Method	EM	F1
ALBERT	87.0	90.2
Two-model Ensemble	87.6	90.6
Retro-Reader	87.8	90.9

Table 6: Comparisons with Equivalent Parameters on the dev set of SQuAD2.0.

Method	SQuAD2.0		NewsQA	
	EM	F1	EM	F1
BERT	78.8	81.7	51.8	62.5
+ CA	78.8	81.7	52.1	62.7
+ MA	78.3	81.4	52.4	62.6
ALBERT	87.0	90.2	57.1	67.5
+ CA	87.3	90.3	56.0	66.3
+ MA	86.8	90.0	55.8	66.1

Table 7: Results (%) with matching interaction methods on the dev sets of SQuAD2.0 and NewsQA.

exists a tradeoff between the HasAns and NoAns accuracies. We see that the final RV that combines E-FV and I-FV shows the best performance, which we select as our final implementation for testing.

We further conduct the experiments on our model performance of the 5,945 unanswerable questions from the SQuAD 2.0 dev set. Results in Table ?? show that our method improves the performance on unanswerable questions by a large margin, especially in the primary F1 and accuracy metrics.

5.2 Comparisons with Equivalent Parameters

When using sketchy reading module for external verification, we have two parallel modules that have independent parameters. For comparisons with equivalent parameters, we add an ensemble of two baseline models, to see if the advance is purely from the increase of parameters. Table ?? shows the results. We see that our model can still outperform two ensembled models. Although the two modules share the same design of the Transformer encoder, the training objectives (e.g., loss functions) are quite different, one for answer span prediction, the other for answerable decision. The results indicate that our two-stage reading modules would be more effective for learning diverse aspects (verification and span prediction) for solving MRC tasks with different training objectives. From the two modules, we can easily find the effectiveness of either the span prediction or answer verification, to improve the modules correspondingly. We believe this design would be quite useful for real-world applications.

5.3 Evaluation on Matching Interactions

Table ?? shows the results with different interaction methods described in §??. We see that merely adding extra layers could not bring noticeable improvement, which indicates

Passage:

Southern California consists of a heavily developed urban environment, home to some of the largest urban areas in the state, along with vast areas that have been left undeveloped. It is the third most populated megalopolis in the United States, after the Great Lakes Megalopolis and the Northeastern megalopolis. Much of southern California is famous for its large, spread-out, suburban communities and use of automobiles and highways. The dominant areas are Los Angeles, Orange County, San Diego, and Riverside-San Bernardino, each of which are the centers of their respective metropolitan areas...

Question:

What are the second and third most populated megalopolis after Southern California?

Answer:

Gold: (no answer)

ALBERT (+TAV): Great Lakes Megalopolis and the

Northeastern megalopolis.

Retro-Reader over ALBERT: $\langle \text{no answer} \rangle$ $score_{has} = 0.03, score_{na} = 1.73, \delta = -0.98$

Table 8: Answer prediction examples from the ALBERT baseline and Retro-Reader.

that simply adding more layers and parameters would not substantially benefit the model performance. The results verified the PrLMs' strong ability to capture the relationships between passage and question after processing the paired input by deep self-attention layers. In contrast, answer verification could still give consistent and substantial advance.

5.4 Comparison of Predictions

To have an intuitive observation of the predictions of Retro-Reader, we give a prediction example on SQuAD2.0 from baseline and Retro-Reader in Table ??, which shows that our method works better at judging whether the question is answerable on a given passage and gets rid of the plausible answer.

6 Conclusion

As machine reading comprehension tasks with unanswerable questions stress the importance of answer verification in MRC modeling, this paper devotes itself to better verifieroriented MRC task-specific design and implementation for the first time. Inspired by human reading comprehension experience, we proposed a retrospective reader that integrates both sketchy and intensive reading. With the latest PrLM as encoder backbone and baseline, the proposed reader is evaluated on two benchmark MRC challenge datasets SQuAD2.0 and NewsQA, achieving new state-of-the-art results and outperforming strong baseline models in terms of newly introduced statistical significance, which shows the choice of verification mechanisms has a significant impact for MRC performance and verifier is an indispensable reader component even for powerful enough PrLMs used as the encoder. In the future, we will investigate more decoder-side problem-solving techniques to cooperate with the strong encoders for more advanced MRC.

Vitae maxime ullam harum, corporis amet quod