



MONASH University

Formal Explainability for Artificial Intelligence in Dynamic Environments

Jime Cuartas Granada

Supervisors:

Alexey Ignatiev

Peter J. Stuckey

Julian Gutierrez

A Progress Review report at
Monash University in 2026
School of Information Technology

Abstract

In dynamic environments, the goal of Artificial Intelligence (AI) is to build intelligent agents capable of addressing sequential decision-making settings. Reinforcement Learning (RL) is a branch of Machine Learning that addresses sequential decision-making by agents to perform tasks. In this context, there are two important challenges for humans to understand decisions made by agents: (1) the sequential decisions are connected, and (2) the agents may use opaque black-box models (e.g., neural networks) for each decision.

Despite the success of RL in sequential decision-making, the lack of transparency in understanding their decisions can make the agents hard to validate. To address the need for transparency, there are efforts to develop Explainable Artificial Intelligence (XAI) and its subfield, Explainable Reinforcement Learning. XAI is a set of methods designed to make AI models easier to comprehend. Despite the importance of Explainable Reinforcement Learning in developing trustworthy intelligent agents, there are gaps in current research to make sequential decision-making explainable.

This project proposes to explain sequential decision-making using formal reasoning. To achieve this goal, the proposal focuses on (1) Formal Explainability for Finite Automata, to address sequential actions in deterministic environments, and (2) Formal Explainability for Reinforcement Learning, where the agent's behaviour is non-interpretable.

Contents

Abstract	i
List of Figures	iii
List of Tables	iv
1 Introduction	1
1.1 Refined scope - problem statement	1
1.2 Contributions to knowledge - achieved and projected	2
2 Progress	3
2.1 Model Proposal: Context-Free Grammar (CFG) Explanations	3
2.2 Motivation: The Gap between Detection and Explanation	4
2.3 Preliminaries	5
2.4 Future Work and Timeline to Completion	12
A An Appendix	13
Bibliography	15

List of Figures

2.1	A Pushdown Automaton \mathcal{A} accepting the language of balanced parentheses. The symbol $\$$ is used as the bottom-of-stack marker.	6
-----	---	---

List of Tables

Chapter 1

Introduction

The deployment of Artificial Intelligence (AI) algorithms has necessitated the need for eXplainability AI (XAI) methods in order to ensure transparency, trust, and accountability. While much of the field has focused on heuristic explanations for opaque models, there is an interest in formal approaches that provide rigorous guarantees about the explanations generated [1, 2].

A fundamental challenge in dynamic environments is explaining sequential decision-making. To address this, we model these processes using Automata, which provide a symbolic and tractable representation of sequential decision functions. This approach allows us to generate formal explanations, why a specific sequence of actions leads to a particular outcome. Automata are widely used in software verification [3], design of communication protocols [4], and syntax parsing in compiler [5]. When a computational model, such as a Finite Automaton (FA) or a Pushdown Automaton (PDA), accepts or rejects an input string, the reasoning behind that decision can be non-trivial. Understanding why a specific input was accepted or rejected is crucial for debugging, and refinement purposes.

This research project investigates the formalization of explanations for sequential decision-making. Having addressed an approach to deliver formal explanations for Finite Automata (FA) in the first stage of this research, and submitting it to a ICALP 2026. I now move to address explanations for Context-Free Languages (CFG) / Pushdown Automata (PDA) decisions.

1.1 Refined scope - problem statement

While standard XAI focuses on feature attribution in classifiers, the "features" in formal languages are sequential and structural. Since the confirmation report, the research

scope has been refined to address three primary gaps:

- **Research Problem 1 (Completed): Explaining Finite Automata.** Finite Automata are often assumed to be interpretable. However, large FA are cognitively inaccessible to humans. I have developed a framework to compute formal explanations for the acceptance and rejection of inputs in FA, providing a rigorous foundation for automaton-based explainability.
- **Research Problem 2: Explaining Pushdown Automata (PDA).** Context-free languages, recognized by PDAs, introduce a stack-based memory that allows to represent more complex behaviors. My objective is the generation of Minimal Contrastive Explanations (CXPs) the minimal sets of modifications required to turn a rejected word into an accepted one; and Minimal Abductive Explanations (AXps) the minimal sets of tokens such that the word is going to be rejected. Finally, there is a lack of quantitative metrics that assign a Feature Attribution Score (“degree of responsibility”) to specific indices in a rejected string.
- **Research Problem 3: Explaining Decisions of Sthocastic models.** How can the Feature Attribution Score be extended to explain outcomes for policies in sthocastic environments? This problem explores sequential decision-making, where actions influence future states under uncertain contexts.

1.2 Contributions to knowledge - achieved and projected

This research provides both theoretical and practical contributions to the field of Computer Science:

Achieved Contributions:

- Development of an theoretical and practical approach to explain Finite Automata decisions.
- A paper submitted to ICALP 2026 titled “A Formal Framework for the Explanation of Finite Automata Decisions”

Projected Contributions:

- Explaining Pushdown Automata decisions: (In Progress) Extending the formal explanation framework to PDAs, which recognize context-free languages. This involves developing algorithms to identify the minimal contrastive explanations (CXPs) and Abductive Explanations (AXPs), and quantifying the contribution of specific tokens to the decision (acceptance or rejection).
- Explaining Decisions of Sthocastic models: Extending the formal explanation framework to Markov Decision Processes (MDPs). The goal is to provide verifiable explanations able to identify the environmental factors or decision points that lead to a particular outcomes.

Chapter 2

Progress

2.1 Model Proposal: Context-Free Grammar (CFG) Explanations

The research has evolved from the study of Finite Automata (FA) to more expressive computational models. While FA provided a baseline for explaining sequential behaviors, many complex problems require the model to “remember” an arbitrary number of previous inputs in the sequence to determine the validity of subsequent inputs

consider the challenge of identifying and explaining syntax errors in programming languages. Code structures are inherently hierarchical; blocks of logic are nested within one another, defined by matching delimiters such as `{ }` in C or `()` in Lisp.

Consider the abstract language $L = \{a^n b^n \mid n \geq 1\}$, which represents a sequence where every ‘a’ must be matched by a corresponding ‘b’.

- **The Limitation:** A standard Finite Automaton (FA) possesses no auxiliary memory. Therefore, it cannot count the number of a ’s to ensure they match the number of b ’s once n exceeds the number of states in the machine.
- **The Explanation Failure:** If an FA were used to validate such a sequence, it would process inputs locally. Upon encountering a mismatch (e.g., $a^4 b^2$), it might reject the sequence, but it lacks the structural context to generate a contrastive explanation such as: “The sequence is invalid because the third or fourth ‘a’ was not closed by a matching ‘b’.” Instead, it can only report a “transition failure” at the specific index, obscuring the root cause of the error.
- **From FA to PDA:** We propose the use of Pushdown Automata (PDAs) to model decision-making processes with memory.

Unlike FA, the addition of a stack allows the description of more complex languages. Here the challenge is how to explain PDA decisions.

2.2 Motivation: The Gap between Detection and Explanation

One well established applications of Context-Free Grammars is in the design of programming languages and compilers. Compilers are highly efficient at detecting when a sequence of tokens fails to belong to a grammar. However, a fundamental question is: are they able of generating a useful *explanation* for why the failure occurred or how to fix it?

Consider the following C code, where a typo has introduced a double opening bracket `{{` in the `for` loop:

```

1 | int main() {
2 |     for(int i=0; i<10; i++){ { // <--- Error: Double bracket
3 |         printf("hello");
4 |     }
5 | }
```

When compiled (e.g., using GCC or Clang), the parser consumes the input until it reaches `<<EOF>>`(end-of-file), finding only then that EOF was not expected, instead there is an incomplete structure, the token `}` is expected before `<<EOF>>`. The resulting error message is:

```

error: expected '}' at end of input
    5 | }
      | ^
```

The Explanation: While the compiler’s output is correct, the file ended while the stack still contained an open brace that was not closed yet. It is *misleading*.

- **Root Cause:** The compiler points to line 5 (the end of the file) as the location of the error. However, the root cause is located at line 2.
- **Lack of Contrastive Reasoning:** A human (or a formal explainer) would identify that the input is “almost correct”. The explanation should not just report a missing symbol at the end, but rather propose a *minimal correction*.

In this case, the minimal correction is not to add a brace at the end, but to remove the redundant opening brace at the loop initialization:

```

1 | int main() {
2 |     for(int i=0; i<10; i++){
3 |         printf("hello");
4 |     }
5 | }
```

This discrepancy motivates the need for our proposed formal framework. We aim to move beyond isolated decisions to explained decisions revealing these minimal set of edits (Contrastive Explanations).

2.3 Preliminaries

Pushdown Automata and Context-Free Grammars

To provide a foundation for the proposed explanation extraction methods, standard definitions and notations for Pushdown Automata and Context-Free Grammar are adopted [6, 7].

Definition 2.1 (Pushdown Automaton). A Pushdown Automaton (PDA) extends the capabilities of a Finite Automaton by incorporating an infinite memory stack. A PDA is formally defined as a 7-tuple $\mathcal{A} = (Q, \Sigma, V, \delta, q^0, v^0, F)$, where:

- Q is a finite set of states.
- Σ is the input alphabet, a finite terminal alphabet.
- V is the stack alphabet, a finite nonterminal alphabet that can be pushed onto or popped from the stack.
- $\delta : Q \times (\Sigma \cup \{\epsilon\}) \times V \rightarrow \mathcal{P}(Q \times V^*)$ ¹ is the transition function. It dictates how the machine transitions between states and modifies the stack based on the current state, input symbol, and the top symbol of the stack.
- $q^0 \in Q$ is the initial state.
- $v^0 \in V$ is the initial pushdownstore symbol.
- $F \subseteq Q$ is the set of accepting states.

A configuration of \mathcal{A} is a triple $c = (q, \gamma, x)$ in $Q \times V^* \times \Sigma^*$. And the automaton moves from c into configuration $c' = (q', \gamma', x')$, denoted as $c \vdash c'$ if:

- $\gamma = v\gamma_1 (v \in V)$ ², $x = ax' (a \in \Sigma)$, $\gamma' = m\gamma_1 (m \in V^*)$ and $(q', m) \in \delta(q, a, v)$, namely “a-move”;
- or $\gamma = v\gamma_1 (v \in V)$, $x = x'$, $\gamma' = m\gamma_1 (m \in V^*)$ and $(q', m) \in \delta(q, \epsilon, v)$, namely “ ϵ -move”.

A *word* is accepted by a pushdown automaton if, starting with an empty stack, there is a path through the automaton such that the automaton stops in an accepting state after the entire string has been read. The *language* recognized by a PDA \mathcal{A} is the set of all accepted words, denoted as $L(\mathcal{A})$.

¹ \mathcal{P} denotes the power set (the set of all its subsets). V^* and Σ^* denote the Kleene closures of the stack and input alphabets, respectively, representing the sets of all finite strings formed by those alphabets.

² $\gamma_1 \in V^*$ represents the remaining symbols on the stack below the top symbol v . Thus, γ represents the full current stack, formed by the top symbol v (to be processed) and the rest of the stack γ_1 .

The following PDA recognizes the language of balanced parentheses (a subset of the Dyck language)³ and is used throughout the document to illustrate the proposed ideas.

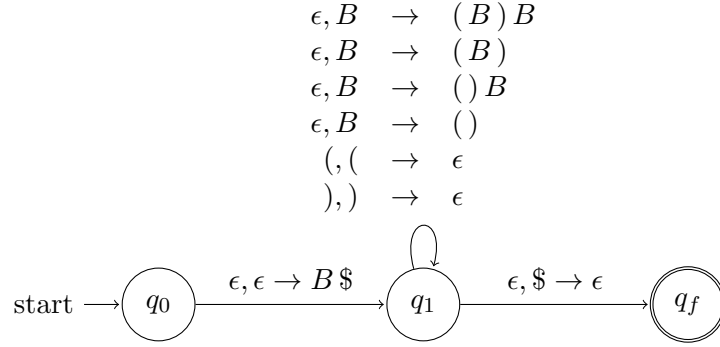


FIGURE 2.1: A Pushdown Automaton \mathcal{A} accepting the language of balanced parentheses. The symbol $\$$ is used as the bottom-of-stack marker.

Example 2.1. Let \mathcal{A} be the PDA that accepts the language generated by G in [Example 2.2](#). The PDA uses a stack to track the depth of nesting, pushing a symbol for every open parenthesis and popping for every closed one. B represents *Balanced* in G .

- *States:* $Q = \{q_0, q_1, q_f\}$, $q^0 = q_0$, $F = \{q_f\}$
- *Input Alphabet:* $\Sigma = \{ (,) \}$
- *Stack Alphabet:* $V = \{ (,), B, \$ \}$, $v^0 = B$
- *Transitions:*
 1. $\delta(q_0, \epsilon, \epsilon) = \{(q_1, B\$)\}$ (Initialize stack with Start Symbol)
 2. $\delta(q_1, \epsilon, B) = \{(q_1, (B)B), (q_1, (B)), (q_1, ()B), (q_1, ())\}$ (Expand B)
 3. $\delta(q_1, (, () = \{(q_1, \epsilon)\}$ (Match input '(' with stack '(')
 4. $\delta(q_1,),)) = \{(q_1, \epsilon)\}$ (Match input ')' with stack ')')
 5. $\delta(q_1, \epsilon, \$) = \{(q_f, \epsilon)\}$ (Accept if bottom marker is reached)

[Figure 2.1](#) illustrates this PDA. The configuration history of \mathcal{A} for the input $"()()"$ is:

$$\begin{array}{ll}
 (q_0, ()(), \epsilon) & \\
 \vdash (q_1, ()(), B\$) & \text{(Initialize)} \\
 \vdash (q_1, ()(), ()B\$) & \text{(Expand } B \rightarrow ()B) \\
 \vdash (q_1,)(),)B\$) & \text{(Match '(')} \\
 \vdash (q_1, (), B\$) & \text{(Match ')')} \\
 \vdash (q_1, (), ()\$) & \text{(Expand } B \rightarrow ()) \\
 \vdash (q_1,),)\$) & \text{(Match '(')} \\
 \vdash (q_1, \epsilon, \$) & \text{(Match ')')} \\
 \vdash (q_f, \epsilon, \epsilon) & \text{(Accept)}
 \end{array} \tag{2.1}$$

³The Dyck language describes a set of strings with balanced and properly nested brackets (e.g., $()$, $[]$, $\{\}$) [7]. The example focuses solely on non-empty sequences of balanced (and).

Definition 2.2 (Context-Free Grammar). A Context-Free Grammar (CFG) is defined as a 4-tuple $G = (V, \Sigma, R, S)$, where:

- V (Variables/Non-terminals) is a finite set of variables (non-terminal symbols).
- Σ (Terminals) is a finite set of terminal symbols, disjoint from V .
- R is a finite set of production rules of the form $A \rightarrow \alpha$, where $A \in V$ describes a variable and $\alpha \in (V \cup \Sigma)^*$ is a string of variables and terminals.
- $S \in V$ is the start variable.

A fundamental equivalence between CFGs and PDAs: a language L is context-free iff there exists a PDA \mathcal{A} such that $L(\mathcal{A}) = L$ [6, 7]. This equivalence allows us to use grammar-based parsing algorithms, such as CYK, to analyze the behavior of PDAs.

Example 2.2. Let $G = (\{Balanced\}, \{(),\}, R, Balanced)$ be the grammar defined by the following production rules R :

$$\begin{aligned}
 Balanced &\rightarrow (Balanced) Balanced && (Rule\ 1) \\
 Balanced &\rightarrow (Balanced) && (Rule\ 2) \\
 Balanced &\rightarrow () Balanced && (Rule\ 3) \\
 Balanced &\rightarrow () && (Rule\ 4)
 \end{aligned} \tag{2.2}$$

This grammar generates the language of properly nested parentheses. For instance, the string “ $()()$ ” can be derivated as follows.

$$\begin{aligned}
 Balanced &\Rightarrow () \mathbf{Balanced} && (Rule\ 3: \text{parentheses and } Balanced) \\
 &\Rightarrow () () && (Rule\ 4: \text{Reduce final 'Balanced' to '()'})
 \end{aligned} \tag{2.3}$$

Transition to Explainability for Pushdown Automata decisions Following the confirmation of candidacy, a natural progression from the study of Finite Automata is to consider computational models with greater expressive power. the research scope has expanded to include Pushdown Automata (PDA).

The CYK-Based Explanation Engine

To extract formal explanations from these Context-Free structures, we have developed a novel adaptation of the CYK (Cocke-Younger-Kasami) algorithm.

While standard CYK is primarily used for parsing, we utilize the resulting triangular parsing table as a search space for explanations.

By analyzing partial trees within the CYK table, for a rejected string w and a PDA \mathcal{A} we can identify:

- *Minimal Abductive Explanations*: minimal subsets of input tokens (leaves in the parse tree) such that is impossible for the start symbol S to be derivable (does not appear in the top cell)
- *Minimal Contrastive Explanations*: minimal set of edits (substitutions) required to “repair” the CYK table such that the start symbol S is derivable (appears in the top cell).

Minimal Contrastive Explanations: In the case of rejection, the Definition environment.

The Grammar \mathcal{A} PCFG is defined as a tuple $G = (V, \Sigma, R, S, P)$, where:

- V is a finite set of non-terminal symbols.
- Σ is a finite set of terminal symbols (the alphabet).
- R is a finite set of production rules.
- $S \in V$ is the start symbol.
- $P : R \rightarrow [0, 1]$ is a probability function such that for each $A \in V$, $\sum_{A \rightarrow \alpha \in R} P(A \rightarrow \alpha) = 1$.

The Rejected Word Let $w = \sigma_1\sigma_2\ldots\sigma_n$ be a string in Σ^* . We say w is rejected if $w \notin L(G)$, where $L(G)$ is the language generated by the grammar.

Definition: Contrastive Set. For a rejected word w of length n , a set of indices $I \subseteq \{1, \dots, n\}$ is a Contrastive Explanation if:

- Feasibility: There exists a word $w' \in L(G)$ such that w and w' differ only at indices $i \in I$. Formally, $\forall j \notin I, \sigma_j = \sigma'_j$.
- Minimality: No proper subset $I' \subset I$ satisfies the feasibility condition.

Example: For $w = ()))$, the index set $I = \{2\}$ is a contrastive explanation because changing index 2 to $)$ results in $w' = (()) \in L(G)$, and the empty set \emptyset is not feasible.

Introducing Probabilistic Preference. In a PCFG, not all accepted words w' are equal. We can rank explanations by the likelihood of the correction they enable.

The Scoring Function For any word $w' \in L(G)$, the score $P(w')$ is the maximum probability among all possible parse trees T that yield w' (Viterbi Algorithm):

$$P(w') = \max_{T \in \text{Trees}(w')} P(T)$$

Optimal Contrastive Explanation Given a rejected word w , an explanation I_1 is probabilistically superior to I_2 if the best correction enabled by I_1 is more likely than that of I_2 . We define the Explanation Weight as:

$$\text{Weight}(I) = \max\{P(w') \mid w' \text{ matches } w \text{ except at indices } I, w' \in L(G)\}$$

The Extended CYK Table

Standard CYK populates a 3D table $T[i, j, A]$, representing the maximum probability that non-terminal A derives the substring from index i to j .

For a word $w = \sigma_1\sigma_2\ldots\sigma_n$ and an index set I :

- **Base Case (Length 1)** For each position $i \in \{1, \dots, n\}$ and each non-terminal A :

- If $i \notin I$ (Fixed):

$$T[i, i, A] = P(A \rightarrow \sigma_i)$$

(If no such rule exists, the probability is 0).

- If $i \in I$ (in exp):

$$T[i, i, A] = \sum_{\sigma \in \Sigma} P(A \rightarrow \sigma)$$

This selects the most likely terminal that A can produce at that “broken” index.

- **Recursive Step (Length $l > 1$)** For each length l from 2 to n , each starting position i from 1 to $n - l + 1$, and each non-terminal A :

$$T[i, l, A] = \max_{A \rightarrow BC \in R} \left(\max_{1 \leq k < l} \{P(A \rightarrow BC) \cdot T[i, k, B] \cdot T[i + k, l - k, C]\} \right)$$

Heatmap for Contrastive Explanations

Theorem: Given a word $w = \sigma_1\sigma_1 \dots \sigma_n$, and grammar G , such that $w \notin L(G)$. Let \mathcal{E} be the set of all contrastive explanations for w . If H is a hitting set of all contrastive explanations \mathcal{E} , then no word w' that keeps the indices in H fixed to their original values in w can be accepted by the grammar. Mathematically:

If $\forall i \in H, \sigma'_i = \sigma_i$, then $w' \notin L(G)$.

Proof by Contradiction

Step 1: Assume the negation. Assume there exists a word $w' \in L(G)$ such that w' agrees with the original rejected word w on all indices in the hitting set H .

$$\forall i \in H : \sigma'_i = \sigma_i$$

Step 2: Let J be the set of indices where w' differs from w ($J \cap H = \emptyset$).

$$J = \{j \mid \sigma'_j \neq \sigma_j\}$$

Step 3: Relate to Contrastive Explanations.

Since $w' \in L(G)$ and it was formed by changing indices J in w , then by definition, J is a “feasible” contrastive set. It must either be a minimal contrastive explanation or a

superset of one.

$$\exists I \in \mathcal{E} \text{ such that } I \subseteq J$$

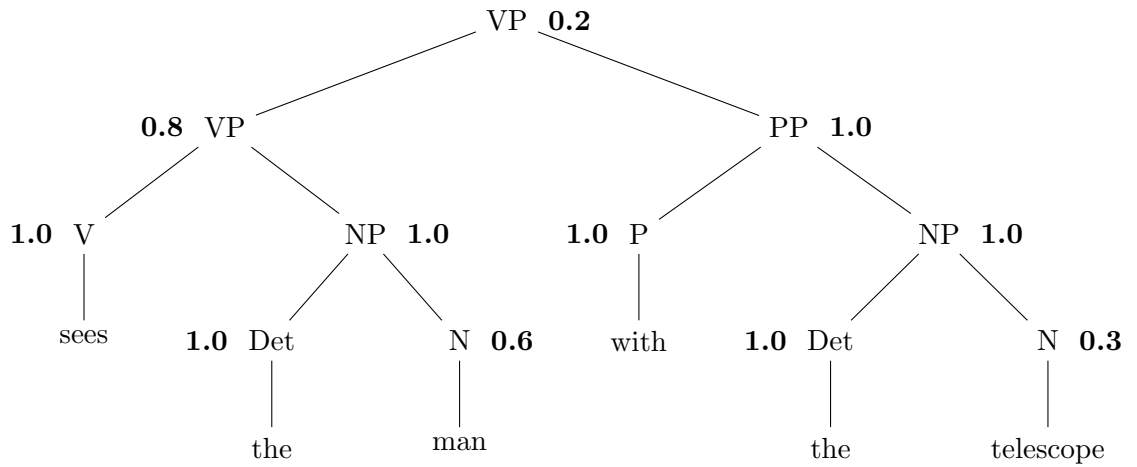
Step 4: The Contradiction.

We know from Step 2 that $J \cap H = \emptyset$. Since $I \subseteq J$, it follows that $I \cap H = \emptyset$. However, by definition, H is a hitting set of \mathcal{E} , meaning it must have a non-empty intersection with every $I \in \mathcal{E}$.

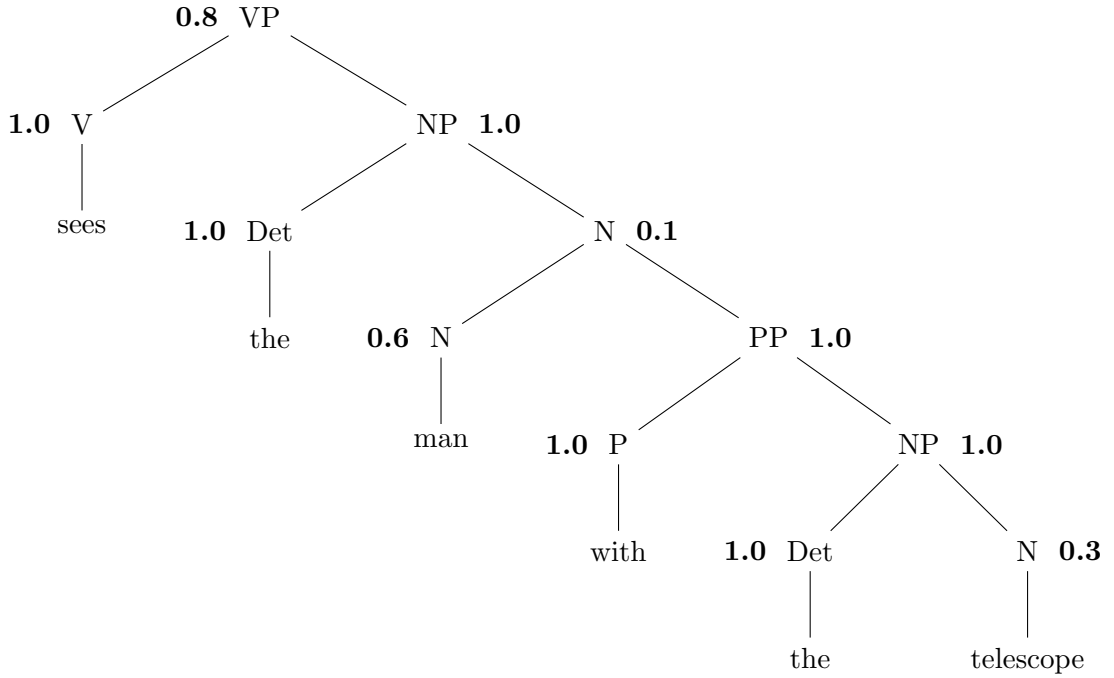
$$H \cap I \neq \emptyset \quad (\text{Contradiction})$$

Conclusion: Our assumption that an accepted word w' exists is false. Therefore, the indices in H effectively “block” all possible paths to acceptance. They are the necessary components of the error.

$$\begin{array}{llllllll} 0.8 & VP \rightarrow V \ NP & 1 & NP \rightarrow Det \ N & 0.1 & N \rightarrow N \ PP & 1 & Det \rightarrow \text{the} & 0.6 & N \rightarrow \text{man} \\ 0.2 & VP \rightarrow VP \ PP & 1 & PP \rightarrow P \ NP & 1 & V \rightarrow \text{sees} & 1 & P \rightarrow \text{with} & 0.3 & N \rightarrow \text{telescope} \end{array}$$



$$P(t1) = 0.6 * 0.8 * 0.2 * 0.3 = 0.0288$$



$$P(t2) = 0.6 * 0.8 * 0.1 * 0.3 = 0.0144$$

$$p(VP, \text{sees the man with the telescope}) = 0.0288 + 0.0144 = 0.0432$$

If “telescope” were “man”

$$P(t3) = 0.6 * 0.8 * 0.2 * 0.6 = 0.0576$$

$$P(t4) = 0.6 * 0.8 * 0.1 * 0.6 = 0.0288$$

$$\text{sees the man with the ?} = P(t1) + P(t2) + P(t3) + P(t4) = 0.1296$$

Given the grammar G and the rejected word $w = \sigma_1\sigma_2\ldots\sigma_n \notin L(G)$. Let $\mathcal{E} = \{I_1, I_2, \ldots, I_k\}$ the collection of all Minimal Contrastive Explanations

Each $I \in \mathcal{E}$: “The error happened exactly in the set of index I ”. and $P(I)$ Is the sum up of probabilities of all possible trees that derivates accepted words $w' = \sigma'_1\sigma'_2\ldots\sigma'_n$ where $\forall_{i \notin I} \sigma'_i = \sigma_i$

$P(A)$: Probability to get an accepted Tree with a minimal exp. Every tree is disjoint.

$$P(A) = P(I_1) + P(I_2) + \cdots + P(I_k)$$

$$P(I | A) = \frac{P(I)}{\sum_{J \in \mathcal{E}} P(J)}$$

$$P(\text{Error in } i | A) = \sum_{I \in \mathcal{E}} P(\text{Error in } i | I, A) P(I | A)$$

$$P(\text{Error in } i \mid I, A) = \mathbb{P}(i \in I) = \begin{cases} 1 & \text{if } i \in I \\ 0 & \text{if } i \notin I \end{cases}$$

2.4 Future Work and Timeline to Completion

Appendix A

An Appendix

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Vivamus at pulvinar nisi. Phasellus hendrerit, diam placerat interdum iaculis, mauris justo cursus risus, in viverra purus eros at ligula. Ut metus justo, consequat a tristique posuere, laoreet nec nibh. Etiam et scelerisque mauris. Phasellus vel massa magna. Ut non neque id tortor pharetra bibendum vitae sit amet nisi. Duis nec quam quam, sed euismod justo. Pellentesque eu tellus vitae ante tempus malesuada. Nunc accumsan, quam in congue consequat, lectus lectus dapibus erat, id aliquet urna neque at massa. Nulla facilisi. Morbi ullamcorper eleifend posuere. Donec libero leo, faucibus nec bibendum at, mattis et urna. Proin consectetur, nunc ut imperdiet lobortis, magna neque tincidunt lectus, id iaculis nisi justo id nibh. Pellentesque vel sem in erat vulputate faucibus molestie ut lorem.

Quisque tristique urna in lorem laoreet at laoreet quam congue. Donec dolor turpis, blandit non imperdiet aliquet, blandit et felis. In lorem nisi, pretium sit amet vestibulum sed, tempus et sem. Proin non ante turpis. Nulla imperdiet fringilla convallis. Vivamus vel bibendum nisl. Pellentesque justo lectus, molestie vel luctus sed, lobortis in libero. Nulla facilisi. Aliquam erat volutpat. Suspendisse vitae nunc nunc. Sed aliquet est suscipit sapien rhoncus non adipiscing nibh consequat. Aliquam metus urna, faucibus eu vulputate non, luctus eu justo.

Donec urna leo, vulputate vitae porta eu, vehicula blandit libero. Phasellus eget massa et leo condimentum mollis. Nullam molestie, justo at pellentesque vulputate, sapien velit ornare diam, nec gravida lacus augue non diam. Integer mattis lacus id libero ultrices sit amet mollis neque molestie. Integer ut leo eget mi volutpat congue. Vivamus sodales, turpis id venenatis placerat, tellus purus adipiscing magna, eu aliquam nibh dolor id nibh. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Sed cursus convallis quam nec vehicula. Sed vulputate neque eget odio fringilla ac sodales urna feugiat.

Phasellus nisi quam, volutpat non ullamcorper eget, congue fringilla leo. Cras et erat et nibh placerat commodo id ornare est. Nulla facilisi. Aenean pulvinar scelerisque eros eget interdum. Nunc pulvinar magna ut felis varius in hendrerit dolor accumsan. Nunc pellentesque magna quis magna bibendum non laoreet erat tincidunt. Nulla facilisi.

Duis eget massa sem, gravida interdum ipsum. Nulla nunc nisl, hendrerit sit amet commodo vel, varius id tellus. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Nunc ac dolor est. Suspendisse ultrices tincidunt metus eget accumsan. Nullam facilisis, justo vitae convallis sollicitudin, eros augue malesuada metus, nec sagittis diam nibh ut sapien. Duis blandit lectus vitae lorem aliquam nec euismod nisi volutpat. Vestibulum ornare dictum tortor, at faucibus justo tempor non. Nulla facilisi. Cras non massa nunc, eget euismod purus. Nunc metus ipsum, euismod a consectetur vel, hendrerit nec nunc.

Bibliography

- [1] João Marques-Silva. Logic-based explainability in machine learning. In Leopoldo E. Bertossi and Guohui Xiao, editors, *Reasoning Web. Causality, Explanations and Declarative Knowledge - 18th International Summer School 2022, Berlin, Germany, September 27-30, 2022, Tutorial Lectures*, volume 13759 of *Lecture Notes in Computer Science*, pages 24–104. Springer, 2022. doi: 10.1007/978-3-031-31414-8_2. URL https://doi.org/10.1007/978-3-031-31414-8_2.
- [2] Adnan Darwiche. Logic for explainable AI. In *LICS*, pages 1–11. IEEE, 2023.
- [3] Christel Baier and Joost-Pieter Katoen. *Principles of model checking*. MIT Press, 2008. ISBN 978-0-262-02649-9.
- [4] Gerard J. Holzmann. The model checker SPIN. *IEEE Trans. Software Eng.*, 23(5):279–295, 1997. doi: 10.1109/32.588521. URL <https://doi.org/10.1109/32.588521>.
- [5] Alfred V. Aho, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison-Wesley series in computer science / World student series edition. Addison-Wesley, 1986. ISBN 0-201-10088-6. URL <https://www.worldcat.org/oclc/12285707>.
- [6] John E. Hopcroft, Rajeev Motwani, and Jeffrey D. Ullman. *Introduction to automata theory, languages, and computation, 3rd Edition*. Pearson international edition. Addison-Wesley, 2007. ISBN 978-0-321-47617-3.
- [7] Jean Berstel and Luc Boasson. Context-free languages. In Jan van Leeuwen, editor, *Handbook of Theoretical Computer Science, Volume B: Formal Models and Semantics*, pages 59–102. Elsevier and MIT Press, 1990. doi: 10.1016/B978-0-444-88074-1.50007-X. URL <https://doi.org/10.1016/b978-0-444-88074-1.50007-x>.