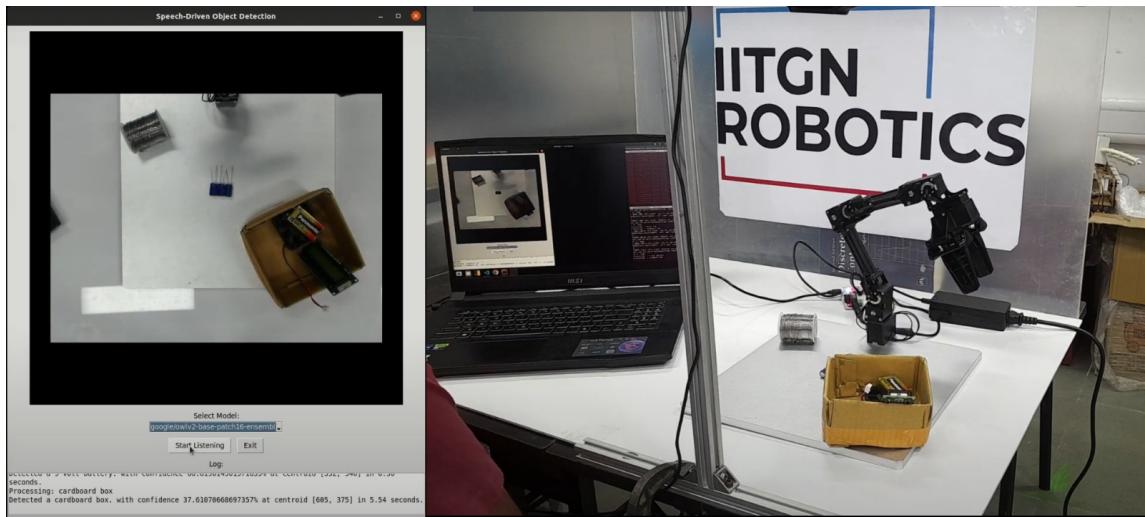


# 1 Manipulation using LLM and Modular CV Model Integration

2  
3 ANONYMOUS AUTHOR(S)  
4



23  
24 Fig. 1. OpenMANIPULATOR-X's workspace, where the experiments were performed (on the right) and the UI of the pipeline (on the  
25 left)

26 Achieving seamless robotic manipulation through voice commands is made possible by advancements in voice recognition, machine  
27 learning, and computer vision. A key strength lies in the modularity of the computer vision models, allowing for the integration of  
28 various zero-shot object recognition models tailored to specific tasks. This adaptability ensures the system can be easily updated or  
29 expanded as new models emerge. The process involves converting voice commands into text, extracting relevant objects using a  
30 Large Language Model, and determining the manipulator's joint angles via an Inverse Kinematics algorithm. This approach enhances  
31 scalability, flexibility, and future-proofing for advanced robotic manipulation.

32  
33 CCS Concepts: • Computing methodologies → Motion path planning; Vision for robotics; Information extraction; Object  
34 detection; Speech recognition; Machine translation.  
35

36 Additional Key Words and Phrases: Robotic Manipulation, Speech-To-Text, Large Language Models, Human-Robot Interaction,  
37 Zero-Shot Object Detection  
38

## 40 ACM Reference Format:

41 Anonymous Author(s). 2018. Manipulation using LLM and Modular CV Model Integration. In *Proceedings of Make sure to enter the*  
42 *correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 17 pages. <https://doi.org/XXXXXXX.XXXXXXX>  
43  
44

45 Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not  
46 made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components  
47 of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on  
48 servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.  
49

50 © 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
51 Manuscript submitted to ACM  
52

## 53    1 Introduction

54    The quest for seamless human-machine interaction has witnessed substantial progress in recent years, driven by  
 55    advancements in voice recognition technologies, machine learning, Large Language Models (LLMs) and computer  
 56    vision. Within the domain of robotic manipulation, enabling a robot to execute complex tasks through simple voice  
 57    commands represents the convergence of these technologies and epitomizes their potential for real-world applications.  
 58

59    A central strength of the proposed system lies in its modular integration of computer vision (CV) models. This  
 60    modularity facilitates the incorporation of diverse zero-shot object detection models, each optimized for specific object  
 61    types or tasks, thereby enhancing the system's flexibility and scalability. Such an approach ensures that the pipeline  
 62    can adapt to various environments and application domains, addressing the limitations of monolithic vision systems.  
 63

64    To demonstrate the practical utility of this research, we propose a versatile application: a robotic assistant designed  
 65    for electrical laboratories and machine workshops. This assistant is capable of performing manipulation tasks on tools  
 66    and components within its configuration space (C-space) based on user voice commands. The manipulator used in this  
 67    study, the OpenMANIPULATOR-X, is a 4-DOF robotic arm capable of precise object manipulation.  
 68

### 69    1.1 Contributions

70    This paper introduces a novel modular pipeline that integrates an LLM, zero-shot object detection models, and inverse  
 71    kinematics for robotic manipulation. The primary contributions of this work include the demonstration of modular  
 72    computer vision (CV) model integration, enabling versatile and scalable object detection across various tasks. Addi-  
 73    tionally, the pipeline is validated through practical experiments conducted in diverse environments, showcasing its  
 74    adaptability and effectiveness in real-world scenarios.  
 75

### 76    1.2 Paper Structure

77    The remainder of this paper is organized as follows:  
 78

- 79    • **Section 2: Methodology** details the components of the proposed system, including speech recognition, LLM  
     80    integration, zero-shot object detection models, inverse kinematics, and text-to-speech processing.
- 81    • **Section 3: Experimental Validation** describes the experimental setup, datasets used, and results obtained across  
     82    different scenarios, such as manipulating tools and components in laboratory environments.
- 83    • **Section 4: Importance of Modularity** emphasizes the advantages of using modular CV models and their impact  
     84    on the system's adaptability and robustness.
- 85    • **Section 5: Challenges Faced** discusses the technical and practical difficulties encountered during the development  
     86    and deployment of the pipeline.
- 87    • **Section 6: Future Scope and Discussion** explores potential enhancements to the system and broader applications  
     88    of the proposed approach.
- 89    • **Section 7: Conclusion** summarizes the findings and implications of the research, highlighting its contributions  
     90    and future prospects.

## 91    2 Methodology

92    We can easily curate the sequence of steps we must take to ensure that the manipulation task is accomplished by a  
 93    mere voice command from the user. So, that directly leads us to the voice recognition part of the pipeline. Once, we  
 94    recognize the voice from the user and convert it into "speech\_text", we feed the same to the LLM model. From which  
 95    Manuscript submitted to ACM  
 96

we would get only specific parts of the sentence which is required for the Zero-shot object recognition CV models (in our case objects present in the sentence). The CV models process the LLM's input and gives us the required coordinates to the manipulator (to perform pick-and-place). We would need an IK algorithm to solve the given coordinates and find the joint angles of the manipulator (here 4 joint angles), so that the manipulation can finally be performed successfully. The above sequence of steps can be simply depicted through the below flowchart.

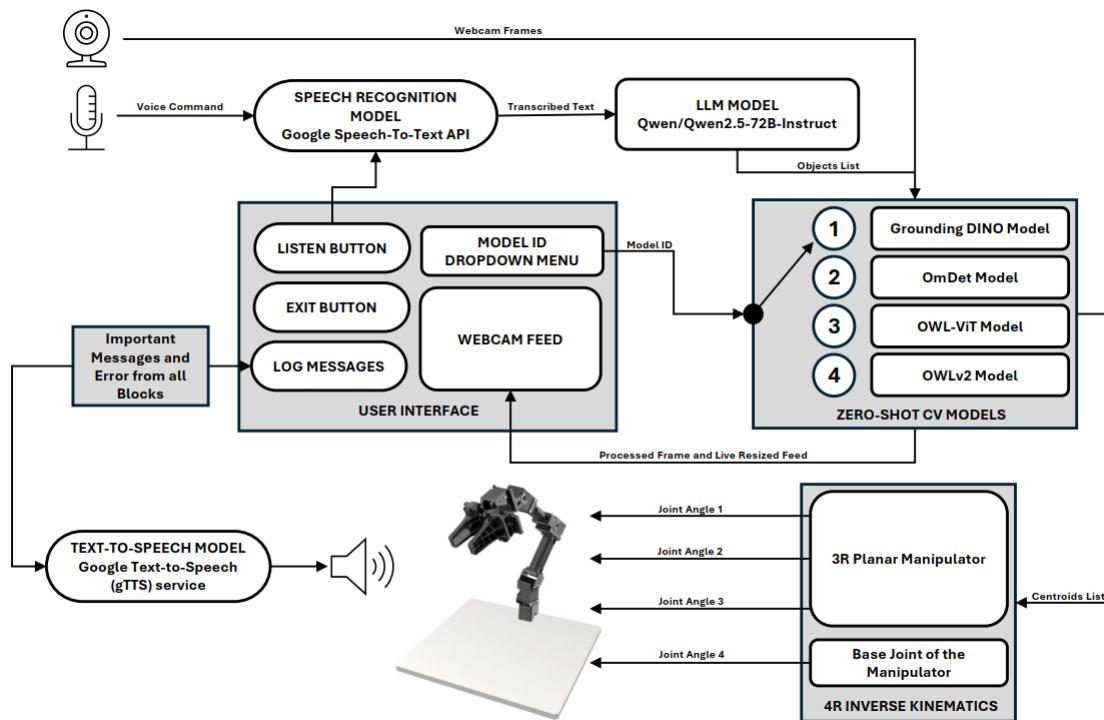


Fig. 2. Flowchart depicting all the blocks present in the pipeline.

From the flowchart we can easily relate to the pseudocode below. The pseudocode is given below for better understanding of the following subsections.

---

**Algorithm 1:** Robotic Manipulation Pipeline

---

**Input:** Speech recognition, LLM, modular CV models, and inverse kinematics modules are initialized.

157 Initialize an empty list `detected_centroids` to store object centroids;

158

159 **Step 1: Initialization;**

160     Announce readiness using text-to-speech (TTS);

161     Prompt the user to select the required CV model;

162     Retrieve the selected model ID using modular CV and LLM modules;

163     Initialize the speech recognizer;

164

165     **while** *True* **do**

166

167         **Step 2: Speech Recognition;**

168         Prompt user to speak via TTS;

169         Capture audio input using a microphone;

170         **if** *Speech is successfully transcribed* **then**

171             Announce successful recognition via TTS;

172         **else**

173             Announce error and prompt the user to try again;

174             **continue** to the next iteration;

175

176         **Step 3: Text Processing with LLM;**

177         Process the `speech_text` to generate actionable commands;

178         **if** *No actionable commands are detected* **then**

179             Announce shutdown via TTS;

180             Terminate the program and **break** from the loop;

181

182         **Step 4: Object Detection Using Zero-shot CV Models;**

183         **foreach** *object* **in** the `ordered_response_list` **from** LLM **do**

184

185             Prepend "a" to the `object_text` for query formatting;

186             Perform object detection using the selected CV model;

187             Retrieve centroid, confidence, `processing_time`, and `processed_frame`;

188             Display the `processed_frame` temporarily;

189             Append the detected centroid to `detected_centroids`;

190

191         **Step 5: Inverse Kinematics and Pick-and-Place;**

192         Transform the pixel coordinates of the first and second detected centroids into world coordinates;

193         Perform pick-and-place operation using the transformed coordinates;

194         Generate a completion message using the LLM and announce it via TTS;

195         **if** *An error occurs during object detection or manipulation* **then**

196             Announce the error via TTS;

197             **continue** to the next iteration;

198

199

200

201

202

203

204

205

---

**209 2.1 Speech Recognition Model**

210 Speech recognition models convert audio input, such as sounds from a microphone or pre-recorded files, into text  
211 by processing audio signals into meaningful features. They use neural networks to match sounds to words (acoustic  
212 models), ensure grammatical correctness with language models, and output the final transcribed text for further use.  
213 Some of Speech Recognition Models are DeepSpeech by Mozilla, Whisper by OpenAI and Google Speech-to-Text API.  
214

215 For our task, we have used Google Speech-To-Text API, which is capable of capturing live audio and processing it in  
216 real-time. It supports over 100 languages, including accents like English(India). While performing the task, the language  
217 was chosen as "en-IN" to account for the Indian English accent.  
218

**219 2.2 Large Language Model - "Qwen/Qwen2.5-72B-Instruct"**

220 Large Language Models (LLMs) are advanced artificial intelligence systems designed to comprehend, generate, and  
221 interact with human language. These models, such as OpenAI's GPT-3 and Google's BERT, are trained on extensive  
222 datasets comprising diverse text, enabling them to grasp context, answer questions, and produce human-like text. Their  
223 versatility makes them suitable for a range of tasks, including translation, summarization, query answering, and content  
224 generation.  
225

226 For this work, we utilize the LLM "Qwen2.5-72B-Instruct," developed by Alibaba Cloud's Qwen team. This model,  
227 part of the Qwen LLM family, is characterized by its architecture of 72.7 billion parameters, which includes 2.7 billion  
228 embedding parameters and an 80-layer deep neural network [8]. The model is deployed in two key scenarios to facilitate  
229 the robotic manipulation task:  
230

- 231 (1) To extract and structure object information from user queries. For instance, given the input query speech\_text:  
232     *"Place the screwdriver on top of the circular plate,"* the model processes the request and returns the objects and  
233     their properties in the order of appearance as a comma-separated list: `["screwdriver", "circular plate"]`.  
234 (2) To generate a grammatically correct confirmation statement. For example, for the query speech\_text: *"Place*  
235     *the screwdriver on top of the circular plate,"* the model produces the statement: *"You have placed the object as*  
236     *requested."*

237 To integrate this model into our system, we utilize OpenAI's openai module. The deployment process involves creat-  
238 ing an API key from Hugging Face and configuring the LLM to respond accurately to the specified tasks. The LLM's output  
239 in scenario (1) is provided to the zero-shot object detection CV models as a structured list - `"ordered_response_list"`,  
240 enabling precise detection and manipulation of objects. Once the manipulation task is completed, the output of scenario  
241 (2) is processed by the text-to-speech module to provide an interactive and natural user experience.  
242

**243 2.3 Zero-shot Object Detection CV Models**

244 Zero-shot object detection (ZSD) is a computer vision task where models detect objects in images without prior training  
245 on those specific object classes. Given an image and a list of candidate classes, ZSD models output bounding boxes and  
246 labels for detected objects. This approach eliminates the need for extensive data annotation and model fine-tuning,  
247 making it efficient for applications like image search, object counting, and tracking. In our pipeline we use four different  
248 ZSD models, which obtain `"ordered_response_list"` from the LLM and they process the objects present in the list  
249 to draw the bounding boxes with certain confidence level or score to eventually calculate the centroid of the specific  
250 object asked for.

All of these models follow the same sequence of steps to return us the centroids of the objects and all of them recognize objects better when only a single object is given to it to detect. Hence, define a function which takes "model\_id", "object\_text", "frame\_end" and "camera\_index" as input parameters and finally returns "centroid", "confidence", "processing\_time" and "processed\_frame" to the user interface.

**2.3.1 Grounding DINO Model (tiny variant).** The Grounding DINO model integrates the DINO architecture with grounded pre-training to enable open-set object detection. Developed by IDEA-Research, it extends traditional closed-set object detection models by incorporating a text encoder, allowing it to detect objects without prior training on specific classes [4]. The tiny variant of the model achieves notable performance, including an average precision (AP) of 52.5 on the COCO (Common Objects in Context) zero-shot benchmark [5]. This capability makes it particularly effective for diverse applications such as image search and object counting. Its robust generalization across domains without requiring additional training data establishes it as a highly capable tool for zero-shot object detection tasks.

The deployment of this model in Python leverages Hugging Face's `transformers` library, specifically utilizing the `AutoProcessor` and `AutoModelForZeroShotObjectDetection` modules. Initially, input image frames captured using OpenCV are converted from BGR to RGB format, as OpenCV processes images in BGR format, whereas most computer vision models, including Grounding DINO, require images in RGB format. After this conversion, the preprocessed image and the corresponding `object_text` query are fed into the processor, which outputs detection results. These results include bounding boxes around detected objects and their respective confidence scores.

To identify the object with the highest detection confidence, a straightforward algorithm iterates through the bounding boxes to locate the one with the highest confidence score. The bounding box with the highest confidence is visually highlighted in blue to distinguish it from other detections.

The centroid of the detected object is computed to assist in robotic manipulation. This is achieved by calculating the geometric center of the bounding box using coordinate geometry. The centroid plays a crucial role in ensuring that the robotic manipulator can effectively grasp the object while maintaining balance during manipulation tasks.

**2.3.2 OmDet Model.** The OmDet model, introduced in the work "*Real-time Transformer-based Open-Vocabulary Detection with Efficient Fusion Head*" by Tiancheng Zhao et al., is a state-of-the-art approach for zero-shot (open-vocabulary) object detection. The model leverages the Swin-Tiny-HF architecture, which features a hierarchical structure and shifted windows to reduce computational complexity. This architectural design enables OmDet to deliver efficient and scalable performance across a range of computer vision tasks. The model is particularly effective in detecting objects without predefined labels, ensuring adaptability to diverse applications. Furthermore, its efficient design, coupled with the Swin-Tiny-HF architecture, provides high performance and fast inference, making it well-suited for real-time scenarios where both speed and accuracy are paramount [9].

The deployment of the OmDet model follows a process similar to that of the Grounding DINO model, with the primary difference being the use of the `OmDetTurboForObjectDetection` module to initialize the model. Once deployed, the process for obtaining results, including the highest confidence score, centroids, bounding boxes, and frame outputs, remains identical to the methodology used for the DINO model.

**2.3.3 OWL-ViT Model.** The OWL-ViT (Vision Transformer for Open-World Localization) is a zero-shot, text-conditioned object detection model proposed by Matthias Minderer et al. The model leverages CLIP (Contrastive Language-Image Pre-Training) as its multi-modal backbone, integrating a ViT (Vision Transformer) for extracting visual features and a causal language model for processing text features. Unlike traditional vision transformers, OWL-ViT removes the final

313 token pooling layer and attaches lightweight classification and bounding box heads to each transformer output token,  
 314 enabling efficient and precise object detection [7].  
 315

316 The term *text-conditioned* refers to the model's ability to utilize text queries to guide the object detection process,  
 317 effectively linking textual descriptions with visual representations. This capability is achieved through the use of  
 318 contrastive loss, a training strategy that maximizes the similarity between paired image and text embeddings. OWL-  
 319 ViT's architecture supports open-vocabulary classification, allowing it to detect objects based on textual descriptions  
 320 without being restricted to predefined labels. The model is trained from scratch and further fine-tuned on detection  
 321 datasets, enhancing its generalization and detection accuracy [7].  
 322

323 The deployment of the OWL-ViT model follows a process similar to that of the Grounding DINO and OmDet models.  
 324 However, for initialization, we utilize the OwlViTProcessor and OwlViTForObjectDetection modules from Hugging  
 325 Face's transformers library. Once deployed, the procedures for obtaining detection results, such as bounding boxes,  
 326 centroids, and confidence scores, remain consistent with the previously described methods.  
 327

328  
 329  
 330 2.3.4 *OWLv2 Model.* The OWLv2 model, short for Open-World Localization version 2, is a zero-shot, text-conditioned  
 331 object detection model proposed in "*Scaling Open-Vocabulary Object Detection*" by Matthias Minderer, Alexey Gritsenko,  
 332 and Neil Houlsby. Building on its predecessor, OWL-ViT, OWLv2 leverages CLIP as its multi-modal backbone, integrating  
 333 a ViT (Vision Transformer) for visual feature extraction and a causal language model for text processing. Unlike OWL-  
 334 ViT, OWLv2 employs the ViT-B/16 Transformer architecture, which captures finer details compared to the ViT-B/32  
 335 used in OWL-ViT [6].  
 336

337 Similar to its predecessor, OWLv2 removes the final token pooling layer and attaches lightweight classification and  
 338 bounding box heads. It achieves open-vocabulary classification by replacing fixed classification layer weights with  
 339 class-name embeddings derived from the text model [6]. OWLv2 introduces improved scalability, enhanced real-time  
 340 performance, and better generalization to unseen objects, making it a significant step forward in open-vocabulary  
 341 detection tasks [6]. These advancements are attributed to its refined training process and enhanced text processing  
 342 pipeline, which together ensure superior efficiency and adaptability for diverse applications requiring zero-shot object  
 343 detection.  
 344

345 While OWLv2 demonstrates notable improvements over OWL-ViT, its deployment process remains unchanged.  
 346 The model is initialized using the Owlv2Processor and Owlv2ForObjectDetection modules from Hugging Face's  
 347 transformers library. Once deployed, the methods for obtaining detection results, including bounding boxes, centroids,  
 348 and confidence scores, are identical to those used for its predecessors.  
 349

## 353 2.4 Inverse Kinematics

354 Inverse kinematics (IK) is a fundamental concept in robotics, used to compute the joint configurations required  
 355 for a manipulator to achieve a specified end-effector position and orientation. Unlike forward kinematics, which  
 356 determines the end-effector position based on given joint angles, IK operates in reverse by specifying the target  
 357 location and calculating the necessary joint angles to reach it. This process often involves solving complex, nonlinear  
 358 equations, frequently requiring iterative techniques such as Jacobian-based methods. IK is indispensable for robotic  
 359 arm manipulation, enabling precise motion and positioning for applications in domains such as manufacturing, medical  
 360 robotics, and autonomous systems.  
 361

365 In practice, IK for robots can be implemented using the MoveIt! framework by providing the Unified Robot De-  
 366 scription Format (URDF) file to the MoveIt! setup assistant. Once configured, the generated package includes con-  
 367 troller launch files that enable position control for the robot using the default IK solver (`KDLKinematicsPlugin`).  
 368 However, in our case, the MoveIt! package's IK solver proved to be inefficient, exhibiting significant delays and  
 369 frequently moving into singular configurations even when the target position was easily reachable by the manipulator  
 370 (OpenMANIPULATOR-X).  
 371



400 Fig. 3. OpenMANIPULATOR-X, a 4-DOF manipulator with 4R configuration. It can be assumed as a 3R planar manipulator affixed on  
 401 top of base revolute joint.

402  
 403 To address these limitations, we developed a custom IK algorithm tailored to the OpenMANIPULATOR-X, a 4-DOF  
 404 manipulator (excluding the 1-DOF gripper). The manipulator is modeled as a combination of a 3R planar manipulator (3-  
 405 DOF) mounted on a revolute joint (1-DOF). This approach enabled us to achieve more efficient and reliable computa-  
 406 tion of joint angles, ensuring precise and stable performance for the manipulation tasks.  
 407

408  
 409 2.4.1 *IK of a 3R Planar Manipulator.* The inverse kinematics (IK) equations for a 3R planar manipulator can be derived  
 410 through geometric and trigonometric principles and are commonly detailed in standard robotics textbooks. These  
 411 equations can also be obtained by applying the laws of vector addition and trigonometric relationships to the diagram  
 412 shown in Figure 4.  
 413

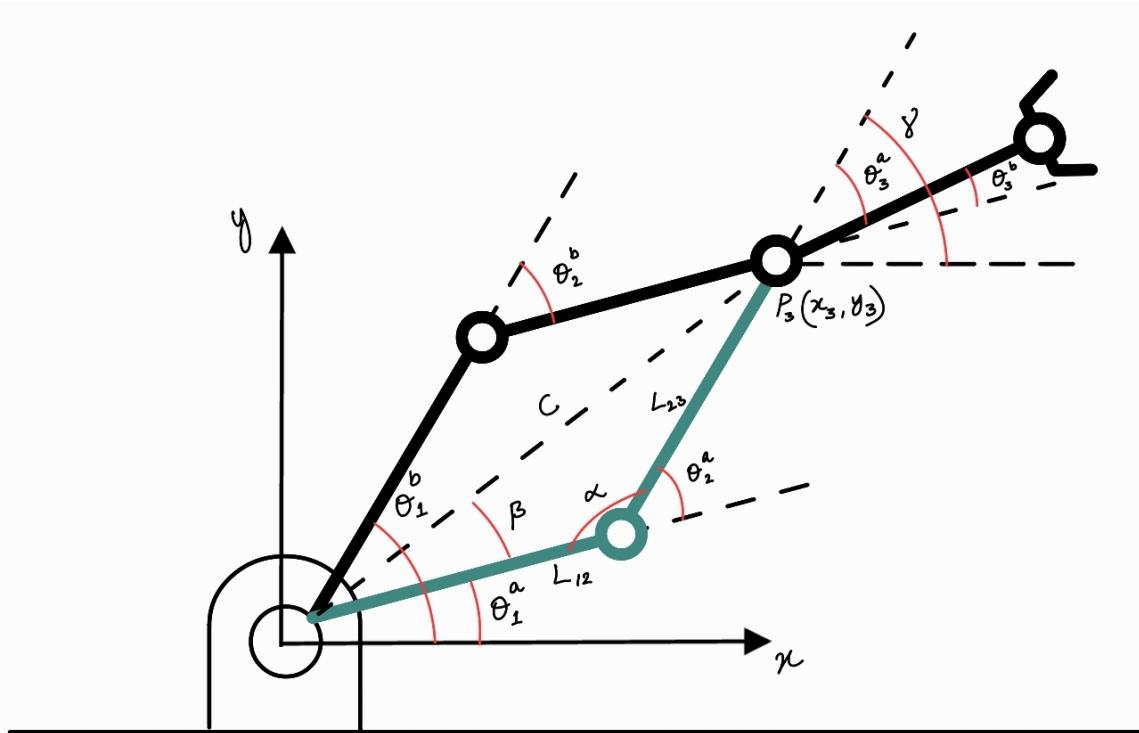
414 Using the triangular law of vector addition and the law of sines, the following equations can be derived:  
 415

417  
418  
419  
420

$$\alpha = \cos^{-1} \left( \frac{x_3^2 + y_3^2 - L_{12}^2 - L_{23}^2}{2L_{12}L_{23}} \right), \quad (1)$$

421  
422  
423  
424

$$\beta = \sin^{-1} \left( \frac{L_{23} \sin \alpha}{\sqrt{x_3^2 + y_3^2}} \right). \quad (2)$$



453 Fig. 4. 3R Planar manipulator diagram with joint angles and link lengths and the two possible orientations for a single position are  
454 differentiated using colors.  
455

456 From the geometry in Figure 4, the position of point  $P_3$  in Cartesian coordinates can be expressed as shown below  
457 [3].

458

$$P_3 = \begin{bmatrix} x_3 \\ y_3 \end{bmatrix} = \begin{bmatrix} x_e - L_{34} \cos(\gamma) \\ y_e - L_{34} \sin(\gamma) \end{bmatrix}. \quad (3)$$

461 Using these relationships, the joint angles can be determined. For most configurations, the manipulator can assume  
462 two orientations: elbow-up (a) or elbow-down (b) [3]. For the purposes of this study, which focus on pick-and-place  
463 tasks, the elbow-down configuration is selected exclusively. The joint angles for this configuration are calculated as  
464 follows:

466  
467

$$\theta_1^b = \tan^{-1} \left( \frac{y_3}{x_3} \right) + \beta, \quad (4)$$

$$\theta_2^b = -(180^\circ - \alpha), \quad (5)$$

$$\theta_3^b = \gamma - \theta_1^b - \theta_2^b. \quad (6)$$

**2.4.2 IK of the Base Joint.** The base joint angle,  $\theta_{\text{base}}$ , can be determined using the coordinates of the end-effector ( $x_e, y_e$ ). The angle is given by:

$$\theta_{\text{base}} = \tan^{-1} \left( \frac{y_e}{x_e} \right). \quad (7)$$

Equations (4), (5), (6), and (7) are integrated to compute all the joint angles required for the OpenMANIPULATOR-X manipulator.

**2.4.3 Coordinate Transformation.** To compute the joint angles, the end-effector coordinates must be transformed from the camera's base frame to the manipulator's base frame, and the pixel coordinates must be converted to meters. This transformation is performed using the following equations:

$$x_m = (y_p - 15) \cdot 0.000625, \quad y_m = (x_p - 435) \cdot 0.000625, \quad (8)$$

where  $x_m$  and  $y_m$  are the coordinates in meters,  $x_p$  and  $y_p$  are the pixel coordinates in the camera frame. The constants 15 and 435 represent the reference pixel coordinates, and the scaling factor 0.000625 converts pixel units to meters. With the transformed coordinates, the joint angles for any position within the manipulator's configuration space (C-space) are computed. Joint control is then performed, which provides faster and more accurate positioning compared to the MoveIt! package's default position control.

## 2.5 Text-to-Speech Model

To provide auditory feedback to the user, we utilized the Google Text-to-Speech (gTTS) service, which converts textual input into speech. The gTTS service supports over 100 languages and accents, offering flexibility for diverse applications. It also allows customization of speech speed through the `slow` and `speed_factor` parameters, enabling control over the pace of the generated speech.

In this study, the text-to-speech model is employed to deliver feedback such as "Speech Recognized!" upon successful speech recognition. Additionally, it provides error messages for scenarios where recognition fails, such as "Could not understand the audio. Please try again!" This integration ensures a user-friendly and interactive system, enhancing the overall experience.

### 3 Experimental Validation

### 3.1 Task Description

We test the working of our pipeline by placing the manipulator in three different scenarios:

- colored blocks of different shapes
  - electrical lab environment
  - mechanical lab environment

As mentioned earlier we give the manipulator natural voice commands and it assists you in the lab by completing the assigned task through manipulation (pick-and-place).

This video [2] (watch in 1.75x) shows the performance of the manipulator under different scenarios and different CV models.

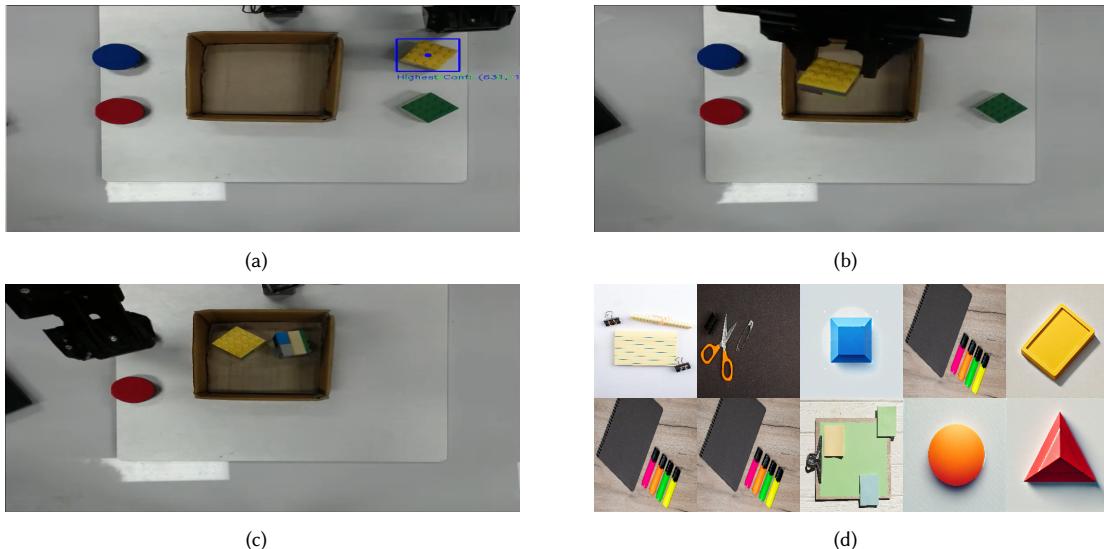
521 **3.2 Dataset Used**

522 It is also well recommended to perform benchmark tests on these CV models to get a clear idea on how well or worse it  
 523 performs under these different scenarios. So, we create a sample dataset of 50 images, which consists of 20 images of  
 524 electrical or electronic components and tools, 20 images of mechanical tools and 10 images of simple shapes of different  
 525 colors and popular stationary. We feed the sample dataset to all of these Zero-shot object detection CV models.  
 526

527 The data accumulation of these CV models can be seen in through the Excel worsheet [1].  
 528

529 **3.3 Experimental Results**

530 *3.3.1 Colored Blocks of Different Shapes.* In the video [2] (watch in 1.75x) we ask the robot to place the yellow square  
 531 object inside the cardboard box followed by asking it to place the green square object, the blue circular object and the  
 532 red circular object inside the cardboard box. We clearly see it perform the manipulations seamlessly without any issues.  
 533



556 Fig. 5. (a) The Robot is searching for a yellow square. (b) The robot places the Yellow Square in the cardboard box. (c) The robot places  
 557 the red square in the cardboard box. (d) A sample image dataset consisting of 20 images of shapes, colors, and stationery items used to  
 558 benchmark the model, via free licensed images (<https://www.freepik.com/>), (<https://www.vecteezy.com/>), and (<https://unsplash.com/>).  
 559  
 560

561 Table 1. Testing Models on Shapes, Colours, and Stationery Items.  
 562

563 <b>Models</b>	564 <b>Accuracy</b>	565 <b>Average Processing Time (s)</b>	566 <b>Mean Confidence (Correct Predictions)</b>
567 <i>Grounding DINO</i>	70.00%	10.489	80.81%
568 <i>Omdet Turbo</i>	80.00%	2.347	68.15%
569 <i>OwlV2</i>	90.00%	15.247	60.43%
570 <i>OwlVit</i>	80.00%	2.157	34.18%

573 After testing the four models on different shapes and colors, *OwlV2* was the most accurate. *OwlVit* emerged to be the  
 574 fastest model, and *Grounding DINO* proved to be the most confident model in predictions.  
 575

576 3.3.2 *Electrical Lab Environment*. In the video [2] (watch in 1.75x) we ask the robot to place all the electrical or electronic  
 577 tools and components in a specific order as shown in the video, inside the cardboard box. We clearly see that it performs  
 578 the manipulations seamlessly without any issues.  
 579



600 Fig. 6. (a) The Robot is searching for a 9V Battery, (b) The Robot is placing a Blue Capacitor in the Cardboard Box, (c) The Robot  
 601 has successfully placed a green LCD display in the Cardboard Box, (d) A sample Image dataset consisting of 20 images of Various  
 602 Electrical Tools and Components used to perform Benchmark Tests on the Computer Vision Models, via free licensed images  
 603 (<https://www.freepik.com/>), (<https://www.vecteezy.com/>), and (<https://unsplash.com/>).  
 604

605 Table 2. Comparison of models for detecting Electrical Tools.  
 606

611 Models	612 Accuracy	613 Average Processing Time (s)	614 Mean Confidence (Correct Predictions)
613 <i>Grounding DINO</i>	614 70.00%	615 13.524	616 80.72%
615 <i>Omdet Turbo</i>	616 95.00%	617 2.158	618 68.10%
617 <i>OwlV2</i>	618 95.00%	619 17.008	620 49.70%
619 <i>OwlVit</i>	620 75.00%	621 1.145	622 28.73%

623 While testing on Electrical Tools, *Omdet Turbo* and *OwlV2* showed the highest accuracy. *OwlVit* was the most  
 624 accurate model, and *Grounding DINO* was the most confident model for predictions.

625     3.3.3 *Mechanical Lab Environment.* In the video [2] (watch in 1.75x) we ask the robot to place all the mechanical tools  
 626     in clockwise order as shown in the video, inside the cardboard box. We clearly see that it performs the manipulations  
 627     seamlessly without any issues.  
 628



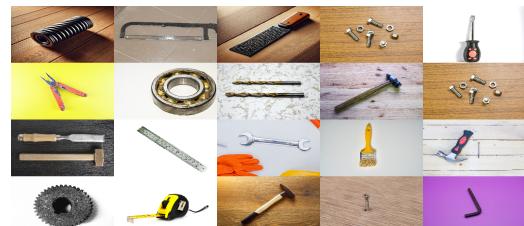
(a)



(b)



(c)



(d)

651     Fig. 7. (a) The model detects a green gear to be placed in the cardboard box. (b) The robot has placed the green gear in the cardboard  
 652     box. (c) The robot has placed the black chisel in the cardboard box. (d) A sample image dataset of 20 images of various mechanical  
 653     tools used to benchmark computer vision models, via free licensed images (<https://www.freepik.com/>), (<https://www.vecteezy.com/>),  
 654     and (<https://unsplash.com/>).

Table 3. Comparison of models to detect Mechanical Tools.

Models	Accuracy	Average Processing Time (s)	Mean Confidence (Correct Predictions)
<i>Grounding DINO</i>	80.00%	14.076	81.25%
<i>Omdet Turbo</i>	75.00%	2.372	60.34%
<i>OwlV2</i>	80.00%	14.883	47.04%
<i>OwlVit</i>	65.00%	1.175	28.44%

673     While testing over Mechanical Tools, *Grounding DINO* and *OwlV2* exhibited the maximum accuracy. *OwlV2* was the  
 674     fastest model and *Grounding DINO* was the most confident model for predictions.  
 675

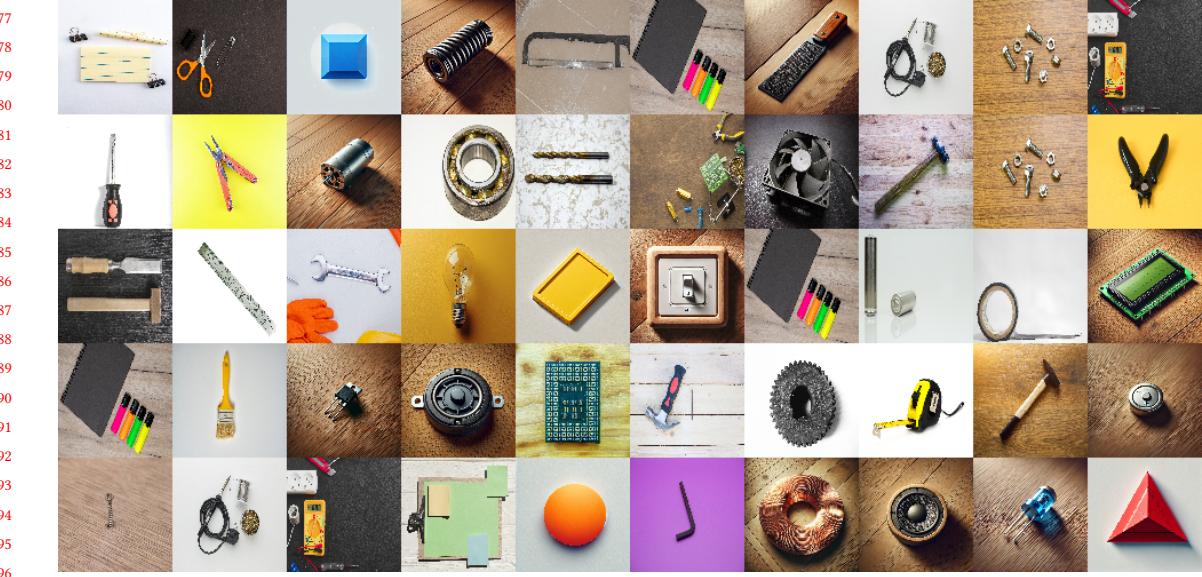


Fig. 8. A sample image dataset consisting of 50 images used to perform Benchmark Tests on the Computer Vision Models to judge the overall performance of the Models, via free licensed images (<https://www.freepik.com/>), (<https://www.vecteezy.com/>), and (<https://unsplash.com/>).

Table 4. Overall comparison of models.

Models	Accuracy	Average Processing Time (s)	Mean Confidence (Correct Predictions)	Confidence-Weighted Accuracy
<i>Grounding DINO</i>	74.00%	13.138	80.97%	59.92%
<i>Omdet Turbo</i>	84.00%	2.282	65.34%	54.88%
<i>OwlV2</i>	88.00%	15.806	50.93%	44.81%
<i>OwlVit</i>	72.00%	1.360	29.84%	21.48%

3.3.4 *Overall Performance of Models.* Overall, *OwlV2* was the most accurate model. *OwlVit* was the fastest model, and *Grounding Dino* was the most confident model.

For the highest accuracy, *OwlV2* is preferred, albeit at the cost of speed and confidence. This is because it uses a deeper and more complex architecture, which increases accuracy but reduces speed. For good confidence and accuracy, *Grounding DINO* excels, making it reliable but slightly slower. Its use of additional attention mechanisms or object-centric processing explains its high confidence but also incurs a speed penalty. The fastest model was *OwlVit*, but it consistently performed poorly in accuracy and confidence, making it ideal only for high-speed applications. This is due to its relatively lighter architecture. *Omdet Turbo* was a well-balanced model, as it uses a moderately heavy architecture, making it suitable for tasks that require a balance between speed and accuracy.

#### 729   **4 Why Modularity of CV Models is Important?**

730   The modularity of computer vision (CV) models is crucial for advancing the development and application of Robot-LLM  
731   systems that integrate language models with visual perception. Modular CV models offer flexibility, scalability, and  
732   adaptability, which are essential in addressing the diverse challenges faced in robotics and multi-modal AI systems.  
733

- 734   (1) Modularity facilitates customization and reuse. By designing CV systems as interchangeable components,  
735   developers can easily tailor specific modules for tasks like object detection, scene segmentation, or gesture  
736   recognition without rebuilding the entire pipeline. This approach enables researchers to reuse pre-trained  
737   modules in different contexts, reducing development time and computational costs.  
738
- 739   (2) Modularity enhances scalability. Complex robotics applications often require incremental updates or additions  
740   to their vision capabilities. A modular architecture allows seamless integration of new functionalities. This  
741   scalability is critical in dynamic environments where robotic systems need to adapt to evolving requirements.  
742
- 743   (3) Modularity supports fault isolation and debugging. In large-scale systems, identifying the source of errors  
744   can be challenging. Modular CV architectures enable developers to isolate and diagnose specific components,  
745   ensuring more efficient troubleshooting and maintenance.  
746

747   If in future, we need to add a new, better zero-shot object detection CV model to our pipeline we just need to add  
748   or change only the "process\_frame" function. Once, that is done we are good to go. This makes our implementation  
749   modular, which is crucial for all the reasons given above. Finally, modularity promotes collaboration across research  
750   domains. Teams can independently develop and improve modules, fostering innovation and rapid prototyping.  
751

#### 752   **5 Challenges Faced**

- 753   • Finding the right LLM which is simple to deploy and works everytime for the same repetitive instruction it has  
754   been given.  
755   • Writing modules for every block in the pipeline (efficient IK, Modular CV, pick-and-place, speech recognition,  
756   user interface, text-to-speech) which would be finally imported in the UI code.  
757   • Overlay of live webcam feed and the display of  
758   "processed\_frame" together in the same window of the user interface.  
759   • Devising a single common function to detect objects and return centroids for all of the zero-shot object detection  
760   models.  
761   • Figuring out to how to divide the inverse kinematics problem among the parts of the OpenMANIPULATOR-X  
762   manipulator, to solve it efficiently and easily.  
763

#### 764   **6 Future Scope and Discussion**

765   The integration of robotics and Large Language Models (LLMs) presents significant potential, given that both fields are  
766   among the fastest-evolving domains in modern technology. The rapid adoption of systems like OpenAI's ChatGPT into  
767   everyday life underscores the transformative impact of LLMs. Extending this capability to interact effectively with  
768   the physical environment through robotic systems offers unparalleled opportunities for societal impact. A technology  
769   pipeline, such as the one proposed in this study, has the potential to bridge the gap between digital intelligence and  
770   real-world application, surpassing the current limitations of LLMs in generating textual content.  
771

772   One of the most impactful applications of such a pipeline is in assistive technologies, particularly for individuals with  
773   disabilities. For instance, the world's blind population could benefit significantly from a system capable of interpreting  
774

781 visual inputs and executing contextual physical tasks. Beyond the blind community, individuals with varying disabilities  
 782 could use this cross-domain innovation to perform everyday activities more efficiently. Our experimentation, which  
 783 employed a standard webcam and a 4-DOF manipulator, highlights the versatility of such systems. The addition of a  
 784 depth-sensing camera, for example, could allow manipulators to perform precise tasks like placing tools directly into a  
 785 user's hands.  
 786

787 The modular design of the pipeline ensures that it remains future-proof. Faster, more accurate CV or LLM modules can  
 788 be seamlessly integrated without requiring a complete redesign. Another exciting avenue lies in industrial automation,  
 789 where modular CV-LLM systems could optimize assembly lines by integrating complex decision-making capabilities  
 790 with robotic dexterity. These advancements underscore the transformative potential of modular, cross-domain pipelines  
 791 in addressing real-world challenges across industries.  
 792

## 793 7 Conclusion

794 The integration of robotics and Large Language Models (LLMs) with modular computer vision (CV) architectures  
 795 represents a transformative step in advancing multi-modal AI systems. Our research highlights the importance of  
 796 modularity in ensuring adaptability, scalability, and interoperability, allowing the pipeline to evolve with advancements  
 797 in both CV and LLM technologies. Through experimentation with a simple setup, we demonstrated the pipeline's  
 798 potential in real-world applications, from assistive technologies for individuals with disabilities to broader industrial  
 799 automation tasks. The modular design ensures future readiness, enabling seamless integration of more accurate and  
 800 efficient models to meet emerging challenges.  
 801

802 This innovative cross-domain approach bridges the gap between digital intelligence and physical interaction, paving  
 803 the way for a wide range of use cases. Ultimately, this research underscores the immense potential of combining LLMs  
 804 with robotics, offering a pathway toward smarter, more inclusive, and highly adaptable systems for diverse real-world  
 805 environments.  
 806

## 807 8 References

- 808 [1] Anonymized. 2024. Data Accumulation of CV Models. Excel Sheet. <https://docs.google.com/spreadsheets/d/1LGY4fcigky1658t3-OpXl4GrfP7dnuGN/edit?gid=1779260325> Accessed: December 15, 2024.
- 809 [2] Anonymized. 2024. Robot\_LLM\_Modular\_CV\_anonymous.mp4. Video. <https://drive.google.com/file/d/1QHUNmvWTMWS5MtPXRNRxhEmTEuAL/view> Accessed: December 15, 2024.
- 810 [3] Juecoree. 2020. Forward and Reverse Kinematics for 3R Planar Manipulator. <https://hive.blog/hive-196387/@juecoree/forward-and-reverse-kinematics-for-3r-planar-manipulator> Accessed: 2024-12-15.
- 811 [4] Allan Kouidri. 2023. Grounding DINO: Leading the Way in Zero-Shot Object Detection. <https://www.ikomia.ai/blog/grounding-dino-zero-shot-detection-explained?form=MG0AV3> Accessed: 2024-12-15.
- 812 [5] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, and Lei Zhang. 2023. Grounding DINO: Marrying DINO with Grounded Pre-Training for Open-Set Object Detection. arXiv:2303.05499 [cs.CV]
- 813 [6] Matthias Minderer, Alexey Gritsenko, and Neil Houlsby. 2023. Scaling Open-Vocabulary Object Detection. arXiv:2306.09683 [cs.CV]
- 814 [7] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, Xiao Wang, Xiaohua Zhai, Thomas Kipf, and Neil Houlsby. 2022. Simple Open-Vocabulary Object Detection with Vision Transformers. *arXiv preprint arXiv:2305.06230* (2022).
- 815 [8] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zhihao Fan. 2024. Qwen2 Technical Report. *arXiv preprint arXiv:2407.10671* (2024).
- 816 [9] Tiancheng Zhao, Peng Liu, Xuan He, Lu Zhang, and Kyusong Lee. 2024. Real-time Transformer-based Open-Vocabulary Detection with Efficient Fusion Head. *arXiv preprint arXiv:2403.06892* (2024).

833 Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884