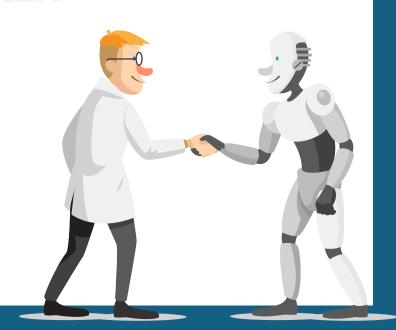
机辅翻译课程·对外经济贸易大学

Week 5: 语料库工具及语料库创建入门

黄婕 2024年10月16日 Class 3







- § 语料库的创建
- § 术语的导出和对齐
- § 本地单语、双语语料库工具
- § 英文词性标注工具
- § 语料库对齐工具
- § 翻译记忆库的创建和导出

任务场景

为了完成一篇关于人工智能的科技文章的翻译,你需要搜索在线资料,然后使用语料库工具来提高翻译效率。具体包括:

- 1) 建立一个双语的语料库
- 2) 筛选出其中的术语,生成术语表
- 3) 基于双语的语料,生成双语翻译记忆库

请思考: 你需要进行什么操作? 过程中用到哪些工具?

本课学习的软件

语料库工具

Sketch Engine

强大的云端语料库工具

AntConc 4.1.1

北外ParaConc (tools)

过程工具

中文分词工具

英文词性标注Tree Tagger

Abbyy Aligner 2.0



1. 语料库基础

语料库是什么? Corpus/Corpora

按照明确的目的和设计要求,根据语言学或翻译学原则,运用科学合理的技术方法,将一定规模的语言文本汇总而成的<u>电子文本库</u>。(管新潮,陶友兰,《语料库与翻译》,2017)

运用计算机技术,按照一定的语言学规则,根据特定的语言研究目的而大规模收集并存储在计算机中的真实语料,这些语料经过<u>一定程度的标注,便于检索</u>,可用于<u>描述研究和实证研究</u>。(王克非,《语料库翻译学探索》,2012)

语料库基础: 分类

通用语料库/专用语料库

文字语料库/ 声音语料库

单语语料库/ 平行语料库/ 对比语料库 用于<u>翻译</u>研究的语料库/ 用于<u>翻译实</u> 践的语料库

常见的语料库

语言学或翻译学研究类语料库

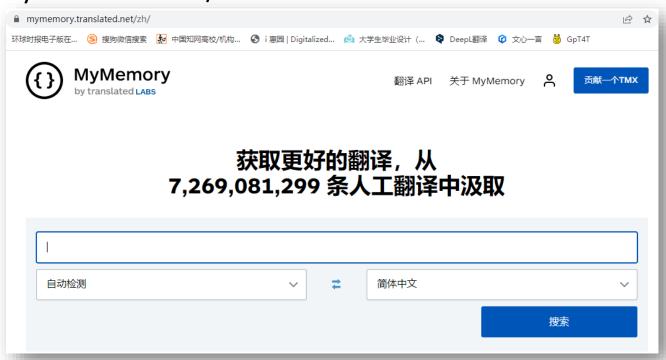
- 布朗(Brown)美国英语语料库(F. Nelson和H. Kucera,世界最早的计算机语料库)
- 上海交通大学科技英语语料库(JDEST, 上海交通大学杨惠中)
- •翻译英语语料库(世界第一个翻译语料库,英国曼彻斯特大学, Mona Baker)
- 英国国家语料库(BNC)
- 美国当代英语语料库(COCA,美国杨百翰大学Mark Davies)
- 现代汉语语料库(中国国家语言文字委员会)
- 通用汉英对应语料库(北京外国语大学,王克非)
- 两岸三地英汉科普历时平行语料库(上海交通大学)
- 英汉医学平行语料库(上海交通大学,管新潮)
- 中国法律法规汉英平行语料库(绍兴文理学院)
- 汉学文史著作英汉平行语料库(山东师范大学,徐彬)

常见的语料库

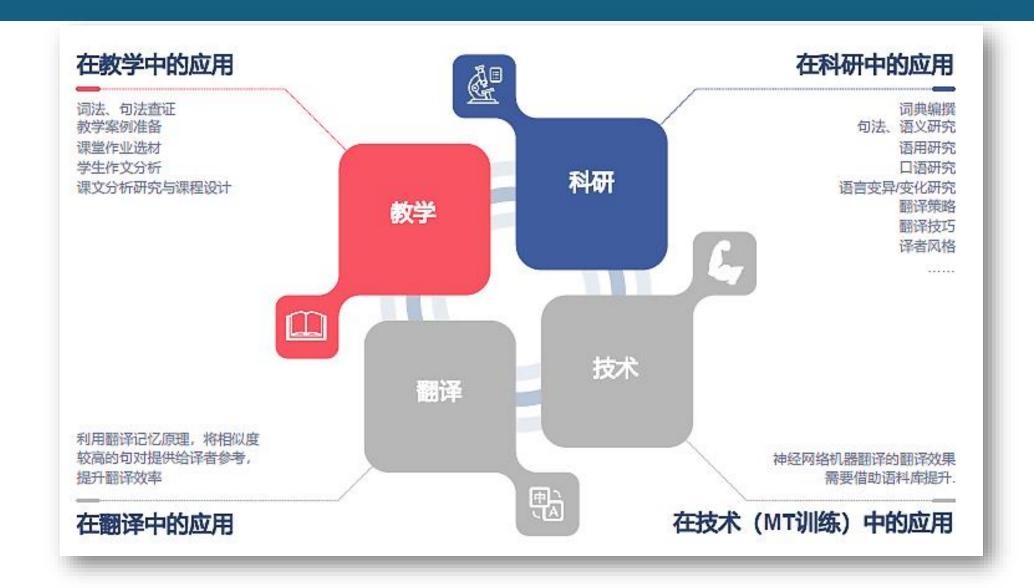
翻译实践应用类语料库

- 欧洲议会平行语料库(12亿单词,21种欧洲语言)
- TAUS Data (700亿单词, 2200个语言对)
- MyMemory (https://mymemory.translated.net/)





语料库的功能



语料库术语和功能

KWIC -- keyword in context

KWIC/concordance search 关键词索引: 特定词/短语的上下文使用案例

Word Frequency lists 词频列表: 词或词组频繁出现的次数

Collocation analysis 搭配分析:某个词的前后搭配关系

WordList 词表:按照字母或词频顺序列出给定文本的词表,统计其文本特征。可以用于翻译准备阶段的术语提取

POS tagging 词性标注: (Part-of-speech tagging) 对语料库内的单词按照语义和

上下文进行标记,即标注其词性

Keyword analysis 关键词分析

Term extraction 术语提取

双语平行语料库的应用

语料是基础

<mark>翻译风格</mark>研究

<mark>句法结构</mark>研究

<mark>跨文化</mark>传播研究、<mark>法律/商务</mark>话语研究、<mark>文化交际和形象</mark>研究、<mark>翻译研究</mark>以及翻译实践

构建<mark>语料</mark>检索、匹配、加工、利用、交易、开放的<mark>系统平台</mark>

制定<mark>专业词表、术语</mark>(单语、双语)

统计<mark>字频、词频</mark>,编写教材

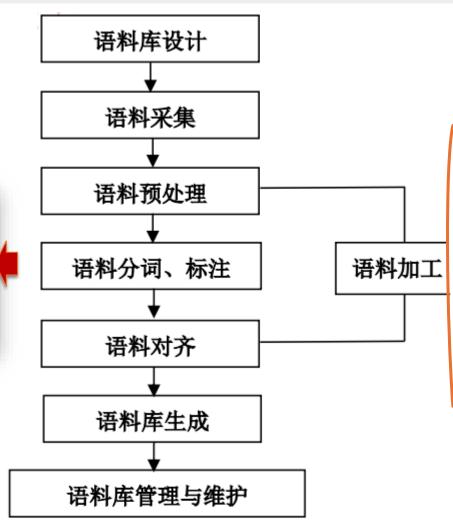
<mark>辅助</mark>写作,辅助翻译

训练机器翻译

语料库建设流程

Corpora	Li	Le	reaches	Benz	Inc
pku	李	乐	到达	奔驰	公司
msr	李乐		到达	奔驰公司	
as	李樂		到達	賓士	公司
cityu	李樂		到達	平治	公司

Table 1: Illustration of different segmentation criteria of SIGHAN bakeoff 2005.



什么是"对齐"?

- 在源语文本和目的语文本具 体单位之间建立的对应关系。
- 可分为词汇、语块、语句、 段落和篇章等层次对齐。

图 1:语料库建设流程图

中国翻译协会,语料库通用技术规范,2018

2. 语料库的创建



任务场景

为了完成一篇关于人工智能的科技文章的翻译,你需要搜索在线资料,然 后使用语料库工具来提高翻译效率。 具体包括:

- 1) 建立一个双语的语料库
- 2) 筛选出其中的术语,生成一份术语对齐表。



先获取语料:

- 1. 客户提供参考文件
- 2. 网络搜索可靠资料

Step I: 获取双语专业文档——搜索WIPO数据库

https://patentscope.wipo.int/



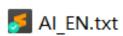
Step I: 获取双语专业文档——保存为本地文档

摘要

(EN) The embodiment of the invention provides a training method of a large language model and a text processing method based on the large language model, and relates to the fields of artificial intelligence, cloud technology, natural language processing, machine learning and the like, in particular to a language model in a pre-training language model. The method comprises the following steps: acquiring a training set and a pre-training language model corresponding to each task in a plurality of natural language processing tasks in the same target field, and acquiring a second feature extraction network corresponding to each task; and repeatedly performing training operation on the second feature extraction network corresponding to the task based on the training set corresponding to the task until a training ending condition is met, obtaining a trained second feature extraction network corresponding to the task, and obtaining the target target based on the pre-training language model and the trained second feature extraction network corresponding to each task. And obtaining a target large language model of the target domain. Based on the method, the accuracy of the text processing result output by the large language model can be improved.

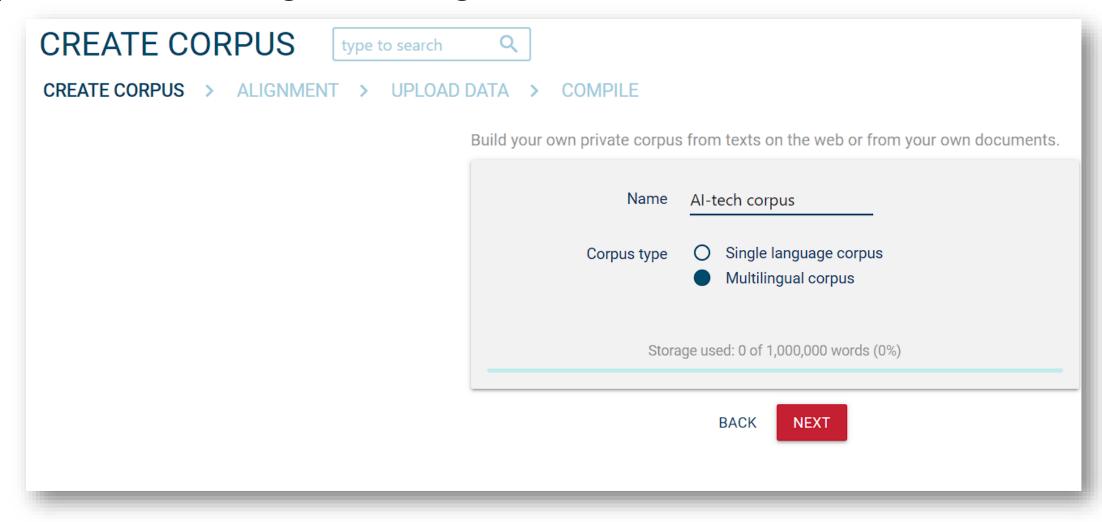
(ZH) 本申请实施例提供了一种大语言模型的训练方法及基于大语言模型的文本处理方法,涉及人工智能、云技术、自然语言处理及机器学习等领域,尤其涉及预训练语言模型中的语言模型。该方法包括:获取同一目标领域的多个自然语言处理任务中每一任务对应的训练集和预训练语言模型,获取每一任务对应的第二特征提取网络,对于每一任务,基于该任务对应的训练集对该任务对应的第二特征提取网络重复执行训练操作,直至满足训练结束条件,得到该任务对应的训练后的第二特征提取网络,基于所述预训练语言模型和各所述任务对应的训练后的第二特征提取网络,得到所述目标领域的目标大语言模型。基于该方法,可以提高大语言模型输出的文本处理结果的准确性。

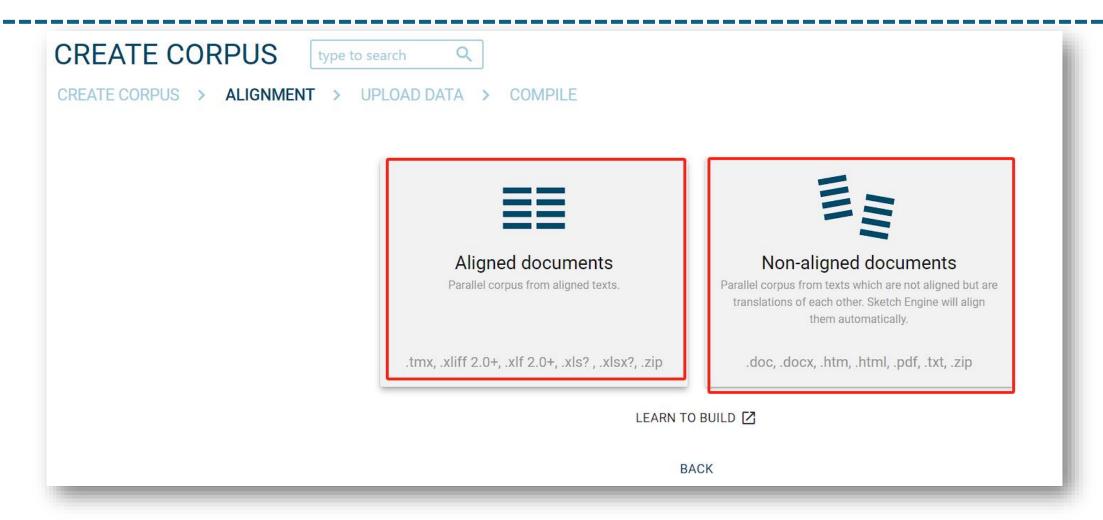




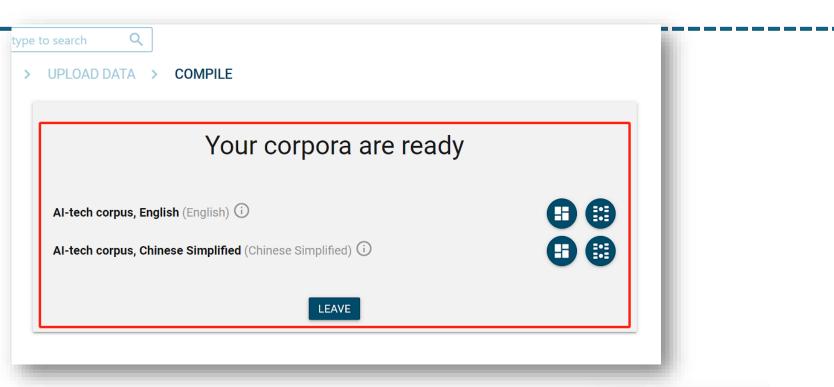


https://auth.sketchengine.eu/#login



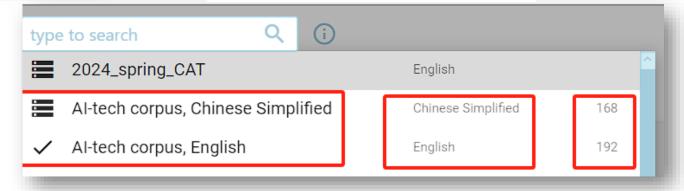




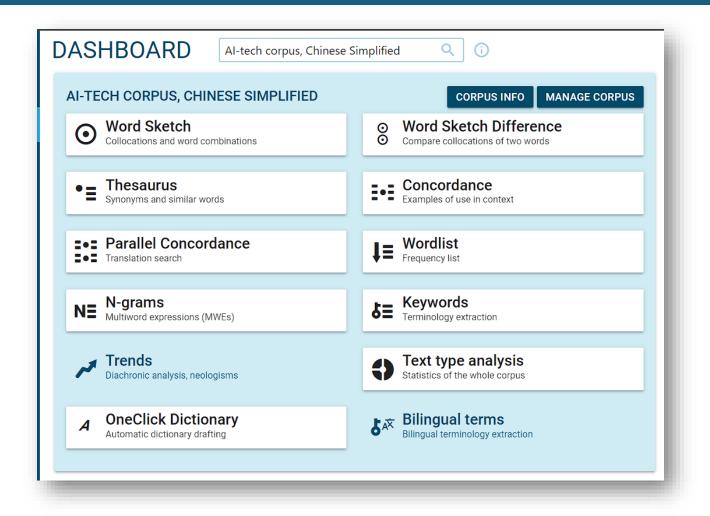


创建成功





Sketch engine 的基本功能



Word sketch 词汇素描:

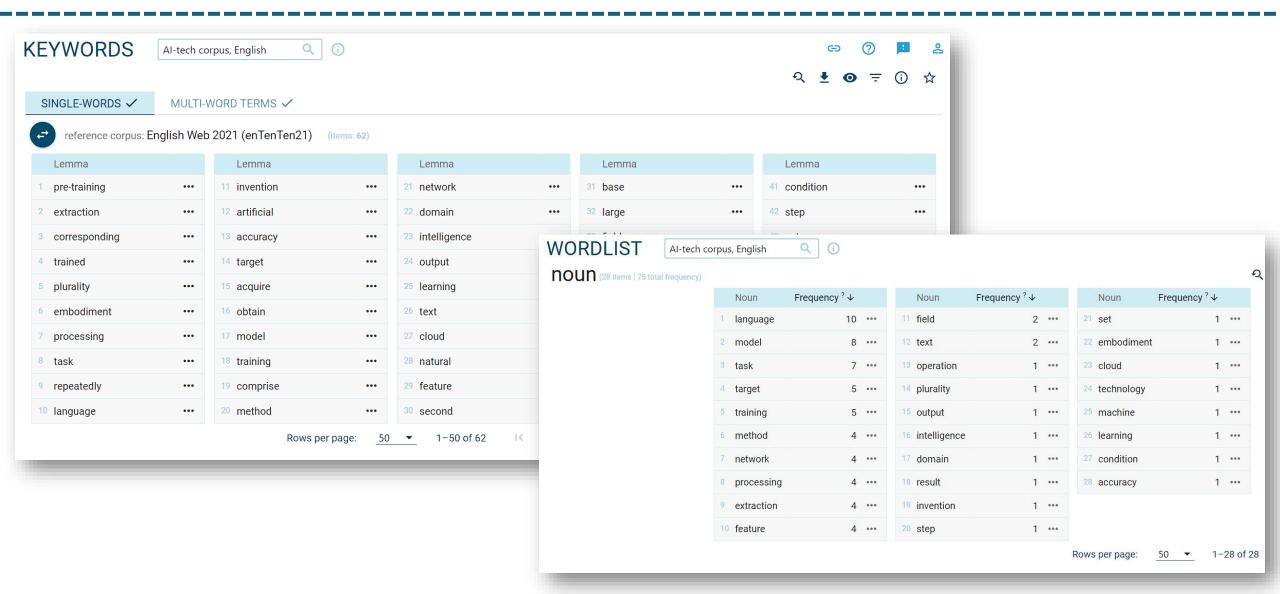
展示一个单词的语法和 搭配行为,包括常见的修饰词、宾语、主语等。

Thesaurus 同义词:

用于查找同义词,特别适用于写作时遇到词穷的情况。

Wordlist 词频列表 Keywords 关键词术语表

Sketch engine: 关键词、词频(自动词性标注)





3. 术语的导出和对齐

Step 3 导出术语、导出双语记忆库文件





Plain text

Without part-of-speech tags and lemmas



Vertical

One token per line with part-ofspeech tags and lemmas



TMX

For aligned multilingual corpora

More settings ▼

CLOSE

Step 3 导出术语、导出双语记忆库文件

vert 文件: 术语的词性标注

```
ai_tech_corpus_english.vert ×
 kdoc id="file32377953" filename="AI EN.txt" parent folder="upload">
 <s>
 The DT the-x
 embodiment NN
                embodiment-n
 of IN of-i
 the DT the-x
 invention NN invention-n
 provides
            VVZ provide-v
    DT a-x
 training
            NN training-n
 method NN method-n
 of IN of-i
                       术语的词性标注+详细信息
 a DT a-x
       JJ large-j
 large
 language
            NN language-n
 model NN model-n
 and CC and-c
 a DT a-x
        NN text-n
 text
 processing NN processing-n
 method
        NN method-n
 based
        VVN base-v
 on IN on-i
 the DT the-x
       JJ large-j
 large
            NN language-n
 language
        NN model-n
 model
 <g/>
```

Step 3 导出术语、导出双语记忆库文件

tmx 文件: 双语翻译记忆库

```
ai_tech_corpus_english.tmx ×
    <?xml version="1.0" encoding="UTF-8" standalone="no" ?>
    <tmx version="1.4"><header /><body>
    <tuv xml:lang="en"><seg>The embodiment of the invention provides a training method of a
    large language model and a text processing method based on the large language model, and
    relates to the fields of artificial intelligence, cloud technology, natural language
    processing, machine learning and the like, in particular to a language model in a
    pre-training language model.</tuv>
   <tuv xml:lang="zh-Hans"><seg>本 申请 实施 例 提供 了 — 种 大 语言 模型 的 训练 方法 及 基于 大
    语言 模型 的 文本 处理 方法 ,涉及 人工 智能 、云 技术 、自然 语言 处理 及 机器 学习 等 领域 ,
    尤其 涉及 预 训练 语言 模型 中 的 语言 模型 。</seg></tuv>
    </tu>
   <tu>>
    <tuv xml:lang="en"><seg>The method comprises the following steps: acquiring a training set
    and a pre-training language model corresponding to each task in a plurality of natural
    language processing tasks in the same target field, and acquiring a second feature
    extraction network corresponding to each task; and repeatedly performing training operation
    on the second feature extraction network corresponding to the task based on the training set
    corresponding to the task until a training ending condition is met, obtaining a trained
    second feature extraction network corresponding to the task, and obtaining the target target
    based on the pre-training language model and the trained second feature extraction network
    corresponding to each task.</seg></tuv>
9 <tuv xml:lang="zh-Hans"><seg>该 方法 包括 : 获取 同一 目标 领域 的 多 个 自然 语言 处理 任务 中
    oldsymbol{a} 一 任务 对应 的 训练 集 和 预 训练 语言 模型 ,获取 oldsymbol{a} 一 任务 对应 的 第二 特征 提取 网络 ,
    对于 每 一 任务 , 基于 该 任务 对应 的 训练 集 对 该 任务 对应 的 第二 特征 提取 网络 重复 执行
    训练 操作 ,直 至 满足 训练 结束 条件 ,得到 该 任务 对应 的 训练 后 的 第二 特征 提取 网络 ,基于
    所述预 训练 语言 模型 和 各 所 述 任务 对应 的 训练 后 的 第二 特征 提取 网络 , 得到 所述 目标 领域
    的 目标 大 语言 模型 。</seg></tuv>
10 </tu>
```

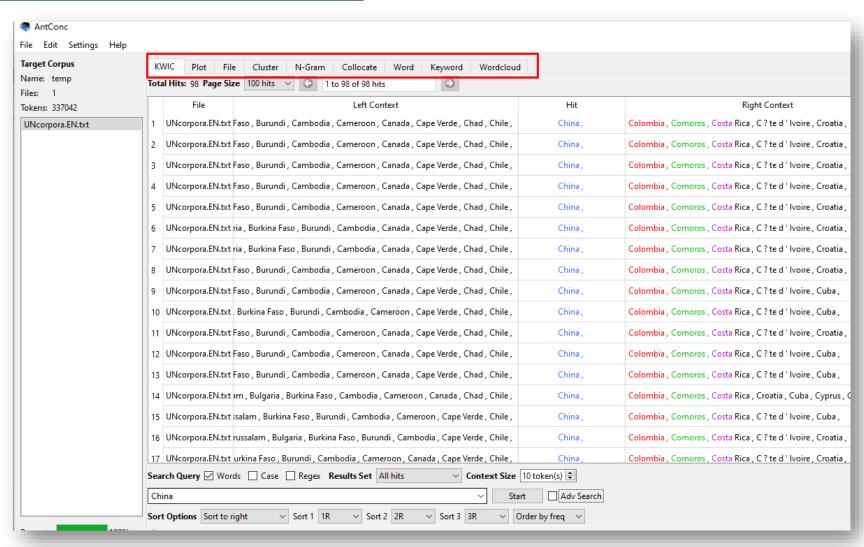
4. 本地单语、双语语料库工具



单语语料库工具: AntConc

http://www.laurenceanthony.net/software/antconc/

UTF-8编码: 通用、国际化



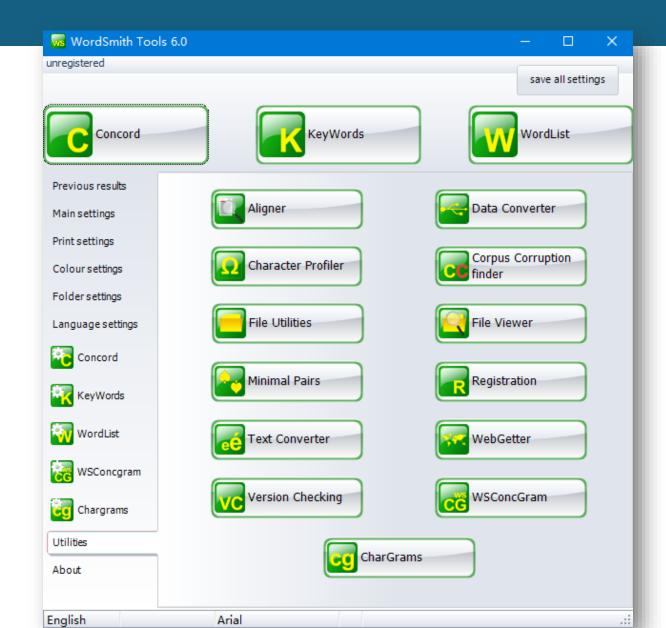
单语语料库工具: WordSmith

优点:

- **1.功能丰富**:词汇检索、词频统计、共现分析等
- **2.界面友好**: 软件界面布局合理, 操作相对简便
- **3.支持大文件处理**:可以处理大型语料库。

缺点:

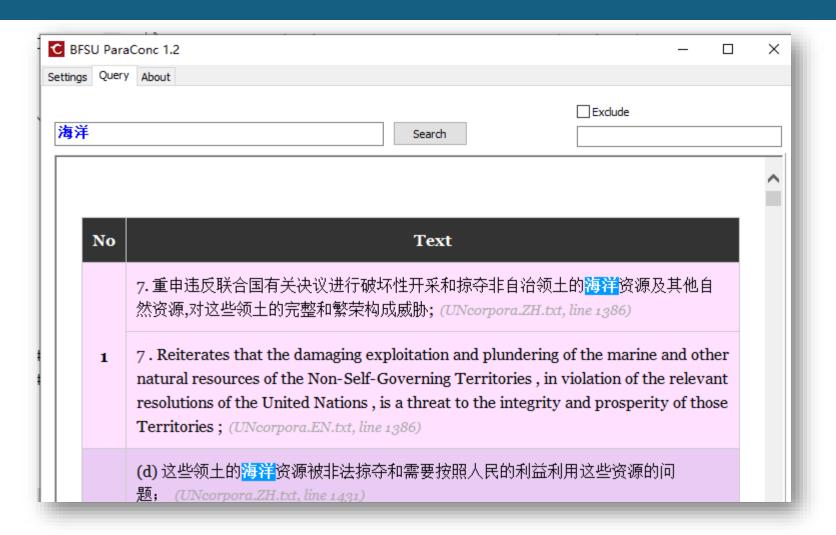
- 1.商业软件:需要付费使用
- 2.功能细节不尽完善:虽然功能丰富,但细节方面可能不如其他免费的AntConc。



双语语料库工具: ParaConc

ANSI编码: 需转换、有局 限性

但是可以创建双语语料库





5. 中文分词、 英文词性标注工具

中文分词的目的

将自然语言分解为最小的单元(词语),为以下活动做准备:

生成词典

多语语料对齐

常见中文分词的工具(Python包)

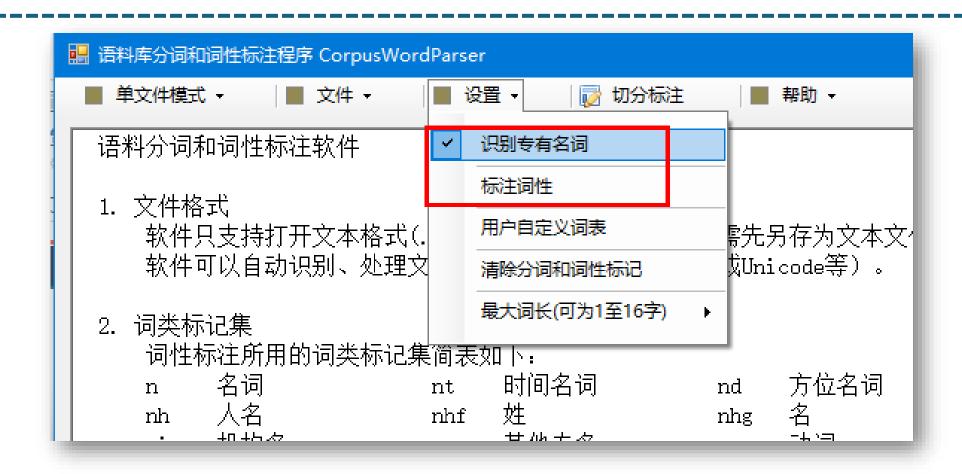
Jieba: 中文分词+词性标注 (热门工具)

SnowNLP: 中文分词+词性 标注、情感分析、文本分类 等

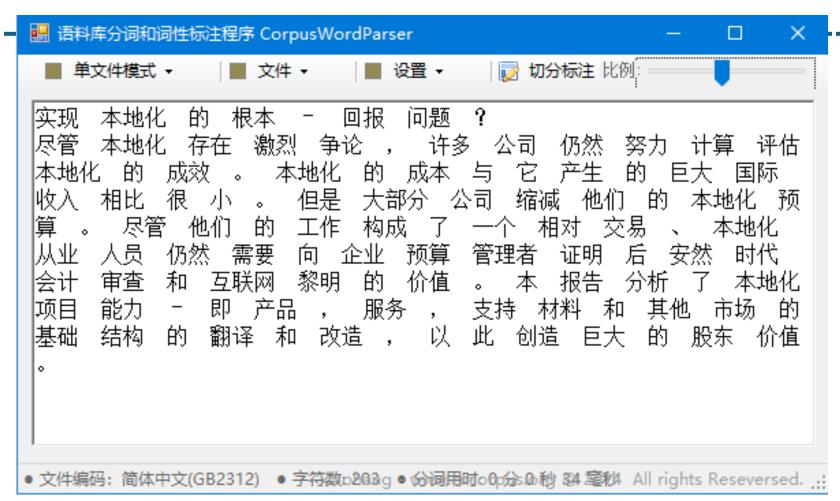
PkuSeg: 多领域分词,个性化的预训练模型。

THULAC: 中文分词+词性标注。清华大学自然语言处理与社会人文计算实验室。

中文自动分词工具: CorpusWordParser.exe



中文自动分词工具: CorpusWordParser.exe



中文语料库在线分词与标注工具,分词后的文件为ANSI编码的TXT文件

英文词性标注工具: Tree Tagger

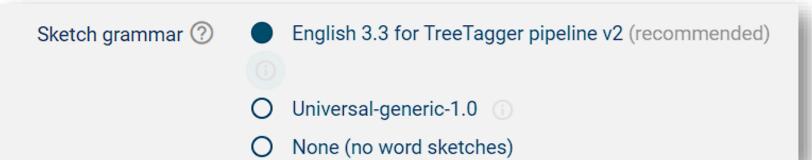


Getting VVG to TO the DT Bottom NP of IN Localization NP - NN What WP Is VBZ the DT Payback NN ? SENT Despite IN compelling JJ arguments NNS for IN localization NN , , many JJ firms NNS are VBP still RB struggling VVG with IN figuring VVG out RP how WRB to TO justify VV the DT effort NN SENT The DT cost NN of IN localization NN happens VVZ to TO be VB very RB small JJ compared VVN to TO the DT big JJ international JJ revenue NN it PP can MD help VV generate VV SENT But_CC most_JJS firms_NNS shortchange_VV their PP\$ localization NN budgets NNS . SENT Even_RB though_IN their_PP\$ work_NN constitutes_VVZ a_DT relative_JJ bargain_NN ,_, localization_NN practitioners NNS still RB have VHP to TO prove VV their PP\$ value NN to TO corporate JJ budgeters NNS mindful JJ of IN post-Enron NN accounting NN scrutiny NN and CC the DT morning-after JJ internet NN malaise NN . SENT

ANSI编码的 TXT文件

英文词性标注工具: Tree Tagger

其他调用 tree tagger 的方式



The following steps are necessary to install the TreeTagger (see below for the <u>Windows version</u>). Download the files by right-clicking on the link. Then select "save file as". All files should be stored in the same directory.

- 1. Download the tagger package for your system (<u>PC-Linux</u>, <u>Mac OS-X (Intel</u>), <u>Mac OS-X (M1</u>), <u>ARM64</u>, <u>ARMHF</u>, <u>ARM-Android</u>, <u>PPC64le-Linux</u>). If you have problems with your Linux kernel version, download this <u>older Linux version</u> and rename it to tree-tagger-linux-3.2.5.tar.gz.
- 2. Download the tagging scripts into the same directory.
- 3. Download the installation script install-tagger.sh.
- 4. Download the parameter files for the languages you want to process.
- 5. Open a terminal window and run the installation script in the directory where you have downloaded the files: sh install-tagger.sh
- 6. Make a test, e.g.

 echo 'Hello world!' | cmd/tree-tagger-english

 or

 echo 'Das ist ein Test.' | cmd/tagger-chunker-german
- 7. You also might want to have a look at my new part-of-speech tagger RNNTagger.

6. 语料库对齐工具



语料对齐

■ 什么是对齐 (Alignment) ?

通过比较和关联源语言文档和目标语言文档创建双语平行语料库的过程。

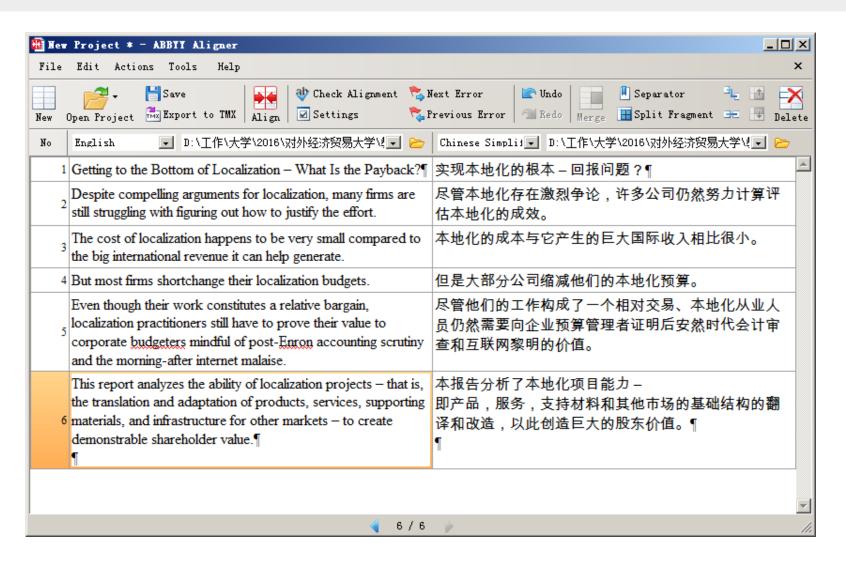
■ 对齐方法

- 使用CAT工具的对齐功能
- 使用特定工具 (Abbyy Aligner, TMXmall的在线对齐)

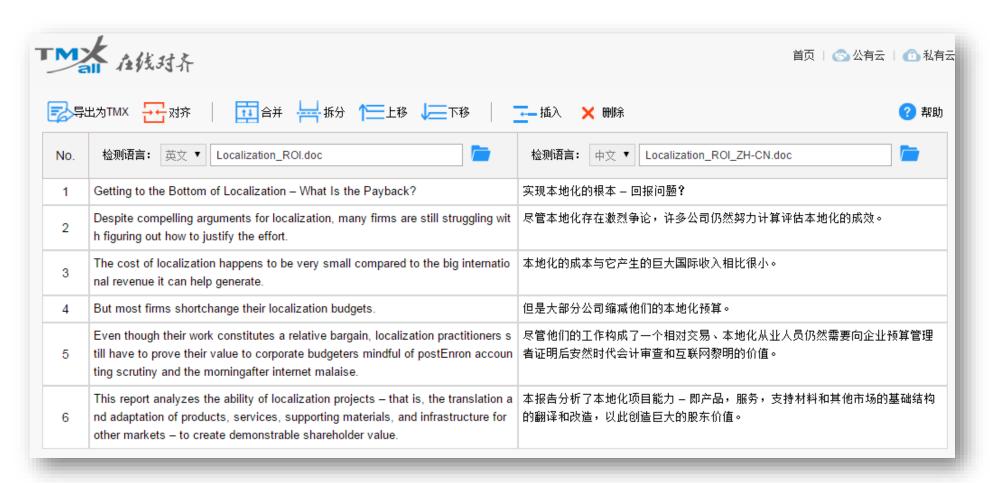
■ 对齐注意事项

- 原文和译文段落数相同
- 手工处理断句问题

使用Abbyy Aligner对齐



使用TMXMall在线对齐



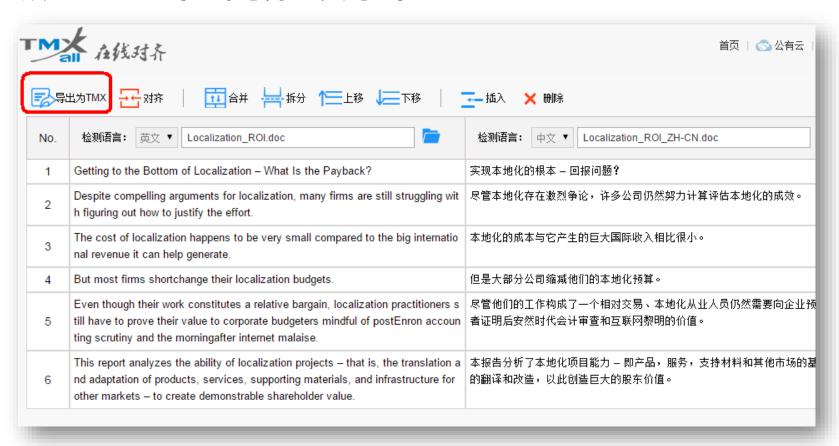
http://www.tmxmall.com



• 7. 翻译记忆库的创建和导出

翻译记忆库的导出

■ 从TMXMall在线对齐工具导出



ChatGPT创建翻译记忆库

将英文和中文以句子为单位对齐成标准TMX格式的翻译记忆库文件,下面第一段是英文原文,第二段是中文译文:

Getting to the Bottom of Localization – What Is the Payback?

Despite compelling arguments for localization, many firms are still struggling with figuring out how to justify the effort. The cost of localization happens to be very small compared to the big international revenue it can help generate. But most firms shortchange their localization budgets. Even though their work constitutes a relative bargain, localization practitioners still have to prove their value to corporate budgeters mindful of post-Enron accounting scrutiny and the morning-after internet malaise. This report analyzes the ability of localization projects — that is, the translation and adaptation of products, services, supporting materials, and infrastructure for other markets — to create demonstrable shareholder value.

实现本地化的根本 - 回报问题?

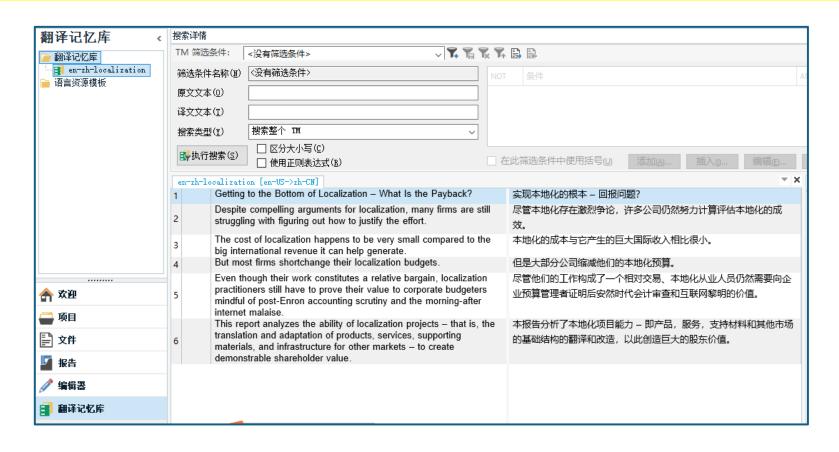
尽管本地化存在激烈争论,许多公司仍然努力计算评估本地化的成效。本地化的成本与它产生的巨大国际收入相比很小。但是大部分公司缩减他们的本地化预算。尽管他们的工作构成了一个相对交易、本地化从业人员仍然需要向企业预算管理者证明后安然时代会计审查和互联网黎明的价值。本报告分析了本地化项目能力 – 即产品,服务,支持材料和其他市场的基础结构的翻译和改造,以此创造巨大的股东价值。

ChatGPT创建翻译记忆库

<tmx version="1.4"> <header creationtool="Translation Memory eXchange" creationtoolversion="1.0"</pre> segtype="sentence" o-tmf="ABCTransMem" adminlang="en-US" srclang="en-US" datatype="PlainText" /> <body> <tu> <tuv xml:lang="en-US"><seg>Getting to the Bottom of Localization – What Is the Payback?</seg></tuv> <tuv xml:lang="zh-CN"><seg>实现本地化的根本 – 回报问题? </seg></tuv> </tu> <tu> <tu> xml:lang="en-US"><seg>Despite compelling arguments for localization, many firms are still struggling with figuring out how to justify the effort.</seg></tuv><tuv xml:lang="zh-CN"><seg>尽管本地化存在激烈争论,许多公司仍然努力计算评 估本地化的成效。</seg></tu>></tu>>tu><tu>xml:lang="en-US"><seg>The cost of localization happens to be very small compared to the big international revenue it can help generate.</seg></tuv> <tuv xml:lang="zh-CN"><seg>本地化的成本与它产生的巨大国际收入相比很小。</seg></tu>></tu> <tu> <tu> <tu> xml:lang="en-US"><seg>But most firms shortchange their localization budgets.</seg></tuv> <tuv xml:lang="zh-CN"><seg>但是大 部分公司缩减他们的本地化预算。</seg></tu>>/tu> <tu> <tu> xml:lang="en-US"><seg>Even though their work constitutes a relative bargain, localization practitioners still have to prove their value to corporate budgeters mindful of post-Enron accounting scrutiny and the morning-after internet malaise.</seg></tuv> <tuv xml:lang="zh-CN"><seg>尽管他们的工作构成了一个相对交易、本地化从业人员仍然需要向企业预算管理者证明后安然时代 会计审查和互联网黎明的价值。</seg></tu>></tu>>tu> <tu> xml:lang="en-US"><seg>This report analyzes the ability of localization projects – that is, the translation and adaptation of products, services, supporting materials, and infrastructure for other markets – to create demonstrable shareholder value.</seg></tuv> <tuv xml:lang="zh-CN"><seg>本报告分析了本地化项目能力-即产品,服务,支持材料和其他市场的基础结构的翻译和改造,以 此创造巨大的股东价值。</seg></tu>>/tu></body></tmx>

ChatGPT创建翻译记忆库

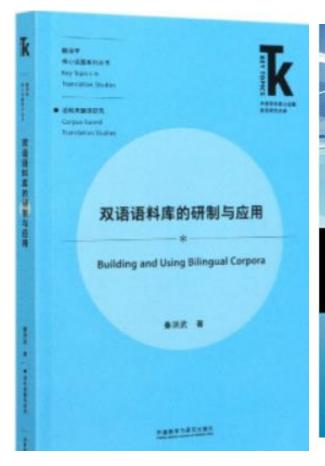
将文本内容保存到TXT文件,修改文件扩展名为TMX,在Trados Studio新建翻译记忆库文件,导入TMX文件。

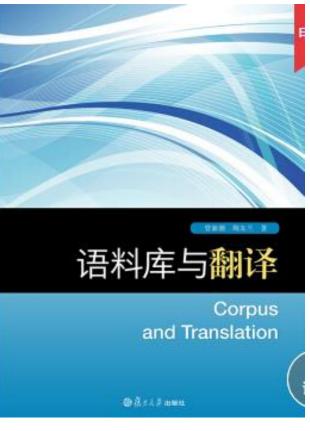


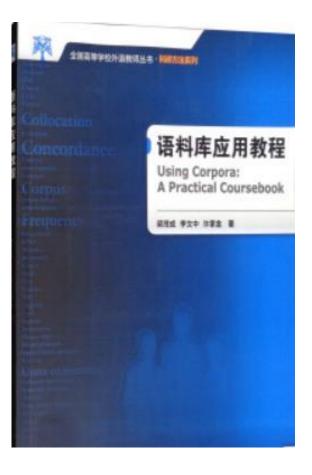
课外阅读材料

Excel文件与TMX文件的相互转换方法

来源:"本地化世界"微信公众号

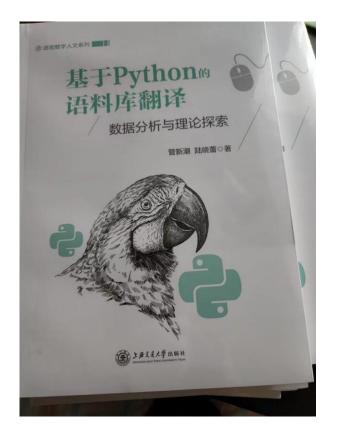


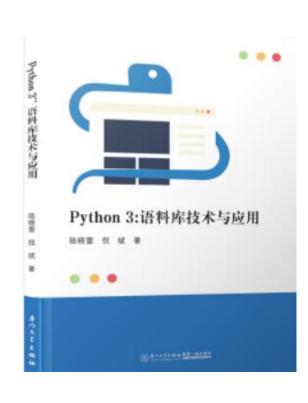


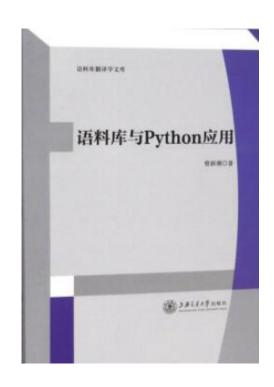


参考材料

Python语料库编程





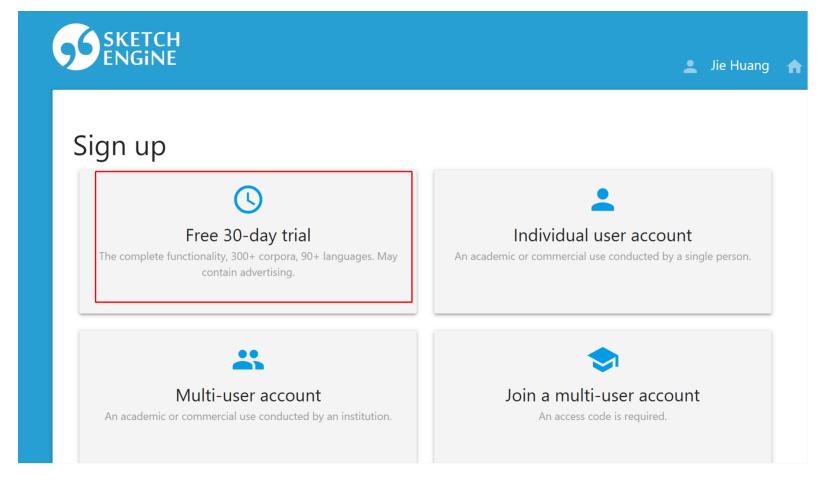


番外: SketchEngine · 注册和使用指南

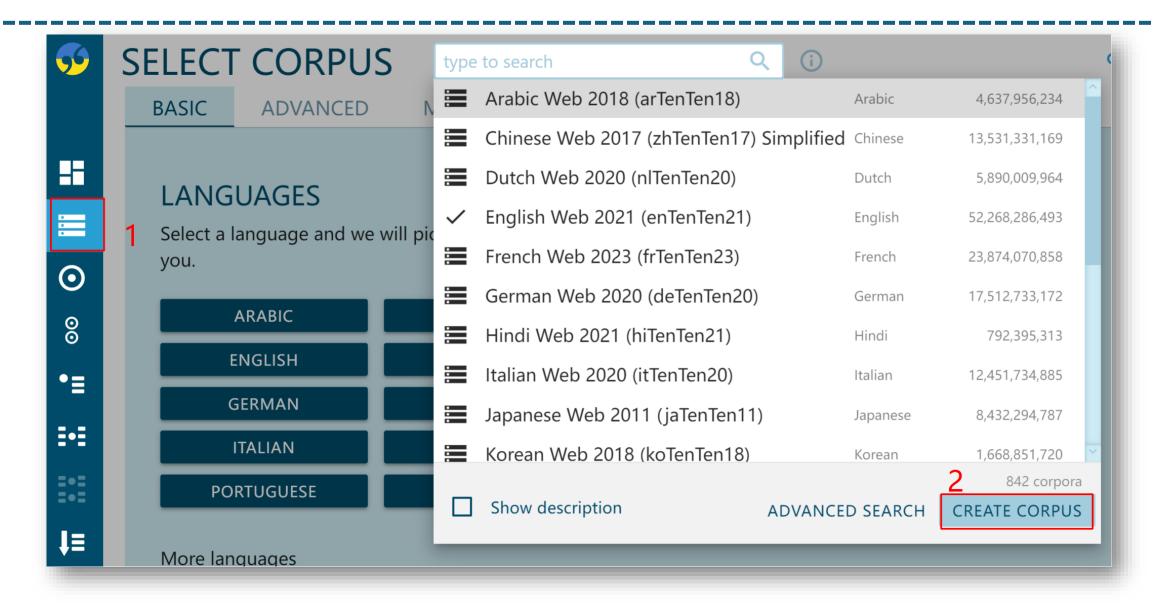


注册 30-day free trial

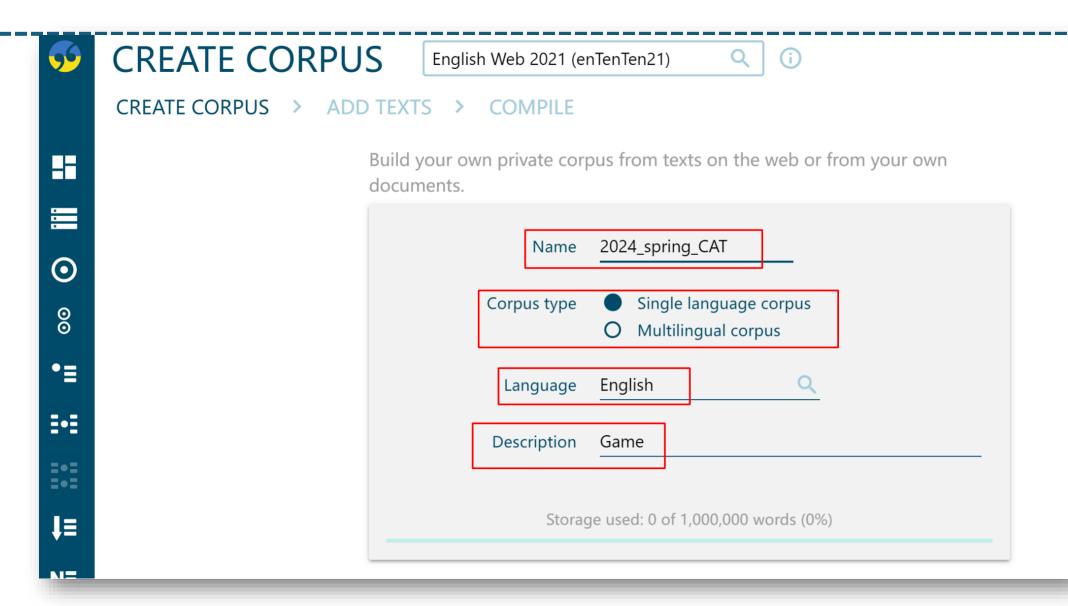
https://auth.sketchengine.eu/#register



创建语料库—create corpus



设置语料库属性



选择语料来源——网页/本地文件

CORPUS: 2024_spring_CAT (English)

CREATE CORPUS > ADD TEXTS > COMPILE



Find texts on the web

Automatically find and download relevant texts



I have my own texts

Upload your own files (.txt, .pdf,...) or paste text

网页搜索关键词:至少3个

Input type



Web search ②



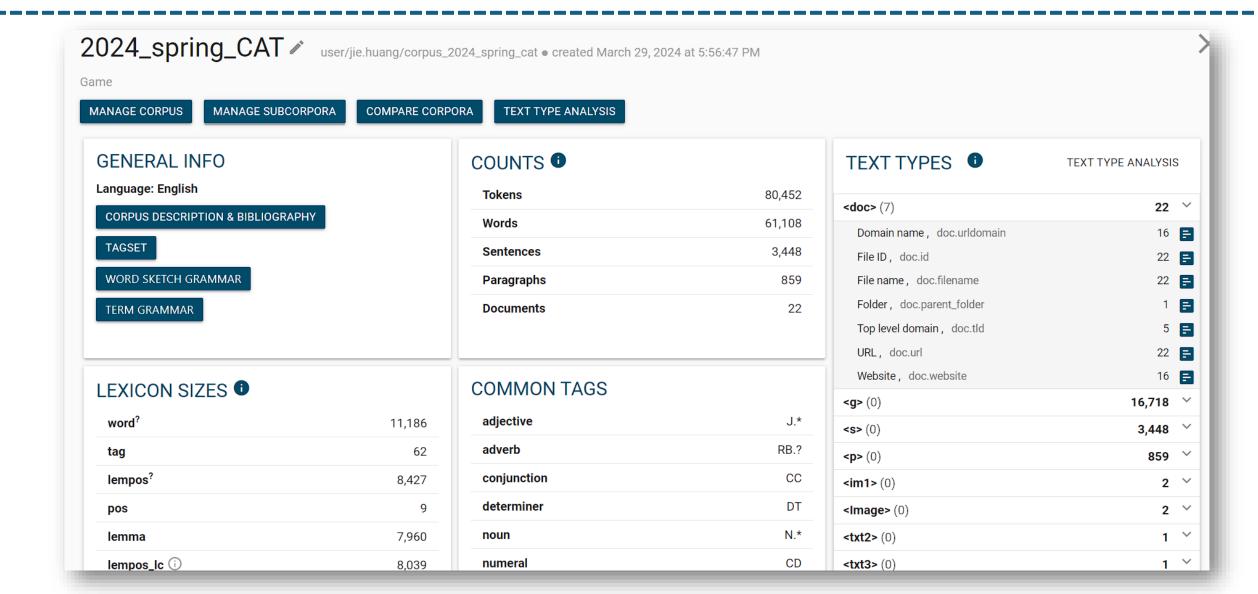
Input some words and phrases that define the topic of the new corpus. Words will be randomly selected and groups of 3 will be sent to the Bing search engine. The web pages that Bing returns will be downloaded and processed into a corpus. Input between 3 and 20 words or phrases.

Hit ENTER after each one.

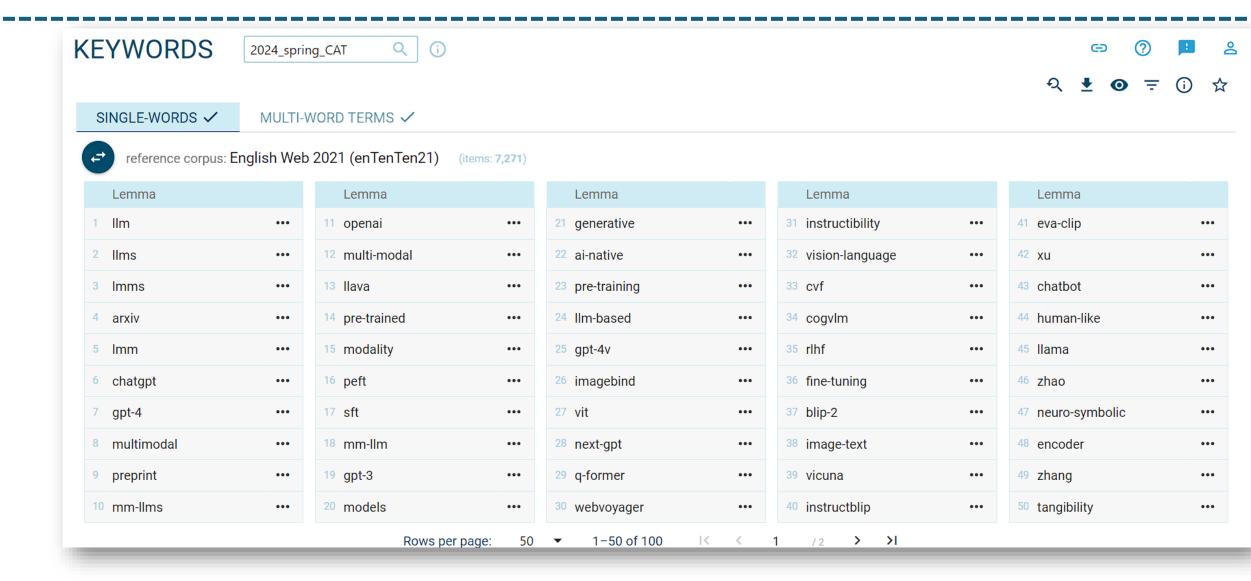
搜索——GO

Select web pages to download The selected web pages will be downloaded. Deselect those that should be skipped. A page may be removed after the downloaded. match your denylist settings, allowlist settings or size restrictions. Check the settings now (Your current selection will be lost.) Filter COLLAPSE ALL **DESELECT VISIBLE EXPAND ALL** SELECT VISIBLE game terminology ● game localization ● video game (23/23 selected) ^ **blog.andovar.com**/games-translation-ultimate-guide **ehlion.com**/magazine/gaming-terminology/ en.wikipedia.org/wiki/Glossary of video game terms **gengo.com**/industry-translation/video-game-translation-services/ **j-entranslations.com**/what-skills-do-i-need-to-be-a-game-translator-part-1-translation-skills/ link.springer.com/chapter/10.1007/978-3-030-42105-2 15 link.springer.com/chapter/10.1007/978-3-030-88292-1 3 **multiplatform.com**/news/demystifying-game-terminology-in-video-game-localization-ptw-s-experience/ research-information.bris.ac.uk/en/publications/terminology-management-in-game-localization **smartcat.com**/blog/game-localization/ ✓ academia.edu/6639017/Challenges in video game localization An integrated perspective <a>I

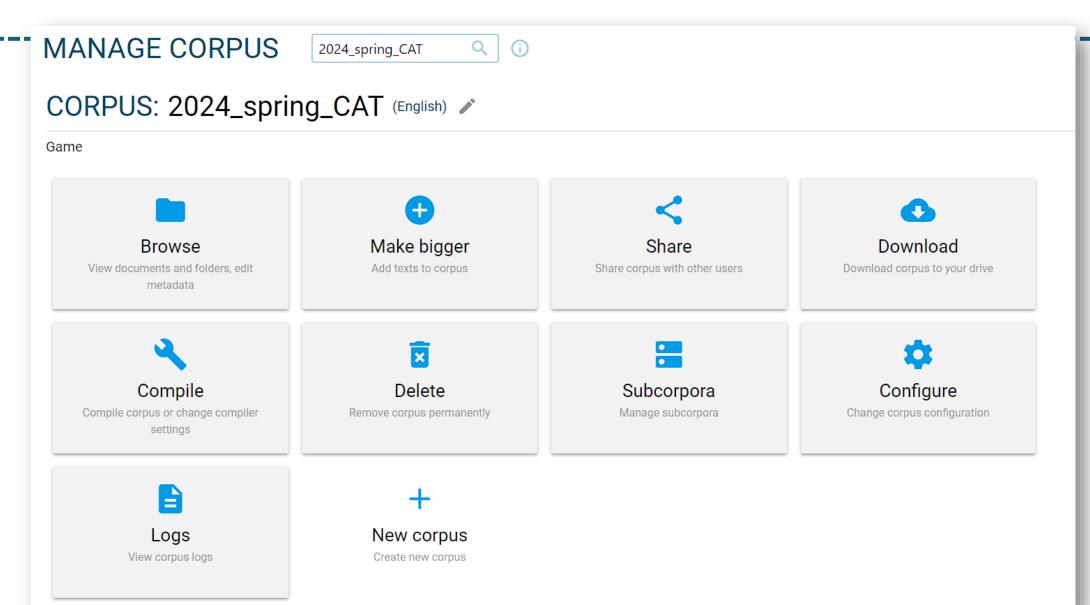
语料库创建成功



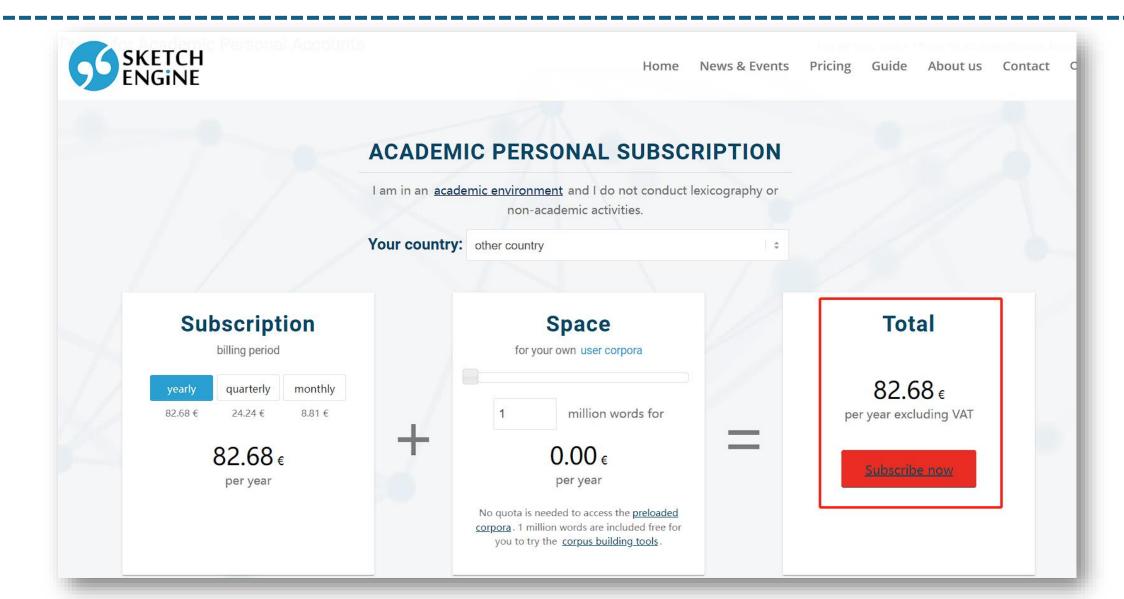
术语列表



语料库管理



Pricing



Sketch engine vs. English-corpora



english-corpora.org/compare-sketchEngine.asp

English-Corpora.org 🤼



companies, including Google, Amazon, Microsoft, IBM, Samsung; Merriam-

₩☆

corpora guides related resources users my account upgrade help

English-Corpora.org and SketchEngine are probably the two largest sites for online corpora. We believe that both sites provide valuable resources for linguists, lexicographers, and language learners and teachers.

The following is a comparison of the two sites, for those who are already family with Sketch Engine, but are new to English-Corpora.org. Admittedly (because this list is at English-Corpora.org), it is probably biased towards English-Corpora.org, and we invite you to look more in depth at what Sketch Engine has to offer as well. Finally, if there is incomplete / incorrect information below, please let us know.

Feature	Sketch Engine	English-Corpora.org
Corpora	- Extremely wide (90+) range of languages, and hundreds of corpora - For English, very large web-based corpora, as well as many other specialized corpora	 Mostly English, as well as some for Spanish and Portuguese For English, perhaps the best suite of corpora for looking at variation: genre-based, historical, and dialectal Largest corpora are iWeb (14 billion words) and NOW (14.6 billion words and growing by ~250 million words each month)
Users / research	- Linguistics and lexicographers, teachers and learners, etc (For those with information on Sketch Engine, please send us more detailed / verifiable information on number of users, researchers, universities with licenses, number of publications, etc)	users - ~300 universities have academic (group) licenses, as well as large government-funded licenses

End