



Utilización de mecanismos de Síntesis de Voz y Reconocimiento de Voz para el desarrollo de aplicaciones sobre dispositivos móviles-celulares.

Julio C. González Moreno.

Carrera de Ingeniería en Telemática, Dpto. de Computación, Facultad de Ciencias y Tecnología,
Universidad Nacional Autónoma de Nicaragua, León. (UNAN – León).
Email: jcgonzalez@ct.unanleon.edu.ni

Recibido: 11 Febrero, 2014

Aceptado: Julio 12, 2014

RESUMEN

El artículo presenta el análisis del diseño y de la arquitectura empleada en el desarrollo de un sistema, bajo la plataforma Android, que por medio de comandos de voz permite la comunicación de las personas mediante mensajes de texto (SMSs).

Se implementan acciones de lectura, recepción, borrado y escritura de SMSs. La idea es permitir que tanto las manos como la vista queden libres y sean útiles en el desarrollo de otras actividades cotidianas como por ejemplo: manejar, hacer ejercicio, cocinar, etc.

La funcionalidad está basada en la creación de una gramática libre de contexto que permitió generar un lenguaje finito para que los usuarios interactúen con su teléfono mediante comandos dictados a través de la voz. Para incorporar dicha funcionalidad se empleó una plataforma software que funcionará sobre dispositivos móviles-celulares y que a su vez permitirá acceder a servicios propios como son el envío y la recepción de SMSs, así como también la posibilidad de incorporar tecnologías proporcionadas por terceros y poder lograr la unificación de diversas funcionalidades y soluciones.

El sistema puede ser empleado por cualquier persona que posea conocimientos básicos de funcionamiento de tecnologías móviles-celulares que emplean Android y que su sistema vocal se encuentre en buenas condiciones para poder emitir, por medio de su voz, los comandos que se encargan de gobernar el sistema.

Palabras Claves: SMS, Android, Síntesis de voz, Reconocimiento de voz, Gramática libre de contexto

1. INTRODUCCION

Los siguientes escenarios: personas con algún tipo de discapacidad visual y el envío controlado de SMSs mientras se conduce un automóvil guardan una estrecha relación. En ambos casos resulta necesario utilizar la tecnología móvil-celular para la comunicación a través de SMSs. El problema aparece porque en el primer escenario el sentido de la vista es nulo o casi nulo y en el segundo escenario el sentido de la vista se encuentra ocupado ya que la concentración visual es crítica. En ambos casos se dificulta, o imposibilita, el envío de SMSs para la comunicación.

El problema planteado en los dos escenarios anteriores podría ser solucionado si se incorporasen características de interacción vocal, referida al intercambio de SMSs, dentro de los dispositivos móviles-celulares. Para dar solución al problema se ha desarrollado un sistema, para dispositivos móviles-celulares con Android, que incorpora dichas características de interacción vocal para el envío y la recepción de SMSs^[1].

La funcionalidad del sistema está basada en tres grandes elementos: el primero se corresponde con la definición de una gramática libre de contexto, la cual permitió representar un lenguaje finito utilizado para consolidar los comandos que son admitidos por el sistema; el segundo elemento se corresponde con la interacción usuario-teléfono mediante comandos de voz a través de un mecanismo de reconocimiento de voz ofrecido por Google y el tercer elemento se corresponde con la interacción teléfono-usuario mediante la síntesis de voz ofrecida por el motor Svox Pico TTS^[2].

2. DESARROLLO

El sistema desarrollado integra una serie de elementos principales que lo conforman, los cuales se detallan en la siguiente figura:

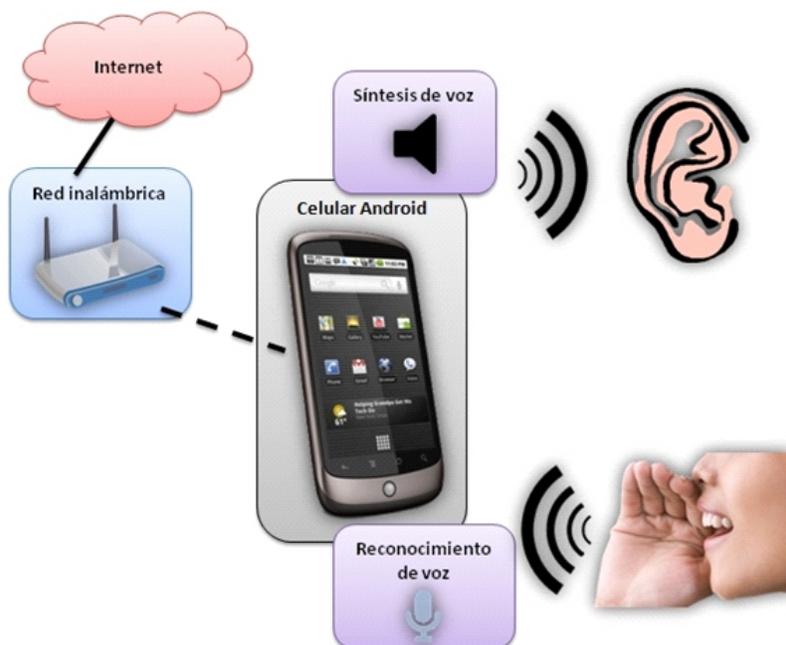


Figura 1. Elementos principales que conforman la solución



Celular Android: Se requiere un dispositivo móvil-celular con Android versión 2.0 o superior, poseer micrófono y altavoces, capacidades de envío y recepción de SMSs y conexión a Internet vía Wifi.

Red inalámbrica: Para ofrecer la funcionalidad del reconocimiento de voz se requiere de acceso a Internet; dicho acceso se puede llevar a cabo a través de una conexión vía WiFi.

Internet: El acceso a Internet ofrecerá el medio de comunicación para interactuar con el servidor de Google; el cual proporciona la funcionalidad necesaria para el procesamiento realizado por el mecanismo reconocimiento de voz^[3].

Síntesis de voz: Corresponde con el módulo del sistema encargado de emitir todos los tipos de mensajes (informativos, de error, etc.) que son generados por el sistema; dichos mensajes son emitidos en formato sonoro (o de audio) a través de una voz artificial generada a través del motor Svox Pico TTS.

Reconocimiento de voz: Hace uso de una implementación proporcionada por Google. Para su correcto funcionamiento requiere de una conexión a Internet.

3. COMPONENTES LÓGICOS

A continuación se detallan cómo y para que fueron utilizados los componentes lógicos que en esencia constituyen gran parte del sistema.

3.1 RECONOCIMIENTO DE VOZ

El proceso de reconocimiento de voz mediante ordenador dota a los sistemas digitales de la capacidad de recibir mensajes orales tomando como entrada la señal acústica de la voz recogida a través de un micrófono y tiene como objetivo final decodificar el mensaje contenido en la onda acústica de la voz para realizar las acciones pertinentes. Para lograr este objetivo es necesario conjugar una gran cantidad de conocimientos acerca del sistema auditivo humano, sobre la estructura del lenguaje humano, la representación del significado de los mensajes y sobre todo el auto-aprendizaje de la experiencia diaria^[4].

Un sistema de reconocimiento de voz debe cumplir con tres tareas esenciales:

- **Pre-procesamiento:** se encarga de convertir la entrada de la voz a un formato que el sistema reconocedor de voz es capaz de procesar.
- **Reconocimiento:** trata de identificar lo que se dijo. En esta tarea se lleva a cabo una traducción de la señal de voz a un texto o cadenas de caracteres específicos.
- **Comunicación:** formado el texto correctamente, este es enviado al sistema (Software/Hardware) que lo requiera.

Para hacer uso del reconocimiento de voz fue necesario emplear el paquete Java *android.speech.RecognizerIntent*. Esto permite ejecutar el mecanismo de reconocimiento de voz para capturar la onda de voz, emitida por el usuario, mediante un flujo de audio y enviarla a un servidor de Google para que la procese y obtener a partir de dicho flujo una representación en cadena de caracteres (texto plano); luego esta representación es devuelta por el servidor de Google hacia el dispositivo móvil-celular en forma de texto el cual ya puede ser utilizado en la lógica de programación para, por ejemplo, proporcionar la interacción con los comandos.



Figura 2. Reconocimiento de voz

3.2 SÍNTESIS DE VOZ

La síntesis de voz es un proceso de sintetización del habla que genera una voz sintética o artificial (no pregrabada), el cual consiste en la producción artificial de habla humana. Un sistema usado con este propósito recibe el nombre de sintetizador de voz y puede implementarse a través de soluciones basadas en software o en hardware. La síntesis de voz se llama a menudo Text To Speech (TTS), en referencia a su capacidad de convertir texto en voz artificial.

La calidad que posea una voz sintética vendrá dada por los siguientes factores:

- **Inteligibilidad:** determina con qué facilidad (o dificultad) es entendida la voz artificial.
- **Naturalidad:** determina en qué medida se asemeja la voz artificial a la voz real de un ser humano.

Un sistema sintetizador de voz, que convierte el texto en voz artificial, se compone de dos partes fundamentales: un *front-end* que toma como elemento de entrada el texto y produce una representación lingüística fonética del mismo; y un *back-end* que retoma como entrada la representación lingüística fonética y produce una forma de onda sintetizada.

El motor TTS empleado por el sistema desarrollado fue *Svox Pico TTS*, el cual es una solución telefónica dedicada para complementar el creciente éxito de la tecnología de texto a voz en el mercado de dispositivos móviles-celulares y permitir que aumente la cantidad de usuarios que disfrutan de una tecnología de manos libres basada en mensajes sonoros emitidos mediante voz sintética. Para hacer uso de dicho motor TTS, fue necesario emplear el paquete Java *android.speech.tts.TextToSpeech*. Lo que se consigue cuando se incorpora dicho paquete es que una vez ejecutada la síntesis de voz, es posible emplear capacidades de habla mediante la generación de voz artificial; para así, comunicarle al usuario las acciones que suceden mientras se hace uso del sistema^[5].



Figura 3. Síntesis de voz

3.3 PROCESAMIENTO DE SMSs

El procesamiento de SMSs corresponde con todas aquellas acciones que implican una interacción directa con SMSs, ya sea para leerlos, borrarlos o actualizarlos. Estas acciones están asociadas y sirven de apoyo para dar respuesta a algunos comandos propios del sistema como son: nuevo, borrar y leer.

Las acciones que fueron implementadas y que implican una interacción directa con SMSs fueron las siguientes:

- **Capturar un nuevo mensaje entrante:** fue necesario determinar cuándo se genera el evento de recepción de un nuevo mensaje para detectar y reaccionar ante dicho evento y de esta manera agregar el nuevo mensaje recibido a la bandeja de entrada y marcarlo como no leído.
- **Enviar un nuevo mensaje hacia un destinatario:** se utiliza para ofrecer la funcionalidad del comando “nuevo” el cual solicita al usuario toda la información necesaria (cuerpo del mensaje y número del destinatario) para la emisión de un nuevo SMSs hacia un determinado destinatario.
- **Borrar mensajes:** accede a las bandejas de entrada o de salida, según corresponda, para realizar la eliminación de SMSs. Los tipos de eliminaciones contempladas fueron: eliminar mensajes entrantes leídos (comando “borrar recibidos”), eliminar el último mensaje entrante leído (comando “borrar”), eliminar todos los mensajes entrantes leídos y no leídos (comando “borrar todo”) y eliminar todos los mensajes salientes (comando “borrar enviados”).
- **Acceso a la bandeja de entrada:** se utiliza para ofrecer la funcionalidad del comando “leer” el cual lee (a través de la síntesis de voz) algún mensaje almacenado en la bandeja de entrada. El usuario debe indicarle al sistema cuál es el número del mensaje que desea leer.
- **Actualizar el estado de un mensaje:** utilizado cuando es necesario acceder a la bandeja de entrada para marcar un mensaje como leído o para marcar un nuevo mensaje recibido como no leído.

3.4 NOTIFICACIONES

Una notificación corresponde con un mensaje informativo que le indica al usuario la ocurrencia de un evento. Dicha notificación es mostrada en la barra de estado del teléfono y permanece ahí hasta que el usuario decide eliminarla.

Las notificaciones en combinación con la síntesis de voz se emplean para informar al usuario la recepción de un nuevo mensaje. Si se recibe un nuevo mensaje entrante mientras el sistema se encuentra en ejecución; el sistema informará de este suceso al usuario a través de un mensaje generado mediante una voz artificial; también se mostrará la ocurrencia de dicho suceso en la barra de estado del teléfono a través de una notificación.

Los tipos de estados de las notificaciones mostradas en la barra de estado del teléfono son:

- **Estado inicial:** se muestra en toda la barra de estado un mensaje informativo que dice: “Mensaje recibido de: <Nombre_emisor>”.
- **Estado informativo:** se muestra en la esquina superior izquierda de la barra de estado, el logo del sistema indicando que se ha recibido un nuevo mensaje.
- **Estado desplegado:** una vez desplegada la barra de estado se muestra el nombre y el cuerpo del mensaje enviado por el emisor.



Figura 4. Notificaciones y estados de las notificaciones

3.5 GRAMÁTICA LIBRE DE CONTEXTO

La lógica de los comandos que rigen el sistema está basada en la definición de un vocabulario de términos finito que indica la sintaxis que se debe seguir para cada uno de los comandos admitidos. El vocabulario de términos finito fue consolidado con el apoyo de una gramática libre de contexto.

Las gramáticas libres de contexto se utilizan para describir la mayoría de los lenguajes de programación, de hecho, la sintaxis de la mayoría de lenguajes de programación está definida mediante gramáticas libres de contexto. Por otro lado, estas gramáticas son suficientemente simples como para permitir el diseño de algoritmos de análisis sintácticos eficientes que, para una cadena de caracteres dada, determinen como puede ser generada desde la gramática.

La sintaxis de cada uno de los comandos admitidos por el sistema está dado por la gramática libre de contexto.



Los comandos admitidos se detallan en la siguiente tabla:

Comando	Descripción
nuevo	Crea nuevo mensaje y lo envía a destinatario
borrar recibidos	Borra mensajes recibidos leídos
borrar enviados	Borra todos los mensajes enviados
borrar todo	Borra todos los mensajes enviados y recibidos (leídos y no leídos)
borrar	Borra último mensaje recibido leído
leer	Leer algún mensaje recibido

TABLA N° 1: Comandos admitidos por el sistema

La gramática libre de contexto propuesta, se diseñó para dar respuesta a la forma de operar de cada uno de los comandos admitidos por el sistema.

La propuesta de gramática libre de contexto es la siguiente:

```

<inicio> := (<comandos> <mensaje> <num> | <contactos> <confirmación>) | (<comandos> <confirmación>) | (<comandos>
<num>)
<comandos> := nuevo | borrar recibidos | borrar enviados | borrar todo | borrar | leer
<mensaje> := <mensaje> (<palMay> | <palMin> <espacio>)
<confirmación> := <afirmativo> | <negativo>
<palMay> := A|B|C|D|E|F|G|H|I|J|K|L|M|N|O|P|Q|R|S|T|U|V|W|X|Y|Z
<palMin> := a|b|c|d|e|f|g|h|i|j|k|l|m|n|o|p|q|r|s|t|u|v|w|x|y|z
<espacio> := " "
<num> ::= 0|1|2|3|4|5|6|7|8|9
<contactos> := buscar
<afirmativo> := si
<negativo> := no

```



Una gramática libre de contexto está compuesta por cuatro elementos:

- 1. Símbolos terminales:** es el conjunto finito de los símbolos que forman las palabras del lenguaje, es decir, son los elementos que no generan nada. Por ejemplo: <afirmativo>, <negativo>
- 2. Símbolos no terminales:** es el conjunto finito de símbolos que permiten representar estados intermedios de la generación de las palabras del lenguaje. Son los elementos del lado izquierdo de una producción, antes de la flecha "→". Por ejemplo: <inicio>, <comandos>
- 3. Producciones:** permiten generar las palabras del lenguaje y por tanto son sentencias que se escriben en la gramática. Por ejemplo: <mensaje> → Hola amigos
- 4. Símbolo inicial:** es el símbolo a partir del que se aplican las reglas de la gramática para obtener las distintas palabras del lenguaje, es decir, es el primer elemento de la gramática. Por ejemplo: <inicio>

La sintaxis de cada uno de los comandos admitidos por el sistema según la gramática libre de contexto propuesta se detalla en la siguiente tabla:

nuevo							
nuevo	< cuerpo del mensaje >	< número del destinatario / buscar >		< si / no >			
borrar recibidos							
borrar recibidos		< si / no >					
borrar enviados							
borrar enviados		< si / no >					
borrar todo							
borrar todo		< si / no >					
borrar							
borrar		< si / no >					
leer							
leer		< número del mensaje a leer >					
recibir nuevo mensaje							
¿Desea leer el mensaje?	< si / no >	lectura del mensaje	¿Desea emitir una respuesta sobre el mensaje recibido?	< si / no >	< cuerpo del mensaje >	¿Seguro de enviar el mensaje?	< si / no >

TABLA N° 2: Sintaxis de cada comando admitido por el sistema

Una vez que el usuario ha emitido el comando; el sistema automáticamente lo irá guiando, a través de una voz artificial, hasta completar la sintaxis del comando emitido.

Al utilizar el comando “nuevo” es posible que el usuario indique el número del destinatario de dos maneras diferentes: la primera es que directamente se dicte el número del destinatario al que se desea enviar el mensaje y la segunda es indicando la palabra “buscar” la cual mostrará una lista con todos los números telefónicos de los contactos almacenados. Esto le permite al usuario seleccionar el número del destinatario al cual será enviado el nuevo mensaje sin que tenga que saberlo de memoria.

Al hacer uso del comando “leer”, el usuario debe indicar un número que represente el mensaje que quiere leer. Es importante destacar que números pequeños corresponderán con mensajes recibidos recientemente y números grandes corresponderán con mensajes antiguos. Si se recibe un nuevo mensaje entrante mientras el sistema se encuentra en ejecución; se informará de este suceso al usuario a través de un mensaje generado mediante una voz artificial y también se mostrará dicho suceso en la barra de estado del teléfono.

Un árbol de análisis sintáctico permite mostrar gráficamente cómo se puede derivar cualquier cadena de un lenguaje a partir del símbolo distinguido de una gramática que genera ese lenguaje. A continuación se muestran algunos ejemplos de árboles de análisis sintáctico según la gramática propuesta.

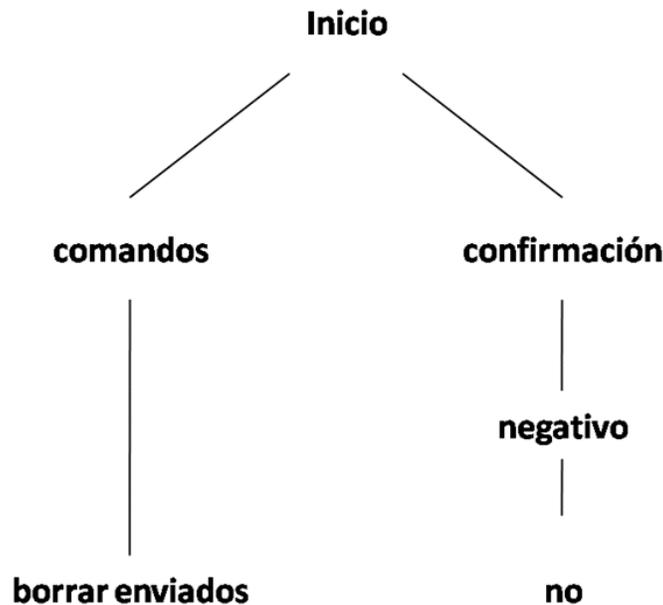


Figura 5. Árbol de análisis sintáctico asociado a la eliminación de todos los mensajes enviados

3.6 OPCIONES CONFIGURABLES

El sistema proporciona una serie de opciones que son configurables y que modifican el comportamiento del mismo. El acceso a la configuración de dichas opciones se puede llevar a cabo al pulsar el botón MENU del teléfono. Dichas opciones son:

- **Volumen multimedia:** permite cambiar la intensidad del volumen multimedia entre los valores: bajo, medio o alto (valor por defecto).
- **Velocidad de voz:** permite seleccionar la velocidad a la que se quiere que la voz artificial emita los mensajes sonoros, es posible establecer los valores: rápida, normal (valor por defecto) o lenta.
- **Idioma:** permite seleccionar el idioma que se quiere que sea empleado por el sistema para su funcionamiento, es posible establecer los valores: español (valor por defecto) e inglés.
- **Acerca de:** muestra una breve descripción y el nombre del autor del sistema.



Figura 6. Opciones configurables del sistema



4. ASPECTOS RELEVANTES

Un aspecto importante a citar es que el usuario puede indicar los comandos con un máximo de cinco intentos dictados mediante su voz, si después de cinco intentos el sistema no ha sido capaz de obtener la información solicitada de forma correcta (según la sintaxis de cada comando) se abortará la acción que se esté desarrollando en ese momento.

Al ejecutar el sistema; éste verifica si el reconocimiento de voz está presente en el teléfono, en caso de no estarlo se notifica dicho suceso por medio del mensaje sonoro: “Detección de voz no presente en el teléfono” y se inhabilita la posibilidad de emplear algún comando de voz; si el mecanismo de reconocimiento de voz se encuentra presente en el teléfono se procederá a verificar si éste se encuentra conectado a alguna red inalámbrica, en caso de no estarlo se notifica dicho suceso a través del mensaje sonoro: “Conexión de red no disponible” y se inhabilita la posibilidad de que el usuario pueda emplear algún comando; en caso de que el teléfono si se encuentre conectado a alguna red inalámbrica se habilitará la posibilidad de que el usuario pueda utilizar el sistema mediante algún comando de voz. Es importante aclarar que cuando se verifica si el teléfono está conectado a alguna red inalámbrica no se verifica si la red inalámbrica a la que se encuentra conectada cuenta con acceso a Internet.

Es posible guardar las preferencias que han sido configuradas por el usuario. Dentro de estas preferencias se encuentran todas aquellas opciones que se muestran al pulsar el botón MENU del teléfono (exceptuando la opción Acerca de...). Una vez que el usuario modifica alguna de las opciones antes citadas su nuevo valor será almacenado para que perdure aún cuando el sistema sea cerrado o el teléfono sea apagado. Una vez que el usuario inicia el sistema, se restablecen los valores almacenados para cada una de las opciones.

Resulta necesario dejar muy claro que el sistema puede presentar un funcionamiento no deseado o incorrecto bajo ciertas circunstancias como pueden ser:

- Palabras complejas dictadas con mucha rapidez.
- El motor TTS empleado no es capaz de realizar entonaciones interrogativas o exclamativas, por lo que suenan igual que las enunciativas.
- El reconocimiento de voz requiere de acceso a Internet para cada frase, palabra u oración dictada, por lo que puede sufrir una leve latencia que dependerá de la velocidad de conexión a Internet.
- El volumen, el eco y el ruido de fondo juegan un papel importante en el reconocimiento de las frases, palabras u oraciones. En algunas situaciones este aspecto puede ser un factor no influyente, ya que el propio micrófono del dispositivo aplica mecanismos que tratan de paliar dichos problemas.



CONCLUSIONES

Una vez que el sistema fue diseñado, codificado y puesto a prueba se obtuvieron las siguientes conclusiones

- La arquitectura y programación empleadas permitieron un funcionamiento autónomo mediante comandos de voz en respuesta a las necesidades originales (envío, recepción, lectura y eliminación de SMSs).
- El diseño sencillo e intuitivo ha permitido que usuarios con capacidades especiales (de la tercera edad, lesionados, con discapacidades visuales, etc.) hagan un uso eficiente y frecuente del sistema.
- Haber incorporado los mecanismos de texto a voz y de reconocimiento de voz de Google permitió el desarrollo de un sistema que emplea dichos mecanismos para la interacción de SMSs a través de la voz.
- La combinación de la tecnología móvil-celular con plataformas de trabajo flexibles y programables permiten potenciar la creación de sistemas que aprovechan los recursos gestionados por un dispositivo móvil-celular.

Lo antes mencionado reafirma que fue posible el diseñado y desarrollado de un sistema sencillo y práctico que resuelve la problemática planteada en el trabajo.

AGRADECIMIENTOS

A Denis Espinoza por brindarme su apoyo y conocimiento cuando más lo necesite. También a las personas, organismos y proyectos que han hecho posible la adquisición de nuevos conocimientos y del desarrollo de este trabajo. A Santiago Molina, por cultivar su conocimiento, educación y cultura sobre mi persona.

REFERENCIAS

1. Chami, F. Java Code Geeks. [Sitio en Internet]. Consultado el 23 de Septiembre de 2011. Disponible en: <http://www.javacodegeeks.com/2010/09/android-text-to-speech-application.html>
 2. Elsey, J. James Elsey - Digital learnings and other such nonsense. [Sitio en Internet]. Consultado el 12 de Febrero de 2011. Disponible en: <http://www.jameselsey.co.uk/blogs/techblog/android-how-to-implement-voice-recognition-a-nice-easy-tutorial/>
 3. oriolpons. desctrl Webs & Mobile Applications. [Sitio en Internet]. Consultado el 25 de Octubre de 2011. Disponible en: <http://www.desctrl.com/blog/2010/10/25/android-check-network-available/>
 4. Murphy, M. L. (2009). *Beginning Android*. Apress, 386 pags.
- Project, A. O. (s.f.). Android Developers. [Sitio en Internet]. Consultado el 20 de Octubre de 2011. Disponible en: <http://developer.android.com/resources/articles/tts.html>