

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。它把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。

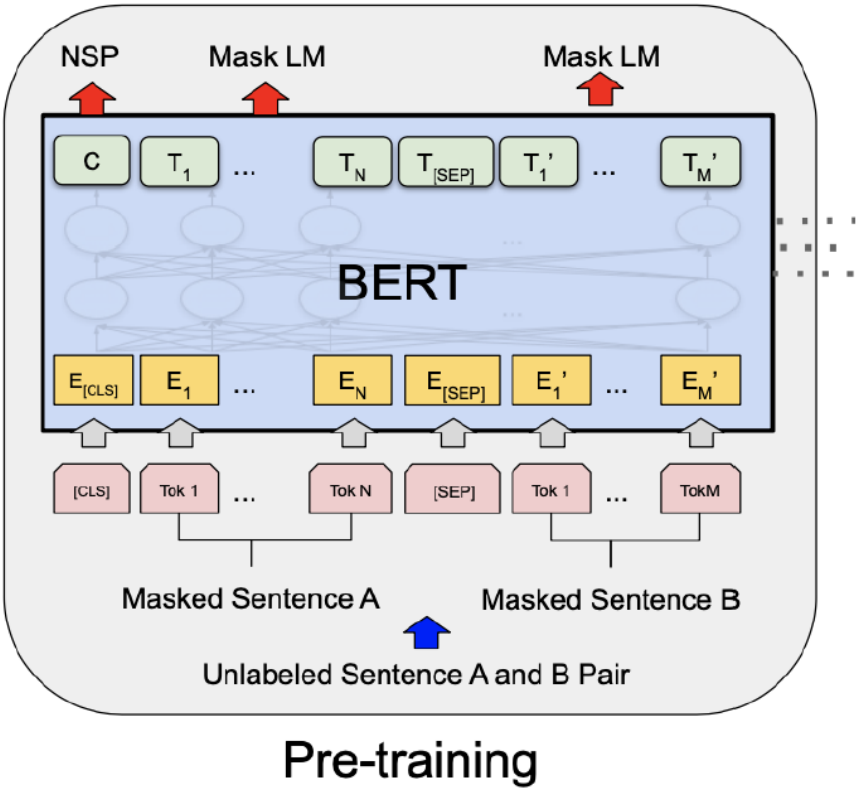
$$\begin{aligned} \text{Micro-Average Precision} &= \frac{\sum TP}{\sum TP + \sum FP} \\ \text{Micro-Average Recall} &= \frac{\sum TP}{\sum TP + \sum FN} \\ \text{Micro-Average F1} &= 2 \times \frac{(\text{Micro Precision} \times \text{Micro Recall})}{(\text{Micro Precision} + \text{Micro Recall})} \end{aligned}$$

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	0.9713613930185667	0.9717125382262997	0.9715153155843004	0.9717125382262997	0.1591310352087021	4.3364
数据集 B	0.9607732624971154	0.9653663548752834	0.9627669359301028	0.9653663548752834	0.17902402579784393	8.1976

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	0.9717125382262997	0.9717125382262997	0.9717125382262997	0.8673289276579453	0.8478130526942544	0.8573962308701384

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	0.965859360425647	0.9673042704626335	0.9661755533216788	0.9673042704626335	0.10091111063957214	3.5747
4	0.9695756273928118	0.9709185943060499	0.9698405783513404	0.9709185943060499	0.09781475365161896	3.5827
6	0.9748397118458895	0.9753113879003559	0.9749971259961178	0.9753113879003559	0.10068126767873764	3.5825
8	0.9729364759002894	0.9735042259786477	0.9731642014794695	0.9735042259786477	0.11649686098098755	3.5887
10	0.9743676182461154	0.9747553380782918	0.974540298619832	0.9747553380782918	0.13155880570411682	3.5785

数据集 B

Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.9634874417593906	0.9647728883219955	0.9636102287355702	0.9647728883219955	0.08292320370674133	7.6235
4	0.9611707205126476	0.9681919642857143	0.9634537970119926	0.9681919642857143	0.08083999902009964	7.7177
6	0.9619109614420617	0.9659775368480725	0.9631737884030217	0.9659775368480725	0.10151440650224686	7.5857
8	0.9607420133460248	0.9666772959183674	0.9630476118906922	0.9666772959183674	0.13081666827201843	7.5832
10	0.9608151673880866	0.9650031887755102	0.9626049702565307	0.9650031887755102	0.15603597462177277	7.5866

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	0.96811 2015577 6729	0.9683438 455657493	0.9682108 01617711 9	0.96834384 55657493	0.15556995 570659637	4.3383
数据集 B	0.96170 5493661 1624	0.9659598 214285714	0.9635556 88824680 2	0.96595982 14285714	0.15483438 968658447	8.1275

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	0.968343 84556574 93	0.968343 84556574 93	0.968343 84556574 93	0.836226 94752648 86	0.835406 76629689 05	0.835655 05209439 84
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	0.96370769 86681025	0.96394016 90391459	0.96373988 33761893	0.96394016 90391459	0.112768314 77880478	3.59 75
4	0.96781548 71484845	0.96880560 49822064	0.96748257 52554453	0.96880560 49822064	0.099622294 30675507	3.57 92
6	0.97051037 10023891	0.97069617 43772242	0.97053512 02725301	0.97069617 43772242	0.112534001 46961212	3.56 79
8	0.97274631 90910541	0.97333741 10320284	0.97298543 07448372	0.97333741 10320284	0.120456494 39096451	3.57 99
10	0.97290976 87048657	0.97350422 59786477	0.97314598 18475224	0.97350422 59786477	0.139481723 30856323	3.58 32

数据集 B

Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.96122561 92641184	0.96793509 07029478	0.96301542 82453001	0.96793509 07029478	0.080131962 89539337	7.49 49
4	0.96121514 59441082	0.96965348 63945578	0.96223684 60374825	0.96965348 63945578	0.079407051 20563507	7.74 77

6	0.96174086 35243314	0.96769593 25396826	0.96372869 30779716	0.96769593 25396826	0.090674847 36442566	7.79 6
8	0.96191942 89796517	0.96650899 94331065	0.96386045 28771143	0.96650899 94331065	0.117915615 43941498	7.70 83
10	0.96180479 87592484	0.96613697 56235828	0.96355845 47151035	0.96613697 56235828	0.130092933 77399445	7.91 48

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	0.96394 0242927 4961	0.9646167 813455657	0.9642436 22737465	0.96461678 13455657	0.18482051 78976059	3.9459
数据集 B	0.96140 6071894 0829	0.9663495 606575964	0.9634749 39687441	0.96634956 06575964	0.16236752 271652222	7.9751

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	0.9646167 813455657	0.9646167 813455657	0.9646167 813455657	0.8273717 706255621	0.80382831 03729272	0.81531084 07512697
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	0.95918495 04729812	0.96146574 73309609	0.95983018 22708686	0.96146574 73309609	0.116816453 6356926	3.08 97
4	0.96398771 51221328	0.96527468 86120996	0.96426760 70302236	0.96527468 86120996	0.110925152 89783478	3.13 4
6	0.96665166	0.96719306	0.96687019	0.96719306	0.120304748	3.14

	36368039	04982206	72485794	04982206	41594696	27
8	0.96622835 21442985	0.96752669 03914591	0.96662462 75754931	0.96752669 03914591	0.133583903 3126831	3.06 76
10	0.96876239 78945292	0.96950066 72597865	0.96908259 48759366	0.96950066 72597865	0.154163137 07828522	3.05 11

数据集 B

Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.96407831 97007398	0.96525120 46485261	0.96390404 02066386	0.96525120 46485261	0.082221254 70638275	9.69 45
4	0.96121962 33374363	0.96908659 29705216	0.96293416 77537446	0.96908659 29705216	0.081010073 42338562	9.34 74
6	0.96147254 17305747	0.96737705 49886621	0.96339996 4029979	0.96737705 49886621	0.092464007 43722916	7.18 56
8	0.96124318 50099584	0.96758964 00226758	0.96360780 31217096	0.96758964 00226758	0.134863793 84994507	6.72 07
10	0.96141629 98781131	0.96671272 67573696	0.96350522 48092323	0.96671272 67573696	0.137260258 19778442	6.67 15

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	0.898053521 4526109	0.916260512 2324159	0.900511550 865519	0.916260512 2324159	0.304275244 474411	4.35 13
数据集 B	0.950169384 2044238	0.961734693 8775511	0.954995042 3377727	0.961734693 8775511	0.138649314 64195251	8.13 59

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	0.9162605 122324159	0.9162605 122324159	0.9162605 122324159	0.65320906 79952062	0.3968061 30054585	0.46904291 833949097
数据集 B	0.9959254	0.9959254	0.9959254	0.40479139	0.4110982	0.39143839

数据集 B	535147393	535147393	535147393	253470026	524602927	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

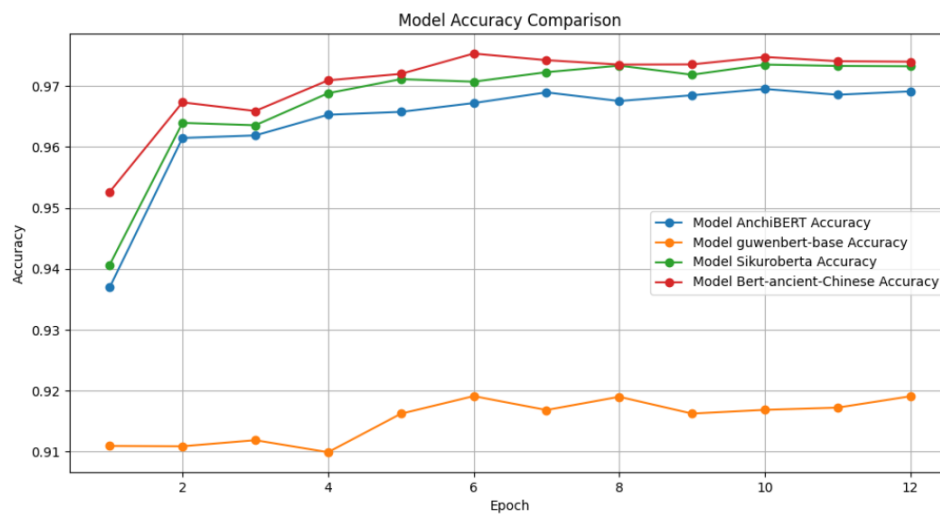
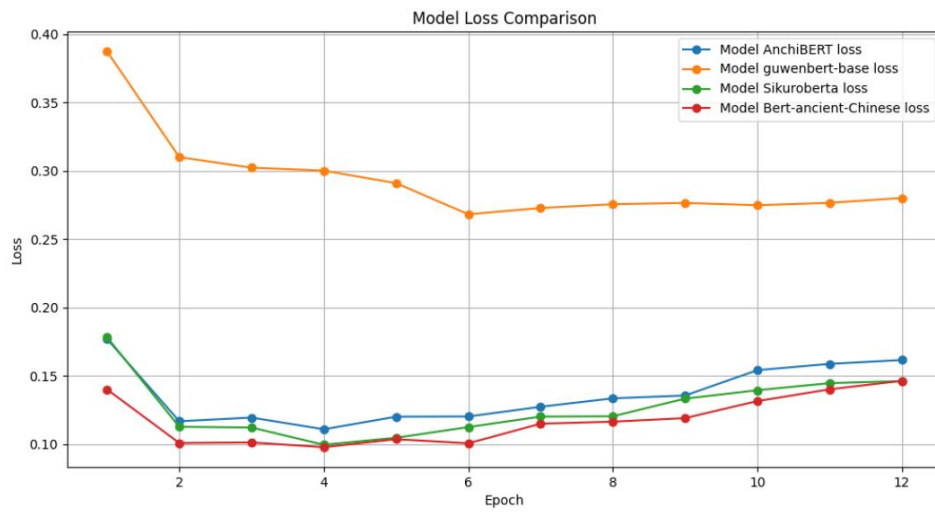
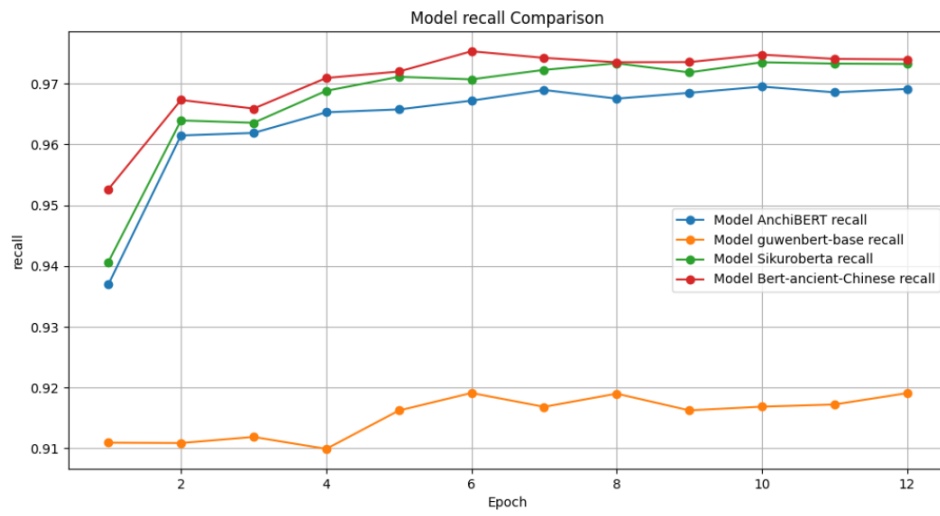
Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.889935668553033	0.9108930160142349	0.887439981938873	0.9108930160142349	0.3099377155303955	3.5996
4	0.8930562664726759	0.9099199288256228	0.8947124462689304	0.9099199288256228	0.30007246136665344	3.5699
6	0.9070822079888535	0.919122553380783	0.8958977944838243	0.919122553380783	0.2681475579738617	3.5973
8	0.9026081384861598	0.9190113434163701	0.9012913717862462	0.9190113434163701	0.27551373839378357	3.5942
10	0.8998557531360152	0.9168705516014235	0.9025673921036375	0.9168705516014235	0.27476179599761963	3.6135

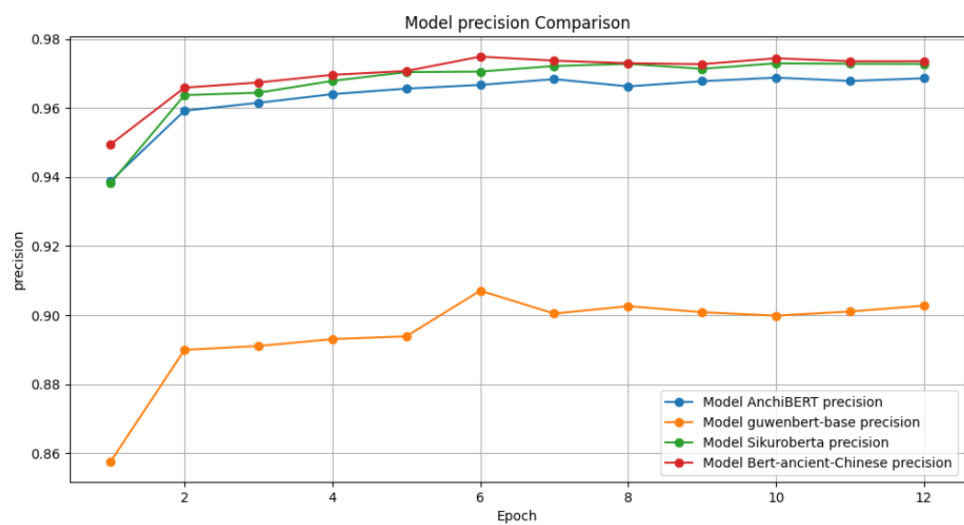
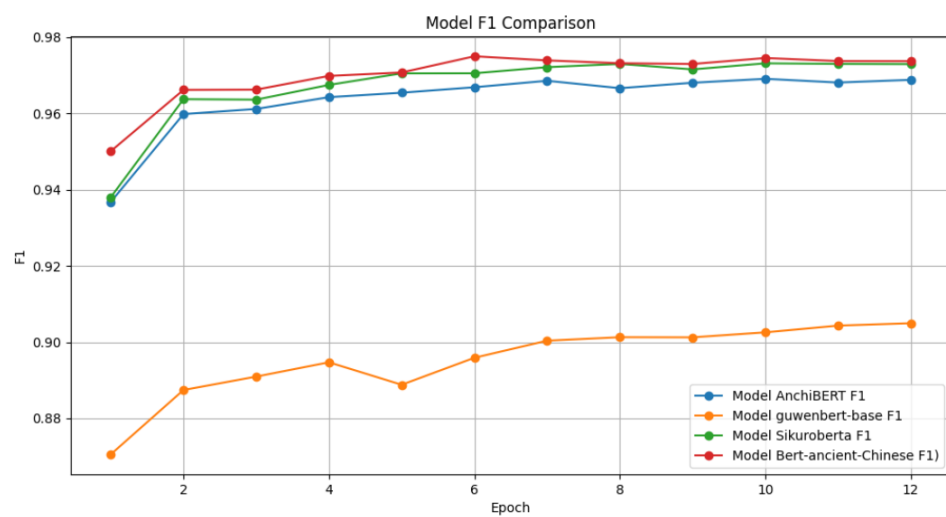
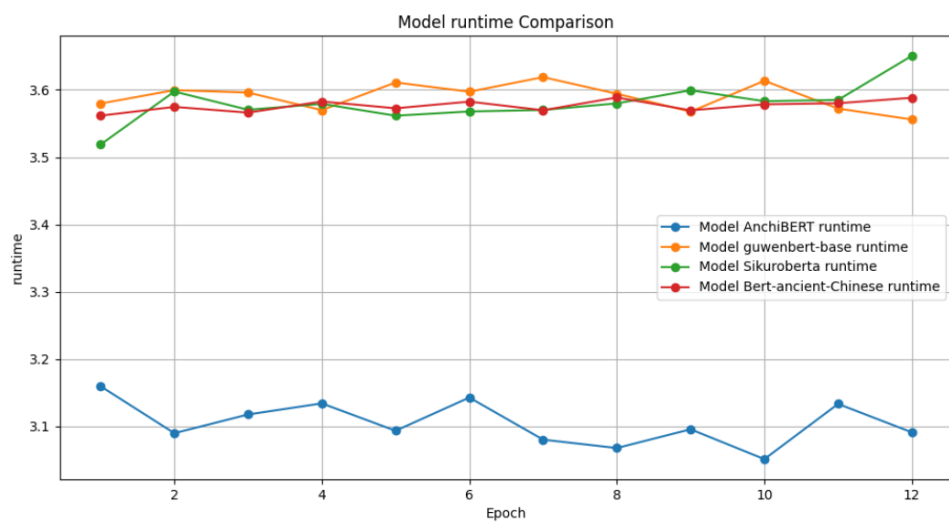
数据集 B

Epoch	Eval Precision	Eval Recall	Eval F1	Eval Accuracy	Eval Loss	Eval Runtime
2	0.9455200607648422	0.9664912840136054	0.9503247965618769	0.9664912840136054	0.14051271975040436	7.7432
4	0.9460385929088998	0.9663407029478458	0.9513311115749682	0.9663407029478458	0.1322239190340042	7.8783
6	0.9501682554545535	0.9662078373015873	0.9526954770627326	0.9662078373015873	0.1294545829296112	7.746
8	0.948499906885771	0.9653309240362812	0.953915135909169	0.9653309240362812	0.13299113512039185	7.5674
10	0.9488682087282262	0.9633645124716553	0.9544498117034105	0.9633645124716553	0.1366061121225357	7.5882

比较分析

数据集 A





guwenbert-base

guwenbert-base 在各方面表现较差。

Bert-ancient-Chinese

Bert-ancient-Chinese 在整个训练中效果最好。

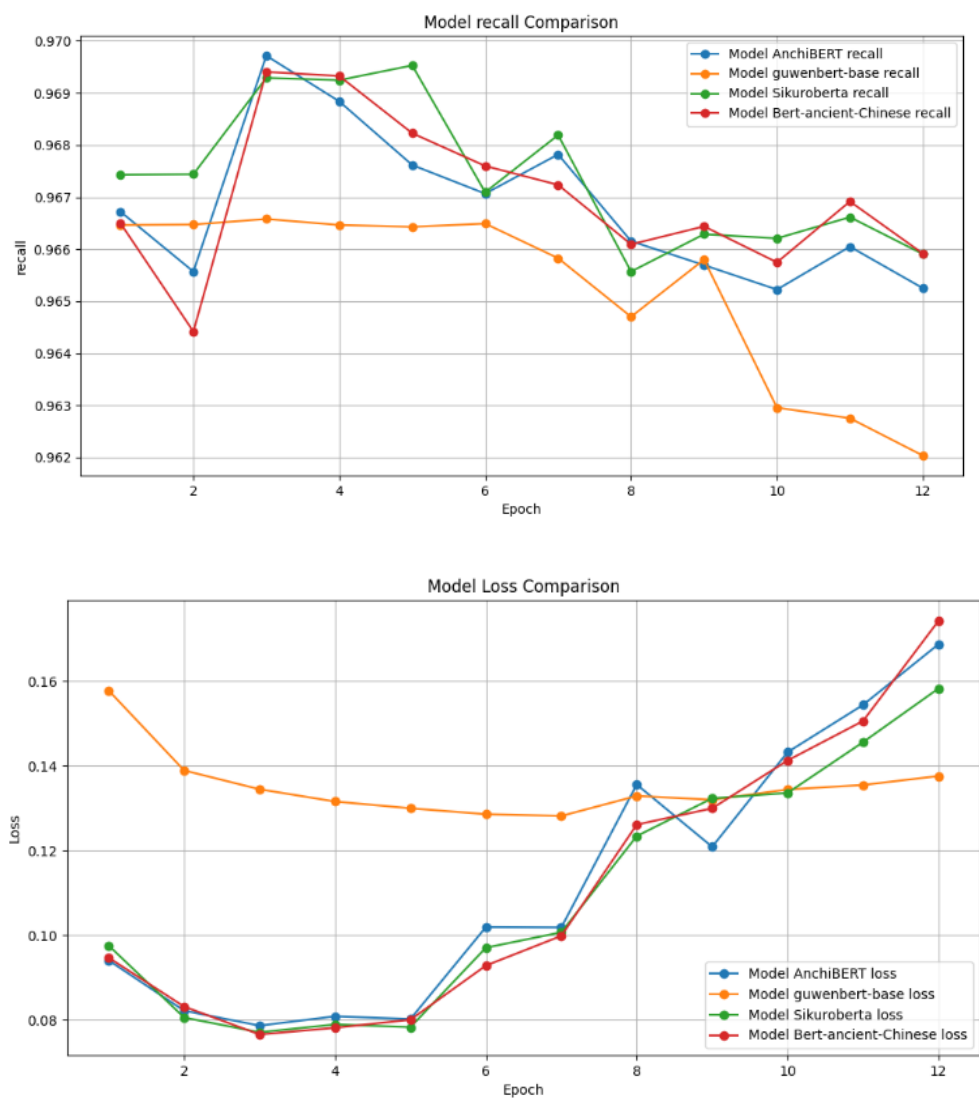
AnchiBERT

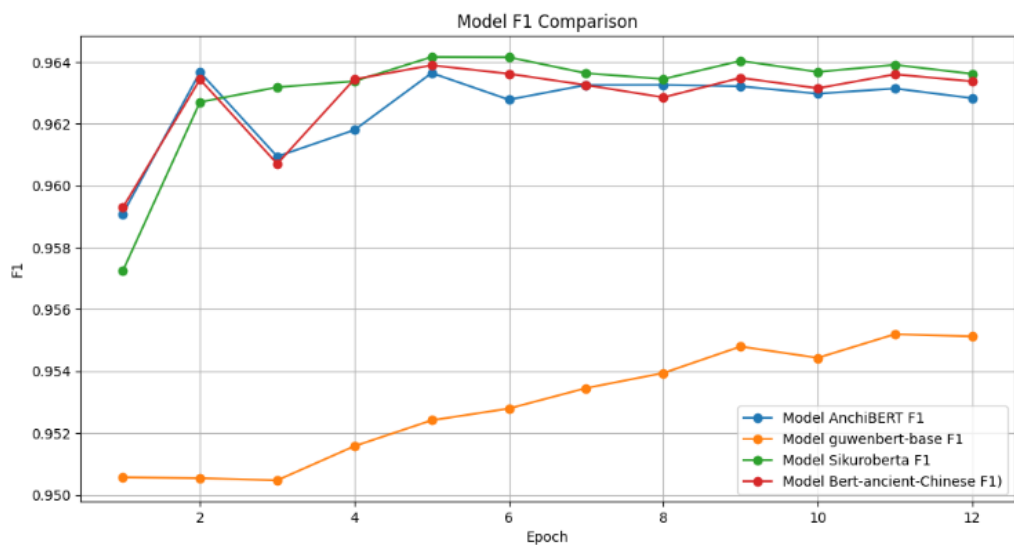
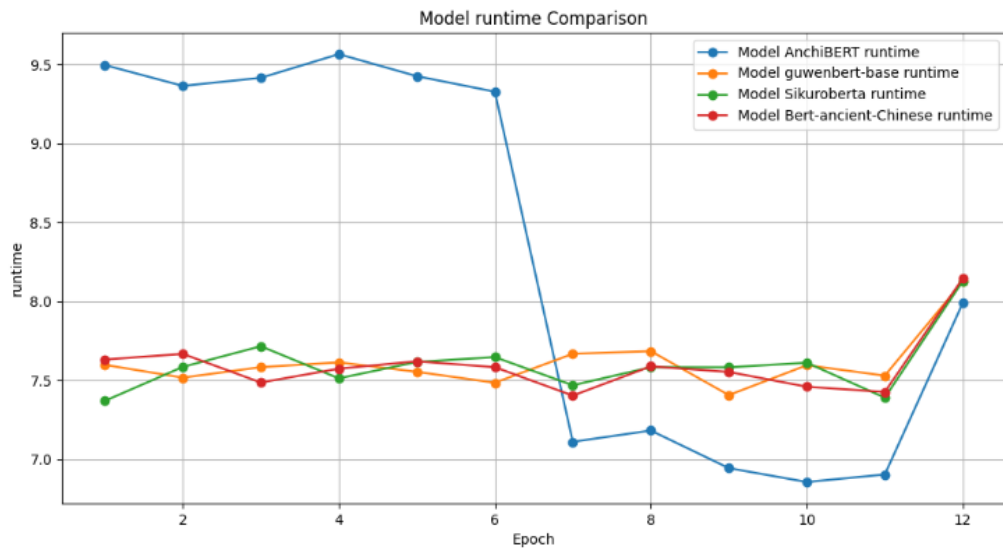
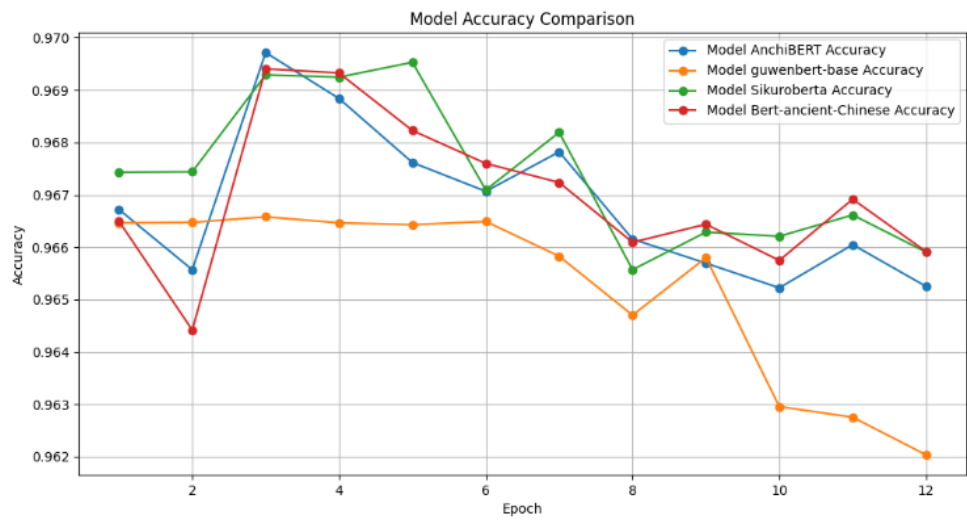
AnchiBERT 在整个训练中不如 SikuRoberta 和 Bert-ancient-Chinese，但是时间最短。

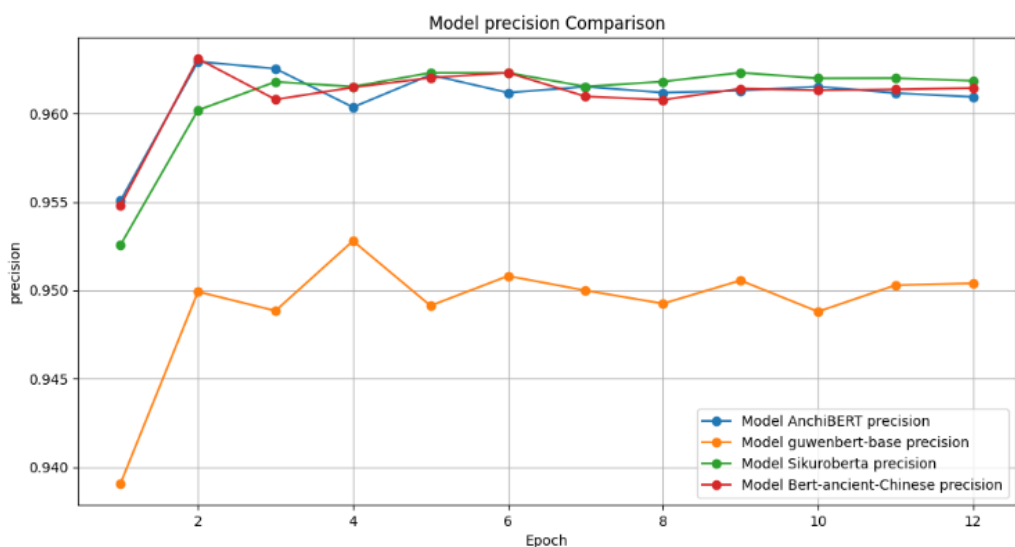
SikuRoberta

SikuRoberta 各项指标稍逊于 Bert-ancient-Chinese，训练效果相对优秀。

数据集 B







guwenbert-base

guwenbert-base 在各方面表现不如其他三个模型。

Bert-ancient-Chinese

Bert-ancient-Chinese 在 precision 和 f1 的表现与 AnchiBERT 和 SikuRoberta 类似，在后续 epoch 中相较于 SikuRoberta 更加平稳，表现相对不错。

AnchiBERT

anchbert 在其他方面表现不错，但在 runtime 中不太稳定。

SikuRoberta

SikuRoberta 在 precision 和 f1 的表现与 AnchiBERT 和 SikuRoberta 类似，表现相对不错。

综合比较

Bert-ancient-Chinese 和 SikuRoberta 表现都相对出色，Bert-ancient-Chinese 在后续 epoch 中更加平稳。

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	0.9717125382262997	0.9717125382262997	0.9717125382262997	0.8673289276579453	0.8478130526942544	0.8573962308701384
SikuR	0.9683438	0.9683438	0.9683438	0.8362269	0.8354067	0.8356550

obert a	45565749 3	45565749 3	45565749 3	47526488 6	66296890 5	520943984
Anchi BERT	0.9646167 81345565 7	0.9646167 81345565 7	0.9646167 81345565 7	0.8273717 70625562 1	0.8038283 10372927 2	0.8153108 407512697
guwe nbert -base	0.9162605 12232415 9	0.9162605 12232415 9	0.9162605 12232415 9	0.6532090 67995206 2	0.3968061 30054585	0.4690429 183394909 7

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	0.9629373 50478468 8	0.9629373 50478468 8	0.9629373 50478468 8	0.4463562 23356897 7	0.3469447 37913422 5	0.3858195 247664353 3
SikuRoberta	0.9640774 5215311	0.9640774 5215311	0.9640774 5215311	0.6133806 35777134 8	0.3840230 41465861 5	0.4477773 397719211 3
AnchiBERT	0.9632177 03349282 3	0.9632177 03349282 3	0.9632177 03349282 3	0.4179623 60384415 6	0.3378474 05702294 1	0.3702823 300980698
guwenbert-base	0.9600029 90430622	0.9600029 90430622	0.9600029 90430622	0.3253865 76678924	0.2004047 49296543	0.2373501 789184325

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，Bert-ancient-Chinese 在数据集 A 中的各项指标都是最高的，说明在整体上取得了最好的平衡性能。 在数据集 B 中指标相差不大。

Macro-Average 指标展示了模型对所有类别同等对待的性能，Bert-ancient-Chinese 在数据集 A 中的各项指标都是最高的，表明其对于少数类别的识别能力较强。 在数据集 B 中指标相差不大。