# End-to-End Prediction of Buffer Overruns from Raw Source Code via Neural Memory Networks

**Min-je Choi** and **Sehun Jeong** and **Hakjoo Oh** and **Jaegul Choo**

Department of Computer Science and Engineering

Korea University, Seoul

{devnote5676, gifaranga, hakjoo_oh, jchoo}@korea.ac.kr

## Abstract

Detecting buffer overruns from a source code is one of the most common and yet challenging tasks in program analysis. Current approaches based on rigid rules and handcrafted features are limited in terms of flexible applicability and robustness due to diverse bug patterns and characteristics existing in sophisticated real-world software programs. In this paper, we propose a novel, data-driven approach that is completely end-to-end without requiring any hand-crafted features, thus free from any program language-specific structural limitations. In particular, our approach leverages a recently proposed neural network model called memory networks that have shown the state-of-the-art performances mainly in question-answering tasks. Our experimental results using source code samples demonstrate that our proposed model is capable of accurately detecting different types of buffer overruns. We also present in-depth analyses on how a memory network can learn to understand the semantics in programming languages solely from raw source codes, such as tracing variables of interest, identifying numerical values, and performing their quantitative comparisons.

## 1 Introduction

Detecting potential bugs in software programs has long been a challenge ever since computers were first introduced. To tackle this problem, researchers in the domain of programming languages developed various techniques called static analysis, which tries to find potential bugs in source codes without having to execute them based on a solid mathematical framework [Cousot and Cousot, 1977]. However, designing a static analyzer is tightly coupled with a particular programming language, and it is mainly based on a rigid set of rules designed by a few experts, considering numerous types of possible program states and bug cases. Thus, even with its slight syntax changes frequently found in real-world settings, e.g., several variants of ANSI C languages, a significant amount of engineering effort is required to make a previously designed analyzer applicable to the other similar languages.

To overcome these limitations, one can suggest data-driven, machine learning-based approaches as the rapid growth of deep neural networks in natural language processing has proved its effectiveness in solving similar problems such as defect predictions. Studies show that deep convolutional neural networks (CNNs) and recurrent neural networks (RNNs) are capable of learning patterns or structures within text corpora such as source codes, so they can be applied to programming language tasks such as bug localization [Lam et al., 2016], syntax error correction [Bhatia and Singh, 2016; Pu et al., 2016], and code suggestion [White et al., 2015].

Despite their impressive performances at detecting syntax-level bugs and code patterns, deep neural networks have shown less success at understanding how data values are transferred and used within source codes. This semantic level of understanding requires not only knowledge on the overall structure but also the capability to track the data values stored in different variables and methods. Although the aforementioned deep learning models may learn patterns and structures, they cannot keep track of how values are changed. This restriction greatly limits their usefulness in program analysis since run-time bugs and errors are usually much more difficult to detect and thus are often treated with greater importance.

In response, we introduce a new deep learning model with the potential of overcoming such difficulties: memory networks [Weston et al., 2015b; Sukhbaatar et al., 2015]. Memory networks are best described as neural networks with external memory 'slots' to store previously introduced information for future uses. Given a question, it accesses relevant memory slots via an attention mechanism and combines the values of the accessed slots to reach an answer. While long short-term memories (LSTMs) and earlier models also have external memories, theirs tend to evolve as longer sequences of information are fed in to the network, thus failing to fully preserve and represent information introduced at earlier stages. Memory networks on the other hand can preserve the given information even during long sequences.

This unique aspect of memory networks makes it and its variant models [Kumar et al., 2016; Henaff et al., 2016] perform exceptionally well at question answering tasks, e.g., the Facebook bAbI task [Weston et al., 2015a], a widely-used QA benchmark set. The structure of these tasks comprises a story, a query, and an answer, from which a model has to predict the correct answer to the task mentioned in the query

by accessing relevant parts of the given story. These tasks are logical questions such as locating an object, counting numbers, or basic induction/deduction. All questions can be correctly answered by referring to appropriate lines of the given story.

We point out that this task setting is in fact similar to that of a buffer overrun analysis that requires the understanding of previous lines in a source code to evaluate whether a buffer access is valid. Both tasks require knowledge not only on how each line works but also on how to select the best relevant information from previous lines. It is this very situation at which our work sets a starting point.

In this study we set the objective as demonstrating a data-driven model free of hand-crafted features and rules, and yet capable of solving tasks with the complexity of buffer overrun analyses. We present how memory networks can be effectively applied to tasks that require the understanding of not only syntactic aspects of a source code but also more complex tasks such as how values are transferred along code lines. We present how our models can learn the concept of numbers and numerical comparison simply by training on such buffer overrun tasks without any additional information. We also introduce a generated source code dataset that was used to compensate for difficulties we faced in our data-driven approach. As far as our knowledge goes, our proposed approach is the first to use deep learning to directly tackle a run-time error prediction task such as buffer overruns.

In Section 2, we cover previous work related to our task. In Section 3, we redefine our tasks, introduce our generated dataset and its purposes, and propose characteristics of the memory network model and how it is applied to this domain. In Section 4, we report experimental results and further discuss the performance of memory networks and notable characteristics it learned during the process. In Section 5 we conclude our work and discuss future work as well as the potential of memory networks for future tasks.

## 2 Related Work

To improve traditional static analysis techniques in the programming language domain, data-driven approaches based on machine learning have been recently studied. Obtaining general properties of a target program, namely, invariants, is one of the prime examples. When concrete data of target programs such as test cases or logs are available, data-driven approaches can be used to identify general properties [Sharma *et al.*, 2012; Sharma *et al.*, 2013b; Sharma *et al.*, 2013a; Sankaranarayanan *et al.*, 2008b; Sankaranarayanan *et al.*, 2008a; Nori and Sharma, 2013], similar to static analysis techniques. This use case is particularly useful when a target program has inherent complexity that makes contemporary static analyzers to compromise either of precision and cost, but is bundled with test cases that can cover most of cases.

Meanwhile, following the upsurge in the developing field of neural computing and deep learning, many models have been applied to natural language texts, especially in identifying language structure and patterns. Socher et al. [Socher *et al.*, 2013] introduced recursive neural networks which parse a sentence into subsections. [Sutskever *et al.*, 2014] proposed

RNNs that learn structures of long text sequences. Source codes of a program can also be seen as a text corpus with its own grammar structure, thus being applicable for such neural network models. [Karpathy *et al.*, 2015] showed that a character-level LSTM trained with Linux kernel codes is capable of detecting features such as brackets or sentence length as well as generating simulated codes that greatly resemble actual ones in syntactic structure. Motivated by such results in the pattern discovery of source codes, several approaches have been taken to solve practical issues in source code analysis. [Pu *et al.*, 2016] and [Bhatia and Singh, 2016] gathered data from programming assignments submitted for a MOOC class to train a correction model which corrects syntax errors in assignments. [Huo *et al.*, 2016] and [Lam *et al.*, 2016] applied attention-based CNN models to detect buggy source codes. [Allamanis *et al.*, 2014] learned coding styles by searching for patterns with neural networks. While these approaches proved that neural networks are capable of detecting patterns within codes, they are limited to detecting only syntax errors or bugs and not the transition of values stored inside variables or functions of a source code program.

Neural networks with external memories have shown better performances in inference or logical tasks compared to contemporary models. Following the introduction of neural Turing machines [Graves *et al.*, 2014] and memory networks [Weston *et al.*, 2015b; Sukhbaatar *et al.*, 2015], many variants of these models were applied to various tasks other than QA tasks such as sentiment analysis, part-of-speech tagging [Kumar *et al.*, 2016], and information extraction from documents [Miller *et al.*, 2016]. Yet, so far there has been no work that applies a memory network-based model to tasks with the complexity of semantic analysis in source codes.

## 3 Model Description

In this section, we first provide the rationale for solving buffer overruns as a QA task. We also introduce a source code-based training dataset that we designed. Lastly, we describe the structure of our model which is based on the memory network [Sukhbaatar *et al.*, 2015] and how it predicts buffer overruns from a source code.

### 3.1 Benchmark Source Code Generation



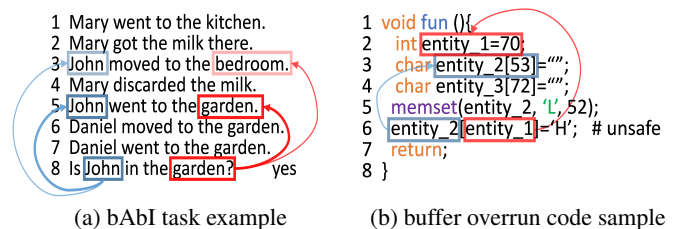(a) bAbI task example          (b) buffer overrun code sample

Figure 1: Comparison of a bAbI and a buffer overrun tasks

**Comparison of bAbI tasks and buffer overruns**. We return to our statement that analyzing buffer overruns is similar to solving bAbI tasks. Consider Fig. 1. The bAbI task shown in Fig. 1(a) is given a story (lines 1-7) and a query (line 8). A solver model understands this task by looking at 'John'
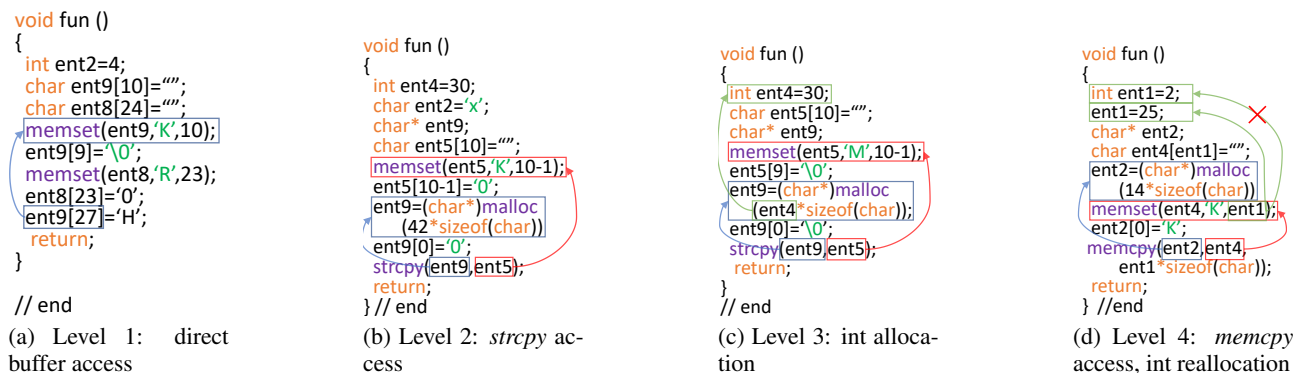
```
void fun ()
{
  int ent2=4;
  char ent9[10]="";
  char ent8[24]="";
  memset(ent9,'K',10);
  ent9[9]='\0';
  memset(ent8,'R',23);
  ent8[23]='0';
  ent9[27]='H';
  return;
}

// end
```

(a) Level 1: direct buffer access

```
void fun ()
{
  int ent4=30;
  char ent2='x';
  char* ent9;
  char ent5[10]="";
  memset(ent5,'K',10-1);
  ent5[10-1]='0';
  ent9=(char*)malloc
    (42*sizeof(char))
  ent9[0]='0';
  strcpy(ent9,ent5);
  return;
} // end
```

(b) Level 2: strcpy access

```
void fun ()
{
  int ent4=30;
  char ent5[10]="";
  char* ent9;
  memset(ent5,'M',10-1);
  ent5[9]='\0';
  ent9=(char*)malloc
    (ent4*sizeof(char));
  ent9[0]='\0';
  strcpy(ent9,ent5);
  return;
} // end
```

(c) Level 3: int allocation

```
void fun ()
{
  int ent1=2;
  ent1=25;
  char* ent2;
  char ent4[ent1]="";
  ent2=(char*)malloc
    (14*sizeof(char))
  memset(ent4,'K',ent1);
  ent2[0]='K';
  memcpy(ent2,ent4,
    ent1*sizeof(char));
  return;
} //end
```

(d) Level 4: memcpy access, int reallocation

Figure 2: Different levels of buffer overrun tasks

and 'where' from the query and then attends the story to find lines related to 'John.' Lines 3 and 5 are chosen as candidates. The model understands from the sequential structure that line 5 contains more recent, thus relevant information. In the end, the model returns the answer 'garden' by combining the query and the information from line 5.

Meanwhile, the task of Fig. 1(b) is to discriminate whether the buffer access made at line 6 is valid. Our analyzer first understands that its objective is to compare the size of the character array *entity_2* and the integer variable *entity_1*. Next, it searches for the length of *entity_2* at line 3, where 53 is allocated to the variable. It also gains knowledge from line 2 that *entity_1* is equivalent to 70. The remaining task is to compare the integer variables 53 and 70 and return an alarm (*unsafe*) if the index exceeds the length of the character array. One can think of lines 1-5 as a story and line 6 as a query, perfectly transforming this problem into a bAbI task.

**Limitations of test suites**. Although test suites such as Juliet Test Suite for C programming language [Boland and Black, 2012] are designed for benchmarking buffer overrun and other program analysis tasks, the data is not diverse enough. Code samples differ by only a small fraction such as a different variable nested in a conditional statement or loop, while a large portion of code appears repeatedly over several samples. A data-driven model will inevitably learn from only the small variations and ignore a large portion of the code where much of the valuable information is stored.

**Program structure**. We tackle this problem of data inadequacy by generating our own training source code dataset.[1] Our dataset adopts buffer access functions and initialization methods from Juliet to maintain at least an equal level of task complexity, while also preserving an underlying structure that makes it applicable for deep learning approaches. Each sample is a void function of 10 to 30 lines of C code and consists of three stages: initialization, allocation, and query. During the initialization stage, variables are initialized as either characters, character arrays, or integers. At the allocation stage, these variables are assigned values using randomly generated integers between 1 to 100. Buffer sizes are allocated to character arrays with *malloc* and *memset* functions. At the

query stage, a buffer access is attempted on one of the allocated character arrays via a direct access on an array index (Fig. 2(a)). We formulate this task into a binary classification problem where an 'unsafe' sign is returned if a character array is accessed with a string or index that exceeds its size.

**Naming conventions**. We assume that a limited number of individual variables appear in each program sample. Each variable is given the name *entity_n* where $n \in \{i | 0 \leq i \leq N_{upper}, i \in Z\}$ and $N_{upper}$ is an integer set by default to 10. Each $n$ is assigned randomly to variables and invariant of their introduced order or data type. One can imagine a situation where an agent (variable) is given a fake ID (entity name) for a particular task (sample). The agent learns to complete the task with that fake ID, then discards it upon task completion, and selects a new one for the subsequent task. In this manner, we can prevent *entities* from learning task-specific knowledge and instead train them as representations of universal variables which can replace any kind of variable that appears in a program. We can easily apply our model to real-life source codes using this naming convention by simply changing the names of newly introduced variables and methods to different *entities*.

**Adding complexity**. Our model has to adapt to more realistic source codes with complex structures. Possible settings that complicate our task include

- selecting only the appropriate variables out of several dummy variables,
- introducing different buffer access methods requiring the comparison of two character arrays such as *strcpy* or *memcpy* functions,
- allocating the sizes of character arrays not with integers but indirectly with previously assigned integer variables,
- reallocating integer variables prior to or after their use in allocating a character array.

We first assign a number of dummy variables to each sample program. Each dummy variable is initialized and allocated in the same manner as the ones actually used in the buffer access. We include the use of *strcpy* (Fig. 2(b)) / *memcpy* (Fig. 2(d)) functions for buffer accesses. We also add cases where character arrays are allocated not directly by integers, but indirectly with additionally introduced integer variables (Fig. 2(c)). Given this setting, the model has to learn to store the integer value allocated to the integer variable first, then
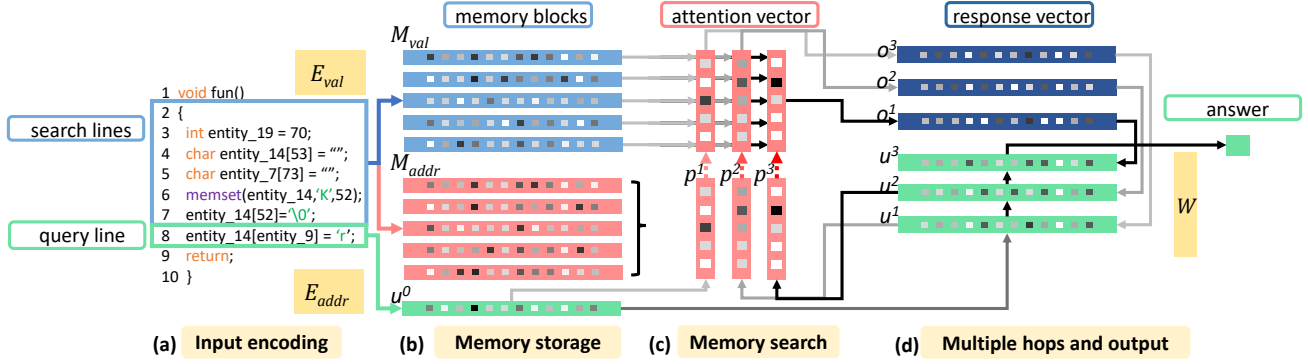
Figure 3: Our proposed memory network-based model for buffer overrun tasks

use that value to obtain the length of the character array. Finally, we add further cases where the additional integer variable *itself* is reallocated (Fig. 2(d)), either before or after it is used to define the character array length. Now the model has to learn to choose whether the previously assigned or reallocated value was used for allocating a character array.

Our generated source codes are equivalent to an expanded version of Juliet test suite in terms of data flow, that is, the flow of data values defined and used within a source code. Compared to codes in Juliet which only define the source and destination variables that will be used for buffer access, ours include dummy variables which are defined and used similarly. The inclusion of various settings such as memory allocation using assigned integer variables instead of raw integers and reassignment of variables also increase data flow. In overall, our source codes provide a tougher environment for models to solve buffer overrun tasks than the existing Juliet dataset as there are more variations to consider.

## 3.2 Model Structure

The overall structure of our model is displayed in Fig. 3.

**Input encoding (Fig. 3(a)).** The memory network takes in as input a program code $X$ consisting of $n$ search lines $X_1, X_2, \cdots, X_n$ and a single buffer access line or query $X_q$. A single program line $X_m$ is a list of words $w_m^1, w_m^2, \cdots, w_m^l$. With $V$ as the vocabulary size, we define $x_m^l$ as the $V$-dimensional one-hot vector representation of a word $w_m^l$. We set an upper limit $N$ for the max number of lines a memory can store, and we pad zeros for the remaining lines if a program is shorter than $N$ lines.

Note that every word in the source code is treated as a word token. This includes not only variable names ($entity$), type definitions ($int$) and special characters, ('[', '*'), but also integers as well. This setting matches our concept of an end-to-end model that does not require explicit parsing. While it is possible to apply parsers to extract numbers and represent them differently from other word tokens, this would contradict our goal of applying a purely data-driven approach. Treating integers as individual word tokens means that our model will not be given any prior information regarding the size differences between numbers, and thus has to learn such numerical concepts by itself. We further discuss this in Section 4.

Next, we compute vector representations for each sentence using its words. Each word is represented in a $d$-dimensional vector using an embedding matrix $E_{val} \in R^{d \times V}$. We also multiply a column vector $l_j$ to the $j$-th word vector for each word to allocate different weights according to word positions. This concept known as position encoding [Sukhbaatar *et al.*, 2015] enables our model to discriminate the different roles of variables when two or more identical words appear in a single sentence. Without such settings, our model may fail to discriminate between source and destination variables such as in a *strcpy* function. The memory representation $m_i$ of line $i$ consisting of $J$ words and the $k$-th element $l_j^k$ of the position encoding vector $l_j \in \mathbb{R}^d$ for word $j$ in the line $i$ are obtained as

$$m_i = \Sigma_{j=1}^t l_j \cdot A x_i^j, \qquad (1)$$

$$l_j^k = (1 - j/J) - (k/d)(1 - 2j/J), \qquad (2)$$

where '·' is element-wise multiplication.

**Memory storage (Fig. 3(b)).** Next, we allocate our encoded sentences $m_i$ into matrices called memory blocks. Fig. 3(b) shows two memory blocks, the memory value block $\left( M_{val} \in R^{N \times d} \right)$ and the memory address block $\left( M_{addr} \in R^{N \times d} \right)$. Each sentence is allocated into one row of memory block, namely a memory slot. $M_{val}$ stores semantical information about the contents of a code line while $M_{addr}$ stores information for locating how much to address each line. For this reason, sentences are encoded using two different word embedding matrices, $E_{val}$ and $E_{addr}$ for $M_{val}$ and $M_{addr}$, respectively.

**Memory search (Fig. 3(c)).** The query is encoded into a representation using $E_{addr}$. We denote the initial query embedding as $u^0$. By computing the inner products between the query embedding and each slot of the memory address block, then applying a softmax function to the resulting vector, we obtain the attention vector $p$ which indicates how related each line is to the query. The $i$-th element of $p$ is obtained as

$$p_i = \text{softmax}\left( \left( u^0 \right)^T M_{addr} \right), \qquad (3)$$

with the softmax function as

$$\text{softmax}(z_i) = e^{z_i} / \Sigma_j e^{z_j}. \qquad (4)$$

The response vector $o$ is computed as in

$$o = \Sigma_i p_i \left( M_{val} \right)_i. \qquad (5)$$

Table 1: Comparison on generated source codes. Inside brackets are the standard deviations

| | level 1 | | | level 2 | | | level 3 | | | level 4 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | acc | F1 | auc | acc | F1 | auc | acc | F1 | auc | acc | F1 | auc |
| CNN | 0.67 | 0.69 | 0.75 | 0.73 | 0.71 | 0.81 | 0.61 | 0.61 | 0.66 | 0.62 | 0.62 | 0.67 |
| | (0.01) | (0.02) | (0.01) | (0.02) | (0.02) | (0.01) | (0.01) | (0.03) | (0.01) | (0.01) | (0.03) | (0.01) |
| LSTM | 0.8 | 0.84 | 0.92 | 0.82 | 0.80 | 0.90 | 0.69 | 0.66 | 0.76 | 0.67 | 0.64 | 0.75 |
| | (0.01) | (0.01) | (0.00) | (0.01) | (0.02) | (0.01) | (0.00) | (0.01) | (0.01) | (0.01) | (0.02) | (0.01) |
| Memory network | 0.84 | 0.84 | 0.92 | 0.86 | 0.85 | 0.93 | 0.83 | 0.83 | 0.90 | 0.82 | 0.82 | 0.90 |
| | (0.01) | (0.01) | (0.01) | (0.01) | (0.02) | (0.01) | (0.02) | (0.02) | (0.02) | (0.02) | (0.02) | (0.02) |

This vector contains information collected over all lines of the memory value block according to their attention weights obtained from the memory address block. This is equivalent to searching the memory for different parts of information with respect to a given query.

**Multiple hops and output (Fig. 3(d))**. The response vector $o$ can be either directly applied to a weight matrix $W$ to produce an output, or added to strengthen the query $u$. In the latter case, the query is updated as in Eq. (6) by simply adding the response vector to the previous query embedding.

$$u^{k+1} = u^k + o^k \qquad (6)$$

We repeat from Eq. (3) to obtain a new response vector. Our model iterates through multiple hops where at each hop the desired information to be obtained from the memory slightly changes. This accounts for situations where a model has to first look for lines where an array is allocated, and then gather information from lines stating the size of the variables used for allocating the array size. The final output is a floating value ranging from 0 (unsafe) to 1 (safe), which we round to the nearest integer to obtain a binary prediction result.

## 4 Experiments

In this section, we present both quantitative and qualitative results of our experiments on model performance and learned characteristics.

### 4.1 Quantitative Evaluation

**Experiment settings**. Our main dataset consists of C-style source codes discussed in Section 3.2. We used a single training set consisting of 10,000 sample programs. We generated four test sets with 1,000 samples each and assigned them levels one to four, with a higher level indicating a more complex condition (see Table 2). Samples ranged from 8 to 33 lines of code, with an average of 16.01. A total of 196 unique words appeared in the training set. A maximum of four dummy variables were added to each sample. We used random integers between 0 and 100 for buffer allocation and access. We conducted every experiment on a Intel (R) Xeon (R) CPU E5-2687W v3 @ 3.10GHz machine equipped with two GeForce GTX TITAN X GPUs. All models were implemented with Tensorflow 0.12.1 using Python 2.7.1 on an Ubuntu 14.04 environment.

**Model Comparison**. We set our memory network to three hops with a memory of 30 lines and the embedding size of $d = 32$. As there has been no previous work on using deep learning models for such tasks, we used existing deep learning models often used for text classification tasks as baselines. That is, we included a CNN for text classification [Kim, 2014] and a two-layer LSTM binary classifier. All models were trained with Adam [Kingma and Ba, 2014] at a learning rate of 1e-2. We used the classification accuracy, F1 score, and the area under the ROC curve (AUC) as performance metrics. We averaged the scores of the ten best cases with the smallest training error.

Table 2: Different levels of test sets

| | Level 1 | Level 2 | Level 3 | Level 4 |
|---|---|---|---|---|
| Direct buffer access | √ | √ | √ | √ |
| Access by *strcpy* / *memcpy* | | √ | √ | √ |
| Allocation by int variable | | | √ | √ |
| Reallocation of int variable | | | | √ |

Performance results shown in Table 2 demonstrate that all models decrease in performance as task levels increase, due to our level assignment. Of the deep learning models, only memory networks performed consistently at a high level on all four level settings with accuracy rates higher than 80%. This is expected since their hops allow them to solve even complex situations such as variable reallocation and different buffer access types. Meanwhile, CNNs failed to complete even the simplest tasks since they cannot capture the sequential information in input sentences and instead apply convolutional filters to words of all regions on an equal basis. Any positional information is discarded.

Memory networks require substantially shorter computing time compared to other models, requiring an average of 0.63s per epoch, while LSTMs and CNNs each require 7.13s and 13.21s. As for the latter models, the number of computations is proportional to the number of words that appear in a code. However, memory networks sum up all words of a sentence to form a single representation, and thus computation time relies on the number of lines instead of individual words. This significantly reduces computation time.

Interestingly, LSTM models also performed well when set to easier tasks. Results show that LSTMs performed comparably to memory networks, even equaling them on Level 1 tasks. However, its performance sharply dropped when used on higher level tasks. This partial success of LSTMs relates to the simple structure of Level 1 tasks. The size of the character array always appears before the index to access, so the model can cheat by comparing the only two numbers that appear within the entire code. This cheating becomes obsolete

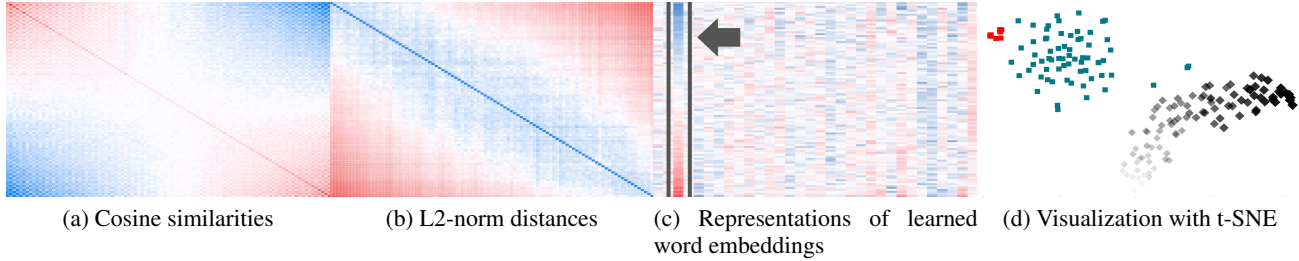| (a) Cosine similarities | (b) L2-norm distances | (c) Representations of learned word embeddings | (d) Visualization with t-SNE |

Figure 4: Visualizations of word embedding vectors of numbers 1-100. Red and blue indicate high and low values, respectively.

as higher-level tasks require knowledge only obtainable by attending previous lines in a stepwise manner.

## 4.2 Qualitative Analysis



Figure 5: Prediction result with attention per hop

We further examine the performance of our memory network model and the steps it takes to obtain a correct answer. We also present visualization results on how our model learns the concepts of numbers and numerical comparison without being explicitly supervised about such tasks.

**Tracking hop-wise results.** In order to prove that our model solves the tasks in our desired manner, that is, by attending and collecting relevant information from different parts of the memory at different hops, we analyze individual prediction cases by inspecting which parts of information our model has obtained from taking each hop.

Fig. 5 displays an example of buffer overrun analysis using our model. We can observe that when given a *strcpy* buffer access as a query, the model's initial attention shifts to the sentence where the destination buffer (*entity_3*) is allocated. The model decides here to next look for *entity_9*, which contains the size used for allocating to *entity_3*. During the next hop it attends the line where the source buffer (*entity_2*) is allocated and obtains data of 99, the size of *entity_2*. At the last hop the memory network visits entity_9 and obtains 69. After the three hops, the destination size 69 is compared with source size 99, and being a smaller number, returns 'unsafe' as a result. The prediction confidence in Fig. 5 indicates how close the predicted value is to the ground answer.

**Numerical concepts automatically learned.** Recall from Section 3 that our model was not given any prior information regarding the notion of quantitative values. Interestingly, our model learned to compare between different numbers. Fig. 4 displays visualization results using only the word embedding vectors corresponding to the 100 numbers.

Figs. 4(a) and (b) display the cosine similarities and the L2-norm distances of all numbers from 1 to 100, with 1 at the topmost left-hand side. The colors observed at the first and third quadrants from both figures show that numbers with large differences are trained to minimize cosine similarities while maximizing L2-norm distances, thus spacing themselves apart. In contrast, similar numbers in the second and fourth quadrants have opposite characteristics, meaning they are similarly placed.

The word embedding vectors of numbers across all $d$ dimensions as seen in Fig. 4(c) further demonstrate a clear sequential order between numbers. The highlighted column forms a strong color spectrum starting from a low value which gradually increases as the corresponding number increases from 1 to 100. As all word embeddings were initialized with random values at the beginning, this spectrum indicates that our model learns by itself to assign such values for comparison purposes.

Last of all, Fig. 4(d) is a t-SNE representation of all word embedding vectors. The black gradation indicates the word embeddings of numbers, with denser colors indicating larger numbers. We notice that they are embedded in a consistent direction in an increasing order. While this again shows how our model learns numerical characteristics, we also discover that dots in red, which correspond to *entities* from Section 3, stand out. As mentioned earlier, *entities* correspond to the variables that appear in source codes as integer variables or character buffers. This implies that our model learns to train word embeddings differently according to their purposes within a code.

## 5 Conclusions and Future Work

In this work, we proposed a memory network-based model for predicting buffer overruns in programming language analysis. Our work is the first to apply a deep learning-based approach to a problem in the field of program analysis that requires both syntactic and semantic knowledge. Performance results show that memory networks are superior to other models in solving buffer overrun tasks across all difficulty levels. We also presented that our model successfully learns the notion of numbers and their quantitative comparisons from merely textual data in an end-to-end setting.

Our work has room to improve in many interesting aspects. We can expand our model to cover different program analysis tasks such as pointer analysis, interval analysis, and flow-

sensitivity analysis, which share similar semantic natures. We can apply advanced variants of memory networks to handle various conditions in source codes such as *if* and *for* statements. Our knowledge of models learning numerical representations can further aid deep learning models compatible with arithmetic and logical reasoning. All of these combined, our work marks a stepping stone to a fully data-driven program analyzer.

# References

[Allamanis *et al.*, 2014] Miltiadis Allamanis, Earl T. Barr, Christian Bird, and Charles Sutton. Learning natural coding conventions. In *Proceedings of the 22Nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*, FSE 2014, pages 281–293, New York, NY, USA, 2014. ACM.

[Bhatia and Singh, 2016] Sahil Bhatia and Rishabh Singh. Automated correction for syntax errors in programming assignments using recurrent neural networks. *CoRR*, abs/1603.06129, 2016.

[Boland and Black, 2012] T. Boland and P. E. Black. Juliet 1.1 c/c++ and java test suite. *Computer*, 45(10):88–90, Oct 2012.

[Cousot and Cousot, 1977] Patrick Cousot and Radhia Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Proceedings of the 4th ACM SIGACT-SIGPLAN symposium on Principles of programming languages*, pages 238–252. ACM, 1977.

[Graves *et al.*, 2014] Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *CoRR*, abs/1410.5401, 2014.

[Henaff *et al.*, 2016] Mikael Henaff, Jason Weston, Arthur Szlam, Antoine Bordes, and Yann LeCun. Tracking the world state with recurrent entity networks. *CoRR*, abs/1612.03969, 2016.

[Huo *et al.*, 2016] Xuan Huo, Ming Li, and Zhi-hua Zhou. Learning unified features from natural and programming languages for locating buggy source code. In *Proceedings of the 25th International Joint Conference on Artificial Intelligence (IJCAI 2016)*, pages 1606–1612, 2016.

[Karpathy *et al.*, 2015] Andrej Karpathy, Justin Johnson, and Li Fei-Fei. Visualizing and understanding recurrent networks. *CoRR*, abs/1506.02078, 2015.

[Kim, 2014] Yoon Kim. Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP 2014)*, pages 1746–1751, 2014.

[Kingma and Ba, 2014] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.

[Kumar *et al.*, 2016] Ankit Kumar, Ozan Irsoy, Peter Ondruska, Mohit Iyyer, James Bradbury, Ishaan Gulrajani, Victor Zhong, Romain Paulus, and Richard Socher. Ask me anything: Dynamic memory networks for natural language processing. In *ICML*, 2016.

[Lam *et al.*, 2016] An Ngoc Lam, Anh Tuan Nguyen, Hoan Anh Nguyen, and Tien N. Nguyen. Combining deep learning with information retrieval to localize buggy files for bug reports. In *Proceedings - 2015 30th IEEE/ACM International Conference on Automated Software Engineering, ASE 2015*, pages 476–481. IEEE, 2016.

[Miller *et al.*, 2016] Alexander H. Miller, Adam Fisch, Jesse Dodge, Amir-Hossein Karimi, Antoine Bordes, and Jason Weston. Key-value memory networks for directly reading documents. *CoRR*, abs/1606.03126, 2016.

[Nori and Sharma, 2013] Aditya V Nori and Rahul Sharma. Termination proofs from tests. In *Proceedings of the 2013 9th Joint Meeting on Foundations of Software Engineering*, pages 246–256. ACM, 2013.

[Pu *et al.*, 2016] Yewen Pu, Karthik Narasimhan, Armando Solar-Lezama, and Regina Barzilay. sk-p: a neural program corrector for moocs. In *Companion Proceedings of the 2016 ACM SIGPLAN International Conference on Systems, Programming, Languages and Applications: Software for Humanity*, pages 39–40, New York, New York, USA, 2016. ACM Press.

[Sankaranarayanan *et al.*, 2008a] Sriram Sankaranarayanan, Swarat Chaudhuri, Franjo Ivančić, and Aarti Gupta. Dynamic inference of likely data preconditions over predicates by tree learning. In *Proceedings of the 2008 international symposium on Software testing and analysis*, pages 295–306. ACM, 2008.

[Sankaranarayanan *et al.*, 2008b] Sriram Sankaranarayanan, Franjo Ivančić, and Aarti Gupta. Mining library specifications using inductive logic programming. In *Proceedings of the 30th international conference on Software engineering*, pages 131–140. ACM, 2008.

[Sharma *et al.*, 2012] Rahul Sharma, Aditya V Nori, and Alex Aiken. Interpolants as classifiers. In *International Conference on Computer Aided Verification*, pages 71–87. Springer, 2012.

[Sharma *et al.*, 2013a] Rahul Sharma, Saurabh Gupta, Bharath Hariharan, Alex Aiken, Percy Liang, and Aditya V Nori. A data driven approach for algebraic loop invariants. In *European Symposium on Programming*, pages 574–592. Springer, 2013.

[Sharma *et al.*, 2013b] Rahul Sharma, Saurabh Gupta, Bharath Hariharan, Alex Aiken, and Aditya V Nori. Verification as learning geometric concepts. In *International Static Analysis Symposium*, pages 388–411. Springer, 2013.

[Socher *et al.*, 2013] Richard Socher, Alex Perelygin, and Jy Wu. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing (EMNLP 2013)*, pages 1631–1642, 2013.

[Sukhbaatar *et al.*, 2015] Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, and Rob Fergus. End-to-end memory networks. In *Advances in Neural Information Processing Systems*, pages 1–11, 2015.

[Sutskever *et al.*, 2014] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, pages 3104–3112, 2014.

[Weston *et al.*, 2015a] Jason Weston, Antoine Bordes, Sumit Chopra, and Tomas Mikolov. Towards ai-complete question answering: A set of prerequisite toy tasks. *CoRR*, abs/1502.05698, 2015.

[Weston *et al.*, 2015b] Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks. In *Int'l Conf. on Learning Representations*, pages 1–15, 2015.

[White *et al.*, 2015] Martin White, Christopher Vendome, Mario Linares-Vásquez, and Denys Poshyvanyk. Toward deep learning software repositories. In *Mining Software Repositories (MSR), 2015 IEEE/ACM 12th Working Conference on*, pages 334–345. IEEE, 2015.