



# **SQUARE: TOWARDS MULTI-DOMAIN AND FEW-SHOT COLLABORATING QUESTION ANSWERING AGENTS**

Prof. Iryna Gurevych

Haritz Puerto

# AGENDA

- 1** New Possibilities
- 2** Meta-QA
- 3** Relevance Feedback In Neural Re-Ranking
- 4** UKP-SQuARE
- 5** Future Work



**SQUARE: TOWARDS MULTI-DOMAIN AND FEW-SHOT COLLABORATING QA AGENTS**

# NEW POSSIBILITIES

# EXPLOSION OF QA DATASETS AND MODELS

How to leverage all this collective effort?

How to study all these models and datasets?

 **Hugging Face** Models 3,942

distilbert-base-cased-distilled-squad  
📄 • Updated Dec 5, 2022 • ↓ 2.36M • ❤️ 86

Datasets 380

 squad

👁️ Preview • Updated Nov 3, 2022 • ↓ 174k • ❤️ 60

# MULTI-SKILL QA

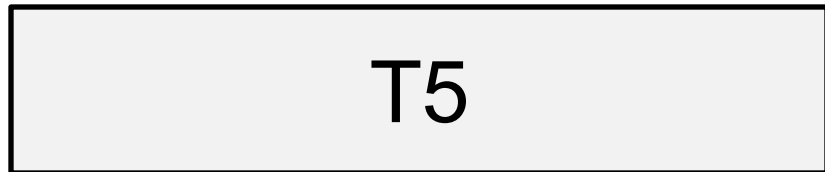
## MULTI-DATASET MODELS

Train a model on many datasets



Extractive QA Datasets

[MultiQA](#) (Talmor & Berant, ACL 2019)

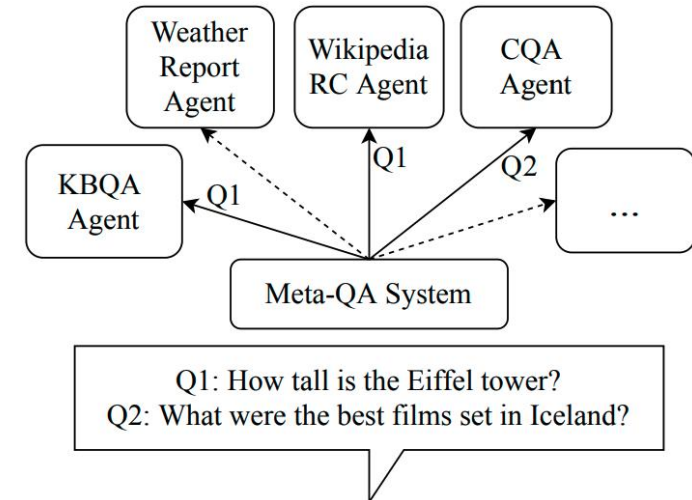


Extractive, Multiple Choice,  
Abstractive, Boolean QA Datasets

[UNIFIEDQA](#) (Khashabi et al., Findings 2020)

## MULTI-AGENT MODELS

Combine many models

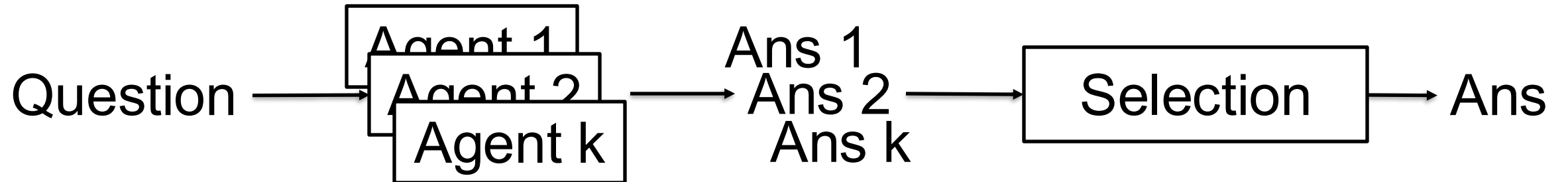


[TWEAC: Transformer with Extendable QA Agent Classifiers](#) (Geigle et al., Arxiv 2021)

**METAQA: COMBINING EXPERT AGENTS FOR MULTI-SKILL QUESTION ANSWERING**

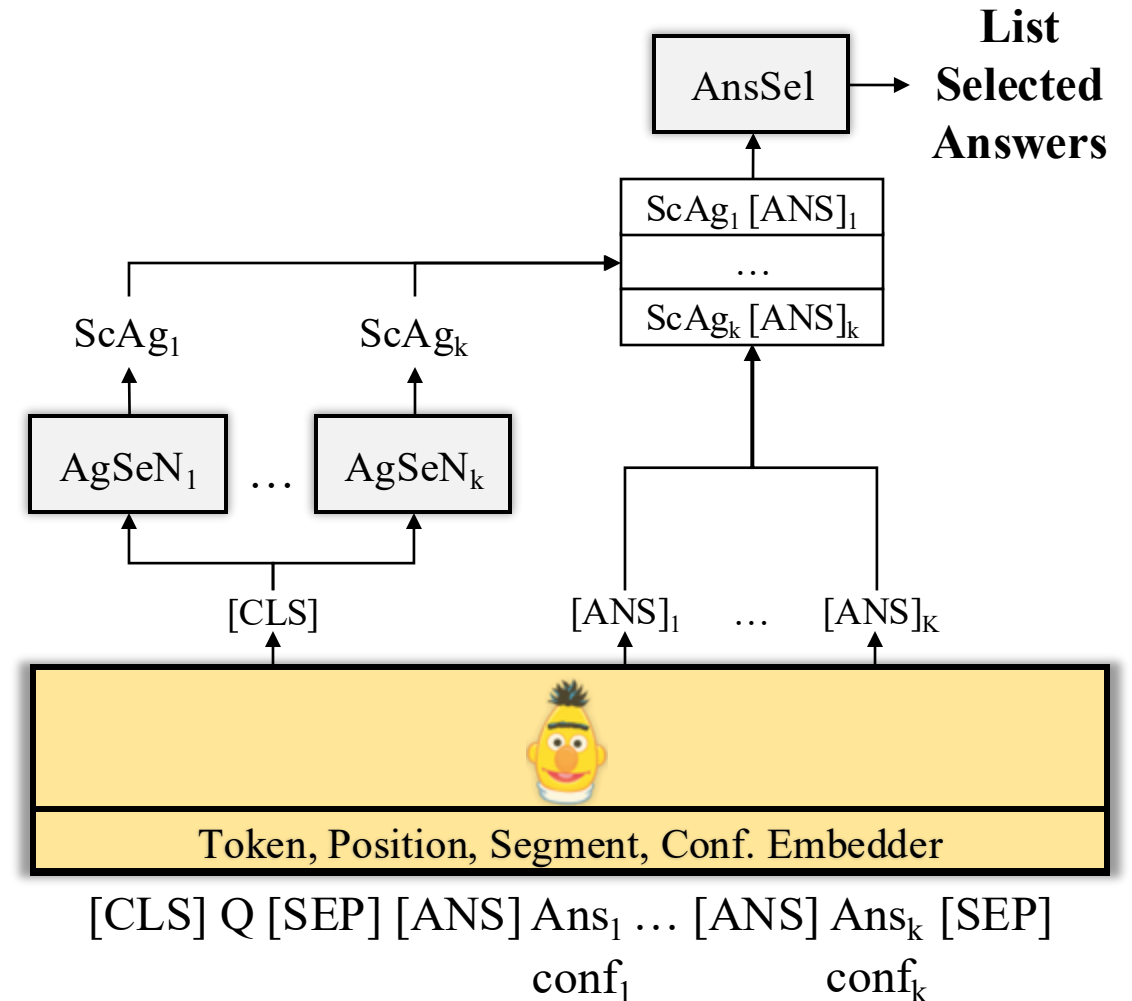
# METAQA

# OVERVIEW



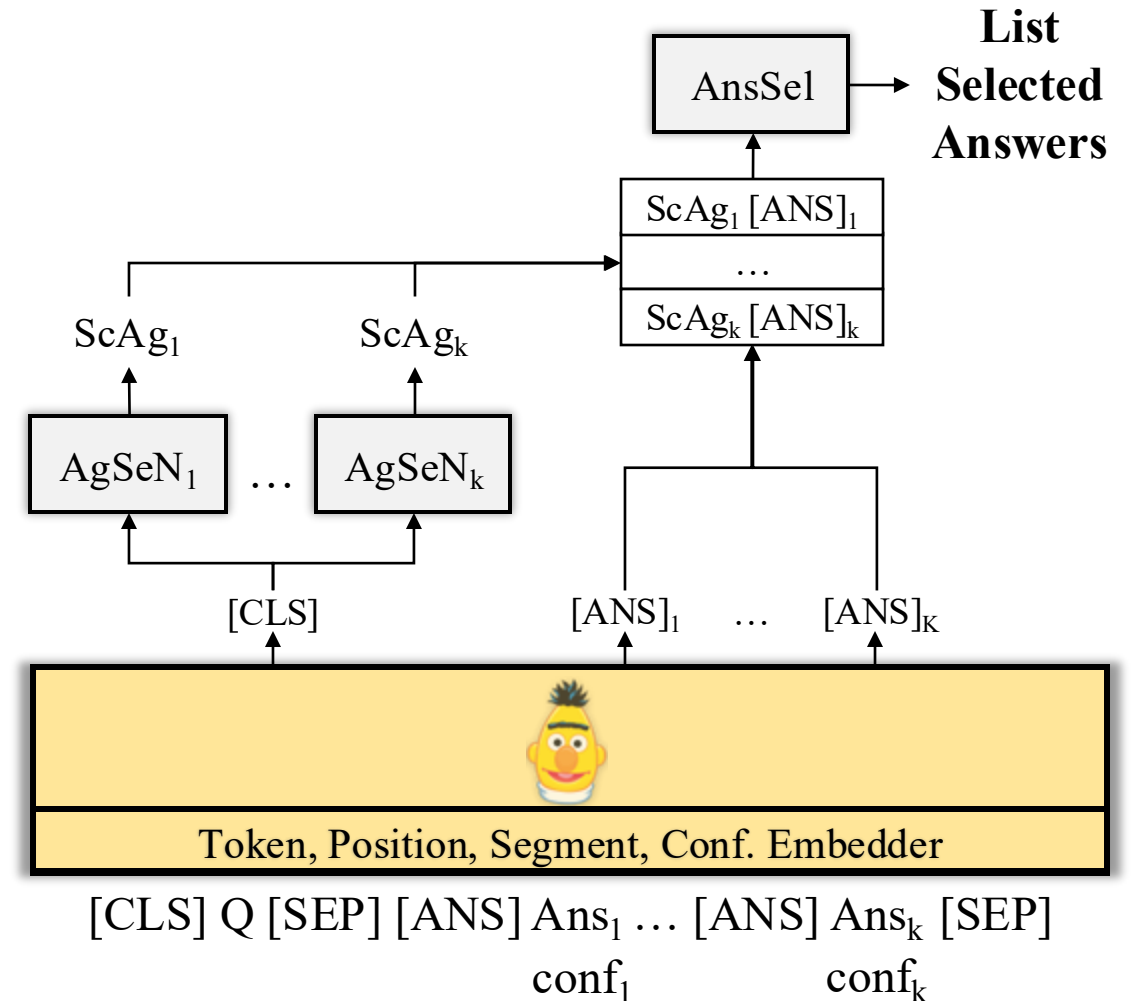
# MODEL

- **Input:** Question, Answers, Confidence Scores
- **Output:** List of selected answers
- **Task:** identify the best answer



# MULTI-TASK

- **Agent Selection**  $\approx$  Identify the in-domain agent
- **Answer Selection**: selects the best answer based on:
  - Domain of the question
  - Answer semantics



# DATASETS

Dataset	Dataset	Dataset
SQuAD (Rajpurkar et al., 2016)	RACE (Lai et al., 2017)	DROP (Dua et al., 2019)
NewsQA (Trischler et al., 2017)	CSQA (Talmor et al., 2019)	NarrativeQA (Kočíský et al., 2018)
HotpotQA (Yang et al., 2018)	BoolQ (Clark et al., 2019)	HybridQA (Chen et al., 2020)
SearchQA (Dunn et al., 2017)	HellaSWAG (Zellers et al., 2019)	Abstractive & MultiModal QA
NQ (Kwiatkowski et al., 2019)	SIQA (Sap et al., 2019)	
TriviaQA-web (Joshi et al., 2017)	Multiple-Choice QA	
QAMR (Michael et al., 2018)		
DuoRC (Saha et al., 2018)		

Extractive QA

1 agent for each dataset



**METAQA: COMBINING EXPERT AGENTS FOR MULTI-SKILL QUESTION ANSWERING**

# RESULTS

# RESULTS

Datasets	MetaQA	UnifiedQA	$\Delta$
Average (all – 16)	<b>76.39</b>	65.9	10.49
Extractive (8)	<b>77.72</b>	72.65	5.07
Multiple Choice (5)	<b>76.77</b>	61.41	15.36

# LIMITATIONS OF MULTI-DATASET MODELS

- Cannot model specific tasks (eg: numerical reasoning)
- Weak in minority domains (eg: Trivia and movie questions)

Dataset	MetaQA	UnifiedQA	$\Delta$
TriviaQA-web	<b>80.63</b>	72.34	8.29
DuoRC	<b>51.24</b>	34.65	16.59
CSQA	<b>78.66</b>	58.43	20.23
HellaSWAG	<b>73.19</b>	36.01	37.18
DROP	<b>73.04</b>	42.45	30.59

# OOD RESULTS

Model	TriviaQA	DuoRC	CSQA	HellaSWAG	SIQA	Δ
UnifiedQA	72.34	34.65	58.43	36.01	<b>61.62</b>	-
OOD MetaQA	<u><b>75.02</b></u>	<u><b>50.51</b></u>	<u><b>58.59</b></u>	<u><b>52.13</b></u>	59.28	-
OOD UnifiedQA	69.33	32.84	50.57	29.35	44.93	-8.12

Best results in bold.

OOD MetaQA outperform in-domain UnifiedQA underlined

# COLLABORATION BETWEEN AGENTS

Dataset	In-Domain Agent	OOD Agent
DuoRC	48%	NewsQA (18%)
TriviaQA	80%	DuoRC (3%)
SearchQA	86%	TriviaQA (5%)

Dataset	Question	In-domain Agent	OOD Agent
DuoRC	Who does Rocky Balboa work for as an enforcer?	<b>Adrian</b>	<b>Tony Gazzo</b> (NewsQA Agent)
TriviaQA-web	Who played the character Mr Chips in the 2002 TV adaptation of Goodbye Mr Chips?	<b>Timothy Carroll</b>	<b>MartinClunes</b> (DuoRC Agent)
SearchQA	This short story, written around 1820, contains the line "If I can but reach that bridge... I am safe"	<b>Legend</b>	<b>Legend of Sleepy Hollow</b> (TriviaQA Agent)

# EFFICIENCY

## Training

- MetaQA is very training efficient
- MetaQA uses 16% of the data used by UnifiedQA

MetaQA learns to match Questions with Answers instead of end-to-end QA

## Inference

Model	Time (s) per question
UnifiedQA	2.15 ± 0.02
MetaQA	7.08 ± 0.16

On SQuARE Platform using CPU



<https://square.ukp-lab.de>

# CONCLUSIONS

- We proposed MetaQA, a system to combine multiple expert agents
- Large performance gains wrt multi-dataset models (in-domain and OOD scenarios)
- Agent Collaboration → key to performance
- MetaQA is highly efficient to train
- Inference time remains reasonable



## **MetaQA: Combining Expert Agents for Multi-Skill Question Answering**

**Haritz Puerto<sup>1</sup>, Gözde Gül Şahin<sup>2</sup>, Iryna Gurevych<sup>1</sup>**

<sup>1</sup>Ubiquitous Knowledge Processing Lab (UKP Lab)

Department of Computer Science and Hessian Center for AI (hessian.AI)  
Technical University of Darmstadt, Germany

<sup>2</sup> Department of Computer Science, Koç University, KUIS AI Lab  
Istanbul, Türkiye





**SQUARE: TOWARDS MULTI-DOMAIN AND FEW-SHOT COLLABORATING QA AGENTS**

# RELEVANCE FEEDBACK IN NEURAL RE-RANKING

# QUERY TYPES [1]

	Navigational / Transactional	Information-Seeking
Example query	“arxiv sentence bert” “acl paper submission”	“origin coronavirus” “neural sentence representation”
# relevant documents	1-few	Many
Scenarios	Navigation Known-Item Retrieval	Scientific Literature Review News Background Case Law Factoid QA

# MOTIVATION

## Challenges

- Unknown Search Domain
- Difficult Query Formulation

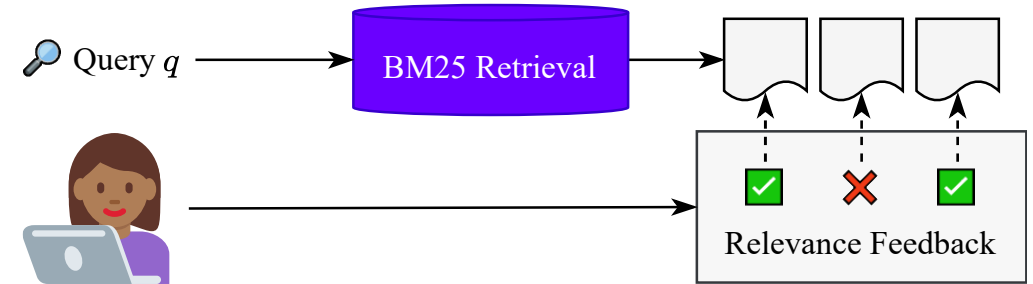
## Observation

- Judging a document is easier than formulating good queries
- Some relevant documents might be already known to the user

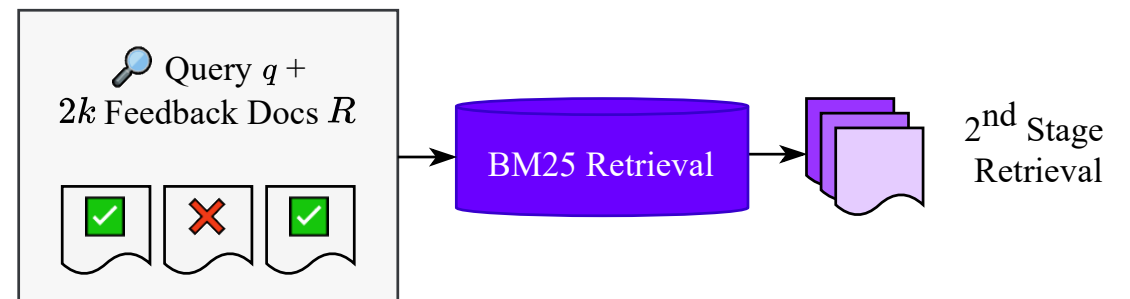
**→ Use relevance feedback to improve search**

# BACKGROUND

- 1) Retrieve Documents using the query
- 2) Obtain feedback on retrieved documents



- 3) Extract "expansion terms" from documents
- 4) Retrieve with query + expansion terms



**How to integrate relevance feedback into neural retrieval?**

# TASK SETUP

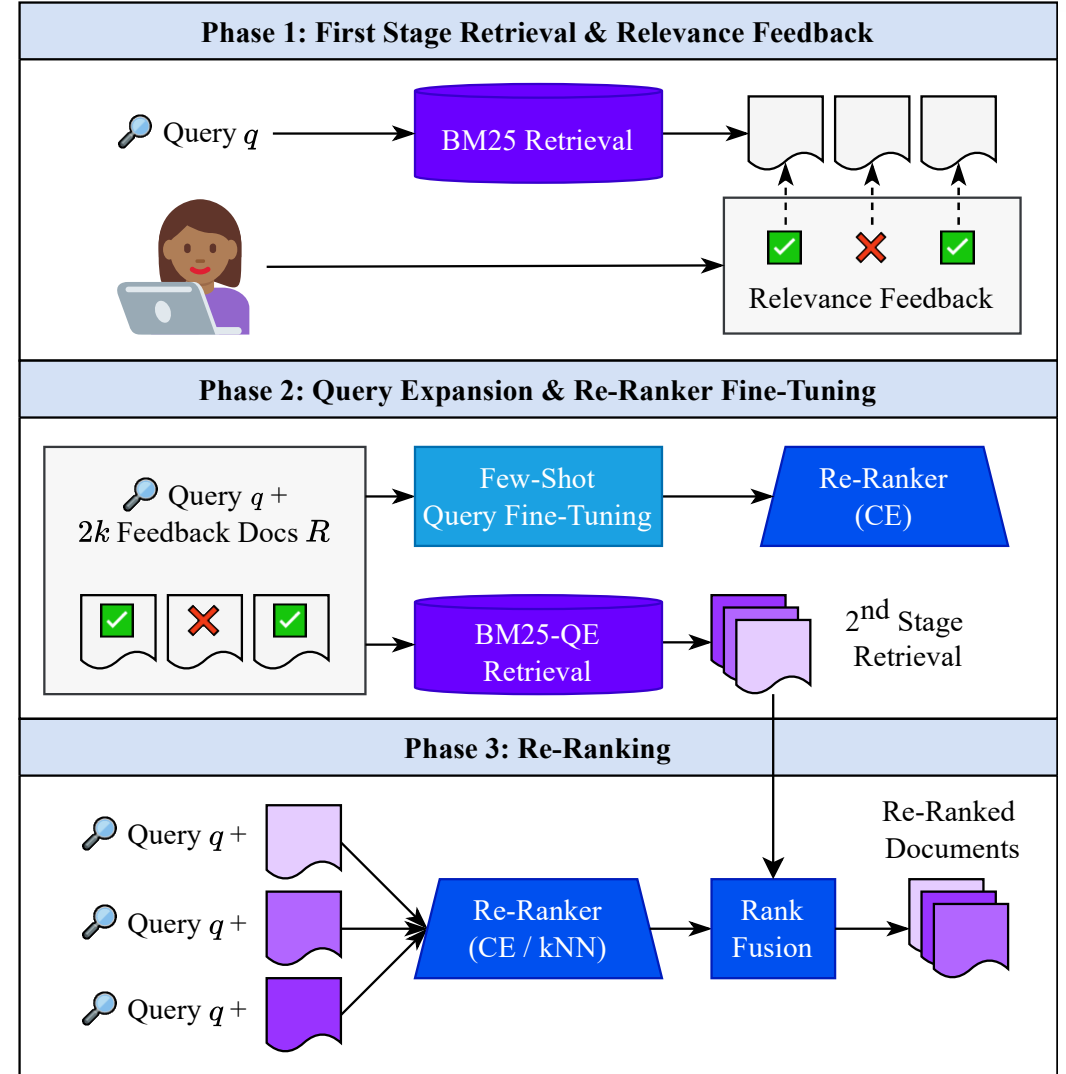
**Goal:** Re-rank documents from 2. Stage Retrieval with neural re-rankers incorporating relevance feedback

## Input

- Query  $q$
- $k$  relevant and non-relevant feedback documents, where  $k \in \{2, 4, 8\}$
- 1000 documents from 2. stage retrieval

## Evaluation

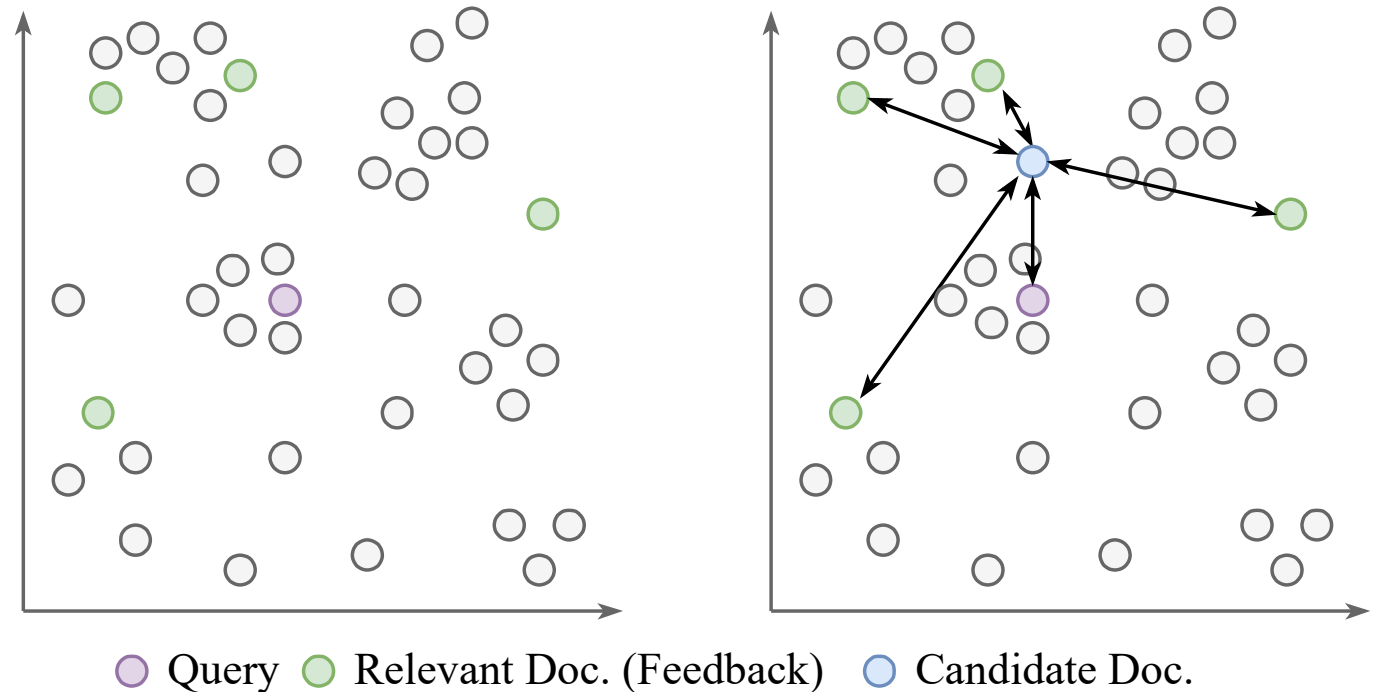
- Re-ranking [nDCG@20]



# METHOD: KNN

- Compute document representations  $d_i \in D$
- Score documents by summing the similarities between the candidate document  $d_i$ , query  $q$  and relevant feedback documents  $d_j \in R^+$

$$s_i = f(d_i, q) + \sum_{d_j \in R^+} f(d_i, d_j)$$



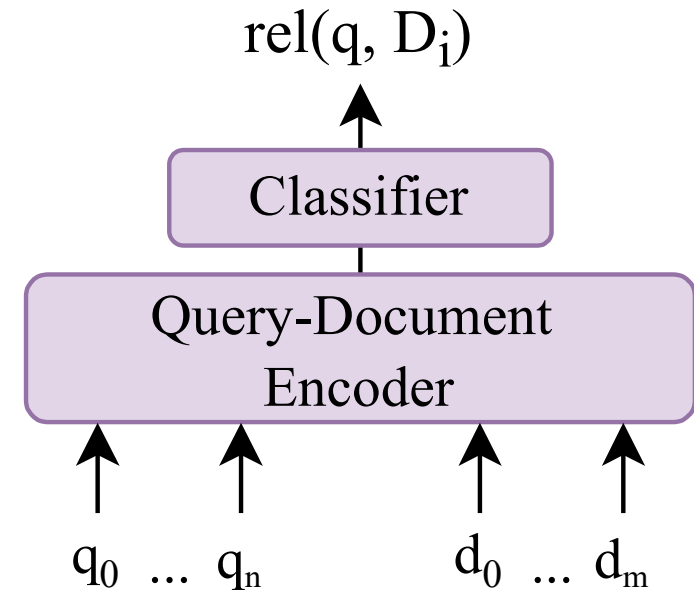
# METHOD: CROSS-ENCODER [2]

## CE Query Fine-Tuning

Fine-Tune bias layers per query on  $2k$  Relevance Feedback Documents

## CE MAML + Query Fine-Tuning

1. Fine-Tune bias layers on in-domain annotations with Meta-Learning to obtain "fast parameters"
2. Fine-Tune per query on on  $2k$  Feedback Documents



# METHOD: RANK-FUSION [3]

**Idea:** Merge rankings of different ranking functions  $h \in H$

**Problem:** Different methods produce different scores, simple adding the scores is biased

=> Use ranks instead of raw scores

$$s_i = \sum_{h \in H} \frac{1}{c + h(d_i)}$$

$h(d_i)$  returns the integer rank that method  $h$  assigns document  $d_i$

$c$  is constant smoothing the impact of top ranked documents

# DATASETS

Dataset	Domain	Docs.	Queries	Judgments
Robust04	News	528k	148	1287 ( $\pm 501$ )
TREC-Covid	Biomedical	191k	50	1370 ( $\pm 323$ )
TREC-News	News	595k	34	259 ( $\pm 82$ )
Touché-2020	Debates	383k	49	50 ( $\pm 7$ )

*Includes only queries with at least 32 relevant documents.*

# RESULTS

Method	Robust	Covid	News	Touché	Avg
BM25-QE (2. Stage Retrieval)	<b>0.496</b>	0.610	<b>0.392</b>	<b>0.271</b>	<b>0.442</b>
kNN	0.443	0.686	0.365	0.174	0.417
CE Zero-Shot	0.415	<b>0.702</b>	0.314	0.176	0.402

=> *BM25-QE is hard to beat!*

# RESULTS

Method	Robust	Covid	News	Touché	Avg
BM25-QE (2. Stage Retrieval)	0.496	0.610	<b>0.392</b>	<b>0.271</b>	<b>0.442</b>
CE Query Fine-Tune	0.484	0.723	0.335	0.198	0.435
CE MAML + Query FT	<b>0.506</b>	<b>0.735</b>	0.314	0.223	<b>0.445</b>

=> Fine-tuning per query helps  
=> CE MAML + Query FT is on par with BM25-QE

# RESULTS

Method	Robust	Covid	News	Touché	Avg
BM25-QE (2. Stage Retrieval)	0.496	0.610	0.392	<b>0.271</b>	0.442
kNN	0.443	0.686	0.365	0.174	0.417
CE Zero-Shot	0.415	0.702	0.314	0.176	0.402
CE Query Fine-Tune	0.484	0.723	0.335	0.198	0.435
CE MAML + Query FT	0.506	0.735	0.314	0.223	0.445
BM25-QE $\cap$ kNN	0.507	0.707	<b>0.412</b>	0.248	0.468
BM25-QE $\cap$ CE MAML + Query FT	<b>0.570</b>	<b>0.740</b>	0.405	<b>0.272</b>	<b>0.497</b>

=> Fusing BM25-QE and neural model is highly effective!



# Incorporating Relevance Feedback for Information-Seeking Retrieval using Few-Shot Document Re-Ranking

**Tim Baumgärtner,<sup>1</sup> Leonardo F. R. Ribeiro,<sup>1\*</sup> Nils Reimers,<sup>2</sup> Iryna Gurevych<sup>1</sup>**

<sup>1</sup>Ubiquitous Knowledge Processing Lab (UKP Lab),  
Department of Computer Science and Hessian Center for AI (hessian.AI),  
Technical University of Darmstadt

<sup>2</sup>cohere.ai

[www.ukp.tu-darmstadt.de](http://www.ukp.tu-darmstadt.de)



**EMNLP  
2022**





**SQUARE: TOWARDS MULTI-DOMAIN AND FEW-SHOT COLLABORATING QA AGENTS**

# UKP-SQUARE

AN ONLINE PLATFORM FOR QA RESEARCH

# AN ECOSYSTEM FOR QA RESEARCH



Common **interface**



Easy **analysis** of the strengths, weaknesses, and biases



Common **explainability** methods



Common **adversarial** attacks



Graph **visualizations**



Running on the **browser** (no configurations, no installations)



Re-use **data sources**



Easy to **deploy** new models



Dozens of **models available**



ACL 2022



ACL 2022




arXiv 2023

Data icons created by Freepik - Flaticon

# EXPLAINABILITY

## Saliency Maps

- Attention (Jain et al., ACL 2020)
- Scaled Attention (Serrano & Smith, ACL 2019)
- Simple Gradient (Simonyan et al., arXiv 2013)
- Smooth Gradients (Smilkov et al., arXiv 2017)
- Integrated Gradients (Sundararajan et al., PMLR 2017)

Saliency Map 

Method: Attention Scaled Attention Simple Gradients Smooth Gradients Integrated Gradients

Showing the top 3 most important words

---

**NewsQA BERT Adapter**

**Question:** what was problem with getting marriage license ?

**Context:** new orleans , louisiana ( cnn ) - two newlyweds are fighting for the dismissal of the justice of the peace who refused them a marriage license because they are of different races . we ' ve retained an attorney , and we ' re in the process of taking the next steps in order to make sure that ( the justice of the peace ) loses his job .

**Answer:** they are of different races.

# ADVERSARIAL ATTACKS

- HotFlip (Ebrahimi et al., ACL 2018)
- Input Reduction (Feng et al., EMNLP 2018)
- Sub-Span
- Top-K

### Attack Methods

Method: HotFlip **Input Reduction** Sub-Span Top K

# Reductions = 10 

---

#### SQuAD 1.1 BERT Adapter

Question: to whom **did the virgin mary** allegedly **appear in 1959 in Lourdes** france ?

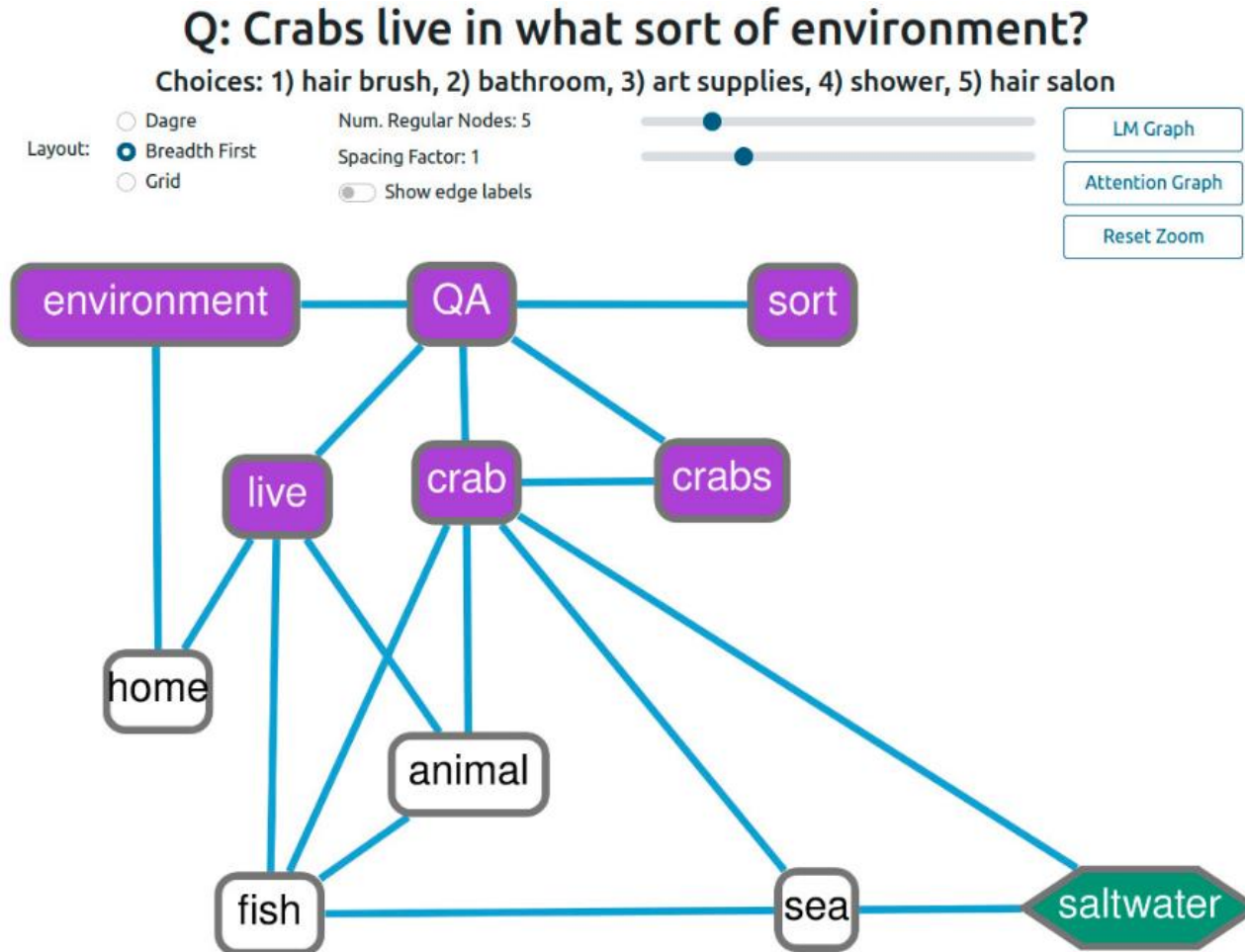
Context: Architecturally, the school has a Catholic character. Atop the Main Building's gold dome is a golden statue of the Virgin Mary. Immediately in front of the Main Building and facing it, is a copper statue of Christ with arms upraised with the legend "Venite Ad Me Omnes". Next to the Main Building is the Basilica of the Sacred Heart. Immediately behind the basilica is the Grotto, a Marian place of prayer and reflection. It is a replica of the grotto at Lourdes, France where the Virgin Mary reputedly appeared to Saint Bernadette Soubirous in 1858. At the end of the main drive (and in a direct line that connects through 3 statues and the Gold Dome), is a simple, modern stone statue of Mary.

New Answer: saint bernadette soubirous      Old Answer: Saint Bernadette Soubirous

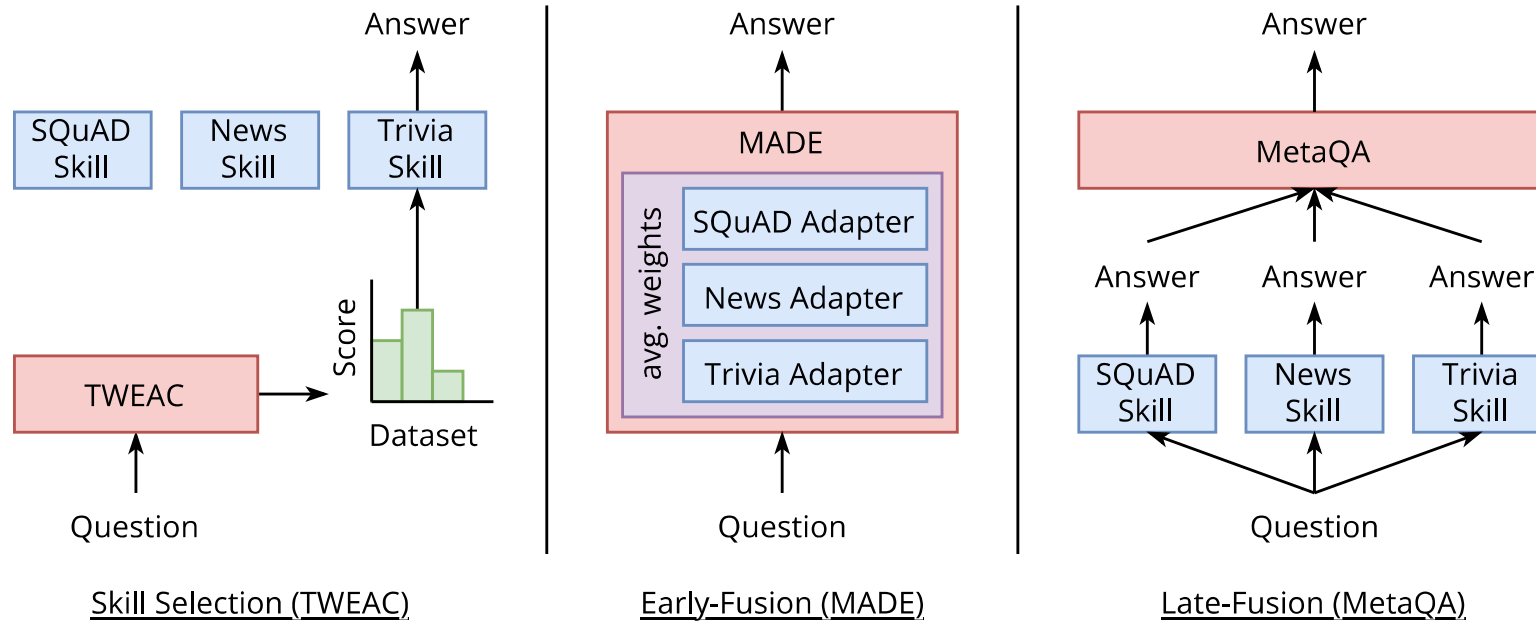
# GRAPH-BASED MODELS

## QA-GNN

- LM + KB for common-sense QA
- KB: ConceptNet
- Graph Viz can help us understand the behavior of the model



# MULTI-AGENT SYSTEMS



- Qualitative analysis of multi-agent systems
- Qualitative comparison against multi-dataset systems



**SQUARE: TOWARDS MULTI-DOMAIN AND FEW-SHOT COLLABORATING QA AGENTS**

**DEMO**

[HTTPS://SQUARE.UKP-LAB.DE](https://square.ukp-lab.de)



**SQUARE: TOWARDS MULTI-DOMAIN AND FEW-SHOT COLLABORATING QA AGENTS**

# **FUTURE WORK**

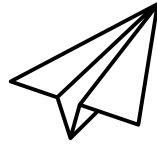


# PROMPTING LANGUAGE MODELS

- Comparing multiple LM in parallel
- Sharing prompts easily with the community
- Augmenting LM
  - Documents
  - External tools (eg: calculators)
  - Knowledge Graphs
  - Other (smaller) models



# THANK YOU!



[gurevych@tu-darmstadt.de](mailto:gurevych@tu-darmstadt.de)



<https://square.ukp-lab.de>



[https://twitter.com/UKP\\_SQuARE](https://twitter.com/UKP_SQuARE)



<https://github.com/UKP-SQuARE>