
Identifying the Context Shift between Test Benchmarks and Production Data

Matthew Groh
MIT Media Lab
groh@mit.edu

Abstract

Machine learning models are often brittle on production data despite achieving high accuracy on benchmark datasets. Benchmark datasets have traditionally served dual purposes: first, benchmarks offer a standard on which machine learning researchers can compare different methods, and second, benchmarks provide a model, albeit imperfect, of the real world. The incompleteness of test benchmarks (and the data upon which models are trained) hinder robustness in machine learning, enable short-cut learning, and leave models systematically prone to err on out-of-distribution and adversarially perturbed data. The mismatch between a single static benchmark dataset and a production dataset has traditionally been described as a dataset shift (or distribution shift with subcategories including covariate shift, prior probability shift, and concept shift). These shifts are simultaneously over-specified with formal definitions for comparing two data samples and under-specified for evaluating the data-generating process that drives the mismatch between data samples. In an effort to clarify how to address the mismatch between test benchmarks and production data, we introduce context shift to describe semantically meaningful changes in the underlying data generation process. Moreover, we identify three methods for addressing context shift that would otherwise lead to model prediction errors: first, we describe how human intuition and expert knowledge can identify semantically meaningful features upon which models systematically fail, second, we detail how dynamic benchmarking – with its focus on capturing the data generation process – can promote generalizability through corroboration, and third, we highlight that clarifying a model’s limitations can reduce unexpected errors. Robust machine learning is focused on model performance beyond benchmarks, and as such, we consider three model organism domains – facial expression recognition, deepfake detection, and medical diagnosis – to highlight how implicit assumptions in benchmark tasks lead to errors in practice. By paying close attention to the role of context in a prediction task, researchers can design more comprehensive benchmarks, reduce context shift errors, and increase generalization performance.

1 Machine Learning Models are Brittle in Production

Dataset benchmarks offer a standard for comparing and evaluating the performance of machine learning models on real-world tasks like object detection (1), handwritten digit recognition (2), image captioning (3), general language understanding (4), affect recognition (5), deepfake detection (6), medical diagnosis (e.g. for skin disease (7), pneumonia (8), critical care (9), etc.), and many other tasks. As a standard for comparison, dataset benchmarks have enabled rapid progress in computer vision and natural language processing.

Despite intentions to create and curate data that match the real-world as closely as possible, the dynamic, high-dimensional, combinatoric complexity of many real-world tasks is often difficult to capture in a single static benchmark. Indeed, the development and evaluation of machine learning

models on benchmarks often suffer from a variety of historical, representational, measurement, aggregation, and evaluation biases (10). These biases can be further exacerbated by deployment biases where the task that a benchmark is intended to measure differs from the real-world task (11). Moreover, data for benchmarks are often collected at scale with minimal oversight (12), which leaves data open to poisoning attacks (13), leakage (14), multiple interpretations (15) and error (16). As a consequence, machine learning models that appear to be approaching (and sometimes surpassing) human-level ability on a test benchmark will often error when shown out-of-distribution (17) data. In other words, the reliance on static test benchmarks as metrics for projecting production performance (18) inflates the accuracy of machine learning model performance and leaves open the questions, “Can you trust your model? Will it work in deployment?” (19)

The meaning of out-of-distribution data depends on a task’s context. Two canonical examples of out-of-distribution data in object detection tasks are images of either a cow on a sandy beach or a camel on a green pasture (20). Today’s commonly used training data rarely contain such animal-environment pairs, and as a result, machine learning models often learn spurious correlations such as cloven hoofed mammals next to sand are camels but the ones next to grass are cows. With *a priori* knowledge of potentially spurious correlations, one approach for addressing this kind of model brittleness is to include auxiliary labels that can serve as a causally-motivated regularization framework (21). However, post hoc model explanations are often ineffective for identifying previously unknown shortcuts (22) (though both explanations via concept traversals (23) and identifying model failures as directions in latent space via contrastive learning where images and natural language are embedded in a shared latent space show promise (24)). In contrast, human intuition can identify many out-of-distribution contexts on which spurious correlations (sometimes called shortcut learning) may occur.

In one of the clearest examples of spurious correlations that lead to the benchmark-production gap, researchers recreated ImageNet (1) and CIFAR-10 (25) with news data and demonstrated that the state-of-the-art models’ performance is significantly lower on the recreated versions of these datasets (26). The benchmark-production gap is particularly salient in this example because these two datasets have been the most commonly used benchmarks for object recognition over the last decade. Recht et al explain that the drop in performance does not appear to be explained by random sampling error, hyperparameter tuning for optimizing performance on the original test set, or obvious changes in semantically meaningful features, but instead, the performance gap appears to arise from subtle changes in the data (26). Object recognition is not as straightforward a task as it might appear at first glance and involves edge cases arising from a variety of contexts.

In complex human-centered machine learning applications, a task’s context involves answers to the following kinds of questions: What is the task? For whom is the task designed? When and where does it take place? Why is it done? Are there any interventions happening that might alter features and labels associated with the task? And how is the task measured? The lack of clear answers to these questions indicates that the model and its evaluation lack generalizability simply because it is not clear to what the model should generalize. Likewise, clear answers to these questions without a corresponding diverse representation in the benchmark dataset to evaluate performance leaves open the question of whether the dataset generalizes to the contexts in which the model is intended to generalize.

As an example of a generalization failure in a human-centered machine learning application, consider facial recognition. In Joy Buolamwini’s and Timnit Gebru’s algorithmic audit of facial recognition benchmarks and classifiers, the authors reveal the most commonly used benchmarks for evaluating facial recognition accuracy were composed of images of people with predominantly light skin. In other words, images of people with dark skin were relatively out-of-distribution (27). Furthermore, the Buolamwini et al 2018 audit presented a new benchmark to evaluate accuracy across intersectional identities. Commercial gender classification models performed extremely accurately in identifying men with light skin (with a maximum error rate of less than 1%) but incorrectly in women with dark skin (with a maximum error rate of 35%) (27). This large accuracy disparity reveals how failures to generalize can be hidden by benchmarks that do not represent the diversity of the real world. Research on machine learning applied to the diagnosis of skin disease reveals a similar story to facial recognition: models trained to classify skin disease based on images of only light or dark skin are more accurate in skin tones closest to the skin tones in the images in which the model was trained (28). These examples corroborate the notion that simply optimizing for predictive accuracy with very large datasets can often misrepresent the true data generating process and lead to systematic errors (29).

In other domains like affect recognition, an out-of-distribution context can be very task specific. For example, spontaneous facial expressions can be out-of-distribution for facial expression benchmarks that primarily contain posed expressions (30). Likewise, images labeled with emotions such as anger or surprise can be out-of-distribution for the same benchmarks where happy and neutral labels are most common (31).

Machine learning models that have been trained on perceptual data are subject to systematic failures on a special case of out-of-distribution data: adversarial perturbations. Adversarial perturbations refer to minor changes in data that do not influence classification of the data by humans but radically alter a model’s classification. As an example, researchers have demonstrated that adding a small sticker to a stop sign can alter the classification of machine learning models’ such that the models incorrectly classify the stop sign as a yield sign (32; 33). Researchers have shown that one can generalize adversarial perturbations by attaching a mainly translucent sticker on the lens of a camera (34). Likewise, researchers have demonstrated that adversarial perturbations can be applied to medical data e.g. noise or rotations in medical images and text substitution in medical notes and reimbursement codes (35). In general, adversarial perturbations demonstrate a lack of model robustness (36), lead to model errors that reasonable humans would rarely make, and open the question: How can we build models that are invariant to the same semantically meaningful features to which humans are invariant? Training robust models with adversarial perturbations is a starting point for aligning model performance more closely with human perceptions (37), but it is often difficult to identify the comprehensive possibility space of adversarial perturbations.

What drives the systematic errors by machine learning models on out-of-distribution data? The next section discusses two perspectives for characterizing the benchmark-production gap: the distribution shift perspective and the context shift perspective. The rest of the paper describes three methods for addressing context shift and considers three case studies of context shift in facial expression recognition, deepfake detection, and medical diagnosis.

2 Systematic Errors Arise from Context Shift and Lead to Distribution Shift

The mismatch between two datasets (e.g. the train and test splits or a test benchmark and production data) has been traditionally described as a dataset shift (38). More recently, machine learning researchers have described the same concept as distribution shift. In order to illustrate the growing attention to and evolving semantics of distribution shift, we present the number of papers on Google Scholar containing both “machine learning” and “distribution shift” (and other sub-components of distribution shift) in Figure 1.

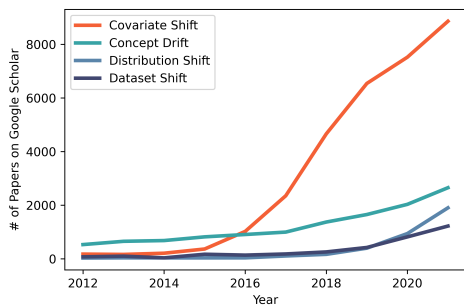


Figure 1: Number of papers on Google Scholar from 2012 to 2021 for search queries combining “machine learning” + the four most common terms for distribution shifts. For context, “machine learning” returns 185,000 articles in 2012 and 245,000 articles in 2021. The terms “prior probability shift” and “concept shift” return 398 and 1,040 papers over all time, respectively, when paired with “machine learning”.

Distribution shift refers to the non-equivalence of the joint distributions between two datasets. Formally, distribution shift describes the following equation $P_1(y, x) \neq P_2(y, x)$ where $P_n(y, x)$ is the joint distribution of labels, y , and covariates, x for a particular dataset, n (39). Based on Moreno et al 2012, the four subcategories of distribution shift include **covariate shift** when the distribution of

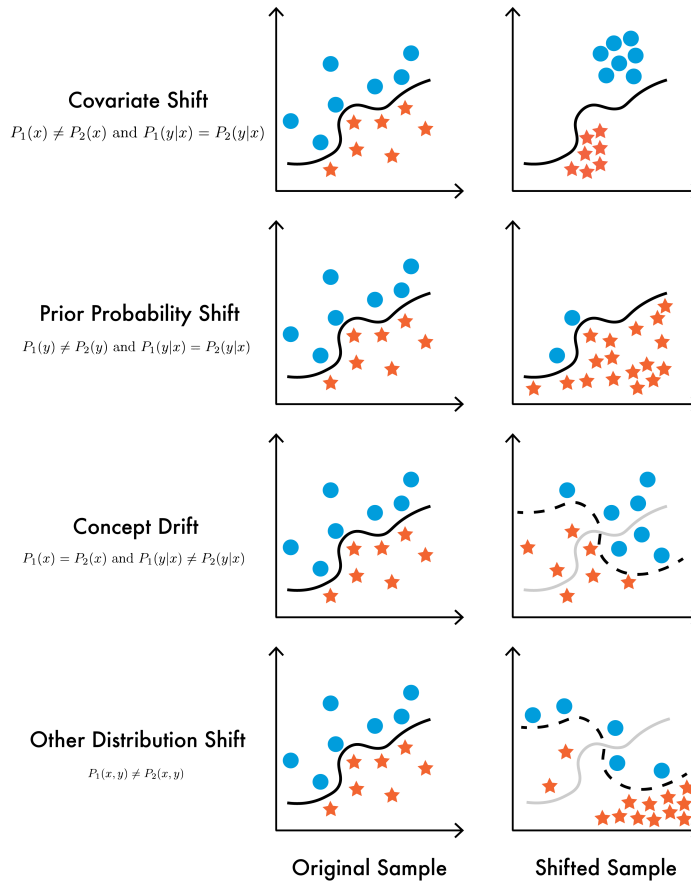


Figure 2: Illustrations of the four kinds of distribution shifts as defined in Moreno et al. 2012 (39). The spatial positions represents the feature space, geometric shapes and colors represents the ground truth label, the solid boundary line represents the learned representation of labels from the original sample, and the dotted boundary line represents the learned representation of labels from the shifted sample. Most real-world distribution shifts involve changes across features, labels, and the relationship between features and labels, and as such would be characterized as “Other Distribution Shift.” The core problem with the conceptual framework of distribution shift is that it is merely a symptom of changes in data-generating processes - how data are created, collected, and curated – but not part of the data-generating process itself. In order to improve model reliability and robustness, researchers need to take into consideration the data generating process.

features changes but everything else remains the same, **prior probability shift** when the distribution of labels changes but everything else remains the same, **concept shift** (more commonly referred to as concept drift) when the distribution of labels conditional on features changes but everything else remains the same, and **other distribution shift** when none of the other three shifts hold and the joint distributions between two datasets is different. We illustrate examples of each shift in Figure 2 to motivate intuition as to how the changes appear. Moreno et al 2012 formally specify the four subcategories of distribution shifts as follows (39):

- Covariate shift: $P_1(x) \neq P_2(x)$ but $P_1(y|x) = P_2(y|x)$
- Prior probability shift: $P_1(y) \neq P_2(y)$ but $P_1(y|x) = P_2(y|x)$
- Concept drift: $P_1(y|x) \neq P_2(y|x)$ but $P_1(x) = P_2(x)$
- Other distribution shift: $P_1(y, x) \neq P_2(y, x)$ where none of the above three shifts applies.

In theory (and within synthetic data), these four subcategories of distribution shift can be disentangled. However, production data, especially in human-centered applications, is subject to changing distributions and is often best characterized by the catch-all “Other distribution shift” sub-category. As such, the fundamental problem with trying to directly address the benchmark-production gap by focusing on distribution shift (or robustness under covariate shift) is the solution focuses only on the symptoms of the changes but not the underlying changes themselves. Distribution shift is downstream of the data generating process, and machine learning researchers has long considered the hidden contexts behind the distribution shift (40). In this paper, we seek to re-direct attention from the perspective of “distribution shift” towards “context shift” which refers to changes in the semantically meaningful features that influence data-generating processes. The solution to addressing context shift involves focusing on how to identify the changes in the creation, collection, and curation of data that lead to distribution shifts. By centering the problem of robustness and generalizability of applied machine learning on context shift, we seek to illustrate the importance of data-centered machine learning (alongside model-centered machine learning) in generating research that produces robust and generalizable models.

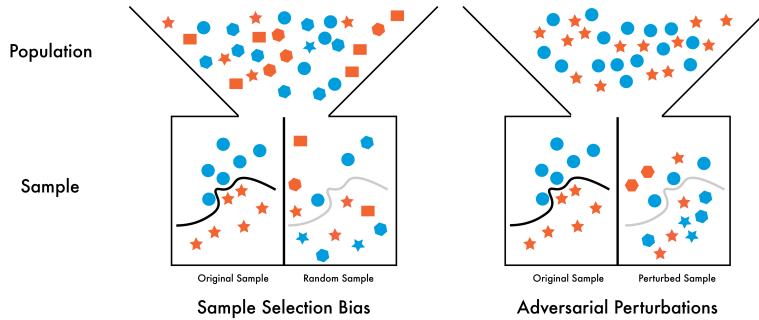


Figure 3: Illustrations of sample selection bias and adversarial perturbations with colors representing the ground truth label, geometric shapes and spatial positions representing the features, the top of the funnel representing the full populations, the bottom of the funnel representing the samples drawn from the population, and the solid boundary line representing the learned representation of labels from the original sample. On the left, the population contains upright stars, rotated stars, hexagons, rectangles, and circles, but the biased original sample only contains circles and stars. The random sample contains much higher diversity of features and relationships between features and labels. As such, the learned representation fails in more than 50% of observations. On the right, the population contains upright stars and blue circles. The original sample contains the same set of features, but the perturbed sample includes both rotated hexagons and stars, which may not be immediately noticeable to humans at first glance. Depending on the rotation, the learned representation misclassifies the perturbed shapes. Both pairs of samples present changes in features and changes in labels conditional on the features, which would make these examples of “Other Distribution Shift.” This figure is intended to provide intuition for where the perspective of distribution shift is inadequate and where the perspective of identifying semantically meaningful features that influence how samples are curated and created can inform approaches for addressing robustness in applications of machine learning.

Rather than focusing on differences in two distributions with disregard for the reasons behind the difference, we suggest researchers consider three concepts that drive context shift. These semantically meaningful shifts include sample selection bias (e.g. the new dataset contains images of people from a demographic not represented in the old dataset), adversarial perturbations (e.g. the new dataset contains noise injections that are imperceptible to human perception but change model performance), or non-stationarity (e.g. the new dataset contains images of smart phones post 2018 but the old dataset only contains flip phones before 2010). While we list non-stationarity separately from sample selection bias, non-stationarity can be considered as a special case of sample selection bias where sample selection bias arises from the inability to sample from features and labels in the future. We present Figure 3 to illustrate sample selection bias and adversarial perturbations, which can be formally described as follows:

- Sample selection bias: $P_1(s) \not\subset P_2(s)$ where s indicates x , y , or $y|x$
- Adversarial perturbations: $P_1(x) \neq P_2(x)$ but $P_1(y|H(x)) = P_2(y|H(x))$ where $H(x)$ represents human perception of the data

Unlike distribution shift, which can be measured between two datasets, context shift can only be fully addressed by learning the entire population’s data distribution, the kinds of changes that are and are not perceptible to humans, and how the population’s data distribution changes over time and space. Outside of artificially constrained spaces like synthetic datasets or games, access to the entire population’s data distribution (or the rules governing the distribution) across space and time is rare. Nevertheless, people generally have intuition and the ability to reason about data distributions of combinatoric contexts that they might never experience. In fact, cognitive science research shows that intuitive reasoning about statistical power analysis begins early in childhood (41).

By addressing the benchmark-production gap problem from the perspective of context shift as opposed to distribution shift, we can consider three approaches for increasing generalizability: human intuition and subject matter expertise in machine learning model development, dynamic benchmarking in the evaluation of machine learning models, and limitations statements that clarify how a machine learning model will generalize.

3 Addressing Robustness with Human Intuition and Expertise

Over the last few years, researchers have been developing data-centered frameworks to offer guidance for breaking down the data generating process into relevant component parts that reveal where context shift may lead to benchmark-production performance gaps. These frameworks include *Data Statements for Natural Language Processing* (42), *The Dataset Nutrition Label* (43), *Model Cards for Model Reporting* (44), *Datasheets for Datasets* (45), *Closing the AI accountability gap* (46), *The Ethical Pipeline for Healthcare Model Development* (47), and *The Clinician and Dataset Shift in Artificial Intelligence* (48). Likewise, meta-frameworks offer guidance for ensuring data documentation frameworks are useful and actionable (49).

As a heuristic for human-centered machine learning applications, teams of conscientious, creative, and skilled model developers, data engineers, and subject matter experts may find it useful to identify a first-order, non-exhaustive list of dimensions on which context shift is likely to occur. This list of dimensions depends largely on the context and the degree to which the data are subjective, representative, and missing (50). In ethnographic interviews with machine learning engineers, researchers find that engineers often address changes in context with “elaborate rule-based guardrails to avoid incorrect outputs” (51). Recent examples of semantically meaningful dimensions that have been demonstrated as useful for evaluating robustness in applied machine learning include skin color in face recognition (27) and dermatology diagnosis (28; 7), background scenery for affect recognition (5), number of people in a video for deepfake detection (52), number of chronic illnesses for algorithmic healthcare risk prediction (53), data artifacts like surgical markings (54) or clinically irrelevant labels (55) for medical diagnosis classification, patients’ self reports of pain for quantifying severity of knee osteoarthritis (56), and image similarity characteristics for pathologists to disambiguate between machine learning and user errors (57).

Knowledge elicitation is not a solved problem, but helpful questions that may guide the identification of potential context shifts in complex, human-centered machine learning applications include (and are not limited to): who are represented in the data and as annotators of the data, when and where is the

data collected, how do social, geographical, temporal, technological, aesthetic, financial incentives and other idiosyncrasies influence the creation of the data, and why the data is curated as it is. Knowledge elicitation has been historically ill-defined in artificial intelligence applications (58), but recent work developing taxonomies for knowledge elicitation helps to formalize the process and increase transparency along the way (59; 60).

Another expert intuition guided approach to closing the benchmark-production gap involves developing test benchmarks with adequate diversity in the data along the contextual dimensions upon which human intuition and expertise suggests model performance is most likely to vary. Recent examples of benchmark datasets working towards this goal are *BREEDS: Benchmarks for Subpopulation Shift* (61) and *WILDS: A Benchmark of in-the-Wild Distribution Shifts* (62), which includes labels for relative contexts and sub-populations for the explicit examination of context shifts.

4 Addressing Robustness with Dynamic Benchmarking

A second approach to addressing the benchmark-production gap is to transform the practice of evaluation from static benchmarks to dynamic benchmarks where models' performance is not evaluated on a single dataset, but rather continually evaluated on datasets produced via well-specified, quality controlled data generation processes. One example of dynamic benchmarking is dynabench (63), which is designed for natural language processing tasks. For general development of dynamic benchmarks, data generation process desiderata should include specifying the following dimensions of a dynamic benchmark:

- **Prediction task:** What are the input features and output labels? For example, inputs may be images and outputs may be lists of objects or inputs may be described more specifically as images of skin lesions photographed by dermatoscopes and outputs may be classifications of benign and malignant by board-certified dermatologists in the United States. It is important to be careful that the task matches the expected goal because unexpected mismatches between tasks and goals are relatively common (64; 65).
- **Ground truth annotation arbitration:** Who has the authority to annotate the data? How do experts differ from crowdworkers or an algorithm (66)? How should the data be annotated? How should inter-annotator disagreement be represented? What categories should be included?
- **Data inclusion and exclusion criteria:** What are the possible data sources? How are data curated from these sources? What is the data distribution of categories and subcategories? What are the quality constraints?
- **Benchmark size and shape:** What is the minimum size of a batch of data to serve as a benchmark? How should benchmarks by different groups for the same task be combined together?

These desiderata enable the development of dynamic benchmarks that further enable quantitative evaluation of model robustness via corroborated accuracy, which is the distribution of accuracy scores across dynamic benchmarks. Rather than simply evaluating a model on a single or a few static test benchmarks, we might consider a well-corroborated model to be one that meets two criteria: first, it is reasonably available for evaluation, and second, all attempts to uncover systematic errors in well-specified contexts reveal no significant accuracy disparities. The practice of dynamic benchmarking could be particularly relevant for addressing the *AI Knowledge Gap* (67) characterized by the disparity between the large number of machine learning models and the small number of studies evaluating these models' performance. Furthermore, dynamic benchmarking can be combined with benchmark task misalignment methodologies (68; 69) to assess how aligned (or misaligned) model predictions are with human annotations and considering diverse examples that bring transparency to the ethical implications and societal impact to model development (70).

The transition from static benchmarks on a particular instance (or set of instances) to dynamic benchmarks on data generation processes defined by explicit desiderata may be useful for addressing the fundamental issue of construct validity that arises in singular, static benchmarks (46).

5 Addressing Robustness by Clarifying a Model’s Limitations

A third approach to reducing the benchmark-production gap is to appropriately specify the contexts in which a model is expected to work via a limitations section (71).

To clarify domain-specific limitations driving the benchmark-production gap, we consider implicit assumptions that lead to a context shift in three real-world computer vision tasks: facial expression recognition, deepfake detection, and medical diagnosis.

6 Case Studies for Addressing Context Shift in Applied Machine Learning

6.1 Facial Expression Recognition

In the field of affective computing, facial expression recognition (FER) is a task to classify human facial expressions with affective labels (72; 31), which can be a useful component in designing human-AI interactions with computational empathy (73; 74; 75). Model-based FER is similar to how humans recognize the emotions of others (called empathic accuracy in affective science (76) and emotion reasoning in developmental psychology (77)) except that FER is based solely on facial expressions, whereas affect recognition can include information about someone’s gestures, language, tone, physiological measurements, and the long-tail of context, which can include factors such as the temperature outside, the social relationship between two individuals, what happened the day before, and more.

Consider an example from relatively recent research (78) where a standard neural network architecture, AlexNet (79), is trained on a large number of images of spontaneous and posed facial expressions to classify images into seven categories (anger, disgust, fear, happiness, sadness, surprise, and neutral) and achieves accuracy scores ranging from 48.6% in SFEW (80) to 56.0% in MMI (81) to 56.1% in DISFA (82) to 61.1% in FER2013 (83) to 77.4% in FERA (84; 85) to 92.2% in CK + (86) to 94.8% in MultiPie (87). While this model’s accuracy is significantly better than random guessing, which would be 14.2%, it varies dramatically depending on the chosen benchmark dataset. How should we interpret a performance gain of 21.9 percentage points on one dataset and an average performance gain of 3.5 percentage points on the other 6 datasets in an alternative network architecture? How should we interpret the model’s ability to achieve higher accuracy scores than non-neural network methods on three of the seven benchmark datasets? What does the distribution of performance tell us about how this model would perform on real-world production data? There is no clear answer to any of these questions, yet an implicit assumption in the well-cited, peer-reviewed publication of this FER paper is the slightly improved performance on several benchmark datasets appears to mark a contribution to the field of facial expression recognition. This assumption has the potential to lead to another more pernicious and mistaken assumption: the role of contextual features for real-world performance can be ignored when assessing the state-of-the-art methodology in applied problems like FER.

Clearly, models can learn facial expression features that map to human annotations of a handful of emotion categories to classify images at significantly better than chance rates. But, it is not reported nor clear how changes in lighting, head pose, occlusion, skin tone, ethnicity, age, gender, and background scenery influence both the model’s performance or human annotations. It is also underexplored how well FER models would perform if humans of diverse cultures annotated these images. Likewise, it is unclear how the model would perform on more fine-grained emotion categories (88) or labels based on affective dimensions like valence, arousal, and dominance. Furthermore, in many real-world settings where people may feign smiles to appease their managers, cry to express joy, or appear neutral to hide a winning poker hand, the perspective of outside observers may be very different than the perspective of close friends or individuals themselves. We highlight these relevant contextual features to highlight the many dimensions in which context shift can occur between test benchmarks and real-world production data. While these are not an exhaustive list of contextual features, these represent intuitive, first-order contexts for conducting algorithmic audits, developing future benchmark datasets with these labeled contexts, and adapting models to handle these dimensions. While researchers build the next version of contextualized dynamic benchmarks, other researchers who are focused on developing models should at the very least include caveats in their papers about the likely contextual dimensions that may affect performance.

6.2 Deepfake Detection

As a second case study of context shift in real-world applications of computer vision, we consider deepfake detection. Deepfakes are videos that have been manipulated to make someone appear to do or say something they have not said (89). These types of manipulation can be qualitatively characterized as face swapping where two people’s faces are swapped, head puppetry where facial landmarks are adjusted to make someone appear to be speaking, and lip-syncing where an individual’s lips are moved in sync with the phonemes from an external audio track (90).

The largest deepfake detection benchmark dataset to date is the Deepfake Detection Competition Dataset (DFDC) (91; 6), which consists of 128,154 videos based on performances by 960 consenting actors representing diversity across sex and ethnicity. However, Groh et al 2022 point out, “Unlike viral deepfake videos of politicians and other famous people, the videos from [this benchmark dataset] have minimal context: These are all 10 [second] videos depicting unknown actors making uncontroversial statements in nondescript locations” (52). This deepfake test benchmark is designed to evaluate algorithmic performance in identifying videos that have (and have not) been manipulated by seven synthetic techniques.

But, the real-world deepfake detection problem is not simply identifying whether one of seven synthetic techniques has been applied to a video. Instead, the real-world problem is identifying videos that have been algorithmically altered to impersonate innocent people and deceive the viewer. This problem is more than just a computer vision problem; it is a deception detection problem that involves both searching for artifacts that reveal that a manipulation has occurred and applying prior knowledge and critical reasoning to assess the likelihood that the video has been fabricated.

The DFDC does not include politicians or any scenes of news conferences or people speaking to a large audience. If we assume that harmful deepfakes will involve these kinds of contexts (like a deepfake of President Volodymyr Zelensky that appeared in March 2022 (92)), then it is important to evaluate models on videos with these kinds of dimensions, such as those from the Presidential Deepfakes Dataset (93; 94) and the Protecting World Leaders against Deepfakes Dataset (95). When Groh et al 2022 examined the leading state-of-the-art for detecting DFDC videos on deepfakes of Kim Jung-un and Vladimir Putin, they found the the leading model predicted a 2% and 8% likelihood these videos are deepfakes. While failure on two examples is only an anecdote, this failure speaks to an important need: diverse test benchmarks that cover the first-order dimensions where human intuition and expertise suggests context shift is most likely to occur.

6.3 Medical Diagnosis

As a third case study of context shift, we consider medical diagnosis in store-and-forward teledermatology settings where clinical data are collected at one site and sent electronically for evaluation at another site. Recent research on machine learning applied to skin disease classification has demonstrated the human expert-level performance of models in a number of specific tasks (96; 97). However, it is unclear how these models will perform on people with dark skin because the first paper does not describe the distribution of ethnicity or skin tone in the evaluation benchmark (96) and the evaluation benchmark in the second paper contains only 2.7% of people with the second darkest of the six Fitzpatrick Skin Types (FSTs) and 1 person with the darkest of the FSTs (97). Given the accuracy disparities that appeared across skin types in facial recognition, expert intuition suggests that systematic errors are likely to also appear in skin disease classifiers.

In fact, empirical research corroborates this intuition (28), and the Diverse Dermatology Images (DDI) dataset (7) reveals that state-of-the-art skin disease classification models make systematically more errors on dark skin than on light skin. The DDI represents a more comprehensive benchmark than previous datasets, and as a result, the DDI exposed errors that should guide and motivate the future development of machine learning models towards more robustness. However, the DDI is not perfectly comprehensive; the dataset is de-identified for privacy reasons and lacks free text clinical notes and other information that physicians would acquire via an in-person examination (7). Given that many skin diseases appear similarly and expert diagnoses are based on clinical history and non-visual features, expert intuition would expect, once again, that systematic errors lurk in the state-of-the-art machine learning models for store-and-forward skin disease classification.

7 Towards Robustness in Applied Machine Learning

Supervised machine learning models are very good at identifying statistical regularities in a given dataset but tend to err on out-of-distribution data that may arise from sample selection bias, adversarial perturbation, or nonstationarity. On the other hand, humans can be quite good at identifying contextual examples of out-of-distribution data. By combining the strengths of machine learning models with human intuition and expertise, early career ancient historians can quickly restore and date ancient texts (98), content moderation teams can more accurately distinguish between real and fake videos (52), and general practitioners can more accurately diagnose skin conditions from images (99) (although AI advice can also mislead experts; see (100; 101; 102; 103; 104)). In fact, initial evidence suggests that human intuition is fairly accurate in predicting model misclassifications on common object detection tasks (105). The integration of machine predictions with human decisions in collaborative decision making systems may be the most immediately effective way to avoid errors from context shift. The three case studies suggest the following advice for applied machine learning researchers:

- **Human intuition and subject matter expertise** can be useful for identifying first-order dimensions where context shift is likely to occur. These dimensions can inform the write-up of a limitations section, the development of a test benchmark, the collection of new data, or changes to model architecture.
- The practice of **dynamic benchmarking** mirrors the real-world more closely than static benchmarking and can enable insights from anywhere into systematic model failures.
- The inclusion of **limitations statements** in peer-reviewed research can increase model generalizability by simply clarifying the contexts in which a model is expected to generalize or not.

Promising future research directions for developing robust machine learning models under distribution shift involve the following iterative process: first, identify missing contexts in test benchmarks, second, collect data that contain those missing contexts, and third, adjust the model accordingly. Researchers can begin to identify missing contexts by collaborating with human experts who may be able to identify first-order drivers of context shift on a task-by-task basis. Similarly, researchers can further identify missing contexts by evaluating models against data generation process desiderata rather than a single or a few datasets.

Finally, one of the most effective solutions for addressing the benchmark-production gap is for researchers to clearly communicate the contexts in which a model has been evaluated and the contexts in which the model's performance is unknown.

References

- [1] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [2] L. Deng, "The mnist database of handwritten digit images for machine learning research [best of the web]," *IEEE signal processing magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- [3] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollár, and C. L. Zitnick, "Microsoft coco captions: Data collection and evaluation server," *arXiv preprint arXiv:1504.00325*, 2015.
- [4] A. Wang, A. Singh, J. Michael, F. Hill, O. Levy, and S. R. Bowman, "Glue: A multi-task benchmark and analysis platform for natural language understanding," *arXiv preprint arXiv:1804.07461*, 2018.
- [5] R. Kosti, J. M. Alvarez, A. Recasens, and A. Lapedriza, "Context Based Emotion Recognition using EMOTIC Dataset," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2019, arXiv: 2003.13401. [Online]. Available: <http://arxiv.org/abs/2003.13401>
- [6] B. Dolhansky, J. Bitton, B. Pflaum, J. Lu, R. Howes, M. Wang, and C. C. Ferrer, "The DeepFake Detection Challenge (DFDC) Dataset," *arXiv:2006.07397 [cs]*, Oct. 2020, arXiv: 2006.07397. [Online]. Available: <http://arxiv.org/abs/2006.07397>

- [7] R. Daneshjou, K. Vodrahalli, R. A. Novoa, M. Jenkins, W. Liang, V. Rotemberg, J. Ko, S. M. Swetter, E. E. Bailey, O. Gevaert *et al.*, “Disparities in dermatology ai performance on a diverse, curated clinical image set,” *Science advances*, vol. 8, no. 31, p. eabq6147, 2022.
- [8] J. Irvin, P. Rajpurkar, M. Ko, Y. Yu, S. Ciurea-Ilicus, C. Chute, H. Marklund, B. Haghgoo, R. Ball, K. Shpanskaya *et al.*, “Chexpert: A large chest radiograph dataset with uncertainty labels and expert comparison,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 590–597.
- [9] A. E. Johnson, T. J. Pollard, L. Shen, L.-w. H. Lehman, M. Feng, M. Ghassemi, B. Moody, P. Szolovits, L. Anthony Celi, and R. G. Mark, “Mimic-iii, a freely accessible critical care database,” *Scientific data*, vol. 3, no. 1, pp. 1–9, 2016.
- [10] R. Srinivasan and A. Chander, “Biases in ai systems,” *Communications of the ACM*, vol. 64, no. 8, pp. 44–49, 2021.
- [11] H. Suresh and J. Gutttag, “A framework for understanding sources of harm throughout the machine learning life cycle,” in *Equity and access in algorithms, mechanisms, and optimization*, 2021, pp. 1–9.
- [12] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, “Imagenet large scale visual recognition challenge,” *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [13] X. Chen, C. Liu, B. Li, K. Lu, and D. Song, “Targeted Backdoor Attacks on Deep Learning Systems Using Data Poisoning,” *arXiv:1712.05526 [cs]*, Dec. 2017, arXiv: 1712.05526. [Online]. Available: <http://arxiv.org/abs/1712.05526>
- [14] S. Kapoor and A. Narayanan, “Leakage and the reproducibility crisis in ml-based science,” *arXiv preprint arXiv:2207.07048*, 2022.
- [15] M. L. Gordon, M. S. Lam, J. S. Park, K. Patel, J. T. Hancock, T. Hashimoto, and M. S. Bernstein, “Jury Learning: Integrating Dissenting Voices into Machine Learning Models,” *arXiv:2202.02950 [cs]*, Feb. 2022, arXiv: 2202.02950. [Online]. Available: <http://arxiv.org/abs/2202.02950>
- [16] C. G. Northcutt, A. Athalye, and J. Mueller, “Pervasive Label Errors in Test Sets Destabilize Machine Learning Benchmarks,” *arXiv:2103.14749 [cs, stat]*, Apr. 2021, arXiv: 2103.14749. [Online]. Available: <http://arxiv.org/abs/2103.14749>
- [17] A. Torralba and A. A. Efros, “Unbiased look at dataset bias,” in *CVPR 2011*. IEEE, 2011, pp. 1521–1528.
- [18] R. L. Thomas and D. Uminsky, “Reliance on metrics is a fundamental challenge for ai,” *Patterns*, vol. 3, no. 5, p. 100476, 2022.
- [19] Z. C. Lipton, “The mythos of model interpretability,” *Communications of the ACM*, vol. 61, no. 10, pp. 36–43, Sep. 2018. [Online]. Available: <https://dl.acm.org/doi/10.1145/3233231>
- [20] M. Arjovsky, L. Bottou, I. Gulrajani, and D. Lopez-Paz, “Invariant Risk Minimization,” Mar. 2020, number: arXiv:1907.02893 arXiv:1907.02893 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1907.02893>
- [21] M. Makar, B. Packer, D. Moldovan, D. Blalock, Y. Halpern, and A. D’Amour, “Causally motivated shortcut removal using auxiliary labels,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022, pp. 739–766.
- [22] J. Adebayo, M. Muelly, H. Abelson, and B. Kim, “Post hoc explanations may be ineffective for detecting unknown spurious correlation,” in *International Conference on Learning Representations*, 2021.
- [23] A. Ghandeharioun, B. Kim, C.-L. Li, B. Jou, B. Eoff, and R. W. Picard, “Dissect: Disentangled simultaneous explanations via concept traversals,” *arXiv preprint arXiv:2105.15164*, 2021.

- [24] S. Jain, H. Lawrence, A. Moitra, and A. Madry, “Distilling model failures as directions in latent space,” *arXiv preprint arXiv:2206.14754*, 2022.
- [25] A. Krizhevsky, V. Nair, and G. Hinton, “The cifar-10 dataset,” *online: <http://www.cs.toronto.edu/kriz/cifar.html>*, vol. 55, no. 5, 2014.
- [26] B. Recht, R. Roelofs, L. Schmidt, and V. Shankar, “Do imagenet classifiers generalize to imagenet?” in *International Conference on Machine Learning*. PMLR, 2019, pp. 5389–5400.
- [27] J. Buolamwini and T. Gebru, “Gender shades: Intersectional accuracy disparities in commercial gender classification,” in *Conference on fairness, accountability and transparency*. PMLR, 2018, pp. 77–91.
- [28] M. Groh, C. Harris, L. Soenksen, F. Lau, R. Han, A. Kim, A. Koochek, and O. Badri, “Evaluating Deep Neural Networks Trained on Clinical Images in Dermatology with the Fitzpatrick 17k Dataset,” in *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. Nashville, TN, USA: IEEE, Jun. 2021, pp. 1820–1828. [Online]. Available: <https://ieeexplore.ieee.org/document/9522867/>
- [29] J. Hullman, S. Kapoor, P. Nanayakkara, A. Gelman, and A. Narayanan, “The worst of both worlds: A comparative analysis of errors in learning from data in psychology and machine learning,” *arXiv preprint arXiv:2203.06498*, 2022.
- [30] D. Dupré, E. G. Krumhuber, D. Küster, and G. J. McKeown, “A performance comparison of eight commercially available automatic classifiers for facial affect recognition,” *PLOS ONE*, vol. 15, no. 4, p. e0231968, Apr. 2020. [Online]. Available: <https://dx.plos.org/10.1371/journal.pone.0231968>
- [31] S. Li and W. Deng, “Deep Facial Expression Recognition: A Survey,” *IEEE Transactions on Affective Computing*, pp. 1–1, 2020. [Online]. Available: <https://ieeexplore.ieee.org/document/9039580/>
- [32] T. B. Brown, D. Mané, A. Roy, M. Abadi, and J. Gilmer, “Adversarial patch,” *arXiv preprint arXiv:1712.09665*, 2017.
- [33] K. Eykholt, I. Evtimov, E. Fernandes, B. Li, A. Rahmati, C. Xiao, A. Prakash, T. Kohno, and D. Song, “Robust physical-world attacks on deep learning visual classification,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1625–1634.
- [34] J. Li, F. Schmidt, and Z. Kolter, “Adversarial camera stickers: A physical camera-based attack on deep learning systems,” in *International Conference on Machine Learning*. PMLR, 2019, pp. 3896–3904.
- [35] S. G. Finlayson, J. D. Bowers, J. Ito, J. L. Zittrain, A. L. Beam, and I. S. Kohane, “Adversarial attacks on medical machine learning,” *Science*, vol. 363, no. 6433, pp. 1287–1289, 2019.
- [36] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, “Adversarial examples are not bugs, they are features,” *Advances in neural information processing systems*, vol. 32, 2019.
- [37] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry, “Robustness may be at odds with accuracy,” *arXiv preprint arXiv:1805.12152*, 2018.
- [38] J. Quiñero-Candela, M. Sugiyama, A. Schwaighofer, and N. D. Lawrence, *Dataset shift in machine learning*, 2008.
- [39] J. G. Moreno-Torres, T. Raeder, R. Alaiz-Rodríguez, N. V. Chawla, and F. Herrera, “A unifying view on dataset shift in classification,” *Pattern recognition*, vol. 45, no. 1, pp. 521–530, 2012.
- [40] G. Widmer and M. Kubat, “Learning in the presence of concept drift and hidden contexts,” *Machine learning*, vol. 23, no. 1, pp. 69–101, 1996.
- [41] M. C. Pelz, K. R. Allen, J. B. Tenenbaum, and L. E. Schulz, “Foundations of intuitive power analyses in children and adults,” *Nature Human Behaviour*, pp. 1–12, 2022.

- [42] E. M. Bender and B. Friedman, “Data statements for natural language processing: Toward mitigating system bias and enabling better science,” *Transactions of the Association for Computational Linguistics*, vol. 6, pp. 587–604, 2018.
- [43] S. Holland, A. Hosny, S. Newman, J. Joseph, and K. Chmielinski, “The dataset nutrition label: A framework to drive higher data quality standards,” *arXiv preprint arXiv:1805.03677*, 2018.
- [44] M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, and T. Gebru, “Model cards for model reporting,” in *Proceedings of the conference on fairness, accountability, and transparency*, 2019, pp. 220–229.
- [45] T. Gebru, J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. D. Iii, and K. Crawford, “Datasheets for datasets,” *Communications of the ACM*, vol. 64, no. 12, pp. 86–92, 2021.
- [46] I. D. Raji, E. M. Bender, A. Paullada, E. Denton, and A. Hanna, “AI and the Everything in the Whole Wide World Benchmark,” *arXiv:2111.15366 [cs]*, Nov. 2021, arXiv: 2111.15366. [Online]. Available: <http://arxiv.org/abs/2111.15366>
- [47] I. Y. Chen, E. Pierson, S. Rose, S. Joshi, K. Ferryman, and M. Ghassemi, “Ethical Machine Learning in Health Care,” *arXiv:2009.10576 [cs]*, Oct. 2020, arXiv: 2009.10576. [Online]. Available: <http://arxiv.org/abs/2009.10576>
- [48] S. G. Finlayson, A. Subbaswamy, K. Singh, J. Bowers, A. Kupke, J. Zittrain, I. S. Kohane, and S. Saria, “The Clinician and Dataset Shift in Artificial Intelligence,” *New England Journal of Medicine*, vol. 385, no. 3, pp. 283–286, Jul. 2021. [Online]. Available: <http://www.nejm.org/doi/10.1056/NEJMc2104626>
- [49] A. Heger, E. B. Marquis, M. Vorvoreanu, H. Wallach, and J. W. Vaughan, “Understanding machine learning practitioners’ data documentation perceptions, needs, challenges, and desiderata,” *arXiv preprint arXiv:2206.02923*, 2022.
- [50] S. Mullainathan and Z. Obermeyer, “Does Machine Learning Automate Moral Hazard and Error?” vol. 107, no. 5, p. 5, 2017.
- [51] S. Shankar, R. Garcia, J. M. Hellerstein, and A. G. Parameswaran, “Operationalizing machine learning: An interview study,” 2022. [Online]. Available: <https://arxiv.org/abs/2209.09125>
- [52] M. Groh, Z. Epstein, C. Firestone, and R. Picard, “Deepfake detection by human crowds, machines, and machine-informed crowds,” *Proceedings of the National Academy of Sciences*, vol. 119, no. 1, p. e2110013119, Jan. 2022. [Online]. Available: <http://www.pnas.org/lookup/doi/10.1073/pnas.2110013119>
- [53] Z. Obermeyer, B. Powers, C. Vogeli, and S. Mullainathan, “Dissecting racial bias in an algorithm used to manage the health of populations,” *Science*, vol. 366, no. 6464, pp. 447–453, Oct. 2019. [Online]. Available: <https://www.sciencemag.org/lookup/doi/10.1126/science.aax2342>
- [54] J. K. Winkler, C. Fink, F. Toberer, A. Enk, T. Deinlein, R. Hofmann-Wellenhof, L. Thomas, A. Lallas, A. Blum, W. Stolz, and H. A. Haenssle, “Association Between Surgical Skin Markings in Dermoscopic Images and Diagnostic Performance of a Deep Learning Convolutional Neural Network for Melanoma Recognition,” *JAMA Dermatology*, vol. 155, no. 10, p. 1135, Oct. 2019. [Online]. Available: <https://jamanetwork.com/journals/jamadermatology/fullarticle/2740808>
- [55] L. Oakden-Rayner, J. Dunnmon, G. Carneiro, and C. Re, “Hidden stratification causes clinically meaningful failures in machine learning for medical imaging,” in *Proceedings of the ACM Conference on Health, Inference, and Learning*. Toronto Ontario Canada: ACM, Apr. 2020, pp. 151–159. [Online]. Available: <https://dl.acm.org/doi/10.1145/3368555.3384468>
- [56] E. Pierson, D. M. Cutler, J. Leskovec, S. Mullainathan, and Z. Obermeyer, “An algorithmic approach to reducing unexplained pain disparities in underserved populations,” *Nature Medicine*, vol. 27, no. 1, pp. 136–140, Jan. 2021. [Online]. Available: <http://www.nature.com/articles/s41591-020-01192-7>

- [57] C. J. Cai, M. C. Stumpe, M. Terry, E. Reif, N. Hegde, J. Hipp, B. Kim, D. Smilkov, M. Wattenberg, F. Viegas, and G. S. Corrado, “Human-Centered Tools for Coping with Imperfect Algorithms During Medical Decision-Making,” in *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems - CHI '19*. Glasgow, Scotland Uk: ACM Press, 2019, pp. 1–14. [Online]. Available: <http://dl.acm.org/citation.cfm?doid=3290605.3300234>
- [58] D. E. Forsythe, “Engineering knowledge: The construction of knowledge in artificial intelligence,” *Social studies of science*, vol. 23, no. 3, pp. 445–477, 1993.
- [59] D. Kerrigan, J. Hullman, and E. Bertini, “A survey of domain knowledge elicitation in applied machine learning,” *Multimodal Technologies and Interaction*, vol. 5, no. 12, p. 73, 2021.
- [60] V. Chen, U. Bhatt, H. Heidari, A. Weller, and A. Talwalkar, “Perspectives on incorporating expert feedback into model updates,” *arXiv preprint arXiv:2205.06905*, 2022.
- [61] S. Santurkar, D. Tsipras, and A. Madry, “BREEDS: Benchmarks for Subpopulation Shift,” Aug. 2020, number: arXiv:2008.04859 arXiv:2008.04859 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/2008.04859>
- [62] P. W. Koh, S. Sagawa, H. Marklund, S. M. Xie, M. Zhang, A. Balsubramani, W. Hu, M. Yasunaga, R. L. Phillips, I. Gao, T. Lee, E. David, I. Stavness, W. Guo, B. A. Earnshaw, I. S. Haque, S. Beery, J. Leskovec, A. Kundaje, E. Pierson, S. Levine, C. Finn, and P. Liang, “WILDS: A Benchmark of in-the-Wild Distribution Shifts,” *arXiv:2012.07421 [cs]*, Jul. 2021, arXiv: 2012.07421. [Online]. Available: <http://arxiv.org/abs/2012.07421>
- [63] D. Kiela, M. Bartolo, Y. Nie, D. Kaushik, A. Geiger, Z. Wu, B. Vidgen, G. Prasad, A. Singh, P. Ringshia *et al.*, “Dynabench: Rethinking benchmarking in nlp,” *arXiv preprint arXiv:2104.14337*, 2021.
- [64] S. Mullainathan and Z. Obermeyer, “On the Inequity of Predicting A While Hoping for B,” *AEA Papers and Proceedings*, vol. 111, pp. 37–42, May 2021. [Online]. Available: <https://pubs.aeaweb.org/doi/10.1257/pandp.20211078>
- [65] S. Kerr, “On the folly of rewarding a, while hoping for b,” *Academy of Management journal*, vol. 18, no. 4, pp. 769–783, 1975.
- [66] M. Groh, C. Harris, R. Daneshjou, O. Badri, and A. Koochek, “Towards transparency in dermatology image datasets with skin tone annotations by experts, crowds, and an algorithm,” *arXiv preprint arXiv:2207.02942*, 2022.
- [67] Z. Epstein, B. H. Payne, J. H. Shen, A. Dubey, B. Felbo, M. Groh, N. Obradovich, M. Cebrian, and I. Rahwan, “Closing the AI Knowledge Gap,” *arXiv:1803.07233 [cs]*, Mar. 2018, arXiv: 1803.07233. [Online]. Available: <http://arxiv.org/abs/1803.07233>
- [68] D. Tsipras, S. Santurkar, L. Engstrom, A. Ilyas, and A. Madry, “From imagenet to image classification: Contextualizing progress on benchmarks,” in *International Conference on Machine Learning*. PMLR, 2020, pp. 9625–9635.
- [69] A. Ilyas, S. M. Park, L. Engstrom, G. Leclerc, and A. Madry, “Datamodels: Understanding predictions with data and data with predictions,” in *International Conference on Machine Learning*. PMLR, 2022, pp. 9525–9587.
- [70] A. Paullada, I. D. Raji, E. M. Bender, E. Denton, and A. Hanna, “Data and its (dis) contents: A survey of dataset development and use in machine learning research,” *Patterns*, vol. 2, no. 11, p. 100336, 2021.
- [71] J. J. Smith, S. Amershi, S. Barocas, H. Wallach, and J. Wortman Vaughan, “Real ml: Recognizing, exploring, and articulating limitations of machine learning research,” in *2022 ACM Conference on Fairness, Accountability, and Transparency*, 2022, pp. 587–597.
- [72] J. F. Cohn and F. De la Torre, “Automated face analysis for affective computing.” 2015.
- [73] R. W. Picard, *Affective computing*, 2000.

- [74] A. Paiva, I. Leite, H. Boukricha, and I. Wachsmuth, “Empathy in virtual agents and robots: A survey,” *ACM Transactions on Interactive Intelligent Systems (TiiS)*, vol. 7, no. 3, pp. 1–40, 2017.
- [75] M. Groh, C. Ferguson, R. Lewis, and R. Picard, “Computational empathy counteracts the negative effects of anger on creative problem solving,” *arXiv preprint arXiv:2208.07178*, 2022.
- [76] W. Ickes, “Empathic accuracy,” *Journal of personality*, vol. 61, no. 4, pp. 587–610, 1993.
- [77] A. L. Ruba and S. D. Pollak, “The Development of Emotion Reasoning in Infancy and Early Childhood,” *Annual Review of Developmental Psychology*, vol. 2, no. 1, pp. 503–531, Dec. 2020. [Online]. Available: <https://www.annualreviews.org/doi/10.1146/annurev-devpsych-060320-102556>
- [78] A. Mollahosseini, D. Chan, and M. H. Mahoor, “Going deeper in facial expression recognition using deep neural networks,” in *2016 IEEE Winter conference on applications of computer vision (WACV)*. IEEE, 2016, pp. 1–10.
- [79] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Advances in neural information processing systems*, vol. 25, 2012.
- [80] A. Dhall, R. Goecke, S. Lucey, and T. Gedeon, “Static facial expression analysis in tough conditions: Data, evaluation protocol and benchmark,” in *2011 IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*. IEEE, 2011, pp. 2106–2112.
- [81] M. Pantic, M. Valstar, R. Rademaker, and L. Maat, “Web-based database for facial expression analysis,” in *2005 IEEE international conference on multimedia and Expo*. IEEE, 2005, pp. 5–pp.
- [82] S. M. Mavadati, M. H. Mahoor, K. Bartlett, P. Trinh, and J. F. Cohn, “Disfa: A spontaneous facial action intensity database,” *IEEE Transactions on Affective Computing*, vol. 4, no. 2, pp. 151–160, 2013.
- [83] I. J. Goodfellow, D. Erhan, P. L. Carrier, A. Courville, M. Mirza, B. Hamner, W. Cukierski, Y. Tang, D. Thaler, D.-H. Lee *et al.*, “Challenges in representation learning: A report on three machine learning contests,” in *International conference on neural information processing*. Springer, 2013, pp. 117–124.
- [84] M. F. Valstar, B. Jiang, M. Mehu, M. Pantic, and K. Scherer, “The first facial expression recognition and analysis challenge,” in *2011 IEEE International Conference on Automatic Face & Gesture Recognition (FG)*. IEEE, 2011, pp. 921–926.
- [85] T. Bänziger and K. R. Scherer, “Introducing the geneva multimodal emotion portrayal (gemep) corpus,” *Blueprint for affective computing: A sourcebook*, vol. 2010, pp. 271–94, 2010.
- [86] P. Lucey, J. F. Cohn, T. Kanade, J. Saragih, Z. Ambadar, and I. Matthews, “The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression,” in *2010 IEEE computer society conference on computer vision and pattern recognition-workshops*. IEEE, 2010, pp. 94–101.
- [87] R. Gross, I. Matthews, J. Cohn, T. Kanade, and S. Baker, “Multi-pie,” *Image and vision computing*, vol. 28, no. 5, pp. 807–813, 2010.
- [88] A. S. Cowen, D. Keltner, F. Schroff, B. Jou, H. Adam, and G. Prasad, “Sixteen facial expressions occur in similar contexts worldwide,” *Nature*, vol. 589, no. 7841, pp. 251–257, Jan. 2021. [Online]. Available: <http://www.nature.com/articles/s41586-020-3037-7>
- [89] D. Boneh, A. J. Grotto, P. McDaniel, and N. Papernot, “Preparing for the age of deepfakes and disinformation.”
- [90] S. Lyu, “DeepFake Detection: Current Challenges and Next Steps,” *arXiv:2003.09234 [cs]*, Mar. 2020, arXiv: 2003.09234. [Online]. Available: <http://arxiv.org/abs/2003.09234>

- [91] B. Dolhansky, R. Howes, B. Pflaum, N. Baram, and C. C. Ferrer, “The Deepfake Detection Challenge (DFDC) Preview Dataset,” *arXiv:1910.08854 [cs]*, Oct. 2019, arXiv: 1910.08854. [Online]. Available: <http://arxiv.org/abs/1910.08854>
- [92] J. Wakefield, “Deepfake presidents used in russia-ukraine war,” Mar 2022. [Online]. Available: <https://www.bbc.com/news/technology-60780142>
- [93] A. Sankaranarayanan, M. Groh, R. Picard, and A. Lippman, “The presidential deepfakes dataset,” 2021.
- [94] M. Groh, A. Sankaranarayanan, and R. Picard, “Human detection of political deepfakes across transcripts, audio, and video,” *arXiv preprint arXiv:2202.12883*, 2022.
- [95] S. Agarwal, H. Farid, Y. Gu, M. He, K. Nagano, and H. Li, “Protecting world leaders against deep fakes.” 2019.
- [96] A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, and S. Thrun, “Dermatologist-level classification of skin cancer with deep neural networks,” *Nature*, vol. 542, no. 7639, pp. 115–118, Feb. 2017. [Online]. Available: <http://www.nature.com/articles/nature21056>
- [97] Y. Liu, A. Jain, C. Eng, D. H. Way, K. Lee, P. Bui, K. Kanada, G. de Oliveira Marinho, J. Gallegos, S. Gabriele, V. Gupta, N. Singh, V. Natarajan, R. Hofmann-Wellenhof, G. S. Corrado, L. H. Peng, D. R. Webster, D. Ai, S. J. Huang, Y. Liu, R. C. Dunn, and D. Coz, “A deep learning system for differential diagnosis of skin diseases,” *Nature Medicine*, vol. 26, no. 6, pp. 900–908, Jun. 2020. [Online]. Available: <http://www.nature.com/articles/s41591-020-0842-3>
- [98] Y. Assael, T. Sommerschild, B. Shillingford, M. Bordbar, J. Pavlopoulos, M. Chatzipanagiotou, I. Androutsopoulos, J. Prag, and N. de Freitas, “Restoring and attributing ancient texts using deep neural networks,” *Nature*, vol. 603, no. 7900, pp. 280–283, Mar. 2022. [Online]. Available: <https://www.nature.com/articles/s41586-022-04448-z>
- [99] A. Jain, D. Way, V. Gupta, Y. Gao, G. de Oliveira Marinho, J. Hartford, R. Sayres, K. Kanada, C. Eng, K. Nagpal, K. B. DeSalvo, G. S. Corrado, L. Peng, D. R. Webster, R. C. Dunn, D. Coz, S. J. Huang, Y. Liu, P. Bui, and Y. Liu, “Development and Assessment of an Artificial Intelligence–Based Tool for Skin Condition Diagnosis by Primary Care Physicians and Nurse Practitioners in Tele dermatology Practices,” *JAMA Network Open*, vol. 4, no. 4, p. e217249, Apr. 2021. [Online]. Available: <https://jamanetwork.com/journals/jamanetworkopen/fullarticle/2779250>
- [100] P. Tschandl, C. Rinner, Z. Apalla, G. Argenziano, N. Codella, A. Halpern, M. Janda, A. Lallas, C. Longo, J. Malvehy, J. Paoli, S. Puig, C. Rosendahl, H. P. Soyer, I. Zalaudek, and H. Kittler, “Human–computer collaboration for skin cancer recognition,” *Nature Medicine*, vol. 26, no. 8, pp. 1229–1234, Aug. 2020. [Online]. Available: <http://www.nature.com/articles/s41591-020-0942-0>
- [101] A. Abeliuk, D. M. Benjamin, F. Morstatter, and A. Galstyan, “Quantifying machine influence over human forecasters,” *Scientific Reports*, vol. 10, no. 1, p. 15940, Dec. 2020. [Online]. Available: <http://www.nature.com/articles/s41598-020-72690-4>
- [102] S. Gaube, H. Suresh, M. Raue, A. Merritt, S. J. Berkowitz, E. Lerner, J. F. Coughlin, J. V. Gutttag, E. Colak, and M. Ghassemi, “Do as AI say: susceptibility in deployment of clinical decision-aids,” *npj Digital Medicine*, vol. 4, no. 1, p. 31, Dec. 2021. [Online]. Available: <http://www.nature.com/articles/s41746-021-00385-9>
- [103] M. Jacobs, M. F. Pradier, T. H. McCoy, R. H. Perlis, F. Doshi-Velez, and K. Z. Gajos, “How machine-learning recommendations influence clinician treatment selections: the example of antidepressant selection,” *Translational Psychiatry*, vol. 11, no. 1, p. 108, Jun. 2021. [Online]. Available: <http://www.nature.com/articles/s41398-021-01224-x>
- [104] M. Vaccaro and J. Waldo, “The effects of mixing machine learning and human judgment,” *Communications of the ACM*, vol. 62, no. 11, pp. 104–110, 2019.

[105] Z. Zhou, M. Nartker, and C. Firestone, “When will ai misclassify? human intuition for machine (mis) perception,” *Journal of Vision*, vol. 20, no. 11, pp. 1325–1325, 2020.

Checklist

1. Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [Yes]
2. Did you describe the limitations of your work? [Yes]
3. Did you discuss any potential negative societal impacts of your work? [Yes]
4. Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]