

# Evaluating how interactive visualizations can assist in finding samples where and how computer vision models make mistakes

Hayeong Song\*  
Georgia Institute of Technology

Gonzalo Ramos†  
Microsoft Research

Peter Bodik‡  
Microsoft Research

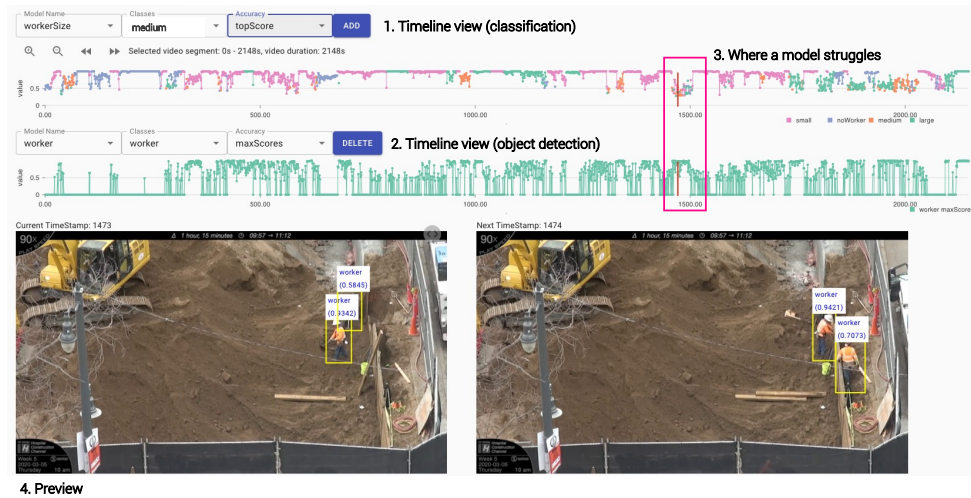


Figure 1: *Sprite*'s timeline view for a model that classifies images based on the size (small, medium, and large) of the worker in view. A user can utilize this view in the following way: **1.** The user can display the worker size classification prediction results in a timeline view, which shows prediction results over time (when a score is high, the model is confident about its prediction). Color coding shows the predicted size classification (label) of the worker for each video frame. **2.** The user can add another timeline view that shows prediction results for worker object detection model. The preview images include bounding boxes where the model detects a worker. **3.** The user can find an area in the timeline (pink rectangle) where the model seems to struggle. For example, the pink rectangle in the top portion of the preview above, the model struggled to classify a sequence of images with similar traits over a short period of time. The view differentiates its borderline prediction scores based on predicted labels, such as small (pink) and large (green). **4.** After reviewing a sequence of images, the user finds out that the model often misclassifies images of medium-sized workers. The user is then able to select meaningful new training images (with medium-sized workers) to improve the model's performance.

## ABSTRACT

Creating Computer Vision (CV) models remains a complex practice, despite their ubiquity. Access to data, the requirement for ML expertise, and model opacity are just a few points of complexity that limit the ability of end-users to build, inspect, and improve these models. Interactive ML perspectives have helped address some of these issues by considering a teacher in the loop where planning, teaching, and evaluating tasks take place. We present and evaluate two interactive visualizations in the context of *Sprite*, a system for creating CV classification and detection models for images originating from videos. We study how these visualizations help *Sprite*'s users identify (evaluate) and select (plan) images where a model is struggling and can lead to improved performance, compared to a baseline condition where users used a query language. We found that users who had used the visualizations found more images across a wider set of potential types of model errors.

**Index Terms:** Human-centered computing—Visualization—Human computer interaction (HCI)—HCI design and evaluation

\*e-mail: hsong300@gatech.edu

†e-mail: goramos@microsoft.com

‡e-mail: peterb@microsoft.com

methods—Usability testing; Human-centered computing—Visualization—Visualization application domains—Visual analytics

## 1 INTRODUCTION

Research in the field of interactive Machine Learning (ML) and teaching [9, 23] has made progress in reducing the complexity of creating ML models by having a human in the loop [2, 9, 10, 16]. This interactive loop can be characterized by curriculum building, knowledge transfer, and model evaluation stages [23]. While defining these teaching tasks is in principle straightforward, it presents challenges. In the case of curriculum building, the existence of unknowns (e.g., examples a model has not seen or features an ML model is blind to) and the presence of bias or obstacles that can get in the way of the task [28]. The transfer of knowledge can be limited to the inputs a learning algorithm permits, and challenges still persist in allowing teachers in the loop to express rich knowledge [21, 29]. Model evaluation is an activity that has a direct relationship with topics in the area of explainable ML (XAI) [1, 25]. In the context of an interactive ML model-building process, the system can only provide feedback and confidence about what it knows.

This paper describes an augmentation to *Sprite* that allows users to train CV models that get applied to images originating from video footage. We focus on the planning and evaluation stages, seeking ways not only to support subject domain experts (people with knowledge germane to the task the model is performing and those who do not necessarily have ML expertise) in efficiently identifying chal-

lenging documents (images) in which the model is not performing well but also to identify the potential reasons why. Our enhancement to *Sprite* supports users facing these challenges by improving the interactive machine teaching process so that a user and the learning system collaborate through interactive visualizations to find images that should be labeled (and added to the training dataset) to improve the CV model. We propose two interactive visualizations, *timeline view* and *scatterplot view*, in the context of *Sprite*, a system to build CV models from video footage.

We conducted a usability study to assess the effects of interactive visualizations. This study looked into how our views—the timeline and scatterplot views—help participants identify (evaluate) and select (plan) images about which a model is struggling to make correct predictions. We tested two conditions: baseline (*Sprite* without our augmentation) and visualization (*Sprite* with our augmentation). With this assessment, our work offers the following contributions. First, we’ve designed interactive visualizations and have studied how they can help users (who are training CV models) more efficiently sample images that contain diverse prediction errors in the context of *Sprite* (see Section 3). These visualizations support users in browsing and finding candidate images and in assessing and contrasting the prediction behavior of one or more models (e.g., classification and object-detection models). Second, we present the results of a study that shows that participants using these interactive visualizations were more efficient in finding candidate images that stumped the model when using a *baseline* condition. Finally, we provide insights, discuss design implications, and suggest future directions to support subject-domain experts during interactive machine teaching of CV models.

## 2 RELATED WORK

### 2.1 Interactive machine learning and teaching

There is a great body of research in the area of interactive machine learning (IML) where a person in the loop “iteratively builds and refines a mathematical model to describe a concept through iterative cycles of input and review” [8,9]. Systems like the Wekinator [11] and NorthStar [17] are but just a small example of IML experiences making ML accessible for end-users. There are different perspectives one can take when thinking about the skills, roles, and goals of people in the loop in interactive machine teaching (IMT) systems. IMT is an opinionated perspective on IML where people in the loop take on the role of a teacher of a particular concept and interact with the learning system in a “teacher-student-type of interaction”. In the IMT perspective, people (teachers) in the loop are engaged in planning and updating the teaching curriculum, explaining knowledge pertinent to the subject domain, and reviewing the learner’s progress while integrating the given knowledge. In our work, we explore IML as the general framework for helping people create CV models and use IMT’s perspective to unpack the activities people partake in while training the model. Our work focuses on the planning and reviewing stages. For a comprehensive list of articles on the topic of IML and IMT, we refer readers to these works [2,9,23,26,31].

### 2.2 Building CV models interactively

Systems and prototypes like Lobe.ai<sup>1</sup> and Google’s teachable machines<sup>2</sup> provide experiences that hide the complexity of directly interacting with a learning algorithm design and parameters and instead, present interfaces where people in the loop only need their subject-domain expertise to teach to the system. While lowering the barrier of entry for end-users, these tools only focus on part of the interactive machine teaching cycle, limiting evaluation and leaving behind tasks such as planning. Current advances in transfer learning make it possible for people to create ML models with less effort

and labels by leveraging pre-existing large models [22,32]. This makes it possible for services like Azure’s Custom Vision [20] to provide people with experiences where they only need to provide a fraction of the labels that otherwise building a robust CV model would take. Our work leverages this type of learning in its own experience, providing people with different opportunities to review a model and select the most relevant labeled data to present next.

## 3 THE *Sprite* SYSTEM

*Sprite* implements an end-to-end interactive ML workflow in the context of creating CV models using images extracted from video footage. A system like this can be used to produce predictive models to help monitor video sources such as security feeds or game streaming for content of interest. While *Sprite* follows the standard structure found in many interactive machine learning and teaching systems, *it is not the main focus of the article*; We write about it only to provide a fuller picture of our work. Our contributions do not focus on *Sprite* itself but, rather, on studying the effects of using interactive visualizations on discovering diverse examples of images about which models made prediction errors. In this section, we discuss the design goals and relevant implementation details of our contribution.

### 3.1 Design Goals

- G1. **Allow easy discovery of diverse prediction error patterns.** Improving sampling means being able to efficiently find diverse samples where a model struggles (i.e., variation of object appearances [12]).
- G2. **Support leveraging and comparing existing CV models to find images with prediction errors.** People are naturally attuned to react to both cognitive and visual dissonances or differences.
- G3. **Help subject-domain experts discover prediction patterns and errors across a video’s timeline.** Predictions over data that has a temporal component (e.g., temporal stability [15,24]) can reveal insightful patterns about a model’s performance, that cannot be observed by inspecting individual predictions in isolation.
- G4. **Increase people’s access to ML building tools by reducing the need for ML literacy.** We aspire to create experiences where subject domain experts can build CV models in efficient ways. Of the many ways to achieve this goal, we seek to leverage the direct agencies that interactive visualization can provide [7,18,34].

**General System Architecture.** *Sprite* is implemented as a web app using backend services that handle tasks, such as model training and applying the model(s) to video frames. *Sprite* performs transfer learning using *Resnet50* [14] for classification tasks, and *YOLOv4-tiny* [3] for object detection tasks. Model evaluation is performed on a prediction cluster that uses multiple GPUs to process images with high throughput. After performing an inference on a specific image, the image metadata in the database is updated with the inference output. All processed video frames are indexed in Azure Cosmos DB store, where users can retrieve relevant images using SQL queries [19] at interactive speeds.

**General System Flow & General UI.** In *Sprite*, users are engaged in the following workflow when they train and evaluate CV models: 1) upload videos (see Figure 2 (a)), 2) decode video to frames (images), 3) sample decoded images that meet a target concept, 4) label samples images (see Figure 2 (d)), 5) train CV models, and 6) evaluate model performance and predictions. To check prediction results, users can see prediction results for the in-process CV model. This paper focuses on stages 3 and 6, where people sample and evaluate decoded video frames to improve the next CV model. We specifically focus on these two phases because they are critical

<sup>1</sup><https://lobe.ai/>

<sup>2</sup><https://teachablemachine.withgoogle.com/>

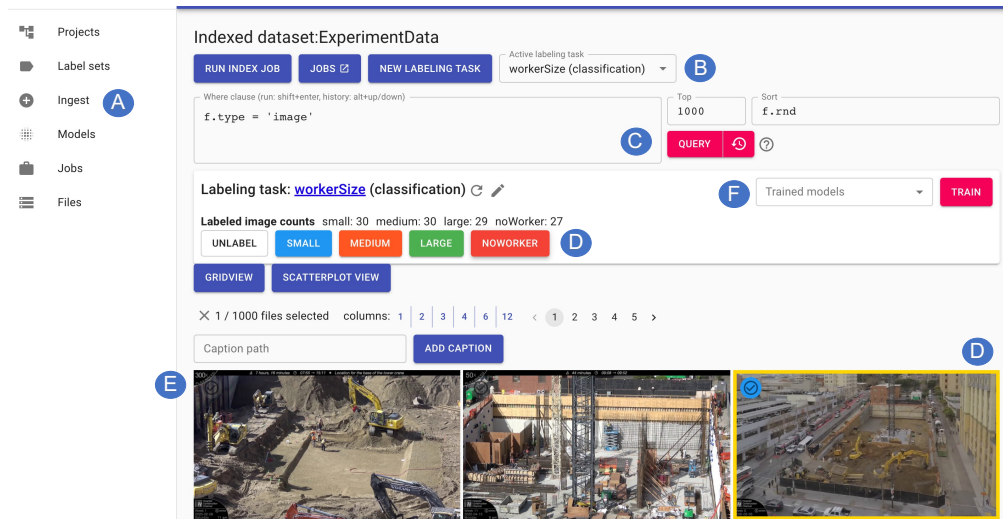


Figure 2: General UI of a *Sprite* system. A-The system’s menu shows where people can ingest (upload) videos into the system. B-Selector for the current model being worked/trained on. C-Main sampling query area to retrieve from the set of decoded images. D-Labeling selector that indicates the current classes of a selected classification model, and that is used to indicate what label to assign to a selected image. E-Grid of retrieved images from a query. F-Control to initiate a training operation and check prediction accuracy.

phases for debugging (evaluating) and improving (planning) a CV model in an interactive machine teaching loop [23].

### 3.2 Interactive Visualizations

We designed the interactive views (design process) based on 1) a review of ML exploration systems, 2) characterization and extracted design alternatives from literature, and 3) discussions on trade-offs on design alternatives. Our design was an iterative process, which involved researchers discussing design alternatives (i.e., confusion matrix [6]), prototyping, and informal usability tests with our target users (e.g., novices in ML to knowledgeable in ML) to refine our design. Following our design goals, we support two main interactive visualizations, which are *timeline view* and *scatterplot view* (see Section 3.1 for design goals). We support these views to help people in the loop to sample a set of images (labeled & unlabeled) where a CV model makes prediction errors (G4). This way, users can interact with visualization to explore different aspects of the data and model, which can allow users to formulate hypotheses on why CV models make mistakes and gain insight into the IMT process. We choose these views based on our design goals, which can help users review images on a large scale and help them understand local and global prediction patterns at a glance.

**Timeline view.** Our first design follows our goals to help people find images with potential errors by leveraging and comparing existing CV models (G2) and to discover prediction patterns/errors across a video’s timeline (G3). A timeline view provides a holistic view of prediction patterns for decoded video frames over time. This allows people to quickly find suspicious patterns such as spikes and dips globally and locally by relying on the temporal consistency in a video (see Figure 1 (1) & (2)). Users can add and stack up to three graphs in a configuration that allows them to check prediction results across different CV models. This feature can help people troubleshoot CV models by comparing multiple models (e.g., object detection and classification models). For example, if a classification model predicted an image to have a worker, but a detection model did not detect any worker, then one of the models is wrong. This disparity between the prediction results of CV models is a strong indicator that an image is worth inspecting. For *classification*, a timeline view displays results with color coding (each prediction class is represented with a different color). This view can help view

global/local patterns over time, as users can see prediction results in groups (see Figure 1 (1)). For a classification model, scores represent the prediction scores for individual classes, *topScore* is the highest score, and *topClass* is the class with the highest score (scores range from 0 - 1). Note that a prediction with a low *topScore* (e.g.,  $0.4 < \text{score} < 0.7$ ) is a sign that the model is not confident in the prediction. For *detection*, the system annotates preview image(s) with bounding boxes and shows scores of what the detection model detected (see Figure 1 (4)). For object detectors, *classes* is a list of detected classes in the image, *counts* is the number of detected objects of each class, and *maxScore* is the maximum score across all detections of a certain class.

**Scatterplot View.** Our follow-up design stems from our goals to allow people to discover diverse error patterns easily (G1) and help people easily leverage and compare existing CV models (G2). In the context of our system, we call this design the scatterplot view (see Figure 3). This view can help users find prediction errors faster because it shows similar images in clusters (based on people-specified metrics) and helps people find corner cases, where two CV models’ prediction results do not agree on. To check prediction, users can select what they want the x- and y-axis to represent to compare prediction results for two models. Users can select a model, class name (if applicable), and accuracy through the corresponding drop-down menus. For example, users can look for outliers to see if there is a disparity between the two models’ prediction results. An example of an outlier or irregularity is an image that has been classified as *noWorker* by a classification model, but where a detection model detected a worker. Then they retrieve an image of interest to inspect, by clicking on a data point on the plot (see Figure 3).

## 4 STUDY DESIGN

We wanted to assess how our proposed interactive visualizations can help users find a diverse set of examples where the model they are building makes prediction errors. In designing our study, we considered allowing participants to complete a few model training cycles to assess the success of the provided support but decided to instead observe the diversity of the examples that participants selected. This is because adding full training cycles would have significantly extended the duration of the study, compromising both human factors such as comfort and fatigue. Our study then consisted

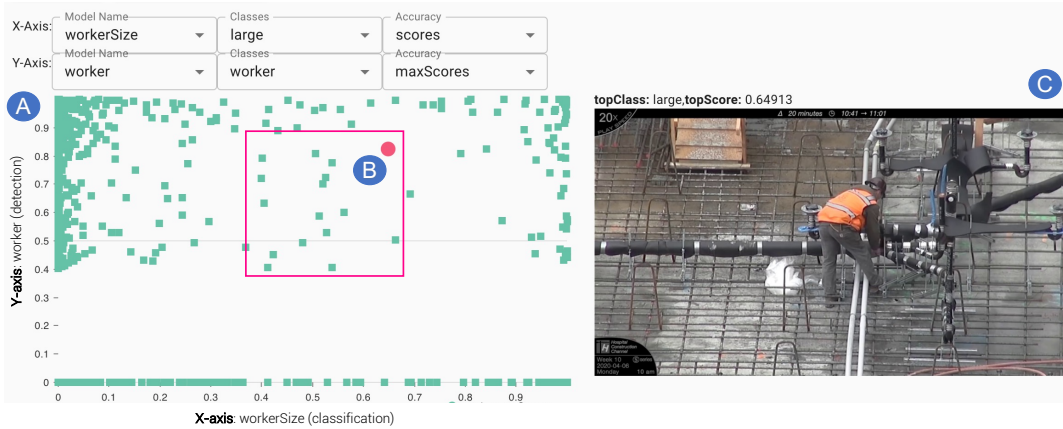


Figure 3: A - Scatterplot view used to inspect the predictions of two CV models: *workerSize* (classification) and *worker* (detection). The X-axis displays worker size classifier prediction results. The Y-axis displays prediction results for a worker detection model. B - The selected data point (the red circle) represents an image where the two models seem to agree. C - The image corresponding to the selected data point is shown along with a prediction label and score for the model currently being evaluated. To troubleshoot a classification model, a user can check data points corresponding to borderline prediction scores (pink area, e.g.  $0.4 < \text{score} < 0.7$ ). Then the user can check samples where a worker was detected (high worker detection scores) and borderline scores for a classification model, as the models seem to struggle (i.e., noisy background). This type of visualization can be helpful to quickly retrieve sub-samples of images that can help improve a model, by assisting users to compare prediction results across semantically related CV models.

of a between-subjects experiment (*baseline* and *visualization*) to understand how our visualizations support people to be more efficient in finding meaningful video frames (images) to improve a CV model. In the *baseline* condition, participants could use a query language to sample images and participants in *visualization* condition could use both query languages and visualization. We randomly divided the 20 participants into two groups of 10 each. After the study’s main task, participants filled out a survey where we collected qualitative information about their experience as well as quantitative scores to measure usability [4] and subjective workload [13]) based on a 5-point Likert survey.

**Participants.** We recruited 20 participants (13 male, 7 female), all of whom work for a large technology company. Participants came from diverse backgrounds that included design, research, and project management. Only 35% of the participants reported that they used ML for CV tasks. Participants with this type of background fit the profile of subject-domain experts, that have various levels of ML and CV model training experience. The study took about 90 minutes to complete. We compensated participants with a \$50 gift card.

**Dataset and Models.** The dataset we used was about a hospital construction site from a Youtube Channel [33] (6 videos, 40 minutes long) were uploaded in *Sprite*. For the user study, we presented participants with a system where a number of CV models were in the process of being trained. The set-up consisted of participants taking over the process of identifying sample images where in-training CV models struggled, which could then be used to fuel a subsequent labeling and training stage. The in-progress CV models we included in the study consisted of two classification models and a detection model. These **classifications** models were, first *view classifier* classifies images depending on the camera zoom level (labels: *closeup*, *medium*, and *full*) in the scene. Second, *workerSize classifier*, classifies images based on worker height compared to image size (labels: *small*, *medium*, *large*, and *noWorker*). A **detection** model detects the presence and location of construction workers in an image. All of the models had above 90% accuracy.

**Procedure.** We conducted the studies using the Teams teleconferencing platform. Participants used their laptops to access a link to the particular experimental condition we assigned them to. The general procedure for the experiment was the same. We first asked

participants to fill out a demographic survey. Then participants watched a 13-minute tutorial video, which explained the main features of *Sprite*, such as how to use queries, read prediction results, and demonstrate examples where the model made prediction errors.

**User Task.** After onboarding and practicing, we asked participants to find images, where a CV model makes prediction errors for a classification model and a detection model. We gave 20 minutes to complete each task. To retrieve images, sometimes participants had to write queries. To allow participants to explore images freely, we provided example queries participants could use and the first author (researcher) facilitated the study and served as *wizard* [27] to alleviate the need for the participants to know query language and to make the contrast between conditions as fair as possible. For the *classification* model, we asked participants to find images, where the *workerSize* model makes a prediction error. For example, the model might have classified images with a *large* worker to a *medium* worker. For the *detection model*, we asked participants to find images where the *worker* model makes prediction errors. The model might have also missed worker objects. For example, an image might have three workers, but a model might have only detected one of them. In the *baseline* condition, we asked participants to find an image where a model made prediction errors using querying language and general UI (Figure 2). For the *visualization* condition, the task was the same and participants could use query languages and interactive visualization for the task. Once participants finished the task, we asked them if they noticed any patterns that might have made a CV model struggle on those images. After the completion of the task, we followed up with a survey and interview that asked about their user experience.

## 5 RESULTS

### User Performance: Identify Images, Where CV Model Struggles.

We asked participants to capture images where the CV models in *Sprite* made prediction errors. We categorized the reasons participants hypothesized a model made an error into 11 error patterns (e.g., images without color (grayscale) or motion blur) using thematic analysis, which can be applied to classification and detection models. For our user performance analysis, we tracked the total number of images participants captured for each error pattern, and



	Baseline median, mean of ranks	Visualization median, mean of ranks	Mann-Whitney U test
<b>*Mental</b>	<b>4,13.2</b>	<b>3,7.8</b>	<b>*U = 23, X = 2, P &lt; 0.05, r = 0.45 (medium effect)</b>
Physical	1.5,10	1.5,11	U = 45, X = -0.34, P > 0.05, r = 0.07 (low effect)
Temporal	3, 11.65	2.5, 9.35	U = 38.5, X = 0.83, P > 0.05, r = 0.18 (low effect)
<b>*Effort</b>	<b>3, 13.3</b>	<b>2.5, 7.7</b>	<b>*U = 22, X = 2.08, P &lt; 0.05, r = 0.45 (medium effect)</b>
Stress	2, 12.1	1, 8.9	U = 34, X = 1.17, P > 0.05, r = 0.26 (low effect)

Table 1: Summary of participants' reported subjective workload. \* denotes significance.

how many people captured each error pattern.

For the *classification* task, participants captured images per error pattern on average was higher in the *visualization* (M = 20.27, SD = 16.14) than in the *baseline* condition (M = 6.36, SD = 7.59). Also, participants captured images during the task that led to a particular error pattern on average was higher in the *visualization* (M = 5.63, SD = 2.46) than in the *baseline* condition (M = 3.63, SD = 2.5). For the *detection* task, participants captured images per error pattern on average was higher in the *visualization* condition (M = 32.45, SD = 32.81) than in the *baseline* (M = 14.45, SD = 15.37). Also, participants captured images during the task that led to a particular error pattern on average was higher in the *visualization* (M = 6.36, SD = 3.13) than in the *baseline* condition (M = 4.27, SD = 2.49).

Our results showed that participants in *visualization* condition found more images that contained prediction errors and more variety of error patterns for both *classification* and *detection* tasks. Additionally, the total number of images captured per error was higher in the *visualization* condition. For example, participants in *visualization* condition found scarcer error patterns (e.g., scenes that were less common in the video or only a few participants identified) that participants in *baseline* condition did not find. This may be because the scatter plot view helped participants find corner cases, where the CV model struggled. Participants in the *visualization* condition found more easier to find error patterns (e.g., low light level, noisy background). These error patterns were easier to find because these scenes were common in the ingested videos. We hypothesize that this is because the timeline view assisted them in capturing similar images that exist within a time window. Similar images likely appear for a certain duration of time in a time-series view, which helped participants to sample those images quickly. This can help sample images that can later be used as new labels to potentially improve the current CV model.

**Survey Results: Usability & Subjective workload.** We measured subjective workload in performing user tasks and usability of two conditions. To analyze the data, we used Mann-Whitney's U test (a non-parametric version of the t-test).

**Subjective workload.** We asked participants to self-report 5 subjective workloads using a Likert scale (1-very low, 5-very high) [13], which were mental demand, physical demand, temporal demand, effort level, and stress level. See Table 1 for a summary.

**Usability.** We collected SUS scores [4] and our results showed that participants in *visualization* (median = 77.25 (good)) condition reported significantly higher usability scores than in the *baseline* (median = 70.25 (okay)) condition. Mann-Whitney U test showed that the mean of ranks for *baseline* was 7.15 and *visualization* was 13.85; U = 16.5, X = -2.49, P < 0.05, r = 0.55 (large effect).

## 6 DISCUSSION

We studied how two interactive visualizations can provide improved support for two fundamental activities (planning and evaluating) within the interactive machine teaching cycle when creating CV models, in the context of *Sprite*, a system that allows subject matter experts to build classifiers and extractors for images originating from videos. We assessed these techniques against the *baseline* condition and found that, first, our visualizations helped participants to better

leverage and compare existing CV models to find images with prediction errors. Second, our visualizations helped participants to more quickly and easily discover prediction patterns and errors across a video's timeline. Finally, participants under *visualization* condition, found more images across a wider set of potential types of prediction errors and reported significantly higher usability scores and significantly lower mental demand and effort levels in performing user tasks. These results suggest that our interactive visualizations provide useful support to the IML process of building CV models, by facilitating evaluation and finding the images that contain diverse prediction errors, which can be used to guide what next actions to take to continuously improve a model's performance. We discuss the following insights:

**To evaluate a model, use it with and compare it to other models.** We hypothesize that using and comparing multiple models can help people find useful samples to improve a CV model's performance. We observed that having access to these visualizations of a model's prediction results led to a better, more diverse set of examples in which the models struggled. Finding these images can directly lead to improving a model (e.g., labeling and potential feature selection). This can assist people in the loop, in the planning stage, to familiarize themselves with the existing dataset and assist them to apply insight and foresight [23, 30] to make the model-in-training better. Also, the ability to use and compare different models opens the door to the notion of what we are calling helper concepts. A helper concept involves using a simple model, using only a small sample as a microcosm representative of the larger sample, and applying the lessons derived from that microcosm to the whole.

**Support people looking at the data globally and locally.** Our observations suggest that our visualizations can lead to insights about a model's performance that go beyond inspecting a single prediction result. By comparing the outcome of the different conditions in our study, we see that participants using interactive visualizations can better assess a model's prediction patterns globally and locally. For example, we observed that the *scatterplot* view not only allowed participants to see clusters that suggested sets where different models have different opinions about the data but also allowed participants to see concentrations of data around prediction values that suggested performance issues with a model. These enhancements can facilitate image exploration and provide a guided search to improve CV models [5].

## 7 CONCLUSION AND FUTURE WORK

The goal of our work has been to improve the experience of users involved in interactive ML and teaching processes of CV models working with images originating from video. We have worked towards this goal by focusing on improving the evaluation and sampling stages of the training loop using interactive visualizations. While our work's contributions and insights exist in the context of a particular type of media-task and application—we suggest our work could be generalized, using helper models to assist in sampling engineering tasks beyond CV. Also, our insights into providing ways for end-users to inquire about model behavior locally, globally, and across time can be advantageous in other scenarios involving data with a time component, such as sensor signals.

## REFERENCES

- [1] A. Adadi and M. Berrada. Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access*, 6:52138–52160, 2018.
- [2] S. Amershi, M. Cakmak, W. B. Knox, and T. Kulesza. Power to the people: The role of humans in interactive machine learning. *Ai Magazine*, 35(4):105–120, 2014.
- [3] A. Bochkovskiy, C.-Y. Wang, and H.-Y. M. Liao. Yolov4: Optimal speed and accuracy of object detection, 2020.
- [4] J. Brooke. "SUS-A quick and dirty usability scale." *Usability evaluation in industry*. CRC Press, June 1996. ISBN: 9780748404605.
- [5] D. H. Chau, A. Kittur, J. I. Hong, and C. Faloutsos. Apollo: making sense of large network data by combining rich user interaction and machine learning. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pp. 167–176, 2011.
- [6] D. Chen, R. K. Bellamy, P. K. Malkin, and T. Erickson. Diagnostic visualization for non-expert machine learning practitioners: A design study. In *2016 IEEE Symposium on Visual Languages and Human-Centric Computing (VLHCC)*, pp. 87–95. IEEE, 2016.
- [7] J. Choo and S. Liu. Visual analytics for explainable deep learning. *IEEE computer graphics and applications*, 38(4):84–92, 2018.
- [8] H. Dang, L. Mecke, and D. Buschek. Ganslider: How users control generative models for images using multiple sliders with and without feedforward information. In *CHI Conference on Human Factors in Computing Systems*, pp. 1–15, 2022.
- [9] J. J. Dudley and P. O. Kristensson. A review of user interface design for interactive machine learning. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 8(2):1–37, 2018.
- [10] J. A. Fails and D. R. Olsen Jr. Interactive machine learning. In *Proceedings of the 8th international conference on Intelligent user interfaces*, pp. 39–45, 2003.
- [11] R. Fiebrink and P. R. Cook. The wekinator: a system for real-time, interactive machine learning in music. In *Proceedings of The Eleventh International Society for Music Information Retrieval Conference (ISMIR 2010)(Utrecht)*, vol. 3, 2010.
- [12] L. Gou, L. Zou, N. Li, M. Hofmann, A. K. Shekar, A. Wendt, and L. Ren. Vatld: a visual analytics system to assess, understand and improve traffic light detection. *IEEE transactions on visualization and computer graphics*, 27(2):261–271, 2020.
- [13] S. G. Hart. Nasa-task load index (nasa-tlx); 20 years later. In *Proceedings of the human factors and ergonomics society annual meeting*, vol. 50, pp. 904–908. Sage publications Sage CA: Los Angeles, CA, 2006.
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition, 2015.
- [15] F. Hohman, K. Wongsuphasawat, M. B. Kery, and K. Patel. Understanding and visualizing data iteration in machine learning. In *Proceedings of the 2020 CHI conference on human factors in computing systems*, pp. 1–13, 2020.
- [16] A. Holzinger, M. Plass, M. Kickmeier-Rust, K. Holzinger, G. C. Crişan, C.-M. Pintea, and V. Palade. Interactive machine learning: experimental evidence for the human in the algorithmic loop. *Applied Intelligence*, 49(7):2401–2414, 2019.
- [17] T. Kraska. Northstar: An interactive data science system. *Proceedings of the VLDB Endowment*, 11(12):2150–2164, 2018.
- [18] S. Liu, X. Wang, M. Liu, and J. Zhu. Towards better analysis of machine learning models: A visual analytics perspective. *Visual Informatics*, 1(1):48–56, 2017.
- [19] Microsoft. Getting started with sql queries, 2021.
- [20] Microsoft. What is azure custom vision?, 2021.
- [21] F. Ng, J. Suh, and G. Ramos. Understanding and supporting knowledge decomposition for machine teaching. In *ACM conference on Designing Interactive Systems (DIS)*, July 2020.
- [22] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Learning and transferring mid-level image representations using convolutional neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1717–1724, 2014.
- [23] G. Ramos, C. Meek, P. Simard, J. Suh, and S. Ghorashi. Interactive machine teaching: a human-centered approach to building machine-learned models. *Human-Computer Interaction*, 35(5-6):413–451, 2020. doi: 10.1080/07370024.2020.1734931
- [24] D. Ren, S. Amershi, B. Lee, J. Suh, and J. D. Williams. Squares: Supporting interactive performance analysis for multiclass classifiers. *IEEE transactions on visualization and computer graphics*, 23(1):61–70, 2016.
- [25] M. T. Ribeiro, S. Singh, and C. Guestrin. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144, 2016.
- [26] P. Y. Simard, S. Amershi, D. M. Chickering, A. E. Pelton, S. Ghorashi, C. Meek, G. Ramos, J. Suh, J. Verwey, M. Wang, et al. Machine teaching: A new paradigm for building machine learning systems. *arXiv preprint arXiv:1707.06742*, 2017.
- [27] A. Steinfeld, O. C. Jenkins, and B. Scassellati. The oz of wizard: simulating the human for interaction research. In *Proceedings of the 4th ACM/IEEE international conference on Human robot interaction*, pp. 101–108, 2009.
- [28] J. Suh, S. Ghorashi, G. Ramos, N.-C. Chen, S. Drucker, J. Verwey, and P. Simard. Anchorviz: Facilitating semantic data exploration and concept discovery for interactive machine learning. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 10(1):1–38, 2019.
- [29] N. Sultanum, S. Ghorashi, C. Meek, and G. Ramos. A teaching language for building object detection models. In *Proceedings of the 2020 ACM Designing Interactive Systems Conference*, pp. 1223–1234, 2020.
- [30] S. Teso, Ö. Alkan, W. Stammer, and E. Daly. Leveraging explanations in interactive machine learning: An overview. *arXiv preprint arXiv:2207.14526*, 2022.
- [31] X. Wu, L. Xiao, Y. Sun, J. Zhang, T. Ma, and L. He. A survey of human-in-the-loop for machine learning. *Future Generation Computer Systems*, 2022.
- [32] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? *arXiv preprint arXiv:1411.1792*, 2014.
- [33] YouTube. Hospital construction, 2021.
- [34] Q.-s. Zhang and S.-C. Zhu. Visual interpretability for deep learning: a survey. *Frontiers of Information Technology & Electronic Engineering*, 19(1):27–39, 2018.