

# Geometry-Enhanced Portion Estimation for Multimodal LLMs

8up.ai Research

research@8up.ai

July 7, 2026

## Abstract

Image-based dietary assessment promises to replace costly, bias-prone manual recalls, but portion estimation remains a major blocker. Multimodal LLMs (MLLMs) recognize a wide range of foods zero-shot in uncontrolled photos, yet they are weak at portion estimation—a gap we measure across the current frontier (Gemini, GPT, and Claude flagships alike). We present a method that *enhances a frozen, commercial MLLM with an accurate portion head*: a small geometry-enhanced network on a frozen DINOv2 backbone with a *structured softmax-ownership volume*, consuming the MLLM’s per-food name, bounding box, and density range—no depth sensor, no MLLM fine-tuning. Evaluated fully open-vocabulary on three real-world benchmarks, the head cuts per-food portion error by **33–41 %** relative to the MLLM alone, outperforms every flagship MLLM’s direct estimates, and surpasses each benchmark’s originally published image-only model at its own reported metric.

## 1 Introduction

Dietary assessment underpins nutrition research and public-health monitoring, but the conventional instruments face a trade-off between burden, detail, and accuracy [13]. A 24-hour dietary recall captures fine-grained intake but asks respondents to remember and report everything they ate, which is burdensome and prone to recall bias. Food-frequency questionnaires are easier to administer—respondents only report how often they eat items from a fixed list—but sacrifice the detail needed to quantify specific foods and portions. Image-based food recognition and portion estimation offers a scalable, objective alternative that aims for detail without the respondent burden [1]. Existing methods, however, are developed on small, controlled datasets and do not transfer to real life. Many tend to report only a single dish-level total (mass, calories, or macronutrients); the lack of *per-food identity and portion size* makes them insufficient for dietary assessment and nutrition research, because critical diet-quality scores such as the Healthy Eating Index (HEI) and specific micronutrients (e.g., iron, sodium, vitamins) cannot be recovered from one aggregate.

Multimodal LLMs (MLLMs), trained on web-scale image-text data, recognize a wide range of foods zero-shot in uncontrolled, real-world photos. This open-vocabulary recognition is exactly what a real-life dietary tool requires, and MLLMs have already been adopted in modern nutrition apps for image-based food recognition and nutrition estimation. However,

MLLMs have been shown to be weak at fine-grained portion estimation, producing large discrepancies in energy and nutrient estimates [6].

We propose to **enhance a commercial MLLM with a dedicated portion head**—keeping its open-vocabulary recognition while making the portion estimates accurate enough for nutrition research on real-life images. Our head treats the MLLM as a frozen upstream detector: for each detected food it returns a name, a bounding box, and a plausible density range. The head then reasons jointly over all detected foods and the image so the foods coordinate to partition shared regions of the plate. From this it estimates each food’s mass through structured prediction. The design uses no depth sensor and does not fine-tune the MLLM itself.

### Contributions.

1. **A geometry-enhanced head that endows a frozen MLLM with accurate portion estimation.** Two components are central: an *image $\times$ food cross-attention transformer* that jointly reasons over all detected foods and the image so foods coordinate rather than double-count, and a *structured softmax-ownership volume* that turns this into per-food mass. It substantially lowers per-food error relative to the MLLM alone, using no depth sensor and no MLLM fine-tuning.
2. **A real-life system suitable for large-scale dietary assessment.** It extracts a broad range of foods from real-world meal photos and estimates per-food portion, which can be mapped to canonical food databases for high-fidelity nutrition information (e.g., macronutrients, micronutrients, and diet-quality scores).

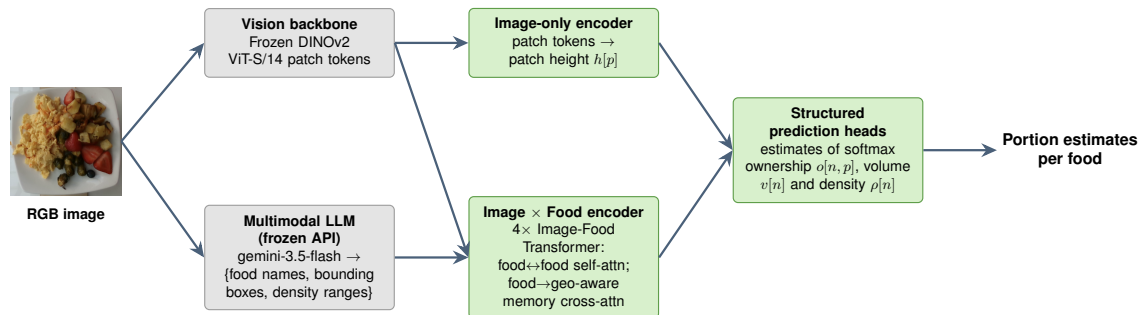
In the remainder of the paper we discuss related work, detail the model and its training, and report its performance on three public food benchmarks.

## 2 Related Work

Image-based food recognition and portion estimation has received considerable attention from both the nutrition-research and computer-vision communities.

**Nutrition Research and Apps** Image-assisted and image-based dietary assessment methods have been developed over the last decade [1]. In recent years, AI-based methods, including convolutional neural networks (CNNs) and multimodal LLMs (MLLMs), have been proposed and evaluated [10, 15]. Although image-based AI has shown great promise, large and inconsistent errors remain a major barrier to its wide adoption in population-scale dietary assessment. Consumer nutrition apps, such as MyFitnessPal and Noom, have adopted image-based food tracking, but a substantial quality gap still separates them from research-grade dietary assessment [6].

**Computer Vision** Recovering mass or volume from a food image is the central hard sub-problem, because a 2D photo discards the 3D information that determines how much food is present. Early monocular approaches regress portion or energy directly from a single RGB image, coupling food classification with portion regression in a multi-task



**Figure 1 | Geometry-enhanced MLLM architecture.** A frozen MLLM (Gemini-3.5-Flash) extracts per-food name, bounding box, and density range from the RGB image; a frozen DINOv2 ViT-S/14 supplies patch tokens. The image $\times$ food encoder and an image-only encoder feed structured heads that produce an exclusive per-patch ownership partition, a per-food density anchored to the MLLM range, and a per-patch height field. Volume estimate for food  $n$  is the sum across all patches  $p$ :  $v_n = \sum_p o_{n,p} h_p$  and per food mass estimate is  $m_n = v_n \cdot \rho_n$ .

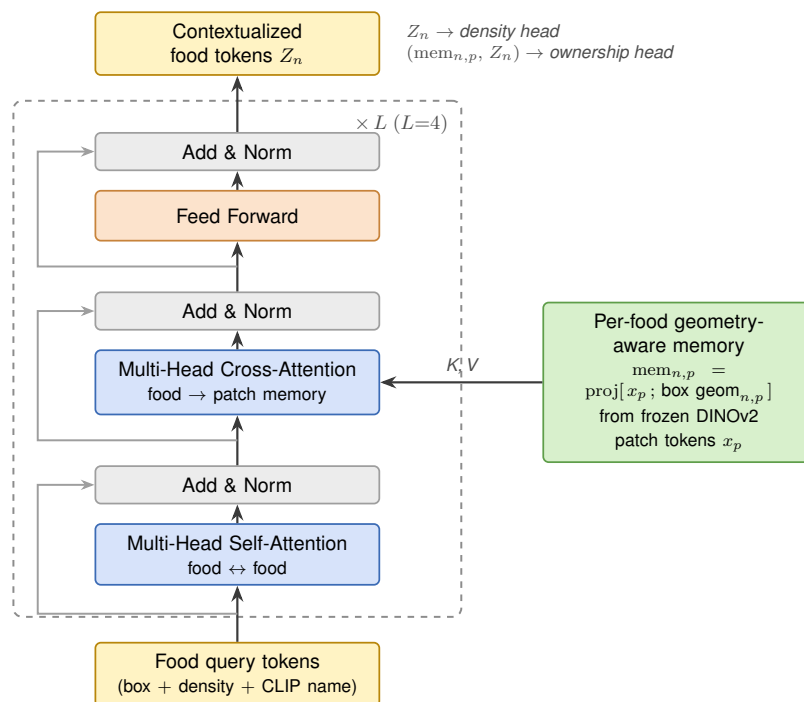
network [5]. To restore the missing geometry, a dominant line adds depth: Nutrition5k [12] supplies a sensor depth map per dish, and DPF-Nutrition [4] predicts depth from the RGB image and fuses it for nutrition estimation. A more recent line makes geometry explicit, reconstructing a 3D food model or point cloud and inferring scale from physical references or contextual objects [3, 14]. These geometry-based methods raise accuracy but depend on depth sensors, multi-view capture, or fragile monocular 3D reconstruction, and operate over a closed food vocabulary. In contrast, we keep a single RGB image with no explicit depth, learning a structured volume from frozen appearance features while an MLLM supplies open-vocabulary food identity.

### 3 Method

Our system pairs a frozen commercial MLLM with a small, trainable geometry-enhanced head. The MLLM supplies open-vocabulary semantics—per-food name, bounding box, and density range—while the head supplies the geometry that the MLLM lacks, turning appearance features into a calibrated per-food portion estimate. Figure 1 gives the overall data flow.

#### 3.1 Multimodal LLM and Vision Backbone

The MLLM is queried as a frozen API with no fine-tuning. We benchmarked five flagship multimodal models and adopt Gemini-3.5-Flash, which pairs strong recognition with the best direct portion accuracy among the flagships we benchmarked (Section 4.5). For each of the  $N$  detected foods it returns the **name**, a **bounding box** enclosing all of that food’s pieces, and a **density range**  $[\rho_{\min}, \rho_{\max}]$ . The name is embedded with a CLIP text encoder [8] and is the dominant semantic signal; the box and density range act as geometric and physical



**Figure 2 | Image×Food encoder block.** Each of the  $L=4$  layers is a Transformer decoder-style block: multi-head *self*-attention over the food query tokens (food↔food), multi-head *cross*-attention from the foods into the per-food geometry-aware patch memory  $\text{mem}_{n,p}$  (keys/values, built once from the frozen DINOv2 patch tokens and per-food box geometry), and a feed-forward sublayer, each wrapped by a residual connection and layer normalization. The stack outputs contextualized per-food tokens  $Z_n$ : the density head reads  $Z_n$  alone, while the ownership head scores each  $(\text{mem}_{n,p}, Z_n)$  pair.

priors.

The vision backbone is a frozen DINOv2 ViT-S/14 [7]. We resize every image—across all datasets—to a fixed 336 px square; with the ViT’s 14 px patch this yields a  $24 \times 24 = 576$  patch-token grid  $\{x_p\}_{p=1}^P$  ( $P=576$ ), cached for speed. Both encoders below read these same frozen tokens; neither the backbone nor the MLLM is updated.

### 3.2 Image×Food Encoder

This encoder produces, for every detected food, a contextualized food embedding and, jointly with a per-food patch memory, the evidence the ownership and density heads consume. Each of its layers is a standard Transformer decoder-style block (Figure 2).

**Inputs.** (i)  $N$  food query tokens  $\{q_n\}$ , each formed by projecting the concatenation of that food’s box features, density features, and CLIP name embedding into the model width  $d$ —the shared hidden dimension of all head components ( $d=256$  throughout); (ii) a *per-food geometry-aware memory*: for food  $n$  and patch  $p$ , the entry  $\text{mem}_{n,p}$  projects the frozen DINOv2

patch token  $x_p$ , concatenated with geometry features describing where  $p$  falls relative to food  $n$ 's bounding box, into width  $d$  (plus a positional embedding). The memory is built once from the  $P$  frozen patch tokens and serves as the keys and values of every layer's cross-attention.

**Layer structure.** The queries pass through  $L=4$  stacked layers. Each layer applies three sublayers, each wrapped with a residual connection and layer normalization:

- **Food↔food self-attention**, where the queries, keys, and values are all the food tokens. This lets foods coordinate so they can partition shared regions of the plate rather than each claiming the same pixels.
- **Food→patch cross-attention**, where food  $n$ 's token attends into its own memory slice  $\text{mem}_{n,\cdot}$ . Crucially the geometry is placed in the attention *values*, not only as an additive attention bias: a bias would merely steer *where* a food looks, whereas putting geometry in the values lets the model read it as *content* and reason over it. We find this choice is load-bearing.
- **A position-wise feed-forward network** over each food token.

**Outputs.** A contextualized per-food token  $Z_n$ . The density head reads  $Z_n$  directly; the ownership head scores the pair  $(\text{mem}_{n,p}, Z_n)$ —how well food  $n$ 's contextualized token matches its geometry-aware view of patch  $p$ —to produce the ownership logits of Section 3.4. The memory is computed once by this encoder and reused unchanged by the ownership head.

### 3.3 Image-only Encoder

This encoder predicts a single positive scalar per patch—the “pile-height” field  $h_p$ —from appearance alone, independent of which foods are present.

**Structure.** The  $P$  DINOv2 patch tokens are linearly projected to width  $d$  and passed through a shallow *patch self-attention* encoder (one transformer layer, queries/keys/values all the projected patch tokens), so each patch can use its spatial neighborhood—a patch in the middle of a tall pile versus at a thin edge. A per-patch MLP then maps each token to one scalar, and a softplus<sup>1</sup> yields the height  $h_p > 0$ .

**Input / output.** Input: the  $P$  frozen patch tokens only (no food tokens, no box, no name). Output: a field  $\{h_p\}$  on the  $24 \times 24$  patch grid. Because it does not depend on the food set, it is a reusable estimate of how much the surface rises above the plate at each patch, shared by all foods when the volume is integrated.

<sup>1</sup>softplus( $x$ ) =  $\log(1 + e^x)$ , a smooth positive-valued activation, so  $h_p > 0$ .

$h_p$  is a **latent height, not predicted depth**. Although  $h_p$  plays the role that a depth map plays in depth-based portion methods, it is deliberately *not* predicted depth. It carries no depth supervision: there is no sensor target and no depth loss anywhere in training, and  $h_p$  is learned end-to-end *only* through the downstream mass objective (Section 3.5)—it becomes whatever makes the integrated volume, times density, match the weighed grams. It is therefore a unitless, latent per-patch weight rather than a calibrated metric surface, and it is defined at the coarse  $24 \times 24$  *patch* resolution, not the dense per-pixel resolution of monocular depth estimators. This is a design choice, not a limitation: explicit single-view depth—sensor depth, pseudo-depth, and depth-supervised height heads—is a tested negative in our setting, whereas a latent height learned under mass supervision keeps the model depth-free while still recovering the missing vertical dimension.

### 3.4 Structured Prediction Heads

The heads take three inputs from the encoders—the contextualized per-food tokens  $Z_n$  and the per-food, per-patch memory  $\text{mem}_{n,p}$  from the image $\times$ food encoder (Section 3.2), and the per-patch height  $h_p$  from the image-only encoder (Section 3.3)—and turn them into a per-food volume, density, and mass.

**Ownership.** Ownership answers, for each patch, “what fraction of this patch belongs to each food?”. We first form a *claim logit* for every food–patch pair and a background logit for every patch,

$$z_{n,p} = \text{MLP}_{\text{own}}([\text{mem}_{n,p}; Z_n]), \quad z_{\text{bg},p} = \text{MLP}_{\text{bg}}(x_p), \quad (1)$$

where  $z_{n,p}$  scores how strongly food  $n$  claims patch  $p$ —from the match between its token  $Z_n$  and the geometry-aware memory  $\text{mem}_{n,p}$ —and  $z_{\text{bg},p}$  scores how much patch  $p$  is plate or empty. A softmax over the  $N$  foods *and* background makes them compete for each patch:

$$o_{n,p} = \frac{\exp(z_{n,p})}{\exp(z_{\text{bg},p}) + \sum_{n'=1}^N \exp(z_{n',p})}, \quad \sum_{n=1}^N o_{n,p} + o_{\text{bg},p} = 1. \quad (2)$$

Hence  $o_{n,p} \in [0, 1]$  is the *share* of patch  $p$  assigned to food  $n$ , and the shares over  $\{1, \dots, N, \text{bg}\}$  partition each patch exactly. This exclusive, soft partition is what makes the downstream volume physical: because a patch’s area cannot be counted twice, foods that overlap must divide the shared region between them, a larger food wins more patches, and background absorbs the plate rather than inflating any food.

**Density.** From each food token  $Z_n$  a per-food log-density is predicted and hinge-anchored to the MLLM’s range  $[\rho_{\min,n}, \rho_{\max,n}]$ , so  $\rho_n$  stays physically plausible while the model calibrates within the range.

**Volume and mass.** The volume of food  $n$  integrates the shared height field over the patches it owns, and mass follows from density:

$$v_n = \sum_p o_{n,p} h_p, \quad m_n = v_n \cdot \rho_n, \quad M_{\text{dish}} = \sum_n m_n. \quad (3)$$

The dish total is simply the sum of per-food masses. We deliberately omit a learned dish-total residual head: it helps dish-total accuracy but overfits the per-food objective.

### 3.5 Training

Queries are **MLLM detections, not ground-truth foods**, and the same detections are used at training and test time; ground-truth (GT) per-food mass is the only supervision target, and no oracle food list is supplied (the per-dataset recognition setup is described in Section 4.1).

**Match-and-mask supervision.** Because detections need not align with GT ingredients, we align detections to GT foods one-to-one—by fuzzy name matching on Nutrition5k, whose ingredients have no ground-truth boxes, and by token- and descriptor-based name matching gated with box IoU where GT boxes exist (FPB, NV-Real). We call the resulting scheme *match-and-mask*: it is the matching-based supervision of set-prediction models such as DETR [2], except that unmatched queries are *masked out* of the loss rather than pushed toward a “no object” target. Concretely, the per-food mass loss is applied *only on matched detections*, while the density hinge anchors every detection (matched or not) to a plausible density; unmatched detections remain as no-loss queries whose appearance-grounded mass is still summed at inference. The backbone stays frozen and we optimize with AdamW; the structured-volume model is warm-started from a simpler direct scalar-volume model (identical to the no-structured-ownership ablation in Section 4), since training the structured volume from scratch is less stable.

**Loss.** The per-food mass objective sums two Smooth-L1 terms with complementary roles, plus the density hinge:

$$\mathcal{L} = \sum_{n \in \text{matched}} \left[ \text{SmoothL1}(\log m_n, \log g_n) + \alpha \text{SmoothL1}(m_n, g_n) \right] + \lambda_\rho \sum_n \text{hinge}(\rho_n, [\rho_{\min, n}, \rho_{\max, n}]), \quad (4)$$

where  $g_n$  is the GT mass,  $m_n$  the predicted mass (Section 3.4),  $\alpha$  the linear-gram weight, and  $\lambda_\rho$  the density-hinge weight. The two mass terms are complementary. The *log-space* term is scale-invariant—it penalizes *relative* error, so a few-gram garnish and a large pile of rice contribute comparable gradients despite spanning two orders of magnitude in grams, which keeps optimization stable across the wide mass range. But scale-invariance also means it under-penalizes *absolute* error on large portions: a 20% miss is only 1 g on a 5 g food yet 60 g on a 300 g one, and the log term weights the two equally. The *linear-gram* term adds a penalty measured directly in grams, restoring pressure on the heavy-portion under-prediction that the log term alone tolerates. We keep  $\alpha$  small—enough linear signal to correct large misses, but not so much that the noisy heavy-portion tail dominates and hurts the many small and medium foods.

## 4 Experiments

We show that a lightweight geometry head on a frozen commercial MLLM substantially improves per-food portion accuracy over the MLLM’s own gram estimates. We focus on *per-food* accuracy—identity-resolved mass—because that is what dietary assessment needs. We evaluate on three real-world datasets (Section 4.1), compare per-food accuracy against the MLLM alone (Section 4.2), compare our portion estimates against each dataset’s originally reported result (Section 4.3), ablate the model’s components on Nutrition5k (Section 4.4), and compare flagship MLLMs as direct portion estimators (Section 4.5).

### 4.1 Datasets

We use three public datasets that provide real meal images with ground-truth *per-food* portions, so the method is exercised beyond any single source.

1. **Nutrition5k** [12] is our primary controlled benchmark: ~5,000 cafeteria dishes with continuous, scale-weighted per-ingredient masses, of which ~3,500 have overhead captures from a fixed RGB-D rig (the subset we use). Its single overhead viewpoint makes it the canonical setting for the footprint-to-mass ambiguity we target; we evaluate on a sealed held-out test set of 502 dishes.
2. **FPB** [9] (Food Portion Benchmark) is a large set of 14,083 photos over 138 food classes of Central Asian cuisine with manually annotated boxes and laboratory-measured component weights. Each dish is photographed from *multiple viewing angles* (a top-down view plus four side angles); we train and test on this mixed-angle image pool rather than a single fixed view, which stresses generalization to varied capture geometry and unfamiliar cuisines. We evaluate on its 2,197-image test split.
3. **NutritionVerse-Real (NV-Real)** [11] is a real consumer-photography set of 889 hand-collected smartphone images spanning 251 dishes, every ingredient individually weighed. Images are taken freehand from random angles, with *no fixed camera-to-food distance* and no scale reference, so it reflects the realistic, uncontrolled phone-capture setting a deployed dietary tool must handle. We evaluate on its 265-image test split.

Although Nutrition5k ships with sensor depth, we deliberately ignore it: such sensors are rarely available in practice, so all our results are image-only (RGB), no depth.

**Protocol.** On each dataset we train the same head—identical architecture, losses, and hyper-parameters—on that dataset’s training split and evaluate on its held-out test split; no test image is used for training or model selection. Recognition is fully open-vocabulary on every dataset: the MLLM names foods freely and never sees a class menu or a per-dish food list, in training or testing. For FPB, whose dishes are regionally unfamiliar, we add a brief cuisine context to the prompt (that the images show Central Asian / Kazakh / Russian / international cafeteria food from a catering service in Kazakhstan). This is a general hint, not a class menu or per-dish food list, so recognition stays open-vocabulary.

## 4.2 Main Results: Portion Estimates Versus MLLMs

We compare per-food portion accuracy—our model versus the MLLM’s own gram estimates—using per-food Mean Absolute Error (MAE) and percentage MAE (PMAE). Error is measured only over matched detection–ground-truth food pairs  $\mathcal{M}$  (see the recognition note below):

$$\text{MAE} = \frac{1}{|\mathcal{M}|} \sum_{(i,j) \in \mathcal{M}} |\hat{m}_i - m_j^*|, \quad \text{PMAE} = 100 \frac{\text{MAE}}{\overline{m^*}},$$

where  $\hat{m}_i$  is a predicted per-food mass,  $m_j^*$  the matched ground-truth mass, and  $\overline{m^*}$  the mean of the matched ground-truth masses  $\{m_j^*\}_{j \in \mathcal{M}}$  (unmatched foods enter neither term). For statistical stability we run our model with three seeds and report the average (Table 1).

**Table 1 | Per-food MAE/PMAE: our model vs. the MLLM alone.** Lower is better. Our model is fully open-vocabulary on every dataset.

Dataset	MLLM alone	Ours	MAE reduction
Nutrition5k	28.1 g / 36.1 %	<b>16.7 g</b> / 21.4 %	−41%
FPB	66.9 g / 31.8 %	<b>43.6 g</b> / 19.3 %	−35%
NV-Real	41.0 g / 33.1 %	<b>27.4 g</b> / 21.0 %	−33%

The geometry head cuts per-food portion error by 33–41 % relative to the MLLM alone, and the gain holds across controlled (Nutrition5k, −41%) and uncontrolled real-world data (FPB, −35%; NV-Real, −33%).

**Recognition coverage (mass-weighted).** Because per-food MAE is defined only on matched foods, missed and spurious detections are excluded; recognition quality therefore bounds *coverage* but is orthogonal to the portion head. Measured as mass-weighted recall, our open-vocabulary recognition covers 87% on Nutrition5k, 88% on NV-Real, and 64% on FPB. FPB’s lower coverage is a recognition-*identity* limit, not a weighing one: without a menu the MLLM *describes* unfamiliar regional dishes instead of naming them (chak-chak becomes “puffed rice cereal”; orama nan becomes “steamed dumplings”).

## 4.3 Portion Estimates Versus the Original Papers

The three source papers report accuracy at different granularities, so we compare our model to each on the metric that paper reports (Table 2), all in the image-only, no-depth setting. The Nutrition5k and NV-Real papers predict only a *dish total*—an image backbone (InceptionV2 and Inception-ResNet, respectively) regresses total dish mass directly, with no per-food breakdown—so for these we compare our dish total, the sum of our per-food masses, against their reported total. The FPB paper is different: its YOLOv12 model regresses a weight per detected food item and reports a per-food MAE, so for FPB we compare at the per-food level instead. We stay depth-free and never regress the total directly; for reference the Nutrition5k paper’s depth/volume model reaches 29.4 g, but it uses depth.

**Table 2 | Portion estimates vs. the original papers** (image-only, no depth). Each dataset is compared on the granularity its paper reports: *dish total MAE* for Nutrition5k and NV-Real, and *per-food MAE* for FPB. Lower is better.

Dataset	Metric	Original paper	Ours
Nutrition5k	dish total MAE	40.4 g / 18.8 %	<b>36.6 g</b> / 18.4 %
NV-Real	dish total MAE	115.9 g / 27.2 %	<b>68.6 g</b> / 16.1 %
FPB	per-food MAE	90.95 g / 40.3 %	<b>43.6 g</b> / 19.3 %

#### 4.4 Ablations

We remove one component at a time from the full model and report per-food MAE, broken down by single-food and multi-food dishes (Table 3); this ablation study is conducted on Nutrition5k only. **(1) No MLLM context:** zeros the per-food name, density, and box geometry, leaving only DINOv2 patches and a generic query, isolating the value of the MLLM priors. **(2) No structured ownership:** replaces the ownership-height volume with direct mass regression, testing whether the exclusive patch partition is load-bearing on mass. **(3) No image-only height:** sets  $h_p=1$ , collapsing volume to footprint area.

**Table 3 | Ablations on Nutrition5k.** One component removed at a time from the full model; per-food MAE (grams, three-seed average) on the sealed held-out test, broken down by single-food dishes, multi-food dishes, and all dishes (total), each with  $\Delta$  vs. the full model. Lower is better.

Variant	Single		Multi		Total	
	MAE	$\Delta$	MAE	$\Delta$	MAE	$\Delta$
Full model	<b>14.2 g</b>	—	<b>17.1 g</b>	—	<b>16.7 g</b>	—
(1) No MLLM context (no name/density/box)	20.6 g	+6.4 g	36.7 g	+19.6 g	34.4 g	+17.7 g
(2) No structured ownership (direct mass regr.)	16.4 g	+2.3 g	18.6 g	+1.5 g	18.3 g	+1.6 g
(3) No image-only height ( $h_p=1$ )	15.4 g	+1.2 g	18.3 g	+1.2 g	17.9 g	+1.2 g

The MLLM per-food context is by far the most load-bearing component: removing name, box, and density more than *doubles* per-food error (16.7 g  $\rightarrow$  34.4 g, +106%). For both single-food and multi-food dishes, the errors go up greatly without per-food context. The error is especially large on multi-food dishes, because without per-food priors the model can only regress toward the average mass of the foods on the plate. The structured softmax-ownership partition (+1.6 g, +10%) and the learned image-only height field (+1.2 g, +7%) each add a smaller, consistent gain across both single- and multi-food dishes.

#### 4.5 Which MLLM? A Flagship Comparison

Our head can sit on any MLLM, so we benchmarked five flagships as *direct* portion estimators—each recognizes foods and estimates per-food grams itself, open-vocabulary, with the same decomposition prompt and matcher—on 100 multi-food dishes ( $\geq 2$  significant

foods, the hard regime of Table 3) from the sealed Nutrition5k test. Each model runs at the minimum reasoning effort it allows, with temperature 0 where supported; all models were accessed through their public APIs in July 2026.

**Table 4 | Flagship MLLMs as direct portion estimators.** Each model recognizes foods and estimates per-food grams directly (open-vocabulary, same decomposition prompt and matcher) on 100 multi-food dishes. Mass cov. = mass-weighted recall; higher is better. Per-food MAE: lower is better. Our model inherits Gemini-3.5-Flash recognition through the pipeline’s extraction prompt (name, box, density—no gram estimation), so its coverage differs slightly from the Gemini-3.5-Flash row, whose detections come from the direct-estimation prompt: same model, different prompt.

Model	Mass cov.	Per-food MAE
Gemini-3.5-Flash	68.9%	24.1 g / 34.7 %
Gemini-3.1-Pro	68.6%	25.8 g / 36.6 %
GPT-5.5	61.4%	26.7 g / 38.7 %
Claude-Fable-5	<b>71.6%</b>	35.7 g / 51.3 %
Claude-Opus-4.8	66.8%	36.2 g / 52.9 %
Our model (head on Gemini-3.5-Flash)	71.1%	<b>16.5 g / 23.7 %</b>

Results are shown in Table 4. Three observations. **(1) Recognition is broadly comparable.** On these hard cases, mass-weighted coverage clusters at 61–72% across all five models, with Claude-Fable-5 highest at 71.6%. **(2) Portion estimation splits the field.** The Gemini models and GPT-5.5 sit at 24.1 g–26.7 g, while both Claude models are near 36 g; the Claude models under-estimate large piles less but pay with larger errors on ordinary portions—a different bias profile, not a uniformly weaker one. **(3) A small trained head beats every flagship.** Our portion estimation head, fed Gemini-3.5-Flash recognition, reaches 16.5 g per-food MAE on the same 100 dishes—32% below the best flagship direct estimate.

## 5 Limitations

Four limitations bound the current system, in decreasing order of impact. First, **single-view portion is fundamentally ambiguous**: a footprint does not determine pile height, so the heavy-portion tail is systematically under-predicted. This is a property of the input, not of the head—we observe it consistently across datasets. Second, **image-only recognition cannot capture invisible foods**: ingredients that leave little or no visual trace—cooking oil, added sugar, salt, sauces mixed into a dish—are missed by any image-based recognizer, yet they carry real mass and a disproportionate share of the nutrients dietary assessment cares about. Third, **the head is trained per dataset**: the architecture, losses, and hyper-parameters are identical everywhere, but each benchmark uses its own weights, and cross-dataset transfer of a single universal weigher has not been demonstrated. Fourth, **we validate mass, not nutrients**: the per-food masses are designed for downstream nutrient conversion, but that conversion is not evaluated here.

## 6 Conclusion and Future Directions

Food *recognition* is now largely an MLLM capability; food *weighing* is not—every flagship we tested under-performs at portion estimation, and larger reasoning models do not close the gap. Our results show this gap is closed not by scale but by a small, task-specific geometry head: consuming only the MLLM’s structured output and frozen DINOv2 features, it cuts per-food portion error by 33–41 % relative to the MLLM alone, beats every flagship’s direct estimates, and surpasses each benchmark’s originally published image-only model—fully open-vocabulary, depth-free, and without touching the MLLM. Recognition and weighing thus decouple cleanly: the MLLM supplies identity and priors, the head supplies geometry.

The end goal is not a better plate estimator but a replacement for the 24-hour dietary recall—the instrument nutrition research actually runs on—and framing future work against that target makes the remaining gaps concrete. **(1) From served portion to consumed intake.** A recall measures what was *eaten*; an image measures what was *served*. Closing the difference requires paired before/after capture (or short video) with leftover subtraction—a natural extension of a per-food weigher, since intake is the difference of two portion estimates over matched foods. **(2) From visible foods to complete intake.** A recall interviewer *probes*: cooking oil, sugar in coffee, sauces mixed in, the composition of a casserole. No image model can see these; the MLLM’s dietary knowledge, used conversationally (brief follow-up questions, recipe inference from dish identity), can substitute for the probe—turning the invisible-ingredient limit from a silent bias into a bounded, queryable one. **(3) From plates to days.** A recall covers all eating occasions; a deployed system must handle missed captures, snacks, and drinks across a full day with one universal weigher (our identical-recipe result suggests one is trainable) and graceful degradation when images are absent. **(4) From benchmark agreement to clinical validity.** Replacing the recall ultimately requires validation the way recalls themselves were validated: against recovery biomarkers and controlled-feeding ground truth, on nutrient intakes rather than grams, across diverse cuisines—where our coverage results show current MLLMs still under-serve underrepresented foods. Per-food mass estimation of the accuracy demonstrated here is, we believe, the enabling component that makes this agenda realistic.

## References

- [1] Carol J. Boushey, Melissa Spoden, Fengqing M. Zhu, Edward J. Delp, and Deborah A. Kerr. New mobile methods for dietary assessment: Review of image-assisted and image-based dietary assessment methods. *Proceedings of the Nutrition Society*, 76(3): 283–294, 2017. doi: 10.1017/S0029665116002913.
- [2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 213–229, 2020.
- [3] Yuhao Chen, Gautham Vinod, Siddeshwar Raghavan, Talha Ibn Mahmud, Bruce Coburn, Jinge Ma, Fengqing Zhu, and Jiangpeng He. Implicit-scale 3D reconstruction

- for multi-food volume estimation from monocular images. In *IEEE Southwest Symposium on Image Analysis and Interpretation (SSIAI)*, pages 61–64, 2026.
- [4] Yuzhe Han, Qimin Cheng, Wenjin Wu, and Ziyang Huang. DPF-Nutrition: Food nutrition estimation via depth prediction and fusion. *Foods*, 12(23):4293, 2023. doi: 10.3390/foods12234293.
- [5] Jiangpeng He, Zeman Shao, Janine Wright, Deborah Kerr, Carol Boushey, and Fengqing Zhu. Multi-task image-based dietary assessment for food recognition and portion size estimation. In *IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*, pages 49–54, 2020.
- [6] Xiao Li, Angela Yin, Hyun Yang Choi, Virginia Chan, Margaret Allman-Farinelli, and Juan Chen. Evaluating the quality and comparative validity of manual food logging and artificial intelligence-enabled food image recognition in apps for nutrition care. *Nutrients*, 16(15):2573, 2024. doi: 10.3390/nu16152573.
- [7] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. DINOv2: Learning robust visual features without supervision. *Transactions on Machine Learning Research (TMLR)*, 2024.
- [8] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, pages 8748–8763, 2021.
- [9] Aibota Sanatbyek, Tomiris Rakhimzhanova, Bibinur Nurmanova, Zhuldyz Omarova, Aidana Rakhmankulova, Rustem Orazbayev, Huseyin Atakan Varol, and Mei Yen Chan. A multitask deep learning model for food scene recognition and portion estimation – the Food Portion Benchmark (FPB) Dataset. *IEEE Access*, 2025. doi: 10.1109/ACCESS.2025.3603287.
- [10] Eleanor Shonkoff, Kelly Copeland Cara, Xuechen (Anna) Pei, Mei Chung, Shreyas Kamath, Karen Panetta, and Erin Hennessy. AI-based digital image dietary assessment methods compared to humans and ground truth: A systematic review. *Annals of Medicine*, 55(2):2273497, 2023. doi: 10.1080/07853890.2023.2273497.
- [11] Chi-en Amy Tai, Matthew Keller, Mattie Kerrigan, Yuhao Chen, Saejith Nair, Pengcheng Xi, and Alexander Wong. NutritionVerse: Empirical study of various dietary intake estimation approaches. In *Proceedings of the 8th International Workshop on Multimedia Assisted Dietary Management (MADiMa)*, pages 11–19, 2023. doi: 10.1145/3607828.3617799.
- [12] Quin Thames, Arjun Karpur, Wade Norris, Fangting Xia, Liviu Panait, Tobias Weyand, and Jack Sim. Nutrition5k: Towards automatic nutritional understanding of generic

- food. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8903–8911, 2021.
- [13] Frances E. Thompson and Amy F. Subar. Dietary assessment methodology. *Nutrition in the Prevention and Treatment of Disease*, pages 5–48, 2017.
- [14] Gautham Vinod, Jiangpeng He, Zeman Shao, and Fengqing Zhu. Food portion estimation via 3D object scaling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 3741–3749, 2024.
- [15] Runze Yan, Hanqi Luo, Jiaying Lu, Darren Liu, Hannah Posluszny, Mehak Preet Dhaliwal, Janice MacLeod, Yao Qin, Carl Yang, Terry J. Hartman, and Xiao Hu. DietAI24 as a framework for comprehensive nutrition estimation using multimodal large language models. *Communications Medicine*, 5:450, 2025. doi: 10.1038/s43856-025-01159-0.