

Question Answering with Hybrid Data and Models

Sanjay Kamath Ramachandra Rao^{1,2}

Supervisors: Brigitte Grau¹, Yue Ma²

¹ LIMSI, CNRS, Université Paris-Saclay, Orsay, France

² LRI, Université Paris-Sud, CNRS, Université Paris-Saclay, France

Project: ANR GoASQ

February 6, 2020





Plan

- 1 Introduction
- 2 State of the art
- 3 Building Domain-Specific Models
- 4 Leveraging Semantic Information
- 5 Conclusion

Question Answering

Who was Heisenberg in breaking bad?



Bryan Cranston

It's officially been ten years since Bryan Cranston first graced our television screens as down-and-out chemistry teacher **Walter White** in Breaking Bad. But his metamorphosis into meth kingpin Heisenberg over five epic seasons still hasn't left our minds. The unforgettable AMC drama premiered on Jan. Jan 20, 2018



[Breaking Bad: Walter White Transformation Into Heisenberg ...](https://time.com/breaking-bad-walter-white-transformation)
<https://time.com/breaking-bad-walter-white-transformation>



Walter White

Fictional character

Walter Hartwell White Sr., also known by his clandestine alias Heisenberg, is a fictional character and the main protagonist of Breaking Bad. He is portrayed by Bryan Cranston. [Wikipedia](#)

Played by: [Bryan Cranston](#)

Question Answering - Types of Questions

Q: What country are Volvo automobiles made in? (Location - Country)

A: Sweden

Q: How tall is Mount McKinley? (Numerical value - Height)

A: 6,190 m

Q: What currency is used in Ukraine? (Currency)

A: Ukrainian hryvnia

Q: Who played the role of Heisenberg in the series Breaking Bad? (Person)

A: Bryan Cranston

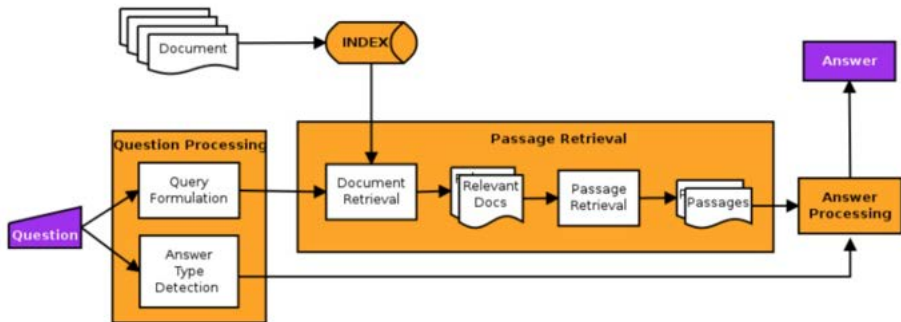
Factoid

Question: Why is ice less dense than water?

Answer Passage: The molecules of water are closer together and constantly moving, whereas the molecules of ice are in a crystal lattice, meaning they're in a rigid formation. When water freezes, the molecules spread out a little more to form the crystal lattice. Since density is mass over volume, and ice has takes up more volume than water, the density of ice is lesser than that of water. Which makes ice float on water.

Non - Factoid

Question Answering Pipeline - The General Approach



- **Question Processing** module analyses questions to detect the Expected Answer Type.
- **Passage retrieval** module uses indexed set of documents to find relevant set of documents and further retrieves a set of relevant paragraphs.
- **Answer Processing** module extracts the answer for the question from the set of relevant paragraphs.

Question Answering Pipeline - Document Retrieval

Question: Who is the President of France?

Document 1: The President of France, officially the President of the French Republic, is the [executive head of state](#) of [France](#) in the [French Fifth Republic](#). In French terms, the presidency is the supreme magistracy of the country.

The powers, functions and duties of prior presidential offices, as well as their relation with the [prime minister](#) and [Government of France](#), have over time differed with the various constitutional documents since the [French Second Republic](#). The president of the French Republic is also the *ex officio* [co-prince of Andorra](#), grand master of the [Legion of Honour](#) and of the [National Order of Merit](#). The officeholder is also honorary proto-canon of the [Basilica of St. John Lateran](#) in Rome (although some have rejected the title in the past). The current president of the French Republic is [Emmanuel Macron](#), who succeeded [François Hollande](#) on 14 May 2017.

Document 2: The presidency of France was first publicly proposed during the [July Revolution](#) of 1830, when it was offered to the [Marquis de Lafayette](#). He demurred in favour of Prince [Louis Philippe](#), who became King of the French. Eighteen years later, during the opening phases of the [Second Republic](#), the title was created for a popularly elected head of state, the first of whom was [Louis-Napoléon Bonaparte](#), nephew of Emperor [Napoleon](#). Bonaparte served in that role until he staged an [auto coup](#) against the republic, proclaiming himself [Napoleon III](#), [Emperor of the French](#).

Document 3: Under the [Third Republic](#) and [Fourth Republic](#), which were [parliamentary systems](#), the office of President of the Republic was a largely ceremonial and powerless one. The Constitution of the [Fifth Republic](#) greatly increased the President's powers. A [1962 referendum](#) changed the constitution, so that the president would be directly elected by universal suffrage and not by the Parliament. In 2000, a [referendum](#) shortened the presidential term from seven years to five years. A maximum of two consecutive terms was imposed after the [2008 constitutional reform](#).

QA Pipeline - Paragraph/Sentence Selection

Question: Who is the President of France?

Paragraph 1: The President of France, officially the President of the French Republic, is the **executive head of state** of France in the **French Fifth Republic**. In French terms, the presidency is the supreme magistracy of the country.

Paragraph 2: The powers, functions and duties of prior presidential offices, as well as their relation with the **prime minister & Government of France**, have over time differed with the various constitutional documents since the **French Second Republic**. The president of the French Republic is also the *ex officio* co-prince of **Andorra**, grand master of the Legion of Honour and of the **National Order of Merit**.

Paragraph 3: The officeholder is also honorary proto-canon of the Basilica of St. John Lateran in Rome (although some have rejected the title in the past).

The current president of the French Republic is Emmanuel Macron, who succeeded François Hollande on 14 May 2017.

Paragraph 4: The presidency of France was first publicly proposed during the **July Revolution** of 1830, when it was offered to the **Marquis de Lafayette**. He demurred in favour of Prince **Louis Philippe**, who became King of the French.

Paragraph 5: Eighteen years later, during the opening phases of the **Second Republic**, the title was created for a popularly elected head of state, the first of whom was **Louis-Napoléon Bonaparte**, nephew of Emperor **Napoleon**. Bonaparte served in that role until he staged an **auto coup** against the republic, proclaiming himself **Napoleon III, Emperor of the French**.

Paragraph 6: Under the **Third Republic** and **Fourth Republic**, which were **parliamentary systems**, the office of President of the Republic was a largely ceremonial and powerless one. The Constitution of the **Fifth Republic** greatly increased the President's powers.

Paragraph 7: A **1962 referendum** changed the constitution, so that the president would be directly elected by universal suffrage and not by the Parliament. In 2000, a **referendum** shortened the presidential term from seven years to five years. A maximum of two consecutive terms was imposed after the **2008 constitutional reform**.

Hurdles using Deep Learning models for QA



- Having enough data (Size)
- Having the right kind of labelled data (Suitable type)
- Building an end-to-end model which does everything (Complexity)
- Generalizing the model to work on all QA tasks (Generalization)

Size and Type of the data

Size

- How large is large enough to use a deep learning algorithm?
- Do we always need a large scale data in a specific domain?
- Can similar datasets from other domains be useful?

Type

- Deep learning based approaches mainly focus on building end-to-end models.
- How can we use semantic features effectively along with neural network models?
- Are synthetic datasets and human annotated datasets comparable?

Complexity and Generalization

Complexity of the model

- Are complex model always performing better than simple ones?
- How to choose a good model to experiment on a new dataset?
- How to choose the required hardware needed for experiments? (Number of GPUs or TPUs required)

Generalization

- Does one model performing better on a dataset perform similarly on others?
- Does it generalize across different data domains?

Timeline of different QA system approaches



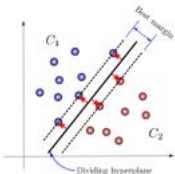
Whole QA



Rule based systems

1960 - 90s

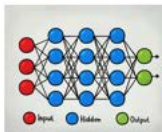
Question Classification, Sentence Selection



Feature based modules

1990-2014

Sentence Selection, Reading Comprehension



Neural Network based modules

2014 - 2018

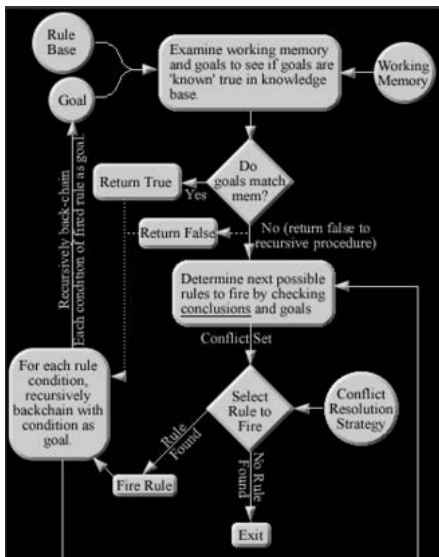
Reading Comprehension



Large Scale Language model based modules

2018 - present

Rule based systems (1960-1990)



Rule based expert systems

- Hard coded rules by experts.
- Term matching module triggers rules and applies actions.
- BASEBALL (1961), LUNAR (1973), SYNTEX, LIFER, and PLANES were some of the systems built.

Rule based expert systems - Limitations

- Hard to create rules.
- Extensive amount of human work is required.
- Systems are not robust.
- Not easy to adapt for expert domains.

Pipeline based systems

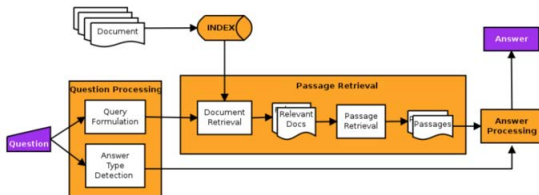
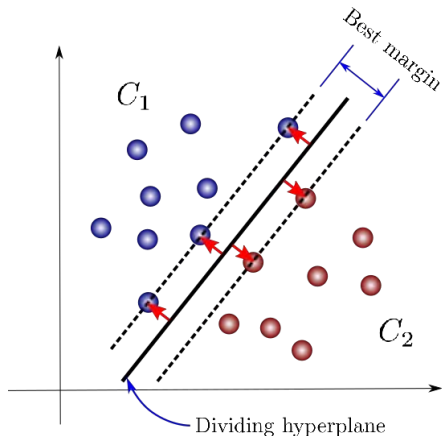


Figure: A typical QA pipeline

Trec QA task

- Trec Question Answering task [Voorhees et al., 2000] since 1999 gave rise to several works which followed pipeline based systems.

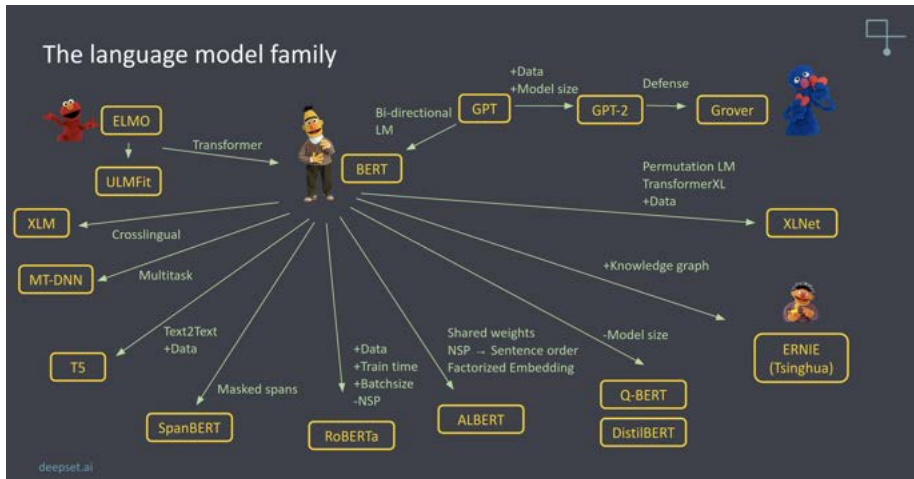
Features based Machine Learning systems (1990-2014)



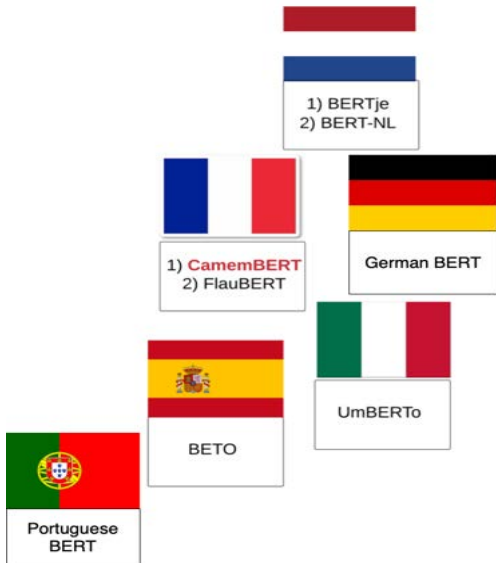
Machine Learning modules - Limitations

- System performance depends mainly on input features.
- Several NLP tools are used for extracting those features.
- Domain expertise is required for feature extraction.
- NLP tools used for pre-processing may contribute to error propagation.

Large Scale Language Models (2018 - Present)



Large Scale Language Models (2018 - Present)



Datasets

	Datasets	Train	Dev	Test
Small Scale	BIOASQ 4b	427	59	161
	BIOASQ 5b	544	75	150
	BIOASQ 6b	685	94	161
Large Scale	SQUAD v1.0	87,599	10,570	9,533
	QUASAR-T	37,012	3,000	3,000

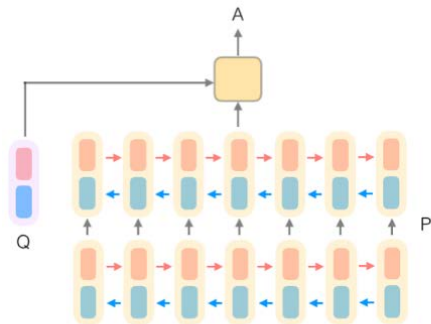
Figure: Large scale and Small scale datasets comparison

Datasets

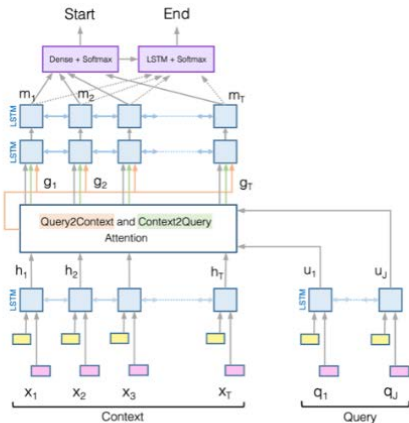
Dataset	Gold Questions	RC eligible
Bioasq 4	486	321 (66%)
Bioasq 5	619	428 (69.1%)
Bioasq 6	779	543 (69.7%)
SQUAD v1.0	1,07,702	1,07,702 (100%)

Figure: Datasets used in our experiments which are suitable for Reading Comprehension (RC) setting, as done by [Weissenborn et al., 2018, Lee et al., 2019]

Choosing a Simple Model



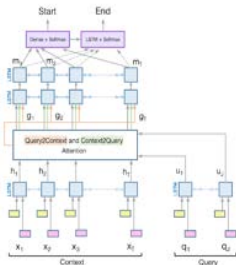
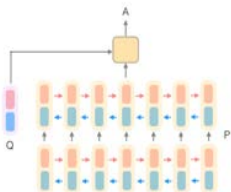
DRQA model



BIDAf model

Figure: DRQA by [Chen et al., 2017] and BIDAf by [Seo et al., 2016]

Choosing a Simple Model



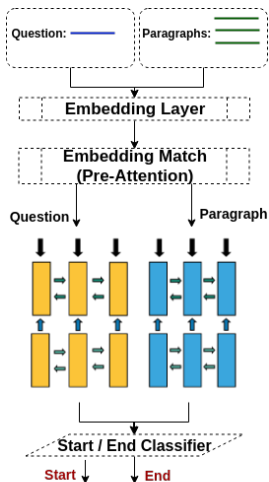
DRQA

- Training time: ~4 hours on a single GPU
- Exact Match score on SQUAD dataset: 69.5%
- Simple model compared to BiDAF
- Published in March 2017

BiDAF

- Training time: ~20 hours a single GPU
- Exact Match score on SQUAD dataset: 67.7%
- Complex model compared to DRQA
- Published in Nov. 2016

DRQA



DRQA

- Question (Q) and Paragraphs (P) both are encoded with GLOVE vectors.
- An attention mechanism is used to map embeddings between Q and P.

$$F_{align}(p_i) = \sum_j a_{i,j} E(q_j) \quad (1)$$

Where $a_{i,j}$ is,

$$a_{i,j} = \frac{\exp(\alpha(E(s_i)) \cdot \alpha(E(q_j)))}{\sum_{j'} \exp(\alpha(E(s_i)) \cdot \alpha(E(q_{j'})))} \quad (2)$$

- Bi-LSTMS are used individually and are connected to two separate bilinear classifiers for Start and End predictions.

$$P_{start}(i) \propto \exp(\mathbf{p}_i \mathbf{W}_s \mathbf{q}) \quad (3)$$

$$P_{end}(i) \propto \exp(\mathbf{p}_i \mathbf{W}_e \mathbf{q}) \quad (4)$$

- Model and implementation by [Chen et al., 2017]

Domain Adaptation using DRQA

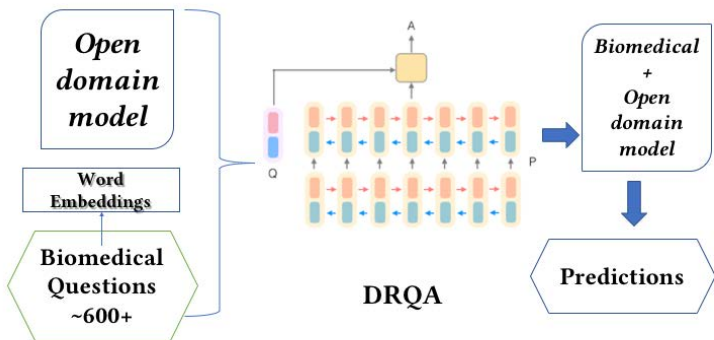
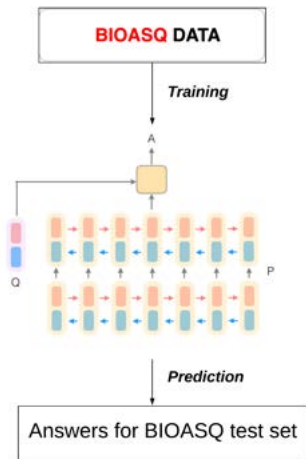
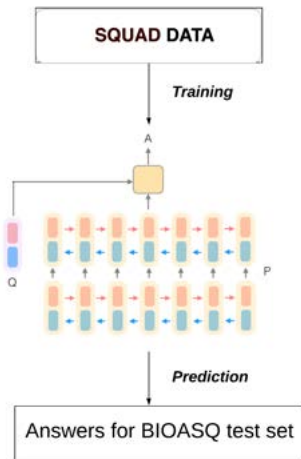


Figure: From open domain towards biomedical domain [Weissenborn et al., 2017]

Pre-training and Finetuning

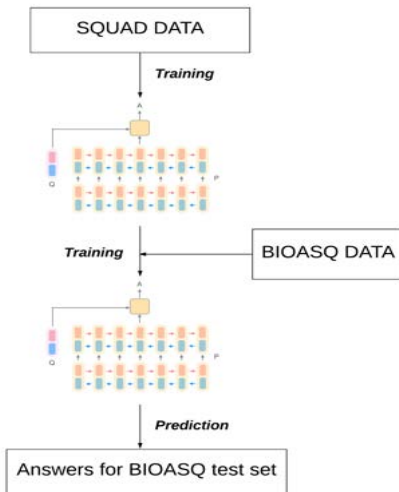


NO PRE-TRAINING



NO FINE-TUNING

Pre-training and Finetuning



PRE-TRAINING + FINE-TUNING

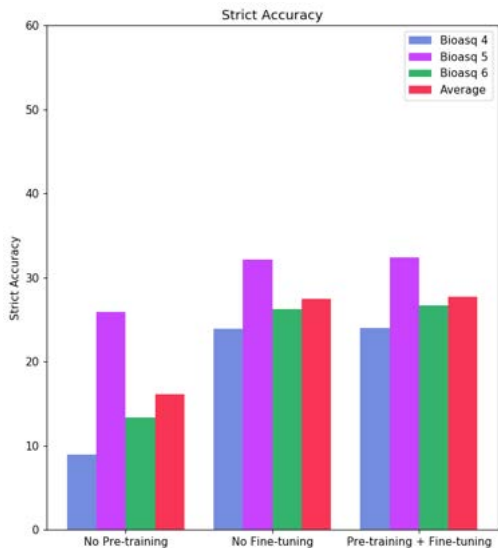


Figure: A study to show the importance of domain adaptation - [Kamath et al., 2019]

Reading Comprehension vs Open QA

Question: Which NFL team represented the AFC at Super Bowl 50?

Answer Passage: Super Bowl 50 was an American football game to determine the champion of the National Football League (NFL) for the 2015 season. The American Football Conference (AFC) champion **Denver Broncos** defeated the National Football Conference (NFC) champion Carolina Panthers 24-10 to earn their third Super Bowl title. The game was played on February 7, 2016, at Levi's Stadium in the San Francisco Bay Area at Santa Clara, California. As this was the 50th Super Bowl, the league emphasized the "golden anniversary" with various gold-themed initiatives, as well as temporarily suspending the tradition of naming each Super Bowl game with Roman numerals (under which the game would have been known as "Super Bowl L"), so that the logo could prominently feature the Arabic numerals 50.

Answer: **Denver Broncos**. Start Offset: 177

Q: What's the capital of Ireland?

P₁: As the capital of Ireland, **Dublin** is...

P₂: Ireland is an island in the North Atlantic. . .

P₃: **Dublin** is the capital of Ireland. Besides, Ottawa is one of famous tourist cities in Ireland and ...

BIOASQ Data (Recap)

Question: Which calcium channels does ethosuximide target?

Answer: **T-type calcium channels**

Snippets:

- Theta rhythms remained disrupted during a subsequent week of withdrawal but were restored with the **T-type channel** blocker ethosuximide.
- Given evidence that chemotherapy-induced neuropathic pain is blocked by ethosuximide, known to block **T-type calcium channels**, we examined if more selective T-type calcium channel blockers....
- The Ca(v)3.2 channel is sensitive to ethosuximide, amlodipine and amiloride.

Figure: Some snippets contain answers and some do not.

Reading Comprehension Format

- **All paragraphs are considered in the OpenQA model setting.**

OPEN-QA model - PSPR

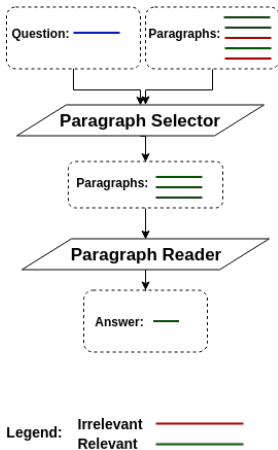


Figure: OpenQA model by [Lin et al., 2018]

OPEN-QA with DRQA

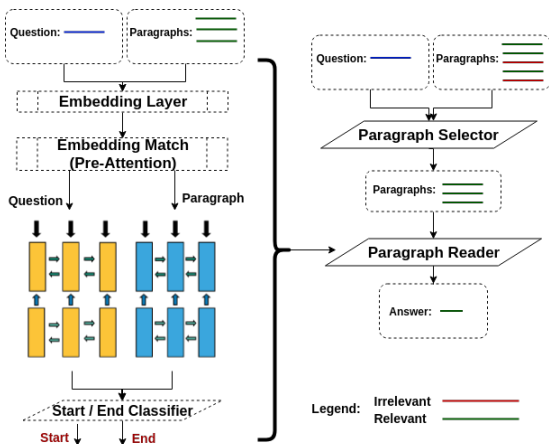


Figure: OpenQA model by [Lin et al., 2018]

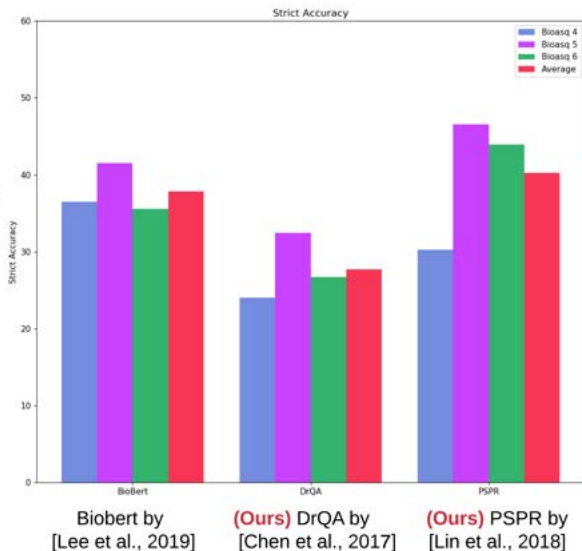
Training Times (Approx):

Biobert:

Bert - 4 days on 1 TPU pod,
Biobert - 23 days on 8 GPUs
and 20 minutes on 1 GPU
(Fine-tuning)

DrQA - 4 hours on 1 GPU

PSPR - 8 hours on 1 GPU



BERT - New SOTA on several QA tasks



BERT - New SOTA on several QA tasks

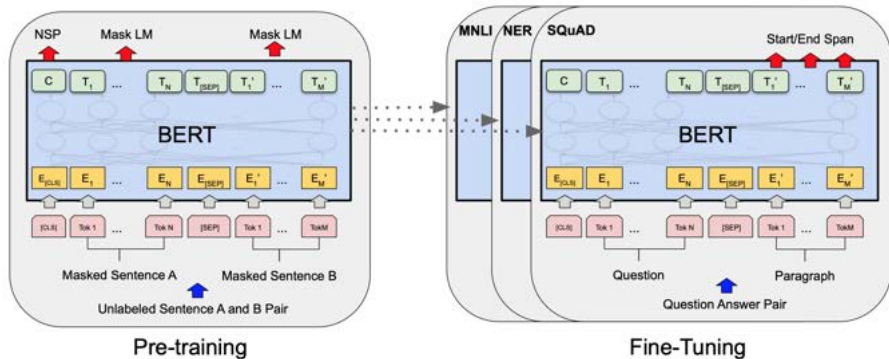


Figure: BERT model by [Devlin et al., 2019]

BIOBERT

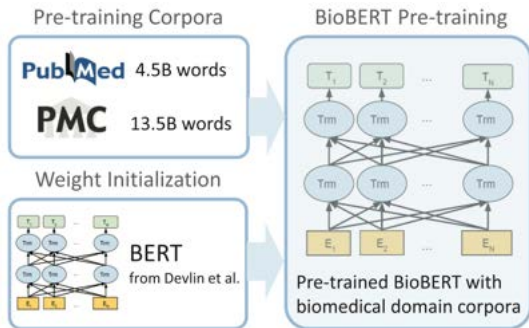


Figure: BIOBERT model by [Lee et al., 2019]

BIOBERT in BIOASQ

- Applied BIOBERT on BIOASQ task data but on document level text snippets.
- Pre-trained with SQUAD data and fine-tuned with BIOASQ data.

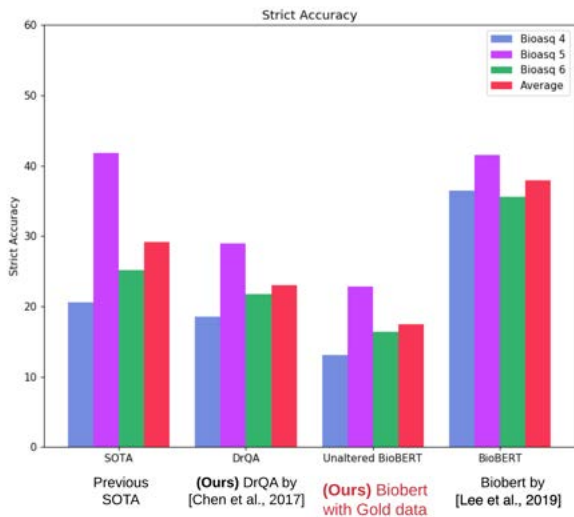


Figure: Modification of the paragraph text results in the variation of performance

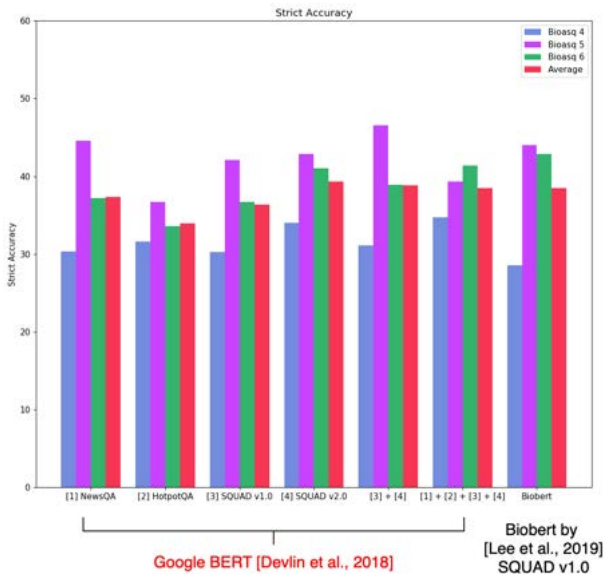
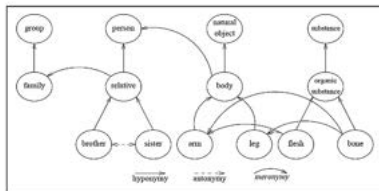
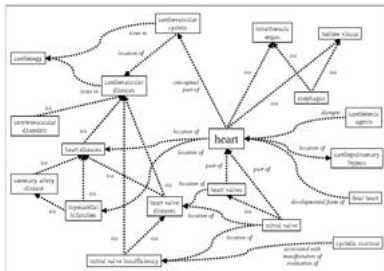


Figure: Experiments on different Pre-training datasets for Biomedical QA task.

Semantic Information

- End-to-end Neural models rely only on input and output to learn.
- Traditional QA pipeline methods rely on features from different sources such as named entities, part of speech tags, question types etc.
- "How can one build models which use best of the both worlds?"



Using Answer Variants

Q: Mutation of which gene is implicated in the Brain-lung-thyroid syndrome?

- Novel **NKX2-1** Frameshift Mutations in Patients with Atypical Phenotypes of the Brain-Lung-Thyroid Syndrome.
- **NKX2-1** mutations in brain-lung-thyroid syndrome: a case series of four patients.
- Brain-lung-thyroid syndrome (BLTS) characterized by congenital hypothyroidism, respiratory distress syndrome, and benign hereditary chorea is caused by **thyroid transcription factor 1 (NKX2-1/TTF1)** mutations.
- The disorder is caused by mutations to the **NKX2.1 (TTF1)** gene and also forms part of the "brain-lung-thyroid syndrome", in which additional developmental abnormalities of lung and thyroid tissue are observed.

A: thyroid transcription factor 1, NKX2-1, TTF1, TTF1

Several answer variants which are syntactically different but semantically represents the same entity are annotated by us.

Annotating Answer Variants

1	Question: Which species of bacteria did the mitochondria originate from?
2	Answer: [u'Biologists agree that the ancestor of mitochondria was an alpha-proteobacterium.']
	Begin
3	Recently, α -proteobacteria have been shown to possess virus-like gene transfer agents that facilitate high frequency gene transfer in natural environments betw
4	This system could have driven the genomic integration of the mitochondrial progenitor and its proto-eukaryote host and contributed to the evolutionary mosaic of ge
	eukaryotic genomes.
	Begin ExactAnswer
5	Although the Alphaproteobacteria are thought to be the closest relatives of the mitochondrial progenitor, there is dispute as to what its particular sister group is
	Begin
6	More detailed phylogenetic analyses with additional Alphaproteobacteria and including genes from the mitochondria of Reclinomonas americana found matches
	the Rickettsiaceae, Anaplasmataceae, and Rhodospirillaceae families.
	Begin ExactAnswer
7	Biologists agree that the ancestor of mitochondria was an alpha-proteobacterium.
	Begin
8	Mitochondria originated by permanent enslavement of purple non-sulphur bacteria.
	Begin ExactAnswer
9	Phylogenetic analyses based on genes located in the mitochondrial genome indicate that these genes originated from within the alpha-proteobacteria.

Manual Annotations

- Annotated using the Brat tool and UMLS meta thesaurus references. 618 questions were annotated manually.
- 3 people from CS background annotated for answer variants.
- Released this annotated dataset publicly

https://zenodo.org/record/1346193#.W3_WUZMzZQI

Annotating Answer Variants

Gold standard answer : MDR - TB

Paragraph 1: Delamanid: a review of its use in patients with multidrug-resistant **tuberculosis**.

Paragraph 2: In conclusion, delamanid is a useful addition to the treatment options currently available for patients with **MDR-TB**.

Paragraph 3:

EXPERT OPINION: Delamanid showed potent activity against drug-susceptible and -resistant Mycobacterium **tuberculosis** in both in vitro and in vivo studies.

Automatic Annotations

- Annotated using the UMLS CUI identifiers from entities detected.
- CUI from gold standard answer and matching CUIs from paragraphs are mapped to find variants.

Annotating Answer Variants

Gold standard answer : **MDR - TB**

Paragraph 1: Delamanid (CUI : C0206526) is used in patients with multidrug-resistant tuberculosis.

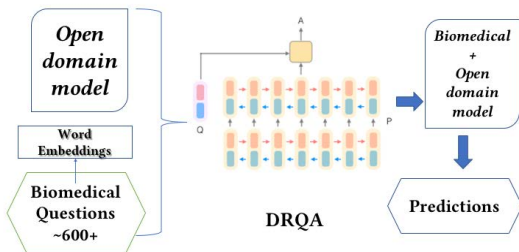
Conclusion, delamanid is a useful addition to the treatment options for patients with MDR-TB. (CUI : C0206526)

Paragraph 3: EXPERT OPINION: Delamanid (CUI : C0206526) has activity against drug-susceptible and -resistant Mycobacterium tuberculosis in both in vitro and in vivo studies. (CUI : C0206526)

Automatic Annotations

- Annotated using the UMLS CUI identifiers from entities detected.
- CUI from gold standard answer and matching CUIs from paragraphs are mapped to find variants.

Answer Variants and Performance Increase



Experiments

- Reading comprehension experiments with BIOASQ Gold Standard data and Annotated data.
- Comparison with Automatic and Manually annotated answers.
- SQUAD dataset for Open domain, BIOASQ dataset for biomedical domain.

Approach 2 - Expected Answer Types (EAT)

Biomedical Domain

Question: What **disease** in Loxapine prominently used for?

Answer: **Schizophrenia**

Expected Answer Type: Disease or Syndrome.

Semantic Group: DISO Disease or Syndrome.

Lexical Answer Type: **disease**

Question: Which **drugs** are utilized to treat amiodarone-induced thyrotoxicosis?

Answer: **Antithyroid drugs**

Expected Answer Type: Chemical Drug

Semantic Group: CHEM Chemicals & Drugs

Lexical Answer Type: **drugs**

Open Domain

Question: Which **actor** played the role of Walter White in the series Breaking Bad?

Answer: **Bryan Cranston**

Expected Answer Type (EAT): HUM - Person, Organisation

Lexical Answer Type (LAT): **actor**

Question: **Where** was the TV show Breaking Bad primarily filmed at?

Answer: **Albuquerque**

Expected Answer Type (EAT): LOC - Location

Lexical Answer Type (LAT): **Where**

Improving Answer Sentence Selection - Using EAT (Expected Answer Type)

-	Method	Question	Sentence
1	Original text	Who is the author of the book, 'The Iron Lady: a biography of Margaret Thatcher'	in 'The Iron Lady,' <i>Young</i> traces the greatest woman political leader since <i>Catherine the Great</i> .
2	Replacement - (Tayyar Madabushi et al., 2018) (EAT Single type)	Who is the author of the book, 'The Iron Lady: a biography of Margaret Thatcher'	in 'The Iron Lady,' <code>max_entity_left</code> traces the greatest woman political leader since <code>entity_left</code> .
3	EAT (Different types)	Who is the author of the book, 'The Iron Lady: a biography of Margaret Thatcher'	in 'The Iron Lady,' <code>max_entity_left</code> traces the greatest woman political leader since <code>entity_hum</code> .
4	EAT (MAX + Different types)	Who is the author of the book, 'The Iron Lady: a biography of Margaret Thatcher'	in 'The Iron Lady,' <code>max_entity_hum</code> traces the greatest woman political leader since <code>entity_hum</code> .

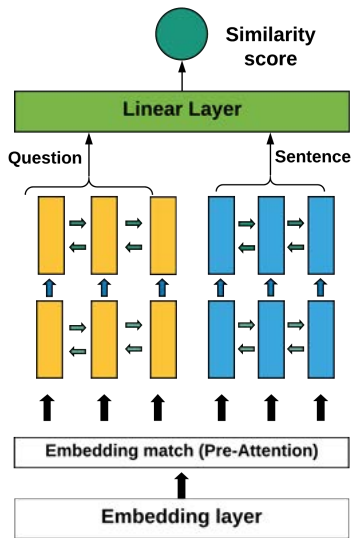


Figure: RNN-S model for Sentence Selection Task

Improving Answer Sentence Selection - Using EAT (Expected Answer Type)

Datasets	Method	Acc.@1	MAP	MRR
TrecQA	Plain words - (Rao et al., 2016)	-	78	83.4
	EAT words - (Tayyar Madabushi et al., 2018)	-	83.6	86.2
	Plain words - RNN-S	78.95	80.24	84.81
	EAT words (single type) - RNN-S	85.26	85.28	89.16
	EAT words (different types) - RNN-S	85.26	85.48	88.11
	EAT words (MAX+different types) - RNN-S	86.32	85.42	88.86

Figure: Results using RNN-S model. EAT (Expected Answer Type) - [Kamath et al., 2019]

Improving Top-1 Accuracy Using Semantic Features

Algorithm	Accuracy
PSPR - Baseline	41.1
Randomforest	42.86
Adaboost	43.23
MLP classifier	42.23

Features removed from the best model	Accuracy	Drop
Best model score	43.23	-
- Maximum value of paragraph probability + Answer presence ratio	41.36	1.87
- Maximum value of paragraph probability	42.23	1
- Summation of Ans Prob. + Answer presence ratio	42.46	0.77
- Answer presence ratio	42.53	0.7
- Summation of answer probability	43.03	0.2
- Answer Probability	43.06	0.17
- Sentence Probability	43.2	0.03

Figure: Results on the QUASAR-T dataset (Open Domain)

Improving Top-1 Accuracy Using Semantic Features

Datasets	Finetune	BIOASQ 4	BIOASQ 5	BIOASQ 6
Baseline	Strict	30.31	46.83	42.79
	Lenient	45.00	52.66	53.41
Adaboost	Strict	39.37	44.00	45.96
	Lenient	45.00	52.66	53.41
Randomforest	Strict	38.75	46.00	46.58
	Lenient	45.00	52.66	53.41
MLP	Strict	34.37	46.00	38.50
	Lenient	45.00	52.66	53.41

Features removed from the best model	Bioasq 5 Acc.	Bioasq 6 Acc.
Best model score	46.00	46.58
- Maximum value of answer probability + Answer overlap	42.66 (3.34)	44.72 (1.86)
- Maximum value of answer probability	43.33 (2.67)	44.72 (1.86)
- Paragraph Probability	44.66 (1.34)	45.96 (0.62)
- LAT Semantic Type	45.33 (0.67)	45.96 (0.62)
- Answer overlap	45.33 (0.67)	45.96 (0.62)

Figure: Results on the BIOASQ dataset (Biomedical Domain)

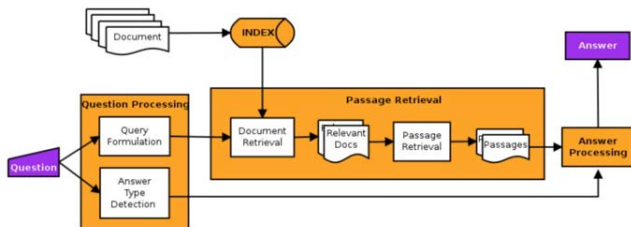




Plan

- 1 Introduction
- 2 State of the art
- 3 Building Domain-Specific Models
- 4 Leveraging Semantic Information
- 5 Conclusion

Whole QA pipeline using different QA models



QA Pipeline using SOTA models

- Can the current State of the Art models replace the pipeline with an end-to-end model?
- How does the performance compare with neural models which use pipeline architecture?

Whole QA pipeline using SOTA models

-	Model	Accuracy
1	SQUAD using DRQA (Chen et al., 2017)	69.5
2	Open QA (BM25) and BERT model (K. Lee et al., 2019)	28.1
3	Open QA using LSTM model (Chen et al., 2017)	27.1
4	Whole QA using BERT model (K. Lee et al., 2019)	26.5

QA Pipeline using SOTA models

- BERT end-to-end model for whole QA pipeline performs **worse** compared to LSTM model DRQA which is end-to-end.
- BM25 for document retrieval and BERT for QA performs better than BERT alone.

Whole QA Pipeline on TrecQA data

Model	Accuracy
Tree kernels methods by (Severyn and Moschitti, 2013)	70.8
Sequence tagging model by (Yao et al., 2013)	67.2
Reading Comprehension mode	58.02
OpenQA pre-trained with QUASAR-T	49.38
OpenQA on TrecQA data only	45.68

QA Pipeline on TrecQA data

- Using the OpenQA model by [Lin et al., 2018] on TrecQA task dataset for Answer Extraction.
- Non-neural models performed better.

Large Scale Language Models - Concerns

Consumption	CO ₂ e (lbs)
Air travel, 1 passenger, NY↔SF	1984
Human life, avg. 1 year	11,023
American life, avg. 1 year	36,156
Car, avg incl. fuel, 1 lifetime	126,000
Training one model (GPU)	
NLP pipeline (parsing, SRL)	39
w/ tuning & experimentation	78,468
Transformer (big)	192
w/ neural architecture search	626,155

Model	Hardware	Power (W)	Hours	kWh-PUE	CO ₂ e	Cloud compute cost
Transformer _{base}	P100x8	1415.78	12	27	26	\$41-\$140
Transformer _{big}	P100x8	1515.43	84	201	192	\$289-\$981
ELMo	P100x3	517.66	336	275	262	\$433-\$1472
BERT _{base}	V100x64	12,041.51	79	1507	1438	\$3751-\$12,571
BERT _{base}	TPUv2x16	—	96	—	—	\$2074-\$6912
NAS	P100x8	1515.43	274,120	656,347	626,155	\$942,973-\$3,201,722
NAS	TPUv2x1	—	32,623	—	—	\$44,055-\$146,848
GPT-2	TPUv3x32	—	168	—	—	\$12,902-\$43,008

Table 1: Estimated CO₂ emissions from training common NLP models, compared to familiar consumption.¹

Some concerns

- Expensive compute required. BERT took 4 days on 16 cloud TPUv3 (64 chips) which costs around \$14,000.
- CO₂ emissions are a side effect which are often neglected into consideration. [Strubell et al., 2019]

Large Scale Language Models - Concerns



Conversation B	
Human:	<i>Hi!</i>
Meena:	<i>Hi! How are you this fine evening?</i>
Human:	<i>I'm doing well. it's morning my time!</i>
Meena:	<i>Good morning! What are your plans for today?</i>
Human:	<i>Get my work done then watch a movie at home in the night</i>
Meena:	<i>Sounds fun! What movie?</i>
Human:	<i>I haven't decided yet, perhaps an Oscar-nominated one. any suggestion?</i>
Meena:	<i>I'm not too familiar with recent Oscar winners, but I do know I wanted to see The Grand Budapest Hotel.</i>
Human:	<i>ah, tell me about that movie!</i>
Meena:	<i>It's from the same director and writer as Moonrise Kingdom. Wes Anderson is pretty well known for making funny, colorful films. I've heard it's pretty good.</i>
Human:	<i>ha! I've searched on the internet and both movies are very colorful!</i>
Meena:	<i>I'm looking forward to it. I love colorful movies.</i>

Some concerns

- A new chatbot model by Google AI [Adiwardana et al., 2020] costs \$1,400,000 to train for 30 days on the TPU cloud.
- Google uses renewable energy sources for cloud, what about others?

Conclusion

- Domain Adaptation
 - From Open domain towards Biomedical domain QA using several techniques.
- Semantic information
 - Explicit use of semantic and structured information can help.
- Choose simple models most of the times.
 - Simple models are better for several reasons.
- Contributions
 - Several SOTA QA models built, modified and experimented.
 - Annotated datasets and codes released publicly.



Publications

- **2019 (Co-authored publication) - Measuring semantic similarity of clinical trial outcomes using deep pre-trained language representations** - Anna Koroleva, Sanjay Kamath, Patrick Paroubek. Journal of Biomedical Informatics: X, Published in October 2019.
- **2019 - How to Pre-Train Your Model? Comparison of Different Pre-Training Models for Biomedical Question Answering.** - Proceedings of the 7th BioASQ Workshop A challenge on large-scale biomedical semantic indexing and question answering. ECMLPKDD, September 2019.
- **2019 - Predicting and Integrating Expected Answer Types into a Simple Recurrent Neural Network Model for Answer Sentence Selection.** - 20th International Conference on Computational Linguistics and Intelligent Text Processing - CICLING 2019, April 2019.
- **2018 - An Adaption of BIOASQ Question Answering dataset for Machine Reading systems by Manual Annotations of Answer Spans.** - Proceedings of the 6th BioASQ Workshop A challenge on large-scale biomedical semantic indexing and question answering. EMNLP, October 2018.
- **2018 - Verification of the Expected Answer Type for Biomedical Question Answering.** - HQA workshop, companion proceedings of the The Web Conference 2018, April 2018.
- **2017 - A Study of Word Embeddings for Biomedical Question Answering.** - 4e édition du Symposium sur l'Ingénierie de l'Information Médicale, November 2017.