
Appendices

Table of Contents

- Appendix A: Additional Background on HMMs
- Appendix B: Additional Details of Experimental Setup
- Appendix C: Details of Benchmark Models
- Appendix D: Additional Synthetic Experiment Results
- Appendix E: Ablations on LLMs
- Appendix F: Spectral Learning HMMs for Prediction Task
- Appendix G: Additional Real World Experiments

A Additional Background on HMMs

In this section, we define in detail the HMM settings we are interested in, including the conditions for Markov chains to converge to unique stationary distributions. Recall that a HMM is characterized by the Markov chain’s *initial state distribution* and its *state transitions*, along with the *emission probabilities* of an observation given the hidden state. With finitely many states and observations, without loss of generality, states take values in $\mathcal{X} = \{1, 2, \dots, M\}$ while observations take values in $\mathcal{O} = \{1, 2, \dots, L\}$. The initial state distribution is denoted as $\boldsymbol{\pi} \in \mathbb{R}^M$ with π_j the probability of starting in state j , the state transitions are describe by the matrix $\mathbf{A} \in \mathbb{R}^{M \times M}$ with elements a_{ij} the probability of transitioning to state j from state i , and the emission matrix $\mathbf{B} \in \mathbb{R}^{M \times L}$ contains b_{jl} the probability of observing l when in hidden state j . The triple $\lambda = (\boldsymbol{\pi}, \mathbf{A}, \mathbf{B})$ completely parameterizes a finite-alphabet HMM.

Let $\{X_1, X_2, \dots\}$ denote a discrete-time Markov chain taking values in \mathcal{X} with transition matrix \mathbf{A} . Let $p_{ij}^{(n)} = \mathbb{P}(X_{t+n} = j | X_t = i)$ denote the n -step transition probability between states $i, j \in \mathcal{X}$. State j is said to be *accessible* from state i if there exists an integer $n \geq 1$ such that $p_{ij}^{(n)} > 0$. A subset $\mathcal{C} \subseteq \mathcal{X}$ is called *irreducible* if every pair of states $i, j \in \mathcal{C}$ is mutually accessible. The *period* of state i is defined as $c(i) = \gcd\{n \geq 1 : p_{ii}^{(n)} > 0\}$, the greatest common divisor of all possible return times. State i is *aperiodic* if $c(i) = 1$. A Markov chain is termed *geometrically ergodic* if it is irreducible and aperiodic, which guarantees convergence to a unique *stationary distribution* $\boldsymbol{\mu} \in \mathbb{R}^M$ satisfying $\boldsymbol{\mu} = \boldsymbol{\mu}\mathbf{A}$. The *mixing rate* $\rho \in [0, 1)$ is such that for all states $i, j \in \mathcal{X}$, there exists a constant $C \geq 0$ for which $|p_{ij}^{(n)} - \mu_j| \leq C\rho^n$ for all $n \geq 1$. For a finite-alphabet HMM, ρ equals λ_2 , the second-largest eigenvalue of \mathbf{A} . We run experiments on a few non-ergodic cases, while the majority of HMMs are with ergodic state transitions to avoid dependence on the initial state.

The *entropy* $H(X)$ of a discrete random variable X is defined as $H(X) = -\sum_{x \in \mathcal{X}} p(x) \log p(x)$. A fundamental property of entropy is that conditioning reduces uncertainty: for any two random variables X and Y , we have $H(X|Y) \leq H(X)$, with equality holding if and only if X and Y are statistically independent [1]. By applying the chain rule of entropy, the joint entropy of a stochastic process can be expressed as $H(X_1, X_2, \dots, X_n) = \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_1)$. For a Markov chain with stationary distribution $\boldsymbol{\mu}$, the *entropy rate* is defined as $H(\mathcal{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, X_2, \dots, X_n) = -\sum_{i,j} \mu_i a_{ij} \log a_{ij}$, which depends solely on the transition matrix \mathbf{A} . We additionally define the entropy of the emission matrix \mathbf{B} as $-\sum_{j,l} \mu_j b_{jl} \log b_{jl}$, which quantifies the average uncertainty in observations given the underlying states. Although the entropy rate of the observation process in a HMM has no known closed-form expression, it can be bounded as $H(O_n | O_{n-1}, \dots, O_1, X_1) \leq H(\mathcal{O}) \leq H(O_n | O_{n-1}, \dots, O_1)$. As \mathbf{A} defines transitions from X_t to X_{t+1} , and \mathbf{B} determines sampling O_t from X_t , the entropies of \mathbf{A} and \mathbf{B} combined help us to control the entropy lower bound of the sampled HMM sequence.

44 B Additional Details of Experimental Setup

45 **Construct \mathbf{A} with specific mixing rate, entropy, and steady state distribution.** For an ergodic
 46 Markov chain that converges to a unique stationary distribution, the stochastic matrix \mathbf{A} can be
 47 decomposed into eigenvalues and eigenvectors with the ordering shown in Figure 1, where $\vec{1} \in \mathbb{R}^M$
 48 is a vector of ones, λ_2 is the second-largest eigenvalue of \mathbf{A} , and μ is the stationary distribution
 49 [2]. We leverage this decomposition to construct \mathbf{A} with predefined λ_2 and μ . To determine
 50 the remaining eigenvalues and eigenvectors, we formulate an optimization problem based on the
 51 following requirements: (1) all entries of \mathbf{A} are non-negative; (2) each row of \mathbf{A} sums to 1; (3)
 52 $\mathbf{U}^{-1}\mathbf{U} = \mathbf{I}$; and (4) all remaining eigenvalues have magnitudes not exceeding λ_2 . The optimization
 53 problem has the following form, where we translate the constraints above into penalty terms.

$$\min_{\lambda_{3:M}, \mathbf{V}_2} \sum_{i,j=1}^M \max\{-a_{ij}, 0\} + \sum_{j=1}^M \left(\left(\sum_{i=1}^M a_{ij} \right) - 1 \right)^2 + \sum_{i,j=1}^M (\mathbf{V}\mathbf{U} - \mathbf{I})_{ij}^2 + \sum_{i=3}^M \max\{\lambda_i - \lambda_2, 0\}$$

$$\text{s.t. } \mathbf{A} = \mathbf{V} \text{diag}(1, \lambda_2, \lambda_3, \dots, \lambda_M) \mathbf{U}, \quad \mathbf{V} = [\mathbf{1} \quad \mathbf{V}_2], \quad \mathbf{U} = \begin{bmatrix} \mu \\ \mathbf{V}_2^\dagger \end{bmatrix}$$

54 This is a nonconvex problem, which we solve using first order methods with pytorch. We randomly
 55 initialize the free variables $\lambda_3, \dots, \lambda_M$ and \mathbf{V}_2 and then run 5000 iterations of Adam with step size
 56 0.01 and default values for other parameters. After the optimizer terminates, we reject instances
 57 which do not satisfy the constraints exactly. By initializing with multiple random seeds, we generate
 58 matrices spanning the desired entropy spectrum.

$$\mathbf{A} = \mathbf{U}^\dagger \mathbf{A} \mathbf{U} = \begin{bmatrix} \vec{1} & \dots \end{bmatrix} \begin{bmatrix} 1 & 0 \\ \lambda_2 & \ddots \\ 0 & \dots \end{bmatrix} \begin{bmatrix} \mu \\ \vdots \end{bmatrix}$$

Figure 1: The singular value decomposition of ergodic unichain Markov matrix \mathbf{A} . The darker shaded region is pre-defined for our controlled experiments. The lighter shaded region is randomly initialized and calculated using a neural network.

59 **Steady state distribution.** We construct steady state distributions with varying skewness using the
 60 Beta distribution with $\alpha = 1$ and different values of β . When $\alpha = 1$ and $\beta = 1$, the resulting steady
 61 state distribution is uniform. As β increases, the distribution becomes increasingly skewed toward
 62 smaller state indices. Unless otherwise specified (Appendix D.3), we use a uniform steady state
 63 distribution as the default configuration.

64 **Entropy for visualizations.** The entropy definitions $H(\mathbf{A})$ and $H(\mathbf{B}, \mu)$ we introduced in Section
 65 2.1 are used for constructing HMM parameters and sampling trajectories. For graphing main paper
 66 Figure 4 (Left), we define normalized entropy considering both matrices:

$$\tilde{H}(\mathbf{A}, \mathbf{B}, \mu) = \frac{H(\mathbf{A}) + H(\mathbf{B}, \mu)}{\log M + \log L}.$$

67 We define $\tilde{H}(\mathbf{A}) = H(\mathbf{A})/\log M$ and $\tilde{H}(\mathbf{B}, \mu) = H(\mathbf{B}, \mu)/\log L$ for the main paper Figure 4
 68 (Middle) and (Right).

69 **T is when LLM converges to Viterbi.** The concept of “convergence”, though intuitive to human eyes,
 70 requires a specific numerical definition for plots like the main paper Figure 4. We define convergence
 71 as the point where two conditions are simultaneously satisfied: (i) the accuracy difference between
 72 Viterbi and LLM is within 0.025, and (ii) LLM achieves at least 95% of Viterbi’s accuracy. We
 73 use both constant and relative thresholds to ensure a strict convergence definition that accounts for
 74 different baseline performance levels across experimental conditions.

75 **Hellinger Distance.** For two discrete probability distributions $P, Q \in \mathbb{R}^L$, the Hellinger distance is
 76 defined as

$$D_{\text{Hellinger}}(P, Q) = \frac{1}{\sqrt{2}} \sqrt{\sum_{i=1}^L (\sqrt{P_i} - \sqrt{Q_i})^2}.$$

77 C Details of Benchmark Models

78 In this section, we provide descriptions and pseudocode for the benchmark models we use (main
79 paper Table 1). The executable code for all methods are included in supplemental materials.

Algorithm 1: Viterbi Algorithm

Input: States $\mathcal{X} = \{1, 2, \dots, M\}$, initial distribution μ , transition matrix \mathbf{A} , emission matrix \mathbf{B} ,
and observation sequence $\{o_1, \dots, o_T\}$.
Output: Most likely state sequence path $= \{x_1, \dots, x_T\}$
Initialization: $\mathbf{P}[0][s] \leftarrow \mu[s] \cdot \mathbf{B}[s][o_1]$ for all $s \in \mathcal{X}$;
Forward recursion: for $t = 1$ to $T - 1$ do
 for $s \in \mathcal{X}$ do
 $\mathbf{P}[t][s] \leftarrow \max_{r \in \mathcal{X}} \{\mathbf{P}[t-1][r] \cdot \mathbf{A}[r][s] \cdot \mathbf{B}[s][o_t]\}$;
 $\mathbf{Q}[t][s] \leftarrow \arg \max_{r \in \mathcal{X}} \{\mathbf{P}[t-1][r] \cdot \mathbf{A}[r][s] \cdot \mathbf{B}[s][o_t]\}$;
 end
end
Backtracking: path[$T - 1$] $\leftarrow \arg \max_{s \in \mathcal{X}} \mathbf{P}[T - 1][s]$;
path[t] $\leftarrow \mathbf{Q}[t + 1][\text{path}[t + 1]]$ for $t = T - 2, \dots, 0$;
return path

80 **Viterbi algorithm.** The Viterbi algorithm is a dynamic programming technique for efficiently finding
81 the most likely sequence of hidden states in a Markov model, given a sequence of observations. It
82 iteratively computes the highest probability path to each state at time t by considering all possible
83 predecessor states at time $t - 1$, their transition probabilities, and the emission probabilities of the
84 current observation. Rather than exhaustively evaluating all M^T possible state sequences, Viterbi
85 maintains only the M most promising paths at each time step, storing both their probabilities and the
86 penultimate states that maximize these probabilities. After computing probabilities for all time steps,
87 the algorithm traces backward from the most probable final state to reconstruct the optimal state
88 sequence. We use the most probable final state and ground-truth \mathbf{A} and \mathbf{B} to calculate the prediction
89 distribution of the next observation.

Algorithm 2: Compute $P(O_{t+1}|O_{t-k:t})$

Input: States $\mathcal{X} = \{1, 2, \dots, M\}$, initial distribution μ , transition matrix \mathbf{A} , emission matrix \mathbf{B} ,
and observation sequence $\{o_{t-k}, \dots, o_t\}$.
Output: Probability of next observation $P(o_{t+1}|o_{t-k:t})$
Forward pass over observation window: $\alpha_{t-k}[s] \leftarrow \mathbf{B}[s][o_{t-k}] \cdot \mu[s]$ for all $s \in \mathcal{X}$;
 $\alpha_i[s] \leftarrow \mathbf{B}[s][o_i] \cdot \sum_{r \in \mathcal{X}} \mathbf{A}[r][s] \cdot \alpha_{i-1}[r]$ for $i = t - k + 1, \dots, t, s \in \mathcal{X}$;
Normalize to get posterior: $P(s|o_{t-k:t}) \leftarrow \frac{\alpha_t[s]}{\sum_{s' \in \mathcal{X}} \alpha_t[s']}$ for all $s \in \mathcal{X}$;
Prediction step: $P(s|o_{t-k:t}) \leftarrow \sum_{r \in \mathcal{X}} \mathbf{A}[r][s] \cdot P(r|o_{t-k:t})$ for all $s \in \mathcal{X}$;
Marginalize over states: $P(o_{t+1}|o_{t-k:t}) \leftarrow \sum_{s \in \mathcal{X}} \mathbf{B}[s][o_{t+1}] \cdot P(s|o_{t-k:t})$;
return $P(o_{t+1}|o_{t-k:t})$

90 **Optimal inference with truncated memory** $P(O_{t+1}|O_{t-k:t})$. The forward-based prediction al-
91 gorithm computes the probability of the next observation in a hidden Markov model by using a
92 three-step approach. First, it calculates the posterior distribution over current hidden states via the
93 forward algorithm, recursively processing the observation window while accounting for transitions
94 and emissions. Second, it projects this belief state forward by applying the transition matrix to com-
95 pute the distribution over next possible states. Finally, it determines $P(o_{t+1}|o_{t-k:t})$ by marginalizing
96 over all possible next states, weighting each by its emission probability.

97 **Baum-Welch algorithm.** The Baum-Welch algorithm is an expectation-maximization method for
98 estimating hidden Markov model parameters. It iteratively alternates between computing state
99 posteriors $\gamma_t(s)$ and transition posteriors $\xi_t(s, r)$ via forward-backward recursion (E-step), and
100 updating parameters to maximize likelihood (M-step): setting the initial distribution to γ_1 , transition

Algorithm 3: Baum-Welch Algorithm

Input: States $\mathcal{X} = \{1, 2, \dots, M\}$, observations $\mathcal{Y} = \{1, 2, \dots, L\}$, observation sequence $\{o_1, \dots, o_T\}$, initial parameters $\mu^{(0)}, \mathbf{A}^{(0)}, \mathbf{B}^{(0)}$, and threshold ϵ .

Output: Refined parameters $\mu, \mathbf{A}, \mathbf{B}$

Initialize: $\mu \leftarrow \mu^{(0)}, \mathbf{A} \leftarrow \mathbf{A}^{(0)}, \mathbf{B} \leftarrow \mathbf{B}^{(0)}, \mathcal{L}_{\text{prev}} \leftarrow -\infty$;

repeat

$\mathcal{L}_{\text{prev}} \leftarrow \mathcal{L}$;

E-Step:

Forward pass: $\alpha_1[s] \leftarrow \mathbf{B}[s][o_1] \cdot \mu[s]$ for all $s \in \mathcal{X}$;

$\alpha_t[s] \leftarrow \mathbf{B}[s][o_t] \cdot \sum_r \mathbf{A}[r][s] \cdot \alpha_{t-1}[r]$ for $t = 2, \dots, T, s \in \mathcal{X}$;

$\mathcal{L} \leftarrow \sum_s \alpha_T[s]$;

Backward pass: $\beta_T[s] \leftarrow 1$ for all $s \in \mathcal{X}$;

$\beta_t[s] \leftarrow \sum_r \mathbf{A}[s][r] \cdot \mathbf{B}[r][o_{t+1}] \cdot \beta_{t+1}[r]$ for $t = T-1, \dots, 1, s \in \mathcal{X}$;

Expected counts: $\gamma_t[s] \leftarrow \frac{\alpha_t[s] \cdot \beta_t[s]}{\mathcal{L}}$ for $t = 1, \dots, T, s \in \mathcal{X}$;

$\xi_t[s][r] \leftarrow \frac{\alpha_t[s] \cdot \mathbf{A}[s][r] \cdot \mathbf{B}[r][o_{t+1}] \cdot \beta_{t+1}[r]}{\mathcal{L}}$ for $t = 1, \dots, T-1, s, r \in \mathcal{X}$;

M-Step: $\mu[s] \leftarrow \gamma_1[s]$ for all $s \in \mathcal{X}$;

$\mathbf{A}[s][r] \leftarrow \frac{\sum_{t=1}^{T-1} \xi_t[s][r]}{\sum_{t=1}^{T-1} \gamma_t[s]}$ for all $s, r \in \mathcal{X}$;

$\mathbf{B}[s][v] \leftarrow \frac{\sum_{t=1}^T \gamma_t[s] \cdot \mathbf{1}(o_t=v)}{\sum_{t=1}^T \gamma_t[s]}$ for all $s \in \mathcal{X}, v \in \mathcal{Y}$;

until $|\mathcal{L} - \mathcal{L}_{\text{prev}}| < \epsilon$;

return $\mu, \mathbf{A}, \mathbf{B}$

101 probabilities to normalized expected transitions, and emission probabilities to normalized observation
102 counts per state. This process continues until the log-likelihood converges, yielding locally optimal
103 parameters that maximize the probability of generating the observed sequence. We use the learned
104 parameters to predict next observation similar to the Viterbi algorithm.

Algorithm 4: n -gram Based Next-Observation Prediction

Input: Observation sequence $O = \{o_1, \dots, o_T\}$, context length $n-1$, smoothing parameter δ

Output: n -gram model for predicting $P(o_t | o_{t-(n-1):t-1})$

Count extraction: $\text{counts}_n, \text{counts}_{n-1} \leftarrow$ empty associative arrays;

for $t = n-1$ **to** $T-1$ **do**

 context $\leftarrow [o_{t-(n-1)}, \dots, o_{t-1}]$;

 Increment $\text{counts}_n[\text{context} \oplus o_t]$ and $\text{counts}_{n-1}[\text{context}]$;

end

Model construction with smoothing: $V \leftarrow$ number of unique symbols in O ;

for each observed context c in counts_{n-1} **do**

for each unique observation o in O **do**

 model[c, o] $\leftarrow \frac{\text{counts}_n[c \oplus o] + \delta}{V \cdot \delta + \text{counts}_{n-1}[c]}$;

end

end

Back-off for unseen contexts: $P_{\text{unif}}(o) \leftarrow \frac{1}{V}$ for all o ;

$P(o_t | o_{t-(n-1):t-1}) \leftarrow \begin{cases} \text{model}[[o_{t-(n-1)}, \dots, o_{t-1}], o_t], & \text{if context observed;} \\ P_{\text{unif}}(o_t), & \text{otherwise} \end{cases}$;

return model

105 **n -gram.** n -gram models provide an elegant, computationally efficient framework for next-observation
106 prediction in Markov chain processes by directly estimating conditional probabilities from observed
107 sequences. These models embody the Markov assumption that $P(O_{t+1} | O_{1:t}) \approx P(O_{t+1} | O_{t-n+2:t})$,
108 making them particularly effective for stochastic processes where future states depend only on a
109 limited history of previous states. For first-order Markov chains, bigram models ($n = 2$) precisely
110 capture the underlying transition dynamics, while higher-order dependencies can be modeled by
111 increasing n .

Algorithm 5: LSTM for Single Sequence Prediction

Input: Observation sequence $O = \{o_1, \dots, o_T\}$, vocabulary size V , embedding dimension d , hidden dimension h , layers L , learning rate α , epochs E
Output: Trained LSTM model for $P(o_{t+1}|o_{1:t})$
Architecture: Initialize embedding layer: $\text{Embedding} : \mathbb{Z} \rightarrow \mathbb{R}^d$;
Initialize LSTM layers: $\text{LSTM} : \mathbb{R}^d \rightarrow \mathbb{R}^h$ with L layers;
Initialize output projection: $\text{Linear} : \mathbb{R}^h \rightarrow \mathbb{R}^V$;
Initialize optimizer with learning rate α ;
Training: for $epoch = 1$ to E do
 for $t = 1$ to $T - 1$ do
 $x \leftarrow O_{1:t}$; // Use all previous observations as context
 $y \leftarrow O_{t+1}$; // Next observation as target
 Update model to maximize $P(y|x)$ via gradient descent;
 end
end
Inference: For prefix $O_{1:t}$, compute $P(o_{t+1}|O_{1:t}) = \text{softmax}(\text{model}(O_{1:t}))$;
return *model*

112 **RNN LSTM.** LSTM networks are specialized recurrent neural architectures designed to model
113 sequential data through memory cells regulated by input, forget, and output gates. These gates
114 control information flow, allowing LSTMs to selectively retain relevant historical patterns while
115 discarding irrelevant information. LSTMs excel at next-observation prediction tasks by capturing
116 both short-term correlations and long-term dependencies in the observation history. The network
117 processes a window of prior observations sequentially, updating its hidden state to encode temporal
118 patterns, then projects this state through a softmax layer to generate a probability distribution over
119 possible next observations—making LSTMs particularly effective for forecasting future values in
120 time series where the prediction depends on complex patterns spanning multiple time scales.

121 D Additional Synthetic Experiment Results

122 This section presents additional results from synthetic experiments. All methods are evaluated using
 123 the average performance over 4,096 sequences, with the exception of LSTM, which is evaluated on
 124 16 sequences due to its high computational cost. Consequently, the LSTM results exhibit higher
 125 variance. Nonetheless, in metrics such as Hellinger distance—which account for the full output
 126 distribution rather than relying solely on the argmax for accuracy—LSTM underperforms compared
 127 to the LLM most of the time.

128 D.1 Varying Entropy of A

129 In this section, we present detailed results on varying the entropy of A matrix over 4/8/16 states and
 emissions, reporting accuracies and Hellinger distances.

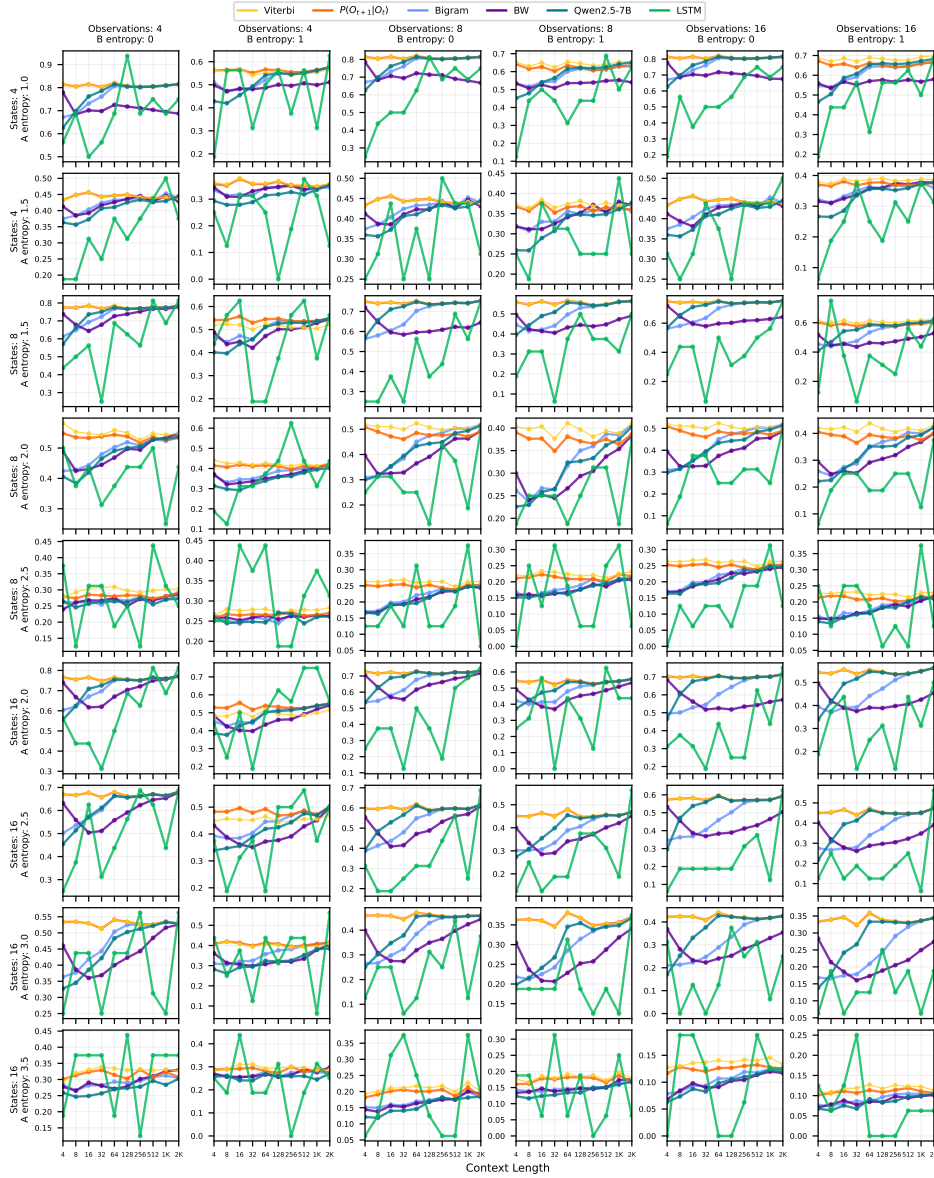


Figure 2: Accuracies of six methods across different A entropy, B entropy, number of states, and number of emissions with $\lambda_2 = 0.75$ and uniform steady state distribution.

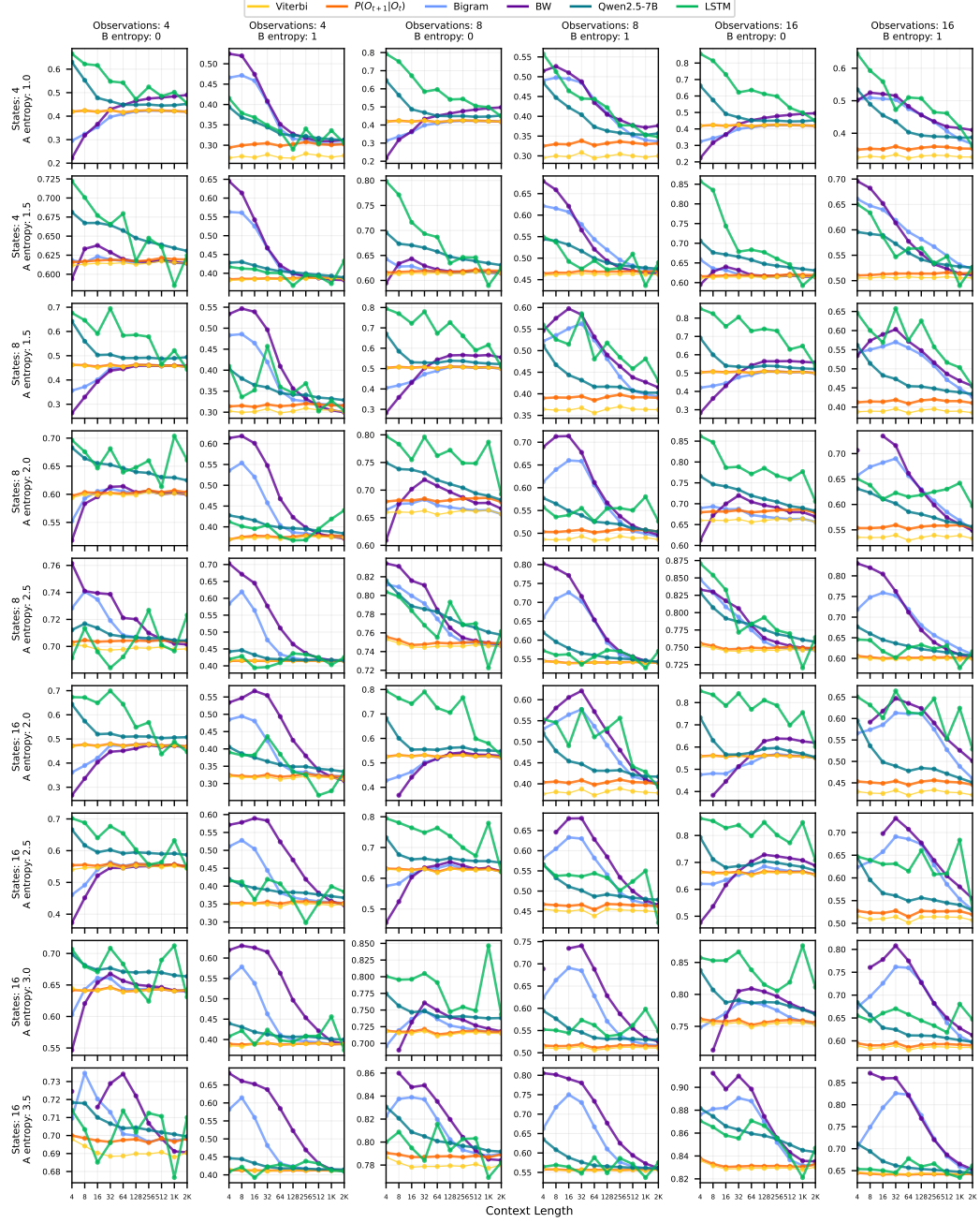


Figure 3: Hellinger distances of six methods across different A entropy, B entropy, number of states, and number of emissions with $\lambda_2 = 0.75$ and uniform steady state distribution.

131 D.2 Varying Mixing Rate of A

132 In this section, we present detailed results on varying the mixing rate (λ_2) of A matrix over 4/8/16
states and emissions, reporting accuracies and Hellinger distances.

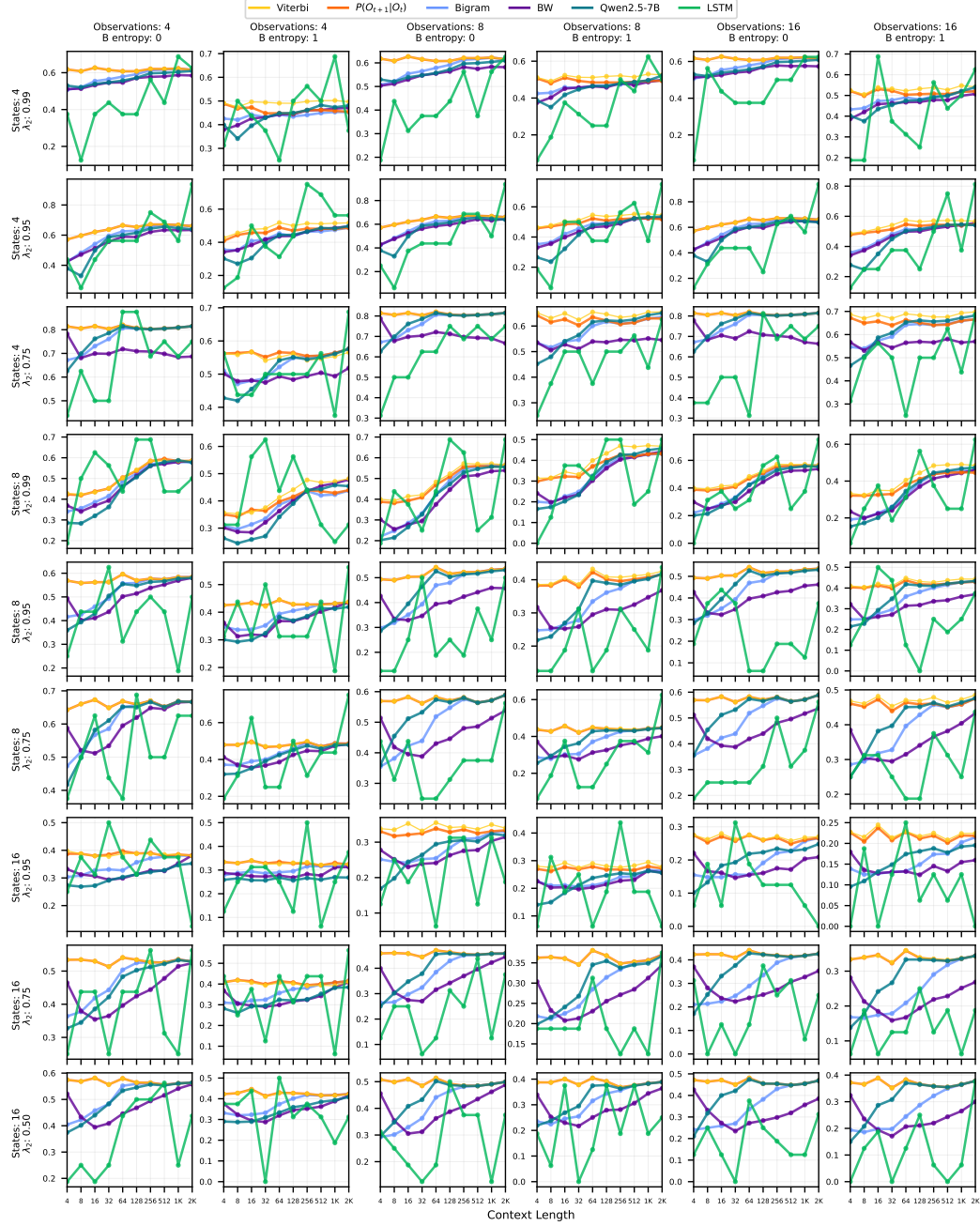


Figure 4: Accuracies of six methods across different mixing rates (λ_2), B entropy, number of states, and number of emissions with uniform steady state distribution and (1, 2, 3) A entropy for (4, 8, 16) states respectively.

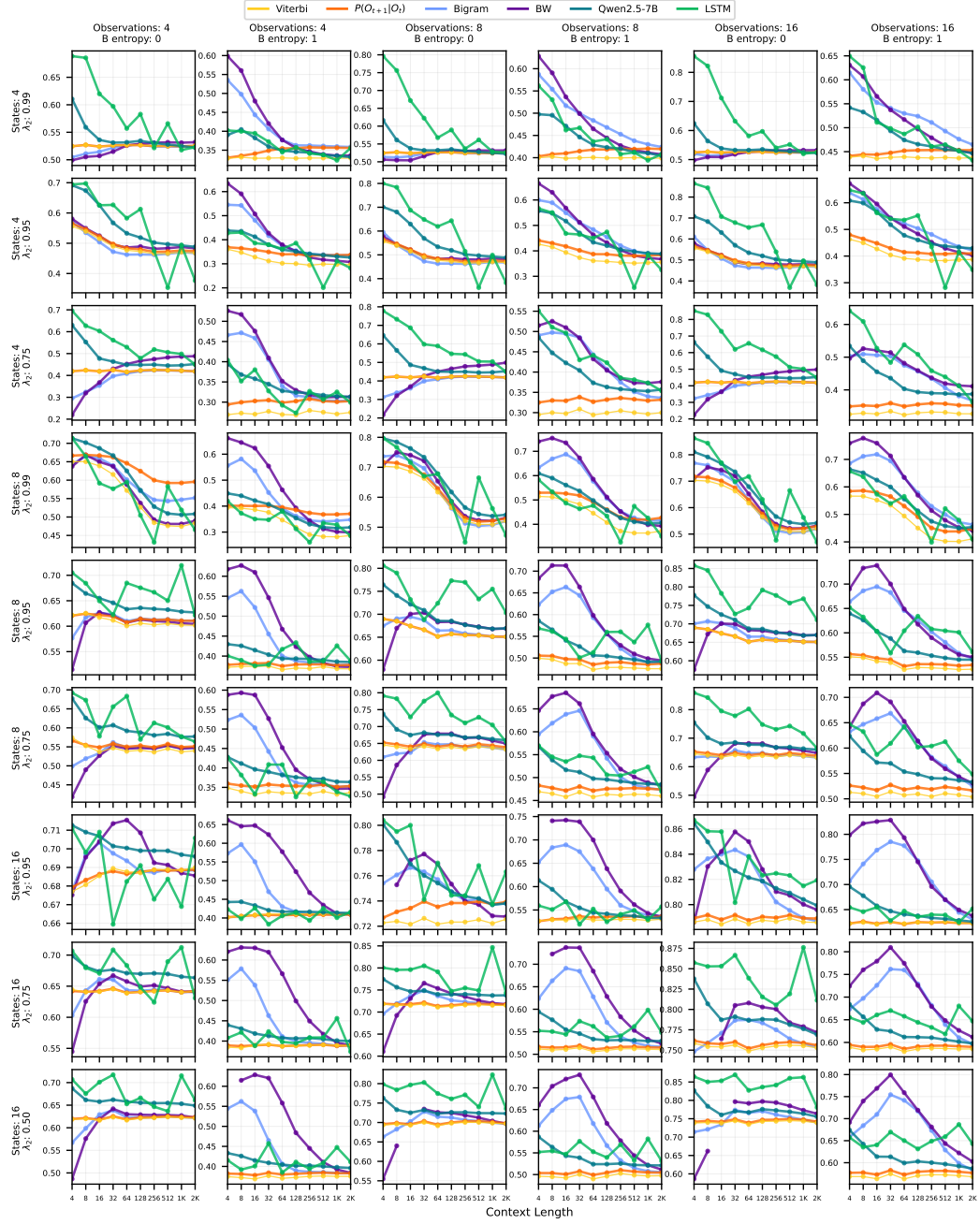


Figure 5: Hellinger distances of six methods across different mixing rates (λ_2), B entropy, number of states, and number of emissions with uniform steady state distribution and (1, 2, 3) A entropy for (4, 8, 16) states respectively.

134 D.3 Varying Steady State Distribution of A

135 In this section, we present detailed results on varying the steady state distributions of A matrix over
 136 4/8/16 states and emissions, reporting accuracies and Hellinger distances. We construct steady states
 137 with different skewness using Beta distribution with $\alpha = 1$. Notably, with $\alpha = 1$ and $\beta = 1$, the
 138 steady state distribution is uniform. As we increase β , the distribution becomes more skewed. We
 139 test with $\beta = 1, 2, 3$, representing uniform, skewed, and very skewed respectively.

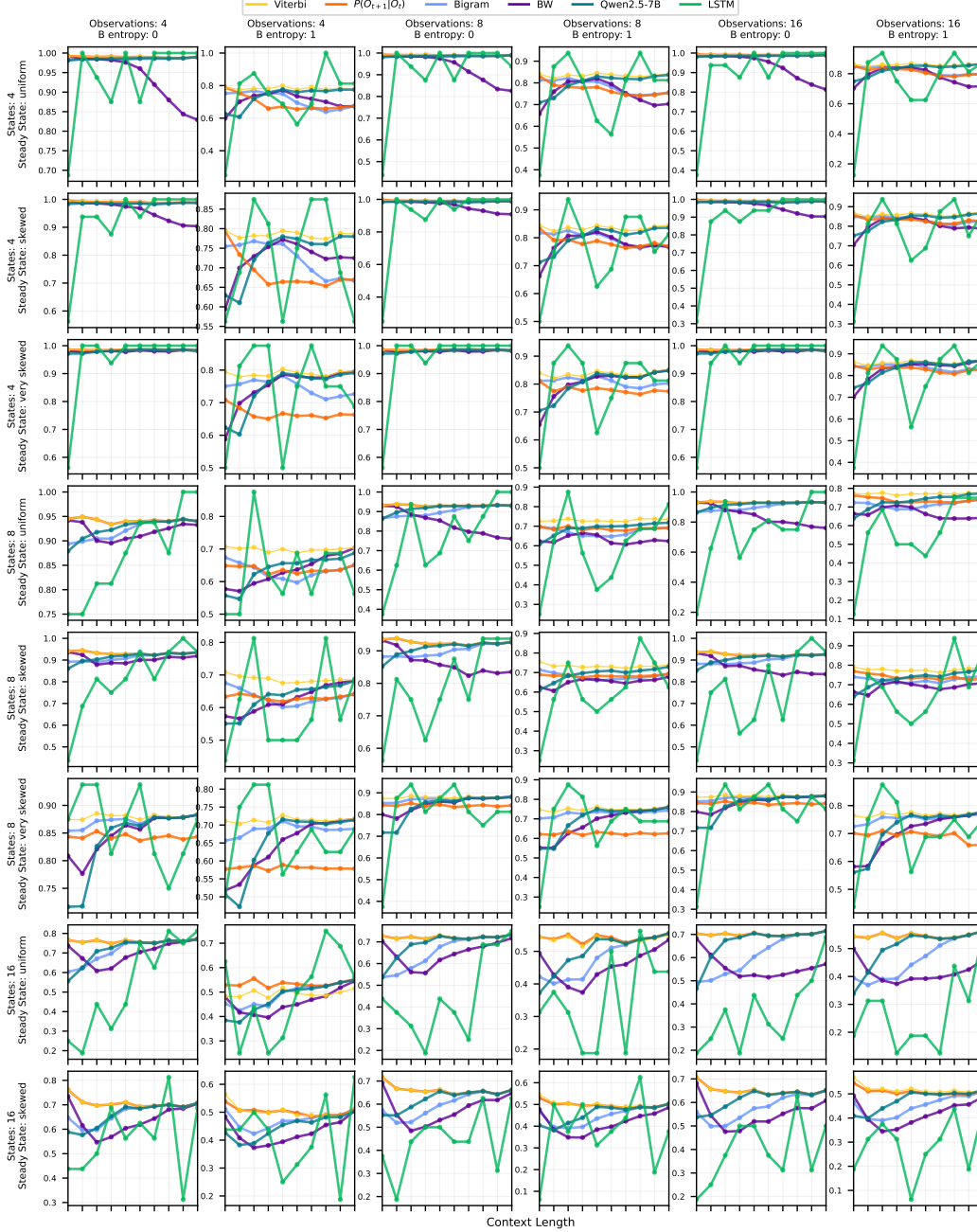


Figure 6: Accuracies of six methods across different steady state distributions, B entropy, number of states, and number of emissions with (0, 0.5, 2) A entropy for (4, 8, 16) states respectively and (0.99, 0.95, 0.75) λ_2 for (4, 8, 16) states respectively.

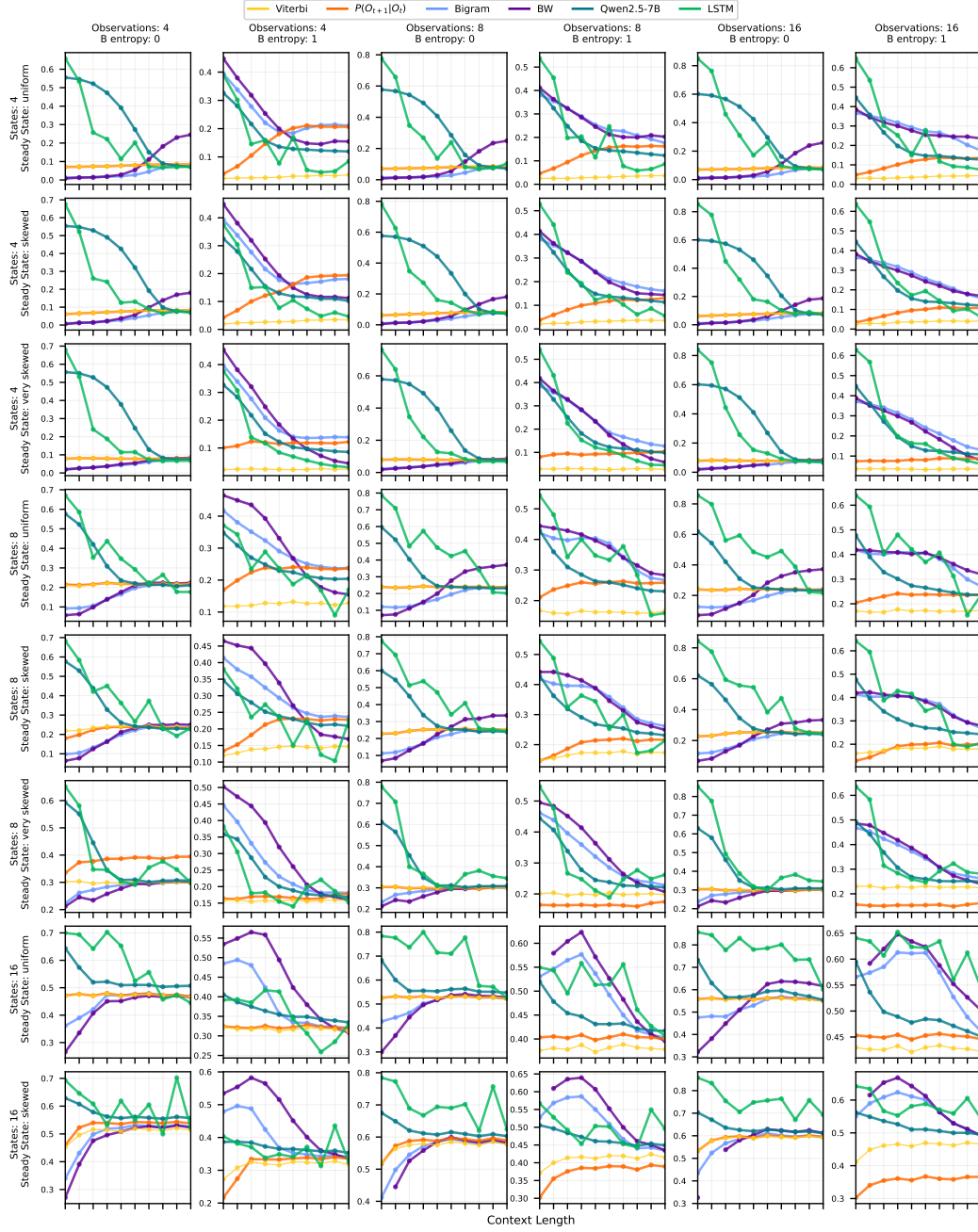


Figure 7: Hellinger distances of six methods across different steady state distributions, B entropy, number of states, and number of emissions with (0, 0.5, 2) A entropy for (4, 8, 16) states respectively and (0.99, 0.95, 0.75) λ_2 for (4, 8, 16) states respectively.

D.4 Discussions

When LLMs fail to converge. While LLMs converge to Viterbi performance efficiently under most HMM parameter settings (scaling trends summarized in Section 3.1), we identify two conditions where convergence fails or proceeds exceptionally slowly. First, when entropy of \mathbf{A} or \mathbf{B} is approaches its maximum ($\log M$ or $\log L$ respectively), the prediction accuracy gap ε at context length 2048 remains substantial. For instance, in the last row of Figure 2 with $M = 16$ and the entropy of \mathbf{A} is 3.5 (near the maximum of $\log 16 = 4$), the LLM (Qwen2.5-7B) exhibits gradual convergence with a persistent gap. Second, when mixing is slow (λ_2 approaches 1), such as in the third-to-last row of Figure 4 with $M = 16$ and $\lambda_2 = 0.95$, a performance gap persists even at maximum context length.

Importantly, the Viterbi algorithm also struggles under these challenging conditions. Under high entropy, as shown in the last row of Figure 2, Viterbi accuracy barely exceeds random prediction (0.25/0.125/0.0625 for $L = 4/8/16$). Under slow mixing, such as in the fourth row of Figure 4 with $M = 8$ and $\lambda_2 = 0.99$, Viterbi algorithm requires context length 512 to achieve peak performance. These results demonstrate that LLM performance degradation under high entropy and slow mixing conditions reflects fundamental limits of stochastic system learnability—arising from random dynamics and long-range dependencies—that affect even optimal inference methods.

Monotonicity of LLM performance with respect to context length. We observe that LLM performance almost always improves monotonically with longer context length—a property notably absent in other learning baselines. Even excluding LSTM from this comparison (due to high variance from averaging over fewer sequences, as discussed in the first paragraph in Appendix D), both Baum-Welch and bigram models lack monotonic convergence behavior. For Baum-Welch, the accuracy graphs (Figures 2, 4, and 6) reveal multiple cases where performance “dips” and recovers, or deteriorates as context length increases. The Hellinger distance graphs (Figures 3, 5, and 7) provide clearer evidence that both BW and bigram exhibit non-monotonic learning patterns. In most cases, LLM Hellinger distance decreases monotonically, while BW and bigram display erratic behavior: sometimes experiencing early-context “bumps”, other times starting very close to the ground truth emission distribution (occasionally even closer than the oracle Viterbi by empirical chance) before gradually converging to statistically sound distributions. Importantly, when BW or bigram achieve lower Hellinger distances, this does not necessarily indicate better performance—the corresponding prediction accuracy graphs often show poor results, highlighting the distinction between distributional similarity and predictive capability.

When (normalized) entropy \tilde{H} is held constant, varying the number of states does not affect the LLM convergence rate. We provide concrete evidence for this claim in Figure 2, where rows 1, 3, and 6 all have the same normalized entropy $\tilde{H}(\mathbf{A}) = 0.5$. Across each column, the Qwen2.5-7B convergence curves for these three rows exhibit nearly identical shapes, demonstrating that the convergence rate depends primarily on normalized entropy rather than absolute state space size.

We emphasize that convergence rate differs from convergence target—the Viterbi performance. While the rate of improvement remains consistent across different state space sizes (when normalized entropy is fixed), larger state spaces result in lower achievable prediction accuracy due to increased task difficulty.

180 E Ablations on LLMs

181 In this section, we provide the results on the families, sizes, and tokenization of the LLMs.

182 E.1 LLM Size

183 We compare Qwen and Llama model families with seven different models. We found that their performances are similar, with slight degradation when the model size is small.

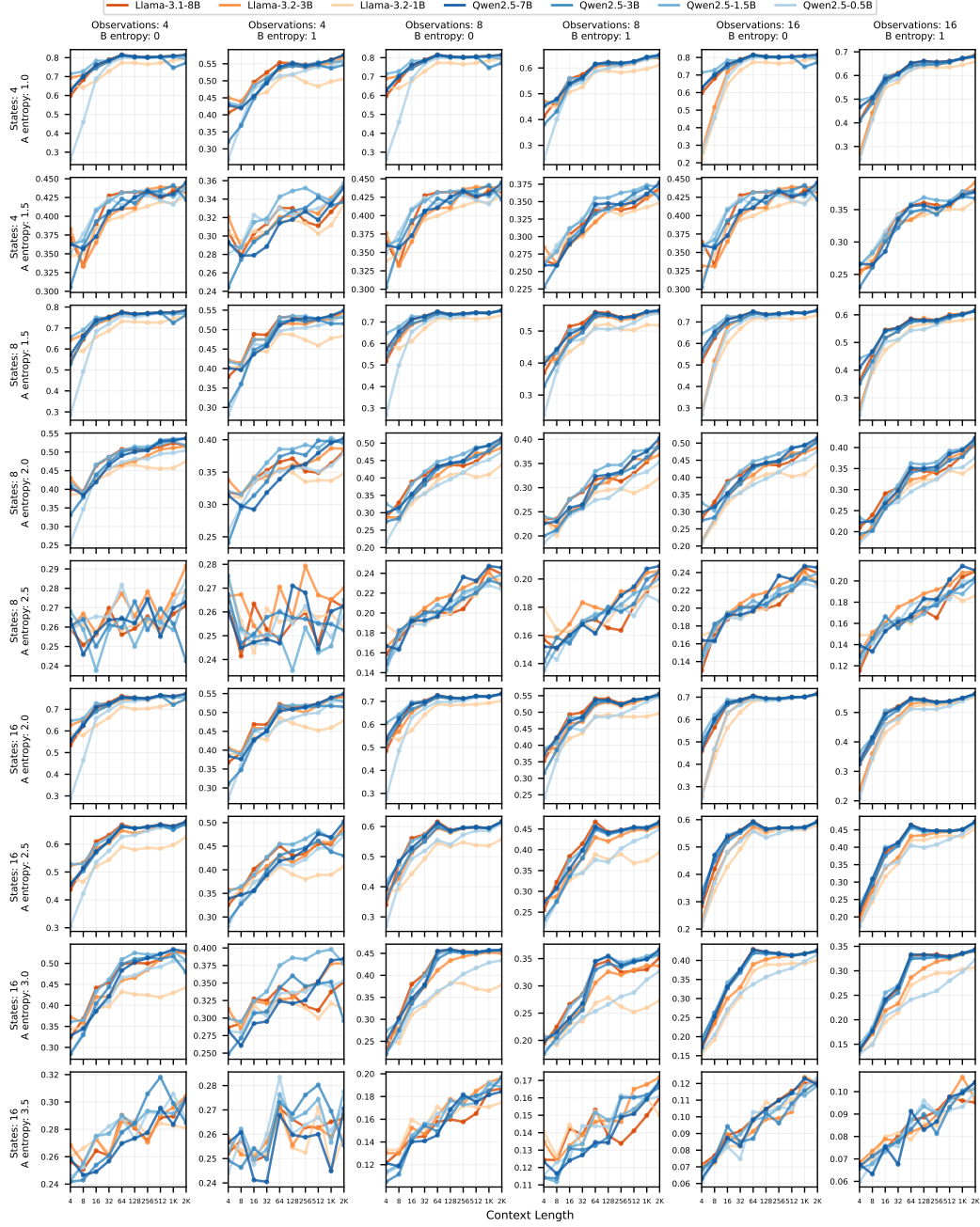


Figure 8: Accuracies of seven models across different **A** entropy, **B** entropy, number of states, and number of emissions with $\lambda_2 = 0.75$ and uniform steady state distribution. Lighter color represents smaller models. The two smallest models from each family have suboptimal performance.

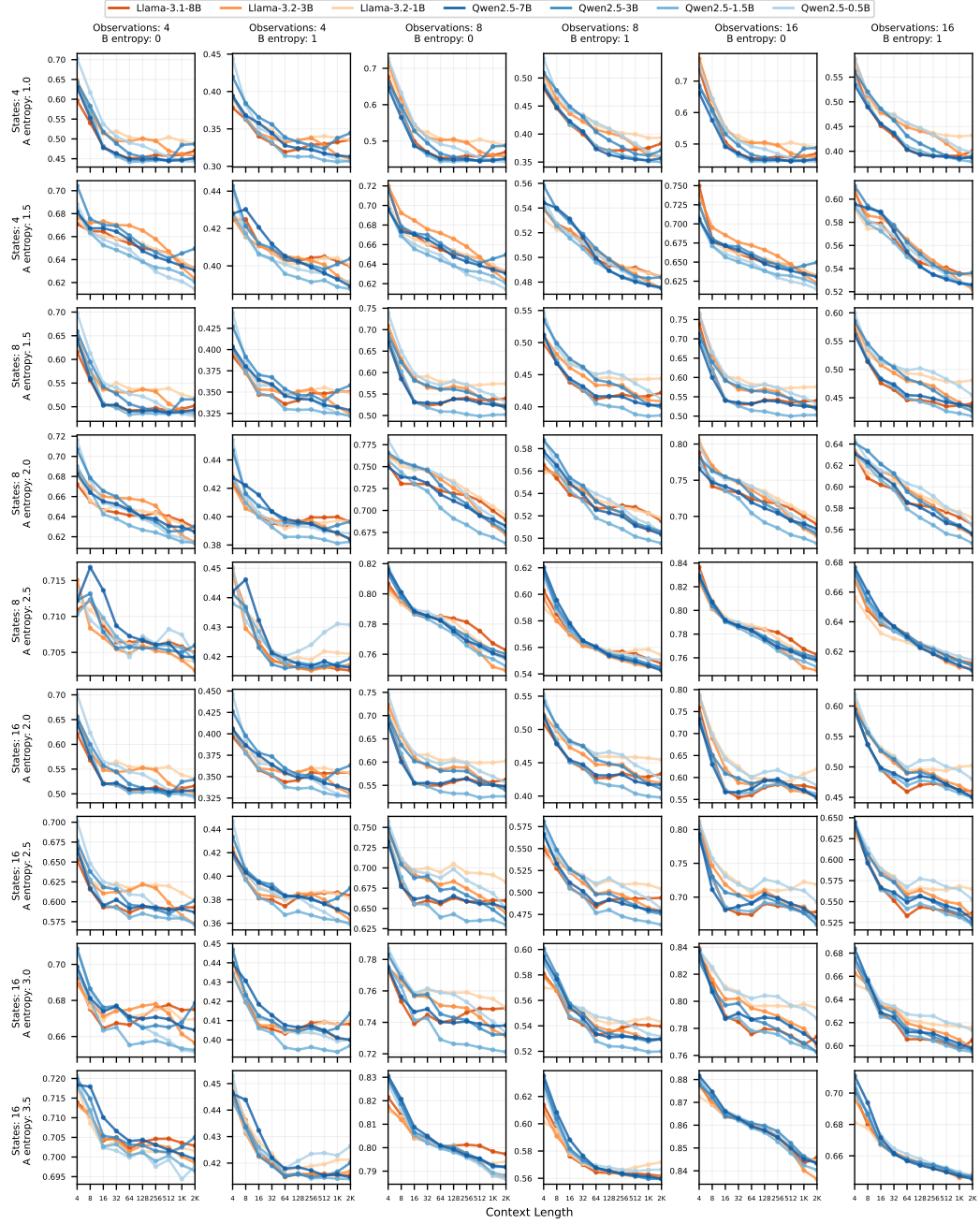


Figure 9: Hellinger distances of seven models across different A entropy, B entropy, number of states, and number of emissions with $\lambda_2 = 0.75$ and uniform steady state distribution. The models converge similarly, especially when entropy is high.

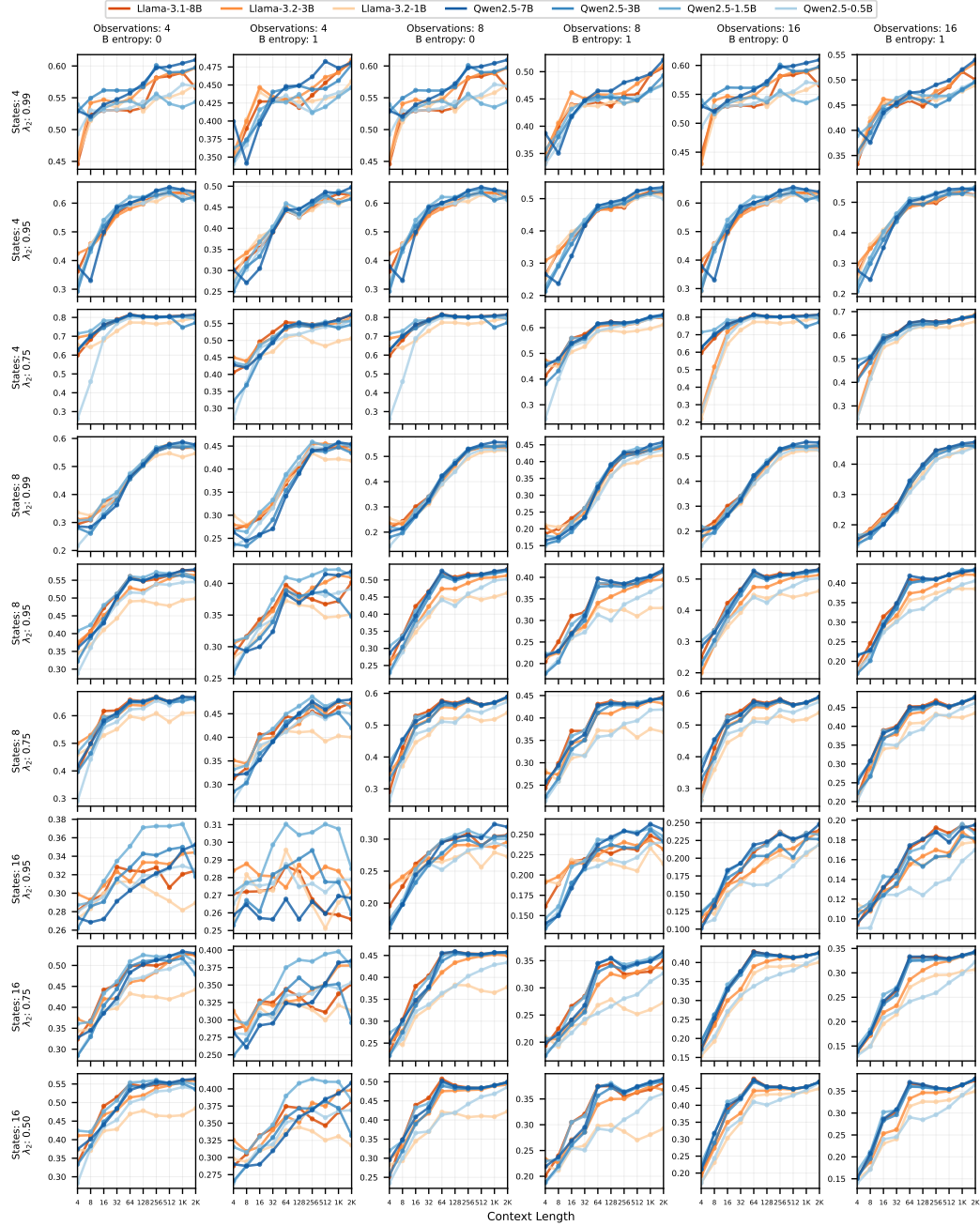


Figure 10: Accuracies of seven models across different mixing rates (λ_2), B entropy, number of states, and number of emissions with uniform steady state distribution and (1, 2, 3) A entropy for (4, 8, 16) states respectively. The two smallest models from each family have suboptimal performance, especially when mixing is fast.

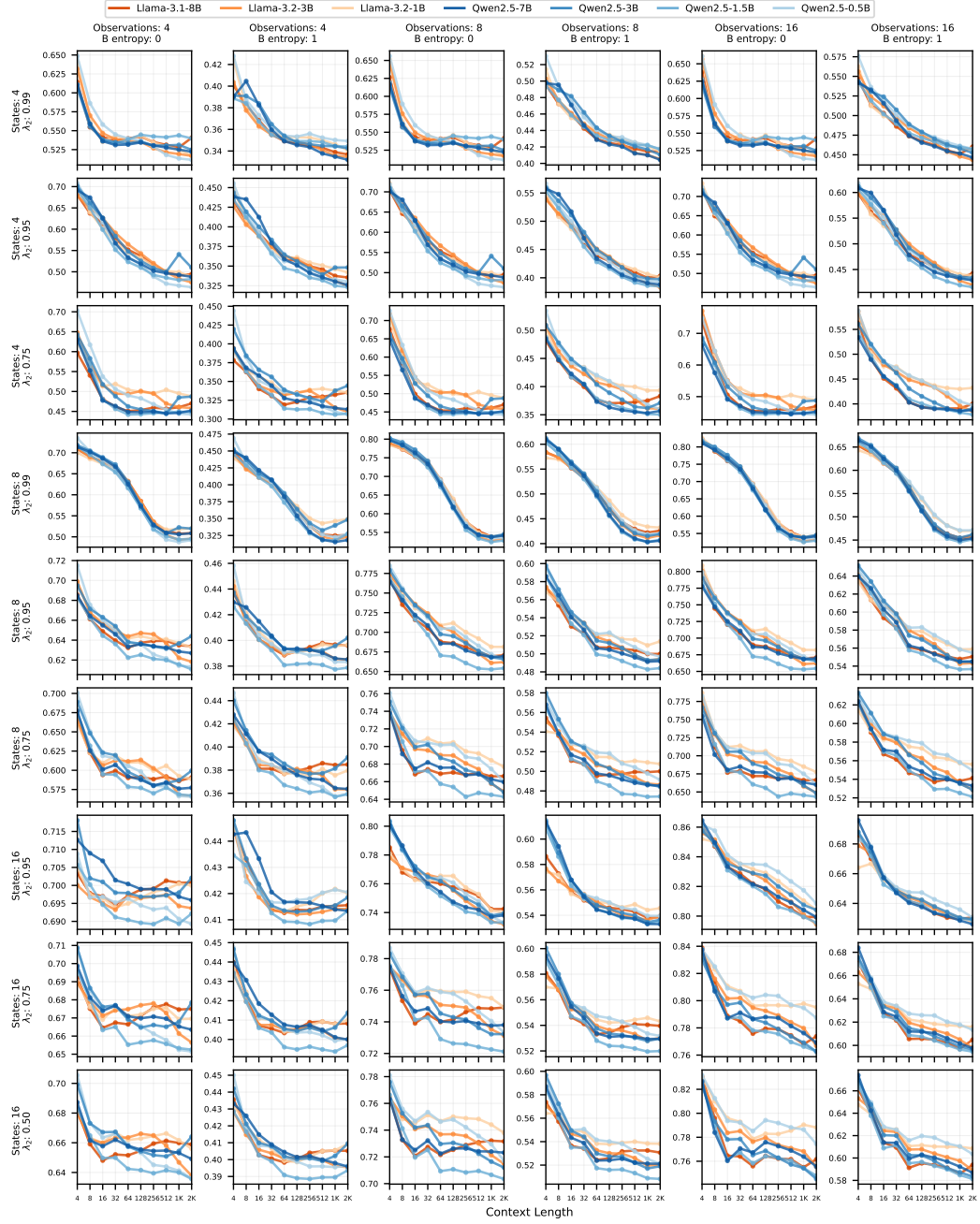


Figure 11: Hellinger distances of seven models across different mixing rates (λ_2), B entropy, number of states, and number of emissions with uniform steady state distribution and (1, 2, 3) A entropy for (4, 8, 16) states respectively. The models converge similarly, especially when mixing is slow.

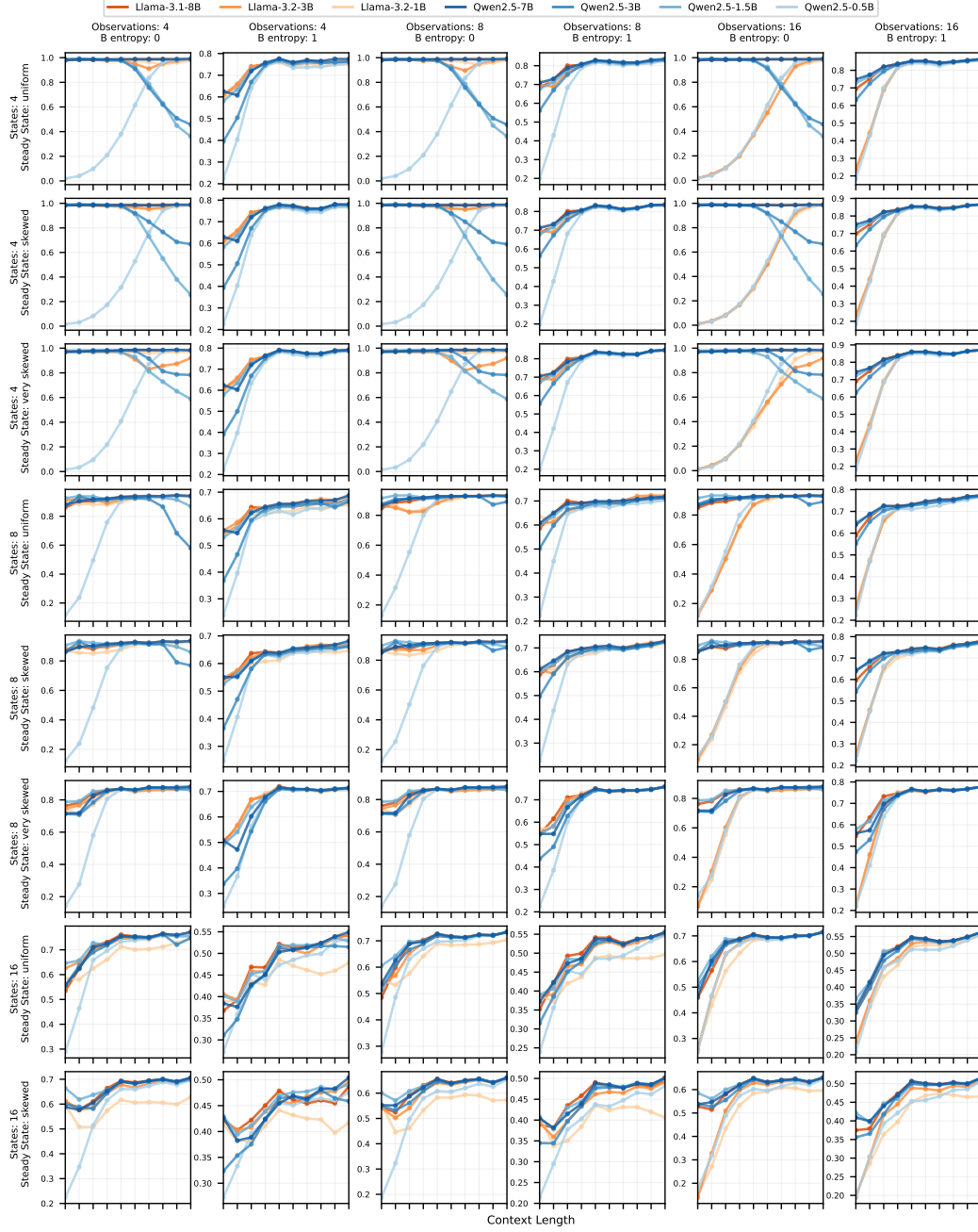


Figure 12: Accuracies of seven models across different steady state distributions, B entropy, number of states, and number of emissions with $(0, 0.5, 2)$ A entropy for $(4, 8, 16)$ states respectively and $(0.99, 0.95, 0.75)$ λ_2 for $(4, 8, 16)$ states respectively. The poor performance observed in smaller models at short context length under low A & B entropy settings may be attributed to the filtering of repeated n -grams during pretraining, as discussed in Appendix E.2.

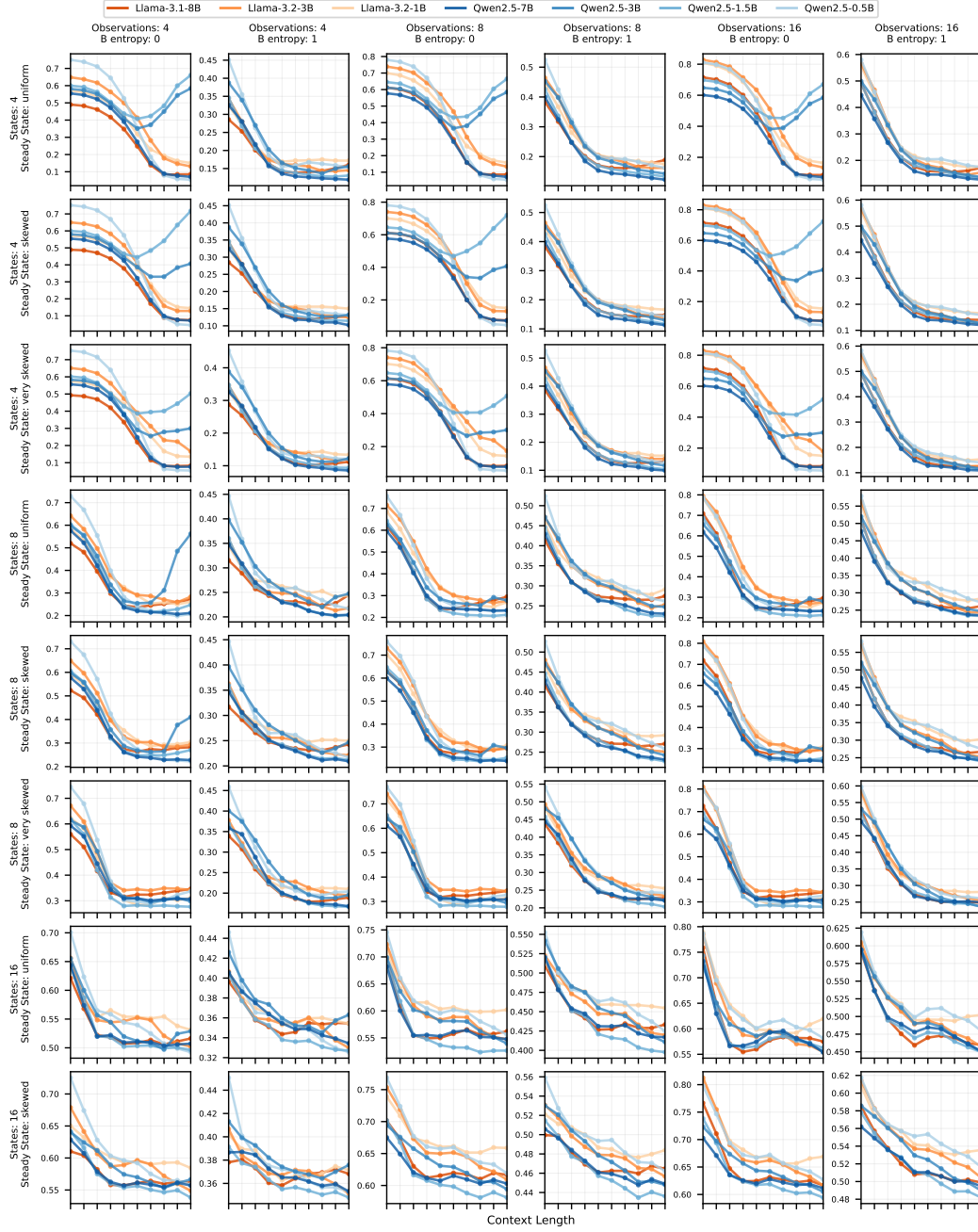


Figure 13: Hellinger distances of seven models across different steady state distributions, B entropy, number of states, and number of emissions with (0, 0.5, 2) A entropy for (4, 8, 16) states respectively and (0.99, 0.95, 0.75) λ_2 for (4, 8, 16) states respectively.

185 E.2 Tokenization

186 In this section, we evaluate three tokenization strategies: **ABC**, which encodes emissions as single
 187 letters; **123**, which encodes them as single digits; and **random**, which maps emissions to random
 188 tokens from the LLM’s tokenizer. For the **random** strategy, we specifically map emissions to special
 189 tokens (!@#\$). All experiments are conducted using the Qwen2.5-1.5B model, and the results are
 190 presented below.

191 We observe that all tokenization methods converge to similar performance levels in terms of accuracy,
 192 with **ABC** converging slightly faster when the entropy of A is large. This suggests that the choice
 193 of tokenization has limited impact on final performance. In our experiments, we adopt the **ABC**
 194 tokenization for maximum performance on the LLM. However, when the entropy of matrix A is
 195 low, **ABC** tokenization exhibits significantly lower initial accuracy and a higher Hellinger distance
 196 with short context length. We hypothesize that this is due to the increased likelihood of repetitive
 197 state sequences early in the sequence—for example, ‘AAAAA...’. During pretraining, such repeated
 198 n-gram patterns are often filtered out, as they could cause loss spikes [7]. As a result, the model may
 199 have limited exposure to these patterns, leading to poor initial performance on such inputs.

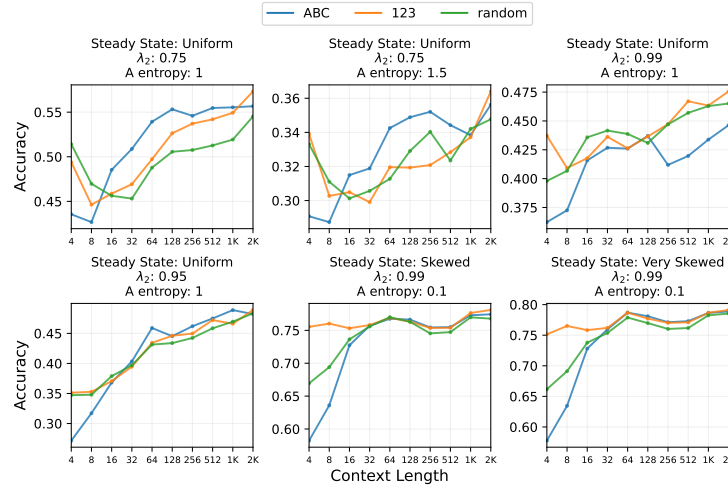


Figure 14: Accuracy of three tokenization methods across different mixing rates (λ_2), A entropy, and steady states with 4 states, 4 emissions, and 1 for B entropy.

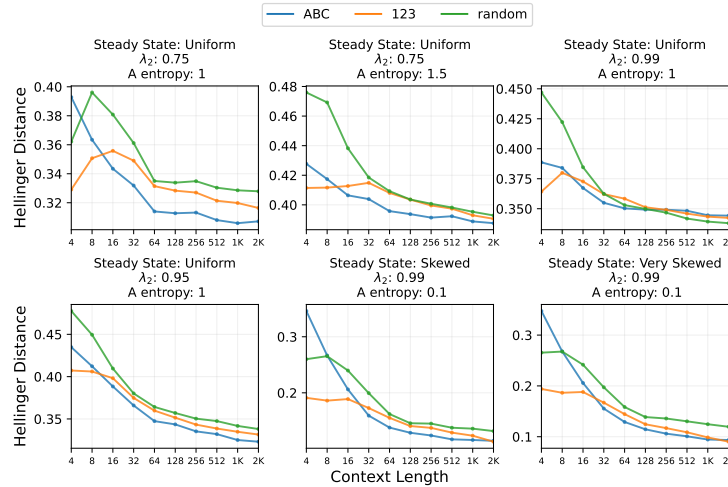


Figure 15: Hellinger distance of three tokenization methods across different mixing rates (λ_2), A entropy, and steady states with 4 states, 4 emissions, and 1 for B entropy.

F Spectral Learning HMMs for Prediction Task

Notations: We use $[\mathbf{X}]_{i,j}$ to denote the element of matrix \mathbf{X} at its i -th row and j -th column. The indicator function $\mathbf{1}_{\{x=i\}}$ is 1 only when $x = i$ and is 0 otherwise. We use $\mathbf{1}_M$ to denote a vector of all 1's with dimension M . We use the notation $[L] = \{1, 2, \dots, L\}$. $\|\cdot\|$ denotes the Frobenius norm for matrices, and depending on the context it denotes ℓ_1 or ℓ_2 norm for vectors.

Algorithm 6: Spectral Learning-Based Prediction

Input: Number of hidden states M , number of observations L , sequence $\{o_1, \dots, o_N\}$

Output: Conditional probability distribution $\hat{P}(O_{N+1}|O_{1:N} = o_{1:N})$

Estimate empirical probabilities: for all combinations $i, j, n \in [L]$ do

$$\begin{aligned} [\hat{\mathbf{P}}_1]_i &\leftarrow \frac{1}{N} \sum_{k=1}^N \mathbf{1}_{\{o_k=i\}}; \\ [\hat{\mathbf{P}}_2]_{i,j} &\leftarrow \frac{1}{N} \sum_{k=1}^N \mathbf{1}_{\{o_k=i, o_{k-1}=j\}}; \\ [\hat{\mathbf{P}}_{3,n}]_{i,j} &\leftarrow \frac{1}{N} \sum_{k=1}^N \mathbf{1}_{\{o_k=i, o_{k-1}=n, o_{k-2}=j\}}; \end{aligned}$$

end

Compute SVD for dimensionality reduction:

$\hat{\mathbf{U}} \leftarrow$ left singular vectors of $\hat{\mathbf{P}}_2$ corresponding to M largest singular values;

Estimate spectral parameters: $\hat{\mathbf{b}}_1 \leftarrow \hat{\mathbf{U}}^\top \hat{\mathbf{P}}_1$;

$\hat{\mathbf{b}}_\infty \leftarrow (\hat{\mathbf{P}}_2^\top \hat{\mathbf{U}})^\dagger \hat{\mathbf{P}}_1$;

for each observation $o \in [L]$ **do**

$\hat{\mathbf{C}}_o \leftarrow \hat{\mathbf{U}}^\top \hat{\mathbf{P}}_{3,o} (\hat{\mathbf{U}}^\top \hat{\mathbf{P}}_2)^\dagger$;

end

Hidden state belief update: $\hat{\mathbf{b}}_1 \leftarrow$ initial belief;

for $\tau = 1$ **to** N **do**

$$\hat{\mathbf{b}}_{\tau+1} \leftarrow \frac{\hat{\mathbf{C}}_{o_\tau} \hat{\mathbf{b}}_\tau}{\hat{\mathbf{b}}_\tau^\top \hat{\mathbf{C}}_{o_\tau} \hat{\mathbf{b}}_\tau};$$

end

Conditional probability prediction: for each possible next observation $o_{N+1} \in [L]$ **do**

$$\hat{P}(O_{N+1} = o_{N+1} | O_{1:N} = o_{1:N}) \leftarrow \frac{\hat{\mathbf{b}}_\infty^\top \hat{\mathbf{C}}_{o_{N+1}} \hat{\mathbf{b}}_{N+1}}{\sum_{k=1}^L \hat{\mathbf{b}}_\infty^\top \hat{\mathbf{C}}_k \hat{\mathbf{b}}_{N+1}};$$

end

return $\hat{P}(O_{N+1} | O_{1:N} = o_{1:N})$

F.1 Preliminaries

For a Markov chain with transition matrix \mathbf{A} , we let $\boldsymbol{\pi} \in \mathbb{R}_+^M$ denote the initial state distribution. We assume that $\boldsymbol{\pi}$ is also the stationary distribution of the Markov chain. This can be achieved by taking samples after a burn-in time which is proportional to $\frac{1}{1-\lambda_2(\mathbf{A})}$. Note that $\boldsymbol{\pi}_t = (\mathbf{A}^t)^\top \boldsymbol{\pi}$ is essentially a convex combination of rows of matrix \mathbf{A}^t , then by triangle inequality, we have $\|\boldsymbol{\pi}_t - \boldsymbol{\pi}_\infty\|_1 \leq \max_{i \in [M]} \|([\mathbf{A}^t]_{i,:})^\top - \boldsymbol{\pi}_\infty\|_1$. Thus, for an ergodic Markov matrix \mathbf{A} , we define the following to quantify the convergence of $\|\boldsymbol{\pi}_t - \boldsymbol{\pi}_\infty\|_1$. For an ergodic Markov matrix $\mathbf{A} \in \mathbb{R}_+^{M \times M}$, let $\tau_{\text{MC}} > 1$ and $\rho_{\text{MC}} \in (\lambda_2(\mathbf{A}), 1)$ be two constants [5, Theorem 4.9] such that

$$\max_{i \in [M]} \|([\mathbf{A}^t]_{i,:})^\top - \boldsymbol{\pi}_\infty\|_1 \leq \tau_{\text{MC}} \rho_{\text{MC}}^t. \quad (1)$$

Furthermore, we define the mixing time of \mathbf{A} as

$$t_{\text{MC}}(\epsilon) := \min \left\{ t \in \mathbb{N} : \max_{i \in [M]} \frac{1}{2} \|([\mathbf{A}^t]_{i,:})^\top - \boldsymbol{\pi}_\infty\|_1 \leq \epsilon \right\}. \quad (2)$$

Note that $\tau(\mathbf{M})$ and τ_{MC} have similar roles except $\tau(\mathbf{M})$ is usually used to study state matrices while τ_{MC} is for Markov matrices. For a square \mathbf{M} , we have $\|\mathbf{M}^k\| \leq \tau(\mathbf{M}) \rho(\mathbf{M})^k$, and for a Markov matrix, we have $\|\mathbf{A}^t - \mathbf{1}_M \boldsymbol{\pi}_\infty^\top\| \leq \tau_{\text{MC}} \rho_{\text{MC}}^t$.

217 F.2 Sample Complexity Analysis

218 In this section, we analyze the sample complexity of spectral learning algorithm (Alg 6) when the
 219 observation sequence is coming from a single trajectory. Our proof builds on [3] by modifying
 220 their analysis in Appendix A to incorporate single trajectory learning. We only present the Sample
 221 complexity analysis here and refer the reader to [3] for the remaining proofs.

222 F.3 Proof of Theorem 1

223 Fix $2 < T < N$, and recall from [3] that $[\mathbf{P}_1]_i = \mathbb{E}[\mathbf{1}_{\{o_T=i\}}]$, $[\mathbf{P}_2]_{i,j} = \mathbb{E}[\mathbf{1}_{\{o_T=i, o_{T-1}=j\}}]$,
 224 $[\mathbf{P}_{3,k}]_{i,j} = \mathbb{E}[\mathbf{1}_{\{o_T=i, o_{T-1}=k, o_{T-2}=j\}}]$, for all $k \in [L]$, when the initial distribution π is the stationary
 225 distribution of the Markov chain. In the following, we will present three different estimators for each
 226 of these quantities and analyze their convergence.

227 • **Estimation of \mathbf{P}_1 :** Let $\bar{N} := \lfloor \frac{N}{T} \rfloor$, and without loss of generality, suppose $\frac{N}{T}$ is an integer. Suppose
 228 $\{o_T^{(1)}, \dots, o_T^{(\bar{N})}\}$ be the i.i.d. samples obtained from \bar{N} independent trajectories of the HMM. We
 229 define the following three estimators of \mathbf{P}_1 ,

$$[\hat{\mathbf{P}}_1]_i = \frac{\sum_{k=1}^N \mathbf{1}_{\{o_k=i\}}}{N} \quad [\hat{\mathbf{P}}_1^{(\ell)}]_i = \frac{\sum_{k=1}^{\bar{N}} \mathbf{1}_{\{o_{kT-\ell}=i\}}}{\bar{N}} \quad [\hat{\mathbf{P}}_1^{(\perp)}]_i = \frac{\sum_{k=1}^{\bar{N}} \mathbf{1}_{\{o_T^{(k)}=i\}}}{\bar{N}}, \quad (3)$$

230 for all $\ell = 0, \dots, T-1$. By triangle inequality, we have

$$\|\hat{\mathbf{P}}_1 - \mathbf{P}_1\| \leq \|\hat{\mathbf{P}}_1 - \hat{\mathbf{P}}_1^{(\perp)}\| + \|\hat{\mathbf{P}}_1^{(\perp)} - \mathbf{P}_1\|. \quad (4)$$

231 [3] showed that, with probability at least $1 - \delta$, we have, $\|\hat{\mathbf{P}}_1^{(\perp)} - \mathbf{P}_1\| \lesssim \sqrt{\frac{\log(1/\delta)}{N}} + \sqrt{\frac{1}{N}}$. In the
 232 following, we will upper bound the term $\|\hat{\mathbf{P}}_1 - \hat{\mathbf{P}}_1^{(\perp)}\|$ by considering entry-wise concentration of
 233 each ℓ -th subtrajectory as follows: We have

$$[\hat{\mathbf{P}}_1^{(\ell)}]_i - [\hat{\mathbf{P}}_1^{(\perp)}]_i = \frac{\sum_{k=1}^{\bar{N}} (\mathbf{1}_{\{o_{kT-\ell}=i\}} - \mathbf{1}_{\{o_T^{(k)}=i\}})}{\bar{N}}. \quad (5)$$

234 First, we observe that $\mathbb{E}[\mathbf{1}_{\{o_{kT-\ell}=i\}} - \mathbf{1}_{\{o_T^{(k)}=i\}}] = 0$. Moreover, $|\mathbf{1}_{\{o_{kT-\ell}=i\}} - \mathbf{1}_{\{o_T^{(k)}=i\}}| \leq 1$,
 235 almost surely. However, the summation in (5) has weakly dependent terms. Therefore, we use the
 236 Bernstein type inequality for a class of weakly dependent and bounded random variables proposed
 237 in [6]. Before that, we need to upper bound the variance of the summation in (5). Observing that
 238 $\mathbb{E}[\hat{\mathbf{P}}_1^{(\ell)}]_i - [\hat{\mathbf{P}}_1^{(\perp)}]_i = 0$, we have,

$$\begin{aligned} \text{Var}([\hat{\mathbf{P}}_1^{(\ell)}]_i - [\hat{\mathbf{P}}_1^{(\perp)}]_i) &:= \mathbb{E} \left[\left([\hat{\mathbf{P}}_1^{(\ell)}]_i - [\hat{\mathbf{P}}_1^{(\perp)}]_i \right)^2 \right], \\ &= \mathbb{E} \left[\left([\hat{\mathbf{P}}_1^{(\ell)}]_i \right)^2 \right] + \mathbb{E} \left[\left([\hat{\mathbf{P}}_1^{(\perp)}]_i \right)^2 \right] - 2 \mathbb{E} \left[[\hat{\mathbf{P}}_1^{(\ell)}]_i [\hat{\mathbf{P}}_1^{(\perp)}]_i \right]. \end{aligned} \quad (6)$$

239 In the following, we will upper bound each term in (6) separately. We begin with,

$$\begin{aligned} \mathbb{E} \left[\left([\hat{\mathbf{P}}_1^{(\ell)}]_i \right)^2 \right] &= \frac{1}{N^2} \mathbb{E} \left[\sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbf{1}_{\{o_{kT-\ell}=i\}} \mathbf{1}_{\{o_{k'T-\ell}=i\}} \right], \\ &= \frac{1}{N^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{E} \left[\mathbf{1}_{\{o_{kT-\ell}=i, o_{k'T-\ell}=i\}} \right], \\ &= \frac{1}{N^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{P}(O_{kT-\ell} = i, O_{k'T-\ell} = i), \\ &= \frac{[\mathbf{B}^\top \boldsymbol{\pi}]_i}{N} + \frac{1}{N^2} \sum_{k=1}^{\bar{N}} \sum_{\substack{k'=1 \\ k' \neq k}}^{\bar{N}} [\mathbf{B}^\top \text{diag}(\boldsymbol{\pi}) \mathbf{A}^{|k-k'|T} \mathbf{B}]_{i,i}. \end{aligned} \quad (7)$$

240 Next, we have,

$$\begin{aligned}
\mathbb{E} \left[\left([\hat{\mathbf{P}}_1^{(\perp)}]_i \right)^2 \right] &= \frac{1}{\bar{N}^2} \mathbb{E} \left[\sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbf{1}_{\{o_T^{(k)}=i\}} \mathbf{1}_{\{o_T^{(k')}=i\}} \right], \\
&= \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{E} \left[\mathbf{1}_{\{o_T^{(k)}=i, o_T^{(k')}=i\}} \right], \\
&= \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{P} \left(o_T^{(k)} = i, o_T^{(k')} = i \right) = \frac{[\mathbf{B}^\top \boldsymbol{\pi}]_i}{\bar{N}} + (\bar{N} - 1) \frac{[\mathbf{B}^\top \boldsymbol{\pi}]_i^2}{\bar{N}}. \quad (8)
\end{aligned}$$

241 Lastly, we have

$$\begin{aligned}
\mathbb{E} \left[\left([\hat{\mathbf{P}}_1^{(\ell)}]_i \right) \left([\hat{\mathbf{P}}_1^{(\perp)}]_i \right) \right] &= \frac{1}{\bar{N}^2} \mathbb{E} \left[\sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbf{1}_{\{o_{kT-\ell}=i\}} \mathbf{1}_{\{o_T^{(k')}=i\}} \right], \\
&= \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{E} \left[\mathbf{1}_{\{o_{kT-\ell}=i, o_T^{(k')}=i\}} \right], \\
&= \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{P} \left(o_{kT-\ell} = i, o_T^{(k')} = i \right) = [\mathbf{B}^\top \boldsymbol{\pi}]_i^2. \quad (9)
\end{aligned}$$

242 Combining (7), (8), and (9) into (6), we get

$$\begin{aligned}
\mathbf{Var} \left([\hat{\mathbf{P}}_1^{(\ell)}]_i - [\hat{\mathbf{P}}_1^{(\perp)}]_i \right) &= \frac{2[\mathbf{B}^\top \boldsymbol{\pi}]_i}{\bar{N}} + \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{\substack{k'=1 \\ k' \neq k}}^{\bar{N}} \left[\mathbf{B}^\top \mathbf{diag}(\boldsymbol{\pi}) \mathbf{A}^{|k-k'|T} \mathbf{B} \right]_{i,i} \\
&\quad - (\bar{N} + 1) \frac{[\mathbf{B}^\top \boldsymbol{\pi}]_i^2}{\bar{N}}, \\
&= \frac{2([\mathbf{B}^\top \boldsymbol{\pi}]_i - [\mathbf{B}^\top \boldsymbol{\pi}]_i^2)}{\bar{N}} + \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{\substack{k'=1 \\ k' \neq k}}^{\bar{N}} \left[\mathbf{B}^\top \mathbf{diag}(\boldsymbol{\pi}) \mathbf{A}^{|k-k'|T} \mathbf{B} \right]_{i,i} \\
&\quad - (\bar{N} - 1) \frac{[\mathbf{B}^\top \boldsymbol{\pi}]_i^2}{\bar{N}}, \\
&= \frac{2(\mathbf{b}_i^\top \boldsymbol{\pi} - (\mathbf{b}_i^\top \boldsymbol{\pi})^2)}{\bar{N}} \\
&\quad + \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{\substack{k'=1 \\ k' \neq k}}^{\bar{N}} \mathbf{b}_i^\top \mathbf{diag}(\boldsymbol{\pi}) \left(\mathbf{A}^{|k-k'|T} - \mathbf{1}_M \boldsymbol{\pi}^\top \right) \mathbf{b}_i, \\
&\lesssim \frac{\mathbf{b}_i^\top \boldsymbol{\pi} - (\mathbf{b}_i^\top \boldsymbol{\pi})^2}{\bar{N}} + \frac{\|\mathbf{b}_i\|^2 \tau_{\text{MC}} \rho_{\text{MC}}^T}{\bar{N}(1 - \rho_{\text{MC}}^T)} \lesssim \frac{\mathbf{b}_i^\top \boldsymbol{\pi} - (\mathbf{b}_i^\top \boldsymbol{\pi})^2}{\bar{N}}, \quad (10)
\end{aligned}$$

243 where \mathbf{b}_i denotes the i -th column of \mathbf{B} and we get the last inequality by choosing,

$$T \gtrsim \log \left(\frac{\|\mathbf{b}_i\|^2 \tau_{\text{MC}}}{(\mathbf{b}_i^\top \boldsymbol{\pi} - (\mathbf{b}_i^\top \boldsymbol{\pi})^2)(1 - \rho_{\text{MC}}^T)} \right) / (1 - \rho). \quad (11)$$

244 Hence, using the Bernstein type inequality for weakly dependent and bounded random variables
245 (Theorem 1 in [6]), together with (10) (11), and the observations we made right after (5), with
246 probability at least $1 - \delta$, we have

$$|[\hat{\mathbf{P}}_1^{(\ell)}]_i - [\hat{\mathbf{P}}_1^{(\perp)}]_i| \lesssim \sqrt{\frac{(\mathbf{b}_i^\top \boldsymbol{\pi} - (\mathbf{b}_i^\top \boldsymbol{\pi})^2)}{\bar{N}} \log \left(\frac{1}{\delta} \right)}. \quad (12)$$

247 Union bounding over all $i \in [L]$, and $\ell \in \{0, 1, \dots, T-1\}$, with probability at least $1 - \delta$, we have

$$\|\hat{\mathbf{P}}_1^{(\ell)} - \hat{\mathbf{P}}_1^{(\perp)}\| \lesssim \sqrt{\frac{\mathbf{1}_L^\top \mathbf{B}^\top \boldsymbol{\pi} - \|\mathbf{B}^\top \boldsymbol{\pi}\|^2}{\bar{N}} \log\left(\frac{LT}{\delta}\right)}, \quad (13)$$

248 given $T \gtrsim \max_{i \in [L]} \left\{ \log\left(\frac{\|\mathbf{b}_i\|^2 \tau_{\text{MC}}}{(\mathbf{b}_i^\top \boldsymbol{\pi} - (\mathbf{b}_i^\top \boldsymbol{\pi})^2)(1 - \rho_{\text{MC}}^T)}\right) \right\} / (1 - \rho)$. This further implies that, with
 249 probability at least $1 - \delta$, the same upper bound also holds for $\|\hat{\mathbf{P}}_1 - \hat{\mathbf{P}}_1^{(\perp)}\|$. Combining this with
 250 (4) and [3], with probability at least $1 - \delta$, we have

$$\|\hat{\mathbf{P}}_1 - \mathbf{P}_1\| \lesssim \sqrt{\frac{\log(1/\delta)}{\bar{N}}} + \sqrt{\frac{1}{\bar{N}}} + \sqrt{\frac{\mathbf{1}_L^\top \mathbf{B}^\top \boldsymbol{\pi} - \|\mathbf{B}^\top \boldsymbol{\pi}\|^2}{\bar{N}} \log\left(\frac{LT}{\delta}\right)}. \quad (14)$$

251 **• Estimation of \mathbf{P}_2 :** Here, we follow a similar line of reasoning as above. We begin with defining
 252 the three estimators of \mathbf{P}_2 as follows,

$$\begin{aligned} [\hat{\mathbf{P}}_2]_{i,j} &= \frac{\sum_{k=1}^N \mathbf{1}_{\{o_k=i, o_{k-1}=j\}}}{N} & [\hat{\mathbf{P}}_2^{(\ell)}]_{i,j} &= \frac{\sum_{k=1}^{\bar{N}} \mathbf{1}_{\{o_{kT-\ell}=i, o_{kT-\ell-1}=j\}}}{\bar{N}}, \\ [\hat{\mathbf{P}}_2^{(\perp)}]_{i,j} &= \frac{\sum_{k=1}^{\bar{N}} \mathbf{1}_{\{o_T^{(k)}=i, o_{T-1}^{(k)}=j\}}}{\bar{N}} \end{aligned} \quad (15)$$

253 Similar to \mathbf{P}_1 , we consider the entry-wise concentration of each ℓ -th subtrajectory as follows,

$$[\hat{\mathbf{P}}_2^{(\ell)}]_{i,j} - [\hat{\mathbf{P}}_2^{(\perp)}]_{i,j} = \frac{\sum_{k=1}^{\bar{N}} \left(\mathbf{1}_{\{o_{kT-\ell}=i, o_{kT-\ell-1}=j\}} - \mathbf{1}_{\{o_T^{(k)}=i, o_{T-1}^{(k)}=j\}} \right)}{\bar{N}}. \quad (16)$$

254 Observing that $\mathbb{E} \left[[\hat{\mathbf{P}}_2^{(\ell)}]_{i,j} - [\hat{\mathbf{P}}_2^{(\perp)}]_{i,j} \right] = 0$, we have,

$$\begin{aligned} \text{Var} \left([\hat{\mathbf{P}}_2^{(\ell)}]_{i,j} - [\hat{\mathbf{P}}_2^{(\perp)}]_{i,j} \right) &= \mathbb{E} \left[\left([\hat{\mathbf{P}}_2^{(\ell)}]_{i,j} - [\hat{\mathbf{P}}_2^{(\perp)}]_{i,j} \right)^2 \right], \\ &= \mathbb{E} \left[\left([\hat{\mathbf{P}}_2^{(\ell)}]_{i,j} \right)^2 \right] + \mathbb{E} \left[\left([\hat{\mathbf{P}}_2^{(\perp)}]_{i,j} \right)^2 \right] - 2 \mathbb{E} \left[[\hat{\mathbf{P}}_2^{(\ell)}]_{i,j} [\hat{\mathbf{P}}_2^{(\perp)}]_{i,j} \right]. \end{aligned} \quad (17)$$

255 In the following, we will upper bound each term in (17) separately. We begin with,

$$\begin{aligned} \mathbb{E} \left[\left([\hat{\mathbf{P}}_2^{(\ell)}]_{i,j} \right)^2 \right] &= \frac{1}{\bar{N}^2} \mathbb{E} \left[\sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbf{1}_{\{o_{kT-\ell}=i, o_{kT-\ell-1}=j\}} \mathbf{1}_{\{o_{k'T-\ell}=i, o_{k'T-\ell-1}=j\}} \right], \\ &= \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{E} \left[\mathbf{1}_{\{o_{kT-\ell}=i, o_{kT-\ell-1}=j, o_{k'T-\ell}=i, o_{k'T-\ell-1}=j\}} \right], \\ &= \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{P} (O_{kT-\ell} = i, O_{kT-\ell-1} = j, O_{k'T-\ell} = i, O_{k'T-\ell-1} = j), \\ &= \frac{\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M}{\bar{N}} + \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{\substack{k'=1 \\ k' \neq k}}^{\bar{N}} \boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{A}^{|k-k'|T-1} \mathbf{D}_{j,i} \mathbf{1}_M, \end{aligned} \quad (18)$$

256 where, given the i -th column \mathbf{b}_i , and the j -th column \mathbf{b}_j of \mathbf{B} , we define

$$\mathbf{D}_{j,i} := \text{diag}(\mathbf{b}_j) \mathbf{A} \text{diag}(\mathbf{b}_i). \quad (19)$$

257 Next, we have

$$\begin{aligned}
\mathbb{E} \left[\left([\hat{\mathbf{P}}_2^{(\perp)}]_{i,j} \right)^2 \right] &= \frac{1}{\bar{N}^2} \mathbb{E} \left[\sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbf{1}_{\{o_T^{(k)}=i, o_{T-1}^{(k)}=j\}} \mathbf{1}_{\{o_T^{(k')}=i, o_{T-1}^{(k')}=j\}} \right], \\
&= \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{E} \left[\mathbf{1}_{\{o_T^{(k)}=i, o_{T-1}^{(k)}=j, o_T^{(k')}=i, o_{T-1}^{(k')}=j\}} \right], \\
&= \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{P} \left(O_T^{(k)} = i, O_{T-1}^{(k)} = j, O_T^{(k')} = i, O_{T-1}^{(k')} = j \right), \\
&= \frac{\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M}{\bar{N}} + (\bar{N} - 1) \frac{(\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M)^2}{\bar{N}}. \tag{20}
\end{aligned}$$

258 Lastly, we have

$$\begin{aligned}
\mathbb{E} \left[\left([\hat{\mathbf{P}}_2^{(\ell)}]_{i,j} \right) \left([\hat{\mathbf{P}}_2^{(\perp)}]_{i,j} \right) \right] &= \frac{1}{\bar{N}^2} \mathbb{E} \left[\sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbf{1}_{\{o_{kT-\ell}^{(k)}=i, o_{kT-\ell-1}^{(k)}=j\}} \mathbf{1}_{\{o_T^{(k')}=i, o_{T-1}^{(k')}=j\}} \right], \\
&= \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{E} \left[\mathbf{1}_{\{o_{kT-\ell}^{(k)}=i, o_{kT-\ell-1}^{(k)}=j, o_T^{(k')}=i, o_{T-1}^{(k')}=j\}} \right], \\
&= \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{k'=1}^{\bar{N}} \mathbb{P} \left(O_{kT-\ell}^{(k)} = i, O_{kT-\ell-1}^{(k)} = j, O_T^{(k')} = i, O_{T-1}^{(k')} = j \right), \\
&= (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M)^2. \tag{21}
\end{aligned}$$

259 Combining (18), (20), and (21) into (17), we get

$$\begin{aligned}
\mathbf{Var} \left([\hat{\mathbf{P}}_2^{(\ell)}]_{i,j} - [\hat{\mathbf{P}}_2^{(\perp)}]_{i,j} \right) &= \frac{2\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M}{\bar{N}} + \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{\substack{k'=1 \\ k' \neq k}}^{\bar{N}} \boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{A}^{|k-k'|T-1} \mathbf{D}_{j,i} \mathbf{1}_M \\
&\quad - (\bar{N} + 1) \frac{(\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M)^2}{\bar{N}}, \\
&= \frac{2(\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M - (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M)^2)}{\bar{N}} \\
&\quad + \frac{1}{\bar{N}^2} \sum_{k=1}^{\bar{N}} \sum_{\substack{k'=1 \\ k' \neq k}}^{\bar{N}} \boldsymbol{\pi}^\top \mathbf{D}_{j,i} \left(\mathbf{A}^{|k-k'|T-1} - \mathbf{1}_M \boldsymbol{\pi}^\top \right) \mathbf{D}_{j,i} \mathbf{1}_M, \\
&\lesssim \frac{2(\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M - (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M)^2)}{\bar{N}} + \frac{\|\boldsymbol{\pi}^\top \mathbf{D}_{j,i}\| \|\mathbf{D}_{j,i} \mathbf{1}_M\| \tau_{\text{MC}} \rho_{\text{MC}}^{T-1}}{\bar{N}(1 - \rho_{\text{MC}}^T)}, \\
&\lesssim \frac{\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M - (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M)^2}{\bar{N}}, \tag{22}
\end{aligned}$$

260 where we get the last inequality by choosing,

$$T \gtrsim 1 + \log \left(\frac{\|\boldsymbol{\pi}^\top \mathbf{D}_{j,i}\| \|\mathbf{D}_{j,i} \mathbf{1}_M\| \tau_{\text{MC}}}{(\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M - (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M)^2)(1 - \rho_{\text{MC}}^T)} \right) / (1 - \rho). \tag{23}$$

261 Hence, using similar line of reasoning as we did in the case of P_1 , with probability at least $1 - \delta$, we
262 have

$$\|\hat{\mathbf{P}}_2^{(\ell)} - \hat{\mathbf{P}}_2^{(\perp)}\| \lesssim \sqrt{\frac{\sum_{i,j=1}^L (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M - (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M)^2)}{\bar{N}} \log \left(\frac{L^2 T}{\delta} \right)}, \tag{24}$$

263 given $T \gtrsim 1 + \max_{i,j \in [L]} \left\{ \log \left(\frac{\|\boldsymbol{\pi}^\top \mathbf{D}_{j,i}\| \|\mathbf{D}_{j,i} \mathbf{1}_M\| \tau_{\text{MC}}}{(\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M - (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M)^2)(1 - \rho_{\text{MC}}^T)} \right) \right\} / (1 - \rho)$. This further implies
 264 that, with probability at least $1 - \delta$, the same upper bound also holds for $\|\hat{\mathbf{P}}_2 - \hat{\mathbf{P}}_2^{(\perp)}\|$. Combining
 265 this with the triangle inequality and [3], with probability at least $1 - \delta$, we have

$$\|\hat{\mathbf{P}}_2 - \mathbf{P}_2\| \lesssim \sqrt{\frac{\log(1/\delta)}{\bar{N}}} + \sqrt{\frac{1}{\bar{N}}} + \sqrt{\frac{\sum_{i,j=1}^L (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M - (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M)^2)}{\bar{N}} \log \left(\frac{L^2 T}{\delta} \right)}. \quad (25)$$

266 **• Estimation of \mathbf{P}_3 :** Here, we follow a similar line of reasoning as above. We begin with defining
 267 the three estimators of \mathbf{P}_3 as follows,

$$\begin{aligned} [\hat{\mathbf{P}}_{3,n}]_{i,j} &= \frac{\sum_{k=1}^N \mathbf{1}_{\{o_k=i, o_{k-1}=n, o_{k-2}=j\}}}{N} & [\hat{\mathbf{P}}_{3,n}^{(\ell)}]_{i,j} &= \frac{\sum_{k=1}^{\bar{N}} \mathbf{1}_{\{o_{kT-\ell}=i, o_{kT-\ell-1}=n, o_{kT-\ell-2}=j\}}}{\bar{N}}, \\ [\hat{\mathbf{P}}_{3,n}^{(\perp)}]_{i,j} &= \frac{\sum_{k=1}^{\bar{N}} \mathbf{1}_{\{o_T^{(k)}=i, o_{T-1}^{(k)}=n, o_{T-2}^{(k)}=j\}}}{\bar{N}} \end{aligned} \quad (26)$$

268 Following the same line of reasoning as we did in the case of \mathbf{P}_2 , with probability at least $1 - \delta$, we
 269 have

$$\begin{aligned} &\|\hat{\mathbf{P}}_{3,n} - \mathbf{P}_{3,n}\| \\ &\lesssim \sqrt{\frac{\log(1/\delta)}{\bar{N}}} + \sqrt{\frac{1}{\bar{N}}} + \sqrt{\frac{\sum_{i,j,n=1}^L (\boldsymbol{\pi}^\top \mathbf{D}_{j,n,i} \mathbf{1}_M - (\boldsymbol{\pi}^\top \mathbf{D}_{j,n,i} \mathbf{1}_M)^2)}{\bar{N}} \log \left(\frac{L^3 T}{\delta} \right)}, \end{aligned} \quad (27)$$

270 provided that,

$$T \gtrsim 2 + \max_{i,j,n \in [L]} \left\{ \log \left(\frac{\|\boldsymbol{\pi}^\top \mathbf{D}_{j,n,i}\| \|\mathbf{D}_{j,n,i} \mathbf{1}_M\| \tau_{\text{MC}}}{(\boldsymbol{\pi}^\top \mathbf{D}_{j,n,i} \mathbf{1}_M - (\boldsymbol{\pi}^\top \mathbf{D}_{j,n,i} \mathbf{1}_M)^2)(1 - \rho_{\text{MC}}^T)} \right) \right\} / (1 - \rho), \quad (28)$$

271 where, given the i -th column \mathbf{b}_i , the j -th column \mathbf{b}_j and the n -th column \mathbf{b}_n of \mathbf{B} , we define

$$\mathbf{D}_{j,n,i} := \text{diag}(\mathbf{b}_j) \mathbf{A} \text{diag}(\mathbf{b}_n) \mathbf{A} \text{diag}(\mathbf{b}_i) \quad (29)$$

272 **• Finalizing the proof:** Theorem 1 follows by repeating the proof of Theorem 7 in [3], with the i.i.d.
 273 estimators replaced by the single trajectory estimators, and the values of ϵ_1 , $\epsilon_{2,1}$ and $\epsilon_{3,x,1}$ replaced
 274 by,

$$\begin{aligned} \epsilon_1 &\lesssim \sqrt{\frac{\log(1/\delta)}{\bar{N}}} + \sqrt{\frac{1}{\bar{N}}} + \sqrt{\frac{\mathbf{1}_L^\top \mathbf{B}^\top \boldsymbol{\pi} - \|\mathbf{B}^\top \boldsymbol{\pi}\|^2}{\bar{N}} \log \left(\frac{LT}{\delta} \right)}, \\ \epsilon_{2,1} &\lesssim \sqrt{\frac{\log(1/\delta)}{\bar{N}}} + \sqrt{\frac{1}{\bar{N}}} + \sqrt{\frac{\sum_{i,j=1}^L (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M - (\boldsymbol{\pi}^\top \mathbf{D}_{j,i} \mathbf{1}_M)^2)}{\bar{N}} \log \left(\frac{L^2 T}{\delta} \right)}, \\ \epsilon_{3,x,1} &\lesssim \sqrt{\frac{\log(1/\delta)}{\bar{N}}} + \sqrt{\frac{1}{\bar{N}}} + \sqrt{\frac{\sum_{i,j,n=1}^L (\boldsymbol{\pi}^\top \mathbf{D}_{j,n,i} \mathbf{1}_M - (\boldsymbol{\pi}^\top \mathbf{D}_{j,n,i} \mathbf{1}_M)^2)}{\bar{N}} \log \left(\frac{L^3 T}{\delta} \right)}, \end{aligned}$$

275 where $\bar{N} = \lfloor \frac{N}{T} \rfloor = \mathcal{O}(N(1 - \lambda_2(\mathbf{A})))$. The proof is completed by upper bounding the Hellinger-
 276 distance in terms of KL-distance.

G Additional Real World Experiments

We design an additional experiment using real-world datasets to validate our findings. We artificially simulate different emission entropy levels for the same underlying hidden transition process by controlling the amount of information included in the observation sequence. Using complete information corresponds to low emission entropy, while limiting information artificially increases emission entropy.

We use the IBL decision-making mice dataset [4]. In our LLM in-context learning experiment, we implement four ablation conditions that vary the information presented in each trial: (i) “choice only”; (ii) “choice reward”; (iii) “stimulus choice”; (iv) “stimulus choice reward”. Note that the baseline GLM-HMM uses all available information as in condition (iv). These ablations describe the same underlying mouse decision-making sequences but with varying levels of environmental state detail.

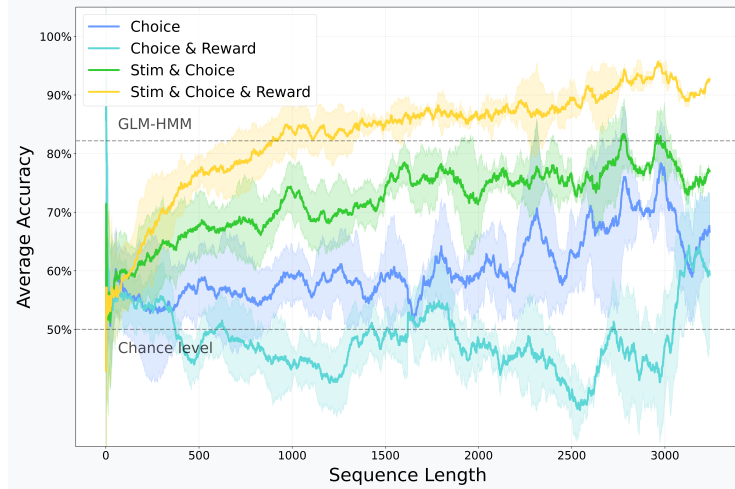


Figure 16: LLM in-context learning prediction accuracy for mice decision-making task with varying types of information in the observed sequences. Each line is averaged over 7 mice, with $1\text{-}\sigma$ error bar. The model we use is Qwen2.5-7B.

The results shown in Figure 16 reveal significant differences across ablation conditions: while “stimulus choice reward” achieves performance exceeding GLM-HMM, “choice reward” is merely at chance level with its convergence trend similar to the synthetic experiments when the transitions or emissions are near random. This demonstrates that accurately modeling mouse decision-making in this task requires both stimulus and reward information.

These findings highlight a broader principle: obtaining appropriate information (corresponding to low emission entropy) is essential for successful task modeling. This experiment parallels real-world experimental design, where scientists must choose which signals to collect when studying task structure. When researchers omit critical information needed to describe a sequence, it can easily lead to incorrect conclusions about the underlying process.

References

- [1] Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. Wiley-Interscience, USA, 2006. ISBN 0471241954.
- [2] Robert G Gallager. Discrete stochastic processes. *Journal of the Operational Research Society*, 48(1):103–103, 1997.
- [3] Daniel Hsu, Sham M Kakade, and Tong Zhang. A spectral algorithm for learning hidden markov models. *Journal of Computer and System Sciences*, 78(5):1460–1480, 2012.
- [4] The International Brain Laboratory, Valeria Aguilon-Rodriguez, Dora Angelaki, Hannah Bayer, Niccolo Bonacchi, Matteo Carandini, Fanny Cazettes, Gaelle Chapuis, Anne K Churchland, Yang Dan, Eric Dewitt, Mayo Faulkner, Hamish Forrest, Laura Haetzel, Michael Häusser, Sonja B Hofer, Fei Hu, Anup Khanal, Christopher Krasniak, Ines Laranjeira, Zachary F Mainen, Guido Meijer, Nathaniel J Miska, Thomas D Mrsic-Flogel, Masayoshi Murakami, Jean-Paul Noel, Alejandro Pan-Vazquez, Cyrille Rossant, Joshua Sanders, Karolina Socha, Rebecca Terry, Anne E Urai, Hernando Vergara, Miles Wells, Christian J Wilson, Ilana B Witten, Lauren E Wool, and Anthony M Zador. Standardized and reproducible measurement of decision-making in mice. *eLife*, 10:e63711, may 2021. ISSN 2050-084X. doi: 10.7554/eLife.63711. URL <https://doi.org/10.7554/eLife.63711>.
- [5] David A Levin and Yuval Peres. *Markov chains and mixing times*, volume 107. American Mathematical Soc., 2017.
- [6] Florence Merlevède, Magda Peligrad, and Emmanuel Rio. Bernstein inequality and moderate deviations under strong mixing conditions. In *High dimensional probability V: the Luminy volume*, volume 5, pages 273–293. Institute of Mathematical Statistics, 2009.
- [7] Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*, 2024.