

Conditional Deontics over Terminals: A Mildly Context-Sensitive Formal Grammar for Constrained Decoding

Anonymous ACL submission

Abstract

With the broad adoption of large language models (LLMs), techniques to guide their generation have become increasingly important, with particularly sophisticated patterns found in constrained decoding frameworks expressed via formal grammars. Recent frameworks such as IterGen expand beyond context-free grammars to achieve some context sensitivity but with drawbacks: expensive backtracking or speculative lookahead, permissive semantics that risk spurious constraint activation, and difficult integration with runtime optimizations like XGrammar. To address these concerns, we propose Conditional Deontics over Terminals (CDoT), a mildly context-sensitive grammar paired with a GPU-accelerated $O(n^3)$ parser that can be employed for constrained decoding with LLMs. We show that CDoT can express useful constraints that improve LLM performance in a traditional text-to-SQL task and a complex logic puzzle involving unit propagation.

1 Introduction

As large language models (LLMs) have become more popular, the importance of techniques to guide their generation has also grown. Techniques for steering LLMs target various model states, including activations, attention patterns, MoE router logits, and output logits. These interventions can be applied uniformly, or conditioned on the generation process itself. For example, CAST (Lee et al., 2025) conditions on the cosine similarity between activations and behavior vectors, FASB (Cheng et al., 2025) applies a neural classifier to activations to make steering decisions, and PASTA (Zhang et al., 2024) conditions on user-selected input spans. In such simple conditioning on the generation process, these methods depend on only a fixed window of past steps or inputs.

Constrained decoding is unique among steering methods in that it explicitly models long-range and

structurally complex patterns – such as programming language specifications, molecular structures, and privacy requirements (Ugare et al., 2025; Loula et al., 2025; Scholak et al., 2021; Poesia et al., 2022) – that cannot be captured by finite-window or input-fixed conditions. Formal grammars are central to this power.

Our contributions to this research direction are:

- We increase the expressive power of constrained decoders by designing a new formal grammar, Conditional Deontics over Terminals (CDoT), that is strictly more powerful than context-free grammars.
- We demonstrate that constrained decoding with LLMs under this new grammar formalism yields improvements on both a traditional text-to-SQL dataset and on a complex logic puzzle involving unit propagation.

2 Design Statement

2.1 Formal Grammars

Formal grammars operate on any finite sequence where the elements are drawn from a finite set (e.g., text, molecules, or even neural network modules). While our experiments apply our CDoT grammar formalism to constrained decoding with LLMs, we will define it with sufficient abstraction that it could be used on sequences other than text. We lean on two alternative interpretations of formal grammars:

Grammar the proto programming language:

Computation theory studies grammar as a fundamental model of computation with the capacity to reach Turing-completeness—meaning any algorithm can be implemented, given an appropriate formal grammar. Following this intuition, we design the "knight-and-knave" experiment, where CDoT is used to implement the unit propagation algorithm to constrain the LLM generation.

Grammar the deduction system: Proposed by Pereira and Warren (1983), and used in PEG (Ford, 2004) and conjunctive grammars (Okhotin, 2001a; Mrykhin and Okhotin, 2023; Barash and Okhotin, 2014; Okhotin, 2013), under this interpretation, instead of saying a sequence $t_0 \dots t_n$ can be parsed into or derived from non-terminal X , we say $t_0 \dots t_n$ has property X . Consequently, grammar rule $X \rightarrow A B$ means a sequence of property X can be obtained by concatenating sequences of property A then B . This interpretation inspires us to blend grammar with deontic logic (von Wright, 1951), from which we inherit the name and the symbols (\diamond , \square). We will also depend on this interpretation to define formal semantics for CDoT.

2.2 Reach for Context-Sensitiveness

Classical constrained decoding use cases, such as JSON schema enforcement, are based on context-free grammars, implying a tree-shaped latent structure. Context-sensitive grammars additionally allow references across subtrees, for example, ensuring that variable references in generated code refer back to valid variable definitions.

Figure 1 illustrates the challenge of handling context-sensitiveness during LLM generation. Assume that we wish to enforce a constraint where the SYMPTOM of a diagnosis constrains which CONDITIONS are allowable. Note that constrained decoding is an incremental process, so the ancestor nodes are unavailable at the time a constraint must be applied. Ignoring ancestor structure and using only precedence, as in PiCARD (Scholak et al., 2021) and partially inherited by IterGen (Ugare et al., 2025), will result in the constraint firing incorrectly for the bottom tree in figure 1. An alternative adopted by IterGen (Ugare et al., 2025) and GenLM Control (Loula et al., 2025) is to allow ancestor structures to be considered by invoking expensive backtracking and/or speculative look-ahead, where ancestor structures are allowed to emerge before constraints are applied. ChopChop (Nagy et al., 2025) employs more powerful formalism to reason over partial syntax trees, but still cannot produce a complete logit mask and relies on backtracking incorrect token instead. There have also been specialized solutions for particular context-sensitive constraints, such as indentation in Python syntax (Melcer et al., 2024).

We will instead design a cross reference mechanism that can be resolved even when ancestor structure is non-determinant.

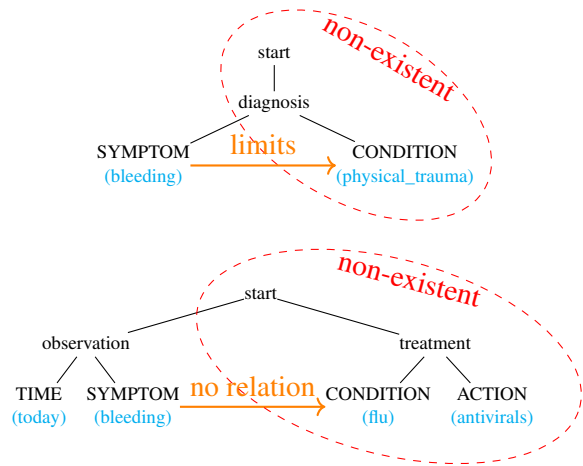


Figure 1: Example of challenges in enforcing a constraint where the SYMPTOM of a diagnosis constrains which CONDITIONS are allowable. Because constrained decoding is incremental, the ancestors of CONDITION are not known when it is being generated, yet the constraint must fire in the top tree but not the bottom.

2.3 Efficiency

Despite a diverse selection of mildly context-sensitive grammars (Joshi and Schabes, 1997; Vijay-Shanker and Weir, 1993; Okhotin, 2001b; Mrykhin and Okhotin, 2023; Barash and Okhotin, 2014), such grammars are often not designed with the goal of efficient incremental parsing in mind. For example, Tree-Adjoining Grammars have an $O(n^8)$ step time-complexity (Schabes, 1991). The recently published grammar with one-sided context (Barash and Okhotin, 2014) within the conjunctive grammar family overlaps with the goal of constrained decoding, but is currently a theoretical construct applied only to toy grammars. We design our formalism with the goal of parsing at $O(n^2)$ step with GPU acceleration, inline with general context-free grammars (Azimov and Grigorev, 2017).

2.4 The Solution

In a chart-based recognizer like Valiant (Valiant, 1975), given a subtree root, it is easy to find both the terminals in this subtree and the immediate children of the subtree root. Therefore, we design our cross-subtree references to have the following form: Given a subtree root, terminals under one of its children (consequent) are conditioned on terminals under a child to the left (antecedent). To provide further expressiveness, we deviate from the traditional practice of treating terminals as opaque objects, instead allowing predicates to be placed

on the consequent terminals. Additionally, Azimov and Grigorev (2017) shows that Valiant recognizer can be GPU accelerated and outperform the high performance general context-free parser GLL.

3 CDoT Example

Before formally defining CDoT, we provide some intuition with a concrete example of the formalism.

CDoT extends context free grammars (CFGs; formally defined in appendix A) with *deontic rules*, which constrain not only the structure of derivations, but also the *obligations* and *permissions* of terminals. A *Conditional Deontics over Terminals (CDoT)* grammar is a pair $G^+ = (G, D)$ where

- $G = (V, \Sigma, R, S)$ is a context-free grammar (possibly with regular-expression terminals),
- D is a finite set of deontic rules, which constrain terminal instances generated by G .

Each deontic rule has an *anchor*, identifying the scope of the rule, and a *deontic function*, which maps terminal instances to a pair $(\psi^\square, \psi^\diamond)$, where ψ^\square is a set of regular expressions (necessary conditions) and ψ^\diamond is a set of regular expressions (permitted conditions).

Non-conditional deontic rules. A non-conditional rule has the form

$$X^{\text{consequent}} \ll t^{\text{consequent}},$$

with $X^{\text{consequent}} \in V$ and $t^{\text{consequent}} \in \Sigma$. When the deontic function of a non-conditional deontic rule returns (N^\square, P^\diamond) , the semantics are that all terminals descended from X must *collectively* satisfy every $r \in N^\square$, and each individual terminal must match at least one $r \in P^\diamond$.

Example (Balanced Meals). Given the grammar:

$meal \rightarrow food\ beverage$
 $food \rightarrow \text{DISH } food \mid \text{DISH}$
 $beverage \rightarrow \text{BEVERAGE } beverage \mid \text{BEVERAGE}.$

The following deontic rule enforces that every meal includes grains, meat, and vegetables (\square), while also allowing optional extras like soup or steamed_egg (\diamond).

```

anchor : food << DISH
def enforce_meal():
    return ({ rice|bread, chicken|fish,
              salad|roasted_vegetable } $\square$ ,
            { rice, bread, chicken, fish,
              salad, roasted_vegetable,
              steamed_egg, soup } $\diamond$ )

```

Conditional deontic rules. A conditional rule has the form

$$t^{\text{ant}} \gg X^{\text{ant}} B \dots D > X^{\text{par}} < X^{\text{con}} \ll t^{\text{con}}$$

where $t^{\text{ant}}, t^{\text{con}} \in \Sigma$, and $B, D, X^{\text{ant}}, X^{\text{con}} \in V$ are variables appearing in a parent production $X^{\text{par}} \rightarrow \dots X^{\text{ant}} B \dots D X^{\text{con}} \dots$. The deontic function of a non-conditional deontic rule takes an antecedent terminal instance t^{ant} as input, and uses it to determine and return constraints (N^\square, P^\diamond) , on the consequent terminal t^{con} .

Example (Food-Drink Pairings). Given the same grammar as the previous example, the following deontic rule enforces that when a food and a beverage are paired within a meal, steaks require red wine, fish requires white wine, and spicy noodles require soda.

```

anchor : DISH >> food > meal < beverage << BEVERAGE
def drink2food(terminal_instance):
    match terminal_instance.text:
        case 'steak':
            return {red_wine} $\square$ , {} $\diamond$ 
        case 'fish':
            return {white_wine} $\square$ , {} $\diamond$ 
        case 'spicy_noodle':
            return {soda} $\square$ , {} $\diamond$ 

```

Notice that, while a single antecedent variable *food* can cover multiple antecedent DISH variables, each invocation of the deontic function *match_drink_to_food* sees only a single instance of DISH. This prevents repetitive calling of the deontic function over different combinations. At the same time, the conjunctive nature of \square and disjunctive nature of \diamond will still naturally combine so that, if both fish and steak are ordered, both red and white wine are requested.

4 Formal Semantics of CDoT

4.1 Basic Notations

Terminal Instance (t_j) is a subspan of the generation process. Each terminal instance is defined by a sequential index (j), a terminal type, and a user-defined payload. In the case of constrained decoding, a terminal instance might have type “VAR_NAME” and payload of token sequence [“start”, “_”, “time”].

Type-based Filtering ($[t_i \dots t_j]_b$) extracts the set of all terminal instances within the span that have terminal type b in sequence $t_i \dots t_j$. Formally,

$$[t_i \dots t_j]_b \triangleq \{t_k \mid i \leq k \leq j \wedge \text{type}(t_k) = b\}$$

250	Predicate (ρ) We write $\rho(t_i)$ to indicate that the	286
251	single terminal instance t_i satisfies property ρ . For	287
252	example, regular expressions are good predicate	288
253	for constrained decoding (i.e. a variable name must	289
254	start with a certain prefix). Alternatively, a neural	290
255	probe or other heuristics might also be predicates.	291
256	Predicate Verifier A predicate verifier is needed	292
257	to confirm a terminal instance satisfies the predi-	293
258	cate, for example, matching regular expressions or	294
259	comparing to a threshold value.	295
260	Predicate Enforcer A predicate enforcer pro-	296
261	duces a new terminal instance that satisfies the predi-	297
262	cate, for example, running a regular expression	298
263	constrainer or cutting off values. An all-permitting	299
264	enforcer is acceptable and appropriate if permit-	300
265	ted is set to universal (e.g., when using a neural	301
266	probe predicate). Enforcers need to be mergable,	302
267	with the merged enforcer producing only t_i s.t.	303
268	$\rho_1(t_i) \vee \rho_2(t_i)$. Enforcers needs to be at least sound,	304
269	if completeness is not achievable.	305
270	4.2 Terminal Specification	306
271	CDoT deviates from traditional formal grammars	307
272	and allows grammar rules to specify properties to	308
273	terminal instances under their scope.	309
274	Terminal Specification (ψ) is a set of predicates.	310
275	Terminal $t_i \dots t_j$ satisfies a specification ψ is writ-	311
276	ten as $t_i \dots t_j \in \psi$. Notice that the \in operator is	312
277	overloaded as specification is a set of predicates.	313
	Necessary Specification (ψ^\square) is the conceptual	314
	“lower bound”. Each predicate in a necessary speci-	315
	fication must be satisfied by at least one consequent	316
	terminal instance, formally	317
	$t_i \dots t_j \in \psi^\square \iff \forall \rho \in \psi^\square : \exists k \in [i, j] : \rho(t_k)$	318
	Permitted Specification (ψ^\diamond) is the conceptual	319
	“upper bound”. Each consequent terminal instance	320
	must satisfy at least one predicate in the permitted	321
	specification, formally	322
	$t_i \dots t_j \in \psi^\diamond \iff \forall k \in [i, j] : \exists \rho \in \psi^\diamond : \rho(t_k)$	323
278	Obligation Implies Can we stipulate that every	324
279	predicate in a necessary specification is automati-	325
280	cally added to the matching permitted specification.	326
281	Aggregating Specifications is achieved by the	327
282	union $\psi_A \cup \psi_B$. Intuitively, the necessary speci-	328
283	fication tightens (more predicates to be satisfied),	329
284	and the permitted specification loosens (more predi-	330
285	cates to choose from).	331
	4.3 Deontic Function	332
	A <i>deontic function</i> maps a single terminal instance	333
	to a token specification. This serves a similar role	
	as functions in recursive descent parsers or IterGen.	
	Invoking unrestricted functions during parsing in-	
	creases both expressiveness and convenience, al-	
	lowing, for example, a query to an external tool or	
	a recursive LLM call during generation.	
	We define the function to take a single termi-	
	nal instance, rather than all antecedent instances to	
	reduce the number of deontic function calls from	
	$O(n^2)$ to $O(n)$. We compute each terminal’s de-	
	ontic result once upon completion, then discard	
	the terminal instance – useful for large payloads	
	such as LLM hidden states. Due to the aggrega-	
	tion behavior of necessary and permitted semantics	
	described in section 4.2, taking a single terminal	
	instance often gets the same result specification.	
	4.3.1 Notations of deontic function	
	Concretely, we write a deontic function as a python	
	function that takes a terminal instance and returns	
	a pair of necessary and permitted specifications	
	as a tuple. Given a deontic function, we write	
	$\square t_i$ or $\diamond t_i$ to indicate the necessary or permitted	
	specification obtained by applying such deontic	
	function to t_i . We introduce shorthands:	
	$\square(t_i \dots t_j) \triangleq \square(t_i) \cup \square(t_{i+1}) \cup \dots \cup \square(t_j)$	
	$\diamond(t_i \dots t_j) \triangleq \diamond(t_i) \cup \diamond(t_{i+1}) \cup \dots \cup \diamond(t_j)$	
	4.4 The Inference Rules	
	Terminal Instance Partitions (F, P, R) Under	
	a CDoT subtree, each terminal instance belongs	
	to one of three partitions. F are free terminal in-	
	stances not under deontic constraints (e.g., keyword	
	<i>SELECT</i> in a SQL query is likely unconstrained).	
	P are the permitted terminal instances under de-	
	ontic constraints and receive permission (\diamond) from	
	the current subtree. R are the remaining terminal	
	instances under deontic constraints, but no permis-	
	sion received from the current subtree (e.g., an	
	invocation of variable from outer scope).	
	Assertion ($A(t_i \dots t_j \mid \langle R, F \rangle)$) is read as	
	$t_i \dots t_j$ has property A , conditioned on the termi-	
	nal instances in $R \subseteq t_i \dots t_j$ receiving permission	
	from ancestors. We say $t_0 \dots t_n$ is in the language	
	of grammar if and only if $START(t_0 \dots t_n \mid \emptyset)$.	
	4.4.1 Context-Free Rules	
	We rewrite context-free rules into Chomsky normal	
	form to include the deontic condition.	

Terminal Rule $X \rightarrow a$ generates a single terminal from a non-terminal. Because it is not under deontic constraint yet, t_i belongs to partition F .

$$\frac{\text{type}(t_i) = a}{X(t_i \mid \langle \emptyset, \{t_i\} \rangle)} \quad (\text{Terminal})$$

Context-free Rule $X \rightarrow A B$ simply aggregates and passes on the deontic condition.

$$\frac{A(t_i \dots t_j \mid \langle R_A, F_A \rangle) \quad B(t_{j+1} \dots t_k \mid \langle R_B, F_B \rangle)}{X(t_i \dots t_k \mid \langle R_A \cup R_B, F_A \cup F_B \rangle)} \quad (\text{Context-free})$$

Since neither rule is conditioned on deontic premises nor adds terminals to the R partition, CDoT behaves identically to a context-free grammar in the absence of deontic rules.

4.4.2 Strict Permitted Rule

The crux of the problem is defining behavior when deontic rules are nested. We begin with a naive but useful semantics termed "strict permitted."

Strict- \diamond rule requires the consequent to strictly satisfy the terminal specification returned by the antecedent's deontic function, disregarding any constraints from ancestor rules or rules nested within the consequent. As a result, all terminals of consequent type (b) under the consequent non-terminal (B) belong to partition P for parent X , and therefore do not appear in F_x or R_x , as shown in the inference rule.

$$\frac{A(t_i \dots t_j \mid \langle R_A, F_A \rangle) \quad B(t_{j+1} \dots t_k \mid \langle R_B, F_B \rangle) \quad \begin{array}{l} [t_{j+1} \dots t_k]_b \in \square([t_i \dots t_j]_a) \\ [t_{j+1} \dots t_k]_b \in \diamond([t_i \dots t_j]_a) \end{array}}{X(t_i \dots t_k \mid \langle R_X, F_X \rangle) \quad \begin{array}{l} R_X = (R_A \cup R_B) \setminus [t_{j+1} \dots t_k]_b \\ F_X = (F_A \cup F_B) \setminus [t_{j+1} \dots t_k]_b \end{array}} \quad (\text{Deontic - Strict } \diamond)$$

Upon closer examination, one may realize that in a system with only context-free and strict- \diamond again keeps $R = \emptyset$ and never conditions on $\langle R, F \rangle$. Thus, a CYK-like parser can be easily implemented if the entire terminal instance sequence is available.

However, decidability issues arise when we seek to parse incrementally while maintaining the valid-prefix property. Because each consequent terminal instance can further serve as an antecedent in descendant subtrees—used as inputs to another deontic function—the parser must determine whether future terminals could yield a terminal specification incompatible with the current one.

4.4.3 Union Permitted Rule

Beside tractability, the strict- \diamond semantic sometimes also does not capture desired expressions when we need to aggregate specifications from multiple antecedent subtrees, for example, variables from different scopes.

Union- \diamond rule While keeping the necessary(\square) behavior unchanged, we allow subtrees to put terminal instances under constraint but defer permission giving to ancestors. This is reflected in removal of \diamond condition in the premise and updated R_X definition. Notice that we are overloading the \setminus operator between terminal instance collection and specification, $(R_B \cup F_B) \setminus \diamond [t_i \dots t_j]_a \triangleq \{t \mid t \in [R_B \cup F_B]_b \wedge t \notin \diamond [t_i \dots t_j]_a\}$.

$$\frac{A(t_i \dots t_j \mid \langle R_A, F_A \rangle) \quad B(t_{j+1} \dots t_k \mid \langle R_B, F_B \rangle) \quad [t_{j+1} \dots t_k]_b \in \square([t_i \dots t_j]_a)}{X(t_i \dots t_k \mid \langle R_X, F_X \rangle) \quad \begin{array}{l} R_X = R_A \cup (R_B \cup F_B \setminus \diamond [t_i \dots t_j]_a) \\ F_X = (F_A \cup F_B) \setminus [t_{j+1} \dots t_k]_b \end{array}} \quad (\text{Deontic - Union } \diamond)$$

5 Parser

Building on past work on GPU accelerated parsers, we construct our parser using two waves of extended Incremental Valiant Recognizers. The detail of parser is shown in [appendix C](#).

Parser Behavior The constrainer is complete, meaning it would not falsely deny a valid token.

However, due to the decidability issue discussed in [section 4.4.2](#), we only achieve eventual soundness: a strict- \diamond deontic rule instance is pruned only upon its first direct violation. At that point, if the grammar does not provide an alternative path, generation terminates. Our Knight-and-Knave grammar ([figure 7](#)) demonstrates how to provide such a fallback. This pattern commonly arises when constraining a thinking trace, where deontic rules are used for early termination of an attempt.

Furthermore, under union- \diamond semantics, consider a prefix where two deontic rule instances satisfy the following conditions: (1) they share at least one potential derivation tree, and (2) each can also belong to derivation trees that exclude the other. In such cases, permission granted in one derivation tree may incorrectly propagate to another, causing cross contamination of permissions and looser than

specified constraint. We demonstrate this behavior in [appendix B](#).

Efficiency By preserving the loop structure of incremental Valiant recognizer and only modify the constant-time loop body, we achieve $O(n^2)$ step complexity. Our current implementation is CPU bound and the token-to-token time is similar to HuggingFace implementation of Llama3 3B models. However, an efficient implementation of incremental Valiant recognizer steps at around 100 μ s with $n = 100$ and $|G^{CNF}| = 3000$ using a single stream multiprocessor (out of > 100) on a Nvidia H100 GPU.

6 Experiments and Results

6.1 Knight and Knave

Knight and Knave is a logic puzzle where each character is either a knight (always tells the truth) or a knave (always lies). The task is to infer each character’s type from their statements. For example, if Alice says “Alice and Bob have the same identity”, and Bob says “Alice is lying”, we must deduce that Alice is a knave and Bob is a knight. [Xie et al. \(2024\)](#) provides a readily usable dataset containing both textual and logical form of puzzles with 2-8 characters.

6.1.1 Preparation: conjunctive normal form and unit prop

Conjunctive Normal Form is a conjunction of *clauses*, each clause being a disjunction of *literals* (a variable or its negation). For example, $(x_1 \vee \neg x_3 \vee x_4) \wedge (\neg x_1 \vee x_3) \wedge (x_1 \vee x_5)$. An arbitrary boolean expression can be converted into equalsatisfiable conjunctive normal form efficiently via the Tseytin transformation ([Tseytin, 1983](#)) with a caveat of adding $O(n)$ auxiliary variables.

Unit Propagation is an algorithm that deduces assignments based on committed assignments of other literals. It exploits the fact that each clause in conjunctive normal form must have at least one literal to be true. Therefore, if all but one literal in a clause are false, the last remaining literal must be true. In the example above, if we assume x_1 is true, then x_3 must also be true, because it is the last remaining option to keep the second clause satisfied.

6.1.2 Constrainer Design

[Figure 7](#) shows the context free grammar we employ for the output of the LLM on the Knight and

Knave problem. It ensures that the model outputs a sequence of guesses. Each guess assigns roles to characters and boolean values to auxiliary variables. The grammar alone will not ensure the desired behavior, so we add the following deontic rules.

To ensure that *sufficient_assignment_seq* is recognized iff an *assignment_seq* commits to at least one literal from each clause is, we add the deontic rule:

```

anchor : sufficient_assignment_seq << ASSIGNMENT
% attr : hidden\_permitted
def correct_guess():
    return ( $\square\{ "x_i \dots x_j"
                \text{for } x_i \vee \dots \vee x_j
                \text{in input\_clauses } \},
            \diamond\{ \} )$ 
```

We add a similar rule for constraining *contradictory_assignment_seq*.

To ensure that a REPEATED_ASSIGNMENT is recognized iff it replicates a previous ASSIGNMENT within this guess, we add the deontic rule:

```

anchor : ASSIGNMENT >> assignment_seq_semi_1
        > repetition_assignment_seq
        < repeating_assignment
        << REPEATED_ASSIGNMENT
def repetition(terminal):
    return ( $\square\{ \}, \diamond\{ \text{terminal} \} )$ 
```

Deontic rules for implementing the unit propagation algorithm are defined for each x_k in the input clauses. We add context-free rules of the form:

$$\text{unit_prop} \rightarrow x_k_condition_semi_2 \quad x_k_assignment$$

We add deontic rules of the form:

```

anchor : x_k_assignment << ASSIGNMENT
% attr : self\_sufficient\_necessary
def x_k_condition():
    return ( $\square\{ "x_k" \}, \diamond\{ \} )$ 
```

For each input clause c that contains x_k , we add a deontic rule:

```

anchor : x_k_condition << ASSIGNMENT
% attr : hidden\_permitted
def x_k_assignment():
    return ( $\square\{ "NEGATE\_LITERAL(x_m)"
                \text{for } x_m \text{ in } c \text{ if } m \neq k \},
            \diamond\{ "x_m"
                \text{if } x_m \text{ exists in any clause}
                \text{and } m \neq k \} )$ 
```

These deontic rules result in the real-time triggering of unit propagation and the real-time detection of good and bad guesses. As the LLM commits to literals, unit propagation will trigger immediately upon conditions being met, and will constrain generation to only entailed literals. Bad guesses will be detected immediately upon their first contradiction and generation will be forced to start a new

Model	Constrained		Unconstrained	
	Thinking	Non-thinking	Thinking	Thinking
	512	512	512	2048
Llama-3.2-1B-Instruct	0.52	0.04	0.04	OOM
Llama-3.2-3B-Instruct	0.65	0.13	0.13	OOM
Qwen/Qwen3-0.6B	0.61	0.07	0.07	0.36
Qwen/Qwen3-1.7B	0.69	0.03	0.03	0.55

Table 1: Knight and Knave test: solve rate of constrained vs. unconstrained models with 512 or 2048 token limits. Llama-3 models do not have a thinking switch, but does exhibit chain of thought-like behavior. Xie et al. (2024) reported 0.14 for Llama-3-8B-instruct in the same setting.

guess. Good guesses will be detected immediately once the formula is solved and generation will be immediately terminated.

6.1.3 Baselines

We compare against only unconstrained models for the lack of alternatives. PiCARD is hard-coded to SQL grammar. IterGen can only backtrack an incorrect generation, while our constraints proactively forces an assignment once it is fixed. Furthermore, it is unclear how many assignments should be backtracked once conflict arises. ChopChop has similar backtracking issue and is not publicly released.

6.1.4 Results

We set LLMs to solve puzzles in a 0-shot setting. Models are provided with both the textual and logical forms of the puzzle in the prompt. The results are shown in table 1. When the output is limited to 512 tokens, constrained models achieve about a 50-point higher solve rates than non-constrained models, and even outperform non-constrained models that are allowed to generate 4 times more tokens.

We also plot the number of output tokens for successful (dark) and failed (light) LLM responses across constrained (blue) and unconstrained (orange) Qwen3-1.7B, shown in figure 2. The constrained model shows a skewed distribution, where a larger portion of puzzles are solved with lesser tokens, and failures don't show up until around 512 tokens. (Qualitatively, we observe that at this point, the constrained model gets stuck in loops.) The non-constrained model forms a flat bell curve distribution, and requires many more output tokens before starting to solve puzzles.

We observe that the constrained model can often solve the puzzle on the first attempt (figure 3). This indicates that the constrainer and LLM are organ-

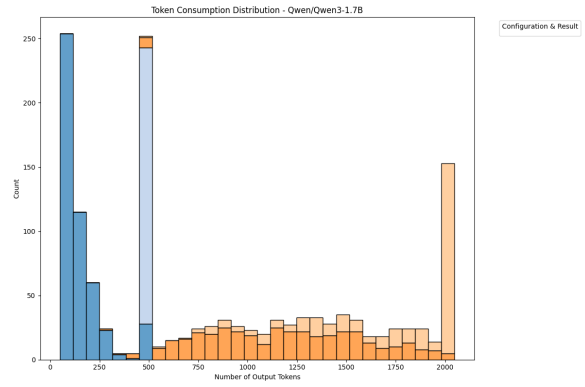


Figure 2: Knight and Knave test: token efficiency of successes (dark) vs failures (light) of constrained (blue) vs. unconstrained (orange) models. Constrained models were cut off at 512 tokens, unconstrained at 2048 tokens.

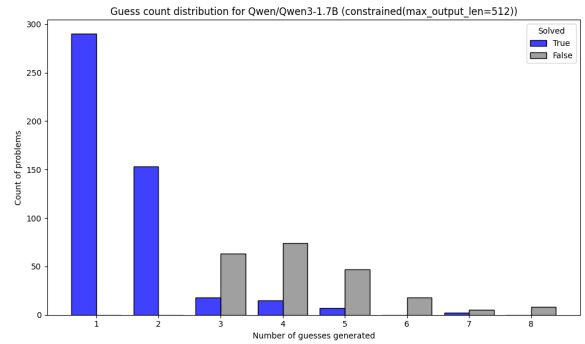


Figure 3: Knight and Knave test: constrained models solve many problems in the first guess. Notice that the solve rate at first guess is significantly higher than random guess success rate, indicating organic LLM/Constrainer collaboration instead of constrainer dominance.

ically collaborating. If the constrainer alone was driving success, we would expect the distribution of guesses to be closer to random uniform, and it would be rare to solve the puzzle on the first guess.

When we break down the solve rate by the size of the puzzle as in figure 6 in appendix D, we observe that the constrained models scale much better and still achieve solve rate of around 0.50, in contrast to under 0.20 for non-constrained models.

6.2 Text-to-SQL

Spider (Yu et al., 2019) is a text-to-sql task, where the model is provided with a database schema and a natural text question, and it is expected to answer the question with an generated SQL query.

We finetune an LLM with the constrainer on during both finetuning and inferencing. For the constrainer, we use a grammar extracted from the SQLite source code as the base context-free gram-

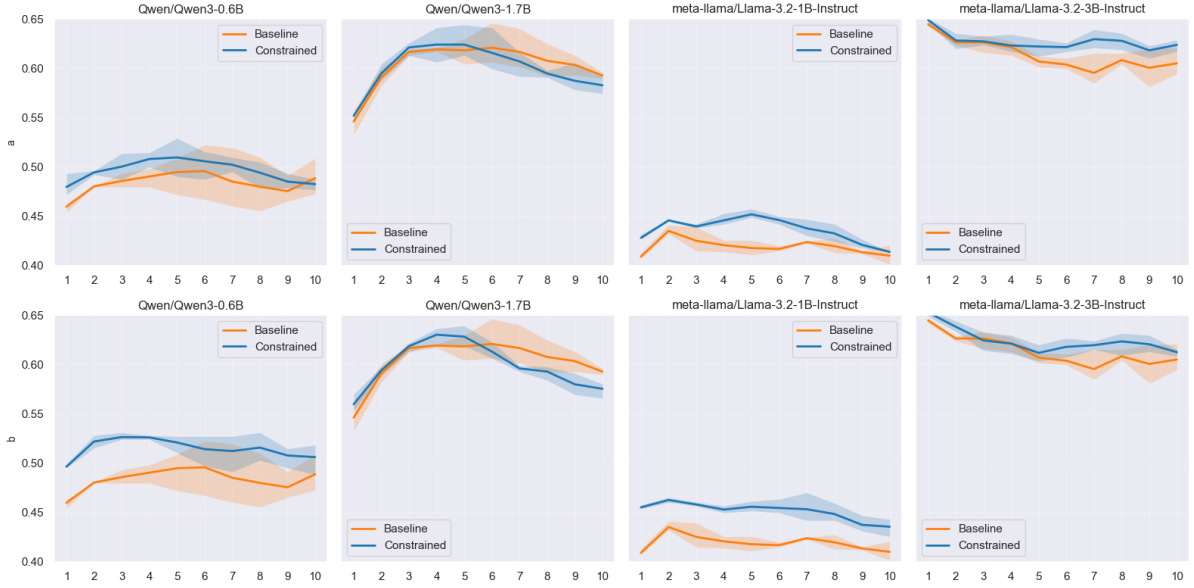


Figure 4: Text-to-SQL test: Execution accuracy of constrained vs non-constrained models over 3 runs and 10 epochs. Row (a) uses selected columns in “SELECT” clause to limit the remaining identifiers to that of databases containing them. Row (b) additionally constrains the “SELECT” clause to identifiers of all databases. Constrained decoding helps three of the four models.

mar, then add constraints that allow only identifiers that come from databases containing the columns in the “SELECT” clause. We also run experiments with the additional constraint of limiting identifiers in the “SELECT” clause to the scope of all databases. Formally:

```

575 anchor : ID >> select_clause > select < from_clause << ID
576 def allow_containing_database (terminal):
577     return (□ {},
578             ◇ {name if name is an
579                identifier of a database
580                 containing terminal})

```

```

581 anchor : select_clause << ID
582 def correct_guess ():
583     return (□ {},
584             ◇ {name if name is a table or
585                column of any database})

```

For both constrained and unconstrained models, we remove tokens from the tokenizer that do not align with SQL grammar terminals, such as “.Name” and “);”. LLM performance is known to suffer when such tokens are not handled by the constrainer (Beurer-Kellner et al., 2024). For the unconstrained model, we add the target database schema to the model input, so the model has the same information on the target database that is available to the constrained model.

The results are shown in figure 4. “Qwen3-0.7B”, “Llama-3.2-1B-Instruct”, and “Llama-3.2-3B-Instruct” benefit from the constrainer with ei-

ther constraint configuration. However, “Qwen3-1.7B” does not benefit from the constrainer.

Baseline is again restricted to unconstrained models. PiCARD was implemented for T5 and is now defunct. Backtracking systems such as IterGen cannot be used during finetuning for its inability to produce logit mask, while our system requires finetuning in SQL generation due to token/terminal alignment issue mentioned above. For reference, IterGen reports improvement from 25.5 to 30.9 in accuracy in a non-finetuned Llama-3.2-1B model, while our finetuned model shows similar delta but at a higher starting point in the 40s.

7 Conclusion

We exploit chart-based parsing to design a novel cross-subtree reference mechanism that achieves mild context-sensitivity. By allowing predicates to be attached to terminals, we expand the expressiveness. Attaching deontic functions to grammar rules improves both expressiveness and usability. Our parser’s efficiency builds on prior work in GPU-accelerated parsing.

We use our formalism to implement the unit propagation algorithm, guiding reasoning traces on Knight-and-Knave puzzles and achieve significant performance gain at reduced token budget. Our constrainer also proves helpful in text-to-SQL.

626 Limitations

627 Although our conainer is efficient with GPU ac-
628 celeration, there are cases where it fails to reject
629 incorrect tokens, resulting in constraints that are
630 looser than specified. Future work could identify
631 restrictions on the underlying grammar or deontic
632 functions that guarantee tightness of constraints.
633 Additionally, our conainer has difficulty hand-
634 dling tokens that span multiple terminals, limiting
635 its usefulness in code generation tasks.

636 References

637 Rustam Azimov and Semyon V. Grigorev. 2017.
638 [Graph parsing by matrix multiplication](#). *CoRR*,
639 abs/1707.01007.

640 Mikhail Barash and Alexander Okhotin. 2014. [An ex-
641 tension of context-free grammars with one-sided con-
642 text specifications](#). *Information and Computation*,
643 237:268–293.

644 Luca Beurer-Kellner, Marc Fischer, and Martin Vechev.
645 2024. [Guiding llms the right way: Fast, non-invasive
646 constrained generation](#). *Preprint*, arXiv:2403.06988.

647 Zifeng Cheng, Jinwei Gan, Zhiwei Jiang, Cong Wang,
648 Yafeng Yin, Xiang Luo, Yuchen Fu, and Qing Gu.
649 2025. [Steering when necessary: Flexible steering
650 large language models with backtracking](#). *Preprint*,
651 arXiv:2508.17621.

652 Bryan Ford. 2004. [Parsing expression grammars: a
653 recognition-based syntactic foundation](#). In *Proceed-
654 ings of the 31st ACM SIGPLAN-SIGACT Symposium
655 on Principles of Programming Languages*, POPL ’04,
656 page 111–122, New York, NY, USA. Association for
657 Computing Machinery.

658 Aravind K. Joshi and Yves Schabes. 1997. *Tree-
659 Adjoining Grammars*, pages 69–123. Springer Berlin
660 Heidelberg, Berlin, Heidelberg.

661 Bruce W. Lee, Inkit Padhi, Karthikeyan Natesan Rama-
662 murthy, Erik Miebling, Pierre Dognin, Manish Na-
663 gireddy, and Amit Dhurandhar. 2025. [Programming
664 refusal with conditional activation steering](#). *Preprint*,
665 arXiv:2409.05907.

666 João Loula, Benjamin LeBrun, Li Du, Ben Lipkin,
667 Clemente Pasti, Gabriel Grand, Tianyu Liu, Yahya
668 Emara, Marjorie Freedman, Jason Eisner, Ryan Cot-
669 terell, Vikash Mansinghka, Alexander K. Lew, Tim
670 Vieira, and Timothy J. O’Donnell. 2025. [Syntactic
671 and semantic control of large language models via
672 sequential monte carlo](#). *Preprint*, arXiv:2504.13139.

673 Daniel Melcer, Nathan Fulton, Sanjay Krishna Gouda,
674 and Haifeng Qian. 2024. [Constrained decoding for
675 fill-in-the-middle code language models via efficient
676 left and right quotienting of context-sensitive gram-
677 mars](#). *Preprint*, arXiv:2402.17988.

Mikhail Mrykhin and Alexander Okhotin. 2023. [The
678 hardest language for grammars with context opera-
679 tors](#). *Theoretical Computer Science*, 958:113829.
680

Shaan Nagy, Timothy Zhou, Nadia Polikarpova, and
681 Loris D’Antoni. 2025. [Chopchop: a programmable
682 framework for semantically constraining the output
683 of language models](#). *Preprint*, arXiv:2509.00360.
684

Alexander Okhotin. 2001a. [Conjunctive grammars](#).
685 *Journal of Automata, Languages and Combinatorics*,
686 6(4):519–535.
687

Alexander Okhotin. 2001b. [Conjunctive grammars](#). *J.
688 Autom. Lang. Comb.*, 6(4):519–535.
689

Alexander Okhotin. 2013. [Conjunctive and boolean
690 grammars: The true general case of the context-free
691 grammars](#). *Computer Science Review*, 9:27–59.
692

Fernando C. N. Pereira and David H. D. Warren. 1983.
693 [Parsing as deduction](#). In *21st Annual Meeting of
694 the Association for Computational Linguistics*, pages
695 137–144, Cambridge, Massachusetts, USA. Associa-
696 tion for Computational Linguistics.
697

Gabriel Poesia, Oleksandr Polozov, Vu Le, Ashish Ti-
698 wari, Gustavo Soares, Christopher Meek, and Sumit
699 Gulwani. 2022. [Synchromesh: Reliable code gen-
700 eration from pre-trained language models](#). *arXiv
701 preprint arXiv:2201.11227*.
702

Yves Schabes. 1991. [The valid prefix property and
703 left to right parsing of Tree-Adjoining Grammar](#). In
704 *Proceedings of the Second International Workshop on
705 Parsing Technologies*, pages 21–30, Cancun, Mexico.
706 Association for Computational Linguistics.
707

Torsten Scholak, Nathan Schucher, and Dzmitry Bah-
708 danau. 2021. [PICARD: Parsing incrementally for
709 constrained auto-regressive decoding from language
710 models](#). In *Proceedings of the 2021 Conference on
711 Empirical Methods in Natural Language Processing*,
712 pages 9895–9901, Online and Punta Cana, Domini-
713 can Republic. Association for Computational Lin-
714 guistics.
715

G. S. Tseitin. 1983. [On the Complexity of Derivation
716 in Propositional Calculus](#), pages 466–483. Springer
717 Berlin Heidelberg, Berlin, Heidelberg.
718

Shubham Ugare, Rohan Gumaste, Tarun Suresh, Gagan-
719 deep Singh, and Sasa Misailovic. 2025. [Itergen: It-
720 erative semantic-aware structured LLM generation
721 with backtracking](#). In *The Thirteenth International
722 Conference on Learning Representations*.
723

Leslie G. Valiant. 1975. [General context-free recogni-
724 tion in less than cubic time](#). *J. Comput. Syst. Sci.*,
725 10:308–315.
726

K Vijay-Shanker and David J. Weir. 1993. [Parsing some
727 constrained grammar formalisms](#). *Computational
728 Linguistics*, 19(4):591–636.
729

G. H. von Wright. 1951. [Deontic logic](#). *Mind*,
730 60(237):1–15.
731

Chulin Xie, Yangsibo Huang, Chiyuan Zhang, Da Yu, Xinyun Chen, Bill Yuchen Lin, Bo Li, Badih Ghazi, and Ravi Kumar. 2024. [On memorization of large language models in logical reasoning](#).

Tao Yu, Rui Zhang, Kai Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir Radev. 2019. [Spider: A large-scale human-labeled dataset for complex and cross-domain semantic parsing and text-to-sql task](#). *Preprint*, arXiv:1809.08887.

Qingru Zhang, Chandan Singh, Liyuan Liu, Xiaodong Liu, Bin Yu, Jianfeng Gao, and Tuo Zhao. 2024. [Tell your model where to attend: Post-hoc attention steering for llms](#). *Preprint*, arXiv:2311.02262.

A Context-Free Grammars and Languages

A *context-free grammar (CFG)*, which is a 4-tuple $G = (V, \Sigma, R, S)$ where:

- V is a finite set of *nonterminal symbols* (also called variables), which represent syntactic categories that can be expanded,
- Σ is a finite set of *terminal symbols*, which are the basic alphabet symbols that appear in the strings of the language, with $V \cap \Sigma = \emptyset$,
- $R \subseteq V \times (V \cup \Sigma)^*$ is a finite set of *production rules* of the form $A \rightarrow \alpha$, where $A \in V$ (a nonterminal) and $\alpha \in (V \cup \Sigma)^*$ (a string of terminals and/or nonterminals),
- $S \in V$ is the designated *start symbol*.

The language generated by G is defined as $L(G) = \{w \in \Sigma^* \mid S \Rightarrow^* w\}$, where \Rightarrow^* denotes the reflexive, transitive closure of the single-step derivation relation \Rightarrow induced by R . A language $L \subseteq \Sigma^*$ is called *context-free* if there exists a CFG G such that $L = L(G)$.

B Union- \diamond Counter Example

Under union- \diamond semantics, consider a prefix where two deontic rule instances satisfy the following conditions: (1) they share at least one potential derivation tree, and (2) each can also belong to derivation trees that exclude the other. In such cases, permission granted in one derivation tree may incorrectly propagate to another, causing cross contamination of permissions.

We illustrate this with two examples: one demonstrating correct pruning behavior, and one demonstrating a failure case.

	t_2	t_3
A	alive (prop. from A_2)	pruned
A_2	alive ($t_2 \in \diamond A_2$)	pruned
A_3	pruned	—
tree 1	alive ($t_2 \in \diamond A_2$)	pruned
tree 2	pruned	—

Table 2: Case 1 (sound and complete): after committing, both trees are pruned.

	t_2	t_3
A	alive (prop. from A_2)	alive (prop. from A_3)
A_2	alive ($t_2 \in \diamond A_2$)	pruned
A_3	alive (prop. from A_4)	alive ($t_3 \in \diamond A_3$)
A_4	alive (prop. from A_3)	alive ($t_3 \in \diamond A_4$)
tree 1	alive ($t_2 \in \diamond A_2$)	pruned
tree 2	alive	alive
tree 3	alive ($t_2 \in \diamond A_4$)	alive ($t_3 \in \diamond A_3$)

Table 3: Case 2 (complete but not sound): propagation keeps some combinations alive that cannot co-occur in any single grammar-valid tree. Red boxes mark components falsely kept alive.

B.1 Case 1: Sound and complete (as strict as specified)

Consider trees in [figure 5](#). Assume a grammar allows tree 1 and tree 2, but does not have tree 3. Although node A is used in two separate trees, it does not cause problem. Tree 2 is pruned after committing t_2 for lack of permission. Tree 1 is kept alive for having A_2 , but pruned after committing t_3 for lack of permission. This is illustrated in [table 2](#).

B.2 Case 2: Complete but not sound (looser than specified)

Now assume a grammar allows tree 1, tree 2, and tree 3. Tree 1 and tree 2 as we know should have been pruned after committing t_2 and t_3 , but instead tree 2 is incorrectly kept alive. This is because node A is kept alive by being tree 1 and 2, node A_3 is kept alive by being in tree 3, collectively all components of tree 2 survives, as shown in [table 3](#) **ERROR**: no tree contains A , A_2 , and A_3 .

C Parser Details

Data Structure ($\mathbf{LF}[tier, row, X \rightarrow A B]$, $\mathbf{TIER}[row, X]$) Chart-like parsers such as CYK use a table indexed by *row*, *column* and non-terminal X . The Incremental Valiant recognizer uses the same table with the alternative index *tier*, *row*, X , where $tier = col + row$.

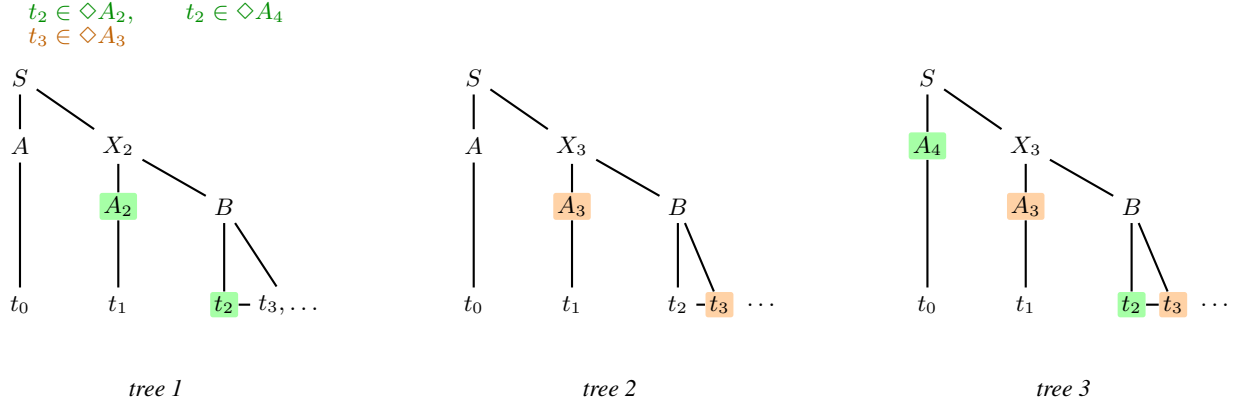


Figure 5: Three trees used in the counter example (highlighted nodes correspond to the committed constraints).

The CDoT parser inherits the boolean matrix $TIER[row, X]$, and the left-foot boolean tensor $\mathbf{LF}[tier, row, X \rightarrow A B]$, with the last dimension indexed by binarized grammar rules.

Necessary (\square) is enforced by denying the parent node from forming until the necessary condition is met. In the case that the underlying context-free grammar allows a finite number of consequent terminals, we assume the prescribed necessary requirement does not exceed this limit and the LLM completes the necessary requirement before the limit is exhausted. If this assumption is violated and the grammar does not provide an alternative path, generation will terminate.

Permitted (\diamond) enforcement allows an antecedent non-terminal A instance to be used in anchoring grammar rule $X \rightarrow A B$ until a violating consequent terminal instance is committed.

C.1 Preprocessing

After a new terminal instance t_n is completed, we need to find left feet that 1) have all necessary conditions **fulfilled** after accepting t_n , 2) can give permission to **enable** t_n . This can be achieved by simple dynamic programming shown as [algorithm 1](#). We omit proof for simplicity.

$$\forall j \leq \text{FULFILLER}[row, col] \Leftrightarrow [t_j \dots t_n]_b \in \square([t_{col} \dots t_{col+row}]_a) \quad (1)$$

$$\text{ENABLER}[row, col] \Leftrightarrow t_n \in \diamond([t_{col} \dots t_{col+row}]_a)$$

C.2 Commit Wave

We update the parser for newly formed terminal instance t_n . We call the persistent state alive left foot, $\mathbf{LF}_n^{\text{ALIVE}}$. This requires a bottom-up and a top-

down pass that can be run in parallel. The pseudo code can be found as [algorithm 2](#) and [algorithm 6](#).

C.3 Probe Wave

Probe wave finds all valid next terminal types and confirms each doentic left foot belongs to a derivation that supports current prefix $t_0 \dots t_n$.

Additionally, to confirm each consequent terminal type have a unconstrained path, we run an additional bottom-up incremental Valiant step with set of grammar rules excluding corresponding deontic anchors. If a terminal type has a unconstrained path, the constrainer receive its context-free definition as specification.

Shown as [algorithm 10](#) and [algorithm 11](#)

C.4 Post Processing

If a antecedent terminal is under a valid deontic left feet, then ψ^\diamond resulting from the corresponding deontic function should be added to enforcer during generation of the consequent terminal ([algorithm 13](#)). This takes a simple top-down dynamic programming with $O(|D| \cdot n^2)$ complexity.

D Knight and Knave by Puzzle Size

E Knight and Knave Constrainer Grammar

Algorithm 1: Deontic Preprocessing

Input: Terminal sequence $t_0 \dots t_n$, deontic rule with antecedent type a and consequent type b , necessary function \square , permitted function \diamond

Output: FULFILLED[row, col], ENABLER[row, col]

// Base case: row 0

for $col \leftarrow 0$ **to** $n - 1$ **do**

if $type(t_{col}) \neq a$ **then**

 FULFILLED[0, col] $\leftarrow +\infty$;

 ENABLER[0, col] $\leftarrow false$;

else

$J \leftarrow \{ j \mid [t_j \dots t_n]_b \in \square(t_{col}) \}$;

$J \leftarrow J \cup \{-\infty\}$;

 FULFILLED[0, col] $\leftarrow \max(J)$;

 ENABLER[0, col] $\leftarrow t_n \in \diamond t_{col}$;

// Inductive step:

for $row \leftarrow 1$ **to** $n - 1$ **do**

$L^{FF} \leftarrow$ FULFILLED[$row-1, :-1$];

$R^{FF} \leftarrow$ FULFILLED[$row-1, 1$];

 FULFILLED[row] $\leftarrow \min(L^{FF}, R^{FF})$;

$L^{EB} \leftarrow$ ENABLER[$row-1, :-1$];

$R^{EB} \leftarrow$ ENABLER[$row-1, 1$];

 ENABLER[row] $\leftarrow L^{EB} \vee R^{EB}$;

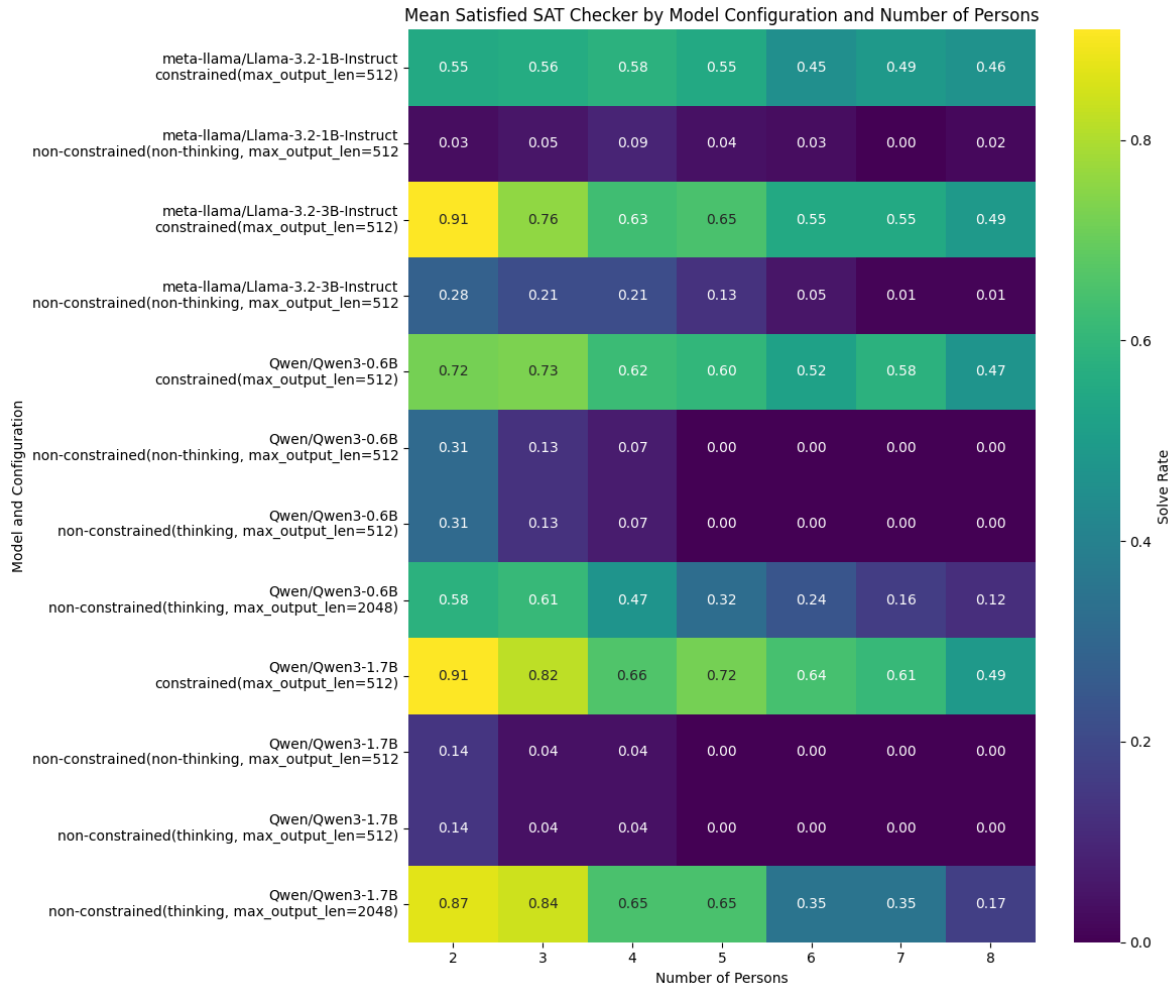


Figure 6: Knight and Knave test: Constrained models outperform unconstrained models across all puzzle sizes.

Algorithm 2: Commit \uparrow : Computing LF_n^{ALIVE} (Bottom-Up)

Input : ENABLER, FULFILLER, LF_{n-1}^{ALIVE} , FIRST, t_n

Output : LF_n^{ALIVE}

// Initialization

$LF_n^{\text{ALIVE}} \leftarrow \text{false}$ everywhere

$\text{PARENT}^{\text{FREE}}, \text{PARENT}^{\text{CONSTR}} \leftarrow \text{false}$ everywhere

$\text{ANCESTOR}^{\text{FREE}}, \text{ANCESTOR}^{\text{CONSTR}} \leftarrow \text{false}$ everywhere

// Base Case: row = 0

foreach terminal rule ($X \rightarrow t_n$) **do**

└ $\text{PARENT}^{\text{FREE}}[0, X] \leftarrow \text{true}$

foreach $V \in \mathcal{N}$ **do**

└ $\text{ANCESTOR}^{\text{FREE}}[0, V] \leftarrow \text{PARENT}^{\text{FREE}}[0, V]$

AddAncestor($\text{ANCESTOR}^{\text{FREE}}[0]$, FIRST)

// Inductive Step

for $row_X \leftarrow 1$ **to** n **do**

┌ **for** $row_B \leftarrow 0$ **to** $row_X - 1$ **do**

└ $tier_A \leftarrow n - row_B - 1$

└ $row_A \leftarrow row_X - row_B - 1$

└ **foreach** rule ($X \rightarrow A B$) **do**

└ **if** $\neg LF_{n-1}^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B]$ **then**

└ **continue**

└ Merge $_{\text{commit}}^{\uparrow}(row_X, row_B, tier_A, row_A, X \rightarrow A B)$

└ AddAncestor($\text{ANCESTOR}^{\text{FREE}}[row_X]$, FIRST)

└ AddAncestor($\text{ANCESTOR}^{\text{CONSTR}}[row_X]$, FIRST)

// Finalization: Create new left feet at tier n

for $row \leftarrow 0$ **to** n **do**

┌ **foreach** rule ($X \rightarrow A B$) **do**

└ $LF_n^{\text{ALIVE}}[n, row, X \rightarrow A B] \leftarrow \text{PARENT}^{\text{FREE}}[row, A] \vee \text{PARENT}^{\text{CONSTR}}[row, A]$

return LF_n^{ALIVE}

Algorithm 3: Merge[↑]_{commit} (Context-Free)

Input: $row_X, row_B, tier_A, row_A, (X \rightarrow A B)$ where $\text{type}(X \rightarrow A B) = \text{CF}$

```
if PARENTFREE[rowB, B] then
  PARENTFREE[rowX, X] ← true
  LFALIVEn[tierA, rowA, X → A B] ← true
if PARENTCONSTR[rowB, B] then
  PARENTCONSTR[rowX, X] ← true
  LFALIVEn[tierA, rowA, X → A B] ← true
if ANCESTORFREE[rowB, B] then
  ANCESTORFREE[rowX, X] ← true
  LFALIVEn[tierA, rowA, X → A B] ← true
if ANCESTORCONSTR[rowB, B] then
  ANCESTORCONSTR[rowX, X] ← true
  LFALIVEn[tierA, rowA, X → A B] ← true
```

Algorithm 4: Merge[↑]_{commit} (Strict Permitted)

Input: $row_X, row_B, tier_A, row_A, (X \rightarrow A B)$ where $\text{type}(X \rightarrow A B) = \text{STRICT}$

```
enabled ← ENABLER[tierA, rowA, X → A B]
fulfilled ← FULFILLER[tierA, rowA, X → A B] ≥ (n - rowB)
// All paths must satisfy ◇ at this scope
if enabled ∧ (ANCESTORFREE[rowB, B] ∨ ANCESTORCONSTR[rowB, B]) then
  ANCESTORCONSTR[rowX, X] ← true
  LFALIVEn[tierA, rowA, X → A B] ← true
if enabled ∧ fulfilled ∧ (PARENTFREE[rowB, B] ∨ PARENTCONSTR[rowB, B]) then
  PARENTCONSTR[rowX, X] ← true
  LFALIVEn[tierA, rowA, X → A B] ← true
```

Algorithm 5: Merge $_{commit}^{\uparrow}$ (Union Permitted)

Input: $row_X, row_B, tier_A, row_A, (X \rightarrow A B)$ where $\text{type}(X \rightarrow A B) = \text{UNION}$
 $enabled \leftarrow \text{ENABLER}[tier_A, row_A, X \rightarrow A B]$
 $fulfilled \leftarrow \text{FULFILLER}[tier_A, row_A, X \rightarrow A B] \geq (n - row_B)$
// Already permitted: propagate without re-checking
if $\text{ANCESTOR}^{\text{CONSTR}}[row_B, B]$ **then**
 $\text{ANCESTOR}^{\text{CONSTR}}[row_X, X] \leftarrow \text{true}$
 $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true}$
if $\text{PARENT}^{\text{CONSTR}}[row_B, B]$ **then**
 $\text{PARENT}^{\text{CONSTR}}[row_X, X] \leftarrow \text{true}$
 $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true}$
// Free t_n entering deontic scope: check \diamond
if $enabled \wedge \text{ANCESTOR}^{\text{FREE}}[row_B, B]$ **then**
 $\text{ANCESTOR}^{\text{CONSTR}}[row_X, X] \leftarrow \text{true}$
 $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true}$
if $enabled \wedge fulfilled \wedge \text{PARENT}^{\text{FREE}}[row_B, B]$ **then**
 $\text{PARENT}^{\text{CONSTR}}[row_X, X] \leftarrow \text{true}$
 $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true}$

Algorithm 6: Commit $_{\downarrow}$: Computing $\text{LF}_n^{\text{ALIVE}}$ (Top-Down)

Input: ENABLER, FULFILLER, $\text{LF}_{n-1}^{\text{ALIVE}}$, FIRST, t_n
Output: $\text{LF}_n^{\text{ALIVE}}$
// Initialization
 $\text{LF}_n^{\text{ALIVE}} \leftarrow \text{false}$ everywhere;
 $\text{CHILD}^{\text{FREE}}, \text{CHILD}^{\text{CONSTR}} \leftarrow \text{false}$ everywhere;
 $\text{ANCESTOR}^{\text{FREE}}, \text{ANCESTOR}^{\text{CONSTR}} \leftarrow \text{false}$ everywhere;
// Base case: Start symbol at row n
foreach $V \in \text{Nonterminals}$ **do**
 $\text{ANCESTOR}^{\text{FREE}}[n, V] \leftarrow (V = S)$;
ADDDDESCENDANT($\text{ANCESTOR}^{\text{FREE}}[n]$, FIRST);
// Inductive step: top-down
for $row_X \leftarrow n$ **to** 1 **do**
 for $row_B \leftarrow 0$ **to** $row_X - 1$ **do**
 $tier_A \leftarrow n - row_B - 1$;
 $row_A \leftarrow row_X - row_B - 1$;
 foreach rule $(X \rightarrow A B)$ **do**
 if $\neg \text{LF}_{n-1}^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B]$ **then continue**;
 MERGE $\downarrow (row_X, row_B, tier_A, row_A, X \rightarrow A B)$;
 ADDDDESCENDANT($\text{ANCESTOR}^{\text{FREE}}[row_X - 1]$, FIRST);
 ADDDDESCENDANT($\text{ANCESTOR}^{\text{CONSTR}}[row_X - 1]$, FIRST);
return $\text{LF}_n^{\text{ALIVE}}$;

Algorithm 7: Merge \downarrow (Context-Free)

Input: $row_X, row_B, tier_A, row_A, (X \rightarrow A B)$ where $\text{type}(X \rightarrow A B) = \text{CF}$

if $\text{CHILD}^{\text{FREE}}[row_X, X]$ **then**

- ┌ $\text{CHILD}^{\text{FREE}}[row_B, B] \leftarrow \text{true};$
- └ $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true};$

if $\text{CHILD}^{\text{CONSTR}}[row_X, X]$ **then**

- ┌ $\text{CHILD}^{\text{CONSTR}}[row_B, B] \leftarrow \text{true};$
- └ $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true};$

if $\text{ANCESTOR}^{\text{FREE}}[row_X, X]$ **then**

- ┌ $\text{ANCESTOR}^{\text{FREE}}[row_B, B] \leftarrow \text{true};$
- └ $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true};$

if $\text{ANCESTOR}^{\text{CONSTR}}[row_X, X]$ **then**

- ┌ $\text{ANCESTOR}^{\text{CONSTR}}[row_B, B] \leftarrow \text{true};$
- └ $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true};$

Algorithm 8: Merge \downarrow (Strict Permitted)

Input: $row_X, row_B, tier_A, row_A, (X \rightarrow A B)$ where $\text{type}(X \rightarrow A B) = \text{STRICT}$

$enabled \leftarrow \text{ENABLER}[tier_A, row_A, X \rightarrow A B];$

$fulfilled \leftarrow \text{FULFILLER}[tier_A, row_A, X \rightarrow A B] \geq (n - row_B);$

// All paths must satisfy \diamond at this scope

if $enabled \wedge (\text{ANCESTOR}^{\text{FREE}}[row_X, X] \vee \text{ANCESTOR}^{\text{CONSTR}}[row_X, X])$ **then**

- ┌ $\text{ANCESTOR}^{\text{CONSTR}}[row_B, B] \leftarrow \text{true};$
- └ $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true};$

if $enabled \wedge fulfilled \wedge (\text{CHILD}^{\text{FREE}}[row_X, X] \vee \text{CHILD}^{\text{CONSTR}}[row_X, X])$ **then**

- ┌ $\text{CHILD}^{\text{CONSTR}}[row_B, B] \leftarrow \text{true};$
- └ $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true};$

Algorithm 9: Merge_↓ (Union Permitted)

Input: $row_X, row_B, tier_A, row_A, (X \rightarrow A B)$ where $\text{type}(X \rightarrow A B) = \text{UNION}$
 $enabled \leftarrow \text{ENABLER}[tier_A, row_A, X \rightarrow A B];$
 $fulfilled \leftarrow \text{FULFILLER}[tier_A, row_A, X \rightarrow A B] \geq (n - row_B);$
// Already constrained: propagate (permission granted below)
if $\text{ANCESTOR}^{\text{CONSTR}}[row_X, X]$ **then**
 $\text{ANCESTOR}^{\text{CONSTR}}[row_B, B] \leftarrow \text{true};$
 $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true};$
if $\text{CHILD}^{\text{CONSTR}}[row_X, X]$ **then**
 $\text{CHILD}^{\text{CONSTR}}[row_B, B] \leftarrow \text{true};$
 $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true};$
// Free t_n entering deontic scope: check \diamond , if satisfied \rightarrow constrained
if $enabled \wedge \text{ANCESTOR}^{\text{FREE}}[row_X, X]$ **then**
 $\text{ANCESTOR}^{\text{CONSTR}}[row_B, B] \leftarrow \text{true};$
 $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true};$
if $enabled \wedge fulfilled \wedge \text{CHILD}^{\text{FREE}}[row_X, X]$ **then**
 $\text{CHILD}^{\text{CONSTR}}[row_B, B] \leftarrow \text{true};$
 $\text{LF}_n^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true};$

Algorithm 10: Probe_↑

Input: b – hypothetical terminal type at position n ; LF^{ALIVE} ; FIRST
Output: $\text{LF}^{\text{REACHABLE}\uparrow}$
 $\text{ANCESTOR}^{\uparrow}[\cdot, \cdot] \leftarrow \text{false};$
 $\text{LF}^{\text{REACHABLE}\uparrow}[\cdot, \cdot, \cdot] \leftarrow \text{false};$
foreach terminal rule $(V \rightarrow b)$ **do**
 $\text{ANCESTOR}^{\uparrow}[0, V] \leftarrow \text{true};$
 $\text{ADDANCESTOR}(\text{ANCESTOR}^{\uparrow}[0], \text{FIRST});$
for $row_X \leftarrow 1$ **to** n **do**
 for $row_B \leftarrow 0$ **to** $row_X - 1$ **do**
 $tier_A \leftarrow n - row_B - 1;$
 $row_A \leftarrow row_X - row_B - 1;$
 foreach rule $(X \rightarrow A B)$ **do**
 if $\text{LF}^{\text{ALIVE}}[tier_A, row_A, X \rightarrow A B] \wedge \text{ANCESTOR}^{\uparrow}[row_B, B]$ **then**
 $\text{ANCESTOR}^{\uparrow}[row_X, X] \leftarrow \text{true};$
 $\text{LF}^{\text{REACHABLE}\uparrow}[tier_A, row_A, X \rightarrow A B] \leftarrow \text{true};$
 $\text{ADDANCESTOR}(\text{ANCESTOR}^{\uparrow}[row_X], \text{FIRST});$
return $\text{LF}^{\text{REACHABLE}\uparrow}$

Algorithm 11: Probe \downarrow

Input: $\mathbf{LF}^{\text{ALIVE}}$, FIRST
Output: $\mathbf{LF}^{\text{REACHABLE}\downarrow}$
ANCESTOR $\downarrow[\cdot, \cdot] \leftarrow \text{false}$;
 $\mathbf{LF}^{\text{REACHABLE}\downarrow}[\cdot, \cdot, \cdot] \leftarrow \text{false}$;
ANCESTOR $\text{uncon}[\cdot, \cdot, \cdot] \leftarrow \text{false}$; // indexed by $[b, \text{row}, V]$
ANCESTOR $\downarrow[n, \text{START}] \leftarrow \text{true}$;
ADDDESCENDANT(ANCESTOR $\downarrow[n]$, FIRST);
foreach *deontic consequent terminal* b **do**
 ANCESTOR $\text{uncon}[b, n, \text{START}] \leftarrow \text{true}$;
 ADDDESCENDANT(ANCESTOR $\text{uncon}[b, n]$, FIRST);
for $\text{row}_B \leftarrow n - 1$ **to** 0 **do**
 for $\text{row}_X \leftarrow \text{row}_B + 1$ **to** n **do**
 $\text{tier}_A \leftarrow n - \text{row}_B - 1$;
 $\text{row}_A \leftarrow \text{row}_X - \text{row}_B - 1$;
 foreach *rule* $(X \rightarrow A B)$ **do**
 if $\mathbf{LF}^{\text{ALIVE}}[\text{tier}_A, \text{row}_A, X \rightarrow A B] \wedge \text{ANCESTOR}\downarrow[\text{row}_X, X]$ **then**
 ANCESTOR $\downarrow[\text{row}_B, B] \leftarrow \text{true}$;
 $\mathbf{LF}^{\text{REACHABLE}\downarrow}[\text{tier}_A, \text{row}_A, X \rightarrow A B] \leftarrow \text{true}$;
 foreach *deontic consequent terminal* b **do**
 if $(X \rightarrow A B)$ *is not anchor of any deontic rule with consequent* b **then**
 if ANCESTOR $\text{uncon}[b, \text{row}_X, X]$ **then**
 ANCESTOR $\text{uncon}[b, \text{row}_B, B] \leftarrow \text{true}$;
 ADDDESCENDANT(ANCESTOR $\downarrow[\text{row}_B]$, FIRST);
 foreach *deontic consequent terminal* b **do**
 ADDDESCENDANT(ANCESTOR $\text{uncon}[b, \text{row}_B]$, FIRST);
return $\mathbf{LF}^{\text{REACHABLE}\downarrow}$

Algorithm 12: AddAncestor / AddDescendant

Procedure ADDANCESTOR(*cell*, FIRST)
 foreach $V \in \text{Nonterminals}$ **do**
 if $\exists W : \text{cell}[W] \wedge \text{FIRST}[V, W]$ **then**
 $\text{cell}[V] \leftarrow \text{true}$;

Procedure ADDDESCENDANT(*cell*, FIRST)
 foreach $W \in \text{Nonterminals}$ **do**
 if $\exists V : \text{cell}[V] \wedge \text{FIRST}[V, W]$ **then**
 $\text{cell}[W] \leftarrow \text{true}$;

Algorithm 13: Collect Permitted Constraints

Input : $LF^{RRA}[\text{tier}, \text{row}, \text{rule}] : \text{bool}$

Input : Terminal sequence $t_0 \dots t_{n-1}$

Input : Deontic rules D

Output : ψ_b^\diamond : permitted specification sent to enforcer

$LF^{RRA}[\text{tier}, \text{row}]$ $\text{Reachable}^\uparrow \wedge \text{Reachable}^\downarrow \wedge \text{Alive}$ for each left foot

$\text{memo}[\text{row}, \text{col}]$ true if and only if underneath a RRA left foot

ψ_b^\diamond permitted specification for next terminal of type b

for each deontic rule with anchor $a \gg A > X < B \ll b$ and permitted deontic function \diamond **do**

```
/* Initialize Memo with RRA */
memo[row, col] ←  $LF^{RRA}[\text{row} + \text{col}, \text{row}, X \rightarrow A B]$ 

/* Propagate Reachability */
for row ←  $(n - 1)$  to 0 do
  memo[row, :] |= memo[row + 1, :] // vertical (↓)
  memo[row, 1:] |= memo[row + 1, :-1] // diagonal (↘)

/* Collect Active Antecedents */
for  $i \leftarrow 0$  to  $n - 1$  do
  if  $\text{type}(t_i) = a$  and memo[0,  $i$ ] then
     $\psi_b^\diamond \leftarrow \psi_b^\diamond \cup \diamond(t_i)$ 

/*  $\psi_b^\diamond$  is sent to enforcer */
```

```
start → bad_guess* good_guess
good_guess → sufficient_assignment_seq semi_3 GOOD_GUESS_TAIL
sufficient_assignment_seq → assignment_seq
bad_guess → contradictory_assignment_seq semi_3 CONTRADICTORY_GUESS_TAIL
           | repetition_assignment_seq semi_3 REPETITION_GUESS_TAIL
contradictory_assignment_seq → assignment_seq
repetition_assignment_seq → assignment_seq_semi_1 repeating_assignment
assignment_seq_semi_1 → assignment_seq semi_1
assignment_seq → assignment_seq semi_1 assignment
                 | GUESS_HEAD assignment
                 | unit_prop
assignment → ASSIGNMENT
ASSIGNMENT → (NAME "->" KK)
           | AUX_VAR
repeating_assignment → REPEATING_ASSIGNMENT
REPEATING_ASSIGNMENT2 → ASSIGNMENT
semi_3 → SEMI_3 "\n"
semi_2 → SEMI_2 "\n"
semi_1 → SEMI_1 "\n"
SEMI_33 → ";"
SEMI_22 → ";"
SEMI_11 → ";"
GUESS_HEAD → "Here is a different and improved guess: {\n"
GOOD_GUESS_TAIL → "}} «DONE»"
CONTRADICTORY_GUESS_TAIL → "}} contradiction\n\n"
REPETITION_GUESS_TAIL → "}} Bad guess\n\n"
NAME → /(Jacob|Noah|Michael|Liam|Ella|... |Ava)/
KK → "knight" | "knave"
```

Figure 7: Context free grammar for LLM generations in the Knight and Knave problem.