

Towards Multi-Language Recipe Personalisation and Recommendation

Niall Twomey*

Cookpad Ltd
Bristol, UK

niall-twomey@cookpad.com

Andrey Ponikar

Cookpad Ltd
Bristol, UK

andrey-ponikar@cookpad.com

Mikhail Fain†

Cookpad Ltd
Bristol, UK

mikhail-fain@cookpad.com

Nadine Sarraf

Cookpad Ltd
Bristol, UK

nadine-sarraf@cookpad.com

ABSTRACT

Multi-language recipe personalisation and recommendation is an under-explored field of information retrieval in academic and production systems. The existing gaps in our current understanding are numerous, even on fundamental questions such as whether consistent and high-quality recipe recommendation can be delivered across languages. Motivated by this need, we consider the multi-language recipe recommendation setting and present grounding results that will help to establish the potential and absolute value of future work in this area. Our work draws on several billion events from millions of recipes, with published recipes and users incorporating several languages, including Arabic, English, Indonesian, Russian, and Spanish. We represent recipes using a combination of normalised ingredients, standardised skills and image embeddings obtained without human intervention. In modelling, we take a classical approach based on optimising an embedded bi-linear user-item metric space towards the interactions that most strongly elicit cooking intent. For users without interaction histories, a bespoke content-based cold-start model that predicts context and recipe affinity is introduced. We show that our approach to personalisation is stable and scales well to new languages. A robust cross-validation campaign is employed and consistently rejects baseline models and representations, strongly favouring those we propose. Our results are presented in a language-oriented (as opposed to model-oriented) fashion to emphasise the language-based goals of this work. We believe that this is the first large-scale work that evaluates the value and potential of multi-language recipe recommendation and personalisation.

*Corresponding author.

† Authors listed alphabetically on last name.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

RecSys '20, September 22–26, 2020, Virtual Event, Brazil

© 2020 Copyright held by the owner/author(s).

ACM ISBN 978-1-4503-7583-2/20/09.

<https://doi.org/10.1145/3383313.3418478>

KEYWORDS

information retrieval, recommendation, personalisation, recipes and food modelling

ACM Reference Format:

Niall Twomey, Mikhail Fain, Andrey Ponikar, and Nadine Sarraf. 2020. Towards Multi-Language Recipe Personalisation and Recommendation. In *Fourteenth ACM Conference on Recommender Systems (RecSys '20)*, September 22–26, 2020, Virtual Event, Brazil. ACM, New York, NY, USA, 6 pages. <https://doi.org/10.1145/3383313.3418478>

1 INTRODUCTION AND RELATED WORK

A plurality of complex factors influence the choices we make when deciding on which recipes to cook. The occurrence of allergens in a recipe will render it inappropriate for some users, the personal commitments of other users will lead to total avoidance of some food categories, and for yet others, embarking on short- and long-term dietary campaigns will disturb otherwise steady and predictable eating habits. Indeed, since membership to any categorisation is non-exclusive and transient, we cannot rely on the existence of a uniformly acceptable recipe ranking across any subgroup. In this paper, we concern ourselves with delivering recipe personalisation models that capture individual preferences.

The formative works in recipe recommendation [11, 12] constructed user-based ingredient preferences from historic recipe interactions and represented recipes by their ingredients. Distances between users and recipes in the ingredient representation space can be used to make recommendations using k Nearest Neighbours (k -NN). This is a classic design pattern in the literature [2, 19, 41]. Because ingredients and recipes are written in un- or semi-structured forms, they are not necessarily amenable to immediate analysis, and recipe normalisation is known to be beneficial for downstream tasks [24, 28]. Thus, early work employed ontologies or knowledge graphs [6, 8, 9], supervised training [28] and scoring methods [32] to extract clean elements. Most of these approaches require a labelled set of ‘canonical’ recipe entities (ingredients, tools, skills) and to date have been evaluated on a single language, which presents a problem for scaling the approaches to a multi-language system.

Extracting the quality and quantity of recipes and ingredients [29, 32] is a key precursor in many application areas of food computing, including healthy recommendation [33]. The multi-modal aspect

of recipes has shown promise in enhancing cooking procedure understanding [40] by using auxiliary data such as video [22, 30] or images [27, 42]. The existing work on leveraging these modalities for recommendation [17, 34, 35] uses established pre-trained image models or specific image features (including measures of sharpness, contrast). Not all of the generic image representation approaches are suitable for recipe image datasets, which typically consist of images along with semi-structured text (title, ingredients, steps). Thus, transfer learning [5], cross-modal training [4, 13, 36], as well as self-supervised training with weak labels [10], have been proposed for building recipe image representation.

The application of recipe recommendation models to more than one language is not unexplored [20, 35], though the scope of these approaches is limited to two languages and is reliant on manual intervention for recipe pre-processing. Consequently, the limits of large-scale multi-lingual recipe modelling are under-explored, despite the existence of several relevant datasets and platforms operating in several languages [14, 21]. State-of-the-art (SOTA) neural translation [1, 39] and multi-language Natural Language Processing (NLP) frameworks [26] offer opportunities for bridging these gaps. To the best of our knowledge, ours is the first comprehensive work delivering intrinsic language-agnostic pipelines for recipe recommendation and personalisation across many languages. Towards this end, our methodological contributions and results are outlined in Sections 2 and 3, and we conclude on the value, limitations and future directions of this research in Section 4.

2 METHODS

2.1 Dataset

We consider only the published recipes with valid titles, ingredient lists, ingredient quantities, and method steps (as well as optional fields cooking duration, serving size, images). We limit ourselves to a single online multi-language recipe platform (Cookpad) as we are unaware of alternative data sources fulfilling our multi-language requirements.

We employ ten different user interaction types, and each is assigned a weight based on its likelihood of indicating future cooking behaviour estimated from proprietary data. Search data are also considered. Cookpad’s search seeks to serve the best *new* recipes to users based on their query. We have access to the queries, search result order and recipe clicks. Data fusion techniques are incorporated to merge the interaction data to the click data arising from the search histories in our analyses.

2.2 Data Representations

Qualitative Features. The purpose of these deterministic features is to provide key insights into recipe complexity, completeness, quality, and regionality. Specifically, we extract the following features: 1. the number of ingredients used; 2. the number of skills used; 3. recipe image; 4. the number of steps; 5. the number of step images; 6. the ratio of steps to step images; 7. the published year; 8. the published month; 9. the published time; 10. the author’s identity; 11. the author’s system ID 12. the cooking time; 13. the number of portions; 14. author’s country.

Normalised Ingredients. Recipe data contains ingredients with quantities in separate fields in various languages. Quantity extraction quality is not uniform across datasets with certain regional traits, while the ingredients themselves are written in raw form. For each language, we split the dataset by spaces, and, ignoring numerical values, extract a set of 200 common quantity tokens. We then remove any ingredients with punctuation or quantity present, and sort the remaining ingredients by frequency, picking the top 1500 common ingredients as the dictionary.

For inference, the raw ingredient string is matched to an ingredient from the dictionary by tokenising it and finding the largest common subset of tokens between a candidate normalised ingredient and the original ingredient. To deal with misspellings we use a threshold on the cosine distance between word vectors trained on a recipe corpus using FastText [3]. Our algorithm is unsupervised and can be applied to all languages with space separation between words.

Normalised Skills. Pre-trained SOTA language models [26] are deployed to detect verbs in recipe sentences. For a given sentence, the cartesian product between the sets of detected verbs and ingredients defines all possible ingredient-verb pairs for that sentence, e.g. ‘slice’ and ‘onion’. We incorporated partial matching and Levenshtein distance to overcome ingredient naming inconsistencies (e.g. when ‘peppers’ in a recipe step refers to ‘red bell peppers’ from the ingredient list). We broadcast ingredient-verb pair extraction over our datasets, extending existing work in skill extraction [42].

Let $C_I(i)$, $C_V(v)$, $C_{I,V}(i,v)$, $C_{I|V}(i|v)$ and $C_{V|I}(v|i)$ denote a family of counting functions for ingredients, verbs, joint ingredient-verb pairs, and conditional ingredient-verb pairs in recipe sentences. The domain of the random variables are $I \in \{i, \neg i\}$ and $V \in \{v, \neg v\}$ (i.e. ‘present’ and ‘absent’). Removing subscripts for improved clarity, counts are normalised to form distributions (e.g. $P(i,v) = C(i,v) / \sum_{i',v'} C(i',v')$), allowing us to calculate Mutual Information (MI) ($\sum_{i,v} P(i,v) \log \frac{P(i,v)}{P(i)P(v)}$). The summands are reformulated in terms of positive ingredient and verb occurrences since only these counts are available, i.e. $P(\neg a) = 1 - P(a)$ and $P(a, \neg b) = (1 - P(b|a))P(a)$. Denoting the MI matrix as $\mathbf{M} \in \mathcal{R}^{|I| \times |V|}$, the expected MI of ingredient i over associated verbs as $\mathbb{E}_{v' \sim P_{V|I=i}} [\mathbf{M}_{i,v'}]$, and the expected MI of verb v over ingredients as $\mathbb{E}_{i' \sim P_{I|V=v}} [\mathbf{M}_{i',v}]$, the data-dependent threshold for the (i,v) -th ingredient-verb pair is defined as

$$\Theta_{i,v} = \alpha \mathbb{E}_{v' \sim P_{V|I=i}} [\mathbf{M}_{i,v'}] + (1 - \alpha) \mathbb{E}_{i' \sim P_{I|V=v}} [\mathbf{M}_{i',v}]$$

where $\alpha \in [0, 1]$ is a hyperparameter (default value 0.5) that balances the relative weight of ingredients and verbs in skill selection. The final set of ingredient-skill pairs is given by $V = \{(i,v) : \mathbf{M}_{i,v} > \Theta_{i,v} \forall i,v\}$.

Text Representations. We explore representing recipe text as a bag of sub-word-units with Term Frequency-Inverse Document Frequency (TF-IDF) embeddings. The dimensionality of the embedding was reduced to 300 using singular value decomposition, and we followed the FastText [3] procedure in sub-word unit selection.

Image Representations. We used self-supervised training to extract image representations [10], extending the method to the multi-lingual setting. We trained a DenseNet-201 [15] model using a variation of TagSpace [38]. According to this approach, TagSpace labels in all languages and images are all embedded in the same shared 300-dimensional space, which leads to similar labels (in the same or different languages) ending up close to each other in the shared embedding space. The labels for training were the top 1000 common unigrams and bigrams extracted from recipe titles per each of the 5 languages. After the model was trained for 40 epochs on a dataset with 1.5M images and 5k labels, we discard the label embeddings and run the CNN on all recipe images. The image representations are extracted from the global average pooling layer after the last convolutional layer.

User Representations. In this work, user profiles derive directly from users’ interaction histories. Let $\mathbf{I} \in \mathbb{R}_+^{N_u \times N_r}$ be the (sparse) user-to-recipe interaction matrix that encodes the interaction importance numerically. Moreover, let $\boldsymbol{\mu}$ be the row-wise normalised interaction matrix, i.e. $\boldsymbol{\mu}_{u,r} = \mathbf{I}_{u,r} / \sum_{r'} \mathbf{I}_{u,r'} \quad \forall u, r$. Finally, let the recipe features be embedded in $\mathbf{X} \in \mathbb{R}^{N_r \times D}$ (c.f. Section 2.2). Using these definitions, the user features are calculated simply by averaging recipe embeddings over interaction history, i.e. $\mathbf{U} = \boldsymbol{\mu}\mathbf{X}$, and $\mathbf{U} \in \mathbb{R}^{N_u \times D}$.

2.3 Behavioural Models

‘Clickability’ Model. We used a 3-layer feed-forward neural network with ReLU non-linear activation units, dropout, and batch normalisation as a cold-start. The model is optimised in a pairwise approach [25] to differentiate between positive (clicked recipes) and negative (recipes not clicked) recipes. Negatives are sampled randomly from among the recipes viewed in the search results but not clicked. We incorporate late context fusion to measure the impact of using query context in ranking. We call this the ‘clickability’ model since it is based on click-through data.

Personalisation Model. LightFM [18] is a framework offering linear Collaborative Filtering (CF), Content Based (CB) and hybrid recommendation models. It is known to be strongly performant in scenarios with sparse and transient data, even for new users with little interaction history [18]. Our objective in this research is to establish strong definitive baselines for multi-language recipe recommendation [7]. Consequently, our future work will develop the assessment of SOTA recommendation models in this application area. In LightFM’s setting, given user and recipe embedding matrices ($\mathbf{E}^U \in \mathbb{R}^{D_U \times K}$ and $\mathbf{E}^R \in \mathbb{R}^{D_R \times K}$), user and recipe embeddings are calculated with $\mathbf{X}^U = \mathbf{U}\mathbf{E}^U$ and $\mathbf{X}^R = \mathbf{X}\mathbf{E}^R$. User-recipe affinity is measured by $S_{u,r} = f(\mathbf{X}_u^U \cdot \mathbf{X}_r^R + \mathbf{b}_u^U + \mathbf{b}_r^R)$ where \mathbf{b}^U and \mathbf{b}^R are user- and recipe-specific biases, and $f(\cdot)$ is a suitable function selected based on the task, e.g. logistic. In order to optimise the embedding matrices and biases, several hyperparameters must be specified (including learning rate, number of iterations, user and recipe regularisation, the loss function, sample weights, feature groups). Owing to the large number of hyperparameters, we take a sequential approach and use Bayesian Optimisation (BO) [23] on a validation set to select these parameters.

2.4 Evaluation Procedures

Table 1: The approximate number of users, recipes and interactions available for analysis.

	Arabic	English	Indonesian	Russian	Spanish	Total
Users	2M	3M	6M	2M	4M	18M
Items	0.6M	0.5M	2M	0.5M	0.5M	4M
Events	0.8B	0.2B	4B	0.5B	1M	7B

Table 1 presents dataset size that is available to us in this work. We stratify interaction data based on event time into four non-overlapping partially ordered sets named profile, train, validation and test, denoted $S_p < S_t < S_v < S_e$. Of particular importance is the ‘profile’ interactions (S_p) since these are exclusively used in creating user profiles (c.f. Section 2.2). Performance is evaluated with Mean Average Precision (mAP) [31] at $k \in \{1, 20\}$. Since in this work we consider a lot of moving parts for our system, we opted to rely on BO in assessment rather than on ablation studies.

3 RESULTS AND DISCUSSION

3.1 Ingredient and Skill Validation

Our ingredient normalisation algorithm picks up non-trivial patterns in the data. For example, the phrase ‘large sweet strawberries’ is normalised to ‘strawberry’, and ‘large sweet potatoes’ gets normalised to ‘sweet potato’. It is noteworthy that these two similar phrases are both correctly normalised in different ways by our unsupervised model. From a quantitative point of view, ingredient normalisation can match and outperform systems based on manually maintained ingredients dictionaries across the five languages considered. We compare performance of our ingredient normaliser to a system built on top of a professionally-maintained proprietary ingredient dictionary provided by Cookpad in Table 2. This shows statistically significant error rates on the dedicated test sets for our normaliser.

Table 2: Normalisation error rates evaluated by native speakers on each language. Statistical significance results in bold.

Language	Arabic	English	Indonesian	Russian	Spanish
Baseline	0.4	0.18	0.18	0.14	0.21
Proposed	0.26	0.10	0.12	0.05	0.09
Reduction	0.35	0.44	0.33	0.66	0.57

Table 3 shows discovered skills from all languages, omitting general skills such as ‘add’ and ‘mix’. Focusing specifically on English, the proposed pairings are of high quality and diversity (i.e. ‘deboning fish’ is not a skill that all cooks will employ). The skill quality on the remaining languages are similar in nature to English, and translations of ‘peel potato’ and ‘boil water’ can be found in the non-English columns of the table. Since skills are unstructured and may be ambiguous or unclear, we designed a small labelled experiment to evaluate definitive skill detection performance. Training data were

Table 3: An example of the discovered skill-ingredient pairs across five different languages.

Arabic		English		Indonesian		Russian		Spanish	
Ing.	Skill	Ing.	Skill	Ing.	Skill	Ing.	Skill	Ing.	Skill
بيض	خفق	onion	slice	air	didihkan	лук	нарезать	cebolla	cortar
دقيق	خلط	fish	debone	telur	aduk	морковь	натереть	agua	hervir
ثوم	قطع	eggs	beat	margarin	panaskan	разрыхлитель	просеять	cebolla	pelar
ماء	غلي	potatoes	peel	gula pasir	aduk	картофель	очистить	ajo	picar
دجاج	قطع	flour	sieve	tepung terigu	sajikan	сыр	посыпать	harina	amasar

acquired using a set of regular expressions and human validation on recipe title, ingredients and steps, and logistic regression models were optimised to detect the chosen skills using a bag of words encoding for each recipe field as features. The dataset contains ≈ 30 times more negatives than positives, yet our method has precision of approximately 0.5. This indicates that skills can be detected with reasonable precision with straightforward approaches.

3.2 Case Study 1: Interactions

Since the clickability model is aimed at users without interaction histories, it is trained with click-through data. Additionally, the text and image embeddings from Section 2.2 were used for recipe representations, and a grid search was used to select the model’s hyperparameters (learning rate, network architecture, dropout). For personalisation models, BO was used to select model hyperparameters but also to specify the model type (from CF, CB, or hybrid models) and recipe representation (from any combination of qualitative, ingredient or skill features). We ran BO for 100 iterations and selected the optimal model based on the performance on validation sets.

Results for both models are presented in Table 4. High consistency of selected models and representations are obtained, providing evidence that supports ingredients and skills in recommendation. Additionally, hybrid models are always selected over CF and CB by BO. The performance gap between clickability and personalisation is due to two main factors. Firstly, since we sample negatives from the recipes that received new interactions during the test period, it is likely that most of these recipes are of high quality (*i.e.* ‘clickable’) making the task more challenging for content-based models. Secondly, the pool of negatives covers a diverse set of cuisines, which makes recommendation easier for models that have learnt users preferences. Hybrid models are always selected by BO which consistently rejected CF and CB alternatives, suggesting that both preference and content are important for the task.

We experimented with a variety of alternate text embeddings to understand the source of our good personalisation performance. We found that normalised ingredient and skill were still constantly selected by BO even when other text embeddings were available for consideration. This indicates that for the recommendation task defined, targeted ingredient and skill representations are more expressive than general text embeddings.

3.3 Case Study 2: Search

In this case study, models are tasked to re-rank candidate recipe lists that were served from search queries. The served search order is known to be significantly biased [16, 37], and consequently we expect them to act as an upper bound on performance. Four rankings are considered in this case study: served order (de-biased), clickability, personalisation and the biased served order (biased). Since recipe publication is a random process and served search order is currently a strong function of recency, we can de-bias served results with randomisation. This establishes an unbiased baseline for evaluating clickability and personalisation models. BO again selects between CF, CB and hybrid models.

Table 5 presents the results of the search re-ranking experiment. We base our clickability results on the context-free model variation. We found that adding the query context to ranking models does not substantially improve performance on these metrics since the candidates presented to the model necessarily encapsulate this context already. Our clickability model out-performs baseline significantly and improved performance is obtained over all languages.

Personalisation out-performs clickability models in all cases, with average mAP@1 improvements of $\approx 20\%$. This is a vital metric in search and measures the proportion of time users engage with the top-ranked recipe. English is the weakest language for search personalisation, though it still out-performs baseline and clickability, and Arabic registers the highest improvement over clickability. The interaction experiment has higher base results than search. Although several factors contribute to this, the key explanation is that re-ranking small sets of (potentially) similar candidates for search is more challenging because candidate diversity is lower.

When evaluating search re-ranking against recipe clicks, we were unable to surpass the strong bias of served order. However, if instead we evaluate performance against other interactions (*e.g.* bookmark, cookplan) the personalisation models out-perform the (biased) served order by $\approx 10\%$. Personalisation models surpassing the strong bias of served search order is noteworthy and highlights the appropriateness of our approach to re-rank search results meaningfully.

3.4 Case Study Summary and Discussion

We tested our models extensively against several popular and competitive baseline methods (including CF and CB) and our proposed approach was exclusively selected by BO in interaction and search case studies. Popular text embedding models were also tested, but,

Table 4: Results of interaction prediction. CF and CB baseline results are not shown since they were rejected by BO.

Language	Model spec.			Random model		Clickability		Personalisation	
	Model	Ings.	Skills	mAP@1	mAP@20	mAP@1	mAP@20	mAP@1	mAP20
Arabic	Hybrid	✓	✓	0.103	0.246	0.126	0.272	0.397	0.459
English	Hybrid	✓	✓	0.104	0.248	0.159	0.309	0.280	0.373
Indonesian	Hybrid	✓	✓	0.090	0.236	0.123	0.272	0.407	0.485
Russian	Hybrid	✓	✓	0.113	0.258	0.120	0.280	0.369	0.451
Spanish	Hybrid	✓	✓	0.099	0.245	0.140	0.286	0.408	0.474

Table 5: Results on search re-ranking. CF and CB baseline results are not shown since they were rejected by BO.

Language	Served (de-biased)		Clickability		Personalisation		Served (biased)	
	mAP@1	mAP@20	mAP@1	mAP@20	mAP@1	mAP20	mAP@1	mAP20
Arabic	0.096	0.220	0.169	0.291	0.220	0.340	0.332	0.415
English	0.100	0.218	0.185	0.314	0.205	0.354	0.273	0.404
Indonesian	0.112	0.234	0.186	0.315	0.212	0.334	0.289	0.413
Russian	0.109	0.242	0.180	0.305	0.208	0.341	0.281	0.382
Spanish	0.108	0.230	0.185	0.313	0.217	0.339	0.287	0.400

disappointingly, these did not increase performance due to averaging effects over long recipe text. This exemplifies the value of targeted recipe representations in recipe recommendation. We focused on reporting qualitative performance measures in this emerging work, and broader measures (including coverage, qualitative) will be factored into more mature future presentations. The prime enabler of our success is the deliberate integration of SOTA language models and targeted ingredient and skill recipe representations.

4 CONCLUSIONS

The express objective of this paper was to develop initial understanding and expectations in multi-language recipe recommendation. Our analysis, using the most extensive dataset available for cooking and recipe recommendation, validates all representations and models with our results suggesting that multi-language recipe recommendation is suitably modelled with the proposed methodology. Despite this early work employing linear models for personalisation, our approach significantly outperforms popular content-based and collaborative baselines. We believe that we have established a strong standard for comparing the absolute value of succeeding multi-language recommendation research, and our future work will expand into three key areas. First, we will deploy our models to production systems and measure the utility of our methods in live experiments. We will then embark on an exploration of sophisticated non-linear neural recommendation frameworks and evaluate their merit. Finally, we will explore end-to-end cross-language recipe recommenders.

REFERENCES

- [1] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural Machine Translation by Jointly Learning to Align and Translate. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7–9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). <http://arxiv.org/abs/1409.0473>
- [2] J. Bobadilla, F. Ortega, A. Hernando, and A. Gutiérrez. 2013. Recommender systems survey. *Knowledge-Based Systems* 46 (2013), 109 – 132. <https://doi.org/10.1016/j.knosys.2013.03.012>
- [3] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics* 5 (2017), 135–146.
- [4] Micael Carvalho, Rémi Cadène, David Picard, Laure Soulier, Nicolas Thome, and Matthieu Cord. 2018. Cross-Modal Retrieval in the Cooking Context: Learning Semantic Text-Image Embeddings. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval* (Ann Arbor, MI, USA) (SIGIR '18). Association for Computing Machinery, New York, NY, USA, 35–44. <https://doi.org/10.1145/3209978.3210036>
- [5] Jing-Jing Chen, Chong-Wah Ngo, Fu-Li Feng, and Tat-Seng Chua. 2018. Deep Understanding of Cooking Procedure for Cross-Modal Recipe Retrieval. In *Proceedings of the 26th ACM International Conference on Multimedia* (Seoul, Republic of Korea) (MM '18). Association for Computing Machinery, New York, NY, USA, 1020–1028. <https://doi.org/10.1145/3240508.3240627>
- [6] Paula Chocron and Paolo Pareti. 2018. Vocabulary Alignment for Collaborative Agents: a Study with Real-World Multilingual How-to Instructions. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*. International Joint Conferences on Artificial Intelligence Organization, 159–165. <https://doi.org/10.24963/ijcai.2018/22>
- [7] Maurizio Ferrari Dacrema, Paolo Cremonesi, and Dietmar Jannach. 2019. Are We Really Making Much Progress? A Worrying Analysis of Recent Neural Recommendation Approaches. In *Proceedings of the 13th ACM Conference on Recommender Systems* (Copenhagen, Denmark) (RecSys '19). Association for Computing Machinery, New York, NY, USA, 101–109. <https://doi.org/10.1145/3298689.3347058>
- [8] Damion M Dooley, Emma J Griffiths, Gurinder S Gosal, Pier L Buttigieg, Robert Hoehndorf, Matthew C Lange, Lynn M Schriml, Fiona S L Brinkman, and William W L Hsiao. 2018. FoodOn: a harmonized food ontology to increase global food traceability, quality control and data integration. *npj Sci. Food* 2, 1 (2018), 23. <https://doi.org/10.1038/s41538-018-0032-6>
- [9] M. A. El-Dosuky, M. Z. Rashad, T. T. Hamza, and A. H. EL-Bassiouny. 2012. Food Recommendation Using Ontology and Heuristics. In *Advanced Machine Learning Technologies and Applications*, Aboul Ella Hassanien, Abdel-Badeeh M. Salem, Rabie Ramadan, and Tai-hoon Kim (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 423–429.
- [10] Mikhail Fain, Andrey Ponikar, Ryan Fox, and Danushka Bollegala. 2019. Dividing and Conquering Cross-Modal Recipe Retrieval: from Nearest Neighbours Baselines to SoTA. [arXiv:1911.12763 \[cs.CV\]](https://arxiv.org/abs/1911.12763)
- [11] Jill Freyne and Shlomo Berkovsky. 2010. Intelligent Food Planning: Personalized Recipe Recommendation. In *Proceedings of the 15th International Conference on Intelligent User Interfaces* (Hong Kong, China) (IUI '10). Association for Computing Machinery, New York, NY, USA, 321–324. <https://doi.org/10.1145/1719970.1720021>
- [12] Jill Freyne and Shlomo Berkovsky. 2010. Recommending Food: Reasoning on Recipes and Ingredients. In *User Modeling, Adaptation, and Personalization*, Paul

- De Bra, Alfred Kobsa, and David Chin (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 381–386.
- [13] Han Fu, Rui Wu, Chenghao Liu, and Jianling Sun. 2020. MCEN: Bridging Cross-Modal Gap between Cooking Recipes and Dish Images with Latent Variable Model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [14] Jun Harashima, Michiaki Ariga, Kenta Murata, and Masayuki Ioki. 2016. A Large-scale Recipe and Meal Data Collection as Infrastructure for Food Research. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, European Language Resources Association (ELRA), Portoroz, Slovenia, 2455–2459. <https://www.aclweb.org/anthology/L16-1389>
- [15] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. 2017. Densely Connected Convolutional Networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2261–2269.
- [16] Thorsten Joachims, Laura Granka, Bing Pan, Helene Hembrooke, and Geri Gay. 2017. Accurately Interpreting Clickthrough Data as Implicit Feedback. *SIGIR Forum* 51, 1, 4–11. <https://doi.org/10.1145/3130332.3130334>
- [17] Y. Kawano, T. Sato, T. Maruyama, and K. Yanai. 2013. Mirurecipe: A mobile cooking recipe recommendation system with food ingredient recognition. In *2013 IEEE International Conference on Multimedia and Expo Workshops (ICMEW)*, 1–2.
- [18] Maciej Kula. 2015. Metadata embeddings for user and item cold-start recommendations. In *CEUR Workshop Proc.*, Vol. 1448. arXiv, N/A, 14–21. arXiv:1507.08439 <http://groupLens.org/datasets/movielens/>
- [19] Jie Lu, Dianshuang Wu, Mingsong Mao, Wei Wang, and Guangquan Zhang. 2015. Recommender system application developments: A survey. *Decision Support Systems* 74 (2015), 12–32. <https://doi.org/10.1016/j.dss.2015.03.008>
- [20] Rui Maia and Joao C. Ferreira. 2018. Context-aware food recommendation system. In *Lect. Notes Eng. Comput. Sci.*, Vol. 2237, 349–356. <https://repositorio.iscte-iul.pt>
- [21] Bodhisattwa Prasad Majumder, Shuyang Li, Jianmo Ni, and Julian McAuley. 2019. Generating Personalized Recipes from Historical User Preferences. In *EMNLP*, 5975–5981. <https://doi.org/10.18653/v1/D19-1613>
- [22] Jonathan Malmaud, Jonathan Huang, Vivek Rathod, Nicholas Johnston, Andrew Rabinovich, and Kevin Murphy. 2015. What's Cookin'? Interpreting Cooking Videos using Text, Speech and Vision. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Denver, Colorado, 143–152. <https://doi.org/10.3115/v1/N15-1015>
- [23] Jonas Mockus. 1989. *Bayesian approach to global optimization: theory and applications*. Vol. 37. Springer, Springer, Dordrecht. <https://doi.org/10.1007/978-94-009-0909-0>
- [24] Donghyeon Park, Keonwoo Kim, Yonggyu Park, Jungwoon Shin, and Jaewoo Kang. 2019. KitcheNette: Predicting and Ranking Food Ingredient Pairings using Siamese Neural Network. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*. International Joint Conferences on Artificial Intelligence Organization, 5930–5936. <https://doi.org/10.24963/ijcai.2019/822>
- [25] Rama Kumar Pasumarthi, Sebastian Bruch, Xuanhui Wang, Cheng Li, Michael Bendersky, Marc Najork, Jan Pfeifer, Nadav Golbandi, Rohan Anil, and Stephan Wolf. 2019. TF-Ranking: Scalable TensorFlow Library for Learning-to-Rank. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (Anchorage, AK, USA) (KDD '19)*. Association for Computing Machinery, New York, NY, USA, 2970–2978. <https://doi.org/10.1145/3292500.3330677>
- [26] Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: A Python Natural Language Processing Toolkit for Many Human Languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics, Online, 101–108. <https://doi.org/10.18653/v1/2020.acl-demos.14>
- [27] Amaia Salvador, Michal Drozdal, Xavier Giro-i Nieto, and Adriana Romero. 2019. Inverse Cooking: Recipe Generation From Food Images. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [28] A. Salvador, N. Hynes, Y. Aytar, J. Marin, F. Ofli, I. Weber, and A. Torralba. 2017. Learning Cross-Modal Embeddings for Cooking Recipes and Food Images. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3068–3076.
- [29] Nuno Silva, David Ribeiro, and Liliana Ferreira. 2019. Information Extraction from Unstructured Recipe Data. In *Proceedings of the 2019 5th International Conference on Computer and Technology Applications (Istanbul, Turkey) (ICCTA 2019)*. Association for Computing Machinery, New York, NY, USA, 165–168. <https://doi.org/10.1145/3323933.3324084>
- [30] Young Chol Song, Iftekhar Naim, Abdullah Al Mamun, Kaustubh Kulkarni, Parag Singla, Jiebo Luo, Daniel Gildea, and Henry Kautz. 2016. Unsupervised Alignment of Actions in Video with Text Descriptions. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence (New York, New York, USA) (IJCAI'16)*. AAAI Press, 2025–2031.
- [31] Wanhua Su, Yan Yuan, and Mu Zhu. 2015. A Relationship between the Average Precision and the Area Under the ROC Curve. In *Proceedings of the 2015 International Conference on The Theory of Information Retrieval (Northampton, Massachusetts, USA) (ICTIR '15)*. Association for Computing Machinery, New York, NY, USA, 349–352. <https://doi.org/10.1145/2808194.2809481>
- [32] Chun-Yuen Teng, Yu-Ru Lin, and Lada A. Adamic. 2012. Recipe Recommendation Using Ingredient Networks. In *Proceedings of the 4th Annual ACM Web Science Conference (Evanston, Illinois) (WebSci '12)*. Association for Computing Machinery, New York, NY, USA, 298–307. <https://doi.org/10.1145/2380718.2380757>
- [33] Thi Ngoc Trang Tran, Müslüm Atas, Alexander Felfernig, and Martin Stettinger. 2018. An overview of recommender systems in the healthy food domain. *J. Intell. Inf. Syst.* 50, 3 (2018), 501–526. <https://doi.org/10.1007/s10844-017-0469-0>
- [34] Christoph Trattner and David Elswiler. 2019. An evaluation of recommendation algorithms for online recipe portals. In *CEUR Workshop Proc.*, Vol. 2439, 24–28. <http://www.librec.net/>
- [35] Christoph Trattner, Dominik Moesslang, and David Elswiler. 2018. On the predictability of the popularity of online recipes. *EPJ Data Sci.* 7, 1 (2018), 20. <https://doi.org/10.1140/epjds/s13688-018-0149-5>
- [36] Hao Wang, Doyen Sahoo, Chenghao Liu, Ee-peng Lim, and Steven C. H. Hoi. 2019. Learning Cross-Modal Embeddings With Adversarial Networks for Cooking Recipes and Food Images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 11572–11581.
- [37] Xuanhui Wang, Nadav Golbandi, Michael Bendersky, Donald Metzler, and Marc Najork. 2018. Position Bias Estimation for Unbiased Learning to Rank in Personal Search. In *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining (Marina Del Rey, CA, USA) (WSDM '18)*. Association for Computing Machinery, New York, NY, USA, 610–618. <https://doi.org/10.1145/3159652.3159732>
- [38] Jason Weston, Sumit Chopra, and Keith Adams. 2014. #TagSpace: Semantic Embeddings from Hashtags. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (Doha, Qatar). Association for Computational Linguistics, 1822–1827. <https://doi.org/10.3115/v1/D14-1194>
- [39] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144* (2016).
- [40] Yoko Yamakata, Shinsuke Mori, and John Carroll. 2020. English Recipe Flow Graph Corpus. In *Proceedings of The 12th Language Resources and Evaluation Conference*. European Language Resources Association, Marseille, France, 5187–5194. <https://www.aclweb.org/anthology/2020.lrec-1.638>
- [41] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. 2019. Deep Learning Based Recommender System: A Survey and New Perspectives. *ACM Comput. Surv.* 52, 1, Article 5 (Feb. 2019), 38 pages. <https://doi.org/10.1145/3285029>
- [42] Yixin Zhang, Yoko Yamakata, and Keishi Tajima. 2019. Categorization of Cooking Actions Based on Textual/Visual Similarity. In *Proceedings of the 5th International Workshop on Multimedia Assisted Dietary Management (Nice, France) (MADiMa '19)*. Association for Computing Machinery, New York, NY, USA, 42–49. <https://doi.org/10.1145/3347448.3357165>