

# 历时演变视角下的古汉语分词：时期嵌入与大规模语料库的应用

柯永红

北京师范大学民俗典籍文字研究中心, 中国文字整理与规范研究中心 北京市新街口外大街19号100875  
yh5555@126.com

## 摘要

古汉语自动分词是古籍数字化和智能化处理的关键环节, 但古汉语在数千年演变过程中呈现出显著的历时性差异, 对构建通用的分词模型构成了严峻挑战。为应对这一挑战, 本研究构建了一个覆盖上古、中古及近代三个主要历史时期的大规模古汉语分词标注语料库, 在此基础上, 本文提出了一种基于时期嵌入 (Period Embedding) 的古汉语历时分词模型 ‘RoBERTa-PeriodEmb-Fusion-CRF’。该模型以预训练语言模型 ‘roberta-classical-chinese-large-char’ 为骨干, 通过引入可学习的时期向量来感知文本的时代背景, 并设计了非线性融合层以有效整合时期信息与上下文语义表示, 最后结合条件随机场 (CRF) 进行序列解码。在构建的历时语料库上的大量实验结果表明, 与不包含时期信息的强基线模型相比, 本文提出的模型在整体分词性能 (F1值达到0.9505) 以及跨时期文本的适应性上均取得了显著提升。本研究不仅验证了显式建模时期信息对于提升古汉语分词效能的重要性, 也为构建高性能、通用的古汉语处理工具提供了有益的思路和数据支持。

**关键词:** 古汉语分词; 历时演变; 时期嵌入; 大规模语料库

## Ancient Chinese Word Segmentation from a Diachronic Perspective: Applications of Period Embedding and a Large-Scale Corpus

Yonghong Ke

Research Center for Folklore, Classics and Chinese Characters,  
Research Center for Collation and Standardization of Chinese Characters  
No.19, Xinjiekouwai St, Haidian District, Beijing, 100875, P.R.China  
yh5555@126.com

## Abstract

Automatic word segmentation is a crucial step for the digitization and intelligent processing of ancient Chinese texts. However, ancient Chinese has undergone significant diachronic variations over millennia, posing severe challenges to building universal segmentation models. To address this issue, this study first presents the construction of a large-scale annotated corpus for ancient Chinese word segmentation, systematically covering three major historical periods: Ancient, Middle, and Early Modern Chinese, which provides a valuable data foundation for diachronic computational linguistic research. Based on this corpus, we propose a period-aware diachronic word segmentation model ‘RoBERTa-PeriodEmb-Fusion-CRF’ for ancient Chinese. The model leverages ‘roberta-classical-chinese-large-char’ as its backbone pre-trained language model. It introduces learnable Period Embeddings to perceive the temporal context of the text and employs a non-linear fusion layer to effectively integrate period information with

contextual semantic representations, followed by a Conditional Random Field (CRF) layer for sequence decoding. Extensive experiments conducted on our diachronic corpus demonstrate that the proposed model achieves significant improvements in overall segmentation performance (F1-score reaching 0.9505 ) and cross-period adaptability compared to strong baseline models awareness. This research not only validates the importance of explicitly modeling period information for enhancing ancient Chinese word segmentation but also offers valuable insights and data support for developing high-performance, universal tools for ancient Chinese language processing.

**Keywords:** Ancient Chinese Word Segmentation , Diachronic Variation , Period Embedding , Large-scale Corpus

## 1 引言

**分词(Word Segmentation)** 作为一项基础且关键的环节, 其准确性直接影响着后续诸如词性标注、命名实体识别、信息抽取、机器翻译乃至文本理解等高级任务的效果。然而, 与现代汉语相比, 古汉语分词面临着更为严峻和独特的挑战。**显著的历时性差异**是古汉语处理的核心难点之一。汉语从上古时期的文言发展到近代时期的白话, 其词汇系统、语法规则乃至文字形态均发生了深刻变化 (王力, 1980)。例如, 同一结构在短语向词转变的过程中存在词和短语两种形态。这种显著的历时差异使得为单一历史时期语料设计的模型难以普适于其他时期, 极大地限制了古汉语分词模型的**通用性和实用价值**。尽管已有部分研究关注古汉语分词, 但多数工作局限于特定断代或小规模语料, 难以系统性地应对古汉语的历时复杂性。

制约古汉语历时分词研究深入发展的一个关键瓶颈在于缺乏大规模、高质量、覆盖多个历史时期的标注语料库。现有公开的古汉语分词资源相对零散, 难以支撑对语言历时演变模式的有效学习和通用模型的稳健训练。为突破这一基础性障碍, 并为古汉语的历时计算语言学研究提供坚实的数据支撑, 本研究构建并细致标注了一个大规模的古汉语分词**历时语料库**。该语料库系统性地选取了自上古、中古至近代三个主要历史时期的代表性典籍, 力求为历时分词模型提供一个可靠的、可用于训练和评估的基础资源。

在构建这一大规模历时语料库的基础上, 本文进一步探索了如何利用显式的时期信息来提升古汉语分词模型的跨时期适应能力。我们认为, 让模型感知并利用文本所属的历史时期特征, 是解决历时性挑战的有效途径。为此, 本文提出了一种面向多时期古汉语文本的时期自适应分词模型。该模型以强大的预训练语言模型 `roberta-classical-chinese-large-char` (Yasuoka, 2022) 为基础编码器, 充分汲取其在古汉语文本上学习到的丰富上下文表示能力。核心在于引入了**时期嵌入 (Period Embedding)** 机制, 将文本的时期信息 (上古/中古/近代) 编码为低维稠密向量; 并设计了**非线性特征融合层 (Non-linear Feature Fusion Layer)**, 以更有效地整合上下文语义信息与时期特定的先验知识, 从而引导模型根据不同的时期背景动态调整其分词决策。最终, 模型采用条件随机场 (Conditional Random Field, CRF) 层进行全局最优序列解码。

本文工作的包括以下四个方面:

- 构建了一个大规模、覆盖上古、中古、近代三个关键历史阶段的古汉语分词标注语料库, 为古汉语历时研究提供了重要的基础数据资源。
- 提出了一种融合时期嵌入与非线性特征融合的古汉语分词模型 (RoBERTa-PeriodEmb-Fusion-CRF), 旨在显式地建模和利用文本的历时信息, 以提升模型的跨时期泛化能力。

©2025 中国计算语言学大会

根据《Creative Commons Attribution 4.0 International License》许可出版

本研究获得了以下项目的支持与资助: 国家社会科学基金项目“上古汉语词标记语料库及应用系统构建研究”(项目编号: 20BYY127), 教育部、国家语委研究基地项目“面向古籍智能化研究和应用的古籍分词语料库建设”(项目编号: ZDI145-86), 以及教育部哲学社会科学研究重大专项“古文字构形系统发展研究与数据库建设”(项目编号: 2024JZDZ048)。

- 在我们构建的历时标注数据集上进行了全面的实验。结果表明，与不包含时期信息的强基线模型相比，本文提出的方法在整体分词性能以及在不同历史时期文本上的表现均有提升，验证了所提方法的有效性。
- 通过细致的实验分析，包括分时期性能对比和案例研究，探讨了时期嵌入机制对模型识别历时语言现象的作用，为理解古汉语的计算建模提供了新的视角。

## 2 相关工作

### 2.1 通用中文分词技术

中文分词（CWS）技术发展经历了多个阶段。

早期使用基于词典和规则的方法。其中，最大匹配法（Maximum Matching, MM）、逆向最大匹配法（Reverse Maximum Matching, RMM）及其双向结合（Bi-directional MM, BMM）是最为经典的技术。这些方法依赖于预先构建的词典，通过贪心策略进行匹配，虽然实现简单，但在处理歧义和未登录词（OOV）方面存在天然的局限性。基于规则的方法试图引入语言学知识来消解歧义，但规则库的构建和维护成本高昂且难以穷尽。

随着机器学习的兴起，基于统计模型的方法，尤其是将分词视为序列标注任务的方法，取得了显著进展。隐马尔可夫模型（HMM）（Rabiner, 1989）在一定程度上提升了性能（Xue, 2003）。条件随机场（Conditional Random Field, CRF）（Lafferty et al., 2001）由于其判别式建模能力，能够利用丰富的上下文特征并克服HMM的观测独立性假设，成为后续很长一段时间内性能最优、应用最广的分词模型之一（Tseng, 2005）。

深度学习技术进一步推动了CWS的发展。循环神经网络（RNN），特别是长短期记忆网络（LSTM）（Hochreiter, 1997）（Chen, 2015）和门控循环单元（GRU），能够有效捕捉文本序列的长期依赖关系。BiLSTM-CRF 架构（Huang et al., 2015）（Lample et al., 2016）（Ma, 2018）结合了双向LSTM的强大上下文编码能力和CRF对输出标签序列的全局优化能力，在包括CWS在内的多种序列标注任务上取得了卓越表现。

近年来，基于Transformer架构的预训练语言模型（PLMs）（Devlin et al., 2019）彻底改变了NLP领域。BERT及其改进模型，如RoBERTa（Liu, 2019）、ERNIE（Sun, 2019）等，通过在海量文本上进行自监督学习，掌握了丰富的语言知识。在CWS任务中，当前最优的范式通常是把这些强大的PLMs作为编码器，在其上连接一个简单的线性分类层或CRF层，然后针对特定任务进行微调（Tian, 2020）。

### 2.2 古汉语分词技术

相较于现代汉语，古汉语自动分词的研究起步较晚，面临着数据稀疏、语言规范与现代汉语差异巨大等挑战。早期的古汉语分词研究也借鉴了现代汉语的处理方法。例如，郭辉等人（郭辉等, 2002）探索了改进的最大匹配法在古汉语文本上的应用。统计模型，特别是条件随机场（CRF），因其在序列标注任务上的优越性能，也被广泛应用于特定古籍或时期的分词任务。陈薇薇和俞士汶（陈薇薇等, 2007）较早地将CRF模型用于古代汉语的自动分词。后续研究者们针对不同时期的文献进行了探索，如梁社会和陈小荷（梁社会等, 2013）对先秦文献《孟子》的分词研究，王晓玉和李斌（王晓玉等, 2017）结合词典信息对中古汉语进行分词，陆文（陆文, 2018）则尝试了对《左传》进行分词。这些研究通常聚焦于某一特定历史时期或特定文体的文献，并取得了一定的进展。第一届古代汉语分词与词性标注国际评测（EvaHan）（李斌等, 2023）的举办，进一步推动了该领域技术的发展，并涌现了多种基于传统统计模型和深度学习的方法。

近年来，深度学习模型也开始应用于古汉语分词。俞敬松等人（俞敬松等, 2020）结合非参数贝叶斯模型和深度学习方法对古文进行分词。Tang等人（Tang, 2022）关注跨时代文本的序列切分问题，提出了一种带有切换记忆机制的方法。尽管取得了这些进展，现有研究的一个主要局限性仍然在于缺乏对汉语历时性系统演变的深入考量和统一建模。多数模型或是为单一时期设计并在小规模数据集上验证，或是在混合数据上进行通用训练而未能充分利用文本的时期信息。这导致其在面对时间跨度较大的古汉语文本时，性能和鲁棒性受到限制。此外，正如引言中所述，大规模、覆盖多时期的标准分词语料库的匮乏，也严重制约了相关研究的深入。

本研究在大规模历时分词语料库的基础上提出融合时期信息的模型，为构建通用的古汉语分词器提供支持。王力先生在《汉语史稿》中指出：“语言的发展既是渐进的，那么，由旧质到



新质的过渡阶段就应该是很长的，它可以历时几百年甚至上千年。”这深刻揭示了语言演变的连续性和渐进性特征，任何分期方案本质上都是一种基于特定标准的人为概括。目前，学术界（特别是汉语史研究领域）在宏观层面将古汉语划分为“上古”、“中古”、“近代”三个主要发展阶段，这一划分是基于对语音、词汇、语法等多个语言层面在这些阶段间发生显著系统性变化的长期观察和归纳，具有较高的概括性和通行度，也为我们理解和研究古汉语的宏观演变脉络提供了一个基础框架。

### 2.3 预训练语言模型在古汉语中的应用

预训练语言模型（PLMs）的出现是自然语言处理领域的一大突破，它们通过在海量无标注文本上进行自监督学习，获得了强大的语言表示能力。针对古汉语独特的语言特性和丰富的文献资源，学术界和工业界相继开发了一系列专门的古汉语PLMs。其中具有代表性的包括：哈工大讯飞联合实验室发布的GuwenBERT (高嗣佳等, 2021); 浙江大学等单位基于四库全书等大规模古籍文献训练的SikuBERT (Fei et al., 2021); 以及由社区贡献者（如Jihuai-wpy）在Hugging Face平台上发布的bert-ancient-chinese<sup>1</sup>等模型，它们为研究者提供了便捷的古汉语PLM资源。本文所采用的基础编码器roberta-classical-chinese-large-char<sup>2</sup>是一个基于RoBERTa架构并在古典中文文本上进行了预训练的字符级大模型，其设计考虑了古汉语的特点，如繁简字处理。

这些古汉语PLMs为各类下游任务（包括分词、词性标注、命名实体识别、古文今译等）提供了坚实的语义表示基础，并在相关评测中展现出优于传统方法的性能。然而，这些预训练模型本身通常并不直接包含或区分细致的历时信息。因此，如何将这些强大的PLMs与显式的时期知识相结合，以提升其在处理具有历时演变特征的古汉语文本时的表现，是本研究关注的一个重要方面。

### 2.4 处理数据异质性与时期演变的相关方法

在自然语言处理中，如何有效处理来自不同来源、不同风格或不同时间段文本所表现出的数据异质性，是一个持续受到关注的核心问题。针对历时语言变化的研究，计算语言学领域也进行了诸多有益的探索。例如，在词义演变检测（Semantic Change Detection）任务中，研究者广泛采用的方法是比较词嵌入在不同时期语料库上的分布差异，以追踪词义的动态漂移 (Kutuzov et al., 2018)。在文本分类、情感分析或机器翻译等任务中，当面临不同领域 (Domain) 或风格 (Style) 的数据时，引入领域嵌入 (Domain Embedding) (Daumé III, 2007) 或使用多任务学习 (Multi-Task Learning, MTL) 框架 (Caruana, 1997) 来学习共享表示和特定表示，是提升模型泛化能力的常见策略。近年来，Adapter模块 (Houlsby et al., 2019) 作为一种参数高效的迁移学习方法，允许在不修改预训练模型主体参数的情况下，为特定任务或领域注入新知识，也为处理数据异质性提供了新的思路。

具体到序列标注任务中融入条件信息（如时期、领域或文体），**时期嵌入 (Period Embedding)** 或更广义的**条件嵌入 (Conditional Embedding)** 是一种直观且被证明有效的方法。其核心思想是将离散的条件变量（例如，文本所属的“上古”、“中古”、“近代”时期）映射为可学习的低维稠密向量，并将这些向量作为额外的特征信息输入到神经网络模型中。这样，模型就能够感知到当前的条件上下文，并据此调整其内部参数或表示，从而对不同条件下的数据做出更具适应性的预测。例如，在神经机器翻译中，研究者曾使用领域标记 (Domain Tags) 的嵌入来指导模型生成更符合特定领域风格的译文 (Kobus et al., 2017)。

本文借鉴了这一思想，将时期嵌入机制应用于古汉语分词任务，系统性地探索其在捕捉和适应古汉语显著历时演变特征方面的具体效用。与简单地为每个时期分别训练独立模型或者在混合数据上训练单一通用模型相比，本文旨在构建一个**统一的、能够根据明确的时期输入动态调整其分词策略的自适应模型**，从而在保证模型通用性的同时，提升其在具体历史时期文本上的分词精度。

## 3 方法

为有效处理古汉语文本的历时性差异，并构建一个能够适应多时期语言特征的通用分词模型，本文提出了一种融合时期嵌入的自适应序列标注方法。模型整体遵循主流的“编码器-解码

<sup>1</sup><https://huggingface.co/Jihuai/bert-ancient-chinese>

<sup>2</sup><https://huggingface.co/KoichiYasuoka/roberta-classical-chinese-large-char>

器”框架，如图1所示。其核心组件包括：基于预训练语言模型的基础编码器、用于感知时期信息的时期嵌入模块、用于整合多源特征的非线性融合层，以及用于全局序列解码的条件随机场（CRF）层。

### 3.1 基础编码器：RoBERTa

我们选用强大的预训练语言模型‘roberta-classical-chinese-large-char’作为模型的基础编码器。该模型基于RoBERTa架构 (Liu, 2019)，并在包含大量古汉语文本的数据上进行了预训练，使其能够有效捕捉古汉语的字词特征和深层上下文语义信息。对于输入的古汉语字符序列  $X = (c_1, c_2, \dots, c_n)$ ，经过RoBERTa编码后，我们得到其对应的上下文相关的隐藏状态序列  $H = (h_1, h_2, \dots, h_n)$ ，其中  $h_i \in \mathbb{R}^{d_{roberta}}$  是第  $i$  个字符的表示向量， $d_{roberta}$  是RoBERTa模型的隐藏层维度。这个序列  $H$  将作为后续特征融合的基础。为了增强模型的泛化能力并减少过拟合，我们在RoBERTa的输出后添加了一个Dropout层。

### 3.2 时期信息感知模块

为了让模型能够感知并适应不同历史时期（本文设定为上古、中古、近代三个时期）的语言特征，我们引入了时期嵌入（Period Embedding）机制。该机制的核心是将离散的时期类别信息转化为模型可以利用的、可学习的连续向量表示。

首先，我们将每个历史时期映射为一个唯一的整数ID，例如：上古（sg） $\rightarrow 0$ ，中古（zg） $\rightarrow 1$ ，近代（jd） $\rightarrow 2$ 。对于一个给定的输入文本，其所属的时期ID  $p \in \{0, 1, \dots, N_p - 1\}$ （其中  $N_p = 3$  为时期总数）作为本模块的输入。

然后，我们采用标准的嵌入层（`torch.nn.Embedding`）定义一个可学习的时期嵌入矩阵  $M_{period} \in \mathbb{R}^{N_p \times d_{period}}$ ，其中时期嵌入维度  $d_{period}$  在本研究中设为64。通过查表操作，可以获得时期ID  $p$  对应的时期嵌入向量  $e_p = M_{period}[p]$ 。采用可学习的嵌入层主要基于以下考虑：其一，它能够使模型根据下游分词任务自动学习和优化时期表示；其二，该嵌入层本身的参数量可控（本研究中为  $3 \times 64 = 192$  个参数），是为模型提供宏观时期上下文信息的一种轻量级且有效的机制。

为了将时期信息融入到序列的每个字符表示中，我们将学习到的时期嵌入向量  $e_p$  扩展（Expand）至与输入字符序列等长，得到时期特征序列  $E'_p = (e_p, e_p, \dots, e_p)$ 。该时期特征序列  $E'_p$  将作为一种补充性的信息源，与RoBERTa编码器输出的上下文语义表示  $H$  在后续的非线性融合层中进行交互，旨在使模型能够根据不同的时期背景动态地调整其分词决策。

### 3.3 非线性特征融合层

简单地将不同来源的特征进行拼接可能不足以让模型充分学习它们之间的复杂交互。为了更有效地整合RoBERTa编码器提供的上下文语义信息  $H$  和时期信息感知模块提供的时期特征  $E'_p$ ，我们设计了一个非线性特征融合层。

首先，我们将上下文表示  $h_i$  和对应位置的时期特征  $e_p$ （来自  $E'_p$ ）进行拼接（Concatenation）操作，得到初始的融合特征向量  $h'_i$ ：

$$h'_i = \text{concat}(h_i, e_p) \quad (1)$$

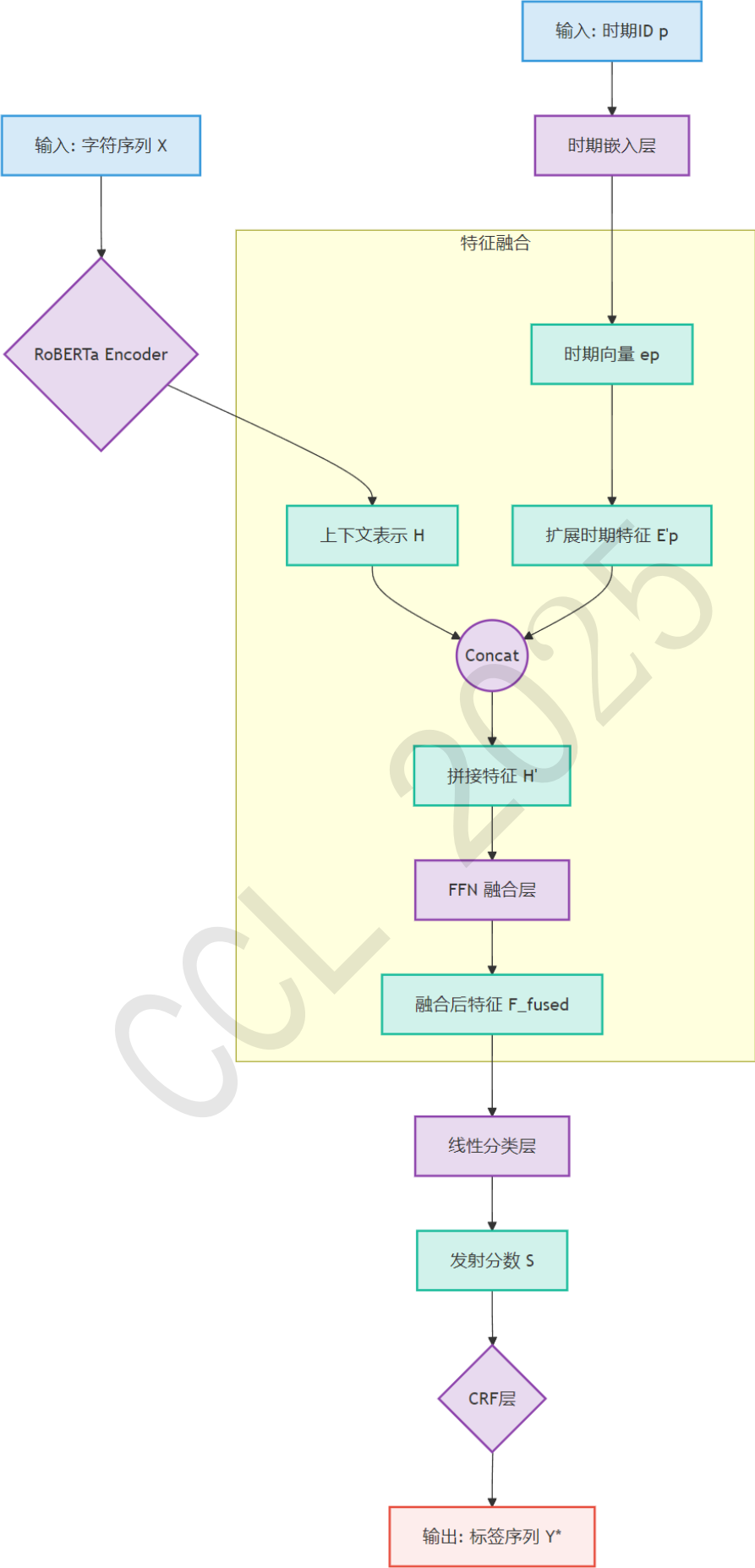
其中  $h'_i \in \mathbb{R}^{d_{roberta} + d_{period}}$ 。

随后，我们将拼接后的特征序列  $H' = (h'_1, h'_2, \dots, h'_n)$  输入到一个前馈神经网络（Feed-Forward Network, FFN）中进行深度的非线性变换和融合。该FFN包含两层线性变换，并使用GELU作为激活函数，同时加入了Dropout和Layer Normalization以增强模型的稳定性和泛化能力。具体计算如下：

$$F^{(1)} = \text{LayerNorm}(\text{GELU}(\text{Linear}_1(H'))) \quad (2)$$

$$F_{fused} = \text{Linear}_2(\text{Dropout}(F^{(1)})) \quad (3)$$

其中， $\text{Linear}_1$  将输入维度从  $d_{roberta} + d_{period}$  映射到一个中间维度  $d_{ffn}$ ， $\text{Linear}_2$  再将维度从  $d_{ffn}$  映射回一个合适的输出维度  $d_{fused}$ （例如，可以设为  $d_{roberta}$  或一个更小的维度，如256或512，以便后续CRF层处理）。 $F_{fused} = (f_1, f_2, \dots, f_n)$  即为最终融合后的特征序列，其中  $f_i \in \mathbb{R}^{d_{fused}}$ 。这种非线性融合方式使得模型能够以更复杂的方式学习上下文信息和时期信息的交互关系，而不是简单地将它们视为独立的附加特征。



第二十四届中国计算语言学大会论文集, 第651页-第665页, 济南, 中国, 2025年8月11日至14日。  
Figure 1: 本文提出的时期自适应古汉语分词模型框架图。  
(c) 2025 中国中文信息学会计算语言学专业委员会

### 3.4 解码层与损失函数

融合后的特征序列 $F_{fused}$ 随后被送入一个线性分类层，以计算每个字符对应各个分词标签（如B, I, E, S, O<sup>3</sup>）的发射分数（Emission Scores）。对于序列中的第 $i$ 个字符的融合特征 $f_i$ ，其发射分数向量 $S_i \in \mathbb{R}^{N_{label}}$ 计算如下：

$$S_i = W_{class}f_i + b_{class} \quad (4)$$

其中 $W_{class} \in \mathbb{R}^{N_{label} \times d_{fused}}$ 和 $b_{class} \in \mathbb{R}^{N_{label}}$ 是线性分类层的可学习参数， $N_{label}$ 是标签的总数（本文中为5）。

为了建模输出标签之间的依赖关系（例如，“B”标签后面更可能跟随“I”或“E”标签），我们采用了条件随机场（Conditional Random Field, CRF）层进行全局序列解码。CRF层引入一个状态转移矩阵 $T \in \mathbb{R}^{N_{label} \times N_{label}}$ ，其中 $T_{jk}$ 表示从标签 $j$ 转移到标签 $k$ 的分数。对于一个给定的标签序列 $Y = (y_1, y_2, \dots, y_n)$ ，其总分定义为发射分数和转移分数之和：

$$\text{score}(X, Y) = \sum_{i=1}^n S_{i, y_i} + \sum_{i=1}^{n-1} T_{y_i, y_{i+1}} \quad (5)$$

在训练阶段，模型通过最大化真实标签序列的对数似然来进行优化。损失函数定义为负对数似然损失：

$$L_{CRF} = -\log P(Y_{true}|X) = -\text{score}(X, Y_{true}) + \log \sum_{Y' \in \mathcal{Y}(X)} e^{\text{score}(X, Y')} \quad (6)$$

其中 $Y_{true}$ 是真实的标签序列， $\mathcal{Y}(X)$ 是输入 $X$ 所有可能的标签序列集合。这个计算可以通过前向-后向算法高效完成。在预测阶段，我们使用维特比算法（Viterbi Algorithm）来寻找具有最高分数的标签序列 $Y^*$ 作为最终的分词结果。

为了稳定训练初期，我们采用了逐步引入CRF损失的策略。在训练的前 $k$ 个轮次（epoch），模型仅使用标准的交叉熵损失（Cross-Entropy Loss）在上述线性分类层的输出上进行优化（忽略标签间的转移）。从第 $k+1$ 轮开始，逐渐增加CRF损失的权重，最终完全使用CRF损失。

## 4 实验

为了全面评估本文提出的融合时期嵌入的古汉语分词模型的有效性，我们进行了一系列实验。本节将详细介绍所使用的实验数据集、评估指标、具体的实验设置以及用于对比的基线模型。

### 4.1 数据集

正如引言中所述，缺乏大规模、覆盖多时期的标注语料是古汉语历时分词研究的关键瓶颈。为解决这一问题，并为本研究提供坚实的数据基础，我们构建了一个大规模历时古汉语分词语料库(DHACWS Corpus)。<sup>4</sup>

#### 4.1.1 时期划分与数据格式

我们选取了能够代表不同历史阶段语言特点的文献，并根据王力先生《汉语史稿》（王力, 1980）等的划分标准，并结合文献的成书年代，将语料划分为三个主要历史时期：

- **上古汉语**：包括《论语》《左传》《诗经》等53部典籍。
- **中古汉语**：包括《抱朴子内篇》《世说新语》《生经》等70部典籍。

<sup>3</sup>本文采用BIESO标注体系，其中B、I、E、S分别表示词语的开始、中间、结束和单字成词，O表示非词语成分（如标点符号）。

<sup>4</sup>本文所使用的历时分词语料库主要基于以下科研项目的支持与成果：国家社会科学基金项目“上古汉语词标记语料库及应用系统构建研究”（项目编号：20BYY127）和教育部、国家语委研究基地型项目“面向古籍智能化研究和应用的古籍分词语料库建设”（项目编号：ZDI145-86）。此外，语料库亦整合了部分公开获取的古汉语资源。关于本语料库（特别是上古部分）的构建理念、方法论和初步成果，可参见“上古汉语分词和词性标记语料的构建”（柯永红, 2024）。



- **近代汉语**: 包括《敦煌变文》《红楼梦》《祖堂集》等19部典籍。

数据集采用BIESO标注体系进行分词标注。其中，B表示词语的开始字，I表示词语的中间字，E表示词语的结束字，S表示单字成词，O表示非词语成分（主要用于标记标点符号）。

处理后的数据保存为jsonlines 格式，每行一个JSON对象，包含三个字段：

- **text**: 原始古汉语文本字符串。
- **labels**: 与text 字符一一对应的BIESO标签列表。
- **period**: 文本所属时期的标识符（"sg", "zg", 或"jd"）。

标注完成的数据如：{"text": "三月，葬蔡平公。", "labels": ["B", "E", "O", "S", "B", "I", "E", "O"], "period": "sg"}

#### 4.1.2 数据统计与划分

经过整理和标注，最终构建的历时古汉语分词语料库的详细统计信息如表1所示。该表展示了整个语料库以及各个时期子集在句子数、字符数、词数方面的规模。

Table 1: 数据集统计信息（训练集/验证集/测试集）				
时期	指标	训练集	验证集	测试集
上古(sg)	句子数	181123	10027	10175
	字符数	3493973	194311	196192
	词数	3101838	172231	174075
中古(zg)	句子数	198780	11090	10931
	字符数	2984421	164012	164037
	词数	2387958	132634	131285
近代(jd)	句子数	241776	13421	13432
	字符数	4952117	275043	273244
	词数	3983623	221348	219917
总体(Overall)	句子数	621679	34538	34538
	字符数	11430511	633366	633473
	词数	9473419	526213	219917

考虑到古汉语封闭语料、数据稀疏的特点，我们将整个语料库按照9:0.5:0.5的比例，随机划分为训练集、验证集和测试集，以保证各数据集中不同时期样本的分布与总体分布一致，从而进行更可靠的模型训练和评估。各划分集的详细统计信息也包含在表1中。

#### 4.2 评估指标

我们采用词级别的精确率（Precision, P）、召回率（Recall, R）和F1值（F1-score）作为评估模型分词性能的主要指标。计算方式基于模型预测的分词结果与人工标注的黄金标准之间的匹配程度。为了更全面地考察模型的性能，我们将从以下几个层面进行评估：

- **整体性能(Overall Performance)**: 在整个测试集上计算P/R/F1。
- **分时期性能(Period-specific Performance)**: 分别在测试集的上古、中古、近代三个子集上独立计算P/R/F1，以考察模型的跨时期适应性。

#### 4.3 实验设置

本文所有实验均基于PyTorch (Paszke et al., 2019)和Hugging Face Transformers (Wolf et al., 2020)库实现。基础编码器采用‘roberta-classical-chinese-large-char’的预训练权重。关键的超参数设置如表2所示。其他未在表中列出的超参数均采用相关库的默认设置或根据验证集上的表现进行调整。



Table 2: 主要超参数设置

超参数(Hyperparameter)	值(Value)
预训练模型(Pre-trained Model)	‘roberta-classical-chinese-large-char‘
最大序列长度(Max Sequence Length)	128
时期嵌入维度(Period embedding dim)	64
RoBERTa输出后Dropout率	继承自模型配置
自定义层Dropout率	继承自模型配置
训练批次大小(Train Batch Size)	24
评估批次大小(Evaluation Batch Size)	48
梯度累积步数(Gradient Accumulation Steps)	4
学习率(Learning Rate)	1e-5
优化器(Optimizer)	AdamW
训练轮数(Number of Epochs)	20

4.4 对比模型

为了验证本文提出的时期自适应分词模型的有效性，我们设置了以下对比模型：

- **RoBERTa-CRF (基线模型)**：该模型直接使用roberta-classical-chinese-large-char 作为编码器，其输出的序列表示直接送入一个线性分类层和CRF层进行分词。该模型不包含任何显式的时期信息或额外的特征融合机制，代表了当前基于强预训练模型的标准分词方法。

本文提出的模型记为**RoBERTa-PeriodEmb-Fusion-CRF**。所有模型均在相同的训练集、验证集和测试集上进行训练和评估，以保证比较的公平性。

5 结果与分析

本节将详细呈现并分析上一节所述实验的评估结果。我们将首先比较本文提出的时期自适应分词模型RoBERTa-PeriodEmb-Fusion-CRF 与基线模型在整体性能上的表现。随后，我们将考察模型在不同历史时期文本上的适应能力，并分析其对未登录词（OOV）的处理效果。最后，通过消融实验和案例分析进一步验证模型各组件的有效性以及时期嵌入的作用机制。

5.1 整体性能比较

为了评估模型的整体分词能力，我们在整个测试集上对RoBERTa-PeriodEmb-Fusion-CRF <sup>(5)</sup> 以及所有对比模型进行了测试。表3展示了各模型在词级别（Word-level）的精确率（P）、召回率（R）和F1值。

Table 3: 不同模型在整体测试集上的分词性能比较（词级别F1值）

模型(Model)	精确率(P)	召回率(R)	F1值(F1)	提升(F1 Δ)
RoBERTa-CRF (基线)	0.8975	0.8735	0.8853	-
<b>OurModel (本文模型)</b>	<b>0.9478</b>	<b>0.9532</b>	<b>0.9505</b>	↑ +7.36%

从表3可以看出，**OurModel**在整体测试集上的F1值达到了0.9505，相比于基线模型RoBERTa-CRF的0.8853，取得了7.36%的相对提升。这一结果初步表明，通过引入时期嵌入和非线性特征融合机制，能够有效提升古汉语分词模型的整体性能。

5.2 分时期性能分析

为了更深入地考察模型在处理不同历史时期文本时的适应能力，我们在测试集的上古（sg）、中古（zg）和近代（jd）三个子集上分别对各模型进行了评估。详细的词级别F1值如表4所示。

<sup>5</sup>下文统一使用**OurModel**指代本文提出的模型。

Table 4: 不同模型在各时期子测试集上的词级别F1值

模型(Model)	上古(sg) F1	中古(zg) F1	近代(jd) F1
RoBERTa-CRF (基线)	0.8794	0.8882	0.9072
<b>OurModel (本文模型)</b>	<b>0.9147</b> (↑ 4.01%)	<b>0.9542</b> (↑ 7.43%)	<b>0.9680</b> (↑ 6.70%)

从表4的对比结果可以观察到以下几点:

- **基线模型的时期差异:** RoBERTa-CRF基线模型在三个时期上的表现存在一定差异, 例如, 在近代汉语上F1值较高, 而在上古汉语上相对较低, 这反映了不同时期语言特征对单一模型的挑战。
- **本文模型的跨时期提升:** 本文提出的RoBERTa-PeriodEmb-Fusion-CRF模型在所有三个时期的子测试集上均取得了优于基线模型的性能。特别是在中古汉语, F1值从基线的0.8882提升至0.9542, 提升了7.43%个百分点。这充分证明了时期嵌入和非线性融合机制在帮助模型感知和适应不同历史时期语言特征方面的有效性, 显著增强了模型的跨时期泛化能力和鲁棒性。
- **时期间的性能趋势:** 尽管所有时期都有提升, 但中古和近代汉语的F1值是提升更为显著, 这可能与中古和近代汉语更接近现代书面语、分词语料更规范或骨干模型等因素有关。

这些结果支持了我们的核心假设: 显式地将时期信息融入模型是能够有效提升古汉语历时分词性能。

### 5.3 词长性能分析

为了进一步探究模型的性能特点, 我们分析了本文提出的**OurModel**在不同长度词语上的识别效果。表5展示了模型在测试集上对1字词、2字词直至更长词语的精确率、召回率和F1值。

Table 5: OurModel在不同长度词语上的性能分析

词长(Word Length)	F1值(F1)	精确率(P)	召回率(R)
1字词(1-char)	0.9877	0.9918	0.9837
2字词(2-char)	0.9511	0.9473	0.9549
3字词(3-char)	0.9412	0.9416	0.9408
4字词(4-char)	0.8698	0.8562	0.8839
5字词(5-char)	0.8691	0.8492	0.8900
6字词(6-char)	0.8211	0.8298	0.8125
7字词(7-char)	0.7375	0.9077	0.6211
≥8字词(≥8-char)	0.8235	0.8750	0.7778

从表中数据可以看出, 模型对单字词和双字词的识别准确率非常高, F1值分别达到了0.9877和0.9511。随着词语长度的增加, 模型的性能呈现出一定的下降趋势。例如, 对于4字词和5字词, F1值分别降至0.8698和0.8691。对于更长的词语(如7字词), F1值进一步下降至0.7375, 其中召回率(0.6211)相对精确率(0.9077)有更明显的降低, 这可能表明模型在识别这些较长且出现频率较低的完整词语边界方面存在一定的困难, 或者倾向于将它们切分的更细。值得注意的是, 对于≥8字的极长词, 虽然样本数少, 但模型仍然取得了一个相对不错的F1值0.8235, 这可能得益于预训练模型的长距离依赖捕捉能力以及CRF层的全局优化。总体而言, 词长分析结果符合预期, 即模型对常见长度的词语(尤其是1-3字词)处理效果优秀, 而对于稀有的长词, 其识别仍具挑战性, 这可能与训练数据中长词样本相对不足有关。

5.4 消融实验分析

为了验证模型中不同组件的有效性，我们进行了消融实验。主要考察去掉时期嵌入模块（退化为RoBERTa-CRF基线）对模型性能的影响。实验结果如表6所示（该表可以与表3合并，或单独展示，这里假设单独展示）。

Table 6: 模型组件消融实验结果（整体测试集词级别F1值）

模型配置(Model Configuration)	F1值(F1)
<b>OurModel (完整模型)</b>	<b>0.9505</b>
- 时期嵌入模块(退化为基线)	0.8853

消融实验结果清晰地展示了时期嵌入的有效性：当从完整模型中移除时期嵌入模块后（即与基线模型RoBERTa-CRF对比），模型性能从0.9505下降至0.8853，下降了7.36个百分点。这再次证明了时期嵌入对于提升模型性能的作用。

5.5 案例分析

为了更直观地展示模型的优势，我们选择了一些典型的分词难例进行分析。我们以双字序列“如今”为例，该序列在不同时期的主要用法（作为单个双字词或两个单字词“如”和“今”）存在显著差异。我们首先统计了“如今”在训练集中不同时期的出现模式，结果如表7所示。

Table 7: “如今”在训练数据中的出现频次统计

双字序列	近代(jd)		中古(zg)		上古(sg)	
	“如今”词	“如”+“今”*	“如今”词	“如”+“今”*	“如今”词	“如”+“今”*
如今	2864	2	52	20	0	6

“如”+“今”指训练数据中“如”和“今”被标注为两个独立词（通常是单字词S）且紧邻出现的次数；“如今”词指“如今”被标注为一个双字词（B-E）的次数。

表7的统计数据显示，“如今”作为一个固定的双字词主要出现在中古和近代文本中，而在上古文本中，它更倾向于被理解为“如”和“今”两个独立的字。

接下来，我们展示了我们的模型（已融合时期嵌入）在包含“如今”序列的测试集样本上的分词结果示例：

- 上古(sg):
  - 预测文本: 后生可畏，焉知来者之不如今也？（《论语·子罕》）
  - 分词结果: 后/ 生/ 可/ 畏/ ， / 焉/ 知/ 来者/ 之/ 不/ 如/ 今/ 也/ ？ （《论语·子罕》）
- 中古(zg):
  - 预测文本: 如今现在十方诸佛，亦为受佛职位诸菩萨说。（《悲华经·卷第一》）
  - 分词结果: 如今/ 现在/ 十方/ 诸/ 佛/ ， / 亦/ 为/ 受/ 佛/ 职/ 位/ 诸/ 菩萨/ 说/ 。 （《悲华经·卷第一》）
- 近代(jd):
  - 预测文本: 如今且说那邢皋门的行止。（《醒世姻缘传·第十六回》）
  - 分词结果: 如今/ 且/ 说/ 那/ 邢皋门/ 的/ 行止/ 。 （《醒世姻缘传·第十六回》）

上述示例清晰地展示了模型根据不同的时期背景对同一字符序列“如今”做出了不同的分词决策。这种区分能力证明了**时期嵌入机制**的有效性。模型在训练过程中，通过学习将输入的时期信息与文本的上下文表示相结合，从而捕捉到了特定词语或结构在不同历史时期的不同用法和组词规律。在预测时，相应的时期嵌入为模型提供了关键的判别信号，引导其做出符合时代背景的分词决策，进而提高了模型处理古汉语历时性差异的能力和整体分词的准确性。

## 6 结论与展望

古汉语自动分词是古籍文献资源深度利用与智能化处理的关键环节。然而，古汉语在漫长的历史发展过程中经历了显著的历时演变，给构建通用、高效的分词模型带来了巨大挑战。针对这一问题，本文基于涵盖了上古、中古及近代三个主要历史时期的代表性文献的分词语料库，提出了一种面向多时期古汉语文本的时期自适应分词模型。该模型以强大的预训练语言模型roberta-classical-chinese-large-char为基础，引入时期嵌入机制来捕获文本的时代背景信息，并通过一个非线性特征融合层有效整合时期特征与深层上下文语义表示，最后结合CRF层进行全局最优解码。

我们在构建的历时分词语料库上进行了一系列全面的实验评估。实验结果充分证明了本文所提方法的有效性：

- 与不包含时期信息的强基线模型相比，融入时期嵌入和非线性融合机制的模型在整体分词性能上取得了明显提升，F1值达到了0.9505。
- 更为重要的是，本文模型在不同历史时期的子测试集上均表现出优于基线模型的性能，证明了时期信息对于提升模型跨时期泛化能力的作用。

通过消融实验和案例分析，我们进一步验证了时期嵌入模块和非线性融合层对模型性能的积极贡献，并直观地展示了模型如何利用时期信息处理具有历时演变特征的语言现象。

本文采用静态时期嵌入与“上古/中古/近代”三时期划分，是在综合考量了当前古汉语标注资源的实际状况、大规模语料标注的可行性以及模型实现的复杂度后作出的选择。作为一项工程性的探索，本研究初步但有力地证实了将宏观时期信息融入先进的预训练语言模型，在提升古汉语自动分词处理跨时期文本的效用方面，展现出可行性与应用潜力。

综上所述，本研究通过构建基础语料资源和提出时期自适应模型，为解决古汉语分词的历时性挑战提供了一个有效的解决方案。研究结果不仅验证了显式建模时期信息对于提升古汉语NLP任务性能的重要性，也为未来构建更智能、更通用的古汉语处理工具积累了有益的经验。

展望未来，本研究仍有进一步拓展和深化的空间：

- **更丰富的历时特征融合：**除了时期ID嵌入，未来可以尝试融入更多能够表征时代特色的语言学特征，如特定句式的出现、特定虚词的用法等。
- **与其他知识的协同：**探索如何将时期信息与词典知识、句法知识等其他外部知识更有效地协同融合，以期在解决OOV问题和复杂歧义方面取得更大突破。
- **探索更先进的条件化模型架构：**例如，可以尝试使用Adapter模块或更复杂的多任务学习框架来更参数高效或更具针对性地融入时期信息。
- **下游任务的应用与评估：**将本文提出的历时分词模型应用于古汉语的词性标注、命名实体识别、信息抽取等下游任务，评估其带来的连锁效益。
- **集成预测：**基于现有数据，训练模型对未知文本尝试所有已知时期的预测，给未知语料打上“上古/中古/近代”标签。

## 参考文献

- 陈薇薇, 俞士汶. 基于条件随机场的古代汉语自动分词研究[J]. 中文信息学报, 2007, 21(6): 7-12.
- 高嘉琦, 赵庆聪. 基于新词发现的古典文学作品分词方法研究[J]. 计算机技术与发展, 2021, 31(9): 178-181.
- 高嗣佳, 赵子帅, 黑兹涵, 李泊林, 李鹏, & 吴辉. (2021). GuwenBERT: 古文预训练语言模型. In 第二十届中国计算语言学大会论文集(CCL 2021) (pp. 1040-1051).
- 郭辉, 苏中义, 王文, 等. 一种改进的MM分词算法[J]. 微型电脑应用, 2002(1): 13-15.
- 柯永红. 上古汉语分词和词性标记语料的构建[J]. 数字人文, 2024(4).



- 李斌, 袁义国, 芦靖雅, 等. 第一届古代汉语分词和词性标注国际评测[J]. 中文信息学报, 2023, 37(3): 46-53+64.
- 梁社会, 陈小荷. 先秦文献《孟子》自动分词方法研究[J]. 南京师范大学文学院学报, 2013(3): 175-182.
- 陆文. 基于条件随机场的《左传》自动分词研究[D]. 南京: 南京农业大学, 2018.
- 钱智勇, 周建忠, 童国平, 等. 基于HMM的楚辞自动分词标注研究[J]. 图书情报工作, 2014, 58(4): 105-110.
- 邱冰, 皇甫娟. 基于中文信息处理的古代汉语分词研究[J]. 微计算机信息, 2008, 24(24): 100-102.
- 唐雪梅, 苏祺, 王军, 等. 基于图卷积神经网络的古汉语分词研究[J]. 情报学报, 2023, 42(6): 740-750.
- 王嘉灵. 以《汉书》为例的中古汉语自动分词[D]. 南京: 南京师范大学, 2014.
- 王力. 汉语史稿[M]. 北京: 中华书局, 1980: 35.
- 王晓玉, 李斌. 基于CRFs和词典信息的中古汉语自动分词[J]. 数据分析与知识发现, 2017, 1(5): 62-70.
- 邢付贵, 朱廷劭. 基于大规模语料库的古文词典构建及分词技术研究[J]. 中文信息学报, 2021, 35(7): 41-46.
- 严顺. 基于CRF的古汉语分词标注模型研究[J]. 江苏科技信息, 2016(8): 10-12.
- 俞敬松, 魏一, 张永伟, 等. 基于非参数贝叶斯模型和深度学习的古文分词研究[J]. 中文信息学报, 2020, 34(6): 1-8.
- Caruana, R. (1997). Multitask Learning. *Machine Learning*, 28(1), 41-75.
- X. Chen, X. Qiu, C. Zhu, P. Liu, and X. Huang. "Long Short-Term Memory Neural Networks for Chinese Word Segmentation." In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 1197-1206, Lisbon, Portugal, Sep. 2015. Association for Computational Linguistics. [Online]. Available: <https://aclanthology.org/D15-1141/>
- Chu, C., Dabre, R., & Kurohashi, S. (2017). An Empirical Study of Domain Adaptation for Neural Machine Translation. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics (EACL 2017), Volume 2: Short Papers* (pp. 533-538). Valencia, Spain: Association for Computational Linguistics.
- Daumé III, H. (2007). Frustratingly Easy Domain Adaptation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics (ACL 2007)* (pp. 256-263). Prague, Czech Republic: Association for Computational Linguistics.
- Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*, Volume 1 (Long and Short Papers) (pp. 4171-4186).
- Fei, H., Li, Y., Wang, S., Li, R., Chen, Y., & Wu, H. (2021). SikuBERT: A Pre-trained Language Model for Ancient Chinese. *arXiv preprint arXiv:2109.07185*.
- S. Gao, Z. Zhao, Z. He, B. Li, P. Li, and H. Wu. "GuwenBERT: 古文预训练语言模型." In *Proceedings of the Twentieth China National Conference on Computational Linguistics (CCL '21)*, pp. 1040-1051, 2021.
- T. Gui, Y. Chen, Q. Zhang, Q. Zhu, and X. Huang. "CNN-based chinese NER with lexicon rethinking." In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI '19)*, pp. 4986-4992, 2019.
- S. Hochreiter and J. Schmidhuber. "Long short-term memory." *Neural Computation*, vol. 9, no. 8, pp. 1735-1780, 1997.
- Houlsby, N., Giurugu, A., Jastrzebski, S., Morrone, B., De Laroussilhe, Q., Gesmundo, A., Attariyan, M., & Gelly, S. (2019). Parameter-Efficient Transfer Learning for NLP. In *Proceedings of the 36th International Conference on Machine Learning (ICML)* (pp. 2790-2799).
- Huang, Z., Xu, W., & Yu, K. (2015). Bidirectional LSTM-CRF models for sequence tagging. *arXiv preprint arXiv:1508.01991*.

- Kobus, C., Crego, J., & Senellart, J. (2017). Domain Control for Neural Machine Translation. In *Proceedings of the Second Conference on Machine Translation (WMT17), Volume 1: Research Papers* (pp. 388–394). Copenhagen, Denmark: Association for Computational Linguistics.
- Kutuzov, A., Arefyev, N., Kuzmenko, E., & Pivovarov, L. (2018). Diachronic Word Embeddings and Semantic Shifts: A Survey. In *Proceedings of the 27th International Conference on Computational Linguistics (COLING)* (pp. 1384–1397).
- Lafferty, J., McCallum, A., & Pereira, F. C. N. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning (ICML '01)* (pp. 282–289).
- Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., & Dyer, C. (2016). Neural architectures for named entity recognition. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL HLT)* (pp. 260–270).
- X. Li, Y. Li, X. Yang, H. Chen, L. Nie, Z. Chen, and C. Sun. "FLAT: Chinese NER using flat-lattice transformer." In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL '20)*, pp. 1529–1539, 2020.
- Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov. "RoBERTa: A robustly optimized BERT pretraining approach." *arXiv preprint arXiv:1907.11692*, 2019.
- I. Loshchilov and F. Hutter. "Decoupled Weight Decay Regularization." In *International Conference on Learning Representations (ICLR)*, 2019. [Online]. Available: <https://arxiv.org/abs/1711.05101>.
- J. Ma, K. Ganchev, and D. Weiss. "State-of-the-art Chinese Word Segmentation with Bi-LSTMs." In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 1256–1261, Brussels, Belgium, 2018. Association for Computational Linguistics. [Online]. Available: <https://aclanthology.org/D18-1460/>
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., & Chintala, S. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, & R. Garnett (Eds.), *Advances in Neural Information Processing Systems 32 (NeurIPS 2019)* (pp. 8024–8035). Curran Associates, Inc. Retrieved from <http://papers.neurips.cc/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html>
- L. Rabiner. "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition." *Proceedings of the IEEE*, vol. 77, no. 2, pp. 257–286, Feb. 1989. doi: 10.1109/5.18626.
- Y. Sun, S. Wang, Y. Li, S. Feng, X. Chen, H. Zhang, X. Tian, D. Zhu, H. Tian, and H. Wu. "ERNIE: Enhanced representation through knowledge integration." In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence (AAAI '19)*, pp. 2304–2311, 2019.
- X. M. Tang and Q. Su. "That Sleep Al the Nyght with Open Ye! Cross-era Sequence Segmentation with Switch-memory." In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7830–7840, Dublin, Ireland, 2022. Association for Computational Linguistics.
- Y. Tian, Y. Song, F. Xia, T. Zhang, and Y. Wang. "Improving Chinese Word Segmentation with Wordhood Memory Networks." In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7634–7644, Online, 2020. Association for Computational Linguistics. [Online]. Available: <https://aclanthology.org/2020.acl-main.686/>
- H. Tseng, P. Chang, G. Andrew, D. Jurafsky, and C. Manning. "A Conditional Random Field Word Segmenter for Sighan Bakeoff 2005." In *Proceedings of the Fourth SIGHAN Workshop on Chinese Language Processing*, pp. 168–171, 2005. [Online]. Available: <https://aclanthology.org/I05-3027/>
- Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Davison, J., Shleifer, S., von Platen, P., Ma, C., Jernite, Y., Plu, J., Xu, C., Le Scao, T., Gugger, S., Drame, M., Lhoest, Q., & Rush, A. M. (2020). Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural*

- Language Processing: System Demonstrations (EMNLP 2020 Demo)* (pp. 38–45). Online: Association for Computational Linguistics. <https://www.aclweb.org/anthology/2020.emnlp-demos.6>. doi:10.18653/v1/2020.emnlp-demos.6
- N. Xue. "Chinese Word Segmentation as Character Tagging." *International Journal of Computational Linguistics & Chinese Language Processing*, vol. 8, no. 1, pp. 29–48, Feb. 2003. [Online]. Available: <https://aclanthology.org/003-4002/>
- J. Yang, Y. Zhang, and S. Wu. "Subword Encoding in Lattice LSTM for Chinese Word Segmentation." In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL HLT '19), Volume 1 (Long and Short Papers)*, pp. 556–561, Minneapolis, Minnesota, 2019. Association for Computational Linguistics. [Online]. Available: <https://aclanthology.org/N19-1053/>
- Yasuoka, K. (2022). *roberta-classical-chinese-large-char*. Hugging Face Model Repository. Retrieved from <https://huggingface.co/KoichiYasuoka/roberta-classical-chinese-large-char>
- Y. Zhang and J. Yang. "Chinese NER using lattice LSTM." In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1123–1133, 2018.