# Hierarchical Ranking Neural Network for Long Document Readability Assessment

Yurui Zheng[a], Yijun Chen[a], Shaohong Zhang[a,*]

*[a]School of Computer Science and Cyber Engineering, Guangzhou University, Guangzhou, 510000, Guangdong, China*

## Abstract

Readability assessment aims to evaluate the reading difficulty of a text. In recent years, while deep learning technology has been gradually applied to readability assessment, most approaches fail to consider either the length of the text or the ordinal relationship of readability labels. This paper proposes a bidirectional readability assessment mechanism that captures contextual information to identify regions with rich semantic information in the text, thereby predicting the readability level of individual sentences. These sentence-level labels are then used to assist in predicting the overall readability level of the document. Additionally, a pairwise sorting algorithm is introduced to model the ordinal relationship between readability levels through label subtraction. Experimental results on Chinese and English datasets demonstrate that the proposed model achieves competitive performance and outperforms other baseline models.

*Keywords:*
Long document, Multidimensional context weights, Ranking model

## 1. Introduction

Automatic Text Readability (ARA) research originated in the early 20th century, aiming to evaluate text reading difficulty and assist educators in recommending appropriate reading materials for learners [1]. Readability assessment approaches are generally classified into three paradigms: human evaluation, co-selection-based analysis, and content-based analysis. Human evaluation involves expert annotation or reader surveys; co-selection methods leverage user interaction data such as reading time or choices [2]; and content-based approaches infer readability using linguistic, syntactic, or semantic features extracted from the text itself. Early studies predominantly relied on experts' subjective evaluations and simple statistical features, such as sentence length and word complexity. However, these approaches suffered from high subjectivity, with evaluation results often varying depending on the evaluators' criteria and purposes. The current mainstream evaluation is based on content-based analysis. In the field of natural language processing (NLP), readability assessment research was initially limited. Early researchers developed readability formulas by designing simple linguistic features [3, 4, 5]. As the number of readability formulas grew, the application of readability analysis expanded from education to diverse industries, including law [6], medicine [7, 8], and government policy to improve document intelligibility.

Interest in this issue among NLP researchers has only emerged in the past two decades. From statistical language models and feature-engineering-based machine learning methods to more recent deep neural networks, a range of methods has been explored for this task [9]. Existing works [10, 11, 12, 13] primarily involve designing diverse language features and employing machine learning models to text classification, such as support vector machines(SVM) or multilayer perceptron.

In recent years, the proliferation of deep neural networks has advanced text readability research, circumventing the need for cumbersome linguistic feature engineering. However, deep neural networks have not outperformed traditional machine learning models by a wide margin. One key challenge is the highly variable text lengths in readability

---

*Corresponding author
*Email addresses:* 846006629@qq.com (Yurui Zheng), 1051575227@qq.com (Yijun Chen), zimzsh@qq.com (Shaohong Zhang )

evaluation datasets [14]—upper-grade texts are typically longer, while lower-grade texts are shorter, a discrepancy that is particularly pronounced in Chinese corpora. Most current models [15, 16, 17, 18] employ pre-trained BERT [19] as word encoders, but BERT's maximum input length is limited to 512 tokens. To address this, the hierarchical attention network (HAN) framework [20] offers a solution: its word-to-sentence-to-document hierarchical structure enables models to capture long-text dependencies more effectively.

Additionally, while the text readability task is commonly framed as a text classification problem, Zheng and Zheng [21] highlights the critical role of context vectors in traditional attention mechanisms. However, prior studies [15] either neglect context vectors or initialize them randomly, a practice that undermines contextual modeling. To address this, we introduce a multi-dimensional context weight vector, which generates sentence-level representation vectors by aggregating context weights and word embeddings. Notably, each text is assigned a unique context weight vector to capture its global semantic characteristics.

Moreover, this paper innovatively proposes a bidirectional text readability assessment framework based on hierarchical modeling. The proposed method first employs document-level readability annotations to infer sentence-level readability labels via a hierarchical model, thereby constructing a sentence-level readability corpus. Subsequently, this sentence corpus is used as auxiliary supervision to enhance the model's forward prediction of document-level readability. This approach not only enables the modeling of implicit sentence-level readability features, but also enhances the model's contextual understanding by leveraging hierarchical information, leading to a more fine-grained readability assessment system.

Due to the orderliness brought by the text difficulty level labels, there are also some works that treat it as a regression task [22] or a ranking task [23]. Because texts in adjacent grades tend to be more similar than distant grades, but do not distinguish them in the classification task. Therefore, we construct the Ranking Model in the tail of the model, which adopts the way of pairwise comparison and label subtraction to learn the order relationship between categories.

Our research contributions can thus be concluded as follows:

- Compared with English, research on Chinese text readability remains relatively limited. We constructed a Chinese corpus based on textbooks from mainland China and developed a feature engineering scheme tailored to the Chinese language domain.

- We designed a hierarchical model inspired by hierarchical neural architectures, which preserves more information from long documents. Meanwhile, we introduced multi-dimensional contextual weighting to guide the attention mechanism in identifying informative words within the input sequence.

- We proposed a bidirectional readability assessment framework, which utilizes sentence-level readability labels to further enhance the model's ability to predict text-level readability in the forward direction.

- At the end of the forward prediction module, we introduced a ranking model that learns the ordinal relations between readability levels by modeling label differences, and outputs the optimal text ranking through pairwise comparisons.

The rest of this paper is structured as follows:

- Section 2 provides a review of related work.

- Section 3 introduces our proposed model.

- Sections 4 and 5 describe the datasets and experimental setup, respectively.

- Section 6 discusses the experimental results.

- Section 7 concludes the paper and outlines future directions.

## 2. Related Work

Early research focused on developing the readability formula, which is a weighted linear function, including Dale-Chall [4], SMOG [5], and Flesch-Kincaid [24]. Using readability formula to evaluate text difficulty is objective and easy to calculate, but it only considers simple text features, and does not consider language features well. Nonetheless, the traditional readability formulation also laid the foundation for later readability research.

In machine learning, researchers train ARA models with the help of classifiers by constructing a large number of linguistic features. Heilman et al. [25] used a combination of lexical features and grammatical features that are derived from subtrees of syntactic parses, while also verifying that ordinal regression models were most effective in predicting reading difficulty. Feng et al. [26] employed SVM and logistic regression to compare and evaluate several sets of explanatory variables - including shallow, language modeling, POS, syntactic, and discourse features, and checked that the judicious combination of various features led to significant improvements over the state of the art. Hancke et al. [27] developed new morphological features and achieved 89.7% accuracy in German readability classification based on these features. Qiu et al. [28] designed 100 factors to systematically evaluate the influence of four levels of linguistic features (shallow, part of speech, syntax and discourse) on the difficulty of predicting texts for L1 Chinese learners, and further selected 22 features with regression significance. Deutsch et al. [29] and Lee et al. [11] similarly leverage various language-driven features combined with simple machine learning models and aided by deep learning models to improve performance.

Deep learning methods are becoming more and more widely used in ARA. Jiang et al. [30] proposed a graph-based readability evaluation method using word coupling, which combines the merits of word frequencies and text features. Azpiazu and Pera [22] present a multiattentive recurrent neural network architecture for automatic multi-lingual readability assessment. This architecture considers raw words as its main input, but internally captures text structure and informs its word attention process using other syntax- and morphology-related datapoints, known to be of great importance to readability. Blaneck et al. [31] studied the ability of ensembles of fine-tuned GBERT and GPT-2-Wechsel models to reliably predict the readability of German sentences.

ARA can similarly be viewed as a regression or ordinal regression task due to the orderliness of the readability labels. Meng et al. [32] proposed a new comprehensive framework that uses a hierarchical self-attention model to analyze document readability. In this model, the goal is to minimize the ordinal regression loss [33]. Lee and Vajjala [23] proposed the first neural pairwise ranking model for ARA and showed the first results of cross-lingual, zero-shot evaluation of ARA using neural models. Z Zeng et al. [34] used soft labels [35] to exploit the ordinal nature of the readability assessment task.In addition to the above label modeling-based methods, Tanaka-Ishii et al. [36] proposed a readability assessment framework based on a sorting mechanism. Instead of directly predicting the absolute readability level, a binary comparator is trained to judge the relative readability between any two texts, and then the text set is sorted by a sorting algorithm. The biggest advantage of this method is that the training data only needs two levels (easy and difficult), which greatly reduces the difficulty of annotation. It is particularly suitable for low-resource language environments where training data is scarce and level annotation is difficult.

In experiments across corpora, Xia et al. [37] applied a generalization method to adapt models trained on larger native corpora to estimate text readability for learners, and explored domain adaptation and self-learning techniques to make use of the native data to improve system performance on the limited L2 data. Madrazo Azpiazu and Pera [38] developed a cross-lingual readability assessment strategy that serves as a means to empirically explore the potential advantages of employing a single strategy (and set of features) for readability assessment in different languages, including interlanguage prediction agreement and prediction accuracy improvement for low-resource languages.

## 3. Methodology

First, we built a set of linguistic features for the Chinese corpus. Then, we proposed an ARA hierarchical model that can be used to evaluate the readability level of sentences in long documents, and introduced a multidimensional context weight vector and a multi-head difficulty embedding matrix(MDEM) into the model, aiming to solve the information loss problem in the traditional attention mechanism and how to reversely predict the readability level of sentences through text. We call this model HHNN-MDEM (Hierarchical Hybrid Neural Network with Multi-Head Embedding Matrix). The overall structure of the model is shown in Fig. 1.
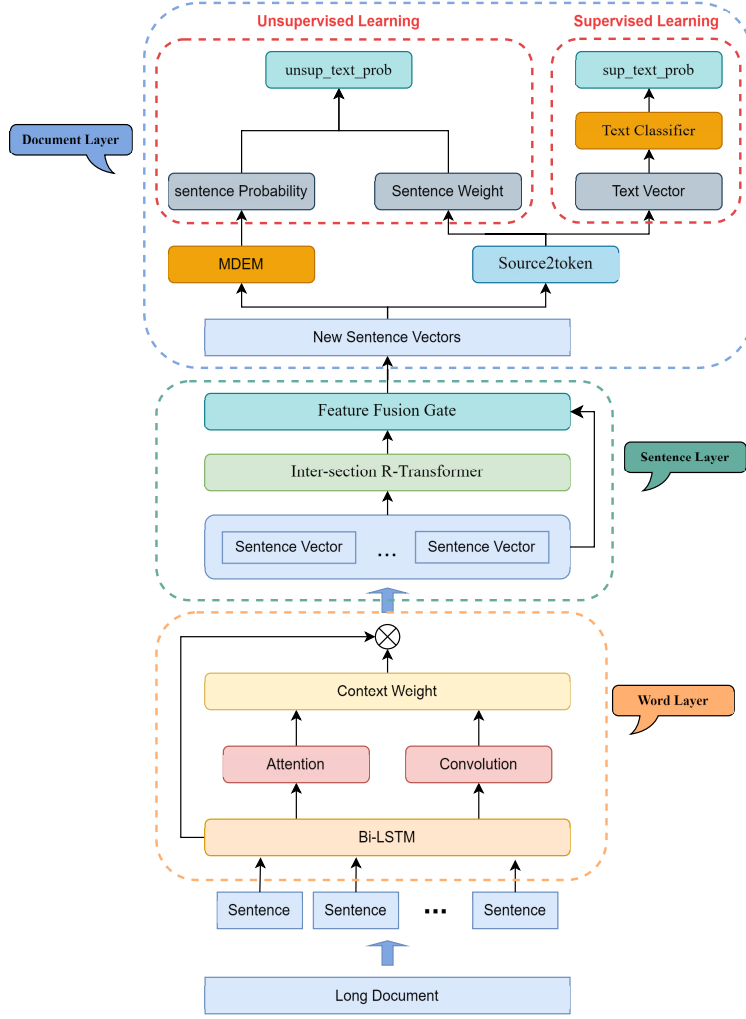
Figure 1: The overall structure of the HHNN-MDEM.

## 3.1. Explicit Features

Compared with English, there are few researches on Chinese readability feature engineering. Therefore, we construct Chinese traditional features, including lexical features, part-of-speech features, discourse features and article cohesion features, with reference to relevant traditional features [28, 39, 40]. In the specific domain of classical Chinese, we construct thematic features and complex semantic features for this purpose. See the appendix for a detailed explanation.

For traditional features in English, explicit feature extraction comes from lingfeat [11], which studies 255 language features. For already existing features, variants were added to expand coverage, including the development of high-level semantic features related to the topic: semantic richness, clarity, and noise. In order to improve the generalization of the experiment, the features related to the dataset were removed in this experiment, that is, the topic knowledge features of WeeBit and OneStopEnglish, a total of 32 features. Since there are 16 Entity Grid Features missing from some of the features generated by lingfeat, 207 explicit features were finally generated.

## 3.2. Proposed Model

As shown in Figure. 1, HHNN-MDEM is a hierarchical semi-supervised learning framework that uses a hybrid supervision and hierarchical consistency training strategy to achieve automatic labeling of sentence readability labels.
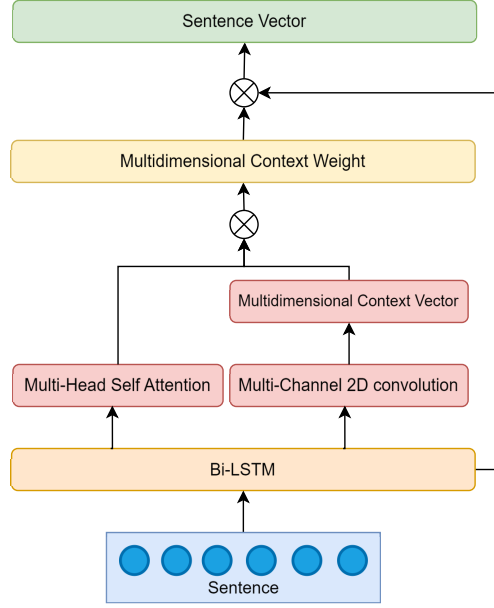
Figure 2: Overall structure of word layer.

The hierarchical learning architecture can be divided into three components, namely the word layer, sentence layer, and document layer.

### 3.2.1. Word Layer

In order for the model to learn more text information from long documents, we divide the text into multiple sentences. Specifically, for a given long document, we represent it as a sequence of sentences $S = \{s_1, \ldots, s_n\}$, where $s_i = \{w_{i,1}, \ldots, w_{i,m}\}$, $i = 1, \ldots, n$ denotes the token sequence of the $i^{th}$ sentence. If necessary, each sentence will be padded and cut to maintain the same length. The Embedding layer encodes each word $w_{i,j}$ into a d-dimensional vector based on word embeddings. The output is an $m \times d$ dimensional matrix $A^i = [\mathbf{e}_1, \ldots, \mathbf{e}_m]$, where $d$ is the embedding dimension and $m$ is the number of words in the $i^{th}$ sentence.

Due to the orderliness of text sequences, we use a bidirectional recurrent neural network, such as LSTM or GRU, to capture the order information between words. Specifically, the word vectors of $A^i$ are sequentially input into the Bi-LSTM. By connecting the forward hidden state and the backward hidden state, a new word-level representation $h_i^t$ is produced, denoted as $\mathbf{h}_i^t = [h_{i,1}^t, \ldots, h_{i,m}^t] \in \mathbb{R}^{m \times d}$, i.e., $\mathbf{h}_i^t = Bi - LSTM(A^i)$.

**Multidimensional context weights.** The multi-head self-attention mechanism is used to learn the interaction between words as the attention coefficient. Then, a multi-channel two-dimensional convolutional network is used to learn the multi-dimensional context vector of the sentence, and the attention coefficient and the multi-dimensional context vector are multiplied as the final multi-dimensional context weight of the word. The representation vector of each sentence is formed by aggregating the multi-dimensional context weights and the words. The overall structure is shown in Fig. 2.

Du et al. [41] proved that the convolution of CNN is precisely the process of calculating the similarity between text and the *attentivesearch templates*. Therefore, the CNN layer extracts the most influential n-gram syntax of different semantic aspects from the text and uses them as context vectors. According to the output of the Bi-LSTM, the local features on $h_i^t$ are extracted using convolutional kernels. The convolution operation involves a set of $k$ convolutional kernels, where each convolutional kernel $Conv \in \mathbb{R}^{l \times d}$ is applied to a window of $l$ words, generating new features $c_j$ from the window of the vector $h_{i,j:j+l-1}^t$, with $j = 1, \ldots, m$ representing the $j^{th}$ word of the $i^{th}$ sentence, as follows,

$$c_j = Conv \cdot h_{i,j:j+l-1}^t + b \tag{1}$$

where is $b$ bias. This convolution kernel is applied to each possible window of the matrix $h^t$ to generate the context

5

vector $\hat{\mathbf{c}} = [c_1, \ldots, c_{m-l+1}]$. By employing k convolutional kernels, the context vectors are concatenated together, resulting in a k-dimensional context vector $C = [\hat{\mathbf{c}}_1, \ldots, \hat{\mathbf{c}}_k] \in \mathbb{R}^{(m-l+1) \times k}$.

Multi-head attention (MHA) [42] allows the model to jointly attend to information from different representational subspaces at different positions. The multi-head self-attention mechanism (MHSA) is a special case of MHA, where the queries ($Q$), keys ($K$), and values ($V$) of the self-attention layer all come from the output of the previous encoder layer, i.e., the inputs are $Q = K = V$. Here, we modify this mechanism by taking $Q$ and $K$ from the output $h^t$, and by calculating the similarity between $Q$ and $K$, we learn the interactions between words to obtain the attention coefficients. The values $V$, on the other hand, are taken from the multidimensional contextual vector $C$. Thus, the multidimensional contextual weights $W_i$ are calculated using the weighted sum of the attention coefficients and the multidimensional contextual vector $C$. The specific calculation is shown as follows,

$$W_i^t = softmax(ReLU(MHSA(\mathbf{h}_i^t, \mathbf{h}_i^t, C))) \tag{2}$$

where $softmax$ refers to the along-column normalization and $ReLU$ is the modified linear unit activation function. Finally, by taking the weighted sum of $h_i^t$ and $W_i^t$, the sentence representation vector $h_i^s$ is obtained, as represented below,

$$\mathbf{h_i^s} = \sum_{j=1}^{m} W_{i,j}^t \odot h_{i,j}^t \tag{3}$$

where $\odot$ represents the element wise product of two matrices. We denote $h^s = [\mathbf{h}_1^s, \ldots, \mathbf{h}_n^s] \in \mathbb{R}^{n \times d}$. According to the multi-dimensional context vector generated by the convolutional network, the interactive word vector of multi-head self-attention is combined with the context vector to generate multi-dimensional context weights. The weights select useful local features from the recurrent layer. This hybrid network retains the advantages of the three models, and each text has its own context vector.

### 3.2.2. Sentence Layer

In order to further capture the long-term dependency of different sentences, we adopt the Inter-section R-Transformer [43], which adopts N layer of transformers and replaces the residual blocks with residual fusion gates. The multi-head self-attention layer MHSA and normalization are firstly applied to the sentence level features as follows,

$$o^s = norm(MHSA(h^s)) \tag{4}$$

where $o^s = [o_1^s, \ldots, o_n^s] \in \mathbb{R}^{n \times d}$, and $norm$ represents the normalization operation.

To combine the local and global context features at the sentence level, the Inter-section R-Transformer uses a residual fusion gate to dynamically merge the input and output of multi-head self-attention. The output sequence of the residual fusion gate $e^s = e_1^s, \ldots, e_n^s$ is computed as follows,

$$G1 = sigmoid(W_{11}^s o^s + W_{12}^s h^s + b_1^s) \tag{5}$$

$$e^s = G1 \odot h^s + (1 - G1) \odot o^s \tag{6}$$

where $W_{11}^s$, $W_{12}^s$ and $b_1^s$ are the parameters of the first residual fusion gate and $sigmoid$ is the activation function. The output of gate $e^s$ is then passed through a fully connected layer $f^s$ followed with a normalization process. The final output of the Inter-section R-Transformer is obtained by applying another residual fusion gate, which can be expressed as,

$$G2 = sigmoid(W_{21}^s norm(f^s(e^s)) + W_{22}^s e^s + b_2^s) \tag{7}$$

$$v^s = G2 \odot e^s + (1 - G2) \odot norm(f^s(e^s)) \tag{8}$$

where $W^{21}$, $W^{22}$, and $b_2^s$ are the parameters of the second residual fusion gate.

We employ feature fusion gate [43] to fuse $h^s$ and $v^s$ to generate the final feature representation of the sentence. The formula of the feature fusion gate is as follows,

$$F = ReLU(W_3^s[h^s, v^s] + b_3^s) \tag{9}$$

$$G = sigmoid(W_4^s[h^s, v^s] + b_4^s) \tag{10}$$

$$u^s = G \odot F + (1 - G) \odot h^s \tag{11}$$

where $W_3^s \in \mathbb{R}^{2d \times d}$, $W_4^s \in \mathbb{R}^{2d \times d}$, $b_3^s \in \mathbb{R}^d$, $b_4^s \in \mathbb{R}^d$ are the learning parameters of the feature fusion gate. $u^s = [u_1^s, \ldots, u_n^s] \in \mathbb{R}^{n \times d}$ is the output of the feature fusion gate.

### 3.2.3. Document Layer

We apply the sentence vector $u^s$ to the "source2token" self-attention [44], explore the dependencies between sentences, and compress $u^s$ into a document vector $d$, which is represented as follows,

$$W^D = softmax(W_1^d ReLU(W_2^d u^s + b_1^d) + b_2^d) \tag{12}$$

$$d = \sum_{i=1}^{n} W_i^D \odot u_i^s \tag{13}$$

where $W_1^d \in \mathbb{R}^{d \times d}$, $W_2^d \in \mathbb{R}^{d \times d}$, $b_1^d \in \mathbb{R}^d$, $b_2^d \in \mathbb{R}^d$ are the learning parameters of the "source2token" self-attention module.

**Supervised Learning.** The objective of the supervised component is to leverage document-level labels as explicit supervision signals to train a text classifier that learns document-level feature representations. This enables the effective transfer of document-level label information to the sentence level, driving the model to learn semantic mappings from documents to sentences and laying a solid foundation for the subsequent automatic generation of sentence-level labels. Specifically, the input text is first encoded by the HHNN model to obtain a document representation $d$. Then, $d$ is passed through a fully connected layer to map it to a probability distribution **r** over readability levels. The supervised training objective is to minimize the cross-entropy loss, as defined by the following formula,

$$\mathcal{L}_{\text{sup}} = -\frac{1}{N} \sum_{i=1}^{N} \sum_{k=1}^{Y} y_{ik} \log(r_{ik}) \tag{14}$$

where $N$ denotes the number of samples, $y_{ik}$ represents the ground-truth label of sample $i$ for class $k$, and $r_{ik}$ denotes the corresponding predicted probability. The supervised component not only trains a document-level classifier but also ensures that the model, optimized through cross-entropy loss, captures explicit semantic information from the document. This helps align the sentence representation space with the semantic structure of readability levels.

**Unsupervised Learning.** The core objective of the unsupervised learning component is to learn sentence-level readability labels through hierarchical propagation of document-level labels. To achieve this, we introduce the Multi-Head Difficulty Embedding Matrix (MDEM) module, which, combined with hierarchical consistency training, enables the automatic generation of sentence-level labels.

The MDEM module learns the relative difficulty scores of sentence representations with respect to different readability levels. At its core, it employs a multi-head attention mechanism to compute the difficulty scores of each sentence under each readability category, thereby producing a difficulty distribution for each sentence. Specifically, the MDEM module performs a matrix multiplication between the sentence representations $u^s \in \mathbb{R}^{n \times d}$ and the multi-head difficulty embedding matrix $M \in \mathbb{R}^{h \times z \times Y}$, where each head learns category-specific difficulty features. The sentence vectors $u^s$ are first projected into multi-head format, as formulated below,

$$u^{s'} = \text{reshape}(u^s) \in \mathbb{R}^{h \times n \times z} \tag{15}$$

here, $h$ denotes the number of attention heads, $n$ is the number of sentences, and $z$ represents the dimension of each input feature, where $z = d/h$. Next, the sentence representations are multiplied with MDEM to obtain the sentence-level scores under each readability category across different heads $a = [a_1, \ldots, a_h] \in \mathbb{R}^{h \times n \times Y}$. Finally, a summation is performed across the multi-head dimension to obtain the aggregated sentence-level score matrix $A$,

$$A = \sum_{i=1}^{h} a_i \tag{16}$$

$A$ reflects the difficulty scores of each sentence across different readability categories.

The prediction result of the sentence label is not provided directly, but the information of the document-level label is propagated to the sentence level through the reverse label, and finally obtained through the MDEM module. In order to ensure the accuracy and consistency of the sentence label, the experiment introduced the KL divergence to measure the probability consistency between the sentence label and the document label. First, the sentence score matrix A is weighted summed to obtain the text score vector $\hat{\mathbf{r}}$, and the weight comes from the "source2token" self-attention mechanism at the document word level. Then, the text score vector $\hat{\mathbf{r}}$ is processed by LogSoftmax and compared with the document prediction probability $\mathbf{d}$ obtained in supervised learning. The probability distribution of the two is constrained by the KL divergence. The specific formula is as follows,

$$\mathcal{L}_{\text{unsup}} = \frac{1}{N} \sum_{i=1}^{N} D_{KL}(\text{LogSoftmax}(\hat{\mathbf{r}}) \| \mathbf{d}) \tag{17}$$

KL divergence is used as the loss function of the unsupervised part to directly correct the weights of the MDEM module and the sentence encoder, improve the distribution alignment capability, and ensure the semantic consistency between sentence-level labels and document labels. It helps the model learn semantic features consistent with document labels from unlabeled sentences and further optimizes the prediction accuracy of sentence labels. It also forms a two-way closed loop of "document label-driven sentence annotation—sentence feature reconstruction document prediction" with supervised learning.

**Hierarchical consistency training.** A hybrid supervision strategy is adopted during the training process. When supervised learning and unsupervised learning are used for joint training, a weighting factor $\lambda$ is used to balance the supervised cross entropy and unsupervised consistency training losses. Formally speaking, the complete training objective formula is as follows,

$$\mathcal{L} = \mathcal{L}_{\text{sup}} + \lambda \mathcal{L}_{\text{unsup}} \tag{18}$$

In the hierarchical consistency training, the experiment adopted three training techniques proposed by Google UDA [45] to ensure that the model establishes a stable and consistent semantic mapping between the document level and the sentence level. These techniques are Training Signal Annealing (TSA), Confidence-based Masking, and Sharpening Predictions.

### 3.3. Forward Text Readability Assessment

After using the sentence corpus obtained by HHNN-MEDM, sentence labels are used to assist in predicting the document-level readability level. Specifically, the structure of the DSDR [46] model is borrowed, as shown in the Fig. 3, and combined with the Ranking Model, DSDRRM is proposed.

### 3.3.1. DSDR

First, the pre-trained BERT model is directly used to train the sentence corpus to obtain an enhanced pre-trained model, which combines the difficulty multi-view representation and multi-view representation fusion, and can effectively apply sentence labels to document-level readability prediction, thereby further improving the accuracy of readability evaluation.

**Sentence-level difficulty-aware pre-training.** Use BERT to perform supervised training on a sentence corpus and build an enhanced pre-training model (EPTM) to learn the difficulty representation capability of sentences.
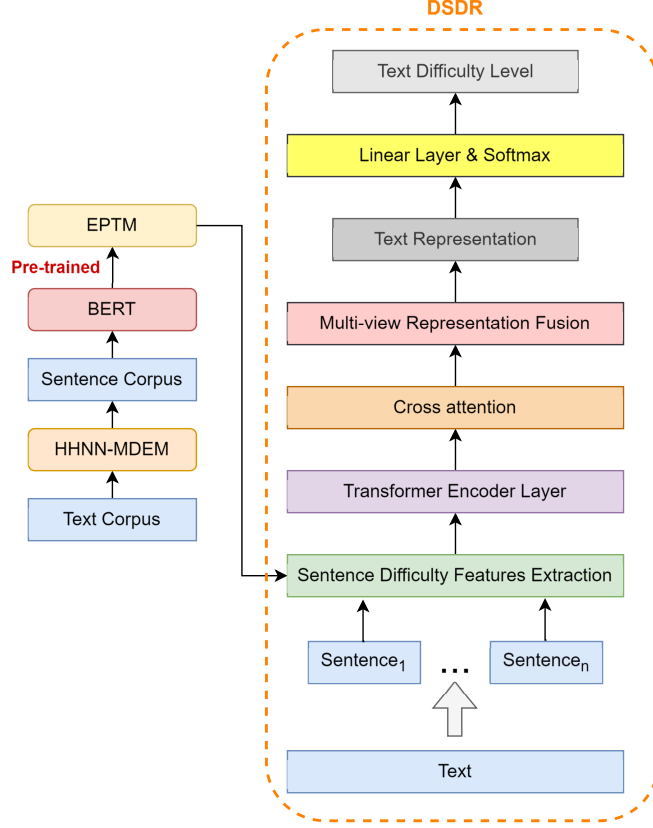
Figure 3: The overall structure of the DSDR.

**Difficulty multi-view representation.** EPTM is used to generate the representation of each sentence, and the context information is supplemented by the Transformer encoder to form a sentence-level context representation. Furthermore, multiple learnable difficulty vectors $C \in \mathbb{R}^{m \times d}$ are introduced, where m is the number of difficulty categories and d is the dimension of the sentence vector. The semantic representation of sentences at different difficulty levels is extracted through the cross-attention mechanism to form a multi-view representation $R = Attention\left(CW^Q, H^t W^k, H^t W^V\right)$.

**Multi-view representation fusion.** The multi-view difficulty representations are fused by average pooling to obtain the document-level representation $T$.

### 3.3.2. Ranking Model

In ARA prediction tasks, the order relationship between categories is important, and there is an obvious order relationship between adjacent difficulty labels, which is difficult to capture by simple classification tasks, because classification tasks usually assume that categories are equally weighted and independent. Therefore, we propose the Ranking Model, which typically requires the model to output a specific category of the corresponding sample for classification tasks. The Ranking Model further attempts to predict the order relationship of categories.

Firstly, based on the original data set, a number of data subsets are constructed according to the process in Fig. 4. Each data subset contains samples corresponding to the number of readability levels. That is, there are as many samples in the data subset as there are classes in the readability level, and there are no samples of the same level in the subset. Then, the samples in each subset are combined by pairwise permutation, and the difference in difficulty labels of the samples is used as the new label. This comparison allows the model to learn the differences between neighboring categories and thus better understand the sequential relationships between categories. In addition, this pairwise comparison method can also greatly expand the amount of data, that is, each data subset can construct $B_Y^{Y-1} = Y \times (Y-1)$ samples. The combined two samples are concatenated $d_{concat} = [d_1; d_2]$ into a fully connected layer
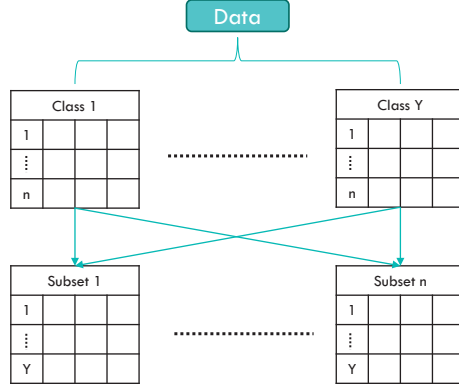
Figure 4: Data subset construction process. Suppose we have readability rank Y ordered categories in our dataset, each with n samples.

Table 1: Statistics of datasets for OSP,CEE, and CTRDG. Avg.Length is the average number of words per text.

| Grade | OSP | | CEE | | CTRDG | |
|---|---|---|---|---|---|---|
| | Texts | Avg.Length | Texts | Avg.Length | Texts | Avg.Length |
| **1** | 189 | 519 | 64 | 140 | 714 | 14.12 |
| **2** | 189 | 654 | 60 | 268 | 1102 | 23.32 |
| **3** | 189 | 809 | 71 | 613 | 1310 | 40.69 |
| **4** | - | - | 67 | 767 | 971 | 85.26 |
| **5** | - | - | 69 | 752 | 1163 | 235.28 |
| **6** | - | - | - | - | 461 | 580.57 |

to act as a classifier. Finally, each data subset was put into the model as a batch to continue training, and cross-entropy loss was used as the training target. By randomly combining samples with samples from different data subsets, a single sample will have multiple prediction results, which increases the error tolerance of each sample. Finally, the final readability level of each sample is obtained by hard voting.

As for the test set, we combined the data of the test set with the data subset of the training set to form a pairwise comparison. By predicting the label difference and adding it to the sample label in the data subset, we finally get the readability grade of the sample in the test set. Similarly, each sample grade in the test set is obtained by hard voting.

## 4. Data

In order to verify the effectiveness of our proposed model, we conduct experiments on five datasets of long documents, including two English datasets and three Chinese datasets. We split the training and test sets with an 8:2 ratio across all datasets. The dataset statistical distribution is shown in Table 1 and Table 2.

**OneStopEnglish(OSP)** [47] is a parallel corpus that can be used for automatic readability assessment and automatic text simplification. The corpus consists of 189 texts, each with three versions (567 texts in total).

**Cambridge(CEE)** [37] is a corpus for L2 learners that contains reading articles from five major Cambridge English examinations (KET, PET, FCE, CAE, CPE). The five exams are aimed at learners at levels A2-C2 of the Common European Frame of Reference.

**CMER** [34] consists of texts from extracurricular reading books for kids and teenagers at China mainland currently on the book market, with a total of 3,395,923 characters, distributed in 2,260 texts in 12 levels.

**CLT** is a dataset of Chinese language textbooks that we collected from primary and secondary school textbooks from multiple publishing houses. All the Chinese textbooks are taken from the first grade of primary school to the third grade of junior high school, with a total of 9 grades.

**CTRDG** [48] is a dataset about the Chinese Proficiency Test (HSK), with a total of 6 levels.

Table 2: Statistics of datasets for CMER, and CLT. Avg.Length is the average number of words per text.

| Grade | CMER | | CLT | |
|---|---|---|---|---|
| | Texts | Avg.Length | Texts | Avg.Length |
| **1** | 218 | 164 | 107 | 112 |
| **2** | 217 | 347 | 181 | 191 |
| **3** | 234 | 604 | 203 | 308 |
| **4** | 229 | 699 | 192 | 429 |
| **5** | 199 | 757 | 171 | 534 |
| **6** | 255 | 775 | 155 | 651 |
| **7** | 221 | 1352 | 91 | 1237 |
| **8** | 204 | 1409 | 85 | 1124 |
| **9** | 187 | 1429 | 61 | 1927 |
| **10** | 100 | 2384 | - | - |
| **11** | 95 | 2418 | - | - |
| **12** | 97 | 2226 | - | - |

## 5. Experimental Setup

This section presents the comparison baselines with our proposed model, the evaluation metrics, and the detailed details of the experimental implementation.

### 5.1. Statistical classification algorithms

This baseline is based on the explicit features of traditional classifiers including **Logistic Regression (LR)**, **Random Forest (RF)**, and **Support Vector Machines (SVM)**. In Section 3.1 we introduce the traditional features about Chinese and English. The model was implemented using the scikit-learn [49] tool, and the hyperparameters were dominated by default Settings.

### 5.2. Neural document classifiers

Such baselines represent another line of previous works that employ variants of neural document models for sentence or document classification.

**Vec2Read** [22] uses static word embeddings, Bi-LSTM, word level and sentence level attention mechanisms. Word Attention version of Vec2Read model is adopted in this experiment. The embedding size and hidden layer size of Bi-LSTM are set to 300 and 128, respectively.

**ReadNet** [32] proposed a new synthesis framework based on transformers that uses a hierarchical self-attention model to analyze document readability. The version of ReadNet model without explicit features is used in this experiment. For article coding, following the settings of the original paper, the number of sentences in each article and the number of words in each sentence were both limited to a maximum of 50, and the number of encoder layers p and q were set to 6. The embedding dimension is d = 100.

**HAN** [20] uses two GRUs, word level and sentence level attention mechanisms to encode word and sentence representations. This experiment uses the same experimental setup as Martinc et al. [15], where the context vector is randomly initialized and the word and sentence embedding sizes are 200 and 100, respectively.

**BERT** [19] is fine-tuned on English and Chinese using bert-base-uncased[1] and bert-base-chinese[2] respectively, with the default learning rate of 2e-5.

**DTRA** [34] uses a hierarchical attention network composed of BERT and Bi-LSTM, combined with word level and sentence level attention mechanisms, and performs model pre-training through soft labels of ordinal regression and predicting pairwise relative text difficulty. Since the paper does not provide the corresponding code and experiment detailed details, the results of the paper are extracted directly.

---

[1] https://huggingface.co/google-bert/bert-base-uncased
[2] https://huggingface.co/google-bert/bert-base-chinese

Table 3: Experimental results of readability evaluation on English and Chinese datasets.

| Dataset | Metrics | LR | RF | SVM | VecRead | ReadNet | HAN | Bert | DTRA | Lite-DTRA | DSDRRM |
|---|---|---|---|---|---|---|---|---|---|---|---|
| OSP | acc | 48.25 | 81.58 | 57.89 | 55.85 | 84.21 | 80.70 | 78.66 | 85.00 | 86.67 | **89.47** |
| | adj-acc | 89.47 | 99.12 | 93.86 | 100.00 | 100.00 | 100.00 | 97.96 | 100.00 | 100.00 | **100.00** |
| | F1 | 47.65 | 81.47 | 57.39 | 56.16 | 84.36 | 79.93 | 78.04 | 84.91 | 86.79 | **89.38** |
| | p | 47.71 | 81.74 | 58.23 | 61.02 | 85.15 | 81.41 | 79.95 | - | - | **89.38** |
| | r | 48.25 | 81.58 | 57.89 | 55.85 | 84.21 | 80.70 | 78.66 | - | - | **89.47** |
| | qwk | 40.05 | 83.81 | 54.56 | 64.60 | 87.50 | 86.19 | 79.59 | - | - | **92.00** |
| CEE | acc | 47.76 | 80.60 | 50.75 | 53.03 | 72.73 | 73.74 | 62.63 | - | - | **83.58** |
| | adj-acc | 77.61 | 94.03 | 76.12 | 88.38 | 95.45 | 92.93 | **97.47** | - | - | 97.01 |
| | F1 | 42.72 | 79.84 | 43.54 | 51.17 | 72.07 | 72.74 | 58.38 | - | - | **83.64** |
| | p | 45.45 | **85.65** | 46.37 | 56.19 | 77.50 | 75.90 | 64.25 | - | - | 83.34 |
| | r | 47.76 | 80.60 | 50.75 | 53.03 | 72.73 | 73.74 | 62.63 | - | - | **83.58** |
| | qwk | 57.18 | 91.27 | 67.14 | 75.12 | 88.25 | 87.19 | 88.68 | - | - | **94.05** |
| CMER | acc | 23.06 | 28.38 | 22.62 | 22.20 | 26.40 | 26.11 | 28.61 | 26.50 | 26.50 | **48.89** |
| | adj-acc | 50.78 | 59.87 | 49.67 | 46.46 | 57.37 | 53.10 | 56.64 | 58.50 | 62.47 | **76.99** |
| | F1 | 20.28 | 27.28 | 17.51 | 19.64 | 26.61 | 24.87 | 26.39 | 25.16 | 22.06 | **48.51** |
| | p | 21.26 | 28.11 | 15.98 | 24.95 | 29.45 | 27.36 | 30.30 | - | - | **50.59** |
| | r | 23.06 | 28.38 | 22.62 | 22.20 | 26.40 | 26.11 | 28.61 | - | - | **48.89** |
| | qwk | 64.29 | 71.80 | 64.11 | 57.97 | 75.37 | 76.60 | 71.20 | - | - | **85.08** |
| CLT | acc | 30.40 | 42.00 | 35.20 | 29.33 | 41.47 | 41.73 | 37.87 | - | - | **46.00** |
| | adj-acc | 76.80 | 81.60 | 78.80 | 64.80 | **85.60** | 84.00 | 75.73 | - | - | 84.80 |
| | F1 | 27.32 | 41.61 | 30.79 | 26.93 | 41.59 | 41.29 | 33.29 | - | - | **46.56** |
| | p | 27.92 | 42.71 | 34.42 | 28.53 | 44.53 | 46.02 | 34.88 | - | - | **50.19** |
| | r | 30.40 | 42.00 | 35.20 | 29.33 | 41.47 | 41.73 | 37.87 | - | - | **46.00** |
| | qwk | 74.98 | 85.17 | 77.69 | 54.85 | **85.22** | 83.40 | 73.59 | - | - | 84.93 |
| CTRDG | acc | 38.95 | 76.07 | 24.19 | 75.78 | 80.70 | 82.18 | 89.43 | - | - | **90.48** |
| | adj-acc | 77.73 | 99.65 | 58.34 | 97.00 | 99.39 | 99.83 | **99.91** | - | - | 99.74 |
| | F1 | 27.94 | 75.96 | 10.43 | 76.04 | 80.66 | 82.18 | 89.18 | - | - | **90.50** |
| | p | 35.48 | 75.97 | 18.71 | 77.70 | 81.26 | 82.34 | 89.77 | - | - | **90.63** |
| | r | 38.95 | 76.07 | 24.19 | 75.78 | 80.70 | 82.18 | 89.43 | - | - | **90.48** |
| | qwk | 61.65 | 94.79 | 2.91 | 92.81 | 95.52 | 96.16 | 97.67 | - | - | **97.84** |

**Lite-DTRA** [34] is a streamlined version of the DTRA model, proposed to reduce the requirements for hardware storage memory. In this version, the pre-trained BERT with frozen parameters is replaced by ALBERT [50], allowing the model to be trained in an end-to-end manner. Similarly, due to the lack of corresponding code and detailed experimental information provided in the paper, the results of the paper are extracted directly.

## 5.3. Training and Evaluation Details

We used the Pytorch [51] framework for our experiments. In the Embedding layer, the output dimension of the other datasets is 400 except for CMER, which has an output dimension of 512. Similarly, the hidden layer of Bi-LSTM, the number of convolutional kernels of CNN and the number of heads of multi-head self-attention h size in CMER are 256, 256, 16 respectively. 200, 200, 8 in the other datasets, respectively. The window size of the convolution kernel is 3. During training, the learning rate is set to 1e-3, Adam [52] is used as the optimizer, and the weights decay to 5e-4. TSA was performed in linear form, the $\beta$ threshold based on confidence masking was set to 0.45, the temperature parameter $\tau$ was set to 0.85, the training rounds were set to 30, and the experimental hyperparameters of the DSDR-MDEM model followed the settings of the DSDR model. For evaluation, we calculated precision (acc), adjacent accuracy (adj-acc), weighted F1 score (F1), precision (p), recall (r), and quadratically weighted kappa (qwk). We repeat each experiment three times and report the average score.

Following previous work on readability evaluation, we use qwk as a primary metric to reflect the ordinal alignment between predicted readability levels and ground-truth labels. Although rank correlation metrics such as Spearman's $\rho$ and Kendall's $\tau$ have been recommended for ordinal tasks [53], qwk provides a widely accepted, label-sensitive alternative that penalizes larger rank discrepancies more severely.

## 6. Experimental Results

We report experimental results on all datasets (Section 6.1). We then present an ablation study (Section 6.2) and a comparison between single dimensional context weights (Section 6.3) and ordinal regression (Section 6.4) with DSDRRM.

### 6.1. Overall Results

The experimental results of all models are summarized in Table 3. Our DSDRRM achieves consistent improvements over the baselines on all datasets, which verifies the effectiveness of our proposed method. First, on the English dataset, our method has slightly lower p and qwk than the baseline model on the CEE dataset, but outperforms all baselines on other indicators on both English datasets, with OSP improving the accuracy by 2.8% and CEE improving the accuracy by 2.98%. To our surprise, on the Chinese dataset, our model improves the accuracy by 22.39% over DTRA on the CMER dataset, and also improves to varying degrees on the CLT and CTRDG datasets. Second, in terms of the order of readability labels, our model achieves great improvements on qwk, which also verifies the effectiveness of the Ranking Model on ordered labels. Finally, surprisingly, RF seems to be more suitable for the evaluation of text readability, and the results on the Chinese dataset are even comparable to the neural network baseline, which illustrates the effectiveness of explicit features designed for Chinese datasets. This also lays the foundation for the future combination of explicit features and neural network features.

In the experiments, it can be observed that there are obvious differences in the accuracy of the DSDRRM model on four different readability datasets. The reason is that the structure of the dataset and the label standardization have an important impact on the model performance. There are differences in the modelability of the language itself. The grammatical rules of English are relatively fixed, and the subject-verb-object structure is clear. The granularity of Chinese language units is fuzzy, the syntax is flexible, and the semantic dependence is long, which makes modeling relatively difficult. In addition, the difference in the number of levels also significantly affects the performance of the model. As the number of classification levels increases, the model needs to make judgments in a finer-grained label space, and the classification difficulty increases accordingly. CMER has a total of 12 levels, which is far more than the 3 levels of OSP and the 6 levels of CTRDG. This increases the difficulty of learning the model, especially when the sample distribution is uneven. The model is more likely to be biased towards the prediction of the middle level, thereby reducing the overall accuracy.

### 6.2. Ablation Study

In order to measure the contribution of multi-dimensional context weight, sentence tag assistance and ranking model to the model, we conducted ablation experiments on OSP and CLT. F1 and qwk were selected as evaluation indicators in the experiment. The experimental results are shown in Table 4.

After removing the multi-dimensional context weight vector, the F1 and qwk of both datasets decreased, verifying that the multi-dimensional context weights contain useful information and can guide the attention model to locate information-rich words from the input sequence, thus playing an important role in the attention mechanism. After removing the sentence label auxiliary step, the performance of the model on the dataset also decreased. This result fully verifies the important role of sentence label assistance in improving document-level readability assessment. The strategy of sentence label-assisted document-level readability assessment provides an effective solution for fine-grained optimization of readability assessment tasks and significantly improves the generalization ability of the model. After removing the Ranking Model module, the OSP dataset has a significant decrease in F1 and qwk indicators. Due to the ordinal characteristics of readability labels, qwk also reflects the practicality of the Ranking Model module for ARA tasks.

### 6.3. Multi vs. Single Context: Comparative Analysis

In order to further verify the effectiveness of multidimensional context weights, we construct a single dimensional context weight for comparison experiments, that is, the same weight is used for each feature dimension in the word vector output by Bi-LSTM. Similar to Section 3.2.1, k convolution kernels are used to capture the output $h_i^t$ of Bi-LSTM to obtain the context information $\hat{c} = [c_1, \ldots, c_{m-l+1}]$. The pooling layer converts text of various lengths into fixed-length vectors. With the pooling layer, we can capture the information of the entire text. Therefore, the average

Table 4: Results of ablation experiments. -Context means removing multi-dimensional context weights. -MDEM means removing sentence tag assistance. -Ranking Model means removing the Ranking Model module.

| Model | OSP | | CLT | |
|---|---|---|---|---|
| | F1 | qwk | F1 | qwk |
| -Context | 88.6 | 90.92 | 45.20 | 84.60 |
| -MDEM | 86.06 | 89.16 | 42.40 | 80.07 |
| -Ranking Model | 87.58 | 90.68 | 45.69 | 84.40 |
| DSDRRM | 89.38 | 92.00 | 46.56 | 84.93 |

Table 5: Experimental results on Multi vs. Single Context. -SDW is represented with single dimensional context weights.

| Model | OSP | | CLT | |
|---|---|---|---|---|
| | F1 | qwk | F1 | qwk |
| DSDRRM-SDW | 89.33 | 92.20 | 45.16 | 83.88 |
| DSDRRM | 89.38 | 92.00 | 46.56 | 84.93 |

pooling operation is applied to $\hat{c}$ to extract the average $a = mean(\hat{c})$. Finally, the sentence representation vector $h_i^s = a \cdot h_i^t$ is obtained by matrix multiplication of $a$ and $h_i^t$.

We conducted experiments using OSP and CLT, and the experimental results are shown in Table 5. Under multi-dimensional context weighting, although the model performance has a small loss in the qwk indicator on the OSP dataset, it has improved to varying degrees on the CLT dataset. In NLP, some tokens are polysemous. Since the traditional attention mechanism calculates the overall weight score of each word based on the word vector, it is impossible to distinguish the meaning of the same word in different contexts. The multi-dimensional weight vector calculates a weight score for each feature of each word, so it can select the feature that best describes the specific meaning of the word in any given context and include this information in the sentence encoding output.

### 6.4. Ranking Model vs. Ordinal Regression

ARA can also be formulated as an ordinal regression task. Given a dataset with Y readability level categories, the document vector d represented by the model is input into a fully connected layer, outputting a readability label vector $r \in \mathbb{R}^Y$. The goal of ordinal regression is to minimize the ordinal regression loss, which is defined as follows:

$$L(r; y) = -log(Sigmoid(\theta_k - r_k) - Sigmoid(\theta_{k-1} - r_{k-1})) \qquad (19)$$

where $k = y$. $r_k$ denotes the $k^{th}$ dimension of the $r$. y is the true label. The threshold parameter $\theta_0, \ldots, \theta_{Y-1}$ is also learned automatically from data. The probability of the current class is calculated by the threshold, and this probability is obtained by comparing the difference between the current class and the neighboring classes. This process is similar to the Ranking Model, understanding the ordinal relationship between categories by subtracting.

Similarly, we conduct experiments with OSP and CLT to further examine the effect of Ranking Model by comparing ordinal regression and multi-class classification. The experimental results are shown in Table 6. No matter classification or ordinal regression, their results are worse than the Ranking model, which verifies the effectiveness of the pairwise comparison ranking algorithm. When comparing classification and ordinal regression, although ordinal regression is worse than classification on the weighted F1 measure, however, in terms of qwk, ordinal regression has a good improvement over classification on the whole. Therefore, it is also verified that the model can learn the differences between adjacent categories by doing subtraction, so as to better understand the sequential relationship between categories.

## 7. Conclusion and Future Work

This paper proposes a deep learning model for readability assessment of long documents. Compared to baseline models, our proposed forward and reverse readability assessment and pairwise sorting algorithms achieve competitive performance across all five datasets. In future work, we can further study the fusion of explicit features and neural

Table 6: Model accuracy based on classification (C), ordinal regression (OR), and Ranking Model(RM).

| Model | OSP | | CLT | |
|---|---|---|---|---|
| | F1 | qwk | F1 | qwk |
| DSDR | 87.58 | 90.68 | 45.69 | 84.40 |
| DSDR-OR | 89.34 | 91.57 | 44.28 | 85.23 |
| DSDR-RM | 89.38 | 92.00 | 46.56 | 84.93 |

network features to improve the performance of readability assessment. At the same time, sentence labels from different corpora can be applied in cross-corpus evaluation, and more domain adaptation techniques can be considered to find the optimal feature set capable of generalizing well to unseen texts.

## Acknowledgments

## References

[1] M. Vogel, C. Washburne, An objective method of determining grade placement of children's reading material, The Elementary School Journal 28 (1928) 373–381.

[2] U. Cop, D. Drieghe, W. Duyck, Eye movement patterns in natural reading: A comparison of monolingual and bilingual reading of a novel, PloS one 10 (2015) e0134008. doi:10.1371/journal.pone.0134008.

[3] R. Flesch, A new readability yardstick., Journal of applied psychology 32 (1948) 221.

[4] E. Dale, J. S. Chall, A formula for predicting readability: Instructions, Educational research bulletin (1948) 37–54.

[5] G. H. Mc Laughlin, Smog grading-a new readability formula, Journal of reading 12 (1969) 639–646.

[6] S. Villata, et al., Plain language assessment of statutes, in: Legal Knowledge and Information Systems: JURIX 2020: The Thirty-third Annual Conference, Brno, Czech Republic, December 9-11, 2020, volume 334, IOS Press, 2020, p. 207.

[7] A. Sare, A. Patel, P. Kothari, A. Kumar, N. Patel, P. A. Shukla, Readability assessment of internet-based patient education materials related to treatment options for benign prostatic hyperplasia, Academic Radiology 27 (2020) 1549–1554. doi:10.1002/lary.23424.

[8] S. Perni, M. K. Rooney, D. P. Horowitz, D. W. Golden, A. R. McCall, A. J. Einstein, R. Jagsi, Assessment of use, specificity, and readability of written clinical informed consent forms for patients with cancer undergoing radiotherapy, JAMA oncology 5 (2019) e190260–e190260. doi:10.1001/jamaoncol.2019.0260.

[9] S. Vajjala, Trends, limitations and open challenges in automatic readability assessment research, arXiv preprint arXiv:2105.00973 (2021). doi:10.48550/arxiv.2105.00973.

[10] S. E. Schwarm, M. Ostendorf, Reading level assessment using support vector machines and statistical language models, in: Proceedings of the 43rd annual meeting of the Association for Computational Linguistics (ACL'05), 2005, pp. 523–530. doi:10.3115/1219840.1219905.

[11] B. W. Lee, Y. S. Jang, J. Lee, Pushing on text readability assessment: A transformer meets handcrafted linguistic features, in: M.-F. Moens, X. Huang, L. Specia, S. W.-t. Yih (Eds.), Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, Online and Punta Cana, Dominican Republic, 2021, pp. 10669–10686. URL: https://aclanthology.org/2021.emnlp-main.834/. doi:10.18653/v1/2021.emnlp-main.834.

[12] H. Hansen, A. Widera, J. Ponge, B. Hellingrath, Machine learning for readability assessment and text simplification in crisis communication: A systematic review, in: Proceedings of the 54th Hawaii International Conference on System Sciences, 2021. doi:10.24251/hicss.2021.277.

[13] I. Pilán, E. Volodina, T. Zesch, Predicting proficiency levels in learner writings by transferring a linguistic complexity model from expert-written coursebooks, in: Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers, 2016, pp. 2101–2111.

[14] W. B. Li, Z. Wang, Y. Wu, A unified neural network model for readability assessment with feature projection and length - balanced loss, arXiv preprint arXiv:2210.10305 (2023). doi:10.48550/arxiv.2210.10305.

[15] M. Martinc, S. Pollak, M. Robnik-Šikonja, Supervised and unsupervised neural approaches to text readability, Computational Linguistics 47 (2021) 141–179. doi:10.1162/COLI_A_00398.

[16] J. M. Imperial, Bert embeddings for automatic readability assessment, in: Proceedings of the International Conference on Recent Advances in Natural Language Processing, 2021, pp. 611–618. doi:10.26615/978-954-452-072-4_069.

[17] B. W. Lee, J. H.-J. Lee, Prompt-based learning for text readability assessment, arXiv preprint arXiv:2302.13139 (2023). doi:10.48550/arxiv.2302.13139.

[18] J. Risch, R. Krestel, Bagging bert models for robust aggression identification, in: Proceedings of the Second Workshop on Trolling, Aggression and Cyberbullying, European Language Resources Association (ELRA), Marseille, France, 2020, pp. 55–61. URL: https://aclanthology.org/2020.trac-1.9/. doi:10.5281/zenodo.3727018.

[19] J. D. M.-W. C. Kenton, L. K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, in: Proceedings of NAACL-HLT, 2019, pp. 4171–4186. doi:10.18653/v1/N19-1423.

[20] Z. Yang, D. Yang, C. Dyer, X. He, A. Smola, E. Hovy, Hierarchical attention networks for document classification, in: Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, 2016, pp. 1480–1489. URL: http://www.aclweb.org/anthology/N16-1174. doi:10.18653/v1/N16-1174.

[21] J. Zheng, L. Zheng, A hybrid bidirectional recurrent convolutional neural network attention - based model for text classification, IEEE Access 7 (2019) 106673–106685. doi:10.1109/access.2019.2932619.

[22] I. M. Azpiazu, M. S. Pera, Multiattentive recurrent neural network architecture for multilingual readability assessment, Transactions of the Association for Computational Linguistics 7 (2019) 421–436. URL: https://aclanthology.org/Q19-1028/. doi:10.1162/tacl_a_00278.

[23] J. Lee, S. Vajjala, A neural pairwise ranking model for readability assessment, in: Findings of the Association for Computational Linguistics: ACL 2022, Association for Computational Linguistics, 2022, pp. 3802–3813. doi:10.18653/v1/2022.findings-acl.300.

[24] J. P. Kincaid, R. P. Fishburne Jr, R. L. Rogers, B. S. Chissom, Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel (1975).

[25] M. Heilman, K. Collins-Thompson, M. Eskenazi, An analysis of statistical models and features for reading difficulty prediction, in: Proceedings of the third workshop on innovative use of NLP for building educational applications, 2008, pp. 71–79.

[26] L. Feng, M. Jansche, M. Huenerfauth, N. Elhadad, A comparison of features for automatic readability assessment, in: Coling 2010: Posters, 2010, pp. 276–284.

[27] J. Hancke, S. Vajjala, D. Meurers, Readability classification for german using lexical, syntactic, and morphological features, in: Proceedings of COLING 2012, 2012, pp. 1063–1080.

[28] X. Qiu, K. Deng, L. Qiu, X. Wang, Exploring the impact of linguistic features for chinese readability assessment, in: Natural Language Processing and Chinese Computing: 6th CCF International Conference, NLPCC 2017, Dalian, China, November 8–12, 2017, Proceedings 6, Springer, 2018, pp. 771–783. doi:10.1007/978-3-319-73618-1_67.

[29] T. Deutsch, M. Jasbi, S. Shieber, Linguistic features for readability assessment, arXiv preprint arXiv:2006.00377 (2020). doi:10.18653/v1/2020.bea-1.1.

[30] Z. Jiang, G. Sun, Q. Gu, T. Bai, D. Chen, A graph-based readability assessment method using word coupling, in: Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, 2015, pp. 411–420. doi:10.18653/v1/D15-1047.

[31] P. G. Blaneck, T. Bornheim, N. Grieger, S. Bialonski, Automatic readability assessment of german sentences with transformer ensembles, arXiv preprint arXiv:2209.04299 (2022).

[32] C. Meng, M. Chen, J. Mao, J. Neville, Readnet: A hierarchical transformer framework for web article readability analysis, in: Advances in Information Retrieval: 42nd European Conference on IR Research, ECIR 2020, Lisbon, Portugal, April 14–17, 2020, Proceedings, Part I 42, Springer, 2020, pp. 33–49. doi:10.1007/978-3-030-45439-5_3.

[33] J. D. Rennie, N. Srebro, Loss functions for preference levels: Regression with discrete ordered labels, in: Proceedings of the IJCAI multidisciplinary workshop on advances in preference handling, volume 1, AAAI Press, Menlo Park, CA, 2005.

[34] J. Zeng, Y. Xie, X. Yu, J. S. Lee, D.-X. Zhou, Enhancing automatic readability assessment with pre-training and soft labels for ordinal regression, in: Findings of the Association for Computational Linguistics: EMNLP 2022, 2022, pp. 4557–4568. doi:10.18653/v1/2022.findings-emnlp.334.

[35] R. Diaz, A. Marathe, Soft labels for ordinal regression, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 4738–4747. doi:10.1109/CVPR.2019.00487.

[36] K. Tanaka-Ishii, S. Tezuka, H. Terada, Sorting texts by readability, Computational linguistics 36 (2010) 203–227. doi:10.1162/coli.09-036-R2-08-050.

[37] M. Xia, E. Kochmar, T. Briscoe, Text readability assessment for second language learners, arXiv preprint arXiv:1906.07580 (2019). doi:https://doi.org/10.18653/v1/W16-0502.

[38] I. Madrazo Azpiazu, M. S. Pera, Is cross-lingual readability assessment possible?, Journal of the Association for Information Science and Technology 71 (2020) 644–656. doi:10.1002/asi.24293.

[39] Y.-T. Sung, J.-L. Chen, Y.-S. Lee, J.-H. Cha, H.-C. Tseng, W.-C. Lin, T.-H. Chang, K.-E. Chang, Investigating chinese text readability: linguistic features, modeling, and validation., Chinese Journal of Psychology (2013).

[40] K. Ma, Z. Liu, L. Yang, N. Sun, Y. Wang, Z. Qiu, Research on the evaluation of the classical chinese difficulty in the compulsory education stage, in: 2022 International Conference on Asian Language Processing (IALP), IEEE, 2022, pp. 353–357.

[41] J. Du, L. Gui, Y. He, R. Xu, X. Wang, Convolution-based neural attention with applications to sentiment classification, IEEE Access 7 (2019) 27983–27992. doi:10.1109/ACCESS.2019.2900335.

[42] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, Advances in neural information processing systems 30 (2017).

[43] Y. Hu, P. Chen, T. Liu, J. Gao, Y. Sun, B. Yin, Hierarchical attention transformer networks for long document classification, in: 2021 International Joint Conference on Neural Networks (IJCNN), IEEE, 2021, pp. 1–7. doi:10.1109/IJCNN52387.2021.9533869.

[44] T. Shen, T. Zhou, G. Long, J. Jiang, S. Pan, C. Zhang, Disan: Directional self-attention network for rnn/cnn-free language understanding, in: Proceedings of the AAAI conference on artificial intelligence, volume 32, 2018. doi:10.48550/arXiv.1709.04696.

[45] Q. Xie, Z. Dai, E. Hovy, T. Luong, Q. Le, Unsupervised data augmentation for consistency training, Advances in neural information processing systems 33 (2020) 6256–6268.

[46] W. Li, Research on text readability evaluation based on neural network model, Master's thesis, Peking University, 2023.

[47] S. Vajjala, I. Lučić, Onestopenglish corpus: A new corpus for automatic readability assessment and text simplification, in: Proceedings of the thirteenth workshop on innovative use of NLP for building educational applications, 2018, pp. 297–304. doi:10.18653/v1/W18-0535.

[48] K. Tan, Y. Lan, Y. Zhang, A. Ding, A chinese text readability classification model based on multi-level linguistic feature fusion, Journal of Chinese Information Processing 38 (2024) 41–52.

[49] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al., Scikit-learn: Machine learning in python, the Journal of machine Learning research 12 (2011) 2825–2830. doi:10.5555/1953048.2078195.

[50] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, R. Soricut, Albert: A lite bert for self-supervised learning of language representations, arXiv preprint arXiv:1909.11942 (2019).

[51] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al., Pytorch: An imperative style, high-performance deep learning library, Advances in neural information processing systems 32 (2019).

[52] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, arXiv preprint arXiv:1412.6980 (2014). doi:10.48550/arXiv.1412.6980.

[53] Y. Ehara, Evaluation of unsupervised automatic readability assessors using rank correlations, in: Proceedings of the 2nd Workshop on Evaluation and Comparison of NLP Systems, 2021, pp. 62–72. doi:10.18653/v1/2021.eval4nlp-1.7.

## Appendix A. Chinese Explicit features Details

Table A.1: Lexical class features

| Categories | Feature name |
|---|---|
| Number of words | Word count |
| | Word count |
| Vocabulary richness | Dissimilar Word Ratio (TTR) |
| | Content word density |
| Word frequency | Log average of real word frequencies |
| | Number of difficult words |
| Vocabulary length | Low stroke number of characters |
| | The number of characters in the stroke |
| | Number of high stroke characters |
| | Average number of strokes per character |
| | Number of two-character words |
| | More than three words |
| Semantic class metrics | Number of content words |
| | Negative words |
| Syntactic metrics | Average number of words in simple sentences |
| | Average number of words in complex sentences |
| | Single sentence ratio |
| | Noun phrase ratio |
| Vocabulary density | Number of function words |
| | Function word density |
| Diversity of words | RTTR |
| | MTLD |

Table A.2: Part of speech feature

| Categories | Feature name |
|---|---|
| Adjective | Percentage of adjectives |
| | Percentage of unique adjectives |
| | Number of unique adjectives |
| | Average number of adjectives per sentence |
| | Average number of unique adjectives per sentence |
| Noun | Percentage of nouns |
| | Percentage of unique nouns |
| | Number of unique nouns |
| | Average number of nouns per sentence |
| | Average number of unique nouns per sentence |
| Verb | Percentage of verbs |
| | Percentage of unique verbs |
| | Number of unique verbs |
| | Average number of verbs per sentence |
| | Average number of unique verbs per sentence |

Table A.3: Discourse feature

| Categories | Feature name |
|---|---|
| Solid density | Number of named entities |
| | Number of unique named entities |
| | Percentage of named entities |
| | Percentage of unique named entities |
| | Average number of named entities per sentence |
| | Number of unique named entities per sentence |

Table A.4: Article cohesion characteristics

| Categories | Feature name |
|---|---|
| Term of reference | Number of pronouns |
| Connective words | Number of connectives |
| | Positive connectives |
| | Negative connectives |