



HAL
open science

A Comparative Study of Glottal Source Estimation Techniques

Thomas Drugman, Baris Bozkurt, T. Dutoit

► **To cite this version:**

Thomas Drugman, Baris Bozkurt, T. Dutoit. A Comparative Study of Glottal Source Estimation Techniques. *Computer Speech and Language*, 2011, 26 (1), pp.20. 10.1016/j.csl.2011.03.003. hal-00770346

HAL Id: hal-00770346

<https://hal.science/hal-00770346>

Submitted on 5 Jan 2013

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

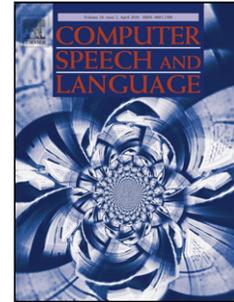
L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Accepted Manuscript

Title: A Comparative Study of Glottal Source Estimation Techniques

Authors: Thomas Drugman, Baris Bozkurt, Thierry Dutoit

PII: S0885-2308(11)00021-0
DOI: doi:10.1016/j.csl.2011.03.003
Reference: YCSLA 487



To appear in:

Received date: 25-3-2010
Revised date: 8-3-2011
Accepted date: 10-3-2011

Please cite this article as: Drugman, T., Bozkurt, B., Dutoit, T., A Comparative Study of Glottal Source Estimation Techniques, *Computer Speech & Language* (2010), doi:10.1016/j.csl.2011.03.003

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

A Comparative Study of Glottal Source Estimation Techniques

Thomas Drugman^a, Baris Bozkurt^b, Thierry Dutoit^a

^a*TCTS Lab, University of Mons, 31 Boulevard Dolez, 7000 Mons, Belgium*

^b*Department of Electrical & Electronics Engineering, Izmir Institute of Technology, Gulbahce Koyu 35430, Urla, Izmir, Turkey*

Abstract

Source-tract decomposition (or glottal flow estimation) is one of the basic problems of speech processing. For this, several techniques have been proposed in the literature. However studies comparing different approaches are almost nonexistent. Besides, experiments have been systematically performed either on synthetic speech or on sustained vowels. In this study we compare three of the main representative state-of-the-art methods of glottal flow estimation: closed-phase inverse filtering, iterative and adaptive inverse filtering, and mixed-phase decomposition. These techniques are first submitted to an objective assessment test on synthetic speech signals. Their sensitivity to various factors affecting the estimation quality, as well as their robustness to noise are studied. In a second experiment, their ability to label voice quality (tensed, modal, soft) is studied on a large corpus of real connected speech. It is shown that changes of voice quality are reflected by significant modifications in glottal feature distributions. Techniques based on the mixed-phase decomposition and on a closed-phase inverse filtering process turn out to give the best results on both clean synthetic and real speech signals. On the other hand, iterative and adaptive inverse filtering is recommended in noisy environments for its high robustness.

Keywords:

Source-tract Separation, Glottal Flow Estimation, Inverse Filtering, Mixed-Phase Decomposition, Voice Quality

Corresponding author: Thomas Drugman, *thomas.drugman@umons.ac.be*

Research Highlights

- CPIF and CCD give the best results on synthetic speech
- NAQ, H1-H2 and HRF are useful glottal parameters
- Glottal flow can be used for voice quality analysis
- Methods are tested on a large speech database

Accepted Manuscript

The authors would like to thank again the reviewers for their time and their fruitful suggestions.

Reviewer #1: The revised paper has obviously been improved. I think that this paper can be published without further review but I recommend that the authors consider the following points.

Section 4.4:

The revised descriptions would still confuse the readers.

In this section, the authors don't show any results about a relationship between the glottal formant parameter and the estimation accuracy. However, the authors say "Throughout our experiments we confirmed for all glottal source estimation techniques the performance degradation as a function of these factors."

I understand the glottal formant affects the estimation accuracy but it is not shown in this experiment.

Therefore, the glottal formant should not be included in "these factors."

The authors had better further revise this section to clarify what the authors evaluated.

So as to clarify this point, the first and second paragraph of section 4.4 (page 17) have been modified:

"The stronger this interference, the more important the time overlap between the minimum-phase component and the maximum-phase response of the next glottal cycle, and consequently the more difficult the decomposition."

And:

"A strong interference then appears with high pitch, and with low F_1 and F_g values. The previous experiments confirmed for all glottal source estimation techniques the performance degradation as a function of F_0 and F_1 . Although we did not explicitly measure the sensitivity of these techniques to F_g in this manuscript, it was confirmed in other informal experiments we performed."

Section 2.1:

iterative => iterative

Thanks.

Reviewer #2: The revised manuscript corrects most of the issues that I pointed out in my review. I believe that the current version can be accepted for publication.

Reviewer #1: This paper presents an evaluation of three typical glottal flow estimation methods: closed-phase inverse filtering (CPIF); iterative and adaptive inverse filtering (IAIF); and complex cepstrum decomposition (CCD). These methods are compared with each other in terms of robustness to additive noise, sensitivity to F_0 , and sensitivity to the vocal tract transfer function in the experiments with speech samples synthesized with the LF model. Interesting properties of each estimation method are shown from the experimental results. Moreover, another experimental evaluation using real speech samples is also conducted. To evaluate the estimation performance, speech samples in three different speaking styles are analyzed, and then the parameters of the glottal flow estimated by the individual estimation methods are compared with each other. The experimental result demonstrates that CCD estimates the most discriminative style-dependent glottal flow compared with the others. This paper is relatively well written. The presented experimental results are interesting and worthwhile to be published.

There are some unclear parts. I recommend that the authors revise the following points to further improve the paper.

*** General remark ***

Discussions in section 4:

It is demonstrated from the presented experimental results which estimation method is robust or sensitive to some conditions. However, it is still unclear the reason why these results are caused: e.g., what is the reason why IAIF is more effective than the others in noisy conditions but it is less effective in clean conditions? Please add more discussions about the possible reasons to interpret the presented results.

Thanks for your comment. We have extended Section 4.1 so as to give more explanations about the reasons for which one method is more robust or sensitive to noise than another one.

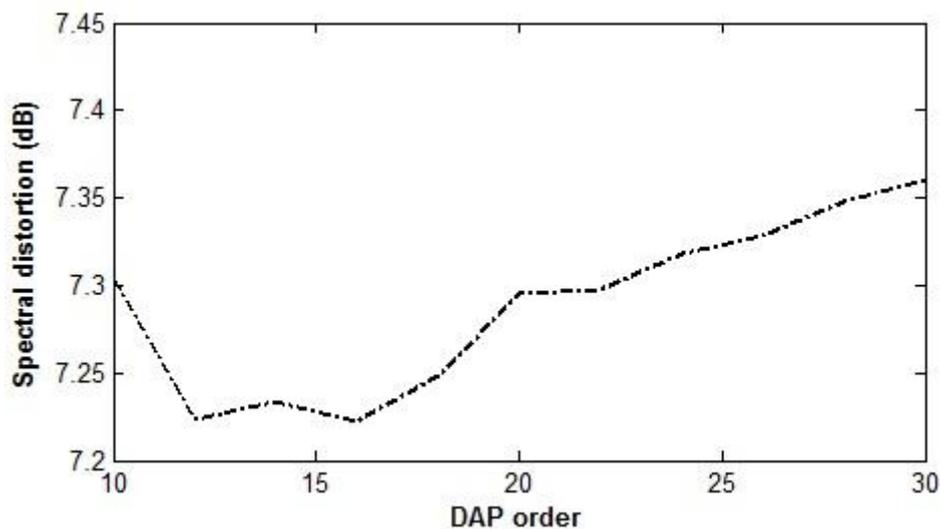
Effect of the order of the inverse filter:

How did you determine the order of the inverse filter in CPIF and IAIF? The authors have to describe it. The order of synthesis filter used for synthesizing speech samples in section 4 should also be described. Does the order of the inverse filter strongly affect the results of CPIF and IAIF shown in figures 4 and 5? Is it possible that significant improvements of the estimation performance are yielded if the order is carefully adjusted? Please add a discussion about the sensitivity to the setting of the order of inverse filtering as well.

For IAIF, as mentioned in the paper (end of Section 2.1.2), we used the implementation from the toolbox available on the TKK Aparat website with its default options. Since the authors of IAIF have probably optimized their method, we used their implementation, without modifying the code.

As for CPIF, the order used for DAP analysis was fixed to 18 ($=F_s/1000+2$), as commonly used in the literature. In order to investigate the improvements that can be obtained by carefully adjusting the analysis order, the following figure exhibits the evolution of the overall Spectral Distortion as a function of the DAP order. It can be observed the best performance is obtained using an order comprised between 12 and 18. However, the resulting improvement is extremely weak (in the range $[10,20]$, only a difference of 0.07dB can be reached by optimizing the order), and insignificant with regard to the gap separating the different methods compared in this work. Furthermore, this

Figure was obtained for synthetic signals, and this order optimization may not hold for real speech signals. More information about the filter order has been provided in the manuscript, in the very end of Section 2.1.1.



Experimental results in section 4.3:

Only four typical vowels are considered in the evaluation. Why did you choose only them? The number of data points shown in figure 6 seems too small. If the first formant frequency is an effective parameter for capturing an impact of the vocal tract transfer function on the glottal flow estimation, the same tendencies as observed in these four vowels would also be observed in the other vowels. It seems better to show more data points if possible. If more data points would make the tendencies invisible, other parameters as well as the first formant frequency have also to be considered.

Thanks for this suggestion. Section 4.3 and Figure 6 have been substantially modified so as to give a more complete answer to these questions.

*** The other comments ***

Explanation about DAP in section 2.1.1:

The 2nd sentence in the 2nd paragraph explains DAP but this explanation is insufficient due to lack of precision. A more precise explanation is required.

Thanks. More information about DAP has been added in the second paragraph of Section 2.1.1.

Experimental conditions in section 4:

Please clarify sampling frequency of synthetic speech samples. In the last sentence in the 1st paragraph, "20.000 experiments" should be "20,000 experiments."

Thanks for this remark. Changes have been made in the manuscript.

Figure 6 in section 4.3:

Is this figure correct? The F1 values should be 336, 378, 466, and 827 Hz but data points at 466 Hz don't seem to be plotted correctly.

Thanks for the correction. We have probably made a mistake when plotting these points. Nevertheless, Figure 6 has now been totally modified.

Section 4.4:

The authors mention that the glottal formant is used as one of the parameters to capture interference. Is this parameter controlled in the additive noise experiment? Please clearly describe this point. Some examples showing the waveform (or spectral) change caused by varying each parameter would be helpful for readers to understand how each parameter affects the glottal flow or the impulse response of the vocal tract transfer function. Please show them if possible.

Changes in the text (third enumeration of Section 4.4) have been brought so as to clarify the effect of the glottal formant on the interference. Due to limitation space, we think that showing the influence of each parameter is impossible for this document. However, we refer to (Doval and d'Alessandro, 2006; or others) which discuss this point more in depth.

Reviewer #2: This is an important paper that aims at comparing 3 different techniques for the estimation of the glottal flow relying solely on the acoustic speech signal. I find the paper in general well written and organized. I also did not find any methodological errors.

My main critique of the work presented in the paper regards the limited number of robustness tests that were conducted when analyzing the performance of the algorithms on synthetic speech. There are many well known problems in voice recordings that affect the estimation of the glottal parameters and the evaluation of the robustness of the algorithms just by adding white Gaussian noise to the speech signal only accounts for a very small part of these problems. There are two types of problems: the ones related with the glottal flow signal and the ones related with the recording conditions. In the first case, the jitter is, in my view, the major challenge for these algorithms especially for the ones that require more than one glottal cycle for each analysis frame. Jitter is a very common effect in real speech signals, even for non-pathological voices. The amount of jitter is very good discriminator for detecting some types of voice pathologies. Another important effect is the case of breathy voices in which there is noise added not to the speech signal but to the glottal waveform. This is an amplitude modulated noise that also affects the robustness of most algorithms. The problems related with recording condition are just partially modeled by the additive noise model that was used. A major difficulty in real cases is the lack of fixed distance from the mouth to the microphone and the presence and non-stationary noise. This is what happens in speech acquired in hospital environment, in which each patient is recorded in different conditions. These problems should at least be referred.

Thank you for this fruitful comment. The end of the introduction of Section 4 (last sentences of the first paragraph) has been modified so as to refer to this issue.

It is true that this study, for simply practical reasons, does not cover the whole diversity of glottal production. This is now specified in the modified text. Our concentration is more on the real speech tests, which are comparatively broad.

Another question is the selection of the voice qualities for the study in real speech. The discrimination between Loud, Modal and Soft qualities is not one of hardest to perform. Using these 3 it would have been easy to predict that the Normalized Amplitude Quotient would be best discriminator. The results of section 5 thus are not very interesting.

This comment is surprising to us since Section 5 contains very large tests with real speech that has not been performed before. We frankly think that it is a very valuable contribution of this paper. We therefore believe that without this test some information would be missing for the readers with comparatively less expertise than Reviewer 2. Besides, it was not so clear to us, before performing the test itself, that NAQ would lead to the best discrimination. At least, we hope that, for the reader who could predict that, the results of Section 5 will give an objective confirmation led on a large dataset of what he intuitively expected.

There are also some points that, in my view, should be improved.

1. In the introduction (page 2), regarding the need for source-filter deconvolution. It is said that:

In many speech processing applications, it is important to separate the contributions from the glottis and the vocal tract. Achieving such a source-filter deconvolution could lead to a distinct characterization and modeling of these two components, as well as to a better understanding of the human phonation process.

This misses the most important reason that is the completely different dynamics of the two contributions.

Thanks. We have modified our introduction so as to account for this remark.

2. In the section regarding the Closed Phase Inverse Filtering technique (page 4-5):

Closed phase refers to the timespan during which the glottis is closed (see Figure 1). During this period, the vocal tract is supposed to be free of any excitation.

That is not exactly true, since the vocal tract is generally considered a system with memory and the effects of excitation open-phase are still present. The reason for using the closed phase is because during that part of the glottal cycle, the effects of the sub-glottal cavities are minimized, and thus providing a better estimate of the vocal tract transfer function.

Thank you for this comment. The beginning of Section 2.1.1 has been modified accordingly.

3. In figure 2 (page 8) the time axis of the time domain graphics are not correct. The impulse response of the vocal tract should start at $n=0$ and not at $n=250$ (more or less), as presented. This is the only way to have the depicted response speech signal.

Thank you. Figure 2 has been modified consequently.

4. In page 10 the definition of the QOQ is not correct:

Avoiding this problem and preferred to the traditional Open Quotient, the Quasi-Open Quotient (QOQ) was proposed as a parameter describing the relative open time of the glottis (Hacki, 1989). It is defined as the ratio between the quasi-open time and the quasi-closed time of the glottis, and corresponds to the timespan (normalized to the pitch period) during which the glottal flow is 50% above the minimum flow.

The QOQ was defined for the EGG signal and not for the flow itself. In this case there is no minimum flow but a variation in conductivity for which there is no zero value. I agree that the measure can be used for the glottal flow but the 50% above minimum does not make sense. It should be the normalized time span in which the flow is above 50% of the difference between the maximum and minimum flow. This is the way the measure is computed in the aparat matlab package that I believed the authors used.

Thank you for this comment. The definition of QOQ has been changed so as to reflect your suggestion.

5. Again in page 10, regarding the normalized amplitude quotient, NAQ:

It is defined as the ratio between the maximum of the glottal flow and the minimum of its derivative, normalized with respect to fundamental frequency. I believe that it is more correct to say that NAQ is normalized with respect to the fundamental period, since it is this values that goes in the denominator of the definition. In this case and in the previous one, I would have found useful to have the definition of both measures expressed as a mathematical equation.

Thanks. The definition of NAQ has been changed according to your suggestion.

6. Regarding the experiments on Synthetic Speech, in page 12, some important details are missing. The first one is the sampling frequency that was used. I assume that the experiments were conducting using the same sampling frequency used in the real speech experiments (16K) of the next section but it must be stated. Also missing is the order of the auto-regressive filter used for synthesis.

This information has been added in the beginning of Section 4.

7. Regarding the measures studied I think that the "determination rate" is a strange one. The correct measure must be the "error rate" on the estimation of NAQ and QOQ. Besides the meaning be more clearly stated, it also change the grading so that higher values are worst. This also makes these two measures compatible with the third one, the spectral distortion. In Fig 4 and Fig 5, the y axis of the first to graphs are "the higher the better" whilst the last one is "the lower the better". This makes the data harder to interpret.

As suggested, the use of the 'determination rate' has been replaced by the 'error rate' in Section 4.

8. In the conclusion (page 20), the sentence following sentence is to strong:

For the first time in the literature, a large real speech corpus was used for testing, without limiting analysis to sustained vowels. It was not the effectiveness of the glottal flow estimation techniques that was used but its ability to discriminate 3 voice qualities. This does not mean that the parameters were correctly estimated, and it is not even a good indication on the relative quality of the 3 methods.

Thanks for this comment. We made this sentence softer in the text. And indeed it is true that these results do not explicitly reflect the performance of the method for the goal of glottal flow estimation. However, it can be expected that, since the actual glottal behavior is known to significantly vary across the different voice qualities, these changes should be highlighted after parametrization of the estimates of the glottal flow. So, although our tests do not explicitly assess the performance of the methods, we think they perform contribute to their implicit assessment, specifically in the context of the unavailability of ground truth data.

Accepted Manuscript

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

A Comparative Study of Glottal Source Estimation Techniques

Thomas Drugman^a, Baris Bozkurt^b, Thierry Dutoit^a

^aTCTS Lab, University of Mons, 31 Boulevard Dolez, 7000 Mons, Belgium

^bDepartment of Electrical & Electronics Engineering, Izmir Institute of Technology, Gulbahce Koyu 35430, Urla, Izmir, Turkey

Abstract

Source-tract decomposition (or glottal flow estimation) is one of the basic problems of speech processing. For this, several techniques have been proposed in the literature. However studies comparing different approaches are almost nonexistent. Besides, experiments have been systematically performed either on synthetic speech or on sustained vowels. In this study we compare three of the main representative state-of-the-art methods of glottal flow estimation: closed-phase inverse filtering, iterative and adaptive inverse filtering, and mixed-phase decomposition. These techniques are first submitted to an objective assessment test on synthetic speech signals. Their sensitivity to various factors affecting the estimation quality, as well as their robustness to noise are studied. In a second experiment, their ability to label voice quality (tensed, modal, soft) is studied on a large corpus of real connected speech. It is shown that changes of voice quality are reflected by significant modifications in glottal feature distributions. Techniques based on the mixed-phase decomposition and on a closed-phase inverse filtering process turn out to give the best results on both clean synthetic and real speech signals. On the other hand, iterative and adaptive inverse filtering is recommended in noisy environments for its high robustness.

Keywords:

Source-tract Separation, Glottal Flow Estimation, Inverse Filtering, Mixed-Phase Decomposition, Voice Quality

1. Introduction

Speech results from filtering the glottal flow by the vocal tract cavities, and converting the resulting velocity flow into pressure at the lips (Quatieri, 2002). In many speech processing applications, it is important to separate the contributions from the glottis and the vocal tract. Achieving such a *source-filter deconvolution* could lead to a distinct characterization and modeling of these two components, as well as to a better understanding of the human phonation process. Such a decomposition is thus a preliminary condition for the study of glottal-based vocal effects, which can be segmental (as for vocal fry), or be controlled by speakers on a separate, supra-segmental layer (as it is the case for the voice quality modifications mentioned in Section 5). Their dynamics is very different from that of the vocal tract contribution, and requires further investigation. Glottal source estimation is then a fundamental problem in speech processing, finding applications in speech synthesis (Cabral et al., 2008), voice pathology detection (Drugman et al., 2009b), speaker recognition (Plumpe et al., 1999), emotion analysis/synthesis (Airas and Alku, 2006), etc.

In this paper, we limit our scope to the methods which perform an estimation of the glottal source contribution directly from the speech waveform. Although some devices such as electroglottographs or laryngographs, which measure the impedance between the vocal folds (but not the glottal flow itself), are informative about the glottal behaviour (Henrich et al., 2004), in most cases the use of such apparatus is inconvenient and only the speech signal is available for analysis. This problem is then a typical case of blind separation, since neither the vocal tract nor the glottal contribution are observable. This also implies that no quantitative assessment of the performance of glottal source estimation techniques is possible on natural speech, as no target reference signal is available.

As one of the basic problems and challenges of speech processing research, glottal flow estimation has been studied by many researchers and various techniques are available in the literature (Walker and Murphy, 2007). However the diversity of algorithms and the fact that the reference for the actual glottal flow is not available often leads to the questionability about relative effectiveness of the methods in real life applications. In most of studies, tests are performed either on synthetic speech or on a few recorded sustained vowels. In addition, very few comparative studies exist (such as (Sturmel et al., 2007)). In this paper, we compare three of the main representative

1
2
3
4
5
6
7
8
9 state-of-the-art methods: closed-phase inverse filtering, iterative and adap-
10 tive inverse filtering, and mixed-phase decomposition. For testing, we first
11 follow the common approach of using a large set of synthetic speech signals
12 (by varying synthesis parameters independently), and then we examine how
13 these techniques perform on a large real speech corpus. In the synthetic
14 speech tests, the original glottal flow is available, so that objective measures
15 of decomposition quality can be computed. In real speech tests the ability of
16 the methods to discriminate different voice qualities (tensed, modal and soft)
17 is studied on a large database (without limiting data to sustained vowels).
18

19
20 The paper is structured as follows. In Section 2 the main state-of-the-art
21 methods for glottal source estimation are reviewed, and the three techniques
22 compared in this study are detailed. Section 3 discusses how the resulting
23 glottal signal can be parametrized both in time and frequency domains. The
24 three methods are evaluated in Section 4 through a wide systematic study
25 on synthetic signals. Their robustness to noise, as well as the impact of
26 the various factors that may affect source-tract separation, are investigated.
27 Section 5 presents decomposition results on a real speech database containing
28 various voice qualities, and shows that the glottal source estimated by the
29 techniques considered in this work conveys relevant information about the
30 phonation type. Finally Section 6 draws the conclusions of this study.
31
32
33
34

35 36 2. Glottal Source Estimation 37

38
39 Glottal flow estimation mainly refers to the estimation of the voiced exci-
40 tation of the vocal tract. During the production of voiced sounds, the airflow
41 arising from the trachea causes a quasi-periodic vibration of the vocal folds
42 (Quatieri, 2002), organized into so-called opening/closure cycles. During the
43 *open phase*, vocal folds are progressively displaced from their initial state due
44 to the increasing subglottal pressure. When the elastic displacement limit
45 is reached, they suddenly return to this position during the so-called *return*
46 *phase*. Figure 1 displays the typical shape of one cycle of the glottal flow
47 (Fig.1(a)) and its time derivative (Fig.1(b)) according to the Liljencrants-
48 Fant (LF) model (Fant et al., 1985). It is often preferred to gather the lip
49 radiation effect (whose action is close to a differentiation operator) with the
50 glottal component, and work in this way with the glottal flow derivative on
51 the one hand, and with the vocal tract contribution on the other hand. It is
52 seen in Figure 1 (bottom plot) that the boundary between open and return
53 phases corresponds to a particular event called the Glottal Closure Instant
54
55
56
57
58

(GCI). GCIs refer to the instances of significant excitation of the vocal tract (Drugman and Dutoit, 2009). Being able to determine their location is of particular importance in so-called pitch-synchronous speech processing techniques, and in particular for a more accurate separation between vocal tract and glottal contributions.

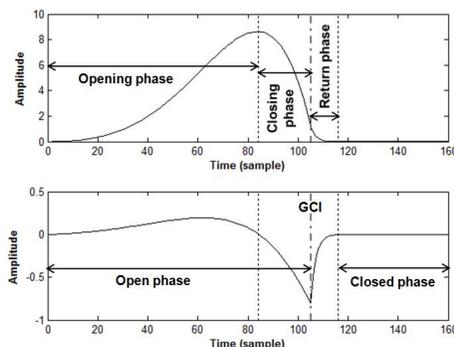


Figure 1: Typical waveforms, according to the Liljencrants-Fant (LF) model, of one cycle of: (top) the glottal flow, (bottom) the glottal flow derivative. The various phases of the glottal cycle, as well as the Glottal Closure Instant (GCI) are also indicated.

The main techniques for estimating the glottal source directly from the speech waveform are now reviewed. Relying on the speech signal alone, as it is generally the case in real applications, allows to avoid the use of intrusive (e.g. video camera at the vocal folds) or inconvenient (e.g. laryngograph) device.

Such techniques can be separated into two classes, according to the way they perform the source-filter separation. The first category (Section 2.1) is based on inverse filtering, while the second one (Section 2.2) relies on the mixed-phase properties of speech.

2.1. Methods based on Inverse Filtering

Most glottal source estimation techniques are based on an inverse filtering process. These methods first estimate a parametric model of the vocal tract, and then obtain the glottal flow by removing the vocal tract contribution via inverse filtering. The methods in this category differ by the way the vocal tract is estimated. In Section 2.1.1 this estimation is computed during the glottal closed phase, while in Section 2.1.2 an iterative and/or adaptive procedure is used. A more extended review of the inverse filtering-based

1
2
3
4
5
6
7
8
9 process for glottal waveform analysis can be found in (Walker and Murphy,
10 2007).

11 12 13 *2.1.1. Closed Phase Inverse Filtering*

14 Closed phase refers to the timespan during which the glottis is closed (see
15 Figure 1). During this period, the effects of the subglottal cavities are min-
16 imized, providing a better way for estimating the vocal tract transfer func-
17 tion. Therefore, methods based on a Closed Phase Inverse Filtering (CPIF)
18 estimate a parametric model of the spectral envelope, computed during the
19 estimated closed phase duration (Wong et al., 1979). The main drawback
20 of these techniques lies in the difficulty in obtaining an accurate determi-
21 nation of the closed phase. Several approaches have been proposed in the
22 literature to solve this problem. In (Veeneman and Bement, 1985), authors
23 use information from the electroglottographic signal (which is avoided in this
24 study) to identify the period during which the glottis is closed. In (Plumpe
25 et al., 1999), it was proposed to determine the closed phase by analyzing the
26 formant frequency modulation between open and closed phases. In (Alku
27 et al., 2009), the robustness of CPIF to the frame position was improved by
28 imposing some dc gain constraints. Besides this problem of accurate deter-
29 mination of the closed phase, it may happen that this period is so short (for
30 high-pitched voices) that not enough samples are available for a reliable filter
31 estimation. It was therefore proposed in (Brookes and Chan, 1994) a tech-
32 nique of multicycle closed-phase LPC, where a small number of neighbouring
33 glottal cycles are considered in order to have enough data for an accurate vo-
34 cal tract estimation. Finally note that an approach allowing non-zero glottal
35 wave to exist over closed glottal phases was proposed in (Deng et al., 2006).

36 In this study, the CPIF-based technique that is used is based on a Dis-
37 crete All Pole (DAP, (Jaroudi and Makhoul, 1991)) inverse filtering process
38 estimated during the closed phase. In order to provide a better fitting of
39 spectral envelopes from discrete spectra (Jaroudi and Makhoul, 1991), the
40 DAP technique aims at computing the parameters of an autoregressive model
41 by minimizing a discrete version of the Itakura-Saito distance (F. Itakura,
42 1970), instead of the time squared error used by the traditional LPC. The use
43 of the Itakura-Saito distance is justified as it is a spectral distortion measure
44 arising from the human hearing perception. The closed phase period is de-
45 termined using the Glottal Opening and Closure Instants (GCIs and GOIs)
46 located by the algorithm detailed in (Drugman and Dutoit, 2009). This algo-
47 rithm has been shown to be effective for reliably and accurately determining
48
49
50
51
52
53
54
55
56
57
58

1
2
3
4
5
6
7
8
9 the position of these events on a large corpus containing several speakers.
10 For tests with synthetic speech, the exact closed phase period is known and
11 is used for CPIF. Note that for high-pitched voices, two analysis windows
12 were used as suggested in (Brookes and Chan, 1994), (Yegnanarayana and
13 Veldhuis, 1998) and (Plumpe et al., 1999). In the rest of the paper, speech
14 signals sampled at 16 kHz are considered, and the order for DAP analysis is
15 fixed to 18 ($=F_s/1000 + 2$, as commonly used in the literature). Through our
16 experiments, we reported that the choice of the DAP order is not critical in
17 the usual range, and that working with an order comprised between 12 and
18 18 leads to sensibly similar results.
19
20
21
22

23 *2.1.2. Iterative and/or Adaptive Inverse Filtering*

24 Some methods are based on iterative and/or adaptive procedures in order
25 to improve the quality of the glottal flow estimation. In (Fu and Murphy,
26 2006), Fu and Murphy proposed to integrate, within the AutoRegressive eX-
27 ogenous (ARX) model of speech production, the LF model of the glottal
28 source. The resulting ARXLF model is estimated via an adaptive and itera-
29 tive optimization (Vincent et al., 2005). Both source and filter parameters are
30 consequently jointly estimated. The method proposed by Moore in (Moore
31 and Clements, 2004) iteratively finds the best candidate for a glottal wave-
32 form estimate within a speech frame, without requiring a precise location of
33 the GCIs. Finally a popular approach was proposed by Alku in (Alku et al.,
34 1992) and called Iterative Adaptive Inverse Filtering (IAIF). This method
35 is based on an iterative refinement of both the vocal tract and the glottal
36 components. In (Alku and Vilkmann, 1994), the same authors proposed an
37 improvement, in which the LPC analysis is replaced by the Discrete All Pole
38 (DAP) modeling technique (Jaroudi and Makhoul, 1991), shown to be more
39 accurate for high-pitched voices.
40
41
42
43
44

45 As a representative technique of this category, the IAIF method proposed
46 by Alku in (Alku et al., 1992) is considered in the rest of this paper. More
47 precisely, we used the implementation of the IAIF method (Airas, 2008) from
48 the toolbox available on the TKK Aparat website ([Online], 2008), with its
49 default options.
50
51

52 *2.2. Mixed-Phase Decomposition*

53 The methods presented in this Section rely on the mixed-phase model
54 of speech (Bozkurt and Dutoit, 2003). According to this model, speech is
55
56
57
58

1
2
3
4
5
6
7
8
9 composed of both minimum-phase (i.e causal) and maximum-phase (i.e anti-
10 causal) components. While the vocal tract impulse response and the glottal
11 *return phase* of the glottal component can be considered as minimum-phase
12 signals, it has been shown in (Doval et al., 2003) that the glottal *open phase*
13 of the glottal flow is a maximum-phase signal. Besides it has been shown in
14 (Gardner and Rao, 1997) that mixed-phase models are appropriate for model-
15 ing voiced speech due to the maximum-phase nature of the glottal excitation.
16 They showed that the use of an anticausal all-pole filter for the glottal pulse
17 is necessary to resolve magnitude and phase information correctly. The key
18 idea of mixed-phase decomposition methods is then to separate minimum
19 from maximum-phase components of speech, where the latter is only due to
20 the glottal contribution.
21
22
23

24 A crucial issue in mixed-phase separation is the weighting window that
25 is applied to the speech signal for short-term analysis. Indeed, since the de-
26 composition is based on phase properties, windowing may have a dramatic
27 influence. It has been shown that GCI-synchronization, as well as the re-
28 spect of some constraints on the window length and function, are essential
29 for guaranteeing a correct decomposition (Drugman et al., 2009a), (Bozkurt
30 et al., 2005). Throughout the rest of this study, we use an appropriate GCI-
31 centered two pitch period-long Blackman window satisfying these conditions.
32
33

34 In previous works, we proposed two approaches achieving such a decom-
35 position: a technique based on the Zeros of the Z-Transform (ZZT, (Bozkurt
36 et al., 2005)), and one based on the Complex Cepstrum Decomposition (CCD,
37 (Drugman et al., 2009a), (Quatieri, 2002)). Both techniques are briefly pre-
38 sented in Sections 2.2.1 and 2.2.2 and depicted in Figure 2. Finally, the
39 methods are shown to be functionally equivalent in Section 2.2.3.
40
41
42

43 2.2.1. Zeros of the Z-Transform (ZZT)

44 For a series of N samples $(x(0), x(1), \dots, x(N - 1))$ taken from a discrete
45 signal $x(n)$, the ZZT representation is defined as the set of roots (zeros)
46 $(Z_1, Z_2, \dots, Z_{N-1})$ of the corresponding z-transform $X(z)$:
47
48
49
50
51
52
53
54
55
56
57
58

$$X(z) = \sum_{n=0}^{N-1} x(n)z^{-n} \quad (1)$$

$$= x(0)z^{-N+1} \prod_{m=1}^{N-1} (z - Z_m) \quad (2)$$

$$= x(0)z^{-N+1} \prod_{k=1}^{M_o} (z - Z_{max,k}) \prod_{k=1}^{M_i} (z - Z_{min,k}) \quad (3)$$

To achieve the ZZT-based decomposition of speech, speech frames are first weighted by a specific window (see above). When computing the ZZT of this signal as in Equation 3, some roots $Z_{max,k}$ fall outside the unit circle. These are due to the maximum-phase (i.e anticausal) component of speech, and are consequently only related to the glottal open phase. On the opposite, roots located inside the unit circle $Z_{min,k}$ are due to the minimum-phase component of speech, i.e mainly to the vocal tract impulse response. Mixed-phase decomposition can then be easily achieved in the ZZT domain, using the unit circle as a discriminant boundary (see Figure 2, third column).

2.2.2. Complex Cepstrum Decomposition (CCD)

The Complex Cepstrum (CC) $\hat{x}(n)$ of a discrete signal $x(n)$ is defined by the following equations (Oppenheim and Schaffer, 1989):

$$X(\omega) = \sum_{n=-\infty}^{\infty} x(n)e^{-j\omega n} \quad (4)$$

$$\log[X(\omega)] = \log(|X(\omega)|) + j\angle X(\omega) \quad (5)$$

$$\hat{x}(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \log[X(\omega)]e^{j\omega n} d\omega \quad (6)$$

where Equations 4, 5, 6 respectively correspond to a Discrete-Time Fourier Transform (DTFT), a complex logarithm and an inverse DTFT (IDTFT). Decomposition in the CC domain arises from the fact that the complex cepstrum $\hat{x}(n)$ of an anticausal (causal) signal is zero for all n positive (negative). Retaining only the negative indexes of the CC makes then it possible to estimate the glottal contribution. The separation in the complex cepstrum

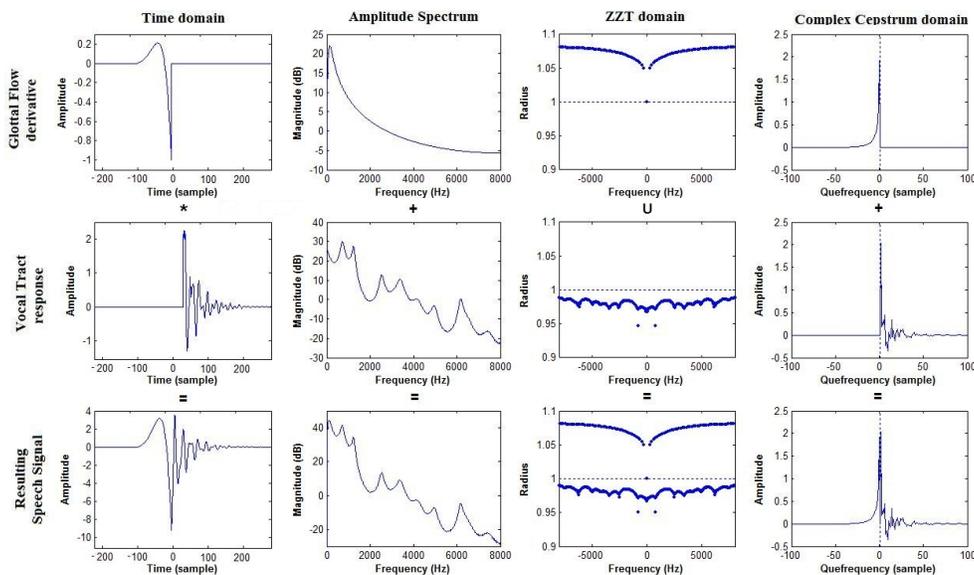


Figure 2: Illustration of mixed-phase decomposition. Rows respectively exhibit the following signals: the glottal flow derivative (top), the vocal tract response (middle), and the resulting speech signal (bottom). Each column corresponds to a domain of representation of these signals: time domain (first column), amplitude spectrum (second column), ZYT representation in polar coordinates (third column), and complex cepstrum domain (fourth column). Interestingly, convolution in the time domain corresponds to the union operator in the ZYT domain and to the addition operator in the complex cepstrum domain. The ZYT and CC domains are suited for achieving mixed-phase decomposition since minimum and maximum-phase components become linearly separable. In the ZYT domain, the unit circle is used as a discriminant boundary, while the quefrency origin is used as a boundary in the complex cepstrum domain.

domain using the quefrency origin as a discriminant boundary is clearly seen in Figure 2, fourth column.

2.2.3. Equivalence between ZYT and CCD

If $X(z)$ is written as in Equation 3, it can be easily shown that the corresponding complex cepstrum can be expressed as (Oppenheim and Schaffer, 1989):

$$\hat{x}(n) = \begin{cases} |x(0)| & \text{for } n = 0 \\ \sum_{k=1}^{M_o} \frac{Z_{max,k}^n}{n} & \text{for } n < 0 \\ \sum_{k=1}^{M_i} \frac{Z_{min,k}^n}{n} & \text{for } n > 0 \end{cases} \quad (7)$$

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

This equation shows the narrow link between the ZZT and the CCD techniques. These two methods can then be seen as two different ways of performing the same operation: separate the minimum and maximum-phase components from a given z-transform $X(z)$. Nevertheless, although functionally equivalent, it has been shown (Pedersen et al., 2010), (Drugman et al., 2009a) that CCD performs much faster (speed is increased between 10 and 100 times for a sampling rate of 16 kHz, depending on the pitch period). This may be explained by the fact that it only relies on FFT and IFFT operations while ZZT requires the factoring of high-order polynomials.

As a method achieving mixed-phase separation, the CCD is considered in the rest of this paper for its higher computational speed. To guarantee good mixed-phase properties (Drugman et al., 2009a), GCI-centered two pitch period-long Blackman windows are used. For this, GCIs were located on real speech using the technique we proposed in (Drugman and Dutoit, 2009). CC is calculated as explained in Section 2.2.2 and FFT is computed on a sufficiently large number of points (typically 4096), which facilitates phase unwrapping.

3. Glottal Source Parametrization

Once the glottal signal has been estimated by any of the aforementioned algorithms, it is interesting to derive a parametric representation of it, using a small number of parameters. Various approaches, both in the time and frequency domains, have been proposed to characterize the human voice source. This Section gives a brief overview of the most commonly used parameters in the literature, since some of them are used in Sections 4 and 5.

3.1. Time-domain features

Several time-domain features can be expressed as a function of time intervals derived from the glottal waveform (Alku, 1992). These are used to characterize the shape of the waveform, by capturing for example the location of the primary or secondary opening instant (Laukkanen et al., 1996), of the glottal flow maximum, etc. The formulation of the source signal in the commonly used LF model (Fant et al., 1985) is based on time-domain parameters, such as the Open Quotient O_q , the Asymmetry coefficient α_m , or the Voice Speed Quotient S_q (Doval and d'Alessandro, 2006). However in most cases these instants are difficult to locate with precision from the glottal flow estimation. Avoiding this problem and preferred to the traditional Open

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Quotient, the Quasi-Open Quotient (QOQ) was proposed as a parameter describing the relative open time of the glottis (Hacki, 1989). It is defined as the ratio between the quasi-open time and the quasi-closed time of the glottis, and corresponds to the timespan (normalized to the pitch period) during which the glottal flow is above 50% of the difference between the maximum and minimum flow. Note that QOQ was used in (Laukkanen et al., 1996) for studying the physical variations of the glottal source related to the vocal expression of stress and emotion. In (Airas and Alku, 2007) various variants of Oq have been tested in terms of the degree by which they reflect phonation changes. QOQ was found to be the best for this task.

Another set of parameters is extracted from the amplitude of peaks in the glottal pulse or its derivative (Gobl and Chasaide, 2003). The Normalized Amplitude Quotient (NAQ) proposed by Alku in (Alku et al., 2002) turns out to be an essential glottal feature. NAQ is a parameter characterizing the glottal closing phase (Alku et al., 2002). It is defined as the ratio between the maximum of the glottal flow and the minimum of its derivative, normalized with respect to the fundamental period. Its robustness and efficiency to separate different types of phonation was shown in (Alku et al., 2002), (Airas and Alku, 2007). Note that a quasi-similar feature, called *basic shape parameter*, was proposed by Fant in (Fant, 1995), where it was qualified as “*most effective single measure for describing voice qualities*”.

In (Plumpe et al., 1999), authors propose to use 7 LF parameters and 5 energy coefficients (defined in 5 subsegments of the glottal cycle) respectively for characterizing the coarse and fine structures of the glottal flow estimation. Finally some approaches aim at fitting a model on the glottal flow estimate by computing a distance in the time domain (Plumpe et al., 1999), (Drugman et al., 2008).

3.2. Frequency-domain features

In the frequency domain, the LF model presents a low-frequency resonance called the *glottal formant* (Doval and d’Alessandro, 2006) (see the amplitude spectrum of the glottal flow derivative in Figure 2, row 1, column 2). Some approaches characterize the glottal formant both in terms of frequency and bandwidth (Drugman et al., 2009a). By defining a spectral error measure, other studies try to match a model to the glottal flow estimation (Ling et al., 2005), (Fant, 1995), (Drugman et al., 2008). This is also the case for the Parabolic Spectrum Parameter (PSP) proposed in (Alku et al., 1997).

1
2
3
4
5
6
7
8
9 An extensively used measure is the $H1-H2$ parameter (Fant, 1995). This
10 parameter is defined as the ratio between the amplitudes of the magnitude
11 spectrum of the glottal source at the fundamental frequency and at the second
12 harmonic (Klatt and Klatt, 1990), (Titze and Sundberg, 1992). It has been
13 widely used as a measure characterizing voice quality (Hanson, 1995), (Fant,
14 1995), (Alku et al., 2009).
15
16

17 For quantifying the amount of harmonics in the glottal source, the Har-
18 monic to Noise Ratio (HNR) and the Harmonic Richness Factor (HRF) have
19 been proposed in (Murphy and Akande, 2005) and (Childers and Lee, 1991).
20 More precisely, HRF quantifies the amount of harmonics in the magnitude
21 spectrum of the glottal source. It is defined as the ratio between the sum
22 of the amplitudes of harmonics, and the amplitude at the fundamental fre-
23 quency (Childers, 1999). It was shown to be informative about the phonation
24 type in (Childers and Lee, 1991) and (Alku et al., 2009).
25
26
27

28 **4. Experiments on Synthetic Speech**

29

30 The first experimental protocol we opted for is close to the one presented
31 in (Sturmel et al., 2007). Decomposition is achieved on synthetic speech
32 signals (sampled at 16 kHz) for various test conditions. The idea is to cover
33 the diversity of configurations one can find in continuous speech by varying all
34 parameters over their whole range. Synthetic speech is produced according
35 to the source-filter model by passing a known sequence of Liljencrants-Fant
36 (LF) glottal waveforms (Fant et al., 1985) through an auto-regressive filter
37 extracted by LPC analysis (with an order of 18) from real sustained vowels
38 uttered by a female speaker. As the mean pitch during these utterances
39 was about 180 Hz, it can be considered that fundamental frequency should
40 not exceed 100 and 240 Hz in continuous speech. For the LF parameters,
41 the Open Quotient Oq and Asymmetry coefficient α_m are varied through
42 their common range (see Table 1). For the filter, 14 types of typical vowels
43 are considered. Noisy conditions are modeled by adding a white Gaussian
44 noise to the speech signal, from almost clean conditions ($SNR = 80dB$) to
45 strongly adverse environments ($SNR = 10dB$). Table 1 summarizes all test
46 conditions, which makes a total of slightly more than 250,000 experiments. It
47 is worth mentioning that the synthetic tests presented in this section focus
48 on the study of non-pathological voices with a regular phonation. Although
49 the glottal analysis of less regular voices (e.g presenting a jitter or a shimmer;
50 or containing an additive noise component during the glottal production, as
51
52
53
54
55
56
57
58

it is the case for a breathy voice) is a challenging issue, this latter problem is not addressed in the present study.

Source			Filter	Noise
Pitch (Hz)	Oq	α_m	Vowel type	SNR (dB)
100:5:240	0.3:0.05:0.9	0.55:0.05:0.8	14 vowels	10:10:80

Table 1: Table of synthesis parameter variation range.

The three source estimation techniques described in Section 2 (CPIF, IAIF and CCD) are compared. In order to assess their decomposition quality, two objective quantitative measures are used (and the effect of noise, fundamental frequency and vocal tract variations to these measures are studied in detail in the next subsections):

- **Error rate on NAQ and QOQ** : An error on the estimation of NAQ and QOQ after source-tract decomposition should be penalized. An example of distribution for the relative error on QOQ in clean conditions is displayed in Figure 3. Many attributes characterizing such a histogram can be proposed to evaluate the performance of an algorithm. The one we used in our experiments is defined as the proportion of frames for which the relative error is higher than a given threshold of $\pm 20\%$. The lower the error rate on the estimation of a given glottal parameter, the better the glottal flow estimation method.

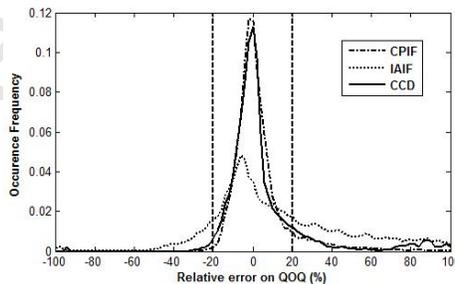


Figure 3: Distribution of the relative error on *QOQ* for the three methods in clean conditions ($SNR = 80dB$). The *error rate* is defined as the percentage of frames for which the relative error is higher than a given threshold of 20% (indicated on the plot).

- **Spectral distortion** : Many frequency-domain measures for quantifying the distance between two speech frames x and y arise from the

speech coding literature. Ideally the subjective ear sensitivity should be formalised by incorporating psychoacoustic effects such as masking or isosonic curves. A simple and relevant measure is the spectral distortion (SD) defined as (Nordin and Eriksson, 2001):

$$SD(x, y) = \sqrt{\int_{-\pi}^{\pi} (20 \log_{10} \left| \frac{X(\omega)}{Y(\omega)} \right|)^2 \frac{d\omega}{2\pi}} \quad (8)$$

where $X(\omega)$ and $Y(\omega)$ denote both signals spectra as a function of normalized angular frequency. In (K. Paliwal, 1993), authors argue that a difference of about 1dB (with a sampling rate of 8kHz) is hardly perceptible. In order to take this point into account, we used the following measure between the spectra of the estimated and reference glottal signals:

$$SD(Estimated, Reference) \approx \sqrt{\frac{2}{8000} \int_{20}^{4000} (20 \log_{10} \left| \frac{S_{Estimated}(f)}{S_{Reference}(f)} \right|)^2 df} \quad (9)$$

An efficient technique of glottal flow estimation is then reflected by low spectral distortion values.

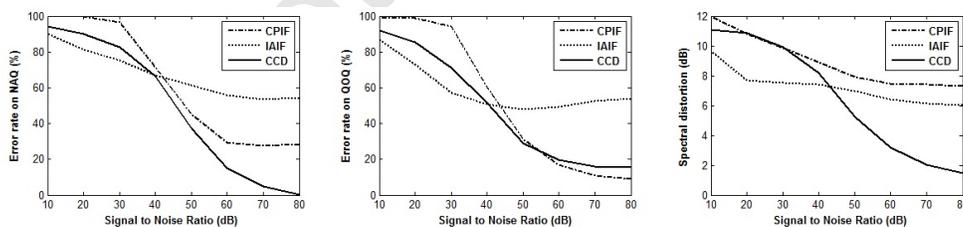


Figure 4: Evolution of the three performance measures (error rate on *NAQ* and *QOQ*, and spectral distortion) as a function of the Signal to Noise Ratio for the three glottal source estimation methods.

Based on this experimental framework, we now study how the glottal source estimation techniques behave in noisy conditions, or with regard to some factors affecting the decomposition quality, such as the fundamental frequency or the vocal tract transfert function.

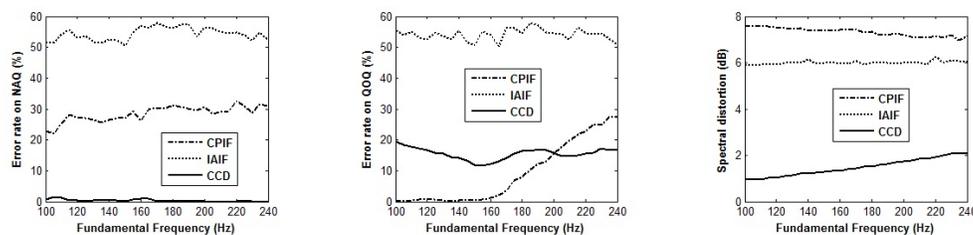


Figure 5: Evolution of the three performance measures as a function of the fundamental frequency for the three glottal source estimation methods.

4.1. Robustness to Additive Noise

As mentioned above, white Gaussian noise has been added to the speech signal, with various SNR levels. This noise is used as a (weak) substitute for recording or production noise but also for every little deviation to the theoretical framework which distinguishes real and synthetic speech. Results according to our three performance measures are displayed in Figure 4. As expected, all techniques degrade as the noise power increases. More precisely, CCD turns out to be particularly sensitive. This can be explained by the fact that a weak presence of noise may dramatically affect the phase information, and consequently the decomposition quality. The performance of CPIF is also observed to strongly degrade as the noise level increases. This is probably due to the fact that noise may dramatically modify the spectral envelope estimated during the closed phase, and the resulting estimate of the vocal tract contribution becomes erroneous. On the contrary, even though IAIF is, in average, the less efficient on clean synthetic speech, it outperforms other techniques in adverse conditions (below 40 dB of SNR). One possible explanation of its robustness is the iterative process it relies on. It can be indeed expected that, although the first iteration may be highly affected by noise (as it is the case for CPIF), the severity of the perturbation becomes weaker as the iterative procedure converges.

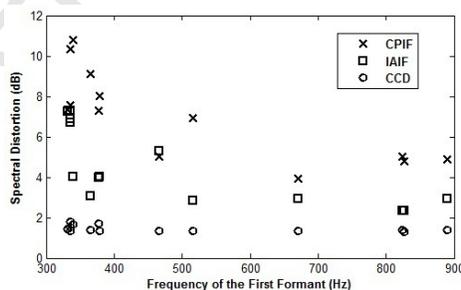
4.2. Sensitivity to Fundamental Frequency

Female voices are known to be especially difficult to analyze and synthesize. The main reason for this is their high fundamental frequency which implies to process shorter glottal cycles. As a matter of fact the vocal tract response has not the time to freely return to its initial state between two glottal sollicitation periods (i.e. the duration of the vocal tract response can be much longer than that of the glottal closed phase). Figure 5 shows the

1
2
3
4
5
6
7
8
9 evolution of our three performance measures with respect to the fundamen-
10 tal frequency in clean conditions. Interestingly, all methods maintain almost
11 the same efficiency for high-pitched voices. Nonetheless an increase of the
12 error rate on QOQ for CPIF, and an increase of the spectral distortion for
13 CCD can be noticed. It can be also observed that, for clean synthetic speech,
14 CCD gives the best results with an excellent determination of NAQ and a
15 very low spectral distortion. Secondly, despite its high spectral errors, CPIF
16 leads to an efficient parametrization of the glottal shape (with notably the
17 best results for the determination of QOQ).

21 4.3. Sensitivity to Vocal Tract

22
23 In our experiments, filter coefficients were extracted by LPC analysis on
24 sustained vowels. Even though the whole vocal tract spectrum may affect the
25 decomposition, the first formant, which corresponds to the dominant poles,
26 generally imposes the longest contribution of its time response. To give an
27 idea of its impact, Figure 6 exhibits, for the 14 vowels, the evolution of the
28 spectral distortion as a function of the first formant frequency F_1 . A general
29 trend can be noticed from this graph: it is observed for all methods that the
30 performance of the glottal flow estimation degrades as F_1 decreases. This will
31 be explained in the next Section by an increasing overlap between source and
32 filter components, as the vocal tract impulse response gets longer. It is also
33 noticed that this degradation is particularly important for both CPIF and
34 IAIF methods, while the quality of CCD (which does not rely on a parametric
35 modeling) is only slightly altered.



51 Figure 6: Evolution, for the 14 vowels, of the spectral distortion with the first formant
52 frequency F_1 .

4.4. Conclusions on Synthetic Speech

Many factors may affect the quality of the source-tract separation. Intuitively, one can think about the *time interference* between minimum and maximum-phase contributions, respectively related to the vocal tract and to the glottal open phase. The stronger this interference, the more important the time overlap between the minimum-phase component and the maximum-phase response of the next glottal cycle, and consequently the more difficult the decomposition. Basically, this interference is conditioned by three main parameters:

- the pitch F_0 , which imposes the spacing between two successive vocal system responses,
- the first formant F_1 , which influences the length of the minimum-phase contribution of speech,
- and the glottal formant F_g , which controls the length of the maximum-phase contribution of speech. Indeed, the glottal formant is the most important spectral feature of the glottal open phase (see the low-frequency resonance in the amplitude spectrum of the glottal flow derivative in Figure 2). It is worth noting that F_g is known (Doval and d’Alessandro, 2006) to be a function of the time-domain characteristics of the glottal open phase (i.e of the maximum-phase component of speech): the open quotient O_q , and the asymmetry coefficient (α_m).

A strong interference then appears with high pitch, and with low F_1 and F_g values. The previous experiments confirmed for all glottal source estimation techniques the performance degradation as a function of F_0 and F_1 . Although we did not explicitly measure the sensitivity of these techniques to F_g in this manuscript, it was confirmed in other informal experiments we performed.

It can be also observed from Figures 4 and 5 that the overall performance through an objective study on synthetic signals is the highest for the complex cepstrum-based technique. This method leads to the lowest values of spectral distortion and gives relatively high rates for the determination of both NAQ and QOQ parameters. The CPIF technique exhibits better performance in the determination of QOQ in clean conditions and especially for low-pitched speech. As for the IAIF technique, it turns out that it gives the worst results in clean synthetic speech but outperforms other approaches in adverse

1
2
3
4
5
6
7
8
9
10
11
12
13
14
noisy conditions. Note that our results corroborate the conclusions drawn in
(Sturmel et al., 2007) where the mixed-phase deconvolution (achieved in that
study by the ZZT method) was shown to outperform other state-of-the-art
approaches of glottal flow estimation.

15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65

5. Experiments on Real Speech

Reviewing the glottal flow estimation literature, one can easily notice that testing with natural speech is a real challenge. Even in very recent published works, all tests are performed only on sustained vowels. In addition, due to the unavailability of a reference for the real glottal flow (see Section 1), the procedure of evaluation is generally limited to providing plots of glottal flow estimates, and checking visually if they are consistent with expected glottal flow models. For real speech experiments, here we will first follow this state-of-the-art experimentation (of presenting plots of estimates for a real speech example), and then extend it considerably both by extending the content of the data to a large connected speech database (including non-vowels), and extending the method to a comparative parametric analysis approach.

In this study, experiments on real speech are carried out on the De7 corpus, a diphone database designed for expressive speech synthesis (Schroeder and Grice, 2003). The database contains three voice qualities (modal, soft and loud) uttered by a German female speaker, with about 50 minutes of speech available for each voice quality (leading to a total of around 2h30). Recordings sampled at 16 kHz are considered. Locations of both GCIs and GOIs are precisely determined from these signals using the algorithm described in (Drugman and Dutoit, 2009). As mentioned in Section 2, an accurate position of both events is required for an efficient CPIF technique, while the mixed-phase decomposition (as achieved by CCD) requires, among others, GCI-centered windows to exhibit correct phase properties.

Let us first consider in Figure 7 a concrete example of glottal source estimation on a given voiced segment ($/aI/$ as in "ice") for the three techniques and for the three voice qualities. In the IAIF estimate, some ripples are observed as if some part of the vocal tract filter contribution could not be removed. On the other hand, it can be noticed that the estimations from CPIF and CCD are highly similar and are very close to the shape expected by the glottal flow models, such as the LF model (Fant et al., 1985). It can be also observed that the abruptness of the glottal open phase around the

GCI is stronger for the loud voice, while the excitation for the softer voice is smoother.

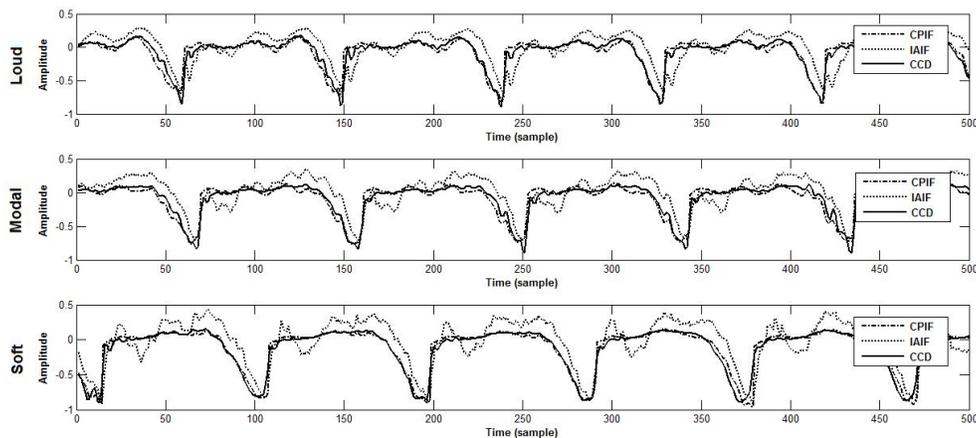


Figure 7: Example of glottal flow derivative estimation on a given segment of vowel (/aI/ as in "ice") for the three techniques and for the three voice qualities: (*top*) loud voice, (*middle*) modal voice, (*bottom*) soft voice.

We investigated whether the glottal source estimated by these techniques conveys information about voice quality. Indeed the glottis is assumed to play an important part for the production of such expressive speech (d'Alessandro, 2006). As a matter of fact we found some differences between the glottal features in our experiments on the De7 database. In this experiment, the NAQ, H1-H2 and HRF parameters described in Section 3 are used. Figure 8 illustrates the distributions of these features estimated by CPIF, IAIF and CCD for the three voice qualities. This Figure can be considered as a summary of the voice quality analysis using three state-of-the-art methods on a large speech database. The parameters NAQ, H1-H2 and HRF have been used frequently in the literature to label phonation types (Alku et al., 2002), (Hanson, 1995), (Childers and Lee, 1991). Hence the separability of the phonation types based on these parameters can be considered as a measure of effectiveness for a particular glottal flow estimation method.

For the three methods, significant differences between the histograms of the different phonation types can be noted. This supports the claim that, by applying one of the given glottal flow estimation methods and by parametrizing the estimate with one or more of the given parameters, one can perform automatic voice quality/phonation type labeling with a much higher success

rate than by random labeling. It is noticed from Figure 8 that parameter distributions are convincingly distinct, except for the IAIF and H1-H2 combination. The sorting of the distributions with respect to vocal effort are consistent and in line with results of other works ((Alku et al., 2002) and (Alku et al., 2009)). Among other things, strong similarities between histograms obtained by CPIF and CCD can be observed. In all cases, it turns out that the stronger the vocal effort, the lower NAQ and H1-H2, and the higher HRF.

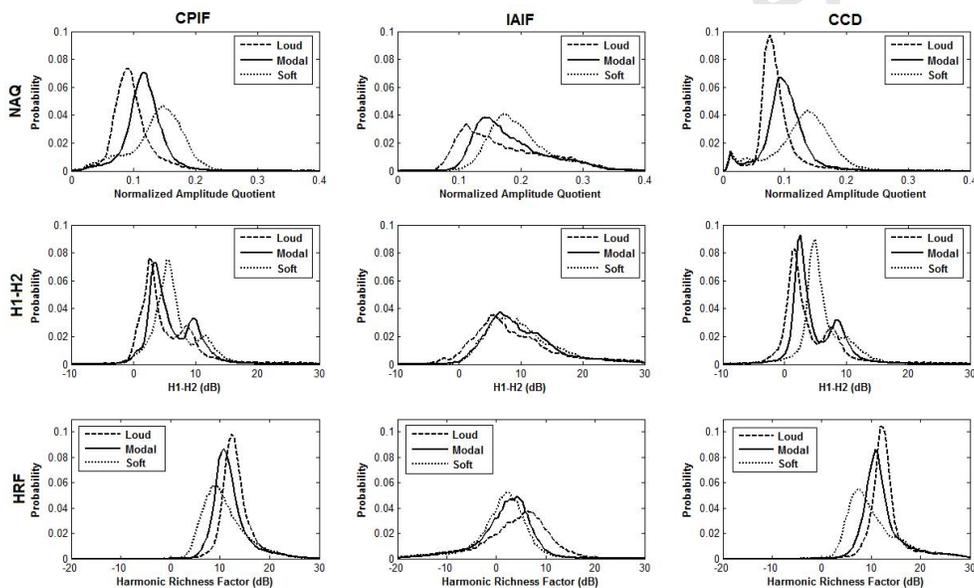


Figure 8: Distributions, for various voice qualities, of three glottal features (from top to bottom: NAQ, H1-H2 and HRF) estimated by three glottal source estimation techniques (from left to right: CPIF, IAIF and CCD). The voice qualities are shown as dashed (loud voice), solid (modal voice) and dotted (soft voice) lines.

Although some significant differences in glottal feature distributions have been visually observed, it is interesting to quantify the discrimination between the voice qualities enabled by these features. For this, the Kullback-Leibler (KL) divergence, known to measure the separability between two discrete density functions A and B , can be used (Lin, 1991):

$$D_{KL}(A, B) = \sum_i A(i) \log_2 \frac{A(i)}{B(i)} \quad (10)$$

Since this measure is non-symmetric (and consequently is not a true distance), its symmetrised version, called Jensen-Shannon divergence, is often preferred. It is defined as a sum of two KL measures (Lin, 1991):

$$D_{JS}(A, B) = \frac{1}{2}D_{KL}(A, M) + \frac{1}{2}D_{KL}(B, M) \quad (11)$$

where M is the average of the two distributions ($M = 0.5 * (A + B)$). Figure 9 displays the values of the Jensen-Shannon distances between two types of voice quality, for the three considered features estimated by the three techniques.

From this figure, it can be noted that NAQ is the best discriminative feature (i.e. has the highest Jensen-Shannon distance between distributions), while H1-H2 and HRF convey a comparable amount of information for discriminating voice quality. As expected, the loud-soft distribution distances are highest compared to loud-modal and modal-soft distances. In seven cases out of nine (three different parameters and three different phonation type couples), CCD leads to the most relevant separation and in two cases (loud-modal separation with NAQ, loud-modal separation with HRF) CPIF provides a better separation. Both Figures 8 and 9 show that the effectiveness of CCD and CPIF is similar, with slightly better results for CCD, while IAIF exhibits clearly lower performance (except for one case: loud-modal separation with HRF).

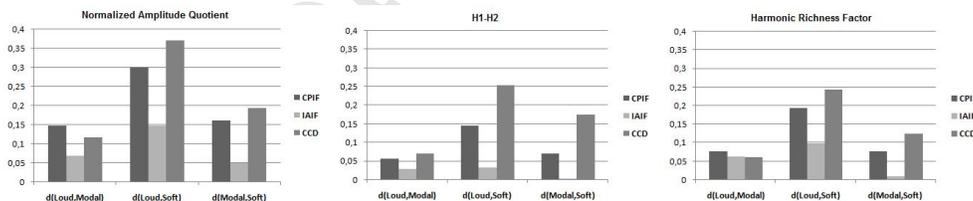


Figure 9: Jensen-Shannon distances between two types of voice quality using (from left to right) the NAQ, H1-H2 and HRF parameters. For each feature and pair of phonation types, the three techniques of glottal source estimation are compared.

6. Conclusion

This study aimed at comparing the effectiveness of the main state-of-the-art glottal flow estimation techniques. For this, detailed tests on both synthetic and real speech were performed. For real speech, a large corpus was

1
2
3
4
5
6
7
8
9 used for testing, without limiting analysis to sustained vowels. Due to the
10 unavailability of the reference glottal flow signals for real speech examples,
11 the separability of three voice qualities was considered as a measure of the
12 ability of the methods to discriminate different phonation types. In synthetic
13 speech tests, objective measures were used since the original glottal flow
14 signals were available. Our first conclusion is that the usefulness of NAQ,
15 H1-H2 and HRF for parameterizing the glottal flow is confirmed. We also
16 confirmed other works in the literature (such as (Alku et al., 2002) and
17 (Alku et al., 2009)) showing that these parameters can be effectively used as
18 measures for discriminating different voice qualities. Our results show that
19 the effectiveness of CPIF and CCD appears to be similar and rather high,
20 with a slight preference towards CCD. However, it should be emphasized here
21 that in our real speech tests, clean signals recorded for text-to-speech (TTS)
22 synthesis were used. We can thus confirm the effectiveness of CCD for TTS
23 applications (such as emotional/expressive TTS). However, for applications
24 which require the analysis of noisy signals (such as telephone applications)
25 further testing is needed. We observed that in the synthetic speech tests, the
26 ranking dramatically changed depending on the SNR and the robustness of
27 CCD was observed to be rather low. IAIF has lower performance in most tests
28 (both in synthetic and real speech tests) but shows up to be comparatively
29 more effective in very low SNR values.
30
31
32
33
34
35
36

37 **Acknowledgment**

38
39 Thomas Drugman is supported by the Belgian Fonds National de la
40 Recherche Scientifique (FNRS). Authors also would like to thank the re-
41 viewers for their fruitful comments.
42
43
44

45 **References**

- 46
47 Airas, M., 2008. Tkk aparat: An environment for voice inverse filtering and
48 parameterization. *Logopedics Phoniatrics Vocology* 33, 49–64.
49
50 Airas, M., Alku, P., 2006. Emotions in vowel segments of continuous speech
51 : Analysis of the glottal flow using the normalised amplitude quotient.
52 *Phonetica* 63, 26–46.
53
54 Airas, M., Alku, P., 2007. Comparison of multiple voice source parameters
55 in different phonation types, in: *Proc. Interspeech*, pp. 1410–1413.
56
57
58

- 1
2
3
4
5
6
7
8
9 Alku, P., 1992. An automatic method to estimate the time-based parameters
10 of the glottal pulseform, in: Proc. ICASSP, pp. 29–32.
11
- 12 Alku, P., Backstrom, T., Vilkmán, E., 2002. Normalized amplitude quotient
13 for parametrization of the glottal flow. *Journal of the Acoustical Society*
14 *of America* 112, 701–710.
15
16
- 17 Alku, P., Magi, C., Yrttiaho, S., Backstrom, T., Story, B., 2009. Closed
18 phase covariance analysis based on constrained linear prediction for glottal
19 inverse filtering. *Journal of the Acoustical Society of America* 125, 3289–
20 3305.
21
22
- 23 Alku, P., Strik, H., Vilkmán, E., 1997. Parabolic spectral parameter - a new
24 method for quantification of the glottal flow. *Speech Communication* 22,
25 67–79.
26
27
- 28 Alku, P., Svec, J., Vilkmán, E., Sram, F., 1992. Glottal wave analysis with
29 pitch synchronous iterative adaptive inverse filtering. *Speech Communica-*
30 *tion* 11, 109–118.
31
32
- 33 Alku, P., Vilkmán, E., 1994. Estimation of the glottal pulseform based on
34 discrete all-pole modeling, in: *Third International Conference on Spoken*
35 *Language Processing*, pp. 1619–1622.
36
- 37 Bozkurt, B., Doval, B., d’Alessandro, C., Dutoit, T., 2005. Zeros of z -
38 transform representation with application to source-filter separation in
39 speech. *IEEE Signal Processing Letters* 12.
40
41
- 42 Bozkurt, B., Dutoit, T., 2003. Mixed-phase speech modeling and formant es-
43 timation, using differential phase spectrums, in: *ISCA ITRW VOQUAL03*,
44 pp. 21–24.
45
46
- 47 Brookes, D., Chan, D., 1994. Speaker characteristics from a glottal airflow
48 model using glottal inverse filtering. *Institute of Acoust.* 15, 501–508.
49
- 50 Cabral, J., Renals, S., Richmond, K., Yamagishi, J., 2008. Glottal spec-
51 tral separation for parametric speech synthesis, in: *Proc. Interspeech*, pp.
52 1829–1832.
53
54
- 55 Childers, D., 1999. *Speech Processing and Synthesis Toolboxes*. Wiley and
56 Sons, Inc.
57
58

- 1
2
3
4
5
6
7
8
9 Childers, D., Lee, C., 1991. Vocal quality factors: Analysis, synthesis, and
10 perception. *Journal of the Acoustical Society of America* 90, 2394–2410.
11
12 d’Alessandro, C., 2006. Voice source parameters and prosodic analysis, in:
13 *Method in empirical prosody research*, pp. 63–87.
14
15
16 Deng, H., Ward, R., Beddoes, M., Hodgson, M., 2006. A new method for
17 obtaining accurate estimates of vocal-tract filters and glottal waves from
18 vowel sounds. *IEEE Trans. on Acoustics, Speech, and Signal Processing*
19 14, 445–455.
20
21
22 Doval, B., d’Alessandro, C., 2006. The spectrum of glottal flow models. *Acta*
23 *acustica united with acustica* 92, 1026–1046.
24
25
26 Doval, B., d’Alessandro, C., Henrich, N., 2003. The voice source as a
27 causal/anticausal linear filter, in: *ISCA ITRW VOQUAL03*, pp. 15–19.
28
29 Drugman, T., Bozkurt, B., Dutoit, T., 2009a. Complex cepstrum-based de-
30 composition of speech for glottal source estimation, in: *Proc. Interspeech*.
31
32
33 Drugman, T., Dubuisson, T., d’Alessandro, N., Moinet, A., Dutoit, T., 2008.
34 Voice source parameters estimation by fitting the glottal formant and the
35 inverse filtering open phase, in: *16th European Signal Processing Confer-*
36 *ence*.
37
38
39 Drugman, T., Dubuisson, T., Dutoit, T., 2009b. On the mutual information
40 between source and filter contributions for voice pathology detection, in:
41 *Proc. Interspeech*.
42
43
44 Drugman, T., Dutoit, T., 2009. Glottal closure and opening instant detection
45 from speech signals, in: *Proc. Interspeech*.
46
47
48 F. Itakura, S.S., 1970. A statistical method for estimation of speech spectral
49 density and formant frequencies. *Electron. Commun. Japan* 53-A, 36–43.
50
51
52 Fant, G., 1995. The lf-model revisited. transformations and frequency domain
53 analysis. *STL-QPSR* 36, 119–156.
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9 Fu, Q., Murphy, P., 2006. Robust glottal source estimation based on joint
10 source-filter model optimization. *IEEE Trans. on Audio, Speech, and Lan-*
11 *guage Processing* 14, 492–501.
12
13
14 Gardner, W., Rao, B., 1997. Noncausal all-pole modeling of voiced speech.
15 *IEEE Trans. Speech and Audio Processing* 5, 1–10.
16
17 Gobl, C., Chasaide, A., 2003. Amplitude-based source parameters for mea-
18 suring voice quality, in: *VOQUAL03*, pp. 151–156.
19
20 Hacki, T., 1989. Klassifizierung von glottisdysfunktionen mit hilfe der elek-
21 troglottographie. *Folia Phoniatica* 41, 43–48.
22
23
24 Hanson, H., 1995. Individual variations in glottal characteristics of female
25 speakers, in: *Proc. ICASSP*, pp. 772–775.
26
27 Henrich, N., d’Alessandro, C., Doval, B., Castellengo, M., 2004. On the
28 use of the derivative of electroglottographic signals for characterization of
29 non-pathological phonation. *J. Acoust. Soc. Am.* 115, 1321–1332.
30
31
32 Jaroudi, A.E., Makhoul, J., 1991. Discrete all-pole modeling. *IEEE Trans.*
33 *on Signal Processing* 39, 411–423.
34
35
36 K. Paliwal, B.A., 1993. Efficient vector quantization of lpc parameters at 24
37 bits/frame. *IEEE Trans. Speech Audio Processing* 1, 3–14.
38
39 Klatt, D., Klatt, L., 1990. Analysis, synthesis and perception of voice qual-
40 ity variations among female and male talkers. *Journal of the Acoustical*
41 *Society of America* 87, 820–857.
42
43
44 Laukkanen, A.M., Vilkman, E., Alku, P., Oksanen, H., 1996. Physical vari-
45 ations related to stress and emotional state: a preliminary study. *Journal*
46 *of Phonetics* 24, 313–335.
47
48
49 Lin, J., 1991. Divergence measures based on the shannon entropy. *IEEE*
50 *Trans. on Information Theory* 37, 145–151.
51
52
53 Ling, Z., Hu, Y., Wang, R., 2005. A novel source analysis method by match-
54 ing spectral characters of lf model with straight spectrum. *Lecture Notes*
55 *in Computer Science* 3784, 441–448.
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9 Moore, E., Clements, M., 2004. Algorithm for automatic glottal waveform
10 estimation without the reliance on precise glottal closure information, in:
11 Proc. ICASSP.
12
13
14 Murphy, P., Akande, O., 2005. Quantification of glottal and voiced speech
15 harmonics-to-noise ratios using cepstral-based estimation, in: Nonlinear
16 Speech Processing Workshop, pp. 224–232.
17
18
19 Nordin, F., Eriksson, T., 2001. A speech spectrum distortion measure with
20 interframe memory, in: Proc. ICASSP, pp. 717–720.
21
22 [Online], 2008. http://aparat.sourceforge.net/index.php/main_page. TKK
23 Aparat Main Page .
24
25
26 Oppenheim, A., Schafer, R., 1989. Discrete-time signal processing. Prentice-
27 Hall.
28
29 Pedersen, C., Andersen, O., Dalsgaard, P., 2010. Separation of mixed-phase
30 signals by zeros of the z-transform - a reformulation of complex cepstrum-
31 based separation by causality, in: Proc. ICASSP.
32
33
34 Plumpe, M., Quatieri, T., Reynolds, D., 1999. Modeling of the glottal
35 flow derivative waveform with application to speaker identification. IEEE
36 Trans. on Speech and Audio Processing 7, 569–586.
37
38
39 Quatieri, T., 2002. Discrete-time speech signal processing. Prentice-Hall.
40
41 Schroeder, M., Grice, M., 2003. Expressing vocal effort in concatenative syn-
42 thesis, in: 15th International Conference of Phonetic Sciences, pp. 2589–
43 2592.
44
45
46 Sturmel, N., d’Alessandro, C., Doval, B., 2007. A comparative evaluation
47 of the zeros of z transform representation for voice source estimation, in:
48 Proc. Interspeech, pp. 558–561.
49
50
51 Titze, I., Sundberg, J., 1992. Vocal intensity in speakers and singers. Journal
52 of the Acoustical Society of America 91, 2936–2946.
53
54
55 Veeneman, D., Bement, S., 1985. Automatic glottal inverse filtering from
56 speech and electroglottographic signals. IEEE Trans. on Acoustics, Speech,
57 and Signal Processing 33, 369–377.
58
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9 Vincent, D., Rosec, O., Chovanel, T., 2005. Estimation of lf glottal source
10 parameters based on an arx model, in: Proc. Interspeech, pp. 333–336.
11
12 Walker, J., Murphy, P., 2007. A review of glottal waveform analysis, in:
13 Progress in Nonlinear Speech Processing, pp. 1–21.
14
15 Wong, D., Markel, J., Gray, A., 1979. Least squares glottal inverse filtering
16 from the acoustic speech waveform. IEEE Trans. on Acoustics, Speech,
17 and Signal Processing 27.
18
19
20 Yegnanarayana, B., Veldhuis, R., 1998. Extraction of vocal-tract system
21 characteristics from speech signals. IEEE Trans. Speech Audio Processing
22 6, 313–327.
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65