

---

# DP-KB: Data Programming with Knowledge Bases Improves Transformer Fine Tuning for Answer Sentence Selection

---

**Nic Jedema**

Alexa AI - Graphiq  
Santa Barbara, CA 93101  
jedem@amazon.com

**Thuy Vu**

Alexa AI - Search  
Manhattan Beach, CA 90266  
thuyvu@amazon.com

**Manish Gupta**

Alexa AI - Graphiq  
Santa Barbara, CA 93101  
manishg@amazon.com

**Alessandro Moschitti**

Alexa AI - Search  
Manhattan Beach, CA 90266  
amosch@amazon.com

## Abstract

While transformers demonstrate impressive performance on many knowledge intensive (KI) tasks, their ability to serve as implicit knowledge bases (KBs) remains limited, as shown on several slot-filling, question-answering (QA), fact verification, and entity-linking tasks. In this paper, we implement an efficient, data-programming technique that enriches training data with KB-derived context and improves transformer utilization of encoded knowledge when fine-tuning for a particular QA task, namely answer sentence selection (AS2). Our method outperforms state of the art transformer approach on WikiQA and TrecQA, two widely studied AS2 benchmarks, increasing by 2.0% p@1, 1.3% MAP, 1.1% MRR, and 4.4% p@1, 0.9% MAP, 2.4% MRR, respectively. To demonstrate our improvements in an industry setting, we additionally evaluate our approach on a proprietary dataset of Alexa QA pairs, and show increase of 2.3% F1 and 2.0% MAP. We additionally find that these improvements remain even when KB context is omitted at inference time, allowing for the use of our models within existing transformer workflows without additional latency or deployment costs.

## 1 Introduction

Transformers are powerful sequence-to-sequence models that leverage the mechanism of self-attention [28] to capture long-range dependencies, such as relationships between words in a natural language text. Compared to sequential models such as recurrent or convolutional networks, they make efficient use of available processing power. Transformer based models such as BERT, XLNet and BART [4, 33, 18, 17] hold state of the art on many natural language processing (NLP) tasks, including next sentence prediction, natural language generation, and natural language inference [23, 17, 10]. In addition to efficiently encoding linguistic information from unlabelled text, their top performance on knowledge-intensive (KI) NLP tasks [21], such as question-answering [24] have led to the hypothesis that transformers also encode relational knowledge, and as such serve as parameterized, implicit knowledge bases (KBs) [22].

However, it has also been shown that transformer knowledge acquisition [22, 25] and subsequent utilization [27, 12] can be uncontrollable, highly context dependent, and tightly coupled to language acquisition. These limitations may impact performance on downstream tasks, including KI tasks

Example 1:	<p><b>Q:</b> How old is Elton John's husband?</p> <p><b>Correct:</b> David Furnish is 57 years old. He was born on October 25, 1962.</p> <p><b>Selected:</b> Elton John and David Furnish became an item after meeting in the early 1990s and in 2005.</p>
Example 2:	<p><b>Q:</b> How many humps on a Camel?</p> <p><b>Correct:</b> The two surviving species of camel are the dromedary, or one-humped camel, which is native to the Middle East and the Horn of Africa; and the Bactrian, or two-humped camel, which inhabits Central Asia.</p> <p><b>Selected:</b> A camel is an even-toed ungulate within the genus <i>Camelus</i>, bearing distinctive fatty deposits known as "humps" on its back.</p>
Example 3:	<p><b>Q:</b> What some legal uses of meth?</p> <p><b>Correct:</b> Although rarely prescribed, methamphetamine hydrochloride is approved by the U.S. Food and Drug Administration (FDA) for the treatment of attention deficit hyperactivity disorder and obesity under the trade name Desoxyn.</p> <p><b>Selected:</b> Methamphetamine, also known as metamfetamine, meth, ice, crystal, glass, tik, N-methylamphetamine, methylamphetamine, and desoxyephedrine, is a psychostimulant of the phenethylamine and amphetamine class of psychoactive drugs.</p>

Table 1: Three QA examples incorrectly predicted by a state-of-the-art transformer answer selection model (TANDA [7]).

like answer sentence selection (AS2) [31]. Table 1 is illustrative of limitations of some of these deficiencies of transformers in precisely leveraging encoded information. Transformer models that show top performance [7, 15] on widely studied benchmarks [32, 31] still fail to classify many QA pairs correctly. In *Example 1*, the model is unable to leverage knowledge of the identity between *Elton John's husband* and *David Furnish*. In *Example 2*, *one-humped* or *two-humped* are not recognizable as quantities pertaining to the uncommonly quantity *humps*. *Example 3* shows the difficulty in reasoning for the a rare prescriptive use of the illicit drug *methamphetamine*. These examples also illustrate relevance of this task as a means to assess impact of deficiencies in transformer knowledge utilization.

In this paper study, we propose an efficient, data-programming approach utilizing a KB that improves performance on answer selection tasks and demonstrate that some of these limitations can be mitigated during fine-tuning with simple data augmentation technique.

A number of recent studies have also studied approaches that aim to improve transformer performance on KI tasks, proposing the use of differentiable knowledge retrievers [8, 13, 16], retrieval-augmented generation (RAG) [13], KB embeddings such as KnowBERT [20] and ERNIE [34], and pre-training on verbalized KBs such as KELM [1]. While these approaches offer promising benefits for transformer knowledge encoding and retrieval, to our knowledge, none of them have been shown to outperform existing state of the art for answer selection, a task that is essential to several question answering services provided by commercial voice assistants. Additionally, each of these approaches is significantly complex and require significant work to leverage in production applications. Our approach, on the other hand, leverages Elasticsearch to tag KB entries in input QA pairs, derives weak-supervision signals from tagged KB entries, and incorporates this context only during fine-tuning. We show that our simple, efficient and data-programming method confers significant performance benefits over the state of the art for answer sentence selection, even when KB context is omitted at inference time.

The main contributions of our work are:

- We show that several limitations in the use of transformers implicit KBs can be overcome using a simple data-programming approach by outperforming state-of-the-art models on several QA tasks:
  1. increasing by 2.0% p@1, 1.3% MAP, 1.1% MRR and 4.4% p@1, 0.9% MAP, 2.4% MRR on WikiQA and TrecQA respectively, two widely used AS2 benchmarks.
  2. increasing by 2.3% F1 and 2.0% MAP on AlexaQA pairs, a proprietary commercial answer classification benchmark.
- We show that KB is not needed at inference time, allowing our trained models to be used as drop-in replacements for existing transformer-based AS2 systems.

## 2 Background

### 2.1 Transformers

The transformer [28] is an architecture for efficiently transforming one sequence into another via *self-attention*, a mechanism that differentially weighs the significance of discrete tokens in an input sequence. Compared to sequentially aligned or convolutional networks such as RNNs and CNNs, transformer models have proven to be extremely effective at efficiently capturing long-range dependencies between words in natural language [4, 18, 17], including some structured knowledge such as the *husband of* relation between named entities [22]. Pre-training transformer models on large, unstructured corpora of unlabelled text [3, 11] allows them to capture linguistic and factual knowledge prior to subsequent fine-tuning on downstream tasks. The suitability of transformer based models such as BERT [4] for this type of transfer learning dramatically increases their reusability, driving state of the art results on many tasks [7, 15, 8, 33, 9] in addition to widespread adoption in industry.

### 2.2 Transformer Limitations as Knowledge Bases

While transformers have demonstrated strong performance on question-answering [25] and fill-in-the-blank cloze tasks [22, 12] without access to external information, the modulation of transformer knowledge acquisition and utilization is limited. Cloze task [22] and question answering [25] probes demonstrate transformer knowledge acquisition is largely uncontrollable and often only results in the acquisition of frequently observed information. Further, transformer recall of factual knowledge on cloze tasks remains tightly bound to learned linguistic representation [12]. In a systematic study [27] on multiple tasks, it was shown that transformers lack robust multi-hop reasoning faculties, are insensitive to adverbial modifiers like "always", "some", and "never", and are unable to robustly compare quantities. For example, while RoBERTa [18] appears able to effectively compare numbers, it is unable to compare when values are given in *ages*. Other studies additionally have shown insensitivity to negation [5], difficulty with misspellings and short, simple sequences [26], and sensitivity to sequence length, punctuation, and subject-verb agreement [2].

### 2.3 Answer Sentence Selection (AS2) and Answer Classification

Answer sentence selection (AS2) consists of ranking answer candidates given a question and one or more answer candidates, while binary answer classification consists of classifying answer candidates as *correct* or *incorrect* given the same input. Both tasks encourage models that leverage encoded knowledge to select the most correct answer and thus may be used to probe model knowledge and reasoning capabilities. Let  $q$  be a question,  $C_q = \{c_1, \dots, c_n\}$  be a set of answer sentence candidates for  $q$ , we define  $\mathcal{R}$  as a ranking function, which orders the candidates in  $C_q$  according to a score,  $p(q, c_i)$ , indicating the probability of  $c_i$  to be a correct answer for  $q$ . Answer sentence selection is performed by taking the highest scoring candidate in  $C_q$ , while binary answer classification is performed by assigning the label of the highest probability class as determined by the ranking function  $\mathcal{R}$ .

Widely used metrics for AS2 performance are mean average precision (MAP) and mean reciprocal rank (MRR), while mean average precision (MAP) and F-score (F1) are commonly used for binary answer classification. To our knowledge, transformer models [7, 15] demonstrate a strong state of the art on the AS2 task.

## 3 Modeling

### 3.1 Datasets

We study popular Answer Sentence Selection datasets to evaluate the benefits of our approach for this task: Answer Sentence Natural Questions (ASNQ) [7], WikiQA [32], and TrecQA [30]. ASNQ is a large scale QA dataset derived from Google’s Natural Questions [14] dataset, with more than  $\sim 84K$  unique questions. The train split of this dataset is used to transfer a pre-trained transformer model to the AS2 task. WikiQA [32] and TrecQA [30] are widely studied benchmark datasets for Answer Sentence Selection with over  $\sim 1.1K$  and  $\sim 1.2K$  unique questions respectively. We utilize the *clean*

versions of both WikiQA and TrecQA, as well as the *TRAIN-ALL* split of TrecQA for fine-tuning. All of these datasets are available under dataset specific licenses that permit their use and distribution for academic purposes.

We additionally evaluate the benefits of our approach in an industry setting using the binary answer classification task using AlexaQA. AlexaQA is a proprietary benchmark dataset that contains  $\sim 107K$  unique questions obtained from de-identified samples of Amazon Alexa QA traffic with correct/incorrect labels assigned by expert annotators. Detailed statistics of each dataset by split are shown in Table 6 in Appendix A.

### 3.2 Dataset Preprocessing

We implement a novel data enrichment pipeline that use an ElasticSearch index of 20.7M item and relation labels obtained by popularity-based filtering of Alexa’s KB. Our pipeline tags KB entries in input text by aggregating the results of three queries on our index. For each word  $w$  in the set of words  $W = \{w_1, \dots, w_n\}$  in the input text, we tag  $w_i$  as a KB entry if:

- $w_i$  is an *exact* match for a label in the index
- $w_i$  is *contained* by a label in the index
- $w_i$  and  $w_i + 1$  is a *quorum* match for a label in the index

Consecutive labels matching the same entry are assumed joined together, matches are sorted for relevance, and the top result is selected as the KB entry. KB meta-data for entries is derived from selected KB properties, such as the *collection* property that indicates classifications such as *celebrity*, *book*, or *album*.

### 3.3 Incorporating KB-derived Context for Transformer Training

Metadata for each entry tagged by the preprocessing pipeline (Section 3.2) is resolved to a textual representation using corresponding KB labels. An example of the JSON produced from this resolution is shown below:

---

```
{
  "text": "David Furnish is 57 years old.",
  "kb_tags": [{
    "kb_id": "e-478772",
    "popularity": 0.981,
    "candidate_title": "David Furnish",
    "candidate_aliases": "David James Furnish, Elton John's husband",
    "collection": "celebrity",
    "relations": "married_to, years_old, birth_date, ... ",
  ]}]
```

---

Inspired by other studies [19, 1] that verbalize structured data for use in language models, we insert the textual representation of KB context directly into model input. This approach may distract the model from attending to the QA pair itself if too much context is added and we thus employ two strategies to prevent this. First, we limit metadata to the *collection* property, whose values include common categories such as "celebrity", "quantity", and "generic drug form".<sup>1</sup> The *collection* property in our KB has many analogous properties in other KBs, for example, the *instance of* relation in Wikidata [6].

Second, we employ a filter that constrains the number of entries from which KB context is added. The *intersection* filter exploits the intuitive hypothesis that correct QA pairs will contain the same KB entries, adding context only if the same entry is tagged in both the question and the answer. For example, this filter adds context for entry *David Furnish* from the QA pair: *Q: how old is Elton John's husband; A: David Furnish is 57 years old* because the question contains *Elton John's husband*, an alias for "David Furnish" in our KB, and the answer contains *David Furnish*. The *intersection*

---

<sup>1</sup>initial experimentation using metadata derived from the *popularity*, *aliases*, and *relations* suggested that the *collection* property was the most effective.

filter excluded context for entries like *57* and *husband*, even though entries for both exist in our KB. We additionally study the *1-best* filter, which selects the KB entry from the answer with the highest *popularity* in our KB as a more lenient alternative. Two strategies of concatenating context to question/answer text are also explored: *append* and *prepend*; in both cases, the model’s special separator token <sup>2</sup> is used to separate the context from question/answer text.

An example of the resulting sequences are shown below:

- **Append:** how old is elton john’s husband <\s> john furnish is 57 years old. he was born on october 25, 1962 <\s> *celebrity* <\s> *celebrity*
- **Prepend:** <\s> *celebrity* <\s> *celebrity* <\s> how old is elton john’s husband <\s> john furnish is 57 years old. he was born on october 25, 1962

### 3.4 Model Architecture

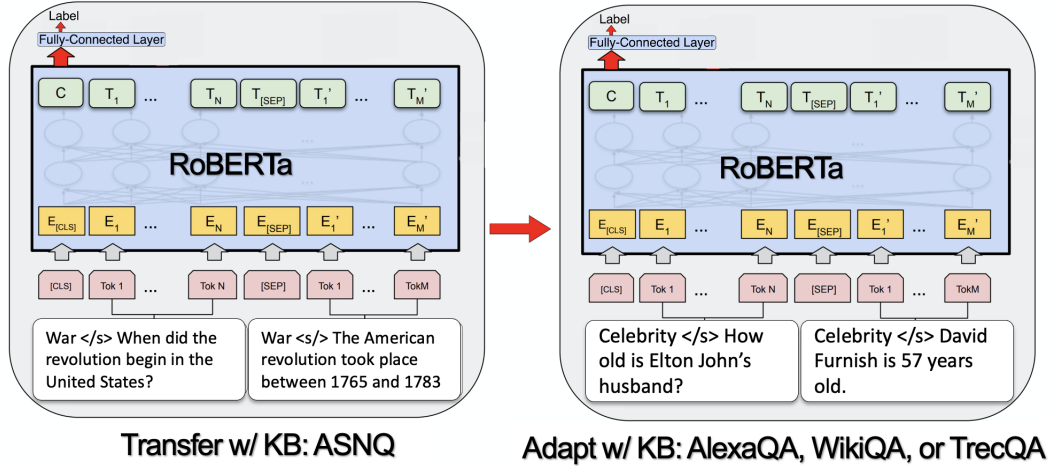


Figure 1: The Transfer-and-Adapt architecture using our approach

Our approach builds upon the Transfer-and-Adapt (TANDA) [7] architecture, the state of the art approach for answer sentence selection, by leveraging KB-derived context to address deficiencies observed in transformer knowledge utilization for this task. As illustrated in figure Figure 1, we transfer a pre-trained RoBERTa-base model [18] to the answer sentence selection task by fine-tuning on ASNQ and adapt the transferred model via further fine-tuning on our target dataset, either WikiQA, TrecQA, or AlexaQA. Training incorporates KB-derived context in both transfer and adapt steps, as discussed below. During inference, we optionally remove KB context so as to evaluate our approach as a drop-in replacement for existing transformer-based AS2 systems.

In order to isolate the benefits of our approach, we reuse the same optimizer, hyper-parameters, and early stopping strategy described in [7] and only alter the sequence length, increasing from 128 to 256 to accommodate additional context. Experiments on ASNQ, WikiQA, and TrecQA use AWS EC2 p3dn.24xlarge hosts, and those on AlexaQA use AWS SageMaker ml.p3.16xlarge notebook instances.

## 4 Results

Performance of KB augmented transformer models for standard fine-tuning (FT) on ASNQ is shown in Table 2. Transfer-and-Adapt performance with KB augmentation is reported for WikiQA, TrecQA, and AlexaQA in Tables 3, 4 and 5 respectively. We indicate the datasets used in Transfer-and-Adapt setting using two arguments, *transfer dataset* → *adapt dataset* with numerics in parentheses indicate

<sup>2</sup>We tried other separator tokens, including "#", ":", and " ", and found the special separator performs marginally better

training epochs. Baseline models - i.e. the RoBERTa base TANDA state-of-the-art set by [7] - are indicated by \* and lack the -KB suffix.

We additionally evaluate a setting in which KB context is omitted at inference time to explore the ability of our approach to modulate transformer knowledge utilization. Results for this setting are reported for each dataset and are indicated by the value of the *Incl. KB at Inference* column.

The results in the tables below demonstrate that:

- KB context improves fine-tuning performance on ASNQ, increasing the p@1, MRR and MAP by 2.9%, 3.0%, and 2.9% after 9 epochs.
- Training with KB context improves on the strong performance set by the state of the art TANDA approach on widely studied benchmarks, increasing the p@1, MRR and MAP by 2%, 1.3%, and 1.1% and 4.4%, 0.9%, and 2.4% on WikiQA and TrecQA respectively.
- The benefits of KB context generalize to our industry setting, increasing the F1 and MAP by 2.3% and 2.0% over the TANDA state of the art, RoBERTa ASNQ(9)→AlexaQA(1), and by .4% and 2.3% over the more challenging baseline, RoBERTa ASNQ(1) → AlexaQA(1).
- Models trained with our approach continue to outperform the TANDA state of the art even when KB context is omitted at inference time; in other words, the benefits of KB context are primarily realized during model training.

Model	KB Approach	Incl. KB at Inference	p@1	MAP	MRR
RoBERTa FT ASNQ(9)*	–	No	.599	.672	.716
RoBERTa FT ASNQ-KB(9)	Append, Intersection	Yes	.627	.696	.737
RoBERTa FT ASNQ-KB(9)	Prepend, Intersection	Yes	.627	<b>.702</b>	<b>.745</b>
RoBERTa FT ASNQ-KB(9)	Prepend, 1 best	Yes	.616	.694	.736
RoBERTa FT ASNQ-KB(9)	Append, Intersection	No	<b>.628</b>	.692	.736
RoBERTa FT ASNQ-KB(9)	Prepend, Intersection	No	.621	<b>.696</b>	<b>.739</b>
RoBERTa FT ASNQ-KB(9)	Prepend, 1 best	No	.617	.693	.735

Table 2: Performance of KB-augmented fine-tuned (FT) transformer models on ASNQ

Model	KB Approach	Incl. KB at Inference	p@1	MAP	MRR
RoBERTa ASNQ(9) → WikiQA(9)*	–	No	.827	.890	.901
RoBERTa ASNQ-KB(9) → WikiQA-KB(9)	Append, Intersection	Yes	.835	.891	.903
RoBERTa ASNQ-KB(9) → WikiQA-KB(9)	Prepend, Intersection	Yes	<b>.847</b>	<b>.903</b>	<b>.913</b>
RoBERTa ASNQ-KB(9) → WikiQA-KB(9)	Prepend, 1-best	Yes	.835	.885	.898
RoBERTa ASNQ-KB(9) → WikiQA-KB(9)	Append, Intersection	No	.835	.892	.902
RoBERTa ASNQ-KB(9) → WikiQA-KB(9)	Prepend, Intersection	No	<b>.843</b>	<b>.895</b>	<b>.907</b>
RoBERTa ASNQ-KB(9) → WikiQA-KB(9)	Prepend, 1-best	No	.839	.887	.900

Table 3: Performance of KB-augmented fine-tuned (FT) transformer models on WikiQA

Model	KB Approach	Incl. KB at Inference	p@1	MAP	MRR
RoBERTa ASNQ(9) → TrecQA(9)*	–	No	.897	.906	.942
RoBERTa ASNQ-KB(9) → TrecQA-KB(9)	Append, Intersection	Yes	.911	.901	.952
RoBERTa ASNQ-KB(9) → TrecQA-KB(9)	Prepend, Intersection	Yes	<b>.926</b>	<b>.914</b>	<b>.960</b>
RoBERTa ASNQ-KB(9) → TrecQA-KB(9)	Prepend, 1-best	Yes	.897	.900	.944
RoBERTa ASNQ-KB(9) → TrecQA-KB(9)	Append, Intersection	No	<b>.941</b>	<b>.915</b>	<b>.966</b>
RoBERTa ASNQ-KB(9) → TrecQA-KB(9)	Prepend, Intersection	No	.911	.901	.955
RoBERTa ASNQ-KB(9) → TrecQA-KB(9)	Prepend, 1-best	No	.926	.905	.959

Table 4: Performance of KB-augmented fine-tuned (FT) transformer models on TrecQA

Model	KB Approach	Incl. KB at Inference	F1	MAP
RoBERTa ASNQ(1) → AlexaQA(1)	–	No	.848	.839
RoBERTa ASNQ(9) → AlexaQA(1)*	–	No	.829	.842
RoBERTa ASNQ-KB(1) → AlexaQA-KB(1)	Append, Intersection	Yes	<b>.852</b>	.860
RoBERTa ASNQ-KB(1) → AlexaQA-KB(1)	Prepend, Intersection	Yes	.850	<b>.862</b>
RoBERTa ASNQ-KB(1) → AlexaQA-KB(1)	Prepend, 1-best	Yes	.850	.858
RoBERTa ASNQ-KB(1) → AlexaQA-KB(1)	Append, Intersection	No	<b>.851</b>	.859
RoBERTa ASNQ-KB(1) → AlexaQA-KB(1)	Prepend, Intersection	No	.850	<b>.861</b>
RoBERTa ASNQ-KB(1) → AlexaQA-KB(1)	Prepend, 1-best	No	.849	.857

Table 5: Performance of KB-augmented fine-tuned (FT) transformer models on AlexaQA. Models transferred for only (1) epoch are shown, since our experiments indicate that further epochs of transfer to ASNQ conveyed marginal benefits for AlexaQA.

## 5 Discussion

### 5.1 Comparing Context Generation Strategies

Results reported in Tables 2, 3, 4 and 5 all demonstrate that our approach outperforms the state of the art approach, even in the more challenging setting where KB context is omitted at inference time. We explain the robustness of our models to the omission of KB context in light of the proportion of each dataset that our approach impacts. The *intersection* filter adds KB context to only 3.38% of the ASNQ dataset, 5.27% of TrecQA, and 8.33% of WikiQA while the *1 best* filter adds context for 31.17% for ASNQ, 51.1% for TrecQA, 40.79% for WikiQA. We hypothesize that the large number of training examples seen without context allows the model to leverage context as a for weak supervision that encourages knowledge utilization and elaborate further in subsection 5.2 below.

These results show that the more intuitive *intersection* filter performs better than the *1 best* filter for both concatenation strategies, despite impacting between significantly less of each dataset. We conclude that the explicit conceptual alignment provided by the *intersection* conveys additional benefits beyond the addition of conceptual keywords provided by the *1 best* filter. The *prepend* strategy outperforms the *append* strategy on all datasets other than TrecQA, a deviation that we attribute to the small size of the TrecQA test set. We explicate these findings in light of the positional invariance the *prepend* strategy - that is, prepend always adds context in the same position in the sequence, whereas *append* does not. As a result, *prepend* models appear better able to attend to context and outperform their *append* counterparts, even though *prepend* models suffer more when context is omitted at inference.

### 5.2 Impact of KB Context

We leverage the three illustrative examples presented in Table 1 to probe the impact of our KB context and its potential to address the previously studied [12, 27] deficiencies of transformers as implicit KBs. Models trained with our approach classify each of these examples correctly, even when KB is omitted at inference, indicating that they may be able to exploit our context to refine their utilization of encoded knowledge. In order to identify the mechanism behind these benefits, we compare the attention of TANDA with that of our best model, *prepend, intersection*, using box plots of attention intensity and bar plots of activate head counts per layer in Appendix B.

**Example 1** requires the model to leverage encoded knowledge in order to make the connection between "husband" and "David Furnish" necessary to recognize that the phrase "is 57 years old" answers the question phrase "how old". Figure 2 presents model attention weights between tokens "how" and "57" and between "husband" and "David", where it can be seen that our approach significantly improves both the quantity of heads attending to these keywords and the intensity of this attention. It is likely that model pre-training has encoded this knowledge, given that the second sentence on David Furnish’s Wikipedia page reads: "He is married to English musician Sir Elton John". Unsurprisingly, changing the question or the answer text to remove this relation - to either "How old is David Furnish" or "Elton John’s husband David Furnish is 57 years old" - produces the correct answer from the TANDA model.

**Example 2** probes transformer ability to robustly recognize that *"one-humped"* and *"two-humped"* are values for the quantity sought by *"how many"* and are related to the subject *"Camel"*. We hypothesize that the KB context *"animal"* added for similar entities during training increases attention on *"camel"* tokens and their modifiers, *"one-humped"* and *"two-humped"* in this case. Figure 3 compares model attention weights of tokens *"many"* and *"Camel"* with the values *"one"* and *"two"* and again demonstrate that our approach significantly increases the intensity of model attention between these terms. Changing the answer to use common numeric values *"the Dromedary Camel has 1 hump...and the Bactrian Camel has 2 humps"* is sufficient for the TANDA model to select the correct answer.

**Example 3** illustrates whether the model is able to connect the adverbial phrase *"some legal uses"* in the question with the phrase *approved...for the treatment of...* in the correct answer. Interestingly, the KB context added for *"meth"* and entities like it is *"generic drug"*, which we hypothesize may encourage attention to relevant terms like *"treatment"* that are not commonly used in context of the subject *"meth"*. Figure 4 shows the weights connecting *"treatment"* with *"uses"* and *"meth"* and further demonstrates the impact of our approach on model attention. We conclude that in some cases, the context itself may provide relevant information that helps the model more effectively utilize uncommon knowledge, like that meth may be used as a medical treatment.

## 6 Conclusion

In this paper, we presented a data-programming approach that enriches transformer training data with KB-derived context, and demonstrate that it beats state of the art approach on several challenging knowledge-intensive question-answering benchmarks such as ASNQ, WikiQA, TrecQA, and Alexa QA. Our findings indicate that our approach addresses some deficiencies of transformer knowledge utilization that negatively impact AS2 performance. We probed the mechanism of our approach with challenging examples that highlight the potential ways in which our KB context may allow transformers to better utilize encoded knowledge. Our method is simple, efficient and task-agnostic, and training benefits remain even when KB context is omitted at inference time. We believe that our approach provides a way to rapidly integrate the benefits of KBs within the deployed inference pipelines utilized in many virtual-assistant workflows.

While we improve on the state of the art approach in AS2, we do acknowledge that our approach may face limitations of its own. While our approach is efficient in that it not require significant pre-training, unlike KB based approaches like KELM, KnowBERT, and ERNIE as well as retrieval oriented approaches like REALM and RAG, it is inefficient in that it likely does not leverage the full richness of our KB. This has the negative consequence that our approach still requires significant task-specific training and thus consumes significant GPU hours and the natural resources used to power them. Further work beyond the data-programming approach that we propose in the direction of more effective transformer architectures that enhance knowledge utilization can lessen this impact and provide models capable of more completely disentangling knowledge and language acquisition.

## References

- [1] AGARWAL, O., GE, H., SHAKERI, S., AND AL-RFOU, R. Knowledge graph based synthetic corpus generation for knowledge-enhanced language model pre-training, 2021.
- [2] CHERNYAVSKIY, A., ILVOVSKY, D., AND NAKOV, P. Transformers: "the end of history" for nlp?, 2021.
- [3] DAI, A. M., AND LE, Q. V. Semi-supervised sequence learning, 2015.
- [4] DEVLIN, J., CHANG, M.-W., LEE, K., AND TOUTANOVA, K. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [5] ETtinger, A. What BERT is not: Lessons from a new suite of psycholinguistic diagnostics for language models. *Transactions of the Association for Computational Linguistics* 8 (2020), 34–48.
- [6] FARDA-SARBAS, M., AND MÜLLER-BIRN, C. Wikidata from a research perspective – a systematic mapping study of wikidata, 2019.



- [7] GARG, S., VU, T., AND MOSCHITTI, A. Tanda: Transfer and adapt pre-trained transformer models for answer sentence selection, 2019.
- [8] GUU, K., LEE, K., TUNG, Z., PASUPAT, P., AND CHANG, M.-W. Realm: Retrieval-augmented language model pre-training, 2020.
- [9] HARDALOV, M., KOYCHEV, I., AND NAKOV, P. Enriched pre-trained transformers for joint slot filling and intent detection, 2020.
- [10] HE, P., LIU, X., GAO, J., AND CHEN, W. DeBERTa: Decoding-enhanced bert with disentangled attention, 2021.
- [11] HOWARD, J., AND RUDER, S. Universal language model fine-tuning for text classification, 2018.
- [12] JIANG, Z., XU, F. F., ARAKI, J., AND NEUBIG, G. How can we know what language models know?, 2020.
- [13] KARPUKHIN, V., OĞUZ, B., MIN, S., LEWIS, P., WU, L., EDUNOV, S., CHEN, D., AND TAU YIH, W. Dense passage retrieval for open-domain question answering, 2020.
- [14] KWIATKOWSKI, T., PALOMAKI, J., REDFIELD, O., COLLINS, M., PARIKH, A., ALBERTI, C., EPSTEIN, D., POLOSUKHIN, I., KELCEY, M., DEVLIN, J., LEE, K., TOUTANOVA, K. N., JONES, L., CHANG, M.-W., DAI, A., USZKOREIT, J., LE, Q., AND PETROV, S. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics* (2019).
- [15] LASKAR, M. T. R., HUANG, J. X., AND HOQUE, E. Contextualized embeddings based transformer encoder for sentence similarity modeling in answer selection task. In *Proceedings of the 12th Language Resources and Evaluation Conference* (Marseille, France, May 2020), European Language Resources Association, pp. 5505–5514.
- [16] LEE, K., CHANG, M.-W., AND TOUTANOVA, K. Latent retrieval for weakly supervised open domain question answering, 2019.
- [17] LEWIS, M., LIU, Y., GOYAL, N., GHAZVININEJAD, M., MOHAMED, A., LEVY, O., STOYANOV, V., AND ZETTLEMOYER, L. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension, 2019.
- [18] LIU, Y., OTT, M., GOYAL, N., DU, J., JOSHI, M., CHEN, D., LEVY, O., LEWIS, M., ZETTLEMOYER, L., AND STOYANOV, V. Roberta: A robustly optimized bert pretraining approach, 2019.
- [19] OGUZ, B., CHEN, X., KARPUKHIN, V., PESHTERLIEV, S., OKHONKO, D., SCHLICHTKRULL, M., GUPTA, S., MEHDAD, Y., AND YIH, S. Unified open-domain question answering with structured and unstructured knowledge, 2020.
- [20] PETERS, M. E., NEUMANN, M., AU2, R. L. L. I., SCHWARTZ, R., JOSHI, V., SINGH, S., AND SMITH, N. A. Knowledge enhanced contextual word representations, 2019.
- [21] PETRONI, F., PIKTUS, A., FAN, A., LEWIS, P., YAZDANI, M., CAO, N. D., THORNE, J., JERNITE, Y., KARPUKHIN, V., MAILLARD, J., PLACHOURAS, V., ROCKTÄSCHEL, T., AND RIEDEL, S. Kilt: a benchmark for knowledge intensive language tasks, 2021.
- [22] PETRONI, F., ROCKTÄSCHEL, T., LEWIS, P., BAKHTIN, A., WU, Y., MILLER, A. H., AND RIEDEL, S. Language models as knowledge bases?, 2019.
- [23] RAFFEL, C., SHAZEER, N., ROBERTS, A., LEE, K., NARANG, S., MATENA, M., ZHOU, Y., LI, W., AND LIU, P. J. Exploring the limits of transfer learning with a unified text-to-text transformer, 2020.
- [24] RAJPURKAR, P., ZHANG, J., LOPYREV, K., AND LIANG, P. Squad: 100,000+ questions for machine comprehension of text, 2016.

- [25] ROBERTS, A., RAFFEL, C., AND SHAZEER, N. How much knowledge can you pack into the parameters of a language model?, 2020.
- [26] SUN, L., HASHIMOTO, K., YIN, W., ASAI, A., LI, J., YU, P., AND XIONG, C. Adv-bert: Bert is not robust on misspellings! generating nature adversarial samples on bert, 2020.
- [27] TALMOR, A., ELAZAR, Y., GOLDBERG, Y., AND BERANT, J. olympics – on what language model pre-training captures, 2020.
- [28] VASWANI, A., SHAZEER, N., PARMAR, N., USZKOREIT, J., JONES, L., GOMEZ, A. N., KAISER, L., AND POLOSUKHIN, I. Attention is all you need, 2017.
- [29] VIG, J. A multiscale visualization of attention in the transformer model, 2019.
- [30] VOORHEES, E., AND TICE, D. *The TREC-8 Question Answering Track Evaluation*. Department of Commerce, National Institute of Standards and Technology, 1999, pp. 77–82.
- [31] WANG, M., SMITH, N., AND MITAMURA, T. What is the jeopardy model? a quasi-synchronous grammar for qa. pp. 22–32.
- [32] YANG, Y., YIH, W.-T., AND MEEK, C. WikiQA: A challenge dataset for open-domain question answering. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing* (Lisbon, Portugal, Sept. 2015), Association for Computational Linguistics, pp. 2013–2018.
- [33] YANG, Z., DAI, Z., YANG, Y., CARBONELL, J., SALAKHUTDINOV, R., AND LE, Q. V. Xlnet: Generalized autoregressive pretraining for language understanding, 2020.
- [34] ZHANG, Z., HAN, X., LIU, Z., JIANG, X., SUN, M., AND LIU, Q. Ernie: Enhanced language representation with informative entities, 2019.

## Checklist

1. For all authors...
  - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
  - (b) Did you describe the limitations of your work? [\[Yes\]](#)
  - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) We discuss in our conclusion the relative inability of our approach to leverage the full richness of our KB and that it thus requiring significant fine-tuning in light of the energy consumed by our approach.
  - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)
2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [\[N/A\]](#)
  - (b) Did you include complete proofs of all theoretical results? [\[N/A\]](#)
3. If you ran experiments...
  - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[No\]](#) The code for the augmentation pipeline is proprietary as well as the resulting augmented datasets, which contain the proprietary information added by our pipeline. We are additionally not able to release the resulting model binaries, which have encoded this additional information.
  - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#) We use the provided splits for the publicly available datasets referenced in A. We give statistics of the 80/10/10 split used for the Alexa QA dataset, though opt not justify this selection. We use the hyperparameters selected by Garg et. al. [7] in order to accurately assess the added benefits of our approach over their state of the art performance.

- (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[No]** We ran each experiment multiple times with different seeds and observed no statistically significant difference in the outcome. This result is consistent with the robustness of the TandA approach studied by Garg et. al. in [7]. We do not include error bars as this requires at least a page of graphs that we omit for brevity.
- (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** 3.4
- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
  - (a) If your work uses existing assets, did you cite the creators? **[Yes]** We cite [7] and [32], creators of the ASNQ and WikiQA datasets respectively.
  - (b) Did you mention the license of the assets? **[Yes]** We reference the licenses of the publicly available datasets used in 3.1
  - (c) Did you include any new assets either in the supplemental material or as a URL? **[No]**
  - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **[No]** We don’t discuss the consent of the Alexa data used in this study because we are not able to release this data.
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[Yes]** We note that the dataset has been stripped of PII by our expert annotators.
- 5. If you used crowdsourcing or conducted research with human subjects...
  - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? **[N/A]**
  - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? **[N/A]**
  - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? **[N/A]**

## A Dataset statistics

The table below shows the distribution of the datasets studied in this by each split, such as *train*, *dev*, and *test*. These demonstrate that our pipeline is able to tag at least one KB entry in each input QA pair, indicating that our simple tagging method is effective at producing KB context.

Dataset	#QA pairs	% w/o KB	#Correct w/ KB	#Incorrect w/ KB
ASNQ Dev	276,809	.020%	1,117	275,692
ASNQ Test	879,594	.036%	3,600	875,672
ASNQ Train	29,987,324	.027%	120,184	29,867,166
WikiQA Dev	1,130	.000%	140	990
WikiQA Test	2,351	.000%	293	2,507
WikiQA Train	8,672	.000%	1,040	7,632
TrecQA Dev	1,117	.000%	205	912
TrecQA Test	1,442	.000%	248	1,194
TrecQA Train	53,417	.000%	6,403	47,011
AlexaQA Dev	26,951	.040%	25,822	1,192
AlexaQA Test	26,965	.000%	25,796	1,169
AlexaQA Train	215,416	.635%	205,070	8,978

Table 6: Dataset Statistics and KB Tag Rate by Split

## B Attention Weight Comparison

In the graphs below, we illustrate the impact of our approach on model attention for the challenging AS2 examples presented in Table 1. We do not add KB context at inference for any of these examples,

opting to visualize the impact of our approach in the more challenging "omit KB" setting. We leverage BertViz [29] to extract model attention weights and quantify model attention between meaningful keywords selected in question and answer texts. Box plots, shown on the left, quantify the intensity of model attention across all layers, while bar plots, shown on the right, quantify the number of heads per layer exhibiting attention weights greater than an arbitrary minimum of 0.1.

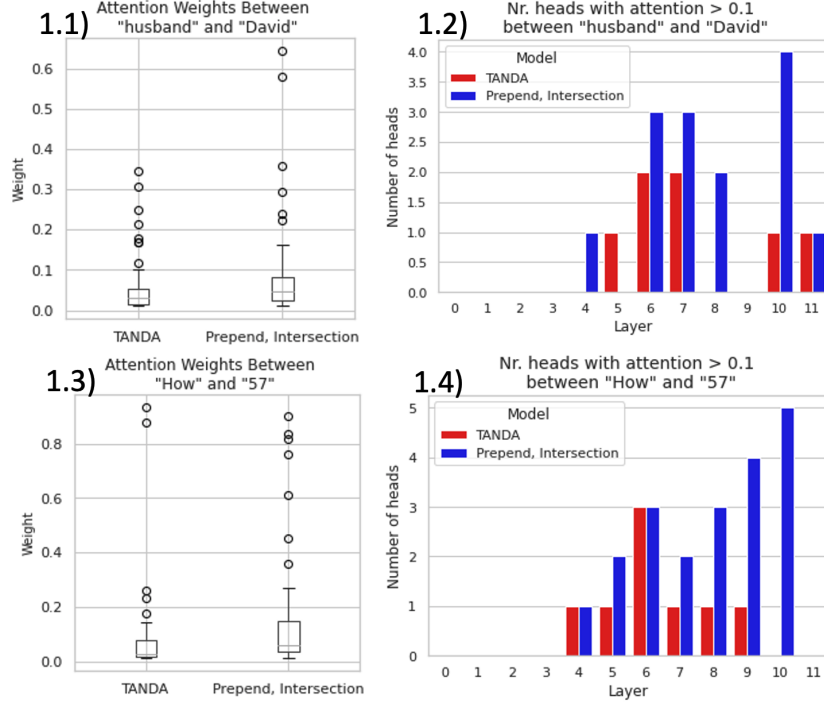


Figure 2: Attention comparison for the correct QA pair *Q: How old is Elton John's husband A: David Furnish is 57 years old. He was born on October 25, 1962*

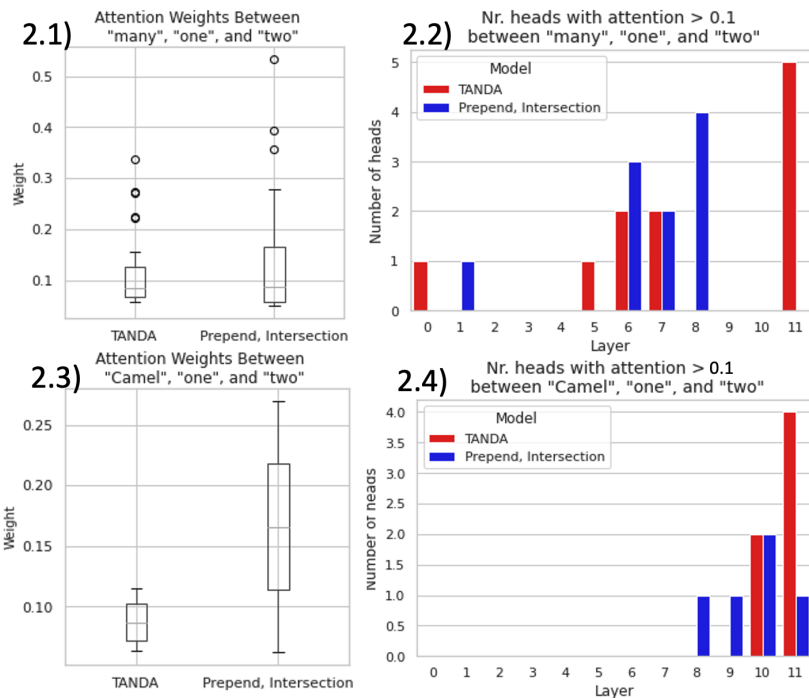


Figure 3: Attention comparison for the correct QA pair *Q: How many humps on a Camel? A: The two surviving species of camel are the dromedary, or one-humped camel, which is native to the Middle East and the Horn of Africa; and the Bactrian, or two-humped camel, which inhabits Central Asia.*

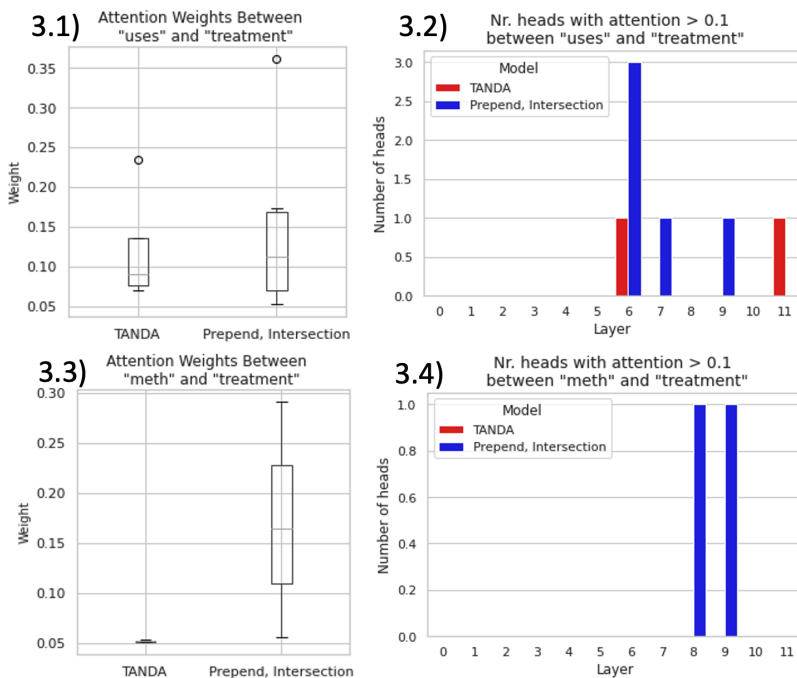


Figure 4: Attention comparison for the correct QA pair *Q: What some legal uses of meth? A: Although rarely prescribed, methamphetamine hydrochloride is approved by the U.S. Food and Drug Administration (FDA) for the treatment of attention deficit hyperactivity disorder and obesity under the trade name Desoxyn.*