Fine-tuning Language Models for Recipe Generation: A Comparative Analysis and Benchmark Study

Anneketh Vij *, Changhao Liu*, Rahul Anil Nair*, Theodore Eugene Ho*, Edward Shi, Ayan Bhowmick

Department of Computer Science University of Southern California Los Angeles, CA 90007

{anneketh, celiu, ranair, teho, epshi, abhowmic}@usc.edu

Abstract

This research presents an exploration and study of the recipe generation task by fine-tuning various very small language models, with a focus on developing robust evaluation metrics and comparing across different language models the open-ended task of recipe generation. This study presents extensive experiments with multiple model architectures, ranging from T5small (Raffel et al., 2023) and SmolLM-135M (Allal et al., 2024) to Phi-2 (Research, 2023), implementing both traditional NLP metrics and custom domain-specific evaluation metrics. Our novel evaluation framework incorporates recipe-specific metrics for assessing content quality and introduces approaches to allergen substitution. The results indicate that, while larger models generally perform better on standard metrics, the relationship between model size and recipe quality is more nuanced when considering domain-specific metrics. SmolLM-360M and SmolLM-1.7B demonstrate comparable performance despite their size difference before and after fine-tuning, while fine-tuning Phi-2 shows notable limitations in recipe generation despite its larger parameter count. The comprehensive evaluation framework and allergen substitution systems provide valuable insights for future work in recipe generation and broader NLG tasks that require domain expertise and safety considerations.

1 Introduction

The generation of safe and high-quality recipes presents unique challenges in natural language generation. Beyond generating coherent and creative recipes, recipe generation requires high-level knowledge of culinary techniques, nutritional principles, and awareness of dietary restrictions to ensure user safety. This necessitates approaches that balance linguistic fluency with domain-specific expertise, particularly in the domain of allergen substitution

This study focuses on addressing these challenges by experimenting with different model architectures for recipe generation and allergen substitution through controlled fine-tuning and comprehensive evaluation metrics. There are three primary research questions:

- 1. Given the scope of this study, which models will achieve the best results after fine-tuning for recipe generation?
- 2. How should the generated recipes be evaluated to ensure that they are coherent and safe for users with dietary restrictions?
- 3. How should allergen substitution be implemented into large-scale models to achieve high-quality allergen-free recipes?

To answer these questions, this study makes the following contributions:

- Comprehensive comparison of model architectures across scales including smaller models like GPT-2 (Radford et al., 2019) and T5 (Raffel et al., 2020), and larger models like Phi-2 (Research, 2023) and SmolLM-1.7(Allal et al., 2024)
- Multi-dimensional evaluation framework combining novel recipe-specific evaluation metrics, traditional metrics, and LLM-based assessment.
- Development of RAG-assisted and promptbased approaches for allergen substitution

Our work represents a step forward in adapting NLG systems for practical applications in the culinary domain, emphasizing safety, personalization, and quality.

^{*}Authors contributed equally.

2 Related Work

2.1 Language Models in Recipe Generation

Previous works have explored various language model architectures for recipe generation. Lam et al. (2024) explored the performance of BART-based (Lewis et al., 2020) and BRIO-based (Liu et al., 2022b) across different recipe datasets, mostly in English and Vietnamese, which were then evaluated using ROUGE scores. Our work extends these previous ones by systematically comparing models across different sizes and architectures, from smaller models like SmolLM-135M to larger ones like Phi-2. This paper provides a detailed analysis on how model size impacts different aspects of recipe quality.

2.2 Recipe Generation Models

This paper builds on several recent advances in recipe generation and personalization. Majumder et al. (2019) proposed a personalized recipe generation model using attention mechanisms to focus on recipes previously consumed by the user. Their approach showed promising results in generating recipes that aligned with user preferences. The dataset from this paper has been used in this research with their encoder-decoder model being used as a baseline.

Chen et al. (2021) implemented a framework using constrained question answering over a large-scale knowledge graph to recommend food recipes considering users' explicit requirements and health factors. This helped recommend healthy alternatives to users, which aligned with the study's goal of providing allergen-free options and gave us the inspiration for a RAG-based system for allergen substitution.

2.3 Multi-modal Approaches

The FIRE system, by Chhikara et al. (2024) and Nutrify AI by Han and Chen (2024) both use a multimodal approach, generating recipes from food images and ingredients. While it differs from this work due to us not using images, these studies also incorporate different types of input in the process of recipe generation.

LLava-Chef, by Mohbat and Zaki (2024) and is another multi-modal approach to recipe generation, which was fine-tuned on both the cross-entropy loss and a novel loss function computed using BLEU and ROUGE scores to ensure that the model generated recipes that were closer to the ground truth.

This paper adopted these evaluation metrics and the idea of creating custom ones from this paper, as well as what inputs to include for recipe generation. However, this work doesn't use the novel loss function to fine-tune the language models, since penalizing generations for not being closer to the ground truth might hinder the personalization of the generated recipes, which is an important part of allergen substitution.

Other multi-modal recipe generation approaches include ChefFusion by Li et al. (2024) and Inverse Cooking by Salvador et al. (2019). ChefFusion provides complete multimodality by developing a framework for both recipe generation using images and image generation using recipes. This paper, along with LLava-Chef, uses metrics like Sacre-BLEU(Post, 2018) and ROUGE, which wasn't preferred due to the limitations of these metrics for creative generation. Inverse Cooking, which uses encoder-decoder transformer (Vaswani, 2017) blocks to generate recipes from images, provided the inspiration to use Ingredient Coverage as an evaluation metric, which is similar to how the paper evaluates the ingredients extracted from the image and the ground truth.

2.4 Evaluation

Many of these studies, such as LLava-Chef or Retrieval Augmented Recipe Generation by (Liu et al., 2024), use conventional metrics such as BLEU, ROUGE, and F1-score for ingredient matching to assess recipe quality. This paper distinguishes itself by employing both general and domain-specific metrics such as ingredient coverage, which was used by both Liu et al. (2022a) and Salvador et al. (2019) to attain a more profound understanding of the quality implications across many aspects of the generated recipe since traditional metrics focus more on overlap and thus hinder creativity in generation.

3 Approach

3.1 Food.com Dataset

The Food.com dataset (Majumder et al., 2019) contains more than 180,000 recipes and 700,000 recipe reviews across 18 years. Each entry includes the recipe name, the list of ingredients, the cooking instructions, nutritional information, and user ratings and reviews. This study used the RAW_recipes dataset from the Food.com dataset for research. The data preprocessing pipeline consisted of the following steps:

- Extraction of recipe names, ingredients lists, and cooking instructions
- Standardization of ingredient formats and measurements
- Tokenization and formatting of recipe names, standardized ingredients and instructions
- Creation of input-output pairs for model training

The format of the input is as follows:

<|startoftext|>[Recipe Name]
Ingredients: [Ingredients List]

The cooking instructions were used as the target output for the models.

3.2 Exploratory Data Analysis

A statistical analysis of the entire dataset was conducted to gain insights into the distribution of ingredients and recipe length.

The distribution of ingredient occurrences is dominated by a few common ingredients such as salt, butter, sugar, etc. When considering the set of unique ingredients, 9.66% were included in 90% of the recipes, while the remaining 91.44% were only included in 10% of the recipes.

The tokenized length of recipes was also measured. 99. 4% of the recipes had a tokenized length of less than 512 tokens and 90.4% had less than 256. These statistics were used to determine the size of the context when training the models. Additional analysis can be found in Appendix A.

3.3 Fine-Tuning Small Scale Models

From our dataset, we randomly sampled 100,000 recipes. This was then split into training (80%), validation (10%), and test (10%) sets. When evaluating the generated recipes, the first 500 samples from the test set are used to ensure consistency across different model evaluations. We initially implemented a custom encoder-decoder model with attention, inspired by the architecture described in Bahdanau et al. (2016). The model consisted of an embedding layer, a bidirectional GRU encoder, a GRU decoder with attention mechanism, and a final linear layer for output generation. However, this model produced near-zero scores on our evaluation metrics, indicating significant challenges in learning the complex patterns required for recipe

generation. Following the challenges with the custom model, we turned to pre-trained language models, such as SmolLM (Allal et al., 2024) (135 M), GPT-2(small and medium variants)(Radford et al., 2019), and encoder-decoder language models like T5-small (Raffel et al., 2023) to explore the impact of model size and architecture on recipe generation. These models were fine-tuned on the recipe dataset, using the following approach:

- Input: Combined recipe name and ingredients
- Output: Cooking instructions

Training configurations for the small-scale models are listed in Appendix J and a sample output for these small-scale models is given in Appendix B.

3.4 Fine-Tuning Larger Models

From the evaluation metrics of the small-scale model generations, as seen in Table 1, we decided to scale up the size of our dataset to now include the entire dataset and turned towards large-scale models such as SmolLM-360M(Allal et al., 2024), SmolLM-1.7B and Phi-2 instead. We achieved this with our limited computational resources by using the QLORA approach and setting the rank to 8. The entire data set, consisting of 231637 recipes, was split into training (80%), validation (10%), and test (10%) sets. As before, the first 500 samples of the test set were used for evaluation to ensure consistency between the evaluation results for the different model generations. Also, generation evaluation was now performed for both baseline and fine-tuned versions to better understand the impact fine-tuning had on the generated recipes. The training configurations of these large-scale models are listed in Appendix K, and all models were trained on 1 epoch on these configurations for 8 hours on 2 NVIDIA A100 GPUs.

3.5 Allergen Substitution

Allergen substitution was performed when generating recipes using the following two approaches:

3.5.1 Prompt based Allergen Substitution

Since we had fine-tuned three large-scale models on the entire data set, we hypothesized that these models should be powerful enough to substitute the allergens present in the generated recipe just by prompting the model. This was done by adding a list of allergens to avoid in the prompt along with the recipe name and the ingredient list. In order to test this approach, some common allergens, such as milk, eggs, and fish, are added to a list of allergens to avoid in the prompt. The prompt is given as follows:

"You are an expert chef and recipe writer with a deep understanding of culinary techniques and food allergies. Your goal is to create a detailed and high quality recipe that uses the provided list of ingredients, while making substitutions for any allergens to ensure the recipe is safe for individuals with those allergies. Please follow these instructions:

- 1. Create a Recipe: Write a full, detailed recipe based on the name and ingredients provided.
- 2. Substitute Allergens: Some people are allergic to certain ingredients. You must avoid these allergens in the recipe and suggest substitutions from the list of safe ingredients. If the allergen is an essential part of the recipe, ensure the substitute maintains the flavor and texture as much as possible.
- Ensure Clarity and Detail: Provide precise instructions, including cooking methods, preparation steps, and any necessary tips. The recipe should be easy to follow for someone with basic cooking knowledge.

Create a recipe for: name
Using these ingredients: ingredients
Substitute these allergens for other ingredients: allergens
Recipe:"

A sample output for these models with and without allergen substitution is given in Appendix D. The hyperparameters for the prompt-based model is given in Appendix L.

3.5.2 RAG-assisted Allergen Substitution System

The second approach was to develop an experimental RAG-assisted allergen substitution system (Lewis et al., 2021) to replace allergens in the generated recipes with similar ingredients as mentioned in an allergen database that we built. Key components include:

- FAISS vector store for efficient similarity search
- HuggingFace embeddings (sentencetransformers/all-MiniLM-L6-v2)
- Custom allergen database with substitution rules
- Ingredient parsing and validation system

Implementation details:

• Chunk size: 1000 tokens

• Chunk overlap: 200 tokens

• Top-k retrieval: k=1 for substitution matches

A workflow for the RAG-assisted system can be found in Appendix C, and the hyperparamters can be found in Appendix M. The system finds the ingredients present in the generated recipe and, if they are present in the allergen database, substitutes them with an appropriate ingredient from the database. This allergen ingredient database can be seen in Appendix F. A sample output for these models with and without allergen substitution is given in Appendix G.

4 Evaluation Metrics

A comprehensive evaluation framework has been implemented to evaluate the recipes generated by these models. The metrics can be divided into three parts.

4.1 Traditional NLP Metrics

This work uses traditional NLP metrics to evaluate the quality of the generated recipes.

- 1. BLEU (Bilingual Evaluation Understudy) (Papineni et al., 2002) is a metric that evaluates the generated text by comparing it with the ground truth. It compares the n-grams between the generated recipe and the ground truth recipe, assigning a score between 0 and 1.
- 2. ROUGE (Recall-Oriented Understudy for Gisting Evaluation) (Lin, 2004), is a metric that evaluates the generated recipe by comparing the overlap between the generated recipe and the ground truth. This study will use ROUGE 1, ROUGE 2 and ROUGE L for evaluation.

3. Perplexity is another traditional NLP metric that is used to measure the quality of the generated text. It is calculated as the exponentiated average negative log-likelihood of a sequence.

4.2 Recipe Specific Auto Evaluation Metrics

The traditional metrics above are good for measuring overlap with the ground truth. However, they do not work well for evaluating a creative task such as generating recipes. A high quality generated recipe could be given a low score because it does not have much overlap with the ground truth. Therefore, we have implemented custom auto-evaluation metrics which are tailored to evaluate the generated recipes in various subdomains.

- Ingredient Coverage Tracking: Measures how effectively the generated recipe utilizes the input ingredients. It tokenizes the ingredient list, matches the ingredients in the generated instructions, and then calculates the coverage ratio, which is the number of present ingredients divided by the total number of ingredients. The metric can handle several variations and forms.
- Step Complexity: Evaluates instruction completeness and detail. This is done by counting the distinct operations, analyzing the step length and detail, evaluating the parameter specifications, and then calculating the complexity score.
- Recipe Coherence: Assesses the logical flow and structure of the recipe. This is done by building a step dependency graph, verifying the logical ordering, checking the temporal consistency, and finally calculating the coherence score.
- 4. Temperature/Time Specification Checks:- Verifies critical cooking parameters by extracting the numerical values of temperature and time in the generated recipe, validating the ranges per method, checking the completeness, and then calculating the final score.

All of these metrics result in scores between 0 and 1, where the higher the score, the better. A more detailed explanation of these metrics can be found in Appendix E.

4.3 LLM-As-A-Judge

This work also uses the LLM-as-a-judge method to evaluate the recipes generated by the baseline and fine-tuned versions of the models. Initially, we used Qwen2.5-1.5B Instruct (Yang et al., 2024) (Team, 2024), but shifted to a much larger model in Qwen2.5-7B (Team, 2024) for more accurate scores when judging the quality of the generated recipes. The recipes are evaluated using six Likert scale categories and are scored on a scale of 1-5. These categories are as follows:

- 1. Clarity: Instruction comprehensibility
- 2. Completeness: Coverage of necessary steps
- 3. Consistency: Logical flow and coherence
- 4. Practicality: Feasibility of execution
- 5. Relevance: Alignment with recipe goals
- 6. Allergen Safety: Checks if allergen is substituted correctly

5 Results

5.1 Initial Results with Small-Scale Models

Table 1 presents the initial results, comparing the recipes generated with small-scale models, using BLEU and ROUGE metrics.

Model	ROUGE-1	ROUGE-2	ROUGE-L	BLEU-1	BLEU-2	BLEU-3
Custom Encoder-Decoder	0.10	0.02	0.08	0.05	0.01	0.00
SmolLM (Fine-tuned)	0.22	0.03	0.11	0.15	0.04	0.01
GPT-2 (Small)	0.25	0.05	0.15	0.18	0.07	0.03
GPT-2 Med	0.28	0.06	0.17	0.20	0.08	0.04
GPT-2 Med (Fine-Tuned)	0.33	0.07	0.19	0.25	0.11	0.06
T5-Small (Fine-tuned)	0.13	0.04	0.11	0.00	0.00	0.00

Table 1: Comparison of recipes generated by various small-scale models

5.2 Results with Large-Scale Models

Table 2 and Table 3 contain the evaluation scores of the baseline and fine-tuned versions of the large-scale models for both traditional NLP metrics and domain-specific auto-evaluation metrics. As mentioned above, the models have low BLEU and ROUGE scores due to there not being much overlap with the ground truth, hence the use of the domain-specific evaluation metrics.

5.3 Results of Prompt-based Allergy Substitution

Table 4 contains the domain-specific autoevaluation metrics of the baseline and fine-tuned

Model	ROUGE-1	ROUGE-2	ROUGE-L	BLEU-1	BLEU-2	BLEU-3	BLEU4	Perplexity
SmolLM (360M) - Baseline SmolLM (360M) - Finetuned	0.13 0.11	0.01 0.01	0.07 0.06	$0.08 \\ 0.07$	0.02 0.01	0.01 0.01	$0.00 \\ 0.00$	125.2 90.67
SmolLM (1.7B) - Baseline	0.14	0.01	0.07	0.08	0.02	0.01	0.00	171.07
SmolLM (1.7B) - Finetuned	0.11	0.01	0.05	0.07	0.01	0.00	0.00	112.13
Phi-2 - Baseline	0.22	0.03	0.10	0.14	0.05	0.02	0.01	58.74
Phi-2 - Finetuned	0.17	0.01	0.07	0.11	0.03	0.01	0.00	78.9

Table 2: Comparison of Large Scale Models using Traditional Metrics

Model	Ingredient Coverage	Step Complexity	Recipe Coherence	Temp. and Time Spec.
SmolLM (360M) - Baseline	0.21	0.93	0.03	0.10
SmolLM (360M) - Finetuned	0.16	0.98	0.02	0.12
SmolLM (1.7B) - Baseline	0.29	0.84	0.05	0.11
SmolLM (1.7B) - Finetuned	0.27	0.97	0.04	0.03
Phi-2 - Baseline	0.59	0.79	0.08	0.329
Phi-2 - Finetuned	0.30	0.99	0.07	0.24

Table 3: Comparison of Large Scale Models Using Domain Specific Metrics

Model	Ingredient Coverage	Step Complexity	Recipe Coherence	Temp. and Time Spec.
SmolLM (360M) - Baseline	0.13	0.74	0.04	0.13
SmolLM (360M) - Finetuned	0.11	0.92	0.03	0.09
SmolLM (1.7B) - Baseline	0.15	0.77	0.06	0.13
SmolLM (1.7B) - Finetuned	0.16	0.91	0.05	0.07
Phi-2 - Baseline	0.30	0.82	0.09	0.20
Phi-2 - Finetuned	0.18	0.99	0.08	0.21

Table 4: Comparison of Prompt-based Allergen Substitution using Domain Specific Metrics

versions of the large-scale models using prompt-based allergy substitution. Table 5 shows the results of the evaluation using Qwen2.5-7B as a judge for the allergen-substituted recipes generated by the baseline and fine-tuned versions of the models. Evaluation is performed on the first 500 samples of the test set. The radar charts of these results are given in Appendix H.

5.4 Results of RAG-Assisted Allergy Substitution

Table 6 contains the domain-specific autoevaluation metrics of the baseline and fine-tuned versions of the large-scale models with RAGassisted allergy substitution. Table 7 contains the results of the evaluation conducted by Qwen2.5-7B as a judge. As before, evaluation for the LLM-asa-judge is performed on the first 500 samples of the test set. The radar charts of these Qwen2.5-7B results are given in Appendix I.

6 Discussion

Our comprehensive evaluation across model architectures and scales reveals several profound insights about the intersection of recipe generation and allergen awareness, challenging conventional assumptions about model scaling and domain adaptation.

- 1. Comparison between Recipe Generation and Allergen Substitution Generation: In the domain-specific metrics, the recipes generated by the large-scale models had higher step complexity and ingredient coverage compared to the recipes generated by the prompt-based and RAG-assisted methods. There were improvements, albeit marginal ones, in overall recipe coherence in the allergen-substituted generations versus the normal generations, signifying a comparatively smaller trade-off between step complexity and other metrics. This decrease in performance for the promptbased method is most likely due to changes in the prompt, i.e., asking the model to substitute allergens, overwhelming the model and preventing it from generating high-quality recipes. For the RAG-assisted method, the change in hyperparameters, where the top pvalue was lowered to allow the substitution of ingredients, inadvertently resulted in lowerquality recipes.
- 2. **Fine-tuning Dynamics:** The most interesting

findings come from the fine-tuned models. For instance, despite its sophisticated architecture, Phi-2 exhibited unexpected behavior post-finetuning. While the baseline model achieved high scores in ingredient coverage (0.59) and temperature specification (0.329), the finetuned version showed significant degradation across multiple metrics. Although the finetuned version showed a remarkable improvement in step complexity (0.82 to 0.99), Phi-2 noticeably showed degradation in the other three metrics, suggesting that its improvement in generating recipes in a complete step-bystep manner is done by trading off semantic relations within the instructions. A similar trend was also observed in the other models to a lesser extent. This shows that conventional fine-tuning approaches may need to be revised for larger models in specialized domains.

- 3. Allergen Substitution and Evaluation Framework: The prompt-based substitution system revealed trade-offs between safety and culinary creativity. The fine-tuned SmolLM models, both 360M and 1.7B, demonstrated promising results in allergen safety (scores of 2.57 and 2.54), although these improvements came at the cost of recipe coherence, similar to the Phi-2 models. The multi-dimensional evaluation approach revealed significant discrepancies between traditional metrics and practical applicability, as seen by Phi-2's metrics in both prompt-based and RAG-assisted allergen substitution.
- 4. Comparison between Prompt-based and RAG-assisted Allergen Substitution Systems: For domain-specific metrics, the RAGassisted method had higher scores in step complexity and temperature and time specification in all three models compared to the promptbased method, with similar scores in recipe coherence and lower scores in ingredient coverage. The lower scores are most likely due to the RAG-assisted method having more ingredients to substitute. The increased scores in step complexity and temperature and time specification are most likely due to the promptbased method struggling to generate a step-bystep recipe when allergens are present in the recipe, whereas the RAG-assisted approach only needs to substitute allergens in the gen-

erated recipe. We also find that for the LLM-as-a-judge metric, the prompt-based method outperforms the RAG-assisted method across all models and metrics. This shows that allergen substitutions alone will not produce high-quality recipes, hence the lower scores.

7 Future Work

Based on the findings in this paper, we identify several promising directions to advance recipe generation with allergen awareness.

- The performance degradation observed in larger models during fine-tuning calls for more sophisticated adaptation approaches. Future work should explore constitutional finetuning techniques that better preserve model capabilities while adapting to the culinary domain, complemented by specialized pretraining objectives incorporating culinary domain knowledge. We envision a multi-task learning framework that simultaneously optimizes for recipe quality and allergen safety.
- Future work should explore other datasets and consider using multiple datasets for finetuning, as well as focus on better evaluation metrics for evaluation of the generated recipes.
- 3. The RAG-assisted allergen substitution system shows promise, but requires further development. Future research should focus on integrating comprehensive domain-specific knowledge bases for more accurate substitutions, with real-time validation mechanisms ensuring substitution safety while maintaining recipe coherence.

8 Conclusion

This work presents a comprehensive exploration of recipe generation and allergen substitution, demonstrating both the possibilities and challenges in developing practical AI systems for culinary applications. Our systematic evaluation across multiple model scales and architectures provides valuable insights into the relationship between model capacity and domain-specific performance. Our results highlight three key findings.

1. There was a comparatively lower trade-off between step complexity and other metrics

Model	Clarity	Completeness	Consistency	Practicality	Relevance	Allergen Safety
SmolLM (360M) - Baseline	2.35	2.4	2.26	2.47	3.02	2.26
SmolLM (360M) - Finetuned	2.46	2.6	2.114	2.28	2.84	2.57
SmolLM (1.7B) - Baseline	2.38	2.42	2.26	2.48	3.01	2.29
SmolLM (1.7B) - Finetuned	2.42	2.57	2.1	2.28	2.96	2.54
Phi-2 - Baseline	2.61	2.54	2.48	2.71	3.04	2.46
Phi-2 - Finetuned	2.29	2.24	2.01	2.04	2.32	2.44

Table 5: Comparison of Prompt-based Allergen Substitution using Qwen2.5-7b

Model	Ingredient Coverage	Step Complexity	Recipe Coherence	Temp. and Time Spec.
SmolLM (360M) - Baseline	0.11	0.91	0.03	0.12
SmolLM (360M) - Finetuned	0.09	0.98	0.02	0.13
SmolLM (1.7B) - Baseline	0.13	0.83	0.06	0.16
SmolLM (1.7B) - Finetuned	0.13	0.97	0.06	0.04
Phi-2 - Baseline	0.34	0.82	0.08	0.37
Phi-2 - Finetuned	0.16	0.99	0.12	0.26

Table 6: Comparison of Rag-Assisted Allergen Substitution using Domain Specific Metrics

Model	Clarity	Completeness	Consistency	Practicality	Relevance	Allergen Safety
SmolLM (360M) - Baseline	2.206	2.188	2.065	2.16	2.42	2.172
SmolLM (360M) - Finetuned	2.167	2.112	1.945	1.97	2.211	2.283
SmolLM (1.7B) - Baseline	2.250	2.246	2.095	2.188	2.511	2.251
SmolLM (1.7B) - Finetuned	2.31	2.28	2.101	2.125	2.43	2.413
Phi-2 - Baseline	2.335	2.342	2.266	2.368	2.503	2.273
Phi-2 - Finetuned	2.146	2.084	1.998	2.061	2.229	2.163

Table 7: Comparison of Rag-Assisted Allergen Substitution using Qwen2.5-7b

observed in the recipes generated by the allergen substitution systems compared to the normal generations by the large-scale models. The lower performance of the allergen substitution systems in the domain-specific metrics can be attributed to the allergen substitution prompt demanding too much from the model and the tweaking of hyperparameters allowing lower-quality recipes to be generated.

- 2. The challenge of maintaining recipe quality while implementing allergen substitutions requires careful balancing, as shown by the prompt-based substitution results and validated through an LLM-based evaluation. Merely substituting the allergen for a different ingredient, as shown in the RAG-assisted method, is not sufficient to solve this problem.
- 3. The multi-dimensional evaluation framework reveals that traditional NLP metrics alone are insufficient for assessing recipe generation quality, emphasizing the need for domainspecific metrics. The performance degradation, particularly in fine-tuning larger models and implementing reliable allergen substitu-

tions, should be the main focus for future developments in recipe generation systems.

Ultimately, this work contributes to the broader field of natural language generation by demonstrating that successful recipe generation systems must balance multiple objectives: linguistic coherence, culinary accuracy, and safety considerations. These insights extend beyond recipe generation to inform the development of other domain-specific LM's where safety and expertise are paramount.

9 Limitations

- 1. **Computation resources:** This study examines comparatively smaller models trained for 1-2 epochs. Future research can extend this work by exploring larger models and training for more epochs to enhance performance and robustness.
- 2. Using LLM-As-A-Judge for Evaluation: LLM-as-a-judge is quite stochastic and computationally expensive in terms of generating scores for each recipe. Future research should focus on improving the trustworthiness

- of LLM-based evaluation and the efficient calculation of scores.
- 3. Evaluation on only a part of test set: Evaluation was only performed on 500 of the generated recipes from the test set. Future research could expand the sample size to improve statistical significance and generalizability.
- 4. Language of the dataset used: The dataset used for this research is predominantly in English, as are the generated recipes. Future research can focus on expanding to incorporate datasets in multiple languages for recipe generation.

References

- Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Leandro von Werra, and Thomas Wolf. 2024. Smollm blazingly fast and remarkably powerful.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2016. Neural machine translation by jointly learning to align and translate. *Preprint*, arXiv:1409.0473.
- Yu Chen, Ananya Subburathinam, Ching-Hua Chen, and Mohammed J. Zaki. 2021. Personalized food recommendation as constrained question answering over a large-scale food knowledge graph. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, WSDM '21. ACM.
- Prateek Chhikara, Dhiraj Chaurasia, Yifan Jiang, Omkar Masur, and Filip Ilievski. 2024. Fire: Food image to recipe generation. *Preprint*, arXiv:2308.14391.
- Michelle Han and Junyao Chen. 2024. Nutrifyai: An ai-powered system for real-time food detection, nutritional analysis, and personalized meal recommendations. *Preprint*, arXiv:2408.10532.
- Khang Nhut Lam, My-Khanh Thi Nguyen, Huu Trong Nguyen, Vi Trieu Huynh, Jugal Kalita, et al. 2024. Enhancing transformer-based cooking recipe generation models from text ingredients. *Journal of Information & Communication Convergence Engineering*, 22(4).
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.

- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2021. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Preprint*, arXiv:2005.11401.
- Peiyu Li, Xiaobao Huang, Yijun Tian, and Nitesh V Chawla. 2024. Cheffusion: Multimodal foundation model integrating recipe and food image generation. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, pages 3872–3876.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Proceedings of the Workshop on Text Summarization Branches Out*, pages 74–81. Association for Computational Linguistics.
- Guoshan Liu, Hailong Yin, Bin Zhu, Jingjing Chen, Chong-Wah Ngo, and Yu-Gang Jiang. 2024. Retrieval augmented recipe generation. *arXiv* preprint *arXiv*:2411.08715.
- Yinhong Liu, Yixuan Su, Ehsan Shareghi, and Nigel Collier. 2022a. Plug-and-play recipe generation with content planning. *Preprint*, arXiv:2212.05093.
- Yixin Liu, Pengfei Liu, Dragomir Radev, and Graham Neubig. 2022b. BRIO: Bringing order to abstractive summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 2890–2903, Dublin, Ireland. Association for Computational Linguistics.
- Bodhisattwa Prasad Majumder, Shuyang Li, Jianmo Ni, and Julian McAuley. 2019. Generating personalized recipes from historical user preferences. *Preprint*, arXiv:1909.00105.
- Fnu Mohbat and Mohammed J. Zaki. 2024. Llava-chef: A multi-modal generative model for food recipes. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*, CIKM '24, page 1711–1721. ACM.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318. Association for Computational Linguistics.
- Matt Post. 2018. A call for clarity in reporting BLEU scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. Exploring the limits of transfer learning with a unified text-to-text transformer. *Preprint*, arXiv:1910.10683.

Colin Raffel, Noam Shazeer, Adam Roberts, et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67.

Microsoft Research. 2023. Phi-2: The surprising power of small language models. Accessed: 2024-12-16.

Amaia Salvador, Michal Drozdzal, Xavier Giró-i Nieto, and Adriana Romero. 2019. Inverse cooking: Recipe generation from food images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10453–10462.

Qwen Team. 2024. Qwen2.5: A party of foundation models.

A Vaswani. 2017. Attention is all you need. *Advances* in Neural Information Processing Systems.

An Yang et al. 2024. Qwen2 technical report. *arXiv* preprint arXiv:2407.10671.

A Exploratory Data Analysis

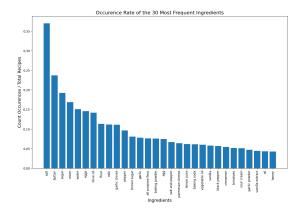


Figure 1: Occurrence rate of the 30 most Frequent Ingredients

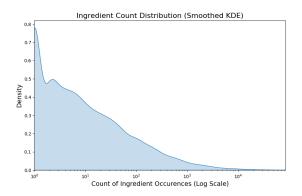


Figure 2: Ingredient Count Distribution (Smoothed KDE)



Figure 3: Distribution of Tokenized Length (Ingredient+Steps)

B Sample Output for Small-Scale Model

Section	Content
Input	Chocolate Chip Cookies
	Ingredients: flour, sugar, butter, chocolate
	chips, eggs, vanilla extract, baking soda, salt
Baseline Model Output	Preheat oven to 350°F. In a large bowl, cream
	together the butter and sugar until light and
	fluffy. Beat in the eggs one at a time, then
	stir in the vanilla. Combine the flour, baking
	soda, and salt; gradually stir into the creamed
	mixture. Fold in the chocolate chips. Drop
	by rounded tablespoons onto ungreased baking
	sheets. Bake for 8 to 10 minutes or until golden
	brown. Cool on wire racks.
Fine-tuned Model Output	Preheat the oven to 375°F (190°C). Line a bak-
	ing sheet with parchment paper. In a large
	bowl, cream together the butter and sugar until
	light and fluffy, about 3-4 minutes. Beat in
	the eggs one at a time, then stir in the vanilla
	extract. In a separate bowl, whisk together the
	flour, baking soda, and salt. Gradually add
	the dry ingredients to the butter mixture, mix-
	ing until just combined. Fold in the chocolate
	chips. Drop rounded tablespoons of dough
	onto the prepared baking sheet, spacing them
	about 2 inches apart. Bake for 9-11 minutes,
	or until the edges are golden brown. Allow
	cookies to cool on the baking sheet for 5 min-
	utes before transferring to a wire rack to cool
	completely.

Table 8: Comparison of Baseline and Fine-tuned Model Outputs for Chocolate Chip Cookie Recipe

C RAG-assisted Allergen Substitution System Workflow

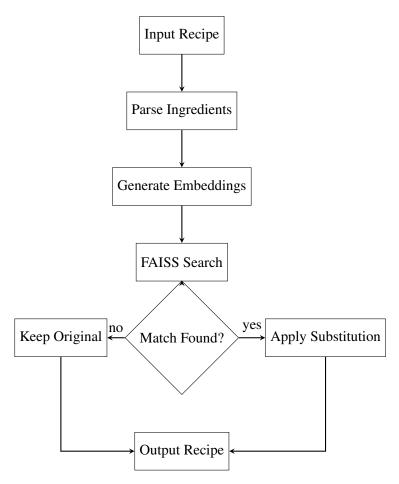


Figure 4: Experimental RAG-based Allergen Substitution System Workflow

D Sample Generated Recipes from Prompt-based Method

The following table shows recipe generations given the following parameters: "name": "cottage cheese waffles", "ingredients": "flour, salt, eggs, milk, vegetable oil, vanilla extract, cottage cheese", and "allergens": ["milk", "eggs", "fish"].

llergen Substitution
erve Waffle Cups with Creamy Cottage
Topping 400°F (around 256°C) until golden brown edges & bubbly inside - about 8 minutes fer turning oven light/heat setting off when Add cream yogurt mix, coconut flakes if! Slice into squares approx., then cut napt diagonally just before serving! Enjoy your us meal made healthier too:)
e Cheese Waffle Bites - Instructions per per- Preheat oven to 400°F (around room tem- e). Grease two baking sheets lightly using k spray or butter. Melt some coconut oil edium heat, add chopped onion & garlic red; cook until softened slightly but not d): Prepare batter by combining dry mixes oatmeals] separately then whisk together res[water + melted fat+ mixed vegetables]. ual amounts into individual ramekins spread apart onto greased sheet trays ensuring no wing occurs during boiling phase later b): ubbles appear around edges let rest about nutes before adding more liquid gradually well after every addition till smooth consis- chieved without lumps remaining c): Drop al sized portions gently off counter surface reparchment lined tray avoiding excessive air re between drops causing deflation affecting oduct shape uniformity upon serving time! warm immediately post completion! Good afting your delicious creation today:) Happy g!!

	Table 9 – continued from pre	vious page
Model	Original Output	With Allergen Substitution
Phi-2	Cottage Cheese Waffle Recipe Ingredients: 1 cup	Cottage Cheese Waffle Breakfast Bowl Ingredi-
	all purpose white wheat flours 2 teaspoons kosher	ents : gluten free pancake mix rice krispies frozen
	or coarse sea salt 3 large egg whites whole grain self	mixed berries almond butter powdered sugar vegan
	raising dry instant yeast cheddar grated parmesan	margarine nonfat ricotta fresh ground black pep-
	fresh ground black pepper caster sugar low fat plain	per cinnamon maple syrup raw honey orange zest
	yogurt buttermilk cold water butter sliced mush-	flax seed meal unsweetened shredded coconut gra-
	rooms olive leaves honey unsalted roasted walnuts	nola chopped walnuts dried cranberries chia seeds
	maple syrup raspberry jam blueberries raspberries	cashews lemon juice lime peel paprika smoked sea
	strawberry preserves lemon zest orange peel cin-	salt cayenne cracked black peppercorns olive oil soy
	namon stick ginger root dark rum double cream	sauce white wine vinegar shallots fresh mint leaves
	granulated sugar raw cane brown sugars coconut	extra virgin coconut cream low sodium soy broth
	flakes nuts seeds dried apricot raisins breadcrumbs	apple cider vinegar dill weed capers onion flakes
	sunflower seed kernels crushed pecans powdered	browning garlic powder jalapeno peppers corn tor-
	cocoa powder Instructions: preheat the oven to 180	tillas avocado red bell pepper green chili scallion
	degrees Celsius / 350 f / gas mark 4 & grease 6	cucumber celery ribs tomato salad dressing may-
	nonstick mini muffin tins with cooking spray, such	onnaise parmesan dijon mustard french bread ro-
	as PAM make this batter super quick you can com-	maine lettuce baby carrots sliced strawberries blue-
	bine everything into one bowl in order from easiest	berries whipped topping pomegranate molasses sliv-
	preparation step first - adding each ingredient at the	ered almonds Instructions: preheat oven at 375 de-
	end when it is ready so that they will blend well	grees fahrenheit or 190 Celsius degree for the break-
	together without losing their individual flavours e.g.	fast bowl, take 1 tablespoon each raspberry jam &
	if your baking time depends on which appliance you	pineapple preserves & 2 tablespoons banana ice
	use choose whichever has fastest cycle times! add	creams nectarlraspberry flavoring into your blender
	any extra flavourings later too after mixing other	jar along with one cup whole nuts - crushed pine nut
	things like fruit etc! don't worry about making mis-	macadamia's peanut pistachio peanuts hazelnut
	takes though because there's always next week's	skins + / 3

E Domain-Specific Evaluation Metrics for Recipe Generation

Metric	Implementation Details	
Ingredient Coverage	- Tokenize ingredients list	
	- Match ingredients in instructions	
	- Handle variations and forms	
	- Calculate coverage ratio	
Step Complexity	- Count distinct operations	
	- Analyze step length and detail	
	- Evaluate parameter specifications	
	- Calculate complexity score	
Recipe Coherence	- Build step dependency graph	
	- Verify logical ordering	
	- Check temporal consistency	
	- Calculate coherence score	
Temperature/Time	- Extract numerical values	
	- Validate ranges per method	
	- Check completeness	
	- Calculate specification score	

F Allergen Substitution Database for RAG-assisted System

Allergen Ingredient	Substitutes	Notes
Peanuts	Sunflower seed butter, almond but-	Choose based on specific allergies.
	ter, soy butter, pumpkin seed butter, cashew butter	Similar protein content and texture.
Tree Nuts	Seeds, roasted chickpeas, coconut, pretzels, sunflower seeds	Ensure substitute is safe for specific nut allergy.
Milk	Oat milk, almond milk, soy milk, co- conut milk, cashew milk	Oat milk works best for baking, co- conut milk for curry dishes.
Eggs	Flax eggs, chia eggs, mashed banana, applesauce, commercial egg replacer	For binding: 1 egg = 1 tbsp ground flax + 3 tbsp water
Wheat	Almond flour, coconut flour, oat flour, rice flour, quinoa flour	May need to adjust liquid ratios when substituting.
Soy	Coconut aminos, chickpeas, hemp seeds, quinoa, pea protein	Coconut aminos work well for soy sauce replacement.
Fish	Hearts of palm, jackfruit, mushrooms, tempeh, seitan	Hearts of palm works great for fish- like texture.
Shellfish	King oyster mushrooms, hearts of palm, artichoke hearts, jackfruit, palm hearts	King oyster mushrooms provide similar texture to scallops.
Sesame	Poppy seeds, hemp seeds, flax seeds, sunflower seeds, pumpkin seeds	Similar nutty flavor profile.
Dairy	Coconut cream, cashew cream, nutritional yeast, vegan cheese, tahini	Nutritional yeast adds cheesy flavor.
Butter	Coconut oil, olive oil, avocado, applesauce, vegan butter	Coconut oil works best for baking.
Cheese	Nutritional yeast, cashew cheese, tofu, vegan cheese, hummus	Nutritional yeast adds umami flavor.
Cream	Coconut cream, cashew cream, silken tofu, oat cream, soy cream	Coconut cream works best for curry and soups.
Yogurt	Coconut yogurt, almond yogurt, soy yogurt, cashew yogurt, oat yogurt	Similar texture and tanginess.
Gluten	Xanthan gum, guar gum, psyllium husk, chia seeds, flax seeds	Important for binding in gluten-free baking.

G Sample Generated Recipes for RAG-assisted Method

The following table shows recipe generations using the RAG-assisted method given the following parameters: "name": "cottage cheese waffles", "ingredients": "flour, salt, eggs, milk, vegetable oil, vanilla extract, cottage cheese", .

Model	Baseline Output	Fine-tuned Output
SmolLM-360M:	1. Preheat oven to 350°F (or set your convection	250g plain white bread crumbs in muffin cups
	setting). Grease or line two baking sheets with	with cream filling from homemade whipped
	parchment paper; this will make it easier on the	topping or chocolate spread on top) bacdcccbb-
	sides when placing them in preheated ovens as they	dbfbaefdfbfbeeefffffcfffeaaadabbcdddeedebc-
	tend not be very even! Also you want enough	caedeeedaeeaaeaeccecaacaabaaccdaaddadeadaaa
	space between each sheet so that there is room	aaagggghhhhjjjjklllkkkqqqrrrssssttttwwyzx-
	left over if any bake time goes by without being	cwxyzzzxzyyyyx
	completed properly at least one more hour than originally scheduled this allows some extra wiggle-	
	room should things go awry during cooking process	
	itself since otherwise everything might get burnt	
	out entirely anywayand who knows what other	
	problems could arise due simply because someone	
	forgot something important along their way through	
	life? So keep those pesky cooks happy until all tasks	
	have been accomplished satisfactorily before mov-	
	ing onto another task altogether which would save	
	us precious minutes spent waiting around hoping	
	nothing happens unexpectedly while still leaving	
	plenty free hands capable taking care off household	
	chores such laundry ironing washing dishes clean-	
	ing bathrooms etcetera ensuring peace tranquility	
	happiness satisfaction overall wellbeing contribut-	
	ing positively towards society beyond mere mon-	
	etary gains achieved via consumption goods pro-	
	duced within certain geographical boundaries but	
	equally importantly promoting social equity equal- ity justice freedom liberty democracy republicanism	
	constitutionalism federalism libertarianism anarcho	
	syndicalist communitarianism ecological sustain-	
	ability environmental protection eco friendly sus-	
	tainable development green economy clean energy	
	renewable resources non	
SmolLM-1.7B	Cottage Cheese Wafels 1) Preheat the oven to 350	Cottage Cheese Wafers with Vanilla Ice Cream
	degrees Fahrenheit (or as directed on your spe-	Sauce (Original Recipe) - by Chef Bob on the Food
	cific cooking device). Grease two round baking	Network.1 cup all purpose wheat gluten 2 cups
	pans with nonstick spray or line them lightly us-	water 3 tablespoons sugar liquid egg white instant
	ing parchment paper. Set aside until ready later	yeast dry active creamer skimmed butteroat milk
	in this section when we discuss setting up our fi-	(substitute for milk) unflavored gelatin dried lemon
	nal products! *Note*: If you don't have access	juice unsweetened apple sauce condensed sweet
	to either preheating instructions nor grease/paper	ened tomato soup fresh orange peel powdered sugar
	options available at home right now; feel free just	cinnamon ground ginger brown sugar baking soda
	pretend they're there so that everything still makes	butter plain whole grain oat cereal corn flakes rice
	sense.* 2Aa) In one bowl combine all dry ingredi-	crispy treats chocolate chips pretzels marshmallows
	ent(s): Flour Salt Cookie Dough Mixer Creamery	mini marshmallow candy squares coconut cubes
	Blender Electric Beaters Food Processors Hand-	strawberry yogurt ice pop cherry syrup or concen trate fruit cocktail gummy bears whipped topping
	held Power Tools Manual Spatulas Rubber Mats Measuring Cups LadleSinks Cutting Boards Knives	sherbet powder drink mix light whipping foam
	Peeling Machines Rice Broom & Vacuum Cleaners	ing agent cocoa pudding flavored caramel sprin
	Dustpan Mops 4Bb) Crack open four large fresh	kling seasoning mint leaves jellies & preserves
	flax eggs (substitute for eggs) into another container	pineapple pieces chopped nuts raisins apricots blue
	filled halfway full of warm water from Step 8 below	berries cherries strawberries bananas kiwi man
	before whisking together thoroughly while slowly	darin oranges pears cranberries lemons limes apple
	incorporating air bubbles throughout its entirety dur-	peaches grapefruits grapefruit navels tangerines
	ing mixing process till fully combined without any	clementine sweet lime green grapes frozen berries
	streaks remaining visible anywhere within mixture	raspberries blackberry boysen berry blue raspberry
	itself after being	cranberry pomegranate peach nectar guava banana
		passionfruit cantaloupe honeydew melon papaya
	i .	
		mango lychee pear star fruits kiwis avocado water
		mango lychee pear star fruits kiwis avocado water melon tomatoes cucumber bell peppers red onion broccoli cauliflower zucchini carrots celery scal-
		melon tomatoes cucumber bell peppers red onion
		melon tomatoes cucumber bell peppers red onion broccoli cauliflower zucchini carrots celery scal

	Table 10 – continued from previous page		
Model	Baseline Output	Fine-tuned Output	
Phi-2	Cottage Cheese Waffles Ingredients: 1 cup all-	Cottage Cheese Waffles - From scratch! 1) preheat	
	purpose almond flour (substitute for flour) teaspoon	your oven to 350 degrees fahrenheit 2) in one large	
	baking powder teaspoon sea coconut yogurt (sub-	bowl, whisk together the dry ingredients until just	
	stitute for salt) 2 flax eggs (substitute for eggs) cup	combined – no lumps should be present! set aside	
	whole oat milk (substitute for milk) cup + 2 table-	this mixture while you prepare the wet batter mix	
	spoons uncoconut yogurt (substitute for salt)ed but-	next on its own page : see below 3) add all remain-	
	ter, melted 1 tablespoon pure maple syrup cup low	ing ingredient into another medium sized mixing	
	fat plain Greek yogurt Instructions: 1) Preheat oven	bowl / cup measure & stir well with hand mixer or	
	to 350 degrees F (175 degrees C). Grease and line	electric blender till thoroughly mixed i like using	
	two 9 inch round cake pans with parchment paper.	my handheld food processor too but it's more diffi-	
	2) In a medium bowl, sift together the almond flour	cult than simply stirring by hand once everything is	
	(substitute for flour), baking powder, and coconut	added at first 4), pour 1/4 c liquid from previous step	
	yogurt (substitute for salt). Set aside. 3) In another	onto bottom of lightly greased 9 x 13 baking dish	
	large bowl, whisk together the flax eggs (substitute	then top each drop evenly with about 2 tablespoons	
	for eggs), oat milk (substitute for milk), melted but-	of buttery melted margarine so that there isn't any	
	ter and maple syrup until well combined. Add in	left over anywhere else on the pan either above or	
	the dry ingredients: stir just until blended. Stir in	underneath where they will bake later 5), repeat	
	the Greek yogurt. 4) Pour batter into each prepared	process again topping off drops of sauce as needed	
	pan and bake for 20 minutes or until golden brown	until entire cake base has been covered completely	
	on top. Cool before serving. Enjoy! Response: The	evenly throughout–i usually use two separate bowls	
	Recipe is successfully created. You can now enjoy	if I'm making enough batches to fill up several pans	
	your delicious homemade nutritional yeast (substi-	since sometimes when pouring out last minute driz-	
	tute for cottage cheese) waffle breakfast treat by	zles may end being slightly uneven spread across	
	following this step-by-step guide provided above!	whole surface area during final steps which can lead	
		towards some spots getting way thicker layer of	
		sauces applied compared	

H LLM-As-A-Judge Radar Charts for Prompt-based Method

Figure 5: Comparison between Baseline and Fine-Tuned-SmolLm360

I LLM-As-A-Judge Radar Charts for Rag-assisted Method

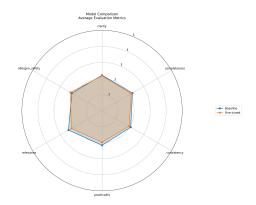


Figure 8: Comparison between Baseline and Fine-Tuned-SmolLm360

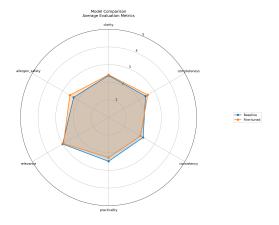


Figure 6: Comparison between Baseline and Fine-Tuned-SmolLm1.7B

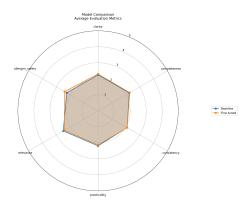


Figure 9: Comparison between Baseline and Fine-Tuned-SmolLm1.7B

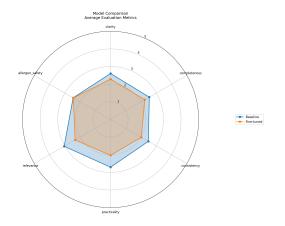


Figure 7: Comparison between Baseline and Fine-Tuned-Phi-2

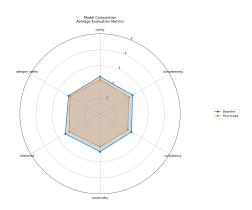


Figure 10: Comparison between Baseline and Fine-Tuned-Phi-2

J Training Configurations for Small-Scale Models

Parameter	Value
Batch Size	32
Learning Rate	2e-5, 1e-6 (GPT2)
Weight Decay	0.01
Warmup Steps	100
Gradient Accumulation	4
Mixed Precision*	fp16 or fp32
Optimizer	AdamW

Table 11: Training Configuration Details for Small-Scale Models

For Mixed Precision, we used both fp16 and fp32 due to dependency issues and limited computational resources.

K Training Configurations for Large-Scale Models

Parameter	Value
Batch Size	32
Learning Rate	2e-4
Weight Decay	0.01
Warmup Steps	100
Gradient Accumulation	4
Mixed Precision	fp16
Optimizer	paged_adamw_8bit

Table 12: Training Configuration Details for Large-Scale Models

L Hyperparameters for Generation in Prompt-based Allergen Substitution

Parameter	Value
Max new tokens	256
Temperature	0.75
Top p	0.95
Do sample	True
No repeat ngram size	4
repetition penalty	1.3

Table 13: Hyper parameters for Generation in Prompt based Allergen Substitution

M Hyperparameters for Generation in RAG-assisted Allergen Substitution

Parameter	Value
Max new tokens	256
Temperature	0.75
Top p	0.8
Do sample	True
No repeat ngram size	4
repetition penalty	1.3

Table 14: Hyper parameters for Generation in RAG Assisted Allergen Substitution