

Semantic Hilbert Space for Text Representation Learning

Benyou Wang, Qiuchi Li, Massimo Melucci
University of Padua
Padua, Italy
wang.qiuchili,melo@dei.unipd.it

Dawei Song
Beijing Institute of Technology
Beijing, China
dawei.song@open.ac.uk

ABSTRACT

Capturing the meaning of sentences has long been a challenging task. Current models tend to apply linear combinations of word features to conduct semantic composition for a bigger-granularity units e.g. phrase, sentence and documents. However, the semantic linearity does not always hold in human language. For instance, the meaning of the phrase "ivory tower" can not be deduced by linearly combining the meanings of "ivory" and "tower". To address this issue, we propose a new framework that models different levels of semantic units (e.g. sememe, word, sentence and semantic abstraction) on a single *Semantic Hilbert Space*, which naturally admits a non-linear semantic composition by means of a complex-valued vector word representation. An end-to-end neural network¹ is proposed to implement the framework in the text classification task, and evaluation results on six benchmarking text classification datasets demonstrate the effectiveness, robustness and self-explanation power of the proposed model. Furthermore, intuitive case studies are conducted to help end users to understand how the framework works.

CCS CONCEPTS

• **Information systems** → *Document structure; Content analysis and feature selection;*

KEYWORDS

text understanding, neural network, quantum theory

ACM Reference Format:

Benyou Wang, Qiuchi Li, Massimo Melucci and Dawei Song. 1997. Semantic Hilbert Space for Text Representation Learning. In *Proceedings of ACM Woodstock conference (WWW '2019)*. ACM, New York, NY, USA, 7 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

¹<https://github.com/wabyking/qnn>

[†] Benyou Wang and Qiuchi Li contribute equally and share the co-first authorship. Qiuchi Li (qiuchili@dei.unipd.it) is the corresponding author.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://www.acm.org).

WWW '2019, May 2019, San Francisco, California USA

© 2016 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

In natural language understanding, it is crucial, yet challenging, to model sentences and capture their meanings. Essentially, most statistical machine learning models [3, 9, 15, 20, 23] are built within a linear bottom-up framework, where words are the basic features adopting a low-dimensional vector representation, and a sentence is modeled as a linear combination of individual word vectors. Such linear semantic composition is efficient, but does not always hold in human language. For example, the phrase "ivory tower", which means "a state of privileged seclusion or separation from the facts and practicalities of the real world", is not a linear combination of the individual meanings of "ivory" and "tower". Instead, it carries a new meaning. We are therefore motivated to investigate a new language modeling paradigm to account for such intricate non-linear combination of word meanings.

Drawing inspiration from the recent findings in the emerging research area of quantum cognition, which suggest that human cognition [1, 2, 8] especially language understanding [6, 7, 29] exhibit certain non-classical phenomena (i.e. quantum-like phenomena), we propose a theoretical framework, named *Semantic Hilbert Space*, to formulate quantum-like phenomena in language understanding and to model different levels of semantic units in a unified space.

In Semantic Hilbert Space, we assume that words can be modeled as microscopic particles in superposition states, over the basic sememes (i.e. minimum semantic units in linguistics), while a combination of word meanings can be viewed as a mixed system of particles. The Semantic Hilbert Space represents different levels of semantic units, ranging from basic sememes, words and sentences, on a unified complex-valued vector space. This is fundamentally different from existing quantum-inspired neural networks for question answering [31, 32] which are based on a real vector space. In addition, we introduce a new semantic abstraction, named as *Semantic Measurements*, which are also embedded in the same vector space and trainable to extract high-level features from the mixed system.

As shown in Fig. 1, the Semantic Hilbert Space is built on the basis of quantum probability (QP), which is the probability theory for explaining the uncertainty of quantum superposition. As quantum superposition requires the use of the complex field, Semantic Hilbert Space has complex

values and operators. In particular, the probability function is implemented by a unique (complex) density operator.

Semantic Hilbert Space adopts a complex-valued vector representation of unit length, where each component adopts an amplitude-phase form $z = re^{i\phi}$. We hereby hypothesize that the amplitude r and complex phase ϕ can be used to encode different levels of semantics such as lexical-level co-occurrence, hidden sentiment polarity or topic-level semantics. When word vectors are combined, even in a simple complex-valued addition form, the resulting expression will entail a non-linear composition of amplitudes and phases, thus indicating a complicated fusion of different levels of semantics. A more detailed explanation is given in Sec. 3. In this way, the complex-valued word embedding is fundamentally different from existing real-valued word embedding. A series of ablation tests indicate that the complex-valued word embedding can increase performance.

The Semantic Hilbert Space is an abstract representation of our approach to modeling language through QP. At the level of implementation, an efficient and effective computational framework is needed to cope with large text collections. To do so, we propose an end-to-end neural network architecture, which provides means for training of the network components. Each component corresponds to a physical meaning of quantum probability with well-defined mathematical constraints. Moreover, each component is easier to understand than the kernels in convolutional neural network and cells in recurrent neural networks.

The network proposed in this paper is evaluated on six benchmarking datasets for text classification and achieves a steady increase over existing models. Moreover, it is shown that the proposed network is advantageous due to its high robustness and self-explanation capability.

2 SEMANTIC HILBERT SPACE

The mathematical foundation of Quantum Theory is established on a Hilbert Space over the complex field. In order to borrow the underlying mathematical formalism of quantum theory for language understanding, it is necessary to build such a Hilbert Space for language representation. In this study, we build a *Semantic Hilbert Space* \mathcal{H} over the complex field. As is illustrated in Fig. 1, multiple levels of semantic units are modeled on this common Semantic Hilbert Space. In the rest of this section, the semantic units under modeling are introduced separately.

We follow the standard *Dirac Notation* for Quantum Theory. A unit vector and its transpose are denoted as a ket $|\mu\rangle$ and a bra $\langle\mu|$, respectively. The inner product and outer product of two unit vectors \vec{u} and \vec{v} are denoted as $\langle u|v\rangle$ and $|u\rangle\langle v|$ respectively.

2.1 Sememes

Sememes are the minimal non-separable semantic units of word meanings in language universals [13]. For example, the word “ironsmith” is composed of sememes “human”, “occupation”, “metal” and “industrial”. We assume that the Semantic Hilbert Space \mathcal{H} is spanned by a set of orthogonal basis $\{|e_j\rangle\}_{j=1}^n$ corresponding to a finite closed set of sememes $\{e_j\}_{j=1}^n$. In the quantum language, the set of sememes are modeled as *basis states*, which is the basis for representing any quantum state. In Fig. 1, the axes of the Semantic Hilbert Space correspond to the set of sememe states, and semantic units with larger granularity are represented on its basis.

2.2 Words

The meaning of a word is a combination of sememes. We adopt the concept of *superposition* to formulate this combination. Essentially, a word w is modeled as a quantum particle in *superposition state*, represented by a unit-length vector in the Semantic Hilbert Space \mathcal{H} , as can be seen in Fig. 1. It can be written as a linear combination of the basis states for sememes:

$$|w\rangle = \sum_{j=1}^n r_j e^{i\phi_j} |e_j\rangle \quad (1)$$

where the complex-valued weight $r_j e^{i\phi_j}$ denotes how much the meaning of word w is associated with the sememe e_j . Here $\{r_j\}_{j=1}^n$ are non-negative real-valued amplitudes satisfying $\sum_{j=1}^n r_j^2 = 1$ and $\phi_j \in [-\pi, \pi]$ are the corresponding complex phases. We could also transfer the complex number in a complex plane as $re^{i\phi} = r \cos \phi + ir \sin \phi$.

It is worth noting that the complex phases $\{\phi_j\}$ are crucial as they implicitly entail the *quantum interference* between words. Suppose two words w_1 and w_2 are associated to weights $r_j^{(1)} e^{i\phi_j^{(1)}}$ and $r_j^{(2)} e^{i\phi_j^{(2)}}$ for the sememe e_j . The two words in combination are therefore at the state e_j with a probability of

$$\left| r_j^{(1)} e^{i\phi_j^{(1)}} + r_j^{(2)} e^{i\phi_j^{(2)}} \right|^2 = \left| r_j^{(1)} \right|^2 + \left| r_j^{(2)} \right|^2 + 2r_j^{(1)} r_j^{(2)} \cos \left(\phi_j^{(1)} - \phi_j^{(2)} \right) \quad (2)$$

where the term $2r_j^{(1)} r_j^{(2)} \cos(\phi_j^{(1)} - \phi_j^{(2)})$ reflects the interference between the two words, where as the classical case corresponds to a particular case $\phi_j^{(1)} = \phi_j^{(2)} = 0$.

2.3 Semantic Compositions

As is illustrated in Fig. 1, we view a word composition (e.g. a sentence) as a bag of words [14], each of which is modeled as a particle in superposition state on the Semantic Hilbert Space \mathcal{H} . To obtain the semantic composition of words, we leverage the concept of *quantum mixture* and formulate the word composition as a mixed system composed of the word

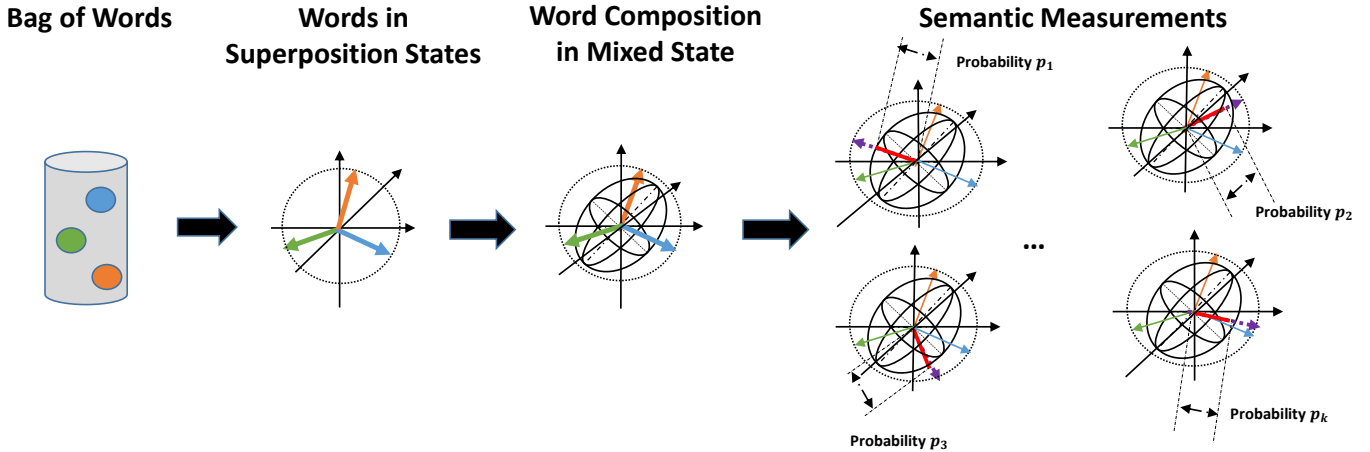


Figure 1: Illustration of Semantic Hilbert Space. The green, blue and orange colors correspond to three different words modeled as quantum particles. The black dotted circle represents the unit ball in the Semantic Hilbert Space. The ellipsoid in solid line refers to the quantum probability distribution defined by the density matrix of the word composition. The purple lines are semantic measurements. The intersections of the ellipsoids and semantic measurements are in thick red lines, the lengths of which correspond to measurement probabilities.

superposition states. The system is in a *mixed state* represented by a n -by- n density matrix ρ on \mathcal{H} , which is positive semi-definite with trace 1. It is computed as follows:

$$\rho = \sum_i p(i) |w_i\rangle \langle w_i|, \tag{3}$$

where $|w_i\rangle$ denotes the superposition state of the i -th word and $p(i)$ is the classical probability of the state $|w_i\rangle$ with $\sum_i p(i) = 1$. It determines the contribution of the word w_i to the overall semantics.

The complex-valued density matrix ρ can be seen non-classical distribution of sememes in \mathcal{H} . Its diagonal elements are real and form a classical distribution of sememes, while its complex-valued off-diagonal entries encode the interplay between sememes, which in turn gives rise to the interference between words. A density matrix assigns a probability value for any state on \mathcal{H} such that the values for any set of orthogonal states sum up to 1 [12]. Hence it is visualized as an ellipsoid in Fig. 1, assigning a quantum probability to a unit vector with the intersection length.

2.4 Semantic Measurements

As a non-classical probability distribution, a sentence density matrix carries rich information and in particular it contains all the information about a quantum system. In order to extract the relevant information to a concrete task from the semantic composition, we build a set of measurements and compute the probability that the mixed system falls onto each of the measurements as a high-level abstraction of the semantic composition.

Suppose our proposed *semantic measurements* are associated with a set of measurement projectors $\{P_i\}_{i=1}^k$. According to the Born’s rule [5], applying the measurement projector P_i onto the sentence density matrix ρ yields the following result:

$$p_i = tr(P_i\rho) \tag{4}$$

Here, we only consider pure states as measurement states, i.e. $P_i = |v_i\rangle \langle v_i|$. Moreover, we ignore the constraints of the measurements states $\{|v_i\rangle\}_{i=1}^k$ (i.e. orthogonality or completeness), but keep them trainable, so that the most suitable measurements can be determined automatically by the data in a concrete task, such as classification or regression. In this way, the trainable semantic measurements can be understood as a similar approach to supervised dimensionality reduction [11], but in a quantum probability framework with complex values.

3 QUANTUM PROBABILITY DRIVEN NETWORK

In order to implement the proposed framework, we further propose an end-to-end neural network on its basis. Fig. 2 shows the architecture of the proposed Quantum Probability Driven Network (QPDN). The embedding layer, composed of a unit complex-valued embedding and a term-weight lookup table, captures the basic lexical features. While the mixture layer is designed to combine the low-level bag-of-word features with an additive complex-valued outer product operation. The measurement layer adopts a set of trainable semantic measurements to extract the higher-level features for the final linear classifier. In the following we will introduce the architecture layer by layer.

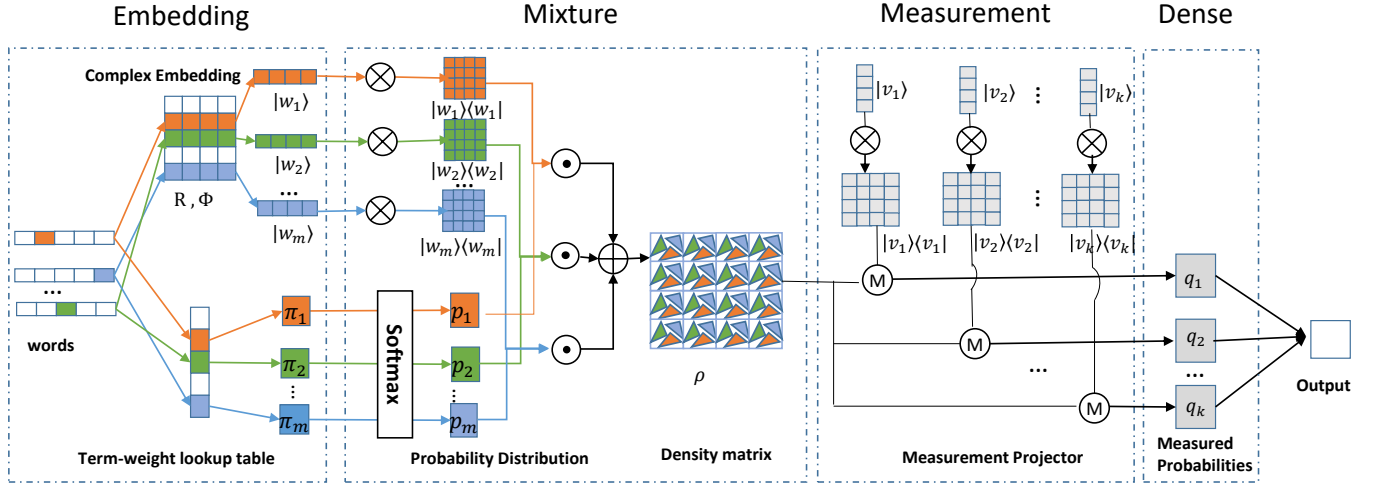


Figure 2: Architecture of Quantum probability-driven Neural Network. \odot means that a matrix multiplies a number with each elements. \oplus refers to a element-wise addition. \otimes denotes a outer production to a vector, \textcircled{M} means a measurement operation according to Eq. 4.

3.1 Embedding Layer

The parameters of the embedding layer are $\{R, \Phi, \Pi\}$, respectively, denoting the amplitude embedding, the phase embedding, and the term-weight lookup table. Eq. 1 expresses a quantum representation as a unit-length, complex-valued vector representation for a word w , i.e. $|w\rangle = [r_1 e^{i\phi_1}, r_2 e^{i\phi_2}, \dots, r_n e^{i\phi_n}]^T$. The term-weight lookup table is used to weight words for semantic combinations, which will be described in the next subsection. During training, word embeddings need to be normalized to unit length after each batch.

This representation allows for a non-linear composition of amplitudes and phases in its mathematical form. Suppose two words w_1 and w_2 are of weights $r_j^{(1)} e^{i\phi_j^{(1)}}$ and $r_j^{(2)} e^{i\phi_j^{(2)}}$ for the j^{th} dimension (corresponding to the j^{th} sememe). The combination of w_1 and w_2 gives rise to a weight $r_j e^{i\phi_j}$ for the j^{th} dimension computed as

$$\begin{aligned} r_j e^{i\phi_j} &= r_j^{(1)} e^{i\phi_j^{(1)}} + r_j^{(2)} e^{i\phi_j^{(2)}} \\ &= \sqrt{|r_j^{(1)}|^2 + |r_j^{(2)}|^2 + 2r_j^{(1)} r_j^{(2)} \cos(\phi_j^{(1)} - \phi_j^{(2)})} \\ &\quad \times e^{i \arctan\left(\frac{r_j^{(1)} \sin(\phi_j^{(1)}) + r_j^{(2)} \sin(\phi_j^{(2)})}{r_j^{(1)} \cos(\phi_j^{(1)}) + r_j^{(2)} \cos(\phi_j^{(2)})}\right)} \end{aligned} \quad (5)$$

Where both r_j and ϕ_j is a non-linear combination of $r_j^{(1)}, r_j^{(2)}, \phi_j^{(1)}$ and $\phi_j^{(2)}$. If the amplitudes and phases are associated to different levels of information, the amplitude-phase representation then naturally gives rise to a non-linear fusion of information.

3.2 Mixture Layer

A sentence is modeled as a density matrix, which is constructed in Sec. 2.3. Instead of using uniform weights in Eq. 3, word-sensitive weights are used for each word, which is commonly used in IR, e.g. inverse document frequency (IDF) as a word-dependent weight in TF-IDF scheme [28].

In order to guarantee the unit trace length for density matrix, the word weights which are from the lookup table in a sentence are normalized to a probability value through a softmax operation: $p(i) = e^{\pi(w_i)} / \sum_j e^{\pi(w_j)}$. Compared to IDF weight, the normalized weight for a specific word in our approach is not static but updated adaptively in the training phase. Even in the inference/test phase, the real term weight i.e. $p(w_i)$ is also not static, but highly depends on the neighbor context words through nonlinear softmax function.

3.3 Measurement Layer

The measurement layer adopts a set of 1-order measurement projectors $\{|v_i\rangle \langle v_i|\}_{i=1}^k$ where $|v_i\rangle \langle v_i|$ is the outer product of its corresponding state in Semantic Hilbert Space $|v_i\rangle$. After each measurement, we can obtain one probability for each measurement state like $q_j = \text{tr}(\rho |v_j\rangle \langle v_j|)$. Finally, we can obtain a vector $\vec{q} = [q_1, q_2, \dots, q_k]$. Similarly to the word vectors which are also represented as unit states, the states $|v_i\rangle$ are also normalized after several batches.

3.4 Dense Layer

The vector \vec{q} in the measurement layer consists of k positive scalar numbers and it is used to infer the label for a given sentence. A dense layer with softmax activation is adopted

Dataset	train	test	vocab.	task	Classes
CR	4K	CV	6K	product reviews	2
MPQA	11k	CV	6K	opinion polarity	2
SUBJ	10k	CV	21k	subjectivity	2
MR	11.9k	CV	20k	movie reviews	2
SST	67k	2.2k	18k	movie reviews	2
TREC	5.4k	0.5k	10k	Question	6

Table 1: Dataset Statistics. (CV means 10-fold cross validation for testing performance.)

after the measurement layer to get a classification probability distribution, i.e. $\hat{y} = \text{softmax}(\hat{q} \cdot W)$. The loss is designed as a cross-entropy loss between \hat{y} and the one-hot label \vec{y} .

4 EXPERIMENTS

Our model is evaluated on 6 datasets for text classification: CR customer review [17], MPQA opinion polarity [30], SUBJ sentence subjectivity [25], MR movie review [25], SST binary sentiment classification [27], and TREC question classification [21]. The statistics of them are shown in Tab. 1.

We compared the proposed QPDN with various models, including Uni-TFIDF, Word2vec, FastText [18] and Sent2Vec [24] as unsupervised representation learning baselines, CaptionRep [15] and DictRep [16] as supervised representation learning baselines, as well as CNN [19] and BiLSTM [10] for advanced deep neural networks. We report the classification accuracy values of these models from the original papers.

We used Glove word vectors [26] with 50,100,200 and 300 dimensions respectively. The amplitude embedding values are initialized by L2-norm, while the phases in complex-valued embedding are randomly initialized in $-\pi$ to π . We searched for the best performance in a parameter pool, which contains a learning rate in $\{1E-3, 1E-4, 1E-5, 1E-6\}$, an L2-regularization ratio in $\{1E-5, 1E-6, 1E-7, 1E-8\}$, a batch size in $\{8, 16, 32, 64, 128\}$, and the number of measurements in $\{5, 10, 20, 50, 100, 200\}$.

The main parameters in our model are R and Φ . Since both of them are $n \times |V|$ in shape, the number of parameters is roughly two times that of fastText [22]. For the other parameters, Π is $|V| \times 1$, $\{|v_i\rangle\}_{i=1}^k$ is $k \times 2n$, while W is $k \times |L|$ with L being the label set. Apart from word embeddings, the model is robust with limited scale at $k \times 2n + n \times |V| + k \times |L|$ for the number of parameters.

The results in Tab. 2 demonstrate the effectiveness of our model, with improved classification accuracies over some strong baseline supervised and unsupervised representation models on most of the datasets except MPQA. In comparison with more advanced models including BiLSTM and CNN, our model generally performs better than BiLSTM with increased accuracy values on the multi-class classification

Table 2: Experimental Results in percentage (%). The best performed value (except for CNN/LSTM) for each dataset is in bold. where \dagger means a significant improvement over FastText.

Model	CR	MPQA	MR	SST	SUBJ	TREC
Uni-TFIDF	79.2	82.4	73.7	-	90.3	85.0
Word2vec	79.8	88.3	77.7	79.7	90.9	83.6
FastText [18]	78.9	87.4	76.5	78.8	91.6	81.8
Sent2Vec [24]	79.1	87.2	76.3	80.2	91.2	85.8
CaptionRep [15]	69.3	70.8	61.9	-	77.4	72.2
DictRep [16]	78.7	87.2	76.7	-	90.7	81.0
Ours: QPDN	81.0\dagger	87.0	80.1\dagger	83.9\dagger	92.7\dagger	88.2\dagger
CNN [19]	81.5	89.4	81.1	88.1	93.6	92.4
BiLSTM [10]	81.3	88.7	77.5	80.7	89.6	85.2

dataset (TREC) and three binary text classification datasets (MR, SST & SUBJ). However, it under-performs CNN on all 6 datasets with a difference of over 2% on 3 of them (MPQA, SST & TREC), probably because that it uses fewer parameters and simpler structures. We argue that QPDN achieves a good balance between effectiveness and efficiency, due to the fact that it outperforms BiLSTM.

5 DISCUSSIONS

This section discusses the power of self-explanation and conducts an ablation test to examine the usefulness of important components of the network, especially the complex-valued word embedding.

Self-explanation Components. As is shown in Tab. 3, all components in our model have a clear physical meaning corresponding to quantum probability, where classical Deep Neural Network (DNN) can not well explain the role each component plays in the network. Essentially, we construct a bottom-up framework to represent each level of semantic units on a uniform Semantic Hilbert Space, from the minimum semantic unit, i.e. sememe, to the sentence representation. The framework is operationalized through superposition, mixture and semantic measurements. On the one hand, the explanation is reflected by well-designed constraints for all the components. On the other hand, some intuitive explanation can be performed on the crucial components of the network i.e. measurements, as shown in Sec. 5.

Ablation Test. An ablation test is conducted to examine how each component influences the final performance of QPDN. In particular, a double-length real word embedding network is implemented to examine the use of complex-valued word embedding, while mean weights and IDF weights are used as alternative word weighting strategies to check the necessity of introducing trainable weights. A set of non-trainable orthogonal projectors and a dense

Table 3: Physical meanings and constraints

Components	DNN	QPDN
Sememe	-	basis vector / basis state $\{w w \in C^n, w _2 = 1, \}$ complete & orthogonal
Word	real vector $(-\infty, \infty)$	unit complex vector / superposition state $\{w w \in C^n, w _2 = 1\}$
Low-level representation	real vector $(-\infty, \infty)$	density matrix / mixed system $\{\rho \rho = \rho^*, tr(\rho) = 1\}$
Abstraction	CNN/RNN $(-\infty, \infty)$	unit complex vector / measurement $\{w w \in C^n, w _2 = 1\}$
High-level representation	real vector $(-\infty, \infty)$	probabilities/ measured probability (0, 1)

Table 4: Ablation Test

Setting	SST	Δ
FastText [18]	0.7880	-0.0511
FastText [18] with double-dimension real word vectors fixed amplitude part but trainable phase part	0.7883	-0.0508
replace trainable weights with fixed mean weights	0.8199	-0.0192
replace trainable weights with fixed IDF weights	0.8303	-0.0088
non-trainable projectors with fixed orthogonal ones	0.8259	-0.0132
replace projectors with dense layer	0.8171	-0.0220
QPDN	0.8221	-0.0170
	0.8391	-

layer on top of the sentence density matrix are implemented to analyze the effect of trainable semantic measurements.

Due to limited space, we only report the ablation test result for SST, which is the largest and hence the most representative dataset. We use 100-dimensional real-valued word vectors and 50-dimensional complex-valued vectors for the models in the ablation test. All models under ablation are comparable in terms of time cost. Tab. 4 shows that each component plays an important role in the QPDN model. In particular, replacing complex embedding with double-dimension real word embedding leads to a 5% drop in performance, which indicates that the complex-valued word embedding is not merely doubling the number of parameters.

The comparison with IDF and mean weights shows that the data-driven scheme gives rise to high-quality word weights. The comparison with non-trainable projectors and directly applying a dense layer on the density matrix shows that trainable measurements bring benefits to the network.

Discriminative Semantic Directions. In order to better understand the well-trained measurement projectors, we obtained the top 10 nearest words in the complex-valued vector space for each trained measurement state (like $|v_i\rangle$), using KD tree [4]. Due to limited space, we take 5 measurements from the trained model for the MR dataset, and select words from the top 10 nearest words to each measurement. As can be seen in Tab. 5, the first measurement is roughly about

Table 5: The learned measurement for dataset MR. They are selected according to nearest words for a measurement vector in Semantic Hilbert Space

Measurement	Selected neighborhood words
1	change, months, upscale, recently, aftermath
2	compelled, promised, conspire, convince, trusting
3	goo, vez, errol, esperanza, ana
4	ice, heal, blessedly, sustains, make
5	continue, warned, preposterousness, adding, falseness

changes over time, the second concerning being motivated or forced to do something. While the third measurement groups uncommon non-English words together. The last two measurements also group words sharing similar meanings. It is therefore interesting to see that relevant words can somehow be grouped together into certain topics during the training process, which may be discriminative for the given task.

6 CONCLUSIONS

In order to better model the non-linearity of word semantic composition, we have developed a quantum-inspired framework that models different granularities of semantic units on the same Semantic Hilbert Space, and implement this framework into an end-to-end text classification network. The network shows a promising performance on 6 benchmarking text datasets, in terms of effectiveness, robustness and self-explanation ability. Moreover, the complex-valued word embedding approach, which inherently achieves the non-linear combination of word meanings, does bring benefits to the classification accuracy in a comprehensive ablation study.

This work is among the first steps to apply the quantum probabilistic framework to text modeling. We believe it is a promising direction. In the future, we would like to further extend this work by considering deeper and more complicated structures such as attention or memory mechanism in language, in order to investigate related quantum-like phenomena on textual data to provide more intuitive insights. Additionally, Semantic Hilbert Space in a tensor space is also worthy to be explored like [32], which may provide more interesting insights for current communities.

ACKNOWLEDGEMENT

This work is supported by the Quantum Access and Retrieval Theory (QUARTZ) project, which has received funding from the European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No. 721321.

REFERENCES

- [1] Diederik Aerts, Liane Gabora, and Sandro Sozzo. 2013. Concepts and Their Dynamics: A Quantum-Theoretic Modeling of Human Thought. *Topics in Cognitive*

- Science* (Sept. 2013). <https://doi.org/10.1111/tops.12042> arXiv: 1206.1069.
- [2] Diederik Aerts and Sandro Sozzo. 2014. Quantum Entanglement in Concept Combinations. *International Journal of Theoretical Physics* 53, 10 (Oct. 2014), 3587–3603. <https://doi.org/10.1007/s10773-013-1946-z>
- [3] Sanjeev Arora, Yingyu Liang, and Tengyu Ma. 2016. A Simple but Tough-to-Beat Baseline for Sentence Embeddings. (Nov. 2016). <https://openreview.net/forum?id=SyK00v5xx>
- [4] Jon Louis Bentley. 1975. Multidimensional binary search trees used for associative searching. *Commun. ACM* 18, 9 (1975), 509–517.
- [5] Max Born. 1926. Zur Quantenmechanik der Stoßvorgänge. *Zeitschrift für Physik* 37, 12 (Dec. 1926), 863–867. <https://doi.org/10.1007/BF01397477>
- [6] Peter Bruza, Kirsty Kitto, Douglas Nelson, and Cathy McEvoy. 2009. Is there something quantum-like about the human mental lexicon? *Journal of Mathematical Psychology* 53, 5 (2009), 362–377.
- [7] Peter D Bruza, Kirsty Kitto, Douglas McEvoy, and Cathy McEvoy. 2008. Entangling words and meaning. (2008).
- [8] Jerome R Busemeyer and Peter D Bruza. 2012. *Quantum models of cognition and decision*. Cambridge University Press.
- [9] Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc Barrault, and Antoine Bordes. 2017. Supervised Learning of Universal Sentence Representations from Natural Language Inference Data. In *EMNLP: Association for Computational Linguistics*, 670–680. <http://aclweb.org/anthology/D17-1070>
- [10] Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc Barrault, and Antoine Bordes. 2017. Supervised Learning of Universal Sentence Representations from Natural Language Inference Data. *arXiv:1705.02364 [cs]* (May 2017). arXiv: 1705.02364.
- [11] Ronald A Fisher. 1936. The use of multiple measurements in taxonomic problems. *Annals of eugenics* 7, 2 (1936), 179–188.
- [12] Andrew M Gleason. 1957. Measures on the closed subspaces of a Hilbert space. *Journal of mathematics and mechanics* (1957), 885–893.
- [13] Cliff Goddard and Anna Wierzbicka. 1994. *Semantic and lexical universals: Theory and empirical findings*. Vol. 25. John Benjamins Publishing.
- [14] Zellig S Harris. 1954. Distributional structure. *Word* 10, 2-3 (1954), 146–162.
- [15] Felix Hill, Kyunghyun Cho, and Anna Korhonen. 2016. Learning Distributed Representations of Sentences from Unlabelled Data. In *NAACL: Association for Computational Linguistics*, San Diego, California, 1367–1377. <http://www.aclweb.org/anthology/N16-1162>
- [16] Felix Hill, Kyunghyun Cho, Anna Korhonen, and Yoshua Bengio. 2016. Learning to Understand Phrases by Embedding the Dictionary. *TACL* 4 (2016), 17–30. <https://tacl2013.cs.columbia.edu/ojs/index.php/tacl/article/view/711>
- [17] Mingqing Hu and Bing Liu. 2014. Mining and Summarizing Customer Reviews. (2014), 10.
- [18] Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759* (2016).
- [19] Yoon Kim. 2014. Convolutional Neural Networks for Sentence Classification. *EMNLP* (2014), 1746–1751. <https://doi.org/10.3115/v1/D14-1181>
- [20] Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2015. Skip-thought Vectors. In *NIPS'15*. MIT Press, Cambridge, MA, USA, 3294–3302.
- [21] Xin Li and Dan Roth. 2002. Learning Question Classifiers. In *COLING: Association for Computational Linguistics*.
- [22] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781* (2013).
- [23] Matteo Pagliardini, Prakhar Gupta, and Martin Jaggi. 2017. Unsupervised Learning of Sentence Embeddings using Compositional n-Gram Features. *arXiv:1703.02507 [cs]* (March 2017). arXiv: 1703.02507.
- [24] Matteo Pagliardini, Prakhar Gupta, and Martin Jaggi. 2018. Unsupervised Learning of Sentence Embeddings Using Compositional n-Gram Features. *NAACL* 1 (2018), 528–540. <https://doi.org/10.18653/v1/N18-1049>
- [25] Bo Pang and Lillian Lee. 2005. Seeing stars: exploiting class relationships for sentiment categorization with respect to rating scales. *Association for Computational Linguistics*, 115–124. <https://doi.org/10.3115/1219840.1219855>
- [26] Jeffrey Pennington, Richard Socher, and Christopher D Manning. 2014. Glove: Global Vectors for Word Representation. In *EMNLP*, Vol. 14. 1532–1543.
- [27] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. 2013. Recursive Deep Models for Semantic Compositionality Over a Sentiment Treebank. *EMNLP* (2013), 1631–1642.
- [28] Karen Sparck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation* 28, 1 (1972), 11–21.
- [29] Benyou Wang, Peng Zhang, Jingfei Li, Dawei Song, Yuexian Hou, and Zhen-guo Shang. 2016. Exploration of quantum interference in document relevance judgement discrepancy. *Entropy* 18, 4 (2016), 144.
- [30] Janyce Wiebe, Theresa Wilson, and Claire Cardie. 2005. Annotating Expressions of Opinions and Emotions in Language. *Language Resources and Evaluation* 39, 2-3 (May 2005), 165–210. <https://doi.org/10.1007/s10579-005-7880-9>
- [31] Peng Zhang, Jiabin Niu, Zhan Su, Benyou Wang, Liqun Ma, and Dawei Song. 2018. End-to-End Quantum-like Language Models with Application to Question Answering. (2018).
- [32] Peng Zhang, Zhan Su, Lipeng Zhang, Benyou Wang, and Dawei Song. 2018. A Quantum Many-body Wave Function Inspired Language Modeling Approach. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. ACM, 1303–1312.