

# Image Analysis for AI-Powered Post Creation

---

## Basic Image Analysis Techniques

Image analysis encompasses several key techniques:

1. **Image Classification:** Identifies what a whole image represents (e.g., "This is a restaurant interior")
2. **Object Detection:** Locates and identifies multiple objects within a single image (e.g., "There's a steak, a glass of wine, and a salad")
3. **Semantic Segmentation:** Classifies each pixel in an image to a specific class (e.g., precisely outlining the boundaries of all objects)
4. **Instance Segmentation:** Combines object detection and semantic segmentation to identify separate instances of objects (e.g., distinguishing between multiple steaks on a table)

For our project, **object detection** is indeed the most appropriate approach, as we need to identify multiple food items, restaurant decor elements, and other objects that would be relevant for caption generation.

## Role of Image Analysis in Our Project

Image analysis serves a critical function in the AI-powered post creation workflow:

1. The system analyzes uploaded images to detect key objects
2. These detected objects are presented to the user
3. Users can select relevant items to focus on while removing irrelevant or incorrectly detected objects
4. The selected objects provide context for generating accurate, relevant captions that match the image content
5. This improves the quality of AI-generated captions by ensuring they reference what's actually in the image

This approach gives users control over what aspects of the image to emphasize in captions while leveraging AI to identify content automatically.

## Implementation Options

Table 1: Implementation Options Comparison

Option	Approach	Pros	Cons	Feasibility

Option	Approach	Pros	Cons	Feasibility
1	Run pre-trained models (YOLO) locally on developer machine	<ul style="list-style-type: none"> <li>• Zero hosting costs</li> <li>• Simplest initial setup</li> <li>• Complete privacy</li> <li>• No internet dependency</li> <li>• Good for proof-of-concept</li> </ul>	<ul style="list-style-type: none"> <li>• Only runs on developer's machine</li> <li>• Team members can't test it</li> <li>• Requires GPU for good performance</li> <li>• Not usable in production</li> </ul>	Low - Limited to single developer testing
2	Deploy entire backend with pre-trained models (YOLO) on cloud servers (AWS EC2, Google Colab)	<ul style="list-style-type: none"> <li>• Complete backend solution</li> <li>• Integrated with all other features</li> <li>• Production-ready environment</li> <li>• Centralized deployment</li> </ul>	<ul style="list-style-type: none"> <li>• Highest cost (full server)</li> <li>• Most complex setup</li> <li>• GPU costs even when not using ML</li> <li>• Requires server management</li> <li>• Overkill for small project</li> </ul>	Low - Excessive overhead for image analysis alone
3	Self-host pre-trained models (YOLO) as API on cloud servers (AWS EC2, Google Colab)	<ul style="list-style-type: none"> <li>• Dedicated ML microservice</li> <li>• Only pay for ML processing</li> <li>• Separates ML from main backend</li> <li>• Reusable across projects</li> <li>• Can optimize server for ML</li> </ul>	<ul style="list-style-type: none"> <li>• Still requires separate GPU server</li> <li>• Deployment complexity</li> <li>• API design and implementation needed</li> <li>• Time-consuming for a small feature</li> </ul>	Medium - Requires significant initial setup

Option	Approach	Pros	Cons	Feasibility
4	Use free third-party image analysis APIs	<ul style="list-style-type: none"> <li>No setup required</li> <li>Free to use</li> <li>Simple API integration</li> </ul>	<ul style="list-style-type: none"> <li>Poor detection quality (see Table 2)</li> </ul>	Low - Inadequate quality (tested)
5	Use paid commercial vision APIs (AWS Rekognition, Google Vision API)	<ul style="list-style-type: none"> <li>Best accuracy for food items</li> <li>No server setup or maintenance</li> <li>Highly reliable services</li> </ul>	<ul style="list-style-type: none"> <li>Usage costs (~\$1-2 per 1,000 images)</li> <li>API key management</li> <li>Data privacy concerns</li> </ul>	High - Best balance of effort vs quality

Table 2: Free API Test Results

Image	Expected Detection	Cloudmersive Results
	Steak, Red wine	cup
	Pasta, Shrimp, Garlic, Herbs, Clam, White wine	wine glass, chair
	Burger	sandwich

Image	Expected Detection	Cloudmersive Results
	Pancake, Berries, Whipped Cream	cake, bowl

## Conclusions

**Option 5 (Paid Commercial APIs)** offers the best balance between implementation effort and quality outcomes. Given the scope of our project, the associated costs would be minimal while delivering the necessary food detection capabilities for generating high-quality captions.

Further research on AI Vision APIs

API Service	Pricing	Free Tier	Food Detection Capabilities	Other Benefits
<b>AWS Rekognition</b>	\$1 per 1,000 images	5,000 images/month (first 12 months)	Restaurant-specific object detection with detailed labels	Robust food and ingredient recognition; Extensive documentation
<b>Google Cloud Vision</b>	\$1.50 per 1,000 images	1,000 units/month	Strong food recognition with ingredient breakdown	Food label detection; Integration with other Google services
<b>BLIP (Replicate)</b>	~\$0.73 per 1,000 images	~100 calls per day	Combined object detection and caption generation	Generates descriptive captions; Single API call for both detection and description

## Next steps

1. Research and compare sample responses from AWS Rekognition, Google Vision API, and BLIP for restaurant food images
2. Set up developer accounts and obtain API keys for the selected service
3. Create a simple API wrapper in our backend to handle image upload and analysis
4. Implement the user interface for displaying and selecting detected objects before caption generation
5. Test the integration with actual restaurant food images to verify detection quality