



# Evaluating Entity Disambiguation and the Role of Popularity in Retrieval-Based NLP

Anthony Chen\*  Pallavi Gudipati  Shayne Longpre   
Xiao Ling  Sameer Singh 

 University of California, Irvine  Apple  
{anthony.chen, sameer}@uci.edu  
{pgudipati, slongpre, xiaoling}@apple.com

## Abstract

Retrieval is a core component for open-domain NLP tasks. In open-domain tasks, multiple entities can share a name, making disambiguation an inherent yet under-explored problem. We propose an evaluation benchmark for assessing the entity disambiguation capabilities of these retrievers, which we call *Ambiguous Entity Retrieval (AmbER) sets*. We define an AmbER set as a collection of entities that share a name along with queries about those entities. By covering the set of entities for polysemous names, AmbER sets act as a challenging test of entity disambiguation. We create AmbER sets for three popular open-domain tasks: fact checking, slot filling, and question answering, and evaluate a diverse set of retrievers. We find that the retrievers exhibit popularity bias, significantly under-performing on rarer entities that share a name, e.g., they are twice as likely to retrieve erroneous documents on queries for the less popular entity under the same name. These experiments on AmbER sets show their utility as an evaluation tool and highlight the weaknesses of popular retrieval systems.<sup>1</sup>

## 1 Introduction

Substantial progress in NLP has been made on “closed” tasks, where queries are paired with relevant documents (Rajpurkar et al., 2016; Dua et al., 2019). However, there is growing interest in “open-domain” tasks, where relevant documents need to be retrieved from a knowledge source before an NLP system can perform reasoning and produce an answer (Chen et al., 2017; Petroni et al., 2021). The open-domain setting better reflects real-world usage for tasks where relevant information is generally not provided (e.g., fact checking).

\*Work started during an internship at Apple.

<sup>1</sup>The AmbER sets used in this paper and the code to generate them are available at <https://github.com/anthonywchen/AmbER-Sets>.

**Q:** Which battle did **Abe Lincoln** fight in?

**A:** World War II

**Wikipedia Documents Ranked by BLINK:**

1. **Abraham.Lincoln**
2. Abraham.Lincoln.in.the.Black.Hawk.War
3. Abraham.Lincoln.(captain)
4. Benjamin.Lincoln
5. Lincoln.Nebraska
6. Lincoln.England

**Q:** What musical instrument does **Abe Lincoln** play?

**A:** Trombone

**Wikipedia Documents Ranked by BLINK:**

1. Abraham.Lincoln
2. John.Wilkes.Booth
3. Abe.(musical)
4. Nebraska
5. Lincoln.Nebraska
6. **Abe.Lincoln.(musician)**

Figure 1: Queries for two entities (**president** & **musician**) with the name “Abe Lincoln”. Retrieving the **gold document** involves disambiguating which “Abe Lincoln” each query is asking about. BLINK performs sub-optimally on the second query, as it ranks the document of the president over the gold document.

Because success hinges on finding relevant documents, open-domain progress has been closely tied to improvements in retrieval systems<sup>2</sup> (Lee et al., 2019; Karpukhin et al., 2020; Lewis et al., 2020b).

A crucial challenge when interacting with a large knowledge source (e.g., Wikipedia) is entity ambiguity, the phenomenon where a single name can map to multiple entities. Resolving this ambiguity is referred to as entity disambiguation and is an important step for effective retrieval. For example, given the query “What musical instrument does Abe Lincoln play?”, documents about the musician should rank higher than other entities with the same name (Figure 1). Although entity disambiguation has been extensively studied in entity linking (Hof-fart et al., 2011; Rao et al., 2013; Sevgili et al.,

<sup>2</sup>For example, replacing the BM25 retriever with DPR on Natural Questions increases exact match by 15 points.

2020) and search (Balog et al., 2010, 2011), in the context of open-domain NLP, it is unclear how good retrieval systems are when faced with queries with ambiguous entities. Evaluating entity ambiguity is challenging because the popularity of entities follows a long-tail (Figure 2) and rare entities are seldom covered in naturally-occurring datasets.

In this paper we introduce AmbER sets, a benchmark for evaluating the entity disambiguation capabilities of retrievers across multiple NLP tasks. Each AmbER set is a collection of Wikidata entities that share a name, and their corresponding queries for specific NLP tasks. For each set, we define the **head** entity as the most popular entity and **tail** entities as the less popular ones. By creating queries for multiple entities that share a name, AmbER sets provide an accurate test of entity disambiguation capabilities of retrievers and help assess the role of entity popularity in disambiguation. We show examples of AmbER sets for the question answering task in Table 1. We automatically create AmbER sets by mining the Wikidata knowledge graph (Vrandečić and Krötzsch, 2014) for relevant names and entities, and leveraging task-specific templates to generate inputs for three tasks: fact checking, slot filling, and question answering (Figure 3). In total, our AmbER sets contain 80k task-specific queries which we align to the Wikipedia snapshot from KILT (Petroni et al., 2021).

We use AmbER sets to conduct a systematic study of various retrieval systems that operate under different principles, such as token overlap and dense embedding similarity. Retrievers perform very differently on AmbER sets in terms of absolute retrieval numbers, with Bootleg (Orr et al., 2020), an entity-linking-based retriever, performing best. Despite these differences, all retrievers exhibit a large degree of popularity bias, underperforming on inputs concerning tail entities. TF-IDF, a token-based retriever, performs about four times worse on tail entity inputs compared to head entity inputs. Even with Bootleg, the best performing retriever, performance on tail entities is still 1.5 times lower than on head entities. Our results on AmbER sets demonstrate that there is significant work to be done on making retrievers robust in handling entity disambiguation.

## 2 AmbER Sets

Retrieving relevant documents from large knowledge sources such as Wikipedia is an important

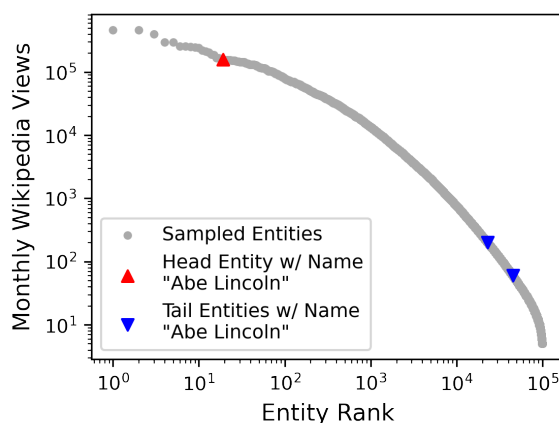


Figure 2: **The Long Tail of Entity Popularity:** Graph of the Wikipedia pageviews (in October 2019) for each Wikidata entity, ranked by popularity. Gray are 100k randomly sampled entities, while red/blue are entities with the name “Abe Lincoln”.

first step in the open-domain pipeline. An inherent problem in working with such sources is entity disambiguation: resolving a name (mention) to an entity in the knowledge source. Entity disambiguation can be challenging because many entities share a name, and the popularity of entities follows a long-tail distribution (Figure 2). Despite the importance of entity disambiguation, it remains an understudied problem for open-domain NLP. We introduce AmbER sets for evaluating entity disambiguation capabilities of retrievers and analyze the role of entity popularity in disambiguation.

### 2.1 What is an AmbER Set?

We first provide an intuition for an AmbER set before concretely defining one. Consider two entities, a president and a musician, both of which have the name “Abe Lincoln” (Figure 1). Now, consider the query “Which battle did Abe Lincoln fight in?” and assume a retriever correctly returns the article about the president for this query. Simply because the correct document was retrieved does not mean a retriever has the ability to disambiguate between the president and the musician, as the president is much more popular. We should only be confident in its ability to disambiguate entities if we *also* pose a query about the less popular musician and the retriever again returns the correct document (as opposed to the document about the president).

Based on this intuition, we define an AmbER set as a collection of queries that satisfy the following:

- **Criteria 1: Polysemous Name:** The queries in an AmbER set are all about entities that share a common name (e.g., Abe Lincoln).

	QID	Input	Answer	Gold Document
AmbER-H	Q517	What wars did <b>Napoleon</b> participate in?	Napoleon Wars	Napoleon
	Q3335909	What sport does <b>Napoleon</b> play?	Rugby	Napolioni_Nalaga
	Q3335909	Which team does <b>Napoleon</b> play for?	Fiji National	Napolioni_Nalaga
	Q117012	What movement did <b>Yoko Ono</b> participate in?	Fluxus	Yoko_Ono
	Q16264827	Which sport does <b>Yoko Ono</b> participate in?	Judo	Yoko_Ono_(judoka)
AmbER-N	Q312	Which industry is <b>Apple</b> in?	Electronics	Apple_Inc.
	Q532100	What is the record label of <b>Apple</b> ?	Page One	Apple_(band)
	Q7714007	Who acted in <b>Apple</b> ?	Ray Shell	The_Apple_(1980_film)
	Q788822	Who is a cast member on <b>Her</b> ?	Steve Zissis	Her_(film)
	Q788822	Who is <b>Her</b> 's screenwriter?	Spike Jonze	Her_(film)
	Q28441308	Who performed <b>Her</b> ?	Aaron Tippin	Her_(song)

Table 1: **Examples of QA AmbER sets.** An AmbER set is a collection of entities that share a name, with instantiated queries for each entity. In this work, we use Wikidata to collect entities (QID). We also create queries for two more tasks, fact checking and slot filling (omitted from this table).

- **Criteria 2: Disparity in Popularity:** An AmbER set contains queries about both the most popular entity for a name (the **head** entity), *e.g.*, the president, and the less popular entities (the **tail** entities), *e.g.*, the musician.
- **Criteria 3: Resolvable Ambiguity:** The content of the query should be sufficient to resolve to the correct entity. The query “Which battle did Abe Lincoln fight in?” satisfies this criteria, because there is only one Abe Lincoln that fought in a war, while “Where was Abe Lincoln born?” does not since it applies to all Abe Lincolns.

We provide examples of AmbER sets for the task of question answering in Table 1.

## 2.2 Open-Domain Tasks

In this work, we create AmbER sets for three tasks: fact checking, slot filling, and question answering (Table 2). We consider these three tasks for three reasons. First, these three set of tasks are diverse in nature. In this work, slot filling is a generation task, question answering is a span selection task, and fact checking is a classification task. Second, the training sets available for each task are quite disparate. The largest fact checking training set, FEVER (Thorne et al., 2018), has 80k instances, while the slot filling dataset, T-REx (Elsahar et al., 2018), has over 2 million instances. The final reason we study these three tasks is that their inputs are short and easy to create.

## 3 Creating AmbER Sets

While AmbER sets can be manually created, doing so can be time-consuming, requiring a human to manually scour a knowledge base for polysemous

Task	Input Instance	Output
FC	John Mayer plays music.	True
SF	Nike [SEP] country	USA
QA	Whose face is on \$100 bill?	Benjamin Franklin

Table 2: Examples for each open-domain NLP task.

names and related entities before manually writing queries for those entities. Instead, we present a pipeline for automatically creating AmbER sets using the Wikidata knowledge graph (Vrandečić and Kröttsch, 2014). In this section, we describe two different *collections* of AmbER sets, and discuss our automatic pipeline for creating AmbER sets.

### 3.1 Two Collections of AmbER Sets

A natural question is “How do retrievers handle entity ambiguity when two entities have the same entity type as opposed when they have different types?”. To answer this question, we create two *collections* of AmbER sets. The first is AmbER-H, a collection of AmbER sets where all entities are humans. The choice to restrict AmbER-H to humans is motivated by the fact that humans have properties that help distinguish themselves from other humans, generally based on occupation. The second is AmbER-N, a collection of AmbER sets where all entities contained are non-humans, and disambiguation of a name is between non-human entities with different entity types. This is because a non-human entity, like a movie, does not generally have a single distinguishing property to distinguish from other movies. This makes it natural to compare non-human entities to other non-human entities with different types. We specify the entity types in each collection in Table 3.

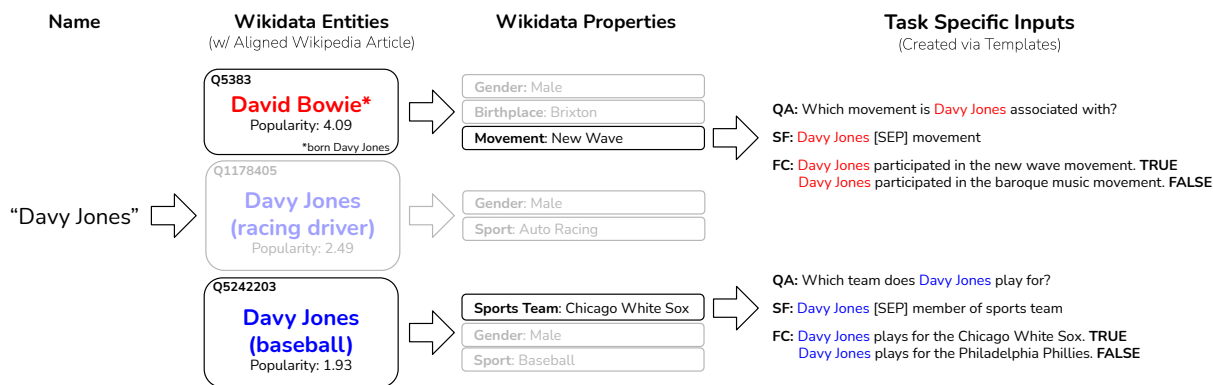


Figure 3: **Automated creation of AmbER sets for three tasks.** We collect sets of entities from Wikipedia that share a name, where the most popular entity is the head entity (in red) and others are tail entities (in blue), along with their properties and associated values. We filter out properties that do not help distinguish entities in the set (gray-ed out), and remove entities that do not have any properties remaining. From the remaining properties, we instantiate queries via templates for three tasks: question answering (QA), slot filling (SF), and fact checking (FC).

### 3.2 Automatic Creation of AmbER Sets

We now describe a pipeline to automatically create AmbER sets for three tasks: fact checking, slot filling, and question answering. We provide a visualization of the pipeline in Figure 3.

**Collecting Names and Entities** We begin by collecting all entity aliases<sup>3</sup> in Wikidata. From these aliases, we filter for those that are shared by multiple Wikidata entities. Each entity in Wikidata is represented by a unique QID. The entities must have an entity type from Table 3 depending on the collection we are collecting AmbER sets for. Each alias and associated entities form the basis for an AmbER set. Within each set, we define the head and tail entities based on the number of Wikipedia page views for the month of October 2019. We filter out AmbER sets where the percentage gap in popularity between the head entity and the most popular tail entity is less than 10% to account for noise in the monthly page views.

**Collecting Distinguishing Properties** We gather properties and associated values for each entity from Wikidata. We only retain properties that are in a specified list (Table 3), as they are useful for resolving ambiguity (*Criteria 3*). We also filter a property if two entities within an AmbER set have that property, ensuring that the remaining properties can be used to disambiguate between entities with the same name. These properties are used to instantiate the queries.

**Aligning Entities to Wikipedia** We use the KILT Wikipedia snapshot (Petroni et al., 2021) as

<sup>3</sup>Aliases are all possible names for an entity.

	Entity Type	Property ( <i>PID</i> )	Percent
AmbER-H	Human	instrument ( <i>P1303</i> )	17.01
		movement ( <i>P135</i> )	2.04
		appears in ( <i>P1441</i> )	0.08
		killed by ( <i>P157</i> )	0.19
		PhD student ( <i>P185</i> )	0.42
		military branch ( <i>P241</i> )	12.22
		sports position ( <i>P413</i> )	12.82
		sports team ( <i>P54</i> )	17.25
		battles or wars ( <i>P607</i> )	12.29
		sport ( <i>P641</i> )	25.68
AmbER-N	Album	performer ( <i>P175</i> )	16.57
		record label ( <i>P264</i> )	7.11
		tracklist ( <i>P658</i> )	0.21
	Business	industry ( <i>P452</i> )	0.65
	City	population ( <i>P1082</i> )	0.24
	Film	cast member ( <i>P161</i> )	27.14
		screenwriter ( <i>P58</i> )	18.28
	Literary Work	author ( <i>P50</i> )	11.13
	Musical Group	record label ( <i>P264</i> )	2.1
	Song	performer ( <i>P175</i> )	4.42
		record label ( <i>P264</i> )	0.62
	TV Series	cast member ( <i>P161</i> )	2.01
		# seasons ( <i>P2437</i> )	1.85
screenwriter ( <i>P58</i> )		0.21	
Written Work	author ( <i>P50</i> )	7.43	

Table 3: **Distinguishing Properties** selected to create queries based on whether they are sufficient to resolve ambiguity. We provide the percent breakdown of how often each property occurs in each AmbER collection.

the knowledge source for AmbER sets for better reproducibility. Each Wikipedia document in KILT has an associated QID. For each entity, we find all Wikipedia documents with that associated QID. After this alignment, we apply a round of filtering on the tuples. For each tuple, we check that the value of the tuple is within the first 350 tokens of the aligned Wikipedia article. If not, we remove

	AmbER-H	AmbER-N
# AmbER Sets	2,093	5,237
Averages per AmbER Set		
...# entities	2.98	2.42
...# entities w/ properties	2.03	2.06
...# properties	2.84	2.64
# Input Queries	23,768	55,216
... Question Answering (QA)	5,942	13,804
... Slot Filling (SF)	5,942	13,804
... Fact checking (FC)	11,884	27,608

Table 4: Statistics of AmbER collections.

the tuple.<sup>4</sup> Aligned Wikipedia articles that contain the tuple value serve as gold documents.

**Instantiating AmbER Instances** Recall that our goal was to create AmbER sets for three tasks: fact checking, slot filling, and question answering. We are able to create queries for all three tasks simultaneously using the collected Wikidata tuples. For question answering and fact checking, we use templates based on properties to instantiate inputs. Three of the authors wrote a template each for each property for the two tasks. Duplicate templates are removed, resulting in an average of 3 question answering templates per property and 2.7 fact checking templates per property. See Appendix B for the complete list of templates.

For slot filling, we create a single input from each Wikidata tuple by concatenating the AmbER set name with the property name, and using the value of the tuple as the answer. For question answering, we also create a single input for each tuple by filling in the template with the AmbER set name and using the value of the tuple as the answer. For fact checking, we create two inputs for each tuple, one claim that is true using the tuple value and one claim that is false. The false claim is created by finding the most popular value for the tuple property that does not match the tuple value<sup>5</sup>.

### 3.3 Dataset Statistics

We provide statistics for AmbER sets in Table 4. On average, each AmbER set has about three entities that share the same name. Of these three entities, on average, only two have properties after filtering. In total, our AmbER sets contain about 80k task-specific input queries.

<sup>4</sup>This reduces the number of tuples for AmbER-H from 17,079 to 5,942 and for AmbER-N from 22,219 to 13,804.

<sup>5</sup>The most popular instrument in Wikidata is piano. Therefore, given the true claim “*Abe Lincoln played the trombone.*”, the false claim would be “*Abe Lincoln played the piano.*”.

### 3.4 Limitations

Since our pipeline is automated and relies on Wikipedia and Wikidata, there are a few limitations worth noting. AmbER sets will be affected by incompleteness of the knowledge source, sometimes resulting ambiguous queries if a property is missing from Wikidata, but answerable from Wikipedia text. For this reason, we only select a few properties for each type (Table 3). Second, even though we author multiple templates for each property, the reliance on these templates limits the syntactic diversity in the queries (not a critical concern, since we are only evaluating existing models). Also, we use Wikipedia page views as a proxy for real-world popularity of entities. Defining popularity in this way may be problematic, as page views for an entity can fluctuate, and may make our pipeline difficult to generalize to other knowledge sources, where this information may not be available.

Several design choices in creating AmbER sets are worth further investigation. We limit AmbER sets to a pre-specified list of entity types and properties to ensure that entities in an AmbER set are distinguishable. This precludes other properties that may be useful in distinguishing entities, reducing the diversity in AmbER sets. Another design choice is we allow any alias in Wikidata to form an AmbER sets, however, not all aliases are canonical ways to refer to the entity. For instance, Shaquille O’Neal has the unusual alias “The Big Cactus”, potentially leading to a somewhat unrealistic query “*What sport did The Big Cactus play?*”. We plan to revisit these design choices in future work.

## 4 Evaluation Setup

**Retrieval Systems** The primary focus of this work is to evaluate entity ambiguity of retrieval systems. We consider four retrievers based on different retrieval paradigms. The first three are TF-IDF, a token-based retriever using sparse embeddings, DPR (Karpukhin et al., 2020), a dense embedding based retriever, and BLINK (Wu et al., 2020), a linker-based retriever which ranks documents based on input entities. These three retrievers have been thoroughly evaluated on a number of open-domain tasks in Petroni et al. (2021) with no obvious winner across tasks. Encouraged by the disambiguation success on rare entities by Orr et al. (2020), we also evaluate a retriever based on Bootleg, another entity linker. We provide additional details about these retrievers in Appendix D.

Collection	Retriever	Fact Checking (FC)				Slot Filling (SF)				Question Answering (QA)			
		All	Head	Tail	∇	All	Head	Tail	∇	All	Head	Tail	∇
AmbER-H	TF-IDF	17.3	28.5	8.2	0.0	18.8	31.9	8.1	0.0	16.7	28.2	7.3	0.1
	DPR	18.1	23.9	13.3	0.1	8.0	11.6	5.1	0.3	13.1	19.6	7.9	1.1
	BLINK	<b>55.9</b>	<b>64.4</b>	<b>49.0</b>	<b>5.6</b>	38.2	57.0	22.9	11.5	31.7	40.5	24.6	6.6
	Bootleg	34.8	43.0	28.2	0.7	<b>56.5</b>	<b>63.9</b>	<b>50.6</b>	<b>25.3</b>	<b>67.2</b>	<b>77.1</b>	<b>59.1</b>	<b>36.1</b>
AmbER-N	TF-IDF	9.4	13.6	4.9	0.0	13.4	21.0	5.2	0.2	13.9	21.7	5.4	0.3
	DPR	<b>36.9</b>	<b>48.0</b>	<b>24.8</b>	<b>4.4</b>	29.9	40.9	18.0	6.0	36.2	49.2	22.2	9.3
	BLINK	11.7	13.9	9.4	0.0	5.7	7.3	3.9	0.7	35.2	44.7	24.9	10.1
	Bootleg	3.5	4.6	2.4	0.0	<b>52.3</b>	<b>61.3</b>	<b>42.5</b>	<b>22.4</b>	<b>59.8</b>	<b>69.5</b>	<b>49.3</b>	<b>29.0</b>

Table 5: **Top-1 retrieval results** on each collection of AmbER sets. We report accuracy@1 results on all instances as well as results on instances about head entities and instances about tail entities. We also report a set-level metric, *all correct* (∇), the percentage of AmbER sets where *all* inputs had the correct document retrieved.

		FC		SF		QA	
		Head	Tail	Head	Tail	Head	Tail
<i>H</i> *	TF-IDF	19.5	67.5	28.2	75.7	27.9	76.1
	DPR	1.2	10.0	2.3	23.8	2.6	27.0
	BLINK	9.8	32.2	14.0	58.2	4.4	27.6
	Bootleg	6.2	24.7	9.3	30.5	3.7	28.7
<i>N</i> *	TF-IDF	10.1	49.9	22.0	76.9	23.0	76.8
	DPR	6.2	32.2	9.1	48.3	8.7	44.0
	BLINK	5.8	22.8	5.1	32.2	5.5	31.9
	Bootleg	7.7	26.1	16.1	36.2	7.8	31.6

\* *H* represents AmbER-*H* and *N* represents AmbER-*N*.

Table 6: **Entity confusion** measures the % of queries the gold document ranks **worse** (lower) than a document for another entity with the same name (*i.e.*, another entity in the AmbER set). Retrievers are four times as likely to exhibit this when dealing tail queries.

**Downstream Models** The dominant approach to open-domain tasks is a two-stage process where a retriever first finds relevant documents, followed by a downstream model that processes these documents to produce an answer. We evaluate the end-to-end performance on AmbER sets by training downstream NLP models on our tasks of interest. For fact checking, we fine-tune a BERT classifier (Devlin et al., 2019) on FEVER (Thorne et al., 2018). For question answering, we fine-tune a RoBERTa model (Liu et al., 2019) on Natural Questions (Kwiatkowski et al., 2019). For slot filling, a generation task, we fine-tune a BART model (Lewis et al., 2020a) on T-Rex (Elsahar et al., 2018). We provide example training instances in Table 2 and additional details on the models in Appendix E. We use the AllenNLP and HuggingFace Transformers library to finetune our downstream models (Gardner et al., 2018; Wolf et al., 2020).

## 5 Results

In this section, we evaluate existing open-domain NLP pipelines using AmbER sets. We also conduct

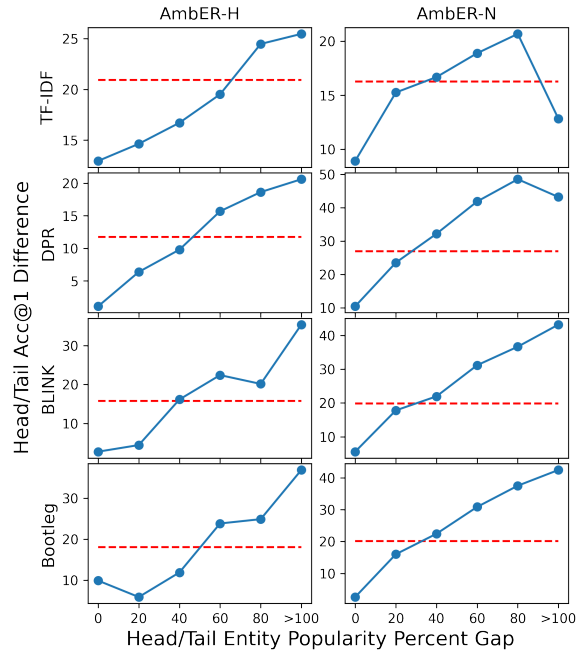


Figure 4: **Popularity Gap vs Retrieval Gap.** We bin QA queries of pairs of head and tail entities based on the popularity gap between the entities. For each bin, we calculate the retrieval accuracy@1 difference on the head and tail queries. Larger popularity gaps tend to lead to a wider gaps in retrieval performance. The red line is retrievers’ performance gaps between head and tail queries on the entire collection.

a user study to evaluate the quality of the queries in the AmbER sets.

**Top Document Retrieval** We report retrieval performance in Table 5 in terms of retriever accuracy@1 (the % of instances where the first retrieved document is the gold document). For each task, we report values on the entire AmbER set (“All”), as well as instances corresponding only to “Head” entities or to “Tail” entities. We also report a metric we call *all correct* (∇), the fraction

Task	System	Results			
		All	Head	Tail	
<i>H</i>	FC	BERT (Oracle)	77.7	73.6	80.3
		BERT + BLINK	59.8	60.1	57.7
	SF	BART (Oracle)	83.9	85.0	83.5
		BART + BLINK	34.4	38.2	32.6
	QA	BERT (Oracle)	71.4	77.7	83.0
		BERT + BLINK	27.5	33.8	22.3
<i>N</i>	FC	BERT (Oracle)	66.6	63.9	69.5
		BERT + DPR	60.9	61.4	60.4
	SF	BART (Oracle)	82.1	80.1	84.3
		BART + DPR	18.6	18.6	18.6
	QA	BERT (Oracle)	83.5	85.1	81.8
		BERT + DPR	26.0	31.3	20.4

Table 7: **End-to-end performance on AmbER sets.** We evaluate systems in an oracle setting, where the gold document is provided, and a retrieval setting, where 20 documents are provided from a retriever.

of AmbER sets in which all queries had the correct document retrieved. All retrievers do better on head entities compared to tail entities. Since BLINK, Bootleg, and DPR are initialized using pre-trained language models, they may have a predisposition towards being biased to more popular entities. However, we find TF-IDF also does better on head entities, perhaps because more popular entities have longer Wikipedia pages, possibly increasing term-frequency scores. Second, there are large discrepancies between a retriever’s performance on different tasks for an AmbER collection. For instance, DPR does substantially worse on slot filling compared to its performance on question answering. This is surprising since queries for all tasks are created from the same set of Wikidata tuples. Finally, we find that retrievers are mostly incorrect on getting all the queries in a set correct, with some receiving a  $\forall$  score of 0 on some tasks. Overall, we find that the Bootleg retriever on average does the best across tasks, however there is significant scope for improvement.

**Entity Confusion** To explicitly evaluate whether retrievers get confused by entities in the same AmbER set, we compute *entity confusion* for retrievers defined as the percentage of queries where the retriever ranks a document for an incorrect entity from the same AmbER set over the gold document (Table 6). We find that across retrievers, tasks, and AmbER collections, entity confusion is twice as high for tail entity inputs. This result indicates that the popularity of an entity for a given name plays a significant role in retrieval performance.

**Effect of Popularity Gap** Since the difference in popularity between the head and tail entities can vary considerably, these results obfuscate the effect of the *size* of the popularity gap. We explore how the gap in popularity between head and tail entities translates to the gaps in performance on their associated queries. For a head entity with popularity  $p_h$  and a tail entity with popularity  $p_t$  from the same AmbER set, we calculate popularity gap,  $\frac{p_h - p_t}{p_t}$ , and bin associated head/tail inputs based on the gap<sup>6</sup>. For each bin, we calculate the difference in accuracy@1 between the head and tail entity queries. Results for QA AmbER sets (Figure 4) show that there is a strong correlation between the popularity gap and the difference in performance.

**End to End Results** We evaluate end to end performance in several evaluation settings with all results provided in Table 7. The metrics used are F1 for slot filling and question answering and accuracy for fact checking. In the “oracle” setting, we directly provide the downstream NLP model the gold document, and find that the gap between head entities and tail entities is fairly small. This suggests that in closed NLP settings, where the gold document is known, entity disambiguation is not a major concern.

In the regular retrieval setting, we provide the model the top 20 documents as ranked by a retrieval system (BLINK and DPR), and find that retrievers still perform better on head entity queries (see Appendix A). The downstream systems that use retrieved documents display a noticeable gap in end-to-end performance between head and tail entity inputs. This is expected, as retrieval systems perform worse on tail entities.

**User Study** AmbER sets are created in a largely automatic process, raising questions about data quality. To address these questions, we conduct a small user study on AmbER sets to evaluate whether the queries are resolvable by humans. We present a query from a QA AmbER set along with three documents for the entities from the same AmbER set, one of which is the gold document. We first ask the user to select the relevant document, then we ask the user to select an answer span from the selected document. In total, we asked 7 subjects to examine about 120 queries across AmbER-*H* and AmbER-*N*, and computed their accuracy in

<sup>6</sup>Bin width of 20%. Queries with a popularity gap higher than 100% are binned into the highest bin.

System	AmbER-H		AmbER-N	
	Doc Acc.	EM	Doc Acc.	EM
TF-IDF	43.3	-	50.3	-
DPR	69.1	-	68.3	-
BLINK	69.1	-	74.1	-
Bootleg	79.6	-	73.1	-
BERT	-	71.8	-	75.5
<b>Human</b>	<b>100</b>	<b>78.8</b>	<b>97.9</b>	<b>77.5</b>

Table 8: **User study on AmbER QA.** Humans are nearly perfect in identifying the correct document for each query (Doc Acc), while existing retrievers frequently fail. When the gold document is provided to downstream NLP models (BERT), they do almost as well as humans in answering the question (EM).

selecting the correct document and answer (Table 8). We also compare retrievers for this task, *i.e.* select from 3 documents for the same queries, and find that humans perform very well on the document selection task compared to retrievers on both sets. We also compare the accuracy of answer selection, and see that the closed domain NLP model (fine-tuned BERT) is as almost accurate as humans on the same set of queries<sup>7</sup>. This further confirms that closed NLP models are not the source of bias towards head entities, but the retrievers are.

## 6 Related Work

**Entity Ambiguity** As previously mentioned, entity ambiguity is when a single name can match multiple entities in a knowledge source. Entity ambiguity has been most studied in the context of entity linking (Rao et al., 2013). To improve disambiguation, entity linkers have included auxiliary information such as entity types (Onoe and Durrett, 2020) and entity descriptions (Logeswaran et al., 2019). A recent thread of work aims to study how language models recall and leverage information about names and entities. Prabhakaran et al. (2019) shows that names can have a measurable effect on the prediction of sentiment analysis systems. Shwartz et al. (2020) demonstrates that pre-trained language models implicitly resolve entity ambiguity by grounding names to entities based on the pre-training corpus. The problem of entity ambiguity also appears implicitly in entity-centric tasks such as determining the semantic relatedness between entities (Hoffart et al., 2012) and entity-oriented

<sup>7</sup>The relatively low answer score is due to artifacts in using EM for QA evaluation, and is consistent with human performance on span selection (Rajpurkar et al., 2016).

search (Balog et al., 2010, 2011). We draw inspiration from these works by studying entity ambiguity in the context of open-domain NLP.

**Popularity Bias** System’s that perform worse on the long-tail suffer from what is known as popularity bias. This problem has been studied extensively in the recommendation systems literature, where recommendation systems are known to often ignore the long-tail of products and instead recommend very popular items (Abdollahpouri et al., 2017; Chen et al., 2020). This has the effect of unfairly hurting users who would prefer these less-popular items (Abdollahpouri et al., 2019; Ciampaglia et al., 2018). We explore popularity bias from the angle of retrieval as opposed to recommendation, and find popularity bias exists in retrieval systems.

**Open-Domain Ambiguity** Ambiguity is an inherent problem when it comes to open-domain reasoning. Min et al. (2020) showed that half of instances sampled from Natural Questions are ambiguous, with multiple correct answers. AmbER sets are similar in that the ambiguity is in terms of the entity in the query, however, in contrast to Natural Questions, AmbER set inputs have been constructed such that the ambiguity is resolvable.

**Challenge Sets** There have been many evaluation sets specifically designed to assess a model’s ability to handle a specific phenomenon (Naik et al., 2018; Zhao et al., 2018; McCoy et al., 2019; Warstadt et al., 2020; Richardson et al., 2020; Jeretic et al., 2020; Ribeiro et al., 2019). Some of these challenge sets, similar to AmbER sets, use templates to generate a large amount of evaluation data quickly (Richardson et al., 2020; McCoy et al., 2019; Ribeiro et al., 2020). AmbER sets can be viewed as a challenge set for assessing open-domain systems’ ability to handle entity ambiguity.

## 7 Conclusion

Entity ambiguity is an inherent problem in retrieval, as many entities can share a name. For evaluating disambiguation capabilities of retrievers, we introduce AmbER sets; an AmbER set is a collection of task-specific queries about entities that share a name, but the queries have sufficient content to resolve the correct entity. We create a broad range of AmbER sets, covering many entity types, with input queries for three open-domain NLP tasks: fact checking, slot filling, and question answering. Our experiments demonstrate the struggles of current



retrievers in handling entity ambiguity. In particular, we find that the popularity of an entity in relation to other entities that share a name plays a significant role during disambiguation. For instance, we find that all tested retrievers are about twice as likely to retrieve erroneous documents when dealing with less popular entities than the most popular entity with the same name. Future goals include improving entity disambiguation capabilities of retrievers, perhaps more directly incorporating ideas from entity linking and coreference resolution. The AmbER sets and the code for the generation pipeline is available at <https://github.com/anthonywchen/AmbER-Sets>.

## Acknowledgements

We would like to thank Jo Daiber, Michael Tu, Russ Webb, Matt Gardner, Robert Logan, Sherry Tongshuang Wu, and the anonymous reviewers for providing valuable feedback for our work. This work is funded in part by the DARPA MCS program under Contract No. N660011924033 with the United States Office Of Naval Research.

## References

- Himan Abdollahpouri, Robin Burke, and Bamshad Mobasher. 2017. [Controlling popularity bias in learning-to-rank recommendation](#). In *Proceedings of the Eleventh ACM Conference on Recommender Systems, RecSys 2017, Como, Italy, August 27-31, 2017*, pages 42–46. ACM.
- Himan Abdollahpouri, Masoud Mansoury, Robin Burke, and Bamshad Mobasher. 2019. [The unfairness of popularity bias in recommendation](#). *arXiv preprint arXiv:1907.13286*.
- K. Balog, Pavel Serdyukov, and Arjen P. de Vries. 2010. [Overview of the trec 2010 entity track](#). In *TREC*.
- K. Balog, Pavel Serdyukov, and Arjen P. de Vries. 2011. [Overview of the trec 2011 entity track](#). In *TREC*.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. [Reading Wikipedia to answer open-domain questions](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879, Vancouver, Canada. Association for Computational Linguistics.
- J. Chen, Hande Dong, Xiao lei Wang, Fuli Feng, Ming-Chieh Wang, and X. He. 2020. [Bias and debias in recommender system: A survey and future directions](#). *arXiv preprint arXiv:2010.03240*.
- Giovanni Luca Ciampaglia, Azadeh Nematzadeh, Filippo Menczer, and Alessandro Flammini. 2018. [How algorithmic popularity bias hinders or promotes quality](#). *Scientific Reports*, 8.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. [DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.
- Hady Elsahar, Pavlos Vougiouklis, Arslan Remaci, Christophe Gravier, Jonathon Hare, Frederique Laforest, and Elena Simperl. 2018. [T-REx: A large scale alignment of natural language with knowledge base triples](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)*, Miyazaki, Japan. European Language Resources Association (ELRA).
- Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke Zettlemoyer. 2018. [AllenNLP: A deep semantic natural language processing platform](#). In *Proceedings of Workshop for NLP Open Source Software (NLP-OSS)*, pages 1–6, Melbourne, Australia. Association for Computational Linguistics.
- Johannes Hoffart, Stephan Seufert, Dat Ba Nguyen, Martin Theobald, and Gerhard Weikum. 2012. [KORE: keyphrase overlap relatedness for entity disambiguation](#). In *21st ACM International Conference on Information and Knowledge Management, CIKM'12, Maui, HI, USA, October 29 - November 02, 2012*, pages 545–554. ACM.
- Johannes Hoffart, Mohamed Amir Yosef, Ilaria Bordino, Hagen Fürstenauf, Manfred Pinkal, Marc Spaniol, Bilyana Taneva, Stefan Thater, and Gerhard Weikum. 2011. [Robust disambiguation of named entities in text](#). In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 782–792, Edinburgh, Scotland, UK. Association for Computational Linguistics.
- Paloma Jeretic, Alex Warstadt, Suvrat Bhooshan, and Adina Williams. 2020. [Are natural language inference models IMPPRESsive? Learning IMPlicature](#)

- and PRESupposition. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8690–8705, Online. Association for Computational Linguistics.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. [Dense passage retrieval for open-domain question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural questions: A benchmark for question answering research](#). *Transactions of the Association for Computational Linguistics*, 7:452–466.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. [Latent retrieval for weakly supervised open domain question answering](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6086–6096, Florence, Italy. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020a. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Patrick S. H. Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020b. [Retrieval-augmented generation for knowledge-intensive NLP tasks](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#). *arXiv preprint arXiv:1907.11692*.
- Lajanugen Logeswaran, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, Jacob Devlin, and Honglak Lee. 2019. [Zero-shot entity linking by reading entity descriptions](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3449–3460, Florence, Italy. Association for Computational Linguistics.
- Tom McCoy, Ellie Pavlick, and Tal Linzen. 2019. [Right for the wrong reasons: Diagnosing syntactic heuristics in natural language inference](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3428–3448, Florence, Italy. Association for Computational Linguistics.
- Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. [AmbigQA: Answering ambiguous open-domain questions](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5783–5797, Online. Association for Computational Linguistics.
- Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. [Stress test evaluation for natural language inference](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 2340–2353, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Yasumasa Onoe and Greg Durrett. 2020. [Fine-grained entity typing for domain independent entity linking](#). In *AAAI*.
- Laurel Orr, Megan Leszczynski, Simran Arora, Sen Wu, Neel Guha, Xiao Ling, and Christopher Ré. 2020. [Bootleg: Chasing the tail with self-supervised named entity disambiguation](#). *arXiv preprint arXiv:2010.10363*.
- Fabio Petroni, Aleksandra Piktus, Angela Fan, Patrick Lewis, Majid Yazdani, Nicola De Cao, James Thorne, Yacine Jernite, Vladimir Karpukhin, Jean Maillard, Vassilis Plachouras, Tim Rocktäschel, and Sebastian Riedel. 2021. [KILT: a benchmark for knowledge intensive language tasks](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2523–2544, Online. Association for Computational Linguistics.
- Vinodkumar Prabhakaran, Ben Hutchinson, and Margaret Mitchell. 2019. [Perturbation sensitivity analysis to detect unintended model biases](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5740–5745, Hong Kong, China. Association for Computational Linguistics.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ questions for machine comprehension of text](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Delip Rao, Paul McNamee, and Mark Dredze. 2013. [Entity linking: Finding extracted entities in a knowledge base](#). In *Multi-source, Multilingual Information Extraction and Summarization*.

- Marco Tulio Ribeiro, Carlos Guestrin, and Sameer Singh. 2019. [Are red roses red? evaluating consistency of question-answering models](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6174–6184, Florence, Italy. Association for Computational Linguistics.
- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. 2020. [Beyond accuracy: Behavioral testing of NLP models with CheckList](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4902–4912, Online. Association for Computational Linguistics.
- Kyle Richardson, H. Hu, L. Moss, and A. Sabharwal. 2020. Probing natural language inference models through semantic fragments. In *AAAI*.
- Ozge Sevgili, Artem Shelmanov, Mikhail V. Arkhipov, Alexander Panchenko, and Christian Biemann. 2020. [Neural entity linking: A survey of models based on deep learning](#). *arXiv preprint arXiv:2006.00575*.
- Vered Shwartz, Rachel Rudinger, and Oyvind Tafjord. 2020. [“you are grounded!”: Latent name artifacts in pre-trained language models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6850–6861, Online. Association for Computational Linguistics.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. [FEVER: a large-scale dataset for fact extraction and VERification](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.
- Denny Vrandečić and M. Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM*, 57:78–85.
- Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. [BLiMP: The benchmark of linguistic minimal pairs for English](#). *Transactions of the Association for Computational Linguistics*, 8:377–392.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45. Association for Computational Linguistics.
- Ledell Yu Wu, F. Petroni, Martin Josifoski, Sebastian Riedel, and Luke Zettlemoyer. 2020. Zero-shot entity linking with dense entity retrieval. In *EMNLP*.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. [Gender bias in coreference resolution: Evaluation and debiasing methods](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20, New Orleans, Louisiana. Association for Computational Linguistics.

## Appendix

### A Top-20 Retrieval Results

We provide results for top-20 retrieval in Table 9. Top-20 retrieval is used for providing documents in the end-to-end evaluation setting. In this setting, retrieval accuracy measures whether a gold document appears in one of the top-20 retrieved documents. Similar to top-1 retrieval, retrievers continue to perform better on head queries.

### B Task Specific Templates

Table 10 contains the templates used to instantiate the task-specific inputs. Templates were written on a per-property basis. We note that many of the properties share templates that are very similar.

### C Computational Resources

All experiments (*e.g.*, training baselines, generating AmbER sets, etc.) were conducted on a machine with 500 GB of RAM, 64 CPUs, and using an NVIDIA TitanRTX with 24 GB of RAM. Retrieval on a collection of AmbER sets takes about 12 hours for the most time-consuming retriever, BLINK. Training a downstream model takes roughly 5 hours and inference on a collection of AmbER sets takes less than 30 minutes.

### D Retriever Details

For BLINK, DPR, and TF-IDF, we use the retriever code in the KILT repository released by Facebook<sup>8</sup>. For Bootleg, we use the code provided by the Hazy Research group<sup>9</sup>.

### E Downstream Model Details

For question answering, we train a RoBERTa-Large model on Natural Questions. We use the negative documents in Natural Questions to train a “no-answer” classifier using the [CLS] token. During inference, we take the highest-scoring span where the answer is not classified as “no-answer”. For slot filling, we train a BART-base model. For each slot filling instance, we train with the top non-gold document retrieved by TF-IDF as a negative document. For this negative document, we train the model to generate a “none” token, and during inference, we take the highest scoring answer that is

not “none”. For fact checking, we train a three-way (*i.e.*, SUPPORTS, REFUTES, NEUTRAL) BERT-base classifier. Similar to slot filling, we train with the top non-gold document retrieved by TF-IDF as a negative document and train the model to classify this negative document as NEUTRAL. During inference, we take the highest scoring prediction that is not NEUTRAL. When training baselines models, we do not tune over hyperparameters and train with a batch size of 32 for 3 epochs.

<sup>8</sup><https://github.com/facebookresearch/KILT>

<sup>9</sup><https://github.com/HazyResearch/bootleg>

Collection	Retriever	Fact Checking				Slot Filling				Question Answering			
		All	Head	Tail	∇	All	Head	Tail	∇	All	Head	Tail	∇
<b>AmbER-H</b>	TF-IDF	65.8	78.5	55.4	26.7	72.0	83.5	62.5	55.6	72.6	82.0	64.8	55.9
	DPR	39.8	51.0	30.6	4.1	26.6	37.0	18.1	6.8	36.1	49.3	25.3	9.6
	BLINK	78.6	82.0	76.0	43.8	73.3	73.9	72.8	64.6	58.8	60.3	57.5	32.2
	Bootleg	<b>96.5</b>	<b>97.6</b>	<b>95.6</b>	<b>93.2</b>	<b>96.6</b>	<b>97.7</b>	<b>95.7</b>	<b>93.6</b>	<b>96.5</b>	<b>97.6</b>	<b>95.6</b>	<b>93.5</b>
<b>AmbER-N</b>	TF-IDF	50.8	57.0	44.1	12.0	46.8	53.4	39.7	35.3	52.0	59.1	44.4	40.7
	DPR	62.3	75.8	47.7	27.8	57.3	71.4	42.0	29.4	63.4	77.9	47.8	37.2
	BLINK	33.5	38.7	27.9	1.3	18.2	21.5	14.6	5.8	74.7	80.6	68.3	53.0
	Bootleg	<b>79.3</b>	<b>80.2</b>	<b>78.4</b>	<b>61.5</b>	<b>89.6</b>	<b>91.9</b>	<b>87.1</b>	<b>85.3</b>	<b>83.8</b>	<b>83.6</b>	<b>84.1</b>	<b>71.1</b>

Table 9: Top-20 retrieval results measuring retrieval accuracy and ∇.

	Property	Question Answering Template	Fact Checking Template
AmbER-H	instrument	Which musical instrument did <i>\$name</i> play? What musical instrument does <i>\$name</i> play? What instrument does <i>\$name</i> play?	<i>\$name</i> plays the <i>\$object</i> . <i>\$name</i> plays the musical instrument <i>\$object</i> . The <i>\$object</i> is played by <i>\$name</i> .
	movement	What movement did <i>\$name</i> participate in? Which movement is <i>\$name</i> associated with? What movement is <i>\$name</i> associated with?	<i>\$name</i> was a member of the <i>\$object</i> movement. <i>\$name</i> participated in the <i>\$object</i> movement. <i>\$name</i> was a part of the <i>\$object</i> movement.
	appears in	What works does the fictional entity <i>\$name</i> appear in? What work is the character <i>\$name</i> present in? Which work was the character <i>\$name</i> in?	<i>\$name</i> is a character in <i>\$object</i> . <i>\$name</i> is a fictional character in <i>\$object</i> . <i>\$object</i> features the fictional character <i>\$name</i> .
	doctoral student	Who were the doctoral students of <i>\$name</i> ? Who are <i>\$name</i> 's doctoral students? Who did <i>\$name</i> advise?	<i>\$name</i> has a doctoral student named <i>\$object</i> . <i>\$name</i> 's doctoral student is <i>\$object</i> . <i>\$name</i> advised their student <i>\$object</i> .
	military branch	What branch of the military does <i>\$name</i> belong to? Which military branch does <i>\$name</i> belong to? What military branch is <i>\$name</i> affiliated with?	<i>\$name</i> is a member of the <i>\$object</i> . <i>\$name</i> belongs to the military branch <i>\$object</i> . <i>\$name</i> belongs to the <i>\$object</i> branch of the military.
	sports position	What is the position that <i>\$name</i> plays? What position does <i>\$name</i> play? Which position does <i>\$name</i> play?	<i>\$name</i> plays the <i>\$object</i> position. <i>\$name</i> plays as a <i>\$object</i> .
	sports team	<i>\$name</i> plays for which team? What team does <i>\$name</i> play for? Which team does <i>\$name</i> play for?	<i>\$name</i> is a player on the <i>\$object</i> . <i>\$name</i> plays for the <i>\$object</i> team. <i>\$name</i> plays for the <i>\$object</i> .
	battles or wars	What were the wars that <i>\$name</i> participated in? Which battle did <i>\$name</i> fight in? Which war did <i>\$name</i> fight?	<i>\$name</i> fought in the <i>\$object</i> . <i>\$name</i> fought in <i>\$object</i> .
	sport	Which sport does <i>\$name</i> participate in? Which sport does <i>\$name</i> play? What sport does <i>\$name</i> play?	<i>\$name</i> plays <i>\$object</i> . <i>\$name</i> plays the sport <i>\$object</i> .
	AmbER-N	performer	Who performs <i>\$name</i> ? Who is the performer of <i>\$name</i> ? Who performed <i>\$name</i> ?
record label		What is the record label of <i>\$name</i> ? What is the record label for <i>\$name</i> ? <i>\$name</i> belongs to which record label?	<i>\$object</i> is the record label for <i>\$name</i> . <i>\$name</i> 's record label is <i>\$object</i> .
tracklist		What song appears in the album <i>\$name</i> ? What song appears on <i>\$name</i> ? What are the tracks in <i>\$name</i> ?	<i>\$name</i> belongs to <i>\$object</i> tracklist. <i>\$object</i> is on the release of <i>\$name</i> . <i>\$object</i> is a song in the <i>\$name</i> tracklist.
industry		Which industry is <i>\$name</i> in? In what industry is <i>\$name</i> ? What is <i>\$name</i> 's industry?	<i>\$name</i> is in the industry of <i>\$object</i> . The company <i>\$name</i> is in the <i>\$object</i> industry. <i>\$name</i> 's industry is <i>\$object</i> .
population		What is the total population of <i>\$name</i> ? What is the population of <i>\$name</i> ? How many people live in <i>\$name</i> ?	The population of <i>\$name</i> is <i>\$object</i> . <i>\$name</i> 's population is <i>\$object</i> . <i>\$name</i> has a population of <i>\$object</i> .
cast member		Who acted in <i>\$name</i> ? Who is a cast member on <i>\$name</i> ? Who starred in <i>\$name</i> ?	<i>\$object</i> was a cast member in <i>\$name</i> . <i>\$object</i> appeared in <i>\$name</i> . <i>\$object</i> acted in <i>\$name</i> .
screenwriter		Who was the screenwriter for <i>\$name</i> ? Who was screenwriter for <i>\$name</i> ? Who is <i>\$name</i> 's screenwriter?	<i>\$name</i> 's screenwriter is <i>\$object</i> . <i>\$object</i> wrote the screenplay of <i>\$name</i> . <i>\$object</i> screenwrote <i>\$name</i> .
# seasons		How many seasons are there in <i>\$name</i> ? How many seasons does <i>\$name</i> have? How many seasons were there in <i>\$name</i> ?	There were <i>\$object</i> seasons in <i>\$name</i> . <i>\$name</i> has <i>\$object</i> seasons.
author		Who is the author of <i>\$name</i> ? Who wrote <i>\$name</i> ? Who authored <i>\$name</i> ?	<i>\$name</i> wrote <i>\$object</i> . <i>\$name</i> is written by <i>\$object</i> . <i>\$object</i> authored <i>\$name</i> .

Table 10: Templates used to instantiate the task-specific inputs.