

Extraction of ingredient names from recipes by combining linguistic annotations and CRF selection

Thierry Hamon
LIM&BIO (EA3969), Université Paris 13
Sorbonne Paris Cité
74, rue Marcel Cachin
93017 Bobigny, France
thierry.hamon@univ-paris13.fr

Natalia Grabar
CNRS UMR 8163 STL
Université Lille 1&3
59653 Villeneuve d'Ascq, France
natalia.grabar@univ-lille3.fr

ABSTRACT

Nutrition and diet have direct and considerable impact on our well-being and health. This field attracts researchers from different areas, such as medicine, nutrition and epidemiology, computer sciences, artificial intelligence and natural language processing (NLP). We process the recipes with NLP methods in order to automatically identify ingredient names within recipes. We propose a hybrid system based on linguistic enrichment of the recipes and selection of the relevant ingredient names with a CRF method. Semantic resources have been specifically built for processing two kinds of information: exact (*e.g.* quantity expressed in grams or liters, durations expressed in minutes or days) and fuzzy (*e.g.* quantities expressed in *chouilla* (*smidgeon*) and *louche* (*ladle*), durations sequenced with *après*, *ensuite*, *alors que* (*the, after that, while*)). The experiments are performed with French-language textual data. The results demonstrate that the proposed method is useful for searching and managing the recipes.

Categories and Subject Descriptors

H.3.1 [Information Storage and Retrieval]: Content Analysis and Indexing—*Indexing methods*; H.3.1 [Information Storage and Retrieval]: Content Analysis and Indexing—*Linguistic processing*; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Information filtering*; H.3.3 [Information Storage and Retrieval]: Information Search and Retrieval—*Selection process*; I.2.7 [Artificial Intelligence]: Natural Language Processing; I.2.7 [Artificial Intelligence]: Text analysis

General Terms

Languages, Documentation, Experimentation, Performance

Keywords

Text mining, Indexing, Conditional Random Fields, Recipes

Permission to make digital or hard copies of all or part 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 components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

CEA 2013 Barcelona, Spain

Copyright 2013 ACM 978-1-4503-2392-5/13/10 \$15.00.

<http://dx.doi.org/10.1145/2506023.2506035>.

1. INTRODUCTION

Nutrition and diet have direct and considerable impact on our well-being and health. For instance, it has been demonstrated that unbalanced diet, such as one based on hamburgers, hot-dogs, fries, chips and sodas, can favor several disorders, *e.g.*, cardiovascular disorders, depression, and even some kinds of cancer [9, 16, 12, 20]. This area has attracted several scientists from different research domains, such as medicine and biology, computer sciences, natural language processing, knowledge representation, robotics and visual recognition. We present here some of the existing studies.

In biomedical domain, the main effort of researchers consists to draw the attention on the fact that the nutrition has immediate and direct impact on health. For instance, in France, the nation-wide network NutriNet-Santé¹ has been built and concentrates today various initiatives related to the nutrition questions in France. Similar initiatives may exist in other countries as well. Beyond the relation between nutrition and disorders, the corresponding studies investigate also nutritional requirements specific to some pathologies, such as diabetes, to the social condition or to the geographical origin of the studied population [7, 1].

In computer sciences and Natural Language Processing (NLP) domains, the object which is addressed is the text of the recipes. The first study of the kind seems to be dedicated to the application of the Epicure system [3] to recipes. The general objective of the system is to generate linguistic description of recipes. For this, both, deep and shallow, syntactic analyses and representations are built using the unification grammar. The system relies on domain knowledge, such as objects (individual, massive, quantified or not), encoded in a dedicated ontology. The objects change their state according to the performed NLP treatments. This NLP system also processes the pronominal and anaphoric constructions. The objective of this study was the implementation of a fine-grained semantic approach for the processing of recipes.

An example of knowledge representation and engineering study has been performed within the recent Computer Cooking Contest². For instance, in 2012, the objectives of this contest were oriented on case-based reasoning. Such orientation may allow to automatically build a case-based reasoning system in order (1) to demonstrate the feasibility of the process using the narrative documents, such as recipes, and (2) to propose new recipes further to the exploitation of the ex-

¹<https://www.etude-nutrinet-sante.fr>

²<http://computercookingcontest.net>

Resource	Size
Ingredients, food	6,070 entries
Ustensils	222 entries
Actions	13,856 entries
Variants	163 entries
Stop words	616 entries
Training set	13,864 occurrences of words
Test set	2,306 occurrences of words

Table 1: Resources built and used for the NLP processing of the recipes.

isting recipes and according to the specifications of the task [5]. In [5], special attention is paid to actions (often verbs) and to their arguments. Moreover, several NLP treatments are applied, such as assignment of grammatical categories (noun, verbs...), searching for pronominal (*it*, *them*) references. 15 recipes are processed in this way.

The recipes represent also a good example of procedural texts. The verbs are also the central points of the recipes [19], around which the terms (term-based approach) or noun phrases (frame-based approach) are extracted and linked to them. The results obtained for 40 recipes show that the frame-based approach may be more efficient, although no proper evaluation is performed.

Among other studies, we can mention multimodal and interdisciplinary methods for the recognition of food from pictures of the refrigerator contents [8], or learning and modeling of cooking gestures using specially designed gloves [15].

We propose to design and apply the NLP resources and methods for the automatic detection and extraction of ingredient names within recipes written in French. Few research works have explored this task [3, 19], and the task is not well described yet. We propose to combine NLP methods and semantic resources (to provide linguistic and semantic annotations of recipes and the first set of extractions and weightings) with machine-learning system (to sort out the extracted ingredient names using their weighting and scoring). Such information can be exploited for managing and searching the recipes. In the following of the paper, we present the material built and exploited (Section 2) and the NLP methods (Section 3) we design and use. We then introduce the experimental design (Section 4) and the obtained results (Section 5). We finish with a conclusion and we draw some perspectives to the presented work (Section 6).

2. MATERIAL

We exploit several kinds of resources (Table 1). The usefulness of these resources is mainly related to the fact that some information usually appear within the ingredient context and may give important information on them:

- List of ingredient names and food is collected from available sources online³⁴⁵⁶; the French part of the UMLS Metathesaurus [14]; and documents from the training set. A difficulty we were faced to consisted to distinguish between ingredients and food (or staples):

³<http://www.bioweight.com/glucides.html>

⁴<http://www.bioweight.com/proteines.html>

⁵<http://www.centre-clauderer.com/acides-bases/femme-2.htm>

⁶<http://les.calories.free.fr/>

ingredients are typically used in the recipes, while staples usually correspond to the recipe result (their extraction may reduce the quality of the system). Nevertheless, some recipes may also require food products, such as pasta, ice-cream, sauces or cheese. To resolve the situation, we have introduced the ambiguity and categorized the concerned products both as food and as ingredients. The whole list contains 6,070 entries. In addition to the categorization of entries as *ingredients* and/or *food*, these are also manually categorized into semantic categories, such as meat, vegetables, fruits, bakery products, fish, seafood, candies.

- List of kitchen utensils is also collected from available resources online⁷⁸. These may appear within the context of ingredients. This list contains 222 entries.
- List of actions (verbs and nouns) provided by the resource Verbaction [6] and further to manually added entries. Actions also appear within the context of ingredients. This resource contains 13,856 entries.
- Recipe variant indicators indicate whether a given recipe contains optional and possible variations, such as the possibility to use one ingredient (*zucchini*) instead of another (*eggplant*). The variant ingredients cannot be considered as relevant for a given recipe. This list contains 163 manually detected variant markers.
- List of French stopwords contains 616 entries. These usually correspond to grammatical words (*le, une, ils... (the, an, they...)*) and should not be considered as possible candidates for ingredients.
- Resources for the detection of quantities of ingredients are specifically built for the purpose of this study. The quantities usually appear within the context of ingredients, such as: *250 g de sucre (250 g of sugar)*, *3 oeufs (3 eggs)*, *une bonne cuillère d'huile (one big spoon of oil)*, *100 + 50 gr de beurre (100 + 50 gr of butter)*, *1/2 l de lait (1/2 l of milk)*, *deux graines de cardamome (two seeds of cardamon)*, *un chouilla de sel (smidgeon of salt)*, *2 louches de bouillon (2 ladles of broth)*, *beaucoup de menthe (a lot of mint)*, etc. The resources build allow detection of such entities. We have distinguished standard quantities, already expressed in grams or liters, or easily convertible to these, and non standard quantities [4], not expressed in grams or liters and not directly convertible to these. In order to be able to process both types of quantities, we propose and apply heuristics and convert the non standard quantities into liters or kilograms. We rely for this on existing resources or converters available online⁹¹⁰. When several equivalences are proposed or when the conversion question has not been addressed yet, we apply our own intuition. For instance, the expressions such as *pincée (pinch)*, *soupçon (suspicion)*, *chouilla (smidgeon)*, *giclée (spurt)*, *goutte (drop)* mean that the

⁷<http://popoblog.unblog.fr/liste-ustensiles-de-cuisine-mise-a-jour-le-130808/>

⁸fr.wikipedia.org/wiki/Ustensile_de_cuisine

⁹<http://www.supertoilette.com/mesures-equivalences-culinaires.html>

¹⁰<http://webcafe.highbb.com/t1827-mesures-metriques-et-nord-americaine#13245>

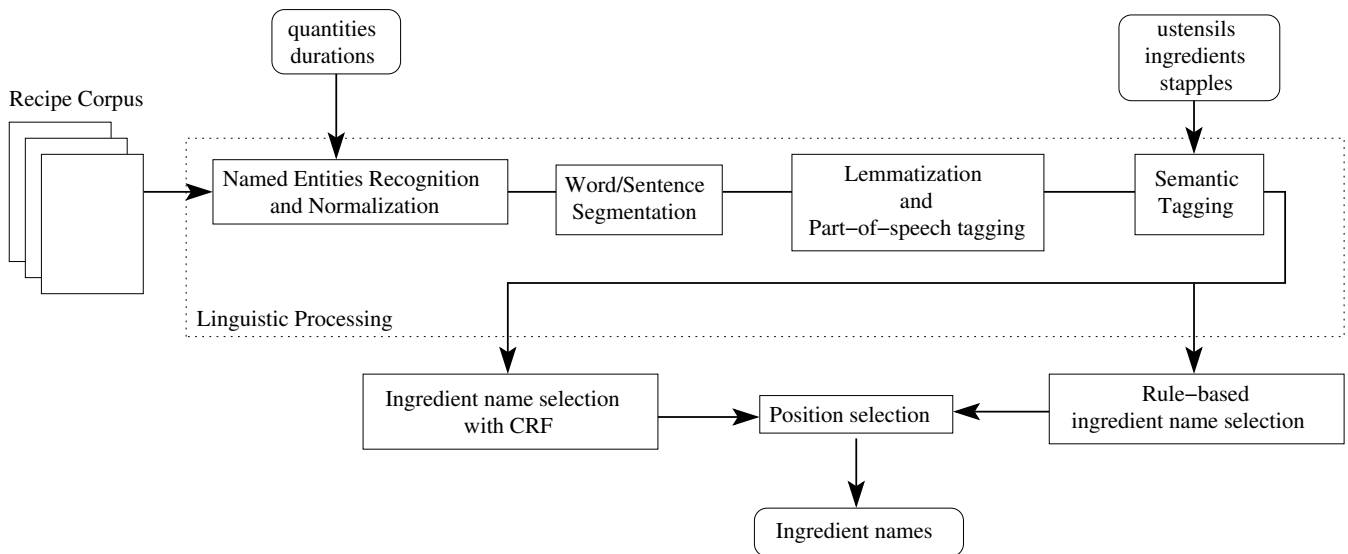


Figure 1: General schema of the method.

recipe contains 1 gram of the corresponding product; expressions such as *louche* (*ladles*), *tube* (*tube*), *assiette creuse* (*bowl*), *poignée* (*handful*) mean that the recipe contains 100 grams of the corresponding product, etc. For normalization, several arithmetic operations are applied (addition, division, multiplication).

- Resources for the detection of durations are also specifically built for the purpose of this study. Sometimes, the duration indications also appear within the context of ingredients. Here again, we distinguish between exact duration quantified and expressed in minutes, hours, days, nights, etc., and fuzzy duration non quantified as such but sequenced with expressions such as *ensuite* (*after that*), *pendant que* (*when*), *puis* (*then*), *alors que* (*while*).

The corpus we work with has been provided by the DEFT 2013 challenge¹¹. It is composed of training (13,864 recipes) and test sets (2,306 recipes). Every recipe contains the following information: title, list of ingredient names, steps required for the preparation of the recipes. The reference data correspond to the ingredient names as indicated by people who wrote the recipes.

3. NLP METHODS

We design and implement two systems: rule-based and machine learning systems. The existing tools involved in the method perform dedicated tasks, which importance is fully gained thanks to the combination of these tools and their results. Prior to the application of these systems, the recipes are pre-processed with the NLP tools. The recipes are treated with TreeTagger [17] for performing the assignment of grammatical categories (POS tagging) and lemmatization (computing the canonic form of words). This allows applying the first normalization to the processed documents. For instance, *pommes* (*apples*) is POS-tagged as plural noun and is lemmatized to *pomme* (*apple*). Optionnally,

the POS-tagging and lemmatization can post-processed with the FLEMM [13], which objective is to verify the grammatical category and lemma, and to add additional morpho-syntactic features. Moreover, terms or entries from lexical resources are also syntactically analyzed with the Y_AT_EA[2] shallow parser. Using this processing, we obtain terms which are syntactically parsed into head and expansion components. For instance, within the term *pomme verte* (*green apple*), *pomme* (*apple*) is the head component and *verte* (*green*) is the expansion component.

The rule-based system is used to recognize and extract the ingredient names and other kinds of semantic information from the recipes, as well as information needed for the machine-learning system. Machine-learning system based on the CRF [10] is applied to data extracted by rule-based system and allows performing additional and contextual filtering of the ingredient names but also selecting the correct normalized form of the ingredient name.

3.1 Rule-based system

The objective of the rule-based system is to perform the recognition of terms (*e.g.*, ingredients, kitchen utensils, food, actions) and of the associated information (mainly quantities and durations) in the recipes. The system functions through three main steps: term extraction (Section 3.1.1), ingredient name weighting (Section 3.1.2), and ingredient name selection (Section 3.1.3).

3.1.1 Extraction of ingredient names and of associated information

Resources described in Section 2 (lists of ingredient names, food, kitchen utensils, actions, recipe variants indicators...) are projected on the recipe text. The projection is done on both, lemmas and inflected forms. We use for this the Perl module `Alvis::TermTagger`¹². The named entities associated to the ingredients (quantities and durations) are also extracted. For instance, we use the resources designed for

¹¹<http://deft.limsi.fr/>

¹²<http://search.cpan.org/~thamon/Alvis-TermTagger/>

the detection of quantities of the ingredients, and of the duration time required for preparing a given recipe. At this step, all the entries of the exploited resources are recognized and extracted, together with the associated information.

3.1.2 Weighting of ingredient names

In order to identify, among all the extracted ingredient names, those which are the most important to the recipe, we weight them with several contextual methods:

- The position of the ingredients is expected to provide important information on their status. We exploit several heuristics, which give importance to different possible positions: (1) The first occurrence of the ingredient name (**position**): we consider that ingredients positioned at the beginning of the recipe should receive a higher weight than those positioned at the end of the recipe; (2) Cosine of the position of the first occurrence of ingredient names (**positionCos**): our hypothesis is that those ingredients which occur at the beginning and at the end of the recipe are the most relevant; (3) Position of the first occurrence of ingredient according to the middle of the recipe (**positionMid**): our hypothesis is that those ingredient names which are close to the middle of the recipe are the most relevant.
- Frequency of the canonic form (lemmatized form) of the ingredient name (**canon**);
- Association of the recognized ingredient names with quantities (**quant**).

At this step, several weights are associated to every ingredient name computed according to the presented methods.

3.1.3 Filtering and selection of ingredient names

Once weighted, the ingredient names can now be scored and filtered. Several criteria are applied for this:

- Computing and consideration of the lexical inclusion between the extracted ingredient names (**filtrInclLex**). Pairs, such as $\{pomme\ verte, pomme\}$ ($\{green\ apple, apple\}$) or $\{thé\ glacé, thé\}$ ($\{ice\ tea, tea\}$) are concerned. In such examples, the short ingredient name is usually the syntactic head of the large ingredient name and is lexically included in it. Lexical inclusion is computed at the syntactic level, in which case it is based upon the syntactic analysis of ingredient names, and at the level of strings, which is useful when the syntactic analysis is not efficient for the detection of neither syntactic relations nor inclusions. Ingredient names which correspond to the head components, (*pomme (apple)* or *thé (tea)*), within a larger ingredient name and which have a higher score are removed: their extended forms, *pomme verte (green apple)* or *thé glacé (ice tea)*, are considered then.
- Grouping the ingredient names according to their canonic forms (**filtrCan**) allows grouping the lemmas and their filtering on their most frequent inflected form.

Several combinations of weights, filters and scores are tested. For instance, we propose to test formulas such as: **position * quant**. When several ingredient names have the same weight, they are scored according to their **canon** frequency within the recipe. At this step, we keep only those ingredient names which satisfy the filtering and selection criteria.

3.2 CRF-based selection of ingredient names

The CRF-based system uses a CRF (Conditional random fields) classifier implementation [11] Wapiti¹³. CRFs are a class of statistical modelling method based on the hidden Markov models. The advantage of the CRFs is that they provide the possibility to process both, the features of the treated items and the features of the neighboring items. In the exploited implementation, the CRFs perform a quasi-newton optimisation with a limited memory [18] implemented through the L-BFGS algorithm. The system is applied on the rule-based system output and performs selection of the ingredient names. The setting of the CRF-based selection is the following:

- We consider the sentences as sequences;
- Each element (one-word and multi-word ingredient names) of the sequences is linguistically annotated with its inflected and lemmatized forms, its grammatical category, its semantic tag, the number of words (1 for one-word and n for multi-word ingredient names), and its co-occurrence with quantities;
- The CRF system has to predict whether the elements are ingredients, and whether the correct elements should correspond to inflected or lemmatized form.

Feature function for a given element and five elements before and after it are the following information associated to them: (1) linguistic and semantic annotation of sequence elements; (2) combination of 4-grams of lemmas and grammatical category, associated with the semantic tag and the number of words of a given element. The output of the CRF system is post-processed in order to select the correct form (inflected or lemmatized) of an ingredient name: it corresponds to the first occurrence of this ingredient within the recipe.

4. EXPERIMENTAL SET-UP

Within the training sets, the semantic tagging is performed with ingredients, ustensils, actions and recipe variants (Section 3.1.1). The lemmatized form, inflected form, or both of them are annotated as positive examples of ingredient names (Section 3.1.3). Finally, the CRF-based selection of the ingredient names is performed (Section 3.2). The experiments performed, as well as features used, are shown in Table 2. Our baseline is the results provided by the rule-based system only (BL), with large ingredient names kept.

The proposed methods have been designed and fitted on the training set. They have been then evaluated on the test set. Two versions of the recipes are analyzed: only text of the recipes which describes the steps of the preparation (without the list of ingredients), and whole narrative text of the recipes (list of ingredients and preparation steps).

Evaluation is performed with the following measures:

- Precision (percentage of correct extractions), Recall (exhaustivity of correct extractions) and F-measure (harmonic mean of Precision and Recall) in their micro and macro versions;
- MAP (mean average precision), defined as the following: $MAP = \frac{1}{N} \sum_{i=1}^N \frac{1}{n_i} \sum_{j=1}^{n_i} P(I_i^j)$, where $P(I_i^j)$ is the not interpolated precision of the ingredient name

¹³<http://wapiti.limsi.fr>

Runs	semantic tagging	annotation	lexical inclusion
RB	ingredients, ustensils	infl_form + can_form	large terms
CRF1	ingredients, ustensils, actions, variant indicator	infl_form + can_form	short and large terms
CRF2	ingredients, ustensils, actions, variant indicator	infl_form	short and large terms
CRF3	ingredients, ustensils, actions, variant indicator	can_form	short and large terms
CRF4	ingredients, ustensils, actions	infl_form	short and large terms
CRF5	ingredients, ustensils	infl_form	short and large terms
CRF6	ingredients, ustensils	infl_form + can_form	short and large terms
CRF7	ingredients, ustensils, actions	infl_form + can_form	short and large terms
CRF8	ingredients, ustensils, actions	can_form	short and large terms
CRF9	ingredients, ustensils	can_form	short and large terms
RB+ingr	ingredients, ustensils	infl_form + can_form	large terms
CRF6+ingr	ingredients, ustensils	infl_form + can_form	short and large terms
CRF6+ingr-Mdl2	ingredients, ustensils	infl_form + can_form	short and large terms

Table 2: Features exploited in different experiments

run	MAP	macro		
		Precision	Recall	F-measure
RB	0.4729	0.5252	0.7318	0.6038
CRF1	0.6036	0.7600	0.6830	0.7160
CRF2	0.6024	0.7615	0.6800	0.7151
CRF3	0.6026	0.7610	0.6812	0.7160
CRF4	0.6041	0.7642	0.6827	0.7170
CRF5	0.6159	0.7540	0.7011	0.7230
CRF6	0.6171	0.7578	0.7010	0.7246
CRF7	0.6053	0.7650	0.6825	0.7174
CRF8	0.6046	0.7621	0.6838	0.7175
CRF9	0.6157	0.7616	0.6980	0.7245
RB+ingr	0.6522	0.5488	0.8600	0.6588
CRF6+ingr	0.7582	0.7704	0.8296	0.7950
CRF6+ingr-Mdl2	0.8394	0.7833	0.9181	0.8421

Table 3: Performance on the test set

I_i^j at the rank j , N is the number of recipes, n_i is the number of ingredient names I_i^j of the recipe R_i .

5. RESULTS AND DISCUSSION

In Table 3, we indicate the results for the identification of the ingredient names in raw narrative recipes.

5.1 Rules-based system

We can observe that the baseline RB provides the highest recall (the number of ingredient names extracted is the largest), although precision is low. This experiment shows the lowest F-measure and MAP values. Parameters which appear to be suitable for the extraction of ingredient names with the rule-based system are: (1) use of lists of ingredient names built from online resources; (2) computing of the weight of ingredient names (Section 3.1.2); (3) computing of the syntactic inclusion among ingredient names; (4) computing of the inclusion among ingredient names at the level of characters; (5) application of FLEMM. Our best results are obtained with inclusions (both syntactic and at the character level), and use of FLEMM. The designed approach often allows extracting correct ingredient names, although their forms (inflected or lemmatized, short or large) and weighted positions may be incorrect. On the whole, it remains difficult to completely reproduce the reference data in the naming of the ingredient names. Among the weighting methods,

the most efficient appears to be: **position * quant**, *i.e.* when position of the first mention of an ingredient name is combined with the co-occurrence of this ingredient with quantity, both exact and fuzzy (these two provide indeed important and complementary indications). The difficulties we are still faced to are the hesitation between short and large forms of the ingredient names ($\{pomme\}$ ($\{pomme\}$), ($\{green\ apple, apple\}$) or $\{thé\ glacé, thé\}$ ($\{ice\ tea, tea\}$)), and the hesitation between inflected and lemmatized forms of the ingredient names ($\{pommes, pomme\}$ ($\{apples, apple\}$), $\{aubergines, aubergine\}$ ($\{eggplants, eggplant\}$)). In the recipes, all of these can occur, although the reference data contain only one of the possible forms. With the rule-based system, it is not always correctly chosen. Another observed difficulty is when the recipe does not provide necessary information. For instance, in some recipes, statements like *Mélanger tous les ingrédients* (*Mix all the ingredients*) is the only information available and does not provide specific information about the ingredients involved. Ingredients may also be underspecified: use of *viande* (*meat*) instead of precise piece (*e.g., steak*). In such cases, the automatic system cannot extract the relevant information.

5.2 Machine-learning system

5.2.1 Text of preparation steps

Machine learning system leads to an improvement of the global performance (runs CRF1 to CRF9), in terms of MAP and F-measure. Precision and recall show then more comparable values, although the recall values become lower than those obtained with the baseline: we loose 3 to 5%. Precision is significantly improved with up to 24 to 25%. More particularly, resources with ingredients and ustensils are sufficient for this task. As we explained in the Method section, the quantification information is used for the weighting of ingredient names: it provides very useful information. We also observed that the use of actions and variant indicators does not improve the results. This is a surprising observation. We expected indeed that, because the actions often co-occur with ingredients (*peel the apples, cut the carrots*), they should give useful indication on the presence of ingredients. It may also be possible that the actions are already used through their lemmatized or inflected forms within the studied windows. Concerning the form of the ingredient names (inflected or lemmatized), it appears that the best

solution is to let the CRF to make the choice by itself: the results show to be best then.

5.2.2 Whole text of recipes

When we use the whole text of the recipes (list of ingredients and description of preparation steps), the results are improved (runs with +ingr). We did several experiments: (1) the model learned on the description of preparation steps is applied to the whole text of recipes, and (2) specific model is built on the whole text of recipes. In both cases, the results are improved, although they are better when second solution is applied. More specifically, recall values are improved which improves the global results.

6. CONCLUSION AND PERSPECTIVES

We presented experiments performed for the automatic identification of ingredient names within raw text of recipes. The experiments have been done with the DEFT 2013 datasets, which gathers over 20,000 recipes. We exploit for this specifically built resources and two kinds of methods (rule-based and CRF-based). Combination of these two methods leads to the improvement of macro F-measure by 0.24 and of MAP by 0.37, and provides an important gain.

Several perspectives are open: (1) better study the weighting and ordering of ingredient names, and use more contextual information; (2) processing misspellings and anaphora; (3) exploiting other semantic relations between the ingredient names; (4) perform additional tests and study the influence of other attributes on the performance; (5) test the extracted ingredient names for other tasks, such as recipe search and classification; (6) propose a system for recipe advising, given the available ingredients and the health condition of users. Moreover, the CRF-based system can be applied directly on the recipes, without their preprocessing by the rule-based system: this can indicate the usefulness of information generated by the rule-based system.

7. REFERENCES

- [1] V. Andreeva, C. Martin, S. Issanchou, S. Hercberg, E. Kesse-Guyot, and C. Méjean. Sociodemographic profiles regarding bitter food consumption. cross-sectional evidence from a general french population. *Appetite*, 67(3):53–60, 2013.
- [2] S. Aubin and T. Hamon. Improving term extraction with terminological resources. In *FinTAL 2006*, number 4139 in LNAI, pages 380–387. Springer, 2006.
- [3] R. Dale. Cooking up referring expressions. In *Annual meeting on Association for Computational Linguistics*, pages 68–75, 1989.
- [4] M. Delamasure. Norme française NF UNM 00-003. système d’unités pifométriques. Technical report, NA, 2007.
- [5] V. Dufour-Lussier, F. Le Ber, J. Lieber, and E. Nauer. Automatic case acquisition from texts for process-oriented case-based reasoning. *Information Systems*, 2013. Comming soon.
- [6] N. Hathout, F. Namer, and G. Dal. An experimental constructional database: the MorTAL project. In P. Boucher, editor, *Morphology book*. Cascadilla Press, Cambridge, MA, 2001.
- [7] J. Jones-Smith, A. Karter, E. Warton, M. Kelly, E. Kersten, H. Moffet, N. Adler, D. Schillinger, and B. Lاراia. Obesity and the food environment: Income and ethnicity differences among people with diabetes: The diabetes study of Northern California (DISTANCE). *Diabetes Care*, 36(1):1200–1208, 2013.
- [8] R. Kamoda, M. Ueda, T. Funatomi, M. Iiyama, and M. Minoh. Grocery re-identification using load balance feature on the shelf for monitoring grocery inventory. In *Proceedings of the Cooking with Computers workshop (CwC)*, pages 8–18, 2012.
- [9] M. Koleva. Nutrition, nutritional habits, obesity, and prevalence of chronic diseases in workers. *Rev Environ Health*, 14(1):21–9, 1999.
- [10] J. D. Lafferty, A. McCallum, and F. C. N. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning, ICML ’01*, pages 282–289, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc.
- [11] T. Lavergne, O. Cappé, and F. Yvon. Practical very large scale CRFs. In *Proceedings the 48th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 504–513. Association for Computational Linguistics, July 2010.
- [12] S. Munsch, T. Michael, E. Biedert, A. Meyer, and J. Margraf. Negative mood induction and unbalanced nutrition style as possible triggers of binges in binge eating disorder (BED). *Eat Weight Disord*, 13(1):22–9, 2008.
- [13] F. Namer. FLEMM : un analyseur flexionnel du français à base de règles. *Traitement automatique des langues (TAL)*, 41(2):523–547, 2000.
- [14] NLM. *UMLS Knowledge Sources Manual*. National Library of Medicine, Bethesda, Maryland, 2011. www.nlm.nih.gov/research/umls/.
- [15] S. Ota, M. Soga, N. Yamamoto, and H. Taki. Design and development of a learning support environment for apple peeling using data gloves. In *Proceedings of the Cooking with Computers workshop (CwC)*, pages 7–12, 2012.
- [16] O. Pinhas-Hamiel, R. Newfield, I. Koren, A. Agmon, P. Lilos, and M. Phillip. Greater prevalence of iron deficiency in overweight and obese children and adolescents. *Int J Obes Relat Metab Disord*, 27(3):416–8, 2003.
- [17] H. Schmid. Probabilistic part-of-speech tagging using decision trees. In *Proceedings of the International Conference on New Methods in Language Processing*, pages 44–49, Manchester, UK, 1994.
- [18] N. N. Schraudolph, J. Yu, and S. Günter. A stochastic quasi-Newton method for online convex optimization. In M. Meila and X. Shen, editors, *Proc. 11th Intl. Conf. Artificial Intelligence and Statistics (AISTATS)*, volume 2 of *Workshop and Conference Proceedings*, pages 436–443, San Juan, Puerto Rico, 2007.
- [19] P. Schumacher, M. Minor, and K. Walter. Extraction of procedural knowledge from the web. In *WWW’12 Companion*, pages 739–747, 2012.
- [20] Z. Zadik. Unbalanced nutrition might be a threat to our endocrine system. *J Pediatr Endocrinol Metab*, 22(4):287–8, 2009.