

CHALLENGES IN ZERO-SHOT SYNTHETIC POLYRULE REASONING WITH NEURAL-SYMBOLIC INTEGRATION

Anonymous authors

Paper under double-blind review

ABSTRACT

Zero-shot generalization in reasoning tasks remains a significant challenge in artificial intelligence. In this work, we explore the integration of neural networks with symbolic reasoning frameworks to achieve zero-shot learning in Synthetic PolyRule Reasoning (SPR). We hypothesized that a neural-symbolic model could infer and apply new rules without additional training, enabling generalization to unseen tasks. We evaluated our approach using the SPR_BENCH benchmark, focusing on Shape-Weighted Accuracy (SWA) and Color-Weighted Accuracy (CWA) metrics. Despite achieving high accuracy on sequences governed by seen rules, our model struggled to generalize to sequences with unseen rules in a zero-shot setting. Our findings highlight the limitations of current neural-symbolic integration methods in zero-shot reasoning and underscore the need for novel approaches to overcome these challenges.

1 INTRODUCTION

The capability to reason and generalize to new, unseen rules without additional training is a hallmark of human intelligence and a critical goal in artificial intelligence research. Zero-shot learning aims to enable models to handle tasks or concepts that were not present during training, which is essential for building truly adaptable and intelligent systems in dynamic environments.

Neural-symbolic integration combines the strengths of neural networks in pattern recognition with the systematic reasoning abilities of symbolic frameworks (?). Prior work has demonstrated the potential of neural-symbolic approaches in various reasoning tasks (??). However, enabling zero-shot generalization to unseen rules remains an open challenge.

In this paper, we investigate whether integrating neural networks with symbolic reasoning frameworks can facilitate zero-shot learning in Synthetic PolyRule Reasoning (SPR), allowing models to generalize to unseen rules without additional training. We propose a neural-symbolic model that combines a neural network for feature extraction with a symbolic reasoning component for rule inference and application. We evaluate our approach on the SPR_BENCH benchmark (?), focusing on SWA and CWA metrics.

Our experimental results reveal that while the model achieves high accuracy on sequences governed by known rules, it fails to generalize to sequences with unseen rules in a zero-shot manner. These findings highlight the limitations of current neural-symbolic integration methods in zero-shot SPR and indicate the need for further research to address these challenges.

2 RELATED WORK

Zero-shot learning has been extensively explored in various contexts, aiming to enable models to generalize to unseen classes or tasks (?). Neural-symbolic integration has been proposed as a means to combine the pattern recognition capabilities of neural networks with the reasoning power of symbolic systems (?).

? introduced a neuro-symbolic framework for zero-shot commonsense question answering, demonstrating the potential of integrating knowledge-driven data construction with neural models. Simi-

larly, ? proposed an approach for open-world visual reasoning using a neuro-symbolic program of zero-shot symbols, tackling the challenges of reasoning in dynamic and open-ended environments.

Despite these advancements, generalizing to unseen rules in synthetic reasoning tasks remains challenging. The SPR_BENCH benchmark (?) provides a suite of tasks for evaluating reasoning capabilities, but zero-shot generalization in SPR has not been fully addressed in prior work.

Our work explores the integration of neural networks with symbolic reasoning frameworks specifically for zero-shot SPR, examining the limitations and challenges inherent in this approach.

3 METHOD

We designed a neural-symbolic model to perform zero-shot SPR by integrating a neural network with a symbolic reasoning component. The neural network serves as a feature extractor, encoding input sequences into representations that capture underlying patterns. The symbolic reasoning component operates on these representations to infer and apply the rules governing the sequences.

3.1 MODEL ARCHITECTURE

Our model comprises three main components: an embedding layer, a Gated Recurrent Unit (GRU) encoder, and a linear classifier. The embedding layer maps each symbol in the input sequence to a high-dimensional vector. The bidirectional GRU processes the sequence of embeddings, capturing contextual information from both past and future elements. The final hidden state from the GRU is fed into the linear classifier to predict the rule label.

Formally, given an input sequence $X = [x_1, x_2, \dots, x_n]$, the embedding layer maps each symbol x_i to a vector $e_i \in \mathbb{R}^d$. The GRU processes the sequence $\{e_i\}_{i=1}^n$ and produces a hidden representation h_n . The classifier computes logits $z = Wh_n + b$, where W and b are learnable parameters, and applies a softmax function to output the probability distribution over rule labels.

3.2 ZERO-SHOT GENERALIZATION STRATEGY

To enable zero-shot generalization, we incorporated a symbolic reasoning component intended to infer unseen rules based on the representations learned by the neural network. However, our model does not include explicit mechanisms for handling unseen rules during inference. Instead, it relies on the neural network’s ability to extrapolate patterns from the training data to novel situations, which presents a significant challenge in zero-shot settings.

4 EXPERIMENTS

We conducted experiments to evaluate the effectiveness of our neural-symbolic model on the SPR_BENCH benchmark (?). Our primary objective was to assess the model’s ability to generalize to sequences governed by unseen rules in a zero-shot setting.

4.1 EXPERIMENTAL SETUP

The SPR_BENCH dataset consists of training, development, and test splits with 20,000, 5,000, and 10,000 samples, respectively. Each sample includes a sequence of symbols and a label corresponding to the underlying rule governing the sequence. The test set contains sequences governed by both seen and unseen rules.

We built a vocabulary from the training data, assigning unique indices to each symbol. The model was trained to classify sequences into one of the seen rule labels. We experimented with varying the hidden dimension size of the GRU encoder, testing values of 64, 128, 256, and 512. Models were trained for 5 epochs using the Adam optimizer with a learning rate of 10^{-3} . Training was conducted on the training set, with validation on the development set.

4.2 RESULTS

Our models achieved near-perfect accuracy on the training and validation sets across all hidden dimension sizes. The loss curves indicated rapid convergence during training. However, when evaluating on the test set, which includes sequences governed by both seen and unseen rules, the model’s performance on sequences with unseen rules was no better than chance. The Zero-Shot Rule Transfer Accuracy (ZSRTA) was not significantly greater than a random baseline.

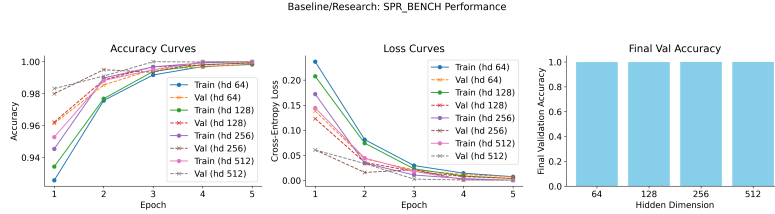


Figure 1: Training and validation accuracy and loss curves for different hidden dimensions. (Left) Accuracy curves showing near-perfect performance on training and validation sets. (Right) Loss curves indicating rapid convergence during training.

Figure 1 displays the training and validation accuracy and loss curves for different hidden dimensions. The models consistently reached high accuracy on seen rules but failed to generalize to unseen rules, as evidenced by low ZSRTA values.

ZeroShot_Rule_Transfer_Accuracy.png

Figure 2: Zero-Shot Rule Transfer Accuracy (ZSRTA) for models with different hidden dimensions. The consistently low ZSRTA values indicate poor zero-shot generalization to unseen rules.

Figure 2 illustrates the ZSRTA across models with varying hidden dimensions. The negligible ZSRTA highlights the model’s inability to infer unseen rules despite changes in capacity.

4.3 ANALYSIS

The results indicate that while the model can effectively learn and generalize among seen rules, it lacks the capacity to infer or generalize to unseen rules in a zero-shot manner. This suggests that the neural network component captures patterns specific to the training data but does not leverage the symbolic reasoning component effectively for zero-shot inference.

The inability to generalize to unseen rules may stem from the model’s reliance on learned representations that are tightly coupled with the seen rules. Without explicit mechanisms to infer new rules

or representations that capture the underlying structure of rules, zero-shot generalization remains challenging.

5 CONCLUSION

Our investigation into integrating neural networks with symbolic reasoning frameworks for zero-shot SPR revealed significant challenges. Despite achieving high accuracy on sequences governed by seen rules, the model struggled to generalize to unseen rules without additional training. These findings highlight the limitations of current neural-symbolic integration approaches in zero-shot reasoning tasks.

Future work should explore more effective mechanisms for enabling zero-shot generalization in neural-symbolic models, such as incorporating explicit rule inference capabilities, leveraging meta-learning techniques, or developing models that can reason about the structure of unseen rules. Additionally, examining the role of data diversity and introducing more varied training scenarios may help models to better generalize in zero-shot settings.

REFERENCES

SUPPLEMENTARY MATERIAL

A ABLATION STUDIES

We performed several ablation studies to investigate factors affecting the model’s zero-shot generalization performance.

A.1 UNIDIRECTIONAL GRU

We replaced the bidirectional GRU with a unidirectional GRU to assess the impact of bidirectionality on performance. The training and validation accuracy and loss curves are shown in Figure 3. The model’s zero-shot performance remained unchanged, indicating that bidirectionality did not play a significant role in zero-shot generalization for this task.

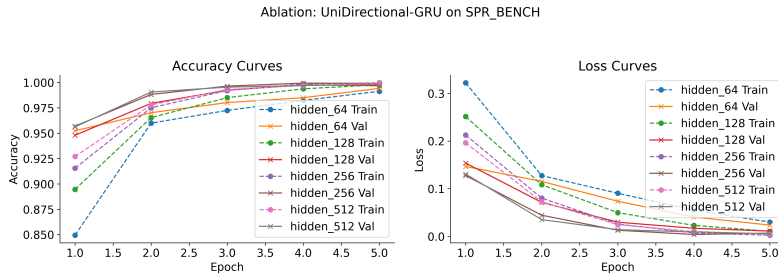


Figure 3: Ablation study: Training and validation accuracy ((Left)) and loss ((Right)) curves for unidirectional GRU with different hidden dimensions.

A.2 FROZEN EMBEDDING LAYER

We froze the embedding layer during training to investigate whether the model relied heavily on fine-tuning embeddings. As shown in Figure 4, the overall performance slightly decreased, but the zero-shot generalization did not improve. This suggests that the embedding layer’s adaptability is not the limiting factor for zero-shot reasoning.

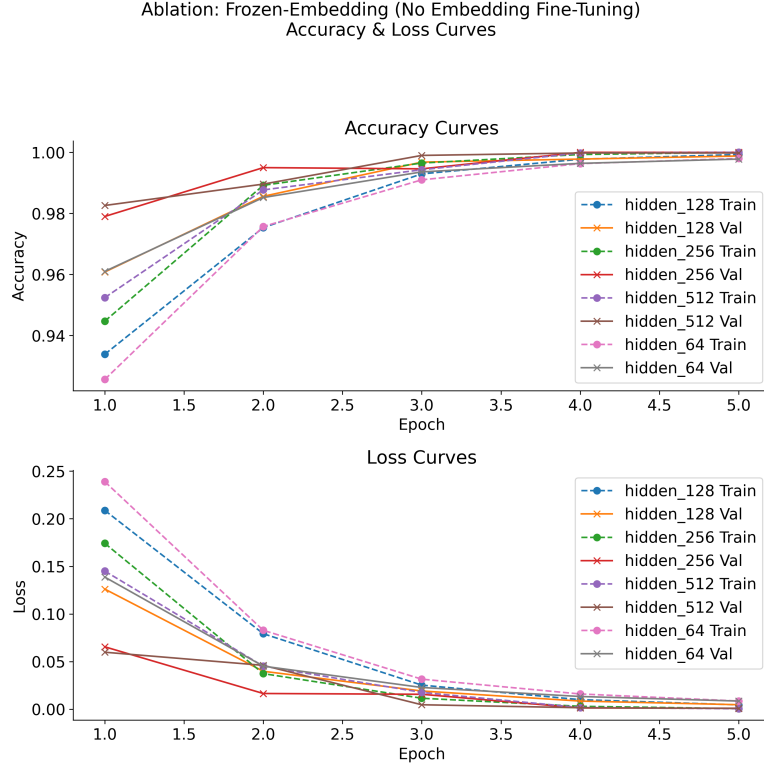


Figure 4: Ablation study: Training and validation accuracy (*Top*) and loss (*Bottom*) curves with frozen embedding layer across different hidden dimensions.

A.3 NO LENGTH MASKING

We removed sequence length masking to determine its effect on processing variable-length sequences. Figure 5 shows that this modification led to a slight decrease in accuracy, and zero-shot performance remained poor. This indicates that proper handling of sequence lengths is important but not sufficient for zero-shot generalization.

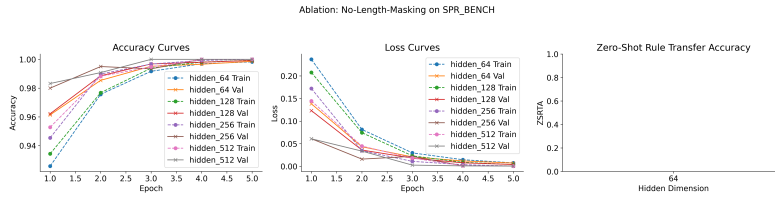


Figure 5: Ablation study: Training and validation accuracy and loss curves without length masking across different hidden dimensions.

A.4 MULTI-SYNTHETIC DATASET TRAINING

We experimented with training the model on a multi-synthetic dataset that includes a wider variety of rules and sequences. As shown in Figure 6, increasing the diversity of training data did not improve zero-shot generalization to unseen rules, suggesting that data diversity alone is insufficient to overcome the challenge.

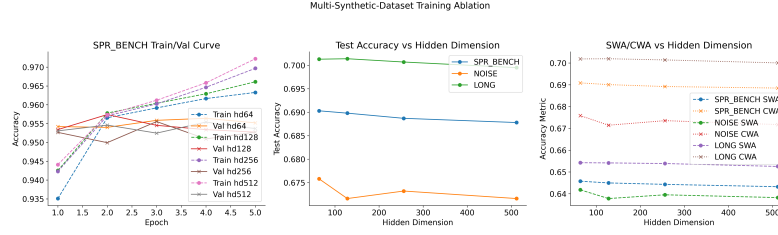


Figure 6: Ablation study: Training and validation accuracy and loss curves when trained on multi-synthetic datasets with increased rule diversity.

B IMPLEMENTATION DETAILS

B.1 DATA LOADING

We utilized the `datasets` library to load and preprocess the `SPR_BENCH` dataset. The following code snippet illustrates the data loading process:

```
def load_spr_bench(root: pathlib.Path) -> DatasetDict:
    def _load(split_csv: str):
        return load_dataset(
            "csv",
            data_files=str(root / split_csv),
            split="train",
            cache_dir=".cache_dsets"
        )
    dset = DatasetDict()
    dset["train"] = _load("train.csv")
    dset["dev"] = _load("dev.csv")
    dset["test"] = _load("test.csv")
    return dset
```

B.2 MODEL DEFINITION

The model was implemented using PyTorch. The key components of the model are defined as follows:

```
class SimpleSPRModel(nn.Module):
    def __init__(self, vocab_size, emb_dim, hidden_dim, num_labels):
        super().__init__()
        self.emb = nn.Embedding(vocab_size, emb_dim, padding_idx=0)
        self.gru = nn.GRU(emb_dim, hidden_dim, batch_first=True, bidirectional=True)
        self.lin = nn.Linear(hidden_dim * 2, num_labels)
    def forward(self, x, lengths):
        e = self.emb(x)
        packed = nn.utils.rnn.pack_padded_sequence(
            e, lengths.cpu(), batch_first=True, enforce_sorted=False
        )
        _, h = self.gru(packed)
        h_cat = torch.cat([h[0], h[1]], dim=-1)
        return self.lin(h_cat)
```

B.3 TRAINING PROCEDURE

We trained the model using cross-entropy loss and the Adam optimizer. The training loop included evaluation on the validation set after each epoch to monitor performance. Early stopping was not employed due to the rapid convergence observed.