
Reuse Bench

Anonymous Author(s)
Affiliation
Address
email

Abstract

1 Ricci is awesome.

2 **1 Introduction**

3 **2 Background**

4 **3 Datasets**

5 **4 Analysis**

6 **5 Conclusion**

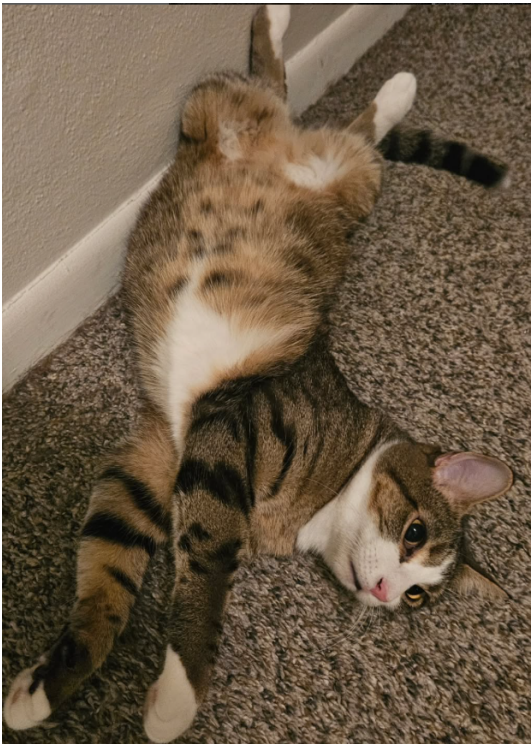


Figure 1: Ricci

7 References

8 A Dataset Summary

9 A.1 Identity

10 Problem Description.

11 Problem Example.

12 Solution Example.

13 Reuse component.

14 Difficulty control.

15 Evaluation.

16 Dataset Info. Dataset size

17 B CNF Conversion

18 **Problem Description.** CNF (Conjunctive Normal Form) conversion tasks typically ask us to start
19 with a propositional logic formula (which may include connectives $\wedge, \vee, \rightarrow, \neg, \dots$) and system-
20 atically transform it into a *conjunction of disjunctions of literals*. The solver is required to *apply*
21 *specific rewriting rules* (e.g., removing implications, moving \neg inward, distributing \vee over \wedge) in
22 a **step-by-step** procedure, ensuring that each intermediate formula is logically equivalent to the
23 previous one.

24 **Problem Example.** A typical assignment might read:

“Convert $(p \rightarrow q) \wedge \neg(r \vee p)$ into CNF. Show each transformation step: remove \rightarrow , push \neg , distribute \vee , etc.”

25 So, the solver must:

- 26 • Rewrite $(p \rightarrow q)$ as $(\neg p \vee q)$,
- 27 • Apply De Morgan’s laws on $\neg(r \vee p)$ to get $(\neg r \wedge \neg p)$,
- 28 • Possibly distribute any \vee over \wedge to finalize CNF form,
- 29 • Verify each step is an *equivalence* transformation.

30 **Solution Example.** Following the example $(p \rightarrow q) \wedge \neg(r \vee p)$:

31 1. **Remove implication:** $p \rightarrow q \equiv \neg p \vee q$. So the formula becomes:

$$(\neg p \vee q) \wedge \neg(r \vee p).$$

32 2. **Push negation inwards:** $\neg(r \vee p)$ by De Morgan’s law is $(\neg r \wedge \neg p)$. So now:

$$(\neg p \vee q) \wedge (\neg r \wedge \neg p).$$

33 3. **Distribute** if needed:

$$(\neg p \vee q) \wedge \neg r \wedge \neg p$$

34 is already in CNF, since it is a conjunction of three clauses: $(\neg p \vee q)$, $(\neg r)$, and $(\neg p)$.

35 Hence the final CNF is:

$$(\neg p \vee q) \wedge (\neg r) \wedge (\neg p).$$

36 **Reuse component.** Once certain rewriting *identities* are established (like $\neg(p \vee q) \equiv \neg p \wedge \neg q$,
37 or $p \rightarrow q \equiv \neg p \vee q$), one can *cite* them without re-deriving at each step. Similarly, if an earlier
38 problem already removed implications from a sub-formula, that sub-formula can be *embedded* in a
39 new formula—saving repeated transformations. This is akin to citing *known lemmas* in mathematics.

40 **Difficulty control.** Ways to scale **difficulty**:

- 41 • **Formula complexity:** more nested connectives, multiple implications, extended negations.
- 42 • **Extra equivalences needed:** dealing with \leftrightarrow , exclusive or, or special connectives (like \oplus).
- 43 • **Size of final distribution:** distributing \vee over \wedge in large expressions can be many steps.
- 44 • **Strict normal forms:** Some tasks might require not just CNF, but also further factorization
- 45 (like a canonical ordering of literals).

46 **Evaluation.** To check correctness:

- 47 • **Equivalence check per step:** each rewrite must be a known logical equivalence (implication
- 48 removal, De Morgan's, double negation, distribution).
- 49 • **Final form** is indeed CNF: a conjunction of disjunctions of literals (no leftover operators in
- 50 the wrong place).
- 51 • **No skipped transformations:** if an entire distribution step is large, ensure the solver shows
- 52 partial expansions or at least references the standard distribution law carefully.

53 A correct derivation ensures the final formula has the same truth conditions as the original, now in a

54 strict CNF layout.

55 **Dataset Info.** Dataset size

56 C Context-free Grammar

57 **Problem Description.** Context-free grammar (CFG) checks often ask the solver to:

- 58 • Verify whether a given string can be derived by a certain CFG (i.e. membership checking).
- 59 • Prove or disprove that two grammars generate the same language.
- 60 • Convert a CFG to a certain normal form (e.g. Chomsky Normal Form), or show step-by-step
- 61 derivations.

62 Such tasks require explicit **step-by-step** reasoning, for example enumerating a parse tree or using a

63 parsing algorithm (like CYK, Earley, LL(1), or LR parsing) in a canonical order.

64 **Problem Example.** A typical instance might be:

“Given the grammar G with start symbol S and productions:

65

$$\begin{array}{lcl} S \rightarrow AB & | & b \\ A \rightarrow aA & | & \epsilon \\ B \rightarrow b \end{array}$$

66

check whether the string aab is in $L(G)$. Show the stepwise derivation or parse steps.”

67 Alternatively, a question could ask to transform G into Chomsky Normal Form, detailing every

68 elimination of ϵ -rules, unit-rules, etc.

69 **Solution Example.** Below is a short membership check using a leftmost derivation for aab :

70 1. Start with S .

71 2. Apply $S \rightarrow AB$:

$$S \Rightarrow AB.$$

72 3. For A , pick $A \rightarrow aA$:

$$AB \Rightarrow aAB.$$

73 4. Again for A , pick $A \rightarrow aA$:

$$aAB \Rightarrow aaAB.$$

74 5. Now $A \rightarrow \epsilon$:

$$aaAB \Rightarrow aaB.$$

75 6. For B , pick $B \rightarrow b$:

$$aaB \Rightarrow aab.$$

76 7. Hence $S \Rightarrow aab$. Thus $aab \in L(G)$.

77 **Reuse component.** Many CFG transformations or parsing steps are reusable:

- 78 • Once we have removed ϵ -productions or converted a subgrammar to CNF, we can *cite* that
79 partial result in a larger grammar.
- 80 • A standard parsing table for a known sub-grammar can be reused in multiple parse checks.
- 81 • Known normal forms or well-known derivation patterns (like $S \rightarrow aSb \mid \epsilon$ for balanced
82 parentheses) can be invoked repeatedly without re-derivation.

83 **Difficulty control.** Ways to scale **difficulty** in CFG tasks:

- 84 • **Grammar complexity:** more nonterminals, more productions, nested recursion.
- 85 • **String length:** from short strings (quick derivations) to longer ones requiring multiple
86 expansions.
- 87 • **Normalization or transformations:** Chomsky Normal Form, Greibach Normal Form, or
88 removing left recursion can be more involved for bigger grammars.
- 89 • **Parsing methods:** e.g. CYK vs. LL(1), each step must follow canonical rules for a unique
90 path.

91 **Evaluation.** To judge correctness:

- 92 • **Verify each production use** in the derivation: every expansion must match a grammar rule
93 exactly.
- 94 • **Check final derived string:** it must equal the target string with no extraneous symbols left.
- 95 • **Parsing correctness:** if using a parser-based approach (LL(1), LR(1)), confirm each step in
96 the parsing table is consistent with the grammar and the next input symbol.
- 97 • **Grammar transformations:** each rewrite (eliminating ϵ , factoring out left recursion) must
98 be verifiable by standard transformations that preserve language equivalence.

99 Any mismatch in expansions or incorrect unit-rule elimination reveals a flawed step. A correct chain
100 of derivations or transformations ensures soundness of the final conclusion.

101 **Dataset Info.** Dataset size

102 D Lambda Calculus

103 **Problem Description.** Lambda calculus reduction tasks typically involve starting with a λ -
104 expression (possibly containing multiple nested abstractions and applications) and *reducing* it via
105 β -reduction rules (and sometimes α -conversion to avoid variable capture). The solver is asked to
106 derive a normal form, or to perform a specific reduction strategy (e.g., *call-by-value*, *call-by-name*, or
107 *leftmost-outermost*). Each intermediate step must be **verifiable** by the standard λ -calculus rules, and
108 a *unique path* is enforced by fixing a particular strategy.

109 **Problem Example.** A typical problem might read:

“Reduce the expression $(\lambda x. x+1)((\lambda y. y+y) 3)$ using call-by-value. Show each step of the reduction.”

110 Here, the solver must:

- 111 • Evaluate the *argument* $(\lambda y. y + y) 3$ first (call-by-value),
- 112 • Perform the actual β -reduction once that argument is evaluated,
- 113 • Arrive at a final numeric form (or a Church numeral, depending on how arithmetic is
114 encoded).

115 **Solution Example.** Below is a short step-by-step reduction (in call-by-value) for

$$(\lambda x. x + 1)((\lambda y. y + y) 3).$$

116 1. Evaluate the argument $(\lambda y. y + y) 3$:

$$(\lambda y. y + y) 3 \longrightarrow 3 + 3 = 6$$

117 (assuming an arithmetic interpretation or a further β -reduction if Church numerals are in
118 play).

119 2. Now substitute the result into $\lambda x. x + 1$:

$$(\lambda x. x + 1) 6 \longrightarrow 6 + 1 = 7.$$

120 3. Final result: 7.

121 Each \rightarrow step references a valid β -reduction (or arithmetic evaluation if representing numbers as
122 Church numerals), consistent with call-by-value ordering.

123 **Reuse component.** Common *mini-reductions* or known sub-expressions are often reused:

- 124 • Once you have shown $\lambda y. (y + y)$ acting on a Church numeral n reduces to $2n$, you can *cite*
125 that directly in larger expressions.
- 126 • Standard combinators (like $\mathbf{I} = \lambda x. x$ or $\mathbf{K} = \lambda x. \lambda y. x$) can be used and referenced as
127 lemmas.

128 Hence, partial λ -calculus derivations can be integrated into bigger tasks much like citing theorems or
129 subroutines.

130 **Difficulty control.** Ways to scale **difficulty**:

- 131 • **Expression size:** More nested applications and abstractions, e.g. multiple λ layers.
- 132 • **Reduction strategy:** (Call-by-name vs. call-by-value vs. normal-order) influences the path
133 of reduction.
- 134 • **Explicit α -conversion:** If variable capture arises, ensuring correct renaming steps can add
135 complexity.
- 136 • **Encoding of data structures:** If the expression uses Church numerals, booleans, or pairs,
137 the solver must recall these encodings and systematically reduce them.

138 **Evaluation.** To verify a solution:

- 139 • **Check each β -reduction:** Ensure correct substitution (no variable capture, no missed
140 occurrences of the replaced variable).
- 141 • **Check the sequence of redexes** matches the required strategy (e.g. leftmost-outermost).
- 142 • **Confirm final normal form:** No further β -reductions are possible (if a normal form is
143 required).
- 144 • If the solution is partial, check that all intermediate steps align with the official reference
145 derivation or a standard λ -calculus engine for correctness.

146 **Dataset Info.** Dataset size

147 E Merge Sort

148 **Problem Description.** Merge Sort is a *divide-and-conquer* sorting algorithm that splits an array
149 (or list) into two halves, recursively sorts each half, then *merges* the two sorted halves into one fully
150 sorted array. A typical Merge Sort reasoning task requires demonstrating each recursive call and
151 intermediate merge step in **step-by-step** detail. This ensures the solver clearly shows how subarrays
152 get sorted and combined.

153 **Problem Example.** A standard prompt might say:

“Apply Merge Sort to the list [4, 1, 3, 9, 7]. Show each major recursive call and how you merge sorted sublists.”

154 The solver is expected to partition the array repeatedly (splitting in half each time), then do step-by-
155 step merges of smaller sorted subarrays until the entire list is sorted.

156 **Solution Example.** A brief illustration using the list $A = [4, 1, 3, 9, 7]$:

157 1. **Split** A into two parts:

$$A_{\text{left}} = [4, 1, 3], \quad A_{\text{right}} = [9, 7].$$

158 2. **Recursively sort** A_{left} :

- 159 • Split $[4, 1, 3]$ into $[4, 1]$ and $[3]$.
- 160 • Sort $[4, 1]$ by splitting into $[4]$ and $[1]$, then merging to get $[1, 4]$.
- 161 • Merge $[1, 4]$ with $[3]$ to get $[1, 3, 4]$.

162 3. **Recursively sort** $A_{\text{right}} = [9, 7]$:

- 163 • Split into $[9]$ and $[7]$, then merge to get $[7, 9]$.

164 4. **Final merge:** merge $[1, 3, 4]$ with $[7, 9]$ step by step:

- 165 • Compare 1 and 7: pick 1 first, etc.
- 166 • Eventually obtain the sorted array $[1, 3, 4, 7, 9]$.

167 **Reuse component.** Once you have demonstrated the subroutine for *merging two sorted lists* in one
168 problem, you can *cite* that process in subsequent tasks. For example, if we already established how to
169 merge $[1, 4]$ and $[3]$, we do not need to re-derive every comparison in a new question about merging
170 $[2, 5]$ and $[1, 9]$. Thus, the merging algorithm acts like a reusable “lemma” or subroutine in future
171 Merge Sort tasks.

172 **Difficulty control.** To tune the **difficulty**:

- 173 • **List size:** A short list of 5 elements vs. a large list of 20+ elements (which yields more
174 recursive calls).
- 175 • **Varied structure:** Repeated elements, negative numbers, or special patterns (e.g., partially
176 sorted lists).
- 177 • **Intermediate details required:** Some tasks only require the final sorted array, others
178 demand each *partial* merge step in detail.
- 179 • **Tie-breaking strategy:** If you want a unique step-by-step path, specify how to handle ties
180 or even how to split uneven arrays (e.g., if 7 elements, which half gets 3 vs. 4).

181 **Evaluation.** To verify correctness:

- 182 • Check each *recursive partition* is done as specified (splitting the list correctly).
- 183 • Confirm each *merge step* is done in ascending order, comparing the smallest elements first.
184 One can re-check the subarrays are indeed sorted before merging.
- 185 • The final output should be the list in sorted order.
- 186 • If the task demands full tracing, ensure each sub-step matches a known reference solution
187 (every partial merge or recursion is in the correct sequence).

188 **Dataset Info.** Dataset size

189 F Quantified Boolean Formulas (QBF)

190 **Problem Description.** A *Quantified Boolean Formula* (QBF) is an extension of SAT where variables
191 are universally (\forall) or existentially (\exists) quantified. A typical QBF has the form

$$Q_1 x_1 \ Q_2 x_2 \ \dots \ Q_n x_n \ \Phi(x_1, x_2, \dots, x_n),$$

192 where each $Q_i \in \{\forall, \exists\}$ and Φ is a propositional formula (often in CNF). The question is whether the
193 quantified formula is *true* or *false*, given the sequence of quantifiers. Solving QBF requires reasoned
194 assignment or proof for each variable in turn, respecting \exists (can choose a witness) and \forall (must hold
195 for all choices).

196 **Problem Example.** A small QBF instance could be:

$$\forall x_1 \exists x_2 \forall x_3 ((x_1 \vee \neg x_2 \vee x_3) \wedge (\neg x_1 \vee x_2)).$$

197 Here, we must check whether for *all* values of x_1 , there is a choice for x_2 such that for *every* choice
198 of x_3 , the formula is satisfied. The solver must systematically analyze each quantifier layer.

199 **Solution Example.** A step-by-step approach might follow:

- 200 1. **Outer quantifier:** $\forall x_1$. Consider $x_1 = \text{true}$ and $x_1 = \text{false}$ separately.
- 201 2. For each x_1 value, check **next quantifier:** $\exists x_2$. Possibly guess or systematically pick an x_2
202 that satisfies the sub-formula for both possibilities of x_3 .
- 203 3. Then **last quantifier:** $\forall x_3$. Evaluate the formula for $x_3 = \text{true}$ and $x_3 = \text{false}$.
- 204 4. If for *every* x_1 one can pick an x_2 that works for *all* x_3 , the formula is *true*. Otherwise, it is
205 *false*.

206 Each sub-step can be turned into a short propositional check (like small SAT tasks), verifying the
207 partial assignments plus universal checks.

208 **Reuse component.** Sub-parts of a QBF often reduce to simpler SAT or lower-level QBF sub-tasks.
209 Once a smaller formula or quantifier structure is resolved in one step, it can be *cited* in later steps.
210 For example, if $\exists x_2$ satisfies a sub-formula ψ regardless of x_3 , that partial assignment is reusable.
211 Also, known lemmas about certain standard QBF forms (like “ $\forall x \exists y : x \neq y$ is always true if the
212 domain is more than 1 bit”) can be cited repeatedly.

213 **Difficulty control.** To scale **difficulty** in QBF:

- 214 • **Number of quantifiers:** short chain (2–3) vs. deeper nesting (5+).
- 215 • **Placement of universal vs. existential variables:** more frequent alternation $\forall \exists \forall \exists$ is
216 typically more complex.
- 217 • **Size of underlying propositional formula:** simpler 2-3 clauses vs. large CNF.
- 218 • **Special structure:** e.g., formula structured in *prenex normal form* with a particular pattern
219 can be easier or harder to solve.

220 **Evaluation.** To judge correctness of a QBF solution:

- 221 • Verify each quantifier step is checked properly: for $\forall x_i$, the solution must hold for both
222 $x_i = \text{true}$ and false (if x_i is Boolean). For $\exists x_i$, there must be *some* setting of x_i making the
223 sub-formula hold for all subsequent universal variables.
- 224 • Validate any sub-formula solves: e.g., each partial assignment is consistent with the original
225 clauses or sub-formulas.
- 226 • If the solver claims *true*, confirm each universal branch truly has an existential witness. If
227 *false*, confirm no existential assignment can satisfy all universal choices.

228 Hence, a thorough *step-by-step* argument or a known QBF solver trace can be used to confirm
229 correctness.

230 **Dataset Info.** Dataset size

231 G Regular Expression

232 **Problem Description.** Regular expression tasks typically involve:

- 233 • *Constructing* a regular expression R that describes a certain language L (e.g., “all strings
234 over $\{a, b\}$ of even length”).
- 235 • *Transforming* an existing expression into a simplified or normal form (e.g., converting $(ab)^*$
236 to a certain canonical representation).
- 237 • *Proving equivalence* of two given regular expressions (or an expression and a DFA), showing
238 that they define the same language.

Such tasks generally require a **step-by-step** application of known *regex* identities or constructive proofs (union, concatenation, Kleene star, etc.), with each step *verifiable* by recognized rewriting or equivalence checks.

Problem Example. A typical problem might say:

“Show that the language of the regular expression $(ab)^*$ is accepted by the following NFA. Then prove their equivalence or con

Or

“Simplify the regular expression $(a \cup \epsilon)^*(a \cup b)$ to a canonical form, using standard regex identities.”

Solution Example. A step-by-step solution might illustrate how to rewrite:

$$(a \cup \epsilon)^*(a \cup b)$$

into an equivalent simpler expression. For instance:

1. Observe $(a \cup \epsilon)^* = (a \cup \epsilon)^*(\epsilon \cup a)$ by commutativity of union. Potentially factor or rewrite in a normal form.

2. Use the identity $R \cup \epsilon \equiv \epsilon \cup R$ which sometimes can be converted to optional forms, or repeated expansions.

3. Carefully apply distribution laws if needed:

$$(XY \cup XZ) = X(Y \cup Z) \quad (\text{for expressions } X, Y, Z).$$

4. Eventually arrive at a more canonical or simpler expression such as $(a^*)(a \cup b) = a^*(a \cup b)$, etc.

Each step references a known equivalence rule (distributive law, absorbing law, or expansion of Kleene star).

Reuse component. Sub-results or partial transformations can be reused:

- If you have already shown $(ab)^* \cup (ab)^*a = (ab)^*$ in some prior step, you can *cite* that result in new proofs.

- A known minimal expression for a sub-language can be embedded within a more complicated expression.

Much like “citing a lemma” in mathematics, once a certain regex identity is established, it can be referenced repeatedly.

Difficulty control. Ways to scale **difficulty** in regular expression tasks:

- **Expression size:** More operators, repeated nesting of Kleene stars, or multiple unions.
- **Target normal forms:** Requiring rewriting to a standard normal form (like GNFA-based transformations) can introduce many steps.
- **Equivalence proofs:** Checking that two large expressions define the same language can be more difficult than simpler rewriting tasks.
- **Combined tasks:** Converting a DFA to a regex or vice versa for big automata is more involved.

Evaluation. To evaluate correctness:

- **Step-by-step identity checks:** Verify each rewrite rule is correctly applied (e.g., distribution, union/concatenation expansions).
- **Equivalence checks:** For final expressions, confirm they represent the same language as the initial expression (possibly by known theorems or by constructing NFAs).
- **No skipped transformations:** Each intermediate form must be justified by a standard regex equivalence law.

Because *regex rewriting* is typically *verifiable* by formal transformations or by mechanical checks (like testing string membership in both expressions), any incorrect step stands out upon closer inspection.

280 **Dataset Info.** Dataset size

281 **H SAT (Boolean Satisfiability Problem)**

282 **Problem Description.** SAT (Boolean satisfiability) problems present a Boolean formula in CNF
283 (Conjunctive Normal Form), or some other equivalent form, and ask whether there exists a truth
284 assignment to the variables that makes the entire formula *true*. For example, the formula might be a
285 collection of clauses (disjunctions of literals) that must all be satisfied simultaneously. A step-by-
286 step derivation typically involves applying a systematic decision procedure or resolution method to
287 demonstrate either a **SAT** assignment or a proof of **UNSAT**.

288 **Problem Example.** A typical SAT instance might look like:

$$\Phi = (x_1 \vee \neg x_2 \vee x_3) \wedge (\neg x_1 \vee x_2) \wedge (x_2 \vee x_3) \wedge \dots$$

289 The question is: *Does a truth assignment to $\{x_1, x_2, x_3, \dots\}$ exist that satisfies all clauses? If so,*
290 *provide one such assignment. Otherwise, show it is unsatisfiable.*

291 **Solution Example.** A valid step-by-step solution might adopt a **DPLL** (Davis-Putnam-Logemann-
292 Loveland) approach or a **resolution** derivation:

- 293 1. **Pick a variable** x_1 . Assume $x_1 = \text{true}$.
- 294 2. Simplify the formula by removing any clauses containing x_1 (since they are satisfied) and
295 removing $\neg x_1$ from any clauses where it appears.
- 296 3. If we reach a conflict (an empty clause), backtrack and assume $x_1 = \text{false}$ instead.
- 297 4. Continue branching on x_2, x_3, \dots in a similar fashion.
- 298 5. If a complete assignment is found that satisfies all clauses, we are done (SAT). If all branches
299 fail, the formula is UNSAT.

300 Each branching step and clause simplification can be checked mechanically.

301 **Reuse component.** When analyzing multiple SAT instances, we can reuse certain sub-formulas or
302 *lemmas* about them. For instance:

- 303 • A known tautology or small sub-formula that is always satisfiable (like $(p \vee \neg p)$).
- 304 • A repeated *core* formula whose satisfiability is already proven in a separate question.

305 By referencing these known components, one can shortcut repeated expansions or branching steps.

306 **Difficulty control.** To scale the **difficulty** of SAT problems:

- 307 • **Number of variables:** Small (3–5) vs. large (20+).
- 308 • **Clause density or structure:** Some formulas have many short clauses, others have fewer
309 but longer clauses.
- 310 • **Special forms:** 2-SAT (polynomial-time solvable) vs. general 3-SAT.
- 311 • **Branching complexity:** Adding symmetrical or *trick* clauses to force more branching in a
312 DPLL approach.

313 **Evaluation.** To evaluate correctness:

- 314 • **Check each branching step:** ensuring any partial assignment is consistent, and each
315 simplified formula is correct.
- 316 • **Verify the final assignment** (for SAT): plugging it into the original formula to confirm
317 every clause is satisfied.
- 318 • **Resolution-based approach:** verify that each derived clause follows from previous ones,
319 culminating in an empty clause (if unsatisfiable).

320 If any step misapplies the branching, resolution, or simplification rules, the solver's derivation is
321 invalid. A valid chain, in contrast, demonstrates soundness and completeness for that specific instance.

322 **Dataset Info.** Dataset size

I Symbolic Differentiation

Problem Description. Symbolic differentiation tasks provide a symbolic function $f(x)$ (possibly involving polynomials, exponentials, products, quotients, or trigonometric functions) and ask the solver to *differentiate* it step by step. The goal is to produce a fully expanded or simplified derivative, following a specific rule sequence (e.g., product rule, chain rule) at each stage. This enforces a **multi-step** solution, and each intermediate step is **verifiable** by checking standard differentiation rules.

Problem Example. A typical question might read:

“Differentiate $f(x) = x^2 \cdot e^x$ with respect to x . Show each rule used (product rule, chain rule) at each step.”

In this scenario, we identify the relevant differentiation rule (the product rule), apply it carefully, and optionally simplify the resulting expression.

Solution Example. A valid step-by-step solution for $f(x) = x^2 e^x$:

1. Let $u = x^2$ and $v = e^x$.

2. By the product rule: $\frac{d}{dx}[u \cdot v] = u' \cdot v + u \cdot v'$.

3. Compute $u' = \frac{d}{dx}(x^2) = 2x$.

4. Compute $v' = \frac{d}{dx}(e^x) = e^x$.

5. Therefore,

$$f'(x) = (2x)e^x + x^2(e^x) = e^x(2x + x^2).$$

One can verify each sub-step by checking the correctness of product rule application and basic derivative formulas.

Reuse component. Often, the same *mini-derivatives* appear across multiple questions. For instance, once we have established that $\frac{d}{dx}(x^n) = nx^{n-1}$ or that $\frac{d}{dx}(e^x) = e^x$, we can cite those directly in subsequent problems. Similarly, if we already derived that $\frac{d}{dx}(\sin(ax)) = a \cos(ax)$ for a constant a , we can reuse it for many tasks. This is analogous to *citing known sub-results* or lemmas, eliminating repeated derivations.

Difficulty control. To scale the **difficulty** of symbolic differentiation problems:

- **Number of factors:** Larger products (e.g., three- or four-term products) require multiple product rule expansions.
- **Composition depth:** Nested chain rules (e.g., $\sin(e^{x^2})$).
- **Mixed functions:** Polynomials, exponentials, trigonometric functions, logarithms, or piecewise definitions.
- **Simplification complexity:** Requiring the final answer in a specific form (fully expanded vs. factored) can add steps.

Evaluation. Each differentiation step is checked against a *known derivative rule* (e.g., product rule, chain rule, quotient rule, basic function derivatives). A solver’s correctness can be assessed by:

- Verifying each intermediate expression matches the declared rule application (e.g., does $u' \cdot v + u \cdot v'$ truly expand to the solver’s next line?).
- Optionally, performing a final check: re-differentiate the solver’s result or substitute numeric values to see if it matches the true derivative at sample points.

Because differentiation is well-defined in terms of standard rules, any incorrect algebra or omitted factor stands out quickly during verification.

Dataset Info. Dataset size

J Trigonometric Identity Transform

Problem Description. Trigonometric transformation problems present an expression involving trigonometric functions (e.g., \sin , \cos , \tan , \cot , etc.), possibly combined with algebraic factors or exponentials, and require rewriting it into a specified *target form* (such as a single trig function, or only \sin and \cos , or a sum-product form). The solver must apply **canonical trigonometric identities** (angle-sum formulas, double-angle formulas, product-to-sum identities, etc.) in a specific stepwise manner, ensuring each transformation is *verifiable*.

Problem Example. A typical question might say:

“Rewrite $\tan(x) \cdot \cos(2x)$ as an expression purely in terms of $\sin(x)$ and $\cos(x)$. Show each trig identity used.”

For instance, one might expand $\cos(2x)$ with $\cos^2(x) - \sin^2(x)$, then rewrite $\tan(x)$ as $\sin(x)/\cos(x)$, and so on, step by step.

Solution Example. Here is a brief sketch of the solution steps for $f(x) = \tan(x) \cdot \cos(2x)$:

1. Rewrite $\tan(x)$:

$$\tan(x) = \frac{\sin(x)}{\cos(x)}.$$

2. Use a double-angle identity on $\cos(2x)$:

$$\cos(2x) = \cos^2(x) - \sin^2(x).$$

3. Combine them:

$$f(x) = \frac{\sin(x)}{\cos(x)} \cdot [\cos^2(x) - \sin^2(x)].$$

4. Distribute if desired:

$$f(x) = \sin(x) \cos(x) - \frac{\sin^3(x)}{\cos(x)}.$$

5. Further simplification may depend on the instructions (e.g., factor out $\sin(x)$, or keep it in that form). Each step references a known identity (double-angle, tan definition).

Reuse component. Once certain identities are proven or used (e.g., $\sin^2(x) + \cos^2(x) = 1$, $\sin(2x) = 2 \sin(x) \cos(x)$, or product-to-sum formulas), the solver can *cite* these sub-results repeatedly. Similarly, any partial expression that has been successfully transformed in a prior question (e.g., rewriting $\tan(x) \cos(2x)$) can be plugged in as a “lemma” in a larger or more complicated expression.

Difficulty control. Ways to scale **difficulty** in trigonometric transformations:

- **More layers of angles:** $\sin(2x)$ vs. $\sin(3x + \pi/4)$, multiple nested sums or differences.
- **Complex compositions:** e.g., rewriting $\sin(\cos(x))$ is typically out of scope, but adding exponentials, $\tan^2(x)$, or $\sec(x)$ can make transformations more involved.
- **Number of steps:** A short expression with one or two identities vs. a longer expression requiring four or five.
- **Specific final form:** Requiring a sum/product form, or purely sin-cos form, or factoring out common terms can add complexity.

Evaluation. Each step is evaluated by checking:

- **Identity correctness:** Did the solver apply the correct double-angle formula, sum-to-product identity, etc.?
- **Algebraic manipulation:** Are multiplications, divisions, and factorizations accurate?
- **Final form:** Does the final expression match the requested format (e.g., no leftover \tan , no powers of \sin , etc., if that was specified)?

Because trigonometric identities are standard, a mismatch at any step indicates an incorrect application or algebra slip. A correct chain of transformations is easily *verifiable* by substituting numeric values to confirm equivalence or by standard identity checks.

K Turing Machine Simulation

Problem Description. Turing Machine (TM) simulation problems present a deterministic TM description (states, tape alphabet, transition function, etc.) along with an initial tape configuration. The solver must step through each configuration of the TM until it halts, then report whether it *accepts* or *rejects* (or possibly provide the final tape contents). Such tasks require **step-by-step** verification of transitions and enforce a **unique reasoning path** because the machine is deterministic.

Problem Example. Consider a TM M with:

- States: $\{q_0, q_1, q_{\text{accept}}\}$, where q_0 is the start state and q_{accept} is the accepting state.
- Tape alphabet: $\{\sqcup, 0, 1\}$, where \sqcup represents a blank symbol.
- Input alphabet: $\{0, 1\}$.
- Transition function (partial example):

$$\delta(q_0, 0) = (q_1, 1, R), \quad \delta(q_0, 1) = (q_{\text{accept}}, 1, R), \quad \dots$$

If the initial tape is $\sqcup \sqcup 0 1 0 \sqcup \dots$ with the head on the first 0, the question is: *Does the machine accept or reject this input?*

To solve, one must enumerate each configuration: (current state, tape contents, head position) and apply δ until reaching either q_{accept} (accept) or a non-existent transition (reject).

Solution Example. A valid step-by-step solution might look like this (abbreviated):

1. **Configuration 1:** $(q_0, \underline{\sqcup 1 0} \dots)$.
 $\delta(q_0, 0) = (q_1, 1, R)$.
 \Rightarrow **Configuration 2:** $(q_1, \underline{\sqcup 1 1 0} \dots)$.
2. **Configuration 2:** $(q_1, \underline{\sqcup 1 1 0} \dots)$.
Suppose $\delta(q_1, 1) = (q_{\text{accept}}, 1, R)$.
 \Rightarrow **Configuration 3:** $(q_{\text{accept}}, \underline{\sqcup 1 1 0} \dots)$.
3. Since the machine is now in q_{accept} , it halts and **accepts** the input.

Each transition is mechanically checkable by looking up the relevant $\delta(q, \text{symbol})$.

Reuse component. A sub-routine or partial sequence of moves can be reused across problems. For instance, if a portion of the TM “rewinds the tape to the left,” and we have already traced that behavior in one question, we can cite that known sub-computation in another task where the same sequence of transitions occurs. This is analogous to “citing theorems” in proofs. Once established, the user can simply reference the known step-by-step transitions (e.g., “Subroutine \mathcal{R} that scans right until a blank symbol is found”) without re-deriving them.

Difficulty control. To scale **difficulty** in Turing Machine simulation tasks:

- **Tape length.** A short tape with 2–3 symbols vs. a longer tape with 10+ symbols.
- **Number/complexity of states.** A TM with 2–3 states vs. 5–10 states (or more intricate transitions).
- **Special sub-procedures.** Inclusion of loops or sub-routines that require multiple passes over the tape.
- **Time complexity.** TMs that run for many steps vs. those that halt quickly.

Larger tapes, more states, or complicated transition functions yield longer simulation traces and more verification steps.

Evaluation. To **evaluate** a solver’s correctness:

- Check whether each transition from configuration i to $i + 1$ matches the TM’s δ .
- Verify the final state (accept/reject) matches the official outcome.
- Optionally, compare the entire configuration trace to a “ground truth” reference to ensure no step is skipped or incorrectly applied.

447 Because the TM is deterministic, if the solver's trace diverges from the official reference at any point,
448 the solver's answer is incorrect.

449 **Dataset Info.**