# Snap, Segment, Deploy: A Visual Data and Detection Pipeline for Wearable Industrial Assistants

Di Wen[1], Junwei Zheng[1], Ruiping Liu[1], Yi Xu[1], Kunyu Peng[1,†], Rainer Stiefelhagen[1]

*Abstract*— Industrial assembly tasks increasingly demand rapid adaptation to complex procedures and varied components, yet are often conducted in environments with limited computing, connectivity, and strict privacy requirements. These constraints make conventional cloud-based or fully autonomous solutions impractical for factory deployment. This paper introduces a mobile-device-based assistant system for industrial training and operational support, enabling real-time, semi-hands-free interaction through on-device perception and voice interfaces. The system integrates lightweight object detection, speech recognition, and Retrieval-Augmented Generation (RAG) into a modular on-device pipeline that operates entirely on-device, enabling intuitive support for part handling and procedure understanding without relying on manual supervision or cloud services. To enable scalable training, we adopt an automated data construction pipeline and introduce a two-stage refinement strategy to improve visual robustness under domain shift. Experiments on our generated dataset, *i.e.*, Gear8, demonstrate improved robustness to domain shift and common visual corruptions. A structured user study further confirms its practical viability, with positive user feedback on the clarity of the guidance and the quality of the interaction. These results indicate that our framework offers a deployable solution for real-time, privacy-preserving smart assistance in industrial environments. We will release the Gear8 dataset and source code upon acceptance.

## I. INTRODUCTION

The complexity of manufacturing assembly tasks is escalating, driven by the increasing diversity of components, intricate assembly procedures, and the growing need for customization in contemporary production settings [1]. While full automation is still unfeasible in many contexts, newly onboarded or rotating workers often struggle with steep learning curves, resulting in extended training periods and a higher incidence of assembly mistakes [2]. To address these challenges, vision-based assistance systems have emerged as a promising solution, offering real-time, step-by-step guidance that helps reduce both task completion time and error rates [3].

However, the development of multimodal, on-device assistants for industrial use remains underexplored. Many existing solutions rely on cloud infrastructure, require extensive manual annotation, or fail to operate under practical constraints such as limited compute, offline usage, and strict data privacy [5], [6], [7].

† Corresponding author

All authors are with the Computer Vision for Human-Computer Interaction Lab (cv:hci), Institute for Anthropomatics and Robotics (IAR), Karlsruhe Institute of Technology, 76131 Karlsruhe, Germany. Emails: {di.wen, junwei.zheng, ruiping.liu, kunyu.peng, rainer.stiefelhagen}@kit.edu,yi.xu@student.kit.edu
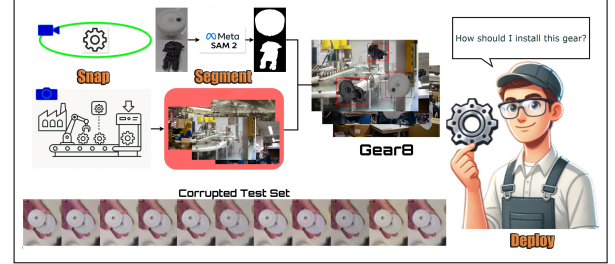
Fig. 1. Overview of our three-stage pipeline: **Snap**, **Segment**, and **Deploy**. In the *Snap* stage, we collect multi-angle part videos and real-world background images from industrial environments. In the *Segment* stage, we extract instance masks using SAM2 [4] and generate a synthetic dataset (Gear8) by compositing parts into factory scenes. In the *Deploy* stage, a lightweight object detector trained on Gear8 is integrated into a wearable assistant system featuring speech input, visual recognition, knowledge retrieval, and audio feedback for interactive on-site guidance. The bottom row shows corrupted test images simulating real-world conditions for robustness evaluation.

We present a mobile-device-based industrial assistance system that leverages computer vision and Large Language Models (LLMs) to deliver real-time, context-aware guidance during assembly tasks. The system comprises a lightweight object detection module for part identification, a speech-to-text interface for natural language input, and a Retrieval-Augmented Generation (RAG) engine [8] that generates part-specific responses, which are conveyed through text-to-speech output. This end-to-end, on-device pipeline supports intuitive and semi-hands-free interaction, eliminating the need for cloud-based processing. To facilitate training without manual annotations, we propose a fully automated data generation pipeline that synthesizes part-level training images using multi-view video captures and compositing over diverse industrial backgrounds. Additionally, to improve the detector's generalization under domain shift, we introduce a two-stage training approach termed Background-Agnostic Refinement (BAR), wherein the model is fine-tuned on plain-background object crops to prioritize object-centric features.

The system is modular, allowing components such as the vision model, knowledge base, or language model to be updated independently for new tasks or deployment settings.

Bridging real-time object detection, vision-language interaction, and mobile industrial AI, our framework offers a deployable solution for contextual task assistance in constrained factory environments.

Extensive experiments on the Gear8 dataset, along with a structured user study, demonstrate the system's effectiveness in improving detection robustness, task efficiency, and overall usability.

Our contributions are summarized as follows:

- We design a real-time, multimodal assistant system for mobile industrial deployment, integrating visual detection, voice interaction, and semantic retrieval for hands-free part recognition and instruction.
- We propose a fully automated dataset generation pipeline that requires no manual annotation and supports diverse industrial backgrounds using consumer-grade equipment.
- We introduce a two-stage training strategy, Background-Agnostic Refinement (BAR), which improves detection robustness under domain shift without modifying model architecture.

## II. RELATED WORK

### A. Synthetic Data Generation for Vision Tasks

Synthetic data has proven crucial for modern vision, initially in object detection and now across segmentation, pose estimation, and beyond. Early work used 3D renders: Peng et al. [9] fine-tuned networks pre-trained on CAD images with minimal real data for PASCAL VOC. Game-engine datasets like "Playing for Data"/"Playing for Benchmarks" [10], [11] and Virtual KITTI [12] provided pixel-perfect urban scenes, while URSA [13] and ProcSy [14] scaled this to millions of driving images.

Domain randomization [15], [16] addresses the sim-to-real gap by randomizing textures, lighting, and viewpoints. GAN-based refinement improves realism: Nogues et al. [17] and Hu et al. [18] apply CycleGAN variants to industrial parts; Lin et al. [19] use multi-task GANs for traffic-sign detection. BigDatasetGAN [20] synthesizes pixel-annotated ImageNet, and InstaGen [21] employs diffusion models for open-vocabulary detection. Copy-paste augmentation remains simple yet effective [22]. Synthetic data has also enabled progress in diverse tasks such as urban-scene parsing [23], remote-sensing detection [24], and medical or robotic vision. Moreover, domain-adaptive detectors [25] and self-supervised learners benefit significantly from combining synthetic and real data.

We build on the copy-paste paradigm by combining factory-specific background images with instance masks extracted from full-view object captures, producing diverse, domain-adapted synthetic images that strike a balance between simplicity and high realism for industrial inspection.

### B. Real-Time Egocentric Vision in Wearable Systems

Assistive technology [26], [27], especially real-time egocentric vision [28], [29], has become increasingly relevant with the rise of smart glasses and edge AI devices. Several studies [30], [31], [32] have demonstrated the deployment of YOLO-based detectors on wearable platforms for tasks like wire harness inspection and redundant part detection in assembly lines. Other works [33], [34], [35], [36] extend this idea to maintenance and defect localization, leveraging egocentric views to support visual overlays and spatial guidance.

To address hardware limitations, researchers have proposed lightweight architectures [37], edge-cloud offloading [38], and wearable edge fusion [39], balancing latency and model performance. Some works further explore interaction modeling and long-term tracking [40], [41], supporting higher-level scene understanding in dynamic environments.

Assistive systems also benefit from real-time egocentric detection. Vision–tactile feedback [42], auditory guidance [43], and multi-modal hazard detection [44] have shown promise in aiding the visually impaired or supporting safety in challenging environments. Additionally, MATERobot [28] enables real-time material and object recognition with tactile feedback to improve environmental awareness for visually impaired individuals. Recent efforts such as ObjectFinder [29] integrate open-vocabulary object detection and navigation to support interactive object search for blind users. Moreover, vision-language systems like EgoVLP [45] and EdgeVL [46] begin to bridge perception and semantic reasoning, though on-device deployment remains constrained. Notably, Vinci [47] introduces a portable egocentric vision-language assistant capable of real-time scene understanding and user interaction on wearable devices.

Unlike prior work relying on proprietary hardware or cloud infrastructure, our system emphasizes low-cost, edge-deployable design for real-time multimodal perception and feedback in practical industrial workflows.

## III. DATASET CONSTRUCTION

Acquiring large-scale labeled datasets in industrial settings poses significant challenges due to factors such as strict confidentiality policies that limit third-party access and the lack of specialized annotation resources or infrastructure in small- and medium-sized enterprises. In light of these constraints, it becomes essential to develop a data generation framework that empowers non-expert factory staff to efficiently produce high-quality visual datasets using readily available consumer-grade devices, without the need for advanced technical training. To meet this objective, we introduce a synthetic dataset of industrial parts specifically designed to reflect the conditions and requirements of real-world deployment. This section outlines our methodology for data capture, the synthetic generation pipeline, and the key properties of the resulting dataset.

### A. Background Collection

Due to restricted factory access, we selected 12 high-resolution, license-free industrial background images from public sources. Although this introduces a domain gap, it imposes a stricter evaluation setting and may underestimate real deployment performance.

### B. Multi-View Part Capture

To simulate realistic conditions, we captured multi-view video sequences of each part using smartphones, without controlled environments or specialized equipment. The intentional use of different imaging devices from those in

deployment introduces a domain gap, enabling robustness evaluation under suboptimal conditions.

From each video sequence $\{F_t\}_{t=1}^T$, a subset of frames is uniformly sampled across diverse viewpoints to maximize coverage while minimizing redundancy. These frames are then used for subsequent instance segmentation and synthetic data generation.

## C. Synthetic Composition

To generate training samples, we first extract instance masks from each video frame $F_t$ using the SAM2 [4] model with center point guidance:

$$M_t = \text{SAM}_2(F_t, c_t), \tag{1}$$

where $c_t$ denotes the center of the target object.

We then synthesize composite images by randomly pasting $k \in [3,5]$ object masks onto a selected background image $B$. Each mask undergoes random affine transformations:

$$m_i' = \text{Affine}(m_i; \theta_i), \quad \theta_i = \{\text{scale}, \text{rotation}\}. \tag{2}$$

To ensure realism, we constrain the pairwise IoU between pasted instances to be below 0.5 and limit the number of instances per category to two per image. Object categories are sampled uniformly to ensure class balance.

The final composite image is defined as:

$$I = B \oplus \{m_1', \ldots, m_k'\}. \tag{3}$$

This composition strategy promotes diversity in object appearance, position, and category, supporting generalizable and robust model training.

## D. Robustness Evaluation via Corruption Simulation

To assess model robustness under deployment constraints, we construct a corrupted test set used exclusively for evaluation. We consider ten types of common perturbations, categorized into blur (motion blur, Gaussian blur), noise (Gaussian noise, ISO noise), and color or lighting variations (HSV modification, color shift, brightness adjustment, and contrast change).

Each clean test image $I$ is corrupted by applying $c_i \in \mathcal{C}$ at a fixed, moderate severity $\alpha_i$:

$$I_i' = c_i(I; \alpha_i), \tag{4}$$

resulting in ten distinct corrupted versions per image, each reflecting a specific type of degradation.

Overall, our dataset construction strategy enables scalable, annotation-free training while ensuring diversity, realism, and robustness evaluation. The intentional domain gap between training and deployment conditions further validates the practical applicability of our approach to real-world industrial environments.
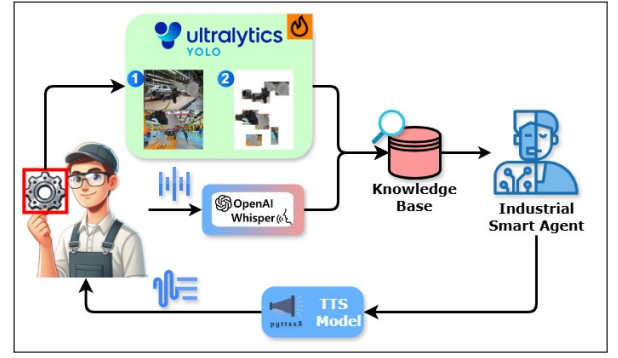


Fig. 2. System pipeline of the proposed industrial assistant integrates YOLO [48]-based detection, Whisper speech recognition, structured knowledge retrieval, and audio feedback. The detector is trained via a two-stage strategy with Background-Agnostic Refinement (BAR) to enhance robustness to domain shifts.

## IV. METHODOLOGY

Our proposed industrial smart assistant system is designed to support real-time part recognition and multimodal user interaction under constrained computational environments. The pipeline consists of two main modules: a lightweight object detection framework enhanced by two-stage refinement, and a retrieval-augmented multimodal interaction system integrating visual and speech inputs.

## A. Background-Agnostic Refinement (BAR)

To improve the generalization ability of lightweight object detectors under domain shifts, we adopt a two-stage training strategy called Background-Agnostic Refinement (BAR).

In the first stage, the detector is trained conventionally on the synthetic dataset containing diverse background conditions. After reaching high validation performance, we proceed to the second stage: the trained detector is re-applied to the original training set to extract high-confidence predictions. Each detected bounding box $b$ is cropped from its original image and placed onto a clean white canvas, eliminating background noise while preserving part-specific visual cues.

The resulting purified pseudo-labeled dataset emphasizes object-intrinsic features and reduces overfitting to contextual artifacts. We fine-tune the detector on this refined set, which empirically enhances robustness and improves detection stability across varying deployment environments.

## B. Retrieval-Augmented Multimodal Interaction

To enable intuitive user interaction with minimal manual effort, we design a retrieval-augmented multimodal system based on smart glasses. The pipeline consists of three modules:

*1) Knowledge Base Construction:* To enable the smart assistant to provide part-specific internal knowledge required in industrial scenarios, we construct a structured knowledge base covering all relevant components.

For each labeled part $x \in \mathcal{D}$, a semantic embedding $\mathbf{v}_x$ is computed using a pretrained SentenceTransformer [49]:

$$\mathbf{v}_x = f_{\text{enc}}(x), \tag{5}$$

where $f_{\text{enc}}(\cdot)$ denotes the encoder function.

All embeddings are indexed using FAISS [50], an efficient library for similarity search. In our setting, FAISS constructs a flat L2 index to enable rapid nearest-neighbor retrieval based on Euclidean distance in the semantic space.

*2) Query Acquisition:* When the user triggers a query through a button press, the system continuously captures video frames $\{I_1, \ldots, I_n\}$ from the wearable device. To ensure that the user's field of view is properly aligned with the target object, detection is performed on each frame individually, and only when valid detections are observed over $N$ consecutive frames does the system proceed to the next step.

The final frame $I_n$ is used for subsequent processing. It is passed through our object detector, producing a set of object predictions. To improve detection robustness, we perform multi-frame fusion on the buffered detections across the last $N$ frames, aggregating predictions based on label consistency and bounding box IoU. For detections with the same label and IoU greater than a predefined threshold $\tau$, the bounding boxes are merged via confidence-weighted averaging:

$$b' = \frac{\sum_{i=1}^{k} w_i b_i}{\sum_{i=1}^{k} w_i}, \qquad (6)$$

where $b_i$ and $w_i$ denote the bounding box and its associated confidence score, respectively. After fusion, redundant detections with high spatial overlap and identical labels are suppressed using a duplicate box filtering strategy, ensuring that each object is represented only once. The resulting fused and filtered detections are then passed through DepthAnything [8] to infer pixel-wise depth maps.

The fused detections are then sorted in ascending order of depth, prioritizing closer objects for subsequent retrieval and interaction. Finally, the top-$K$ nearest objects (with $K = 3$) are selected to construct the query candidates.

Formally, the output of this phase is:

$$\{(l_j, b_j, d_j)\}_{j=1}^{K}, \qquad (7)$$

where $l_j$ denotes the object label, $b_j$ the bounding box, and $d_j$ the estimated depth for the $j$-th selected object.

*3) Semantic Retrieval and Knowledge-Augmented Response Generation:* During runtime, after the top-$K$ detected objects $\{(l_j, b_j, d_j)\}_{j=1}^{K}$ are selected based on depth ranking, each label $l_j$ is embedded into a vector representation $\mathbf{v}_{l_j}$ using a pretrained SentenceTransformer. These embeddings are used to query the FAISS index to retrieve the nearest database entries $x^*$, providing part-specific context $C$ for response generation.

Given the user's transcribed query $q$ and the retrieved context $C$, a language model (LLM) generates a natural language response tailored to the user's question:

$$\text{Answer} = \text{LLM}(q, C). \qquad (8)$$

The generated answer is synthesized into speech via a text-to-speech (TTS) engine and played back to the user through the wearable device, enabling hands-free interaction.

## V. EXPERIMENTS

### A. Dataset

We construct the *Gear8* dataset following the pipeline in Section IV-A. It consists of two visually similar first-stage reduction gearboxes, each containing two gear types and a pair of cover and housing components. To synthesize the training data, we composited segmented parts onto 12 industrial backgrounds sourced from license-free image platforms, selected to cover diverse industrial contexts. This process yields 4,000 training images and 1,000 validation images. The test set comprises 196 real-world images capturing actual assembly scenarios and background environments. To simulate deployment conditions in practical systems, we further applied 10 types of realistic visual corruptions, resulting in an augmented test set of 2,156 images.

### B. Implementation Details

All experiments are conducted using PyTorch 2.6.0 with CUDA 12.4. The YOLOv11n [48] model is fine-tuned on the *Gear8* dataset, with second-stage training triggered once validation mAP50 exceeds 95.0. Multi-frame predictions are aggregated using confidence-weighted IoU merging. We set the detection confidence threshold to 0.4 and use a minimum of 5 consecutive frames with valid detections to confirm object presence. IoU threshold for merging is set to 0.5, with a minimum vote count of 3. Depth estimation is performed using the DepthAnything (vitl14) [51] model. Semantic embeddings are generated using SentenceTransformer (all-MiniLM-L6-v2) [49] and queried via FAISS [50] over a structured component database. A local Phi-3-mini-4k-instruct [52] language model is used for dialogue generation. Whisper-small [53] handles speech recognition, and pyttsx3 [54] is used for text-to-speech synthesis. Voice input is recorded with a sample rate of 16kHz and a duration of 8 seconds. Real-time interaction is supported through depth-aware sorting of detected components and response synthesis based on retrieved metadata.

### C. Results on Object Detection

Table I presents the performance comparison under different combinations of training strategies and inference augmentations. The baseline model, YOLOv11n, shows limited accuracy due to its ultra-lightweight architecture. Applying Background-Agnostic Refinement (BAR) during training leads to consistent improvements across all metrics, particularly mAP@0.5:0.95, which rises from 0.06 to 0.17. This improvement stems from BAR's ability to decouple foreground object learning from background distractions. By refining the model with synthetic variations of the same objects on diverse backgrounds, BAR encourages the detector to focus on intrinsic object features rather than context-dependent cues, leading to better generalization under domain shifts.

Among inference-time methods, Test-Time Augmentation (TTA) yields the greatest gain. TTA aggregates predictions from spatially transformed inputs to improve robustness against scale and localization variations. When combined with BAR, it improves mAP@0.5 to 0.30 and

mAP@0.5:0.95 to 0.20, demonstrating stronger generalization across object sizes and configurations.

| Method | mP | mR | mAP@0.5 | mAP@0.5:0.95 |
|---|---|---|---|---|
| YOLOv11n (Baseline) | 0.24 | 0.21 | 0.20 | 0.06 |
| w/ BAR | 0.28 | 0.25 | 0.26 | 0.17 |
| w/ BAR+SAHI | 0.23 | 0.23 | 0.20 | 0.06 |
| w/ BAR+TTA | 0.35 | 0.26 | 0.30 | 0.20 |
| w/ BAR+TTA+SAHI | 0.31 | 0.26 | 0.25 | 0.16 |

In contrast, SAHI [55], which slices input images into overlapping tiles for fine-grained detection, shows minimal or even negative effect. This is likely because our dataset primarily consists of small, well-centered objects in uncluttered settings. Tiling fragments the global context and introduces redundant detections, which ultra-lightweight models like YOLOv11n [48] are less capable of reconciling effectively. Notably, combining TTA and SAHI leads to a drop in performance compared to TTA alone, suggesting that the spatial fragmentation from SAHI interferes with the consistency benefits introduced by TTA.

Overall, BAR and TTA complement each other well in our setting, while SAHI proves less effective due to the dataset characteristics and model capacity. Based on this evaluation, we deploy the BAR+TTA-enhanced model in our industrial smart assistant system to ensure accurate, real-time performance under deployment constraints.

### D. Qualitative Results

To further evaluate the performance of our framework, we present qualitative detection results for the object detection module under both clean and corrupted conditions, as shown in Figure 3.

In the first row, accurate detections on uncorrupted test images demonstrate that the object detector is capable of handling real-world challenges such as partial occlusions and high inter-class similarity, indicating its ability to perform fine-grained recognition under clean conditions.

In contrast, the second row illustrates failure cases under clean conditions. Columns one to three show misclassifications of similar-looking parts—such as different versions of the gearbox cover, gearbox housing, and gears—particularly under specific viewpoints. These errors may be attributed to the limited capacity of the lightweight detector employed in our system, which may struggle to resolve subtle visual differences. Column four further demonstrates a failure under severe occlusion, where critical part features are no longer visible, leading to missed detections.

The third row depicts corruption-induced failures (as defined in Section III-D). Although the corresponding clean images were correctly detected, moderate perturbations lead to a notable drop in performance. In particular, ISO noise in column two results in complete detection failure, while imperceptible corruptions in column four cause misclassification into similar-looking parts. These observations confirm

the sensitivity of the detection module to input distribution shifts and further validate the need for robustness-aware evaluation in real-world deployment settings.

### E. User Study

Following preliminary experiments on object detection using our constructed dataset, we further evaluate the system in real-world scenarios with a focus on user experience in wearable settings. Given that usability and interaction quality are critical for practical deployment, we conduct a structured user study consisting of task-based evaluations and a post-task questionnaire to assess system effectiveness, usability, and user satisfaction.

*Organization :* To evaluate the effectiveness and usability of the system, we conducted a user study involving 9 participants (6 male, 3 female) in real-world assembly scenarios. For the evaluation, we selected components from 8 different object categories present in the training set, with varying instance counts for each class. Each participant interacted with the system for 30 minutes. To compare performance with and without system assistance, participants were asked to complete two rounds of part assembly: one without any system support, and another strictly following the system's step-by-step guidance. The time taken for each round was recorded separately. During the assisted phase, participants could also interact with the system to inquire about part-specific knowledge. After both rounds, all participants completed an anonymous questionnaire session. This included the NASA Task Load Index (NASA-TLX) [56] to assess cognitive workload, along with six custom-designed questions targeting the perceived usefulness, clarity, and reliability of our system.

*Performance:* To assess system performance from a user perspective, we designed six custom questions covering object recognition, instruction accuracy, task guidance, efficiency, voice clarity, and speech recognition. Participants rated each item on a 10-point Likert scale, with higher scores indicating better experience. As shown in Figure 4, the system received high ratings in task-related aspects, including *Task Guidance* and *Efficiency* (both 8.0), as well as in *Voice Clarity* and *Speech Recognition* (8.0). *Instruction Accuracy* also performed well (7.0), suggesting consistent and understandable system output. The lower score in *Object Recognition* (6.0) is attributed to the use of a lightweight YOLOv11n model (Table I) with only 2.6M parameters, necessary for real-time performance on mobile devices. Additional domain gaps caused by resolution differences between training and deployment environments further impacted detection stability. Despite mitigation through multi-frame fusion, recognition remains constrained by computational limitations. These results indicate that the proposed system offers reliable and efficient assistance in real-world assembly tasks, with potential for further improvement as hardware capabilities advance.

*Time Efficiency and Error Prevention:* To evaluate the impact of the smart assistant on operational performance, we recorded task completion times for each participant with
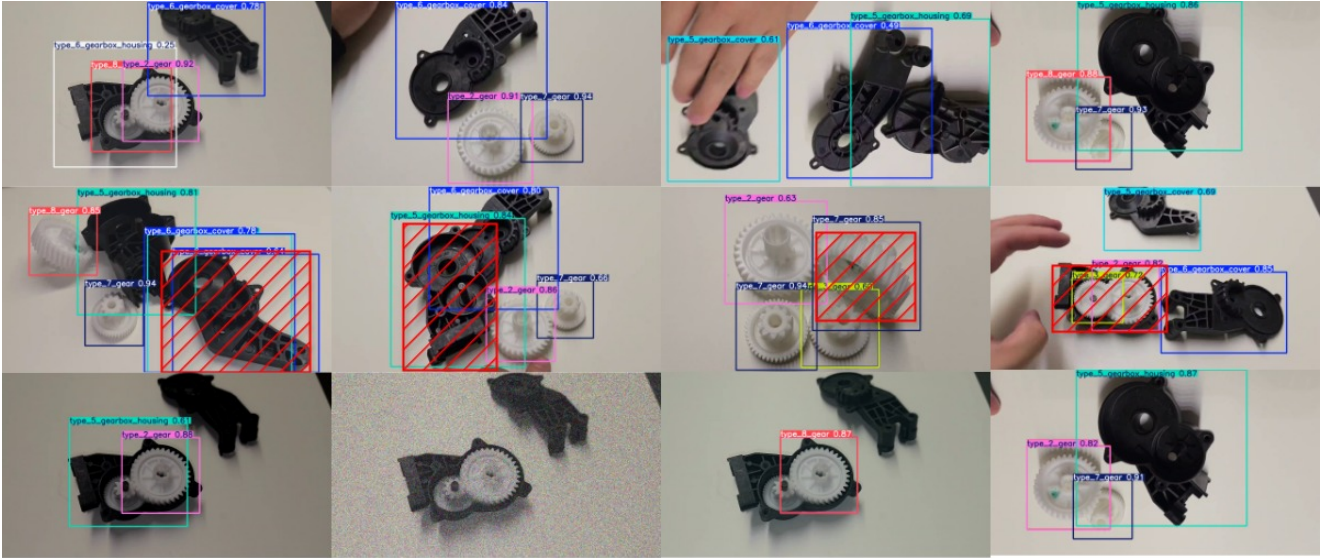
Fig. 3. Qualitative detection results on the Gear8 dataset. The first row shows successful object detection cases on the uncorrupted test set, where all components are correctly identified. The second row displays failure cases on the same test set, where missed or incorrect detections are explicitly highlighted using red bounding boxes with hatched fill. The third row illustrates corruption-induced failures, as defined in Section III-D. All results are visualized using a confidence threshold of 0.25.

and without system support. As shown in Table 4, the average completion time decreased from 542.3 s to 219.2 s, demonstrating a substantial improvement in task efficiency. An exception was observed for participant ID 7, who completed the task in 93 s unaided but took 354 s with guidance. This likely reflects prior task familiarity, where the system's structured workflow introduced additional overhead. Such outliers are rare and do not affect the overall trend.

In addition to improved efficiency, the system also enhanced operational safety. In the unassisted condition, a total of 5 parts were damaged due to incorrect interference fits, resulting in surface abrasion on shaft components. No such incidents occurred under system guidance, indicating that the assistant effectively prevents critical assembly errors and reduces potential production costs.

*Cognitive Load:* To evaluate the cognitive load introduced by our wearable system, we employed the NASA Task Load Index (NASA-TLX) [56], a widely used measure of subjective workload. Participants rated six dimensions: Mental, Physical, Temporal Demand, Performance, Effort, and Frustration. Figure 4 presents the average scores across all participants. The overall score was 48.0, close to the global median (range: 6.21–88.50, $N = 1173$ [57]), indicating moderate workload. The system is effective without excessive strain. Among all factors, *Frustration* (57) and *Mental Demand* (55) were relatively high. Likely due to the effort required to interact with a multimodal system, especially when precise part names are needed. Occasional language model hallucinations, verbose outputs, or extended response times may also have contributed. In contrast, the lowest score was for *Physical Demand* (37), confirming that the semi-hands-free interface minimizes physical effort. Overall, the system provides a positive user experience with acceptable cognitive load with potential for further improvement through interaction refinement.

*Overall Rating:* Participants rated their likelihood to recommend the Smart Assistant system to others on a 0–10 scale. The average score was 6.11 (SD = 1.62), indicating moderate satisfaction. Individual scores ranged from 4 (lowest) to 8 (highest), highlighting variability in user perceptions and experiences.

*Limitations:* The current system is subject to two main limitations. First, due to safety and privacy constraints in industrial environments, the system runs fully offline on resource-constrained wearable devices. This requires lightweight models to meet real-time constraints, limiting overall capacity and accuracy. Second, the devices used for data collection and deployment differ in resolution and sensor characteristics, introducing a domain gap that may affect detection performance during testing. These issues, however, are not fundamental limitations of the proposed pipeline. In practice, they can be addressed by aligning the deployment hardware with the data acquisition setup, and by integrating higher-performance servers within the factory's local network. Such improvements would allow the adoption of more capable models while preserving the modular structure and deployment flow of the system, leading to further gains in recognition accuracy and interaction quality.

## VI. CONCLUSION

We introduced a semi-hands-free smart assistant for industrial assembly, tailored for offline use, limited compute, and strict privacy. The system integrates lightweight object detection, depth estimation, and retrieval-augmented dialogue in a mobile-friendly pipeline. To avoid manual labeling, we developed an automated data generation process and a two-stage training method (BAR) to boost detection robustness. Experiments on the Gear8 dataset show improved accuracy with high efficiency, while a user study confirmed gains in task speed, reduced part damage, and lower cognitive load. Despite current hardware and domain shift limitations,
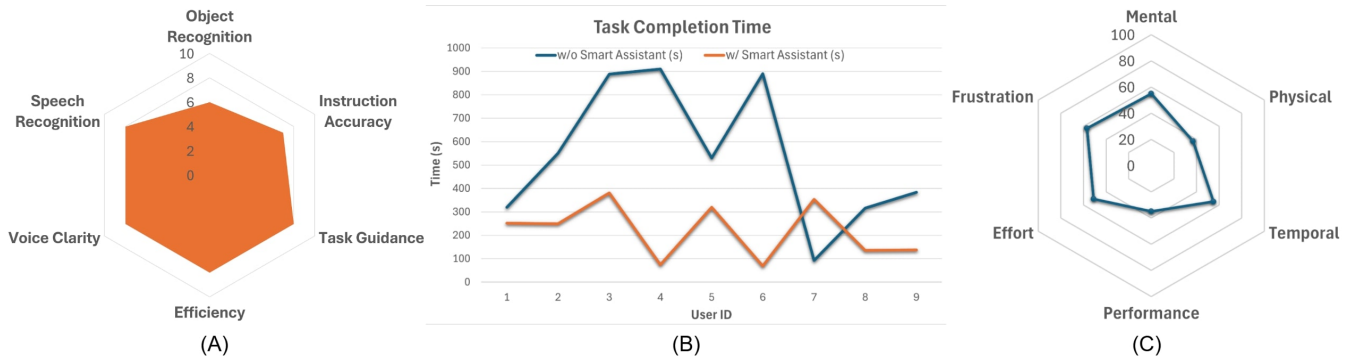
Fig. 4. User-study results for the Smart Assistant system: (a) average scores on the customized functionality-assessment questionnaire (0–10, higher is better); (b) task-completion time with and without the Smart Assistant for nine participants; (c) average NASA-TLX workload scores (0–100, lower is better).

future updates with local servers and aligned data capture can address these. Our approach provides a scalable, privacy-preserving solution for intelligent industrial assistance.

## ACKNOWLEDGMENT

## REFERENCES

[1] Matteo Capponi, Riccardo Gervasi, Luca Mastrogiacomo, and Fiorenzo Franceschini. Assembly complexity and physiological response in human-robot collaboration: Insights from a preliminary experimental analysis. *Robotics and Computer-Integrated Manufacturing*, 2024.

[2] Luca Gualtieri, Maximilian Öhler, Andrea Revolti, and Patrick Dallasega. A visual management and augmented-reality-based training module for the enhancement of short and long-term procedural knowledge retention in complex machinery setup. *Computers & Industrial Engineering*, 2024.

[3] Lea M Daling and Sabine J Schlittmeier. Effects of augmented reality-, virtual reality-, and mixed reality–based training on objective performance measures and subjective evaluations in manual assembly tasks: a scoping review. *Human factors*, 2024.

[4] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, et al. Sam 2: Segment anything in images and videos. *arXiv preprint arXiv:2408.00714*, 2024.

[5] Jinjie Liu, Guolong Liu, Huan Zhao, Junhua Zhao, Jing Qiu, and Zhao Yang Dong. Real-time industrial carbon emission estimation with deep learning-based device recognition and incomplete smart meter data. *Engineering Applications of Artificial Intelligence*, 127:107272, 2024.

[6] Vahid Gharibvand, Mohammad Karimzadeh Kolamroudi, Qasim Zeeshan, Zeki Murat Çınar, Saeid Sahmani, Mohammed Asmael, and Babak Safaei. Cloud based manufacturing: A review of recent developments in architectures, technologies, infrastructures, platforms and associated challenges. *The International Journal of Advanced Manufacturing Technology*, 131(1):93–123, 2024.

[7] Naval Kishore Mehta, Shyam Sunder Prasad, Sumeet Saurav, Ravi Saini, and Sanjay Singh. Iar-net: A human-object context guided action recognition network for industrial environment monitoring. *IEEE Transactions on Instrumentation and Measurement*, 2024.

[8] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474, 2020.

[9] Xingchao Peng, Baochen Sun, Karim Ali, and Kate Saenko. Learning deep object detectors from 3d models. In *Proceedings of the IEEE international conference on computer vision*, pages 1278–1286, 2015.

[10] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, pages 102–118. Springer, 2016.

[11] Stephan R Richter, Zeeshan Hayder, and Vladlen Koltun. Playing for benchmarks. In *Proceedings of the IEEE international conference on computer vision*, pages 2213–2222, 2017.

[12] Adrien Gaidon, Qiao Wang, Yohann Cabon, and Eleonora Vig. Virtual worlds as proxy for multi-object tracking analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4340–4349, 2016.

[13] Matt Angus, Mohamed ElBalkini, Samin Khan, Ali Harakeh, Oles Andrienko, Cody Reading, Steven Waslander, and Krzysztof Czarnecki. Unlimited road-scene synthetic annotation (ursa) dataset. In *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pages 985–992. IEEE, 2018.

[14] Samin Khan, Buu Phan, Rick Salay, and Krzysztof Czarnecki. Procsy: Procedural synthetic dataset generation towards influence factor studies of semantic segmentation networks. In *CVPR workshops*, volume 3, page 4, 2019.

[15] Josh Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *2017 IEEE/RSJ international conference on intelligent robots and systems (IROS)*, pages 23–30. IEEE, 2017.

[16] Jonathan Tremblay, Aayush Prakash, David Acuna, Mark Brophy, Varun Jampani, Cem Anil, Thang To, Eric Cameracci, Shaad Boochoon, and Stan Birchfield. Training deep networks with synthetic data: Bridging the reality gap by domain randomization. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pages 969–977, 2018.

[17] Fernando Camaro Nogues, Andrew Huie, and Sakyasingha Dasgupta. Object detection using domain randomization and generative adversarial refinement of synthetic images. *arXiv preprint arXiv:1805.11778*, 2018.

[18] Jiwei Hu, Feng Xiao, Qiwen Jin, Guangpeng Zhao, and Ping Lou. Synthetic data generation based on rdb-cyclegan for industrial object detection. *Mathematics*, 11(22):4588, 2023.

[19] Y Lin, K Suzuki, H Takeda, and K Nakamura. Generating synthetic training data for object detection using multi-task generative adversarial networks. *ISPRS Annals of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, 2:443–449, 2020.

[20] Daiqing Li, Huan Ling, Seung Wook Kim, Karsten Kreis, Sanja Fidler, and Antonio Torralba. Bigdatasetgan: Synthesizing imagenet with

pixel-wise annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21330–21340, 2022.

[21] Chengjian Feng, Yujie Zhong, Zequn Jie, Weidi Xie, and Lin Ma. Instagen: Enhancing object detection by training on synthetic dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14121–14130, 2024.

[22] Debidatta Dwibedi, Ishan Misra, and Martial Hebert. Cut, paste and learn: Surprisingly easy synthesis for instance detection. In *Proceedings of the IEEE international conference on computer vision*, pages 1301–1310, 2017.

[23] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3234–3243, 2016.

[24] Jacob Shermeyer, Thomas Hossler, Adam Van Etten, Daniel Hogan, Ryan Lewis, and Daeil Kim. Rareplanes: Synthetic data takes flight. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 207–217, 2021.

[25] Yuhua Chen, Wen Li, Christos Sakaridis, Dengxin Dai, and Luc Van Gool. Domain adaptive faster r-cnn for object detection in the wild. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3339–3348, 2018.

[26] Xin Jiang, Junwei Zheng, Ruiping Liu, Jiahang Li, Jiaming Zhang, Sven Matthiesen, and Rainer Stiefelhagen. @bench: Benchmarking vision-language models for human-centered assistive technology. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2025.

[27] Ruiping Liu, Jiaming Zhang, Kunyu Peng, Junwei Zheng, Ke Cao, Yufan Chen, Kailun Yang, and Rainer Stiefelhagen. Open scene understanding: Grounded situation recognition meets segment anything for helping people with visual impairments. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1857–1867, 2023.

[28] Junwei Zheng, Jiaming Zhang, Kailun Yang, Kunyu Peng, and Rainer Stiefelhagen. Materobot: Material recognition in wearable robotics for people with visual impairments. In *2024 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2303–2309. IEEE, 2024.

[29] Ruiping Liu, Jiaming Zhang, Angela Schön, Karin Müller, Junwei Zheng, Kailun Yang, Kathrin Gerling, and Rainer Stiefelhagen. Objectfinder: Open-vocabulary assistive system for interactive object search by blind people. *arXiv preprint arXiv:2412.03118*, 2024.

[30] Shuo Li, Mingzhe Yuan, Wenhong Wang, Feidao Cao, Hongyan Shi, Yuhang Zhang, and Xiangpu Meng. Enhanced yolo-and wearable-based inspection system for automotive wire harness assembly. *Applied Sciences*, 14(7):2942, 2024.

[31] Qijie Zhao, Yaohui Kong, Shaojie Sheng, and Junjun Zhu. Redundant object detection method for civil aircraft assembly based on machine vision and smart glasses. *Measurement Science and Technology*, 33(10):105011, 2022.

[32] José Silva, Pedro Coelho, Luzia Saraiva, Paulo Vaz, Pedro Martins, and Alfonso López-Rivero. Validating the use of smart glasses in industrial quality control: a case study. *Applied Sciences*, 14(5):1850, 2024.

[33] Shaohan Wang, Sakib Ashraf Zargar, and Fuh-Gwo Yuan. Augmented reality for enhanced visual inspection through knowledge-based deep learning. *Structural Health Monitoring*, 20(1):426–442, 2021.

[34] Shufei Li, Pai Zheng, and Lianyu Zheng. An ar-assisted deep learning-based approach for automatic inspection of aviation connectors. *IEEE Transactions on Industrial Informatics*, 17(3):1721–1731, 2020.

[35] Ana Malta, Mateus Mendes, and Torres Farinha. Augmented reality maintenance assistant using yolov5. *Applied Sciences*, 11(11):4758, 2021.

[36] Kyeong-Beom Park, Minseok Kim, Sung Ho Choi, and Jae Yeol Lee. Deep learning-based smart task assistance in wearable augmented reality. *Robotics and Computer-Integrated Manufacturing*, 63:101887, 2020.

[37] Junjie Jiang, Zhile Yang, Chengke Wu, Yuanjun Guo, Meng Yang, and Wei Feng. A compatible detector based on improved yolov5 for hydropower device detection in ar inspection system. *Expert Systems with Applications*, 225:120065, 2023.

[38] Alessandro Farasin, Francesco Peciarolo, Marco Grangetto, Elena Gianaria, Paolo Garza, et al. Real-time object detection and tracking in mixed reality using microsoft hololens. In *VISIGRAPP 2020-Proceedings of the 15th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications*, volume 4, pages 165–172. SciTePress, 2020.

[39] Iason Karakostas, Aikaterini Valakou, Despoina Gavgiotaki, Zinovia Stefanidi, Ioannis Pastaltzidis, Grigorios Tsipouridis, Nikolaos Kilis, Konstantinos C Apostolakis, Stavroula Ntoa, Nikolaos Dimitriou, et al. A real-time wearable ar system for egocentric vision on the edge. *Virtual Reality*, 28(1):44, 2024.

[40] Mingzhen Huang, Xiaoxing Li, Jun Hu, Honghong Peng, and Siwei Lyu. Tracking multiple deformable objects in egocentric videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1461–1471, 2023.

[41] Shaowei Liu, Subarna Tripathi, Somdeb Majumdar, and Xiaolong Wang. Joint hand motion and interaction hotspots prediction from egocentric videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3282–3292, 2022.

[42] Yiwen Chen, Junjie Shen, and Hideyuki Sawada. A wearable assistive system for the visually impaired using object detection, distance measurement and tactile presentation. *Intell. Robot*, 3:420–435, 2023.

[43] Sanjeev U Rao, Swaroop Ranganath, TS Ashwin, Guddeti Ram Mohana Reddy, et al. A google glass based real-time scene analysis for the visually impaired. *IEEE Access*, 9:166351–166369, 2021.

[44] Zhongze Ma, Chenjie Zhang, and Pengcheng Jiao. Underwater smart glasses: A visual-tactile fusion hazard detection system. *Iscience*, 27(4), 2024.

[45] Kevin Qinghong Lin, Jinpeng Wang, Mattia Soldan, Michael Wray, Rui Yan, Eric Z Xu, Difei Gao, Rong-Cheng Tu, Wenzhe Zhao, Weijie Kong, et al. Egocentric video-language pretraining. *Advances in Neural Information Processing Systems*, 35:7575–7586, 2022.

[46] Kaiwen Cai, Zhekai Duan, Gaowen Liu, Charles Fleming, and Chris Xiaoxuan Lu. Self-adapting large visual-language models to edge devices across visual modalities. In *European Conference on Computer Vision*, pages 301–318. Springer, 2024.

[47] Yifei Huang, Jilan Xu, Baoqi Pei, Yuping He, Guo Chen, Mingfang Zhang, Lijin Yang, Zheng Nie, Jinyao Liu, Guoshun Fan, et al. An egocentric vision-language model based portable real-time smart assistant. *arXiv preprint arXiv:2503.04250*, 2025.

[48] Rahima Khanam and Muhammad Hussain. Yolov11: An overview of the key architectural enhancements. *arXiv preprint arXiv:2410.17725*, 2024.

[49] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019.

[50] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. 2024.

[51] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10371–10381, 2024.

[52] Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.

[53] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International conference on machine learning*, pages 28492–28518. PMLR, 2023.

[54] nateshmbhat. pyttsx3. https://github.com/nateshmbhat/pyttsx3, 2024. Python text-to-speech library.

[55] Fatih Cagatay Akyon, Sinan Onur Altinuc, and Alptekin Temizel. Slicing aided hyper inference and fine-tuning for small object detection. In *2022 IEEE international conference on image processing (ICIP)*, pages 966–970. IEEE, 2022.

[56] Sandra G Hart and Lowell E Staveland. Development of nasa-tlx (task load index): Results of empirical and theoretical research. In *Advances in psychology*, volume 52, pages 139–183. Elsevier, 1988.

[57] Rebecca A Grier. How high is high? a meta-analysis of nasa-tlx global workload scores. In *Proceedings of the human factors and ergonomics society annual meeting*, volume 59, pages 1727–1731. Sage Publications Sage CA: Los Angeles, CA, 2015.