



## Motivation & Contribution

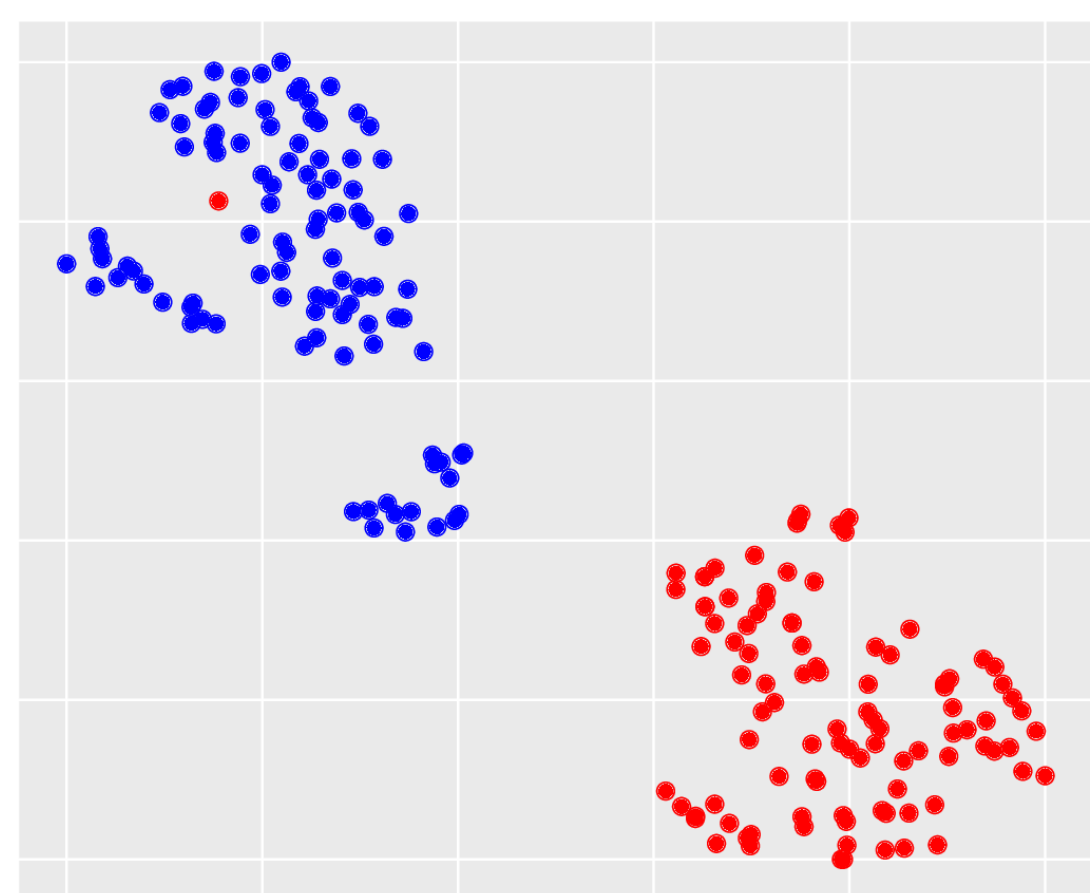
### Motivation

- Recent success in MRC relies on large-scale annotated in-domain data (e.g., SQuAD)
- Directly adapting models from source domain to low-resource target domain performs poorly due to domain shift

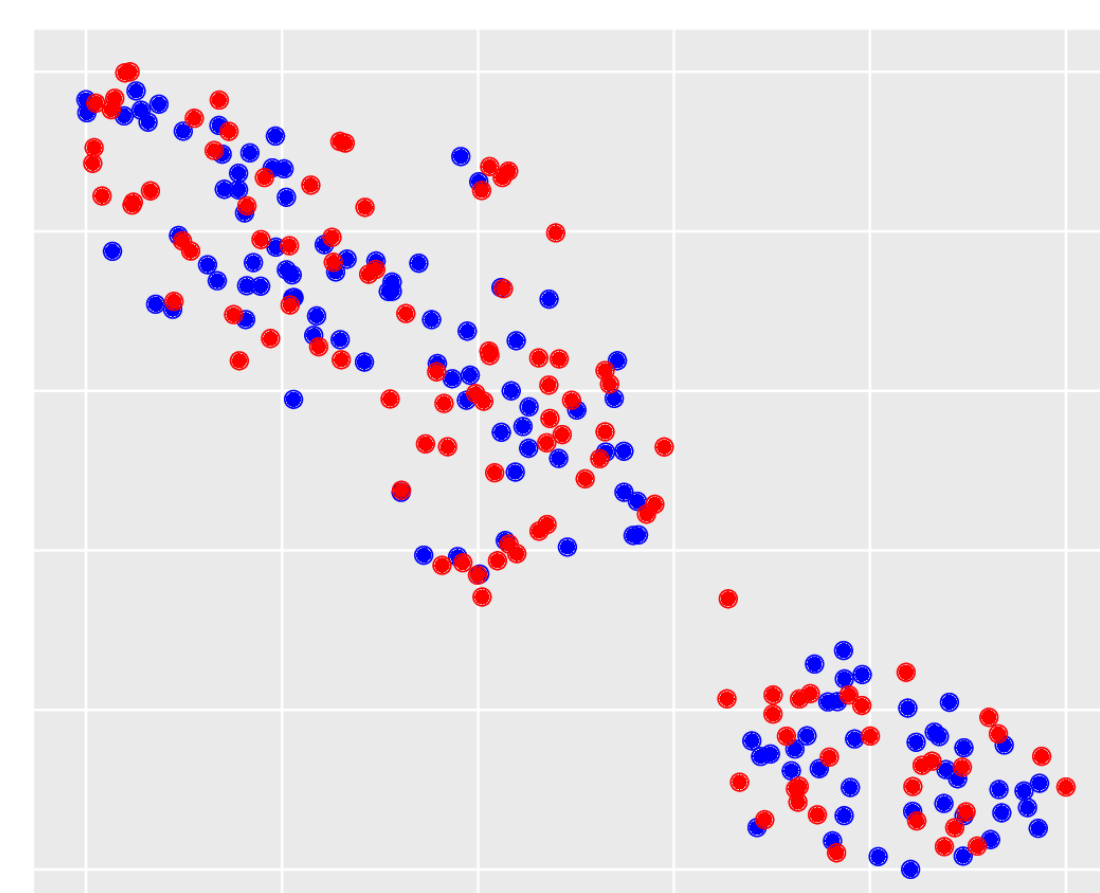
### Contribution

- Unsupervised Domain Adaptation by generating pseudo data on target domain and learning *domain-invariant* representations through adversarial learning

T-SNE plot of encoded feature representations

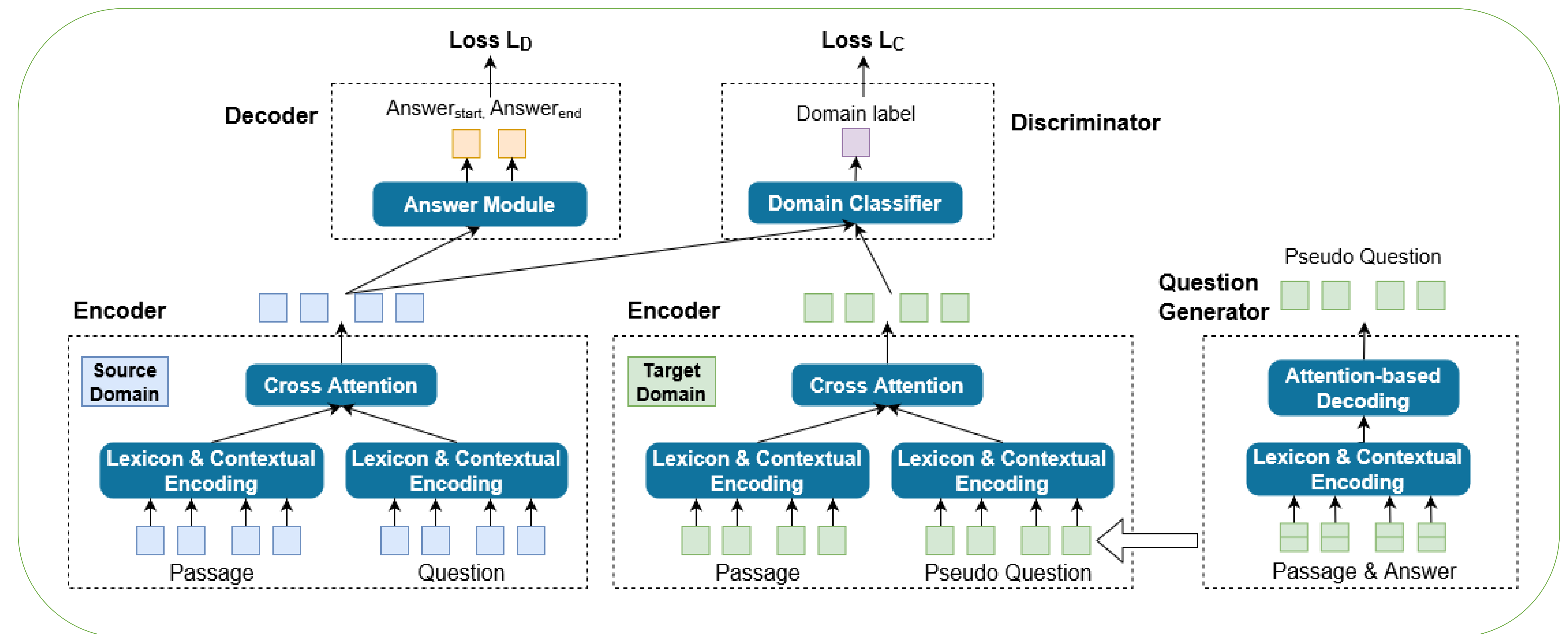


Without domain adaptation



With domain adaptation

## AdaMRC Framework



- Question Generator (QG)**: using passage and answer (extracted by NER) as input for generating pseudo questions in the target domain
- Encoder & Decoder**: source domain and target domain share the same encoder & decoder (i.e., MRC Module)
- Discriminator**: an MLP as domain classifier. A gradient reverse layer is used for gradient backpropagation
- Training**: pseudo target domain data is used for encoder and discriminator training, but *not* for decoder training, because the synthetic data could be noisy

## Training Algorithm

### Algorithm 1 AdaMRC training procedure.

- Input**: source domain labeled data  $S = \{p^s, q^s, a^s\}$ , target domain unlabeled data  $T = \{p^t\}$
- Train the MRC model  $\theta^s = (\theta_e^s, \theta_d^s)$  on source domain  $S$ ;
- Train the QG model  $\theta_{QG}$  on source domain  $S$ ;
- Generate  $T_{gen} = \{p^t, q^t, a^t\}$  using the QG model;
- Initialize  $\theta = (\theta_e, \theta_d, \theta_c)$  with  $\theta^s$ ;
- for** epoch  $\leftarrow 1$  to  $\#epochs$  **do**
- Optimize  $\theta$  on  $S \cup T_{gen}$ . Each minibatch is composed with  $k_s$  samples from  $S$  and  $k_t$  samples from  $T_{gen}$ ;
- end for**
- Output**: Model with the best performance on the target development set  $\theta^*$ .

## Experimental Results

Dataset	Domain
SQuAD (v1.1)	Wiki
NewsQA	News
MS MARCO (v1)	Web

- Main results are based on Stochastic Answer Network (SAN)
- AdaMRC consistently improves performance over baselines
- Direct data augmentation and fine-tuning (SynNet) hurts performance
- Question generation is effective (margin with “AdaMRC with GT questions” is relatively small)
- Generalizable to other datasets and other MRC models with consistent performance gain

Method	EM/F1
<b>SQuAD <math>\rightarrow</math> NewsQA</b>	
SAN	36.68/52.79
SynNet + SAN	35.19/49.61
AdaMRC	<b>38.46/54.20</b>
AdaMRC with GT questions	39.37/54.63
<b>NewsQA <math>\rightarrow</math> SQuAD</b>	
SAN	56.83/68.62
SynNet + SAN	50.34/62.42
AdaMRC	<b>58.20/69.75</b>
AdaMRC with GT questions	58.82/70.14
<b>SQuAD <math>\rightarrow</math> MS MARCO (BLEU-1/ROUGE-L)</b>	
SAN	13.06/25.80
SynNet + SAN	12.52/25.47
AdaMRC	<b>14.09/26.09</b>
AdaMRC with GT questions	15.59/26.40
<b>MS MARCO <math>\rightarrow</math> SQuAD</b>	
SAN	27.06/40.07
SynNet + SAN	23.67/36.79
AdaMRC	<b>27.92/40.69</b>
AdaMRC with GT questions	27.79/41.47

- Compatible with pre-trained language models

Method	EM/F1
SAN	36.68/52.79
AdaMRC + SAN	<b>38.46/54.20</b>
SAN + ELMo	39.61/55.18
AdaMRC + SAN + ELMo	<b>40.96/56.25</b>
BERT <sub>base</sub>	42.00/58.71
AdaMRC + BERT <sub>base</sub>	<b>42.59/59.25</b>

- Can be extended to semi-supervised setting

Ratio (%Labeled data)	SAN	AdaMRC+ SAN
0%	36.68/52.79	<b>38.46/54.20</b>
5%	47.61/62.69	<b>48.50/63.17</b>
10%	48.66/63.32	<b>49.64/63.94</b>
20%	50.75/64.80	<b>51.14/65.38</b>
50%	53.24/67.30	<b>53.34/67.30</b>
100%	<b>56.48/69.14</b>	56.29/68.97