# What You See is What You Grasp: User-Friendly Grasping Guided by Near-eye-tracking

Shaochen Wang*, Wei Zhang*, Zhangli Zhou*, Jiaxi Cao, Ziyang Chen, Kang Chen, Bin Li, Zhen Kan

*Abstract*—This paper presents a next-generation human-robot interface that can infer and realize the user's manipulation intention via sight only. Our system integrates near-eye-tracking with robotic manipulation to enable user-specified actions, such as grasping or pick-and-place operations. We develop a head-mounted near-eye-tracking device that tracks the user's eyeball movements in real-time to identify the user's visual attention and enable sight-guided manipulation. To improve grasping performance, we introduce a transformer-based grasp model that uses stacked attention blocks to extract hierarchical features, expanding channel volumes while squeezing feature map resolutions. Experimental validation demonstrates that our system can effectively assist users to complete manipulation tasks through their eyes, which holds great potential for an assistive robot that leverages gaze interaction to aid individuals with upper limb disabilities or the elderly in their daily lives.

## I. Introduction

Robotic systems [1], [2], powered by recent advances in artificial intelligence, have become increasingly prevalent in both industrial and daily life. To provide more daily assistance, robots are no longer limited to performing simple and repetitive tasks but are designed to understand human intention for better helping users. In particular, it is essential to build bridges between humans and machines to transfer human intentions to robots.

Traditional manipulation methods [3] often involves using a joystick to deliver user-specified actions. Unfortunately, this can be challenging for elderly individuals, particularly those with upper limb disabilities. Recent advancements in wearable technology have demonstrated the potential in robotic assistance systems such as brain-computer interface (BCI) [4]. However, invasive BCI devices require microelectrodes to be surgically implanted into the human cerebral cortex, which can be risky. Furthermore, non-invasive BCI devices are susceptible to noise interference and are often expensive. Therefore, there is a pressing need to develop a new human-robot interface that is both safe and user-friendly.

The majority of human sensory input is acquired through the eyes. In fact, over 80% of human information is obtained through vision [5], which motivates the development of eye-based robotic assistive systems to enable manipulation using sight. Robotic assistive systems that incorporate eye-tracking technology have the potential to revolutionize various applications, including surgical diagnosis, rehabilitation, and academic research. For individuals with physical disabilities,

eye-tracking offers a natural interface that can bridge the gap between humans and robots, since the user's vision is typically not affected by their disability. However, despite its potential benefits, eye-tracking has not been widely adopted in practice, largely due to the difficulty of accurately modeling the motion of the eye to capture the gaze point. Currently, many eye-tracking robotic assistive systems [6] rely on fixed cameras and are primarily built for desktop use only.

This paper introduces a next-generation human-robot interface that can infer and execute user's manipulation intention using only sight. We have developed a system that integrates near-eye-tracking and robotic manipulation to enable user-specified actions, including grasping and pick-and-place tasks. To achieve sight-guided manipulation, we have designed a head-mounted near-eye-tracking device that tracks eyeball movements in real-time, enabling us to identify the user's visual attention. Additionally, we have developed a transformer-based global grasp detection framework that enhances the robot's sensing capabilities, facilitating successful execution of user-specified manipulations. Our framework utilizes self-attention to model the long-term spatial dependencies among pixels, and a feature fusion pyramid to merge multi-scale features from each stage, thereby determining the final grasping pose. Experimental validation has demonstrated that the eye-tracking system has low gaze estimation error, and the grasping system performs well on multiple grasping datasets.

The contribution of this paper can be summarized as follows:

- We build a sight-based robotic assistive system for user-specified manipulations. Our system includes a head-mounted device that enables real-time intention inference and a grasping subsystem that utilizes self-attention mechanisms for improved visual grasping.
- We propose a novel human-robot interface that enables more natural and instinctive manipulation by utilizing eye-tracking only.
- Extensive experiments demonstrate the effectiveness of the developed robotic assistive system for manipulation tasks.

## II. Related Work

Robotic manipulation [7], [8], [9] is a fundamental skill that has found widespread applications in manufacturing, industry, and medical operations. Vision-based grasping techniques have been extensively investigated by researchers. Lenz et al. [10] were the first to utilize deep learning to detect grasping rectangles. Redmon et al. [11] employed a convolutional

\* Contribute equally
The authors are with the School of Information Science, University of Science and Technology of China, Hefei, 230026, China.

neural network (CNN) to regress the grasping pose that the robot can execute. Additionally, Morrison et al. [12] designed a generative grasping CNN (GG-CNN) that uses depth input to generate antipodal grasp candidates. With the recent advances in artificial intelligence, the new generation of robots is expected to understand human intention through interactions with users, rather than being limited to low-level tasks such as grasping.

Assistive robotic arms have become increasingly popular among users with upper limb impairments in their daily lives, such as grasping objects and pouring water. While joysticks are often used to control these robots, they can be difficult to use, especially for elderly or upper limb disabled users. In contrast, sight is a natural way for people with physical disabilities or mobility and speech impairments to interact with robotic systems. Eye-tracking technology, which has been used for almost a century, has advanced significantly in recent years and can now be used to manipulate robotic devices by following human gaze.

Note that the attention mechanism mimics the way how human vision works. Since sight can indicate human attention, Hollenstein et al. [13] enhanced the performance of the annotation model by using human sight, and demonstrated that the semantic information embedded in the sight can be well utilized by the entity model. In computer vision tasks, Karessli et al. [14] improved the classification accuracy of zero-shot tasks by introducing human sight as an auxiliary task. In vision-language tasks, Sugano et al. [15] assisted the image caption tasks with human sight annotation information. In addition, eye-tracking technology has been used in augmented reality [16], mixed reality [17], and deep learning [15]. These successful applications motivate our research of using sight to enable human-robot manipulation.

## III. METHOD

### A. System Overview

This section presents a friendly human-robot interaction by incorporating near-eye-tracking with robotic manipulation. The developed robotic assistive system is integrated with a head-mounted eye-tracking device so that the user's intention can be conveniently captured. That is, the user can use the sight to control a robotic arm to manipulate and grasp objects. Fig. 1 illustrates the pipeline of how the human intentions are perceived through eye gaze and translated to actions that the robot can execute. The method involves a combination of three steps: i) A head-mounted eye-tracking device integrated with low-cost stereo cameras measures the user's gaze direction. The biological model of human eyeball is incorporated with computer vision to identify the pupil orientation and locate the gaze coordinates of the eyes in 3D space. ii) In parallel, a hierarchical transformer visual model is developed to extract effective features for grasping, where the attention performs global perception. A feature pyramid inside the transformer network gathers features from each stage for multi-scale sensing in order to generate the final grasping configuration. iii) The information from the two subsystems is fused and the gaze point with human attention is filtered against the grasping



Fig. 1. The whole system pipeline. The pink section indicates the grasp detection subsystem and the green section is the eye-tracking module. The bottom yellow area is the fusion module, where objects of interest are selected for grasping based on the user's gaze.

quality heatmap to obtain the grasping pose parameters for the desired object.

### B. Eye-tracking

The human eye is a highly sophisticated optical instrument that allows light entering the eye to be imaged on the retina through a series of reflections and refractions. The refractive system of eyes consists of the cornea, vitreous humor, and lens, which forms multiple refractive surfaces at their interfaces due to the different refractive indices of each part, making the optical system extremely complex. In practice, to simplify the model, the refraction of the cornea during imaging of pupil is ignored and the cornea is modelled as a sphere with the same curvature at all points. We employ the eyeball model presented by [18], where the eye is modelled as two nested ellipsoids. The larger ellipsoid is the ocular and the smaller is the cornea. The corneal surface is modelled as a rotating ellipsoidal plane. The pupil is the channel for light to enter the eye and the direction of pupil directly indicates the rotation of eyeball and sight. The central pupil coordinate is the most important feature in sight tracking. Our pupil central localization is split into two steps: coarse localization and refined localization. The coarse localization mainly uses radical symmetry transformations to quickly locate points in the pupil area and eliminate invalid cases such as blink frames. Refined localization is carried out by edge extraction using Canny operation, followed by edge filtering and ellipse fitting of edge segments to obtain the pupil's precise position.

According to the coordinates of the pupil centre obtained by processing the eye image, based on the camera imaging calibration theory [19], the 3D coordinates of the pupil center satisfy the following relationship

$$
\begin{bmatrix} u_p \\ v_p \\ 1 \end{bmatrix} = \frac{1}{P_z} A \begin{bmatrix} P_x \\ P_y \\ P_z \end{bmatrix} = \frac{1}{P_z} \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} P_x \\ P_y \\ P_z \end{bmatrix},
$$
(1)

Fig. 2. Illustration of the near-eye-tracking.



Fig. 3. Overview of the transformer based grasp detection model.

where $[u_p, v_p]$ is the gaze point estimated on the near-eye image, $[P_x, P_y, 1]$ is the gaze point in the world coordinate system. $P_z$ and $A$ are the coordinate transformation matrix and the camera internal parameter matrix, respectively, which can be implemented by OpenCV [20]. To simplify the computation, the distance between the image plane and the camera is usually taken to be 1. Therefore, the 3D coordinates of the pupil are

$$
\begin{bmatrix} P_x \\ P_y \\ P_z \end{bmatrix} = \begin{bmatrix} P_x \\ P_y \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{u_p - c_x}{f_x} \\ \frac{u_p - c_y}{f_y} \\ 1 \end{bmatrix}. \tag{2}
$$

Similarly, the 3D coordinates of the centre of spot can also be found. Thus, the direction vectors $p_1, p_2$ are calculated by the connection between the camera optical center and the pupil center. The direction vectors $g_1, g_2$ are derived from the lines between the camera center and the corneal spot center. Since the ocular optic axis, reflected corneal spot, and the center of eye camera are coplanar, the normal vectors of the two planes shown in Fig. 2 are calculated through

$$
\begin{cases} \mathbf{n_1} = \frac{\mathbf{p_1} \times \mathbf{g_1}}{\|\mathbf{p_1} \times \mathbf{g_1}\|} \\ \mathbf{n_2} = \frac{\mathbf{p_2} \times \mathbf{g_2}}{\|\mathbf{p_2} \times \mathbf{g_2}\|}. \end{cases} \tag{3}
$$

The optical axis of the eye is the intersection of two planes, so the direction of the optical axis $l_{OA}$ can be obtained from the normal vectors of the two planes through equation $\frac{n_1 \times n2}{\|n_1 \times n2\|}$, after which the corresponding coordinates of the other positions of the optical axis can be calculated.

Our near-eye assistance system as a whole is shown in Fig. 5, where tiny eye cameras are installed on a head-mounted frame to capture high-resolution near-eye images. As shown in Fig. 2, for each eye, two near-eye cameras are used to capture the corneal reflected spot and calculate the corneal spherical center and 3D optical axis for the coordinates of the imaging.

At the same time, a scene image with the same view as the user is acquired by a scene camera mounted on the same frame. The system builds on the spherical cornea model and applies multiple near-eye cameras to resolve the corneal center and optical axis of the eyes in 3D space, detects the pupil center, and fits the line of sight. The line of sight is solved by analyzing the eyeball and applying optical principles. The pupil's 3D coordinates are obtained by detecting the centres of pupil in the image from near-eye cameras. Since the optical axis of eyes passes through the center of pupil and the sphere of the cornea, the line of sight can be obtained by calculating the line between the two coordinates based on the geometric model of the eye [18].

### C. Grasp Detection

Compared to object detection, the grasp detection is generally made of small rectangles and is more sensitive to positions and rotation angles. To provide a more global understanding of the model, as opposed to convolutional kernels with fixed receptive fields, stacked transformer layers are adopted as backbone to gradually extract coarse-grained to fine-grained feature representations. As shown in Fig. 3, the input image $\mathcal{I} \in \mathbb{R}^{W \times H}$, where W and H are the width and height of the image $\mathcal{I}$. At first, the image is split into non-overlapping patches through a conv projection layer. Each patch in the image is treated as a word token. Similar to [21], there are four successive blocks to extract semantic rich information, and each stage contains a patching merging layer and a swin transformer layer. Each block is composed of patch merging and swin transformer layer. Patching merging is functionally similar to the pooling mechanism in CNN, which is designed to reduce the resolution of an image while increasing the number of channels of features. The dimension of features of each stage is shown in Fig. 3. The foundation of swin transformer layer is the attention mechanism. Input features

are linearly transformed to obtain query, key, and value. Then the self-attention is computed as

$$\text{Attention}(Q, K, V) = \text{SoftMax}(\frac{QK^T}{\sqrt{d}})V, \qquad (4)$$

where $\sqrt{d}$ is the scale factor. Swin transformer layer performs self-attention within a local window, greatly reducing the computational complexity. Meanwhile, the shifted window is applied to model the global relationships. The whole computation flow is as follows:

$$
\begin{aligned}
\hat{\mathbf{x}}^l &= \text{W-MSA}\left(\text{LN}\left(\mathbf{x}^{l-1}\right)\right) + \mathbf{x}^{l-1}, \\
\mathbf{x}^l &= \text{MLP}\left(\text{LN}\left(\hat{\mathbf{x}}^l\right)\right) + \hat{\mathbf{x}}^l, \\
\hat{\mathbf{x}}^{l+1} &= \text{SW-MSA}\left(\text{LN}\left(\mathbf{x}^l\right)\right) + \mathbf{x}^l, \\
\mathbf{x}^{l+1} &= \text{MLP}\left(\text{LN}\left(\hat{\mathbf{x}}^{l+1}\right)\right) + \hat{\mathbf{x}}^{l+1}
\end{aligned}
\qquad (5)
$$

The feature $\mathbf{x}^{l-1}$ from last layer enters the W-MSA module via the layerNorm layer, and there exists a residual connection between each module. After that, it goes through SW-MSA layer in a similar way. Here W-MSA represents the window multi-head attention layer and SW-MSA indicates shifted window multi-head attention layer. One motivation for using swin transformer as backbone network is that it maintains both global and local perception. Meanwhile, it reduces computational complexity compared to vanilla self-attention.

A feature fusion pyramid is designed at the bottom of model in Fig. 3 to collect features from each stage in the backbone network. Features from different stages are aggregated through a feature pyramid for a multiscale fusion of contextual information. The feature fusion module uses concatenation to fuse these features that learn semantic and spatial contextual information. The network finally outputs the grasping quality head, grasping width head, and grasping angle head via $1 \times 1$ convolutional kernels. Each output head is the same size as the original input. For grasping quality head, each parameter in the quality output takes a value between 0 and 1, indicating the probability of a successful grasp in the corresponding position in the image. The angle head consists of two parts: $cos2\theta$ and $sin2\theta$, and the resulting angle is calculated by $\frac{1}{2}arctan\frac{cos2\theta}{sin2\theta}$. Afterwards, all outputs of grasping quality head are searched for the point with the highest grasping quality as the grasping center, as well as its corresponding gripper rotation angle and width.

The loss function $\mathcal{L}$ is defined as $\mathcal{L} = w_1\mathcal{L}_{pose} + w_2\mathcal{L}_{angle} + w_3\mathcal{L}_{pose}$. For each component of the loss function, $\mathcal{L}_i$ is the mean square error between the corresponding value of the model and the ground truth. $w1, w2$, and $w3$ are the relevant weight factors. For instance, the first term of $\mathcal{L}$ is defined as $\mathcal{L}_{pose} = \sum_{i=1}^{N} \|\tilde{G}_i - G_i^*\|^2$, where $\tilde{G}_i$ is the output of the grasp quality head and $G_i^*$ is the corresponding ground truth.

## IV. EXPERIMENT

### A. Dataset and Implementation Details

The Cornell [10] grasping dataset is used to evaluate the effectiveness of our grasp detection model. Each image in the dataset has been taken with a center crop of $224 \times 224$. The full grasp detection model is implemented in Pytorch, running on a single NVIDIA GTX 3090 GPU. The batch size is set to 32 and we use AdamW optimizer with a learning rate of $1e$-4.

**Evaluation Criteria.** Following the standards in [10], [8], [22], the grasping rectangle metric is used to evaluate the grasping results. A predicted grasp is treated as positive if it meets the following two criteria.

i) The discrepancy between the rotation angle of the predicted grasp and the ground truth does not exceed $30°$.

ii) The Jaccard index, defined in (6), of the predicted grasp and the ground truth is greater than $0.25$,

$$J\left(\mathcal{R}^*, \mathcal{R}\right) = \frac{|\mathcal{R}^* \cap \mathcal{R}|}{|\mathcal{R}^* \cup \mathcal{R}|}, \qquad (6)$$

where $\mathcal{R}$ is the predicted grasping rectangle region and $\mathcal{R}^*$ is the ground truth. $\mathcal{R}^* \cap \mathcal{R}$ is an intersection of these two areas and $\mathcal{R}^* \cup \mathcal{R}$ is the union of these two regions.



Fig. 5. The developed head-mounted eye-tracking device.

### B. Grasping Results

Given the desired grasping coordinates, the robot inverse kinematics are utilized to realize the desired trajectory. The detailed comparisons with other methods on the Cornell dataset are listed in Table II. Although the performance gap among the state-of-the-art methods is moderate, our model achieves the best performance. For instance, our transformer-based grasping model achieves an accuracy of 96.28 if only depth images are used as input, and 98.86 for the RGB-D as input. In particular, our model directly predicts the grasping quality, angle, and width of grasping rectangles, obviating the requirements for design anchors for different targets.

To test whether our model can be generalized to new scenes, objects are arranged in new positions with different orientations. We divide the objects into three categories, including objects that appear in the dataset, objects that are similar in the dataset, and objects that have never been seen before. In each category, there are at least four objects. Each type of object is grasped several times, and the number of successful grasps are recorded. The detailed grasping results are shown in Table I. In Table I, we can see that it performs well for tools seen in the dataset, and shows good generalization to similar objects. For unseen objects and complex scenes, our method also provides a decent improvement in grasp detection accuracy.

(a) Visual attention        (b) Gazed-based physical grasping

Fig. 4. (a) Visualization of the scene images and the human attention heatmaps captured by our proposed method. From left to right, the human gaze points are the blue bottle, eyeglass box, banana, and apple, respectively. (b) Experiment for real-time gaze-based grasping.

TABLE I
EXPERIMENTAL ROBOTIC GRASPING SUCCESS RATE FOR DIFFERENT OBJECTS.

| Seen Objects | | | Familiar Objects | | | Unseen Objects | | |
|---|---|---|---|---|---|---|---|---|
| Objects | Detected | grasp (%) | Objects | Detected | grasp (%) | Objects | Detected | grasp (%) |
| Mouse | 15 / 15 | 13 / 15 | Orange | 14 / 15 | 12 / 15 | Scissor | 12 / 15 | 11 / 15 |
| Remote Control | 15 / 15 | 13 / 15 | Staples Box | 15 / 15 | 12 / 15 | Toothpaste Box | 12 / 15 | 12 / 15 |
| Apple | 14 / 15 | 12 / 15 | Knife | 12 / 15 | 11 / 15 | Razor | 13 / 15 | 11 / 15 |
| Pencil | 14 / 15 | 13 / 15 | Screwdriver | 12 / 15 | 14 / 15 | Toy | 12 / 15 | 9 / 15 |
| Average | 96 % | 85 % | Average | 88 % | 81.6 % | Average | 81.6 % | 71.6 % |

TABLE II
THE ACCURACY ON CORNELL GRASPING DATASET.

| Authors | Algorithm | Accuracy (%) |
|---|---|---|
| Jiang [23] | Fast Search | 60.5 |
| Lenz [10] | SAE, struct. reg. | 73.9 |
| Redmon [11] | AlexNet, MultiGrasp | 88.0 |
| Wang [22] | Two-stage closed-loop | 85.3 |
| Asif [24] | STEM-CaRFs | 88.2 |
| Kumra [8] | ResNet-50x2 | 89.2 |
| Morrison [12] | GG-CNN | 73.0 |
| Guo [25] | ZF-net | 93.2 |
| Zhou [26] | FCGN, ResNet-101 | 97.7 |
| Karaoguz [27] | GRPN | 88.7 |
| Asif [28] | GraspNet | 90.2 |
| | GraspFormer-D | 96.28 |
| Our | GraspFormer-RGB | 97.72 |
| | GraspFormer-RGB-D | **98.86** |

### C. System Design

The system consists of two main components, the eye-tracking module and grasping module. The robot employed in our experiments is a Franka Emika Panda robot. An RGB-D camera is fixed on the robot's gripper. The used camera is RealSense D435i. The panda robot has a parallel finger and its active range is limited to 10cm with a maximum loading capacity of no more than 3kg.

The eye-tracking hardware includes four eye cameras associated with infrared light source, one scene camera, multi-channel video acquisition circuits, and a head-mounted frame. The near-eye cameras are positioned below the human eyes to take high-resolution images of eyes, and the near-infrared light around the camera provides a corneal reflective spot used to accurately resolve the 3D eye features under the assumption of an aspheric corneal model.

### D. Limitations

Since properties such as corneal refraction and corneal asphericity can introduce additional parameters which are hard to calibrate, it is challenging to directly solve for pupil centre coordinates by optical principles. To mitigate the challenge, a simplified and approximated eye model is used in our eye-tracking module. Specifically, the refraction of the cornea is ignored, whereas the refractive index of human cornea is approximately 1.376. The curvature of the corneal surface is also overlooked, which is modelled as a sphere in our model. Note that in reality the curvature is not constant at all points on the surface of the eyes. Such approximations can lead to discrepancies in the gaze estimation. In the grasping subsystem, our model does not perform well when grasping transparent objects, since the RealSense camera does not yet provide a good depth for such objects. It is found in the experiment that the objects with complex surfaces or smooth materials are likely to slip out of the grippers during grasping.

### V. CONCLUSIONS AND DISCUSSIONS

TABLE III
TIME DELAY ANALYSIS OF GAZE INTERACTION-BASED GRASPING ASSISTIVE SYSTEM.

| Setup | | Time Delay |
|---|---|---|
| Image Acquisition | | 0.01 ms |
| Image Preprocessing | | 0.2 ms |
| Gaze Point Estimation | | 3 ms |
| grasp detection | GraspFormer Tiny | 47 ms |
| | GraspFormer Small | 56 ms |
| | GraspFormer Base | 88 ms |
| Total Time | GraspFormer Tiny | 50.21ms |
| | GraspFormer Small | 59.21ms |
| | GraspFormer Base | 91.21ms |

Fig. 6. Images of the eye captured by near-eye cameras.

In this work, we provide a contactless human-robot interface that enables robotic manipulation using sight. The proposed framework leverages a head-mounted eye-tracking device to automatically locate the object that human pays attention to. Once the user's gaze point is obtained, a transformer-based grasp model takes advantages of global perception to identify the user's region of interest. The results show that the developed gaze-based robotic arm is capable of moving objects or grasping desired objects by near-eye-tracking. Ablation studies demonstrate that our eye-tracking can achieve reasonably decent tracking accuracy. The associated head-mounted device provides a low-cost design and meets real-time requirements.

Compared to using other human-machine interfaces, our eye interaction is more flexible and user-friendly. Eye-driven human-robot interaction serves as a novel framework and shows great potential in various applications. Future research will consider extending the current work to more challenging environments (e.g., outdoor environments) for more complex tasks (e.g., pouring water, serving food).

## REFERENCES

[1] J. B. Sol, "Effective grasping enables successful robot-assisted dressing," *Sci. Robotics*, vol. 7, no. 65, 2022.

[2] S. Wang, Z. Zhou, and Z. Kan, "When transformer meets robotic grasping: Exploits context for efficient grasp detection," *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 8170–8177, 2022.

[3] R. Rahman, M. S. Rahman, and J. R. Bhuiyan, "Joystick controlled industrial robotic system with robotic arm," in *IEEE International Conference on Robotics, Automation, Artificial-intelligence and Internet-of-Things (RAAICON)*, 2019, pp. 31–34.

[4] F. Sun, W. Zhang, J. Chen, H. Wu, C. Tan, and W. Su, "Fused fuzzy petri nets: A shared control method for brain-computer interface systems," *IEEE Trans. Cogn. Dev. Syst.*, vol. 11, no. 2, pp. 188–199, 2019.

[5] L. D. Rosenblum, *See what I'm saying: The extraordinary powers of our five senses.* WW Norton & Company, 2011.

[6] Y.-S. L.-K. Cio, M. Raison, C. L. Ménard, and S. Achiche, "Proof of concept of an assistive robotic arm control using artificial stereovision and eye-tracking," *IEEE Trans. Neural Syst. and Rehabilitation Engineering*, vol. 27, no. 12, pp. 2344–2352, 2019.

[7] S. S. Groothuis, S. Stramigioli, and R. Carloni, "Lending a helping hand: toward novel assistive robotic arms," *IEEE Robot. Autom. Magazine*, vol. 20, no. 1, pp. 20–29, 2013.

[8] S. Kumra and C. Kanan, "Robotic grasp detection using deep convolutional neural networks," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst*, 2017, pp. 769–776.

[9] S. Wang, Z. Zhou, H. Wang, Z. Li, and Z. Kan, "Unsupervised representation learning for visual robotics grasping," in *International Conference on Advanced Robotics and Mechatronics (ICARM)*, 2022, pp. 57–62.

[10] I. Lenz, H. Lee, and A. Saxena, "Deep learning for detecting robotic grasps," *Int. J. Robotics Res.*, vol. 34, no. 4-5, pp. 705–724, 2015.

[11] J. Redmon and A. Angelova, "Real-time grasp detection using convolutional neural networks," in *Proc. IEEE Int. Conf. Robot. Autom.*, 2015, pp. 1316–1322.

[12] D. Morrison, P. Corke, and J. Leitner, "Learning robust, real-time, reactive robotic grasping," *Int. J. Robotics Res.*, vol. 39, no. 2-3, pp. 183–201, 2020.

[13] N. Hollenstein and C. Zhang, "Entity recognition at first sight: Improving ner with eye movement information," in *NAACL-HLT (1)*, 2019, pp. 1–10.

[14] N. Karessli, Z. Akata, B. Schiele, and A. Bulling, "Gaze embeddings for zero-shot image classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 4525–4534.

[15] Y. Sugano and A. Bulling, "Seeing with humans: Gaze-assisted neural image captioning," *arXiv preprint arXiv:1608.05203*, 2016.

[16] J.-Y. Lee, H.-M. Park, S.-H. Lee, T.-E. Kim, and J.-S. Choi, "Design and implementation of an augmented reality system using gaze interaction," in *Int. Conf. on Information Science and Applications*, 2011, pp. 1–8.

[17] F. Bruno, L. Barbieri, and M. Muzzupappa, "A mixed reality system for the ergonomic assessment of industrial workstations," *Int. Journal on Interactive Design and Manufacturing (IJIDeM)*, vol. 14, no. 3, pp. 805–812, 2020.

[18] T. Nagamatsu, Y. Iwamoto, J. Kamahara, N. Tanaka, and M. Yamamoto, "Gaze estimation method based on an aspherical model of the cornea: surface of revolution about the optical axis of the eye," in *Proc. of Symposium on Eye-Tracking Research & Applications*, 2010, pp. 255–258.

[19] Z. Zhang, "A flexible new technique for camera calibration," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 22, no. 11, pp. 1330–1334, 2000.

[20] G. Bradski and A. Kaehler, "Opencv," *Dr. Dobb's journal of software tools*, vol. 3, p. 120, 2000.

[21] L. Z. et al., "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proc. IEEE Int. Conf. Comput. Vision*, 2021, pp. 10 012–10 022.

[22] Z. Wang, Z. Li, B. Wang, and H. Liu, "Robot grasp detection using multimodal deep convolutional neural networks," *Advances in Mechanical Engineering*, vol. 8, no. 9, p. 1687814016668077, 2016.

[23] Y. Jiang, S. Moseson, and A. Saxena, "Efficient grasping from rgbd images: Learning using a new rectangle representation," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2011, pp. 3304–3311.

[24] U. Asif, M. Bennamoun, and F. A. Sohel, "Rgb-d object recognition and grasp detection using hierarchical cascaded forests," *IEEE Trans. on Robotics*, vol. 33, no. 3, pp. 547–564, 2017.

[25] D. Guo, F. Sun, H. Liu, T. Kong, B. Fang, and N. Xi, "A hybrid deep architecture for robotic grasp detection," in *Proc. IEEE Int. Conf. Robot.Automat.*, 2017, pp. 1609–1614.

[26] X. Zhou, X. Lan, H. Zhang, Z. Tian, Y. Zhang, and N. Zheng, "Fully convolutional grasp detection network with oriented anchor box," in *Proc. IEEE Int. Conf. Intell. Robots Syst.*, 2018, pp. 7223–7230.

[27] H. Karaoguz and P. Jensfelt, "Object detection approach for robot grasp detection," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2019, pp. 4953–4959.

[28] U. Asif, J. Tang, and S. Harrer, "Graspnet: An efficient convolutional neural network for real-time grasp detection for low-powered devices." in *IJCAI*, vol. 7, 2018, pp. 4875–4882.