# Aligning Text and Phonemes for Speech Technology Applications Using an EM-Like Algorithm

R.I. DAMPER*

*Image, Speech and Intelligent Systems (ISIS) Research Group, School of Electronics and Computer Science, University of Southampton, Southampton SO17 1BJ, UK*

rid@ecs.soton.ac.uk

Y. MARCHAND

*Institute for Biodiagnostics (Atlantic), National Research Council Canada, Neuroimaging Research Laboratory, 1796 Summer Street, Suite 3900, Halifax, Nova Scotia, Canada B3H 3A7*

J.-D. S. MARSTERS AND A.I. BAZIN

*Image, Speech and Intelligent Systems (ISIS) Research Group, School of Electronics and Computer Science, University of Southampton, Southampton SO17 1BJ, UK*

**Abstract.**    A common requirement in speech technology is to align two different symbolic representations of the same linguistic 'message'. For instance, we often need to align letters of words listed in a dictionary with the corresponding phonemes specifying their pronunciation. As dictionaries become ever bigger, manual alignment becomes less and less tenable yet automatic alignment is a hard problem for a language like English. In this paper, we describe the use of a form of the expectation-maximization (EM) algorithm to learn alignments of English text and phonemes, starting from a variety of initializations. We use the British English Example Pronunciation (BEEP) dictionary of almost 200,000 words in this work. The quality of alignment is difficult to determine quantitatively since no 'gold standard' correct alignment exists. We evaluate the success of our algorithm indirectly from the performance of a pronunciation by analogy system using the aligned dictionary data as a knowledge base for inferring pronunciations. We find excellent performance—the best so far reported in the literature. There is very little dependence on the start point for alignment, indicating that the EM search space is strongly convex. Since the aligned BEEP dictionary is a potentially valuable resource, it is made freely available for research use.

**Keywords:**    text-to-speech synthesis, string alignment, dynamic programming, EM algorithm, pronunciation by analogy

## 1. Introduction

The requirement commonly arises in speech technology and natural language processing to align two linear, symbolic representations of the same linguistic entity. One important example, which forms the focus of this paper, is the alignment of the textual (orthographic or spelling) and phonemic (pronunciation) representations of isolated words (of English, in this work). The necessity to align text and phonemes arises in, for instance, inferring the complete form of spelling-pronunciation word pairs from elliptical entries in a dictionary (Lawrence and Kaye, 1986) and adding new entries to the pronunciation dictionary that provides a mapping between sub-word models and language models in automatic speech recognition (Knill and Young,

*To whom all correspondence should be addressed.

150    *Damper et al.*

**42** 1997, p. 48). But as (Jansche, 2001) writes: "The prob-
**43** lem of finding a good alignment has not received its
**44** due attention in the literature".

**45**    Two examples from the domain of text-to-
**46** speech (TTS) synthesis suffice to motivate the search
**47** for powerful automatic alignment techniques.

**48** 1. In (supervised) training of neural networks to per-
**49**    form spelling-to-sound conversion, as in the well-
**50**    known NETtalk and NETspeak of Sejnowski and
**51**    Rosenberg (1987) and McCulloch et al. (1987) re-
**52**    spectively, it is necessary to associate each letter
**53**    of an input word with a target output phoneme. In
**54**    both works, alignment was done manually, but this
**55**    is time-consuming, error-prone, and limits the size
**56**    of datasets that can be used for training. As speech
**57**    synthesis becomes ever more data-driven (Damper,
**58**    2001) using ever larger dictionaries and corpora
**59**    (Young and Bloothooft, 1997), so manual alignment
**60**    becomes less and less tenable and the need for au-
**61**    tomatic alignment methods increases.
**62** 2. Increasingly in recent years, an approach known
**63**    as *pronunciation by analogy* (PbA) has been used
**64**    in TTS synthesis to derive pronunciations for un-
**65**    known words, i.e., those not listed in the system
**66**    dictionary (Dedina and Nusbaum, 1991; Sullivan
**67**    and Damper, 1993; Pirrelli and Federici, 1994; Pir-
**68**    relli and Federici, 1995; Federici et al., 1995;
**69**    Damper and Eastmond, 1996; Yvon, 1996a; Yvon,
**70**    1996b; Damper and Eastmond, 1997; Bagshaw,
**71**    1998; Damper et al., 1999; Pirrelli and Yvon, 1999;
**72**    Marchand and Damper, 2000; Sullivan, 2001).
**73**    PbA assembles pronunciations for such (unknown)
**74**    words from partial matches to the (known) words
**75**    listed in the dictionary—a process that requires
**76**    each letter of every word in the dictionary to be
**77**    aligned with a corresponding phoneme in contigu-
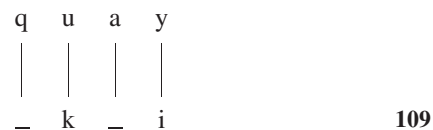**78**    ous, one-to-one fashion.

**79**    However, automatic alignment is a difficult prob-
**80** lem. Much of the difficulty arises because of the lack
**81** of regularity ('consistency' and 'transparency') in the
**82** English writing system. By 'consistency', we mean
**83** that the same letter always corresponds to the same
**84** phoneme. In fact, English is notorious for the lack
**85** of consistency in its spelling-to-sound correspondence
**86** (Venezky, 1965; Carney, 1994) at the level of single
**87** letters. For instance, the letter *c* is pronounced /s/ in
**88** *cider* but /k/ in *cat*. On the other hand, the /k/ sound
**89** of *kitten* is written with a letter *k*. By 'transparency',

**90** we mean that a single letter corresponds to a single
**91** phoneme (Henderson, 1984, p. 17) and vice versa.

**92**    The lack of consistency in English orthography is
**93** problematic for alignment since any given letter can
**94** potentially align with (i.e., correspond to) many differ-
**95** ent phonemes. To illustrate the problems that arise from
**96** lack of transparency, consider the word (*quay*, /ki/), for
**97** which a reasonable alignment might be:

q    u    a    y
|    |    |    |

k    _    i    _                                    **98**

**99**    This word is not unusual for English in having fewer
**100** phonemes than letters, necessitating the insertion of
**101** 'null phonemes' in the transcription if a one-to-one
**102** mapping is to be maintained. Such null symbols are
**103** entirely 'artificial' in that they play no role in speci-
**104** fying the pronunciation; their only purpose is to main-
**105** tain the one-to-one correspondence between letters and
**106** phonemes. Yet it is not clear precisely where the null
**107** letters should be placed, since the following is also a
**108** reasonable alignment:

q    u    a    y
|    |    |    |

_    k    _    i                                    **109**

**110**    This example illustrates a key aspect of the lack
**111** of transparency in that letter combinations frequently
**112** correspond to a single phoneme—a form of con-
**113** text dependency. Such letter combinations have been
**114** called "functional spelling units" (Venezky, 1970; Colt-
**115** heart, 1984). Examples of functional spelling units are
**116** *th* → /ð/ as in *that*, *ch* → /tʃ/ as in *church*, and *qu* → /k/
**117** as in this example of *quay*. Unfortunately, any of the let-
**118** ters of the functional spelling unit could plausibly align
**119** with the corresponding phoneme, with the others corre-
**120** sponding to nulls, leading to a degree of indeterminacy.
**121**    More rarely, there are fewer letters than phonemes
**122** in a word of English. Examples are (*six*, /sɪks/) and
**123** (*sex*, /sɛks/) in which the single letter *x* maps to the
**124** two phonemes /ks/, so that 'null letters' may have to
**125** be introduced to maintain a one-to-one mapping. Se-
**126** jnowski and Rosenberg (1987) actually invented 'new'
**127** phonemes (/K/, /X/ and /#/) in NETtalk to avoid intro-
**128** ducing null letters. As with null phonemes, the prob-
**129** lem arises as to exactly where the nulls should be

130 placed. Worse yet, both problems—null letters and null
131 phonemes—can occur in the same word, as in the case
132 of (*axe*, /aks/) for which a reasonable alignment is:

$$
\begin{array}{cccc}
a & x & \_ & e \\
| & | & | & | \\
133 \quad a & k & s & \_
\end{array}
$$

134 So the simple-minded presumption that the same num-
135 ber of letters and of phonemes implies a one-to-one
136 mapping is mistaken in this case.
137 These examples illustrate that there is no canoni-
138 cally correct alignment of text and phonemes in every
139 case, nor should we expect this, since the process is
140 essentially a computational convenience lacking any
141 sound linguistic or theoretical basis. The alignment
142 problem is especially severe for languages like En-
143 glish and French whose writing systems are 'deep', i.e.,
144 they display a complex relation between spelling and
145 sound lacking consistency and transparency, unlike the
146 'shallow' orthographies of Finnish or Serbian for exam-
147 ple, where the correspondence is mostly if not entirely
148 consistent and transparent (Coltheart, 1978; Liberman
149 et al., 1980; Katz and Feldman, 1981; Turvey et al.,
150 1984; Sampson, 1985). Indeed, (Abercrombie, 1981,
151 p. 209) describes the English spelling-to-sound system
152 as "... one of the least successful applications of the
153 Roman alphabet."
154 As one last illustration of the complexities of
155 spelling-sound correspondence in English, consider the
156 word (*made*, /meɪd/):

$$
\begin{array}{cccc}
m & a & d & e \\
| & | & | & | \\
157 \quad m & eɪ & d & \_
\end{array}
$$

158 Here, the final *e* aligns with a null phoneme, yet
159 it does not seem natural to view *de* as a functional
160 spelling unit in this case. Removing the *e* yields the
161 word (*mad*, /mad/), so that it acts as a 'marking'
162 (Venezky, 1970), signifying that the preceding vowel
163 is lengthened or dipthongized: /a/ becomes /eɪ/. This
164 contrasts with the final *e* of *axe*, which has no such
165 marking effect, further illustrating the inconsistent and
166 partly-arbitrary nature of the English spelling sys-
167 tem. Markings in English, whereby a final letter af-
168 fects the sound of a medial vowel letter, can be very
169 long range, as in the well-known example word pairs
170 *photograph*/*photography* and *telegraph*/*telegraphy*

171 (Chomsky and Halle, 1968). They can be seen as an
172 interaction of the lack of consistency and transparency,
173 both of which—as we have seen—complicate the pro-
174 cess of alignment.
175 Given these difficulties, it is clear that the automatic
176 alignment of text and phonemes is not a straightforward
177 matter. In the remainder of this paper, we develop an
178 approach to alignment based on ideas originally found
179 in Luk and Damper (1991, 1992, 1993, 1996), but us-
180 ing much-improved algorithms. Although imperfect,
181 our earlier methods have in fact been used by other au-
182 thors (e.g., Parfitt and Sharman, 1991; Jansche, 2001),
183 reflecting the widespread need for a good alignment
184 algorithm.

## 2.  Alignment by Dynamic Programming    185

186 Dynamic programming (Bellman, 1957; Kruskal,
187 1983) offers a simple and powerful way to align text
188 and phonemes on the assumption that we have some
189 knowledge of the probability of a particular letter map-
190 ping to a particular phoneme. In this work, knowl-
191 edge about letter-phoneme mappings will be compiled
192 in an 'association' matrix, $\mathbf{A}$, of dimension $L \times P$,
193 where $L$ is the size of the letter inventory (i.e., 26)
194 and $P$ is the size of the phoneme inventory (which
195 is 44 here). The dynamic programming (DP) princi-
196 ple asserts that the global solution to a path-finding
197 problem can be found by a sequence of locally-optimal
198 steps; in other words, no local non-optimality can con-
199 tribute to a globally-optimal solution. This principle is
200 well-known and widely-used in computational linguis-
201 tics and speech technology, forming for instance the
202 basis of the CYK parsing algorithm (Hopcroft et al.,
203 2001, pp. 298–301) and the Viterbi algorithm (Viterbi,
204 1967; Forney, 1973; Neuhoff, 1975), used in various
205 guises in speech recognition, speech synthesis, and text
206 processing.
207 The process of aligning text and phonemes for a spe-
208 cific word can be cast as a path-finding problem by
209 building a table, or $\mathbf{B}$ matrix, indexed by the letters of
210 the word's spelling and the phonemes of its pronun-
211 ciation. This is illustrated for the word (*phase*, /feɪz/)
212 in Fig. 1(a). The entries in this matrix are to be inter-
213 preted as degrees of 'association' between each letter
214 and each phoneme. The procedure for inferring these
215 entries is detailed in later sections. (The values seen
216 here are taken from one iteration of an actual run of our
217 algorithm.) Note that we have added word delimiters
218 (# and $ for letter and phoneme domains respectively),

152    *Damper et al.*

**219** with an association of 0 for (#, \$). This is done to allow
**220** the DP algorithm to align the leading letter or phoneme
**221** of a word with a null; otherwise the first letter would
**222** always align with the first phoneme. The 'best' align-
**223** ment of letters and phonemes is then defined by the path
**224** from the top-left entry of the matrix to the bottom-right
**225** that maximizes the accumulation of association values
**226** along this path.

**227**    To find this best alignment, we introduce two new
**228** matrices **C** and **D**. Matrix **C** is a table of accumulated
**229** associations, such that each entry is the maximum accu-
**230** mulated association up to that point in the table (i.e., up
**231** to that point in the alignment). Matrix **D** holds pointers
**232** indicating the precursor cell from which the DP algo-
**233** rithm moved to each cell. The **C** and **D** matrices are
**234** filled left-to-right, top-to-bottom using some appropri-
**235** ate form of simple recursive maximization equation.
**236** At the end of the process, the **C** matrix holds the max-
**237** imum accumulated association for the complete word
**238** in its bottom right cell, and the best alignment can be
**239** found by tracing pointers back from the bottom right
**240** cell of the **D** matrix.

   In this work, we have used the implementation of DP
due to Needleman and Wunsch (1970), since it is
simple, well-known and performed very satisfactorily
in preliminary, exploratory investigations. The spe-
cific form of the recursive maximization equation for
a given word $w$ is:

$$C_{i,j} = \max \left\{ \begin{array}{l} C_{i-1,j-1} + B_{i,j}, \\ C_{i-1,j} - \delta, \\ C_{i,j-1} - \delta \end{array} \right\} \quad \begin{array}{l} 1 \leq i \leq |l_w| \\ 1 \leq j \leq |p_w| \end{array} \tag{1}$$

**241** where $|l_w|$ and $|p_w|$ are the lengths of word $w$ in terms
**242** of letters and phonemes (including delimiters) respec-
**243** tively, and $\delta$ is some suitably chosen penalty term,
**244** which here is set to 0.

**245**    Figure 1(b) shows the **C** and **D** matrices found for
**246** the word (*phase*, /feɪz/) with the associations tabu-
**247** lated in Fig. 1(a). For ease of illustration, the two
**248** matrices are shown superimposed. If the maximiza-
**249** tion chose the $C_{i-1,j-1} + B_{i,j}$ argument, correspond-
**250** ing to a diagonal move in the **B** and **C** matrices, the
**251** entry in the **D** matrix is "↘". If the maximization
**252** chose the $C_{i-1,j}$ argument, corresponding to a ver-
**253** tical move in the **B** and **C** matrices, the entry in the
**254** **D** matrix is "↓", corresponding to alignment of a let-
**255** ter with a null phoneme. If the maximization chose the
**256** $C_{i,j-1}$ argument, corresponding to a horizontal move in
**257** the **B** and **C** matrices, the entry in the **D** matrix is "→",
**258** corresponding to alignment of a phoneme with a null

**259** letter. The "ε" in the top left cell indicates the start
**260** for the DP alignment from which no back-tracing is
**261** possible. The maximal association (or DP score) for
**262** the word is align(*phase*) = 71446. By tracing point-
**263** ers back from the bottom right entry, the alignment
**264** is found as:

```
p   h   a   s   e
|   |   |   |   |
_   f   eɪ  z   _
```
**265**

**266**    Note that the dynamic programming handles con-
**267** text dependency (e.g., letter group *ph* acts here as
**268** a functional spelling unit) in an implicit manner,
**269** since at each step of the maximization, Eq. (1),
**270** we consider moves from the three possible pre-
**271** cursors (cells $(i-1, j-1)$, $(i-1, j)$, and $(i, j-1)$)
**272** of cell $(i, j)$. At the same time, the very strong
**273** $a \rightarrow$ /eɪ/ and $s \rightarrow$ /z/ associations of 23098 and 45788
**274** respectively in Fig. 1 act as 'anchors' for the DP
**275** alignment.

**276**    It only remains to find the **A** matrix and thereafter
**277** we can align any word in the dictionary. This is done

|   | \$ | f | eɪ | z | \$ |
|---|---|---|---|---|---|
| # | 0 | 0 | 0 | 0 | 0 |
| p | 0 | 9 | 0 | 0 | 0 |
| h | 0 | 2580 | 27 | 35 | 0 |
| a | 0 | 42 | 23098 | 937 | 0 |
| s | 0 | 79 | 3 | 45788 | 0 |
| e | 0 | 947 | 1732 | 2641 | 0 |
| # | 0 | 0 | 0 | 0 | 0 |

(a)

|   | \$ | f | eɪ | z | \$ |
|---|---|---|---|---|---|
| # | 0, ε | 0, → | 0, → | 0, → | 0, → |
| p | 0, ↓ | 9, ↘ | 9, → | 9, → | 9, → |
| h | 0, ↓ | 2580, ↘ | 2580, → | 2580, → | 2580, → |
| a | 0, ↓ | 2580, ↓ | 25678, ↘ | 25678, → | 25678, → |
| s | 0, ↓ | 2580, ↓ | 25678, ↓ | 71446, ↘ | 71446, → |
| e | 0, ↓ | 2580, ↓ | 25678, ↓ | 71446, ↓ | 71446, ↘ |
| # | 0, ↓ | 2580, ↓ | 25678, ↓ | 71446, ↓ | 71446, ↘ |

(b)

*Figure 1.* (a) Example matrix of letter-phoneme associations (**B** matrix) for the word (*phase*, /feɪz/). The word is delimited by # and \$ in the letter and phoneme domains respectively. See text for explanation of entries. (b) Table of cumulative associations found by dynamic programming, together with the production or 'move' from the precursor cell that maximizes this value. This table can be viewed as a superposition of **C** and **D** matrices (see text).

278 using a form of the EM algorithm, which is the subject
279 of the next section.

## 3.  Estimating Associations with the EM Algorithm

282 The expectation-maximum (EM) algorithm is an iter-
283 ative approach to the solution of maximum-likelihood
284 estimation problems when there are data missing from
285 the set of observations and/or the likelihood function
286 cannot be easily differentiated to find its maxima. Al-
287 though the basic idea had appeared in the literature
288 previously (e.g., Hartley, 1958; Baum, 1972), the term
289 "EM algorithm" was coined by Dempster et al. (1977).
290 A useful introduction is provided by Moon (1996); an
291 excellent survey and treatment of recent developments
292 is given by McLachlan and Krishnan (1997).
293     The EM algorithm interleaves two steps, starting
294 from initial, assumed values for the missing data:

295 1. the $E$-step, in which the expected value of the like-
296    lihood is found with respect to the unknown values,
297    using the current estimate of the parameters, condi-
298    tioned on the observations.
299 2. the $M$-step, in which this expectation is maximized
300    to yield a new set of parameters.

301     The E- and M-steps are iterated with each iteration
302 guaranteed to increase the likelihood until we con-
303 verge to a local maximum of the likelihood function.
304 Convergence is proved by Dempster et al. (1977) and
305 Wu (1983) among others. Like other optimisation tech-
306 niques that find local maxima by gradient ascent, the
307 particular local maximum found in general depends on
308 the start point of the iteration—i.e., the assumed initial
309 values of the missing data.
310     In the specific case of letter-phoneme alignment, the
311 observed data are the words listed in the dictionary
312 in terms of their paired spellings/pronunciations. The
313 missing data are the parameters describing the proba-
314 bilistic correspondence between words and letters that
315 underlie the alignment process and that are compiled
316 into matrix $\mathbf{A}$. As mentioned in Section 4 below, we
317 maximize not the likelihood for word $w$ at iteration $k$
318 but the maximal DP score (as described in the previous
319 section) given the association matrix from the itera-
320 tion. Hence, the process must start with an association
321 matrix $\mathbf{A}^0$ initialized with some appropriate values.
322     The simplest way to obtain $\mathbf{A}^0$ is the *naïve* initial-
323 ization, found as follows. Processing each word of the

324 dictionary in turn, every time a letter $l$ and a phoneme $p$
325 appear in the same word, *irrespective of relative po-*
326 *sition*, the corresponding element $a_{lp}^0$ of $\mathbf{A}^0$ is incre-
327 mented. After the first pass through the dictionary, each
328 element $a_{lp}^0$ contains a count of the number of times let-
329 ter $l$ and phoneme $p$ appear in the same word. This is
330 not of course to say that a specific $l$ and $p$ do align; the
331 rationale is that they can *only* align if they occur in the
332 same word. Although we do not expect this to give a
333 very good estimate of $\mathbf{A}$, an initial alignment can be
334 attempted from $\mathbf{A}^0$.
335     Once we have this (imperfect) alignment, we can per-
336 form a second pass through the dictionary to produce a
337 new and better association matrix $\mathbf{A}^1$ with elements $a_{lp}^1$
338 that count the number of times letter $l$ and phoneme $p$
339 appear *at the same (aligned) position*, $i$. At this first
340 iteration, nulls are now introduced into the dictionary
341 as a consequence of the DP matching so that letters
342 can associate with null phonemes and phonemes can
343 associate with null letters. Although these nulls obvi-
344 ously affect the counts of letter-phoneme associations,
345 they are not themselves entered as part of the updated
346 matrix $\mathbf{A}^1$. They are omitted because to do so worked
347 far better than including nulls. If we include nulls in
348 the set of letters and phonemes at the EM stage, we are
349 effectively building in an unnatural tendency for align-
350 ments to exploit nulls, because of their cumulative high
351 scoring over a variety of situations. Hence, we restrict
352 the role of the nulls to the DP matching stage.
353     Proceeding as above, a new set of candidate align-
354 ments can now be produced and scored, a new 'best'
355 alignment again selected, and $\mathbf{A}^1$ updated to $\mathbf{A}^2$. Fur-
356 ther iterations can then be used to improve the align-
357 ments, and the estimates of the association matrix,
358 until convergence.
359     By its use of a step in which expectations of new cor-
360 respondences are computed (using the current estimate
361 of the correspondences conditioned on the dictionary
362 data) followed by a maximization step, this can be seen
363 as an EM-like algorithm.

## 4.  Issues with the Alignment Algorithm

365 Many interesting issues arise with respect to alignment
366 based on the EM and DP algorithms. In this section,
367 we briefly discuss the more important of them.
368     As a form of gradient ascent procedure, convergence
369 is to a local maximum that in general depends upon the
370 start point, i.e., the matrix $\mathbf{A}^0$. One possible start point
371 uses the simple naïve approach of the previous section.

154    *Damper et al.*

Intuitively, this has the disadvantage of allowing any letter to associate with any phoneme, no matter that one might appear at the beginning of a long word and the other at the end. Hence, an attractive possibility is to weight the entries $a_{lp}^0$ inversely according to the difference of the position indices of the $l$ and $p$ symbols. For example, the position-index difference between letter $h$ and phoneme /z/ of (*phase*, /feɪz/) is $|2 - 3| = 1$. Various weighting schemes could be envisaged. Yet another possibility is to use the manual alignments devised for training NETtalk (Sejnowski and Rosenberg, 1987) or NETspeak (McCulloch et al., 1987) to obtain $\mathbf{A}^0$. (In this latter case, the counts entered into $\mathbf{A}^0$ *will* have taken account of nulls.) Further, Black et al. (1998) have described a similar algorithm to ours in which they specify a set of "allowables", i.e., letters and phonemes that can plausibly associate on the basis of prior intuitive knowledge of letter-phoneme correspondences. This can be used to define binary values for $a_{lp}^0$ (which become continuous on subsequent EM iterations). One of the major aims of this paper was to evaluate the wide variety of possibilities for initialization (see Section 5.2).

One very important issue is evaluating quantitatively the effectiveness of any alignment algorithm. However, this is difficult since there is no canonically correct 'gold standard' alignment in all cases (see Introduction). Scoring on the basis of human judgement is likely to be subjective and inconsistent between judges and is, in any case, not practical for the sort of very large dictionaries that we wish to use. Although it is possible (and indeed sensible) to have a human expert check obvious problem cases (e.g., *axe*, *know*, *phase*, . . . ), and we did in fact do this during program development, it does not amount to a full and thorough evaluation, giving a global summary figure of merit. Thus, we have decided to assess our alignment results indirectly according to the number of words correctly transcribed by a pronunciation by analogy (PbA) system. For this purpose, we have used the PbA system of 2000.

Another issue is what we have previously called the 'harmonization' of the different phoneme inventories used by different researchers and/or dictionary compilers (Damper et al., 1999). Thus, if we wish to use the NETtalk manual alignment to estimate $\mathbf{A}^0$ in order to align a dictionary such as BEEP (see below), we must have some way of mapping the different sets of phonemes used by the different dictionaries onto a common set. Because our goal is to align BEEP, we obviously choose the BEEP symbols as the common set. Tables 1 and 2 show the harmonization scheme

*Table 1.* Harmonization scheme used to map the NETtalk phoneme set onto the BEEP set.

| NETtalk | BEEP | as in . . . | IPA |
|---|---|---|---|
| a | aa | f<u>a</u>ther | ɑ |
| b | b | <u>b</u>et | b |
| c | ao | b<u>ou</u>ght | ɔ |
| d | d | <u>d</u>ime | d |
| e | ey | b<u>a</u>ke | eɪ |
| f | f | <u>f</u>in | f |
| g | g | <u>g</u>uess | g |
| h | hh | <u>h</u>ead | h |
| i | iy | p<u>ea</u>t | i |
| k | k | <u>k</u>itten | k |
| l | l | <u>l</u>et | l |
| m | m | <u>m</u>et | m |
| n | n | <u>n</u>et | n |
| o | ow | b<u>oa</u>t | oʊ |
| p | p | <u>p</u>et | p |
| r | r | <u>r</u>ed | r |
| s | s | <u>s</u>et | s |
| t | t | <u>t</u>est | t |
| u | uw | l<u>u</u>te | u |
| v | v | <u>v</u>est | v |
| w | w | <u>w</u>et | w |
| x | ax | <u>a</u>bout | ə |
| y | y | <u>y</u>et | j |
| z | z | <u>z</u>oo | z |
| A | ay | b<u>i</u>te | aɪ |
| C | ch | <u>ch</u>in | tʃ |
| D | dh | <u>th</u>is | ð |
| E | eh | b<u>e</u>t | ɛ |
| G | ng | si<u>ng</u> | ŋ |
| I | ih | b<u>i</u>t | ɪ |
| J | jh | <u>g</u>in | ʤ |
| K | k s | se<u>x</u>ual | k ʃ |
| L | l | bott<u>le</u> | ɫ |
| M | m | abys<u>m</u> | (ə)m |
| N | n | butto<u>n</u> | (ə)n |
| O | oy | b<u>oy</u> | ɒɪ |
| Q | k w | <u>qu</u>est | k w |
| R | er | b<u>ir</u>d | ɜ |
| S | sh | <u>sh</u>in | ʃ |
| T | th | <u>th</u>in | θ |
| U | uh | b<u>oo</u>k | ʊ |
| W | aw | b<u>ou</u>t | aʊ |

(*Continue on next page.*)

*Table 1.* (*Continue*).

| NETtalk | BEEP | as in . . . | IPA |
|---------|------|-------------|-----|
| X | k s | se<u>x</u> | k s |
| Y | y uw | c<u>u</u>te | j u |
| Z | zh | lei<u>s</u>ure | ʒ |
| @ | ae | b<u>a</u>t | a |
| ! | t s | na<u>z</u>i | t s |
| # | g z | e<u>x</u>amine | ɡ z |
| + | w aa | bourge<u>ois</u> | w ɑ |
| * | w | <u>wh</u>ack | ʍ |
| ^ | ah | b<u>u</u>t | ʌ |

*Table 2.* Harmonization scheme used to map the NETspeak phoneme set onto the BEEP set.

| NETspeak | BEEP | as in . . . | IPA |
|----------|------|-------------|-----|
| A | ax | <u>a</u>bout | ə |
| B | b | <u>b</u>et | b |
| D | d | <u>d</u>ime | d |
| E | eh | b<u>e</u>t | ɛ |
| F | f | <u>f</u>in | f |
| G | g | <u>g</u>uess | ɡ |
| H | hh | <u>h</u>ead | h |
| I | ih | b<u>i</u>t | ɪ |
| J | jh | <u>g</u>in | ʤ |
| K | k | <u>k</u>itten | k |
| L | l | <u>l</u>et | l |
| M | m | <u>m</u>et | m |
| N | n | <u>n</u>et | n |
| O | oh | st<u>o</u>ck | ɒ |
| P | p | <u>p</u>et | p |
| R | r | <u>r</u>ed | r |
| S | s | <u>s</u>et | s |
| T | t | <u>t</u>est | t |
| U | ah | b<u>u</u>t | ʌ |
| V | v | <u>v</u>est | v |
| W | w | <u>w</u>et | w |
| Y | y | <u>y</u>et | j |
| Z | z | <u>z</u>oo | z |
| AA | ae | b<u>a</u>t | a |
| AI | ey | b<u>a</u>ke | eɪ |
| AR | aa | f<u>a</u>ther | ɑ |
| AW | ao | b<u>ou</u>ght | ɔ |
| CH | ch | <u>ch</u>in | tʃ |
| DH | dh | <u>th</u>is | ð |
| EE | iy | p<u>ea</u>t | i |
| EI | ea | <u>air</u> | ɛə |
| ER | er | b<u>ir</u>d | ɜ |
| EY | ih | d<u>e</u>spite | ɪ |
| GZ | g z | e<u>x</u>amine | ɡ z |
| IA | ia | <u>ear</u> | ɪə |
| IE | ay | b<u>i</u>te | aɪ |
| KH | k sh | an<u>xi</u>ous | k ʃ |
| KS | k s | se<u>x</u> | k s |
| KW | k w | <u>qu</u>est | k w |
| NG | ng | si<u>ng</u> | ŋ |
| OA | ow | b<u>oa</u>t | oʊ |
| OI | oy | b<u>oy</u> | ɒɪ |

(*Continue on next page.*)

422 used to map the NETtalk and NETspeak phoneme sets
423 onto BEEP. Note that BEEP uses a phoneme inventory
424 of 44 symbols (excluding the null phoneme), whereas
425 the NETtalk and NETspeak inventories are both of size 51
426 (again excluding the null phoneme).
427 The symbols listed in the 'NETtalk' column of Ta-
428 ble 1 are those in the file downloaded from http://
429 www.speech.cs.cmu.edu/comp.speech and
430 not the ones tabulated in Appendix A of 1987. The
431 downloaded file includes a symbol '+' which is not
432 listed in the paper and excludes a symbol '|' which is
433 listed in the paper. In general, harmonization can never
434 be an exact process, because of idiosyncratic choice of
435 phoneme inventories by the different individual com-
436 pilers of the transcribed dictionaries, which often re-
437 flect dialectal differences. For instance, Sejnowski and
438 Rosenberg (1987) use the same symbol /a/ to transcribe
439 both the a vowel in *father* and the ɒ vowel in *stock*, as
440 these are probably the same vowel for their dialect of
441 American English. So the mapping from NETtalk to
442 BEEP symbols is not one-to-one. We can only try to
443 achieve the most consistent mapping according to our
444 intuitions.
445 A final issue is that the EM algorithm is properly a
446 probabilistic algorithm. We experimented with various
447 normalizations, corresponding to various probabilistic
448 models, but none performed as well as using simple
449 (unnormalized) frequency counts directly from the as-
450 sociation matrix **A**. Hence, all results presented here
451 use this formulation. This is the reason we refer to our
452 algorithm as "EM-like". The effect of using unnormal-
453 ized counts (rather than proper probabilities) on con-
454 vergence is unknown but, as we shall see, this did not
455 prove to be an issue in practice.

156     *Damper et al.*

*Table 2.*   (*Continue*).

| NETspeak | BEEP | as in … | IPA |
|---|---|---|---|
| OO | uh | b<u>oo</u>k | ʊ |
| OU | aw | b<u>ou</u>t | aʊ |
| SH | sh | <u>sh</u>in | ʃ |
| TH | th | <u>th</u>in | θ |
| UL | l | bott<u>le</u> | ɫ |
| UR | ua | m<u>oor</u> | ωə |
| UU | uw | l<u>u</u>te | u |
| YU | y uw | c<u>u</u>te | j u |
| ZH | zh | lei<u>s</u>ure | ʒ |

## 5. Results

In this section, we report the results of using our algorithm to align a large dictionary.

### 5.1. BEEP Dictionary

Our algorithm has been tested by using it to align BEEP: the British English Example Pronunciation dictionary. BEEP is publically accessible and can be downloaded from `http://www.speech.cs.cmu.edu/comp.speech`. It is typical of the size and content of the on-line dictionaries used for current speech technology applications. BEEP was constructed by amalgamating several public domain dictionaries to yield a large composite. The version used here contained 257,033 words. Note that there has been no strong quality control in constructing BEEP. Consequently, it contains several erroneous word entries (e.g., INDISPUTABLE for *indissoluble*, UNDILAPIDATED for *undiluted*) and transcriptions (e.g., for *abnegation*). Those that we discovered have been removed but we certainly cannot guarantee to have found all errors. We also removed all words with multiple pronunciations for conformity with the evaluation protocol in Marchand and Damper (2000). This gives a dictionary with 198,632 entries in all.

### 5.2. Initializations

The following initializations were used:

- naïve;
- a weighted scheme with $W = \beta/(1 + |d|)$ where $d$ is the letter-phoneme position-index difference, and $\beta$ is a heuristic scaling set to 40 for the results reported here;
- the NETtalk manual alignment (20,009 words);
- the NETspeak manual alignment (16,280 words);
- various random alignments.

### 5.3. Convergence

The convergence criterion was that there was no change as between $\mathbf{A}^k$ and $\mathbf{A}^{k-1}$.

Figure 2 shows the convergence behavior for the NETtalk initialisation. The quantity graphed is the total DP score for the whole dictionary at the end of iteration $k$, i.e., $S_k = \sum_{i=1}^{198,632} \text{align}^k(w_i)$. Note that convergence requires that the $\mathbf{A}$ matrix is unchanged between iterations, $\mathbf{A}^k = \mathbf{A}^{k-1}$, which (because nulls are not included in the $\mathbf{A}$ matrix) is not quite the same as the total DP score remaining unchanged, $S_k = S_{k-1}$. The total DP score at the zeroth iteration, $S_0$, is very low in this case, because only the 20,009 words of the originally-aligned NETtalk dictionary can be scored.

Figure 3 shows convergence behavior for two different initializations, excluding the total DP score at the zeroth iteration, $S_0$. This gives a clearer view of the convergence for the NETtalk initialization than does Fig. 2 where the very low value of $S_0$ swamps the trend. For the naïve initialization, it is not really sensible to depict $S_0$ anyway since the dramatic overcounting of associations (every letter is counted $|p_w|$ times and every phoneme is counted $|l_w|$ times) produces a very high score that is effectively meaningless. For both initializations, most of the improvement takes place between the first and second iterations. This was found to be a general characteristic of the results. For all initializations, convergence was achieved in between 5 to 8 iterations. The manual alignment of the NETtalk dictionary, even though it is much smaller than BEEP, shows a clear benefit in terms of a higher score at iteration 1 together with faster convergence. The score at convergence, $S_C$, was remarkably consistent across the various initializations, suggesting that the search problem is strongly convex. The best value obtained was $S_C = 8.579 \times 10^{10}$ for the NETtalk initialization whereas the worst value was $S_C = 8.473 \times 10^{10}$ for one of the random initializations. Generally, the random initialization values were slightly lower than the others.

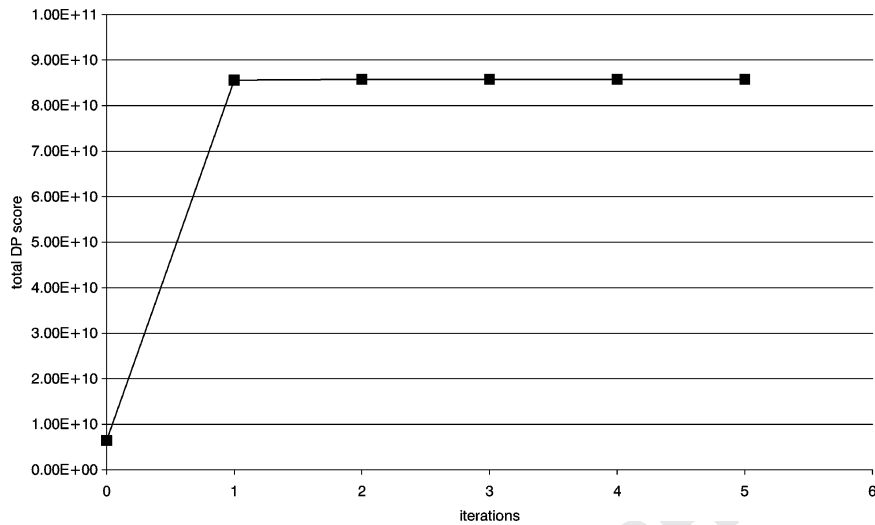Aligning Text and Phonemes for Speech Technology     157



*Figure 2.*   Convergence behavior of the alignment algorithm for the NETtalk initialization.
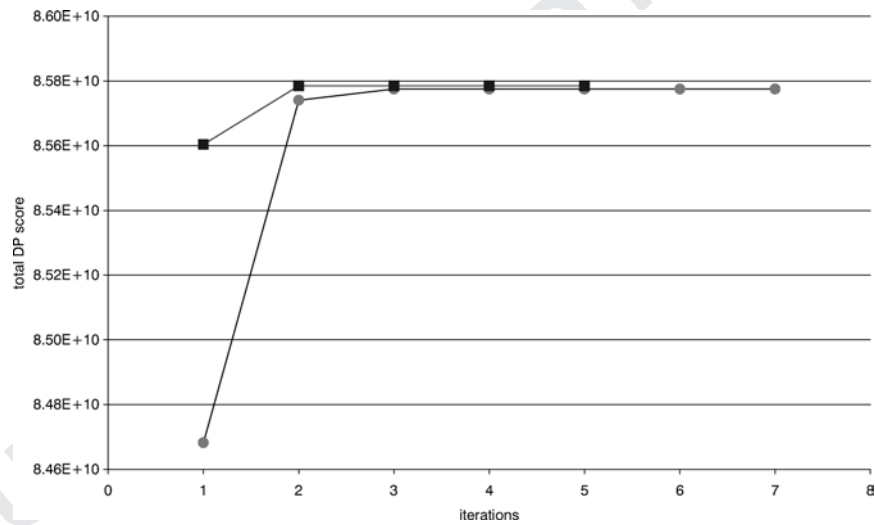


*Figure 3.*   Convergence behavior of the alignment algorithm for two different initializations. Rectangles: NETtalk initialization; Circles: naïve initialization.

### 5.4.   *Analysis of Association Matrices*

Figures 4(a) and (b) show the association matrices for the naïve initialization initially, $\mathbf{A}^0$, and at convergence, $\mathbf{A}^7$. The larger association values in Fig. 4 are a consequence of the overcounting mentioned above. As expected, the matrix is considerably less random (i.e., peakier) at convergence. Quantitatively, the (negative) entropy of the $\mathbf{A}^0$ matrix was 8.84 bits whereas that of the converged matrix was 5.24 bits; these figures compare with 10.13 bits for the equiprobable case. En-couragingly, the strongest peaks at convergence, corresponding to the major letter-phoneme associations, are also among the strongest peaks in $\mathbf{A}^0$, indicating that the naïve initialization, albeit very simple, still provides an effective start point for our algorithm.

There is a wealth of information about letter-phoneme correspondences in English to be gleaned from the $\mathbf{A}$ matrix obtained at convergence. Since nulls are introduced into the aligned dictionary only at the DP matching stage (see Section 3) and do not figure in the $\mathbf{A}$ matrix, they are not considered explicitly in
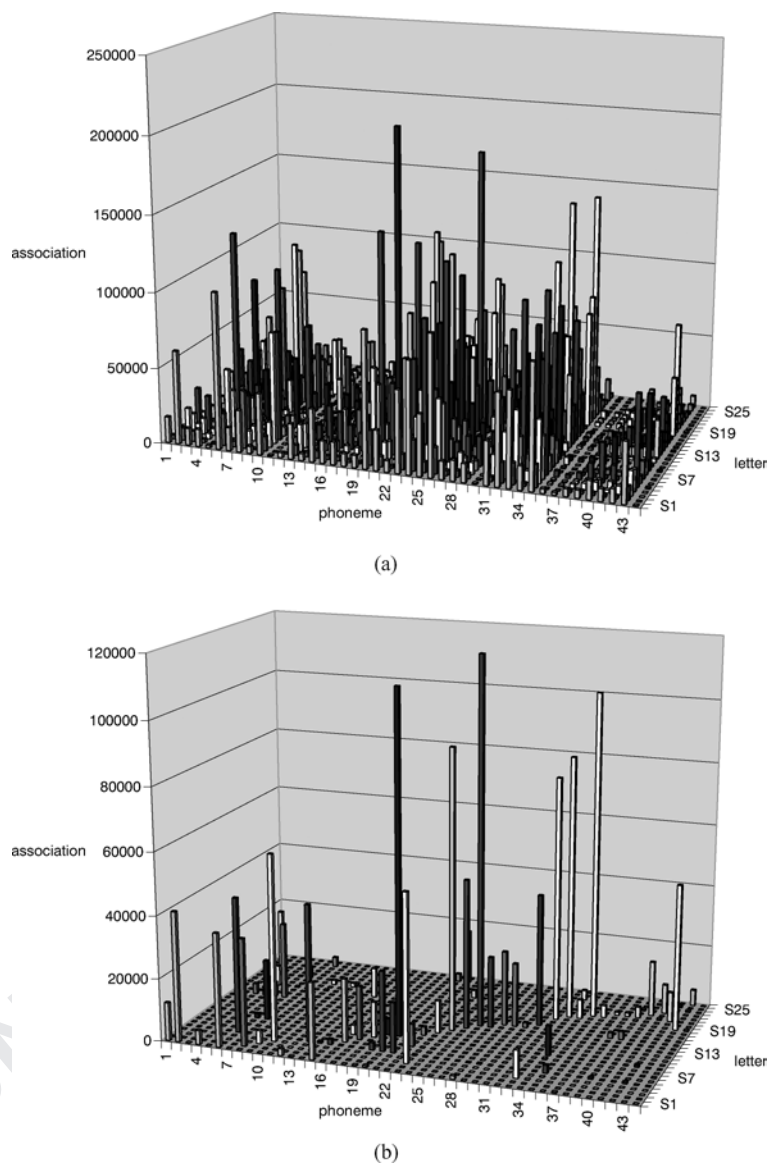
158    *Damper et al.*



*Figure 4.*    Association matrices for the naïve initialization both initially, $\mathbf{A}^0$, and at convergence, $\mathbf{A}^7$.

the remarks that follow. With this proviso, the commonest correspondence overall was $n \rightarrow$ /n/. The commonest letter participating in correspondences is $i$, which occurs 148,913 times in the matrix. This is slightly surprising as the commonest letter overall is $e$. The apparent discrepancy is explained by the number of times letter $e$ participates in a functional spelling unit such as $ea$ and so aligns with null (with the letter $a$ aligning with the vowel phoneme). The least common letter participating in correspondences is $q$, which occurs just 17 times. Again, $q$ almost invariably occurs in a $qu$ functional spelling unit, with $q$ aligning with a null phoneme, which reduces its count in the matrix. The commonest phoneme is /ɪ/ at 138,176 occurrences, which can be understood from the frequency with which letter $i$ occurs and the fact that $i \rightarrow$ /ɪ/ is a very common correspondence (at 109,508 occurrences). Schwa, /ə/, is relatively less common than /ɪ/ at 190,975 occurrences. Intuitively, one might expect schwa to be the commonest vowel, but it is perhaps more likely than /ɪ/ to align with a null letter. Of the letters, $o$ displays most variability in its association with

**574** phonemes, with no less than nine correspondences with
**575** a frequency count of 1000 or more. The least variabil-
**576** ity is shown by letter *m*, which almost always associates
**577** with phoneme /m/. Schwa displays easily the most vari-
**578** ability in its association with letters, participating in
**579** five correspondences (with letters *a*, *e*, *i*, *o* and *u*) with
**580** a count greater than 1000. The least variable phoneme
**581** was /ŋ/, which associated with letter *n* in all but just
**582** 2 cases.

**583** *5.5.  Assessing Alignment Performance Using PbA*

**584** As previously stated, alignment results were assessed
**585** using PbA. Each word was removed from the dictio-
**586** nary and a pronunciation determined from the word's
**587** spelling by analogy with all other words. The March-
**588** hand and Damper PbA system uses multiple (actually
**589** five) criteria to select between candidate pronuncia-
**590** tions to find the 'best'. There is, however, a problem in
**591** that PbA was designed to transcribe text in which there
**592** will obviously be no null letters. Yet here, null letters
**593** have been added to the alignments of many words. Our
**594** first step, then, has been to ignore any words with null
**595** letters, reducing the number of words to be tested from
**596** 198,632 to approximately 177,000. (The number varies
**597** with the exact initialization used.)  This is an obvious
**598** simplification of the problem, but should nonetheless
**599** yield interesting insights.
**600**    Table 3 shows results obtained (for words without
**601** null letters) in terms of words and phonemes correctly
**602** pronounced for each of the initializations used. Several
**603** different random initializations were used, but results
**604** were very similar and so figures for one only are tab-
**605** ulated here. In each case, we show the results for the
**606** best single scoring criterion of the five, for the best
**607** combination, and when all five are combined. Note
**608** that 10100 in the column heading indicates that scor-
**609** ing strategies 1 and 3 as described by Marchand and
**610** Damper (2000, pp. 207–208) provided the best com-
**611** bination performance for all initializations. Although
**612** space precludes a full description of our PbA method-
**613** ology, we mention that strategy 1 takes the product of
**614** arc frequencies along the shortest path in the pronun-
**615** ciation lattice, whereas strategy 3 counts the number
**616** of identical pronunciations having the same shortest
**617** path length. Strategy 1 is relatively popular in PbA
**618** (e.g., Damper and Eastmond, 1997) whereas we are
**619** not aware that any other researchers have ever used
**620** strategy 3, which interestingly turns out to be best per-
**621** forming single strategy overall.

*Table 3.*  Results when alignment of the BEEP dictionary is assessed by the performance of a pronunciation by analogy system, for various initializations. Words with null letters in their alignments have been ignored at this stage.

| | Best Single | Best Combination | All 5 |
|---|---|---|---|
| NAïVE | 00100 | 10100 | 11111 |
| Words (%) | 85.84 | 87.32 | 85.96 |
| Phonemes (%) | 97.52 | 97.78 | 97.57 |
| *W* WEIGHTED | 00100 | 10100 | 11111 |
| Words (%) | 85.87 | 87.36 | 86.00 |
| Phonemes (%) | 97.60 | 97.85 | 97.65 |
| NETTALK | 00100 | 10100 | 11111 |
| Words (%) | 86.00 | 87.41 | 86.05 |
| Phonemes (%) | 97.59 | 97.83 | 97.63 |
| NETSPEAK | 00100 | 10100 | 11111 |
| Words (%) | 86.01 | 87.48 | 86.11 |
| Phonemes (%) | 97.64 | 97.89 | 97.70 |
| RANDOM | 00100 | 10100 | 11111 |
| Words (%) | 85.87 | 87.38 | 85.69 |
| Phonemes (%) | 97.51 | 97.78 | 97.57 |

**622**    The figures in Table 3 are remarkably consistent, in-
**623** dicating that the particular initialization used does not
**624** have a dramatic effect. This is in spite of our attempts
**625** to restart the algorithm from a variety of very differ-
**626** ent points, suggesting that the search space is strongly
**627** convex. It is worth noting, however, that as a con-
**628** sequence of the large dictionary size (approximately
**629** 177,000 words) the difference between the best Best
**630** Combination of 87.48% (for the NETtalk initialization)
**631** and the worst Best Combination of 87.32% (for the
**632** naïve initialization) is in fact marginally significant at
**633** the 5% level (binomial test, $z = 2.026$, $p \sim 0.021$).

**634**    The best PbA performance is found for NETspeak
**635** but initializing alignment with the NETspeak dictionary
**636** actually produced a slightly lower total DP score at con-
**637** vergence than initializing with NETtalk. In other words,
**638** the total DP score at convergence is a good but not per-
**639** fect indicator of PbA performance. Examination of the
**640** final alignments revealed that these were strongly sim-
**641** ilar; there were typically somewhere between 10 and
**642** 100 different alignments only between one initializa-
**643** tion and another. Most often, differences were due to
**644** the specific placement of nulls in words having many
**645** silent letters (e.g., *bourgeoisie*, *heavyweight*, *mem-*
**646** *oirs*). Frequently, these were words of foreign (French)
**647** origin.

648    This is certainly among the best performance fig-
649 ures ever reported on English letter-phoneme conver-
650 sion, in terms of word-level accuracy on a large dic-
651 tionary. Previously (Damper et al., 1999), we obtained
652 71.8% words correct using PbA on a much smaller
653 dictionary—the 16,280 manually-aligned words used
654 by McCulloch et al. (1987) to train NETspeak. (It
655 should be noted, however, that BEEP uses a smaller
656 phoneme inventory of 44 symbols than the 51 used
657 in the NETspeak dictionary, making for a somewhat
658 easier problem.) A further observation is that using
659 all five strategies does not give best performance, as
660 it did for our earlier work with smaller dictionaries
661 (Marchand and Damper, 2000). In assessing perfor-
662 mance, however, we must remember that we have sim-
663 plified the problem by ignoring words with nulls, which
664 arguably gives a too optimistic view of the present re-
665 sults. However, even under the maximally pessimistic
666 assumption that PbA were to get *all* the words with
667 null letters wrong, the 85.8% words correct for best
668 single strategy, naïve start point, would fall to 76.1%—
669 still a very respectable result on such a sizable
670 dictionary.

671    To gain further insight into this issue, PbA was used
672 to produce pronunciations for all 198,632 words includ-
673 ing those with null letters in their alignment, treating the
674 latter as a legitimate input symbol (even though it never
675 could be in practice). Results for the best combination
676 averaged 82.3% words correct, showing that high ac-
677 curacy is potentially achievable if only 'missing' nulls
678 in the PbA input could be appropriately introduced.

679    **6. Discussion and Conclusions**

680 We have described a form of the EM algorithm, used
681 with dynamic programming to align a dictionary of
682 word spellings and their pronunciations. Such align-
683 ment problems commonly occur in speech technology
684 and natural language processing. The issues that arise
685 in solving this important problem have been detailed
686 and discussed. The quality of the obtained alignment
687 has been assessed using pronunciation by analogy to
688 derive pronunciations for all words in the dictionary
689 from their spelling, using the aligned data as a knowl-
690 edge base. Since the EM algorithm is effectively a gra-
691 dient ascent procedure prone to finding local maxima,
692 alignment has been performed from a variety of ini-
693 tializations, or start points. Results are judged to be
694 extremely encouraging, and are relatively insensitive
695 to a wide variety of start points. This indicates that the

696 search space is strongly convex and, hence, that local
697 maxima are not a practical problem.

698    Our work has several similarities with that of Ristad
699 and Yianilos (1998). This is perhaps not surprising as
700 they take the topic of stochastic transduction as their
701 motivation, whereas the ideas reported in this paper
702 had their early expression in our own work, which led
703 to the use of stochastic transduction to solve problems
704 in TTS conversion, including letter-phone alignment
705 (Luk and Damper, 1996, 1998). Ristad and Yianilos
706 also use dynamic programming in conjunction with
707 the EM algorithm to learn edit distances between two
708 strings. Since the string edit operations of insertion and
709 deletion can be interpreted as the introduction of nulls
710 into one string or another—either the word's spelling
711 or its pronunciation—there is clearly a strong relation
712 between the two pieces of work. As Jansche (2001)
713 writes: "The problem of letter-to-sound conversion is
714 very similar to the problem of modeling pronunciation
715 variation". However, although Ristad and Yianilos con-
716 sider the problem of pronunciation modelling in speech
717 technology, they do not consider alignment problems
718 as such.

719    This work represents the most comprehensive study
720 to date of letter-phoneme alignment, at the same time
721 achieving what is probably the best reported perfor-
722 mance on the difficult task of letter-phoneme conver-
723 sion of unknown words of English. Since the aligned
724 BEEP dictionary is a potentially valuable resource,
725 the version obtained from the NETspeak initialization
726 (which produced best performance on letter-phoneme
727 conversion) is made freely available for research
728 use at `http://festvox.org/packed/data/`
729 `damper`. Since our software has wide applicability, we
730 are also working to provide an on-line facility at which
731 researchers can submit dictionaries for alignment.

**References**                                                    732

Abercrombie, D. (1981). Extending the Roman alphabet: Some or-    733
    thographic experiments of the past four centuries. In R.E. Asher   734
    and E. Henderson (Eds.), *Towards a History of Phonetics*. Edin-   735
    burgh, UK: Edinburgh University Press, pp. 207–224.              736
Bagshaw, P.C. (1998). Phonemic transcription by analogy in text-  737
    to-speech synthesis: Novel word pronunciation and lexicon com-    738
    pression. *Computer Speech and Language*, *12*(2):119–142.        739
Baum, L.E. (1972). An inequality and associated maximization tech- 740
    nique in statistical estimation for probabilistic functions of Markov   741
    processes. In *Inequalities III: Proceedings of the Third Symposium   742
    on Inequalities*, Los Angeles, CA, pp. 1–8.                       743
Bellman, R. (1957). *Dynamic Programming*. Princeton, NJ: Prince-  744
    ton University Press.                                            745

Black, A.W., Lenzo, K., and Pagel, V. (1998). Issues in building general letter-to-sound rules. In *Proceedings of 3rd European Speech Communication Association (ESCA)/COCOSDA International Workshop on Speech Synthesis*. Jenolan Caves, Australia, pp. 77–80.

Carney, E. (1994). *A Survey of English Spelling*. London, UK: Routledge.

Chomsky, N. and Halle, M. (1968). *The Sound Pattern of English*. New York, NY: Harper and Row.

Coltheart, M. (1978). Lexical access in simple reading tasks. In G. Underwood (Ed.), *Strategies of Information Processing*. New York: Academic Press, pp. 151–216.

Coltheart, M. (1984). Writing systems and reading disorders. In L. Henderson (Ed.), *Orthographies and Reading: Perspectives from Cognitive Psychology, Neuropsychology and Linguistics*. London, UK: Lawrence Erlbaum Associates, pp. 67–79.

Damper, R.I. (Ed.) (2001). *Data-Driven Methods in Speech Synthesis*. Dordrecht, The Netherlands: Kluwer Academic Publishers.

Damper, R.I. and Eastmond, J.F.G. (1996). Pronouncing text by analogy. In *Proceedings of 16th International Conference on Computational Linguistics*. Copenhagen, Denmark, Vol 2, pp. 268–273.

Damper, R.I. and Eastmond, J.F.G. (1997). Pronunciation by analogy: Impact of implementational choices on performance. *Language and Speech*, 40(1):1–23.

Damper, R.I., Marchand, Y., Adamson, M.J., and Gustafson, K. (1999). Evaluating the pronunciation component of text-to-speech systems for English: A performance comparison of different approaches. *Computer Speech and Language*, 13(2):155–176.

Dedina, M.J. and Nusbaum, H.C. (1991). PRONOUNCE: A program for pronunciation by analogy. *Computer Speech and Language*, 5(1):55–64.

Dempster, A.P., Laird, N.M., and Rubin, D.B. (1977). Maximum-likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 39(1):1–38.

Federici, S., Pirrelli, V., and Yvon, F. (1995). Advances in analogy-based learning: False friends and exceptional items in pronunciation by paradigm-driven analogy. In *Proceedings of International Joint Conference on Artificial Intelligence (IJCAI'95) Workshop on New Approaches to Learning for Natural Language Processing*, Montreal, Canada, pp. 158–163.

Forney, G. D. (1973). The Viterbi algorithm. In *Proceedings of the IEEE*, vol. 61, no. 3, pp. 268–278.

Hartley, H. (1958). Maximum likelihood estimation from incomplete data. *Biometrics 14*, 174–194.

Henderson, L. (1984). Writing systems and reading processes. In L. Henderson (Ed.), *Orthographies and Reading: Perspectives from Cognitive Psychology, Neuropsychology and Linguistics*. London, UK: Lawrence Erlbaum Associates, pp. 11–24.

Hopcroft, J.E., Motwani, R., and Ullman, J.D. (2001). *Introduction to Automata Theory, Languages, and Computation*, 2nd ed. Boston, MA: Addison-Wesley.

Jansche, M. (2001). Re-engineering letter to sound rules. In *Second Meeting of the North American Chapter of the Association for Computational Linguistics, NAACL 2001*, Pittsburg, PA. Paper N01-1015 in on-line archive at http://acl.ldc.upenn.edu/N/N01/.

Katz, L. and Feldman, L.B. (1981). Linguistic coding in word recognition: comparisons between a deep and a shallow orthography. In A.M. Lesgold and C.A. Perfetti (Eds.), *Interactive Processes in Reading*. Hillsdale, NJ: Lawrence Erlbaum Associates, pp. 85–106.

Knill, K. and Young, S. (1997). Hidden Markov models in speech and language processing. See Young and Bloothooft (1997), pp. 27–68.

Kruskal, J.B. (1983). An overview of sequence comparison. In D. Sankoff and J. B. Kruskal (Eds.), *Time Warps, String Edits and Macromolecules: The Theory and Practice of Sequence Comparison*. Reading, MA: Addison-Wesley, pp. 11–44.

Lawrence, S.G.C. and Kaye, G. (1986). Alignment of phonemes with their corresponding orthography. *Computer Speech and Language*, 1(2):153–165.

Liberman, I., Liberman, A., Mattingly, I. and Shankweiler, D. (1980). Orthography and the beginning reader. In J. Kavanagh and R. Venezky (Eds.), *Orthography, Reading and Dyslexia*. Baltimore, OH: University Park Press, pp. 137–153.

Luk, R.W.P. and Damper, R.I. (1991). A novel approach to inferring letter-phoneme correspondences. In *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP'91*. Toronto, Canada, Vol 2, pp. 741–744.

Luk, R.W.P. and Damper, R.I. (1992). Inference of letter-phoneme correspondences by delimiting and dynamic time warping techniques. In *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP'92*. San Francisco, CA, Vol 2, pp. II.61–II.64.

Luk, R.W.P. and Damper, R.I. (1993). Inference of letter-phoneme correspondences with pre-defined consonant and vowel patterns. In *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing, ICASSP'93*. Minneapolis, MN, Vol 2, pp. II.203–II.206.

Luk, R.W.P. and Damper, R.I. (1996). Stochastic phonographic transduction for English. *Computer Speech and Language*, 10(2):133–153.

Luk, R.W.P. and Damper, R.I. (1998). Computational complexity of a fast Viterbi decoding algorithm for stochastic letter-phoneme transduction. *IEEE Transactions on Speech and Audio Processing*, 6(3):217–225.

Marchand, Y. and Damper, R.I. (2000). A multistrategy approach to improving pronunciation by analogy. *Computational Linguistics*, 26(2):195–219.

McCulloch, N., Bedworth, M., and Bridle, J. (1987). NETspeak—a re-implementation of NETtalk. *Computer Speech and Language*, 2(3/4):289–301.

McLachlan, G.J. and Krishnan, T. (1997). *The EM Algorithm and Extensions*. New York, NY: John Wiley.

Moon, T.K. (1996). The expectation-maximization algorithm. *IEEE Signal Processing Magazine*, 13(6):47–60.

Needleman, S.B. and Wunsch, C.D. (1970). An efficient method applicable to the search for similarities in the amino acid sequences of two proteins. *Journal of Molecular Biology*, 48(3):444–453.

Neuhoff, D.L. (1975). The Viterbi algorithm as an aid in text recognition. In *IEEE Transactions on Information Theory*, IT-21:222–226.

Parfitt, S.H. and Sharman, R.A. (1991). A bidirectional model of English pronunciation. In *Proceedings of 2nd European Conference on Speech Communication and Technology, Eurospeech'91*. Genova, Italy, Vol 2, pp. 800–804.

Pirrelli, V. and Federici, S. (1994). On the pronunciation of unknown words by analogy in text-to-speech systems. In *Proceedings of the*

162    *Damper et al.*

*Second Onomastica Research Colloquium*. London, UK, pp. 43–50.

Pirrelli, V. and Federici, S. (1995). You'd better say nothing than something wrong: Analogy, accuracy and text-to-speech applications. In *Proceedings of 4th European Conference on Speech Communication and Technology, Eurospeech'95*. Madrid, Spain, Vol 1, pp. 855–858.

Pirrelli, V. and Yvon, F. (1999). The hidden dimension: A paradigmatic view of data-driven NLP. *Journal of Experimental and Theoretical Artificial Intelligence*, 11(3):391–408.

Ristad, E.S. and Yianilos, P.M. (1998). Learning string-edit distance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(5):522–532.

Sampson, G. (1985). *Writing Systems*. London, UK: Hutchinson.

Sejnowski, T.J. and Rosenberg, C.R. (1987). Parallel networks that learn to pronounce English text. *Complex Systems*, 1(1):145–168.

Sullivan, K.P.H. (2001). Analogy, the corpus and pronunciation. See 2001, pp. 45–70.

Sullivan, K.P.H. and Damper, R.I. (1993). Novel-word pronunciation: A cross-language study. *Speech Communication*, 13(3–4):441–452.

Turvey, M.T., Feldman, L.B. and Lukatela, G. (1984). The Serbo-Croatian orthography constrains the reader to a phonologically analytic strategy. In L. Henderson (Ed.) *Orthographies and Reading: Perspectives from Cognitive Psychology, Neuropsychology and Linguistics*. London, UK: Lawrence Erlbaum Associates, pp. 81–89.

Venezky, R.L. (1965). *A Study of English Spelling-to-Sound Correspondences on Historical Principles*. Ann Arbor, MI: Ann Arbor Press.

Venezky, R.L. (1970). *The Structure of English Orthography*. The Hague, The Netherlands: Mouton.

Viterbi, A.J. (1967). Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE Transactions on Information Theory*, IT-13(2):260–269.

Wu, C.F.J. (1983). On the convergence properties of the EM algorithm. *Annals of Statistics*, 11(1):95–103.

Young, S. and G. Bloothooft (Eds.) (1997). *Corpus-Based Methods in Language and Speech Processing*. Dordrecht, The Netherlands: Kluwer Academic Publishers.

Yvon, F. (1996a). Grapheme-to-phoneme conversion using multiple unbounded overlapping chunks. In *Proceedings of Conference on New Methods in Natural Language Processing (NeMLaP-2'96)*. Ankara, Turkey, pp. 218–228.

Yvon, F. (1996b). *Prononcer par Analogie: Motivations, Formalisations et Évaluations*. PhD thesis, ENST, Paris, France.