## Alignment of Phonemes with their Corresponding Orthography

S.G.C.Lawrence and G.Kaye

Speech Research, IBM UK Scientific Centre,
Athelstan House, St Clement Street, WINCHESTER, SO23 9DR

### Introduction

The original motivation for developing this algorithm arose during the automated production of a computer-based dictionary of phonemic transcription based on the *Collins English Dictionary* [1], using a copy of the photo-typesetting tape kindly supplied by Wm. Collins and Co.. Collins's dictionary gives the phonemic transcription of all headwords, and follows the usual lexicographic conventions for representing inflected and derived forms within a headword entry, i.e. by giving partial spellings, and in some cases the pronunciation of the inflectional or derivational suffix. Where the inflected or derived form is considered to be regular no pronunciation is given and the user must apply rules given in the introduction to the dictionary to determine the pronunciation. The alignment algorithm has been used to combine the dictionaries pronunciation of a root with the pronunciation of the suffix (also given in the dictionary) to produce the pronunciation of the inflected or derived form, and for checking all pronunciations so derived. It is a mechanism for combining partial pronunciations and not a means of producing a phonemic transcription *ab initio*. Typical examples of the use of the algorithm in building the dictionary entries are:-

1. The entry for **go** is

   **go** (/gəʊ/) *vb* **goes, go+ing, went, gone**

   For our purposes pronunciations for **goes** and **going** are derived by applying the rules given in the introduction to the dictionary. The plausibility of the pronunciation so produced is checked by applying the algorithm to the result. If the pronunciation and the orthography cannot be aligned the pronunciation is considered invalid and must be supplied manually. The irregular forms **went** and **gone** appear as separate head words and so do not need to have their pronunciations derived.

2. The algorithm can be used to create complete phonemic transcriptions from partial transcriptions. For example the dictionary entry

   **aba+cus** /æbəkəs/ pl. **-ci** /-saɪ/

   The '+' denotes a syllable boundary, and in this case also marks the place at which **cus** is to be replaced by **ci**. There is no such marker in the phonemic transcription. It is necessary to know where the syllable boundary is in the transcription so that /kəs/ can be replaced by /saɪ/. The algorithm is used for this. After alignment **aba** is found to correspond with /æbə/ and **cus** with /kəs/. The syllable boundary is between /æbə/ and /kəs/. The end of the orthography can be changed to **ci** and the end of the phonemic transcription to /saɪ/ thus producing the generated form **abaci** /æbəsaɪ/.

ALIGNMENT OF PHONEMES WITH THEIR CORRESPONDING ORTHOGRAPHY

3. Where a phonemic transcription is given in *Collins English Dictionary* the lexical stress is indicated as part of it. Words without phonemic transcriptions have stress marked in the orthography. It is therefore necessary to transfer the stress marking to the generated phonemic transcription. The algorithm is used to check the plausibility of the generated phonemic transcription and also to determine the placement of the stress markers. Thus for the word **sug'gester** the phonemic transcription /sedʒesta/ was generated by appending /a/ to the transcription /sedʒest/ of **suggest**. The pronunciation and orthography of **sug'gest** are aligned by the algorithm as:

```
      s   u   g'g   e   s   t
 /    s   ə   dʒ    ɛ   s   t        /
```

this can be used to infer that the stress should be placed on /dʒ/ (using the convention that when the lexical stress in the orthography is between two gemminate consonants, it is placed before the corresponding consonant in the phonemic transcription), giving the derived dictionary entry **suggester** /sə'dʒestə/.

4. Where a dictionary entry contains alternative orthographies for a head word, but only one phonemic transcription, the algorithm is used to decide if the transcription is plausible for all orthographies. In *Collins English Dictionary* are the following entries.

   **abridgement** /əbrɪdʒmənt/ or **abridgment**

   **abutment** /əbʌtmənt/ or **abuttal**

The pronunciation of **abridgement** and **abridgment** are the same. However, **abuttal** is not pronounced /əbʌtment/, nor does any English word contain **al** pronounced as /mənt/. The algorithm indicates that the phonemic transcription of **abutment** is not plausible for **abuttal**, and so its pronunciation must be supplied manually.

## The Algorithm

The alignment algorithm is driven by a table of RP phonemes and their corresponding orthography. A different table is required for each accent of English. The present one was deduced with the aid of *Walker's Rhyming Dictionary of the English Language* [2], and Collins English Dictionary. The following is a small section of the table and gives example words. The complete table is given later (see "The Table of Correspondences").

| /-/ | ⇔ | e | have |
| /ɛ/ | ⇔ | e | when |
| /ʃ/ | ⇔ | e | before |
| /iː/ | ⇔ | e | be |
| /ə/ | ⇔ | e | government |
| /eɪ/ | ⇔ | e | crochet |
| /ɛə/ | ⇔ | e | whereas |
| /iˑ/ | ⇔ | e | recipe |
| /ɒ/ | ⇔ | e+n | gendarme |

The symbol ⇔ denotes *corresponds with.*

ALIGNMENT OF PHONEMES WITH THEIR CORRESPONDING ORTHOGRAPHY

When a correspondence occurs only in a specific following orthographic context, this is shown separated from the grapheme by a + (e.g. ə corresponds to /ɒ/ only in the context of a following n). The grapheme following the + must be separately aligned.

Where a grapheme can be considered to have no pronunciation then the symbol /-/ is used to represent the *silent phoneme*. The following multi-graphic phonemes are treated as a single unit and are each stored internally as a single code point:- ɪəʊ, eʊə, iː, ɪ, ɑː, ɔː, uː, əː, eɪ, aɪ, ɔɪ, eʊ, aʊ, ɛə.

To improve the robustness of the algorithm, an attempt has been made to minimise the number of *silent* graphemes. There are five cases in the table, and all have high frequencies of occurrence in the LOB corpus [3]. The most common of these graphemes is ə, occurring largely in a morph or word-final position. Except for these five, all other potentially *silent* graphemes have been combined with one or more adjacent graphemes.

For example:

| | | | |
|---|---|---|---|
| /m/ | ⇔ | mb | climb |
| /s/ | ⇔ | ps | psychology |
| /w/ | ⇔ | wh | which |
| /f/ | ⇔ | ph | phone |
| /ə/ | ⇔ | eur | amateur |
| /ʊə/ | ⇔ | ur | sure |
| /ɔː/ | ⇔ | ar | war |
| /əː/ | ⇔ | or | work |

Such combinations are usually well founded for one of the following reasons.

1. They follow a long established orthographic convention e.g., th ⇔ /θ,ð/, ch ⇔ /tʃ/, sh ⇔ /ʃ/.

2. They are phonetically motivated based on other accents of English. One example is wh ⇔ /w/ which in Scottish accents is realised as [ʍ] and may be represented phonemically as /hw/ or /ʍ/. Another example is post vocalic r which in RP usually results in lengthening of the vowel, whilst in rhotic accents, such as American English, produces rhoticization of the vowel. It was decided to append the r to the preceding vowel when the r does not precede a vowel phoneme as in award /əwɔːd/ hence a silent r cannot occur after a vowel. However when the r is intervocalic as in boring /bɔːrɪŋ/, the r corresponds to /r/.

Those grapheme clusters which cannot be justified for the above reasons are either idiosyncrasies of English spelling, or a consequence of a foreign loan word. They usually have a low frequency of occurrence.

The alignment of -tion presents some difficulties. As a matter of principle the alignment t ⇔ /ʃ/, i ⇔ /-/, o ⇔ /ə/, n ⇔ /n/ was excluded, leaving two other possibilities. The alignment chosen is t ⇔ /ʃ/, io ⇔ /ə/ instead of ti ⇔ /ʃ/, o ⇔ /ə/. For whilst the consonant alternation in word sequences such as permit, permission, permissive, is governed by the /t/ in the underlying phonemic form, it is less consistent for the purposes of this algorithm, to group ti rather than io, for word sequences such as the following:

ALIGNMENT OF PHONEMES WITH THEIR CORRESPONDING ORTHOGRAPHY

```
      p  er  m  i  t                    d  i  r  e  c  t
 /    p  ə   m  ı  t  /              /   d  ı  r  ɛ  k  t /
      p  er  m  i  ss  i  o  n          d  i  r  e  c  t  i  o  n
 /    p  ə   m  ı  ʃ   ə  n  /       /   d  ı  r  ɛ  k  ʃ  ə  n  /
      p  er  m  i  ss  i  v  e          d  i  r  e  c  t  i  v  e
 /    p  ə   m  ı  s   ı  v  -  /    /   d  ı  r  ɛ  k  t  ı  v  -  /
```

The algorithm matches the word and its phonemic transcription against the table of correspondences. This table is stored in the computer such that the longest grapheme clusters are scanned first. The initial objective is to match the largest number of graphemes that produces the corresponding phonemes. When a match is found a space is inserted after the last of the matched graphemes and also their corresponding phonemes. If a match is not found then a scan of the graphemes which correspond to the silent phoneme is performed. If a correspondence is found then a '-' is inserted into the phonemic transcription string. If no correspondence is found then there is either an error in the table or in the phonemic transcription and therefore no alignment is possible. If the correspondence consists of multiple graphemes or multiple phonemes then a note is kept of the position in the word and phonemic transcription where the cluster occurs. If, during the left to right scan of the word and transcription, no correspondence can be found and there has been a previous correspondence with more than one grapheme or phoneme then the length of the matched graphemic or phonemic cluster is decreased by one, and a search is made of the correspondence table with the length of the matching string limited to the length of the shortened cluster.

If the cluster is reduced to a single grapheme or phoneme and no match found, then there is either an error in the correspondence table or the phonemic transcription.

The following example illustrates the algorithm. For **create** ⇔ /kri:eɪt/, the first attempt at aligning the orthography with the corresponding phonemes produces

| | | | |
|---|---|---|---|
| /k/ | ⇔ | c | |
| /r/ | ⇔ | r | |
| /i:/ | ⇔ | ea | eg. **each** |
| /-/ | ⇔ | t | eg. **often** |
| /eɪ/ | ⇔ | e | eg. **crochet** |

with

/t/ unaligned

The **ea** cluster is split and a second attempt to match the graphemes with the phonemes is attempted starting with the **e** that follows the **r**. Although /eɪ/ is a cluster that could be split, it is not because the RP phoneme /e/ occurs only in the diphthong /eɪ/. The splitting of a cluster into non-RP phonemes is prevented by using the input conversion table described above that causes such clusters to be treated as single entities.

The second attempt produces

| /k/ | ⇔ | c | |
|-----|---|---|---|
| /r/ | ⇔ | r | |
| /iː/ | ⇔ | e | eg. be. |
| /eɪ/ | ⇔ | a | eg. made. |
| /t/ | ⇔ | t | |
| /-/ | ⇔ | e | eg. have. |

which gives

```
    c   r   e    a    t    e
  / k   r   iː   eɪ   t   -/
```

## Testing the Algorithm

The algorithm has been tested against 32,049 words of a phonemically tagged version of the *Lancaster-Oslo/Bergen Corpus* (see "The Table of Correspondences").

31,773 words were aligned without error.

276 words could not be aligned because they fell into one of the following categories.

* Words containing diacritics eg. **attaché**. This could be avoided if the appropriate correspondence were to be added to the table.

* Words given in Collins with transcriptions containing non-RP phonemes eg. **Truffaut** /tryfo/. It was decided to keep a table of only RP phoneme-to-grapheme correspondences.

* Partial transcriptions. eg. **Sunbury-on-Thames** /sʌnbəriˈ/. It seems odd that the dictionary entry does not give the pronunciation of **Thames** as well as that of **Sunbury**.

## Alignment Examples

The following examples illustrate the performance of the algorithm.

```
a d v e n t  ur e          b ɔr s t a l           e ss e x
/ e d v c n tʃ e - /        / b ɔː s t e l /        / c s ɪ ks /
a d v i  s e                c a m e r a             e  v e
/ e d v ɑ z - /             / k æ m ə r e /         /iː v - /
a ff e c t i o n            c o n v e n t i o n      e x  a m i  n e
/ ə f  c kʃ e n /            / k e n v c nʃ e n /     /ɪ gz æ mɪ  n - /
a n g l i c a n             c y c l e               f o n d
/ æ ŋ g l ɪ k e n /          / s aɪ k ǀ - /          /f ɒ n d/
a s l ee p                  d i r e c t i o n s      h i  th e r t o
/ e s l iː p /              / d ɪ r c kʃ e n z /      / h ɪ ꞵ e t uː /
a wer d                     d i r t y               h o r i  z ɒn
/ e w ɔː d /                / d ɜː t i· /            / h e r ɑ z ŋ /
b a r e                     d i sh                  h u rr y
/ b cə - /                  / d ɪ ʃ /               / h ʌ r i· /
```

```
i d e n t i c a l          p l a t f o r m            s u r r e n d e r
/aɪ d e n t ɪ k |  /        /p l æ t f ɔ: m/           /s ə r e n d ə  /
l e ss on                  p o r t u gue s e           u l t i m a t e l y
/l e s n  /                /p ɔ: t j u g i: z - /      /ʌ l t ɪ mɪ t - l iˑ  /
mea t                      p o ss e ss i o n           v a gu e
/mi: t  /                  /p ə z ɛ ʃ e n/             /v eɪ g - /
n o t a b l e              p r a i s e                 v i s i t o r
/n əʊ t e b |  - /         /p r eɪ z - /               /v ɪ z i t ə  /
n o u n                    p r e s t i ge              v o c a l
/n aʊ n /                  /p r ɛ s t i: ʒ - /         /v eʊ k |  /
o b s er v er              s e tt l e me n t           wer e n ' t
/ə b z ɜ: v ə  /           /s ɛ t l | - me n t /       /w ɜ: - n - t /
p e n al t y               s p a r e                   wor th y
/p ɛ n | t iˑ /            /s p eə - /                 /w ɜ: ð iˑ /
```

## The Table of Correspondences

The table that follows shows the phoneme-to-grapheme correspondences together with some statistics gathered by using the algorithm on a lexicon. Only part of the table is given. The complete table will be published at a later date. The lexicon used is those 32.049 words in the *LOB Corpus* that are not homographs and whose pronunciation has been derived from *Collins English Dictionary*. The frequencies of the words are taken from the *LOB Corpus*. The table is in the order in which it is scanned when matching the orthography with the phonemic transcriptions.

The **Phone** column shows the phonemes that correspond to the graphemes in the column entitled **Graph**. Characters other than IPA symbols indicate the following :

- The graphemes in the **Graph** column do not correspond with a phoneme (i.e. are silent).

* If the grapheme **r** follows the graphemes in the **Graph** column in a word and the phoneme /r/ does not follow the phonemes in the **Phone** column in the transcription then the grapheme **r** forms part of the correspondence.

In the **Graph** column, characters other than graphemes are used to indicate a condition which has to be satisfied for there to be a correspondence. The characters and their meaning are as follows :

+ The graphemes following the ' + ' must follow the graphemes that precede the ' + ' for the correspondence to occur.

> The graphemes following the ' > ' must not follow the graphemes that precede the ' > ' for the correspondence to occur.

- The graphemes following the '-' must precede the graphemes that precede the '-' for the correspondence to occur.

< The graphemes following the ' < ' must not precede the graphemes that precede the ' < ' for the correspondence to occur.

# Proceedings of The Institute of Acoustics

ALIGNMENT OF PHONEMES WITH THEIR CORRESPONDING ORTHOGRAPHY

. May be substituted for one of the characters . , / ' or a word end.

The **Example** column gives the most frequently occurring word in the lexicon in which the correspondence occurs. The **# Text** column gives the cumulative sum of the frequencies of occurrence of those words in the corpus (LOB), from which the lexicon was produced, in which the correspondence occurs. If the correspondence does not occur in the lexicon (i.e. 0 in this column), it is listed because it occurs in words with RP phonemic transcriptions in *Collins English Dictionary*. The **# Lex** column gives the number of times the correspondence occurred in the lexicon.

| Phone | Graph | Example | # Text | # Lex | Phone | Graph | Example | # Text | # Lex |
|---|---|---|---|---|---|---|---|---|---|
| - | ʔed | often | 579 | 32 | ɑː | e+m | contretemps | 1 | 1 |
| - | s+. | cannes | 0 | 0 | ɑː | e+n | genre | 4 | 1 |
| - | r-r | referred | 130 | 8 | ɑː˙ | e | sergeant | 58 | 10 |
| - | h | john | 1748 | 174 | ɑː | au | laughed | 213 | 17 |
| - | e>r | have | 139867 | 7594 | ɑː | ard | boulevards | 1 | 1 |
| aɪə˙ | y | byron | 65 | 8 | ɑː | al | half | 483 | 32 |
| aɪə˙ | oi | choir | 10 | 3 | ɑː | agh | omagh | 0 | 0 |
| aɪə | Iro+n | iron | 79 | 7 | ɑː˙ | ae | stranraer | 6 | 2 |
| aɪə˙ | ie | society | 800 | 45 | ɑː˙ | aa | afrikaans | 15 | 7 |
| aɪə˙ | ia | trial | 546 | 86 | ɑː˙ | a | are | 24710 | 1548 |
| aɪə˙ | i+r | required | 1627 | 126 | biː | b+. | b | 378 | 22 |
| aɪ | y | by | 10824 | 331 | b | pb | cupboard | 22 | 4 |
| aɪ | uy | buy | 148 | 8 | b | bb | abbey | 410 | 144 |
| aɪ | igh | might | 4335 | 181 | b | b | by | 54975 | 4016 |
| aɪ | ia | vial | 0 | 0 | diː | d+. | d | 191 | 25 |
| aɪ | i | i | 34269 | 2096 | dʌbljuː | w+. | w. | 96 | 10 |
| aɪ | ey | eyes | 550 | 12 | d | ld | could | 0 | 0 |
| aɪ | eigh | height | 83 | 4 | d | dh | oldham | 14 | 2 |
| aɪ | ei | either | 544 | 40 | d | dd | added | 1674 | 154 |
| aɪ | ay | ramayana | 11 | 5 | d | d | and | 153411 | 8977 |
| aɪ | ai | cairo | 38 | 17 | dʒeɪ | j+. | j. | 145 | 8 |
| aɪ | ae | maestro | 2 | 1 | dʒiː | g+. | g. | 207 | 20 |
| aɪ | a | bandaranaike | 2 | 1 | dʒ | jj | hajji | 0 | 0 |
| aʊə˙ | ow | dowry | 1 | 1 | dʒ | j | just | 4912 | 405 |
| aʊə˙ | ou | our | 1654 | 18 | dʒ | gi | religious | 378 | 36 |
| aʊə˙ | au | gaur | 0 | 0 | dʒ | gg | suggested | 402 | 15 |
| aʊ | ow | now | 5940 | 184 | dʒ | g | general | 11034 | 1028 |
| aʊ | ough | plough | 25 | 7 | dʒ | dj | adjustment | 129 | 27 |
| aʊ | ou | out | 10918 | 464 | dʒ | dg | knowledge | 1057 | 123 |
| aʊ | eo | macleod | 49 | 2 | dʒ | d | soldiers | 197 | 11 |
| aʊ | aw | haworth | 1 | 1 | dʒ | ch | norwich | 23 | 5 |
| aʊ | au | strauss | 29 | 11 | eɪtʃ | h+. | h. | 161 | 13 |
| aʊ | aou | caoutchouc | 0 | 0 | eɪ | ez | laissez | 0 | 0 |
| aʊ | ao | laos | 26 | 3 | eɪ | eyo | eyot | 0 | 0 |
| æ | i | meringue | 2 | 2 | eɪ | ey | they | 3900 | 44 |
| æ | ai | laing's | 4 | 4 | eɪ | es | demesnes | 1 | 1 |
| æ | ach | drachm | 0 | 0 | eɪ | eigh | eight | 444 | 38 |
| æ | aa | naali | 0 | 0 | eɪ | ei | reign | 125 | 30 |
| æ | a'a | ma'am | 5 | 1 | eɪ | ee | beethoven | 27 | 7 |
| æ | a | and | 94264 | 4539 | eɪ | ed | passepied | 0 | 0 |
| ɔːr | r+. | r.m.s. | 17 | 9 | eɪ | eagh | castlereagh | 0 | 0 |
| ɔː | ɪ+. | r. | 169 | 11 | eɪ | ea | great | 1169 | 21 |
| ɑː | ia | acciaccatura | 0 | 0 | eɪ˙ | e | crochet | 200 | 56 |
| ɑː˙ | ea | heart | 303 | 25 | | | | | |

ALIGNMENT OF PHONEMES WITH THEIR CORRESPONDING ORTHOGRAPHY

| Phone | Graph | Example | # Text | # Lex | | Phone | Graph | Example | # Text | # Lex |
|---|---|---|---|---|---|---|---|---|---|---|
| eɪ | ay | may | 8574 | 284 | | eɪ | aigh | straight | 134 | 7 |
| eɪ | au | gauge | 22 | 4 | | eɪ | ai | main | 6515 | 481 |
| eɪ | ao | gaol | 17 | 3 | | eɪ | ae | phaeton | 8 | 5 |
| eɪ | ais | renaissance | 17 | 5 | | eɪ | a | made | 33492 | 2930 |

## Summary

The aligned forms of words and their phonemic transcriptions, together with statistics on frequency of rule use, are seen as invaluable aids in developing a set of text-to-phoneme rules, where the ordering of rules, based on frequency of occurrence is critical. The aligned forms may also be useful in speech recognition systems.

The table of correspondences may be used in the development of aids for the teaching of English spelling. From the table can be deduced the most likely orthography for a given sound.

## References

[1] 'Collins English Dictionary', William Collins, Glasgow, (1979).

[2] J. Walker, 'The Rhyming Dictionary of the English Language', Routledge & Keegan Paul, London, (1924).

[3] 'LOB (Lancaster-Oslo/Bergen) Corpus'. Department of English, University of Oslo, Oslo, (1978).

## Acknowledgements