# NwQM: A neural quality assessment framework for Wikipedia

**Bhanu Prakash Reddy**[*1]**,Sasi Bhushan**[*2]**,Soumya Sarkar**[*3]**,Animesh Mukherjee**[4]

IIT Kharagpur, India[2,3,4], Adobe Research, India[1]

soumya015@iitkgp.ac.in[3],guda@adobe.com[1]
sasibhushan3@gmail.com[2], animesh@cse.iitkgp.ac.in[4]

## Abstract

Millions of people irrespective of socio-economic and demographic backgrounds, depend on Wikipedia articles everyday for keeping themselves informed regarding popular as well as obscure topics. Articles have been categorized by editors into several quality classes, which indicate their reliability as encyclopedic content. This manual designation is an onerous task because it necessitates profound knowledge about encyclopedic language, as well navigating circuitous set of *wiki* guidelines. In this paper we propose **N**eural **w**ikipedia **Q**uality **M**onitor (NwQM), a novel deep learning model which accumulates signals from several key information sources such as *article text*, *meta data* and *images* to obtain improved Wikipedia article representation. We present comparison of our approach against a plethora of available solutions and show 8% improvement over state-of-the-art approaches with detailed ablation studies.

## 1 Introduction

Wikipedia is one of the most prominent sources of free information in the world today. Since reliable and advertisement free material is mostly behind pay-walled sources, a huge volume of global population is directed toward this extraordinary crowd sourced platform for information ranging from history, politics, pop-culture to even scientific topics (Horta Ribeiro et al., 2020).

Although Wikipedia has grown significantly in terms of volume and veracity over the last decade, the quality of articles is not uniform (Warncke-Wang et al., 2015). The quality of Wikipedia articles is monitored through a rating system where each article is assigned one of several class indicators. Some of the *major* article categories are **FA**, **GA**, **B**, **C**, **Start** and **Stub**. Most complete

and dependable content is annotated by an FA (*aka featured article*) tag while lowest quality content is annotated with a Stub tag. The intention behind this elaborate scheme is to notify editors regarding current state of the article and extent of effort needed for escalating to encyclopedic standards[1].

There exist several guidelines which direct editors in annotating articles into respective classes. Some of the traits of a FA article are engaging and comprehensive prose with neutral point of view and verifiable claims. It must also rigorously follow the style manual, i.e. the page structure. Further the content should be stable, i.e. devoid of edit warring. Understanding compliance with these guidelines often require detailed knowledge about language usage as well as domain knowledge about Wikipedia page layout and style principles. Often it is nontrivial to discern qualifying differences between articles which merits their ratings without inculcating personal biases.

Consider the wikipage of two prominent US presidents *Abraham Lincoln* (GA) and *John F. Kennedy* (B). Both pages are indistinguishable in-terms of coverage and engagement, however on closer assessment it is apparent that President Kennedy's page contains unattributed opinion such as the statement *This crisis brought the world closer to nuclear war than at any point before or after ...* in the *Cuban Missile Crisis* section. It also has vague quantifiers such as *some questioned, some crtics, somewhat successful etc.* Similarly, if we look at wikipages of historical figures *Akhenaten* (B) and *Cleopatra* (FA) it is difficult to discern their quality just from the content. A deep dive into the individual *talk pages* reveal that the former page has unresolved content issues and disputes which justifies the given rating. Hence it is a difficult task to manually judge language specific nuances

---

[1] wiki/Wikipedia:WikiProject Wikipedia/Assessment

present in the main page text, topic level disputes manifesting in the talk pages as well as section layout and image positioning before deciding a correct rating. This is also apparent from the quality statistics[2] which shows that only $0.09\%, 0.5\%$ of the $\sim 6M$ English Wikipedia articles have a FA and a GA tag respectively.

Current approaches for automatic quality assessment by *Wikimedia foundation*[3] use handcrafted features from main article text for classifying quality class (Halfaker and Geiger, 2019). Analogous approaches exist (Dang and Ignat, 2016) which attempt to automatically generate features from main text using deep learning models such as *doc2vec* (Le and Mikolov, 2014). Other approaches use deep sequence models such as BIL-STM (Shen et al., 2017) as well as combining representations obtained with additional modalities like image (Shen et al., 2019b). The principal focus of existing works have been concentrated on main article text. However one of the key sources of metadata about an article, i.e., the corresponding *talk page* has been ignored. Talk pages contain crucial information concerning stability of a page. They also hold evidence whether discussion threads encompassing topics are decisive. Besides, representation of main page text using sequence models cannot capture high level semantic signals such as whether the introduction section is a summary, whether the coverage of topics is polarized or the information is redundant or the wording is convoluted.

In this paper we propose **N**eural **w**ikipedia **Q**uality **M**onitor (**NwQM**) which integrates information from multiple sources, i.e., main page text, metadata and html rendering resulting in improved quality assessment. We use bidirectional contextual representation (Devlin et al., 2018) for encoding article text. However, contrary to document representation using BERT (Adhikari et al., 2019), which is not adequate for large text documents, we first segment articles organically based on sections. We fine tune on each section text individually followed by a summarization layer which preserves the sequential nature of sections. Similar representation of atomic units and summarisation is applied on talk pages. We also obtain images from the raw markup using Imagekit[4]. These images are further embedded in a vector space using Inception V3 (Szegedy

et al., 2016). Inception V3 is pre-trained on Imagenet [5] and we fine tune on our dataset to cater to our task. Our experiments show that combining information sources from these sources leads to improved result eclipsing current state of the art (Shen et al., 2019b) by $8\%$.

Our main contributions are enumerated below.

1. We propose a multimodal framework from quality assessment of Wikipedia articles which leverages contextual representation obtained from bidirectional transformers and supports conditional summarization.

2. To the best of our knowledge this is the first work which utilizes meta pages, i.e., talk pages as an additional signal for this task. All code, sample data and image embeddings related to the paper are made available[6,7] to promote reproducible research.

## 2  Related work

Automatic article assessment is one of the key research agendas of the Wikimedia foundation[8]. One of the preliminary approaches (Halfaker and Taraborelli, 2015) seeking to solve this problem extracted structural features such as presence of infobox, references, level 2 headings etc. as indicators of the article quality. Other approaches explored distributional representation as well as sequence models (Dang and Ignat, 2016; Shen et al., 2017, 2019b). (Zhang et al., 2018) attempted to solve this problem by formulating features capturing dynamic nature of the articles. A complementary direction of exploration has been put forward by (Li et al., 2015; de La Robertie et al., 2015) where correlation between article quality and structural properties of co-editor network and editor-article network has been exploited.

This task can also be solved by exploring the rich literature of document classification. One of the characteristics of Wikipedia articles are that these are long documents hence sequence models can suffer from (Atkinson, 2018) catastrophic forgetting. (Yang et al., 2016) proposed to solve this by leveraging a hierarchical organization of documents. Further improvements have been demonstrated by employing bidirectional transformers (Adhikari et al., 2019; Ostendorff et al., 2019). However to the best of our knowledge no

---

[2] `wiki/Wikipedia:Good article statistics`
[3] `wikimediafoundation.org`
[4] `pypi.org/project/imgkit`

[5] `www.image-net.org/`
[6] `https://github.com/sasibhushan3/NwQM_EMNLP`
[7] `https://zenodo.org/record/4066405`
[8] `www.mediawiki.org/wiki/ORES`

| Class | Article count |
|-------|---------------|
| FA | 3589 |
| GA | 5900 |
| B | 5900 |
| C | 5900 |
| Start | 5900 |
| Stub | 5900 |
| Total | 33089 |

Table 1: Wikipedia dataset of articles with respective *talk* pages.

previous work have used metadata about article pages as source of additional signals. Also we investigate information fusion from multimodal sources for generating improved article representation.

## 3 Dataset

Wikimedia foundation stores all data for its multilingual wikiprojects in the form of Wikidumps[9]. We downloaded first 100 English Wikipedia dumps which are 7z archived *xml* files. The combined uncompressed size of these files is $\sim 8TB$ and it contains a sample of $\sim 6M$ English Wikipedia text from the first version to last update as of June 2019. Because of limitation of space we did not go though the entire English Wikipedia archive. Each uncompressed file is approximately of $80GB$ size and have random samples of approximately $\sim 5k$ Wikipedia pages. We parsed each file using mediawiki xml parser[10] and in one linear scan we tried to locate if the main page text and talk page text is in the same dump xml file. More specifically if we encounter the main Wikipedia article *Cleopatra*, we remember it in a dictionary and seek to locate *Talk:Cleopatra* in future scans or vice versa. However it is entirely possible that *Talk:Cleopatra* is not present at all in the currently encountered xml file and may be present somewhere else; in such case we ignore that article. If we can locate both main and talk pages of the same article we save it for reference. We include in our corpus maximum number of articles for **GA**, **B**, **C**, **Start** and **Stub** while maintaining equality. For **FA** articles, we included entire extracted corpus, because such articles are scarce.

Although this process is naïve, yet we manage to extract moderately balanced number of datapoints for each class. We would further like to note that few other public datasets exits for this task. The first version is made available by the Wikimedia foundation[11] which has $30K$ datapoints with approximately $5k$ pages in each of the 6 classes. (Shen et al., 2017) has pointed out that this dataset contains many noisy datapoints, e.g., empty pages labeled as FA class. Other datasets are made available by (Shen et al., 2019a,c; Warncke-Wang et al., 2015). Our investigation shows that none of the former datasets contain meta pages which prompts us to extract this novel data ourselves. Also some of the datasets have uniform article length across classes, which is often not the case in reality. For example *Dodo* and *Grey-necked wood rail* are both FA articles with very different article length. Considering uniform article length may lead to overfitting. We used stratified random sample of $80\%$ data for training and $10\%$ each for validation and testing. The overall distribution of the articles with respective classes in our dataset is enumerated in Table 1.

## 4 Proposed solution

In this section we present a detailed description of our multimodal approach. We start this section by explaining the various notations that we use in the subsequent sections of the paper. We then proceed to first describe representation mechanisms of explicit signals about article quality, obtained from main text, i.e., presence of bias, claim verifiability, coherent wording, citations etc. We further elaborate on the mechanisms employed, to capture implicit quality indicators such as article stability, collaborative nature of editors obtained from article talk text as well as visual renderings of documents to capture how well the article follows the style manual. Further, we describe the different ways in which we combine the information to predict the quality of the page. We present the overall architecture of our model in Figure 1.

### 4.1 Preprocessing

We first pre-process the Wikipedia articles and talk pages to convert the text from Wiki Markup Language[12] to plain English text format using a python text crawler[13]. We then replace the meta content in

---

[9] https://dumps.wikimedia.org
[10] https://pypi.org/project/mwxml

[11] analytics.wikimedia.org/published/datasets
[12] en.wikipedia.org/wiki/Wikipedia:Wiki_Markup_Language
[13] github.com/attardi/wikiextractor

the page such as infobox, level 1 section headings, level 2 section heading, internal wikilink, external link, inline reference, footnote template, image template, quotation template and categories into special tokens which act as additional features using mediawiki parser[14]. We use this pre-processed text in the subsequent models.

## 4.2 NwQM overview

Our approach is inspired by the hierarchical document representation approach proposed in (Yang et al., 2016), designed to capture signals from multiple levels of document organization. Wikipedia pages have an organic structure, i.e., words form sentences, sentences form paragraphs, paragraphs form sections and sections form a page. We build page representation by first generating section embeddings, followed by a suitable summarization. We use BERT (Devlin et al., 2018) for section representation and bidirectional GRU with self attention for summarization. Explicitly our encoder has two levels in the hierarchy, however implicitly by employing bidirectional transformers (Vaswani et al., 2017) with special tokens, we can aggregate contextual features across multiple levels.

## 4.3 Fine tuning BERT

We fine tune the BERT (Devlin et al., 2018) model on the extracted Wikipedia pages for classifying them into the 6 quality categories. We pass the pre-processed textual content of the Wikipedia page to the BERT model and pass the [CLS] token's representation to a dense followed by softmax layers to classify the page. We follow the fine-tuning strategy proposed by (Sun et al., 2019) for long input sequences to overcome the limitations of BERT in handling inputs more than 512 tokens. We concatenate the first 128 and the last 384 tokens of the page. For the statistical observations on advantage of this approach please refer to (Sun et al., 2019). We use the BERT's tokenizer with additional meta content tokens explained in Section 4.1. We fine tune the model end-to-end (110M parameters). Once the model is finetuned we freeze the weights, thus in subsequent steps of the training the parameters are not updated.

## 4.4 BERT section encoder

We use previously obtained fine tuned BERT for generating section representations. It has been

shown that fine tuned BERT performs extremely well in subject verb agreement task (Goldberg, 2019; Clark et al., 2019), hence it is able to capture long range semantic dependency. Besides, since it is pre-trained with next sentence prediction objective, it remembers context across multiple sentences. Also since 95% of the sections have less than 512 words, we use the pre-trained BERT model with maximum number of input tokens, to directly encode the sections, hence *collapsing* the hierarchy starting from words upto sections in the encoding process. We collapsed the hierarchy up to sections but do not extend further due to several drawbacks in the BERT model discovered through probing (Liu et al., 2019; Si et al., 2019) which state that very long sequences spanning multiple sentences leads to incorrect comprehension. We take the final hidden state $h$ of the first token [CLS] as the aggregate representation of the section i.e. $\mathbf{Sc}_x$. (See Figure 1).

## 4.5 Conditional summarizer

We next generate the page representation ($\mathbf{D}_{p_i}$, see Figure 1) from the sequence of sections using attention based bidirectional GRU encoder as discussed in the previous section. The inputs $i_x$ in equations (1) to (6) correspond to the section representations obtained from the BERT models. The $\alpha_x$ correspond to the attention weights, the final output $o_n$ is the page representation $\mathbf{D}_{p_i}$ (see Figure 1). To aggregate the sequences at section level $Sc_x, x \in [1, n]$, we use the bidirectional GRU (Chung et al., 2014) which provides representations of the document text by summarizing information from both directions. We concatenate the forward and backward hidden states $h_x$ and feed the hidden states to the self attention module (Lin et al., 2017; Bahdanau et al., 2014).

$$\overrightarrow{h}_x = \overrightarrow{GRU}(i_x),\ x \in [1, L] \tag{1}$$

$$\overleftarrow{h}_x = \overleftarrow{GRU}(i_x),\ x \in [L, 1] \tag{2}$$

$$h_x = [\overrightarrow{h}_x, \overleftarrow{h}_x] \tag{3}$$

$$u_x = \sigma(W_i h_x + b_i) \tag{4}$$

$$\alpha_x = \frac{exp(u_x^T u_i)}{\sum_{x=1}^{L} exp(u_x^T u_i)} \tag{5}$$

$$o_n = \sum_x \alpha_x u_x \tag{6}$$

## 4.6 Talk page encoder

We split the talk page content into sentences and pass each sentence through the Google Universe
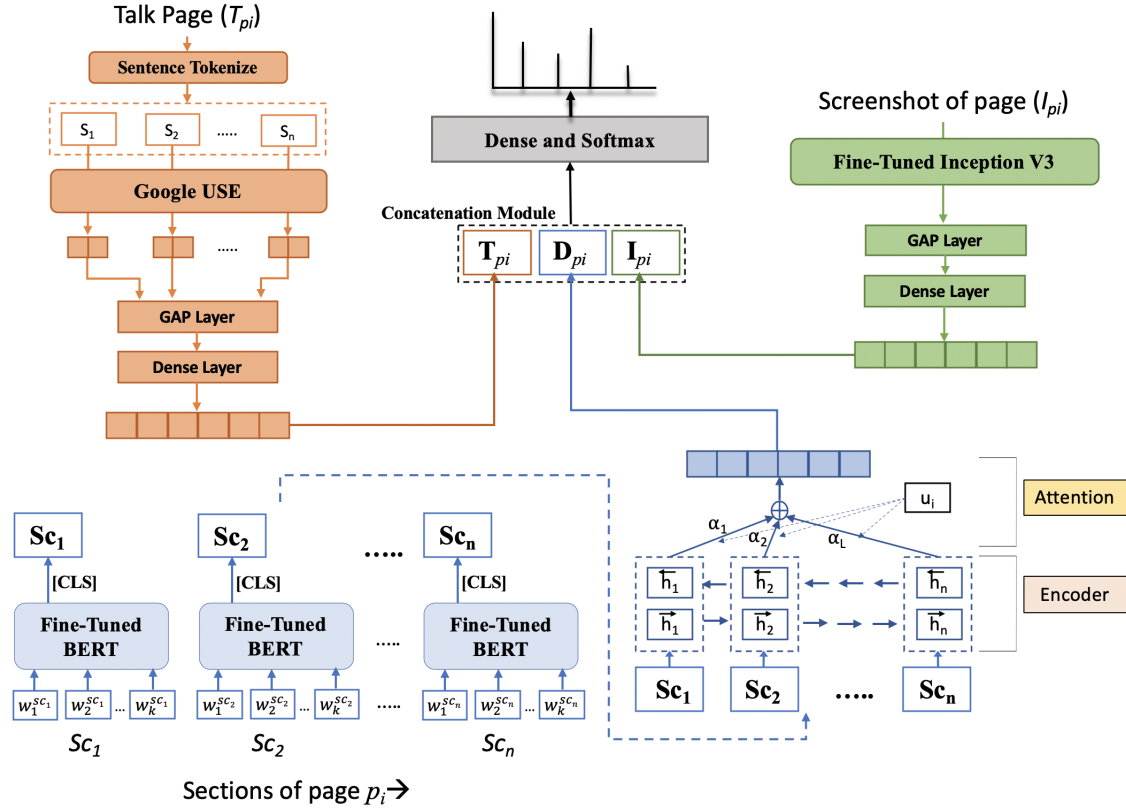
Figure 1: The overall pipeline of **NwQM**. We have taken cutoff for the number of tokens in sequence i.e $k = 512$ and no. of sequence $n = 16$. Parameters of fine-tuned BERT, Inception V3 are fixed.

Sentence Encoder model (Cer et al., 2018) to obtain a 512 dimensional representation of the sentence. These sentence embeddings are then passed through a global average pooling layer (Lin et al., 2013) followed by a dense layer to get a 200 dimensional final representation of the talk page ($\mathbf{T}_{p_i}$, see Figure 1). The motivation behind this approach is to obtain an aggregate representation of the discourse in talk pages.

### 4.7 Fine tuned Inception V3

Shen et al. (2019b) shows the effectiveness of using the visual rendering (screenshot image) of a page to predict the quality of a document in a multimodal setup. To embed the visual rendering we use the Inception V3 (Szegedy et al., 2016) model. Similar to the BERT model, to learn better representations, we fine tune the Inception V3 model. We follow the setting proposed in (Shen et al., 2019b) to fine-tune the Inception V3 model. We resize the images from varied high dimensions to a standard low dimension. We then fine tune the Inception v3 model to classify the screenshots of the pages into the 6 quality classes. We do an end-to-end training for fine tuning the Inception V3 model. We flatten the

final representation from the last convolution layer, stack the dropout and global average 2D pooling layers, and classify using softmax layer. The hidden flat layer output (2048 dimension) of the fine tuned Inception V3 model is the final representation of the visual rendering of the page. The input to the fine tuned Inception model is the screenshot of the wikipage, and the output is a visual embedding of the page ($\mathbf{I}_{p_i}$, see Figure 1). We integrate this information with the text and talk content of the page and report the improved performance of the final model.

### 4.8 Concatenation module

In this subsection we illustrate the concatenation of the information from various modules representing the features from text and image modalities to assess the quality of a page. We experiment with different modes of concatenation presented in (Reimers and Gurevych, 2019) $(u, v)$, $(u, v, |u - v|)$, $(u, v, |u - v|, u * v)$, $(u, v, u * v)$, $(|u-v|, u*v)$, $(|u-v|)$, $(u*v)$, and choose the best performing strategy $(u, v, |u - v|)$ to concatenate the vectors $u$ and $v$. We experiment with various combinations of modules, i.e., conditional docu-

ment summarizer (fine tuned BERT encoder with GRU summarizer), talk page encoder and visual rendering of the page (fine tuned Inception V3) and tabulate the best results in Table 2.

## 4.9 Model configurations

We set the hidden states $(\overrightarrow{h_x}, \overleftarrow{h_x})$ of the GRU units of the section encoder to 100. Therefore the final page representation $(\mathbf{D}_{p_i})$ is a 200 dimensional vector. Since 90% of the data has number of sections in a page less than 16 respectively, we limit their maximum size of a page to 16 sections. For sequences less than the specified length, we left-pad using $\overrightarrow{0}$. We load the pre-trained weights of BERT-base model from TensorflowHub[15].

We use the implementation of Google Universal Sentence Encoder available at TensorFlow Hub[16]. We use the nltk library[17] to tokenize the pre-processed talk page content into sentences. We train the hierarchical content encoder for 10 epochs with a learning rate of 0.001 and batches of 16 using Adam optimizer (Kingma and Ba, 2014). For fine tuning the BERT models, we use the Adam optimizer with a learning rate of 2e-5. We empirically set the number of training epochs to 4. To fine tune the Inception V3 model, we again use Adam optimizer with a learning rate 1e-4 and train for 20 epochs. We employ the categorical cross-entropy as loss function for all the models and train using batches of size 16. For all the joint models, we set the learning rate to 0.001, batch size to 32, number of epochs for training to 40. For classification, we use dense layers followed by softmax layer. We further utilize dropout probability of 0.5 in the dense layers. Prior fine-tuning of individual units reduces explosive training time, common in end-to-end models

## 5  Experiment

In this section we evaluate **NwQM** against several existing approaches.

**Baselines.**   We provide a brief outline of the competing methods in the following.

- **ORES** (Halfaker and Geiger, 2019) is a machine learning service made available by Wikimedia foundation through a RESTful HTTP interface serving prediction about target articles. It uses handcrafted features along with gradient boosted machine as classifier.
- **DOC2VEC** (Dang and Ignat, 2016) proposed the first application of deep neural networks into quality assessment task where they employed distributional representation of documents (Le and Mikolov, 2014) without using manual features.
- **BILSTM+** (Shen et al., 2017) is a hybrid model, where textual content of the Wikipedia articles are encoded using a BILSTM model. The hidden representation captured by the sequence model is further augmented with handcrafted features and the concatenated feature vector is used for final classification
- **H-LSTM** (Zhang et al., 2018) is an edit history based approach where every version of an article is represented by 17 dimensional handcrafted features. Hence an acticle with $k$ versions will be represented by $k \times 17$ matrix. This $k$ length sequence is passed through a stacked LSTM for final representation used in classification.
- **M-BILSTM** (Shen et al., 2019b) proposed a multimodal information fusion approach where embeddings obtained from both article text as well as html rendering of the article webpage is used for final classification.
- **DOCBERT** (Adhikari et al., 2019) proposed this method to generate document representations using bidirectional transformers (Devlin et al., 2018). The primary idea is filtering the representation obtained from the CLS token using a fully connected layer which translates 768 dimensional encoding to class distribution using a softmax layer. This architecture is further fine tuned end-to-end for the respective document classification task.
- **HAN** (Yang et al., 2016) proposed a hierarchical approach which iteratively constructs a document vector by coalescing important words into sentence vectors and subsequently combining important sentences vectors to obtain the document vectors. A convenient consequence of this approach is that it is suitable for large documents like Wikipedia articles.

**Result**   We evaluate **NwQM** against existing solutions for automatic quality assessment and tabulate the obtained results in Table 2. Since our classes are roughly balanced, we opt to report accuracy as the metric for evaluation. Some of the

---

[15]https://tfhub.dev/google/bert_uncased_L-12_H-768_A-12/1
[16]https://tfhub.dev/google/universal-sentence-encoder/
[17]https://www.nltk.org/api/nltk.tokenize.html

approaches that we compare our model with are ORES, DOC2VEC, BILSTM+, M-LSTM, H-LSTM. We achieve an improvement of 8% compared to the best performing baseline. This we believe is a considerable leap for a 6 class classification task.

We also compare our model against novel document classification approaches, i.e., DOCBERT and HAN because inherently Wikipedia quality assessment problem is closely related to document categorization. To this purpose, we compare the document classification approaches with textual content of the Wikipedia main article as well as concatenated version of the main article and talk pages denoted as DOCBERT-wT, HAN-wT respectively. We obtain at most 5% improvement against the existing approaches. Note that apart from our model, including the talk page meta data in the document classification models also considerably enhances their respective performances.

**NwQM** is constituted of concatenated representation from article text, talk and image and therefore it is important to look at how individual components perform independently. We evaluate **NwQM** without signals from image (**NwQM**-w/oI), without talk (**NwQM**-w/oT) and with solely the main article text, i.e., without any secondary and tertiary signals from image and metadata (**NwQM**-w/oTI). Our experiment show that the combined approach (**NwQM**) obtains the best result. Without talk and image we land in a drop of accuracy of 1.3% and 5% respectively validating our hypothesis that extracting signals from external sources can serve fruitful in this task. We also compare against representation from talk pages and fine tuned image embeddings individually for the sake of completeness and our results show significant drop in accuracy compared to the combined approach.
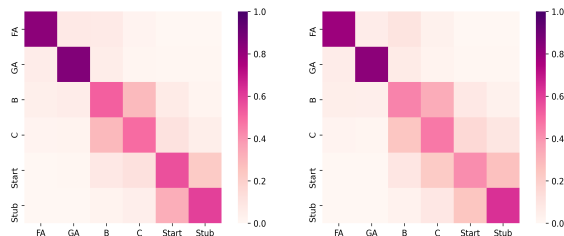
| Model | Accuracy |
|---|---|
| ORES (Halfaker and Geiger, 2019) | 43.21 |
| DOC2VEC wRF (Dang and Ignat, 2016) | 44.01 |
| DOC2VEC wLR (Dang and Ignat, 2016) | 49.33 |
| BILSTM+ (Shen et al., 2017) | 54.5 |
| H-LSTM (Zhang et al., 2018) | 53.05 |
| M-BILSTM (Shen et al., 2019b) | 58.47 |
| DOCBERT (Adhikari et al., 2019) | 57.66 |
| DOCBERT-wT (Adhikari et al., 2019) | 59.87 |
| HAN (Yang et al., 2016) | 56.35 |
| HAN-wT (Yang et al., 2016) | 57.48 |
| NwQM | 63.23 |
| NwQM-w/oI | 59.95 |
| NwQM-w/oT | 62.37 |
| NwQM-w/oTI | 59.10 |
| TALK | 37.95 |
| INCEPTION V3 | 52.96 |

Table 2: Results obtained from different models. The best result is highlighted in green, the best result among document classification models is highlighted in red and the best result among the state-of-the-art quality assessment models is highlighted in blue.
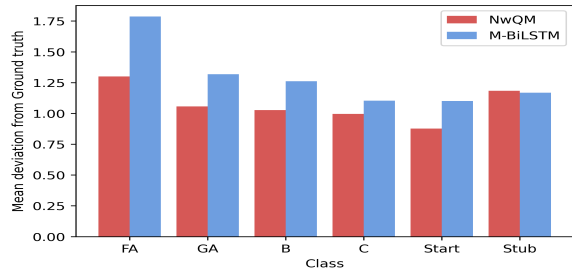


Figure 3: Mean absolute distance of the misclassifications from the true labels for individual classes. (Best viewed in color)

## 6 Discussion and qualitative analysis

In this section we perform a deep dive into the predictions obtained by our model and we contrast it with those from closest known competitor M-BILSTM (Shen et al., 2019b) for quality assessment. **Confusion matrix**: We first tabulate the confusion matrix obtained by **NwQM** on test data in Figure 2. Results show that most of the misclassifactions made by **NwQM** are on average between very similar classes. This is due to the inherent ordinal nature of the classes in this dataset. More specifically, there is an increasing degree of quality in pages from Stub to FA, though the gradient may not be smooth. Thus FA, GA articles have overlapping guidelines which is very different from other pages. Likewise B, C and Start, Stub have mutually overlapping guidelines. **NwQM** can successfully recover this structure, hence misclassifications have occurred between closer classes. However, com-



Figure 2: Left panel shows confusion matrix obtained by **NwQM**; Right panel shows confusion matrix obtained by M-BILSTM (Best viewed in color) .
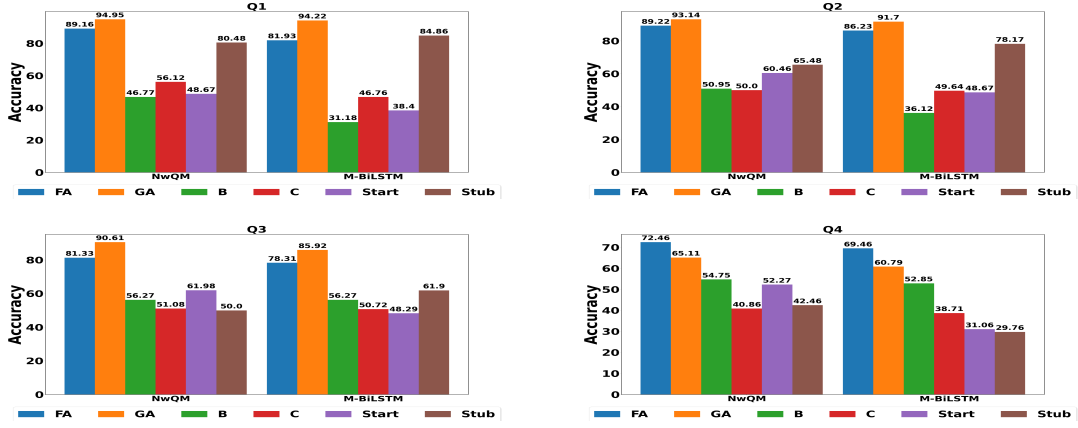
Figure 4: Accuracy per class for test data segregated into quartiles $Q_1, Q_2, Q_3, Q_4$ with respect to main article text length with $Q_1$ smallest and $Q_4$ largest. (Best viewed in color)
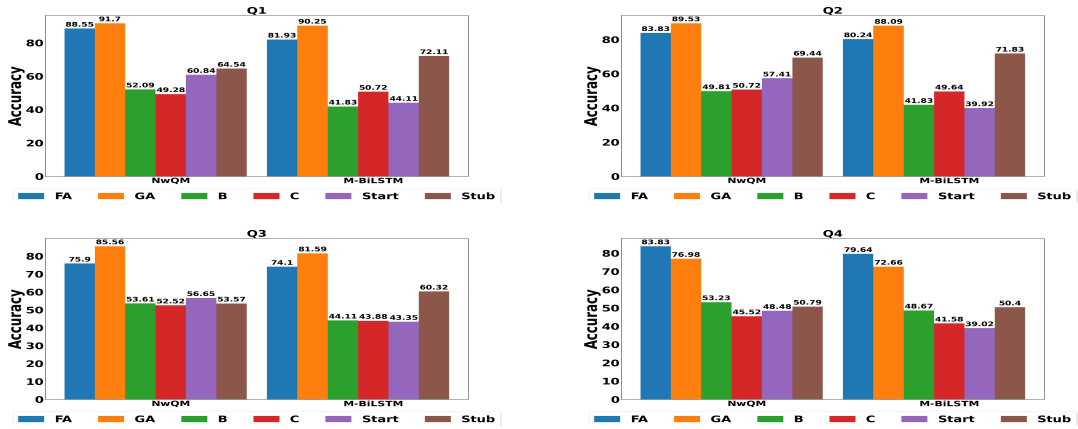


Figure 5: Accuracy per class for test data segregated into quartiles $Q_1, Q_2, Q_3, Q_4$ with respect to article talk length with $Q_1$ smallest and $Q_4$ largest. (Best viewed in color)

pared to M-BILSTM (see Figure 2) **NwQM** can capture class specific features significantly better thus showing lower mistakes for certain datapoints especially in case of Start, B and C classes.

**Distance of wrong predictions from ground truth**: In cases of wrong predictions, we calculate how far individual models are from ground truth labels. We transform each class into integers with Stub transformed to 0 and FA transformed to 5. We then find the mean absolute distance of the incorrect predictions from the true labels for each individual class. The results are tabulated in Figure 3 where we observe that for **NwQM** the absolute distance is consistently lower. The results are statistically significant with a $p$-value of 1e-5 using Stuart-Maxwell test for multiclass classification (Sun and Yang, 2008).

**Effect of article length**: We further ascertain how the main text and talk page length of an plays a role

in prediction, for respective models. Articles with detailed coverage of topics and multiple discussion threads may not be nominated among high quality articles. This is because even if the coverage is diverse, the language of the article may be biased, convoluted, there could be unverified claims or disaccord among editors. Deep learning models such as BERT have often been shown to rely on surface forms (Ettinger, 2020) as shotcuts for classification instead of semantic understanding. We investigate whether **NwQM** is relying spurious signals like article length in final classification. For example, article *Kauri Gum* (GA) and *Guar Gum* (Start) have very similar coverage, but are distant in terms of the quality criteria. Similarly, article *Frog cake* (GA) and *Sugar* (GA) have very different coverage but the same quality tag. For further investigation we rank predictions with respect to article and talk page length and divide our predictions into 4 quan-

tiles. We investigate the accuracy of **NwQM** in the individual quantiles starting from the smallest quantile, i.e., $Q_1$ to the largest, i.e., $Q_4$. Accuracy scores for test set ordered by main article text are illustrated in Figure 4 and those ordered based on talk pages are shown in Figure 5. Results indicate that **NwQM** is not biased toward length and irrespective of the quantiles we obtain improved prediction accuracy almost always compared to the closest baseline.

**t-SNE plots**: Finally, we visualize the representations obtained by **NwQM** and M-BILSTM using t-SNE (Maaten and Hinton, 2008) scatter plots (see Figure 6). The degree of separation obtained by **NwQM** is much better which translates to the improved accuracy.
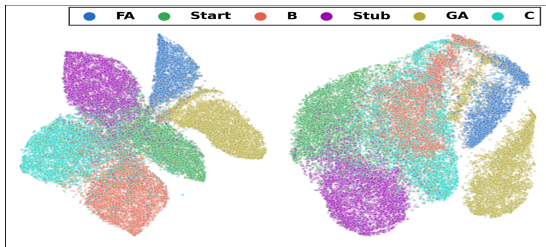


Figure 6: Left panel shows t-SNE visualization for **NwQM**; right panel shows visualization for M-BISTM. (Best viewed in color)

**Interpreting the effect of the different modalities**: One of the critical issues in deep learning models is interpretability. In order to ascertain whether different sources of signals are indeed contributing toward the final prediction task, we leverage the model agnostic evaluation tool LIME (Ribeiro et al., 2016). We generate representations for text, talk and images, i.e., $D_{p_i}, T_{p_i}, I_{p_i}$ respectively (see Figure 1) from our learned model. These embeddings are passed through classification layers comprising dense layer and softmax layer for prediction. This network takes concatenated input of $D_{p_i}, T_{p_i}, I_{p_i}$ and outputs the classification probabilities for each test instance. We evaluate this black box neural network using LIME. For data points on the test set we identify top 500 features contributing to the outcome of the highest class probability. We further calculate the average contribution from each modality toward the respective classes. More specifically for every test page instance, we identify the top 500 contributing features as per LIME. Each feature can be contributed by any one of the three modalities. We compute the mean of features scores per modality in the top 500 for a page and

then aggregate that over all pages. The results are tabulate in Figure 7. Our results show that signals from $D_{p_i}$ play most important overall role in prediction of quality. $I_{p_i}$ play almost equal role in prediction of all the classes. Interestingly, for high quality pages, $T_{p_i}$, i.e., the talk page information contributes in classification higher than text and image information. Talk also contributes for low quality pages such as *Start, Stub*. We speculate that since high quality pages have larger discussion archives and low quality pages have very low discussion threads $T_{p_i}$ plays significant role in distinguishing these classes. We stress that these interpretability results are the prime insights that these paper neatly establishes.
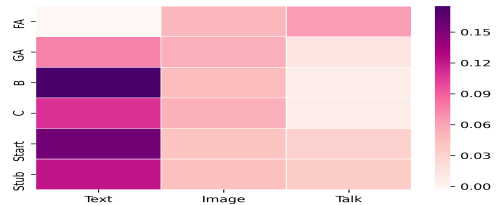


Figure 7: (Best viewed in color)

## 7 Conclusion

In this paper we proposed a novel multimodal deep learning based model **NwQM** for quality assessment of English Wikipedia articles. Our model combines signals from article text, meta pages and image rendering to construct an improved document representation. We evaluate it against several existing approaches and obtain at most $8\%$ improvement compared to the state-of-the-art method. For a 6 class classification this leap in accuracy is notable. We also perform extensive investigation of the different components of our model to understand their individual utility. We perform in-depth qualitative analysis of the obtained predictions and contrast them with the closest baseline.

To the best of our knowledge this is the first work which combines several aspects of information available for Wikipedia articles and, in particular, the talk page dynamics toward quality assessment. We also showcase the utility of fine tuned bidirectional transformers toward document classification especially when combined with niche platform specific signals. We believe our work opens up the necessity of further investigation pertaining to careful information fusion techniques for downstream tasks.

# References

Ashutosh Adhikari, Achyudh Ram, Raphael Tang, and Jimmy Lin. 2019. Docbert: Bert for document classification. *arXiv preprint arXiv:1904.08398*.

David C Atkinson. 2018. Charlottesville and the alt-right: a turning point? *Politics, Groups, and Identities*, 6(2):309–315.

Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.

Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, et al. 2018. Universal sentence encoder. *arXiv preprint arXiv:1803.11175*.

Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*.

Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. 2019. What does bert look at? an analysis of bert's attention. *arXiv preprint arXiv:1906.04341*.

Quang Vinh Dang and Claudia-Lavinia Ignat. 2016. Quality assessment of wikipedia articles without feature engineering. In *Proceedings of the 16th ACM/IEEE-CS on Joint Conference on Digital Libraries*, pages 27–30.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Allyson Ettinger. 2020. What bert is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Transactions of the Association for Computational Linguistics*, 8:34–48.

Yoav Goldberg. 2019. Assessing bert's syntactic abilities. *arXiv preprint arXiv:1901.05287*.

Aaron Halfaker and R Stuart Geiger. 2019. Ores: Lowering barriers with participatory machine learning in wikipedia. *arXiv preprint arXiv:1909.05189*.

Aaron Halfaker and Dario Taraborelli. 2015. Artificial intelligence service "ores" gives wikipedians x-ray specs to see through bad edits. *Wikimedia Blog*.

Manoel Horta Ribeiro, Kristina Gligorić, Maxime Peyrard, Florian Lemmerich, Markus Strohmaier, and Robert West. 2020. Sudden attention shifts on wikipedia following covid-19 mobility restrictions. *arXiv*, pages arXiv–2005.

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.

Baptiste de La Robertie, Yoann Pitarch, and Olivier Teste. 2015. Measuring article quality in wikipedia using the collaboration network. In *2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, pages 464–471. IEEE.

Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *International conference on machine learning*, pages 1188–1196.

Xinyi Li, Jintao Tang, Ting Wang, Zhunchen Luo, and Maarten De Rijke. 2015. Automatically assessing wikipedia article quality by exploiting article–editor networks. In *European Conference on Information Retrieval*, pages 574–580. Springer.

Min Lin, Qiang Chen, and Shuicheng Yan. 2013. Network in network. *arXiv preprint arXiv:1312.4400*.

Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. 2017. A structured self-attentive sentence embedding. *arXiv preprint arXiv:1703.03130*.

Nelson F Liu, Matt Gardner, Yonatan Belinkov, Matthew Peters, and Noah A Smith. 2019. Linguistic knowledge and transferability of contextual representations. *arXiv preprint arXiv:1903.08855*.

Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605.

Malte Ostendorff, Peter Bourgonje, Maria Berger, Julian Moreno-Schneider, Georg Rehm, and Bela Gipp. 2019. Enriching bert with knowledge graph embeddings for document classification. *arXiv preprint arXiv:1909.08402*.

Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv preprint arXiv:1908.10084*.

Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. " why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1135–1144.

Aili Shen, Daniel Beck, Bahar Salehi, Jianzhong Qi, and Timothy Baldwin. 2019a. Modelling uncertainty in collaborative document quality assessment. In *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*, pages 191–201.

Aili Shen, Jianzhong Qi, and Timothy Baldwin. 2017. A hybrid model for quality assessment of wikipedia articles. In *Proceedings of the Australasian Language Technology Association Workshop 2017*, pages 43–52.

Aili Shen, Bahar Salehi, Timothy Baldwin, and Jianzhong Qi. 2019b. A joint model for multimodal document quality assessment. In *2019 ACM/IEEE Joint Conference on Digital Libraries (JCDL)*, pages 107–110. IEEE.

Aili Shen, Bahar Salehi, Jianzhong Qi, and Timothy Baldwin. 2019c. Feature-guided neural model training for supervised document representation learning. In *Proceedings of the The 17th Annual Workshop of the Australasian Language Technology Association*, pages 47–51.

Chenglei Si, Shuohang Wang, Min-Yen Kan, and Jing Jiang. 2019. What does bert learn from multiple-choice reading comprehension datasets? *arXiv preprint arXiv:1910.12391*.

Chi Sun, Xipeng Qiu, Yige Xu, and Xuanjing Huang. 2019. How to fine-tune bert for text classification? In *China National Conference on Chinese Computational Linguistics*, pages 194–206. Springer.

Xuezheng Sun and Zhao Yang. 2008. Generalized mcnemar's test for homogeneity of the marginal distributions. In *SAS Global forum*, volume 382, pages 1–10.

Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. 2016. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.

Morten Warncke-Wang, Vladislav R Ayukaev, Brent Hecht, and Loren G Terveen. 2015. The success and failure of quality improvement projects in peer production communities. In *Proceedings of the 18th ACM Conference on Computer Supported Cooperative Work & Social Computing*, pages 743–756.

Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. 2016. 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*, pages 1480–1489.

Shiyue Zhang, Zheng Hu, Chunhong Zhang, and Ke Yu. 2018. History-based article quality assessment on wikipedia. In *2018 IEEE International Conference on Big Data and Smart Computing (BigComp)*, pages 1–8. IEEE.