



# MIT&AIS-2026

**27-29 квітня 2026**  
**Харків-Яремче 2026**

**2 Міжнародна науково-практична конференція**  
**«Сучасні інформаційні технології**  
**та системи штучного інтелекту»**  
**(MIT&AIS-2026)**

## **Матеріали** **конференції**

**Частина 1**

**Харків 2026**

Міністерство освіти та науки України  
Національна академія наук України  
Координаційна рада НАН України з питань штучного інтелекту  
Харківський національний університет радіоелектроніки  
Харківський національний університет імені В.Н. Каразіна  
Карпатський національний університет імені Василя Стефаника  
Північно-Східний координаційний науковий центр з питань штучного інтелекту  
Інститут кібернетики імені В.М. Глушкова НАН України  
Університет технологій в Лодзі  
Університет Павла Йозефа Шафарика в Кошице

# **СУЧАСНІ ІНФОРМАЦІЙНІ ТЕХНОЛОГІЇ ТА СИСТЕМИ ШТУЧНОГО ІНТЕЛЕКТУ MIT&AIS-2026**

## **Матеріали**

**2-ї Міжнародної науково-практичної конференції**

**Частина 1**

**27-29 квітня 2026 р.  
Харків – Яремче, Україна**

**Харків 2026**

# Mutual Information Preference Optimization for Robust Multi-Modal Recipe Generation

Maksym Shaposhnyk <sup>a</sup> and Serhii Minukhin <sup>a</sup>

<sup>a</sup> *Simon Kuznets Kharkiv National University of Economics, Nauky Ave, 9A, Kharkiv, 61166, Ukraine*

## Abstract

This study evaluates the impact of Mutual Information Preference Optimization (MIPO) as a corrective layer within a hybrid vision-language architecture. Rather than introducing a new standalone framework, the research modifies an existing multimodal pipeline by integrating MIPO to bridge the operational gap between a DenseNet-121 ensemble and Llama 3.1 8B. The central hypothesis—that LLMs can act as autonomous semantic filters—was tested through contrastive alignment, which synchronizes CNN-derived visual features with the textual latent space.

Experimental results on the Food-101 dataset validate this modification, demonstrating that the system can successfully suppress false-positive detections without a complete retraining of the visual backbone. By filtering out incongruous artifacts through preference optimization, the modified architecture achieved a 60,8% reduction in semantic hallucinations. This confirms the viability of using LLMs for real-time error correction in specialized domains, such as personalized dietetics, where output fidelity is a critical requirement.

## Keywords

Food Computing; Llama 3.1; DenseNet-121; Mutual Information Preference Optimization; visual grounding.

## 1. Introduction

The evolution of Food Computing has reached a critical juncture, transitioning from rudimentary image classification to the sophisticated, autonomous synthesis of culinary content [1, 2]. However, the integration of Computer Vision (CV) and Large Language Models (LLMs) remains hampered by a persistent challenge: the propagation of errors across modalities. In typical hybrid architectures, the LLM acts as a passive recipient of visual tokens derived from CNN backbones [3]. When these visual features are noisy or misclassified, the language model—lacking an inherent verification mechanism—extrapolates these errors into semantic hallucinations [1, 3, 4].

This phenomenon is particularly detrimental in the context of personalized nutrition, where the inclusion of a "hallucinated" ingredient can render a recipe not only useless but potentially hazardous for users with strict dietary constraints.

The core of this research lies in testing the hypothesis that the "modality gap" can be bridged without a resource-intensive retraining of the visual module. A modification of the standard hybrid pipeline by implementing Mutual Information Preference Optimization (MIPO) [5]. The objective is to transform the LLM from a passive generator into an active semantic filter. By maximizing the conditional mutual information between the visual context and the generated text [5], the system learns to prioritize contextually grounded ingredients over stochastic artifacts.

## 2. Problem Analysis

The primary limitation of contemporary vision-to-recipe pipelines lies in the asymmetric reliability of their constituent modules. While convolutional neural networks (CNNs), such as DenseNet-121 [6], exhibit high recall in broad object detection, they are inherently susceptible to textural ambiguity—

where visual similarities between disparate ingredients (e.g., the visual grain of "fish" vs. the texture of "onions") trigger false-positive classifications [1, 7].

In a standard hybrid configuration, the Large Language Model (LLM) operates as a stochastic decoder of these visual tokens. Because the interface between the CNN and the LLM is typically unidirectional and lacks a verification layer, the LLM is forced to reconcile contradictory or erroneous visual inputs within a coherent narrative. This lack of semantic agency results in "hallucinations" [8]: the model generates instructions for ingredients that are visually suggested but contextually illogical [4].

Empirical observation suggests that even when a CNN backbone achieves a high Recall (e.g., 0.91), the resulting Precision at the recipe level remains suboptimal. This discrepancy occurs because the LLM lacks the mechanism to challenge the visual input. To address this, a shift from Passive Translation to Active Semantic Filtering is justified.

### 3. Objectives

The primary aim of this research is to adapt the hybrid architecture originally proposed in [1] by evolving it into a system that uses a Large Language Model (LLM) as a "semantic filter." This modification is designed to autonomously rectify visual classification errors and ensure high grounding (contextual fidelity) in generated recipes, addressing the reliability gaps identified in [1, 4].

To achieve this aim, the following objectives must be addressed:

1. Synthesize contrastive pairs via hallucination simulation.
2. Formulate the MIPO loss function.
3. Execute supervised fine-tuning (SFT) of the Llama 3.1 8B model [9].

### 4. Methodology and Empirical Results

The input image is processed through two parallel streams based on the DenseNet-121 architecture. DenseNet-121 was selected as the visual backbone due to its high efficiency in feature propagation and its ability to mitigate the vanishing-gradient problem through dense connectivity.

The first stream is dedicated to global scene recognition. It classifies the input into one of 101 distinct dish categories. This branch utilizes a Softmax activation layer paired with a Categorical Cross-Entropy loss function. The primary objective of this stream is to provide a high-level semantic anchor (e.g., "Pancakes") to guide subsequent recipe generation.

The second stream operates in parallel to identify a set of specific ingredients within the image. Unlike the classification branch, this stream treats ingredient detection as a multi-label classification problem. To account for the simultaneous presence of multiple components, a Sigmoid activation function is applied. This allows the model to assign independent probabilities to each potential ingredient, facilitating the detection of complex compositions.

The core mechanism of Mutual Information Preference Optimization (MIPO) is the systematic construction of a preference dataset  $D$ , formulated as follows [4] (1):

$$D = \{(x, y_w, y_l)\}, \quad (1)$$

where  $x$  —the input prompt containing the "noisy" ingredient list generated by the CNN backbone. This represents the raw, unverified visual detections that may contain errors,  $y_w$  — the reference (ground-truth) recipe. This serves as the target sequence, characterized by high semantic grounding and accurate ingredient representation.,  $y_l$  — the hallucinated recipe. This sequence is intentionally constructed to include contextually incongruous ingredients (hallucinations) derived from the CNN's false-positive detections.

The model optimization is performed by minimizing the following MIPO loss function (2), as established in [5,8] (2):

$$L_{MIPO}(\pi_\theta; \pi_{ref}) = -E_{(x, y_w, y_l) \sim D} \left[ \log \sigma \left( \beta \log \frac{\pi_\theta(y_w | x)}{\pi_{ref}(y_w | x)} - \beta \log \frac{\pi_\theta(y_l | x)}{\pi_{ref}(y_l | x)} \right) \right], \quad (2)$$

where  $\pi_\theta$  — denotes the policy (probability distribution) of the Llama 3.1 8B model being fine-tuned,  $\pi_{ref}$  — represents the frozen reference model, ensuring the optimization does not deviate excessively from the base linguistic capabilities,  $\beta$  — a hyperparameter that controls the strength of the penalty for deviations from the reference model, effectively modulating the "conservatism" of the semantic filter.

The experimental validation was conducted using the Food-101 [10] dataset. The training phase for the visual modules (the DenseNet-121 classification and recognition streams) required 8 hours of computation on an NVIDIA RTX 3090 GPU.

To implement the MIPO mechanism, Low-Rank Adaptation (LoRA) was employed. This parameter-efficient fine-tuning technique enabled the adaptation of the Llama 3.1 8B model in just 2 hours, demonstrating the framework's efficiency under constrained computational resources [1, 4]. By freezing the original weights and injecting trainable rank-decomposition matrices, we maintained the model's vast linguistic foundation while specializing its "reasoning" for culinary semantic filtering. This optimization ensures high-fidelity hallucination pruning without the catastrophic forgetting typically associated with full-parameter retraining.

The comparative analysis between the baseline hybrid architecture (CNN + LLM) and the proposed MIPO-enhanced hybrid system is presented in Tables 1 and 2, highlighting the shift toward superior factual grounding.

**Table 1**  
Quantitative Performance Comparison on Food-101 with visual grounding.

Category	Image ID	BLEU	ROGUE-1	ROGUE-2	ROGUE-L	Cosine similarity	Ingredient hallucination count
pancakes	2606759.jpg	0.162	0.583	0.260	0.397	0.897	5
tacos	2937937.jpg	0.07	0.501	0.117	0.258	0.89	4
fried_calamari	2262574.jpg	0.07	0.44	0.111	0.259	0.89	2
spaghetti_carbonara	3542036.jpg	0.128	0.523	0.157	0.26	0.798	7
sashimi	2787824.jpg	0.001	0.149	0.025	0.08	0.79	3
bread_pudding	3337424.jpg	0.053	0.412	0.142	0.263	0.781	11
frozen_yogurt	3456727.jpg	0.036	0.405	0.12	0.222	0.781	1
crème_brulee	2319223.jpg	0.004	0.274	0.086	0.148	0.757	2
cheese_plate	3430261.jpg	0.035	0.372	0.034	0.162	0.791	1
lobster_bisque	3301541.jpg	0.181	0.569	0.212	0.319	0.884	4

**Table 2**  
Quantitative Performance Comparison on Food-101 with visual grounding and MIPO.

Category	Image ID	BLEU	ROGUE-1	ROGUE-2	ROGUE-L	Cosine similarity	Ingredient hallucination count
pancakes	2606759.jpg	0.152	0.625	0.214	0.439	0.832	0
tacos	2937937.jpg	0.011	0.3031	0.071	0.177	0.718	2
fried_calamari	2262574.jpg	0.077	0.499	0.164	0.299	0.887	2
spaghetti_carbonara	3542036.jpg	0.077	0.48	0.171	0.277	0.696	0
sashimi	2787824.jpg	0.048	0.345	0.07	0.241	0.838	1
bread_pudding	3337424.jpg	0.069	0.441	0.131	0.281	0.648	6
frozen_yogurt	3456727.jpg	0.039	0.400	0.096	0.242	0.812	1
crème_brulee	2319223.jpg	0.051	0.459	0.127	0.224	0.728	3
cheese_plate	3430261.jpg	0.013	0.339	0.042	0.174	0.758	1
lobster_bisque	3301541.jpg	0.018	0.371	0.134	0.182	0.776	2

The most critical metric in this study—Ingredient Hallucination Count—showed a dramatic reduction across the dataset. The total number of hallucinations across the ten sampled categories dropped from 46 in the baseline system to 18 with MIPO intervention, representing a 60.8% overall

reduction in semantic errors. Specifically, in the "Pancakes" and "Spaghetti Carbonara" scenarios, the MIPO mechanism achieved a zero-hallucination state, effectively pruning all contextually incongruous tokens generated by the visual backbone. While the baseline system frequently included "fish" or "onion" in pancake recipes due to textural similarities detected by the CNN, the MIPO-enhanced Llama 3.1 model successfully identified these as stochastic artifacts and suppressed them in favor of semantically valid ingredients.

A key observation from the data is the stabilization of Cosine Similarity alongside the reduction of hallucinations. In categories such as "Sashimi" and "Frozen Yogurt," an increase in similarity scores (to 0.838 and 0.812, respectively) is observed, signaling a higher degree of semantic grounding. This indicates that preference optimization effectively bridges the multimodal gap.

However, it is noteworthy that in some cases (e.g., "Tacos"), aggressive filtering of hallucinated ingredients led to a slight decrease in linguistic overlap metrics such as BLEU and ROUGE. This phenomenon suggests that the MIPO-filtered system prioritizes biological and culinary logic over literal adherence to the noisy CNN prompt. By sacrificing "incorrect" textual labels, the system ensures that the final recipe is operationally reliable and safe for the user, resolving the hallucination bottleneck that plagues traditional multimodal relay architectures. This emphasizes that semantic precision is more vital than mere lexical imitation.

## 5. Declaration on Generative AI

During the preparation of this manuscript, the authors used generative AI and AI-assisted tools, specifically Grammarly and Reverso Context, to refine grammar, correct orthography, and optimize language. Following the AI-assisted processing, the authors personally reviewed, verified, and edited the content to ensure academic accuracy. The authors maintain full responsibility for the integrity, factual correctness, and original contributions of the final publication.

## 6. Conclusions

The empirical validation of the proposed architecture confirms that semantic correction of visual noise can be successfully performed without modifying the underlying classifier backbones. By integrating the Mutual Information Preference Optimization (MIPO) mechanism, significant improvements in system reliability were achieved, particularly in high-complexity culinary scenes. The most definitive metric of success is the Ingredient Hallucination Count, which saw an aggregate reduction of 60.8% (from 46 to 18 total errors).

The baseline system struggled significantly with high-texture complexity dishes; for instance, "Bread Pudding" initially triggered 11 hallucinated ingredients, likely due to visual ambiguities in the crumb structure being misinterpreted as discrete ingredients. Under MIPO, this was reduced to 6. Even more striking is the achievement of a "zero-hallucination state" in the "Pancakes" and "Spaghetti Carbonara" categories. In these instances, the MIPO-enhanced Llama 3.1 effectively acted as a logic-based gatekeeper, discarding erroneous tokens like "fish" or "onion" that the visual backbone had incorrectly suggested based on textural similarities.

Beyond error reduction, the quality of semantic alignment—measured via Cosine Similarity—demonstrated notable stabilization. While the baseline model occasionally produced high similarity through sheer word volume, MIPO achieved higher precision. In the "Sashimi" category, cosine similarity rose from 0.79 to 0.838, and "Frozen Yogurt" increased from 0.781 to 0.812. This upward trend suggests that when the model is forced to prioritize culinary logic over noisy visual prompts, the resulting output aligns more closely with the true semantic essence of the dish.

A nuanced finding in the data is the divergence between linguistic overlap metrics (BLEU/ROUGE) and semantic accuracy. In the "Tacos" category, we observed a drop in BLEU scores from 0.07 to 0.011. While a traditional interpretation might view this as a performance regression, a granular audit reveals that the MIPO mechanism was aggressively pruning hallucinated ingredients. By sacrificing literal adherence to a noisy prompt (which often contains hallucinated words that may happen to match incorrect reference labels), the system prioritizes operational safety. This confirms that in specialized

domains like automated recipe generation, semantic integrity and the suppression of false positives are more critical than maintaining high linguistic overlap with potentially flawed visual detections.

This confirms the viability of using LLMs for real-time error correction in specialized domains [8].

Furthermore, the successful application of LoRA for this task demonstrates that such high-level alignment can be achieved with minimal computational overhead. Beyond immediate performance gains, this research provides a robust technical solution to the "blind trust" problem inherent in modular multimodal systems. Future work could explore integrating MIPO into continual food learning frameworks [11] to maintain performance as recipe datasets expand.

## 7. References

- [1] S. Minukhin, M. Shaposhnyk, A Hybrid Approach to Visually Oriented Generation of Culinary Recipes Based on Convolutional Neural Networks and Large Language models, Herald of Khmelnytskyi National University. Technical sciences (2026) 418–434. doi:10.31891/2307-5732-2026-363-57
- [2] A. Salvador, et al., Inverse cooking: Recipe generation from food images, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR '19, 2019, pp. 10445–10454. doi: 10.1109/CVPR.2019.01070.
- [3] Z. Ji, et al., Survey of hallucination in natural language generation, ACM Comput. Surv. 55 (2023) 1–38. doi:10.1145/3571730.
- [4] R. Rafailov, et al., Direct preference optimization: Your language model is secretly a reward model, arXiv preprint arXiv:2305.18290, 2023. doi:10.48550/arXiv.2305.
- [5] H. Nam, et al., Maximizing mutual information between user-contexts and responses improve LLM personalization with no additional data, arXiv preprint arXiv:2603.19294, 2026. doi:10.48550/arXiv.2603.19294.
- [6] G. Huang, et al., Densely connected convolutional networks, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, CVPR '17, 2017. doi:10.48550/arXiv.1608.06993.
- [7] W. Min, et al., Large Scale Visual Food Recognition, IEEE Trans. Pattern Anal. Mach. Intell. 45 (2023) 9932–9949. doi:10.1109/TPAMI.2023.3237871.
- [8] C. Jang, Modulated Intervention Preference Optimization (MIPO): Keep the Easy, Refine the Difficult, arXiv preprint arXiv:2409.17545, 2024. doi:10.48550/arXiv.2409.17545.
- [9] S. E. Repede, R. Brad, LLaMA 3 vs. State-of-the-Art Large Language Models: Performance in Detecting Nuanced Fake News, Computers 13 (2024) 292. doi:10.3390/computers13110292.
- [10] L. Bossard, M. Guillaumin, L. Van Gool, Food-101 – Mining Discriminative Components with Random Forests, in: D. Fleet, T. Pajdla, B. Schiele, T. Tuytelaars (Eds.), Computer Vision – ECCV 2014, volume 8694 of Lecture Notes in Computer Science, Springer, Cham, 2014, pp. 446–461. doi:10.1007/978-3-319-10599-4\_29.
- [11] X. Wu, B. Zhu, F. Han, P. Jiao, J. Chen, Dual-LoRA and Quality-Enhanced Pseudo Replay for Multimodal Continual Food Learning, arXiv preprint arXiv:2511.13351, 2025. doi:10.48550/arXiv.2511.13351.