# Compiling an Online Dictionary based on Field Data:
## The Case of Kelabit Utilizing TEI/XML, XSLT and ChatGPT

**Yasuka Fukaya**

Kyushu University/JSPS (Japan Society for the Promotion of Science)

(yasuka.fukaya@gmail.com)

# Abstract

**Introduce Kelabit and ODK project**

*Online Dictionary of Kelabit

**How NLP assists in this project**

ChatGPT helped the author create...

XSLT (to transform XML to HTML)

JavaScript (to add research function)
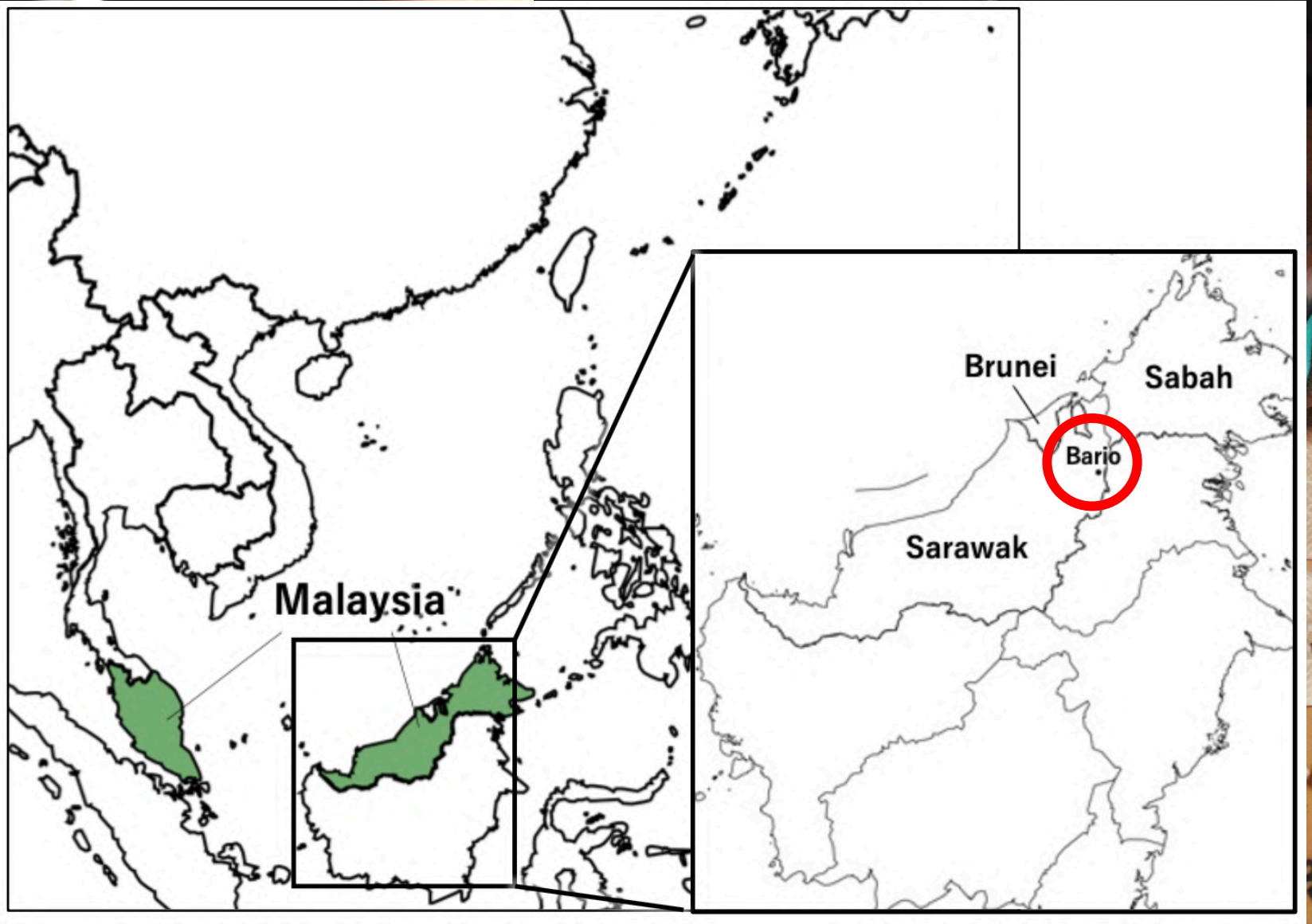
# Structure of Presentation

Part 1: Background

Part 2: How NLP were used in the ODK project
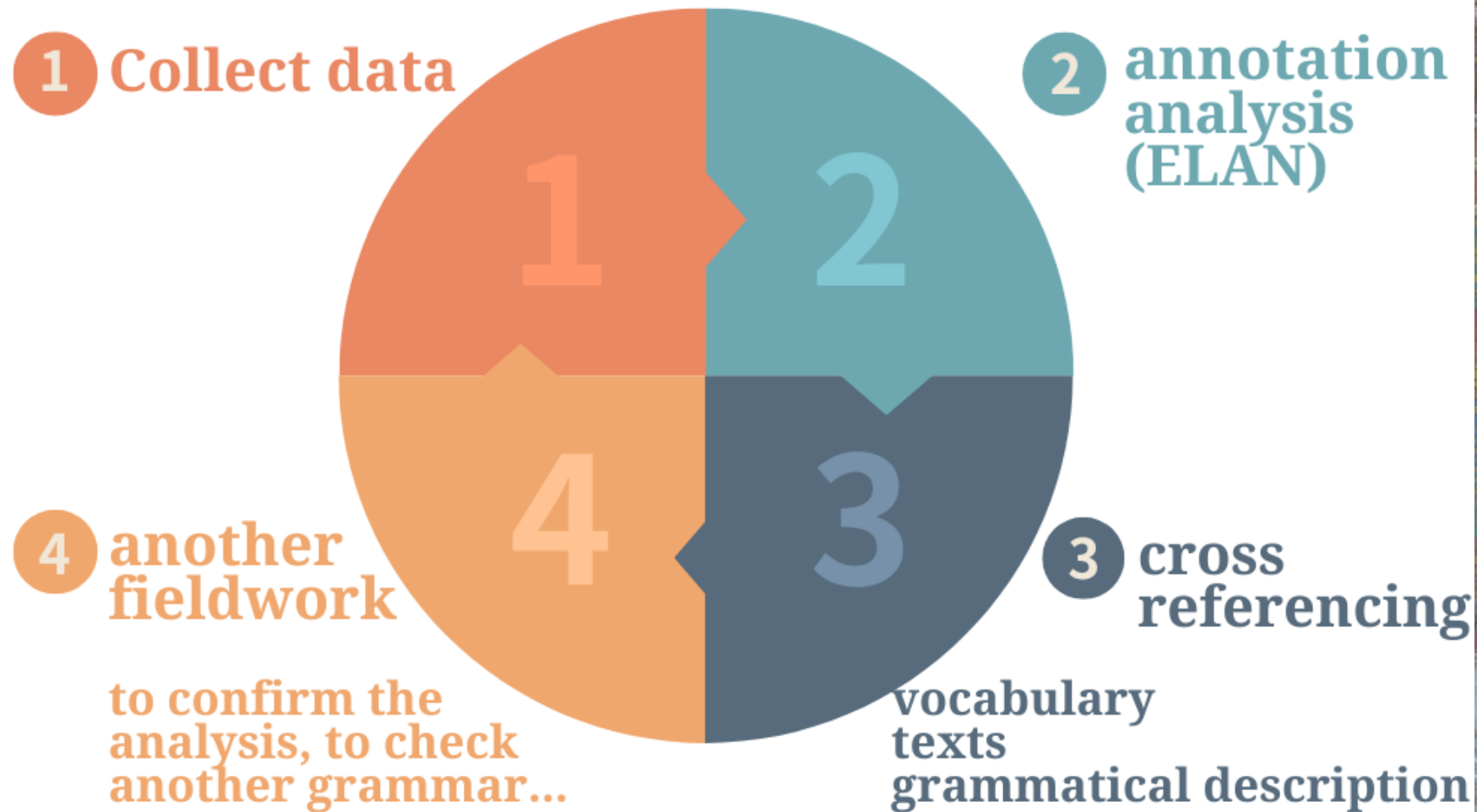
# Part 1: Background
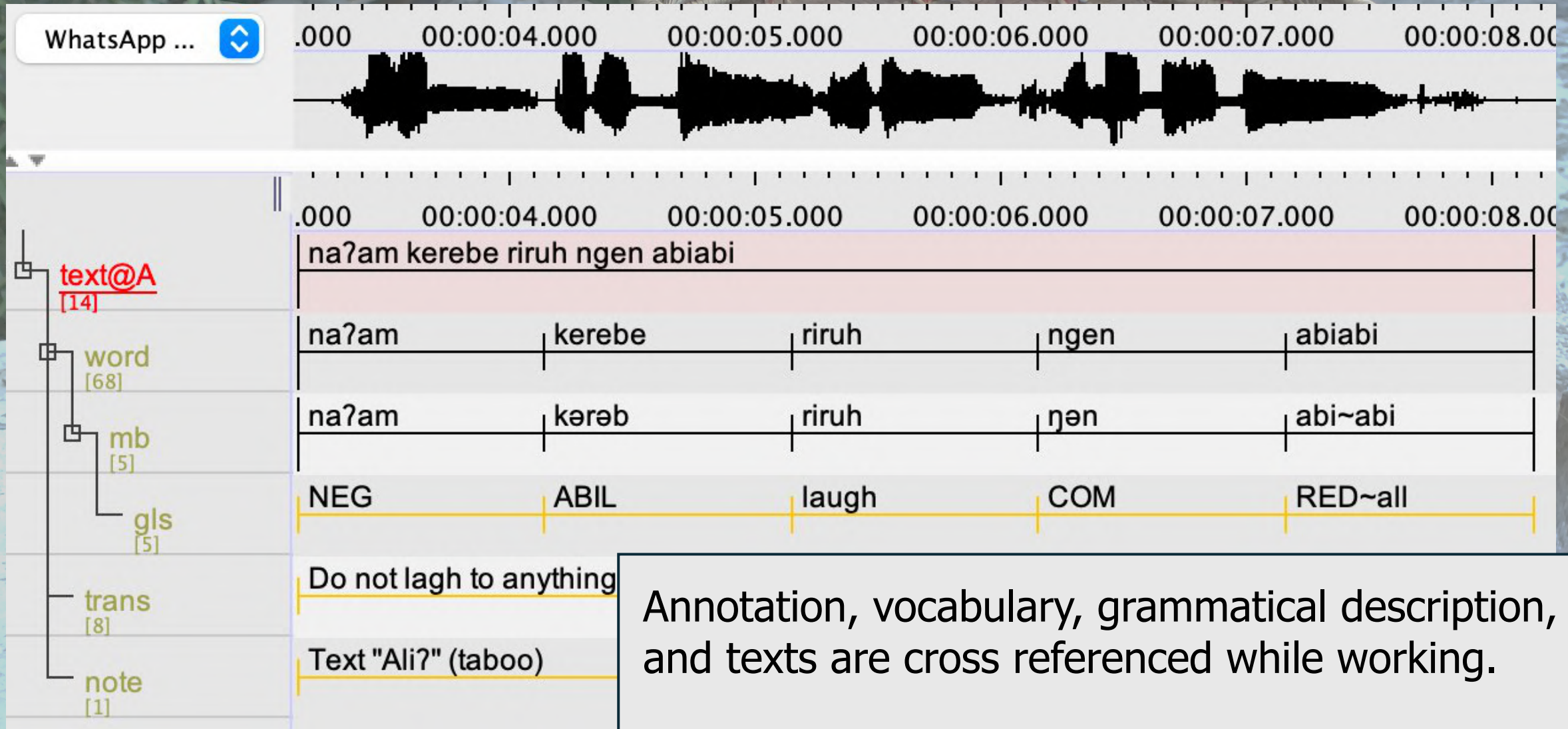
**Bario, Sarawak, Malaysia**

Kelabit

Brunei
Sabah
Bario
Sarawak
Malaysia

[kaɾuh kəlabit]
kaɾuh  Kelabit

# Annotation work using ELAN



Annotation, vocabulary, grammatical description, and texts are cross referenced while working.

We have lexical data for a dictionary…

**Needs of speakers' community**

- **language inheritance**

  dictionary
  children's book
  textbook for grown-ups
  ...

grammar description

making a dictionary

# Part 2:
# How NLP were used in the ODK project?

# The process of creating ODK

- **The process of making lexical data for the dictionary**

Raw data → transcription (ELAN) **Done while grammar writing**

→ **transformation to TEI/XML (Oxygen)**

- **The process of getting the output in different formats**

XSLT(ChatGPT) → XML to HTML (Oxygen)

→ search function in JavaScript (ChatGPT)

# ELAN --> TEI/XML (OxygenXMLEditor)



WhatsApp Ptt 2023-...9 at 14.07.02_FY.eaf

Kelabit_lexicon.xml

| lexicon | entry | sense | grammatical-category |

```xml
1  <?xml version="1.0" encoding="UTF-8" standal
2  <lexicon xmlns:xsi="http://www.w3.org/2001/X
3      <header>
4          <name>Kelabit_lexicon</name>
5          <language>Kelabit</language>
6      </header>
7      <entry id="e_df7e5c06-facf-4c0b-a552-f5d
8          <lexical-unit>na?am</lexical-unit>
9          <morph-type>stem</morph-type>
10         <sense id="s_06c06f78-2104-433c-87e1
11             <grammatical-category>adv</gramm
12             <gloss>NEG</gloss>
13         </sense>
14     </entry>
```

**&lt;entry&gt;**
   **&lt;lexical-unit&gt;...**
   **&lt;morph-type&gt;...**
   **&lt;sense&gt;**
       **&lt;grammatical-category&gt;...**
       **&lt;gloss&gt;...**
   **&lt;/sense&gt;**
**&lt;/entry&gt;**

**NOT suitable for ODK**

# TEI/XML for ODK (TEI Guideline Ch. 10 Dictionaries)

https://tei-c.org/release/doc/tei-p5-doc/en/html/DI.html

```
70  <entry n="1" xml:id="AAlistentry1" class="entry">
71      <form>
72          <orth>uluh</orth>
73          <pron notation="ipa">uluh</pron>
74      </form>
75      <gramGrp>
76          <pos>n</pos>
77      </gramGrp>
78      <sense>
79          <sense type="meaning" n="">
80              <def xml:lang="jp">頭</def>
81              <def xml:lang="en">head</def>
82          </sense>
83          <sense>
84              <cit type="example">
85                  <quote id="W1.S1">neh ngudeh uluh mudih?</quote>
86                  <morphs>
87                      <morph id="W1.S1.morph1">neh</morph>
```

**TEI/XML is made by hand with Oxygen XML editor.**
- IPA
- part of speech
- meaning
- example
  - with gloss
- other information

# The process of creating ODK

- **The process of making lexical data for the dictionary**
  Raw data → transcription (ELAN)
  → transformation to TEI/XML (Oxygen)

- **The process of getting the output in different formats**
  XSLT(ChatGPT) → XML to HTML (Oxygen)
  → search function in JavaScript (ChatGPT)

# Creating XSLT with ChatGPT

ChatGPT recommended XSLT over CSS.

```
<?xml version="1.0" encoding="UTF-8"?>
<TEI >...</TEI>
```

**Using this XML code, how can I make XSLT with following features?**

In <entry> elements, <orth> is bolded style and <pron> is blue Times New Roman font. <gramGrp> is font in light blue shading. ...

To create an XSLT stylesheet with the desired features for the provided XML code, you can use the following XSLT template:

xml                                         コードをコピーする

```
<?xml version="1.0" encoding="UTF-8"?>
<xsl:stylesheet version="1.0" xmlns:xsl="http://www.w3.org/1999/XSL,
    xmlns:tei="http://www.tei-c.org/ns/1.0">
```

16

# XML to HTML

# Add a search function made by ChatGPT

```
<div
    style="max-width: 550px; margin: 0px auto; x-column-count: 2; x-column-gap: 40px;">
    <h1>Kelabit Dictionary</h1><br />

    <!-- 検索フォームの追加 -->
    <div class="search-box">
        <form id="search-form">
            <input type="text" id="search-input"
                <input type="submit" value="検索
        </form>
    </div>
    <br />

    <xsl:apply-templates select=".//tei:entry"/>
</div>
<script>
    document.getElementById("search-form").addEv
    event.preventDefault(); <!-- リロードしない --

    var searchKeyword = document.getElementById("search-input").value.toLowerCase(); <!-- 検索キーワードを小文
    var entries = document.querySelectorAll(".entry");

    entries.forEach(function(entry) {
    var heading = entry.querySelector(".orth").textContent.toLowerCase(); <!-- 見出し語<orth>を小文字に-->
    var senses = entry.querySelectorAll(".sense"); <!-- <sense>のすべてを取得?-->

    if (heading.includes(searchKeyword)) {
```

**Add JavaScript to the .xsl file**
　・search function

　・jump to the entry

I learned where to put <script>, and how it works from ChatGPT.
It did not work well, so I read other general websites on JavaScript.

18

**Present version**
  limited words and function

**Future goal**
  Better way (xml, xsl…)
  Photos and sounds

https://yasukafukaya.weebly.com/kelabit-dictionary.html

19

I would like to express my special thanks of gratitude to my language consultants and Kelabit community who supports me.

# Terima kasih mula'-mula'!