

Hey AI, Can You Solve Complex Tasks by Talking to Agents?

Anonymous ACL submission

Abstract

Training giant models from scratch for each complex task is resource- and data-inefficient. To help develop models that can leverage existing systems, we propose a new challenge: Learning to solve complex tasks by communicating with existing agents (or models) in natural language. We design a synthetic benchmark, COMMAQA, with three complex reasoning tasks (explicit, implicit, numeric) designed to be solved by communicating with existing QA agents. For instance, using text and table QA agents to answer questions such as "Who had the longest javelin throw from USA?". We show that black-box models struggle to learn this task from scratch (accuracy under 50%) even with access to each agent's knowledge and gold facts supervision. In contrast, models that learn to communicate with agents outperform black-box models, reaching scores of 100% when given gold decomposition supervision. However, we show that the challenge of learning to solve complex tasks by communicating with existing agents *without relying on any auxiliary supervision or data* still remains highly elusive. We will release COMMAQA, along with a compositional generalization test split, to advance research in this direction.

1 Introduction

A common research avenue pursued these days is to train monolithic language models with billions of parameters to solve every language understanding and reasoning challenge. In contrast, humans often tackle complex tasks by breaking them down into simpler sub-tasks, and solving these by interacting with other people or automated agents whose skill-sets we are familiar with. This approach allows us to learn to solve new complex tasks quickly and effectively, by building upon what's already known. Can AI systems learn to do the same?

To facilitate research in this direction, we propose a new reasoning challenge and a benchmark called COMMAQA where, in addition to the usual

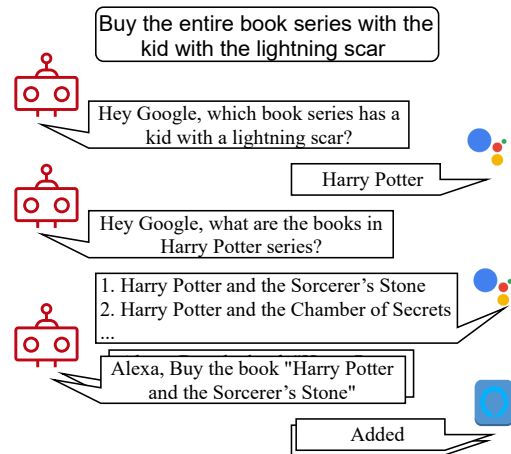


Figure 1: Motivating example for a setup where a system is expected to learn to accomplish goals by interacting with agents via a natural language interface.

end-task supervision, one has access to a set of pre-defined AI agents with examples of their natural language inputs.¹ Importantly, the target end-task is designed to be too difficult for current models to learn based only on end-task supervision. The goal is instead to build models that learn to solve the target task by decomposing it into sub-tasks solvable by these agents, and interacting with these agents in natural language to do so.

As a motivating example, consider the interaction depicted in Figure 1 where a system is asked to buy a book series with a certain property. The system breaks this goal down, using agent-1 (here Google Assistant) to identify the referenced book series as well as the list of books in that series, and then using agent-2 (here Amazon Alexa) to make the purchase. While both of these agents interact with the system in natural language, they have different and complementary skill sets,² rely on privately held knowledge sources, and have been built at an enormous cost. At the same time, neither agent by itself can accomplish the original goal.

¹Our benchmark will be released upon publication.

²but not necessarily mutually exclusive skills

065 An alternative to building such a system that in- 116
066 teracts with existing agents is to teach all requisite 117
067 sub-tasks and skills to a large black-box system, 118
068 say via multi-task learning (Khashabi et al., 2020; 119
069 Gupta et al., 2021). This, however, not only wastes 120
070 time and resources, but is often also infeasible. For
071 example, agents such as Google Assistant and Ope-
072 nAI GPT-3 use private knowledge resources and
073 are computationally expensive to train even once.
074 It would thus be nearly impossible to build a single
075 system with the capabilities of both of these agents.

076 We note that agents need not be sophisticated 121
077 AI assistants. An agent may simply be a previ- 122
078 ously developed question-answering (QA) model, 123
079 a math module, a function of textual input, an im- 124
080 age captioning system—anything the community 125
081 already knows how to build. The goal is to *learn to* 126
082 *leverage existing agents for more complex tasks.* 127

083 To enable the development of general systems 128
084 for this task, we identify the minimal inputs that 129
085 must be assumed for the task to be learnable— 130
086 training data for the complex task, existing agents 131
087 that together can solve the complex task, and ex- 132
088 amples of valid questions that can be asked of 133
089 these agents (capturing the agents’ capabilities). 134
090 We build a new synthetic benchmark dataset called 135
091 COMMAQA (Communicating with agents for QA), 136
092 containing three complex multihop QA tasks (in- 137
093 volving Explicit, Implicit, and Numeric reasoning) 138
094 and four input QA agents that can solve these tasks. 139

095 COMMAQA is not yet another multi-hop reading 140
096 comprehension dataset. It is designed to facilitate 141
097 the development of a new family of techniques that 142
098 teach systems to communicate with a wide variety 143
099 of agents to solve different types of complex tasks. 144

100 We demonstrate that black-box models struggle 145
101 on COMMAQA even when provided with auxil- 146
102 iary data, such as domain-relevant agent knowl- 147
103 edge. On the other hand, a model that leverages 148
104 the agents (Khot et al., 2021) can achieve very high 149
105 accuracy but relies on auxiliary supervision (de- 150
106 composition annotations). While it is possible to 151
107 identify valid decompositions using just the end- 152
108 task labels, the search space is extremely large and 153
109 naïve approaches, as we show, help only with one 154
110 of the datasets. COMMAQA thus serves as a new 155
111 challenge for the NLP community. 156

112 **Contributions:** We (1) propose a new challenge 157
113 of learning to solve complex tasks by communicat- 158
114 ing with agents; (2) develop a synthetic multi-hop 159
115 QA dataset COMMAQA with three reasoning types; 160

(3) provide auxiliary training data and a composi- 116
117 tional generalization test set; (4) demonstrate the 118
119 challenging nature of COMMAQA for black-box 120
121 models; and (5) show the promise of compositional 122
123 models that learn to communicate with agents. 124

2 Related Work 121

122 **Semantic Parsing** typically focuses on mapping 123
124 language problems to executable symbolic repre- 125
126 sentation based on a pre-defined grammar (Krish- 127
128 namurthy et al., 2017; Chen et al., 2020). Similar 129
130 ideas are also found in the area of program syn- 131
132 thesis (Gulwani, 2011; Desai et al., 2016). These 133
134 goals, like ours, seek to simplify complex prob- 135
136 lems into simpler executable forms, without relying 137
138 on explicit intermediate annotation (Clarke et al., 138
139 2010; Berant et al., 2013). We, however, diverge 140
141 from this line by seeking agent communication in 142
143 free-form language, not bound to any pre-specified 144
145 set of operations or domain specific languages. 146

147 **Multi-hop QA** focuses on reasoning with mul- 148
149 tiple facts. Despite the development of many 150
151 datasets (Khashabi et al., 2018; Yang et al., 2018; 152
153 Khot et al., 2020; Geva et al., 2021) and mod- 154
155 els (Min et al., 2019b; Pan et al., 2021), exist- 156
157 ing benchmarks often contain single-hop short- 158
159 cuts (Min et al., 2019a), resulting in brittle mod- 159
160 els (Gardner et al., 2020) and little progress to- 160
161 wards true multi-hop reasoning (Trivedi et al., 161
162 2020). Additionally these datasets often contain 162
163 sub-problems not solvable by existing models, fur- 163
164 ther disincentivising the development of composi- 164
165 tional models (Khot et al., 2021). 165

166 **Question Decomposition** is used to solve multi- 167
168 hop QA but the resulting models (Talmor and Be- 168
169 rant, 2018; Min et al., 2019b; Perez et al., 2020; 169
170 Khot et al., 2021) are often dataset-specific, rely on 170
171 decomposition annotations, and limited to one or 171
172 two QA agents. To address these limitations, our 172
173 proposed challenge covers three dataset types and 173
174 four agents. Additionally, models are expected to 174
175 learn to decompose the task by interacting with the 175
176 agents, rather than relying on human annotations. 176

177 **Synthetic Reasoning Challenges** have recently 177
178 been proposed (Lake and Baroni, 2018; Sinha et al., 178
179 2019; Clark et al., 2020) to help systematically 179
180 identify the weaknesses of existing models and in- 180
181 spire modeling innovation (Liu et al., 2021). Our 181
182 new tasks are unique and focus on simulating com- 182
183 plex agent interaction to motivate the development 183
184 of decomposition-based modeling approaches. 184
185 185

3 Challenge Task Definition

We formalize the new challenge task of *learning to talk with agents to solve complex tasks*. To ensure generality of solutions, we identify minimal inputs for the task to be well-defined and learnable.

First we must define f_i , the agents or models that solve simpler sub-tasks.³ Minimally, we need to define the space of valid inputs \mathcal{L}_i for each agent f_i , i.e., how can they be invoked. For a system to identify the appropriate agent for each sub-task, we also need to define the capabilities of each agent. Since these agents are often defined for natural language tasks, the space of inputs captures the capabilities of these agents too. For instance, "Buy the book 'Harry Potter and the Sorcerer's Stone'" captures the Alexa agent's capability of buying books. Instead of complex formal specifications of the agent's capabilities, we use natural language inputs as a rich and convenient representation.

Next, we need a target task \mathcal{T} that can be solved via a composition of the capabilities of $\{f_i\}$.⁴ Finally, to pose this as a machine learning problem, we need training data $\mathcal{D} = \{(x_k, y_k)\}_{k=1}^N$ for \mathcal{T} . Since collecting annotations for complex tasks can be difficult, \mathcal{D} is expected to be relatively small. Models must therefore use the available agents, instead of learning the complex task from scratch.

Given these pre-requisites, we can define the challenge task as follows:

Challenge: Learn a model to solve a complex task \mathcal{T} , given only:

- Training dataset $\mathcal{D} = \{(x_k, y_k)\}_{k=1}^N$ for \mathcal{T} ;
- Agents $\{f_1, \dots, f_m\}$ that can help solve \mathcal{T} ;
- Examples from the space \mathcal{L}_i of valid inputs for each agent f_i that captures its capabilities.

One example of this challenge is answering multi-hop questions given two agents: an open-domain TextQA agent f_1 and an open-domain TableQA agent f_2 . Agent f_1 can use large textual corpora to answer questions such as "Who directed Kill Bill?". Agent f_2 can use tables (e.g., Filmography tables) to answer questions such as "List the movies directed by Quentin Tarantino". Finally, the training data \mathcal{T} for the complex task would contain examples such as ("What movies has the director of Kill Bill appeared in?", ["Reservoir Dogs", ...,]).

³As mentioned earlier, we use *agents* to refer interchangeably to models, assistants, or functions that take free-text as input and produce free-text as output.

⁴Existing datasets lack this requirement, making it impossible to focus only on the agent communication aspect.

Auxiliary Information. Apart from the above minimal pre-requisites, in some cases we may be able to obtain additional training supervision, or additional data about the internals of the agents. We emphasize that such auxiliary information may not always be available (e.g., when using a proprietary agents such as Alexa). A general-purpose system should thus ideally learn to solve this challenge without relying on it. However, it's acceptable for initial methods to use some auxiliary signals as stepping stones toward more general systems.

We consider two kinds of such information—*auxiliary supervision* for the complex task's training examples $(x_k, y_k) \in \mathcal{D}$, and *auxiliary data* about the agents $\{f_i\}$ themselves (not tied to \mathcal{D}).

For auxiliary supervision, we consider having access to annotated decomposition \mathcal{D}_k of a complex task training input x_k into valid inputs for various agents. We also consider annotated gold facts \mathcal{F}_k that could be used to answer x_k .

For auxiliary data, we consider having access to the training data used to build the agents, or the underlying knowledge base \mathcal{K}_i used by them (and possibly even a question-specific relevant subset \mathcal{K}_{ik}). In the example above, \mathcal{K}_i would be equivalent to the entire text and table corpora used by the agents, and \mathcal{K}_{ik} could be the texts and tables relevant to the *movie* domain. Such information can be used to train a stronger black-box model on the end-task, e.g. fine-tuning on the agent's training data first or using the gold facts to identify relevant context. These approaches that circumvent the agents are not the target of our dataset, but we nevertheless evaluate them to highlight their limits.

In summary, we have the following potential auxiliary information as stepping stones:

Auxiliary Supervision for $(x_k, y_k) \in \mathcal{D}$:

- Gold Decomposition \mathcal{D}_k for x_k
- Gold Knowledge \mathcal{F}_k for x_k

Auxiliary Data for agents $\{f_i\}$:

- Training data $\mathcal{D}_{f_i} = \{(u_{ij}, v_{ij})\}_{j=1}^M$ for agent f_i , where $u_{ij} \in \mathcal{L}_i$ and $v_{ij} = f_i(u_{ij})$
- Complete knowledge resource \mathcal{K}_i used by f_i , or a manageable subset $\mathcal{K}_{ik} \subset \mathcal{K}_i$ containing \mathcal{F}_k

4 Dataset: COMMAQA Benchmark

We next propose a new benchmark dataset COMMAQA that enables the development of models that can learn to communicate with existing agents. Specifically, we provide a collection of *three syn-*

| Operator | Pseudo-code | Example |
|---------------|---|---|
| select | return $f_i(q(a))$ | #1=[23, 35] q="Which is largest value in #1?" $f_i = \text{mathqa}$ → 35 |
| project | return $[(x, f_i(q(x))) \text{ for } x \text{ in } a]$ | #1=[Jordan, Johnson] q="What were the lengths of throw by #1?" $f_i = \text{textqa}$ → [(Jordan, [23, 34]), (Johnson, [45, 56])] |
| projectValues | return $[(k, f_i(q(v))) \text{ for } (k, v) \text{ in } a]$ | #1=[(Jordan, [23, 34]), (Johnson, [45, 56])] q="Which is largest value in #1?" $f_i = \text{mathqa}$ → [(Jordan, 34), (Johnson, 56)] |
| filter | return $[x \text{ for } x \text{ in } a \text{ if } f_i(q(x))]$ | #1=[23, 34, 56] q="Is #1 greater than 50?" $f_i = \text{mathqa}$ → [56] |
| filterValues | return $[(k, v) \text{ for } (k, v) \text{ in } a \text{ if } f_i(q(v))]$ | #1=[(Jordan, 34), (Johnson, 56)] q="Is #1 greater than 50?" $f_i = \text{mathqa}$ → [(Johnson, 56)] |

Table 1: Compositional Operators used in this work to transform structured answers into queries answerable by an agent. The operator takes the agent f_i , a structured answer a (we use the answer index, e.g. #1, to refer to any answer), and a query with a placeholder as inputs and executes the pseudo-code shown here.

thetic datasets where each question is answerable by talking to simple QA agents. Note that we are not proposing a new class of questions but a new dataset for the proposed challenge task.

We choose QA as the underlying task and use QA agents for this challenge because the question-answer format can capture a broad range of tasks (Gardner et al., 2019) while also naturally surfacing the capability of each agent. For instance, the question "What are the key frames in v ?" describes a capability of the invoked agent (namely, identifying key frames), in addition to the specific inputs. We next describe our framework for building COMMAQA, which we believe can be extended to other complex tasks, e.g., video summarization.

4.1 Agent Definition

To define the i -th agent, we build a knowledge base that captures its internal knowledge resource \mathcal{K}_i . We use natural language question templates to define the set of questions that this agent can answer over this internal knowledge. For example, given a KB with relations such as "directed(x, y)", the agent would answer questions based on the template: "Who directed the movie $_$?"

Knowledge Base, \mathcal{K}_i . To build the knowledge base, we define a KB schema as a set of binary relations between entity types, e.g., director(movie, person). We build a list of entity names that belong to each entity type. To avoid potential conflicts with the LM’s pre-training knowledge, all entity names are generated non-existent words.⁵

Rather than building a static and very large KB, we sample *a possible world* independently for each question, by sub-sampling entities for each entity type and then randomly assigning the KB relations

between these entities. This prevents memorization of facts across the train and test splits, which in the past has led to over-estimation of QA model performance (Lewis et al., 2021). This also encourages models to learn proper multi-hop reasoning using the agents, rather than memorizing answers.

Examples of Valid Inputs. To define the space of valid inputs for each agent f_i , we define a set of question templates that can be answered by it over \mathcal{K}_{ik} (e.g., Who directed $_$?). We construct questions corresponding to a relation in both directions, e.g., "Who all directed $_$?" and "For which movies was $_$ a director?". To emulate redundancy in natural language, we specify multiple phrasings for the same question. We use these templates to generate examples of valid inputs in \mathcal{L}_i by grounding them with entities of the appropriate entity type (e.g., Who directed Kill Bill?).

To ensure generalization to a broad set of tasks, we do not limit the questions to only single span answers. Depending on the question, the agent can produce answers as a single string (span, boolean or a number), a list of strings (e.g., "Which movies did Spielberg direct?"), or a map (e.g., "What are the states and their capitals in USA?").

Implementation. To answer the question, agents convert questions into queries against the internal knowledge (based on the templates) which we implement as a symbolic function (written in Python), instead of a model. While a language model might be able to generalize to out-of-distribution variations in language, its behavior can be often unpredictable. By implementing the agents as pattern-based functions, we ensure that the resulting systems would stay within the language constraints of each agent and generalize to restricted language models. Additionally, this enables faster develop-

⁵<https://www.thisworddoesnotexist.com/>

ment of approaches without spending resources on running a large-scale LM for each agent.

4.2 Complex Task Definition

Given the space of valid input questions for each agent, we construct training examples for the complex task using templated theories. These theories consist of a complex question template and a composition rule expressed as a sequence of questions asked to appropriate agents. For example,

```
"What movies have the directors from $1 directed?"
#1 = [textqa] "Who is from the country $1?"
#2 = [tableqa] "Which movies has #1 directed?"
```

Composition Operators. While this simple theory would work for single span answers, these agents often return list or map answers. Even within this simple example, there can be multiple directors from a given country and this list can not be directly fed to the tableqa model, i.e., "Which movies has [...] directed?". This problem gets even more challenging with complex structures. E.g., maintaining a map structure while operating on the values of the map (see 3rd row in Table 1).

To handle the different answer structures, we define a special set of compositional operators in Table 1. These operators take agent f_i , a structured answer a , and a query with a placeholder as inputs, and execute a set of queries (as defined by the pseudo-code in Table 1) against f_i . These operators are inspired by QDMR (Wolfson et al., 2020), but modified to be actually *executable*. E.g., the "project" operator in QDMR: "return directors of #1?" does not specify how to execute this query whereas our operation (project) [textqa] "Who are the directors of #1?" specifies how to use the TextQA model and #1 to generate a map.

We also define a set of agent-independent data structure transformations in Table 2, e.g., convert a map into a list of its keys. Since longer chains of reasoning are prone to more errors (Fried et al., 2015; Khashabi et al., 2019), we don't model these simple transformations as additional reasoning steps. Instead, we concatenate compositional operators with transformations to create about 20 new, combined operators such that transformations can be applied after an operation in a single step, e.g., project_values operation performs the project operation followed by the Values transformation.

Given these operators, the final theory for the above example would look like:

| Transf. | Procedure |
|---------|--|
| FLAT | Flatten list of lists into a single list |
| UNIQUE | Return the unique items from a list |
| KEYS | Return the list of keys from a map |
| VALUES | Return the list of values from a map |

Table 2: Simple transformations that modify the output data structure. These transformations can be chained together with an operation, e.g., PROJECT_VALUES.

```
"What movies have the directors from $1 directed?"
#1 = (select) [textqa] "Who is from the country $1?"
#2 = (project_values_flat_unique) [tableqa] "Which movies
has #1 directed?"
```

Building Examples. Given a KB schema, question templates for each agent, and theories, we can now build examples for the complex task (Fig. 2). We first sample a possible world based on the KB schema. We assign each relation to one of the agents based on which agents are likely to answer such questions, i.e., only this agent would answer questions about this relation. This captures multimodality of knowledge, e.g., movie awards might be described in text or a table, but a person's birth date is likely described in text. When a relation can be captured by knowledge in multiple modalities, it is assigned to one of them per KB. This emulates the challenging setting where a model must interact with multiple agents to find the answer.⁶ We use the templated theories to construct questions by grounding placeholders. We select m valid questions⁷ for each KB such that each theory has the same number of examples across the dataset.

As a stepping stone towards solving the full challenge and to evaluate the limits of current baselines given oracle supervision, we provide auxiliary information as mentioned in Sec. 3. Specifically, we provide the gold decomposition \mathcal{D}_k for each example x_k using the same language as the theories (see Fig. 3). We verbalize each relation to create the underlying knowledge resource \mathcal{K}_{ik} used by the agent f_i (e.g., relation director(M, P) is converted into "M was a movie directed by P" or "movie: M ; director: P" depending on the agent assigned to this relation). While our KB and resulting facts are intentionally simple to show the limitations of black-box models, such verbalization may not always be possible with larger KBs. For each training example, we collect the facts used by each agent in the decomposition and treat these as gold facts \mathcal{F}_k .

⁶With real questions and agents, models may be able to avoid this by just memorizing the agents.

⁷has a non-empty answer and up to five answer spans

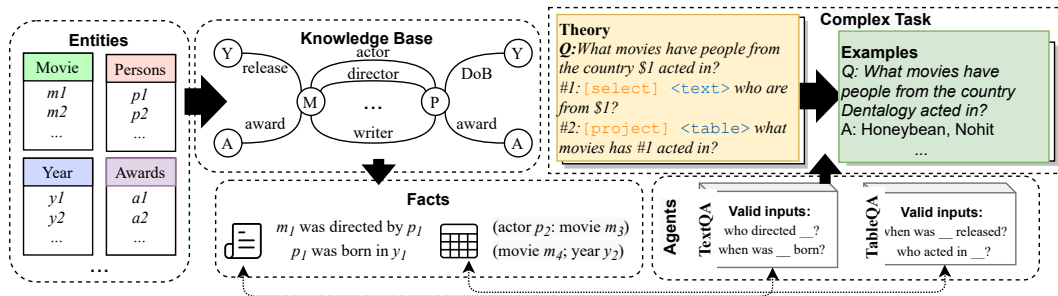


Figure 2: High-level schema of our dataset construction process. We use a list of entities and a KB schema to generate a list of facts. The QA agents operate over these facts to answer a set of pre-determined questions that form the examples of valid inputs from \mathcal{L}_i . We define multiple complex question templates and a corresponding *theory* that can be used to answer them. We then ground these question templates (i.e. sample \$1) to create complex questions and use the agents to generate the answers.

4.3 COMMAQA Dataset

We use the above framework to build three datasets capturing three challenges in multi-hop reasoning.

COMMAQA-E: Explicit Decomposition. This dataset consists of multi-hop questions from the movie domain where the reasoning needed to answer the question is Explicitly described in the question itself (Yang et al., 2018; Ho et al., 2020; Trivedi et al., 2021). For example, "What awards have the movies directed by Spielberg won?". We use a TextQA and TableQA agent where certain relations can be either be expressed in text or table (more details in App. Fig. 5).

COMMAQA-I: Implicit Decomposition. This dataset consists of multi-hop questions where the reasoning needed is Implicit (Khot et al., 2020; Geva et al., 2021), for example, "Did Aristotle use a laptop?". Inspired by such questions in StrategyQA (Geva et al., 2021), we create this dataset using three agents (TextQA, KBQA and MathQA) with just two question styles: (1) "What objects has __ likely used?" and (2) "What objects has __ helped make?". However each question has three possible strategies depending on the context (see App. Fig. 6 for more details). This is a deliberate choice as similar sounding questions can have very different strategies in a real world setting, e.g., "Did Steve Jobs help develop an Iphone?" vs. "Did Edison help develop the television?".

COMMAQA-N: Numeric Decomposition. This dataset consists of Numeric (also referred to as discrete) reasoning questions (Dua et al., 2019; Amini et al., 2019) requiring some mathematical operation, in addition to standard reasoning. For example, "Who threw javelins longer than 5 yards?". We create this dataset in the sports domain with TextQA, TableQA and MathQA agents (more details in App. Fig. 7).

Dataset Statistics. The final dataset⁸ consists of the three QA sub-datasets described above (key statistics in Table 7 in the App.). We have 10K total examples in each dataset with 80%/10%/10% train/dev/test split. To prevent models from guessing answer spans, we introduce more distractors by sampling a large number of facts for COMMAQA-E and COMMAQA-I. This results in a larger number of facts in the KB (~ 170) and larger length of the KB in these two datasets (~ 2500 tokens). Since COMMAQA-N can have derived answers from numeric reasoning and has longer chains (avg #steps 4.7 vs. 2.7 in COMMAQA-E), we do not need a large number of distractor facts (80 facts/KB).

Metrics. The answer y_k to each question x_k in COMMAQA is an unordered list of single-word entities.⁹ By the design of the dataset, a model that performs the desired reasoning should be able to output y_k correctly, barring entity permutation. Hence, we use *exact match accuracy* as the metric.¹⁰ (see appendix for a softer metric, F1 score)

5 Experiments

We next evaluate state-of-the-art models on COMMAQA, with and without auxiliary information.

5.1 Models

Models with Access to Agent Knowledge: Given access to the facts associated with each (train or test) question x_k , i.e., *each agent's domain-relevant knowledge* \mathcal{K}_{ik} , the facts can be concatenated to create a context and frame the challenge as a reading comprehension (RC) task.¹¹ We train two stan-

⁸released under CC BY license

⁹Although not in the current dataset, entities in the unordered list y_k may be repeated, i.e., we have a multi-set.

¹⁰Our implementation uses "exact match" in the DROP multi-span evaluator, which accounts for entity reordering.

¹¹We reiterate that it is often unreasonable to expect access to \mathcal{K}_i and especially \mathcal{K}_{ik} . This model tries to solve

| | | |
|--|--|-----------|
| What awards have movies written by people born in 1905 won? (select) [text] Who were born in the year 1905? (project_values_flat_unique) [table] What movies has #1 written? (project_values_flat_unique) [table] Which awards were given to #2? | A: ["Gigafuna"] A: ["Pneumodendron", "Pipesia", "Riften"] A: ["Pludgel", "Dessication", "Pianogram"] | CommaQA-E |
| What objects has Calcid helped to make? (select) [text] Calcid is the founder of which companies? (project_values_flat_unique) [text] #1 produces which materials? (project_values_flat_unique) [text] Which objects use #2 as a material? | A: ["Dufferate"] A: ["comander"] A: ["chickenpot", "yaki"] | CommaQA-I |
| Who threw discuses shorter than 51.8? (select) [text] Who threw discus? (project) [text] What were the lengths of the discus throws by #1? (projectValues) [math_special] What is the smallest value among #2? (filterValues_keys) [math_special] Is #3 less in value than 51.8? | A: ["Lobsteroid", "Karfman", "Terbaryan", ...] A: [["Lobsteroid", "65.6", "46.0"], ["Karfman", ...]] A: [["Lobsteroid", 46.0], ["Karfman", 51.8], ...] A: ["Lobsteroid", ...] | CommaQA-N |

Figure 3: Sample Decomposition Annotations for example questions in COMMAQA. We denote the composition operators using the format (operation) [agent] "question".

dard black-box models, T5-L (Raffel et al., 2020) and UnifiedQA-L (Khashabi et al., 2020),¹² to generate answers¹³ given the question and context.

Models with Fact Supervision: If, in addition to access to the underlying knowledge \mathcal{K}_{ik} , we also have the auxiliary supervision for the gold facts \mathcal{F}_k , we can use this annotation to train a model to first retrieve a small subset of relevant facts from \mathcal{K}_{ik} (see App. D.1 for details). Since the context is shorter, we also train a T5-3B model¹⁴ on this task.

Models with Decomposition Supervision: Given decomposition supervision, we can develop a model that actually solves the task of communicating with the agents. Specifically, we use the Text Modular Network (TMN) framework (Khot et al., 2021) that trains a NextGen model that communicates with the agents. This model is trained to produce the next question (including operation and agent) in a decomposition chain, given the questions and answers so far, which is then executed against the agent to produce the answer for the current step. Additionally this framework samples multiple questions at each step of the chain to search¹⁵ for the most likely chain of reasoning. We refer to this model as TMN-S when we use this search and TMN-G when we greedily select the most likely question at each step.

Models with No Supervision: Finally, we develop a baseline approach that directly targets the challenge task without relying on any auxiliary information. To this end, we generate the training data for NextGen via distant supervision. Specifically, we perform a brute-force search where we

sample l questions at each step for up to o steps.¹⁶ The operations are chosen randomly but we only consider the applicable operations (e.g., "select" for the first step). We use lexical overlap between the questions in the examples of valid inputs and the complex question to avoid wasteful random sampling.¹⁷ We assume all chains that lead to the gold answer represent valid decompositions, and use them to build the training dataset for TMNs. We refer to these generated decompositions as $\hat{\mathcal{D}}_k$. More details are deferred to Appendix B.

5.2 Results

Table 3 reports the accuracy of these four classes of models on the COMMAQA dataset.

Black-box models struggle on COMMAQA: Due to the large number of distractors, black-box models struggle to learn the task across all three datasets with average accuracy below 20. The extremely low performance on COMMAQA-E is especially notable, given that the reasoning needed for each question is explicitly described. While these models are able to solve similar datasets (Yang et al., 2018), the low scores on our synthetic dataset with more distractors indicates that they are still unable to truly learn this kind of reasoning.

Fact annotations help but insufficient: The models trained on shorter context (obtained by relying on gold fact training annotation) are able to take advantage of the reduced number of distractors, improving their score to about 45 pts across all datasets. However, even with the larger 3B model, there is no noticeable improvement, indicating 45 pts being roughly a ceiling for these models.

COMMAQA can be solved by talking to the agents: The TMN model trained on the gold de-

COMMAQA without invoking agents, which deviates from the purpose of our benchmark dataset. Nevertheless, we conduct experiments in this setting for completeness.

¹²We use T5 models as they can handle longer contexts.

¹³We alphabetically sort answers for a deterministic order.

¹⁴T5-11B performed worse than or same as the 3B model.

¹⁵Score is the sum log likelihood of the generated questions.

¹⁶ o is set based on the length of the rules in each dataset, i.e., $o = 3$ for COMMAQA-E, $o = 4$ for I, $o = 7$ for N.

¹⁷We also found random generally performed worse.

| Model | Aux. Info | E | I | N | Avg. |
|-------|---------------------------------------|-------|-------|-------|-------|
| T5-L | $\{\mathcal{K}_{ik}\}$ | 0.9 | 10.2 | 35.4 | 15.5 |
| UQA-L | $\{\mathcal{K}_{ik}\}$ | 1.0 | 10.2 | 39.0 | 16.7 |
| T5-L | $\mathcal{F}_k, \{\mathcal{K}_{ik}\}$ | 42.2 | 49.4 | 44.7 | 45.4 |
| UQA-L | $\mathcal{F}_k, \{\mathcal{K}_{ik}\}$ | 40.1 | 49.7 | 43.4 | 44.4 |
| T5-3B | $\mathcal{F}_k, \{\mathcal{K}_{ik}\}$ | 42.3 | 49.9 | 43.4 | 46.2 |
| TMN-G | \mathcal{D}_k | 75.4 | 36.0 | 100.0 | 70.5 |
| TMN-S | \mathcal{D}_k | 100.0 | 100.0 | 100.0 | 100.0 |
| TMN-S | $\hat{\mathcal{D}}_k(l=5)$ | 0.0 | 0.0* | 0.0 | 0.0 |
| TMN-S | $\hat{\mathcal{D}}_k(l=10)$ | 17.0 | 0.0* | 0.0 | 5.7 |

Table 3: Accuracy of models trained and tested separately on the 3 datasets. Last column reports average accuracy across the datasets (weighed equally). **TOP** half: Black-box models struggle even when given the domain-relevant KB \mathcal{K}_{ik} . Using the additional fact supervision \mathcal{F}_k helps these models, but their accuracy remains below 50%. **BOTTOM** half: TMN models with auxiliary decomposition supervision \mathcal{D}_k can solve all tasks with search ("TMN-S"), showing promise. Solving the full challenge by generating training data $\hat{\mathcal{D}}$ helps on COMMAQA-E but does not result in any valid decomposition (indicated by *) on COMMAQA-I.

composition annotations can solve this task. This experiment shows that COMMAQA is noise-free, unambiguous, and solvable by a model that learns to talk to the agents (as designed). Note that greedily selecting the next question results in much lower performance on the two datasets (E and I) that have multiple decompositions for the same question.

Naïve Search is insufficient: In the final two rows, we evaluate the brute-force search approach to generate training data for TMNs. For COMMAQA-I, we don't find even a single chain that leads to the gold answer, resulting in no training data. With COMMAQA-E and COMMAQA-N, we do find valid decompositions for a subset of the questions (cf. Table 6 in the Appendix for statistics), but they are insufficient to train an effective NextGen model. Expanding the search to $l=20$ helps achieve near 100% accuracy on COMMAQA-E (with $\sim 700K$ agent calls). However, we don't observe any gains on COMMAQA-I and COMMAQA-N with even 2M agent calls (see App. C).

5.3 Compositional Generalization

We also design compositional generalization test sets COMMAQA-E^{CG} and COMMAQA-N^{CG}. Specifically we create questions using novel composition of queries that have been seen during training but never together in this form. For instance, we create a new question "What awards have the di-

rectors of the __ winning movies received?", given that the model was trained on questions such as "What awards have the actors of the __ winning movies received?", "What movies have the directors from __ directed?", and "What movies have people from the country __ acted in?".

| Model | Aux. Info | COMMAQA-E ^{CG} | COMMAQA-N ^{CG} |
|-------|---------------------------------------|-------------------------|-------------------------|
| T5-L | $\mathcal{F}_k, \{\mathcal{K}_{ik}\}$ | 37.0 | 2.0 |
| T5-3B | $\mathcal{F}_k, \{\mathcal{K}_{ik}\}$ | 39.2 | 23.8 |
| TMN-S | \mathcal{D}_k | 79.4 | 97.6 |
| TMN-S | $\hat{\mathcal{D}}_k(l=10)$ | 16.2 | 0.0 |

Table 4: Lower accuracy on compositional generalization test sets. TMN-S with decomposition supervision still outperforms other models.

As shown in Table 4, all models exhibit a drop in accuracy relative to their score in Table 3, but the compositional model trained on gold decomposition still outperforms black-box models. Our error analysis of TMN-S on COMMAQA-E identified this key issue: While TMN-S learns to generalize, it generates questions outside the space of valid agent inputs (e.g., "Who are the directors in the movie __?" vs. "Which movies has __ directed?").

6 Closing Remarks

We motivated a new challenge task of solving complex task by communicating with existing AI agents. Developing approaches for this challenge, we argue, can result in more generalizable and efficient models. Towards this goal, we introduced a new benchmark dataset COMMAQA which involves multi-hop questions with three multi-hop reasoning challenges, all solvable by composing four QA agents. Experiments with state-of-art language models indicated that they struggle to solve COMMAQA, even when provided with agents' internal knowledge. In contrast, a model that is able to learn to communicate with the agents, albeit using annotated decompositions, is able to solve this task. These results point to the need for and the potential of such approaches, but without reliance on auxiliary annotations, to solve complex tasks.

COMMAQA is only one instantiation of our overall framework. One can extend it in many ways, such as using LMs to enrich lexical diversity, emulating the behavior of imperfect real-world agents that even attempt to answer out-of-scope questions, and diversifying to other reasoning types such as Boolean questions where using distant supervision is even harder (Dasigi et al., 2019).

612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665

References

Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. 2019. MathQA: Towards interpretable math word problem solving with operation-based formalisms. In *NAACL*.

Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv:2004.05150*.

Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. Semantic Parsing on Freebase from Question-Answer Pairs. In *EMNLP*.

Xinyun Chen, Chen Liang, Adams Wei Yu, Denny Zhou, Dawn Song, and Quoc V. Le. 2020. Neural symbolic reader: Scalable integration of distributed and symbolic representations for reading comprehension. In *ICLR*.

Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. Transformers as soft reasoners over language. *IJCAI*.

James Clarke, Dan Goldwasser, Ming-Wei Chang, and Dan Roth. 2010. Driving semantic parsing from the world’s response. In *CoNLL*, pages 18–27.

Pradeep Dasigi, Matt Gardner, Shikhar Murty, Luke Zettlemoyer, and Eduard Hovy. 2019. Iterative search for weakly supervised semantic parsing. In *NAACL-HLT*.

Aditya Desai, Sumit Gulwani, Vineet Hingorani, Nidhi Jain, Amey Karkare, Mark Marron, and Subhajit Roy. 2016. Program synthesis using natural language. In *Proceedings of the 38th International Conference on Software Engineering*, pages 345–356.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*.

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 *NAACL*.

Daniel Fried, Peter A. Jansen, Gus Hahn-Powell, Mihai Surdeanu, and Peter E. Clark. 2015. Higher-order lexical semantic models for non-factoid answer reranking. *TACL*, 3:197–210.

Matt Gardner, Yoav Artzi, Jonathan Berant, Ben Bogin, Sihao Chen, Dheeru Dua, Yanai Elazar, Ananth Gotumukkala, Nitish Gupta, Hanna Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Eric Wallace, Ally Quan Zhang, and Ben Zhou. 2020. Evaluating models’ local decision boundaries via contrast sets. In *Findings of EMNLP*.

Matt Gardner, Jonathan Berant, Hannaneh Hajishirzi, Alon Talmor, and Sewon Min. 2019. Question answering is a format; when is it useful? *ArXiv*, abs/1909.11291. 666
667
668
669

Matt Gardner, Joel Grus, Mark Neumann, Oyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke S. Zettlemoyer. 2017. AllenNLP: A deep semantic natural language processing platform. *arXiv preprint arXiv:1803.07640*. 670
671
672
673
674
675

Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. Did Aristotle Use a Laptop? A Question Answering Benchmark with Implicit Reasoning Strategies. *TACL*. 676
677
678
679

Sumit Gulwani. 2011. Automating string processing in spreadsheets using input-output examples. *ACM Sigplan Notices*, 46(1):317–330. 680
681
682

Tanmay Gupta, Amita Kamath, Aniruddha Kembhavi, and Derek Hoiem. 2021. Towards general purpose vision systems. *ArXiv*, abs/2104.00743. 683
684
685

Xanh Ho, A. Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. Constructing a multi-hop qa dataset for comprehensive evaluation of reasoning steps. In *COLING*. 686
687
688
689

Daniel Khashabi, Erfan Sadeqi Azer, Tushar Khot, Ashish Sabharwal, and Dan Roth. 2019. On the possibilities and limitations of multi-hop reasoning under linguistic imperfections. *arXiv: Computation and Language*. 690
691
692
693
694

Daniel Khashabi, Snigdha Chaturvedi, Michael Roth, Shyam Upadhyay, and Dan Roth. 2018. Looking beyond the surface: a challenge set for reading comprehension over multiple sentences. In *NAACL*. 695
696
697
698

Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and Hannaneh Hajishirzi. 2020. UnifiedQA: Crossing format boundaries with a single QA system. In *Findings of EMNLP*. 699
700
701
702
703

Tushar Khot, Peter Clark, Michal Guerquin, Paul Edward Jansen, and Ashish Sabharwal. 2020. QASC: A dataset for question answering via sentence composition. In *AAAI*. 704
705
706
707

Tushar Khot, Daniel Khashabi, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2021. Text modular networks: Learning to decompose tasks in the language of existing models. In *NAACL*. 708
709
710
711

Jayant Krishnamurthy, Pradeep Dasigi, and Matt Gardner. 2017. Neural semantic parsing with type constraints for semi-structured tables. In *EMNLP*. 712
713
714

Brenden Lake and Marco Baroni. 2018. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks. In *ICML*, pages 2873–2882. 715
716
717
718

| | | | |
|-----|--|--|-----|
| 719 | Patrick Lewis, Pontus Stenetorp, and Sebastian Riedel. | Shleifer, et al. 2020. Transformers: State-of-the-art natural language processing. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45. | 773 |
| 720 | 2021. Question and answer test-train overlap in open-domain question answering datasets. In <i>EACL</i> . | | 774 |
| 721 | | | 775 |
| 722 | | | 776 |
| 723 | Nelson F Liu, Tony Lee, Robin Jia, and Percy Liang. | Tomer Wolfson, Mor Geva, Ankit Gupta, Matt Gardner, Yoav Goldberg, Daniel Deutch, and Jonathan Berant. 2020. Break it down: A question understanding benchmark. <i>TACL</i> . | 778 |
| 724 | 2021. Can small and synthetic benchmarks drive modeling innovation? a retrospective study of question answering modeling approaches. <i>arXiv preprint arXiv:2102.01065</i> . | | 779 |
| 725 | | | 780 |
| 726 | | | 781 |
| 727 | | | |
| 728 | 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. <i>arXiv:1907.11692</i> . | Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In <i>EMNLP</i> . | 782 |
| 729 | | | 783 |
| 730 | | | 784 |
| 731 | | | 785 |
| 732 | | | 786 |
| 733 | Sewon Min, Eric Wallace, Sameer Singh, Matt Gardner, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2019a. Compositional questions do not necessitate multi-hop reasoning. In <i>ACL</i> . | Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. 2020. Big Bird: Transformers for longer sequences. <i>NeurIPS</i> , 33. | 787 |
| 734 | | | 788 |
| 735 | | | 789 |
| 736 | | | 790 |
| 737 | Sewon Min, Victor Zhong, Luke S. Zettlemoyer, and Hannaneh Hajishirzi. 2019b. Multi-hop reading comprehension through question decomposition and rescoring. In <i>ACL</i> . | | 791 |
| 738 | | | |
| 739 | | | |
| 740 | | | |
| 741 | Liangming Pan, Wenhua Chen, Wenhan Xiong, Min-Yen Kan, and William Yang Wang. 2021. Unsupervised multi-hop question answering by question generation. In <i>NAACL</i> . | | |
| 742 | | | |
| 743 | | | |
| 744 | | | |
| 745 | Ethan Perez, Patrick Lewis, Wen-tau Yih, Kyunghyun Cho, and Douwe Kiela. 2020. Unsupervised question decomposition for question answering. In <i>EMNLP</i> . | | |
| 746 | | | |
| 747 | | | |
| 748 | | | |
| 749 | Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>Journal of Machine Learning Research</i> , 21:1–67. | | |
| 750 | | | |
| 751 | | | |
| 752 | | | |
| 753 | | | |
| 754 | | | |
| 755 | Koustuv Sinha, Shagun Sodhani, Jin Dong, Joelle Pineau, and William L Hamilton. 2019. CLUTRR: A diagnostic benchmark for inductive reasoning from text. In <i>EMNLP</i> . | | |
| 756 | | | |
| 757 | | | |
| 758 | | | |
| 759 | Alon Talmor and Jonathan Berant. 2018. The web as a knowledge-base for answering complex questions. In <i>NAACL</i> . | | |
| 760 | | | |
| 761 | | | |
| 762 | H. Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2020. Is multihop QA in DiRe condition? Measuring and reducing disconnected reasoning. In <i>EMNLP</i> . | | |
| 763 | | | |
| 764 | | | |
| 765 | | | |
| 766 | H. Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2021. MuSiQue: Multi-hop questions via single-hop question composition. <i>ArXiv</i> , abs/2108.00573. | | |
| 767 | | | |
| 768 | | | |
| 769 | | | |
| 770 | Thomas Wolf, Julien Chaumond, Lysandre Debut, Victor Sanh, Clement Delangue, Anthony Moi, Pierric Cistac, Morgan Funtowicz, Joe Davison, Sam | | |
| 771 | | | |
| 772 | | | |

792 **A Multiple Answers in a Question**

793 If a question refers to multiple answers, e.g. "Is #3
 794 a part of #2?", the operator execution is unclear. To
 795 handle such cases, the operator must specify the
 796 answer to operate over as a parameter. E.g. (**fil-**
 797 **ter(#3)) [mathqa]** "Is #3 a part of #2?" would filter
 798 the answers in #3 whereas (**filter(#2)) [mathqa]** "Is
 799 #3 a part of #2?" would filter the answers in #2.

800 **B Search Approach**

801 We describe the approach used to build the training
 802 data \hat{D} using the simple search technique in more
 803 detail. To generate the space of possible decomposi-
 804 tions, for each question, we first select f operations
 805 from the list of valid operations in Table 5. We only
 806 consider these operations as these are the only oper-
 807 ators needed for COMMAQA. Note that even with
 808 this restricted set of operators, models struggle on
 809 COMMAQA-I and COMMAQA-N. Additionally,
 810 we only consider the select operation for the first
 811 step. For all subsequent steps, we only consider
 812 replacements of $__$ with a previous answer index.

813 To select the questions, we first simplify the
 814 space of inputs by converting the questions into
 815 Fill-In-The-Blank (FITB) questions by removing
 816 the named entities. E.g "Who was born in 1991?"
 817 is changed to "Who was born in $__$?". This is also a
 818 necessary step as the operators need questions with
 819 placeholders to handle structured answers. At every
 820 step, we expand this pool of questions by replacing
 821 the blanks with entities in the complex question
 822 and any answer index from the previous steps (e.g.
 823 #1, #2 in the third step of a decomposition). To
 824 avoid wasteful sampling, we use lexical overlap
 825 between questions in this expanded question pool
 826 and the input question to identify the top g most
 827 relevant questions. The agent associated with each
 828 question is tracked throughout this process.

829 In the end, we consider the cross product be-
 830 tween the f operations and g questions to produce
 831 $l = f \times g$ total questions at each steps. These l
 832 questions are then executed using the appropriate
 833 agent and only the successful questions (i.e. an-
 834 swered by the agent) are considered for the next
 835 step. This is the key reason why the search space
 836 is much smaller than l^o for o reasoning steps.

837 Table 6 presents the overall statistics of the
 838 search approach.

```

select filter
filterValues_keys
filter(__)
filterValues(__)_keys
project
projectValues
projectValues_flat
projectValues_flat_unique
project_values_flat
project_values_flat_unique
    
```

Table 5: Set of operations considered in the search ap-
 proach. $__$ can be replaced by any of the answer indices
 from the previous steps to create a new operation.

| Dataset | NumQs/ Step | Models calls | Num +ve chains | Num qs w/ +ve chains | Dev Acc |
|-----------|----------------|--------------|-------------------|-------------------------|---------|
| CommaQA-E | 5 | 70801 | 246 | 242 | 0 |
| CommaQA-E | 10 | 116595 | 456 | 421 | 17 |
| CommaQA-E | 15 | 541816 | 1325 | 870 | 32.8 |
| CommaQA-E | 20 | 683168 | 2505 | 1669 | 98.9 |
| CommaQA-I | 5 | 81325 | 0 | 0 | 0 |
| CommaQA-I | 10 | 123202 | 0 | 0 | 0 |
| CommaQA-I | 15 | 1149762 | 0 | 0 | 0 |
| CommaQA-I | 20 | 1525736 | 0 | 0 | 0 |
| CommaQA-N | 5 | 94481 | 40 | 27 | 0 |
| CommaQA-N | 10 | 351178 | 46 | 27 | 0 |

Table 6: Statistic of the search-based approach for dif-
 ferent values of l (NumQs/Step). While we get few +ve
 chains for COMMAQA-N, it is not sufficient to train an
 effective model

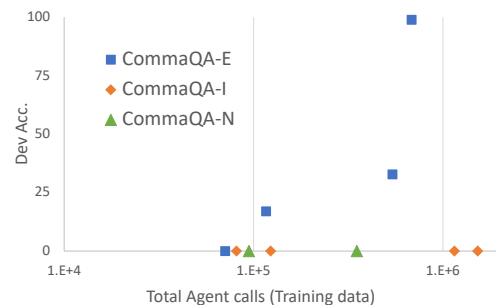


Figure 4: With an order of magnitude increase in
 search space, we can achieve close to 100% accu-
 racy on COMMAQA-E. However COMMAQA-I and
 COMMAQA-N need smarter search strategies to gen-
 erate useful training supervision.

| | COMMAQA | | |
|-------------------------------|---------|--------|--------|
| | E | I | N |
| #questions | 10K | 10K | 10K |
| #theories | 6 | 6 | 6 |
| #steps per theory | 2.7 | 3.2 | 4.7 |
| #entity types | 7 | 13 | 5 |
| #relations | 11 | 16 | 4 |
| #templates in \mathcal{L}_i | 42 | 68 | 30 |
| #entities per answer | 3.21 | 3.29 | 1.36 |
| #KB facts per KB | 169.4 | 175.7 | 80 |
| #T5tokens per KB | 2252.9 | 2540.9 | 1513.4 |
| #Gold facts per qn | 7.5 | 6.9 | 15.4 |

Table 7: Statistics of COMMAQA. All per-question and per-KB statistics are averages.

| EM/ F1 scores | CommaQA | | |
|-------------------------|-------------|----------------|-------------|
| | E | I | N |
| Full Context | | | |
| T5-Large | 0.9 / 30.12 | 10.2 / 25.4 | 35.4 / 38.4 |
| UQA-Large | 1.00 / 30.0 | 10.2 / 25.75 | 39.0 / 41.4 |
| Using Gold Facts | | | |
| T5-Large | 42.2 / 75.5 | 49.9 / 65.5 | 44.7 / 45.3 |
| UQA-Large | 40.1 / 75.3 | 49.7 / 65.8 | 43.4 / 44.8 |
| T5-3B* | 42.3 / 75.7 | 49.9 / 65.6 | 43.4 / 45.3 |
| Decompositions | | | |
| TMN-G | 75.4 / 75.4 | 36 / 36 | 100 / 100 |
| TMN-S | 100 / 100 | 100 / 100 | 100 / 100 |
| TMN-S (l=5) | 0.0 / 0.0 | 0.0 / 0.0 | 0.0 / 0.0 |
| TMN-S (l=10) | 17.0 / 17.1 | 0.0 / 0.0 | 0.0 / 0.0 |

Table 8: EM / F1 scores on the test set using the baseline approaches

C Search Cost vs Accuracy

One could always exhaustively search for *all* possible decompositions to reproduce the gold decompositions for all the questions. But this would be computationally highly expensive as each call to the agent would often invoke a large-scale LM or a complex AI assistant. To characterize the computational cost of these approaches, we extend the search parameter to include $l=15$ and $l=20$ (capped at 5M agent calls) and compute the accuracy of the TMN-S model trained on the resulting dataset (shown in Fig. 4). We can achieve close to 100% accuracy on COMMAQA-E when the search is sufficiently exhaustive (about 700K model calls) mainly due to the shorter rules and the lexical signal. COMMAQA-I and COMMAQA-N, on the other hand, even with an order of magnitude increase in the number of agent calls, we don't observe any increase in the model accuracy.

D Black Box Models

We train the T5 models on each of the three datasets to generate the answer given the question and facts. We format the input sequence as `<concatenated facts> Q: <question> A:.` Since many of the answers can be multiple spans, we sort¹⁸ and concatenate them into a single string with '+' as the separator. As noted in Table 7, the verbalized facts can result in a context over 2K tokens long. We trained T5-Large models on A100 80G GPUs and RTX8000s to train on such a long context. Transformers designed for longer documents (Beltagy et al., 2020; Zaheer et al., 2020) would be able to handle such contexts more efficiently but generally under-perform due to sparse attention. Hence we don't evaluate them here.

For all T5-based models, model tuning was standardly performed using a random hyper-parameter search in the style of Devlin et al. (2019) using the public huggingface implementation (Wolf et al., 2020); model selection was done based on the highest EM accuracy on the development sets. We specifically experimented with learning rates in the range of $(1e-3f$ to $5e-5f)$ using both Adam and Adagrad optimizers and generally found the settings comparable to the original T5 pre-training parameters (Raffel et al., 2020) to be optimal (Adafactor, lr=0.001, 10 epochs, 0-1000 warmup steps, gradient accumulation was used extensively in place of batching to fit long sequences into GPU memory). The optimal T5-3B models and T5-L for full context on COMMAQA-E were trained with lr=5e-5. All other models were trained with a lr of 1e-3. We will release the complete list of optimal hyper-parameters along with the code.

D.1 Models with Fact Supervision

To select the relevant facts, we train a RoBERTa-Large (Liu et al., 2019) model on the gold facts and select the top-scoring facts to produce a shorter context that fits in 512 tokens. The RoBERTa model was trained using the AllenNLP library (Gardner et al., 2017) with the standard parameters used for RoBERTa – learning rate of 2e-5, triangular LR scheduler with 10% warmup steps, gradient clipping at 1.0, batch size of 16, 5 epochs of training with patience of 3 epochs. We didn't observe a noticeable difference in score with a random pa-

¹⁸To ensure a deterministic order, we sort the answers in alphabetical order.

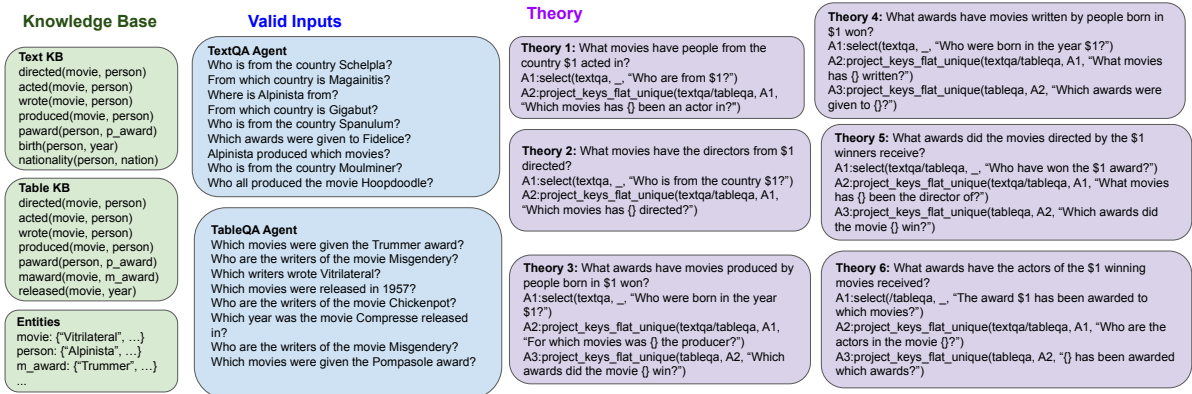


Figure 5: Example KB, space of valid inputs and theory used to construct COMMAQA-E

906 parameter search, so kept these parameters constant.
907 The model was trained to score each fact independ-
908 dently on the train set and the best model was se-
909 lected based on the accuracy on the dev set. The
910 model was then evaluated on the facts from the
911 train, dev and test set to produce the shorter context
912 for all three sets. The facts were sorted based on
913 the model’s scores and the top-scoring facts were
914 added to the context till the number of tokens did
915 not exceed 512 tokens (white-space splitting).

916 E TMN

917 To train the NextGen model for TMNs, we use
918 the same parameters as the prior work (Khot et al.,
919 2021). We train a T5-Large model as the NextGen
920 model using a batch size of 64, lr of 5e-6, 5 epochs
921 and warmup of 1000 steps in all our experiments.
922 We used the public huggingface implementation
923 (Wolf et al., 2020) to train this model. During in-
924 ference, we use a beam size of 10 and select 5
925 questions at each step. We use nucleus sampling
926 with p=0.95 and k=10. For greedy search, we use
927 the same parameters but select one question at each
928 step. We use the sum log likelihood of each gener-
929 ated question as the score of the reasoning chain.
930 All the code will be released on publication for
931 reproducibility.

| | KB | Facts | Valid Inputs for Agents |
|-----------|---|--|---|
| ConceptKB | <pre>studied(occupation2, field) graduate(field2, occupation) isa(device, obj)</pre> | <pre>(Airpipe ; Isa ; haystone). (Working as kreuse ; HasPrerequisite ; Studying googolome). (Misigram ; Isa ; chikor). (Misigram device ; Isa ; pistarmen object). (Study metatoun ; MotivatedByGoal ; Work as kreuse).</pre> | <pre>What occupation do people who study scampot work in? What would be the field of study for someone working as a matularch? Which field have people working as zorgion graduated from? What devices are types of teeplemole? What is the device Pomorpha a type of? Which devices are of the type gastrat? What object is Pludgel a type of?</pre> |
| TextKB | <pre>dob(person, year) dod(person, year) occupation(person, occupation) field(person, field) invent(obj, year) usedo(obj, occupation2) usedf(obj, field2) founded(person, company) invented(person, tech) developed(company, device) manufactures(company, material) usedin(tech, device) contains(material, obj)</pre> | <pre>When studying kinneticket, saltcoat would be used. todou material is needed to make vetto. stretchwork is often used by people working as bartery. Carpoon device was developed based on the vout technology. Triclops studied chasmogon in college. flawpack was first invented in the year 1943. gambilla was invented in 2005. Kapod studied duriel in college. nosecutter is commonly used in the field of blaubrudin. Chaudelaire died in 1980. chickenshaw was invented in 1940. Dentology works as a scritigraphy. flawpack was first invented in the year 1989. Stoptite was born in 1937. chickenspaw material is needed to make stretchwork. Terbaryan was developed by the Coathanger company.</pre> | <pre>Which company produces the material topboard? Who have founded the company Moderexample? Monocytotyping is the founder of which companies? What is Loisy's occupation? When was cursaire invented? Which year was teeplemole invented? Which technologies has Kapod developed? Polyhoney is the inventor of which technologies? Which materials does Gutskin produce? What would be the occupation of someone using demiplane? What does Teinteen work as? What is Triclops's field of study? Which company produces the material enovasion? Who have developed the technology coule? herbalife is used by people in which field of study?</pre> |

Complex Questions (and Theory)

```
QC: What objects has Loisy likely used?
[select] <text> What is Loisy's field of study? A: ["cougarism", "nightslash"]
[project_flat_unique] <kb> What is the occupation of people who study #1? A: ["nephewskin", "skirtsicine"]
[project_flat_unique] <text> Which objects are used by a #2? A: ["cannolium", "microallocation", "tenderstiltskin", "monovacuum"]

QC: What objects has Triclops helped to make?
[select] <text> Triclops is the founder of which companies? A: ["Mechanicism"]
[project_flat_unique] <text> Which devices has #1 developed? A: ["Terbaryan"]
[project_flat_unique] <kb> What object is #2 a type of? A: ["vetto"]

QC: What objects has Stoptite helped to make?
[select] <text> Which technologies has Stoptite developed? A: ["thralline"]
[project_flat_unique] <text> #1 technology is used in which devices? A: ["Cabaretilonite"]
[project_flat_unique] <kb> What object is #2 a type of? A: ["cavata", "piperfish"]

QC: What objects has Kapod helped to make?
[select] <text> Which companies has Kapod founded? A: ["Superglitch"]
[project_flat_unique] <text> #1 produces which materials? A: ["fannyxist"]
[project_flat_unique] <text> Which objects use #2 as a material? A: ["epicanoine"]

QC: What objects has Minimiseris likely used?
[select] <text> What does Minimiseris work as? A: ["infilng", "glodome"]
[project_flat_unique] <kb> Which field have people working as #1 graduated from? A: ["kernwood", "kinneticket"]
[project_flat_unique] <text> What objects are used in the study of #2? A: ["pistarmen", "dactylin", "pilefork", "enableness"]

QC: What objects has Duriel likely used?
[select] <text> When did Duriel die? A: ["1928"]
[select] <text> Which invented objects are mentioned? A: ["legault", "stoptite", "stridery", "hydrallium", "...", "waxbox"]
[project] <text> Which year was #2 invented? A: [{"legault", ["1997"]}, {"stoptite", ["1991"]}, {"stridery", ["1921"]}, {"hydrallium", ["1993"]}, {"waxbox", ["1971"]}]]
[filterValues(#3)_keys] <math_special> Is #3 smaller than #1? A: ["stridery", "pistarmen"]
```

Figure 6: Example KB, space of valid inputs and theory used to construct COMMAQA-I

| | KB | Facts | Valid Inputs for Agents |
|---------|---|---|--|
| TableKB | nation(personj, nation) nation(persond, nation) | Athlete: Gigabut ; Nation: Besprit; Sport: Javelin. athlete: Fidelice ; country: Coathanger; sport: Javelin Throw. Athlete: Jimayo ; Nation: Tremolophore; Sport: Discus. athlete: Jungdowda ; country: Epicuratorion; sport: Discus Throw. | Which country does Metrix play for? Who are the discus throwers from Premericy? Which country is Entine from? Who are the discus throwers from Waxseer? Which country does Thym play for? Which country is Queness from? |
| TextKB | threwj(personj, length) threwd(persond, lengthd) | Mossia hurled the javelin to a distance of 87.2. Insimetry registered a throw of 85.6 in the javelin event. Undercabin registered a discus throw of 50.0. Diaquam registered a throw of 88.4 in the javelin event. Darecline registered a discus throw of 48.4. Vitule hurled the javelin to a distance of 66.4. Karmacogram threw the discus to a distance of 69.6. Sequinodactyl hurled the javelin to a distance of 70.2. | What were the lengths of the javelin throws by Predigime? Who was a discus thrower for 55.0? Who threw the discus for 67.6? Who threw the javelin for 67.2? Who was a javelin thrower for 93.0? Who was a discus thrower for 60.0? Who threw the javelin for 67.2? Who performed discus throws? |

Complex Questions (and Theory)

| |
|--|
| <p>QC: Who threw javelins longer than 89.6?</p> <pre>[select] <text> Who performed javelin throws? A: ["Jungdowda", "Prostigma", "Biopsie", "Thym", "Coacheship", "Knebbit", "Lowrise", "Sealt", "Seeper", "Entine", "Queness", "Cutthrough"] [project_zip] <text> What lengths were #1's javelin throws? A: [{"Jungdowda", ["71.2", "66.0", "73.6"]}, {"Prostigma", ["64.6"]}, {"Biopsie", ["77.6", "93.0"]}, {"Thym", ["87.0", "89.4", "86.8"]}, {"Coacheship", ["92.2", "72.2"]}, {"Knebbit", ["71.8", "84.0", "64.8", "75.8"]}, {"Lowrise", ["64.0", "82.8"]}, {"Sealt", ["68.6"]}, {"Seeper", ["65.6"]}, {"Entine", ["67.0"]}, {"Queness", ["91.2"]}, {"Cutthrough", ["80.8", "89.6", "79.4"]} [project_values] <math_special> max(#2) A: [{"Jungdowda", 73.6}, {"Prostigma", 64.6}, {"Biopsie", 93.0}, {"Thym", 89.4}, {"Coacheship", 92.2}, {"Knebbit", 84.0}, {"Lowrise", 82.8}, {"Sealt", 68.6}, {"Seeper", 65.6}, {"Entine", 67.0}, {"Queness", 91.2}, {"Cutthrough", 89.6}] [filter_keys(#3)] <math_special> is_greater(#3 89.6) A: ["Biopsie", "Coacheship", "Queness"]</pre> |
| <p>QC: How many discus throws were shorter than 48.0?</p> <pre>[select] <text> Who threw discus? A: ["Zayage", "Endography", "Dewbar", "Skullard", "Cabaretilonite", "Terbaryan", "Siligar", "Triclops", "Polyparity", "Cheapnose", "Flumph"] [project_flat] <text> What lengths were #1's discus throws? A: ["72.4", "54.4", "55.8", "66.8", "46.0", "70.8", "50.0", "59.4", "51.6", "70.0", "48.0", "45.0", "72.2", "66.2", "58.0", "65.6", "48.4", "61.8", "66.6", "44.0", "56.4", "50.2", "68.2", "47.2"] [filter(#2)] <math_special> is_smaller(#2 48.0) A: ["46.0", "45.0", "44.0", "47.2"] [select] <math_special> count(#3) A: 4</pre> |
| <p>QC: Who threw discuses shorter than 45.0?</p> <pre>[select] <text> Who threw discus? A: ["Dewbar", "Biscus", "Whime", "Dumasite", "Blumen", "Colorectomy", "Guazepam", "Metatoun", "Siligar", "Lechpin", "Sahaki", "Barbrauch", "Noosecutter", "Pompasole"] [project_zip] <text> What were the lengths of the discus throws by #1? A: [{"Dewbar", ["65.2", "44.0", "72.0"]}, {"Biscus", ["72.4", "73.6"]}, {"Whime", ["44.8", "65.0"]}, {"Dumasite", ["58.8"]}, {"Blumen", ["44.4", "54.6"]}, {"Colorectomy", ["53.6", "60.0"]}, {"Guazepam", ["52.8", "65.8"]}, {"Metatoun", ["46.8", "54.4", "51.4"]}, {"Siliga", ["59.4"]}, {"Lechpin", ["62.6"]}, {"Sahaki", ["48.6"]}, {"Barbrauch", ["45.0", "52.6"]}, {"Noosecutter", ["69.6"]}, {"Pompasole", ["64.0"]} [project_values] <math_special> min(#2) A: [{"Dewbar", 44.0}, {"Biscus", 72.4}, {"Whime", 44.8}, {"Dumasite", 58.8}, {"Blumen", 44.4}, {"Colorectomy", 53.6}, {"Guazepam", 52.8}, {"Metatoun", 46.8}, {"Siligar", 59.4}, {"Lechpin", 62.6}, {"Sahaki", 48.6}, {"Barbrauch", 45.0}, {"Noosecutter", 69.6}, {"Pompasole", 64.0}] [filter_keys(#3)] <math_special> is_smaller(#3 45.0) A: ["Dewbar", "Whime", "Blumen"]</pre> |
| <p>QC: What was the gap between the longest and shortest discus throws by Honeywax?</p> <pre>[select] <text> What lengths were Honeywax's discus throws? A: ["48.0", "59.8", "50.6"] [select] <math_special> max(#1) A: 59.8 [select] <math_special> min(#1) A: 48.0 [select] <math_special> diff(#2 #3) A: 11.8</pre> |
| <p>QC: What was the gap between the longest and shortest javelin throws by athletes from Misapportionment?</p> <pre>[select] <table> Who are the javelin throwers from Misapportionment? A: ["Zekkobe", "Featsaw", "Tantor"] [project_flat] <text> What lengths were #1's javelin throws? A: ["79.0", "67.8", "89.6", "80.4", "89.4", "79.6", "87.8"] [select] <math_special> max(#2) A: 89.6 [select] <math_special> min(#2) A: 67.8 [select] <math_special> diff(#3 #4) A: 21.8</pre> |
| <p>QC: What was the gap between the best javelin throws from Haystone and Pistarmen?</p> <pre>[select] <table> Which javelin throwers are from the country Haystone? A: ["Modiparity", "Polyacrylate", "Sequinodactyl"] [project_flat] <text> What lengths were #1's javelin throws? A: ["89.6", "75.2", "85.4", "67.8", "76.4", "68.4"] [select] <math_special> max(#2) A: 89.6 [select] <table> Who are the javelin throwers from Pistarmen? A: ["Crowdstrike"] [project_flat] <text> What were the lengths of the javelin throws by #4? A: ["66.0", "85.6"] [select] <math_special> max(#5) A: 85.6 [select] <math_special> diff(#3 #6) A: 4.0</pre> |

Figure 7: Example KB, space of valid inputs and theory used to construct COMMAQA-N