



# EVALITA 2011

*Evaluation of NLP and Speech Tools for Italian*

EVALITA 2011 is the third evaluation campaign of Natural Language Processing and Speech tools for Italian, supported by the NLP working group of AI\*IA (*Italian Association for Artificial Intelligence*) and AISV (*Italian Association of Speech Science*)

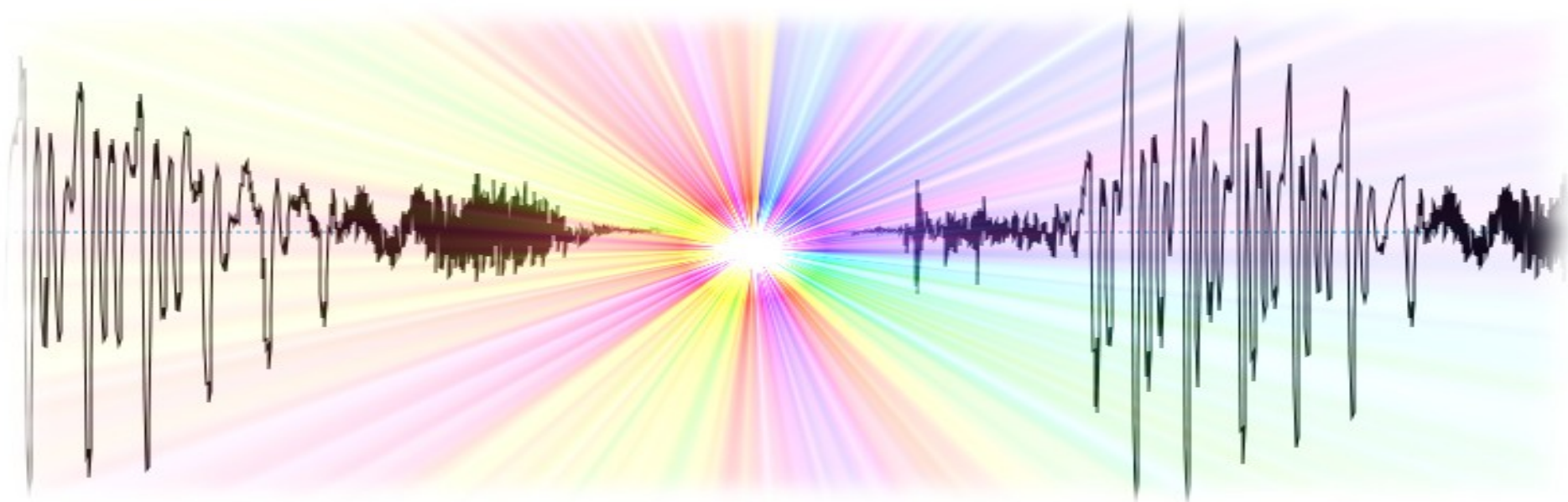
<http://www.evalita.it/2011>

## The SPPAS participation to Evalita 2011

Brigitte Bigi



# ma che cos' è SPPAS?

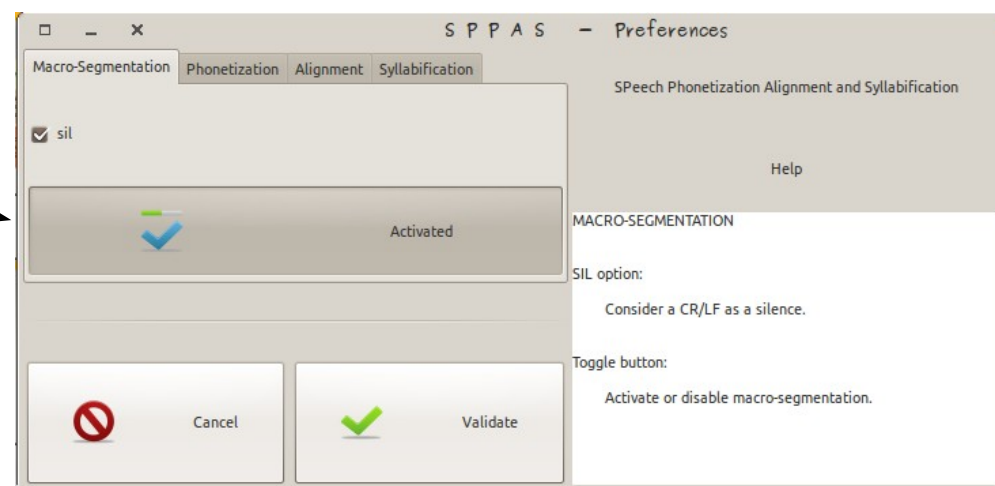
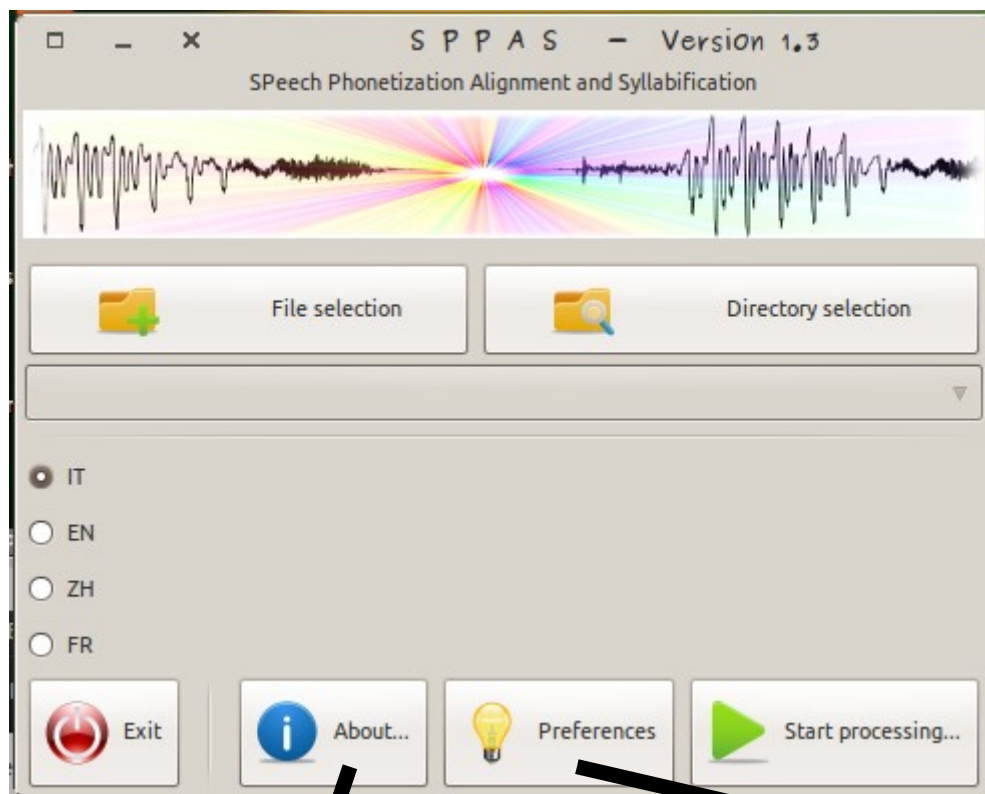


**SP**eech **P**honetization **A**lignment and **S**yllabification

# Main description

- A new tool to produce automatically annotations which includes utterance, word, syllabic and phonemic segmentations from a recorded speech sound and its transcription.
- Language-independent.
- Currently designed for French, English, Italian and Chinese and there is an easy way to add other languages.
- Distributed under GPL license.

# SPPAS Screenshots



# Evalita task

- "Forced Alignment on Spontaneous Speech":
  - ✓ Phone segmentation;
  - ✓ Word segmentation.
- Data: ✓ Closed; · Opened.
- Dialogues, map-tasks:
  - 3h30 speech;
  - 15% phones: pauses, filled-pauses, garbage.
- SPPAS Forced-Alignment is 2 sub-tasks:
  - phonetization + alignment

# Phonetization

- The process of representing sounds by phonetic signs.
- There are two general ways to construct a phonetization process:
  - rule based systems (with rules based on inference approaches or proposed by expert linguists);
  - dictionary based solutions which consist of storing a maximum of phonological knowledge in a lexicon.

# SPPAS phonetization

- SPPAS uses the dictionary-based approach.
- The phonetization is the equivalent of a sequence of dictionary-look-ups:
  - Input transcription needs to be word-segmented
  - It is supposed that all words of the transcription are mentioned in the pronunciation dictionary.
- A specific phone to represent filled pauses.

# Dictionary

- The dictionary contains a set of possible pronouciations of words, including accents as "*perchè*" pronounced as /b e r k e/, and reduction phenomena as /p e k/.
- Made of:
  - 390k words;
  - 5k variants.
- From:
  - Festival + Evalita training corpus.
- Manual corrections

```

264377 PESAVO [PESAVO] p e s a v o
264378 PESCA [PESCA] p E s k a
264379 PESCA(2) [PESCA] p e s k a
264380 PESCADOR [PESCADOR] p e s k a d o r
264381 PESCAGGI [PESCAGGI] p e s k a d z i
264382 PESCAGGIO [PESCAGGIO] p e s k a d z o
264383 PESCAI [PESCAI] p e s k a i
264384 PESCAIA [PESCAIA] p e s k a j a

```

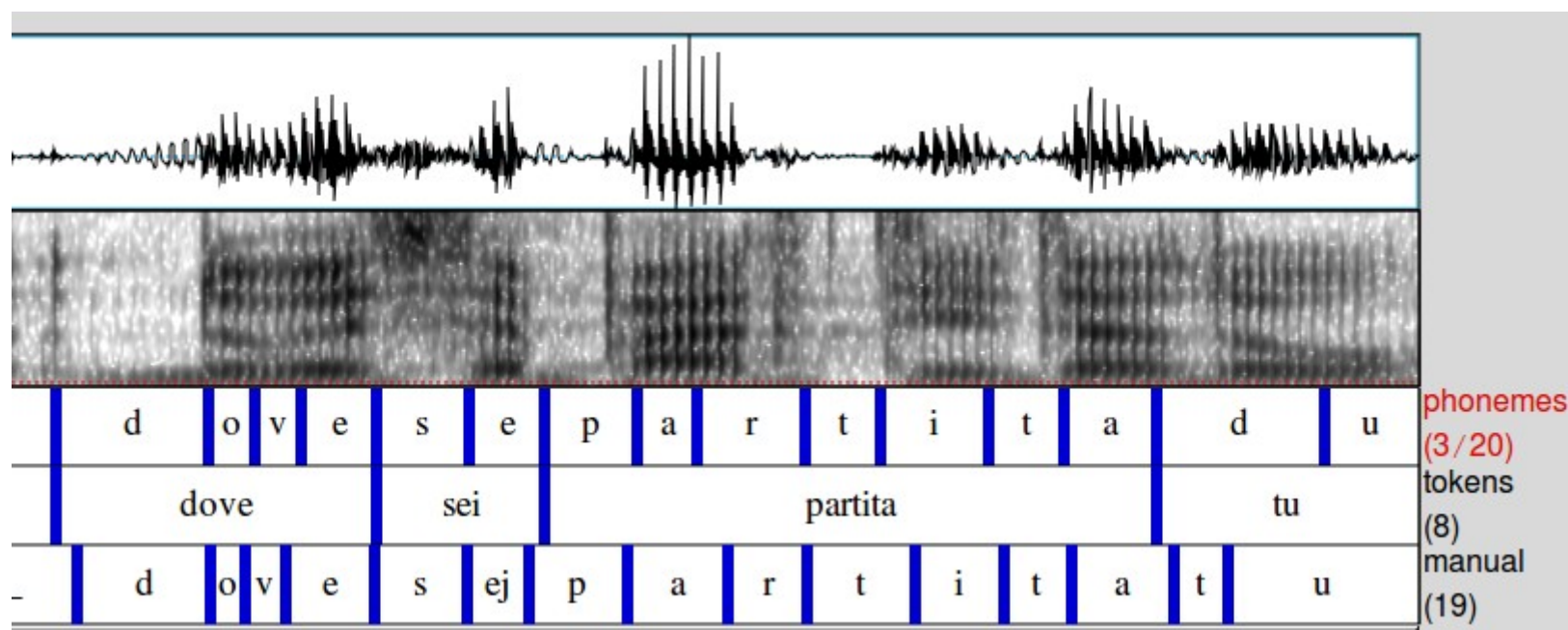


# Phonetization: variants

- No rules are applied:
  - all possibilities are proposed to the aligner.
- For example, the sentence "*del mio frigorifero*" will produce the following list of words with associated pronunciations:
  - d.e.l | d.E.l
  - m.jo | m.i
  - f.r.i.g.o.r.i.f.e.r.o | f.r.i.g.o.r.i.f.E.r.o |  
f.r.i.g.o.r.i.f.e.r | f.r.i.g.o.r.i.f.e.r.0

# Alignment

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



# SPPAS Alignment

- Alignment in SPPAS is based on the Julius Speech Recognition Engine (SRE):
  - A finite state grammar that describes sentence patterns to be recognized;
  - An acoustic model.
- The alignment task is a 2-steps process:
  - the first one choose the phonetization;
  - the second one perform the segmentation.

# Grammar

- A grammar: constraints on what the SRE can expect as input. It is a list of words that the SRE listens for. Each word has a set of associated list of phonemes.

Dictionary	Grammar
0 w_0 d e l	0 2 1 0 1
0 w_0 d E l	1 1 2 0 0
1 w_1 m jo	2 0 3 0 0
1 w_1 m i	3 -1 -1 1 0
2 w_2 f r i g o r i f e r o	
2 w_2 f r i g o r i f E r o	
2 w_2 f r i g o r i f e r	
2 w_2 f r i g o r i f e r 0	

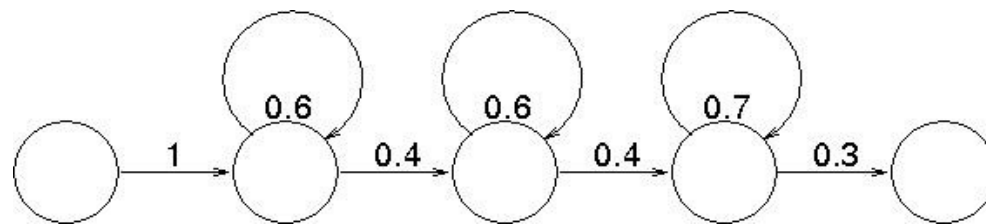
# Acoustic Model

- HMM, 5-states

- Triphones

- Trained with HTK:

- from the proposed phonetized transcription, without using the phonetic time-alignment;
- using 16 bits, 16000 hz wav files;
- Mel-frequency cepstrum coefficients (MFCC) along with their first and second derivatives were extracted from the speech in the standard way: MFCC\_D\_N\_Z\_0.



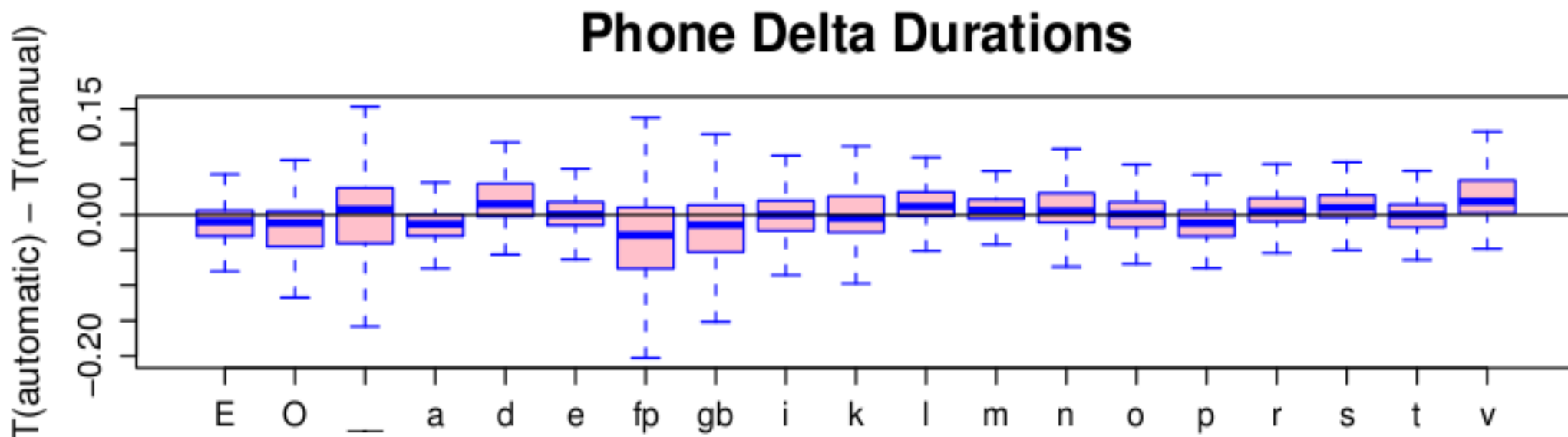
# Results

- Development corpus:
  - 200 utterances; 12 min 04 sec.
  - 2373 words, 6282 phonemes, including:
    - 689 “\_” (pauses);
    - 246 “#” (garbage);
    - both represent 14.88%.
- Evaluate separately:
  - Alignment;
  - Phonetization obtained after the alignment.

# Results alignment

- The availability of our system to align the good phoneme sequence.
- On the basis of the *manual phonetization*.
- Sclite using the time-alignment option :
  - a correct rate of 89.8%, with 7.6% substitutions, 2.6% deletions and 2.6% insertions.

# Results alignment



- Pauses, filled pauses and garbages: greatest ranges
- Vowels: automatic shorter than manual
- Consonants: automatic higher than manual



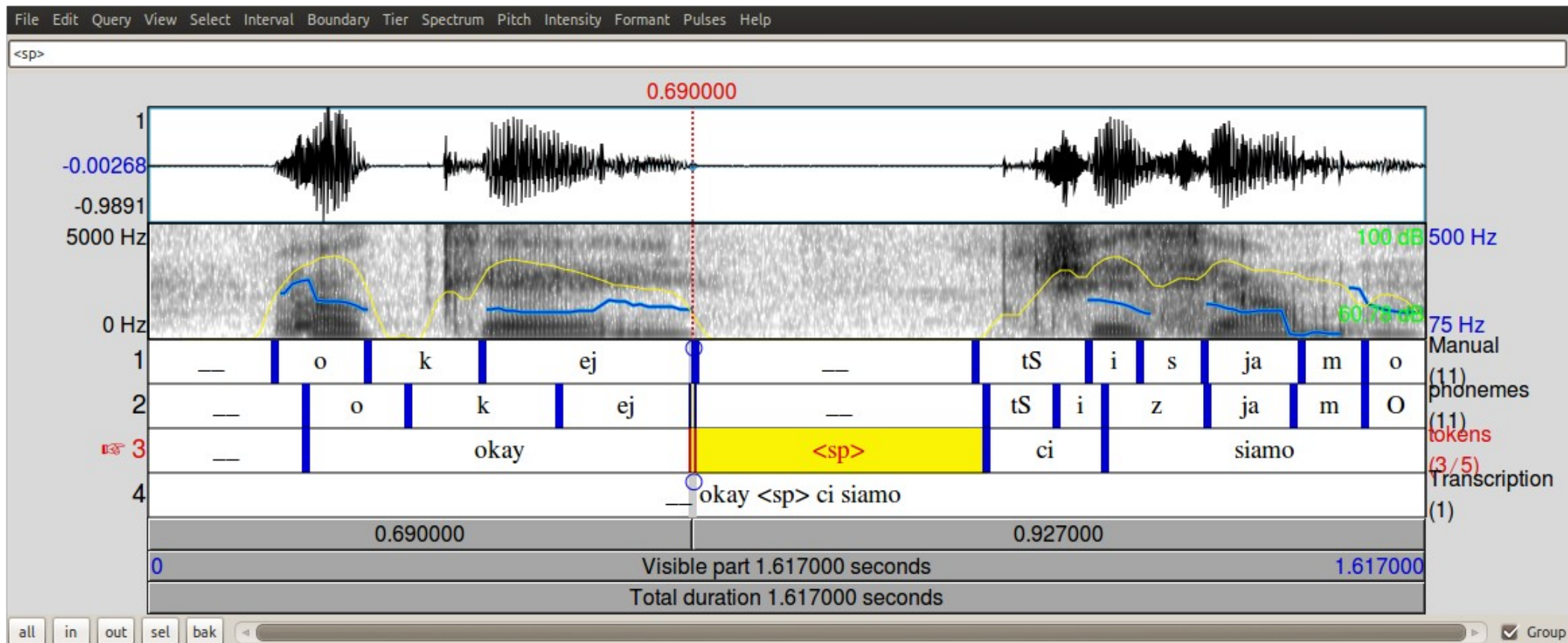
# Results phonetization

- The availability of our system to propose the expected phoneme sequence.
- Sclite *without* using the time-alignment option :
  - a correct rate of 89.5%, with 8.1% substitutions, 2.3% deletions and 6.9% insertions.
- Most frequent errors are due to the garbage manual annotation. Example: *bravissimo a questo*
  - Automatic: b r a v i s i m o a k w e s t o
  - Manual: b r # s # k w e s t o
    - 5 insertions, 2 substitutions!

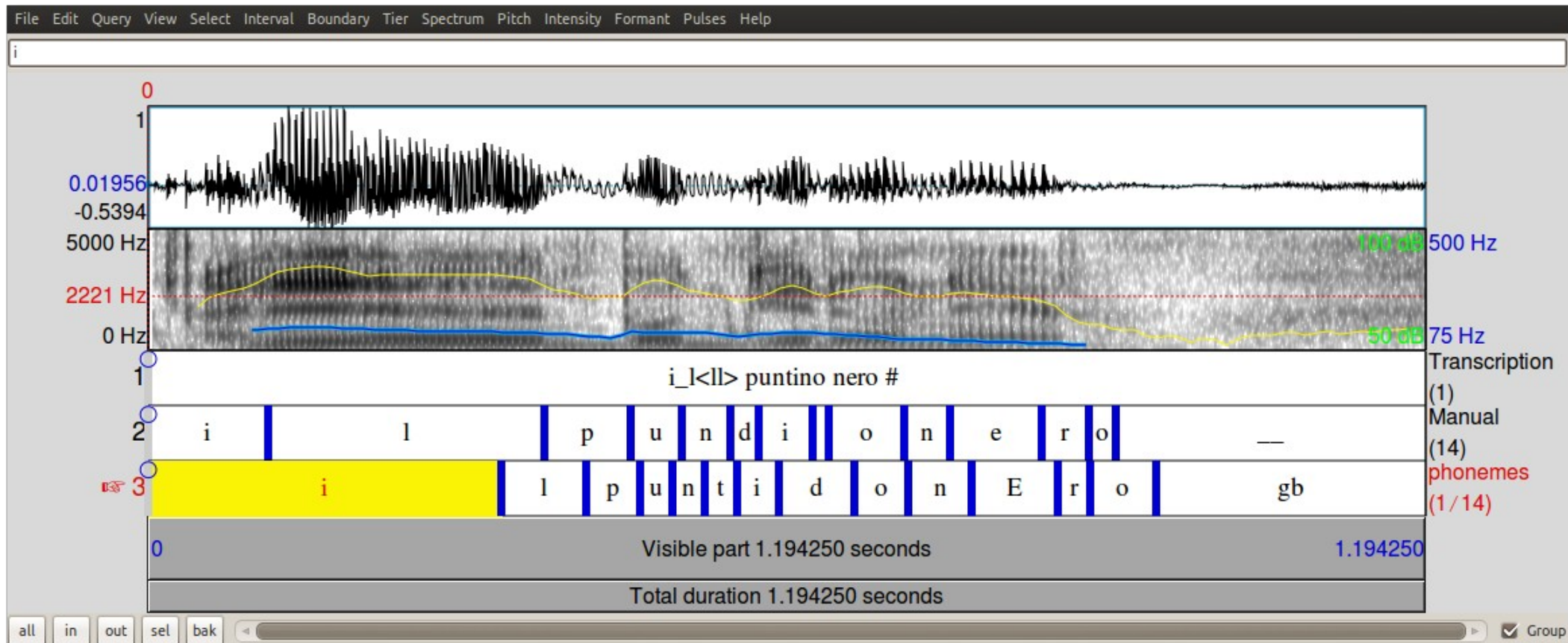
# SPPAS final results

- Official results estimated using sclite:
  - 88.4% good phoneme alignments:
    - This score contains both phonetizations and alignments errors.
  - 96.7% good word alignments.

# Example



# Example



# Conclusion

- SPPAS: a tool to perform the forced-alignment task during the Evalita 2011 campaign, on Italian map-task dialogues.
- SPPAS was not specifically devoted to Italian: it can deal with various languages: French, English, Chinese.
- <http://www.lpl-aix.fr/~bigi/sppas/>

