

# A Comprehensive Analysis of the Quantum-like Approach for Integrating Syntactic and Semantic Information

Anonymous ACL submission

## Abstract

Transformers have proved effectiveness in understanding and deciphering the intricate context of languages. This success is achieved by those models that lack explicit modeling of syntactic structures, which were hypothesised by decades of computational linguistic research to be necessary for logical text understanding. In this work, we present a comprehensive analysis of syntactic and semantic context integration by proposing Compressed Phrase Embedding and adopting quantum-like methods for text classification. We first introduce Compressed Phrase Embedding (ComPhE) by integrating syntactic parsing and semantic contextual information. We test those with two types of quantum-like approaches, 1) quantum-like input processing (DisCoWord) and 2) quantum-like attention (QSA), and discuss the contribution of compressed phrase syntactic and semantic integration towards the model performance on different text classification benchmarks.

## 1 Introduction

Transformer has dramatically enhanced the performance of tasks in NLP, but they only focus on capturing the contextual semantics of the language, even though the syntactic structure of the sentences are also an essential aspect of a language. Previous studies (Clark et al., 2019; Htut et al., 2019; Mareček and Rosa, 2019; Rogers et al., 2020; Sharma et al., 2022; Zheng and Liu, 2023; Ma et al., 2023) have shown that pretrained transformer-based models like BERT can implicitly capture partial syntactic information. However, this substantial information is a by-product of learning language semantics via self-supervised learning, not by explicitly trying to understand the syntactic structure. Several recent studies (Hu et al., 2020; Du et al., 2023; Zhu et al., 2022b; Li et al., 2022; Zhang and Li, 2022) try to make use of Graph-based Neural Networks (GNNs) to inject such syntactic information into Transformer

explicitly. However, GNN research can lead to higher time and space complexity. Few works in machine translation apply the simple token grouping (Bharadwaj and Shevade, 2022) and attention masks (Hou et al., 2022), and there is still room for improvement. Hence, how to incorporate the syntactic and semantic aspects effectively into transformers is also unsettled. Recently, gradient-based training of quantum circuits has been successfully adopted to generate joint distributions over multiple aspects and variables (Delgado and Hamilton, 2022; Zhu et al., 2022a).

This paper presents a comprehensive analysis of syntactic and semantic context integration by proposing Compressed Phrase Embedding and adopting quantum-like methods for text classification. At first, we introduce Compressed Phrase Embedding (ComPhE) by integrating syntactic parsing and semantic contextual information. To integrate the syntactic and semantic aspects of inputs into transformers effectively, we then apply and test two types of quantum-like approaches: quantum-like input processing and quantum-like attention. Firstly, for the quantum-like input processing, we adopt the concept from DisCoCat, which is a computing framework for compositional sentence meaning (Clark et al., 2008; O’Riordan et al., 2020; Lorenz et al., 2023; Kartsaklis et al., 2021). We introduce a new quantum-based representation called DisCoWord to help ComPhE better gather semantic contextual information from the syntactic aspect. Secondly, for the quantum-like attention, inspired by the available quantum modelling of transformer structure (Di Sipio et al., 2022; Cherrat et al., 2022; Shi et al., 2023), we propose a Quantum Self-Attention (QSA), which uses the measurements of the quantum qubits to compute the attention scores. The main contributions are as follows: We investigate the benefit of quantum-like approaches for syntactic and semantic integration by proposing the ComPhE, and adopting quantum-like components

083 on several text classification benchmarks. Note  
084 that we introduce all new components, DisCoWord  
085 and QSA, to efficiently apply to text classifications,  
086 inspired by the existing quantum-like models.

## 087 2 Method<sup>1</sup>

088 We propose Compressed Phrase Embedding (Com-  
089 PhE), which merges the token embeddings within  
090 each phrase after splitting the input sentences into  
091 phrases using constituency parsing. This approach  
092 aims at syntactic and semantic integration, i.e.,  
093 gathering semantic contextual information accord-  
094 ing to syntactic parsing. Then, we apply ComPhE  
095 with two quantisation approaches: Firstly, we aug-  
096 ment the vanilla token embeddings by adding Dis-  
097 CoWord, another word embedding that uses quan-  
098 tisation to combine both syntactic and semantic  
099 context information. Secondly, we adopt quantisa-  
100 tion to improve the efficiency of the self-attention  
101 module. Specifically, we use quantisation in the at-  
102 tention score calculation to encode  $2^n$  dimensional  
103  $Q/K$  vectors into  $n$  qubits quantum states and use  
104 quantum states to model the semantics. We adopt  
105 the joint learning ability of quantum circuits since  
106 those have been adopted to generate joint distribu-  
107 tions over multiple aspects and variables.

### 108 2.1 Compressed Phrase Embedding

109 We propose Compressed Phrase Embedding (Com-  
110 PhE) that utilises constituency parsing to gather  
111 semantic contextual information from the syntactic  
112 aspect. Based on the constituent tree of the in-  
113 puts, we design two phrasing methods for splitting  
114 the words into phrases. The first is Top-to-Bottom  
115 (T2B), which generates phrases whose tag range is  
116 all Penn Treebank II Constituent Tags. The second  
117 is Bottom-to-Top (B2T), which only generates NP,  
118 PP, and VP phrases. In detail, T2B goes through the  
119 tree except for leaf nodes in a level order traversal  
120 from the second to the bottom layer and checks the  
121 subtree starting from each no-leaf node. If there is  
122 no non-root node with the same tag as the subtree’s  
123 root, the subtree is determined as a phrase. B2T  
124 goes through the tree except for leaf nodes in a  
125 level order traversal starting from the bottom to the  
126 second layer, splitting all of the minimal NP and  
127 PP. All of the other words will be seen as the VP  
128 phrases. In addition, both methods merge the con-  
129 secutive single words as a new phrase. We compare

<sup>1</sup>The Overview of Model and Test Architecture can be found in Appendix A.1

130 the two phrasing methods to demonstrate which  
131 syntactic aspects are more critical to gathering se-  
132 mantic contextual information. After splitting, we  
133 apply stemming and stopword removal in phrases.  
134 Then we tokenise the words in phrases. Finally, we  
135 obtain ComPhE by summing the token embeddings  
136 phrase-wise.

### 137 2.2 DisCoWord Representation

138 We introduce a new quantum-based representation  
139 called DisCoWord to help ComPhE better gather  
140 semantic contextual information from the syntac-  
141 tic aspect. DisCoWord can represent the seman-  
142 tic meanings of a word according to its contex-  
143 tual grammar structure, i.e., the pregroup grammar  
144 (Lambek, 1999) of the word in a context. This  
145 approach is inspired by DisCoCat (Lorenz et al.,  
146 2023), a distributional compositional categorical  
147 model which uses pregroup grammar to compute  
148 the meaning of words in a quantum way and can  
149 learn the grammar-based contextual representations  
150 for words. Based on DisCoCat, we use pretrained  
151 word embeddings to initialise the parameters of  
152 the quantum circuit representing the word, which  
153 make the representations hold both the syntactic  
154 and semantic information based on their contex-  
155 tual grammar structure and pretrained word embed-  
156 dings. See more details in Appendix A.2.

### 157 2.3 Quantum Self-Attention

158 We propose Quantum Self-Attention (QSA), which  
159 uses the measurements of the quantum qubits to  
160 compute the attention scores, instead of using dot-  
161 product between  $Q$  and  $K$  in a head. In more  
162 detail, we design a quantum circuit containing  $n$   
163 input qubits and  $m$  output qubits to help with that.  
164 First, we use a quantum feature map to transform  
165 each  $Q/K$  vector belonging to a head to a quantum  
166 state of the  $n$  qubits. After that, we apply param-  
167 eterised quantum gates on the  $n$  input qubits and  
168 the  $m$  output qubits, of which the parameters will  
169 be trained in the training process. In the end, we  
170 use measurement values of the  $m$  output qubits for  
171 each  $Q/K$  vector to compute the attention score.

## 172 3 Experiment Setup

173 We articulate how to evaluate our ComPhE with  
174 DisCoWord and QSA on text classification.

### 175 3.1 Datasets

176 We use four widely used text classification bench-  
177 mark datasets, including Movie Reviews (MR)

Statistics	MR	Twitter	SST-2	OffensEval
Split	7108/3554	8000/2000	6920/1821	11916/1324
Doc.	10.78/4.96	5.32/3.21	9.91/5.01	8.57/6.14
Emb.	9.94/26.71	7.20/21.70	9.87/26.82	7.84/19.64
Dropped	172	913	128	573

Table 1: The summary statistics of datasets. The split represents the Train/Test split. Doc. and Emb. represent document and embedding length, with average/standard deviation values. We dropped the documents which failed to be transformed into diagrams or with the word having lengthy states ( longer than 256 ).

(Pang and Lee, 2005), OffensEval (Zampieri et al., 2019), SST-2 (Socher et al., 2013) and Twitter<sup>2</sup> in our experiments. These datasets are binary classified, and their text lengths are short enough for running quantum computing on classical computers with the quantum simulation software. All the datasets are preprocessed by lowering case, stemming, removing punctuation and removing stopwords. The statistics of the four datasets and the DisCoWord embeddings for each dataset are presented in Table 1.

### 3.2 DisCoWord Training

For DisCoWord training, we use BobcatParser<sup>3</sup>, the state-of-the-art statistical Combinatory Categorical Grammar (CCG) parser (Clark, 2021). We choose the same ansatz (Hadfield et al., 2019) for each grammar type (Lambek, 2008) used in (Lorenz et al., 2023) with 1 qubit for resource saving. Then we use the Simultaneous Perturbation Stochastic Approximation (SPSA) (Spall, 1998) algorithm and CrossEntropy loss for optimisation. The SPSA algorithm is an efficient gradient approximation method only using the value of the object function, i.e., applying the random perturbation on the parameters and calculating the approximated gradient. Thus, it can deal with quantum simulation optimization, which is challenging to calculate gradients directly and usually has noise. The hyperparameters of the SPSA algorithm are referred to as the one from (Spall, 1998). As for pretrained word embeddings, we use glove-wiki-gigaword-50, glove-twitter-25, fasttext-wiki-news-subwords-300, and fasttext-twitter-100 (Camacho-Collados et al.) to initialise the parameters of the quantum circuit representing the word and choose the best DisCoWord according to the test accuracy. See more details in Appendix B.2.

<sup>2</sup>A built-in dataset in NLTK (Bird et al., 2009) library.

<sup>3</sup>we use the BobcatParser implemented in lambeq (Kartsaklis et al.) library.

### 3.3 Quantum Self-Attention Design

For the quantum circuit we used in the self-attention module, the number of input qubits  $n$  is determined by the number of attention heads, the dimension of  $Q/K$  vectors, and the future map. In our experiments, we use AmplitudeEmbedding feature map (Jaeger, 2007; Möttönen et al., 2005) to encode the  $2^n$  dimensional  $Q/K$  vector to the quantum state of the  $n$  input qubits. To save hardware resources, the number of output qubits  $m$  is set to 1, i.e., the last input qubit is set as the output qubit. In that case, each  $Q/K$  vector will be transformed into a measurement value. We then apply the block including a  $R_X(\theta)$  gate, a  $R_Y(\theta)$  gate and a  $CNOT$  gate on every two qubits from top to bottom, where  $\theta$  is the randomly initialised trainable parameter. Finally, we apply the  $Pauli - Z$  gate (DiVincenzo, 1998) as the measurement operator on the output qubit. The details of the quantum circuit can be found in Appendix B.3.

### 3.4 Method Verification

We verify our method in the text classification task on the four aforementioned datasets. We use Self-Attentive Encoder (SAE) constituency parser proposed by (Kitaev and Klein, 2018) for ComPhE. We use AdamW (Loshchilov and Hutter) optimiser and CrossEntropy loss in training. We use early stopping to get the test accuracy as the evaluation metric. We have three variations for testing: ComPhE (Vanilla), which applies ComPhE on Vanilla Transformer. ComPhE (DisCoWord), which applies ComPhE on Vanilla Transformer with DisCoWord augmented input. ComPhE (QSA), which applies ComPhE on the QSA augmented Transformer with Vanilla Transformer input.

## 4 Results

### 4.1 Performance Evaluation

We compare the performance of baselines and our ComPhE with variations on four text classification datasets. From Table 2, we can see that our ComPhE variations are better than all baselines, which proves the ability to improve text classification. More specifically, using ComPhE (Vanilla) is better than the Vanilla Transformer and most baselines, demonstrating the effectiveness of integrating semantic context information from the syntactic aspect. ComPhE (DisCoWord) does not improve the accuracy except on the Tweet dataset.

Methods	MR	Tweet	SST-2	OffensEval
TFIDF+LR	75.5	68.4	80.1	77.6
CNN-Rand	70.9	99.1	75.3	72.7
CNN-Pretrained	72.0	93.7	74.9	70.0
LSTM-Rand	66.4	92.5	67.8	62.5
LSTM-Pretrained	71.5	87.0	69.6	64.4
Vanilla Transformer	74.8	99.8	73.2	76.9
ComPhE (Vanilla)	75.1	99.8	80.1	<b>78.0</b>
ComPhE (DisCoWord)	74.1	<b>99.9</b>	78.5	77.9
ComPhE (QSA)	<b>75.9</b>	<b>99.9</b>	<b>80.8</b>	73.7

Table 2: Overall performance comparison with the baselines and the ComPhE variations. The ComPhE variations are better than the baselines overall; TFIDF+LR and CNN-Rand are competitive. Note we mainly focus on comparing with Vanilla Transformer and other classical methods to demonstrate the effectiveness of our approach. We mainly focus on Transformer variants because our proposed methods can be added to other Transformer-based models.

We believe the reason is that the quality of DisCoWord is not good because we use as few qubits as possible and split the training process due to the limited hardware resources. ComPhE (QSA) gets the best performance, 75.9 in MR, 99.9 in Tweet and 80.8 in SST-2. We draw the following inferences from the results. **a)** Evolution to quantum states can replace dot-product self-attention. **b)** OffensEval aims to identify offensive documents from English tweets, which may place a higher demand on text comprehension. Therefore, the simple quantum circuit used in the experiment does not convert the classical vectors of this dataset into quantum states well. As for the Tweet dataset, whose average document length is the shortest, it may be too easy to learn by the model, resulting in the close performances from the three ComPhE variations.

## 4.2 Compressed Phrasing Analysis

To better understand the impact of the phrasing methods after constituency parsing, we apply our ComPhE with T2B and B2T phrasing methods, respectively. Table 3 shows that both B2T and T2B methods outperform Vanilla Transformer on the four datasets, and are very close while B2T is slightly better. However, when using T2B to get ComPhE, the average token length of the input is 17.6% to 31.4% less than when using B2T<sup>4</sup>. Therefore, we believe using the T2B method will be more advantageous in processing long text.

<sup>4</sup>The details of input token length statistics can be found in Appendix B.4

Method	MR	Tweet	SST-2	OffensEval
w/o-phrasing	74.8	<b>99.8</b>	73.2	76.9
ComPhE (B2T)	<b>75.1</b>	<b>99.8</b>	<b>80.1</b>	77.5
ComPhE (T2B)	74.8	<b>99.8</b>	79.5	<b>78.0</b>

Table 3: Performance comparison between ComPhE test results using different phrasing methods. w/o-phrasing represents the Vanilla Transformer.

## 4.3 Positional Encoding Analysis

We also test the impact of positional embedding for ComPhE input handling. As can be seen in Table 4, Rel. PE is better on the OffensEval dataset, while Abs. PE (sum) is better on MR and SST-2 datasets. We believe that this discrepancy is caused by differences in the datasets. Offenseval’s data is derived from Twitter, so its text is less rigorous than that of MR and SST-2. Therefore, the relative positional embedding, added at the attention layer and can supply phrase position information, indicates that phrase position information is more important than token position information in the Twitter text. As for the absolute positional embedding, which is added to the token embeddings and influences the gathering process and the quality of the ComPhE, is better than relative positional embedding on the MR and SST-2 datasets. Since the Tweet is the dataset with the shortest average document length, positional embedding may be useless. Therefore, these positional embeddings’ performances are close on the Tweet dataset.

PE	MR	Tweet	SST-2	OffensEval
w/o-PE	72.7	<b>99.9</b>	77.9	78.2
Abs. PE (sum)	<b>75.1</b>	99.8	<b>80.1</b>	78.0
Abs. PE (cat)	72.2	<b>99.9</b>	76.8	77.3
Rel. PE	74.4	99.8	78.7	<b>78.9</b>

Table 4: The positional embedding analysis results. All the experiments use ComPhE (Vanilla).

## 5 Conclusion

We propose Compressed Phrase Embedding, which integrates syntactic parsing and semantic contextual information, and apply it with DisCoWord and QSA on four text classification tasks. Our ComPhE, which gathers semantic contextual information using syntactic parsing, can help understand the text. Both DisCoWord and QSA can enhance the performance of ComPhE. Hence, we hope that ComPhE with quantum-like approaches will be a good reference for integrating syntactic and semantic information and introducing quantum machine learning in text classification tasks.



## 327 Limitations

328 The limitation of our work comes from the hard-  
329 ware resources for running quantisation parts.  
330 Since we use the quantum simulation software in-  
331 stead of the quantum computer, the memory usage  
332 increases exponentially, i.e., the space complexity  
333 is  $O(2^n)$  where  $n$  is the number of qubits. There-  
334 fore, we can only use as few qubits as possible in  
335 DisCoWord training and QSA, which limits the  
336 performance and parameter searching. In addition,  
337 we do not analyse the impact of different syntactic  
338 parsing methods, which is left to future work.

## 339 References

340 Ville Bergholm, Josh Izaac, Maria Schuld, Christian  
341 Gogolin, Shahnawaz Ahmed, Vishnu Ajith, M So-  
342 haib Alam, Guillermo Alonso-Linaje, B Akash-  
343 Narayanan, Ali Asadi, et al. 2018. Pennylane: Au-  
344 tomatic differentiation of hybrid quantum-classical  
345 computations. *arXiv preprint arXiv:1811.04968*.

346 Shikhar Bharadwaj and Shirish Shevade. 2022. Effi-  
347 cient constituency tree based encoding for natural  
348 language to bash translation. In *Proceedings of the*  
349 *2022 Conference of the North American Chapter of*  
350 *the Association for Computational Linguistics: Hu-*  
351 *man Language Technologies*, pages 3159–3168.

352 Steven Bird, Ewan Klein, and Edward Loper. 2009. *Nat-*  
353 *ural language processing with Python: analyzing text*  
354 *with the natural language toolkit*. " O'Reilly Media,  
355 Inc."

356 Jose Camacho-Collados, Yerai Doval, Eugenio  
357 Martinez-Cámara, Luis Espinosa-Anke, Francesco  
358 Barbieri, and Steven Schockaert. Learning cross-  
359 lingual embeddings from twitter via distant supervi-  
360 sion.

361 El Amine Cherrat, Iordanis Kerenidis, Natansh Mathur,  
362 Jonas Landman, Martin Strahm, and Yun Yvonna Li.  
363 2022. Quantum vision transformers. *arXiv preprint*  
364 *arXiv:2209.08167*.

365 Kevin Clark, Urvashi Khandelwal, Omer Levy, and  
366 Christopher D Manning. 2019. What does bert look  
367 at? an analysis of bert's attention. In *Proceedings*  
368 *of the 2019 ACL Workshop BlackboxNLP: Analyzing*  
369 *and Interpreting Neural Networks for NLP*, pages  
370 276–286.

371 Stephen Clark. 2021. Something old, something new-  
372 grammar based ccg parsing with transformer models.

373 Stephen Clark, Bob Coecke, and Mehrnoosh Sadrzadeh.  
374 2008. A compositional distributional model of mean-  
375 ing. In *Proceedings of the Second Quantum Interac-*  
376 *tion Symposium (QI-2008)*, pages 133–140. Citeseer.

Andrea Delgado and Kathleen E Hamilton. 2022. Un-  
supervised quantum circuit learning in high energy  
physics. *Physical Review D*, 106(9):096006. 377  
378  
379

Riccardo Di Sipio, Jia-Hong Huang, Samuel Yen-Chi  
Chen, Stefano Mangini, and Marcel Worrying. 2022.  
The dawn of quantum natural language processing.  
In *ICASSP 2022-2022 IEEE International Confer-*  
*ence on Acoustics, Speech and Signal Processing*  
*(ICASSP)*, pages 8612–8616. IEEE. 380  
381  
382  
383  
384  
385

David P DiVincenzo. 1998. Quantum gates and circuits.  
*Proceedings of the Royal Society of London. Series A:*  
*Mathematical, Physical and Engineering Sciences*,  
454(1969):261–276. 386  
387  
388  
389

Mingzhe Du, Mouad Hakam, See-Kiong Ng, and  
Stéphane Bressan. 2023. Constituency-informed  
and constituency-constrained extractive question an-  
swering with heterogeneous graph transformer. In  
*Transactions on Large-Scale Data-and Knowledge-*  
*Centered Systems LIII*, pages 90–106. Springer. 390  
391  
392  
393  
394  
395

Stuart Hadfield, Zihui Wang, Bryan O'gorman,  
Eleanor G Rieffel, Davide Venturelli, and Rupak  
Biswas. 2019. From the quantum approximate opti-  
mization algorithm to a quantum alternating operator  
ansatz. *Algorithms*, 12(2):34. 396  
397  
398  
399  
400

Shengyuan Hou, Jushi Kai, Haotian Xue, Bingyu  
Zhu, Bo Yuan, Longtao Huang, Xinbing Wang, and  
Zhouhan Lin. 2022. Syntax-guided localized self-  
attention by constituency syntactic distance. *arXiv*  
*e-prints*, pages arXiv–2210. 401  
402  
403  
404  
405

Phu Mon Htut, Jason Phang, Shikha Bordia, and  
Samuel R Bowman. 2019. Do attention heads in  
bert track syntactic dependencies? *arXiv preprint*  
*arXiv:1911.12246*. 406  
407  
408  
409

Ziniu Hu, Yuxiao Dong, Kuansan Wang, and Yizhou  
Sun. 2020. Heterogeneous graph transformer. In  
*Proceedings of the web conference 2020*, pages 2704–  
2710. 410  
411  
412  
413

Gregg Jaeger. 2007. *Quantum information*. Springer. 414

Dimitri Kartsaklis, Ian Fan, Richie Yeung, Anna Pear-  
son, Robin Lorenz, Alexis Toumi, Giovanni de Fe-  
lice, Konstantinos Meichanetzidis, Stephen Clark,  
and Bob Coecke. lambeq-an efficient high-level  
python library for quantum nlp. 415  
416  
417  
418  
419

Dimitri Kartsaklis, Ian Fan, Richie Yeung, Anna Pear-  
son, Robin Lorenz, Alexis Toumi, Giovanni de Fe-  
lice, Konstantinos Meichanetzidis, Stephen Clark,  
and Bob Coecke. 2021. lambeq: An efficient high-  
level python library for quantum nlp. *arXiv e-prints*,  
pages arXiv–2110. 420  
421  
422  
423  
424  
425

Nikita Kitaev and Dan Klein. 2018. **Constituency pars-**  
**ing with a self-attentive encoder**. In *Proceedings*  
*of the 56th Annual Meeting of the Association for*  
*Computational Linguistics (Volume 1: Long Papers)*,  
pages 2676–2686, Melbourne, Australia. Association  
for Computational Linguistics. 426  
427  
428  
429  
430  
431

432	Joachim Lambek. 1999. Type grammar revisited.	Shangshang Shi, Zhimin Wang, Jiaxin Li, Yanan Li,	488
433	In <i>Logical Aspects of Computational Linguistics:</i>	Ruimin Shang, Haiyong Zheng, Guoqiang Zhong,	489
434	<i>Second International Conference, LACL'97 Nancy,</i>	and Yongjian Gu. 2023. A natural nisq model of	490
435	<i>France, September 22-24, 1997 Selected Papers 2,</i>	quantum self-attention mechanism. <i>arXiv e-prints,</i>	491
436	pages 1–27. Springer.	pages arXiv–2305.	492
437	Joachim Lambek. 2008. <i>From Word to Sentence: a</i>	Richard Socher, Alex Perelygin, Jean Wu, Jason	493
438	<i>computational algebraic approach to grammar.</i> Poli-	Chuang, Christopher D Manning, Andrew Y Ng, and	494
439	metrica sas.	Christopher Potts. 2013. Recursive deep models for	495
440	Zuchao Li, Kevin Parnow, and Hai Zhao. 2022. Incorporating rich syntax information in grammatical error	semantic compositionality over a sentiment treebank.	496
441	correction. <i>Information Processing &amp; Management,</i>	In <i>Proceedings of the 2013 conference on empirical</i>	497
442	59(3):102891.	<i>methods in natural language processing,</i> pages	498
443		1631–1642.	499
444	Robin Lorenz, Anna Pearson, Konstantinos Meichanetzidis, Dimitri Kartsaklis, and Bob Coecke. 2023. Qnlp in practice: Running compositional models of meaning on a quantum computer. <i>Journal of Artificial Intelligence Research,</i> 76:1305–1342.	James C Spall. 1998. Implementation of the simultaneous perturbation algorithm for stochastic optimization. <i>IEEE Transactions on aerospace and electronic systems,</i> 34(3):817–823.	500
445			501
446			502
447			503
448			
449	Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In <i>International Conference on Learning Representations.</i>	Marcos Zampieri, Shervin Malmasi, Preslav Nakov, Sara Rosenthal, Noura Farra, and Ritesh Kumar. 2019. Predicting the type and target of offensive posts in social media. In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers),</i> pages 1415–1420.	504
450			505
451			506
452	Weicheng Ma, Brian Wang, Hefan Zhang, Lili Wang, Rolando Coto-Solano, Saeed Hassanpour, and Soroush Vosoughi. 2023. Improving syntactic probing correctness and robustness with control tasks. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers),</i> pages 402–415.		507
453			508
454			509
455			510
456			511
457			
458			
459	David Mareček and Rudolf Rosa. 2019. From balustrades to pierre vinken: Looking for syntax in transformer self-attentions. In <i>Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP,</i> pages 263–275.	Yue Zhang and Zhenghua Li. 2022. Csyngec: Incorporating constituent-based syntax for grammatical error correction with a tailored gec-oriented parser. <i>arXiv e-prints,</i> pages arXiv–2211.	512
460			513
461			514
462			515
463			
464			
465	Mikko Möttönen, Juha J Vartiainen, Ville Bergholm, and Martti M Salomaa. 2005. Transformation of quantum states using uniformly controlled rotations. <i>Quantum Information &amp; Computation,</i> 5(6):467–473.	Jianyu Zheng and Ying Liu. 2023. What does chinese bert learn about syntactic knowledge? <i>PeerJ Computer Science,</i> 9.	516
466			517
467			518
468			
469	Lee J O’Riordan, Myles Doyle, Fabio Baruffa, and Venkatesh Kannan. 2020. A hybrid classical-quantum workflow for natural language processing. <i>Machine Learning: Science and Technology,</i> 2(1):015011.	Elton Yechao Zhu, Sonika Johri, Dave Bacon, Mert Esencan, Jungsang Kim, Mark Muir, Nikhil Murgai, Jason Nguyen, Neal Pisenti, Adam Schouela, et al. 2022a. Generative quantum learning of joint probability distribution functions. <i>Physical Review Research,</i> 4(4):043092.	519
470			520
471			521
472			522
473			523
474	Bo Pang and Lillian Lee. 2005. Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In <i>Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics (ACL’05),</i> pages 115–124.	Fangyi Zhu, Lok You Tan, See-Kiong Ng, and Stéphane Bressan. 2022b. Syntax-informed question answering with heterogeneous graph transformer. In <i>Database and Expert Systems Applications: 33rd International Conference, DEXA 2022, Vienna, Austria, August 22–24, 2022, Proceedings, Part I,</i> pages 17–31. Springer.	524
475			525
476			526
477			527
478			528
479	Anna Rogers, Olga Kovaleva, and Anna Rumshisky. 2020. A primer in BERTology: What we know about how BERT works. <i>Transactions of the Association for Computational Linguistics,</i> 8:842–866.		529
480			530
481			531
482			
483	Rishab Sharma, Fuxiang Chen, Fatemeh Fard, and David Lo. 2022. An exploratory study on code attention in bert. In <i>Proceedings of the 30th IEEE/ACM International Conference on Program Comprehension,</i> pages 437–448.	<b>A Methodology Details</b>	532
484			
485			
486			
487			
		<b>A.1 Architecture Overview</b>	533
		Our architecture overview is shown in Figure 1.	534
		<b>A.2 DisCoWord Supplementary</b>	535
		Specifically, we follow DisCoCat to transform sentences into parameterised quantum circuits according to their string diagrams produced by the combinatory categorial grammar (CCG) parser, as shown in Figure 2 and Figure 3. In this case, the word is	536
			537
			538
			539
			540

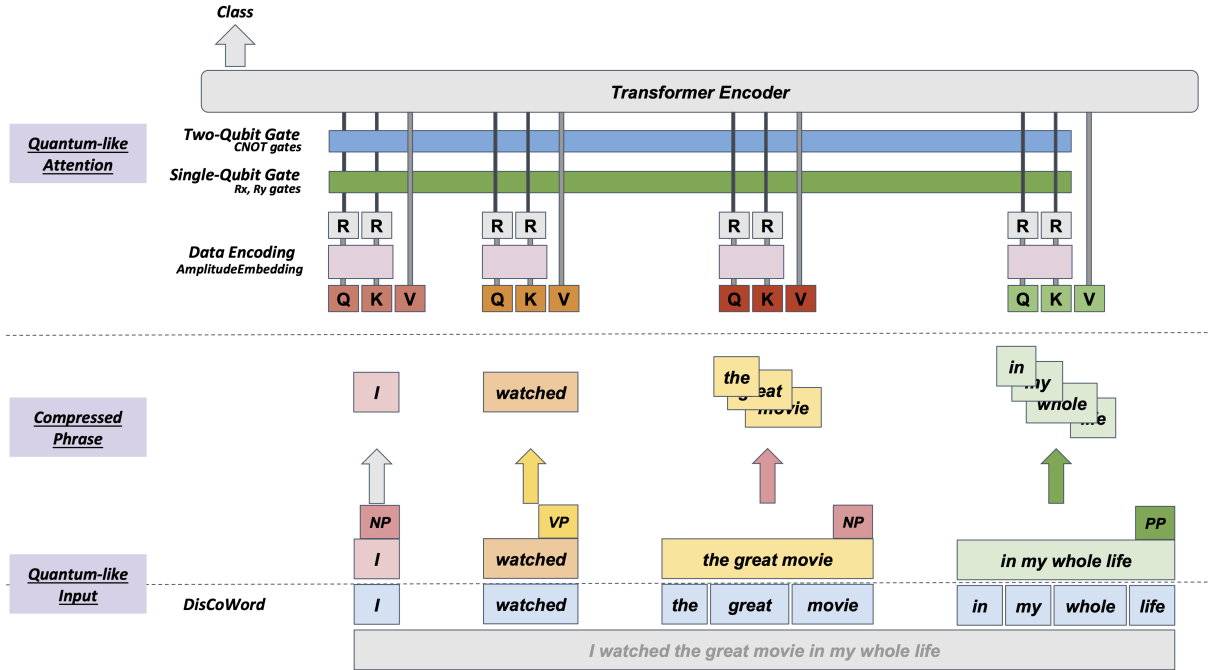


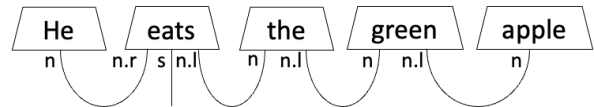
Figure 1: The overview of model architecture

541 represented by the quantum state of the qubits of  
 542 its pregroup grammar, named word state. However,  
 543 DisCoCat randomly initialises the quantum circuit  
 544 parameters, which means the word states are not the  
 545 semantics of these words. They can only represent  
 546 the ‘meaning’ in the specific sentence classification  
 547 task and can not be used in other NLP tasks. We  
 548 use pretrained word embeddings to initialise these  
 549 parameters so that the word states can hold both the  
 550 syntactic and semantic information based on their  
 551 contextual grammar structure and pretrained word  
 552 embeddings. After training, we evaluate the word  
 553 states. Note that the evaluated word states are complex  
 554 vectors. To use them in later task classification  
 555 experiments, we concatenate the real parts and image  
 556 parts as word embeddings, named DisCoWord  
 557 representation.

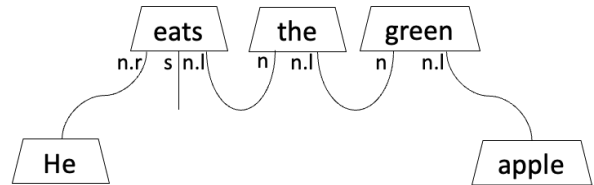
## 558 B Experiments Details

### 559 B.1 Computational Resource Utilization

560 We use four A100 GPUs in our work. It takes 20-  
 561 50 GPU hours to train DisCoWord representation.  
 562 For training our Transformer variants, it takes 5-90  
 563 GPU hours to train the Transformer (QSA) and 1-  
 564 10 GPU hours for other variants. The reason that  
 565 training Transformer (QSA) takes more time is due  
 566 to the use of PennyLane (Bergholm et al., 2018)  
 567 for quantum simulation.



(a) String diagram before bending noun.



(b) String diagram after bending noun.

Figure 2: String diagrams of the sentence, where  $n$ ,  $n.r$ ,  $n.l$ ,  $s$  are the grammar types (Lambek, 2008) of words, the types under a word form its pregroup grammar (Lambek, 1999).

### B.2 DisCoWord Training Details

568 Due to the hardware resource limitation, we split  
 569 each dataset into subsets and train the DisCoCat  
 570 with them separately. In addition, there are three  
 571 processings on the word state: **a)** If it is longer than  
 572 256 dimensions, we drop the word state evaluations  
 573 to save memory. If it is less than 256 dimensions,  
 574 we convert the states into 256 dimensions and then  
 575 conduct a zero-padding. **b)** If the specific pregroup  
 576 grammar appears in the testing set but not in the  
 577 training set, we take the mean of a word’s states to  
 578 deal with the case. For example, the word ‘like’ has  
 579 three pregroup grammars representing its adjective,  
 580

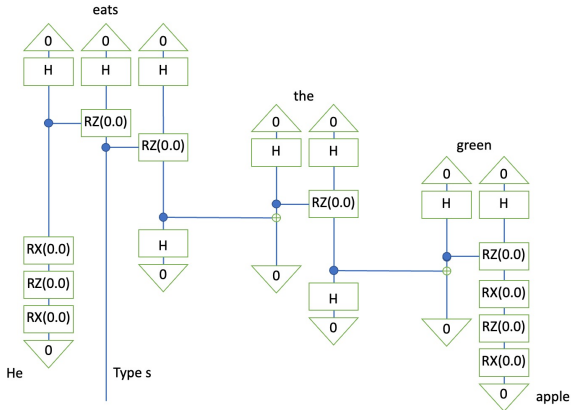


Figure 3: Quantum circuit of the sentence. The quantum state of the qubits belonging to a word represents its word state.

noun, and conjunction meaning respectively, but the conjunctive ‘like’ only appears in the testing set. We take the mean of the word states of the adjective ‘like’ and the nominal ‘like’ as the state of the conjunctive ‘like’. c) Since the evaluated states are complex vectors, we concatenate the real-valued and image components for the Transformer.

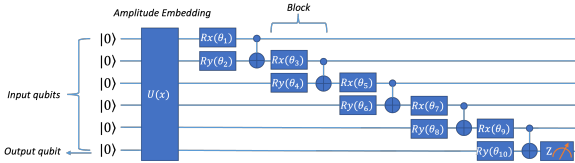


Figure 4: The quantum circuit example.

### B.3 Quantum Self-Attention Design

An example of the quantum circuit for QSA is shown in Figure 4. In the Figure,  $x$  is the  $Q/K$  vector of a head, and  $U(x)$  represents the AmplitudeEmbedding feature map.  $\theta_1, \theta_2 \dots \theta_{10}$  are the trainable parameters. In this circuit, the number of input qubits  $n$  is 6, which implies that the dimensionality  $2^n$  of  $Q/K$  vector of a head is 64.

Method	MR	Tweet	SST-2	Offenseval
w/o-phrasing	21.35	15.27	19.6	23.57
ComPhE (B2T)	9.1	6.72	8.43	10.26
ComPhE (T2B)	6.43	5.54	6.06	7.04

Table 5: The statistics of input token length using different phrasing methods.

### B.4 Statistics of Phrasing Methods

As aforementioned, phrasing methods will reduce the input token length. Here we list the statistics of the input token length using different phrasing methods in Table 5.