# TmxSQL

## Managing Translation Memories with SQL

Steve Braich

https://github.com/stevebpdx/TmxSql

# Managing Translation Memories

**Background:**

- Translation memories (TMs) are saved translations that can be leveraged later for future translations.

- They are usually saved in an XML file format known as TMX (Translation Memory Exchange)

- There are virtually no open source TM desktop management tools that have continued support.

- Translation memories are a very valuable asset for training machine translation models.

# Why the need for TM Management?

- Translators use saved translations so they can reduce future translation work (reference, fuzzy matching)

- Localization Engineers use TMs when they develop Enterprise TMS software

- Machine Translation Engineers use TMs when they train machine translation models

**Summary:**

Translation memories are a valuable asset for translators and organizations. Effective management of TMs is critical to leveraging their potential and reducing cost with expanding an organization's global footprint.

# What is a TMX

- TMX stands for Translation Memory eXchange, an XML file format

```
 5  />
 6  <tu srclang="EN-US" tuid="61">
 7      <tuv xml:lang="EN-US">
 8          <seg>I don't speak French very well.</seg>
 9      </tuv>
10      <tuv xml:lang="FR" changedate="20140210T125940Z" changeid="service_acme" creationdat
11          <prop type="x-Context">Travel</prop>
12          <prop type="x-Source">http://www.fodors.com/</prop>
13          <seg>Je ne parle pas très bien français.</seg>
14      </tuv>
15  </tu>
16  <tu srclang="EN-US" tuid="79">
17      <tuv xml:lang="EN-US">
18          <seg>Do you speak English?</seg>
19      </tuv>
20      <tuv xml:lang="FR" changedate="20140210T071608Z" changeid="service_acme" creationdat
21          <prop type="x-Context">Travel</prop>
22          <prop type="x-Source">http://www.fodors.com/</prop>
23          <seg>Parlez-vous anglais?</seg>
24      </tuv>
```

- For more information on the TMX file format:
  https://en.wikipedia.org/wiki/Translation_Memory_eXchange
  https://www.gala-global.org/tmx-14b

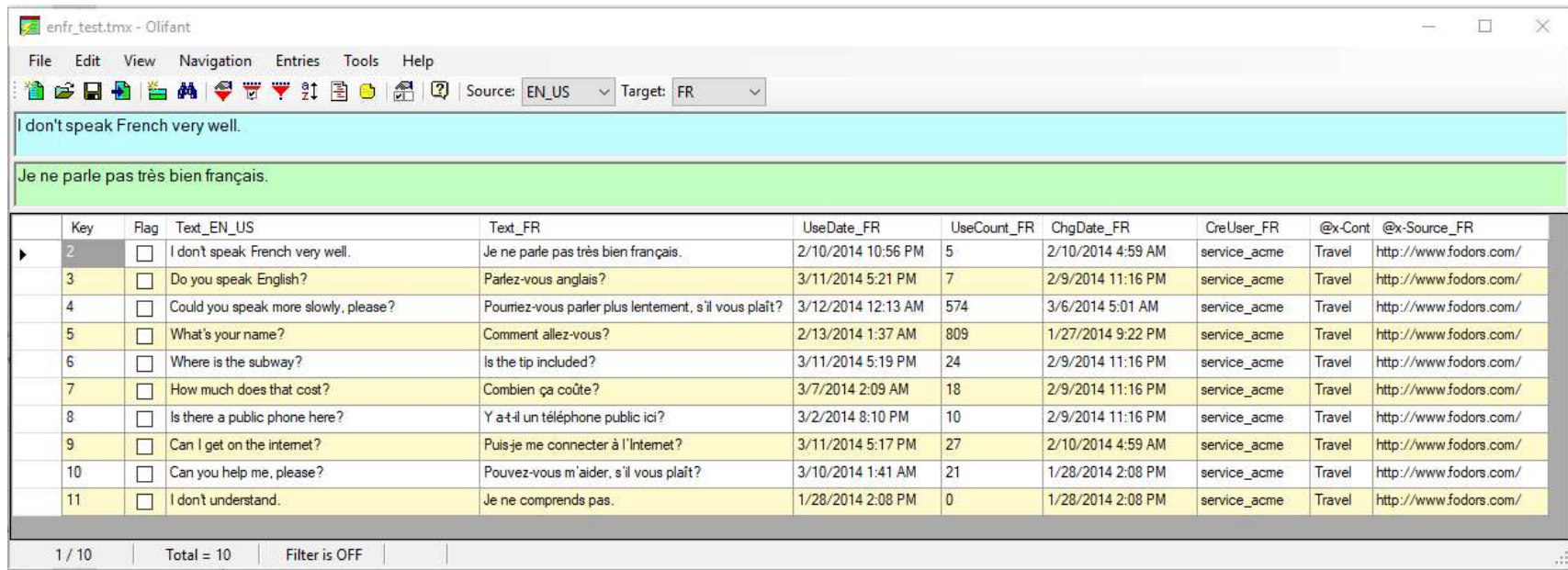# Types of Translation Memory Management

There are two basic types of TM Management:

| Translation Management System (TMS) | Enterprise end-to-end solution. Very few free open source offerings. Often include desktop CAT tools. |
| --- | --- |
| Computer Assisted Translation Tools (CAT) | Various types of tools that deal with a specific stages of the translation process. Usually a desktop tool that will allow translators without direct access to a TMS to do work offline and submit their translations. |

# Open Source Desktop TM Management Tools

## Okapi Framework - Olifant:

- http://okapi.sourceforge.net/Release/Olifant/Help/

- https://bitbucket.org/okapiframework/olifant/overview

- https://github.com/OkapiFramework/okapi-olifant

# Olifant Limitations:

- Hasn't been maintained in four years
- Custom properties (meta-data) are not strongly typed which causes problems with sorting and filtering
- Doesn't allow you to manage a repository of TMs.  Just one single TMX file at a time.
- Limited import/export features
- Limited features for advanced users like localization engineers.

# Other open source or free solutions:

- Google Translator's Toolkit
  https://translate.google.com/toolkit

- Heartsome
  https://github.com/heartsome/tmxeditor8

# Latest Olifant Development

- Allows management of multiple TMX files in a centralized repository
- In 'ALPHA' phase
- Last update was four years ago
- https://bitbucket.org/okapiframework/olifant

# Introducing: TmxSql

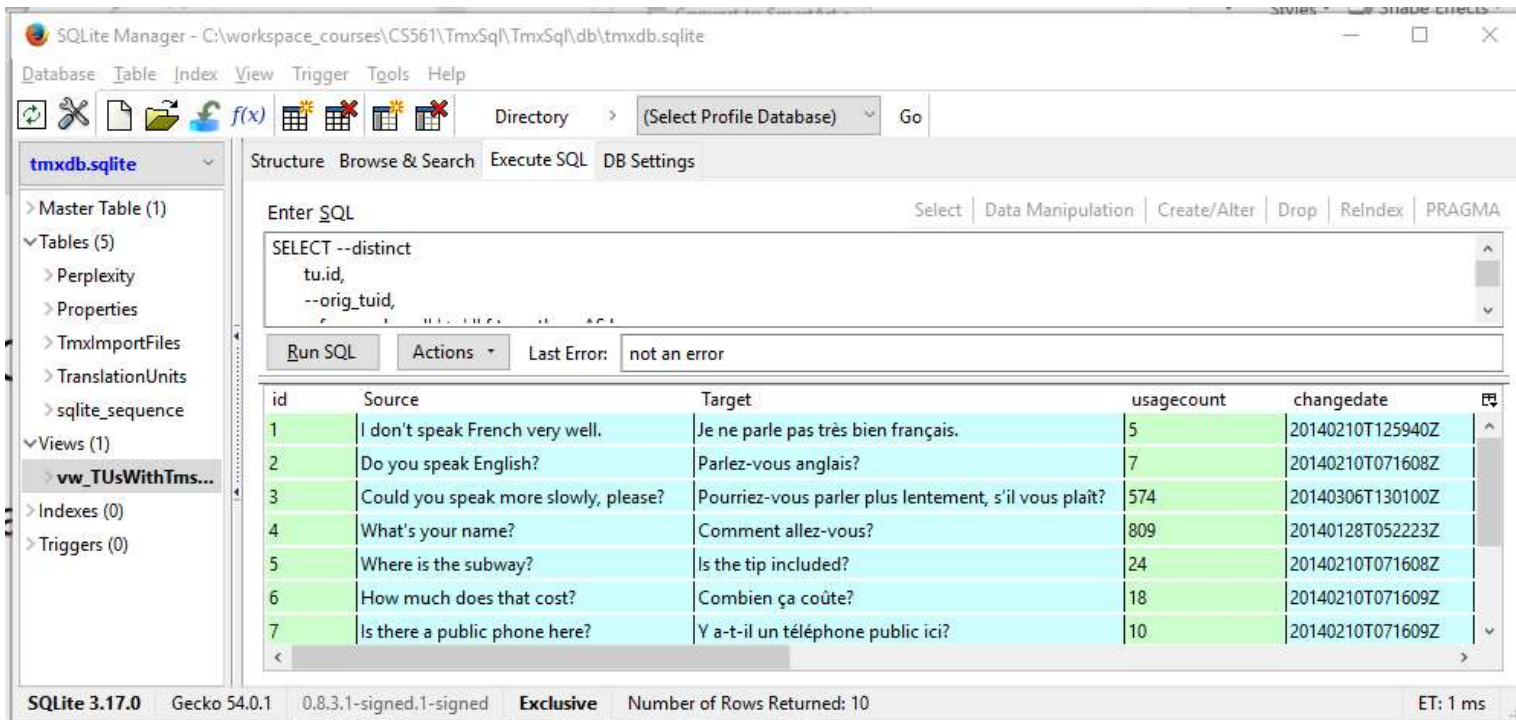TmxSql uses a python script to import TMX files into a single Sqlite database where engineers can easily manage TMs using SQL.

# TmxSql Features:

- You can centralize your TMs in a single repository
- Allows localization engineers to use SQL to manage TMs
- Imports custom properties from TMX files
- Other than installing a Sqlite and a python script, no other software is required

# TmxSql Limitations:

- Only works with bilingual corpora, not multilingual corpora
- Not suitable for translators or users that are not localization engineers
- Issues dealing with ambiguities related to custom properties in TMX

# Self Assessment / Lessons Learned:

1. Defining the project goal was/is difficult
2. Spent a lot of time trying to join the Okapi Framework team.
3. Tried to get a too many birds with one stone:
   - A potential client wanted to pay me to work on something similar
   - I tried to make that work open source
   - Tried to make it part of the Okapi Framework
   - And… to use the Okapi Framework project as my class project
4. I spent less time adding new features to my project and more time on point #3.
5. Failed to drum interest by potential contributors and users.

# Interested in Contributing?

TmxSql:

Steve Braich

stevebpdx@gmail.com

https://github.com/stevebpdx/TmxSql

Okapi Framework:

https://bitbucket.org/okapiframework