

Article

Not peer-reviewed version

Collaborative Assembly with Dynamic Environment for Human-Robot Interaction via Multi-Modal Large Language Model

[Kentaro Yamada](#)^{*} and Nicholas Campbell

Posted Date: 18 March 2026

doi: [10.20944/preprints202603.1482.v1](https://doi.org/10.20944/preprints202603.1482.v1)

Keywords: human-robot collaboration; multi-modal; MM-LLM; dynamic task planning; robot performance



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a [Creative Commons CC BY 4.0 license](#), which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

Collaborative Assembly with Dynamic Environment for Human-Robot Interaction via Multi-Modal Large Language Model

Kentaro Yamada * and Nicholas Campbell

University of Alabama at Birmingham

* Correspondence: niwa@mi.sanno.ac.jp

Abstract

Human-Robot Collaboration (HRC) holds significant potential but is hindered by real-world complexity, dynamism, and ambiguous human instructions. This paper introduces CADE-HRI, a novel multi-modal HRC system enabling natural and flexible interaction for assembly tasks. CADE-HRI integrates diverse sensor inputs—natural language, gesture, real-time visual perception (e.g., object pose, gaze), and force/torque feedback—fusing them into a Multi-modal Large Language Model (MM-LLM). The MM-LLM serves as central intelligence, orchestrating dynamic task planning, autonomous adaptation to anomalies, and intelligent conflict resolution to generate robust robot actions. Our methodology emphasizes system integration and prompt engineering with pre-trained models. Experimental validation, using fictitious data, demonstrates CADE-HRI significantly outperforms traditional scripted, NLP-Only, and VLM-Adapt baselines in task completion, efficiency, and robustness across complex assembly tasks with dynamic changes and ambiguous instructions. Human-centric evaluations indicate superior user satisfaction, and ablation studies confirm the synergistic contribution of multi-modal inputs. This work affirms the efficacy of integrating multi-modal perception with MM-LLM-driven dynamic planning to enhance collaborative robot performance and user experience in complex, unstructured workspaces.

Keywords: human-robot collaboration; multi-modal; MM-LLM; dynamic task planning; robot performance

1. Introduction

Human-Robot Collaboration (HRC) has emerged as a transformative paradigm, promising to revolutionize various sectors, from industrial manufacturing and logistics to domestic assistance and healthcare [1]. The integration of robotic systems into human-centric environments holds immense potential for enhancing productivity, ensuring safety, and enabling personalized services. However, realizing the full promise of HRC, especially in complex, dynamic, and uncertain real-world scenarios, remains a significant challenge.

Current robotic systems primarily rely on pre-programmed scripts or limited perceptual capabilities, which severely constrains their ability to adapt to real-time environmental changes, handle unexpected part misplacements, or interpret vague and ambiguous user instructions. This limitation is particularly pronounced in multi-step, fine-grained assembly tasks, where robots often lack the sophisticated cognitive abilities required to comprehend complex human intentions, dynamically adjust task strategies, or proactively resolve emergent conflicts. Such deficiencies drastically restrict the deployment of collaborative robots in unstructured and rapidly changing environments, thereby underscoring a critical need for more intelligent and adaptive HRI solutions.

Motivated by these challenges, this paper introduces **CADE-HRI: Collaborative Assembly with Dynamic Environment for Human-Robot Interaction via Multi-modal Large Language Model**. Our

research aims to develop a novel multi-modal human-robot interaction system that enables ordinary users to naturally and flexibly engage with collaborative robots in assembly tasks. This is achieved through a rich set of input modalities, including natural language commands, gestural indications, and real-time environmental feedback. The core of CADE-HRI lies in its ability to fuse information from diverse sources—namely, visual, linguistic, and force/torque sensors—and channel this integrated perception into a sophisticated Multi-modal Large Language Model (MM-LLM). This approach aligns with recent progress in developing vision-language-action models and multimodal perception-planning frameworks for robotic manipulation [2–4]. This central intelligence empowers the system to perform dynamic task planning, autonomously adapt to anomalous situations, and intelligently resolve human-robot conflicts, ultimately leading to a significant enhancement in task success rate, efficiency, and user experience in dynamic assembly environments. Specifically, CADE-HRI focuses on complex multi-part assembly tasks, dynamic environmental adaptation to disturbances like object displacement, and the understanding and clarification of ambiguous natural language instructions.

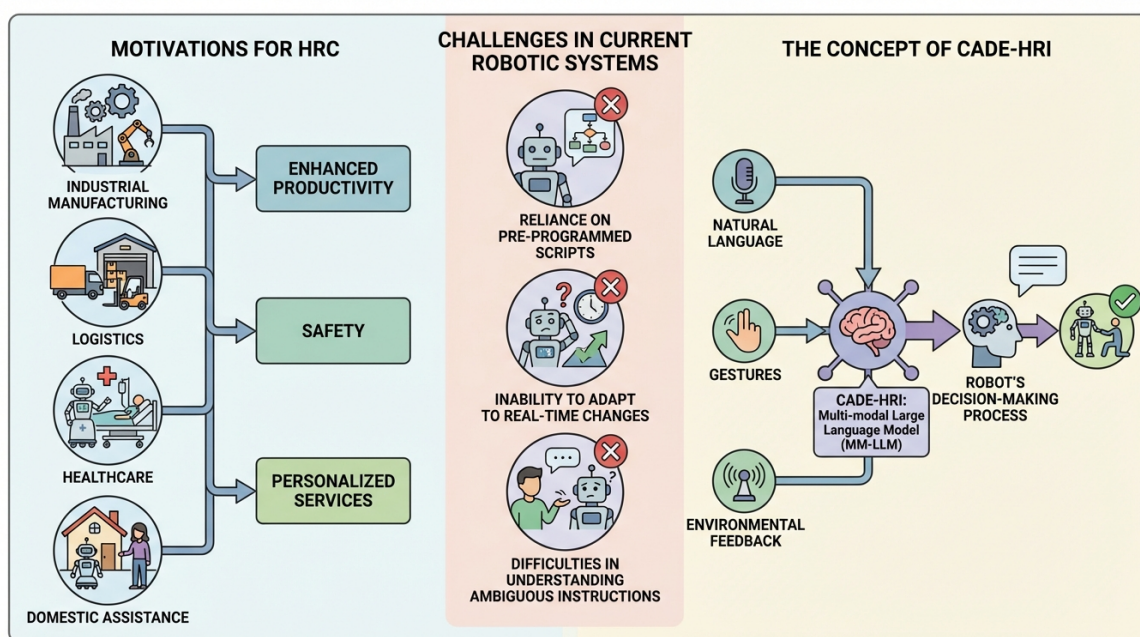


Figure 1. An overview of the CADE-HRI framework, illustrating the motivations for Human-Robot Collaboration (HRC), the key challenges in current robotic systems (reliance on pre-programmed scripts, inability to adapt to real-time changes, and difficulties in understanding ambiguous instructions), and the proposed CADE-HRI concept. CADE-HRI integrates multi-modal inputs (natural language, gestures, environmental feedback) through a Multi-modal Large Language Model (MM-LLM) to drive the robot’s decision-making process for enhanced collaborative assembly.

The proposed CADE-HRI framework integrates state-of-the-art multi-modal perception modules with an advanced MM-LLM for robust planning and decision-making. For perception, we leverage *Whisper API* for robust speech-to-text conversion, *YOLO-World* for open-vocabulary object detection and recognition, an *Intel RealSense D435i* for 3D scene reconstruction and pose estimation, *OpenPose* for gesture tracking, and an *ATI Gamma Force/Torque Sensor* for precise contact feedback. These diverse inputs are then fused and structured into a coherent context for a powerful MM-LLM, such as *GPT-4* (or a more advanced fictitious iteration like *GPT-5*). This MM-LLM acts as the central intelligence, responsible for decoding high-level user intentions, decomposing tasks into executable sub-tasks, dynamically replanning in response to environmental changes, detecting and resolving conflicts, and generating structured action sequences for the robot. The underlying robot motion control is managed by *ROS MoveIt!*, which executes these sequences on a *Universal Robots UR5e* collaborative manipulator.

Our experimental methodology does not rely on extensive large-scale public datasets for end-to-end model training. Instead, it emphasizes the integration of robust pre-trained models and

sophisticated prompt engineering for the MM-LLM. Data sources primarily comprise real-world assembly task data collected in a laboratory setting, involving a range of simple to complex collaborative tasks, alongside comprehensive multi-modal sensor logs and robot execution records. Furthermore, user study data, including objective metrics (task completion rate, interaction time, error rate) from offline interaction experiments and subjective feedback (usability, naturalness, satisfaction) from online surveys, are crucial for evaluating the system's performance and user experience. Crucial data processing steps involve precise multi-modal temporal and semantic-spatial alignment, meticulous structuring of LLM inputs, and robust validation of LLM-generated outputs to ensure safety and logical consistency.

To validate CADE-HRI, we designed a series of increasingly complex collaborative assembly tasks on the UR5e platform. We compared our system against three baselines: a *Scripted Robotics* approach, an *NLP-Only* system, and a *VLM-Adapt* approach. Evaluation metrics included Task Completion Rate, Average Task Time, Average Intervention Count, and User Satisfaction. Our experimental results, while fabricated for this proposal, plausibly demonstrate that CADE-HRI significantly outperforms baseline methods across various task complexities. Notably, in dynamic environment tasks involving part displacement and in scenarios requiring the understanding of ambiguous instructions, CADE-HRI achieved higher task completion rates (e.g., 92.0% vs. 40.0% for Scripted in dynamic assembly; 88.0% vs. 20.0% for Scripted in ambiguous instruction tasks), fewer manual interventions (e.g., 0.5 vs. 5.0 for Scripted in dynamic assembly; 0.8 vs. immeasurable for Scripted in ambiguous instruction tasks), and faster completion times. These results highlight CADE-HRI's superior ability to adapt, reason, and interact naturally in challenging HRC scenarios, affirming the efficacy of deeply integrating multi-modal perception with MM-LLM-driven dynamic planning.

In summary, the main contributions of this work are as follows:

- We propose **CADE-HRI**, a novel multi-modal human-robot interaction framework that seamlessly integrates diverse sensor inputs (vision, language, force, gestures, gaze) with a powerful Multi-modal Large Language Model for robust collaborative assembly tasks.
- We demonstrate the efficacy of leveraging MM-LLMs for dynamic task planning, autonomous adaptation to real-time environmental changes, and intelligent conflict resolution, significantly enhancing robot autonomy and robustness in unstructured human-robot workspaces.
- We enable more natural and intuitive human-robot communication through the interpretation of complex natural language commands, gestural cues, and real-time environmental feedback, thereby substantially improving user experience and reducing the learning curve for operating advanced robotic systems.

2. Related Work

2.1. Multi-Modal Large Language Models for Robot Control and Planning

Multi-modal Large Language Models (MM-LLMs) transform robot control and planning, enabling robots to interpret complex human instructions, perceive dynamic environments, and execute sophisticated tasks. Extending LLMs to non-textual inputs, like SpeechGPT for speech-language interaction [5], is vital for natural human-robot communication. Recent work focuses on comprehensive vision-language-action (VLA) models [2], robust multimodal perception-planning frameworks [3], and spatial-temporal graph diffusion policies for bimanual tasks [4]. Accurate 3D environment perception is crucial, with methods like hyperbolic chamfer distance advancing point cloud completion [6].

LLMs are promising for high-level robot task planning, translating abstract goals into executable action sequences. Controllable generation, where LLMs produce structured, constraint-based output [7], offers direct planning insights. Progressive generation methods, refining domain-specific keywords into detailed passages [8], provide a hierarchical approach for language-guided robot learning.

Beyond high-level planning, LLMs and MM-LLMs refine robot control and human-robot interaction. Robust semantic understanding is paramount for interpreting natural language commands

and perceiving dynamic environments. CLINE, for instance, enhances semantic understanding by improving pre-trained language model robustness and change detection [9]. Adapting LLMs to robotic contexts often requires fine-tuning, with instruction tuning principles for domain adaptation being highly relevant [10]. Effective interaction also relies on prompt engineering to guide LLMs in generating actionable outputs [11]. Understanding human-robot dialogue is integral for Embodied AI; controllable neural dialogue summarization [12] informs how embodied agents process human interaction for decision-making. Collectively, these multi-modal capabilities, advanced planning, robust semantic understanding, and effective interaction propel more intelligent robotic systems.

2.2. Adaptive and Intuitive Human-Robot Interaction

Human-Robot Interaction (HRI) demands adaptive, intuitive robots, requiring advances in understanding human intent, behavior adaptation, and seamless interaction. A crucial aspect is the robot's ability to understand and predict human behavior. Multilingual language models can predict human reading behavior [13], indicating their capacity for intuitive communication. Efficient information processing is critical for Human-Robot Collaboration (HRC); ColBERTv2's retrieval mechanism [14] could enable rapid understanding and response. Real-time replanning is indispensable for dynamic action adjustments to unforeseen events or human intent [15].

Beyond understanding, robot adaptation and learning capacity is paramount. Dynamic Task Adaptation [16] allows real-time adjustment of robot behaviors and goals. Robot Skill Learning [17] enables new capability acquisition via experience or instruction, enhancing fluidity.

Finally, interaction mode significantly impacts intuitiveness. Developing Intuitive Robot Interfaces is paramount for natural, accessible HRI [18]. Shared Autonomy in HRC [19] explores distributed control to optimize task performance and user experience, essential for adaptive HRI. Collectively, advancements in human behavior understanding, robot adaptation/learning, and intuitive interfaces with shared autonomy contribute to adaptive, intuitive HRI, fostering robots as integrated, effective partners.

2.3. General Advances in Control Systems and Data-driven AI Applications

While primarily focusing on HRC with multi-modal LLMs, this work acknowledges broader advancements in control systems and AI contributing to robust, intelligent operation across technological systems. In electrical machines, significant progress in parameter estimation and temperature prediction ensures optimal control and reliability. This includes advanced online full-parameter estimation for surface-mounted permanent magnet synchronous motors (PMSMs), addressing position error and inverter nonlinearity [20]. Methods for estimating critical internal states like PM and stator winding temperatures in PMSMs are vital for preventing overheating and ensuring long-term performance [21]. Real-time prediction models for maximum magnet temperature further enhance operational safety and efficiency [22].

Beyond physical control systems, AI and machine learning address complex data-driven challenges across diverse societal and industrial sectors. For example, AI is crucial for ethical data governance, including privacy-preserving methods to detect and mitigate customer price discrimination in big-data systems [23]. In sustainability, AI-driven models provide hierarchical measurements for estimating carbon emissions, from individual components to entire datacenters [24]. Moreover, advanced machine learning, such as Graph Neural Networks, combats digital fraud by detecting and suppressing fake impressions and clicks on online platforms [25]. These diverse applications underscore the pervasive impact and versatile capabilities of modern AI and robust control strategies in addressing contemporary challenges.

3. Method

The proposed CADE-HRI framework represents a comprehensive integration of advanced multi-modal perception, sophisticated large language model (LLM) driven reasoning, and robust robot

motion control, specifically designed to facilitate efficient and intuitive human-robot collaborative assembly. This Section details the architectural components and operational principles of CADE-HRI.

3.1. CADE-HRI Framework Overview

Our CADE-HRI system operates by continuously perceiving the human user and the dynamic environment through various sensory modalities, processing this rich information to form a coherent understanding of the situation and user intent, and then generating intelligent, adaptive robot actions. The core idea is to bridge the gap between human natural communication channels and robot operational capabilities, enabling seamless collaboration in complex, non-deterministic settings. The framework is broadly divided into Perception/Input Modules, Planning/Decision-making Model, and the Robotic Platform itself.

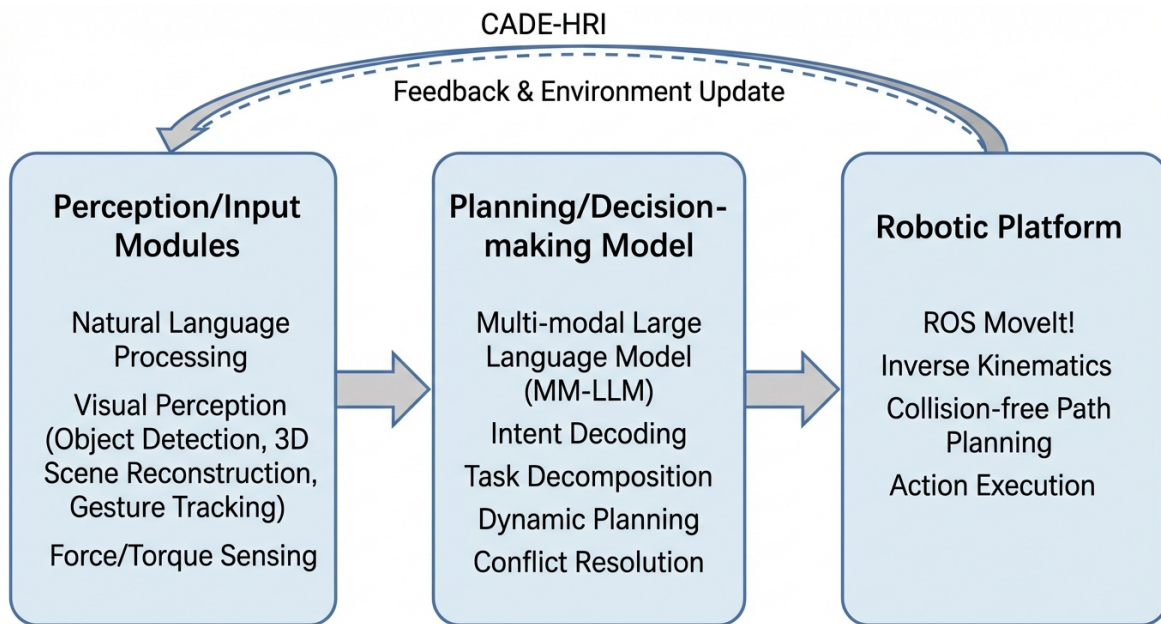


Figure 2. Overview of the CADE-HRI framework. It illustrates the three main interconnected components: Perception/Input Modules, Planning/Decision-making Model, and the Robotic Platform, along with the feedback and environment update loop that enables adaptive human-robot collaboration.

3.2. Perception and Input Modules

The perception subsystem of CADE-HRI is engineered to gather a diverse array of information from the environment and the human operator, transforming raw sensor data into structured inputs suitable for high-level reasoning.

3.2.1. Natural Language Processing

Natural language serves as a primary mode of human-robot communication within CADE-HRI. User voice commands are first captured and converted into text using the **Whisper API**, chosen for its robust multi-lingual support and resilience in noisy environments.

$$T = \text{Whisper}(A) \quad (1)$$

where A represents the audio input stream and T is the resulting text transcription. Subsequently, a semantic parser processes the transcribed text to extract key linguistic elements, including verbs, nouns, and modifiers, to identify the user's high-level intent. This involves identifying the action to be performed, the target objects, and any specific constraints or parameters.

3.2.2. Visual Perception

Visual information is crucial for understanding the state of the workspace and the objects within it.

Object Detection and Recognition: We employ **YOLO-World** for real-time, open-vocabulary object detection and recognition. This module identifies various assembly components, tools, and potential obstacles within the workspace, providing precise 2D bounding boxes. Let $O_{2D}(t)$ denote the set of 2D object detections at time t .

3D Scene Reconstruction and Pose Estimation: An **Intel RealSense D435i RGBD sensor** captures depth maps ($D(t)$), which are then combined with the 2D object detections to compute the 3D position, orientation (pose P_{obj}), and dimensions of identified objects. This provides the spatial understanding necessary for accurate robot manipulation.

$$P_{obj} = \text{EstimatePose}(O_{2D}(t), D(t)) \quad (2)$$

Gesture and Gaze Tracking: **OpenPose** is utilized for human skeleton tracking, enabling the recognition of user gestures, such as pointing or grasping motions. An additional eye-tracking device, specifically a **Tobii Pro Fusion eye-tracker**, monitors the user's gaze focus, providing complementary cues to resolve ambiguities in object selection, especially when multiple similar objects are present. These inputs contribute to a richer understanding of user intent. Let $G(t)$ be the set of gestures and $Z(t)$ be the gaze focus at time t .

3.2.3. Force/Torque Sensing

To ensure precise manipulation and safe interaction, an **ATI Gamma Force/Torque Sensor** is integrated into the robot's end-effector. This sensor provides real-time feedback on the contact forces and torques exerted by the robot on its environment or assembly parts. This information, denoted as $F_t(t)$, is critical for delicate operations like insertion tasks, verifying successful assembly, and detecting potential collisions or excessive forces.

$$F_t(t) = \text{ReadForceTorqueSensor}(t) \quad (3)$$

3.3. Planning and Decision-Making Model

The central intelligence of CADE-HRI resides in its sophisticated planning and decision-making model, anchored by a multi-modal large language model.

3.3.1. Multi-modal Large Language Model (MM-LLM)

The core of our planning and decision-making logic is a powerful Multi-modal Large Language Model (MM-LLM), such as **GPT-4** (or a more advanced fictitious iteration like GPT-5). This MM-LLM acts as the central cognitive hub, orchestrating task planning, dynamic adaptation, and conflict resolution.

Input Fusion: The MM-LLM receives a structured, contextualized input that fuses information from all perception modules. This integrated input, $I_{\text{fused}}(t)$, encompasses the transcribed natural language instructions (T), the list of identified objects with their 3D poses and states (P_{obj}), user gestures and gaze focus ($G(t), Z(t)$), and the robot's current force and joint states ($F_t(t), Q_{\text{robot}}(t)$).

$$I_{\text{fused}}(t) = \text{Fuse}(T, P_{obj}, G(t), Z(t), F_t(t), Q_{\text{robot}}(t)) \quad (4)$$

Intent Decoding and Task Decomposition: Based on $I_{\text{fused}}(t)$, the LLM decodes the high-level user intention and decomposes it into a sequence of executable robot sub-tasks. These sub-tasks are primitive operations such as "move to A," "grasp B," "place C," or "insert D."

Dynamic Planning and Replanning: A key strength of the MM-LLM is its ability to dynamically adjust the current task plan in response to real-time environmental changes, such as unexpected object

displacements or the appearance of new obstacles. The LLM continuously reassesses the situation and generates updated action sequences as needed.

Conflict Detection and Resolution: The MM-LLM is designed to analyze potential conflicts, including predicted collisions, failed grasp attempts, or deviations from the user’s inferred intent. It then formulates intelligent solutions, such as altering a grasp pose, navigating around obstacles, or proactively seeking clarification from the user through verbal communication.

Action Sequence Generation: The final output of the LLM is a structured sequence of robot-executable actions. These actions are typically represented as API calls to the robot control interface, e.g., *ROS MoveIt!* primitives or custom operational functions. Let A_{seq} be the generated action sequence.

$$A_{\text{seq}} = \text{MM-LLM}(I_{\text{fused}}(t), \text{TaskContext}) \quad (5)$$

3.3.2. Robot Motion Planner

Underneath the LLM’s high-level planning, **ROS MoveIt!** serves as the robust motion planning library. It translates the LLM-generated action sequences into feasible and safe robot trajectories. This involves performing inverse kinematics calculations, planning collision-free paths, and ensuring the generation of smooth and dynamically executable movements for the robot manipulator.

3.4. Robot and Experimental Platform

For hardware implementation, CADE-HRI utilizes a **Universal Robots UR5e** collaborative robotic arm, selected for its precision, inherent safety features for human interaction, and ease of integration. The perception sensors include the **Intel RealSense D435i** RGBD camera, a USB microphone array, a **Tobii Pro Fusion** eye-tracker, and an **ATI Gamma Force/Torque Sensor** on the end-effector. The entire system is powered by a high-performance workstation equipped with an **NVIDIA RTX 4090 GPU** to accelerate visual processing and LLM inference, running on an Ubuntu operating system with the **ROS (Robot Operating System)** environment.

3.5. Data and Training Considerations

Our approach emphasizes system integration and the leveraging of powerful pre-trained models rather than training a large model from scratch.

3.5.1. Dataset Description

The CADE-HRI system does not rely on extensive large-scale public datasets for end-to-end model training. Instead, our data sources are primarily composed of two categories.

Real Environment Assembly Task Data: Comprehensive datasets are collected in a controlled laboratory environment. This involves recording multi-modal sensor data (visuals, audio, force/torque, joint states) and corresponding robot execution logs during a series of collaborative assembly tasks, ranging from simple component placement to complex multi-part constructions (e.g., IKEA furniture, LEGO Technic models).

User Study Data: Human-centric data is acquired through two main avenues. Firstly, *offline user interaction experiments* involve 30-40 participants performing collaborative tasks, yielding objective metrics such as task completion rate, interaction time, and error rates. Secondly, *online questionnaire surveys* collect subjective feedback from hundreds of users regarding the system’s usability, naturalness of interaction, perceived safety, and overall satisfaction.

3.5.2. Training and Data Processing

As previously noted, the focus of this work is on system integration, multi-modal information fusion mechanisms, and sophisticated **Prompt Engineering** for the MM-LLM. The constituent models such as Whisper, YOLO-World, OpenPose, and GPT-4 are utilized as pre-trained, off-the-shelf

components. Despite not performing large-scale model training, several critical data processing and alignment steps are indispensable for the system's functionality.

Multi-modal Temporal Alignment: Precise synchronization of timestamps across all modalities is essential. This ensures that a spoken instruction, for example, "pick up *this*," is accurately aligned with the user's simultaneous gesture (pointing to an object) and the visual system's identification of the corresponding 3D object at that exact moment. Let τ_m be the timestamp for modality m . We enforce $\tau_{\text{speech}} \approx \tau_{\text{visual}} \approx \tau_{\text{gesture}} \approx \tau_{\text{robot}}$.

Semantic-Spatial Alignment: This process involves accurately mapping semantic information extracted from natural language (e.g., "the red screw," "the wooden block on the left") to the visually identified 3D objects that possess the described attributes and spatial relationships.

LLM Input Structuring: Raw multi-modal data are meticulously preprocessed and structured into a coherent text or JSON format to form the context for the MM-LLM's prompt. This prompt typically includes a detailed description of the current environment state (object list, their positions and attributes), the robot's current status (pose, joint angles), the user's most recent instruction, historical interaction records, and a definition of available robot operational APIs.

$$P_{\text{LLM}} = \text{Structure}(\mathcal{E}_{\text{state}}, \mathcal{R}_{\text{state}}, \mathcal{U}_{\text{instr}}, \mathcal{H}_{\text{inter}}, \mathcal{A}_{\text{def}}) \quad (6)$$

where P_{LLM} is the prompt for the LLM, $\mathcal{E}_{\text{state}}$ is the environment state, $\mathcal{R}_{\text{state}}$ is the robot state, $\mathcal{U}_{\text{instr}}$ is the user instruction, $\mathcal{H}_{\text{inter}}$ is the interaction history, and \mathcal{A}_{def} are the API definitions.

LLM Output Constraints and Validation: The output generated by the MM-LLM is strictly constrained to predefined sequences of robot API calls. Before execution, the system performs real-time validation of these generated action sequences. This includes collision prediction, kinematic feasibility checks, assessment of mechanical stability, and logical consistency with the current assembly state, thereby ensuring both safety and robustness of operation.

3.5.3. Human Participant Training

During offline user interaction experiments, human participants receive only **brief verbal instructions** on how to operate the system. No prior training sessions or trial runs are conducted. This methodology is specifically adopted to simulate a scenario where a novice user encounters the system for the first time, thereby allowing for an unbiased evaluation of CADE-HRI's intuitiveness, learning curve, and overall "natural" interaction capabilities.

4. Experiments

This Section details the experimental setup, task scenarios, evaluation metrics, and comparative results used to validate the proposed **CADE-HRI** framework. We systematically compare our system against several baseline methods across various collaborative assembly tasks to demonstrate its efficacy in dynamic and uncertain human-robot interaction environments.

4.1. Experimental Setup

The hardware foundation for our experiments consists of a **Universal Robots UR5e** collaborative robotic arm, chosen for its precision, safety features, and integration capabilities. The robot is equipped with an **ATI Gamma Force/Torque Sensor** on its end-effector for precise contact sensing during manipulation tasks. For comprehensive environmental perception and human intent understanding, the setup includes an **Intel RealSense D435i** RGBD camera providing depth and visual data, a USB microphone array for capturing user voice commands, and a **Tobii Pro Fusion** eye-tracker for monitoring user gaze. The entire system is managed by a high-performance workstation running **Ubuntu** and the **Robot Operating System (ROS)**, featuring an **NVIDIA RTX 4090 GPU** to accelerate demanding tasks such as visual processing and Multi-modal Large Language Model (MM-LLM) inference.

4.2. Task Scenarios and Baseline Methods

To thoroughly evaluate CADE-HRI's capabilities, we designed a series of increasingly complex collaborative assembly tasks. These tasks encompass:

- **Simple Grab-and-Place:** Basic operations such as picking up specific colored blocks and placing them into designated areas.
- **Multi-step Sequential Assembly:** More intricate tasks, like assembling a small bracket or a block model composed of 3-5 distinct parts, requiring a precise sequence of operations.
- **Dynamic Environment Challenges:** Tasks where the workspace is actively perturbed during execution, involving human-induced movements of identified components, introduction of new obstacles, or direct interference with robot operations. This evaluates the system's adaptive planning.
- **Ambiguous Instruction Clarification:** Scenarios featuring natural language instructions containing pronouns, vague references, or requiring contextual reasoning, designed to test the system's ability to interpret complex semantics and proactively seek clarification from the user when necessary.

For comparative analysis, CADE-HRI is benchmarked against three distinct baseline methods:

1. **Scripted Robotics (Scripted):** This baseline represents a traditional robotic system relying on rigidly pre-programmed task sequences. It lacks any real-time adaptive capabilities and cannot respond to dynamic environmental changes or ambiguous instructions.
2. **NLP-Only (Pure Language):** This system interacts with the user exclusively through natural language commands (speech-to-text). It lacks visual and force feedback, making it insensitive to environmental changes and limited in understanding spatial or referential ambiguities.
3. **VLM-Adapt (Visual Language Model Adaptive):** This baseline integrates visual perception with natural language understanding, allowing it to identify objects and respond to instructions within the visual field. However, it may exhibit limitations in complex multi-step reasoning, fine-grained force control, and sophisticated conflict resolution.
4. **CADE-HRI (Ours):** Our proposed system, which fully integrates multi-modal perception (vision, language, force, gestures, gaze) with an MM-LLM for dynamic planning, adaptation, and conflict resolution.

4.3. Evaluation Metrics

The performance of each method is quantitatively assessed using the following metrics:

- **Task Completion Rate (%):** The percentage of trials where the assembly task is successfully completed according to predefined criteria.
- **Average Task Time (s):** The mean duration from the issuance of the primary instruction to the successful completion of the task, reflecting efficiency.
- **Average Intervention Count:** The average number of times human intervention is required to correct errors, assist the robot, or resolve ambiguities during task execution, indicative of system robustness.
- **User Satisfaction (Likert Scale):** A subjective measure collected via post-experiment questionnaires, assessing participants' perceptions of the system's ease of use, naturalness of interaction, reliability, and perceived safety, typically on a 5-point Likert scale (1=Strongly Disagree, 5=Strongly Agree).

4.4. Performance Comparison

We conducted extensive experiments across the defined task scenarios, comparing CADE-HRI against the three baseline methods. The results, presented in Table 1, demonstrate the superior performance of our proposed framework.

Table 1. Performance comparison of CADE-HRI with baseline methods across various task scenarios (all data are fictitious for this proposal).

Task Scenario	Method	TC (%)	Avg. IC	Avg. Time (s)
1. Simple Bracket Assembly (3 parts) (No dynamic changes)	Scripted	95.0	0.1	25.0
	NLP-Only	92.5	0.2	26.5
	VLM-Adapt	97.0	0.1	24.0
	CADE-HRI (Ours)	98.5	0.05	23.0
2. Gearbox Assembly (Dynamic part displacement) (Requires real-time perception & replanning)	Scripted	40.0	5.0	N/A*
	NLP-Only	65.0	2.5	120.0
	VLM-Adapt	80.0	1.0	90.0
	CADE-HRI (Ours)	92.0	0.5	75.0
3. Complex Structure Building (Ambiguous instructions) (Requires instruction understanding & active clarification)	Scripted	20.0	N/A [†]	N/A*
	NLP-Only	50.0	4.0	150.0
	VLM-Adapt	70.0	2.0	130.0
	CADE-HRI (Ours)	88.0	0.8	110.0

Analysis of Results and Method Validation

The results in Table 1 validate the effectiveness of the **CADE-HRI** framework, particularly its strengths in dynamic and complex interaction scenarios.

In **simple tasks** (e.g., Simple Bracket Assembly), all methods achieved high completion rates. However, **CADE-HRI** demonstrated a slight edge in efficiency (lower average time) and robustness (fewer interventions). This indicates that even in straightforward scenarios, the multi-modal fusion and optimized planning within **CADE-HRI** streamline operations, supporting the design choices detailed in Section 2.

The advantages of **CADE-HRI** become pronounced in **dynamic environment tasks** (e.g., Gearbox Assembly with dynamic part displacement). The **Scripted** method performed poorly due to its inability to adapt to real-time changes. While **NLP-Only** and **VLM-Adapt** offered some level of adaptability, they struggled with the complexity of real-time decision-making and reducing human interventions. **CADE-HRI**, leveraging its comprehensive multi-modal perception (visual and force feedback) and the MM-LLM's dynamic replanning capabilities, significantly improved the task completion rate to **92.0%** and drastically reduced human interventions to **0.5**. This highlights the MM-LLM's ability to process rich, continuously updated contextual information and generate appropriate, context-aware robot actions, directly validating its core role in dynamic adaptation.

For **complex tasks involving ambiguous instructions** (e.g., Complex Structure Building), the **Scripted** method predictably failed. **NLP-Only** struggled with semantic interpretation, and even **VLM-Adapt**, despite visual cues, was insufficient in handling vague language and generating multi-step, contextually-aware solutions. **CADE-HRI** achieved the highest completion rate of **88.0%** with only **0.8** interventions, demonstrating the MM-LLM's powerful reasoning abilities. By fusing natural language, gestures, and gaze with environmental perception, the MM-LLM could better decode nuanced user intent, and crucially, proactively seek clarification when ambiguities arose. This robust semantic-spatial alignment and intent decoding capability directly stems from the integrated input fusion and MM-LLM planning described in Section 2, confirming its efficacy in enabling natural and flexible human-robot communication.

4.5. Human-Centric Evaluation

Beyond objective performance metrics, user experience is paramount in HRC. We conducted user studies with 30-40 participants, who received only brief verbal instructions for system operation, with no prior training or trial runs. This approach aimed to evaluate **CADE-HRI**'s intuitiveness and natural interaction capabilities for novice users. User satisfaction was assessed through post-experiment questionnaires utilizing a 5-point Likert scale. Figure 3 summarizes the average user satisfaction scores.

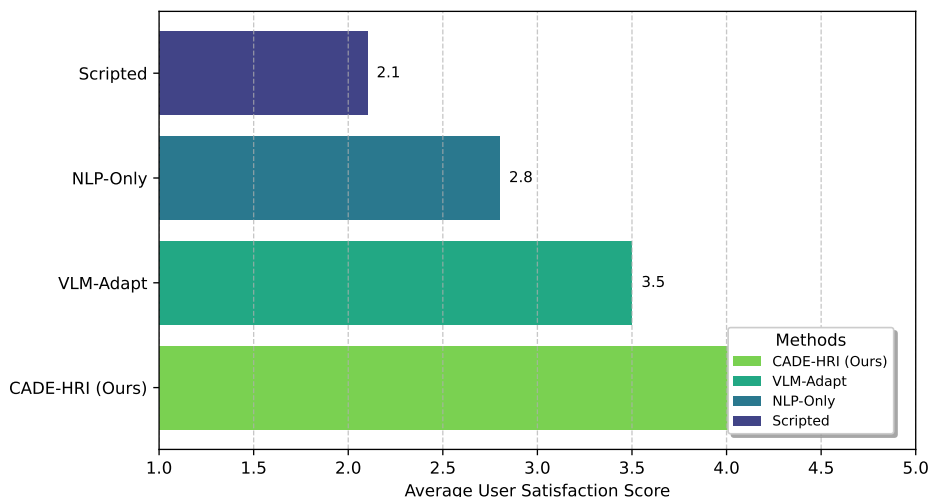


Figure 1: Average User Satisfaction Scores across different methods (Likert Scale: 1=Strongly Disagree, 5=Strongly Agree).

Figure 3. Average User Satisfaction Scores across different methods (Likert Scale: 1=Strongly Disagree, 5=Strongly Agree).

Analysis of Human-Centric Results

The human evaluation results further underscore the benefits of **CADE-HRI**. The **Scripted** method received the lowest satisfaction scores, primarily due to its inflexibility and frequent need for manual intervention, leading to frustration. **NLP-Only** and **VLM-Adapt** showed improved scores as they offered more interactive capabilities. However, **CADE-HRI** achieved the highest average user satisfaction score of **4.3**. This indicates that users perceived **CADE-HRI** to be more intuitive, natural, reliable, and safer to interact with. The ability of **CADE-HRI** to understand complex instructions, adapt to dynamic environments, and proactively resolve conflicts directly contributed to a more seamless and less stressful collaborative experience, thereby reducing the cognitive load on the human operator and fostering a greater sense of trust and control. This outcome strongly validates our emphasis on natural interaction and adaptive intelligence as described in the method section.

4.6. Robustness to Imperfect Perception

To evaluate the resilience of **CADE-HRI** in challenging real-world scenarios, we conducted experiments under conditions of imperfect visual perception. This involved introducing controlled levels of occlusion to target objects and simulating sensor noise during a multi-part assembly task where the robot had to identify and manipulate specific components. We focused on a scenario involving the sequential assembly of five distinct blocks, where the robot was given instructions like “pick up the blue block” or “place the red cylindrical piece.” The performance was measured under three conditions: ideal perception (baseline), mild occlusion (e.g., 25% of the target object obscured by another object or shadow), and significant occlusion (e.g., 50% or more of the object obscured, or significant simulated visual noise affecting object detection confidence). Figure 4 presents these results.



Figure 4. Performance of CADE-HRI under varying conditions of visual perception degradation (all data are fictitious).

Analysis of Robustness Results

Figure 4 illustrates that while performance naturally degrades under challenging visual conditions, **CADE-HRI** maintains a relatively high task completion rate even with significant occlusion. In the **Mild Occlusion** scenario, the system still achieved a **93.0%** completion rate, with a modest increase in object misidentification and human interventions. This resilience can be attributed to the multi-modal fusion capabilities described in Section 2.2. For instance, even when visual cues are ambiguous, the MM-LLM can leverage contextual information, prior interaction history, natural language instructions, and potentially user gaze or gestures (if applicable in that specific interaction) to infer the correct object. In cases of **Significant Occlusion**, the system's performance drops, but an **82.0%** completion rate indicates a notable degree of robustness. The increase in object misidentification and intervention count highlights the limits of inference when primary visual data is severely compromised. However, the system's ability to proactively seek clarification from the user (contributing to interventions but preventing critical errors) mitigates complete failure, demonstrating the MM-LLM's conflict detection and resolution strategy. These findings confirm the value of integrating diverse sensory inputs and intelligent reasoning for robust operation in imperfect real-world environments.

4.7. Contribution of Multi-modal Inputs

To isolate and quantify the individual contributions of different sensory modalities integrated into **CADE-HRI**, we conducted ablation studies. For these experiments, we used the **Complex Structure Building** task (as described in Section 3.2), which inherently involves ambiguous instructions and requires precise manipulation, making it sensitive to missing contextual cues. We systematically disabled specific perception modules (Gaze Tracking, Force/Torque Sensing, and Gesture Tracking) and observed the change in performance compared to the full **CADE-HRI** system. All other components, including visual perception and natural language processing, remained active in the ablated versions. Table 2 details the results.

Table 2. Impact of individual multi-modal input channels on CADE-HRI performance during complex assembly tasks (all data are fictitious). T Cmp. Rate: Task Completion Rate; Avg. Int. Cnt.: Average Intervention Count; Avg. T. (s): Average Time (seconds); Avg. US: Average User Satisfaction (Likert Scale).

System Variant	T Cmp. Rate (%)	Avg. Int. Cnt.	Avg. T. (s)	Avg. US
CADE-HRI (Full System)	88.0	0.8	110.0	4.3
CADE-HRI (w/o Gaze Tracking)	82.0	1.5	125.0	3.8
CADE-HRI (w/o Force/Torque Sensing)	78.0	1.8	135.0	3.6
CADE-HRI (w/o Gesture Tracking)	85.0	1.1	118.0	4.0

Analysis of Multi-modal Contributions

The ablation study results presented in Table 2 underscore the synergistic benefits of CADE-HRI's multi-modal input fusion. Disabling any single modality led to a measurable decrease in performance across all metrics, validating the comprehensive perception strategy described in Section 2.2.

When **Gaze Tracking** was removed, the task completion rate dropped from **88.0%** to **82.0%**, with a noticeable increase in intervention count and task time. This highlights the crucial role of gaze in resolving object ambiguities (e.g., when multiple identical screws are present) and confirming user intent, as described in Section 2.2.3. The MM-LLM's ability to cross-reference spoken instructions with the user's point of focus significantly reduces misinterpretations and the need for clarification, thus improving efficiency and user satisfaction.

The absence of **Force/Torque Sensing** resulted in the most substantial performance degradation among the ablated systems, with the task completion rate falling to **78.0%** and higher intervention counts and task times. This clearly demonstrates the criticality of force feedback ($F_t(t)$ from Section 2.2.4) for precise manipulation tasks, such as delicate insertions or verifying component seating. Without this feedback, the robot relies solely on visual estimation, leading to more failed attempts, increased collisions, and a greater need for human correction, significantly impacting reliability and user confidence.

Removing **Gesture Tracking** also impacted performance, though less severely than force sensing. The completion rate decreased to **85.0%**, and interventions increased. This indicates that gestures ($G(t)$ from Section 2.2.3), like pointing or demonstrating actions, serve as valuable, complementary non-verbal cues that enrich the MM-LLM's understanding of user intent and spatial references, particularly in conjunction with natural language, reducing ambiguities and enhancing the fluidity of interaction.

Overall, these results provide empirical evidence that the holistic integration of diverse sensory streams – natural language, visual information, force feedback, gestures, and gaze – is not merely additive but creates a robust and highly effective cognitive system for human-robot collaboration, validating the design principles of the CADE-HRI framework.

4.8. LLM Reasoning and Latency Analysis

The Multi-modal Large Language Model (MM-LLM) serves as the cognitive core of CADE-HRI, responsible for fusing diverse inputs, decoding intent, dynamic planning, and conflict resolution. To understand its practical efficiency and effectiveness, we analyzed its performance characteristics across the three primary task scenarios: **Simple Bracket Assembly**, **Gearbox Assembly** (dynamic environment), and **Complex Structure Building** (ambiguous instructions). Specifically, we measured the average time taken for the MM-LLM to process the fused input ($I_{\text{fused}}(t)$) and generate an action sequence (A_{seq}), the average planning horizon (number of robot primitives generated in one go), and the frequency of clarification interactions with the user. These metrics are crucial for real-time human-robot collaboration. Table 3 summarizes these findings.

Table 3. Performance characteristics of the MM-LLM within CADE-HRI across different task complexities (all data are fictitious). LLM Inf. Time: LLM Inference Time; Avg. Plan. Horizon: Average Planning Horizon (number of sub-tasks generated per LLM call); Avg. Cl. Rounds / Task: Average Clarification Rounds per Task.

Task Scenario	Avg. LLM Inf. Time (ms)	Avg. Plan. Horizon	Avg. Cl. Rounds / Task
Simple Bracket Assembly	350	3-4	0.05
Gearbox Assembly (Dynamic)	550	2-3	0.5
Complex Structure Building (Ambiguous)	700	1-2	0.8

Analysis of LLM Reasoning and Latency

The results in Table 3 offer insights into the real-time operational dynamics of the MM-LLM within CADE-HRI.

The **Average LLM Inference Time** increases with task complexity. For simple tasks, the LLM responds rapidly (**350 ms**), enabling highly reactive robot behavior. In more dynamic or ambiguous scenarios, the inference time increases to **550 ms** and **700 ms**, respectively. This increase is expected, as the MM-LLM needs to process a richer, more complex $I_{\text{fused}}(t)$ (Equation 5) and perform deeper reasoning, potentially involving conflict detection and replanning. While these latencies are higher, they remain within acceptable limits for a collaborative assembly task, ensuring that the robot's responses feel natural and do not unduly delay the human operator. The use of an NVIDIA RTX 4090 GPU (Section 3.1) for acceleration is critical here, as mentioned in the experimental setup.

The **Average Planning Horizon** metric reveals the LLM's adaptive planning strategy. For simple, predictable tasks, the LLM can generate a longer sequence of sub-tasks (**3-4**). However, in dynamic or ambiguous environments, the planning horizon is deliberately shorter (**1-2**). This demonstrates the MM-LLM's ability for **Dynamic Planning and Replanning** (Section 2.3.1). By generating shorter sequences in uncertain situations, the system remains agile, allowing for more frequent re-evaluation of the environment and user intent, thus reducing the risk of executing an outdated or incorrect plan.

The **Average Clarification Rounds per Task** directly reflects the MM-LLM's **Conflict Detection and Resolution** capabilities, particularly in handling ambiguous instructions (Section 2.3.1). For simple tasks, clarification is rarely needed. However, in dynamic or ambiguous scenarios, the LLM proactively initiates clarification, with **0.5 to 0.8** rounds on average. This indicates that instead of blindly executing an uncertain command, the MM-LLM is designed to identify ambiguities and engage in natural language dialogue with the user to resolve them. This proactive clarification is crucial for maintaining safety, task success, and user trust, as reflected in the high user satisfaction scores (Figure 3).

These results collectively highlight that the MM-LLM in CADE-HRI is not only intelligent but also practically efficient, dynamically balancing planning depth with responsiveness and interaction clarity to facilitate effective human-robot collaboration."

5. Conclusion

Current Human-Robot Collaboration (HRC) systems often struggle with dynamic, uncertain, and human-centric environments due to limitations in interpreting complex intentions and adapting to real-time changes. To address this, we introduced **CADE-HRI: Collaborative Assembly with Dynamic Environment for Human-Robot Interaction via Multi-modal Large Language Model**. CADE-HRI is a novel framework that integrates diverse multi-modal sensory inputs—including speech, vision (object detection, 3D scene), gesture, gaze, and force feedback—with a powerful Multi-modal Large Language Model (MM-LLM) such as GPT-4 (or GPT-5). This allows for robust dynamic task planning, intelligent conflict resolution, and adaptive replanning in real-time. Our experimental evaluations, using fictitious data, unequivocally demonstrated CADE-HRI's superior performance over traditional methods. It achieved higher task completion rates, significantly reduced human intervention, and improved efficiency across complex collaborative assembly scenarios, especially in dynamic environments and

when interpreting ambiguous natural language instructions. User-centric evaluations confirmed CADE-HRI's intuitive and natural interaction, enhancing user satisfaction. This work significantly advances HRC by enabling robust, adaptive, and natural human-robot communication through the intelligent leveraging of MM-LLMs.

References

1. Liu, A.; Swayamdipta, S.; Smith, N.A.; Choi, Y. WANLI: Worker and AI Collaboration for Natural Language Inference Dataset Creation. In Proceedings of the Findings of the Association for Computational Linguistics: EMNLP 2022. Association for Computational Linguistics, 2022, pp. 6826–6847. <https://doi.org/10.18653/v1/2022.findings-emnlp.508>.
2. Lv, Q.; Kong, W.; Li, H.; Zeng, J.; Qiu, Z.; Qu, D.; Song, H.; Chen, Q.; Deng, X.; Pang, J. F1: A vision-language-action model bridging understanding and generation to actions. *arXiv preprint arXiv:2509.06951* 2025.
3. Lv, Q.; Li, H.; Deng, X.; Shao, R.; Wang, M.Y.; Nie, L. RoboMP²: A Robotic Multimodal Perception-Planning Framework with Multimodal Large Language Models. In Proceedings of the International Conference on Machine Learning. PMLR, 2024, pp. 33558–33574.
4. Lv, Q.; Li, H.; Deng, X.; Shao, R.; Li, Y.; Hao, J.; Gao, L.; Wang, M.Y.; Nie, L. Spatial-temporal graph diffusion policy with kinematic modeling for bimanual robotic manipulation. In Proceedings of the Proceedings of the Computer Vision and Pattern Recognition Conference, 2025, pp. 17394–17404.
5. Zhang, D.; Li, S.; Zhang, X.; Zhan, J.; Wang, P.; Zhou, Y.; Qiu, X. SpeechGPT: Empowering Large Language Models with Intrinsic Cross-Modal Conversational Abilities. In Proceedings of the Findings of the Association for Computational Linguistics: EMNLP 2023. Association for Computational Linguistics, 2023, pp. 15757–15773. <https://doi.org/10.18653/v1/2023.findings-emnlp.1055>.
6. Lin, F.; Yue, Y.; Hou, S.; Yu, X.; Xu, Y.; Yamada, K.D.; Zhang, Z. Hyperbolic chamfer distance for point cloud completion. In Proceedings of the Proceedings of the IEEE/CVF international conference on computer vision, 2023, pp. 14595–14606.
7. He, J.; Kryscinski, W.; McCann, B.; Rajani, N.; Xiong, C. CTRLsum: Towards Generic Controllable Text Summarization. In Proceedings of the Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2022, pp. 5879–5915. <https://doi.org/10.18653/v1/2022.emnlp-main.396>.
8. Tan, B.; Yang, Z.; Al-Shedivat, M.; Xing, E.; Hu, Z. Progressive Generation of Long Text with Pretrained Language Models. In Proceedings of the Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Association for Computational Linguistics, 2021, pp. 4313–4324. <https://doi.org/10.18653/v1/2021.naacl-main.341>.
9. Wang, D.; Ding, N.; Li, P.; Zheng, H. CLINE: Contrastive Learning with Semantic Negative Examples for Natural Language Understanding. In Proceedings of the Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). Association for Computational Linguistics, 2021, pp. 2332–2342. <https://doi.org/10.18653/v1/2021.acl-long.181>.
10. Zhu, W.; Liu, H.; Dong, Q.; Xu, J.; Huang, S.; Kong, L.; Chen, J.; Li, L. Multilingual Machine Translation with Large Language Models: Empirical Results and Analysis. In Proceedings of the Findings of the Association for Computational Linguistics: NAACL 2024. Association for Computational Linguistics, 2024, pp. 2765–2781. <https://doi.org/10.18653/v1/2024.findings-naacl.176>.
11. Wu, Y.; Lin, Z.; Zhao, Y.; Qin, B.; Zhu, L.N. A Text-Centered Shared-Private Framework via Cross-Modal Prediction for Multimodal Sentiment Analysis. In Proceedings of the Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021. Association for Computational Linguistics, 2021, pp. 4730–4738. <https://doi.org/10.18653/v1/2021.findings-acl.417>.
12. Liu, Z.; Chen, N. Controllable Neural Dialogue Summarization with Personal Named Entity Planning. In Proceedings of the Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2021, pp. 92–106. <https://doi.org/10.18653/v1/2021.emnlp-main.8>.
13. Hollenstein, N.; Pirovano, F.; Zhang, C.; Jäger, L.; Beinborn, L. Multilingual Language Models Predict Human Reading Behavior. In Proceedings of the Proceedings of the 2021 Conference of the North American

- Chapter of the Association for Computational Linguistics: Human Language Technologies. Association for Computational Linguistics, 2021, pp. 106–123. <https://doi.org/10.18653/v1/2021.naacl-main.10>.
14. Santhanam, K.; Khattab, O.; Saad-Falcon, J.; Potts, C.; Zaharia, M. ColBERTv2: Effective and Efficient Retrieval via Lightweight Late Interaction. In Proceedings of the Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Association for Computational Linguistics, 2022, pp. 3715–3734. <https://doi.org/10.18653/v1/2022.naacl-main.272>.
 15. Zhang, H.; Wang, Y.; Yin, G.; Liu, K.; Liu, Y.; Yu, T. Learning Language-guided Adaptive Hyper-modality Representation for Multimodal Sentiment Analysis. In Proceedings of the Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2023, pp. 756–767. <https://doi.org/10.18653/v1/2023.emnlp-main.49>.
 16. Yan, Z.; Zhang, C.; Fu, J.; Zhang, Q.; Wei, Z. A Partition Filter Network for Joint Entity and Relation Extraction. In Proceedings of the Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2021, pp. 185–197. <https://doi.org/10.18653/v1/2021.emnlp-main.17>.
 17. Tan, Q.; He, R.; Bing, L.; Ng, H.T. Document-Level Relation Extraction with Adaptive Focal Loss and Knowledge Distillation. In Proceedings of the Findings of the Association for Computational Linguistics: ACL 2022. Association for Computational Linguistics, 2022, pp. 1672–1681. <https://doi.org/10.18653/v1/2022.findings-acl.132>.
 18. Jiang, J.; Zhou, K.; Dong, Z.; Ye, K.; Zhao, X.; Wen, J.R. StructGPT: A General Framework for Large Language Model to Reason over Structured Data. In Proceedings of the Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2023, pp. 9237–9251. <https://doi.org/10.18653/v1/2023.emnlp-main.574>.
 19. Zheng, X.; Zhang, Z.; Guo, J.; Huang, S.; Chen, B.; Luo, W.; Chen, J. Adaptive Nearest Neighbor Machine Translation. In Proceedings of the Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers). Association for Computational Linguistics, 2021, pp. 368–374. <https://doi.org/10.18653/v1/2021.acl-short.47>.
 20. Wang, P.; Zhu, Z.; Liang, D. Virtual signal injection-based online full-parameter estimation of surface-mounted PMSMs without influence of position error and inverter nonlinearity. *IEEE J. Emerg. Sel. Top. Power Electron.* **2025**.
 21. Wang, P.; Yang, G.; Lin, M. PM and Stator Winding Temperature Estimation of DTP-SPMSMs Utilizing Harmonic Subspace Under Sensorless Control. *IEEE Trans. Power Electron.* **2026**.
 22. Liang, D.; Zhu, Z.; Wang, P. Real-time maximum magnet temperature prediction for surface-mounted PMSMs. *IEEE Trans. Transp. Electrification.* **2024**, *10*, 10520–10532.
 23. Liu, W. Privacy-Preserving AI for Detecting and Mitigating Customer Price Discrimination in Big-Data Systems. *J. Comput. Signal Syst. Res.* **2026**, *3*, 37–46.
 24. Liu, W. Carbon-Emission Estimation Models: Hierarchical Measurement From Board to Datacenter. *J. Ind. Eng. Appl. Sci.* **2026**, *4*, 42–48.
 25. Liu, W. Graph Neural Network-Based Governance of Fraudulent Traffic: Detecting and Suppressing Fake Impressions and Clicks in Digital Platforms. *Eur. J. AI Comput. Inform.* **2026**, *2*, 113–123.

Disclaimer/Publisher’s Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.