

**(12) STANDARD PATENT**  
**(19) AUSTRALIAN PATENT OFFICE**

(11) Application No. **AU 2014200407 B2**

(54) Title  
**Method for Voice Activation of a Software Agent from Standby Mode**

(51) International Patent Classification(s)  
**G10L 15/22** (2006.01) **G10L 25/78** (2013.01)  
**G10L 15/26** (2006.01) **G10L 99/00** (2013.01)  
**G10L 15/28** (2013.01)

(21) Application No: **2014200407** (22) Date of Filing: **2014.01.24**

(30) Priority Data

(31) Number	(32) Date	(33) Country
<b>DE102013001219.8</b>	<b>2013.01.25</b>	<b>DE</b>

(43) Publication Date: **2014.08.14**

(43) Publication Journal Date: **2014.08.14**

(44) Accepted Journal Date: **2019.09.19**

(71) Applicant(s)  
**inodyn NewMedia GmbH;Lothar Pantel**

(72) Inventor(s)  
**Pantel, Lothar**

(74) Agent / Attorney  
**Lothar Pantel, c/o George PO Box 665, Berwick, VIC, 3806, AU**

(56) Related Art  
**US 2012/0034904 A1**  
**US 7996228 B2**

## ABSTRACT

### Method for Voice Activation of a Software Agent from Standby Mode

A method for voice activation of a software agent from a standby mode. An audio recording (2) is buffered in an audio buffer (6) and at the same time, the audio recording is input to a secondary voice recognition process (7) which is economical in terms of energy and has an increased false positive rate. When a keyword is recognized, a primary voice recognition process (8) is activated from an inactive state, which converts the audio buffer to text and inputs it to a dialog system (9) which analyzes as to whether there is a relevant question made by the user. If this is the case, the user gets an acoustic reply (3), and if this is not the case, the dialog system and the primary voice recognition process immediately return to the inactive state and transfer the control to the secondary voice recognition process.

Use figure 2.



## **METHOD FOR VOICE ACTIVATION OF A SOFTWARE AGENT FROM STANDBY MODE**

### **CROSS-REFERENCE TO RELATED APPLICATION**

[0001] This application claims priority from German Patent Application DE102013001219.8, filed January 25, 2013, the entire disclosure of which is expressly incorporated herein by reference.

### **TECHNICAL FIELD**

[0002] The present invention relates to the field of voice recognition, in particular to voice-based activation of processes.

### **BACKGROUND OF THE INVENTION**

[0003] Voice recognition, that is, the conversion of acoustic speech signals to text, concretely, the conversion to a digital text representation by means of character encoding, is known. It is possible to control systems without haptic operation. The methods and systems of US patent 8,260,618 and US patent 7,953,599 describe how devices can be controlled or also activated by voice.

[0004] Owing to their small size, the ergonomics of smartphones, i.e. mobile telephones with computer functionality, is very restricted when they are operated by touchscreen. An alternative is personal assistant systems where the smartphone can be controlled with voice commands, in part also with natural speech without special control commands. A known example is the "Siri" system in the "iPhone" from Apple (source: <http://www.apple.com>). A personal assistant system can be an independent application ("app") on the smartphone or be integrated in the operating system. Voice recognition, interpretation and reaction can be done locally on the hardware of the smartphone. But because of the greater processing power an Internet-based server network ("in the cloud") is normally used, with which the personal assistant system communicates, i.e. compressed voice or sound recordings are sent to the server or server network and the verbal reply generated by voice synthesis is streamed back to the smartphone.

[0005] Personal assistant systems are a subset of software agents. There are various options for interaction: e.g. retrieval of facts or knowledge, status updates in social networks or dictation of emails. In most cases, a dialog system (or a so-called chatbot) is used for the personal assistant system which operates partly with semantic analysis or approaches from artificial intelligence to simulate a virtually realistic conversation about a topic.

[0006] Another example of a personal assistant is the system designated as "S voice" on the "Galaxy S III" smartphone from Samsung (source: <http://www.samsung.com>). This product has the option of waking up the smartphone from a standby or sleep state, namely by means of a voice command, without touching the touchscreen or any key. For this purpose the user can store a spoken phrase in the system settings which is used for waking up. "Hi Galaxy" has been factory set. The user must explicitly activate the acoustic monitoring and again deactivate it later because the power consumption would be too great for a day-long operation. According to the manufacturer, the system is provided for situations in which manual operation is not an option, e.g. while driving. By way of example, the driver gives the verbal command "Hi Galaxy", to which, depending on the setting, the "S voice" replies with the greeting: "What would you like to do?" Only now, in a second step, and after the user has already lost productive time due to his first command and waiting for the wake up time - including the greeting - he can actually ask e.g. "What is the weather like in Paris?"

[0007] By storing a limited number of further phrases in the control panel very simple actions can be activated by voice. By means of the command "take a picture" the camera app could be started. It is, however, not possible to ask the smartphone or rather the "S voice" complex questions or request complex actions from the smartphone, as long as the system is in the standby or sleep state. A question such as "Will I need a raincoat in Paris the day after tomorrow?", cannot be answered by the system from the standby or sleep state in spite of the acoustic monitoring. It has to be explicitly awakened for this purpose.

[0008] The voice activation technology used in the "Galaxy S III" smartphone is from Sensory Inc. (source: <http://www.sensoryinc.com>). The manufacturer emphasizes the extremely low false positive rate on acoustic monitoring by means of their "TrulyHandsFree" technology. "False positive" means falsely interpreting other noise as a phrase and the undesired initiation of the trigger. The manufacturer restricts his descriptions to the sequential process during which the device is first brought to life by means of a keyword, only then to be controlled via further commands. Quote: "TrulyHandsFree can be always-on and listening for dozens of keywords that will bring the device to life to be controlled via further voice commands." No other procedure is disclosed.

## **SUMMARY OF THE INVENTION**

[0009] The object underlying the present invention is to provide a method or system which permits asking a software agent, which is in a standby or sleep state, questions via "natural" voice, whereby the system should reply or respond without further interposed interaction steps.

[0010] According to the present invention, the object mentioned above is attained by means of the features of independent claims 1 and 13. Advantageous embodiments, possible alternatives, and optional functionalities are specified in the dependent claims. According to the claims, a system for voice activation of a software agent from a standby mode may comprise at least one microphone, at least one output device, at least one audio buffer and a hardware infrastructure that is able to execute a primary voice recognition process, a secondary voice recognition process and a dialog system process, wherein said hardware infrastructure may be configured to:

- a) capture audio data by means of said at least one microphone,
- b) continually buffer said audio data in said audio buffer, such that said audio buffer always contains the audio data of a time period up to the present,
- c) input said audio data to said secondary voice recognition process, which, on recognizing a keyword or a phrase triggers the start or the activation of said primary voice recognition process,
- d) wherein said primary voice recognition process converts the entire or most recent content of said audio buffer as well as subsequent live audio data to text, said text being handed over to said dialog system process which likewise starts or is activated and which analyzes the content of said text as to whether it contains a question, a message or a command made by the user to said software agent, in which case, if it is answered in the affirmative, said dialog system process triggers an appropriate action or generates an appropriate reply and contacts the user via said output device and
- e) otherwise, if there is only background noise in said audio data or if said text does not contain any relevant or any evaluable content, said primary voice recognition process and said dialog system process return to the inactive state or terminate and again return the control to said secondary voice recognition process.

[0011] For instance, a software agent or a personal assistant system is in a power-saving standby mode or sleep state, the ambient noise - for example voice - picked up by one or more microphones being digitized and continually buffered in an audio buffer, so that the audio buffer constantly contains the ambient noises or voice from the most recent past. Apart from that, the digitized ambient noise or voice that is picked up by the microphone (or several microphones) is input without significant delay to an energy saving secondary voice recognition process, which, on recognition of a keyword or a phrase from a defined keyword- and phrase-catalog, starts a primary voice recognition process or activates it from an inactive or sleep state.

[0012] The more energy-intensive, primary voice recognition process now converts either the entire audio buffer or the most recent part starting at a recognized voice pause (which typically characterizes the beginning of a question phrase) into text, the primary voice recognition process then seamlessly continuing the conversion of the live transmission from the microphone. The text generated via voice recognition, from the audio buffer as well as from the subsequent live transmission, is input to a dialog system (or chatbot). This dialog system process is likewise started or activated from a sleep state or inactive state.

[0013] The dialog system analyzes the content of the text as to whether it contains a question, a message, and/or a request made by the user to the software agent or to the personal assistant system, for example, by means of semantic analysis.

[0014] If a request or a topic is recognized in the text, which the software agent or personal assistant system is competent for, an appropriate action is initiated by the dialog system, or an appropriate reply is generated and communicated to the user via an output device (e.g. loudspeaker and/or display). The software agent or personal assistant is now in full regular operation and interacting with the user.

[0015] However, if the analyzed text (from the audio buffer and the subsequent live transmission) does not contain any relevant or evaluable content, by way of example, when the text string is empty or the dialog system cannot recognize any sense in the word arrangement, the dialog system process and the primary voice recognition process is immediately returned to the sleep state or terminated in order to save power. The control then again returns to the secondary voice recognition process which monitors the surrounding noise or the voice for further keywords or phrases.

## **BRIEF DESCRIPTION OF THE DRAWINGS**

[0016] Further features, advantages, and possible applications will be apparent from the description of the drawings. All described and/or illustrated features, alone or in any combination, independent of the synopsis in individual claims, constitute the subject matter of the invention.

FIG. 1 shows a smartphone with microphone and loudspeaker on which a personal assistant may run as software.

FIG. 2 is a data flow diagram of a basic method.

FIG. 3 is a schematic diagram of the time flow of a process on a time axis  $t$ ; the keyword in the center of the text sample is "what".

FIG. 4 shows an embodiment in which the primary voice recognition process (executed on a processor) as well as the secondary voice recognition process (implemented as a hardware circuit) are located in the local terminal.

FIG. 5 shows an embodiment in which the primary voice recognition process as well as the secondary voice recognition process are executed on the same single core or multi-core processor.

FIG. 6 shows an embodiment in which the secondary voice recognition process is located in the local terminal, and in which the primary voice recognition process is executed on the processor of a server that is connected via a network.

FIG. 7 is a flowchart of an example method; the method supports, inter alia, the recognition of the beginning and end of a sentence, and the recognition of irrelevant audio recordings.

## DETAILED DESCRIPTION OF THE INVENTION

[0017] A terminal can be a mobile computer system or a stationary, cable-based computer system. The terminal is connected to a server via a network and communicates according to the client-server model. Mobile terminals are connected to the network via radio. Typically, the network is the Internet.

[0018] FIG. 1 depicts a smartphone which represents the terminal 1. The software of a personal assistant system runs on this terminal 1. The terminal 1 has a device for digital audio recording and reproduction, typically, one or more microphones 2 and one or more loudspeakers 3 together with the corresponding A/D-converter 5 and D/A-converter circuits. During regular full operation, the digital audio recording 11 (ambient noise or voice) is input to a primary voice recognition process 8. Depending on the embodiment, the primary voice recognition process 8 can be realized in software or as a hardware circuit. In addition, depending on the embodiment, the primary voice recognition process 8 can be located in the local terminal 1 or on a server 28, the digital audio recording then being continually transmitted via the network 29 to the server 28.

[0019] A typical embodiment uses the server 28 for the primary voice recognition process 8, said primary voice recognition process 8 being implemented in software.



[0020] The primary voice recognition process 8 is a high-grade voice recognition technique, which converts the acoustic information to text 13 as completely as possible during the dialog with the user and typically uses the entire supported vocabulary of the voice recognition system. This operating state is designated as full operation. Prior or after the dialog with the user, the terminal 1 can switch to a sleep state or standby mode to save energy.

[0021] Apart from voice recognition for full operation, the system has a second voice recognition process for the sleep state or standby mode. This secondary voice recognition process 7 is optimized for a low consumption of resources and, depending on the embodiment, can likewise be implemented in software or as a hardware circuit. When designed as hardware, attention should be paid to low power consumption, and when implemented in software, attention should be paid to a low demand on resources, like the processor or RAM. Depending on the embodiment, the secondary voice recognition process 7 can be realized on the local terminal 1 or on the server 28, the digital audio recording 11 then being transmitted to the server 28. In a power-saving embodiment the voice recognition in standby mode is done on the local terminal 1, the secondary voice recognition process 7 being realized as a FPGA (field programmable gate array) or as an ASIC (application specific integrated circuit) and optimized for low power consumption.

[0022] In order for a low consumption of resources by the secondary voice recognition process 7 to be possible, it has a very limited vocabulary. The secondary voice recognition process 7 can thus only understand a few words or short segments from idiomatic expressions (phrases).

[0023] These keywords 18 and phrases should be selected such that they contain the typical features when contacting or asking a question to the personal assistant system. The selected keywords 18 and phrases need not necessarily be at the beginning of a sentence. For example all keywords 18 and phrases to infer a question are suitable: e.g. "do you have", "have you got", "are there", "do I need", "do I have".

[0024] With reference to FIG. 2, in the standby mode, all incoming audio signals 11 are buffered in an audio buffer 6 for a certain time. Random-Access Memory (RAM) may be used for this purpose. If the secondary voice recognition process 7 is located in the terminal 1, the audio buffer 6 should also be located in the terminal 1. If the standby voice recognition is server-based, the audio buffer 6 should also be managed by the server 28.

[0025] As soon as the secondary voice recognition process 7 recognizes a potentially relevant keyword 18 or a phrase, e.g. "do you know", it arranges the temporary wakeup 12 of the primary voice recognition process 8 and a switch to full operation takes place. The content 21 of the audio buffer 6 is now handed over to the primary voice recognition process 8.

[0026] In one embodiment, the audio buffer 6 is located in the RAM of terminal 1. If the primary voice recognition process 8 is also located on the terminal 1, accessing the audio buffer 6 in the RAM will be sufficient. If the primary voice recognition process 8 is executed on the server 28, the content 21 of the audio buffer 6 is now transferred to the server 28 via the network 29.

[0027] The primary voice recognition process 8 now has the past of a potential conversation available via the audio buffer 6. The primary voice recognition process 8 must be able to process the audio data 11 with high priority: The objective is to promptly empty the audio buffer 6 in a timely way in order to again process live audio data 22 as soon as possible; see FIG. 3 and the corresponding list of reference numerals. The result of the primary voice recognition process 8 is the spoken text 13 from the recent past up to the present.

[0028] This text 13 is now input to the dialog system 9 which, by means of semantic analysis or also artificial intelligence, analyzes to what extent a query to the personal assistant system actually exists. It is also possible that the keyword 18 recognized by the secondary voice recognition process 7 does no longer appear in the current text 13 because the voice recognition during full operation (primary voice recognition process 8) is of a higher quality and the secondary voice recognition process 7 was therefore wrong.

[0029] In all cases in which the audio recording 21 (located in the audio buffer 6) and the subsequent live audio data 22 turns out to be irrelevant, the dialog system 9 arranges an immediate return to the standby mode, in particular if there is only background noise or if the meaning of the text 13 is not recognized by the dialog system 9; see the flowchart in FIG. 7 and the corresponding list of reference numerals.

[0030] If the dialog system 9, however, concludes that the question, message, or request contained in the audio buffer 6 is relevant, the terminal 1 remains in full operation and the dialog system 9 will interact with the user. As soon as there are no more queries or messages from the user, the terminal 1 may again switch to standby mode and thus transfers control to the secondary voice recognition process 7.

[0031] Additional embodiments are described in the following. Alternatives or optional functions are also mentioned in some cases:

[0032] In one embodiment, after recognizing a keyword 18 or a phrase by the secondary voice recognition process 7, first of all the audio buffer 6 is scanned for the beginning of the sentence with the question, message, or request. In most cases, as illustrated in FIG. 3, it can be assumed that there is a short fraction of time without voice (that is to say with relative silence with respect to the ambient noise) before the beginning of a sentence because most people make a short pause 16 when they want to give the personal assistant a concrete, well formulated question, message or request.

[0033] In order to find the beginning of a sentence the audio buffer 6 is scanned backward in time starting at the position in time of the recognized keyword 18 or phrase until a period is found that can be interpreted as a silence 16. Typically, this period with the speech pause 16 should have a duration of at least one second. As soon as such a position with a relative silence 16 is found and thus the probable beginning of a sentence is established, the subsequent content 17 of the audio buffer 6 is then handed over to the primary voice recognition process 8, which is started or activated next to generate the text 13.

[0034] If during the evaluation of the text 13 the dialog system 9 does not recognize any meaning in the text 13, possibly because the beginning of the sentence was incorrectly interpreted, there can be a second, optional step: The entire content 21 of the audio buffer 6 can be converted to text 13 together with the subsequent live transmission 22 and be analyzed by the dialog system 9.

[0035] If it is not possible to localize a position of relative silence 16 in the entire audio buffer 6 then probably there is no question, message, or request to the personal assistant system, but interfering noise or a conversation between people. In this case, as shown in the flowchart in FIG. 7, there is no need to start or activate the primary voice recognition process 8.

[0036] In order for a user not to have to wait excessively long for a reply or action, it is advantageous that after activation 12 via a keyword 18 or via phrase, the primary voice recognition process 8 is executed with high priority and completed in a short time, as illustrated by means of the dotted lines 23 and 24 in FIG. 3.

[0037] Since according to the present invention, a full-fledged voice recognition is realized by the primary voice recognition process 8, the secondary voice recognition process 7 can

have an increased false positive rate when recognizing keywords 18 or phrases. That is to say the trigger 12 of the secondary voice recognition process 7 may react very sensitive and during monitoring the ambient noise, overlooking a keyword 18 or phrase is extremely rare. If other noises or other words are falsely interpreted as keywords 18 or phrases, these errors are then corrected by the primary voice recognition process 8: As soon as the faulty trigger 12 is recognized, the primary voice recognition process 8 is immediately terminated or deactivated again.

[0038] The highly reduced recognition performance of the secondary voice recognition process 7 makes it possible to design it as especially energy saving; by way of example, as software running on a slow clocked processor with low power consumption, or on a digital signal processor that is likewise optimized for low power consumption. An FPGA or an ASIC, or, in general, an energy saving hardware circuit 25 is suitable, too. (See FIG. 4)

[0039] With reference to FIG. 5, in case the primary voice recognition process 8 as well as the secondary voice recognition process 7 is running on the local hardware 1, they can both run on the same single core or multi-core processor 27, the secondary voice recognition process 7 running in an especially resource conserving mode of operation with low memory requirements and low power consumption.

[0040] Alternatively, the primary voice recognition process 8 and the dialog system 9 may run on an external server 28 or on a server network, as shown in FIG. 6. In this connection, the entire content 21 or the most recent content 17 of the audio buffer 6, and subsequently also the live transmission 22 is transferred to the server 28 or server network via a network 29 or radio network. Typically, the network 29 is the Internet.

[0041] With continued reference to FIG. 6, after a voice activation 12 triggered by the secondary voice recognition process 7 a latency or transmission delay will occur as soon as the content 17 of the audio buffer 6 has to be transferred via the network 29 to the server 28 or server network, so that the primary voice recognition process 8 and the dialog system 9 can evaluate the content. In order to prevent such a transmission delay, an "anticipatory standby mode" can be used: As soon as the presence of a user is detected, the "anticipatory standby mode" transfers the content 21 of the audio buffer 6 and the ensuing live transmission 22 of the ambient noise or voice to the external server 28 or server network. The audio data 11 are temporarily stored there, so that in the event of a voice activation 12, the primary voice recognition process 8 can access the audio data 11 almost without latency.

[0042] Furthermore, in the "anticipatory standby mode", the secondary voice recognition process 7 can optionally intensify the monitoring of the ambient noise for keywords 18 or phrases.

[0043] The presence of a user can be assumed when there are user activities; by way of example, input via a touchscreen 4 or movements and changes in the orientation of the terminal 1, which are detected by means of acceleration- and position-sensors. It is likewise possible to recognize changes in brightness by means of a light sensor, to recognize changes in position by means of satellite navigation (e.g. GPS), and to perform face recognition by means of a camera.

[0044] Keywords 18 and/or phrases supported by the secondary voice recognition process 7 may be stored in a keyword- and phrase-catalog and may include:

- Question words and question phrases: e.g. "who has", "what", "how is", "where is", "are there", "is there", "are there", "do you know", "can one".
- Requests and commands: By way of example: "Please write an email to Bob". The phrase "write an email" will be recognized. Another example: "I would like to take a picture". The phrase "take a picture" will be recognized.
- Nouns referring to topics on which there is information in the database of the dialog system: e.g. "weather", "appointment", "deadline", "football", "soccer".
- Product names, nicknames and generic terms for a direct address of the personal assistant system. Examples of generic terms: "mobile", "mobile phone", "smartphone", "computer", "navigator", "navi".

[0045] Using a product name as a keyword has the advantage that compared to a catalog with question words, the frequency at which the system unnecessarily changes to full operation can be reduced. When using a product name, it can be assumed that the personal assistant system is in charge. Example: "Hello, <product name>, please calculate the square root of 49", or "What time is it, <product name>?"

[0046] In an advantageous embodiment, the keyword- and phrase-catalog can be modified by the user. If the voice activation is done via the product name or a generic term, the user could, for example, define a nickname for the terminal 1 as a further, alternative keyword.

[0047] The user could also delete some keywords or phrases from the catalog, e.g. if the personal assistant system should report less frequently or only in relation to certain topics.

[0048] As soon as the secondary voice recognition process 7 has recognized a keyword 18 or a phrase, the user has to wait for a few moments until the primary voice recognition process 8 and the dialog system 9 have generated a reply or response. Therefore, in a further embodiment, on recognition of a keyword 18 or phrase by the secondary voice recognition process 7, an optical, acoustic and/or haptic signal is output to the user, for example, a short beep through the loudspeaker 3, a vibration of the terminal 1, an indication on the display 4 or by turning on the backlight of the display 4. The user is then informed that his/her query has reached the terminal 1. At the same time, this type of signaling is only minimally disturbing in case the keyword 18 or the phrase was erroneously recognized. In this case, if no relevant or evaluable content can be recognized in the audio buffer 6 or in the resulting text 13, it is advantageous to output a further optical, acoustic or haptic signal that is conveniently different from the first signal, by way of example, a double beep (first high, then low) or by turning off the backlight of the display 4 that had previously been turned on.

[0049] In another embodiment, the personal assistant system can distinguish different voices or speakers, so that only questions, messages, and requests coming from an entitled person are replied by the dialog system 9, by way of example, only questions by the user. As the primary voice recognition process 8 has a considerably greater recognition performance, only this process may be able to distinguish different speakers by their voice, whereas the secondary voice recognition process 7 may not be able to distinguish different speakers.

[0050] Given a keyword 18 or phrase spoken by a still unidentified speaker, the secondary voice recognition process 7 will arrange the execution of the primary voice recognition process 8. The primary voice recognition process 8 recognizes from the speaker's voice whether he/she is entitled to use the personal assistant system. If a corresponding entitlement is not available, the primary voice recognition process 8 terminates itself or returns to the inactive state, and the control is again passed to the secondary voice recognition process 7. During this procedure, the dialog system 9 can remain in the inactive or sleep state.

[0051] In an optional embodiment, the dialog system 9 takes the context of a conversation into consideration: A conversation between people is monitored and a keyword 18 or a phrase from the keyword- and phrase-catalog appears in the conversation (e.g. "soccer"), so that the primary voice recognition process 8 and the dialog system 9 is started or activated. The dialog system 9 checks if it is competent for the content 21, 22 of the current

conversation, in particular, whether a question, message, or request was made to the personal assistant system. If the dialog system 9 is not in charge, the dialog system 9 stores the context and/or topic and/or keywords or phrases for later reference and returns to the sleep state together with the primary voice recognition process 8. If the dialog system 9 is again started or activated by another keyword 18 or phrase (e.g. "who") at a later time, the previously stored information can be considered as a context. In accordance with the above example, the question "Who won the match today?" can be replied with the soccer results of the current match day.

[0052] It is also possible to repeatedly perform a voice recognition within the primary voice recognition process 8. In the first instance, the voice recognition could be done with an especially quick algorithm that reduces the user's waiting time. In case the resulting text 13 is not valid for the dialog system 9 or cannot be evaluated, the content in the audio buffer 6 can again be converted to text 13 by means of one or more different voice recognition methods, which e.g. are particularly resistant to background noise.

#### **LIST OF REFERENCE NUMERALS**

- 1 Smartphone (Terminal)
- 2 Microphone
- 3 Loudspeaker
- 4 Display / Touchscreen
- 5 Analog-Digital Converter (A/D)
- 6 Audio Buffer
- 7 Secondary Voice Recognition Process
- 8 Primary Voice Recognition Process
- 9 Dialog System
- 10 Analog Microphone Signals
- 11 Digital Audio Signals
- 12 Activation Signal (Trigger) After Recognizing A Keyword
- 13 Text (Digital Representation by Means of Character Coding)
- 14 Reply or Response of the Dialog System
- 15 Audio Recording of the Previously Spoken Sentence in the Audio Buffer
- 16 Audio Recording of the Speech Pause (Silence)
- 17 Audio Recording of the Current Sentence (First Part) in the Audio Buffer
- 18 Recognized Keyword or Phrase
- 19 Live Transmission of the Current Sentence (Second Part)

- 20 Start of the Dialog System
- 21 Audio Data of the Most Recent Past in the Audio Buffer
- 22 Live Transmission of the Audio Data
- 23 Processing Delay Relative to the Beginning of the Sentence
- 24 Reduced Processing Delay at the End of the Sentence
- 25 Hardware Circuit (Digital Signal Processor, FPGA or ASIC)
- 26 Main Processor
- 27 Single Core or Multi-Core Processor with Power Saving Function
- 28 Server or Server Network
- 29 Network (e.g. Radio Network, Internet)
- 30 Digitize Microphone Signals via A/D Converter
- 31 Buffer Live Audio Data in the Audio Buffer
- 32 Execute Secondary Voice Recognition Process with Live Audio Data
- 33 Keyword or Phrase Found?
- 34 Scan Audio Buffer Backward for a Speech Pause
- 35 Was the Speech Pause Found?
- 36 Start/Activate Primary Voice Recognition Process and Dialog System
- 37 Apply Primary Voice Recognition Process to Audio Buffer, Starting at the Speech Pause
- 38 Apply Primary Voice Recognition Process to New Live Audio Data
- 39 Speech Pause at the End of Sentence Found?
- 40 Analyze the Text of the Sentence by means of the Dialog System
- 41 Does the Text Contain A Relevant Question, Message, or Command?
- 42 Generate Reply or Activate Action/Response (Full Regular Operation)
- 43 Are there Further Questions/Commands by the User? (Full Regular Operation)
- 44 Terminate/Deactivate Primary Voice Recognition Process and Dialog System



**CLAIMS**

1. A method for voice activation of a software agent from a standby mode, comprising:
  - a) providing at least one microphone, at least one output device, at least one audio buffer and a hardware infrastructure that is able to execute a primary voice recognition process, a secondary voice recognition process and a dialog system process,
  - b) capturing audio data by means of said at least one microphone,
  - c) continually buffering said audio data in said audio buffer, such that said audio buffer always contains the audio data of a time period up to the present,
  - d) inputting said audio data to said secondary voice recognition process, which, on recognizing a keyword or a phrase triggers the start or the activation of said primary voice recognition process,
  - e) wherein said primary voice recognition process converts the entire or most recent content of said audio buffer as well as subsequent live audio data to text, said text being handed over to said dialog system process which likewise starts or is activated and which analyzes the content of said text as to whether it contains a question, a message or a command made by the user to said software agent, in which case, if it is answered in the affirmative, said dialog system process triggers an appropriate action or generates an appropriate reply and contacts the user via said output device and
  - f) otherwise, if there is only background noise in said audio data or if said text does not contain any relevant or any evaluable content, said primary voice recognition process and said dialog system process return to the inactive state or terminate and again return the control to said secondary voice recognition process.
2. The method according to claim 1, further comprising scanning said audio buffer backwards, beginning at the position in time of the recognized keyword or phrase until a period is found that can be interpreted as a speech pause, the most recent content of said audio buffer, beginning at the position with the recognized speech pause, being handed over to said primary voice recognition process.
3. The method according to claim 2, wherein said primary voice recognition process remains in the inactive state, if no speech pause is found in said audio buffer in a range beginning at said position in time of the recognized keyword or phrase up to the oldest entries.

4. The method according to any one of claims 1 to 3, wherein said secondary voice recognition process has a false positive rate on recognition of keywords and/or phrases higher than said primary voice recognition process, while the interplay between said secondary voice recognition process and said primary voice recognition process corrects false positive error of said secondary voice recognition process.
5. The method according to any one of claims 1 to 4, wherein said secondary voice recognition process consumes less power than said primary voice recognition process.
6. The method according to any one of claims 1 to 5, wherein said primary voice recognition process and said dialog system process are executed on an external server or on a server network, the entire or the most recent content of said audio buffer being transferred via a network and/or radio network to said server or server network.
7. The method according to claim 6, further comprising switching said software agent to an anticipatory standby mode as soon as the presence of the user is detected by means of a sensor, while, during said anticipatory standby mode, the entire or the most recent content of said audio buffer and/or live audio data is continually transferred via said network to said external server or server network and buffered there, whereby, in case of voice activation said primary voice recognition process can access the buffered audio data almost latency-free.
8. The method according to any one of claims 1 to 7, further comprising intensifying the monitoring of said audio data for keywords and/or phrases by said secondary voice recognition process as soon as the presence of the user is detected by means of a sensor.
9. The method according to claim 7 or 8, wherein said sensor is:
  - a) a user interface for user input  
and/or
  - b) an acceleration- and/or position-sensor measuring movement or changes in position  
and/or
  - c) a light sensors measuring changes in the brightness  
and/or
  - d) a satellite navigation sensor measuring changes in position  
and/or
  - e) a camera for face recognition.

10. The method according to any one of claims 1 to 9, further comprising outputting an optical, acoustic and/or haptic signal to the user by means of an output device as soon as a keyword or a phrase is recognized by said secondary voice recognition process.
11. The method according to any one of claims 1 to 10, wherein said primary voice recognition process can distinguish different speakers by their voice by means of an acoustic model, and wherein said secondary voice recognition process cannot distinguish different speakers, whereby said secondary voice recognition process triggers the execution of said primary voice recognition process as soon as a keyword or a phrase from any speaker is detected by said secondary voice recognition process, said primary voice recognition process establishing from the speaker's voice whether the speaker is entitled to utilize said software agent and if there is no entitlement, said primary voice recognition process is terminating or returning to the inactive state, and again passing on the control to said secondary voice recognition process.
12. The method according to any one of claims 1 to 11, wherein, in the event that said dialog system process is not competent for a question, message or command in said audio data, converted to text by said primary voice recognition process, said dialog system process stores the context and/or the topic and/or the keywords or phrases on a storage means, and the stored information is taken into consideration at least on one of the subsequent reactivations of said dialog system process.
13. A system for voice activation of a software agent from a standby mode, comprising at least one microphone, at least one output device, at least one audio buffer and a hardware infrastructure that is able to execute a primary voice recognition process, a secondary voice recognition process and a dialog system process, wherein said hardware infrastructure is configured to:
  - a) capture audio data by means of said at least one microphone,
  - b) continually buffer said audio data in said audio buffer, such that said audio buffer always contains the audio data of a time period up to the present,
  - c) input said audio data to said secondary voice recognition process, which, on recognizing a keyword or a phrase triggers the start or the activation of said primary voice recognition process,
  - d) wherein said primary voice recognition process converts the entire or most recent content of said audio buffer as well as subsequent live audio data to text, said text being handed over to said dialog system process which likewise starts or is activated and which

analyzes the content of said text as to whether it contains a question, a message or a command made by the user to said software agent, in which case, if it is answered in the affirmative, said dialog system process triggers an appropriate action or generates an appropriate reply and contacts the user via said output device and

e) otherwise, if there is only background noise in said audio data or if said text does not contain any relevant or any evaluable content, said primary voice recognition process and said dialog system process return to the inactive state or terminate and again return the control to said secondary voice recognition process.

14. The system according to claim 13, wherein:

- a) said at least one microphone, said at least one output device and said at least one audio buffer are part of a local terminal, said secondary voice recognition process being executed on said local terminal,
- b) said primary voice recognition process and said dialog system process are executed on an external server or server network, the entire or most recent content of said audio buffer, as well as subsequent live audio data, being transferred via a network and/or radio network to said server or server network.

15. The system according to any one of claims 13 or 14, wherein said secondary voice recognition process consumes less power than said primary voice recognition process.

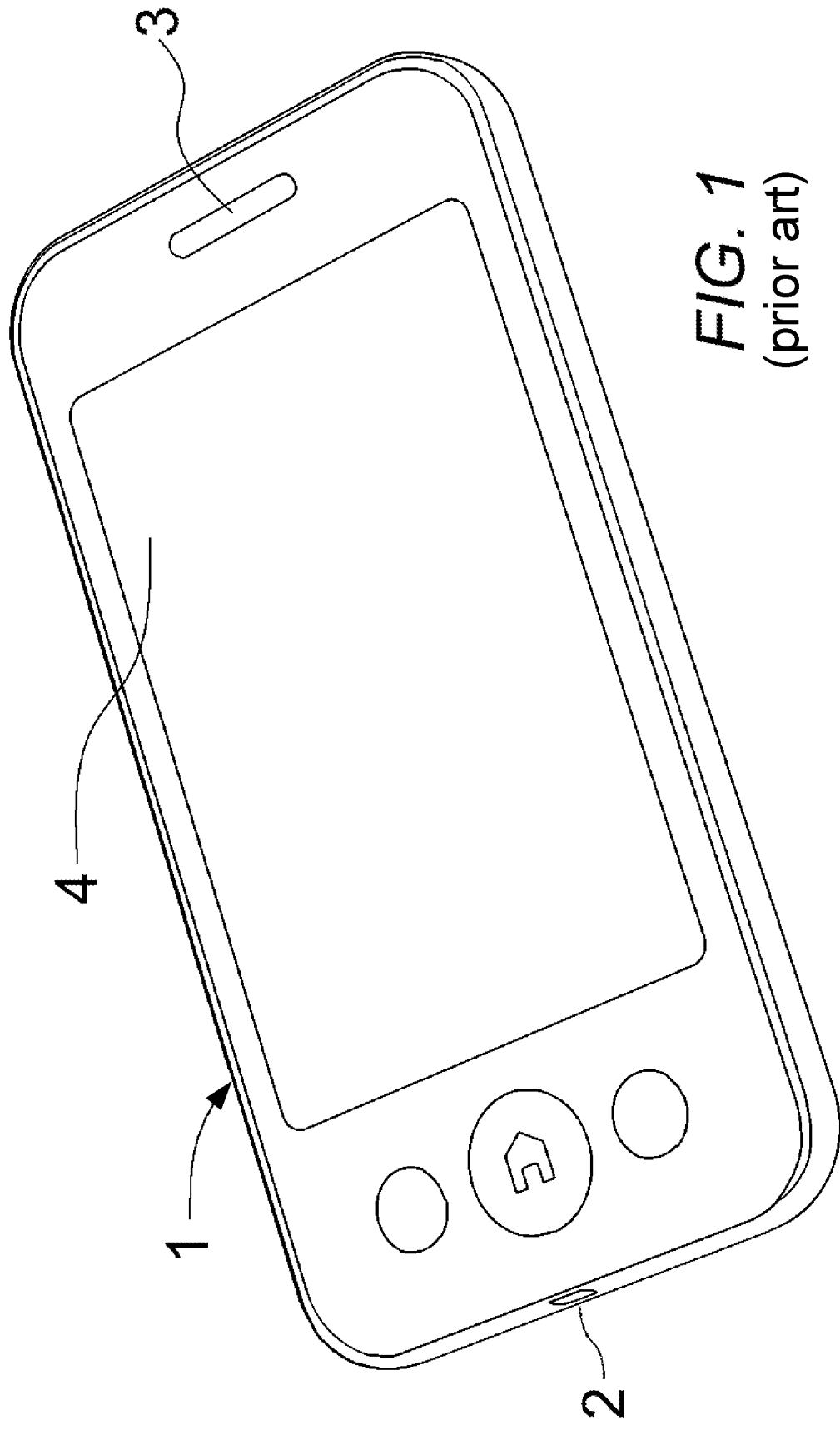
16. The system according to any one of claims 13 to 15, wherein: if noise or words are falsely interpreted as a keyword or phrase by said secondary voice recognition process, these errors are at least partially corrected by said primary voice recognition process.

17. The system according to any one of claims 13 to 16, wherein said software agent is a personal assistant system.

18. The system according to any one of claims 13 to 17, wherein said keyword or phrase is a product name or a nickname.

19. The system according to any one of claims 13 to 18, wherein said output device is a loudspeaker.

20. The system according to any one of claims 13 to 19, wherein said hardware infrastructure is configured to output an optical signal to the user by means of a second output device as soon as a keyword or phrase is recognized by said secondary voice recognition process.



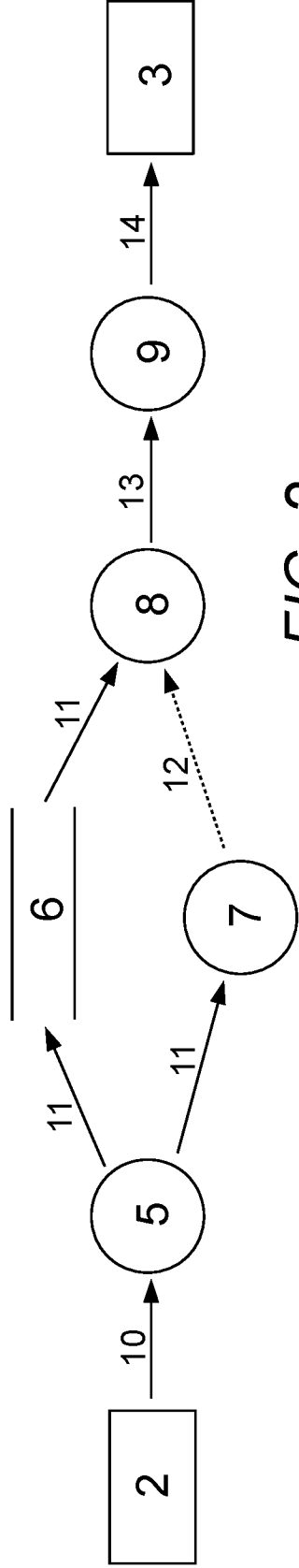


FIG. 2

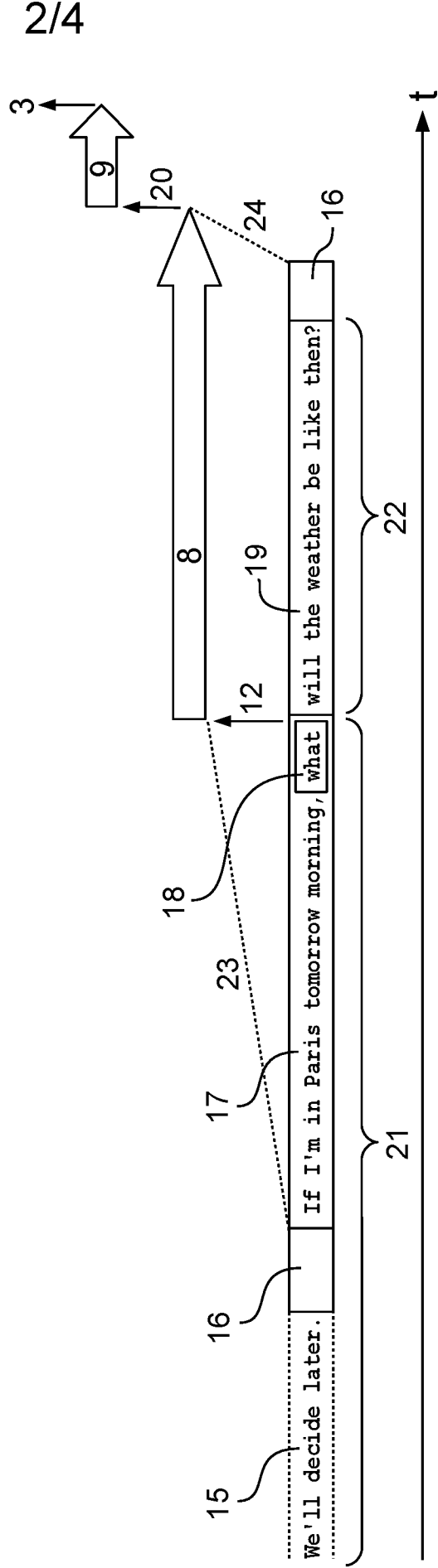


FIG. 3

3/4

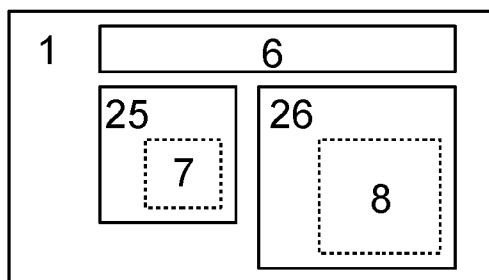


FIG. 4

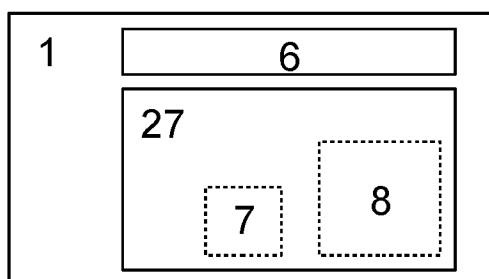


FIG. 5

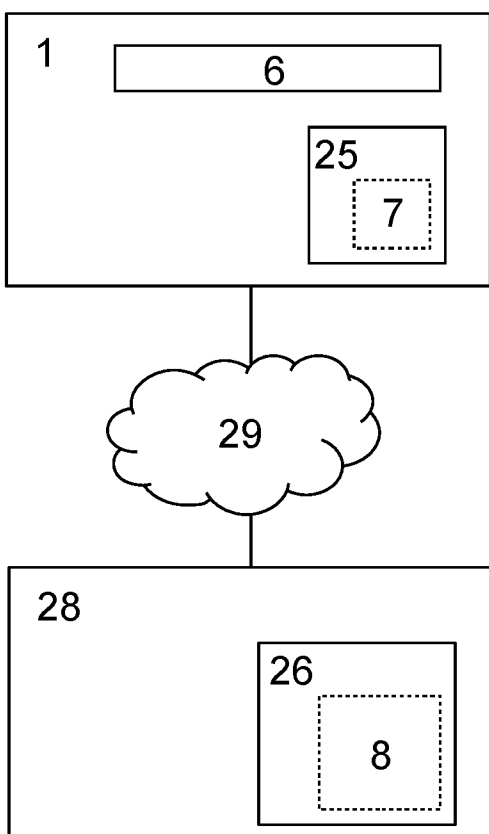


FIG. 6

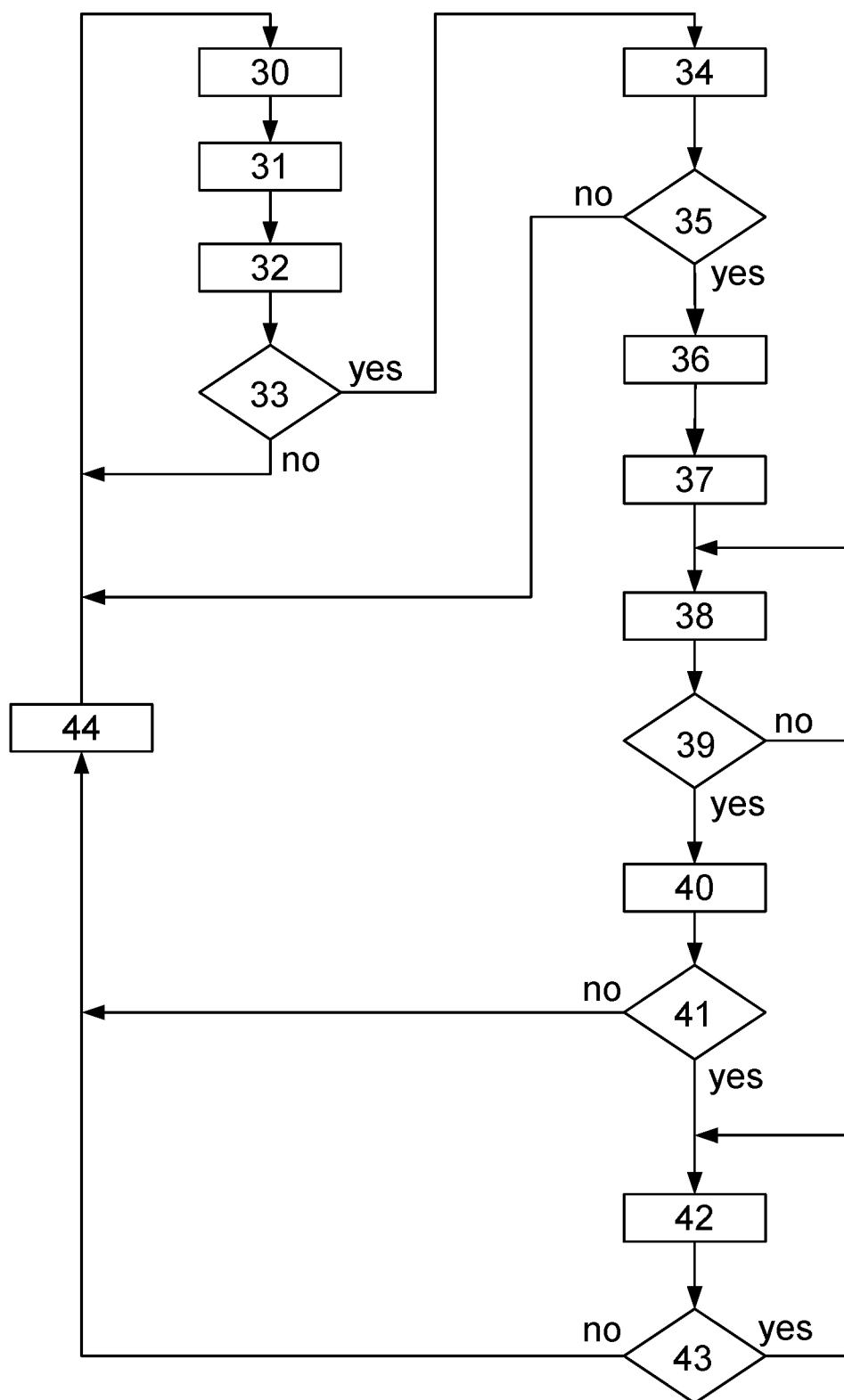


FIG. 7