

## データ内容について

はじめに

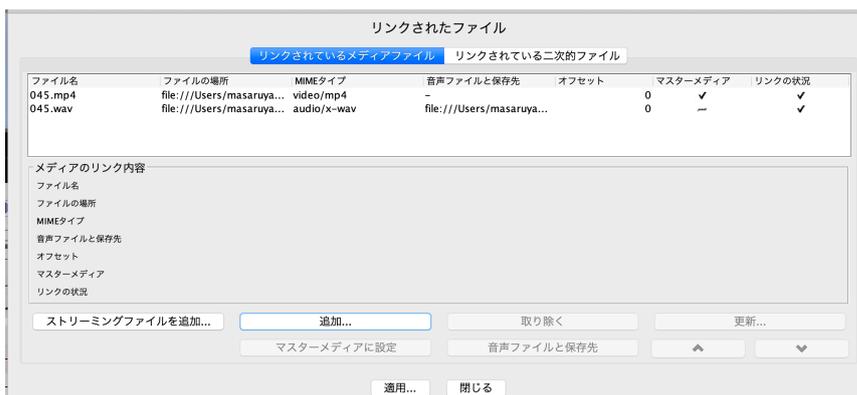
各フォルダに格納されている3種類のファイル（.eaf, .wav, .mp4）を、フリーソフト ELAN を使って開きます。ELAN を事前にインストールしておいてください。

注意：

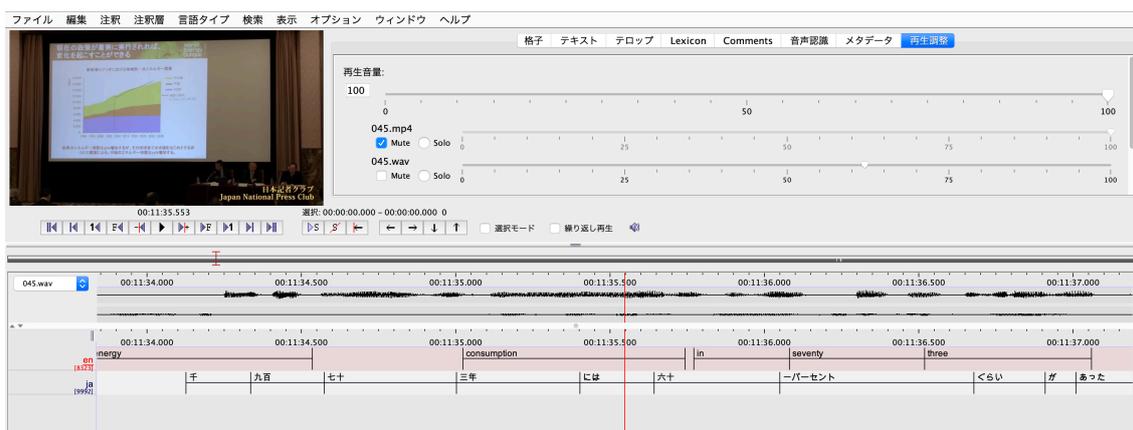
- ✚ フォルダ名は、「S001」や「C001」のような通し番号になっています。「S」は同時通訳、「C」は逐次通訳を意味します。
- ✚ フォルダ内に格納されているファイル名は、フォルダ名とは異なります。例えば、「S001」フォルダ内の3つのファイル名は、「11.eaf」「11.wav」「11.mp4」です。
- ✚ フォルダ名は「収録番号」、ファイル名は「作業番号」（データ作成上の部品番号のようなもの）になります。これらの対応については、別添の「収録データ一覧.xlsx」を御覧ください。

## 設定方法

1. まず、ELAN から.eaf ファイルを開きます。
2. 音声ファイル（.wav）と動画ファイル（.mp4）を設定します。[編集]-[リンクファイル]を選択します。
3. 下のダイアログボックスが表示されたら、[追加]ボタンをクリックして、.eaf ファイルと同じフォルダ内にある、.wav と.mp4 を選択します（2つのファイルそれぞれに対して追加を行う）。追加できたら、[適用]ボタンを押します。



4. 正しく選択されると、下のように動画と音声波形、文字情報が表示されるので、確認してください。



ELAN 上で再生するときには、[再生調整]上の mp4 の音声の[Mute]をチェックしておいたほうが、音がクリアに聞けます（上の場合は、045.mp4 の下の Mute ボタンをチェック）。

## 免責事項

- 本データの書き起こし（トランスクリプト）は、英語と日本語ともに同梱した.wav ファイルを音声認識ソフト（speech to text）にかけて自動生成したものです。使用したシステムは IBM Watson Speech to Text と Speechmatics の 2 種類<sup>1</sup>です。自動生成した後、誤認識によるエラーを、人手によって修正しています。3 人の異なる作業者により 3 度の確認作業を行っていますが、基本的には自動生成テキストをできるだけ活かす方針としているため、2 種類の音声認識システムの記述の不統一や表記揺れ等がございます。研究の目的に合わせ、ご自身で修正や調整をしてご使用ください。
- 音声データ（.wav）はステレオ録音されており、英語と日本語の音声は独立したチャンネルに収録されています。書き起こしのテキストも、ELAN 上で、英語（en）と日本語（ja）のように、異なる注釈層に記録しています。
- 書き起こしの「単位」に関する基本ルールは以下の通りです。「単位」とは「開始」と「終了」の時間が記録されている言葉のかたまりを表します。ELAN 上では、「注釈」の 1 つの区切りが入っている言葉です。

英語の場合は単語が単位になります。例えば、「Thank you very much」という表現は「Thank」「you」「very」「much」のそれぞれの単語を単位として、それぞれの発話時間（開始と終了の時間）が記録されています。

日本語の場合は、わかち書き、もしくはそれより小さい形態素になります。例えば、わかち書きでは、「中国の」「エネルギー需要は」となり、形態素では「中国」「の」「エネルギー」「需要」「は」となります。

※日本語の書き起こしの「単位」は統一されていないため、研究の目的に合わせて調整してください。

<sup>1</sup> <https://www.ibm.com/watson/jp-ja/developercloud/speech-to-text.html>  
<https://www.speechmatics.com/>

- 言い淀み (hesitation) については、「%HESITATION」「 D\_ 」または「 [ah] または [ああ] (角括弧で囲われている) のように記載しています。日本語の場合は「 D\_アア 」のように記していることもあります。しかし、音声認識ツールによっては、言い淀みが削除される場合があるため、すべてが網羅されているわけではありません。また、英語のフィラーや日本語の冗語の一部は、自動音声認識で「不要」と判断され、書き起こされていない場合があります。明らかな false start や言い直しは、[ ]で示していますが、これもすべてを網羅しているわけではありません。こうした現象についての研究を予定されている場合は、ご自身でデータの確認と追加・修正をお願いします。
- 表記の不統一。以下については、表記にバラツキがございます。
  - 漢字・ひらがなの使用。例：「いただいて」「頂いて」
  - 数字の表記。例：アラビア数字と漢数字。例：1995年、千九百九十五年
  - 不要な空白。例：「現役 の」(同一単位内に空白がある)
  - 人名や固有名詞。例：「ふじた」と「藤田」が混在する場合があります
  - どの漢字が使用されるのか不明な場合。例：「あずま」「吾妻」「東」
- その他、正しく書き起こしされていない箇所もあるかもしれませんが、その場合は適宜、修正・調整してください。書き起こしテキストは、ELAN から[ファイル]-[別ファイル形式保存]で任意の形式 (スプレッドシートなど) に書き出すこともできます。

以上

2020年4月13日