



① Veröffentlichungsnummer: 0 612 059 A2

# EUROPÄISCHE PATENTANMELDUNG

(21) Anmeldenummer: 93120010.9 (51) Int. Cl.<sup>5</sup>: **G10L** 3/00, G10L 3/02

2 Anmeldetag: 11.12.93

(12)

Priorität: 23.12.92 DE 4243831

(43) Veröffentlichungstag der Anmeldung: **24.08.94 Patentblatt 94/34** 

Benannte Vertragsstaaten: **DE FR GB** 

Anmelder: Daimler-Benz Aktiengesellschaft
Postfach 80 02 30

D-70546 Stuttgart (DE)

Erfinder: Linhard, Klaus, Dr.-Kasernstrasse 43 D-89231 Neu-Ulm (DE)

Vertreter: Amersbach, Werner, Dipl.-Ing. AEG Aktiengesellschaft Postfach 70 02 20 D-60591 Frankturt (DE)

- (54) Verfahren zur Laufzeitschätzung an gestörten Sprachkanälen.
- © Die Erfindung betrifft ein Verfahren zur Geräuschreduktion in einem Spracherkennungssystem. Es werden die Phasen von zumindest zwei gestörten Signalen geschätzt. Die Phasenschätzung und der für die Geräuschreduktion erforderliche Phasenausgleich wird im Frequenzbereich durchgeführt. Die Hintergrundstörung und das Einschwingverhalten des Raumes werden ständig mitgeschätzt.

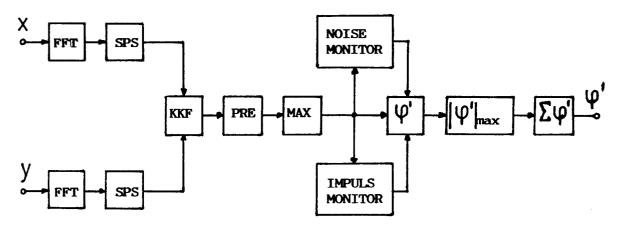


FIG.1

Die Erfindung betrifft ein Verfahren nach dem Oberbegriff des Patentanspruchs 1.

Ein derartiges Verfahren findet Verwendung bei automatischen Spracherkennungssystemen oder für Freisprechanlagen z.B. in Büroräumen, Kraftfahrzeugen etc..

Gestörte Sprache ist besser erfaßbar, wenn sie mit zwei oder mehreren Kanälen aufgezeichnet wird. Der Mensch benutzt zwei Kanäle, seine beiden Ohren. Durch eine psychoakustische Nachverarbeitung wird bei ihm die Richtung des Sprechers ermittelt und die Hintergrundstörung ausgeblendet. Bei technischen Geräten können zwei oder mehrere Kanäle zur Aufzeichnung verwendet werden. Diese Signale können dann mit einer digitalen Signalverarbeitung aufbereitet werden.

Ein wesentlicher Aspekt der mehrkanaligen Verarbeitung ist die Schätzung des Laufzeitunterschiedes der einzelnen Kanäle. Ist der Laufzeitunterschied bekannnt, kann die Richtung des Schallereignisses (Sprecher) ermittelt werden. Die Signale der einzelnen Kanäle können entsprechend laufzeitkorrigiert und weiterverarbeitet werden. Werden z.B. nicht korrigierte Signale zu einem Summensignal zusammengefaßt, können sich einzelne spektrale Anteile des Signals durch Interferenz verstärken, dämpfen oder auslöschen.

Ein Verfahren zur automatischen Ermittlung der Laufzeitunterschiede zweier Mikrofone ist aus einer Veröffentlichung von M. Schlang, ITG-Fachtagung 1988, Bad Nauheim S. 69-73 bekannt. Es arbeitet im Zeitbereich. Jedoch ist dieses Verfahren bei starken Störungen nicht anwendbar.

Der Erfindung liegt deshalb die Aufgabe zugrunde ein Verfahren zur Laufzeitschätzung für ein Spracherkennungssystem anzugeben, das auch bei starken Hintergrundgeräuschen anwendbar ist, für ein Mehrkanalübertragungssystem geeignet ist und zeit- und kostensparend arbeitet.

Die Aufgabe wird gelöst durch die im kennzeichnenden Teil des Patentanspruchs 1 angegebenen Merkmale. Vorteilhafte Ausgestaltungen und/oder Weiterbildungen sind den Unteransprüchen zu entnehmen.

Die Erfindung wird anhand eines Ausführungsbeispiels beschrieben unter Bezugnahme auf schematische Zeichnungen.

In FIG. 1 wird anhand eines Blockschaltbilds die Phasenschätzung erläutert.

FIG. 2 gibt für ein Fahrgeräusch von 140km/h eine Darstellung der Größen  $S_B$ ,  $S_I$ ,  $S_N$  und g in Abhängigkeit von der Zeit an.

In der vorliegenden Erfindung wird ein 2-kanaliger Laufzeitausgleich vorgestellt. Die Erweiterung auf mehrere Kanäle ist mit dem entsprechenden Mehraufwand leicht möglich. Der Laufzeitausgleich ist ein Teil der Signalvorverarbeitung einer mehrkanaligen Geräuschreduktion, die z.B. für einen Spracherkenner im Fahrzeug verwendet werden kann.

Die Laufzeit wird im Frequenzbereich ermittelt. Dies ermöglicht eine einfache Laufzeitkorrektur durch die Multiplikation des Spektrums mit der neuen Phase und führt zu einem geringen Rechenaufwand.

Die Sprach- und Geräuschaufnahmen zur Entwicklung und Bewertung des vorliegenden Verfahrens wurden in einem Fahrzeug mit zwei Mikrofonen durchgeführt. Die Störung ist das Fahrgeräusch bei verschiedenen Fahrsituationen.

Mit dem erfindungsgemäßen Verfahren werden im Frequenzbereich die Phasen an einer Anzahl von Maxima der Kreuzkorrelation bestimmt. Die Hintergrundstörung und das Einschwingverhalten des Raumes werden ständig mitgeschätzt. Die einzelnen Phasenwerte werden nur zu Beginn eines Einschwingvorgangs verarbeitet und wenn das Hintergrundgeräusch um einen gewissen Faktor überschritten wird. Bei der Weiterverarbeitung der Phasenwerte wird eine lineare Phasenbeziehung vorausgesetzt und die Varianz der Schätzung wird bei der Glättung der Werte mitberücksichtigt. Die Berücksichtigung des Einschwingvorgangs des Raumes fuhrt dazu, daß nur bei starken Energieanstiegen der Sprache eine Phasenschätzung stattfindet. Sofort zu Beginn des Wortes steht ein neuer Phasenschätzwert zur Verfügung. Der Einfluß von Reflexionen wird vermindert. Durch die Berücksichtigung des Hintergrundgeräuschs ist das Verfahren für den praktischen Einsatz z.B. im Fahrzeug gut geeignet. Anhand eines Blockschaltbildes in FIG. 1 wird der Verfahrensablauf der Phasenschätzung näher erläutert.

Die Mikrofonsignale x und y werden in den Frequenzbereich transformiert (FFT, Fast Fourier Transformation). Die Transformationslänge wird zu N=256 gewählt. Es ergaben sich die transformierten Segment  $X_I$  (i) und  $Y_I$  (i). I bezeichnet den Blockindex der Segmente, i die diskrete Frequenz (i = 0,1,2,...,N-1). Die Segmente sind halb überlappt und werden mit einem Hanning Fenster gewichtet. (Die Abtastrate der Signale x und y beträgt 12 kHz.)

Im Frequenzbereich wird der Langzeitmittelwert des Betragsspektrums subtrahiert (SPS, spektrale Substraktion). Die Phase der Signale wird nicht verändert. Das Störgeräusch wird reduziert. Es ergeben sich die Schätzwerte  $\widehat{X}$  und  $\widehat{Y}$ . Die SPS ist ein Standardverfahren und kann hier in einer einfachen Version eingesetzt werden. Sind nur geringe Störungen vorhanden, kann auf die SPS ganz verzichtet werden.

Mit der Glättungskonstante  $\beta$  wird das Störspektrum S<sub>nn</sub>(i) geschätzt. Das Störspektrum wird normiert und subtrahiert. I bezeichnet den Blockindex, i die diskrete Frequenz. Als Glättungskonstante wird z.B.  $\beta_1$  =

0.03 verwendet.

$$\hat{S}_{nn,l}(i) = (l-\beta_l)\hat{S}_{nn,l-1}(i) + \beta_l|X_l(i)|^2$$
 (1)

5

$$|\hat{x}_{1}(i)| = |x_{1}(i)| - \frac{\hat{s}_{m,1}(i)}{|x_{1}(i)|}$$
 (2)

10

$$\hat{X}_{1}(i) = [1 - \frac{\hat{S}_{nn}, 1^{(i)}}{|X_{1}(i)|^{2}} \times_{I}(i)$$
(3)

15

Für den zweiten Kanal Y gelten die entsprechenden Gleichungen.

Aus den geschätzten Werten  $\hat{X}$  und  $\hat{X}$  wird der Betrag der Kreuzleistungsdichte  $B_{XY,l}$  berechnet. Der Bereich  $(N_u, N_o)$  liegt z.B. zwischen 300 und 1500Hz  $(N_u = 6, N_o = 31, bei N = 256)$ . Dabei gilt

$$S_{xy,I}(i) = (I-\alpha)S_{xy,I-I}(i) + \alpha \widehat{X}_I(i)\widehat{Y}_I^*(i); N_U \le i \le N_0$$
 (4)

$$B_{xy,l}(i) = |S_{xy,l}(i)|$$
 (5)

25

Als Glättungskonstante  $\alpha$  wird z.B.  $\alpha$  = I gewählt. Werte  $\alpha \ll$  I sind nicht sinnvoll.

Mit einer Präemphase können höhere Frequenzen angehoben werden. Dies ist dann vorteilhaft, wenn das Sprachsignal und das Störsignal bei höheren Frequenzen eine geringere Leistung aufweisen. Die Werte der Kreuzleistung B<sub>xy</sub> (i) können z.B. im Bereich 300 bis 1500 Hz um 10dB linear an-steigend angehoben werden. Die Präemphase kann aber auch schon durch die Mikrofoncharakteristik vorgegeben sein.

Aus den Werten  $B_{xy}$  (i) werden M Maxima bestimmt und summiert. Es können z.B. M = 8 verwendet werden. Es wird ein aktueller Schätzwert

35

45

$$s_{B,1} = \frac{1}{M} \sum_{M=1}^{M} B_{xy,1}(i)$$
 (6)

40 bestimmt.

Über einen Impulsmonitor wird eine "simulierte Impulsantwort"  $S_l$  berechnet. Das Einschwingverhalten des umgebenden Raumes auf plötzliche energiestarke Schallereignisse (Sprache) wird hiermit grob simuliert (z.B. wird  $\gamma = 0.1$  gewählt). Die Glättung des Phasenwerts "vom Wortanfang in das Wort hinein" ist mit  $\gamma$  einstellbar.

$$S_{I,1} = (I - \gamma)S_{I,I-1} + \gamma S_{B,I}$$
 (7)

Außerdem wird über einen Geräuschmonitor eine adaptive Glättungskonstante h berechnet. Mit dieser Glättungskonstanten ergibt sich ein Schätzwert  $S_N$  für die Störung. Wurde zuvor eine spektrale Substraktion (SPS) durchgeführt, ist  $S_N$  ein Schätzwert für die Reststörung. Für die Glättungskonstante  $h_o$  gilt z.B.  $h_o$  = 0.03

$$h_1 = h_0 \xrightarrow{S_{N,1-1}} S_{N,1-1+B,1}$$
 (8)

$$S_{N,L} = (1 - h_I)S_{N,I-1} + h_IS_{B,I}$$
 (9)

Die Phase der gestörten Signale wird aus den Real- und Imaginärteilen von  $S_{xy}$  berechnet. Die Phase wird nur an den M zuvor bestimmten Maxima berechnet.

$$\varphi_{1}(i) = \arctan \frac{Im[S_{xy,1}(i)]}{Re[S_{xy,1}(i)]} ; für Re > 0$$
 (10)

und

5

10

25

40

50

$$\varphi_{1}(i) = \pi - \arctan \frac{-\operatorname{Im}[S_{xy,1}(i)]}{\operatorname{Re}[S_{xy,1}(i)]}; \quad \text{sonst}$$
 (11)

20 Daraus ergibt sich der Phasenanstieg:

$$\varphi_{1}(i) = \frac{\varphi_{1}(i)}{i} \tag{12}$$

Mit der Länge der Fouriertransformation N und der max. zulässigen Verschiebung um n Taps ergibt sich (N = 256):

30 
$$|\phi'|_{\text{max}} = |n|_{N}^{2\pi}$$
 (13)

Übersteigt der Phasenanstieg  $|\phi'|$  an einem der Maxima  $|\phi'|_{max}$ , so wird dieser Wert  $\phi'$  nicht weiterverwendet. Es wird eine adaptive Glättungskonstante g berechnet:

$$g_{1} = \frac{g_{0}(S_{B,1} - S_{I,1})}{S_{I,1}}$$
 (14)

 $g_l \le g_{maX}$  (15)

$$g_{max} = 0.25; g_0 = 0.25$$
 (16)

45 Der aktuelle Wert S<sub>B</sub> muß um den Faktor c größer sein als die simulierte Impulsantwort S<sub>I</sub>

$$S_{B,I} \ge cS_{I,I}$$
;  $c = 2$  (17)

sonst gilt:

$$g_1 = 0$$
 (18)

Der aktuelle Wert  $S_{B}$  muß um den Faktor d größer sein als das Restrauschen  $S_{N}$ 

$$S_{B,I} \ge dS_{N,I}; d = 3$$
 (19)

sonst gilt ebenfalls

$$g_1 = 0$$
 (20)

Ist Gl. (17) oder Gl. (19) nicht erfüllt, d.h. gilt g = O, so kann die Phasenschätzung abgebrochen werden. Es gilt der alte Phasenschätzwert.

Für alle

$$|\phi'|(i)| \le |\phi| |max$$
 (21)

gilt:

10

5

$$m_{\varphi',1} = \frac{1}{M'} \sum_{\varphi'_1(i)} \varphi'_1(i) \qquad (22)$$

15

$$s^2 \varphi', 1 = \frac{1}{M'} \sum_{i=1}^{M'} (\varphi'_1(i))^2$$
 (23)

20

30

35

Von den ursprünglichen M Maxima werden wegen Gl. (21) nur M' für die Gl. (22, 23) verwendet. Ist die Anzahl M' der für die Summen gültigen Werte  $\phi$  kleiner als  $M_{min}$ , gilt der geschätzte Phasenanstieg als zu unsicher oder außerhalb des Nutzbereichs (z.B.  $M_{min}$  = 6, bei M = 8). Die Phasenschätzung wird dann nicht aktualisiert und das Verfahren hier abgebrochen. Es gilt der alte Phasenschätzwert.

Es wird die Varianz der Schätzung berechnet:

$$\sigma^{2}_{\varphi',1} = s^{2}\varphi'_{11} - m^{2}\varphi'_{11}$$
 (24)

Als maximale Varianz wird

$$\sigma^2_{\text{max}} = |\phi'|^2_{\text{max}} \qquad (25)$$

verwendet.

Entsprechend der Varianz wird die Glättungskonstante g gewichtet. Bei einer großen Streuung gilt:

40

$$g_1$$
: = 0.09 \*  $g_1$ ; für 0,2 $\sigma^2_{max}$  <  $\sigma^2_{\phi'_1 1}$  <  $\sigma^2_{max}$  (26)

45

Bei einer mittleren Streuung gilt:

$$g_1$$
: = 0.3\*  $g_1$ ; für 0.02 $\sigma^2_{\text{max}} \le \sigma^2_{\phi',1} \le 0.2\sigma^2_{\text{max}}$  (27)

Bei sehr geringer Streuung gilt:

55

50

$$g_1 := g_1; \text{ für } \sigma^2_{\psi_1^{'}1} < 0.02\sigma^2_{\text{max}}$$
 (28)

Entsprechend den Gl. 19 - 22 wird g in der Regel nur am Wortanfang größer Null sein. Dabei muß die Energie des Wortes größer sein als die Energie des Restgeräusches und der simulierten Impulsantwort. Mit der Variablen j wird die aufeinanderfolgende Anzahl für g > 0 gezählt. Entsprechend gilt für die Glättung:

$$j=2: 
 m_{\varphi',1} = \frac{(m_{\varphi',1} + m_{\varphi',1-1})}{2} 
 (30)$$

$$\widetilde{\varphi}'_{1} = (1 - 1, 5g_{1}) \widetilde{\varphi}'_{1-1} + 1, 5 g_{1}^{m}_{\varphi'_{1}}$$
 (31)

(29)

Wird z.B. infolge einer Störung die Bedingung g > 0 nur einmal in Folge erfüllt, wird die Phasenschätzung nicht aktualisiert. Eine Aktualisierung der Phasenschätzung erfolgt nur dann, wenn g > 0 mindestens 2-mal in Folge erfüllt wird.

Ein Beispiel für die Zwischengrößen  $S_B$ ,  $S_I$ ,  $S_N$ ' und g und die daraus abgeleitete Phasenschätzung zeigt, FIG. 2. Dabei wird das Wort "Senderwahl" gesprochen und das Fahrgeräusch bei 140km/h addiert. Das Verfahren wird, wie oben angegeben, verwendet. Der Phasenschätzwert ist in Abtastwerten n angegeben. Mit der Größe  $S_I$  wird der "Sprachimpuls" teilweise verdeckt und so nur bei starken Energieanstiegen eine Schätzung erlaubt ( $S_B$  muß  $S_I$  um den Faktor 2 übersteigen). Die Schätzung der Reststörung  $S_N$  ermöglicht eine größere Robustheit gegenüber Geräuschen ( $S_B$  muß  $S_N$  um den Faktor 3 übersteigen).

# Patentansprüche

10

15

20

45

50

- Verfahren zur Laufzeitschätzung bei dem Laufzeitunterschiede von geräuschgestörten Signalen von zumindest zwei Sprachkanälen mittels einer Kreuzkorrelation bestimmt werden, dadurch gekennzeichnet,
  - daß im Frequenzbereich die Phasenwerte von zumindest zwei Signalen über eine bestimmte Anzahl von Maxima der Kreuzleistungsdichte ermittelt werden und deren Phasenverschiebung bestimmt wird, und
  - daß der erforderliche Phasenausgleich ebenfalls im Frequenzbereich durchgeführt wird.
  - 2. Verfahren nach Anspruch I, dadurch gekennzeichnet, daß Hintergrundstörungen und das Einschwingverhalten des Raumes bei der Bestimmung der Phasenwerte ständig mitgeschätzt werden.
  - 3. Verfahren nach Anspruch 2, dadurch gekennzeichnet, daß das Hintergrundgeräusch über einen Geräuschmonitor geschätzt wird, und daß ein neuer Phasenwert lediglich dann ermittelt wird, wenn der Schätzwert des Hintergrundgeräusches um einen bestimmten Faktor überschritten wird.
- Verfahren nach Anspruch 2, dadurch gekennzeichnet, daß das Einschwingverhalten des umgebenden Raumes über einen Impulsmonitor derart geschätzt wird, daß lediglich bei starkem Energieanstieg in den Signalen ein neuer Phasenschätzwert ermittelt wird.

- 5. Verfahren nach einem der vorhergehenden Ansprüche, dadurch gekennzeichnet, daß eine lineare Laufzeitverschiebung der Signale angenommen wird.
- 6. Verfahren nach einem der vorhergehenden Ansprüche, da-durch gekennzeichnet, daß eine Glättung des Phasenwertes vom Wortanfang in das gesprochene Wort hinein durchgeführt wird, und daß die Varianz der Schätzung bei der Glättung der Phasenwerte mitberücksichtigt wird.
  - 7. Verfahren nach einem der vorhergehenden Ansprüche, dadurch gekennzeichnet,

5

10

15

20

25

30

45

50

55

- daß zumindest zwei Mikrofonsignale x, y mittels einer FFT (Fast Fourier Tansformation) in den Frequenzbereich transformiert werden,
- daß durch spektrale Substraktion aus den transformierten Signalen die Schätzwerte  $\widehat{X}$ ,  $\widehat{Y}$  bestimmt werden,
- daß aus den geschätzten Werten  $\widehat{X}$ ,  $\widehat{Y}$  der Betrag der Kreuzleistungsdichte  $B_{xy}$  bestimmt wird,
- daß die Maxima der Kreuzleistungsdichte bestimmt werden, und daß aus einer bestimmten Anzahl Maxima der Kreuzleistungsdichte B<sub>xy</sub> ein aktueller Wert S<sub>B</sub> für die gestörten Signale ermittelt wird, daß abhängig vom aktuellen Wert S<sub>B</sub> die Phasen φ der gestörten Signale ermittelt werden und damit der Phasenanstieg φ' bestimmt wird,
- daß der Phasenanstieg φ'geglättet wird, indem über einen Impulsmonitor ein simulierter Sprachimpuls S<sub>I</sub> mit dem aktuellen Wert S<sub>B</sub> der gestörten Signale gekoppelt wird, derart, daß eine
  erneute Phasenschätzung lediglich dann durchgeführt wird, wenn ein starker Energieanstieg des
  Mikrofonsignals registriert wird, und
- daß mit einem Geräuschmonitor ein Schätzwert S<sub>N</sub> für die Hintergrundgeräuschstörung ermittelt wird und mit dem aktuellen Wert S<sub>B</sub> der gestörten Signale gekoppelt wird, derart, daß eine erneute Phasenschätzung lediglich dann durch geführt wird, wenn vom Signal die Hintergrundstörung deutlich überschritten wird.
- 8. Verfahren nach Anspruch 7, dadurch gekennzeichnet, daß ein maximaler Phasenanstieg |\phi'|\_{max} für die Phase an den einzelnen Maxima vorgebbar ist und eine erneute Phasenschätzung lediglich dann durchgeführt wird, wenn der Phasenanstieg um mindestens M' der M Maxima den maximalen Anstieg |\phi'|\_{max} nicht überschreitet.
- 9. Verfahren nach Anspruch 7, dadurch gekennzeichnet, daß die Varianz der Phasenanstiege an den einzelnen Maxima bei der zeitlichen Glättung des Phasenanstiegs berücksichtigt wird.
- **10.** Verfahren nach den Ansprüchen 7 bis 9, dadurch gekennzeichnet, daß eine erneute Phasenschätzung lediglich dann durchgeführt wird, wenn die Bedingungen für einen gültigen Phasenanstieg zeitlich mehrfach in Folge auftreten.
- 11. Verfahren nach einem der vorhergehenden Ansprüche, dadurch gekennzeichnet, daß die gestörte Sprache auf mehr als zwei Sprachkanälen aufgenommen wird und daß die Laufzeitunterschiede der einzelnen Kanäle geschätzt werden.

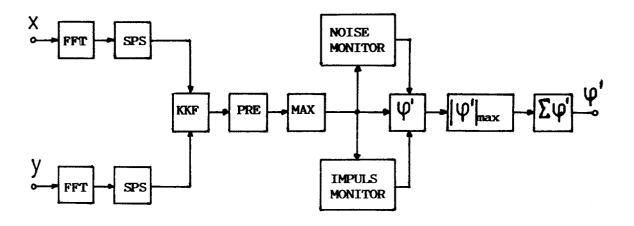


FIG.1

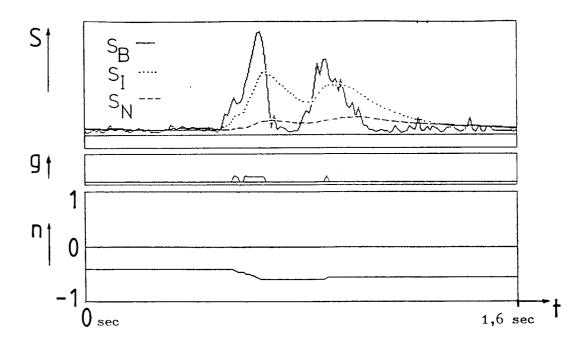


FIG. 2