

Healthy Cooking with Large Language Models, Supervised Fine-Tuning, and Retrieval Augmented Generation

Andrea Morales-Garzón¹, Oscar A. Rocha, Sara Benel Ramirez, Gabriel Tuco Casquino², Alberto Medina³

¹Dept. of Computer Science and Artificial Intelligence, University of Granada

²Universidad Católica de Santa María, Perú

³ETSII-UPM, Universidad Politécnica de Madrid

Correspondence: amoralesg@decsai.ugr.es

Abstract

In response to the growing demand of the global society to adopt healthy habits, this paper presents the design, development and validation of a Spanish culinary recipe dataset focused on healthy nutrition that includes representative dishes from the gastronomy of Spanish-speaking countries. We also evaluate the dataset using Gemma 2B, Supervised Fine-Tuning and Retrieval Augmented Generation methodologies, showing its use to solve concrete problems related to nutrition.

1 Introduction

Food is essential in human development, and maintaining proper eating habits prevents diet-related risk factors and diseases. According to the World Health Organization, adopting healthy eating habits helps prevent malnutrition and diseases such as cancer, diabetes and obesity. Previous works propose scrapping a corpus of recipes from specialized cooking websites and their following homogenization (Li et al., 2022; Majumder et al., 2019; Marin et al., 2021; Salvador et al., 2017; Yagcioglu et al., 2018), or extending existing resources, tackling data curation and dataset extension (Bieñ et al., 2020). Despite their relevance, there is a notable shortage of Spanish resources to address nutrition-related computational problems. While large language models (LLMs) are a commonly used tool for providing culinary knowledge, their limitations, such as being error-prone, monotonous in similar scenarios, and potentially dangerous when dealing with diseases or allergies, underscore the need for a more accurate and reliable tool (Niszczota and Rybicka, 2023). We aim to alleviate these drawbacks: the need for a large recipe corpus in Spanish and the inefficient use of LLMs to improve the model’s responsiveness. We present a corpus of 20,447 recipes from Spanish-speaking countries to provide open resources to address societal ques-

tions and make them accessible to the Spanish-speaking community. This dataset, called *RecetasDeLaAbuela*¹, includes recipes belonging to the gastronomy of different Spanish-speaking countries obtained from multiple recipe websites while taking into account possible geographical and linguistic biases. To our knowledge, this is the largest dataset of recipes of Spanish-speaking origin available for free use. In addition, we tested its quality through an application that allows for nutritional and culinary context and food queries. For this, we use Gemma 2B and the latest Supervised Fine-Tuning (SFT) 4bit + Retrieval Augmented Generation (RAG) methodologies (Lewis et al., 2020).

2 Methodology

Our methodology is structured as follows: (1) Design for corpus creation, (2) Collection of recipes by web scraping, (3) Data curation and unification, (4) Corpus generation with instructions, (5) SFT training on lightweight 4-bit LLMs, (6) RAG with FAISS/LangChain and (7) Development and deployment of a demo in Gradio. Fig. 1 shows a diagram of the applied methodology.

3 Recipe corpus

Web scraping. We scraped several web pages on Hispanic/Latin American and international cuisine written in Spanish to create the dataset. We collected a total of 20,447 recipes using the Newspaper3k, Scrapy and BeautifulSoup libraries from the repository Frorozcoloa/ChatCocina², an open-source project for recipe extraction. In the extraction, attributes such as name, ingredients, preparation steps, duration, category, context or description, rating and votes, diners and difficulty of each recipe were obtained. In some recipes,

¹<https://huggingface.co/datasets/somosnlp/RecetasDeLaAbuela>

²<https://github.com/Frorozcoloa/ChatCocina>

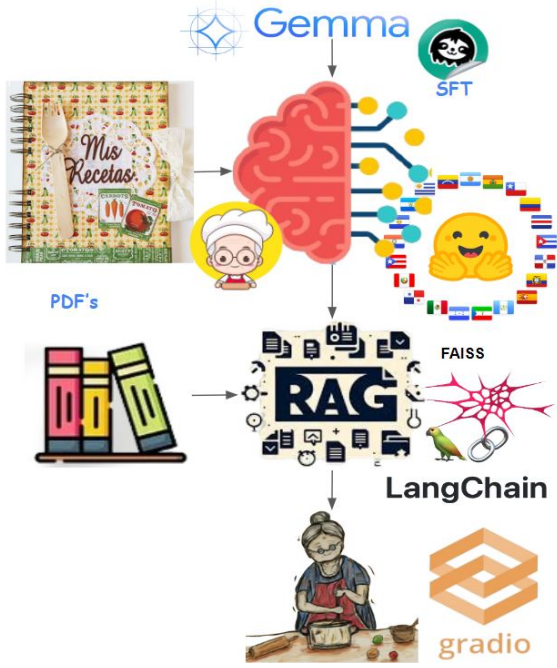


Figure 1: Methodology (LLMs SFT 4bit +RAG)

we calculated the country of origin by analyzing the keywords of the recipe (name of the country, gentilism and keywords of ancient cultures such as Aztec, Inca, gaucho, Guarani, Mapuche or Bolivarian). In adverse cases, we solved it by querying using the Together.AI API³ with *Mixtral-8x7B-Instruct-v0.1*. It is essential to highlight that some recipes lack values in specific attributes in the final corpus due to differences in structure or data in their origins.

Data preprocessing and homogenization. We performed data curation on specific fields of the dataset. For “*Ingredients*” and “*Steps*”, we have separated quantities from metric units (e.g., 1ml → 1 ml). For “*Name*” and “*Country*”, we have implemented a filter for duplicate recipes with the same name and same country. We have homogenized the corpus fields to facilitate the acquisition of statistics and visualization. We have used the HH:MM format for “*Duration*”. Finally, for “*Country*”, we have implemented the ISO format with the country codes.

Statistics and biases. We have obtained statistics through WordCloud visualizations, making use of TF-IDF to get three types of results:

(1) *Most used ingredients by country.* For example, Fig. 2 shows a Wordcloud with the most used ingredients in Mexican recipes.

(2) *Geographic biases.* The geographic representa-

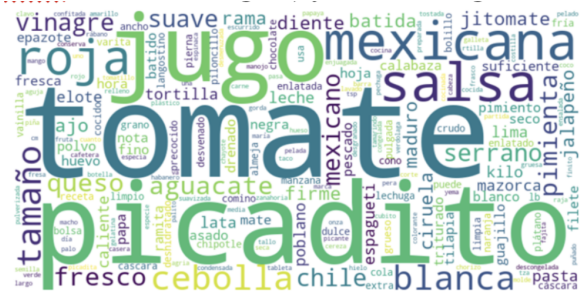


Figure 2: Most used ingredients in Mexico



Figure 3: Most used terms in recipe titles

tion of the dataset is closely linked to the nationality of those involved in the project. There is a higher proportion of recipes from Spain, and to a lesser extent from Mexico and Peru. Fig. 4 shows the heatmap of the number of recipes by country, where the higher the number, the darker the color. We have observed very few recipes from countries such as Panama or Paraguay since the original web pages do not focus on these countries.

(3) *Most present terms in recipe titles.* Some terms that refer to an origin may not necessarily be closely related to the recipe’s origin. Fig. 3 shows the WordCloud generated from the recipe names.

Creation of the instruction corpus. We have created the instruction corpus in two phases. First, we curated the initial corpus. We removed records with no data (nulls, blanks and line breaks) in the “*Ingredients*” and “*Steps*” columns and then concatenated the information from all columns into one. Secondly, we created a corpus of synthetic instructions using the LLM *Genstruct* in the *distilabel* environment. We used this LLM to generate question-answer pairs from the previously aggregated information. The main advantage of *Genstruct* is that it can massively simulate questions from a human user in a very varied and natural way instead of always asking for ingredients or preparation steps of different recipes, e.g., “Tell me a vegetarian dish” or “Recommend me a healthy chicken-based dish”.

³<http://together.ai>



Figure 4: Distribution of recipes per country

4 Fine-tuning of LLMs y RAG

SFT with Gemma. The SFT training has been performed using the unsloth/gemma-2b-bnb-4bit model since 4-bit quantization offers a superior effort/quality ratio (2.4x faster and 58% less VRAM vs. the Gemma 7B LLM). Mistral 7B (slower) and TinyLlama 1.1B (slightly faster but less accurate) have also been tested. A UTF-8 subset of 2.5k recipes of 3 attributes forming the question/answer pair (Recipe name vs ingredients and preparation steps) is extracted from the RecetasDeLaAbuel@ corpus. The SFT training lasts approximately two hours (equivalent to 8 epochs and 1550 steps) on HuggingFace Nvidia T4 medium (8 vCPU, 30 GV RAM, 16GB VRAM), and we obtained the RecetasDeLaAbuela5k model⁴.

SFT Hiperparameters. We trained the model using CUDA 12.1 Pytorch 2.2.2 with maximum input tokens 4096, LoRA $r=16$, $\alpha=16$, no dropout/bias/rslo/LoftQ, two processes, two batches, four gradient accumulation steps, five warmup steps, $2e-4$ learning rate, adamw_8bit optimizer, 0.01 weight drop rate and linear scheduler. A total of 19M parameters are trained for Gemma 2B. The maximum peak SFT/total memory reaches 12/14.4Gb (83%/98%).

Training results. We saved the SFT training every 10 steps in Weights & Biases. Fig. 5 shows a plot with the model loss, which is decreasing (1.8 at 200 steps and 1.3 at 1500 steps) with ± 0.05 shallow ripple. We used BERTScore (Zhang et al., 2019) to evaluate the resemblance of the original recipes to those generated by the model, obtaining precision: 0.67, recall: 0.71 and f1: 0.69.

Environmental impact. Experiments were performed using HuggingFace (AWS) in the sa-east-1 region, which has a carbon efficiency of 0.2 kg CO₂

⁴https://huggingface.co/somosnlp/RecetasDeLaAbuela5k_gemma-2b-bnb-4bit

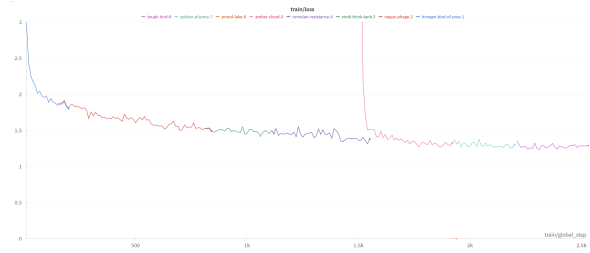


Figure 5: Model loss visualization

eq/kWh. A cumulative of 50 hours of computation was performed on HW type T4 (TDP of 70W). Total estimated emissions are 0.7 kg eq. CO₂, obtained through the ML CO₂ Impact website (Lacoste et al., 2019).

5 Results

The trained models can be tested in the Gradio demos RecetasDeLaAbuel@⁵ and ComeBien⁶. The user can enter their question about a recipe and add either a nutritional context (e.g., “You are an AI expert on cooking and nutrition”) or a culinary context extracted from cookbook PDFs using RAG and FAISS LangChain.

6 Conclusions and future work

The RecetasDeLaAbuel@ corpus has been successfully generated, trained and validated. With it, we have contributed to creating an open-source strategy to compile the most extensive corpus of recipes from Spanish-speaking countries. In future work, we will extend the experimentation with Mistral and TinyLlama and use MoE and RAG of traditional recipe books.

Acknowledgments

This project was developed during the international Spanish NLP Hackathon Somos600M organized by SomosNLP; we thank the organizers and sponsors, especially SomosNLP and HuggingFace, for GPU credits and model endpoints. We also thank to Tomás Vergara Browne for contributing with the English translation of this paper.

References

Michał Bień, Michał Gilski, Martyna Maciejewska, Wojciech Taisner, Dawid Wisniewski, and Agnieszka

⁵https://huggingface.co/spaces/somosnlp/RecetasDeLaAbuela_Demo

⁶https://huggingface.co/spaces/somosnlp/ComeBien_Demo

- Lawrynowicz. 2020. Recipenlg: A cooking recipes dataset for semi-structured text generation. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 22–28.
- Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. 2019. Quantifying the carbon emissions of machine learning. *arXiv preprint arXiv:1910.09700*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Ming Li, Lin Li, Qing Xie, Jingling Yuan, and Xiaohui Tao. 2022. Mealrec: a meal recommendation dataset. *arXiv preprint arXiv:2205.12133*.
- Bodhisattwa Prasad Majumder, Shuyang Li, Jianmo Ni, and Julian McAuley. 2019. Generating personalized recipes from historical user preferences. *arXiv preprint arXiv:1909.00105*.
- Javier Marin, Aritro Biswas, Ferda Ofli, Nicholas Hynes, Amaia Salvador, Yusuf Aytar, Ingmar Weber, and Antonio Torralba. 2021. Recipe1m+: A dataset for learning cross-modal embeddings for cooking recipes and food images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(1):187–203.
- Paweł Niszczoła and Iga Rybicka. 2023. The credibility of dietary advice formulated by chatgpt: robo-diets for people with food allergies. *Nutrition*, 112:112076.
- Amaia Salvador, Nicholas Hynes, Yusuf Aytar, Javier Marin, Ferda Ofli, Ingmar Weber, and Antonio Torralba. 2017. Learning cross-modal embeddings for cooking recipes and food images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3020–3028.
- Semih Yagcioglu, Aykut Erdem, Erkut Erdem, and Nazli Ikizler-Cinbis. 2018. Recipeqa: A challenge dataset for multimodal comprehension of cooking recipes. *arXiv preprint arXiv:1809.00812*.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Aprendiendo a cocinar de manera saludable con Large Language Models, Supervised Fine Tuning y Retrieval Augmented Generation

Andrea Morales-Garzón¹, Oscar A. Rocha, Sara Benel Ramirez, Gabriel Tuco Casquino², Alberto Medina³

¹Dept. of Computer Science and Artificial Intelligence, University of Granada

²Universidad Católica de Santa María, Perú

³ETSII-UPM, Universidad Politécnica de Madrid

Correspondence: amoralesg@decsai.ugr.es

Abstract

En respuesta a la creciente demanda de la sociedad global en adoptar hábitos saludables, este artículo presenta el diseño, desarrollo y validación de un dataset de recetas culinarias en español enfocado a nutrición saludable que incluye platos representativos de la gastronomía de países hispanohablantes. Además, evaluamos el dataset utilizando Gemma 2B, junto con las metodologías Supervised Fine-Tuning y Retrieval Augmented Generation, mostrando su uso para resolver problemas concretos relacionados con la nutrición.

1 Introducción

La alimentación constituye un pilar fundamental en el desarrollo humano, y mantener unos hábitos alimenticios adecuados es esencial para prevenir factores de riesgo asociados con la dieta y enfermedades directamente vinculadas con la alimentación. Según la Organización Mundial de la Salud, adoptar hábitos alimentarios saludables ayuda a prevenir la desnutrición y enfermedades como el cáncer, la diabetes y la obesidad, entre otras. Trabajos previos proponen escrapear un corpus de recetas de páginas web especializadas, junto con su posterior homogeneización (Li et al., 2022; Majumder et al., 2019; Marin et al., 2021; Salvador et al., 2017; Yagcioglu et al., 2018), o extender recursos existentes en términos de curado de datos y extensión del dataset (Bieñ et al., 2020). A pesar de su relevancia, hay una escasez notable de recursos en español que permitan abordar problemas computacionales relacionados con la nutrición. Los modelos grandes de lenguaje (LLMs) son una herramienta útil para proporcionar conocimiento culinario. Sin embargo, su uso es propenso a errores, se comportan de forma monótona en escenarios similares, y pueden ser peligrosos cuando tratamos con enfermedades o alergias (Niszczota and Rybicka, 2023). Nuestro objetivo es paliar estos

inconvenientes: la falta de corpus de recetas en español, y su aprovechamiento ineficiente con LLMs para mejorar la respuesta del modelo. Presentamos un corpus de 20,447 recetas de países hispanohablantes con el propósito de proporcionar recursos abiertos para abordar interrogantes de la sociedad y hacerlos accesibles a la comunidad hispanohablante. Este dataset, llamado *RecetasDeLaAbuela*¹, incluye recetas pertenecientes a la gastronomía de distintos países hispanohablantes obtenida de múltiples páginas web de recetas, a la vez que se tiene en cuenta posibles sesgos geográficos y lingüísticos. Hasta donde sabemos, este es el mayor dataset de recetas de origen hispanohablante disponible para su libre uso. Además, probamos su calidad a través de una aplicación que permite incluir contexto nutricional y culinario y realizar consultas sobre alimentación. Para ello, usamos Gemma 2B, junto con las últimas metodologías Supervised Fine-Tuning (SFT) 4bit + Retrieval Augmented Generation (RAG) (Lewis et al., 2020).

2 Metodología

La metodología se estructura siguiendo los siguientes pasos: (1) Diseño para creación del corpus, (2) Recopilación recetas mediante web scraping, (3) Curado y unificación de los datos, (4) Generación del corpus con instrucciones, (5) Entrenamiento SFT sobre LLMs ligeros de 4 bits, (6) RAG con FAISS/LangChain y (7) Desarrollo y despliegue de una demo en Gradio. La Fig. 1 muestra un diagrama de la metodología aplicada.

3 Corpus de recetas

Web scraping. Para la creación del dataset se realizó Web Scraping en diferentes páginas web sobre cocina hispano/latinoamericana e internacional redactada en español. En total se recopilaron

¹<https://huggingface.co/datasets/somosnlp/RecetasDeLaAbuela>



Figure 1: Metodología (LLMs SFT 4bit +RAG)

20,447 recetas mediante las librerías Newspaper3k, Scrapy y BeautifulSoup a partir del repositorio Frorozcoloa/ChatCocina² (proyecto open-source para extracción de recetas). En la extracción se obtuvieron atributos como nombre, ingredientes, pasos de preparación, duración, categoría, contexto o descripción, valoración y votos, comensales y dificultad de cada receta. En algunas recetas, el país de procedencia se ha calculado mediante análisis de palabras clave de la receta (nombre de países, gentilicios y palabras clave de culturas antiguas como azteca, inca, gaucho, guaraní, mapuche o bolivariano) y en caso negativo se ha resuelto preguntando mediante la API Together.AI³ al modelo Mixtral-8x7B-Instruct-v0.1. Es importante resaltar que algunas recetas carecen de valores en ciertos atributos en el corpus final debido a diferencias en su estructura o datos en sus orígenes.

Preprocesamiento y homogeneización. Se realizó un curado de datos en campos específicos del dataset. Para “Ingredientes” y “Pasos” se implementó una separación de cantidades con unidades métricas. Ejemplo: 1ml -> 1 ml. Para “Nombre” y “País” se implementó un filtro de recetas duplicadas que tuvieran mismo nombre y estuvieran asociadas a un mismo país. Hemos realizado una homogeneización de campos del corpus para facilitar las tareas de obtención de estadísticas y visualización. Para “Duración” se utilizó el formato HH:MM, y para “País” se implementó el formato ISO_A3 con los códigos de los países.

Estadísticas y sesgos. Las estadísticas se obtu-

²<https://github.com/Frorozcoloa/ChatCocina>

³<http://together.ai>

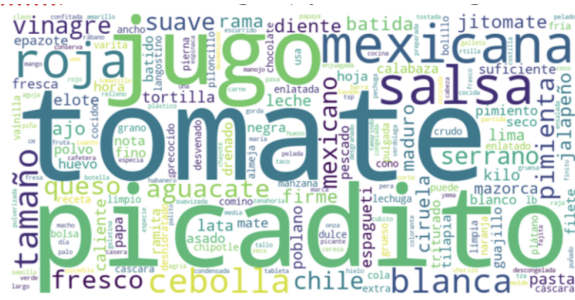


Figure 2: Ingredientes más usados en México



Figure 3: Términos más usados en las recetas

vieron mediante visualizaciones de WordCloud, haciendo uso de TF-IDF para obtener tres tipos de resultados:

(1) *Ingredientes más utilizados por país.* Por ejemplo, Fig. 2 muestra un Wordcloud con los ingredientes más utilizados en recetas mexicanas.

(2) *Sesgos geográficos.* Con respecto a la representación por país que tiene el dataset, que está estrechamente vinculado a la nacionalidad de los involucrados en el proyecto. Existe mayor proporción de recetas de España, y en menor medida de México y Perú. La Fig. 4 muestra el mapa de calor de la cantidad de recetas por país, donde a mayor cantidad, más oscuro es el color. Se observan muy pocas recetas de países como Panamá o Paraguay puesto que las páginas web originales no se centran en estos países.

(3) *Términos más presentes a la hora de definir los nombres de las recetas.* Algunos términos que hacen referencia a un origen, no necesariamente deben de estar estrechamente relacionados al origen de la receta. El WordCloud generado a partir de los nombres de las recetas se muestra en la Fig. 3.

Creación del Corpus de Instrucciones. La creación del corpus de instrucciones se ha realizado en 2 fases. Primero, depuramos corpus inicial. Se eliminaron los registros sin datos (nulos, espacios vacíos y saltos de línea) en las columnas de “Ingredientes” y “Pasos” y luego se concatenó la información de todas las columnas en una sola. A



Figure 4: Cantidad de recetas por país

continuación, creamos un corpus de instrucciones sintéticas mediante el LLM Genstruct en el entorno distilabel. Este LLM generó pares de preguntas-respuestas a partir de la información agregada previamente. La ventaja de Genstruct es que es capaz de simular de forma masiva preguntas de un usuario humano de forma muy variada y natural en vez de preguntar siempre por los ingredientes o pasos de preparación de distintas recetas, p. ej., “Dime un plato vegetariano” o “Recomiéndame un plato saludable basado en pollo”.

4 Fine-tuning de LLMs y RAG

SFT con Gemma. El entrenamiento SFT se ha realizado tomando como base el modelo gemma-2b-bnb-4bit de sloth puesto que la cuantización de 4 bits ofrece un ratio esfuerzo/calidad superior (2.4x más rápido y 58% menos VRAM frente al LLM Gemma 7b). Se ha probado también Mistral 7B (más lento) y TinyLlama 1.1B (un poco más rápido pero menos preciso). Se extrae del corpus RecetasDeLaAbuel@ un subconjunto UTF-8 de 2.5k recetas de 3 atributos que forman el par pregunta/respuesta (Nombre de la receta vs ingredientes y pasos de preparación). El entrenamiento SFT dura 2h aproximadamente (equivalente a 8 épocas y 1550 pasos) en HuggingFace Nvidia T4 medium (8 vCPU, 30 GV RAM 16GB VRAM) y se obtiene el modelo RecetasDeLaAbuela5k⁴.

SFT Hiper-parámetros. El entrenamiento se realiza bajo CUDA 12.1 Pytorch 2.2.2 con un máximo de tokens de entrada 4096, LoRA r=16, =16 sin dropout/bias/rsloa/LoftQ, 2 procesos, 2 batches, 4 pasos de acumulación de gradiente, 5 pasos de calentamiento, tasa de aprendizaje 2e-4, optimizador adamw_8bit, tasa de caída de pesos de 0.01 y planificador lineal. Para Gemma 2B se entrenan un

⁴https://huggingface.co/somosnlp/RecetasDeLaAbuela5k_gemma-2b-bnb-4bit

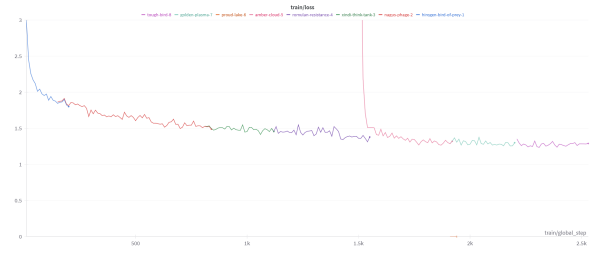


Figure 5: Visualización de la pérdida del modelo

total de 19M de parámetros. El máximo pico de memoria SFT/total alcanza 12/14.4Gb (83%/98%). **Resultados del entrenamiento.** El entrenamiento SFT se guarda cada 10 pasos en Wandb. La Fig. 5 muestra una gráfica con la pérdida del modelo, la cual es decreciente (1.8 a los 200 pasos y 1.3 a los 1500 pasos) con rizado ± 0.05 muy bajo. Utilizamos BERTScore (Zhang et al., 2019) para evaluar el parecido de las recetas originales con las generadas por el modelo, obteniendo precisión: 0.67, recall: 0.71 y f1: 0.69.

Impacto medioambiental. Los experimentos se realizaron utilizando HuggingFace (AWS) en la región sa-east-1, que tiene una eficiencia de carbono de 0.2 kg CO₂ eq/kWh. Se realizó un acumulado de 50 horas de cómputo en HW tipo T4 (TDP de 70W). Las emisiones totales estimadas son 0.7 kg eq. CO₂., obtenidas a través de la web ML CO₂ Impact (Lacoste et al., 2019).

5 Resultados

Los modelos entrenados se pueden probar en las demos Gradio RecetasDeLaAbuel@⁵ y ComeBien⁶. El usuario puede introducir su pregunta sobre una receta y añadir ya sea un contexto nutricional (“Eres una IA especialista en cocina y nutrición”) o un contexto culinario extraído de PDFs de libros de cocina mediante RAG y FAISS LangChain.

6 Conclusiones y trabajo futuro

El corpus RecetasDeLaAbuel@ ha sido generado, entrenado y validado con éxito. Se ha iniciado un trabajo open-source de compilación del corpus más extenso de recetas hispanoamericanas. Como trabajo futuro, ampliaremos la experimentación con pruebas en Mistral, TinyLlama, y utilizando MoE junto con RAG de libros tradicionales de cocina.

⁵https://huggingface.co/spaces/somosnlp/RecetasDeLaAbuela_Demo

⁶https://huggingface.co/spaces/somosnlp/ComeBien_Demo

Acknowledgments

Este proyecto se ha desarrollado en el Hackaton de NLP en Español Somos600M organizado por SomosNLP. Agradecemos a toda la organización y a sus patrocinadores, y en especial a SomosNLP y HuggingFace por facilitar los créditos de cómputo para GPU y recursos *endpoints* de los modelos para inferencia. También queremos agradecer a Tomás Vergara Browne por su colaboración en la traducción del artículo a inglés.

References

- Michał Bień, Michał Gilski, Martyna Maciejewska, Wojciech Taisner, Dawid Wisniewski, and Agnieszka Lawrynowicz. 2020. Recipenlg: A cooking recipes dataset for semi-structured text generation. In *Proceedings of the 13th International Conference on Natural Language Generation*, pages 22–28.
- Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. 2019. Quantifying the carbon emissions of machine learning. *arXiv preprint arXiv:1910.09700*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.
- Ming Li, Lin Li, Qing Xie, Jingling Yuan, and Xiaohui Tao. 2022. Mealrec: a meal recommendation dataset. *arXiv preprint arXiv:2205.12133*.
- Bodhisattwa Prasad Majumder, Shuyang Li, Jianmo Ni, and Julian McAuley. 2019. Generating personalized recipes from historical user preferences. *arXiv preprint arXiv:1909.00105*.
- Javier Marin, Aritro Biswas, Ferda Ofli, Nicholas Hynes, Amaia Salvador, Yusuf Aytar, Ingmar Weber, and Antonio Torralba. 2021. Recipe1m+: A dataset for learning cross-modal embeddings for cooking recipes and food images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(1):187–203.
- Paweł Niszczoła and Iga Rybicka. 2023. The credibility of dietary advice formulated by chatgpt: robot-diets for people with food allergies. *Nutrition*, 112:112076.
- Amaia Salvador, Nicholas Hynes, Yusuf Aytar, Javier Marin, Ferda Ofli, Ingmar Weber, and Antonio Torralba. 2017. Learning cross-modal embeddings for cooking recipes and food images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3020–3028.
- Semih Yagcioglu, Aykut Erdem, Erkut Erdem, and Nazli Ikizler-Cinbis. 2018. Recipeqa: A challenge dataset for multimodal comprehension of cooking recipes. *arXiv preprint arXiv:1809.00812*.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.