

Beyond Memorization: Testing LLM Reasoning on Unseen Theory of Computation Tasks

Anonymous ACL submission

Abstract

Large language models (LLMs) have demonstrated strong performance on formal language tasks, yet whether this reflects genuine symbolic reasoning or pattern matching on familiar constructions remains unclear. We introduce a benchmark for deterministic finite automata (DFA) construction from regular languages, comprising factual knowledge questions, seen construction problems from public sources, and two types of unseen problems: hand-crafted instances with multiple interacting constraints and systematically generated problems via Arden’s theorem. Models achieve perfect accuracy on factual questions and 84-90% on seen tasks. However, accuracy drops sharply on unseen problems (by 30-64%), with failures stemming from systematic misinterpretation of language constraints, incorrect handling of Kleene-star semantics, and a failure to preserve global consistency. We evaluate a three-stage hint protocol that enables correction of shallow errors but does not reliably resolve globally inconsistent or structurally flawed automata. Our analysis across multiple prompting strategies (direct, Chain-of-Thought, Tree-of-Thought) reveals that errors persist regardless of prompting approach, exposing a fundamental gap between LLMs’ ability to generate syntactically plausible DFAs and their capacity for semantically correct formal reasoning.

1 Introduction

Large language models (LLMs) have demonstrated remarkable performance on diverse reasoning benchmarks, from mathematical problem-solving (Lewkowycz et al., 2022; Welleck et al., 2021; Azerbayev et al., 2023) to code generation (Wu et al., 2022). However, a fundamental question remains unresolved: do these models perform genuine symbolic reasoning, or do they primarily rely on pattern matching over memorized examples? Recent work reveals persistent failures

on tasks requiring structured symbolic manipulation (Katz et al., 2025; Yue et al., 2024), suggesting that strong benchmark performance may not reflect robust reasoning capabilities.

We address this question through the lens of *deterministic finite automata (DFA) construction from regular languages*, a core problem in the Theory of Computation (ToC). Moreover, this also depicts the task of lexical analyser which are solved by using tools like Lex, Flex with limitations (Aho et al., 2006). This task offers unique advantages as a reasoning probe: (1) correctness is formally verifiable through exhaustive testing, (2) solutions require multi-step symbolic manipulation with global consistency constraints, (3) the space of possible problem instances is combinatorially vast, and (4) DFAs represent the simplest non-trivial computational model, ensuring that failures cannot be attributed to problem complexity or ambiguous specifications. We focus on prompting-based evaluation (rather than fine-tuning) to reflect practical LLM usage. Critically, while existing ToC benchmarks (Golesteanu et al., 2024; Zahraei and Asgari, 2024) evaluate factual knowledge and proof verification, they do not systematically control for *memorization versus compositional generalization*, a model may succeed by recalling similar problems from training data rather than reasoning from first principles.

To isolate genuine reasoning capability, we introduce a carefully designed benchmark with three components: (1) a *knowledge-checking dataset* verifying mastery of DFA definitions and properties, (2) a *seen construction dataset* comprising 90 publicly available DFA problems, and (3) an *unseen construction dataset* with 180 novel problems generated via two complementary approaches. The first approach, *mathematical art*, manually constructs problems with multiple interacting constraints, forbidden substrings, and narrative-based specifications (e.g., encoding

chess moves). The second approach, *mathematical engineering*, systematically generates problems via Arden’s theorem (Sipser, 2013), producing structurally complex regular expressions unlikely to appear in training data. This seen/unseen split enables us to measure the performance gap attributable to memorization versus reasoning.

We evaluate frontier LLMs: GPT-5.1, Gemini-2.5-Flash, and Grok-4.1-fast-reasoning – across multiple prompting strategies including Chain-of-Thought (CoT) and Tree-of-Thought (ToT) (Wang and Zhou, 2024; Yao et al., 2023). We also introduce a *three-stage hint protocol* that progressively reveals construction errors, enabling us to assess whether LLMs can self-correct when guided.

Main Findings. Our results reveal a stark dissociation between knowledge and reasoning: all models achieve 100% accuracy on factual questions and 84–90% on seen construction tasks, but accuracy drops sharply on unseen problems (20.67–59.12% under direct prompting, representing 30–64 % point drops). Detailed error analysis shows systematic failure modes: incorrect simplification of Kleene star semantics, failure to preserve constraints under concatenation, and introduction of spurious states. Critically, these failures persist across all prompting strategies, and the hint protocol primarily corrects shallow errors while leaving globally inconsistent automata uncorrected.

Contributions. This work makes the following contributions: (i) **Novel benchmark:** We introduce the first DFA construction benchmark with systematic seen/unseen splits, comprising 50 knowledge, 90 seen, and 180 unseen problems (60 hand-crafted, 120 via Arden’s theorem). (ii) **Controlled memorization study:** By evaluating structurally similar seen/unseen pairs, we provide the first evidence that LLM success on ToC tasks primarily reflects memorization rather than compositional reasoning. (iii) **Comprehensive prompting evaluation:** We evaluate CoT, ToT (with four construction methods: direct, minimization, derivative-based, Thompson’s algorithm), and a novel hint-based self-correction protocol. (iv) **Systematic failure taxonomy:** Through detailed analysis of 500+ incorrect DFAs, we identify six recurring failure modes (derivative normalization errors, constraint composition failures, etc.) that reveal fundamental limitations in symbolic state tracking. All datasets, prompts, evaluation code, and model outputs are released at

<https://anonymous.4open.science/r/dfa-llm-evaluation-B82D/> to support reproducibility and future research on formal reasoning in LLMs.

2 Seen Dataset and Task Formulation

2.1 Task Definition

We evaluate LLMs on the task of constructing deterministic finite automata (DFAs) from formal language specifications. Given a target language L specified either as a regular expression (RE) or natural language description, the model must produce a DFA $D = \langle Q, \Sigma, \delta, q_0, F \rangle$ that recognizes exactly L . Here, Q is a finite set of states, Σ is the input alphabet, $\delta : Q \times \Sigma \rightarrow Q$ is the transition function, $q_0 \in Q$ is the start state, and $F \subseteq Q$ is the set of accepting states (Sipser, 2013). A DFA D accepts string $w = w_1w_2 \cdots w_n \in \Sigma^*$ if there exists a sequence of states s_0, s_1, \dots, s_n such that (i) $s_0 = q_0$, (ii) $s_i = \delta(s_{i-1}, w_i)$ for all $1 \leq i \leq n$, and (iii) $s_n \in F$. We say D recognizes language L if $L = \{w \in \Sigma^* \mid D \text{ accepts } w\}$. A language is *regular* if there exists a DFA that recognizes it. Correct DFA construction requires models to: (1) parse and interpret formal language specifications, (2) identify the minimal state structure capturing all constraints, (3) design transitions ensuring acceptance of all and only valid strings, and (4) maintain global consistency across all states and symbols. Critically, there exist exponentially many invalid DFAs for any language, and small errors in state semantics or transition assignments can invalidate the entire construction.

Problem Instances. Each problem specifies the target language in one of two formats: (a) *Natural language:* $L_1 = \{\text{Construct a DFA over } \{a, b\} \text{ that accepts all strings in which the third-to-last symbol from the end must be 'a'}\}$. Here, if the problem format is given in natural language, it is not possible to design DFA using Lex, Flex (Aho et al., 2006). (b) *RE:* $L_1 = (a+b)^*a(a+b)(a+b)$. Both formats specify the same language (Figure 1(a)). Models must output a complete DFA specification including states, transitions, start state, and accepting states. We focus on *minimal* DFAs (fewest states) when possible, as minimality requires identifying semantically equivalent states – a reasoning challenge beyond simply enumerating valid transitions. However, during checking the performance of LLMs, we consider only correct DFA (particularly, not minimal DFA).

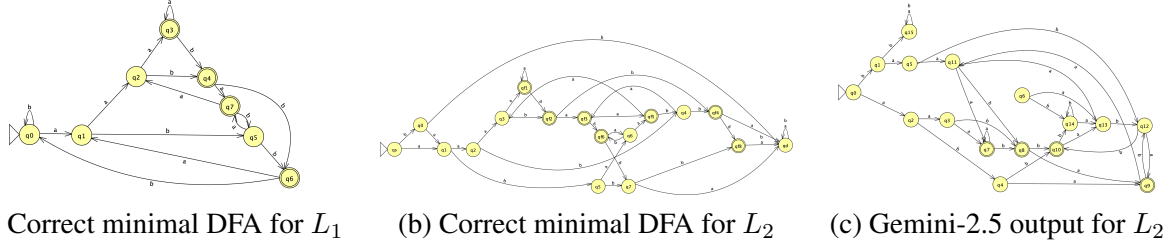


Figure 1: (a): Minimum DFA for language L_1 ; (b) & (c): Comparison of correct minimal DFA for unseen language L_2 (see b) and representative incorrect output (see c). Notation: \circ = state, \triangleright = start state, \odot = accepting state. Transition $q_0 \xrightarrow{a} q_1$ indicates that on input a from state q_0 , the DFA transitions to state q_1 (Please zoom for better readability).

Dataset	Prompting	GPT-5.1	Grok-4.1	Gemini-2.5
Knowledge	Zero-Shot	100%	100%	100%
Seen DFA	Zero-Shot	84.2%	89.5%	85.0%

Table 1: Results of knowledge and seen DFA datasets.

2.2 Knowledge and Seen DFA Dataset

Before evaluating construction ability, we verify that models possess foundational knowledge of automata theory. This dataset comprises 50 questions: 25 multiple-choice and 25 short-answer questions covering DFA definitions, RE properties, and the relationship between DFAs and nondeterministic FAs (NFAs). All questions are drawn from standard ToC textbooks and publicly available materials. **Results:** All models achieve 100% accuracy (Table 1), confirming complete mastery of relevant definitions and properties. Consequently, errors on construction tasks cannot be attributed to missing conceptual knowledge.

The seen dataset comprises 90 DFA construction problems collected from online university problem sets, textbooks, and publicly accessible ToC resources. For each problem, we verified that both the problem statement and solution DFA are publicly available online, ensuring these represent patterns likely encountered during pretraining.

Language Characteristics. The seen dataset includes standard patterns such as: (i) Suffix/prefix constraints (e.g., “ends with ab ”); (ii) Counting modulo k (e.g., “even number of a ’s”); (iii) Position-based constraints (e.g., “ 3^{rd} -last symbol a ”); and (iv) Boolean combinations (e.g., “contains aa or bb ”).

Evaluation Protocol. We evaluate models using zero-shot direct prompting (no worked examples). Models are queried via official APIs with temperature set to 0 for deterministic decoding. Each model receives the problem specification and must output a complete DFA in structured JSON format

specifying states, alphabet, transitions, start state, and accepting states (Appendix D).

Results. Table 1 shows that all models achieve strong performance: Grok-4.1-fast-reasoning achieves the highest success rate (89.5%), followed by Gemini-2.5-Flash (85.0%) and GPT-5.1 (84.2%). For the example language L_1 , all three models produce the correct minimal DFA. These results suggest that models can successfully construct DFAs for *familiar* problem patterns.

The Memorization Question. High accuracy on seen problems does not imply genuine reasoning capability. Models may succeed by retrieving similar examples from training data rather than performing compositional symbolic manipulation. To isolate reasoning from memorization, we next introduce the unseen construction dataset with carefully controlled novelty.

2.3 The Memorization Question

Strong performance on the seen dataset (84–90% accuracy) might suggest that LLMs possess robust DFA construction capabilities. However, this conclusion is premature: models may succeed by retrieving memorized solution patterns rather than performing genuine symbolic reasoning. To distinguish these explanations, we must evaluate performance on *structurally novel* problems that require compositional generalization. **Motivating Example:** Consider two closely related languages: L_1 (**seen**): Accepts all strings over $\{a, b\}$ where the third-to-last symbol is ‘a’; L_2 (**unseen**): Accepts all strings over $\{a, b\}$ where (i) the fourth-to-last symbol is ‘a’, and (ii) substring ‘bb’ does not appear before any ‘a’. Language L_2 extends L_1 with one additional constraint (forbidding ‘bb’ before ‘a’) and a minor modification to the position constraint (4^{th} -last instead of 3^{rd} -last). Both languages require similar reasoning – tracking symbol positions while enforcing order-

ing constraints – yet L_2 is constructed to be absent from public problem sets and textbook solutions.

Empirical Evidence. Despite achieving 100% accuracy on L_1 , *all models fail on L_2 under direct prompting*. Figure 1(c) shows Gemini-2.5-Flash’s output: while the constructed DFA accepts some valid strings (e.g., ‘aaaa’, ‘aab’), it also accepts invalid strings such as ‘aaa’ (violates position constraint), and ‘aabbaba’ (violates both constraints). Similar systematic errors occur for GPT-5.1 and Grok-4.1-fast-reasoning (Appendix A). This failure is particularly revealing because L_1 and L_2 differ only in *constraint composition*, not in fundamental reasoning requirements. The models correctly construct DFAs for simpler positional constraints (seen in training) but fail when these constraints are combined with ordering restrictions – suggesting success on seen tasks reflects pattern retrieval rather than robust symbolic reasoning.

Research Question. These observations motivate our central research question: *To what extent does LLM performance on formal reasoning tasks depend on memorization of training examples versus compositional symbolic reasoning?* Answering this question requires systematic evaluation on unseen problems that control for structural novelty while preserving task and reasoning requirements.

Unseen Dataset Design. To enable this controlled evaluation, we construct an unseen DFA dataset using two complementary approaches: (i) **Mathematical Art (60 problems):** Manually designed problems with multiple interacting constraints, forbidden substrings, positional patterns, and narrative-based specifications (e.g., encoding chess openings, maze navigation). These problems require creative constraint combination absent from standard curricula. (ii) **Mathematical Engineering (120 problems):** Systematically generated problems via Arden’s theorem (Sipser, 2013). We construct random NFAs, derive their accepted languages through algebraic elimination, and use the resulting regular expressions (often highly nested and non-standard) as problem specifications. This approach ensures structural diversity and scalability.

All unseen problems are manually verified to be absent from public problem repositories, textbooks, and online course materials. We categorize problems by difficulty (easy, medium, hard) based on the number of constraints (Art) or nesting depth (Engineering), enabling fine-grained analysis of where reasoning breaks down. Detailed con-

struction procedures and representative examples are provided in Section 3.

Evaluation Framework. Beyond measuring accuracy on unseen problems, we evaluate multiple prompting strategies (CoT, ToT with four construction branches) and introduce a three-stage hint protocol to assess self-correction capability. This comprehensive evaluation isolates whether failures stem from initial misinterpretation, inability to maintain symbolic consistency, or fundamental reasoning deficits that persist even with corrective guidance.

3 Unseen DFA Construction Dataset

We employ two generation strategies: *manual constraint composition* (60 problems) for creative problem design, and *systematic Arden’s theorem inversion* (120 problems) for scalable generation with guaranteed structural diversity.

3.1 Manual Constraint Composition

This approach extends seen DFA patterns through controlled increases in constraint complexity. Following the design principles illustrated by the $L_1 \rightarrow L_2$ transformation (Section 2.3), we systematically combine multiple constraints that rarely co-occur in standard curricula: (i) **Product constructions:** Multiple conjunctive/disjunctive constraints requiring state-space products (e.g., “strings starting with ‘aba’ OR ‘bab’ AND ending with their reverse”). (ii) **Interacting constraints:** Independent conditions that cannot be verified locally (e.g., “fourth-to-last symbol is ‘a’ AND substring ‘bb’ never precedes ‘a’” – language L_2). (iii) **Structural restrictions:** Positional patterns and forbidden substrings (e.g., “XOR of first three bits equals final bit”; “every 4-bit window’s product is divisible by 4”). (iv) **Narrative encodings:** Real-world scenarios requiring formalization (e.g., chess opening sequences encoded as moves; fruit-mixing recipes as syrup combinations; see Appendix B for additional examples). All 60 problems were manually verified to be absent from the top 100 Google search results, standard textbooks (Sipser, 2013; Hopcroft et al., 2001), and popular online repositories (GitHub, StackOverflow, course websites). This approach yields high-quality problems testing creative constraint integration, but scalability is limited by human effort.

3.2 Generation via Arden's Theorem

To address scalability while ensuring structural novelty, we employ a reverse-engineering approach: generate random NFAs, derive their accepted languages algebraically using Arden's theorem, and use the resulting regular expressions as problem specifications. This systematic procedure guarantees that each generated problem has a unique structure determined by random NFA topology rather than memorized patterns. Following is the generation procedure: (i) **Random NFA construction:** Random NFA with 5 – 8 states, 2 – 3 accepting states, and random transition density 0.3 – 0.6. The nondeterministic transition function $\delta : Q \times \Sigma \rightarrow 2^Q$ permits multiple successor states per input symbol. (ii) **Language derivation via Arden's theorem:** Apply Arden's theorem (Sipser, 2013) – if regular expressions P and Q satisfy $R = Q + RP$ and P does not contain ϵ , then $R = QP^*$ – to eliminate states iteratively and derive a closed-form RE. (iii) **Minimal DFA construction:** Convert the NFA to a DFA via subset construction to obtain the ground-truth solution.

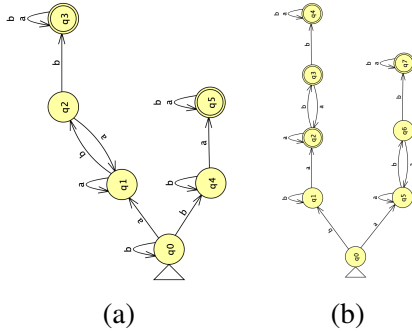


Figure 2: (a): Random NFA over $\{a, b\}$. (b): Minimal DFA recognizing the derived language L_3 .

Detailed Example: Deriving Language L_3 . We now demonstrate the complete derivation process using the random NFA shown in Figure 2(a).

Step 1: State equations from NFA. Each state's language is expressed recursively based on incoming transitions of Figure 2(a):

$$q_5 = q_5a + q_5b + q_4a \quad (1)$$

$$q_4 = q_0b + q_4b \quad (2)$$

$$q_3 = q_3a + q_3b + q_2b \quad (3)$$

$$q_2 = q_1b \quad (4)$$

$$q_1 = q_0a + q_1a + q_2a \quad (5)$$

$$q_0 = \epsilon + q_0b \quad (6)$$

Here q_0 is the start state (accepting ϵ), and q_3, q_5 are accepting states. Our goal is to derive closed-

form REs for q_3 and q_5 in terms of the input alphabet $\{a, b\}$ only.

Step 2: Solve for q_0 using Arden's theorem. Equation (6) has the form $R = Q + RP$ with $R = q_0$, $Q = \epsilon$, and $P = b$. Since b does not contain ϵ , Arden's theorem yields: $q_0 = \epsilon \cdot b^* = b^*$

Step 3: Solve for q_1 . Substitute $q_2 = q_1b$ from Equation (4) into Equation (5):

$$\begin{aligned} q_1 &= q_0a + q_1a + q_2a = q_0a + q_1a + (q_1b)a \\ &= q_0a + q_1(a + ba) \end{aligned}$$

Applying Arden's theorem with $R = q_1$, $Q = q_0a$, $P = (a + ba)$: $q_1 = q_0a(a + ba)^* = b^*a(a + ba)^*$.

Following a similar application of Arden's theorem, we get the following REs for the final states: $q_3 = b^*a(a + ba)^*bb(a + b)^*$ and $q_5 = b^*bb^*a(a + b)^*$

Final language: Therefore, the language recognized by the NFA is: $L_3 = q_3 \cup q_5$. Hence,

$$L_3 = b^*a(a + ba)^*bb(a + b)^* + b^*bb^*a(a + b)^*$$

The resulting expression exhibits deep nesting (concatenations of unions containing Kleene stars) and non-standard structure unlikely to match textbook examples. Figure 2(b) shows the minimal DFA for L_3 , obtained by converting the NFA via subset construction and state minimization. Using this procedure, we generated 120 problems spanning a range of structural complexities. Each random NFA seed produces a distinct RE, ensuring comprehensive coverage of expression patterns while maintaining verifiable correctness through algorithmic DFA construction¹.

3.3 Dataset Statistics and Difficulty level

We categorize all 180 unseen problems by difficulty based on structural complexity: (i) **Easy (41 problems, 22.8%)**: Simple expressions with shallow nesting (e.g., $L_4 = b(a + b)^*ab - 1$ Kleene star, 1 union, 3 concatenations, minimal DFA has ≤ 5 states). (ii) **Medium (65 problems, 36.1%)**: Moderate compositionality with 2 – 4 nested operators (e.g., $L_5 = (a^*b + ((a^*a + a^*b)b^*)a + (((a^*a + a^*b)b^*)b^*)b)a^*$; minimal DFA has 6–12 states). (iii) **Hard (74 problems, 41.1%)**: Deep nesting requiring global consistency tracking (e.g., $L_6 = (((a^*b)(a + b((\epsilon + a)(ba)^*(a + bb) + b)b^*a)^*)b((\epsilon + a)(ba)^*(a + bb) + b))b^*$; minimal DFA has ≥ 13 states). For manually composed

¹All datasets, file formats, and additional implementation details are provided in Appendix B.

Dataset Component	Count	Easy	Med	Hard
Unseen construction	180	41	65	74
Manual composition	60	15	22	23
Arden generation	120	26	43	51

Table 2: Dataset statistics and difficulty categorization.

problems, difficulty reflects the number of interacting constraints (1–2 for easy, 3–4 for medium, 5+ for hard). For Arden-generated problems, difficulty corresponds to RE nesting depth and the number of states in the DFA. Table 2 summarizes the complete benchmark. The difficulty distribution ensures comprehensive evaluation across complexity levels, with a slight bias toward hard problems (41%) to stress-test reasoning capabilities, for detailed classification see Appendix C.

Validation of Novelty. To ensure problems are truly unseen, we: (1) manually searched for exact and structurally similar matches in top Google results and online repositories, (2) verified that Arden-generated regular expressions do not match patterns in standard textbooks, and (3) confirmed that manual problems combine constraints in non-standard ways. While we cannot guarantee complete absence from all pretraining corpora, the controlled generation process and verification steps provide strong evidence of novelty.

4 Evaluation Methodology

Here, all models are accessed via official APIs with temperature set to 0 for deterministic decoding, and all experiments use identical prompt templates to ensure fair comparison². We evaluate five prompting configurations spanning zero-shot, explicit reasoning, multi-branch exploration, and guided self-correction:

Direct Input-Output (Zero-Shot). Models receive only the problem specification (RE or natural language) and must output a complete DFA with no intermediate reasoning or worked examples (Kojima et al., 2022). This baseline isolates pure construction capability without scaffolding.

Chain-of-Thought (CoT). Models are instructed to reason step-by-step before producing the final DFA (Wang and Zhou, 2024). For DFA construction, this naturally decomposes into: (1) interpreting the language specification, (2) identifying required states and their semantic roles, and (3) de-

²Additional details regarding determinism controls, decoding settings, and the runtime environment are provided in Appendices I and J.

signing transitions ensuring correct acceptance/rejection. For example, given $L_7 = (a + b)^*a(a + b)^*$, CoT prompting encourages models to explicitly reason: “State q_0 tracks strings without ‘a’; state q_1 (accepting) tracks strings with at least one ‘a’; transition $\delta(q_0, a) = q_1$ handles first ‘a’; etc.”.

Chain-of-Thought (One-Shot). Extends CoT by providing a single worked example (problem + solution + reasoning trace) before the target problem. This tests whether explicit demonstration reduces ambiguity in reasoning structure.

Tree-of-Thought (ToT). Decomposes DFA construction into four distinct reasoning branches corresponding to standard automata-theoretic methods (Yao et al., 2023). Models are prompted to explore multiple construction approaches in parallel: (i) **Direct (Intuitive):** Interpret the regular expression semantically, identify necessary states based on language properties, and design transitions directly. Tests high-level symbolic reasoning without algorithmic scaffolding. (ii) **Minimization-Based:** Construct an initial DFA (possibly with redundant states), then merge equivalent states via partition refinement. Tests whether models understand state equivalence beyond local transition correctness. (iii) **Derivative-Based:** Apply Brzozowski’s method (Sipser, 2013), where each state corresponds to a distinct derivative of the regular expression with respect to input prefixes. Tests symbolic algebraic manipulation and normalization of equivalent expressions. (iv) **Thompson’s Construction:** Follow the algorithmic pipeline: convert regular expression to ϵ -NFA via structural recursion, then determinize via subset construction. Tests procedural correctness on multi-stage formal algorithms. Lexical analyser tools Lex, Flex follow this approach to automatically create DFA from RE with limitations (Aho et al., 2006). These branches span fundamentally different reasoning styles (semantic interpretation, equivalence optimization, algebraic manipulation, algorithmic execution), enabling fine-grained diagnosis of where models succeed or fail. All prompt templates are provided in Appendix D.

Hint-Based Self-Correction Protocol. For problems answered incorrectly under direct prompting, we evaluate whether models can self-correct when provided structured feedback. We introduce a three-stage hint protocol with increasing levels of guidance: **Stage 1: Counterexample Feedback.** The model receives concrete counterexamples exposing errors in its DFA: **False negatives:**

Strings in L that the DFA rejects; **False positives:** Strings not in L that the DFA accepts. Counterexamples are selected to cover all major structural paths through the automaton, ensuring errors cannot be fixed by local patches. **Stage 2: Error Localization.** If errors persist, the model is informed *which reasoning stage* is incorrect (language interpretation, state design, or transition) along with a high-level description (e.g., ‘Your DFA incorrectly simplifies the Kleene star constraint’). **Stage 3: Explicit Error Disclosure.** If the model still fails, it receives the exact error location and nature (e.g., ‘Missing transition: $\delta(q_2, a) \rightarrow q_4$ to handle *aba*’). This progressive disclosure tests whether failures stem from shallow misinterpretation (correctable with counterexamples), incomplete reasoning (correctable with localization), or fundamental inability to maintain symbolic consistency (uncorrectable even with explicit guidance), Appendix D contains hint templates.

Validation and Grading. All DFA outputs are validated through a two-stage pipeline combining automated behavioural testing and manual structural inspection: (i) **Automated Validation:** A custom validator checks: (a) *syntactic correctness*: valid JSON schema, well-formed state/transition structure; (b) *totality*: every state has exactly one outgoing transition per input symbol; and (c) *behavioural equivalence*: the DFA accepts the same language as the ground-truth RE. Behavioural equivalence is tested via exhaustive enumeration of all strings up to length 6 and random sampling of 2000 strings of length 7 – 15. Any mismatch triggers explicit counterexample logging. (ii) **Manual Inspection:** For all DFAs, at least two independent reviewers from the research team manually verify: (a) correct interpretation of the language specification, (b) meaningful state semantics aligned with the RE, and (c) valid transition logic. All grading is conducted in a model blinded manner to prevent bias. A DFA is marked correct if and only if it passes both tests (see Appendices E and F for validation codes, evaluation scripts and execution scripts).

5 Results and Analysis

Results. Table 3 reports DFA construction success rates across datasets, prompting strategies, and models. (i) **Seen vs. Unseen:** In contrast with the results of seen dataset (Section 2.2), performance drops sharply on the *Unseen DFA Con-*

Dataset	Prompting	GPT-5.1	Grok-4.1	Gemini-2.5
Seen DFA	Zero-Shot	84.2%	89.5%	85.0%
Unseen DFA	Zero-Shot	20.67%	59.12%	29.33%
	CoT	16.67%	51.90%	23.33%
	CoT (One-Shot)	19.33%	55.87%	28.10%
	ToT	24.10%	–	54.00%

Table 3: DFA construction success rates across datasets, prompting strategies, and models. “–” denotes inference-time timeouts.

struction dataset across all models. Under direct prompting, Grok-4.1-fast-reasoning performs best, followed by Gemini-2.5-Flash and GPT-5.1. This corresponds to a degradation of 30–64 % relative to seen tasks. Crucially, the task formulation is identical across seen and unseen settings. The observed performance gap therefore isolates a failure of *compositional generalization* rather than task understanding. (ii) **CoT Effects:** CoT consistently degrades performance on unseen DFA construction across all models. Here, explicit step-by-step reasoning increases effective state-space complexity and amplifies the impact of early semantic errors. To assess whether examples mitigate this failure mode, we evaluate *CoT (One-Shot)* prompting. Providing a single worked example improves performance relative to standard CoT. (iii) **ToT Effects:** ToT yields the strongest performance on complex unseen instances. Gemini-2.5, GPT-5.1 improve by 24.67, 3.43 % over direct prompting and 30.67, 7.43 % over CoT prompting respectively. In contrast, Grok-4.1-fast-reasoning frequently fails to return an output within the fixed inference-time budget³.

Analysis. Across models and prompting strategies, the dominant failure modes arise from difficulties in maintaining *globally consistent symbolic structure*, rather than from a lack of formal knowledge. To illustrate, we analyze following representative outputs under the four ToT framework methods (which subsumes Direct and CoT). (i) **Thompson Construction:** Here, LLMs generally succeed on moderately complex REs with limited nesting because this method is fully algorithmic. However, for complex expressions involving deep nesting, the intermediate ϵ -NFA grows rapidly in size. This *state explosion* acts as a reason behind the failure of LLMs. Importantly, these failures do not reflect a misunderstanding of method, but rather a limitation in reliably managing large,

³A detailed analysis of timeout behaviour, inference cost, and deployability trade-offs is provided in Section 7.

densely connected intermediate structures. (ii) **Direct Construction:** Here, LLMs attempt to construct DFAs by informally reasoning about the language constraints. This setting exhibits following recurring failure modes. (a) *Partial or oversimplified language interpretation.* To produce a correct DFA, the first most logically important step demands correct understanding about the language. However, there are many evidences where LLMs lack in this step (e.g., for $L_8 = a^*(a+b)(a+b)^*a$, LLM develops wrong understanding ‘ending with a’, ignoring presence of $(a+b)$). The same is also applicable for ‘Kleene star’. (b) *Inconsistency between a model’s stated understanding and its final construction.* (e.g., starts with correct understanding ‘strings ending with ab ’, but returns answer for ‘ending with ba ’). It indicates difficulty in maintaining symbolic invariants across construction. (c) *Failure to preserve constraints under concatenation.* (e.g., fails to handle L_8 , but capable to handle unit components, namely a^* , $(a+b)$, $(a+b)^*$, and a). (d) *Over-acceptance of strings outside the target language.* Models often validate their constructions using a finite set of accepted strings while ignoring rejection cases. In addition, LLMs often introduce unreachable or redundant states, which further invalidate the answer ⁴. (iii) **Minimization-Based Construction.** Here, most errors originate before minimization is applied and propagates during the algorithmic minimization approach. (iv) **Derivative-Based Construction.** LLMs frequently fail to normalize semantically equivalent RE derivatives, leading to the creation of spurious DFA states. (e.g., for $L_9 = b(a+b)^*ab^5$ the derivative $(a+b)^*ab + \varepsilon$ is a correct intermediate step reflecting nullability. However, LLMs frequently replace ε with $(a+b)^*$, yielding the incorrect intermediate RE $(a+b)^*ab + (a+b)^*$. This RE normalizes to $(a+b)^*$, thus discards the original language constraints). Importantly, this failure shows an inability to perform basic algebraic simplification and language containment reasoning.

Hint-Based Protocol Analysis. The hint-based protocol is designed to evaluate whether LLMs can recover from incorrect DFA constructions when provided with structured guidance. Table 4 summarizes model performance across difficulty levels and hint stages. (i) **Easy:** Most errors are

Stage	GPT-5.1			Grok-4.1			Gemini-2.5		
	Easy	Medium	Difficult	Easy	Medium	Difficult	Easy	Medium	Difficult
First Hint	60.13%	10.13%	30.76%	85.25%	39.20%	52.93%	57.45%	49.13%	31.79%
Second Hint	15.50%	16.10%	15.38%	14.75%	27.67%	7.83%	42.55%	12.47%	6.88%
Final Hint	24.37%	36.17%	-	-	13.33%	-	-	12.37%	16.31%
Not Solved	-	37.60%	53.86%	-	19.80%	39.24%	-	26.03%	45.02%

Table 4: Hint-based DFA construction performance.

resolved after the *first hint* across all models, indicating that failures at this level primarily arise from superficial misinterpretations rather than fundamental reasoning errors. (ii) **Medium:** Here, recovery is distributed across multiple hint stages rather than concentrated at a single level of guidance. Moreover, a non-trivial fraction remain incorrect even after the final hint. (iii) **Difficult:** A large proportion of problems remain unsolved even after all hints are provided, particularly for GPT-5.1 and Gemini-2.5-Flash. The persistence of unsolved cases suggests that, for complex REs, guidance is insufficient to overcome underlying limitations in symbolic state-transition tracking, and global automaton structure. Overall, the hint-based analysis demonstrates that while all models can leverage guidance to correct shallow errors, their ability to recover degrades rapidly as DFA complexity increases. Hints primarily assist in resolving local misunderstandings, but they do not reliably enable correction of globally inconsistent or structurally flawed automata.

6 Concluding Remarks

This work presents a systematic evaluation of LLMs on DFA construction tasks, with an emphasis on symbolic correctness rather than surface-level pattern matching. We find that while current LLMs reliably reproduce familiar constructions, their performance degrades substantially on compositionally complex inputs. Our analysis shows that these failures are systematic rather than random. These limitations persist across prompting strategies indicating structural weaknesses in symbolic reasoning rather than prompt-specific shortcomings. The hint-based framework demonstrates that limited self-correction is possible for simpler errors; however, deeper semantic inconsistencies are rarely resolved. By introducing a principled benchmark with controlled novelty and by identifying recurring failure modes, this work provides a foundation for future research on improving symbolic robustness through improved evaluation methodologies, targeted training, or hybrid symbolic–neural approaches.

⁴Appendix G contains a gallery of representative errors.

⁵see Appendix H for model-generated output.

718	7 Limitations	
719	This study has several limitations that should be	
720	considered when interpreting the results.	
721	Model access and evaluation constraints. All	
722	experiments were conducted via public APIs or	
723	standard model interfaces with fixed inference-	
724	time budgets. As a result, some models—most	
725	notably Grok-4.1-fast-reasoning under Tree-of-	
726	Thought prompting—frequently failed to return	
727	outputs within the allotted budget. Such cases	
728	were treated as unsuccessful attempts, reflecting	
729	practical interface and deployability constraints	
730	rather than definitive limitations of the underlying	
731	reasoning capabilities.	
732	Token and state-space scalability. DFA con-	
733	struction from complex regular expressions often	
734	induces large intermediate representations, par-	
735	ticularly under Chain-of-Thought and Tree-of-	
736	Thought prompting. Token budget exhaustion and	
737	implicit state-space explosion limit the reliability	
738	of these methods for highly nested or composi-	
739	tional expressions, even when the underlying con-	
740	struction is theoretically well defined.	
741	Finite validation horizon. Behavioral equiva-	
742	lence between generated DFAs and ground-truth	
743	regular expressions is tested using exhaustive enu-	
744	meration only up to a bounded string length,	
745	supplemented by randomized testing for longer	
746	strings. While this approach provides strong em-	
747	pirical assurance, it does not constitute a formal	
748	proof of equivalence for all possible strings. As	
749	a result, rare counterexamples beyond the tested	
750	length range may remain undetected, potentially	
751	leading to overestimation of construction correct-	
752	ness.	
753	Dataset format constraints. All DFA construc-	
754	tion datasets are released in PDF format, contain-	
755	ing problem statements, reference transition ta-	
756	bles, and DFA diagrams. This design choice pri-	
757	oritizes human interpretability and preserves ped-	
758	agogical structure, but it requires manual or semi-	
759	automated parsing for downstream reuse and lim-	
760	its direct machine-readability.	
761	Scope of tasks and formalisms. The experi-	
762	ments focus exclusively on regular expressions,	
763	DFA/NFA construction, and Arden’s theorem-	
764	-based inverse synthesis. Consequently, the re-	
765	ported results may not generalize to richer au-	
766	tomata models, non-regular languages, or other	
	formal systems. Accordingly, our findings should	767
	not be generalized to richer formalisms such as	768
	context-free grammars, pushdown automata, or	769
	Turing-complete models.	770
	Model variability across interfaces. Identical	771
	prompts issued via APIs and web interfaces	772
	can yield different outputs due to undocumented	773
	system-level differences. All reported results cor-	774
	respond to API-based evaluations and should not	775
	be interpreted as exact replicas of web interface	776
	behavior.	777
	Absence of fine-tuning or adaptive prompting.	778
	No model-specific fine-tuning, adaptive decoding,	779
	or prompt optimization was applied beyond the	780
	fixed prompting protocols described in the paper.	781
	While this ensures controlled comparability across	782
	models, it may underestimate achievable perfor-	783
	mance under specialized or tuned evaluation set-	784
	tings.	785
	These limitations primarily affect the external	786
	validity of our conclusions, but do not alter the ob-	787
	served relative performance trends across models	788
	and prompting strategies.	789
	8 Ethical Considerations	790
	This work evaluates the capability of large	791
	language models (LLMs) to perform formal	792
	language-theoretic reasoning tasks, specifically	793
	deterministic finite automaton (DFA) construction	794
	and related inverse problems. The study does not	795
	involve human subjects, personal data, or sensitive	796
	information.	797
	Dataset sourcing. The datasets used in this	798
	study consist exclusively of symbolic, mathemati-	799
	cal objects, including regular expressions, DFAs,	800
	NFAs, transition tables, and multiple-choice or	801
	true/false questions.	802
	The Knowledge Checking dataset and the Seen	803
	DFA Construction dataset were sourced from pub-	804
	licly available educational materials and standard	805
	automata theory problem sets accessible online.	806
	These materials are commonly used for teaching	807
	and assessment and do not carry restrictive li-	808
	censes.	809
	The Unseen DFA Construction dataset, includ-	810
	ing the Mathematical Engineering and Mathemat-	811
	ical Art subsets, was manually constructed by the	812
	authors and was not sourced from online reposi-	813
	tries. All datasets are released with an explicit data	814
	license in the accompanying repository.	815

All newly created datasets are released under the Creative Commons Attribution 4.0 International (CC BY 4.0) license. Publicly sourced educational datasets are provided for research and evaluation purposes only, with original copyrights retained by their respective owners. Full license text and terms of use are provided in the accompanying repository.

Model usage and evaluation. All evaluated models were accessed exclusively through their official public APIs under standard usage conditions. No attempts were made to extract proprietary model internals, training data, or hidden reasoning traces. Prompts were designed solely to elicit task performance and did not request personal, sensitive, or copyrighted information.

Reproducibility and transparency. To promote transparency and responsible research practice, we release the full set of prompts, experimental scripts, validation code, and representative outputs. Automated validation and targeted human inspection were used to reduce the risk of reporting spurious or misleading results.

Potential risks and misuse. The tasks studied are purely formal and mathematical in nature, and the released artifacts are not expected to enable harmful applications or deployment-facing systems. A plausible risk lies in misinterpretation or overgeneralization of the results, for example if observed failures or successes are taken as evidence of general reasoning ability or deficiency across broader linguistic or real-world domains. We explicitly limit our claims to regular-language formalisms and symbolic automata construction and caution against extrapolation beyond these settings.

Use of AI assistants. AI assistants were used in a limited capacity during the development and debugging of auxiliary code. They were not used to generate datasets, ground-truth solutions, experimental results, analyses, or conclusions, and did not contribute to the scientific claims of this work.

Overall, we identify no significant ethical risks associated with this study and adhere to ACL guidelines on data sourcing, transparency, and responsible reporting.

References

- Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. 2006. *Compilers: Principles, Techniques, and Tools (2nd Edition)*. Addison-Wesley Longman Publishing Co., Inc., USA.
- Zhangir Azerbayev, Bartosz Piotrowski, Hailey Schoelkopf, Edward W. Ayers, Dragomir Radev, and Jeremy Avigad. 2023. [Proofnet: Autoformalizing and formally proving undergraduate-level mathematics](#). *Preprint*, arXiv:2302.12433.
- Matei A. Golesteanu, Garrett B. Vowinkel, and Ryan E. Dougherty. 2024. Can ChatGPT pass a theory of computing course? In *Proceedings of the ACM Virtual Global Computing Education Conference*, pages 33–38. Association for Computing Machinery.
- John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. 2001. *Introduction to Automata Theory, Languages, and Computation*, 2 edition. Addison-Wesley.
- Michael Katz, Harsha Kokel, and Sarath Sreedharan. 2025. [Seemingly simple planning problems are computationally challenging: The countdown game](#). *Preprint*, arXiv:2508.02900.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. [Large language models are zero-shot reasoners](#). *Preprint*, arXiv:2205.11916.
- Aitor Lewkowycz, Anders Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, Yuhuai Wu, Behnam Neyshabur, Guy Gur-Ari, and Vedant Misra. 2022. Solving quantitative reasoning problems with language models. In *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS)*.
- Michael Sipser. 2013. *Introduction to the Theory of Computation*, 3 edition. Course Technology.
- Xuezhi Wang and Denny Zhou. 2024. [Chain-of-thought reasoning without prompting](#). *Preprint*, arXiv:2402.10200.
- Sean Welleck, Jiacheng Liu, Ronan Le Bras, Hannaneh Hajishirzi, Yejin Choi, and Kyunghyun Cho. 2021. [Naturalproofs: Mathematical theorem proving in natural language](#). *Preprint*, arXiv:2104.01112.
- Yuhuai Wu, Albert Q. Jiang, Wenda Li, Markus N. Rabe, Charles Staats, Mateja Jamnik, and Christian Szegedy. 2022. Autoformalization with large language models. In *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS)*.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate

916 problem solving with large language models. In
 917 *Proceedings of the 37th Conference on Neural In-*
 918 *formation Processing Systems (NeurIPS)*.

919 Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng,
 920 Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu
 921 Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao
 922 Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan
 923 Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang,
 924 and 3 others. 2024. MMMU: A massive multi-
 925 discipline multimodal understanding and reasoning
 926 benchmark for expert AGI. In *Proceedings of the*
 927 *IEEE/CVF Conference on Computer Vision and Pat-*
 928 *tern Recognition (CVPR)*.

929 Pardis Sadat Zahraei and Ehsaneddin Asgari. 2024.
 930 **Turingq: Benchmarking AI comprehension in the-**
 931 **ory of computation.** *Preprint*, arXiv:2410.06547.

A Appendix : Model-Specific DFA Outputs 932

This appendix reports the exact DFA constructions 934
 produced by each evaluated model for the seen 935
 task (L_1) and the unseen task (L_2). All outputs are 936
 shown exactly as returned by the models, without 937
 post-processing or correction. 938

A.1 Seen Problem (L_1) 939

A.1.1 GPT-5.1 Output for L_1 940

```

  941 {
  942   "states": [
  943     "q0",
  944     "q1",
  945     "q2",
  946     "q3",
  947     "q4",
  948     "q5",
  949     "q6",
  950     "q7"
  951   ],
  952   "alphabet": [
  953     "a",
  954     "b"
  955   ],
  956   "start_state": "q0",
  957   "accept_states": [
  958     "q4",
  959     "q5",
  960     "q6",
  961     "q7"
  962   ],
  963   "transitions": {
  964     "q0": {
  965       "a": "q1",
  966       "b": "q0"
  967     },
  968     "q1": {
  969       "a": "q3",
  970       "b": "q2"
  971     },
  972     "q2": {
  973       "a": "q5",
  974       "b": "q4"
  975     },
  976     "q3": {
  977       "a": "q7",
  978       "b": "q6"
  979     },
  980     "q4": {
  981       "a": "q1",
  982       "b": "q0"
  983     },
  984     "q5": {
  985       "a": "q3",
  986       "b": "q2"
  987     },
  988     "q6": {
  989       "a": "q5",
  990       "b": "q4"
  991     },
  992     "q7": {
  993       "a": "q7",
  994       "b": "q6"
  995     }
  996   }
  997 }
  998
  
```

A.1.2 Grok-4-1-fast-reasoning Output for L_1 1000

```

  1001 {
  1002   "states": [
  1003     "q0",
  1004     "q1",
  1005     "q2",
  1006     "q3",
  1007     "q4",
  1008
  
```

```

1009     "q5",
1010     "q6",
1011     "q7"
1012 ],
1013 "alphabet": [
1014   "a",
1015   "b"
1016 ],
1017 "start_state": "q0",
1018 "accept_states": [
1019   "q4",
1020   "q5",
1021   "q6",
1022   "q7"
1023 ],
1024 "transitions": {
1025   "q0": {
1026     "a": "q1",
1027     "b": "q0"
1028   },
1029   "q1": {
1030     "a": "q3",
1031     "b": "q2"
1032   },
1033   "q2": {
1034     "a": "q5",
1035     "b": "q4"
1036   },
1037   "q3": {
1038     "a": "q7",
1039     "b": "q6"
1040   },
1041   "q4": {
1042     "a": "q1",
1043     "b": "q0"
1044   },
1045   "q5": {
1046     "a": "q3",
1047     "b": "q2"
1048   },
1049   "q6": {
1050     "a": "q5",
1051     "b": "q4"
1052   },
1053   "q7": {
1054     "a": "q7",
1055     "b": "q6"
1056   }
1057 }
1058 }

```

A.1.3 Gemini-2.5-flash Output for L_1

```

1060 {
1061   "states": [
1062     "q0",
1063     "q1",
1064     "q2",
1065     "q3",
1066     "q4",
1067     "q5",
1068     "q6",
1069     "q7"
1070 ],
1071 "alphabet": [
1072   "a",
1073   "b"
1074 ],
1075 "start_state": "q0",
1076 "accept_states": [
1077   "q4",
1078   "q5",
1079   "q6",
1080   "q7"
1081 ],
1082 "transitions": {
1083   "q0": {
1084     "a": "q1",
1085     "b": "q0"
1086   },
1087   "q1": {
1088     "a": "q2",
1089     "b": "q3"
1090   },
1091   "q2": {
1092     "a": "q4",

```

```

    "b": "q5"
  },
  "q3": {
    "a": "q6",
    "b": "q7"
  },
  "q4": {
    "a": "q4",
    "b": "q5"
  },
  "q5": {
    "a": "q6",
    "b": "q7"
  },
  "q6": {
    "a": "q2",
    "b": "q3"
  },
  "q7": {
    "a": "q1",
    "b": "q0"
  }
}

```

A.2 Unseen Problem (L_2)

A.2.1 GPT-5.1 Output for L_2

```

{
  "states": [
    "q0",
    "q1",
    "q2",
    "q3",
    "q4",
    "q5",
    "q6",
    "q7",
    "q8",
    "q9",
    "q10",
    "q11",
    "q12",
    "q13",
    "q14",
    "q15",
    "q16",
    "q17",
    "q18",
    "q19",
    "q20",
    "q21",
    "q22",
    "q23",
    "q24",
    "q25",
    "q26",
    "q27",
    "q28",
    "q29",
    "q30",
    "q31",
    "qDead"
  ],
  "alphabet": [
    "a",
    "b"
  ],
  "start_state": "q0",
  "accept_states": [
    "q24",
    "q25",
    "q26",
    "q27",
    "q28",
    "q29",
    "q30",
    "q31"
  ],
  "transitions": {
    "q0": {
      "a": "q1",
      "b": "q2"
    },

```

```

1179 "q1": {
1180   "a": "q3",
1181   "b": "q4"
1182 },
1183 "q2": {
1184   "a": "q1",
1185   "b": "qDead"
1186 },
1187 "q3": {
1188   "a": "q5",
1189   "b": "q6"
1190 },
1191 "q4": {
1192   "a": "q3",
1193   "b": "q7"
1194 },
1195 "q5": {
1196   "a": "q8",
1197   "b": "q9"
1198 },
1199 "q6": {
1200   "a": "q5",
1201   "b": "q10"
1202 },
1203 "q7": {
1204   "a": "q6",
1205   "b": "q11"
1206 },
1207 "q8": {
1208   "a": "q12",
1209   "b": "q13"
1210 },
1211 "q9": {
1212   "a": "q8",
1213   "b": "q14"
1214 },
1215 "q10": {
1216   "a": "q9",
1217   "b": "q15"
1218 },
1219 "q11": {
1220   "a": "q10",
1221   "b": "q16"
1222 },
1223 "q12": {
1224   "a": "q17",
1225   "b": "q18"
1226 },
1227 "q13": {
1228   "a": "q12",
1229   "b": "q19"
1230 },
1231 "q14": {
1232   "a": "q13",
1233   "b": "q20"
1234 },
1235 "q15": {
1236   "a": "q14",
1237   "b": "q21"
1238 },
1239 "q16": {
1240   "a": "q15",
1241   "b": "q22"
1242 },
1243 "q17": {
1244   "a": "q23",
1245   "b": "q24"
1246 },
1247 "q18": {
1248   "a": "q17",
1249   "b": "q25"
1250 },
1251 "q19": {
1252   "a": "q18",
1253   "b": "q26"
1254 },
1255 "q20": {
1256   "a": "q19",
1257   "b": "q27"
1258 },
1259 "q21": {
1260   "a": "q20",
1261   "b": "q28"
1262 },
1263 "q22": {
1264   "a": "q21",
1265   "b": "q29"

```

```

},
"q23": {
  "a": "q23",
  "b": "q30"
},
"q24": {
  "a": "q23",
  "b": "q31"
},
"q25": {
  "a": "q24",
  "b": "q31"
},
"q26": {
  "a": "q25",
  "b": "q31"
},
"q27": {
  "a": "q26",
  "b": "q31"
},
"q28": {
  "a": "q27",
  "b": "q31"
},
"q29": {
  "a": "q28",
  "b": "q31"
},
"q30": {
  "a": "q30",
  "b": "q31"
},
"q31": {
  "a": "q30",
  "b": "q31"
},
"qDead": {
  "a": "qDead",
  "b": "qDead"
}
}

```

```

1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308

```

A.2.2 Grok-4-1-fast-reasoning Output for L_2

```

{
  "states": [
    "q0",
    "q1",
    "q2",
    "q3",
    "q4",
    "q5",
    "q6",
    "q7",
    "q8",
    "q9",
    "q10",
    "q11",
    "q12",
    "q13",
    "q14",
    "q15",
    "q16",
    "q17",
    "q18",
    "q19",
    "q20",
    "q21",
    "q22",
    "q23",
    "q24",
    "q25",
    "q26",
    "q27",
    "q28"
  ],
  "alphabet": [
    "a",
    "b"
  ],
  "start_state": "q0",
  "accept_states": [
    "q13",

```

```

1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351

```

1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438

```
"q14",
"q15",
"q16",
"q17",
"q18",
"q19",
"q20"
],
"transitions": {
  "q0": {
    "a": "q1",
    "b": "q2"
  },
  "q1": {
    "a": "q3",
    "b": "q4"
  },
  "q2": {
    "a": "q5",
    "b": "q12"
  },
  "q3": {
    "a": "q6",
    "b": "q7"
  },
  "q4": {
    "a": "q8",
    "b": "q9"
  },
  "q5": {
    "a": "q10",
    "b": "q11"
  },
  "q6": {
    "a": "q13",
    "b": "q14"
  },
  "q7": {
    "a": "q15",
    "b": "q16"
  },
  "q8": {
    "a": "q17",
    "b": "q18"
  },
  "q9": {
    "a": "q19",
    "b": "q20"
  },
  "q10": {
    "a": "q21",
    "b": "q22"
  },
  "q11": {
    "a": "q23",
    "b": "q24"
  },
  "q12": {
    "a": "q12",
    "b": "q12"
  },
  "q13": {
    "a": "q13",
    "b": "q14"
  },
  "q14": {
    "a": "q15",
    "b": "q16"
  },
  "q15": {
    "a": "q17",
    "b": "q18"
  },
  "q16": {
    "a": "q19",
    "b": "q20"
  },
  "q17": {
    "a": "q21",
    "b": "q22"
  },
  "q18": {
    "a": "q23",
    "b": "q24"
  },
  "q19": {
    "a": "q25",
```

```
"b": "q26"
},
"q20": {
  "a": "q27",
  "b": "q28"
},
"q21": {
  "a": "q13",
  "b": "q14"
},
"q22": {
  "a": "q15",
  "b": "q16"
},
"q23": {
  "a": "q17",
  "b": "q18"
},
"q24": {
  "a": "q19",
  "b": "q20"
},
"q25": {
  "a": "q21",
  "b": "q22"
},
"q26": {
  "a": "q23",
  "b": "q24"
},
"q27": {
  "a": "q25",
  "b": "q26"
},
"q28": {
  "a": "q27",
  "b": "q28"
}
}
```

1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478

A.2.3 Gemini-2.5-flash Output for L_2

```
{
  "states": [
    "q0",
    "q1",
    "q2",
    "q3",
    "q4",
    "q5",
    "q6",
    "q7",
    "q8",
    "q9",
    "q10",
    "q11",
    "q12",
    "q13",
    "q14",
    "q15"
  ],
  "alphabet": [
    "a",
    "b"
  ],
  "start_state": "q0",
  "accept_states": [
    "q7",
    "q8",
    "q9",
    "q10"
  ],
  "transitions": {
    "q0": {
      "a": "q2",
      "b": "q1"
    },
    "q1": {
      "a": "q5",
      "b": "q15"
    },
    "q2": {
      "a": "q3",
      "b": "q4"
    },
```

1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524

1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578

```
"q3": {
  "a": "q7",
  "b": "q8"
},
"q4": {
  "a": "q9",
  "b": "q10"
},
"q5": {
  "a": "q11",
  "b": "q12"
},
"q6": {
  "a": "q13",
  "b": "q14"
},
"q7": {
  "a": "q7",
  "b": "q8"
},
"q8": {
  "a": "q9",
  "b": "q10"
},
"q9": {
  "a": "q11",
  "b": "q12"
},
"q10": {
  "a": "q13",
  "b": "q14"
},
"q11": {
  "a": "q7",
  "b": "q8"
},
"q12": {
  "a": "q9",
  "b": "q10"
},
"q13": {
  "a": "q11",
  "b": "q12"
},
"q14": {
  "a": "q13",
  "b": "q14"
},
"q15": {
  "a": "q15",
  "b": "q15"
}
}
```

1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596

B Appendix : Datasets

This appendix documents the datasets used in our experiments, including their composition, sizes, formats, and public availability. Detailed dataset construction procedures, design motivations, and validation protocols are described in the main paper.

B.1 Overview

We evaluate large language model (LLM) performance using four datasets spanning factual knowledge, seen DFA construction, and two categories of unseen DFA construction tasks. All datasets are grounded in formal language theory and deterministic finite automata. Ground-truth answers are verified through manual construction and automated validation procedures described in the main paper.

B.2 Dataset Composition and Sizes 1597

- **Knowledge Checking Dataset** (50 questions): 1598
1599
This dataset consists of multiple-choice and true/false questions assessing factual knowledge of core Theory of Computation concepts, including regular languages, deterministic and nondeterministic finite automata, and DFA minimization. Each question admits a single unambiguous correct answer. 1600
1601
1602
1603
1604
1605
1606
- **Seen DFA Construction Dataset** (90 questions): 1607
1608
This dataset contains regular expression-to-DFA and language-to-DFA construction problems whose structural patterns closely resemble commonly available textbook exercises and publicly accessible online examples. For each problem, we additionally provide explicit evidence demonstrating the availability of structurally identical or equivalent questions on the public internet, serving as proof of prior exposure. 1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
- **Unseen DFA Construction Dataset** (180 questions): 1619
1620
This dataset is designed to evaluate generalization beyond memorized patterns and is intentionally constructed to avoid overlap with common instructional examples. It is divided into two subsets that differ in their generative principles: 1621
1622
1623
1624
1625
1626
 - **Mathematical Art Subset** (60 questions): 1627
1628
This subset consists of problems generated by combining multiple interacting constraints, including structural restrictions on automata, conditional dependencies between symbols, and narrative-based task formulations. The resulting languages are intentionally highly irregular and non-canonical, designed to stress symbolic consistency, constraint integration, and long-horizon reasoning. 1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
To design the languages in this unseen subset, we employ the following principled construction strategies: 1639
1640
1641
 - (a) **Use of product construction.** Multiple constraints are combined using product-style constructions (e.g., repeated use of conjunctions and disjunctions to increase constraint complexity); for example, 1642
1643
1644
1645
1646
1647

1648	$L_a = \{$ Construct a DFA	$L_e = \{$ Construct a DFA	1700
1649	that accepts the set of	over $\{0,1\}$ that accepts	1701
1650	all strings over $\{a,b\}$	a string w such that	1702
1651	where it starts with	the digit obtained by	1703
1652	either 'aba' or 'bab'	the XOR operation of	1704
1653	and ends with 'bab' or	the first three digits	1705
1654	'aba' respectively. }	is equal to the digit	1706
1655	$L_b = \{$ Construct a DFA	present at the end of	1707
1656	that accepts the set of	the string w . }	1708
1657	all strings over $\{a,b\}$	$L_f = \{$ Construct a DFA	1709
1658	where 'ab' should be	over $\{0,1\}$ such that the	1710
1659	followed by 'ba' and	string length is greater	1711
1660	'aa' should be followed	than or equal to 4, and	1712
1661	by 'bb' and the total	in each substring of	1713
1662	count of 'a's and 'b's	length 4 the first two	1714
1663	is even. }	bits are interpreted	1715
1664	(b) Multiple interacting constraints.	as a decimal number and	1716
1665	Independent constraints (e.g., prefix	the last two bits are	1717
1666	conditions paired with counting or	interpreted as another	1718
1667	exclusion constraints) are imposed	decimal number, and	1719
1668	simultaneously, preventing purely	their multiplication	1720
1669	local reasoning from yielding a cor-	is divisible by 4. }	1721
1670	rect construction; for example,	(d) Narrative-based tasks. Narrative-	1722
1671	$L_c = \{$ Construct a DFA	driven problems are formulated us-	1723
1672	over $\{a,b\}$ that accepts	ing real-world or game-based sce-	1724
1673	all strings in which the	narios, which must be translated into	1725
1674	fourth-last symbol from	precise symbolic automaton con-	1726
1675	the end must be 'a' and	straints; for example,	1727
1676	the substring 'bb' does	$L_g = \{$ In chess, consider	1728
1677	not appear before any	the Ruy-Lopez opening.	1729
1678	'a'. }	States correspond to	1730
1679	$L_d = \{$ Construct a DFA	positions of pieces	1731
1680	over $\{a,b,c\}$ such that	on the board. The	1732
1681	there exists a block of	input alphabet is	1733
1682	A's of length no more	defined as $\Sigma =$	1734
1683	than 5, immediately	$\{pxyi\}$, where p denotes	1735
1684	followed by a block	the player (white or	1736
1685	of B's that is twice	black), x denotes the	1737
1686	the length of the A	piece, y the square	1738
1687	block; the character	alphabet, and i the	1739
1688	C may appear anywhere	square number. A move	1740
1689	before or after these	such as a black knight	1741
1690	two consecutive blocks	moving to c6 is written	1742
1691	with $ C \equiv 0 \pmod{3}$, and	as bkc6. Assuming	1743
1692	no C appears between the	standard chess rules	1744
1693	two AB blocks. }	with alternating turns	1745
1694	(c) Structural restrictions. Structural	and a fixed initial	1746
1695	constraints such as adjacency rules,	board configuration,	1747
1696	forbidden substrings, or positional	construct a minimized	1748
1697	patterns are enforced, requiring pre-	DFA whose accepting	1749
1698	cision tracking of symbol relation-	paths correspond to	1750
1699	ships; for example,	the Ruy-Lopez opening	1751



Figure 3: Initial board configuration before the Ruy-Lopez opening sequence.



Figure 4: Illustration of the Ruy-Lopez opening sequence considered in L_g .

1752 achieved in exactly five
 1753 moves, excluding the two
 1754 initial moves. }
 1755 $L_h = \{$ You are a fruit
 1756 juice manufacturing
 1757 company owner with
 1758 three possible flavors:
 1759 sweet, sour, and
 1760 mild. Five syrups are
 1761 available: 0-Orange,
 1762 1-Lemon, 2-Mango,
 1763 3-Banana, and 4-Grape.
 1764 The flavor combinations
 1765 are defined as follows:
 1766 Sweet = Mango + Grape;
 1767 Sour = Orange + Lemon +
 1768 Grape; Mild = Banana +
 1769 Mango. Construct a DFA

1770 that accepts exactly
 1771 those strings that
 1772 constitute a particular
 1773 flavor according to
 1774 these rules. }

1775 Across these strategies, we include con-
 1776 crete instances such as maze-solving
 1777 analogies, real-world scenario encod-
 1778 ings, and composite logical require-
 1779 ments to illustrate their effect on DFA
 1780 construction difficulty. Each strategy
 1781 is represented by multiple example lan-
 1782 guages in the dataset.

1783 – **Mathematical Engineering Subset**
 1784 (120 questions):

1785 This subset is constructed using Ar-
 1786 den’s theorem as a core generative
 1787 tool. Regular expressions are derived
 1788 from systematically designed DFA
 1789 and NFA transition structures, yielding
 1790 symbolically precise but nonstandard
 1791 DFA construction problems that em-
 1792 phasize algebraic manipulation and
 1793 formal reasoning rather than surface
 1794 familiarity.

1795 **B.3 Data Format**

1796 All DFA construction datasets (seen and unseen)
 1797 are provided in **PDF format**. Each problem in-
 1798 stance includes:

- 1799 • The problem statement (regular expression or
 1800 formal language specification),
- 1801 • The corresponding DFA state transition table,
- 1802 • A visual DFA diagram used for reference and
 1803 verification.

1804 The Knowledge Checking dataset is provided in
 1805 the same PDF format and includes both questions
 1806 and explicitly labeled correct answers.

1807 **B.4 Ground Truth and Validation**

1808 All ground-truth DFAs and corresponding answers
 1809 were manually constructed by the authors and val-
 1810 idated through formal reasoning and automated
 1811 behavioral equivalence testing, as described in the
 1812 main paper. Dataset construction and validation
 1813 were conducted over a period of more than three
 1814 months in the context of advanced Theory of Com-
 1815 putation coursework. The process can be summa-
 1816 rized as follows:

- 1817 • **Course Context:** Dataset development was car-
 1818 ried out within two advanced Theory of Compu-
 1819 tation courses:

- **CSE 406 (Theory of Computation):** 193 undergraduate students.
- **CSE 525 (Advanced Theory of Computation):** 107 students, comprising both undergraduate and graduate students.

In total, instructional materials involved contributions across multiple course offerings from **300 students**.

- **Nature of Student Contributions:** Student contributions occurred strictly within the normal scope of coursework and assessment activities. No student-generated data was collected, analyzed, or annotated specifically for the purposes of this research, and students were not treated as research participants.
- **Instructional Review and Validation:** Instructional materials underwent systematic review and validation involving:
 - **Five teaching assistants** responsible for intermediate checking, consistency verification, and instructional review.
 - **Two faculty members** who performed final checking and validation.
 - A dedicated **team of research assistants** supporting curation, standardization, and cross-verification.

All final problem instances and ground-truth solutions included in the datasets were subsequently curated, standardized, and independently verified by the authors to ensure correctness, consistency, and suitability for research use.

For the Seen DFA Construction dataset, the repository additionally includes supporting documentation demonstrating the presence of structurally equivalent questions in publicly accessible online sources. This serves as explicit evidence of prior availability and confirms that these instances were not newly introduced for this work.

B.5 Availability and Licensing

All datasets, prompts, experimental scripts, validation code, evidence of internet availability for seen questions, and representative sample outputs are publicly released at:

<https://anonymous.4open.science/r/dfa-llm-evaluation-B82D/>

The repository includes:

- Dataset PDFs for all four datasets,
- Machine-readable metadata files describing dataset properties,
- Prompt templates and experimental scripts,

- Automated validation and evaluation code,
- Proof of public availability for seen dataset questions,
- Sample model outputs and validation reports.

All datasets are released under an explicit data license included in the repository, permitting research use and reproducibility in accordance with ACL guidelines.

C Appendix : Difficulty Labeling Criteria

C.1 Mathematical Engineering Dataset

For the regular-expression-based dataset, difficulty labels are assigned using a quantitative *structural complexity score* derived directly from the syntax of each regular expression. The score combines the following five measurable properties:

- *Maximum nesting depth*, capturing hierarchical structure and long-range dependencies;
- *Number of union operators* (`|`), reflecting branching and nondeterminism;
- *Number of Kleene stars* (`*`), indicating unbounded repetition;
- *Implicit concatenations*, representing sequential composition complexity;
- *Expression length*, capturing overall description size (log-scaled).

These features are linearly combined into a single scalar score. Regular expressions are then partitioned into *easy*, *medium*, and *difficult* categories using dataset-wide tertiles of this score. This procedure yields adaptive and reproducible difficulty labels without the use of manually tuned thresholds. The full implementation used in our experiments is shown below.

```
import math
import json

# ----- Regex Analysis -----

def analyze_regex(regex):
    length = len(regex)

    union = 0
    star = 0
    concat = 0

    nesting = 0
    max_nesting = 0

    prev = None

    for c in regex:
        if c == '(':
            nesting += 1
            max_nesting = max(max_nesting,
nesting)
        elif c == ')':
            nesting -= 1
        elif c == '|':
            union += 1
        elif c == '*':
```

```

1932         star += 1
1933         elif prev and prev not in '(|' and c
1934 not in '|)':
1935             concat += 1
1936
1937         prev = c
1938
1939     return {
1940         "length": length,
1941         "union": union,
1942         "star": star,
1943         "concat": concat,
1944         "nesting": max_nesting
1945     }
1946
1947 # ----- Difficulty Score -----
1948
1949 def difficulty_score(m):
1950     return (
1951         1.5 * m["nesting"] +
1952         1.0 * m["union"] +
1953         0.5 * m["star"] +
1954         0.2 * m["concat"] +
1955         math.loglp(m["length"])
1956     )
1957
1958 # ----- Dataset Classification -----
1959
1960 def classify_dataset(entries):
1961     analyzed = []
1962
1963     for e in entries:
1964         metrics = analyze_regex(e["regex"])
1965         score = difficulty_score(metrics)
1966         analyzed.append((e, metrics, score))
1967
1968     scores = sorted(score for _, _, score in
1969 analyzed)
1970
1971     q1 = scores[int(0.33 * len(scores))]
1972     q2 = scores[int(0.66 * len(scores))]
1973
1974     output = []
1975     for e, metrics, score in analyzed:
1976         if score <= q1:
1977             difficulty = "easy"
1978         elif score <= q2:
1979             difficulty = "medium"
1980         else:
1981             difficulty = "difficult"
1982
1983         entry = e.copy()
1984         entry["metrics"] = metrics
1985         entry["score"] = round(score, 2)
1986         entry["difficulty"] = difficulty
1987         output.append(entry)
1988
1989     return output
1990
1991 # ----- Example Usage -----
1992
1993 if __name__ == "__main__":
1994     dataset = [
1995         {
1996             "id": "r001",
1997             "regex":
1998                 "((b*a)|(b*a(a|b)*(a|b)))a*(a|b)a*aa*",
1999             "alphabet": ["a", "b"]
2000         }
2001     ]
2002
2003     classified = classify_dataset(dataset)
2004     print(json.dumps(classified, indent=2))
2005

```

C.2 Mathematical Art Dataset

The *Mathematical Art* dataset consists of natural-language DFA construction tasks characterized by multiple interacting constraints. Unlike regular-expression-based problems, these tasks are not defined by formal syntax alone; difficulty is therefore

assigned based on semantic structure rather than surface form.

Specifically, difficulty labels are assigned as follows:

- **Easy and medium** tasks are distinguished based on the number of explicit constraints and the extent of their interaction (e.g., positional constraints, counting constraints, and forbidden substrings).
- **Difficult** tasks involve narrative-based formulations, such as grid navigation, game scenarios, or real-world analogies. These problems require multi-step interpretation, implicit state tracking, and translation from informal descriptions into formal automata.

This labeling scheme reflects the additional reasoning burden imposed by narrative grounding and constraint integration, which consistently leads to higher error rates across evaluated models.

D Appendix: Prompt Templates

This appendix provides the exact prompts used in all experiments. Prompts are reported verbatim to ensure reproducibility.

D.1 Prompts Used for Regular Expression to DFA Construction

D.1.1 Direct Input–Output Prompt

```

SYSTEM:
You are an expert in formal languages and
automata.

USER:
Given the following regular expression and
alphabet:

REGEX: {{REGEX}}
ALPHABET: {{ALPHABET}}

Task:
Construct a correct deterministic finite
automaton (DFA) that recognizes exactly
the language denoted by the regular expression.

Constraints:
- The DFA must be total (every state has exactly
one transition per symbol).
- Use short state names: q0, q1, q2, ...
- Do NOT include explanations, reasoning,
derivations, or intermediate steps.
- Do NOT include any text outside the JSON object.

OUTPUT ONLY a single JSON object matching EXACTLY
this schema:

{
  "states": ["q0", "q1", "..."],
  "alphabet": ["a", "b", "..."],
  "start_state": "q0",
  "accept_states": ["q1", "..."],
  "transitions": {
    "q0": { "a": "q1", "b": "q0" },
    "q1": { "a": "q1", "b": "q2" }
  }
}

```

2078

D.1.2 Chain-of-Thought (CoT) Prompt

2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140

```

SYSTEM:
You are an expert in formal languages and finite automata. When asked to produce a DFA transition table, THINK step-by-step internally (Chain-of-Thought) but DO NOT output any reasoning. Output only a single valid JSON object that exactly matches the schema described below.

USER:
Given the following regular expression and alphabet:

REGEX: {{REGEX}}
ALPHABET: {{ALPHABET}}

Task:
1. Internally (chain-of-thought) derive a correct deterministic finite automaton (DFA) that recognizes the language of the regular expression.
2. Do not output any intermediate reasoning or explanation.
3. OUTPUT ONLY a single JSON object that contains the DFA transition table and nothing else. The JSON must match the exact schema (keys and types) below.

Required JSON schema (exact keys and types):
{
  "states": ["q0", "q1", ...],
  "alphabet": ["a", "b", ...],
  "start_state": "q0",
  "accept_states": ["q1", ...],
  "transitions": {
    "q0": { "a": "q1", "b": "q0" },
    "q1": { "a": "q1", "b": "q2" }
  }
}

Formatting rules and constraints (must follow exactly):
- The JSON object must contain only the five keys above and nothing else.
- The DFA must be total: for every state and for every symbol in "alphabet", there must be exactly one target state.
- Use short state names like "q0", "q1", "q2", ...
- All values must be valid JSON types.
- Do not include comments, trailing commas, or extra text.
- If construction fails, output the following machine-readable failure JSON:
  {"error":"cannot_construct","reason":"<one-line_reason>"}

Edge-case guidance (zero-shot):
- Output ONLY the JSON object (or the failure JSON above).
- No diagnostic text before or after the JSON.
- The experiment script retries up to three times.

```

2142

D.1.3 Chain-of-Thought One-Shot Prompt

2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162

```

SYSTEM:
You are an expert in formal languages and finite automata. When asked to produce a DFA transition table, THINK step-by-step internally (Chain-of-Thought) but DO NOT output any reasoning. Output only a single valid JSON object that exactly matches the schema described below.

USER:
Below is an example (one-shot) followed by the real input. Read the example carefully it demonstrates exact JSON schema, formatting, and totality rules. After the example, produce the DFA JSON for the real input only.

=== EXAMPLE (one-shot) ===
Regular expression: (a|b)*
Alphabet: ["a","b"]

```

```

Correct DFA (JSON only):
{
  "states": ["q0"],
  "alphabet": ["a", "b"],
  "start_state": "q0",
  "accept_states": ["q0"],
  "transitions": {
    "q0": { "a": "q0", "b": "q0" }
  }
}

=== END EXAMPLE ===

Now the actual input:

REGEX: {{REGEX}}
ALPHABET: {{ALPHABET}}

Task:
1. Internally (chain-of-thought) derive a correct deterministic finite automaton (DFA) that recognizes the language of the regular expression.
2. Do not output any intermediate reasoning or explanation.
3. OUTPUT ONLY a single JSON object that contains the DFA transition table and nothing else. The JSON must match the exact schema (keys and types) below.

Required JSON schema (exact keys and types):
{
  "states": ["q0", "q1", ...],
  "alphabet": ["a", "b", ...],
  "start_state": "q0",
  "accept_states": ["q1", ...],
  "transitions": {
    "q0": { "a": "q1", "b": "q0" },
    "q1": { "a": "q1", "b": "q2" }
  }
}

Formatting rules and constraints (must follow exactly):
- The JSON object must contain only the five keys above and nothing else.
- The DFA must be total: for every state and for every symbol in "alphabet", there must be exactly one target state.
- Use short state names like "q0", "q1", "q2", ... (no spaces, no special characters).
- All values must be valid JSON types (arrays, strings, objects). Do not include comments or trailing commas.
- Do not include any text before or after the JSON object (no backticks, no code fences, no extra explanation).
- If you cannot produce a correct DFA, DO NOT output {}. Instead, attempt to produce a compact but valid DFA. Only if you truly cannot construct any DFA, output a machine-readable failure JSON:
  {"error":"cannot_construct","reason":"<one-line_reason>"}
  (This is allowed but discouraged - prefer producing a compact DFA.)
- If the DFA is large, produce a compact encoding using short state names (q0,q1,q2,...) and include every transition for each symbol.

Edge-case guidance (one-shot):
- The above example demonstrates the EXACT output shape and formatting (including spacing/newlines are not important, but content must be valid JSON).
- The experiment script will retry up to 3 times if the first output is invalid. OUTPUT ONLY the JSON object (or the small error JSON) no additional text.

IMPORTANT: OUTPUT ONLY the JSON object (or the small error JSON). Do not output any other text.

```

2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245

2247 **D.1.4 Tree-of-Thought (ToT) Direct-Branch**
 2248 **Prompt**

2249 SYSTEM:
 2250 You are an expert in formal languages and finite
 2251 automata. When asked to produce a DFA
 2252 transition table, THINK step-by-step
 2253 internally (Chain-of-Thought) but DO NOT
 2254 output any reasoning. Output only a single
 2255 valid JSON object that exactly matches the
 2256 schema described below.
 2257
 2258 USER:
 2259 Given the following regular expression and
 2260 alphabet:
 2261
 2262 REGEX: {{REGEX}}
 2263 ALPHABET: {{ALPHABET}}

2266 Task:
 2267 1. Internally derive a correct deterministic
 2268 finite automaton (DFA) that recognizes
 2269 exactly the language denoted by the regular
 2270 expression.
 2271 2. Use any sound formal reasoning internally, but
 2272 do not output any intermediate steps.
 2273 3. OUTPUT ONLY a single JSON object that contains
 2274 the DFA transition table and nothing else.
 2275
 2276 Required JSON schema (exact keys and types):
 2277 {
 2278 "states": ["q0", "q1", ...],
 2279 "alphabet": ["a", "b", ...],
 2280 "start_state": "q0",
 2281 "accept_states": ["q1", ...],
 2282 "transitions": {
 2283 "q0": { "a": "q1", "b": "q0" },
 2284 "q1": { "a": "q1", "b": "q2" }
 2285 }
 2286 }
 2287
 2288 Formatting rules and constraints:
 2289 - Output only the JSON object, with no extra text.
 2290 - The DFA must be total.
 2291 - Use short state names q0, q1, q2, ...
 2292 - Do not include comments, explanations, or
 2293 trailing commas.
 2294 - If construction is difficult, still output a
 2295 compact but valid DFA rather than an empty
 2296 object.

2298 **D.1.5 Tree-of-Thought (ToT) Minimization**
 2299 **Prompt**

2300 SYSTEM:
 2301 You are an expert in deterministic finite
 2302 automata. Internally aim to construct a
 2303 correct DFA that recognizes the given
 2304 regular expression and is close to minimal
 2305 in the number of states. You may internally
 2306 apply minimization techniques such as state
 2307 equivalence or partition refinement. Do NOT
 2308 output any reasoning.
 2309
 2310 USER:
 2311 Given the following regular expression and
 2312 alphabet:
 2313
 2314 REGEX: {{REGEX}}
 2315 ALPHABET: {{ALPHABET}}

2318 Task:
 2319 1. Internally construct a correct DFA for the
 2320 regular expression.
 2321 2. Internally minimize the DFA if possible.
 2322 3. Ensure the DFA is total.
 2323 4. OUTPUT ONLY the final minimized DFA as a JSON
 2324 object matching the schema below.
 2325
 2326 Required JSON schema:
 2327 {
 2328 "states": ["q0", "q1", ...],
 2329 "alphabet": ["a", "b", ...],
 2330 "start_state": "q0",

```
"accept_states": ["q1", ...],
"transitions": {
  "q0": { "a": "q1", "b": "q0" },
  "q1": { "a": "q1", "b": "q2" }
}
}
```

2331
 2332
 2333
 2334
 2335
 2336
 2337
 2338
 2339
 2340
 2341
 2342

2343 Formatting constraints:
 2344 - Output only the JSON object.
 2345 - No explanations or comments.
 2346 - DFA must be total.

2343 **D.1.6 Tree-of-Thought (ToT) Derivative**
 2344 **Method Prompt**

2345 SYSTEM:
 2346 You are an expert in formal language theory.
 2347 Internally construct a DFA using regular
 2348 expression derivatives (Brzozowski or
 2349 Antimirov derivatives). Each DFA state
 2350 corresponds to a derivative of the original
 2351 regular expression. Do NOT output any
 2352 intermediate reasoning.
 2353
 2354 USER:
 2355 Given the following regular expression and
 2356 alphabet:
 2357
 2358 REGEX: {{REGEX}}
 2359 ALPHABET: {{ALPHABET}}

2361 Task:
 2362 1. Internally compute the set of regular
 2363 expression derivatives with respect to the
 2364 alphabet.
 2365 2. Internally build the corresponding DFA from
 2366 these derivatives.
 2367 3. Ensure the DFA is total.
 2368 4. OUTPUT ONLY the final DFA as a JSON object
 2369 matching the exact schema below.
 2370
 2371 Required JSON schema:
 2372 {
 2373 "states": ["q0", "q1", ...],
 2374 "alphabet": ["a", "b", ...],
 2375 "start_state": "q0",
 2376 "accept_states": ["q1", ...],
 2377 "transitions": {
 2378 "q0": { "a": "q1", "b": "q0" },
 2379 "q1": { "a": "q1", "b": "q2" }
 2380 }
 2381 }
 2382
 2383 Constraints:
 2384 - Output only valid JSON.
 2385 - DFA must be total.
 2386 - No explanatory text.
 2387 - Use compact state naming.

2390 **D.1.7 Tree-of-Thought (ToT) Thompson's**
 2391 **Construction Prompt**

2392 SYSTEM:
 2393 You are an expert in automata theory. Internally
 2394 follow a Thompson-style construction: first
 2395 convert the regular expression to an -NFA,
 2396 then apply subset construction to obtain a
 2397 DFA, and finally ensure the DFA is total. Do
 2398 NOT output any intermediate structures or
 2399 reasoning.
 2400
 2401 USER:
 2402 Given the following regular expression and
 2403 alphabet:
 2404
 2405 REGEX: {{REGEX}}
 2406 ALPHABET: {{ALPHABET}}

2408 Task:
 2409 1. Internally construct an -NFA using Thompsons
 2410 construction.
 2411 2. Internally determinize the NFA using subset
 2412 construction.
 2413

2414 3. Internally make the DFA total by adding a sink
 2415 state if necessary.
 2416 4. OUTPUT ONLY the final DFA as a single JSON
 2417 object matching the schema below.
 2418

2419 Required JSON schema:
 2420 {
 2421 "states": ["q0", "q1", ...],
 2422 "alphabet": ["a", "b", ...],
 2423 "start_state": "q0",
 2424 "accept_states": ["q1", ...],
 2425 "transitions": {
 2426 "q0": { "a": "q1", "b": "q0" },
 2427 "q1": { "a": "q1", "b": "q2" }
 2428 }
 2429 }
 2430

2431 Formatting rules:
 2432 - Output ONLY the JSON object.
 2433 - DFA must be total.
 2434 - No explanations, comments, or extra text.
 2435 - Use short state names q0, q1, q2, ...

2437 **D.2 Hint-Based Framework Prompts**

2438 This section lists the hints used to evaluate
 2439 whether models can correct their own errors un-
 2440 der guided feedback.

2441 **D.2.1 Initial Prompt (No Hints)**

2442 Construct a deterministic finite automaton (DFA)
 2443 that recognizes exactly the following
 2444 language over the alphabet {a, b}:
 2445
 2446 $L = (b)^*(a)(a + ba)^*(bb)(a + b)^* + (b)^*b(b)^*a(a +$
 2447 $b)^*$
 2448
 2449 Requirements:
 2450 - The DFA must be total.
 2451 - Use short state names (q0, q1, q2, ...).
 2452 - Specify the start state, accept states, and
 2453 transitions.
 2454 - Output only the final DFA.
 2455

2457 **D.2.2 Hint 1: Counterexamples protocol**

2458 You are given the following language:
 2459
 2460 $L = (b)^*(a)(a + ba)^*(bb)(a + b)^* + (b)^*b(b)^*a(a +$
 2461 $b)^*$
 2462
 2463 The DFA you previously constructed does not
 2464 correctly recognize this language.
 2465
 2466 In particular, the following strings belong to L
 2467 but are rejected by your DFA:
 2468 - ba
 2469 - bba
 2470 - baa
 2471
 2472 Revise or reconstruct the DFA so that it accepts
 2473 all valid strings in L while continuing to
 2474 reject invalid strings. Output only the
 2475 corrected DFA.
 2476

2478 **D.2.3 Hint 2: Error localization**

2479 You are given the following language:
 2480
 2481 $L = (b)^*(a)(a + ba)^*(bb)(a + b)^* + (b)^*b(b)^*a(a +$
 2482 $b)^*$
 2483
 2484 Your revised DFA still contains an error.
 2485
 2486 At least one transition does not correctly
 2487 reflect the language definition, leading to
 2488 incorrect acceptance or rejection of some
 2489 strings.
 2490

Re-examine the DFA structure and transition assignments, and provide a corrected DFA that recognizes exactly L. Output only the DFA.

D.2.4 Final Hint: Explicit error disclosure

You are given the following language:
 $L = (b)^*(a)(a + ba)^*(bb)(a + b)^* + (b)^*b(b)^*a(a + b)^*$
 The remaining error arises from an incomplete interpretation of the regular expression.
 In particular, the Kleene-star subexpression $(a + ba)^*$ permits arbitrary repetitions and combinations that must be fully captured by the automaton.
 As a consequence of this incomplete interpretation, your DFA incorrectly accepts the string:
 aabab
 Reconstruct the DFA with a correct interpretation of all subexpressions and their interactions. Output only the final corrected DFA.

E Appendix: Evaluation and Validation Code

Automated Validation Pipeline. The automated validation pipeline provides a principled and reproducible mechanism for assessing the correctness of LLM-generated DFAs. It combines (i) strict structural validation, ensuring syntactic correctness, totality, and well-formed transitions, with (ii) behavioral equivalence testing against the ground-truth regular expression. Behavioral validation is performed using exhaustive enumeration of all strings up to a fixed length and is supplemented with randomized sampling of longer strings, enabling detection of both local and global acceptance errors. This design ensures that reported correctness reflects semantic language equivalence rather than surface-level plausibility.

The validator is model-agnostic and deterministic, producing identical results given the same DFA and regular expression. It records explicit counterexamples whenever a mismatch is detected, supporting transparent error analysis and independent verification. All validation artifacts, including per-instance reports and aggregated summaries, are saved to disk to facilitate reproducibility and post-hoc inspection. At the same time, the validation procedure has inherent limitations. Exhaustive testing is bounded by a maximum string length, and while

2491
 2492
 2493
 2494
 2495
 2496
 2497
 2498
 2499
 2500
 2501
 2502
 2503
 2504
 2505
 2506
 2507
 2508
 2509
 2510
 2511
 2512
 2513
 2514
 2515
 2516
 2517
 2518
 2519
 2520
 2521
 2522
 2523
 2524
 2525
 2526
 2527
 2528
 2529
 2530
 2531
 2532
 2533
 2534
 2535
 2536
 2537
 2538
 2539
 2540
 2541
 2542
 2543
 2544
 2545
 2546
 2547
 2548
 2549
 2550
 2551
 2552

Parameter	GPT-5.1	Gemini-2.5	Grok-4.1
API endpoint	Public	Public	Public
Model variant	GPT-5.1	Gemini-2.5-Flash	Grok-4.1-fast-reasoning
Temperature	0.0	0.0	0.0
Max output tokens	4000	4000	4000
Reasoning toggle	Not exposed	Not exposed	Not exposed
Prompt template	Identical	Identical	Identical
Retries	3	3	3
Timeout	120s	120s	120s

Table 5: API parameters used across model providers. No provider-specific reasoning or deliberation modes were enabled.

random sampling extends coverage, formal equivalence between a DFA and a regular expression is undecidable via finite testing alone. As a result, DFAs classified as correct should be interpreted as *likely correct* within the tested bounds rather than formally proven equivalent. Additionally, the validator assumes the correctness of the regular expression semantics as interpreted by the host regex engine, which may differ from theoretical automata semantics in edge cases. These limitations are mitigated in the main evaluation by complementary human expert inspection, but they remain important considerations when interpreting automated results.

E.1 Automated DFA Validation Pipeline

```
#!/usr/bin/env python3
"""
validate_dfa_outputs.py

Automated DFA validation pipeline used in
experiments.

Validation protocol:
1. Schema validation (JSON structure + required
fields).
2. Totality validation (every state has one
transition per symbol).
3. Behavioral equivalence testing against the
ground-truth regular expression:
- Exhaustive enumeration up to
MAX_EXHAUSTIVE_LEN.
- Large-scale randomized testing for longer
strings.

Outputs:
- Per-DFA detailed validation reports (JSON).
- Aggregated summary table (CSV + JSON).

This script is model-agnostic and
dataset-agnostic.
"""

import os
import json
import re
import random
import csv
from itertools import product
from typing import List, Dict, Tuple

# =====
# CONFIGURATION (EDIT HERE)
# =====

TABLES_DIR = "outputs/tables"
RAW_DIR = "outputs/raw"
VALID_DIR = "outputs/validation"
```

```
MAX_EXHAUSTIVE_LEN = 6 # exhaustive
testing for strings of length 0..6
N_RANDOM = 2000 # number of random
test strings
MAX_RANDOM_LEN = 15 # maximum length of
random strings
MAX_COUNTEREXAMPLES = 100 # cap stored
counterexamples per DFA
RANDOM_SEED = 42 # reproducibility

# =====
# INITIALIZATION
# =====

os.makedirs(VALID_DIR, exist_ok=True)
random.seed(RANDOM_SEED)

# =====
# UTILITIES
# =====

def load_json(path: str):
    with open(path, "r", encoding="utf-8") as f:
        return json.load(f)

def save_json(path: str, obj):
    with open(path, "w", encoding="utf-8") as f:
        json.dump(obj, f, indent=2,
ensure_ascii=False)

def is_valid_schema(dfa: Dict) -> Tuple[bool,
str]:
    required = {"states", "alphabet",
"start_state", "accept_states",
"transitions"}
    if not isinstance(dfa, dict):
        return False, "not_json_object"
    if not required.issubset(dfa.keys()):
        return False, "missing_required_keys"
    if not isinstance(dfa["states"], list):
        return False, "states_not_list"
    if not isinstance(dfa["alphabet"], list):
        return False, "alphabet_not_list"
    if dfa["start_state"] not in dfa["states"]:
        return False, "invalid_start_state"
    if not isinstance(dfa["accept_states"], list):
        return False, "accept_states_not_list"
    if not isinstance(dfa["transitions"], dict):
        return False, "transitions_not_dict"
    return True, "ok"

def check_totality(dfa: Dict):
    missing = []
    invalid_targets = []

    for s in dfa["states"]:
        if s not in dfa["transitions"]:
            missing.append({"state": s, "reason":
"no_transition_block"})
            continue
        for a in dfa["alphabet"]:
            if a not in dfa["transitions"][s]:
                missing.append({"state": s,
"symbol": a})
            else:
                tgt = dfa["transitions"][s][a]
                if tgt not in dfa["states"]:
```

2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566

2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610

2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
2638
2639
2640
2641
2642
2643
2644
2645
2646
2647
2648
2649
2650
2651
2652
2653
2654
2655
2656
2657
2658
2659
2660
2661
2662
2663
2664
2665
2666
2667
2668
2669
2670
2671
2672
2673
2674
2675
2676
2677
2678

```

2679     invalid_targets.append({"state": s,
2680     "symbol": a, "target": tgt})
2681
2682     return missing, invalid_targets
2683
2684
2685 def simulate_dfa(dfa: Dict, word: str) -> bool:
2686     cur = dfa["start_state"]
2687     for ch in word:
2688         if ch not in dfa["alphabet"]:
2689             return False
2690         cur = dfa["transitions"].get(cur,
2691         {}).get(ch)
2692         if cur is None:
2693             return False
2694     return cur in dfa["accept_states"]
2695
2696
2697 def generate_exhaustive_strings(alphabet:
2698     List[str], max_len: int) -> List[str]:
2699     strings = [""]
2700     for L in range(1, max_len + 1):
2701         for tup in product(alphabet, repeat=L):
2702             strings.append("".join(tup))
2703     return strings
2704
2705
2706 def generate_random_strings(alphabet: List[str],
2707     n: int, max_len: int) -> List[str]:
2708     out = set()
2709     while len(out) < n:
2710         L = random.randint(0, max_len)
2711         out.add("".join(random.choice(alphabet)
2712         for _ in range(L)))
2713     return list(out)
2714
2715 # =====
2716 # MAIN VALIDATION LOOP
2717 # =====
2718
2719 def main():
2720     table_files = [
2721         os.path.join(TABLES_DIR, f)
2722         for f in os.listdir(TABLES_DIR)
2723         if f.endswith(".json")
2724     ]
2725
2726     summary = []
2727
2728     for table_path in table_files:
2729         base = os.path.splitext(os.path.basename(
2730         table_path))[0]
2731         print(f"[VALIDATING] {base}")
2732
2733         report = {
2734             "id": base,
2735             "schema_valid": False,
2736             "schema_error": None,
2737             "totality_missing": [],
2738             "totality_invalid_targets": [],
2739             "tested_strings": 0,
2740             "counterexamples": [],
2741             "verdict": ""
2742         }
2743
2744         try:
2745             dfa = load_json(table_path)
2746         except Exception as e:
2747             report["schema_error"] =
2748             f"json_load_error: {e}"
2749             report["verdict"] = "invalid_json"
2750             save_json(os.path.join(VALID_DIR,
2751             base + "_report.json"), report)
2752             continue
2753
2754             ok, reason = is_valid_schema(dfa)
2755             report["schema_valid"] = ok
2756             report["schema_error"] = None if ok else
2757             reason
2758
2759             if not ok:
2760                 report["verdict"] = "schema_invalid"
2761                 save_json(os.path.join(VALID_DIR,
2762                 base + "_report.json"), report)
2763                 continue
2764
2765             missing, invalid = check_totality(dfa)
2766             report["totality_missing"] = missing

```

```

report["totality_invalid_targets"] =
invalid

# Locate matching raw file for regex
regex = None
for rf in os.listdir(RAW_DIR):
    if rf.startswith(base):
        raw =
load_json(os.path.join(RAW_DIR, rf))
        regex = raw.get("regex")
        break

if regex is None:
    report["verdict"] =
    "no_ground_truth_regex"
    save_json(os.path.join(VALID_DIR,
    base + "_report.json"), report)
    continue

try:
    pattern = re.compile(f"^{(regex)}$")
except Exception as e:
    report["verdict"] =
    f"regex_compile_error: {e}"
    save_json(os.path.join(VALID_DIR,
    base + "_report.json"), report)
    continue

exhaustive = generate_exhaustive_strings(
    dfa["alphabet"], MAX_EXHAUSTIVE_LEN)
random_tests =
generate_random_strings(dfa["alphabet"],
    N_RANDOM, MAX_RANDOM_LEN)
tests = exhaustive + [s for s in
    random_tests if s not in exhaustive]

for w in tests:
    dfa_accept = simulate_dfa(dfa, w)
    regex_accept = pattern.fullmatch(w)
    is not None
    report["tested_strings"] += 1

    if dfa_accept != regex_accept:
        report["counterexamples"].append({
            "string": w,
            "dfa_accepts": dfa_accept,
            "regex_accepts": regex_accept
        })
        if len(report["counterexamples"])
        >= MAX_COUNTEREXAMPLES:
            break

    if missing or invalid:
        report["verdict"] = "not_total"
    elif report["counterexamples"]:
        report["verdict"] = "incorrect"
    else:
        report["verdict"] = "likely_correct"

    save_json(os.path.join(VALID_DIR, base +
    "_report.json"), report)

    summary.append({
        "id": base,
        "schema_valid": ok,
        "totality_ok": not (missing or
        invalid),
        "tested_strings":
        report["tested_strings"],
        "num_counterexamples":
        len(report["counterexamples"]),
        "verdict": report["verdict"]
    })

# =====
# WRITE SUMMARY TABLES
# =====

with open(os.path.join(VALID_DIR,
    "summary.csv"), "w", newline="",
    encoding="utf-8") as f:
    writer = csv.DictWriter(
        f,
        fieldnames=["id", "schema_valid",
        "totality_ok",
        "tested_strings",
        "num_counterexamples", "verdict"]

```

```

2765
2766
2767
2768
2769
2770
2771
2772
2773
2774
2775
2776
2777
2778
2779
2780
2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807
2808
2809
2810
2811
2812
2813
2814
2815
2816
2817
2818
2819
2820
2821
2822
2823
2824
2825
2826
2827
2828
2829
2830
2831
2832
2833
2834
2835
2836
2837
2838
2839
2840
2841
2842
2843
2844
2845
2846
2847
2848
2849
2850
2851

```

```

)
writer.writeheader()
for row in summary:
    writer.writerow(row)

save_json(os.path.join(VALID_DIR,
    "summary.json"), summary)

print("[DONE] DFA validation completed.")

if __name__ == "__main__":
    main()

```

F Appendix: Experiment Execution Code

Experiment Execution Code. We provide the full experiment execution code used to query language models and collect DFA outputs. The script is model-agnostic and was used uniformly across all evaluated models and prompting strategies, differing only in the model identifier and prompt templates. It enforces deterministic decoding, robust retry logic, strict JSON extraction, and complete logging of raw outputs, parsed DFAs, and metadata. All runs write outputs regardless of success or failure, enabling reproducible analysis of both correctness and failure modes.

Strengths and Limitations. The execution pipeline ensures reproducibility through deterministic decoding and uniform evaluation settings, and robustness through retries and structured output validation. However, it does not attempt to correct or post-process invalid DFAs, and relies on downstream automated and human validation for semantic correctness. Timeout events and malformed outputs are treated as unsuccessful attempts, reflecting practical deployment constraints rather than purely theoretical capability.

F.1 Experiment Execution Code

```

#!/usr/bin/env python3
"""
run_experiment_tot.py

General Tree-of-Thought (ToT) experiment runner
for LLM-based DFA construction.

Features:
- Model-agnostic (OpenAI / API-based LLMs)
- Deterministic decoding (temperature = 0)
- Robust JSON extraction and retry logic
- Strict output logging (raw, parsed, metadata)
- Always writes outputs (success or failure)
- Reproducible and configurable via constants

This script is used uniformly across models and
prompting strategies by
changing only the MODEL identifier and prompt
templates.
"""
import os
import json
import time
import re

```

```

from datetime import datetime
from typing import Optional
from openai import OpenAI, OpenAIError

# =====
# CONFIGURATION (MODEL-AGNOSTIC)
# =====

MODEL = os.environ.get("LLM_MODEL", "gpt-5.1")
TEMPERATURE = 0.0
MAX_TOKENS = 4000

DATASET_PATH = "data/analysis_data.json"

PROMPT_FILES = {
    "direct": "prompts/tot_direct.txt",
    "minimal": "prompts/tot_minimal.txt",
    "derivative": "prompts/tot_derivative.txt",
    "thompson": "prompts/tot_thompson.txt",
}

OUTPUT_ROOT = "outputs/tot_experiments"

RETRIES = 3
RETRY_SLEEP = 1.5
RATE_LIMIT_SLEEP = 0.3

# =====
# INIT
# =====

client = OpenAI()

TIMESTAMP =
    datetime.utcnow().strftime("%Y%m%dT%H%M%SZ")
OUTPUT_RAW_DIR = os.path.join(OUTPUT_ROOT,
    MODEL, "raw")
OUTPUT_TABLE_DIR = os.path.join(OUTPUT_ROOT,
    MODEL, "tables")
OUTPUT_META_DIR = os.path.join(OUTPUT_ROOT,
    MODEL, "meta")

for d in [OUTPUT_RAW_DIR, OUTPUT_TABLE_DIR,
    OUTPUT_META_DIR]:
    os.makedirs(d, exist_ok=True)

PROJECT_ROOT =
    os.path.abspath(os.path.join(os.path.
    dirname(__file__), ".."))

# =====
# UTILITIES
# =====

def now_ts():
    return datetime.utcnow().strftime("%Y%m%dT%H%M%SZ")

def safe_filename(s: str) -> str:
    return re.sub(r"[^0-9A-Za-z._-]", "_",
        s)[:200]

def load_json(path):
    with open(path, "r", encoding="utf-8") as f:
        return json.load(f)

def load_prompt(path):
    with open(path, "r", encoding="utf-8") as f:
        return f.read()

def save_json(path, obj):
    with open(path, "w", encoding="utf-8") as f:
        json.dump(obj, f, indent=2,
            ensure_ascii=False)

def try_parse_json(s: str) -> Optional[dict]:
    if not s:
        return None
    s = re.sub(r"````(?:json)?", "", s,
        flags=re.IGNORECASE).strip("` \n")
    try:
        return json.loads(s)
    except Exception:
        pass
    m = re.search(r"\{[\s\S]+\}", s)
    if m:

```

```

3004     try:
3005         return json.loads(m.group(0))
3006     except Exception:
3007         pass
3008     return None
3009
3010 # =====
3011 # TOKEN DISPATCH
3012 # =====
3013
3014 def completion_kwargs():
3015     if MODEL.startswith("gpt-5"):
3016         return {"max_completion_tokens":
3017                 MAX_TOKENS}
3018     return {"max_tokens": MAX_TOKENS}
3019
3020 # =====
3021 # MODEL ACCESS CHECK
3022 # =====
3023
3024 def check_model_access():
3025     try:
3026         r = client.chat.completions.create(
3027             model=MODEL,
3028             messages=[{"role": "user", "content":
3029                       "ping"}],
3030             temperature=0.0,
3031             **completion_kwargs()
3032         )
3033         actual = getattr(r, "model", None)
3034         if actual is None or not
3035         actual.startswith(MODEL.split("-")[0]):
3036             raise RuntimeError(f"Model mismatch:
3037                               requested={MODEL}, actual={actual}")
3038     except OpenAIError as e:
3039         raise RuntimeError(f"Model access failed:
3040                             {e}")
3041
3042 # =====
3043 # PROMPTING
3044 # =====
3045
3046 def build_messages(template, regex, alphabet):
3047     filled = (
3048         template
3049         .replace("{REGEX}", regex)
3050         .replace("{ALPHABET}",
3051                 json.dumps(alphabet))
3052     )
3053     return [
3054         {
3055             "role": "system",
3056             "content": (
3057                 "You are an expert in formal
3058                 languages and automata. "
3059                 "Internally reason as needed, but
3060                 OUTPUT ONLY the DFA JSON."
3061             )
3062         },
3063         {"role": "user", "content": filled}
3064     ]
3065
3066 def call_model(messages):
3067     r = client.chat.completions.create(
3068         model=MODEL,
3069         messages=messages,
3070         temperature=TEMPERATURE,
3071         **completion_kwargs()
3072     )
3073     text = r.choices[0].message.content or ""
3074     meta = {
3075         "requested_model": MODEL,
3076         "actual_model": getattr(r, "model", None),
3077         "finish_reason":
3078             r.choices[0].finish_reason,
3079         "usage": getattr(r, "usage", None)
3080     }
3081     return text, meta
3082
3083 # =====
3084 # CORE LOOP
3085 # =====
3086
3087 def run_single_branch(entry, branch,
3088                      prompt_template):
3089     regex_id = entry["id"]
3090     regex = entry["regex"]

```

```

alphabet = entry.get("alphabet", [])
ts = now_ts()

messages = build_messages(prompt_template,
                           regex, alphabet)

raw_text, meta, parsed = "", None, None
attempts = 0

while attempts < RETRIES:
    attempts += 1
    try:
        raw_text, meta = call_model(messages)
    except Exception as e:
        meta = {"error": str(e)}
        time.sleep(RETRY_SLEEP)
        continue

    parsed = try_parse_json(raw_text)
    if parsed:
        break

    messages.append({
        "role": "user",
        "content": "OUTPUT ONLY the DFA JSON.
    No explanations."
    })
    time.sleep(RETRY_SLEEP)

base =
safe_filename(f"{regex_id}_{branch}_{ts}")

save_json(os.path.join(OUTPUT_RAW_DIR, base +
".json"), {
    "id": regex_id,
    "branch": branch,
    "regex": regex,
    "alphabet": alphabet,
    "raw_output": raw_text,
    "attempts": attempts,
    "meta": meta,
    "model": MODEL,
    "timestamp": ts
})

save_json(os.path.join(OUTPUT_META_DIR, base
+ "_meta.json"), meta)

table_obj = parsed if parsed else {
    "error": "GENERATION_FAILED",
    "reason": "Invalid or non-JSON output",
    "regex_id": regex_id,
    "branch": branch,
    "model": MODEL,
    "attempts": attempts
}

save_json(os.path.join(OUTPUT_TABLE_DIR, base
+ ".json"), table_obj)

# =====
# MAIN
# =====

def main():
    check_model_access()
    dataset =
        load_json(os.path.join(PROJECT_ROOT,
DATASET_PATH))
    prompts = {
        k: load_prompt(os.path.join(PROJECT_ROOT,
v))
        for k, v in PROMPT_FILES.items()
    }

    for entry in dataset:
        for branch, tmpl in prompts.items():
            run_single_branch(entry, branch, tmpl)
            time.sleep(RATE_LIMIT_SLEEP)

if __name__ == "__main__":
    main()

```

```

3091
3092
3093
3094
3095
3096
3097
3098
3099
3100
3101
3102
3103
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
3132
3133
3134
3135
3136
3137
3138
3139
3140
3141
3142
3143
3144
3145
3146
3147
3148
3149
3150
3151
3152
3153
3154
3155
3156
3157
3158
3159
3160
3161
3162
3163
3164
3165
3166
3167
3168
3169
3170
3171
3172

```

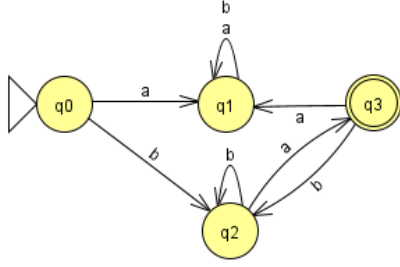


Figure 5: DFA generated via direct construction for $L = b(a | b)^*ab$.

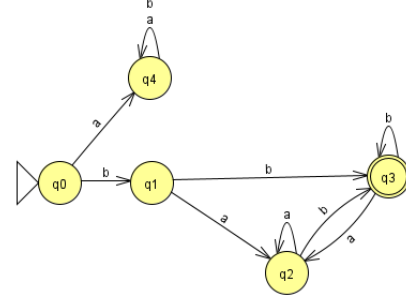


Figure 6: Derivative-based DFA generated for $L = b(a | b)^*ab$.

G Appendix: Common Problems Identified from LLM-Generated DFAs

This appendix summarizes recurring error patterns observed in deterministic finite automaton (DFA) constructions produced by large language models (LLMs) across different prompting strategies. Each problem type is illustrated using a concrete example and a corresponding DFA diagram. The goal is diagnostic: to characterize systematic construction failures rather than isolated errors.

Summary of Identified Problem Types

Problem 1: Disregard for Kleene Star Semantics

Figure 5 shows a DFA generated via direct construction for the language $L = b(a | b)^*ab$. The model correctly identifies the terminal pattern ab but fails to account for the Kleene star $(a | b)^*$, resulting in a DFA that accepts strings ending in ba rather than enforcing the required suffix ab . This error reflects an incomplete semantic interpretation of unbounded repetition under the Kleene star.

Problem 2: Violation of Brzowski Derivative Semantics

Construction method: Derivative-based construction. Figure 6 shows a derivative-based DFA generated for the language $L = b(a | b)^*ab$. The construction fails to normalize semantically equivalent Brzowski derivatives and incorrectly reintroduces consumed prefixes. As a result, distinct residual languages are treated as separate states, violating the formal semantics of regular expression derivatives and yielding an incorrect DFA.

Problem 3: Errors in Initial DFA Construction

Construction method: Minimization-based construction. Figure 7 shows a DFA produced prior to minimization for the language $L = b(a | b)^*ab$. The initial DFA is incorrectly constructed due to

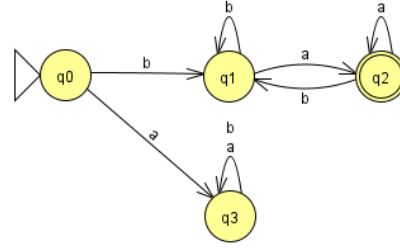


Figure 7: DFA produced prior to minimization for $L = b(a | b)^*ab$.

faulty state semantics and transition assignments. Subsequent minimization amplifies these errors by merging states that should remain distinct, resulting in an invalid minimized DFA.

Problem 4: Over-Acceptance of Strings

Construction method: Derivative-based construction. Figure 8 shows an over-accepting DFA generated for the language

$$L = ((a | b)^*(aa^*bb^* | b(b | ab)^*aaa^*bb^*) (aa^*bb^*)^* | (a | b)^*b(b | ab)^*a)aa^*.$$

The constructed DFA accepts strings that do not belong to the target language: formally, there exists a string $w \in \Sigma^*$ such that $w \in L(\text{DFA})$ while $w \notin L$. This over-acceptance arises when syntactically similar but semantically distinct derivatives are merged without establishing true language equivalence. In this example, strings of the form bb^*a are incorrectly accepted.

Problem 5: Failure to Preserve Constraints under Concatenation

Construction methods: Direct and Thompson constructions. Figure 9 shows a DFA that fails to preserve constraints under concatenation for the language

$$L = ((a^*ba^*b | (a^*ba^*b | a^*a)(b | a)^*a)a^*.$$

Table 6: Summary of Common DFA Construction Problems

Problem	Identified Issue
Problem 1	Disregard for Kleene star semantics
Problem 2	Violation of Brzowski derivative semantics
Problem 3	Errors in initial DFA construction before minimization
Problem 4	Over-acceptance of strings outside the target language
Problem 5	Failure to preserve constraints under concatenation
Problem 6	Introduction of redundant or unreachable states

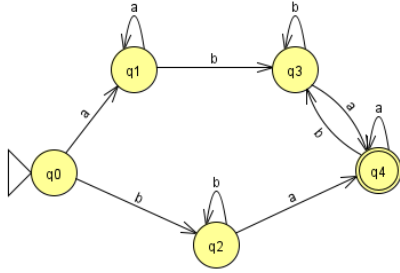


Figure 8: Over-accepting DFA generated for the target language.

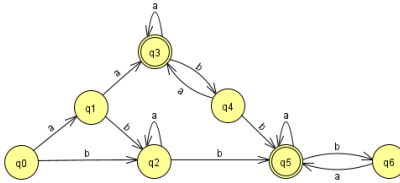


Figure 9: DFA failing to preserve constraints under concatenation.

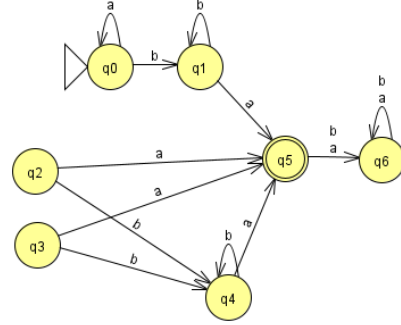


Figure 10: DFA containing redundant and unreachable states.

icates that the model introduces states corresponding to syntactic components of the regular expression without performing reachability or usefulness analysis, resulting in redundant and semantically irrelevant states.

For an overall view, Figure 11 shows the common mistakes for the Language L_9 under the ToT methods.

H Appendix : Derivative-Based DFA Output for L_9

This appendix reports the exact output produced by the derivative-based DFA construction for language L_9 . The content is shown verbatim as returned by the model, without normalization, correction, or post-processing.

H.0.1 Model Output (Derivative Construction for L_9)

```
{
  "method": "derivative",
  "derivatives": {
    "q0": {
      "regex": "b(a|b)*ab",
      "nullable": false
    },
    "q1": {
      "regex": "(a|b)*ab",
      "nullable": false
    },
    "q2": {
      "regex": "b(a|b)*ab+(a|b)*b",
      "nullable": false
    },
    "q3": {
      "regex": "(a|b)*ab+(a|b)*",

```

The constructed DFA correctly enforces individual sub-constraints such as a^*ba^*b and a^*a in isolation. However, when these constraints are composed via concatenation, global boundary conditions are not preserved.

Problem 6: Introduction of Redundant or Unreachable States

Construction methods: Direct, Thompson, derivative-based, and minimization-based constructions.

Figure 10 shows a DFA containing redundant and unreachable states for the language

$$L = ((a^*b | (a^*a)a^*b | (a^*b)b^*b | ((a)^*ba | (a^*a)a^*ba | (a^*b)b^*ba) | (aa | b)^*a)(a | b) | (a^*b)b^*a.$$

The constructed DFA contains multiple states that are unreachable from the start state and therefore do not contribute to language recognition. In some examples, some of these states are accepting despite having no incoming paths. This indi-

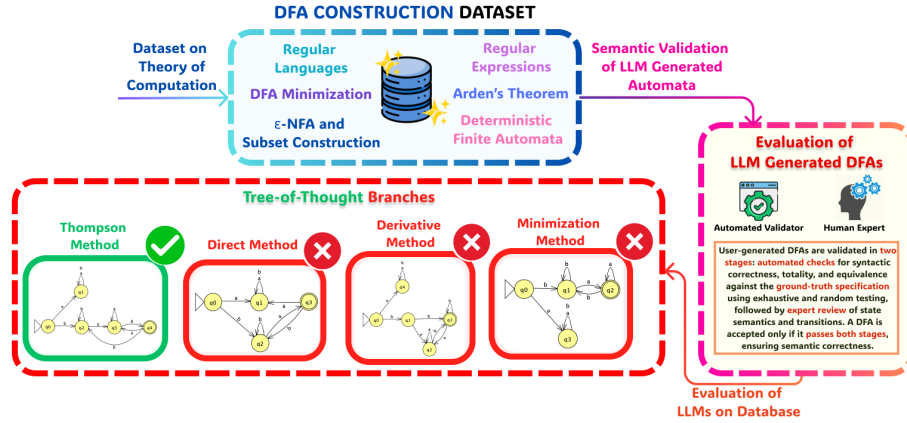


Figure 11: DFA construction dataset and evaluation pipeline, with representative DFA outputs generated by ToT prompting following construction strategies for the language $L_9 = b(a + b)^*ab$.

```

3288     "nullable": true
3289   },
3290   "q4": {
3291     "regex": "",
3292     "nullable": false
3293   }
3294 },
3295 "dfa": {
3296   "states": [
3297     "q0",
3298     "q1",
3299     "q2",
3300     "q3",
3301     "q4"
3302   ],
3303   "alphabet": [
3304     "a",
3305     "b"
3306   ],
3307   "start_state": "q0",
3308   "accept_states": [
3309     "q3"
3310   ],
3311   "transitions": {
3312     "q0": {
3313       "a": "q4",
3314       "b": "q1"
3315     },
3316     "q1": {
3317       "a": "q2",
3318       "b": "q3"
3319     },
3320     "q2": {
3321       "a": "q2",
3322       "b": "q3"
3323     },
3324     "q3": {
3325       "a": "q2",
3326       "b": "q3"
3327     },
3328     "q4": {
3329       "a": "q4",
3330       "b": "q4"
3331     }
3332   }
3333 }
3334 }

```

I Appendix: Determinism and Randomness Control

All experiments were designed to be as deterministic and reproducible as possible, subject to the constraints of black-box LLM APIs.

Decoding determinism. For all models and prompting strategies, decoding temperature was fixed to 0.0, and no stochastic sampling parameters were modified. This ensures that, for a given prompt and model endpoint, the generation process is deterministic to the extent supported by the provider.

Prompt determinism. Each prompt template is fixed and version-controlled in the repository. For Tree-of-Thought and Chain-of-Thought settings, a fixed, strategy-specific prompt template is reused across all runs. No adaptive prompt modification or dynamic hint generation is performed during a single model invocation.

Retry policy. To mitigate transient API failures and formatting errors, each query is retried up to a fixed number of attempts using an identical prompt. Retries are only triggered when the output is invalid (e.g., non-JSON or schema-violating) and do not introduce additional randomness or alternative prompts.

Validation randomness control. Behavioral equivalence testing uses a hybrid validation strategy consisting of: (i) exhaustive enumeration of all strings up to a fixed maximum length, and (ii) randomized testing over longer strings. All randomized testing is performed with a fixed random seed, ensuring that validation results are fully reproducible.

Non-determinism across API calls. Despite the above controls, repeated API calls to the same model may still produce different outputs

3373 due to undocumented provider-side nondetermin-
3374 ism (e.g., model updates, inference-time optimiza-
3375 tions, or distributed serving). Such variation is
3376 treated as an inherent property of deployed LLM
3377 systems rather than experimental noise.

3378 **Reporting policy.** All reported results corre-
3379 spond to the actual outputs returned by the API
3380 during evaluation and are not post-selected. Time-
3381 outs, formatting failures, and invalid outputs are
3382 logged explicitly and counted according to the
3383 evaluation protocol described in the main paper.

3384 **J Appendix: Runtime Environment**

3385 All experiments were executed using custom
3386 Python scripts on a local workstation. The imple-
3387 mentation was written in Python (version 3.10+)
3388 and relied exclusively on standard libraries and
3389 official model APIs. No proprietary or undocu-
3390 mented tooling was used.

3391 All LLM interactions were performed via API
3392 calls using deterministic decoding settings (tem-
3393 perature set to 0). Experiments were conducted
3394 sequentially with explicit rate-limiting and retry
3395 logic to ensure stability and reproducibility. Each
3396 model invocation was stateless: no conversational
3397 context was carried across calls. This includes
3398 the hint-based framework, where the full problem
3399 specification was re-provided at every stage.

3400 Model outputs, metadata (including token usage
3401 and finish reasons), parsed DFA tables, and val-
3402 idation reports were logged to disk in structured
3403 JSON format. The runtime environment enforced
3404 strict output handling: results were recorded for
3405 every attempt, including malformed outputs, time-
3406 outs, or generation failures.

3407 Validation and evaluation were performed of-
3408 fline using a separate automated pipeline that
3409 exhaustively and randomly tested DFA behavior
3410 against the ground-truth regular expressions. All
3411 randomness used in validation (e.g., random string
3412 generation) was controlled via fixed random seeds.

3413 No fine-tuning, model-side configuration
3414 changes, or system-level optimizations were
3415 applied. Differences in performance therefore
3416 reflect intrinsic model behavior under identical
3417 runtime conditions rather than environmental
3418 variability.