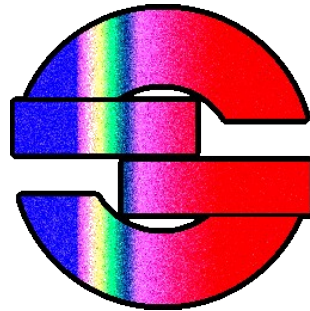# SPPAS: Automatic Phonetic Annotation of SPeech

**Brigitte Bigi**

**Keywords:**

Phonetizaton          Automatic          Speech

Syllabification     Segmentation

Alignment          Prosody

# What SPPAS 1.4.9 can do?

- Automatic annotations:

  - **Momel/INTSINT**: Modelisation of Mélodie

  - **IPUs segmentation**: utterance level segmentation

  - **Tokenization**: text normalization

  - **Phonetization**: grapheme to phoneme conversion

  - **Alignment**: phonetic segmentation

  - **Syllabification**: group phonemes into syllables

- Components:

  - Wav Player

  - Manual transcription based on IPUs segmentation
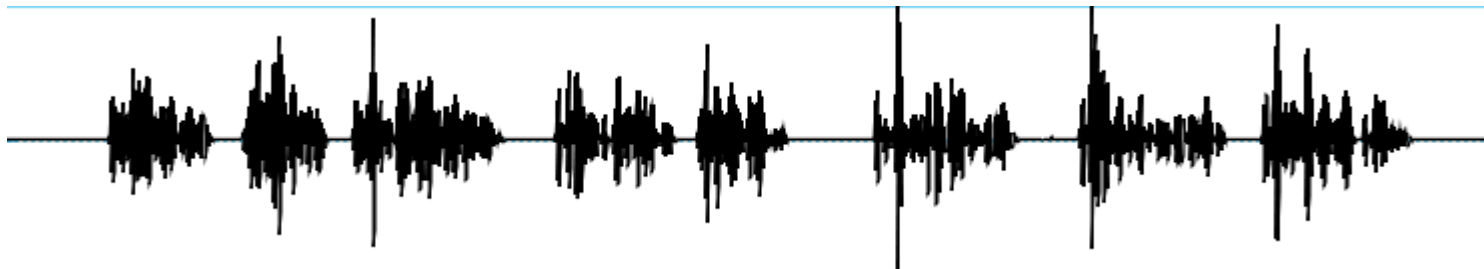
  - Get tiers information and Filter tiers

# Key-points

- A tool dedicated to computer scientists **and** linguists

- Language-independent algorithms
  - Resources for French, English, Italian, Chinese, Taiwanese, partially Vietnamese
  - There is an easy way to add other languages

- GNU Public License

# Automatic annotations: inputs
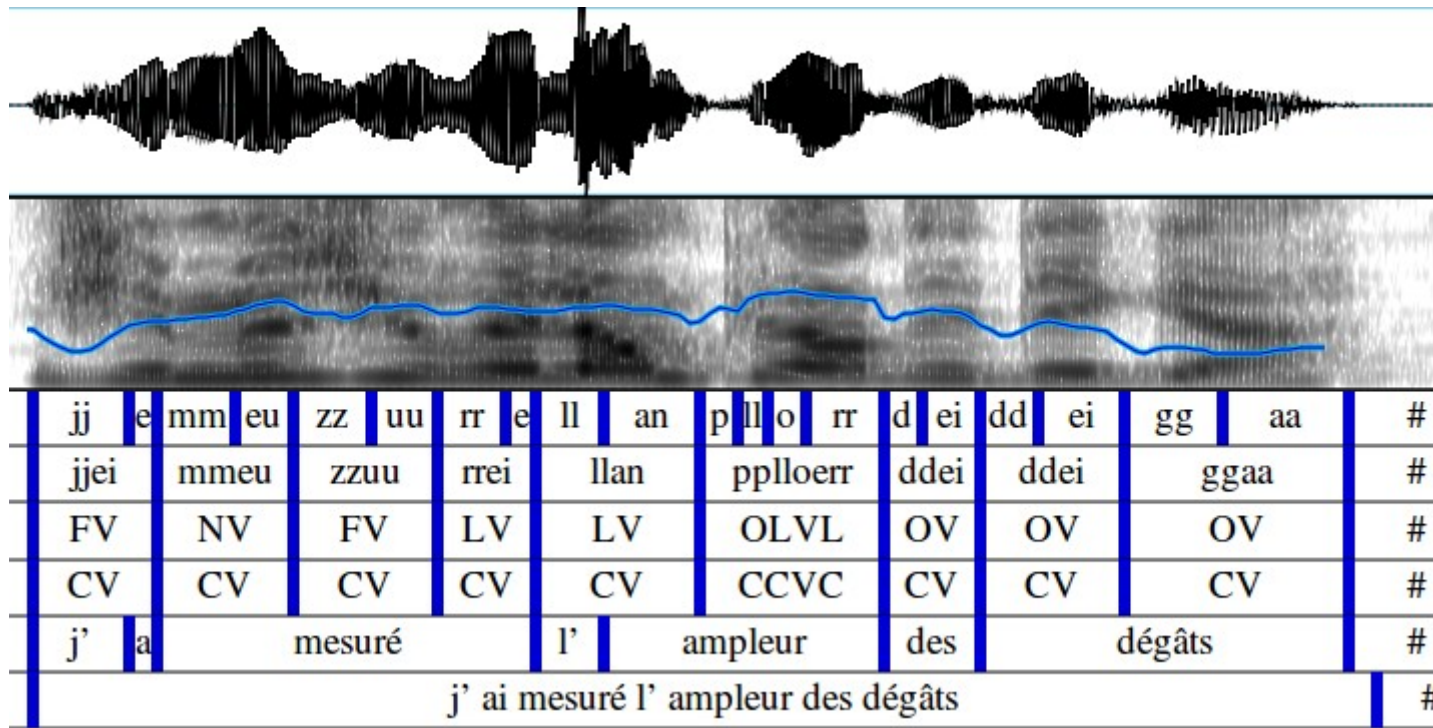
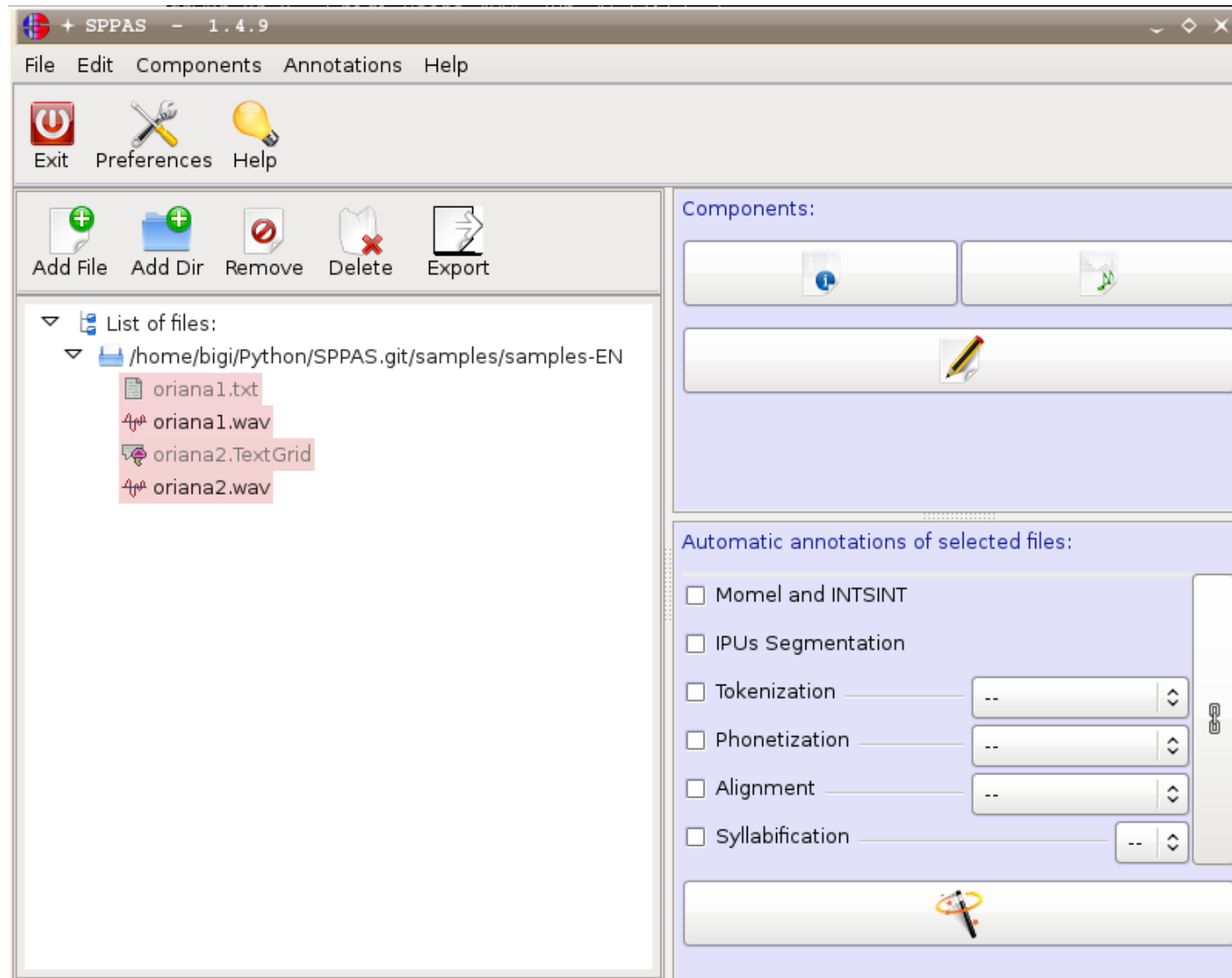- Speech signal: wav file



- Transcription: txt or TextGrid



assis sur le mur du jardin potager
j' ai mesuré l' ampleur des dégâts
les choux avaient été entièrement dévorés par les limaces
le potager était complètement dévasté
et ressemblait à un terrain en friche
mais pourquoi est-ce_que j' ai pas pensé à mettre du tue_limaces
au point où j' en suis si je m' écoutais je ferais tout cimenter
comme ça j' aurais une belle cour intérieure et plus de soucis

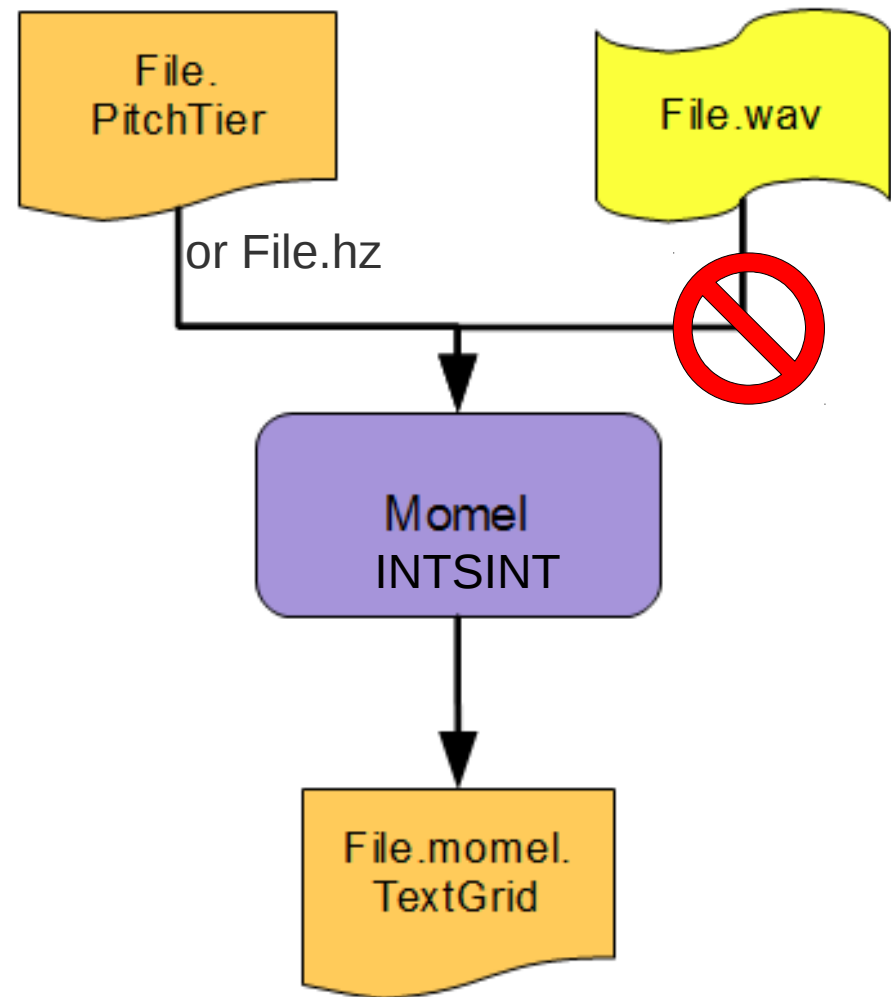# Automatic annotations: outputs

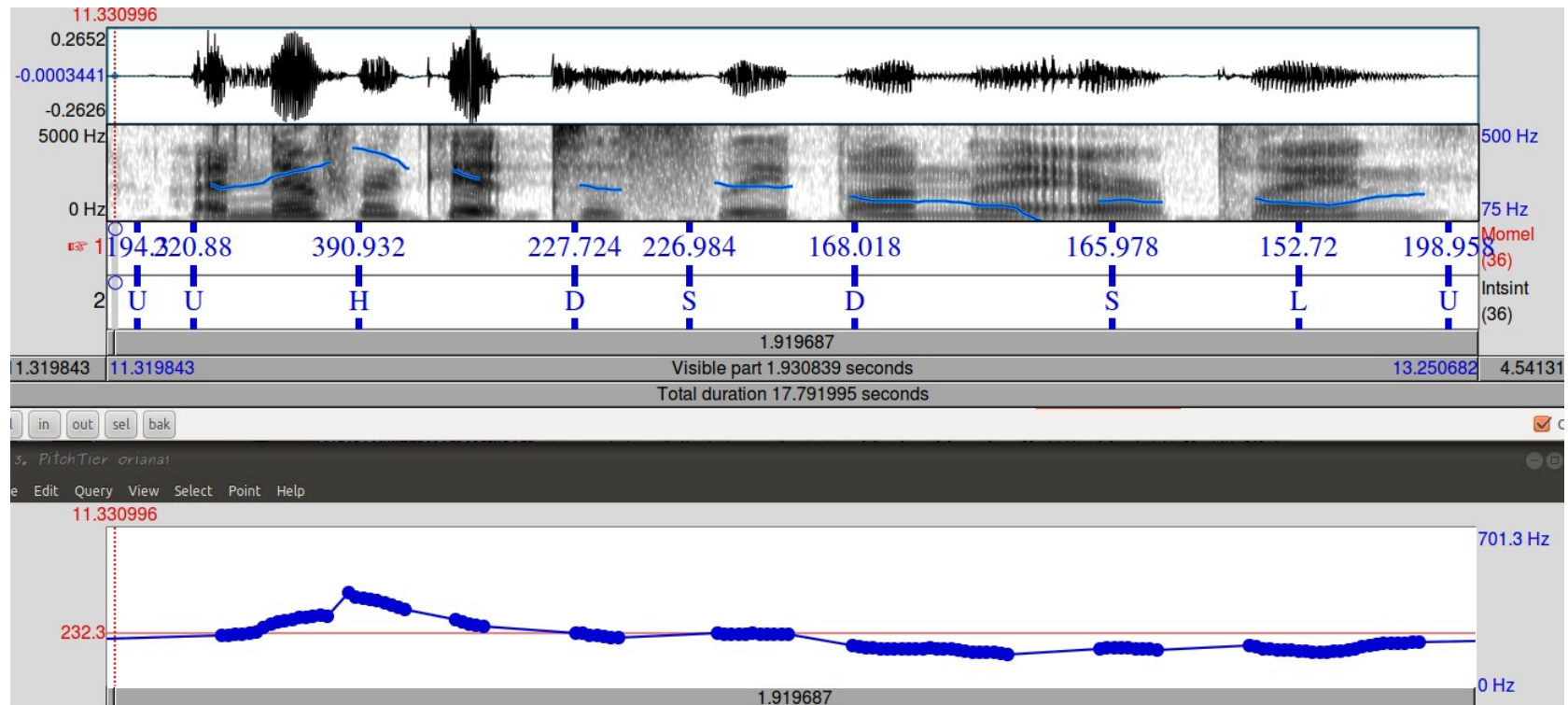- A set of TextGrid files

# Screenshot

# Momel and INTSINT

- SPPAS implements Momel and INTSINT: Daniel Hirst

- A file with pitch values is required (one value each 10ms).

# Momel/INTSINT: example
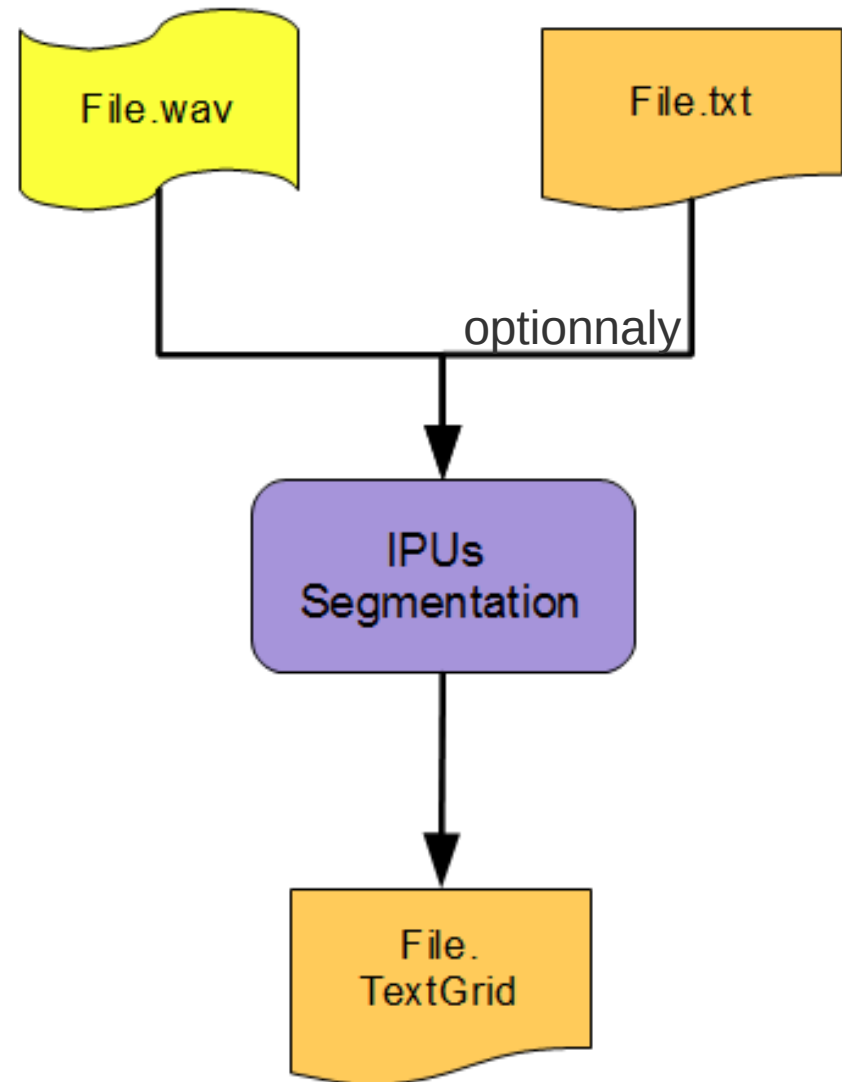
- Output: a TextGrid file with 2 tiers
  - Momel targets (pitch values)
  - INTSINT annotation of these targets

# IPUs segmentation

- Inter-Pausal Units segmentation

- The algorithm computes a heuristics based on the detection of silences, by using:

  - volume

  - min silence duration

  - min speech duration



File.wav    File.txt

optionnaly
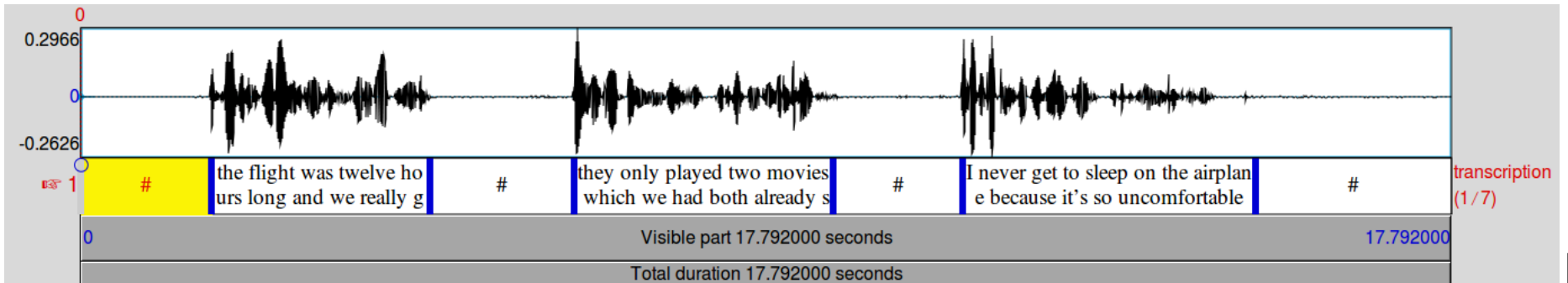
IPUs Segmentation

File. TextGrid

# IPUs segmentation: example

Transcription: silences are indicated by newlines or '#'

the flight was twelve hours long and we really got bored
they only played two movies which we had both already seen
I never get to sleep on the airplane because it's so uncomfortable

# Tokenization

- Process of normalizing text:
    - Remove punctuation, comments etc.
    - Convert numbers to letters:
        - 2 → deux
    - Convert characters to the lower form
    - Segment into words/tokens:
        - parce que → parce_que

- An algorithm as language independent as possible (dict-based)
    - But... *bien que* tu sois là, je pars / c'est *bien que* tu sois là.

File .TextGrid

Tokenization ← L.vocab

File-tokens .TextGrid

# Phonetization

- Process of representing sour
  with phonetic signs

- The phonetization is the
  equivalent of a sequence of
  dictionary look-ups.

- Phonetic variants:

  - no rules are applied, all
    possibilities are stored



File-tokens TextGrid → Phonetization ← L.dict → File-phon TextGrid

# Phonetization: example
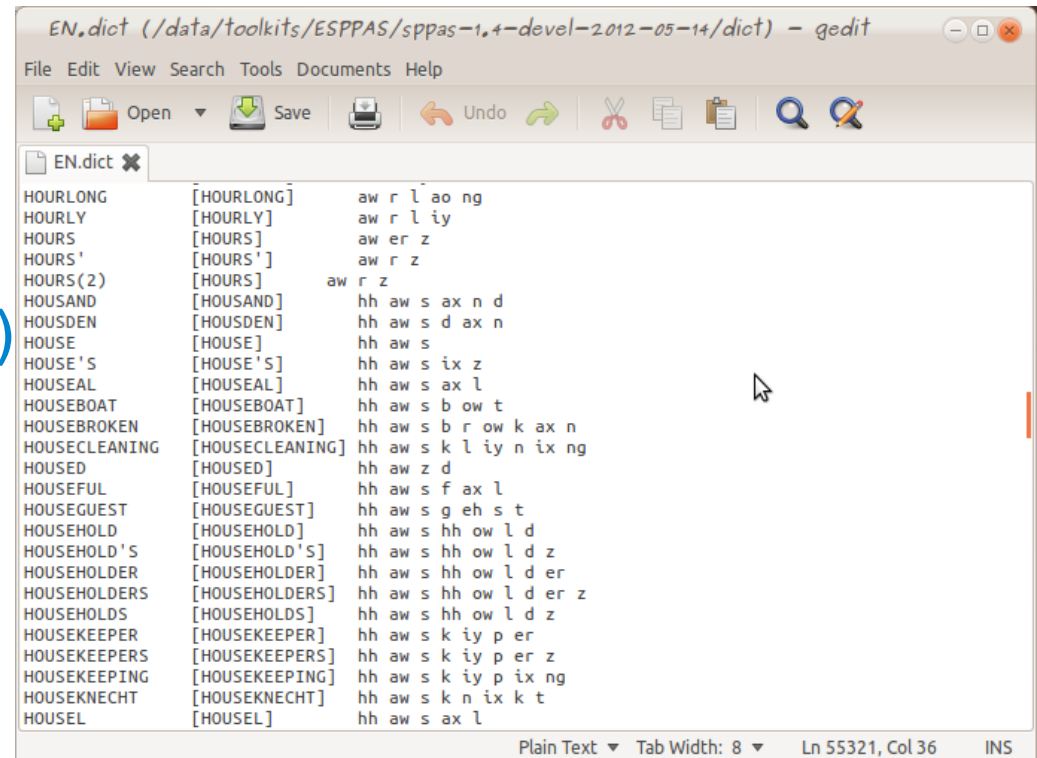
- Resources:

  - a dictionary

    (HTK-ASCII format)

EN.dict (/data/toolkits/ESPPAS/sppas-1.4-devel-2012-05-14/dict) – gedit

File   Edit   View   Search   Tools   Documents   Help

Open ▼   Save

EN.dict

```
HOURLONG        [HOURLONG]         aw r l ao ng
HOURLY          [HOURLY]           aw r l iy
HOURS           [HOURS]            aw er z
HOURS'          [HOURS']           aw r z
HOURS(2)        [HOURS]       aw r z
HOUSAND         [HOUSAND]          hh aw s ax n d
HOUSDEN         [HOUSDEN]          hh aw s d ax n
HOUSE           [HOUSE]            hh aw s
HOUSE'S         [HOUSE'S]          hh aw s ix z
HOUSEAL         [HOUSEAL]          hh aw s ax l
HOUSEBOAT       [HOUSEBOAT]        hh aw s b ow t
HOUSEBROKEN     [HOUSEBROKEN]      hh aw s b r ow k ax n
HOUSECLEANING   [HOUSECLEANING]  hh aw s k l iy n ix ng
HOUSED          [HOUSED]           hh aw z d
HOUSEFUL        [HOUSEFUL]         hh aw s f ax l
HOUSEGUEST      [HOUSEGUEST]       hh aw s g eh s t
HOUSEHOLD       [HOUSEHOLD]        hh aw s hh ow l d
HOUSEHOLD'S     [HOUSEHOLD'S]      hh aw s hh ow l d z
HOUSEHOLDER     [HOUSEHOLDER]      hh aw s hh ow l d er
HOUSEHOLDERS    [HOUSEHOLDERS]     hh aw s hh ow l d er z
HOUSEHOLDS      [HOUSEHOLDS]       hh aw s hh ow l d z
HOUSEKEEPER     [HOUSEKEEPER]      hh aw s k iy p er
HOUSEKEEPERS    [HOUSEKEEPERS]     hh aw s k iy p er z
HOUSEKEEPING    [HOUSEKEEPING]     hh aw s k iy p ix ng
HOUSEKNECHT     [HOUSEKNECHT]      hh aw s k n ix k t
HOUSEL          [HOUSEL]           hh aw s ax l
```

Plain Text ▼   Tab Width: 8 ▼      Ln 55321, Col 36      INS

```
the flight was twelve hours long and we really got bored
```
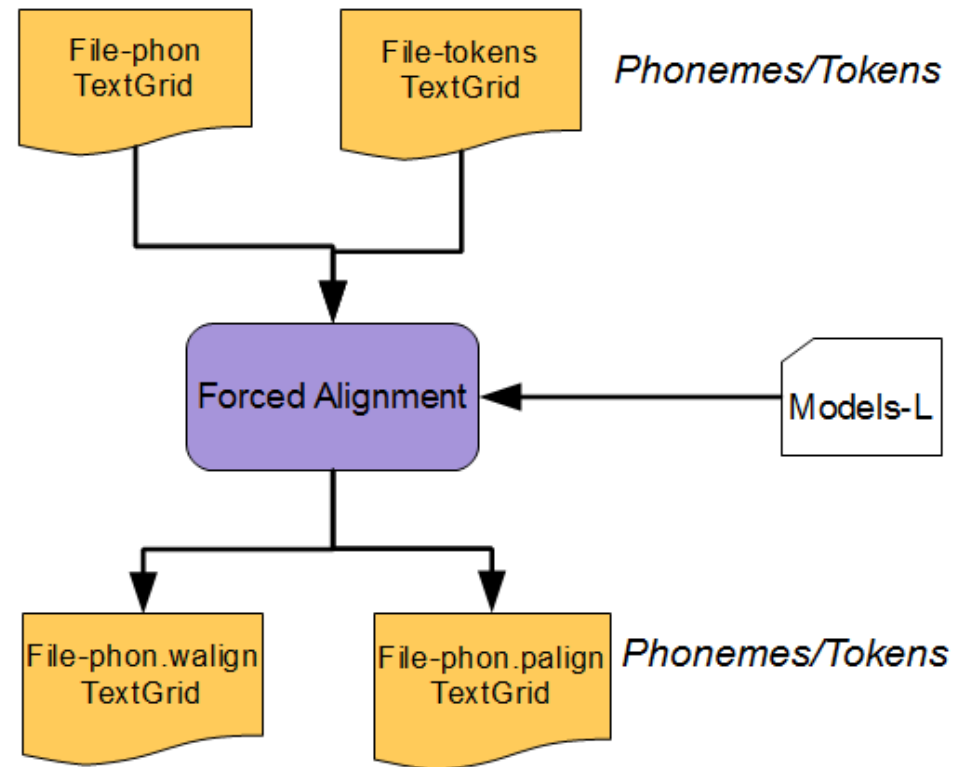
is phonetized as follow:

```
dh.ax|dh.ah|dh.iy    f.l.ay.t    w.aa.z|w.ah.z|w.ax.z|w.ao.z    t.w.eh.l.v
aw.er.z|aw.r.z   l.ao.ng   ae.n.d|ax.n.d   w.iy   r.ih.l.iy|r.iy.l.iy   g.aa.t
b.ao.r.d
```
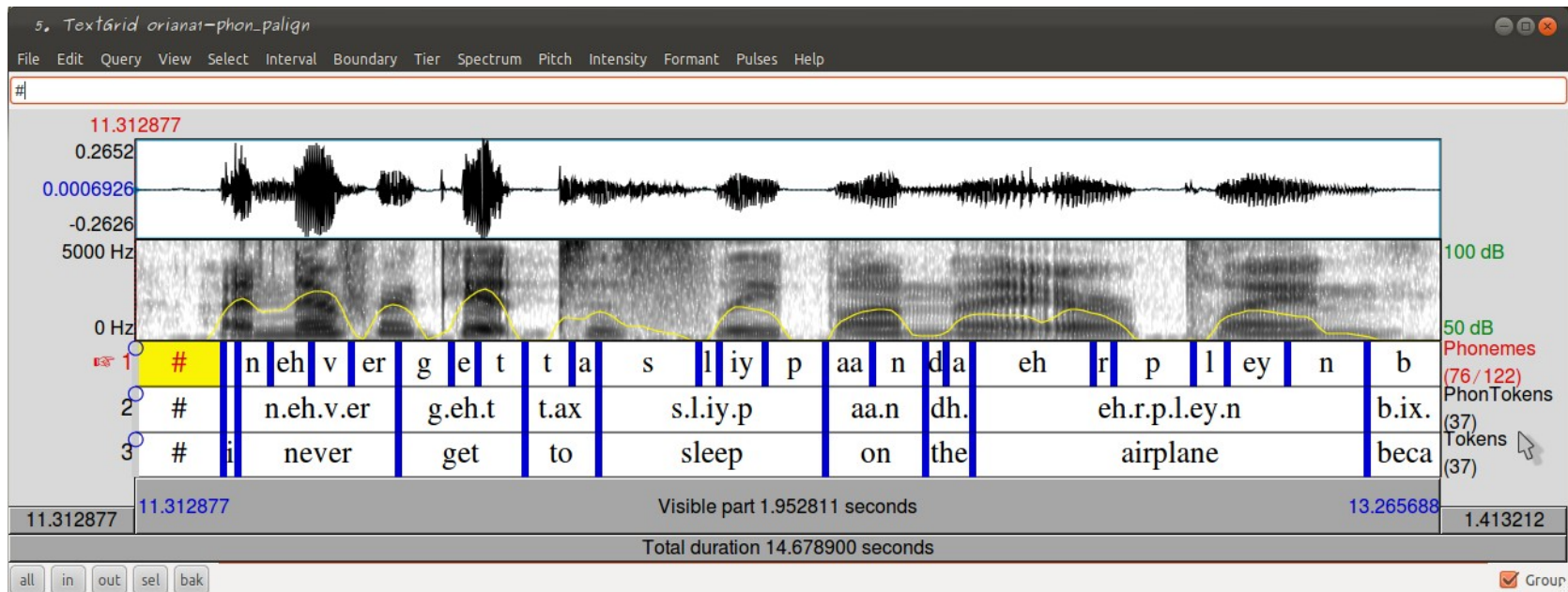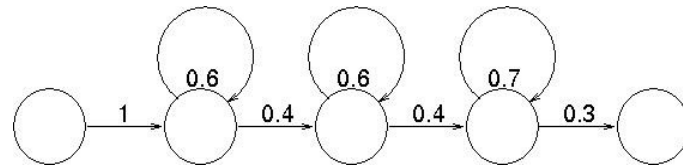
# Alignment

- A time-matching between a given speech utterance along with a phonetic representation of the utterance

- Forced-alignment in SPPAS is based on the **Julius** Speech Recognition Engine



- The alignment task is a 2-step process:

  - the first one: choose the phonetization;
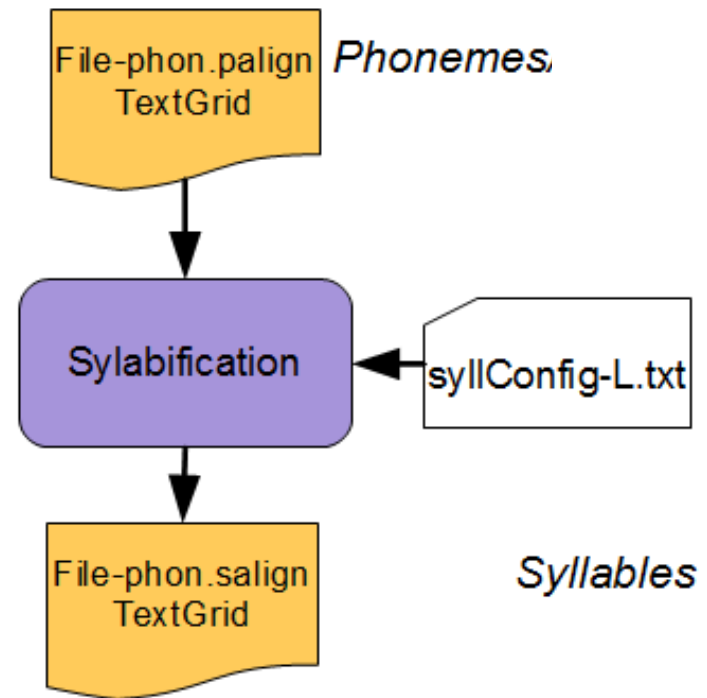
  - the second one: perform the segmentation.

# Alignment: example

- Resources:

  - A finite state grammar that describes sentence patterns to be recognized (created by SPPAS);

  - An acoustic model.

# Syllabification

- Development of a Rule-Based System for automatic syllabification of phonemes' strings

- The syllabification is based on 2 principles:

  - a syllable contains a vowel, and only one;

  - a pause is a syllable boundary.

V C C V
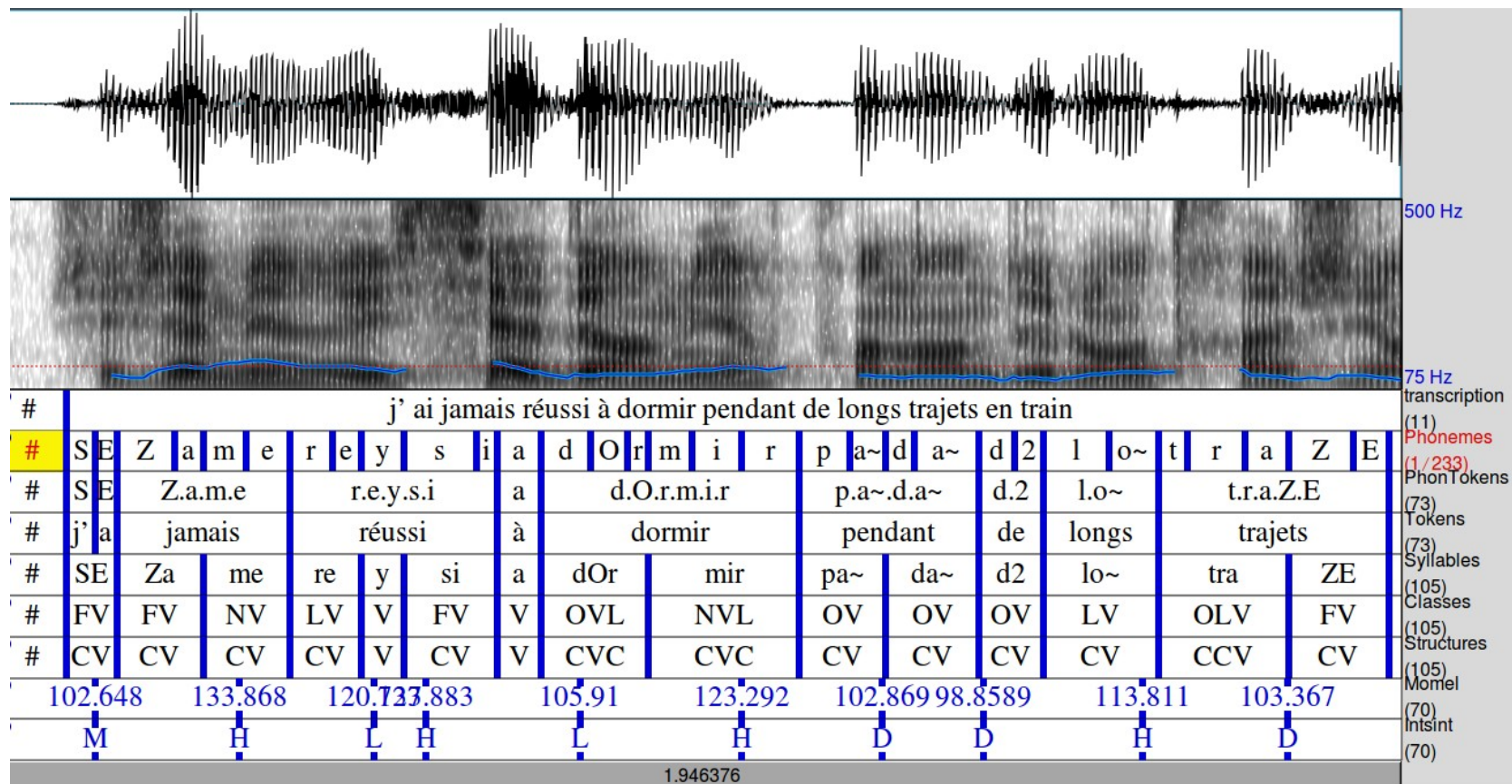
File-phon.palign TextGrid — *Phonemes*

Sylabification ← syllConfig-L.txt

File-phon.salign TextGrid — *Syllables*

# Syllabification: example

- Resources (FR and IT):

  - a configuration file with the phoneme set, the classes and all rules

# Resources summary

| | FR | IT | ZH | EN |
|---|---|---|---|---|
| Dictionary : Number of entries | 350k words and 300k variants | 390k words and 5k variants | 88k words (350 syllables) | 121k words and 10k variants |
| Acoustic model: Data to train | Triphones - 7h30 CID +30min read | Triphones - 3h30 map-task | Monophones - 90min read | Triphones See voxforge.org |

SLDR forge     Evalita 2011     Eurom1     CMU dictionary

# A few words about technical stuff...

- The transcription encoding must correspond to that of SPPAS dictionary:

  - UTF-8 for French, Chinese or Italian,

  - us-ascii for English.

- The transcription and the audio files must have the same name (except for the extension)

Recorded input speech files are **mono wav** files only.
Other file formats are not supported.

*SPPAS* verifies if the wav file is 16 bits and 16000 Hz sample rate.
Otherwise it automatically converts to this configuration using sox.

# About

- URL: http://www.lpl-aix.fr/~bigi/sppas/

- Supported by the Equipex ORTOLANG

- SPPAS can achieve a set of automatic phonetic annotations of speech; results are depending on...

  - The resources quality;

  - The input wav quality;

  - The transcription quality...

# References

[1] B. Bigi, C. Meunier, I. Nesterenko, R. Bertrand. *Automatic detection of syllable boundaries in spontaneous speech.* Language Resource and Evaluation Conference (LREC), pp 3285-3292, La Valetta, Malte, 2010.

[2] Brigitte Bigi (2011). *A Multilingual Text Normalization Approach.* 2nd Less-Resourced Languages workshop, 5th Language  Technology Conference, Poznàn (Poland). 2011.

[3] B. Bigi. *The SPPAS participation to Evalita 2011.* Working Notes of EVALITA 2011, Rome, Italy, ISSN: 2240-5186, January 2012.

[4] B. Bigi, D. Hirst. *SPeech Phonetization Alignment and Syllabification: a tool for the automatic analysis of speech prosody.* Speech Prosody, Shanghai, China, May 2012, .

[5] B. Bigi.  *SPPAS: a tool for SPeech Phonetization Alignment and Syllabification.* Language Resource and Evaluation Conference (LREC), Istanbul, Turkey, May 2012. .

[6]. Bigi. *Forced Alignment on Spontaneous Speech for Italian: the SPPAS tool.* B. Magnini et al. (Eds.): EVALITA 2012, LNCS 7689, pp. 312--321. Springer, Heidelberg.

http://www.lpl-aix.fr/~bigi/sppas/