

# TREC Complex Answer Retrieval Overview

Laura Dietz\*, Ben Gamari, Jeff Dalton, Nick Craswell

Homepage: <http://trec-car.cs.unh.edu>

Google group mailinglist: TREC-CAR

## Abstract

This notebook gives an overview of activities, datasets, and results of the second year of TREC Complex Answer Retrieval. We lay out the tasks offered and how provided datasets are automatically derived from Wikipedia and TQA. Manual relevance assessments are created by NIST. We describe the details of the assessment procedures, inter-annotator agreement, and statistics. Nine teams submitted runs exploring interactions of entities and passages, neural as well as traditional retrieval methods. We see that combining traditional methods with learning-to-rank can outperform neural methods, even when many training queries are available.

## 1 Introduction

Generating information objects, sub-document retrieval, answer aggregation, retrieving material for conversational agents are important challenges raised in the SWIRL 2012 and 2018 workshops on frontiers, challenges, and opportunities for information retrieval [1, 3]. These challenges share several commonalities: We desire answers to complex information needs, and wish to find them in a single and well-organized page. Such a page may not yet exist, therefore it needs to be synthesized from multiple information sources. Advancing the state of the art in this area is the goal of this TREC Complex Answer Retrieval track.

Algorithms that can automatically author a complex answer in response to a query would benefit users that investigate new and unfamiliar topics. They would improve information access in mobile environments with restricted interaction capabilities. In contrast to extensive work on finding the best short answer, the CAR track focuses on the retrieval of longer answers—especially answers that cover a range of different subtopics. We envision answers to be composed of multiple text fragments from multiple sources, recycling information about related topics, but selected to highlight insightful connections.

Retrieving high-quality comprehensive answers is challenging as it is not sufficient to choose a lower rank-cutoff with the same techniques as for short answers. Instead, we need new approaches for finding and organizing relevant information units of a complex answer space.

Many examples of manually created complex answers exist on the Web: [howstuffworks.com](http://howstuffworks.com), travel guides, fanzines, or educational text books. These are collections of articles where each article constitutes a long answer to an information need expressed by the title of the article. Ideally, by reading an article, users will gain new information about the topic. We can measure this information gain by testing how well the text enables users to answer questions about the topic.

The fundamental task of collecting references, facts, and opinions into a single point-of-entry has traditionally been a manual process. We envision that automated information retrieval systems can relieve users from a large amount of manual work through sub-document retrieval, consolidation and organization. Ultimately, the goal is to retrieve synthesized *information* rather than documents.

## 2 A Worked Example

To motivate a brief example, consider a user wondering about how coffee preparation techniques lead to different tastes. With this intention in mind, she enters the query “Coffee preparation”. A possible answer

---

\*dietz@cs.unh.edu

could look as follows:

## Coffee Preparation

### Grinding

*Arabic coffee and Turkish coffee require that the grounds be almost powdery in fineness, finer than can be achieved by most burr grinders. Pounding the beans with a mortar and pestle can pulverize the coffee finely enough.*

*The fineness of the grind strongly affects brewing. Brewing methods that expose coffee grounds to heated water for longer require a coarser grind than faster brewing methods.*

### Steeping

*The softer flavors come out of the coffee first and the more bitter flavors only after some time, so a large brew will tend to be both stronger and more bitter.*

*Another variation is cold brew coffee, sometimes known as “cold press.” Cold water is poured over coffee grounds and allowed to steep for eight to twenty-four hours. This process produces a very strong concentrate which can be stored in a refrigerated, airtight container for up to eight weeks.*

Here “grinding” and “steeping” are two important facets. After reading this article, the user is able to tell that the taste is affected by the fineness of the ground, the time of exposure, and the temperature.

Of course, one might envision other responses that would satisfy the information need equally well. While this example was taken from Wikipedia<sup>1</sup> it should be possible to identify such information from a Web collection with passage retrieval, consolidation, and organization.

## 3 Task Description

While the long-term goal of this track is to retrieve complex answers without any more information than the given complex topic, in this year we focus on a simpler task, where both the topic and an appropriate outline is provided as a query. An example outline is given in Figure 1. We offer two tasks: passage and entity.

**Passage Task:** Given an outline for complex topic  $Q$ , retrieve for each of its sections  $H_i$ , a ranking of relevant passages  $S$ .

**Entity Task:** Given outline for complex topic  $Q$ , retrieve for each of its sections  $H_i$ , a ranking of relevant entities  $E$ . Each entity is to be supported with a passage  $S$  that motivates the why the entity is relevant for the query.

The passage  $S$  is taken from the provided passage corpus. The entity  $E$  refers to an entry in the provided knowledge base. We define a passage or entity as relevant if the passage content or entity is appropriate to be mentioned an article about the topic  $Q$ .

## 4 TREC CAR Data Set (v2.1)

The 2018 Complex Answer Retrieval track uses topics and outlines that are extracted from English Wikipedia (XML dump from June 2018) and the Text Book Question Answering (TQA) dataset<sup>2</sup>. We refer to these subsets as Wiki-18 and TQA in the following. Paragraphs, entities and training data is extracted from previous years’s English Wikipedia (XML dump from Dec 20th, 2016; referred to as Wiki-16). Wikipedia articles are split into the outline of sections and the contained paragraphs.

<sup>1</sup>[https://en.wikipedia.org/wiki/Coffee\\_preparation](https://en.wikipedia.org/wiki/Coffee_preparation)

<sup>2</sup>Available at <http://data.allenai.org/tqa/>

**MUST be mentioned:**

Many water-saving devices (such as low-flush toilets) that are useful in homes can also be useful for business water saving. Other water-saving technology for businesses includes:

**CAN be mentioned:**

Recycling one gallon of paint could save 13 gallons of water, 1 quart of oil, and 250,000 gallons of water pollution, 13.74 pounds of , save enough energy to power the average home for 3 hours, or cook 6 meals in a microwave oven, or blow dry someone’s hair 27 times.

**Roughly on TOPIC but non-relevant:**

Dual piping is a system of plumbing installations used to supply both potable and reclaimed water to a home or business. Under this system, two completely separate water piping systems are used to deliver water to the user. This system prevents mixing of the two water supplies, which is undesirable, since reclaimed water is usually not intended for human consumption.

Paragraph IDs:

left: dbcef592762b4711012041f6bdf1bdf7cb5a521

center: f26730da3b7860c727411480b08ae6466dcc9a54

right: 21e6e381383e392cb7d1432200c51c095cdf3fbe

Figure 2: Example passages and relevance for “Protecting the Water Supply / Saving Water at Home” (Query ID “tqa:protecting%20the%20water%20supply/Saving%20Water%20at%20Home”).

All paragraphs from all articles are gathered and deduplicated to form the paragraph corpus. Due to a bug fix in the Wikipedia parser, the paragraph collection for Y2 is larger and cleaner than in Y1.

Each section outline, such as depicted in Figure 1, is a description of a complex topic. By keeping the information which paragraph originates from which article and section, we have a means of providing (automatic) training data for the passage retrieval task. By preserving hyperlinks inside paragraphs that point to other Wikipedia pages (also known as entity links), we have a means of providing training data for the entity retrieval task.

Title: Protecting the Water Supply

1. Rationing Water
2. Reducing Water Pollution
3. Saving Water in Irrigation
4. Conserving Water
5. Water Treatment
6. What You Can Do
7. Saving Water at Home
8. Controlling Water Pollution

Figure 1: Example outline for TQA topic.

Figure 3 depicts the set of Wiki-18 pages selected for the benchY2test set. This allows us to study how well Wikipedia content can be recycled to populate articles on new and unseen topics. With the release of allButBenchmark, all Wiki-16 pages are made available. However, taking paragraphs and outlines from different collections poses significant challenges. Only a small subset of these Wiki-18 pages in contained paragraphs available in the provided paragraph collection (derived from Wiki-16). An unfortunate consequence is that the “automatic passage evaluation” option used Y1 (Section 4.2) cannot be applied to passage runs in this year’s evaluation.

After filtering and processing procedures described in Section 4.1, several datasets for training and evaluation are derived. The size of the datasets is given in Table 1. The paragraph collection contains 29,678,367 unique paragraphs.

### 4.1 Data Set Creation Pipeline

The TREC Complex Answer Retrieval benchmark (v2.1) is derived from Wikipedia so that complex topics are chosen from articles on open information needs, i.e., not people, not organizations, not events, etc. However, any paragraph or entity on Wikipedia is a legal paragraph/entity for the retrieval task even if a person entity or a paragraph from an article on an event. The data set creation process (similar to the v1.5 data) is as follows:

1. Mediawiki format of each article in the Wikipedia dump is parsed, preserving paragraph boundaries,

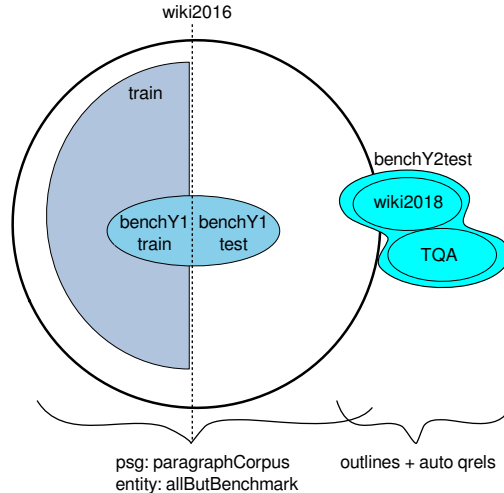


Figure 3: Y1 and Y2 datasets. Manual Qrels for benchY1test were collected in the previous year’s evaluation.

intra-Wikipedia hyper links, and section hierarchy. The TQA collection was obtained preprocessed in JSON format.

2. Name information from redirect and disambiguation pages and category information added to the article page. Redirects in hyperlinks are resolved. Templates, talk pages, portals, disambiguation, redirect, and category pages are discarded. (Disambiguation, redirects, and category information only provided in training set).
3. While articles about query topics in benchY1train, benchY1test, benchY2test, and test200 are withheld, all remaining articles in Wiki-16 are released as a legal set for entity retrieval (**allButBenchmark**).
4. Articles tagged with categories that indicate people, organizations, music, books, films, events, and lists are discarded.
5. Sections with headings that do not contain prose are discarded, for example external links, references, bibliography, notes, gallery etc.
6. Each article is separated into (1) the outline of section headings and (2) paragraphs.
7. The set of paragraphs across all of Wiki-16 are collected, and unique paragraph IDs are derived through SHA256 hashes on the text content (ignoring links).
8. The paragraphs are further deduplicated with min hashing using word embedding vectors provided by GloVe (with 50 dimensions). For each set of duplicates, one representative paragraph is chosen.
9. The collection representative paragraphs is released as the **paragraphCorpus**.
10. Articles are rewritten, replacing paragraphs that have duplicates with the representative paragraph.
11. The set of articles is further filtered to remove images, sections with very long (>100 characters), and very short headings (<3 letters). Articles with less than three remaining sections are discarded.

Table 1: Data set sizes in terms of articles, section, and automatic positive assessments. (\*\*) used in this years’ evaluation.

	benchY1train	benchY1test	train	benchY2test **	
	auto	manual	auto	manual	auto
number of articles (complex topics)	117	133	285,924	27	65
hierarchical sections (queries)	1,816	2,125	2,180,868	269 / 271	976
total positive paragraphs assessments	4,530	5,820	5,276,624	3,788	N/A
total positive entity assessments	13,031	15,085	12,310,616	3,173	17,044

	Paragraph		Entity	
	Automatic	Manual	Automatic	Manual
train	yes	-	yes (original entity)	-
benchmarkY1train	yes	-	yes (TagMe)	-
benchmarkY1test	yes	yes	yes (TagMe)	no
benchmarkY2test Wiki-18	no, too small	yes	yes (TagMe)	yes
benchmarkY2test TQA	no	yes	yes (TagMe)	yes

Table 2: Availability of manual and automatic assessments for different benchmarks.

12. The set of Wiki-16 articles is split into training and holdout data, training data is further split into five folds. To ensure uniform distribution and reproducibility, these decisions are made based on the SipHash of the article title.
13. The five folds of the training data, with separated outlines and paragraphs and extracted automatic qrels are made available as **train**.
14. The set of pages used for benchmarks used in Y1 (**benchY1train**, **benchY1test**, **test200**) were re-processed with the new Wikipedia parser, manual judgments were translated to new paragraph IDs, the re-released.
15. A manual selection of articles in Wiki-18 and the TQA corpus were only released as outlines as **benchY2test.public**. Official contributed runs were submitted on these topics. A complete **benchmarkY2test** is released after the TREC workshop.<sup>3</sup> It contains:
  - Manual ground truth for paragraphs (qrels)
  - Automatic and manual ground truth for entities (qrels)
  - Original articles

## 4.2 Automatic Ground Truth

Two kinds of ground truth signals are collected: automatic and manual. For each, we release true paragraphs and true entities. While the manual ground truth is assessed after participants submit runs, the automatic ground truth is derived along with the dataset from the Wikipedia/TQA dump.

For the benchmarkY1 datasets, the automatic ground truth is derived as follows

- If a paragraph is contained in the page/section it is defined as relevant, and non-relevant otherwise.
- If the page/section contains an entity link, then the (link target) entity it is defined as relevant, and non-relevant otherwise. For benchmarkY1train, benchmarkY1test, and benchmarkY2test, the TagMe [5] entity linker was used. For the much larger “train” collection, entity links that were manually added by Wikipedia editors were used.

Qrels are derived with several different levels:

- Hierarchical: Only content of leaf sections is considered.
- Article: All content of the page is considered (independent of the section). However, due to a mistake in processing, the lead text is missing.
- Top-level: Only relevance for top-levels sections is provided. All content in this section or a child section is considered relevant.
- Tree: For all page titles and headings, all content in the subtree is considered relevant. Tree qrels also contain article-level relevance assessments, although these are not part of the TREC CAR evaluation.

The benchmarkY2 dataset is constructed from pages outside the Wiki-16 dump (depicted in Figure 3). Only a small fraction of paragraphs on Wiki-18 pages already existed in Wiki-16 on a page with a different name. The paragraph sets from TQA and Wiki-16 are disjoint. Thus, the automatic evaluation procedure for paragraphs, used in Y1, is not applicable to the Y2 dataset.

The automatic entity ground truth is available on benchmarkY2test as it is constructed from entity links and does not rely on an overlap in paragraphs.

An overview for which benchmarks which kind of ground truth is available is given in Table 2.

## 5 Submission

For the passage ranking task, participants were asked to submit a ranking of paragraph IDs per heading in the outlines of benchY2test. For the entity ranking task, participants were asked to submit a ranking of entity IDs per heading in the outlines of benchY2test. To support assessing entity relevance, participants were asked to provide provenance to annotate each entity ID with the ID of a paragraph that explains why the entity is relevant for the corresponding section heading.

Participants were allowed to consider all headings in the outline at once, use external resources such as knowledge graphs, entity linking tools, pre-trained word embeddings, and any of the provided TREC CAR data sets. **The participants were not allowed to directly use a dump of Wikipedia**, as this would allow them to look up the paragraphs on the page—the information used in the automatic ground truth.

Each participating team was allowed to submit up to three runs to the passage task and three runs to the entity task. Nine teams participated in this second year of the track.

## 6 Assessment of the Manual Ground Truth

For each heading of query outlines, the top five paragraph IDs or entity IDs of participant-contributed runs were merged to build the assessment pools. Additionally, paragraphs and entities that are relevant according to the automatic ground truth were added to the pool for verification. In the previous year, complex topics were only partially judged (to obtain a larger topical variety). In contrast, this year, all sections of complex topics were assessed. This yielded manual assessments for 65 complex topics (31 TQA and 34 Wiki-18). One section was annotated by all assessors in order to measure inter-annotator agreement across the six NIST assessors (see Section 6.2).

For the passage task, the assessor is presented with the complex topic (page title) and the headings in the outline, followed by a randomized list of paragraphs from the assessment pool. In the case of the entity task, the list displayed the canonical entity names together with the provenance paragraph if given. As not all participants submitted provenance, the list also displays an entry of the canonical entity name together with first paragraph from the entity’s Wikipedia pages as provenance. This information was intended to support the assessment process. However, the first paragraph of the entity’s Wikipedia page turned out to be generally not relevant. As only one team submitted provenance, assessors has no choice by to resort to world-knowledge.

Assessors were asked to envision writing a Wikipedia article on the given complex topic. A graded assessment scale was used based on how importantly the paragraph/entity should be mentioned in this section of the article, using grades as follows.

- MUST be mentioned
- SHOULD be mentioned
- CAN be mentioned
- Non-relevant, but roughly on TOPIC of the page
- NO, non-relevant
- Trash

The grade “trash” is assigned to paragraphs/entities of low quality which therefore would not be relevant for any topic imaginable.

---

<sup>3</sup><http://trec-car.cs.unh.edu/datareleases/>

Table 3: Assessment scale for manual assessments. Horizontal line: Cutoff for positive/negative assessments.

	binary scale	manual scale	manual lenient scale
MUST be mentioned	1	3	5
SHOULD be mentioned	1	2	4
CAN be mentioned	1	1	3
Non-relevant, but roughly on TOPIC	0	0	2
NO, non-relevant	0	-1	0
Trash	0	-2	-2

Table 4: Manual assessment: Grade histogram and distribution across both passage and entity task.

	annotator1	annotator2	annotator3	annotator4	annotator5	annotator6	Total %
<b>Trash</b>	40	118	0	14	5	2	1
<b>No</b>	1213	980	810	1124	801	1066	42
<b>Topic</b>	439	674	613	361	770	946	27
<b>Can</b>	241	489	469	261	212	115	13
<b>Should</b>	210	304	181	305	140	472	11
<b>Must</b>	145	214	65	288	191	51	7

## 6.1 Label distribution

Six assessors created 13,310 passage annotations on a total of 269 topic sections for passages. For 271 topics sections, entity assessments were created with 8415 assessments in total.

The grade histogram per annotator and the overall grade distribution is given in Table 4 (discrepancies due to merging and cleaning). We notice that only a third of all assessments are graded as relevant, while an additional third were annotated as being on topic.

72% of passages (62% of entities) in the assessment pool were marked as non-relevant. This demonstrates that the task is feasible, but challenging.

## 6.2 Inter-annotator agreement

One section (both passage and entity) was selected for annotation by all assessors to measure inter-annotator agreement in the middle of the assessment period.

We measure inter-annotator agreement using Cohen’s  $\kappa$  for pairwise comparison and Fleiss’  $\kappa$  across all annotators. We analyze agreement on the derived binarized assessment and graded assessment. Subtle differences between neighboring grades are often not reliable. Therefore, we consider the case of graded

Table 5: Statistics for anual assessments after merging and cleaning (as used for results).

	Passage			Entity		
	Total%	TQA	Wiki-18	Total%	TQA	Wiki-18
<b>Trash</b>	1	25	131	0	1	2
<b>No</b>	43	2325	3364	25	642	1480
<b>Topic</b>	28	1601	2076	37	1215	1902
<b>Can</b>	12	1211	445	13	602	450
<b>Should</b>	10	1009	374	13	630	458
<b>Must</b>	6	400	349	12	585	448

Table 6: Inter annotator agreement according to Cohen’s  $\kappa$  on graded scale, counting grades that are “off by one” as agreement.

(a) Passage assessment.

	annotator1	annotator2	annotator3	annotator4	annotator5	annotator6
annotator1		0.687	0.430	0.781	0.628	0.449
annotator2	0.687		0.568	0.632	0.644	0.871
annotator3	0.430	0.568		0.409	0.422	0.859
annotator4	0.781	0.632	0.409		0.407	0.740
annotator5	0.628	0.644	0.422	0.407		0.512
annotator6	0.449	0.871	0.859	0.740	0.512	

(b) Entity assessment.

	annotator1	annotator2	annotator3	annotator4	annotator5	annotator6
annotator1		0.342	0.542	0.708	0.705	0.134
annotator2	0.342		0.526	0.096	0.500	0.514
annotator3	0.542	0.526		0.422	0.716	0.457
annotator4	0.708	0.096	0.422		0.643	0.467
annotator5	0.705	0.500	0.716	0.643		0.408
annotator6	0.134	0.514	0.457	0.467	0.408	

assessment where assessments that differ by no more than one grade step, e.g., grades “SHOULD” and “MUST”, are also counted as agreements for both  $p_0$  and  $p_e$ . We call this graded evaluation “off by one”, for which results are displayed in in Table 6.

Inspecting Cohen’s  $\kappa$ , we find that pair-wise agreement is relatively similar across all pairs of assessors. In other words, there is no “odd one out” which speaks to the quality of NIST’s assessment procedures. As expected, the agreement for binarized paragraph judgments (Fleiss  $\kappa = 0.500$ ) is higher than for graded judgments (Fleiss  $\kappa = 0.304$ ). For entity judgments, we find that Annotator 2 was an outlier, with a disproportionate high number of non-relevant assessments. After removing assessment from Annotator 2, binarized agreement is Fleiss  $\kappa = 0.416$ , and graded Fleiss  $\kappa = 0.374$ . This may sound small, yet it is comparable to previous work [2]. However, once neighboring grades are counted as agreement (“off by one”), the inter-annotator agreement is even higher agreement than on binarized assessments.

We conclude that, aside from subtle nuances in the grading scale, assessors agree on the whether the passage or entity should be included in the article on the complex topic.

### 6.3 Annotation Time

each, yielding a total of 240 hours. Including breaks and training, the average annotation time per passage or entity judgments depends on the annotator and ranges between 17 and 64 seconds (median: 30 seconds).

### 6.4 Interaction between Manual and Automatic Assessments

Figure 4 depicts differences between the sets of relevant passages and entities according to manual and automatic relevant data. Difficulties arise since the test queries in benchY2test are taken from Wiki-18, but the provided paragraph corpus and legal set of entities (allButBenchmark) are derived from Wiki-16. While automatic entity assessments can be extended by applying an entity linking tool (Figure 4, bottom, left), the paragraphs cannot be automatically aligned. (We experimented to re-align paragraphs using ROUGE, but were not convinced by the results.)

We asked assessors to annotate paragraphs (and entities) from the original page in addition to participant contributed paragraphs (and entities). See Figure 4, top. The motivation is to display positive examples to the assessors to retain high standards even when contributed runs were not containing relevant entries. Only paragraphs (and entities) that were both manually assessed and contained in the paragraph corpus (and legal set of entities) were used in the evaluation.



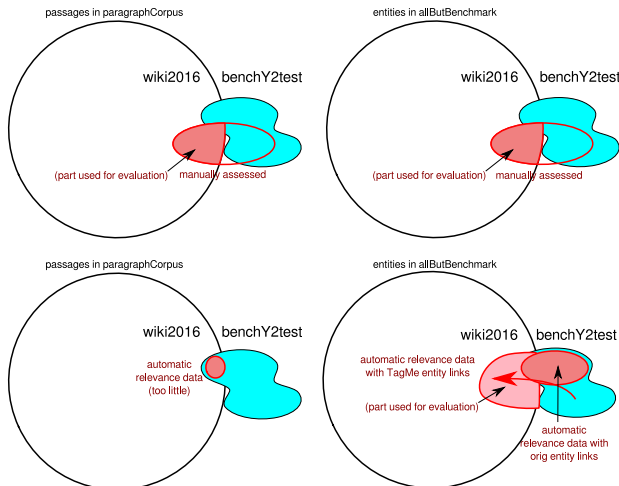


Figure 4: Manual (top) versus automatic (bottom) benchmarks for passages (left) and entities (right). The black circle marks the paragraph/entity collection, the blue area the gold articles, and the red circle marks the set of assessed paragraphs/entities. The light pink area depicts entities determined as relevant via TagMe entity links. The automatic passage collection has too few shared entries and is therefore not included in this study.

## 7 Participant Submitted Runs

Nine teams contributed runs. Most methods are based on Lucene’s BM25 ranking model as a candidate method. Some methods used entity linking, bigrams, pre-trained word vectors, and other forms of query expansion. Table 7 gives an overview over the kinds of methods contributed by participating teams. Below, detailed descriptions of submitted runs provided by teams:

- The `guir` team submitted two passage retrieval runs. The first run (`guir`) is an approach based off prior work [6], in which an ad-hoc neural ranking architecture is modified for the CAR task by incorporating heading frequency statistics from the training data and by incorporating separate matching phases for each heading in the query. The second run (`guir-exp`) uses the same approach, and adds the top-scoring query expansion terms for each heading as determined by a learned match-gating mechanism. The models were trained on the automatic relevance judgments from the train dataset, and validated using the manual relevance judgments on the benchmarkY1-test dataset.
- NYU: Lucene is the underlying retrieval engine. We train 20 query reformulators on random disjoint partitions of the training set as in Nogueira et. al., 2018 [7]. For each query, each reformulator produces a list of ranked documents. We re-rank the union of these 20 lists using a simple ranking model that scores each query-document pair using a feed-forward neural network whose input is the concatenation of the average word embeddings of the query and document. To further improve the performance of the system, we train an ensemble of 9 ranking models whose network architectures are randomly chosen. For each query, we re-rank the union of the 9 lists produced by the 9 ranking models using the best ranking model in the ensemble. All the models are trained on the first 4 folds of TREC-CAR v2.1 queries and the last fold is used for dev/hyperparameter tuning. Y1 test queries are used as test queries. We use paragraphCorpus v2.0 to retrieve documents for train, dev, and test queries.
- CG: We developed a deep learning model for information retrieval TREC Complex Answer Retrieval (CAR) task. We used an attention based sequence-to-sequence model to first translate content passages to outlines. Then across all extracted outlines, use sentence embedding to rank outlines based on the available query. We used attention based bidirectional LSTM model for encoder and decoder layers. In order to capture rare words while limiting our dictionary size, we used Byte Pair Encoding subword units to tokenize sentences. The main advantage of using seq2seq model is to preform main inference

	guir	NYU-DL	CGNADAA	UTDHILTRI	uogTr	CUIS	UMass	TREMA-UNH	DWS-UMA
Pre-trained word embeddings	X	X		X			X		X
Neural network technology	X	X	X	X					
Learning to rank	X	X	X		X		X		
Uses allButBenchmark					X	X	X	X	
Hierarchy in outlines					X		X		
Entity linking					X	X	X	X	
BM25		X		X		X		X	X
SDM					X	X	X	X	
Language model	X				X		X	X	
Uses query expansion	X	X			X	X	X	X	
Passage Task	X	X	X	X	X	X	X	X	X
Entity Task					X	X		X	X

Table 7: Participant contributed runs in CAR Year 2.

computation offline and only use the model to re-rank all outlines based on the query in inference time which is much faster compare to recent neural IR models. Using attention based model also provides position-dependent information required to assess the relevance of a snippet of a document to a given query. Attention signals illustrate term dependencies between query and given passage. The model is trained on TREC V2 data set which has 50% of Wikipedia Articles and test on TREC benchmarkY1 which is official evaluation topics for the TREC CAR task.

- UTD: We extended our approach from last year to create the TRANSformer Complex Answer Paragraph Retrieval (TRANS-CAPAR) system to perform complex answer retrieval consisting of the following five modules: (1) The Paragraph Indexing Module creates a searchable index of paragraphs from Wikipedia articles; (2) The Query Processing Module processes a Wikipedia article outline into a set of queries - one for each section of the outline; (3) The Paragraph Search Module searches each query against the paragraph index, resulting in a list of relevant paragraphs for each section in the article outline; (4) The Feature Extraction Module is used to extract features from each paragraph; (5) The Paragraph Ranking Module produces a separate ranking of the retrieved paragraphs for each section. To calculate relevance scores for each paragraph we use a neural relevance model that combines dynamic IR features with two semantic matching networks that capture complementary relevance signals. One is based on the Transformer sequence-to-sequence model and the other uses the cosine similarity matrix. The model was trained on train.v2.0 dataset and validated using the test200.v2.0 dataset.
- uog: Our 2018 runs study entity-aware expansion models tailored to the TREC CAR task. In particular, we use fine-grained expansion features based on heading components of the TREC CAR query topics. We employ query entity linking, entity retrieval, and performed expansion over diverse matching vocabularies (words, entity IDs, aliases). For entity retrieval we experimented with using feedback on the paragraph collection. All methods parameters (hyper-parameters and LTR models) used Y1Train (hierarchical qrels) and runs were selected by performance on Y1Test (tree qrels).
- CUIS: Team CUIS provided two passage runs and three entity runs. The system consists of two stages. The first stage chooses top 1000 candidate passages based on Lucene’s BM25 method. The second stage reranks these candidate passages using a Markov random field based model where unigrams,

bigrams and concepts induced by query terms from different query sections are considered. Besides, the system incorporates the Wikipedia article information or query entity mentions based on a Dirichlet prior smoothed language model. The run “CUIS-F150” incorporates the Wikipedia article information. The run "CUIS-MX5" incorporates both Wikipedia article and query entity mentions. The system is trained on benchmarkY1-train v2.0 dataset. The entity runs are derived by replacing the paragraph id with the containing Wikipedia article id.

- UMass (entityEmbedLambdaMart): For the lambdamart run we use the benchmarkY1-train.v2.0 as the training set, benchmarkY1-test-public.v2.0 as the validation set and benchmarkY2test-public.v2.1.1 evaluation topics for the submitted run. In this run, we learned a joint entity-word embedding representation based on the Wikipedia corpus. In TREC CAR topics, each query topic consists of three subtopics: Root subtopic (R), Intermediate subtopic (I) and leaf subtopic (H). We retrieve a set of documents with three baseline methods: SDM (Sequential Dependence Model), RM3 and query likelihood with the subtopic combinations of R-H, R-I-H and R-I-H, respectively. Furthermore, We represent each document and query based on their entity embeddings. Each query has fine-grained subtopic word-vector representations as well as the complete topic representation. To be more specific, we represent an entity by the average embedding vectors of entities only in Root, entities only in Leaf (H), and all of the entities in the topic. We use cosine similarity between the document vector representation and each query representation as a feature in LambdaMart learning-to-rank model as well as the retrieval scores from the base retrieval model.
- TREMA-UNH (UNH): We provided three passage runs and three entity runs that were all based on combination of low-level input runs such as BM25, Query likelihood, SDM, RM3 and Entity Context model [4] with combinations trained with coordinate ascent Learning-to-Rank optimized for MAP. We also experimented with a new “Learning-to-Walk” methods for supervised graph walks. Our best performing passage run is a combination of BM25, Query Likelihood with and without RM3. All combinations were trained on benchmarkY1train, and the best three methods were selected on benchmarkY1test. In our notebook also includes evaluation results on Y1 benchmarks.
- DWS-UMA: Our Trec-CAR submission is a simple unsupervised method. At query time we perform semantic query expansion in combination with term specificity boosting on a Lucene Index. Our model first represents the query by its lemmatized query terms. In the next step, for each term, the query is expanded by including the top k nearest neighbors from a semantic word embedding space. The expanded query is executed against a Lucene Index with BM25. Query terms at lower levels in the outline, i.e., more specific query terms, are boosted and receive a higher weight. For our submission we used a pre-trained embedding space and the value for k is tuned on benchmarkY1-train.

## 7.1 Assessment Interface and Fixed Mistakes in Submitted Runs

The assessment interface was populated by a) pooling the top 5 of all submitted runs and b) paragraphs and entities on the original article (i.e., the automatic ground truth).

Several submitted runs contained mistakes. We fixed those mistakes post-hoc, but we were unable to do so before the assessment. The following issues were fixed post-hoc

- Team NYU assigned all ranked items the score of 1.0, all information was contained in the rank information. We derived corrected runs turning rank information into scores. As a result random 5 elements of their ranking were assessed.
- Several teams submitted illegal entity IDs. Despite recommended otherwise, participants created entity IDs from page names by replacing spaces with %20. This will not address non-ASCII characters such as accents or umlauts. We derived corrected runs post-hoc. As a result, for affected runs entities containing accents or umlauts were not assessed.
- Team CG submitted rankings that only contained three passages. As a result, fewer elements were assessed for CG than other teams.

## 8 Results

The official evaluation<sup>4</sup> of participant-contributed is conducted with respect to four standard TREC evaluation measures, R-Precision (RPrec), Mean-average Precision (MAP), Reciprocal Rank (MRR), and Normalize Discounted Cumulative Gain (NDCG). Of these measures only NDCG considers the graded scale, for all other methods the positive/negative cutoff indicated in Table 3 are used.

### 8.1 Passage Task

We evaluate participant-contributed passage runs on manual assessments, since automatic assessments are not available. As described in Section 6.4, while the manual assessments included paragraphs within and outside the paragraph collection, we evaluate participating runs only on passages that are included in the paragraph corpus.

Results for the benchY2test passage retrieval task are presented in Figure 5 on the manual graded scale, and a lenient variant of the manual graded scale. Standard error bars and paired-t-test with respect to the best performing method are given for reference. All analyses across all measures are painting the same picture. Acknowledging consistent patterns in the results for different metrics, here the ranking of passage methods by across-the-board performance (tied methods on the same rank):

1. uog-heading-rh-sdm, UNH-p-l2r
2. uog-linear-lrt-hier, UNH-p-sdm, UNH-p-mixed
3. entityEmbedLambdaMart (UMass), guir-exp, UTDHLTRI2
4. guir, uog-linear-raw-expansion, NYU-XL-f, CUIS-MX5
5. NYU-L-f, NYU-M-f, CUIS-F150
6. DWS-UMASemQueryExp, DWS-UMASemQueryExp20, DWS-UMASemQueryExp30
7. CG-Seq2Seq

To study whether the differences are due to better performance on easy queries, difficult queries, or overall, we include divide the set of all annotated topic sections into percentiles ranging from easy to difficult according to the best performing method. The results are presented in Figure 9 and show a consistent ranking of methods for difficult and easy queries (interquartile ranges 25%-50%, 50%-75%, and 75%-95%).

Separating results for queries originating from Wiki-18 and TQA (Figures 5d and 5e), demonstrates that methods performing well on Wiki-18, also perform well on TQA. However, the difference between methods is less pronounced for the Wiki-18 subset. TQA queries seem to be slightly easier, which is probably because the TQA outlines contain fewer sections.

### 8.2 Entity Task

We evaluate participant-contributed entity runs on automatic and manual assessments. While the manual assessments included entities within and outside the legal entity set (allButBenchmark), we evaluate participating runs only on entities that are included in legal set (cf. Section 6.4) We further evaluate on the automatic benchmark derived with TagMe entity links on ground truth pages. (The previous method of using only entity links included manually by Wikipedia editors does not apply to TQA articles.)

Results for the benchY2test entity retrieval task are presented in Figure 6 on the manual graded scale, lenient variant of the manual graded scale, and automatic assessment based on TagMe entity links, including error bars and paired-t-tests. Once more, all analyses across different measures are providing a coherent picture, resulting in the following ranking of entity methods (tied methods on the same rank):

1. UNH-e-L2R, UNH-e-mixed, UNH-e-graph, uog-rf-ent
2. uog-linear-ltr-hier-ent, uog-heading-rh-sdm-ent, DWS-UMA-AspQLrm, DWS-UMA-EntAspBM25none

---

<sup>4</sup>Official results available <http://trec-car.cs.unh.edu/results/>

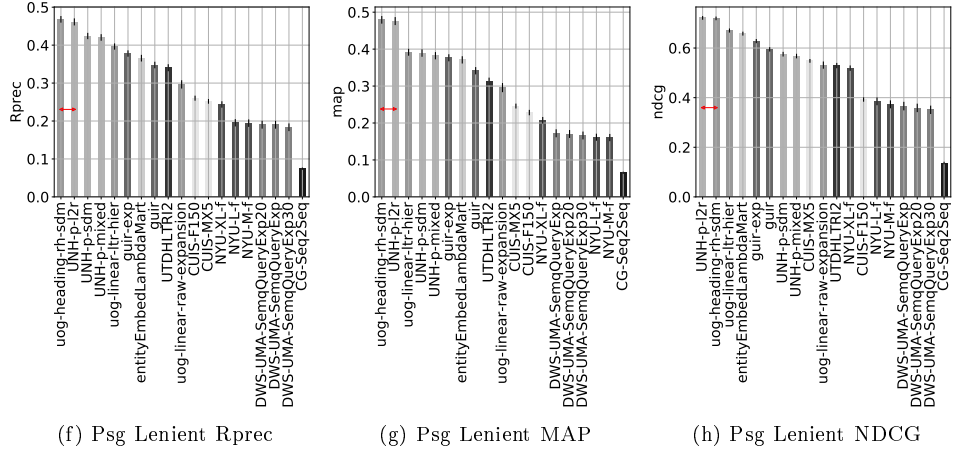
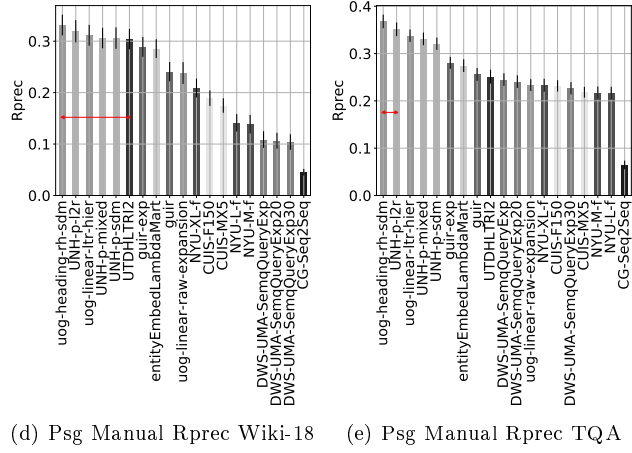
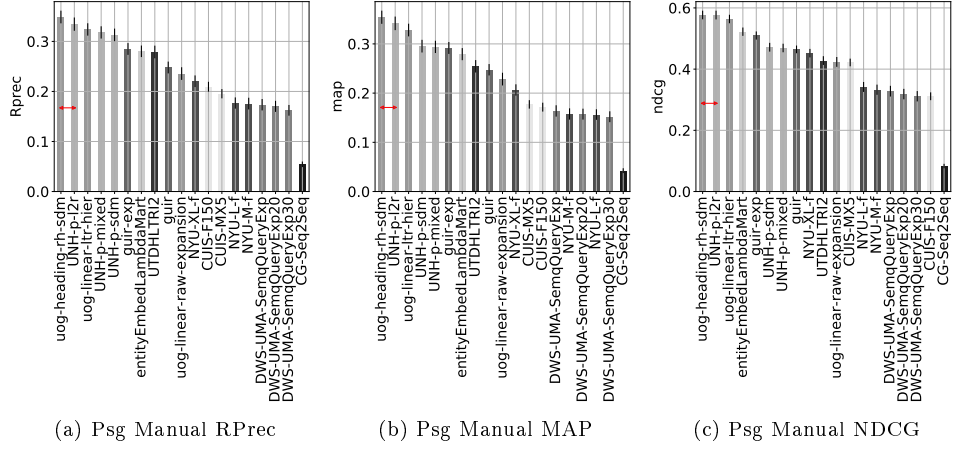
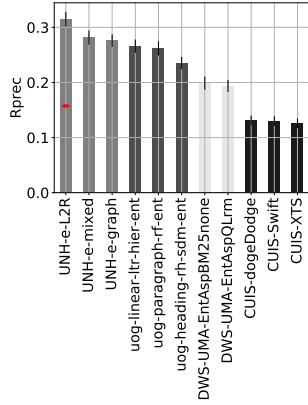
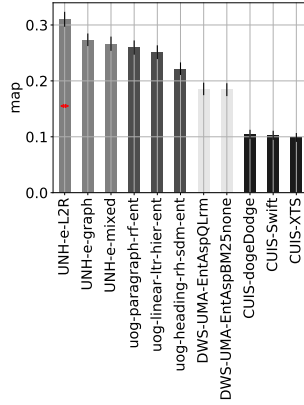


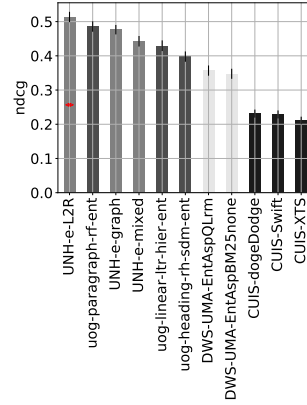
Figure 5: Results of contributed passage runs under the manual ground truth. Lenient is based on the manual graded scale, but counting TOPIC as relevant. The red arrow marks systems for which no significant difference to the best system could be detected.



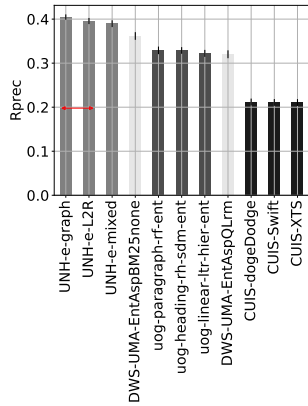
(a) Entity Manual Rprec



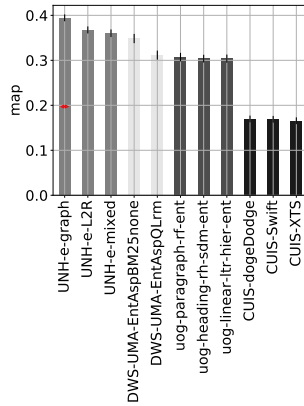
(b) Entity Manual MAP



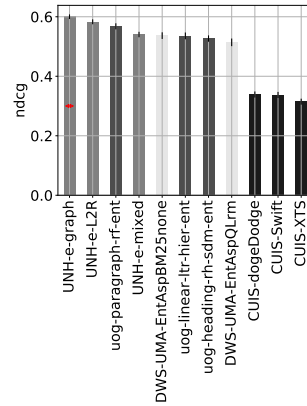
(c) Entity Manual NDCG



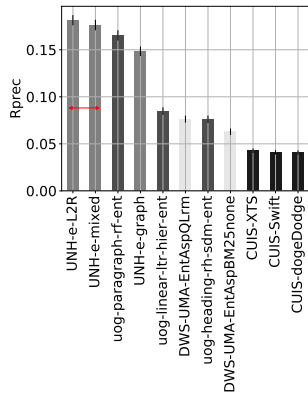
(d) Entity Lenient Rprec



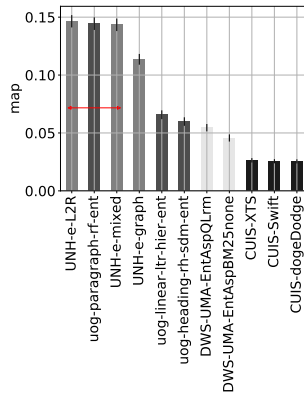
(e) Entity Lenient MAP



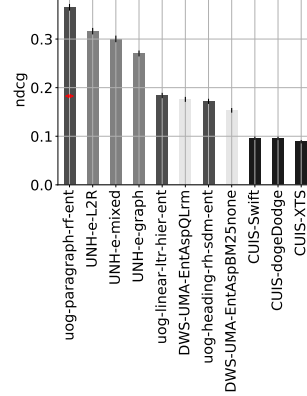
(f) Entity Lenient NDCG



(g) Entity Automatic Rprec



(h) Entity Automatic MAP



(i) Entity Automatic NDCG

Figure 6: Results of contributed entity runs under automatic and manual ground truth. Lenient is based on the manual graded scale, but counting TOPIC as relevant. The red arrow marks systems for which no significant difference to the best system could be detected.

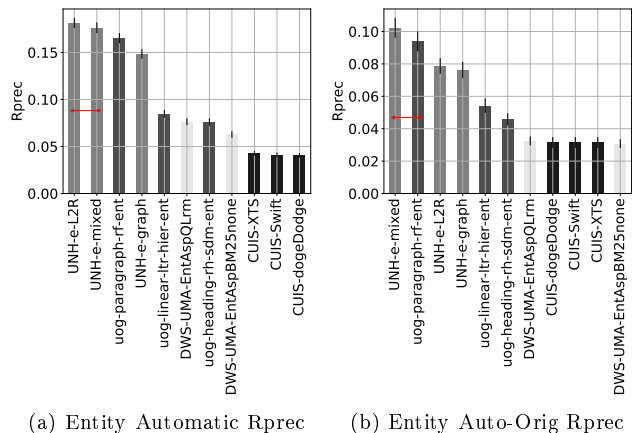


Figure 7: Automatic evaluation using entity links from TagMe versus evaluation using links in Wikipedia source (Auto-Orig, only available for Wiki-18). The red arrow marks systems for which no significant difference to the best system could be detected.

### 3. CUIS-dodgeDodge, CUIS-Swift, CUIS-XTS

Comparing results under the automatic entity benchmark (using TagMe) with the automatic entity benchmark derived from entity links included by Wikipedia editors (Auto-Orig, provided as training data for the “train” benchmark), we see a similar pattern emerging, but some systems swap ranks.

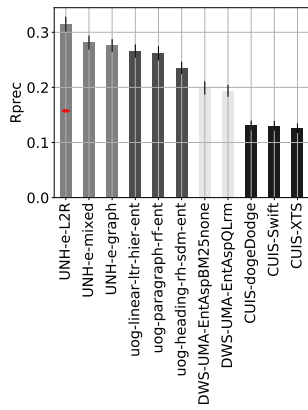
When analyzing queries originating from Wiki-18 and TQA presented in Figure 8, we observe less sharp distinctions on Wiki-18 (with two or three systems tied for rank 1), where on TQA, the method UNH-e-L2R is consistently leading with significant difference to other methods.

## 9 Conclusion

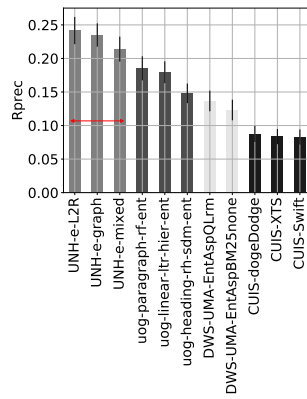
In contrast to the previous year of TREC CAR, where neural network methods were dominating the leaderboard, in this second year we see that learning-to-rank with unsupervised retrieval models such as BM25, SDM and query expansion significantly outperformed state-of-the-art neural methods. We see that systems that perform well in passage retrieval tasks, are also performing well in the entity retrieval task. Regarding benchmark construction for this project. In previous year, we confirmed that automatically derived relevance data for passages agrees with human-created benchmarks on the ranking of systems. In this year we further show that automatically derived relevance data for entities agrees with human-created benchmarks. This means, that we have an effective test bed for method development in TREC CAR, when the outline is provided. In the next year of TREC CAR, we will move beyond population of existing outlines and leave it to participants to also identify a suitable ordering of paragraphs to automatically populate a complete article, given only a suitable title.

## Acknowledgement

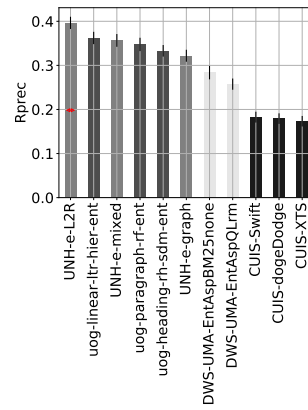
We express our gratitude for many suggestions of several experts in the field, who helped to make this track successful. We thank the University of New Hampshire for providing computational resources and web servers. We are deeply thankful for Ellen Voorhees’ experience, patience, and persistence in running the assessment process. Finally we thank all our participants for their contributions.



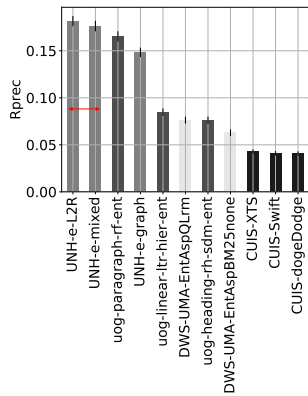
(a) Entity Manual Rprec



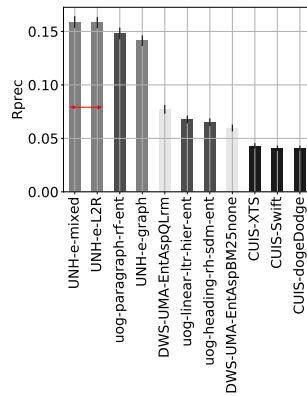
(b) Entity Manual Rprec Wiki-18



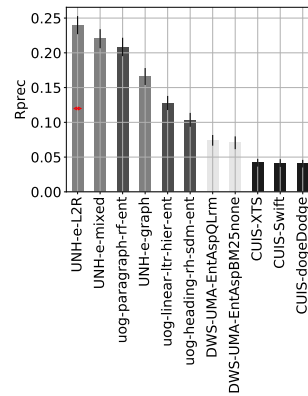
(c) Entity Manual Rprec TQA



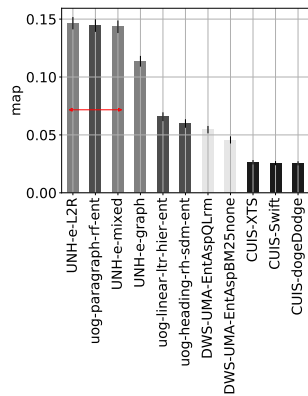
(d) Entity Automatic Rprec



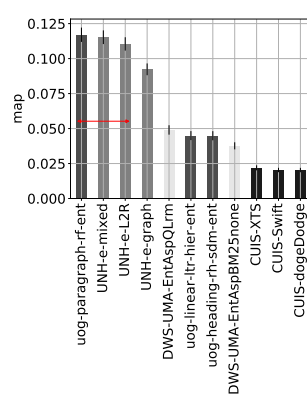
(e) Entity Auto Rprec Wiki-18



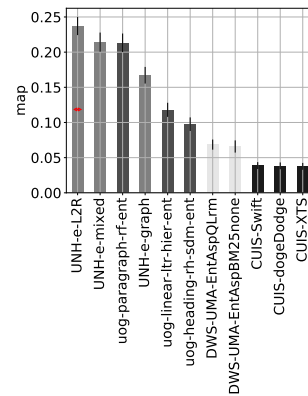
(f) Entity Auto Rprec TQA



(g) Entity Automatic MAP



(h) Entity Auto MAP Wiki-18



(i) Entity Auto MAP TQA

Figure 8: Entity performance on Wiki-18 versus TQA. The red arrow marks systems for which no significant difference to the best system could be detected.



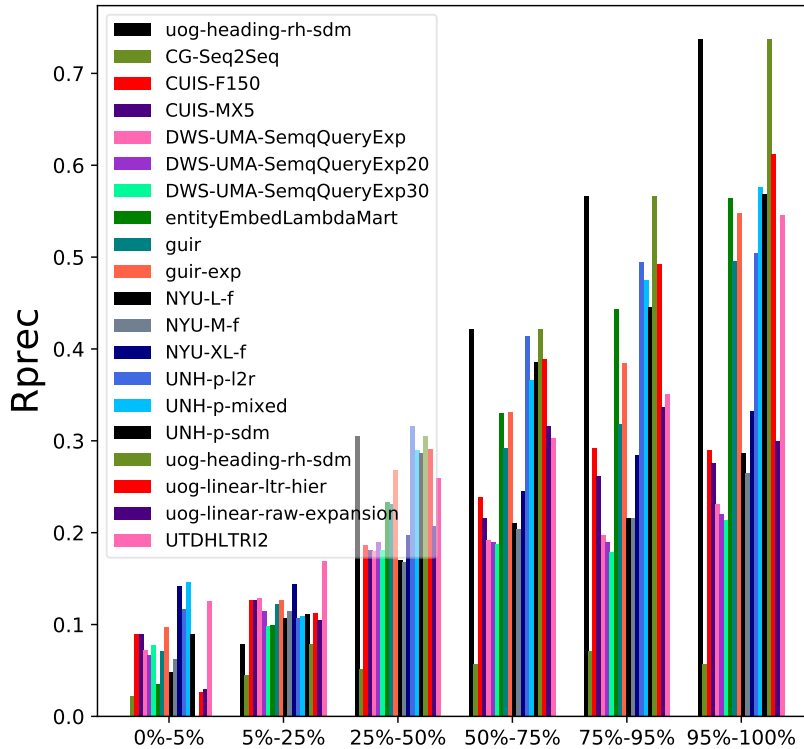


Figure 9: Performance according to manual truth on difficulty percentiles according to best performing method.

## References

- [1] James Allan, W. Bruce Croft, Alistair Moffat, and Mark Sanderson. Frontiers, challenges, and opportunities for information retrieval: Report from SWIRL 2012 the second strategic workshop on information retrieval in Lorne. *SIGIR Forum*, 46(1):2–32, 2012. doi: 10.1145/2215676.2215678. URL <http://doi.acm.org/10.1145/2215676.2215678>.
- [2] Omar Alonso and Stefano Mizzaro. Using crowdsourcing for trec relevance assessment. *Information Processing & Management*, 48(6):1053–1066, 2012.
- [3] J Shane Culpepper, Fernando Diaz, and Mark D Smucker. Research frontiers in information retrieval: Report from the third strategic workshop on information retrieval in lorne (swirl 2018). In *ACM SIGIR Forum*, volume 52, pages 34–90. ACM, 2018.
- [4] Jeffrey Dalton, Laura Dietz, and James Allan. Entity query feature expansion using knowledge base links. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*, pages 365–374. ACM, 2014.
- [5] Paolo Ferragina and Ugo Scaiella. Tagme: on-the-fly annotation of short text fragments (by wikipedia entities). In *Proceedings of the 19th ACM international conference on Information and knowledge management*, pages 1625–1628. ACM, 2010.
- [6] Sean MacAvaney, Andrew Yates, Arman Cohan, Luca Soldaini, Kai Hui, Nazli Goharian, and Ophir Frieder. Characterizing question facets for complex answer retrieval. In *Proceedings of the 41st International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2018.
- [7] Rodrigo Nogueira, Jannis Bulian, and Massimiliano Ciaramita. Learning to coordinate multiple reinforcement learning agents for diverse query reformulation. 2018.