# A Probabilistic Approach to Pronunciation by Analogy

Janne V. Kujala*     Aleksi Keurulainen†

September 21, 2011

### Abstract

The relationship between written and spoken words is convoluted in languages with a deep orthography such as English and therefore it is difficult to devise explicit rules for generating the pronunciations for unseen words. Pronunciation by analogy (PbA) is a data-driven method of constructing pronunciations for novel words from concatenated segments of known words and their pronunciations. PbA performs relatively well with English and outperforms several other proposed methods. However, the best published word accuracy of 65.5% (for the 20,000 word NETtalk corpus) suggests there is much room for improvement in it.

Previous PbA algorithms have used several different scoring strategies such as the product of the frequencies of the component pronunciations of the segments, or the number of different segmentations that yield the same pronunciation, and different combinations of these methods, to evaluate the candidate pronunciations. In this article, we instead propose to use a probabilistically justified scoring rule. We show that this principled approach alone yields better accuracy (66.21% for the NETtalk corpus) than any previously published PbA algorithm. Furthermore, combined with certain ad hoc modifications motivated by earlier algorithms, the performance climbs up to 66.6%, and further improvements are possible by combining this method with other methods.

## 1 Introduction

Pronunciation by analogy (PbA) was originally proposed by Glushko (1979) as a psychological model of how humans pronounce words. He suggested that instead of retrieving a single pronunciation from memory for known words and using abstract spelling-to-sound rules for pseudowords, humans pronounce both words and pseudowords using a similar process based on analogy to similar words.

The first concrete definition of PbA was given by the PRONOUNCE computer program of Dedina and Nusbaum (1991). The program uses four components: 1) a *lexical database* to store the pronunciations of known words, 2) a *matcher* to list segments of the input word that match segments of known words in the lexical database, a 3) *pronunciation lattice* to arrange the segments and their pronunciations within the lexical database into a graph-like data structure representing possible pronunciations of the input word, and 4) a *decision function* to select the best pronunciation for the input word among the different candidate pronunciations. The same basic structure has been used in PbA implementations since then, but the different components have gone through various improvements.

---

*Address: Department of Mathematical Information Technology, University of Jyväskylä, Finland. Email address: `jvk@iki.fi`

†Address: Agora Center, University of Jyväskylä, Finland

PbA is currently one of the most accurate methods for automatically generating the pronunciations of unseen words. However, all PbA algorithms yield several candidate pronunciations and the selection of the best one is generally based on ad hoc methods. Indeed, Damper and Marchand (2006) state that little or no theory exists to guide the selection and the desire for a more probabilistic formulation was already apparent in Damper and Eastmond (1997) and Damper and Marchand (1998). In the lack of a theoretically preferable approach, an important idea leading to the current state of the art is the multi-strategy approach of Marchand and Damper (2000) which combines several evaluation methods in the decision function to squeeze out the best accuracy scores (Damper and Marchand, 2006; Polyákova and Bonafonte, 2008, 2009).

In this article, we instead propose to use a probabilistically justified evaluation function for choosing the best candidate pronunciation. We show that a single, principled method can yield better results than any of the previously published algorithms. Furthermore, our analysis of the best previously published single strategy, the PFSP method of Polyákova and Bonafonte (2008), suggests the use of what we call "magic root" in this article, a certain power function applied in the evaluation function. Although it is difficult to justify this theoretically, we show that it can be used to further improve the accuracy of our otherwise principled algorithm.

In the rest of this section, we describe the four basic components of a PbA system in more detail, review the relevant literature, and present the results of our reimplementation of the best published PbA algorithms. In Section 2, we present and evaluate our simple probabilistically justified evaluation method using non-overlapping segments; in Section 3, we generalize the method for overlapping segments and consider the "magic root" mentioned above for improving performance. Then, we conclude.

## 1.1 Lexical Database

The lexical database is assumed to consist of a list of *aligned* word-pronunciation pairs $(x, y) = (x_1 \ldots x_l, y_1 \ldots y_l)$ so that there is a one-to-one correspondence of the letters (or graphemes) $x_i$ to the phonemes $y_i$. Thus, the English word *that* should ideally be aligned as $(x_1, x_2, x_3) = (th, a, t)$ and $(y_1, y_2, y_3) = (ð, æ, t)$, where we have used the IPA phonetic symbols for the phonemes. To facilitate direct comparison to other PbA work, we have for now used the NETtalk corpus which inserts null phonemes /-/ (e.g., *th* is transcribed as /ð-/) and represents certain combinations as single phonemes (e.g., /ks/ replaced by /X/) so as to yield single-letter-to-phoneme alignment. However, all the algorithms can be naturally generalized to use many-to-many alignment.

## 1.2 Matcher

Dedina and Nusbaum (1991) defined the matcher as follows. Given an input word, all the words in the lexical database are iterated through. For each lexical entry, the matching procedure starts with the input word and the lexical entry left-aligned and then slides the shorter word to the right one letter at a time until the two words are right aligned. At each step, all substrings of letter positions where the two strings match are extracted together with the corresponding pronunciation of the lexical entry and stored in the pronunciation lattice (described below).

Damper and Marchand (1998) generalized this procedure to so called *full matching*, where the matching is done over the range of all possible overlaps of the two strings so that all common substrings are extracted regardless of their positions within the words. This matching procedure is in fact equivalent to the approach described by Sullivan and Damper (1993), where the frequencies of all substrings of the lexicon together with their pronunciations are first enumerated and stored in a dictionary, and then, for each input word,

the frequencies of its substrings are looked up in the precomputed dictionary to form the pronunciation lattice.

Dedina and Nusbaum (1991) and Damper and Marchand (1998) match the beginning of a word only to beginnings of other words and the end of a word only to the ends of other words. Sullivan and Damper (1993) investigate the effect of this restriction. In their approach, this amounts to storing a word as *#word#*, padded with silent beginning and end characters # that match the beginnings and ends of other words. In most cases, they obtained the best results by using the word boundary characters (see the reference for details) and this result was also confirmed by our experiments.

## 1.3    Pronunciation Lattice

In Dedina and Nusbaum (1991) and Damper and Marchand (1998), the pronunciation of a word is always constructed from segments that overlap by one letter and have matching pronunciations for the overlap letter. The nodes of the pronunciation lattice represent the overlap letters, identified by their pronunciation and position in the target word. In addition, special Start and End nodes represent the beginning and end of the word (corresponding to the silent word boundary null phonemes). Each arc represents a matched segment of the target word between and including the overlap letters. The arc is labeled by the pronunciation of the segment in the matched word and the frequency of this pronunciation of the segment in all matched words. Each path from Start to End represents one possible pronunciation for the word.

Yvon (1996) generalized the one character overlap to an unbounded overlap, with the nodes of the pronunciation lattice given by all matched substrings of the input word identified by their letter positions and pronunciations in matched words. This results in pronunciations given by several multiply overlapping chunks, for example, a five letter word could be segmented into three segments at letter positions 1–3, 2–4, and 3–5 with all three segments overlapping the letter position 3.

Sullivan and Damper (1993) define the pronunciation lattice differently, with the nodes representing the junctures between letters. In this case, any path from Start to End represents one possible pronunciation consisting of non-overlapping segments. They label the arcs with probability-like "preference values" instead of plain frequencies.

## 1.4    Decision Function

In their original article, Dedina and Nusbaum (1991) use a scoring heuristic that chooses the shortest path (i.e., the pronunciation with the smallest number of segments) and in case of a tie, chooses the path with the greatest sum of the arc frequencies.

Damper and Eastmond (1997) also use shortest paths as the primary heuristic but argue that a product of the frequencies is intuitively better as it is closer to a probabilistic formulation. As a secondary heuristic they use the maximum product of arc frequencies. Sullivan and Damper (1993) also use the maximum product of preference values over all paths as one of their heuristics.

Damper and Eastmond (1997) define a single heuristic strategy called "total product" (TP'), which sums the products of arc frequencies over *all* paths that yield the same pronunciation. Damper and Marchand (1998) consider a variation of the total product rule, summing the products of arc frequencies over all *shortest* paths yielding the same pronunciation (TP). Based on the published accuracy scores, the results for this variation appear to be somewhat better than for summing over all paths. Damper and Marchand (1998) also consider another variant called "weighted total product" (WTP), where the product of arc frequencies is divided by the product of segments lengths before summing, which was found to yield a slight improvement.

In his method (SMPA) using unbounded overlap, Yvon (1996) replaced the shortest-paths primary heuristic with a one taking into account the varying overlap. His primary heuristic is equivalent to choosing paths with the maximum average length of the segments. The secondary heuristic uses the frequency data to break ties as in Dedina and Nusbaum (1991) (we assume this means the sum of absolute frequencies).

In their seminal article, Marchand and Damper (2000) proposed a method for combining different strategies. They use the shortest paths as the primary heuristic but define several different secondary strategies. First, they rank the candidates using each of the component strategies and award points to the candidates according to a certain formula that gives the most points to the highest ranking candidate and a decreasing number of points for each lower ranking candidate. The scores of the component strategies are then combined by taking their sum or product. The product rule was found to yield slightly better results than the sum rule.

## 1.5 Reimplementation of Best Published Algorithms

A common problem well-known in the PbA literature is the difficulty of replicating published results. Differences in implementation details typically yield different results for reimplemented algorithms. Therefore, we have reimplemented the state of the art PbA algorithms described in the literature to facilitate fair comparisons to our new methods.

Several ideas have been proposed in the literature, but so far the best published PbA results are obtained from segmentations with overlap of one and the shortest-paths primary heuristic with various secondary heuristics. As a starting point, Table 1 shows the lower and upper bounds for accuracy for the NETtalk corpus given the primary heuristic of shortest paths.

In all published PbA works, the best performance has been obtain by combinations of the following component strategies for evaluating the candidate pronunciations as a secondary heuristic:

1. Maximum product of arc frequencies (PF)

2. Minimum standard deviation of arc lengths (SDPS)

3. Maximum number of candidates with the same pronunciation (FSP)

4. Minimum sum of the Hamming distances[1] of the candidate pronunciation from other candidates (NDS)

5. Maximum smallest arc frequency of the path (WL)

6. Maximum weighted product of arc frequencies (WPF): weighting by dividing each arc frequency by the number of different pronunciations for the letter sequence of the arc

7. Maximum frequency of the first arc of the path (SF)

8. Maximum frequency of the last arc of the path (SL)

9. Maximum length of the longest arc of the path and in case of a tie, maximum frequency of the maximum frequency arc (SLN)[2]

---

[1]The Hamming distance is the number of positions where two strings of equal length differ from each other.

[2]This definition appears to yield results closer to Polyákova and Bonafonte (2008) although their description seems to literally indicate maximum frequency of the longest arc as the secondary heuristic.

|  | Text-to-Speech | | Speech-to-Text | |
|---|---|---|---|---|
|  | Words (%) | Phones (%) | Words (%) | Letters (%) |
| Lower bound | 44.25 | (84.39) | 59.25 | (91.62) |
| Upper bound | 85.42 | (96.68) | 89.76 | (98.08) |

Table 1: Lower and upper bounds for accuracy given the primary heuristic of shortest paths and *overlap of one letter*. The lower bound is obtained by unit score, which corresponds to tie-braking by random guessing. The upper bound is given by the theoretically best possible secondary heuristic, which is obtained by minimizing the Levenshtein edit distance to the correct pronunciation (this implies that the correct pronunciation is chosen when available and otherwise the phoneme error is minimized).

10. Maximum sum over letter positions of the number of other candidates having the same phoneme multiplied by the frequency of the arc containing the phoneme (SSPF)[3]

11. Maximum sum of the geometric means of the products of arc frequencies over all candidates with the same pronunciation (PFSP)

Strategies 1–5 were proposed by Marchand and Damper (2000) who evaluated all $2^5 - 1 = 31$ combinations of the first five strategies; the best result was obtained as a combination of all five. Polyákova and Bonafonte (2008, 2009) extended the set by strategies 6–11 and evaluated all $2^{11} - 1 = 2047$ combinations. The best results were obtained by certain combinations not including all strategies. Table 2 shows the accuracy results of our reimplementation of these methods.

## 1.6   Evaluation Methodology

As in many previous works, we use the so called leave-one-out method of evaluation. That is, the accuracy of the algorithm for a given word is evaluated by teaching all other words in the corpus and then testing the generated pronunciation for the excluded word. After generating the pronunciation, we strip out all null phonemes from it as well as from the correct pronunciation before evaluation. This improves performance as sometimes the algorithms generate the correct pronunciation but with null phonemes at different places.[4]

The *word accuracy* of the algorithm is computed as the proportion of correct generated pronunciations over all words in the corpus. If for a given word there are several candidate pronunciations with the same score, we avoid arbitrary tie-braking by computing the accuracy of the word as the number of correct candidates with the best score divided by the total number of candidates with the best score.

The *phoneme accuracy* is computed by taking the Levenshtein edit distance[5] of the generated pronunciation from the correct pronunciation and dividing it by the number of phonemes in the correct pronunciation (not counting any silent null phonemes). If there are several candidates with the best score, then, as for word accuracy, the phoneme accuracy is averaged over all of them.

When using overlapping segments, certain words are left without pronunciation. This so called *silence problem* occurs for example for the word *anecdote*, which does not get a pronunciation as the substring *cd* does not occur in any other word. We solve the silence

---

[3]It is unclear how to define the frequency of the containing arc for an overlap letter as it belongs to two arcs (or no arcs, depending on the definition). We have included overlap letters twice, once for each arc.

[4]The stripping of null-phonemes is not done (or at least not reported) in the literature, but it seems like a more logical way of evaluating the algorithms as the null phonemes are only used for technical purposes.

[5]The Levenshtein distance is the minimum number of edits (changing, adding, or removing one character) required to change one string to the other.

| Evaluation method | | Text-to-Speech | | Speech-to-Text | |
|---|---|---|---|---|---|
| Combination | Name | Words (%) | Phones (%) | Words (%) | Letters (%) |
| 10000 | PF | 59.31 | (89.31) | **73.91** | **(94.78)** |
| 01000 | SDPS | 46.38 | (84.76) | 60.58 | (91.69) |
| 00100 | FSP | **62.46** | (89.90) | 72.41 | (94.42) |
| 00010 | NDS | 60.81 | **(89.96)** | 71.83 | (94.41) |
| 00001 | WL | 55.91 | (88.48) | 68.53 | (93.81) |
| 10101 | | 65.46 | (91.02) | **75.90** | (95.18) |
| 11111 | | **65.62** | **(91.14)** | 75.52 | (95.10) |
| 00000100000 | WPF | 59.25 | (89.18) | 74.06 | (94.82) |
| 00000010000 | SF | 58.28 | (88.84) | 68.00 | (93.67) |
| 00000001000 | SL | 51.16 | (86.84) | 69.58 | (93.80) |
| 00000000100 | SLN | 55.83 | (88.07) | 70.45 | (94.15) |
| 00000000010 | SSPF | 62.02 | (90.17) | 73.43 | (94.69) |
| 00000000001 | PFSP | **64.00** | **(90.57)** | **75.77** | **(95.21)** |
| . . . | | | | | |
| 00101100011 | | 65.55 | (91.12) | 76.09 | (95.24) |
| 00111111001 | | 65.46 | (91.11) | 76.10 | (95.23) |
| 10111111001 | | 64.94 | (90.98) | 76.12 | (95.23) |
| 00101001011 | | 65.31 | (91.03) | 76.12 | (95.22) |
| 00101100001 | | 65.35 | (90.99) | **76.19** | **(95.26)** |
| . . . | | | | | |
| 00111010011 | | 66.05 | **(91.25)** | 75.71 | (95.15) |
| 01101000001 | | 66.12 | (91.18) | 75.91 | (95.17) |
| 01100010001 | | 66.13 | (91.22) | 75.59 | (95.12) |
| 00100010001 | | 66.14 | (91.18) | 75.55 | (95.13) |
| 00101000001 | | **66.14** | (91.20) | 76.00 | (95.22) |
| – | TP | 61.76 | (89.96) | 74.91 | (94.99) |
| – | WTP | 62.02 | (90.04) | **74.97** | **(95.01)** |
| – | SMPA | **64.59** | **(90.55)** | 74.95 | (94.95) |

Table 2: Accuracy scores of our reimplementation of the best published PbA methods. First are shown the five strategies of Marchand and Damper (2000) followed by the best combinations for speech-to-text and text-to-speech. Then are shown the six additional strategies of Polyákova and Bonafonte (2008,2009) followed by the best five combinations for speech-to-text and text-to-speech. For comparison, the last three rows also show the performance of our reimplementation of the total product and weighted total product methods of Damper and Marchand (1998) and the SMPA method of Yvon (1996).

problem when it occurs by allowing one segment boundary without overlap at any position and then restricting the resulting segmentations (e.g., by requiring minimal number of segments) as usual. This solution gives a pronunciation for all words in the NETtalk corpus, but for the speech-to-text direction, leaves 12 words (0.06%) without spelling.[6]

In all our experiments, we use the NETtalk corpus, which is the most widely used corpus in PbA work (Damper and Eastmond, 1996, 1997; Damper and Marchand, 1998; Marchand and Damper, 2000; Yvon, 1996; Damper and Marchand, 2006; Soonklang et al., 2008; Polyákova and Bonafonte, 2008, 2009). It is based on a 20,009-word corpus from the 1974 edition of Webster's Pocket Dictionary, and was used in the NETtalk project in 1987 to train a neural network. The lexicon is manually aligned by Sejnowski and Rosenberg (1987) on a one-to-one basis. As in (Marchand and Damper, 2000), we removed all homophones (or homographs for the speech-to-text direction) and one-letter words from the corpus, leaving 19,595 words (or 19,545 words for the speech-to-text direction).[7] The speech-to-text direction is evaluated analogously to the text-to-speech direction, we simply swap the words and their pronunciations.

# 2 A Probabilistic Approach to Pronunciation by Analogy

In this section, we describe our simple probabilistic approach to PbA. We present first a method based on the pronunciation lattice of Sullivan and Damper (1993), which yields candidate pronunciations with non-overlapping segments, as it is more straightforward to define probabilities in this case. In Section 3, for faster computation and better results, we generalize the method to overlapping segments as used by among others Dedina and Nusbaum (1991) and Damper and Marchand (1998).

## 2.1 Training Data

We assume that a lexical database is given, consisting of aligned pairs $(x, y) = (x_1 \ldots x_l, y_1 \ldots y_l)$ of words and their pronunciations. Each word is assumed to be padded with the silent space characters $\#$ at both ends.

We consider all substrings $(x_a \ldots x_b, y_a \ldots y_b)$, $1 \leq a \leq b \leq l$, of every word-pronunciation pair $(x, y)$ in the corpus and compute the frequencies of these substrings over the entire list of words. For each possible written substring $x$, the frequencies of different corresponding pronunciations $y$ are then normalized as

$$\hat{p}(y \mid x) = \frac{\text{(number of times } x \text{ is pronounced as } y \text{ in the training data)}}{\text{(number of times } x \text{ appears in the training data)} + 1} \tag{1}$$

to yield a probability distribution $\hat{p}(y \mid x)$. This is the training data for our PbA algorithm. We add one to the denominator to take into account the fact that it is possible that the segment $x$ has a novel pronunciation in the unseen word. This is not necessary, but it is probabilistically more consistent and slightly increases the accuracy of our algorithm.

---

[6]We could generalize to allow two boundaries without overlap if we cannot get a pronunciation with one, but as we are mostly interested in the text-to-speech direction, we leave it for further work.

[7]The numbers of words quoted by Marchand and Damper (2000) are one less than ours, which we assume is due to an error (last word listed twice) in the version of the lexicon cited by them. We use the lexicon in the version that is freely downloadable from http://archive.ics.uci.edu/ml/datasets/Connectionist+Bench+(Nettalk+Corpus)

## 2.2 Segmentations of a Word

Given a written word $x$, we consider different *segmentations* $s = (s_1, \ldots, s_n)$ of the word, that is, lists of segment lengths that sum to the length $l$ of the word. For any given segmentation $s$ of the word $x$, we can compute a probability distribution of possible pronunciations as

$$p(y \mid x, s) = \prod_{k=1}^{n} \hat{p}(y_k \mid x_k), \tag{2}$$

where $x_k$ and $y_k$ denote the segments of the written word and its pronunciation.

The choice of the segmentation greatly affects the resulting distribution of its possible pronunciations. For example, for the string $l$, the distribution of its pronunciations in the training data might by given by the probability function[8]

$$
\begin{aligned}
/l/ &\mapsto 0.94 \\
/\ / &\mapsto 0.06 \text{ (silent)}.
\end{aligned}
$$

For the string $y$, the distribution might be

$$
\begin{aligned}
/\imath/ &\mapsto 0.78 \\
/a\imath/ &\mapsto 0.22.
\end{aligned}
$$

Thus, the string $ly$ segmented in two pieces would have the following distribution of pronunciations:

$$
\begin{aligned}
/l\imath/ &\mapsto 0.94 \cdot 0.78 = 0.7332, \\
/la\imath/ &\mapsto 0.94 \cdot 0.22 = 0.2068, \\
/\ \imath/ &\mapsto 0.06 \cdot 0.78 = 0.0468, \\
/\ a\imath/ &\mapsto 0.06 \cdot 0.22 = 0.0132.
\end{aligned}
$$

However, if one considers $ly$ as one segment, the statistics might be

$$
\begin{aligned}
/l\imath/ &\mapsto 0.96 \\
/la\imath/ &\mapsto 0.04.
\end{aligned}
$$

Thus, the distributions for segmentations with a small number of segments will generally be less spread out than the ones with a higher number of segments (assuming, of course, that there are high order correlations between the letters of the words of the language).

## 2.3 Generating the Pronunciation

To generate the pronunciation, we consider all possible segmentations $s = (s_1, \ldots, s_n)$ of the input word $x$ simultaneously. Each segmentation is considered as a potential *model* for having generated the true pronunciation of the word. Hence, assuming a prior probability distribution $p(s)$ over all segmentations, we can average our pronunciation distribution over all these models:

$$p(y \mid x) = \sum_s p(s)p(y \mid x, s). \tag{3}$$

Assuming that one of the segmentations is the correct generating model, the most likely true pronunciation for the written word $x$ can then be output as

$$\hat{y} := \arg\max_y p(y \mid x), \tag{4}$$

where $\arg\max_y$ denotes the value of $y$ that maximizes the expression.

---

[8]For simplicity, we do not include the novel pronunciation case of (1) in the examples.

|              | Text-to-Speech |            | Speech-to-Text |             |
| ------------ | :------------: | :--------: | :------------: | :---------: |
|              | Words (%)      | Phones (%) | Words (%)      | Letters (%) |
| Lower bound  | 18.34          | (73.79)    | 60.53          | (89.18)     |
| Upper bound  | 91.18          | (98.32)    | 94.12          | (99.02)     |

Table 3: Lower and upper bounds for accuracy as in Table 1 given the primary heuristic of shortest paths for segmentations *without overlap*.

## 2.4 The Choice of the Prior Distribution of Segmentations

A good choice for $p(s)$ is a distribution that assigns equal probabilities to all feasible segmentations with the minimum number of segments. Here a segmentation of the written word is considered feasible if all of its segments can be found in the training data (i.e., it corresponds to a path through the pronunciation lattice). As we use a uniform prior distribution, the model average (3) is equivalent to simply summing (collating) the estimated probabilities of all candidates corresponding to the same pronunciation.

The restriction to the minimum number of segments yields good performance in practice since it yields maximally concentrated probability distributions. The same restriction is also applied in the best-performing previously published PbA algorithms as described in Sec. 1.3. Another choice we have considered is a uniform distribution over so called "minimal segmentations", that is, segmentations such that you cannot get a feasible segmentation by removing a segment boundary. However, in our intitial experiments, this did not improve performance.

Table 3 shows the upper and lower bounds for accuracy for the NETtalk corpus given the restriction to minimum number of segments for the prior distribution. As expected, the bounds are somewhat wider than for the more restrictive overlap of one shown in Table 1.

## 2.5 Results and Discussion for Non-Overlapping Segments

Our probabilistically justified evaluation method corresponds to a *collated* product of estimated probabilities, where "estimated probabilities" are defined by Eq. (1) and "collation" means that the values of all candidates with the same pronunciation are summed together. Roughly speaking, this method is a combination of certain ideas already suggested in the literature:

1. The product of arc frequencies has been used as the evaluation function by Sullivan and Damper (1993) and Damper and Eastmond (1997) with the justification of it being closer to a probabilistic definition than the sum as used by Dedina and Nusbaum (1991) in their original program.

2. Collation has been used in the total product methods of Damper and Eastmond (1997) and Damper and Marchand (1998) and it is also present in Strategy 3 of Marchand and Damper (2000), the best of the five component methods used therein, and in Strategy 11 of Polyákova and Bonafonte (2008), the best of the 11 methods therein (see Sec. 1.5 for details).

3. Our estimated probabilities are similar to normalized frequencies as used by Sullivan and Damper (1993) in one of their definitions of "preference values", but we add one to the denominator in order to better estimate the true probabilities.

Our PROB method combines all three of the above and Table 4 shows that this probabilistically coherent method performs better than any other combination of normalization type and collation or no collation. It also performs better than any of the five component strategies of Marchand and Damper (2000) and better than any of the six new strategies

Text-to-Speech

| Evaluation method | No collation | | With collation | |
|---|---|---|---|---|
| | Words (%) | Phones (%) | Words (%) | Phones (%) |
| Product of | | | | |
|    absolute frequencies | 47.37 | (85.72) | 53.03 | (87.26) |
| normalized frequencies | 54.33 | (87.63) | 62.85 | (90.12) |
| estimated probabilities | 56.05 | (88.21) | **63.80** | **(90.51)** |

Speech-to-Text

| Evaluation method | No collation | | With collation | |
|---|---|---|---|---|
| | Words (%) | Letters (%) | Words (%) | Letters (%) |
| Product of | | | | |
|    absolute frequencies | 69.83 | (93.92) | 71.96 | (94.36) |
| normalized frequencies | 71.67 | (94.35) | 75.70 | (95.14) |
| estimated probabilities | 73.17 | (94.58) | **76.23** | **(95.22)** |

Table 4: Results for non-overlapping segmentations with minimal number of segments. The candidate pronunciations are evaluated by the product of absolute frequencies, normalized frequencies, or estimated probabilities as defined in Eq. (1). Without collation, the candidate with the highest value is chosen; with collation, the values of all candidates with the same pronunciation are summed together first and the the pronunciation with the highest sum is chosen. The method of collated estimated probabilities (results shown in bold) corresponds to our probabilistically justified method PROB and yields the unambiguously best results of those shown in the table.

of Polyákova and Bonafonte (2008) except for Strategy 11 (we will get to this strategy and why it works so well later in Sec. 3.6) and the generalization of our method for overlapping segments will outperform even Strategy 11. Thus, our theoretically justified approach appears to work well also in practice.

Indeed, normalization very clearly improves performance. Why it does so can be illustrated by pronouncing the word *recent* using all other words of the NETtalk corpus as the training set. The algorithm gives five different segmentations of two segments each and a total of 60 pronunciation candidates for the word. When absolute frequencies are considered, the best candidate is the pronunciation rEkxnt (/rɛkənt/), built from the segments *#rec* and *ent#*. The absolute frequencies for these segments are 29 and 218, respectively, and their product is 6322. However, when considering the normalized frequencies of the segment pronunciations, the values are 0.392 and 0.729, showing that neither segment has a particularly stable or unambiguous pronunciation.

When normalized frequencies are used, the best candidate pronunciation is ris-Nt (/ris-ņt/), formed from the segments *#r* and *ecent#*, which have absolute frequencies of 903 and 2 but relative frequencies of 1.0 and 1.0, respectively. In other words, the lexicon contains no words in which either segment is pronounced in any other way, so the segment pronunciations are very stable.

Collation also categorically improves performance in our experiments and is justified by model averaging as discussed in Sec. 2.3.

# 3 A Probabilistic Approach for Overlapping Segments

In many previous PbA algorithms, the pronunciations are generated using overlapping segments. This avoids some of the potential problems at the segment boundaries, in particular, it avoids certain impossible pronunciations at the segment boundary such as two consecutive stop consonants. It also limits the number of possible segmentations and therefore yields faster computation. We can also generalize our algorithm to use overlapping segments.

## 3.1 Product Rule (PROD)

Given a segmentation $s = ((a_1 : b_1), \ldots, (a_n : b_n))$, where the segments $(a_k : b_k) := \{a_k, \ldots, b_k\}$ cover all indices of the word $x$, but can overlap each other arbitrarily, we can define

$$p(y \mid x, s) \propto \prod_{k=1}^{n} \hat{p}(y_{a_k:b_k} \mid x_{a_k:b_k}) \tag{5}$$

and use this definition with the same decision function (4) as with non-overlapping segments. However, the expression on the right is no longer a properly normalized distribution as we require that the pronunciations of the overlapping parts of the segments match each other.[9] Thus, the actual distribution resulting from this definition can be considered as being generated by a procedure that samples the pronunciation for each segment from their respective distributions and then outputs the pronunciation for the full word only if the overlapping parts happen to match each other. This will bias the distribution of the overlapping part towards the most common pronunciations.

## 3.2 Conditional Probability Rules (CONDR, CONDL, CONDRL, CONDALL)

A better justified definition that avoids the bias of the PROD rule mentioned above is given by

$$p(y \mid x, s) = \prod_{k=1}^{n} \hat{p}(y_{a_k:b_k} \mid x_{a_k:b_k}, y_{a_1:b_1}, \ldots, y_{a_{k-1}:b_{k-1}}), \tag{6}$$

where for any sequence $\mathbf{y} = y_{a_1:b_1}, \ldots, y_{a_n:b_n}$, we define the estimated conditional probability expression used above as

$$\hat{p}(y_{a:b} \mid x_{a:b}, \mathbf{y}) = 0 \tag{7}$$

if $\mathbf{y}$ does not agree with $y_{a:b}$ at some index, as

$$\hat{p}(y_{a:b} \mid x_{a:b}, \mathbf{y}) = 1 \tag{8}$$

if $\mathbf{y}$ fully covers and agrees with $y_{a:b}$, and as

$$\hat{p}(y_{a:b} \mid x_{a:b}, \mathbf{y}) =$$
$$\frac{(\text{number of times } x_{a:b} \text{ is pronounced as } y_{a:b} \text{ in the training data})}{(\text{number of times } x_{a:b} \text{ appears with its pronunciation agreeing with } \mathbf{y}) + 1} \tag{9}$$

otherwise.

Thus, we sample the pronunciation for each segment from the conditional distribution given the pronunciations of any overlapping parts of the previous segments. This is of course

---

[9]The symbol $\propto$ means that the right side should be divided by a normalization constant to yield a properly normalized distribution

not symmetric, that is, the resulting distribution depends on the order of the segments. However, if one considers all orders (CONDALL) or just the left-to-right and right-to-left orders (CONDRL) of the segments as separate models, then symmetry is restored. In addition to the symmetric CONDALL and CONLR rules, we have also implemented and tested the left-to-right order (CONDR) and right-to-left order (CONDL) rules.

## 3.3 CONDF Rule

A third way to handle the overlaps is as follows. Let use denote by $y_{\text{overlap}}$ the pronunciations of the overlapping letters of a candidate pronunciation $y$. Then, the CONDF rule is given by

$$p(y \mid x, s) = \prod_{k=1}^{n} \hat{p}(y_{a_k:b_k} \mid x_{a_k:b_k}, y_{\text{overlap}}). \tag{10}$$

That is, we consider each segmentation and each configuration of the overlap letters as a separate model and normalize the probabilities given the pronunciations of the overlap letters.

## 3.4 Example Calculations of Conditional Probabilities

We shall use the word *longevity* to demonstrate the different normalization schemes used. Using a fixed one-letter overlap, the word can be segmented into four different segmentations with a minimal number of segments and these segmentations involve altogether ten different segments:

$$\#longe + evi + ity\#$$
$$\#longe + ev + vity\#$$
$$\#long + ge + evity\#$$
$$\#lon + nge + evity\#$$

After removing the word itself from the corpus, the ten segments have the following pronunciations and frequencies:

| | |
|---|---|
| #longe | lcGg-: 1 |
| #long | lcG-: 4, lanJ: 2, lcGg: 1 |
| #lon | lcG: 5, lan: 2, lon: 1 |
| nge | nJ-: 54, nJx: 18, Gg--: 12, nJE: 9, nJi: 9, G--: 6, NJ-: 3, Ggx: 1, n-i: 1 |
| ge | J-: 284, Jx: 105, JE: 80, Ji: 40, Z-: 26, g-: 19, --: 16, gE: 11, -x: 8, |
| | gx: 6, JI: 4, gA: 3, gi: 3, Ze: 2, -i: 1, Ja: 1, Za: 1, gI: 1, gY: 1, ge: 1 |
| ev | Ev: 83, Iv: 50, iv: 36, -v: 24, xv: 15, Ef: 1 |
| evi | ivi: 10, Evx: 8, IvA: 7, Ev-: 3, IvI: 3, ivA: 3, xvI: 3, -vI: 2, iv-: 2, |
| | ivI: 2, -v-: 1, Evi: 1, IvY: 1, ivY: 1, ivy: 1, xvA: 1, xvY: 1 |
| evity# | Evxti: 2 |
| vity# | vxti: 22 |
| ity# | xti: 421, Iti: 2 |

Despite the numerous different pronunciations that the middle segments have, only few of them can occur in our pronunciation candidates since the overlap between segments forces the pronunciations to overlap as well. For the same reason the two-segment segmentation #longe + evity# is not valid; the only pronunciations for the two segments do not have a common phoneme overlapping on the letter *e*.

The conditionality introduced by the overlap must also be taken into consideration when normalizing the frequencies. There are up to four different ways of normalizing each segment's frequency, since both the left and right side of the segment might or might not

have a pronunciation fixed for the overlapping part. This in turn is dependent on the order of assigning pronunciations for the segments in a segmentation; in a segmentation with $n$ segments there are $n!$ different orderings, which in our example case will mean six different orderings for the word *longevity*.

To demonstrate the different evaluation methods using conditionally normalized frequencies, let us consider the segmentation #lon + nge + evity# and its pronunciation lanJEvxti. The first segment #lon has the beginning marker # and will never overlap from its left end. Similarly, the last segment evity# with the ending marker will never overlap from its right end.

The absolute frequencies of the segment pronunciations (lan, nJE and Evxti) are 2, 9 and 2 respectively. The unconditionally normalized frequency for a segment pronunciation is the absolute frequency divided by the sum of frequencies of all candidate pronunciations plus one as shown in Eq. (1).
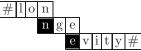
Conditional probabilities for the segments' pronunciations are derived by dividing the pronunciation's absolute frequency by the sum (plus one) of frequencies of all candidate pronunciations for the segment that conform to the overlaps already fixed by the neighboring segments as shown in Eq. (9).

For our example segmentation and pronunciation, the normalized frequencies and conditional probabilities are:

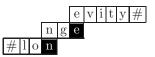|                            | lan | nJE   | Evxti |
|----------------------------|-----|-------|-------|
| unconditional normalization| 2/9 | 9/114 | 2/3   |
| left overlap fixed         |     | 9/92  | 2/3   |
| right overlap fixed        | 2/4 | 9/10  |       |
| both overlaps fixed        |     | 9/10  |       |

In our conditional probability rules we define five evaluation methods that utilize different orderings of the segments:

- CONDR: The pronunciation is constructed in a left-to-right order beginning with the leftmost segment. The first segment is normalized unconditionally and the construction continues to the right so that each consecutive segment's left overlap is fixed by the previous segment's pronunciation (black indicates a letter with a fixed pronunciation):



  Thus, the resulting estimated probability of our example pronunciation is $(2/9)(9/92)(2/3) \approx 0.0145$.

- CONDL: The pronunciation is constructed in a right-to-left order beginning with the rightmost segment. The first segment is normalized unconditionally and the construction continues to the left so that each consecutive segment's right overlap is fixed by the previous segment's pronunciation:



  The resulting estimated probability of our example pronunciation is $(2/3)(9/10)(2/4) = 0.3$.

- CONDLR: This method takes the average of the results of the CONDR and CONDL methods to make the evaluation symmetrical, yielding an estimated probability of $\approx (0.0145 + 0.3)/2 \approx 0.157$

| Evaluation method | Text-to-Speech | | Speech-to-Text | |
|---|---|---|---|---|
| | Words (%) | Phones (%) | Words (%) | Letters (%) |
| 10000 | 59.31 | (89.31) | **73.91** | **(94.78)** |
| 00100 | **62.46** | (89.90) | 72.41 | (94.42) |
| 00010 | 60.81 | **(89.96)** | 71.83 | (94.41) |
| 10101 | 65.46 | (91.02) | **75.90** | (95.18) |
| 11111 | **65.62** | **(91.14)** | 75.52 | (95.10) |
| 00000000001 | **64.00** | **(90.57)** | 75.77 | **(95.21)** |
| 00101100001 | 65.35 | (90.99) | **76.19** | **(95.26)** |
| 00111010011 | 66.05 | **(91.25)** | 75.71 | (95.15) |
| 00101000001 | **66.14** | (91.20) | 76.00 | (95.22) |
| PROB | 63.80 | (90.51) | 76.23 | (95.22) |
| PROD | 63.88 | (90.58) | 75.75 | (95.14) |
| CONDF | **66.21** | **(91.13)** | 76.26 | (95.25) |
| CONDR | 63.26 | (90.14) | 75.80 | (95.11) |
| CONDL | 65.98 | (91.12) | 76.05 | (95.19) |
| CONDRL | 65.50 | (90.98) | 76.30 | (95.27) |
| CONDALL | 65.66 | (91.02) | **76.39** | **(95.29)** |

Table 5: The results of our probabilistic algorithms for overlapping segments compared to the best results of previously published algorithms from Table 2. The probabilistically questionable PROD method is also shown for comparison.

- CONDALL: The pronunciation is assembled in all $n!$ orders and the average of their estimated probabilities is taken, yielding $\approx 0.022$. (Here the averaging corresponds to assuming a uniform prior distribution over the $n! = 6$ possible models).

- CONDF: This method first fixes the overlap letters and then normalizes each segment with both ends fixed (except at the beginning and end of the word):



In this case, the estimated probability of our example pronunciation is $(2/4)(9/10)(2/3) = 0.3$.

Summed over all segmentations, the CONDL and CONDF methods in fact yield the best score (probability) for the correct pronunciation while the PROD method, the product of unconditional probabilities, does not.

## 3.5 Results and Discussion for Overlapping Segments

For overlapping segments, a strictly probabilistic evaluation is less straightforward than for non-overlapping segments, because there is no obvious way of defining the probability distribution of the overlapping segments. This is probably the reason why the approach of Sullivan and Damper (1993), the one coming closest to a probabilistic formulation (and using non-overlapping segments), has been given up in later work, as overlapping segments are used exclusively in the most recent work. However, as our results will show, our probabilistically justified evaluation rules are in fact better than taking the direct product of segment probabilities. Thus, it appears that also for overlapping segments, our probabilistically coherent theory works in practice.

Table 5 shows the results of our probabilistic evaluation methods compared to the best previously published algorithms. The CONDF method performs better than any

previously published algorithm including all combinations of the 11 component strategies of Marchand and Damper (2000) and Polyákova and Bonafonte (2008). In fact, all of our probabilistically justified left-right symmetric (i.e., excluding the asymmetric CONDR and CONDL methods) methods perform better than any previously published single strategy.

For the text-to-speech direction, the asymmetric CONDR and CONDL strategies yield very different results, and interestingly, it is the reverse right-to-left order of the segments that results in higher accuracy. For the speech-to-text direction, the difference between CONDR and CONDL is much smaller, but still in the same direction. The symmetric CONDRL strategy, which includes both orders of the segments, yields more consistent performance and is likely to maintain good performance across different lexicons better. The CONDALL strategy, which includes all $n!$ orders of the segments offers a slight improvement over CONDRL, and may be the best candidate for a general method, but has a somewhat higher computational cost for long words.

## 3.6 Strategy 11 of Polyákova and Bonafonte (2008) and the Magic Root

The only previously published single scoring strategy that works better than PROB (our simplest probabilistic approach using non-overlapping segments) is Strategy 11 of Polyákova and Bonafonte (2008). This is the collated geometric mean of the arc frequencies. As we have shown, collation categorically increases performance so it is not surprising that this method works better than just the product of arc frequencies. However, what is surprising is that the geometric mean (i.e., the $n$th root of the product of the $n$ arc frequencies) actually works better than the product, contrary to any probabilistic justification.

What strikes as odd in using the geometric mean is that the degree of the root involved depends on the number of segments in the candidate pronunciation. This would perhaps be justifiable if the different candidate pronunciations had different numbers of segments (so as to bring the products on the same scale), but as the primary strategy already limits the candidates to those with minimum number of segments, it is difficult to justify why an effectively different rule should be used depending on the number of segments. In our leave-one-out evaluation with the NETtalk corpus, the pronunciations for most words can be constructed from two or three segments; thus, the question arises if it would be better to instead use two or three or some other constant as the degree of the root in the method.

Table 6 shows the results of the product rule with different constant roots and using the varying $kt$th root of Polyákova and Bonafonte (2008). It turns out that a constant root is indeed better and that the accuracy increases further by replacing the absolute frequencies by the estimated probabilities as defined by Eq. (1). Still, the highest accuracy obtained is less than that of CONDF, the best of our probabilistically justified methods shown in Table 5.

These results lead one to wonder if this "magic root" could be used to improve the performance also in our probabilistically justified methods. This is indeed the case as shown by the results of Table 7. Taking a root of the estimated probabilities (before any summing) appears to improve the accuracy of all of our probabilistic methods, except for the speech-to-text direction of the CONDF method. The optimum root appears to vary a bit depending on the method and the lexicon, but based on these experiments, the cubic root seems like a good compromise that works well for all methods. The CONDF method, which performed best, benefits the least from the root. The asymmetric CONDL method obtains the best word accuracy (66.61%) of all methods when combined with the cubic root.

Text-to-Speech

| Root | Product of absolute frequencies | | Product of normalized frequencies | | Product of estimated probabilities | |
|---|---|---|---|---|---|---|
| | Words (%) | Phones (%) | Words (%) | Phones (%) | Words (%) | Phones (%) |
| NC | 59.31 | (89.31) | 60.81 | (89.81) | 60.15 | (89.63) |
| $n$ | 64.00 | (90.57) | 65.45 | (91.00) | 65.56 | (91.04) |
| 1 | 61.76 | (89.96) | 63.66 | (90.54) | 63.88 | (90.58) |
| 2 | 63.81 | (90.51) | 65.22 | (90.93) | 65.41 | (90.99) |
| 3 | 64.98 | (90.85) | 65.84 | (91.12) | 65.99 | (91.17) |
| 4 | 65.40 | (90.96) | 66.01 | (91.16) | 66.11 | (91.20) |
| 5 | 65.58 | (90.99) | **66.06** | **(91.18)** | **66.15** | **(91.21)** |
| 6 | 65.69 | (91.01) | 66.05 | (91.17) | 66.13 | (91.19) |
| 7 | 65.75 | (91.02) | 66.02 | (91.15) | 66.10 | (91.17) |
| 8 | **65.83** | **(91.04)** | 66.01 | (91.13) | 66.09 | (91.15) |
| 9 | 65.82 | (91.04) | 66.00 | (91.13) | 66.11 | (91.15) |
| 10 | 65.83 | (91.04) | 65.99 | (91.12) | 66.11 | (91.15) |

Speech-to-Text

| Root | Product of absolute frequencies | | Product of normalized frequencies | | Product of estimated probabilities | |
|---|---|---|---|---|---|---|
| | Words (%) | Letters (%) | Words (%) | Letters (%) | Words (%) | Letters (%) |
| NC | 73.91 | (94.78) | 74.37 | (94.92) | 74.51 | (94.92) |
| $n$ | 75.77 | (95.17) | 76.09 | (95.23) | 76.25 | (95.25) |
| 1 | 74.91 | (94.99) | 75.50 | (95.12) | 75.76 | (95.14) |
| 2 | 75.68 | (95.17) | 75.98 | (95.21) | 76.12 | (95.23) |
| 3 | **75.99** | **(95.20)** | **76.14** | **(95.25)** | 76.26 | **(95.26)** |
| 4 | 75.97 | (95.19) | 76.12 | (95.24) | **76.27** | (95.26) |
| 5 | 75.98 | (95.19) | 76.04 | (95.22) | 76.18 | (95.24) |
| 6 | 75.94 | (95.18) | 75.96 | (95.20) | 76.11 | (95.22) |
| 7 | 75.91 | (95.17) | 75.96 | (95.19) | 76.09 | (95.21) |
| 8 | 75.91 | (95.16) | 75.91 | (95.17) | 76.05 | (95.19) |
| 9 | 75.88 | (95.15) | 75.88 | (95.16) | 76.03 | (95.18) |
| 10 | 75.87 | (95.14) | 75.86 | (95.16) | 75.98 | (95.17) |

Table 6: Accuracy obtained for the product of absolute frequencies, normalized frequencies, or estimated probabilities as the secondary evaluation function (as in Table 4) for shortest paths with overlap of one. The first row (NC) shows the results without collation. The remaining rows show the results with collation and with the indicated root taken of the product before collation (without collation the root does not have any effect). Here $n$ is the number of segments in the candidate pronunciation and so the degree $n$ root of the product of absolute frequencies corresponds to the PFSP rule of Polyákova and Bonafonte (2008). However, these results indicate that a constant root is in fact better than the varying root. Also, consistent with our other results, normalization increases accuracy.

| Evaluation method | Text-to-Speech | | Speech-to-Text | |
| --- | --- | --- | --- | --- |
| | Words (%) | Phones (%) | Words (%) | Letters (%) |
| PROB | 63.80 | (90.51) | 76.23 | (95.22) |
| PROB$^{1/2}$ | 65.25 | (90.98) | 76.66 | (95.36) |
| PROB$^{1/3}$ | 65.60 | (91.14) | **76.83** | **(95.41)** |
| PROB$^{1/4}$ | 65.81 | **(91.20)** | 76.74 | (95.38) |
| PROB$^{1/5}$ | 65.87 | (91.19) | 76.69 | (95.37) |
| PROB$^{1/6}$ | **65.92** | (91.19) | 76.65 | (95.37) |
| PROB$^{1/7}$ | 65.89 | (91.19) | 76.58 | (95.35) |
| PROB$^{1/8}$ | 65.88 | (91.17) | 76.51 | (95.33) |
| PROB$^{1/9}$ | 65.88 | (91.16) | 76.51 | (95.33) |
| PROB$^{1/10}$ | 65.88 | (91.15) | 76.52 | (95.33) |
| CONDF | 66.21 | (91.13) | **76.26** | **(95.25)** |
| CONDF$^{1/2}$ | 66.29 | (91.15) | 76.17 | (95.21) |
| CONDF$^{1/3}$ | 66.28 | (91.14) | 76.07 | (95.19) |
| CONDF$^{1/4}$ | **66.31** | **(91.14)** | 76.02 | (95.17) |
| CONDF$^{1/5}$ | 66.27 | (91.11) | 76.00 | (95.16) |
| CONDR | 63.26 | (90.14) | 75.80 | (95.11) |
| CONDR$^{1/2}$ | 64.76 | (90.61) | 75.94 | (95.17) |
| CONDR$^{1/3}$ | 65.20 | (90.77) | **76.00** | **(95.17)** |
| CONDR$^{1/4}$ | **65.41** | **(90.84)** | 75.95 | (95.15) |
| CONDR$^{1/5}$ | 65.46 | (90.85) | 75.93 | (95.14) |
| CONDL | 65.98 | (91.12) | 76.05 | (95.19) |
| CONDL$^{1/2}$ | 66.52 | (91.31) | 76.20 | (95.23) |
| CONDL$^{1/3}$ | **66.61** | **(91.33)** | 76.21 | (95.24) |
| CONDL$^{1/4}$ | 66.54 | (91.31) | 76.17 | (95.23) |
| CONDL$^{1/5}$ | 66.51 | (91.29) | 76.12 | (95.22) |
| CONDRL | 65.50 | (90.98) | 76.30 | (95.27) |
| CONDRL$^{1/2}$ | 66.24 | (91.20) | **76.36** | **(95.27)** |
| CONDRL$^{1/3}$ | **66.33** | **(91.23)** | 76.32 | (95.26) |
| CONDRL$^{1/4}$ | 66.28 | (91.19) | 76.23 | (95.24) |
| CONDRL$^{1/5}$ | 66.26 | (91.16) | 76.16 | (95.22) |
| CONDALL | 65.66 | (91.02) | 76.38 | (95.29) |
| CONDALL$^{1/2}$ | 66.30 | (91.23) | **76.39** | **(95.28)** |
| CONDALL$^{1/3}$ | **66.39** | **(91.25)** | 76.36 | (95.27) |
| CONDALL$^{1/4}$ | 66.35 | (91.22) | 76.26 | (95.24) |
| CONDALL$^{1/5}$ | 66.31 | (91.18) | 76.17 | (95.21) |

Table 7: Experiments applying the "magic root" suggested by the good performance of the PFSP method of Polyákova and Bonafonte (2008) to our probabilistically justified methods. The degree of the applied root is shown as the inverse exponent of the method name.

# 4   Conclusion

The current state of the art in pronunciation by analogy is well summarized by the following passage from Damper and Marchand (2006):

> Automatic pronunciation of words from their spelling alone is a hard computational problem, especially for languages like English and French where there is only a partially consistent mapping from letters to sound. Currently, the best known approach uses an inferential process of analogy with other words listed in a dictionary of spellings and corresponding pronunciations. However, the process produces multiple candidate pronunciations and little or no theory exists to guide the choice among them.

In this article, we have have presented a simple probabilistically justified approach to choosing the best of several candidate pronunciations. Our principled approach outperforms all previously published PbA algorithms (the best one by a small margin). Thus, we can obtain or exceed state of the art performance with a strictly theoretically justified evaluation method without any ad hoc modifications.

Our probabilistic approach differs from the more conventional methods by a series of relatively small changes, each of which makes the method closer to a strictly probabilistic formulation. We have shown that each of these changes improves performance. Thus, even if one prefers not to adopt our strictly probabilistic approach, one can still use the following results yielded by our work in any PbA algorithms:

1. Accuracy of a PbA algorithm can usually be improved by collation, that is, by summing together the values of all candidates with the same pronunciation.

2. Instead of using absolute frequencies of the pronunciations, performance will typically increase by using normalized frequencies, or even better, estimated probabilities as defined by Eqs. (1) and (9), instead.

3. Some improvements of accuracy can be expected by replacing the frequencies or estimated probabilities of the segment pronunciations by their cubic root or some other root with degree $> 1$. (This works in practice although we do not have a theoretical justification for it).

Outside PbA literature, there are also other, statistically justified approaches for text-to-speech conversion, see for example Chen (2003), Bisani and Ney (2008), and Jiampojamarn and Kondrak (2009). Based on reported experiments on the NETtalk corpus, the state of the art accuracy in these methods is actually better than in PbA; indeed, the best quoted word accuracy for the NETtalk corpus is 69% in Bisani and Ney (2008). However, these methods are far more complicated, both conceptually and computationally than PbA algorithms; the programs run for several hours and typically require setting aside a subset of words for tuning certain global parameters of the algorithm. The PbA algorithms remain attractive due to their simplicity and comparatively good accuracy.

# Acknowledgement

# References

Maximilian Bisani and Hermann Ney. Joint-sequence models for grapheme-to-phoneme conversion. *Speech Communication*, 50:434–451, 2008.

Stanley F. Chen. Conditional and joint models for grapheme-to-phoneme conversion. In *Proceedings of European Conference on Speech Communication and Technology*, pages 2033–2036, Brighton, UK, 2003.

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

R. I. Damper and Y. Marchand. Information fusion approaches to the automatic pronunciation of print by analogy. *Information Fusion*, 7:207–230, 2006.

R.I. Damper and Y. Marchand. Improving pronunciation by analogy for text-to-speech applications. In *Proceedings of 3rd European Speech Communication Association (ESCA)/COCOSDA International Workshop on Speech Synthesis*, 1998.

Robert I. Damper and John F. G. Eastmond. Pronouncing text by analogy. In *Proceedings of the 16th conference on Computational linguistics*, pages 268–273, 1996.

Michael J. Dedina and Howard C. Nusbaum. PRONOUNCE: a program for pronunciation by analogy. *Computer Speech and Language*, 5:55–64, 1991.

Robert J. Glushko. The organization and activation of orthographic knowledge in reading aloud. *Journal of Experimental Psychology: Human Perception and Performance*, 5: 674–691, 1979.

Sittichai Jiampojamarn and Grzegorz Kondrak. Online discriminative training for grapheme-to-phoneme conversion. In *Proceedings of INTERSPEECH 2009*, pages 1303–1306, 2009.

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

Tatyana Polyákova and Antonio Bonafonte. Further improvements to pronunciation by analogy. In *Actas de las V Jornadas en Tecnología del Habla*, pages 149–152, Bilbao, Spain, 2008.

Tatyana Polyákova and Antonio Bonafonte. New strategies for pronunciation by analogy. In *ICASSP '09*, 2009.

Terrence Sejnowski and Charles Rosenberg. Parallel networks that learn to pronounce english text. *Complex Systems*, 1:145–168, 1987.

Tasanawan Soonklang, Robert I. Damper, and Yannick Marchand. Multilingual pronunciation by analogy. *Natural Language Engineering*, 1:1–21, 2008.

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

François Yvon. Grapheme-to-phoneme conversion using multiple unbounded overlapping chunks. In *Proceedings of the Conference on New Methods in Natural Language Processing, NeMLaP '96*, pages 218–228, Ankara, Turkey, 1996.