(12) UK Patent Application (19)GB (11)2493849 (13)A

(43)Date of A Publication 20.02.2013

(21) Application No: 1214706.2

(22) Date of Filing: 17.08.2012

(30) Priority Data:
(31) 13213609 (32) 19.08.2011 (33) US

(71) Applicant(s):
The Boeing Company
(Incorporated in USA - Illinois)
100 North Riverside Plaza, Chicago,
Illinois 60606-2016, United States of America

(72) Inventor(s):
George A Velius
David A Wheland

(74) Agent and/or Address for Service:
Carpmaels & Ransford
One Southampton Row, LONDON, WC1B 5HA,
United Kingdom

(51) INT CL:
G10L 17/02 (2013.01)

(56) Documents Cited:
US 20110311144 A1     US 20050060153 A1

(58) Field of Search:
INT CL G06F, G10L
Other: WPI EPODOC TXTE

(54) Title of the Invention: **Methods and systems for speaker identity verification**
Abstract Title: **A system for speaker identity verification**

(57) A system for confirming that a subject isthe source of spoken audio and the identity of the subject providing the spoken audio is described. The system includes at least one motion sensor operable to capture physical motion of at least one articulator (vocal chords, lips, tongue, cheek, jaw, mouth, teeth, trachea etc) that contributes to the production of speech, at least one acoustic signal sensor to receive acoustic signals, and a processing device comprising a memory and communicatively coupled to the at least one motion sensor and the at least one acoustic signal sensor. The processing device is programmed to correlate physical motion data with acoustical signal data to uniquely characterize the subject for purposes of verifying the subject is the source of the acoustical signal data and the identity of the subject.
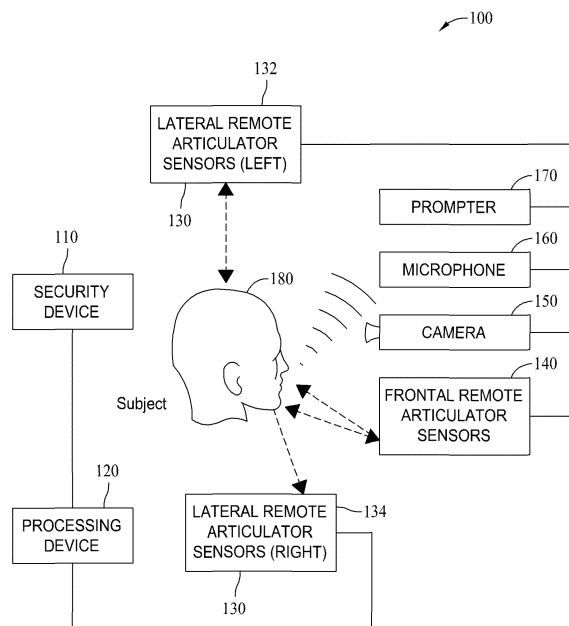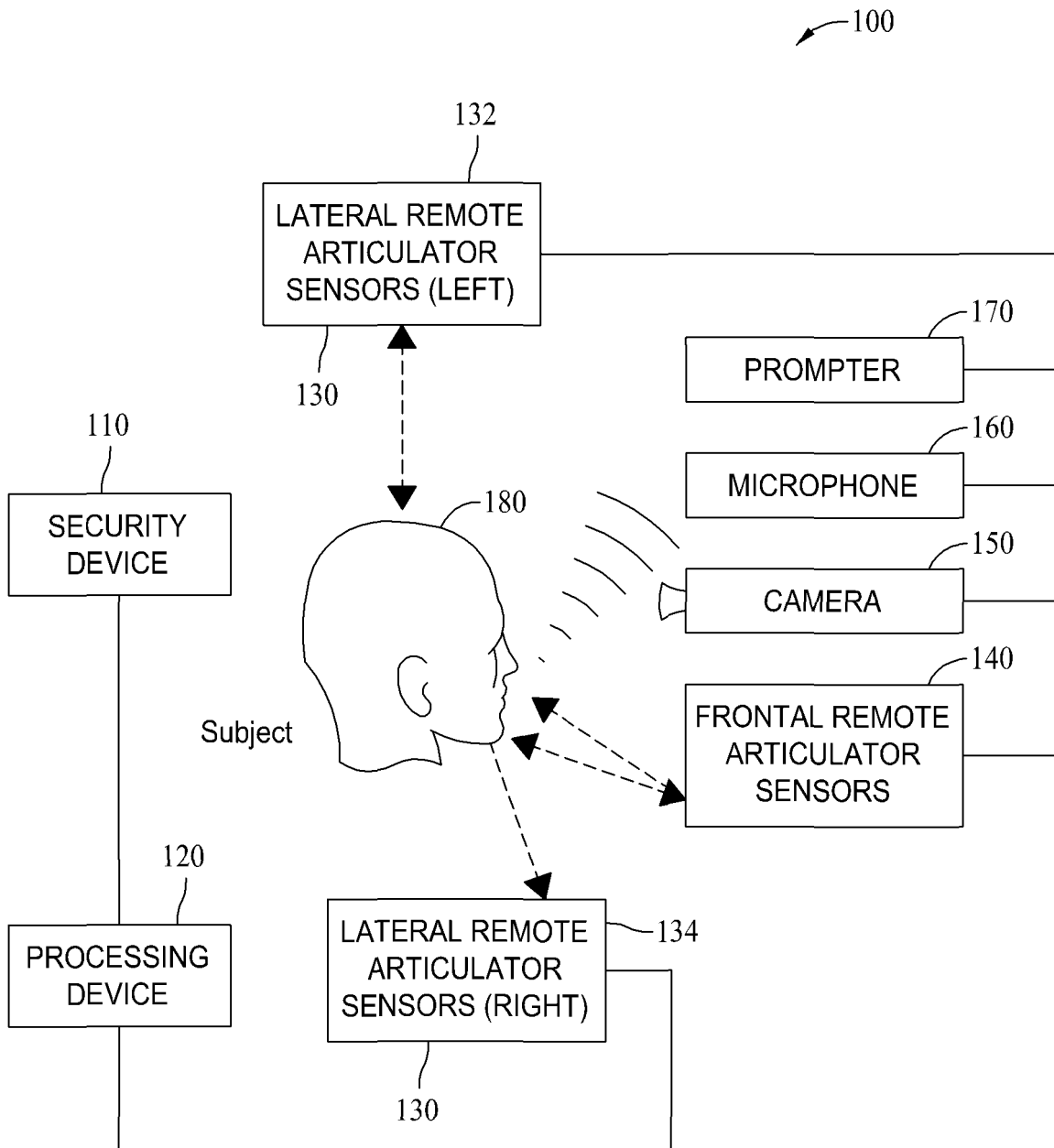


FIG. 1

100

132

LATERAL REMOTE
ARTICULATOR
SENSORS (LEFT)

170

PROMPTER

130

160

MICROPHONE

110

180

150

SECURITY
DEVICE

CAMERA

Subject

140

FRONTAL REMOTE
ARTICULATOR
SENSORS

120

134

PROCESSING
DEVICE

LATERAL REMOTE
ARTICULATOR
SENSORS (RIGHT)

130

FIG. 1

200

SUBJECT PRESENTS
SELF TO SYSTEM TO
MAKE IDENTITY CLAIM — 202

HEAD OF SUBJECT IS
LOCATED — 204

PROMPT SUBJECT
FOR SPEECH — 206

ACQUIRE SIGNALS
AS SUBJECT SPEAKS
INCLUDING PHYSICAL
ARTICULATOR MOTION
MEASUREMENTS — 208

ANALYZE SIGNALS TO
DETERMINE IF SPEECH
WAS RECORDED OR
PHYSICALLY SPOKEN — 210

ANALYZE SPEECH TO
DETERMINE IF
IDENTITY CLAIM IS
VALID — 212

FIG. 2

300

304  306  308

| PROCESSOR UNIT | MEMORY | PERSISTENT STORAGE |

302

| COMMUNICATIONS UNIT | INPUT/OUTPUT UNIT | DISPLAY |

310  312  314

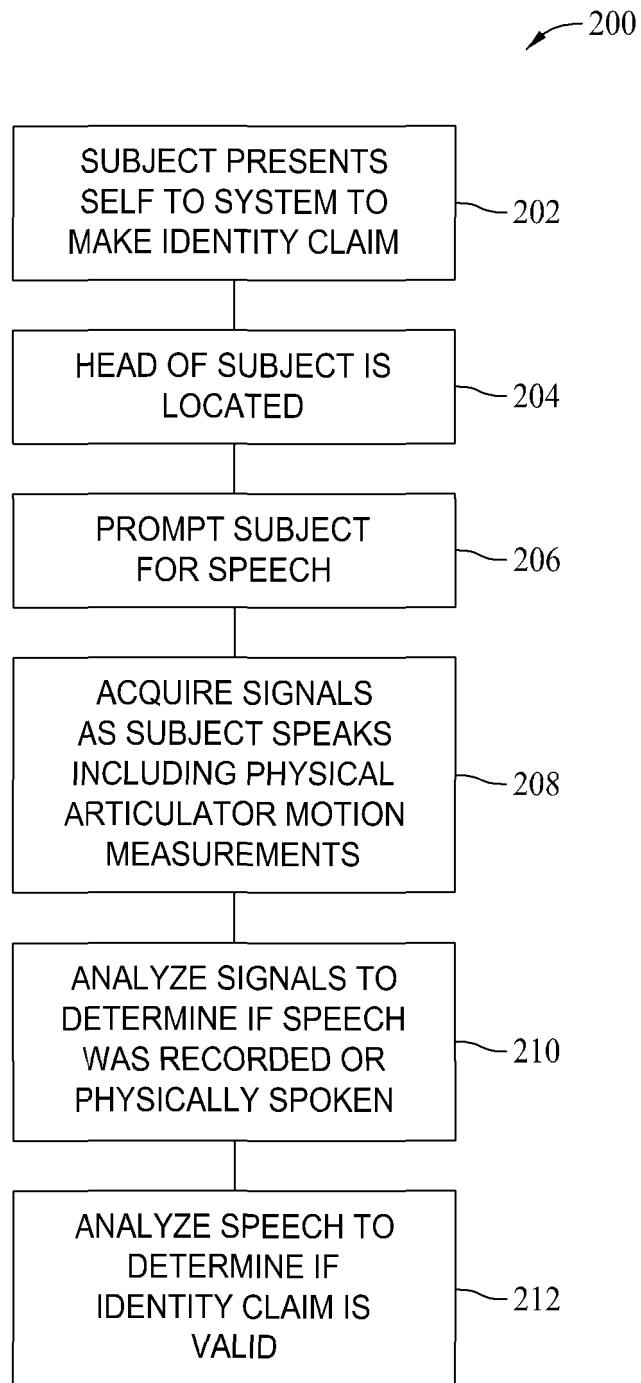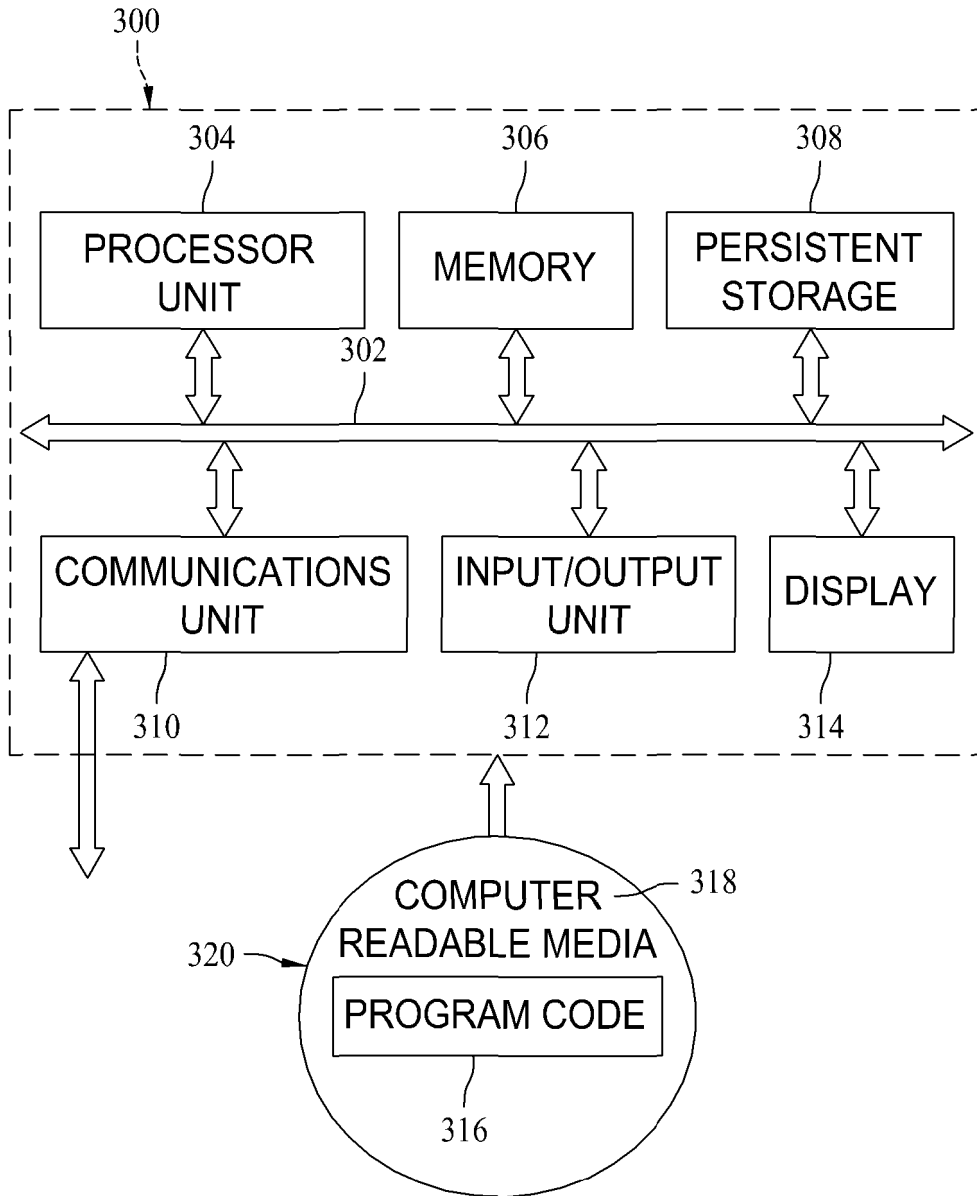COMPUTER READABLE MEDIA ⟋ 318

320

PROGRAM CODE

316

FIG. 3

# METHODS AND SYSTEMS FOR SPEAKER
# IDENTITY VERIFICATION

## BACKGROUND

[0001] The field of the disclosure relates generally to security systems that utilize voice recognition, and more specifically, to methods and systems for speaker identity verification.

[0002] Security is a major concern for many and draws significant investment from both military and commercial endeavors. Identity verification is a key capability for many aspects of security. Speech-based identity verification is one of the least objectionable biometrics. However, an open issue with this biometric, computerized Speaker Identity Verification (SIV) is the threat of replay attacks.

[0003] A speech signal is quite complex and every natural utterance is unique. By focusing on speaker-specific characteristics of the speech signal, individual speakers can be distinguished. In producing speech, the human body performs many complex and coordinated tasks. When the result of all these tasks is reduced to a one-dimensional audio signal, the result is easily captured and replayed with high fidelity.

[0004] Previous attempts to address the replay attack issue have focused on determining the 'liveness' of the source by referring to additional signal data, for example, a two dimensional video of the speaker's face. However, modern video capture and replay technology makes such liveness testing also susceptible to spoofing.

## BRIEF DESCRIPTION

[0005] In one aspect, a system for confirming that a subject is the source of spoken audio and the identity of the subject providing the spoken audio is provided. The system includes at least one motion sensor operable to capture physical motion of at least one articulator that contributes to the production of speech, at least

one acoustic signal sensor to receive acoustic signals, and a processing device comprising a memory and communicatively coupled to the at least one motion sensor and the at least one acoustic signal sensor. The processing device is programmed to correlate physical motion data with acoustical signal data to confirm that the subject is the source of the acoustical signal data for purposes of uniquely characterizing the subject and determining the identity of the subject.

[0006] In another aspect, a method for verifying that a subject is the source of spoken audio is provided. The method includes receiving physical vibration information and acoustic speech information for the spoken audio via a plurality of sensors, the information purportedly from the subject, correlating the physical vibration information with the acoustic speech information, analyzing the correlated information with stored physical vibration information and acoustic speech information associated with the subject, and determining whether the correlation is within a defined threshold for verifying that the subject is the source of spoken audio and identifying the subject.

[0007] In still another aspect, a security system is provided that includes a security device, at least one motion sensor operable to capture physical motion of at least one articulator that contributes to the production of speech, at least one acoustic signal sensor to receive acoustic signals, and a processing device having a memory. The processing device is communicatively coupled to the security device, the at least one motion sensor and the at least one acoustic signal sensor and is programmed to correlate physical motion data received from the at least one motion sensor with previously stored physical motion data and correlate acoustical signal data received from the at least one acoustic signal sensor with previously stored acoustical signal data to uniquely characterize the subject for purposes of verifying the subject is the source of the acoustical signal data and verifying the identity of the subject for providing access to the security device.

[0008] The features, functions, and advantages that have been discussed can be achieved independently in various embodiments or may be

combined in yet other embodiments further details of which can be seen with reference to the following description and drawings.

BRIEF DESCRIPTION OF THE DRAWINGS

[0009] Figure 1 is a block diagram of a system configured with sensors to establish a correlation between the motion of internal and external human articulators and a captured speech signal.

[0010] Figure 2 is a flow chart illustrating a process for verifying the identity of a subject.

[0011] Figure 3 is a diagram of a data processing system that might be found in the processing device of Figure 1.

DETAILED DESCRIPTION

[0012] Ultimately, to ensure that a speech signal is not artificially reproduced by an impostor, the physical source of the acoustic signal must be determined. The described embodiments thwart such replay attacks by confirming the presence of a functioning human vocal apparatus that produced the spoken utterances which then uniquely characterize the speaker. More specifically, modern remote sensing technology is used to monitor the human speech production system with sufficient precision to show conclusively that motion of human articulatory members is indeed the source of the acoustic speech signal used to verify an identity. As used herein, the human speech production system includes, but is not limited to, a plurality of vocal chords, cheeks, tongue, lips, trachea, etc.

[0013] As described above, the verification of an individual's claimed identity using their unique speech characteristics has been vulnerable to replay attacks. The described embodiments address such issues through measuring internal and external members of a speaker's vocal articulators and correlating physical motion with concurrently measured acoustic signals of the speaker's speech. Particularly, externally measurable vibrations, internal motion of the tongue/epiglottis/teeth/jaw, and/or internal/external vibration of vocal chords, cheek,

jaw, lips, and trachea are all potential sources of data that might be correlated with an acoustic signal, as further described below.

[0014]  In one embodiment, technical effects of the methods, systems, and computer-readable media described herein to confirm the identity of a source include at least one of: (a) sensor to capture physical vibration data, (b) an audio device to capture acoustic speech data, and (c) an analysis unit to correlate the vibration and speech data and determine whether the correlated data is within a defined threshold.  Further technical effects include (d) simultaneously receiving physical vibration and acoustic speech information from an audio source, (e) correlating the information based upon potential delays in receiving the information from the sources, and (f) determining whether the correlated data is within a defined threshold.

[0015]  Figure 1 is a block diagram of a computerized system 100 embodiment that uses remote sensing techniques to establish a sufficiently high correlation of the motion of both internal and external human articulators with a captured speech signal such that the physical source of the captured speech signal can be conclusively determined.  In one embodiment, system 100 operates to confirm the identity claim of an individual and works in cooperation with an authentication system, such as security device 110, and do so in a manner resistant to both live and mechanical impersonation.  In one example, security device 110 is an access control device, such as a door into a secure or classified work area.  In another example, security device 110 is associated with an airport security terminal or other airport security device.  Computerized system 100 shown in Figure 1 is but one embodiment of a system that can be utilized to perform the processes described herein.  Other system configurations may be utilized as well.

[0016]  Computerized system 100 includes a processing system 120 which receives data from lateral remote articulator sensors 130 (left 132, right 134), frontal remote articulator sensors 140, a camera 150, and a microphone 160.  A prompter 170 is controlled by processing system 120 and prompts the subject 180 with words or phrases to be spoken.  The speaking of the words and/or phrases is

captured by the sensors 130, 140, and one or more of camera 150 and microphone 160.

[0017] In one embodiment, sensors 130 and 140 include one or more micro-power pulse radars which are positioned to determine vocal chord closures and cheek oscillations. In short, the sensors 130, 140 combine to measure the actions of internal and external human articulators that contribute to the production of speech, for example, the vocal chords, tongue, lips, cheeks, and trachea. Each of these articulators vibrate or oscillate in some fashion during the production of speech, and such vibrations and oscillations can be captured and analyzed to ensure that the speech is not recorded and being generated by the subject 180. Correlation of the vibrations and oscillations allows for the modeling, by system 100, of physical positions of the articulators in time.

[0018] In one example, sensors 130 are laser devices that can capture movements of the cheeks and sensor 140 is a micro-power radar that captures movements of one or more of the tongue, lips, trachea and vocal chords. Processing device 120 ensures that captured vibrations and oscillations from the various sensors are in temporal alignment. To further illustrate, Figure 2 is a flowchart 200 that illustrates one process that might be performed by the system 100 of Figure 1. Initially, system 100 is presented 202 with an individual (i.e. subject 180) that wishes to make an identity claim. System 100 operates to locate 204 the head of the subject 180. In system 100 the head location function is performed, for example, by camera 150 along with processing device 120 to determine the position the head of the subject with respect to at least one motion sensor. In embodiments, camera 150 incorporates one or more of visual, infrared (IR), ultrasonic, or electromagnetic (EM) sensors. System 100, through operation of processing device 120 and prompter 170 prompts 206 subject 180 for speech. As the subject 180 speaks, system 100 acquires 208 signals including the subject's speech and physical articulator motion measurements, for example, cheek, throat, and vocal chord oscillations, and head, jaw, tongue, teeth, and lip position and movement. System 100 performs 210 an analysis to determine if the measured physical motions, potentially including oscillations, indicate, for example, a production mechanism for a recorded speech signal. If the performed 210

analysis indicates that an actual person is the source of the speech, a second analysis is performed to determine 212 if the subject's identity claim is valid. System 100 operates as programmed based on the performed 210 analysis and the determination 212.

[0019] The described embodiments utilize dynamic models of internal and external articulatory members, as well as speech-induced oscillations of external members, as the basis for liveness testing. Further, the described embodiments suggest that external articulator motion, such as lip and jaw movements, can be mimicked by an impostor while an alternate audio source, such as a hidden mechanical speaker, replays audio recordings of the individual authorized to access the security device 110. However, system 100 by measuring physical dynamics of internal articulators associated with the subject as well as by measuring external oscillations of articulatory members, counteracts any fabricated system used in an impersonation attempt. Specifically, any fabricated system would have to approach the complexity of the human vocal mechanism including vocal chords, tongue, cheeks, jaw, and mouth, which is highly unlikely.

[0020] Computerized system 100 uses remote sensing techniques to establish such a high correlation of the motion of both internal and external human articulators with a captured speech signal that no spoofing and replay attacks will succeed. As stated above, an objective of the described embodiments is to confirm the identity claim of an individual cooperating with an authentication system (e.g., security device 110).

[0021] With synchronous capture of both acoustic and motion data, a degree of temporal alignment between motion events and their corresponding acoustic signal events can be computed. Thresholds or decision rules for accepting the articulator production mechanism as the source of the acoustic signal must be empirically determined. In various embodiments, a degree of temporal alignment of motion events may also be computed between captured data and previously captured data. For example, system 100 may be programmed to verify that vibrations and oscillations captured by at least one motion sensor are in temporal alignment with

previously captured vibrations and oscillations stored in memory. In various embodiments, system 100 may be programmed to verify that speech information captured by an acoustic sensor is in temporal alignment with previously captured speech information stored in memory. In various embodiments, system 100 may verify both physical vibration information and acoustic speech information are in temporal alignment with previously stored physical vibration and acoustic speech information.

[0022] A number of techniques from different engineering disciplines are available as candidate solutions for measuring articulator motion. For example, for cheek oscillations, a micro-powered pulse radar is plausible. For head, jaw, and lip motion data, a real-time 3D camera 150 using light-phase and volumetric imaging techniques is used. For internal tissue motion of the tongue, glottis, and epiglottis, either radiofrequency-induced thermoacoustic tomography or low-power x-ray diffraction are used.

[0023] As used herein, an element or step recited in the singular and proceeded with the word "a" or "an" should be understood as not excluding plural elements or steps unless such exclusion is explicitly recited. Furthermore, references to "one embodiment" of the present invention or the "exemplary embodiment" are not intended to be interpreted as excluding the existence of additional embodiments that also incorporate the recited features.

[0024] Turning now to Figure 3, a diagram of processing system 100 is depicted in accordance with an illustrative embodiment. In this illustrative example, data processing system 300 includes communications fabric 302, which provides communications between processor unit 304, memory 306, persistent storage 308, communications unit 310, input/output (I/O) unit 312, and display 314. Communication unit 310 is utilized for communications with peripheral devices such as sensors 130, 140, camera 150, microphone 160 and prompter 170.

[0025] Processor unit 304 serves to execute instructions for software that may be loaded into memory 306. Processor unit 304 may be a set of one or more

processors or may be a multi-processor core, depending on the particular implementation. Further, processor unit 304 may be implemented using one or more heterogeneous processor systems in which a main processor is present with secondary processors on a single chip. As another illustrative example, processor unit 304 may be a symmetric multi-processor system containing multiple processors of the same type.

[0026] Memory 306 and persistent storage 308 are examples of storage devices. A storage device is any piece of hardware that is capable of storing information either on a temporary basis and/or a permanent basis. Memory 306, in these examples, may be, for example, without limitation, a random access memory or any other suitable volatile or non-volatile storage device. Persistent storage 308 may take various forms depending on the particular implementation. For example, without limitation, persistent storage 308 may contain one or more components or devices. For example, persistent storage 308 may be a hard drive, a flash memory, a rewritable optical disk, a rewritable magnetic tape, or some combination of the above. The media used by persistent storage 308 also may be removable. For example, without limitation, a removable hard drive may be used for persistent storage 308.

[0027] Communications unit 310, in these examples, provides for communications with other data processing systems or devices. In these examples, communications unit 310 is a network interface card. Communications unit 310 may provide communications through the use of either or both physical and wireless communication links.

[0028] Input/output unit 312 allows for input and output of data with other devices that may be connected to data processing system 300. For example, without limitation, input/output unit 312 may provide a connection for user input through a keyboard and mouse. Further, input/output unit 312 may send output to a printer. Display 314 provides a mechanism to display information to a user.

[0029] Instructions for the operating system and applications or programs are located on persistent storage 308. These instructions may be loaded into

memory 306 for execution by processor unit 304. The processes of the different embodiments may be performed by processor unit 304 using computer implemented instructions, which may be located in a memory, such as memory 306. These instructions are referred to as program code, computer usable program code, or computer readable program code that may be read and executed by a processor in processor unit 304. The program code in the different embodiments may be embodied on different physical or tangible computer readable media, such as memory 306 or persistent storage 308.

[0030] Program code 316 is located in a functional form on computer readable media 318 that is selectively removable and may be loaded onto or transferred to data processing system 300 for execution by processor unit 304. Program code 316 and computer readable media 318 form computer program product 320 in these examples. In one example, computer readable media 318 may be in a tangible form, such as, for example, an optical or magnetic disc that is inserted or placed into a drive or other device that is part of persistent storage 308 for transfer onto a storage device, such as a hard drive that is part of persistent storage 308. In a tangible form, computer readable media 318 also may take the form of a persistent storage, such as a hard drive, a thumb drive, or a flash memory that is connected to data processing system 300. The tangible form of computer readable media 318 is also referred to as computer recordable storage media. In some instances, computer readable media 318 may not be removable.

[0031] Alternatively, program code 316 may be transferred to data processing system 300 from computer readable media 318 through a communications link to communications unit 310 and/or through a connection to input/output unit 312. The communications link and/or the connection may be physical or wireless in the illustrative examples. The computer readable media also may take the form of non-tangible media, such as communications links or wireless transmissions containing the program code.

[0032] In some illustrative embodiments, program code 316 may be downloaded over a network to persistent storage 308 from another device or data

processing system for use within data processing system 300. For instance, program code stored in a computer readable storage medium in a server data processing system may be downloaded over a network from the server to data processing system 300. The data processing system providing program code 316 may be a server computer, a client computer, or some other device capable of storing and transmitting program code 316.

[0033] The different components illustrated for data processing system 300 are not meant to provide architectural limitations to the manner in which different embodiments may be implemented. The different illustrative embodiments may be implemented in a data processing system including components in addition to or in place of those illustrated for data processing system 300. Other components shown in Figure 3 can be varied from the illustrative examples shown.

[0034] As one example, a storage device in data processing system 300 is any hardware apparatus that may store data. Memory 306, persistent storage 308 and computer readable media 318 are examples of storage devices in a tangible form.

[0035] In another example, a bus system may be used to implement communications fabric 302 and may be comprised of one or more buses, such as a system bus or an input/output bus. Of course, the bus system may be implemented using any suitable type of architecture that provides for a transfer of data between different components or devices attached to the bus system. Additionally, a communications unit may include one or more devices used to transmit and receive data, such as a modem or a network adapter. Further, a memory may be, for example, without limitation, memory 306 or a cache such as that found in an interface and memory controller hub that may be present in communications fabric 302.

[0036] The above described embodiments are directed to methods and systems that capture data related to the motion of both internal and external human articulators using remote sensing techniques. The captured data, for example, cheek oscillations captured with a micro-powered pulse radar, is used to establish a

correlation between the motion of human articulators and a captured speech signal such that spoofing and replay attacks on a security device do not succeed.

[0037] The embodiments described herein sometimes refer to "previously stored" and/or "previously captured" data such as audio data and data relating to articulator movements. Previously captured/stored data may refer to either or both of sensor data from the same session or access-attempt and/or reference data from earlier sessions or access attempts. Previously captured/stored data may be stored in memory 306, persistent storage 308, or computer readable media 318 and may include data representative of physical motion of vocal articulators for a specific speaker for correlation with data received from one or more motion sensors. The correlation of vibration data is sometimes speaker-specific. For example, a first speaker of audio may produce vibration data that produces a correlation of 80%, while other speakers may produce a correlation of 90%. Advantageous system performance is achieved when such statistical reference data is available for selecting the threshold of acceptable correlation on a per-session basis. However, using sensor data from a same session can still provide an indicator of "liveness", for example, by showing that tongue position is consistent with cheek vibrations during voiced phonemes.

[0038] This written description uses examples to disclose various embodiments, which include the best mode, to enable any person skilled in the art to practice those embodiments, including making and using any devices or systems and performing any incorporated methods. The patentable scope is defined by the claims, and may include other examples that occur to those skilled in the art. Such other examples are intended to be within the scope of the claims if they have structural elements that do not differ from the literal language of the claims, or if they include equivalent structural elements with insubstantial differences from the literal languages of the claims.

WHAT IS CLAIMED IS:

1. A system for confirming that a subject is a source of spoken audio and the identity of the subject providing the spoken audio, said system comprising:

at least one motion sensor operable to capture physical motion of at least one articulator that contributes to the production of speech;

at least one acoustic signal sensor to receive acoustic signals; and

a processing device comprising a memory and communicatively coupled to said at least one motion sensor and said at least one acoustic signal sensor, said processing device programmed to correlate physical motion data with acoustical signal data to confirm that the subject is the source of the acoustical signal data for purposes of uniquely characterizing the subject and determining the identity of the subject.

2. The system according to Claim 1 wherein said at least one motion sensor comprises at least one micro-power pulse radar.

3. The system according to Claim 2 wherein said at least one micro-power pulse radar is positioned and operable to capture physical movement of at least one of vocal chords, lips, tongue, cheek, jaw, mouth, teeth, and trachea of the subject.

4. The system according to Claim 1 further comprising at least one secondary motion sensor operable to capture physical motion of at least one additional articulator that contributes to the production of speech.

5. The system according to Claim 4 wherein said at least one secondary motion sensor comprises a micro-power pulse radar operable for capturing vibrations of at least one cheek of the subject.

6.  The system according to Claim 1 further comprising at least one prompter communicatively coupled to said processing device operable to solicit a specific word or phrase from the subject.

7.  The system according to Claim 1 wherein said memory comprises data representative of physical motion data of vocal articulators for a specific subject for correlation with data received from said at least one motion sensor.

8.  The system according to Claim 1 wherein said at least one motion sensor comprises at least one of a lateral remote articulator sensor and a frontal remote articulator sensor operable for capturing at least one of vibrations and oscillations produced during the production of speech and configured for positioning with respect to a head of the subject.

9.  The system according to Claim 1 wherein said processing device is programmed to model, in time, physical positions of speech articulators provided by the subject.

10.  The system according to Claim 1 wherein said processing device is programmed to verify that vibrations and oscillations captured by said at least one motion sensor are in temporal alignment with previously captured vibrations and oscillations stored in said memory.

11.  The system according to Claim 1 further comprising a camera communicatively coupled to said processing device and operable, with said processing device, for accurately positioning a head of the subject with respect to said at least one motion sensor.

12.  The system according to Claim 1 wherein said at least one acoustic signal sensor comprises a microphone, said microphone and said processing device operable for capturing an audible speech pattern of the subject for comparison with a speech pattern stored in said memory.

13. The system according to Claim 1 wherein said at least one motion sensor comprises at least one of a radiofrequency-induced thermoacoustic tomography device and a low-power x-ray diffraction device.

14. A method for verifying that a subject is the source of spoken audio, said method comprising:

receiving physical vibration information and acoustic speech information for the spoken audio via a plurality of sensors, the information purportedly from the subject;

correlating the physical vibration information with the acoustic speech information;

analyzing the correlated information with stored physical vibration information and acoustic speech information associated with the subject; and

determining whether the correlation is within a defined threshold for verifying that the subject is the source of spoken audio and identifying the subject.

15. The method according to Claim 14 wherein receiving physical vibration information and acoustic speech information comprises receiving data from at least one micro-power pulse radar positioned to capture physical movement of at least one internal articulator for the subject that contributes to the production of speech.

16. The method according to Claim 14 wherein receiving physical vibration information and acoustic speech information comprises receiving data from at least one micro-power pulse radar positioned to capture vibrations and oscillation of at least one of vocal chords, lips, tongue, cheek, jaw, mouth, teeth, and trachea of the subject.

17. The method according to Claim 14 wherein correlating the physical vibration information and acoustic speech information comprises verifying that the received physical vibration information and acoustic speech information are

in temporal alignment with previously stored physical vibration information and acoustic speech information associated with the subject in a memory.

18.     The method according to Claim 14 wherein correlating the physical vibration information and acoustic speech information comprises comparing a speech pattern captured using a microphone with a speech pattern stored in said memory.

19.     The method according to Claim 14 further comprising using information captured with a camera to accurately position a head of the subject with respect to at least one motion sensor.

20.     A security system comprising:

a security device;

at least one motion sensor operable to capture physical motion of at least one articulator that contributes to the production of speech;

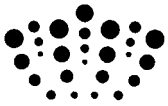at least one acoustic signal sensor to receive acoustic signals; and

a processing device comprising a memory and communicatively coupled to said security device, said at least one motion sensor and said at least one acoustic signal sensor, said processing device programmed to:

correlate physical motion data received from said at least one motion sensor with previously stored physical motion data; and

correlate acoustical signal data received from said at least one acoustic signal sensor with previously stored acoustical signal data to uniquely characterize the subject for purposes of verifying the subject is the source of the acoustical signal data and verifying the identity of the subject for providing access to said security device.

21. A system substantially as described herein, with reference to and as shown in figures 1 and/or 3.


22. A method substantially as described herein, with reference to figure 2.

**INTELLECTUAL**
PROPERTY OFFICE

| | | | |
|---|---|---|---|
| **Application No:** | GB1214706.2 | **Examiner:** | Mrs Hannah Sylvester |
| **Claims searched:** | 1-22 | **Date of search:** | 8 December 2012 |

## Patents Act 1977: Search Report under Section 17

### Documents considered to be relevant:

| Category | Relevant to claims | Identity of document and passage or figure of particular relevance |
|---|---|---|
| X | 1, 14 and 20 | US2005/060153 A1 (UNIV CALIFORNIA [US]) see claims 1 and 6 for quick reference and the whole document. |
| A | - | US2011/311144 A1 (MICROSOFT CORP [US]) |

### Categories:

| | | | |
|---|---|---|---|
| X | Document indicating lack of novelty or inventive step | A | Document indicating technological background and/or state of the art. |
| Y | Document indicating lack of inventive step if combined with one or more other documents of same category. | P | Document published on or after the declared priority date but before the filing date of this invention. |
| & | Member of the same patent family | E | Patent document published on or after, but with priority date earlier than, the filing date of this application. |

### Field of Search:

Search of GB, EP, WO & US patent documents classified in the following areas of the UKC$^X$ :

| |
|---|
| |

Worldwide search of patent documents classified in the following areas of the IPC

| |
|---|
| G06F; G10L |

The following online and other databases have been used in the preparation of this search report

| |
|---|
| WPI EPODOC TXTE |

### International Classification:

| Subclass | Subgroup | Valid From |
|---|---|---|
| None | | |