



(11) **EP 3 016 314 A1**

(12) **EUROPEAN PATENT APPLICATION**

(43) Date of publication:
04.05.2016 Bulletin 2016/18

(51) Int Cl.:
H04L 9/00 (2006.01) **G06F 21/32 (2013.01)**
H04L 9/32 (2006.01) **G10L 17/12 (2013.01)**

(21) Application number: **14461584.6**

(22) Date of filing: **28.10.2014**

(84) Designated Contracting States:
AL AT BE BG CH CY CZ DE DK EE ES FI FR GB GR HR HU IE IS IT LI LT LU LV MC MK MT NL NO PL PT RO RS SE SI SK SM TR
Designated Extension States:
BA ME

(72) Inventors:
• **Galka, Jakub**
31-621 Krakow (PL)
• **Grzywacz, Marcin**
43-602 Jaworzno (PL)
• **Samborski, Rafal**
31-303 Krakow (PL)

(71) Applicant: **Akademia Gorniczo-Hutnicza im. Stanisława Staszica w Krakowie**
30-059 Krakow (PL)

(74) Representative: **Eupatent.pl**
ul. Zeligowskiego 3/5
90-752 Lodz (PL)

(54) **A system and a method for detecting recorded biometric information**

(57) A method for detecting a recorded biometric information, the method comprising preparing (201) a spectrogram $S_{n,k}$ of the recorded speech biometric information, wherein the spectrogram $S_{n,k}$ is an $N \times K$ matrix, where $S_{n,k}$ is a log-amplitude of the spectrogram at a given time frame n and frequency bin k ; the method being characterized in that it further comprises the steps of: normalizing (202) the spectrogram by its mean spectral value; subjecting (203) the spectrogram to high-pass digital filtering; extracting (204) a local maxima pair; and calculating (205) a similarity by comparing the output of the local maxima pair extraction with a previously stored data.

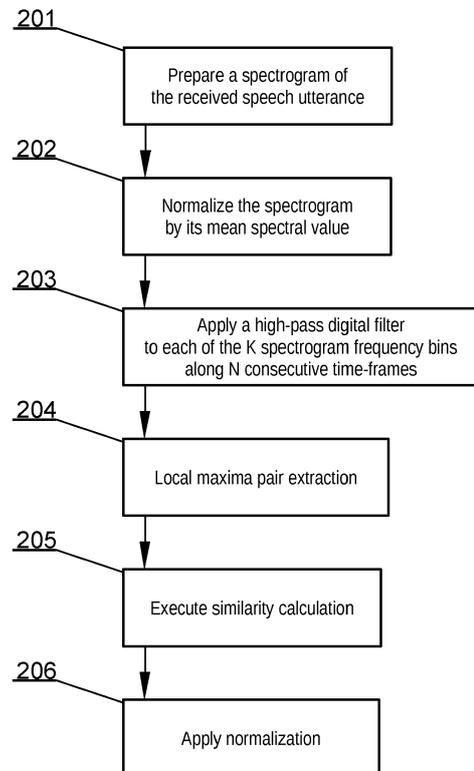


Fig. 2

EP 3 016 314 A1

Description

TECHNICAL FIELD

5 [0001] There are presented herein a system and a method for detecting recorded biometric information, which are useful in particular for detection of playback attack.

BACKGROUND

10 [0002] A playback attack (also known as a replay attack) is commonly defined as a form of a network attack, in which a valid data transmission is maliciously or fraudulently repeated or delayed. One of the biggest threats in biometric speaker verification systems are playback (replay) attacks, in which a previously recorded passphrase is played back by an unprivileged person in order to gain access.

15 [0003] There is a known concept of biometric speaker verification, the aim of which is to accept or reject the identity of the speaker based on the sample of their voice. Telephone-based automatic speaker verification by use of telephone channel has been evaluated. Despite the fact that such systems perform very well, consumers and organizations still have their doubts when it comes to high-security applications (e.g. e-banking). One of the prevailing arguments against voice biometry concerns the common passphrase text-dependent systems, in which the passphrase uttered by the speaker does not change from one login attempt to another login attempt.

20 [0004] In such cases it could be possible to break into such systems by playing-back a recording, which was obtained earlier using a microphone or any other eavesdropping method (e.g. malicious mobile software). This type of attack is called a playback attack and can be relatively easily performed by a person experienced with signal processing.

25 [0005] One of the solutions to this problem is to use a text-prompted system (in which the user is asked to say a randomly selected phrase for each access attempt), although such systems are more sensitive to other types of attacks (such as concatenation of previously recorded sounds) and, due to the fact that the system cannot use lexical knowledge in its assessment, achieve higher error rates as compared to text-dependent solutions.

30 [0006] A publication "A Playback Attack Detector for Speaker Verification Systems" by Wei Shang and Maryhelen Stevenson of Department of Electrical and Computer Engineering University of New Brunswick (ISCCSP 2008, Malta, 12-14 March 2008) discloses a playback attack detector (PAD), which can be mobilized in guarding speaker verification systems against playback attacks. To detect playback attacks, the PAD uses a feature set called a peakmap, which includes frame and FFT bin numbers of the five highest spectral peaks from each of the voiced frames in an utterance. During the detection, the peakmap of an incoming recording is first extracted and then compared to those of all the other recordings that are stored at the system end. Each comparison will yield a similarity score that represents the level of similarity between the two recordings. The incoming recording is declared to be a playback recording if its maximum similarity score is above a threshold.

35 [0007] A US patent application US2010131279 discloses a method for controlling user access to a service available in a data network and/or to information stored in a user database, in order to protect stored user data from unauthorized access, wherein the method comprises the following steps: input of a user's speech sample to a user data terminal, processing of the user's speech sample in order to obtain a prepared speech sample as well as a current voice profile of the user, comparison of the current voice profile with an initial voice profile stored in an authorization database, and output of an access-control signal to either permit or refuse access, taking into account the result of the comparison step, such that the comparison step includes a quantitative similarity evaluation of the current and the stored voice profiles as well as a threshold-value discrimination of a similarity measure thereby derived, and an access-control signal that initiates permission of access is generated only if a prespecified similarity measure is not exceeded.

40 [0008] A US patent application US2013006626 discloses a system, which includes a login process controller; a speech recognition module; a speaker verification module; a speech synthesis module; and a user database. Responsive to a user-provided first verbal answer to a first verbal question, the first verbal answer is converted to text and compared with data previously stored in the user database. The speech synthesis module provides a second question to the user, and responsive to a user-provided second verbal answer to the second question, the speaker verification module compares the second verbal answer with a voice print of the user previously stored in the user database and validates that the second verbal answer matches a voice print of the user previously stored in the user database. There is also disclosed a method of logging in to the telecommunications system and a computer program product for logging in to the telecommunications system.

45 [0009] It would be advantageous to improve the security of text-dependent, biometric speaker verification systems against playback attacks. It would be also advantageous to provide an improved and resource-effective system and method for detecting recorded biometric information.

SUMMARY

[0010] There is presented a method for detecting a recorded biometric information, the method comprising: preparing a spectrogram $S_{n,k}$ of the recorded speech biometric information, wherein the spectrogram $S_{n,k}$ is an $N \times K$ matrix, where $S_{n,k}$ is a log-amplitude of the spectrogram at a given time frame n and frequency bin k ; wherein the method further comprises the steps of: normalizing the spectrogram by its mean spectral value; subjecting the spectrogram to high-pass digital filtering; extracting a local maxima pair; and calculating a similarity by comparing the output of the local maxima pair extraction with a previously stored data.

[0011] Preferably, step of preparing a spectrogram includes Hamming windowing, wherein the windows overlap.

[0012] Preferably, the Hamming window is from 20 ms to 100 ms and the overlap is from 25% to 75%.

[0013] Preferably, the step of normalizing the spectrogram comprises calculating a normalized spectrogram

$$S_{n,k}^* = S_{n,k} - \mu_S, \quad (1)$$

[0014] wherein for $n = 1, \dots, N$ and $k = 1, \dots, K$:

$$\mu_S = \frac{1}{N * K} \sum_{n=1}^N \sum_{k=1}^K S_{n,k} \quad (2)$$

[0015] Preferably, the high-pass digital filter is applied to each of the K spectrogram frequency bins along N consecutive time-frames:

$$H(z) = \frac{1 - z^{-1}}{1 - \alpha z^{-1}} \quad (3)$$

[0016] Preferably, the step of executing a local maxima pair extraction comprises the steps of: (a) for each i -th local spectral maximum, extracting its coordinates (n_i, k_i) according to a differential criterion; (b) pruning the obtained set of peak coordinates (n_i, k_i) in order to eliminate insignificant peaks; (c) setting all the other frequency bins to a zero-value; (d) eliminating peak candidates by masking selected candidates with the envelope $en(k)$ obtained by a convolution of the spectral frame with a Gaussian window; (e) for each peak, inside the target region, combining this seed peak with other peaks from the target region into pairs; (f) representing the set of such pairs by a 4-column matrix, where each row corresponds to one pair of maxima for all the pairs found in the spectrogram wherein each row consists of 4 integer values: (k_i, n_i) - coordinates of the seed peak; (dk_i, dn_i) - frequency and time shifts between peaks paired with the seed peak; (g) sorting the matrix row-wise and removing duplicated rows as well as the first column of the matrix.

[0017] Preferably, the step of executing a similarity calculation preserves the time-ordering of the maxima pairs and each comparison results in a specific number of identical rows, wherein this number creates the hits vector U of M length, wherein M is the number of recordings of a particular speaker in a database.

[0018] Preferably, the method further comprises the step of normalizing the similarity calculation output.

[0019] Preferably, the normalization is based on L_C cohort cores.

[0020] There is also presented a computer program comprising program code means for performing all the steps of the computer-implemented method as described above when said program is run on a computer, as well as a computer readable medium storing computer-executable instructions performing all the steps of the computer-implemented method as described above when executed on a computer.

[0021] There is also presented a system for detecting pre-recorded biometric information, the system comprising: a data bus communicatively coupled to a memory; a controller communicatively coupled to the data bus and configured to prepare a spectrogram $S_{n,k}$ of the recorded speech biometric information, wherein the spectrogram $S_{n,k}$ is an $N \times K$ matrix, wherein $S_{n,k}$ is a log-amplitude of the spectrogram at a given time frame n and frequency bin k ; a spectrogram normalization module configured to normalize the spectrogram by its mean spectral value; a high-pass filter module configured to high-pass filter the spectrogram $S_{n,k}$; a maxima pair extraction module configured to extract a local maxima

pair; the controller being further configured to calculate a similarity by comparing the output of the local maxima pair extraction with a previously stored data.

BRIEF DESCRIPTION OF FIGURES

[0022] These and other objects of the invention presented herein are accomplished by providing a system and a method for detecting recorded biometric information. Further details and features of the presented method and system, their nature and various advantages will become more apparent from the following detailed description of the preferred embodiments shown in a drawing, in which:

Fig. 1 presents a diagram of the system for detecting recorded biometric information;

Fig. 2 presents a diagram of the method for detecting recorded biometric information;

Figs. 3A-B present examples of a spectrogram with its maxima pairs plotted for a reference pattern and for a playback attack pattern.

NOTATION AND NOMENCLATURE

[0023] Some portions of the detailed description which follows are presented in terms of data processing procedures, steps or other symbolic representations of operations on data bits that can be performed on computer memory. Therefore, a computer executes such logical steps thus requiring physical manipulations of physical quantities.

[0024] Usually these quantities take the form of electrical or magnetic signals capable of being stored, transferred, combined, compared, and otherwise manipulated in a computer system. For reasons of common usage, these signals are referred to as bits, packets, messages, values, elements, symbols, characters, terms, numbers, or the like.

[0025] Additionally, all of these and similar terms are to be associated with the appropriate physical quantities and are merely convenient labels applied to these quantities. Terms such as "processing" or "creating" or "transferring" or "executing" or "determining" or "detecting" or "obtaining" or "selecting" or "calculating" or "generating" or the like, refer to the action and processes of a computer system that manipulates and transforms data represented as physical (electronic) quantities within the computer's registers and memories into other data similarly represented as physical quantities within the memories or registers or other such information storage.

[0026] A computer-readable (storage) medium, such as referred to herein, typically may be non-transitory and/or comprise a non-transitory device. In this context, a non-transitory storage medium may include a device that may be tangible, meaning that the device has a concrete physical form, although the device may change its physical state. Thus, for example, non-transitory refers to a device remaining tangible despite a change in state.

DETAILED DESCRIPTION

[0027] Fig. 1 presents a diagram of a system for detecting recorded biometric information. The system may be realized using dedicated components or custom made FPGA or ASIC circuits. The system comprises a data bus 101 communicatively coupled to a memory 104. In addition, other components of the system are communicatively coupled to the system bus 101, so that they may be managed by a controller 105.

[0028] The memory 104 may store a computer program or a plurality of programs to be executed by the controller 105 in order to execute steps of the method for detecting recorded biometric information. The memory may further store digital speech data to be tested against a playback attack, as well as intermediate and final processing results.

[0029] The system comprises a spectrogram normalization module 102 configured to normalize a spectrogram of the analyzed speech by its mean spectral value. This process corresponds to step 202 of the procedure shown in Fig. 2.

[0030] The system further comprises a high-pass digital filter 103 configured such that the filter is applied to each of the K spectrogram frequency bins along N consecutive time-frames. This process corresponds to step 203 of the procedure shown in Fig. 2.

[0031] The next module of the system is a maxima pair extraction module 106 configured to execute local maxima pair extraction according to step 204 of the procedure shown in Fig. 2.

[0032] The system further comprises a results normalization module 107 configured to execute similarity calculation and results normalization according to steps 205 and 206 of the procedure shown in Fig. 2.

[0033] Fig. 2 presents a diagram of a method for detecting recorded biometric information. The method starts at step 201 from preparing a spectrogram of the received speech utterance, which is a preferred form of representation of a signal that allows for features analysis. For this purpose, a window function may be applied. In signal processing, a window function (also known as an apodization function or tapering function) is a mathematical function that is zero-valued outside of a chosen interval.

[0034] Preferably, a 512-point FFT (Fast Fourier transform) with 64ms Hamming windowing and 32ms overlapping is

used to obtain the spectral features of the recording. Other windowing times may be used as well, for example from 48 ms to 80 ms. The overlapping is optional and it serves to improve efficiency, other overlapping ranges are also possible. For simpler implementations, overlapping may be avoided. These values have been determined experimentally to provide the optimal spectral resolution for the extraction of the most significant features.

5 **[0035]** The windowing need not be Hamming windowing and other approaches may be used, which does not alter the operating principle but only affects system efficiency.

[0036] Let the spectrogram S be a $N \times K$ matrix, where $S_{n,k}$ is a log-amplitude of the spectrogram at a given time frame n and frequency bin k (a bin refers to a given frequency while a window may comprise a plurality of bins).

10 **[0037]** In order to reduce border-effects introduced in the next steps of the procedure, the spectrogram is normalized by its mean spectral value at step 202.

$$15 \quad S_{n,k}^* = S_{n,k} - \mu_S, \quad (1)$$

wherein

$$20 \quad \mu_S = \frac{1}{N * K} \sum_{n=1}^N \sum_{k=1}^K S_{n,k} \quad (2)$$

25

for $n = 1, \dots, N$ and $k = 1, \dots, K$.

[0038] At step 203, there is subsequently applied a high-pass digital filter to each of the K spectrogram frequency bins along N consecutive time-frames to increase the contrast of the spectrogram and to improve maxima extraction efficiency.

30 The closer the pole is to $z = 1$, the less spectral maxima are found. During the performance test, the value $\alpha = 0.98$ was found to produce the best results (nevertheless a preferred range of 0.98 to 1 may be used).

$$35 \quad H(z) = \frac{1 - z^{-1}}{1 - \alpha z^{-1}} \quad (3)$$

[0039] The filtering discards slowly changing values (numerous local maxima in proximity to each other. Decrease of this parameter will affect local maxima extraction as there will be fewer of them due to smoothing of the spectrogram).

[0040] Next, at step 204, there is executed a local maxima pair extraction. The method of the local maxima pair extraction may, for example, be based on the method published in: A. L. chun Wang, An industrial-strength audio search algorithm, in: Proceedings of the 4th International Conference on Music Information Retrieval, 2003.

45 **[0041]** For each i -th local spectral maximum, its coordinates (n_i, k_i) are extracted according to the following differential criterion:

$$50 \quad (n_i, k_i) : S_{n_i, k_i} > \max\{S_{n_i, k_i-1}, S_{n_i, k_i+1}, S_{n_i-1, k_i} + d, S_{n_i+1, k_i}\} \quad (4)$$

wherein d is a temporal peak-masking decay coefficient. The obtained set of peak coordinates (n_i, k_i) is then pruned in order to eliminate insignificant peaks. That is, for each n -th frame, only the top-5 spectral maxima are left in the set.

[0042] All the other frequency bins are set to a zero-value. The resulting sparse spectrogram is then processed further, as the elimination of peak candidates is performed by way of masking selected candidates with the envelope $en(k)$ obtained by a convolution of the spectral frame with a Gaussian window. The width of the spreading Gaussian window can be adjusted during system evaluation. Only peak candidates higher than the envelope:

wherein $p = 1.1$ is an overmasking factor, are eventually selected as the set of spectral maxima to be used for

utterance parameterization. The value of p has been determined by experimentation.

[0043] Having the set of local maxima of the spectrogram, inside the target region R_i

$$\{(n_i, k_i) : S_{n,k} > e_n(k)/\rho\}, \quad (5)$$

$$R_i(n, k) = \{(n, k) : n_i < n < n_i + r \wedge k_i - p < k < k_i + p\} \quad (6)$$

of each peak, this seed peak is then combined with other peaks from the target region into pairs.

[0044] Parameters r and p determine the size of the target region along the time and frequency dimension respectively. The maximum number of pairs originating from each seed peak is limited by selecting the closest peak candidates.

[0045] The set of such pairs is represented by a 4-column matrix G , wherein each row corresponds to one pair of maxima for all the pairs found in the spectrogram. Each row consists of 4 integer values: (k_i, n_i) - coordinates of the seed peak; (dk_i, dn_i) - frequency and time shifts between peaks paired with the seed peak. The matrix is sorted row-wise and duplicated rows are removed.

[0046] The first redundant column is removed as well. Each row of the pair-matrix G is 20 bits in size and it is called marker for further reference. The marker is a concatenation of 8 bits of k_i , 6 bits of dk_i and 6 bits of dn_i . The matrix G of the original utterance is ready to be saved in the system's database and to be used in the scoring phase. An example of a spectrogram (a) with its maxima pairs plotted (b) is presented in Fig. 3A and Fig. 3B.

[0047] Figs. 3A-B present spectrograms for a reference pattern and a playback attack pattern respectively, i.e. examples of the authentic and playback utterances. There are presented spectrograms with peaks for: a) authentic utterance, c) playback utterance. There are also presented maxima pairs constellations for: b) authentic utterance, d) playback utterance. Solid lines on b), d) show similarity between authentic recording and its played-back version. The graphic shows simplified utterances' features sets to maintain the clarity of the figure. The object is to find the similarities. The similarities measured as common pairs of maxima of histograms are marked with solid lines (one line - one pair). On both recordings, the solid lines are in the same places (in contrast to dashed lines, which indicate differences).

[0048] Subsequently, at step 205, there is executed similarity calculation. In this step the pair-matrix G_{test} of the tested utterance is compared with all the pair-matrices of a particular target, stored during previous verification attempts (or during the step of setting up a biometric speaker verification access to a given system).

[0049] The comparison is conducted in a way to preserve the time-ordering of the maxima pairs. Each comparison results in a specific number of identical rows. This number creates the hits vector U having the length of M , wherein M is the number of recordings of a particular speaker in a database. In order to compute the similarity of a test recording to reference patterns, an L coefficient may be computed. Each $L(i)$ is a ratio of a count of the same markers in the test recording and a reference recording (i) to a count of all markers in the reference recording (i).

[0050] Therefore, each $L(i)$ is a set of baseline system scores. The ratio of each element of U to the number of rows of its corresponding reference G is calculated, wherein Y_i is the number of rows in G . For example, if a comparison is made between the tested utterance's pair-matrix G_{test} with the reference pair-matrix G_A of the recording A from the speaker's database, the ratio of the number of obtained identical rows to the number of all rows G_A is calculated. The PAD algorithm returns the value L_{max} , which is the maximum score of the baseline detection scores vector L . The obtained value is then used as a baseline similarity

$$L(i) = \frac{U(i)}{Y_i}, \quad (7)$$

score to test the playback attack hypothesis. The baseline similarity score is meant as the raw score without normalization.

[0051] In order to gain accuracy and robustness, system normalization techniques are used at step 206. The main goal of detection score normalization is to increase the system's accuracy by making the scoring space independent from the varying testing, training, and other operational conditions. Applying score normalization in the PAD algorithm is used to improve its robustness and accuracy as well.

[0052] Let L_C be a so-called cohort-score vector of baseline detection scores L , wherein the maximum score L_{max} is

excluded from the original score set. Normalization approaches based on L_C cohort cores may for example utilize techniques disclosed in: Y. Zigel, A. Cohen, On cohort selection for speaker verification., in: INTERSPEECH, ISCA, 2003 (URL <http://dblp.uni-trier.de/db/conf/interspeech/interspeech2003.html>).

[0053] There may be applied different methods of cohort-based score normalization, such as:

$$L_{Cnorm1} = \frac{L_{max}}{\max(L_C)}, \quad (8)$$

$$L_{Cnorm2} = L_{max} - \mu_{L_C}, \quad (9)$$

$$L_{Cnorm3} = \frac{L_{max}}{\mu_{L_C}}. \quad (10)$$

[0054] In the case of cohort normalization, to normalize scores using T-normalization (C. Barras, J. Gauvain, Feature and score normalization for speaker verification of cellular data, in: Acoustics, Speech, and Signal Processing, 2003. Proceedings. (ICASSP '03). 2003 IEEE International Conference on, Vol. 2, 2003, pp. II-49-52 vol.2. doi:10.1109/ICASSP.2003.1202291) there may be used a cohort of L_C scores. The size of the L_C cohort is limited to 30, the closest to the L_{max} value, elements only. T-normalization is performed using the mean μ_T and standard deviation σ_T parameters of the L_C cohort:

$$L_{Tnorm}(i) = \frac{L(i) - \mu_T}{\sigma_T + \epsilon} \quad (11)$$

where $\epsilon = 0.00431$ when $\sigma = 0$ is used as a backup (precalculated off-line from test where $\sigma \neq 0$) and $\epsilon = 0$ otherwise.

[0055] Another method of score normalization is Z-normalization (R. Auckenthaler, Score Normalization for Text-Independent Speaker Verification Systems, Digital Signal Processing 10 (1-3) (2000) 42-54. doi:10.1006/dspr.1999.0360), where the mean μ_Z and the standard deviation σ_Z parameters, calculated for each of the speaker's reference recordings, are obtained by scoring the recording against a set of randomly chosen reference recordings of other speakers:

The main advantage of the Z-norm is that the parameters μ_Z and σ_Z can be pre-calculated offline.

$$L_{Znorm}(i) = \frac{L(i) - \mu_Z}{\sigma_Z + \epsilon}, \quad (12)$$

[0056] Combinations of different normalization methods, such as ZT-norm (first a normalization using T-norm, next using Z-norm), TZ-norm and the normalization described by the formula which combines the properties of T- and Z-normalization simultaneously, may also be applied.

$$L_{Xnorm}(i) = \frac{L(i) - \frac{\mu_T + \mu_Z}{2}}{\frac{1}{4} \sqrt{\sigma_T^2 + \sigma_Z^2}}, \quad (13)$$

[0057] As it turns out from comparison of the test results of the presented PAD system and the expectations of the market, it is possible to use the presented solution in real-world applications. This method, due to spectral-pairs-based detection, can be used to detect artificially prepared (tailored) playback utterances as well, with the help of concatenation or other, more advanced speech modification methods.

[0058] The presented method and system allow for biometric speaker verification and prevention of playback attack, which is very useful in authorization systems. Therefore, the presented method and system provide a useful, concrete and tangible result.

[0059] Due to the fact that speech data are processed and transformed in order to execute biometric speaker verification and prevention of playback attack and to the fact that data processing takes place in a dedicated, programmed machine, the machine or transformation test is fulfilled and the idea is not abstract.

[0060] It can be easily recognized, by one skilled in the art, that the aforementioned method for detecting recorded biometric information may be performed and/or controlled by one or more computer programs. Such computer programs are typically executed by utilizing the computing resources in a computing device. Applications are stored on a non-transitory medium. An example of a non-transitory medium is a non-volatile memory, for example a flash memory or volatile memory, for example RAM. The computer instructions are executed by a processor. These memories are exemplary recording media for storing computer programs comprising computer-executable instructions performing all the steps of the computer-implemented method according to the technical concept presented herein.

[0061] While the method and system presented herein has been depicted, described, and has been defined with reference to particular preferred embodiments, such references and examples of implementation in the foregoing specification do not imply any limitation on the presented method or system. It will, however, be evident that various modifications and changes may be made thereto without departing from the broader scope of the technical concept. The presented preferred embodiments are exemplary only, and are not exhaustive of the scope of the technical concept presented herein.

[0062] Accordingly, the scope of protection is not limited to the preferred embodiments described in the specification, but is only limited by the claims that follow.

Claims

1. A method for detecting a recorded biometric information, the method comprising:

- preparing (201) a spectrogram $S_{n,k}$ of the recorded speech biometric information, wherein the spectrogram $S_{n,k}$ is an $N \times K$ matrix, where $S_{n,k}$ is a log-amplitude of the spectrogram at a given time frame n and frequency bin k ;

the method being **characterized in that** it further comprises the steps of:

- normalizing (202) the spectrogram by its mean spectral value;
 - subjecting (203) the spectrogram to high-pass digital filtering;
 - extracting (204) a local maxima pair; and
 - calculating (205) a similarity by comparing the output of the local maxima pair extraction with a previously stored data.

2. The method according to claim 1 wherein the step of preparing a spectrogram includes Hamming windowing, wherein the windows overlap.

3. The method according to claim 2 wherein the Hamming window is from 20 ms to 100 ms and the overlap is from 25% to 75%.

4. The method according to claim 1 wherein the step of normalizing the spectrogram comprises calculating a normalized spectrogram

$$S_{n,k}^* = S_{n,k} - \mu_S, \quad (1)$$

5

wherein for $n = 1, \dots, N$ and $k = 1, \dots, K$:

$$\mu_S = \frac{1}{N * K} \sum_{n=1}^N \sum_{k=1}^K S_{n,k} \quad (2)$$

15

5. The method according to claim 4, wherein the high-pass digital filter is applied to each of the K spectrogram frequency bins along N consecutive time-frames:

20

$$H(z) = \frac{1 - z^{-1}}{1 - \alpha z^{-1}} \quad (3)$$

25

6. The method according to claim 1 wherein the step of executing a local maxima pair extraction comprises the steps of:

30

- (a) for each i-th local spectral maximum, extracting its coordinates (n_i, k_i) according to a differential criterion;
- (b) pruning the obtained set of peak coordinates (n_i, k_i) in order to eliminate insignificant peaks;
- (c) setting all the other frequency bins to a zero-value;
- (d) eliminating peak candidates by masking selected candidates with the envelope $en(k)$ obtained by a convolution of the spectral frame with a Gaussian window;
- (e) for each peak, inside the target region, combining this seed peak with other peaks from the target region into pairs;
- (f) representing the set of such pairs by a 4-column matrix, where each row corresponds to one pair of maxima for all the pairs found in the spectrogram wherein each row consists of 4 integer values: (k_i, n_i) - coordinates of the seed peak; (dk_i, dn_i) - frequency and time shifts between peaks paired with the seed peak; and
- (g) sorting the matrix row-wise and removing duplicated rows as well as the first column of the matrix.

40

7. The method according to claim 1, wherein the step of executing a similarity calculation preserves the time-ordering of the maxima pairs and each comparison results in a specific number of identical rows, wherein this number creates the hits vector U of M length, wherein M is the number of recordings of a particular speaker in a database.

45

8. The method according to claim 1, further comprising the step of normalizing (206) the similarity calculation (205) output.

9. The method according to claim 8, wherein the normalization is based on L_C cohort cores.

50

10. A computer program comprising program code means for performing all the steps of the computer-implemented method according to any of claims 1-9 when said program is run on a computer.

11. A computer readable medium storing computer-executable instructions performing all the steps of the computer-implemented method according to any of claims 1-9 when executed on a computer.

55

12. A system for detecting pre-recorded biometric information, the system comprising:

- a data bus (101) communicatively coupled to a memory (104);
- a controller (105) communicatively coupled to the data bus (101) and configured to prepare a spectrogram

EP 3 016 314 A1

$S_{n,k}$ of the recorded speech biometric information, wherein the spectrogram $S_{n,k}$ is an $N \times K$ matrix, wherein $S_{n,k}$ is a log-amplitude of the spectrogram at a given time frame n and frequency bin k ;

- a spectrogram normalization module (102) configured to normalize the spectrogram by its mean spectral value;

5

- a high-pass filter module (103) configured to high-pass filter the spectrogram $S_{n,k}$;

- a maxima pair extraction module (106) configured to extract a local maxima pair;

- the controller (105) being further configured to calculate a similarity by comparing the output of the local maxima pair extraction with a previously stored data.

10

15

20

25

30

35

40

45

50

55

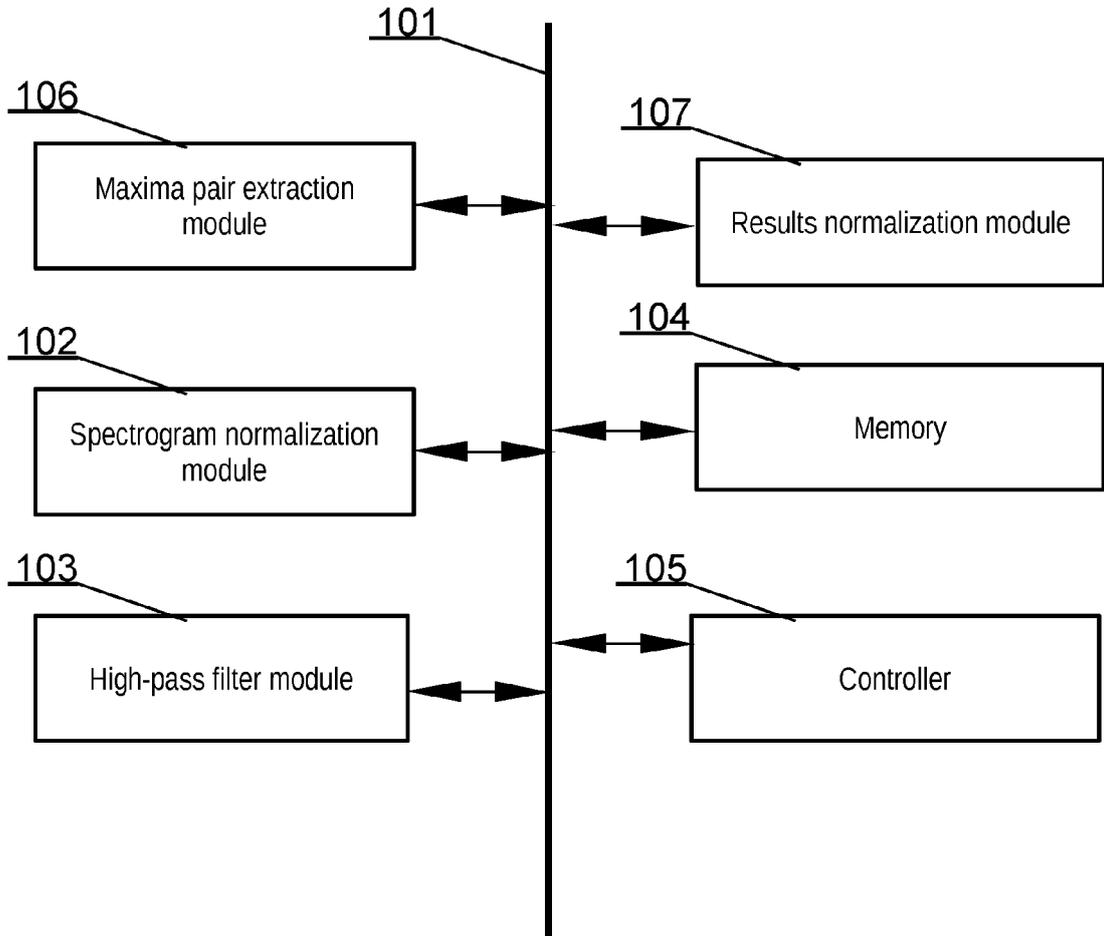


Fig. 1

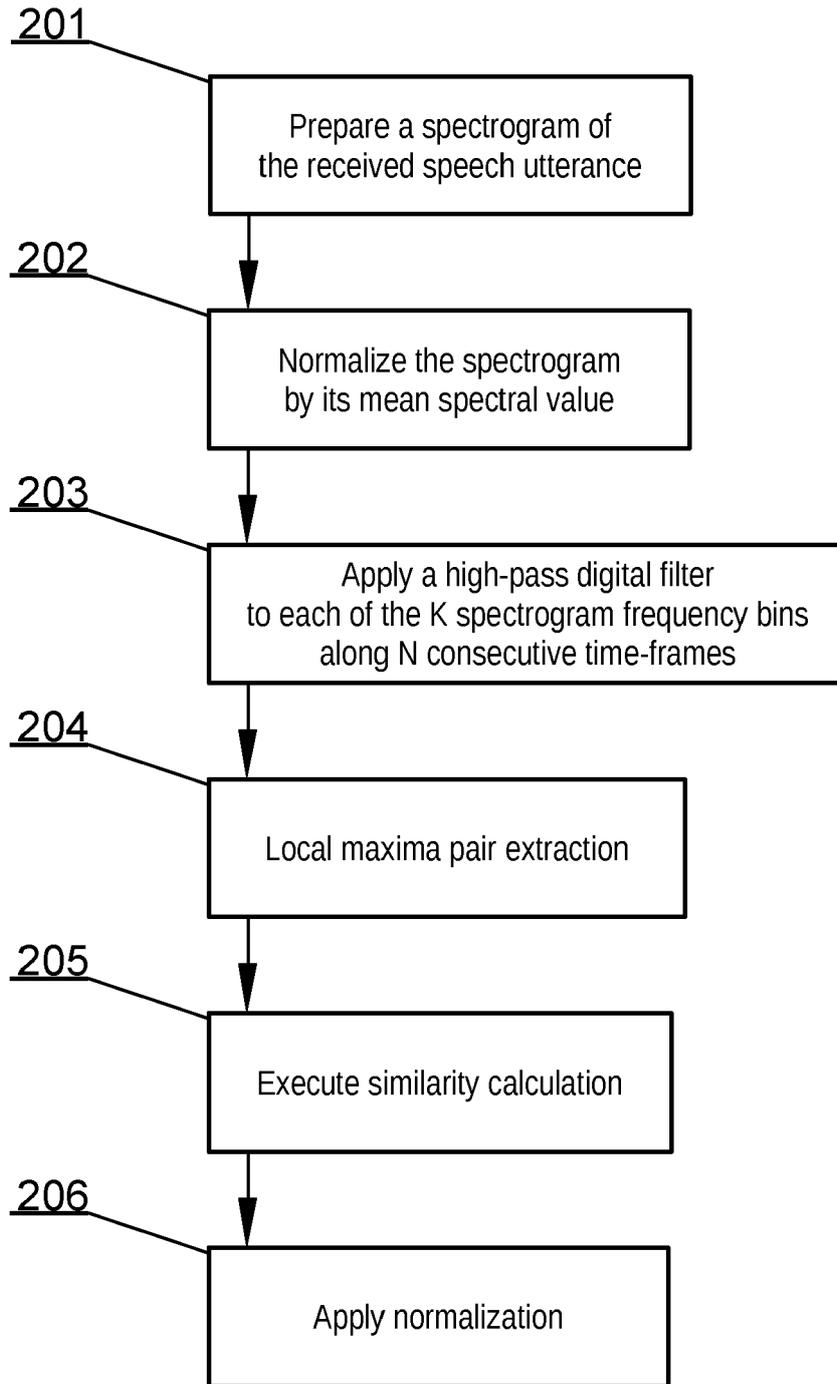


Fig. 2

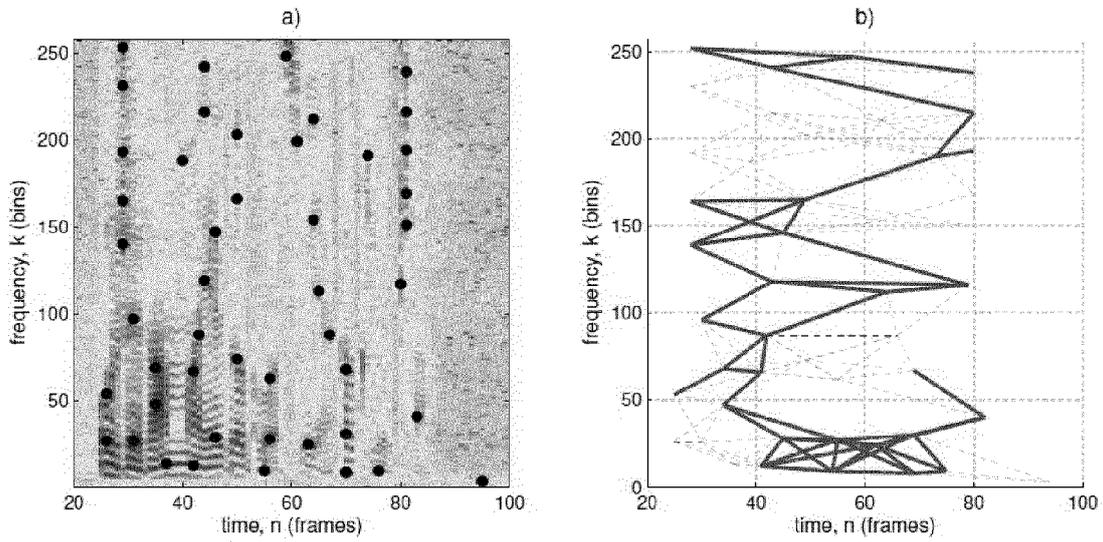


Fig. 3A

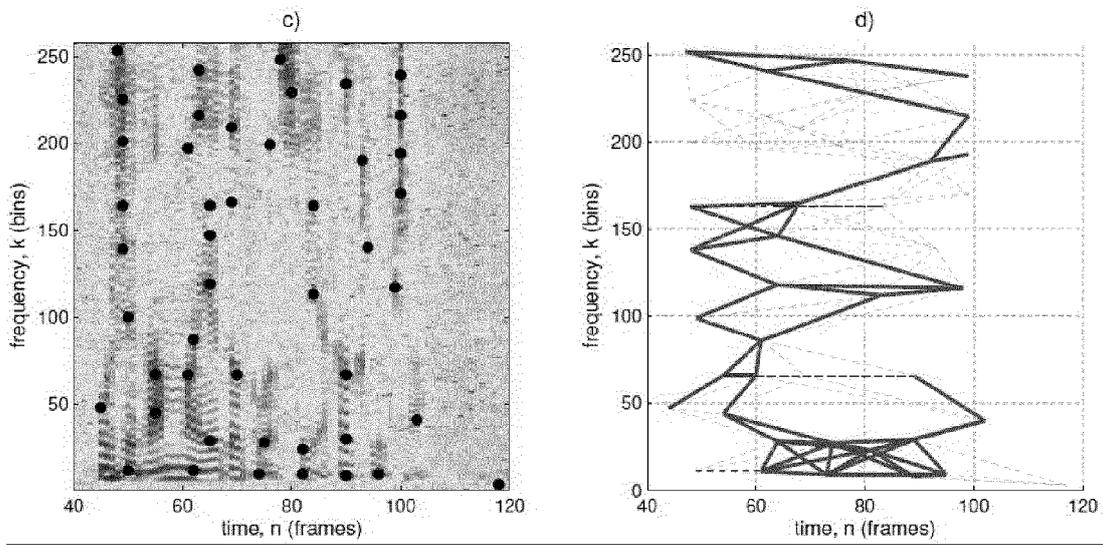


Fig. 3B



EUROPEAN SEARCH REPORT

Application Number
EP 14 46 1584

5

10

15

20

25

30

35

40

45

50

55

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (IPC)
X,P	JAKUB GALKA ET AL: "Playback attack detection for text-dependent speaker verification over telephone channels", SPEECH COMMUNICATION, vol. 67, 1 March 2015 (2015-03-01), pages 143-153, XP055182502, ISSN: 0167-6393, DOI: 10.1016/j.specom.2014.12.003 * paragraph [0011] * -----	1-12	INV. H04L9/00 G06F21/32 H04L9/32 G10L17/12
A,D	WEI SHANG ET AL: "A playback attack detector for speaker verification systems", COMMUNICATIONS, CONTROL AND SIGNAL PROCESSING, 2008. ISCCSP 2008. 3RD INTERNATIONAL SYMPOSIUM ON, IEEE, PISCATAWAY, NJ, USA, 12 March 2008 (2008-03-12), pages 1144-1149, XP031269241, ISBN: 978-1-4244-1687-5 * paragraph [0011] - paragraph [001V] * -----	1-12	TECHNICAL FIELDS SEARCHED (IPC)
A	ZHI-FENG WANG ET AL: "Channel pattern noise based playback attack detection algorithm for speaker recognition", MACHINE LEARNING AND CYBERNETICS (ICMLC), 2011 INTERNATIONAL CONFERENCE ON, IEEE, 10 July 2011 (2011-07-10), pages 1708-1713, XP031966488, DOI: 10.1109/ICMLC.2011.6016982, ISBN: 978-1-4577-0305-8 * paragraph [0002] * ----- -/--	1-12	H04L G06F G10L
The present search report has been drawn up for all claims			
Place of search Munich		Date of completion of the search 13 April 2015	Examiner Bec, Thierry
CATEGORY OF CITED DOCUMENTS X : particularly relevant if taken alone Y : particularly relevant if combined with another document of the same category A : technological background O : non-written disclosure P : intermediate document T : theory or principle underlying the invention E : earlier patent document, but published on, or after the filing date D : document cited in the application L : document cited for other reasons & : member of the same patent family, corresponding document			

EPO FORM 1503 03/02 (P04C01)



EUROPEAN SEARCH REPORT

Application Number
EP 14 46 1584

5

10

15

20

25

30

35

40

45

50

55

DOCUMENTS CONSIDERED TO BE RELEVANT				
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (IPC)	
A	<p>WU ZHIZHENG ET AL: "Synthetic speech detection using temporal modulation feature", 2013 IEEE INTERNATIONAL CONFERENCE ON ACOUSTICS, SPEECH AND SIGNAL PROCESSING (ICASSP); VANCOUCER, BC; 26-31 MAY 2013, INSTITUTE OF ELECTRICAL AND ELECTRONICS ENGINEERS, PISCATAWAY, NJ, US, 26 May 2013 (2013-05-26), pages 7234-7238, XP032508956, ISSN: 1520-6149, DOI: 10.1109/ICASSP.2013.6639067 [retrieved on 2013-10-18] * paragraph [0111] *</p> <p>-----</p>	1-12		
A	<p>CORREIA M J ET AL: "Preventing converted speech spoofing attacks in speaker verification", 2014 37TH INTERNATIONAL CONVENTION ON INFORMATION AND COMMUNICATION TECHNOLOGY, ELECTRONICS AND MICROELECTRONICS (MIPRO), MIPRO, 26 May 2014 (2014-05-26), pages 1320-1325, XP032622913, DOI: 10.1109/MIPRO.2014.6859772 [retrieved on 2014-07-17] * paragraph [0111] - paragraph [0014] *</p> <p>-----</p>	1-12		TECHNICAL FIELDS SEARCHED (IPC)
A,D	<p>EP 2 192 576 A1 (VOICE TRUST AG [DE]) 2 June 2010 (2010-06-02) * paragraph [0037] - paragraph [0048] *</p> <p>-----</p>	1-12		
A	<p>US 6 084 977 A (BORZA STEPHEN J [CA]) 4 July 2000 (2000-07-04) * column 4, line 31 - column 7, line 36 *</p> <p>-----</p>	1-12		
The present search report has been drawn up for all claims				
Place of search Munich		Date of completion of the search 13 April 2015	Examiner Bec, Thierry	
<p>CATEGORY OF CITED DOCUMENTS</p> <p>X : particularly relevant if taken alone Y : particularly relevant if combined with another document of the same category A : technological background O : non-written disclosure P : intermediate document</p> <p>T : theory or principle underlying the invention E : earlier patent document, but published on, or after the filing date D : document cited in the application L : document cited for other reasons & : member of the same patent family, corresponding document</p>				

EPO FORM 1503 03/82 (P04/C01)

**ANNEX TO THE EUROPEAN SEARCH REPORT
ON EUROPEAN PATENT APPLICATION NO.**

EP 14 46 1584

5 This annex lists the patent family members relating to the patent documents cited in the above-mentioned European search report.
The members are as contained in the European Patent Office EDP file on
The European Patent Office is in no way liable for these particulars which are merely given for the purpose of information.

13-04-2015

10	Patent document cited in search report	Publication date	Patent family member(s)	Publication date
15	EP 2192576 A1	02-06-2010	DE 102008058883 A1 EP 2192576 A1 US 2010131279 A1	27-05-2010 02-06-2010 27-05-2010
20	US 6084977 A	04-07-2000	NONE	
25				
30				
35				
40				
45				
50				
55				

EPO FORM P0459

For more details about this annex : see Official Journal of the European Patent Office, No. 12/82

REFERENCES CITED IN THE DESCRIPTION

This list of references cited by the applicant is for the reader's convenience only. It does not form part of the European patent document. Even though great care has been taken in compiling the references, errors or omissions cannot be excluded and the EPO disclaims all liability in this regard.

Patent documents cited in the description

- US 2010131279 A [0007]
- US 2013006626 A [0008]

Non-patent literature cited in the description

- A Playback Attack Detector for Speaker Verification Systems. **WEI SHANG ; MARYHELEN STEVENSON**. ISCCSP 2008. Department of Electrical, 12 March 2008 [0006]
- **A. L. CHUN WANG**. An industrial-strength audio search algorithm. *Proceedings of the 4th International Conference on Music Information Retrieval*, 2003 [0040]
- **C. BARRAS ; J. GAUVAIN**. Feature and score normalization for speaker verification of cellular data. *Acoustics, Speech, and Signal Processing, 2003. Proceedings. (ICASSP '03). 2003 IEEE International Conference on*, 2003, vol. 2, II-49-52 [0054]
- **R. AUCKENTHALER**. Score Normalization for Text-Independent Speaker Verification Systems. *Digital Signal Processing*, 2000, vol. 10 (1-3), 42-54 [0055]