

Interpretable Proof Generation via Iterative Backward Reasoning

Anonymous ACL submission

Abstract

We present IBR, an Iterative Backward Reasoning to solve the proof generation task on rule-based Question Answering (QA), where models are required to reason over a series of textual rules and facts to find out the related proof path and derive the final answer. We handle the limitations of existed works in two folds: 1) enhance the interpretability of reasoning procedure with detailed tracking, by predicting nodes and edges in the proof path iteratively backward from the question; 2) promote the efficiency and accuracy via reasoning on the elaborate representations of nodes and history path, without any intermediate texts that may introduce external noise during proof generation. There are three main modules in IBR, QA and proof strategy prediction to obtain the answer and offer guidance for the following procedure; parent node prediction to determine a node in the existing proof that a child node will link to; child node prediction to find out which new node will be added to the proof. Experiments on both synthetic and paraphrased datasets demonstrate that IBR has a better in-domain performance as well as cross-domain transferability than state-of-the-art models.

1 Introduction

Endowing machines with reasoning capabilities is a longstanding problem (Newell and Simon, 1956a) in the field of AI. Though existing tasks such as multi-hop QA (Yang et al., 2018; Welbl et al., 2018) or logical-reasoning QA (Yu et al., 2020; Dua et al., 2019) impose a higher requirement on the reasoning capabilities, they usually just request for an answer without reasoning the procedure that make it interpretable. Recently, Clark et al. (2020) proposed new datasets and tasks for interpretable reasoning. Given a question, along with a set of facts (plain statements) and rules (implication relationships) that are expressed in natural language, there are two tasks: 1) predicting the binary an-

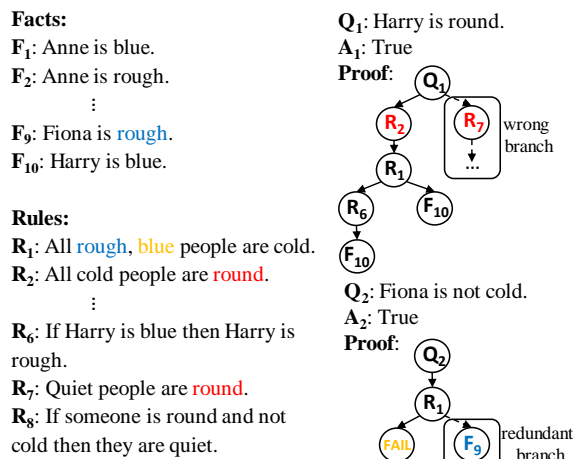


Figure 1: Illustration of generating proof iteratively. Regarding the proof path as a graph, and using the question as the initial node, other nodes and edges will be added step by step. (The gold proof is the obtained path in a reverse order exclude the question). The main challenges are wrong (cannot derive the answer) or redundant (can derive the answer, but the path is longer than the optimal one) branches may be involved at each step.

swer; 2) generating the proof path behind this answer. Large-scale pretrained models have shown strong performance on the first subtask in the early work (Liu et al., 2019), but there still remain challenges for the second one.

Several approaches have been proposed to simultaneously address the two subtasks. PROVER (Saha et al., 2020) and PROBR (Sun et al., 2021) tried to construct the reasoning path at once, where two classifiers are used to determine whether each node or edge is involved in the proof path respectively based on corresponding encoded representations. But they lack interpretability on tracking the detailed reason for selecting each step. To make proof generation more interpretable, Proofwriter (Tafjord et al., 2021) and EVR (Liang et al., 2021) decompose complex reasoning over the question into multiple simple procedures, resulting in iterative and intermediate processes with the help of intermediate texts in backward

062 order. Nevertheless, both of them suffer from
063 efficiency and external errors issues. The reason
064 is that they both require a large searching space,
065 as they perform on the whole inferable texts
066 and ignore the structure information from the
067 history path have already derived. Moreover, the
068 generation of intermediate text is costly and may
069 introduce extra noise propagation.

070 Inspired by the top-down AMR parsing (Cai
071 and Lam, 2019), we present **Iterative Backward**
072 **Reasoning (IBR)** for better proof generation. It gener-
073 erates proof path iteratively starting from the core
074 component for QA, i.e. question, making the pro-
075 cess more interpretable with trackable intermediate
076 states. Regarding a higher efficiency and accuracy,
077 along with two challenges mentioned in Figure 1,
078 the proof generation module of IBR simplifies the
079 intermediate process of reasoning as well as avoids
080 the unnecessary search for a possible unsuitable
081 branch. To add a new node and edge to the path,
082 there are two steps in IBR for each iteration: 1)
083 finding out the next parent node, i.e. one existing
084 rule or fact in the parsed history path that a new
085 node will become its child; 2) determine which
086 rule or fact that will be the new child node and
087 added to the path. Equipped with question-aware
088 representations from a pretrained encoder, along
089 with structure-aware node and path features, our
090 model can choose the optimal endpoint. It accom-
091 plishes reasoning with the highest possibility to
092 obtain a correct subsequent proof path based on
093 relevant features, getting rid of intermediate texts
094 while utilizing history path to avoid redundancy on
095 all possible texts as previous iterative works.

096 In addition, to make IBR applicable for samples
097 with incomplete proof path, which are abandoned
098 in former models like EVR (Liang et al., 2021), we
099 employ a proof strategy predictor to output proof
100 types. The predictions will be integrated into the
101 later proof generation actions, making the process
102 more controllable under different conditions.

103 We validate our approach on several datasets that
104 are widely used in previous studies (i.e. DU0-DU5,
105 Birds-Electricity, and ParaRules) spanning differ-
106 ent supervision settings (i.e. fully-supervised, few-
107 shot, and zero-shot). Experimental results show
108 that, compared to existing strong baselines, IBR
109 can achieve the best overall performance of proof
110 generation and comparable answer prediction accu-
111 racy, along with noticeable generalization capabil-
112 ity. Extensive analyses show that 1) the improve-

113 ments come from our elaborately designed iterative
114 and simplified proof generation modules, and 2)
115 the latency could be significantly improved, con-
116 firming our claims.

2 Related Work 117

Question answering and reasoning. Endowing
118 machines to do reasoning over explicit knowledge
119 is a primitive task (Newell and Simon, 1956b).
120 Early works tried to solve it by converting texts into
121 logic forms (Newell and Simon, 1956b; Musen and
122 van der Lei, 20p). But such kinds of approaches
123 can be affected by the error propagation caused by
124 semantic parsing (Zettlemoyer and Collins, 2012;
125 Berant et al., 2013; Berant and Liang, 2014).
126

127 Lately, question answering (QA) is employed as
128 an important task for machine reasoning. Numer-
129 ous datasets were proposed, including synthesized
130 data (Weston et al., 2016), comprehension on natu-
131 ral texts (Rajpurkar et al., 2016; Joshi et al., 2017;
132 Fisch et al., 2019) or more complex relationship
133 reasoning (Tafjord et al., 2019; Lin et al., 2019).
134 There are also multi-hop QA (Yang et al., 2018;
135 Welbl et al., 2018) or logical QA datasets (Yu et al.,
136 2020; Liu et al., 2020) where textual rules need to
137 be inferred implicitly. Plenty of studies try to solve
138 these problems via neural networks and achieve re-
139 markable performance (Joshi et al., 2020; Yu et al.,
140 2018; Shao et al., 2020). Nevertheless, nearly all of
141 them only focus on the prediction of final answers
142 and neglect the acquisition of interpretable proofs.

Proof generation. NLProlog (Weber et al., 2019)
143 first employ logic programming to search a proof
144 then predict the answer in Multi-hop QA. Re-
145 cently, Clark et al. (2020) first propose new rule-
146 based QA datasets for this line of research that
147 include proof path, and present RuleTaker to an-
148 swer questions. Saha et al. (2020) argue that pro-
149 ducing answer proofs makes models more reliable
150 and propose PROVER, a transformer-based model
151 that predicts all nodes and edges within the proof
152 at once using their embeddings. PROBER (Sun
153 et al., 2021) further improves this framework us-
154 ing probabilistic graph modeling more variables.
155 There is also an increasing interest in solving proof
156 generation iteratively. EVR (Liang et al., 2021)
157 splits the question into sub-questions, using gener-
158 ated intermediate texts to guide proof step by step.
159 ProofWriter (Tafjord et al., 2021) employs a similar
160 idea but instead using intermediate textual conclu-
161 sions and a more powerful T5-11B model (Raffel
162

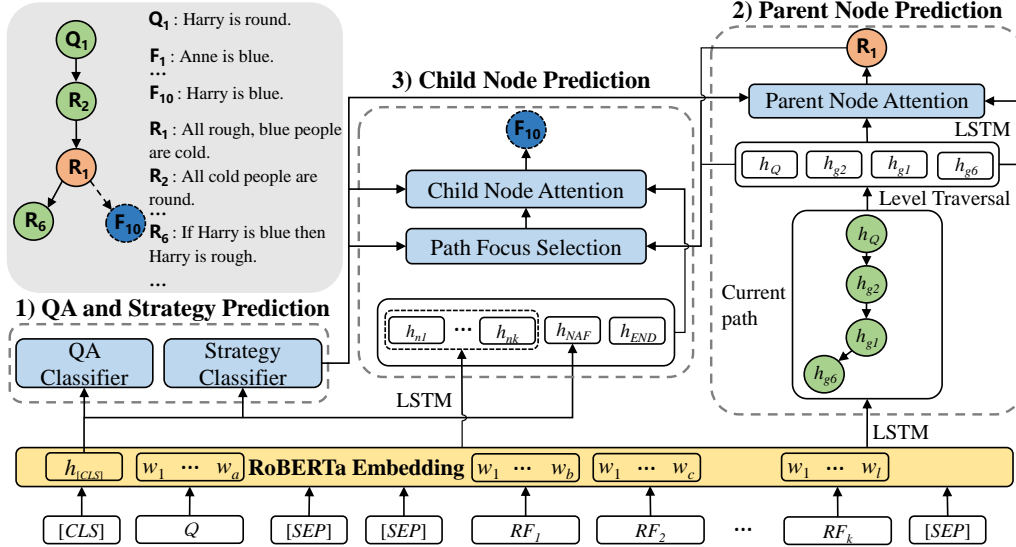


Figure 2: The Model architecture of IBR. 1) is only used once at the start, then 2) and 3) will be applied iteratively to generate the whole proof. It also illustrates the state when adding F_{10} into the proof (F: facts, R: rules).

et al., 2020) for generation, which makes it hard to reproduce. IBR is also an iterative model, being more interpretable than at-once models. It gets rid of intermediate texts and directly uses various representations to finish each reasoning step, improving efficiency and effectiveness.

3 Methodology

3.1 Task Definition

We first formulate the proof generation task as follows. Given a tuple (C, Q, A, P) , where $C = \{RF_i\}$ is the contexts containing several textual rules and facts RF , Q is the question, $A \in \{True, False\}$ is the answer, and P indicates the proof path for the detailed reasoning procedure to derive A , our goal is twofold: 1) predicting the answer A , and 2) generating the proof path P . Taking **DU0-DU5** (Clark et al., 2020) dataset as example, P is a single-directed acyclic graph having the shortest path to derive A . P can start from one or multiple nodes but must end in one node that directly entails or contradicts Q . A node in P can be a fact, a rule, or a special NAF node (Negation As Failure)¹. Edges between nodes indicate the start nodes can be used to prove the end nodes during reasoning. Proofs in the dataset can be roughly classified according to their **strategies** S to prove the question: (1)*Proof*: the question can be directly proven to be True or False using the given C and NAF; (2) *Fail-*

¹A start node when the negation condition in the next node has no corresponding fact nor rule node, it will be considered as true. E.g., there is no item in C related to “Anne is big”, its negation “Anne is not big” will be considered as true.

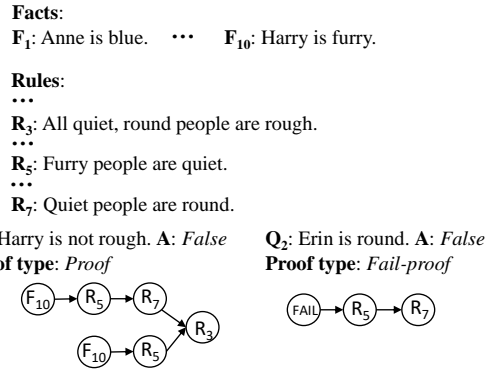


Figure 3: Examples of *Proof* and *Fail-proof* strategies.

Proof: the question cannot be explicitly deduced barely using C as some key information is missed, a positive statement is judged as **False** while a negative statement as **True** under such cases (Figure 3).

3.2 Overview

The proposed Iterative Backward Reasoning (IBR) model will take Q as the initial node and produce proof path P backward, from the end node to start node. Two actions are included at each iteration: (1) **Predicting the new parent node**, i.e. a node in the derived proof path where a child node will be added (except the first step that only Q exists); (2) **Predicting the child node**, i.e. the fact or rule in C that will be the child for the selected parent node. After each iteration, a new node and an associated edge will be added. After obtaining the whole reasoning path, we remove Q and reverse all edges to get the final proof P .

The Figure 2 illustrates our IBR model, which can be divided into three main modules, (1) **QA**

and Strategy Prediction, (2) Parent Node Prediction, and (3) Child Node Prediction. In order to make the question Q can fully interact context C (facts and rules) and obtain better representations, IBR uses pretrained RoBERTa (Liu et al., 2019) as the backbone network. The input of RoBERTa is the concatenation of the question Q and the context $C = \{RF_i\}$, separated by special $[SEP]$ tokens, denoted as $[CLS] Q [SEP] [SEP] C [SEP]$.

IBR only uses the QA prediction and strategy prediction modules once at first to predict the answer A and the strategy (refer to §3.1, where the latter one will result in different proof generation procedures. In order to improve the reasoning efficiency as well as accuracy, instead of using generated intermediate texts (Liang et al., 2021; Tafjord et al., 2021), all possible nodes (rules and facts) are represented by node embeddings in IBR. The initial state of the proof is only the representation of the question h_Q , then the rest of the reasoning path will be constructed backward from it.

Samples with *Fail-Proof* strategy differs from ones with *Proof*, because their proofs are usually short without sub-branch, and only consist of rules due to lacking essential supporting facts. To take the advantage of such a property distinction and extend the applicability compared to former models (Liang et al., 2021) that can not generate proofs for *Fail-Proof* samples, we apply different actions in modules (1) and (2) depending on the output from strategy prediction.

3.3 QA and Strategy Prediction Module

This module predicts the answer A to the question Q and the corresponding strategy S of proof P . Since the representation of $[CLS]$ token from pretrained models is proven to have the capability of modeling the whole input, we use it as the input feature for both predictions as they condition the global information. The encoded $[CLS]$ by RoBERTa, $h_{[CLS]}$ is passed to a linear layer and the softmax function σ for answer and strategy classification respectively,

$$P_{QA} = \sigma(f_{QA}(h_{[CLS]})),$$

$$P_{Strategy} = \sigma(f_{Strategy}(h_{[CLS]})).$$

Here, f_{QA} and $f_{Strategy}$ indicate the linear layer for QA classification and strategy classification, respectively. P_{QA} and $P_{Strategy}$ are binary-class probability values, the former one for values in $A \in \{True, False\}$ while the later one for values in $S \in \{Proof, Fail-proof\}$.

3.4 Parent Node Prediction Module

This module predicts which node in the current reasoning path is going to be the next parent node that a new child node will link to. To better represent the sequential information of each possible node (fact or rule), an LSTM (Hochreiter and Schmidhuber, 1997) is used to further encode the token-level embedding from RoBERTa. The hidden state in the last step will be used as the textual representation h_{gi} of a possible parent node RF_i .

In addition, selecting a node from the existing proof path also needs global and structural modeling on the history path. To make this procedure more convenient and involving the reasoning order information, the path will be regarded as a tree structure and nodes are reordered by level traversal from top to down. Since Q will always be the root node of the tree, e.g., if Q have two children RF_1 and RF_3 , and RF_1 has a child RF_2 , the reordered representation sequence is $[h_Q, h_{g1}, h_{g3}, h_{g2}]$. We then utilize another LSTM model to encode the reordered representation sequence of the current reasoning path obtained before, extracting the overall state of the path, which is the hidden state h_g at the last time step in this LSTM.

A parent node attention based on the Transformer attention (Vaswani et al., 2017), taking h_g and the representation sequence of current path $\mathbf{H}_p = [h_Q, h_{g1} \dots h_{gt}]$ as input, i.e.

$$\text{Att}(h_g, \mathbf{H}_p) = \sigma(f_Q(h_g)(f_K(\mathbf{H}_p))^T/d), \quad (1)$$

is used to obtain the weights of all possible parents nodes, where f_Q and f_K indicate linear layers, σ is a softmax function, and d is the dimension of h_g . As we discussed in §3.2, different operations will be employed for corresponding strategy types of proofs. 1) If the predicted proof strategy is *Proof*, we select the node with the highest weight as the parent node RF_p . 2) If the predicted proof strategy is *Fail-proof*, we will use the last node in the current path, i.e. h_{gt} in \mathbf{H}_p , as the parent node RF_p , because no sub-branch is included in such proofs.

3.5 Child Node Prediction Module

This module predicts which node is added to the proof path and linked to the parent node RF_p obtained before. To obtain the representations of candidate child nodes, similar to §3.4, we apply another LSTM model to the encoded RoBERTa embeddings and get h_{ni} for RF_i . Since we discussed a special NAF node in §3.1 which may contain

information from the whole context, we utilize a linear layer f_{NAF} to transform the $[CLS]$ token embedding $h_{[CLS]}$ into its representation h_{NAF} . Moreover, we initialize the representation h_{END} for a special END node, indicating the proof generation process will finish here.

During selecting the new child node, we need to consider not only the knowledge of the history path but also the state of the parent node. To better model such relationships, we propose a Path Focus Selection module to generate relevant features before predicting the child node. We employ a 2-layer Transformer model along with a LSTM model, first encode the representations of node sequence \mathbf{H}_p from Parent Node Prediction respectively, then fuse their hidden state via a linear layer f_U ,

$$h_F = f_U([\text{Trans}(h_{gp}, \mathbf{H}_p, \mathbf{H}_p); \text{LSTM}(\mathbf{H}_p)]). \quad (2)$$

Here, h_{gp} is the representation of the selected parent node in §3.4, f_U is the linear layer for feature fusing, while $[\cdot; \cdot]$ stands for concatenation. q, k, v in $\text{Trans}(q, k, v)$ indicate the inputs corresponding to Query, Key, and Value for a transformer model, and only the hidden state in the last time step is remained in both Trans and LSTM. LSTM used here is a supplementary knowledge source for a better representation. Such an operation results in a feature h_F that is aware of both the history proof path and the parent node that a child will link to.

This feature h_F will then be used in the Child Node Attention to calculate the attention weights on all possible child nodes. Particularly, an attention model same as Eq. 1 is applied on h_F and a series of child node representations obtained before $\mathbf{H}_c = [h_{n1} \dots h_{nk}, h_{NAF}, h_{END}]$, and the attention weights will be $\text{Att}(h_F, \mathbf{H}_c)$. It contains all facts and rules in the context, and the special NAF node as well as END node.

Similar to §3.4, we also apply different actions according to our predicted proof strategies before. (1) If the strategy is Proof, we will select the child node with the highest attention weight from all candidates as the new node in the proof path. (2) If the strategy is *Fail-proof*, since RF_p is the last node during reasoning and this procedure is a first-order logical under such a situation, there is no need to make complex modeling on the derived path. Therefore, we directly use its parent node representation h_{gp} rather than encoded state from Transformer in Eq. 2 to get h_F . But LSTM will be remained to maintain some basic modeling capa-

bility on the path. In child node attention, we mask all fact nodes and select the one with the highest weight among the remaining nodes, because such kinds of proofs usually only contain rules and such masking can avoid extra errors.

Moreover, it is worth noting that fact nodes have a higher priority than rule nodes during inference, due to the clearer subject information in facts.

3.6 Training and Inference

The whole model will be trained via binary cross-entropy losses from all three above modules jointly,

$$L = L_{QA} + L_{Parent} + L_{Child} + \alpha * L_{Strategy}.$$

L_{QA} and $L_{Strategy}$ correspond to the loss of QA prediction and strategy prediction, respectively. α is a hyperparameter to reweigh the influence of $[CLS]$ token. L_{Parent} is the loss for parent node prediction, where the cross-entropy is calculated between the attention weight vector and a one-hot vector indicating the gold parent node. While L_{Child} is in a similar way on child node prediction. Note that samples labeled as *Fail-proof* strategy are not involved in the training of parent node prediction, as all their proof path are chains and the new parent node is always the last node added to the path. So learning on these data may introduce model bias.

During inference, IBR will first make predictions on the answer A and strategy S , then generate parent node and child node iteratively, until the special END node is predicted as the new child node.

4 Experiments

Following former studies (Saha et al., 2020; Sun et al., 2021), we evaluate our IBR² on three datasets and four settings including fully supervised learning, few-shot learning, zero-shot learning, and generalization.

4.1 Setup

Datasets. Experiments are conducted on three datasets raised by Clark et al. (2020)³, where we use the same test split as previous works for fair comparison:

- **DU0-DU5:** Five synthesized datasets created by translating hand-crafted rules and formal language to natural language. It is divided by the highest depth of proof, where DU stands for "Depth Upto" (DU=0,1,2,3,5). Data in higher DU values will also contain samples with lower depth. Note that proofs

²Refer to Appendix A.1 for implementation details.

³More details are given in Appendix A.2

in DU0 only have one supporting or opposing fact. All related results are reported on DU5 test split.

- **Bird-Electricity**: It is a test-only dataset that contains samples about birds and electric circuits. It is generated in the same way as DU0-DU5, but are in different domains as DU0-DU5.

- **ParaRules**: This dataset consists of 40k questions expressed in paraphrased natural language based on synthetic data, which is created by crowdsourcing. Multiple facts get together in one statement here rather than separated in DU0-DU5.

Baselines. We consider the following baselines⁴.

- **RuleTaker (RT)** (Clark et al., 2020): a RoBERTa based model that only predicts answers.

- **PROVER (PV)** (Saha et al., 2020): a method that treats the proof as a graph and predicts all its nodes and edges at once, also using RoBERTa model as the backbone, same as IBR.

- **PROBR(PB)** (Sun et al., 2021): it improves PROVER by introducing the probabilistic graph that jointly consider the answer, nodes and edges.

- **EVR** (Liang et al., 2021): an iterative model that predicts the next proof item by generating textual sub-questions based on logical operator. Note that this model is **not applicable** for samples whose **proof strategy is *Fail-proof*** discussed in §3.1, so we make comparison with it separately.

Metrics. We closely follow previous works to evaluate the model performance with answer prediction (QA) accuracy and proof generation (PA) accuracy. Some samples may have multiple gold proofs. A generated proof is considered correct, as long as its nodes and edges match exactly with the nodes and the edges in any of the gold proofs. Full Accuracy (FA) is also included, where a sample is regarded as correct only both the predicted answer and proof are correct.

4.2 Results on Fully Supervised Scenario

We train IBR on the training split of the DU5 dataset and evaluate on the test split of DU5. We compare the performance of IBR with baselines except for EVR in Table 1, while with EVR in Table 2 where only partial test split is included, excluding samples whose proof strategy is *Fail-proof*. Because EVR always fails on these samples (EVR on these excluded samples is given in Appendix A.5).

Obviously, IBR achieves the best proof generation accuracy (PA) as well as full accuracy (FA)

⁴Results of baselines are obtained from the original papers or by running the released code.

	D	0	1	2	3	4	5	all
	Cnt	6299	4434	2915	2396	2134	2003	20192
QA	RT	100	98.4	98.4	98.8	99.2	99.8	99.2
	PV	100	99.0	98.8	99.1	98.8	99.3	99.3
	PB	100	99.9	99.9	100	100	100	99.9
	IBR	100	99.2	99.2	98.9	99.3	99.6	99.4
PA	PV	98.4	93.2	84.8	80.5	72.5	65.1	87.1
	PB	98.4	94.3	86.1	82.0	76.1	72.2	88.8
	IBR	99.5	95.6	93.0	90.7	86.5	81.7	93.5
FA	PV	98.4	93.1	84.8	80.5	72.4	65.1	87.1
	PB	98.4	94.3	86.1	82.0	76.1	72.2	88.8
	IBR	99.5	95.6	92.9	90.7	86.5	81.6	93.5

Table 1: Performance of different models on varying proof depth (D) when trained on DU5. Cnt: sample count, RT:RuleTaker, PV: PROVER, PB: PROBR.

	D	0	1	2	3	4	5	all
	Cnt	1934	1934	1934	1934	1934	1934	11604
QA	EVR	99.4	99.3	96.9	93.3	88.9	88.3	94.4
	IBR	100	99.3	99.6	99.3	99.6	99.5	99.5
PA	EVR	95.8	92.5	87.7	79.3	77.3	68.8	83.6
	IBR	98.8	96.4	94.7	92.2	88.7	83.6	92.4
FA	EVR	95.8	92.5	87.7	79.3	77.3	68.8	83.6
	IBR	98.8	96.3	94.6	92.2	88.7	83.5	92.3

Table 2: Performance of IBR compared to EVR on a partial test split of DU5(exclude Fail-proof samples) after training on train split of DU5.

among all baseline models, on samples with every depth. Our model also shows a greater advantage on samples with the deeper proof path, e.g. 81.7% vs. 72.2 on PA when depth is 5, illustrating the superiority of iterative models on complex proof path. Besides, despite not being the best on answer accuracy (QA), our model only has very narrow gaps with the best one, which proves that IBR is still a comprehensive model covering both subtasks. When compare to EVR, also an iterative model, IBR shows a significantly stronger performance on all metrics, benefiting from our elaborate two-fold reasoning process at each step.

4.3 Results on Few-shot Scenario

We also explore the performance of IBR when training using fewer samples, ranging from 10k to 30k to all the examples (70k) in DU5. The comparison between our model, PROVER(PV), and PROBR(PB) is shown in Table 3, in all three metrics. Our model significantly has the best proof generation performance than the other two baselines under all cases, due to the iterative architecture requiring less global modeling capability and thus fewer training samples. Although PB shows a

Data	QA			PA			FA		
	PV	PB	IBR	PV	PB	IBR	PV	PB	IBR
70k	99.3	99.9	99.4	87.1	88.8	93.5	87.1	88.8	93.5
30k	97.8	99.9	98.3	72.5	86.8	89.8	72.4	86.8	89.7
10k	87.1	99.9	94.3	44.0	72.4	75.7	42.7	72.3	75.4

Table 3: Few-shot performance comparison among IBR, PROVER (PV), and PROBR (PB) on the full test split of DU5 after trained on partial DU5 samples.

Data	QA		PA		FA	
	EVR	IBR	EVR	IBR	EVR	IBR
70k	94.4	99.5	83.6	92.4	83.6	92.3
30k	95.7	99.4	84.4	88.2	84.4	88.1
10k	96.2	97.9	82.8	71.2	82.8	70.8

Table 4: Few-shot performance comparison among EVR and IBR on partial test split of DU5 (without *Fail-proof* samples) after trained on partial DU5 samples.

promising answer prediction accuracy under few-shot settings, the performance of IBR is close to it while better than PV, e.g., 94.3 vs. 87.1 under 10k. In addition, in Table 4, we also compare with EVR under the same settings but using a different test set that excludes *Fail-proof* samples. EVR outperforms IBR under the 10k setting for proof generation, but IBR is stronger if more training samples are available.

4.4 Results on Zero-shot Scenario

We further test the out-of-domain performance of IBR against baselines on Birds-Electricity dataset to evaluate their robustness, where B1 and B2 are two sets from the birds domain, and E1-E4 are four sets from the electricity domain. Results are shown in Table 5 and Table 6. Note that *Fail-proof* samples are still not involved in the comparison for EVR. Overall, our IBR achieves 2.5% promotion in PA while an equivalent result on QA, compared to PROVER. Despite being the best one on QA, PROBR is also defeated by IBR on both PA and FA. In addition, our model shows more improvement on the hardest E3 and E4 subsets, which further verifies its robustness. When it comes to EVR, we can find its cross-domain capability is relatively weak as it sees a significant drop on PA, and IBR is superior to it without any doubt. Because the cross-domain generation for intermediate texts is much harder, our usage of high-level node features to finished reasoning can alleviate this challenge.

4.5 Generalization Ability

Generalize to higher depths. Following the previous work (Sun et al., 2021), we test the generaliza-

Test	B1	B2	E1	E2	E3	E4	all	
Cnt	40	40	162	180	624	4224	5270	
QA	RT	97.5	100.0	96.9	98.3	91.8	76.7	80.1
	PV	95.0	95.0	100.0	100.0	89.7	84.8	86.5
	PB	100.0	100.0	100.0	100.0	98.2	95.6	96.3
	IBR	100.0	97.5	100.0	100.0	89.2	84.1	86.0
PA	PV	92.5	95.0	95.1	91.7	72.3	80.6	80.7
	PB	100.0	100.0	97.5	93.3	79.3	77.7	79.3
	IBR	100.0	100.0	95.6	94.4	80.2	82.4	83.2
FA	PV	92.5	95.0	95.1	91.7	71.8	80.6	80.5
	PB	100.0	100.0	97.5	93.3	79.3	77.7	79.3
	IBR	100.0	97.5	95.6	94.4	78.2	82.4	82.9

Table 5: Zero-shot performance comparison among RuleTakers(RT), PROVER(PV), and PROBR(PB) on Birds-Electricity dataset after training on DU5.

Test	B1	B2	E1	E2	E3	E4	all	
Cnt	28	28	72	90	312	1206	1736	
QA	EVR	67.8	64.2	83.3	80.0	76.2	83.8	81.6
	IBR	100.0	96.4	100.0	100.0	92.9	100.0	98.6
PA	EVR	32.1	35.7	58.3	50.0	45.5	70.3	63.1
	IBR	100.0	100.0	91.6	91.1	91.3	95.2	94.3
FA	EVR	32.1	32.1	58.3	50.0	45.5	70.3	63.1
	IBR	100.0	96.4	91.6	91.1	87.1	95.2	93.5

Table 6: Zero-shot performance comparison among EVR and IBR on partial Birds-Electricity dataset (exclude *Fail-proof* samples) after training on DU5.

tion ability of IBR by first training the model on the training splits of DU0, DU1, DU2, and DU3, then test them on the test split of DU5 with deeper proof path respectively⁵. Results are shown in Table 7. We notice that all models suffer performance degeneration especially the depth of training set is lower because it is hard for the model to learn complex reasoning based on simple proof path. However, IBR still realizes the best performance in terms of PA and FA, especially on **DU3**, where it gets 4.2% PA/FA promotion to PROBR and even outperforms PROVER trained on the whole **DU5** data. These observations again prove that iterative approaches can better learn the detailed reasoning step by step, obtaining a better generalization capability than at-once models.

Generalize to complex language. We also test whether IBR can be applied to samples where questions and statements are expressed in more human-like natural language. Following Clark et al. (2020), we train models on the combined training partitions of DU3 and ParaRules and test them on the ParaRules test set. Table 8 demonstrates that our

⁵We remove the position embedding in path focus selection to proceed to this test, see Appendix A.1 for details

Data	QA			PA			FA		
	PV	PB	IBR	PV	PB	IBR	PV	PB	IBR
DU0	68.7	56.9	53.5	44.4	50.7	47.0	42.8	41.3	47.0
DU1	73.7	97.7	73.1	63.8	63.9	64.6	61.9	63.9	64.5
DU2	89.6	99.9	89.6	72.6	74.5	76.3	72.3	74.4	76.2
DU3	98.6	99.9	98.6	79.1	83.2	87.4	79.1	83.2	87.4
DU5	99.3	99.9	99.4	87.1	88.8	93.5	87.1	88.8	93.4

Table 7: Performance of generalization ability between PROVER (PV), PROBR (PB), and IBR when testing on the test split of DU5, after trained on DU0, DU1, DU2, DU3, and DU5, respectively.

	D	0	1	2	3	4	all
	Cnt	2968	2406	1443	1036	142	8008
QA	PV	99.7	98.6	98.2	96.5	88.0	98.4
	PB	99.8	99.7	99.9	99.8	100	99.8
	IBR	99.9	98.8	97.5	96.3	88.7	98.4
PA	PV	99.5	98.0	88.9	90.0	76.1	95.4
	PB	99.5	98.0	88.9	90.1	82.4	95.6
	IBR	99.8	98.8	91.1	89.0	75.3	95.9
FA	PV	99.4	97.3	88.7	89.9	76.1	95.1
	PB	99.4	98.0	88.9	90.1	82.4	95.5
	IBR	99.7	98.1	90.9	89.0	75.3	95.7

Table 8: Performance on ParaRules test set, after trained on combined D3+ParaRules training partitions, including PROVER (PV), PROBR (PB), and IBR.

534 model sees a slight promotion on PA/FA while a
535 similar accuracy as PROVER on QA, indicating
536 that IBR has good applicability when reasoning on
537 more complicated and natural texts.

538 5 Analysis

539 5.1 Ablation Study

540 To explore the effects between different compo-
541 nents in our model, we consider the following abla-
542 tions: 1) IBR +Gold-Parent: given the true parent
543 node during inference to explore the accuracy of
544 child node prediction; 2)IBR +Gold-Child: given
545 the true child node to verify the accuracy of parent
546 node prediction; 3) *w/o* QA: removing QA task
547 in loss to check its impact on proof generation; 4)
548 *w/o* node LSTM: using mean pooling rather than
549 LSTM encoding to get the representations of nodes;
550 5) *w/o* focus LSTM: Removing the supplementary
551 LSTM in path focus selection.

552 Through results on the whole DU5 test split
553 given in Table 9, it can be seen that giving either
554 gold parent nodes or child nodes can benefit the per-
555 formance especially giving gold children, signifi-
556 cating our parent node prediction achieves promising
557 accuracy while the prediction of child nodes can be
558 improved. Moreover, IBR can still learn to gener-

Models	QA	PA	FA
IBR	99.4	93.5	93.5
IBR +Gold-Parent	99.4	95.6	95.3
IBR +Gold-Child	99.4	99.6	99.3
<i>w/o</i> QA	-	93.7	-
<i>w/o</i> node LSTM	99.5	93.2	93.2
<i>w/o</i> focus LSTM	99.6	92.6	92.4

Table 9: Results of ablation studies on DU5 dataset.

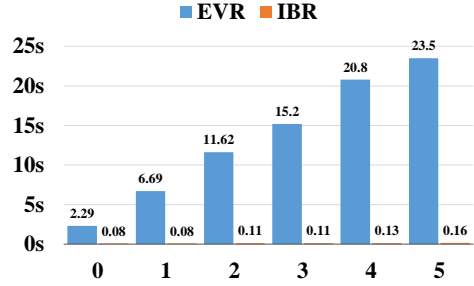


Figure 4: Per-sample inference runtime (in second) of EVR and IBR, on DU5 dataset with varying depths.

559 ate proofs without supervision from answers, while
560 LSTM encoders attribute to a better representation
561 of the node or the path that has been derived.

562 5.2 Latency Analysis

563 To demonstrate the computational efficiency of
564 IBR, we compare the per sample inference time
565 of IBR to EVR, also an iterative proof generation
566 model, on the test split of DU5. Both models are
567 tested on one NVIDIA Tesla-V100 GPU with
568 the same batch size. As shown in Figure 4, our IBR
569 could achieve up to $\times 146.9$ speedup compared with
570 EVR, benefiting from our reasoning based on node
571 and path features rather than intermediate texts. It
572 is also noticeable that the runtime of EVR grows
573 linearly with depth, while such an effect is slight
574 on our model. Because EVR needs to infer on all
575 contexts at every step, but IBR uses a simplified
576 parent node prediction based on the derived path.

577 6 Conclusion

578 This paper presents IBR, a proof generation model
579 via iterative backward reasoning for rule-based QA
580 tasks. We equip the reasoning procedure with de-
581 tailed tracking by predicting nodes and edges in the
582 proof path iteratively backward from the question,
583 and allow the model to reason on the elaborate rep-
584 resentations of nodes and history path. Our model
585 is more interpretable than previous at-once mod-
586 els, and more effective and efficient than former
587 iterative models. Experiments also demonstrate
588 the superiority of IBR on proof generation under
589 various settings.

590
591
592
593
594
595
596
597

598
599
600
601
602
603

604
605
606
607
608
609
610
611

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

References

Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. [Semantic parsing on Freebase from question-answer pairs](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1533–1544, Seattle, Washington, USA. Association for Computational Linguistics.

Jonathan Berant and Percy Liang. 2014. [Semantic parsing via paraphrasing](#). In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1415–1425, Baltimore, Maryland. Association for Computational Linguistics.

Deng Cai and Wai Lam. 2019. [Core semantic first: A top-down approach for AMR parsing](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3799–3809, Hong Kong, China. Association for Computational Linguistics.

Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2020. [Transformers as soft reasoners over language](#). In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pages 3882–3890. ijcai.org.

Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. 2019. [DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2368–2378, Minneapolis, Minnesota. Association for Computational Linguistics.

Adam Fisch, Alon Talmor, Robin Jia, Minjoon Seo, Eunsol Choi, and Danqi Chen. 2019. [MRQA 2019 shared task: Evaluating generalization in reading comprehension](#). In *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pages 1–13, Hong Kong, China. Association for Computational Linguistics.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. [Long short-term memory](#). *Neural computation*, 9(8):1735–1780.

Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2020. [SpanBERT: Improving pre-training by representing and predicting spans](#). *Transactions of the Association for Computational Linguistics*, 8:64–77.

Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke Zettlemoyer. 2017. [TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension](#). In *Proceedings of the 55th Annual Meeting of*

the Association for Computational Linguistics (Volume 1: Long Papers), pages 1601–1611, Vancouver, Canada. Association for Computational Linguistics.

Zhengzhong Liang, Steven Bethard, and Mihai Surdeanu. 2021. [Explainable multi-hop verbal reasoning through internal monologue](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1225–1250, Online. Association for Computational Linguistics.

Kevin Lin, Oyvind Tafjord, Peter Clark, and Matt Gardner. 2019. [Reasoning over paragraph effects in situations](#). In *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pages 58–62, Hong Kong, China. Association for Computational Linguistics.

Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2020. [Logiqa: A challenge dataset for machine reading comprehension with logical reasoning](#). In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, pages 3622–3628. ijcai.org.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#). *ArXiv preprint*, abs/1907.11692.

Mark A. Musen and Johan van der Lei. 20p. [Of brittleness and bottlenecks: Challenges in the creation of pattern-recognition and expert-system models](#). *ArXiv preprint*, abs/p.

Allen Newell and Herbert Simon. 1956a. [The logic theory machine—a complex information processing system](#). *IRE Transactions on information theory*, 2(3):61–79.

Allen Newell and Herbert A. Simon. 1956b. [The logic theory machine—a complex information processing system](#). *IRE Trans. Inf. Theory*, 2:61–79.

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](#). *Journal of Machine Learning Research*, 21(140):1–67.

Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ questions for machine comprehension of text](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.

Swarnadeep Saha, Sayan Ghosh, Shashank Srivastava, and Mohit Bansal. 2020. [PProver: Proof generation for interpretable reasoning over rules](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 122–136, Online. Association for Computational Linguistics.

702	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.	In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.	759 760 761 762
703			
704			
705			
706			
707			
708			
709	Changzhi Sun, Xinbo Zhang, Jiangjie Chen, Chun Gan, Yuanbin Wu, Jiase Chen, Hao Zhou, and Lei Li. 2021. Probabilistic graph reasoning for natural proof generation . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 3140–3151, Online. Association for Computational Linguistics.	Adams Wei Yu, David Dohan, Minh-Thang Luong, Rui Zhao, Kai Chen, Mohammad Norouzi, and Quoc V. Le. 2018. Qanet: Combining local convolution with global self-attention for reading comprehension . In <i>6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings</i> . OpenReview.net.	763 764 765 766 767 768 769 770
710			
711			
712			
713			
714			
715			
716	Oyvind Tafjord, Bhavana Dalvi, and Peter Clark. 2021. ProofWriter: Generating implications, proofs, and abductive statements over natural language . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 3621–3634, Online. Association for Computational Linguistics.	Weihaoyu Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. Reclor: A reading comprehension dataset requiring logical reasoning . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	771 772 773 774 775 776
717			
718			
719			
720			
721			
722	Oyvind Tafjord, Matt Gardner, Kevin Lin, and Peter Clark. 2019. QuaRTz: An open-domain dataset of qualitative relationship questions . 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 5941–5946, Hong Kong, China. Association for Computational Linguistics.	Luke S Zettlemoyer and Michael Collins. 2012. Learning to map sentences to logical form: structured classification with probabilistic categorial grammars . <i>ArXiv preprint</i> , abs/1207.1420.	777 778 779 780
723			
724			
725			
726			
727			
728			
729			
730	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need . In <i>Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA</i> , pages 5998–6008.		
731			
732			
733			
734			
735			
736			
737	Leon Weber, Pasquale Minervini, Jannes Münchmeyer, Ulf Leser, and Tim Rocktäschel. 2019. NLProlog: Reasoning with weak unification for question answering in natural language . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 6151–6161, Florence, Italy. Association for Computational Linguistics.		
738			
739			
740			
741			
742			
743			
744	Johannes Welbl, Pontus Stenetorp, and Sebastian Riedel. 2018. Constructing datasets for multi-hop reading comprehension across documents . <i>Transactions of the Association for Computational Linguistics</i> , 6:287–302.		
745			
746			
747			
748			
749	Jason Weston, Antoine Bordes, Sumit Chopra, and Tomás Mikolov. 2016. Towards ai-complete question answering: A set of prerequisite toy tasks . In <i>4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings</i> .		
750			
751			
752			
753			
754			
755	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 .		
756			
757			
758			

A Appendix

A.1 Implementation Details

Parameter	Value
Training Epochs	8
Optimizer	AdamW
Batch Size	16
RoBERTa Learning rate	1e-5
QA and Strategy Pre Learning rate	1e-5
Parent Node Pre Learning rate	2e-4
Child Node Pre Learning rate	5e-4
All LSTM Learning rate	1e-3
Dropout Rate	0.1
LSTM hidden state for parent node and child node encoding	1024
LSTM hidden state for path encoding in parent node prediction	1024
Transformer hidden state in path focus selection	1024
LSTM hidden state in path focus selection	256
Seed	42

Table 10: Implementation details of IBR.

We implement our model based on PyTorch along with Huggingface-Transformers toolkit⁶. We use RoBERTa_{Large} model⁷ as our backbone encoder to generate token-level representations. Table 10 shows the implementation details of IBR, including learning rates for different modules. All linear layers used in our model have one layer. The model trained after 8 epochs will be used in the evaluation. Each epoch training takes 2 hours. We select these hyper-parameters according to tuning them empirically based on the performance on the validation set of DU0-DU5. All experiments are run on NVIDIA Tesla-V100 GPUs. The performance of IBR fluctuates within one point. We release our code anonymously at <https://anonymous.4open.science/r/IBR-ECC8>.

A.2 Dataset Details

We will introduce the details of three datasets used in our experiment, all of them are firstly applied in rule-based QA and proof generation tasks in Clark et al., 2020.

⁶<https://github.com/huggingface/transformers>

⁷<https://huggingface.co/roberta-large>

Split	D	Num	Fail-proof Num	Proof Num	Avg. Node
Train	0	21,359	14,597	6,762	0.62
	1	15,380	8,618	6,762	1.82
	2	10,112	3,350	6,762	3.37
	3	8,389	1,627	6,762	4.98
	4	7,456	694	6,762	6.90
	5	6,987	225	6,762	9.26
	all	69,683	29,111	40,572	3.35
Test	0	6,299	4,365	1,934	0.59
	1	4,434	2,500	1,934	1.77
	2	2,915	981	1,934	3.36
	3	2,396	462	1,934	4.99
	4	2,134	200	1,934	6.98
	5	2,003	69	1,934	9.47
	all	20,181	8,577	11,604	3.33

Table 11: The statistics of train and test split in DU5 dataset. *Fail-proof* and *Proof* indicate different proof strategies we discussed in §3.1. Avg.Node indicates the average node number in proof path.

DU0-DU5: A series of synthesized datasets where rules and facts are all generated via manually designed logical programming, while questions are generated by combining random logical operations among them. Data are divided into 5 subsets according to their maximum reasoning depth (D) in the proof path, D = 0, 1, 2, 3, 5. There are 100k questions in each subset, where 70k / 10k / 20k samples in the train / validation / test partition respectively. D = 0 means that the question can be proven directly using a fact in contexts. In our experiment in §4, we only use the data from DU5 for testing because it covers all possible depths, while the train set is the train split in DU5 except §4.5, where we use train split from DU0, DU1, DU2 and DU3 for training. We provide some statistics of DU5 in Table 11.

Birds-Electricity: It is a set of data that only contains 5k test samples for the evaluation of robustness and out-of-domain performance of models. The Bird data only requires reasoning up to depth 1 and 2 (B1 and B2), while Electricity data have reasoning depths ranging from 1 to 4. Both of them include new vocabulary that is not included in DU0-DU5.

ParaRules: A more challenging dataset contains paraphrased samples on the synthesized ones via crowdsourcing. It has 40k questions against about 2k theories. The statements are expressed in a more natural way, posing a discrepancy between DU0-DU5. It has 28k / 4k / 8k samples in the train / validation / test split respectively. In §4.5, we combine it with the extensive DU3 for training,

resulting in a train set containing 119k samples.

A.3 Possible Limitations of Our Model

Since our strategy prediction module and operations corresponding to different strategies in node prediction modules are specially designed for the current datasets, we may need to redesign some specific operations if some new proof types are included in new datasets to reach the best performance. But we believe our architecture will still take effect without modification.

A.4 Strategy Accuracy of IBR

D	Cnt	Strategy Accuracy
0	6299	99.9
1	4434	99.1
2	2915	99.3
3	2396	99.0
4	2134	99.2
5	2003	99.7
All	20192	99.4

Table 12: Strategy accuracy of IBR on test split of DU5 after training on training split of DU5.

We provide the strategy prediction accuracy on DU5 in Table 12. It proves that IBR is also well able to make the predictions on the proof strategies. This is partly due to RoBERTa’s powerful representation capability. On the other hand, there is a certain connection between the answer to the question and the strategy, and there are some common elements at the semantic representation level that can be learned together.

A.5 Performance of EVR and IBR on *Fail-proof* Samples

As we have discussed in §4.2, EVR (Liang et al., 2021) is not applicable for samples containing *Fail-proof* proofs, because it cannot obtain proper intermediate questions to proceed correct following reasoning. Here, we compare our model with EVR on these samples in DU0-DU5, as illustrated in Table 13. Although EVR can achieve promising performance on answer prediction (QA) for these samples, it cannot generate any correct proof path under such cases, which have already been discussed in its original paper.

	D	0	1	2	3	4	5	all
	Cnt	4365	2500	981	462	200	69	8577
QA	EVR	99.7	99.1	98.9	99.1	98.5	100	99.4
	IBR	100	99.1	98.3	97.6	96.5	100	99.3
PA	EVR	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	IBR	99.8	95.0	89.5	84.4	65.5	28.9	95.0
FA	EVR	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	IBR	99.8	95.0	89.5	84.4	65.5	28.9	95.0

Table 13: The performance of EVR and IBR on the partial test split of DU5 that only contains samples whose proofs strategies are *Fail-proof*.

A.6 Proof Generation samples

We provide some proof generation samples in Figure 5 for a better understanding of this task, where questions, all contexts, and the proof path generated by our IBR are given (all consistent with the given labels).

Rules:

- R₁: If someone is nice and kind then they like the bear.
- R₂: If someone sees the dog and they eat the bear then the bear is cold.
- R₃: If someone is big then they eat the cat.
- R₄: If someone is big then they do not see the rabbit.
- R₅: If someone is not big and they do not eat the dog then the dog is cold.
- R₆: If someone is cold then they like the rabbit.
- R₇: If someone likes the rabbit then they see the dog.
- R₈: If the dog eats the cat then the dog is kind.
- R₉: If someone likes the dog and they do not eat the cat then the dog eats the bear.

Facts:

- F₁: The bear eats the cat.
- F₂: The bear eats the rabbit.
- F₃: The cat eats the dog.
- F₄: The cat eats the rabbit.
- F₅: The cat likes the bear.
- F₆: The cat sees the rabbit.
- F₇: The dog is round.
- F₈: The dog likes the bear.
- F₉: The dog likes the cat.
- F₁₀: The dog sees the bear.
- F₁₁: The rabbit eats the bear.
- F₁₂: The rabbit is big.
- F₁₃: The rabbit is cold.
- F₁₄: The rabbit is not kind.
- F₁₅: The rabbit does not like the cat.
- F₁₆: The rabbit sees the bear.

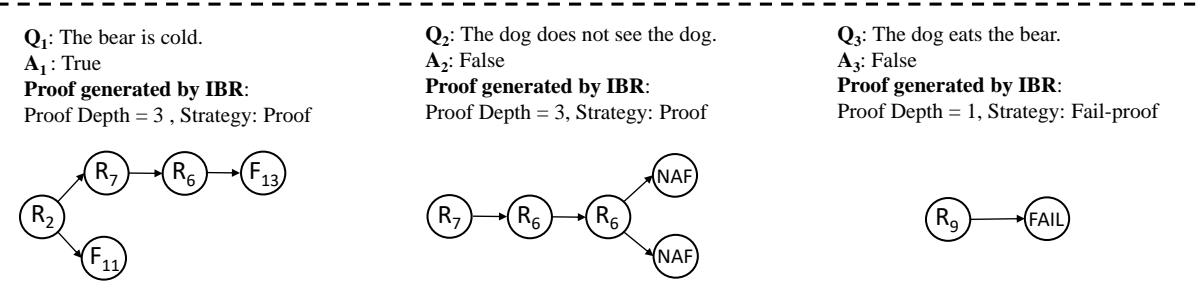


Figure 5: Some proof cases generated by IBR, along with all contexts and questions, including two proof strategies, *Proof* and *Fail-proof*.