

# TOWARDS A MATHEMATICAL PROOF OF THE SPEECH ALIGNMENT PARADOX

David Sündermann<sup>1,2</sup>, Jaka Smrekar<sup>3</sup>, Harald Höge<sup>1</sup>

<sup>1</sup>Siemens Corporate Technology, Munich, Germany

<sup>2</sup>Technical University of Catalonia, Barcelona, Spain

<sup>3</sup>University of Ljubljana, Ljubljana, Slovenia

david@suendermann.com jaka.smrekar@fmf.uni-lj.si harald.hoege@siemens.com

## ABSTRACT

Applying a recently presented text-independent speech alignment technique based on unit selection to the training of a voice conversion system suggested that the more training data was available, the less speaker-specific information was learned. This paradoxical effect contradicts the experience we have from other corpus-based applications as speech recognition, synthesis or translation. There, the performance usually gains with increasing amount of data. In this paper, we investigate this paradox by means of objective tests and derive a mathematical model of the underlying stochastic process.

## 1. INTRODUCTION

In several speech processing applications (e.g. in speech recognition [1], speaker identification [2], or speech data mining [3, 4]), we have to find a time alignment between speech samples, usually generated by different speakers. Mainly, the texts underlying the compared speech samples is identical, which allows for applying dynamic time warping [1] to the problem. If the underlying text is known, forced alignment [5] can be performed, which may lead to more accurate results.

However, certain applications require the alignment of utterances, which are not parallel. Here, we face the text-independent alignment task. Recently, we presented a technique based on unit selection, which was used for text-independent voice conversion training [6] and later extended to cross-language voice conversion [7].

When compared to text-dependent alignment (dynamic time warping), the achieved speech quality of the voice-converted speech was improved by means of the novel technique, whereas the similarity to the target speaker decreased. Table 1 shows the results of a subjective evaluation reported in [7]. As common metrics, for both overall speech quality

|                  | MOS <sub>Q</sub><br>(quality) | MOS <sub>S</sub><br>(similarity) |
|------------------|-------------------------------|----------------------------------|
| text-dependent   | 3.3                           | 2.4                              |
| text-independent | 3.5                           | 2.0                              |
| source voice     | 4.7                           | 1.6                              |

**Table 1.** Results of a subjective evaluation on the application of speech alignment to voice conversion: overall speech quality (MOS<sub>Q</sub>) and similarity to the target (MOS<sub>S</sub>)

and similarity to the target, a mean opinion score [8] on a five-point scale (1 for bad to 5 for excellent) was applied.

As informal listening tests suggested, both effects, the quality boost and the similarity score loss, increased with increasing amount of training data. This paper is to study this paradox focusing on the similarity effect, which can be described by objective criteria, rather than the speech quality, whose objective investigation is still a hard problem [9].

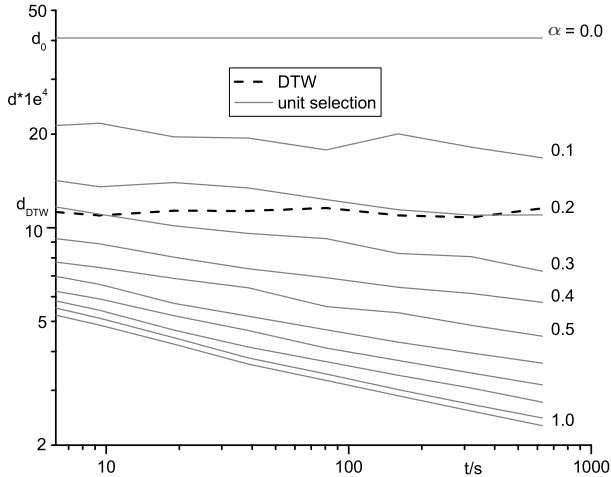
## 2. TEXT-INDEPENDENT SPEECH ALIGNMENT BASED ON UNIT SELECTION

We consider two arbitrary speech samples to be aligned. At first, they are broken down into frames<sup>1</sup>. Now, the frames are encoded leading to two sequences of feature vectors<sup>2</sup> representing source and target speech,  $x_1^M$  and  $y_1^N$ . To perform the alignment, from the latter, vectors are to be selected and joined to a sequence  $\tilde{y}_1^M$  that optimally corresponds to the source sequence. This is done by taking two criteria into account:

<sup>1</sup>In our study, we utilized pitch-synchronous time frames produced by the Praat tool [10], since this allows for using standard pitch modification techniques to change prosodical properties of speech in the framework of voice conversion. However, all the following considerations also apply to constant frame lengths as mostly used in speech recognition.

<sup>2</sup>Here, we use line spectral frequencies; in other applications, one would certainly prefer other types as mel frequency cepstral coefficients or linear predictive coefficients, cf. [11].

This work has been partially funded by the European Union under the integrated project TC-Star - Technology and Corpora for Speech to Speech Translation - <http://www.tc-star.org>.



**Fig. 1.** Text-independent speech alignment: average distance between corresponding source and target feature vectors  $d$  depending on the amount of data and the trade-off parameter  $\alpha$ .

- The distance between source and corresponding target features (*target cost*) is minimum (optimal correspondence).
- The distance to the neighbors of the corresponding target feature vector (*concatenation cost*) is minimum (optimal naturalness). This criterion is supposed to select naturally smooth segments<sup>3</sup> from the target feature vector sequence  $y_1^M$ .

Mostly, these optima do not coincide, and we must get by with a compromise between both: We search for the minimum of the weighted sum of target and concatenation cost for each source feature vector:

$$\tilde{y}_1^M = \arg \min_{y_1^M} \sum_{m=1}^M \left\{ \alpha S(y_m - x_m) + (1 - \alpha) S(y_{m-1} - y_m) \right\}. \quad (1)$$

Here,  $S(w)$  is the Euclidean distance

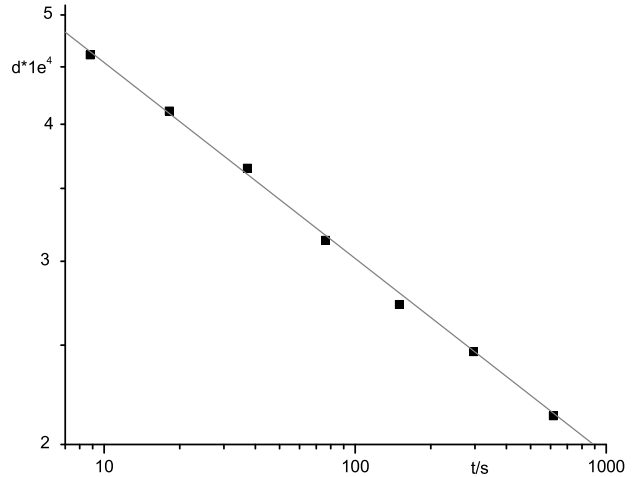
$$S(w) = \sqrt{w'w} \quad (2)$$

and  $0 \leq \alpha \leq 1$  is a weight influencing the trade-off between target and concatenation cost.

### 3. EXPERIMENTAL EVIDENCE OF THE SPEECH ALIGNMENT PARADOX

As already argued in Section 1, we want to limit the investigations on the speech alignment paradox to the similarity of

<sup>3</sup>or *units*; that is, where the term *unit selection* comes from. This paradigm is well-known from concatenative speech synthesis, where optimal speech units are selected and concatenated, cf. [12].



**Fig. 2.** Special case  $\alpha = 1$ .

the aligned speech samples. We claimed that the more data was available, the less speaker-specific information could be extracted for the application to voice conversion. An explanation of this effect is that the units, which are selected to minimize the Euclidean distance to the target become more and more similar to the latter, the more data is available to select from.

To investigate this effect, we want to use the mean Euclidean distance between the aligned feature vector sequences as an objective measure:

$$d = \frac{1}{M} \sum_{m=1}^M S(\tilde{y}_m - x_m).$$

Now, we want to look at the dependence of the increasing similarity, i.e. decreasing  $d$  value, on the amount of data available. In doing so, we also have to take the trade-off parameter  $\alpha$ , see Eq. 1 into account. We conducted experiments using the evaluation corpus of the project TC-Star [13], which consists of about 10 minutes of speech of two female and two male British English voices. Independent of the voice combinations to be aligned, we got very similar outcomes. As an example, we display the results of a female-male voice combination in Figure 1 in double logarithmic representation. We observe that independent of the trade-off parameter  $\alpha$ , the values of  $d$  almost constantly decrease<sup>4</sup>. To simplify matters, in the following, we look at the special case  $\alpha = 1$ , the respective diagram is shown in Figure 2.

For the considered amounts of data, our test samples are almost located on a straight line in double logarithmic representation. Consequently, the relation between  $d$  and  $t$  can

<sup>4</sup>except for  $\alpha = 0$ , which does not lead to a useful alignment, since no target costs are considered

be approximated by<sup>5</sup>:

$$\log d = c - b \log t \quad \text{with } b > 0 ;$$

exponentiation yields

$$d = e^{c-b \log t} = e^c e^{\log t^{-b}} = at^{-b} \quad \text{with } a, b > 0 . \quad (3)$$

If we assume the validity of Eq. 3 also for amounts of data beyond the experiment's scope, we get the limit

$$\lim_{t \rightarrow \infty} d = \lim_{t \rightarrow \infty} at^{-b} = 0 . \quad (4)$$

This means, for very large amounts of data, the aligned speech samples become very similar to each other (for the limit case even identical), which provides evidence for the speech alignment paradox. Unfortunately, the speech alignment algorithm based on unit selection is very computationally expensive (cf. [7]); to process 400 seconds of speech, the computation took more than 80 hours on a 3GHz Intel Xeon machine. Thus, currently, we are not able to massively increase the amount of data. This is the main reason for describing the paradox by mathematical means as done in the next section.

#### 4. TOWARDS A MATHEMATICAL PROOF OF THE SPEECH ALIGNMENT PARADOX

Although the empirical investigations of Section 3 were confirmed by several experimental cycles, doubts arose on the validity of the limit value shown in Eq. 4, as it could be interpreted as follows:

**If there is enough speech data available, an arbitrary utterance of an arbitrary voice can be produced only by selecting and concatenating units from this data.**

However, the crucial point in the statement is the word *enough*. Applying the parameters  $a = 6.8$  and  $b = 0.18$  determined on the data of Figure 2 to Eq. 3, we estimated the required amount of data for several degrees of similarity, cf. Table 2. We see that the amount of necessary data extremely grows when the mean distance between source and aligned target feature vectors becomes smaller and soon exceeds the limits of the technical possible.

Nonetheless, since the validity of the statement phrased above could be of high interest to the speech processing community, in the following, we will investigate the alignment technique's behavior for very large amounts of data by mathematical means.

##### 4.1. Speech as a Mixture of Gaussians

As introduced in Section 2, we describe the processed speech by sequences of feature vectors, whose statistical charac-

<sup>5</sup>in the following equations, we use the normalized time  $t := \frac{t}{s}$  to avoid confusion

| $d$ | $t$                         | disk space |
|-----|-----------------------------|------------|
| 5   | 5.6 s                       | 174 kB     |
| 2   | 900 s = 15 min              | 27 MB      |
| 1   | $4.2 \cdot 10^4$ s = 11.7 h | 1.3 GB     |
| 0.5 | $2.0 \cdot 10^6$ s = 22.8 d | 59 GB      |
| 0.2 | $3.2 \cdot 10^8$ s = 10.3 a | 9.2 TB     |

**Table 2.** Required amount of data ( $t$ ) for certain degrees of similarity ( $d$ ) and the corresponding hard disk space necessary for storing a 16kHz/16bit PCM version of the data

teristics are very often described by means of the Gaussian mixture model – in literature, we find applications of this model to speech recognition [14], language identification [15], voice conversion [16], speaker recognition [17], speaking rate estimation [18], and gender classification [19], and more.

The success of the Gaussian mixture model in these speech processing fields also suggests its application to the investigation of the speech alignment paradox.

In order to keep things manageable, we strongly reduce the degrees of freedom for our first investigations as follows:

- We set the number of Gaussian mixture densities to  $K = 1^6$ .
- We reduce the dimensionality of the feature vectors to  $D = 1$  (w.l.o.g.).
- We assume identical covariance matrices for the feature vector sequences to be aligned, i.e., for  $D = 1$ , we have the standard deviation  $\sigma$ .

##### 4.2. The A-Priori Alignment

When we have a look at the speech samples without performing any alignment, we can determine an a-priori value for the mean vector distance  $d$ : the expected value  $E_1(d)^7$ . The latter is the expected distance between the two normally distributed random vectors  $x$  and  $y$

$$E_1(d) = \int_{-\infty}^{\infty} E_1(d|x) \mathcal{N}(x|\mu_x, \sigma) dx , \quad (5)$$

where  $E_1(d|x)$  is the expected value of  $d$  if  $x$  is fixed and  $\mathcal{N}(x|\mu_x, \sigma)$  is the probability density function of a normal

<sup>6</sup>Hence, for these considerations, there is no need for using the term *mixture* when referring to the model. Interestingly, in particular applications to voice conversion, it turns out that the optimal choice for  $K$  is small anyway: [20] reported  $K = 6$ , in [21] we find  $K = 4$ , and in [22] we even go down to  $K = 1$  for particular cases of text-independent speech alignment.

<sup>7</sup>The subindex 1 is due to the fact that this expected value is a special case of that described in Section 4.3.

distribution. In the following, we use the *standard* normal distribution

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$

and modify Eq. 5 accordingly

$$E_1(d) = \frac{1}{\sigma} \int_{-\infty}^{\infty} E_1(d|x) f\left(\frac{x - \mu_x}{\sigma}\right) dx. \quad (6)$$

Now, we calculate the expected value of  $d$  if  $x$  is fixed

$$\begin{aligned} E_1(d|x) &= \frac{1}{\sigma} \int_{-\infty}^{\infty} |x - y| f\left(\frac{y - \mu_y}{\sigma}\right) dy \\ &= \frac{1}{\sigma} \int_{-\infty}^x (x - y) f\left(\frac{y - \mu_y}{\sigma}\right) dy \\ &\quad - \frac{1}{\sigma} \int_x^{\infty} (x - y) f\left(\frac{y - \mu_y}{\sigma}\right) dy \\ &= \left[ \sigma f\left(\frac{y - \mu_y}{\sigma}\right) + (x - \mu_y) \Phi\left(\frac{y - \mu_y}{\sigma}\right) \right]_{y=-\infty}^x \\ &\quad - \left[ \sigma f\left(\frac{y - \mu_y}{\sigma}\right) + (x - \mu_y) \Phi\left(\frac{y - \mu_y}{\sigma}\right) \right]_{y=x}^{\infty} \\ &= 2\sigma f\left(\frac{x - \mu_y}{\sigma}\right) + (x - \mu_y) \left[ 2\Phi\left(\frac{x - \mu_y}{\sigma}\right) - 1 \right], \end{aligned} \quad (7)$$

where  $\Phi(x)$  is the standard normal cumulative density function, thus we have  $\frac{d\Phi(x)}{dx} = f(x)$ . By inserting the result into Eq. 6, we get

$$\begin{aligned} E_1(d) &= 2 \int_{-\infty}^{\infty} f\left(\frac{x - \mu_y}{\sigma}\right) f\left(\frac{x - \mu_x}{\sigma}\right) dx \\ &\quad + 2 \int_{-\infty}^{\infty} \frac{x - \mu_y}{\sigma} \Phi\left(\frac{x - \mu_y}{\sigma}\right) f\left(\frac{x - \mu_x}{\sigma}\right) dx \\ &\quad - \int_{-\infty}^{\infty} \frac{x - \mu_y}{\sigma} f\left(\frac{x - \mu_x}{\sigma}\right) dx \\ &= T_1 + T_2 + T_3. \end{aligned}$$

For  $T_1$  and  $T_3$ , we have straightforward solutions

$$\begin{aligned} T_1 &= \frac{1}{\pi} \int_{-\infty}^{\infty} e^{-\frac{(x - \mu_y)^2 + (x - \mu_x)^2}{2\sigma^2}} dx \\ &= \frac{1}{\pi} \int_{-\infty}^{\infty} e^{-\frac{(2x - \mu_x - \mu_y)^2 + (\mu_y - \mu_x)^2}{4\sigma^2}} dx \\ &= 2f\left(\frac{\mu_y - \mu_x}{\sqrt{2}\sigma}\right) \int_{-\infty}^{\infty} f\left(\frac{2x - \mu_x - \mu_y}{\sqrt{2}\sigma}\right) dx \\ &= 2f\left(\frac{\mu_y - \mu_x}{\sqrt{2}\sigma}\right) \left[ \frac{\sigma}{\sqrt{2}} \Phi\left(\frac{2x - \mu_x - \mu_y}{\sqrt{2}\sigma}\right) \right]_{x=-\infty}^{\infty} \\ &= \sqrt{2}\sigma f\left(\frac{\mu_y - \mu_x}{\sqrt{2}\sigma}\right); \end{aligned}$$

$$\begin{aligned} T_3 &= - \left[ (\mu_x - \mu_y) \Phi\left(\frac{x - \mu_x}{\sigma}\right) - \sigma f\left(\frac{x - \mu_x}{\sigma}\right) \right]_{x=-\infty}^{\infty} \\ &= \mu_y - \mu_x, \end{aligned}$$

whereas  $T_2$  requires a more complex derivation, which we omit here only giving the final result<sup>8</sup>

$$T_2 = \sqrt{2}\sigma f\left(\frac{\mu_y - \mu_x}{\sqrt{2}\sigma}\right) + 2(\mu_y - \mu_x) \left[ \Phi\left(\frac{\mu_y - \mu_x}{\sqrt{2}\sigma}\right) - 1 \right]$$

yielding the searched expected value of  $d$  (in the following, we use  $\delta = \mu_y - \mu_x$ , the difference between the distribution means)

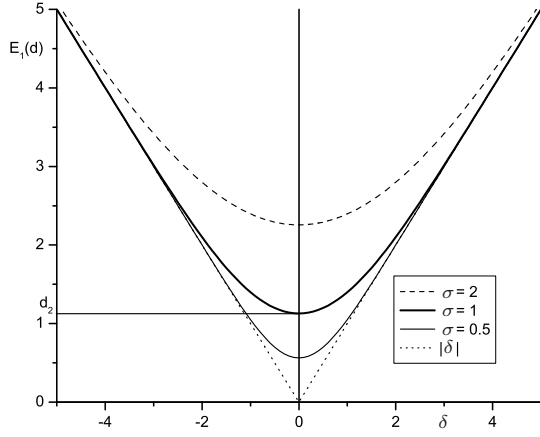
$$E_1(d) = 2\sqrt{2}\sigma f\left(\frac{\delta}{\sqrt{2}\sigma}\right) + 2\delta \Phi\left(\frac{\delta}{\sqrt{2}\sigma}\right) - \delta. \quad (8)$$

Figure 3 shows  $E_1(d)$  as a function of  $\delta$  for  $\sigma \in \{0.5, 1, 2\}$  and indicates the lower bound of  $E_1(d)$ , which is given by the limit

$$\lim_{\delta/\sigma \rightarrow \pm\infty} E_1(d) = |\delta|. \quad (9)$$

This limit can also be calculated using Eq. 6: When  $\frac{\delta}{\sigma}$  approaches infinity, the deviation of  $x$ 's distribution function becomes infinitely small as compared with its mean's distance to  $y$ 's mean. Consequently, the normal distribution

<sup>8</sup>The authors would be happy to provide the proof to everybody, who is interested.



**Fig. 3.** Expected value of the mean distance between two feature vectors  $E_1(d)$  as a function of the difference of the distribution means  $\delta$  and of the standard deviation  $\sigma$ .

can be replaced by the Dirac delta function  $\Delta$  yielding

$$\begin{aligned}
\lim_{\delta/\sigma \rightarrow \pm\infty} E_1(d) &= \lim_{\delta/\sigma \rightarrow \pm\infty} \frac{1}{\sigma} \int_{-\infty}^{\infty} E_1(d|x) f\left(\frac{x - \mu_x}{\sigma}\right) dx \\
&= \lim_{\delta/\sigma \rightarrow \pm\infty} \frac{1}{\sigma} \int_{-\infty}^{\infty} E_1(d|x) \Delta\left(\frac{x - \mu_x}{\sigma}\right) dx \\
&= \lim_{\delta/\sigma \rightarrow \pm\infty} \frac{1}{\sigma} \int_{-\infty}^{\infty} E_1(d|\sigma\xi + \mu_x) \Delta(\xi) \sigma d\xi \\
&= \lim_{\delta/\sigma \rightarrow \pm\infty} E_1(d|\mu_x) \quad (10)
\end{aligned}$$

Applying Eq. 7 yields

$$\lim_{\delta/\sigma \rightarrow \pm\infty} E_1(d) = \lim_{\delta/\sigma \rightarrow \pm\infty} 2\sigma f\left(\frac{\delta}{\sigma}\right) + \delta \left[ \Phi\left(\frac{\delta}{\sigma}\right) - 1 \right] = |\delta|. \quad (11)$$

Interestingly, we found that for the special case  $\delta = 0$  and  $\sigma = 1$ , Eq. 8 becomes the closed form solution of what in statistical process control is referred to as the constant  $d_2^9$ , whose value 1.1284 determined by numerical means is given in textbooks on process control, e.g. [23]. We obtain

$$d_2 = \frac{2}{\sqrt{\pi}}.$$

### 4.3. An Attempt of Finding a Closed-Form Solution

After studying the a-priori alignment, we want to extend the above approach to the unit selection-based alignment.

<sup>9</sup>the expected distance between two instances of a standard normally distributed random process

The principal difference to the former is the minimization in Eq. 1 to search for the optimal sequence of vectors. This directly affects the expected value of  $d$  given  $x$  (cf. Eq. 7), since now we do not have a normal distribution as probability density function of  $y$  but the more complicated term  $p_N(y|x)$

$$E_N(d|x) = \int_{-\infty}^{\infty} |x - y| p_N(y|x) dy. \quad (12)$$

Here,  $N$  denotes the number of feature vectors in the target feature vector sequence  $y_1^N$ , which serves as a pool we select appropriate units from, see Section 2. Again, we assume these vectors to be normally distributed with the parameters  $\mu_y$  and  $\sigma$  and independent of each other.

For each possible  $y$ , we calculate the probability density of the  $n^{\text{th}}$  target feature vector being equal to  $y$  and closest to  $x$ . The sum over all of these vectors from 1 to  $N$  yields the searched density  $p_N(y|x)$ .

To be more detailed: The probability density of the  $n^{\text{th}}$  vector being equal to  $y$  is

$$P_n = \frac{1}{\sigma} f\left(\frac{y - \mu_y}{\sigma}\right).$$

The probability of the  $n^{\text{th}}$  vector being closest to  $x$  means that the distance to all other vectors  $y_\nu$  for  $\nu \in \{1, \dots, N\}$ ,  $\nu \neq n$  is greater than that to  $y_n$  or, given  $y_n = y$ , that  $|y_\nu - x| > |y - x|$

$$\begin{aligned}
Q_n &= p\left(\bigwedge_{\substack{\nu=1 \\ \nu \neq n}}^N |y_\nu - x| > |y - x|\right) \\
&= \prod_{\substack{\nu=1 \\ \nu \neq n}}^N p(|y_\nu - x| > |y - x|) \\
&= p(|\psi - x| > |y - x|)^{N-1} \\
&= \begin{cases} p(\psi < y \vee \psi > 2x - y)^{N-1} & \text{for } y < x \\ p(\psi > y \vee \psi < 2x - y)^{N-1} & \text{otherwise} \end{cases} \\
&= \begin{cases} \left(\Phi\left(\frac{y - \mu_y}{\sigma}\right) + 1 - \Phi\left(\frac{2x - y - \mu_y}{\sigma}\right)\right)^{N-1} & \text{for } y < x \\ \left(1 - \Phi\left(\frac{y - \mu_y}{\sigma}\right) + \Phi\left(\frac{2x - y - \mu_y}{\sigma}\right)\right)^{N-1} & \text{otherwise} \end{cases}
\end{aligned}$$

Here,  $\psi$  is a  $y$ -like distributed random variable replacing  $y_\nu$  for  $\nu \in \{1, \dots, N\}$ ,  $\nu \neq n$ . Accordingly, Eq. 12 becomes (also cf. Eq. 7)

$$E_N(d|x) = \int_{-\infty}^{\infty} |x - y| \sum_{n=1}^N (P_n Q_n) dy \quad (13)$$

$$= \frac{N}{\sigma} \int_{-\infty}^x (x-y) f\left(\frac{y-\mu_y}{\sigma}\right) \left(1 + \Phi\left(\frac{y-\mu_y}{\sigma}\right) - \Phi\left(\frac{2x-y-\mu_y}{\sigma}\right)\right)^{N-1} dy$$

$$- \frac{N}{\sigma} \int_x^{\infty} (x-y) f\left(\frac{y-\mu_y}{\sigma}\right) \left(1 - \Phi\left(\frac{y-\mu_y}{\sigma}\right) + \Phi\left(\frac{2x-y-\mu_y}{\sigma}\right)\right)^{N-1} dy$$

As already mentioned in footnote 7, for  $N = 1$ , this becomes identical to Eq. 7.

One can show that Eq. 13 can be simplified to the problem of solving

$$\int f^m(x) \Phi^n(x-\delta) dx \text{ for } m \in \{1, 2\} \text{ and } n \in \{0, 1, 2, \dots\},$$

whose closed-form solution we do not know for  $n > 1$ .

However, if we consider large values of  $\frac{\delta}{\sigma}$  (cf. Eq. 9), we get a very exact approximation of the searched expected value of  $d$  as already discussed in Eq. 10 by

$$\lim_{\delta/\sigma \rightarrow \pm\infty} E_N(d) = \lim_{\delta/\sigma \rightarrow \pm\infty} E_N(d|\mu_x).$$

Together with Eq. 13, where we substitute  $y$  by  $z + \mu_x$ , we have

$$\lim_{\delta/\sigma \rightarrow \pm\infty} E_N(d) = \quad (14)$$

$$= \lim_{\delta/\sigma \rightarrow \pm\infty} \frac{N}{\sigma} \left[ - \int_{-\infty}^0 z f\left(\frac{z-\delta}{\sigma}\right) \left(1 + \Phi\left(\frac{z-\delta}{\sigma}\right) - \Phi\left(\frac{-z-\delta}{\sigma}\right)\right)^{N-1} dz \right.$$

$$\left. + \int_0^{\infty} z f\left(\frac{z-\delta}{\sigma}\right) \left(1 - \Phi\left(\frac{z-\delta}{\sigma}\right) + \Phi\left(\frac{-z-\delta}{\sigma}\right)\right)^{N-1} dz \right]$$

For  $\frac{\delta}{\sigma} \rightarrow \infty$ , the first integral of this equation becomes zero, since

$$\lim_{\delta/\sigma \rightarrow \infty} f\left(\frac{z-\delta}{\sigma}\right) = 0 \text{ for } z < 0. \quad (15)$$

Furthermore, in the remaining integral, we have

$$\lim_{\delta/\sigma \rightarrow \infty} \Phi\left(\frac{-z-\delta}{\sigma}\right) = 0 \text{ for } z > 0,$$

and, taking into account the observation in Eq. 15, we are allowed to extend the remaining integral's lower limit to  $-\infty$  since this adds zero. Consequently, for  $\frac{\delta}{\sigma} \rightarrow \infty$ , we can express Eq. 14 by

$$\lim_{\delta/\sigma \rightarrow \infty} E_N(d) =$$

$$\lim_{\delta/\sigma \rightarrow \infty} \frac{N}{\sigma} \int_{-\infty}^{\infty} z f\left(\frac{z-\delta}{\sigma}\right) \left(1 - \Phi\left(\frac{z-\delta}{\sigma}\right)\right)^{N-1} dz.$$

Applying the above steps to the case  $\frac{\delta}{\sigma} \rightarrow -\infty$  and using the relations  $\Phi(x) = 1 - \Phi(-x)$  and  $f(x) = f(-x)$ , we have

$$\lim_{\delta/\sigma \rightarrow \pm\infty} E_N(d) =$$

$$\lim_{\delta/\sigma \rightarrow \pm\infty} \frac{N}{\sigma} \int_{-\infty}^{\infty} z f\left(\frac{z-|\delta|}{\sigma}\right) \left(1 - \Phi\left(\frac{z-|\delta|}{\sigma}\right)\right)^{N-1} dz.$$

Substituting  $z$  by  $|\delta| - \sigma\xi$  yields<sup>10</sup>

$$\lim_{\delta/\sigma \rightarrow \pm\infty} E_N(d) =$$

$$= \frac{N}{\sigma} \int_{-\infty}^{\infty} (|\delta| - \sigma\xi) f(-\xi) (1 - \Phi(-\xi))^{N-1} (-\sigma d\xi)$$

$$= N \int_{-\infty}^{\infty} (|\delta| - \sigma\xi) f(\xi) \Phi(\xi)^{N-1} d\xi$$

$$= N|\delta| \left[ \frac{\Phi(\xi)^N}{N} \right]_{\xi=-\infty}^{\infty} - \sigma N \int_{-\infty}^{\infty} \xi f(\xi) \Phi(\xi)^{N-1} d\xi$$

$$= |\delta| - \sigma\mu(N). \quad (16)$$

The structure of this formula gives a qualitative overview about some of the expected value's characteristics. We have the term  $|\delta|$ , which is independent of the standard deviation  $\sigma$  and a term, which is a constant with respect to  $|\delta|$  but linearly depends on  $\sigma$ .

[24] gives solutions to  $\mu(N)$  for  $N \in \{1, \dots, 5\}$ , but so far, for  $k > 5$ , we did not succeed in finding a closed form. Table 3 gives some example values of  $\mu(N)$ , and Figure 4 shows a plot of  $E_N(d)$  as a function of  $\delta$  for several values of  $N$ . We see that for large  $\delta$ , the graphs approach  $|\delta| - \sigma\mu(N)$  as derived in Eq. 16, and for  $N = 1$ , one obtains the special case discussed in Eqs. 9 and 11.

Although we are not able to find a solution to  $\mu(N)$  for an arbitrary  $N$ , we can derive partial results exploiting the symmetries of  $f$  and  $\Phi$

<sup>10</sup>As the quotient  $\frac{\delta}{\sigma}$  disappears here, we can remove the limit in the following expressions.

| $N$       | $\mu(N)$ | closed form of $\mu(N)$  |
|-----------|----------|--|
| 1         | 0        | 0  |
| 2         | 0.56     | $\frac{1}{\sqrt{\pi}}$   |
| 3         | 0.85     | $\frac{3}{2\sqrt{\pi}}$  |
| 4         | 1.03     | $\frac{6}{\sqrt{\pi^3}} \arctan \sqrt{2}$                          |
| 5         | 1.16     | $\frac{15}{\sqrt{\pi^3}} \arctan \sqrt{2} - \frac{5}{2\sqrt{\pi}}$ |
| 10        | 1.54     |  |
| 100       | 2.51     |  |
| 1 000     | 3.24     |  |
| 1 000 000 | 4.86     |  |

**Table 3.** The offset constant  $\mu(N)$  for different values of  $N$ .

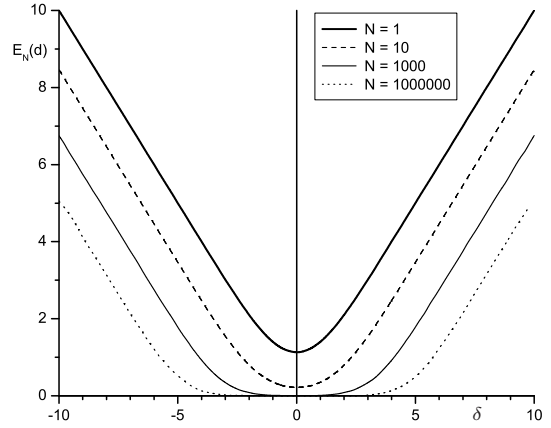
$$\begin{aligned}
\mu(N) &= N \int_{-\infty}^{\infty} x f(x) \Phi(x)^{N-1} dx \\
&= -N \int_{-\infty}^{\infty} \xi f(\xi) (1 - \Phi(\xi))^{N-1} d\xi \\
&= -N \sum_{k=0}^{N-1} (-1)^k \binom{N-1}{k} \int_{-\infty}^{\infty} \xi f(\xi) \Phi(\xi)^k d\xi \\
&= -\sum_{k=0}^{N-1} (-1)^k \binom{N}{k+1} \mu(k+1) \\
&= \sum_{\kappa=1}^{N-1} (-1)^\kappa \binom{N}{\kappa} \mu(\kappa) + (-1)^N \mu(N).
\end{aligned}$$

This finally yields

$$\mu(N) = \frac{1}{2} \sum_{k=1}^{N-1} (-1)^k \binom{N}{k} \mu(k) \text{ for } N \in \{1, 3, \dots\}.$$

This formula enables us to recursively compute  $\mu(N)$  from  $\mu(1), \dots, \mu(N-1)$ ; unfortunately, it holds only for odd  $N$ , so we would not be able to find a general statement unless we find a description for even  $N$ .

However, there is a way to study the behaviour of  $\mu(N)$



**Fig. 4.** Expected value of the minimum distance between a source feature and  $N$  target feature vectors  $E_N(d)$  as a function of the difference of the distribution means  $\delta$  and the number of available target feature vectors  $N$ ;  $\sigma = 1$ .

when  $N$  approaches infinity

$$\begin{aligned}
\lim_{N \rightarrow \infty} \mu(N) &= \lim_{N \rightarrow \infty} N \int_{-\infty}^{\infty} x f(x) \Phi(x)^{N-1} dx \\
&= \lim_{N \rightarrow \infty} N \left[ \int_{-\infty}^0 x f(x) \Phi(x)^{N-1} dx \right. \\
&\quad \left. + \int_0^{\infty} x f(x) \Phi(x)^{N-1} dx \right]. \quad (17)
\end{aligned}$$

We know that  $f(x)$  is an even function, and  $x$  is odd, so  $xf(x)$  is also odd.  $\Phi(x)$  is a strictly positive, monotonous, and bounded function, hence we know that  $\int_{-\infty}^0 xf(x)\Phi(x)^{N-1} dx < 0$  and  $\int_0^{\infty} xf(x)\Phi(x)^{N-1} dx > 0$ . Consequently, both terms become smaller, if we replace  $g_N(x) = \Phi(x)^{N-1}$  by a function  $h_N(x)$ , which is greater than the former for  $x < 0$  and smaller for  $x > 0$

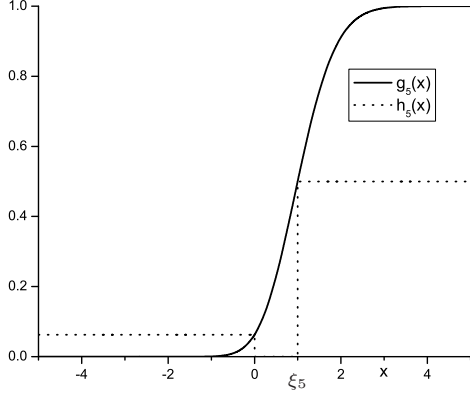
$$h_N(x) = \begin{cases} g_N(0) = \frac{1}{2^{N-1}} & \text{for } x \leq 0 \\ 0 & \text{for } 0 < x \leq \xi_N \\ \frac{1}{2} & \text{for } \xi_N < x \end{cases}$$

Here,  $\xi_N > 0$ <sup>11</sup> is the position, where  $g_N(x)$  becomes  $\frac{1}{2}$

$$\xi_N = \Phi^{-1}\left(2^{\frac{1}{1-N}}\right). \quad (18)$$

Figure 5 shows an example of the functions  $g_N(x)$  and  $h_N(x)$ .

<sup>11</sup>This relation is only true for  $N > 2$ , which is, however, no additional constraint, since we want to investigate the limit for  $N \rightarrow \infty$ .



**Fig. 5.** Example of the functions  $g_N(x)$  and  $h_N(x)$  for  $N = 5$ ;  $\xi_5 = 0.998$ , cf. Eq. 18.

Applying the definition of  $h_N(x)$  to Eq. 17 yields

$$\begin{aligned}
\lim_{N \rightarrow \infty} \mu(N) &\geq \lim_{N \rightarrow \infty} N \int_{-\infty}^{\infty} x f(x) h_N(x) dx \\
&= \lim_{N \rightarrow \infty} \left[ \frac{N}{2^{N-1}} \int_{-\infty}^0 x f(x) dx \right. \\
&\quad \left. + \frac{N}{2} \int_{\xi_N}^{\infty} x f(x) dx \right] \\
&= \lim_{N \rightarrow \infty} \left[ \frac{-N}{2^{N-1} \sqrt{2\pi}} + \frac{N f(\xi_N)}{2} \right] \\
&= \frac{1}{2} \lim_{N \rightarrow \infty} N f(\xi_N). \tag{19}
\end{aligned}$$

Using Eq. 18, we can express  $N$  as a function of  $\xi_N$

$$N = 1 - \frac{\log 2}{\log \Phi(\xi_N)}, \tag{20}$$

and we observe that when  $N$  approaches infinity, also  $\xi_N$  approaches infinity. Consequently, we are allowed to rewrite Eq. 19 as

$$\begin{aligned}
\lim_{N \rightarrow \infty} \mu(N) &\geq \frac{1}{2} \lim_{\xi_N \rightarrow \infty} f(\xi_N) \left( 1 - \frac{\log 2}{\log \Phi(\xi_N)} \right) \\
&= -\frac{\log 2}{2} \lim_{\xi_N \rightarrow \infty} \frac{f(\xi_N)}{\log \Phi(\xi_N)}.
\end{aligned}$$

Application of l'Hôpital's rule produces

$$\begin{aligned}
\lim_{N \rightarrow \infty} \mu(N) &\geq -\frac{\log 2}{2} \lim_{\xi_N \rightarrow \infty} \frac{-\xi_N f(\xi_N)}{\frac{1}{\Phi(\xi_N)} f(\xi_N)} \\
&= \frac{\log 2}{2} \lim_{\xi_N \rightarrow \infty} \xi_N \Phi(\xi_N) \\
&= \infty.
\end{aligned}$$

Hence, the limit value for  $\mu(N)$  is infinity if  $N$  approaches infinity. However, this means that for very large  $N$ , the approximation Eq. 16 is not useful, since the expected value of  $d$  is non-negative. Consequently, when  $N$  approaches infinity, we must not apply the simplifications derived in Eq. 10, but have to consider the original definition of  $d$ 's expected value (cf. Eq. 6)

$$\lim_{N \rightarrow \infty} E_N(d) = \lim_{N \rightarrow \infty} \frac{1}{\sigma} \int_{-\infty}^{\infty} E_N(d|x) f\left(\frac{x - \mu_x}{\sigma}\right) dx,$$

where  $E_N(d|x)$  is declared in Eq. 13. Several substitutions and the application of Lebesgue's dominated convergence theorem considering the fact that  $E_1(d)$  is finite leads to the result<sup>12</sup>

$$\lim_{N \rightarrow \infty} E_N(d) = 0.$$

This can be regarded as a proof of a special case of the speech alignment paradox taking the above formulated conditions into account.

## 5. CONCLUSION

For the highly computational complexity of text-independent speech alignment based on unit selection, we were not able to investigate the speech alignment paradox by means of very large amounts of data. This was the reason for applying a mathematical model describing two speech samples by means of Gaussian mixture models. For a special case, we could derive a mathematical proof of the paradox. Future work is to focus on the generalization of this paper's investigations.

## 6. REFERENCES

- [1] C. Myers and L. Rabiner, "A Comparative Study of Several Dynamic Time-Warping Algorithms for Connected Word Recognition," *Bell System Technical Journal*, vol. 60, no. 7, 1981.
- [2] S. Furui, "Cepstral Analysis Technique for Automatic Speaker Verification," *IEEE Trans. on Acoustics, Speech and Signal Processing*, vol. 29, no. 2, 1981.
- [3] D. Berndt and J. Clifford, "Using Dynamic Time Warping to Find Patterns in Time Series," in *Proc. of the AAAI Workshop on Knowledge Discovery in Databases*, Seattle, USA, 1994.

<sup>12</sup>Again, the authors would be happy to provide the details of this proof on inquiry.



- [4] C. Ratanamahatana and E. Keogh, "Everything you Know about Dynamic Time Warping Is Wrong," in *Proc. of the Workshop on Mining Temporal and Sequential Data*, Seattle, USA, 2004.
- [5] S. Young, P. Woodland, and W. Byrne, *The HTK Book, Version 1.5*, Cambridge University, Cambridge, UK, 1993.
- [6] D. Sündermann, H. Höge, A. Bonafonte, H. Ney, A. Black, and S. Narayanan, "Text-Independent Voice Conversion Based on Unit Selection," in *Proc. of the ICASSP'06*, Toulouse, France, 2006.
- [7] D. Sündermann, H. Höge, A. Bonafonte, H. Ney, and J. Hirschberg, "Text-Independent Cross-Language Voice Conversion," in *Proc. of the Interspeech'06*, Pittsburgh, USA, 2006.
- [8] "Methods for Subjective Determination of Transmission Quality," Tech. Rep. ITU-T Recommendation P.800, ITU, Geneva, Switzerland, 1996.
- [9] S.-H. Chen, S.-J. Chen, and C.-C. Kuo, "Perceptual Distortion Analysis and Quality Estimation of Prosody-Modified Speech for TD-PSOLA," in *Proc. of the ICASSP'06*, Toulouse, France, 2006.
- [10] P. Boersma, "Praat, a System for Doing Phonetics by Computer," *Glott International*, vol. 5, no. 9/10, 2001.
- [11] H. Ye and S. Young, "Perceptually Weighted Linear Transformations for Voice Conversion," in *Proc. of the Eurospeech'03*, Geneva, Switzerland, 2003.
- [12] A. Hunt and A. Black, "Unit Selection in a Concatenative Speech Synthesis System Using a Large Speech Database," in *Proc. of the ICASSP'96*, Atlanta, USA, 1996.
- [13] H. Höge, "Project Proposal TC-STAR - Make Speech to Speech Translation Real," in *Proc. of the LREC'02*, Las Palmas, Spain, 2002.
- [14] R. Duda and P. Hart, *Pattern Classification and Scene Analysis*, Wiley, New York, USA, 1973.
- [15] M. Zissman, "Automatic Language Identification Using Gaussian Mixture and Hidden Markov Models," in *Proc. of the ICASSP'93*, Minneapolis, USA, 1993.
- [16] Y. Stylianou, O. Cappé, and E. Moulines, "Statistical Methods for Voice Quality Transformation," in *Proc. of the Eurospeech'95*, Madrid, Spain, 1995.
- [17] D. Reynolds, "Automatic Speaker Recognition Using Gaussian Mixture Models," *Lincoln Laboratory Journal*, vol. 8, no. 2, 1995.
- [18] R. Faltlhauser, T. Pfau, and G. Ruske, "On-Line Speaking Rate Estimation Using Gaussian Mixture Models," in *Proc. of the ICASSP'00*, Istanbul, Turkey, 2000.
- [19] S. Tranter and D. Reynolds, "Speaker Diarisation for Broadcast News," in *Proc. of the Odyssey 2004 Speaker and Language Recognition Workshop*, Toledo, Spain, 2004.
- [20] A. Kain and M. Macon, "Design and Evaluation of a Voice Conversion Algorithm Based on Spectral Envelope Mapping and Residual Prediction," in *Proc. of the ICASSP'01*, Salt Lake City, USA, 2001.
- [21] D. Sündermann, A. Bonafonte, H. Ney, and H. Höge, "A Study on Residual Prediction Techniques for Voice Conversion," in *Proc. of the ICASSP'05*, Philadelphia, USA, 2005.
- [22] D. Sündermann, A. Bonafonte, H. Ney, and H. Höge, "A First Step Towards Text-Independent Voice Conversion," in *Proc. of the ICSLP'04*, Jeju Island, South Korea, 2004.
- [23] D. Montgomery, *Introduction to Statistical Quality Control*, Wiley, New York, USA, 1996.
- [24] H. David and H. Nagaraja, *Order Statistics*, Wiley, New York, USA, 2003.