

## 数据集

数据集 A 取自《新五代史》《北史》等古籍，训练集，验证集和测试集分别为 21128, 2831, 2148, 平均句长为 26, 采用 BIOES, 共有 5 种实体类型, [O, B-NOUN\_BOOKNAME, I-NOUN\_BOOKNAME, B-NOUN-OTHER, I-NOUN-OTHER]

数据集 B 取自《春秋谷梁传》等古籍，训练集，验证集和测试集分别为 7528, 865, 877, 平均字长为 85, 采用 BIOES, 共有 13 种实体类型, [O, B-PER, I-PER, B-ORG, I-ORG, B-LOC, I-LOC, B-JOB, I-JOB, B-WAR, I-WAR, B-BOO, I-BOO]

实验数据下载:

<https://github.com/jizijing/C-CLUE>

<https://github.com/Ethan-yl/CCLUE>

## 评价指标

### 准确率(Accuracy)

Accuracy 是从整体上衡量模型的性能，即模型预测正确的样本占全部样本的比例。它的计算公式如下：

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

其中，TP (True Positive) 表示真正例，即模型预测为正例且实际上也为正例的样本数量；FP (False Positive) 表示假正例，即模型预测为正例但实际上为负例的样本数量；TN (True Negative) 表示真负例，即模型预测为负例且实际上也为负例的样本数量；FN (False Negative) 表示假负例，即模型预测为负例但实际上为正例的样本数量。

### 精确率(Precision)

Precision 关注的是模型预测为正例的样本中有多少是真正的正例。换句话说，它衡量的是模型预测的正例中有多少是正确的，计算公式为：

$$Precision = \frac{TP}{TP + FP}$$

其中，TP (True Positive) 表示真正例，即模型预测为正例且实际上也为正例的样本数量；FP (False Positive) 表示假正例，即模型预测为正例但实际上为负例的样本数量。

### F1 值(F1 Score)

F1 分数是精确率和召回率的调和平均值，用于衡量模型的准确性。它的计算公式如下：

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

其中，精确率 (Precision) 是指模型预测为正例的样本中真正为正例的比例，召回率 (Recall) 是指所有真正的正例样本中被模型预测为正例的比例。

### 召回率(Recall)

召回率是指在所有真正的正例样本中，被模型正确预测为正例的比例。它的计算公式如下：

$$Recall = \frac{TP}{TP + FN}$$

其中，TP（True Positive）表示真正例，即模型预测为正例且实际上也为正例的样本数量；FN（False Negative）表示假负例，即模型预测为负例但实际上为正例的样本数量。

**Micro-Average**

Micro-Average 把所有类别的结果汇总起来计算平均值。它将所有类别的贡献视为等同，因此对于样本量大的类别，Micro-Average 更加敏感。

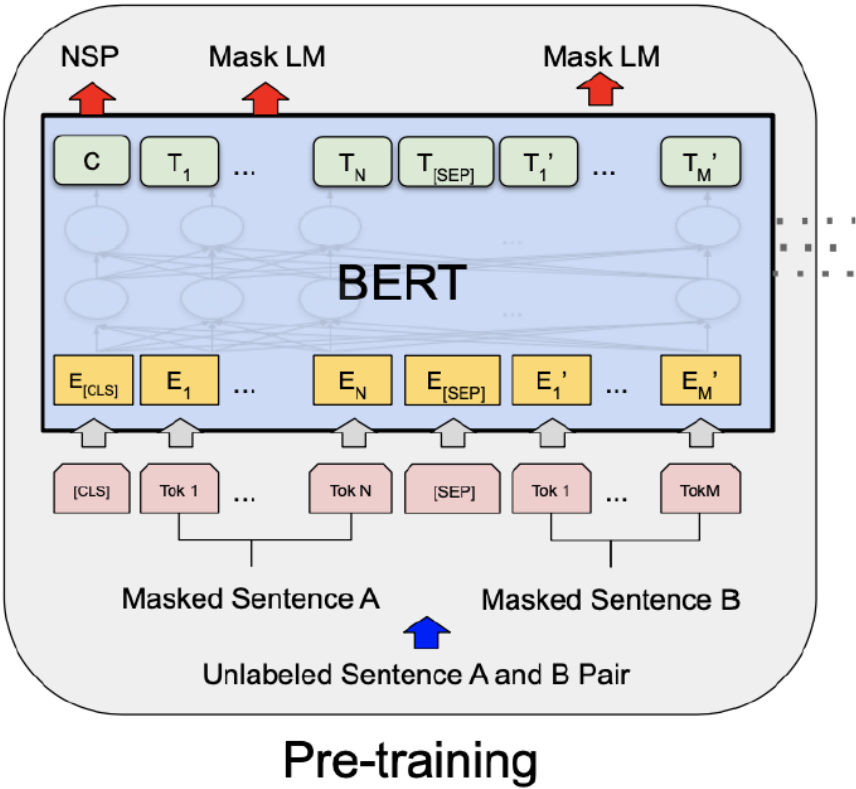
$$\begin{aligned} \text{Micro-Average Precision} &= \frac{\sum TP}{\sum TP + \sum FP} \\ \text{Micro-Average Recall} &= \frac{\sum TP}{\sum TP + \sum FN} \\ \text{Micro-Average F1} &= 2 \times \frac{(\text{Micro Precision} \times \text{Micro Recall})}{(\text{Micro Precision} + \text{Micro Recall})} \end{aligned}$$

**Macro-Average**

Macro-Average 分别计算每个类别的性能指标，然后计算这些指标的算术平均值。它给予每个类别同等的权重，无论类别的样本量大小，因此对于样本量小的类别更加敏感。

$$\begin{aligned} \text{Macro-Average Precision} &= \frac{1}{N} \sum \frac{TP_i}{TP_i + FP_i} \\ \text{Macro-Average Recall} &= \frac{1}{N} \sum \frac{TP_i}{TP_i + FN_i} \\ \text{Macro-Average F1} &= \frac{1}{N} \sum 2 \times \frac{(\text{Precision}_i \times \text{Recall}_i)}{(\text{Precision}_i + \text{Recall}_i)} \end{aligned}$$

模型  
BERT



BertModel 部分包含了 BERT 的所有主要组件，包括词嵌入层、位置嵌入层、令牌类型嵌入层、层归一化层和 Dropout 层。

BertEncoder 部分负责处理输入序列，通过多个 BERT 层来进行特征提取。每个 BERT 层都包含注意力机制和前馈神经网络。

BertForTokenClassification 模型的最后一部分是一个线性分类器，它将 BERT 编码器的输出映射到目标类别的概率分布上。

- 1.Transformer 架构：** BERT 建立在 Transformer 模型的基础上，这是一种使用自注意力机制（Self-Attention Mechanism）的深度学习神经网络。Transformer 允许模型在处理序列数据时同时关注序列中的所有位置，而不是像传统的循环神经网络 (RNN) 或卷积神经网络 (CNN) 那样逐步处理。
- 2.预训练策略：** BERT 采用了无监督的预训练策略，通过大规模的语言模型预训练来学习丰富的语义表示。该模型通过对大量文本数据进行“遮蔽语言模型” (Masked Language Model, MLM) 任务的预训练，使得模型能够理解词汇和语法结构，并捕捉单词之间的关系。
- 3.双向性：** BERT 在预训练时考虑了双向信息，即使用上下文信息来理解每个词的语义。这种双向性有助于模型更好地理解文本中的语境和关联，提高了对上下文相关性的捕捉。
- 4.Fine-tuning：** 预训练后，BERT 模型可以通过微调 (fine-tuning) 来适应特定的下游任务，如命名实体识别、情感分析等。这种能力使得 BERT 在各种 NLP 任务中都表现出色，无需从零开始训练新的模型。
- 5.Contextual Embeddings：** BERT 生成的词向量是上下文相关的，每个词的表示取决于整个输入句子的上下文，而不是简单地从固定的嵌入中获取。这种上下文敏感的嵌入有助于更准确地捕捉语义信息。

实验结果

Bert-ancient-Chinese

训练最终结果

数据集	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
数据集 A	0.9730221685457736	0.9734486209964412	0.9732050340527741	0.9734486209964412	0.15202167630195618	4.6471
数据集 B	0.9892003099895879	0.9911334325396826	0.9900857673216377	0.9911334325396826	0.025596966966986656	10.0205

数据集	Micro-Average Precision	Micro-Average Recall	Micro-Average F1	Macro-Average Precision	Macro-Average Recall	Macro-Average F1
数据集 A	0.9734486209964412	0.9734486209964412	0.9734486209964412	0.8710822709848423	0.8494438276449081	0.8599516197788132

数据集 B	0.9911334	0.9911334	0.9911334	0.24364633	0.21699121	0.21620752
	325396826	325396826	325396826	330184382	039347588	290443754

模型介绍和分析

这是一个基于 BERT 的模型，专门训练用于古汉语处理。它适用于古文的语义理解和文本生成，旨在提高对古汉语的处理能力。

训练过程分析

数据集 A

Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.9651031456399267	0.9666370106761566	0.9651928655745571	0.9666370106761566	0.10154145956039429	4.6365
4	0.9678706559288105	0.9669984430604982	0.9673165804611824	0.9669984430604982	0.10396654903888702	4.5978
6	0.9714789822369363	0.9719472864768683	0.9715817412063057	0.9719472864768683	0.09987952560186386	4.6208
8	0.9718535415569596	0.9717804715302492	0.9717633388042668	0.9717804715302492	0.13383987545967102	4.6113
10	0.9732551487821443	0.973810053380783	0.9734712820803293	0.973810053380783	0.13705913722515106	4.6348

数据集 B

Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.98443603515625	0.9921875	0.9882965686274511	0.9921875	0.029304953292012215	7.6138
4	0.9874243224761378	0.9916648951247166	0.9894932891321785	0.9916648951247166	0.0257867444306612	7.5513
6	0.9880146702085819	0.9918774801587301	0.989184204109248	0.9918774801587301	0.024214763194322586	7.4568
8	0.9886410337332172	0.991328302154195	0.9898839745068818	0.991328302154195	0.024256916716694832	7.5093
10	0.9895343898980461	0.9911777210884354	0.9902227660139274	0.9911777210884354	0.025142531841993332	7.4757

SikuRoberta

训练最终结果

数据集	Eval	Eval Recall	Eval F1	Eval	Eval Loss	Eval
-----	------	-------------	---------	------	-----------	------

	Precision			Accuracy		Runtime
数据集 A	0.97104 9537254 9688	0.9716136 565836299	0.9712887 62010626 7	0.97161365 65836299	0.15085685 25314331	3.5724
数据集 B	0.98963 9677599 6501	0.9911068 594104309	0.9903303 96788302 3	0.99110685 94104309	0.02541288 547217846	7.8184

数据集	Micro-Average Precision	Micro-Average Recall	Micro-Average F1	Macro-Average Precision	Macro-Average Recall	Macro-Average F1
数据集 A	0.971613 65658362 99	0.971613 65658362 99	0.971613 65658362 99	0.864802 89745921 76	0.842991 83475158 23	0.853425 21643909 03
数据集 B	0.991106 85941043 09	0.991106 85941043 09	0.991106 85941043 09	0.231851 15660763 794	0.230711 27190351 204	0.221999 41287605 457

### 模型介绍和分析

这个模型是 Siku Quanshu（四库全书）的基础上训练的 RoBERTa 变体。它对古典文献中的文本理解有很强的能力，适合处理古籍和文献分析任务。

### 训练过程分析

#### 数据集 A

Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.96083311 8499789	0.96160475 97864769	0.96014807 41181836	0.96160475 97864769	0.118305407 46450424	3.56 16
4	0.96757056 9968203	0.96805493 77224199	0.96746086 79832772	0.96805493 77224199	0.103166840 9705162	3.55 9
6	0.96917842 68414622	0.96994550 71174378	0.96946201 4966379	0.96994550 71174378	0.109539739 78757858	3.57 63
8	0.97179422 33651111	0.97228091 63701067	0.97200139 51812252	0.97228091 63701067	0.130345419 049263	3.57 41
10	0.97108461 94442453	0.97150244 6619217	0.97126974 69725421	0.97150244 6619217	0.137427598 23799133	3.55 55

#### 数据集 B

Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.98443603 515625	0.9921875	0.98829656 86274511	0.9921875	0.029090946 540236473	5.43 47
4	0.98735217	0.99180661	0.98954351	0.99180661	0.025082616	5.27

	26932424	84807256	36377729	84807256	13547802	25
6	0.98752439 2029138	0.99191291 09977324	0.98926084 41941434	0.99191291 09977324	0.023758890 107274055	5.23 86
8	0.98878774 15659296	0.99158517 57369615	0.99002233 63677141	0.99158517 57369615	0.024264896 288514137	5.46 21
10	0.98936904 56522269	0.99123086 73469388	0.99020889 12819328	0.99123086 73469388	0.024774728 34289074	5.23 45

## AnchiBERT

### 训练最终结果

数据集	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
数据集 A	0.96772 6471010 9994	0.9684719 75088968	0.9680270 67726180 2	0.96847197 5088968	0.16520337 760448456	3.0608
数据集 B	0.98981 3124692 8284	0.9911068 594104309	0.9903980 65962140 8	0.99110685 94104309	0.02563157 305121421 8	7.6369

数据集	Micro-Average Precision	Micro-Average Recall	Micro-Average F1	Macro-Average Precision	Macro-Average Recall	Macro-Average F1
数据集 A	0.9684719 75088968	0.9684719 75088968	0.9684719 75088968	0.8426885 820191009	0.82060433 51382434	0.83101526 89098269
数据集 B	0.9911068 594104309	0.9911068 594104309	0.9911068 594104309	0.2512589 632357055	0.24764968 296968584	0.23697334 205426301

### 模型介绍和分析

AnchiBERT 是针对古代汉语的 BERT 变体，提供了对古汉语语料的更深入的理解。它优化了 BERT 架构以适应古汉语的特殊需求，适合古文翻译和解析。

### 训练过程分析

#### 数据集 A

Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.96003681 25761395	0.96202179 71530249	0.96050850 30338348	0.96202179 71530249	0.116751648 48566055	3.04 96
4	0.96362474 02560569	0.96421819 39501779	0.96360568 10176072	0.96421819 39501779	0.116235248 74448776	3.07 19

6	0.96702017 12841807	0.96783251 77935944	0.96734084 21844694	0.96783251 77935944	0.116239771 2469101	3.06 06
8	0.96688934 23942839	0.96844417 25978647	0.96726425 57197699	0.96844417 25978647	0.143406972 2890854	3.05 03
10	0.96825936 09860043	0.96897241 99288257	0.96852979 74658393	0.96897241 99288257	0.151419401 16882324	3.06 45

数据集 B

Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.98612550 86400601	0.99209892 29024944	0.98896008 64359019	0.99209892 29024944	0.027748765 42389393	5.07 94
4	0.98649524 49481542	0.99194834 18367347	0.98906720 28096879	0.99194834 18367347	0.024718163 534998894	4.82 41
6	0.98849344 37504811	0.99170032 59637188	0.98993204 85624267	0.99170032 59637188	0.024498600 512742996	4.98 91
8	0.98926487 99954938	0.99128401 36054422	0.99021190 89568483	0.99128401 36054422	0.024361521 005630493	5.38 6
10	0.98922827 63694249	0.99119543 65079365	0.99011360 42309379	0.99119543 65079365	0.025130053 982138634	4.86 51

## guwenbert-base

训练最终结果

数据集	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
数据集 A	0.97582891 07635478	0.97620106 76156584	0.97598140 42908822	0.97620106 76156584	0.1170888915 6579971	2.75 89
数据集 B	0.99591997 07856833	0.99592545 35147393	0.99586389 48914785	0.99592545 35147393	0.0130802895 87378502	6.87 21

数据集	Micro-Average Precision	Micro-Average Recall	Micro-Average F1	Macro-Average Precision	Macro-Average Recall	Macro-Average F1
数据集 A	0.97620106 76156584	0.97620106 76156584	0.97620106 76156584	0.88416424 31251063	0.86803984 26337851	0.87581740 96034843

数据集 B	0.99592545 35147393	0.99592545 35147393	0.99592545 35147393	0.40479139 253470026	0.41109825 24602927	0.39143839 70521158
-------	------------------------	------------------------	------------------------	-------------------------	------------------------	------------------------

模型介绍和分析

guwenbert-base 是一个专门为古文设计的 BERT 模型。它利用大量古文语料进行训练，目标是提高古文文本的理解和处理能力。

训练过程分析

数据集 A

Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.96616543 25093485	0.96741548 04270463	0.96662958 24622826	0.96741548 04270463	0.110058732 3307991	2.69 51
4	0.96761796 2053998	0.96880560 49822064	0.96764848 92367837	0.96880560 49822064	0.103007376 19400024	2.66 98
6	0.97278685 40697165	0.97264234 87544484	0.97263619 29489298	0.97264234 87544484	0.094608858 2277298	2.66 26
8	0.97489096 85689006	0.97486654 80427047	0.97477620 81255394	0.97486654 80427047	0.112511061 1319542	2.65 97
10	0.97490283 35985271	0.97503336 29893239	0.97491338 71389762	0.97503336 29893239	0.115471594 03562546	2.72 61

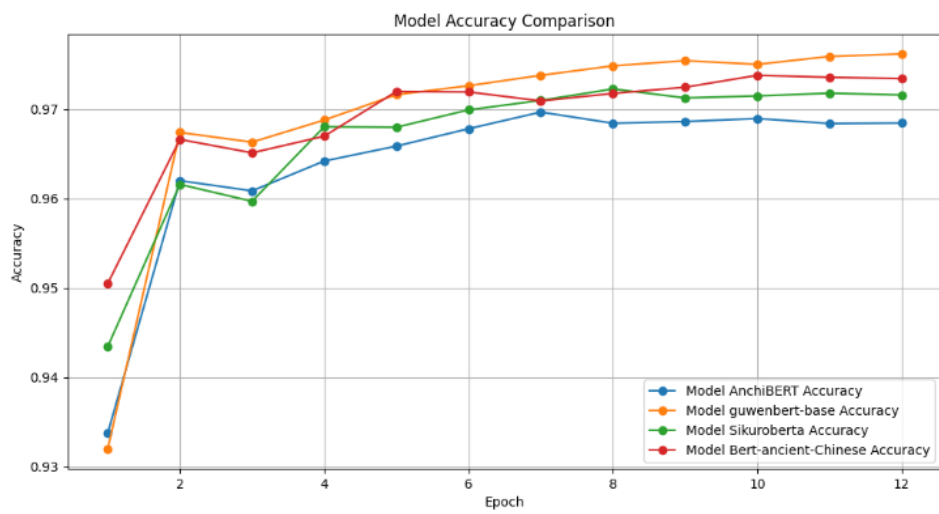
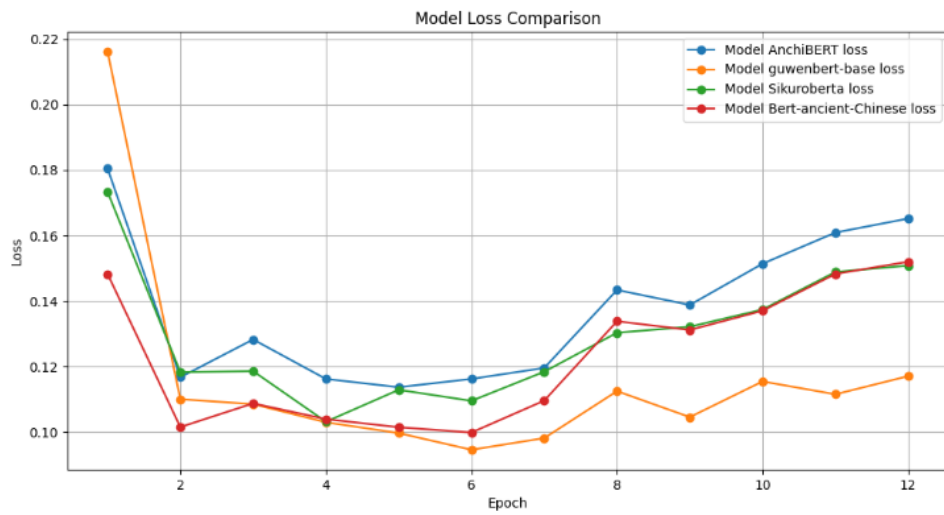
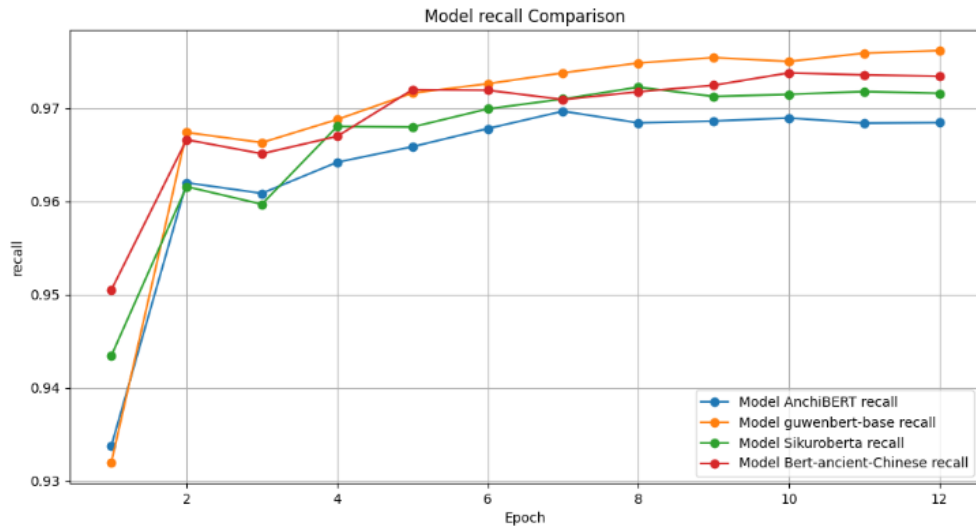
数据集 B

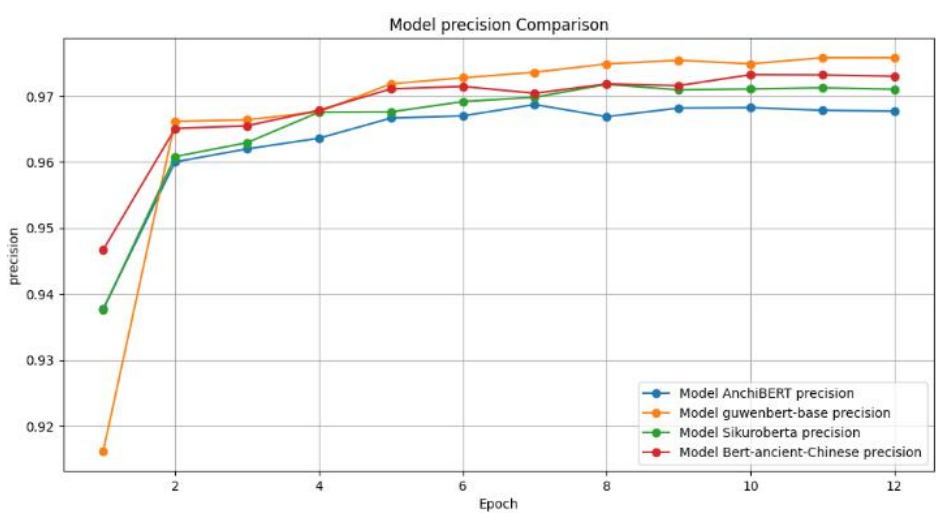
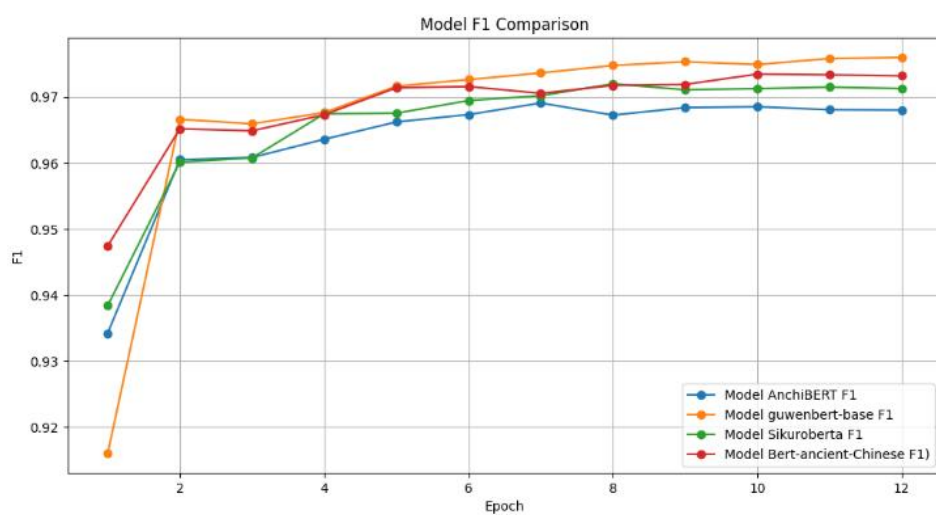
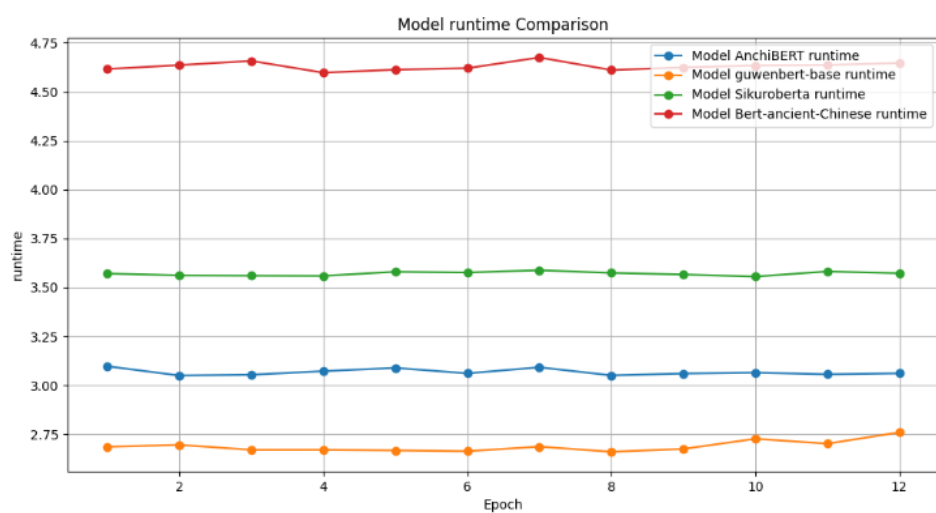
Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.99313401 95989583	0.99490681 68934241	0.99359175 61585891	0.99490681 68934241	0.020061865 44895172	4.01 95
4	0.99412462 81408016	0.99546485 26077098	0.99449116 69147691	0.99546485 26077098	0.016725214 198231697	4.06 48
6	0.99547706 1732981	0.99561543 36734694	0.99525561 96658517	0.99561543 36734694	0.014112563 803792	4.27 64
8	0.99588148 85623937	0.99625318 87755102	0.99586504 30433779	0.99625318 87755102	0.012881815 433502197	4.10 83
10	0.99601992 91173533	0.99594316 89342404	0.99591807 74967511	0.99594316 89342404	0.013141417 875885963	4.61 96

比较分析

数据集 A







## guwenbert-base

guwenbert-base 在各方面表现优越。f1、precision、accuracy、recall 在 1-2 个 epoch 之间

有显著提升，在后续 epoch 中表现优秀，处于最高位。loss 在 1-2 个 epoch 之间有显著下降，在后续 epoch 中表现优秀，处于最低位。训练时间在四个模型中最短。

### Bert-ancient-Chinese

Bert-ancient-Chinese 在整个训练中相对优秀。f1、precision、accuracy、recall 在 1-2 个 epoch 之间变化最小，在后续 epoch 中表现稳定，处于第二高位。loss 在 1-2 个 epoch 之间也变化最小，在后续 epoch 中表现平稳，处于第二低位。但训练时间在四个模型中最长。

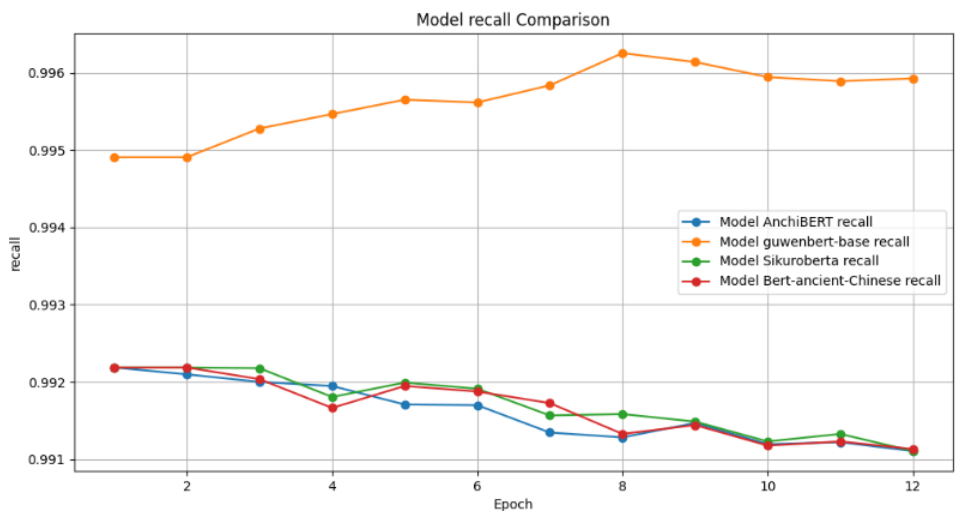
### AnchiBERT

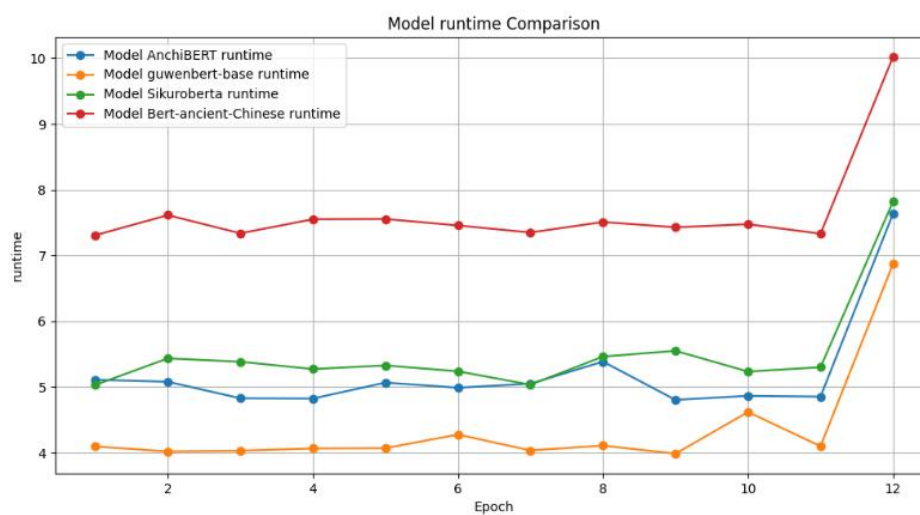
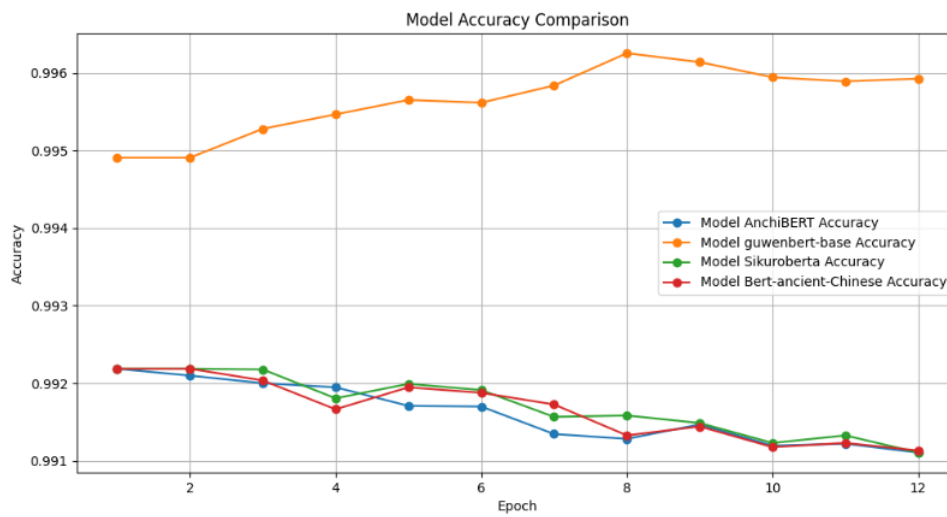
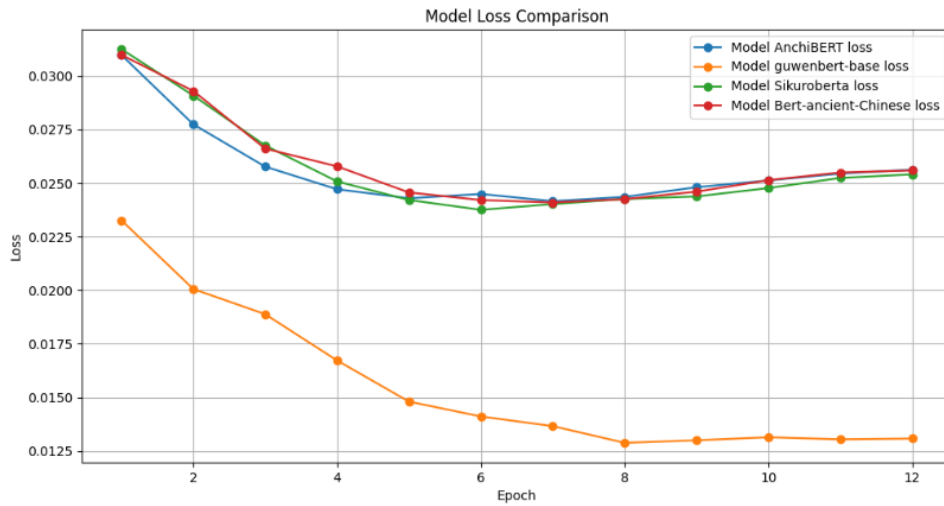
AnchiBERT 在整个训练中相对差。f1、precision、accuracy、recall 在 1-2 个 epoch 之间变化较小，在后续 epoch 中表现稳定，处于最低位。loss 在 1-2 个 epoch 之间也变化较小，在后续 epoch 中表现平稳，处于最高位。训练时间在四个模型中是第二短。

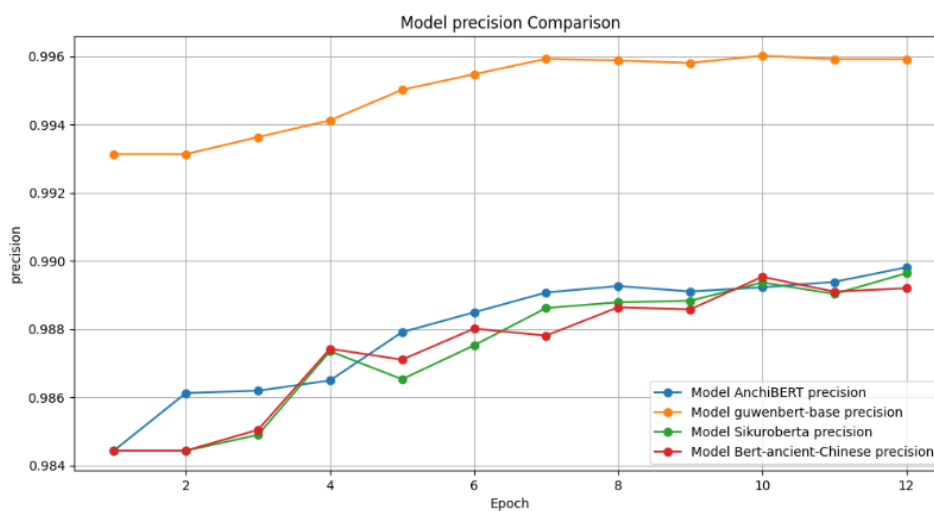
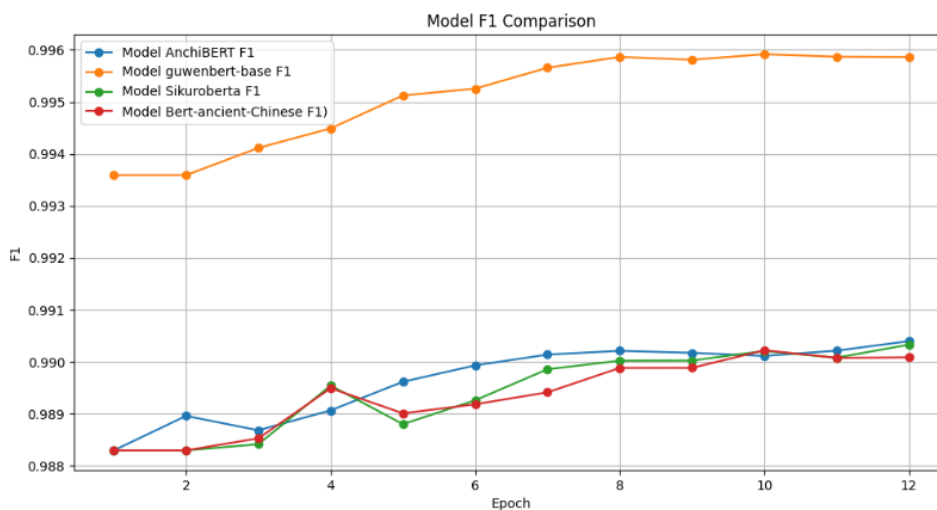
### SikuRoberta

SikuRoberta 与 AnchiBERT 在整个训练中类似，各方面比 AnchiBERT 相对较好，但训练时间处在第二长。

## 数据集 B







## guwenbert-base

guwenbert-base 在各方面表现优越。Recall、accuracy、f1 在训练中逐步提升，precision 在 1-6 个 epoch 之间有显著上升，在后续 epoch 中较为平稳，处于最高位。loss 在 1-6 个 epoch 之间有显著下降，在后续 epoch 中较为平稳，处于最低位。训练时间在四个模型中最短。

## Bert-ancient-Chinese

Bert-ancient-Chinese 在各方面都明显不如 guwenbert-base。训练时间在四个模型中最长。

## AnchiBERT

AnchiBERT 在各方面都明显不如 guwenbert-base，Recall、accuracy、f1、precision、loss 同 Bert-ancient-Chinese 的相似。训练时间在四个模型中处在第二短。

## SikuRoberta

SikuRoberta 在各方面都明显不如 guwenbert-base, Recall、accuracy、f1、precision、loss 同 Bert-ancient-Chinese 的相似。训练时间在四个模型中处在第二长。

综合比较

guwenbert-base 表现相对出色，且在整个训练过程中稳定增长。

实验整体分析

数据集 A

模型	Micro-Average Precision	Micro-Average Recall	Micro-Average F1	Macro-Average Precision	Macro-Average Recall	Macro-Average F1
Bert-ancient-Chinese	0.9734486209964412	0.9734486209964412	0.9734486209964412	0.8710822709848423	0.8494438276449081	0.8599516197788132
SikuRoberta	0.9716136565836299	0.9716136565836299	0.9716136565836299	0.8648028974592176	0.8429918347515823	0.8534252164390903
AnchiBERT	0.968471975088968	0.968471975088968	0.968471975088968	0.8426885820191009	0.8206043351382434	0.8310152689098269
guwenbert-base	0.9762010676156584	0.9762010676156584	0.9762010676156584	0.8841642431251063	0.8680398426337851	0.8758174096034843

数据集 B

模型	Micro-Average Precision	Micro-Average Recall	Micro-Average F1	Macro-Average Precision	Macro-Average Recall	Macro-Average F1
Bert-ancient-Chinese	0.9911334325396826	0.9911334325396826	0.9911334325396826	0.24364633330184382	0.21699121039347588	0.21620752290443754
SikuRoberta	0.9911068594104309	0.9911068594104309	0.9911068594104309	0.23185115660763794	0.23071127190351204	0.22199941287605457
AnchiBERT	0.9911068594104309	0.9911068594104309	0.9911068594104309	0.2512589632357055	0.24764968296968584	0.23697334205426301
guwenbert-base	0.9959254535147393	0.9959254535147393	0.9959254535147393	0.40479139253470026	0.4110982524602927	0.3914383970521158

**Micro-Average 指标**反映了所有类别的总体性能，这意味着在考虑样本不平衡的情况下，guwenbert-base 在整体上取得了最好的平衡性能。

**Macro-Average 指标**展示了模型对所有类别同等对待的性能，guwenbert-base 表现最佳，表明其对于少数类别的识别能力较强。