

# A Human-ML Collaboration Framework for Improving Video Content Reviews

Meghana Deodhar, Xiao Ma, Yixin Cai, Alex Koes, Alex Beutel, Jilin Chen

Google

USA

{mdeodhar,xmaa,yixincai,koes,alexbeutel,jilinc}@google.com

## ABSTRACT

We deal with the problem of localized in-video taxonomic human annotation in the video content moderation domain, where the goal is to identify video segments that violate granular policies, e.g., community guidelines on an online video platform. High quality human labeling is critical for enforcement in content moderation. This is challenging due to the problem of information overload - raters need to apply a large taxonomy of granular policy violations with ambiguous definitions, within a limited review duration to relatively long videos. Our key contribution is a novel human-machine learning (ML) collaboration framework aimed at maximizing the quality and efficiency of human decisions in this setting - human labels are used to train segment-level models, the predictions of which are displayed as "hints" to human raters, indicating probable regions of the video with specific policy violations. The human verified/corrected segment labels can help refine the model further, hence creating a human-ML positive feedback loop. Experiments show improved human video moderation decision quality, and efficiency through more granular annotations submitted within a similar review duration, which enable a 5-8% AUC improvement in the hint generation models.

## KEYWORDS

human computation, machine learning, video content moderation, ranking

### ACM Reference Format:

Meghana Deodhar, Xiao Ma, Yixin Cai, Alex Koes, Alex Beutel, Jilin Chen. 2022. A Human-ML Collaboration Framework for Improving Video Content Reviews. In *Proceedings of Human-in-the-Loop Data Curation Workshop at the 31st ACM International Conference on Information and Knowledge Management (CIKM '22)*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

## 1 INTRODUCTION

The importance of content moderation on online video platforms such as TikTok, YouTube or Instagram is growing [16, 29]. These platforms strive to accurately detect the presence of policy violations within the video, which drive enforcement actions, e.g., the video can be taken down. Given the complexity of this problem,

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).

CIKM '22, October 17–22, 2022, Atlanta, GA

© 2022 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

content moderation relies heavily on human judgement and employs large teams of content moderators to perform reviews. Since human annotations directly lead to high stakes decisions, such as content take downs, the quality of the annotations is critical.

For content moderation decisions there is a growing need for transparency in detected policy violations to provide feedback to content creators [20]. This motivates a in-video taxonomic annotation task, where the goal is to provide *localized* and *fine-grained policy-specific* annotations, i.e., both the time regions (video segments) and the exact policies violated, which inform downstream video content moderation decisions. To cover the spectrum of potential violations, policy playbooks typically contain hundreds of fine-grained policies. This large space of policy violations can be organized as a taxonomy of broad categories such as Profanity, Violence, Nudity, etc., each of which contains several granular violations. For instance, Violence could include a range of granular classes such as animal abuse or graphic violence in video games. The class definitions are complex, ambiguous and often require nuanced judgment to apply, e.g., graphic violence. New policy classes may be added over time as well, e.g., Covid anti-vaccination. Moreover, there is a class imbalance issue - some egregious violations may be very rare.

Our goal is to maximize the quality and efficiency of the complex, granular, localized policy annotations task, hence leading to the correct video level enforcement decision. We achieve this by tackling the key issue of "information overload" faced by raters in providing high quality annotations, where 1) the sheer volume of videos on large online video platforms means raters only have limited review time per video; and 2) the large taxonomy of policies makes it hard for raters to recall the complete set of all granular violations for every video region they watch in the limited review duration.

We propose a human-ML collaboration framework to maximize human ratings quality and efficiency by addressing "information overload". We train models on granular rater annotations to predict policy violations, which are then combined with innovative front-end elements in the rating tool to provide "hints" to assist raters. We borrow from information retrieval literature and use ranking mechanisms for identifying the most useful and succinct set of hints. In experiments, we show that this enables raters to efficiently label policy violations more correctly and comprehensively. The human interactions with model hints pave the way for leveraging human feedback to improve the underlying ML models.

## 2 RELATED WORK

The crowd sourcing literature is very rich in the application of human annotations to perform a variety of tasks such as text processing [8, 18, 22], audio transcription [26], taxonomy creation [12], and social media analysis [14, 41]. Although there is existing literature on video annotation, it is primarily focused on identifying actions or labeling entities easily distinguishable by humans using visual information only, e.g., high jump, thunderstorms. The primary goal of these tasks is to create large datasets for facilitating Machine Learning/Perception applications [36, 40], e.g., the YouTube-8m dataset [1]. This is very different from our set up, where raters annotate granular, ambiguously defined policies using multi-modal signals - video, audio and text, from the transcript, and the video title and description. The recent emergence of crowd sourcing literature on content moderation primarily covers textual content such as user comments [9, 24], however there has been little focus on video moderation tasks.

Human-ML collaboration is an emerging area of research with two main categories of work:

**ML-Assisted Human Labeling:** ML-assistance through predictions and explanations has been used to improve the quality of human decisions in several domains [7, 25, 30], including content moderation [9, 24]. Crossmod [9], for instance, uses a model trained on historic cross-community moderation decisions to enable Reddit human moderators to find more violations. Interactive ML-assistance, which we leverage in Section 3.1.2, is used by Bartolo et al. [6] to assist human annotators to develop adversarial examples for improving a natural language question answering model. ML-assistance has been shown to also improve the efficiency of the human labeling task [3, 5, 9]. Our work aims to exploit both the human labeling quality and efficiency benefits.

**Improving ML-models Through Human Annotations:** Human annotations are useful in constructing hybrid human-ML systems that leverage the complementary strengths of both to improve the performance of ML models [4, 37]. Existing work on active learning [32, 39] shows that strategically sampling data points can reduce human workload, but the purpose is to improve machine learning models instead of assisting raters. Recent work on explainable active learning (XAL) [15] has called for better designing for the human experience in the human-AI interface.

Our work shows that it is possible to achieve model improvements and assist human raters, bridging the gap between ML-assisted human labeling and active learning. The novelty of our approach is that the models are re-trained continuously on the output of the human annotation task, which they provide assistance for, constructing a positive feedback loop between humans and models. In this collaborative framework, we have the opportunity to improve both modeling and human rater performance.

Information overload, which we encounter in our content moderation setting, is a well studied problem that reduces the effectiveness of a human’s decision making ability [11, 13]. To address this, we build on the intuition that humans find it easier to verify or correct suggestions rather than produce new annotations from scratch [8, 19]. Our ML-assistance proposal strives to select the most informative but succinct ML-based "hints" to surface to raters by drawing on the information retrieval and ranking literature.

ML-based ranking has been shown to reduce information overload effectively in electronic messaging [27] and social media [10, 23]. We draw inspiration from the learning to rank idea [21] to reduce information overload for raters.

## 3 PROPOSED HUMAN-ML COLLABORATION FRAMEWORK

The main contribution of this paper is the human-ML collaboration framework visualized in Figure 1. We use the predictions of ML models to provide assistance to human raters and evaluate the effectiveness of different user interfaces for the ML-assistance. Since the policy violation prediction task is hard for ML models, the feedback from human raters is useful to improve the models. ML-assistance enables raters to provide segment level annotations more efficiently leading to more ground truth to train/update the ML models. Additionally, we can enable raters to interact with the ML hints (accept/reject), providing direct feedback to refine the model, establishing a positive human-ML feedback loop.

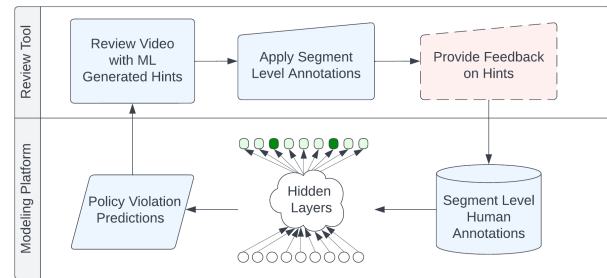


Figure 1: Human-ML collaboration set up.

### 3.1 ML-Assisted Human Reviews

As discussed in the introduction, raters face an information overload problem due to a combination of granular, complex policy definitions and limited time per video. Often, raters need to use their judgement to decide what parts of the video to watch to identify violating segments, resulting in fleeting violations being missed, or inconsistencies across raters watching different sections of the same video.

It is intuitive that raters would benefit from pointers to likely unsafe regions within a video, labeled with the exact policies being violated. Even if not completely precise, this will enable them to optimize their review bandwidth by focusing on potentially more relevant regions, making them less likely to miss violations. We achieve this by training per policy ML models and transforming their predictions into "hints" provided to raters, described in more detail in Section 3.1.2. We tune the models to be high recall to minimize uncaught violations, while relying on human judgement to improve the precision of the labeled violations.

**3.1.1 Segment Level Model Training.** To train segment level policy violation models we frame the following modeling problem - given multi-modal features for a fixed length video segment, predict whether the segment contains specific policy violations. We generate training datasets by extracting per-frame visual and audio features from the human labeled violation region. The visual

and audio features are the dense embedding layers of standard convolutional network image semantic similarity [35] and audio classification models [17] respectively. Based on empirical evidence, we select flat concatenation to aggregate the frame features over a segment, versus average/max pooling. The final model we train is a multi-label DNN model with the aggregated frame-level visual and audio features as input, where each label corresponds to a fine grained policy violation. We use MultiModal Versatile Networks [2] during model training to learn a better representation for audio and visual features for our classification task, which further improves model performance. We use a sliding window approach to utilize the trained model to generate prediction scores per policy violation for a fixed length segment starting at each frame of the video. Using a window of frame size  $n$  and stride of 1 frame, we produce model scores for segments with start and end frames  $[0 \text{ to } n-1]$ ,  $[1 \text{ to } n]$ , and so on until the end of the video, padding with empty features to fit the segment length for the last  $n$  frames.

**3.1.2 Techniques to Provide ML-Assistance.** We proceed to develop ways to use model predictions to most effectively assist human reviews, and provide details on two different designs (V1 and V2).

**V1 Hints: Continuous Line Graphs.** For video annotations, it is standard to display the video itself with playback controls and additional information in the form of a timeline [38]. For V1, we display the ML predictions as a line graph across the entire timeline of the video. The user interface is demonstrated in Figure 2. Raters can examine the line graphs and jump to the point in the video where a peak (policy violation) occurs.

While we have model predictions for hundreds of granular policies, due to visual clutter, we only display plots for a small subset of the most frequent policies. Raters also don't have the ability to provide feedback to improve model predictions.

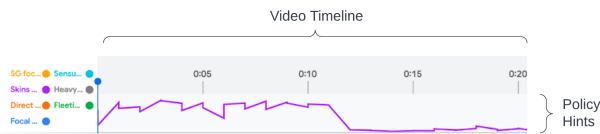


Figure 2: V1 Continuous Line Graph for Specific Policy Hints

**V2 Hints: Towards a Scalable and Interactive-ML Assistance UI.** In V2, we borrow elements from recommender systems [31, 34] to develop a more *scalable* and *interactive* interface (see Figure 3), where we pre-populate video segments that may contain policy violations detected by machine learning models.

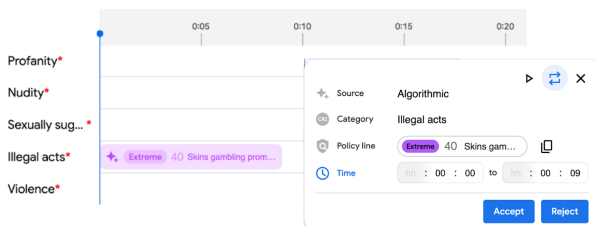


Figure 3: V2 Pre-Populated Segments from ML Models

To generate video segments, as shown in Figure 3, we introduce an algorithm to binarize continuous model scores per policy into discrete segments and use a ranking algorithm to recommend the most useful segments to raters. For simplicity, we used a threshold-based algorithm. The threshold selection constitutes a tradeoff between the precision and recall, where precision captures the utility of the predicted segments to raters and recall captures the comprehensive coverage of all video violations. The algorithm chooses a threshold maximizing recall, while maintaining a minimum precision (40% based on user studies). The regions of the video where the ML scores are above this chosen threshold are displayed as predicted policy violating segments to raters. Several heuristics are then applied to maximize segment quality, e.g., we merge segments that are close to avoid visual clutter (<3% of the whole video length apart).

Finally, to reduce information overload, we borrow from the learning to rank concept [21] to rank candidate segments and limit the number of displayed segments. The ranking algorithm prioritizes segments based on the max score across segment frames, and egregiousness of the predicted policies. We then select the top  $N$  segments to display to raters, with  $N$  selected through user studies. We pre-populate each ML suggested segment in the video timeline UI as seen in Figure 3. Raters can choose to accept or reject the suggested segments. These logged interactions can be used to provide feedback to improve the ML models.

## 3.2 Human Feedback to ML-Models

The ambiguous and fluid policy definitions along with a changing distribution of videos on online platforms poses a challenge for building robust models to accurately predict policy violations for providing ML-assistance. We hence continuously need more data and human feedback to improve the models. We show that ML-assistance increases the total number of segment labels submitted vs. no assistance. The labels in turn can serve as new ground truth for continuously re-training the models and improving performance. Additionally, with the reject button shown in V2 we collect clean negative labels; earlier, we had only "weak" negatives from videos where no violations were annotated<sup>1</sup>. Further exploiting the human feedback in combination with active learning strategies is an area of future work.

## 4 EXPERIMENTS

Our proposed methodology is evaluated using raters from the author's organization that regularly perform video content moderation reviews on live queues. Raters are separated into two pools: experts and generalists, with 150 and 400 raters respectively. The expert pool is a group of more experienced quality assurance (QA) raters with demonstrably better decision quality over a long period of time. Since their focus is QA, they don't have fixed productivity targets as generalists do. They can hence spend more time reviewing each video comprehensively, leading to higher decision quality. In our experiment setup, each video in an evaluation dataset is independently reviewed by 1 expert and 2 generalist raters. We use the labels from expert raters on the dataset as the ground truth to evaluate the 2 sets of generalist rater decisions.

<sup>1</sup>Even if no violations were annotated, they could still be present in video segments the rater did not watch, hence the negative labels inferred are weak/noisy.

## 4.1 Experimental Setup

**4.1.1 Datasets.** Our two evaluation datasets contain videos viewed on a large online video platform: (1) "Live traffic" dataset: Sampled from live traffic, hence containing a very low proportion of policy violating videos; (2) "Affected slice" dataset, sampled from live traffic and filtered to only videos with ML-hints present, containing 13-20% policy violating videos.

**4.1.2 Metrics.** Our human ratings quality metric is calculated at the video-level and conveys the correctness of the final, binary content moderation decision, e.g., take down or not. We compute the precision (P), recall (R), and disagreement rate for each of the 2 sets of generalist's video-level decisions. We consider the expert decision as ground truth and report the averaged values across the 2 sets. On live traffic datasets, we use P/R over standard inter-rater disagreement metrics due to the high class imbalance.

For rater efficiency, we measure: (1) percentage of policy violating videos where raters provide segment annotations. (2) number of segment annotations submitted by raters per video (3) average review duration per video.

## 4.2 Results

**4.2.1 Rater Quality Improvements.** We conduct experiments on both "live traffic" and "affected slice" datasets, with the baseline as a review process without ML-hints. Tables 1 and 2 compare the ratings quality metrics of our proposed V1 (line plot of model scores) ML-assistance treatment relative to the baseline, and evaluate the incremental benefit of the V2 (pre-populated segments) treatment over V1, in the V1 + V2 vs. V1 row.

Treatment	Precision	Recall	# Videos
V1 vs. Baseline	+9.82%	+1.37%	3456
V1 + V2 vs. V1	+9.97%	+5.64%	2914

**Table 1: Relative impact on live traffic (Dataset 1)**

Treatment	Precision	Recall	Disagreement%	# Videos
V1 vs. Baseline	+4.24%	+3.64%	-15.71%	3319
V1 + V2 vs. V1	+7.02%	+14.30%	-32.27%	682

**Table 2: Relative impact on affected slice (Dataset 2)**

From the live traffic results in Table 1, we see that V1 shows an improvement in precision and recall over the baseline, driven by the improvement on the affected slice as seen in Table 2.

We also see large rater quality gains of V2 over V1 on both live traffic and affected slice datasets. The segmentation and ranking algorithms in V2 allows us to overcome the scalability limitation of V1 and expand the number of granular policies covered by model hints from 7 to 18. Specifically for violence related violations, we see a 35% relative recall gain over V1 by expanding policies with ML-hints from 2 to 9. The V2 design can be scaled to cover hundreds of policies in future versions by dynamically surfacing the most relevant violating segments, further improving recall.

**4.2.2 Rater Efficiency Improvements.** We observe reduced review duration on policy violating videos with V1 hints vs. without, with more efficiency benefits on longer videos as expected; -14% on videos longer than 10 minutes and -20% on videos longer than 30 minutes. With V2, relative to V1, we see a 3% increase in review duration, but it is traded off by a 9% increase in the percentage

of policy violating videos with exact segments annotated, and a 24% increase in the number of segment annotations submitted per video.

**4.2.3 Interactive ML-Assistance Metrics.** Isolating the precision of the ML-assisted segments, we see raters accepting 35% of ML generated hints, which is in line with the 40% precision constraint we chose when converting model scores into discrete segments.

**4.2.4 Model Quality Improvements.** Since the introduction of V1 hints, we see significant model performance improvements with more human labels collected within a 3 month period on specific policy areas.

Policy Area	AUCPR	# Positive Labels
Sexually Suggestive	+5.9%	+12.7%
Nudity	+4.9%	+12.7%
Illegal Acts	+8.6%	+12.1%

**Table 3: Model Quality Improvements**

## 5 DISCUSSION

One of the potential risks of our proposed human-ML collaboration framework is automation bias [33], where a rater's over-reliance on ML-assistance can result in (i) blind-spots due to humans missing violations and (ii) raters accepting model hints without verification. Our video-level ratings quality evaluation metrics are robust to this since the ground truth comes from expert (QA) raters who review videos comprehensively, looking beyond ML-hints. In practice, we observe little evidence of both (i) and (ii). 56% of the violation segments submitted by raters in the V2 setup are organically created, i.e., don't overlap with pre-populated hint segments. The segment acceptance rate is 35%, aligned with our segmentation model precision tuning point of 40%, indicating that raters are verifying and rejecting false positive hints at the expected rate. We could mitigate the risk of (ii) further by enforcing that at least some percentage of hint segments/video is actually watched or by surfacing the model's confidence in the predicted hint to raters. To ensure robust evaluation of model quality, the AUC improvements in Section 4.2.4 are evaluated on a set of labels collected without model generated segments.

## 6 FUTURE WORK

For content moderation to scale to the size of online platforms, it is necessary to take model-based enforcement action. We would like to explore the relation between improved ground truth and improvement of automated, model based enforcement. Leveraging active learning strategies in combination with utilizing rater feedback on model generated segments to show further quality improvements in the models is another open area of research. Finally, we will explore multi-armed bandits to balance active learning based *exploration* for model improvement with model *exploitation* for providing high quality ML-assistance [28].

This paper used content moderation as the test bed for our human-ML collaboration proposal. However, it is a more generalized framework that applies to the problem of granular, localized video annotation encountered in various other industry applications such as identifying products/brands in videos to inform the placement of relevant ads, which we would like to explore further.

## REFERENCES

- [1] Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Apostol (Paul) Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan. 2016. YouTube-8M: A Large-Scale Video Classification Benchmark. In *arXiv:1609.08675*. <https://arxiv.org/pdf/1609.08675v1.pdf>
- [2] Jean-Baptiste Alayrac, Adrià Recasens, Rosalia Schneider, Relja Arandjelović, Jason Ramapuram, Jeffrey De Fauw, Lucas Smaira, Sander Dieleman, and Andrew Zisserman. 2020. Self-Supervised MultiModal Versatile Networks. <https://doi.org/10.48550/ARXIV.2006.16228>
- [3] Samreen Anjum, Ambika Verma, Brandon Dang, and Danna Gurari. 2021. Exploring the Use of Deep Learning with Crowdsourcing to Annotate Images. *Human-Computer Interaction* 8 (07 2021), 76–106. <https://doi.org/10.15346/hc.v8i2.121>
- [4] Ines Arous, Jie Yang, Mourad Khayati, and Philippe Cudré-Mauroux. 2020. *Open-Crowd: A Human-AI Collaborative Approach for Finding Social Influencers via Open-Ended Answers Aggregation*. Association for Computing Machinery, New York, NY, USA, 1851–1862. <https://doi.org/10.1145/3366423.3380254>
- [5] Zahra Ashktorab, Michael Desmond, Josh Andres, Michael Muller, Narendra Nath Joshi, Michelle Brachman, Aabhas Sharma, Kristina Brimjoin, Qian Pan, Christine T. Wolf, Evelyn Duesterwald, Casey Dugan, Werner Geyer, and Darrell Reimer. 2021. AI-Assisted Human Labeling: Batching for Efficiency without Overreliance. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW1, Article 89 (apr 2021), 27 pages. <https://doi.org/10.1145/3449163>
- [6] Max Bartolo, Tristan Thrush, Sebastian Riedel, Pontus Stenetorp, Robin Jia, and Douwe Kiela. 2021. Models in the Loop: Aiding Crowdworkers with Generative Annotation Assistants. *CoRR* abs/2112.09062 (2021). arXiv:2112.09062 <https://arxiv.org/abs/2112.09062>
- [7] Emma Beede, Elizabeth Baylor, Fred Hersch, Anna Iurchenko, Lauren Wilcox, Paison Ruamviboonsuk, and Laura M. Vardoulakis. 2020. *A Human-Centered Evaluation of a Deep Learning System Deployed in Clinics for the Detection of Diabetic Retinopathy*. Association for Computing Machinery, New York, NY, USA, 1–12. <https://doi.org/10.1145/3313831.3376718>
- [8] Michael S. Bernstein, Greg Little, Robert C. Miller, Björn Hartmann, Mark S. Ackerman, David R. Karger, David Crowell, and Katrina Panovich. 2010. Soylent: A Word Processor with a Crowd Inside. In *Proceedings of the 23rd Annual ACM Symposium on User Interface Software and Technology* (New York, New York, USA) (UIST '10). Association for Computing Machinery, New York, NY, USA, 313–322. <https://doi.org/10.1145/1866029.1866078>
- [9] Eshwar Chandrasekharan, Chaitrali Gandhi, Matthew Wortley Mustelier, and Eric Gilbert. 2019. Crossmod: A Cross-Community Learning-Based System to Assist Reddit Moderators. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 174 (nov 2019), 30 pages. <https://doi.org/10.1145/3359276>
- [10] Jilin Chen, Rowan Nairn, Les Nelson, Michael Bernstein, and Ed Chi. 2010. Short and tweet: experiments on recommending content from information streams. In *Proceedings of the SIGCHI conference on human factors in computing systems*. 1185–1194.
- [11] Yu-Chen Chen, Rong-An Shang, and Chen-Yu Kao. 2009. The effects of information overload on consumers' subjective state towards buying decision in the internet shopping environment. *Electronic Commerce Research and Applications* 8, 1 (2009), 48–58.
- [12] Lydia B Chilton, Greg Little, Darren Edge, Daniel S Weld, and James A Landay. 2013. Cascade: Crowdsourcing taxonomy creation. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*. 1999–2008.
- [13] Martin J Eppler and Jeanne Mengis. 2008. The concept of information overload—a review of literature from organization science, accounting, marketing, mis, and related disciplines (2004). *Kommunikationsmanagement im Wandel* (2008), 271–305.
- [14] Antigoni Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. 2018. Large scale crowdsourcing and characterization of twitter abusive behavior. In *Twelfth International AAAI Conference on Web and Social Media*.
- [15] Bhavya Ghai, Q. Vera Liao, Yunfeng Zhang, Rachel Bellamy, and Klaus Mueller. 2020. Explainable Active Learning (XAL): An Empirical Study of How Local Explanations Impact Annotator Experience. <https://doi.org/10.48550/ARXIV.2001.09219>
- [16] Kayla Gogarty. 2022. Hate speech and misinformation proliferate on Meta products, with 13,500 policy violations documented in the past year alone. <https://www.mediamatters.org/facebook/hate-speech-and-misinformation-proliferate-meta-products-13500-policy-violations>.
- [17] Shawn Hershey, Sourish Chaudhuri, Daniel P. W. Ellis, Jort F. Gemmeke, Aren Jansen, R. Channing Moore, Manoj Plakal, Devin Platt, Rif A. Saurous, Bryan Seybold, Malcolm Slaney, Ron J. Weiss, and Kevin W. Wilson. 2016. CNN Architectures for Large-Scale Audio Classification. *CoRR* abs/1609.09430 (2016). arXiv:1609.09430 <http://arxiv.org/abs/1609.09430>
- [18] Chang Hu, Benjamin B. Bederson, Philip Resnik, and Yakov Kronrod. 2011. MonoTrans2: A New Human Computation System to Support Monolingual Translation. 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, 1133–1136. <https://doi.org/10.1145/1978942.1979111>
- [19] Hengdong Hu, Lingxi Xie, Zewei Du, Richang Hong, and Qi Tian. 2020. One-bit Supervision for Image Classification. (2020). <https://doi.org/10.48550/ARXIV.2009.06168>
- [20] Shagun Jhaver, Amy Bruckman, and Eric Gilbert. 2019. Does Transparency in Moderation Really Matter? User Behavior After Content Removal Explanations on Reddit. *Proc. ACM Hum.-Comput. Interact.* 3, CSCW, Article 150 (nov 2019), 27 pages. <https://doi.org/10.1145/3359252>
- [21] Alexandros Karatzoglou, Linas Baltrunas, and Yue Shi. 2013. Learning to rank for recommender systems. In *Proceedings of the 7th ACM Conference on Recommender Systems*. 493–494.
- [22] Aniket Kittur, Ed H. Chi, and Bongwon Suh. 2008. Crowdsourcing User Studies with Mechanical Turk. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Florence, Italy) (CHI '08). Association for Computing Machinery, New York, NY, USA, 453–456. <https://doi.org/10.1145/1357054.1357127>
- [23] Ksenia Koroleva and Antonio José Bolufé Röhrler. 2012. Reducing information overload: Design and evaluation of filtering & ranking algorithms for social networking sites. (2012).
- [24] Vivian Lai, Samuel Carton, Rajat Bhatnagar, Q. Vera Liao, Yunfeng Zhang, and Chenhao Tan. 2022. Human-AI Collaboration via Conditional Delegation: A Case Study of Content Moderation. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems* (New Orleans, LA, USA) (CHI '22). Association for Computing Machinery, New York, NY, USA, Article 54, 18 pages. <https://doi.org/10.1145/3491102.3501999>
- [25] Vivian Lai and Chenhao Tan. 2019. On Human Predictions with Explanations and Predictions of Machine Learning Models: A Case Study on Deception Detection. In *Proceedings of the Conference on Fairness, Accountability, and Transparency* (Atlanta, GA, USA) (FAT\* '19). Association for Computing Machinery, New York, NY, USA, 29–38. <https://doi.org/10.1145/3287560.3287590>
- [26] Walter Lasecki, Christopher Miller, Adam Sadilek, Andrew Abumoussa, Donato Borrello, Raja Kushalnagar, and Jeffrey Bigham. 2012. *Real-Time Captioning by Groups of Non-Experts*. Association for Computing Machinery, New York, NY, USA, 23–34. <https://doi.org/10.1145/2380116.2380122>
- [27] Robert M Losee Jr. 1989. Minimizing information overload: the ranking of electronic messages. *Journal of Information Science* 15, 3 (1989), 179–189.
- [28] James McInerney, Benjamin Lacker, Samantha Hansen, Karl Higley, Hugues Bouchard, Alois Gruson, and Rishabh Mehrotra. 2018. Explore, exploit, and explain: personalizing explainable recommendations with bandits. In *Proceedings of the 12th ACM conference on recommender systems*. 31–39.
- [29] Sophie Mellor. 2022. TikTok slammed for videos sharing false information about Russia's war on Ukraine. <https://fortune.com/2022/03/21/tiktok-misinformation-ukraine/>.
- [30] Junwon Park, Ranjay Krishna, Pranav Khadpe, Li Fei-Fei, and Michael Bernstein. 2019. AI-Based Request Augmentation to Increase Crowdsourcing Participation. *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing* 7, 1 (Oct. 2019), 115–124. <https://ojs.aaai.org/index.php/HCOMP/article/view/5282>
- [31] Ivens Portugal, Paulo Alencar, and Donald Cowan. 2018. The use of machine learning algorithms in recommender systems: A systematic review. *Expert Systems with Applications* 97 (2018), 205–227.
- [32] Burr Settles. 2009. *Active Learning Literature Survey*. Computer Sciences Technical Report 1648. University of Wisconsin–Madison.
- [33] Linda J Skitka, Kathleen L Mosier, and Mark Burdick. 1999. Does automation bias decision-making? *International Journal of Human-Computer Studies* 51, 5 (1999), 991–1006.
- [34] Kirsten Swearingen and Rashmi Sinha. 2002. Interaction design for recommender systems. In *Designing Interactive Systems*, Vol. 6. Citeseer, 312–334.
- [35] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. 2015. Rethinking the Inception Architecture for Computer Vision. <https://doi.org/10.48550/ARXIV.1512.00567>
- [36] Milo Z Trujillo, Mauricio Gruppi, Cody Buntain, and Benjamin D Horne. 2022. The MeLa BitChute Dataset. In *Proceedings of the International AAAI Conference on Web and Social Media*, Vol. 16. 1342–1351.
- [37] Jennifer Wortman Vaughan. 2018. Making Better Use of the Crowd: How Crowdsourcing Can Advance Machine Learning Research. *Journal of Machine Learning Research* 18, 193 (2018), 1–46. <http://jmlr.org/papers/v18/17-234.html>
- [38] Carl Vondrick, Deva Ramanan, and Donald Patterson. 2010. Efficiently scaling up video annotation with crowdsourced marketplaces. In *European Conference on Computer Vision*. Springer, 610–623.
- [39] Yi Yang, Zhigang Ma, Feiping Nie, Xiaojun Chang, and Alexander G Hauptmann. 2015. Multi-class active learning by uncertainty sampling with diversity maximization. *International Journal of Computer Vision* 113, 2 (2015), 113–127.
- [40] Hang Zhao, Antonio Torralba, Lorenzo Torresani, and Zhicheng Yan. 2017. HACs: Human Action Clips and Segments Dataset for Recognition and Temporal Localization. <https://doi.org/10.48550/ARXIV.1712.09374>
- [41] Arkaitz Zubiaga, Maria Liakata, Rob Procter, Kalina Bontcheva, and Peter Tolmie. 2015. Crowdsourcing the annotation of rumourous conversations in social media. In *Proceedings of the 24th international conference on World Wide Web*. 347–353.