

MATERobot: Material Recognition in Wearable Robotics for People with Visual Impairments

Junwei Zheng^{1,*}, Jiaming Zhang^{1,2,*}, Kailun Yang^{3,†}, Kunyu Peng¹, and Rainer Stiefelhagen¹

Abstract—People with Visual Impairments (PVI) typically recognize objects through haptic perception. *Knowing objects and materials before touching* is desired by the target users but under-explored in the field of human-centered robotics. To fill this gap, in this work, a wearable vision-based robotic system, MATERobot, is established for PVI to recognize materials and object categories beforehand. To address the computational constraints of mobile platforms, we propose a lightweight yet accurate model MATEViT to perform pixel-wise semantic segmentation, simultaneously recognizing both objects and materials. Our methods achieve respective 40.2% and 51.1% of mIoU on COCOStuff-10K and DMS datasets, surpassing the previous method with +5.7% and +7.0% gains. Moreover, on the field test with participants, our wearable system reaches a score of 28 in the NASA-Task Load Index, indicating low cognitive demands and ease of use. Our MATERobot demonstrates the feasibility of recognizing material property through visual cues and offers a promising step towards improving the functionality of wearable robots for PVI. The source code has been made publicly available at [MATERobot](#).

I. INTRODUCTION

In 2020, it was estimated that approximately 43 million individuals were living with blindness. By 2050, this number is expected to rise significantly to 61 million [1]. Therefore, it is imperative to facilitate the development of assistive systems for helping PVI. In recent years, considerable progress has been witnessed in the field of human-centered assistive technology, such as systems for navigation [2], object localization [3], indoor understanding [4], and path orientation [5].

Material recognition is often a challenging task for PVI, who typically recognize objects through touch [6], [7]. However, material recognition is under-explored in the domain of assistive technology. Therefore, it is essential to develop a wearable system that can assist PVI in recognizing the materials of objects without touching them, *i.e.*, a contact-free material recognition system.

This work was supported in part by the Ministry of Science, Research and the Arts of Baden-Württemberg (MWK) through the Cooperative Graduate School Accessibility through AI-based Assistive Technology (KATE) under Grant BW6-03, in part by BMBF through a fellowship within the IFI programme of DAAD, in part by the Helmholtz Association Initiative and Networking Fund on the HAICORE@KIT and HOREKA@KIT partition, and in part by Hangzhou SurImage Technology Company Ltd.

¹The authors are with the Institute for Robotics and Anthropomatics, Karlsruhe Institute of Technology, Karlsruhe 76131, Germany.

²The author is also with the Department of Engineering Science, University of Oxford, Oxford OX1 3PJ, UK.

³The author is with the School of Robotics and the National Engineering Research Center of Robot Visual Perception and Control Technology, Hunan University, Changsha 410082, China.

*Equal contribution.

†Corresponding author: Kailun Yang. (E-mail: kailun.yang@hnu.edu.cn.)

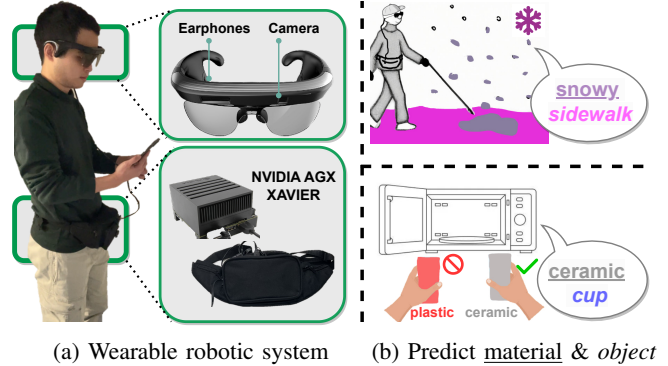


Fig. 1: MATERobot, (a) wearable robotics, can assist (b) material semantic segmentation (e.g., snowy, ceramic) and general *object* semantic segmentation (e.g., *sidewalk*, *cup*).

In this work, we present the wearable **MATERobot** system to cover both object and material recognition. As shown in Fig. 1a, the system consists of a pair of smart glasses with an RGB-D camera, a pair of bone-conducting earphones, and a portable image processor inside a small bag. In Fig. 1b, the MATERobot can recognize materials (indicated by underline, e.g., snowy) and general *objects* (indicated in italics, e.g., *sidewalk*). These two predictions can form feedback of object categories with material properties (e.g., “snowy sidewalk”). This information is conveyed to the blind user in the form of speech through bone-conduction headphones.

Due to computationally complex self-attention operations in vision transformers [8], it is, however, hard to deploy a resource-intensive ViT-based model on portable platforms. To address this, we propose an efficient material segmentation ViT-based model, namely MATEViT, which includes a *Learnable Importance Sampling (LIS)* strategy to maintain only the informative tokens for the material segmentation, so as to reduce the computational cost. Thanks to the LIS strategy, a resource-friendly MATEViT model is obtained, enabling the deployment of vision transformer on wearable robotic devices that have limited computational resources.

Apart from the model efficiency, to enlarge the model capacity, we introduce a *Multi-gate Mixture-of-Experts (MMoE)* method to combine the aforementioned general image semantic segmentation and material semantic segmentation on a single model. Compared to the previous method [5] using a straightforward dual-head structure, our MMoE is used to construct a high-performance and efficient multi-task learning architecture. The feature tokens from input images are forwarded to respective gates and experts to

extract informative features for task-relevant decoder heads that generate final semantic segmentation masks.

In order to endow MATERobot with robust perception capability, including general object and material semantic segmentation, we perform model training on COCOStuff-10K [9] and DMS [10] datasets, both contain more than 10K training samples. Through extensive experiments, our small model obtains 40.2% and 51.1% of mIoU scores, surpassing the previous multi-task learning baseline [5] by absolute +5.7% and +7.0% on COCOStuff-10K and DMS, respectively. For single-task learning on materials segmentation, *i.e.*, on DMS, our model reaches state-of-the-art performance, having +8.1% mIoU gains compared to previous CNN counterpart [10]. To verify the practicability of our MATERobot for recognizing material categories in real-world scenarios, we conduct a task-oriented user study with six blindfolded participants. On the post-study questionnaires, our system obtains respective 28 and 77 scores regarding the evaluation criteria of NASA-Task Load Index (NASA-TLX) and System Usability Scale (SUS), which indicates the ease of use and the usability of our MATERobot in practical scenarios. In summary, our main contributions are:

- For the first time, we integrate material recognition into assistive technology and built a wearable robot system, *i.e.*, MATERobot. This system empowers PVI to achieve contactless long-distance perception similar to that of sighted people.
- We propose an efficient MATEViT model to enable the deployment of resource-intensive ViT-based counterparts on resource-constrained mobile platforms by using a LIS strategy.
- A MMoE method is designed to simultaneously perform object and material semantic segmentation in one unified model.
- Through a task-oriented user study and qualitative analyses, we gain valuable insights into designing wearable material recognition systems for PVI.

II. RELATED WORK

A. Wearable Assistive System

With the tremendous capability revealed by computer vision algorithms, vision-based wearable assistance systems [11]–[15] are becoming increasingly applicable. A vision-based navigation system [2], calculating precise positions and orientations, is proposed to help PVI stay on track while walking, and it can recognize unexpected dynamic obstacles. Due to the COVID-19 pandemic, an object-finding algorithm is introduced in [3] to build a robotic cane system. In [4], a lightweight system with a solid-state LiDAR sensor is proposed for indoor detection and avoidance. In a previous work [5], [16], to cover the segmentation of transparent objects, a dual-head model is deployed on a wearable device. However, only limited recognizable materials are delivered in previous wearable assistance systems. In order to help blind users obtain a more comprehensive and humanized experience on material recognition, we design a wearable

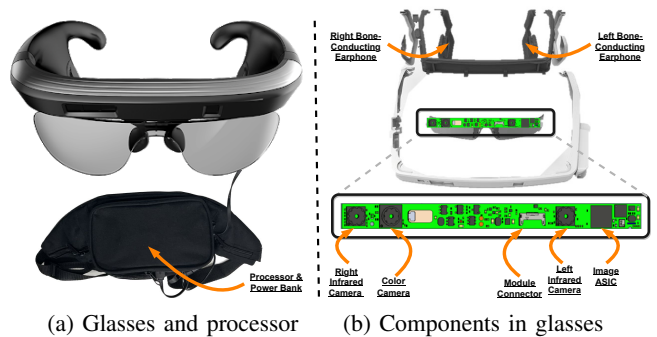


Fig. 2: Hardware components in the MATERobot system.

robotic system in this work, for the first time, which can not only recognize conventional object categories, such as *cups*, but also further recognize the material of the object, such as *plastic cups*, delivering contactless object recognition.

B. Material Semantic Segmentation

Recently, Vision Transformer [8], [17] is proposed to utilize the self-attention operation in transformer layers to extract non-local features from a sequence of image patches, yielding an alternative backbone solution compared to convolutional counterparts [18]–[20]. In DMS [10], a model based on ResNet [20] is used to address dense material segmentation. In contrast to DMS [10], we adopt Vision Transformer as the backbone, *e.g.*, ViT [8], to perform material semantic segmentation, with the aim of extracting long-distance dependencies between image patches, since the long-range contextual information is crucial for robust representation of diverse materials [5]. However, due to the high computational demand of self-attention operation, there is still a bottleneck when deploying a plain vision transformer on resource-constrained mobile platforms, *e.g.* mobile robots and wearable devices. In this work, we propose a novel importance sampling method to reduce the number of tokens and maintain only the informative ones for material segmentation to enable the deployment of plain vision transformers on wearable devices.

III. MATEROBOT: A WEARABLE ROBOTIC SYSTEM

A. Hardware Component

As shown in Fig. 2, there are three main hardware components in our MATERobot, including a pair of KRVision smart vision glasses, a portable processor, and a power bank inside a waist bag. Inside the smart glasses, there is an RGB-Depth camera RealSense R200 and a pair of bone-conduction headphones. For the concept of human-friendly design, there are three advantages of using bone-conduction earphones, which are comfortable to wear, clean and hygienic, and keep in touch with the outside world. Maintaining awareness of ambient sounds is especially important for PVI. To ensure higher portability of the system, we tried different processors and chose the smaller NVIDIA AGX Xavier due to its applicable energy efficiency and inference capabilities. Furthermore, a power bank with high energy capacity is selected to provide the system with up to 6 hours of battery life, which

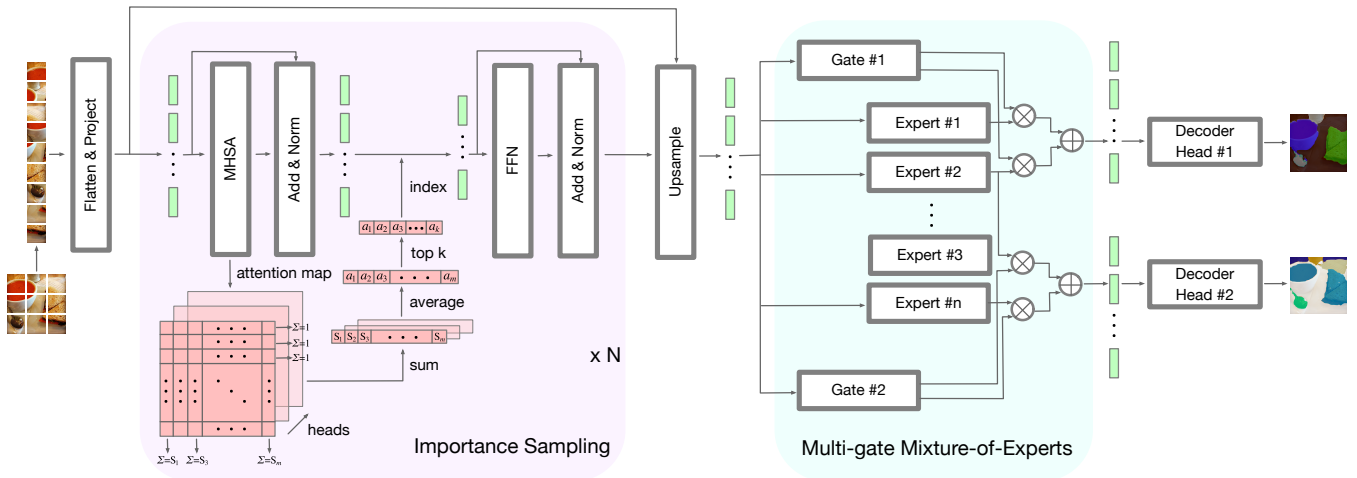


Fig. 3: **Architecture of MATEViT** in a Learnable Importance Sampling (LIS) strategy to reduce the computational complexity, and with a Multi-gate Mixture-of-Experts (MMoE) layer to perform dual-task segmentation (*i.e.*, #1-Object and #2-Material segmentation).

greatly reduces the battery life anxiety of traditional wearable devices [4]. Through the above hardware components, a more portable MATERobot and better user experience can be delivered to PVI when performing contactless material recognition in real-world scenarios.

B. User Interaction

Between the system and the user, we design an easy-to-use interface for PVI. First of all, in order to give users timely information feedback, we adopt an adjustable feedback frequency. For example, if users set a larger interval in a more familiar environment, *e.g.*, home, they will get concise information. If they explore an unknown space, setting a high frequency can obtain more information. Besides, between the pixel-wise image segmentation results to the auditory output to users, only the detected object and material located in the middle of the input image will be selected. Their pre-defined texts are used to generate speeches jointly in the form of “material objects”, such as “ceramic cups” or “metal forks”.

C. MATEViT Model

Apart from the hardware components, we further detail the efficient model designed for MATERobot. To equip mobile platforms with plain vision transformers, we propose MATEViT, which has ViT [8] with LIS as the backbone, followed by an upsampling layer, a MMoE layer and two decoder heads. The architecture is shown in Fig. 3. Note that our method is flexible to include more tasks by adding decoder heads. Mixed by two datasets, each input sample is first encoded by the efficient ViT backbone and then fed into the upsampling layer. One gate in the MMoE layer corresponds to one task, receiving the data sample that is only relevant to the task. Depending on the output of the gate, different selected experts are activated to learn meaningful latent representations, which are finally decoded by the task-relevant decoder head. Different from the training process, one data sample is fed into all gates synchronously during inference, and latent representations from selected experts

are then decoded by expert corresponding decoder heads, producing predictions for all tasks.

D. Learnable Importance Sampling

To make plain ViT [8] more lightweight and feasible in real-world applications, we propose a *Learnable Importance Sampling* strategy for material semantic segmentation, yielding an efficient MATEViT. Our approach does not require an additional class token in forward pass compared to EViT [21]; therefore, it further reduces the model complexity with high-resolution inputs. As illustrated in Fig. 3, all image patches are flattened and projected into tokens, forming Queries (\mathbf{Q}), Keys (\mathbf{K}), and Values (\mathbf{V}). Multi-Head Self-Attention (MHSA) is then calculated as:

$$MHSA(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{Softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{C}}\right)\mathbf{V} \quad (1)$$

where $\mathbf{Q} \in \mathbb{R}^{N \times C}$, $\mathbf{K} \in \mathbb{R}^{N \times C}$ and $\mathbf{V} \in \mathbb{R}^{N \times C}$ are query, key and value matrices; N , C are token number and dimension. Since softmax is introduced in attention map calculation, the summation of each value in rows is equal to 1. However, the result does not always equal to 1 when summing up all values in columns, indicating the importance of each token from an image. Based on this observation, we first calculate the summation in columns and then average the importance vectors among all heads. A fixed number of top k values are selected and tokens are downsampled according to the indices of these k values after *Add & Norm*, which stand for a residual link [20] and layer normalization [22]. Fig. 3 illustrates the whole LIS process in detail. Following [23], the downsampled tokens are then sent to a Feed-Forward Network (FFN) followed by another *Add & Norm*. Since the token number is reduced after the importance sampling, an upsampling layer is then applied, which is a standard transformer decoder block [23]. Thanks to LIS, a plain ViT-based model is expedited and better qualified for real-world mobile applications.

E. Multi-gate Mixture-of-Experts

Since the performance of a wearable robotic system is relevant to its model capacity, Mixture-of-Experts [24] is proposed for enlarging model capacity while maintaining invariant model complexity. However, the usage of MoE is less discussed on wearable systems [5]. For the first time, we adopt the MoE method to perform complementary training of both general object and material segmentation, *i.e.*, the former prevents the latter from overfitting and vice versa. More specifically, we adopt a MMoE layer in our model for high efficiency and performance. Fig. 3 shows the detail of the MMoE layer for multi-task learning. Specifically, one gate is responsible for one task to select the experts. Similar to [24], the resulting selection vector $G(\mathbf{x})$ can be described as:

$$G(\mathbf{x}) = \text{Softmax}(\text{TopM}(H(\mathbf{x}), m)) \quad (2)$$

$$H(\mathbf{x}) = \mathbf{x} \cdot \mathbf{W}_g + N(\mathbf{x}) \quad (3)$$

$$N(\mathbf{x}) = \text{Normal}() \cdot \text{Softplus}(\mathbf{x} \cdot \mathbf{W}_{noise}) \quad (4)$$

$$\text{TopM}(\mathbf{v}, m)_i = \begin{cases} v_i & \text{if } v_i \text{ ranks top } m \\ -\infty & \text{otherwise} \end{cases} \quad (5)$$

where \mathbf{x} , \mathbf{W}_g , and \mathbf{W}_{noise} are token, gate matrix, and noise matrix, respectively. The noise term $N(\mathbf{x})$ is utilized for load balance, which is the multiplication of the standard normal distribution sampling $\text{Normal}()$ and the outcome of $\text{Softplus}(\cdot)$. During the training, for every token in one image, only one same gate is activated to produce the selection vector. According to the indices of the top m values, m experts are selected and the token is only fed into the m experts. The output of the MMoE layer is a weighted sum of the top m values in the selection vector and their corresponding outcomes from the m experts. A task-relevant decoder head [25] is then applied to transform all output tokens from MMoE into a prediction mask. Note that we also employ the load and importance balancing loss following [24] besides the task loss. During the inference, every token from one image is fed into all gates synchronously. The resulting weighted sums from selected experts are then decoded by the corresponding decoder heads shown in Fig. 3. As the model is trained jointly on two tasks, the knowledge absorbed from object segmentation plays an important role in the high performance of material recognition, which provides PVI with accurate material information in their daily life.

IV. EXPERIMENTS

A. Settings and Datasets

Settings. We implement the model with PyTorch 1.12.1 and CUDA 11.6. The learning rate is initialized as 0.01 and it is scheduled by a cosine annealing strategy [26]. SGD with momentum 0.9 is adopted as the optimizer. We initialize our efficient ViT backbone with a pre-trained plain ViT [8] and keep other layers of the model randomly initialized. Data augmentations like random horizontal flipping, random resize with a ratio 0.5-2, and random cropping to 512×512 are used during training. Note that we **do not** use other tricks such

TABLE I: Results of single-task learning on DMS dataset val/test set. GFLOPs are measured in size of 512×512 .

Method	Backbone	GFLOPs	MParams	Pixel Acc (%)	mIoU (%)
PSPNet [19]	MobileNetV2	06.12	02.13	66.5 / 66.3	26.1 / 25.9
DeepLabV3 [28]	MobileNetV2	07.52	02.58	68.4 / 68.1	29.8 / 29.7
DeepLabV3+ [18]	MobileNetV2	07.85	03.11	69.3 / 69.2	30.1 / 29.9
LR-ASPP [29]	MobileNetV3-s	05.41	01.14	66.8 / 66.5	26.5 / 26.4
LR-ASPP [29]	MobileNetV3	08.78	03.29	70.7 / 70.2	30.3 / 29.9
SeaFormer [30]	SeaFormer-L	06.50	14.00	73.1 / 72.8	43.4 / 41.8
DMS [10]	ResNet-50	-	-	73.1 / 72.9	43.5 / 42.0
MATEViT (ours)	ViT-Tiny	04.30	07.78	76.9 / 76.8	45.3 / 44.1
MATEViT (ours)	ViT-Small	15.54	28.79	79.6 / 79.4	51.0 / 50.1
<i>w.r.t.</i> DMS				+6.5 / +6.5	+7.5 / +8.1

as OHEM, auxiliary loss, and class-weighted loss for a fair comparison to other methods. We train our model with a batch size of 4 for 200 epochs on four 1080Ti GPUs.

Datasets. We adopt COCOStuff-10K [9] and DMS [10] for general object and material segmentation, respectively. The COCOStuff-10K dataset [9] has 9000/1000 images for training/testing. We conduct experiments following the implementation of mmsegmentation [27] with 171 categories. The DMS dataset [10] has respective 21857/9057/9152 images for training/validation/testing with 46 categories.

B. Results on Material Segmentation

To verify the proposed method for material segmentation, we conduct experiments of material segmentation on the DMS dataset [10]. Results are reported in Table I. It can be observed that our model using the ViT-Tiny backbone still has the lowest computation expense (4.30 GFLOPs) with high performance (44.1% in mIoU), and the ViT-Small variant outperforms other methods in both pixel accuracy and mIoU. Furthermore, Fig. 4 presents the per-class IoU of all material categories. It is worth noting that our ViT-Small variant achieves performance gains in all 46 categories, especially those are relevant for assisting PVI, *e.g.*, *fire* (gain +20.2%), *snow* (gain +21.3%), *plastic* (gain +22.6%), and *ceramic* (gain +25.5%). The impressive pre-study test scores obtained in our evaluations serve as a testament to the effectiveness and reliability of our proposed system in assisting PVI in accurately recognizing and distinguishing materials.

C. Results on Multi-task Learning

To deliver more information to PVI and to perform complementary training, we further conduct experiments on multi-task learning based on MATEViT, covering both object and material segmentation at once. Table II illustrates the quantitative results. Compared to Trans4Trans MiT-B0 [5], our model with ViT-Tiny requires less computation expense (-35.0% GFLOPs), while reaching higher performance on both datasets. Additionally, compared to the Trans4Trans MiT-B2 variant [5], it becomes evident that our ViT-Small variant has a higher mIoU on both COCOStuff-10K (gain +5.7%) and DMS (gain +7.0%). More important, lower computational complexity (*i.e.*, GFLOPs) is intuitive to reflect the high efficiency of the model running on the mobile

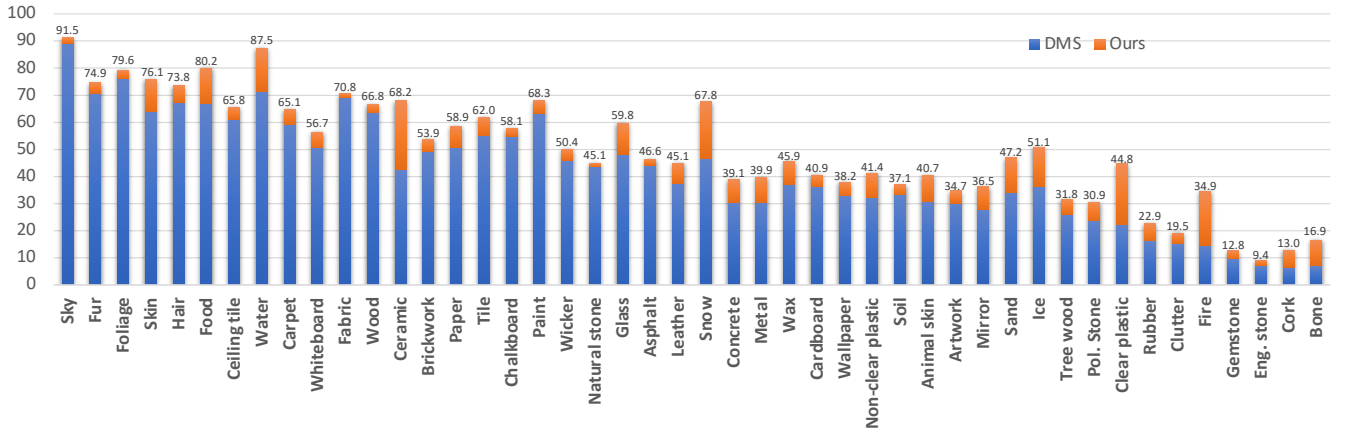


Fig. 4: **Per-class IoU (%) of all material categories.** The blue bar shows the category IoU (%) of the baseline DMS [10], while the orange part shows the performance gains (%) of our proposed method.

TABLE II: **Results (mIoU) of multi-task learning on COCOStuff-10k and DMS test sets.** GFLOPs @ 512×512.

Method	Backbone	GFLOPs	MParams	COCOStuff (%)	DMS (%)
Trans4Trans [5]	MiT-B0	12.24	04.16	27.7	37.1
Trans4Trans [5]	MiT-B1	18.98	14.26	30.6	41.3
Trans4Trans [5]	MiT-B2	27.00	25.30	34.5	44.1
MATEViT (ours)	ViT-Tiny	07.95	11.62	32.7	45.1
MATEViT (ours)	ViT-Small	22.08	37.10	40.2	51.1
<i>w.r.t.</i> Trans4Trans				+5.7	+7.0

platform and, therefore, to improve the user experience of the system. Compared to the single-task results in Table I, multi-task learning brings further improvements. The reason for the gains is two-fold: (1) the MMoE layer enlarges the model capacity; (2) when applying MMoE to both object segmentation and material segmentation, the former prevents the latter from overfitting and vice versa.

D. Ablation Study

To fully understand the proposed components, an ablation study is conducted, as shown in Table III. The baseline model is Segmenter ViT-Tiny [25]. Learned from Table III, replacing ViT-Tiny with ViT-Small boosts the performance (43.2%→49.2% in mIoU). After applying LIS, the mIoU continuously increases (49.2%→50.1%). We then add the MoE layer to our model, leading to an even higher mIoU of 50.7%. It can be observed that our MMoE model achieves the best performances on all metrics in both object segmentation and material segmentation tasks compared to the baseline model, *i.e.*, 6.5% and 3.7% boosts in pixel accuracy, 8.9% and 7.9% boosts in mIoU. Through the extensive pre-study experiments, the effectiveness of the proposed method can be comprehensively proved.

TABLE III: **Ablation study on COCOStuff-10K and DMS.** All values are calculated with test sets.

Method	COCOStuff-10K		DMS	
	Pixel Acc (%)	mIoU (%)	Pixel Acc (%)	mIoU (%)
Segmenter [25]	65.0	31.3	76.9	43.2
+ ViT-Small	69.1	38.2	79.0	49.2
+ ViT-Small + LIS	70.2	38.9	79.4	50.1
+ ViT-Small + LIS + MoE	70.8	39.4	79.9	50.7
+ ViT-Small + LIS + MoE + MMoE (MATEViT)	71.5	40.2	80.6	51.1

E. Qualitative Analysis

Four groups of scenarios related to the daily life of PVI are shown in Fig. 5. The upper-left group describes a scenario where PVI are having their meals. The *hot dog* colored orange is perfectly recognized by our system in the second image, and it is tagged with a food property from the third image. The group in the lower left corner shows a scenario where blind people walk in a park. According to the predictions in the second and third images, blind people using our system know additionally there is an asphaltic road ahead. In the bottom-right case, the entrance is not recognized if only object recognition is performed, however, the material segmentation result can provide supplementary information to find the doors, which can further improve the mobility accessibility. As a result, our system can help PVI better understand their surroundings, and to support them to make correct interactions with the environment.

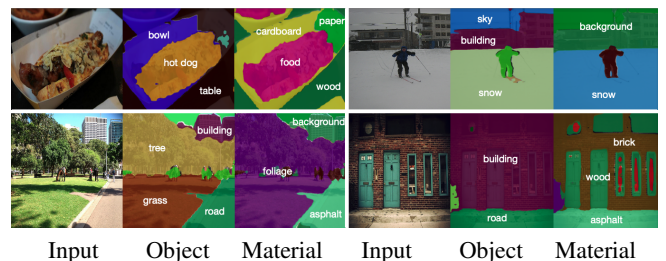


Fig. 5: **Visualization** of both general object segmentation and material segmentation. From left to right in each group: RGB input, object segmentation, and material segmentation.

V. USER STUDY

Following detailed pre-study experiments based on datasets, for wearable robotic systems, a critical factor that needs to be addressed is whether the system can provide a positive user experience in real-world situations. To know about that, we conduct a user study in a structured manner, including task-oriented testing and a questionnaire session.



Fig. 6: Incidences of participants on the user study.

A. Organization

To verify the system’s usability, we organize a user study with six blindfolded participants in real-world testing scenarios. We select seven categories of materials that are common in daily life, *i.e.*, *fabric*, *foliage*, *glass*, *metal*, *paper*, *plastic*, and *wood*. To conduct a comparison between without and with using our system, there are two rounds of material recognition, *i.e.*, **contact** and **contactless** round. The contact round is to recognize by touch, while the contactless round is to recognize by using our system. Note that, to conduct a fair comparison, the order of the objects in two rounds is randomized. The participants are required to name the material after touching an object. The reaction time from touching the object until naming the material is recorded by the organizer. After two rounds of material recognition, all participants take part in an anonymous questionnaire session. The questionnaires regarding the NASA-Task Load Index (NASA-TLX) and System Usability Scale (SUS) are filled out by all participants. Besides, there is space for participants to write down their open comments.

B. Results and Discussion

Cognitive load. To learn about the cognitive load of our wearable system, NASA-TLX is a simple and effective method for cognitive load measurement. We first calculate the average score of every factor among all participants, then average the scores of all six factors, resulting in a final NASA-TLX value of 28. According to [31], this value illustrates the workload caused by our system is in the 20th percentile of global workload scores from 6.21 to 88.50 among 1173 observations, which can assist users without too much burden. From Fig. 7, we notice that the effort value is relatively smaller than the rest five factors, meaning that our system is user-friendly.

Accuracy. The correctness of the recognition is defined as the accuracy. In both rounds of material recognition, the accuracy achieves 98%, indicating our system is useful and valuable in real-world applications.

Efficiency. The time from touching an object to naming the material is defined as the reaction time. We utilize the average reaction time of all materials to evaluate the efficiency. The average reaction time in the contactless round is 3.11s (± 0.21), while the contact one is 3.97s (± 0.34). Without any haptic perception, our wearable system can perform a faster recognition than the contact-based perception, which indicates our MATERobot is feasible and reliable to provide fast assistance in recognizing material properties. More importantly, contactless object and material recognition can provide psychological safety for PVI by allowing them to identify objects without touching them.

Usability. Apart from NASA-TLX, we verify the usability of our system with SUS. Our system scores 77 out of 100, which is a relatively high score. According to Bangor *et al.* [32], who analyzed 2324 surveys from 206 studies, “the best quarter of studies range from 78.51 to 93.93”. Therefore, we find that our system is useful for recognizing not only general objects but also their material properties.

User comments. We analyze the open comments made by users during the testing and from the post-study questionnaire session. The insights are reported below: (1) 66.7% participants were amazed by the fast response of the MATERobot system, which is one of the reasons why they would like to use the system. (2) The user experience with the system was impressive for all participants. They found our system useful and helpful in the daily lives of visually impaired people. (3) Some participants suggest that the voice feedback of the glasses should constantly inform the user of the detected objects and materials. We plan to improve this in future work.

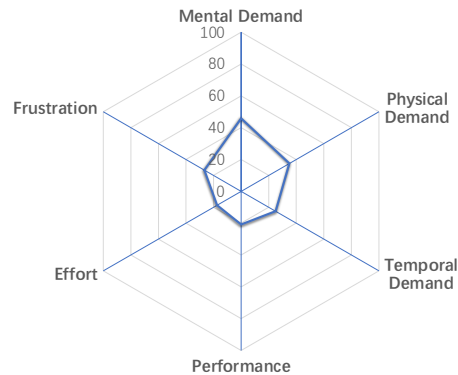


Fig. 7: **Average NASA-TLX score** of every factor among all participants. Values range from 0 to 100, lower is better. Our system requires low cognitive load, with a score of 28.

Limitations. Due to the challenges posed by the pandemic, it was difficult to find a sufficient number of participants to perform a user study on the developed system. The system was tested by six blindfolded participants. While the field test provides insights into the ability of our system, the results cannot be considered representative of the experience of real blind users since PVI could be better at contact perception than sighted people. For future work, we plan to conduct studies involving users who are blind or visually impaired.

VI. CONCLUSION

In this work, we look into semantic material understanding for helping visually impaired people via a wearable robotic system MATERobot. We put forward MATEViT, which unifies general object and material segmentation via an MMoE architecture, whose efficiency is enhanced via LIS to make plain-ViT models suitable for mobile applications. The proposed model is ported to our established assistive MATERobot system designed for supporting PVI. Extensive experiments on DMS and COCOStuff-10K datasets and a user study demonstrate the effectiveness and usefulness of our recognition system.

REFERENCES

- [1] M. J. Burton *et al.*, “The lancet global health commission on global eye health: vision beyond 2020,” *The Lancet Global Health*, vol. 9, no. 4, pp. e489–e551, 2021.
- [2] P.-J. Duh, Y.-C. Sung, L.-Y. F. Chiang, Y.-J. Chang, and K.-W. Chen, “V-eye: A vision-based navigation system for the visually impaired,” *IEEE Transactions on Multimedia*, vol. 23, pp. 1567–1580, 2021.
- [3] S. Agrawal, M. E. West, and B. Hayes, “A novel perceptive robotic cane with haptic navigation for enabling vision-independent participation in the social dynamics of seat choice,” in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 9156–9163.
- [4] H. Liu, R. Liu, K. Yang, J. Zhang, K. Peng, and R. Stiefelhagen, “HIDA: Towards holistic indoor understanding for the visually impaired via semantic instance segmentation with a wearable solid-state LiDAR sensor,” in *2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2021, pp. 1780–1790.
- [5] J. Zhang, K. Yang, A. Constantinescu, K. Peng, K. Müller, and R. Stiefelhagen, “Trans4Trans: Efficient transformer for transparent object and semantic scene segmentation in real-world navigation assistance,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 10, pp. 19 173–19 186, 2022.
- [6] M. Paterson, ““seeing with the hands”: Blindness, touch and the enlightenment spatial imaginary,” *British Journal of Visual Impairment*, vol. 24, no. 2, pp. 52–59, 2006.
- [7] R. Klatzky and S. Lederman, “Object recognition by touch,” in *Blindness and Brain Plasticity in Navigation and Object Perception*, 2007, pp. 197–220.
- [8] A. Dosovitskiy *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *International Conference on Learning Representations (ICLR)*, 2021.
- [9] H. Caesar, J. Uijlings, and V. Ferrari, “COCO-stuff: Thing and stuff classes in context,” in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 1209–1218.
- [10] P. Upchurch and R. Niu, “A dense material segmentation dataset for indoor and outdoor scene parsing,” in *European Conference on Computer Vision (ECCV)*, vol. 13668, 2022, pp. 450–466.
- [11] K. Yang, L. M. Bergasa, E. Romera, X. Huang, and K. Wang, “Predicting polarization beyond semantics for wearable robotics,” in *2018 IEEE-RAS International Conference on Humanoid Robots (Humanoids)*, 2018, pp. 96–103.
- [12] H. Wang, R. K. Katzschmann, S. Teng, B. Araki, L. Giarré, and D. Rus, “Enabling independent navigation for visually impaired people through a wearable vision-based feedback system,” in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 2017, pp. 6533–6540.
- [13] A. Aladren, G. López-Nicolás, L. Puig, and J. J. Guerrero, “Navigation assistance for the visually impaired using RGB-D sensor with range expansion,” *IEEE Systems Journal*, vol. 10, no. 3, pp. 922–932, 2016.
- [14] W. Ou *et al.*, “Indoor navigation assistance for visually impaired people via dynamic SLAM and panoptic segmentation with an RGB-D sensor,” in *International Conference on Computers Helping People with Special Needs (ICCHP)*, 2022, pp. 160–168.
- [15] R. Liu *et al.*, “Open scene understanding: Grounded situation recognition meets segment anything for helping people with visual impairments,” in *2023 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2023, pp. 1849–1859.
- [16] J. Zhang, K. Yang, A. Constantinescu, K. Peng, K. Müller, and R. Stiefelhagen, “Trans4Trans: Efficient transformer for transparent object segmentation to help visually impaired people navigate in the real world,” in *2021 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2021, pp. 1760–1770.
- [17] E. Xie, W. Wang, Z. Yu, A. Anandkumar, J. M. Alvarez, and P. Luo, “SegFormer: Simple and efficient design for semantic segmentation with transformers,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021, pp. 12 077–12 090.
- [18] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, “Encoder-decoder with atrous separable convolution for semantic image segmentation,” in *European Conference on Computer Vision (ECCV)*, vol. 11211, 2018, pp. 833–851.
- [19] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, “Pyramid scene parsing network,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 6230–6239.
- [20] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [21] Y. Liang, C. Ge, Z. Tong, Y. Song, J. Wang, and P. Xie, “Not all patches are what you need: Expediting vision transformers via token reorganizations,” in *International Conference on Learning Representations (ICLR)*, 2022.
- [22] J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” *arXiv preprint arXiv:1607.06450*, 2016.
- [23] A. Vaswani *et al.*, “Attention is all you need,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, 2017, pp. 5998–6008.
- [24] N. Shazeer *et al.*, “Outrageously large neural networks: The sparsely-gated mixture-of-experts layer,” in *International Conference on Learning Representations (ICLR)*, 2017.
- [25] R. Strudel, R. Garcia, I. Laptev, and C. Schmid, “Segmenter: Transformer for semantic segmentation,” in *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 7242–7252.
- [26] I. Loshchilov and F. Hutter, “SGDR: Stochastic gradient descent with warm restarts,” in *International Conference on Learning Representations (ICLR)*, 2017.
- [27] M. Contributors, “MMSegmentation: Openmmlab semantic segmentation toolbox and benchmark,” <https://github.com/open-mmlab/mms Segmentation>, 2020.
- [28] L.-C. Chen, G. Papandreou, F. Schroff, and H. Adam, “Rethinking atrous convolution for semantic image segmentation,” *arXiv preprint arXiv:1706.05587*, 2017.
- [29] A. Howard *et al.*, “Searching for MobileNetV3,” in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 1314–1324.
- [30] Q. Wan, Z. Huang, J. Lu, G. Yu, and L. Zhang, “SeaFormer: Squeeze-enhanced axial transformer for mobile semantic segmentation,” in *International Conference on Learning Representations (ICLR)*, 2023.
- [31] R. 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*, vol. 59, no. 1, 2015, pp. 1727–1731.
- [32] A. Bangor, P. T. Kortum, and J. T. Miller, “An empirical evaluation of the system usability scale,” *International Journal of Human-Computer Interaction*, vol. 24, no. 6, pp. 574–594, 2008.