

# Stochastic Answer Networks for Machine Reading Comprehension

Xiaodong Liu<sup>†</sup>, Yelong Shen<sup>†</sup>, Kevin Duh<sup>‡</sup> and Jianfeng Gao<sup>†</sup>

<sup>†</sup> Microsoft Research, Redmond, WA, USA

<sup>‡</sup> Johns Hopkins University, Baltimore, MD, USA

<sup>†</sup>{xiaodl, yeshen, jfgao}@microsoft.com <sup>‡</sup>kevinduh@cs.jhu.edu

## Abstract

We propose a simple yet robust stochastic answer network (SAN) that simulates multi-step reasoning in machine reading comprehension. Compared to previous work such as ReasoNet which used reinforcement learning to determine the number of steps, the unique feature is the use of a kind of stochastic prediction dropout on the answer module (final layer) of the neural network during the training. We show that this simple trick improves robustness and achieves results competitive to the state-of-the-art on the Stanford Question Answering Dataset (SQuAD), the Adversarial SQuAD, and the Microsoft Machine Reading Comprehension Dataset (MS MARCO).

## 1 Introduction

Machine reading comprehension (MRC) is a challenging task: the goal is to have machines read a text passage and then answer any question about the passage. This task is an useful benchmark to demonstrate natural language understanding, and also has important applications in e.g. conversational agents and customer service support. It has been hypothesized that difficult MRC problems require some form of multi-step synthesis and reasoning. For instance, the following example from the MRC dataset SQuAD (Rajpurkar et al., 2016) illustrates the need for synthesis of information across sentences and multiple steps of reasoning:

*Q:* What collection does **the V&A Theator & Performance galleries** hold?

*P:* **The V&A Theator & Performance galleries** opened in March 2009. ... **They** hold the UK’s biggest national collection of

material about live performance.

To infer the answer (the underlined portion of the passage  $P$ ), the model needs to first perform coreference resolution so that it knows “**They**” refers “**V&A Theator**”, then extract the subspan in the direct object corresponding to the answer.

This kind of iterative process can be viewed as a form of multi-step reasoning. Several recent MRC models have embraced this kind of multi-step strategy, where predictions are generated after making multiple passes through the same text and integrating intermediate information in the process. The first models employed a predetermined fixed number of steps (Hill et al., 2016; Dhingra et al., 2016; Sordani et al., 2016; Kumar et al., 2015). Later, Shen et al. (2016) proposed using reinforcement learning to dynamically determine the number of steps based on the complexity of the question. Further, Shen et al. (2017) empirically showed that dynamic multi-step reasoning outperforms fixed multi-step reasoning, which in turn outperforms single-step reasoning on two distinct MRC datasets (SQuAD and MS MARCO).

In this work, we derive an alternative multi-step reasoning neural network for MRC. During training, we fix the number of reasoning steps, but perform stochastic dropout on the answer module (final layer predictions). During decoding, we generate answers based on the average of predictions in all steps, rather than the final step. We call this a stochastic answer network (SAN) because the stochastic dropout is applied to the answer module; albeit simple, this technique significantly improves the robustness and overall accuracy of the model. Intuitively this works because while the model successively refines its prediction over multiple steps, each step is still trained to generate the same answer; we are performing a kind of stochastic ensemble over the model’s successive predic-

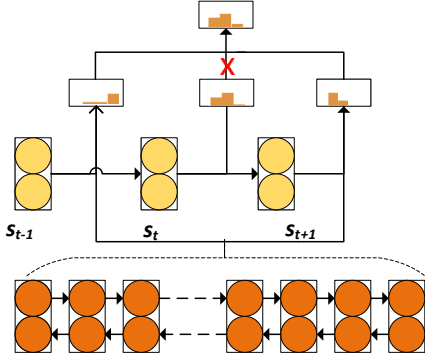


Figure 1: Illustration of “stochastic prediction dropout” in the answer module during training. At each reasoning step  $t$ , the model combines memory (bottom row) with hidden states  $s_{t-1}$  to generate a prediction (multinomial distribution). Here, there are three steps and three predictions, but one prediction is dropped and the final result is an average of the remaining distributions.

tion refinements. Stochastic prediction dropout is illustrated in Figure 1.

## 2 Proposed model: SAN

The machine reading comprehension (MRC) task as defined here involves a question  $Q = \{q_0, q_1, \dots, q_{m-1}\}$  and a passage  $P = \{p_0, p_1, \dots, p_{n-1}\}$  and aims to find an answer span  $A = \{a_{start}, a_{end}\}$  in  $P$ . We assume that the answer exists in the passage  $P$  as a contiguous text string. Here,  $m$  and  $n$  denote the number of tokens in  $Q$  and  $P$ , respectively. The learning algorithm for reading comprehension is to learn a function  $f(Q, P) \rightarrow A$ . The training data is a set of the query, passage and answer tuples  $\langle Q, P, A \rangle$ .

We now describe our model from the ground up. The main contribution of this work is the answer module, but in order to understand what goes into this module, we will start by describing how  $Q$  and  $P$  are processed by the lower layers. Note the lower layers also have some novel variations that are not used in previous work. As shown in Figure 2, our model contains four different layers to capture different concept of representations. The detailed description of our model is provided as follows.

**Lexicon Encoding Layer.** The purpose of the first layer is to extract information from  $Q$  and  $P$  at the word level and normalize for lexical vari-

ants. A typical technique to obtain lexicon embedding is concatenation of its word embedding with other linguistic embedding such as those derived from Part-Of-Speech (POS) tags. For word embeddings, we use the pre-trained 300-dimensional GloVe vectors (Pennington et al., 2014) for the both  $Q$  and  $P$ . Following Chen et al. (2017), we use three additional types of linguistic features for each token  $p_i$  in the passage  $P$ :

- 9-dimensional POS tagging embedding for total 56 different types of the POS tags.
- 8-dimensional named-entity recognizer (NER) embedding for total 18 different types of the NER tags. We utilized small embedding sizes for POS and NER to reduce model size. They mainly serve the role of coarse-grained word clusters.
- A 3-dimensional binary *exact* match feature defined as  $f_{exact\_match}(p_i) = \mathbb{I}(p_i \in Q)$ . This checks whether a passage token  $p_i$  matches the original, lowercase or lemma form of any question token.
- Question enhanced passages word embeddings:  $f_{align}(p_i) = \sum_j \gamma_{i,j} g(GloVe(q_j))$ , where  $g(\cdot)$  is a 280-dimensional single layer neural network  $ReLU(W_0x)$  and  $\gamma_{i,j} = \frac{\exp(g(GloVe(p_j)) \cdot g(GloVe(q_i)))}{\sum_{j'} \exp(g(GloVe(p_i)) \cdot g(GloVe(q_{j'})))}$  measures the similarity in word embedding space between a token  $p_i$  in the passage and a token  $q_j$  in the question. Compared to the *exact* matching features, these embeddings encode *soft* alignments between similar but not-identical words.

In summary, each token  $p_i$  in the passage is represented as a 600-dimensional vector and each token  $q_j$  is represented as a 300-dimensional vector.

Due to different dimensions for the passages and questions, in the next layer two different bidirectional LSTM (BiLSTM) (Hochreiter and Schmidhuber, 1997) may be required to encode the contextual information. This, however, introduces a large number of parameters. To prevent this, we employ an idea inspired by (Vaswani et al., 2017): use two separate two-layer position-wise Feed-Forward Networks (FFN),  $FFN(x) = W_2 ReLU(W_1x + b_1) + b_2$ , to map both the passage and question lexical encodings into the same number of dimensions. Note that this FFN has fewer

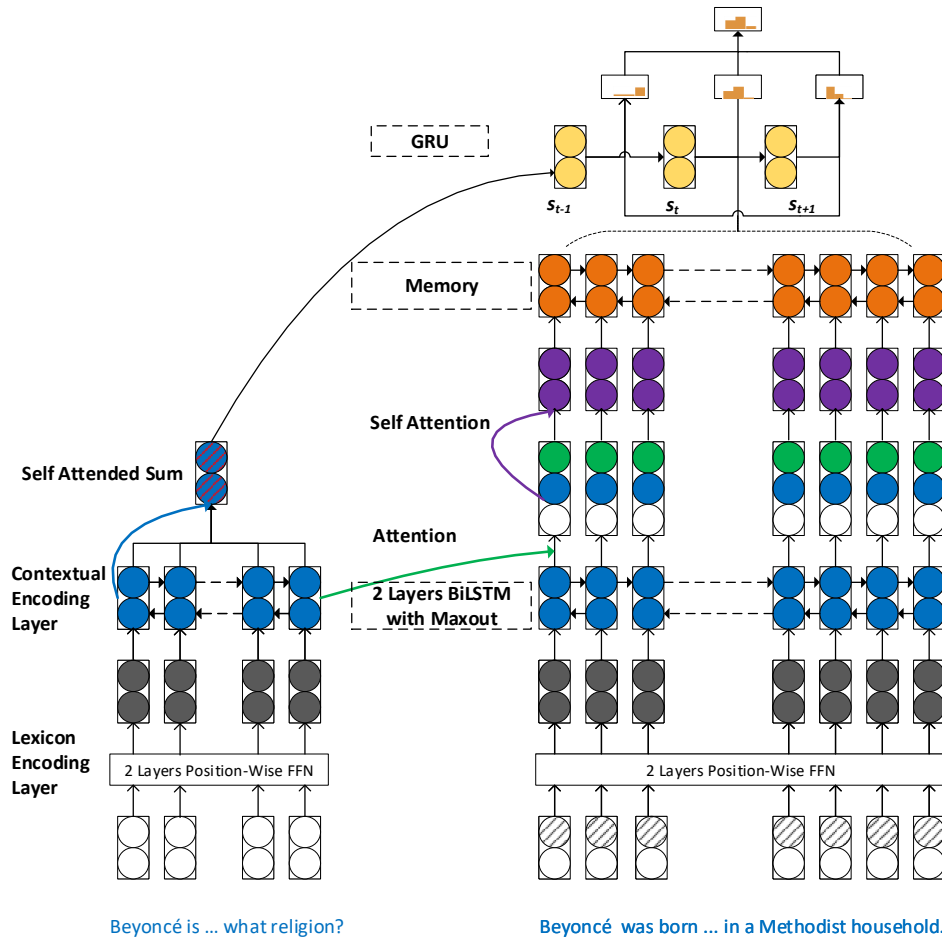


Figure 2: **Architecture of the SAN for Reading Comprehension:** The first layer is a lexicon encoding layer that maps words to their embeddings independently for the question (left) and the passage (right): this is a concatenation of word embeddings, POS embeddings, etc. followed by a position-wise FFN. The next layer is a context encoding layer, where a BiLSTM is used on the top of the lexicon embedding layer to obtain the context representation for both question and passage. In order to reduce the parameters, a maxout layer is applied on the output of BiLSTM. The third layer is the working memory: First we compute an alignment matrix between the question and passage using an attention mechanism, and use this to derive a question-aware passage representation. Then we concatenate this with the context representation of passage and the word embedding, and employ a self attention layer to re-arrange the information gathered. Finally, we use another LSTM to generate a working memory for the passage. At last, the fourth layer is the answer module, which is a GRU that outputs predictions at each state  $s_t$ .

parameters compared to a BiLSTM. Thus, we obtain the final lexicon embeddings for the tokens in  $Q$  as a matrix  $E^q \in \mathbb{R}^{d \times m}$  and tokens in  $P$  as  $E^p \in \mathbb{R}^{d \times n}$ .

**Contextual Encoding Layer.** Both passage and question use a shared two-layers BiLSTM as the contextual encoding layer, which projects the lexicon embeddings to contextual embeddings. We concatenate a pre-trained 600-dimensional CoVe vectors<sup>1</sup> (McCann et al., 2017) trained on German-English machine translation dataset, with

<sup>1</sup><https://github.com/salesforce/cove>

the aforementioned lexicon embeddings as the final input of the contextual encoding layer, and also with the output of the first contextual encoding layer as the input of its second encoding layer. To reduce the parameter size, we use a maxout layer (Goodfellow et al., 2013) at each BiLSTM layer to shrink its dimension. By a concatenation of the outputs of two BiLSTM layers, we obtain  $H^q \in \mathbb{R}^{2d \times m}$  as representation of  $Q$  and  $H^p \in \mathbb{R}^{2d \times n}$  as representation of  $P$ , where  $d$  is the hidden size of the BiLSTM.

**Memory Generation Layer.** In the memory

generation layer, We construct the working memory, a summary of information from both  $Q$  and  $P$ . First, a dot-product attention is adopted like in (Vaswani et al., 2017) to measure the similarity between the tokens in  $Q$  and  $P$ . Instead of using a scalar to normalize the scores as in (Vaswani et al., 2017), we use one layer network to transform the contextual information of both  $Q$  and  $P$ :

$$C = dropout(f_{attention}(\hat{H}^q, \hat{H}^p)) \in \mathbb{R}^{m \times n} \quad (1)$$

$C$  is an attention matrix. Note that  $\hat{H}^q$  and  $\hat{H}^p$  is transformed from  $H^q$  and  $H^p$  by one layer neural network  $ReLU(W_3x)$ , respectively. Next, we gather all the information on passages by a simple concatenation of its contextual information  $H^p$  and its question-aware representation  $H^q \cdot C$ :

$$U^p = concat(H^p, H^q C) \in \mathbb{R}^{4d \times n} \quad (2)$$

Typically, a passage may contain hundred of tokens, making it hard to learn the long dependencies within it. Inspired by (Lin et al., 2017), we apply a self-attended layer to rearrange the information  $U^p$  as:

$$\hat{U}^p = U^p drop_{diag}(f_{attention}(U^p, U^p)). \quad (3)$$

In other words, we first obtain an  $n \times n$  attention matrix with  $U^p$  onto itself, apply dropout, then multiply this matrix with  $U^p$  to obtain an updated  $\hat{U}^p$ . Instead of using a penalization term as in (Lin et al., 2017), we dropout the diagonal of the similarity matrix forcing each token in the passage to align to other tokens rather than itself.

At last, the working memory is generated by using another BiLSTM based on all the information gathered:

$$M = BiLSTM([U^p; \hat{U}^p]) \quad (4)$$

where the semicolon mark ; indicates the vector/matrix concatenation operator.

**Answer module.** There is a Chinese proverb that says: “wisdom of masses exceeds that of any individual.” Unlike other multi-step reasoning models, which only uses a single output either at the last step or some dynamically determined final step, our answer module employs all the outputs of multiple step reasoning. Intuitively, by applying dropout, it avoids a “step bias problem” (where models places too much emphasis one particular step’s predictions) and forces the model to produce good predictions at every individual step. Further,

during decoding, we reuse *wisdom of masses* instead of *individual* to achieve a better result. We call this method “stochastic prediction dropout” because dropout is being applied to the final predictive distributions.

Formally, our answer module will compute over  $T$  memory steps and output the answer span. This module is a memory network and has some similarities to other multi-step reasoning networks: namely, it maintains a state vector, one state per step. At the beginning, the initial state  $s_0$  is the summary of the  $Q$ :  $s_0 = \sum_j \alpha_j H_j^q$ , where  $\alpha_j = \frac{\exp(w_4 \cdot H_j^q)}{\sum_{j'} \exp(w_4 \cdot H_{j'}^q)}$ . At time step  $t$  in the range of  $\{1, 2, \dots, T - 1\}$ , the state is defined by  $s_t = GRU(s_{t-1}, x_t)$ . Here,  $x_t$  is computed from the previous state  $s_{t-1}$  and memory  $M$ :  $x_t = \sum_j \beta_j M_j$  and  $\beta_j = softmax(s_{t-1} W_5 M)$ . Finally, a bilinear function is used to find the begin and end point of answer spans at each reasoning step  $t \in \{0, 1, \dots, T - 1\}$ .

$$P_t^{begin} = softmax(s_t W_6 M) \quad (5)$$

$$P_t^{end} = softmax([s_t; \sum_j P_{t,j}^{begin} M_j] W_7 M). \quad (6)$$

From a pair of begin and end points, the answer string can be extracted from the passage. However, rather than output the results (start/end points) from the final step (which is fixed at  $T - 1$  as in Memory Networks or dynamically determined as in ReasoNet), we utilize all of the  $T$  outputs by averaging the scores:

$$P^{begin} = avg([P_0^{begin}, P_1^{begin}, \dots, P_{T-1}^{begin}]) \quad (7)$$

$$P^{end} = avg([P_0^{end}, P_1^{end}, \dots, P_{T-1}^{end}]) \quad (8)$$

Each  $P_t^{begin}$  or  $P_t^{end}$  is a multinomial distribution over  $\{1, \dots, n\}$ , so the average distribution is straightforward to compute.

During training, we apply stochastic dropout to before the above averaging operation. For example, as illustrated in Figure 1, we randomly delete several steps’ predictions in Equations 7 and 8 so that  $P^{begin}$  might be  $avg([P_1^{begin}, P_3^{begin}])$  and  $P^{end}$  might be  $avg([P_0^{end}, P_3^{end}, P_4^{end}])$ . The use of averaged predictions and dropout during training improves robustness.

Our stochastic prediction dropout is similar in motivation to the dropout introduced by (Srivastava et al., 2014). The difference is that theirs

is dropout at the intermediate node-level, whereas ours is dropout at the final layer-level. Dropout at the node-level prevents correlation between features. Dropout at the final layer level, where randomness is introduced to the averaging of predictions, prevents our model from relying exclusively on a particular step to generate correct output. We used a dropout rate of 0.4 in experiments.

### 3 Experiment Setup

**Dataset:** We evaluate on the Stanford Question Answering Dataset (SQuAD) (Rajpurkar et al., 2016). This contains about 23K passages and 100K questions. The passages come from approximately 500 Wikipedia articles and the questions and answers are obtained by crowdsourcing. The crowdsourced workers are asked to read a passage (a paragraph), come up with questions, then mark the answer span. All results are on the official development set, unless otherwise noted.

Two evaluation metrics are used: Exact Match (EM), which measures the percentage of span predictions that matched any one of the ground truth answer exactly, and Macro-averaged F1 score, which measures the average overlap between the prediction and the ground truth answer.

**Implementation details:** The spaCy tool<sup>2</sup> is used to tokenize the both passages and questions, and generate lemma, part-of-speech and named entity tags. We use 2-layer BiLSTM with  $d = 128$  hidden units for both passage and question encoding. The mini-batch size is set to 32 and Adamax (Kingma and Ba, 2014) is used as our optimizer. The learning rate is set to 0.002 at first and decreased by half after every 10 epochs. We set the dropout rate for all the hidden units of LSTM, and the answer module output layer to 0.4. To prevent degenerate output, we ensure that at least one step in the answer module is active during training.

### 4 Results

The main experimental question we would like to answer is whether the stochastic dropout and averaging in the answer module is an effective technique for multi-step reasoning. To do so, we fixed all lower layers and compared different architectures for the answer module:

1. Standard 1-step: generate prediction from  $s_0$ , the first initial state.

<sup>2</sup><https://spacy.io>

2. 5-step memory network: this is a memory network fixed at 5 steps. We try two variants: the standard variant outputs result from the final step  $s_{T-1}$ . The averaged variant outputs results by averaging across all 5 steps, and is like SAN without the stochastic dropout.
3. ReasoNet<sup>3</sup>: this answer module dynamically decides the number of steps and outputs results conditioned on the final step.
4. SAN: proposed answer module that uses stochastic dropout and prediction averaging.

The main results in terms of EM and F1 are shown in Table 1. We observe that SAN achieves 76.235 EM and 84.056 F1, outperforming all other models. Standard 1-step model only achieves 75.139 EM and dynamic steps (via ReasoNet) achieves only 75.355 EM. SAN also outperforms a 5-step memory net with averaging, which implies averaging predictions is not the only thing that led to SAN’s superior results; indeed, stochastic prediction dropout is an effective technique.

The K-best oracle results is shown in Figure 3. The K-best spans are computed by ordering the spans according to their probabilities  $P^{begin} \times P^{end}$ . We limit K in the range 1 to 4 and then pick the span with the best EM or F1 as oracle. SAN also outperforms the other models in terms of K-best oracle scores. Impressively, these models achieve human performance at  $K = 2$  for EM and  $K = 3$  for F1.

Finally, we compare our results with other top models in Table 2. Note that all the results in Table 2 are taken from the published papers. We see that SAN is very competitive in both single and ensemble settings (ranked in second) despite its simplicity. Note that the best-performing model (Peters et al., 2018) used a large-scale language model as an extra contextual embedding, which gave a significant improvement (+4.3% dev F1). We expect significant improvements if we add this to SAN in future work.

<sup>3</sup>The ReasoNet here is not an exact re-implementation of (Shen et al., 2017). The answer module is the same as (Shen et al., 2017) but the lower layers are set to be the same as SAN, 5-step memory network, and standard 1-step as described in Figure 2. We only vary the answer module in our experiments for a fair comparison.

Answer Module	EM	F1
Standard 1-step	75.139	83.367
Fixed 5-step with Memory Network (prediction from final step)	75.033	83.327
Fixed 5-step with Memory Network (prediction averaged from all steps)	75.256	83.215
Dynamic steps (max 5) with ReasoNet	75.355	83.360
Stochastic Answer Network (SAN ), Fixed 5-step	<b>76.235</b>	<b>84.056</b>

Table 1: **Main results**—Comparison of different answer module architectures. Note that SAN performs best in both Exact Match and F1 metrics.

Ensemble model results:	Dev Set (EM/F1)	Test Set (EM/F1)
BiDAF + Self Attention + ELMo (Peters et al., 2018)	-/-	<b>81.003/87.432</b>
<b>SAN (Ensemble model)</b>	78.619/85.866	79.608/86.496
AIR-FusionNet (Huang et al., 2017)	-/-	78.978/86.016
DCN+ (Xiong et al., 2017)	-/-	78.852/85.996
M-Reader (Hu et al., 2017)	-/-	77.678/84.888
Conductor-net (Liu et al., 2017b)	74.8 / 83.3	76.996/84.630
r-net (Wang et al., 2017)	77.7/83.7	76.9/84.0
ReasoNet++ (Shen et al., 2017)	75.4/82.9	75.0/82.6
<i>Individual model results:</i>		
BiDAF + Self Attention + ELMo(Peters et al., 2018)	-/-	<b>78.580/85.833</b>
<b>SAN (single model)</b>	<b>76.235/84.056</b>	76.828/84.396
AIR-FusionNet(Huang et al., 2017)	75.3/83.6	75.968/83.900
RaSoR + TR (Salant and Berant, 2017)	-/-	75.789/83.261
DCN+(Xiong et al., 2017)	74.5/83.1	75.087/83.081
r-net(Wang et al., 2017)	72.3/80.6	72.3/80.7
ReasoNet++(Shen et al., 2017)	70.8/79.4	70.6/79.36
BiDAF (Seo et al., 2016)	67.7/77.3	68.0/77.3
Human Performance	80.3/90.5	82.3/91.2

Table 2: Test performance on SQuAD. Results are sorted by Test F1.

## 5 Analysis

### 5.1 How robust are the results?

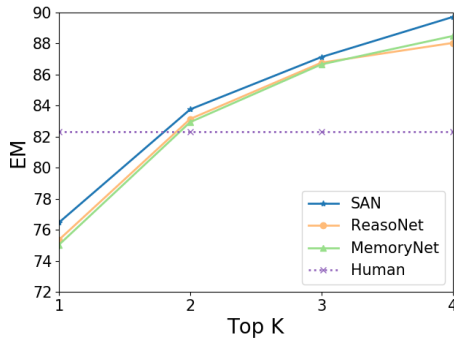
We are interested in whether the proposed model is sensitive to different random initial conditions. Table 3 shows the development set scores of SAN trained from initialization with different random seeds. We observe that the SAN results are consistently strong regardless of the 10 different initializations. For example, the mean EM score is 76.131 and the lowest EM score is 75.922, both of which still outperform the 75.355 EM of the Dynamic step ReasoNet in Table 1.<sup>4</sup>

We are also interested in how sensitive are the results to the number of reasoning steps, which

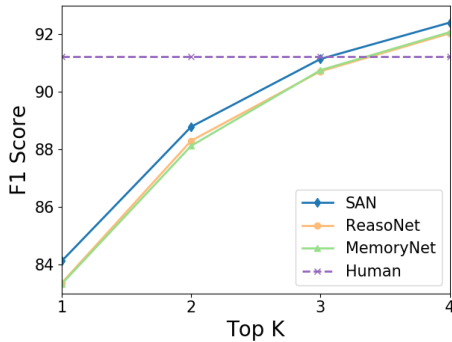
<sup>4</sup>Note the Dev EM/F1 scores of ReasoNet in Table 1 do not match those of ReasoNet++ in Table 2. While the answer module is the same architecture, the lower encoding layers are different.

is a fixed hyper-parameter. Since we are using dropout, a natural question is whether we can extend the number of steps to an extremely large number. Table 4 shows the development set scores for  $T = 1$  to  $T = 10$ . We observe that there is a gradual improvement as we increase  $T = 1$  to  $T = 5$ , but after 5 steps the improvements have saturated. In fact, the EM/F1 scores drop slightly, but considering that the random initialization results in Table 3 show a standard deviation of 0.142 and a spread of 0.426 (for EM), we believe that the  $T = 10$  result does not statistically differ from the  $T = 5$  result. In summary, we think it is useful to perform some approximate hyper-parameter tuning for the number of steps, but it is not necessary to find the exact optimal value.

Finally, we test SAN on two Adversarial SQuAD datasets, AddSent and AddOneSent (Jia and Liang, 2017), where the passages contain



(a) EM comparison on different systems.



(b) F1 score comparison on different systems.

Figure 3: K-Best Oracle results

auto-generated adversarial distracting sentences to fool computer systems that are developed to answer questions about the passages. For example, AddSent is constructed by adding sentences that look similar to the question, but do not actually contradict the correct answer. AddOneSent is constructed by appending a random human-approved sentence to the passage.

We evaluate the single SAN model (i.e., the one presented in Table 2) on both AddSent and AddOneSent. The results in Table 5 show that SAN achieves the new state-of-the-art performance and SAN’s superior result is mainly attributed to the multi-step answer module, which leads to significant improvement in F1 score over the Standard 1-step answer module, i.e., +1.2 on AddSent and +0.7 on AddOneSent.

## 5.2 Is it possible to use different numbers of steps in test vs. train?

For practical deployment scenarios, prediction speed at test time is an important criterion. Therefore, one question is whether SAN can train with, e.g.  $T = 5$  steps but test with  $T = 1$  steps. Table 6 shows the results of a SAN trained on  $T = 5$  steps, but tested with different number of steps. As ex-

Seed#	EM	F1	Seed#	EM	F1
<b>Seed 1</b>	76.24	84.06	Seed 6	76.23	83.99
Seed 2	76.30	<b>84.13</b>	Seed 7	<b>76.35</b>	84.09
Seed 3	<b>75.92</b>	83.90	Seed 8	76.07	83.71
Seed 4	76.00	83.95	Seed 9	75.93	<b>83.85</b>
Seed 5	76.12	83.99	Seed 10	76.15	84.11

Mean: 76.131, Std. deviation: 0.142 (EM)

Mean: 83.977, Std. deviation: 0.126 (F1)

Table 3: **Robustness of SAN (5-step) on different random seeds for initialization:** best and worst scores are boldfaced. Note that our official submit is trained on seed 1.

SAN	EM	F1	SAN	EM	F1
1 step	<b>75.38</b>	<b>83.29</b>	6 step	75.99	83.72
2 step	75.43	83.41	7 step	76.04	83.92
3 step	75.89	83.57	8 step	76.03	83.82
4 step	75.92	83.85	9 step	75.95	83.75
5 step	<b>76.24</b>	<b>84.06</b>	10 step	76.04	83.89

Table 4: **Effect of number of steps:** best and worst results are boldfaced.

pected, the results are best when  $T$  matches during training and test; however, it is important to note that small numbers of steps  $T = 1$  and  $T = 2$  nevertheless achieve strong results. For example, prediction at  $T = 1$  achieves 75.58, which outperforms a standard 1-step model (75.14 EM) as in Table 1 that has approximate equivalent prediction time.

## 5.3 How does the training time compare?

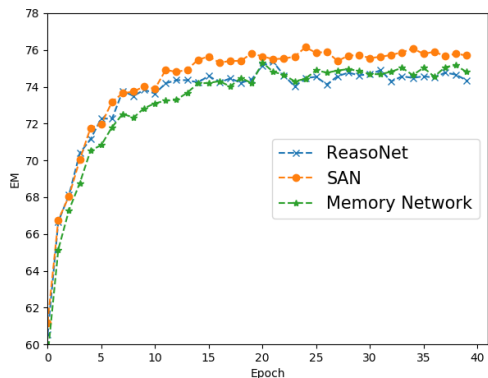
The average training time per epoch is comparable: our implementation running on a GTX Titan X is 22 minutes for 5-step memory net, 30 minutes for ReasoNet, and 24 minutes for SAN. The learning curve is shown in Figure 4. We observe that all systems improve at approximately the same rate up to 10 or 15 epochs. However, SAN continues to improve afterwards as other models start to saturate. This observation is consistent with previous works using dropout (Srivastava et al., 2014). We believe that while training time per epoch is similar between SAN and other models, it is recommended to train SAN for more epochs in order to achieve gains in EM/F1.

Single model:	AddSent	AddOneSent
LR (Rajpurkar et al., 2016)	23.2	30.3
SED <sub>T</sub> (Liu et al., 2017a)	33.9	44.8
BiDAF (Seo et al., 2016)	34.3	45.7
jNet (Zhang et al., 2017)	37.9	47.0
ReasonNet(Shen et al., 2017)	39.4	50.3
RaSoR(Lee et al., 2016)	39.5	49.5
Mnemonic(Hu et al., 2017)	<b>46.6</b>	56.0
QANet(Yu et al., 2018)	45.2	55.7
Standard 1-step in Table 1	45.4	55.8
SAN	<b>46.6</b>	<b>56.5</b>

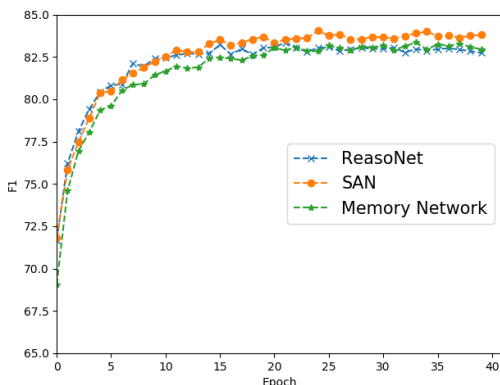
Table 5: Test performance on the adversarial SQuAD dataset in F1 score.

$T =$	EM	F1	$T =$	EM	F1
1	75.58	83.86	4	76.12	83.98
2	75.85	83.90	5	76.24	84.06
3	75.98	83.95	10	75.89	83.88

Table 6: Prediction on different steps  $T$ . Note that the SAN model is trained using 5 steps.



(a) EM



(b) F1

Figure 4: Learning curve measured on Dev set.

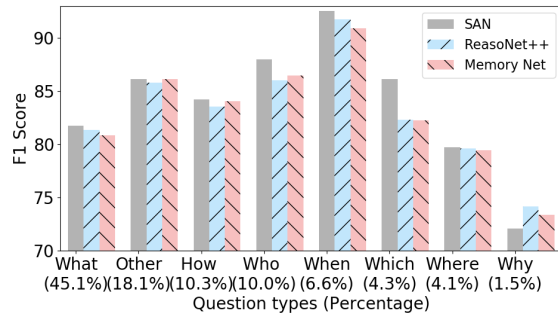


Figure 5: Score breakdown by question type.

#### 5.4 How does SAN perform by question type?

To see whether SAN performs well on a particular type of question, we divided the development set by questions type based on their respective Wh-word, such as “who” and “where”. The score breakdown by F1 is shown in Figure 5. We observe that SAN seems to outperform other models uniformly across all types. The only exception is the Why questions, but there is too little data to derive strong conclusions.

#### 5.5 Experiments results on MS MARCO

MS MARCO (Nguyen et al., 2016) is a large scale real-word RC dataset which contains 100,100 (100K) queries collected from anonymized user logs from the Bing search engine. The characteristic of MS MARCO is that all the questions are real user queries and passages are extracted from real web documents. For each query, approximate 10 passages are extracted from public web documents. The answers are generated by humans. The data is partitioned into a 82,430 training, a 10,047 development and 9,650 test tuples. The evaluation metrics are BLEU(Papineni et al., 2002) and ROUGE-L (Lin, 2004) due to its free-form text answer style. To apply the same RC model, we search for a span in MS MARCO’s passages that maximizes the ROUGE-L score with the raw free-form answer. It has an upper bound of 93.45 BLEU and 93.82 ROUGE-L on the development set.

The MS MARCO dataset contains multiple passages per query. Our model as shown in Figure 2 is developed to generate answer from a single passage. Thus, we need to extend it to handle multiple passages. Following (Shen et al., 2017), we take two steps to generate an answer to a query  $Q$  from  $J$  passages,  $P^1, \dots, P^J$ . First, we run SAN on ev-



<i>SingleModel</i>	ROUGE	BLEU
ReasonNet++(Shen et al., 2017)	38.01	38.62
V-Net(Wang et al., 2018)	45.65	-
Standard 1-step in Table 1	42.30	42.39
SAN	<b>46.14</b>	<b>43.85</b>

Table 7: **MS MARCO devset results.**

ery  $(P^j, Q)$  pair, generating  $J$  candidate answer spans, one from each passage. Then, we multiply the SAN score of each candidate answer span with its relevance score  $r(P^j, Q)$  assigned by a passage ranker, and output the span with the maximum score as the answer. In our experiments, we use the passage ranker described in (Liu et al., 2018)<sup>5</sup>. The ranker is trained on the same MS MARCO training data, and achieves 37.1 p@1 on the development set.

The results in Table 7 show that SAN outperforms V-Net (Wang et al., 2018) and becomes the new state of the art<sup>6</sup>.

## 6 Related Work

The recent big progress on MRC is largely due to the availability of the large-scale datasets (Rajpurkar et al., 2016; Nguyen et al., 2016; Richardson et al., 2013; Hill et al., 2016), since it is possible to train large end-to-end neural network models. In spite of the variety of model structures and attention types (Bahdanau et al., 2015; Chen et al., 2016; Xiong et al., 2016; Seo et al., 2016; Shen et al., 2017; Wang et al., 2017), a typical neural network MRC model first maps the symbolic representation of the documents and questions into a neural space, then search answers on top of it. We categorize these models into two groups based on the difference of the answer module: single-step and multi-step reasoning. The key difference between the two is what strategies are applied to search the final answers in the neural space.

A single-step model matches the question and document only once and produce the final answers. It is simple yet efficient and can be trained using the classical back-propagation algorithm, thus it is adopted by most systems (Chen et al., 2016; Seo et al., 2016; Wang et al., 2017; Liu et al., 2017b; Chen et al., 2017; Weissenborn et al., 2017;

<sup>5</sup>It is the same model structure as (Liu et al., 2018) by using softmax over all candidate passages. A simple baseline, TF-IDF, obtains 20.1 p@1 on MS MARCO development.

<sup>6</sup>The official evaluation on MS MARCO on test is closed, thus here we only report the results on the development set.

Hu et al., 2017). However, since humans often solve question answering tasks by re-reading and re-digesting the document multiple times before reaching the final answers (this may be based on the complexity of the questions/documents), it is natural to devise an iterative way to find answers as multi-step reasoning.

Pioneered by (Hill et al., 2016; Dhingra et al., 2016; Sordoni et al., 2016; Kumar et al., 2015), who used a predetermined fixed number of reasoning steps, Shen et al (2016; 2017) showed that multi-step reasoning outperforms single-step ones and dynamic multi-step reasoning further outperforms the fixed multi-step ones on two distinct MRC datasets (SQuAD and MS MARCO). But these models have to be trained using reinforcement learning methods, e.g., policy gradient, which are tricky to implement due to the instability issue. Our model is different in that we fix the number of reasoning steps, but perform stochastic dropout to prevent step bias. Further, our model can also be trained by using the back-propagation algorithm, which is simple and yet efficient.

## 7 Conclusion

We introduce Stochastic Answer Networks (SAN), a simple yet robust model for machine reading comprehension. The use of stochastic dropout in training and averaging in test at the answer module leads to robust improvements on SQuAD, outperforming both fixed step memory networks and dynamic step ReasonNet. We further empirically analyze the properties of SAN in detail. The model achieves results competitive with the state-of-the-art on the SQuAD leaderboard, as well as on the Adversarial SQuAD and MS MARCO datasets. Due to the strong connection between the proposed model with memory networks and ReasonNet, we would like to delve into the theoretical link between these models and its training algorithms. Further, we also would like to explore SAN on other tasks, such as text classification and natural language inference for its generalization in the future.

## Acknowledgments

We thank Pengcheng He, Yu Wang and Xinying Song for help to set up dockers. We also thank Pranav Samir Rajpurkar for help on SQuAD evaluations, and the anonymous reviewers for valuable discussions and comments.

## References

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. *International Conference on Learning Representations (ICLR2015)*.
- Danqi Chen, Jason Bolton, and Christopher D. Manning. 2016. A thorough examination of the cnn/daily mail reading comprehension task. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Berlin, Germany, pages 2358–2367. <http://www.aclweb.org/anthology/P16-1223>.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to answer open-domain questions. In *Association for Computational Linguistics (ACL)*.
- Bhuvan Dhingra, Hanxiao Liu, William W Cohen, and Ruslan Salakhutdinov. 2016. Gated-attention readers for text comprehension. *arXiv preprint arXiv:1606.01549*.
- Ian J Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron Courville, and Yoshua Bengio. 2013. Maxout networks. *arXiv preprint arXiv:1302.4389*.
- Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. 2016. The goldilocks principle: Reading children’s books with explicit memory representations. *ICLR*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9(8):1735–1780.
- Minghao Hu, Yuxing Peng, and Xipeng Qiu. 2017. Mnemonic reader for machine comprehension. *arXiv preprint arXiv:1705.02798*.
- Hsin-Yuan Huang, Chenguang Zhu, Yelong Shen, and Weizhu Chen. 2017. Fusionnet: Fusing via fully-aware attention with application to machine comprehension. *arXiv preprint arXiv:1711.07341*.
- Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. pages 2021–2031.
- Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Ankit Kumar, Ozan Irsoy, Jonathan Su, James Bradbury, Robert English, Brian Pierce, Peter Ondruska, Ishaan Gulrajani, and Richard Socher. 2015. Ask me anything: Dynamic memory networks for natural language processing. *CoRR* abs/1506.07285. <http://arxiv.org/abs/1506.07285>.
- Kenton Lee, Tom Kwiatkowski, Ankur Parikh, and Dipanjan Das. 2016. Learning recurrent span representations for extractive question answering. *arXiv preprint arXiv:1611.01436*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries.
- Zhouhan Lin, Minwei Feng, Cicero Nogueira dos Santos, Mo Yu, Bing Xiang, Bowen Zhou, and Yoshua Bengio. 2017. A structured self-attentive sentence embedding. *arXiv preprint arXiv:1703.03130*.
- Rui Liu, Junjie Hu, Wei Wei, Zi Yang, and Eric Nyberg. 2017a. Structural embedding of syntactic trees for machine comprehension. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. pages 815–824.
- Rui Liu, Wei Wei, Weiguang Mao, and Maria Chikina. 2017b. Phase conductor on multi-layered attentions for machine comprehension. *arXiv preprint arXiv:1710.10504*.
- Xiaodong Liu, Kevin Duh, and Jianfeng Gao. 2018. Stochastic answer networks for natural language inference. *arXiv preprint arXiv:1804.07888*.
- Bryan McCann, James Bradbury, Caiming Xiong, and Richard Socher. 2017. Learned in translation: Contextualized word vectors. *arXiv preprint arXiv:1708.00107*.
- Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. Ms marco: A human generated machine reading comprehension dataset. *arXiv preprint arXiv:1611.09268*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*. Association for Computational Linguistics, pages 311–318.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Doha, Qatar, pages 1532–1543. <http://www.aclweb.org/anthology/D14-1162>.
- M. E. Peters, M. Neumann, M. Iyyer, M. Gardner, C. Clark, K. Lee, and L. Zettlemoyer. 2018. Deep contextualized word representations. *ArXiv e-prints*.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text pages 2383–2392. <https://aclweb.org/anthology/D16-1264>.

- Matthew Richardson, Christopher J.C. Burges, and Erin Renshaw. 2013. MCTest: A challenge dataset for the open-domain machine comprehension of text. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Seattle, Washington, USA, pages 193–203.
- S. Salant and J. Berant. 2017. Contextualized Word Representations for Reading Comprehension. *ArXiv e-prints* .
- Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, and Hannaneh Hajishirzi. 2016. Bidirectional attention flow for machine comprehension. *arXiv preprint arXiv:1611.01603* .
- Yelong Shen, Po-Sen Huang, Jianfeng Gao, and Weizhu Chen. 2016. Reasonet: Learning to stop reading in machine comprehension. *arXiv preprint arXiv:1609.05284* .
- Yelong Shen, Xiaodong Liu, Kevin Duh, and Jianfeng Gao. 2017. An empirical analysis of multiple-turn reasoning strategies in reading comprehension tasks. *arXiv preprint arXiv:1711.03230* .
- Alessandro Sordani, Philip Bachman, Adam Trischler, and Yoshua Bengio. 2016. Iterative alternating neural attention for machine reading. *arXiv preprint arXiv:1606.02245* .
- Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* 15(1):1929–1958.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762* .
- Wenhui Wang, Nan Yang, Furu Wei, Baobao Chang, and Ming Zhou. 2017. Gated self-matching networks for reading comprehension and question answering. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. volume 1, pages 189–198.
- Y. Wang, K. Liu, J. Liu, W. He, Y. Lyu, H. Wu, S. Li, and H. Wang. 2018. Multi-Passage Machine Reading Comprehension with Cross-Passage Answer Verification. *ArXiv e-prints* .
- Dirk Weissenborn, Georg Wiese, and Laura Seiffe. 2017. Fastqa: A simple and efficient neural architecture for question answering. *arXiv preprint arXiv:1703.04816* .
- Caiming Xiong, Victor Zhong, and Richard Socher. 2016. Dynamic coattention networks for question answering. *arXiv preprint arXiv:1611.01604* .
- Caiming Xiong, Victor Zhong, and Richard Socher. 2017. Dcn+: Mixed objective and deep residual coattention for question answering. *arXiv preprint arXiv:1711.00106* .
- Adams Wei Yu, David Dohan, Minh-Thang Luong, Rui Zhao, Kai Chen, Mohammad Norouzi, and Quoc V. Le. 2018. Qanet: Combining local convolution with global self-attention for reading comprehension.
- Junbei Zhang, Xiaodan Zhu, Qian Chen, Lirong Dai, and Hui Jiang. 2017. Exploring question understanding and adaptation in neural-network-based question answering. *arXiv preprint arXiv:1703.04617* .