

# Semantics-Aware Inferential Network for Natural Language Understanding

Shuailiang Zhang<sup>1,2,3</sup>, Hai Zhao<sup>1,2,3,\*</sup>, Junru Zhou<sup>1,2,3</sup>

<sup>1</sup>Department of Computer Science and Engineering, Shanghai Jiao Tong University

<sup>2</sup>Key Laboratory of Shanghai Education Commission for Intelligent Interaction and Cognitive Engineering, Shanghai Jiao Tong University, Shanghai, China

<sup>3</sup>MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, Shanghai, China  
zsl123@sjtu.edu.cn, zhaohai@cs.sjtu.edu.cn

## Abstract

For natural language understanding tasks, either machine reading comprehension or natural language inference, both semantics-aware and inference are favorable features of the concerned modeling for better understanding performance. Thus we propose a Semantics-Aware Inferential Network (SAIN) to meet such a motivation. Taking explicit contextualized semantics as a complementary input, the inferential module of SAIN enables a series of reasoning steps over semantic clues through an attention mechanism. By stringing these steps, the inferential network effectively learns to perform iterative reasoning which incorporates both explicit semantics and contextualized representations. In terms of well pre-trained language models as front-end encoder, our model achieves significant improvement on 11 tasks including machine reading comprehension and natural language inference.

## 1 Introduction

Recent studies (Zhang et al., 2020; Mihaylov and Frank, 2019; Sun et al., 2019; Zhang et al., 2019, 2018) have shown that introducing extra common sense knowledge or linguistic knowledge into language representations may further enhance the concerned natural language understanding (NLU) tasks that latently have a need of reasoning ability, such as natural language inference (NLI) (Wang et al., 2019; Bowman et al., 2015) and machine reading comprehension (MRC) (Rajpurkar et al., 2018; Koisk et al., 2018). (Zhang et al., 2020) propose incorporating explicit semantics as a well-formed linguistic knowledge by concatenating the pre-trained language model embedding with semantic role labeling embedding, and obtains significant gains on

Corresponding author. This paper was partially supported by National Key Research and Development Program of China (No. 2017YFB0304100) and Key Projects of National Natural Science Foundation of China (U1836222 and 61733011).

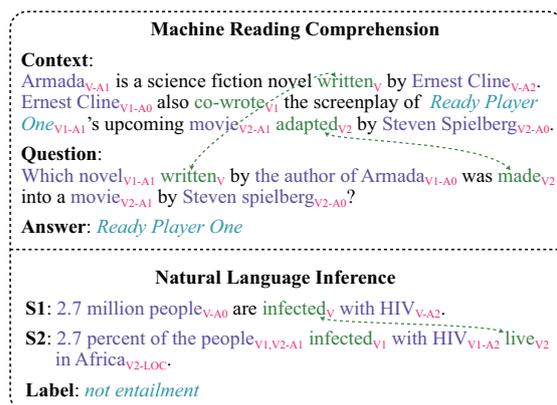


Figure 1: Examples in MRC and NLI with necessary semantic annotations. The connected predicates have important arguments to predict the answer.

the SNLI (Bowman et al., 2015) and GLUE benchmark (Wang et al., 2019). (Mihaylov and Frank, 2019) use semantic information to strengthen the multi-head self-attention model, and achieves substantial improvement on NarrativeQA (Koisk et al., 2018). In this work, we propose a Semantics-Aware Inferential Network (SAIN) to refine the use of semantic structures by decomposing text into different semantic structures for compositional processing in inferential network.

Questions in NLU tasks are usually not compositional, so most existing inferential networks (Weston et al., 2014; Yu et al., 2019) input the same text at each reasoning step, which is not efficient enough to perform iterative reasoning. To overcome this problem, we use semantic role labeling to decompose the text into different semantic structures which are referred as different semantic representations of the sentence (Khashabi et al., 2018; Mihaylov and Frank, 2019).

Semantic role labeling (SRL) over a sentence is to discover *who did what to whom, when and why* with respect to the central meaning (usually verbs) of the sentence and present semantic relationship

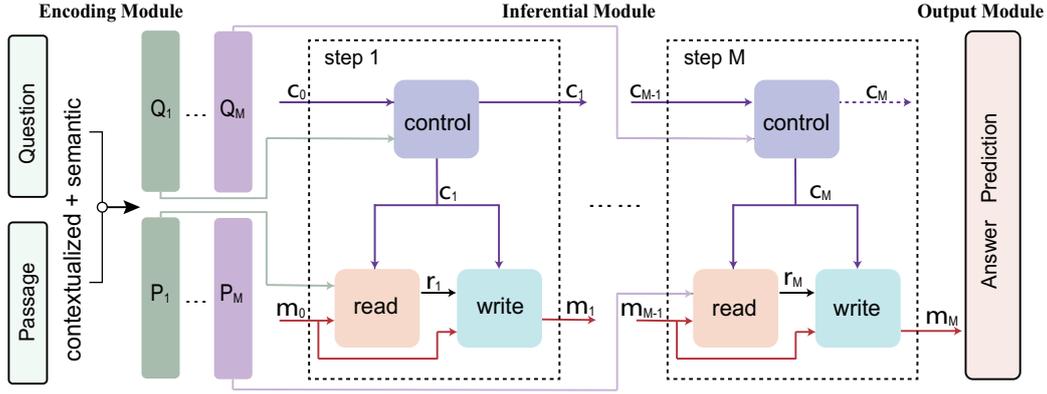


Figure 2: Overview of the framework. Here we only show the inputs and outputs of the first step and last step. The encoding module outputs  $M$  semantic representations that integrate both the contextualized and semantic embedding. The model attends to  $Q_i$  and  $P_i$  in step  $i$ . The final memory state  $m_M$  is used to predict the answer.

as predicate-argument structure, which naturally matches the requirements of MRC and NLI tasks, because questions in MRC are usually formed with *who, what, how, when, why* and verbs in NLI play an important role to determine the answer. Furthermore, when there are several predicate-argument structures in one sentence, there come multiple contextual semantics. Previous neural models are usually with little consideration of modeling these multiple semantic structures which could be critical to predict the answer.

In Figure 1, to correctly answer the MRC question, the model needs to recognize that *the author of Armada is Ernest Cline* firstly, and then knows that *Ernest Cline’s novel Ready Player One was made into a movie by Steven Spielberg*, which requires iteratively reasoning over the two predicates *written* and *made* because they have very similar arguments with the corresponding predicates *written* and *adapted* in the context. For the NLI example, if the model recognizes the predicate *infected* as the central meaning in S2 and ignores the true central word *live*, it probably makes wrong prediction *entailment* because S1 also has a similar structure predicated on *infected*. So it may be helpful to refine the use of semantic clues by integrating all the semantic information into the inference.

We are motivated to model these semantic structures by presenting SAIN, which consists of a set of reasoning steps. In SAIN, each step attends to one predicate-argument structure and can be viewed as a cell consisting of three units: control unit, read unit and write unit, that operate over dull *control* and *memory* hidden states. The cells are recursively connected, where the result of the previous step acts as the context of next step. The interaction between

the cells is regulated by structural constraints to perform iterative reasoning in an end-to-end way.

This work will focus on two typical NLU tasks, natural language inference (SNLI (Bowman et al., 2015), QNLI (Rajpurkar et al., 2016), RTE (Bentivogli et al., 2009) and MNLI (Williams et al., 2018)) and machine reading comprehension (SQuAD (Rajpurkar et al., 2016, 2018) and MRQA (Fisch et al., 2019)). Experiment results indicate that our proposed model achieves significant improvement over the strong baselines on these tasks and obtains the state-of-the-art performance on SNLI and MRQA datasets.

## 2 Approach

The model framework is shown in Figure 2. Our model includes: 1) contextualized encoding module which obtains the joint representation of the pre-trained language model embedding and semantic embedding. 2) inferential module which consists of a set of recurrent reasoning steps/cells, where each step/cell attends to one predicate-argument structure of one sentence. 3) output module which predicts the answer based on the final memory state of the inferential module.

### 2.1 Task Definition

For MRC task, given a passage ( $\mathbf{P}$ ) and a question ( $\mathbf{Q}$ ), the goal is to predict the answer from the given passage. For NLI task, given a pair of sentences, the goal is to judge the relationship between their meanings, such as entailment, neutral and contradiction. Our model will be introduced with the background of MRC task, and the corresponding NLI implementation of our model can be regarded as a simplified case of the MRC, considering that

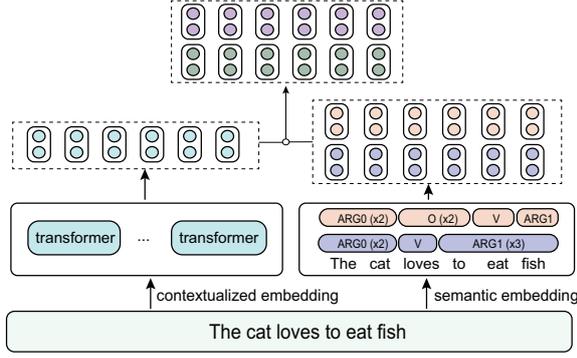


Figure 3: Different semantic representations of one sentence combined by contextualized embedding and semantic embedding.

passage and question in MRC task correspond to two sentences in NLI task.

## 2.2 Semantic Role Labeling

Semantic role labeling (SRL) is generally formulated as multi-step classification subtasks in pipeline systems to identify the semantic structures. There are a few of formal semantic frames, including FrameNet (Baker et al., 1998) and PropBank (Palmer et al., 2005), which generally present the semantic relationship as predicate-argument structure. When several argument-taking predicates are recognized in one sentence, we obtain multiple semantic representations of the sentence. For example, given the context sentence in Figure 3 with target predicates *loves* and *eat*, there are two semantic structures labeled as follows,

$[The\ cat]_{ARG0} [loves]_V [to\ eat\ fish]_{ARG1}$ .

$[The\ cat]_{ARG0} [loves\ to]_O [eat]_V [fish]_{ARG1}$ .

where ARG0, ARG1 represents the argument role 0, 1 of the predicate V, respectively.

## 2.3 Contextual Encoding

**Semantic Embedding** Given the sentence  $X = \{x_1, \dots, x_n\}$  with  $n$  words and  $m$  predicates ( $m = 2$  in Figure 3), there come  $m$  corresponding labeled SRL sequences  $\{L_1, L_2, \dots, L_m\}$  with length  $n$ . Note this is done in data preprocessing and these labels are not updated with the following modules. These semantic role labels are mapped into vectors in dimension  $d_w$  where each sequence  $L_i$  is embedded as  $E^{s_i} = \{e_1^i, \dots, e_n^i\} \in R^{n \times d_w}$ .

**Contextualized Embedding** With an adopted contextualized encoder, the input sequence  $X = \{x_1, \dots, x_n\}$  is embedded as  $E^w = \{e_1, \dots, e_{n_s}\} \in R^{n_s \times d_s}$ , where  $d_s$  is hidden state size of the encoder and  $n_s$  is the tokenized sequence length.

**Joint embedding** Note that the input sequence may be tokenized into subwords. Then the tokenized sequence of length  $n_s$  is usually longer than the SRL sequence of length  $n$ . To align these two sequences, we extend the SRL sequence to length  $n_s$  by assigning the subwords the same label with original word<sup>1</sup>. The aligned contextualized and semantic embeddings are then concatenated as the joint embedding<sup>2</sup> for the sequence  $E^{X_i} = [E^{s_i}; E^w] \in R^{n_s \times d}$ , where  $d = d_s + d_w$ .

Different sentences have various numbers of predicate-argument structures, here we set the maximum number as  $M$  for ease of calculation<sup>3</sup>. So for MRC, the passage and question both have  $M$  encoded representations where  $E^P = \{E^{P_1}, \dots, E^{P_M}\} \in R^{M \times |P| \times d}$  and  $E^Q = \{E^{Q_1}, \dots, E^{Q_M}\} \in R^{M \times |Q| \times d}$ , where  $|P|, |Q|$  are the subwords numbers of passage and question.

## 2.4 Inferential Network

The inferential module performs explicit multi-step reasoning by stringing together  $M$  cells, where each attends to one semantic structure of the sentence. Each cell has three operation units: control unit, read unit and write unit, iteratively aggregating information from different semantic structures.

For MRC, each reasoning step attends to one semantic structure of each sentence from passage and question, respectively. So passage  $E^{P_i} = \{p_{i,1}, \dots, p_{i,|P|}\}$  and question  $E^{Q_i} = \{q_{i,1}, \dots, q_{i,|Q|}\}$  are the input sequences for step  $i$ . Besides, we use biLSTM to get the overall question representation  $bq_i = [\overrightarrow{q_{i,1}}; \overleftarrow{q_{i,|Q|}}] \in R^{2d}$ .

**Reasoning Cell** The reasoning cell is a recurrent cell designed to capture information from different semantic structures. For each step  $i = 1, \dots, M$  in the reasoning process, the  $i^{th}$  cell maintains two hidden states: **control**  $c_i$  and **memory**  $m_i$ , with dimension  $d$ . The control  $c_i$  retrieves information from  $E^{Q_i}$  by calculating a soft attention-based weighted average of the question words. The memory  $m_i$  holds the intermediate results from the reasoning process up to the  $i^{th}$  step by integrating the preceding hidden state  $m_{i-1}$  with the new

<sup>1</sup>For example, if  $x_j$  is tokenized into three subwords  $\{x_{j_1}, x_{j_2}, x_{j_3}\}$ , then  $E^{s_i} = \{e_1^i, \dots, e_j^i, \dots, e_n^i\}$  is extended to  $E^{s_i} = \{e_1^i, \dots, e_{j_1}^i, e_{j_2}^i, e_{j_3}^i, \dots, e_n^i\}$

<sup>2</sup>We also tried summation and multiplication, but experiments show that concatenation is the best.

<sup>3</sup>The sentences without enough number of semantic structures are padded to  $M$  structures where all the labels are assigned to  $O$ . For example, the sentence in Figure 3 is padded as  $[The\ cat\ loves\ to\ eat\ fish]_o$ .

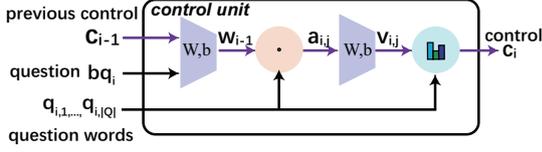


Figure 4: The control unit.

information  $r_i$  retrieved from the passage  $E^{P_i}$ .

There are three units in each cell: control unit, read unit and write unit, which work together to perform iterative reasoning. The control unit retrieves the information from the question, updating the control hidden state  $c_i$ . The read unit extracts relevant information from the passage and outputs extracted information  $r_i$ . The write unit integrates  $c_i$  and  $r_i$  into the memory  $m_{i-1}$ , producing a new memory  $m_i$ . In the following, we give the details of these three units. All the vectors are of dimension  $d$  unless otherwise stated.

The **control unit** (Figure 4) attends to the  $i^{\text{th}}$  semantic structure of the question  $E^{Q_i}$  at step  $i$  and updates the control state  $c_i$  accordingly. Firstly, we combine the overall question representation  $bq_i$  and preceding reasoning operation  $c_{i-1}$  into  $w_i$  through a linear layer. Subsequently, we calculate the similarity between  $w_i$  and each question word  $q_{i,j}$ , and pass the result through a softmax layer, yielding an attention distribution over the question words. Finally, we sum the words over this distribution to get the new control  $c_i$ . The calculation details are as follows:

$$\begin{aligned} w_i &= W^{d \times 2d} [c_{i-1}, bq_i] + b^d \\ a_{i,j} &= W^{1 \times d} (w_i \odot q_{i,j}) + b^1 \\ v_{i,j} &= \text{Softmax}(a_{i,j}), j = 1, \dots, |Q| \\ c_i &= \sum_{j=1}^{|Q|} v_{i,j} \cdot q_{i,j} \end{aligned}$$

where  $W^{d \times 2d}$ ,  $W^{1 \times d}$ ,  $b^d$  and  $b^1$  are learnable parameters,  $|Q|$  is the subwords numbers of question.

The **read unit** (Figure 5) inspects the  $i^{\text{th}}$  semantic structure of the passage  $E^{P_i}$  at step  $i$  and retrieves the information  $r_i$  to update the memory. Firstly, we compute the interaction between every passage word  $p_{i,p}$  and the memory  $m_{i-1}$ , resulting in  $I_{i,p}$  which measures the relevance of the passage word to the preceding memory. Then,  $I_{i,p}$  and  $p_{i,p}$  are concatenated and passed through a linear transformation, yielding  $\hat{I}_{i,p}$  which considers both

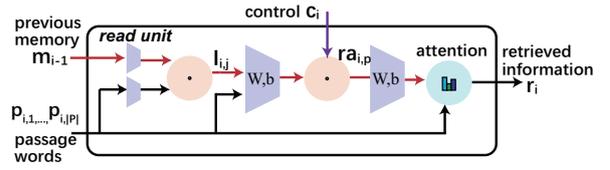


Figure 5: The read unit.

the new information from  $E^{P_i}$  and the information related to the prior intermediate result. Finally, aiming to retrieve the information relevant to the question, we measure the similarity between  $\hat{I}_{i,p}$  and  $c_i$  and pass the result through a softmax layer which produces an attention distribution over the passage words. This distribution is used to get the weighted average  $r_i$  over the passage. The calculation is detailed as follows:

$$\begin{aligned} I_{i,p} &= [W_1^{d \times d} m_{i-1} + b_1^d] \odot [W_2^{d \times d} p_{i,p} + b_2^d] \\ \hat{I}_{i,p} &= W^{d \times 2d} [I_{i,p}, p_{i,p}] + b^d \\ r_{i,p} &= W^{d \times d} (c_i \odot \hat{I}_{i,p}) + b^d \\ r_{v_{i,p}} &= \text{Softmax}(r_{i,p}), p = 1, \dots, |P| \\ r_i &= \sum_{p=1}^{|P|} r_{v_{i,p}} \cdot p_{i,p} \end{aligned}$$

where all the  $W$  and  $b$  are learnable parameters,  $|P|$  is the subwords numbers of passage.

The **write unit** (Figure 6) is responsible for integrating the information retrieved from the read unit  $r_i$  with the preceding memory  $m_{i-1}$ , guided by the  $i^{\text{th}}$  reasoning operation  $c_i$  from the question. Specifically, a sigmoid gate is used when combining the previous memory state  $m_{i-1}$  and the new memory candidate  $m_i^r$ . The calculation details are as follows:

$$\begin{aligned} m_i^r &= W^{d \times 2d} [r_i, m_{i-1}] + b^d \\ \hat{c}_i &= W^{1 \times d} c_i + b^1 \\ m_i &= \sigma(\hat{c}_i) m_{i-1} + (1 - \sigma(\hat{c}_i)) m_i^r \end{aligned}$$

## 2.5 Output Module

For MRC, the output module predicts the final answer to the question based on the set of memory states  $\{m_1, \dots, m_M\}$  produced by the inferential

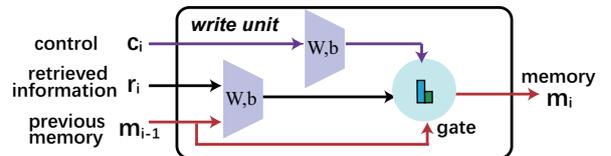


Figure 6: The write unit.

module. For MRC, we calculate the similarity between the  $i^{\text{th}}$  memory  $m_i \in R^d$  and each passage word  $p_{i,p}$  in  $i^{\text{th}}$  semantic passage representation  $E^{P_i}$ , resulting in  $\hat{E}^{P_i}$ ,  $i = 1, \dots, M$ . We concatenate  $\hat{E}^{P_1}, \dots, \hat{E}^{P_M}$  as the final passage representation  $\hat{E}^P \in R^{|P| \times Md}$  which is then passed to a linear layer to get the start and end probability distribution  $p_s, p_e$  on each position. Finally, a cross entropy loss is computed:

$$\begin{aligned} mp_{i,p} &= m_i \cdot p_{i,p} \\ \hat{E}^{P_i} &= [mp_{i,1}, \dots, mp_{i,|P|}] \in R^{|P| \times d} \\ E &= [\hat{E}^{P_1}, \dots, \hat{E}^{P_M}] \in R^{|P| \times Md} \\ [p_s, p_e] &= EW^{Md \times 2} \in R^{|P| \times 2} \\ \text{Loss} &= \frac{1}{2}CE(p_s, y_s) + \frac{1}{2}CE(p_e, y_e) \end{aligned}$$

where  $y_s$  and  $y_e$  are the true start and end probability distribution.  $p_s, p_e, y_s$  and  $y_e$  are all with size  $R^{|P|}$ .  $CE(\cdot)$  indicates the cross entropy function.

For NLI, the final memory state  $m_M$  is directly passed to a linear layer to produce the probability distribution over the labels:  $p = m_M W^{d \times N} \in R^N$ . Cross entropy is used as the metric:  $\text{Loss} = CE(p, y)$ , where  $N$  is the number of labels.  $p \in R^N$  is the predicted probability distribution over the labels and  $y \in R^N$  is the true label distribution.

### 3 Experiments

#### 3.1 Data and Task Description

**Machine Reading Comprehension** We evaluate our model on extractive MRC such as SQuAD (Rajpurkar et al., 2018) and MRQA<sup>4</sup> (Fisch et al., 2019) where the answer is a span of the passage. MRQA is a collection of existing question-answering related MRC datasets, such as SearchQA (Dunn et al., 2017), NewsQA (Trischler et al., 2017), NaturalQuestions (Kwiatkowski et al., 2019), TriviaQA (Joshi et al., 2017), etc. All these datasets as shown in Table 1 are transformed into SQuAD style where given the passage and question, the answer is a span of the passage.

**Natural Language Inference** Given a pair of sentences, the target of natural language inference is to judge the relationship between their meanings, such as entailment, neutral and contradiction. We evaluate on 4 diverse datasets, including Stanford Natural Language Inference (SNLI) (Bowman et al., 2015), Multi-Genre Natural Language Inference (MNLI) (Williams et al., 2018), Question

Dataset	#train	#dev	P	Q
NewsQA	74,160	4,211	599	8
TriviaQA	61,688	7,785	784	16
SearchQA	117,384	16,980	749	17
HotpotQA	72,928	5,904	232	22
NaturalQA	104,071	12,836	153	9

Table 1: Statistics of MRQA datasets. #train and #dev are the number of examples in train and dev set.  $|\cdot|$  denotes the average length in tokens.

Natural Language Inference (QNLI) (Rajpurkar et al., 2016) and Recognizing Textual Entailment (RTE) (Bentivogli et al., 2009).

#### 3.2 Implementation Details

To obtain the semantic role labels, we use the SRL system of (He et al., 2017) as implemented in AllenNLP (Gardner et al., 2018) that splits sentences into tokens and predicts SRL tags such as *ARG0*, *ARG1* for each verb. We use *O* for non-argument words and *V* for predicates. The dimension of SRL embedding is set to 30 and performance does not change significantly when setting this number to 10, 50 or 100. The maximum number of predicate-argument structures (reasoning steps)  $M$  is set to 3 or 4 for different tasks.

Our model framework is based on the Pytorch implementation of transformers<sup>5</sup>. We use Adam as our optimizer with initial learning rate 1e-5 and warm-up rate of 0.1. The batch size is selected in {8, 16, 32} with respect to the task. The total parameters vary from 355M (total steps  $M = 1$ ) to 362M ( $M = 7$ ), increasing 20M to 27M parameters compared to BERT (335M).

#### 3.3 Overall Results

Our main comparison models are the BERT baselines (BERT (Devlin et al., 2019) and SpanBERT (Joshi et al., 2019)) and SemBERT (Zhang et al., 2020). SemBERT improves the language representation by concatenating the BERT embedding and semantic embedding, where embeddings from different predicate-argument structures are simply fused as one semantic representation by using one linear layer. We compare our model to these baselines on 11 benchmarks including 5 MRQA datasets, 4 NLI tasks and 2 SQuAD datasets in Tables 2 and 3.

<sup>4</sup><https://github.com/mrqa/MRQA-Shared-Task-2019>.

<sup>5</sup><https://github.com/huggingface/transformers>.

	NewsQA	TriviaQA	SearchQA	HotpotQA	NaturalQA	(Avg.)
MTL <sub>base</sub> (Fisch et al., 2019)	66.8	71.6	76.7	76.6	77.4	73.8
MTL <sub>large</sub> (Fisch et al., 2019)	66.3	74.7	79.0	79.0	79.8	75.8
CLER (Takahashi et al., 2019)	69.4	75.6	79.0	79.8	79.8	76.7
BERT <sub>large</sub> (Joshi et al., 2019)	68.8	77.5	81.7	78.3	79.9	77.3
HLTC (Su et al., 2019)	72.4	76.2	79.3	80.1	80.6	77.7
SemBERT* (Zhang et al., 2020)	69.1	78.6	82.4	78.6	80.3	77.8
SpanBERT (Joshi et al., 2019)	73.6	83.6	84.8	83.0	82.5	81.5
BERT* <sub>base</sub>	66.2	71.5	77.0	75.0	77.5	73.4
BERT* <sub>large</sub>	69.2	77.4	81.5	78.2	79.4	77.2
SpanBERT*	73.0	83.1	83.5	82.5	81.9	80.8
Our Models						
SAIN <sub>BERT<sub>base</sub></sub>	67.9	72.3	77.8	77.4	78.6	74.8
SAIN <sub>BERT<sub>large</sub></sub>	72.1	80.1	83.4	79.4	82.0	79.4
SAIN <sub>SpanBERT</sub>	74.2	84.5	84.4	83.4	82.7	<b>81.9</b>

Table 2: Performance (F1) on five MRQA tasks. Results with \* are our implementations. Avg indicates the average score of these datasets. All these results are from single models.

Model	MNLI-m/mm		QNLI	RTE	SNLI	(Avg.)	SQuAD 1.1		SQuAD 2.0	
	Acc	Acc	Acc	Acc	Acc		EM	F1	EM	F1
BERT <sub>base</sub>	84.6	83.4	89.3	66.4	90.7	82.9	80.8	88.5	77.1*	80.3*
BERT <sub>large</sub>	86.7	85.9	92.7	70.1	91.1	85.3	84.1	90.9	80.0	83.3
SemBERT <sub>base</sub>	84.4	84.0	90.9	69.3†	91.0*	83.9	-	-	-	-
SemBERT <sub>large</sub>	87.6	86.3	<b>94.6</b>	70.9†	91.6	86.2	84.5*	91.3*	80.9	83.6
Our Models										
SAIN <sub>BERT<sub>base</sub></sub>	84.9	85.0	92.1	72.0	91.2	85.1	82.2	89.3	79.4	82.0
SAIN <sub>BERT<sub>large</sub></sub>	<b>87.7</b>	<b>87.3</b>	94.5	<b>73.9</b>	<b>91.7</b>	<b>87.1</b>	<b>85.4</b>	<b>91.9</b>	<b>82.8</b>	<b>85.4</b>

Table 3: Experiment results on MNLI, QNLI, RTE, SNLI and SQuAD 1.1, SQuAD 2.0. The results of BERT and SemBERT are from (Devlin et al., 2019) and (Zhang et al., 2020). † indicates the results of SemBERT without random restarts and distillation. Results of SQuAD are tested on development sets. Results with \* are our implementations. Avg indicates the average score of these datasets. All these results are from single models.

**SAIN vs. BERT/SpanBERT baselines** Compared to BERT (Devlin et al., 2019), our model achieves general improvements, including 2.1% (79.4 vs. 77.3), 1.8% (87.1 vs. 85.3), 1.6% (88.7% vs. 87.1 %) average improvement on 5 MRQA, 4 NLI and 2 SQuAD datasets. Our model also outperforms other BERT based models CLER (Takahashi et al., 2019) and HLTC (Su et al., 2019) on MRQA. We also compare with SpanBERT (Joshi et al., 2019) on MRQA datasets and our model outperforms this baseline by 0.4% (81.9 vs. 81.5) in average F1 score. To the best of our knowledge, we achieve state-of-the-art performance on MRQA (dev sets) and SNLI.

**SAIN vs. SemBERT** Our SAIN outperforms SemBERT on all tasks, including 1.6% (79.4 vs. 77.8), 0.9% (87.1 vs. 86.2) and 1.3% (86.4 vs.

	RTE	SQuAD 1.1	SQuAD 2.0
SAIN	<b>73.4</b>	<b>91.9</b>	<b>85.4</b>
w/o IM	71.2 (-2.3)	90.6 (-1.3)	82.5 (-1.9)
w/o SI	71.5 (-1.9)	90.1 (-1.8)	83.1 (-2.3)
w/o IR	72.0 (-1.4)	90.9 (-1.0)	83.2 (-2.2)

Table 4: Ablation study on RTE, SQuAD 1.1 and SQuAD 2.0 (F1). We use BERT<sub>large</sub> as contextual encoder here. The definition of IM, SI and IR is detailed in Section 3.4.

85.1) average improvement on MRQA, NLI and SQuAD datasets. We attribute the superiority of our SAIN to its more refined use of semantic clues in terms of inferential network rather than SemBERT which simply encodes all predicate-argument structures into one embedding.

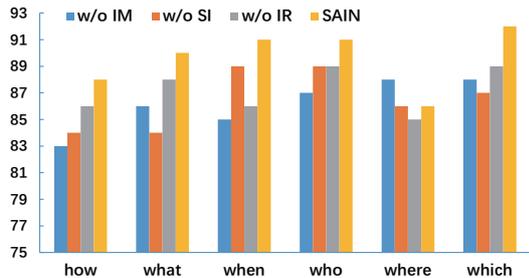


Figure 7: Performance on different question types, tested on the SQuAD 1.1 development set. BERT<sub>base</sub> is used as contextual encoder here. The definition of IM, SI and IR is detailed in Section 3.4.

### 3.4 Ablation Study

To evaluate the contribution of key components in our model, we perform ablation studies on the RTE and SQuAD 2.0 dev sets as shown in Table 4. Here we focus on these components: (1) the whole inferential module (IM); (2) the semantic information (SI); (3) iterative reasoning (IR) that different reasoning cells attend to different predicate-argument structures. To evaluate their contribution, we perform experiments respectively by: (1) IM: removing the inferential module and simply combining the BERT embedding with semantic embeddings from different predicate-argument structures; (2) SI: removing all the semantic embeddings; (3) IR: combining all semantic embeddings from different predicate-argument structures as one and every reasoning step taking the same semantic embedding.

As displayed in Table 4, the ablation on all evaluated components results in performance drop which indicates that all the key components (the inferential module, the semantic information and iterative reasoning process) are indispensable for the model. Particularly, the ablation on iterative reasoning proves that it is necessarily helpful that the model attends to different predicate-argument structures in different reasoning steps.

Furthermore, Figure 7 shows the ablation results on different question types, tested on sampled examples from SQuAD 1.1. The full SAIN model outperforms all other ablation models on all question types except the *where* type questions, which again proves that integrating the semantic information of (*who did what to whom, when and why*) contributes to boosting the performance of MRC tasks where questions are usually formed with *who, what, how, when and why*.

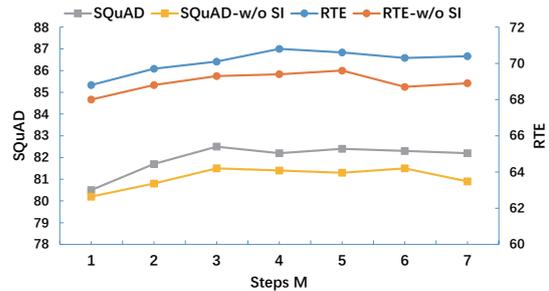


Figure 8: Results on the dev sets of SQuAD 2.0 and RTE when selecting different reasoning steps  $M$ . We use BERT<sub>base</sub> as contextual encoder here. SQuAD/RTE-w/o SI indicates the results without using any semantic information.

### 3.5 Influence of Semantic Information

To further investigate the influence of semantic information, Figure 8 shows the performance comparison of whether to use the semantic information with different numbers of reasoning steps  $M$  (from 1 to 7). The highest performance is achieved when  $M$  is set to 3 on SQuAD, 4 on RTE. The results indicate that semantic information consistently contributes to the performance increase, although the inferential network is strong enough.

To investigate influence of the accuracy of the labeler, we randomly tune specific proportion [0, 20%, 40%] of labels into random error ones. The scores of SQuAD 2.0 and RTE are respectively [85.4, 83.2, 82.6] and [73.4, 71.8, 71.2], which indicate that the model benefits from high-accuracy labeler but can still maintain the performance even using some noisy labels.

### 3.6 Influence of Inferential Mechanism

To obtain better insight into the underlying reasoning processes, we study the visualization of the attention distributions during the iterative computation, and provide examples in Table 5 and Figure 9. Table 5 shows a relatively complex question that is correctly answered by our model, but wrongly predicted by SemBERT (Zhang et al., 2020). In this example, there is misleading contextual similarity between words “*store and transmit*” in sentence S1 and “*transport and storage*” in the question which may lead the model to wrong answer in S1, such as “*fuel*” by SemBERT. To overcome this misleading, the model needs to recognize the central connection predicates “*demand*” and “*requires*” between the question and passage, then extract the correct answer “*special training*” in S2.

**Passage:** (S1) *Steel pipes and storage vessels used to store and transmit both gaseous and liquid oxygen will act as a fuel*; (S2) *and therefore the design and manufacture of oxygen systems requires special training to ensure that ignition sources are minimized.*

**Question:** *What does the transport and storage demand for safety in dealing with oxygen?*

**Golden Answer:** *special training*

**SemBERT:** *fuel*    **SAIN:** *special training*

Table 5: One example that is correctly predicted by SAIN, but wrongly predicted by SemBERT.

Figure 9 shows how our model retrieves information from different semantic structures of the question in each reasoning step. The model first focuses on the word “*what*”, working to retrieve a noun. Then it focuses on the arguments “*transport*” and “*storage*” in step 2 but gets around these words in step 3, and attends to the second verb phrase “*dealing with oxygen*”, taking the model’s attention away from sentence S1. Finally, the model focuses on the main meaning of the question: “*demand for security*” and predicts the correct answer “*special training*” in sentence S2, with respect to the semantic similarity between words “*demand for safety*” and “*requires to ensure*”. This example intuitively explains why our model benefits from the iterative reasoning where each step only attends to one semantic representation.

## 4 Related Work

**Semantic Information for MRC** Using semantic information to enhance the question answering system is one effective method to boost the performance. (Narayanan and Harabagiu, 2004) first stress the importance of semantic roles in dealing with complex questions. (Shen and Lapata, 2007) introduce a general framework for answer extraction which exploits semantic role annotations in the FrameNet (Baker et al., 1998) paradigm. (Yih et al., 2013) propose to solve the answer selection problem using enhanced lexical semantic models. More recently, (Zhang et al., 2020) propose to strengthen the language model representation by fusing explicit contextualized semantics. (Mihaylov and Frank, 2019) apply linguistic annotations to a discourse-aware semantic self-attention encoder which is employed for reading comprehension on narrative texts. In this work, we propose to

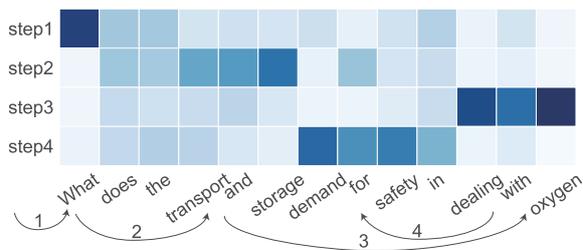


Figure 9: Transformation of attention distribution at each reasoning step, showing how the model iteratively retrieves information from the question.

use inferential model to recurrently retrieve different predicate-argument structures, which presents a more refined way using semantic clues and thus is essentially different from all previous methods.

**Inferential Network** To support inference in neural network, existing models either rely on structured rule-based matching methods (Sun et al., 2018) or multi-layer memory networks (Weston et al., 2014; Liu and Perez, 2017), which either lack end-to-end design or no prior structure to subtly guide the reasoning direction.

Another related works are on Visual QA, aiming to answer the question with regards to the given image. In particular, (Santoro et al., 2017) propose a relation net but restricted the relational question such as comparison. Later, for compositional question, (Hudson and Manning, 2018) introduce an iterative network that separates memory and control to improve interpretability. Our work leverages such separate design, dedicating to inferential NLI and MRC tasks, where the questions are usually not compositional.

To overcome the difficulty of applying inferential network into general NLU tasks, and passingly refine the use of multiple semantic structures, we propose SAIN which naturally decomposes text into different semantic structures for compositional processing in inferential network.

## 5 Conclusion

This work focuses on two typical NLU tasks, machine reading comprehension and natural language inference by refining the use of semantic clues and inferential model. The proposed semantics-aware inferential network (SAIN) is capable of taking multiple semantic structures as input of an inferential network by closely integrating semantics and reasoning steps in a creative way. Experiment results on 11 benchmarks, including 4 NLI tasks and 7 MRC tasks, show that our model outperforms

all previous strong baselines, which consistently indicate the general effectiveness of our model<sup>6</sup>.

## References

- Collin F. Baker, Charles J. Fillmore, and John B. Lowe. 1998. [The Berkeley FrameNet project](#). In *36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics (ACL)*, pages 86–90.
- Luisa Bentivogli, Ido Dagan, Hoa Trang Dang, Danilo Giampiccolo, and Bernardo Magnini. 2009. The Fifth PASCAL Recognizing Textual Entailment Challenge. In *In Proc Text Analysis Conference (TAC09)*.
- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 632–642.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)*, pages 4171–4186.
- Matthew Dunn, Levent Sagun, Mike Higgins, V. Ugur Güney, Volkan Cirik, and Kyunghyun Cho. 2017. [SearchQA: A new q&a dataset augmented with context from a search engine](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Adam Fisch, Alon Talmor, Robin Jia, Minjoon Seo, Eunsol Choi, and Danqi Chen. 2019. [MRQA 2019 shared task: Evaluating generalization in reading comprehension](#). In *Proceedings of 2nd Machine Reading for Reading Comprehension (MRQA) Workshop at EMNLP*, pages 1–13.
- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke Zettlemoyer. 2018. [AllenNLP: A deep semantic natural language processing platform](#). In *Proceedings of Workshop for NLP Open Source Software (NLP-OSS)*, pages 1–6.
- Luheng He, Kenton Lee, Mike Lewis, and Luke Zettlemoyer. 2017. [Deep semantic role labeling: What works and what’s next](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 473–483.
- Drew A. Hudson and Christopher D. Manning. 2018. Compositional attention networks for machine reasoning. In *Proceedings of the 7th International Conference on Learning Representations (ICLR)*.
- Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2019. [Spanbert: Improving pre-training by representing and predicting spans](#). In *CoRR*, volume abs/1907.10529.
- Mandar Joshi, Eunsol Choi, Daniel S. Weld, and Luke Zettlemoyer. 2017. [TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Daniel Khashabi, Tushar Khot, Ashish Sabharwal, and Dan Roth. 2018. [Question Answering as Global Reasoning over Semantic Abstractions](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*.
- Tom Koisk, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gbor Melis, and Edward Grefenstette. 2018. [The NarrativeQA Reading Comprehension Challenge](#). In *Transactions of the Association for Computational Linguistics (TACL)*, volume 6, page 317328.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural questions: a benchmark for question answering research](#). In *Transactions of the Association of Computational Linguistics (TACL)*.
- Fei Liu and Julien Perez. 2017. [Gated End-to-End Memory Networks](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 1–10.
- Todor Mihaylov and Anette Frank. 2019. [Discourse-aware semantic self-attention for narrative reading comprehension](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing, (EMNLP-IJCNLP)*, pages 2541–2552.
- Srini Narayanan and Sanda Harabagiu. 2004. [Question answering based on semantic structures](#). In *Proceedings of the 13th International Conference on Computational Linguistics (COLING)*, pages 693–701.
- Martha Palmer, Daniel Gildea, and Paul Kingsbury. 2005. [The proposition bank: An annotated corpus of semantic roles](#). In *Computational Linguistics (CL)*, pages 71–106.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. [Know what you don’t know: Unanswerable questions for SQuAD](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 784–789.

<sup>6</sup>Our model can be easily adapted to other language models such as ALBERT, which is left for future work.

- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2383–2392.
- Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Timothy Lillicrap. 2017. A simple neural network module for relational reasoning. In *Advances in Neural Information Processing Systems 30*, pages 4967–4976.
- Dan Shen and Mirella Lapata. 2007. [Using semantic roles to improve question answering](#). In *Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing (EMNLP-CoNLL)*, pages 12–21.
- Dan Su, Yan Xu, Genta Indra Winata, Peng Xu, Hyeondey Kim, Zihan Liu, and Pascale Fung. 2019. [Generalizing question answering system with pre-trained language model fine-tuning](#). In *Proceedings of the 2nd Workshop on Machine Reading for Question Answering (EMNLP)*, pages 203–211.
- Yawei Sun, Gong Cheng, and Yuzhong Qu. 2018. [Reading Comprehension with Graph-based Temporal-Casual Reasoning](#). In *Proceedings of the 27th International Conference on Computational Linguistics (COLING)*, pages 806–817.
- Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, and Hua Wu. 2019. [ERNIE: enhanced representation through knowledge integration](#). In *CoRR*, volume abs/1904.09223.
- Takumi Takahashi, Motoki Taniguchi, Tomoki Taniguchi, and Tomoko Ohkuma. 2019. [CLER: Cross-task learning with expert representation to generalize reading and understanding](#). In *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pages 183–190.
- Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordani, Philip Bachman, and Kaheer Suleman. 2017. NewsQA: A Machine Comprehension Dataset. In *Proceedings of the 2nd Workshop on Representation Learning for NLP*, pages 191–200.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *The International Conference on Learning Representations (ICLR)*.
- Jason Weston, Sumit Chopra, and Antoine Bordes. 2014. Memory Networks. In *Proceedings of the 3rd International Conference on Learning Representations (ICLR)*.
- Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. [A broad-coverage challenge corpus for sentence understanding through inference](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)*, pages 1112–1122.
- Wen-tau Yih, Ming-Wei Chang, Christopher Meek, and Andrzej Pastusiak. 2013. [Question Answering Using Enhanced Lexical Semantic Models](#). In *Proceedings of the 44th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1744–1753.
- Jianxing Yu, Zhengjun Zha, and Jian Yin. 2019. [Inferential machine comprehension: Answering questions by recursively deducing the evidence chain from text](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 2241–2251.
- Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. [ERNIE: Enhanced language representation with informative entities](#). In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1441–1451.
- Zhuosheng Zhang, Yuwei Wu, Zuchao Li, and Hai Zhao. 2018. [Explicit Contextual Semantics for Text Comprehension](#). In *Proceedings of the 33rd Pacific Asia Conference on Language, Information and Computation (PACLIC 33)*.
- Zhuosheng Zhang, Yuwei Wu, Hai Zhao, Zuchao Li, Shuailiang Zhang, Xi Zhou, and Xiang Zhou. 2020. [Semantics-aware BERT for Language Understanding](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)*.