

1 Psycho-acoustics Inspired Automatic Speech Recognition

2 Gianpaolo Coro^{a,1,2,*}, Fabio Valerio Massoli^a, Antonio Origlia^b, Francesco Cutugno^b

3 ^a*Istituto di Scienza e Tecnologie dell'Informazione "Alessandro Faedo" – CNR, Pisa, Italy*

4 ^b*Università degli Studi di Napoli Federico II, Napoli, Italy*

5 Abstract

6 Understanding the human spoken language recognition process is still a far scientific goal. Nowadays,
7 commercial automatic speech recognisers (ASRs) achieve high performance at recognising clean speech,
8 but their approaches are poorly related to human speech recognition. They commonly process the phonetic
9 structure of speech while neglecting supra-segmental and syllabic tracts integral to human speech recogni-
10 tion. As a result, these ASRs achieve low performance on spontaneous speech and require enormous costs
11 to build up phonetic and pronunciation models and catch the large variability of human speech. This paper
12 presents a novel ASR that addresses these issues and questions conventional ASR approaches. It uses alter-
13 native acoustic models and an exhaustive decoding algorithm to process speech at a syllabic temporal scale
14 (100-250 ms) through a multi-temporal approach inspired by psycho-acoustic studies. Performance com-
15 parison on the recognition of spoken Italian numbers (from 0 to 1 million) demonstrates that our approach
16 is cost-effective, outperforms standard phonetic models, and reaches state-of-the-art performance.

17 *Keywords:* Automatic Speech Recognition, Deep Learning, Long Short Term Memory, Convolutional
18 Neural Networks, Factorial Hidden Markov Models, Hidden Markov Models, Speech, Psycho-acoustics,
19 Syllables

20 1. Introduction

21 Spoken language recognition in human beings is natural, robust, and effective. People recognise each
22 other's words also in situations of high background noise and reverberation using complex multi-channel
23 information processing (Hawkins and Smith, 2001). However, emulating and understanding these mecha-
24 nisms goes beyond our current technological capabilities (Pieraccini, 2012; Markowitz, 2015). Nevertheless,

*Corresponding author

Email addresses: coro@isti.cnr.it (Gianpaolo Coro), massoli@isti.cnr.it (Fabio Valerio Massoli),
antonio.origlia@unina.it (Antonio Origlia), cutugno@unina.it (Francesco Cutugno)
Preprint submitted to Computers & Electrical Engineering

Telephone Number: +39 050 315 8210

²Fax Number: +39 050 621 3464

May 18, 2021

25 automatic speech recognition has evolved in the last decades to propose high-performance commercial prod-
26 ucts to the large public (Li et al., 2015; Mustafa et al., 2019). The approaches of modern automatic speech
27 recogniser (ASRs) are poorly related with human speech recognition processes, and building ASRs requires
28 significant economic investments (CBInsights, 2019). Indeed, high costs depend mainly on the manual
29 preparation of large corpora of audio samples annotated at multiple levels (usually from sentence to pho-
30 netic levels), pronunciation models, and grammars. As a result, the implementation of high-performance
31 ASRs is usually bounded to few and large corporations and their applications are confined to simple sen-
32 tence transcribers or interactive voice responders. These ASRs usually neglect acoustic indicators such as
33 intonation and supra-segmental tracts that are essential in human spoken dialogues, because these would
34 be expensive in terms of modelling and data preparation. Consequently, modern ASRs have still issues at
35 recognising spontaneous and conversational speech with a high accuracy (Szaszák et al., 2016; Sahu et al.,
36 2018; Naing and Pa Pa, 2018; Knill et al., 2019).

37 The classic ASR architecture we took as a reference to build our ASR, is made up of four main processes
38 (Figure 1): (i) feature extraction, (ii) acoustic unit recognition, (iii) language model, and (iv) decoding. The
39 first process extracts real numbered vectors of acoustic features out of an audio file. The temporal scale
40 of these features is usually strictly sub-segmental and reflects the search for stationary spectral conditions
41 as much as possible. The second process - acoustic unit recognition - extends the temporal scope of the
42 speech chain processing. The acoustic units used in most ASRs are classically related with the acoustic
43 characteristics of the phonotactic distribution of co-articulatory processes. Indeed, most ASRs use contextual
44 phones/tri-phones as acoustic units. This approach embeds coarse assumptions on the internal dynamics of
45 the speech signal and attempts to model a form of contextual prediction of acoustic phenomena related
46 with the nature of connected speech. Consequently, a classic ASR includes a large (~cubic) combination of
47 elementary acoustic phonetic models. A further process estimates the likelihoods of these combined models
48 to a segment of speech signal. State-of-the-art ASRs use between ~2,000 and ~10,000 hours of speech
49 to train these models (CMUSphinx, 2017; Mwiti, 2019). However, these acoustic models are not robust
50 enough to manage the reduction processes that are frequent also in clear speech and that are more easily
51 treated with a syllable-based approach (Greenberg, 1996; Ostendorf et al., 1996; Cutugno et al., 2018).
52 The third process - language model - calculates the joint probability of a sequence of words and guides

53 the ASR search among alternative words during the recognition. This model uses a grammar that specifies
54 the permissible structures of the language. Statistical grammars (e.g. N-grams, Dunning (1994)) are used
55 to model complex languages and are automatically learned from large textual corpora. On the basis of a
56 grammar, the language model assigns a higher probability to more likely word sequences. As a prerequisite,
57 this process requires specifying all possible allowed words (the *lexicon*) and all different ways in which the
58 used units of speech can be combined to build these words (*pronunciation models*). The language model
59 is dependent on the particular lexicon and dialogue context the ASR is meant to manage. State-of-the-art
60 large-vocabulary ASRs include a 10^6 order of magnitude words, and the language model is trained on tens
61 of gigabytes of reference texts (Huang et al., 2001; Google, 2019). The fourth process - decoding - combines
62 acoustic models with the language model to produce the most probable transcription of an input audio file. In
63 most ASRs, acoustic models are implemented as Hidden Markov Models (HMMs, Markov (1913); Rabiner
64 and Juang (1986)), and the decoding process is strictly dependent on their state-based nature, where initial
65 and final states establish the acoustic boundaries between consecutive speech units (Young et al., 1989). This
66 type of HMM is a double stochastic model defined by (i) a finite set of states (ii) a transition matrix between
67 consecutive states, (iii) a set of *emission* probability densities for each state to be associated with the vectors
68 of acoustic features at a certain time, and (iv) a set of initial-state probability densities. Alternative models
69 have been proposed to enhance HMM performance, for example Factorial HMMs (FHMMs, Ghahramani
70 and Jordan (1996)) use sets of HMMs all with the same number of states and independent of each other,
71 except for the fact that the emission probability of one state of an HMM depends also on the states of the other
72 HMMs (Logan and Moreno, 1998). Most state-of-the-art ASRs use Deep Learning models to (i) classify
73 speech-units, (ii) model HMMs emission densities, and (iii) extract acoustic features (Cosi, 1998; Cosi and
74 Hosom, 1999; Ahad et al., 2002; Abdel-Hamid et al., 2014; Hinton et al., 2012; Swietojanski et al., 2014).
75 In particular, Deep Neural Networks (DNNs) are commonly used to model emission probabilities (Povey
76 et al., 2011; Pan et al., 2012; Cosi, 2015) and in some cases are replaced by Recurrent Neural Networks and
77 Long Short-Term Memory models (LSTMs) (Sak et al., 2014; Soltau et al., 2016; Senior et al., 2015; Qu
78 et al., 2017). LSTMs model longer-term dependencies between the elements of the input sequence (Bengio
79 et al., 1994; Hochreiter, 1991; Massoli et al., 2019) and have demonstrated high performance when used to
80 classify single phonemes and syllables (Sak et al., 2014; Senior et al., 2015; Soltau et al., 2016; Qu et al.,

81 2017).

82 End-to-end ASRs are valid alternative architectures and can reach state-of-the-art performance (Rao
83 et al., 2017; Zhang et al., 2017; Chiu et al., 2018; Weng et al., 2018; Watanabe et al., 2018; Zeghidour et al.,
84 2018a,b; Jaitly et al., 2019; Sainath et al., 2019). These systems jointly learn all ASR components in one
85 integrated approach, which reduces training and decoding time. However, they require an amount of training
86 data that is by far higher than what is required by classic ASR architectures (Graves et al., 2006; Graves and
87 Jaitly, 2014; Novoa et al., 2018; Audhkhasi et al., 2019).

88 In this paper, a comparison between ASRs using both conventional and non-conventional approaches
89 is presented. In particular, a novel approach for an ASR is proposed (Figure 2) that uses several possible
90 alternative acoustic models. Each time the ASR is instantiated, one among four models is used for acoustic
91 unit modelling. Three of these models (FHMM, CNN, LSTM) are inspired by studies on the involvement
92 of psycho-acoustic related features of human speech recognition, i.e. the multi-temporal processing of the
93 speech signal at syllabic and phonetic levels (Greenberg, 1996; Jenkins and Strange, 1999; Hawkins and
94 Smith, 2001; Malaia and Wilbur, 2019). Our study complies with the idea that although speech recognition
95 in humans and machines is implemented in different ways, they should compute the speech signal in a similar
96 way (Marr, 1982). Other studies have investigated this similarity at a computational level, to build ASRs that
97 accounted for the high variability of acoustic realisations of lexical representations, speaker independence,
98 and new-word recognition (Scharenborg et al., 2005). For example, the Shortlist and SpeM ASRs addressed
99 these properties by pursuing the hypothesis that human speech processing separates pre-lexical (abstract
100 phonological representations before processing) and lexical levels (Norris, 1994; Scharenborg et al., 2005).
101 Shortlist-B proposed a further pre-processing of the speech signal to reflect the characteristics of human
102 pre-lexical processing (Norris and McQueen, 2008). However, these ASRs still worked at a phonetic-scale
103 (i.e. with phonetic base units) and were mainly conceived for Hidden Markov Models-based acoustic units.
104 Instead, our ASR uses syllabic-scale units and different acoustic model implementations and embeds a new
105 decoding algorithm that is independent of the acoustic model used. The proposed acoustic models address
106 syllable-related dynamics inspired by multi-temporal processing studies. Our results show that these models
107 - especially two involving deep learning models - can use a limited training material to gain performance
108 that is comparable with that of state-of-the-art systems on a non-trivial recognition task. Furthermore, our

109 ASR can outperform standard-approach ASRs built upon the same training material. One drawback is its
110 higher computational complexity, which requires using parallel or distributed processing for operational
111 applications. As benchmark experiments, the recognition of spoken Italian digits (0-9) and numbers ranging
112 from 0 to 1 million (excluded) from telephone-quality recordings were used. Digit recognition was used
113 to compare ASR performance on controlled speech with a simple grammar and a low variability in the
114 utterance of the syllables. Instead, the 0-1 million number experiment was used to test ASR performance
115 with a non-trivial language model and with noisy audio that potentially included features of spontaneous
116 speech (omissions, uncertain speech, false starts, etc.). In these experiments, when our ASR used LSTM
117 syllabic acoustic models through an exhaustive decoding algorithm it had comparable performance with the
118 state-of-the-art Google speech-to-text service (Google, 2019) although it was trained with just one hour of
119 speech samples. Overall, our experiment is a preliminary approach to open the way for questioning base
120 ASR components and thus to provide a cost- and resource-effective solution to build ASRs.

121 Our results support the hypothesis that involving psycho-acoustic and supra-segmental information in
122 an ASR, through the modelling of long and short term dynamics, likely increases its performance. This is
123 an important topic impacting many different situations where general-purpose ASRs may not be applica-
124 ble. First of all, ASRs based on DNNs are challenging for low-resource languages, which may be cut-off
125 from a number of speech interfaces. Domain-specific recognition is also a challenge as it often poses strict
126 constraints to ASRs. For example, pathological speech depends on the effect that a disease may have on
127 human voice and requires strong ASR customisation. Further, domain-specific applications may use words
128 or expressions that are not modelled by general-purpose systems, and sensitive data may not be sent to third
129 parties for transcription. From the point of view of the open source community and of small enterprises,
130 it is important to have the option not to depend on large companies to include speech-to-text capabilities
131 in their applications. In general, a psycho-acoustically motivated solution provides significant adaptation
132 capabilities and flexibility.

133 Overall, the main research question addressed by this paper is: *Can cognitive and psycho-acoustic the-*
134 *ories on the syllable's role in human speech recognition inspire effective syllabic models and ASR architec-*
135 *tures?*

136 This paper is organised as follows: Section 2 describes the assumptions, the material, and the models

137 used in our ASR and alternative baseline ASRs. Section 3 reports the performance comparison between all
138 ASRs on the recognition of syllables, digits, and numbers. Section 4 discusses the results and draws the
139 conclusions.

140 **2. Material and Methods**

141 *2.1. The Base Unit of Speech*

142 The base unit of speech is the minimal form of acoustic information around which human spoken lan-
143 guage recognition is organised (Massaro, 1972). Indeed, the assumption that just one base unit exists is an
144 exemplification of automatic modelling, since linguistic studies have instead indicated that language is or-
145 ganised around a combination of units with different temporal ranges (Greenberg, 1996). Generally, human
146 speech recognition uses several units with different time-scales, each containing coherent information at a
147 given linguistic or paralinguistic level, and likely processes these units concurrently (Hawkins and Smith,
148 2001). In automatic speech recognition, the base unit of speech is usually modelled as one unit (i) having
149 a high number of manifestations, (ii) spectrally defined, and (iii) allowing the implementation of computa-
150 tionally efficient algorithms. The following subsections describe two base units commonly used in ASRs:
151 Phonemes³ and syllables.

152 *2.1.1. Phoneme*

153 Spoken language continuum is still commonly represented by a string of phonetic symbols (e.g. the
154 International Phonetic Alphabet). This representation "hides" co-articulatory transitions and partial supra-
155 segmental labelling (mainly word stress) (Ostendorf, 1999), but allows representing an entire language using
156 a large combination of few tens of symbols. The pronunciation of the string of symbols varies from person
157 to person and from word to word. Generally, a phonetic symbol in the IPA alphabet is associated to a set
158 of reference spectral frequencies (fundamental and formant frequencies) in its stationary section, and all
159 transitional and dynamic spectral variations are assumed to be at the head and the tail of this section. These
160 complex dynamics depend on the variable shape of the vocal tract and the possible activity of the vocal cords

³For brevity, in our model descriptions we will improperly use the term "phoneme" to both indicate classes of speech sounds (phonemes) and their realisations (phones).

161 during the production of the sound. The identification of a phoneme from a portion of the signal spectrum
162 requires catching the exact period in which the vocal tract has a defined structure that produces a stationary
163 signal, and capturing expected transitions toward the following speech sound. ASRs apply an iterative
164 window of ~10 ms, running on the speech signal to capture both transitions and dynamics. Consequently,
165 the identification of a speech segment as a given unit requires using statistical models that take into account
166 its spectral context and the large variability of the phoneme across phonotactic contexts and speakers. Most
167 ASRs use tri-phones as base units. However, this assumption neglects a large amount of information with
168 higher time range contained in the spoken realisation of syllables and words (Fujimura, 1975; Yule and
169 Bernini, 1997).

170 2.1.2. Syllable

171 One empirical definition of *phonetic syllable* (or *pseudo-syllable*, Martin (2010)), is reported in D'Alessandro
172 and Mertens (1995):

173 “[...] a continuous voiced segment of speech organised around one local loudness peak, and
174 possibly preceded and/or followed by voiceless segments.”

175 This definition is application-oriented and is useful in automatic segmentation processes. However, while
176 keeping the term *phonetic syllable*, we adopt a more precise definition by Roach (Roach, 2000, p. 70) that
177 better accounts for co-articulation dynamics:

178 “[...] consisting of a centre which has little or no obstruction to airflow and which sounds
179 comparatively loud; before and after that centre [...] there will be greater obstruction to airflow
180 and/or less loud sound.”

181 Thus, a syllable can also be seen as a 100-250 ms segment of signal constructed around a high energy peak
182 (*nucleus*), possibly preceded by an increasing energy slope (*onset*) and followed by a tail of decreasing
183 energy (*coda*).

184 Syllables are probably the units around which human speech production developed (MacNeilage and
185 Davis, 2000). Several studies have highlighted the importance of syllables in human speech perception,

186 because syllables can be perceived in a spoken word also when they are not actually uttered (*mirage ef-*
187 *fect*) (Fujimura, 1994; Warren et al., 1996; Arnal et al., 2016). However, in "Categories" Aristotle already
188 observed the vague nature of syllables, i.e. although each syllable is heard as separated from the other,
189 generally they do not have defined boundaries. Indeed, disfluencies and reduction processes, observed also
190 in clear but connected speech, can cause segment cancellation and indeterminacy of clear syllabic bound-
191 aries (Greenberg, 1999). However, the rhythmical structure is always preserved, which means that more
192 prominent units are usually less reduced and thus guarantee the preservation of speech-chain intelligibility
193 (Cutugno et al., 2012). From a psycho-acoustic point of view, a syllable contains much more information
194 than the sequence of sounds constituting it (Wu et al., 1998; Kahn, 2015). Different brain activation pat-
195 terns have been observed in human subjects hearing sequences of syllables or single syllables alternatively
196 (Peeva et al., 2010; Rong et al., 2018). These experiments have also highlighted specific activation patterns
197 in different brain areas, corresponding to multiple temporal scales of phonetic, syllabic, and supra-syllabic
198 lengths.

199 One drawback of using syllable as a base unit in ASR, is that it is difficult to numerically describe all
200 syllabic-scale ($\sim 100 - 250$ ms) speech properties that psycho-acoustic studies have indicated as related
201 with human speech recognition robustness to speaker differences and adverse environmental conditions
202 (Kingsbury et al., 1998; Greenberg, 1996). Also, syllable boundaries lack a consistent psycho-acoustic
203 and linguistic definition that makes acoustic model specification non-uniquely defined (Wu et al., 1998;
204 Huang et al., 2001). Indeed, most syllabic-scale features are related with prosody, energy contour, and
205 slow modulations (around 4 Hz). Several studies have demonstrated that incorporating this information in
206 syllabic acoustic models can increase the performance of an ASR (Cutugno et al., 2005; Coro, 2008; Baby
207 and Hamme, 2015; Batliner and Möbius, 2019). However, these studies have also highlighted that it is not
208 convenient to re-use standard ASR algorithms and assumptions when using syllabic base units (Wu et al.,
209 1998; Chang, 2002; Pinson and Pinson, 2019). Generally, most syllabic ASRs either represent syllables as
210 sequences of phonetic acoustic models or build one acoustic model per syllable (or demi-syllable) while
211 using phonetic features extracted from ~ 10 ms signal windows.

212 From a speech-processing operational point of view, pseudo-syllables bring more advantages than tri-
213 phones. The automatic segmentation of a speech signal into pseudo-syllabic segments is facilitated by the

214 correlation of these units with the modulations of sonority movements. In every group of sounds, there
215 are as many syllables as clear relative peaks of sonority (Jespersen, 1905). Furthermore, speech-intensity
216 change is correlated with tonal speech perception (House, 1996) and is typically maximum between pseudo-
217 syllables' onsets and nuclei. These properties allow detecting tonal units automatically (Cutugno et al.,
218 2002; D'Anna and Cutugno, 2003; D'Anna and Petrillo, 2003). Moreover, the acoustic correlates of pseudo-
219 syllables can be used to model pitch movements and produce effective pitch contour stylisation, especially
220 over signal segments with complex spectral content (Origlia et al., 2013). These characteristics allow to build
221 automatic pseudo-syllable classification and segmentation algorithms based on sound-intensity and spectral
222 entropy analysis, overcoming common issues related with sonorant consonants with high intensity (e.g.
223 nasals sounds) (Origlia and Cutugno, 2016). Moreover, automatic emotion detection and tracking models
224 can be more efficient using pseudo-syllabic-scale analyses, instead of phonetic-scale analyses, by harnessing
225 nuclei's spectral richness and extracting information on speech rate and style (Origlia et al., 2014).

226 Overall, pseudo-syllables are (i) widely used in psycho-acoustic studies as the base unit of analysis,
227 (ii) correlated with observable neural activity, (iii) phonetically describable through specific intensity and
228 spectral patterns, (iv) automatically detectable by computationally efficient algorithms, and (v) based on
229 clear phonetic templates. Thus, using pseudo-syllables in ASR allows including human-related patterns and
230 automatically identifying the signal segments that should be extracted and annotated to train the acoustic
231 models and define word pronunciation models. The Italian part of the corpus used in this paper (Section 2.2)
232 was also annotated at the pseudo-syllable level to foster experiments that could explore these operational
233 advantages. For all these reasons, in this paper acoustic models are based on pseudo-syllables, although the
234 term *syllabic model* is used for simplicity.

235 2.2. *Lexicon, Language Model and Acoustic Features*

236 The experiment reported in this paper uses a controlled lexicon to test the performance and the properties
237 of different acoustic models and decoding algorithms. This lexicon was selected to require short preparation,
238 analysis, and development times, and also to produce a non-trivial language model that included sufficiently
239 varied speech and some characteristics of spontaneous and large-vocabulary contexts. Based on these re-
240 quirements and following the suggestions of other works (Wu et al., 1998; Chang, 2002), the range of

241 numbers between 0 and 999,999 - hereinafter indicated as [0,1M) - was selected as a benchmark vocabulary.
242 The Speecon Italian corpus (Siemund et al., 2000; ELRA, 2019) includes sentences in this numeric range,
243 recorded at 16kHz from 400 different speakers with telephonic quality (with $25dB \pm 3dB$ signal-to-noise
244 ratio). Moreover, Speecon includes annotations for numbers and digits (i.e. from 0 to 9) at the phonetic,
245 syllabic, and sentence levels. A total of 42 syllables (Table 1) and 19 phonemes - plus two silence models
246 - and ~220 syllabic combinations were sufficient to build up a language model for the lexicon of numbers
247 in [0,1M). Although the number of phonemes involved is not far from that of the whole Italian language
248 (~32), 42 syllables are just a subset of the thousands of syllables of Italian. However, even with these sylla-
249 bles, spoken long numbers present characteristics of spontaneous speech, e.g. omissions, false starts, dialect
250 inflexions, and uncertain speech.

251 In our experiment, the Speecon recordings were divided into 80-20% training and test sets within a cross-
252 validation process and did not include the same speakers. The Speecon syllabic-level annotations allowed
253 to extract ~65 minutes of speech to train acoustic models and ~13 minutes to test their performance. Word-
254 level annotations allowed to prepare recordings to test ASRs' performance on numbers (~140 minutes) and
255 digits (~55 minutes). The language model for [0,1M) was built using the CMUCLMTK toolkit v7 (CMU,
256 2019) as a statistical model trained with syllabic mono-grams, bi-grams, and tri-grams, and had a non-trivial
257 perplexity of 9.8. As a training set for the language model, the linguistic syllabic subdivisions of all numbers
258 in [0,1M) were used, plus their syllabic transcriptions in Speecon. These transcriptions report the syllables
259 actually uttered in long numbers and thus simulate spoken sentence alterations due to continuous speech,
260 which in turn allows building a more realistic syllabic language model. Finally, back-off probabilities were
261 used to account for non-observed syllabic concatenations.

262 As acoustic features, 13 Mel-frequency cepstral coefficients (MFCCs, Davis and Mermelstein (1980))
263 were used, with delta and double-delta features, for a total number of 39 features extracted from ~10 ms
264 windows sliding over the signal with 50% overlap (5 ms). MFCCs are standard features used in ASR (Sahu
265 et al., 2018). They are extracted out of the application of a filter bank based on the mel scale, which simulates
266 the response of the human auditory system to speech frequencies. Although other types of mel scale-based
267 features could be used (Tyagi and Wellekens, 2005; Parcollet et al., 2018; Kim et al., 2019), MFCCs were
268 the set of features that all ASRs involved in our experiments could use. Thus, MFCCs allowed to measure

269 performance differences that depended on the architectures and acoustic models rather than on the signal
270 representation. Furthermore, MFCCs typically allow to use a lower number of spectral features (typically
271 13 per window) than alternative methods (e.g. Mel-filter bank energies, which typically require 40 features
272 per window) (Paliwal, 1999), which was beneficial to avoid overfitting issues with our limited training set.

273 Although our proposed acoustic models were syllabic (i.e. our ASR used syllables as base units), they
274 were committed to extracting syllabic information from sequences of phonetic-scale features. Indeed, alter-
275 native features using syllabic-scale windows directly (100-250 ms) do not have the same consistency and
276 robustness as MFCCs for speech recognition (Kingsbury et al., 1998; Tyagi et al., 2003; Baby and Hamme,
277 2015).

278 2.3. *Speech Decoding*

279 The decoding process used by our reference ASRs (e.g. Token Passing, Young et al. (1989)) relies on the
280 alignment between sequences of HMM states and the audio signal and makes use of the initial and final states
281 to estimate phonetic and syllabic boundaries. However, some of the syllabic acoustic models proposed in this
282 paper are not made up of sequences of states, thus decoders conceived to work with HMMs could not be used.
283 For this reason, a new decoding algorithm was used that was independent of the nature of the incorporated
284 syllabic acoustic model used. The algorithm described in Coro et al. (2007) (hereinafter named *exhaustive*
285 *Viterbi*) fitted our scopes because it uses syllabic acoustic models as black-boxes. This algorithm optimises
286 the alignment of each acoustic model to the signal because it calculates the likelihoods of the models to
287 all possible sub-sequences of the acoustic features extracted from the audio signal, i.e. it tests all possible
288 alignments of the models to the signal. This approach increases the performance also of standard syllabic
289 HMMs with respect to other decoding algorithms and has also been used in a commercial ASR (D’Anna
290 et al., 2009). In particular, the algorithm calculates the conditional probability distribution $P(W|X)$ of a
291 sequence of n syllables $W = w_1w_2..w_n$ given a sequence of T features $X = x_1x_2..x_n$ extracted from the
292 audio signal, where T is the length of the audio signal. During the calculation, the algorithm combines
293 a syllable-based language model with syllabic acoustic models to find the optimal sequence of syllables
294 W^* associated to X . The output of the algorithm is thus the sequence of syllables that is most probably
295 associated with the audio signal. Through the pronunciation models it is possible to associate lexicon words

296 to the optimal sequence of syllables and produce the orthographic transcription of the audio.

297 Formally, the algorithm produces the following optimal solution (the demonstration is reported in Coro
298 et al. (2007)):

$$P(W^*|X) = \operatorname{argmax}_{m \in Syl} \{f(m, T) \cdot E(m)\}$$

299 where Syl is the complete set of N syllables included in the language model, $E(m)$ is the probability
300 of model m to be an ending syllable, and $f(m, t)$ is the solution to the sub-problem of unit alignment in the
301 time interval $[1, t]$, defined as

$$f(m, t) = \max \left\{ \begin{array}{l} P(X_1^t | m) \cdot \pi(m) \\ \max_{1 \leq t^* < t, n \in Syl} \{f(n, t^*) \cdot P(m|n)^\gamma \cdot P(X_{t^*+1}^t | m)\} \end{array} \right\}$$

302 where γ is the language model's weight. Starting from time T , a backtracking process follows the
303 definition of $P(W^*|X)$ to find the best alignment between the models and the signal. In particular, ac-
304 cording to the definition of $f(m, t)$, the algorithm efficiently tests all possible alignments of all models
305 to all segments of the audio signal, and thus optimises the models' recognition accuracy. However, one
306 drawback is that it requires a pre-calculation of all models' likelihoods to all subsets of observations, i.e.
307 $P(X_{t_i}^{t_j} | m) \forall m \in Syl, 0 \leq t_i \leq T, 0 \leq t_j \leq T$. In particular, the algorithm first computes the V matrix:

$$V = \begin{pmatrix} P(X_1^1 | m) & P(X_1^2 | m) & \dots & P(X_1^{T-1} | m) & P(X_1^T | m) \\ 0 & P(X_2^2 | m) & \dots & P(X_2^{T-1} | m) & P(X_2^T | m) \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 0 & P(X_T^T | m) \end{pmatrix}$$

308 and then the backtracking procedure rapidly reconstructs the optimal solution. The overall algorithm's
309 complexity is $O(T^2 N^2 C_l)$, i.e. it is quadratic in T and N , and also depends on the complexity of all
310 likelihood calculations C_l by the acoustic models. The complexity of the algorithm can be reduced by intro-
311 ducing constraints on the minimum and maximum $t_j - t_i$ difference, and through a beam search strategy in
312 the $f(m, t)$ calculations that filters out all likelihoods falling under a certain threshold. This strategy strongly

313 reduces the number of non-zero elements in the V matrix and thus reduces computational time. In our exper-
 314 iment, forcing $100\text{ ms} \leq t_j - t_i \leq 250\text{ ms}$ and relative likelihood ≥ 0.5 made our ASR return results in short
 315 time without losing performance. Further, since each element of the V matrix is independent of the other,
 316 the matrix calculation can be parallelised to (linearly) reduce decoding time (Section 4). The main differ-
 317 ence between the exhaustive Viterbi algorithm and Token Passing is that the former treats acoustic models as
 318 black boxes. Indeed, exhaustive Viterbi is independent of the acoustic model implementation used and only
 319 requires likelihood calculations from these models. Instead, Token Passing is strongly based on the assump-
 320 tion that acoustic models are made up of sequences of states and that the transition from one model to another
 321 can occur only from a final state to an initial state. This assumption strongly reduces the computational com-
 322 plexity of the decoding strategy. In particular, Token Passing defines the *minimum alignment cost*, between
 323 the vectors X_1^t and a sequence of model states ending in state j , as $s_j(t) = \min_{i \in \text{all states}} \{s_i(t-1) + p_{ij}\} + d_{ij}$.
 324 With p_{ij} being a transition cost that is given either by the model state-transition matrix or by the language
 325 model (when j is an initial state and $\{i\}$ are the final states of other models). The optimal sequence of states
 326 is the one having the minimum cost $S = \min_j \{s_j(T)\}$. Using a bi-gram language model, the computational
 327 complexity of the unit-to-signal alignment is between $O(TN \log(N)C_l)$ and $O(TN^2C_l)$, where N is the
 328 number of connected units (e.g. syllables) and C_l is the complexity of the likelihood calculation of the state-
 329 based model used. Instead, the models managed by exhaustive Viterbi can be non-state-based, which is the
 330 main reason for its higher computational complexity but also for its higher flexibility.

331 2.4. Hidden Markov Models

332 Hidden Markov Models (HMMs) are the most used choice for acoustic modelling (Figure 1-a). Given a
 333 sequence of acoustic features X , they estimate the conditional probability distribution $P(X|S)$ of X given a
 334 sequence of states $S = s_1, s_2, \dots, s_T$. Based on this definition, the Viterbi algorithm (Viterbi, 1967) efficiently
 335 estimates the likelihood of an HMM to X as the conditional probability $P(X|S^*)$ of the sequence of states
 336 S^* that maximises $P(X|S)$ (i.e. the one most likely associated to X). An HMM that models a phoneme is
 337 trained (e.g. through the Baum-Welch algorithm) on many examples of acoustic features for that phoneme,
 338 in order to model the inter-speaker and inter-word variability of that phoneme. Likewise, a syllabic HMM
 339 is trained on the acoustic features of a syllable (Figure 2-a), i.e. concatenations of phonetic-scale features.

340 Often, ASRs use concatenations of phonetic HMMs to build up di-phone or tri-phone models that are re-
341 trained to better assess inter-model transitions (CMUSphinx, 2017). Modern ASRs need thousands of hours
342 of annotated material to train phonetic acoustic models for large-vocabulary applications. For the experiment
343 reported in this paper, HMM implementations from JAHMM (Francois, 2019), KALDI (Povey et al., 2011),
344 and CMUSphinx (Lamere et al., 2003) were used to implement syllabic and tri-phonetic HMMs.

345 For decades, ASRs have used HMMs with Gaussian mixtures (GMMs) to model emission densities
346 (Huang et al., 2001). However, current state-of-the-art ASRs use DNNs to model emission densities (Figure
347 1-b) as $\text{softmax}(o_i(X_t))$, where $o_i(X_t)$ is the value of the activation function in the output layer of the
348 node corresponding to state i (Yu and Deng, 2015). This type of ASR is currently used in many domains and
349 reaches state-of-the-art performance (Serizel and Giuliani, 2017; Ravanelli and Omologo, 2017; Maas et al.,
350 2017; Novoa et al., 2018; Patel et al., 2018; Smit et al., 2018; Chao et al., 2019; Mao et al., 2019). The used
351 DNNs are typically multi-layer perceptrons with many layers (~ 7), with the training phase initialised by a
352 pre-training algorithm. KALDI provides two main implementations of HMM-DNNs, one using Restricted
353 Boltzmann Machines for pre-training and Stochastic Gradient Descent for training (HMM-DNN-nnet1),
354 and the other one using Natural Gradient for Stochastic Gradient Descent and Parameter Averaging (HMM-
355 DNN-nnet2).

356 2.5. Factorial Hidden Markov Models

357 An FHMM is made up of a set of HMMs, all with the same number of states, and inter-dependent
358 emission probabilities usually modelled as multi-variate Gaussians (Figure 2-b) (Logan and Moreno, 1998).
359 FHMMs are particularly suited for speech processing, in particular to model concurrent and overlapping
360 dynamics that are generated by multiple and loosely-coupled processes, as those present in a speech signal
361 (Ghahramani and Jordan, 1996; Gael et al., 2009; Florian et al., 2011). Multi-temporal ASRs have used this
362 property to model the syllabic and phonetic structures contained in ~ 200 ms speech segments. In particular,
363 the transition probability distributions of syllabic FHMM acoustic models with 2 parallel HMMs have high-
364 lighted the presence of two inter-linked syllabic-scale and phonetic-scale dynamics (Coro, 2008). These are
365 likely responsible for the higher performance of FHMMs with respect to HMMs in syllable modelling. FH-
366 MMs have been used in ASRs with the aim to include results from psycho-acoustic studies on overlapping

367 speech dynamics (Logan and Moreno, 1998; Virtanen, 2006; Tu et al., 2016). For the experiment reported in
368 this paper, FHMMs were implemented in Java by porting and optimising the original Matlab implementation
369 by Z. Ghahramani (Ghahramani, 2002).

370 2.6. Deep Learning models

371 Deep Learning (DL) models leverage a multi-layered structure to extract information from raw input
372 data. DNNs are conventional DL models where each layer of the network is assumed to produce an internal
373 representation of the input (*feature map*), with deeper layers producing higher levels of information abstrac-
374 tion. The increasing computational power of graphic processing units (GPUs) has allowed introducing DL
375 models in a vast number of domains, e.g. from computer vision (Krizhevsky et al., 2012; Massoli et al.,
376 2020; Girshick, 2015) to natural language processing (Deng and Liu, 2018; Ortis et al., 2019). In our ex-
377 perimental campaign, two different DL models were implemented - with PyTorch (Paszke et al., 2019) - as
378 acoustic syllabic models (without embedding them in an HMM): a Convolutional Neural Network (CNN,
379 LeCun et al. (1995)) and a Long Short-Term Memory model (LSTM, Hochreiter and Schmidhuber (1997)).

380 2.6.1. Convolutional Neural Network

381 DNNs (which include CNNs) process signal segments in a “static” way, i.e. like they were images (Coro,
382 2004; Coro et al., 2019). Normally, they do not model time as an internal parameter and this limitation neg-
383 atively affects their performance in automatic speech recognition with respect to other time-explicit models.
384 In order to account for this issue, a CNN was built to model pseudo-syllables using a multi-temporal anal-
385 ysis within a convolutional stage (Figure 2-c). This model uses the following operations: Each unit of the
386 convolutional layer is computed by means of multiplications between the input data and a matrix (*kernel*),
387 whose optimal values and size were assessed during the training phase. The kernel size corresponds to the
388 size of the input that is convolved with the kernel (*receptive field*) so that each convolution only looks at a
389 small portion of the input. Through the use of small receptive fields, convolutional layers are generally able
390 to extract and combine *local* information from the input data, i.e. information contained in segments (of
391 speech signal, in our case) with a predefined length. Our CNN used four 1D convolutional layers - each with
392 a different kernel size - that corresponded to different filters and windows on the syllabic signal. The size of
393 each window represents the time scale processed by each convolutional operation. During the convolution,

394 a window stride of one sample maximises the capture of local relations through the signal. In summary, our
395 CNN analyses a syllabic speech signal at multiple time scales through a multi-window processing. After
396 the convolutional step, vector pooling and stacking operations are followed by a flattening operation that
397 projected all the resulting feature vectors (feature maps) on a new 1D vector, whose optimal length was esti-
398 mated during the training phase. This vector is input to a fully-connected (FC) neural network layer, whose
399 optimal size was estimated during the training phase. A rectified linear unit (ReLU) activation function
400 ($\max(0, x)$) is applied to each node of this layer to reduce the vanishing gradient problem (Bengio et al.,
401 1994; Kapur, 2020) and favour generalisation capability (Mishkin et al., 2017; Novak et al., 2018). In order
402 to reduce the risk of data overfitting due to the small amount of training data available, the dropout technique
403 (Srivastava et al., 2014) was used on the FC layer. Dropout statistically excludes some nodes of the FC layer
404 from one training session and re-introduces these nodes with their original weights after the non-dropped
405 nodes connections have been trained. At each training step, a new set of nodes is selected to be dropped.
406 Finally, during the inference phase, each node’s output is multiplied by a dropout probability to account
407 for their possibly missed training steps. Overall, this procedure simulates an ensemble of a high number
408 of different models whose output is eventually averaged at inference time. The last stage of our CNN is a
409 classification layer, i.e. another FC layer with 44 neurons, one for each unit to recognise. This layer allows
410 classifying the acoustic features of a syllabic signal as one among the syllabic units reported in Table 1.
411 Indeed, a softmax function applied to the layer’s outputs makes the CNN overall simulate a posterior prob-
412 ability density $P(W|X)$ of each syllable W given the input vector X (Muller, 2014). In turn, this reduces
413 the complexity of the decoding algorithm, because the probabilities of all syllables for a signal segment are
414 calculated just after one propagation of the input through the network. The described architecture came after
415 testing a large number of alternative architectures, including chained windows and deeper networks. It was
416 the architecture using the lowest number of parameters and gaining the highest performance on the tasks
417 reported in this paper.

418 2.7. Long Short-Term Memory Model

419 LSTMs are naturally suited to process observation sequences and time series (Hochreiter and Schmidhu-
420 ber, 1997), because they consist of one computational unit that is iteratively used to process the observations

421 of an input time series (processing *steps*). The unit uses *gating* mechanisms that process a temporal flow
422 of data while controlling the retain and the release of memorised information. Since an LSTM is suited to
423 simulate a posterior probability density $P(W|X)$ of all acoustic units W given the input features series X ,
424 it cannot directly replace an HMM (which calculates likelihood). Furthermore, the processing steps cannot
425 be treated as the sequence of states of an acoustic HMM and thus cannot be used in classic speech decoding
426 processes.

427 In this paper, an LSTM model was implemented to classify pseudo-syllabic acoustic units directly and
428 was later combined with a speech decoding algorithm, which was able to harness its multi-temporal pro-
429 cessing of the speech signal. Our LSTM model’s unit processes one vector of the input time series (i.e. one
430 window of acoustic features) at a time. It uses a standard unit characterised by one *forget gate*, one *input*
431 *gate*, and one *output gate* (Figure 2-d), all implemented as single-layered neural networks. Within the unit,
432 the cell state c_t is a Real-valued vector that roughly stores the “long-term” memory of the model, whereas the
433 hidden state h_t is the output vector of the LSTM unit that manages “short-term” memory. At each processing
434 step, the LSTM unit receives the current input vector of acoustic features, and the cell and hidden states of
435 the previous unit. The unit outputs a new cell state and a new hidden state. All gates receive the current
436 unit input vector and the previous hidden state as input. As a first operation, the cell state of the previous
437 processing step is multiplied by the output of the *forget gate*, i.e. a neural network with sigmoid activation
438 function with range $[0,1]$, where 0 represents a complete blockading (forget) of an input element and 1 a
439 complete pass (remember). Another process point-wise multiplies the output of a sigmoid-activated neural
440 network (*input gate*) by the output (*proposed cell state*) of a tanh-activated neural network. The output of
441 this process is summed to the output of the forget gate in order to establish which part of the information re-
442 tained by the forget gate should be updated. This result is the unit’s cell state that is passed to the next LSTM
443 processing step. The hidden state is calculated by first passing the cell state to a tanh function (to re-scale its
444 values in $[-1,1]$) and then multiplying this result with the output of another sigmoid-activated neural network
445 (*output gate*). Overall, this final step roughly decides what portion of the cell state is produced as the output
446 of the LSTM unit. As a final step, our LSTM-based syllabic acoustic model uses the last processing step’s
447 output as an input to a classification layer, whose input size is equal to the hidden state size. Similarly to
448 the CNN model, a softmax function is applied to the output of this classification layer in order to simulate

449 the posterior probability density $P(W|X)$ of each syllable W given the input vector X , and to reduce the
450 decoding algorithm complexity.

451 2.8. Baseline Speech Recognisers

452 Several instances of our ASR were produced to assess its performance depending on the acoustic model
453 used. Each instance used one among the four supported models. It is worth noting that some of these models
454 are more suited for recognising entire syllables, but their performance could be positively or negatively
455 affected by the combination with the decoding algorithm (Section 3.3). In particular, syllabic HMMs were
456 enabled in our ASR architecture to measure the performance enhancement that our decoding algorithm
457 would bring to a classic model. Similarly, FHHMs were used to evaluate the performance in word and
458 sentence recognition of a naturally suited model for single pseudo-syllable recognition (Coro et al., 2007).
459 Finally, the CNN and LSTM acoustic models were used to evaluate the performance gained by our psycho-
460 acoustic inspired models in word and sentence recognition.

461 CMUSphinx (Lamere et al., 2003) and KALDI (Povey et al., 2011) were used as reference ASRs. These
462 systems use tri-phonetic HMM-GMMs and HMM-DNNs respectively within a reference ASR architecture,
463 and were trained with an open source reference corpus suited for our recognition tasks. The Italian VoxForge
464 corpus (VoxForge, 2012) was used to train phonetic and tri-phonetic HMMs with 19 hours of speech that
465 involved regional inflexions. Although the dimension of VoxForge is generally not sufficient to build a
466 high-performance large-vocabulary ASR, it was sufficient to build high-performance baseline ASRs for
467 spoken digits and numbers. In particular, CMUSphinx and KALDI were trained through the following
468 operations: (i) Pronunciation models for words uttered in the VoxForge recordings were taken from the
469 large database of Cosi (2015); (ii) phonetic acoustic models (and tri-phones) from the pronunciation models
470 were aligned to the recordings through automatic alignment processes; (iii) the language models described in
471 Section 2.2 were integrated to produce two recognisers, one for digit recognition and another one for number
472 recognition. Furthermore, in the single-syllable recognition task (Section 3.2) the phonetic transcriptions
473 from Cosi (2015) were used to model the pronunciations of the 42 syllables of Table 1.

474 As a second baseline ASR, the Google Speech-to-Text cloud service was used (Google, 2019). This
475 HMM-DNN based ASR, trained with thousands of hours of speech, has top-level performance and high re-

476 sponse efficiency, and is used by almost all Google technology. Google constantly improves its performance
477 after periodically collecting users' data and revising acoustic, pronunciation, and language models. This
478 ASR uses phonetic transcriptions of 10 times the words of an entire language dictionary, for each of the 120
479 languages supported, and also includes a context-specific adaptation process that is able to resize the gram-
480 mar and to optimise the transcription according the language context (Peters et al., 2011; Ballinger et al.,
481 2011; Google, 2019). For example, on number recognition tasks the Google service is able to report the
482 numerical form of the uttered number (e.g. "one hundred three" is reported as 103), while insertions, false
483 starts, and other non-numerical words are deleted if not uttered clearly (Google, 2019). Google Speech-to-
484 Text (Dec. 2019 version) was used as a reference state-of-the-art ASR. Indeed, comparing the performance
485 of our method with that of the Google ASR may not be an optimal choice, because of the different train-
486 ing corpora used (voices, data size, data preparation, etc.) and the lack of details about the Google ASR's
487 architecture. Nevertheless, the Google's context-specific adaptation feature and the relatively small testing
488 context (numbers) reasonably allow to use the comparison as a proxy for a quality assessment of our ASR.

489 **3. Results**

490 This section reports a performance comparison between the models described in the previous section.
491 Accuracy is used as a comparison metric, defined as

$$Accuracy = \frac{\textit{number of correctly recognised units} - \textit{number of over - inserted units}}{\textit{total number of units in the manual transcription}}.$$

492 In order to make comparisons consistent, the interpretation of "unit" in the accuracy formula changed
493 according to the test case. In fact, the compared ASRs had heterogeneous architectures and used different
494 speech units and output types. For example, the Google ASR reported the entire recognised sentence with
495 numerical symbols (e.g. "1" for a digit and "1350" for a number). Furthermore, the other ASRs used either
496 tri-phones or pseudo-syllables. In this context, a comparison could be consistent only at a final orthographic
497 transcription level. Thus, in the single-syllable recognition task, accuracy was calculated on the number of
498 correctly transcribed orthographic syllables. On digit and number recognition tasks, accuracy was calculated

499 on the orthographic transcription of the entire sentence.

500 *3.1. Acoustic Model Topologies*

501 The machine-learning models reported in the Section 2, were trained to recognise the 44 units reported
502 in Table 1. Multiple parametrisations and implementations of the models were tested. Eventually, the
503 topologies and implementations gaining the highest performance were selected for the comparison. This
504 operation required testing thousands of parameter combinations.

505 Optimal HMMs for GMM-based tri-phone models were produced with KALDI and used 5 states and
506 32 mixtures, whereas HMM-DNN phonetic models used 7 hidden layers in the DNN. Out of these HMMs,
507 syllables were represented as concatenations of phonetic HMMs. Optimal syllabic HMMs were produced
508 with JAHMM and had 7 states and 39 Gaussian mixtures. Optimal FHMMs used 2 HMMs with 7 states
509 each and one multivariate Gaussian emission density.

510 Regarding the deep learning models, cross-validation was used to find optimal parameters and topologies
511 of the CNN and the LSTM acoustic models. Specifically, the optimal CNN topology used windows of 48,
512 80, 96 and 112 ms to create 64-length feature maps after convolution, and was made up of two FC layers
513 (one hidden layer and one output layer). The feature map was optimally flattened to 1,280 elements, the
514 dropout probability was 20%, and the optimal size of the first FC layer was 300. The final FC classification
515 layer had 44 neurons, one for each syllable to recognise. Moreover, the importance of introducing non-
516 linearity in the output of this FC layer through ReLU was tested: All models were also trained without
517 ReLU, and a performance degradation up to 4% was observed, which confirms the positive contribution of
518 this transformation and the need to include it in the acoustic model.

519 The optimal LSTM was a mono-directional model with one hidden state with 1000 neurons and a final
520 classification layer with 44 neurons, one for each syllable to recognise. The training phases of both models
521 used the Adam optimiser (Kingma and Ba, 2014) with cross entropy loss criterion and a learning rate of
522 $1.e^{-3}$, reduced of 5 times whenever loss reached a plateau.

523 As for ASR configuration, the language model of CMUSphinx was trained with mono-grams, bi-grams,
524 and tri-grams. This ASR used HMM phonetic models with 32 Gaussian mixtures, trained with 19 hours of
525 annotated recordings from the VoxForge corpus. Our ASR used the exhaustive Viterbi algorithm described in

526 Section 2.3 alternatively combined with HMM, FMM, CNN, and LSTM syllabic acoustic models. The ASR
527 used a language model based on the bi-grams prepared for CMUSphinx. Finally, the Google Speech-to-Text
528 service was used through a Java client that streamed audio files and collected transcriptions.

529 *3.2. Syllable Recognition*

530 Acoustic models' performance was first compared on the recognition of the syllables and silence units of
531 Table 1, without the interference of the language model and the decoding process. This performance com-
532 parison (Table 2-a) showed that our LSTM outperforms phonetic HMMs by 8.91% absolute accuracy and
533 the second optimal model (HMM-DNN-nnet2) by 2.63%. A Chi-squared test confirmed that this discrep-
534 ancy was highly significant with our test set size (with p-value of non-significant discrepancy null hypothesis
535 lower than 0.0001) (NIST, 2018). The accuracy of our LSTM increased non-linearly with the vector length
536 of the LSTM hidden state (Figure 3), which indicated that a ~1000 length was really required to model the
537 complexity and variability of the syllables. Interestingly, our multi-temporal CNN had comparable perfor-
538 mance with HMM-DNN models, and the HMM-DNN model using Natural Gradient had a slightly higher
539 performance than the other HMM-DNN implementation. FHMMs and syllabic HMMs outperformed pho-
540 netic HMMs, in agreement with other studies (Logan and Moreno, 1998; Coro et al., 2007; Gael et al., 2009),
541 but had lower performance than the deep learning models. Since HMM-DNN-nnet2 was the second optimal
542 model, it was selected to be used in the KALDI ASR for the comparison on digit and number recognition
543 tasks.

544 *3.3. Digit Recognition*

545 Although digits are made up of a maximum of two non-silence syllables, spoken digit recognition in-
546 volves the issue of aligning sequences of silence models, short pauses, and (one or two) syllables to the
547 signal. Thus, a performance comparison on digit recognition highlighted how much the slight misalignment
548 of syllabic models to the uttered syllables influenced word recognition. In this case, accuracy was calculated
549 on the recognition of entire words directly, especially for a fair comparison with the Google Speech-to-Text
550 service. The context-specific adaptation of the Google ASR made the reported comparison meaningful be-
551 cause it restricted the grammar to the particular task and deletes non-numeric insertions that were not loudly
552 uttered.

553 Also, due to the moderately-high signal-to-noise ratio, no model reached 100% performance on this
554 "simple" task (Table 2-b). The exhaustive Viterbi algorithm optimised the alignment of the various acoustic
555 models to the speech signal and made syllabic HMMs and FHMMs outperform a standard-approach ASR.
556 Performance was generally high for all ASRs, but the accuracy discrepancy between the CNN and the
557 LSTM models (4%) was higher than in the syllable recognition case (2.76%). Indeed, the CNN model was
558 sensitive to syllable alterations (e.g. stretching and reduction) since it processed signal segments as they
559 were static images. Generally, the difference between the CNN and the LSTM models in accounting for
560 syllabic alterations is more and more evident as long as speech tends to be continuous and spontaneous.
561 For example, alterations of "kwa ttro" as "kwa tro" and of "o tto" as "o to" are more probable within long
562 numbers but also exist with digits. Overall, the high performance of the LSTM-based ASR indicated that our
563 LSTM was a suitable acoustic model for an ASR, and the recogniser also slightly outperformed the Google
564 service (98% vs 97.5% accuracy). A Chi-squared test confirmed that this discrepancy was significant (with
565 p-value of non-significant discrepancy null hypothesis lower than 0.05). Finally, the KALDI ASR using
566 the best tri-phone models (HMM-DNN-nnet2) gained 1.1% higher relative accuracy than the CNN-based
567 model, but lower relative accuracy than the Google ASR (2.5%) and the LSTM-based ASR (3%).

568 *3.4. Number Recognition*

569 The performance comparison on the recognition of [0,1M) numbers further highlighted the differences
570 between the ASRs (Table 2-c). The main difference with respect to the digit recognition case was the higher
571 performance of CMUSphinx with respect to the CNN-, FHMM-, and syllabic HMM-ASRs. This enhance-
572 ment was due to the higher amount of training material used to build the CMUSphinx ASR, and also to
573 the lower flexibility of the other models to work on the more continuous and spontaneous speech of the
574 uttered numbers, which presents a large variability due to omissions, false starts, and uncertain speech. In
575 this context, the acoustic syllabic structures can be very different from those of the training set. Models
576 like CNN, FHMMs, and Syllabic HMMs would require more training material to handle this structural vari-
577 ability. In particular, our CNN model was re-adapted from image processing and does not fully capture the
578 unfolding of information in time and its variability across the training set. Instead, the phonetic CMUSphinx
579 and the KALDI ASRs had comparable performance with Google (3.4% and 0.7% relative accuracy, respec-

580 tively) due to a training material suited for the task. In particular, KALDI demonstrated the high quality and
581 performance that HMM-DNN-based ASRs can reach.

582 Interestingly, when our ASR used the LSTM acoustic model, it outperformed all other ASRs. Indeed, our
583 ASR had a 3.7% higher relative accuracy than the Google ASR and a 7% higher accuracy than CMUSphinx.
584 A Chi-squared test confirmed that these discrepancies were significant (with p-value of non-significant dis-
585 crepancy null hypothesis lower than 0.001). The generalisation capability of the LSTM-based ASR and its
586 flexibility to account for syllable alterations was impressive. The LSTM had optimal performance in all
587 presented cases, and the acoustic models used only ~1 hour of training material, which was much lower than
588 the 19 hours used for CMUSphinx and KALDI and the thousands of Google.

589 *3.5. Issues with large-vocabulary speech recognisers*

590 A one-million-number sentence-set was used instead of a large vocabulary because building a large
591 vocabulary speech recogniser (LVSR) requires solving other additional research questions that were out of
592 our scope. Generally, it is nearly impossible to build a state-of-the-art LVSR for a low-resource language
593 like Italian using publicly available corpora for acoustic and language model training. To better highlight
594 this aspect (and also produce a reference for our future studies), we trained and compared several LVSRs -
595 based on KALDI and CMUSphinx - using alternative open (or low-cost) textual and audio corpora (Table
596 3). The aim of this comparison was principally to highlight some intrinsic practical difficulties in building
597 LVSRs.

598 We compared recognition performance on a 15-minute corpus extracted from the Italian VoxForge corpus
599 (VoxForge, 2012) that was not used during ASR training. The textual corpora used for (4-gram) language
600 model training included: (i) the "Italian Web corpus" (itWaC), made up of texts collected from the Internet
601 and including 1.5 billion words (Baroni et al., 2009); (ii) Paisà, a large and expert-revised collection of Italian
602 texts from the Internet containing ~250 million tokens (Lyding et al., 2014); (iii) CLEF, a large collection
603 of Italian national newspaper articles from the 90's containing ~1 million words overall (CLEF, 2020); and
604 (iii) the Italian Content Annotation Bank (I-CAB), which contains 525 local (Trento province) newspaper
605 articles with ~180,000 words overall (Magnini et al., 2006). The audio corpora used were VoxForge (~20
606 hours) and APASCI (~2 hours), whose audio was based on the same spoken text. The performance across

607 multiple textual corpora was reported only for CMUSphinx for simplicity to highlight performance decrease
608 across the corpora. The following difficulties emerged, which depended on the lack of great effort (and
609 money) investment in data collection and cleaning:

- 610 1. The performance gap between the Google ASR and the other LVSRs was very high (from -25.46% to
611 -49.75% word accuracy);
- 612 2. Using uncontrolled large textual corpora (e.g. itWaC) may end in lower performance because text
613 from social networks introduces too much noise in the language model and is unsuited for spoken
614 dialogues;
- 615 3. Combining different textual corpora (e.g. Paisà+CLEF) may end in lower performance because of too
616 different language structures (e.g. Internet v.s. newspapers);
- 617 4. Data cleaning included in Paisà made this corpus the optimal choice to train the language model, but
618 required greater effort by the corpus producers;
- 619 5. Generally, using many hours of speech (i.e. VoxForge instead of APASCI) and deep learning mod-
620 els for training acoustic models increases performance, but combining different audio corpora can
621 decrease performance probably because of practical audio-transcription inconsistencies between the
622 corpora;
- 623 6. Smaller vocabularies (e.g. CLEF and I-CAB) - even containing thousands of words - may not be
624 sufficient to gain high performance.

625 Thus, selecting and preparing optimal textual and speech corpora for LVSRs is complex and effort-
626 demanding, especially for low-resource languages. Investigating these issues was outside of this paper's
627 scope, which instead aims at introducing new acoustic models and a new decoding algorithm and comparing
628 them with a state-of-the-art ASR on a common vocabulary. However, our future experiments will investigate
629 the above issues because they call for new ways to achieve state-of-the-art performance with less training
630 material and new decoding strategies that optimally use the language model.

631 In summary, a one-million-number benchmark sentence-set was used because it corresponded to a non-
632 trivial language model that did not depend on the used training textual corpus and was reasonably comparable
633 with the one used by a reference state-of-the-art ASR. Furthermore, although the lexicon required a short

634 preparation phase, the speech included several features of spontaneous and large-vocabulary speech (e.g.
635 variability, omissions, uncertain speech, and false starts).

636 **4. Discussion and Conclusions**

637 *4.1. Summary*

638 In this paper, novel syllabic acoustic models and a new ASR have been described and compared with
639 state-of-the-art alternatives. On the single-syllable recognition task, the deep-learning models showed very
640 high performance. FHMMs gained higher performance than syllabic and phonetic HMMs, likely because
641 they recognised both syllabic- and phonetic-scale dynamics associated with different transition speeds in the
642 two parallel HMMs (Coro et al., 2007). Also, our multi-temporal CNN model forced a multi-scale analysis
643 (from phonetic to syllabic scales) and gained high performance. However, this model was not able to fully
644 capture the information contained in feature modulations and transitions in the digit and number recognition
645 tasks. Thus, its performance decreased when the acoustic structures of the modelled syllables were not
646 preserved. The properly trained CMUSphinx and KALDI ASRs reached a very high performance on the
647 [0,1M) number recognition task, but still lower than the cutting edge technology of the Google ASR.

648 In the presented experiment, our LSTM acoustic model reported the highest performance both when used
649 alone and when combined with a decoding algorithm that optimised its prediction capability. In particular,
650 on syllable recognition - i.e. without the presence of the exhaustive Viterbi decoder - the LSTM model
651 outperformed the other models. The performance remained optimal also on digit and number recognition, i.e.
652 when the LSTM was combined with the exhaustive decoding algorithm. This was not the case of the CNN
653 and the other acoustic models, which lost accuracy with respect to the baseline systems as the recognition
654 task became more and more difficult. Thus, the LSTM model both outperformed the other acoustic models
655 and was optimally used by the decoding algorithm. This property indicates that the decoding algorithm
656 was able to use this model at best, although the LSTM performance was already optimal by itself. The
657 LSTM model explicitly accounts for the unfolding of information in time, similarly to HMMs, but also
658 models both long- and short-term information. This likely corresponds to modelling high-rate and low-rate
659 dynamics within one syllable, in agreement with psycho-acoustic studies. At the same time, this behaviour
660 also overcomes the issue of modelling inter-syllabic variability from a small training set, which affected the

661 CNN model’s performance. Our results indicate that the LSTM probably learned this variability from the
662 training set and thus was able to manage a more continuous-like speech. Finally, the separation between the
663 speakers in the training and test sets reduced the potential artefact that the model was trained on the same
664 corpus the test set belonged to.

665 4.2. Using our approach with larger vocabularies and other languages

666 Extending our approach to an LVSR principally requires the availability of annotations of the syllables
667 actually spoken in the utterances (pseudo-syllables) that allow to develop an effective syllabic language
668 model. The use of the Speecon Italian corpus in the presented experiment was mainly driven by the avail-
669 ability of this information. Given the generality of our approach, the presented results are likely valid for
670 all languages currently managed by the reference ASRs of Figure 1, as long as pseudo-syllables are used
671 for acoustic modelling. Indeed, pseudo-syllables ensure the stability of the syllable structure and increase
672 acoustic-model performance (Section 2.1.2). Nevertheless, our future work will test the new proposed ASR
673 on other languages. It is worth noting that the obtained results are compliant with those reported by other
674 studies for English. The Google ASR and HMM-DNN-based ASRs can reach over 99% accuracy on clean
675 English spoken digits (Li et al., 2015) and ~97% accuracy on noisy speech (with a ~15dB signal-to-noise
676 ratio) (Milde and Köhn, 2018). On a task to recognise ~30,000 English numbers (Cole et al., 1995), per-
677 formance can range around a 93% word-recognition accuracy (Greenberg, 1997; Wu et al., 1998; Dimi-
678 trakakis and Bengio, 2011). As a general reference, with clean dialogue speech the Google ASR can reach
679 ~93% word-recognition accuracy on a ~5000 vocabulary of English words (Novoa et al., 2018), and ~63%
680 word-recognition accuracy on a 7.5 million vocabulary of English words (Kimura et al., 2019). On the
681 same million-word vocabulary, an ASR based on KALDI and HMM-DNN acoustic models can reach ~63%
682 word-recognition accuracy but is much more sensible to audio noise than the Google ASR (Kimura et al.,
683 2019).

684 Differently from the *soft alignment* process used in end-to-end models (Wang et al., 2019), our ASR
685 aligns acoustic models to the signal exhaustively and explicitly, and can re-use statistical language models
686 of standard ASRs. Overall, with respect to end-to-end models, our ASR presents a clear separation - as
687 modules - between the phases of feature extraction, language and acoustic modelling, and decoding. This

688 property allows improving the ASR by substituting alternative processes to these modules. One drawback of
689 our ASR is its high computational complexity that mainly depends on the pre-calculation of a large number
690 of likelihoods (Section 2.3). However, this complexity does not compromise efficiency for practical usages
691 of the ASR: Using a parallel implementation of the decoding algorithm on 8 cores with HMM syllabic
692 models, the recognition of a spoken number requires averagely ~5 seconds on a machine with an Intel
693 Core i7-7700HQ CPU and 16GB of Random Access Memory. Parallelising the computation on multiple
694 cores or machines would make computational time manageable also if a large number of syllables were
695 involved, e.g. a large vocabulary (~500 syllables) would require ~5s on a distributed computation using ~100
696 cores/machines on a cloud computing platform. Furthermore, the search space of the decoding algorithm
697 could be drastically reduced through "islands' recognition", i.e. by focusing the process on those portions
698 of the speech signal that are (i) acoustically relevant (prominent) compared to the surrounding units, (ii)
699 pronounced with reasonable accuracy, and (iii) more clearly recognisable (Ludusan et al., 2011). We will
700 explore also this research direction in the future.

701 4.3. Concluding remarks

702 In summary, a completely new ASR architecture has been presented. Our approach's main novel char-
703 acteristics are the inspiration by studies on the multi-temporal processing of speech in human beings and
704 the use of pseudo-syllables instead of tri-phones as acoustic models. These features have produced high-
705 quality results compared to the Google ASR and the KALDI tri-phonetic HMM-DNN ASR. Furthermore,
706 our approach has the technical advantage that it can be applied to new acoustic models and can re-use sta-
707 tistical language models of classic ASRs. The reported results suggest that taking into account the results
708 of psycho-acoustic studies - i.e. including also non-phonetic dynamics - and questioning the standard-used
709 ASR approaches may produce effective solutions. This observation positively answers to our original re-
710 search question. Furthermore, our results show that the high performance of our ASR and acoustic models
711 is likely due to the use of pseudo-syllables instead of tri-phones on the reported recognition tasks, i.e. a per-
712 ceptual syllable definition has direct benefits for both the acoustic models and the ASR. One open question
713 is if the multi-temporal processing included in our acoustic models was crucial to increase ASR perfor-
714 mance. Indeed, our CNN explicitly modelled multi-temporal processing but did not gain top performance.

715 In contrast, the LSTM model implicitly accounted for different time-scale dynamics and gained very high
716 performance.

717 **Acknowledgements**

718 The authors want to thank Daniela Burba for proofreading and correcting the paper. Coro, Origlia and
719 Cutugno want to thank and remember Prof. Renata Savy for having inspired the foundations of this research.

720 **Software**

721 Our ASR is Open Source to allow for comparisons and verification. The source codes of the decoding
722 algorithm and of the HMM-based syllabic acoustic models, and the DL trained models are available on the
723 GitHub at

724 <https://github.com/gianpaolocoro/AutomaticSpeechRecognitionResearch>

725 The source code for training the deep learning acoustic models is available at

726 <https://github.com/fvmassoli/deep-acoustic-modeling>

727 **References**

728 Abdel-Hamid, O., Mohamed, A.r., Jiang, H., Deng, L., Penn, G., Yu, D., 2014. Convolutional neural networks for
729 speech recognition. *IEEE/ACM Transactions on audio, speech, and language processing* 22, 1533–1545.

730 Ahad, A., Fayyaz, A., Mehmood, T., 2002. Speech recognition using multilayer perceptron, in: *IEEE Students Confer-*
731 *ence, ISCON'02. Proceedings., IEEE.* pp. 103–109.

732 Arnal, L.H., Poeppel, D., Giraud, A.L., 2016. A neurophysiological perspective on speech processing in “the neurobi-
733 ology of language”, in: *Neurobiology of language.* Elsevier, pp. 463–478.

734 Audhkhasi, K., Saon, G., Tüske, Z., Kingsbury, B., Picheny, M., 2019. Forget a bit to learn better: Soft forgetting for
735 ctc-based automatic speech recognition. *Proc. Interspeech 2019* , 2618–2622.

736 Baby, D., Hamme, H.V., 2015. Investigating modulation spectrogram features for deep neural network-based automatic
737 speech recognition, in: *Sixteenth Annual Conference of the International Speech Communication Association*, pp.
738 2479–2483.

739 Ballinger, B.M., Schalkwyk, J., Cohen, M.H., Allauzen, C.G.L., Riley, M.D., 2011. Speech to text conversion. US
740 Patent App. 12/976,972.

741 Baroni, M., Bernardini, S., Ferraresi, A., Zanchetta, E., 2009. The wacky wide web: a collection of very large linguisti-
742 cally processed web-crawled corpora. *Language resources and evaluation* 43, 209–226.

743 Batliner, A., Möbius, B., 2019. Prosody in automatic speech processing.

744 Bengio, Y., Simard, P., Frasconi, P., et al., 1994. Learning long-term dependencies with gradient descent is difficult.
745 *IEEE transactions on neural networks* 5, 157–166.

746 CBInsights, 2019. How Big Tech Is Battling To Own The \$ 49B Voice Market. [https://www.cbinsights.com/
747 research/facebook-amazon-microsoft-google-apple-voice/](https://www.cbinsights.com/research/facebook-amazon-microsoft-google-apple-voice/).

748 Chang, S., 2002. A syllable, articulatory-feature, and stress-accent model of speech recognition. Ph.D. thesis. University
749 of California, Berkeley.

750 Chao, G.L., Chan, W., Lane, I., 2019. Speaker-targeted audio-visual models for speech recognition in cocktail-party
751 environments. arXiv preprint arXiv:1906.05962 .

752 Chiu, C.C., Sainath, T.N., Wu, Y., Prabhavalkar, R., Nguyen, P., Chen, Z., Kannan, A., Weiss, R.J., Rao, K., Gonina,
753 E., et al., 2018. State-of-the-art speech recognition with sequence-to-sequence models, in: 2018 IEEE International
754 Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE. pp. 4774–4778.

755 CLEF, 2020. The clef initiative corpus. [http://www.clef-initiative.eu/web/clef-initiative/
756 home](http://www.clef-initiative.eu/web/clef-initiative/home).

757 CMU, 2019. The Carnegie Mellon University CLM Toolkit. [https://sourceforge.net/projects/
758 cmuspinx/files/cmuclmtk/0.7/](https://sourceforge.net/projects/cmuspinx/files/cmuclmtk/0.7/).

759 CMUSphinx, 2017. Training an acoustic model for CMUSphinx. [https://cmuspinx.github.io/wiki/
760 tutorialam/](https://cmuspinx.github.io/wiki/tutorialam/).

761 Cole, R.A., Noel, M., Lander, T., Durham, T., 1995. New telephone speech corpora at csu, in: Fourth European
762 Conference on Speech Communication and Technology, pp. 1–4.

763 Coro, G., 2004. Automatic speech recognition: A syllabic approach.
764 <https://sites.google.com/site/gianpaolocoro/ricerca/tesi-di-laurea>.

765 Coro, G., 2008. A step forward in multi-granular automatic speech recognition. Ph.D. thesis. University of Naples,
766 Federico II, Naples, Italy.

767 Coro, G., Cutugno, F., Caropreso, F., 2007. Speech recognition with factorial-HMM syllabic acoustic models, in: Eighth
768 Annual Conference of the International Speech Communication Association (Interspeech), pp. 870–873.

769 Coro, G., Masetti, G., Bonhoeffer, P., Betcher, M., 2019. Distinguishing violinists and pianists based on their brain
770 signals, in: Tetko, I.V., Kůrková, V., Karpov, P., Theis, F. (Eds.), *Artificial Neural Networks and Machine Learning –*
771 *ICANN 2019: Theoretical Neural Computation*, Springer International Publishing, Cham. pp. 123–137.

772 Cosi, P., 1998. Auditory modeling and neural networks, in: *International Summer School: Speech Processing, Recogni-*
773 *tion and Artificial Neural Networks*. Paper available from [http://www.csrif.pd.cnr.it/Papers/PieroCosi/cp-IIASS98.](http://www.csrif.pd.cnr.it/Papers/PieroCosi/cp-IIASS98.pdf)
774 pdf, Citeseer. p. 235–258.

775 Cosi, P., 2015. A kaldi-dnn-based asr system for italian, in: *2015 International Joint Conference on Neural Networks*
776 *(IJCNN)*, IEEE. pp. 1–5.

777 Cosi, P., Hosom, J.P., 1999. Hmm/neural network-based system for italian continuous digit recognition, in: *Proceedings*
778 *of the 14th International Congress of Phonetic Sciences (ICPhS '99)*, Citeseer. pp. 1669–1672.

779 Cutugno, F., Coro, G., Petrillo, M., 2005. Multigranular scale speech recognizers: Technological and cognitive view, in:
780 Bandini, S., Manzoni, S. (Eds.), *AI*IA 2005: Advances in Artificial Intelligence*, Springer Berlin Heidelberg, Berlin,
781 Heidelberg. pp. 327–330.

782 Cutugno, F., D'Anna, L., Petrillo, M., Zovato, E., 2002. APA: Towards an automatic tool for prosodic analysis, in:
783 *Speech Prosody 2002, International Conference*, pp. 231–234.

784 Cutugno, F., Leone, E., Ludusan, B., Origlia, A., 2012. Investigating syllabic prominence with conditional random
785 fields and latent-dynamic conditional random fields, in: *Thirteenth Annual Conference of the International Speech*
786 *Communication Association*, pp. 2402–2405.

787 Cutugno, F., Origlia, A., Schettino, V., 2018. 7 syllable structure, automatic syllabification and reduction phenom-
788 ena. *Rethinking Reduction: Interdisciplinary Perspectives on Conditions, Mechanisms, and Domains for Phonetic*
789 *Variation* 25, 205.

790 D'Alessandro, C., Mertens, P., 1995. Automatic pitch contour stylization using a model of tonal perception. *Computer*
791 *Speech and Language* 9, 257–288.

792 D'Anna, L., Cutugno, F., 2003. Segmenting the speech chain into tone units: human behaviour vs automatic process,
793 in: Proceedings of The XVth International Congress of Phonetic Sciences (ICPhS), pp. 1233–1236.

794 D'Anna, L., Petrillo, M., 2003. Sistemi automatici per la segmentazione in unità tonali, in: Atti delle XIII Giornate di
795 Studio del Gruppo di Fonetica Sperimentale (GFS), pp. 285–290.

796 Davis, S., Mermelstein, P., 1980. Comparison of parametric representations for monosyllabic word recognition in
797 continuously spoken sentences. IEEE transactions on acoustics, speech, and signal processing 28, 357–366.

798 Deng, L., Liu, Y., 2018. Deep Learning in Natural Language Processing. Springer.

799 Dimitrakakis, C., Bengio, S., 2011. Phoneme and sentence-level ensembles for speech recognition. EURASIP Journal
800 on Audio, Speech, and Music Processing 2011, 1–17.

801 Dunning, T., 1994. Statistical identification of language. Computing Research Laboratory, New Mexico State University
802 Las Cruces, NM, USA.

803 D'Anna, L., Coro, G., Cutugno, F., 2009. EVALITA 2009: Abla srl Participant Report, in: EVALITA 2009 Speech
804 Recognition Challenge, pp. 1–6. URL: `\url{http://www.evalita.it/2009/proceedings}`.

805 ELRA, 2019. Italian speecon database. [http://catalogue.elra.info/en-us/repository/browse/
806 ELRA-S0213/](http://catalogue.elra.info/en-us/repository/browse/ELRA-S0213/).

807 Florian, B., Sepp, K., Joshua, H., Richard, H., 2011. Hidden markov models in the neurosciences. Hidden Markov
808 Models, Theory and Applications , 169.

809 Francois, J.M., 2019. JAHMM: An implementation of hidden Markov models in Java. [https://github.com/
810 KommuSoft/jahmm](https://github.com/KommuSoft/jahmm).

811 Fujimura, O., 1975. Syllable as a unit of speech recognition. IEEE Transactions on Acoustics, Speech, and Signal
812 Processing 23, 82–87.

813 Fujimura, O., 1994. Syllable timing computation in the c/d model, in: Third International Conference on Spoken
814 Language Processing (ICLPS 1994), Yokohama, Japan, pp. 519–522.

815 Gael, J.V., Teh, Y.W., Ghahramani, Z., 2009. The infinite factorial hidden markov model, in: Advances in Neural
816 Information Processing Systems, pp. 1697–1704.

817 Ghahramani, Z., 2002. Matlab implementation of Factorial Hidden Markov Models. [http://mlg.eng.cam.ac.](http://mlg.eng.cam.ac.uk/zoubin/software.html)
818 [uk/zoubin/software.html](http://mlg.eng.cam.ac.uk/zoubin/software.html).

819 Ghahramani, Z., Jordan, M.I., 1996. Factorial hidden markov models, in: *Advances in Neural Information Processing*
820 *Systems*, pp. 472–478.

821 Girshick, R., 2015. Fast r-cnn, in: *Proceedings of the IEEE international conference on computer vision*, pp. 1440–1448.

822 Google, 2019. Cloud Speech-to-Text Features Description. [https://cloud.google.com/](https://cloud.google.com/speech-to-text/)
823 [speech-to-text/](https://cloud.google.com/speech-to-text/).

824 Graves, A., Fernández, S., Gomez, F., Schmidhuber, J., 2006. Connectionist temporal classification: labelling un-
825 segmented sequence data with recurrent neural networks, in: *Proceedings of the 23rd international conference on*
826 *Machine learning*, pp. 369–376.

827 Graves, A., Jaitly, N., 2014. Towards end-to-end speech recognition with recurrent neural networks, in: *International*
828 *conference on machine learning*, PMLR. pp. 1764–1772.

829 Greenberg, S., 1996. Understanding speech understanding: Towards a unified theory of speech perception, in: *Proceed-*
830 *ings of the ESCA Tutorial and Advanced Research Workshop on the Auditory Basis of Speech Perception*, Keele,
831 England. pp. 1–8.

832 Greenberg, S., 1997. On the origins of speech intelligibility in the real world, in: *Robust Speech Recognition for*
833 *Unknown Communication Channels*, pp. 1–11. URL: `\url{http://http.icsi.berkeley.edu/ftp/`
834 `global/pub/speech/papers/escarsr97-origins.pdf}`.

835 Greenberg, S., 1999. Speaking in shorthand—a syllable-centric perspective for understanding pronunciation variation.
836 *Speech Communication* 29, 159–176.

837 Hawkins, S., Smith, R., 2001. Polysp: A polysystemic, phonetically-rich approach to speech understanding. *Italian*
838 *Journal of Linguistics* 13, 99–188.

839 Hinton, G., Deng, L., Yu, D., Dahl, G., Mohamed, A.r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Kingsbury, B.,
840 et al., 2012. Deep neural networks for acoustic modeling in speech recognition. *IEEE Signal processing magazine*
841 29.

- 842 Hochreiter, S., 1991. Untersuchungen zu dynamischen neuronalen netzen. Diploma, Technische Universität München
843 91.
- 844 Hochreiter, S., Schmidhuber, J., 1997. Long short-term memory. *Neural computation* 9, 1735–1780.
- 845 House, D., 1996. Differential perception of tonal contours through the syllable, in: *Proc. of ICSLP*, pp. 2048–2051.
- 846 Huang, X., Acero, A., Hon, H.W., Foreword By-Reddy, R., 2001. *Spoken language processing: A guide to theory,*
847 *algorithm, and system development.* Prentice hall PTR.
- 848 Jaitly, N., Zhang, Y., Chan, W., 2019. Very deep convolutional neural networks for end-to-end speech recognition. US
849 Patent 10,510,004.
- 850 Jenkins, J.J., Strange, W., 1999. Perception of dynamic information for vowels in syllable onsets and offsets. *Perception*
851 *& psychophysics* 61, 1200–1210.
- 852 Jespersen, O., 1905. *Lehrbuch der phonetik.* *Indogermanische Forschungen* 18, 594–594.
- 853 Kahn, D., 2015. *Syllable-based generalizations in English phonology.* Routledge.
- 854 Kapur, R., 2020. The Vanishing Gradient Problem. [https://ayearofai.com/
855 rohan-4-the-vanishing-gradient-problem-ec68f76ffb9b](https://ayearofai.com/rohan-4-the-vanishing-gradient-problem-ec68f76ffb9b).
- 856 Kim, C., Kumar, M., Kim, K., Gowda, D., 2019. Power-law nonlinearity with maximally uniform distribution criterion
857 for improved neural network training in automatic speech recognition, in: *2019 IEEE Automatic Speech Recognition*
858 *and Understanding Workshop (ASRU), IEEE.* pp. 988–995.
- 859 Kimura, T., Nose, T., Hirooka, S., Chiba, Y., Ito, A., 2019. Comparison of speech recognition performance between
860 kaldi and google cloud speech api, in: *Pan, J.S., Ito, A., Tsai, P.W., Jain, L.C. (Eds.), Recent Advances in Intelligent*
861 *Information Hiding and Multimedia Signal Processing,* Springer International Publishing, Cham. pp. 109–115.
- 862 Kingma, D.P., Ba, J., 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* .
- 863 Kingsbury, B.E., Morgan, N., Greenberg, S., 1998. Robust speech recognition using the modulation spectrogram. *Speech*
864 *communication* 25, 117–132.
- 865 Knill, K.M., Gales, M.J.F., Manakul, P.P., Caines, A.P., 2019. Automatic grammatical error detection of non-native
866 spoken learner english, in: *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal*
867 *Processing (ICASSP),* pp. 8127–8131.

- 868 Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. Imagenet classification with deep convolutional neural net-
869 works, in: Pereira, F., Burges, C.J.C., Bottou, L., Weinberger, K.Q. (Eds.), *Advances in Neural Information*
870 *Processing Systems 25*. Curran Associates, Inc., pp. 1097–1105. URL: [http://papers.nips.cc/paper/](http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf)
871 [4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf](http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf).
- 872 Lamere, P., Kwok, P., Gouvea, E., Raj, B., Singh, R., Walker, W., Warmuth, M., Wolf, P., 2003. The cmu sphinx-4
873 speech recognition system, in: *IEEE Intl. Conf. on Acoustics, Speech and Signal Processing (ICASSP 2003)*, Hong
874 Kong, pp. 2–5.
- 875 LeCun, Y., Bengio, Y., et al., 1995. Convolutional networks for images, speech, and time series. *The handbook of brain*
876 *theory and neural networks* 3361, 1995.
- 877 Li, J., Deng, L., Haeb-Umbach, R., Gong, Y., 2015. *Robust automatic speech recognition: a bridge to practical applica-*
878 *tions*. Academic Press.
- 879 Logan, B., Moreno, P., 1998. Factorial hmms for acoustic modeling, in: *Proceedings of the 1998 IEEE International*
880 *Conference on Acoustics, Speech and Signal Processing, ICASSP'98* (Cat. No. 98CH36181), IEEE. pp. 813–816.
- 881 Ludusan, B., Origlia, A., Cutugno, F., 2011. On the use of the rhythmogram for automatic syllabic prominence detection,
882 in: *Twelfth Annual Conference of the International Speech Communication Association*, pp. 2413–2416.
- 883 Lyding, V., Stemle, E., Borghetti, C., Brunello, M., Castagnoli, S., Dell'Orletta, F., Dittmann, H., Lenci, A., Pirrelli,
884 V., 2014. The paisa'corpus of italian web texts, in: *9th Web as Corpus Workshop (WaC-9)@ EACL 2014, EACL*
885 *(European chapter of the Association for Computational Linguistics)*. pp. 36–43.
- 886 Maas, A.L., Qi, P., Xie, Z., Hannun, A.Y., Lengerich, C.T., Jurafsky, D., Ng, A.Y., 2017. Building dnn acoustic models
887 for large vocabulary speech recognition. *Computer Speech & Language* 41, 195–213.
- 888 MacNeilage, P.F., Davis, B.L., 2000. On the origin of internal structure of word forms. *Science* 288, 527–531.
- 889 Magnini, B., Pianta, E., Girardi, C., Negri, M., Romano, L., Speranza, M., Lenzi, V.B., Sprugnoli, R., 2006. I-cab: the
890 italian content annotation bank., in: *LREC, Citeseer*. pp. 963–968.
- 891 Malaia, E.A., Wilbur, R.B., 2019. Syllable as a unit of information transfer in linguistic communication: The entropy
892 syllable parsing model. *Wiley Interdisciplinary Reviews: Cognitive Science* .

893 Mao, S., Tao, D., Zhang, G., Ching, P., Lee, T., 2019. Revisiting hidden markov models for speech emotion recognition,
894 in: ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE.
895 pp. 6715–6719.

896 Markov, A.A., 1913. An example of statistical investigation in the text of 'eugene onyegin' illustrating coupling of tests
897 in chains, in: Proc. of the Academy of Sciences of St. Petersburg, Russia, pp. 153–162.

898 Markowitz, J.A., 2015. Robots that talk and listen: technology and social impact. de Gruyter Berlin.

899 Marr, D., 1982. Vision: A computational investigation into the human representation and processing of visual informa-
900 tion, henry holt and co. Inc., New York, NY 2.

901 Martin, P., 2010. Prominence detection without syllabic segmentation, in: Proc. of Speech Prosody [Online], pp. 1–4.
902 URL: <http://speechprosody2010.illinois.edu/papers/102010.pdf>.

903 Massaro, D., 1972. Perceptual images processing time and perceptual units in auditory perception. Psychological
904 Review 2, 124–145.

905 Massoli, F.V., Amato, G., Falchi, F., 2020. Cross-resolution learning for face recognition. Image and Vision Computing
906 , 103927.

907 Massoli, F.V., Carrara, F., Amato, G., Falchi, F., 2019. Detection of face recognition adversarial attacks. arXiv preprint
908 arXiv:1912.02918 .

909 Milde, B., Köhn, A., 2018. Open source automatic speech recognition for german, in: Speech Communication; 13th
910 ITG-Symposium, VDE. pp. 1–5.

911 Mishkin, D., Sergievskiy, N., Matas, J., 2017. Systematic evaluation of convolution neural network advances on the
912 imagenet. Computer Vision and Image Understanding 161, 11–19.

913 Muller, M.F.K.R., 2014. Estimating a-posteriori probabilities using stochastic network models, in: Proceedings of the
914 1993 Connectionist Models summer school, Psychology Press. p. 324.

915 Mustafa, M.K., Allen, T., Appiah, K., 2019. A comparative review of dynamic neural networks and hidden markov
916 model methods for mobile on-device speech recognition. Neural Computing and Applications 31, 891–899.

917 Mwiti, D., 2019. A 2019 Guide for Automatic Speech Recognition. [https://heartbeat.fritz.ai/
918 a-2019-guide-for-automatic-speech-recognition-f1e1129a141c](https://heartbeat.fritz.ai/a-2019-guide-for-automatic-speech-recognition-f1e1129a141c).

919 Naing, S.H.M., Pa Pa, W., 2018. Automatic speech recognition on spontaneous interview speech, in: 16th International
920 Conference on Computer Applications 2018 (ICCA 2018), Yangon, Myanmar, pp. 1–5.

921 NIST, 2018. SCTL, the NIST Scoring Toolkit. <https://github.com/usnistgov/SCTL>.

922 Norris, D., 1994. Shortlist: A connectionist model of continuous speech recognition. *Cognition* 52, 189–234.

923 Norris, D., McQueen, J.M., 2008. Shortlist b: a bayesian model of continuous speech recognition. *Psychological review*
924 115, 357.

925 Novak, R., Bahri, Y., Abolafia, D.A., Pennington, J., Sohl-Dickstein, J., 2018. Sensitivity and generalization in neural
926 networks: an empirical study. arXiv preprint arXiv:1802.08760 .

927 Novoa, J., Wuth, J., Escudero, J.P., Fredes, J., Mahu, R., Yoma, N.B., 2018. Dnn-hmm based automatic speech recogni-
928 tion for hri scenarios, in: Proceedings of the 2018 ACM/IEEE International Conference on Human-Robot Interaction,
929 pp. 150–159.

930 Origlia, A., Abete, G., Cutugno, F., 2013. A dynamic tonal perception model for optimal pitch stylization. *Computer*
931 *Speech & Language* 27, 190–208.

932 Origlia, A., Cutugno, F., 2016. Combining energy and cross-entropy analysis for nuclear segments detection., in:
933 INTERSPEECH, pp. 2958–2962.

934 Origlia, A., Cutugno, F., Galatà, V., 2014. Continuous emotion recognition with phonetic syllables. *Speech Communi-*
935 *cation* 57, 155–169.

936 Ortis, A., Farinella, G.M., Battiato, S., 2019. An overview on image sentiment analysis: Methods, datasets and
937 current challenges, in: Proceedings of the 16th International Joint Conference on e-Business and Telecommuni-
938 cations, ICETE 2019 - Volume 1: DCNET, ICE-B, OPTICS, SIGMAP and WINSYS, Prague, Czech Republic,
939 July 26-28, 2019., pp. 296–306. URL: <https://doi.org/10.5220/0007909602900300>, doi:10.5220/
940 0007909602900300.

941 Ostendorf, M., 1999. Moving beyond the ‘beads-on-a-string’ model of speech, in: Proc. IEEE ASRU Workshop, pp.
942 79–84.

943 Ostendorf, M., Digalakis, V.V., Kimball, O.A., 1996. From hmm’s to segment models: A unified view of stochastic
944 modeling for speech recognition. *IEEE Transactions on speech and audio processing* 4, 360–378.

945 Padrell-Sendra, J., Martín-Iglesias, D., Diaz-de Maria, F., 2006. Support vector machines for continuous speech recog-
946 nition, in: 2006 14th European Signal Processing Conference, IEEE. pp. 1–4.

947 Paliwal, K.K., 1999. On the use of filter-bank energies as features for robust speech recognition, in: ISSPA'99. Pro-
948 ceedings of the Fifth International Symposium on Signal Processing and its Applications (IEEE Cat. No. 99EX359),
949 IEEE. pp. 641–644.

950 Pan, J., Liu, C., Wang, Z., Hu, Y., Jiang, H., 2012. Investigation of deep neural networks (dnn) for large vocabulary con-
951 tinuous speech recognition: Why dnn surpasses gmms in acoustic modeling, in: 2012 8th International Symposium
952 on Chinese Spoken Language Processing, IEEE. pp. 301–305.

953 Parcollet, T., Zhang, Y., Morchid, M., Trabelsi, C., Linarès, G., De Mori, R., Bengio, Y., 2018. Quaternion convolutional
954 neural networks for end-to-end automatic speech recognition. arXiv preprint arXiv:1806.07789 .

955 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga,
956 L., et al., 2019. Pytorch: An imperative style, high-performance deep learning library, in: Advances in Neural
957 Information Processing Systems, pp. 8024–8035.

958 Patel, T., Krishna, D., Fathima, N., Shah, N., Mahima, C., Kumar, D., Iyengar, A., 2018. Development of large vocabu-
959 lary speech recognition system with keyword search for manipuri., in: Interspeech, pp. 1031–1035.

960 Peeva, M.G., Guenther, F.H., Tourville, J.A., Nieto-Castanon, A., Anton, J.L., Nazarian, B., Alario, F.X., 2010. Distinct
961 representations of phonemes, syllables, and supra-syllabic sequences in the speech production network. *Neuroimage*
962 50, 626–638.

963 Peters, J., Matusov, E., Meyer, C., Klakow, D., 2011. Topic specific models for text formatting and speech recognition.
964 US Patent 8,041,566.

965 Pieraccini, R., 2012. *The voice in the machine: building computers that understand speech*. MIT Press.

966 Pinson, M.B., Pinson, D.T., 2019. Syllable based automatic speech recognition. US Patent App. 16/031,637.

967 Povey, D., Ghoshal, A., Boulianne, G., Burget, L., Glembek, O., Goel, N., Hannemann, M., Motlicek, P., Qian, Y.,
968 Schwarz, P., et al., 2011. The kaldi speech recognition toolkit, in: IEEE 2011 workshop on automatic speech recog-
969 nition and understanding, IEEE Signal Processing Society. pp. 1–4. IEEE Catalog No.: CFP11SRW-USB.

970 Qu, Z., Haghani, P., Weinstein, E., Moreno, P., 2017. Syllable-based acoustic modeling with ctc-smbr-lstm, in: 2017
971 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU), pp. 173–177. doi:10.1109/ASRU.
972 2017.8268932.

973 Rabiner, L.R., Juang, B., 1986. A tutorial on hidden markov models. *IEEE ASSP Magazine* 3, 4–16.

974 Rao, K., Sak, H., Prabhavalkar, R., 2017. Exploring architectures, data and units for streaming end-to-end speech recog-
975 nition with rnn-transducer, in: 2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU),
976 IEEE. pp. 193–199.

977 Ravanelli, M., Omologo, M., 2017. Contaminated speech training methods for robust dnn-hmm distant speech recogni-
978 tion. arXiv preprint arXiv:1710.03538 .

979 Roach, P., 2000. *English Phonetics and Phonology. A Practical Course*. Cambridge University Press.

980 Rong, F., Isenberg, A.L., Sun, E., Hickok, G., 2018. The neuroanatomy of speech sequencing at the syllable level. *PLoS*
981 *one* 13, e0196381.

982 Sahu, P., Dua, M., Kumar, A., 2018. Challenges and issues in adopting speech recognition, in: *Speech and Language*
983 *Processing for Human-Machine Communications*. Springer, pp. 209–215.

984 Sainath, T.N., Pang, R., Rybach, D., He, Y., Prabhavalkar, R., Li, W., Visontai, M., Liang, Q., Strohman, T., Wu, Y.,
985 et al., 2019. Two-pass end-to-end speech recognition. arXiv preprint arXiv:1908.10992 .

986 Sak, H., Senior, A., Beaufays, F., 2014. Long short-term memory recurrent neural network architectures for large
987 scale acoustic modeling, in: *Fifteenth annual conference of the international speech communication association*, pp.
988 338–342.

989 Scharenborg, O., Norris, D., Ten Bosch, L., McQueen, J.M., 2005. How should a speech recognizer work? *Cognitive*
990 *Science* 29, 867–918.

991 Senior, A., Sak, H., Shafran, I., 2015. Context dependent phone models for lstm rnn acoustic modelling, in: 2015 IEEE
992 *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4585–4589. doi:10.1109/
993 *ICASSP.2015.7178839*.

994 Serizel, R., Giuliani, D., 2017. Deep-neural network approaches for speech recognition with heterogeneous groups of
995 speakers including children. *Natural Language Engineering* 23, 325–350.

996 Siemund, R., Höge, H., Kunzmann, S., Marasek, K., 2000. Speecon-speech data for consumer devices, in: LREC, Cite-
997 seer. pp. 329–333. URL: \url{"http://www.lrec-conf.org/proceedings/lrec2000/pdf/63.
998 pdf"}.

999 Smit, M.P., Virpioja, S., Kurimo, M., 2018. Advances in subword-based hmm-dnn speech recognition across languages.
1000 Submitted to Language Resources and Evaluation 29.

1001 Soltau, H., Liao, H., Sak, H., 2016. Neural speech recognizer: Acoustic-to-word lstm model for large vocabulary speech
1002 recognition. arXiv preprint arXiv:1610.09975 .

1003 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R., 2014. Dropout: a simple way to prevent
1004 neural networks from overfitting. The journal of machine learning research 15, 1929–1958.

1005 Swietojanski, P., Ghoshal, A., Renals, S., 2014. Convolutional neural networks for distant speech recognition. IEEE
1006 Signal Processing Letters 21, 1120–1124.

1007 Szaszák, G., Tündik, M.Á., Beke, A., 2016. Summarization of spontaneous speech using automatic speech recognition
1008 and a speech prosody based tokenizer., in: 8th International Conference on Knowledge Discovery and Information
1009 Retrieval (KDIR 2016), Porto, Portugal, pp. 221–227.

1010 Tu, Y.H., Du, J., Dai, L.R., Lee, C.H., 2016. A speaker-dependent deep learning approach to joint speech separation and
1011 acoustic modeling for multi-talker automatic speech recognition, in: 2016 10th International Symposium on Chinese
1012 Spoken Language Processing (ISCSLP), IEEE. pp. 1–5.

1013 Tyagi, V., McCowan, I., Misra, H., Boulard, H., 2003. Mel-cepstrum modulation spectrum (mcms) features for robust
1014 asr, in: 2003 IEEE Workshop on Automatic Speech Recognition and Understanding (IEEE Cat. No. 03EX721), IEEE.
1015 pp. 399–404.

1016 Tyagi, V., Wellekens, C., 2005. On desensitizing the mel-cepstrum to spurious spectral components for robust speech
1017 recognition, in: Proceedings.(ICASSP'05). IEEE International Conference on Acoustics, Speech, and Signal Process-
1018 ing, 2005., IEEE. pp. I–529.

1019 Virtanen, T., 2006. Speech recognition using factorial hidden markov models for separation in the feature space, in:
1020 Ninth International Conference on Spoken Language Processing, pp. 89–92.

1021 Viterbi, A., 1967. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. IEEE
1022 transactions on Information Theory 13, 260–269.

- 1023 VoxForge, 2012. VoxForge Free Speech Recognition Corpora. <http://www.voxforge.org/>.
- 1024 Wang, D., Wang, X., Lv, S., 2019. An overview of end-to-end automatic speech recognition. *Symmetry* 11, 1018.
- 1025 Warren, R.M., Healy, E.W., Chalikia, M.H., 1996. The vowel-sequence illusion: Intrasubject stability and intersubject
1026 agreement of syllabic forms. *The Journal of the Acoustical Society of America* 100, 2452–2461.
- 1027 Watanabe, S., Hori, T., Karita, S., Hayashi, T., Nishitoba, J., Unno, Y., Soplin, N.E.Y., Heymann, J., Wiesner, M., Chen,
1028 N., et al., 2018. Espnet: End-to-end speech processing toolkit. *arXiv preprint arXiv:1804.00015* .
- 1029 Weng, C., Cui, J., Wang, G., Wang, J., Yu, C., Su, D., Yu, D., 2018. Improving attention based sequence-to-sequence
1030 models for end-to-end english conversational speech recognition., in: *Interspeech*, pp. 761–765.
- 1031 Wu, S.L., Kingsbury, E., Morgan, N., Greenberg, S., 1998. Incorporating information from syllable-length time scales
1032 into automatic speech recognition, in: *Proceedings of the 1998 IEEE International Conference on Acoustics, Speech*
1033 *and Signal Processing, ICASSP'98 (Cat. No. 98CH36181), IEEE*. pp. 721–724.
- 1034 Wu, S.L., Kingsbury, E.D., Morgan, N., Greenberg, S., 1998. Incorporating information from syllable-length time scales
1035 into automatic speech recognition, in: *Proceedings of the 1998 IEEE International Conference on Acoustics, Speech*
1036 *and Signal Processing, ICASSP '98 (Cat. No.98CH36181), pp. 721–724 vol.2.*
- 1037 Young, S.J., Russell, N., Thornton, J., 1989. Token passing: a simple conceptual model for connected speech recognition
1038 systems. *Cambridge University Engineering Department Cambridge*.
- 1039 Yu, D., Deng, L., 2015. Deep neural network-hidden markov model hybrid systems, in: *Automatic Speech Recognition*.
1040 Springer, pp. 99–116.
- 1041 Yule, G., Bernini, G., 1997. *Introduzione alla linguistica. Il mulino*.
- 1042 Zeghidour, N., Usunier, N., Synnaeve, G., Collobert, R., Dupoux, E., 2018a. End-to-end speech recognition from the
1043 raw waveform. *arXiv preprint arXiv:1806.07098* .
- 1044 Zeghidour, N., Xu, Q., Liptchinsky, V., Usunier, N., Synnaeve, G., Collobert, R., 2018b. Fully convolutional speech
1045 recognition. *arXiv preprint arXiv:1812.06864* .
- 1046 Zhang, Y., Chan, W., Jaitly, N., 2017. Very deep convolutional networks for end-to-end speech recognition, in: *2017*
1047 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE*. pp. 4845–4849.

1048 **Authors' short bio**

1049 Gianpaolo Coro is a Physicist with a Ph.D. in Computer Science with a focus on Automatic Speech Recognition.
1050 His research focuses on Artificial Intelligence, Data Mining, Cloud Computing, and Open Science paradigm applied to
1051 the domains of Ecology and Natural Language Processing.

1052 Fabio Valerio Massoli is a PostDoc at the AIMH lab of ISTI-CNR. He has a Ph.D. in High Energy Physics from Uni-
1053 versity of Bologna, in collaboration with the Columbia University (NY), with a thesis on Dark Matter search. Currently,
1054 his research interests include deep learning, supervised and unsupervised learning, generative models, and quantum
1055 theory and technologies.

1056 Antonio Origlia took his PhD in 2013 with a thesis on on Affective Computing, focusing on emotional speech
1057 analysis with robotics applications. Then, he concentrated on Human-Computer Interaction topics, mainly focusing
1058 on applications for Cultural Heritage also involving the use of speech. His work mainly concentrates on probabilistic
1059 dialogue systems and their use in advanced applications developed using game engines.

1060 Francesco Cutugno is associate professor of Computational Linguistics and Human Machine Interaction at Univer-
1061 sity Federico II of Naples, Italy. From 2013 to 2018 he has been the President of the Italian Speech Sciences Association.
1062 His main research interests are in the fields of acoustic phonetics; computational linguistics; automatic spoken dialogue
1063 systems, technology applications in the cultural heritage sector.

di	dje	do	due	dze	kwa	kwan	kwe	kwin	la	lle
mi	nno	no	o	ran	ro	se	sei	sil	sp	ssan
sse	ta	ti	to	tre	tren	tSa	tSen	tSi	tSin	tSo
ttan	tte	tto	ttor	ttro	tu	u	un	van	ve	ven

Table 1: Overall set of 42 pseudo-syllables involved in our experiment, plus two silence annotations: "sil" indicates a long silence ($\geq 200ms$), whereas "sp" indicates a shorter pause.

Model Name	Accuracy (%)
a - Syllable Recognition	
LSTM	93.01
HMM-DNN-nnet2	90.38
CNN	90.25
HMM-DNN-nnet1	89.91
FHMMs	86.53
Syllabic HMMs	85.79
Phonetic HMMs	84.10
b - Digit Recognition	
LSTM + Exhaustive Viterbi	98.00
Google Speech-to-Text	97.50
KALDI - HMM-DNN-nnet2	95.06
CNN + Exhaustive Viterbi	94.00
FHMMs + Exhaustive Viterbi	93.30
Syllabic HMMs + Exhaustive Viterbi	92.00
CMUSphinx	87.74
c - Number Recognition	
LSTM + Exhaustive Viterbi	85.00
Google Speech-to-Text	81.81
KALDI - HMM-DNN-nnet2	81.20
CMUSphinx	79.00
CNN + Exhaustive Viterbi	76.60
FHMMs + Exhaustive Viterbi	72.00
Syllabic HMMs + Exhaustive Viterbi	70.00

Table 2: Performance comparison between alternative speech recognition models on the recognition of (a) the 44 units involved in our corpus of data, (b) spoken numbers from 0 to 9 (digits), (c) spoken numbers between 0 and 999,999.

ASR Engine	Corpus for LM	Corpus for AM	Word Accuracy (%)
Google ASR	Google	Google	89.10
KALDI - HMM-DNN-nnet2	Paisà	VoxForge	63.64
KALDI - HMM-DNN-nnet2	Paisà	VoxForge+APASCI	67.00
CMUSphinx	Paisà	VoxForge	51.58
CMUSphinx	Paisà	VoxForge+APASCI	54.41
CMUSphinx	Paisà + I-CAB	VoxForge	49.70
CMUSphinx	Paisà + CLEF + I-CAB	VoxForge	49.90
CMUSphinx	Paisà + CLEF	VoxForge	49.90
CMUSphinx	CLEF	VoxForge	42.60
CMUSphinx	CLEF + I-CAB	VoxForge	42.00
CMUSphinx	itWaC	VoxForge	44.00
CMUSphinx	I-CAB	VoxForge	34.40
KALDI - HMM-DNN-nnet2	Paisà	APASCI	57.95
CMUSphinx	Paisà	APASCI	39.35

Table 3: Performance comparison between several large-vocabulary automatic speech recognisers at the variation of the corpora used for language model (LM) and acoustic model (AM) training: the Google speech-to-text service (Google ASR), KALDI with deep neural network emission densities used in acoustic models (KALDI - HMM-DNN-nnet2), and the Gaussian-mixture based CMUSphinx.

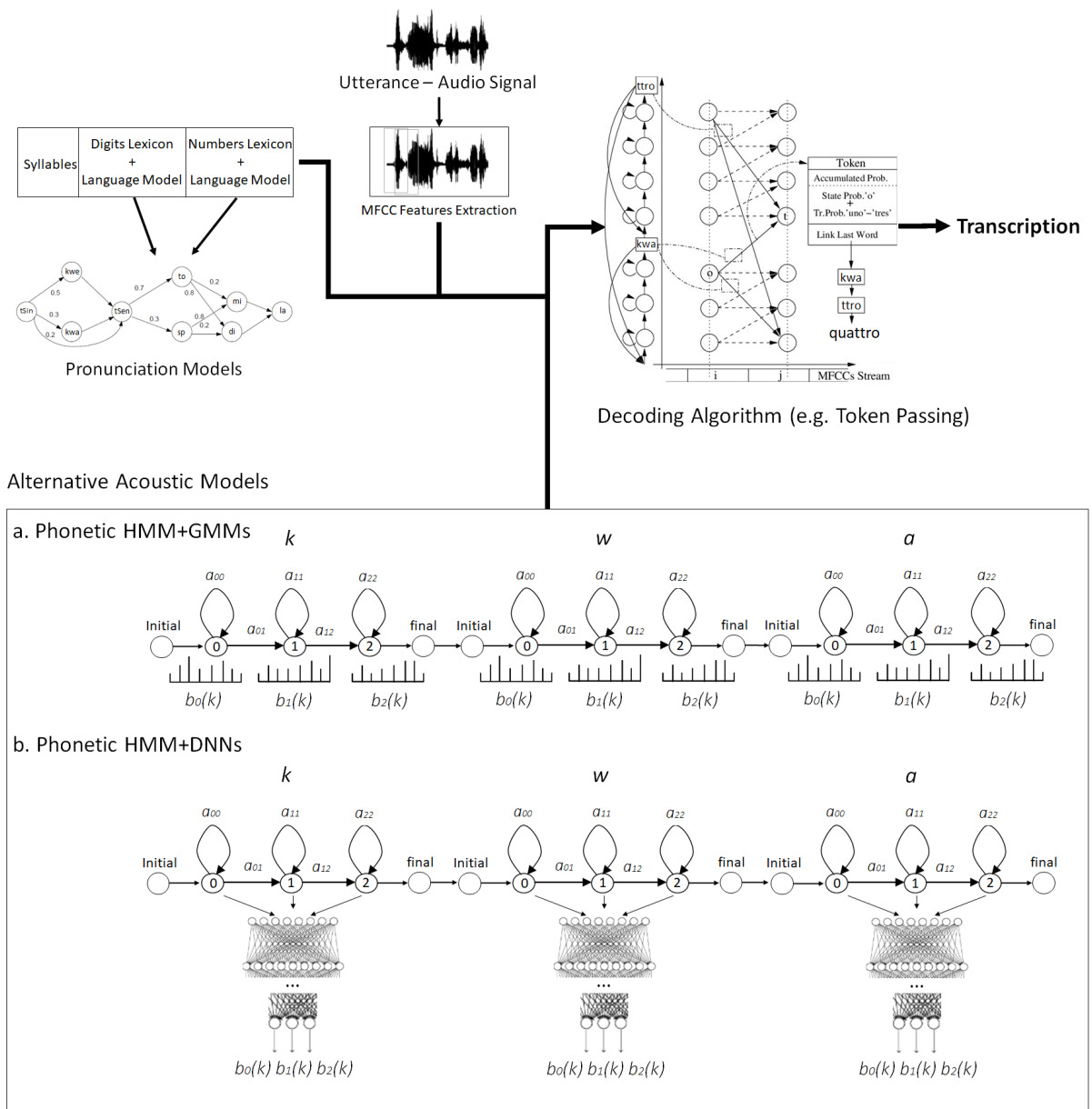
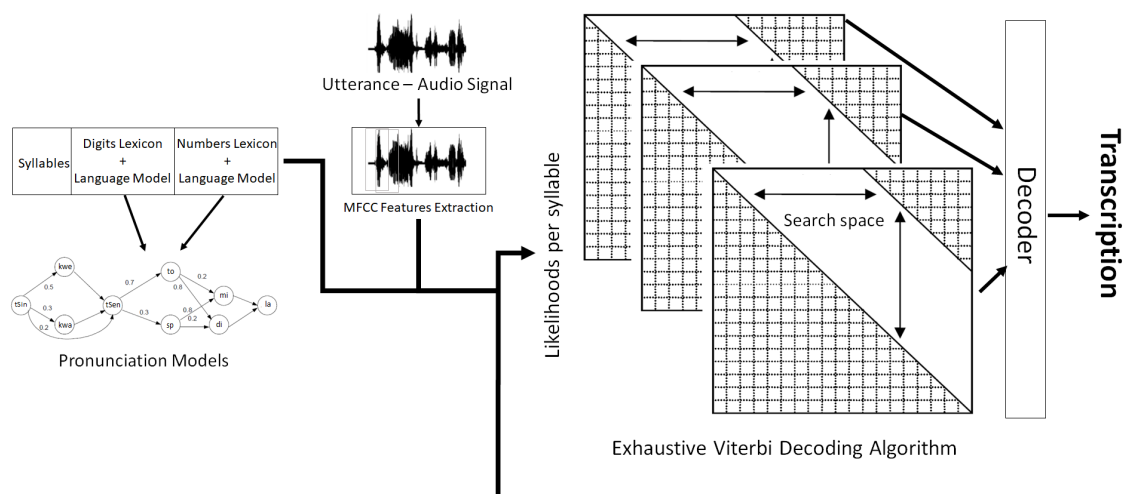


Figure 1: Diagram of a standard ASR with alternative acoustic models: a) tri-phonetic HMMs using GMM emission probabilities, b) tri-phonetic HMMs using a DNN to simulate emission probabilities. The Token Passing schema is adapted from Padrell-Sendra et al. (2006).



Alternative Acoustic Models: only one model among these is used in an ASR instance

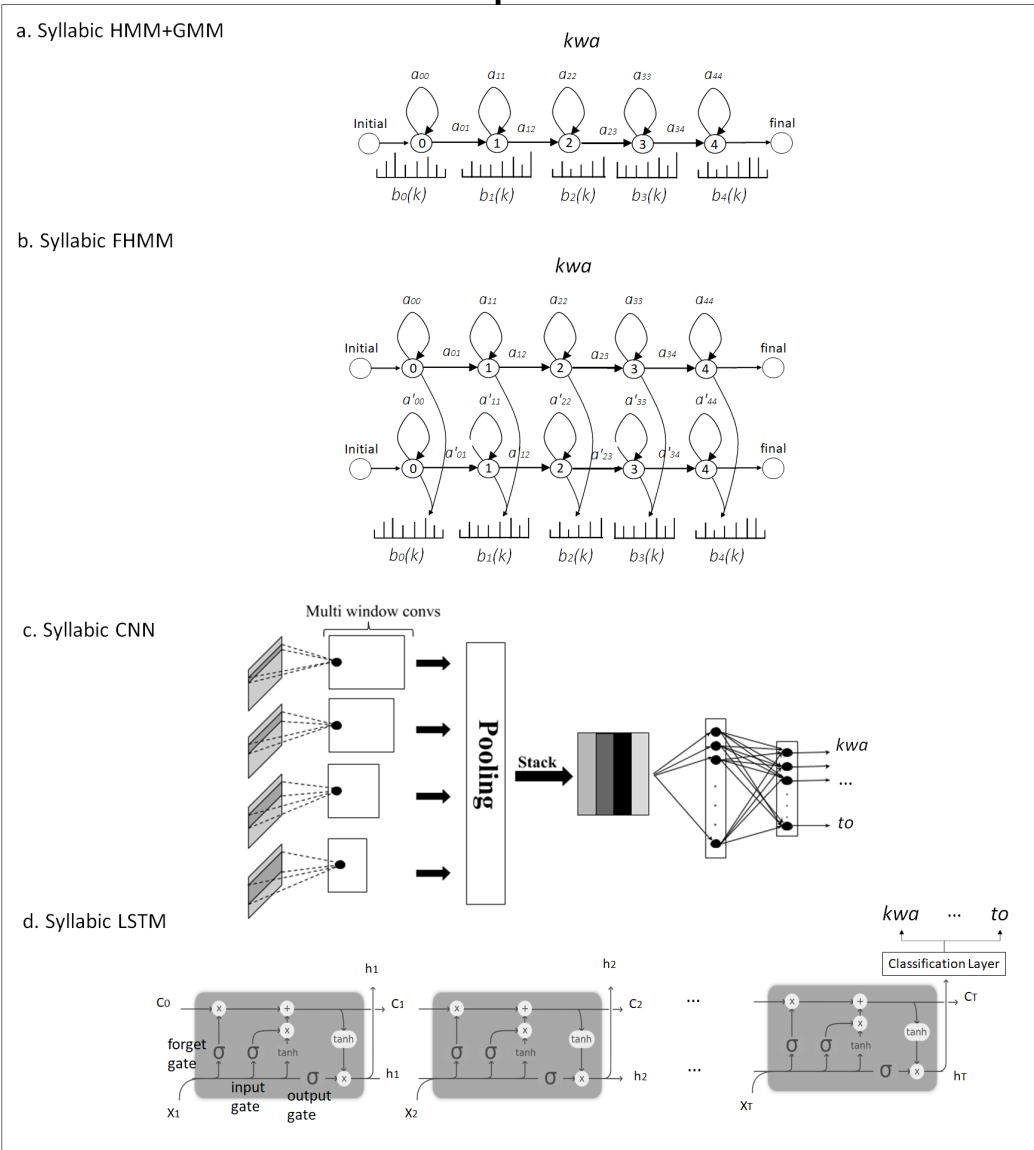


Figure 2: Diagram of our ASR with acoustic models used alternatively (only one in an ASR instance): a) syllabic HMMs using GMM emission probabilities, b) syllabic Factorial HMMs, c) Convolutional Neural Network using multi-temporal windows, with one output neuron for each syllable, and d) Long Short Term Memory model, with one output for each syllable.

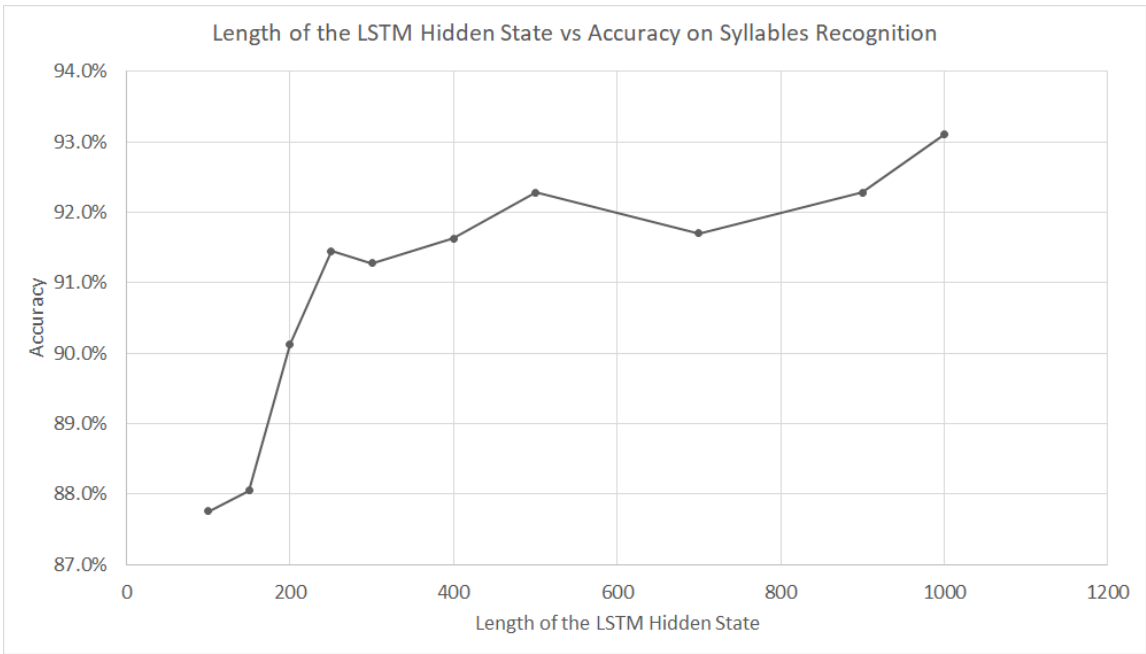


Figure 3: Variation of the accuracy of our LSTM model on the recognition of syllables of numbers between 0 and 999,999, with respect to the LSTM hidden-state length.