

Ugarit: Translation Alignment Visualization

Tariq Yousef

NLP Group, Institute for Computer Science, University Of Leipzig, Germany

E-mail: tariq@informatik.uni-leipzig.de

Abstract—In this paper we introduce Ugarit a web-based tool for manual translation alignment of parallel texts, the aim was to build a user interface to create training data in form of translation pairs to be used later for an automatic translation alignment system at word/phrase level, the tool now is widely used as reading environment for parallel texts and a languages learning tool. The paper gives a short overview of the visualization techniques used to present the alignment results and shows how the translation graph derived from the aligned translation pairs.

Index Terms—Translation Alignment, Data Visualization, Manual Alignment, Translations Graph.

I. INTRODUCTION

Translation alignment is one of the most popular applications of Digital Humanities and Natural Language Processing. It is defined as the comparison of two or more texts in different languages called parallel texts to find which units of the source and target texts correspond together [1]. we can distinguish between various levels of translation alignment according to the text unit to be aligned, it can be document, paragraph, sentence, or word/phrase alignment. Bilingual text alignment is an essential task in statistical Machine translation, several automatic approaches has bees developed to perform the alignment at different levels such as [16] for sentence level alignment and IBM models for word alignment. Most of them use unsupervised statistical methods to create alignment probabilities distribution between the source and target text units.

The accuracy of the automatic methods varies according to various factors such as text type, text length, size of the corpus, and translation accuracy. In contrast, Manual alignment produces more accurate results since it is performed by scholars and experts, but it is expensive in term of time and resources. The lack of user-friendly annotation tools for translation alignment and the need for accurate training data to perform the automatic alignment drove us to create Ugarit. Originally, Ugarit was developed to collect training data for the implementation of statistical machine translation system for historical languages, mainly Ancient Greek, Latin, and Persian, for which few to none aligned data sets exist. Ideally, historical languages are closed systems with a finite number of words and very limited change in the foreseeable future. Therefore, it should be possible to create adequately efficient automated

methods of statistical machine alignment based on a relatively small training data set.

The development of Ugarit started in 2017 at the AvH Chair for digital humanities at university of Leipzig, under supervision of Prof. Gregory Crane and collaboration with Chiara Palladino and Maryam Foradi.

Ugarit is crowd-sourcing project enables users to create translation alignments at word level, and the resulting translation pairs can be re-used in future machine or human translations or to create dynamic lexica and translation memory.

Since the tool was made public, the number and variety of languages included by the users has steadily increased and has gone far beyond the original intent: at the moment this paper is being written, 36 languages are included in Ugarit, and there are 295 active users, and about 23,500 parallel texts. In next section we review similar and related works, then we explain how the manual alignment can be done in Ugarit, next we show the various visualization techniques used by Ugarit to visualize the translations alignments and the dynamic lexicon search results, then we present the translations graph. Finally, we discuss the possible improvements and new features we intend to add to the next release of Ugarit.

II. RELATED WORKS

The Blinker Project [17] has developed the first annotation tool for manual text alignment to align different versions of the Bible in French and English at word level. Alpheios [4] can be considered as the first public web-based manual alignment tool for translation alignment. Alpheios [3] offers two options for the visualization of aligned texts, the side-by-side approach and the interlinear text approach.

However, there are several available softwares for automatic text alignment at sentence level such as **YouAlign**¹, and Moore’s bilingual sentence aligner [9]. Giza++ [5] uses IBM models to create automatic alignments at word level.

Visualization of aligned texts was the subject of interest and research in the recent years, many tools have been developed for this purpose. Some tools such as JuXtacommuns [13] and Versioning Machine use side-by-side approach to visualize the alignment between texts in the same language at word level. Other tools like TRAViz [11] and Stemmaweb [10] use variant graph to represent the alignment between text variants.

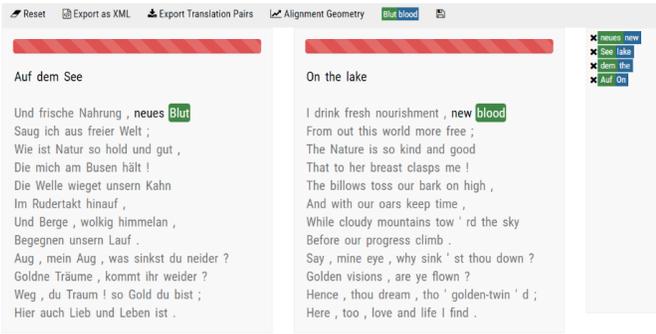


Fig. 1: Ugarit Texts Alignment Editor

Alignment Map [20] uses the parallel view to visualize translation alignment, it shows the base text and its translation, the blue connecting lines indicate alignment, the thickness of the line reflects the length of the aligned segments in words. Itéal [12] introduced the Meso reading to visualize the alignment between parallel texts at line and word level. Other tools like CATView [6] and iAligner [14] use the tabular view approach to show the positions of similarity and difference between parallel texts.

III. WORKFLOW

When we designed Ugarit Interface we tried to make it easy to use even for users who has no experience with annotation tools. Users need to create account on Ugarit and sign in to be able to start a new alignment. Users have the possibility to upload their own corpus in plain text format, or use the canonical text service (CTS) URNs to import texts from the Perseus Digital Library CTS ²repositories [19]. then the languages must be selected.

Texts will be tokenized and prepared for the next step, the alignment process is designed to be as simple as possible, we tried to minimize the mouse clicks needed to create and save the aligned tokens. To align a translation pair, users have to select the word/phrase from the original language and then select the corresponding tokens in the translated text (Fig 1). The paired tokens will be saved automatically when user starts to align new pair. On the right side of the editor, users can see all aligned pairs and have the option to edit or delete any pair. Ugarit editor allows user to create all types of alignment one-to-one, one-to-many, many-to-one, and many-to-many alignment. Resulting pairs are automatically stored in the database, and can then be exported in XML or tabular format. Users are asked to provide some information about the texts such as title, translator, description,.. etc, Users can also decide whether the alignment can be publicly visible on the website or not. Once user saves the alignment, it will appear on the home page in the "New Alignments" panel.

IV. VISUALIZATION TECHNIQUES

Visualization of translation alignment plays a great role to understand and interpret the relation between parallel texts

²cts.perseids.org

and how the aligned patterns are relevant. Ugarit uses different visualization approaches to give the users an overview of the languages and texts hosted by Ugarit and to visualize aligned texts and translation pairs in the search results.

A. Languages Graph

The graph is placed on the home page of the tool ³ to give users a quick overview of the languages currently hosted and how they are related to each other as you can see in Fig 2.

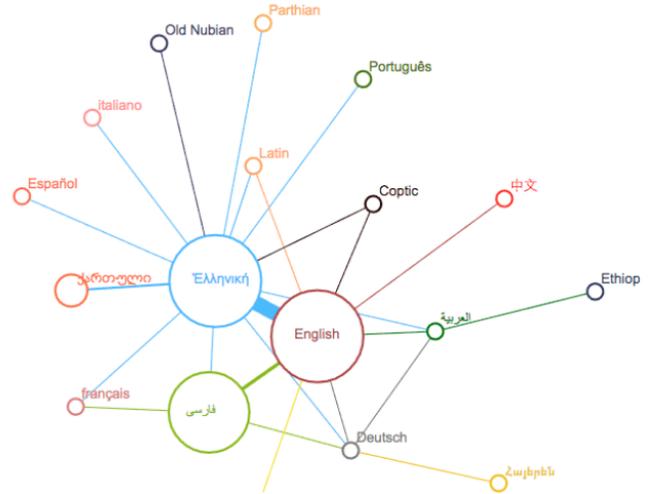


Fig. 2: Languages Graph.

Each vertex represents a language, the size of the vertex reflects the number of texts in this language in the database, each language is assigned a different color, vertices are labelled with the language names in the original form. The connection between two vertices means that there are translations alignment between texts in these two languages, and the thickness of the line reflects the number of aligned translation pairs. In Figure 2 we notice the thick blue line between English and Ancient Greek, because Ugarit contains a huge amount of English-Ancient Greek automatic aligned texts at word level provided by Perseus Digital Library [5].

Clicking on a vertex will load all texts in the language it represents. Clicking on the a line between two vertices will load all aligned texts in these two languages

B. Aligned Texts

Side-by-side visualization is the simplest and most used technique to display the alignment of parallel texts at different levels.



Fig. 3: Two texts side-by-side Visualization

³Ugarit.ialigner.com

