

# GIẢI PHÁP XÂY DỰNG NHANH NGỮ LIỆU SONG NGỮ TỪ TED

*Hoàng Khuê<sup>1</sup>, Lương An Vinh<sup>2</sup>*

<sup>1</sup>Trường Đại học Khoa học Tự nhiên, ĐHQG-HCM

<sup>2</sup>Khoa Công nghệ Thông tin, Đại học Công nghệ Sài Gòn

<sup>1</sup>[hkhueedu@gmail.com](mailto:hkhueedu@gmail.com), <sup>2</sup>[vinh.luongan@stu.edu.vn](mailto:vinh.luongan@stu.edu.vn)

## Tóm tắt:

Hiện nay, ngữ liệu song ngữ đóng vai trò quan trọng đối với việc nghiên cứu và xử lý ngôn ngữ tự nhiên. Tuy nhiên, việc tự xây dựng hoặc mua ngữ liệu song ngữ có sẵn với chất lượng tốt sẽ tốn nhiều chi phí, đặc biệt đối với các cặp ngôn ngữ nghèo tài nguyên như tiếng Việt (vi) – Quốc Tế Ngữ (eo). Vì vậy, trong bài viết này, chúng tôi sẽ trình bày ý tưởng sử dụng tiếng Anh, một ngôn ngữ giàu tài nguyên, làm trung gian để xây dựng ngữ liệu song ngữ cho bất kì cặp ngôn ngữ nào từ dữ liệu của TED<sup>1</sup>:

TED  $\xrightarrow{\text{tài dữ liệu en-L2}}$  {en-vi;en-eo}  $\xrightarrow{\text{\{vi-en\} so khớp \{en-eo\} với trung gian en}}$  {vi-eo}

Bằng giải pháp trên, chúng tôi đã xây dựng được một kho ngữ liệu song ngữ gồm 7 cặp ngôn ngữ: vi-fr; vi-zh\_cn; vi-ja; vi-ko; vi-de; vi-ru; vi-eo, với tổng số 791056 cặp câu.

Từ khóa: ngữ liệu, song ngữ, TED

---

<sup>1</sup> TED là tổ chức phi lợi nhuận hoạt động với mục tiêu truyền bá các ý tưởng thông qua các buổi nói chuyện ngắn (18 phút đồ lại). TED bắt đầu từ năm 1984 như một hội nghị quy tụ các chủ đề về công nghệ giải trí và thiết kế. Ngày nay thì các chủ đề mở rộng ra mọi lĩnh vực từ khoa học đến kinh doanh bằng hơn 100 ngôn ngữ. <https://www.ted.com/about/our-organization>

# BUILDING QUICKLY PARALLEL CORPUS FROM TED

*Hoàng Khuê<sup>1</sup>, Lương An Vinh<sup>2</sup>*

<sup>1</sup>University of Science, VNU-HCM

<sup>2</sup>Saigon Technology University, STU

[hkhueedu@gmail.com](mailto:hkhueedu@gmail.com), [vinh.luongan@stu.edu.vn](mailto:vinh.luongan@stu.edu.vn)

## Abstract:

Nowadays, parallel corpora play an important role in researching and processing natural language. However, building or buying an available qualified parallel corpus may be costly, especially low resource language pairs like Vietnamese (vi) – Esperanto (eo). There fore, in this article, we will present the idea of using English, a rich resource language, as an intermediate language for building parallel corpus for any language pairs from TED<sup>i</sup> data.

$$\text{TED} \xrightarrow{\text{download } \{en-L2\}} \{en-vi;en-eo\} \xrightarrow[\text{with } en \text{ is intermediate}]{\{vi-en\} \text{ matched } \{en-eo\}} \{vi-eo\}$$

With this solution, we built a parallel corpus consisting of seven language pairs: vi-fr; vi-zh\_cn; vi-ja; vi-ko; vi-de; vi-ru; vi-eo, with 791.056 pair sentences in total.

Key word: corpus, parallel, TED

---

<sup>i</sup> TED is a nonprofit devoted to spreading ideas, usually in the form of short, powerful talks (18 minutes or less). TED began in 1984 as a conference where Technology, Entertainment and Design converged, and today covers almost all topics — from science to business to global issues — in more than 100 languages. <https://www.ted.com/about/our-organization>