

Lipz File Format v1.0

By Erik Hermansen

Overview

The Lipz file format is a simple description of lip animation and subtitles corresponding to an audio file for use in games and other software that performs character animation.

Here's a quick example of a Lipz file to give you the gist:

```
{
  "text" : "Now, boys and girls, I'd like you to put away your readersland get re
ady for a very special treat.",
  "visemes" : "ABBBBCABBFGBBCCAABBCDFFGBAAAABAACDEFFCCBBBABBBCCDB|CCCBBAAGGBBCCDDD
BBAACCBBDDBBCCAAABCCC---",
  "viseme_type" : "toonboom",
  "fps" : 24
}
```

The Lipz file above is associated with one audio file containing the spoken dialogue it describes. The game uses information in the Lipz file to show subtitles and synchronize lip animation of the character as the audio file plays.

The `text` element contains text for the game to display on the screen as subtitles. There is a special "segment delimiter" character (`|`) that shows a point in the text that can be used to segment the text into separate parts for subtitle display. In this case, first a subtitle with "Now, boys and girls, I'd like you to put away your readers" would be shown, followed by a subtitle with "and get ready for a very special treat."

The `visemes` element contains codes that indicate specific visemes (lip frames) to show at different points in time. The segment delimiter (`|`) appears here to mark the point in the audio that corresponds to the same point in the `text` element marked by a segment delimiter.

`viseme_type` and `fps` are used to interpret the codes in the `visemes` element correctly.

Goals of Lipz Format

- Cover use cases for character animation in games including:

- Synchronizing lip frames (visemes) to audio
- Subtitles with localization
- Easy to write parsing code in game software
- Extensible to cover different viseme types, e.g. mouth charts
- Ability to attach events to specific times in dialogue that game can use to trigger character animation

Non-Goals

- *Does not* contain settings and information used by animation authoring or generation tools.
- *Does not* store motion-capture or facial position information.
- *Does not* define animation beyond events triggered by character speech.

Association with Audio File

A Lipz file is associated with one audio file that is expected to contain spoken dialogue. There is no concept of multiple characters or tracks inside of a Lipz file, so typically, the Lipz file will have information for just one character's speech.

The association of the Lipz file to an audio file is managed externally (not specified by Lipz file). An easy convention is to name the Lipz file to match the audio file, e.g. if you have an audio file called `samson-hello.wav` then you could name its corresponding Lipz file `samson-hello.lipz`. And whatever code you have to open up and play audio files, would also look for a Lipz file with a common filename but `.lipz` file extension.

Visemes

Visemes (often mistakenly described as "phonemes") are the visual expression corresponding to a spoken part of dialogue. For example, when you say the word "bowl", your lips will come close together for the "B" sound (a phoneme), and that visual expression is a viseme.









In this version of the Lipz format, three viseme types are defined:

- **Blair** - The classic 9-viseme mouth chart from Disney animator Preston Blair.
- **Toon Boom** - This viseme type is based on the Toon Boom software's popular 7-viseme mouth chart. It's essentially the same as Preston Blair's mouth chart, but separate "L" and "U" visemes aren't included.
- **RMS** - This viseme type represents an amplitude average (RMS) for each frame of animation. If the audio is louder, the mouth will be open wider. If the audio is silent (or below a certain threshold, at least) the mouth will be closed.









The hyphen (`-`) viseme code is reserved across all viseme types to indicate a resting mouth position, or in

other words, the character is not actively speaking. It differs from visemes that specifically have the mouth closed, because the character may have a resting mouth position that is not closed. Imagine, for example, an energetic politician that grins constantly while in public.










Toon Boom Viseme Codes

Code	Example Image	Description	Code	Example Image	Description
-		Resting	D		A and I
A		M, B, and P	E		O
B		Most consonants	F		W and Q
C		E	G		F and V

Blair Viseme Codes

Code	Example Image	Description	Code	Example Image	Description
-		Resting	C		Most consonants
A		A and I	W		W and Q
E		E	M		M, B, and P
O		O	L		L
U		U	F		F and V

RMS Viseme Codes

Code	Example Image	Description	Code	Example Image	Description
-		Resting	5		
0		Closed	6		
1		Minimally open	7		
2			8		
3			9		Fully open
4					

File Format Definition

The Lipz file format uses JSON. It is UTF-8 encoded, so may contain Unicode characters. Though JSON supports a hierarchical structure, all elements are declared flat under the root element. Elements that can be used are described in the table below. All elements are optional, but the usefulness of the file to a game relies on at least some of the elements to be present.

Element	Type	Constraints	Description
<code>fps</code>	<i>number</i>	Between 1 and 1000	Frame rate of viseme specification expressed in frames per second. Used to interpret values of <code>viseme</code> elements, where one character equals one frame. Note that the game frame rate is not tied to this value. If unspecified, the value is supplied externally, e.g. a known constant value in game code or a general settings file.
<code>text</code>	<i>string</i>	Segment delimiter count must match <code>visemes</code> .	Dialogue text that corresponds to audio. This form is used when the Lipz file contains text in just one language. Text may be segmented with use of the <code> </code> character (see more about this in "Segmenting" section). Custom meta-information not intended for display can be stored with the text inside of curly braces (<code>{</code> and <code>}</code>).
<code>text_</code> <code>//</code> <code>-</code> <code>cc</code>	<i>string</i>	Segment delimiter count must match <code>visemes</code> .	Dialogue text that corresponds to audio with the additional specification of a language locale. This form is used when the Lipz file contains text for more than one language. The language locale consists of a lower-case, 2-character, ISO 639-1 language code followed by a hyphen (<code>-</code>) followed by a lower-case, 2-character, ISO 3166-1 country code. Example: <code>text_es-mx</code> (Mexican Spanish). Text may be segmented with use of the <code> </code> character (see more about this in "Segmenting" section). Custom meta-information not intended for display can be stored with the text inside of curly braces (<code>{</code> and <code>}</code>).

<code>visemes</code>	<i>string</i>	Constrained according to <code>viseme_type</code> .	Each character represents the viseme for one frame of animation. The exception is the <code> </code> character which delineates segments, and does not count for a frame.
<code>viseme_type</code>	<i>string</i>	"blair", "toonboom", or "rms"	The type of viseme encoding used by the <code>visemes</code> element. If unspecified, the value is supplied externally, e.g. a known constant value in game code or a general settings file.

Segmenting

Segmenting is optional. It can be used to correlate ranges of frames to text or events. With segmenting, the following is possible:

- Subtitles displayed on a screen can be timed to correspond to the playback of matching dialogue from an audio file.
- Algorithms can be devised for displaying subtitles in an automated and localization-optimized way, e.g. "display however many segments fit on screen" will handle large differences in translation text length better than just "display segment X on screen".
- Animated character emotions and gestures can be timed to correspond to an exact moment of playing dialogue.

Segments are always delineated with the `|` character. When used within the `visemes` element, they indicate the exact frame of a segment's boundary. In the example below, the first segment lasts 20 frames, and is then followed by the second segment.

```
{
  "visemes" : "125643356734563-----|9723399421---"
}
```

If text is specified that also contains segmenting characters (`|`) then the segments defined in the text should be interpreted as corresponding to the segments defined in the viseme. This also implies timing values from the visemes correspond to the text. In the example below, the game code that loads the Lipz file will know that the "It was you!" text begins on the 20th frame.


```
{  
  "visemes" : "125643356734563-----l9723399421---",  
  "text" :    "Let me think...!It was you!"  
}
```

Imagine a character on the screen who thinks carefully for a moment, ("Let me think...") and then screams an unexpected accusation ("It was you!"). For dramatic effect, it may be better for the game to show the second segment of text only after the corresponding dialogue has been played. Segmenting gives you this level of control.

Adding to this example, it would be extra-dramatic to have the character point angrily at the accused at the moment the second segment begins. For this, we can add meta-information into the text.

```
{  
  "visemes" : "125643356734563-----l9723399421---",  
  "text" :    "Let me think...!{POINTS BUTLER}It was you!"  
}
```

For the meta-code to do anything in the game, you'd need to write handling that looked for "POINTS BUTLER" in a currently playing segment to trigger a pointing-at-butler animation.

Knowing which Viseme to Show

The game code will perform the tasks of loading and playing an audio file containing dialogue. At the same time it does this, it should also load the Lipz file that corresponds to the audio file.

Later, when the game code plays an audio file, there should be some facility of learning the current play position within the audio file. At any point where a speaking character's lips should be updated, a function like the one below can be called to learn the current viseme that should be used for display. The game code can then use whatever facility is provided for animation to set the character's lips to match the current viseme.

```

//Returns one character code of the current viseme corresponding to playing dialogue
audio.
function getCurrentViseme(
    playPositionSecs, //Float of seconds elapsed since dialogue audio began play
ing.
    visemes,          //String loaded from "visemes" element of Lipz file.
    fps) {            //Frame rate used for interpreting length of each frame in
"visemes".

    var unsegVisemes = visemes.replace(/\|/g, ''); //Removes any segment delimiter "
|" characters.
    var frameNo = Math.floor( playPositionSecs * fps );
    if (frameNo < 0 || frameNo >= visemes.length) { //Frame# is out of bounds.
        return '-'; //Return "resting mouth" viseme as a safe default.
    } else { //Frame# is valid.
        return unsegVisemes.charAt(frameNo);
    }
}

```