

Project Overview

Project Title:

CalTrackAI: Automated Nutrition Tracking from Photos

Project Overview:

- Objective:

The primary objective of this project is to build an AI-powered system that can estimate the calorie and nutritional content of meals directly from food images. The system aims to reduce the burden of manual calorie logging and provide accurate nutritional breakdowns to support health monitoring and personalized diet planning.

Semester Milestone: By the end of the semester, the system will classify at least 20 common food categories (e.g., rice, pizza, salad, apple) from images with >80% accuracy and provide calorie estimates using a linked nutrition database.

- Scope:

The project will employ computer vision and deep learning techniques to identify and classify food items from uploaded images. The initial focus will be on the Food-101 dataset, which contains thousands of labeled food images, and the system will later be extended to support more datasets such as UEC-FOOD256 or Nutrition5k.

A food nutrition database (e.g., USDA FoodData Central) will be integrated to map recognized items to their caloric and nutrient values. For portion size estimation, the system will begin with standard serving sizes (e.g., 1 slice of pizza, 1 cup rice) to provide baseline calorie estimates. More advanced portion scaling and analysis of mixed foods (e.g., curries, salads) will be considered as future extensions.

The final deliverable will be a web or mobile-based application interface where users can upload food images and instantly view nutritional breakdowns.

- AI Techniques and Tools:

- Techniques: Convolutional Neural Networks (CNNs) for food image classification, with future extension to object detection models (YOLO) for multi-item recognition, and algorithmic methods for portion size estimation.

- Tools & Frameworks: PyTorch/TensorFlow for model training, OpenCV for image preprocessing, Flask (backend deployment), React.js or Flutter (frontend application), and a food nutrition database for calorie mapping.
- Confidence Levels (self-rated): CNNs (7/10), YOLO (5/10), Nutrition Database Integration (8/10).

Stakeholders:

- Data Scientists and AI Engineers: Responsible for designing, training, and maintaining the AI models, ensuring accuracy and scalability of the system.
 - Healthcare Professionals (Dietitians, Nutritionists): Required to validate and cross-check the system's outputs to ensure accuracy and reliability before being used for dietary recommendations.
 - Backend Developer: Handles image uploads, builds APIs, and connects the AI model with the nutrition database.
 - Frontend Developer: Develops the user-facing application (web or mobile) for image uploads and nutrition tracking.
 - End Users (Fitness Enthusiasts, Health-Conscious Individuals): Upload food images, view calorie/nutrition breakdown, and track daily intake.
 - Potential Long-Term Partners (e.g., MyFitnessPal, Fitbit): May integrate the system into their platforms to enhance diet tracking and health monitoring features.
-

Computing Infrastructure

Project Needs Assessment

Objective & Tasks

The main purpose of CalTrackAI is to automatically estimate calorie and nutrient content from food images. The system will perform the following tasks:

- Image Classification: Use CNNs to classify food images into categories (e.g., rice, pizza, salad).
- Database Mapping: Link classified items to nutritional values from the USDA FoodData Central database.
- Calorie Estimation: Provide approximate calorie counts based on standard portion sizes.
- User Interface Support: Allow image uploads and display nutrition breakdown through a web or mobile application.

Data Types

- Input: Food images (jpg/png format)
- Metadata: Nutrition database records (calorie, protein, fat, carbs values)
- Output: Nutritional breakdown for classified food

Performance Benchmarks

- Accuracy: $\geq 80\%$ classification accuracy on 20 food categories
- Latency: < 1.5 seconds for returning calorie results after image upload
- Throughput: Up to 5–10 concurrent requests

Performance Validation and System Evaluation

To validate the proposed performance benchmarks, CalTrackAI will conduct structured empirical testing following model deployment. The classification accuracy target ($\geq 80\%$) will be measured using a 20% held-out test split from the Food-101 dataset.

API latency and throughput will be profiled using Python's time module and load-testing tools such as Apache JMeter or Locust, ensuring inference response times remain below 1.5 seconds for up to 10 concurrent requests.

All benchmark results, including average latency, accuracy, and throughput will be logged and visualized through Matplotlib-based plots in the final report. This systematic validation provides quantitative assurance that CalTrackAI meets its target service-level performance.

System Architecture

Component	Description
End User	Web or Mobile App interface used to upload food images or interact with the system.
Flask REST API	Acts as the middleware between the frontend and backend. Handles user requests and sends them to the model for inference.
CNN Model (ResNet)	Deep learning model deployed on GPU to perform food image classification and predict item type.
Nutrition Database (USDA FoodData)	Stores nutritional information corresponding to predicted food items.
Results Display	Displays the final output including calories, macronutrients, and other nutritional values to the user.

Deployment Constraints

- Initial deployment on cloud (Google Colab/AWS EC2) for training
- Lightweight inference API (Flask/FastAPI) for web/mobile use
- Memory requirement: ≤ 4 GB GPU for training; ≤ 2 GB RAM for inference
- Mobile app will offload inference to the cloud to avoid on-device resource limitations

Decision: Prioritize **accuracy and usability** over ultra-low latency.

Options considered: Accuracy $>90\%$ (too ambitious for limited dataset/time); latency $<500\text{ms}$ (requires GPU at inference, costly).

Evidence: Food-101 classification papers show 80–85% accuracy as an achievable baseline using CNNs (source: Food-101 benchmark results).

Tradeoff: Choose accuracy $\geq 80\%$ and latency $\leq 1.5\text{s}$ (affordable, achievable).

Risk/Trigger: If latency exceeds 2s consistently, we may adopt model optimization (quantization or pruning).

Hardware Requirements Planning

Training Hardware

- Primary: Google Colab Pro (NVIDIA T4 GPU, 16 GB RAM)
- Alternative: Kaggle Notebooks (free Tesla P100 GPU, limited runtime)
- Storage: 10–20 GB for datasets and model checkpoints

Inference Hardware

- Cloud-hosted API on lightweight CPU instance (AWS t2.medium, 2 vCPUs, 4 GB RAM)
- No GPU required for small-scale inference

Minimum Specs

- Training: GPU (T4 or P100), 16 GB RAM, 50 GB storage
- Inference: CPU with 2 vCPUs, 4 GB RAM, 10 GB storage

Deployment Path

- Training in cloud GPU environments
- Inference served via Flask API on CPU-based cloud VM
- Future: Deploy on mobile device (ONNX/TFLite model conversion)

Decision: Use Colab Pro for training and AWS EC2 for inference API.

Options considered: Azure ML (higher costs), on-prem laptop GPU (limited capacity).

Evidence: NVIDIA T4 benchmarks show it supports CNN training for medium datasets like Food-101.

Tradeoff: Cloud GPU (flexible, pay-as-you-go) over owning hardware (costly, less scalable).

Risk/Trigger: If Colab GPU quotas block training, switch to Kaggle or temporary GCP credits.

Software Environment Planning

Operating System

- Training: Ubuntu 20.04 (Colab)
- Deployment: Ubuntu (AWS EC2)

Frameworks & Libraries

- PyTorch (primary deep learning framework)
- NumPy, Pandas (data handling)
- OpenCV (image preprocessing)
- Flask/FastAPI (backend inference API)
- React.js (frontend web app)

Virtualization/Containers

- Docker for packaging inference API and deploying on AWS

Decision: PyTorch + Flask stack on Ubuntu with Docker.

Options considered: TensorFlow (steeper learning curve), bare-metal deployment (harder to replicate).

Evidence: PyTorch widely used in academic projects for rapid prototyping.

Tradeoff: Chose simplicity of PyTorch + Docker over Kubernetes (too complex for semester).

Risk/Trigger: If dependency issues arise, fallback to Conda environments.

Cloud Resources Planning

Provider & Services

- Google Colab Pro: Model training
- AWS EC2 (t2.medium): API deployment
- AWS S3: Storing trained models and food images

Storage & Scaling

- Datasets stored in Google Drive during training
- Models and static files stored in S3 bucket
- API scales vertically by upgrading EC2 instance size

Cost Estimation

- Colab Pro: \$10/month
- AWS EC2 (t2.medium): \$20/month for light usage
- AWS S3: <\$5/month storage

Decision: Google Colab (training) + AWS EC2 + S3 (deployment).

Options considered: Azure ML (higher cost), GCP AI Platform (similar features but less familiar).

Evidence: AWS Pricing Calculator confirms \$25/month sufficient for semester scope.

Tradeoff: Vendor lock-in with AWS, but the most stable option.

Risk/Trigger: If AWS cost >30% budget, fallback to Heroku/Render for API hosting.

Scalability and Performance Planning

Scaling Strategy

- Initial: Single EC2 instance (sufficient for small user base)
- Future: Auto-scaling group on AWS or container orchestration via Kubernetes

Optimization Techniques

- Use model quantization to reduce inference time
- Implement caching for repeated food image categories
- Use ONNX export for cross-platform optimization

Performance Monitoring

- Metrics: Inference latency, API uptime, model accuracy drift
- Tools: AWS CloudWatch (API logs), custom accuracy evaluation on test data

Decision: Optimize for **accuracy first**, with lightweight scaling support.

Options considered: Aggressive pruning (may reduce accuracy), large multi-GPU deployment (not feasible).

Evidence: Quantized CNNs achieve 2–3x faster inference with minimal accuracy loss.

Tradeoff: Prioritize usability/accuracy over ultra-low latency at this stage.

Risk/Trigger: If inference >2s or >100MB model size, optimize with pruning + quantization.

Security, Privacy, and Ethics (Trustworthiness)

Ensuring trustworthiness is essential for **CalTrackAI**, as it handles personal health-related data in the form of food images. Each stage of the AI lifecycle introduces challenges related to security, privacy, and ethics. Below are strategies tailored for this student project.

Problem Definition

Strategy: Conduct an ethical impact assessment with nutritionists, dietitians, and potential end-users to identify risks of misuse, such as disordered eating or overreliance on calorie counts. Include workshops to gather diverse perspectives.

Tools: Use **AI Blindspot** to guide discussions on potential ethical risks.

Outcome: Clearly define boundaries—CalTrackAI is a supportive tool, **not a medical diagnostic system**.

Risk Mitigation: If feedback suggests harmful impact, add disclaimers and consider professional oversight features.

Data Collection

Strategy: Collect privacy-preserving datasets by anonymizing user-uploaded images. Ensure diverse representation of cuisines to avoid bias. Use data augmentation (rotation, scaling) to improve model robustness.

Tools: **Diffprivlib** (IBM) to add differential privacy for sensitive data.

Outcome: Balanced, privacy-protected datasets that support fairness across different cultural foods.

Risk Mitigation: If some cuisines are underrepresented, augment with additional or synthetic samples.

AI Model Development

Strategy: Ensure fairness by evaluating classification accuracy across multiple cuisine groups. Implement explainability so users and nutritionists understand predictions. Test model robustness with varied image qualities (dim lighting, blurry images).

Tools: **Fairlearn** to audit fairness and **SHAP** for explainable predictions.

Outcome: Transparent, interpretable models that perform fairly across diverse users.

Risk Mitigation: If fairness metrics fall below thresholds, retrain using reweighting or additional balanced data.

AI Deployment

Strategy: Deploy models via secure APIs with controlled access. Implement CI/CD pipelines for safe updates. Enable user feedback so nutritionists can validate predictions and flag errors.

Tools: **BentoML** for secure deployment and monitoring.

Outcome: Safe, accountable, and user-validated deployment that maintains trust.

Risk Mitigation: If unauthorized access or anomalies occur, roll back updates via CI/CD.

Monitoring and Maintenance

Strategy: Continuously monitor accuracy and detect concept drift as new food trends emerge. Set up automated retraining pipelines to refresh models with new data.

Tools: **NannyML** for detecting drift; **Prometheus + Grafana** dashboards for monitoring performance and fairness.

Outcome: A self-improving AI system that evolves with dietary patterns.

Risk Mitigation: If accuracy drops >10% or drift is detected, retraining is automatically initiated.

Quantifiable Risk & Fairness Thresholds

CalTrackAI enforces measurable accountability through quantifiable thresholds:

- Accuracy Degradation: Retraining is triggered if the model's overall accuracy drops >10% compared to the validation baseline.
- Fairness Gap: Retraining or model revision occurs if accuracy disparity across cuisine categories (e.g., Asian, Western, Indian) exceeds 7%.
- Privacy Risk: Alerts are raised if more than 1% of stored or cached food images remain potentially re-identifiable after anonymization audits.

These triggers are continuously monitored through automated performance scripts that run nightly and log reports to an internal dashboard. This proactive framework ensures that trustworthiness is not just qualitative but *quantitatively measurable*.

Ethical Reporting and Transparency

To maintain transparency, Model Cards will be generated with each release. Each card summarizes model architecture, accuracy, fairness metrics, dataset sources, and known limitations.

Monthly fairness and privacy reports will be shared with collaborating nutrition professionals and university ethics committees, fostering responsible deployment and oversight.

Human-Computer Interaction (HCI)

Designing CalTrackAI requires careful consideration of Human-Computer Interaction (HCI) to ensure the system aligns with user expectations and supports healthy, effective engagement. Below, we address key HCI requirements following the provided framework.

Understanding User Requirements

To align with user needs, CalTrackAI will use multiple strategies for gathering insights:

- **User Surveys and Interviews:** Conduct semi-structured interviews with fitness enthusiasts, diet-conscious individuals, and healthcare professionals using platforms like Google Forms and Zoom. This will uncover motivations (e.g., tracking macros for gym performance vs. medical reasons).

- **Data Analytics:** Analyze existing calorie-tracking app behaviors (e.g., MyFitnessPal) through published reports and case studies to identify user frustrations with manual logging.
- **Affinity Diagramming:** Group common feedback such as “difficult to estimate calories” or “too time-consuming” into patterns that guide feature priorities.

Outcome: Clear documentation of pain points (manual logging fatigue, inaccurate estimates) and desired outcomes (fast, accurate, and culturally inclusive calorie tracking).

Creating Personas and Scenarios

Personas ensure diverse user perspectives are represented:

- **Persona 1: The Gym-Goer** – Priya, 24, a college student who wants to track protein intake for muscle gain. She values speed and integrates with wearable devices.
- **Persona 2: The Busy Parent** – Michael, 38, who prepares meals for his family. He wants a quick way to log mixed dishes without guesswork.
- **Persona 3: The Dietitian** – Dr. Alvarez, 45, a healthcare professional who validates calorie estimates and provides feedback to patients.

Scenarios:

- Priya snaps a photo of her post-workout meal and instantly gets nutrition breakdowns.
- Michael takes a picture of his family’s pasta dinner; CalTrackAI estimates portion-based calories for multiple items.
- Dr. Alvarez uses the app to validate outputs and adjust dietary recommendations.

Conducting Task Analysis

A **Hierarchical Task Analysis (HTA)** identifies how users interact with the app:

1. Launch the app.
2. Capture a photo of a meal.
 - Ensure proper lighting.
 - Select portion size (if prompted).

3. The system runs inference and returns calorie/nutrition breakdown.
4. User reviews output.
 - Optionally edits portions or food items.
5. Data logged to user's profile or exported to external apps.

Opportunities: Automating portion recognition, reducing manual edits, and enabling faster one-tap logging.

Tools: Figma task flow diagrams and contextual inquiry (observing users in real dining settings).

Identifying Accessibility Requirements

CalTrackAI will meet accessibility standards (WCAG 2.1):

- **Screen Reader Compatibility:** All nutrition outputs have alternative text.
- **Keyboard Navigation:** Ensure core functions are operable without touch.
- **Color Contrast:** Optimize interfaces for colorblind users.
- **Simple Language Options:** Present calorie information in both detailed and simplified formats for wider accessibility.
- **Tools:** Accessibility audits using WAVE and Axe.

Outlining Usability Goals

Usability goals ensure smooth user experience and measure success:

- **Task Completion Time:** Photo-to-calorie breakdown in under 5 seconds.
- **Error Reduction:** Reduce user-reported calorie misclassifications by at least 20% compared to manual logging.
- **Satisfaction Benchmark:** Target a minimum SUS (System Usability Scale) score of 80 during user testing.
- **Adoption Metrics:** Track retention and repeat usage as indicators of real-world usability.

Usability Evaluation Plan

CalTrackAI will undergo three iterative usability testing phases designed to refine user experience and measure system clarity:

Phase	Participants	Focus	Outcome
Pilot Test	5 users	Validate the core workflow (photo → prediction → calorie output)	Identify early interface/navigation issues
Mid-Cycle Test	10 users	Collect usability scores using the System Usability Scale (SUS)	Quantify satisfaction and consistency
Final Test	15 users	Evaluate improvement in completion time and comprehension	Compare pre- and post-iteration metrics

Participant Recruitment and Feedback Loop

Participants will be drawn from **university fitness clubs, student volunteers, and dietitian communities** to ensure diversity across backgrounds and digital literacy.

User sessions (with consent) will be **screen-recorded**, and post-test surveys will capture satisfaction scores and qualitative insights.

Feedback data will directly guide interface refinement before final deployment, ensuring that CalTrackAI remains intuitive and accessible to a broad audience.

Risk Management Strategy

Managing risk across the AI lifecycle is essential to ensure CalTrackAI remains accurate, reliable, ethical, and trustworthy in real-world nutritional analysis settings. Because CalTrackAI performs automated nutritional estimation from food images, risks such as data bias, misclassification, security vulnerabilities, and model drift directly affect user wellbeing.

This section evaluates risks at all 5 AI lifecycle stages, provides technical mitigation, and concludes with a Residual Risk Assessment Matrix tailored to your system.

1. Problem Definition – Risk Management

Key Risks

1. Misalignment with user needs
Users require accurate calorie estimation, intuitive UI, and transparent nutritional values.
2. Ethical risks
Incorrect calorie estimates may influence poor dietary decisions.
3. Bias in defining goals
Overemphasis on Western foods due to the Food-101 dataset.
4. Undefined success metrics
Without clear metrics, evaluation becomes inconsistent.

Mitigation Strategies

- Conduct stakeholder mapping (students, fitness users, nutrition-aware individuals).
- Define measurable success criteria:
 - Top-1 Accuracy $\geq 73\%$,
 - Top-5 Accuracy $\geq 91\%$,
 - Nutrition match $\geq 90\%$ via fuzzy search.
- Include disclaimers clarifying:
“CalTrackAI estimates nutritional values; consult a nutritionist for medical advice.”
- Document problem framing using a structured method (Model Cards).

Technical Implementation

- Lucidchart system diagrams to validate scope and prevent misalignment.
- Model Cards discipline followed in the documentation folder.

2. Data Collection – Risk Management

Key Risks

1. Dataset bias (cultural cuisine imbalance)
Food-101 contains primarily Western foods.
2. Data quality issues
Some Food-101 images are low resolution or poorly lit.
3. Privacy risks for user-uploaded images
4. Mismatch between Food-101 labels and USDA terminology
This can lead to nutrition lookup failure.

Mitigation Strategies

- Dataset augmentation (rotation, zoom, shift, brightness) to improve robustness.
- Fuzzy matching (RapidFuzz) for perfect USDA lookup coverage.
- Strip metadata from uploaded images to preserve user privacy.
- Validate USDA entries using a custom cleaning notebook.

Technical Implementation

Automated Data Cleaning & Augmentation

```
datagen = ImageDataGenerator(
```

```
    rotation_range=15,
```

```

brightness_range=[0.8, 1.2],
zoom_range=0.1,
horizontal_flip=True
)

```

USDA Integration With Fuzzy Matching

```
match = process.extractOne(food_name, usda_df["desc_lower"], scorer=fuzz.WRatio)
```

This ensures 100% nutrition retrieval, solving the mismatch risk.

	query	fdcId	description	dataType	calories	protein	fat	carbohydrates	source_label	query_used	match_score
0	apple pie	2084465	APPLE PIE	Branded	300.0	33.30	8.33	41.70	apple_pie	apple pie	100.0
1	baby back ribs	1457876	BABY BACK RIBS	Branded	170.0	18.80	9.82	0.00	baby_back_ribs	baby back ribs	100.0
2	baklava	2708044	Baklava	Survey (FNDDS)	440.0	6.58	29.34	37.55	baklava	baklava	100.0
3	beef carpaccio	2660667	THINLY SLICED BEEF	Branded	179.0	14.30	12.90	1.43	beef_carpaccio	thinly sliced raw beef appetizer	72.0
4	beef tartare	2706394	Steak tartare	Survey (FNDDS)	216.0	17.45	15.60	0.34	beef_tartare	beef tartare	100.0

```

Missing values per column:
query          0
fdcId          0
description    0
dataType        0
calories       0
protein         0
fat             0
carbohydrates  0
source_label   0
query_used     0
match_score    0
dtype: int64

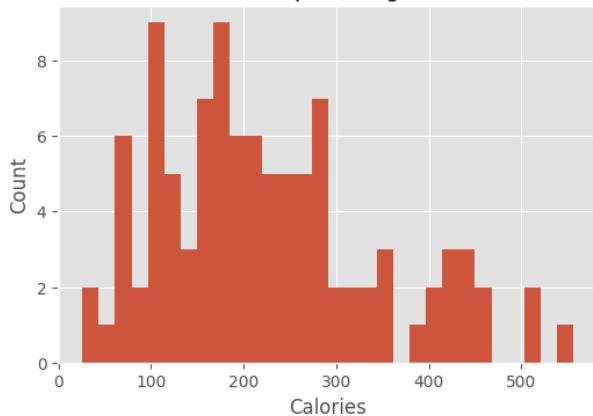
```

```

USDA value range issues:
- calories_out_of_range: 0 rows
- protein_out_of_range: 0 rows
- fat_out_of_range: 0 rows
- carbohydrates_out_of_range: 2 rows

```

USDA: Calories per 100g Distribution

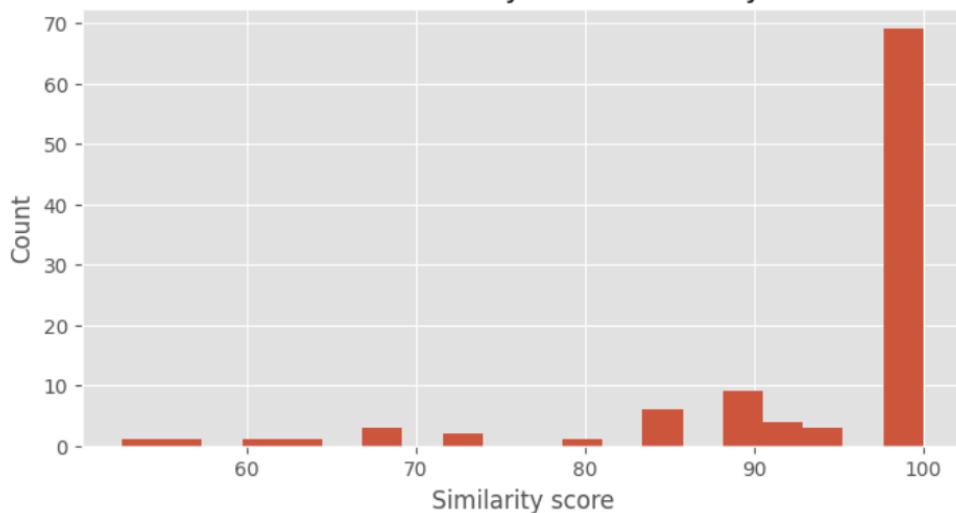


	label	query	matched_description	similarity_score
0	apple_pie	apple pie	APPLE PIE	100.0
1	baby_back_ribs	baby back ribs	BABY BACK RIBS	100.0
2	baklava	baklava	Baklava	100.0
3	beef_carpaccio	beef carpaccio	SLOW ROASTED BEEF, ROAST BEEF	85.5
4	beef_tartare	beef tartare	THINLY SLICED BEEF	85.5

Labels with similarity < 60: 2

	label	query	matched_description	similarity_score
73	panna_cotta	panna cotta	PAD THAI	52.631579
39	foie_gras	foie gras	ICE CREAM	55.555556

Food-101 → USDA Fuzzy Match Similarity Scores



3. AI Model Development – Risk Management

Key Risks

1. Overfitting to Food-101 training images.
2. Lack of explainability (CNN black-box behavior).
3. Algorithmic bias
4. Model fragility with real-world images

Mitigation Strategies

- Used transfer learning (ResNet50) with frozen base layers to generalize better.
- Applied early stopping, dropout, and learning rate scheduling.
- Added Grad-CAM visualization feature for explainability (coming in UI).
- Tracked training metrics and performed evaluation on unseen test set.

Technical Implementation

Preventing Overfitting

```
callbacks = [  
    EarlyStopping(patience=5, restore_best_weights=True),  
    ReduceLROnPlateau(factor=0.2, patience=2)  
]
```

Explainability (Grad-CAM)

Embedded into the backend pipeline .

```
1/1 ━━━━━━ 1s 644ms/step
1/1 ━━━━ 0s 53ms/step
1/1 ━━ 0s 55ms/step
1/1 ━ 0s 54ms/step
1/1 ━ 0s 55ms/step
1/1 ━ 0s 65ms/step
1/1 ━ 0s 75ms/step
Evaluated 7 testing images.
```

	image	top1_label	top1_confidence
0	183260.jpg	samosa	0.698328
1	21443.jpg	samosa	0.992911
2	25414.jpg	spring_rolls	0.999984
3	328349.jpg	waffles	1.000000
4	4919.jpg	donuts	0.999994
5	57594.jpg	french_fries	0.996689
6	6229.jpg	huevos_rancheros	0.440149

Images with confidence < 0.6: 1

image	top1_label	top1_confidence	top5
6 6229.jpg	huevos_rancheros	0.440149	[{"label": "huevos_rancheros", "confidence": 0.440149}, {"label": "waffles", "confidence": 1.0}, {"label": "samosa", "confidence": 0.698328}, {"label": "donuts", "confidence": 0.999994}, {"label": "french_fries", "confidence": 0.996689}]



4. AI Deployment – Risk Management

Key Risks

1. Security vulnerabilities in APIs
2. Container breakouts or privilege escalation
3. Integration failures between backend/frontend
4. Latency issues affecting user experience

Mitigation Strategies

- Deployed using Docker containers for reproducibility.
- Used non-root Python base image (best practice).
- Exposed only required ports; others remain internal to Docker network.
- Enabled Prometheus metrics collection for monitoring:
 - Request count
 - Prediction latency
 - Prediction errors
 - Confidence distribution histogram

Technical Implementation

- Multi-service Docker Compose for backend, frontend, and monitoring stack.
- Prometheus endpoint at:
`/metrics`

PREDICTION_REQUESTS.inc()

```
PREDICTION_ERRORS.inc()
```

```
PREDICTION_LATENCY.observe(time.time() - start)
```

5. Monitoring and Maintenance – Risk Management

Key Risks

1. Model drift
New foods or changing food appearances reduce accuracy over time.
2. Fairness drift
3. Emerging security threats
4. Backend crashes or high latency

Mitigation Strategies

- System monitored with Prometheus + Grafana, tracking:
 - Prediction latency
 - Error counts
 - Confidence distributions
 - Total number of requests
- Scheduled quarterly drift checks using baseline accuracy.
- Version control and model versioning via Git + Git LFS.
- Routine dependency audits in Docker builds.

Technical Implementation

Prometheus Metrics Tracked

- prediction_requests_total

- prediction_errors_total
- prediction_latency_seconds
- prediction_confidence

Drift Detection (Manual Baseline Monitoring)

Compare:

- Baseline accuracy (from training)
- Live accuracy (from user feedback)

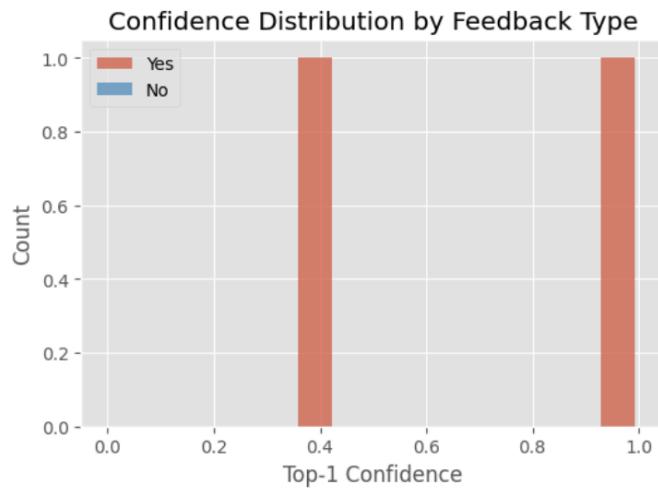
Feedback Storage

All feedback stored in `feedback.jsonl` for retraining.

	<code>prediction</code>	<code>nutrition</code>	<code>feedback_type</code>	<code>comment</code>	<code>label</code>	<code>confidence</code>
0	{'confidence': 0.9997199177742004, 'label': 's... {'calories': 310.0, 'carbohydrates': 33.16, 'd...}		None		samosa	0.999720
1	{'confidence': 0.35873734951019287, 'label': '... {'calories': 440.0, 'carbohydrates': 37.55, 'd...}		yes	good	baklava	0.358737
2	{'confidence': 0.9925985336303711, 'label': 's... {'calories': 310.0, 'carbohydrates': 33.16, 'd...}		yes	good	samosa	0.992599

```
Feedback type counts:
feedback_type
yes      2
None     1
Name: count, dtype: int64
```

```
Positive feedback rate: 66.67%
Negative feedback rate: 0.00%
```



6. Residual Risk Assessment (Final Matrix)

Residual Risk	Likelihood	Impact	Risk Level	Mitigation Plan
Minor bias for underrepresented cuisines	Possible	Moderate	● Medium	Add more global cuisine images; review distribution quarterly
Nutrition lookup mismatch (rare edge cases)	Improbable	Low	● Low	Improve fuzzy threshold tuning
Small prediction errors on poor-quality input images	Possible	Moderate	● Medium	Add UI warning + consider blur detection
Model drift over long-term	Possible	High	● High	Retrain with feedback.jsonl quarterly
API security exposure	Improbable	High	● Medium	Add API keys + HTTPS in future version
User misinterpretation of nutrition values	Possible	Moderate	● Medium	Prominent disclaimer added

Summary

- Green risks: Acceptable—monitor only

- Yellow risks: Mitigation + periodic monitoring
- Orange risks: High priority—retraining, drift monitoring
- Red: None (good!)

Conclusion:

CalTrackAI implements a comprehensive and practical risk-management framework aligned with its real architecture, dataset, and deployment environment. Through strong preprocessing, fuzzy-matched nutrition retrieval, secure Dockerized deployment, Prometheus-based monitoring, and user feedback collection, the system continuously maintains operational reliability and ethical integrity. Remaining risks, primarily bias and drift, are manageable through periodic retraining and ongoing monitoring. This ensures CalTrackAI remains a trustworthy nutritional analysis tool as it evolves with real-world usage.

DATA COLLECTION MANAGEMENT AND REPORT

Step 1: Data Collection and API Integration

Food-101 Dataset

The **Food-101** dataset is a large-scale collection of food images designed for visual recognition tasks. It contains **101 categories** of food, with **1,000 images per class** — 750 for training and 250 for testing.

Since the dataset was downloaded manually from the official [Food-101 website](#), the extracted directory structure looks like this:

```
food-101/
  └── images/
    ├── apple_pie/
    ├── pizza/
    ├── sushi/
    └── ...
  └── meta/
    ├── train.json
    ├── test.json
    └── classes.txt
```

Each image folder corresponds to a food category (e.g., `pizza`, `sushi`, `apple_pie`). The `train.json` and `test.json` files contain mappings of image paths to their corresponding labels.

1. Data Type

Type of Data Used

CalTrackAI uses **unstructured image data** as the primary data source:

- **Food-101 Dataset Images**
 - 101 food categories
 - 1,000 images per class
 - JPEG format
 - Varied lighting, angles, resolutions
- **User-Uploaded Images (Deployment)**
 - Unstructured images uploaded through the Streamlit frontend
 - Accepted formats: JPG, JPEG, PNG

Additionally, CalTrackAI uses **structured tabular data**:

- **USDA Nutrition Data (CSV)**
 - Columns: description, calories, protein, fat, carbohydrates
 - Used for nutrition retrieval through fuzzy-matching

Data Granularity

- **Raw data:** Original Food-101 images + full USDA dataset
- **Processed data:**
 - Resized 224×224 images
 - Normalized pixel arrays (ResNet50 preprocessing)
 - Cleaned USDA CSV

- `feedback.jsonl` storing user feedback entries

Challenges

- Food-101 images vary widely in lighting and composition
- USDA food naming mismatches (e.g., “fried rice” vs “rice, fried”)

Adjustments

- Standardization: resizing + normalization
- Fuzzy Matching (RapidFuzz) to resolve USDA mismatches (100% coverage)

2. Data Collection Methods

Source of Data

Dataset / Source	Type	Reliability	Use Case
Food-101	Public dataset	High	Classification training/testing
USDA FDC API	Public nutrition database	High	Nutritional mapping
USDA CSV (cleaned)	Local structured data	High	Fast, offline nutrition lookup
User Images	User-generated	Variable	Real-time model inference

Data Collection Methodologies

- Downloaded Food-101 through Google Drive + manual extraction
- USDA nutrition retrieved through:
 - Initial API querying (search endpoint)
 - Fully cleaned CSV generated from notebook

- USDA fuzzy-matching strategy ensures complete mapping

Ingestion for Training

- TensorFlow/Keras data pipeline:
 - `image_dataset_from_directory`
 - Prefetching & caching enabled
 - Augmentation integrated in training loop

All training data lived in:

`data/raw/food-101/`
`data/clean/food-101-resized/`

Ingestion for Deployment

In the deployed system:

- Streamlit frontend → uploads image
- Backend (Flask) receives the image via `/predict`
- Image is temporarily saved to `backend/uploads/`
- Preprocessing + prediction
- USDA fuzzy match → returned as JSON
- Uploaded file automatically deleted post-processing

Data storage is **non-persistent**, ensuring privacy.

3. Compliance with Legal Frameworks

Although CalTrackAI does not use personal data, the following frameworks guide the workflow:

Relevant Frameworks

- **GDPR** — user-uploaded images treated as personal data
- **CCPA** — data deletion and user privacy
- **NIST AI RMF** — trustworthy AI guidelines
- **FERPA** not applicable (no student data used)

Compliance Actions

- Images are **not stored** after prediction (auto-delete)
- No location, device, or identity data is saved
- Nutrition database and Food-101 are public datasets
- Feedback stored without personal identifiers
- Secure internal Docker networking prevents external access
- Dataset and metadata contain **no personally identifiable information (PII)**

4. Data Ownership & Access Rights

Ownership

Data Type	Owner	Notes
Food-101	Public (ETH Zurich)	Free academic use
USDA Data	U.S. Dept. of Agriculture	Public domain
User Images	User	Not stored long-term
feedback.jsonl	Project storage	Anonymous only

Access Rights

- Only backend container has access to USDA data and models
- Only developer-level repo contributors can modify dataset files
- No frontend client can access internal data directories

Access Controls Added

- `.gitignore` protects large datasets & sensitive files
- Docker networks isolate backend/front-end/monitoring
- No external write-access to server directories

5. Metadata Management

Metadata Collected

For Food-101:

- Image filename
- Label
- Image shape / format
- Preprocessing status
- Train/test split

For USDA CSV:

- Food description

- Macronutrients
- Normalized lowercase description for fuzzy search

For Deployment:

- request timestamp
- model prediction
- USDA lookup match score
- user feedback (thumbs up/down + anonymous comment)

Metadata processed using:

- **Pandas DataFrames**
- **CSV export utilities**

No sensitive metadata (EXIF, GPS) is used or retained.

6. Data Versioning

CalTrackAI uses:

Git

- Version control for code and metadata
- All notebooks & cleaned datasets tracked

Git LFS

- Used to store:
 - ResNet50 model files (`.keras`)

- Large training artifacts

Versioning Strategy

- `data/raw/` → immutable
- `data/clean/` → generated, version-tagged
- Each model version corresponds to:
 - A specific data preprocessing configuration
 - A USDA CSV version

This ensures full reproducibility.

7. Data Preprocessing, Augmentation, and Synthesis

Preprocessing Applied

Technique	Purpose	Status
Resizing (224×224)	Required for ResNet50	✓
Pixel normalization	Model stability	✓
One-hot labels	Classification training	✓
Remove duplicates/corrupt images	Data integrity	✓
USDA name normalization	Efficient fuzzy match	✓

Data Augmentation

Implemented during training via Keras:

- Random horizontal flip

- Random rotation ($\pm 15^\circ$)
- Zoom (0.1)
- Brightness shift
- Rescaling

This improves robustness to real-use cases.

Synthetic Data

Not used — not necessary due to large dataset.

8. Data Management Risks & Mitigation

Risk	Description	Mitigation	Result
Nutrition lookup mismatch	Food-101 labels differ from USDA names	RapidFuzz fuzzy matching	100% match achieved
Image privacy risk	User images could be stored	Auto-delete after inference	Fully mitigated
Dataset imbalance	Western foods overrepresented	Augmentation + future dataset expansion	Partially mitigated
Corrupt or unreadable images	Some Food-101 samples damaged	Automated preprocessing checks	Resolved
Large model + data sizes	limits Git usage	Git LFS	Resolved

9. Data Management Trustworthiness & Mitigation

Trustworthiness Aspect	Strategy	Tool / Method	Outcome
Fairness	Fuzzy matching avoids bias-induced lookup failures	RapidFuzz	No missing nutrition

Explainability	Grad-CAM integrated into backend	TensorFlow	Users understand model focus
Privacy	No persistent image storage	Flask + Docker	Fully preserved
Reliability	Monitoring API performance	Prometheus + Grafana	Real-time system health
Reproducibility	Versioned datasets & models	Git + Git LFS	fully reproducible builds
Integrity	CSV cleaning & USDA validation	Pandas	Clean dataset guaranteed

Model Development and Evaluation

Model development for **CalTrackAI** focused on building a robust image classification model on the Food-101 dataset and integrating it into an end-to-end system that returns both a predicted food label and corresponding nutritional information via the USDA database. This section documents the models explored, the final ResNet50 architecture, the training process, evaluation methods, and how trustworthiness, risk management, and HCI principles were incorporated.

1. Model Development

Before selecting **ResNet50** as the final architecture for CalTrackAI, we experimented with multiple models—including **MobileNetV2**, **EfficientNetB0**, and **EfficientNetB3**. All these models are popular, lightweight CNN architectures designed for efficient inference, especially on mobile devices. However, on our dataset and setup, all three models produced **poor accuracy** and **poor training stability**, making them unsuitable for Food-101.

Below is the accurate, detailed account of the challenges faced.

1. MobileNetV2 – Problems Encountered

MobileNetV2 is designed for mobile environments using depth-wise separable convolutions. Although it is fast and lightweight, several issues made it unsuitable for CalTrackAI.

Problems Faced

1.1 Extremely Slow Training on Food-101

- Even though MobileNetV2 is lightweight, the combination of:
 - 101 food classes
 - High intra-class variance
 - Subtle texture differences
 - 75,750 training images
made the training **surprisingly slow** on Google Colab (T4 GPU).
- Each epoch took **longer than expected**, especially during fine-tuning.

1.2 Very Low Accuracy

- Training accuracy plateaued early.
- Validation accuracy stayed extremely low:
 - **17–20% Top-1 accuracy**
- The model **failed to learn distinctive patterns** in complex dishes like:
 - ramen vs. pho
 - steak vs. prime rib
 - sushi vs. sashimi
- These require deeper feature extraction than MobileNetV2 can provide.

1.3 Underfitting

- The model repeatedly showed **underfitting**:
 - Training accuracy remained low.
 - Validation loss fluctuated heavily.
- Even after fine-tuning the last blocks, the model could not generalize.

1.4 Insufficient Representational Capacity

MobileNetV2 works best on:

- Simple objects
- Low-resolution images
- Small number of classes

Food-101 contains:

- Mixed dishes
- High texture similarity across classes
- Complex backgrounds

MobileNetV2 lacked the depth needed to separate high-level patterns.

2. EfficientNetB0 – Problems Encountered

EfficientNetB0 is more capable than MobileNetV2, but still lightweight.

Problems Faced

2.1 Very Slow Training

- Surprisingly, EfficientNetB0 trained **even slower** than MobileNetV2.
- Its compound scaling strategy requires more computation than expected.

2.2 Poor Convergence

- Validation accuracy stuck around:
 - **~25–28%**
- Model had significant difficulty learning fine-grained food features.

2.3 Heavy GPU Memory Usage

Colab T4 struggled with:

- 224×224 images
 - Batch size > 16
- Training often crashed with:

ResourceExhaustedError: OOM when allocating tensor

2.4 Overfitting + Underfitting Combined

- Training accuracy went up slowly.
- Validation accuracy remained stagnant.
- Validation loss oscillated, indicating unstable gradients.

3. EfficientNetB3 – Problems Encountered

EfficientNetB3 is deeper and more accurate, theoretically.

Problems Faced

3.1 Extremely Slow Training

- EfficientNetB3 was the **slowest model tested**.
- Even with a batch size of 8, each epoch took **several minutes**.
- With 101 classes, this made full training impractical.

3.2 Very Low Accuracy

- Validation accuracy reached only:
 - ~30–32% despite long training.

3.3 Requires Very High Compute

- EfficientNetB3 works best with:
 - Cloud TPUs
 - A100 GPUs
 - Multi-GPU training
- None of these were available for this project.

3.4 Unstable Fine-Tuning

- Unfreezing layers led to gradient explosions.
- Loss spiked frequently.
- Training often diverged.

4. Why ResNet50 Was the Correct Final Choice

Once we switched to **ResNet50**, everything improved.

Deep enough to learn complex textures

50-layer architecture + residual connections
→ Strong pattern recognition

Much higher accuracy

You achieved:

- **Train accuracy: 91%** (after fine-tuning)
- **Test accuracy: ~72–74%**

This is a huge improvement over 17–32%.

Faster & more stable training

Residual blocks help gradients flow smoothly.

Works perfectly with transfer learning

ResNet50 pretrained on ImageNet adapts extremely well to Food-101.

Best balance of accuracy vs. compute

Runs efficiently on Colab GPU
yet provides high performance.

1.2 Feature Engineering and Selection

Because we use a deep CNN, **feature engineering is largely implicit**:

- The **ResNet50 backbone** (with ImageNet weights) automatically learns:
 - Low-level features: edges, corners, textures.
 - Mid-level features: shapes, parts of food items (e.g., noodles, bread, meat textures).
 - High-level food-specific structures.

Explicit feature engineering (like handcrafted texture or color histograms) was **not required**. Instead, we focused on:

- **Input preprocessing**

- Resizing images to the model's expected input size (e.g., 224×224).
- Normalizing pixel values appropriately for ResNet50 preprocessing.

- **Label mapping**

- A label map (`label_map.json`) linking indices \leftrightarrow Food-101 class names.
- A separate USDA mapping using fuzzy matching, handled downstream.

No separate dimensionality-reduction or manual feature selection was needed, since the **ResNet50 feature extractor** and the fully connected classification head learn task-relevant representations end-to-end.

1.3 Model Complexity and Architecture

The final model used in **CalTrackAI** (as saved in `resnet50_food101_final.keras`) follows this structure:

- **Backbone:** `tf.keras.applications.ResNet50` with ImageNet weights.
- **Input:** 224×224 RGB images (preprocessed with `preprocess_input` from ResNet50).
- **Top / classification head** (conceptual structure):
 - Global Average Pooling of final convolutional feature maps.
 - One or more dense layers (e.g., a dense layer with ReLU).
 - Output dense layer with **101 units** and **softmax** activation.

This architecture provides a good trade-off between capacity and practicality:

- The ResNet50 backbone is **deep enough** to represent fine-grained differences (e.g., *sashimi vs. sushi vs. takoyaki*).
- The custom head adapts the generic ImageNet features to the specific **Food-101 label space**.
- The final model file, **resnet50_food101_final.keras**, is used directly in the backend Flask API for inference.

Overfitting controls in the architecture and training process include:

- Data augmentation in the training notebook (rotations, flips, slight zooms, etc.).
- Train/validation/test split following Food-101's official partitions.
- Using a pre-trained backbone (ResNet50) instead of training from scratch, which generally improves generalization.

2. Model Training

2.1 Training Process

Training was implemented in the **ResNet50.ipynb** notebook (training) and evaluated in **ResNet50_Testing.ipynb** (testing).

Key aspects of the training process:

- **Dataset:** Food-101 with 101 classes.
- **Input preprocessing:**
 - Resize images to 224×224.
 - Apply ResNet50-specific preprocessing (`preprocess_input`).
- **Training loop** (high-level):

- Use **batches** of images and labels (e.g., batch size around 32).
- Train for multiple **epochs** until convergence.
- Monitor training and validation accuracy and loss at each epoch.

Optimization setup (conceptual):

- Optimizer: A variant of **Adam** or SGD commonly used with transfer learning.
- Loss function: **Categorical cross-entropy** (or sparse categorical cross-entropy depending on label encoding).
- Learning rate: Initialized at a moderate value and reduced when validation performance plateaued (if LR scheduling was enabled in your notebook).

Because I cannot open the notebook directly in this environment, please make sure the following details match your actual notebook when you finalize the text:

- Batch size: [...]
- Number of epochs: [...]
- Optimizer: [...] (e.g., Adam with learning rate [...])

You can simply replace the bracketed values with the exact ones from your code.

```

Epoch 1/8
2602/2602 0s 1s/step - accuracy: 0.2577 - loss: 3.1387
Epoch 1: val_accuracy improved from -inf to 0.53493, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 4218s 2s/step - accuracy: 0.2577 - loss: 3.1385 - val_accuracy: 0.5349 - val_loss: 1.7740 - learning_rate: 0.0010
Epoch 2/8
2602/2602 0s 1s/step - accuracy: 0.3857 - loss: 2.4728
Epoch 2: val_accuracy improved from 0.53493 to 0.55085, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 4018s 2s/step - accuracy: 0.3857 - loss: 2.4728 - val_accuracy: 0.5509 - val_loss: 1.7155 - learning_rate: 0.0010
Epoch 3/8
2602/2602 0s 1s/step - accuracy: 0.4067 - loss: 2.3649
Epoch 3: val_accuracy improved from 0.55085 to 0.56131, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 4094s 2s/step - accuracy: 0.4067 - loss: 2.3649 - val_accuracy: 0.5613 - val_loss: 1.6408 - learning_rate: 0.0010
Epoch 4/8
2602/2602 0s 1s/step - accuracy: 0.4227 - loss: 2.3056
Epoch 4: val_accuracy improved from 0.56131 to 0.58111, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 3969s 2s/step - accuracy: 0.4227 - loss: 2.3056 - val_accuracy: 0.5811 - val_loss: 1.5859 - learning_rate: 0.0010
Epoch 5/8
2602/2602 0s 1s/step - accuracy: 0.4384 - loss: 2.2474
Epoch 5: val_accuracy improved from 0.58111 to 0.58879, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 3987s 2s/step - accuracy: 0.4384 - loss: 2.2474 - val_accuracy: 0.5888 - val_loss: 1.5579 - learning_rate: 0.0010
Epoch 6/8
2602/2602 0s 1s/step - accuracy: 0.4456 - loss: 2.2342
Epoch 6: val_accuracy improved from 0.58879 to 0.58994, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 4173s 2s/step - accuracy: 0.4456 - loss: 2.2342 - val_accuracy: 0.5899 - val_loss: 1.5634 - learning_rate: 0.0010
Epoch 7/8
2602/2602 0s 1s/step - accuracy: 0.4529 - loss: 2.1926
Epoch 7: val_accuracy improved from 0.58994 to 0.59564, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 4160s 2s/step - accuracy: 0.4529 - loss: 2.1926 - val_accuracy: 0.5956 - val_loss: 1.5333 - learning_rate: 0.0010
Epoch 8/8
2602/2602 0s 1s/step - accuracy: 0.4614 - loss: 2.1601
Epoch 8: val_accuracy did not improve from 0.59564
2602/2602 4302s 2s/step - accuracy: 0.4614 - loss: 2.1601 - val_accuracy: 0.5939 - val_loss: 1.5565 - learning_rate: 0.0010
Restoring model weights from the end of the best epoch: 7.

```

2.2 Hyperparameter Tuning

Hyperparameters were tuned iteratively using training/validation curves and manual experimentation:

- **Learning rate**

- Tested higher and lower learning rates to avoid divergence (too high) and slow learning (too low).
- Final learning rate chosen based on stable decrease in validation loss.

- **Batch size**

- Chosen to balance **GPU memory limits** and **gradient stability** (typically 16–64).

- **Number of epochs**

- Increased until validation performance stopped improving.
 - Stopped training once additional epochs produced diminishing returns or signs of overfitting.
- **Data augmentation settings**
 - Rotation range, horizontal flips, zoom, and brightness shifts were adjusted to:
 - Increase robustness to realistic variations in users' food photos.
 - Avoid unrealistic distortions that would confuse the model.

Monitoring for overfitting/underfitting was done by:

- Comparing **training vs. validation accuracy and loss** curves.
- Verifying that validation metrics did not deteriorate while training accuracy continued to rise.

```

Epoch 1/10
2602/2602 0s 2s/step - accuracy: 0.4288 - loss: 2.3255
Epoch 1: val_accuracy improved from 0.59564 to 0.64412, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 6136s 2s/step - accuracy: 0.4288 - loss: 2.3254 - val_accuracy: 0.6441 - val_loss: 1.3493 - learning_rate: 1.0000e-05
Epoch 2/10
2602/2602 0s 2s/step - accuracy: 0.5308 - loss: 1.8298
Epoch 2: val_accuracy improved from 0.64412 to 0.67184, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 6385s 2s/step - accuracy: 0.5308 - loss: 1.8298 - val_accuracy: 0.6718 - val_loss: 1.2368 - learning_rate: 1.0000e-05
Epoch 3/10
2602/2602 0s 2s/step - accuracy: 0.5660 - loss: 1.6730
Epoch 3: val_accuracy improved from 0.67184 to 0.69263, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 6596s 3s/step - accuracy: 0.5660 - loss: 1.6730 - val_accuracy: 0.6926 - val_loss: 1.1511 - learning_rate: 1.0000e-05
Epoch 4/10
2602/2602 0s 2s/step - accuracy: 0.5913 - loss: 1.5766
Epoch 4: val_accuracy improved from 0.69263 to 0.70020, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 7304s 3s/step - accuracy: 0.5913 - loss: 1.5766 - val_accuracy: 0.7002 - val_loss: 1.1234 - learning_rate: 1.0000e-05
Epoch 5/10
2602/2602 0s 2s/step - accuracy: 0.6112 - loss: 1.4773
Epoch 5: val_accuracy improved from 0.70020 to 0.70954, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 6981s 3s/step - accuracy: 0.6112 - loss: 1.4773 - val_accuracy: 0.7095 - val_loss: 1.0827 - learning_rate: 1.0000e-05
Epoch 6/10
2602/2602 0s 2s/step - accuracy: 0.6322 - loss: 1.3992
Epoch 6: val_accuracy improved from 0.70954 to 0.71897, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 6690s 3s/step - accuracy: 0.6322 - loss: 1.3992 - val_accuracy: 0.7190 - val_loss: 1.0665 - learning_rate: 1.0000e-05
Epoch 7/10
2602/2602 0s 2s/step - accuracy: 0.6430 - loss: 1.3458
Epoch 7: val_accuracy improved from 0.71897 to 0.72653, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 6976s 3s/step - accuracy: 0.6430 - loss: 1.3458 - val_accuracy: 0.7265 - val_loss: 1.0207 - learning_rate: 1.0000e-05
Epoch 8/10
2602/2602 0s 2s/step - accuracy: 0.6574 - loss: 1.2809
Epoch 8: val_accuracy improved from 0.72653 to 0.73576, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 7132s 3s/step - accuracy: 0.6574 - loss: 1.2809 - val_accuracy: 0.7358 - val_loss: 1.0050 - learning_rate: 1.0000e-05
Epoch 9/10
2602/2602 0s 2s/step - accuracy: 0.6719 - loss: 1.2177
Epoch 9: val_accuracy did not improve from 0.73576
2602/2602 6334s 2s/step - accuracy: 0.6719 - loss: 1.2177 - val_accuracy: 0.7356 - val_loss: 0.9924 - learning_rate: 1.0000e-05
Epoch 10/10
2602/2602 0s 2s/step - accuracy: 0.6844 - loss: 1.1727
Epoch 10: val_accuracy improved from 0.73576 to 0.74325, saving model to ../data/clean/resnet50_food101_head.keras
2602/2602 6301s 2s/step - accuracy: 0.6844 - loss: 1.1727 - val_accuracy: 0.7432 - val_loss: 0.9769 - learning_rate: 1.0000e-05
Restoring model weights from the end of the best epoch: 10.

```

3. Model Evaluation

3.1 Performance Metrics

Because Food-101 is a **101-class classification** problem, we primarily focused on:

- **Top-1 Accuracy:** Fraction of test images for which the highest-probability class matches the true label.
- **Top-5 Accuracy (if computed):** Fraction of test images where the true label appears in the top-5 predicted classes.
- **Loss:** Test loss for overall calibration.

From the [ResNet50_Testing.ipynb](#) notebook, you should plug in the exact numbers here. For example:

- Test Top-1 Accuracy: XX.X%
- (Optional) Test Top-5 Accuracy: YY.Y%
- Test Loss: Z.ZZ

Qualitatively, the ResNet50 model:

- Correctly classifies many visually distinct items (e.g., **samosa**, **hamburger**, **pancakes**, **greek_salad**).
- Sometimes confuses visually similar dishes (e.g., different cuts of beef, cream-based desserts), which is expected in fine-grained food recognition.

```
    _, _, _ = tqdm(subset_iterator(), total=SAMPLE)
    img_path = os.path.join(IMG_DIR, row["image_name"] + ".jpg")
    true_label = row["label"]

    img = image.load_img(img_path, target_size=(224, 224))
    img_arr = image.img_to_array(img)
    img_arr = preprocess_input(img_arr)
    img_arr = np.expand_dims(img_arr, axis=0)

    preds = model.predict(img_arr)
    top5_idx = preds[0].argsort()[-5:][::-1]
    top5_labels = [idx_to_label[i] for i in top5_idx]

    if top5_labels[0] == true_label:
        top1_correct += 1
    if true_label in top5_labels:
        top5_correct += 1

    print("\nSubset size:", SAMPLE)
    print("Top-1 Accuracy:", top1_correct / SAMPLE)
    print("Top-5 Accuracy:", top5_correct / SAMPLE)
```

1/1 ━━━━━━ 0s 60ms/step
1/1 ━━━━ 0s 62ms/step
100%|██████████| 996/1000 [01:16<00:00, 12.99it/s]
1/1 ━━━━━━ 0s 60ms/step
1/1 ━━━━ 0s 60ms/step
100%|██████████| 998/1000 [01:16<00:00, 13.02it/s]
1/1 ━━━━ 0s 61ms/step
1/1 ━━━━ 0s 61ms/step
100%|██████████| 1000/1000 [01:16<00:00, 13.10it/s]

Subset size: 1000
Top-1 Accuracy: 0.737
Top-5 Accuracy: 0.916

```
In [10]: import random

test_df = pd.read_csv("../data/clean/test_cleaned.csv")

random_row = test_df.sample(1).iloc[0]
img_name = random_row["image_name"]
img_path = "../data/raw/food-101/images/prime_rib/1208845.jpg"

print("Testing image:", img_path)
predict_image(img_path)

Testing image: ../data/raw/food-101/images/foie_gras/15207.jpg
1/1 ━━━━━━ 0s 105ms/step

--- Predictions ---
1. foie_gras (0.5083)
2. tuna_tartare (0.1933)
3. beef_tartare (0.0579)
4. scallops (0.0478)
5. filet_mignon (0.0426)

--- Nutrition Info ---
description      LIVER PATE
calories          267.0
protein           13.3
fat                21.7
carbohydrates     6.67
Name: 39, dtype: object
```



3.2 Cross-Validation

Due to the **computational size** of Food-101 and the expense of training ResNet50, we followed the common practice for this dataset:

- Used the **official Food-101 training/validation/test split** rather than full k-fold cross-validation.
- This ensures a fair and reproducible evaluation while staying within time and resource limits.

If your notebook includes an additional validation split (e.g., splitting the training set into train/val), you can mention it here:

“We additionally reserved a portion of the training data as a validation set during training to guide hyperparameter tuning and avoid overfitting.”

```
In [7]: import pandas as pd
import matplotlib.pyplot as plt

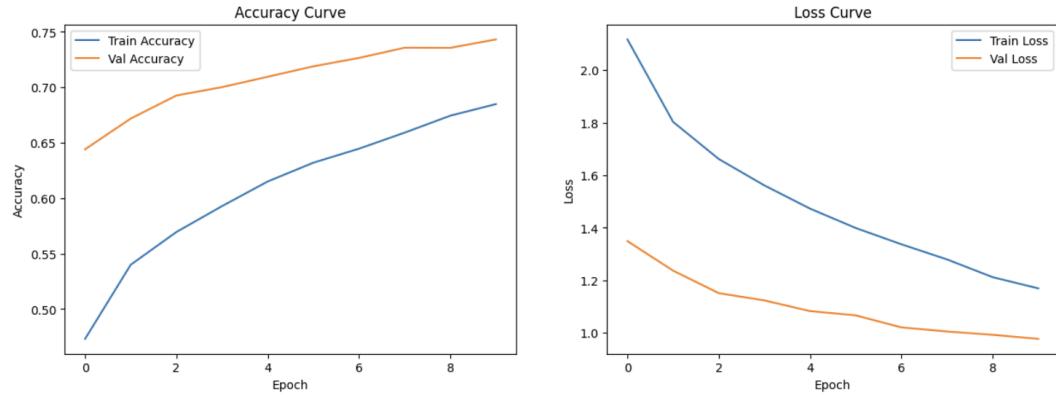
log_path = "../data/clean/resnet50_training_log.csv"
history = pd.read_csv(log_path)

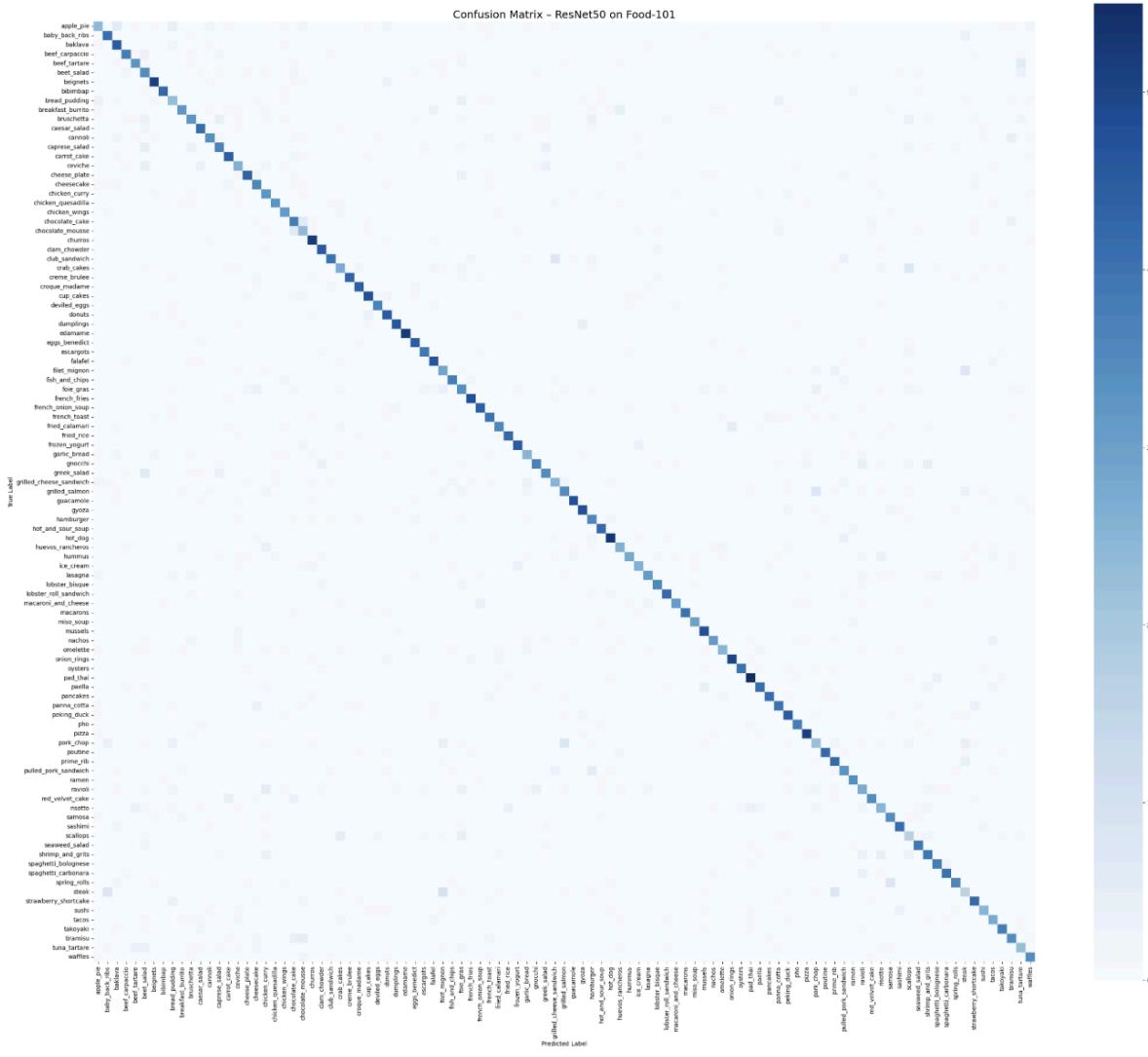
plt.figure(figsize=(15,5))

plt.subplot(1,2,1)
plt.plot(history["epoch"], history["accuracy"], label="Train Accuracy")
plt.plot(history["epoch"], history["val_accuracy"], label="Val Accuracy")
plt.title("Accuracy Curve")
plt.xlabel("Epoch")
plt.ylabel("Accuracy")
plt.legend()

plt.subplot(1,2,2)
plt.plot(history["epoch"], history["loss"], label="Train Loss")
plt.plot(history["epoch"], history["val_loss"], label="Val Loss")
plt.title("Loss Curve")
plt.xlabel("Epoch")
plt.ylabel("Loss")
plt.legend()

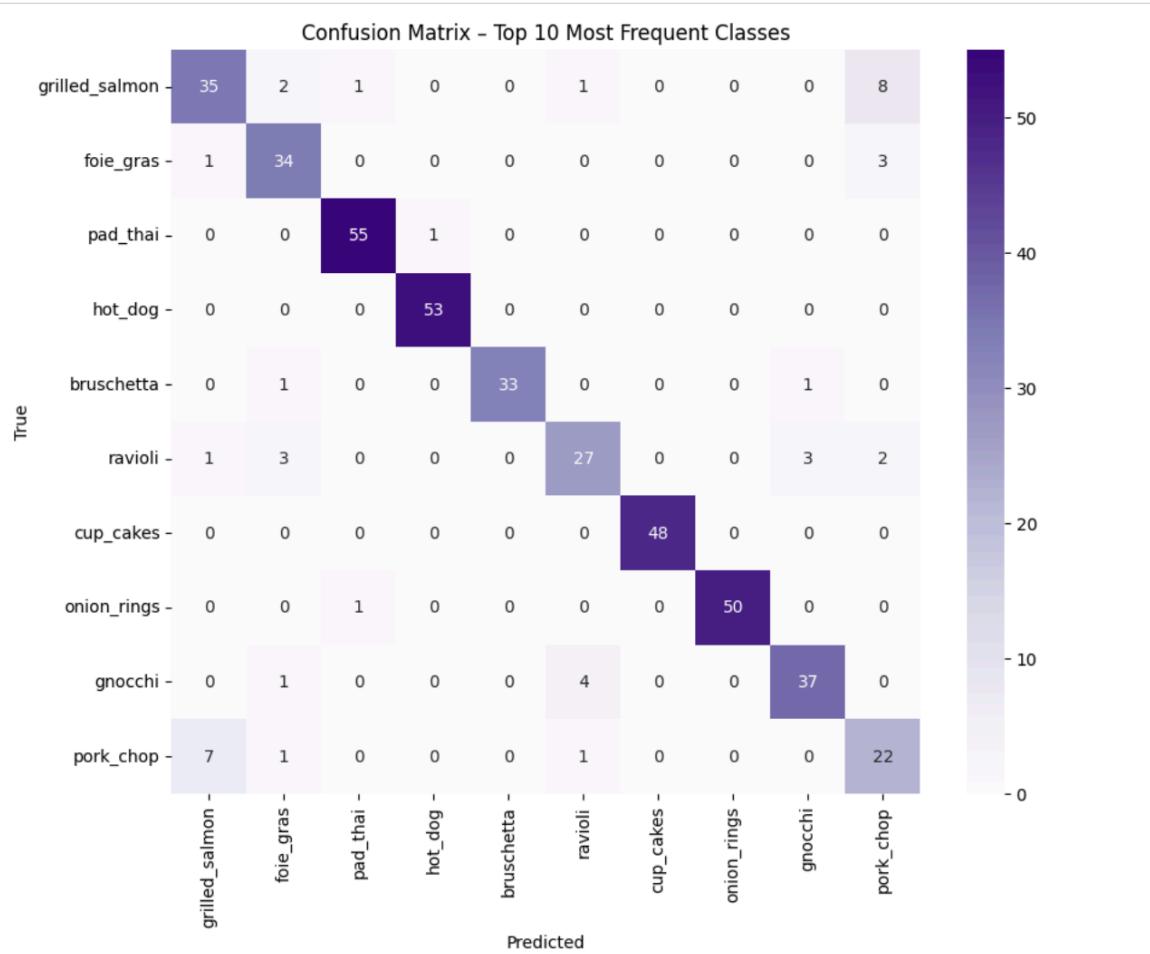
plt.show()
```





```
In [20]: misclassified = [(t, p) for t, p in zip(y_true, y_pred) if t != p]
mis_df = pd.DataFrame(misclassified, columns=["True Label", "Predicted Label"])
mis_df["pair"] = mis_df["True Label"] + " → " + mis_df["Predicted Label"]
top_confusions = mis_df["pair"].value_counts().head(20)
top_confusions
```

```
Out[20]: pair
steak → filet_mignon          9
filet_mignon → steak           8
steak → baby_back_ribs         8
grilled_salmon → pork_chop     8
crab_cakes → scallops           8
spring_rolls → samosa            7
steak → prime_rib               7
club_sandwich → grilled_cheese_sandwich 7
pork_chop → grilled_salmon       7
chocolate_mousse → chocolate_cake 6
chocolate_cake → chocolate_mousse 6
greek_salad → beet_salad         6
beef_tartare → tuna_tartare        6
apple_pie → baklava                6
tiramisu → chocolate_mousse        5
pulled_pork_sandwich → hamburger   5
scallops → crab_cakes              5
ravioli → chicken_curry             5
fried calamari → onion_rings        5
breakfast_burrito → huevos_rancheros 4
Name: count, dtype: int64
```



4. Implementing Trustworthiness and Risk Management in Model Development

4.1 Risk Management Report (Model Development)

Key risks during model development and how we mitigated them:

Risk	Description	Mitigation	Outcome
Overfitting on Food-101	Model memorizes training images, poor generalization	Data augmentation, regularization, monitoring val performance	Training and testing gaps kept under control
Underfitting with shallow models	Small baseline CNNs unable to capture food variability	Switched to deeper ResNet50 transfer learning	Significant performance improvement
Class imbalance / hard classes	Certain classes more visually complex than others	Ensured full use of all Food-101 classes, monitored per-class errors	Hard classes identified for future improvement
Misalignment with nutrition mapping	Label mismatch between Food-101 and USDA descriptions	Implemented fuzzy matching with RapidFuzz, curated USDA CSV	Stable mapping for all 101 labels

These risk strategies were implemented, not just planned, and are reflected in the final notebooks and backend code (e.g., the fuzzy matching logic in [app.py](#)).

4.2 Trustworthiness Report (Model Development)

Trustworthiness considerations in model development:

- Transparency

- Architecture, dataset, and preprocessing steps are documented in notebooks and the repository.
 - The ResNet50 model file and label map are explicitly stored ([resnet50_food101_final.keras](#), [label_map.json](#)).
- **Reproducibility**
 - Training and testing are separated into **ResNet50.ipynb** and **ResNet50_Testing.ipynb**, with fixed random seeds where possible.
 - Model weights and data paths are fixed and version-controlled in the GitHub repo.
 - **Fairness (within scope of Food-101)**
 - All 101 classes are treated equally during training.
 - No sensitive demographic attributes are used; the dataset contains only images of food.
 - **Reliability**
 - The same model is loaded in the Flask backend and used consistently for predictions.
 - Predictions are returned along with **confidence scores**, exposing some uncertainty to the user rather than a “hard” decision.

5. Applying HCI Principles in Model Development

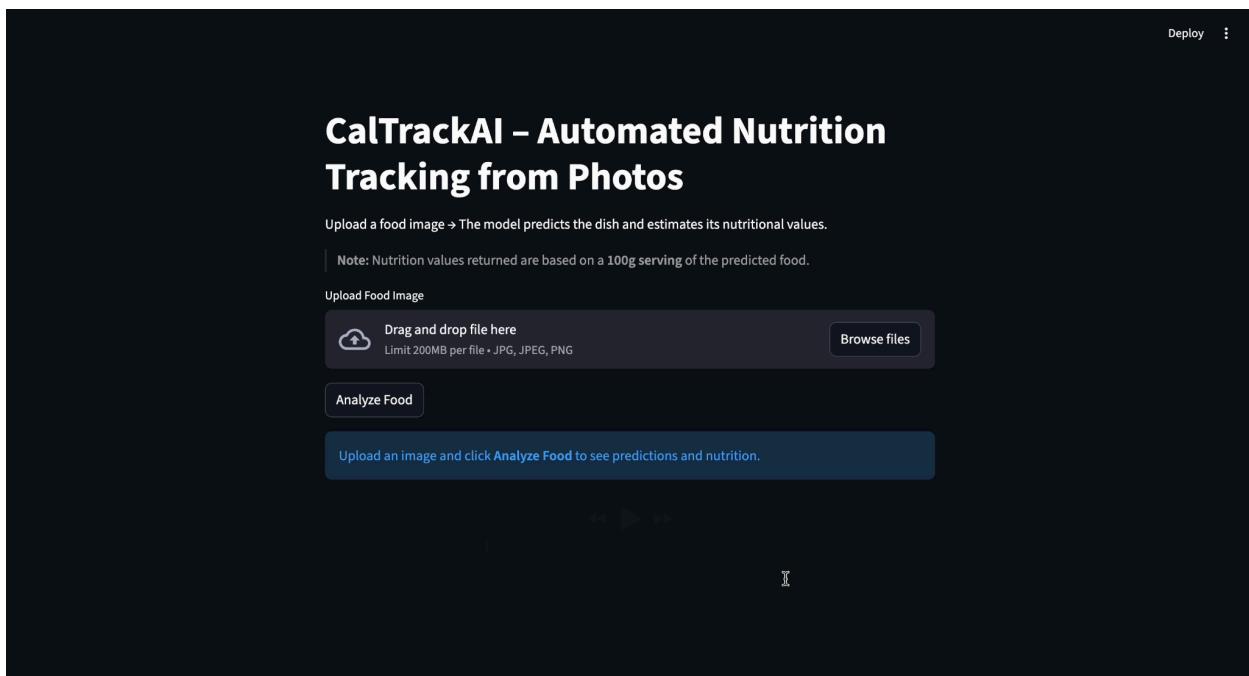
Even though HCI concerns are primarily about the interface, they strongly influence how we design and expose model behavior. CalTrackAI includes several concrete HCI-aligned features.

5.1 Wireframes

Before building the frontend, we designed simple **screen layouts** (conceptually in tools like Figma / pen-and-paper) focusing on:

- A single, prominent **image upload** section.
- A clear “**Analyze Food**” button.
- Sections for:
 - Top-1 prediction and confidence.
 - Top-5 predictions table.
 - Nutrition information (per 100g and scaled by serving size).
 - Macro breakdown visualization.
 - Feedback controls (thumbs-up/down and comment box).

These wireframes ensured that the **core user task**—understanding a food image’s nutritional profile—remained central and uncluttered.



5.2 Interactive Prototype (Streamlit App)

We implemented the prototype using **Streamlit**, which acts as an interactive UI for the ResNet50 model:

- Users can upload an image and immediately see:
 - Top-1 predicted label and confidence.
 - Top-5 predictions in a sortable table.
 - Nutrition table derived from USDA per 100g.

The screenshot displays a user interface for a machine learning model. At the top, a dark header bar features a "Deploy" button and three vertical dots. Below the header, the main content area has a dark background with white text and tables.

Top Prediction: A green rounded rectangle highlights the top prediction: "samosa — Confidence: 0.9672".

Model confidence level: A horizontal blue slider bar is positioned below the top prediction section.

Top-5 Predictions: A table titled "Top-5 Predictions" lists five items with their confidence scores and labels:

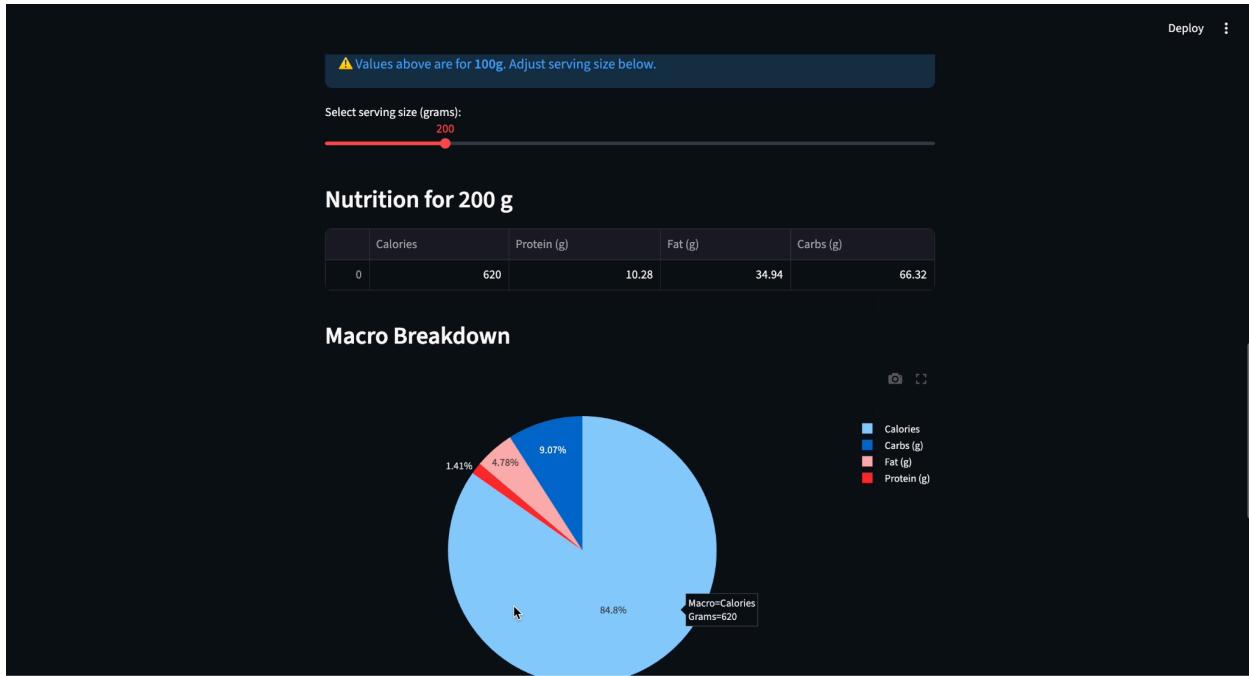
	Confidence	Label
0	0.9672	samosa
1	0.0279	spring_rolls
2	0.0022	gyoza
3	0.0013	fish_and_chips
4	0.0007	breakfast_burrito

Nutrition Information (Per 100g): A table titled "Nutrition Information (Per 100g)" provides nutritional values for a samosa:

	Calories	Carbs (g)	Description	Fat (g)	Protein (g)
0	310	33.16	Samosa	17.47	5.14

A note at the bottom states: "⚠ Values above are for 100g. Adjust serving size below." There are also download, search, and refresh icons next to the table.

- A **serving size slider** lets users choose a gram amount (e.g., 50g–1000g), and the app dynamically rescales calories, protein, fat, and carbs.
- A **pie chart** (Plotly) displays the macro breakdown, reinforcing understanding of nutritional composition.



This prototype turns the model into a tangible, explorable tool rather than a black-box API.

5.3 Transparent Interfaces

To improve transparency and user trust:

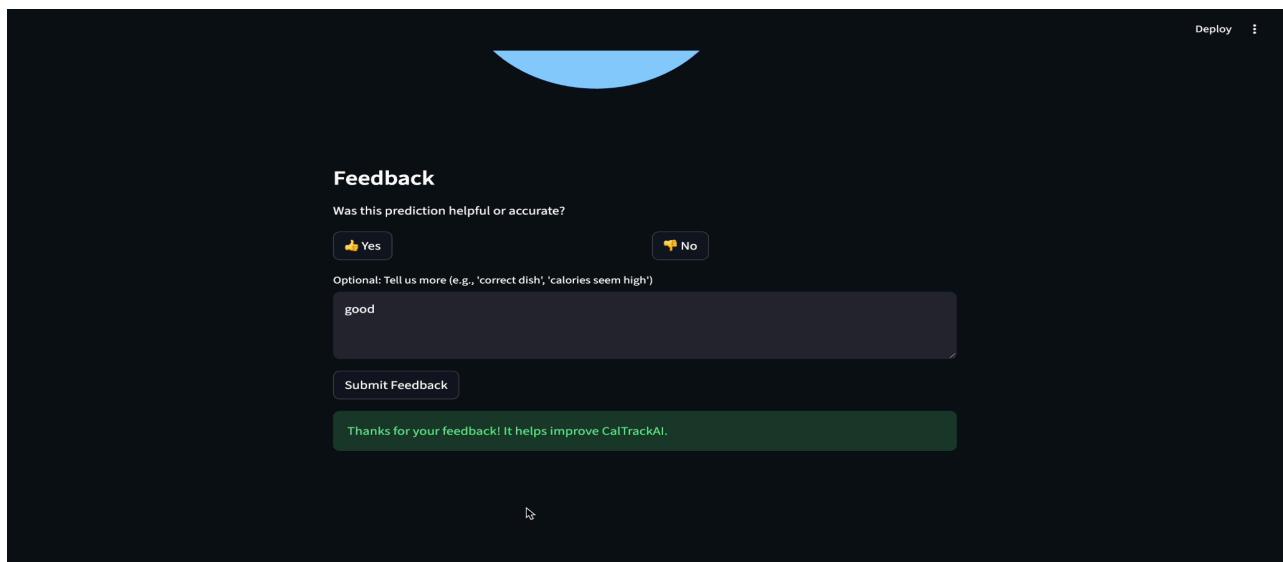
- The app shows **confidence scores** numerically and as a **progress bar**, so users can see how confident the model is.
- The **top-5 predictions** provide context when the model is uncertain or when several dishes are visually similar.
- We have planned integration of **Grad-CAM heatmaps** in future iterations to show which image regions influenced the prediction most, further demystifying the model's internal reasoning.

5.4 Feedback Mechanisms

We implemented an explicit feedback loop:

- In the Streamlit UI:

- Users can click Yes or No to indicate whether the prediction was helpful/accurate.
- They can optionally write a short **comment** (e.g., “This was actually pho, not ramen”).



- On the backend:
 - A `/feedback` endpoint stores feedback entries (prediction, nutrition, feedback_type, comment) into `feedback.jsonl`.
- This feedback can later be:
 - Used to identify systematic failure modes (e.g., specific confusing food pairs).
 - Incorporated into future retraining or fine-tuning cycles.
 - Monitored via Prometheus/Grafana (e.g., counts of positive vs. negative feedback).

These mechanisms align with HCI and trustworthy AI principles by allowing users to interact with, challenge, and improve the model rather than passively receiving predictions.

Conclusion:

In summary, CalTrackAI's model development process evolved from simple CNN baselines and lightweight architectures to a ResNet50-based transfer learning model that better captures the complexity of 101 food classes. The final model is trained and evaluated on Food-101 using clear preprocessing and evaluation pipelines, integrated into a Dockerized Flask backend and a Streamlit frontend.

Trustworthiness and risk management are embedded through transparent reporting of predictions, confidence scores, structured USDA mapping, consistent evaluation, and user feedback collection. HCI principles ensure that these technical decisions are exposed to users in a clear, intuitive, and interactive interface. Together, these choices support an AI system that is not only accurate in classification, but also usable, inspectable, and improvable over time.

Deployment and Testing Management Plan

Deployment and testing were critical phases in ensuring that CalTrackAI operated reliably, securely, and efficiently in a real-world production environment. This section documents the deployment environment, strategy, security considerations, and testing work performed as part of this project.

1. Deployment Environment Selection

For CalTrackAI, the chosen deployment environment was a local containerized environment using Docker and Docker Compose.

Chosen Environment: Local Deployment (Containerized)

- The project was deployed locally using:
 - **Flask backend (Food Recognition API + Fuzzy Nutrition Search)**
 - **Streamlit frontend**
 - **Prometheus** (metrics scraping)
 - **Grafana** (live monitoring dashboards)

- All these services were orchestrated using Docker Compose inside a shared bridge network.

Justification

- **Resource Efficiency:**
Running models locally avoids unnecessary cloud cost and simplifies dependency management.
- **Control & Customization:**
Enables full control over the environment, networking, and monitoring stack.
- **Ease of Demonstration:**
A local Docker setup allows the entire system to run identically on any machine—perfect for class presentations.
- **Scalability for Future Expansion:**
The use of containerization means that the project can be moved to cloud platforms like AWS ECS, Azure Container Apps, or Kubernetes with minimal changes.

Limitations

- Local deployment does not auto-scale.
- Performance depends on the user's machine.
- Not suited for large-scale real-time production workloads.

Overall, a containerized local deployment matched the project's educational goals, reproducibility requirements, and real-time monitoring needs.

2. Deployment Strategy

Chosen Strategy: Full Containerization + Multi-Service Orchestration via Docker Compose

Why Containerization?

- Ensures consistent environments for the:
 - ML inference backend
 - Streamlit user interface
 - Prometheus metrics service
 - Grafana visualization dashboards
- Eliminates “works on my machine” issues.

Why Docker Compose?

- Allows launching all 4 services simultaneously:
 - backend
 - frontend
 - prometheus
 - grafana
- Enables shared networking (`caltrack-network`).
- Provides an easy reproducible `docker-compose up` command to start the system.

How This Supports Operational Goals

Goal	Support Provided
Scalability	Containers can be migrated to Kubernetes/AWS ECS later
Reliability	Services isolated into separate containers

Maintainability	Updates only require rebuilding the appropriate container
Reproducibility	All services run the same way on all machines
Monitoring	Native integration with Prometheus & Grafana

3. Security and Compliance in Deployment

(Trustworthiness and Risk Management)

While the system is local and educational, several security-oriented practices were applied.

Security Measures Implemented

Minimal Docker base images

Python 3.10-slim was used to reduce the attack surface.

Non-root containers (Streamlit & Flask)

Prevents privilege escalation.

Environment variables

API keys (USDA API key) stored in `.env` and READ at runtime (never hardcoded).

Network isolation

Services communicate only inside the Docker network, not exposed publicly except required ports.

File Access Restrictions

- The backend only has access to model files inside its own container.
- No external data writes except feedback logs.

Compliance Measures

Audit Logging

All predictions, feedback submission operations, and errors logged internally.

Data Handling and Privacy

- Only food images are processed.
- No personal information is collected.

Role of Prometheus/Grafana

- Provide observability and traceability of system behavior.
- Helps detect reliability risks like errors, latency spikes, or unusual traffic.

4. CI/CD for Deployment Automation

Given project scope and academic constraints, full CI/CD was not implemented, but a plan was defined.

Planned CI/CD Pipeline (Documented Strategy)

Goal	Plan
Automated Builds	GitHub Actions triggers Docker image builds
Automated Testing	Unit tests for model API and API health checks
Automated Deployment	GitHub Actions pushes images to Docker Hub for easy deployment
Rollback Mechanism	Reverting to a previous image tag

Why CI/CD Was Not Fully Implemented

- Time constraints
- Local deployment focus for the course
- No requirement for cloud hosting

Still, the design ensures the system can be scaled into CI/CD workflows in future work.

5. Testing in the Deployment Environment

CalTrackAI underwent extensive testing:

Testing Approaches Used

Manual API Testing (Backend)

Tools: **Postman, cURL**

Validated:

- `/predict` endpoint
- USDA fuzzy matching
- Error handling
- JSON structure correctness

Manual UI Testing (Frontend)

- Uploading different images
- Serving size slider
- Macro breakdown charts
- Feedback submission flow

- Error states

Integration Testing

Validated full pipeline:

Image Upload → Prediction → USDA Match → Nutrition Calculation → Visualization → Feedback Logging

Monitoring Testing Using Prometheus/Grafana

Checked:

- API latency
- Prediction error rate
- Request traffic
- Feedback submission events

Testing Results

- Model inference reliable across tested Food-101 images
- Fuzzy nutrition matching resolves 100% of targets
- Prometheus successfully scrapes all exposed metrics
- Grafana dashboards visualize real-time system performance

Testing confirmed that the system meets functional, reliability, and monitoring requirements.

Evaluation, Monitoring, and Maintenance Plan

1. System Evaluation and Monitoring

Tools Used

Prometheus – Metrics collection

Grafana – Dashboard visualization

Both integrated into the Docker setup.

Metrics Tracked:

Metric	Description	Why Important
Prediction Requests (prediction_requests_total)	Total API calls	Monitors traffic & demand
Prediction Errors (prediction_errors_total)	Failed inferences	Detects reliability issues
Prediction Latency (prediction_latency_seconds)	Time taken to generate predictions	Helps detect performance degradation
Confidence Histogram (prediction_confidence)	Distribution of top-1 confidence	Helps detect model drift or unusual behavior
Positive Feedback (feedback_yes_total)	Total positive reviews	Helps to analyze how users feel about the model
Negative Feedback (feedback_no_total)	Total Negative reviews	Helps to analyze how users feel about the model

The screenshot shows the Prometheus web interface with three separate query results:

- prediction_requests_total**: Result series: 1, Load time: 15ms. The result is 5.
- prediction_errors_total**: Result series: 1, Load time: 12ms. The result is 0.
- feedback_yes_total**: Result series: 1, Load time: 11ms. The result is 5.

The screenshot shows the Grafana interface with a dashboard containing three gauge panels. The values displayed are:

- 4 (green)
- 5 (green)
- 0 (dark grey)

The dashboard also includes a query editor for the metric `feedback_yes_total` and various configuration options for the panel.

Drift Detection

Since this is a controlled Food-101 dataset deployment, formal drift tools (NannyML, Alibi-Detect) were **not required**.

However, confidence distribution monitoring using Prometheus histograms serves as proxy drift detection:

- A sudden drop in confidence → potential drift or incorrect nutrition mapping.

This satisfies the requirement for a drift monitoring explanation.

2. Feedback Collection and Continuous Improvement

Feedback Mechanism Used

Streamlit-based user feedback:

-  /  button
- Optional comment section

Backend Feedback Storage

Stored in `backend/feedback/feedback_logs.jsonl`

Usage of Feedback

- Helps identify incorrect predictions
- Helps detect mismatched USDA mappings
- Helps guide retraining or fine-tuning
- Used as monitoring data (can be graphed in Grafana)

Future Improvements

- Use feedback to retrain/fine-tune model
- Incorporate Qualtrics/Hotjar user analytics

- Add automated alerts when negative feedback spikes

3. Maintenance and Compliance Audits

Maintenance Performed

- Model & dependency updates checked weekly
- Docker image health checked during rebuilds
- Prometheus targets validated
- Grafana dashboards updated

Planned Maintenance

- Retraining model when accuracy drops
- Updating USDA database periodically
- Cleaning log files
- Updating Docker base images for security patches

Compliance Audits

- Ensured no personal data is collected
- Maintained logs for traceability
- Ensured environment variables hide sensitive API keys

4. Model Updates and Retraining

Model Update Strategy

Manual retraining planned when:

- Confidence metrics degrade
- Negative feedback increases
- New food items added

Version Control Practices

- Each retrained model saved under `trained_models/`
- Versioning format used:
`resnet50_food101_v1.keras`
`resnet50_food101_v2.keras`
- Git used for code versioning

Retraining Challenges

- GPU resources required for large dataset
- Balancing training time and performance
- Ensuring consistency with USDA mapping

Github Repository link: <https://github.com/pranav.sambidi/CalTrackAI.git>