

Reasoning Like Program Executors

Anonymous ACL submission

Abstract

Reasoning over natural language is a long-standing goal for the research community. However, studies have shown that existing language models are inadequate in reasoning. To address the issue, we present POET, a new pre-training paradigm. Through pre-training language models with programs and their execution results, POET empowers language models to harvest the reasoning knowledge possessed in program executors via a data-driven approach. POET is conceptually simple and can be instantiated by different kinds of programs. In this paper, we show three empirically powerful instances, i.e., POET-Math, POET-Logic, and POET-SQL. Experimental results on six benchmarks demonstrate that POET can significantly boost model performance on natural language reasoning, such as numerical reasoning, logical reasoning, and multi-hop reasoning. Taking the DROP benchmark as a representative example, POET improves the F_1 metric of BART from 69.2% to 80.6%. Furthermore, POET shines in giant language models, pushing the F_1 metric of T5-11B to 87.6% and achieving a new state-of-the-art performance on DROP. POET opens a new gate on reasoning-enhancement pre-training and we will make our code, models, and data publicly available to facilitate future research.

1 Introduction

Recent breakthroughs in pre-training illustrate the power of pre-trained Language Models (LM) on a wide range of Natural Language (NL) tasks. Pre-training on self-supervised tasks, such as autoregressive language modeling (Brown et al., 2020) and masked language modeling (Devlin et al., 2019; He et al., 2021) using large amounts of NL sentences, boosts the language understanding of models by a large margin (Wang et al., 2018a). However, existing pre-training paradigms have primarily focused on language modeling and paid little

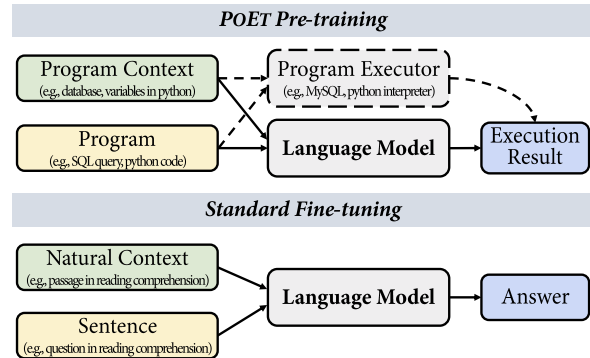


Figure 1: Given a program context and a program as input, POET pre-trains LMs to output the execution result. After fine-tuning on downstream tasks, POET can boost LMs on reasoning-required scenarios. Explanations about program context, program, program executor and execution result can be found in § 3. More examples of natural context and sentence are in Table 1.

attention to advanced *reasoning* capabilities (Table 1). As a result, though reaching near-human performance on several tasks, pre-trained LMs are still far behind expectation in reasoning-required scenarios, such as numerical reasoning (Wallace et al., 2019; Ravichander et al., 2019) and logical reasoning (Yu et al., 2020; Liu et al., 2020). This observed deficiency calls for the development of general-purpose pre-training approaches suitable for learning reasoning skills.

In light of this, we conceive a new pre-training paradigm, POET (**P**rogram **E**xecutor), to boost various reasoning skills over NL sentences by pre-training LMs with the task of *program execution*. As illustrated in Figure 1, with a *program* (e.g., SQL query) and its associated *program context* (e.g., database) as input, the model receives automatic supervision from an established program executor (e.g., MySQL) and learns to produce correct *execution result*. We believe that when LMs imitate program execution procedures, they could potentially learn the reasoning knowledge that humans adopted to create the associated program

Type	Example	Dataset	Task
Numerical	Question: What is the difference in casualty numbers between Bavarian and Austrian? Passage: [DOC] The popular uprising included large areas of . . .	DROP	Reading (RC) Comprehension
Logical	Conclusion: One employee supervises another who gets more salary than himself. Fact: [DOC] David, Jack and Mark are colleagues in a company. David supervises Jack, and Jack supervises Mark. David gets more . . .	LogiQA	Reading (RC) Comprehension
Multi-hop	Question: At which university does the biographer of John Clare teach English Literature? Passage: [DOC] John Clare : John Clare was an English poet . . . [DOC] CMS College Kottayam : The CMS College is one . . .	HotpotQA	Reading (RC) Comprehension
Hybrid	Question: What was the percentage change in gaming between 2018 and 2019? Context: [TAB] Server products and cloud services 32,622 26,129 . . . [DOC] Our commercial cloud revenue, which includes Office . . .	TAT-QA	Question Answering (QA)
Quantitative	Hypothesis: Teva earns \$7 billion a year. Premise: After the deal closes, Teva will generate sales of about \$7 billion a year, the company said.	EQUATE	Natural Language Inference (NLI)

Table 1: The demonstration of five representative reasoning types. Listed are the types, the example questions, the representative dataset and their corresponding tasks. [DOC] and [TAB] indicates the start of a passage and a semi-structured table respectively. Here we regard **Question**, **Conclusion** and **Hypothesis** as *sentence*, and **Passage**, **Fact**, **Context** and **Premise** as *natural context* in Figure 1.

065 executor, and transfer the reasoning capability to
066 NL sentences. This reveals the key hypothesis of
067 POET: *program executors are crystallized knowl-*
068 *edge of human reasoning, and such knowledge can*
069 *be transferred to natural language via pre-training.*

070 While it is extremely difficult to obtain large
071 amounts of clean natural language sentences con-
072 taining clear evidence of reasoning, thanks to the
073 artificial and compositional nature of programming
074 languages, synthesized programs can be made ar-
075 bitrarily complicated but readily available on any
076 scale. These merits greatly facilitate the construc-
077 tion of a high-quality pre-training corpus, address-
078 ing most of unresolved shortcomings in previ-
079 ous reasoning-enhancement pre-training. In other
080 words, POET differs from existing pre-training
081 paradigms relying on noisy NL data. In summary,
082 our contribution is three-fold:

- 083 • We propose POET, a new pre-training
084 paradigm for boosting reasoning capability
085 of language models by imitating program ex-
086 ecutors. Along with this paradigm, we present
087 three exemplary across-program POET instan-
088 tiations for various reasoning capabilities.
- 089 • We show with quantitative experiments that
090 the reasoning ability our models obtains from
091 POET pre-training is transferable to broader
092 natural language scenarios. On six reasoning-
093 focused downstream tasks, POET enables
094 general-purpose language models to achieve
095 comparable or even better performance than
096 previous state-of-the-art specialized models.
- 097 • We carry out comprehensive analytical stud-
098 ies on POET and summarize some insightful

099 findings in our pre-training. We hope these in-
100 sights would shed light on the future research
101 of reasoning like program executors.

2 Related Work 102

103 Since we focus on reasoning over natural language,
104 our work is closely related to previous works which
105 also concentrate on *reasoning skills* in NL tasks.
106 Regarding methods to inject reasoning skills into
107 LMs, our method is related to two lines of work
108 contributing to the topic: the line of *specialized*
109 *models* and the line of *pre-training*. Last, our work
110 is also related to *program execution* since we use
111 program executors in our pre-training.

112 **Reasoning Skills** The literature focuses on rea-
113 soning skills including numerical reasoning (Dua
114 et al., 2019), multi-hop reasoning (Yang et al.,
115 2018), reasoning in hybrid context (Chen et al.,
116 2020b; Zhu et al., 2021) and logical reasoning (Liu
117 et al., 2020; Yu et al., 2020). Our work concentrates
118 on improving the above reasoning skills, leaving
119 the other reasoning abilities such as commonsense
120 reasoning (Zellers et al., 2018; Talmor et al., 2019;
121 Bhagavatula et al., 2020) for future work.

122 **Reasoning via Specialized Models** Early works
123 typically design specialized models and augment
124 them into LMs for different types of questions (Dua
125 et al., 2019; Andor et al., 2019; Hu et al., 2019;
126 Ding et al., 2019). Taking Hu et al. (2019) as
127 an example, they first predicted the answer type
128 of a given question (e.g., “how many”), and then
129 adopted the corresponding module (e.g., count
130 module) to predict the answer. Although these

131 methods work well on a specific dataset, it is chal- 181
132 lenging for them to scale to complex reasoning 182
133 scenarios (Chen et al., 2020c). Differently, our 183
134 work follows the line of reasoning via pre-training, 184
135 which enjoys better scalability. 185
186

136 **Reasoning via Pre-training** This line of work 187
137 focuses on the continued pre-training of LMs using 188
138 large-scale data which involves reasoning. The pre- 189
139 training data are generally NL text, which are either 190
140 crawled from Web with distant supervision (Deng 191
141 et al., 2021), generated by a model-based gener- 192
142 ator (Asai and Hajishirzi, 2020), or synthesized 193
143 via human-designed templates (Geva et al., 2020; 194
144 Yoran et al., 2021; Campagna et al., 2020; Wang 195
145 et al., 2021). However, large-scale high-quality 196
146 textual data involving reasoning are difficult to col- 197
147 lect (Deng et al., 2021). Meanwhile, as the com- 198
148 plexity of desired reasoning operations increases, 199
149 synthesizing high-quality (e.g., fluent) NL sen- 200
150 tences becomes more challenging. Different from 201
151 the above pre-training methods relying on NL data, 202
152 our pre-training is performed on programs. These 203
153 programs can be synthesized at any scale with high- 204
154 quality and rich-diversity, and thus are much easier 205
155 to collect than NL sentences. 206

156 **Program Execution** We present a framework to 207
157 leverage program executors to train LMs, and thus 208
158 our work is close to recent works on learning a neu- 209
159 ral program executor. In this line, the most related 210
160 work to ours is Liu et al. (2021), which revealed 211
161 the possibility of SQL execution on helping table 212
162 pre-training. Different from them mainly focus- 213
163 ing on table-related tasks, we present a general- 214
164 ized approach to include Math, Logic, and SQL, as 215
165 well as their applications on many different natural 216
166 language downstream tasks. Other related stud- 217
167 ies include learning program executors on visual 218
168 question answering (Andreas et al., 2016), read- 219
169 ing comprehension (Gupta et al., 2019; Khot et al., 220
170 2020), knowledge base question answering (Ren 221
171 et al., 2021) and 3D rendering (Tian et al., 2019). 222
172 These works mainly focus on learning a neural 223
173 network to represent the program executor, while 224
174 ours focuses on transferring the knowledge of pro- 225
175 gram executor to downstream tasks via pre-training. 226
176 Other lines of research did not leverage models as 227
177 neural program executors, but instead leveraging 228
178 program execution in inference as a reliable sanity 229
179 guarantee for generated programs by pruning non-
180 executable candidates (Wang et al., 2018b; Chen

et al., 2019, 2021). Others have also noticed that
when a target program is sequential, execution of
the partially generated program provides reliable
guidance towards the final gold output (Odena et al.,
2020; Ellis et al., 2019; Chen et al., 2019; Sun et al.,
2018; Zohar and Wolf, 2018).

3 Reasoning Like Program Executors

Reasoning is the process where deduction and in-
duction are sensibly applied to draw conclusion
from premises or facts (Scriven, 1976). As a
supreme feature of intelligence, humans apply rea-
soning across modalities. Taking numerical rea-
soning as an example, humans can tell how many
chocolates are consumed from a math word prob-
lem description, or from a real-world event where
a mother gets off work and finds the choco-can
empty, aside standing their guilty-looking kids with
brownish stains on their faces. Through detach-
ment of information from their superficial modality
and symbolic abstraction, humans manage to unify
input formats and condense their numerical reason-
ing knowledge into one executable symbolic sys-
tem – This is the origin of an arithmetic program
executor. If a model can master these reasoning
skills by imitating program executors, we believe
in the possibility of transferring those reasoning
skills to different modalities. In our case, we ex-
pect language models to transfer reasoning to NL
related tasks. Given this motivation, we discuss
fundamental components of POET in the rest of
this section, and present three concrete instantia-
tions of our framework in § 4.

Program refers to a finite sequence of symbols
which can be understood and executed by machines.
For example, a program can be a logical form (e.g.,
Prolog), a piece of code (e.g., Python), or a math ex-
pression. Compared with NL sentences, programs
are more formal. Each well-established program
follows a specific set of syntax rules and can thus
be synthesized in a systematic way. The generaliz-
ability of POET framework is free from assumption
and derived from the set of syntax rules on which a
program complies. In POET, as long as a program
returns meaningful output to reflect its computa-
tional procedure, it is an acceptable program.

Program Context is the environment in which
a program is running, which holds numerous vari-
ables accessible to the program. These variables
serve as pivot points that anchor program context

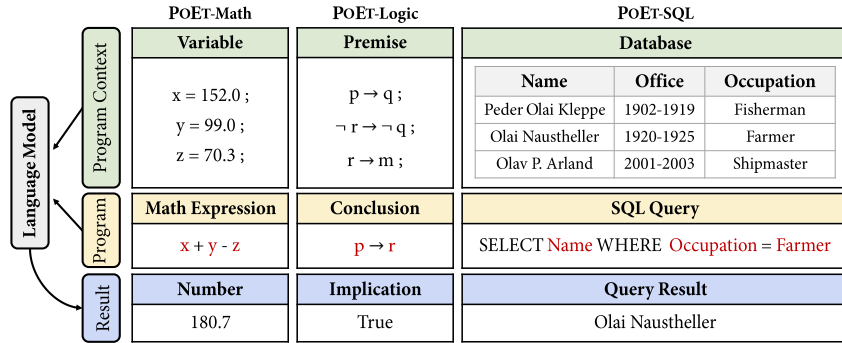


Figure 2: The illustration of three instantiations of POET to inject different kinds of reasoning skills, including POET-Math, POET-Logic and POET-SQL. The red text indicates the variables read by the program.

with the program. In the same sense, the question and the passage in reading comprehension hold a similar relationship. This suggests a natural analogy between the program to program context and the sentence to natural context in Figure 1.

Program Executor is a black-box software that can execute a given program within the program context. An example could be the Python interpreter that executes each line of code, with its specific input data structures as program context. For POET, program executors play the role of teachers to educate student (i.e., LMs) on reasoning knowledge they contain. POET expects program executors to deterministically execute an input program with respect to a specific program context.

Execution Result is obtained from the program executor, given a program and program context as input. It is much analogous to the answer part in NL downstream tasks. The execution result is the primary observable data reflecting the intermediate reasoning process, and serves as the supervision provided by the program executor.

4 Instantiations of POET

Along with the POET paradigm, we manifest three exemplary across-program POET instantiations (Figure 2), named POET-Math, POET-Logic and POET-SQL, for injecting numerical, logical and integrated reasoning capabilities into LMs.

4.1 POET-Math for Numerical Reasoning

The POET-Math (Left in Figure 2) aims at injecting numerical reasoning skills into LMs. Specifically, POET-Math is designed to boost the basic arithmetic skills (i.e., addition and subtraction) of LMs on downstream tasks. This arithmetic skill aligns with requirements to answer questions centered on

addition / subtraction between two numbers, such as “What is the difference in casualty numbers between Bavarian and Austrian?”.

Pre-training Task Given several floating-point variables as the program context and a math expression only involving addition/ subtraction as the program, the pre-training task of POET-Math is to *calculate the math expression*. Taking the leftmost example from Figure 2, receiving the concatenation of the program and the program context as the input, POET-Math is trained to output the number 180.7. Considering the output can be an arbitrary number, the encoder-decoder model (Lewis et al., 2020) is more suitable for this pre-training task.

Pre-training Corpus Each example in the corpus contains a math expression containing up to 2 operators and 3 variables, and a program context which contains at most 30 floating-point variables¹. The mathematical addition and subtraction operators are denoted by + and -, respectively. The values of variables vary from 0.0 to 1000.0. By random generation, we synthesize 4 million examples as the pre-training corpus for POET-Math.

4.2 POET-Logic for Logical Reasoning

The POET-Logic (Mid in Figure 2) aims at injecting logical reasoning (e.g., necessary conditional reasoning) skills into LMs. For example, taking the facts “Only if the government reinforces basic education can we improve our nation’s education to a new stage. In order to stand out among other nations, we need to have a strong educational enterprise.” as premises, POET-Logic is intended to help LMs identify whether the conclusion “In order to stand out among nations, we should reinforce basic education” is necessarily implied.

¹More discussion can be found in Appendix § C.

Pre-training Task Given a few first-order logic premise statements as the program context and one conclusion statement as the program, the pre-training task of POET-Logic is to identify *if the program is necessarily implied from the program context*. The execution result, i.e., the implication relationship between the program and the program context, is either `True` or `False`. Since the output is binary, an encoder-only model (Liu et al., 2019) is sufficient to perform this pre-training task.

Pre-training Corpus Each example in the corpus contains several premise statements and a conclusion statement. Initially, the statement collection for each example is empty. To produce it, we first allocate 5 Boolean variables (e.g., p and q in Figure 2) and randomly sample at most 8 pairs from their pairwise combinations. For each sampled pair (p, q) , we randomly select a statement from the set $\{p \rightarrow q, p \rightarrow \neg q, \neg p \rightarrow \neg q, \neg p \rightarrow q\}$ and add it to the collection. Once the statement collection is prepared, we randomly select a statement as the conclusion statement (i.e., program) and the rest as the premise statements (i.e., program context). Last, we employ Z3 (De Moura and Bjørner, 2008), the well-known satisfiability modulo theory solver, as our program executor to obtain the implied result. Finally, we synthesize 1 million examples as the pre-training corpus for POET-Logic, and nearly 16% examples correspond to `True`.

4.3 POET-SQL for Integrated Reasoning

POET-Math and POET-Logic each focus on one specific reasoning skill. Different from them, POET-SQL allows LMs to master different reasoning skills simultaneously via integrated reasoning.

Pre-training Task Given a database as the program context and a SQL query as the program, the pre-training task of POET-SQL is to mimic the *query result generation*. Since the encoder-decoder LMs can generate arbitrary tokens, they are well suited for the task. On the other hand, encoder-only models have insufficient expressiveness to produce out-of-context query results. To allow them to benefit from the SQL execution, we tailor the task into a *query result selection* task for encoder-only models, which only utilizes query results that can be found in the database. More specifically, the task requires encoder-only models to perform an IO sequence tagging process to find the query results in the database. Here the tag `I` is for golden tokens in the query results, while `O` is for other tokens.

Pre-training Corpus Each example in the corpus contains a SQL query, a database and a query result. Notably, following Liu et al. (2021), each database is flattened into a sequence when it is fed into LMs. Meanwhile, to avoid databases being too large to fit into memory, we randomly drop the rows of large databases until their flattened sequences contains less than 450 tokens. For the query result generation task, we follow the same corpus construction strategy as described in Liu et al. (2021). Concretely, by instantiating SQL templates from SQUALL (Shi et al., 2020) over databases provided by WIKISQL (Zhong et al., 2017), 5 million examples are synthesized for pre-training. For the query result selection task, the pre-training corpus is constructed in a similar way as above, except that only the examples whose query results are suitable for encoder-only are retained. This filtering results in a corpus containing nearly 2 million examples.

5 Experiments & Analysis

To verify the effectiveness of our POET framework on boosting the reasoning capabilities of LMs, we first apply our method on top of several backbone models, including encoder-only models and encoder-decoder models. Then we conduct experiments on six typical reasoning benchmark datasets and compare POET models with previous state-of-the-art (SOTA) methods. Last, we perform a detailed pre-training analysis to demonstrate key insights with respect to each part in our framework.

5.1 Backbone Models

RoBERTa (Liu et al., 2019), one of the most popular LMs, is elected as the backbone in encoder-only LMs. We mark the RoBERTa model trained under POET as POET-X_{RoBERTa}, where X is either Logic or SQL. BART (Lewis et al., 2020) is chosen as the backbone in encoder-decoder LMs. We mark the BART model trained under POET as POET-X_{BART}, where X is either Math or SQL. Meanwhile, to explore whether our approach is simultaneously effective for much larger LMs, we also apply our framework to T5-11B (Raffel et al., 2020), the largest publicly available language model.

5.2 Experimental Datasets

We perform experiments on different datasets including DROP (Dua et al., 2019), HotpotQA (Yang et al., 2018), TAT-QA (Zhu et al., 2021), EQUATE (Ravichander et al., 2019) and

(a) The experimental results of PoET-Math.

Models	DROP [♡] (EM)	DROP [♡] (F1)
BART-Large	66.2	69.2
PoET-Math _{BART}	75.2 (+9.0)	78.1 (+8.9)

(b) The experimental results of PoET-Logic.

Models	LogiQA (EM)
RoBERTa-Large	36.7
PoET-Logic _{RoBERTa}	38.9 (+2.2)

(c) The experimental results of PoET-SQL.

Models	DROP [♡]		HotpotQA [♡]		TAT-QA [♡]		SVAMP	EQUATE
	EM	F1	EM	F1	EM	F1	EM	EM
BART-Large	66.2	69.2	65.6	78.9	38.8	46.7	12.4	62.6
PoET-SQL _{BART}	77.7 (+11.5)	80.6 (+11.4)	66.5 (+0.9)	79.7 (+0.8)	41.5 (+2.7)	49.6 (+2.9)	33.5 (+21.1)	66.5 (+3.9)
RoBERTa-Large	78.1	85.3	67.6	81.1	55.2	62.7	–	64.2
PoET-SQL _{RoBERTa}	79.8 (+1.7)	87.4 (+2.1)	68.7 (+1.1)	81.6 (+0.5)	59.1 (+3.9)	65.9 (+3.2)	–	67.5 (+3.3)
T5-11B	83.5	85.9	71.4	84.5	–	–	52.9	–
PoET-SQL _{T5}	85.2 (+1.7)	87.6 (+1.7)	71.5* (+0.1)	84.4* (-0.1)	–	–	57.4 (+4.5)	–

Table 2: The main experimental results of different backbone models on test sets and dev sets (♡) of datasets² with or without our proposed PoET paradigm. The results of PoET are significantly better than the original LMs ($p < 0.05$), except for those marked by *. PoET-SQL/Math_{BART}, PoET-SQL/Logic_{RoBERTa} and PoET-SQL_{T5} are pre-trained from BART-Large, RoBERTa-Large and T5-11B respectively under the PoET paradigm. We verify the performance of PoET-SQL_{T5} on partial datasets considering our computation budget. Note the performance of RoBERTa-Large and PoET-SQL_{RoBERTa} are evaluated on the subset of DROP where the answer is span(s).

LogiQA (Liu et al., 2020). Table 1 shows examples of these datasets and highlights their corresponding reasoning types. More details can be found in Appendix § B. Furthermore, SVAMP (Patel et al., 2021), the challenging diagnostic dataset for probing *numerical reasoning*, is employed in our experiments to test the generalization capability of our fine-tuned models on DROP. Our models are evaluated on its addition and subtraction subsets.

5.3 Implementation Details

We implement our models based on transformers (Wolf et al., 2020), fairseq (Ott et al., 2019) and DeepSpeed³. Hyperparameters during pre-training and fine-tuning are provided in Appendix § E.

Passage Retrieval in HotpotQA Since the total length of the original passages in HotpotQA is too long to fit into memory, we train a classifier to filter out top-3 passages, as done in previous work (Deng et al., 2021). Specifically, a RoBERTa-Large model is fine-tuned to discriminate if an input passage is required to answer the question. The Hits@3 score of the classifier on HotpotQA is 97.2%.

Numerical Design in DROP and SVAMP As noticed by previous works, sub-word tokenization methods such as byte pair encoding (Sennrich et al., 2015) potentially undermines the arithmetic abil-

ity of models. Instead, the character-level number representation is argued to be a more effective alleviation (Wallace et al., 2019). Additionally, the reverse decoding of numbers is proposed as a better way of modelling arithmetic carry (Geva et al., 2020). Therefore, we employ these design strategies on DROP and SVAMP.

5.4 Methods Comparison

In this section, we compare our models with original LMs and previous state-of-the-art methods.

5.4.1 Comparing to Original LMs

Applying LMs to Different Datasets For any encoder-decoder LM (e.g., BART), we treat all datasets as generative tasks and fine-tune it directly to generate answers. As for the encoder-only LM (e.g., RoBERTa), the fine-tuning strategies on different datasets are slightly different. (i) On **DROP**, we cast the span selection task as a sequence tagging problem following Segal et al. (2020). (ii) On **TAT-QA**, we in-place substitute the RoBERTa-Large encoder in TAGOP (Zhu et al., 2021) with our PoET-SQL_{RoBERTa} to verify its effectiveness, and keep the rest of the components unchanged. (iii) On **HotpotQA**, we train two classifiers independently to predict the start and end positions of the answer span, as done in Devlin et al. (2019). (iv) On **EQUATE**, we train a classifier to perform sequence classification on concatenated premise-hypothesis pairs. Notably, we follow the official setup to train LMs on the MNLI dataset (Williams et al.,

² We compare our models with baselines on dev sets of partial datasets since their test sets are not publicly available.

³ <http://github.com/microsoft/DeepSpeed>

2018) and evaluate their zero-shot performance on EQUATE. (v) On **LogiQA**, we train a classifier to perform binary classification on concatenated question-option-context pairs, as suggested in Liu et al. (2020). (vi) On **SVAMP**, the encoder-only model is not suitable since the answers are out-of-context. On all datasets, our models are evaluated with official evaluation metrics EM and F1.

Experimental Results Table 2 presents a performance comparison between POET models and their vanilla versions without POET. Across all instances, we observe significant performance increment on downstream tasks requiring corresponding reasoning skills. Specifically, (a) POET-Math boosts numerical reasoning ability of BART, bringing in 9.0% EM gain on DROP; (b) POET-Logic improves logical reasoning skill of RoBERTa, resulting in a 2.2% EM improvement on LogiQA; (c) POET-SQL equips popular encoder-only and encoder-decoder models with an integrated package of reasoning skills, effectively improving their performance on five benchmark datasets. As a highlighted example, POET-SQL_{BART} obtains 11.5% (DROP) and 21.1% (SVAMP) improvements on EM, compared with the vanilla BART.

Since POET pre-training is carried purely on program context (Figure 2), whereas all downstream tasks are on natural context, our hypothesis that reasoning capability is transferable from program executors to NL scenarios gets verified. Another interesting observation is that POET also shines in giant LMs. As reflected from the results, T5-11B obtains noticeable performance gains on both DROP (1.7% EM) and SVAMP (4.5% EM).

5.4.2 Comparing to Previous SOTA

Baseline Setup We summarize the baseline methods in short below, and refer readers to their papers for more details. (i) On **DROP**, we include two families of models for comparison: specialized models such as NumNet(+) (Ran et al., 2019), MTMSN (Hu et al., 2019), NeRd (Chen et al., 2020c), QDGAT (Chen et al., 2020a) and language models such as GenBERT (Geva et al., 2020) and PReaM (Yoran et al., 2021). (ii) Similarly, on **HotpotQA** (Distractor), specialized model baselines include DFGN (Qiu et al., 2019), SAE (Tu et al., 2020), C2F Reader (Shao et al., 2020) and the SOTA model HGN (Fang et al., 2020). The language model baselines consist of BERT (Devlin et al., 2019), SpanBERT (Joshi et al., 2020)

Dataset	Models	EM	F1	
<i>Specialized Models</i>				
DROP [♡]	NumNet	64.9	68.3	
	MTMSN (BERT)	76.7	80.5	
	NeRd (BERT)	78.6	81.9	
	NumNet+ (RoBERTa)	81.1	84.4	
	QDGAT (RoBERTa)	<u>84.1</u>	<u>87.1</u>	
	<i>Language Models</i>			
	GenBERT (BERT)	68.8	72.3	
	PReasM (T5)	69.4	72.3	
	PoET-Math _{BART}	75.2	78.1	
	PoET-SQL _{BART}	77.7	80.6	
PoET-SQL+Math _{BART}	78.0	80.9		
PoET-SQL _{T5}	85.2	87.6		
<i>Specialized Models</i>				
HotpotQA [♡]	DFGN	55.7	69.3	
	SAE (BERT)	67.7	80.8	
	C2F Reader (RoBERTa)	68.0	81.2	
	HGN (RoBERTa)	<u>69.2</u>	<u>82.2</u>	
	<i>Language Models</i>			
	BERT	59.1	73.4	
	ReasonBERT (RoBERTa-Base)	64.8	79.2	
	PoET-SQL _{BART}	66.5	79.7	
	SpanBERT (BERT)	67.4	81.2	
	PoET-SQL _{RoBERTa}	68.7	81.6	
PoET-SQL _{T5}	71.5	84.4		
<i>Specialized Models</i>				
TAT-QA [♡]	TAPAS	18.9	26.5	
	NumNet+ V2	38.1	48.3	
	TAGOP (RoBERTa)	<u>55.2</u>	<u>62.7</u>	
	TAGOP (PoET-SQL _{RoBERTa})	59.1	65.9	
<i>Language Models</i>				
EQUATE	BERT	51.8	–	
	GPT	55.8	–	
	Q-REAS	60.7	–	
	PoET-SQL _{BART}	<u>66.5</u>	–	
	PoET-SQL _{RoBERTa}	67.5	–	
<i>Specialized Models</i>				
LogiQA	Co-Matching Network	33.9	–	
	PoET-LogiC _{RoBERTa}	<u>38.9</u>	–	
	DAGN (RoBERTa)	39.3	–	

Table 3: The comparison of our models with previous SOTA methods on test sets and dev sets (♡) of different datasets. LMs used by all baselines are in Large size, except for clarification. Bold and underlined numbers indicate the best and second-best results, respectively.

and ReasonBERT (Deng et al., 2021). (iii) On **TAT-QA**, we adopt the official baselines, including TAPAS (Herzig et al., 2020), NumNet+ V2 and the SOTA model TAGOP (Zhu et al., 2021). (iv) On **EQUATE**, we compare our methods with BERT (Devlin et al., 2019), GPT (Radford et al., 2019) and Q-REAS (Ravichander et al., 2019). (v) On **LogiQA**, we compare our methods with Co-Matching Network (Wang et al., 2018c) and the SOTA model DAGN (Huang et al., 2021).

Experimental Results Table 3 lists all experimental results of baselines and our models on different datasets. As seen, our model generally achieves the best or second-best results over different reasoning skills, showing its strong performance. Meanwhile, POET that utilizes a mix of two different programs (i.e., POET-SQL+Math_{BART}) achieves

Settings	PoET-SQL	PoET-Math
BART-Large	66.2/69.2	66.2/69.2
PoET Models	77.7/80.6	75.2/78.1
w.r.t. Reasoning	67.1/70.4	61.2/64.4
w.r.t. Program	76.9/79.7	–
w.r.t. Program Context	–	67.4/70.5
w.r.t. Program Executor	66.1/69.3	–
w.r.t. Execution Result	15.8/17.8	11.2/12.2

Table 4: The DROP EM/F₁ of PoET-SQL_{BART} and PoET-Math_{BART} with respect to each part in PoET.

a slightly better performance than SQL alone. Furthermore, compared with other reasoning-enhanced LMs, PoET-SQL_{BART} surpasses them by a large margin, demonstrating the effectiveness of our proposed program execution pre-training. For example, compared with PReasM initialized from T5-Large, PoET-SQL_{BART} initialized from BART-Large exceeds it by 8.3%. Finally, along with our proposed PoET framework, PoET-SQL_{T5} tops on the challenging benchmark DROP, revealing the great potential of LMs on reasoning scenarios.

5.5 Pre-training Analysis

In this section, we conduct pre-training analysis with respect to (w.r.t.) each part presented in § 3 to explore their key insights. We carry all feasible pre-training variants of PoET-SQL and PoET-Math, and then fine-tune them on DROP for performance comparison. All results are shown in Table 4.

w.r.t. Reasoning Although the reasoning knowledge in program executors has been proven to boost downstream tasks, we do not know under what conditions such knowledge would be helpful. To explore it, for PoET-SQL, we ablate all SQL queries containing numbers from its pre-training corpus, while for PoET-Math, we pre-train it to execute multiplication/division instead of addition/subtraction. The poor performance of PoET-SQL and PoET-Math variants indicate that it is important to maintain alignment between the reasoning skills involved in the pre-training tasks and the ones required for downstream tasks.

w.r.t. Program As stated before, PoET does not make assumption on syntax rules a program is built upon. To verify it, we randomly map all SQL reserved keywords to the 100 lowest frequency tokens in the BART vocabulary. Results suggest that even such “broken” syntax rules hardly reduce reasoning capability transferability, demonstrating the generality and adaptability of PoET.

w.r.t. Program Context In Figure 1, there is a natural analogy between the program to program context and the sentence to natural context, suggesting the necessity of program context on the reasoning transferability. To verify that, we employ the variant of PoET-Math where there is a variable-free program and an empty program context. Taking the example of PoET-Math in Figure 2, the program is transformed into $152.0 + 99.0 - 70.3$. One can see that there is a dramatic performance drop in the variant compared to PoET-Math_{BART}, verifying the importance of program context.

w.r.t. Program Executor The key hypothesis of PoET is that the program executor is crucial for our pre-training. To verify that, we ablate the program executor in PoET-SQL_{BART} and instead carry out a SQL language modeling pre-training. Practically, we mask each input SQL query in the pre-training corpus of PoET-SQL using the strategy adopted in BART (Lewis et al., 2020), and pre-train BART to output the associated complete SQL query given the masked SQL query and the database. The resulting scarce performance gain suggests what truly brings LMs reasoning ability is the program executor.

w.r.t. Execution Result Since the execution result serves as the supervision for LMs to learn from program executors, it must be a correct result. To verify that, we corrupt the correctness in variants of PoET-Math and PoET-SQL by randomly pairing the execution results of one example with the program and program context of another example. The extremely poor performance suggests that an incorrect pre-training corpus can cause significant damage to the reasoning ability of LMs.

6 Conclusion

We introduce PoET, a new pre-training paradigm for boosting reasoning capability of language models via imitating program executors. Experimental results on six datasets demonstrate that PoET can significantly boost existing language models on several reasoning skills, including numerical, logical and multi-hop reasoning. Our best language model under PoET can reach a comparable or better performance than state-of-the-art methods. Finally, we unveil key factors that make PoET successful. In the future, we hope our analysis could inspire more transference of reasoning knowledge from program executors to models.

References

- 609 Daniel Andor, Luheng He, Kenton Lee, and Emily
610 Pitler. 2019. [Giving BERT a calculator: Finding op-](#)
611 [erations and arguments with reading comprehension.](#)
612 In *Proceedings of the 2019 Conference on Empirical*
613 *Methods in Natural Language Processing and the*
614 *9th International Joint Conference on Natural Lan-*
615 *guage Processing (EMNLP-IJCNLP)*, pages 5947–
616 5952, Hong Kong, China. Association for Computa-
617 tional Linguistics.
618
- 619 Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and
620 Dan Klein. 2016. Neural module networks. In *Pro-*
621 *ceedings of the IEEE conference on computer vision*
622 *and pattern recognition*, pages 39–48.
- 623 Akari Asai and Hannaneh Hajishirzi. 2020. [Logic-](#)
624 [guided data augmentation and regularization for con-](#)
625 [sistent question answering.](#) In *Proceedings of the*
626 *58th Annual Meeting of the Association for Computa-*
627 *tional Linguistics*, pages 5642–5650, Online. As-
628 sociation for Computational Linguistics.
- 629 Chandra Bhagavatula, Ronan Le Bras, Chaitanya
630 Malaviya, Keisuke Sakaguchi, Ari Holtzman, Han-
631 nah Rashkin, Doug Downey, Wen tau Yih, and Yejin
632 Choi. 2020. [Abductive commonsense reasoning.](#) In
633 *International Conference on Learning Representa-*
634 *tions*.
- 635 Tom Brown, Benjamin Mann, Nick Ryder, Melanie
636 Subbiah, Jared D Kaplan, Prafulla Dhariwal,
637 Arvind Neelakantan, Pranav Shyam, Girish Sastry,
638 Amanda Askell, Sandhini Agarwal, Ariel Herbert-
639 Voss, Gretchen Krueger, Tom Henighan, Rewon
640 Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu,
641 Clemens Winter, Chris Hesse, Mark Chen, Eric
642 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
643 Jack Clark, Christopher Berner, Sam McCandlish,
644 Alec Radford, Ilya Sutskever, and Dario Amodei.
645 2020. [Language models are few-shot learners.](#) In
646 *Advances in Neural Information Processing Systems*,
647 volume 33, pages 1877–1901. Curran Associates,
648 Inc.
- 649 Giovanni Campagna, Agata Foryciarz, Mehrad Morad-
650 shahi, and Monica Lam. 2020. [Zero-shot transfer](#)
651 [learning with synthesized data for multi-domain dia-](#)
652 [logue state tracking.](#) In *Proceedings of the 58th An-*
653 *nuual Meeting of the Association for Computational*
654 *Linguistics*, pages 122–132, Online. Association for
655 Computational Linguistics.
- 656 Kunlong Chen, Weidi Xu, Xingyi Cheng, Zou Xi-
657 aochuan, Zhang Yuyu, Le Song, Taifeng Wang,
658 Yuan Qi, and Wei Chu. 2020a. [Question directed](#)
659 [graph attention network for numerical reasoning](#)
660 [over text.](#) pages 6759–6768.
- 661 Shuang Chen, Qian Liu, Zhiwei Yu, Chin-Yew Lin,
662 Jian-Guang Lou, and Feng Jiang. 2021. [ReTraCk:](#)
663 [A flexible and efficient framework for knowledge](#)
664 [base question answering.](#) In *Proceedings of the*
665 *59th Annual Meeting of the Association for Compu-*
666 *tational Linguistics and the 11th International Joint*
Conference on Natural Language Processing: Sys-
tem Demonstrations, pages 325–336, Online. Asso-
ciation for Computational Linguistics.
- 670 Wenhu Chen, Hanwen Zha, Zhiyu Chen, Wenhan
671 Xiong, Hong Wang, and William Yang Wang. 2020b. [HybridQA: A dataset of multi-hop question answering over tabular and textual data.](#) In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1026–1036, Online. Association for Computational Linguistics.
- 677 Xinyun Chen, Chen Liang, Adams Wei Yu, Denny
678 Zhou, Dawn Song, and Quoc V. Le. 2020c. [Neural](#)
679 [symbolic reader: Scalable integration of distributed](#)
680 [and symbolic representations for reading compre-](#)
681 [hension.](#) In *International Conference on Learning*
682 *Representations*.
- 683 Xinyun Chen, Chang Liu, and Dawn Song. 2019. [Execution-guided neural program synthesis.](#) In *International Conference on Learning Representations*.
- 687 Pradeep Dasigi, Nelson F. Liu, Ana Marasović,
688 Noah A. Smith, and Matt Gardner. 2019. [Quoref:](#)
689 [A reading comprehension dataset with questions re-](#)
690 [quiring coreferential reasoning.](#) In *Proceedings of*
691 *the 2019 Conference on Empirical Methods in Nat-*
692 *ural Language Processing and the 9th International*
693 *Joint Conference on Natural Language Processing*
694 *(EMNLP-IJCNLP)*, pages 5925–5932, Hong Kong,
695 China. Association for Computational Linguistics.
- 696 Leonardo De Moura and Nikolaj Bjørner. 2008. Z3:
697 An efficient smt solver. In *International conference*
698 *on Tools and Algorithms for the Construction and*
699 *Analysis of Systems*, pages 337–340. Springer.
- 700 Xiang Deng, Yu Su, Alyssa Lees, You Wu, Cong Yu,
701 and Huan Sun. 2021. [ReasonBERT: Pre-trained to](#)
702 [reason with distant supervision.](#) In *Proceedings of*
703 *the 2021 Conference on Empirical Methods in Natu-*
704 *ral Language Processing*, pages 6112–6127, Online
705 and Punta Cana, Dominican Republic. Association
706 for Computational Linguistics.
- 707 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
708 Kristina Toutanova. 2019. [BERT: Pre-training of](#)
709 [deep bidirectional transformers for language under-](#)
710 [standing.](#) In *Proceedings of the 2019 Conference*
711 *of the North American Chapter of the Association*
712 *for Computational Linguistics: Human Language*
713 *Technologies, Volume 1 (Long and Short Papers)*,
714 pages 4171–4186, Minneapolis, Minnesota. Associ-
715 ation for Computational Linguistics.
- 716 Ming Ding, Chang Zhou, Qibin Chen, Hongxia Yang,
717 and Jie Tang. 2019. [Cognitive graph for multi-hop](#)
718 [reading comprehension at scale.](#) In *Proceedings of*
719 *the 57th Annual Meeting of the Association for Com-*
720 *putational Linguistics*, pages 2694–2703, Florence,
721 Italy. Association for Computational Linguistics.

722	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 <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.	779
723		780
724		781
725		782
726		
727		783
728		784
729		785
730		786
731		787
732	Kevin Ellis, Maxwell I. Nye, Yewen Pu, Felix Sosa, Josh Tenenbaum, and Armando Solar-Lezama. 2019. Write, execute, assess: Program synthesis with a REPL . In <i>Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada</i> , pages 9165–9174.	788
733		789
734		790
735		791
736		
737		792
738		793
739		794
740	Yuwei Fang, Siqi Sun, Zhe Gan, Rohit Pillai, Shuo-hang Wang, and Jingjing Liu. 2020. Hierarchical graph network for multi-hop question answering . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 8823–8838, Online. Association for Computational Linguistics.	795
741		796
742		797
743		798
744		799
745		800
746		
747	Mor Geva, Ankit Gupta, and Jonathan Berant. 2020. Injecting numerical reasoning skills into language models . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 946–958, Online. Association for Computational Linguistics.	801
748		802
749		803
750		804
751		805
752		806
753		807
754		808
755		
756	Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced bert with disentangled attention . In <i>International Conference on Learning Representations</i> .	809
757		810
758		811
759		812
760	Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Eisenschlos. 2020. TaPas: Weakly supervised table parsing via pre-training . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4320–4333, Online. Association for Computational Linguistics.	813
761		814
762		815
763		816
764		817
765		
766		818
767	Minghao Hu, Yuxing Peng, Zhen Huang, and Dongsheng Li. 2019. A multi-type multi-span network for reading comprehension that requires discrete reasoning . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 1596–1606, Hong Kong, China. Association for Computational Linguistics.	819
768		820
769		821
770		822
771		823
772		824
773		825
774		826
775		827
776		828
777		829
778		
		830
		831
		832
		833
		834

835			
836		<i>Language Technologies</i> , pages 2080–2094, Online. Association for Computational Linguistics.	
837	Lin Qiu, Yunxuan Xiao, Yanru Qu, Hao Zhou, Lei Li, Weinan Zhang, and Yong Yu. 2019. Dynamically fused graph network for multi-hop reasoning . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6140–6150, Florence, Italy. Association for Computational Linguistics.		
844	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.		
848	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(140):1–67.		
854	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.		
860	Qiu Ran, Yankai Lin, Peng Li, Jie Zhou, and Zhiyuan Liu. 2019. NumNet: Machine reading comprehension with numerical reasoning . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2474–2484, Hong Kong, China. Association for Computational Linguistics.		
868	Abhilasha Ravichander, Aakanksha Naik, Carolyn Rose, and Eduard Hovy. 2019. EQUATE: A benchmark evaluation framework for quantitative reasoning in natural language inference . In <i>Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)</i> , pages 349–361, Hong Kong, China. Association for Computational Linguistics.		
876	Hongyu Ren, Hanjun Dai, Bo Dai, Xinyun Chen, Michihiro Yasunaga, Haitian Sun, Dale Schuurmans, Jure Leskovec, and Denny Zhou. 2021. Lego: Latent execution-guided reasoning for multi-hop question answering on knowledge graphs . In <i>ICML</i> .		
881	Michael Scriven. 1976. <i>Reasoning</i> . New York: McGraw-Hill.		
883	Elad Segal, Avia Efrat, Mor Shoham, Amir Globerson, and Jonathan Berant. 2020. A simple and effective model for answering multi-span questions . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 3074–3080, Online. Association for Computational Linguistics.		
	Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units . <i>CoRR</i> , abs/1508.07909.		890 891 892
	Nan Shao, Yiming Cui, Ting Liu, Shijin Wang, and Guoping Hu. 2020. Is Graph Structure Necessary for Multi-hop Question Answering? In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7187–7192, Online. Association for Computational Linguistics.		893 894 895 896 897 898 899
	Tianze Shi, Chen Zhao, Jordan Boyd-Graber, Hal Daumé III, and Lillian Lee. 2020. On the potential of lexico-logical alignments for semantic parsing to SQL queries . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 1849–1864, Online. Association for Computational Linguistics.		900 901 902 903 904 905 906
	Shao-Hua Sun, Hyeonwoo Noh, Sriram Somasundaram, and Joseph Lim. 2018. Neural program synthesis from diverse demonstration videos . In <i>International Conference on Machine Learning</i> , pages 4790–4799. PMLR.		907 908 909 910 911
	Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.		912 913 914 915 916 917 918 919 920
	Yonglong Tian, Andrew Luo, Xingyuan Sun, Kevin Ellis, William T. Freeman, Joshua B. Tenenbaum, and Jiajun Wu. 2019. Learning to infer and execute 3d shape programs .		921 922 923 924
	Ming Tu, Kevin Huang, Guangtao Wang, Jing Huang, Xiaodong He, and Bowen Zhou. 2020. Select, answer and explain: Interpretable multi-hop reading comprehension over multiple documents . In <i>The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020</i> , pages 9073–9080. AAAI Press.		925 926 927 928 929 930 931 932 933 934 935
	Eric Wallace, Yizhong Wang, Sujian Li, Sameer Singh, and Matt Gardner. 2019. Do NLP models know numbers? probing numeracy in embeddings . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 5307–5315, Hong Kong, China. Association for Computational Linguistics.		936 937 938 939 940 941 942 943 944
	Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018a.		945 946

947	GLUE: A multi-task benchmark and analysis platform for natural language understanding . In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 353–355, Brussels, Belgium. Association for Computational Linguistics.	1005
948		1006
949		
950		
951		
952		
953	Chenglong Wang, Kedar Tatwawadi, Marc Brockschmidt, Po-Sen Huang, Yi Xin Mao, Oleksandr Polozov, and Rishabh Singh. 2018b. Robust text-to-sql generation with execution-guided decoding. <i>ArXiv</i> , abs/1807.03100.	
954		
955		
956		
957		
958	Shuohang Wang, Mo Yu, Jing Jiang, and Shiyu Chang. 2018c. A co-matching model for multi-choice reading comprehension . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 746–751, Melbourne, Australia. Association for Computational Linguistics.	
959		
960		
961		
962		
963		
964		
965	Siyuan Wang, Wanjun Zhong, Duyu Tang, Zhongyu Wei, Zhihao Fan, Daxin Jiang, Ming Zhou, and Nan Duan. 2021. Logic-driven context extension and data augmentation for logical reasoning of text. <i>arXiv preprint arXiv:2105.03659</i> .	
966		
967		
968		
969		
970	Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122, New Orleans, Louisiana. Association for Computational Linguistics.	
971		
972		
973		
974		
975		
976		
977		
978		
979	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	
980		
981		
982		
983		
984		
985		
986		
987		
988		
989		
990		
991	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.	
992		
993		
994		
995		
996		
997		
998		
999	Ori Yoran, Alon Talmor, and Jonathan Berant. 2021. Turning tables: Generating examples from semi-structured tables for endowing language models with reasoning skills . <i>CoRR</i> , abs/2107.07261.	
1000		
1001		
1002		
1003	Weihaoyu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. Reclor: A reading comprehension dataset requiring logical reasoning. In <i>International Conference on Learning Representations (ICLR)</i> .	1005
1004		1006
	Rowan Zellers, Yonatan Bisk, Roy Schwartz, and Yejin Choi. 2018. SWAG: A large-scale adversarial dataset for grounded commonsense inference . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 93–104, Brussels, Belgium. Association for Computational Linguistics.	1007
		1008
		1009
		1010
		1011
		1012
		1013
	Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2SQL: Generating structured queries from natural language using reinforcement learning. <i>arXiv</i> , abs/1709.00103.	1014
		1015
		1016
		1017
	Fengbin Zhu, Wenqiang Lei, Youcheng Huang, Chao Wang, Shuo Zhang, Jiancheng Lv, Fuli Feng, and Tat-Seng Chua. 2021. TAT-QA: A question answering benchmark on a hybrid of tabular and textual content in finance . <i>CoRR</i> , abs/2105.07624.	1018
		1019
		1020
		1021
		1022
	Amit Zohar and Lior Wolf. 2018. Automatic program synthesis of long programs with a learned garbage collector . <i>CoRR</i> , abs/1809.04682.	1023
		1024
		1025
	A PoET-SQL for Integrated Reasoning	1026
	Table 5 presents seven typical SQL types and their representative SQL programs. We believe that the main reason SQL queries involve integrated reasoning is that they are complex enough to encompass a wide variety of computational procedures. For example, the arithmetic type covers part of the numerical reasoning capability, while the nested type roughly simulates the multi-hop procedure by recursively querying information on the database.	1027
		1028
		1029
		1030
		1031
		1032
		1033
		1034
		1035
	B Dataset Details	1036
	Table 6 presents some statistics about our experimental datasets. Below we introduce each dataset in detail.	1037
		1038
		1039
	DROP A reading comprehension benchmark to measure <i>numerical reasoning</i> ability over a given passage (Dua et al., 2019). It contains three subsets of questions: <i>span</i> , <i>number</i> , and <i>date</i> , each of which involves a lot of numerical operations. Unlike traditional reading comprehension datasets such as SQuAD (Rajpurkar et al., 2016) where answers are always a single span from context, several answers in the <i>span</i> subset of DROP contains multiple spans. The <i>number</i> and <i>date</i> answers are mostly out of context and need generative-level expressiveness.	1040
		1041
		1042
		1043
		1044
		1045
		1046
		1047
		1048
		1049
		1050
		1051

Type	Example SQL Program
Arithmetic	SELECT [COL] ₁ - [COL] ₂
Superlative	SELECT MAX([COL] ₁)
Comparative	SELECT [COL] ₁ WHERE [COL] ₂ > [VAL] ₂
Aggregation	SELECT COUNT([COL] ₁)
Intersection	SELECT [COL] ₁ WHERE [COL] ₂ = [VAL] ₂ AND [COL] ₃ = [VAL] ₃
Union	SELECT [COL] ₁ WHERE [COL] ₂ = [VAL] ₂ OR [COL] ₃ = [VAL] ₃
Nested	SELECT [COL] ₁ WHERE [COL] ₂ IN (SELECT [COL] ₂ WHERE [COL] ₃ = [VAL] ₃)

Table 5: The seven typical SQL types corresponding to numerical reasoning (**Top**) and multi-hop reasoning (**Bottom**). Listed are the type and the example SQL programs. [COL] and [VAL] represent the table column and the table cell value, respectively.

Dataset	Train		Dev	
	# Questions	# Docs	# Questions	# Docs
DROP	77,409	5,565	9,536	582
HotpotQA	90,564	90,564	7,405	7,405
TAT-QA	13,215	2,201	1,668	278
SVAMP	–	–	726	726
EQUATE	–	–	9,606	9,606
LogiQA	6,942	6,942	868	868

Table 6: The statistics of our experimental datasets.

HotpotQA An extractive reading comprehension dataset that requires models to perform *multi-hop reasoning* over different passages (Yang et al., 2018). It contains two settings (i) *Distractor*: reasoning over 2 gold paragraphs along with 8 similar distractor paragraphs and (ii) *Full wiki*: reasoning over customized retrieval results from full Wikipedia passages. We experiment with its distractor setting since retrieval strategy is beyond our focus in this work.

TAT-QA A question answering benchmark to measure reasoning ability over *hybrid* context, i.e., passages and tables (Zhu et al., 2021). It is curated by combing paragraphs and tables from real-world financial reports. According to the source(s) the answers are derived from, the dataset can be divided into three subsets: *Table*, *Text* and *Table-Text(both)*.

EQUATE The first benchmark dataset to explore *quantitative reasoning* under the task of natural language inference (Ravichander et al., 2019). As a test-only dataset, it requires fine-tuned models on MNLI to perform *zero-shot* natural language inference tasks over quantitative statements described in (premise, hypothesis) pairs to reach final entailment decisions.

Models	EM	F1
BART-Large	66.2	69.2
POET-Math _{BART} with 0 irrelevant variable	71.5	74.5
POET-Math _{BART} with 10 irrelevant variables	74.6	77.5
POET-Math _{BART} with 30 irrelevant variables	75.2	78.1

Table 7: The DROP performance with different numbers of irrelevant variables in POET-Math_{BART} pre-training.

LogiQA A multi-choice reading comprehension dataset that evaluates the *logical reasoning* ability, whose questions are designed by domain experts (Liu et al., 2020). It contains four types of logical reasoning, including categorical reasoning, disjunctive reasoning, conjunctive reasoning and conditional reasoning.

SVAMP A challenging math word problem dataset (Patel et al., 2021). It is designed specifically to hack models who leverage spurious patterns to perform arithmetic operations without true understanding of context. We only keep addition and subtraction problems in accordance with our pre-training coverage.

C Variables Design in POET-Math

In the pre-training task of POET-Math, we regard several floating-point variables as the program context. These variables include necessary variables (i.e., variables required by the program) and irrelevant variables. The irrelevant variables exist to make the program context closer to the natural context which generally contains irrelevant sentences. For example, given the program $a + b$ and the program context $a = 1; b = 2; c = 3; d = 4;$, variables c and d are what we refer to as irrelevant variables. This is motivated by the fact that passages are usually full of irrelevant information regarding a specific question in NL downstream tasks. In this section, we explore impacts on pre-training effectiveness brought by numbers of irrelevant variables. Empirically, we experiment on pre-training with 0, 10, 30 irrelevant variables. The total length of 30 irrelevant variables approaches the maximum input length of pre-trained LMs, and thus we do not try more settings.

The experimental results are shown in Table 7. As observed, (i) models can still learn numerical reasoning during pre-training where the program context is free from irrelevant variables, though less effective. (ii) the setting of 30 irrelevant variables

Dataset	Train		Dev	
	# Questions	# Docs	# Questions	# Docs
SQuAD v1.0	77,409	5,565	9,536	582
MNLI	392,702	392,702	9,815	9,815
QuoRef	19,399	3,771	2,418	454

Table 8: POET on NL understanding experiment dataset statistics.

brings BART-Large more performance improvement than the setting of 10 irrelevant variables. Considering there are plenty of lengthy passages in the DROP dataset, we therefore hypothesize that the noise level brought by irrelevant variables in the program context during pre-training should be made closer with the counterpart in the natural context during fine-tuning.

D NL Understanding Performance

Since the program context used in pre-training differs much from the natural context used in downstream tasks, a reasonable concern immediately follows: whether POET pre-training improves reasoning ability at the sacrifice of natural language understanding (NLU) ability of LMs? To investigate the concern, we evaluate POET models on representative benchmarks without emphasis on advanced reasoning skills, also covering the task of reading comprehension (RC) and natural language inference (NLI).

Dataset We fine-tune POET-SQL_{RoBERTa} on (i) SQuAD v1.0: (Rajpurkar et al., 2016): one of the most classical single-span selection RC benchmarks measuring understanding over natural language context; (ii) MNLI (Williams et al., 2018): a large-scale NLI dataset measuring cross-domain and cross-genre generalization of NLU. Notably, our model is evaluated on the *matched* setting for the purpose of simplicity. (iii) QuoRef (Dasigi et al., 2019): A Wikipedia-based multi-span selection RC benchmark with a special emphasis on coreference resolution. All dataset Statistics are shown in Table 8.

Implementation Details (i) On SQuAD, we cast the span selection task as a sequence tagging problem following Segal et al. (2020). (ii) On MNLI-matched, we train both models to perform sequence classification on concatenated premise-hypothesis pairs. (iii) On **Quoref**, we cast the span(s) selection task as an IO sequence tagging problem following Segal et al. (2020).

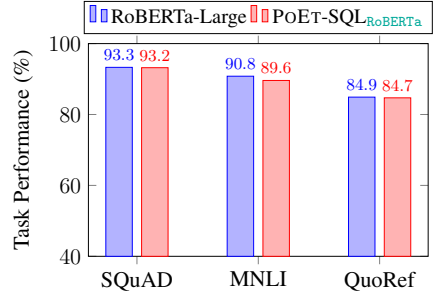


Figure 3: The performance comparison between RoBERTa-Large and POET-SQL_{RoBERTa} on representative NLU tasks. On SQuAD and QuoRef, we compare F_1 , whereas on MNLI we compare Accuracy.

Results As can be observed from performance comparison between POET-SQL_{RoBERTa} and vanilla RoBERTa shown in Figure 3, across all three experimented NLU-focused datasets, POET-SQL_{RoBERTa} performance are almost identical from counterparts of vanilla version. These negligible drops of performance suggest that reasoning capability can be transferred from program execution pre-training to NL downstream tasks, without the expense of LMs’ intrinsic understanding of language.

E Implementation Details

E.1 Pre-training Details

By default, we apply AdamW as pre-training optimizer with default scheduling parameters in fairseq. The coefficient of weight decay is set as 0.05 to alleviate over-fitting of pre-trained models. Additionally, we employ fp16 to accelerate the pre-training.

POET-Math The pre-training procedure lasts for 10,000 steps with a batch size of 512. After the warm up in the first 2000 steps, the learning rate arrives the peak at 3×10^{-5} during pre-training.

POET-Logic The pre-training procedure lasts for 5,000 steps with a batch size of 512. After the warm up in the first 1000 steps, the learning rate arrives the peak at 3×10^{-5} during pre-training.

POET-SQL For POET-SQL_{BART} and POET-SQL_{RoBERTa}, the pre-training procedure lasts for 50,000 steps with a batch size of 512. After the warm up in the first 5000 steps, the learning rate arrives the peak at 3×10^{-5} during pre-training. To save memory, each example in the pre-training corpus could at most contains 512 tokens. For POET-SQL_{T5}, the pre-training procedure lasts for 20,000 steps with a batch size of 512. After the warm

Models	Number	Span	Spans	Date	Total
<i>Previous Systems</i>					
MTMSN (BERT)	81.1	82.8	62.8	69.0	80.5
NumNet+ (RoBERTa)	83.1	86.8*	<u>86.8*</u>	63.9	84.4
QDGAT (RoBERTa)	86.2	88.5*	88.5*	67.5	<u>87.1</u>
GenBERT	75.2	74.5	24.2	56.4	72.3
PREasM	64.4	86.6	78.4	77.7	72.3
<i>Original LMs</i>					
RoBERTa-Large	–	86.4	79.9	–	–
BART-Large	63.6	79.6	74.6	62.1	69.2
T5-11B	83.2	<u>90.2</u>	85.8	84.9	85.8
<i>POET Models</i>					
POET-SQL _{RoBERTa}	–	88.2	83.1	–	–
POET-SQL _{BART}	78.9	84.5	79.6	71.9	80.6
POET-SQL _{T5}	<u>85.2</u>	92.4	86.6	<u>84.4</u>	87.6

Table 9: Breakdown of model F_1 score by answer types on the dev set of DROP. Some works only report overall span type performance (marked by *), and single-span is non-separable from multi-span performance. Bold and underlined numbers indicate the best and second-best results, respectively.

up in the first 2000 steps, the learning rate arrives the peak at 1×10^{-5} during pre-training. The maximum input length in each example is truncated to 384 tokens to increase the batch size.

E.2 Fine-tuning Details

By default, we apply AdamW as fine-tuning optimizer with default scheduling parameters on all datasets. To ensure statistical significance, all fine-tuning procedures are run with three random seeds, except for T5-11B and POET-SQL_{T5} due to the limit of computation budgets.

DROP POET-SQL_{RoBERTa} and RoBERTa-Large are trained with the subset of questions marked as “span” from the DROP dataset. Since a gold answer may occur multiple times in the passage, we optimize over the sum of negative log probability for all possibly-correct IO sequences where each one of gold answers is included at least once, as done in Segal et al. (2020). The fine-tuning procedure runs up to 25,000 steps with a batch size of 64, with the learning rate of 7.5×10^{-6} . As for BART-Large (and POET-SQL_{BART}, POET-Math_{BART}, the same below) and T5-11B (and POET-SQL_{T5}, the same below), they are trained with the whole DROP dataset. For BART-Large, the fine-tuning procedure runs up to 20,000 steps with a batch size as 128 and a learning rate as 3×10^{-5} . For T5-11B, due to the computational budget, the fine-tuning procedure only lasts for 10,000 steps with a batch size of 32, and the learning rate is 1×10^{-5} .

TAT-QA In the experiment of TAT-QA, we employ the official implementation and the default

hyperparameters provided in TAGOP⁴. The fine-tuning procedure runs up to 50 epochs with a batch size of 48. For modules introduced in TAGOP, the learning rate is set as 5×10^{-4} , while for RoBERTa-Large (and POET-SQL_{RoBERTa}), the learning rate is set as 1.5×10^{-5} .

HotpotQA The fine-tuning procedure runs up to 30,000 steps with a batch size of 64. The learning rate is 1×10^{-5} . Overlong inputs are truncated to 512 tokens for both RoBERTa-Large (and POET-SQL_{RoBERTa}), T5-11B (and POET-SQL_{T5}) and BART-Large (and POET-SQL_{BART}).

EQUATE The fine-tuning procedure runs up to 20,000 steps on MNLI with a batch size of 128 for both RoBERTa-Large (and POET-SQL_{RoBERTa}) and BART-Large (and POET-SQL_{BART}), with learning rate is 1×10^{-5} . After fine-tuning, models are directly evaluated on EQUATE.

LogiQA In the experiment of LogiQA, we employ the open-source implementation and the default hyperparameters provided in ReClor⁵ (Yu et al., 2020) to fine-tune RoBERTa-Large (and POET-SQL_{RoBERTa}). The fine-tuning procedure runs up to 10 epochs with a batch size of 24. The learning rate is set as 1×10^{-5} .

F Fine-grained Results

DROP In Table 9 we report model F_1 scores by question type on DROP. Comparing three POET pre-trained models with their vanilla versions, we observe that: (i) POET-SQL_{BART} outperforms the

⁴<https://github.com/NExTplusplus/TAT-QA>

⁵<https://github.com/yuweihao/reclor>

Models	RTE-Q	NewsNLI	RedditNLI	NR ST	AWPNLI	Average
<i>Previous Systems</i>						
MAJ	57.8	50.7	58.4	33.3	50.0	50.4
BERT	57.2	72.8	49.6	36.9	42.2	51.8
GPT	68.1	72.2	52.4	36.4	50.0	55.8
Q-REAS	56.6	61.1	50.8	63.3	71.5	60.7
<i>Original LMs</i>						
BART-Large	68.1	76.2	65.0	53.7	49.7	62.6
RoBERTa-Large	69.3	<u>75.5</u>	<u>65.6</u>	60.1	<u>50.7</u>	64.2
<i>POET Models</i>						
POET-SQL _{BART}	<u>72.3</u>	75.2	64.8	70.7	49.5	<u>66.5</u>
POET-SQL _{RoBERTa}	75.3	<u>75.5</u>	68.1	<u>69.2</u>	50.5	67.5

Table 10: The EM performance of different models on all subsets of the EQUATE benchmark. Bold and underlined numbers indicate the best and second-best results, respectively.

	Table	Text	Table-Text	Total
	EM / F ₁	EM / F ₁	EM / F ₁	EM / F ₁
Arithmetic	50.1 / 50.1	43.8 / 50.0	55.6 / 55.6	51.5 / 51.5
Counting	66.7 / 66.7	- / -	90.0 / 90.0	81.3 / 81.3
Spans	67.4 / 80.6	54.2 / 80.8	79.2 / 84.8	71.4 / 82.6
Span	68.4 / 68.4	51.2 / 76.0	76.2 / 77.8	61.9 / 74.6
Total	56.5 / 58.0	51.1 / 75.0	69.0 / 70.7	59.1 / 65.9

Table 11: The EM performance of TAGOP (POET-SQL_{RoBERTa}) with respect to answer types and sources on the dev set of TAT-QA.

vanilla BART-large with a wide margin in all types of questions, i.e. *number* (15.3%), *date* (9.8%), *span* (around 5%). (ii) POET-SQL_{RoBERTa} only deals with span selection questions, and obtain 1.9%, 3.2% gain on *span*, *spans* questions, respectively. (iii) For the giant POET-SQL_{T5}, we also observe 2% improvement on *number* questions, 2.2% on *span* and 0.8% on *spans* questions. These model-agnostic performance boost on DROP reveals the extra numerical reasoning knowledge models learned from SQL program executors.

vanilla models, suggesting that POET does not harm intrinsic abilities of language models.

TAT-QA Table 11 shows the detailed experimental results of TAGOP (POET-SQL_{RoBERTa}). Considering that the pre-training of POET-SQL_{RoBERTa} is only performed on table-like texts (i.e., the flatten sequence of databases), it is highly non-trivial for our model to generalize to such a hybrid scenario containing both tables and passages, again illustrating the transferability of reasoning capabilities.

EQUATE Table 10 presents performance breakdown by subsets of EQUATE (Ravichander et al., 2019), where we compare POET-SQL_{BART} and POET-SQL_{RoBERTa} with their vanilla versions and previous baselines. For both models, we observe around 10% acc improvement on the *NR ST* subset, where **numerical comparison and quantifiers** are especially emphasized. Stable performance improvement was also observed in both pre-trained models on the *RTE-Q* subset, where **arithmetics and ranges** are primary focus. Interestingly, POET-SQL_{RoBERTa} alone demonstrate improvement on *RedditNLI* (emphasizes approximation and verbal quantitative reasoning) subset. Performance on other subsets are approximately comparable between POET pre-trained models and