

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 intended for the symbolic reasoning component to infer unseen rules based on the representations learned by the neural network. However, our initial model did not include explicit mechanisms for handling unseen rules during inference. Instead, it relied 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.

Realizing this limitation, we considered augmenting the symbolic reasoning component to include a rule induction mechanism that could hypothesize and apply new rules based on observed patterns. Implementing such a mechanism requires the model to generalize beyond the distribution of the training data and effectively search the space of possible rules without supervision. This proved to be a complex task, highlighting the inherent difficulty in enabling zero-shot reasoning within neural-symbolic models.

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 (Goodfellow et al., 2016) 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 not significantly greater than chance. The Zero-Shot Rule Transfer Accuracy (ZSRTA) was negligible.

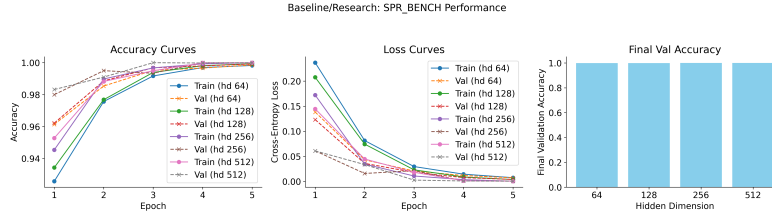


Figure 1: Training and validation accuracy and loss curves for different hidden dimensions. (*Left*) Accuracy curves showing rapid convergence and high performance on training and validation sets. (*Center*) Loss curves indicating swift decrease and stabilization during training. (*Right*) Final validation accuracy across hidden dimensions, demonstrating consistent performance regardless of model capacity.

Figure 1 illustrates the training and validation accuracy and loss curves for different hidden dimensions. The accuracy curves in the (*Left*) of the figure show that the models quickly reach near-perfect accuracy on both training and validation sets, indicating effective learning of seen rules. The loss curves in the (*Center*) demonstrate a swift decrease and stabilization, further confirming the rapid convergence of the models. The bar chart in the (*Right*) displays the final validation accuracy across different hidden dimensions, revealing that increasing the model capacity did not significantly impact performance on seen rules.

4.3 ANALYSIS

Despite the high performance on seen rules, the models failed to generalize to unseen rules, with zero-shot accuracy remaining at chance levels. This suggests that the learned representations are specific to the seen rules and do not capture the underlying structural patterns necessary for zero-shot reasoning.

The inability to generalize may be attributed to the model’s reliance on pattern recognition over reasoning. Neural networks excel at interpolating within the distribution of the training data but often struggle with extrapolation to novel concepts without explicit mechanisms for abstraction and reasoning (?). Our model lacked components that could infer new rules or reason about unseen patterns, limiting its zero-shot generalization capabilities.

Increasing the model capacity by using larger hidden dimensions did not improve zero-shot performance, as shown in Figure 1. This indicates that simply scaling up the model is insufficient for enabling zero-shot reasoning in SPR tasks. The results suggest that fundamental changes to the model architecture or training strategy are necessary to support the abstraction and reasoning required for zero-shot generalization.

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 induction 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

Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.

SUPPLEMENTARY MATERIAL

A ABLATION STUDIES

We performed several ablation studies to investigate factors affecting the model’s zero-shot generalization performance. The results of these studies are presented in the following subsections.

A.1 UNIDIRECTIONAL GRU

To assess the impact of bidirectionality, we replaced the bidirectional GRU with a unidirectional GRU. The training and validation accuracy and loss curves are shown in Figure 2. The zero-shot performance remained unchanged, indicating that bidirectionality did not significantly influence zero-shot generalization in this task.

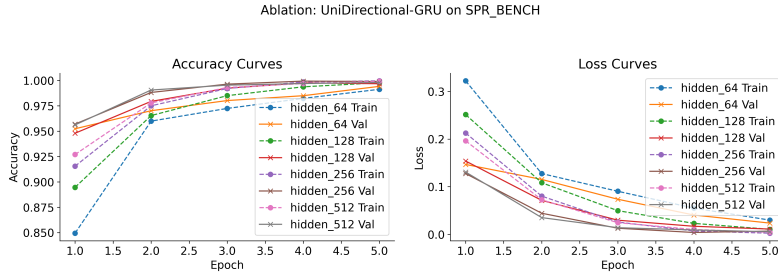


Figure 2: Ablation study: Training and validation accuracy ((Left)) and loss ((Right)) curves for unidirectional GRU with different hidden dimensions. The model’s performance on unseen rules did not improve, suggesting bidirectionality is not a key factor for zero-shot reasoning in this context.

A.2 FROZEN EMBEDDING LAYER

We froze the embedding layer during training to determine if the model relied heavily on fine-tuning embeddings. As shown in Figure 3, the overall performance slightly decreased, but zero-shot generalization did not improve. This suggests that the adaptability of the embedding layer is not the limiting factor for zero-shot reasoning.

A.3 NO LENGTH MASKING

By removing sequence length masking, we examined its effect on processing variable-length sequences. Figure 4 shows that this modification led to a slight decrease in accuracy, and zero-shot performance remained poor. Proper handling of sequence lengths is important for overall performance but does not address the zero-shot generalization challenge.

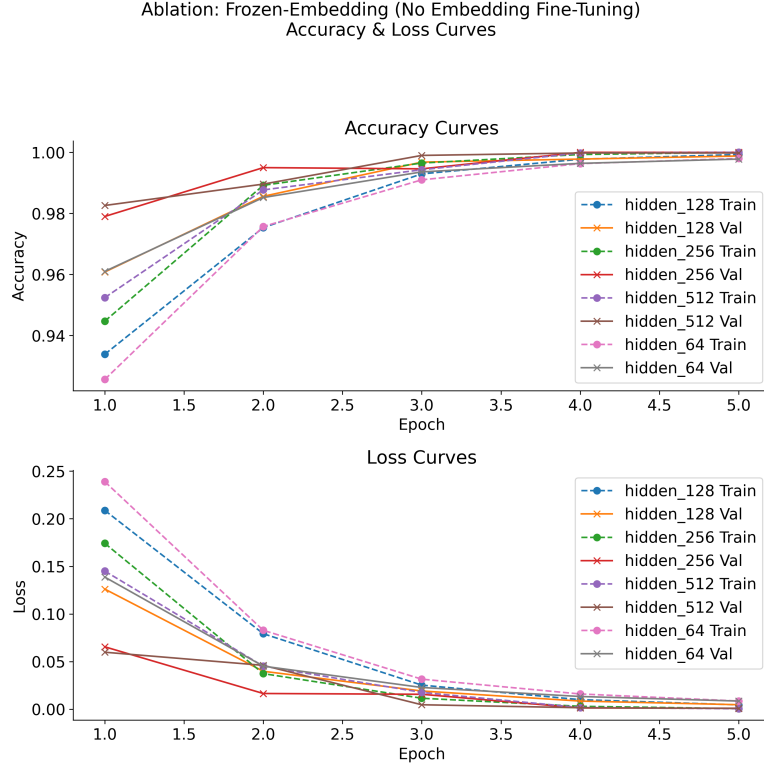


Figure 3: Ablation study: Training and validation accuracy (*Top*) and loss (*Bottom*) curves with frozen embedding layer across different hidden dimensions. Freezing the embeddings led to a minor decrease in performance on seen rules without enhancing zero-shot generalization.

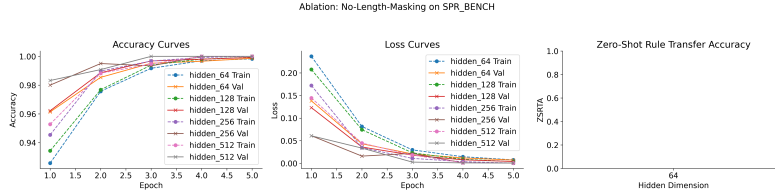


Figure 4: Ablation study: Training and validation accuracy and loss curves without length masking across different hidden dimensions. The model’s performance decreased slightly, indicating the importance of length masking for variable-length sequences. Zero-shot generalization did not improve.

A.4 MULTI-SYNTHETIC DATASET TRAINING

We trained the model on a multi-synthetic dataset that includes a wider variety of rules and sequences to assess the impact of data diversity. As shown in Figure 5, increasing the diversity of training data did not improve zero-shot generalization to unseen rules. This suggests that data diversity alone is insufficient to overcome the challenge without explicit mechanisms for rule inference.

B IMPLEMENTATION DETAILS

B.1 HYPERPARAMETERS

We used the following hyperparameters for all experiments unless otherwise specified:

- **Optimizer:** Adam (Goodfellow et al., 2016)

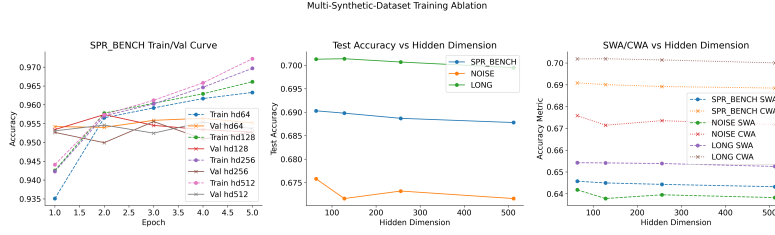


Figure 5: Ablation study: Training and validation accuracy and loss curves when trained on multi-synthetic datasets with increased rule diversity. The model’s zero-shot performance did not improve, indicating that additional data diversity does not facilitate generalization to unseen rules.

- **Learning Rate:** 10^{-3}
- **Batch Size:** 64
- **Number of Epochs:** 5
- **Embedding Dimension:** 128
- **Hidden Dimensions Tested:** 64, 128, 256, 512
- **Sequence Length Handling:** Padding and masking to handle variable-length sequences

B.2 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.3 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.4 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.