# DrawMon: A Distributed System for Detection of Atypical Sketch Content in Concurrent Pictionary Games

Nikhil Bansal
Centre for Visual Information
Technology
International Institute of Information
Technology Hyderabad
Hyderabad, INDIA
nikhil.bansal@research.iiit.ac.in

Kartik Gupta
Centre for Visual Information
Technology
International Institute of Information
Technology Hyderabad
Hyderabad, INDIA
kartik.gupta0204@gmail.com

Kiruthika Kannan
Centre for Visual Information
Technology
International Institute of Information
Technology Hyderabad
Hyderabad, INDIA
kiruthika.k@research.iiit.ac.in

Sivani Pentapati
Centre for Visual Information
Technology
International Institute of Information
Technology Hyderabad
Hyderabad, INDIA
pentapatisivani27@gmail.com

Ravi Kiran Sarvadevabhatla
Centre for Visual Information
Technology
International Institute of Information
Technology Hyderabad
Hyderabad, INDIA
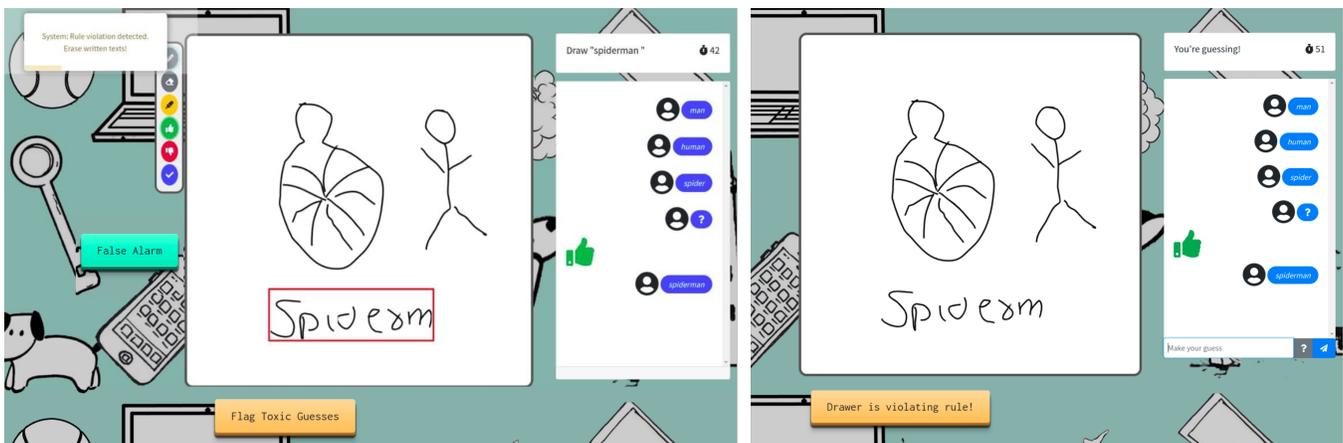ravi.kiran@iiit.ac.in

Figure 1: Screenshots of our game portal showing Drawer (left) and Guesser (right) activity during a Pictionary game. In this case, the Drawer has violated the game rules by writing text ('Spiderm') on the canvas. An automatic alert notifying the player (see top left of screenshot) and identifying the text location (red box on canvas) is generated by our system DrawMon.

## ABSTRACT

Pictionary, the popular sketch-based guessing game, provides an opportunity to analyze shared goal cooperative game play in restricted communication settings. However, some players occasionally draw atypical sketch content. While such content is occasionally relevant in the game context, it sometimes represents a rule violation and impairs the game experience. To address such situations in a timely and scalable manner, we introduce DrawMon, a novel distributed framework for automatic detection of atypical sketch content in concurrently occurring Pictionary game sessions. We build specialized online interfaces to collect game session data and annotate atypical sketch content, resulting in AtyPict, the first ever atypical sketch content dataset. We use AtyPict to train CanvasNet, a deep neural atypical content detection network. We utilize CanvasNet as a core component of DrawMon. Our analysis of post deployment game session data indicates DrawMon's effectiveness for scalable monitoring and atypical sketch content detection. Beyond Pictionary, our contributions also serve as a design guide for customized atypical content response systems involving shared and interactive whiteboards. Code and datasets are available at https://drawm0n.github.io.

## CCS CONCEPTS

• **Human-centered computing** → **Interactive systems and tools**;
• **Computing methodologies** → *Artificial intelligence*; *Scene anomaly detection.*

## KEYWORDS

whiteboard, shared interaction, anomaly detection, Pictionary, deep network, dataset, sketch

## 1 INTRODUCTION

Shared digital whiteboards are becoming increasingly popular in educational and workplace settings as a natural mechanism for collaboration and communication [1, 8, 14, 22, 26, 36]. The sharing aspect offers tremendous scope for interaction and a richer session experience. Unfortunately, shared whiteboards sometimes present situations where malicious participants draw controversial content [47]. Such activities tend to impair the collective experience of participants. Therefore, it is important to have scalable mechanisms for efficiently identifying and responding to such activities.

The popular social sketching game of Pictionary™, which we employ as a use case in this paper, also presents scenarios involving atypical sketched content. Pictionary is a wonderful example of cooperative game play to achieve a shared goal in communication-restricted settings [16, 40, 48]. The game consists of a time-limited episode involving two players - a Drawer and a Guesser. The Drawer is tasked with conveying a given target phrase to a counterpart Guesser by sketching on a whiteboard [15]. The larger the number of target phrases correctly identified and the earlier the phrases are identified from the drawn sketch, greater the number of points accrued for the participating players. The rules of Pictionary forbid the Drawer from writing text on the whiteboard. This is usually not an issue when players are physically co-located. In the anonymized, web-based version of the game, however, the Drawer may cheat by writing text related to the target word on the digitally shared whiteboard, thus violating the rules. Intervention is possible by physically monitoring game sessions. However, such manual intervention is impractical and not scalable to an online setting involving a large number of multiple concurrent game sessions. Providing user interface options for player-triggered flagging of rule violation is another possibility. But such mechanisms are not completely reliable since the Guesser benefits from the content written on the canvas and does not have real incentive to use the flagging mechanism.

Apart from malicious game play, atypical sketch content can also exist in non-malicious, benign scenarios. For instance, the Drawer may choose to draw arrows and other such icons to attract the Guesser's attention and provide indirect hints regarding the target word (see Fig. 2). Accurately localizing such activities can aid statistical learning approaches which associate sketch-based representations with corresponding target words [42]. Considering both malicious and benign scenarios, the broad requirement is for a framework which can respond to a variety of atypical whiteboard sketch content in a reliable, comprehensive and timely manner. To this end, we make the following contributions:

- We introduce ATYPICT - the first ever dataset of atypical whiteboard content.
- We introduce DRAWMON, a distributed system for sketch content-based alert generation (Sec. 5). We analyze sessions with DRAWMON deployed for Pictionary setting and demonstrate its effectiveness (Sec. 6.2).

Although presented in a Pictionary game context, our contributions serve as a design guide for developing response frameworks involving shared and interactive whiteboards. For code, models and additional details, visit the project page https://drawm0n.github.io

## 2 RELATED WORK

**Detecting and Flagging Anomalous Gameplay:** Some approaches employ a diverse mix of techniques for detecting cheating in online games [21, 46]. Dinh et al. [10] use hand-crafted game features and unsupervised machine learning approaches for offline detection of anomalous behavior. In our work, we introduce automatic deep learning based detection and flagging of anomalous gameplay in Pictionary. However, our system is also designed to detect secondary non-anomalous canvas entities which can potentially aid statistical understanding of canvas contents.

**Sketch datasets:** Existing sketch datasets (e.g. TU-Berlin [11], Sketchy [41], QuickDraw [23]) have been created primarily in the context of sketch *object* recognition problem – assign a categorical label to a hand-drawn sketch. The category labels correspond to objects (nouns). Therefore, these datasets lack abstract sketches which tend to be drawn when words from other parts of speech (verbs, adjectives) are provided as targets. Existing datasets are also unnatural because they do not include canvas actions such as erase strokes or location emphasis. Also, no intermediate guess words are associated with sketched content. For a similar reason, these datasets do not contain atypical activities unlike the dataset we introduce. Sarvadevabhatla et al. [42] explore neural network based generation of human-like guesses, but for pre-drawn object sketches. However, they do not accommodate interactivity and non-sketch drawing canvas activities (e.g. erase, pointing emphasis). The Kondate dataset [34] contains on-line handwritten patterns of text, figures, tables, maps, diagrams etc. The OHFCD dataset [2] pertains to online handwritten flowcharts. Although challenging in their own way, these datasets are considerably more structured than our setting. Additionally, they share the sketch datasets' shortcoming of being too cleanly curated because actions such as erase are absent. As a unique aspect, our combination of a game setting and a time limit unleashes greater diversity and creativity, causing sketches in our dataset to be more spontaneous and less homogeneous compared to existing datasets.

**Detecting canvas items:** Recognizing atypical activities can be thought of as a stroke segmentation problem wherein each sketch stroke is labelled as either belonging to an atypical class or the default class (drawing). Stroke segmentation has been employed

for labelling parts in object sketches either from stroke sequence information [25, 37, 50, 53] or within an image canvas [28, 51].

Recognizing atypical sketch content can also be posed as an object detection problem. In this case, the objective is to obtain 2-D spatial bounding boxes enclosing sketch strokes corresponding to the atypical content. We adopt this approach because it is faster and more amenable to near real-time operation compared to segmentation. Handwritten text is the most common atypical sketch content class in Pictionary. Hence, it is reasonable to consider approaches solely designed for text detection in domains such as outdoor scenes and documents [3, 9, 20, 29, 31, 33, 56]. Similarly, detection-based approaches have been proposed for mixed graphic structures [13, 24, 43]. However, graphic elements in these scenarios are more structured compared to our Pictionary setting. **Pictionary-like guessing games:** Borrowing terminology from the seminal work of von Ahn and Dabbish [48], Pictionary can be considered an 'inversion game' with full transparency. Riberio and Igarashi [39] employ a sketching-based interactive guessing game to progressively learn visual models of objects. A review of Pictionary-like word guessing games involving drawing can be found in the work by Sarvadevabhatla et al. [42]. In general, most of the existing works are confined to idealized toy settings [18], with some not even containing any sketching aspect [6, 12]. Unlike what we propose in this paper, they do not discuss the possibility of atypical content.

## 3 DATA COLLECTION

### 3.1 Game Sessions

Our browser-based game portal (Figures 1,5) is compatible with mouse and touch inputs, scalable and can handle up to 50 multiple concurrent Pictionary sessions. Consent is obtained and game instructions are provided when a player accesses the system for the first time. Players are assigned random names and paired randomly as Drawers and Guessers. The targets provided to the Drawers are sampled from a dictionary of 200 guess phrases. We re-emphasize that the target phrases can be nouns (e.g. airplane, bee, chair), verbs (e.g. catch, call, hang) or adjectives (e.g. happy, lazy, scary). To ensure uniform coverage across the dictionary, the probability of a guess phrase being selected for a session is inversely proportional to the number of times it has been selected for elapsed sessions. The game has a time limit of 120 seconds. The game ends when the Guesser enters a word deemed 'correct' by the Drawer or when the time limit is reached.

For the Guesser, a text box is provided for entering guess phrases. For the Drawer, the interface provides a canvas with tools to draw, erase and highlight locations (via a time-decaying spatial animation 'ping') for emphasis (see Fig. 1). In addition, 👍 and 👎 buttons enable Drawer to provide 'hot/cold' feedback on guesses. A question (❓) button is provided to the Guesser for conveying that the canvas contents are not informative and confusing. The canvas strokes are timestamped and stored in Scalable Vector Graphic (SVG) format for efficiency. In addition to canvas strokes (drawing and erasure related), secondary feedback activities mentioned previously (👎,👍, ❓, highlight) are also recorded with timestamps as part of the game session.

| Sketch Content Type class | Number of occurrences | Number of sessions containing | Number of target phrases containing |
|---|---|---|---|
| *Text* | **2419** | **478** | **180** |
| Individual letter | 2244 | 460 | 178 |
| Running hand | 175 | 103 | 81 |
| *Numbers* | **331** | **73** | **28** |
| *Circles* | **110** | **90** | **67** |
| *Iconic* | **750** | **377** | **147** |
| Arrow | 497 | 292 | 129 |
| Question mark | 158 | 116 | 78 |
| Miscellaneous | 95 | 54 | 37 |

**Table 1: Statistics of atypical sketch content categories in game sessions.**

Via our portal, we successfully gathered 3220 timestamped episodes of diverse, realistic game play involving a total of 479 participants in a large age range (14 years to 60 years) and educational demographics (middle and high school students, graduate and undergraduate university students and working professionals). Please refer to project page for sample videos of game sessions, architectural overview of the game play system and plots with additional game session statistics.

### 3.2 Atypical Data

An atypical sketch content instance can be thought of as a subsequence of sketch curves relative to the larger sequence of curves that comprise the game session. We first describe the categories of atypical content usually encountered in Pictionary sessions:

- *Text*: Drawer directly writes the target word or hints related to the target word on the canvas.
- *Numerical*: Drawer writes numbers on canvas.
- *Circles*: Drawers often circle a portion of the canvas to emphasize relevant or important content.
- *Iconic*: Other items used for emphasizing content and abstract compositional structures include drawing a question mark, arrow and other miscellaneous structures (e.g. double-headed arrow, tick marks, addition symbol, cross) and striking out the sketch (which usually implies negation of the sketched item).

Examples can be viewed in Fig. 2. It is important to remember that we consider only *Text* writing as a rule violation in Pictionary. Other categories mentioned above are atypical but their presence is not considered a violation of game rules.

To annotate atypical content, we use our custom-designed, browser-based annotation and visualization tool dubbed CanvasDash (see Fig. 3) - please refer to project page for details.

Using the described annotation procedure, we obtain our atypical Pictionary sketch dataset AtyPict. The occurrence statistics of atypical sketch categories across game sessions can be viewed in Table 1. Representative visual examples can be viewed in Fig. 2. Although we had earlier defined atypical content in terms of curve subsequences, the illustrations in Fig. 2 show that the content can have a defined 2-D spatial extent and context relative to the canvas. For ease of processing, we consider this latter interpretation. In

**Figure 2: Some examples of atypical sketch content in Pictionary game sessions are shown as canvas screenshots. The content instances span text, numbers, question marks, arrows, circles and other icons (e.g. tick marks, addition symbol) categories - refer to Sec 3.2 for details.**
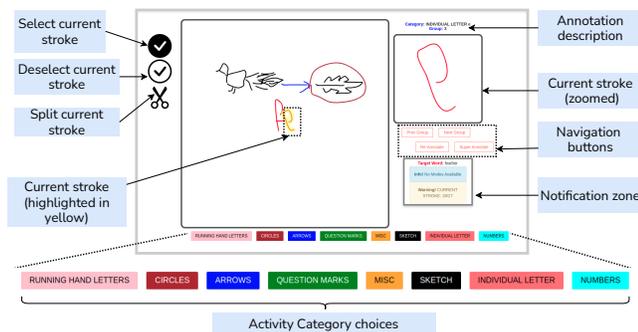


**Figure 3: An illustration of annotation using the CANVAS-DASH interface.**

other words, we consider atypical content instances to be category-labelled 2-D spatial patterns sketched over time on the canvas. In terms of Fig. 2, detecting such content therefore corresponds to accurately detecting and categorizing 2-D spatial extents of entities shown color-coded on the canvas.

## 4 CANVASNET

An effective approach for detecting atypical sketch instances needs to tackle the diversity in scale and appearance of various categories - a glance at Fig. 2 makes this amply clear. In addition, the approach needs to utilize spatial context and be robust to the presence of similar looking yet semantically distinct regular (sketch) canvas content. To meet these requirements, we cast the problem as image-based object detection. In our case, the drawing canvas containing accumulated sketch strokes is the image. Any atypical content instances (e.g. *Text*) present are considered spatially localized 2-D objects to be detected. For the object detection, we design a novel deep neural network which we dub CANVASNET. Before delving into

details of CANVASNET, we first describe the data setup employed for its training and evaluation.

### 4.1 Data Preparation

As mentioned previously, the canvas elements for a given game session are represented as timestamped SVG curve elements. Each SVG element is either a drawing stroke or an erasure stroke. We group drawing strokes into subsequences which are separated by erase strokes. The curves are converted to a 2-D point sequence representation and adaptively downsampled into line-based strokes using Ramer–Douglas–Peucker [38] algorithm ($\epsilon = 2$). The strokes are rendered on a $512 \times 512$ 2-D canvas with a stroke thickness of 4 for the purpose of data annotation and representation. Note that a drawing stroke either belongs to one of the atypical classes (Sec. 3.2) or is a normal sketch stroke. The spatial extents of labelled stroke subsequences are used to automatically generate ground-truth data for training and evaluation of CANVASNET (Sec. 4). Examples of ground-truth bounding boxes for atypical categories can be seen as solid (non-dashed/non-dotted) rectangles in Fig. 7.

**Data Augmentation:** The number of game sessions containing atypical instances are considerably smaller compared to the total number of game sessions. To increase the amount of data available for deep network training in a realistic manner, we first isolate atypical instance stroke subsequences. We sample from this set and add the resulting subsequences to other sessions which share the same target phrase, but do not contain any atypical content. The atypical content subsequences are also randomly rotated and localized carefully. This ensures they are spatially disjoint from strokes of the reference game sessions (which originally lack such atypical entities) – examples can be viewed in project page.

### 4.2 CANVASNET Deep Network

Inspired by the success of deep networks which attempt to detect text in photos [29, 32, 44], we adopt a similar efficient approach for
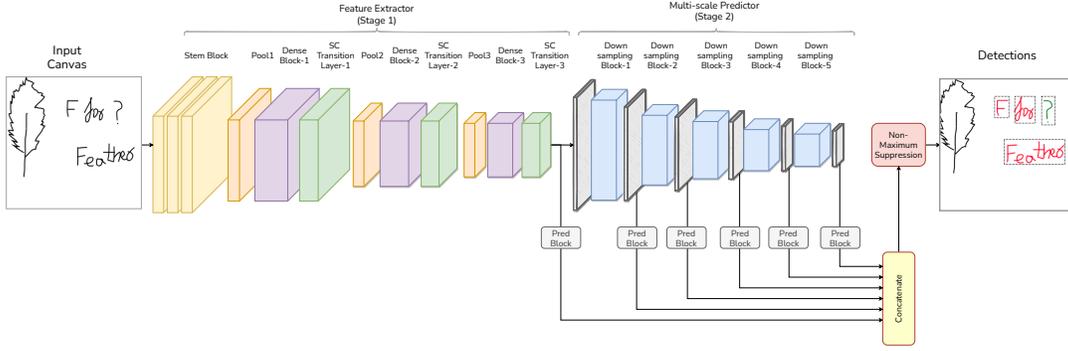
**Figure 4: Architecture of CanvasNet deep neural network. Refer to Sec. 4.2 for details.**

our CanvasNet deep network to detect atypical sketch instances on a drawing canvas. CanvasNet consists of two stages.

**Feature Extractor:** The first stage consists of a stem block containing three $3 \times 3$ unit separable convolution layers [5]. Our choice of seperable convolution layers is motivated by the reduction in number of parameters and operations involved. The first layer in stem block uses a stride of 2 to downsample the input features. The stem block is then followed by a repeating three segment structure consisting of (i) a $2 \times 2$ max pooling layer (ii) a 6 layered dense block [19] with a growth rate of 48. Each of the 6 layers consists of a $1 \times 1$ separable convolution followed by a $3 \times 3$ separable convolution (iii) a $1 \times 1$ separable convolution layer denoted as SC Transition layer. This three segment structure is then repeated three times with similar parameters - see project page for additional details.

**Multi-scale Predictor:** Atypical content (e.g. handwritten text, arrows) can occupy varying amounts of drawing canvas area, to detect them we use multi-scale predictor for prediction on multiple scales of feature maps [29, 44]. For our setting, we use a customized multi-scale predictor for both handwritten text and non-text classes. The second stage sub-network is responsible for generating multi-scale feature maps and generating predictions over each of the feature maps. The output of third segment structure of the Feature Extractor is considered the first scale of the multi-scale feature map. The other feature map scales are obtained as outputs of successive downsampling blocks applied to the Feature Extractor's output (Fig. 4) - see project page for additional details. The feature maps are passed individually through a prediction block. The multi-scale prediction features are concatenated and non-maximal suppression is applied to generate the final bounding box predictions.

*Prediction block:* This consists of a $3 \times 5$ separable convolution layer, followed by a fully connected layer comprising the prediction. The rectangular filter dimensions ($3 \times 5$) used in the block ensure that elongated objects can be detected reliably.

*Anchor boxes with vertical offsets:* Among the atypical object categories, words have larger aspect ratios and range of box orientation. Therefore, we set anchor aspect ratios to $1, 2, 3, 4, 5, 1/2, 1/3, 1/5$ with $\pm 0.25$ as the vertical offset.

**Optimization:** We formulate the loss function for CanvasNet as a combination of a classification loss $L_{cls}$ and a bounding-box localization loss $L_{loc}$:

$$L = \alpha \left( \frac{1}{N} \sum_{i=1}^{N} L_{cls}(P_i, G_i) \right) + \left( \frac{1}{M} \sum_{i=1}^{N} \sum_{j=1}^{c} G_{ij} * L_{loc}(B_i, B_i^{gt}) \right) \quad (1)$$

where $G$ is the ground truth label matrix ($G_{ij} = 1$ if $i$-th anchor box belongs to atypical category $j$, else $G_{ij} = 0$), $P$ is the predicted confidence score matrix ($P_{ij}$ indicates the confidence score that $i$-th anchor box belongs to category $j$), this means $P_i$ is a row of confidence scores for $i$-th box to belong to various atypical categories, similarly $G_i$ is a row where values for all atypical category is 0 except for one(to which $i$-th box belongs). In the above loss function formulation, $B_i^{gt}$ and $B_i$ denotes the ground truth and predicted offsets for the $i$-th anchor box. $N$ is the total number of anchor boxes, $M$ is the total number of ground truth anchor boxes not belonging to background class, $c$ represents the number of atypical categories. We use focal loss[30] for $L_{cls}$ which prioritizes a sparse set of hard examples and prevents large number of negatives from overwhelming the detector. For localisation task, we adopt Distance-IoU Loss [54]:

$$\mathcal{L}_{DIoU} = 1 - IoU + \frac{\rho^2 \left( \mathbf{b}, \mathbf{b}^{gt} \right)}{c^2}$$

where $b$ and $b^{gt}$ denote the central points of predicted and ground truth bounding box, $\rho^2 \left( \mathbf{b}, \mathbf{b}^{gt} \right)$ gives the square of the Euclidean distance between them, and $c$ is the length of the diagonal of the smallest enclosing box covering the two bounding boxes.

## 5 DRAWMON

Consider a scenario with multiple online Pictionary game sessions in progress. We require a framework for automatic and concurrent monitoring of these game sessions for any atypical activities (e.g. a rule violation such as writing text on canvas). Such a framework needs to be reliable, scalable and time-efficient. To meet these requirements, we propose DrawMon - a distributed alert generation system (see Fig. 6). Each game session is managed by a central Session Manager which assigns a unique session id (Fig. 5). For a given session, whenever a sketch stroke is drawn, the accumulated canvas content (i.e. strokes rendered so far) is tagged with session id and relayed to a shared Session Canvas Queue. For efficiency, the canvas content is represented as a lightweight Scalable Vector Graphic (SVG) object. The contents of the Session Canvas Queue

**Table 2: Results for atypical content detection**

**(a) CanvasNet performance for atypical content classes. IoU=0.5 refers to detection threshold. Refer to Sec. 6.1 for details.**

| Atypical Content Category | IoU=0.5 | |
|---|---|---|
| | mAP | mAR |
| Text | 0.58 | 0.80 |
| Number | 0.44 | 0.61 |
| Icon | 0.55 | 0.68 |
| Circle | 0.72 | 0.85 |

**(b) Performance comparison with baselines. mAP = mean Average Precision, mAR = mean Average Recall, #Parameters = the number of trainable weights in the corresponding deep network in millions, ADT = average detection time per image in milliseconds.**

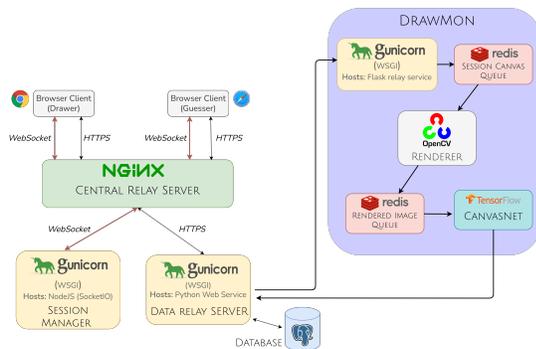| Method | Text only | | Multiclass | | # Parameters | ADT |
|---|---|---|---|---|---|---|
| | mAP | mAR | mAP | mAR | M=million | (m.sec) |
| **CanvasNet** | **0.78** | **0.90** | **0.41** | **0.53** | 1.90 M | 35 |
| BiLSTM+CRF [7] | 0.06 | 0.04 | 0.02 | 0.03 | 0.01 M | 85 |
| SketchsegNet+[37] | 0.56 | 0.32 | 0.04 | 0.11 | 3.90 M | 21 |
| Tiny-YOLOv4[49] | 0.26 | 0.65 | 0.31 | 0.51 | 5.88 M | 40 |
| TextBoxes++[29] | 0.41 | 0.65 | 0.25 | 0.39 | 29.31 M | 51 |
| DSOD[44] | 0.43 | 0.66 | 0.25 | 0.40 | 17.49 M | 52 |
| CRAFT [3] | 0.47 | 0.69 | 0.17 | 0.30 | 1.18 M | 34 |



**Figure 5: System architecture for Pictionary game setup and DrawMon (Sec. 5).**

are dequeued and rendered into corresponding $512 \times 512$ binary images by Distributed Rendering Module in a distributed and parallel fashion. The rendered binary images tagged with session id are placed in the Rendered Image Queue. The contents of Rendered Image Queue are dequeued and processed by Distributed Detection Module. Each Detection module consists of our custom-designed deep neural network CanvasNet which processes the rendered image as input. CanvasNet outputs a list of atypical activities (if any) along with associated meta-information (atypical content category, 2-D spatial location).

The outputs from multiple distributed CanvasNet instances within the Distributed Detection Module are routed to the Alert Generator Module. An activity Record Table within this module records information related to ongoing game sessions and atypical content instances. This table is analyzed with respect to a Rule Base sub-module which generates appropriate alerts and relays them to the appropriate game sessions. Since rule violations are of predominant interest, other atypical content alerts can be filtered out. Incoming alerts are finally displayed on the game session user interface (UI) - see Fig. 1.

Also note that two manually-controlled mechanisms related to alert generation exist within the game UI. The Guesser player can press a button labelled 'Drawer is violating rule!'. This simply

generates the alert (but does not highlight the canvas location where violation occurs). On the Drawer's side, the button 'False Alarm' can be used to dismiss false positive alerts (see Fig. 1). In the current deployment, we utilize the *Text* detection variant of CanvasNet to detect text writing event on canvas.

## 6 EXPERIMENTS AND RESULTS

We first describe the experiments and results for atypical content detection. Following standard machine learning protocols, we divide data into training, validation and test splits. For each target phrase, the sessions containing atypical content are randomly split in the ratio 70 (train) : 15 (validation) : 15 (test). Since we perform data augmentation on atypical content-free sessions, we divide such sessions in the aforementioned ratio as well. The respective data splits are combined to obtain the final groups.

### 6.1 Atypical content detection

**Baselines:** All along, our approach for detecting atypical activities treats the canvas as a 2-D image. In effect, the game session is considered to be a video-like frame sequence of 2-D canvas images. We also consider alternate approaches wherein the game session is processed as a sequence of curves. Each curve is labelled either as a regular sketch stroke or one associated with an atypical content category. Briefly, the baselines we consider are the following: *BiLSTM+CRF [7]* - This classifies each stroke in a sketch sequence as one of the atypical classes using a bidirectional Long Short Term Memory (BiLSTM) neural network [17] and Conditional Random Field (CRF) [45]. The input to this model is a set of hand-crafted features of the strokes. *SketchSegNet+ [37]* - This classifies each point in a sketch sequence as one of the atypical classes using bidirectional LSTM. For image based models, we train appropriately modified versions of two state-of-the-art generic object detectors – *DSOD [44]* and *Tiny-YOLOv4* [49]. We also train modified versions of two popular scene text detection models – *TextBoxes++ [29]* and *CRAFT [3]*. Please see project page for architectural details of baselines.

**Evaluation Protocol:** We conduct evaluation using two protocols. In the first protocol, we consider models trained to detect all atypical content classes. To score performance, we use the standard object
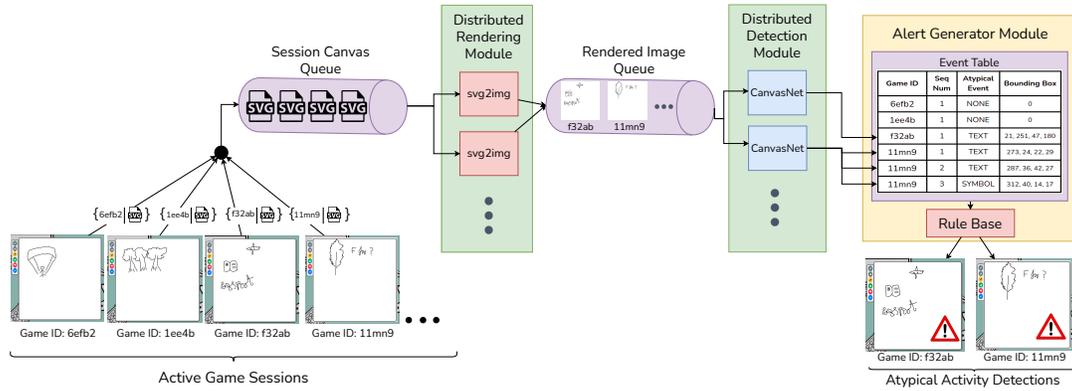
**Figure 6: A pictorial overview of DRAWMON - our distributed atypical sketch content response system (Sec. 5). Also see Fig. 5 for additional architectural details.**

detection measures – mean-average-precision (mAP) and mean-average-recall (mAR) [55]. These measures are typically reported on a $[0, 1]$ scale – larger the better. mAP and mAR are reported at an Intersection-over-Union (IoU) threshold of 0.5. In other words, an overlap of 50% or greater between predicted bounding box and ground-truth bounding box is deemed correct (assuming predicted category label also matches).

**Training:** For training CANVASNET, we employ Adam optimizer [27] with a mini-batch size of 8, the exponential decay rate for 1*st* and 2*nd* moment estimates set to 0.9 and 0.999 respectively. We stop training after 50 epochs. Classification loss weight ($\alpha$) is set to 1000 for quick convergence. The training takes approximately 4.5 hours on two GTX 1080Ti 11GB GPUs. We use *mish* [35] activation function since it provides fast convergence and improved performance compared to the standard *relu* activation. For training the text-only model, we used *hard mining training regime* wherein false-positive samples having IoU overlap value with ground truth in the range $[0.45, 0.55]$ are chosen for mining. We observed significant increase in *mAP*, *mAR* compared to regular training regime due to this regime. For training the multiclass model, we used mini-batch re-sampling with class balancing to counteract the presence of imbalanced per-class sample distribution.

**Detection Results:** CANVASNET's performance for all the atypical categories can be viewed in Table 2a. Fig. 7 depicts examples of CANVASNET detections, including some failure cases. We conducted ablation experiments to determine the relative importance of our architectural and optimization choices. Details of the ablation configurations and results can be viewed in the project page. The comparative evaluation of CANVASNET with baseline approaches is summarized in Table 2b. Note that the comparison also includes compute-related aspects (number of trainable parameters in the approaches and average detection time).

Table 2a shows CANVASNET's performance for various atypical content categories. The consistent depictions of *Circle* enables good performance for the category. Detecting isolated numbers is slightly more challenging. Empirically, we observed that sketched content resembling letters (e.g. a mountain sketch) or numerals (e.g. vertical bar groups) accounted for most of the false positive

detections. Text spanning a significantly large extent of the canvas and unusually oriented numbers accounted for majority of the missed detections (false negatives). Performance scores for ablative variants of CANVASNET are included in project page.

From Table 2b, we see that CANVASNET clearly outperforms a variety of baseline approaches (Sec. 6.1). This is predominantly due to the carefully considered architectural and optimization choices in designing CANVASNET. The results also illustrate the superiority of image-based approaches compared to the sketch stroke processing approaches (*BiLSTM +CRF, SketchSegNet+*). Keeping the rule-violation detection scenario in mind, we also trained variants designed to detect the single class *Text*. As the 'Text only' column in Table 2b shows, CANVASNET remains the best performer. Consequently, we utilize this model variant as part of DRAWMON in our game deployment scenario (Sec. 5). From the table (column named 'Parameters'), we also note that CANVASNET achieves its superior performance despite containing a smaller number of parameters relative to most of the baselines.

## 6.2 DrawMon User Study Experiments

To quantify the efficacy of DRAWMON, we analyzed game session data with DRAWMON deployed to detect text. The canvas contents are relayed to DRAWMON every 1 second. We deployed 4 CANVASNET instances within the Distributed Detection Module on two 2080Ti GPUs alongside 16 worker processes for svg to image conversion. The combined peak usage of GPU RAM was 20 GB while peak CPU RAM usage was 15 GB. 23 participants (11 male, 12 female) in the age group $19 - 25$ (mean=20.8, std.=3.1), recruited using social media and from the institution's student pool, participated in the study. Each session had an average duration of 47.5 sec (std.= 37.4) with the maximum being 120 seconds. Over the study period, the maximum number of concurrent game sessions managed by DRAWMON was 4. From the resulting set of 145 game sessions, 69 contained atypical text activities. During the sessions, we recorded timestamped alerts from DRAWMON, false alarm notifications by the Drawer player and rule violation notifications from the Guesser player. The results from the study are summarized in Table 3. To determine DRAWMON's throughput, we measured two quantities. The first, processing time (p-time), is the average elapsed time between
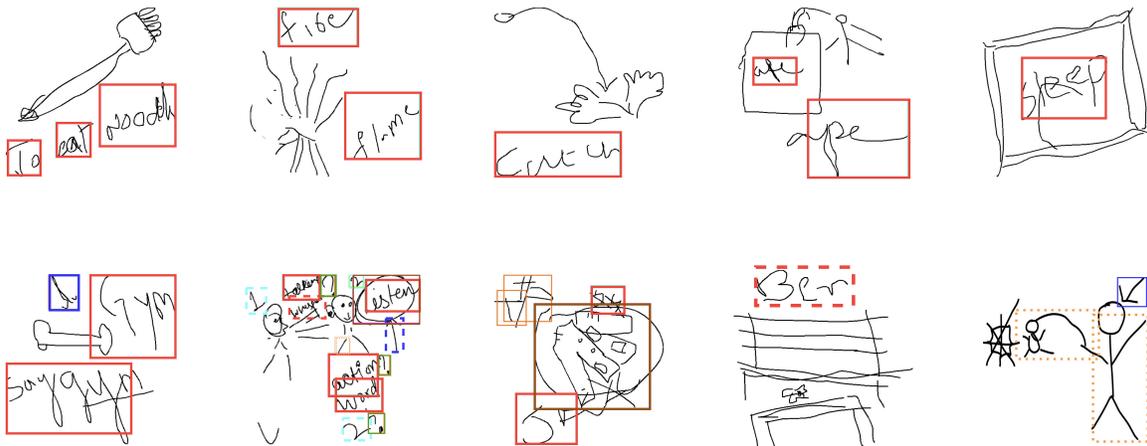
**Figure 7: Examples of atypical content detection by CanvasNet. False negatives are shown as dashed rectangles and false positives as dotted rectangles. Color codes are: text, numbers, question marks, arrows, circles and other icons (e.g. tick marks, addition symbol).**

**Table 3: User study statistics with DrawMon deployed.**

| Game Event Type | Count |
| --- | --- |
| (True Positive) DrawMon generates 'Rule Violation' alert. Drawer doesn't press 'False Alarm' button. | 62 |
| (False Positive) DrawMon generates 'Rule Violation' alert. Drawer presses 'False Alarm' button. | 32 |
| (False Negative) No 'Rule Violation' alert. Guesser presses 'Drawer is violating rule' button. | 6 |

the canvas representation being sent to DrawMon and receiving an alert. In case no alerts were generated, the timestamp corresponding to end of CanvasNet processing was considered. From our data, p-time was 0.4s. The other measurement was the maximum number of concurrently active sessions (n-sess) – this was 4. Defining the effective throughput rate to be tpr = p-time/n-sess, we obtain an average processing rate of 10 items per second.

**Results and Analysis:** The results from DrawMon deployment user study are summarized in Table 3. From the table, we see that a significant fraction of DrawMon generated alerts are valid (see 'True Positives'). From the results, DrawMon's precision is 0.66 while recall is 0.91. Post the user study, we conducted a brief survey with Likert-type questions on a 1 to 5 scale with 5 being the best. 'Q: How responsive was DrawMon to valid rule violations?' : The average score was 3.63 (s.d.=0.74), indicating reasonably high system throughput despite multiple concurrent sessions. This is also supported by the recorded throughput rate (tpr) mentioned previously in this section. 'Q: How was the overall game experience?' : The score was 3.91 (s.d.=0.60), suggesting a positive session experience and satisfaction with rule violation detection and response mechanisms.

User study plots and sample videos of game sessions with Draw-Mon in action can be viewed in project page.

### 6.3 Application Scenarios

**Application Scenarios:** Although we have used Pictionary as a use case scenario, we expect *DrawMon* to be suitable for other shared and interactive whiteboard scenarios. For instance, in a writing related setting, the notion of atypical categories can be the exact opposite of Pictionary scenario: text on canvas would be routine while drawings might be considered abnormal. This can be tackled by appropriate data labelling, for e.g. using our CanvasDash annotation tool, and subsequently retraining CanvasNet deep network. In another scenario, consider participants grouped into teams for a collaborative scene drawing task [8, 52]. DrawMon, using a CanvasNet configured for sketched scene recognition [57], can alert the instructor on progress and task completion. For this task, a CanvasNet instance trained to recognize individual objects and iconic components from our dataset (e.g. arrows, 'addition mark') could also be included as additional detection component for expanding the detection capability.

## 7 CONCLUSION AND FUTURE WORK

DrawMon is a distributed framework for monitoring multiple shared interactive whiteboards for detecting atypical content. We use a Pictionary-like sketching game as the use case scenario. Draw-Mon is enabled by a number of equally important lateral contributions - (i) CanvasDash - an intuitive dashboard UI for annotation and visualization (ii) AtyPict - a first of its kind dataset for atypical sketch content (iii) CanvasNet - a deep neural network for atypical content detection. Together, these reusable contributions create the possibility of developing similar frameworks for other shared and interactive whiteboard scenarios. Apart from atypical content detection, we expect our game session dataset to be a valuable resource in itself for analyzing player characteristics and strategies in communication restricted non-adversarial games [16]. In addition, we plan to develop practical AI agents which can mimic human Pictionary players in a more interactive, realistic manner compared to existing non-interactive works [4, 18].

# REFERENCES

[1] Salvatore Andolina, Hendrik Schneider, Joel Chan, Khalil Klouche, Giulio Jacucci, and Steven Dow. 2017. Crowdboard: augmenting in-person idea generation with real-time crowds. In *Proceedings of the 2017 ACM SIGCHI Conference on Creativity and Cognition*. 106–118.

[2] A. Awal, Guihuan Feng, H. Mouchère, and C. Viard-Gaudin. 2011. First experiments on a new online handwritten flowchart database. In *Electronic Imaging*.

[3] Youngmin Baek, Bado Lee, Dongyoon Han, Sangdoo Yun, and Hwalsuk Lee. 2019. Character Region Awareness for Text Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

[4] Ayan Kumar Bhunia, Ayan Das, Umar Riaz Muhammad, Yongxin Yang, Timothy M. Hospedales, Tao Xiang, Yulia Gryaditskaya, and Yi-Zhe Song. 2020. Pixelor: A Competitive Sketching AI Agent. So you think you can beat me?. In *SIGGRAPH Asia*.

[5] François Chollet. 2017. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1251–1258.

[6] Christopher Clark, Jordi Salvador, Dustin Schwenk, Derrick Bonafilia, Mark Yatskar, Eric Kolve, Alvaro Herrasti, Jonghyun Choi, Sachin Mehta, Sam Skjonsberg, et al. 2021. Iconary: A Pictionary-Based Game for Testing Multimodal Communication with Drawings and Text. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. 1864–1886.

[7] Amirali Darvishzadeh, Thomas Stahovich, Amir Feghahati, Negin Entezari, Shaghayegh Gharghabi, Reed Kanemaru, and Christian Shelton. 2019. CNN-BLSTM-CRF Network for Semantic Labeling of Students' Online Handwritten Assignments. In *2019 International Conference on Document Analysis and Recognition (ICDAR)*. IEEE, 1035–1040.

[8] Stephanie M Davidson, Nicole M Benson, and Scott R Beach. 2021. Drawn Together: a Curriculum for Art as a Tool in Training. *Academic Psychiatry* 45, 3 (2021), 382–387.

[9] Linjie Deng, Yanxiang Gong, Yi Lin, Jingwen Shuai, Xiaoguang Tu, Yuefei Zhang, Zheng Ma, and Mei Xie. 2019. Detecting multi-oriented text with corner-based region proposals. *Neurocomputing* 334 (Mar 2019), 134–142. https://doi.org/10.1016/j.neucom.2019.01.013

[10] Phai Vu Dinh, Thanh Nguyen Nguyen, and Quang Uy Nguyen. 2016. An empirical study of anomaly detection in online games. In *2016 3rd National Foundation for Science and Technology Development Conference on Information and Computer Science (NICS)*. 171–176. https://doi.org/10.1109/NICS.2016.7725645

[11] Mathias Eitz, James Hays, and Marc Alexa. 2012. How Do Humans Sketch Objects? *ACM Trans. Graph. (Proc. SIGGRAPH)* 31, 4 (2012), 44:1–44:10.

[12] Alaaeldin El-Nouby, Shikhar Sharma, Hannes Schulz, Devon Hjelm, Layla El Asri, Samira Ebrahimi Kahou, Yoshua Bengio, and Graham W. Taylor. 2019. Tell, Draw, and Repeat: Generating and Modifying Images Based on Continual Linguistic Instruction. In *The IEEE International Conference on Computer Vision (ICCV)*.

[13] Eyad Elyan, Laura Jamieson, and Adamu Ali-Gombe. 2020. Deep learning for symbols detection and classification in engineering drawings. *Neural networks : the official journal of the International Neural Network Society* 129 (2020), 91–102.

[14] Dikai Fang, Huahu Xu, Xiaoxian Yang, and Minjie Bian. 2019. An augmented reality-based method for remote collaborative real-time assistance: from a system perspective. *Mobile Networks and Applications* (2019), 1–14.

[15] Nicolas Fay, Bradley Walker, and Nik Swoboda. 2017. Deconstructing Social Interaction: The Complimentary Roles of Behaviour Alignment and Partner Feedback to the Creation of Shared Symbols. *Proceedings of the 39th Annual Meeting of the Cognitive Science Society* (2017), 26–29.

[16] Katy Ilonka Gero, Zahra Ashktorab, Casey Dugan, Qian Pan, James Johnson, Werner Geyer, Maria Ruiz, Sarah Miller, David R. Millen, Murray Campbell, Sadhana Kumaravel, and Wei Zhang. 2020. Mental Models of AI Agents in a Cooperative Game Setting. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 1–12. https://doi.org/10.1145/3313831.3376316

[17] Sepp Hochreiter and Jürgen Schmidhuber. 1997. LSTM can solve hard long time lag problems. *Advances in neural information processing systems* (1997), 473–479.

[18] Forrest Huang, Eldon Schoop, David Ha, and John Canny. 2020. Scones: towards conversational authoring of sketches. In *Proceedings of the 25th International Conference on Intelligent User Interfaces*. 313–323.

[19] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4700–4708.

[20] Max Jaderberg, Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2014. Reading Text in the Wild with Convolutional Neural Networks. arXiv:1412.1842 [cs.CV]

[21] Albert B. Jeng and Chia Ling Lee. 2013. A Study on Online Game Cheating and the Effective Defense. In *Recent Trends in Applied Artificial Intelligence*, Moonis Ali, Tibor Bosse, Koen V. Hindriks, Mark Hoogendoorn, Catholijn M. Jonker, and Jan Treur (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 518–527.

[22] Mads Møller Jensen, Roman Rädle, Clemens N Klokmose, and Susanne Bodker. 2018. Remediating a design tool: implications of digitizing sticky notes. In

[23] Jonas Jongejan, Henry Rowley, Takashi Kawashima, Jongmin Kim, and Nick Fox-Gieg. 2016. The Quick, Draw! AI Experiment. *Mountain View, CA* (2016).

[24] F. D. Julca-Aguilar and N. S. T. Hirata. 2018. Symbol Detection in Online Handwritten Graphics Using Faster R-CNN. In *2018 13th IAPR International Workshop on Document Analysis Systems (DAS)*. 151–156. https://doi.org/10.1109/DAS.2018.79

[25] Kurmanbek Kaiyrbekov and Metin Sezgin. 2019. Stroke-based sketched symbol reconstruction and segmentation. arXiv:1901.03427 [cs.GR]

[26] Matthew Kam, Jingtao Wang, Alastair Iles, Eric Tse, Jane Chiu, Daniel Glaser, Orna Tarshish, and John Canny. 2005. Livenotes: A System for Cooperative and Augmented Note-Taking in Lectures. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 531–540. https://doi.org/10.1145/1054972.1055046

[27] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[28] Lei Li, Hongbo Fu, and Chiew-Lan Tai. 2018. Fast Sketch Segmentation and Labeling with Deep Learning. arXiv:1807.11847 [cs.GR]

[29] Minghui Liao, Baoguang Shi, and Xiang Bai. 2018. Textboxes++: A single-shot oriented scene text detector. *IEEE transactions on image processing* 27, 8 (2018), 3676–3690.

[30] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2017. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*. 2980–2988.

[31] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. 2016. SSD: Single Shot MultiBox Detector. *Lecture Notes in Computer Science* (2016), 21–37. https://doi.org/10.1007/978-3-319-46448-0_2

[32] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. 2016. Ssd: Single shot multibox detector. In *European conference on computer vision*. Springer, 21–37.

[33] Yuliang Liu, Sheng Zhang, Lianwen Jin, Lele Xie, Yaqiang Wu, and Zhepeng Wang. 2019. Omnidirectional Scene Text Detection with Sequential-free Box Discretization. arXiv:1906.02371 [cs.CV]

[34] T. Matsushita and M. Nakagawa. 2014. A Database of On-Line Handwritten Mixed Objects Named "Kondate". In *2014 14th International Conference on Frontiers in Handwriting Recognition*. 369–374. https://doi.org/10.1109/ICFHR.2014.68

[35] Diganta Misra. 2020. Mish: A Self Regularized Non-Monotonic Activation Function. arXiv:1908.08681 [cs.LG]

[36] Anja Perlich and Christoph Meinel. 2018. Cooperative note-taking in psychotherapy sessions: An evaluation of the therapist's user experience with tele-board MED. In *2018 IEEE 20th International Conference on e-Health Networking, Applications and Services (Healthcom)*. IEEE, 1–6.

[37] Yonggang Qi and Zheng-Hua Tan. 2019. Sketchsegnet+: An end-to-end learning of rnn for multi-class sketch semantic segmentation. *IEEE Access* 7 (2019), 102717–102726.

[38] Urs Ramer. 1972. An iterative procedure for the polygonal approximation of plane curves. *Computer Graphics and Image Processing* 1, 3 (1972), 244–256. https://doi.org/10.1016/S0146-664X(72)80017-0

[39] Andre Ribeiro and Takeo Igarashi. 2012. Sketch-Editing Games: Human-Machine Communication, Game Theory and Applications. In *Proceedings of the 25th Annual ACM Symposium on User Interface Software and Technology* (Cambridge, Massachusetts, USA) (UIST '12). Association for Computing Machinery, New York, NY, USA, 287–298. https://doi.org/10.1145/2380116.2380154

[40] Melissa J. Rogerson, Martin R. Gibbs, and Wally Smith. 2018. Cooperating to Compete: The Mutuality of Cooperation and Competition in Boardgame Play. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, New York, NY, USA, 1–13. https://doi.org/10.1145/3173574.3173767

[41] Patsorn Sangkloy, Nathan Burnell, Cusuh Ham, and James Hays. 2016. The Sketchy Database: Learning to Retrieve Badly Drawn Bunnies. *ACM Transactions on Graphics (proceedings of SIGGRAPH)* (2016).

[42] R. K. Sarvadevabhatla, S. Surya, T. Mittal, and R. V. Babu. 2020. Pictionary-Style Word Guessing on Hand-Drawn Object Sketches: Dataset, Analysis and Deep Network Models. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42, 1 (2020), 221–231. https://doi.org/10.1109/TPAMI.2018.2877996

[43] Bernhard Schäfer and Heiner Stuckenschmidt. 2019. Arrow R-CNN for flowchart recognition. In *2019 International Conference on Document Analysis and Recognition Workshops (ICDARW)*, Vol. 1. IEEE, 7–13.

[44] Zhiqiang Shen, Zhuang Liu, Jianguo Li, Yu-Gang Jiang, Yurong Chen, and Xiangyang Xue. 2017. Dsod: Learning deeply supervised object detectors from scratch. In *Proceedings of the IEEE international conference on computer vision*. 1919–1927.

[45] Charles Sutton and Andrew McCallum. 2006. An introduction to conditional random fields for relational learning. *Introduction to statistical relational learning* 2 (2006), 93–128.

[46] Luan Bui The and Van Nguyen Khanh. 2010. GameGuard: A Windows-Based Software Architecture for Protecting Online Games against Hackers. In *Proceedings*

*Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. 1–12.

of the 2010 Symposium on Information and Communication Technology (Hanoi, Vietnam) (SoICT '10). Association for Computing Machinery, New York, NY, USA, 171–178. https://doi.org/10.1145/1852611.1852643

[47] John Tunell. 2018. Classification of offensive game-emblem drawings using Convolutional Neural Networks and transfer learning. Master's thesis. Uppsala University.

[48] Luis von Ahn and Laura Dabbish. 2008. Designing Games with a Purpose. Commun. ACM 51, 8 (Aug. 2008), 58–67. https://doi.org/10.1145/1378704.1378719

[49] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. 2020. Scaled-YOLOv4: Scaling Cross Stage Partial Network. arXiv preprint arXiv:2011.08036 (2020).

[50] F. Wang, S. Lin, H. Wu, H. Li, R. Wang, X. Luo, and X. He. 2019. SPFusionNet: Sketch Segmentation Using Multi-modal Data Fusion. In 2019 IEEE International Conference on Multimedia and Expo (ICME). 1654–1659. https://doi.org/10.1109/ICME.2019.00285

[51] Wenhai Wang, Enze Xie, Xiaoge Song, Yuhang Zang, Wenjia Wang, Tong Lu, Gang Yu, and Chunhua Shen. 2020. Efficient and Accurate Arbitrary-Shaped Text Detection with Pixel Aggregation Network. arXiv:1908.05900 [cs.CV]

[52] Blake Williford, Matthew Runyon, Wayne Li, Julie Linsey, and Tracy Hammond. 2020. Exploring the Potential of an Intelligent Tutoring System for Sketching Fundamentals. In Proceedings of the 2020 CHI Conference on Human Factors in

Computing Systems. Association for Computing Machinery, New York, NY, USA, 1–13. https://doi.org/10.1145/3313831.3376517

[53] Lumin Yang, Jiajie Zhuang, Hongbo Fu, Kun Zhou, and Youyi Zheng. 2020. SketchGCN: Semantic Sketch Segmentation with Graph Convolutional Networks. arXiv:2003.00678 [cs.CV]

[54] Zhaohui Zheng, Ping Wang, Wei Liu, Jinze Li, Rongguang Ye, and Dongwei Ren. 2020. Distance-IoU loss: Faster and better learning for bounding box regression. In Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34. 12993–13000.

[55] Dingfu Zhou, Jin Fang, Xibin Song, Chenye Guan, Junbo Yin, Yuchao Dai, and Ruigang Yang. 2019. Iou loss for 2d/3d object detection. In 2019 International Conference on 3D Vision (3DV). IEEE, 85–94.

[56] Xinyu Zhou, Cong Yao, He Wen, Yuzhi Wang, Shuchang Zhou, Weiran He, and Jiajun Liang. 2017. East: an efficient and accurate scene text detector. In Proceedings of the IEEE conference on Computer Vision and Pattern Recognition. 5551–5560.

[57] Changqing Zou, Qian Yu, Ruofei Du, Haoran Mo, Yi-Zhe Song, Tao Xiang, Chengying Gao, Baoquan Chen, and Hao Zhang. 2018. Sketchyscene: Richly-annotated scene sketches. In Proceedings of the european conference on computer vision (ECCV). 421–436.