

机器阅读理解的发展及预训练技术的应用

Zhuosheng Zhang

zhangzs@sjtu.edu.cn

<https://bcmi.sjtu.edu.cn/~zhangzs>

Outline

- ❖ Machine Reading Comprehension

 - ❖ Background, Development, Paradigm

- ❖ Methodology

 - ❖ Two-stage Solving Architecture

 - ❖ Traditional Matching Networks

 - ❖ Pre-trained Language Models

- ❖ Frontiers

 - ❖ Techniques

 - ❖ Tasks

 - ❖ Applications

Outline

- ❖ Machine Reading Comprehension

 - ❖ Background, Development, Paradigm

- ❖ Techniques

 - ❖ Two-stage Solving Architecture

 - ❖ Traditional Matching Networks

 - ❖ Pre-trained Language Models

- ❖ Frontiers

 - ❖ Techniques

 - ❖ Tasks

 - ❖ Applications

Introductions to MRC

There are two categories of branches in natural language processing (NLP)

- Core/fundamental NLP
 - Language model/representation
 - Linguistic structure parsing/analysis
 - Morphological analysis/word segmentation
 - Syntactic/semantic/discourse parsing
 - ...
- Application NLP
 - Machine Reading Comprehension (MRC)
 - Text Entailment (TE) or Natural Language Inference (NLI)
 - SNLI, GLUE
 - QA/Dialogue
 - Machine translation
 - ...

Introductions to MRC

- Aim: teach machines to read and comprehend human languages
- Form: find the accurate Answer for a Question according to a given Passage (document).

- Types

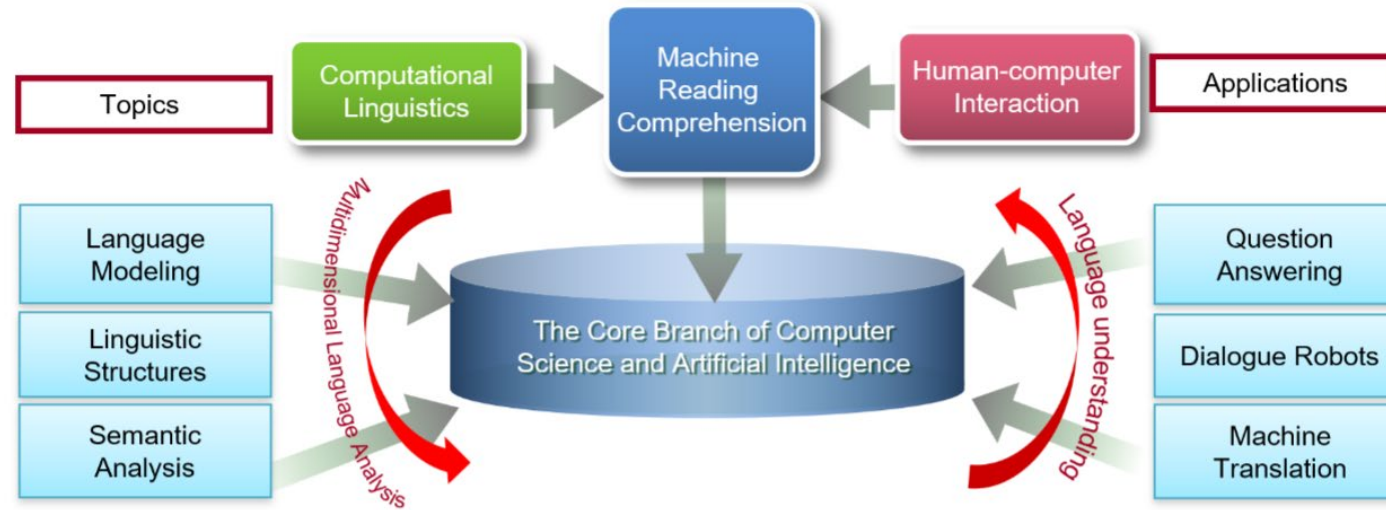
- Cloze-style
- Multi-choice
- Span extraction
- Free-form

- Before 2015

- MCTest
- ProcessBank

- After 2015

- CNN/Daily Mail
- Children Book Test
- WikiReading
- LAMBADA
- SQuAD
- Who did What
- NewsQA
- MS MARCO
- TriviaQA
- CoQA
- QuAC
-



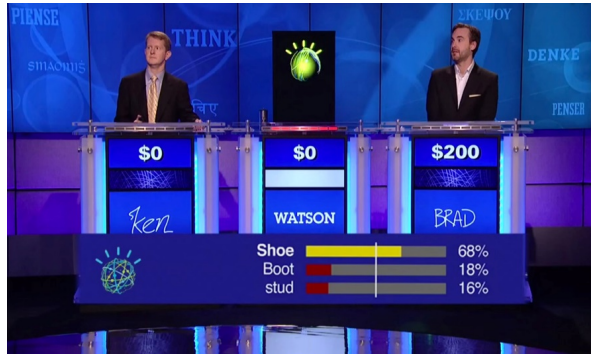
From shared task to **leaderboard**

Introductions to MRC

Cloze-style	from CNN (Hermann et al. 2015)
Context	(@entity0) – a bus carrying members of a @entity5 unit overturned at an @entity7 military base sunday , leaving 23 @entity8 injured , four of them critically , the military said in a news release . a bus overturned sunday in @entity7 , injuring 23 @entity8 , the military said . the passengers , members of @entity13 , @entity14 , @entity15 , had been taking part in a training exercise at @entity19 , an @entity21 post outside @entity22 , @entity7 . they were departing the range at 9:20 a.m. when the accident occurred . the unit is made up of reservists from @entity27 , @entity28 , and @entity29 , @entity7 . the injured were from @entity30 and @entity31 out of @entity29 , a @entity32 suburb . by mid-afternoon , 11 of the injured had been released to their unit from the hospital . pictures of the wreck were provided to the news media by the military . @entity22 is about 175 miles south of @entity32 . e-mail to a friend
Question Answer	bus carrying @entity5 unit overturned at _____ military base @entity7
Multi-choice	from RACE (Lai et al. 2017)
Context	Runners in a relay race pass a stick in one direction. However, merchants passed silk, gold, fruit, and glass along the Silk Road in more than one direction. They earned their living by traveling the famous Silk Road. The Silk Road was not a simple trading network. It passed through thousands of cities and towns. It started from eastern China, across Central Asia and the Middle East, and ended in the Mediterranean Sea. It was used from about 200 B, C, to about A, D, 1300, when sea travel offered new routes, It was sometimes called the world ' s longest highway. However, the Silk Road was made up of many routes, not one smooth path. They passed through what are now 18 countries. The routes crossed mountains and deserts and had many dangers of hot sun, deep snow, and even battles. Only experienced traders could return safely.
Question Answer	The Silk Road became less important because _____. A.it was made up of different routes B.silk trading became less popular C.sea travel provided easier routes D.people needed fewer foreign goods

Span Extraction	from SQuAD (Rajpurkar et al. 2016)
Context	Robotics is an interdisciplinary branch of engineering and science that includes mechanical engineering, electrical engineering, computer science, and others. Robotics deals with the design, construction, operation, and use of robots, as well as computer systems for their control, sensory feedback, and information processing. These technologies are used to develop machines that can substitute for humans. Robots can be used in any situation and for any purpose, but today many are used in dangerous environments (including bomb detection and de-activation), manufacturing processes, or where humans cannot survive. Robots can take on any form, but some are made to resemble humans in appearance. This is said to help in the acceptance of a robot in certain replicative behaviors usually performed by people. Such robots attempt to replicate walking, lifting, speech, cognition, and basically anything a human can do.
Question Answer	What do robots that resemble humans attempt to do? replicate walking, lifting, speech, cognition
Free-form	from DROP (Dua et al. 2019)
Context	The Miami Dolphins came off of a 0-3 start and tried to rebound against the Buffalo Bills. After a scoreless first quarter the Dolphins rallied quick with a 23-yard interception return for a touchdown by rookie Vontae Davis and a 1-yard touchdown run by Ronnie Brown along with a 33-yard field goal by Dan Carpenter making the halftime score 17-3. Miami would continue with a Chad Henne touchdown pass to Brian Hartline and a 1-yard touchdown run by Ricky Williams. Trent Edwards would hit Josh Reed for a 3-yard touchdown but Miami ended the game with a 1-yard touchdown run by Ronnie Brown. The Dolphins won the game 38-10 as the team improved to 1-3. Chad Henne made his first NFL start and threw for 115 yards and a touchdown.
Question Answer	How many more points did the Dolphins score compare to the Bills by the game's end? 28

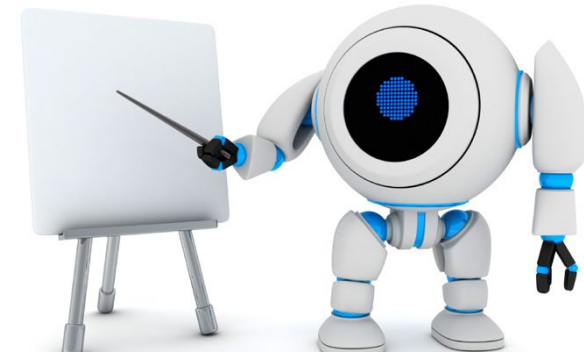
Applications



Question Answering



Dialogue System



Intelligent Teacher



Fake News Identifier



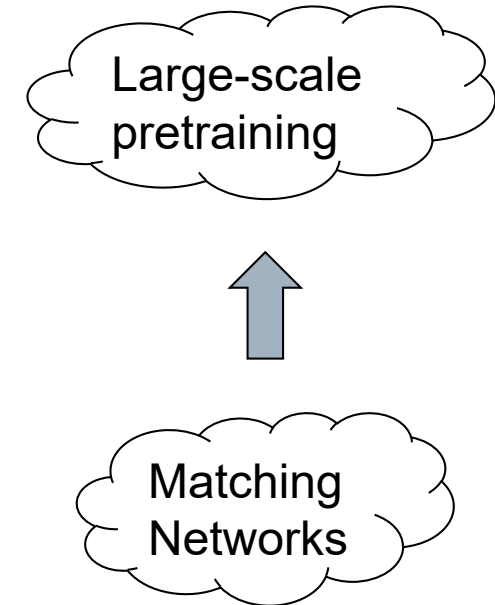
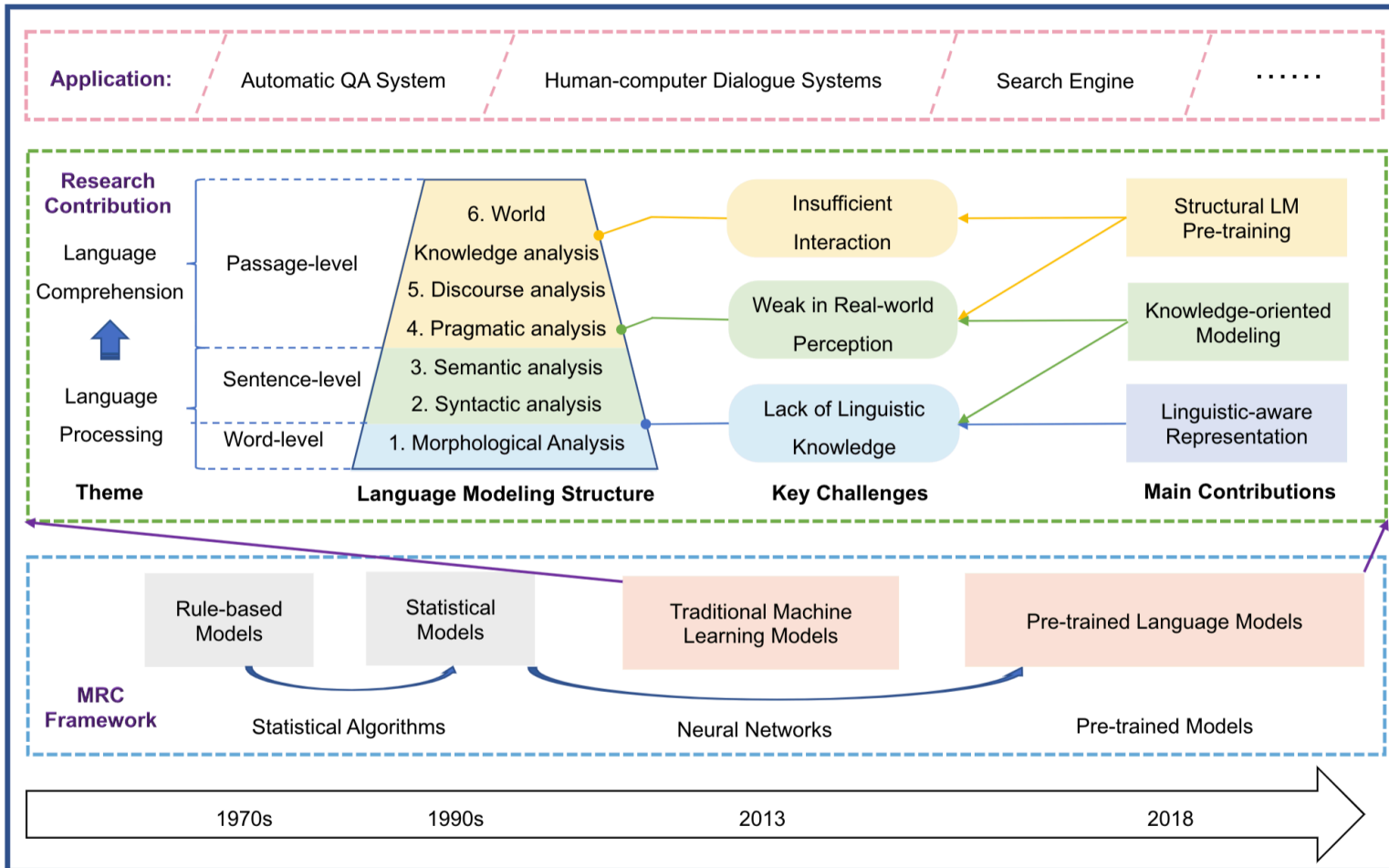
Legal Advisor



Medical Diagnosis

The Boom of MRC researches

- ❑ The burst of deep neural networks, especially attention-based models
- ❑ The evolution of pre-trained language models (large-scale pre-training and task-specific)

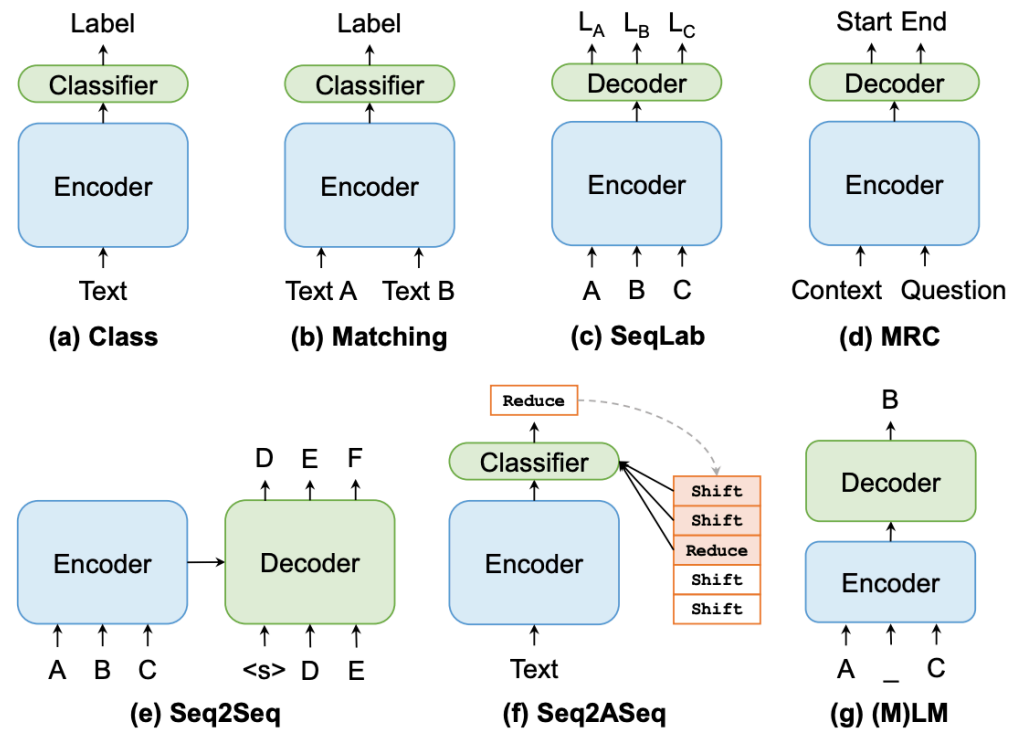
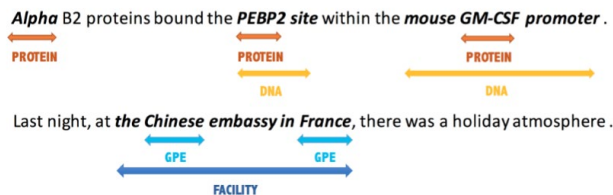


MRC as Paradigm

- ❑ MRC has great inspirations to the NLP tasks.
- **strong capacity** of MRC-style models
- unifying different tasks as **MRC formation**
- ❑ Generalized to other NLP tasks by reformulating them into the MRC format.

Example: nested named entity recognition

*Questoin: Find **XXX** in the text.*



[1] Sun, Tianxiang, et al. Paradigm Shift in Natural Language Processing. 2021.

[2] MCCANN, Bryan, et al. The natural language decathlon: Multitask learning as question answering. *arXiv:1806.08730*, 2018.

[3] KESKAR, Nitish Shirish, et al. Unifying Question Answering, Text Classification, and Regression via Span Extraction. *arXiv:1904.09286*, 2019.

[4] KESKAR, Nitish Shirish, et al. Unifying Question Answering, Text Classification, and Regression via Span Extraction. *arXiv:1904.09286*, 2019.

[5] LI, Xiaoya, et al. Entity-Relation Extraction as Multi-Turn Question Answering. ACL 2019. p. 1340-1350.

[6] LI, Xiaoya, et al. A Unified MRC Framework for Named Entity Recognition. ACL 2020.

Sources



Leaderboards

- SQuAD v1.1/2.0
- RACE
- CoQA
- QuAC
- DREAM
- MuTual
- ShARC
- ...



Venues

- AI/ML: NeurIPS, IJCAI, AAAI, etc.
- NLP/CL: ACL, EMNLP, COLING, etc.



Surveys

- Chen et al, 2018. Neural Reading Comprehension and Beyond
- Liu et al, 2019. Neural machine reading comprehension: Methods and trends
- Zhang et al, 2020. Machine Reading Comprehension: The Role of Contextualized Language Models and Beyond

SQuAD

HomeExplore 2.0Explore 1.1

What is SQuAD?

Stanford Question Answering Dataset (SQuAD) is a reading comprehension dataset, consisting of questions posed by crowdworkers on a set of Wikipedia articles, where the answer to every question is a segment of text, or span, from the corresponding reading passage, or the question might be unanswerable.

SQuAD2.0 combines the 100,000 questions in SQuAD1.1 with over 50,000 unanswerable questions written adversarially by crowdworkers to look similar to answerable ones. To do well on SQuAD2.0, systems must not only answer questions when possible, but also determine when no answer is supported by the paragraph and abstain from answering.

Explore SQuAD2.0 and model predictions

SQuAD2.0 paper (Rajpurkar & Jia et al. '18)

SQuAD 1.1, the previous version of the SQuAD dataset, contains 100,000+ question-answer pairs on 500+ articles.

Explore SQuAD1.1 and model predictions

SQuAD1.0 paper (Rajpurkar et al. '16)

Getting Started

We've built a few resources to help you get started with the dataset.

Download a copy of the dataset (distributed under the CC BY-SA 4.0 license):

Training Set v2.0 (40 MB)

Dev Set v2.0 (4 MB)

To evaluate your models, we have also made available the evaluation script we will use for official evaluation, along with a sample prediction file that the script will take as input. To run the evaluation, use `python evaluate-v2.0.py <path_to_dev-v2.0> <path_to_predictions>`.

Leaderboard

SQuAD2.0 tests the ability of a system to not only answer reading comprehension questions, but also abstain when presented with a question that cannot be answered based on the provided paragraph.

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1	Retro-Reader on ALBERT (ensemble) Shanghai Jiao Tong University Jan 10, 2020	90.115	92.580
2	ALBERT + DAAF + Verifier (ensemble) PINGAN Omni-Sinitic Nov 06, 2019	90.002	92.425
3	ALBERT (ensemble model) Google Research & TTIC https://arxiv.org/abs/1909.11942 Sep 18, 2019	89.731	92.215
4	ALBERT+Entailment DA (ensemble) CloudWalk Dec 08, 2019	88.761	91.745
5	Retro-Reader on ALBERT (single model) Shanghai Jiao Tong University Jan 19, 2020	88.107	91.419
5	XLNet + DAAF + Verifier (ensemble) PINGAN Omni-Sinitic Jul 22, 2019	88.592	90.859
5	albert+verifier (single model) Ping An Life Insurance Company AI Team Nov 22, 2019	88.355	91.019
6	[alber_m_transfor] (single model) QIANXIN Jan 15, 2020	88.186	90.939
6	ALBERT+Entailment DA Verifier (single model) CloudWalk Dec 08, 2019	87.847	91.265
6	ALBERT (single-model) huohua Jan 08, 2020	88.050	91.036
6	ALBERT + SFVerifier (single model) Senseforth AI Research https://www.senseforth.ai/ Jan 07, 2020	88.197	90.830
6	ALBERT (single model) Google Research & TTIC https://arxiv.org/abs/1909.11942 Sep 18, 2019	88.107	90.902

Outline

- ❖ Machine Reading Comprehension

 - ❖ Background, Development, Paradigm

- ❖ Techniques

 - ❖ Two-stage Solving Architecture

 - ❖ Traditional Matching Networks

 - ❖ Pre-trained Language Models

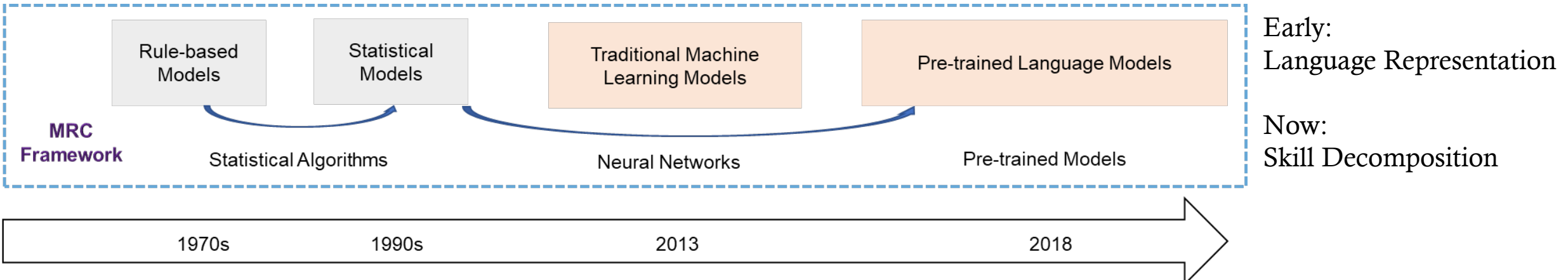
- ❖ Frontiers

 - ❖ Techniques

 - ❖ Tasks

 - ❖ Applications

<p> 1.1 Introduction 1.2 Background 1.3 Objectives 1.4 Scope 1.5 Methodology 1.6 Results 1.7 Conclusion 1.8 References 1.9 Appendix 1.10 Index 1.11 Glossary 1.12 Abbreviations 1.13 Acronyms 1.14 Footnotes 1.15 Endnotes 1.16 References 1.17 Appendix 1.18 Index 1.19 Glossary 1.20 Abbreviations 1.21 Acronyms 1.22 Footnotes 1.23 Endnotes 1.24 References 1.25 Appendix 1.26 Index 1.27 Glossary 1.28 Abbreviations 1.29 Acronyms 1.30 Footnotes 1.31 Endnotes 1.32 References 1.33 Appendix 1.34 Index 1.35 Glossary 1.36 Abbreviations 1.37 Acronyms 1.38 Footnotes 1.39 Endnotes 1.40 References 1.41 Appendix 1.42 Index 1.43 Glossary 1.44 Abbreviations 1.45 Acronyms 1.46 Footnotes 1.47 Endnotes 1.48 References 1.49 Appendix 1.50 Index 1.51 Glossary 1.52 Abbreviations 1.53 Acronyms 1.54 Footnotes 1.55 Endnotes 1.56 References 1.57 Appendix 1.58 Index 1.59 Glossary 1.60 Abbreviations 1.61 Acronyms 1.62 Footnotes 1.63 Endnotes 1.64 References 1.65 Appendix 1.66 Index 1.67 Glossary 1.68 Abbreviations 1.69 Acronyms 1.70 Footnotes 1.71 Endnotes 1.72 References 1.73 Appendix 1.74 Index 1.75 Glossary 1.76 Abbreviations 1.77 Acronyms 1.78 Footnotes 1.79 Endnotes 1.80 References 1.81 Appendix 1.82 Index 1.83 Glossary 1.84 Abbreviations 1.85 Acronyms 1.86 Footnotes 1.87 Endnotes 1.88 References 1.89 Appendix 1.90 Index 1.91 Glossary 1.92 Abbreviations 1.93 Acronyms 1.94 Footnotes 1.95 Endnotes 1.96 References 1.97 Appendix 1.98 Index 1.99 Glossary 1.100 Abbreviations 1.101 Acronyms 1.102 Footnotes 1.103 Endnotes 1.104 References 1.105 Appendix 1.106 Index 1.107 Glossary 1.108 Abbreviations 1.109 Acronyms 1.110 Footnotes 1.111 Endnotes 1.112 References 1.113 Appendix 1.114 Index 1.115 Glossary 1.116 Abbreviations 1.117 Acronyms 1.118 Footnotes 1.119 Endnotes 1.120 References 1.121 Appendix 1.122 Index 1.123 Glossary 1.124 Abbreviations 1.125 Acronyms 1.126 Footnotes 1.127 Endnotes 1.128 References 1.129 Appendix 1.130 Index 1.131 Glossary 1.132 Abbreviations 1.133 Acronyms 1.134 Footnotes 1.135 Endnotes 1.136 References 1.137 Appendix 1.138 Index 1.139 Glossary 1.140 Abbreviations 1.141 Acronyms 1.142 Footnotes 1.143 Endnotes 1.144 References 1.145 Appendix 1.146 Index 1.147 Glossary 1.148 Abbreviations 1.149 Acronyms 1.150 Footnotes 1.151 Endnotes 1.152 References 1.153 Appendix 1.154 Index 1.155 Glossary 1.156 Abbreviations 1.157 Acronyms 1.158 Footnotes 1.159 Endnotes 1.160 References 1.161 Appendix 1.162 Index 1.163 Glossary 1.164 Abbreviations 1.165 Acronyms 1.166 Footnotes 1.167 Endnotes 1.168 References 1.169 Appendix 1.170 Index 1.171 Glossary 1.172 Abbreviations 1.173 Acronyms 1.174 Footnotes 1.175 Endnotes 1.176 References 1.177 Appendix 1.178 Index 1.179 Glossary 1.180 Abbreviations 1.181 Acronyms 1.182 Footnotes 1.183 Endnotes 1.184 References 1.185 Appendix 1.186 Index 1.187 Glossary 1.188 Abbreviations 1.189 Acronyms 1.190 Footnotes 1.191 Endnotes 1.192 References 1.193 Appendix 1.194 Index 1.195 Glossary 1.196 Abbreviations 1.197 Acronyms 1.198 Footnotes 1.199 Endnotes 1.200 References</</p>
--



Two-stage Solving Architecture

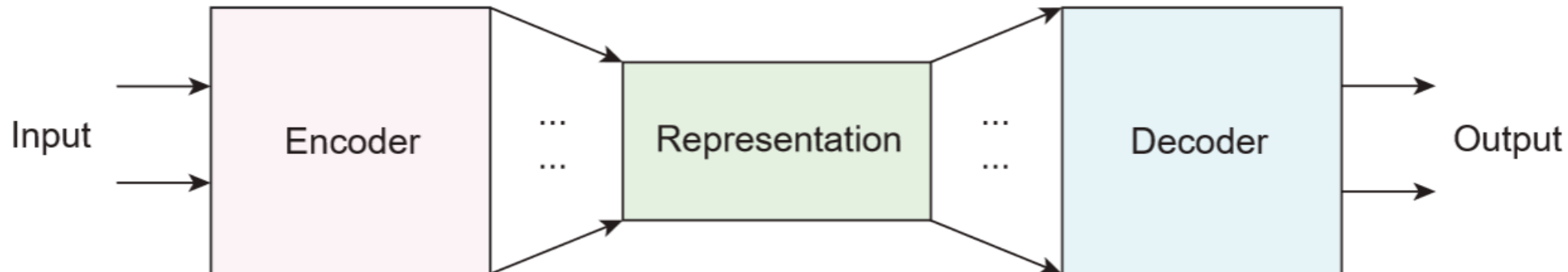
Inspired by **Dual process theory** of cognition psychology:

The cognitive process of human brains potentially involves two distinct types of procedures:

- **contextualized perception** (reading): gather information in an implicit process
- **analytic cognition** (comprehension): conduct the controlled reasoning and execute goals

Standard MRC system:

- building a PrLM as **Encoder**;
- designing ingenious mechanisms as **Decoder** according to task characteristics.



Encoder

❑ Multiple Granularity Features

- Language Units: word, character, subword.
- Salient Features: Linguistic features, such as part-of-speech, named entity tags, semantic role labeling tags, syntactic features, and binary Exact Match features.

❑ Structured Knowledge Injection (Transformer/GNN)

- Linguistic Structures
- Commonsense

❑ Contextualized Sentence Representation

- Embedding pretraining

Encoder (salient features)

SemBERT: Semantics-aware BERT

Passage

- *...Harvard was a founding member of the Association of American Universities in 1900. James Bryant Conant led the university through the Great Depression and World War II and began to reform the curriculum and liberalize admissions after the war. The undergraduate college became coeducational after its 1977merger with Radcliffe College.....*

Question

- *What was the name of the leader through the Great Depression and World War II?*

Semantic Role Labeling (SRL)

- *[James Bryant Conant]_{ARG0} [led]_{VERB} [the university]_{ARG1} through [the Great Depression and World War II]_{ARG2}*

Answer

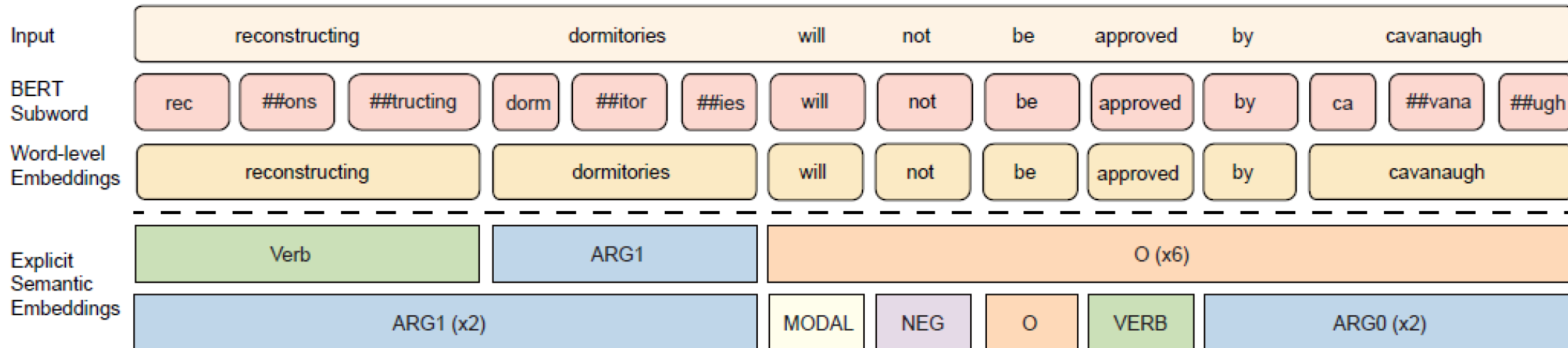
- *James Bryant Conant*

Problem: Who did what to whom, when and why?

Encoder (salient features)

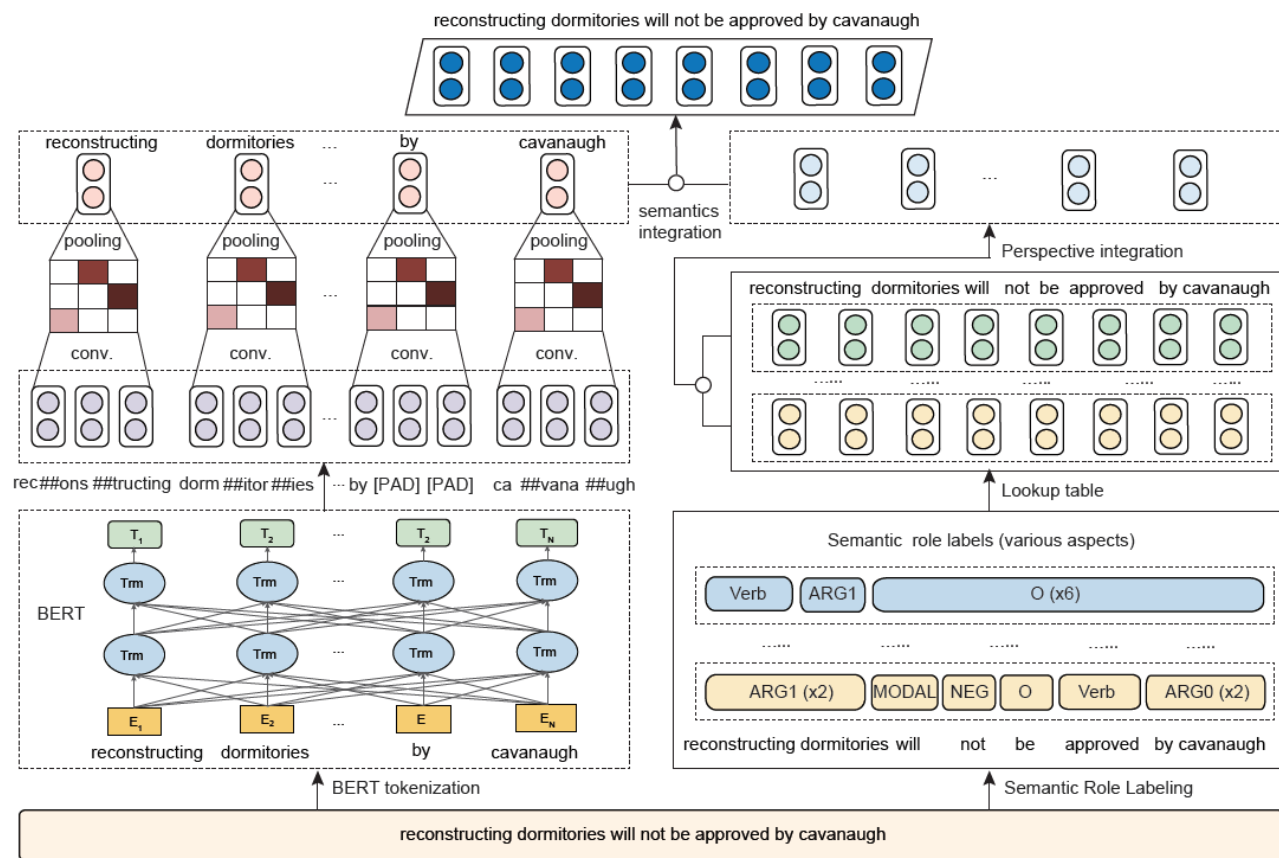
SemBERT: Semantics-aware BERT

- ELMo & BERT: only take **Plain contextual** features
- SemBERT: introduce **Explicit contextual Semantics**, **Deeper representation?**
 - Semantic Role Labeler + BERT encoder



Encoder (salient features)

SemBERT: Semantics-aware



Method	Classification		Natural Language Inference			Semantic Similarity			Score
	CoLA	SST-2	MNLI	QNLI	RTE	MRPC	QQP	STS-B	-
	(mc)	(acc)	m/mm(acc)	(acc)	(acc)	(F1)	(F1)	(pc)	-
Leaderboard (September, 2019)									
ALBERT	69.1	97.1	91.3/91.0	99.2	89.2	93.4	74.2	92.5	89.4
RoBERTa	67.8	96.7	90.8/90.2	98.9	88.2	92.1	90.2	92.2	88.5
XLNET	67.8	96.8	90.2/89.8	98.6	86.3	93.0	90.3	91.6	88.4
In literature (April, 2019)									
BiLSTM+ELMo+Attn	36.0	90.4	76.4/76.1	79.9	56.8	84.9	64.8	75.1	70.5
GPT	45.4	91.3	82.1/81.4	88.1	56.0	82.3	70.3	82.0	72.8
GPT on STILTs	47.2	93.1	80.8/80.6	87.2	69.1	87.7	70.1	85.3	76.9
MT-DNN	61.5	95.6	86.7/86.0	-	75.5	90.0	72.4	88.3	82.2
BERT _{BASE}	52.1	93.5	84.6/83.4	-	66.4	88.9	71.2	87.1	78.3
BERT _{LARGE}	60.5	94.9	86.7/85.9	92.7	70.1	89.3	72.1	87.6	80.5
Our implementation									
SemBERT _{BASE}	57.8	93.5	84.4/84.0	90.9	69.3	88.2	71.8	87.3	80.9
SemBERT _{LARGE}	62.3	94.6	87.6/86.3	94.6	84.5	91.2	72.8	87.8	82.9

GLUE 实验结果

Model	EM	F1	Model	Dev	Test
#1 BERT + DAE + AoA†	85.9	88.6	In literature		
#2 SG-Net†	85.2	87.9	DRCN (Kim et al. 2018)	-	90.1
#3 BERT + NGM + SST†	85.2	87.7	SJRC (Zhang et al. 2019)	-	91.3
U-Net (Sun et al. 2018)	69.2	72.6	MT-DNN (Liu et al. 2019)†	92.2	91.6
BMR + ELMo + Verifier (Hu et al. 2018)	71.7	74.2	Our implementation		
BERT _{LARGE}	80.5	83.6	BERT _{BASE}	90.8	90.7
SemBERT _{LARGE}	82.4	85.2	BERT _{LARGE}	91.3	91.1
SemBERT _{BASE}	84.8	87.9	SemBERT _{BASE}	91.2	91.0
SemBERT _{LARGE}			SemBERT _{LARGE}	92.3	91.6

SQuAD 实验结果

SNLI 实验结果

SNLI: The **best** among all submissions.

<https://nlp.stanford.edu/projects/snli/>

SQuAD2.0: The **best** among all the published work.

GLUE: substantial gains over all the tasks.

Decoder

❑ Matching Network

- Attention Sum, Gated Attention, Self-matching, Attention over Attention, Co-match Attention, Dual Co-match Attention, etc.

❑ Fine-grained Reasoning Network

- Decouple the context into multiple elements and measure the relationships for reasoning

❑ Answer Pointer

- Pointer Network for span prediction
- Reinforcement learning based self-critical learning to predict more acceptable answers

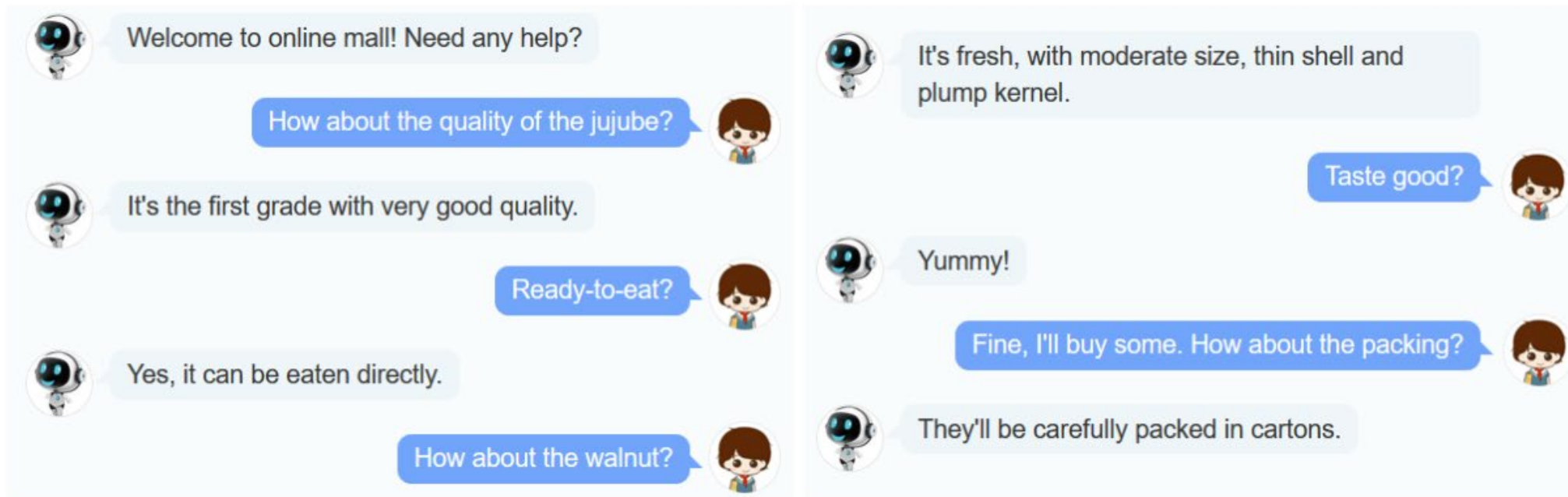
❑ Answer Verifier

- Threshold-based answerable verification
- Multitask-style verification
- External parallel verification

❑ Answer Type Predictor for multi-type MRC tasks

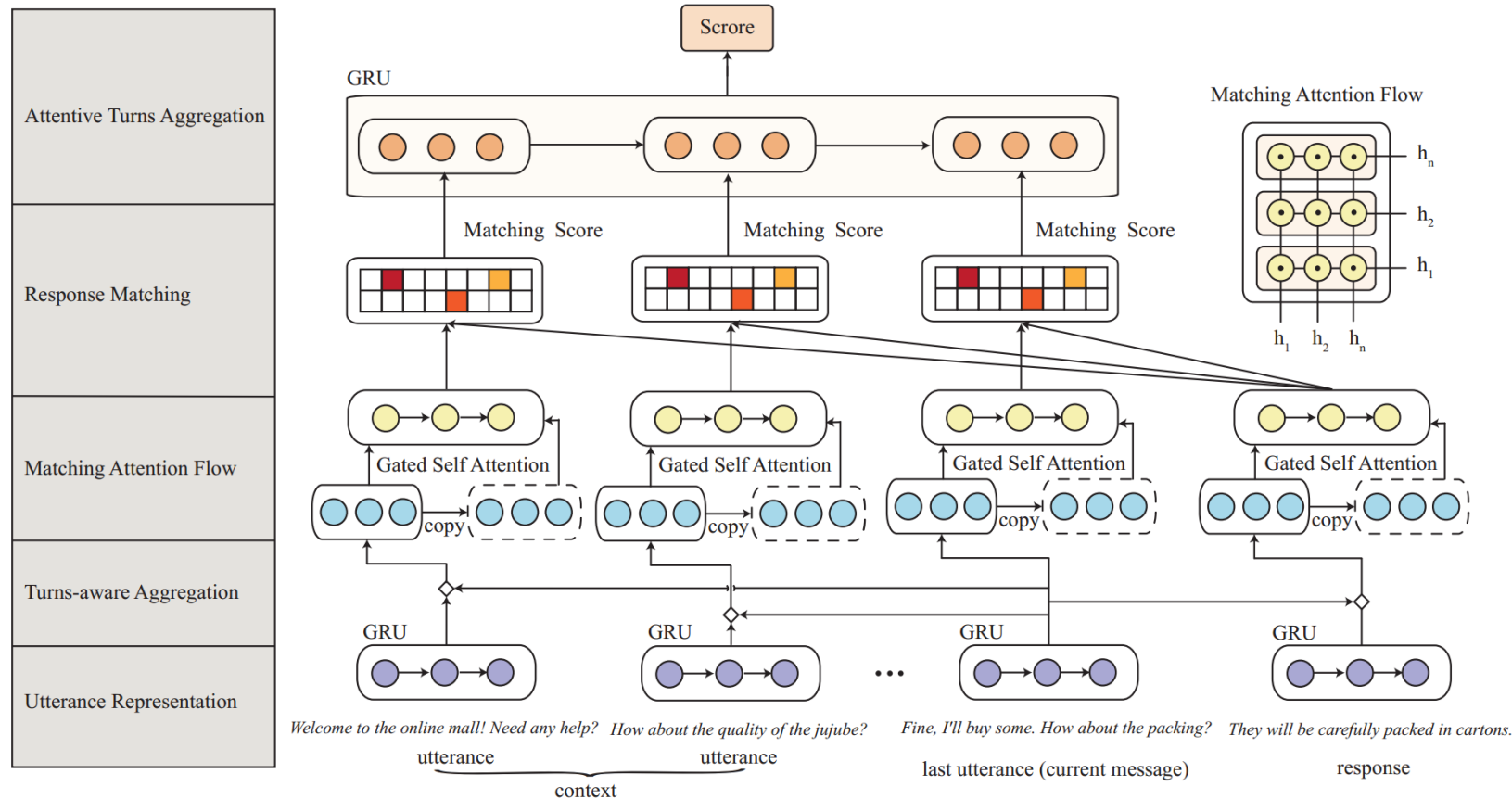
Decoder (Deep Utterance Aggregation)

- ❑ Challenge: **long utterances, multiple intentions, topic shift**, etc.
- ❑ Aim: recognize the **key information** from complex dialogue history
- ❑ Solution: deep utterance aggregation framework (**DUA**)
- ❑ Corpus: a new **E-commerce Dialogue Corpus**



Decoder (Deep Utterance Aggregation)

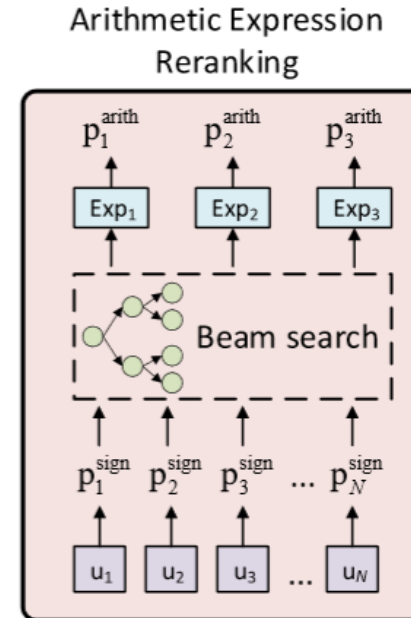
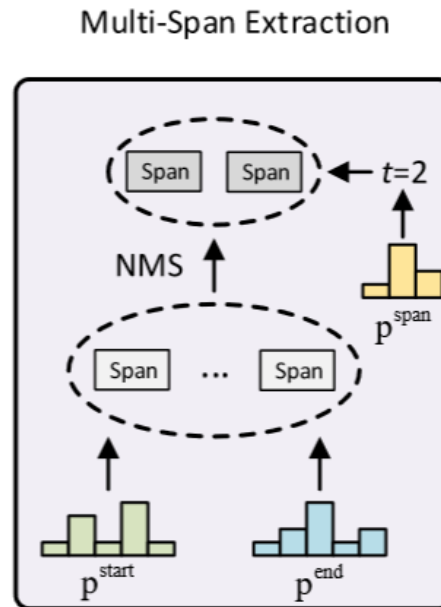
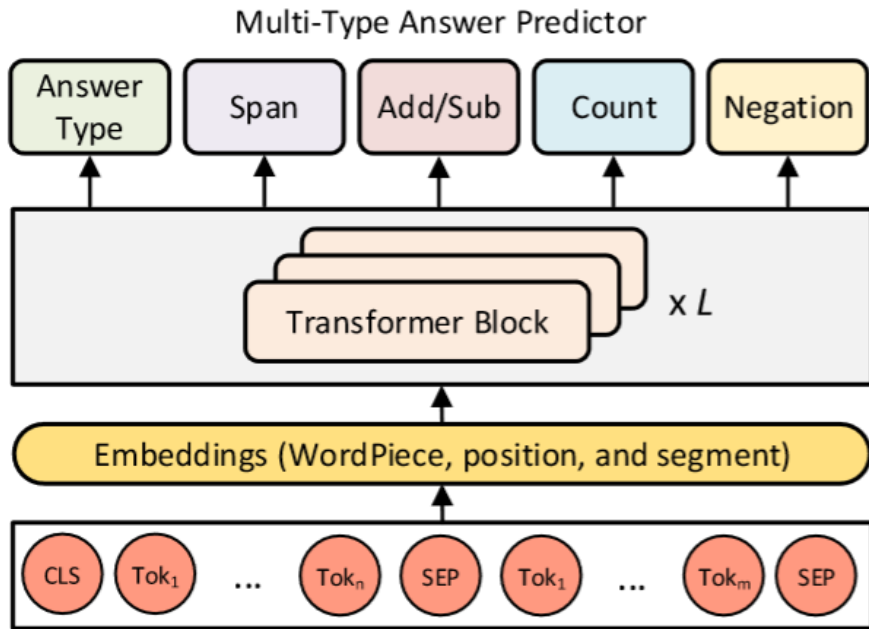
- ❑ Capture the main information in each utterance (**self attention**, first introduced)
- ❑ Model the **information flow through the utterances** in dialogue history
- ❑ Match the relationship **between utterance and candidate response**



Highlight the importance of
the last utterance.

Decoder

□ Answer Type Predictor for multi-type MRC tasks



(MTMSN model from Hu et al., 2019)

Hu, Minghao, et al. A multi-type multi-span network for reading comprehension that requires discrete reasoning. EMNLP-IJCNLP 2019.

Outline

- ❖ Machine Reading Comprehension

 - ❖ Background, Development, Paradigm

- ❖ Techniques

 - ❖ Two-stage Solving Architecture

 - ❖ Traditional Matching Networks

 - ❖ Pre-trained Language Models

- ❖ Frontiers

 - ❖ Techniques

 - ❖ Tasks

 - ❖ Applications

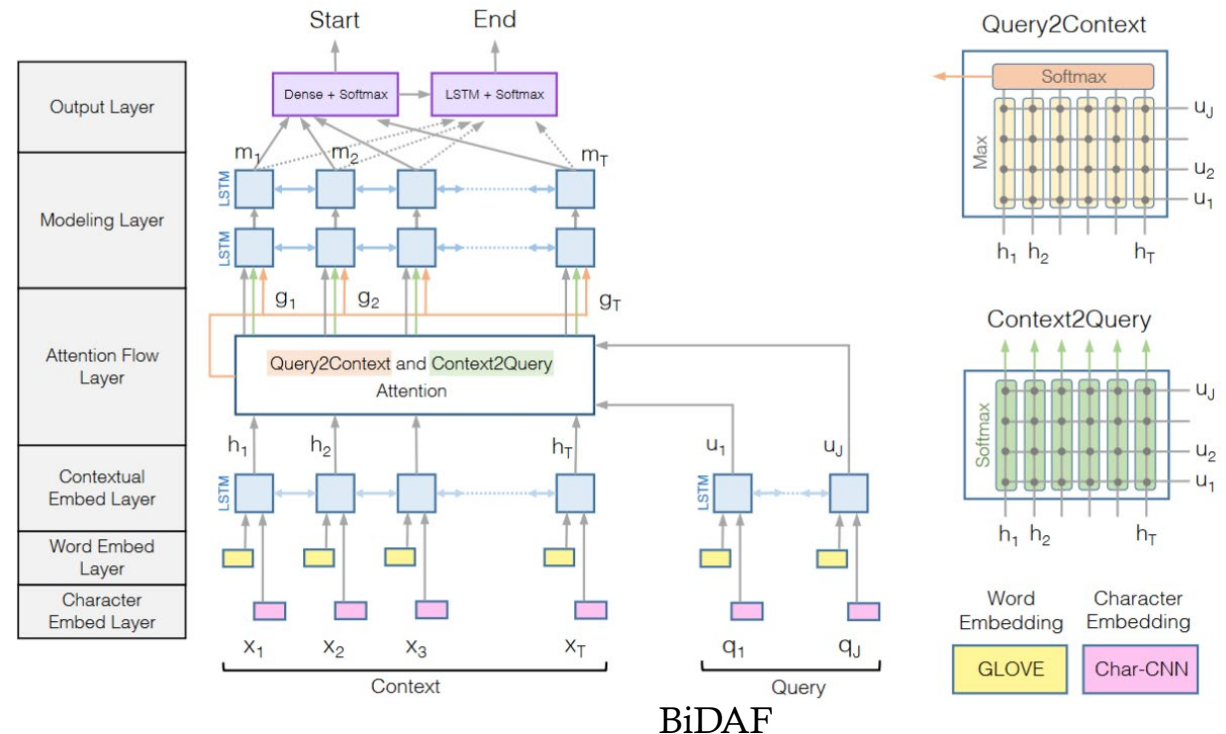
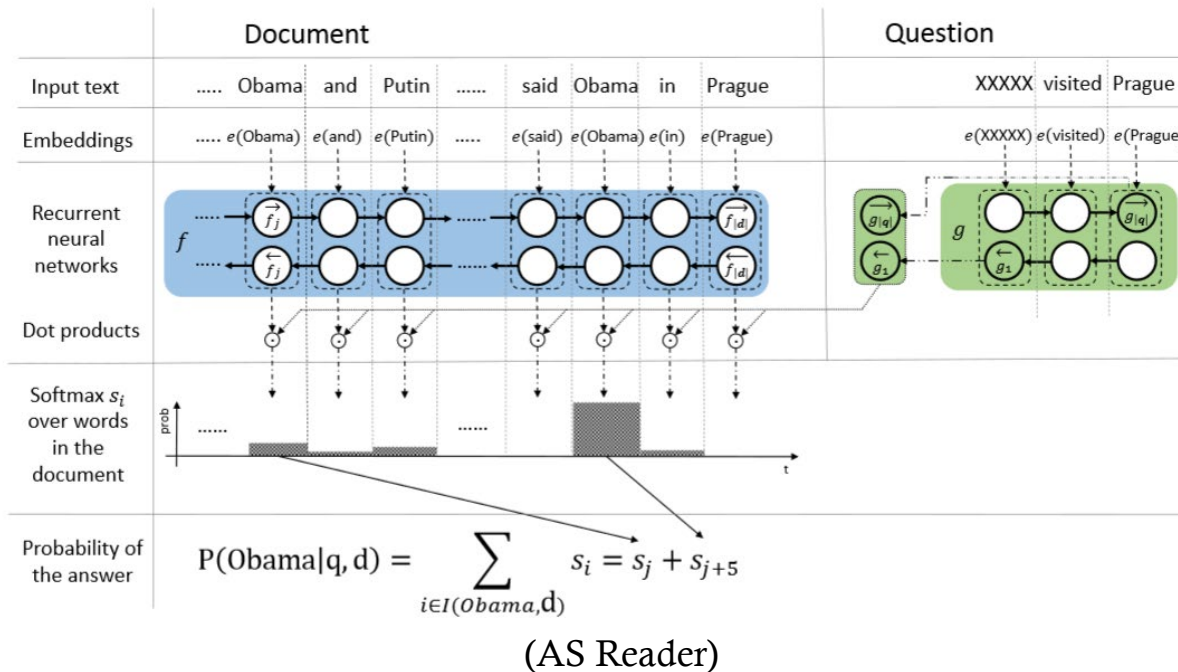
Stage 1: Traditional Matching Networks

❑ Matching Network:

- Attention Sum, Gated Attention, Self-matching, Attention over Attention, BiDAF, etc.

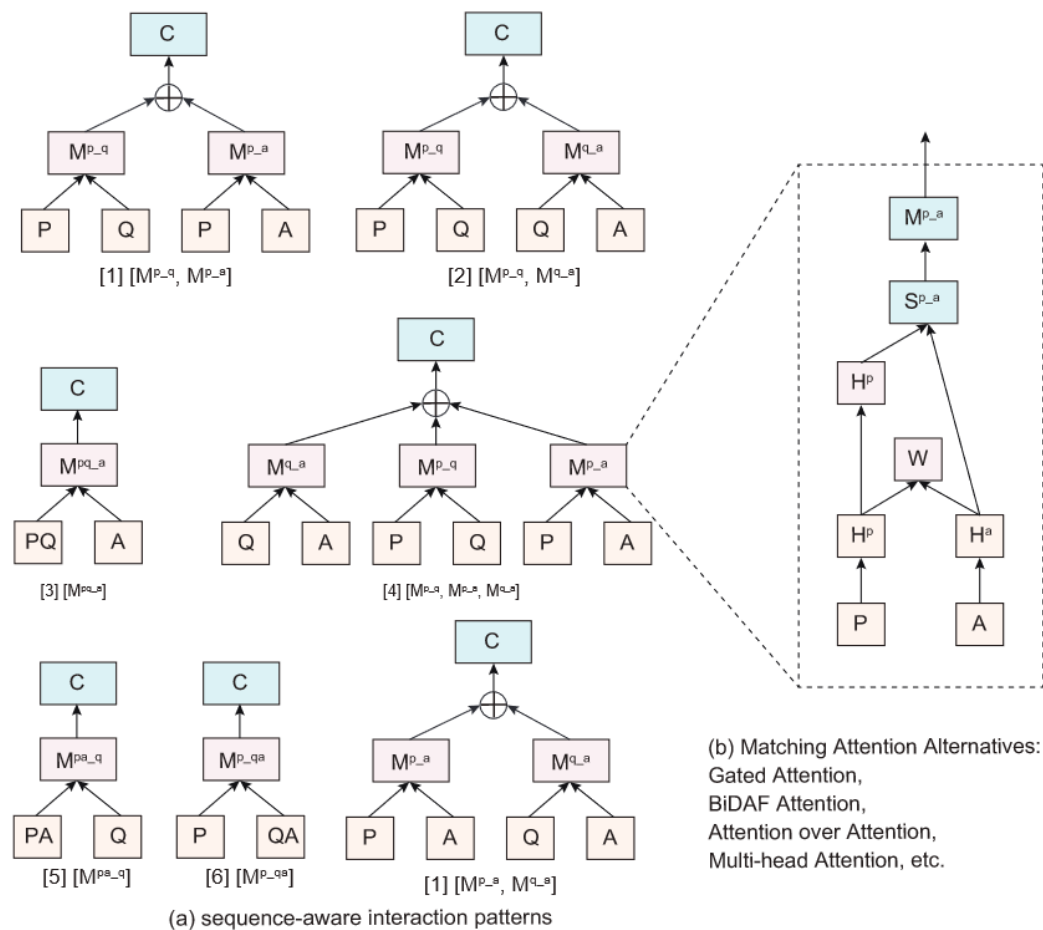
❑ Attention weights: sum, dot, gating, etc.

❑ Attention Direction: question-aware, passage aware, self-attention, bidirectional, etc.



❑ Attention Granularity : word-level, sequence-level, hierarchical, etc.

Stage 1: Traditional Matching Networks



Method	Att. Type	CNN val	CNN test	DailyMail val	DailyMail test
Attentive Reader (Hermann et al. 2015)	UA	61.6	63.0	70.5	69.0
AS Reader (Kadlec et al. 2016)	UA	68.6	69.5	75.0	73.9
Iterative Attention (Sordoni et al. 2016)	UA	72.6	73.3	-	-
Stanford AR (Chen, Bolton, and Manning 2016)	UA	73.8	73.6	77.6	76.6
GARader (Dhingra et al. 2017)	UA	73.0	73.8	76.7	75.7
AoA Reader (Cui et al. 2017)	BA	73.1	74.4	-	-
BiDAF (Seo et al. 2017)	BA	76.3	76.9	80.3	79.6

Model	Matching	M	H	RACE
Human Ceiling Performance (Lai et al. 2017)		95.4	94.2	94.5
Amazon Mechanical Turk (Lai et al. 2017)		85.1	69.4	73.3
HAF (Zhu et al. 2018a)	$[M^{P-A}; M^{P-Q}; M^{Q-A}]$	45.0	46.4	46.0
MRU (Tay, Tuan, and Hui 2018)	$[M^{P-Q-A}]$	57.7	47.4	50.4
HCM (Wang et al. 2018a)	$[M^{P-Q}; M^{P-A}]$	55.8	48.2	50.4
MMN (Tang, Cai, and Zhuo 2019)	$[M^{Q-A}; M^{A-Q}; M^{P-Q}; M^{P-A}]$	61.1	52.2	54.7
GPT (Radford et al. 2018)	$[M^{P-Q-A}]$	62.9	57.4	59.0
RSM (Sun et al. 2019b)	$[M^{P-QA}]$	69.2	61.5	63.8
DCMN (Zhang et al. 2019a)	$[M^{PQA}]$	77.6	70.1	72.3
OCN (Ran et al. 2019a)	$[M^{P-Q-A}]$	76.7	69.6	71.7
BERT _{large} (Pan et al. 2019b)	$[M^{P-Q-A}]$	76.6	70.1	72.0
XLNet (Yang et al. 2019c)	$[M^{P-Q-A}]$	85.5	80.2	81.8
+ DCMN+ (Zhang et al. 2020a)	$[M^{P-Q}; M^{P-O}; M^{Q-O}]$	86.5	81.3	82.8
RoBERTa (Liu et al. 2019c)	$[M^{P-Q-A}]$	86.5	81.8	83.2
+ MMM (Jin et al. 2019a)	$[M^{P-Q-A}]$	89.1	83.3	85.0
ALBERT (Jin et al. 2019a)	$[M^{P-Q-A}]$	89.0	85.5	86.5
+ DUMA (Zhu, Zhao, and Li 2020)	$[M^{P-QA}; M^{QA-P}]$	90.9	86.7	88.0
Megatron-BERT (Shoeybi et al. 2019)	$[M^{P-Q-A}]$	91.8	88.6	89.5

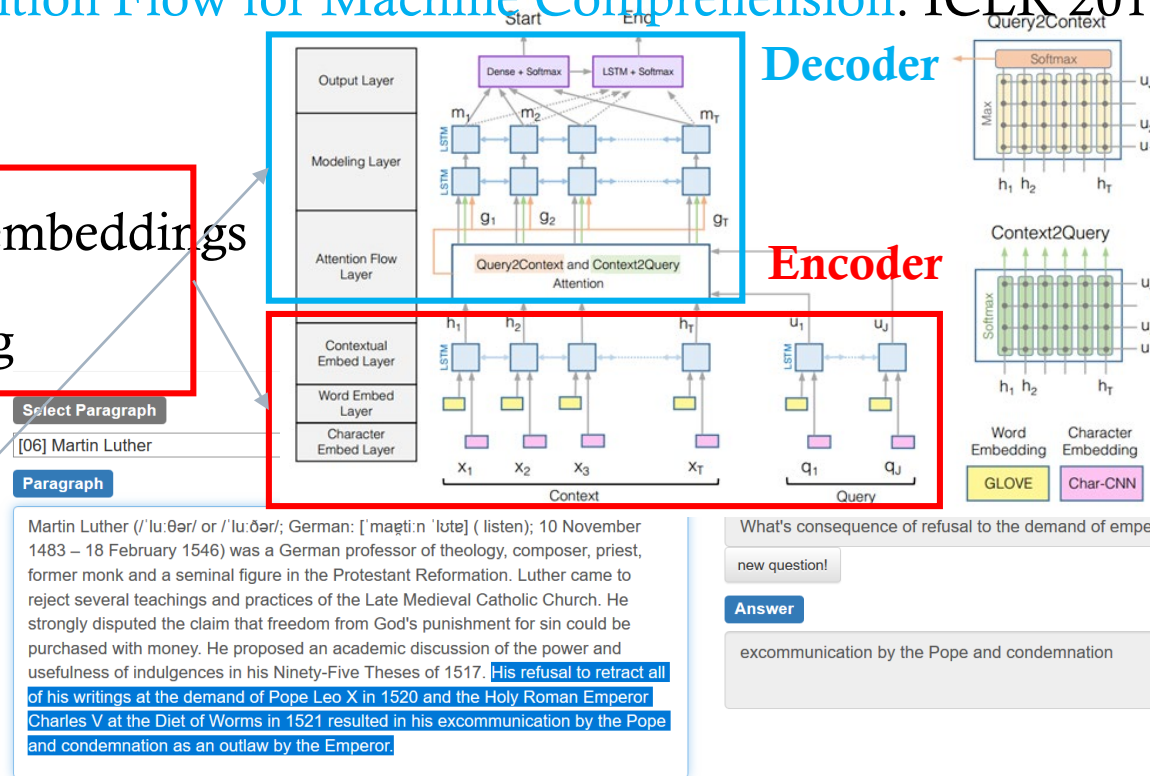
Typical Architecture

□ BiDAF

- Minjoon Seo, Aniruddha Kembhavi, Ali Farhadi, Hannaneh Hajishirzi. 2017. [Bidirectional Attention Flow for Machine Comprehension](#). ICLR 2017.

Hierarchical structure:

- Word + Char level embeddings
- Contextual encoding
- Attention modules
- Answer prediction



Reading Strategies & Data Augmentation

Reading Strategy based on human reading patterns

- Learning to skim text
- Learning to stop reading
- Retrospective reading
- Back and forth reading, highlighting, and self-assessment

Data Augmentation

- Combining various MRC datasets as training data augmentation
- Multi-tasking
- Automatic question generation, such as back translation and synthetic generation

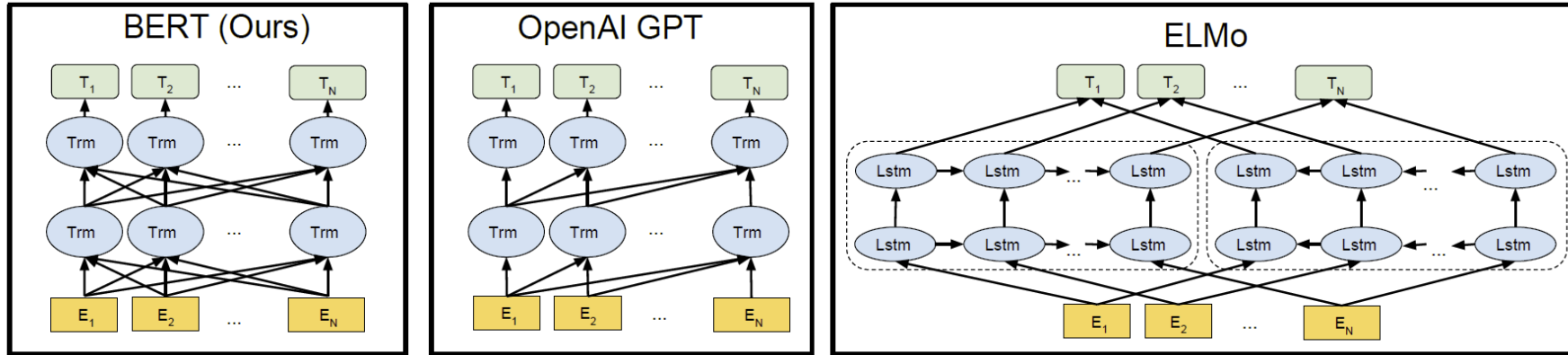
[1] Yu, Adams Wei, et al. Learning to skim text. ACL 2017.

[2] Shen, Yelong, et al. Reasonet: Learning to stop reading in machine comprehension. KDD 2017.

[3] Zhang, Zhuosheng, et al. Retrospective reader for machine reading comprehension. AAAI 2021.

[4] Sun, Kai, et al. Improving machine reading comprehension with general reading strategies. NAACL 2019.

Stage 2: Pre-trained Language Models



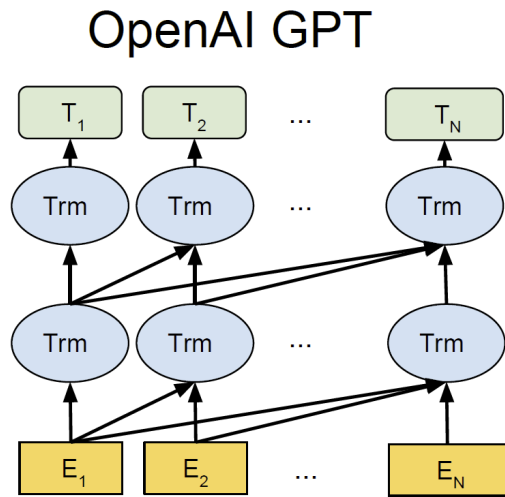
ELMo: Embedding from Language Models

GPT: Generative Pre-Training

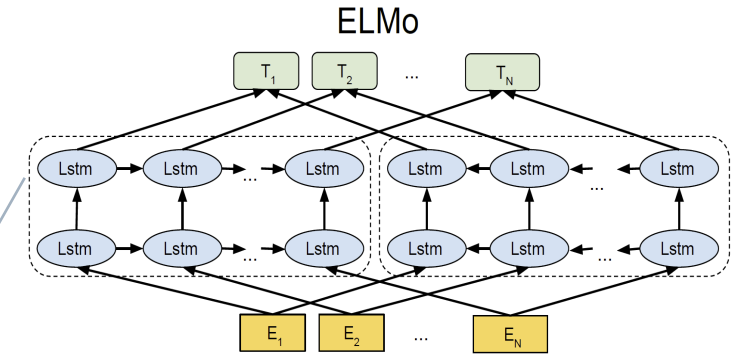
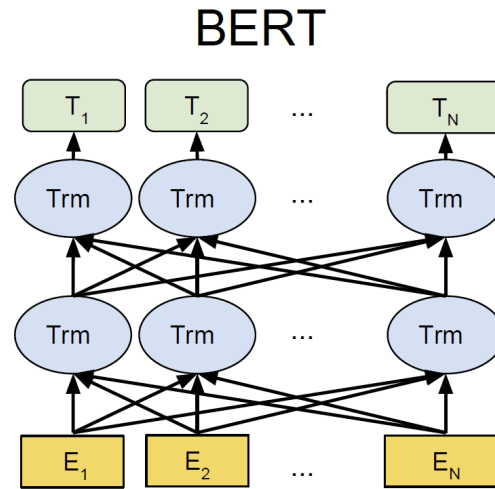
BERT: Pre-training of Deep Bidirectional Transformers

- [1] Peters, Matthew E., et al. Deep contextualized word representations. NAACL-HLT. 2018.
- [2] Radford, Alec, et al. Improving language understanding by generative pre-training. (2018).
- [3] Devlin, Jacob, et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. NAACL-HLT. 2019.

From GPT、ELMo、 Word2Vec to BERT

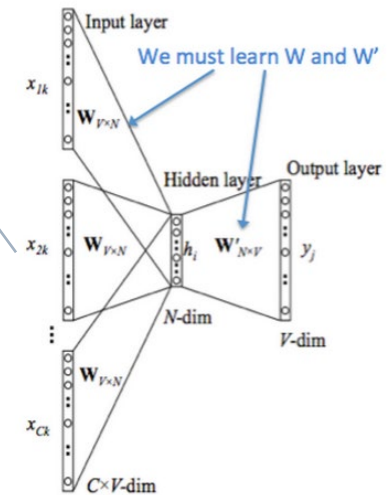


Bi-directional
(idea)



Transformer as feature extractor

Bi-directional language modeling
(method)



BERT

BERT - Bidirectional Encoder Representations from Transformers

Huge Parameters:

BERT base: $L=12$, $H=768$, $A=12$, Total Parameters=110M

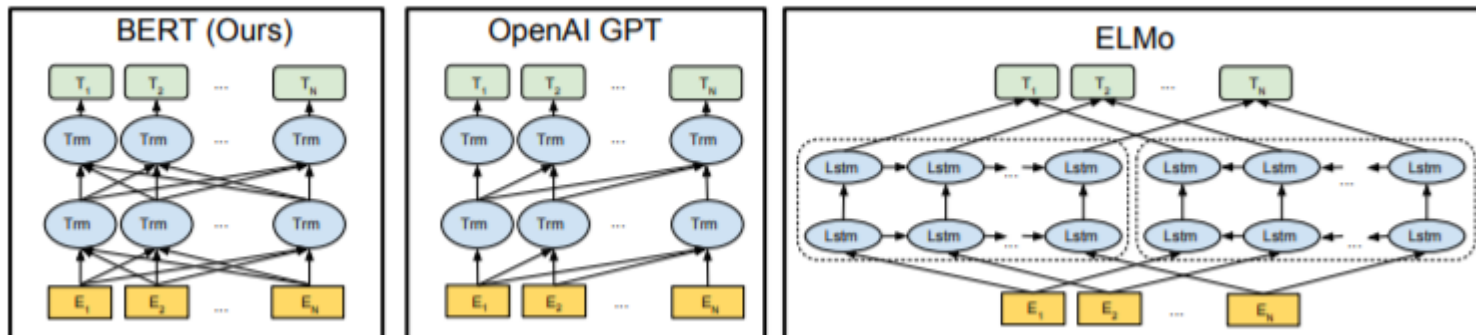
BERT large: $L=24$, $H=1024$, $A=16$, Total Parameters=340M

(L -transformer blocks, H - dimension of hidden state, A – self-attention heads)

Large corpus: BooksCorpus (800M words) + English Wikipedia (2,500M words)

Computing power: BERT base 16 TPU*4 day BERT large 64 TPU *4 day

BERT vs GPT vs ELMo



Devlin, Jacob, et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. NAACL-HLT. 2019.

BERT Pre-training

Task #1: Masked LM

replace the chosen words with [MASK]
then predict it
Not always replace the word with [MASK]

- 80% of the time: Replace the word with the [MASK] token, e.g., my dog is hairy → my dog is [MASK]
- 10% of the time: Replace the word with a random word, e.g., my dog is hairy → my dog is apple
- 10% of the time: Keep the word unchanged, e.g., my dog is hairy → my dog is hairy. The purpose of this is to bias the representation towards the actual observed word.

Task #2: Next Sentence Prediction

[CLS] sentence A [SEP] sentence B
[SEP]
50% of the time B is the actual next sentence that follows A, and 50% of the time it is a random sentence from the corpus

Input = [CLS] the man went to [MASK] store [SEP]

he bought a gallon [MASK] milk [SEP]

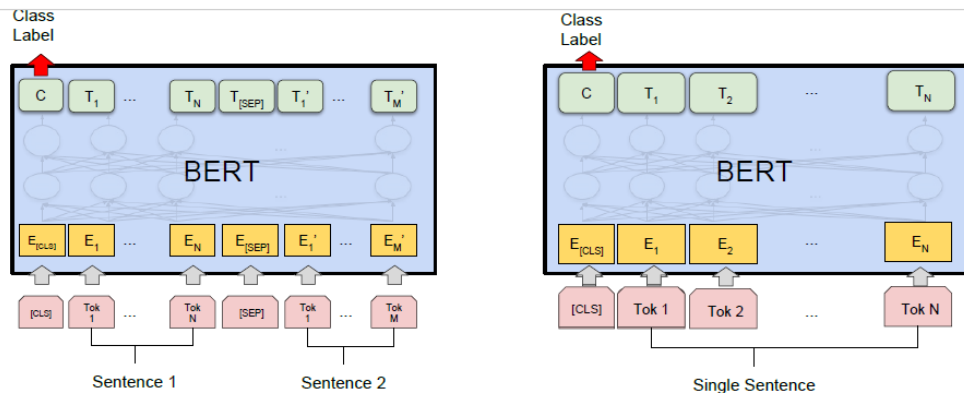
Label = IsNext

Input = [CLS] the man [MASK] to the store [SEP]

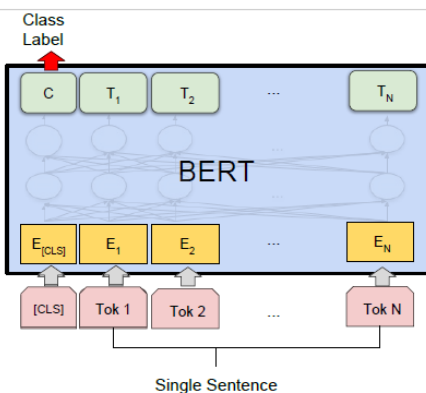
penguin [MASK] are flight ##less birds [SEP]

Label = NotNext

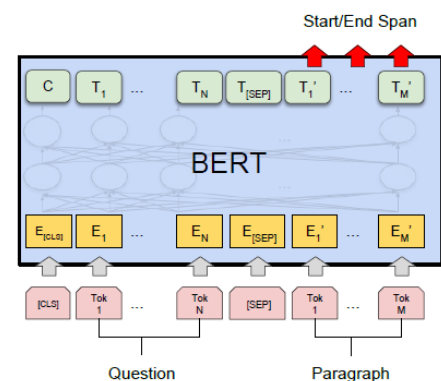
BERT Fine-tuning



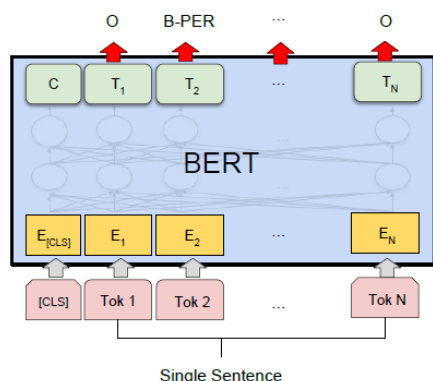
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



(b) Single Sentence Classification Tasks:
SST-2, CoLA



(c) Question Answering Tasks:
SQuAD v1.1



(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

Table 1: GLUE Test results, scored by the GLUE evaluation server. The number below each task denotes the

System	Dev EM	Dev F1	Test EM	Test F1
Leaderboard (Oct 8th, 2018)				
Human	-	-	82.3	91.2
#1 Ensemble - nlnet	-	-	86.0	91.7
#2 Ensemble - QANet	-	-	84.5	90.5
#1 Single - nlnet	-	-	83.5	90.1
#2 Single - QANet	-	-	82.5	89.3
Published				
BiDAF+ELMo (Single)	-	85.8	-	-
R.M. Reader (Single)	78.9	86.3	79.5	86.6
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5
Ours				
BERT _{BASE} (Single)	80.8	88.5	-	-
BERT _{LARGE} (Single)	84.1	90.9	-	-
BERT _{LARGE} (Ensemble)	85.8	91.8	-	-
BERT _{LARGE} (Sgl.+TriviaQA)	84.2	91.1	85.1	91.8
BERT _{LARGE} (Ens.+TriviaQA)	86.2	92.2	87.4	93.2

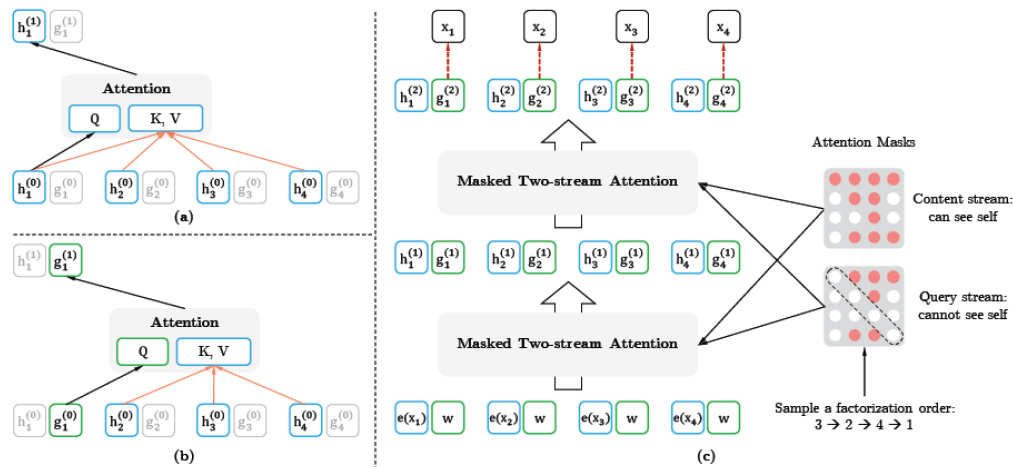
Table 2: SQuAD results. The BERT ensemble is 7x

System	Dev F1	Test F1
ELMo+BiLSTM+CRF	95.7	92.2
CVT+Multi (Clark et al., 2018)	-	92.6
BERT _{BASE}	96.4	92.4
BERT _{LARGE}	96.6	92.8

Table 3: CoNLL-2003 Named Entity Recognition results. The hyperparameters were selected using the Dev set, and the reported Dev and Test scores are averaged over 5 random restarts using those hyperparameters.

XLNet: Token Permutation + Two-stream Attention

- Token permutation + Two-stream Attention
 - Using autoregressive mechanism to overcome the shortcomings of BERT (Masked LM)
 - Permute the tokens in the sentence, and make the LM predictions



Training corpus:

- 13G: BooksCorpus + English Wikipedia
- 16G: Giga5
- 19G: ClueWeb 2012-B
- 78G: Common Crawl

- Computation: 512 TPU v3, 500K steps, batch size = 2048, 2.5 days

Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Ruslan Salakhutdinov, Quoc V. Le. 2019. XLNet: Generalized Autoregressive Pretraining for Language Understanding, NeurIPS 2019.

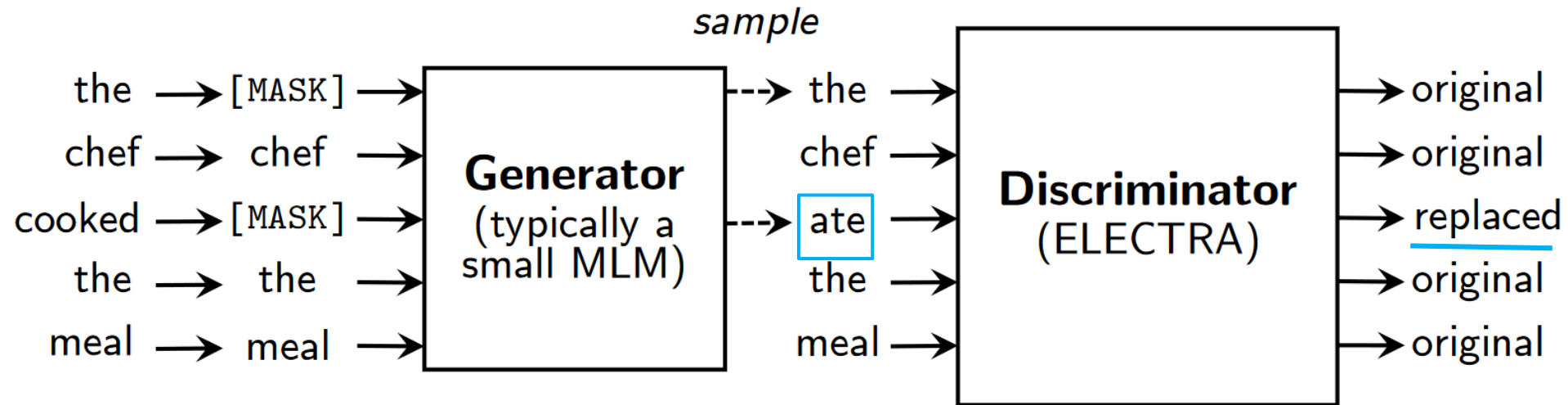
ALBERT: Sentence Order Prediction

- Three improvements:
 - Modify the Embedding (E) and hidden states (H) into the dimension $H \gg E$, instead of $E=H$ in BERT
 - Use full layer parameter sharing, including all forward networks and attention weights (significantly reduce the model size)
 - Modify the sentence training objective (NSP) of BERT to sentence order prediction (SOP)

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, Radu Soricut. ALBERT: A Lite BERT for Self-supervised Learning of Language Representation. *ICLR* 2020.

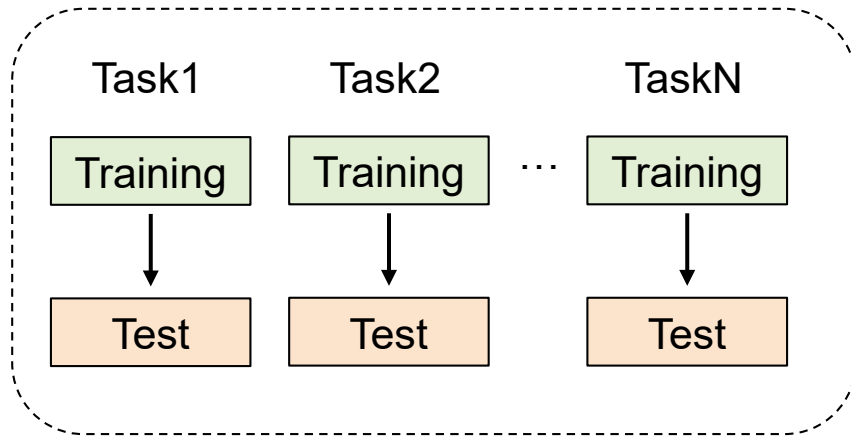
ELECTRA

- Predicts whether each token in the corrupted input was replaced by a generator sample or not.



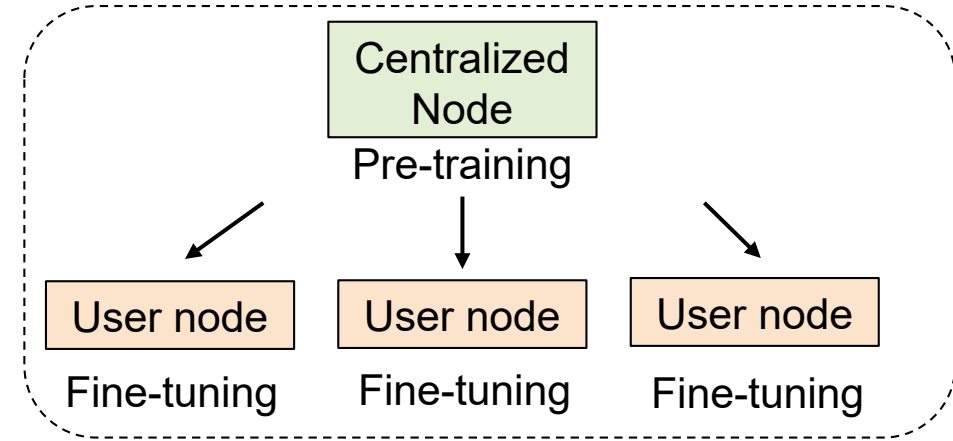
Kevin Clark, Minh-Thang Luong, Quoc V. Le, Christopher D. Manning. ELECTRA: Pre-training Text Encoders as Discriminators Rather Than Generators. *ICLR* 2020.

PrLM: New Paradigm



Previous

Each user trains individual machine learning models for each task.



Now

The central node trains the generalized language model (pre-training) and provides the nearly completed model for users as the standard module for task-specific fine-tuning.

Individual
training



Centralized pre-training + individual fine-tuning

Extreme case: GPT3 gives predictions directly after pre-training, eliminating the fine-tuning process

From Language Models to Language Representation

- ❑ MRC and other application NLP need a full **sentence encoder**,
 - Deep contextual information is required in MRC
 - Word and sentence should be represented as embeddings.
 - ❑ Model can be trained in a style of n -gram language model
 - ❑ So that there comes the **language representation** which includes
 - Contextual encoder (**model architecture**)
 - n -gram language model (**training object**)
 - **Training methods**
- The representation for each word depends on the entire context in which it is used, **dynamic embedding**.

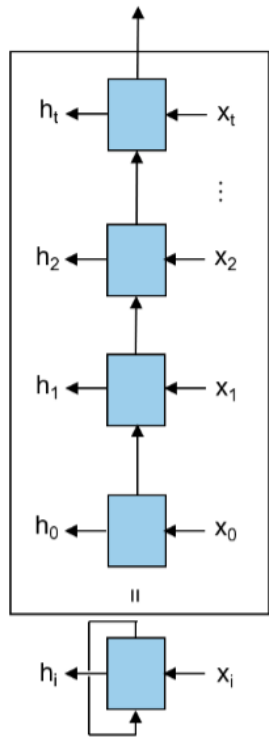
LM Contextualization:
Sentence -> Encoder -> Repr.

The Elements of PrLMs

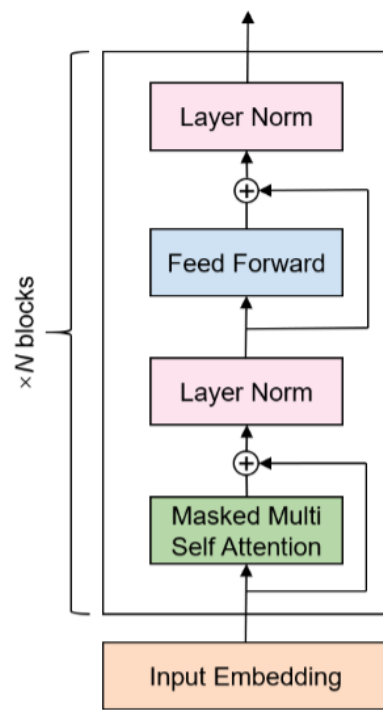
- ❑ Encoder architectures
 - RNN/Transformer/...
- ❑ Training objectives
 - (Autoregressive / denoising) task construction
- ❑ Sampling (training) methods

Architectures of PrLMs

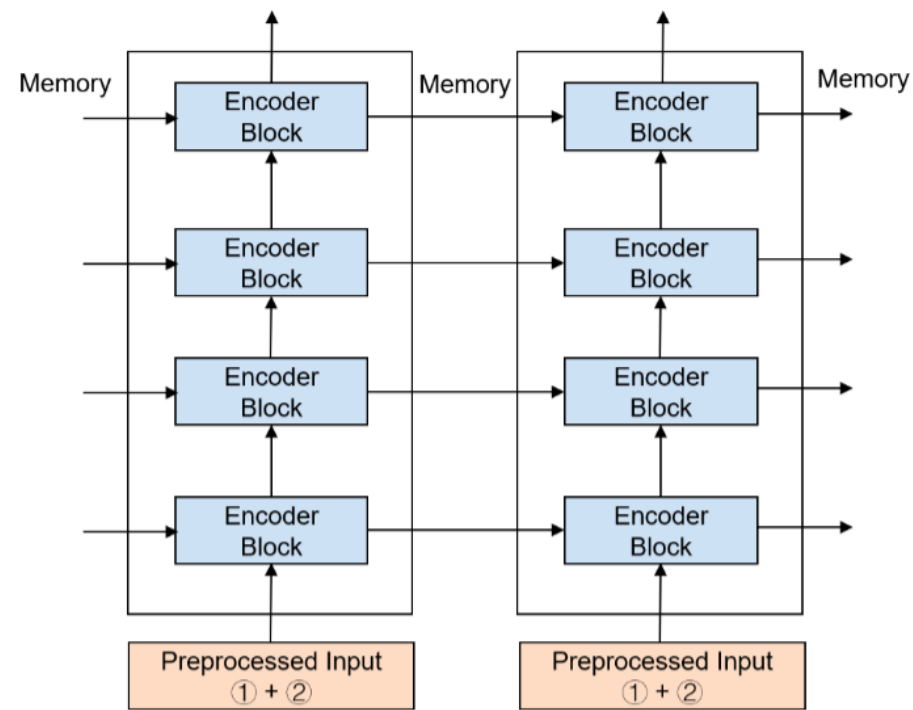
- RNN: GRU/LSTM
- Transformer
- Transformer-XL



(a) RNN



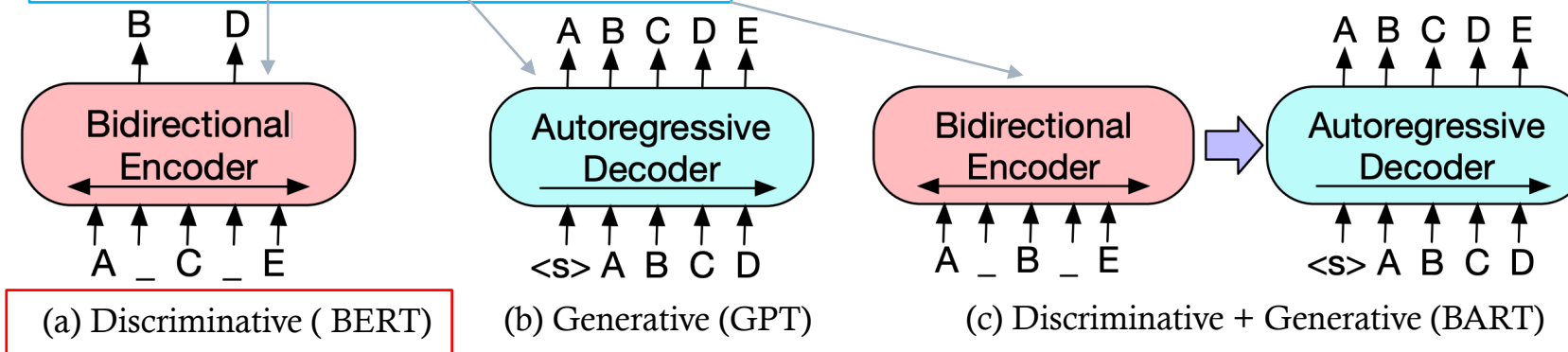
(b) Transformer



(c) Transformer-XL

Training Objectives

- ❑ Constructing the training samples with **generalized autoregressive** method
- ❑ Discriminative vs. Generative
 - **Discriminative**: Predict the corrupted tokens (BERT, ALBERT, ELECTRA, etc)
 - Useful for discriminative tasks like span-based MRC
 - **Generative**: Predict the complete sentence via Decoder (GPT 1-3, etc)
 - Helpful for generative tasks like machine translation
 - **Discriminative + Generative**: Predict the complete sentence via Decoder (BART)



Training Methods (Denoising)

- ❑ LM is an **automatic denoising encoder** in language
- ❑ Manually constructing different levels of corrupted units of natural language text
- ❑ Masking units:

- Subword
- Word
- Span
- Entity
- Etc.

❑ ➔ Edit Operations

- deletion
- addition
- permutation/reordering
- replacement

	word	sentence
deletion	Masking	NSP
replacement		
addition		
permutation	XLNet	SOP

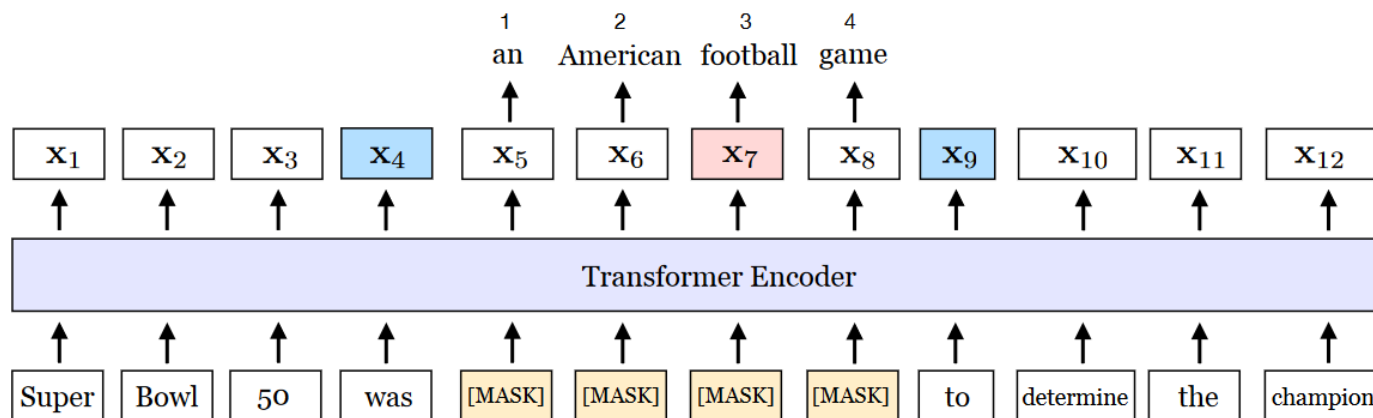
❑ Training strategies:

- direct prediction
- generative-discriminative (ELECTRA)

BERT_{WWM} vs. SpanBERT

- BERT_{WWM} : whole word masking
- SpanBERT
 - Mask continues spans
 - Span boundary objective

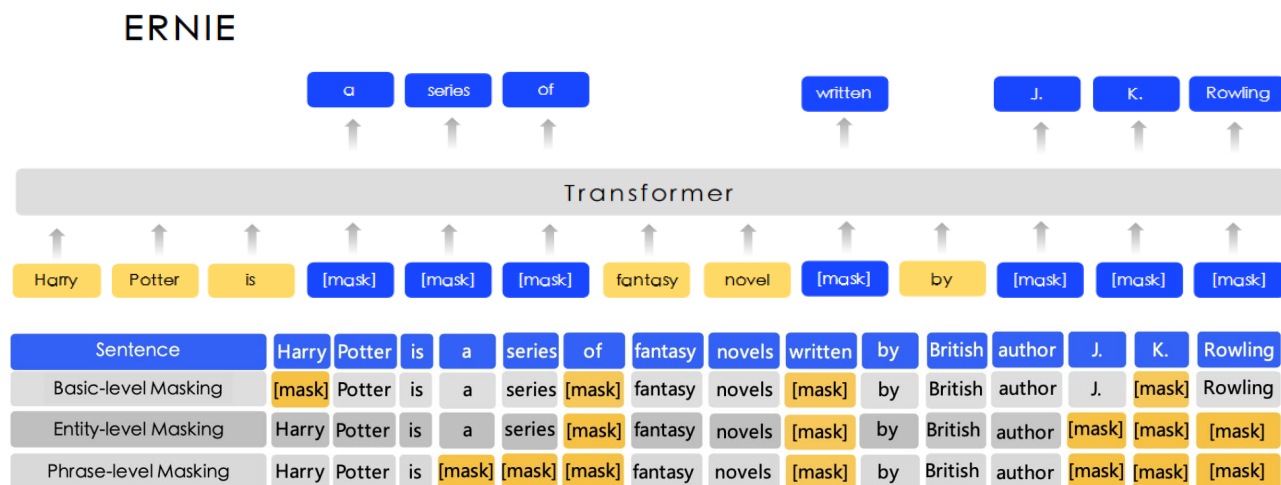
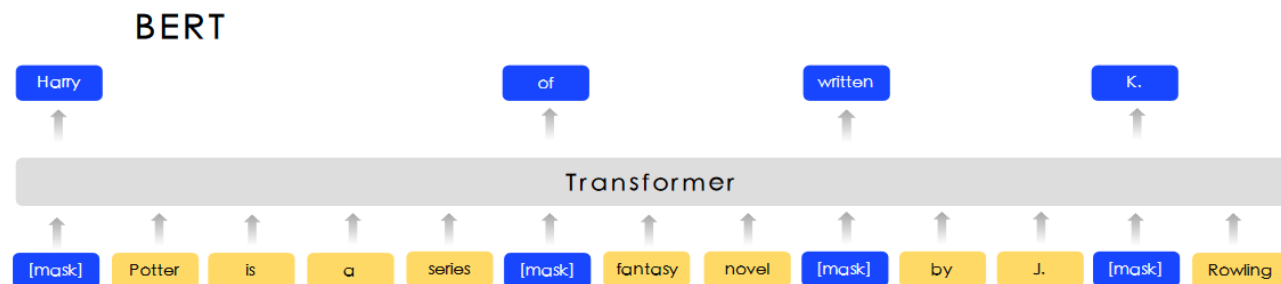
$$\begin{aligned}\mathcal{L}(\text{football}) &= \mathcal{L}_{\text{MLM}}(\text{football}) + \mathcal{L}_{\text{SBO}}(\text{football}) \\ &= -\log P(\text{football} \mid \mathbf{x}_7) - \log P(\text{football} \mid \mathbf{x}_4, \mathbf{x}_9, \mathbf{p}_3)\end{aligned}$$



Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, Omer Levy. 2020. SpanBERT: Improving Pre-training by Representing and Predicting Spans. TACL.

Masking Knowledge Units: ERNIE

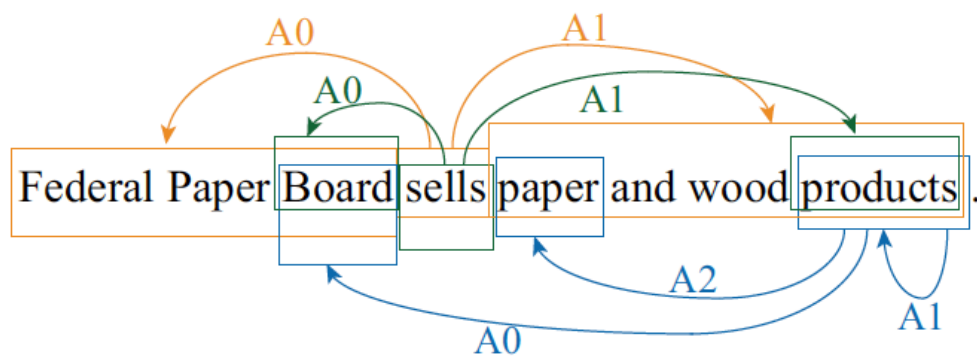
- Knowledge-enhanced masking: **entities** + **phrases**



Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, Hua Wu. ERNIE: Enhanced Representation through Knowledge Integration. ACL 2020.

Linguistic Mask: LIMIT-BERT

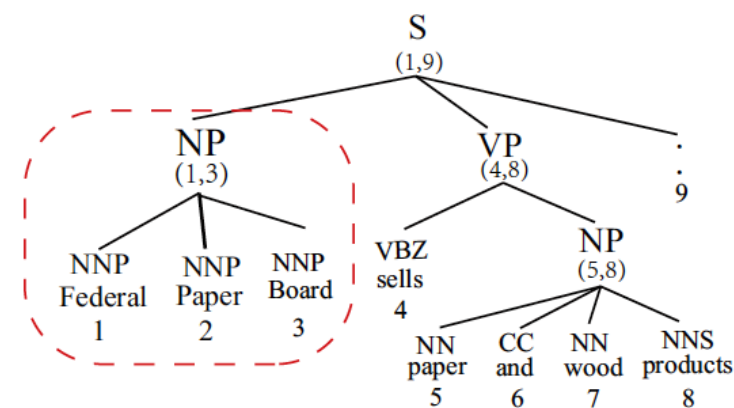
- Mask Strategy: syntactic and semantic masks
- Multitask Learning: improve the modeling performance of language model with linguistic tasks.



Span and Dependency SRL

federal paper board [MASK] paper and wood [MASK] .

(a) Semantic Phrase Masking.



Constituent Syntactic Tree

[MASK] [MASK] [MASK] sells paper and wood products .

(b) Syntactic Phrase Masking.

Derivative of PrLM

□ Embedding Units

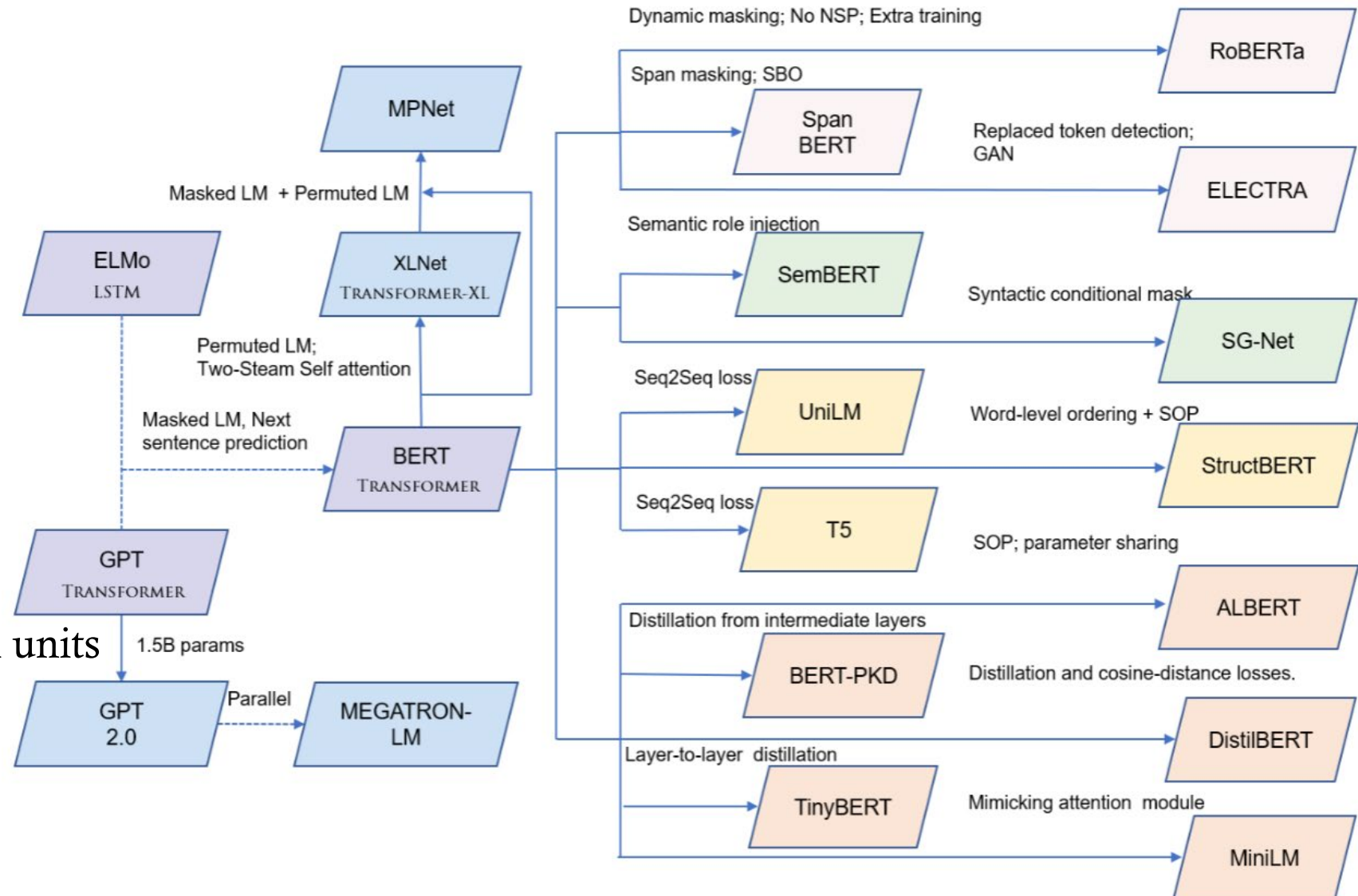
- Character
- Subword
- Word

□ Masked Units

- Subword
- Word/Span/Sentence
- Knowledge pieces
- Statistically meaningful units

□ Sequence Prediction

- Sentence relevance
- Sentence order

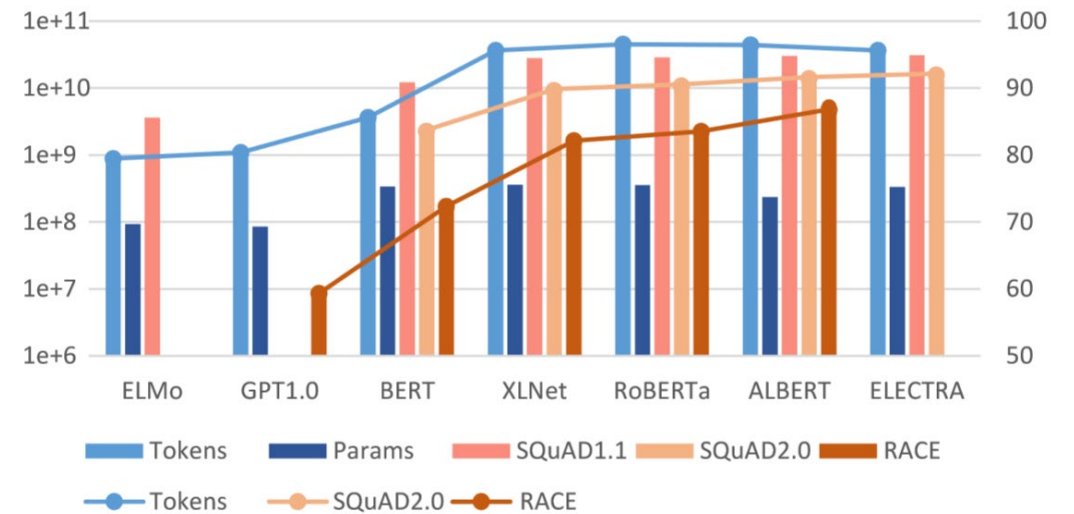


PrLMs greatly boost the benchmark of current MRC

Models	Encoder	EM	F1	↑ EM	↑ F1
Human (Rajpurkar, Jia, and Liang 2018)	-	82.304	91.221	-	-
Match-LSTM (Wang and Jiang 2016)	RNN	64.744	73.743	-	-
DCN (Xiong, Zhong, and Socher 2016)	RNN	66.233	75.896	1.489	2.153
Bi-DAF (Seo et al. 2017)	RNN	67.974	77.323	3.230	3.580
Mnemonic Reader (Hu, Peng, and Qiu 2017)	RNN	70.995	80.146	6.251	6.403
Document Reader (Chen et al. 2017)	RNN	70.733	79.353	5.989	5.610
DCN+ (Xiong, Zhong, and Socher 2017)	RNN	75.087	83.081	10.343	9.338
r-net (Wang et al. 2017)	RNN	76.461	84.265	11.717	10.522
MEMEN (Pan et al. 2017)	RNN	78.234	85.344	13.490	11.601
QANet (Yu et al. 2018)*	TRFM	80.929	87.773	16.185	14.030
<hr/>					
CLMs					
ELMo (Peters et al. 2018)	RNN	78.580	85.833	13.836	12.090
BERT (Devlin et al. 2018)*	TRFM	85.083	91.835	20.339	18.092
SpanBERT (Joshi et al. 2020)	TRFM	88.839	94.635	24.095	20.892
XLNet (Yang et al. 2019c)	TRFM-XL	89.898	95.080	25.154	21.337

Models	Encoder	SQuAD 2.0	↑ F1	RACE	↑ Acc
Human (Rajpurkar, Jia, and Liang 2018)	-	91.221	-	-	-
GPT _{v1} (Radford et al. 2018)	TRFM	-	-	59.0	-
BERT (Devlin et al. 2018)	TRFM	83.061	-	72.0	-
SemBERT (Zhang et al. 2020b)	TRFM	87.864	4.803	-	-
SG-Net (Zhang et al. 2020c)	TRFM	87.926	4.865	-	-
RoBERTa (Liu et al. 2019c)	TRFM	89.795	6.734	83.2	24.2
ALBERT (Lan et al. 2019)	TRFM	90.902	7.841	86.5	27.5
XLNet (Yang et al. 2019c)	TRFM-XL	90.689	7.628	81.8	22.8
ELECTRA (Clark et al. 2019c)	TRFM	91.365	8.304	-	-

Method	Tokens	Size	Params	SQuAD1.1 Dev	SQuAD1.1 Test	SQuAD2.0 Dev	SQuAD2.0 Test	RACE
ELMo	800M	-	93.6M	85.6	85.8	-	-	-
GPT _{v1}	985M	-	85M	-	-	-	-	59.0
XLNet _{large}	33B	-	360M	94.5	95.1*	88.8	89.1*	81.8
BERT _{large}	3.3B	13GB	340M	91.1	91.8*	81.9	83.0	72.0†
RoBERTa _{large}	-	160GB	355M	94.6	-	89.4	89.8	83.2
ALBERT _{xxlarge}	-	157GB	235M	94.8	-	90.2	90.9	86.5
ELECTRA _{large}	33B	-	335M	94.9	-	90.6	91.4	-



- Knowledge from large-scale corpora
- Deep architectures

Correlations Between MRC and PrLM

MRC and PrLM are **complementary** to each other.

MRC serves as an appropriate testbed for language representation, which is the focus of PrLMs.

The progress of PrLMs greatly promotes MRC tasks, achieving impressive gains of model performance.

The initial applications of PrLMs. The concerned NLU task can also be regarded as a special case of MRC

	NLU			MRC	
	SNLI	GLUE	SQuAD1.1	SQuAD2.0	RACE
ELMo	✓	✗	✓	✗	✗
GPT _{v1}	✓	✓	✗	✗	✓
BERT	✗	✓	✓	✓	✗
RoBERTa	✗	✓	✓	✓	✓
ALBERT	✗	✓	✓	✓	✓
XLNet	✗	✓	✓	✓	✓
ELECTRA	✗	✓	✓	✓	✗

Interpretability of Human-parity Performance

- ❑ What kind of **knowledge** or **reading comprehension skills** the systems have grasped?
- ❑ For MRC model side
 - overestimated ability of MRC systems that do not necessarily provide **human-level** understanding
 - unprecise **benchmarking** on the existing datasets.
 - suffers from **adversarial attacks**
- ❑ For PrLM encoder side:
 - good at linguistic notions of **syntax** and **coreference**.
 - struggles with challenging **inferences** and role-based **event prediction**
 - obvious failures with the meaning of **negation**
- ❑ Decomposition of Prerequisite Skills
 - decompose the skills required by the dataset and take skill-wise evaluations
 - provide more explainable and convincing benchmarking of model capacity

Outline

- ❖ Machine Reading Comprehension

 - ❖ Background, Development, Paradigm

- ❖ Techniques

 - ❖ Two-stage Solving Architecture

 - ❖ Traditional Matching Networks

 - ❖ Pre-trained Language Models

- ❖ Frontiers

 - ❖ Techniques

 - ❖ Tasks

 - ❖ Applications

New Frontiers

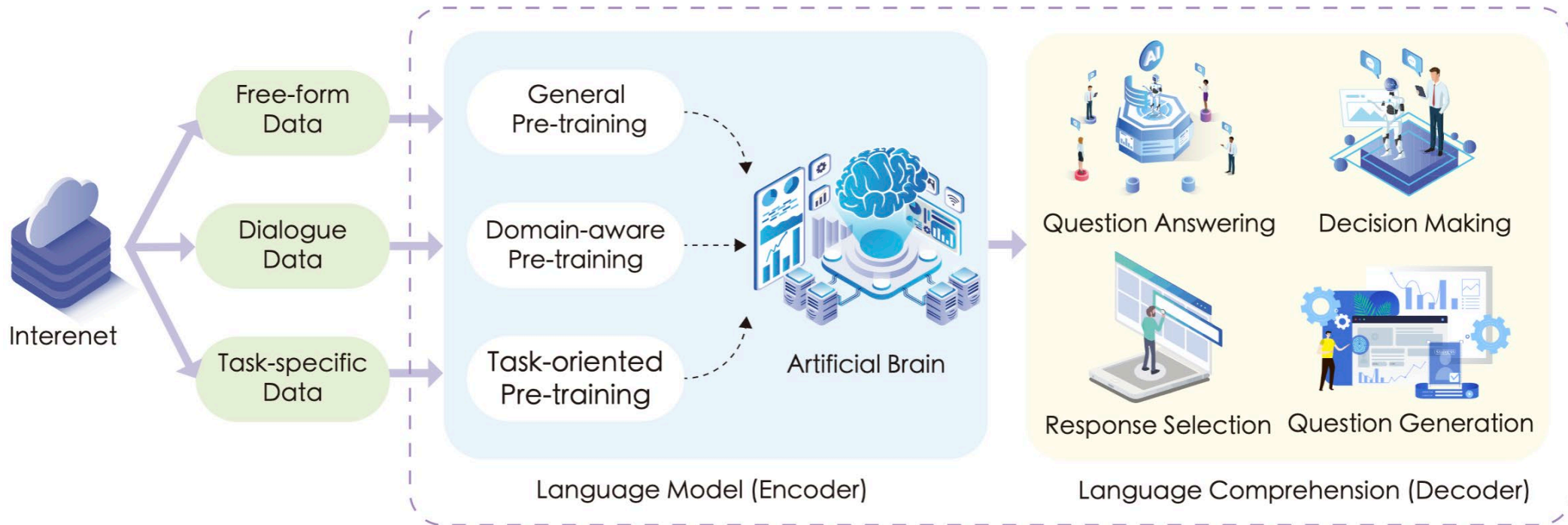
- ❑ Techniques
 - Domain/Task-adaptive Pre-training
 - Graph-aware Knowledge Structure Modeling
- ❑ Tasks
 - Multi-turn Dialogue Comprehension
 - Logical Reasoning
 - Commonsense Reasoning
- ❑ Applications
 - Open-domain QA
 - Multilingual, Multimodal, Multitask

New Frontiers

- ❑ Techniques
 - Domain/Task-adaptive Pre-training
 - Graph-aware Knowledge Structure Modeling
- ❑ Tasks
 - Multi-turn Dialogue Comprehension
 - Logical Reasoning
 - Commonsense Reasoning
- ❑ Applications
 - Open-domain QA
 - Multilingual, Multimodal, Multitask

Domain/Task-adaptive Pre-training

- General-purpose Pre-training (e.g., mask language modeling)
- Domain-aware Pre-training (e.g., science, news, medical domains)
- Task-oriented Pre-training (e.g., dialogue/discourse structure modeling)

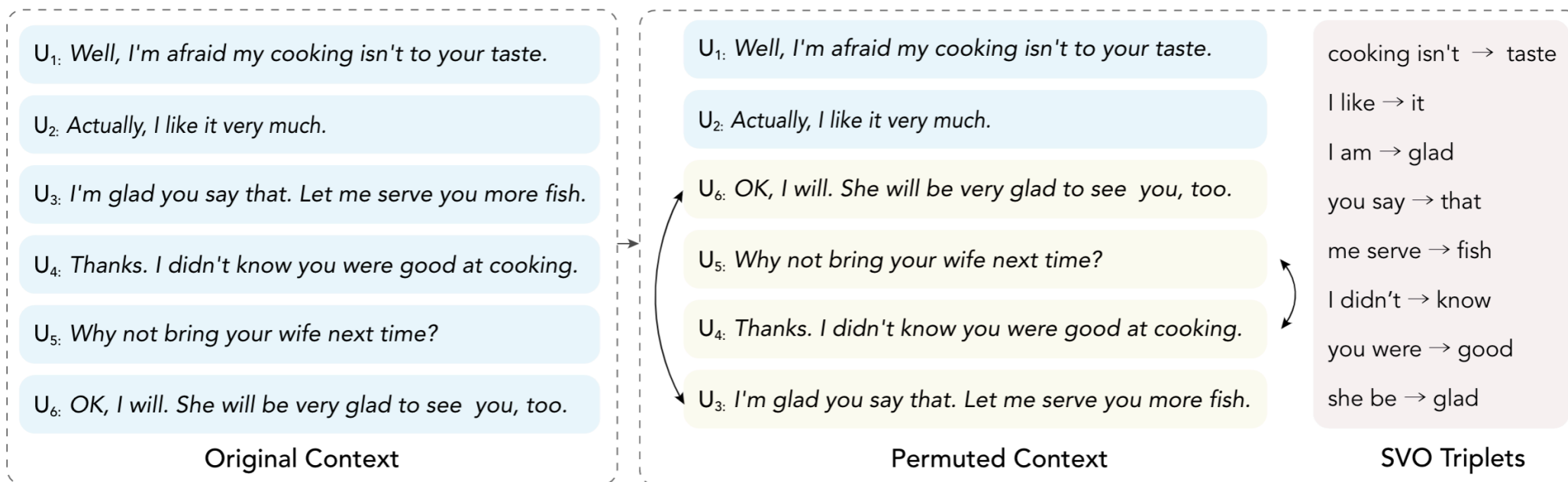


Dialogue-aware Pre-training (SPIDER)

❑ SPIDER: Structural Pre-trained Dialogue Reader

- **sentence backbone regularization:** improve the factual correctness of SVO triples
- **utterance order restoration:** predicts the order of the permuted utterances

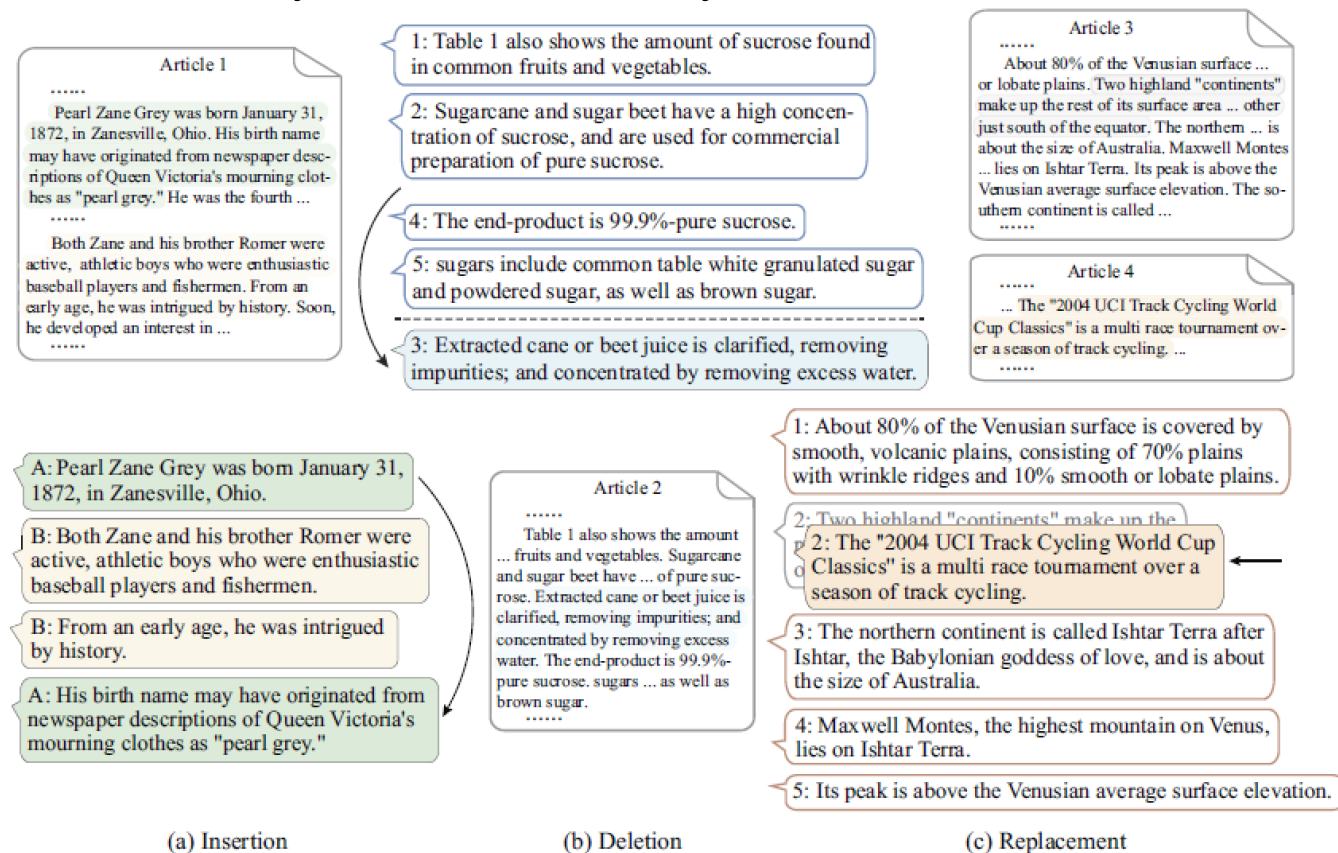
❑ Efficiently and explicitly model the coherence among utterances and the key facts in utterances



Dialogue-oriented Pre-training

- Simulate the conversation features on general plain text to learn dialogue related features including speaker awareness, continuity and consistency:

- Insertion
- Deletion
- Replacement



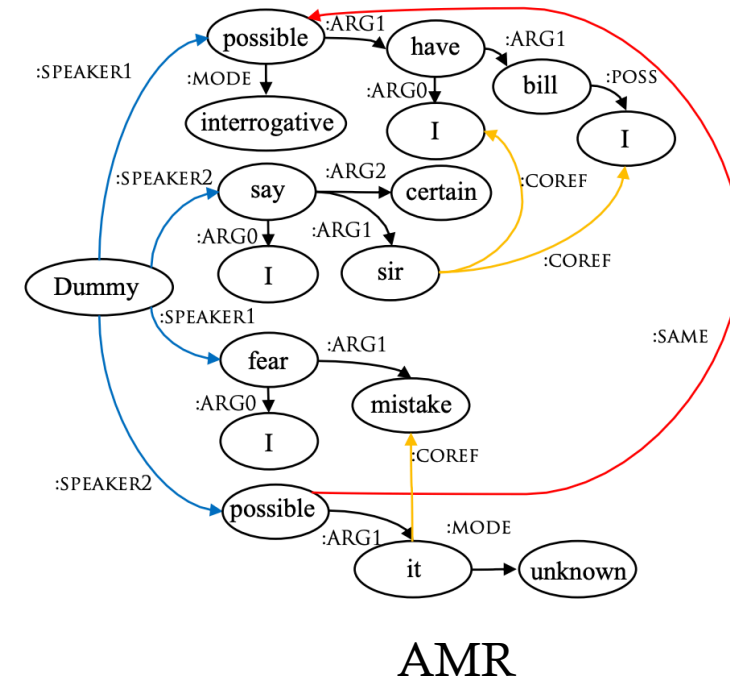
Graph-aware Knowledge Structure Modeling

□ Technical trend: Graph Neural Network (GNN)

- Injecting extra commonsense from knowledge graphs
- Modeling entity relationships
- Graph-attention can be considered as a special case of self-attention

□ Application Scenarios

- Entity linking and coreference modeling
- Dialogue discourse structure
- Abstract meaning representation (AMR)

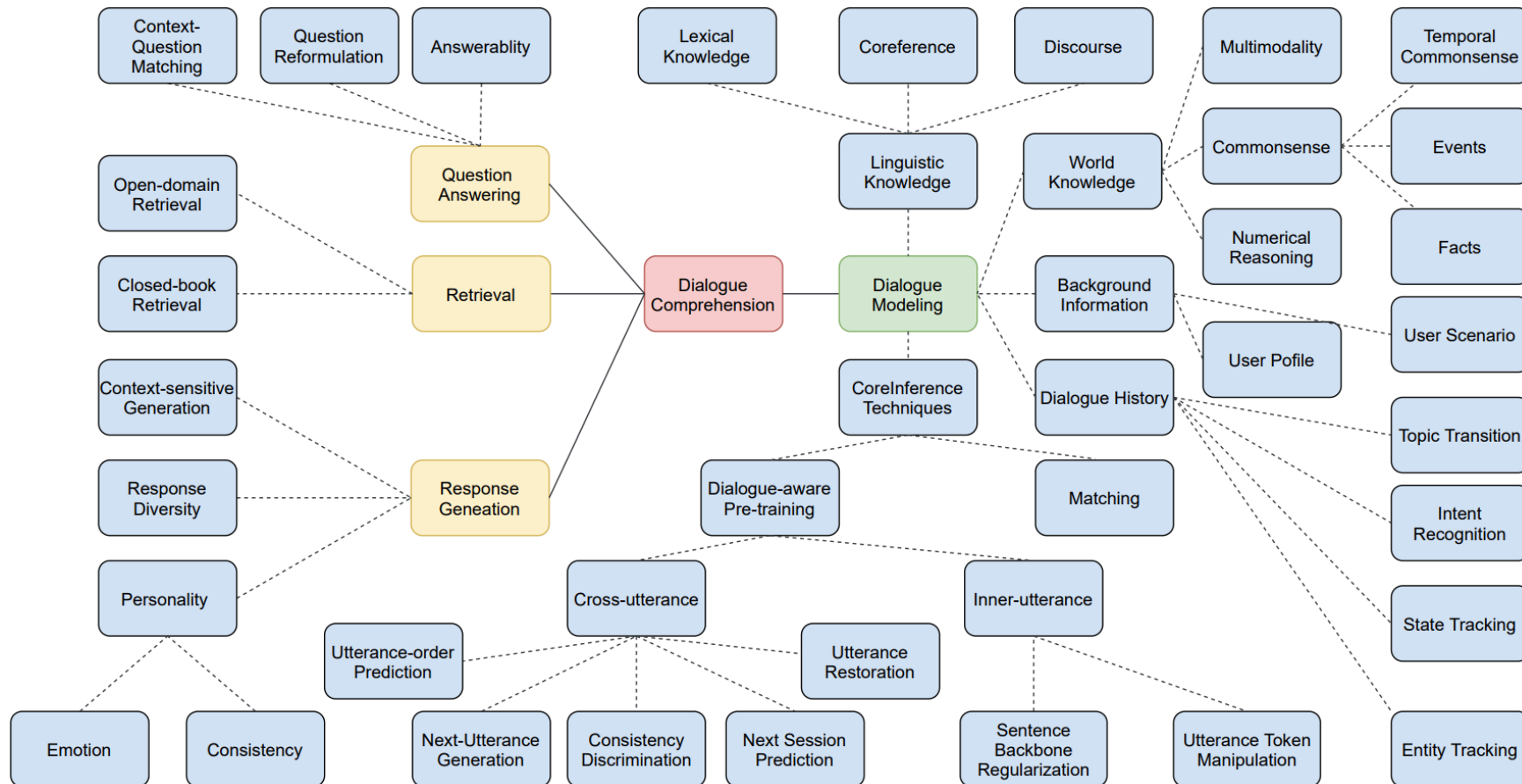


New Frontiers

- ❑ Techniques
 - Domain/Task-adaptive Pre-training
 - Graph-aware Knowledge Structure Modeling
- ❑ Tasks
 - Multi-turn Dialogue Comprehension
 - Logical Reasoning
 - Commonsense Reasoning
- ❑ Applications
 - Open-domain QA
 - Multilingual, Multimodal, Multitask

Multi-turn Dialogue Comprehension

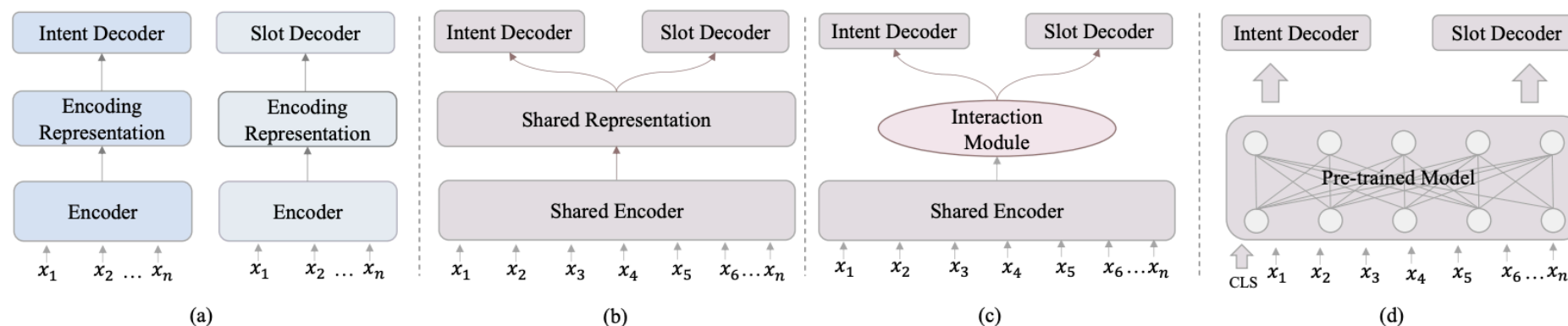
- A multi-turn conversation is intuitively associated with spoken (as opposed to written) language
- Interactive: involves multiple speakers, intentions, topics, thus the utterances are full of transitions.



Other Dialogue Tasks Requiring Comprehension

❑ Spoken Language Understanding

- aims to capture the semantics of user queries
- a core component in task-oriented dialog system



❑ Dialogue Summarization

- Condense the original dialogue into a shorter version covering salient information
- Help people quickly capture the highlights

Fact-driven Logical Reasoning

□ Task: Logical Reasoning

- Challenges: entity-aware commonsense, perception of facts or events.
- Logical supervision is rarely available during language model pre-training.

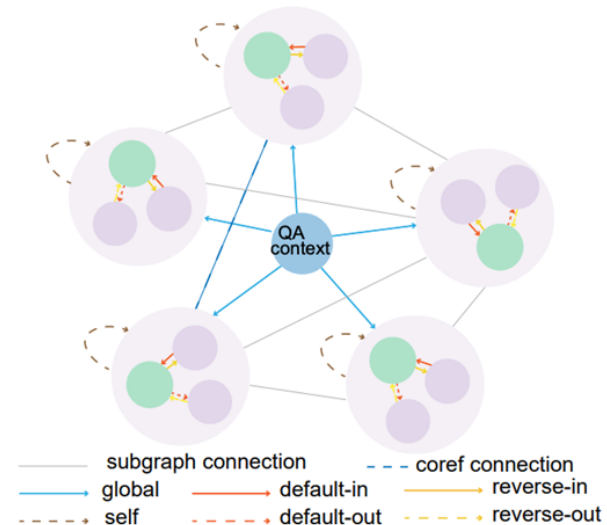
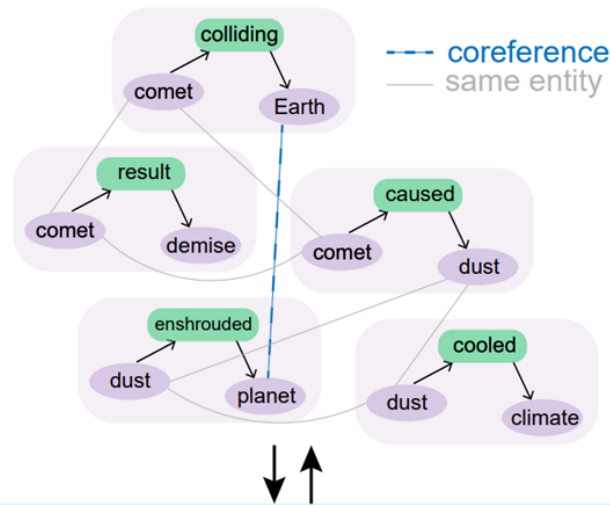
Question	Passage	Answer
<div>Example 1</div> <p>From this we know</p>	Xiao Wang is taller than Xiao Li, Xiao Zhao is taller than Xiao Qian, Xiao Li is shorter than Xiao Sun, and Xiao Sun is shorter than Xiao Qian.	✓ A. Xiao Li is shorter than Xiao Zhao. B. Xiao Wang is taller than Xiao Zhao. C. Xiao Sun is shorter than Xiao Wang. D. Xiao Sun is taller than Xiao Zhao.
<div>Example 2</div> <p>Which one of the following statements, most seriously weakens the argument?</p> A large enough comet colliding with Earth could have caused a cloud of dust that enshrouded the planet and cooled the climate long enough to result in the dinosaurs' demise .	A. Many other animal species from same era did not become extinct at the same time the dinosaurs did. B. It cannot be determined from dinosaur skeletons whether the animals died from the effects of a dust cloud. C. The consequences for vegetation and animals of a comet colliding with Earth are not fully understood. ✓ D. Various species of animals from the same era and similar to them in habitat and physiology did not become extinct.

Fact-driven Logical Reasoning

- ❑ Natural logic units would be the group of backbone constituents of the sentence such as subject, verb and object that cover both global and local knowledge pieces.
- ❑ Design pre-training strategies by restoring facts after masking the units inside a fact and the relations between facts

A large enough comet colliding with Earth could have caused a cloud of dust that enshrouded the planet and cooled the climate long enough to result in the dinosaurs' demise.

comet colliding → Earth
comet caused → dust
dust enshrouded → planet
dust cooled → climate
comet result → demise



Which one of the following, most seriously weakens the argument?

Various species of animals from the same era as dinosaurs and similar to them ... did not become extinct when the dinosaurs did.

Commonsense Reasoning

□ Resources (in natural language)

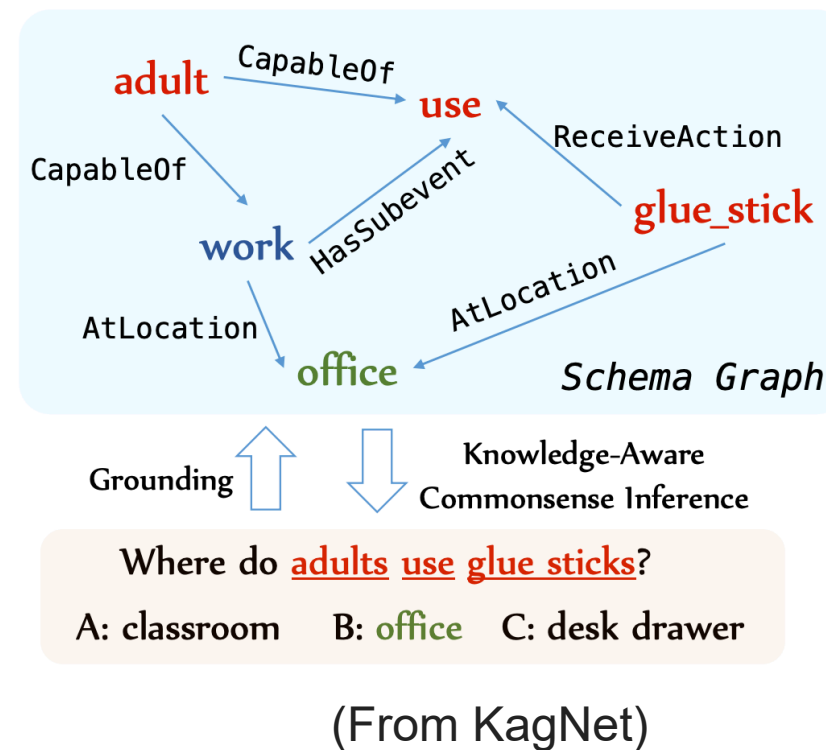
- ConceptNet: semantic knowledge in natural language form
- ATOMIC: knowledge of cause and effect

□ Injecting commonsense into neural networks

- Inserting into the texts
- Attention-based interaction
- Multi-task learning

□ Temporal commonsense

- Understand temporal relations: order, duration, frequency, ..., of events

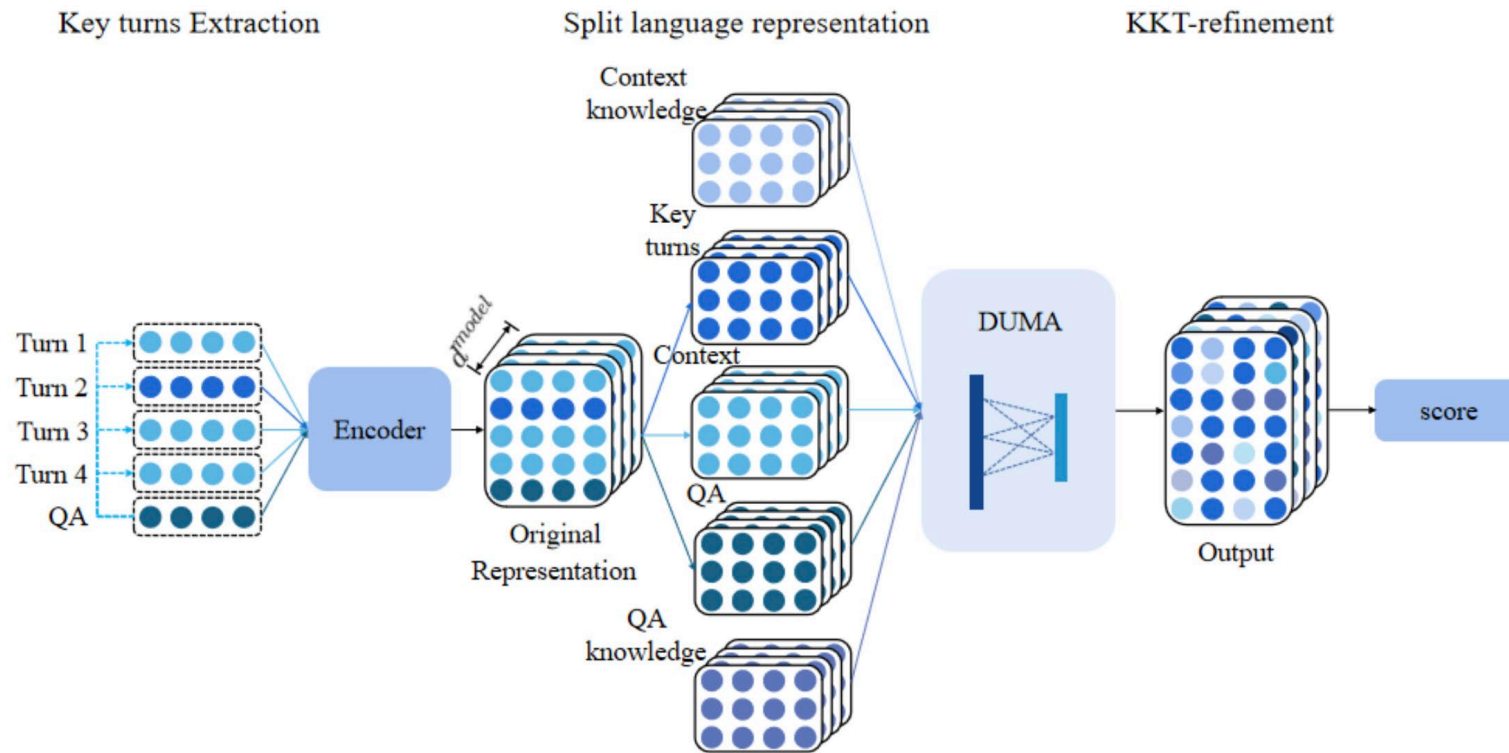


[1] Lin, Bill Yuchen, et al. KagNet: Knowledge-Aware Graph Networks for Commonsense Reasoning. EMNLP 2019.

[2] <https://homes.cs.washington.edu/~msap/acl2020-commonsense/>

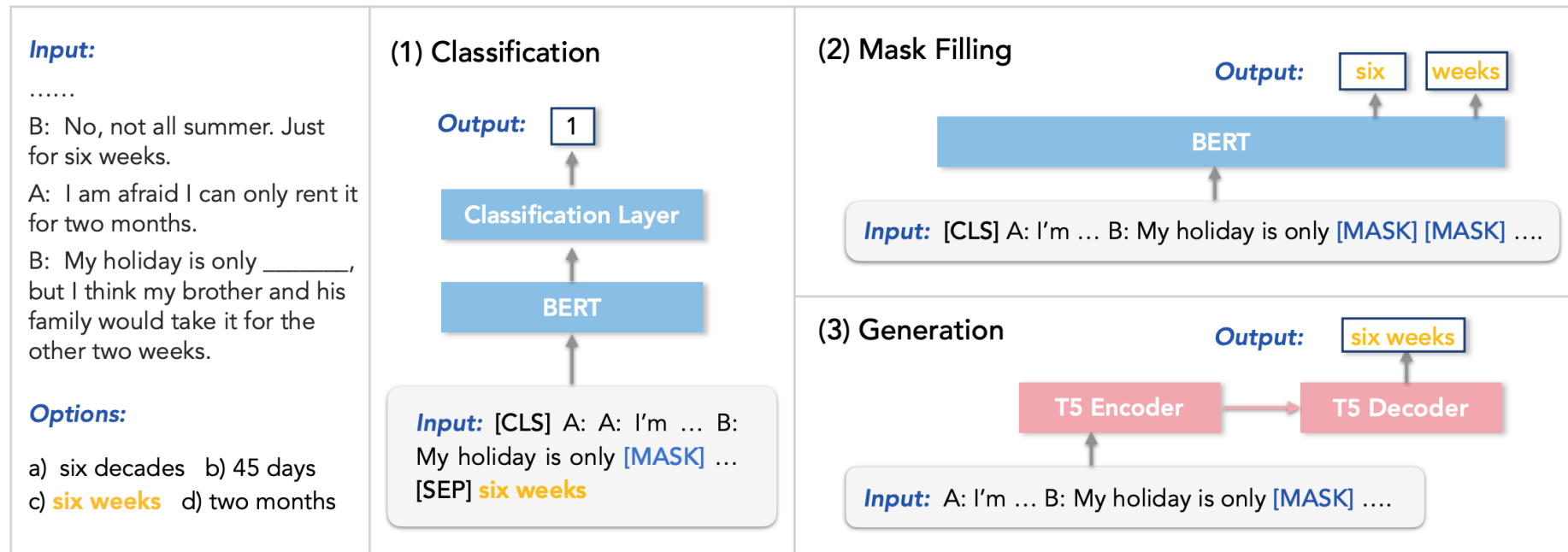
Commonsense Reasoning (KKT)

- ❑ Retrieve Relevant Knowledge from ConceptNet
- ❑ Filter the informative knowledge and use the selected knowledge to enhance the context



Temporal Commonsense

- ❑ Understand temporal relations: order, duration, frequency, ..., of events
- ❑ Humans can easily answer these questions (97.8% accuracy)
- ❑ The best model variant (T5-large with in-domain training) struggles on this challenge set (73%)



New Frontiers

- ❑ Techniques
 - Domain/Task-adaptive Pre-training
 - Graph-aware Knowledge Structure Modeling
- ❑ Tasks
 - Multi-turn Dialogue Comprehension
 - Logical Reasoning
 - Commonsense Reasoning
- ❑ Applications
 - Open-domain QA
 - Multilingual, Multimodal, Multitask

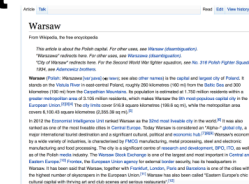
Open-Domain QA

- ❑ Reading Comprehension = Document-level Modeling + QA
- ❑ Open-Domain QA= Open-Domain Reading Comprehension = Open-Domain Document Modeling + QA
 - Machine Reading Comprehension over the whole internet
- ❑ Typical architecture
 - Traditional Retriever-Reader architecture
 - Dense Retrieval vs. BM25
 - Span extraction based on the retrieved documents
- ❑ Next-generation Search Engine

Q: How many of Warsaw's inhabitants spoke Polish in 1933?

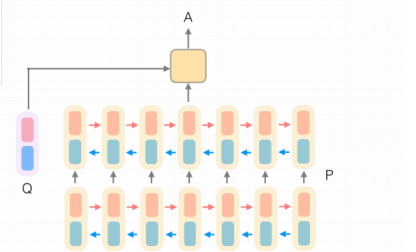


Document Retriever



Document Reader

833,500



Chen, Danqi, et al. 2017. Reading wikipedia to answer open-domain questions. ACL 2017.

Open-Domain QA: DPR

□ Dense Passage **Retriever** (DPR)

- maps any text passage to a fixed dimension of real-valued vectors
- builds an index for all the passages that we will use for retrieval.

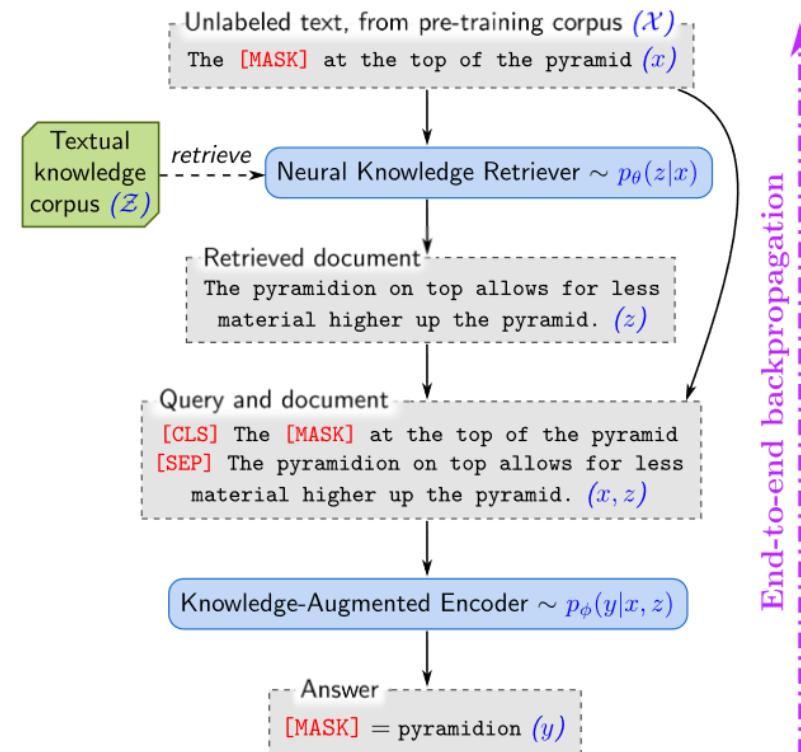
Training	Retriever	Top-20					Top-100				
		NQ	TriviaQA	WQ	TREC	SQuAD	NQ	TriviaQA	WQ	TREC	SQuAD
None	BM25	59.1	66.9	55.0	70.9	68.8	73.7	76.7	71.1	84.1	80.0
Single	DPR	78.4	79.4	73.2	79.8	63.2	85.4	85.0	81.4	89.1	77.2
	BM25 + DPR	76.6	79.8	71.0	85.2	71.5	83.8	84.5	80.5	92.7	81.3
Multi	DPR	79.4	78.8	75.0	89.1	51.6	86.0	84.7	82.9	93.9	67.6
	BM25 + DPR	78.0	79.9	74.7	88.5	66.2	83.9	84.4	82.3	94.1	78.6

Table 2: Top-20 & Top-100 retrieval accuracy on test sets, measured as the percentage of top 20/100 retrieved passages that contain the answer. *Single* and *Multi* denote that our Dense Passage Retriever (DPR) was trained using individual or combined training datasets (all the datasets excluding SQuAD). See text for more details.

Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, Wen-tau Yih. Dense Passage Retrieval for Open-Domain Question Answering. EMNLP 2020.

Open-Domain QA: REALM

- ❑ Two stages: Knowledge Retrieval + Language Modeling
- ❑ Retrieve and attend over documents from a large corpus such as Wikipedia
- ❑ Training Strategies:
 - Only mask “knowledge” tokens (entities, dates, etc.)
 - Add a special empty documents beyond the top-k ones
 - Avoid duplication of pre training documents and knowledge base documents
 - Warmup task: Inverse Cloze Task, retrieve the original document for the sentence



Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, Ming-Wei Chang. REALM: Retrieval-Augmented Language Model Pre-Training. ICML 2020.

Multilingual, Multimodal, Multitask

❑ Multitask

- Training with various types of MRC corpus [1]

❑ Multilingual/Cross-lingual

- Languages other than English are not well-addressed due to the lack of data [2,3]

❑ Multimodal Semantic Grounding

- jointly modeling diverse modalities will be potential research interests [4]
- beneficial for real-world applications, e.g., online shopping and E-commerce customer support
- **Key problem:** 1) the **role** of multimodal features and 2) **when and how** to involve? [5]

[1] MRQA: Workshop on Machine Reading for Question Answering

[2] Cui, Yiming, et al. Cross-Lingual Machine Reading Comprehension. EMNLP 2019.

[3] Anthony Ferritto, Sara Rosenthal, Mihaela Bornea, Kazi Hasan, Rishav Chakravarti, Salim Roukos, Radu Florian, Avirup Sil. A Multilingual Reading Comprehension System for more than 100 Languages. COLING 2020 (Demos).

[4] Hao Tan, Mohit Bansal. Vokenization: Improving Language Understanding with Contextualized, Visual-Grounded Supervision. EMNLP 2020.

[5] Zhuosheng Zhang, Haojie Yu, Hai Zhao, Masao Utiyama. Which Apple Keeps Which Doctor Away? Colorful Word Representations with Visual Oracles. TASLP. 2021.

Conclusion

- ❑ MRC boosts the progress from language **processing** to **understanding**
- ❑ The rapid improvement of MRC systems greatly benefits from the **progress of PrLMs**
- ❑ The theme of MRC is gradually moving from **shallow text matching** to **cognitive reasoning**

Our Survey Papers:

[1] **Machine Reading Comprehension: The Role of Contextualized Language Models and Beyond**

Paper Link: <https://arxiv.org/abs/2005.06249>

[2] **Advances in Multi-turn Dialogue Comprehension: A Survey**

Paper Link: <https://arxiv.org/abs/2103.03125>

Our codes are publicly available at: <https://github.com/cooelf>

Thank You !