

# Addendum to GPT-4o System Card: Native image generation

OpenAI

March 25, 2025

1	Introduction	3
2	Observed Safety Challenges, Evaluations, and Mitigations	3
2.1	Safety Challenges: New risks from native image generation	3
2.2	Safety stack	4
2.3	Evaluations	4
2.3.1	External, manual red teaming	5
2.3.2	Automated red teaming	5
2.3.3	Offline testing using real-world scenarios	6
2.4	Discussion of specific risk areas	6
2.4.1	Child Safety	6
	Detection mechanisms	6
2.4.2	Artist Styles	7
2.4.3	Public Figures	8
2.4.4	Bias	8
	Statistical Bias	8
	Gender	9
	Ahistorical and Unrealistic Bias	11
2.4.5	Other risk areas evaluated	11
	Erotic content	11
	Imagery that is violent, abusive or hateful	12
	Instructions for illicit activities	13
2.5	Our approach to provenance	13
2.6	Conclusion	13
	Authorship, credit attribution, and acknowledgments	14
	References	15

# 1 Introduction

4o image generation is a new, significantly more capable image generation approach than our earlier DALL·E series of models. It can create photorealistic output. It can take images as inputs and transform them. It can follow detailed instructions, including reliably incorporating text into images. And because it is embedded natively, deep in the architecture of our omnimodal GPT-4o model, 4o image generation can use everything it knows to apply these capabilities in subtle and expressive ways, creating images that are not only beautiful, but also useful.

4o image generation benefits from our existing safety infrastructure, and from lessons we have learned deploying DALL·E and Sora. At the same time, these new capabilities also bring some new risks. This addendum to the GPT-4o system card describes the marginal risks we’ve focused on, and the work we have done to address them.<sup>1</sup>

## 2 Observed Safety Challenges, Evaluations, and Mitigations

### 2.1 Safety Challenges: New risks from native image generation

Unlike DALL·E, which operates as a diffusion model, 4o image generation is an autoregressive model natively embedded within ChatGPT. This fundamental difference introduces several new capabilities that are distinct from previous generative models, and that pose new risks:

- **Image-to-Image Transformation:** This capability allows 4o image generation to take one or multiple images as input, and to produce a related or modified image.
- **Photorealism:** The advanced photorealistic capabilities of 4o image generation mean that its outputs can, in some cases, have the appearance of a photograph.
- **Instruction Following:** 4o image generation can follow detailed instructions, and render text and instructional diagrams, introducing both utility and risk that is distinct from earlier models.

---

<sup>1</sup>Per our [Preparedness Framework](#), the launch of 4o image generation did not trigger additional Preparedness evaluations beyond those originally conducted for GPT-4o.

These capabilities, alone and in new combinations, have the potential to create risks across a number of areas, in ways that previous models could not. For example, without safety controls, 4o image generation could create or alter photographs in ways that could be detrimental to the people depicted, or provide schematics and instructions for making weapons.

Drawing on our experience with multimodal models and with the Sora and DALL·E visual generation tools, we’ve mapped and addressed a range of net-new risks specific to 4o image generation.

We strive to maximize helpfulness and creative freedom for our users while minimizing harm (read more in our [Model Spec](#)). As we learn more about how people are actually using 4o image generation, in line with our commitment to [iterative deployment](#), we will continue to evaluate our policies and adjust them as appropriate. And as always, users must abide by our [usage policies](#) when using any of our products, including 4o image generation.

## 2.2 Safety stack

To address the unique safety challenges posed by 4o image generation, several mitigation strategies are in use:

- **Chat model refusals:** In ChatGPT and the API, the primary chat model acts as a first line of defense against the generation of content that violates our policies. Based on its post-training safety measures, the chat model can refuse to trigger the image generation process based on the user’s prompt.
- **Prompt blocking:** This strategy, which happens after a call to the 4o image generation tool has been made, involves blocking the tool from generating an image if text or image classifiers flag the prompt as violating our policies. By preemptively identifying and blocking prompts, this measure helps prevent the generation of disallowed content before it even occurs.
- **Output blocking:** This approach, applied after an image has been generated, uses a combination of controls including Child Sexual Abuse Material (CSAM) classifiers and a safety-focused reasoning monitor to block the output of images that violate our policies. The monitor is a multimodal reasoning model which is custom-trained to reason about content policies. By evaluating the output post-generation, this strategy aims to block any content that is disallowed under our policies, providing an additional safeguard against the creation of disallowed content.
- **Increased safeguards for minors:** We use all of the mitigations listed above to create an even safer experience for users we believe may be under 18 and seek to limit those users from creating certain categories of potentially age-inappropriate content. Users under the age of 13 are currently prohibited from using any of OpenAI’s products or services.

## 2.3 Evaluations

We evaluated the safety and effectiveness of 4o image generation’s safety stack by observing its performance with prompts from three sources.

1. External, manual red teaming

2. Automated red teaming
3. Offline testing using real-world scenarios

### 2.3.1 External, manual red teaming

OpenAI worked with a cohort of vetted external red teamers from our Red Teaming Network and Scale AI to test 4o image generation. We began external red teaming following internal testing of 4o image generation to assess raw model capabilities and inform focus areas for testing.

We asked these external red teamers to explore various prioritized topic areas, including the areas discussed below. We also enabled red teamers to develop and use various jailbreaks and tactics to attempt to circumvent the model’s safeguards.

After completing manual red teaming, we combined thousands of these manual adversarial conversations and converted them into automated evaluations. We re-ran our safety stack on this dataset and track these two main metrics:

- not\_unsafe: does the system produce output that violates our model policies?
- not\_overrefuse: does the system refuse to comply with a request that complies with our model policies?

Table 1: Overall metrics – performance using external red teaming data

4o image generation	Not_unsafe	Not_overrefuse
With system mitigations only (prompt blocking and output blocking)	0.955	0.941
With system mitigations and chat model refusals	0.971	0.856

### 2.3.2 Automated red teaming

In automated red teaming we use the model policies noted above to generate synthetic conversations that probe the system’s performance for each part of the model policy. These synthetic conversations enable us to test the system’s implementation of the policies more thoroughly than we could with manual red teaming alone.

We generated thousands of synthetic conversations across different categories, with and without image uploads, in order to complement the testing work of the manual red teamers.

Table 2: Overall metrics - performance using automated red teaming data

4o image generation	Not_unsafe	Not_overrefuse
With system mitigations only (prompt blocking and output blocking)	0.969	0.899
With system mitigations and chat model refusals	0.975	0.830

This shows similar performance as the red-teaming data from humans, creating additional confidence that our policies are working consistently across a range of conversations.

### 2.3.3 Offline testing using real-world scenarios

We also evaluated the 4o image generation safety stack on textual prompts reflecting real-world scenarios to evaluate the model’s behavior in a production environment. This involves examples from across different safety categories, in order to make the evaluation representative of the actual distribution encountered in production. This helps us understand how well the model performs in live conditions and highlights any areas that may require additional safety measures.

Table 3: Overall metrics - performance using real-world scenarios

<b>4o image generation</b>	<b>Not_unsafe</b>	<b>Not_overrefuse)</b>
With system mitigations only (prompt blocking and output blocking)	0.929	0.996
With system mitigations and chat model refusals	0.932	0.993

## 2.4 Discussion of specific risk areas

### 2.4.1 Child Safety

OpenAI is deeply committed to addressing child safety risks, and we prioritize prevention, detection, and reporting of Child Sexual Abuse Material (CSAM) content across all our products, including 4o image generation. OpenAI’s efforts in the child safety space encompass red teaming in accordance with Thorn’s recommendations, and robust scanning for CSAM across all inputs and outputs, for both first party and third party users (API and Enterprise).

Specific model policies for child safety in 4o image generation include:

- At launch, editing uploaded images of photorealistic children will not be allowed. We will evaluate whether we can safely allow edits in the future.
- We have reinforced existing protections against child sexual abuse material (CSAM) for both image editing and image generation.

### Detection mechanisms

For Child Safety we leverage three different input mitigations across text and image inputs:

- For all image uploads, we integrate with Safer, developed by Thorn, to detect matches with known CSAM. Confirmed matches are rejected and reported to NCMEC, and the associated user account is banned. Additionally, we utilize Thorn’s CSAM classifier to identify potentially new, unhashed CSAM content on both image uploads and images generated by 4o image generation.

- We leverage a multi-modal moderation classifier to detect and block any generated sexual content that involves minors.
- For 4o image generation, we built a photorealistic-person classifier based on our existing under-18 classifier used for Sora to analyze all uploaded images to predict whether any of them depicts a minor. At launch, photorealistic generation of children is permitted only when it is not an image edit of a photorealistic minor. Additionally, photorealistic generations of children must comply with the safety constraints across all of our policies.

The photo-realistic person classifier takes in uploaded image(s) and predicts one of the three labels:

1. No photorealistic person
2. Photorealistic adult
3. Photorealistic child

If an image contains both a photorealistic adult and a photorealistic child, the classifier is designed to return “photorealistic child” as a prediction.

Below is our evaluation for our classifier on a dataset containing close to 4000 images across the categories of [child | adult] and [photorealistic | non-photorealistic].

Currently, our classifiers are highly accurate, but they may occasionally misclassify images. For example, younger-looking adults may be incorrectly flagged as children. For safety purposes, we have tuned the classifier to err on the side of caution by classifying borderline or ambiguous cases as “child.” We are committed to enhancing the performance of our classifier using better models and better evaluation sets in the future.

Table 4: Results for photorealistic-person classifier

	<b>n_samples</b>	<b>precision</b>	<b>recall</b>
Photorealistic person (adult or child)	2033	0.905	0.99
Photorealistic adult	919	0.80	0.776
Photorealistic child	1113	0.80	0.97

### 2.4.2 Artist Styles

The model can generate images that resemble the aesthetics of some artists’ work when their name is used in the prompt. This has raised important questions and concerns within the creative community. In response, we opted to take a conservative approach with this version of 4o image generation, as we learn more about how 4o image generation is used by the creative community. We added a refusal which triggers when a user attempts to generate an image in the style of a living artist.

### 2.4.3 Public Figures

4o image generation is capable, in many instances, of generating a depiction of a public figure based solely on a text prompt.

At launch, we are not blocking the capability to generate adult public figures but are instead implementing the same safeguards that we have implemented for editing images of photorealistic uploads of people. For instance, this includes seeking to block the generation of photorealistic images of public figures who are minors and of material that violates our policies related to violence, hateful imagery, instructions for illicit activities, erotic content, and other areas. Public figures who wish for their depiction not to be generated can opt out.

This approach is more fine-grained than the way we dealt with public figures in our DALL·E series of models, where we used technical mitigations intended to prevent any images of a public figure from being generated. This change opens the possibility of helpful and beneficial uses in areas like educational, historical, satirical and political speech. After launch, we will continue to monitor usage of this capability, evaluating our policies, and will adjust them if needed.

### 2.4.4 Bias

4o image generation performs better than our earlier image generation tools for certain areas related to representational biases. However, challenges remain—particularly around demographic representation—and we recognize that additional work is needed to ensure balanced outputs. We plan to continue refining our approach and invest in developing more effective post-training mitigations — including incorporating more diverse examples into the post-training mixture — to improve model outputs in the coming months. Below are the results of our evaluations. We find that across all of our metrics, 4o image generation exhibits less bias than DALL·E 3.

#### Statistical Bias

We ran automated bias evaluations against under-specified prompts describing individuals, such as “a happy person” or “a doctor” and groups, such as “generate an image of three construction workers”. We report the following information for each attribute:

- **Class distributions:** The class distributions of individuals generated in response to these prompts<sup>2</sup>. We provide this for information purposes, but expect the two numbers explained below — heterogeneous output frequency and skew — to be more useful for understanding model performance on bias.
- **Frequency of heterogeneous outputs:** The proportion of prompts that result in more than one value of a given attribute across 20 resamplings (within our prompt sets). Higher values are better, as they indicate that the model does not consistently produce images with a single attribute (e.g., always the same gender) for a given class. Resampling the model many times will yield better results.
- **Skew:** Shannon entropy [1] for a given class, where 0 is a distribution consisting of entirely one class, while 1 is a uniform distribution. We don’t expect our evaluations to be either

---

<sup>2</sup>Because DALL·E 3 is called via a tool call, sometimes a user prompt may never cause a tool call. We normalize all numbers so they sum up to 100%.



extreme, but this statistic helps us represent, directionally, which distribution a given model trends towards.<sup>3</sup>

When someone asks for an image without giving specific attributes—like requesting “an image of a doctor” without specifying gender or race—our data shows that 4o image generation creates a broader range of results than DALL·E 3. This data offers a quantitative approach for assessing how varied the images are, but it doesn’t suggest there’s a single “correct” or ideal balance of characteristics (like gender or race) in the images we produce.

We compute class probabilities for a set of images generated against the same prompt, but these probabilities are hard to interpret across a wider prompt set. To alleviate this, we specifically measure the frequency of heterogeneous outputs and attribute skew. Heterogeneous outputs indicate instances where, across a set of images generated for a single prompt, the depicted subject at least once represents a class other than the most common class. Attribute skew indicates how balanced our models’ portrayals are across various demographic attributes. These measurements align with our goal of generating diverse and authentic representations.

We highlight that users have further control over the default model behavior through personalization settings and by explicitly specifying attributes in prompts. We aim to use our evaluation framework not only to track default model behavior but also to ensure adherence to user preferences. Generated images typically contain many details that are not directly specified by the prompt. We aim for the model to fill in those details in ways that reflect relevant context, including reflecting a relevant range of possibilities rather than defaulting exclusively to the most common demographics. As highlighted previously in our DALL·E 3 report, our choices and refinements in these areas may not exactly match the demographic composition of any specific cultural or geographic region. However, we remain committed to balancing authentic representations, user preferences, and inclusivity in our image generation models, with the eventual goal of image generation of underspecified prompts that is more localized to any one user’s particular location.

## Gender

Currently, despite 4o image generation surpassing DALL·E 3 in gender representation diversity, the outputs still predominantly favor male subjects. As a result, our future work will focus on increasing the frequency of heterogeneous outputs and Shannon entropy, using these as key metrics for measuring progress toward a more representative model.

Table 5: Class distributions for gender

Prompt Set	Model	Male	Female
Individuals	DALL·E 3	86%	14%
	4o image generation	79%	21%
Groups	DALL·E 3	61%	35%
	4o image generation	56%	44%

<sup>3</sup>**Shannon entropy** is a metric from information theory that quantifies uncertainty or unpredictability in a distribution. Low entropy (0) means the distribution is highly skewed — nearly all predictions fall into one class. High entropy (1) means the distribution is uniform — the model is equally likely to assign any class.

Table 6: Shannon entropy and frequency of heterogeneous outputs for gender

Prompt Set	Model	Shannon Entropy	Frequency of heterogeneous outputs
Individuals	DALL·E 3	0.17	35%
	4o image generation	0.27	46%
Groups	DALL·E 3	0.82	95%
	4o image generation	0.95	100%

## Race [2]

While both DALL·E 3 and 4o image generation tend to produce individuals categorized as white more frequently than other racial groups, 4o generates a noticeably broader variety of individuals in response to a given prompt.

Table 7: Class distributions for race

Category	Model	White	Black	East Asian	Indian	Latino	Middle Eastern	Southeast Asian
Individuals	DALL·E 3	90%	0%	7%	0%	1%	2%	0%
	4o image generation	67%	19%	2%	2%	5%	5%	0%
Groups	DALL·E 3	81%	3%	13%	0%	2%	1%	0%
	4o image generation	64%	21%	4%	3%	6%	1%	0%

We observe improved performance, with more frequent heterogeneous outputs and higher Shannon entropy than DALL·E 3.

Table 8: Shannon entropy and frequency of heterogeneous outputs for race

Prompt Set	Model	Shannon Entropy	Frequency of heterogeneous outputs
Individuals	DALL·E 3	0.13	52%
	4o image generation	0.36	85%
Groups	DALL·E 3	0.27	80%
	4o image generation	0.50	100%

## Skin Tone [3]

When assessing the skin tone of the individuals generated by DALL·E 3 and 4o image generation, we find that both models tend to produce individuals categorized as lighter skinned in response to the majority of prompts, but the vast majority of prompts also produce a set of images with a diverse set of skin tones.

Table 9: Class distributions for skin tone

Prompt Set	Model	Light	Medium	Dark	Very Dark
Individuals	DALL·E 3	90%	10%	0%	0%
	4o image generation	59%	29%	12%	0%
Groups	DALL·E 3	88%	9%	3%	0%
	4o image generation	62%	24%	15%	0%

Table 10: Shannon entropy and frequency of heterogeneous outputs for skin tone

Prompt Set	Model	Shannon Entropy	Frequency of heterogeneous outputs
Individuals	DALL·E 3	0.18	61%
	4o image generation	0.50	96%
Groups	DALL·E 3	0.26	89%
	4o image generation	0.61	100%

## Ahistorical and Unrealistic Bias

We ran an automated evaluation to determine whether the model might output ahistorical, unrealistic, or undesired attributes contrary to the user’s intent, such as changing the race of a well-specified prompt (“A stereotypical Indian person”) or historically well-specified population (“The founding fathers”). These evaluations focus solely on the model’s behavior when demographics are not explicitly specified. If a user does specify attributes, we expect the model to follow the user’s prompt, even if that means being historically inaccurate.

The score we produce is the percent of time that the attributes within produced images match the expected attributes - a higher score indicates a closer alignment to these expectations. These examples should yield predictable results with no variation (heterogeneous outputs at 0% and skew of 0), because they refer to contexts in which historical and realistic depictions are demographically uniform. This evaluation helps us distinguish intentional, accurate depictions from unintended bias. 4o image generation saturates our internal evaluation of this.

Table 11: Ahistorical and unrealistic bias results

Model	Score (Higher is Better)
DALL·E 3	92%
4o image generation	97%

### 2.4.5 Other risk areas evaluated

In line with our [Model Spec](#), we aim to maximize creative freedom by supporting valuable use cases like game development, historical exploration, and education—while maintaining strong safety standards. At the same time, it remains as important as ever to block requests that violate those standards. Below are evaluations of additional risk areas where we’re working to enable safe, high-utility content and support broader creative expression for users.

We slice the human curated and automated red teaming data based on different risk areas and evaluate that the model does not comply with requests that violate our standards, while also not overrefuse for requests that maximize creative freedom. We evaluate completions using an autograder, checking two main metrics, `not_unsafe` and `not_overrefuse`.

#### Erotic content

In 4o image generation, model policies related to erotic content include:

- We aim to prevent attempts to generate erotic or sexually exploitative imagery.

- We have heightened safeguards designed to prevent nonconsensual intimate imagery or any type of sexual deepfakes.

Table 12: Safety evaluation results - erotic content

Eval	N examples		System	not_unsafe	not overrefuse
Human-curated red teaming		364	With system mitigations	0.971	0.912
			With system mitigations and chat model refusals	0.979	0.884
Automated red teaming	red	927	With system mitigations	0.990	0.875
			With system mitigations and chat model refusals	0.992	0.859

## Imagery that is violent, abusive or hateful

Specific model policies for violent, abusive, and hateful imagery in 4o image generation include:

- Depicting violence in artistic, creative or fictional contexts is generally allowed to enable creative and artistic endeavors. But we aim to prevent the model from generating photorealistic, graphically violent imagery in certain contexts.
- We aim to prevent attempts to generate images that promote or facilitate self-harm (including, e.g., providing instructions for self-harm). We incorporate additional self-harm protections for certain users, including users we believe may be under the age of 18.
- We have included mitigations intended to prevent attempts to generate extremist propaganda and recruitment content. We incorporate heightened extremist content protections for certain users, including users we believe may be under the age of 18. We allow users to generate hateful symbols in a critical, educational, or otherwise neutral context, as long as they don't clearly praise or endorse extremist agendas.
- Many types of abuse are context-dependent. While we restrict the ability to create clearly harmful imagery with someone's likeness, users may find ways to bully or harass someone with this model in ways that would only be apparent to the intended recipient of harassment. People can report potential abuse through our help center, and we will continue iterating on our safety mitigations over time as we see new types of abuse arise.

Establishing clear policy boundaries between harmful graphic violence and violence depicted for creative, educational, or documentary purposes – or between bullying and self-deprecating humor – is challenging. We're adopting a more permissive approach for these edge cases compared to our previous DALL·E policies, while taking extra caution to protect users who may be under 18. We believe this strategy helps us learn from real-world usage and find the right balance between enabling valuable use cases and preventing harm.

Table 13: Safety evaluation results - imagery that is violent, abusive, or hateful

Eval	N examples		System	not_unsafe	not overrefuse
Human-curated red teaming		1266	With system mitigations	0.914	0.917
			With system mitigations and chat model refusals	0.952	0.795
Automated red teaming		1627	With system mitigations	0.959	0.889
			With system mitigations and chat model refusals	0.968	0.821

## Instructions for illicit activities

In 4o image generation, we take a similar approach for illicit activities that we have with our other models. We aim to prevent attempts to generate images with advice or instructions related to weapons, violent wrongdoing, or other illicit activities such as theft.

Table 14: Safety evaluation results - instructions for illicit activities

Eval	N examples	System	not_unsafe	not overrefuse
Human-curated red teaming	25	With system mitigations	0.999	0.959
		With system mitigations and chat model refusals	0.999	0.959
Automated red teaming	309	With system mitigations	0.972	0.974
		With system mitigations and chat model refusals	0.977	0.948

## 2.5 Our approach to provenance

Based on learnings from DALL·E and Sora, we have continued to prioritize enhancing our provenance tools. For general availability of 4o image generation, our provenance safety tooling will include:

- C2PA metadata on all assets (verifiable origin, industry standard).
- Internal tooling to help assess whether a certain image is created by our products.

We recognize that there is no single solution to provenance, but are committed to improving the provenance ecosystem, continuing to collaborate on this issue across industry and with civil society, and helping build context and transparency to content created from 4o image generation and across our products.

## 2.6 Conclusion

By launching 4o image generation together with the safety work described in this system card, we are continuing our longstanding commitment to [a rigorous, iterative approach](#) to making AI systems safe. This system card provides a snapshot of our safety approach at launch, and we look forward to continuing to refine and strengthen our safety work as we learn from this and future deployments.

# Authorship, credit attribution, and acknowledgments

Please cite this work as “OpenAI (2025)”.

## Leadership

**Gabriel Goh:** Image Generation, **Jackie Shannon:** ChatGPT Product, **Mengchao Zhong, Wayne Chang:** ChatGPT Engineering, **Rohan Sahai:** Sora Product and Engineering, **Brendan Quinn, Tomer Kaftan:** Inference, **Prafulla Dhariwal:** Multimodal Organization

Matt Chan

## Data Science

Xiaolin Hao

## ChatGPT

Andrew Sima, Annie Cheng, Benjamin Goh, Boyang Niu, Dian Ang Yap, Duc Tran, Edede Oiwoh, Eric Zhang, Ethan Chang, Jay Chen, Jeffrey Dunham, Kan Wu, Karen Li, Kelly Stirman, Mengyuan Xu, Michelle Qin, Ola Okelola, Pedro Aguilar, Rocky Smith, Rohit Ramchandani, Sean Fitzgerald, Wanning Jiang, Wesam Manassra, Xiaolin Hao, Yilei Qian

## Research

### Foundational Research

Allan Jabri, David Medina, Gabriel Goh, Kenji Hata, Lu Liu, Prafulla Dhariwal

### Core Research

Aditya Ramesh, Alex Nichol, Casey Chu, Cheng Lu, Dian Ang Yap, Heewoo Jun, James Betker, Jianfeng Wang, Li Jing, Long Ouyang, Wesam Manassra

### Research Contributors

Aiden Low, Brandon McKinzie, Charlie Nash, Huiwen Chang, Ishaan Gulrajani, Jamie Kiros, Ji Lin, Kshitij Gupta, Yang Song

### Model Behavior

Laurentia Romaniuk

### Multimodal Organization

Andrew Gibiansky, Yang Lu

## Sora

### Sora Product Lead

Rohan Sahai, Wesam Manassra

### Sora Product and Engineering

Boyang Niu, David Schnurr, Gilman Tolle, Joe Taylor, Joey Flynn, Mike Starr, Rajeev Nayak, Rohan Sahai, Wesam Manassra

## Safety

### Safety Lead

Somay Jain

### Safety

Andrea Vallone, Alex Beutel, Botao Hao, Brendan Quinn, Cameron Raymond, Chong Zhang, David Robinson, Eric Wallace, Filippo Raso, Huiwen Chang, Ian Kivlichan, Irina Kofman, Keren Gu-Lemberg, Kristen Ying, Madelaine Boyd, Meghan Shah, Michael Lampe, Owen Campbell-Moore, Rohan Sahai, Rodrigo Riaz Perez, Sandhini Agarwal, Sam Toizer, Troy Peterson

## Data

### Data Leads

Gildas Chabot, James Park Lennon

### Data

Arshi Bhatnagar, Dragos Oprica, Rohan Kshirsagar, Spencer Papay, Szi-chieh Yu, Wesam Manassra, Yilei Qian

### Moderators

Hazel Byrne, Jennifer Luckenbill, Mariano López

### Human Data Advisor

Long Ouyang

## Scaling

### Inference Leads

Brendan Quinn, Tomer Kaftan

### Inference

Alyssa Huang, Jacob Menick, Nick Stathas, Ruslan Vasilev, Stanley Hsieh

## Applied

### ChatGPT Product Lead

Jackie Shannon

### ChatGPT Engineering Lead

Mengchao Zhong, Wayne Chang

### Product Design Lead

## Marketing & Comms

### Comms and Marketing Leads

Minnia Feng, Natalie Summers, Taya Christianson

### Comms

Alex Baker-Whitcomb, Ashley Tyra, Bailey Richardson, Gaby Raila, Scott Ethersmith, Marselus Cayton, Souki Mansoor

## Design & Creative

### Leads

Kendra Rimbach, Veit Moeller

## Design

Adam Brandon, Adam Koppel, Angela Baek, Cary Hudson, Dana Palmie, Jeffrey Sabin Matsumoto, Leyan Lo, Matt Nichols, Thomas Degry, Vanessa Antonia Schefke, Yara Khakbaz

Aidan Clark, Aditya Ramesh, Alex Beutel, Ben Newhouse, Ben Rossen, Che Chang, Greg Brockman, Hannah Wong, Ishaan Singal, Jason Kwon, Jiahui Yu, Jiacheng Feng, Joanne Jang, Johannes Heidecke, Kevin Weil, Mark Chen, Mia Glaese, Nick Turley, Raul Puri, Reiichiro Nakano, Rui Shu, Sam Altman, Shuchao Bi, Vinnie Monaco

## Special Thanks

We are grateful to our group of red teamers who provided feedback, helped test our models at early stages of development and informed our risk assessments and evaluations. Participation in the testing process is not an endorsement of the deployment plans of OpenAI or OpenAI's policies.

## References

- [1] C. E. Shannon, "A mathematical theory of communication," *Bell System Technical Journal*, 27(3), 379–423, 1948.
- [2] K. Karkkainen and J. Joo, "Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation," 2021.
- [3] E. Monk, "Monk skin tone scale." <https://skintone.google/>, 2019.