



MASTER IN
COMPUTER
SCIENCE

Automating high-quality translations for Mobile Apps

Thesis subtitle

Master Thesis

Wanzenried & Stefan
from

Bern BE, Switzerland

Philosophisch-naturwissenschaftliche Fakultät
der Universität Bern

Month and Year

Prof. Dr. Philippe Cudr-Mauroux, eXascale Infolab, University of
Fribourg, Switzerland, Supervisor

Roman Prokofyev, eXascale Infolab, University of Fribourg,
Switzerland, Assistant

Abstract

Every day, hundreds of mobile applications are added to stores such as Google Play or AppStore. Many of them are intended to be used internationally, and thus require translation of the interface. At the same time, many more mobile apps already available for download in these stores. Leveraging the translation bases of existing applications, we could immediately provide high-quality translations for new apps without need to go through a human-translation process.

This project aims to extract and parse translations of existing applications to see if they can be used to translate new ones. A prototype application allows to translate existing apps from a given source to a target language. Translations are generated with different algorithms in statistical machine translation.

Contents

1	Introduction	4
2	Problem Analysis	5
2.1	Importance of Translating Mobile Apps	5
2.2	Translation Services	7
2.3	Statistical Machine Translation	7
2.4	Challenges	7
3	Related Work	8
3.1	Phrase-Based Machine Translation	8
3.1.1	Model	8
3.1.2	Learning a Phrase Translation Table	9
3.1.3	Decoding	11
3.2	Deep Learning Neural Networks	12
4	Data Analysis	13
4.1	Apps	13
4.2	Languages	14
4.3	Top Terms	15
5	Data Extraction, Preparation and Storage	16
5.1	Extract Translations	16
5.2	Preprocess Translations	17
5.2.1	Sanitizing	17
5.2.2	Tokenization and Truecasing	17
5.3	Data storage	17
5.3.1	Parallel and Monolingual Corpus	17
5.3.2	Apache Solr	18
5.3.2.1	Setup and Schema	18
5.3.2.2	Queries	19
6	Translation Process	20
6.1	Baseline System	20
6.1.1	Algorithm	20
6.1.2	Example	21
6.2	Moses	21
6.3	Tensorflow	21

<i>CONTENTS</i>	3
7 Web Application	22
7.1 Architecture	22
7.2 Frontend	22
8 Evaluation	23
8.1 BLEU Score	23
8.2 T	23

1

Introduction

Android, iPhone and Windows smart phones are available in countries all over the world. If an app is intended to be used internationally, translating it into multiple languages can be important to attract users. Users may only use the app if its text is available in their native language. Offering multiple languages is therefore crucial for paid apps to increase sales. However, translating an app requires additional effort. From a developers perspective, all strings must be isolated in order to replace them with different values for each supported language. Furthermore, translating text itself is a difficult task which usually requires the work of a professional translator.

The main idea of this thesis is to collect a lot of existing apps and extract their translations. Are we able to provide high quality translations with the help of Statistical Machine Translation (SMT)? To reduce the scope of different mobile operating systems and languages, we focus on automatic translations of Android apps from English to French and German. From now on, the term *app* in this thesis refers to a Android application. We have chosen French and German as target languages for the translations because they are natively spoken in Switzerland, which makes understanding the produced results easier. SMT is the dominant approach in automating translations and used by services like “Google Translate”¹ and Microsoft’s “Bing Translator”². In this thesis, we use different SMT algorithms to produce translations and compare their results, both manually and automatically. A prototype web application is built which serves as a playground to translate strings and apps.

¹<https://translate.google.com>

²<https://www.bing.com/translator>

2

Problem Analysis

This chapter answers the question why it is important to offer multiple languages for apps.

2.1 Importance of Translating Mobile Apps

Google encourages developers to localize their apps on their localization checklist.¹

Android and Google Play offer you a worldwide audience for your apps, with an addressable user base that's growing very rapidly in countries such as Japan, Korea, India, Brazil, and Russia. We strongly encourage you to localize as it can maximize your apps distribution potential resulting in ratings from users around the world.

The developer should first identify the countries where the app will be distributed based on the overall market size and opportunity, app category, local pricing etc. According to the identified countries, the supported languages are determined. Ideally an app supports multiple languages already before the first version is available in the app store. Users may install the app only once and forget about it, if they are not able to understand it.

Since 2013, Google offers an (human-based) app translation service to developers.² Several developers who participated at the pilot program shared their results after translating their apps:

- The developers of Zombie Ragdoll³ used this tool to launch their new game simultaneously in 20 languages in August 2013. When they combined app translation with local marketing campaigns, they found that 80% of their installs came from non-English-language users.

¹Android localization checklist: <http://developer.android.com/distribute/tools/localization-checklist.html>

²Android Developers Blog post about Googles translation service: <http://android-developers.blogspot.ch/2013/11/app-translation-service-now-available.html>

³Game Zombie Ragdoll in the Google Playstore: <https://play.google.com/store/apps/details?id=com.rvappstudios.zombieragdoll>

- Dating app SayHi⁴ Chat expanded into 13 additional languages using the App Translation Service. They saw 120% install growth in localized markets and improved user reviews of the professionally translated UI.
- The developer of card game G4A Indian Rummy⁵ found that the App Translation Service was easier to use than their previous translation methods, and saw a 300% increase with user engagement in localized apps.

In 2012, “App Annie”⁶ (former called “Distimo”), published a study where they analysed the impact of translations to iPhone apps.⁷ They found out that the use of native language is still limited in several countries, for example Italy, Russia or Brazil. Native-only apps are mostly used in Asian countries (see Figure 2.1). The study followed 200 iPhone apps that introduced native language support in a market and observed their growth. Only one week after the local language support was introduced with an app update, downloading increased by 128% and revenue by 26%.

THE IMPORTANCE OF NATIVE LANGUAGES

Proportion of Free Downloads and Revenue by Language for iPhone

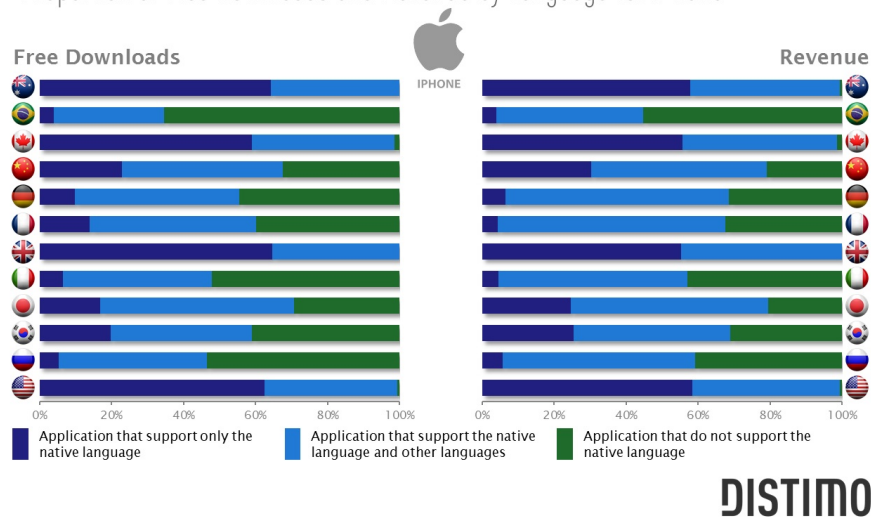


Figure 2.1: Screenshot from the study “The impact of app translations” by App Annie showing the distribution of apps supporting native and other languages for the 12 top languages.

⁴Dating app SayHi! in the Google Playstore: <https://play.google.com/store/apps/details?id=com.unearby.sayhi>

⁵Card game Indian Rummy <https://play.google.com/store/apps/details?id=org.games4all.android.games.indianrummy.prod>

⁶<https://www.appannie.com>

⁷App Annie’s study The impact of app translations: <https://web.archive.org/web/20150327084808/http://www.distimo.com/publications>

2.2 Translation Services

2.3 Statistical Machine Translation

Machine translation has a long history, but over the last decade or two, its evolution has taken on a new direction - a direction that is mirrored in other subfields of natural language processing. This new direction is grounded in the premise that language is so rich and complex that it could never be fully analyzed and distilled into a set of rules, which are then encoded into a computer program. Instead, the new direction is to develop a machine that discovers the rules of translation automatically from a large corpus of translated text, by pairing the input and output of the translation process, and learning from the statistics over the data.[2]

SMT translation models are usually trained on a large parallel corpus which is freely available on the web⁸. The BLEU metric is then used to measure the quality of the produced translations. [3] There exists a recurring translation task of the “Workshop on statistical machine translation”, which focuses on translating european language pairs by improving existing systems⁹.

In this thesis, we compare two different SMT algorithm and evaluate how they perform in the context of translating short sentences of mobile apps. The models are therefore applied on a much smaller data set than usual.

2.4 Challenges

- Quality of existing translations
- Where and how to obtain large amount of apps
- Hardware to train the translation models
- How many data is needed to produce good results?
- How can we measure good translations?

⁸<http://www.statmt.org/europarl/>

⁹<http://www.statmt.org/wmt15/translation-task.html>

3

Related Work

This chapter provides an overview of Statistical Machine Translation (SMT). We consider two different models. The first one is called “Statistical Phrase-Based Translation” and was introduced in 2003 by Philipp Koehn et al. [1][2] The second, more recent model, uses Deep Learning Neural Networks to perform translations. [4]

3.1 Phrase-Based Machine Translation

3.1.1 Model

Given the task to translate a sentence from a foreign language $f = f_1, f_2, f_3, \dots, f_m$ into an english sentence $e = e_1, e_2, e_3, \dots, e_n$, the translation model is based on the following probabilistic equation:

$$P(e|f)$$

To find the best english translation, the following equation is used (Bayes rule).

$$\begin{aligned} e_{best} &= \operatorname{argmax}_e P(e|f) \\ &= \operatorname{argmax}_e \frac{P(f|e)P(e)}{P(f)} \\ &= \operatorname{argmax}_e P(e|f)P(e) \end{aligned}$$

where $P(e)$ is called the *language model*¹ and $P(f|e)$ is the *translation model*. The translation model is further decomposed:

$$p(\bar{f}_1^I | \bar{e}_1^I) = \prod_{i=1}^I \phi(\bar{f}_i^I | \bar{e}_i^I) d(\text{start}_i - \text{end}_{i-1} - 1)$$

where ϕ is called the *phrase translation probability* and d is the *reordering probability*.

¹https://en.wikipedia.org/wiki/Language_model

The foreign sentence f is broken up into I phrases \bar{f}_i . Each foreign phrase is translated into an English phrase \bar{e}_i .

Reordering is handled by a *distance based reordering model*. We define $start_i$ as the position of the first word of the foreign input phrase that translates to the i th English phrase, and end_i as the position of the last word of that foreign phrase. The reordering distance is the number of words skipped when taking foreign words out of sequence (see figure 3.1).

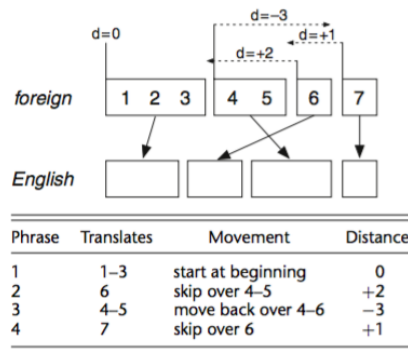


Figure 3.1: Distance based reordering [2]

d is turned into a proper probability distribution by applying an exponentially decaying cost function $d(x) = \alpha^{|x|}$ with $\alpha \in [0, 1]$. This means that movements of phrases over large distances are more expensive than shorter movements or no movement at all.

3.1.2 Learning a Phrase Translation Table

Building a phrase translation table involves three stages:

1. **Word alignment** The words of each sentence from the parallel corpus of the foreign and English language must be aligned in order to extract phrase pairs. The alignment is typically done by using one of the IBM models.² Figure 3.2 shows an example of word alignment for a top term from our corpus.
2. **Extraction of phrase pairs** Extract all phrase pairs of parallel sentences that are consistent with word alignment. A phrase pair (\bar{f}, \bar{e}) is consistent with an alignment A , if all words f_1, \dots, f_n in \bar{f} that have alignment points in A have these with words e_1, \dots, e_n in \bar{e} and vice versa.
3. **Scoring phrase pairs** Assign probabilities to phrase translations.

²IBM word alignment models: https://en.wikipedia.org/wiki/IBM_alignment_models

	Mettre	à	jour	les	services	Google	Play
Update							
the							
Google							
Play							
Services							

Figure 3.2: Word alignment between a top term of our parallel English/French corpus.

	Mettre	à	jour	les	services	Google	Play
Update							
the							
Google							
Play							
Services							

Figure 3.3: Extract phrase pairs consistent with word alignment.

We are able to extract the following phrase pairs from the example sentence:

- Update the Google Play Services / Mettre à jour les services Google Play
- Update / Mettre à jour
- Update the / Mettre à jour les
- Google Play Services / services Google Play
- The Google Play Services / les services Google Play
- Google Play / Google Play

The last step is to estimate the phrase translation probabilities by the relative frequency:

$$\phi(\bar{f}|\bar{e}) = \frac{\text{count}(\bar{e}, \bar{f})}{\sum_{\bar{f}_i} \text{count}(\bar{e}, \bar{f}_i)}$$

These probabilities are put in the phrase translation table, together with the English and foreign phrase (see Figure 3.1).

Foreign Phrase	$\phi(\bar{f} \bar{e})$
paramètres	0.9
réglages	0.36
configuration	0.15

Table 3.1: Example of a phrase translation table containing the estimated probabilities for translating the English word “Setting” to French.

3.1.3 Decoding

Decoding is the process of finding a translation for a given input sentence. Recall the model to compute the translation probability:

$$e_{best} = \underset{e}{\operatorname{argmax}} \prod_{i=1}^I \phi(\bar{f}_i|\bar{e}_i) d(\operatorname{start}_i - \operatorname{end}_{i-1} - 1) p_{LM}(e)$$

- **Phrase translation** Pick phrase \bar{f}_i to be translated as phrase \bar{e}_i by looking up the score from the phrase translation table.
- **Reordering** The previous phrase ended in end_{i-1} , the current phrase starts at start_i . Compute $d(\operatorname{start}_i - \operatorname{end}_{i-1} - 1)$
- **Language model** For a n -gram model, we need to keep track of the last $n - 1$ words. The score is computed for every added word w_i : $p_{LM}(w_i|w_i - (n - 1), \dots, w_{i-1})$

The algorithm reads the input sentence from left to right and creates a new hypothesis by picking any phrase translation at a time. The score is computed incrementally for each partial hypothesis. This is repeated recursively until all hypotheses have been expanded. A hypothesis covering all input words cannot be expanded further and forms and end point in the search graph (Figure 3.4). To get the best translation, we pick the completed hypothesis with the highest probability score and backtrack through the search graph.

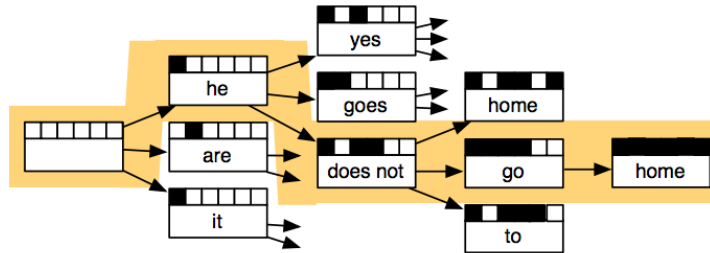


Figure 3.4: Example: Created hypothesis by translating the German sentence “er geht ja nicht nach hause” to English. The squares on top indicate a coverage vector of translated German words (filled black if covered). [2]

This process creates an exponential number of hypothesis. To reduce the search space, the following methods are applied:

- **Recombination** Two hypothesis paths lead to two matching hypotheses, e.g. have the same English words in the output. The worse hypothesis is dropped.
- **Stack Pruning** Bad hypotheses are removed early. Hypothesis that have translated the same number of input words are put into a stack, so that they are comparable. The number of hypotheses in a stack is limited, for example by keeping only k hypotheses according to a score. This score is a combination of the partial probability score and a future cost estimation of the remaining sentence.

3.2 Deep Learning Neural Networks

4

Data Analysis

4.1 Apps

A total of 698 free Android apps were collected from various categories¹, ensuring a good mix of different translations. Most of them were chosen randomly beside the top apps, mostly from Google. Figure 4.1 shows the top ten apps according to the number of available english translations. Note that the counts indicate the number of effective translations after eliminating uninteresting data during pre-processing the corpus (see chapter 5). Most of the apps have similar counts for all the three languages. Surprisingly, some apps provide more translations for german and french than english. Different counts can happen because of:

1. Not all strings were translated to all languages. Some english words remain the same in the target language and don't need to be translated, e.g. E-Mail.
2. The XML files containing the translations are out of sync, e.g. the developer forgot to add a french translation.
3. Some translations were abandoned during pre-processing.

34 apps did not isolate their strings in XML files (we assume that they are hardcoded in the programming code). This leaves a effective number of 664 apps where we could extract translations. However, 241 apps among them were not multilingual and therefore not useful for building translation systems. Unfortunately, apps are not marked as multi-language in the app store, so we only know if there are translations available after downloading and extracting the data.

¹Click on Categories to see a list of available categories at: <https://play.google.com/store/apps?hl=en>

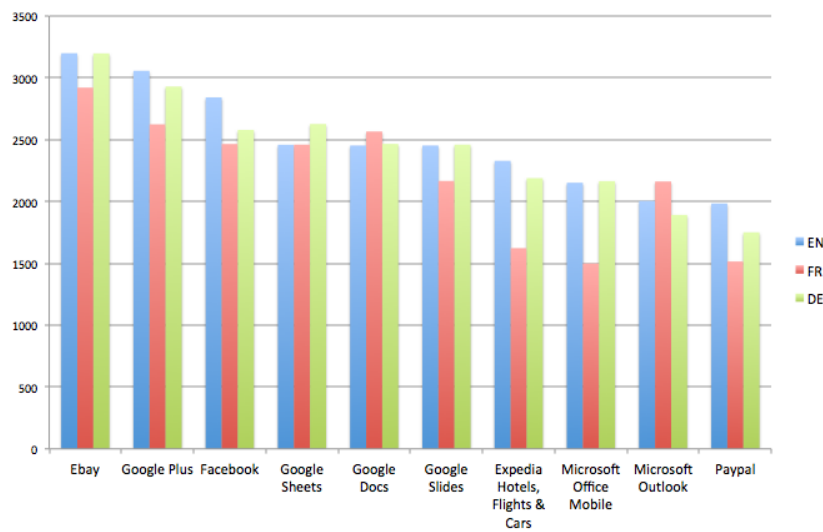


Figure 4.1: Top apps by the number of english translations

4.2 Languages

We collected a total of 423 multilingual apps, offering its content at least for two languages. Figure 4.2 shows the number of apps for the ten most used languages. French and German, where we focus on this thesis, are both in the top 4. The top apps are translated in many more languages. Overall, 117 different languages were available.

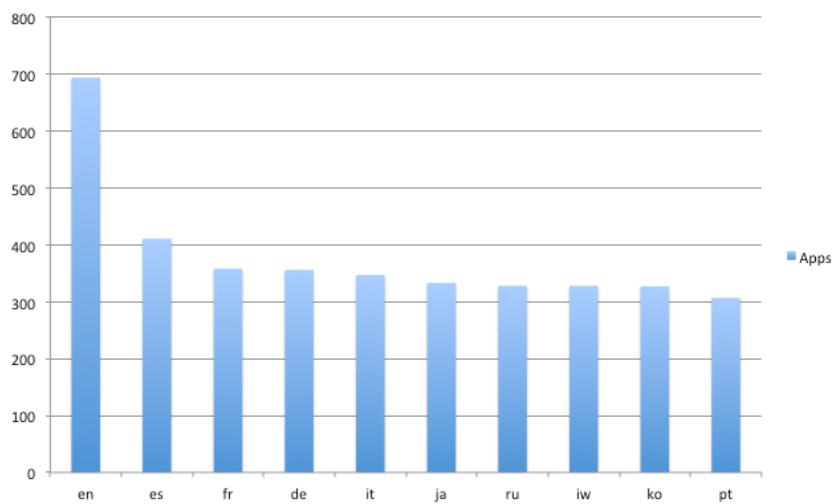


Figure 4.2: Number of apps available in the top ten languages.

4.3 Top Terms

EN	Count	FR	Count	DE	Count
Cancel	1894	Annuler	1067	Abbrechen	886
Done	1230	Supprimer	745	Google Play-Dienste aktivieren	730
Settings	1102	Connexion	745	Google Play-Dienste installieren	728
Search	884	Mettre jour	486	Anmelden	664
Delete	704	Rechercher	485	Lschen	664
Enable Google Play services	692	Paramtres	450	Aktualisieren	652
Get Google Play services	690	Mettre jour les services Google Play	370	Einstellungen	569
Log Out	556	Activer services google Play	368	Fertig	548
Update	545	Activer les services Google Play	366	Weiter	484
Share	541	Installer les services Google Play	366	Schliessen	428

Table 4.1: Top ten terms for english, french and german

5

Data Extraction, Preparation and Storage

5.1 Extract Translations

An Android app has the form of a single binary file (suffix .apk). The XML files containing the translations are packed inside the .apk file. We used a software called “Apktool”¹ to reverse engineer existing apps and extract the needed XML files.

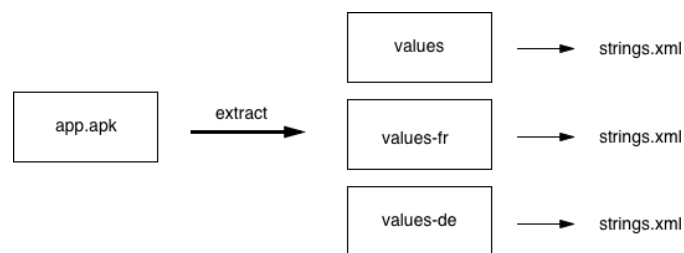


Figure 5.1: For each language, the translations are stored in a XML file

Inside the `strings.xml` files, translations are stored key-value based. Each string has a unique key which holds the translation value for each language. Here is an example how an English XML file looks like:

```
<?xml version="1.0" encoding="utf-8"?>
<resources>
  <string name="close_app">Close Application</string>
  <string name="back">Back</string>
  ...
</resources>
```

¹<http://ibotpeaches.github.io/Apktool/>

5.2 Preprocess Translations

This task involves two parts. First of all, the translations are sanitized from unwanted information such as HTML tags. Some strings are omitted at this stage because they do not provide any meaningful value to the translation process, for example URLs or string placeholders. Secondly, the sanitized strings are tokenized so that they can be used by an SMT system.

5.2.1 Sanitizing

The following rules are used to get rid of unwanted translations or parts of it:

- Trim the string from newlines (`\n`), tabs (`\t`) or carriage return (`\r`)
- Strip any HTML tags
- Omit strings that are URLs (starting with `http` or `www`)
- Remove strings placeholders like `%s` or `%1$s`
- Ignore words with less than 3 characters
- Finally, the resulting string is trimmed again from spaces and must contain at least one alphanumeric character

5.2.2 Tokenization and Truecasing

Tokenization is required to separate words from punctuation. Each token is separated by a space character: “Hi, my name is John.” becomes “Hi , my name is John .”

After tokenisation, truecasing was applied to our text corpus.² In contrast to lowercasing, truecasing determines the proper capitalization of words. This is done by first analyzing all data and then apply the most likely capitalization. It is especially useful for words starting a sentence, that are written uppercase in many languages. For example the word “This” starting a sentence will become “this”, because it is usually written lowercase.

5.3 Data storage

5.3.1 Parallel and Monolingual Corpus

After pre-processing the translations, they are written to text files, one sentence per line. The corpus is split into parallel and monolingual data. A parallel corpus (often also called bitext) describes sentence aligned text-files of one language pair. All sentences must be aligned so that line *x* of the translation source language corresponds to line *x* of the target language. The sentences are shuffled before writing to get a random mix of translations of different apps. Most SMT system use this parallel data to train their translation model, so the parallel corpus can be reused multiple times.

In addition, data of any language is also stored monolingual. All sentences of one language are written into a single text file. This is for example used to build language models, which are an important part of a phrase-based translation model (see chapter 3.1).

²<https://en.wikipedia.org/wiki/Truecasing>

5.3.2 Apache Solr

We needed a simple way to store all translations, preferably with additional meta data, to perform data analysis and calculate statistics. Apache Solr³ is used for this purpose. Solr is an open source text search platform built on top of Apache Lucene⁴. It allows to store the translations with additional meta data, such as the translation key and a unique app ID. Due to its way of indexing text data, Solr is very fast when querying for translation strings. Furthermore, built in tools allow to get nice statistics, for example the top terms grouped by language.

Solr is also used as Baseline translation system, where we translate strings with a simple algorithm directly from the collected data in Solr. (see chapter xy).

5.3.2.1 Setup and Schema

Solr offers so called “Cores” where each core manages a separate index, schema and configuration⁵. In our setup we created one core per language and one document per translation string. The schema is identical for each core (language) and consist of the following fields:

- **id** Unique ID used by Solr to identify a document.
- **app_id** Each app must have a unique app-ID.
- **key** The translation key from the XML file.
- **value** The translation value corresponding to the above key.
- **value_lc** Again the translation value, additionally indexed with a special start- and ending delimiter. This allows to search for exact values.

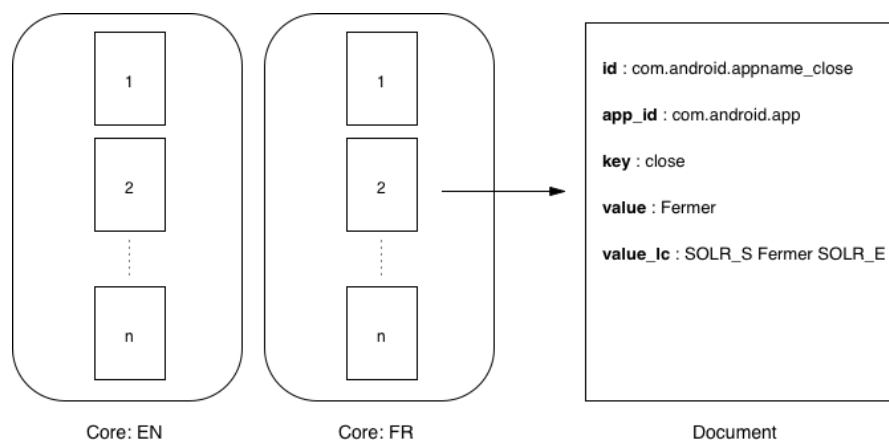


Figure 5.2: Translations are stored in a document inside a separate core per language

The data in the `value` and `value_lc` fields is stored tokenized and lowercased. The Solr Standard Tokenizer⁶ splits text into tokens, treating whitespace and punctuation as delimiters. Delimiter characters are discarded.

³<http://lucene.apache.org/solr/>

⁴<https://lucene.apache.org/core/>

⁵<https://cwiki.apache.org/confluence/display/solr/Solr+Cores+and+solr.xml>

⁶<https://cwiki.apache.org/confluence/display/solr/Tokenizers#Tokenizers-StandardTokenizer>

5.3.2.2 Queries

Solr offers a REST based interface to query the database. There is also a built in web application which allows to inspect results. Here is an example query where we search for the string “Log Out” in the English core⁷. The maximum amount of results is set to 500 rows and the return type is JSON:

```
http://diufpc114.unifr.ch:8983/solr/en/select?q=value:"Log%20Out"&rows=500&wt=json&indent=true
```

This query returns a total of 458 results. Note that there are not only exact matches returned, the term “Log Out” can occur as substring in a longer sentence (see figure xy). By default, Solr does not offer the possibility to return only exact matches. However, we can simulate this behaviour with the help of our `value_lc` field, which indexed the string together with a special start- and end delimiter (SOLR_S and SOLR_E):

```
http://diufpc114.unifr.ch:8983/solr/en/select?q=value_lc:"SOLR_S%20LogOut%20SOLR_E"&rows=500&wt=json&indent=true
```

This query returns only 396 results. Note that the case of the search term doesn’t matter and whitespace and punctuation are ignored. Executing the above query with the term “log Out.” returns the same results.

```
"response": {
  "numFound": 458,
  "start": 0,
  "docs": [
    {
      "id": "com.flaregames.gargoyles_com_facebook_loginview_log_out_action",
      "app_id": "com.flaregames.gargoyles",
      "key": "com_facebook_loginview_log_out_action",
      "value": "Log Out",
      "value_lc": "SOLR_S Log Out SOLR_E"
    },
    {
      "id": "com.instagram.android_must_log_out_one_click_login",
      "app_id": "com.instagram.android",
      "key": "must_log_out_one_click_login",
      "value": "You must log out in order to login using this link.",
      "value_lc": "SOLR_S You must log out in order to login using this link. SOLR_E"
    },
    ...
  ]
}
```

⁷You must be connected to the LAN of the University of Fribourg in order to access the referenced domain. Also note that the parameters should be fully url-encoded, this was omitted in this report for better readability.

6

Translation Process

This chapter describes the different translation systems that were used to translate apps.

6.1 Baseline System

The “Baseline System” makes direct use of the data stored in Solr and produces translations with a simple algorithm. Remember how the data is organized in Solr (see section 5.3.2): Each translation is stored in a document with a fixed schema and each document belongs to one core. To lookup a translation, the following procedure is used:

1. Search the string to translate in the source language
2. Parse the result documents and extract the `app_id` and `key` fields
3. Lookup translations in the target language by querying the target core with the collected `app_ids` and `keys`.
4. Count the number of unique translations found. The best translation is the value that was mostly used.

6.1.1 Algorithm

Given a string $S = [s_1, s_2, \dots, s_n]$ with words s_i in the source language, we want to produce a string $T = [t_1, t_2, \dots, t_n]$ in the target language. The algorithm recursively tries to translate the longest substrings of S by looking up translations for the substrings in Solr. This process is repeated until there are left only single words s_i to translate.

1. Check if we find any translation for S directly. If so, return the translation and exit.
2. Recursively start to translate substrings of S . Let l be the number of words in S at any recursion level. We check for translations of any possible substrings \bar{S}_i of length $l - 1$. If we find translations, we insert them into T at the position where the substring started. Note that if we find translations for multiple substrings, we translate the ones having the highest total count.

3. The above step is repeated until there are left only single words to translate (or in other words, substrings of length 1). If we cannot find a translation at this point, the source word s_i is inserted in T .

The above

6.1.2 Example

6.2 Moses

6.3 Tensorflow

7

Web Application

The main purpose of the web application is to offer a simple way of testing the different implemented translation systems (Baseline with Solr, Moses and Tensorflow). However, the application could already be used by app developers to translate their XML translation files.

The interface allows to translate either plain strings or a XML translation file from one language to another. Each decoder offers some specific settings where the user can influence the decoder.

7.1 Architecture

7.2 Frontend

8

Evaluation

8.1 BLEU Score

8.2 T

Bibliography

- [1] Philipp Koehn Franz, Franz Josef Och, and Daniel Marcu. Statistical phrase-based translation. pages 127–133, 2003.
- [2] Philipp Koehn. *Statistical Machine Translation*. Cambridge University Press, 2010.
- [3] Kishore Papineni, Salim Roukos, Todd Ward, and Wei jing Zhu. Bleu: a method for automatic evaluation of machine translation. pages 311–318, 2002.
- [4] Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks, 2014.