(19) World Intellectual Property Organization International Bureau





(43) International Publication Date 27 December 2002 (27.12.2002)

PCT

(10) International Publication Number WO 02/103680 A2

(51) International Patent Classification7: G10L 17/00

(21) International Application Number: PCT/GB02/02726

(22) International Filing Date: 13 June 2002 (13.06.2002)

(25) Filing Language: English

(26) Publication Language: **English**

(30) Priority Data:

0114866.7 19 June 2001 (19.06.2001) GB 60/302,501 2 July 2001 (02.07.2001)

(71) Applicant (for all designated States except US): SE-CURIVOX LTD [GB/GB]; 336 Perth Road, Dundee DD2 1EQ (GB).

(72) Inventor; and

(75) Inventor/Applicant (for US only): SAPELUK, Andrew, Thomas [GB/GB]; 4 Burn Street, Dundee DD3 0LA (GB).

MURGITROYD & COMPANY; Scotland House, 165-169 Scotland Street, Glasgow G5 8PL (GB).

- (81) Designated States (national): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NO, NZ, OM, PH, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VN, YU, ZA, ZM, ZW.
- (84) Designated States (regional): ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

Published:

without international search report and to be republished upon receipt of that report

For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.



(54) Title: "SPEAKER RECOGNITION SYSTEMS"

(57) Abstract: Speaker recognition (identification and/or verification) methods and systems, in which speech models for enrolled speakers consist of sets of feature vectors representing the smoothed frequency spectrum of each of a plurality of frames and a clustering algorithm is applied to the feature vectors of the frames to obtain a reduced data set representing the original speech sample, and wherein the adjacent frames are overlapped by at least 80 %. Speech models of this type model the static components of the speech sample and exhibit temporal independence. An identifier strategy is employed in which modelling and classification processes are selected to give a false rejection rate substantially equal to zero. Each enrolled speaker is associated with a cohort of a predetermined number of other enrolled speakers and a test sample is always matched with either the claimed identity or one of its associated cohort. This makes the overall error rate of the system dependent only on the false acceptance rate, which is determined by the cohort size. The false error rate is further reduced by use of multiple parallel modelling and/or classification processes. Speech models are normalised prior to classification using a normalisation model derived from either the test speech sample or one of the enrolled speaker samples (most preferably from the claimed identity enrolment sample).

"Speaker Recognition Systems" 1 2 3 The present invention relates to systems, methods 4 and apparatus for performing speaker recognition. 5 6 Speaker recognition encompasses the related fields 7 of speaker verification and speaker identification. The main objective is to confirm the claimed 8 identity of a speaker from his/her utterances, known 9 as verification, or to recognise the speaker from 10 11 his/her utterances, known as identification. Both 12 use a person's voice as a biometric measure and assume a unique relationship between the utterance 13 14 and the person producing the utterance. This unique relationship makes both verification and 15 identification possible. Speaker recognition 16 17 technology analyses a test utterance and compares it to a known template or model for the person being 18 recognised or verified. The effectiveness of the 19 system is dependent on the quality of the algorithms 20 21 used in the process.

2

Speaker recognition systems have many possible 1 2 applications. In accordance with a further aspect 3 of the present invention, speaker recognition 4 technology may be used to permanently mark an electronic document with a biometric print for every 5 6 person who views or edits the content. This produces an audit trail identifying all of the users and the 7 8 times of access and modification. As the user mark is biometric it is very difficult for the user to 9 10 dispute the authenticity of the mark. 11 Other biometric measures may provide the basis for 12 13 possible recognition systems, such as iris scanning, finger printing and facial features. These measures 14 15 all require additional hardware for recording whereas speaker recognition can be used with any 16 voice input such as over a telephone line or using a 17 18 standard multi-media personal computer with no modification. The techniques can be used in 19 conjunction with other security measures and other 20 21 biometrics for increased security. From the point of 22 view of a user the operation of the system is very 23 simple. 24 25 For example, when an on-line document is requested 26 the person requiring access will be asked to give a 27 sample of their speech. This will be a simple prompt from the client software 'please say this phrase..." 28 29 or something similar. The phrase uttered will then 30 be sent to a database server or to a speech recognition server, via any data network such as an 31 32 intranet, to be associated with the document and

3

stored as the key used to activate the document at 1 2 that particular time. A permanent record for a document can therefore be produced, over time, 3 providing an audit trail for the document. The 4 speaker authentication server may maintain a set of 5 6 templates (models) for all currently enrolled 7 persons and a historical record of previously 8 enrolled persons. 9 Speaker recognition systems rely on extracting some 10 11 unique features from a person's speech. turn depends on the manner in which human speech is 12 produced using the vocal tract and the nasal tract. 13 14 For practical purposes, the vocal tract and nasal tract can be regarded as two connected pipes, which 15 can resonate in a manner similar to a musical 16 17 instrument. The resonances produced depend on the diameter and length of the pipes. In the human 18 speech production mechanism, these diameters and to 19 20 some extent the length of the pipe sections can be modified by the articulators, typically the 21 positions of the tongue, the jaw, the lips and the 22 soft palate (velum). These resonances in the 23 spectrum are called the formant frequencies. 24 25 are normally around four formant frequencies in a 26 typical voice spectrum. 27 As with musical instruments, sound will only be 28 29 produced when a constriction of the airflow occurs 30 causing either vibration or turbulence. In human speech, the major vibrations occur when the 31 constriction occurs at the glottis (vocal cords). 32

4

When this happens, voiced speech is produced, 1 2 typically vowel-like sounds. When the constriction is in the mouth, caused by the tongue or teeth, a 3 turbulence is produced, (a hissing type of sound) 4 and the speech produced is called a fricative, 5 6 typified by "s", "sh", "th" etc. From an 7 engineering point of view, this is similar to a 8 source signal (the result of the constriction) being applied to a filter which has the general 9 characteristics (i.e. the same resonances) of the 10 vocal tract and the resulting output signal is the 11 speech sound. True speech is produced by 12 dynamically varying the positions of the 13 14 articulators. 15 All existing speaker recognition systems perform 16 17 similar computational steps. They operate by creating a template or model for an enrolled 18 speaker. The model is created by two main steps 19 applied to a speech sample, namely spectral analysis 20 and statistical analysis. Subsequent recognition of 21 an input speech sample is performed by modelling the 22 input sample (test utterance) in the same way as 23 during speaker enrolment, and pattern/classification 24 25 matching of the input model against a database of enrolled speakers. Existing systems vary in the 26 approach taken when performing some or all of these 27 28 steps. In conventional (industry standard) systems, the spectral analysis is either Linear Predictive 29 Coding (LPC)/Cepstral analysis ("LPCC") or FFT/sub-30 banding. This is followed by a statistical analysis 31 technique, usually a technique called Hidden Markov 32

5

1 Modelling (HMM), and the classification step is a 2 combination of a match against the claimed speaker 3 model and against an "impostor cohort" or "world model" (i.e. a set of other speaker models). 4 5 6 To allow efficient processing of speech samples, all 7 speaker recognition systems use time slices called 8 frames, where the utterance is split into frames and 9 each frame is processed in turn. Frames may or may 10 not be of equal size and may or may not overlap. An 11 example of a typical time signal representation of a 12 speech utterance divided into frames is illustrated 13 in Fig. 1 of the accompanying drawings. A generic 14 speaker recognition system is shown in block diagram form in Fig. 2, illustrating a test utterance being 15 16 processed, through an input filter 10, a spectral 17 analysis (LPCC) stage 12 and a statistical analysis 18 (HMM) stage 14, followed by score normalisation and speaker classification 16, by thresholding, 19 20 employing a database 18 of speaker models (enrolled 21 speaker data-set), before generating a decision as 22 to the identity of the speaker (identification) or the veracity of the speaker's claimed identity 23 24 (verification). 25 26 Such systems have a number of disadvantages or 27 limitations. Firstly, conventional spectral 28 analysis techniques produce a limited and incomplete feature set and therefore poor modelling. Secondly, 29 30 HMM techniques are "black-box" methods, which 31 combine good performance with relative ease of use, 32 but at the expense of transparency. The relative

6

importance of features extracted by the technique

PCT/GB02/02726

WO 02/103680

1

2 are not visible to the designer. Thirdly, the nature of the HMM models do not allow model-against-3 model comparisons to be made effectively. 4 5 Accordingly, important structural detail contained 6 within the enrolled speaker data-set cannot be 7 analysed and used effectively to improve system performance. Fourthly, HMM technology uses temporal 8 9 information to construct the model and is therefore 10 vulnerable to mimics, who impersonate others' voices by temporal variations in pitch etc. Fifthly, the 11 12 world model/impostor cohort employed by the system cannot easily be optimised for the purpose of 13 14 testing an utterance by a claimed speaker. 15 16 The performance of a speaker recognition system 17 relies on the fact that when a true speaker 18 utterance is tested against a model for that speaker it will produce a score, which is lower than a score 19 20 that is produced when an impostor utterance is 21 tested against the same model. This allows an 22 accept/reject threshold to be set. Consecutive 23 tests by the true speaker will not produce identical 24 scores. Rather, the scores will form a statistical distribution. However, the mean of the true speaker 25 distribution will be considerably lower than the 26 27 means of impostor distributions tested against the 28 same model. This is illustrated in Fig. 3, where 25 scores are plotted for each of eight speakers, 29 30 speaker 1 being the true speaker. It can be seen 31 from Fig. 3 that the scores of some speakers are

1	closer to the true speaker than others and can be
2	problematic.
3	
4	The present invention relates to improved speaker
5	recognition methods and systems which provide
6	improved performance in comparison with conventional
7	systems. In various aspects, the invention provides
8	improvements including but not limited to: improved
9	spectral analysis, transparency in its statistical
10	analysis, improved modelling, models that can be
11	compared allowing the data-set structure to be
12	analysed and used to improve system performance,
13	improved classification methods and the use of
14	statistically independent/partially independent
15	parallel processes to improve system performance.
16	
17	The invention further embraces computer programs for
18	implementing the methods and systems of the
19	invention, data carriers and storage media encoded
20	with such programs, data processing devices and
21	systems adapted to implement the methods and
22	systems, and data processing systems and devices
23	incorporating the methods and systems.
24	
25	The various aspects and preferred features of the
26	invention are defined in the Claims appended hereto.
27	
28	Embodiments of the invention will now be described,
29	by way of example only, with reference to the
30	accompanying drawings, in which:
31	

1 Fig. 1 is a time signal representation of an example

- of a speech utterance divided into frames;
- 3 Fig. 2 is a block diagram of a generic, prior art
- 4 speaker recognition system;
- 5 Fig. 3 is a plot of speaker recognition score
- 6 distributions for a number of speakers tested
- 7 against one of the speakers, obtained using a
- 8 conventional speaker recognition system;
- 9 Fig. 4 is a block diagram illustrating a first
- 10 embodiment of the present invention;
- 11 Fig. 5 is a block diagram illustrating a second
- 12 embodiment of the present invention;
- 13 Fig. 6 is a block diagram illustrating a third
- 14 embodiment of the present invention;
- Fig. 7 is a block diagram illustrating a further
- 16 embodiment of a speaker recognition system in
- 17 accordance with the present invention;
- 18 Fig. 8(a) is a time signal representation of an
- 19 example of a speech utterance divided into frames
- and Fig. 8(b) shows the corresponding frequency
- 21 spectrum and smoothed frequency spectrum of one
- 22 frame thereof;
- 23 Fig. 9 illustrates the differences between the
- 24 frequency spectra of two mis-aligned frames;
- 25 Fig. 10 shows the distribution of accumulated frame
- 26 scores plotted against their frequency of
- 27 occurrence;
- Fig. 11(a) shows the same accumulated score
- 29 distributions as Fig. 3 for comparison with Fig.
- 30 11(b), which shows corresponding accumulated score
- 31 distributions obtained using a speaker recognition
- 32 system in accordance with the present invention;

9

- 1 Fig. 12 illustrates the results of model against
- 2 model comparisons as compared with actual test
- scores, obtained using a system in accordance with
- 4 the present invention;
- 5 Fig. 13 illustrates the distribution of speaker
- 6 models used by a system in accordance with the
- 7 present invention in a two-dimensional
- 8 representation of a multi-dimensional dataspace;
- 9 Fig. 14 illustrates the use of an impostor cohort as
- 10 used in a system in accordance with the present
- 11 invention;
- 12 Fig. 15 is a block diagram illustrating a
- 13 normalisation process in accordance with one aspect
- of the present invention;
- 15 Fig. 16 is a block diagram illustrating an example
- of wide area user authentication system in
- 17 accordance with the present invention;
- 18 Fig. 17 is a block diagram illustrating the
- 19 corruption of a speech signal by various noise
- 20 sources and channel characteristics in the input
- 21 channel of a speaker recognition system;
- Figs. 18 and 19 illustrate the effects of noise and
- 23 channel characteristics on test utterances and
- enrolment models in a speaker recognition system;
- 25 and

29

- 26 Fig. 20 illustrates a channel normalisation method
- in accordance with one aspect of the present
- 28 invention.

The present invention includes a number of aspects

- 31 and features which may be combined in a variety of
- 32 ways in order to provide improved speaker

WO 02/103680

10

PCT/GB02/02726

recognition (verification and/or identification) 1 Certain aspects of the invention are 2 systems. concerned with the manner in which speech samples 3 are modelled during speaker enrolment and during 4 subsequent recognition of input speech samples. 5 Other aspects are concerned with the manner in which 6 7 input speech models are classified in order to reach a decision regarding the identity of the speaker. 8 further aspect is concerned with normalising speech 9 signals input to speaker recognition systems 10 11 (channel normalisation). Still further aspects concern applications of speaker recognition systems. 12 13 Referring now to the drawings, Figs. 4 to 6 14 illustrate the basic architectures used in systems 15 embodying various aspects of the invention. It will 16 17 be understood that the inputs to all of the embodiments of the invention described herein are 18 19 digital signals comprising speech samples which have previously been digitised by any suitable means (not 20 shown), and all of the filters and other modules 21 referred to are digital. 22 23 In Fig. 4, a speech sample is input to the system 24 via a channel normalisation module 200 and a filter 25 24. Instead of or in addition to this "front-end" 26 normalisation, channel normalisation may be 27 performed at a later stage of processing the speech 28 sample, as shall be discussed further below. 29 sample would be divided into a series of frames 30 prior to being input to the filter 24 or at some 31 other point prior to feature extraction. In some 32

embodiments, as discussed further below, a noise 1 signal 206 may be added to the filtered signal (or 2 3 could be added prior to the filter 24). The sample data are input to a modelling (feature extraction) 4 module 202, which includes a spectral analysis 5 6 module 26 and (at least in the case of speech sample 7 data being processed for enrolment purposes) a statistical analysis module 28. The model (feature 8 set) output from the modelling module 202 comprises 9 a set of coefficients representing the smoothed 10 frequency spectrum of the input speech sample. 11 During enrolment of a speaker, the model is added to 12 13 a database of enrolled speakers (not shown). During recognition of an input speech sample, the model 14 (feature set) is input to a classification module 15 110, which compares the model (feature set) with 16 models selected from the database of enrolled 17 speakers. On the basis of this comparison, a 18 decision is reached at 204 so as to identify the 19 speaker or to verify the claimed identity of the 20 speaker. The channel normalisation of the input 21 sample and the addition of the noise signal 206 22 comprise aspects of the invention, as shall be 23 described in more detail below, and are preferred 24 features of all implementations of the invention. 25 26 In some embodiments, channel normalisation may be applied following spectral analysis 26 or during the 27 classification process, rather than being applied to 28 29 the input speech sample prior to processing as shown in Figs. 4 to 6. Novel aspects of the modelling and 30 classification processes in accordance with other 31

1

12

aspects of the invention will also be described in 2 more detail below. 3 Other aspects of the invention involve various types 4 of parallelism in the processing of speech samples 5 for enrolment and/or recognition. 6 7 In Fig. 5, the basic operation of the system is the 8 same as in Fig. 4, except that the output from the 9 modelling module 202 is input to multiple, parallel 10 classification processes 110a, 110b ... 110n, and 11 the outputs from the multiple classification 12 processes are combined in order to reach a final 13 decision, as shall be described in more detail 14 below. In Fig. 6, the basic operation of the system 15 is also the same as in Fig. 4, except that the input 16 sample is processed by multiple, parallel modelling 17 processes 202a, 202b ... 202n (typically providing 18 slightly different feature extraction/modelling as 19 described further below), possibly via multiple 20 filters 24a, 24b ... 24n (in this case the noise 21 signal 206 is shown being added to the input signal 22 upstream of the filters 24a, 24b ... 24n), and the 23 outputs from the multiple modelling processes are 24 input to the classification module 110, as shall 25 also be described in more detail below. These types 26 of multiple parallel modelling processes are 27 preferably applied to both enrolment sample data and 28 test sample data. 29 30 Multiple parallel modelling processes may also be 31 combined with multiple parallel classification 32

13

processes; e.g. the input to each of the parallel 1 classification processes 110a-n in Fig. 5 could be 2 the output from multiple parallel modelling 3 processes as shown in Fig. 6. 4 5 6 Various aspects of the invention will now be 7 described in more detail by reference to the 8 modelling, classification and normalisation processes indicated in Figs. 4 to 6. 9 10 11 MODELLING 12 The spectral analysis modules 26, 26a-n may apply 13 similar spectral analysis methods to those used in 14 conventional speaker recognition systems. 15 Preferably, the spectral analysis applied by the 16 17 modules 26a-n is of a type that, for each frame of the sample data, extracts a set of feature vectors 18 19 (coefficients) representing the smoothed frequency 20 spectrum of the frame. This preferably comprises LPC/Cepstral (LPCC) modelling, producing an 21 increased feature set which models the finer detail 22 of the spectra, but may include variants such as 23 delta cepstral or emphasis/de-emphasis of selected 24 25 coefficients based on a weighting scheme. coefficients may alternatively be obtained by other 26 means such as or Fast Fourier Transform [FFT] or by 27 use of a filter bank. 28 29 The complete sample is represented by a matrix 30 consisting of one row of coefficients for each frame 31 of the sample. For the purposes of the preferred 32

14

embodiments of the present invention, these matrices 1 2 will each have a size of the order of 1000 (frames) x 24 (coefficients). In conventional systems, a 3 single first matrix of this type, representing the 4 complete original signal, would be subject to 5 6 statistical analysis such as HMM. 7 As will be understood by those skilled in the art, 8 the LP transform effectively produces a set of 9 filter coefficients representing the smoothed 10 frequency spectrum for each frame of the test 11 utterance. The LP filter coefficients are related 12 to Z-plane poles. The Cepstral transform has the 13 effect of compressing the dynamic range of the 14 smoothed spectrum, de-emphasising the LP poles by 15 moving them closer to the Z-plane origin (away from 16 the real frequency axis at $z=e^{jw}$). The Cepstral 17 transform uses a log function for this purpose. 18 will be understood that other similar or equivalent 19 techniques could be used in the spectral analysis of 20 the speech sample in order to obtain a smoothed 21 frequency spectrum and to de-emphasise the poles 22 thereof. This de-emphasis produces a set of 23 coefficients which when transformed back into the 24 time domain are less dynamic and more well balanced 25 (the cepstral coefficients are akin to a time signal 26 27 or impulse response of the LP filter with deemphasised poles). The log function also transforms 28 multiplicative processes into additive processes. 29 30 The model derived from the speech sample may be 31 regarded as a set of feature vectors based on the 32

PCT/GB02/02726

WO 02/103680

32

frequency content of the sample signal. When a 1 2 feature vector based on frequency content is extracted from a signal, the order of the vector is 3 important. If the order is too low then some 4 important information may not be modelled. To avoid 5 this, the order of the feature extractor (e.g. the 6 7 number of poles of an LP filter) may be selected to 8 be greater than the expected order. However, this in itself causes problems. Poles which match 9 resonances in the signal give good results, whilst 10 the other resulting coefficients of the feature 11 vector will model spurious aspects of the signal. 12 Accordingly, when this vector is compared with 13 14 another model or reference, the distance measure computed may be unduly influenced by the values of 15 those coefficients which are modelling spurious 16 17 aspects of the signal. The distance measure (score) which is returned will thus be inaccurate, possibly 18 19 giving a poor score for a frame which in reality is 20 a good match. 21 In accordance with one aspect of the invention, this 22 23 problem can be obviated or mitigated by adding a noise signal n(t) (206 in Figs. 4 - 6) having known 24 25 characteristics to the speech signal s(t) before the signal is input to the modelling process (i.e. the 26 input signal = s(t)+n(t)). The same noise signal 27 would be used during enrolment of speakers and in 28 subsequent use of the system. The addition of the 29 known noise signal has the effect of forcing the 30 "extra" coefficients (above the number actually 31

required) to model a known function and hence to

WO 02/103680

16

PCT/GB02/02726

give consistent results which are less problematic 1 during model/test vector comparison. 2 particularly relevant for suppressing the effect of 3 noise (channel noise and other noise) during 4 "silences" in the speech sample data. This problem 5 may also be addressed as a consequence of the use of 6 massively overlapping sample frames discussed below. 7 8 As previously mentioned, in order to allow efficient 9 processing of speech samples all speaker recognition 10 systems use time slices called frames, so that the 11 utterance is split into a sequence of frames and 12 each frame is processed in turn. The frames may or 13 may not be of equal size and they may overlap. 14 Models generated by speaker recognition systems thus 15 comprise a plurality of feature sets (vectors 16 corresponding to sets of coefficients) representing 17 a plurality of frames. When models are compared in 18 conventional speaker recognition systems it is 19 necessary to align corresponding frames of the 20 respective models. Different utterances of a given 21 phrase will never be exactly the same length, even 22 when spoken by the same person. Accordingly, a 23 difficulty exists in correctly aligning frames for 24 25 comparison. 26 Conventional systems convert the frames into a 27 spectral or smoothed spectral equivalent as shown in 28 Figs. 8(a) (showing a time signal divided into 29 frames) and 8(b) (showing the corresponding 30 frequency spectrum and smoothed frequency spectrum 31 of one of the frames of Fig. 8(a)). The systems 32

1 then perform further transformations and analysis 2 (such as Cepstral transformation, Vector 3 Quantisation, Hidden Markov Modelling (HMM) and 4 Dynamic Time Warping (DTW)) to obtain the desired result. Frame boundaries can be allocated in many 5 ways, but are usually measured from an arbitrary 6 7 starting point estimated to be the starting point of the useful speech signal. To compensate for this 8 arbitrary starting point, and also to compensate for 9 10 the natural variation in the length of similar sounds, techniques such as HMM and DTW are used when 11 12 comparing two or more utterances such as when building models or when comparing models with test 13 utterances. The HMM/DTW compensation is generally 14 done at a point in the system following spectral 15 analysis, using whatever coefficient set is used to 16 represent the content of a frame, and does not refer 17 18 to the original time signal. The alignment precision is thus limited to the size of a frame. 19 In addition, these techniques assume that the 20 alignment of a particular frame will be within a 21 fixed region of an utterance which is within a few 22 frames of where it is expected to lie. 23 24 introduces a temporal element to the system as the estimated alignment of the current frame depends on 25 the alignment of previous frames, and the alignment 26 27 of subsequent frames depends on the alignment of the present frame. In practice, this means that a 28 particular frame, such as a frame which exists 200 29 ms into an utterance, will in general only be 30 compared with other frames in the 200 ms region of 31

the model or of other utterances being used to

1 construct a model. This approach derives from

2 <u>speech</u> recognition methods (e.g. speech-to-text

3 conversion), where it is used to estimate a phonetic

18

PCT/GB02/02726

- 4 sequence from a series of frames. The present
- 5 applicants believe that this approach is
- 6 inappropriate for <u>speaker</u> recognition, for the
- 7 following reasons.

8

WO 02/103680

- 9 A. Most seriously, the conventional approach
- 10 provides only crude alignment of frames. The
- 11 arbitrary allocation of starting points means that
- it will generally not be possible to obtain accurate
- alignment of the starting points of two respective
- 14 frames, so that even two frames giving a "best
- match" may have significantly different spectral
- 16 characteristics, as illustrated in Fig. 9.

17

- 18 B. Secondly, the conventional approach relies on
- 19 the temporal sequence of the frames and bases
- 20 speaker verification on spectral characteristics
- 21 derived from temporally adjacent frames.

- 23 In accordance with a further aspect of the
- invention, the present enrolment modelling process
- involves the use of very large frame overlaps, akin
- to convolution, to avoid problems arising from frame
- 27 alignment between models (discussed at A. above) and
- 28 to improve the quality of the model obtained. This
- 29 technique is applied during speaker enrolment in
- order to obtain a model, preferably based on
- 31 repeated utterances of the enrolment phrase. By
- 32 massively overlapping the frames, the resulting

WO 02/103680

19

PCT/GB02/02726

1 model effectively approaches a model of all possible 2 alignments, with relatively small differences between adjacent frames, thereby providing good modelling of patterns. Preferably, the frame overlap 4 5 is selected to be at least 80%, more preferably it 6 is in the range 80% to 90%, and may be as high as 7 95%. 8 9 The frames are transformed into representative 10 coefficients using the LPCC transformation as 11 described above, so that each utterance employed in the reference model generated by the enrolment 12 13 process is represented by a matrix (typically having a size of the order of 1000 frames by 24 14 15 coefficients as previously described). There might 16 typically be ten such matrices representing ten 17 utterances. A clustering or averaging technique such as Vector Quantisation (described further 18 19 below) is then used to reduce the data to produce the reference model for the speaker. This model 20 21 does not depend on the temporal order of the frames, addressing the problems described at B. above. 22 23 24 Preferred embodiments of the present invention 25 combine the massive overlapping of frames described 26 above with Vector Quantisation or the like as described below. This provides a mode of operation 27 28 which is quite different from conventional HMM/DTW 29 systems. In such conventional systems, all frames 30 are considered equally valid and are used to derive 31 a final "score" for thresholding into a yes/no 32 decision, generally by accumulating scores derived

20

1 by comparing and matching individual frames. The 2 validity of the scores obtained is limited by the 3 accuracy of the frame alignments. 4 5 In accordance with this aspect of the present invention, the reference (enrolment) models 6 7 represent a large number of possible frame 8 alignments. Rather than matching individual frames 9 of a test utterance with individual frames of the 10 reference models and deriving scores for each 11 matched pair of frames, this allows all frames of 12 the test utterance to be compared and scored against 13 every frame of the reference model, giving a 14 statistical distribution of the frequency of occurrence of frame score values. "Good" frame 15 16 matches will yield low scores and "poor" frame matches will yield high scores (or the converse, 17 18 depending on the scoring scheme). A test utterance 19 frame tested against a large number of reference 20 models will result in a normal distribution as 21 illustrated in Fig. 10. Most frame scores will lie 22 close to the mean and within a few standard deviations therefrom. Because of the massive 23 24 overlapping of frames in the reference models, the score distributions will include "best matches" 25 26 between accurately aligned corresponding frames of 27 the test utterance and reference models. When a test utterance from a particular speaker is tested 28 29 against the reference model for that speaker, the 30 distribution will thus include a higher incidence of 31 very low scores. This ultimately results in "true speaker" scores being consistently low due to some 32

21

parts of the utterance being easily identified as 1 2 originating from the true speaker while other parts less obviously from the true speaker are classified 3 by being from the general population. Impostor 4 5 frame scores will not produce low scores and will be classified as being from the general population. 6 7 8 That is, in accordance with this aspect of the invention, the reference models comprise sets of 9 coefficients derived for a plurality of massively 10 11 overlapping frames, and a test utterance is tested by comparing all of the frames of the test utterance 12 with all of the frames of the relevant reference 13 14 models and analysing the distribution of frame scores obtained therefrom. 15 16 17 The massive overlapping of frames applied to speech 18 samples for enrolment purposes may also be applied 19 to input utterances during subsequent speaker 20 recognition, but this is not necessary. 21 The use of massive overlaps in the enrolment sample 22 23 data is also beneficial in dealing with problems arising from noise present in periods of silence in 24 25 the sample data. Such problems are particularly significant for text-independent speaker recognition 26 systems. The existence of silences may or may not 27 cause problems for an individual model or verification 28 attempt, but they will cause deterioration in the 29 overall system performance. The question is therefore 30 31 how do we remove this completely or minimise the adverse effect. The use of massive frame overlaps in 32

1 the present invention contains an inherent solution.

22

2 Consider the equations, which describe averaging the

3 frame spectra (discussed in more detail below),

4
$$\overline{s(\omega)} = \frac{1}{N} \sum_{n} s_n(\omega) = \frac{1}{N} \sum_{n} (ss_n(\omega) \times sd_n(\omega))$$

5

6 =

$$7 \qquad \frac{1}{N}(ss_1(\omega) \times sd_1(\omega)) + (ss_2(\omega) \times sd_2(\omega)) + \dots + (ss_N(\omega) \times sd_N(\omega))$$

8

9 =

$$ss(\omega) \times \frac{1}{N} (sd_1(\omega) + sd_2(\omega) +sd_2(\omega) + sd_N(\omega))$$

11 It can be seen that the static parts (ss) average to

 $ss(\omega)$ and that individual frames have the spectra

 $ss_n(\omega) \times sd_n(\omega)$. Consider however the spectra of two

14 added frames,

15 $(ss_1(\omega) \times sd_1(\omega)) + (ss_2(\omega) \times sd_2(\omega)) = ss(\omega) \times (sd_1(\omega) + sd_2(\omega))$

16 we have the steady part multiplied by a new spectra

 $sd_1(\omega) + sd_2(\omega)$. But since it is to be reduced by

averaging, and it is also dynamic or variable in

19 nature, the new spectra should behave in exactly the

20 same way as a randomly extracted frame. The

21 implication of this is that frames could be randomly

22 added together with minimal effect on performance.

23 This observation is not entirely true since we can

24 have the case of valid speech frames added to

25 silence frames in which the net result is a valid

speech frame. This in fact results in an improvement

in performance, as we are no longer including

28 unwanted silences in the modelling.

23

1 If a typical signal with some minor silence problems 2 has time frames randomly added, the silences would be eliminated but the signal would appear to have 3 undergone major corruption. However the present 4 5 invention using massively overlapped frames still 6 functions. Interestingly, the implication of this is 7 that channel echoes have no effect and can be It also underlines the fact that the 8 ignored. 9 preferred operating modes of the present invention 10 extract the static parts of the spectra to a larger extent than conventional verifiers (as discussed 11 12 further below). The addition of frames in this way has substantially the same effect as adding coloured 13 14 noise to prevent unwanted modelling as discussed 15 above. 16 17 In accordance with another aspect, the invention 18 uses clustering or averaging techniques such as Vector Quantisation applied by the modules 28, 28a-n 19 20 in a manner that differs from statistical analysis 21 techniques used in conventional speaker recognition 22 systems. 23 24 Preferably, the system of the present invention uses a Vector Quantisation (VQ) technique in processing 25 the enrolment sample data output from the spectral 26 27 analysis modules 26, 26a-n. This is a simplified 28 technique, compared with statistical analysis techniques such as HMM employed in many prior art 29 30 systems, resulting in transparent modelling 31 providing models in a form which allow model-32 against-model comparisons in the subsequent

1	classification stage. Also, VQ as deployed in the
2	present invention does not use temporal information,
3	making the system resistant to impostors.
4	
5	The VQ process effectively compresses the LPCC
6	output data by identifying clusters of data points,
7	determining average values for each cluster, and
8	discarding data which do not clearly belong to any
9	cluster. This results in a set of second matrices
LO	of second coefficients, representing the LPCC data
L1	of the set of first matrices, but of reduced size
L2	(typically, for example, 64 $ ext{x}$ 24 as compared with
L3	1000 x 24).
L4	
L 5	The effects of the use of LPCC spectral analysis and
L6	clustering/averaging in the present invention will
L7	now be discussed.
L8	
L9	The basic model assumes that spectral magnitude is
20	useful and that the phase may be disregarded. This
21	is known to apply to human hearing and if it was not
22	applied to a verifier the system would exhibit
23	undesirable phase related problems, such as
24	sensitivity to the distance of the microphone from
25	the speaker. Further assume that the spectral
26	information of a speech sample can be regarded as
27	consisting of two parts a static part $ss(\omega)$ and a
28	dynamic part $sd(\omega)$ and that the processes are
29	multiplicative. It is also assumed that the dynamic
30	part is significantly larger than the static part.
31	$s(\omega) = ss(\omega) \times sd(\omega)$

WO 02/103680

25

PCT/GB02/02726

As, by definition, the static part is fixed it is 2 the more useful as a biometric as it will be related to the static characteristics of the vocal tract. 3 4 This will relate the measure to some fixed physical characteristic as opposed to $sd(\omega)$ which is related 5 6 to the dynamics of the speech. 7 8 The complete extraction of $ss(\omega)$ would give a 9 biometric which exhibits the properties of a 10 physical biometric, i.e. cannot be changed at will and does not deteriorate over time. Alternatively 11 12 the exclusive use of $sd(\omega)$ will give a biometric 13 which exhibits the properties of a behavioral biometric, i.e. can be changed at will and will 14 deteriorate over time. A mixture of the two should 15 exhibit intermediate properties but as $sd(\omega)$ is much 16 larger than $ss(\omega)$ it is more likely that a 17 18 combination will exhibit the properties of $sd(\omega)$; 19 i.e. behavioral. 20 21 As with all frequency representations of a signal 22 the assumption is that the time signal exists from $-\infty$ to $+\infty$ which clearly is not physically possible. 23 24 In practice all spectral estimates of a signal will 25 be made using a window, which exists for a finite 26 period of time. The window can either be rectangular 27 or shaped by a function (such as a Hamming window). 28 29 The use of a rectangular window amounts to simply taking a section of a signal in the area of interest 30 31 and assuming that it is zero elsewhere. This

1

26

technique is common in speech processing in which

2 the sections of signal are called frames; Fig. 1 shows a time signal with the frames indicated. 3 4 5 The frames can be shaped using an alternate window. Interestingly, the major effect of windowing is a 6 7 spreading of the characteristic of a particular frequency to its neighbours, a kind of spectral 8 averaging. This effect is caused by the main lobe; 9 in addition to this the side lobes produce spectral 10 oscillations, which are periodic in the spectrum. 11 The present system later extracts the all-pole 12 Linear Prediction coefficients, which have the 13 intended effect of spectral smoothing and the extra 14 smoothing, caused by the windowing, is not seen as a 15 major issue. However, the periodic side lobe 16 effects might be troublesome if the window size was 17 inadvertently changed. This however can be avoided 18 by good housekeeping. 19 20 Given that we can divide the time signal into frames 21 the spectral characteristics for frames 1 to N can 22 be represented as 23 $s_1(\omega) = ss_1(\omega) \times sd_1(\omega)$; $s_2(\omega) = ss_2(\omega) \times sd_2(\omega)$; 24 25 $s_n(\omega) = ss_n(\omega) \times sd_n(\omega)$; 26 $s_{N}(\omega) = ss_{N}(\omega) \times sd_{N}(\omega)$ 27 28 But by definition 29 $ss(\omega) = ss_1(\omega) = ss_2(\omega) = ss_3(\omega)$ 30 $= ss_{N}(\omega)$ 31

27

On first impressions to extract $ss(\omega)$ would seem to

2 be possible using an averaging process,

$$\overline{s(\omega)} = \frac{1}{N} \sum_{n} s_n(\omega) = \frac{1}{N} \sum_{n} (ss_n(\omega) \times sd_n(\omega))$$

4 =

$$5 \qquad \frac{1}{N}(ss_1(\omega) \times sd_1(\omega)) + (ss_2(\omega) \times sd_2(\omega)) + \dots (ss_N(\omega) \times sd_N(\omega))$$

6 =

$$7 \qquad \operatorname{ss}(\omega) \times \frac{1}{N} (\operatorname{sd}_1(\omega) + \operatorname{sd}_2(\omega) + \dots \cdot \operatorname{sd}_2(\omega) \dots + \operatorname{sd}_N(\omega)) = \operatorname{ss}(\omega) \times \operatorname{U}(\omega)$$

8 where,

9
$$U(\omega) = \frac{1}{N} (sd_1(\omega) + sd_2(\omega) +sd_2(\omega)..... + sd_N(\omega))$$

10 If the frames had independent spectral

11 characteristics (each resulting from random process)

then $U(\omega)$ would tend to white noise, i.e. would have

a flat spectrum so that $\overline{s(\omega)}$ could be extracted by

14 smoothing the spectrum. This would most likely be

the case if N were very large $\rightarrow \infty$. Given the linear

16 nature of the time domain - frequency domain - time

17 domain transformations a similar analysis could have

18 been described in the time domain.

19

22

20 For real world conditions it cannot be assumed that

N would be large in the sense that the frames have

independent spectral characteristics. It is

23 important to remember that this would require N to

24 be large under two conditions:

1. During model creation

26 2. During a verification event

27

28

1 Failure to comply during either would potentially

- 2 cause a system failure (error), however a failure in
- 3 1 is the more serious as it would remain a potential
- 4 source of error until updated, whereas a problem in
- 5 2 is a single instance event.

6

- 7 If $U(\omega)$ cannot be guaranteed to converge to white
- 8 noise, what can be done to cope with the situation?
- 9 First consider that:
- 1. $U(\omega)$ will be a variable quantity
- 2. When smoothed across the frequency spectrum it
- would ideally be flat; i.e. the smoothed
- 13 version $Usm(\omega) = 1$
- 14 3. $U(\omega)$ is the truncated sum of the speech frames
- the number of which would ideally tend to
- 16 infinity.

17

18 Considering the equation

$$\overline{s(\omega)} = ss(\omega) \times \frac{1}{N} \sum_{n} sd_{n}(\omega)$$

- The summation part tending to a flat spectrum is not
- 21 an ideal performance measure, if we return to the
- 22 frame based equivalent:

$$\overline{s(\omega)} = \frac{1}{N} \sum_{n} (ss_{n}(\omega) \times sd_{n}(\omega))$$

24 If we take the logarithms of the frames:

$$\frac{1}{N} \sum_{n} \log((ss_n(\omega) \times sd_n(\omega)))$$

$$= \frac{1}{N} \sum_{n} \left[\log(ss_n(\omega)) + \log(sd_n(\omega)) \right] = \log(ss(\omega)) + \frac{1}{N} \sum_{n} \log(sd_n(\omega))$$

29

1

 $= lss(\omega) + lsd(\omega)$

3 it can be seen that the relationship between the

4 static and dynamic parts is now additive. Because

5 the relationship between the time domain and the

6 frequency domain is linear a transformation from

7 frequency to time gives:

8
$$lss(\omega) + lsd(\omega) \rightarrow cs(\tau) + cd(\tau) = c(\tau)$$

9

In signal processing $c(\tau)$ is known as the Cepstral

11 transformation of s(t) as discussed previously.

12

13 In general cepstral analysis consists of

time domain \rightarrow frequency domain \rightarrow log(spectrum) \rightarrow time domain

The Cepstral transformation has been used in speech

16 analysis in many forms.

17

18 As discussed above, in our current usage we create

19 the Cepstral coefficients for the frames and extract

20 the static part,

21
$$\frac{1}{N}\sum_{n}c_{n}(\tau) = \frac{1}{N}\sum_{n}\left(cs_{n}(\tau) + cd_{n}(\tau)\right) = cs(\tau) + \frac{1}{N}\sum_{n}cd_{n}(\tau)$$

22

23 Ideally the length of the speech signal would be

long enough so that the dynamic part was completely

25 random and the mean would tend to zero. This would

leave the static part cs(t) as our biometric

27 measure. However, we have a number of problems to

28 overcome.

29

30 1. How do we handle the imperfect nature of the

31 sum-to-zero

PCT/GB02/02726

30

WO 02/103680

1 2. channel variation 2 3. endpointing 4. additive noise 3 4 Referring to the imperfect nature of the sum-to-5 6 zero, the nature of the Cepstral coefficients are 7 such that they decay with increasing time and have the appearance of an impulse response for stable 8 9 systems. This means that the dynamic range of each coefficient is different and they are in general in 10 descending order. 11 12 It can be shown that the differences between the 13 average coefficients of a test sample and the frame 14 15 coefficient values for the true speaker model and the frame coefficient values of an impostor model 16 are not large and a simple summation over all of the 17 utterance frames to produce a distance score will be 18 difficult to threshold in the conventional manner. 19 20 21 If we consider the two difficult problems associated with this methodology together rather than 22 separately the answer to the problem is revealed. To 23 re-emphasise, the two points of difficulty are, 24 1. the utterances will never be long enough for 25 the mean of the dynamic part to converge to 26 27 zero 2. the differences between the true speaker and 28 29 the impostors will be small and difficult to threshold. 30

31

32 Consider two speakers with models based upon

4

6

8

10

13

15

$$1 \qquad \overline{c(\tau)} = \frac{1}{N} \sum_{n} c_n(\tau) = \frac{1}{N} \sum_{n} (cs_n(\tau) + cd_n(\tau)) = cs(\tau) + \frac{1}{N} \sum_{n} cd_n(\tau)$$

so that the models are $m1(\tau)$ and $m2(\tau)$, where,

$$3 \qquad ml(\tau) = \overline{cl(\tau)} = \frac{1}{N} \sum_{n} cl_{n}(\tau) = \frac{1}{N} \sum_{n} (csl_{n}(\tau) + cdl_{n}(\tau)) = csl(\tau) + \frac{1}{N} \sum_{n} cdl_{n}(\tau)$$

5 = $csl(\tau) + el(\tau)$; where $el(\tau)$ is the error

7 In vector form the models are

 $9 m1 = \begin{bmatrix} csl_1 + el_1 \\ csl_2 + el_2 \\ \bullet \\ csl_p + el_p \end{bmatrix} and m2 = \begin{bmatrix} cs2_1 + e2_1 \\ cs2_2 + e2_2 \\ \bullet \\ cs2_p + e2_p \end{bmatrix}$

11 A test utterance from speaker 1 expressed in the

12 same form will be

14 $T1 = \begin{bmatrix} csl_1 + Tel_1 \\ csl_2 + Tel_2 \\ \bullet \\ csl_p + Tel_p \end{bmatrix}$

16 using a simple distance measure true speaker

17 distance is

18
$$dl = |ml - Tl| = \begin{bmatrix} csl_1 + el_1 \\ csl_2 + el_2 \\ \bullet \\ csl_p + el_p \end{bmatrix} - \begin{bmatrix} csl_1 + Tel_1 \\ csl_2 + Tel_2 \\ \bullet \\ csl_p + Tel_p \end{bmatrix} = |el - Tel|$$

19 impostor distance is

32

PCT/GB02/02726

1

WO 02/103680

3

Assuming that the convergence of the dynamic parts 4 5 of the models is good (i.e. that the error vectors are small compared to the static vectors) then in 6 general d1<d2. This is simply stating that the 7 8 models built represent the enrolled speaker (a condition that can easily checked during enrolment 9 using the data available at that time). 10 Interestingly, if e1 and e2 are small compared to 11 the test signal error Tel the distances become 12 independent of e1 and e2. The condition under which 13 14 the test error will be large when compared to the model error is during text-independent test 15 conditions. This shows that if the dynamic 16 17 components of the enrolment speech samples are minimised in the enrolment models then such models 18

21

19 20

The errors e1 and e2 above are average model
construction errors; the actual errors are on a
frame by frame basis and will have a distribution
about the mean. This distribution could be modelled
in a number of ways the simplest being by use of a
standard clustering technique such as k-means to
model the distribution. The use of k-means

can provide a good basis for text-independent

speaker recognition

33

1 clustering is also known in other forms as Vector

- 2 Quantisation (VQ) and is a major part of the Self
- 3 Organising Map (SOM) also known as the Kohonen

4 Artificial Neural Network.

5

- 6 The system just described where a test utterance is
- 7 applied to two models and the closest chosen is a
- 8 variant of identification. In the above case if
- 9 either speaker 1 or speaker 2, the enrolled
- 10 speakers, claim to be themselves and are tested they
- 11 will always test as true and so the False Rejection
- 12 Rate FRR =0. If an unknown speaker claims to be
- either speaker1 or speaker2 he will be classified as
- one or the other, so there is a 1/2 chance of
- 15 success and hence a False Acceptance Rate FAR =50%.
- 16 If an equal number of true speaker tests and random
- impostor tests were carried out, we can calculate an
- overall error rate as (FRR+FAR)/2 = (0+0.5)/2=25%

19

- 20 It is obvious that the number of models (the cohort)
- 21 against which the test utterance is tested will have
- 22 an effect on the FAR and it will reduce as the
- cohort increases. It can be shown that the accuracy
- of recognition under these conditions is asymptotic
- 25 to 100% with increasing cohort size, since FRR=0,
- 26 but as the accuracy is

27

28 accuracy =
$$100 - (FRR + FAR)\frac{100}{2} = 100 - (FRR + \frac{1}{cohort_size})\frac{100}{2}$$

29

it is in more general terms asymptotic to 100-FRR.

WO 02/103680

It is worth observing at this point that the FRR and 2 FAR are largely decoupled: the FRR is fixed by the quality of the model produced and the FAR is fixed 3 4 by the cohort size. It is also worth observing that 5 to halve the error rate we need to double the cohort 6 size e.g. for 99% accuracy the cohort is 50, for 7 99.5% accuracy the cohort is 100, for 99.75% accuracy the cohort is 200. As the cohort increases 8 9 the computational load increases and in fact doubles 10 for each halving of the error rate. As the cohort 11 increases to very large numbers the decoupling of the FRR and FAR will break down and the FRR will 12 13 begin to increase. 14 Rather than continually increasing the cohort size 15 16 in an attempt to reduce the FAR to a minimum another 17 approach is needed. The approach, in accordance with one aspect of the invention, is to use parallel 18 19 processes (also discussed elsewhere in the present description), which exhibit slightly different 20 21 impostor characteristics and are thus partially statistically independent with respect to the 22 23 identifier strategy. The idea is to take a core 24 identifier which exhibits the zero or approximately 25 zero FRR and which has a FAR that is set by the cohort size. The front end processing of this core 26 identifier is then modified slightly to reorder the 27 distances of the cohort member models from the true 28 speaker model. This is done while maintaining the 29 30 FRR~0 and can be achieved by altering the spectral 31 shaping filters 24a-24n (see Fig. 7), or by altering

34

PCT/GB02/02726

35

PCT/GB02/02726

WO 02/103680

the transformed coefficients, such as by using 1 2 delta-ceps etc. 3 When an enrolled speaker uses the system the test 4 signal is applied to all of the processes in 5 parallel but each process has a FRR~0 and the 6 7 speaker will pass. When an unknown impostor uses the system he will pass each individual process with a 8 probability of 1/cohort size. However with the 9 parallel processes we have introduced conditional 10 probabilities. That is, if an impostor passes 11 process1 what is the likelihood of him passing the 12 modified process2 as well etc. Although the 13 probability of an impostor passing all of the 14 processes is not that of the statistically 15 independent case of 16 statistically independent result = process $prob^{no_of_processes}$ 17 it does however reduce with the addition of 18 processes. It can be shown that for a given process 19 FAR value, the overall accuracy of the system 20 increases with the number of processes. 21 22 Where multiple parallel processes are used in this 23 way, the scheme for matching a test sample against a 24 claimed identity may require a successful match for 25 each process or may require a predetermined 26 proportion of successful matches. 27 28 The combined use of massive sample frame overlaps 29 and Vector Quantisation (or equivalent) in building 30 enrolment models in accordance with the present 31 invention provides particular advantages. 32

36

massive overlapping is applied at the time of 1 constructing the models, although it could also be 2 applied at the time of testing an utterance. The 3 technique involves using a massive frame overlap, 4 typically 80-90%, to generate a large possible 5 number of alignments; the frames generated by the 6 7 alignments are then transformed into representative coefficients using the LPCC transformation to 8 produce a matrix of coefficients representing all of 9 the alignments. This avoids conventional problems of 10 11 frame alignment. The matrix is typically of the size no of frames by LPCC order, for example 1000x24. 12 This is repeated for all of the utterances used in 13 constructing the model, typically 10, Giving 10 14 matrices of 1000x24. Vector Quantisation is then 15 used to reduce the data to produce a model for the 16 17 speaker. This has the effect of averaging the frames so as to reduce the significance of the dynamic 18 19 components of the sampled speech data as discussed The resulting model does not take cognisance 20 above. of the frame position in the test utterance and is 21 hence not temporal in nature. This addresses the 22 23 problem of temporal dependency. 24 The combined use of VQ and massive frame overlapping 25 produces an operation mode which is different from 26 conventional systems based upon HMM/DTW. In HMM/DTW 27 all frames are considered to be equally valid and 28 are used to form a final score for thresholding into 29 a yes/no decision. In the present invention every 30 31 row (frame) of the test sample data is tested

against every row of the enrolment model data for

Τ	the claimed speaker and the associated impostor
2	cohort. For each row of the test sample data, a
3	best match can be found with one row of the
4	enrolment model, yielding a test score for the test
5	sample against each of the relevant enrolment
6	models. The test sample is matched to the enrolment
7	model that gives the best score. If the match is
8	with the claimed identity, the test speaker is
9	accepted. If the match is with an impostor the test
10	speaker is rejected.
11	
12	The present system, then, uses LPCC and VQ modelling
13	(or similar/equivalent spectral analysis and
14	clustering techniques), together with massive
15	overlapping of the sample frames, to produce the
16	reference models for each enrolled speaker, which
17	are stored in the database. In use of the system, ar
18	input test utterance is subjected to similar
19	spectral analysis to obtain an input test model
20	which can be tested against the enrolled speaker
21	data-set. Advantageously, this approach can be
22	applied so as to obtain a very low False Rejection
23	Rate (FRR), substantially equal to zero. The
24	significance of this is discussed further below.
25	
26	Parallel Modelling
27	
28	As previously discussed, the performance of speaker
29	recognition systems in accordance with the invention
30	can be improved by using multiple parallel processes
31	to generate the model.
32	

38

PCT/GB02/02726

Referring now to Fig. 7 of the drawings, one 1 preferred embodiment of a speaker recognition system 2 3 employing parallel modelling processes in accordance with one aspect of the invention comprises an input 4 channel 100 for inputting a signal representing a 5 6 speech sample to the system, a channel normalisation process 200 as described elsewhere, a plurality of 7 parallel signal processing channels 102a, 102b ... 8 102n, a classification module 110 and an output 9 channel 112. The system further includes an 10 enrolled speaker data-set 114; i.e. a database of 11 speech models obtained from speakers enrolled to use 12 13 the system. The speech sample data is processed in parallel by each of the processing channels 102a-n, 14 the outputs from each of the processing channels is 15 input to the classification module 110, which 16 communicates with the database 114 of enrolled 17 speaker data, and a decision as to the identity of 18 the source of the test utterance is output via the 19 20 output channel 112. 21 Each of the processing channels 102a-n comprises, in 22 series, a spectral shaping filter 24a-n, an 23 (optional) added noise input 206a-n, as described 24 elsewhere, a spectral analysis module 26a-n and a 25 26 statistical analysis module 28a-n. The outputs from each of the statistical analysis modules 28a-n is 27 input to the classification module 110. 28 29 The spectral shaping filters 24a-n comprise a bank 30 of filters which together divide the utterance 31 signal into a plurality of overlapping frequency 32

39

1 bands, each of which is then processed in parallel 2 by the subsequent modules 26a-n and 28a-n. 3 number of processing channels, and hence the number of frequency bands, may vary, with more channels 4 5 providing more detail in the subsequent analysis of 6 the input data. Preferably, at least two channels 7 are employed, more preferably at least four 8 channels. The filters 24a-n preferably constitute a low-pass or band-pass or high-pass filter bank. 9 10 bandwidth of the base filter 24a is selected such that the False Rejection Rate (FRR) resulting from 11 12 subsequent analysis of the output from the first 13 channel 102a is zero or as close as possible to 14 zero. The subsequent filters 24b-n have 15 incrementally increasing bandwidths that 16 incrementally pass more of the signal from the input 17 channel 100. The FRR for the output from each 18 channel 102a-n is thus maintained close to zero whilst the different channel outputs have slightly 19 20 different False Acceptance (FA) characteristics. Analysis of the combined outputs from the channels 21 22 102a-n yields a reduced overall FA rate (a claimed 23 identity is only accepted if the outputs from all of 24 the channels are accepted) with a FRR close to zero. The significance of this is discussed further below. 25 26 The use of multiple frequency bands improves upon 27 28 conventional single-channel spectral analysis, 29 increasing the size of the feature vectors of 30 interest in the subsequent statistical analysis.

40

It will be understood that different types of 1 2 parallel processing may be employed in the modelling process in order to provide multiple feature sets 3 modelling different (related or unrelated) aspects 4 of the input speech sample and/or alternative models 5 of similar aspects. Banks of filters of other types 6 7 in addition to or instead of low pass filters might be employed. Different types or variants of 8 spectral and/or statistical analysis techniques 9 might be used in parallel processing channels. 10 Parallel statistical analyses may involve applying 11 different weighting values to sets of feature 12 coefficients so as to obtain a set of slightly 13 deviated models. 14 15 It will be understood that the architecture 16 illustrated in Fig. 7 may be used for both obtaining 17 enrolment models for storing in the database 114 and 18 19 for processing test speech samples for testing against the enrolment models. Each enrolment model 20 may include data-sets for each of a plurality of 21 enrolment utterances. For each enrolment utterance, 22 there will be a matrix of data representing the 23 output of each of the parallel modelling processes. 24 Each of these matrices represents the 25 clustered/averaged spectral feature vectors. 26 sample data is subject to the same parallel spectral 27 28 analysis processes, but without clustering/averaging, so that the test model data 29 comprises a matrix representing the spectral 30 analysis data for each of the parallel modelling 31 processes. When a test model is tested against an 32

PCT/GB02/02726

1

41

enrolment model, the test matrix representing a

PCT/GB02/02726

2 particular modelling process is tested against enrolment matrices generated by the same modelling 3 4 process. 5 6 CLASSIFICATION 7 8 The nature of the reference models obtained by the modelling techniques described above is such that 9 they lend themselves to direct model against model 10 comparisons. This enables the system to employ an 11 identifier strategy in which each enrolment model is 12 13 associated with an impostor cohort. That is, for the reference model of each enrolled speaker 14 ("subject"), there is an impostor cohort comprising 15 a predetermined number of reference models of other 16 enrolled speakers, specific to that subject and 17 which has a known and predictable relationship to 18 the subject's reference model. These predictable 19 20 relationships enable the performance of the system to be improved. Fig. 11(a) shows the results 21 obtained by a conventional speaker recognition 22 system, similar to Fig. 3, comparing scores for an 23 input utterance tested against reference data for 24 25 eight speakers. Speaker 1 is the true speaker, but the scores for some of the other speakers are 26 sufficiently close to reduce significantly the 27 degree of confidence that the system has identified 28 the correct speaker. Fig. 11(b) shows equivalent 29 results obtained using a system in accordance with 30 the present invention. It can be seen that the 31 results for speaker 1 are much more clearly 32

1 distinguished from the results of all of the other 2 speakers 2 to 8. 3 4 The speaker modelling method employed in the 5 preferred embodiments of the present invention is 6 inherently simpler (and, in strict mathematical 7 terms, cruder) than conventional techniques such as HMM and possible alternatives such as gaussian 8 9 mixture models. However, the present applicants 10 believe that the conventional use of "tight" 11 statistical methods is inherently flawed and result 12 in poor "real world" performance, and that, 13 surprisingly, the relatively simpler statistical 14 methods of the present invention are much more 15 effective in practice. As previously noted, the 16 temporal nature of HMM makes it susceptible to 17 mimics, a problem which is avoided by the present 18 invention. Further, the models of the present 19 invention are ideally suited to enable analysis of 20 the structure of the enrolled speaker data-set by 21 model against model testing. 22 23 The ability to perform model against model 24 comparisons by using the present speaker models provides two particular advantages. Firstly, this 25 26 provides the ability to identify the most relevant 27 impostors in the enrolled speaker data-set (i.e. those which are close to and uniformly distributed 28 around a particular model) and to produce an 29 30 effective and predictable speaker normalisation 31 mechanism. VQ modelling involves choosing the size 32 of the model; i.e. choosing the number of

42

PCT/GB02/02726

43

coefficients ("centres"). Once this has been done, 1 2 the positions of the centres can be moved around until they give the best fit to all of the enrolment 3 4 data vectors. This effectively means allocating a centre to a cluster of enrolment vectors, so each 5 centre in the model represents a cluster of 6 7 information important to the speaker identity. 8 9 The model against model tests make it possible to 10 predict how an enrolled speaker, or claimed identity, will perform against the database both in 11 the broad sense and in an area local (in the system 12 13 dataspace) to the claimed identity. Fig. 12 illustrates the results of testing reference models 14 for speakers 2 to 8 against the reference model for 15 16 speaker 1. The ellipses show the model against model results whilst the stars show actual scores 17 for speaker utterances tested against model 1. 18 19 can be seen that the model against model tests can be used to predict the actual performance of a 20 particular speaker against a particular reference 21 22 model. The model against model results tend to lie at the bottom of the actual score distributions and 23 therefore indicate how well a particular impostor 24 will perform against model 1. This basic approach 25 of using model against model tests to predict actual 26 performance is known as such. As described further 27 below, this approach may be extended in accordance 28 with one aspect of the present invention to guard 29 30 particular models against impostors using individually selected, statistically variable 31 32 groupings.

PCT/GB02/02726

1							
2	The second advantage derived from model against						
3	model testing is the ability to predict the						
4	performance of a test utterance against some or, if						
5	need be, all of the enrolled speaker models. This						
6	enables a virtually unlimited number of test						
7	patterns to be used to confirm an identity, which is						
8	not possible with conventional systems.						
9							
10	In addition, the model against model test results						
11	may be used to assemble a specific impostor cohort						
12	for use with each reference model. This allows						
13	accurate score normalisation and also allows each						
14	model to be effectively "guarded" against impostors						
15	by using a statistically variable grouping which is						
16	selected for each enrolled speaker. This is						
17	illustrated by Fig. 13. Each reference model can be						
18	regarded as a point in a multi-dimensional						
19	dataspace, so that "distances" between models can be						
20	calculated. Fig. 13 illustrates this idea in two						
21	dimensions for clarity, where each star represents a						
22	model and the two-dimensional distance represents						
23	the distance between models.						
24							
25	It can be seen that the distribution of speaker						
26	models is not uniform, so that a world-model based						
27	normalisation technique will not operate equally						
28	well for all speaker models. It can also be seen						
29	that some speaker models can be relatively close to						
30	one another, which implies that there is potential						
31	for impostors to successfully impersonate enrolled						
32	speakers. For each speaker model, these issues can						

be resolved by creating a specific cohort of

1

45

2 impostors around the subject model. This simplifies normalisation and creates a guard against impostors. 3 This is illustrated in Fig. 14, which shows, in a 4 similar manner to Fig. 13, a subject model 5 6 represented by a circle, members of an impostor cohort represented by stars, and a score for an 7 impostor claiming to be the subject, represented by 8 an "x". The impostor score is sufficiently close to 9 10 the subject model to cause recognition problems. However, because the speaker data-set enables 11 prediction of how the true subject speaker will 12 perform against the models of the impostor cohort, 13 this information can be used to distinguish the 14 impostor x from the true subject, by testing the 15 16 impostor against the models of the cohort members as well as against the true subject model. That is, it 17 can be seen that the impostor utterance x is closer 18 19 to some of the cohort members than would be expected for the true subject, and further away from others 20 than expected. This would indicate an impostor 21 22 event and result in the impostor utterance being rejected as a match for the true subject. 23 24 25 This provides the basis for a two stage recognition process which firstly rejects impostors who are 26 27 clearly not the claimed speaker followed, where necessary, by a more detailed process applied to 28 utterances which are close enough to possibly be the 29 30 claimed speaker. 31

46

1 In certain applications of speaker verification 2 systems, it is important to minimise the possibility 3 of "false rejections"; i.e. instances in which the 4 identity claimed by a user is incorrectly rejected 5 as being false. In accordance with one aspect of 6 the invention, an "identifier strategy" is employed 7 which provides very low false rejections, whilst also providing predictable system performance and 8 9 minimising problems associated with the use of 10 thresholds in accepting or rejecting a claimed 11 identity. 12 13 In accordance with this strategy, the database of enrolled speakers (the "speaker space") is 14 15 partitioned; e.g. so that each speaker enrolled in 16 the system is assigned to a cohort comprising a fixed number N of enrolled speakers, as described 17 The speaker classification module of the 18 above. system (e.g. the module 110 in the system of Fig. 4) 19 operates such that the input test utterance is 20 21 compared with all of the members of the cohort 22 associated with the identity claimed by the speaker, and the test utterance is classified as 23 24 corresponding to that member of the cohort which 25 provides the best match. That is, the test 26 utterance is always matched to one member of the cohort, and will never be deemed not to match any 27 member of the cohort. If the cohort member to which 28 29 the utterance is matched corresponds to the claimed 30 identity, then the claimed identity is accepted as true. If the utterance is matched to any other 31

47

1 member of the cohort then the claimed identity is 2 rejected as false. 3 4 The modelling and classification processes can be tuned such that the proportion of false rejections 5 is effectively zero (FR = 0%) (as discussed above); 6 7 i.e. there is substantially zero probability that a speaker will be wrongly identified as a member of 8 the cohort other than the claimed identity. This is 9 10 facilitated by the use of model against model comparisons such that a match is not based simply 11 upon the test utterance being matched against the 12 13 single closest model, but also on the basis of its relationship to other members of the cohort. Where 14 the cohort is of a fixed size N, the maximum 15 possible proportion of false acceptances 16 FA = 100/N % and the total average error rate 17 = (FA + FR)/2 = 50/N %. If the cohort size N is 20, 18 the error rate is thus 2.5 %; i.e. an accuracy of 19 97.5 %. If the cohort size is fixed, the system is 20 scalable to any size of population while maintaining 21 a fixed and predictable error rate. That is, the 22 accuracy of the system is based on the size of the 23 cohort and is independent of the size of the general 24 population, making the system scalable to very large 25 populations. Accuracy can be improved by increasing 26 the cohort size, as long as the false rejection rate 27 28 does not increase significantly. 29 30 This strategy does not rely on the use of thresholds to determine a result, but thresholds could still be 31 used to reduce false acceptances; i.e. once a test 32

48

1 utterance has been matched to the claimed identity 2 using the foregoing strategy, thresholds could be applied to determine whether the match is close 4 enough to be finally accepted. 5 6 As indicated above, the selection of an impostor 7 cohort associated with a particular enrolment model 8 may involve the use of algorithms so that the 9 members of the impostor cohort have a particular 10 relationship with the enrolment model in question. 11 In principle, this may provide a degree of optimisation in the classification process. 12 13 However, it has been found that a randomly selected 14 impostor cohort performs equally well for most 15 practical purposes. The most important point is 16 that the cohort size should be predetermined in 17 order to give predictable performance. The impostor 18 cohort for a particular enrolment model may be selected at the time of enrolment or at the time of 19 20 testing a test utterance. 21 22 Parallel Classification 23 24 The performance of a speaker recognition system in 25 accordance with the invention may be improved by the use of multiple parallel classification processes. 26 27 Generally speaking, such processes will be 28 statistically independent or partially independent. This approach will provide multiple classification 29 30 results which can be combined to derive a final 31 result, as illustrated in Fig. 5.

1	In one example, using the identifier strategy
2	described above, the same test utterance may be
3	tested against a number of different cohorts, or
4	against different enrolment phrases, or combinations
5	thereof. Where multiple cohorts are employed, each
6	cohort will give a result with a false rejection
7	rate of essentially zero (FR = 0%) and a false
8	acceptance rate $FA = 100/N$ % as before. The overall
9	false acceptance rate for n cohorts of equal size
10	will be
11	$FA = 100*M/N^n$ % and the average error rate
12	= $50*M/N^n$ %, where M is a coefficient having a value
13	greater than 1 and representing the effect of the
14	processes not being entirely statistically
15	independent. That is, with 2 cohorts and a cohort
16	size of 20, the average error rate will be 0.125*M %
17	as compared with 2.5 % for a single cohort as
18	described above. Thresholds may also be applied to
19	further improve accuracy as previously described.
20	
21	Other types of partially statistically independent
22	processes may be employed in the modelling process,
23	the classification process or both as previously
24	discussed. Besides the examples previously given, a
25	single utterance may be divided into parts and
26	processed separately.
27	
28	NORMALISATION
29	
30	A further problem encountered with conventional
31	speaker recognition systems is that system
32	performance may be affected by differences between

50

PCT/GB02/02726

speech sampling systems used for initial enrolment 1 2 and subsequent recognition. Such differences arise from different transducers (microphones), soundcards 3 4 In accordance with a further aspect of the 5 present invention, these difficulties can be obviated or mitigated by normalising speech samples 6 on the basis of a normalisation characteristic which 7 is obtained and stored for each sampling system (or, 8 possibly, each type of sampling system) used to 9 10 input speech samples to the recognition system. Alternatively (preferably), the normalisation 11 characteristic can be estimated "on the fly" when a 12 13 speech sample is being input to the system. normalisation characteristic(s) can then be applied 14 to all input speech samples, so that reference 15 16 models and test scores are independent of the characteristics of particular sampling systems. 17 Alternatively or additionally, in accordance with a 18 19 further aspect of the invention a normalisation process can be applied at the time of testing test 20 sample data against enrolment sample data. 21 22 A normalisation characteristic is effectively a 23 transfer function of the sampling system and can be 24 25 derived, for example, by inputting a known reference signal to the sampling system, and processing the 26 sampled reference signal through the speech 27 recognition system. The resulting output from the 28 recognition system can then be stored and used to 29 30 normalise speech samples subsequently input through the same sampling system or the same type of 31 32 sampling system.

51

1	
2	Alternatively, as illustrated in Fig. 15, a speech
3	signal S(f) which has been modified by the transfer
4	function C(f) of an input channel 300 can be
5	normalised on the fly by inputting the modified
6	speech signal $S(f)*C(f)$ to an estimating module 302,
7	which estimates the transfer function C(f) of the
8	channel 300, and to a normalisation module 304, and
9	applying the inverse of the estimated transfer
10	function $1/C(f)$ to the normalisation module, so that
11	the output from the normalisation module closely
12	approximates the input signal S(f). The estimator
13	module 302 creates a digital filter with the
14	spectral characteristics of the channel 300 and the
15	inverse of this filter is used to normalise the
16	signal. For example, the inverse filter can be
17	calculated by determining the all-pole filter which
18	represents the spectral quality of a sample frame.
19	The filter coefficients are then smoothed over the
20	frames to remove as much of the signal as possible,
21	leaving the spectrum of the channel $(C(f))$. The
22	estimate of the channel spectrum is then used to
23	produce the inverse filter 1/C(f). This basic
24	approach can be enhanced to smooth the positions of
25	the poles of the filters obtained for the frames,
26	with intelligent cancellation of the poles to remove
27	those which are known not to be concerned with the
28	channel characteristics.
29	
30	Depending on the nature of the transfer

31 function/normalisation characteristic, the

normalisation process can be applied to the speech

<u>.</u>	sample prior to processing by the speaker
2	recognition system or to the spectral data or to the
3	model generated by the system.
4	
5	A preferred method of channel normalisation, in
6	accordance with one aspect of the invention, is
7	applied to the test model data and the relevant
8	enrolment models at the time of testing the test
9	sample against the enrolment models.
10	
11	The overall effect of the channel characteristics on
12	a speech signal could be described as
13	$\hat{s}(\omega) = ss(\omega) \times sd(\omega) \times cc(\omega)$
14	where $\hat{s}(\omega)$ is the estimate of the speakers
15	characteristics, $cc(\omega)$ is the channel characteristic
16	or changed channel characteristic as appropriate,
17	and the speech signal is treated as comprising a
18	static part and a dynamic part as before. Ideally
19	the unwanted channel characteristic can be estimated
20	and removed. In practice the removal can be achieved
21	in the time domain, frequency domain or a
22	combination. They both achieve the same effect, that
23	is to estimate $cc(\omega)$ and remove it using some form of
24	inverse filter or spectral division. If $\hat{c}c(\omega)$ is the
25	estimate of the spectrum of the unwanted channel
26	then we would calculate
27	$\frac{\hat{s}(\omega)}{\hat{c}c(\omega)} = ss(\omega) \times sd(\omega) \times \frac{cc(\omega)}{\hat{c}c(\omega)} \approx s(\omega)$

53

1 If the estimation of the channel characteristic is

2 good $\frac{cc(\omega)}{\hat{c}c(\omega)} \approx 1$ and our estimate of the speech is

3 good with the unwanted spectral shaping removed.

4 This would normally be implemented using a algorithm

5 based on the FFT.

6

7 An alternative implementation is to model the

8 channel characteristic as a filter, most likely in

9 the all-pole form,

10

11
$$h(z) = \frac{z^{N}}{z^{N} + a_{N-1}z^{N-1} + \dots + a_{0}}$$

12

13 This is the most basic form of the ARMA and would

normally be extracted from the time signal directly,

15 possibly using Linear Prediction.

16

17 A similar normalisation could be carried out on the

18 Cepstral representation.

19

20 In the Cepstral domain the speech signal is

21 represented as

$$c(\tau) = cs(\tau) + cd(\tau)$$

and the speech signal modified by the unwanted

24 channel characteristics is

$$\hat{c}(\tau) = cs(\tau) + cd(\tau) + cc(\tau)$$

26 It can be seen that in this case we have an additive

27 process rather than a product. But it should also be

28 remembered that both cs and cc are static and we may

need to remove one cc without removing the other.

54

T	it is important to consider the context in which we						
2	would wish to remove the signal cc and their						
3	different conditions (enrolled model, database						
4	derived cohort, test speaker etc.).						
5							
6	Figure 16 illustrates various sources of corruption						
7	of a speech sample in a speaker recognition system.						
8	The input speech signal s(t) is altered by						
9	environmental background noise, b(t), the recording						
10	device bandwidth, r(t), electrical noise and channel						
11	crosstalk, t(t), and transmission channel bandwidth,						
12	c(t), so that the signal input to the recognition						
13	system is an altered signal $v(t)$. The system is						
14	easier to analyse in the frequency domain and the						
15	signal at the verifier is:						
16	$v(\omega) = ((s(\omega) + b(\omega)).r(\omega) + t(\omega)).c(\omega)$ eq1						
17							
18	At the verifier we can define two conditions, when						
19	the person is speaking and when he is not. Resulting						
20	in two equations,						
21	$\mathbf{v}(\omega) = ((\mathbf{s}(\omega) + \mathbf{b}(\omega)).\mathbf{r}(\omega) + \mathbf{t}(\omega)).\mathbf{c}(\omega)$						
22	and						
23	$\mathbf{v}(\omega) = ((0 + \mathbf{b}(\omega)).\mathbf{r}(\omega) + \mathbf{t}(\omega)).\mathbf{c}(\omega)$						
24							
25	First consider the simplified problem as it applies						
26	to the systems in accordance with the present						
27	invention; assume that b(t)=t(t)=0						
28	$v(\omega) = s(\omega).r(\omega).c(\omega) = s(\omega).h(\omega)$						
29	where h() is the combined channel spectral						
30	characteristic,						

 $h(\omega) = r(\omega).c(\omega)$

55

1 $v(\omega) = s(\omega).h(\omega) = ss(\omega).sd(\omega).h(\omega)$ 2 3 The cohort models are selected from the database of speakers recorded using the same channel (b) and the 4 5 true speaker model is recorded using a different channel (a). The test speaker can either be the true 6 7 speaker or an impostor and will be recorded using a 8 third channel (c). Figure 17 shows this 9 diagrammatically. Fig. 18 shows the same thing expressed in the alternate form using the Cepstral 10 coefficients. It should be remembered that the 11 values of the signal components as represented in 12 Figs 17 and 18 are averages corresponding to the 13 summations of sample frame data. 14 15 Consider the claimed identity model, which was built 1,6 17 from, $v_1(\tau) = cs_1(\tau) + cd_1(\tau) + h_2(\tau)$ 18 eq2 and the cohort models which were built from, 19 $v_m(\tau) = cs_m(\tau) + cd_m(\tau) + h_h(\tau)$ 20 eq3 21 The problem for the verifier is that there are two 22 23 different channels used in the identifier and if we assume the difference between them is 24 $hd(\tau) = h_a(\tau) - h_b(\tau)$ 25 $h_a(\tau) = h_h(\tau) + hd(\tau)$ 26 or 27 28 then the claimed identity model referred to the cohort channel (b) will be 29 $v_1(\tau) = cs_1(\tau) + cd_1(\tau) + h_2(\tau) = cs_1(\tau) + cd_1(\tau) + h_2(\tau) + hd(\tau)$ 30 $v_1(\tau) = (cs_1(\tau) + hd(\tau)) + cd_1(\tau) + h_h(\tau)$ 31 and

56

1 2 it can be seen that the mean of the static part of 3 the claimed identity model has been shifted by the difference between the channels and will cause an 4 5 error if the true speaker is tested using channel-b 6 if the situation is not corrected. Similar problems 7 involving false acceptances using channel-a will 8 also occur. 9 10 One method of addressing this problem is to remove 11 the mean from the claimed identity model, but a 12 simple removal of the mean would at first glance 13 produce, 14 $v_1(\tau) = cd_1(\tau)$ where the static part of the speaker model has also 15 been removed. However, examining equation 1 (the 16 17 system model including additive noise) 18 $v(\omega) = ((s(\omega) + b(\omega)).r(\omega) + t(\omega)).c(\omega)$ if we consider the case during which the speaker 19 20 pauses, $s(\omega) = 0$ $v(\omega) = (b(\omega).r(\omega) + t(\omega)).c(\omega)$ 21 then 22 and $v(\omega) = n(\omega).c(\omega)$ 23 where $n(\omega)$ is a noise signal. 24 25 In cepstral form this would be 26 $v(\tau) = n(\tau) + c(\tau) = sn(\tau) + dn(\tau) + c(\tau)$ where as before sn is the static part of the noise 27 28 and dn is the result of the summation of the dynamic

29

30

part.

57

PCT/GB02/02726

WO 02/103680

29

30 31

1 The average of a model constructed from this would 2 be 3 $\operatorname{sn}(\tau) + \operatorname{c}(\tau)$ 4 where sn is any steady state noise such as an 5 interference tone and c is the channel. 6 7 Considering again equation1 (the claimed identity model build conditions) 8 $v_1(\tau) = cs_1(\tau) + cd_1(\tau) + h_2(\tau)$ 9 10 this was the noise free case, adding a steady state noise gives, 11 $v_1(\tau) = cs_1(\tau) + cd_1(\tau) + h_a(\tau) + sn(\tau)$ 12 13 If we constructed the speaker pause model for this 14 case we would get $\operatorname{sn}(\tau) + h_{a}(\tau)$ 15 16 and using this to remove the mean results in 17 $\mathbf{v}_1(\tau) = \mathbf{c}\mathbf{s}_1(\tau) + \mathbf{c}\mathbf{d}_1(\tau)$ 18 This gives us a model unbiased by the channel. A 19 similar process could be applied to each model 20 whereby it has the channel bias removed by its own silence model. The test speaker could be similarly 21 22 treated, i.e. its silence model is used to remove the channel effects. 23 24 The removal (reduction) of the channel 25 characteristics using the silence model as described 26 above requires suitable channel noise and perfect 27 detection of the silence parts of the utterance. As 28

these cannot be guaranteed they need to be mitigated

we will include some of the claimed identity speaker

(for instance, if the silence includes some speech

58

static speech and inadvertently remove it). 1 Fortunately they can be dealt with in one simple 2 modification to the process: the cohort models 3 should all be referred to the same silence model. 4 5 That is, if we re-add the silence average of the 6 claimed identity model to all of the models in the 7 cohort (including the claimed identity model). This 8 refers all of the models to the same mean 9 $sn(\tau) + h_a(\tau) \,.$ This normalisation is also applied to the 10 test model, thereby referring all of the models and 11 the test utterance to the same reference point. In 12 effect we choose a reference channel and noise 13 condition and refer all others to it. 14 15 This is illustrated diagrammatically in Fig. 19, 16 which shows the Cepstral coefficients of the test 17 utterance together with the claimed identity model 18 and the cohort models 1 to m being input to the 19 classifier 110. A "silence model" or "normalisation 20 model" 400 derived from the claimed identity 21 enrolment data is used to normalise each of these 22 before input to the classifier, so that the actual 23 inputs to the classifier are a normalised test 24 utterance, normalised claimed identity model and 25 normalised cohort models. Ideally, the 26 normalisation model 400 is based on data from 27 periods of silence in the claimed identity enrolment 28 sample as discussed above, but it could be derived 29 from the complete claimed identity enrolment sample. 30 In practical terms, the normalisation model 31 comprises a single row of Cepstral coefficients, 32

59

each of which is the mean value of one column (or 1 2 selected members of one column) of Cepstral coefficients from the claimed identity model. These 3 mean values are used to replace the mean values of 4 each of the sets of input data. That is, taking the 5 test utterance as an example, the mean value of each 6 column of the test utterance Cepstral coefficients 7 8 is subtracted from each individual member of that 9 column and the corresponding mean value from the normalisation model is added to each individual 10 11 member of the column. A similar operation is applied to the claimed identity model and each of 12 13 the cohort models. 14 It will be understood that the normalisation model 15 16 could be derived from the claimed identity model or from the test utterance or from any of the cohort 17 models. It is preferable for the model to be 18 19 derived from either the claimed identity model or the test utterance, and it is most preferable for it 20 21 to be derived from the claimed identity model. 22 normalisation model could be derived from the "raw" 23 enrolment sample Cepstral coefficients or from final model after Vector Quantisation. That is, it could 24 be derived at the time of enrolment and stored along 25 with the enrolment model or it could be calculated 26 when needed as part of the verification process. 27 Generally, it is preferred that a normalisation 28 29 model is calculated for each enrolled speaker at the 30 time of enrolment and stored as part of the enrolled speaker database. 31 32

60

1 These normalisation techniques can be employed with 2 various types of speaker recognition systems but are 3 advantageously combined with the speaker recognition 4 systems of the present invention. 5 6 Speaker recognition systems in accordance with the invention provide improved real world performance 7 for a number of reasons. Firstly, the modelling 8 techniques employed significantly improve separation 9 10 between true speakers and impostors. This improved modelling makes the system less sensitive to real 11 world problems such as changes of sound system 12 (voice sampling system) and changes of speaker 13 characteristics (due to, for example, colds etc.). 14 Secondly, the modelling technique is non-temporal in 15 nature so that it is less susceptible to temporal 16 voice changes, thereby providing longer persistence 17 of speaker models. Thirdly, the use of filter pre-18 processing allows the models to be used for variable 19 bandwidth conditions; e.g. models created using high 20 fidelity sampling systems such as multimedia PCs 21 will work with input received via reduced bandwidth 22 input channels such as telephony systems. 23 24 It will be understood that the preferred methods in 25 accordance with the present invention are inherently 26 suited for use in text-independent speaker 27 recognition systems as well as text-dependent 28 29 systems. 30 31 SYSTEMS

61

PCT/GB02/02726

The invention thus provides the basis for flexible, 1 2 reliable and simple voice recognition systems operating on a local or wide area basis and 3 4 employing a variety of communications/input channels. Fig. 16 illustrates one example of a wide 5 area system operating over local networks and via 6 7 the Internet, to authenticate users of a database system server 400, connected to a local network 402, 8 such as an Ethernet network, and, via a router 404, 9 to the Internet 406. A speaker authentication 10 system server 408, implementing a speaker 11 12 recognition system in accordance with the present invention, is connected to the local network for the 13 purpose of authenticating users of the database 400. 14 15 Users of the system may obviously be connected directly to the local network 402. More generally, 16 users at sites such as 410 and 412 may access the 17 18 system via desktop or laptop computers 414, 416 equipped with microphones and connected to other 19 local networks which are in turn connected to the 20 Internet 406. Other users such as 418, 420, 422 may 21 22 access the system by dial-up modem connections via the public switched telephone network 424 and 23 24 Internet Service Providers 426. 25 26 **IMPLEMENTATION** 27 The algorithms employed by speaker recognition 28 systems in accordance with the invention may be 29 30 implemented as computer programs using any suitable programming language such as C or C++, and 31 32 executable programs may be in any required form

1	including stand alone applications on any					
2	hardware/operating system platform, embedded code in					
3	DSP chips etc. (hardware/firmware implementations),					
4	or be incorporated into operating systems (e.g. as					
5	MS Windows DLLs). User interfaces (for purposes of					
6	both system enrolment and subsequent system access)					
7	may similarly be implemented in a variety of forms,					
8	including Web based client server systems and Web					
9	browser-based interfaces, in which case speech					
10	sampling may be implemented using, for example,					
11	ActiveX/Java components or the like.					
12						
13	Apart from desktop and laptop computers, the system					
14	is applicable to other terminal devices including					
15	palmtop devices, WAP enabled mobile phones etc. via					
16	cabled and/or wireless data/telecommunications					
17	networks.					
18						
19	APPLICATIONS					
20						
21	Speaker recognition systems having the degree of					
22	flexibility and reliability provided by the present					
23	invention have numerous applications. One					
24	particular example, in accordance with a further					
25	aspect of the present invention, is in providing an					
26	audit trail of users accessing and/or modifying					
27	digital information such as documents or database					
28	records. Such transactions can be recorded,					
29	providing information regarding the date/time and					
30	identity of the user, as is well known in the art.					
31	However, conventional systems do not normally verify					
32	or authenticate the identity of the user.					

63

PCT/GB02/02726

1							
2	Speaker recognition, preferably using a speaker						
3	recognition system in accordance with the present						
4	invention, may be used to verify the identity of a						
5	user whenever required; e.g. when opening and/or						
6	editing and/or saving a digital document, database						
7	record or the like. The document or record itself						
8	may be marked with data relating to the speaker						
9	verification procedure, or such data may be recorded						
10	in a separate audit trail, providing a verified						
11	record of access to and modification of the						
12	protected document, record etc. Unauthorised users						
13	identified by the system will be denied access or						
14	prevented from performing actions which are						
15	monitored by the system.						
16							
17	Improvements and modifications may be incorporated						
18	without departing from the scope of the invention as						
19	defined in the appended claims.						

1	C.1	a	i.	m	9
_	~=		=		=

2

3 1. A method of processing speech samples to obtain

64

- 4 a model of a speech sample for use in a speaker
- 5 recognition system, comprising:
- dividing the speech sample into a plurality of
- 7 frames;
- 8 for each frame, obtaining a set of feature
- 9 vectors representing the smoothed frequency spectrum
- 10 of the frame;
- 11 applying a clustering algorithm to the feature
- vectors of the frames to obtain a reduced data set
- 13 representing the original speech sample;
- wherein the adjacent frames are overlapped by
- 15 at least 80%.

16

- 17 2. The method of claim 1, wherein the adjacent
- frames are overlapped by less than 95%.

19

- 20 3. The method of claim 1, wherein the adjacent
- frames are overlapped by an amount in the range 80%
- 22 to 90%.

23

- 4. The method of any preceding claim, wherein the
- 25 clustering algorithm comprises a Vector Quantisation
- 26 algorithm or a k-means algorithm.

- 28 5. The method of any preceding claim, wherein the
- 29 set of feature vectors representing the smoothed
- 30 frequency spectrum of the frame is obtained by means
- of Linear Predictive Coding/Cepstral analysis [LPCC]

65

1 or Fast Fourier Transform [FFT] or by use of a 2 filter bank. 3 The method of any preceding claim, further 4 including storing the model of the speech sample and 5 6 the identity of the speaker in a database of enrolment models of speakers enrolled in a speaker 7 8 recognition system. 9 The method of claim 6, wherein each enrolment 10 7. model comprises a plurality of speech sample models 11 12 representing a plurality of different utterances. 13 The method of claim 6 or claim 7, wherein each 14 8. 15 enrolment model comprises a plurality of speech sample models representing the same utterance 16 modelled using a plurality of parallel modelling 17 18 processes. The method of any of claims 6 to 8, further

19

20 21 including associating the model of the speech sample 22 with a cohort comprising a predetermined number of other speakers enrolled in the speaker recognition 23 24 system.

25

31

The method of any of claims 6 to 9, including 26 processing a second speech sample to obtain a test 27 model of the second speech sample for testing 28 against said database of enrolment models, wherein 29 processing said second speech sample comprises: 30 dividing the second speech sample into a

32 plurality of frames; and

PCT/GB02/02726

66

speakers;

32

WO 02/103680

for each frame, obtaining a set of feature 1 2 vectors representing the smoothed frequency spectrum of the frame. 3 4 The method of claim 10, wherein the set of 5 feature vectors representing the smoothed frequency 6 spectrum of the frame of the second sample is 7 obtained by means of Linear Predictive 8 Coding/Cepstral analysis [LPCC] or Fast Fourier 9 Transform [FFT] or by use of a filter bank. 10 11 The method of claim 10 or claim 11, wherein the 12 test model comprises a plurality of speech sample 13 models representing the same utterance modelled 14 using a plurality of parallel modelling processes. 15 16 17 The method of any of claims 10 to 12, wherein the identity of the speaker of the second speech 18 19 sample is tested by testing the test model against the enrolment model for the claimed identity and the 20 associated cohort as defined in claim 9. 21 22 A method of speaker recognition in which a 23 plurality of speakers to be recognised by a speaker 24 recognition system are enrolled by storing an 25 enrolment model for each speaker in a database of 26 enrolled speakers, the enrolment model representing 27 28 at least one speech sample from that speaker, wherein: 29 each enrolled speaker is associated with a 30 31 cohort of a predetermined number of other enrolled

32

a test speech sample from a speaker claiming to 1 2 be one of the enrolled speakers is modelled and tested, using a classification process, against the 3 enrolment model of the claimed speaker and the 4 enrolment models of the associated cohort; and 5 6 the classification process always matches the 7 test model with either the claimed speaker or one of 8 the associated cohort such that a false acceptance 9 rate of the system is determined by the cohort size. 10 11 The method of claim 14, wherein modelling processes used for modelling the enrolled speaker 12 speech samples and the test speech sample and/or 13 classification processes used for testing the test 14 model against the enrolment models are selected to 15 provide a false rejection rate substantially equal 16 17 to zero, so that an overall error rate of the system is determined substantially only by the false 18 19 acceptance rate. 20 The method of claim 14 or claim 15, wherein the 21 test model is tested using multiple parallel 22 23 classification processes and the test model is matched with an enrolment model only if at least a 24 25 predetermined number of the parallel classification processes produces a match with that enrolment 26 model, so as to reduce the false acceptance rate of 27 the system for a given cohort size. 28 29 The method of claim 16, wherein the enrolment 30 models and test model are each obtained using 31

multiple parallel modelling processes and the

67

PCT/GB02/02726

68

1 parallel classification processes compare the 2 results of the parallel modelling processes applied 3 to the test speech sample with corresponding results of the parallel modelling processes applied to the 4 5 enrolment speech samples. 6 The method of any one of claims 8, 12 or 17, 7 8 wherein the parallel modelling processes comprise at 9 least one of: different frequency banding applied to the 10 11 speech samples; 12 different spectral modelling applied to the 13 speech samples; and different clustering applied to the feature 14 15 vectors representing the speech samples. 16 The method of claim 16, wherein the parallel 17 19. 18 classification processes comprise testing the test model against different cohorts of enrolled 19 20 speakers. 21 The method of claim 16, wherein the parallel 22 classification processes comprise testing the test 23 model against different utterances represented by 24 the enrolment models. 25 26 27 The method of any of claims 14 to 20, wherein 21. the enrolment models and test model are obtained 28 29 using the method of any of claims 1 to 13. 30

A method of normalising speech models in a 31

32 speaker recognition system of the type in which

69

- speech samples are input to the system via different
- 2 input channels having different channel
- 3 characteristics, and wherein a test model
- 4 representing a test sample is tested, using a
- 5 classification process, against a set of enrolment
- 6 models representing speech samples from speakers
- 7 enrolled in the system, comprising deriving a
- 8 normalisation model from the test speech sample or
- 9 from one of the enrolment speech samples and using
- the normalisation model to normalise the test model
- and the enrolment models against which the test
- model is to be tested prior to testing the
- 13 normalised test model against the normalised
- 14 enrolment models.

15

- 16 23. The method of claim 22, wherein the
- 17 normalisation model is derived from the enrolment
- speech sample for the identity claimed for the test
- 19 speech sample.

20

- 21 24. The method of claim 23, wherein the
- 22 normalisation model is derived from the enrolment
- 23 model for the identity claimed for the test speech
- 24 sample.

- 26 25. The method of any of claims 22 to 24, wherein
- 27 the speech samples are divided into a plurality of
- frames, a set of feature vectors are obtained
- 29 representing the smoothed frequency spectrum of each
- frame, and the normalisation model is obtained by
- 31 calculating the mean values of sets of feature
- 32 vectors from at least some of said frames of the

70

speech sample from which the normalisation model is derived. 3 The method of claim 25, wherein the frames used 4 for deriving the normalisation model are frames 5 corresponding to periods of silence in the speech 6 7 sample from which the normalisation model is derived. 8 9 The method of claim 25 or claim 26, wherein the 10 11 test model and enrolment models are normalised by replacing mean values of the feature vectors of the 12 test model and the enrolment models with the 13 corresponding mean values from the normalisation 14 model. 15 16 17 The method of any of claims 22 to 27, wherein 18 the speech samples are processed using the method of 19 any of claims 1 to 13. 20 The method of any of claims 14 to 21, wherein 21 the test model and enrolment models are normalised 22 23 prior to classification using the method of any of claims 22 to 28. 24 25 The method of any preceding claim, wherein 26 speech samples are input to a speaker recognition 27 28 system via an input channel having a transfer function which modifies the speech sample data, 29 comprising estimating the transfer function of said 30

input channel and normalising the modified speech

71

sample data using the inverse of said estimated

2 transfer function.

3

4 31. A speaker recognition system comprising data

5 processing and storage means adapted to implement

6 the method of any of claims 1 to 30.

7

8 32. A computer program comprising symbolic code for

9 instructing a computer to execute the method of any

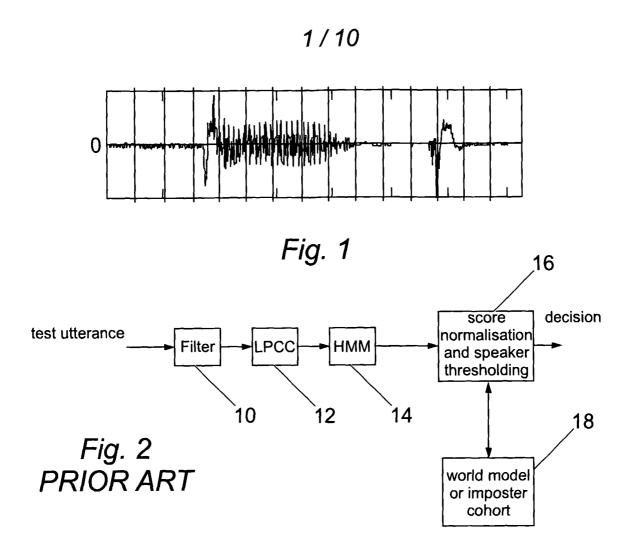
10 of claims 1 to 30.

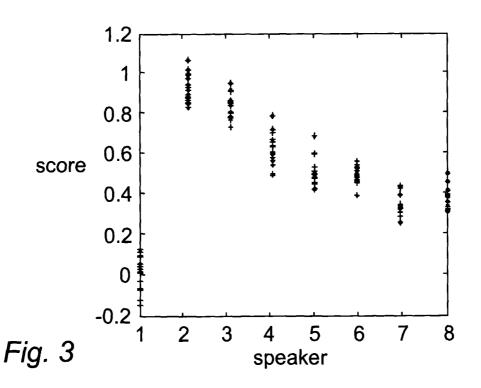
11

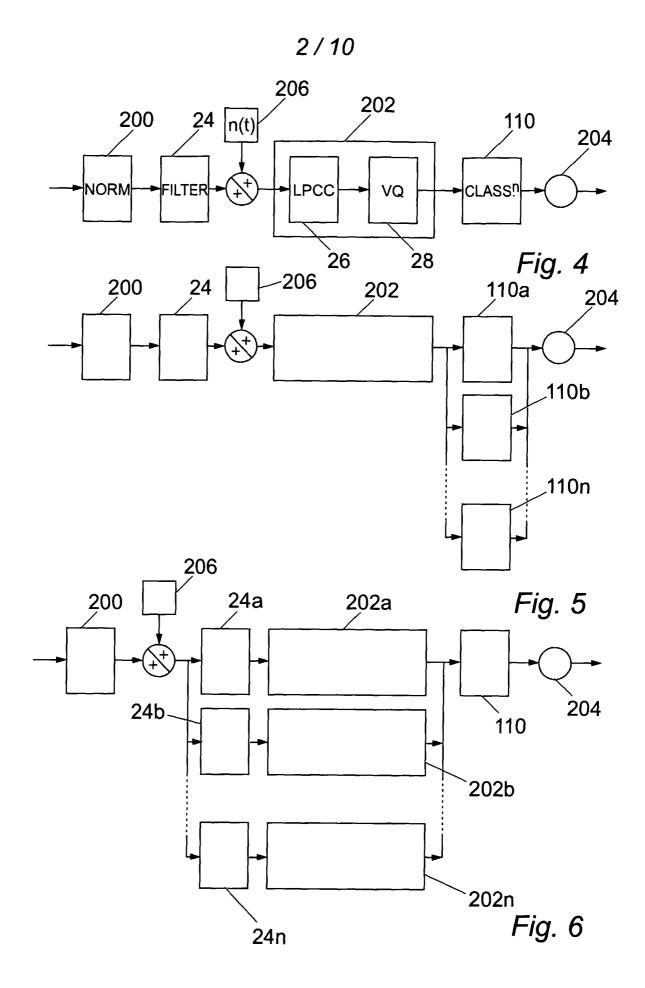
12 33. A data carrier encoded with a computer program

comprising symbolic code for instructing a computer

to execute the method of any of claims 1 to 30.







3/10

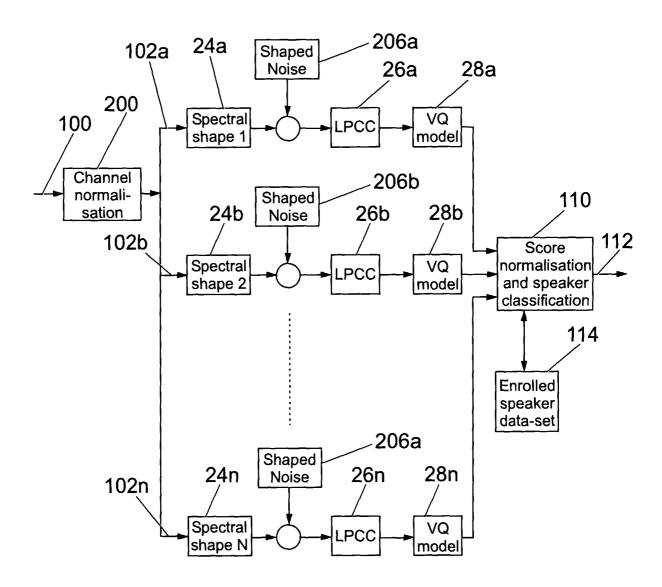
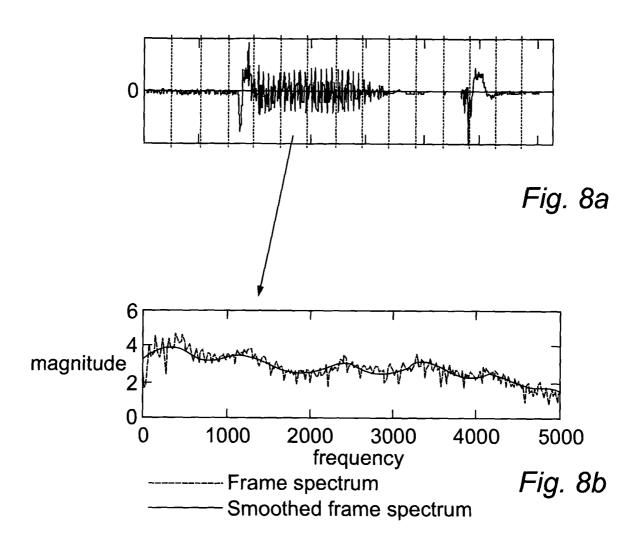
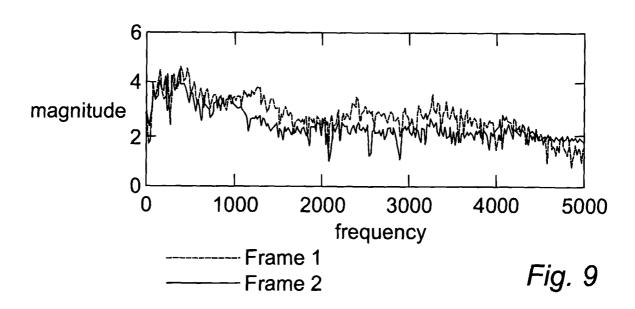


Fig. 7





5/10

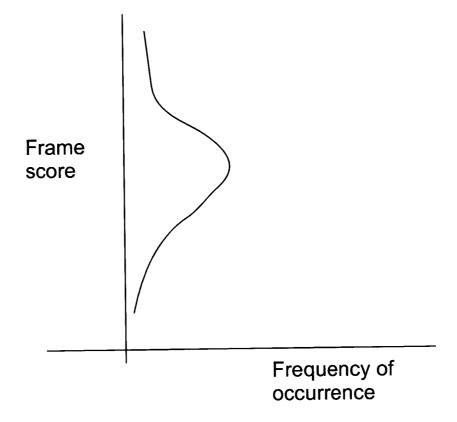
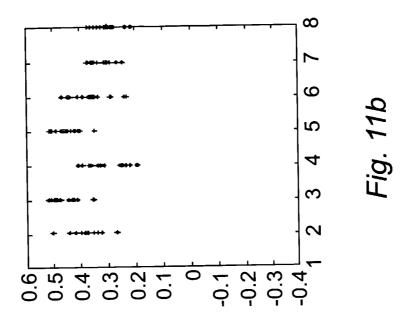
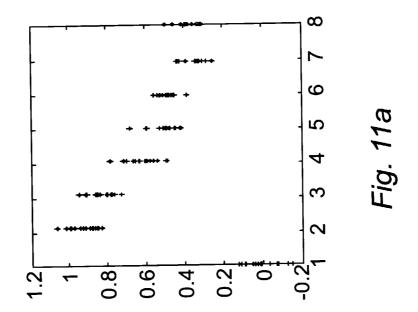
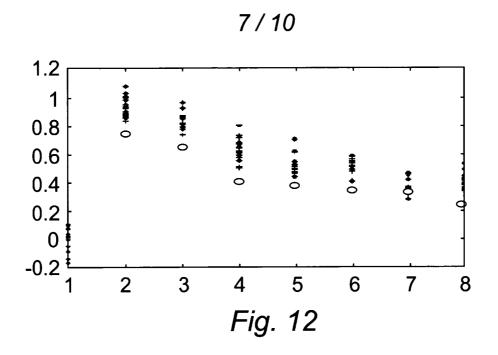


Fig. 10







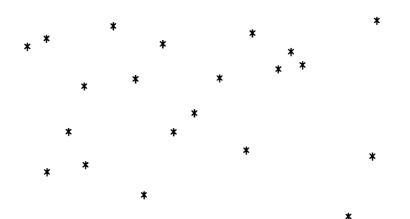


Fig. 13

*

*

*

*

*

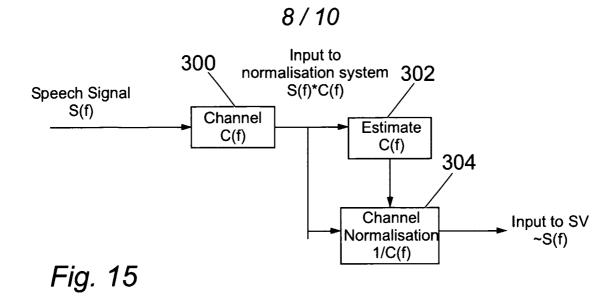
*

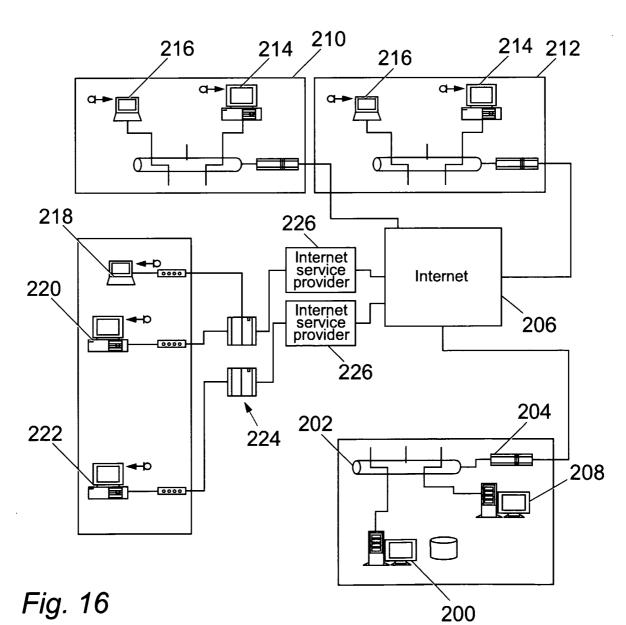
*

*

Fig. 14

*





9/10

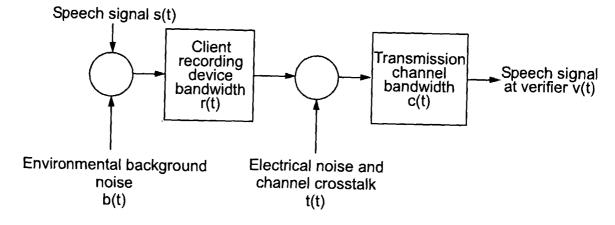


Fig. 17

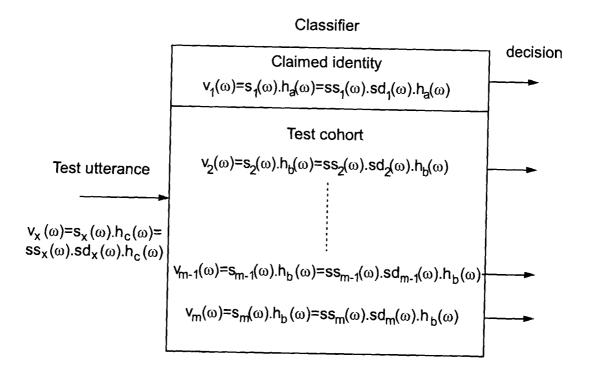
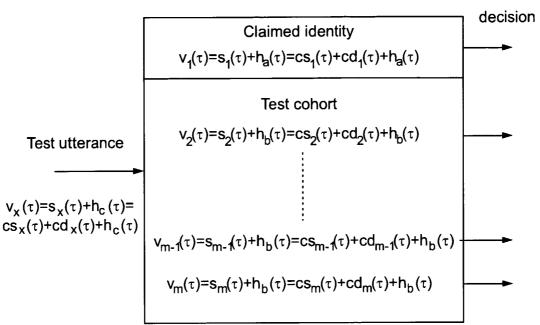


Fig. 18

10/10

Classifier



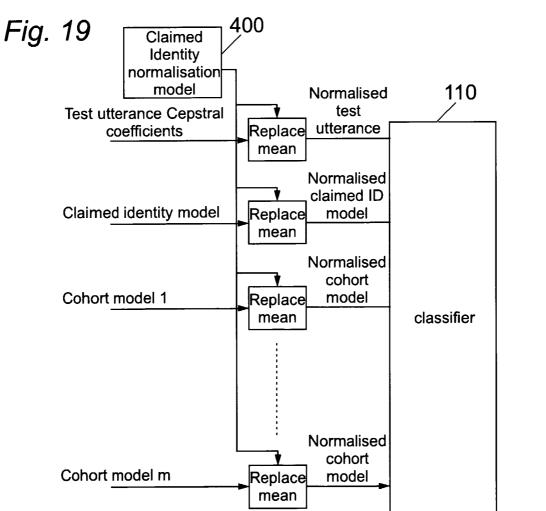


Fig. 20