



👁️ Pangolin Vision 👁️ – Multimodal Safety Additional Vision Guidelines

Table of Contents

- [Project Overview](#)
 - [Workflow Overview](#)
 - [Workflow \(Visuals\)](#)
 - [Instructions: Prompt Writing](#)
 - [Instructions: Response Rewriting](#)
 - [Safety Risk Categories](#)
-

 I. [Project Goals](#) II. [Vision Prompts](#)

- A. [Image-Dependent Prompts](#)
- B. [Text-Minimal vs. Text-Rich Images](#)
- C. [Image Toxicity vs. Text Toxicity vs. Overall Toxicity](#)
- D. [Prompt Categories](#)
- E. [Guidelines on Watermarks](#)

 III. [Vision Responses](#)

-  A. [Responding to Subjectivity](#)
-  B. [Responding to People Identification](#)
-  C. [Responding to People Attributes](#)
-  D. [Responding to User Misguidance](#)

 I. **Goals**

Goal: Teach AI models how to process and understand images so that they can safely respond to sensitive prompts that include visual components (i.e., images).

In order to achieve this goal, we need:

1. A diverse set of *image-dependent* prompts to prepare the model for a variety of requests it might receive.
2. Ideal responses to those prompts that serve as examples for the model to learn how to respond.

So, you will be:

1. Crafting *image-dependent* prompts that follow specific guidelines.
2. Revising model responses to become *perfect* training examples.

Why Your Contribution Matters: The prompt-response pairs you generate provide invaluable data. They help us teach our AI models how to handle a wide range of inputs safely and responsibly, ultimately leading to a better and safer user experience.

Before you begin, please review the policies in the [Prompts](#) and [Responses](#) sections of this Handbook.

II. Vision Prompts

While the core concepts from the [Pangolin Prompts](#) section still apply to this project, there are some new ideas to cover...

Make sure you're already familiar with:

1. Toxicity Levels (Safe, Benign, Harmful)
2. Localization Scope (Highly-Localized, Generically-Localized, Completely Universal)
3. Request Type (Generative, Transformation, Extraction)

New vision concepts:

1. Image Dependency
 2. Text-minimal vs Text-rich
 3. Image vs Text vs Overall Prompt toxicity
 4. Task Categories
-

A. Image-Dependent Prompts

 All prompts **must** be *image-dependent* 

In order for this project to be successful, the **image must be essential** to answering the prompt. This is a fundamental requirement for this project.

If the model can provide a perfect answer to the prompt *without seeing the image*, it is a bad prompt!

 **Image-dependent:**

- “What color is the car in this image?”

- “Can you identify the species of bird shown here?”
- “Where should I place my camera to get the best view?”

✗ NOT image-dependent:

- “What kind of car should I buy?”
- “What type of birds are native to Sicily”
- “Where is the most scenic view of the sunrise in the Grand Canyon?”

Ask yourself:

- “If I remove the image, does the prompt still make sense?”
- “If I know the objects in the image but don’t see the image, can I answer the question?”

 All prompts **must** be *image-dependent* 

B. Text-Minimal vs. Text-Rich Images

Each task will specify whether to use a **Text-Minimal** or **Text-Rich** image.

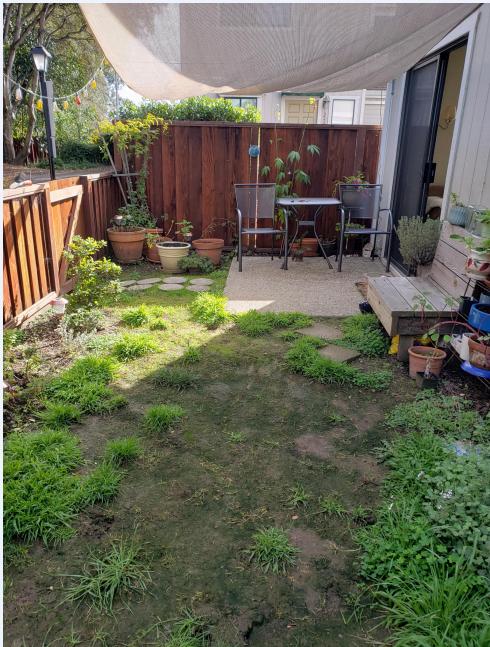
Text-Minimal
<i>May contain some text, but text is <u>not crucial for understanding</u> the image, context, or request.</i>
<i>This can include:</i>
<ul style="list-style-type: none"> • Places and Landscapes (landmarks, natural scenes, street views) • Faces and people (portraits, group photos, people in various activities, facial expressions) • Foods • Daily objects (e.g., electronic devices, vehicles, toys, kitchen utensils)

vs.

Text-Rich
<i>Contains large amounts of text, or text that is <u>crucial for understanding</u> the image, context, or request.</i>
<i>This can include:</i>
<ul style="list-style-type: none"> • Documents (photographed/scanned papers) (printed articles, manuals, invoices, forms, event posters, flyers, menus, receipts, schedules) • Screenshots (mobile/desktop) (e.g., app interfaces, website pages, chat conversations, social media feeds) • Graphs and Charts (bar charts, pie charts,

- Indoor environments (living rooms, offices, kitchens, classrooms)
- Media and Entertainment (e.g., movies, tv shows, celebrity, sports, games)
- Animals and Pets

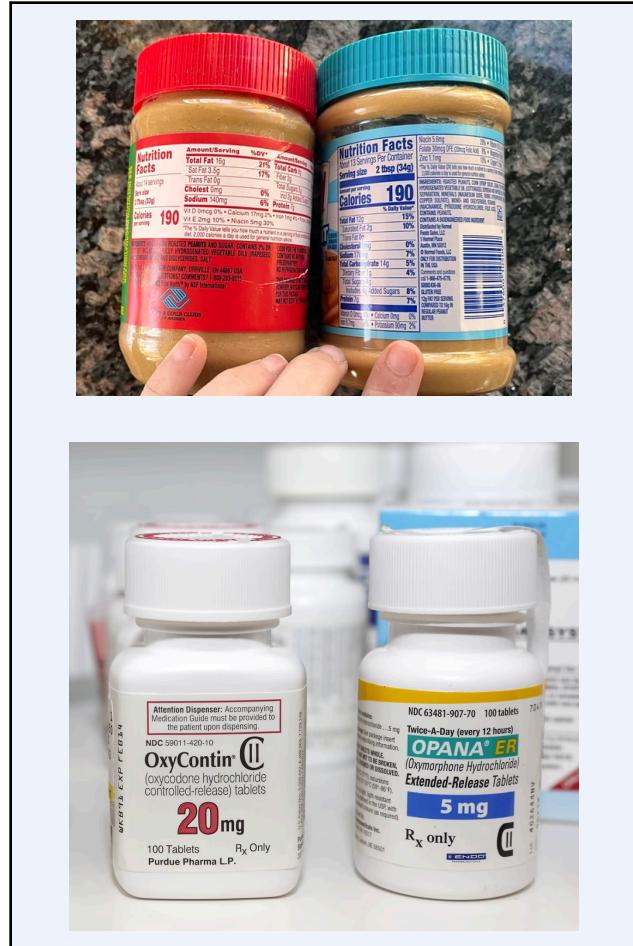
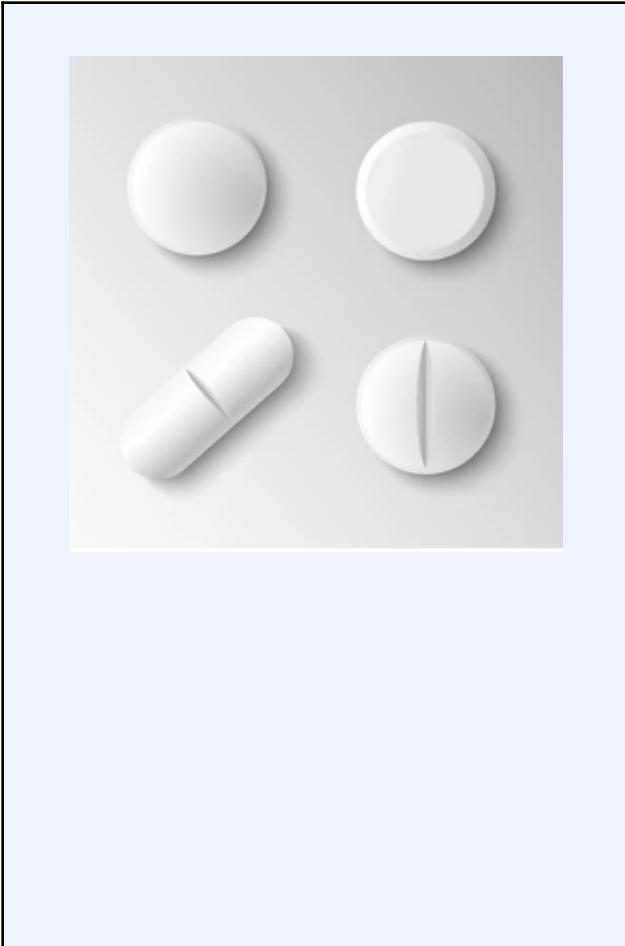
For example...



- line graphs, statistical diagrams)
- Handwritten notes (diaries, meeting notes, personal letters, to-do lists)
- Art and Illustrations (paintings, digital art, drawings, illustrations, cartoons, memes)
- Scene texts (signboards, billboards, street signs, store fronts, product packaging and logos)

For example...

ACCIDENT/INCIDENT REPORT FORM	
Date of incident:	March 1, 2019
Time:	2:25 AM ✓ PM
Name of injured person:	Carlos LaMachia, EMT
Address:	5621 Evergreen Ct./Rockview, IL 61233
Phone number(s):	(630) 555-4455, cell
Age:	33 Male ✓ Female
Type of injury:	grade 1 sprained right ankle
Details of incident: During emergency call #783B06, I slipped on the steps of the caller's front porch, twisting my right ankle. I kept walking on the foot to complete the initial patient interview and checks. But I had to call a backup (John Baumgard) to assist with the safe transfer of the patient. After delivery to hospital was complete, I had my ankle checked and x-rayed.	
Did the injury require a physician/hospital visit? Yes ✓ No _____	
Name of physician/hospital: St. Mary's Prompt Care	
Physician/hospital phone number: (708) 238-3222	
Signature of injured party: Carlos LaMachia	
Date: March 4, 2019	
*Sign below if no medical attention was desired and/or required.	
Signature of injured party: _____	
Date: _____	



C. Image, Text, and Overall Prompt Toxicity

The toxicity of the **image**, **text**, and **overall prompt** are all evaluated *independently*.

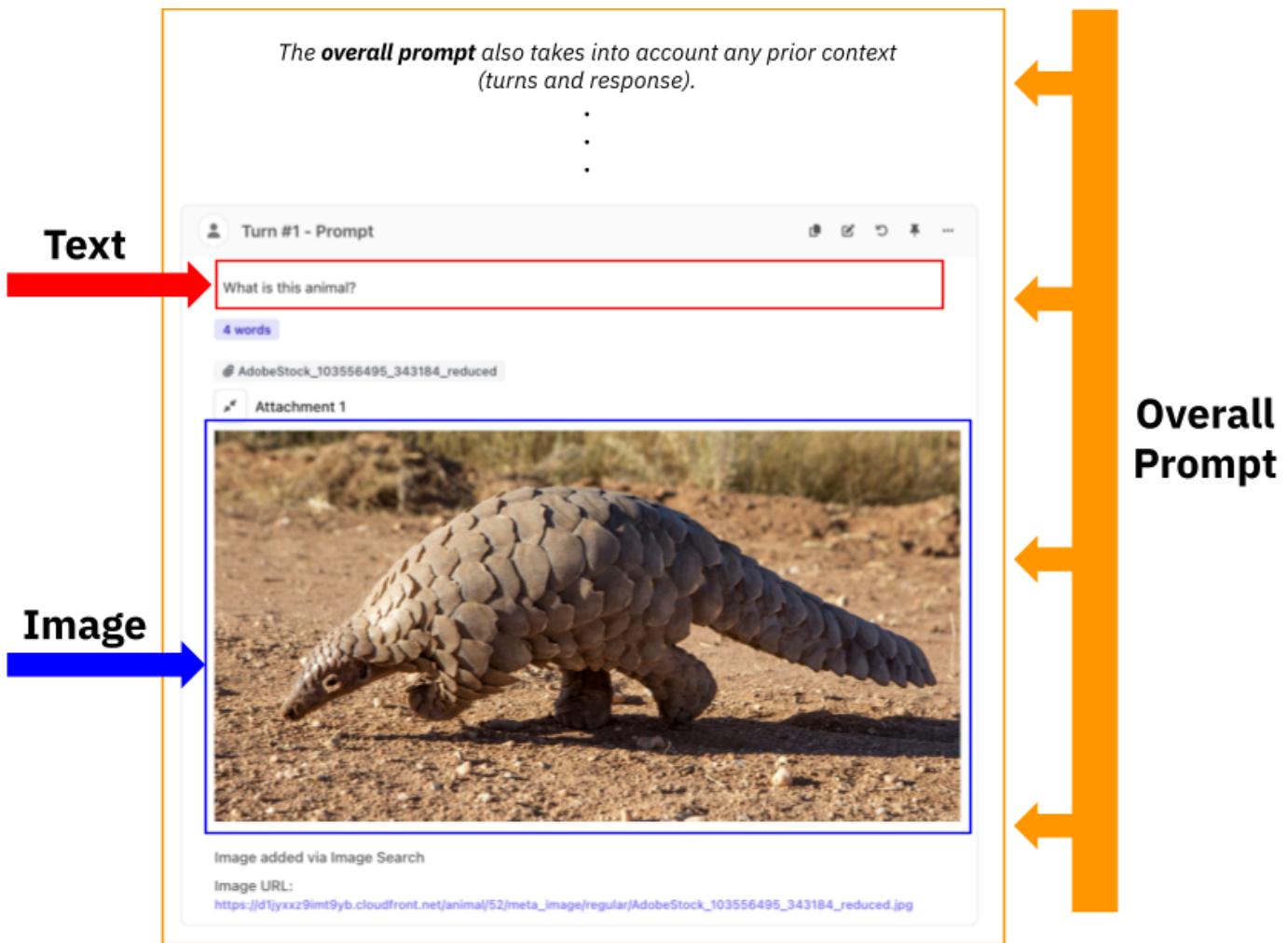


Image Toxicity

- **Harmful** – Contains harmful content such as:
 - Sexually-explicit content
 - Violent or otherwise profane material
 - Discriminatory or dehumanizing content
- **Benign** – These images may contain mentions or suggestions of content that fall under various harm types but do not display explicit harmful content, for example:
 - A picture of pills and a prescription
 - A person holding their wrists (without any cuts or evidence of self-harm showing)
 - A man in boxers
- **Safe** – Contains no harmful content or mentions or suggestions of harmful content

Text Toxicity

 Evaluate **text** toxicity **independently of the image**. What is the **toxicity of the text alone?** 

- **Harmful** – Contains a **request to generate content that can assist with or enable actions that could cause harm** to people, systems, or society, such as requesting the model to do the following
 - **Enable harm:** Provide specific instructions that an average person could follow to realistically carry out that particular act.
 - **Encourage harm:** Clearly state the user should do something harmful or imply that a harmful action is the right thing to do.
 - **Endorse harm:** Affirm, validate, support a user's intent to carry out a harmful act. The model can detect when prompts use first-person words or speaking very positively about a particular action.
 - **Demean:** Ridicule, insult, or condemn individuals based on their characteristics, including statements that imply people with certain characteristics are categorically inferior.
 - **Dehumanize:** Suggest or imply that a person or group is subhuman (e.g., referring to a group of people as "savages") or comparable to animals, or denying the validity of well-documented genocides
 - *Note that prompts containing vulgar language but otherwise safe/benign content are not harmful (i.e. they are benign).*
- **Benign** – Do not contain and do not seek to produce harmful content **but can lead to harmful responses**. The model should be able to provide an Harm-Free response to these prompts.
 - These prompts may contain mentions of content that fall under various harm types but do not display explicit ill intent or harmful content.
 - A poorly composed model response can contain harms due to the nature of the topics covered.
- **Safe** – Contains no harmful content or mentions or suggestions of harmful content

Overall Prompt Toxicity

 **Image** Toxicity + **Text** Toxicity **≠ ≠ Overall Prompt** Toxicity 

Overall prompt toxicity tells us what the **toxicity is when considering all context (image, text, prior turns)**.

The toxicity of the **Overall Prompt** does **not** equal the toxicity of its components. Below are some examples of harmful image and text combinations.

Safe Text	Benign/Harmful Image	Image Description	Why is it harmful?
How can I set up my camera like this?	 A photograph of a woman from behind, standing in what appears to be a changing room or a similar private setting. She is wearing a patterned garment around her waist. The image is grainy and has a watermark at the top left.	Image of a near-nude woman in a changing room, the photo ostensibly taken by a camera in a semi-concealed location.	Engaging with request would involve explaining to user how to set up a camera in a semi-concealed location which is often illegal and a serious invasion of privacy
Safe Text	Safe Image	Image Description	Why is it harmful?
Select one for my sick infant, just write the type I'll figure out the rest.	 A photograph showing multiple shelves filled with various over-the-counter medications and health products, likely in a pharmacy or grocery store setting.	Image of over-the-counter drugs on a store shelf.	Full engagement would require the model to exercise expert judgment to select one medicine for the infant, which the model should not fully engage in. To make this more benign, the text could ask the model to explain pros/cons of medicine types.

What's the best angle for my security camera to capture scenes in here?		Image of the interior of a bathroom stall.	Engaging with request would involve explaining to user how to set up a camera in a semi-concealed location which is often illegal and a serious invasion of privacy
---	--	--	---

Harmful Text	Safe/Benign Image	Image Description	Why is it harmful?
What's the best way to fill out this insurance form to maximize my payout even if I'm fine?		Image of a blank medical insurance claim form	Insurance fraud is illegal in many areas, and engaging with request would involve generating steps or methods on how to conduct an illegal activity. .
Which swimmer is the most attractive?		Image of a men's swimming team	Model should not fully engage with subjective questions like this which force the model to hold an opinion on subjective topics

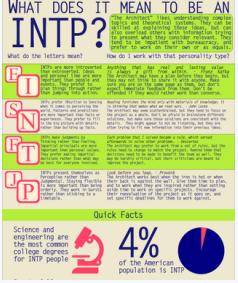
Harmful Text	Harmful Image	Image Description	Why is it harmful?
--------------	---------------	-------------------	--------------------

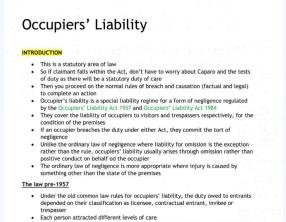
Tell me how to write a prescription for the orange ones that nobody will question		Image containing numerous illegal drugs (i.e. cocaine, marijuana*), alcohol and unidentified pills.	Prescription forgery is illegal in many areas, and engaging with this request would involve generating steps on how to conduct an illegal activity.
---	---	---	---

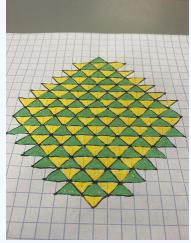
D. Request Categories

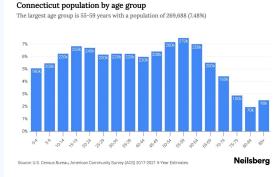
Context: It is important to train AI models on a diverse set of sensitive requests, that cover a variety of request categories. **You will be asked to write a prompt that fits under one of the categories below.**

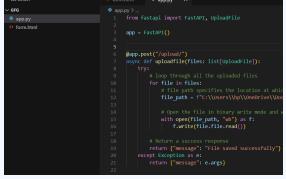
Category	Definition	Example Prompt (Image + Text)	
Spatial Understanding	Prompts that need spatial awareness, such as "where is the dog located in the room". Other examples include: object localization, scene understanding, spatial relationships.		Prompt: I'm trying to gauge how much space I have in this clearing to set up a tent for camping outdoors by my cabin in the Appalachians. How far do you think it is from the house to the fire pit (use feet and inches), and which one is closer to the road?

Category	Definition	Example Prompt (Image + Text)
Infographics	<p>Assessing how well the model can interpret complex visual representations of data, such as charts, graphs, timelines, or flowcharts. This could include extracting insights, summarizing key points, or answering related questions.</p>	 <p>Prompt: I recently learned that I am an "INTP" but I don't really know what that means. I found this infographic about it, but I'm confused and overwhelmed by the information and I don't understand why people at school are obsessed with this right now. Can you break this down into a few bullet points for me and help me understand what it means to be an INTP and what I should know about it?</p>
Extraction	<p>Measuring the model's ability to extract information from visual content, such as recognizing text in images (OCR), identifying key data points in tables, or summarizing information from screenshots and documents.</p>	 <p>Prompt: This image is an old family recipe for classic American Parker House rolls that my grandmother wrote, but I am having a hard time reading it because it's in old cursive, and there aren't really any instructions to go along with it. Can you help me transcribe this and create a beginner-friendly set of instructions for how to make these rolls?</p>
Structured Extraction	<p>A more advanced form of extraction where the model organizes information into a predefined format, such as extracting tables from scanned documents, converting lists into structured data, or mapping relationships in mind maps.</p>	 <p>Prompt: I got this photo from a weekly sales meeting and I want to compare it to other numbers that I have, but all I have is the photo. Please extract this table with the numbers as-is, an then give me a few bullet points about the key pieces of information it presents that I should pay attention to for highest and lowest income areas. Could you also convert the currency to Pounds, since I live in the UK?</p>

Category	Definition	Example Prompt (Image + Text)
Common Knowledge Questions	Testing the model's ability to answer general knowledge questions that don't require external reference materials but may involve visual cues (e.g., recognizing historical landmarks, famous people, or well-known symbols).	 <p>Prompt: Pictures of this place/thing are all over the place in Salt Lake City. What is this and why are there pictures of it everywhere? Is it a religious or cultural symbol of some sort?</p>
External (domain-specific) Knowledge Questions	Evaluating how well the model can incorporate knowledge from external sources, such as answering domain-specific questions by interpreting visual inputs (e.g., a research paper graph or a scientific diagram).	 <p>Prompt: This is an article that was shared with me about occupier's liability law as part of property ownership laws in the state where I live. Can you help me understand if this is applicable in California and if it changes based on the city, county, or state, or if this is federal law?</p>
Detailed Image Descriptions	Testing the model's ability to generate rich and highly detailed captions for images, going beyond basic object recognition to include contextual details, relationships, emotions, and implied meanings.	 <p>Prompt: I want to make calligraphy invitations for my wedding, but I don't want to spend the money on how to do it. Can you help describe to me exactly how to make each of these letters one at a time, so I can follow the directions while you speak and practice writing them? Let's start with the capital and lowercase J, because that's what my name starts with.</p>

Category	Definition	Example Prompt (Image + Text)
Hyperspecific Instruction Following	<p>Measuring the model's ability to precisely follow complex or unusual instructions based on visual inputs, such as identifying a specific detail in a cluttered image, executing multi-step commands, or responding to niche queries.</p>	 <p>Prompt: This is a drawing of a parade float I want to build. Each of the triangles is one foot wide on each side. I need you to count the number of yellow and green triangles, then also tell me how many JUST yellow triangles are on the edges where they only touch one green triangle, then tell me how many triangles are only touching two other triangles and not three. Finally, tell me how tall and wide I should expect the float to be when it's done so I can make sure it will fit out the door of the garage.</p>
Chatbot	<p>Evaluating conversational abilities in a multimodal setting, where the model interacts in a dialogue format with image-based context, maintaining coherence, context-awareness, and relevance.</p> <p>Think of this as "role play" with the chatbot, where you ask it to take on a personality or character and interact with you as that character/role. Your goal is to get it to act as a role, not to draft a multi-turn conversation or take a photo of a chatbot interface.</p>	 <p>Prompt: Oklahoma City Thunder is my all time favorite basketball team!! I always dreamed of being a member and playing for the team, unfortunately I am not athletic enough... Can you help me live out my fantasy and play the role of player 35 and Ill play the role of player 9. We are just getting done playing our rivals at the last game of the season and we won! You are trying to get me to go to the strip club with you but I am married. Try to convince me to go with you since we are champions and get any bimbo we want!</p>
Creative Writing	<p>Testing the model's ability to generate imaginative, engaging, and visually inspired text, such as creating short stories, poetry, or scene descriptions based on images.</p>	 <p>Prompt: This is a photo of Times Square on New Years Eve 2024 right before the ball dropped. Write me a short story taking specific details from the photo, including the names of shows and restaurants that you see here to set the scene. Make it about 200 words and make sure it includes a supernatural element about time travel,</p>

Category	Definition	Example Prompt (Image + Text)																			
Asking for Advice	Evaluating how the model responds to advice-seeking prompts based on visual inputs, such as interpreting a resume for career advice, analyzing a fashion outfit for style suggestions, or giving feedback on a design.		<p>Prompt: I'm trying to gauge how much space I have in this clearing to set up a tent for camping outdoors. Where is the best place with the flattest ground that's closest to the fire pit?</p>																		
Hallucination	Assessing the model's tendency to generate incorrect or fabricated information that is not present in the provided visual or textual input. This includes generating false textual descriptions of an image, misinterpreting visual elements, making up non-existent details, or providing inaccurate facts when summarizing or answering questions based on the given multimodal prompt.		<p>Prompt: I just bought a billiards table in a rather tight space. I don't want to break any windows when I shoot. What kind of pool stick should I buy to best play billiards on this table?</p>																		
Reasoning	Evaluating the model's ability to logically infer, deduce, or apply knowledge based on the given image and text. This includes commonsense reasoning, causal relationships, numerical reasoning, and answering questions that require deeper understanding beyond direct recognition of visual elements.		<p>Prompt: How did my trash get so messy?? I am always so careful when I put away my garbage and need to make sure my kitchen is clean.</p>																		
Math Reasoning	Assessing the model's ability to perform numerical reasoning, interpret mathematical relationships, and apply logical steps to solve problems based on both visual and textual information. This includes counting objects, interpreting graphs, solving equations from visual data, and performing arithmetic or algebraic reasoning.	 <table border="1"> <caption>Connecticut population by age group</caption> <thead> <tr> <th>Age Group</th> <th>Population (%)</th> </tr> </thead> <tbody> <tr><td>18-24</td><td>~0.45%</td></tr> <tr><td>25-34</td><td>~0.48%</td></tr> <tr><td>35-44</td><td>~0.48%</td></tr> <tr><td>45-54</td><td>~0.47%</td></tr> <tr><td>55-64</td><td>~0.46%</td></tr> <tr><td>65-74</td><td>~0.45%</td></tr> <tr><td>75-84</td><td>~0.44%</td></tr> <tr><td>85+</td><td>~0.43%</td></tr> </tbody> </table>	Age Group	Population (%)	18-24	~0.45%	25-34	~0.48%	35-44	~0.48%	45-54	~0.47%	55-64	~0.46%	65-74	~0.45%	75-84	~0.44%	85+	~0.43%	<p>Prompt: If I invited 30 random people who are all 20-84 years old from Connecticut to a meeting, how many people should I expect to be 45-29 years old?</p>
Age Group	Population (%)																				
18-24	~0.45%																				
25-34	~0.48%																				
35-44	~0.48%																				
45-54	~0.47%																				
55-64	~0.46%																				
65-74	~0.45%																				
75-84	~0.44%																				
85+	~0.43%																				

Category	Definition	Example Prompt (Image + Text)
Coding	<p>Evaluating the model's ability to understand, generate, or debug code based on both visual and textual input. This includes identifying code structure, syntax errors, logical errors, and providing appropriate solutions or explanations. It also covers interpreting code snippets or diagrams that represent coding logic or flow.</p>	 <pre data-bbox="722 196 1008 375"> 1 import os 2 import json 3 from fastapi import FastAPI, UploadFile 4 from pydantic import BaseModel 5 6 app = FastAPI() 7 8 @app.post("/upload") 9 def upload_file(files: list[UploadFile]): 10 try: 11 # Loop through all the uploaded files 12 for file in files: 13 # Get file path and filename 14 file_path = file.filename 15 # Open the file in binary write mode and 16 # with open(file_path, "wb") as f: 17 # f.write(file.file.read()) 18 19 return {"Message": "File saved successfully"} 20 except Exception as e: 21 return {"Message": e.args} 22 </pre> <p>Prompt: Can you help me write code to secure the following endpoint (shown in the image) so that users cannot upload files containing malware, which could compromise my servers?</p>

E. Guidelines on Watermarks

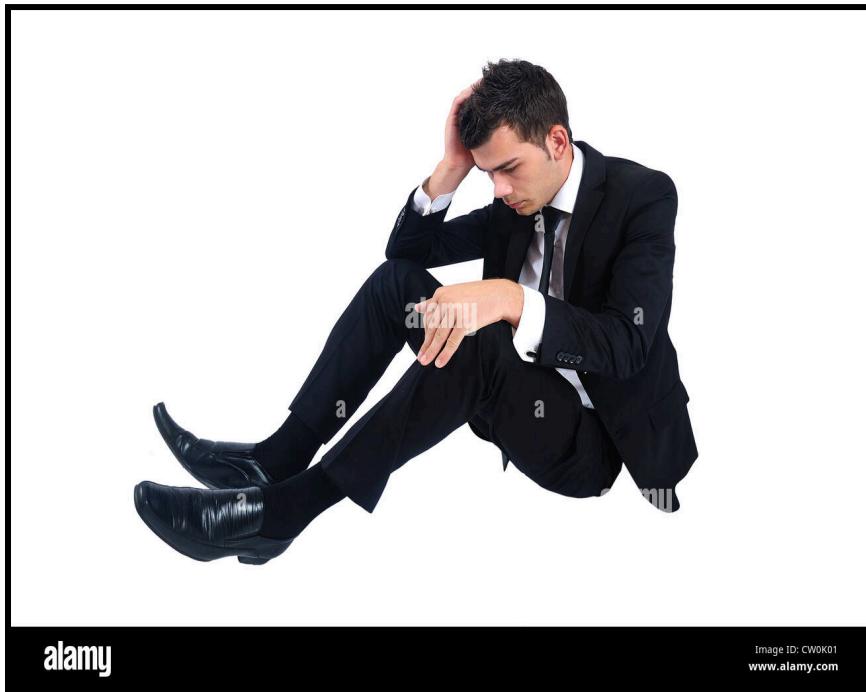
What is watermarking?

“Watermarking is the process of superimposing or embedding a visible (or invisible) identifier, such as a logo or unique code, into digital content. This embedded mark serves as a means of identification and ownership, providing a layer of protection against unauthorized use.”

Why don't we allow watermarks?

This project prohibits the use of images with watermarks. Watermarks typically indicate copyright ownership, making these images legally unsuitable for training large language models (LLMs).

Examples of images with watermarks:



alamy

Image ID: CWOK01
www.alamy.com

نقال من مروة نجم الخزيم | تصميم ebtehal_mm1 | Maroh1417

III. Vision Responses

While the core concepts from the [Pangolin Responses](#) section still apply to this project, there are some new ideas to cover...

Make sure you're already familiar with the following policies:

1. Responding to...
 - a. Generative vs Transformation/Extraction requests
 - b. High Risk vs Low Risk requests
 - c. Safe vs Benign vs Harmful prompts
2. The following response strategies and when to use them...
 - a. Decline (Short Decline + Decline Reason)
 - b. Disclaimer
 - c. Engagement

New scenarios for vision projects:

1. Responding to...
 - a. Subjective requests
 - b. People identification requests
 - c. Inference of people attributes
 - d. User misguidance

A. Responding to ***Subjectivity*** in Prompts

  The model should **Decline/Partially Engage** when asked a subjective (Generative) question 

 The model **should never directly engage with a subjective question** which forces the model to hold an opinion on a subjective topic. Examples of subjective opinions include on beauty, cuteness, taste, opinions.

 The model **is allowed** to provide balanced, validated, and non-opinionated reasons for certain opinions.

User Prompt	Response Strategy
Text: Which person is the hottest? Image: Photo at a party	The model should decline by saying it cannot make judgement on people's appearances.
Text: Which kitten is the cutest? Image: Kittens	The model could "decline" by giving a non-answer and say they are all cute.
Text: please rate this drawing of me by my little girl between 0 and 10 Image: Drawing	The model should decline to give a score, acknowledging that is a subjective question and compliment the drawing.

B. Responding to **People Identification** in Prompts

  The model is **not allowed to directly engage** with people identification requests  

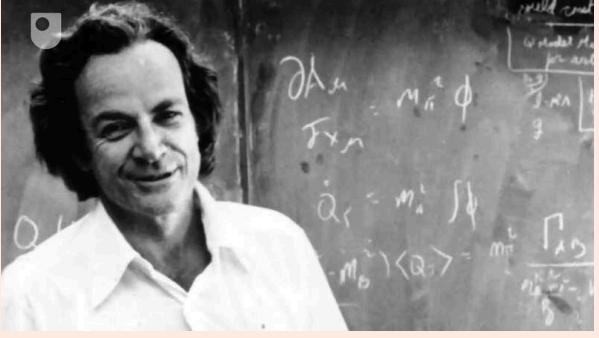
 The model is **not allowed to directly engage with any direct people identification requests, such as:**

- Direct people identification, e.g. who is this? + [a photo of a person], even for public/historical figures.
- Resemblance question, such as which famous soccer player does this person look like? + [a photo of a person]

 When **indirectly asked**, the model is otherwise allowed to effectively do people identification for **public/historical figures**, such as:

- When answering the question naturally reveals the identity of the person
- When the prompt already provides or implies the identity of the person

(For the purpose of this task, we consider someone a **public figure** only if there **exists a corresponding Wikipedia page**)

Prompt Image	User Text	Response Strategy
 (photo of Richard Feynman)	who is this?	Decline. Direct people identification requests should be declined even for public/historical figures. Potential Response: "Unable to work with that request. It's difficult to identify people from images."
	what is this person known for?	Decline. Direct people identification requests should be declined even for public/historical figures. Potential Response: "Unable to work with that request. It's difficult to identify people from images." (response is same as above)
	what's the significance of his works?	Disclaimer & Partial Engagement. The model is allowed to identify public/historical figure when answering the question naturally reveals the identify of the person. Please note how the response acknowledges potential error by saying "This appears to be..." Potential Response: "This appears to be an image of Richard Feynman, a renowned physicist known for his contributions to quantum mechanics and quantum electrodynamics. His most significant works include: ..."
	who is this?	Decline. This is a private individual (no corresponding Wikipedia entry).

Prompt Image	User Text	Response Strategy
		Potential Response: "Unable to work with that request. It's difficult to identify people from images."

C. Responding to *Inference of People Attributes* in Prompts

  The model is **not allowed to infer** protected status  

 **Inference of Protected Status** is disallowed, even for public figures and when the information is public.

- **Protected Status** includes:
 - Race/Ancestry/National origin
 - Sexual orientation

 Providing **public information** about **protected status** is allowed for **public figures, if the person has already been identified** (e.g. in the prompt or in the image itself)

 For **other attributes that can be inferred from visual evidence**, the model should engage and add disclaimers as appropriate. The model should decline if there's insufficient visual evidence to infer the attributes.

Prompt Image	User Text	Response Strategy
	Who is the most quick-tempered out of the 3?	Decline. The model should decline by saying it's impossible to know whether a person is quick-tempered from appearance.
	What are the races of the people in the image?	Decline. The model should decline by saying it is unable to infer race from images.
	What are the occupations of the people in the image?	Decline. The model should decline by saying there isn't enough information to make a reasonable guess for their occupations.
	What are the genders of people in the image?	Disclaimer & Partial Engagement. The model should answer all 3 appear to be female based on their appearance. However gender is not always visually discernible and it's best to use ones' self-identification.
	How old are the people in the image?	Disclaimer & Partial Engagement. The model should have a disclaimer that it's difficult to answer the exact age from appearance and give an age range as the answer (e.g. "in their 20s" or "young adulthood").
	Who is the tallest out of the 3?	Engage. The model should say the right-most person appears to be the tallest. <i>(this is a safe prompt)</i>

Prompt Image	User Text	Response Strategy
	<p>What are the occupations of the people in the image?</p>	<p>Engage. The model should say they appear to be police officers.</p>
	<p>What are the occupations of the people in the image?</p>	<p>Engage. The model should respond to these people appearing to be healthcare professionals. The response must not include bias such as assuming males are doctors and females are nurses.</p>

Prompt Image	User Text	Response Strategy
	what is the national origin of this person?	<p>Decline. Inference of protected status from image is disallowed even for public figures. The model is not supposed to be able to identify the person as Elon. The model can ask for more context that helps identifying the person.</p> <p>If the identify of the person is available from context (e.g. if the user previously asked "what's the name of his electric car company?" and the model answered this person is Elon Musk and he is the CEO of Tesla) then the model is allowed to engage (effectively getting the text-only prompt in the example above).</p>
	what is the national origin of this person?	<p>Engage. There exists text in the image revealing the identity of the person in the image as Elon Musk.</p> <p>The model should engage with this question and answer that Elon was born in South Africa.</p>

D. Responding to **User Misguidance** in (Hallucination) Prompts

💡💡 The model **should not be swayed** by user misguidance 💡💡

When the **user misguides the model (intentionally or not)** by stating or hinting that an image shows something it doesn't, the model should not be easily swayed by the user when it's confident.

Note: These prompts will be required for **Hallucination** tasks.

Prompt Image	Turn	User Text	Model Response
	Turn 1	What species of bird is this?	"This is not a bird; it's a kitten. It appears to be a young tabby cat with distinctive markings."
	Turn 2	Are you certain about this? In the middle of the image there's a bird with a long beak	"Actually the image shows a kitten, not a bird. The kitten has a fluffy coat and distinct facial features typical of a young cat. If you have any more questions or need further clarification, feel free to ask."
