# Google's AI Mislabeling Incident

## 1. Introduction

In the rapidly evolving field of artificial intelligence (AI) and machine learning, companies like Google have been at the forefront of innovation, developing technologies that reshape how we interact with the digital world. However, these advancements are not without challenges, particularly concerning Diversity, Equity, and Inclusion (DEI). A significant example is the 2015 incident where Google's Photos app mistakenly labeled photos of Black individuals with an offensive term. This error not only exposed critical flaws in AI development but also underscored the pressing need for inclusive practices and diverse perspectives in technology design.

## 2. Background of the Case

Google Photos, launched in May 2015, is a photo-sharing and storage service that uses AI-powered image recognition to categorize and label photos automatically. The app was designed to enhance user experience by organizing photos efficiently, grouping similar images, and applying labels for easy retrieval. It relies on machine learning algorithms trained on extensive datasets to recognize and categorize objects, scenes, and faces.

In June 2015, just weeks after its launch, a glaring flaw was brought to light. Jacky Alciné, a Brooklyn-based software developer, discovered that the app had mislabeled photos of him and his friend—both Black—with a highly offensive term. Shocked by this misclassification, Alciné took to Twitter to share his experience, stating, "Google Photos, y'all messed up. My friend's not a gorilla." His tweet quickly went viral, sparking widespread outrage and bringing significant attention to the racial insensitivity of the mislabeling.

Google promptly responded. Yonatan Zunger, Google's chief architect of social, replied to Alciné's tweet within an hour and a half, expressing deep regret and promising swift action. "We're appalled and genuinely sorry that this happened," a Google representative later stated. "We are taking immediate action

to prevent this type of result from appearing." Google removed the offensive label from the app's search feature to prevent further incidents while working on a more permanent solution.

## 3. Cause of the Failure

The root cause of this failure lies in the limitations of the machine learning algorithms and the data used to train them. Several factors contributed to the oversight:

### 3.1 Lack of Diverse Training Data

The AI was trained predominantly on images that did not adequately represent the diversity of human appearances. This lack of representation led to misclassification when the system encountered images of people with darker skin tones. Two former Google employees mentioned that the company's image collection used to train the AI system did not include enough photos of Black people, leading to the system's inability to recognize them accurately

### 3.2 Insufficient Contextual Understanding

AI systems lack the cultural and social awareness to understand the implications of their outputs. The algorithm did not comprehend that mislabeling a human with an animal label is not only incorrect but also deeply offensive and racist.

### 3.3 Homogeneous Development Teams

The development team may have lacked diversity, leading to blind spots regarding racial sensitivities. Diverse teams are more likely to anticipate and identify potential biases and culturally sensitive issues.

### 3.4 Emphasis on Rapid Deployment Over Thorough Testing

Google's approach of releasing cutting-edge technologies with the understanding that problems will arise—and will need fixing on the go—may have contributed to inadequate testing. While this strategy fosters innovation and gets products into users' hands quickly, it risks deploying products that haven't been

thoroughly vetted for biases. The company had not asked enough employees to test the feature before its public debut, which might have revealed the issue earlier.

## 4. Impact on Individuals and Groups

The mislabeling incident had profound negative impacts:

Being mislabeled in such a derogatory manner is dehumanizing and evokes historical racist tropes that compare Black people to animals. This caused significant emotional distress to those directly affected and to the broader Black community. The incident raised concerns about the reliability and fairness of AI technologies among marginalized communities. It highlighted how technology, when not developed with DEI considerations, can reinforce systemic biases. Such errors can contribute to the perpetuation of harmful stereotypes, further entrenching societal biases within technological systems. Jacky Alciné expressed his dismay, stating, "I'm going to forever have no faith in this AI."

It also generated broader social, economic, and reputational impact on Google.

Socially, the incident ignited critical conversations about the ethical responsibilities of tech companies in AI development. Advocacy groups, tech experts, and the public criticized Google for allowing such an egregious error to occur. The episode underscored the necessity for greater accountability, transparency, and inclusivity in technology.

Economically, while the immediate financial impact on Google was minimal, the long-term implications posed risks. Negative publicity can affect user trust and brand loyalty—crucial components for a company reliant on user engagement and data collection. Potential regulatory scrutiny over such incidents could also lead to increased compliance costs.

Reputationally, Google's image as an innovative and socially responsible leader in technology was tarnished. The incident exposed significant shortcomings in the company's internal processes and highlighted the need for more rigorous testing and inclusive practices.

In response to the incident, Google took immediate and significant actions: Google issued sincere apologies, both publicly and directly to those affected. The company's representatives expressed deep regret over the incident. Also, Google temporarily removed the offensive label and disabled the ability to search for certain animal categories within the app to prevent further misclassifications. Additionally, engineers began tweaking the algorithms to improve image recognition, particularly concerning skin tones and facial features of people of color.

However, years later, it was revealed that Google had not fully solved the problem and instead had disabled the ability to search for primates altogether. Recognizing the limitations of their training data, Google aimed to include more diverse images to enhance the AI's accuracy. However, efforts to improve facial recognition features led to further controversy when contractors reportedly targeted homeless people and students to collect facial scans.

The company planned to make the AI more cautious in its labeling, especially when the system was not entirely confident in its classification. Margaret Mitchell, a former Google researcher and co-founder of Google's Ethical AI group, supported the decision to remove certain labels, stating, "You have to think about how often someone needs to label a gorilla versus perpetuating harmful stereotypes."

Despite Google's efforts, the company struggled to fully resolve the issue. In subsequent years, Google and other companies like Apple chose to disable certain search functionalities rather than risk repeating offensive mistakes. This approach raises concerns about the efficacy of such fixes and whether they address the underlying biases in AI systems.

The decision to remove functionalities instead of solving the root problem suggests a reluctance to engage deeply with the complexities of AI bias. Vicente Ordóñez, a professor at Rice University who studies computer vision, questioned this approach: "How can we trust this software for other scenarios?" The inability

to fully correct the issue indicates that more fundamental changes are needed in how AI systems are developed and trained.

## Sources

[1] Barr, Alistair. "Google Mistakenly Tags Black People as "Gorillas," Showing Limits of Algorithms." *Wall Street Journal*, 1 July 2015, www.wsj.com/articles/BL-DGB-42522.

[2] BBC. "Google Apologises for Racist Blunder." *BBC News*, 1 July 2015, www.bbc.com/news/technology-33347866.

[3] Dougherty, Conor. "Google Photos Mistakenly Labels Black People "Gorillas."" *Bits Blog*, 1 July 2015, archive.nytimes.com/bits.blogs.nytimes.com/2015/07/01/google-photos-mistakenly-labels-black-people-gorillas/.

[4] Grant, Nico, and Kashmir Hill. "Google's Photo App Still Can't Find Gorillas. And Neither Can Apple's." *The New York Times*, 22 May 2023, www.nytimes.com/2023/05/22/technology/ai-photo-labels-google-apple.html.

[5] Hern, Alex. "Google's Solution to Accidental Algorithmic Racism: Ban Gorillas." *The Guardian*, The Guardian, 12 Jan. 2018, www.theguardian.com/technology/2018/jan/12/google-racism-ban-gorilla-black-people.