

MPRC: A Multilingual Parallel Reference Corpus

The last decades have seen a growing interest in use of translation corpora for language learning, translation training e.g. Bowker (2003) and corpus-based translation studies e.g. Baker (1996). Multilingual parallel corpora “consist of original texts (i.e. texts originally written in a given language), and/or translations from other languages” (Mikhailov & Cooper 2016). Most available translation corpora dealing with (Arabic, English, Chinese) languages are bitexts focusing on the alignment of two languages only (i.e. English-Arabic, English-Chinese, or Arabic-Chinese) e.g. Wang (2005). None of the available multilingual corpora deals with aligning multi-targets of parallel Arabic, English and Chinese data e.g. Tiedemann (2012). Thus, this paper introduces a new resource for use in translator training and for research on contrastive and comparative linguistics. It describes the current status of the ongoing project called Multilingual Parallel Reference Corpus (MPRC). The resource itself is a subcorpus in a conglomerate of corpora, most of them learner translator corpora, some of them reference corpora. The paper aims to present the information about the corpus design and compilation process.

One of the aims of this project is to provide comprehensive standardized metadata (i.e. author, translators, source text, target texts, and the translator’s preface). It comes in response to Saudi Arabia decision to include the Chinese language in the curriculum in schools and universities (2019). It should inform textbooks compilers of Chinese in Saudi Arabia and provide them with interesting naturalistic examples as the languages involved in this project (Arabic, English, Chinese) are not cognate and belong to different language families. MPRC is constructed to be a reference corpus to a parallel learner corpus that intends to include future productions of the learners of Chinese in Saudi universities and schools. The corpus consists of sample published texts from different genres and translation directions. Texts have been classified according to the source language and genre in a balanced way. The current corpus consists of 5 million words and manually aligned trilingual sentences. The following table indicates information about the available data in the corpus.

Table 1. MPRC preprocessing corpus

Number of Files	Source Language	Source Texts Tokens (Approx.)	First Target Tokens (Approx.)	Second Target Tokens (Approx.)
15	Arabic	600,000	455,000	365,000
15	English	590,000	555,000	490,000
15	Chinese	400,000	600,000	500,000

MPRC developed a user-friendly interface and easy to use corpus analysis tools. MPRC concordancing tool allows for retrieval of all the available examples of any query in context. The concordancing tool allows for the retrieval of the extended context of any query up to 10 sentences from the text.

فنادى بي طلب الشجر يوماً فلم أرخ عليهما حتى ناما فحكيت لهما غورقهما
 فوجدتهما نالمين فكرت أن أوقظهما وأن أعيق قلوبهما أهلاً أو ماله، فليئتُ -
 والقدح على يدي- أنتظر استيقاظهما حتى نزلَّ الفجر والصبية يتصاعون عند
 قدمي فاستيقظا فشرىا غورقهما.

One day, I went far away in search of grazing and could not come back until they had slept. When I milked as usual and brought the drink I found them both asleep. I hated to disturb them and also disliked to give milk to my children before them. My children were crying out of hunger at my feet but I awaited with the bowl in my hand for them to wake up. When they awoke at dawn, they drank milk.

有一天我赶着羊群去远处吃树叶，没有按时回来，直到二老已入睡。我为老人煮好奶子，发现二老仍熟睡，我既不愿唤醒，又不愿在老人之前让家人先喝。于是我双手端着奶子等待老人醒来。一直到黎明，孩子们饿得围在我脚前哭闹。老人终于醒了，喝了奶子。

اللهم إن كنت فعلت ذلك ابتغاء وجهك ففَرِّجْ عنا ما نحن فيه من هذه الصخرة،
 فدفعرت شيئاً لا يستطيعون الخروج منه.

O Allah! If I did so to seek Your Pleasure, then deliver us from the distress caused by the rock! The rock moved slightly but they were unable to escape.

安拉啊！如果我这样做只是为了求得你喜悦，就求你挪开堵住我们的这块岩石吧！石头果然挪动了一下，但仍不能出去。

مَعْقُوقٌ عَلَيْهِ [Al-Bukhari and Muslim] - 两大圣训集

AUTHOR TRANSLATOR SOURCE TEXT TARGET TEXT SECOND TRANSLATOR SECOND TARGET TEXT PREFACE

فنادى بي طلب الشجر يوماً فلم أرخ عليهما حتى ناما فحكيت لهما غورقهما
 فوجدتهما نالمين فكرت أن أوقظهما وأن أعيق قلوبهما أهلاً أو ماله، فليئتُ -
 والقدح على يدي- أنتظر استيقاظهما حتى نزلَّ الفجر والصبية يتصاعون عند
 قدمي فاستيقظا فشرىا غورقهما.

One day, I went far away in search of grazing and could not come back until they had slept. When I milked as usual and brought the drink I found them both asleep. I hated to disturb them and also disliked to give milk to my children before them. My children were crying out of hunger at my feet but I awaited with the bowl in my hand for them to wake up. When they awoke at dawn, they drank milk.

有一天我赶着羊群去远处吃树叶，没有按时回来，直到二老已入睡。我为老人煮好奶子，发现二老仍熟睡，我既不愿唤醒，又不愿在老人之前让家人先喝。于是我双手端着奶子等待老人醒来。一直到黎明，孩子们饿得围在我脚前哭闹。老人终于醒了，喝了奶子。

Author	Translator	Target text	Second Translator	Second target text
Name: Abu Zakaria Nawawi Native Language: Arabic Other Spoken Languages: null Nationality: null Level of Education: null Speciality: null Gender: null	Name: Ibrahim Ma'Rouf Native Language: null Other Spoken Languages: null Nationality: null Level of Education: null Speciality: null Gender: null	Title: The Meadows of the Prophet Years of publication: 2000 Language: English Genre: null Sub genre: Book Publisher: Dar Al-Manar Country: Country Link: Link Error: Error Positive practice: Positive Practice	Name: AbuAbdullah Alsenie Native Language: Arabic Other Spoken Languages: null Nationality: null Level of Education: null Speciality: null Gender: null	Title: 利雅得圣训集 Years of publication: Year of publication Language: Chinese Genre: Relegious Sub genre: Book Publisher: Alwarraq Country: Country Link: Link Error: Error Positive practice: Positive Practice

AUTHOR TRANSLATOR SOURCE TEXT TARGET TEXT SECOND TRANSLATOR SECOND TARGET TEXT PREFACE

Figure 1. MPRC Annotation Layers and Concordancing tool

Additionally, MPRC developed a multilingual collocation tool that can extract formulaic sequences that can be bigrams or trigrams (see Figure 3).

QUERY TOOL 1	QUERY TOOL 2	QUERY TOOL 3
تشيخ	old	老
كبير	man	头
All ▾	All ▾	All ▾

Figure 3. MPRC Multilingual collocation tool

This paper addresses MPRC design criteria and compilation process of data gathered in the first phase of project. The main aim of this project is to make data accessible for teaching and research.

References

- Baker, M. (1996). Corpus-based translation studies: The challenges that lie ahead. In H. Somers (Ed.), *Terminology, LSP and translation. Studies in language engineering in honour of Juan C. Sager* (pp. 175–186). Amsterdam: Benjamins.
- Bowker, L., and Peter, B. (2003). “Student Translation Archive and Student Translation Tracking System: Design, Development and Application.” In *Corpora in Translator Education*, edited by Federico Zanettin, Silvia Bernardini, and Dominic Stewart: 103–119. Manchester: St Jerome Publishing.
- Mikhailov & Cooper. (2016). *Corpus linguistics for translation and contrastive studies: A guide for research*. Routledge. Corpus Linguistics Guides. London & New York: Routledge
- Tiedemann, J. (2012). Parallel data, tools and interfaces in OPUS. In *Proceedings of the 8th International Conference on Language Resources and Evaluation (LREC '12)* (pp. 2214–2218). Istanbul: European Language Research Association.
- Wang, L. (2005). E-C Concord. Hong Kong Institute of Education. <http://ecconcord.ied.edu.hk/paraconc/index.htm>.

Dear Reem Alfuraih,

We thank you for your submission to the 15th Teaching and Language Corpora (TaLC) Conference Corpus Linguistics International Conference 2022 which will take place at the [University of Limerick](#) in conjunction with [Mary Immaculate College](#).

We are delighted to inform you that your paper (details given below) has been accepted for this year's conference. Comments from the Reviewer Panel can now be found linked to your submission in the original Abstract Submission Portal. This portal will re-open on 11th April 2022 for one week only to facilitate any edits/revisions suggested by reviewers.