

# PsyLingXAV: A Psycholinguistics Design Framework for XAI in Automated Vehicles

Ashkan Y. Zadeh<sup>1,\*</sup>, Xiaomeng Li<sup>1</sup>, Andry Rakotonirainy<sup>1</sup>, Ronald Schroeter<sup>1</sup> and Sebastien Glaser<sup>1</sup>

<sup>1</sup>Centre for Accident Research & Road Safety – Queensland (CARRS-Q), Queensland University of Technology, Australia

## Abstract

The increasing deployment of Automated Vehicles (AVs) necessitates effective Explainable Artificial Intelligence (XAI) models to enhance user trust through transparent decision-making. While these models employ diverse techniques, they often prioritize algorithmic interpretability over human cognitive and linguistic processes. This paper introduces PsyLingXAV, a framework that integrates psycholinguistic principles into AV explanations, leveraging the psychology of language to optimize communication with users. By grounding explanation design in cognitively and linguistically informed strategies, PsyLingXAV tailors outputs to the driving context, advancing AV explainability. This work paves the way for more intuitive, human-centered explanations in AVs, with future studies poised to validate its impact.

## Keywords

Automated Vehicles, Psycholinguistics, Explainable Artificial Intelligence, Natural Language Processing

## 1. Introduction

Automated Vehicles (AVs) offer substantial potential to improve road safety, traffic efficiency [1], and sustainability [2]. However, as AVs directly affect human safety and depend on robust user trust, providing clear explanations of their decisions and behaviors is vital for fostering trust, promoting user acceptance, and facilitating effective human-machine interaction [3]. Such transparency enables users to understand system capabilities and develop confidence in automated driving functions [4]. Yet, this trust is undermined by the opaque decision-making of AI-driven systems [4]. AI models, such as Machine Learning (ML) and Deep Learning (DL), which underpin AV functionality, often operate as black boxes, obscuring their reasoning and making it challenging for users to comprehend how decisions are reached [5]. This lack of transparency presents significant challenges for Explainable Artificial Intelligence (XAI), a growing research field aimed at improving AI interpretability [3]. Despite increasing attention to XAI, no universally accepted framework for AV explainability has emerged [6], due to several factors: (a) a lack of interdisciplinary consensus on what constitutes a meaningful explanation [7], (b) the impracticality of developing a generic, application-agnostic XAI model [8], and (c) insufficient research into context-dependent, user-centric explanations [9, 10, 11]. This research addresses two critical gaps identified in the literature:

### 1.1 The need for application-specific eXplainable AVs (XAV)

AV decision-making relies on ML/DL models [9] that manage either the entire driving task [12, 13] or only specific components [14]. Unlike generic AI applications, AVs operate in dynamic, safety-critical environments where decisions must be rapid, reliable, and contextually appropriate [3, 12]. Current XAI methods, dominated by algorithmic processes such as SHAP [15], provide general algorithmic transparency but are insufficient for the unique demands of AVs. These methods often fail to address domain-specific challenges, such as interpreting sensor data under varying road conditions or balancing

*Late-breaking work, Demos and Doctoral Consortium, colocated with the 3rd World Conference on eXplainable Artificial Intelligence: July 09–11, 2025, Istanbul, Turkey*

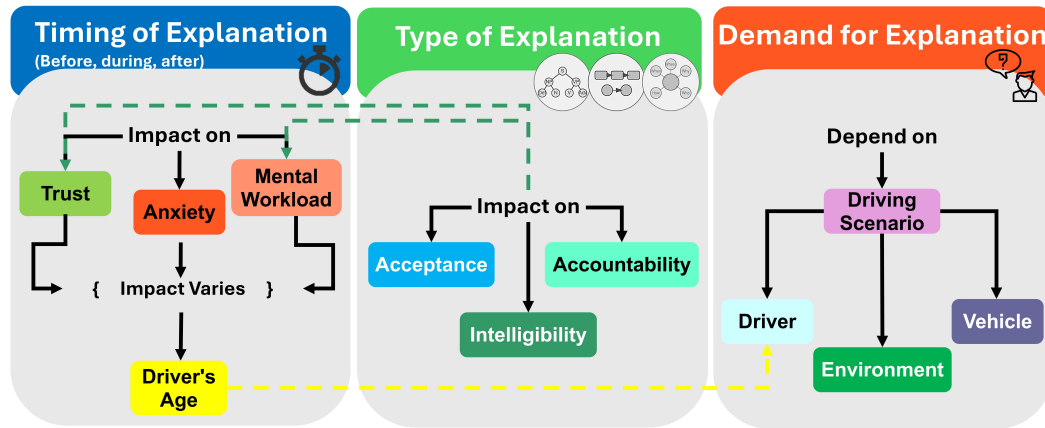
\*Corresponding author.

✉ ashkan.zadeh@hdr.qut.edu.au (A. Y. Zadeh); xiaomeng.li@qut.edu.au (X. Li); r.andry@qut.edu.au (A. Rakotonirainy); r.schroeter@qut.edu.au (R. Schroeter); sebastien.glaser@qut.edu.au (S. Glaser)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).





**Figure 2:** Key Factors Affecting the Effectiveness of Explanations in AVs: Insights from User-Centered Studies.

## 2. Related Work

### 2.1 User-Centred Studies for Explanations in AVs

Studies emphasize that timing is critical for effective explanations in AVs. Proactive (pre-action) explanations—delivered before an AV executes an action—build driver trust, reduce anxiety, and lower mental workload more effectively than post-action explanations [25, 20, 26]. This feedforward approach helps drivers to respond appropriately, though its impact varies by age: younger and middle-aged drivers show lower trust when AVs seek permission pre-action, while older drivers favor it [26]. Post-action explanations, however, consistently undermine trust, especially for older and middle-aged groups, suggesting timing should align with driver demographics [26].

While pre-action explanations are generally preferred, their necessity depends on driving context or scenario. Research shows explanations are more critical in emergencies or near-crashes than routine driving, with demand shaped by both scenario and driver traits [22]. Regarding explanation types, “why-only” explanations (e.g., “The car slowed to avoid a pedestrian”) elicit the most positive emotions, and both “why-only” and “why+how” explanations (e.g., “The car slowed to avoid a pedestrian by braking”) boost acceptance over no explanation [25]. Yet, “why+how” types enhance safety most, despite raising cognitive load, while contrastive “why not” explanations (e.g., “The car didn’t swerve due to a truck nearby”) improve intelligibility and accountability in crises [23]. Trust also hinges on risk perception: simple explanations work best as risk rises, but at extreme levels, no explanation may suffice [21, 22]. A key trust factor is anthropomorphism—attributing human-like traits to AVs [27]. Speech-based explanations increase perceived anthropomorphism and trust compared to non-speech interfaces [28]. Figure 2 captures how timing, demand, and explanation type shape human factors in AV trust, drawing on the aforementioned user-centered research.

### 2.2 Psycholinguistic perspectives on sentence processing

Sentence processing occurs in three main stages: conceptual, syntactic, and phonological. The conceptual stage focuses on forming the intended meaning of a sentence by identifying events and their participants, known as thematic roles [29]. Each component of a sentence serves a specific function, and thematic roles are essential in shaping our understanding of sentences. They provide a structured framework for linking semantic meaning to syntactic structure, making them fundamental to language comprehension [30]. The agent (who performs the action) and the patient (who receives the action) are common roles, but others exist. The syntactic stage structures the sentence by arranging key components such as the subject, verb, and object, though assigning roles can be complex. Finally, the phonological stage prepares the sentence for speech by determining syllables, stress patterns, and

intonation, ensuring smooth articulation [29, 31].

Psycholinguists believe sentence comprehension involves syntax (structure) and semantics (meaning), but they debate whether this occurs in one or two stages [29]. The two-stage model suggests syntax is processed first, then meaning, while the constraint-based model (one stage) argues both happen simultaneously. A “Garden-path” sentence initially leads the reader or listener to an incorrect interpretation before requiring reanalysis due to an unexpected structure or meaning [31, 32, 33]. For example, in “*While Sarah bathed her baby played on the floor,*” we first assume “baby” is the object of “bathed,” but “played” disrupts this, requiring reinterpretation [34]. The garden path model is typically associated with the two-stage approach. Research indicates people read ambiguous sentences more slowly and often reread them [35], which supports the garden path model, where syntax is analyzed before meaning. Because language is processed rapidly and working memory is limited, we rely on heuristics—quick mental shortcuts—that usually work but sometimes cause errors, as with garden path sentences.

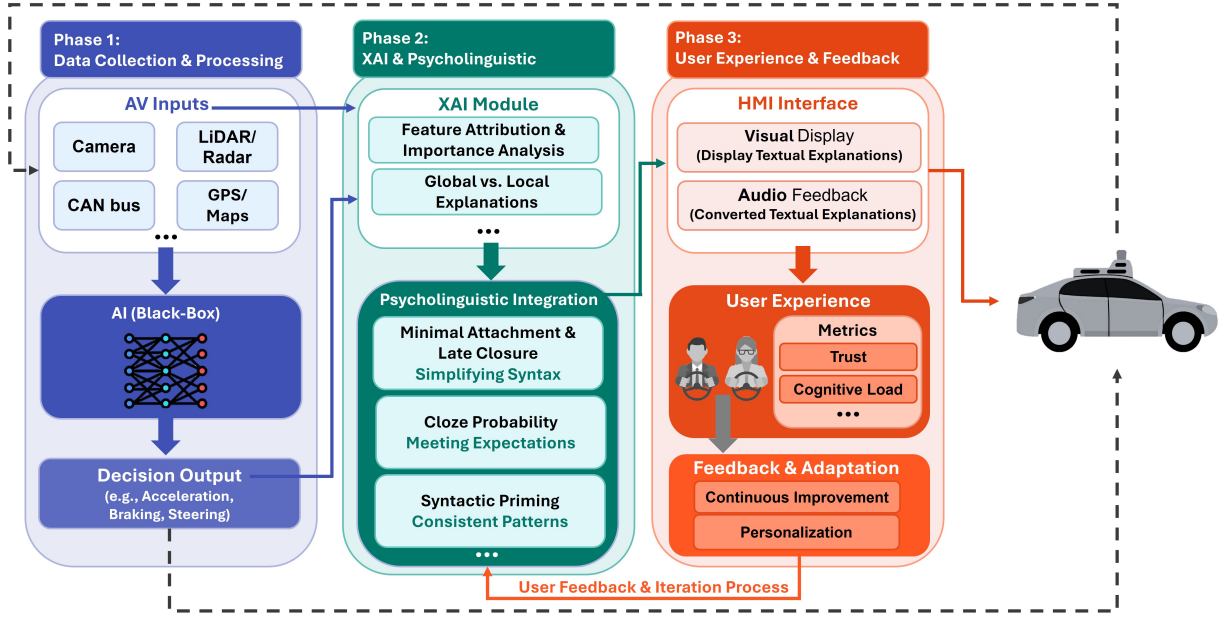
The garden path model explains sentence structure assignment using two heuristics. One is “Late Closure,” which means we keep adding new words to the current phrase unless there is strong evidence to start a new one. This often leads to misinterpretations because our brains tend to assume a subject-verb-object structure [29, 31, 32, 33]. Even when we detect a parsing error, the initial interpretation can linger. Research shows people who hear ambiguous sentences often recall them incorrectly later, demonstrating we rely on a “good enough” approach in everyday communication, where speed can outweigh perfect accuracy [36]. Late closure can also cause early misinterpretation by prematurely “closing” a structure [29, 33]. Some sentences contain complex grammar, like reduced relative clauses, which can be difficult to process. When a word is syntactically ambiguous, the brain tends to pick the more common structure, sometimes leading us down the wrong garden path [29].

“Minimal Attachment” is another garden path heuristic, which assumes the simplest possible sentence structure [32, 33]. This strategy can lead to initial interpretations that later need correction if they clash with the intended meaning [29]. Minimal attachment assumes a prepositional phrase attaches to the main verb rather than to an object noun phrase, simplifying processing but occasionally causing misunderstandings. There are two types of attachment: high attachment which links the phrase to the verb, and low attachment which links it to the object [37]. In syntactic theory, the verb is structurally higher than the object, so high attachment is simpler [29]. However, context also strongly influences how people interpret ambiguous sentences, guiding them to the correct attachment [29, 32, 33].

“Syntactic Priming” is the tendency to repeat a previously encountered sentence structure [38]. Studies show repeated structures are processed more easily, especially if the verb is also repeated [29, 32, 39, 40]. During comprehension, listeners often anticipate upcoming words while constructing sentence structures, a phenomenon known as “Cloze Probability” [41]. Research using the visual world paradigm reveals that visual context helps predict words [42, 43]. Event-Related Potential (ERP) studies show that unexpected words trigger an N400 response, basically a negative brain wave peaking at 400 milliseconds after a stimulus; in this context, it means the brain detects a meaning mismatch [44], reflecting semantic difficulty, while unexpected sentence structures elicit a P600 response, basically a positive brain wave peaking at 600 milliseconds after a stimulus; in this context, it means the brain struggles with sentence structure [45], indicating syntactic processing issues. Expectation thus plays a central role in language comprehension: the brain continuously makes predictions about upcoming words and structures, enabling efficient speech processing despite noise and distractions [46].

### 2.3 The Linkage of Psycholinguistics and Explanation Design in AVs

As shown in Section 2.1 and Figure 2, the time of provision, demand requirements, and different types of AV explanations shape driver trust, mental workload, and other related factors. Psycholinguistic insights from section 2.2 reveal how human can process these explanations in natural language comprehension. The garden path model suggests complex syntactic structures may mislead drivers, requiring reinterpretation that raises cognitive load and anxiety during driving [31]. Heuristics like minimal attachment and late closure show that drivers may favor simple interpretations, and mismatches with



**Figure 3:** PsyLingXAV: A Psycholinguistics Design Framework for XAI in Automated Vehicles.

expectations may erode trust [29].

Demand for explanations ties to language anticipation. The constraint-based model, blending syntax and meaning, implies unclear or context-lacking explanations can confuse drivers, especially in high-stakes scenarios [29]. Experienced drivers, leveraging cloze probability, may need less detail, while novices demand clear, structured explanations to align with their cognitive adaptability [41]. Explanation type affects acceptance and accountability. Syntactic priming indicates familiar structures boost intelligibility by easing effort, while N400 and P600 responses flag processing spikes from unexpected phrasing, potentially lowering acceptance [39, 44]. Misrecall of ambiguous sentences underscores the need for concise, well-timed explanations to ensure accurate understanding [36]. Thus, psycholinguistics presents a valuable opportunity to enhance the effectiveness of AV explanations by offering a user-aligned approach that improves trust, reduces anxiety, and can optimize cognitive load in driver-AV interactions.

### 3. PsyLingXAV Conceptual Framework

Figure 3 illustrates the PsyLingXAV framework, a theoretically grounded approach that integrates psycholinguistics with XAI to transform raw decision outputs from AV systems into cognitively optimized, linguistically structured explanations for the AV’s HMI. The framework operates through three interconnected phases.

In **Phase 1**, the AV collects multisensory inputs, which are processed by an AI black-box model to generate decision outputs. In **Phase 2**, the XAI module extracts critical decision information and refines it using a psycholinguistics module, implemented through natural language processing techniques or large language models. To ensure rapid comprehension during driving, PsyLingXAV employs predefined explanation templates that simplify syntactic structures based on minimal attachment and late closure heuristics, avoiding complexity associated with garden-path sentences. The system also leverages cloze probability and constraint-based parsing to align explanations with user expectations, and syntactic priming to maintain consistent sentence patterns, thereby reducing cognitive load and minimizing trust erosion. The implementation of this phase can follow a similar approach to RAG-Driver [47], a retrieval-augmented, multi-modal large language model that employs in-context learning



to generate explanation patterns similar to those found in the BDD-X dataset [48]. In **Phase 3**, the HMI delivers these explanations through visual displays (text) or audio outputs, enabling real-time driver comprehension. For example, a raw AV output such as *“There is no obstacle detected ahead. The car is initiating acceleration.”* would be translated into a psycholinguistically optimized explanation: *“I’m speeding up since ahead is clear.”* This transformation enhances syntactic simplicity, preserves familiar sentence patterns, and supports faster driver understanding. The framework also integrates a continuous **Feedback & Iteration** process, collecting passive behavioral data (e.g., reaction times, steering corrections, braking) and active user inputs (e.g., survey responses, verbal feedback). These insights allow dynamic adaptation of explanation syntax, timing, and detail to individual user profiles while maintaining psycholinguistic alignment.

## 4. Limitations and future Work

Although the PsyLingXAV framework is based on established psycholinguistic principles, its effectiveness in driving scenarios requires empirical validation. A user study with a VR-based driving simulator will expose participants to three conditions: (1) no explanation, (2) GPT-generated explanations, and (3) psycholinguistically optimized explanations (PsyLingXAV). This controlled setup will allow direct comparison of explanation styles under realistic but safe conditions. Success will be evaluated through trust, cognitive workload, comprehension, and user preference metrics. Findings will inform framework refinements, and future work will explore multimodal extensions, including spoken outputs, with additional research needed to align auditory explanations with psycholinguistic principles.

## 5. Conclusion

This paper introduced PsyLingXAV, a psycholinguistic framework for XAI in AVs, designed to enhance AV explanation effectiveness by aligning them with users’ cognitive and linguistic processes. PsyLingXAV integrates insights from user-centered studies on AV explainability, grounded in robust psycholinguistic principles. While further empirical validation via real-world testing and human-centered evaluations is needed, PsyLingXAV offers a promising approach to improving AV explainability. It lays a foundation for advancing human-centered explanations in AVs through psycholinguistics, paving the way for more intuitive user-AV interactions.

## 6. Acknowledgments

This research was supported partially by the Australian Research Council Discovery projects funding scheme (DP220102598).

## Declaration on Generative AI

The authors used GenAI for grammar check and language assistance during the preparation of this work. They reviewed, edited, and take full responsibility for the final content.

## References

- [1] Z. Mehraban, A. Y. Zadeh, H. Khayyam, R. Mallipeddi, A. Jamali, Fuzzy adaptive cruise control with model predictive control responding to dynamic traffic conditions for automated driving, *Engineering Applications of Artificial Intelligence* 136 (2024) 109008.
- [2] A. Y. Zadeh, H. Khayyam, R. Mallipeddi, A. Jamali, Integrated intelligent control systems for eco and safe driving in autonomous vehicles, *IEEE Transactions on Intelligent Transportation Systems* (2024).

- [3] A. M. Nascimento, L. F. Vismari, C. B. S. T. Molina, P. S. Cugnasca, J. B. Camargo, J. R. de Almeida, R. Inam, E. Fersman, M. V. Marquezini, A. Y. Hata, A systematic literature review about the impact of artificial intelligence on autonomous vehicle safety, *IEEE Transactions on Intelligent Transportation Systems* 21 (2019) 4928–4946.
- [4] S. Tekkesinoglu, A. Habibovic, L. Kunze, Advancing explainable autonomous vehicle systems: A comprehensive review and research roadmap, *ACM Transactions on Human-Robot Interaction* 14 (2025) 1–46.
- [5] W. J. Von Eschenbach, Transparency and the black box problem: Why we do not trust ai, *Philosophy & Technology* 34 (2021) 1607–1622.
- [6] C. Rudin, Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead, *Nature machine intelligence* 1 (2019) 206–215.
- [7] A. Preece, D. Harborne, D. Braines, R. Tomsett, S. Chakraborty, Stakeholders in explainable ai, *arXiv preprint arXiv:1810.00184* (2018).
- [8] P. Linardatos, V. Papastefanopoulos, S. Kotsiantis, Explainable ai: A review of machine learning interpretability methods, *Entropy* 23 (2020) 18.
- [9] D. Calvaresi, A. Najjar, M. Schumacher, K. Främling, Explainable, Transparent Autonomous Agents and Multi-Agent Systems: First International Workshop, EXTRAAMAS 2019, Montreal, QC, Canada, May 13–14, 2019, Revised Selected Papers, volume 11763, Springer Nature, 2019.
- [10] I. Lage, E. Chen, J. He, M. Narayanan, B. Kim, S. Gershman, F. Doshi-Velez, An evaluation of the human-interpretability of explanation, *arXiv preprint arXiv:1902.00006* (2019).
- [11] A. Explainable, The basics policy briefing, Available at [royalsociety.org/ai-interpretability](https://royalsociety.org/ai-interpretability) (2019).
- [12] S. Kuutti, R. Bowden, Y. Jin, P. Barber, S. Fallah, A survey of deep learning applications to autonomous vehicle control, *IEEE Transactions on Intelligent Transportation Systems* 22 (2020) 712–733.
- [13] S. Hecker, D. Dai, L. Van Gool, End-to-end learning of driving models with surround-view cameras and route planners, in: *Proceedings of the european conference on computer vision (eccv)*, 2018, pp. 435–453.
- [14] L. Claussmann, M. Revilloud, D. Gruyer, S. Glaser, A review of motion planning for highway autonomous driving, *IEEE Transactions on Intelligent Transportation Systems* 21 (2019) 1826–1848.
- [15] H. A. Tahir, W. Alayed, W. U. Hassan, A. Haider, A novel hybrid xai solution for autonomous vehicles: Real-time interpretability through lime–shap integration, *Sensors* 24 (2024) 6776.
- [16] K. Panduru, J. Walsh, et al., Exploring the unseen: A survey of multi-sensor fusion and the role of explainable ai (xai) in autonomous vehicles, *Sensors (Basel, Switzerland)* 25 (2025) 856.
- [17] T. Miller, Explanation in artificial intelligence: Insights from the social sciences, *Artificial intelligence* 267 (2019) 1–38.
- [18] D. Wang, Q. Yang, A. Abdul, B. Y. Lim, Designing theory-driven user-centric explainable ai, in: *Proceedings of the 2019 CHI conference on human factors in computing systems*, 2019, pp. 1–15.
- [19] Q. Meteier, M. Capallera, L. Angelini, E. Mugellini, O. A. Khaled, S. Carrino, E. De Salis, S. Galland, S. Boll, Workshop on explainable ai in automated driving: a user-centered interaction approach, in: *Proceedings of the 11th International Conference on Automotive User Interfaces and Interactive Vehicular Applications: Adjunct Proceedings*, 2019, pp. 32–37.
- [20] N. Du, J. Haspiel, Q. Zhang, D. Tilbury, A. K. Pradhan, X. J. Yang, L. P. Robert Jr, Look who’s talking now: Implications of av’s explanations on driver’s trust, av preference, anxiety and mental workload, *Transportation research part C: emerging technologies* 104 (2019) 428–442.
- [21] T. Ha, S. Kim, D. Seo, S. Lee, Effects of explanation types and perceived risk on trust in autonomous vehicles, *Transportation research part F: traffic psychology and behaviour* 73 (2020) 271–280.
- [22] Y. Shen, S. Jiang, Y. Chen, E. Yang, X. Jin, Y. Fan, K. Campbell, To explain or not to explain: A study on the necessity of explanations for autonomous vehicles, *arxiv*, *arXiv preprint arXiv:2006.11684* (2020).
- [23] D. Omeiza, H. Web, M. Jirotko, L. Kunze, Towards accountability: Providing intelligible explanations in autonomous driving, in: *2021 IEEE Intelligent Vehicles Symposium (IV)*, IEEE, 2021, pp. 231–237.
- [24] A. Garnham, *Psycholinguistics (PLE: Psycholinguistics): Central Topics*, Psychology Press, 2013.

- [25] J. Koo, J. Kwac, W. Ju, M. Steinert, L. Leifer, C. Nass, Why did my car just do that? explaining semi-autonomous driving actions to improve driver understanding, trust, and performance, *International Journal on Interactive Design and Manufacturing (IJIDeM)* 9 (2015) 269–275.
- [26] Q. Zhang, X. J. Yang, L. P. Robert Jr, Drivers' age and automated vehicle explanations, *Sustainability* 13 (2021) 1948.
- [27] N. Epley, A. Waytz, J. T. Cacioppo, On seeing human: a three-factor theory of anthropomorphism., *Psychological review* 114 (2007) 864.
- [28] Y. Forster, F. Naujoks, A. Neukum, Increasing anthropomorphism and trust in automated driving functions by adding speech output, in: *2017 IEEE intelligent vehicles symposium (IV)*, IEEE, 2017, pp. 365–372.
- [29] D. Ludden, *The psychology of language: an integrated approach*, Sage Publications, 2015.
- [30] C. Anderson, *Essentials of linguistics*, McMaster University, 2018.
- [31] T. A. Harley, *The psychology of language: From data to theory*, Psychology press, 2013.
- [32] M. J. Traxler, *Introduction to psycholinguistics: Understanding language science* (2011).
- [33] P. Warren, *Introducing psycholinguistics*, Cambridge University Press, 2013.
- [34] N. D. Patson, E. S. Darowski, N. Moon, F. Ferreira, Lingering misinterpretations in garden-path sentences: evidence from a paraphrasing task., *Journal of Experimental Psychology: Learning, Memory, and Cognition* 35 (2009) 280.
- [35] L. Frazier, K. Rayner, Making and correcting errors during sentence comprehension: Eye movements in the analysis of structurally ambiguous sentences, *Cognitive psychology* 14 (1982) 178–210.
- [36] I. Ivanova, M. J. Pickering, H. P. Branigan, J. F. McLean, A. Costa, The comprehension of anomalous sentences: Evidence from structural priming, *Cognition* 122 (2012) 193–209.
- [37] M. J. Pickering, J. F. McLean, H. P. Branigan, Persistent structural priming and frequency effects during comprehension., *Journal of Experimental Psychology: Learning, Memory, and Cognition* 39 (2013) 890.
- [38] J. K. Bock, Syntactic persistence in language production, *Cognitive psychology* 18 (1986) 355–387.
- [39] K. M. Tooley, M. J. Traxler, T. Y. Swaab, Electrophysiological and behavioral evidence of syntactic priming in sentence comprehension., *Journal of Experimental Psychology: Learning, Memory, and Cognition* 35 (2009) 19.
- [40] K. Segaert, G. Kempen, K. M. Petersson, P. Hagoort, Syntactic priming and the lexical boost effect during sentence production and sentence comprehension: An fmri study, *Brain and language* 124 (2013) 174–183.
- [41] E. W. Wlotko, K. D. Federmeier, Age-related changes in the impact of contextual strength on multiple aspects of sentence comprehension, *Psychophysiology* 49 (2012) 770–785.
- [42] M. K. Tanenhaus, M. J. Spivey-Knowlton, K. M. Eberhard, J. C. Sedivy, Integration of visual and linguistic information in spoken language comprehension, *Science* 268 (1995) 1632–1634.
- [43] M. Tanenhaus, Sentence comprehension, *Handbook of Perception and Cognition* 11 (1995).
- [44] L. Hunt III, S. Politzer-Ahles, L. Gibson, U. Minai, R. Fiorentino, Pragmatic inferences modulate n400 during sentence comprehension: Evidence from picture–sentence verification, *Neuroscience Letters* 534 (2013) 246–251.
- [45] L. Osterhout, P. J. Holcomb, Event-related brain potentials elicited by syntactic anomaly, *Journal of memory and language* 31 (1992) 785–806.
- [46] E. Gibson, L. Bergen, S. T. Piantadosi, Rational integration of noisy evidence and prior semantic expectations in sentence interpretation, *Proceedings of the National Academy of Sciences* 110 (2013) 8051–8056.
- [47] J. Yuan, S. Sun, D. Omeiza, B. Zhao, P. Newman, L. Kunze, M. Gadd, Rag-driver: Generalisable driving explanations with retrieval-augmented-in-context learning in multi-modal large language model, *arXiv preprint arXiv:2402.10828* (2024).
- [48] J. Kim, A. Rohrbach, T. Darrell, J. Canny, Z. Akata, Textual explanations for self-driving vehicles, *Proceedings of the European Conference on Computer Vision (ECCV)* (2018).