



Europäisches Patentamt
European Patent Office
Office européen des brevets



(11) **EP 1 515 305 B1**

(12) **EUROPEAN PATENT SPECIFICATION**

(45) Date of publication and mention
of the grant of the patent:
08.02.2006 Bulletin 2006/06

(51) Int Cl.:
G10L 15/06^(2006.01) G10L 15/20^(2006.01)

(21) Application number: **04019236.1**

(22) Date of filing: **13.08.2004**

(54) **Noise adaption for speech recognition**

Rauschadaptierung zur Spracherkennung

Adaptation au bruit pour la reconnaissance de la parole

(84) Designated Contracting States:
DE GB

(30) Priority: **12.09.2003 JP 2003321648**

(43) Date of publication of application:
16.03.2005 Bulletin 2005/11

(73) Proprietors:
• **Furui, Sadaoki**
Tokyo 154-0014 (JP)
• **NTT DoCoMo, Inc.**
Tokyo 100-6150 (JP)

(72) Inventors:
• **Furui, Sadaoki**
Tokyo 154-0014 (JP)
• **Zhang, Zhipeng**
c/o Intellectual Property Dept.
Chiyoda-ku
Tokyo 100-6150 (JP)
• **Horikoshi, Tsutomu**
c/o Intellectual Property Dept.
Chiyoda-ku
Tokyo 100-6150 (JP)

• **Sugimura, Toshiaki**
c/o Intellectual Property Dept.
Chiyoda-ku
Tokyo 100-6150 (JP)

(74) Representative: **Sparing Röhl Henseler**
Patentanwälte
European Patent Attorneys
Rethelstrasse 123
40237 Düsseldorf (DE)

(56) References cited:
• **ZHIPENG ZHANG ET AL: "A Tree-Structured Clustering Method Integrating Noise and SNR for Piecewise Linear-Transformation-Based Noise Adaptation" PROC. IEEE INTERNATIONAL CONFERENCE ON ACOUSTICS, SPEECH, AND SIGNAL PROCESSING, ICASSP 2004, vol. 1, 17 May 2004 (2004-05-17), pages 981-984, XP002297864 MONTREAL, QUEBEC, CANADA**
• **ZHIPENG ZHANG ET AL: "Tree-Structured Noise-Adapted HMM Modeling for Piecewise Linear-Transformation-Based Adaptation" PROC. EUROSPEECH 2003, 1 September 2003 (2003-09-01), pages 669-672, XP002297865 GENEVA, SWITZERLAND**

Note: Within nine months from the publication of the mention of the grant of the European patent, any person may give notice to the European Patent Office of opposition to the European patent granted. Notice of opposition shall be filed in a written reasoned statement. It shall not be deemed to have been filed until the opposition fee has been paid. (Art. 99(1) European Patent Convention).

Description

[0001] The invention relates to a noise adaptation system of speech model, a noise adaptation method, and a noise adaptation program that use noisy speech to be recognized to adapt a clean speech model generated by modeling features of speech by means of a Hidden Markov Model (HMM) so that the recognition rate for the noisy environment can be improved.

[0002] A tree-structure piecewise linear transformation approach is described in an article entitled "Effects of tree-structure clustering in noise adaptation using piecewise linear transformation" by Zhipeng Zhang et al. (Proceedings of 2002 Autumn Meeting of the Acoustical Society of Japan, pp.29-30). According to the approach described in the article, noise is clustered, a tree-structure noisy speech model space is generated based on the result of the clustering, a speech feature parameter of input noisy speech to be recognized is extracted, an optimum model is selected from the tree-structure noisy speech model space, and linear transformation is applied to the selected model so as to increase the likelihood of the selected model, thereby improving the accuracy of input speech.

[0003] Another approach is described in an article entitled "Study on tree-structure clustering in noise adaptation using piecewise linear transformation" by Zhipeng Zhang et al. (2003 Spring Meeting of the Acoustical Society of Japan, pp.37-38), in which noise characteristics are sequentially and hierarchically divided to generate a tree structure of a noise-added speech model. In this approach, noise-added speech is first clustered by signal-to-noise ratio (hereinafter abbreviated to SNR) and then a tree-structure model is provided for each SNR condition to generate a tree-structure noisy speech model space.

[0004] FIG. 6 shows an example of the tree-structure noisy speech model. In FIG. 6, a tree-structure noisy speech model is provided for each of three SNR conditions. In FIG. 6, a tree-structure model for SNR = 5 dB is indicated by K1, a tree-structure model for SNR = 10 dB is indicated by K2, and a tree-structure model for SNR = 15 dB is indicated by K3. The top node (root) of each tree-structure model K1 - K3 represents a clean speech model. Higher levels of each tree structure represent global features of noise characteristics and lower levels represent local features.

[0005] Described in Japanese Patent Laid-Open No. 2002-14692 (FIGS. 2 and 3 and Abstract, in particular) is a technology in which a large number of noise samples are clustered beforehand, acoustic models are generated on the basis of the samples, and noise selected through clustering is added to learning data, thereby enabling efficient learning with a small number of noise samples to achieve high recognition performance.

[0006] Japanese Patent Laid-Open No. 2002-91484 (Abstract, in particular) described a technology in which a language model is generated for each tree-structure

cluster, which is used for speech recognition.

[0007] Japanese Patent Laid-Open No. 2000-298495 (Abstract and Claim 2, in particular) describes combining a number of tree structures to form a new tree structure.

[0008] In the approach in "Study on tree-structure clustering in noise adaptation using piecewise linear transformation" cited above, input noisy speech to be recognized is analyzed to extract a feature parameter string and an optimum model is selected from a tree-structure noisy speech model space. Linear transformation is applied to the selected optimum model to maximize the likelihood. Accordingly, this approach has a drawback that recognition involves a two-step search: an optimum model is first selected under each SNR condition and then the best model is selected from among all SNR models. Problems here are the difficulty of dealing with noisy speech with varying SNR and high costs of computing the conditions.

[0009] None of the technologies described in the above-sited documents can solve these problems.

[0010] An object of the present invention, as defined by the appended independent claims, is to provide a noise adaptation system, a noise adaptation method, and a noise adaptation program for speech recognition that can readily deal with noisy speech with varying SNR and can minimize computation costs by generating a speech model with a single-tree-structure into which noise and SNR are integrated.

[0011] According to an aspect of the invention, there is provided a noise adaptation system of speech model for adapting a speech model for any noise to speech to be recognized in a noisy environment, the speech model being learned by using noise data stored in a noise database and clean speech data, the system comprising: a clustering means of clustering all noise data stored in the noise database; a speech model means generating means of generating a single-tree-structure noisy speech model space based on the result of the clustering performed by the clustering means; a parameter extracting means of extracting a speech feature parameter of input noisy speech to be recognized; a selecting means of selecting an optimum model from the tree-structure noisy speech model space generated by the speech model space generating means; and a linear transformation means of applying linear transformation to the model selected by the selecting means so that the model provides a further increased likelihood. The single-tree-structure noisy speech model space generated as described above allows noisy speech with varying SNR to be readily dealt with and the computation cost to be saved.

[0012] According to a further aspect of the invention, there is provided the noise adaptation system of speech model according to the first aspect, wherein the clustering means generates the noise-added speech by adding the noise to the speech in accordance with a signal-to-noise ratio condition, subtracts the mean value of speech cepstral of the generated noise-added speech, generates a Gaussian distribution model of each of pieces of gen-

erated noise-added speech, and calculates the likelihood between the pieces of noise-added speech to generate a likelihood matrix to provide a clustering result. This allows noise-added speech to be clustered.

[0013] According to an additional aspect of the invention, there is provided the noise adaptation system according to first or second aspect, wherein the selecting means selects a model that provides the highest likelihood for the speech feature parameter extracted by the parameter extracting means. By selecting the model that provides the highest likelihood, the accuracy of speech recognition can be increased. The selecting means may select a model by searching the tree-structure noisy speech model space downward from the highest level to the lowest level. By searching the tree structure from the highest level to the lowest level, an optimum model can be selected.

[0014] The linear transformation means of the noise adaptation system may perform the linear transformation on the basis of the model selected by the selecting means to increase the likelihood. By performing the linear transformation, the likelihood can be maximized.

[0015] According to a still further aspect of the invention, there is provided a noise adaptation method for adapting a speech model for any noise to speech to be recognized in a noisy environment, the speech model being learned by using noise data stored in a noise database and clean speech data, the method comprising: a clustering step of clustering all noise-added speech data stored in the noise database; a speech model space generating step of generating a single-tree-structure noisy speech model space based on the result of the clustering performed in the clustering step; a parameter extracting step of extracting a speech feature parameter of input noisy speech to be recognized; a selecting step of selecting an optimum model from the tree-structure noisy speech model space generated in the speech model space generating step; and a linear transformation step of applying linear transformation to the model selected in the selecting step so that the model provides a further increased likelihood. The single-tree-structure noisy speech model space allows noisy speech with varying SNR to be readily dealt with and the computation cost to be saved.

[0016] According to another aspect of the invention, there is provided a noise adaptation program for adapting a speech model for any noise to speech to be recognized in a noisy environment, the speech model being learned by using noise data stored in a noise database and clean speech data, the program comprising: a clustering step of clustering all noise-added speech data stored in the noise database; a speech model space generating step of generating a single-tree-structure noisy speech model space based on the result of the clustering performed in the clustering step; a parameter extracting step of extracting a speech feature parameter of input noisy speech to be recognized; a selecting step of selecting an optimum model from the tree-structure noisy speech model

space generated in the speech model space generating step; and a linear transformation step of applying linear transformation to the model selected in the selecting step so that the model provides a further increased likelihood.

The single-tree-structure noisy speech model space allows noisy speech with varying SNR to be readily dealt with and the computation cost to be saved.

[0017] In effect, according to the invention, all pieces of noise data in a noise database (hereinafter abbreviated to DB) are used to cluster noise-added speech into a single-tree-structure based on every SNR condition. A noise-added speech space is partitioned in a tree structure according to SNRs and noise characteristics, and sound feature parameter strings of input noisy speech to be recognized are extracted. Then, an optimum model is selected from the tree-structure model space on the basis of the feature parameter string and liner transformation is applied to this model.

[0018] The single-tree-structure into which noise and SNR are integrated is generated to allow the most likely noise-added speech model to be learned. Thus, a high recognition accuracy can be achieved. Furthermore, the approach of the present invention does not require selecting an optimum model under each individual SNR condition. Instead, the approach of the present invention involves just a one-step search through which the best model among all SNR models is selected. Therefore, noisy speech with varying SNR can be readily dealt with and the computation costs can be saved.

[0019] According to the invention, noisy speech with varying SNR can be readily dealt with and the computation costs can be saved by using a single-tree-structure noisy speech model space.

[0020] Noise-added speech can be clustered by adding noise to speech according to signal-to-noise conditions, subtracting the mean value of speech cepstral of each of the pieces of generated noise-added speech, generating a Gaussian distribution model of each of the pieces of noise-added speech, and calculating the likelihood between the pieces of noise-added speech to generate a likelihood matrix.

[0021] An improved accuracy of speech recognition can be achieved by selecting a model that provides the highest likelihood for an extracted speech feature parameter.

[0022] An optimum model can be selected by searching the tree-structure noisy speech model space from the highest level to the lowest level for an optimum model.

[0023] The likelihood can be maximized by performing linear transformation on the basis of the selected model so as to increase the likelihood.

[0024] The invention will now be described in connection with preferred embodiments as shown in the drawings.

[0025] FIG. 1 is a flowchart of a process performed by a noise adaptation system of speech model according to the invention.

[0026] FIG. 2 is a block diagram showing a configura-

tion of a noise adaptation system of speech model according to one embodiment of the invention.

[0027] FIG. 3 is a functional block diagram in which components shown in FIG. 2 are rearranged in accordance with operation flow in the system.

[0028] FIG. 4 is a conceptual diagram showing a process for selecting an optimum model in a tree-structure noisy speech model space in the system.

[0029] FIG. 5 shows a word accuracy achieved by using a speech HMM adapted by the system.

[0030] FIG. 6 is a conceptual diagram showing a process for selecting an optimum model in a tree-structure noisy speech model space used in a conventional noise adaptation system of speech model.

[0031] According to the invention, a noisy speech model space is generated as a tree structure by using SNR and sound quality. To generate the noisy speech model space, a noise database is used to add noise to clean speech according to every SNR condition to produce noise-added speech. Then, noise characteristics are represented as a single-tree-structure to provide a model, in which higher levels of the tree structure represent global features of noise characteristics and lower levels represent local features. An optimum piecewise space of noise can be selected by following the tree structure downward from the root in top-down fashion to select an optimum model.

[0032] Because noise-added speech is consistently used both in clustering and model learning processes, the noise-added speech model that provides the highest likelihood can be learned and an improved accuracy of recognition can be achieved. Configuration of the present system

[0033] A configuration for implementing the above-described process will be described with reference to FIG. 2, which is a block diagram showing one embodiment of the noise adaptation system according to the invention. As shown in FIG. 2, the noise adaptation system according to the embodiment includes a tree-structure-model storage 1, a feature extraction unit 2, a speech recognition unit 3, a model selection and determination unit 4, a model linear transformation adaptation unit 5, and a recognition result storage 6. The present system is realized as terminal apparatus, mobile terminal, server computer, personal computer, and other equipment comprising above units and storages.

[0034] The tree-structure-model storage 1 stores a noise-added speech HMM which has been built as a single-tree-structure based on a result of clustering of noise-added speech.

[0035] The feature extraction unit 2 analyzes speech data inputted to it and transforms it into feature vectors.

[0036] The speech recognition unit 3 applies a Viterbi algorithm to the time series feature vector time transformed from the input speech data to obtain a model sequence that provides the highest likelihood function.

[0037] The model selection and determination unit 4 selects an optimum model that provides an optimum

model that provides the highest likelihood from among models stored in the tree-structure-model storage 1.

[0038] The model linear transformation adaptation unit 5 applies linear transformation to the model selected by the model selection and determination unit 4 so as to maximize its likelihood.

[0039] The recognition result storage 6 stores speech recognition results. Operation of the system

[0040] Operation of the system having the structure described above will be described with reference to FIGS. 1 and 3. FIG. 3 is a functional block diagram in which the components 1 - 6 shown in FIG. 2 are rearranged according to the flow of operation in the system. FIG. 1 is a flowchart of a process performed by the system.

[0041] The process for performing speech recognition in the system follows steps S1 to S9 as described below.

[0042] Step S1 (the step of generating noise-added speech): Every piece of noise data stored in a noise database is used to add noise to clean speech according to every SNR condition to generate noise-added speech.

[0043] Step S2 (the step of subtracting the mean of noise-added speech): CMS (Cepstral Mean Subtraction) is applied to noise-added speech generated at step S1. CMS is a technique for subtracting the mean of speech cepstral. That is, the mean cepstral of value of all frames of speech data in a certain interval is calculated and the mean value is subtracted from the vector of each frame. The cepstral is the Fourier transform of the logarithm of a power spectrum obtained by Fourier transform. The CMS is described in a document entitled "Furui: Cepstral Analysis Technique For Automatic Speaker Verification, IEEE Transaction on Acoustical Speech and Signal Processing, Vol. ASSP - 29, pp.254-272, 1981."

[0044] Step S3 (the step of generating a noise-added speech model): A Gaussian mixture model (GMM) of each noise-added speech is generated by means of the Baum-Welch algorithm. Baum-Welch algorithm is a repetitive approach to getting closer to an optimum value, starting from an appropriate initial value. The Baum-Welch algorithm is described in a document entitled "Speech recognition with probabilistic model" by Seiichi Nakagawa (Institute of Electronics, Information and Communication Engineers, 1988).

[0045] Step S4 (the step of clustering noise-added speech): The GMM is used to calculate the likelihoods between pieces of noise-added speech to generate a likelihood matrix. A SPLIT method based on the likelihood matrix is used to serially cluster the noise-added speech. In the SPLIT method, clusters that provide the largest distortion are split sequentially. Consequently, any number of clusters can be generated. The result of clustering can be obtained fully automatically simply by giving the number of clusters. The SPLIT method is described in a Speech Committee document by Sugamura et al. (S82-64, 1982).

[0046] Step S5 (application to piecewise linear transformation adaptation): A tree-structure clustering result

of the noise-added speech is provided by the step S4. The clustering result is stored in the tree-structure model storage 1. The clustering result is a single-tree-structure model in which noise and SNR are integrated. Also, the clustering result represents features in tree-structure form; global features of the noise-added speech are represented at a higher level of the tree structure and local features of the speech are represented at a lower level.

[0047] The clustering result stored in the tree-structure-model storage 1 is applied to piecewise linear transformation. The piecewise linear transformation is described in the above-cited article "Effects of tree-structure clustering in noise adaptation using piecewise linear transformation" by Zhipeng Zhang et al. In particular, steps S6 to S9 described below are performed.

[0048] Step S6 (the step of extracting feature quantities): The feature extraction unit 2 extracts feature quantities from noise-added speech data to be recognized. In the feature quantity extraction, LPC (Linear Prediction Coding) analysis is applied to each frame of inputted speech data to obtain time series feature parameter vectors such as a cepstral or _ cepstral, as a feature parameter sequence.

[0049] Step S7 (selecting an optimum model): The step of selecting an optimum model will be described with reference to FIG. 4. The node (root) at the top of FIG. 4 represents a clean speech model. Under the root, there are N models, SNR-1 to SNR-N. The N models SNR-1 to SNR-N represent models learned from speech generated by adding all types of noise under all SNR conditions.

[0050] Child nodes below them represent models learned from speech data generated by adding some selected types of noise depending on the clustering result. At the bottom of the tree structure are models learned from speech that are generated by adding only a certain single type of noise. Thus, global noise characteristics are represented at the higher level of the tree structure and local noise characteristics are represented at the lower level.

[0051] Unlike the approach in the above-cited article "Study on tree-structure clustering in noise adaptation using piecewise linear transformation" by Zhipeng Zhang et al. (see FIG. 6), the approach of the present invention does not require selecting an optimum model under each individual SNR condition. Instead, it requires only one-step search in which the best model among all SNR models is selected.

[0052] Returning to FIG. 1, to perform recognition, the likelihood of a given clean model at root is first calculated by using the feature parameter sequence obtained at step S4. This is performed by the speech recognition unit 3 shown in FIG. 1.

[0053] Then, the speech recognition unit 3 uses the models below the root to calculate the likelihoods. The likelihood values thus calculated are used by the model selection and determination unit 4 to select an optimum model. In particular, this is achieved by following the fol-

lowing procedure. Models providing likelihoods higher than that of the clean model at root are reserved. Then the models at the child nodes below them are used to calculate the likelihoods under these SNR conditions.

5 The likelihoods of two child node models are compared with that of the parent node. If a child node model provides the highest likelihood, the likelihoods of the child node models below that node are calculated. On the other hand, if the likelihood of the parent node is higher than those of the child node model, then no further calculation is performed and the parent node is determined as an optimum node.

10 **[0054]** In FIG. 4, the search paths are represented by solid lines. The calculation can be repeated to find an optimum space. Furthermore, the likelihood of the highest-likelihood models under different SNR conditions are compared one another to determine the model providing the highest likelihood among them is selected as the optimum model in the entire noisy speech space. In the example shown in FIG. 4, the fourth node provides the highest likelihood under condition SNR-1. Under SNR-N condition in FIG. 4, the fifth node provides the highest likelihood. The likelihoods of the highest-likelihood models under different SNR conditions are compared with one another to select the model that provides the highest likelihood among the highest-likelihood nodes.

25 **[0055]** Step S8 (linear regression): The model linear transformation adaptation unit 5 applies Maximum Likelihood Linear Regression (hereinafter abbreviated to MLLR) to the selected model so as to provide a further improved likelihood. The MLLR is described in a document entitled "Mean and variance adaptation within the MLLR framework" (M.J.F Gales et al., Computer Speech and Language, pp.249-264, 1996). In particular, a phoneme sequence resulting from recognition is used to estimate a linear transformation matrix on the basis of an maximum likelihood criterion and the mean value and variances of HMM Gaussian distribution are adapted by linear transformation (linear regression).

30 **[0056]** Step S9 (re-recognition): When outputting the result of speech recognition, the speech recognition unit 3 uses the model obtained at step S8 to perform re-recognition and the re-recognition result is stored in the recognition result storage 6.

35 **[0057]** In a noise adaptation system of the present invention, as has been described, all pieces of noise data in a noise database are used to add noise to speech under every SNR condition and learn a noise-added speech model. The distance between all noise models under the SNR conditions are calculated and the noise-added speech is clustered. Based on the result of the noise-added speech clustering, a speech model having a tree structure is generated. Thus, a tree-structure model into which noise and SNR are integrated can be provided and a tree-structure noisy speech model space is generated. In the feature extraction step, an input noisy speech to be recognized is analyzed to extract a feature parameter sequence and the likelihoods of HMMs are

compared with one another to select an optimum model from the tree-structure noisy speech model space. Linear transformation is applied to the model selected from the noisy speech model space so as to provide a further improved likelihood.

[0058] In summary, according to the present invention, every piece of noise-added speech data stored in a noise database is used to add noise to clean speech under every SNR condition to generate noise-added speech (step S1 in FIG. 1). The noise-added speech is clustered to form a single-tree-structure noise-added speech model space. In the noise-added speech model space, each piece of noise belonging to each tree-structure node is added to the clean speech to generate a noise-added speech model (step S3 in FIG. 1). The likelihoods are calculated in the noise-added-speech tree structure model space (step S4 in FIG. 1) and the tree structure is followed downward from the top to select an optimum model (step S7 in FIG. 1). Based on model parameters of an adaptation speech model sequence thus selected, linear transformation is performed so as to maximize the likelihood (step S8 in FIG. 1).

[0059] In effect, according to the invention, a single-tree-structure noise-added speech model space is generated into which noise and SNR are integrated, instead of tree-structure noise-added speech model spaces for individual SNRs. Consequently, Noisy speech with varying SNR can be readily dealt with and the computation cost can be saved.

[0060] The noise-added speech is used not only in the model learning process but also in clustering. Because noise-added speech is consistently used both in clustering and model learning, the most likely noise-added speech model can be learned. As a result, a higher accuracy of recognition can be achieved.

Example

[0061] Effects of recognition of noisy dialog speech that was performed by the present system have been examined. An example of the experiments will be described below.

[0062] A speech HMM used in the experiments is a shared-state, speaker-independent context-dependent phoneme HMM produced by using tree-based clustering. A total of 25 dimensions are used as feature quantities: MFCC (Mel Frequency Cepstral Coefficients) 12 and the first derivative of log power. A "mel frequency" is a value based on the sensitivity of the human ear and often used for representing the level of audibility of a sound. MFCC is generated as follows: discrete Fourier transform is applied to acoustic wave data and the resulting value is transformed into its logarithmic expression. Then inverse discrete Fourier transform is applied to the logarithm to produce a waveform, which is sampled at predetermined intervals. The sampled value is MFCC.

[0063] Effects of the present system will be described below with reference to FIG. 5. FIG. 5 shows a word

accuracy (baseline) achieved by using a given speech HMM and a word accuracy (of the inventive method) achieved by using a speech HMM adapted by the system of the present invention. The vertical axis in FIG. 5 represents the word accuracy (%) and the horizontal axis represents SNR (dB). Indicated on the horizontal axis are SNRs of 5, 10, and 15 dB. The half-tone dot meshing bars in FIG. 5 represent the baseline accuracies and the striped bars represent the accuracies of the present system.

[0064] It can be seen from the results shown in FIG. 5 that the method according to the present invention is more effective than the conventional method. In this example, the word error rate of the present system is lower than the baseline by 56%. That is, the present invention can provide an improved speech recognition accuracy.

(Speech model noise adaptation method)

[0065] The following noise adaptation method is implemented in the noise adaptation system described above. The method is a noise adaptation method for adapting a speech model for any noise that has been learned by using noise data stored in a noise database and clean speech data to speech to be recognized in a noisy environment. The method comprises a clustering step (corresponding to steps S1 to S4 in FIG. 1) of clustering all pieces of noise-added speech data stored in the noise database; a speech model space generating step (corresponding to step S5 in FIG. 1) of generating a single-tree-structure noisy speech model space on the basis of the result of clustering at the clustering step; a parameter extracting step (corresponding to step S6 in FIG. 1) of extracting a speech feature parameter of input noisy speech to be recognized; a selecting step (corresponding to step S7 in FIG. 1) of selecting an optimum model from the tree-structure noisy speech model space generated at the speech model space generating step; and a linear transformation step (corresponding to step S8 in FIG. 1) of applying linear transformation to the model selected at the selecting step so as to provide a further improved likelihood.

[0066] Noisy speech with varying SNR can be readily dealt with and the computation cost can be saved by performing this method and using the single-tree-structure noisy speech model space for speech recognition.

Noise adaptation program of speech model

[0067] A program for performing the process shown in FIG. 1 can be provided and used to control a computer to provide the same effects as those described above. The program is a noise adaptation program for speech recognition that controls a computer to adapt a speech model for any noise that has been learned by using all pieces of noise data stored in a noise database and clean speech data to speech to be recognized in a noisy environment. The program comprises a clustering step (cor-

responding to steps S1 to S4 in FIG. 1) of clustering all pieces of noise-added speech data stored in the noise database; a speech model space generating step (corresponding to step S5 in FIG. 1) of generating a single-tree-structure noisy speech model space on the basis of the result of clustering at the clustering step; a parameter extracting step (corresponding to step S6 in FIG. 1) of extracting a speech feature parameter of input noisy speech to be recognized; a selecting step (corresponding to step S7 in FIG. 1) of selecting an optimum model from the tree-structure noisy speech model space generated at the speech model space generating step; and a linear transformation step (corresponding to step S8 in FIG. 1) of applying linear transformation to the model selected at the selecting step so as to provide a further improved likelihood.

[0068] Noisy speech with varying SNR can be readily dealt with and the computation cost can be saved by executing this program on a computer and using the single-tree-structure noisy speech model space for speech recognition.

[0069] A storage medium for storing the program may be a semiconductor memory, a magnetic disk, an optical disk, or any of other storage media, which are not shown in FIG. 1.

[0070] Automatic speech recognition systems in general can function well under laboratory conditions but their performances drop in real applications. One problem in real-world applications is reduction in performance of recognition of speech containing noise or music in the background. The present invention can solve this noise problem and improve the accuracy of recognition of noise-added speech.

Claims

1. A noise adaptation system of speech model for adapting a speech model for any noise to speech to be recognized in a noisy environment, said speech model being learned by using noise data stored in a noise database which is used to add noise to clean speech according to every SNR condition, and clean speech data, said system comprising:

clustering means for clustering all noise data stored in said noise database;
speech model space generating means for generating a single-tree-structure noisy speech model space based on the result of the clustering performed by said clustering means;
parameter extracting means for extracting a speech feature parameter of input noisy speech to be recognized;
selecting means for selecting an optimum model from the tree-structure noisy speech model space generated by said speech model space generating means; and

linear transformation means for applying linear transformation to the model selected by the selecting means so that the model provides a further increased likelihood.

2. The system of claim 1, wherein said clustering means generates said noise-added speech by adding said noise to said speech in accordance with a signal-to-noise ratio condition, subtracts the mean value of speech cepstral of the generated noise-added speech, generates a Gaussian distribution model of each of pieces of generated noise-added speech, and calculates the likelihood between the pieces of noise-added speech to generate a likelihood matrix to provide a clustering result.
3. The system of claim 1 or 2, wherein said selecting means selects a model that provides the highest likelihood for the speech feature parameter extracted by said parameter extracting means.
4. The system of claim 3, wherein said selecting means selects a model by searching said tree-structure noisy speech model space downward from the highest level to the lowest level.
5. The system of any one of the claims 1 to 4, wherein said linear transformation means performs the linear transformation on the basis of the model selected by said selecting means to increase the likelihood.
6. A noise adaptation method of speech model for adapting a speech model for any noise to speech to be recognized in a noisy environment, said speech model being learned by using noise data stored in a noise database which is used to add noise to clean speech according to every SNR condition, and clean speech data, said method comprising:
 - a clustering step of clustering all noise-added speech data stored in said noise database;
 - a speech model space generating step of generating a single-tree-structure noisy speech model space based on the result of the clustering performed in said clustering step;
 - a parameter extracting step of extracting a speech feature parameter of input noisy speech to be recognized;
 - a selecting step of selecting an optimum model from the tree-structure noisy speech model space generated in said speech model space generating step; and
 - a linear transformation step of applying linear transformation to the model selected in the selecting step so that the model provides a further increased likelihood.
7. A noise adaptation program for speech recognition

for adapting a speech model for any noise to speech to be recognized in a noisy environment, said speech model being learned by using noise data stored in a noise database which is used to add noise to clean speech according to every SNR condition, and clean speech data, said program comprising:

a clustering step of clustering all noise-added speech data stored in said noise database;
 a speech model space generating step of generating a single-tree-structure noisy speech model space based on the result of the clustering performed in said clustering step;
 a parameter extracting step of extracting a speech feature parameter of input noisy speech to be recognized;
 a selecting step of selecting an optimum model from the tree-structure noisy speech model space generated in said speech model space generating step; and
 a linear transformation step of applying linear transformation to the model selected in the selecting step so that the model provides a further increased likelihood.

Patentansprüche

1. Rauschanpassungssystem für ein Sprachmodell, um ein Sprachmodell für beliebiges Rauschen an Sprache anzupassen, die in einer verrauschten Umgebung erkannt werden soll, wobei das Sprachmodell unter Verwendung von Rauschdaten, die in einer Rausch-Datenbank gespeichert sind, die verwendet wird, um zu reiner Sprache entsprechend jeder SNR-Bedingung Rauschen hinzuzufügen, und unter Verwendung von reinen Sprachdaten gelernt wird, wobei das System umfasst:

Sammelmittel, um alle Rauschdaten, die in der Rausch-Datenbank gespeichert sind, zu sammeln;
 Sprachmodellraum-Erzeugungsmittel, um anhand des Ergebnisses des durch die Sammelmittel ausgeführten Sammelns einen Raum für ein verrauschtes Sprachmodell mit Einzelbaumstruktur zu erzeugen;
 Parameterextraktionsmittel, um einen Sprachmerkmalsparameter von zu erkennender eingegebener verrauschter Sprache zu extrahieren;
 Auswahlmittel, um aus dem durch die Sprachmodellraum-Erzeugungsmittel erzeugten Raum für ein verrauschtes Sprachmodell mit Baumstruktur ein optimales Modell auszuwählen; und
 Lineartransformationsmittel, die auf das durch die Auswahlmittel ausgewählte Modell eine lineare Transformation anwenden, damit das Mo-

dell eine weiter erhöhte Wahrscheinlichkeit ergibt.

2. System nach Anspruch 1, bei dem die Sammelmittel die rauschangereicherte Sprache durch Hinzufügen des Rauschens zu der Sprache in Übereinstimmung mit einer Rauschabstandsbedingung erzeugen, den Sprach-Cepstral-Mittelwert der erzeugten rauschangereicherten Sprache subtrahieren, ein Gaußsches Sprachmodell jedes Teils der erzeugten rauschangereicherten Sprache erzeugen und die Wahrscheinlichkeit zwischen den Teilen der rauschangereicherten Sprache berechnen, um eine Wahrscheinlichkeitsmatrix zu erzeugen, um ein Sammelergebnis zu schaffen.
3. System nach Anspruch 1 oder 2, bei dem die Auswahlmittel ein Modell auswählen, das die höchste Wahrscheinlichkeit für den durch die Parameterextraktionsmittel extrahierten Sprachmerkmalsparameter schafft.
4. System nach Anspruch 3, bei dem die Auswahlmittel ein Modell auswählen, indem sie den Raum für ein verrauschtes Sprachmodell mit Baumstruktur von der höchsten Ebene abwärts zur niedrigsten Ebene durchsuchen.
5. System nach einem der Ansprüche 1 bis 4, bei dem die Lineartransformationsmittel die lineare Transformation anhand des durch die Auswahlmittel ausgewählten Modells ausführen, um die Wahrscheinlichkeit zu erhöhen.
6. Rauschanpassungsverfahren für ein Sprachmodell, um ein Sprachmodell für beliebiges Rauschen an Sprache anzupassen, die in einer verrauschten Umgebung erkannt werden soll, wobei das Sprachmodell unter Verwendung von Rauschdaten, die in einer Rausch-Datenbank gespeichert sind, die verwendet wird, um zu reiner Sprache entsprechend jeder SNR-Bedingung Rauschen hinzuzufügen, und unter Verwendung von reinen Sprachdaten gelernt wird, wobei das Verfahren umfasst:
 einen Sammelschritt zum Sammeln von allen rauschangereicherten Sprachdaten, die in der Rausch-Datenbank gespeichert sind;
 einen Sprachmodellraum-Erzeugungsschritt zum Erzeugen eines Raums für ein verrauschtes Sprachmodell mit Einzelbaumstruktur anhand des Ergebnisses des in dem Sammelschritt ausgeführten Sammelns;
 einen Parameterextraktionsschritt zum Extrahieren eines Sprachmerkmalsparameters von zu erkennender eingegebener verrauschter Sprache;
 einen Auswahlsschritt zum Auswählen eines op-

timalen Modells aus dem in dem Sprachmodellraum-Erzeugungsschritt erzeugten Raum für ein verrauschtes Sprachmodell mit Baumstruktur; und
einen Lineartransformationsschritt zum Anwen- 5
den einer linearen Transformation auf das im Auswahlsschritt ausgewählte Modell, damit das Modell eine weiter erhöhte Wahrscheinlichkeit ergibt.

7. Rauschanpassungsprogramm für die Spracherken- 10
nung, um ein Sprachmodell für beliebiges Rauschen an Sprache anzupassen, die in einer verrauschten Umgebung erkannt werden soll, wobei das Sprachmodell unter Verwendung von Rauschdaten, die in einer Rausch-Datenbank gespeichert sind, die verwendet wird, um zu reiner Sprache entsprechend jeder SNR-Bedingung Rauschen hinzuzufügen, und unter Verwendung von reinen Sprachdaten gelernt wird, wobei das Programm umfasst:

einen Sammelschritt, um alle in der Rausch-Datenbank gespeicherten rauschangereicherten Sprachdaten zu sammeln;
einen Sprachmodellraum-Erzeugungsschritt, 25
um anhand des Ergebnisses des im Sammelschritt ausgeführten Sammelns einen Raum für ein verrauschtes Sprachmodell mit Einzelbaumstruktur zu erzeugen;
einen Parameterextraktionsschritt, um einen Sprachmerkmalsparameter von zu erkennen- 30
der eingegebener verrauschter Sprache zu extrahieren;
einen Auswahlsschritt, um ein optimales Modell aus dem im Sprachmodellraum-Erzeugungsschritt erzeugten Raum für ein verrauschtes Sprachmodell mit Baumstruktur auszuwählen; und
einen Lineartransformationsschritt, um eine li- 40
neare Transformation auf das im Auswahlsschritt ausgewählte Modell anzuwenden, damit das Modell eine weiter erhöhte Wahrscheinlichkeit ergibt.

Revendications

1. Un système d'adaptation au bruit d'un modèle vocal pour l'adaptation d'un modèle vocal pour tout bruit à des paroles à reconnaître dans un environnement bruyant, ledit modèle vocal étant assimilé en utilisant des données sur le bruit stockées dans une base de données sur le bruit qui est utilisée pour ajouter du bruit aux paroles claires dans toutes les conditions de rapport signal sur bruit SNR, et des données vo- 50
cales claires, ledit système comprenant:

des moyens de regroupement pour regrouper

toutes les données de bruit stockées dans ladite base de données sur le bruit;
des moyens de création d'un espace des modèles vocaux pour générer un espace de modèles vocaux bruités à structure arborescente unique basé sur le résultat du regroupement réalisé par lesdits moyens de regroupement;
des moyens d'extraction de paramètres pour extraire un paramètre caractéristique des paroles à partir des paroles bruitées en entrée à reconnaître;
des moyens de sélection pour sélectionner un modèle optimal à partir de l'espace des modèles vocaux bruités à structure arborescente généré par lesdits moyens de création d'un espace de modèles vocaux; et
des moyens de transformation linéaire pour appliquer une transformation linéaire au modèle sélectionné par les moyens de sélection de manière que le modèle fournisse une vraisemblance encore améliorée.

2. Le système de la revendication 1, dans lequel lesdits moyens de regroupement génèrent lesdites paroles à bruit ajouté en ajoutant ledit bruit auxdites paroles en accord avec une condition de rapport signal sur bruit, soustraient la valeur moyenne des paroles Cepstral des paroles à bruit ajouté générées, génèrent un modèle à distribution Gaussienne de chacun des morceaux de paroles à bruit ajouté générées, et calculent la vraisemblance entre les morceaux de paroles à bruit ajouté pour générer une matrice de vraisemblance en vue de fournir un résultat de regroupement.
3. Le système de la revendication 1 ou 2, dans lequel lesdits moyens de sélection sélectionnent un modèle qui fournit la plus haute vraisemblance pour le paramètre caractéristique des paroles extrait par lesdits moyens d'extraction de paramètres.
4. Le système de la revendication 3, dans lequel lesdits moyens de sélection sélectionnent un modèle en cherchant vers le bas dans ledit espace des modèles vocaux bruités à structure arborescente depuis le niveau le plus haut jusqu'au niveau le plus bas.
5. Le système de l'une quelconque des revendications 1 à 4, dans lequel lesdits moyens de transformation linéaire réalisent la transformation linéaire sur la base du modèle sélectionné par lesdits moyens de sélection pour améliorer la vraisemblance.
6. Une méthode d'adaptation au bruit d'un modèle vocal pour l'adaptation d'un modèle vocal pour tout bruit à des paroles à reconnaître dans un environnement bruyant, ledit modèle vocal étant assimilé en utilisant des données sur le bruit stockées dans

une base de données sur le bruit qui est utilisée pour ajouter du bruit aux paroles claires dans toutes les conditions de rapport signal sur bruit SNR, et des données vocales claires, ledit système comprenant:

5

une étape de regroupement regroupant toutes les données vocales à bruit ajouté stockées dans ladite base de données sur le bruit;

une étape de création d'un espace des modèles vocaux générant un espace de modèles vocaux bruités à structure arborescente unique basé sur le résultat du regroupement réalisé dans ladite étape de regroupement;

10

une étape d'extraction d'un paramètre extrayant un paramètre caractéristique des paroles à partir des paroles bruitées en entrée à reconnaître;

15

une étape de sélection sélectionnant un modèle optimal à partir de l'espace des modèles vocaux bruités à structure arborescente généré dans ladite étape de création d'un espace de modèles vocaux; et

20

une étape de transformation linéaire appliquant une transformation linéaire au modèle sélectionné dans l'étape de sélection de manière que le modèle fournisse une vraisemblance encore améliorée.

25

7. Un programme d'adaptation au bruit pour reconnaissance vocale pour l'adaptation d'un modèle vocal pour tout bruit à des paroles à reconnaître dans un environnement bruyant, ledit modèle vocal étant appris en utilisant des données sur le bruit stockées dans une base de données sur le bruit qui est utilisée pour ajouter du bruit aux paroles claires dans toutes les conditions de rapport signal sur bruit SNR, et des données vocales claires, ledit système comprenant:

30

une étape de regroupement regroupant toutes les données vocales à bruit ajouté stockées dans ladite base de données sur le bruit;

40

une étape de création d'un espace des modèles vocaux générant un espace de modèles vocaux bruités à structure arborescente unique basé sur le résultat du regroupement réalisé dans ladite étape de regroupement;

45

une étape d'extraction d'un paramètre extrayant un paramètre caractéristique des paroles à partir des paroles bruitées en entrée à reconnaître;

une étape de sélection sélectionnant un modèle optimal à partir de l'espace des modèles vocaux bruités à structure arborescente généré dans ladite étape de création d'un espace de modèles vocaux; et

50

une étape de transformation linéaire appliquant une transformation linéaire au modèle sélectionné dans l'étape de sélection de manière que le modèle fournisse une vraisemblance encore améliorée.

55

FIG. 1

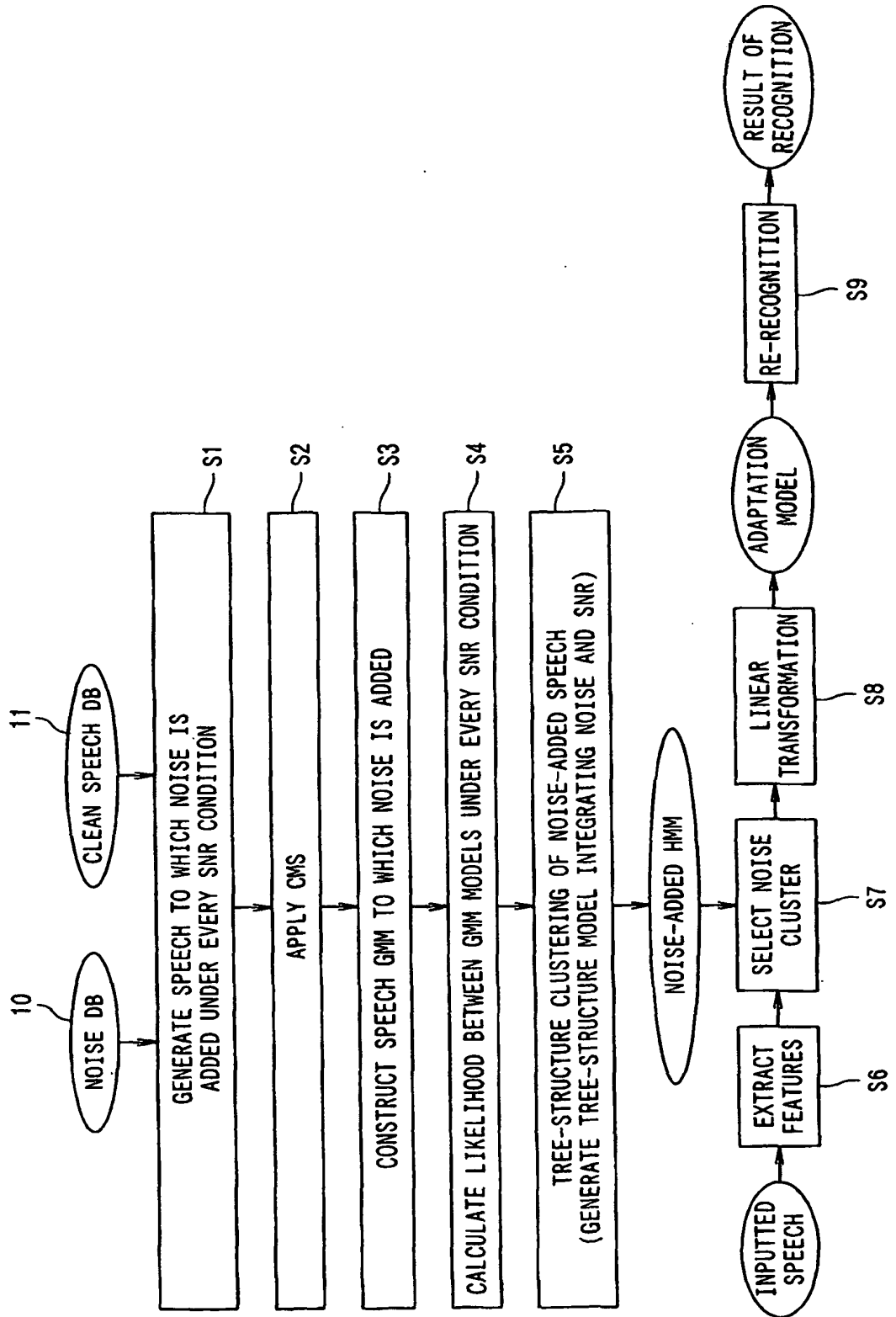


FIG. 2

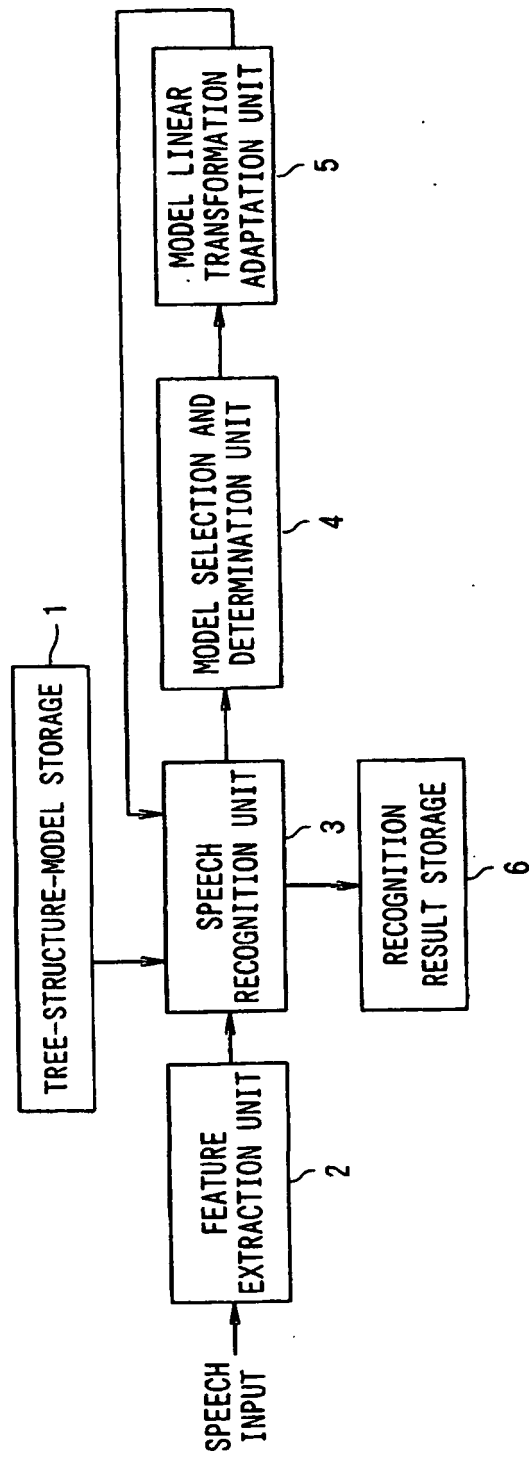


FIG. 3

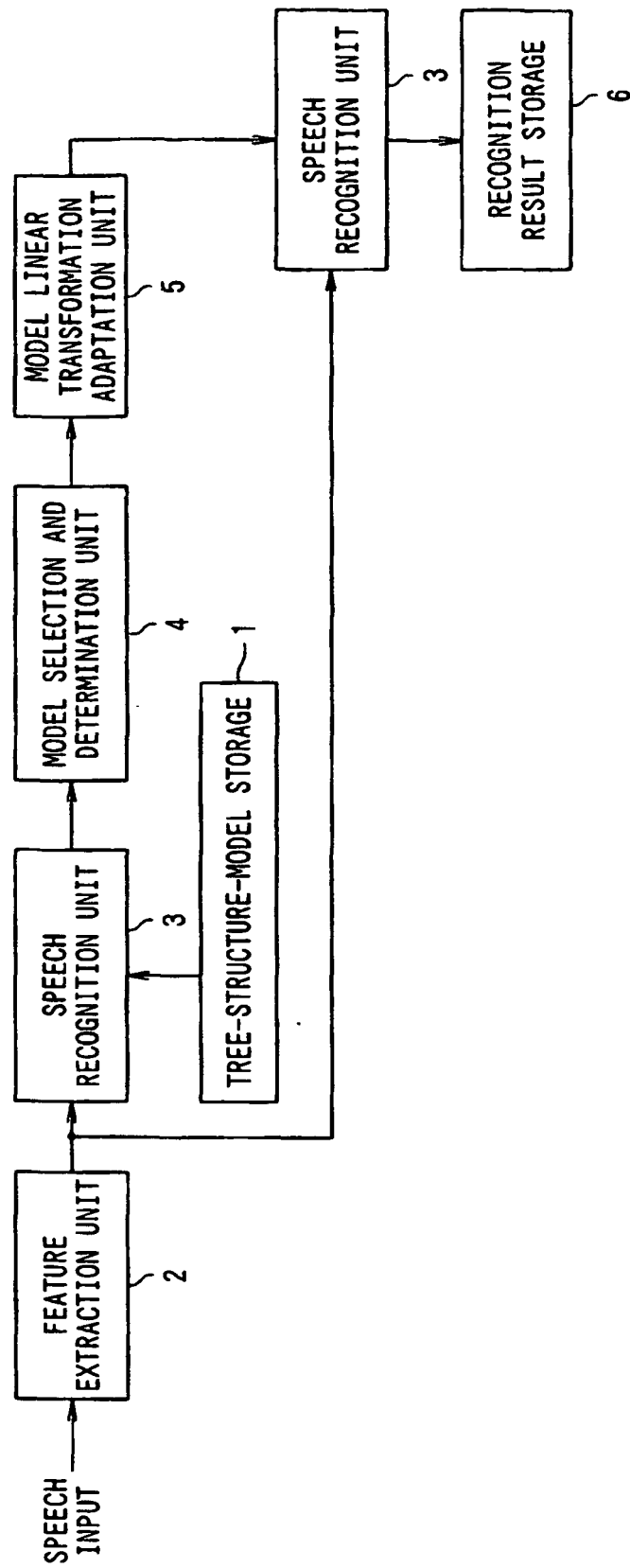


FIG. 4

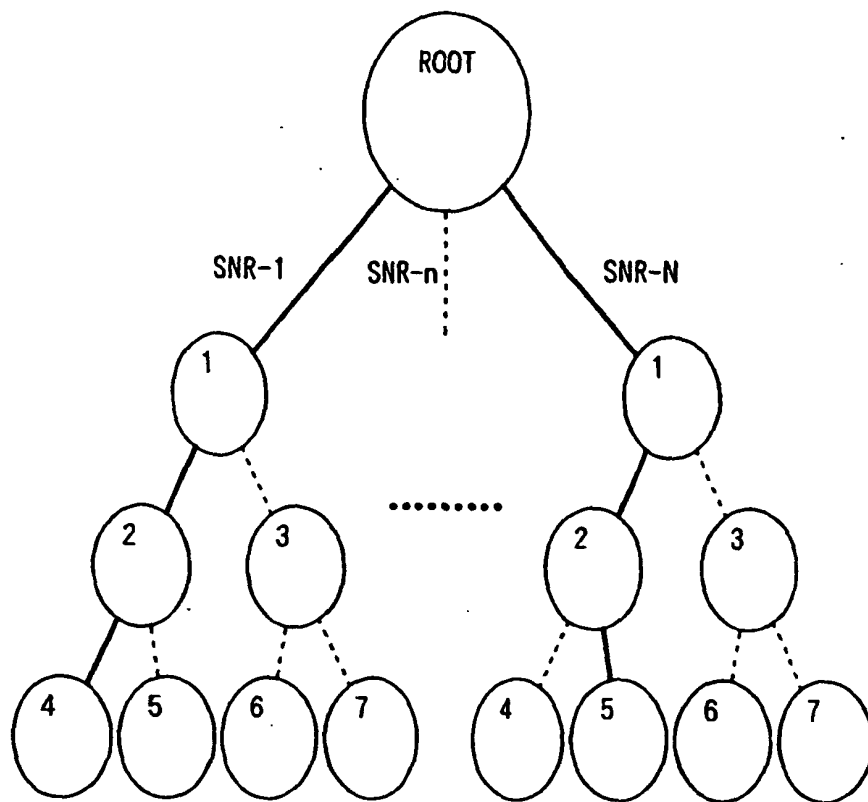


FIG. 5

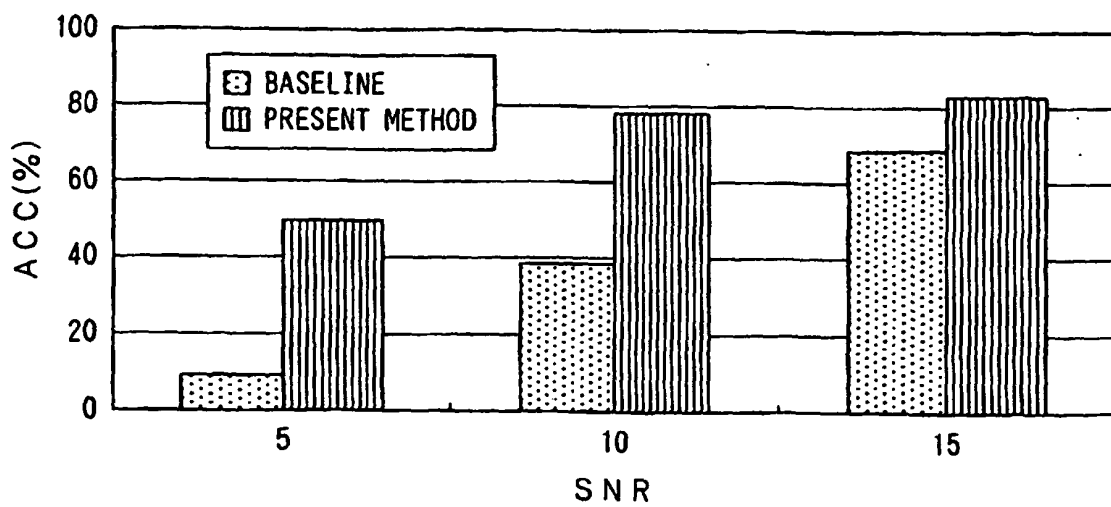


FIG. 6

