# Exhibit H

This procedure, which computes $P(\chi \mid \hat{\theta}^{(j)})$, is the one called by the extremization subroutine for each step in the extremization task. The process ends once the proposed precision for $\hat{r}_0^{(j)}$, $\hat{r}_1^{(j)}$, $\cdots$, $\hat{r}_p^{(j)}$ is obtained and the values for $\hat{a}_1^{(j)}$, $\hat{a}_2^{(j)}$, $\cdots$, $\hat{a}_p^{(j)}$ and $\hat{\sigma}^{2(j)}$ are calculated from (4) and (5).

It is well known that iterative extremization procedures (like quasi-Newton, etc.) perform much better if they are given good initial estimates. These initial estimates may be obtained from [10]

$$\hat{r}_l^{(0)} = \frac{1}{N - l + 1} \sum_{k=l}^{N} x_k x_{k-l} \qquad p \geq l \geq 0. \qquad (7)$$

The Toeplitz matrix $\hat{\Gamma}^{(j)}$, with values of $\hat{r}_l^{(j)}$ obtained from (7), may not be positive definite, and hence, the computation of $\hat{\sigma}^{2(j)}$ from (5) may give a negative value for $\hat{\sigma}^{2(j)}$ for $|\hat{\Gamma}^{(j)}|$; but if that is the case, we may force a positive value $\hat{\sigma}^{2(j)}$ by giving an increased value to $\hat{r}_0^{(j)}$ (changing it to $\hat{r}_0'^{(j)}$) so that the use of $\hat{\Gamma}^{(j)}$ in (5) gives a positive value to $\hat{\sigma}^{2(j)}$ and $|\hat{\Gamma}^{(j)}|$. Once the iteration is close to the true ML value, this will not be a problem anymore.

### III. Conclusion

In many cases, it is necessary to have a computationally efficient algorithm for obtaining an ML estimate of process parameters, and that may compensate for the fact that what is actually obtained is only an approximation to the true ML estimates. In other cases, it is desirable that the estimate be as near as possible to the true optimal value even at the cost of a less efficient algorithm. In this correspondence we show that the true ML estimates of Gaussian zero-mean AR process parameters can be calculated and give the steps toward their obtainment using standard subroutines for calculating extrema. The function that is maximized is the exact likelihood of the data in contrast to current approximate methods. If a fast approximation is needed, methods such as the one in [4] for good approximate ML can be used to begin with, and then this method can improve the estimate as much as necessary.

We would like to thank an anonymous reviewer for his comments and suggestions.

### References

[1] G. Box and G. Jenkins, *Time Series Analysis—Forecasting and Control.* San Francisco, CA: Holden-Day, 1970.
[2] J. Gentler and C. Banyasz, "A recursive (on-line) maximum likelihood identification method," *IEEE Trans. Automat. Contr.*, vol. AC-19, pp. 816–820, Dec. 1974.
[3] K. Astrom, "Maximum likelihood and prediction error methods," *Automatica*, vol. 17, pp. 551–574, Sept. 1980.
[4] S. Kay, "Recursive maximum likelihood estimation of autoregressive processes," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-31, pp. 56–65, Feb. 1983.
[5] G. Goodwin and R. Payne, *Dynamic System Identification and Experiment Design Analysis.* New York: Academic, 1977.
[6] B. Friedlander, "A recursive maximum likelihood algorithm for ARMA spectral estimation," *IEEE Trans. Inform. Theory*, vol. IT-28, July 1982.
[7] R. Fletcher, "Fortran subroutines for minimization by quasi-Newton methods," Rep. R7125 AERE, Harwell, England, June 1972.
[8] S. Kirkpatrick, C. Gellatt, Jr., and M. Vecchi, "Optimization by simulated annealing," IBM T. J. Watson Res. Center, Yorktown Heights, NY, 1982.
[9] J. Burg, "A new analysis technique for time series data," NATO Advanced Study Inst. on Signal Processing with Emphasis on Underwater Acoust., Aug. 12–23, 1968; also in D. Childers, Ed. *Modern Spectrum Analysis.* New York: IEEE Press, 1978, pp. 42–48.
[10] T. Ulrych and T. Bishop, "Maximum entropy spectral analysis and autoregressive decomposition," *Rev. Geophys. and Space Phys.*, Feb. 1975; also in D. Childers, Ed. *Modern Spectrum Analysis.* New York: IEEE Press, 1978.

# Adaptive Silence Deletion for Speech Storage and Voice Mail Applications

CHEONG K. GAN AND ROBERT W. DONALDSON

*Abstract*—An algorithm which uses two adaptive amplitude thresholds and zero-crossing rate was used to delete nonspeech material from speech digitally encoded and then decoded using PCM, adaptive differential PCM, and adaptive delta modulation. Typically, compression rates of 35 percent, resulted. Subject evaluations were used to assess reconstructed speech quality, which improved significantly when absolute silence on playback was replaced with prerecorded background noise.

### I. Introduction

A relatively simple algorithm has been developed to delete from speech waveforms the nonessential acoustic material loosely referred to as "silence." Classification of each 10 ms acoustic segment is based on two adaptive amplitude thresholds, zero-crossing rate, and a minimum speech-segment-duration requirement. Applications of our algorithm include storage/playback or packet transmission, where the objective is to delete as much acoustic material as possible subject to adequate reconstructed speech quality.

Other work and applications related to but different from ours include: TASI/DSI (time-assignment speech interpolation/digital speech interpolation) where percent silence deletion is limited by the need to avoid excessively high talkspurt frequencies [1], [2]; speech recognition where precise delineation of speech/nonspeech boundaries is essential [3]; and synchronous transmission with buffering, where the instananeous speech sampling rate is dynamically adjusted in accordance with buffer space available, with the maximum buffer capacity chosen to balance reconstructed speech quality against excessive delay [4].

### II. Silence Deletion Algorithm

Speech data for our work were digitized at a sampling rate of 8 kHz using 12-bit digital-to-analog conversion following prefiltering by a 75–3000 Hz Butterworth filter. Short-time amplitude magnitude sum $A$ and zero-crossing rate $Z$ were obtained every 10 ms (every 80 samples) using a 100-sample window width:

$$A = \frac{1}{100} \sum_{i=1}^{100} |x_i| \qquad (1)$$

$$Z = \sum_{i=1}^{100} \left[ 1 - \left( x_i x_{i-1} / |x_i x_{i-1}| \right) \right]. \qquad (2)$$

Algorithm development was based on a 30 s speech sample, taken from the following cassette recording of a prepared lecture read by an American male:

> It is claimed that young children up to the age of about seven or eight years are incapable of grasping the abstract fundamental that number and volume remain constant even through changes in the outward appearance of the object. For thirty years the work of Piaget and his colleagues in Geneva has profoundly influenced the education of the young child.

The sample included many weak fricatives and other weak consonants (which are difficult to differentiate from background noise) and silent intervals ranging from 10 ms for intraword silence to 2.5 s for intersentence silence. The background noise level was relatively high.

Amplitude level proved to be very useful in silence/speech discrimination, particularly when the threshold to detect silence-to-speech transitions exceeded the threshold for speech-to-silence transitions and when these thresholds adapted to local background noise level variations. However, like others [2], [3], we found that amplitude alone will not always distinguish speech from silence; we used the relatively high $Z$ values of weak consonants to differentiate these from background noise.

Our silence deletion algorithm appears in Fig. 1. The important parameters are: 1) ZSIL, the zero-crossing rate boundary between speech and silence; 2) MINSP, the minimum number of contiguous 10 ms segments needed for any and each of these to be classified as speech; 3) AON, the amplitude threshold multiplier for detecting silence-to-speech transitions; and 4) AOFF, the amplitude threshold multiplier for identifying speech-to-silence transitions. To adapt to background noise level variations, the actual amplitude threshold is obtained by multiplying $T$ ($T$ = AON or AOFF) by AVG, the local average of the 10 most recent silence-period $A$ values. Because some segments near the silence-to-speech transition can have $A$ values up to twice the average and corrupt it, we allowed only segments for which $A <$ ACRIT*AVG to update AVG, with ACRIT = 1.28.

During initialization, AVG was selected such that AOFF*AVG = 50 on a 0–2048 scale. Following classification of the first silence segment, AVG and all 10 AVG FIFO stack values were set this silence segment's $A$ value. Following subsequent silence segments, AVG was updated as explained in the previous paragraph. This initialization and first updating of AVG avoided the necessity of requiring that the first few acoustic segments be silence.

Immediately after setting $T$ = AON, the first segment with either $A$ or $Z$ above threshold was classified as speech only if $A >$ AON*AVG or $Z >$ ZSIL for each of the subsequent MINSP segments. Short acoustic bursts otherwise classified as speech were thus deleted, without quality loss. Following classification of a segment as speech, $T$ = AOFF was initiated or maintained. Setting $T$ = AON immediately following a silence prior to a weak stop consonant can result in the stop's classification as silence, since AON > AOFF. Premature activation of AON is therefore avoided by requiring 6 contiguous silence segments (60 ms) before setting $T$ = AON.

We varied algorithm parameter values to obtain maximum silence deletion, subject to maintaining each of the following 12 phrases as speech: "ch" and "ren" in children, "th" (the), "ge" (age), "se" (seven), "s" (years, abstract, constant), "ble" (incapable), "ge" (Piaget), "ation" (education), and "ch" in child. Averaged over two different recording levels, we obtained 35 percent silence deletion with ZSIL = 36–38, AON = 2.60–3.02, AOFF = 1.80–2.57, and MINSP = 40 ms (4 segments).

Our algorithm deleted silence from speech encoded, and subsequently decoded using PCM, ADPCM (adaptive differential PCM), and CVSDM (continuously variable slope delta modulation). Unlikely codewords were used to mark the beginning, duration, and end of each silence interval, with less than 0.5 percent overhead. Our ADPCM codec used a fixed previous-sample predictor with gain 0.8. The uniform quantizer step size increased rapidly in response to slope-overload noise, but decayed more slowly in granular noise. Third-order polynomial interpolation yielded 16, 24, and 32 kHz sampling for CVSD modulation, chosen for its simplicity, universality and IC availability. In Fig. 2, $a$ = 0.95, and $z(n)$ = 0.7 $z(n-1)$ + 115.2 if $c(n)$ = $c(n-1)$ = $c(n-2)$; otherwise, $z(n)$ = 0.7 $z(n-1)$ + 0.6. These CVSDM parameters provide good waveform tracking at high signal levels, rather high SQNR (signal-to-quantizing noise ratio) at low signal levels, good syllabic adaptation, and moderate step size error recovery.
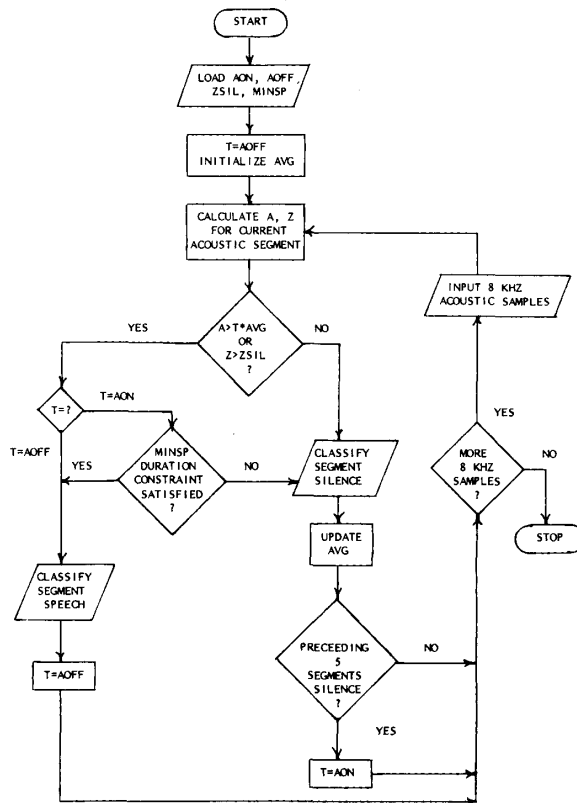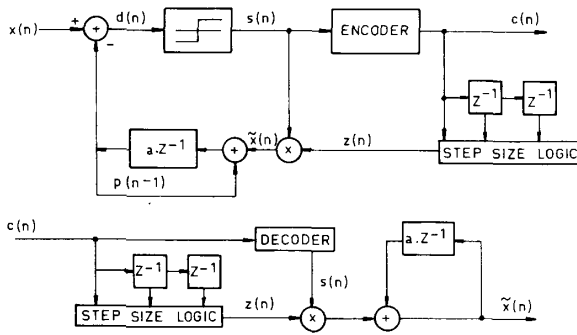


Fig. 1. Silence deletion algorithm.



Fig. 2. CVSD modulation system ($Z^{-1}$ denotes one-sample delay).

Percent silence deletion increased with decreasing data rates, since low-level speech was increasingly decoded as silence. Typical percent silence values were: 37, 36, 35, and 42 for 12-bit uniform PCM, 8-bit, A-law, PCM, 6-bit ADPCM, and 24 kbit/s CVSDM, respectively. Restriction of the maximum playback silence interval to 2.56 s, and some temporary adaptation loss following silence was virtually unnoticeable.

III. SUBJECTIVE EVALUATIONS

As test materials we used the lecture tape described earlier (source $A$), a radio newscast read by a Canadian male (source $B$), and a telephone conversation (recorded from an AM radio channel) between a Canadian male talkshow host and a male caller with a

TABLE I
SUBJECTIVE TEST FORMAT

| | |
|---|---|
| 1 — Original Sample A | 34 — Original Sample A |
| 2 — Original Sample B | 35 — Original Sample B |
| 3 — Original Sample C | 36 — Original Sample C |
| 4 ─┐ | 37 ─┐ |
| . ├─Random Ordering of | . ├─Random Ordering of |
| . │ Processed Samples | . │ Processed Samples |
| . │ | . ┘ |
| 22 — Duplicate of 37 | 55 — Duplicate of 4 |
| . ─┐ | . ─┐ |
| . ├─Random Ordering of | . ├─Random Ordering of |
| . ┘ Processed Samples | . ┘ Processed Samples |
| 33 — Duplicate of 4 | 67 — Duplicate of 37 |

TABLE II
RATINGS FOR SELECTED TEST SAMPLES

| Test Sample Description | Subject Number 1 2 3 4 5 6 7 8 9 10 | Subjective Ratings Mean | Subjective Ratings Variance |
|---|---|---|---|
| Original A ┌ 1 | 4 4 5 4 5 3 4 4 2  3 | 3.8 | 0.9 |
|          └ 34 | 4 4 5 4 3 4 4 4 2  4 | 3.8 | 0.7 |
| Original B ┌ 2 | 5 4 5 4 5 4 4 5 3  3 | 4.2 | 0.7 |
|          └ 35 | 5 3 5 4 4 5 4 5 3  4 | 4.2 | 0.7 |
| Original C ┌ 3 | 3 2 4 3 4 2 3 2 2  1 | 2.6 | 0.9 |
|          └ 36 | 5 4 4 3 2 3 3 4 2  2 | 3.2 | 1.0 |
| Duplicates ┌ 4 | 2 1 3 2 4 1 2 2 3  1 | 2.1 | 0.9 |
|          │ 33 | 3 1 3 3 2 1 2 3 1  2 | 2.1 | 0.8 |
|          └ 55 | 3 2 4 3 3 2 2 3 2  2 | 2.6 | 0.7 |
| Duplicates ┌ 22 | 3 1 5 4 3 4 3 3 1  2 | 2.9 | 1.2 |
|          │ 37 | 4 2 3 3 2 4 3 2 2  3 | 2.8 | 0.7 |
|          └ 67 | 3 2 4 2 2 1 3 3 2  3 | 2.5 | 0.8 |

strong European accent (source $C$). From each source, a 10-s 12-bit sample was obtained. Each of these original samples was coded and subsequently decoded using $A$-law PCM with $N = 2, 4$, or 6 quantization levels, ADPCM with $N = 3, 5$, or 7 levels, and CVSDM at 16, 24, or 32 kbit/s. Twenty-seven distinct speech samples resulted, and each was processed using our silence deletion algorithm with ZSIL = 36, MINSP = 40, AON = 3, and AOFF = 2.5. Deleted acoustic material was replaced on playback with absolute silence. These processed samples, together with the 27 samples coded but not subject to silence deletion processing, constituted 54 test samples. To these were added three samples processed by our silence deletion algorithm, but with a copy of the sample's own background noise (approximately 20 dB signal/noise ratio) replacing absolute silence on playback. These 57 samples were randomly ordered on a two-sided casette as indicated in Table I. Samples 1–3 and 34–36, inclusive, were the original 10-s 12-bit unprocessed samples, presented for listener orientation. Samples 4 and 37 were presented three times during the test to provide a check on listener consistency. Each of 10 untrained subjects was given a Sony Walkman WM-4 cassette player, stereo headphones, a score sheet, and printed instructions as follows:

"The purpose of this subjective listening test is to assess the effect of silence deletion and the subsequent insertion on speech. Please state the degree of acceptability of test samples as recorded messages, on a scale of 1 to 5. 1 means that a sample is unacceptable and 5 denotes the highest degree of acceptance. Rate each sample on its own rather than by comparison with other samples. The test consists of two 10-minute segments recorded on sides 1 and 2 of the accompanying cassette. You are advised to take a 5-minute break between the two segments. There are 67, 10-second test samples recorded and numbered in sequence with a slight pause after each for recording your score. For your information and orientation, the first 3 samples on each side of cassette are original speech samples; the remainder being a random ordering of various processed speech samples."

Table II displays the subjects' scores for each of the three original unprocessed samples, and for samples 4 and 37. The average scores of 3.63 for the originals is well below the maximum, with the radio conversation lowest at 2.9. The original recordings, while not of excellent quality, are representative. Table II indicates good listener group consistency based on average scores for the same samples.

Table III enables determination of the subjective effects of silence deletion. There are two important conclusions: 1) averaged over all speech samples, replacement of deleted acoustic material with absolute silence on playback causes perceived quality to drop significantly, to an average of 2.44 from 3.51 without silence deletion; and 2) insertion of background noise in place of deleted acoustic material during playback tends to restore perceived speech quality toward its presilence-deletion values. The first conclusion follows from comparing the columns labeled "unedited" and "silence-edited" in Table III, and the difference between these in the right-hand column. The second follows from comparing the differ-

TABLE III
SUBJECTIVE RATING COMPARISONS

| Test Sample Class | *Unedited | Silence-edited | Noise-edited | Silence-minus Noise-edited or Unedited |
|---|---|---|---|---|
| (1) Overall Average | 3.51 | 2.44 | | -1.07 |
| (2) Speech Sample Origin | | | | |
| Sample A | 3.59 | 2.32 | | -1.27 |
| Sample B | 3.92 | 2.67 | | -1.25 |
| Sample C | 3.02 | 2.32 | | -0.70 |
| (3) Coding Scheme | | | | |
| A-law PCM | 3.53 | 2.52 | | -1.01 |
| ADPCM | 3.33 | 2.42 | | -0.91 |
| CVSDM | 3.67 | 2.37 | | -1.30 |
| (4) Noise-vs.-Silence-edited | | | | |
| Sample A | 4.00 | 2.30 | | -1.70 |
| 32 kbps CVSDM | 4.00 | | 3.40 | -0.60 |
| Sample B | 3.80 | 2.90 | | -0.90 |
| 7-bit ADPCM | 3.80 | | 3.40 | -0.40 |
| Sample C | 3.00 | 3.00 | | 0.00 |
| 6-bit A-Law PCM | 3.00 | | 3.30 | +0.30 |

*Not subject to silence deletion processing.

ence between the "silence-edited" and "noise-edited" columns in Table III.

By itself, digital coding/decoding using PCM, ADPCM, or CVSDM did not degrade the speech quality very much, on average, as is seen by comparing the "unedited" column in Table III to the "mean subjective rating" column for the three original samples in Table II. For example, compare Sample $A$'s 3.59 value in Table III to the average of 3.8 for Original 12-bit Samples in Table II. Similar comparison involving Table II and the "unedited" column area (4) of Table III shows minimal degradations for CVSDM, ADPCM, and A-law PCM coding at relatively high bit rates. Among the three coding schemes, degradation from deletion of acoustic material was highest CVSDM, probably because the parameters were chosen for good waveform tracking rather than enhanced SQNR. Among the three speech samples, degradation for silence-editing was largest for Sample $A$ whose high background noise level contrasted most during transitions to and from silence during playback.

Further details of this work appear elsewhere [5].

REFERENCES

[1] H. H. Lee and C. K. Un, "A study of the on–off characteristics of conversational speech," *IEEE Trans. Commun.*, vol. COM-34, pp. 630–637, June 1986.

927

[2] P. D. Drago, A. M. Molinari, and F. C. Vagliani, "Digital dynamic speech detectors," *IEEE Trans. Commun.*, vol. COM-26, pp. 140–145, Jan. 1978.

[3] P. de Souza, "A statistical approach to the design of an adaptive self-normalizing silence detector," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-31, pp. 678–684, June 1983.

[4] D. H. Cho and K. C. Un, "Hybrid companding delta modulation with silence deletion," in *Proc. IEEE Global Commun. Conf.*, Miami, FL, Nov. 1982, pp. 1340–1344.

[5] C. K. Gan, "Efficient speech storage via compression of silence periods," M.A.Sc. thesis, Univ. British Columbia, Dep. Elec. Eng., Dec. 1984.

# A Note on "Wigner Distribution for Finite Duration or Band-Limited Signals and Limiting Cases"

## F. HLAWATSCH

*Abstract*—A recent paper by Cohen[1] investigates the Wigner distribution of signals consisting of two finite-duration, nonoverlapping signal components. In this note, we present an alternative analysis which is considerably simplified as compared to Cohen's.

## I. INTRODUCTION

The Wigner distribution (WD) [1] is a time-frequency signal representation which has been shown to be useful for the analysis of time-varying and transient signals (see, e.g., [2] and [3]). Apart from being theoretically attractive, the WD is also the basis for practical time-frequency representations (like the *pseudo-Wigner distribution* or the well-known *spectrogram*) which meet, to a large extent, the requirements encountered in signal analysis applications.

The cross WD (CWD) of two signals $x(t)$, $y(t)$ is defined by

$$W_{x,y}(t,f) = \int_{-\infty}^{\infty} x(t + \tau/2) \, y^*(t - \tau/2) \, e^{-j2\pi f \tau} \, d\tau,$$

$$(1.1)$$

where $t$ and $f$ denote time and frequency, respectively. The CWD may also be expressed, in a similar way, using the signals' spectra $X(f)$, $Y(f)$,

$$W_{x,y}(t,f) = \int_{-\infty}^{\infty} X(f + \nu/2) \, Y^*(f - \nu/2) \, e^{j2\pi t \nu} \, d\nu.$$

$$(1.2)$$

Because of this symmetry of time and frequency, all results derived for, e.g., the time domain apply in the frequency domain as well. We shall discuss a *finite-support property* of the CWD to illustrate this important symmetry; this property will also be used in the subsequent development.

Suppose that the signals $x(t)$ and $y(t)$ are zero outside time intervals $[t_{x1}, t_{x2}]$ and $[t_{y1}, t_{y2}]$, respectively. It is then easily verified that, for all $\tau$, $x(t + \tau/2) \, y^*(t - \tau/2) = 0$ for $t$ outside the interval $[\bar{t}_1, \bar{t}_2]$, where

$$\bar{t}_1 = \frac{t_{x1} + t_{y1}}{2}, \qquad \bar{t}_2 = \frac{t_{x2} + t_{y2}}{2}. \qquad (1.3)$$

[1]L. Cohen, *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-35, pp. 796–806, June 1987.

Inserting into (1.1), it follows that also

$$W_{x,y}(t,f) = 0 \quad \text{for } t \text{ outside } [\bar{t}_1, \bar{t}_2]. \qquad (1.4)$$

Now by virtue of CWD's time-frequency symmetry, an analogous finite-support property exists with respect to the frequency domain: if the signals $x(t)$ and $y(t)$ are band-limited to frequency bands $[f_{x1}, f_{x2}]$ and $[f_{y1}, f_{y2}]$, respectively, then

$$W_{x,y}(t,f) = 0 \quad \text{for } f \text{ outside } [\bar{f}_1, \bar{f}_2], \qquad (1.5)$$

with

$$\bar{f}_1 = \frac{f_{x1} + f_{y1}}{2}, \qquad \bar{f}_2 = \frac{f_{x2} + f_{y2}}{2}. \qquad (1.6)$$

In most applications, the WD of a single signal is of primary interest. The (auto) WD $W_x(t,f)$ of a single signal $x(t)$ is defined by letting $y = x$ in (1.1), so that

$$W_x(t,f) \triangleq W_{x,x}(t,f). \qquad (1.7)$$

Properties (1.4) and (1.5) then reduce to the well-known finite-support properties of the (auto) WD: if the signal $x(t)$ is zero outside a time interval $[t_1, t_2]$, then the WD of $x(t)$ is zero outside the same interval (or, to be more precise, outside the corresponding strip in the $(t, f)$-plane). An analogous result, of course, again exists with respect to the frequency domain.

## II. WIGNER DISTRIBUTION OF TWO-BURST SIGNALS

In a recent paper by Cohen,[1] the WD of a signal $z(t)$ consisting of two finite-duration, nonoverlapping signal segments ("bursts") is considered. To be more specific, this "two-burst signal" is defined by

$$z(t) = \begin{cases} 0, & -\infty < t \le t_{x1} \\ x(t), & t_{x1} < t < t_{x2} \\ 0, & t_{x2} \le t \le t_{y1} \\ y(t), & t_{y1} < t < t_{y2} \\ 0, & t_{y2} \le t < \infty, \end{cases} \qquad (2.1)$$

which is illustrated by Fig. 1.

The paper contains an exhaustive enumeration of 19 different cases which are distinguished by the relative positions of $t_{x1}$, $t_{x2}$, $t_{y1}$, and $t_{y2}$. In each of these cases, different intervals of the time axis are distinguished, and on each of these intervals, the WD of $z(t)$ is expressed as (generally) a sum of at most three integrals which, in our framework, can be identified as WD's of $x(t)$ or $y(t)$ and/or CWD's of $x(t)$, $y(t)$ or $y(t)$, $x(t)$.

The present note is based on our belief that Cohen's discussion is unnecessarily complicated and difficult to read. This is a direct consequence of the comparatively great number of different cases and time intervals which are discussed *separately*. Moreover, the different cases are defined by multiple inequalities which are cumbersome to interpret. As we show in the following, the WD of a two-burst signal (2.1) can be analyzed in a far simpler way by a general treatment which encompasses *all* cases, irrespective of the relative positions of the interval boundaries $t_{x1}$ through $t_{y2}$. Our discussion will be based on WD's quadratic superposition principle (occurrence of *WD interference terms*) and the interference terms' geometrical properties—characteristics of the WD which are of great consequence in practical WD applications and which are discussed in more detail in [4] and [5].

We first note that the two-burst signal $z(t)$ is restricted to the time interval $[t_{x1}, t_{y2}]$. From WD's finite-support property, it thus follows that the WD of $z(t)$ is equally zero outside this interval. On the other hand, the WD is not identically zero in the gap $[t_{x2}, t_{y1}]$ between the two bursts. We shall now show that this is a consequence of WD's interference property.