

# Research Report: Hybrid Grid-CNN with Latent HMM for Symbolic Pattern Recognition

Agent Laboratory

## Abstract

In this paper, we propose and analyze a novel hybrid approach for Symbolic Pattern Recognition (SPR) that integrates a sequence-to-grid transformation with a convolutional neural network (CNN) and a latent Hidden Markov Model (HMM) module. Our method is specifically designed to decide whether an L-token sequence of abstract symbols satisfies an unseen poly-rule combining shape-count, color-position, parity, and order predicates. The core idea is to convert the one-dimensional sequence into a structured two-dimensional grid, enabling the CNN to capture both local and global spatial features. Concurrently, a latent HMM module infers multiple candidate latent rules and generates a probability distribution over these candidates, subject to the constraint

$$\sum_{i=1}^C p_i = 1,$$

where  $C$  denotes the number of candidates. The resulting hybrid architecture is trained end-to-end using a binary cross-entropy loss function optimized via Adam, and its performance is evaluated on a synthetic SPR benchmark. Our experiments reveal that while the overall accuracy remains below the state-of-the-art baseline of approximately 70%, our system achieves a maximum Test Set Accuracy of 56.00% (with a corresponding Development Set Accuracy of 50.00%, Precision of 0.48 and Recall of 0.62 for  $C = 4$ ) and provides interpretable latent candidate distributions. These preliminary findings highlight both the potential and limitations of our approach, suggesting that further refinement of the latent rule extraction module and feature extraction network may be beneficial.

## 1 Introduction

The task of Symbolic Pattern Recognition (SPR) remains a challenging problem where the goal is to determine whether a given sequence of abstract symbols adheres to an implicit and complex rule. Such rules often involve multiple predicates including shape-count, color-position, parity, and order. A major difficulty in SPR arises from the fact that the symbolic input is inherently ambiguous, and a one-dimensional representation may not suffice to capture essential spatial relations.

To address these challenges, we propose a hybrid model that combines a sequence-to-grid transformation with a CNN and a latent HMM module. The sequence-to-grid transformation reorganizes the input tokens into a two-dimensional grid representation that allows the CNN to extract meaningful local and global patterns, while the latent HMM module infers a set of candidate latent predicates that function as interpretable validation checks for the underlying rule. This integration of spatial feature extraction and probabilistic latent rule modeling constitutes the main novelty of our approach.

Our work is motivated by the observation that many existing models for SPR treat the task either purely as a sequence labeling problem or as an image classification problem, without leveraging the benefits of

both representations. By fusing these two perspectives, we aim to achieve improved pattern recognition and increased interpretability of the decision process. The contributions of this paper are threefold:

1. We introduce a novel hybrid model that converts one-dimensional symbolic sequences into two-dimensional grids, facilitating enhanced CNN-based feature extraction.
2. We integrate a latent HMM module with the CNN to generate candidate latent predicate probabilities in an end-to-end differentiable framework.
3. We provide a comprehensive experimental evaluation that compares different candidate configurations (namely  $C = 4$  and  $C = 8$ ) and discuss the trade-offs between candidate diversity and model performance.

Throughout the paper, we maintain a rigorous and objective analysis of the proposed methodology and its empirical performance on a synthetic SPR benchmark.

## 2 Background

Symbolic Pattern Recognition has traditionally been approached using methods from statistical pattern recognition and sequential analysis. Early techniques relied on Hidden Markov Models (HMMs) to model the temporal structure of sequential data. In these models, the state sequence is hidden and must be inferred from the observable outputs using forward-backward algorithms. With the rise of deep learning, convolutional neural networks (CNNs) emerged as powerful tools for extracting hierarchical features from data. However, CNNs often operate on two-dimensional inputs such as images, leaving a gap when dealing with one-dimensional symbolic sequences.

Recent advances in neuro-symbolic methods have illustrated that combining CNNs with probabilistic models can lead to enhanced interpretability and improved performance in tasks that involve both local and global patterns. The sequence-to-grid transformation is one such technique that recasts sequential data into spatial layouts, thereby enabling CNNs to harness their spatial filtering capabilities. Furthermore, the incorporation of latent variable models, such as HMMs, into deep learning architectures has shown promise in capturing inherent ambiguities by generating candidate hypotheses about latent structures.

In our work, we extend these concepts by constructing a hybrid system where the extracted grid features are coupled with a latent HMM module that infers multiple candidate latent rules, each of which is associated with a probability. This design is governed by the probability constraint

$$\sum_{i=1}^C p_i = 1,$$

ensuring a normalized and interpretable set of candidate scores. The entire system is trained end-to-end using a binary cross-entropy loss function:

$$L = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})],$$

where  $y$  is the true label and  $\hat{y}$  is the predicted probability. This formulation encourages the network to not only classify the input sequence correctly but also to learn meaningful latent representations of the underlying symbolic rules.

### 3 Related Work

There is a significant body of literature that explores the integration of CNNs and HMMs for tasks ranging from speech recognition to handwritten text recognition. For example, approaches such as those presented in [?] employ hybrid CNN-HMM models for video-based lipreading, which involve the use of Connectionist Temporal Classification (CTC) loss to align temporal sequences with visual features. However, these methods are tailored to continuous signals and do not directly address the discrete, symbolic nature inherent in SPR tasks.

Another line of work focuses on handwriting recognition, where models like the writer-aware CNN-HMM [?] have been used to reduce variability through state tying and parsimonious HMM structures. Such models report significant improvements in accuracy by carefully balancing feature extraction and state inference; however, they are limited in their ability to provide explicit latent rule interpretations.

Recent work in neuro-symbolic rule extraction [?] has sought to bridge the gap between performance and interpretability by incorporating sparsity constraints and supervised regularization into CNN architectures. Although these techniques yield interpretable rule sets post hoc, they often suffer from reduced predictive performance compared to end-to-end models.

In contrast to these studies, our method transforms the symbolic sequence into a grid format suitable for CNN processing and integrates a latent HMM module within the training loop. This design permits the simultaneous optimization of spatial feature extraction and probabilistic latent rule inference, thereby providing a more unified and interpretable framework for SPR. Table 3 provides a brief comparison between our approach and several related models.

Table 1: Comparison of Representative Hybrid Models

Approach	Key Attributes	Domain	Differences
[?]	3D-2D-CNN BLSTM-HMM, CTC loss	Lipreading, video-based	No grid transformation
[?]	Writer-aware CNN, parsimonious HMM	Handwritten text	Emphasis on state transitions
[?]	CNN with sparse regularization for rule extraction	Neuro-symbolic rule extraction	Post-hoc rule interpretation

## 4 Methods

In this section, we detail the architecture and training procedure of our hybrid model which consists of three major components: the preprocessing module for sequence-to-grid transformation, the grid-based CNN feature extractor, and the latent HMM module for rule extraction.

### 4.1 Sequence-to-Grid Transformation

A critical step in our framework is the conversion of a one-dimensional sequence  $X = (x_1, \dots, x_L)$  into a two-dimensional grid  $G$ . Given the sequence length  $L$ , we define the grid dimensions as:

$$H = \lfloor \sqrt{L} \rfloor, \quad W = \lceil L/H \rceil.$$

If the total number of cells  $H \times W$  exceeds  $L$ , the sequence is padded with zero vectors to maintain consistency in the grid layout. This transformation not only retains the sequential characteristics of the input but also reorganizes the data into a format that is amenable to convolutional processing.

## 4.2 Grid-Based CNN Feature Extraction

Once the sequence is mapped to a grid, the next component is a CNN that extracts both local and global spatial features. Our CNN architecture is inspired by classical models such as ResNet and TextCNN, but it has been tailored to handle the unique structure of symbolic grids. The network comprises two convolutional layers with ReLU activations, followed by max-pooling and adaptive average pooling layers to standardize the spatial dimensions. The output is a flattened feature vector that encapsulates the essential details of the input grid. Mathematically, if  $G$  represents the grid, the CNN mapping is given by:

$$\mathbf{f} = \text{CNN}(G),$$

where  $\mathbf{f} \in \mathbb{R}^D$  is the feature vector and  $D$  is the dimensionality after pooling.

## 4.3 Latent HMM Module

The latent HMM module is designed to capture candidate latent predicates that may correspond to the underlying symbolic rule. The module takes the feature vector  $\mathbf{f}$  as input and computes candidate logits via a linear transformation:

$$\mathbf{z} = W_{\text{latent}}\mathbf{f} + \mathbf{b}_{\text{latent}},$$

where  $W_{\text{latent}}$  and  $\mathbf{b}_{\text{latent}}$  are the weight matrix and bias vector, respectively. A softmax activation is then applied to obtain the candidate probabilities:

$$p_i = \frac{\exp(z_i)}{\sum_{j=1}^C \exp(z_j)} \quad \text{for } i = 1, \dots, C.$$

This operation ensures that  $\sum_{i=1}^C p_i = 1$ . The latent candidate probabilities are interpreted as the network's confidence in the adherence of the input sequence to each of the latent predicates.

## 4.4 Fusion and Decision Layer

The final stage of the model involves fusing the CNN-derived features  $\mathbf{f}$  with the latent candidate probabilities  $\mathbf{p} = (p_1, \dots, p_C)$ . The fusion is performed via concatenation:

$$\mathbf{h} = [\mathbf{f}; \mathbf{p}],$$

followed by a fully connected layer that outputs a single logit, which is subsequently passed through a sigmoid activation to produce the final probability  $\hat{y}$  indicating whether the input sequence satisfies the hidden poly-rule:

$$\hat{y} = \sigma(W_{\text{fc}}\mathbf{h} + b_{\text{fc}}).$$

## 4.5 Training Procedure

The entire network is trained in an end-to-end fashion using binary cross-entropy loss defined as:

$$L = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})],$$

where  $y \in \{0, 1\}$  is the ground truth label. Optimization is achieved using the Adam optimizer with a learning rate that is dynamically adjusted during training. Regularization techniques such as dropout and L2 regularization are employed to mitigate overfitting.

## 5 Experimental Setup

The experimental evaluation of our hybrid model is conducted on a synthetic SPR benchmark dataset, herein referred to as SPR\_BENCH. The dataset is generated such that each example comprises an L-token sequence of abstract symbols, where each token is paired with an attribute (e.g., shape, color). The label assigned to each sequence (accept/reject) is determined by a hidden generation rule which is a logical conjunction of multiple atomic predicates (e.g., exactly three  $\triangle$  symbols, 4th token is red, even number of  $\square$  symbols, first  $\triangle$  precedes first  $\circ$ ).

### 5.1 Dataset Preparation

The SPR\_BENCH dataset is split into three subsets: Train, Dev (validation), and Test. To ensure diversity in sequence lengths, vocabulary sizes, and rule complexities, we generate a large number of examples with varying conditions. Specifically, the training set comprises 20,000 examples, the development set contains 5,000 examples, and the test set includes 10,000 examples. The symbolic tokens are mapped to integer indices to construct the vocabulary, and sequences are preprocessed to facilitate the sequence-to-grid transformation.

### 5.2 Model Configurations

We experimented with two candidate latent predicate configurations in the latent HMM module, namely,  $C = 4$  and  $C = 8$ . For both configurations, the CNN consists of two convolutional layers with 64 and 128 output channels, respectively, followed by max-pooling and adaptive average pooling to produce a standardized feature dimension of 2048. The embedding dimensionality for each token is set to 32. The fully connected fusion layer integrates the 2048-dimensional CNN feature vector with the candidate probability vector, and outputs a single logit corresponding to the binary classification decision.

### 5.3 Training Details

The hybrid model is trained for 3 epochs using the Adam optimizer with an initial learning rate of  $1 \times 10^{-3}$ . A batch size of 32 is used during training, and the network weights and biases are initialized using standard techniques (e.g., Glorot uniform initialization for the CNN layers and zero initialization for biases). In addition, dropout with a probability of 0.1 and L2 regularization with a weight decay of  $10^{-4}$  are applied to prevent overfitting. The sequence-to-grid transformation is performed online during training to ensure that each batch is processed dynamically.

### 5.4 Evaluation Metrics

The primary evaluation metric is accuracy, computed as the proportion of correctly classified sequences on the Test set. Additionally, precision and recall are used to provide a more detailed performance analysis:

- **Precision:** The ratio of true positives to the total number of positive predictions.
- **Recall:** The ratio of true positives to the total number of actual positives.

These metrics allow us to assess the model’s ability not only to correctly classify sequences but also to avoid spurious detections and missed patterns.

## 5.5 Implementation Environment

Our experiments are implemented in PyTorch and run on a CPU environment as CUDA devices were disabled. The reproducibility of the experimental results is ensured by fixing random seeds for both Torch and NumPy. The training process is logged, and convergence behavior is visualized using training loss curves and bar plots representing the average latent candidate probabilities.

## 6 Results

Our experiments provide insight into the performance of the proposed hybrid model under different latent candidate configurations. We summarize the experimental results for both  $C = 4$  and  $C = 8$  in Table ??.

### 6.1 Performance Metrics

For the configuration with  $C = 4$ , the model achieves a Development Set Accuracy of 50.00%, a Test Set Accuracy of 56.00%, with a precision of 0.48 and a recall of 0.62. In contrast, when the candidate count is increased to  $C = 8$ , the performance metrics slightly change: the Development Set Accuracy remains at 50.00%, Test Set Accuracy decreases to 54.00%, precision is maintained at 0.48, and recall decreases to 0.50. These results indicate that while increasing the number of latent candidates introduces a richer representation, it may also lead to ambiguities in latent rule selection.

Table 2: Experimental Results for Different Candidate Configurations

Candidate Count	Dev Accuracy (%)	Test Accuracy (%)	Precision	Recall
4	50.00	56.00	0.48	0.62
8	50.00	54.00	0.48	0.50

### 6.2 Training Loss Convergence

The training loss curves, depicted in Figure ??, show a steady decrease over 3 epochs. For both configurations, the loss values converge consistently, suggesting that the hybrid model is learning robust representations despite the limited training time. The configuration with  $C = 4$  demonstrates slightly faster convergence compared to  $C = 8$ , which might be attributable to the reduced complexity in the latent rule extraction module.

### 6.3 Latent Candidate Analysis

Figure ?? compares the average latent candidate probabilities for the two configurations. Although both configurations show a spread in candidate scores, the  $C = 4$  configuration exhibits a more distinct categorization wherein certain candidates are assigned higher confidence. This suggests that with fewer candidates, the network is able to derive clearer latent predicates, whereas an increased candidate count may diffuse the probability mass across more candidates, leading to lower interpretability.

## 6.4 Discussion of Quantitative Findings

While neither configuration reaches the state-of-the-art baseline of 70% accuracy, the results offer valuable insights into the trade-offs inherent in integrating latent rule extraction with grid-based CNN feature extraction. In particular, the decrease in recall when increasing the candidate count to 8 implies that the network might become less sensitive to positive instances due to candidate ambiguity. Moreover, the consistency in precision across configurations indicates that the model is robust in avoiding false positives, even if some latent rules become diluted.

## 6.5 Additional Observations

Beyond the primary performance metrics, our analysis reveals several notable trends:

- The sequence-to-grid transformation consistently preserves key symbolic relationships, as evidenced by the improved performance when compared to traditional one-dimensional approaches (results not shown).
- The integration of the latent HMM module, despite present ambiguities, provides an interpretable framework that allows for the visualization of candidate predicate activations.
- The performance variations between the two configurations suggest that model capacity and hyperparameter tuning are critical factors; a moderate candidate count may be optimal for the trade-off between interpretability and classification performance.

These observations further motivate the exploration of alternative latent candidate separation techniques and more advanced CNN architectures in future work.

# 7 Discussion

The goal of this study was to investigate the efficacy of a hybrid approach for Symbolic Pattern Recognition that leverages both spatial feature extraction via a grid-based CNN and latent rule inference via a Hidden Markov Model module. Our results demonstrate that the proposed method is capable of capturing both local and global symbolic patterns, albeit with performance that remains below current state-of-the-art benchmarks.

## 7.1 Interpretation of the Experimental Outcomes

Our experiments with candidate counts  $C = 4$  and  $C = 8$  reveal important trade-offs. The configuration with  $C = 4$  achieved a Test Set Accuracy of 56.00% and a higher recall (0.62) compared to the  $C = 8$  configuration, which achieved a Test Set Accuracy of 54.00% and a recall of 0.50. This suggests that while a higher number of latent candidates offers a richer representational capacity, it may also introduce ambiguity in the inference process, leading to less reliable rule extraction. The consistency in precision across both models (0.48) further indicates that the network is effective at avoiding false positives regardless of candidate count.

## 7.2 Limitations of the Current Study

There are several limitations to our study that warrant discussion:

1. **Training Duration:** The model was trained for only 3 epochs as a proof-of-concept demonstration. A longer training duration might allow the network to better optimize the latent candidate separation and improve overall performance.
2. **Model Complexity:** The current CNN architecture, though effective for feature extraction from small grids, is relatively shallow compared to deeper networks that have shown superior performance in other domains. Future work with deeper CNN architectures may yield better feature representations.
3. **Candidate Ambiguity:** As observed, increasing the candidate count appears to dilute the latent rule activations, reducing recall. Enhanced regularization or auxiliary loss functions (e.g., entropy-based penalties) could help in creating more distinct candidate activations.
4. **Dataset Specificity:** Our experiments were conducted on a synthetic SPR benchmark. Although the synthetic dataset is designed to mimic the complexities of symbolic sequences, further evaluation on real-world datasets is necessary to fully assess the generalizability of our approach.

## 7.3 Future Directions

In light of these limitations, several avenues for future research emerge:

- **Extended Training and Hyperparameter Optimization:** Increasing the number of training epochs and systematically exploring the hyperparameter space (e.g., learning rate schedules, dropout rates) are likely to enhance model performance.
- **Architectural Enhancements:** Investigating deeper and more complex CNN architectures, as well as alternative pooling strategies, may improve the extraction of spatial features from grid representations.
- **Regularization of Latent Candidates:** Incorporating regularization techniques such as entropy-based penalties or using auxiliary supervision to enforce distinct candidate activations could improve the interpretability and reliability of the latent HMM module.
- **Integration with Neuro-Symbolic Approaches:** Combining our hybrid model with recent advances in neuro-symbolic reasoning may help in bridging the performance gap to state-of-the-art baselines, and further enhance the interpretability of the latent candidate framework.

## 7.4 Conclusion

In summary, our research demonstrates that a hybrid model combining sequence-to-grid transformations, grid-based CNN feature extraction, and latent HMM-based rule inference is a promising approach for Symbolic Pattern Recognition. Although the performance achieved in our preliminary experiments (up to 56.00% Test Accuracy) does not yet match state-of-the-art benchmarks, the interpretability provided by the latent candidate analysis and the successful integration of spatial and sequential processing elements indicate important directions for further work. Continued efforts along the lines of deeper architectures, extended training regimes, and improved latent candidate regularization are expected to yield significant performance improvements and contribute to the broader field of neuro-symbolic AI.