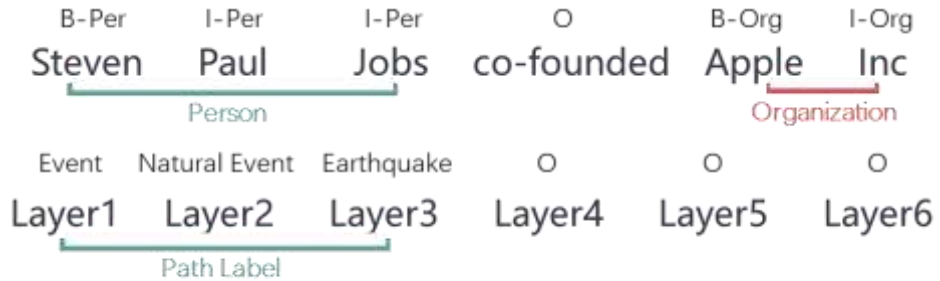




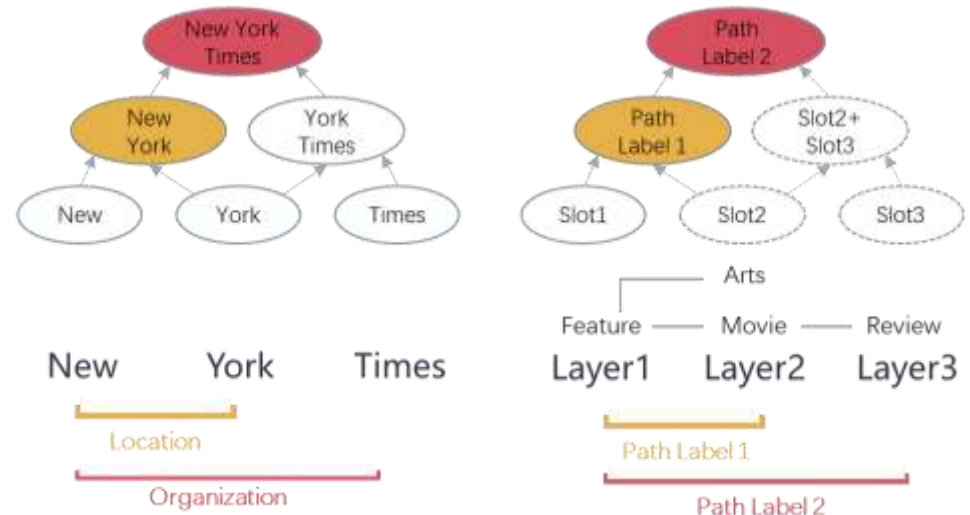
# Motivation

## Support 1:

[CLS] sentence x [SEP] [V1] [MASK] [V2] [MASK] ... [V8] [MASK] [SEP]



(a) Sequence labeling for flat NER and single-path HTC task



(b) Span-based classification for nested NER and multi-path HTC task

## Motivation

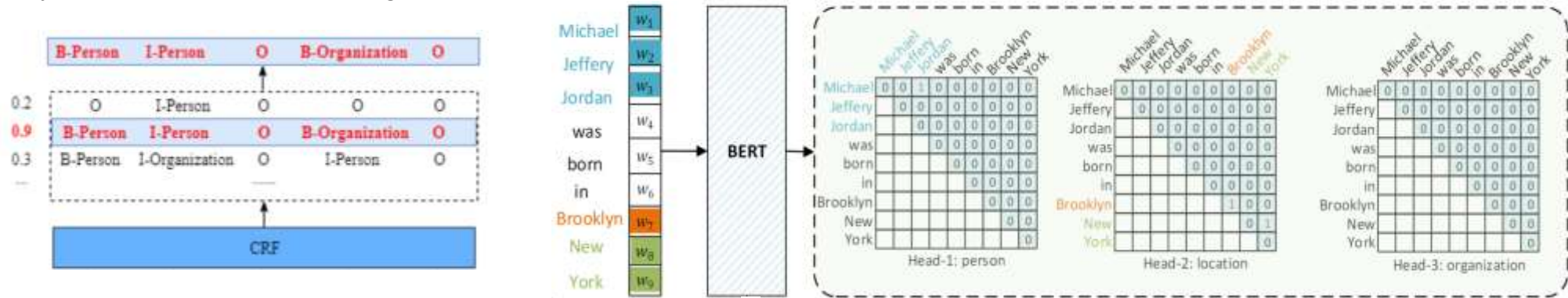
### Support 2:

Louis David watched a game of the Chinese men's basketball team at the Yanyuan Garden of Peking University.

[B-PER, E-PER, O, O, O, O, O, B-ORG, I-ORG, I-ORG, E-ORG, O, O, B-LOC, E-LOC, O, B-ORG, I-ORG, I-ORG, E-ORG]

- Entity Index: B-Beginning, I-Intermediate, E-Ending, O-Others
- Entity Label: PER-Person, ORG-Organization, LOC-Location

Michael Jeffery Jordan was born in Brooklyn New York





## Contributions

1. By considering the characteristics of hierarchical dependencies, we are **the first to adopt NER methods for modeling two prompt-based HTC tasks**, providing a novel perspective for hierarchical-related work.
2. We employ Global Pointer and CRF designed for nested and flat entities to model both multipath and single-path HTC problems, **improving path consistency** in the results through a simple and effective way.
3. We evaluate our method on three popular datasets: Web-of-Science (WOS), NYTimes (NYT), and RCV1-V2. Extensive experiments demonstrate that our method achieves **significant improvements**.



## Overall Framework

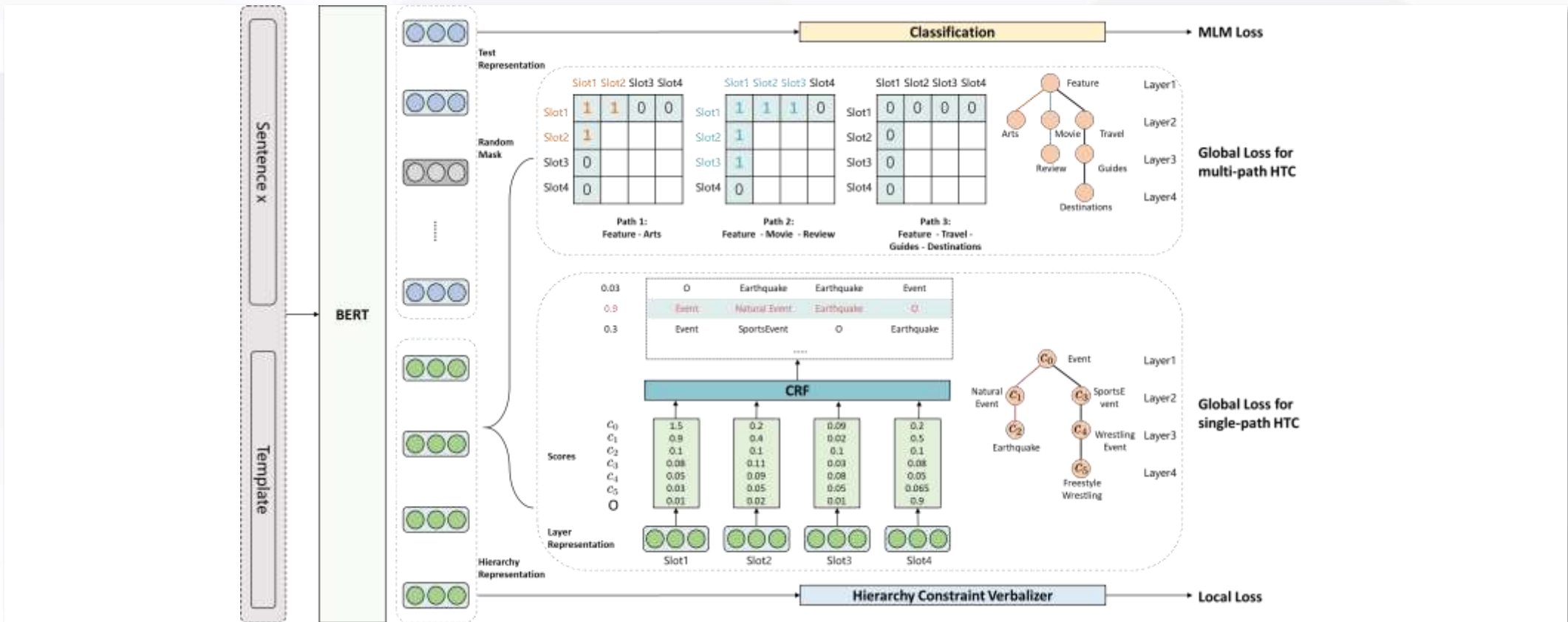
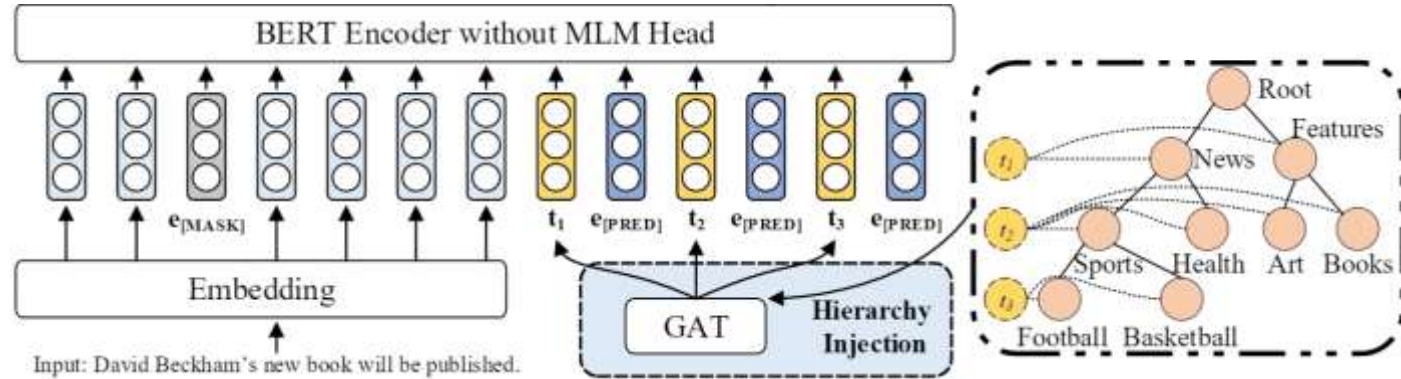


Figure 2: Model structure used in our method. We consider single-path HTC and multi-path HTC as flat and nested NER, respectively. For the multi-path task, we employ span-based Global Pointer method. In the case of single-path HTC, it is treated as a sequence labeling problem. For global-level information, we follow the approach in HPT.

## Global Hierarchy-aware Structure



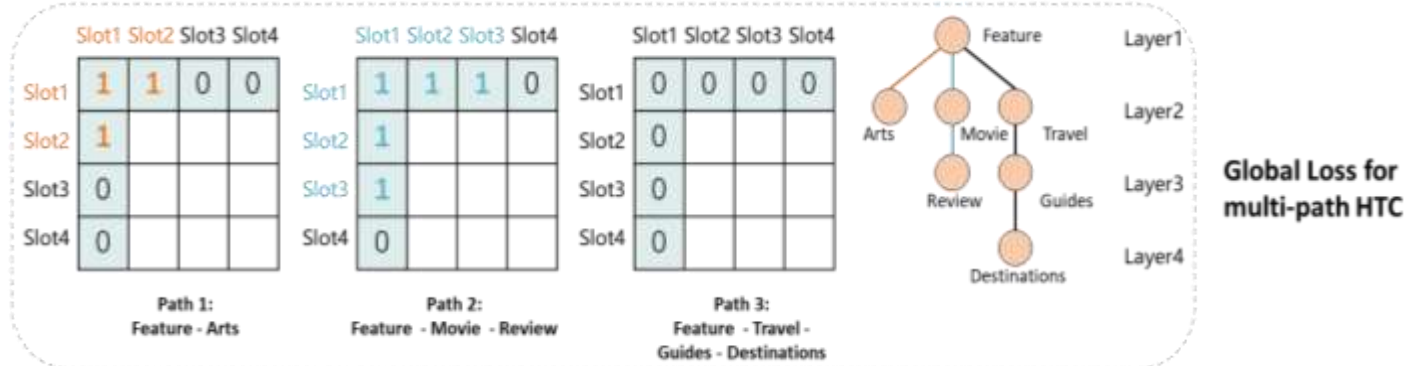
$$\mathbf{Emb} = [\mathbf{X}; \mathbf{T}] = [\mathbf{x}_1, \dots, \mathbf{x}_n, \mathbf{t}_1, \mathbf{e}_{\text{hie}}^1, \dots, \mathbf{t}_L, \mathbf{e}_{\text{hie}}^L]$$

$$\mathbf{G}_v^{(l+1)} = \text{ReLU} \left( \sum_{u \in \mathcal{N}(v) \cup \{v\}} \frac{1}{c_v} \mathbf{W}^{(l)} \mathbf{G}_u^{(l)} \right)$$

$$\mathbf{t}'_i = \mathbf{t}_i + \mathbf{G}_{t_i}^K$$

$$\mathbf{H} = \text{BERT}(\mathbf{Emb}) = [\mathbf{h}_1, \dots, \mathbf{h}_n, \mathbf{h}_{t_1}, \mathbf{h}_{\text{hie}}^1, \dots, \mathbf{h}_{t_L}, \mathbf{h}_{\text{hie}}^L]$$

# Local Hierarchy-aware Structure for Multi-path HTC



$$\mathbf{H} = \text{BERT}(\mathbf{Emb}) = [\mathbf{h}_1, \dots, \mathbf{h}_n, \mathbf{h}_{t_1}, \mathbf{h}_{hie}^1, \dots, \mathbf{h}_{t_L}, \mathbf{h}_{hie}^L]$$

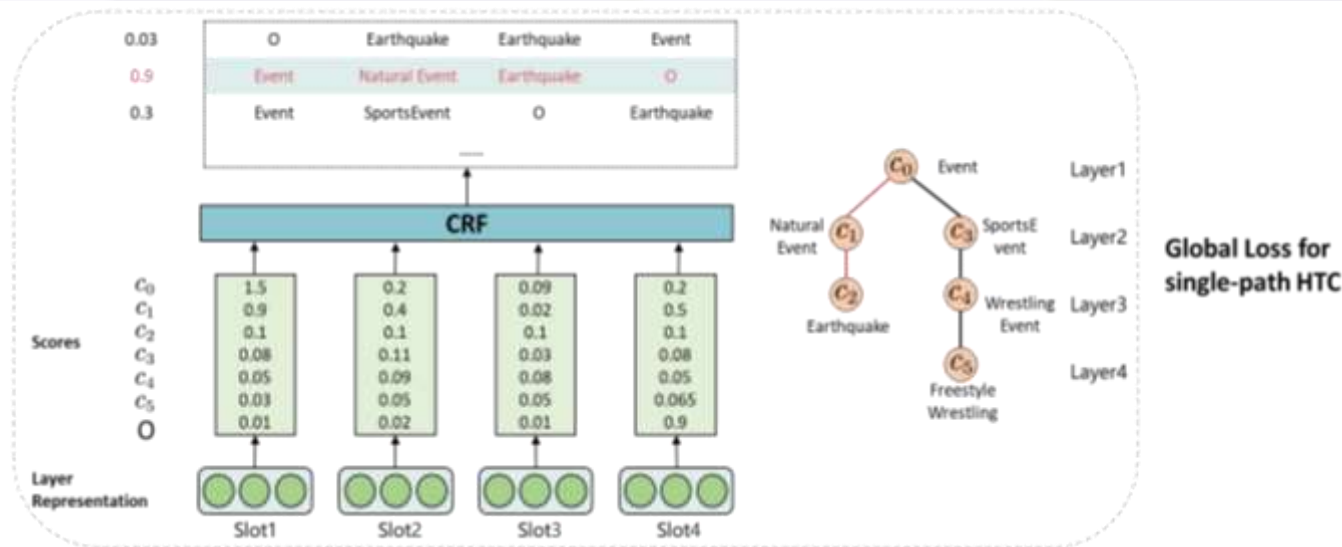
$$\mathbf{q}_{i,\alpha} = W_{q,\alpha} \mathbf{h}_{hie}^i + b_{q,\alpha} \quad Q_\alpha = [\mathbf{q}_{1,\alpha}, \mathbf{q}_{2,\alpha}, \dots, \mathbf{q}_{L,\alpha}]$$

$$\mathbf{k}_{j,\alpha} = W_{k,\alpha} \mathbf{h}_{hie}^j + b_{k,\alpha} \quad K_\alpha = [\mathbf{k}_{1,\alpha}, \mathbf{k}_{2,\alpha}, \dots, \mathbf{k}_{L,\alpha}]$$

$$s_\alpha(i, j) = (\mathcal{M}_i \mathbf{q}_{i,\alpha})^\top (\mathcal{M}_j \mathbf{k}_{j,\alpha})$$

$$= \mathbf{q}_{i,\alpha}^\top \mathcal{M}_{j-i} \mathbf{k}_{j,\alpha}$$

# Local Hierarchy-aware Structure for Single-path HTC



$$\mathbf{x}_{hie} = [x_{hie}^1, x_{hie}^2, \dots, x_{hie}^L] \quad \mathbf{y} = \{y_1, y_2, \dots, y_L\}$$

$$\Pr(\mathbf{y} | \mathbf{x}_{hie}) = \frac{\exp(\text{Score}(\mathbf{x}_{hie}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{P}} \exp(\text{Score}(\mathbf{x}_{hie}, \mathbf{y}'))}$$

$$\text{Score}(\mathbf{x}_{hie}, \mathbf{y}) = \sum_{i=1}^L Em_{i,y_i} + \sum_{i=0}^L Tr_{y_i,y_{i+1}} \quad \text{Score}(\mathbf{x}_{hie}, \mathbf{y}) = \sum_{i=1}^L l_{hie}^i [y_i] + Tr_{y_i,y_{i+1}}$$

## Objective Function

$$\mathcal{L} = \mathcal{L}_{MLM} + \lambda_1 \mathcal{L}_{Global} + \lambda_2 \mathcal{L}_{Local}$$

$$\mathcal{L}_{MLM} = -\log \frac{e^{s_t}}{\sum_{i=1}^V e^{s_i}} = \log \left( 1 + \sum_{i=1, i \neq t}^V e^{s_i - s_t} \right)$$

$$\mathcal{L}_{Global} = \sum_{m=1}^L \left( \log \left( 1 + \sum_{i \in \mathcal{U}_m^{neg}} e^{s_i} \right) + \log \left( 1 + \sum_{i \in \mathcal{U}_m^{pos}} e^{-s_i} \right) \right)$$

$$\mathcal{L}_{Local} = \sum_{\alpha \in P} \left( \log \left( 1 + \sum_{(i,j) \in \mathcal{Q}_\alpha} e^{s_\alpha(i,j)} \right) + \log \left( 1 + \sum_{(i,j) \in \mathcal{Q}_\alpha} e^{-s_\alpha(i,j)} \right) \right) \quad \text{multi-path}$$

$$\mathcal{L}_{Local} = -\log(\Pr(\mathbf{y} | \mathbf{x}_{hie})) = \log \left( \sum_{\mathbf{y}' \in P} \exp^{\text{Score}(\mathbf{x}_{hie}, \mathbf{y}')} \right) - \text{Score}(\mathbf{x}_{hie}, \mathbf{y}) \quad \text{single-path}$$



# Experiments

## Main Results and Consistency Experiments:

Model	WOS (Depth 2)		RCV1-V2 (Depth 4)		NYT (Depth 8)		Average	
	Micro-F1	Macro-F1	Micro-F1	Macro-F1	Micro-F1	Macro-F1	Micro-F1	Macro-F1
<b>Hierarchy-Aware Models</b>								
TextRCNN (Zhou et al., 2020)	83.55	76.99	81.57	59.25	70.83	56.18	78.65	64.14
HiAGM (Zhou et al., 2020)	85.82	80.28	83.96	63.35	74.97	60.83	81.58	68.15
HTCInfoMax (Deng et al., 2021)	85.58	80.05	83.51	62.71	-	-	-	-
HIMatch (Chen et al., 2021)	86.20	80.53	84.73	64.11	-	-	-	-
<b>Pretrained Language Models</b>								
BERT (Wang et al., 2022a)	85.63	79.07	85.65	67.02	78.24	65.62	83.17	70.57
BERT+HiAGM (Wang et al., 2022a)	86.04	80.19	85.58	67.93	78.64	66.76	83.42	71.63
BERT+HTCInfoMax (Wang et al., 2022a)	86.30	79.97	85.53	67.09	78.75	67.31	83.53	71.46
BERT+HIMatch (Chen et al., 2021)	86.70	81.06	86.33	68.66	-	-	-	-
HGCLR (Wang et al., 2022a)	87.11	81.20	86.49	68.31	78.86	67.96	84.15	72.49
HITIN (Zhu et al., 2023)	87.19	81.57	86.71	<b>69.95</b>	79.65	69.31	84.52	73.61
HITIN†	86.75	81.18	86.72	68.94	79.40	68.42	84.29	72.85
HPT (Wang et al., 2022b)	87.16	81.93	87.26	69.53	80.42	70.42	84.95	73.96
HPT†	86.93	81.50	87.39	69.11	80.59	70.45	84.97	73.69
NERHTC (Ours)	<b>87.42</b> <sup>†0.40</sup>	<b>81.93</b> <sup>†0.43</sup>	<b>87.50</b> <sup>†0.11</sup>	69.76 <sup>†0.63</sup>	<b>80.97</b> <sup>†0.38</sup>	<b>70.99</b> <sup>†0.54</sup>	<b>85.30</b> <sup>†0.33</sup>	<b>74.23</b> <sup>†0.54</sup>

Table 2: The experimental results (%) of our proposed method comparing to previous models on three datasets. Best results are in boldface. Our re-implementation scores are marked by "†". "†" indicates the improvement of our model compared with the second-best result within our reproduction.

Method	WOS			
	PMicro-F1	PMacro-F1	CMicro-F1	CMacro-F1
BERT	79.96	78.40	85.43	79.37
HITIN	81.06	79.23	86.45	80.76
HPT	80.69	79.03	86.57	80.85
NERHTC	<b>81.41</b>	<b>79.52</b>	<b>87.14</b>	<b>81.36</b>

Table 4: Consistency experiments of path-based (P-metric) and path-constrained (C-metric) evaluation metrics on WOS.

Method	RCV1-V2		NYT	
	CMicro-F1	CMacro-F1	CMicro-F1	CMacro-F1
BERT	85.68	66.96	78.05	64.62
HITIN	86.48	68.07	78.45	66.79
HPT	86.95	68.15	79.51	68.38
NERHTC	<b>86.99</b>	<b>68.46</b>	<b>80.11</b>	<b>69.42</b>

Table 5: Consistency experiments of path-constrained (C-metric) evaluation metrics on RCV1 and NYT.

## Other Experiments

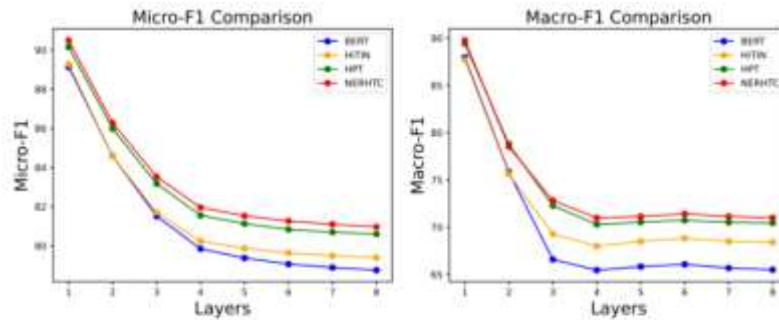


Figure 3: Preference study on label granularity of NYT based on layer increment.

Ablation Models	Micro-F1	Macro-F1
NERHTC	<b>80.97</b>	<b>70.99</b>
<i>r.m.</i> ROPE	80.80	70.40
<i>r.m.</i> MLM loss	80.67	70.40
<i>r.p.</i> BCE loss	80.70	70.10
<i>r.m.</i> GAT	80.82	70.25

Table 3: Ablation study results on NYT. *r.m.* stands for remove. *r.p.* stands for replaced with.

Models	Acc.
Llama 1-7B	<b>22.08</b>
<i>w/ demo</i>	19.05
Llama 2-7B	16.28
<i>w/ demo</i>	16.78

Table 6: The accuracy of using large models for inference on the first-level labels of WOS. *demo* represents demonstration. *w/* stands for with. The results are the average of three experiments

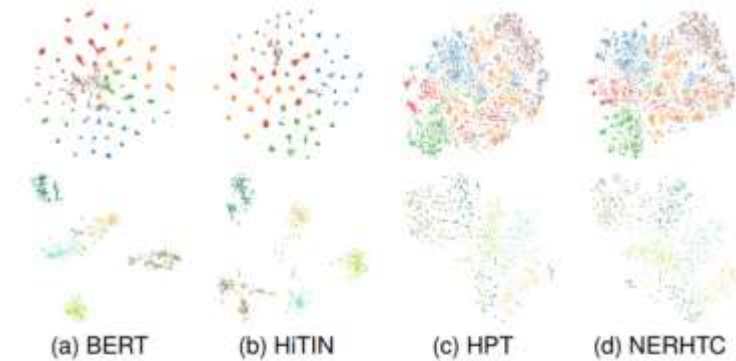


Figure 4: T-SNE visualization of the label representations on WOS. Images in the first row display feature clusters for the first-level labels, while the second row is for the sub-labels of ECE. Dots of the same color belong to the same category.

**Thank You!**