

The Dawn of Quantum Natural Language Processing

Riccardo Di Sipio¹, Jia-Hong Huang², Samuel Yen-Chi Chen³, Stefano Mangini⁴, Marcel Worring²

¹Ceridian HCM Inc., ²University of Amsterdam, ³Brookhaven National Laboratory, ⁴University of Pavia

1. Research Background

In early 2020, the British company called Cambridge Quantum Computing announced the introduction of a “meaning-aware” Quantum NLP model, i.e., a mathematical theory able to splice the semantic information of words with the syntactic structure of a sentence [1]. The bold claim is based on the observation that syntactic structures such as Noam Chomsky’s Context-free grammars or Categorical Compositional Distributional (DisCoCat) can be formulated in the framework of quantum physics [1]. While the concept of a Hilbert space for NLP may seem far-fetched at first, an intuitive explanation goes as follows: NLP-as-we-know-it is rooted in classical statistics (for example, word embeddings are vectors in a \mathbb{R}^d space), but classical statistics itself has its own limits. The field of physics called quantum mechanics is described by the mathematics of quantum statistics, which extends classical statistics by representing objects with matrices of complex numbers. Even if we are not claiming that words are made of particles, it does happen by accident that the kind of statistics needed to describe human language is substantially the same of quantum physics. Coincidentally, we are also entering a time in which rudimentary and still quite noisy devices able to carry out computations based on qubits rather than digital bits of information are becoming more and more accessible to the wider public (including small businesses, start-ups and individuals). Thus, it makes sense in this decade to explore the possibility that quantum computing may give a boost to natural language processing.

2. Introduction

In recent times, large pre-trained neural network models have been successfully used to achieve state-of-the-art performance in computer vision (CV) and natural language processing (NLP). In particular, the field of NLP is seeing an exponential expansion in terms of real-life applications such as machine translation, document classification, interaction with a chatbot but also network size as Transformer-based models seem to outperform long-standing architectures such as Long Short-Term Memory (LSTM) recurrent neural networks. However, the main issue with this marvellous kind of neural networks is that the appalling size of parameters (in the order of hundreds of billions in the case of OpenAI’s GPT-3) usually requires a training dataset of huge dimension such as the complete Wikipedia corpus in several languages, and a massive computer cluster to carry out the training. On the other hand, the past years have witnessed the rise of quantum computing both in terms of hardware development and implementation of algorithms on such platforms. In quantum computing, computations are carried over the equivalent of bits called qubits, which notoriously can handle information in a non-binary state thanks to a property of quantum systems called superposition. A quantum circuit is a series of operations applied to qubits that change their state, e.g., by changing their relative phase or angles in the quantum Hilbert space. Qubits can be represented geometrically with a so called Bloch sphere, so an operation onto a qubit corresponds to a rotation of the quantum state vector in this virtual space. One key concept of quantum algorithms is that, because of the intrinsic nature of qubits, certain calculations can be carried out with smaller complexity compared to the classical equivalent (this property is known as super-polynomial speedup). This is especially true in the case of chemistry, where quantum algorithms are used to predict electronic structures, ground-state energies, and the spatial configuration of proteins. In fact, quantum computers are extremely efficient in cases where CPUs may take literally forever, such as in the simulation of the Ferredoxin molecule. Practical applications are more likely to be a hybrid of classical and quantum operations since the currently available quantum computers are only small to intermediate scale. Under this hybrid paradigm, most of the calculations are carried out on classical computers, while the quantum devices are deployed to solve the computational tasks for which certain quantum advantages can be demonstrated. This hybrid approach is not too different from what has been done in the past decade with GPUs. Thus, the main idea behind Quantum Machine Learning (QML) is to replace parts of a neural network (e.g., linear layers) with a quantum counterpart. A comparison between different learning models is shown in Table 1. Finally, as classical machine learning does not end with just neural network, also quantum machine learning find its strength in a variety of methods, such as quantum support vector machines. While it is still too early to claim that quantum computing has taken over, there are some areas where it can give an advantage as for example in drug discovery or finance. We argue that one field that so far has been poorly explored in QML is NLP.

Classical	Quantum	Hybrid
bits	qubits	bits + qubits
DNN, GBDT	VQC	DNN + VQC
Broad applications	Quantum advantage	Broad applications with Quantum Advantage
Hardware-ready	Limited resources	Quantum hardware only where necessary

Table 1: Comparison between different learning models.

3. Document Classification with Variational Quantum Circuits

In the first example that we want to discuss, a deep neural network is used to provide a sentence embedding. A final layer (“head”) maps these embeddings to a probability vector of dimension (tokens, nclasses). In a classical machine learning setting, this is easily achievable by the means of a linear layer with softmax activation. In the QML framework, the calculation is carried by first putting the input qubits into the initial state (e.g., a string such as 010010), then they are entangled between each other, rotated by arbitrary angles, and finally observed (“measured”). The goal of the training is to find the rotation angles mentioned above that optimize some cost function such as the mean squared error (MSE) between the output and some pre-determined labels. This kind of quantum circuit is known as Variational Quantum Circuit (VQC). See Figure 1 for VQC architecture visualization. However, a VQC cannot change the dimensionality of the input but only rotate the state of the qubits. Hence, the VQC is “dressed”, i.e., it is sandwiched between two classical linear layers to match the dimensionality. A classical layer “squeezes” the input to match the number of qubits, and another classical layer “bloats” the output to match the dimension of hidden vectors.

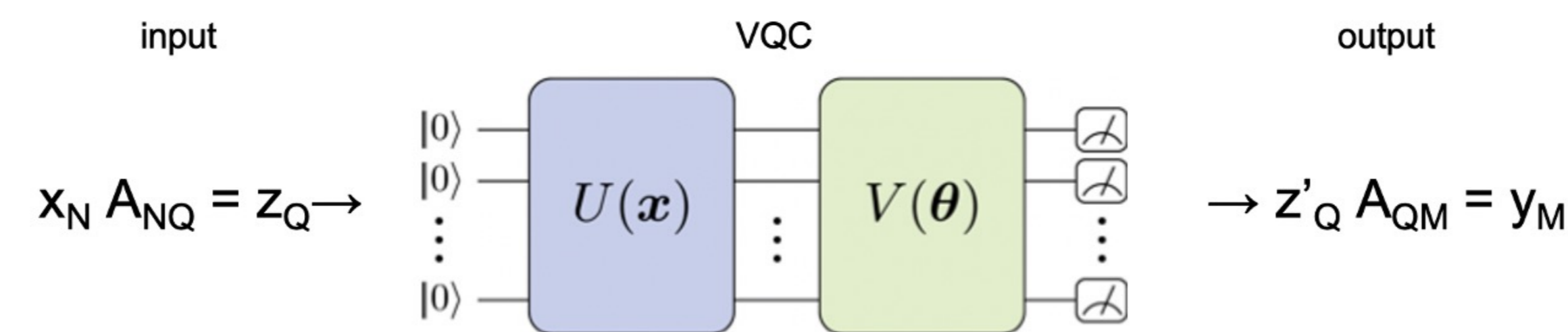


Figure 1: Generic architecture for variational quantum circuits. $U(x)$ is the quantum routine for encoding the (classical) input data x into a quantum state, and $V(\theta)$ is the variational circuit block with tunable parameters θ . A quantum measurement over some or all of the qubits follows.

4. Quantum-enhanced Long Short-Term Memory

As documents are usually presented as sequences of words, historically one of the most successful techniques to manipulate this kind of data has been the Recurrent Neural Network (RNN) architecture, and in particular a variant called Long Short-Term Memory (LSTM). The “trick” that makes these networks so popular in the analysis of sequential data is a combination of “memory” and “statefulness” that help with identifying which components of the input are relevant to compute the output. While the mathematics is quite thick as we will see later on, it suffices to say for now that LSTMs allowed machines to perform translations, classification and intent detection with state-of-the-art accuracy until the advent of Transformer networks. Still, it is interesting at least from an educational point of view to dig into LSTMs to see what good quantum computing may bring to the field. For a more thorough discussion, refer to “Quantum LSTM” and “Recurrent Quantum Neural Network”.

To begin with, let’s review the inner workings of a LSTM. We assume that the input is composed of a sequence of t time steps (e.g., words), each represented by a N -dimensional feature vector (in practical applications N can be large, e.g., 512). Also, the network stores a hidden array of vectors h and a state vector c that are updated for each element of the input sequence. For example, if we are only interested in a summary of the sequence (e.g., in a sentiment analysis), the last element of the array h will be returned. Instead, if we are interested in a representation of each element (e.g., to assign to each word a part-of-speech tag such as noun, verb, etc), we want to have access to every element of h . The calculation can be summarized in the formulas below:

$$\begin{aligned}
 f_t &= \sigma(W_f \cdot v_t + b_f) & (1) \\
 i_t &= \sigma(W_i \cdot v_t + b_i) & (2) \\
 \tilde{C} &= \tanh(W_C \cdot v_t + b_C) & (3) \\
 c_t &= f_t \odot c_{t-1} + i_t \odot \tilde{C} & (4) \\
 o_t &= \sigma(W_o \cdot v_t + b_o) & (5) \\
 h_t &= o_t \odot \tanh(c_t) & (6)
 \end{aligned}$$

Parameter	Value
epochs	300
vocab size	5
number of tags	3
embedding dim	8
hidden dim	6
no. of quantum layers	1
no. of qubits in VQC	4
Total number of weights	199

Table 2: Parameters used to define the quantum-enhanced LSTM network.

, where v_t is a concatenation of the input element at step t and the hidden state at step $t-1$, i.e. $v_t = [h_{t-1}, x_t]$. In the machine learning lingo, borrowed from electric circuit analysis, f_t is called forget gate, it is the input gate, c_t is the update gate and o_t is the output gate. The matrices W_f , W_i , W_c and W_o and the bias vectors b_f , b_i , b_c and b_o are the parameters that have been learned during the supervised training and implement the part of the calculation called linear dense layer that we want to replace with the quantum equivalent. As is often the case, a non-linearity is introduced by the application of sigmoid (σ) and hyperbolic tangent (\tanh) functions to the output of these four dense layers, which effect is to determine whether a part of the input has to be considered (values close to 1) or ignored (values close to 0). Figure 2 shows a graphical representation of the information flow and operations carried out inside a LSTM.

Notably, the linear layers of LSTM change the dimensionality of the input tensor, something that is not possible with qubit rotations. As described before, the VQC layers have to be “dressed”, i.e., sandwiched between two classical linear layers to match the dimensionality. This is applied to f_t , i_t , c_t and o_t separately. Thus, overall there is still a sizeable number of classical parameters to be learned during the training. To test the setting described above, an intuitive but not trivial example is provided: a part-of-speech tagging task. In an example, two sentences (“The dog ate the apple” and “Everybody read that book”) have been annotated with POS tags. For example, the first sentence is [“DET”, “NN”, “V”, “DET”, “NN”]. The LSTM will output the hidden array of vectors $[h_0, h_1, h_2, h_3, h_4]$, one for each word. A dense layer “head” with softmax activation is attached to the LSTM’s outputs to calculate the probability that each word may be a determinant, noun or verb. We trained the two networks (classical and quantum LSTM) for 300 epochs each. The parameters used to define the quantum network are described in Table 2. As shown in Figure 3, the cross-entropy loss decreases as a function of the training epoch, and after 150 epochs both networks are able to tag correctly the two sentences. Due to the complexity of the simulation of the quantum circuit, it took approximately 10 minutes to finish the training, to be compared to a mere 10 seconds for the classical case. Also, it seems that the quantum LSTM needed to be trained for longer epochs to achieve 100% accuracy. At the end of the 300 epochs, the loss of the classical network is 0.053 while that of the quantum network is 0.056. It is worth noticing that the quantum LSTM uses less than half the parameters of the purely classical one (199 vs. 477), while retaining the same overall performances.

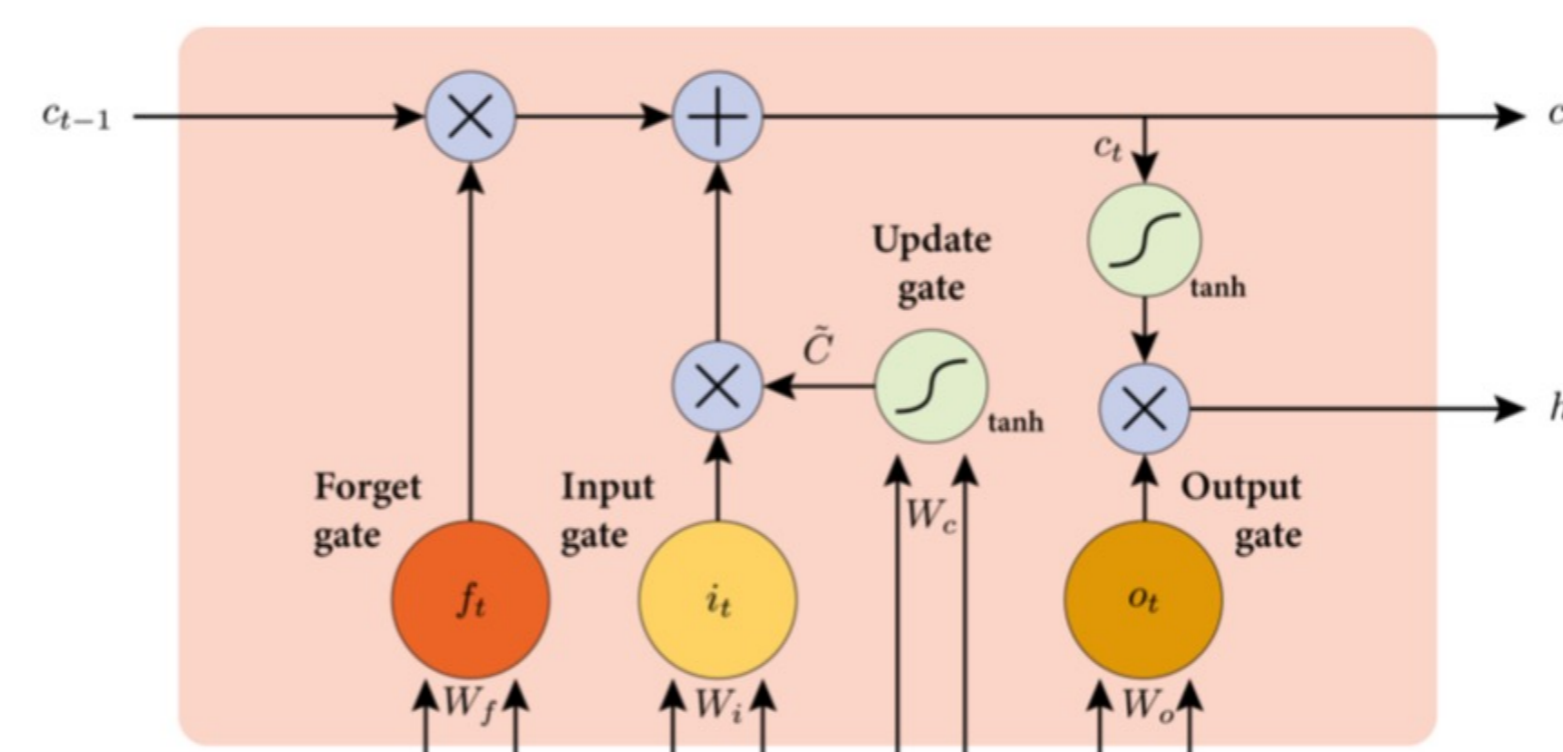


Figure 2: Graphical representation of the information flow and operations carried out inside a LSTM neural network layer. In the quantum-enhanced LSTM, each of the classical operation W_f , W_i , W_c , and W_o is replaced by a hybrid quantum-classical component which includes a VQC sandwiched between classical layers.

5. Toward A Quantum Transformer

As mentioned in the introduction, the Transformer architecture revolutionized the analysis of sequential data, and in particular that of human-written documents. Transformers are a neural network architecture that is optimized to analyze sequential data on highly-parallel devices such as GPUs and TPUs. Differently from recurrent networks, Transformers do not have “memory” but are still able to perform the trick by a combination of position-dependent embeddings (i.e., words embeddings are supplemented by another set of vectors that depend on the position the word has in the sentence, and on the index of the embedding dimension) and attention (i.e., figuring out which parts of the inputs are relevant to calculate the output). At the core of any Transformer sits the so-called Multi-Headed Attention. The idea is to apply three different linear transformations W_Q , W_K , and W_V to each element of the input sequence to transform each word embedding into some other internal representation states called Query (Q), Key (K), and Value (V). These states are then passed to the function that calculates the attention weights, which is simply defined as:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (7)$$

To promote the Transformer from the classical to quantum real, one can replace the linear transformations W_Q , W_K , and W_V with VQCs. A preliminary experiment for sentiment analysis made use of the IMDB dataset for binary classification. We used PennyLane version 0.13.0 and its plugin to perform simulations of quantum processes with the Qulacs library, which is a VQC simulator that runs on a GPU. The parameters of the hybrid network are described in Table 3. It took about 100 hours to train the classifier for a single epoch. It is clear that unless a more sophisticated training method is devised, model development and parameter tuning is prohibitive with the hardware currently available to the general public.

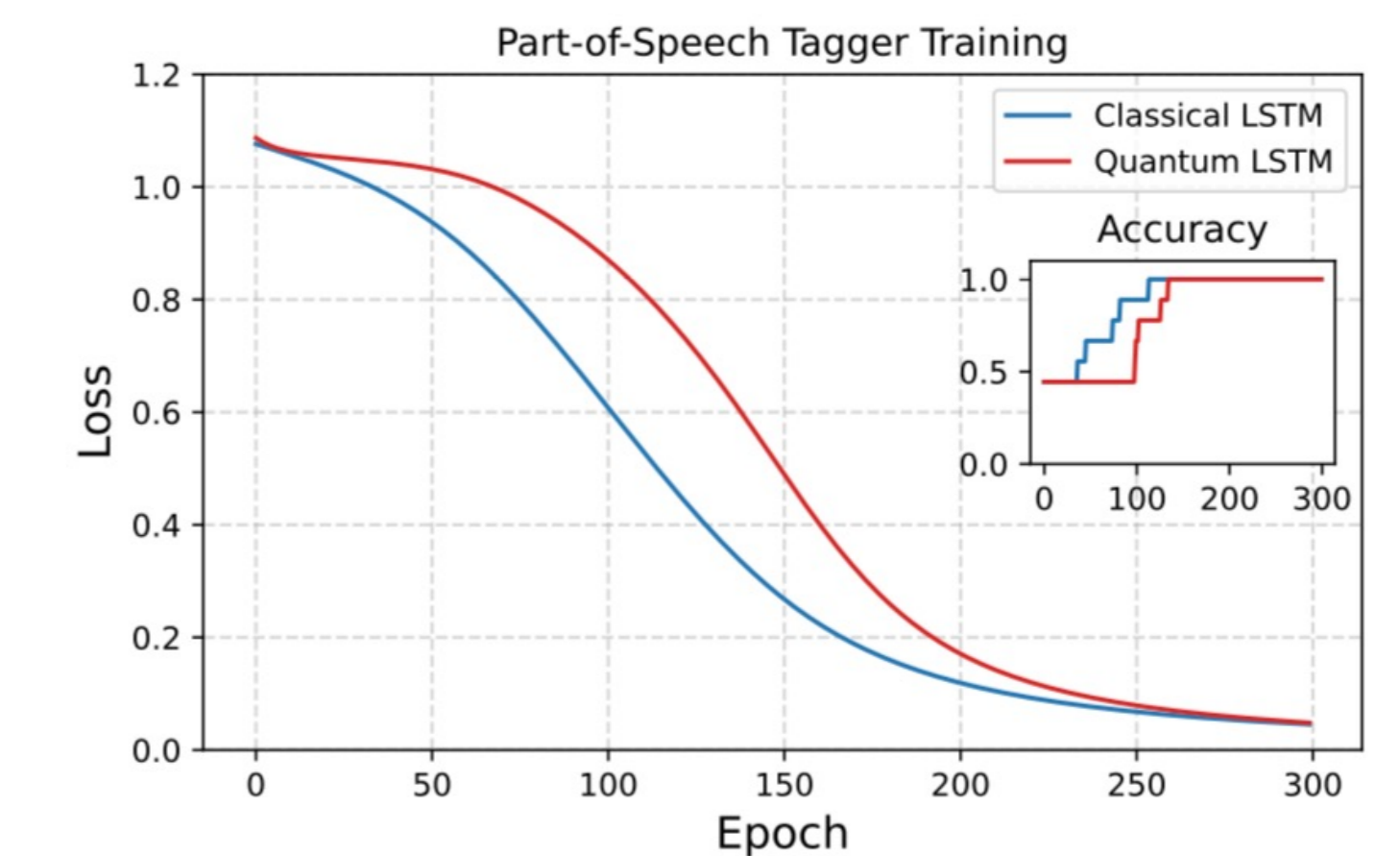


Figure 3: Cross-entropy loss and multi-class accuracy as a function of the training epoch for the classical and quantum LSTM Part-of-Speech taggers.

Parameter	Value
batch size	32
epochs	1
vocab size	50,000
embedding dim	8
max seq length	64
feed-forward net dim	8
drop-out rate	0.1
no. of transformer blocks	1
no. of transformer heads	2
no. of quantum layers	1
no. of qubits in transformer blocks	2
no. of qubits in feed-forward net dim	2
Total number of weights	$\mathcal{O}(150,000)$

Table 3: Parameters used to define the proposed quantum-enhanced transformer network.

6. Conclusion

It is a legitimate question whether it would be possible to devise a fully-quantum Transformer, i.e., a quantum circuit acting upon qubits in a fashion similar to a combination self-attention and positional embedding. While we cannot offer a definite answer here, we argue that most of the building blocks of the basic Transformer have a corresponding quantum operation, as arithmetic operations and exponentiation can be implemented as quantum circuits, or that they can be cast in form of variational quantum circuits. It is the sheer size of the computation that is likely beyond the limits of current quantum computers, although we foresee that this situation may change rapidly over the course of this decade. The end of this paper is a good place to circle back to the original question of what a Transformer actually is, and why it is so successful at interpreting human language. While a number of interpretations have been given in terms of computations on graphs, in due course quantum computing may offer alternative ways to apply attention to sequential data, or provide some completely different conceptual framework altogether. Only time will tell.

[1] Coecke B, de Felice G, Meichanetzidis K, Toumi A. 2021. “How to make qubits speak.” arXiv preprint arXiv:2107.06776.

Contact Email: riccardo.disipio@gmail.com or j.huang@uva.nl

GitHub: <https://github.com/rdisipio/qlstm> ; <https://github.com/rdisipio/qlstm>