



## Introduction

We produce a System to answer questions correctly given paragraph context from SQuAD 2.0. The result would be the span of text or N/A if there is no answer in the paragraph. We use BiDAF as baseline, BERT-based architecture as the core, L1 regularization and other architecture changes on BERT. Ensembling method is also applied for improvement, which combines multiple models into a more robust Question Answering system.

## Data

We use SQuAD 2.0 as the reading comprehension data set. Every answerable SQuAD question has three answers provided.

Dataset has been split into:

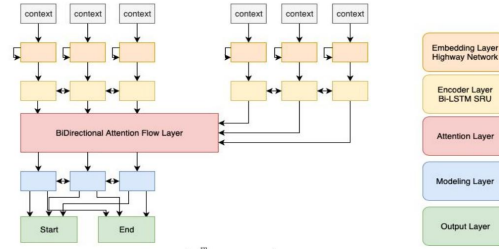
- Train Set: 129941 examples
- Dev Set: 6078 examples
- Test Set: 5291 examples

## Experiment

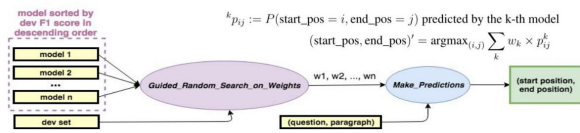
ID	Experiment Name	BERT Model	Pretrained model	Number of Epochs	Learning Rate	Batch Size	Max Sequence Length	Max Sentence Length	Dev F1	Dev EM
1	ori_msmarcoqa1212_330-5_gpt01	BERT-Base-Cased		6	3e-5	12	245	$\lambda = 1e^{-1}$	77.206	73.078
2	ori_msmarcoqa1212_330-5_gpt02	BERT-Base-Cased		4	3e-5	12	245		77.186	73.661
3	ori_msmarcoqa1212_330-5_gpt01-l1e-2	BERT-Base-Cased		4	3e-5	24	140	$\lambda = 1e^{-2}$	76.76	73.955
4	ori_msmarcoqa1212_330-5_gpt01-l1e-3	BERT-Base-Cased		4	3e-5	24	140	$\lambda = 1e^{-3}$	76.666	73.824
5	ori_msmarcoqa1212_330-5_gpt01	BERT-Base-Cased		6	3e-5	12	245		75.925	72.343
6	ori_msmarcoqa1212_330-5_gpt01-l1e-4_ensemble	BERT-Base-Cased		4	3e-5	24	140	$\lambda = 1e^{-4}$	75.889	73.682
7	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		4	3e-5	24	140	$\lambda = 1e^{-4}$	75.705	73.001
8	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		4	3e-5	12	245	$\lambda = 1e^{-4}$	75.871	73.685
9	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		5	3e-5	12	245	$\lambda = 1e^{-4}$ , add one layer	75.334	71.685
10	ori_msmarcoqa1212_330-5_gpt01-ensemble	BERT-Base-Cased		6	3e-5	12	245		75.011	71.272
11	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		4	3e-5	24	140		74.679	71.915
12	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		4	3e-5	48	200		74.633	71.272
13	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		4	3e-5	12	245		74.546	71.092
14	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		4	3e-5	12	200		74.536	71.288
15	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		5	1e-5	12	245		73.885	70.229
16	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		4	3e-5	6	425		73.725	70.7
17	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		4	3e-5	12	128		73.638	71.442
18	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		4	3e-5	12	135		73.202	69.984
19	ori_msmarcoqa1212_330-5_gpt01-l1e-4	BERT-Base-Cased		4	3e-5	48	60		72.954	70.811
20	baseline_sru	GloVe		10	0.1	SRL			64.68	
21	baseline_ensemble	GloVe		20	0.1	LSTM		Baseline	61.588	57.199
22	Guided Random Search for Weighted Average								79.944	77.081

## Models

1. Baseline: BiDAF
2. Baseline: SRU + BiDAF



3. L1 Regularization:  $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_1$
4. Go "deeper": Add one more fully-connected layer to the output of BERT
5. Ensembling: Guided Random Search for Weighted Average Ensembling



## Results

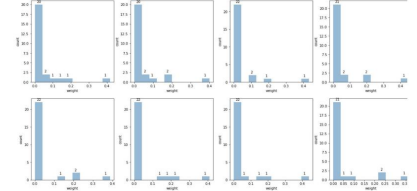
**L1 Regularization:** Train BERT with L1 regularization on weights of output classification and varies the coefficient. Increasing regularization strength helps improve the F1 and EM score.

Regularization Coefficient	Dev F1	Dev EM
0	74.679	71.915
1e-4	75.705	73.001
1e-3	76.666	73.824
1e-2	76.76	73.955

**2. Ensembling:**

Ensembling Model	Dev F1	Dev EM
Guided Random Search for Weighted Average	79.9	77.08
	44	1

## Analysis



This plot visualizes the weights for the top 8 Ensembling models in one run(100 iters) of Guided Random Search of weights by plotting the distribution in histograms. Most of the models only bears a weight in the order of 0.001.

## Conclusions

After training 28 BiDAF-based, BERT-based models, and ensemble them with two algorithms, we push test F1 score to 78.841 and Test EM to 76.1010.

## Future Work

For future work, we would combine the BERT and BiDAF together, which means that we replace BiDAF's GloVe word embedding with BERT last layer's output as as contextual word embedding. Hopefully we can improve our performance more with this idea.

## References

Devlin, Jacob, et al. "Bert: Pre-training of deep bi-directional transformers for language understanding." arXiv preprint arXiv:1810.04805 (2018).