

# Gesture-aware Interactive Machine Teaching with In-situ Object Annotations

Zhongyi Zhou, Koji Yatani

Interactive Intelligent Systems Lab., The University of Tokyo  
Tokyo, Japan  
{zhongyi,koji}@iis-lab.org

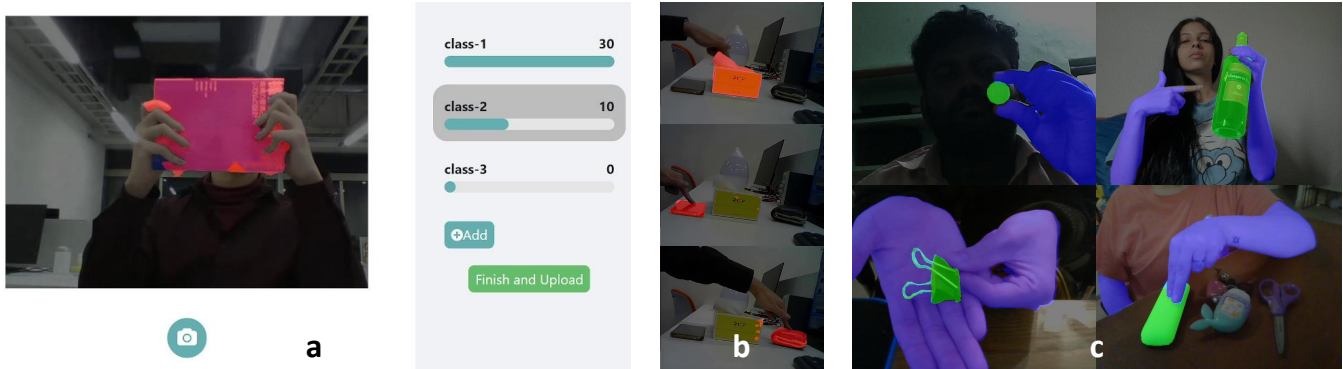


Figure 1: (a): The teaching interface of our vision-based Interactive Machine Teaching system, LookHere. LookHere provides a segmentation mask (an object highlight) on the object guided by users’ deictic gestures in real time during teaching. This segmentation mask is used for model training as additional information for training classifiers. (b): Users’ deictic gestures guides in-situ object annotations. (c): Example images in our *HuTics* dataset that enables the implementation of LookHere. *HuTics* includes 2040 labeled images that capture how 170 people use deictic gestures to present an object.

## ABSTRACT

Interactive Machine Teaching (IMT) systems allow non-experts to easily create Machine Learning (ML) models. However, existing vision-based IMT systems either ignore annotations on the objects of interest or require users to annotate in a post-hoc manner. Without the annotations on objects, the model may misinterpret the objects using unrelated features. Post-hoc annotations cause additional workload, which diminishes the usability of the overall model building process. In this paper, we develop LookHere, which integrates in-situ object annotations into vision-based IMT. LookHere exploits users’ deictic gestures to segment the objects of interest in real time. This segmentation information can be additionally used for training. To achieve the reliable performance of this object segmentation, we utilize our custom dataset called *HuTics*, including 2040 front-facing images of deictic gestures toward various objects by 170 people. The quantitative results of our user study showed that participants were 16.3 times faster in creating a model with our system compared to a standard IMT system with a post-hoc annotation process while demonstrating comparable accuracies. Additionally, models created by our system

showed a significant accuracy improvement ( $\Delta mIoU = 0.466$ ) in segmenting the objects of interest compared to those without annotations.

## CCS CONCEPTS

• **Human-centered computing** → **Interactive systems and tools**; • **Computing methodologies** → *Computer vision*; *Machine learning*.

## KEYWORDS

Interactive machine teaching, deictic gestures, in-situ annotation, dataset

## ACM Reference Format:

Zhongyi Zhou, Koji Yatani. 2022. Gesture-aware Interactive Machine Teaching with In-situ Object Annotations. In *The 35th Annual ACM Symposium on User Interface Software and Technology (UIST '22)*, October 29–November 2, 2022, Bend, OR, USA. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3526113.3545648>

## 1 INTRODUCTION

Interactive Machine Teaching (IMT) [40, 49] aims to enhance users’ teaching experience during the creation of Machine Learning (ML) models. IMT systems are primarily designed for non-ML-experts, and allow such users to provide training data through demonstrations. Vision-based IMT (V-IMT) systems utilize cameras to capture users’ demonstrations. For example, in Teachable Machine [5] users can create a computer vision classification model by showing different views of each object (class) to a camera. Despite its low

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
*UIST '22, October 29–November 2, 2022, Bend, OR, USA*

© 2022 Association for Computing Machinery.  
ACM ISBN 978-1-4503-9320-1/22/10...\$15.00  
<https://doi.org/10.1145/3526113.3545648>

burden for providing training samples, existing work [63] revealed that an ML model trained through a V-IMT system might recognize an object by using visual features unrelated to it. For example, even if a user performs a demonstration of a book, the model may use visual features in the background. Failure to address this error properly could result in the degraded performance of the model when it is deployed to real applications. Therefore, users should have the capability to specify portions of an image that a model should emphasize in learning to achieve reliable classification.

One approach to address this issue is to perform annotations on the objects of interest [46] and feed them into the model as well. Advances in annotation tools reduce user workload by simplifying necessary interactions to clicks [37, 51] or sketches [42, 59]. Despite their lowered burden, existing annotation tools are not well tailored toward V-IMT systems, and users thus have to perform annotations in a post-hoc manner. This would degrade the overall experience of V-IMT systems [50]. Annotation approaches that are more deeply integrated into V-IMT thus need to be explored.

To this end, this work examines V-IMT systems that can integrate object annotations into the teaching process. We observe that when users are doing demonstrations for teaching, they may hold or point to the object of interest. These deictic gestures in demonstrations thus are indicative of what visual features a model should focus on. Therefore, this work focuses on the integration of annotations by leveraging deictic gestures that humans naturally perform during the teaching process. Note that, in this paper, we use the term deictic gestures to represent a wide range of gestures whose purpose is to indicate the object of interest [11, 45] while it typically represents pointing gestures in HCI research [31].

Our V-IMT system, called LookHere<sup>1</sup>, embeds in-situ object annotations inferred from users' deictic gestures into the teaching process. LookHere provides real-time visualizations, named *object highlights*, on what portions of the given frame the system is considering as the region of the target object (the red mask in Figure 1a). Depending on the deictic gestures users are performing, our system infers different object regions (Figure 1b). To achieve this gesture-aware object segmentation, we created *HuTics*, a dataset consisting of 2040 images collected from 170 people that include various deictic gestures and objects with segmentation mask annotations (Figure 1c). Our technical evaluation shows that our object highlights can achieve the accuracy of 0.718 (mean Intersection over Union; *mIoU*) and can run at 28.3 fps. Our user evaluation confirms that participants were able to build accurate models while being liberated from post-hoc manual object annotations.

This work offers the following contributions:

- A vision-based IMT system, LookHere, which integrates in-situ object annotations guided by users' deictic gestures into the teaching process,
- The development of real-time object highlights, which offers users feedback on the object region inferred from their deictic gestures,

- The *HuTics* dataset<sup>2</sup>, which contains 2040 labeled images from 170 participants interacting with various target objects using deictic gestures, and
- Our evaluations that confirm LookHere's benefits through quantitative and qualitative results.

## 2 RELATED WORK

### 2.1 Interactive Machine Teaching

Machine Teaching is a term that has been used by both HCI [25, 44, 49] and ML [64, 65] communities with different definitions. To avoid conflicts, we utilized the term of Interactive Machine Teaching defined by Ramos *et al.* [40]: IMT is “an IML (*interactive machine learning*) process in which the human-in-the-loop takes the role of a teacher, and their goal is to create a machine-learned model”. This definition emphasizes user experience rather than mathematical challenges, and is well aligned with the scope of this work.

IMT systems were typically designed for non-ML-experts to build their own ML models without requiring technical knowledge and skills [16, 17]. Existing work conducted qualitative studies with ML novice users and presented user requirements and design guidelines for IMT systems [20, 43, 44]. For example, Fiebrink *et al.*'s work [18] suggested informing users of “*where and how the model was likely to make mistakes*” so that users can systematically assess the benefits and risks in their applications. Yin *et al.* [56] found that non-ML-experts users evaluated the model only based on its accuracy, and they were not often aware of the potential unreliability of the model when it was used in another application.

Zhou and Yatani [63] enhanced the model assessment process by visualizing the image regions that were highly weighed for predictions. They further found that simple teaching without further fine-grained annotations [5] could cause unexpected failures. Effective approaches for specifying the portions of the image in the teaching process are still under-explored.

This work introduces a V-IMT system that exploits users' deictic gestures to present objects for identifying the region which a model should focus on for learning. This interface design can solve the issue of accidental use of unrelated visual features by ML models by exploiting interactions people would normally perform during the teaching phase.

### 2.2 Interactive Annotations

One standard approach to addressing accidental use of unrelated visual features is to provide annotations (e.g., a segmentation mask over the target object) and inform a system of where a model should focus. While offering useful information, annotation is generally a tedious manual task. Drawing a polygon-based contour on an object [46] is a common approach to generating a segmentation mask, but this is generally very time-consuming. By incorporating computer vision methods, research has demonstrated different ways to reduce input from users [33], including clicks [37, 51], sketches [42, 59], and mouse drags [8].

As these annotation tools are not specifically designed for the integration into V-IMT systems, users would have to use them in a post-hoc manner. This does not thus fully exploit the

<sup>1</sup>The source code is available at <https://github.com/zhongyi-zhou/GestureIMT>

<sup>2</sup>The dataset is available at <https://zhongyi-zhou.github.io/GestureIMT/>.

user interaction that occurred in the teaching phase for inferring segmentation masks on target objects. Instead of proposing another annotation approach, our work utilizes users’ deictic gestures toward target objects when they are performing demonstrations to a camera. In this manner, our system achieves in-situ object annotations while teaching in V-IMT systems.

### 2.3 Interactions Using Deictic Gestures

Prior research found that infants already have an ability to perform and interpret hand gestures [6, 36]. Inspired by this inherent human capability, HCI research has developed various interfaces using deictic gestures [54, 55]. One of the earliest work in this space is “Put-that-there” [1], in which users can manipulate virtual objects through a combination of deictic gestures and natural languages. Interface applications of deictic gestures also include drone manipulations [7], Human-Robot Interaction (HRI) [39, 45] and commutations in Mixed Reality [3]. Sauppe *et al.* [45] demonstrated a human-like robot that can perform deictic gestures, and found that these gestures can contribute to improving communicative accuracy in interactions with users.

Besides deictic gestures, research has investigated different aspects of hand-object interactions. By assuming that the object under humans’ manipulations would follow 1-DOF movements, Hartanto *et al.* [23] created a method to segment the object and classify the type of object motions (pure displacement motion by the prismatic joint or pure rotational motion by the revolute joints). Other work built datasets of hand-object interactions [4, 34, 48], and aimed to derive data-driven approaches for recognizing these interactions. Lee and Kacorri [34] created the TEgO dataset and a system for people with visual impairments to recognize a pre-defined set of daily-life objects and assist interaction with them.

This work extends the application of deictic gestures to V-IMT systems and allows users to perform in-situ annotations while teaching in real time. More importantly, LookHere advances the generalizability by removing those constraints in prior work (e.g., pre-defined object categories [34] or specific motions associated with holding [23]).

## 3 RESEARCH CHALLENGES AND QUESTIONS

### 3.1 Challenges in Existing V-IMT Systems

After reviewing the existing V-IMT systems, the authors summarized our perceived challenges in the following two aspects:

**C1. ML models created through simplified processes supported by V-IMT can be unreliable because they may learn features unrelated to target objects.** One major shortcoming of V-IMT is that ML models created through such systems may unpredictably attend to unrelated objects, which is aligned with the findings by Zhou and Yatani’s work [63]. To simplify the creation of ML models for non-experts [16, 49], V-IMT systems typically only ask users to perform several demonstrations to the camera [5, 19]. During teaching, the computer not only captures the object to be classified, but also other unrelated backgrounds or objects. A model thus may consider those unrelated features as critical components of the target objects while users expect that it would only capture the features on the target objects. This discrepancy

could result in unexpected inaccuracy when the model is brought to actual use. This can greatly degrade the usability of the created model and affect users’ trust [56] toward it.

**C2. Post-hoc annotations can diminish the overall usability of V-IMT systems.** A naive approach to solve the aforementioned issue is to ask users to specify what they want to be included in models (i.e., annotate the image regions of the objects). Existing work has successfully simplified data annotations [37, 52], but these interfaces are mostly designed for more professional use [32]. Furthermore, creating a reliable ML model usually requires the user to provide many samples per class. Performing annotations on many images repeatedly in addition to teaching through V-IMT systems can thus be overwhelming to non-expert users [50]. This also can discourage casual use of ML, which many V-IMT systems envision.

While formative user studies could further confirm these challenges, we decided not to execute them as they are already well explained in the existing literature. Our user evaluation results presented in Section 8 also confirm these challenges well.

### 3.2 Research Questions

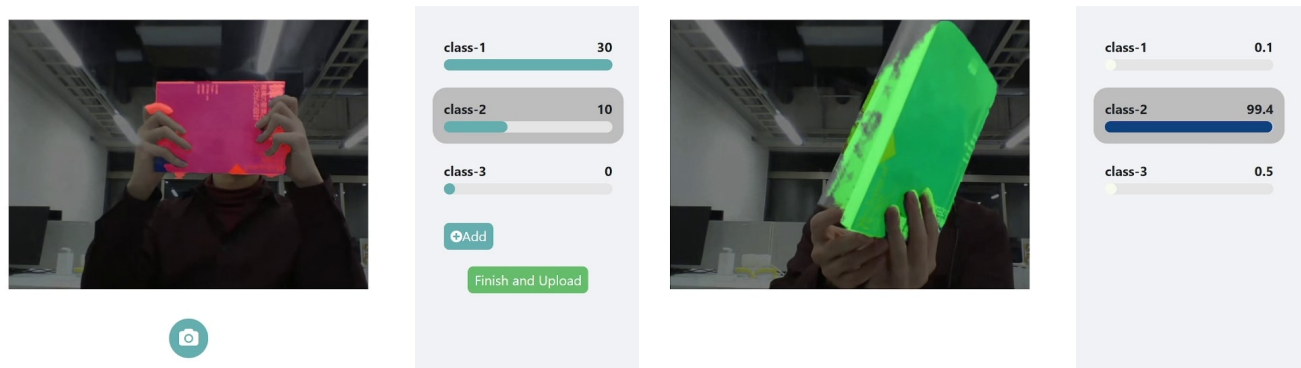
This work explores approaches to solve both challenges by *integrating annotations into the teaching process*. To achieve this integration, we exploit how people interact with objects of interest using *deictic gestures* when they perform demonstrations to a camera. For example, they may hold the object with both hands or may point to the object with their index fingers. Such human behaviors are known in prior HCI research [30] that led to diverse gesture-based applications [34, 38]. Therefore, we hypothesized that such deictic gestures would be an important cue for in-situ annotation. Accordingly, we derive the following two research questions to be answered through this research:

- RQ1.** *How can users’ deictic gestures toward objects of interest be utilized for annotations during teaching?*
- RQ2.** *Can such in-situ annotations inferred from users’ deictic gestures reduce the overall teaching workload while maintaining the model accuracy?*

RQ1 asks for technical approaches for leveraging humans’ deictic gestures for the integration and the corresponding implementation. RQ2 investigates the efficacy of such gesture-aware annotation methods. Our design and implementation of LookHere in the following content explore RQ1, and our evaluation study answers RQ2 using multiple metrics (i.e., time consumption, model accuracies and subjective workload).

## 4 LOOKHERE

Our V-IMT system, LookHere, considers users’ gestures to objects for building accurate ML models. Unlike existing workflows in V-IMT, LookHere directly integrates the annotation process into the teaching process. More specifically, LookHere includes a function called *object highlights* to inform which part of the camera view the system is considering as the region of the object to be learned. In the assessment phase, LookHere also supports a model assessment process by providing a similar visualization, allowing the user to assess whether the trained model attends to the correct features.



(a) Teaching Interface.

(b) Model Assessment Interface.

**Figure 2: The screenshots of LookHere. (a) In this teaching interface, real-time object highlights are provided. The number of samples per class is presented on the right side of the view; (b) In this model assessment interface, the saliency map visualizations for the prediction of a specified class (i.e., class 2 in this example) are shown along with the prediction confidence score. This feedback informs users of what visual features in a given frame a model is weighed for predictions.**



**Figure 3: Highlights are overlaid on different objects depending on users’ deictic gestures.**

Besides these two features explained in this section, the architecture and interaction walkthrough are similar to existing V-IMT systems. In our current implementation, users can train a multi-class classifier (i.e., classifying different objects). To define a class, the user first selects the corresponding class (see the top-right corner of Figure 2a). Then, they can perform demonstrations of the object to the camera, and the system captures the frame when the user clicks a camera button. The number of frames collected for each class is presented as a bar graph. After finishing teaching for the three classes, users may either click the “Add” button to include more classes or the “Finish and Upload” button to finish the teaching session. Appendix A.1 shows our detailed configurations in the ML process after teaching.

#### 4.1 Object Highlights and In-situ Object Annotation

During the teaching process, users can receive visual feedback about which portion of the camera view LookHere is currently considering as the region of the objects of interest. As is shown in Figure 2a, our system infers the object region based on deictic gestures users are performing (e.g., holding or pointing to an object for teaching). Users may simply change how to perform gestures to express different target objects, as shown in Figure 3. LookHere incorporates a gesture-aware algorithm (see details in the next section) to achieve this adaptive highlight on objects.

Another advantage of providing this highlight in real time during the teaching process is to help users avoid including erroneous demonstrations. Users can easily opt out of such frames by not

clicking the camera button. In this manner, LookHere takes a mixed-initiative approach [26] for teaching.

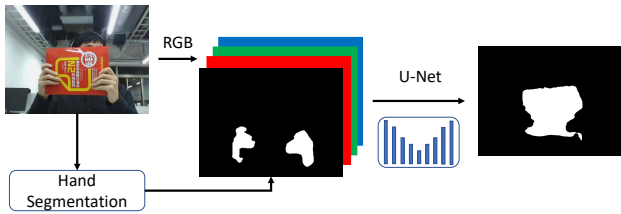
When the user records the current frame by a button click, the system stores the RGB image as well as the inferred object segmentation mask. Both data are used for model training. In this manner, LookHere achieves in-situ object annotations during the teaching process.

#### 4.2 Model Assessment with Saliency Map Visualizations

After the teaching phase, LookHere offers the model assessment mode like other V-IMT systems [5, 19]. However, unlike these systems, LookHere provides saliency map visualizations for users to confirm whether the created model is considering appropriate visual features. Figure 2b illustrates an example of the view in this assessment phase. The interface presents two visualizations for the users: bar graphs to present confidence score distributions (Figure 2b right) and real-time saliency map visualizations (Figure 2b left). The confidence score shows how confidently the model considers that the current frame belongs to the corresponding class. In the example of Figure 2b, the model is 99.4% confident that the object in the frame belongs to its class 2 (which is configured as a “book” class in Figure 2a). Real-time saliency map visualizations then help users understand which portion of the frame the model considers as the object of interest (a book in this example). Existing work [60, 63] leveraged CAM methods to present such visualization while we introduce a new method for more accurate visualizations by utilizing the object segmentation masks originally generated for object highlights (see more details for Section 5.2).

### 5 IMPLEMENTATION

The current prototype of LookHere is implemented as a web-based interface, and most of the computations are executed at the back end. We use WebRTC to synchronize the video between the interface and server for real-time image processing. The two key features presented in Section 4 are supported by two technical components:



**Figure 4: The generation process of object highlights.** LookHere first performs a hand segmentation with the given RGB image. The system then feeds both the RGB image and segmentation mask into U-Net, which predicts a segmentation mask of the object guided by deictic gestures.

gesture-aware object highlights and joint training. We explain the details of the implementation of these components in this section.

### 5.1 Gesture-aware Object Highlights

Figure 4 summarizes the workflow of our gesture-aware object highlight algorithm. The algorithm first applies a hand segmentor on the input image and predicts a hand segmentation mask. It then feeds both the original RGB image and the hand segmentation mask into U-Net [41], which outputs a segmentation mask of the object that is referred to by the users’ deictic gestures.

**5.1.1 Hand Segmentation.** We utilize Li *et al.*’s algorithm [35] trained on the LIP dataset [21] to perform real-time hand segmentation. The LIP dataset parses a person into 20 body parts and garments (e.g., “left-leg”, “gloves” and “pants”), and we regard the segmentation result of “left-arm” and “right-arm” as the portion of hands. We note that the definition of “arm” in the LIP dataset includes both arms and hands that are not covered by clothes or gloves. We notice that the publicly-available model provided by the authors of the LIP dataset is not suitable because it utilizes resnet-101 backbone [24], which is a very deep CNN architecture and is not executable in real time. Therefore, we re-design their methods based on resnet-18, a much lighter model with the same encoding approach. We then tested this light-weighted model on the LIP dataset. The result *mIoU* accuracy of the light model is 0.621, and that of the original model using resnet-101 is 0.680. This demonstrates that our light model for real-time uses can still achieve comparable accuracy to the original deep model.

**5.1.2 Object Highlights.** As explained in Section 4.1, object highlights offer immediate feedback on what portions of the image frame the model to be trained should focus on. To avoid losing the generalizability of V-IMT, LookHere should be able to segment the object of interest in an object-agnostic manner. To tackle this challenge, we feed the RGB image concatenated with the hand segmentation mask inferred from the hand segmentor (Section 5.1.1) into a recognition model as shown in Figure 4. Intuitively, this hand segmentation mask carries the information of what objects in the frame users are specifically referring to in their demonstrations.

The current implementation uses U-Net [41] as the encoder-decoder architecture. It performs the best as well as the fastest among four commonly-used segmentation model architectures (see Appendix A.3 for detailed data). The network uses the EfficientNet [53] backbone, a design toward high computation



(a) CAM ( $\Lambda = 0$ ). (b)  $\Lambda = 1$ . (c)  $\Lambda = 0.718$

**Figure 5: Visual comparison of saliency maps with different settings of  $\Lambda$ .** The parameter  $\Lambda$  in Equ. 1 controls the weight balance between the results by CAM and our trained model.

efficiency. To train this U-Net, we use our own dataset which we will explain in Section 6.

### 5.2 Joint Classification and Segmentation for Saliency Map Visualizations

Saliency map visualizations are useful for users to understand what specific portions of a given image are weighed more in their ML models. Existing work [60, 63] created saliency maps of a classification model through CAM methods [47, 61]. CAM methods are primarily used for simple classification models trained by the dataset without segmentation masks. Unlike existing V-IMT systems, our training data accompany the object segmentation masks inferred during the teaching phase. We thus devise a new model training approach for LookHere to exploit this unique information resource to achieve more accurate saliency maps.

LookHere identifies the areas to be highlighted by saliency maps through solving a classification and segmentation problem jointly. This means that our backend model predicts a class as well as infers the segmentation of the object of interest at the same time. More specifically, we train the model through a joint loss function ( $l_{joint}$ ), which is a weighted sum of classification loss ( $l_{cls}$ ) and segmentation loss ( $l_{seg}$ ):  $l_{joint} = l_{cls} + \lambda \cdot l_{seg}$ .  $\lambda$  is a trade-off weight that determines the relative importance between the classification loss and segmentation loss. In our current prototype, we set  $\lambda$  to 1, making both of them equally important in the training process.

While the segmentation masks originally created for object highlights can be useful for training our backend model for saliency maps as we discussed above, they may also contain some errors because the generated mask is not always perfect. Such errors may lead to degradation in the accuracy of segmentation inference for saliency maps. To eliminate this effect, we introduce another parameter ( $\Lambda$ ) to control the balance between the inference results by our backend model and CAM methods:

$$\Lambda \cdot Out + (1 - \Lambda) \cdot CAM \tag{1}$$

*Out* represents the segmentation output of the our backend model and *CAM* is the CAM inference result. A larger  $\Lambda$  value means that the system weighs more on our inference result for the output for saliency maps.

We found that taking such trade-off in consideration can greatly improve the accuracy of our saliency maps in some challenging cases. Figure 5 illustrates the effect of  $\Lambda$  in a case where a user is holding a plastic bottle. The saliency map visualization can be quite erroneous when we only use the results of CAM (Figure 5a). This approach would include regions that are not related to the object of interest. On the other hand, when we only use the prediction by our backend model, the result tends to be overly conservative

(Figure 5b). One reason of this issue is over-fitting. In this example, we deliberately used different backgrounds for training and testing. As the bottle in this example was transparent, the model might have included (or overfit) some visual features of the background during training. Such features would not appear when the background was changed when being tested, and this could thus explain why our model can be very conservative.

By choosing an appropriate value for  $\Lambda$ , the saliency map can visualize the object region more precisely (Figure 5c). We chose  $\Lambda$  value to be the accuracy of our object highlights in our current implementation and technical evaluation (i.e., 0.718 using EfficientNet-b0 backbone). It is out of our scope to investigate how to achieve optimization on this parameter.

## 6 DEICTIC GESTURE DATASET

### 6.1 Motivation of Data Collection

As explained in the previous section, the backend model for object highlights needs training data of how people perform deictic gestures to objects to a camera. Among existing related human-object datasets [13, 14, 34, 48], TEgO [34] is the one that best fits our task. TEgO includes 5758 labeled egocentric images of hand-object interactions. For each image, there is a hand segmentation mask and a point-level annotation of the object location, which is not immediately sufficient for our purpose (object segmentation). We therefore attempted to infer the segmentation mask of the object by emulating a click-based interactive segmentation method [52]. We then manually inspected all the generated results and removed data samples where the inferred segmentation masks were completely inaccurate. This constitutes our customized dataset with automatically-synthesized object segmentation masks, called TEgO-Syn ( $n=5232$ ).

The trained network using TEgO-Syn achieved  $mIoU=0.895$  on the testing set, showing a seemingly-promising result. Appendix A.2 provides our detailed training configurations. To further evaluate the robustness of the network in real applications, we experimented with this model with images where various objects were presented through different deictic gestures. Our observations showed that the model was not robust enough which we will further confirm in Section 6.3. We then summarized three main reasons why TEgO still cannot fit our target task:

- **A limited set of gestures.** All data in TEgO were collected from two participants, which is insufficient to cover how different people interact with the object using gestures.
- **A limited set of objects.** TEgO-Syn includes 5232 images of 19 objects. Training on a small set of objects repeatedly enables the model to over-fit the features of these specific objects, which is harmful to our target task, i.e., object-agnostic segmentation.
- **Egocentric images.** The images in the TEgO dataset are taken from the egocentric view. Our system uses a front-facing camera, which is a common configuration in V-IMT [5].

### 6.2 HuTics Dataset

To address the three issues above, we created our own dataset. We recruited crowd-workers on Amazon Mechanical Turk, aiming

to enhance the diversity of the dataset.<sup>3</sup> In each task, the worker needed to upload 12 images in total that clearly showed how they would use deictic gestures to express the references to objects. For collecting a diverse set of images from each worker, we first classified deictic gestures into four categories based on Sauppe *et al.*'s taxonomy [45]: pointing, presenting, touching and exhibiting. We then asked the workers to take three different photos for each gesture category. We also provided example pictures to clarify our expectations to the workers.

We collected 2040 qualified images from 170 crowd-workers (M: 99; F: 71) in total. The average age of the workers was 34 ( $SD: 9.2$ ). Example unqualified submissions included images that were highly blurry or where no gesture was involved at all. The crowd-workers spent 15 minutes on average to complete the task, and we paid each participant 2 dollars. We then recruited another five people on our local crowdsourcing platform to annotate object segmentation masks on the collected images. On average, each annotation worker labeled 408 images, and we compensated them with approximately 78 dollars on average in our local currency. During the annotation, the workers used AnnoFab [28], an online polygon-based tool, to label the segmentation masks.

Figure 6 presents example images with the annotated object segmentation masks. Unlike TEgO, our dataset contains a wide range of objects, deictic gestures, backgrounds, and environmental conditions. Table 1 summarizes a comparison of *HuTics* v.s. TEgO-Syn.

### 6.3 Performance of Object Highlights on HuTics

We used the data from 80% of the participants in *HuTics* (i.e., 1632 images from 136 people) for training and 20% for testing. We trained our algorithm using the same configuration above, and the network achieves  $mIoU=0.718$  and 0.806 using EfficientNet-b0 and EfficientNet-b3 backbone on the testing set, respectively. Running on one GTX 2080Ti GPU, our implementation of the algorithm was able to reach 28.3 fps and 24.0 fps with the EfficientNet-b0 and EfficientNet-b3 backbone, respectively. For comparison, we trained another model with the same network architecture using TEgO-Syn and tested with images in *HuTics*. The accuracy of that model was  $mIoU=0.368$ , much lower than that of the same network using *HuTics* for training. This significant accuracy drop from 0.895 (tested on TEgO-Syn) further confirms our observations discussed in Section 6.1.

Figure 7 shows a visual comparison of the results between the networks trained on TEgO-Syn and *HuTics*. Each example in Figure 7 are the one in our testing set that has the closest IoU values (0.366 and 0.719) to the corresponding mean IoU values (0.368 and 0.718). We therefore use the model trained on *HuTics* in our current prototype implementation.

## 7 USER STUDY

We conducted a comparative user study to evaluate how LookHere could improve the experience of V-IMT in terms of time cost for teaching, accuracy performance on models created, and subjective user workload.

<sup>3</sup>We received IRB approval for this data collection at our university.

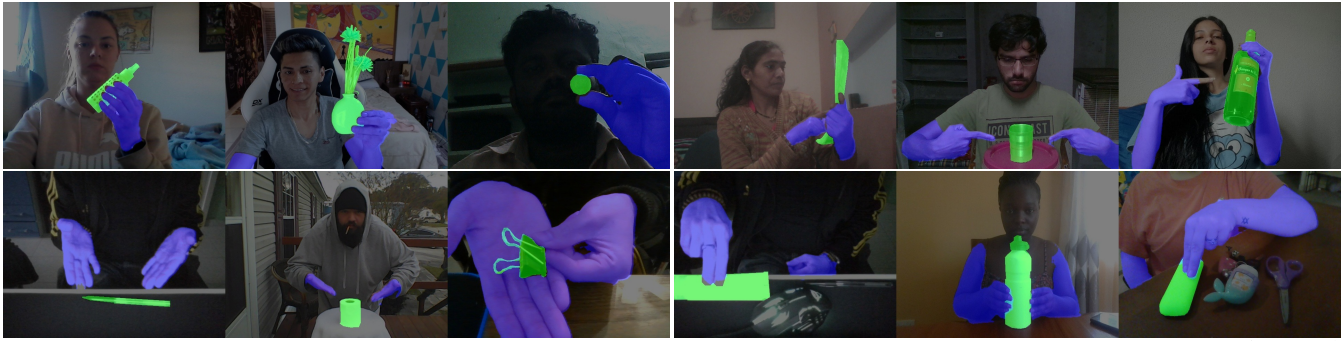
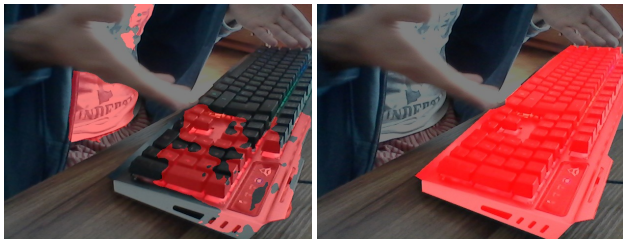


Figure 6: Example images in *HuTics* dataset. *HuTics* covers four kinds of deictic gestures to objects: exhibiting (top-left), pointing (top-right), presenting (bottom-left) and touching (bottom-right). The hands and objects of interest are highlighted in blue and green, respectively.

Table 1: Comparison of *HuTics* and *TEgO-Syn*.

	# of Participants	Object Types	View	# of Images	Object Annotation	Target Task
<i>HuTics</i>	170	Uncontrolled	Front-facing	2040	Segmentation mask	Object-agnostic segmentation specified by gestures
<i>TEgO-Syn</i>	2	Controlled	Egocentric	5232	Point-based	Object recognition for people with visual impairments



Prediction.

Ground truth.

(a) An example with the model trained with *TEgO-Syn* ( $IoU=0.366$ ).



Prediction.

Ground truth.

(b) An example with the model trained with *HuTics* ( $IoU=0.719$ ).

Figure 7: Visual comparison of predictions by the models trained with the two datasets (*TEgO-Syn* and *HuTics*).

### 7.1 Interface Conditions

Besides *LookHere*, we included the following three interface conditions to represent existing V-IMT and object annotation methods.

- *NaïveIMT*: This represents the most common design in current V-IMT systems [5, 19]. In this condition, participants would only perform object demonstrations during the teaching phase. Participants would not have an opportunity to specify which regions of the recorded images would represent the object for

a given class. We implemented this naïve IMT system based on the source code of Zhou and Yatani [63] available online.

- *Contour*: In addition to the teaching process with the naïve IMT system, this condition would involve a manual annotation procedure in a post-hoc manner. In this condition, participants would be asked to perform contour-based annotations. This annotation style is widely used in IMT systems for medical purposes [2]. We used *AnnoFab* [28] for post-hoc contour-based annotations in this study.
- *Click*: The third reference condition included a click-based annotation method [52]. We decided to include this condition as the annotation process would be more lightweight than a contour-based approach. We used *RITM* [52] as the click-based annotation tool.

All these three reference conditions involve the teaching process using the naïve IMT system. To shorten the overall study time, we decided to ask participants to perform teaching under the two conditions of *NaïveIMT* and *LookHere*. After this teaching task, participants were then asked to perform annotations under the two conditions of *Contour* and *Click* using the data recorded under the *NaïveIMT* condition. In this manner, we liberated the participants from performing the same tasks repeatedly with *NaïveIMT* for the *Contour* and *Click* conditions.

We counter-balanced both the condition order of *NaïveIMT* and *LookHere* and that of *Contour* and *Click* across participants. The order of tasks of teaching and annotation was fixed (the teaching process was the first).

### 7.2 Evaluation Metrics

**7.2.1 Teaching and Annotation Time.** We measured how long it took for participants to finish the model creation process under each interface condition. Specifically, we recorded the teaching/annotation time from when the participants started

uploading/annotating the first sample to when they finished the last (30th) sample.

**7.2.2 Model Accuracy.** We measured both classification accuracy and segmentation accuracy (i.e., mean Intersection over Union or *mIoU*) of the created models. We randomly used 80% of the data for training and the rest for testing. For classification accuracy, we utilized cross-condition validation. Specifically, we tested three models trained by data collected from *naïveIMT* on the data collected from LookHere, and vice versa. In terms of the object segmentation accuracy, we only performed cross validation for data collected from LookHere because there were no ground-truth segmentation annotations in LookHere to validate models created from *naïveIMT*. This ensured that LookHere gained no advantage over the three comparative conditions.

**7.2.3 NASA-TLX.** NASA Task Load Index (TLX) [22] is a standard metric for perceived workload. We included this to understand how different conditions could affect the experience of creating ML models with different configurations of IMT systems.

### 7.3 Procedure

At the beginning of the study, we told participants that their goal was to create four AI models to classify and detect objects by the given systems. For the teaching tasks, we allowed participants to use any object available in our experimental space. They were also welcome to bring their own belongings in the study. We did not limit the set of objects to be used in this experiment in order not to lose the validity of the study. After they had explanations about the two teaching methods (*NaïveIMT* and LookHere) and became comfortable with using both, they were asked to create three classes for classification, and generate 30 images for each class through the given teaching method. After completing teaching with the two methods, participants were given an opportunity to take a break.

We next moved to the annotation tasks with the two interfaces (*Contour* and *Click*). We provided explanations about these two tools, and participants were given practice time to become comfortable with using them. Participants were then asked to annotate all the images they captured under the *NaïveIMT* condition. They were instructed to perform each annotation task as fast and accurately as possible. Participants were allowed to take a break between the two sets of tasks (i.e., using the two annotation tools).

Participants were asked to fill in NASA-TLX questionnaires after finishing each of the two task sessions (teaching and annotation). In this manner, we ensured that participants remembered their experience of the conditions. When participants were rating NASA-TLX for *Contour* and *Click*, we explicitly asked them to consider their overall workload for the combination of the teaching method with *NaïveIMT* and the given annotation approach.

After completing both the teaching and annotation tasks, we conducted semi-structured interviews with participants. We first interviewed them about their overall experience and perceived benefits and shortcomings of the methods used in the study. This helped us collect their immediate use experience of each condition without being biased by the performance of the resulting models (i.e., the accuracy of the created models). We next offered our model assessment interface (Figure 2b) for all the four resulting

**Table 2: The mean values and standard deviations of the task completion time and accuracies (classification and segmentation) across the four interface conditions.**

		LookHere		<i>NaïveIMT</i>	<i>Click</i>	<i>Contour</i>
		$\Lambda=1$	$\Lambda=0.718$			
Time [s]		104 (44)		67 (10)	1,197 (228)	1,483 (407)
Acc.	Cls.	0.824 (0.158)		0.847 (0.190)	0.880 (0.141)	0.833 (0.159)
	Seg.	0.578 (0.233)	0.605 (0.153)	0.139 (0.095)	0.716 (0.167)	0.732 (0.151)

models. Participants were allowed to freely use them and check whether their created model would function accurately. We then interviewed them about how they perceived their four models and would characterize them differently.

The whole study takes approximately 3.5 hours on average. We offered each participant compensation of approximately 40 USD in a local currency at the end of the experiment.

### 7.4 Apparatus

We set the video frame rate to be 24 fps across all the conditions. We used the EfficientNet-b0 backbone in our object highlights (i.e., the lightest model), aiming to understand the effectiveness of such feedback even under the least accurate setting.

### 7.5 Participants

We recruited 12 non-expert participants (P1 – P12) for this study. None of them had experience in studying or working in the fields related to AI or ML. Eight of them were female, and the rest were male. The age of participants ranged from 23 to 28.

## 8 RESULTS

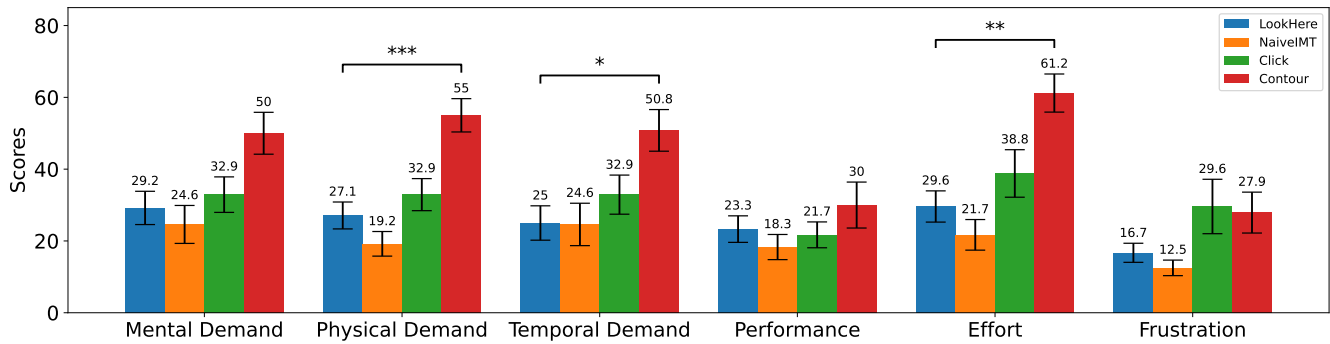
### 8.1 Quantitative Results

Table 2 presents the mean completion time and model accuracy in the study. With respect to the task completion, the *Contour* and *Click* conditions exhibited much longer time than the other two conditions (11.5 and 14.3 times than the LookHere condition, respectively). One-way repeated-measure ANOVA found that the factor of the conditions was significant ( $F(3, 33)=134.02$ ,  $p<.001$ ,  $\eta^2=.92$ ). We then used Scheffe’s multiple comparison procedure to compare the take completion time under the LookHere condition against the three reference conditions. We found that the completion time under the LookHere condition was significantly shorter than those under the *Contor* and *Click* conditions (both  $p < .001$ ). This result clearly suggests that LookHere successfully removed the effort for object annotations.

We next looked into the accuracies of the models created under the four conditions. As shown in Table 2, the mean accuracies for classification (predicting the correct class for the given image from the three classes defined by each participant) did not show large differences. Our one-way repeated-measure ANOVA did not find a significant effect of the interfaces ( $F(3, 33)=.460$ ,  $p=.712$ ,  $\eta^2=.04$ ).

We further examined the segmentation accuracies. As shown in Table 2, the accuracy under the *NaïveIMT* condition was clearly lower than those with the other methods. One-way





**Figure 8: The results of the mean NASA-TLX scores across the six subscales. The error bars represent the standard errors. The observed significant differences are indicated with \*, \*\* or \*\*\* ( $p < .05$ ,  $p < .01$  and  $p < .001$ , respectively).**

repeated-measure ANOVA found a significant effect of the interface conditions ( $F(3, 33)=105.98$ ,  $p < .001$ ,  $\eta^2=.91$ ). Scheffe’s multiple comparison procedure revealed a significant difference between LookHere and NaiveIMT ( $p < .001$ ). This result confirms that the models created with data collected under the NaiveIMT condition did not necessarily weigh the visual features in the objects of interest, implying potential unreliability in actual use.

We also compared the segmentation predictions on the same trained model with two different  $\Lambda$  values: 1 and 0.718 (the default configuration in our current implementation). While the accuracy was improved by 0.027 with the value of 0.718, this difference was not significant. Future research should investigate how  $\Lambda$  should be configured to achieve the best performance, but this result implies that the combination of CAM results and the inference by the backend object segmentation model could offer improvements.

We next examined the NASA-TLX results. Figure 8 shows the mean values of the raw NASA-TLX subscales. One-way repeated-measure ANOVA on each subscale found significant effects by the conditions in mental demand ( $F(3, 33)=7.37$ ,  $p < .001$ ,  $\eta^2=.40$ ); physical demand ( $F(3, 33)=16.27$ ,  $p < .001$ ,  $\eta^2=.60$ ); temporal demand ( $F(3, 33)=7.86$ ,  $p < .001$ ,  $\eta^2=.42$ ); effort ( $F(3, 33)$ ,  $p < .001$ ,  $\eta^2=.57$ ); and frustration ( $F(3, 33)=3.23$ ,  $p < .05$ ,  $\eta^2=.23$ ). No significant result was found in performance ( $F(3, 33)=1.55$ ,  $p=.22$ ,  $\eta^2=.12$ ). Post-hoc Scheffe’s procedure revealed significant differences in physical demand, temporal demand and effort between LookHere and Contour ( $p < .001$ ,  $p < .05$ , and  $p < .01$ , respectively). These results suggest that the contour-based annotation method significantly impacted the user experience negatively.

In summary, the quantitative results show that LookHere was able to achieve the best balance of task completion time and model accuracy. We further looked into how different user experience of the four interface conditions was through our qualitative results.

## 8.2 Qualitative Results

We transcribed the interviews and extracted quotes that were related to user experience and opinions about the four interfaces tested. We then performed the open coding approach to categorize the quotes and derive them in a bottom-up manner.

**8.2.1 Burden for post-hoc annotations.** Six participants (P1, P4, P5, P9, P10 and P12) explicitly mentioned that post-hoc annotations

were tedious and reduced the perceived usability of an overall V-IMT system. While nine participants preferred the Click annotation method to Contour, all participants agreed that both approaches were “time-consuming”.

*“[A good process] shouldn’t contain the annotation process because it is the most time-consuming one and requires lots of effort. On the contrary, these two (NaiveIMT and LookHere) are very comfortable to use because there is only one step.”* [P4]

All participants considered LookHere as “efficient” because it does not involve explicit post-hoc annotations. This was clear from the task completion time and NASA-TLX results, and our qualitative data were also indicative. In particular, P1 appreciated that LookHere combined the teaching and annotation process:

*“It can greatly improve the user experience in terms of not only time consumption but also the sense of satisfaction.”* [P1]

Participants were also satisfied with the accuracy of their created models achieved through LookHere. Despite the simplified teaching experience, they could not notice the accuracy difference between LookHere and Contour.

*“I prefer to use [LookHere]. First, its accuracy is good, and it’s easy to use ... It is a user-friendly design, not requiring much effort and time.”* [P10]

*“In terms of effort and performance, [LookHere] is definitely a cost-effective choice ... Speaking of [Contour], it requires much effort, but its result is not that good, probably similar to [LookHere]. It makes me feel that it is not worthwhile.”* [P6]

**8.2.2 Uncertainty in teaching with NaiveIMT.** Participants expressed their concerns about whether the model created with the NaiveIMT approach correctly interpreted their teaching.

*“Because you can’t find its focus, as a user, you can’t confirm whether it (the computer) understands my idea.”* [P8]

*“[NaiveIMT] is very convenient to use but I am afraid that the performance would be bad.”* [P3]

On the other hand, object highlights shown in LookHere offered our participants more confidence that the regions of the objects would be considered more.

*“[Different from NaiveIMT, LookHere] is simple, and it also has visualizations. It can let users keep well informed whether the object is recognized [by the computer].”* [P5]

In case object highlights were out of place, participants adjusted their deictic gestures until they were well overlaid onto the objects. This offered a sense of control as P12 commented:

*“On one hand, the procedure is simple, and on the other hand, [LookHere] itself has already drawn that pattern (object highlights). [Even though it sometimes has errors,] I can change some positions [of the object] and it can [successfully] capture this [object] ... It provides a sense of control. Unlike [NaiveIMT], I do not know what it captures [within each image].” [P12]*

**8.2.3 Limitations.** Participants pointed out limitations of LookHere, and some further provided suggestions on how we can improve the current prototype. For example, P11 raised an issue that users could not interact with LookHere using bimanual interaction since one hand is required to manipulate the mouse, clicking on the camera icon in Figure 2a.

Additionally, P4 pointed out that there was a lack of further teaching/clarification support when object highlights fail. P4 further suggested that V-IMT should integrate more functions so that users can better correct object highlights in erroneous cases, rather than passively avoid teaching these samples.

*“In terms of [LookHere], is it possible to utilize the Click function there? For example, when I hold something, [if object highlights are erroneous at this moment,] can I tell the computer which region I want it to recognize [by clicking on the object]? ... In the current design, I can only change the position (of the object) to adjust it (the highlight), and this makes me feel quite inactive (i.e., not in good control of object highlights).” [P4]*

## 9 DISCUSSION

As mentioned in Section 3.2, this work aims to (1) explore technical solutions for the integration of object annotations into the teaching process (i.e., RQ1) and (2) study the effectiveness of the solution (i.e., RQ2). We answer RQ1 through our implementation of LookHere as explained in Section 4 and 5. Our evaluation results further show that LookHere can effectively reduce users’ workload during the model creation process while maintaining similar model accuracies, answering RQ2.

Despite its effectiveness, we also found several drawbacks of our systems that limited user experience in practice. In the following content, we share our insights about how future work can improve our system and how to extend our research questions to exploit more human interactions to achieve in-situ annotations in IMT.

Additionally, our dataset, HuTics, which is designed for the gesture-aware object-agnostic segmentation task, is one important contribution of this work. We further show that our approach and dataset can also be used to support other HCI projects by demonstrating several example applications.

### 9.1 Depth-aware Object Highlights

Despite the effectiveness of object highlights to support efficient teaching, we still observed typical erroneous cases that remain to be addressed. Figure 9 shows an example where *IoU* was low (more examples can be found in Appendix B). The person in this figure is pointing at an object on her head, but our algorithm incorrectly highlights the clock in the background. Our object



(a) Prediction. (b) Ground truth.  
Figure 9: A failure case of object highlights.

highlights also tend to fail in cases where a person is pointing at an object at a distance (e.g., buildings or furniture that are not close to hands). These failure cases are mainly caused by the current implementation that uses 2D hand segmentation features without a 3D understanding of the scene. A future system may consider obtaining a richer set of information through 3D scene reconstruction [12], 3D hand pose estimation [4, 27] from RGB images, or directly use depth cameras [29]. Future research should further study how to simplify the aforementioned feature extractors to be used in real time for V-IMT systems or how to use depth sensors to support V-IMT systems [57, 58].

### 9.2 Voice Input and In-situ Correction

As mentioned in Section 8.2.3, we observed several limitations of our interface design. To enable bimanual interactions with the object, future research can investigate how to use technologies like voice input or facial expression recognition to replace a button click. For example, when users want the system to sample the current frame, they can simply say “collect” or smile to the system while performing bimanual deictic gestures.

In addition, as mentioned in Section 8.2.3, future systems should study how to enable users to *actively* correct object highlights when they observe prediction errors. Although the segmentation annotation in LookHere allows users to choose appropriate frames for teaching, the role of users in this Human-AI collaboration is relatively passive. When users observe the failure case of object highlights, they should be given an opportunity to *actively* correct the error [50], instead of *passively* avoiding those data. Allowing such in-situ correction initiated by users can further empower the ability of IMT systems to “leverage human capabilities and human knowledge beyond labels” [40], achieving better human-AI collaboration.

### 9.3 Other Modalities and Privacy Issues

While this paper focuses on studying how to use deictic gestures to enable in-situ annotations, they are not the only human interaction that future IMT research can exploit. As we discussed in Section 9.1, our gesture-aware annotation approach may not function with some deictic gestures users may perform. More importantly, humans also innately perform other interactions as a cue of objects of interest. For example, future research can study how to use gaze tracking technologies to capture the object of interest that is difficult to hold by hand (e.g., buildings or scenery). While examining other modalities is out of scope of this work, future work on this aspect is encouraged.

Despite the benefits from collecting fine-grained annotation by sensing additional human interactions, such systems without proper designs may cause severe privacy issues. We therefore encourage future research to study how to balance privacy protection and the benefits from in-situ annotations in IMT.

#### 9.4 Applications of the Object-agnostic Segmentation Model Trained on *HuTics*

One important contribution of this work is our object-agnostic segmentation model and its dataset, *HuTics*. Although our original objective was to enable LookHere, we envision that our model can be used for a broader range of applications, not limited to IMT research.

**9.4.1 Intelligent Virtual Background.** Using our model, developers can create an intelligent virtual background used in online meeting systems which is aware of the object users are trying to present to others. Segmentation algorithms for virtual backgrounds do not typically consider the behavior of object presentation. Therefore, virtual backgrounds often hide the objects held by users, diminishing user experience in certain scenarios. Our model can address this issue and show the object held by the user while preserving virtual backgrounds to support a better communication experience.

**9.4.2 Gesture-guided Portrait Mode.** Portrait mode in recent smartphones allows users to have a focus effect (e.g., blurring the background to highlight a person in the foreground). Using our model, such a portrait mode may create a focus effect on the objects held by users intelligently. Our object-agnostic segmentation model thus has a potential to enrich user experience of photo shooting with smart devices.

**9.4.3 Supports for People with Visual Impairments.** Prior work demonstrated an assistive technology for people with visual impairments by recognizing objects held by users. While it only recognized 19 objects constrained by the dataset used in that project, future assistive systems may develop a more generalizable approach by using our model trained on *HuTics*. They can first locate an object held by users using our object-agnostic segmentation model, and then apply a classification model trained on large-scale datasets that cover thousands of objects (e.g., 1000 classes in ImageNet [15]), achieving the goal of recognizing various objects for supporting activities of people with visual impairments.

## 10 CONCLUSION

This work demonstrates LookHere, a V-IMT system that allows users to annotate objects in real time during the teaching phase by exploiting users' deictic gestures. We build our own dataset (*HuTics*), consisting of 2040 front-facing images of deictic gestures and objects to achieve our implementation. Our user study results show that LookHere successfully removed substantial user effort on post-hoc manual annotations. However, the models created through LookHere did not show significant differences in their accuracies compared to those using the data with manual annotations.

## ACKNOWLEDGMENTS

A part of this research was supported by the NII CRIS collaborative research program jointly managed by NII CRIS and LINE Corporation. Co-Design Future Society Fellowship also supports the first author of this paper. We thank Xiang 'Anthony' Chen and all the reviewers for their insightful feedback on this paper. We also thank Anran Xu for his help in our data collection.

## REFERENCES

- [1] Richard A. Bolt. 1980. "Put-That-There": Voice and Gesture at the Graphics Interface. *SIGGRAPH Comput. Graph.* 14, 3 (July 1980), 262–270. <https://doi.org/10.1145/965105.807503>
- [2] Dimitrios Bounias, Ashish Singh, Spyridon Bakas, Sarthak Pati, Saima Rathore, Hamed Akbari, Michel Bilello, Benjamin A Greenberger, Joseph Lombardo, Rhea D Chitalia, et al. 2021. Interactive Machine Learning-Based Multi-Label Segmentation of Solid Tumors and Organs. *Applied Sciences* 11, 16 (2021), 7488.
- [3] Minghao Cai, Soh Masuko, and Jiro Tanaka. 2018. Gesture-Based Mobile Communication System Providing Side-by-Side Shopping Feeling. In *Proceedings of the 23rd International Conference on Intelligent User Interfaces Companion (Tokyo, Japan) (IUI '18 Companion)*. Association for Computing Machinery, New York, NY, USA, Article 2, 2 pages. <https://doi.org/10.1145/3180308.3180310>
- [4] Zhe Cao, Ilija Radosavovic, Angjoo Kanazawa, and Jitendra Malik. 2021. Reconstructing Hand-Object Interactions in the Wild. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. 12417–12426.
- [5] Michelle Carney, Barron Webster, Irene Alvarado, Kyle Phillips, Noura Howell, Jordan Griffith, Jonas Jongejan, Amit Pitaru, and Alexander Chen. 2020. Teachable Machine: Approachable Web-Based Tool for Exploring Machine Learning Classification. In *Extended Abstracts of the 2020 CHI Conference on Human Factors in Computing Systems (Honolulu, HI, USA) (CHI EA '20)*. Association for Computing Machinery, New York, NY, USA, 1–8. <https://doi.org/10.1145/3334480.3382839>
- [6] Maria Cristina Caselli. 1990. Communicative gestures and first words. In *From gesture to language in hearing and deaf children*. Springer, 56–67.
- [7] Jessica R. Cauchard, Jane L. E. Kevin Y. Zhai, and James A. Landay. 2015. Drone & Me: An Exploration into Natural Human-Drone Interaction. In *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing (Osaka, Japan) (UbiComp '15)*. Association for Computing Machinery, New York, NY, USA, 361–365. <https://doi.org/10.1145/2750858.2805823>
- [8] Chia-Ming Chang, Chia-Hsien Lee, and Takeo Igarashi. 2021. Spatial Labeling: Leveraging Spatial Layout for Improving Label Quality in Non-Expert Image Annotation. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (Yokohama, Japan) (CHI '21)*. Association for Computing Machinery, New York, NY, USA, Article 306, 12 pages. <https://doi.org/10.1145/3411764.3445165>
- [9] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. 2017. Rethinking Atrous Convolution for Semantic Image Segmentation. <https://doi.org/10.48550/ARXIV.1706.05587>
- [10] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. 2018. Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*.
- [11] Herbert H Clark. 2005. Coordinating with each other in a material world. *Discourse studies* 7, 4-5 (2005), 507–525.
- [12] Manuel Dahnert, Ji Hou, Matthias Niessner, and Angela Dai. 2021. Panoptic 3D Scene Reconstruction From a Single RGB Image. In *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, K. Nguyen, P. S. Liang, J. W. Vaughan, and Y. Dauphin (Eds.), Vol. 34. Curran Associates, Inc., 8282–8293. <https://proceedings.neurips.cc/paper/2021/file/46031b3d04de90994ca317a7c55c4289-Paper.pdf>
- [13] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Antonino Furnari, Jian Ma, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. 2021. Rescaling Egocentric Vision: Collection, Pipeline and Challenges for EPIC-KITCHENS-100. *International Journal of Computer Vision (IJCV)* (2021). <https://doi.org/10.1007/s11263-021-01531-2>
- [14] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. 2018. Scaling Egocentric Vision: The EPIC-KITCHENS Dataset. In *European Conference on Computer Vision (ECCV)*.
- [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*. 248–255. <https://doi.org/10.1109/CVPR.2009.5206848>
- [16] John J. Dudley and Per Ola Kristensson. 2018. A Review of User Interface Design for Interactive Machine Learning. *ACM Trans. Interact. Intell. Syst.* 8, 2, Article 8

- (jun 2018), 37 pages. <https://doi.org/10.1145/3185517>
- [17] Jerry Alan Fails and Dan R. Olsen. 2003. Interactive Machine Learning. In *Proceedings of the 8th International Conference on Intelligent User Interfaces* (Miami, Florida, USA) (IUI '03). Association for Computing Machinery, New York, NY, USA, 39–45. <https://doi.org/10.1145/604045.604056>
- [18] Rebecca Fiebrink, Perry R. Cook, and Dan Trueman. 2011. Human Model Evaluation in Interactive Supervised Learning. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Vancouver, BC, Canada) (CHI '11). Association for Computing Machinery, New York, NY, USA, 147–156. <https://doi.org/10.1145/1978942.1978965>
- [19] Jules Françoise, Baptiste Caramiaux, and Téo Sanchez. 2021. Marcelle: Composing Interactive Machine Learning Workflows and Interfaces. In *The 34th Annual ACM Symposium on User Interface Software and Technology* (Virtual Event, USA) (UIST '21). Association for Computing Machinery, New York, NY, USA, 39–53. <https://doi.org/10.1145/3472749.3474734>
- [20] Yolanda Gil, James Honaker, Shikhar Gupta, Yibo Ma, Vito D'Orazio, Daniel Garijo, Shruti Gadewar, Qifan Yang, and Neda Jahanshad. 2019. Towards Human-Guided Machine Learning. In *Proceedings of the 24th International Conference on Intelligent User Interfaces* (Marina del Rey, California) (IUI '19). Association for Computing Machinery, New York, NY, USA, 614–624. <https://doi.org/10.1145/3301275.3302324>
- [21] Ke Gong, Xiaodan Liang, Dongyu Zhang, Xiaohui Shen, and Liang Lin. 2017. Look Into Person: Self-Supervised Structure-Sensitive Learning and a New Benchmark for Human Parsing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [22] Sandra G Hart and Lowell E Staveland. 1988. Development of NASA-TLX (Task Load Index): Results of empirical and theoretical research. In *Advances in psychology*. Vol. 52. Elsevier, 139–183.
- [23] Richard Sahala Hartanto, Ryoichi Ishikawa, Menandro Roxas, and Takeshi Oishi. 2020. Hand-Motion-guided Articulation and Segmentation Estimation. In *2020 29th IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*, 807–813. <https://doi.org/10.1109/RO-MAN47096.2020.9223433>
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [25] Jonggi Hong, Kyungjun Lee, June Xu, and Hernisa Kacorri. 2020. Crowdsourcing the Perception of Machine Teaching. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) (CHI '20). Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi.org/10.1145/3313831.3376428>
- [26] Eric Horvitz. 1999. Principles of Mixed-Initiative User Interfaces. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Pittsburgh, Pennsylvania, USA) (CHI '99). Association for Computing Machinery, New York, NY, USA, 159–166. <https://doi.org/10.1145/302979.303030>
- [27] Lin Huang, Jianchao Tan, Ji Liu, and Junsong Yuan. 2020. Hand-Transformer: Non-Autoregressive Structured Modeling for 3D Hand Pose Estimation. In *Computer Vision – ECCV 2020*, Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm (Eds.), Springer International Publishing, Cham, 17–33.
- [28] Kurusugawa Computer Inc. 2022. *AnnoFab*. Retrieved Apr 2, 2022 from <https://annofab.com/>
- [29] Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew Davison, and Andrew Fitzgibbon. 2011. KinectFusion: Real-Time 3D Reconstruction and Interaction Using a Moving Depth Camera. In *Proceedings of the 24th Annual ACM Symposium on User Interface Software and Technology* (Santa Barbara, California, USA) (UIST '11). Association for Computing Machinery, New York, NY, USA, 559–568. <https://doi.org/10.1145/2047196.2047270>
- [30] Hernisa Kacorri, Kris M. Kitani, Jeffrey P. Bigham, and Chieko Asakawa. 2017. People with Visual Impairment Training Personal Object Recognizers: Feasibility and Challenges. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems* (Denver, Colorado, USA) (CHI '17). Association for Computing Machinery, New York, NY, USA, 5839–5849. <https://doi.org/10.1145/3025453.3025899>
- [31] Maria Karam and m. c. schraefel. 2005. *A Taxonomy of Gestures in Human Computer Interactions*. Project Report. <https://eprints.soton.ac.uk/261149/>
- [32] Yoni Kasten, Dolev Ofri, Oliver Wang, and Tali Dekel. 2021. Layered Neural Atlases for Consistent Video Editing. *ACM Trans. Graph.* 40, 6, Article 210 (dec 2021), 12 pages. <https://doi.org/10.1145/3478513.3480546>
- [33] Michael Laielli, James Smith, Giscard Biamby, Trevor Darrell, and Bjoern Hartmann. 2019. LabelAR: A Spatial Guidance Interface for Fast Computer Vision Image Collection. In *Proceedings of the 32nd Annual ACM Symposium on User Interface Software and Technology* (New Orleans, LA, USA) (UIST '19). Association for Computing Machinery, New York, NY, USA, 987–998. <https://doi.org/10.1145/3332165.3347927>
- [34] Kyungjun Lee and Hernisa Kacorri. 2019. Hands Holding Clues for Object Recognition in Teachable Machines. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (CHI '19). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300566>
- [35] Peike Li, Yunqiu Xu, Yunchao Wei, and Yi Yang. 2020. Self-Correction for Human Parsing. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2020), 1–1. <https://doi.org/10.1109/TPAMI.2020.3048039>
- [36] Andrew N Meltzoff. 1995. Understanding the intentions of others: re-enactment of intended acts by 18-month-old children. *Developmental psychology* 31, 5 (1995), 838.
- [37] Eric N Mortensen and William A Barrett. 1998. Interactive segmentation with intelligent scissors. *Graphical models and image processing* 60, 5 (1998), 349–384.
- [38] Siyou Pei, Alexander Chen, Jaewook Lee, and Yang Zhang. 2022. Hand Interfaces: Using Hands to Imitate Objects in AR/VR for Expressive Interactions. In *CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (CHI '22). Association for Computing Machinery, New York, NY, USA, Article 429, 16 pages. <https://doi.org/10.1145/3491102.3501898>
- [39] Gabriella Pizzuto and Angelo Cangelosi. 2019. Exploring Deep Models for Comprehension of Deictic Gesture-Word Combinations in Cognitive Robotics. In *2019 International Joint Conference on Neural Networks (IJCNN)*, 1–7. <https://doi.org/10.1109/IJCNN.2019.8852425>
- [40] Gonzalo A. Ramos, Christopher Meek, Patrice Y. Simard, Jina Suh, and Soroush Ghorashi. 2020. Interactive machine teaching: a human-centered approach to building machine-learned models. *Hum. Comput. Interact.* 35, 5-6 (2020), 413–451. <https://doi.org/10.1080/07370024.2020.1734931>
- [41] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*. Springer, 234–241.
- [42] Carsten Rother, Vladimir Kolmogorov, and Andrew Blake. 2004. "GrabCut": Interactive Foreground Extraction Using Iterated Graph Cuts. *ACM Trans. Graph.* 23, 3 (aug 2004), 309–314. <https://doi.org/10.1145/1015706.1015720>
- [43] Téo Sanchez, Baptiste Caramiaux, Jules Françoise, Frédéric Bevilacqua, and Wendy E. Mackay. 2021. How Do People Train a Machine? Strategies and (Mis)Understandings. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW1, Article 162 (apr 2021), 26 pages. <https://doi.org/10.1145/3449236>
- [44] Téo Sanchez, Baptiste Caramiaux, Pierre Thiel, and Wendy E. Mackay. 2022. Deep Learning Uncertainty in Machine Teaching. In *27th International Conference on Intelligent User Interfaces* (Helsinki, Finland) (IUI '22). Association for Computing Machinery, New York, NY, USA, 173–190. <https://doi.org/10.1145/3490099.3511117>
- [45] Allison Sauppé and Bilge Mutlu. 2014. Robot Deictics: How Gesture and Context Shape Referential Communication. In *Proceedings of the 2014 ACM/IEEE International Conference on Human-Robot Interaction* (Bielefeld, Germany) (HRI '14). Association for Computing Machinery, New York, NY, USA, 342–349. <https://doi.org/10.1145/2559636.2559657>
- [46] Boris Sekachev, Nikita Manovich, Maxim Zhiltsov, Andrey Zhavoronkov, Dmitry Kalinin, Ben Hoff, TOSmanov, Dmitry Kruchinin, Artyom Zankevich, Dmitry Sidnev, Maksim Markelov, Johannes222, Mathis Chenuet, a andre, telenachs, Aleksandr Melnikov, Jijoong Kim, Lion Ilouz, Nikita Glazov, Priya4607, Rush Tehrani, Seungwon Jeong, Vladimir Skubriev, Sebastian Yonekura, vugia truong, zliang7, lizhming, and Tritin Truong. 2020. *opencv/cvat: v1.1.0*. <https://doi.org/10.5281/zenodo.4009388>
- [47] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-CAM: Visual Explanations From Deep Networks via Gradient-Based Localization. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*.
- [48] Dandan Shan, Jiaqi Geng, Michelle Shu, and David F. Fouhey. 2020. Understanding Human Hands in Contact at Internet Scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [49] Patrice Y. Simard, Saleema Amershi, David M. Chickering, Alicia Edelman Pelton, Soroush Ghorashi, Christopher Meek, Gonzalo Ramos, Jina Suh, Johan Verwey, Mo Wang, and John Wernsing. 2017. Machine Teaching: A New Paradigm for Building Machine Learning Systems. <https://doi.org/10.48550/ARXIV.1707.06742>
- [50] Alison Smith-Renner, Ron Fan, Melissa Birchfield, Tongshuang Wu, Jordan Boyd-Graber, Daniel S. Weld, and Leah Findlater. 2020. *No Explainability without Accountability: An Empirical Study of Explanations and Feedback in Interactive ML*. Association for Computing Machinery, New York, NY, USA, 1–13. <https://doi.org/10.1145/3313831.3376624>
- [51] Konstantin Sofiiuk, Iliia Petrov, Olga Barinova, and Anton Konushin. 2020. f-brs: Rethinking backpropagating refinement for interactive segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8623–8632.
- [52] Konstantin Sofiiuk, Iliia A. Petrov, and Anton Konushin. 2021. Reviving Iterative Training with Mask Guidance for Interactive Segmentation. [arXiv:2102.06583 \[cs.CV\]](https://arxiv.org/abs/2102.06583)
- [53] Mingxing Tan and Quoc Le. 2019. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, 6105–6114. <https://proceedings.mlr.press/v97/tan19a.html>
- [54] Tijana Vuletic, Alex Duffy, Laura Hay, Chris McTeague, Gerard Campbell, and Madeleine Grealy. 2019. Systematic literature review of hand gestures used in

- human computer interaction interfaces. *International Journal of Human-Computer Studies* 129 (2019), 74–94. <https://doi.org/10.1016/j.ijhcs.2019.03.011>
- [55] Jacob O. Wobbrock, Meredith Ringel Morris, and Andrew D. Wilson. 2009. User-Defined Gestures for Surface Computing. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems (CHI '09)*. Association for Computing Machinery, New York, NY, USA, 1083–1092. <https://doi.org/10.1145/1518701.1518866>
- [56] Ming Yin, Jennifer Wortman Vaughan, and Hanna Wallach. 2019. Understanding the Effect of Accuracy on Trust in Machine Learning Models. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland UK) (CHI '19). Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3290605.3300509>
- [57] Hao Zhang, Zi-Hao Bo, Jun-Hai Yong, and Feng Xu. 2019. InteractionFusion: Real-Time Reconstruction of Hand Poses and Deformable Objects in Hand-Object Interactions. *ACM Trans. Graph.* 38, 4, Article 48 (jul 2019), 11 pages. <https://doi.org/10.1145/3306346.3322998>
- [58] Hao Zhang, Yuxiao Zhou, Yifei Tian, Jun-Hai Yong, and Feng Xu. 2021. Single Depth View Based Real-Time Reconstruction of Hand-Object Interactions. *ACM Trans. Graph.* 40, 3, Article 29 (jul 2021), 12 pages. <https://doi.org/10.1145/3451341>
- [59] Jing Zhang, Xin Yu, Aixuan Li, Peipei Song, Bowen Liu, and Yuchao Dai. 2020. Weakly-Supervised Salient Object Detection via Scribble Annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [60] Wencan Zhang, Mariella Dimiccoli, and Brian Y Lim. 2022. Debiased-CAM to Mitigate Image Perturbations with Faithful Visual Explanations of Machine Learning. In *CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (CHI '22). Association for Computing Machinery, New York, NY, USA, Article 182, 32 pages. <https://doi.org/10.1145/3491102.3517522>
- [61] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. 2016. Learning Deep Features for Discriminative Localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [62] Zongwei Zhou, Md Mahfuzur Rahman Siddiquee, Nima Tajbakhsh, and Jianming Liang. 2018. Unet++: A nested u-net architecture for medical image segmentation. In *Deep learning in medical image analysis and multimodal learning for clinical decision support*. Springer, 3–11.
- [63] Zhongyi Zhou and Koji Yatani. 2021. Enhancing Model Assessment in Vision-Based Interactive Machine Teaching through Real-Time Saliency Map Visualization. In *The Adjunct Publication of the 34th Annual ACM Symposium on User Interface Software and Technology* (Virtual Event, USA) (UIST '21). Association for Computing Machinery, New York, NY, USA, 112–114. <https://doi.org/10.1145/3474349.3480194>
- [64] Xiaojin Zhu. 2015. Machine teaching: An inverse problem to machine learning and an approach toward optimal education. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 29.
- [65] Xiaojin Zhu, Adish Singla, Sandra Zilles, and Anna N. Rafferty. 2018. An Overview of Machine Teaching. <https://doi.org/10.48550/ARXIV.1801.05927>

## A IMPLEMENTATION DETAILS AND EXTRA TECHNICAL RESULTS

### A.1 Configurations of Machine Learning Process in LookHere

Each image captured in the front end is fixed at the size of  $480 \times 640$  (height  $\times$  width). We chose U-Net [41] with EfficientNet-b0 backbone [53] to be the machine learning model at the back end taught by users. During the training stage, LookHere uses Adam optimizer and performs fine-tuning on the model pretrained on ImageNet [15] for 50 epochs. The batch size is four, and the learning rate is  $1e-4$ .

Note that we only use the encoder of the model for the *NaiveIMT* condition (see details in Section 7.1) because there is no ground truth data of segmentation masks in this condition, which is necessary for training the decoder. We used CAM [61] to predict saliency maps using this classification model.

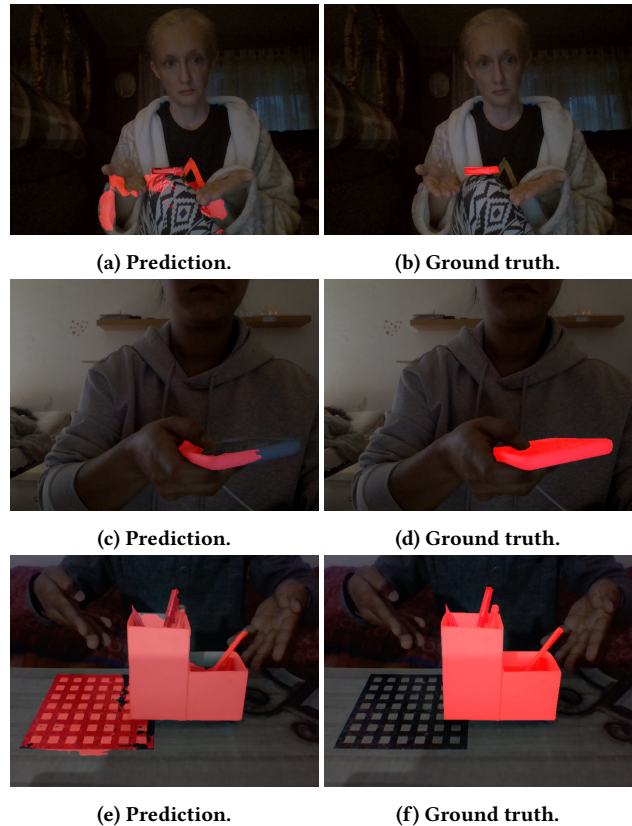


Figure 10: Additional examples of object highlight prediction failure.

Table 3: Performance comparison of the four models for our object highlight.

	U-Net	UNet++	DeepLabV3	DeepLabV3+
<i>mIoU</i>	<b>0.718</b>	0.704	0.698	0.704
fps	<b>28.3</b>	24.7	22.5	27.8

### A.2 Configurations of Training Object Highlights

We fine-tuned the network of object highlights pretrained on ImageNet using the Adam optimizer for 100 epochs. The learning rate maintains  $1e-4$  in the first 25 epochs, and drops exponentially for the subsequent 50 epochs until it reach  $1e-5$  at epoch 75. It then maintains the learning rate of  $1e-5$  in the last 25 epochs. Because there is no validation set, we simply reported the accuracy based the model achieved at the end of the training without the early stop operation. We set the batch size to be 4.

We note that this training process is for object highlights in LookHere. The previous section explains how LookHere trains the model created by users (through demonstrations of objects).

### A.3 Architecture Selections

We chose U-Net [41], UNet++ [62], DeepLabV3 [9] and DeepLabV3+ [10] for comparison because all of them are widely used different segmentation tasks and showed good performance. Table 3 shows the results in which we compared *mIoU* and FPS

of the trained models. The results show that U-Net was the most accurate as well as the fastest. Note that all of the architecture used EfficientNet-b0 as the backbone.

## **B OBJECT HIGHLIGHT PREDICTION FAILURE**

Figure 10 shows additional examples in which our predictions of object highlights have large discrepancy with their ground truth.