

UNIVERSIDAD POLITÉCNICA DE MADRID
DEPARTAMENTO DE INGENIERÍA ELECTRÓNICA
ESCUELA TÉCNICA SUPERIOR DE INGENIEROS
DE TELECOMUNICACIÓN



TESIS DOCTORAL

ESTRATEGIAS PARA LA MEJORA DE LA
NATURALIDAD Y LA INCORPORACIÓN DE
VARIEDAD EMOCIONAL A LA
CONVERSIÓN TEXTO A VOZ EN
CASTELLANO

JUAN MANUEL MONTERO MARTÍNEZ
Ingeniero de Telecomunicación

Director de la Tesis
JOSÉ MANUEL PARDO MUÑOZ
Doctor Ingeniero de Telecomunicación

Madrid, 2003

Datos de la Tesis Doctoral

Autora:

JUAN MANUEL MONTERO MARTÍNEZ

Título:

ESTRATEGIAS PARA LA MEJORA DE LA NATURALIDAD Y LA
INCORPORACIÓN DE VARIEDAD EMOCIONAL A LA CONVERSIÓN
TEXTO A VOZ EN CASTELLANO

Director:

Dr. JOSÉ MANUEL PARDO MUÑOZ

Departamento:

Departamento de Ingeniería Electrónica de la Escuela Técnica Superior de
Ingenieros de Telecomunicación de la Universidad Politécnica de Madrid.

Fecha de lectura:

14 de Noviembre de 2003

Calificación:

Sobresaliente Cum Laude por unanimidad del tribunal

Número de colegiado/asociado: A09073

Resumen de la tesis doctoral

1. Contexto de la tesis doctoral

La Tesis está enmarcada en el campo de las tecnologías del habla, concretamente en los de la Conversión Texto a Voz (CTV) y el Procesamiento de Lenguaje Natural (PLN), cuyo objetivo fundamental es la conversión de un texto en habla por parte de máquinas, alcanzando una inteligibilidad y una naturalidad que la hagan indistinguible del habla humana.

La importancia de la comunicación oral es evidente en la vida de los seres humanos. En todas las sociedades, incluso en las más primitivas, la comunicación oral existe y está basada en mecanismos acústicos, sintácticos y semánticos complejos, con independencia del nivel tecnológico alcanzado por la sociedad en cuestión.

En las últimas décadas se han realizado grandes esfuerzos en el área de la síntesis del habla. De hecho, una gran cantidad de laboratorios en el ámbito nacional e internacional han concentrado sus esfuerzos en la consecución de sistemas cada vez más complejos y eficientes, que no sólo hacen uso de las características acústicas de la voz, sino también de las particularidades sintácticas e incluso semánticas de cada lengua.

La tecnología actual es capaz de convertir texto en voz con una alta tasa de inteligibilidad, aunque su grado de naturalidad no sea tan alto como desearíamos: no podemos imitar el amplio espectro de cadencias, melodías y cualidades que cubre la voz humana. Por lo general las voces sintéticas podían ser catalogadas como monótonas o incluso aburridas: nuestros ordenadores carecían hasta ahora de capacidad para transmitirnos emociones, para adaptar la voz a diferentes estilos de locución (formales o informales), para engañarnos con su naturalidad.

A medida que las tecnologías del habla se han ido implicando cada vez más como parte integral de aplicaciones prácticas en escenarios reales (como servicios de obtención de información por línea telefónica [5] [4] [3] [18] [19] [35] [36], control de dispositivos e instrumental en coches, asistentes personales (PDAs), robots o agentes virtuales animados dotados de personalidad propia [11] [16], etc.), se ha hecho patente la necesidad de desarrollar sistemas automáticos de conversión texto-habla naturales y dotados de la variedad que nos caracteriza a los seres humanos.

Vivimos una época en la que se está dando un gran auge en los estudios teórico-prácticos de la llamada inteligencia social o inteligencia emotiva, la que se encarga de controlar con inteligencia las propias emociones, reconocer las emociones de los demás y reaccionar empáticamente a las mismas. No es descabellado pensar en que esa misma inteligencia social debería gobernar las futuras aplicaciones de intercomunicación hombre-máquina, haciéndolas más y más amigables (*user-friendly*), tanto si son presenciales como si son telefónicas o telemáticas [23] [31] [33]. Para ello, deberíamos de dotar a los sintetizadores de una voz más diversa y humana: un usuario habitual del sistema o un usuario con problemas que dialoga con el sistema, y tras repetidos intentos no consigue acceder a la información que precisa, deben ser tratado de un modo especial, como lo haría un experto humano.

Tres son las grandes líneas que se deben abordar para avanzar hacia la consecución del ambicioso objetivo final, dos de las cuales se abordan en esta tesis:

- Mejorar el procesamiento automático lingüístico-prosódico que, a partir de texto, debe recoger la información sintáctica y semántica del mismo y aprovecharla para

generar una voz más natural desde el punto de vista de su ritmo y su entonación. Colateralmente, las técnicas empleadas son de gran utilidad en comprensión automática de habla en sistemas de diálogo [30] [28] [6].

- Analizar y generar variantes de una voz artificial, en especial incorporar emociones y actitudes que la doten de personalidad y empatía [1] [13] [29]. Esto resulta de utilidad en estudios forenses [32] [21].
- Mejorar la capacidad de imitación del timbre de cada uno de los locutores apropiados para las aplicaciones más habituales.

En esta tesis se ha tratado una parte importante de los dos primeros puntos.

2. Objetivos de la tesis doctoral

El objetivo de la tesis doctoral ha sido profundizar en el estudio de diversas estrategias para incorporar naturalidad y variedad emocional en la conversión texto a voz en castellano, haciendo especial énfasis en el procesado lingüístico orientado al modelado de la prosodia, en el modelado de la frecuencia fundamental dentro de un dominio restringido y en el análisis, modelado y síntesis de voz con emociones.

Los sistemas de síntesis de habla sobre los que se ha realizado este estudio comprenden las principales tecnologías del estado del arte, comenzando con el sistema multilingüe de síntesis por formantes de la empresa sueca Telia-Infovox [15] hasta el sistema registrado Boris de síntesis por concatenación, desarrollado por el autor de la tesis y otros miembros del Grupo de Tecnología del Habla del Departamento de Ingeniería Electrónica de la Universidad Politécnica de Madrid. En la actualidad, ambos sistemas funcionan en aplicaciones en tiempo real, motivo por el cual, a lo largo de toda la tesis, se ha optado por trabajar con técnicas que no requieran de un incremento prohibitivo de la carga computacional y de memoria del sistema, de modo que pudieran ser incorporadas al sistema final para permitir su funcionamiento en aplicaciones reales.

3. Procesamiento de Lenguaje Natural

En el capítulo dedicado a las investigaciones en procesado lingüístico del texto se comienza describiendo en detalle los corpora empleados en la experimentación (que incluyen 54 millones de palabras procedentes de 2 años completos de artículos del periódico El Mundo), tanto en normalización como en etiquetado. La técnica desarrollada en normalización emplea reglas de experto, con muy buenos resultados tanto en precisión como en cobertura (>85%), destacando el empleo de reglas de silabificación para la detección precisa de palabras extranjeras (>99%). La cobertura del sistema desarrollado alcanza un nivel casi insuperable (>99,8%), lo cual confirma la calidad del trabajo llevado a cabo.

Al afrontar la desambiguación gramatical, se comparan tres técnicas: reglas de experto (tasa inferior al 98%), aprendizaje automático de reglas (tasa por debajo del 99%) y modelado estocástico (por encima del 99%), obteniéndose los mejores resultados con esta última técnica, debido a su capacidad de procesar más adecuadamente textos fuera del dominio de entrenamiento.

Finalmente se aborda el análisis sintáctico por medio de gramática de contexto libre como un proceso en dos fases: una primera sintagmática (*shallow parsing*) y una segunda relacional básica, a fin de maximizar la cobertura del análisis. Para la resolución de las ambigüedades que nos permiten alcanzar gran cobertura se adapta el principio de mínima longitud de descripción con notables resultados. Con las gramáticas desarrolladas se alcanza una tasa de cobertura y precisión superior al 96% en

el caso del análisis sintagmático y superiores al 87% en el sintáctico, los mejores resultados publicados en castellano para unas tareas tan complejas.

4. Modelado prosódico en dominio restringido

Para el modelado de F0 en un dominio restringido se emplean perceptrones multicapa. En una primera etapa se describe y evalúa una nueva técnica de diseño de base de datos basada en un algoritmo voraz moderado mediante subobjetivos intermedios. Esta novedosa combinación de voracidad y moderación permite alcanzar precisiones superiores al 95% para los numerosos rasgos que relacionados con la predicción de la prosodia [10] [9].

La exhaustiva experimentación con los diversos parámetros de predicción, la configuración de la red y las subdivisiones de la base de datos ocupa la mayor parte del capítulo, destacando la aportación de un parámetro específico del dominio restringido (el número de la frase portadora del texto que sintetizar) junto a otros más clásicos (la acentuación, el tipo de grupo fónico y la posición en el mismo), además del análisis sobre cómo agrupar las grabaciones para obtener el mejor modelado posible [8] [2] [22] [27].

5. Análisis y síntesis de voz con emociones

El capítulo dedicado a la voz emotiva comienza detallando el proceso de creación de una nueva voz castellana masculina en síntesis por formantes con modelo mejorado de fuente (reglas y metodología), evaluando las posibilidades de personalización de voz que ofrece. La voz desarrollada, por su calidad, fue adoptada por el sintetizador multilingüe de Infovox en castellano [14].

Para trabajar con voz con emociones se diseña, graba y etiqueta una base de datos de voz en la que un actor simula tristeza, alegría, sorpresa, enfado y también una voz neutra. Por medio de técnicas paramétricas (modelo de picos y valles en tono, y multiplicativo en las duraciones) se analiza prosódicamente la base de datos y se establece una primera caracterización de la voz en las distintas emociones. Empleando como base la voz personalizable se desarrolla el sistema completo de conversión texto a voz con emociones y se evalúa, destacando la rápida adaptación de los usuarios en cuanto a la identificación de la emoción expresada [12]. Finalmente se experimenta con síntesis por concatenación y síntesis por copia, llegando a las siguientes conclusiones: la voz sorprendida se identifica prosódicamente, las características segmentales son las que caracterizan al enfado en frío; y, finalmente, la tristeza y la alegría son de naturaleza mixta, hallazgo pionero que fue posteriormente confirmado por posteriores investigadores en diversos países y diversas lenguas [7].

6. Conclusiones

Procesado lingüístico automático

Se han probado tres técnicas de desambiguación contextual gramatical:

1. una basada en reglas manuales (cuyo coste de desarrollo no se ve compensado con una tasa adecuadamente elevada);

2. otra basada en reglas inferidas automáticamente (que supere la dificultad de generar manualmente las reglas). En esta técnica de aprendizaje de reglas los resultados han sido excelentes, aunque bastante dependientes del dominio de entrenamiento (pudiéndose reducir la tasa desde un 99% a menos de un 96%), debido al sobreentrenamiento y al aprendizaje de reglas no generales;
3. otra basada en modelado estocástico, constatándose el superior comportamiento de esta última técnica al realizar ensayos fuera de dominio. Al emplear la técnica estocástica con diccionarios no adaptados a un dominio concreto sin probabilidades, se ha alcanzado una cobertura del 99,89% en textos de un dominio distinto al dominio de entrenamiento, comparable a los mejores sistemas en castellano y que (a pesar de no tener probabilidades) supera significativamente en precisión los resultados de un sistema léxico basado en probabilidades, aunque a costa de una mayor ambigüedad media (>96% en el primer candidato). En la desambiguación, resulta también significativa la mejora debida al tratamiento de las locuciones, dado que los modelos probabilísticos basados en categorías no son capaces de modelar bien contextos amplios y con pocos ejemplos.

Se ha adaptado un sistema de análisis por medio de gramáticas de contexto libre, desarrollando y evaluando con éxito una gramática robusta de dominio general en dos niveles, uno sintagmático y otro relacional. Cabe destacar el empleo de reglas de corte para reducir el número de análisis posibles (con sólo un 0,35% de imprecisión), la aplicación de reglas de concordancia como filtrado posterior al análisis y el uso de un criterio muy simple de número mínimo de segmentos para elegir el mejor análisis (sin necesidad de información probabilística). En el primer nivel sintagmático los resultados han sido excelentes (cobertura y precisión superiores al 96%, comparables a los mejores resultados en castellano, aunque sobre distinto corpus), a pesar de que haya un 1% de errores debidos al etiquetado previo. En el segundo nivel relacional los resultados son prometedores (cobertura y precisión superiores al 87%).

En el nivel léxico, se ha experimentado con diccionarios adaptados al dominio, con diccionarios generales y con diccionarios extranjeros, destacando la aportación de los dos primeros tipos. También se ha constatado la necesidad de incluir reglas robustas para etiquetar palabras fuera de vocabulario, experimentándose con mejoras significativas el empleo de reglas de experto basadas en las terminaciones de las palabras. En este sentido se ha ensayado el empleo de varios conjuntos de reglas manuales de experto procedentes de otros proyectos (completadas con algunas nuevas reglas), aunque la inclusión de los diccionarios ha obligado a filtrarlas y adaptarlas, incrementando su precisión de un 77,5% a un 98,88% (aplicada a un 24,8% de las palabras desconocidas). Se ha incorporado un conjugador verbal de gran cobertura basado en un paradigma sencillo pero efectivo; a pesar de la sobregeneración de este módulo, no se han producido importantes incrementos en el número de etiquetas por palabra.

Trabajando en el nivel de palabra, se han estudiado los distintos tipos de palabras no estándar y la manera de detectar su presencia en un texto, de manera que sea posible procesar no sólo un corpus convenientemente preparado, sino textos obtenidos directamente del dominio sin supervisión. Se ha creado y evaluado un normalizador de texto basado en diccionarios genéricos y diccionarios especializados y en reglas de experto empotradas. La precisión global alcanzada fue del 98,41%, mientras que la precisión sobre las palabras no estándar fue siempre superior al 85% sobre un corpus de evaluación de textos periodísticos, destacando el 96% en nombres propios simples. Para la detección de palabras y nombres propios extranjeros se ha probado un sencillo

método basado en las reglas de silabificación del castellano, cuya precisión supera el 99,5%, aunque cubre pocos casos.

Modelado de F0 en dominio restringido

Se ha estudiado el modelado mediante perceptrones multicapa, destacando la significativa importancia que adquieren la información sobre “la frase portadora” y el “signo de puntuación final del grupo fónico”. El parámetro “número de frase portadora” introduce mejoras significativas; a pesar de que en las condiciones de grabación se intentó aislar el elemento variable de su frase portadora por medio de pausas obligatorias. El parámetro más importante, el “signo de puntuación final del grupo fónico”, es muy relevante porque permite distinguir elementos variables con cadencias y elementos variables que generalmente presentan anticadencias en las grabaciones.

Se ha constatado la importancia de codificar la información inventanada sobre el acento de cada sílaba, así como su situación inicial o final en el grupo fónico. Se han ensayado varias codificaciones alternativas para la misma información sin conseguir superar la tasa. El tamaño óptimo de la ventana es dependiente de la tarea, aunque tiene relación con el tamaño de los grupos fónicos y los datos disponibles.

Se ha desarrollado y evaluado un nuevo método de diseño de bases de datos: por medio de un nuevo algoritmo voraz moderado por medio de subobjetivos parciales, se ha conseguido resumir una gran base de datos con una precisión superior al 95 %, de acuerdo con amplio espectro de vectores prosódico y segmentales.

Es igualmente importante estudiar cómo agrupar las frases en subdominios, (realizar una correcta agrupación de las grabaciones de acuerdo con su prosodia, proponiendo un modelado individual para algunas grabaciones), aunque las diferencias encontradas no han sido significativas.

Parámetros secundarios a la hora de modelar han resultado ser el tamaño del grupo fónico en sílabas o en palabras, la pertenencia de la sílaba a una palabra función o su situación en posición final de palabra. Apenas aportan mejoras; y si las aportan nunca es significativamente ni en todos los subdominios.

Se ha empleado una estrategia de experimentación no exhaustiva, que al ser comparada con la búsqueda exhaustiva del óptimo, ha mostrado su validez.

Análisis y síntesis de voz con emociones

Se han realizado experimentos para determinar la naturaleza segmental o prosódica de las emociones simuladas: en una evaluación pionera confirmada por posteriores experimentos de otros grupos de investigación, hemos mezclado los difonemas y la prosodia de diferentes emociones para concluir la naturaleza segmental del enfado amenazante del actor de nuestra base de datos y la naturaleza prosódica en el caso de la sorpresa; para la alegría y la tristeza se ha revelado como de naturaleza mixta, en parte segmental en parte prosódica.

Se ha creado un sistema completo de conversión texto a voz en castellano, con una nueva voz configurable con emociones: para ello se ha empleado síntesis basada en formantes en castellano, con capacidad de personalización evaluada con usuarios. Los parámetros de personalización han sido elegidos de manera que permitan implementar las emociones como un caso particular de personalización dinámica. Por lo que hemos podido ver, los resultados globales resultan prometedores, puesto que los humanos parecen adaptarse con rapidez a la voz sintética con emociones, y el periodo de

adaptación podría resultar breve y por lo tanto satisfactorio, especialmente en el caso de la tristeza (siendo peor para el enfado). Es cierto que hubo problemas de inteligibilidad en algunos contextos fonéticos, pero los resultados son perfectamente aceptables, incluso si la voz evaluada no resulta totalmente natural. De acuerdo con los resultados de evaluación, la voz sintética con emociones desarrollada en el proyecto VAESS es comparable al estado de la cuestión en otros idiomas a nivel mundial. Aunque en los planteamientos iniciales del proyecto VAESS se consideró que eran independientes los módulos prosódico y segmental del sintetizador, la conclusión de nuestro trabajo dista mucho de ser esta: debemos señalar que las trayectorias de los formantes pueden provocar, al incrementarse la velocidad de elocución, ruidos de naturaleza pseudo-oclusiva, que es necesario limar uno a uno, modificando la reglas segmentales.

Creación de la primera base de datos de habla emotiva simulada en castellano: Está orientada a síntesis prosódica y al análisis de la prosodia en párrafos y frase cortas mediante técnicas paramétricas, y dio lugar a un modelado diferencial de cada emoción respecto a la voz neutra y su evaluación en experimentos de copy-synthesis.

7. Aplicación práctica e interés industrial

A medida que las tecnologías del habla se han ido implicando, cada vez más, como parte integral de aplicaciones prácticas en escenarios reales, se ha hecho patente la necesidad de desarrollar sistemas de conversión texto a voz dotados de gran naturalidad.

Los resultados de esta Tesis han sido exitosamente aplicados en productos como:

- el robot del proyecto nacional Urbano/Ivanhoe (destinado a hacer de guía en el Museo de las Ciencias de Valencia),
- el sintetizador comercial multilingüe de la empresa Telia-Infovox (proyecto europeo VAESS),
- el sintetizador del GTH (registrado y comercializado por la UPM bajo el nombre de Boris, y vendido a empresas nacionales e internaciones como Digra o Ayllón),
- los sistemas de atención telefónica automática de la empresa Natural Vox (proyecto de mejora de su voz femenina).

Relación cronológica de las publicaciones del autor

ARTÍCULOS EN REVISTAS CON ÍNDICE DE IMPACTO EN JCR

1. “*Sesgos cognitivos en el reconocimiento de expresiones emocionales de voz sintética en la alexitimia*” (F. Martínez-Sánchez, J.M. Montero, J. de la Cerra) en *Psicothema* 14(2), pp. 344-349 (ISSN: 0214-9915) 2002 (**Impact Factor en “JCR Social Sciences Edition 2002”: 1,098**)
2. “*Selection of the Most Significant Parameters for Duration Modeling in a Spanish Text-To-Speech System Using Neural Networks*” (R. De Córdoba, J.M. Montero, J. Gutiérrez-Arriola, J.A. Vallejo, E. Enríquez, J.M. Pardo) en *Computer Speech & Language* Vol 16 Number 2 pp. 183-203 (ISSN: 0885-2308) April 2002 (**Impact Factor en “JCR Science Edition 2003”: 0,541**).

ARTÍCULOS EN OTRAS REVISTAS

3. “*Knowledge Combining Methodology for Dialogue Design in Spoken Language Systems*” (Rubén San-Segundo, Juan M. Montero, Javier Macías, Javier Ferreiros y José M. Pardo) en *International Journal of Speech Technology* Vol. 8(1) pp. 45-66 enero 2005 (ISSN: 1381-2416).
4. “*Medidas de confianza en sistemas de diálogo*” (R. San-Segundo, J. Macías, J.M. Montero, J. Ferreiros, R. Córdoba, J.M. Pardo) en *Procesamiento del Lenguaje Natural*, nº 33 pp. 95-102 (ISSN: 1135-5948) Sept. 2004.
5. “*Plataforma de generación semiautomática de sistemas de diálogo multimodales y multilingües: Proyecto GEMINP*” (L.F. D’Haro, R. Córdoba, I. Ibarz, R. San-Segundo, J.M. Montero, J. Macías-Guarasa, J. Ferreiros, J.M. Pardo) en *Procesamiento del Lenguaje Natural*, nº 33 pp. 103-110 (ISSN: 1135-5948) Sept. 2004
6. “*Sistema de comprensión de comunicaciones habladas para el control de tráfico aéreo del proyecto INVOCA*” (V. Sama Rojo, F. Fernández Martínez, J. Ferreiros López, J. Macías-Guarasa, R. De Córdoba, J. M. Montero Martínez, J. Colas Pasamontes, E. Campos Palarea, J. M. Pardo Muñoz) en *Procesamiento del Lenguaje Natural*, nº 31 pp.337-338 (ISSN: 1135-5948) Sept. 2003.

CAPÍTULOS DE LIBROS INTERNACIONALES

7. “*The role of pitch and tempo in Spanish emotional speech: towards concatenative synthesis*” (Juan Manuel Montero, Juana Gutiérrez-Arriola, Ricardo de Córdoba, Emilia Enríquez, José Manuel Pardo) incluido en “*Improvements in speech synthesis*” de los editores Eric Keller y Gerard Bailey, A. Monahan, J. Terken, M. Huckvale (ISBN 0-471-49985-4) pp. 246-251 editado por John Wiley & Sons, Ltd. en el año 2002.
8. “*Application of neural networks to duration modelling in a Spanish text-to-speech system*” (Ricardo de Córdoba, Juan Manuel Montero, José Manuel Pardo) incluido en “*Advances in Systems Engineering, Signal Processing and Communications*” pp. 244-247 editado por WSEAS Press (ISBN 960 8052 696). en el año 2002.

OTRAS PUBLICACIONES INTERNACIONALES REVISADAS

9. “*Parameter Selection for Prosodic Modeling in a Restricted-Domain Spanish Text-to-Speech System*” (J. M. Montero, J. Macías-Guarasa, R. De Córdoba, J. Gutiérrez-Arriola, J. M. Pardo, R. San Segundo) en *Proceedings of World Automation Congress (WAC 2004)* pp. 155-160 (ISBN: 1-889335-23-1) Sevilla.
10. “*ANN F0 Modelling for Female-Voice Synthesis in Spanish: restricted and non-restricted domains*” (J.M. Montero, L.F. d’Haro, R. de Córdoba, J.A. Vallejo, J. Gutiérrez-Arriola, J.M. Pardo) en *Proceedings of the XVth International Congress of Phonetic Sciences* pp. 563-566. Agosto de 2003, Barcelona. (ISBN: 1-876346-48-5) editado en el año 2003.
11. “*ANESTTE: a writer’s assistant for a specific purpose language*” (J.M Montero, M.M. Duque) en *University centre for Computer corpus REserarch on Language Technical Papers Volumen 16 - Special Issue: Proceedings of Corpus Linguistics 2003 Conference* de los editores Dawn Archer, Paul Rayson, Andrew Wilson and Tony McEnery. pp. 544-551. Marzo de 2003, Lancaster (ISBN: 1-86220-131-5) editado en el año 2003.

12. *"Restricted-Domain Female-Voice Synthesis in Spanish: from Database Design to ANN Prosodic Modelling"* (J.M. Montero, R. Córdoba, J.A. Vallejo, J. Gutiérrez-Arriola, J.M. Pardo) en *Proceedings of the International Conference on Spoken Language Processing'2000* pp. 621-624 (ISBN 7-80150-114-4/G.18) editado en el año 2000.
13. *"Development of an emotional speech synthesiser in Spanish"* (J.M. Montero, J. Gutiérrez-Arriola, J. Colás, J. Macías, E. Enriquez, J.M. Pardo) en *Eurospeech'99 Proceedings* pp. 2099-2102 (ISSN: 1018-4074) Budapest 1999.
14. *"Analysis and Modelling of Emotional Speech in Spanish"* (J.M. Montero, J. Gutiérrez-Arriola, J. Colás, E. Enriquez, J.M. Pardo) en *actas de XIVth International Congress of Phonetic Sciences Vol. II* pp. 957-960. Agosto de 1999, San Francisco.
15. *"Emotional Speech Synthesis: from speech database to TTS"* (J.M. Montero, J. Gutiérrez-Arriola, S. Palazuelos, E. Enriquez, S. Aguilera, J.M. Pardo) en *Proceedings of the International Conference on Spoken Language Processing '98* (ISBN 1-876346-17-5) pp. 923-925 en el año 1998.
16. *"Generating Gestures from Speech"* (R. San-Segundo, J.M. Montero, J. Macías-Guarasa, R. de Córdoba, J. Ferreiros, J.M. Pardo) en *INTERSPEECH 2004 – ICSLP* (ISSN 1225-441x), Korea 2004 pp. 1817-1820.
17. *"Analysis of Parameter Importance in Speaker Identity"* (J. Gutiérrez-Arriola, J.M. Montero, R. Córdoba, J.M. Pardo) Artículo publicado en el *ITRW on Voice Quality: functions, analysis and synthesis (VoQual)*, pp. 103-106. (ISSN: 1680-8908). Ginebra, 2003.
18. *"Methodology for Dialogue Design in Telephone-Based Spoken Dialogue Systems: a Spanish Train Information System"* (R. San-Segundo, J.M. Montero, J. Colás, J. Gutiérrez-Arriola, J.M. Ramos, J.M. Pardo) en *Proceedings of EUROSPEECH'01* pp. 2165-2168 (ISBN: 87-90834-09-7) en Aalborg septiembre de 2001.
19. *"An Interactive Directory Assistance Service for Spanish with Large-Vocabulary Recognition"* (R. Córdoba, R. San-Segundo, J.M. Montero, J. Colás, J. Ferreiros, J. Macías-Guarasa, J.M. Pardo) en *Proceedings of EUROSPEECH'01* pp. 1279-1282 (ISBN: 87-90834-09-7), en Aalborg septiembre de 2001.
20. *"A New Multi-speaker Formant Synthesizer that Applies Voice Conversion Techniques"* (J.M. Gutiérrez, J.M. Montero, J.A. Vallejo, R. de Córdoba, R. San Segundo and J.M. Pardo.) en *Proc. Eurospeech 2001*, pp 357-360 (ISBN 87-90834-09-7) Aalborg (Denmark) Sept. 2001.
21. *"New Rule-Based and Data-Driven Strategy to Incorporate Fujisaki's F0 Model to a Text-To-Speech System in Castillian Spanish"* (J. Gutiérrez-Arriola, J.M. Montero, D. Saiz, J.M. Pardo). en *Proceedings of the International Conference on Acoustics and Signal Processing ICASSP' 2001* pp. 821-824 (ISBN: 0-7803-7043-0) en Salt Lake City 2001.
22. *"Duration Modeling in a Restricted-Domain Female-Voice Synthesis in Spanish Using Neural Networks"*, (R. Córdoba, J.M. Montero, J. Gutiérrez-Arriola, J.M. Pardo) en *Proceedings of ICASSP' 2001* pp. 793-796 (ISBN: 0-7803-7043-0) en el año 2001.
23. *"Sistema de información ferroviaria por teléfono: propuesta de una metodología de diseño de gestores de diálogo"* (R. San-Segundo, J.M. Montero, J. Ferreiros, J. Macías-Guarasa y J.M. Pardo) en *actas de "Second International Workshop on Spanish Language Processing and Language Technologies (SLPLT-2)"* (ISBN: 84-8454-095-2) pp.241-245 Septiembre 2001, Jaén.
24. *"Designing confirmation mechanisms and error recover techniques in a railway information system for Spanish"* (R. San-Segundo, J.M. Montero, J. Ferreiros, R. Córdoba, J.M. Pardo.) en *Proceedings of the 2nd SIGdial Workshop on Discourse and Dialogue* pp.: 136-139. Aalborg, Denmark. 1-2 Sept. 2001.
25. *"A telephone-based railway information system for Spanish: development of a methodology for spoken dialogue design"* (R. San-Segundo, J.M. Montero, J.M. Gutiérrez, A. Gallardo, J.D. Romeral and J.M. Pardo) en *Proceedings of the 2nd SIGdial Workshop on Discourse and Dialogue* pp.: 140-148. Aalborg, Denmark. 1-2 Sept. 2001
26. *"Stress Assignment in Spanish Proper Names"* (R. San Segundo, J.M. Montero, R. Córdoba, J.M. Gutiérrez-Arriola) en *Proceedings of the International Conference on Spoken Language Processing'2000* pp. 346-349 (ISBN 7-80150-114-4/G.18) editado en Pekín en el año 2000.

27. "Automatic Modeling of Duration in a Spanish Text-to-Speech System Using Neural Networks" (R. Córdoba, J. A. Vallejo, J.M. Montero, J. Gutiérrez-Arriola, M.A. López, J.M. Pardo) en *Eurospeech '99 Proceedings* (ISSN: 1018-4074) en Budapest en el año 1999 (4 páginas).
28. "On the limitations of Stochastic Conceptual Finite-State Language Models for Speech Understanding" (J. Colás, J. Ferreiros, J.M. Montero, J.M. Pardo) en *Proceedings of the International Conference on Spoken Language Processing '98* (ISBN 1-876346-17-5) pp. 2239-2242 en el año 1998.
29. "Voice Conversion Based on Parameter Transformation" (Gutiérrez-Arriola J.M., Y.S. Hsiao, J.M. Montero, J.M. Pardo, D.F. Childers) en *Proceedings of the International Conference on Spoken Language Processing '98* (ISBN 1-876346-17-5) pp. 987-990 en el año 1998.
30. "An alternative and flexible approach in robust information retrieval systems" (J. Colás, J.M. Montero, J. Ferreiros, J.M. Pardo) en *Eurospeech '97 Proceedings* (ISSN: 1018-4074) pp. 2683-2686 en el año 1997.

ARTÍCULOS EN OTROS CONGRESOS NACIONALES

31. "Proyecto GEMINI: Plataforma avanzada de generación y ejecución de aplicaciones de diálogo hombre-máquina" (R. Córdoba, L. F. D'Haro, F. Fernández, V. Sama, J. M. Montero, R. San-Segundo, J. Macías-Guarasa, J. Ferreiros y J. M. Pardo) en *XIV Jornadas Telecom I+D. Barcelona - Madrid*. (ISBN:), editado en el año 2004.
32. "¿Podemos imitar la voz de una persona? Técnicas de conversión de hablante" (J. Gutiérrez-Arriola, J.M. Montero, R. de Córdoba, J.M. Pardo) en las *Actas del II Congreso de la Sociedad Española de Acústica Forense (SEAF)*. Barcelona pp. 35-46. Abril de 2003.
33. "Generación semiautomática de aplicaciones de diálogo multimodales: proyecto Gemini" (R. Córdoba, L.F. D'Haro, J.M. Montero, J. Ferreiros, J. Macías-Guarasa, y J.D. Romeral) en *XIII Jornadas Telecom I+D. Barcelona - Madrid*. (ISBN: 84 89315 28 0), editado en el año 2003.
34. "Entorno para el desarrollo de aplicaciones multimedia con síntesis y reconocimiento de voz" (R. San-Segundo, J.M. Montero, J. Colás, J. Ferreiros, R. Córdoba, A. Gallardo, J. Macías-Guarasa, J.M. Gutiérrez, J. Pastor, J.M. Pardo) en *X Jornadas Telecom I+D. Barcelona - Madrid*. noviembre de 2000. Madrid (7 páginas).
35. "Optimización de un servicio automático de páginas blancas por teléfono: proyecto IDAS" (R. Córdoba, R. San-Segundo, J. Colás, J.M. Montero, J. Ferreiros, J. Macías-Guarasa, A. Gallardo, J.M. Gutiérrez, J.M. Pardo) *X Jornadas Telecom I+D. Barcelona - Madrid*. Artículo publicado en *X Jornadas Telecom I+D. Barcelona - Madrid*. (ISBN: 84-607-1397-0), editado en el año 2000 (8 páginas).
36. "Servidores Vocales Interactivos: Desarrollo de un servicio de paginas blancas por teléfono con reconocimiento de voz (proyecto IDAS: Interactive telephone-based Directory Assistance Service)" (R. San-Segundo, J. Colás, J.M. Montero, R. Córdoba, J. Ferreiros, J. Macías-Guarasa, A. Gallardo, J.M. Gutiérrez, J. Pastor y J.M. Pardo). en *IX Jornadas Telecom I+D. Barcelona - Madrid*. (ISBN: 84-7653-730-1), editado en el año 1999 (6 páginas).

Otros méritos relacionados con la tesis doctoral

Premios

- Premio al Mejor artículo en una sesión oral en el apartado de "Aplicaciones y Tecnología Multimedia" en las IX Jornadas Telecom I+D. Barcelona - Madrid 1999 (coautor).
- Premio al Mejor artículo en una sesión oral en el apartado de "Servicios, aplicaciones y contenidos multimedia" en las X Jornadas Telecom I+D. Barcelona - Madrid 2000 (coautor).

Cursos y seminarios impartidos

- Seminario "Curso de Tecnologías Lingüísticas: Preguntar al ordenador. Las aplicaciones de los sistemas de diálogo" (del 12 al 16 de Julio de 2004) organizado por la Fundación Duques de Soria Ponencia "La incorporación de emociones a los sistemas de diálogo"

- Curso del Servicio de Promoción Educativa: Fronteras: II Curso: La variación sociolingüística. Enfoque contrastivo, conferencia sobre “Lingüística y tecnología: síntesis y reconocimiento del habla” (Universidad de Murcia). Noviembre de 2000.

Revisor o chairman de revistas y congresos internacionales

- Revisor de revistas internacionales (Computers & Electrical Engineering, Intelligent Automation And Soft Computing) y congresos internacionales (TELEC-2004, EISTA-2005)
- Chairman de sesión en congreso internacional *ICPhS 2003* en Barcelona.

Referencias al trabajo por parte de otros investigadores

Citas a J.M. Montero et al 1998 (ICSLP'1998)

- Louis ten Bosch 2003 "Emotions, speech and the ASR framework" en *Speech Communication*, vol 40, pp 213-225.
- Gábor Olaszky 2000 “The Prosody Structure of Dialogue Components in Hungarian” en *International Journal of Speech Technology* 3 (3-4): 165-176
- Ilona Koutny, Gábor Olaszky y Péter Olaszi 2000 “Prosody Prediction from Text in Hungarian and its Realization in TTS Conversion” en *International Journal of Speech Technology* 3 (3-4): 187-200
- J.M. Sosa 1999 La entonación del español Ed. Cátedra
- Jianhua Tao 2004 "Context Based Emotion Detection from Text Input" en *ICSLP 2004*, pp. 1120-1123
- Marc Schröder 2001 “Emotional Speech Synthesis: A Review” en *Proceedings of Eurospeech* pp.561-564.
- Louis ten Bosch 2000 “Emotions: What Is Possible In The Asr Framework” en *Proceedings of the ISCA Workshop on Speech and Emotion* pp. 189-194.
- Nestor Garay, Julio Abascal, Luis Gardezabal 2002 “Mediación emocional en sistemas de Comunicación Aumentativa y Alternativa” en *Inteligencia Artificial. Revista Iberoamericana de Inteligencia Artificial*, nº 16 (ISSN: 1137-3601)
- Marc Schröder 1999 “Can Emotions Be Synthesized without Controlling Voice Quality?” en *PHONUS* 4, pp. 37-55 *Research Report of the Institute of Phonetics, University of the Saarland*
- Tom Brøndsted, Thomas Dorf Nielsen, Sergio Ortega 1999 “Classification of Emotional Attitudes in Pet-directed Speech” en *Proceedings DALF (The Danish Society for Computational Linguistics)*.
- Tom Brøndsted, Thomas Dorf Nielsen, Sergio Ortega 1999 “Affective MultiModal Interaction with a 3D Agent” en *The Eighth International Workshop on the Cognitive Science of Natural Language Processing*, pp. 102-109
- "Emotion Control of Chinese Speech Synthesis in Natural Environment", Jianhua Tao. pp. 2349-2352 *Eurospeech 2003*
- P.N.Girija and M.Neeraja 2003 “Intonation knowledge for declarative Tegulu sentences” en *Proceedings of Oriental COCOSDA2003* ISBN:981-04-9596-X

Citas a J.M. Montero et al 1999 (ICPhS)

- Murtaza Bulut, Shrikanth Narayanan y Lewis Johnson 2004 “Synthesizing expressive speech: overview, challenges and open questions” en el libro “*Text to Speech Synthesis. New Paradigms and Advances*” de Prentice Hall.
- Murtaza Bulut, Shrikanth S. Narayanan, Ann K. Syrdal 2002 “Expressive Speech Synthesis Using A Concatenative Synthesizer” en *Proceedings of ICSLP: sesión WeC28o (oral)*
- Gendrot Cédric 2002 “Ouverture de la glotte. Fo, intensité et simulations émotionnelles : le cas de la joie, la colère, la surprise, la tristesse et la neutralité” en *XXIVèmes Journées d’Etude sur la Parole JEP2002*

- Philippe Boula de Mareüil, Philippe Célérier, Jacques Toen 2002 "Generation of Emotions by a Morphing Technique in English, French and Spanish" en Proceedings of Speech Prosody. Aix-en-Provence, pp. 187-190
- Marc Schröder and Martine Grice "Expressing vocal effort in concatenative synthesis" en Proceedings of ICPhS'2003 pp. 2589-2592.
- Asa Abelin & Jens Allwood 2000 "Cross Linguistic Interpretation of Emotional Prosody" en Proceedings of the ISCA Workshop on Speech and Emotion pp.110-113
- D Ververidis, C Kotropoulos "A State of the Art Review on Emotional Speech Databases" en Proc. 1st Richmedia Conference
- Ignasi Iriondo, Francesc Alías , Javier Melenchón and M. Angeles Llorca "Modeling and Synthesizing Emotional Speech for Catalan Text-to-Speech Synthesis" en Affective Dialogue Systems: Tutorial and Research Workshop, ADS 2004, pp. 197 – 208
- J. Llisterra 2001, "La conversión de texto en habla" en Quark, Ciencia, Medicina, Comunicación y Cultura 21 pp. 79-89

Citas a J.M. Montero et al 1999b (Eurospeech)

- Felix Burkhardt and Walter F. Sendlmeier 2000 "Verification of Acoustical Correlates of Emotional Speech using Formant-Síntesis" en Proceedings of the ISCA Workshop on Speech and Emotion pp. 151-156 y en "Speech and Signals - Aspects of Speech Synthesis and automatic speech recognition. . Dedicated to Wolfgang Hess on his 60th Birthday", W.F.Sendlmeier (Ed.), Frankfurt: Hector (Forum Phonicum, 69). pp. 27-39.
- Felix Burkhardt "Simulation emotionaler Sprechweise mit Sprachsyntheseverfahren" tesis doctoral 2000

Citas a J.M. Montero et al 2000 (ICSLP)

- X Sevillano, F. Alías , P. Barnola, J.C. Socoró 2003 "ICA-based hierarquical text classification for multi-domain text-to-speech synthesis" en Proceedings of ICASSP2004 pp. V-697-700
- F. Alías, X Sevillano, P. Barnola, J.C. Socoró 2003 "Arquitectura para conversión texto-habla multidominio" en Sociedad Española para el Procesamiento del Lenguaje Natural (SEPLN); Procesamiento del Lenguaje Natural, N° 31, Septiembre 2003, pp. 83-90. ISSN: 1135-5948;

Citas a otros artículos

- R. San-Segundo, J. M. Montero, J. Gutiérrez et al (SIGDial 2001 Workshop) **citado en** "Issue-based dialogue management" de Staffan Larsson. (Tesis Doctoral. Department of linguistics Göteborg University, Sweden, 2002); **citado en** David Schlangen "Causes and Strategies for requesting Clarification in Dialogue" en 5th SIGdial Workshop on Discourse and Dialogue, Cambridge, 2004, pp. 136-143, ISBN: 1-932432-27-2; **citado en** Hiroya Murao, Nobuo Kawaguchi, Shigeki Matsubara, Yukiko Yamaguchi Y Yasuyoshi Inagakiz "Example-based Spoken Dialogue System using WOZ System Log" en SIGDIAL 2003; **citado en** "Automatic Design of Spoken Dialogue Systems" K. Scheffler (Tesis Doctoral Cambridge University Engineering Department)
- R. Córdoba, R., San-Segundo, J.M., Montero et al (Eurospeech 2001). **citado en** F. Torres, E. Sanchis, E. Segarra "Development of a stochastic dialog manager driven by semantics", en European Conference on Speech Communication and Technology (Eurospeech) 2003, pp. 605-608, ISSN: 1018-4074; **citado en** K. Georgila, K. Sgarbas, A. Tsopanoglou, N. Fakotakis "A Speech-Based Human-Computer Interaction System for Automating Directory Assistance Services" en International Journal of Speech Technology 6 (2): 145-159, April 2003; **citado en** K Georgila, N Fakotakis, G Kokkinakis "Large Vocabulary Search Space Reduction Employing Directed Acyclic Word Graphs and Phonological Rules" en International Journal of Speech Technology 5 (4): 355-370, November 2002;
- R. de Córdoba, J.A. Vallejo, J.M. Montero et al (Eurospeech 1999) **citado en** "Segmental Durations Predicted With a Neural Network", de J.P. Teixeira, D. Freitas. en European Conference on Speech Communication and Technology (Eurospeech) 2003, pp. 169-172, ISSN: 1018-4074; **citado en** "Segmental duration control with asymmetric causal retro-causal neural networks", de C. Erdem, H.G. Zimmermann. en 4th ISCA ITRW on Speech Synthesis (SSW-4), 2001, paper 119; **citado en** João Paulo Teixeira y Diamantino Freitas "Evaluation of a Segmental Durations Model for TTS" en

Lecture Notes in Computer Science ISSN: 0302-9743 Volume 2721 / 2003 pp. 40 – 48; **citado en** Eva Navas, Inmaculada Hernandez y Juan Marıa Sanchez “Modelo de duracion para conversion de texto a voz en euskera” en SELPN 2002 pp. 1-15;

- J.M. Gutierrez-Arriola, Y.S. Hsiao, J.M. Montero et al. (ICSLP 1998) **citado en** “Voice Conversion Methods for Vocal Tract and Pitch Contour Modification”, de Oytun Turk, Levent M. Arslan, pp. 2845-2848 ;
- J. M. Gutierrez-Arriola, J. M. Montero, et al (ICASSP 2001) **citado en** en “Inversion Of Model For Natural-Sounding Speech Synthesis” de Pierluigi Salvo Rossi, Francesco Palmieri, Francesco Cutugno en ICASSP 2003, I-520-523; **citado en** P.S. Rossi, F. Palmieri y F. Cutugno “A method for automatic extraction of Fujisaki-model parameters” en Speech Prosody 2002; **citado en** Solimar de S. Silva and Sergio L. Netto “Closed-Form Estimation Of The Amplitude Commands In The Automatic Extraction Of The Fujisaki’s Model” en ICASSP 2004;
- San-Segundo, R. et al “Servidores Vocales Interactivos: Desarrollo de un Servicio de Paginas Blancas por Telefono con Reconocimiento de Voz” en IX Jornadas Telecom I+D, Barcelona-Madrid (1999) **citado en** Jordi lvarez, Victoria Arranz, Nuria Castel, Montserrat Civit “Linguistic and Logical Tools for an Advanced Interactive Speech System in Spanish” en Lecture Notes in Computer Science (International Conference on Industrial and Engineering Applications of Artificial Intelligence and Expert Systems, IEA/AIE) 2001 pp. 519-528;

Proyectos a los que se asocia la tesis:

Proyectos internacionales con financiacion publica

1. **Proyecto VAESS:** “*Voices, Attitudes and Emotions in Speech Synthesis*”, cuyo Investigador Responsable fue Dr. Jose Manuel Pardo Munoz, y que fue financiado por la Union Europea. Referencia: TIDE TP 1174. Duracion: 1994-97 Puesto: **Responsable Tecnico / Investigador.**
2. **Proyecto GEMINI:** *Generic Environment for Multilingual Interactive Natural Interfaces*, cuyo Investigador Responsable es Dr. Jose Manuel Pardo Munoz, y que es financiado por la Union Europea. Referencia: IST-2001-32343 Duracion: 2002-2004. Puesto: **Investigador.**
3. **Proyecto IDAS:** *Interactive telephone-based Directory Assistance Services*, cuyo Investigador Responsable fue Dr. Jose Manuel Pardo Munoz, y que fue financiado por la Union Europea (IV programa marco “*Telematic applications*”) y CICYT. Referencia: LE4-8315. Duracion: 1998-2000. Puesto: **Investigador.**

Proyectos nacionales con financiacion publica

1. **Proyecto ROBINT:** *Tecnologıa del Habla* cuyo Investigador Responsable es Dr. Juan Manuel Montero Martınez, y que es financiado por el MEC. Duracion: 13/13/2004-13/12/2006. Referencia: DPI2004-07908-C02-02. Puesto: **Investigador Principal.** Participantes: DIE ETSIT-UPM, DISAM-ETSII-UPM, Museo Prncipe Felipe-CAC S.A.
2. **Proyecto IVANHOE:** *Interfaces vocales con robot anfitron para visitas presenciales o remotas a exposiciones* cuyo Investigador Responsable es Dr. Jose Manuel Pardo Munoz, y que es financiado por el MCyT. Duracion: 28/12/2001- 27/12/2004. Referencia: DPI2001-3652-C02-02. Puesto: **Investigador**
3. **Accion especial “Red temtica en Tecnologıas del Habla”** cuyo Investigador Responsable es Dr. Antonio Jose Rubio Ayuso, y que es financiado por la MCyT. Duracion: 11/6/2003- 10/6/2005. Referencia: TIC2002-11271-E. Puesto: **Investigador.**
4. **Proyecto SAITE:** *Tecnologıa para el desarrollo de Servicios Avanzados de Informacion Telefonica*, cuyo Investigador Responsable fue Dr. Jose Manuel Pardo Munoz, y que es financiado por la Union Europea a traves de Fondos FEDER Referencia: 2FD1997-1062-C02. Duracion: 1/12/1999-30/11/2001. Puesto: **Investigador**
5. **Proyecto DEMSTENES:** “*Hacia la naturalidad en sntesis de habla a partir de texto*”, cuyo Investigador Responsable fue Dr. Santiago Aguilera Navarro, y que fue financiado por la CICYT. Referencia: TIC 95-0147. Duracion: 1995-1998 Puesto: **Investigador**

Aplicacion prctica e inters industrial

Propiedad intelectual registrada

- Coautor de la obra SW registrada: "Boris. Conversor texto voz del GTH", M-001714/2003 (Registro de la Propiedad Intelectual, Madrid). Autores: Ricardo de Córdoba Herralde, José Manuel Pardo Muñoz, Juan Manuel Montero Martínez y Juana María Gutiérrez Arriola. Entidad titular de los derechos: UPM. Este producto ha sido vendido a diversas empresas nacionales e internacionales (Digra, Ayllón...).

Contratos con empresas

1. **Proyecto Mejora de calidad de síntesis de voz femenina**, cuyo Investigador Responsable fue Dr. José Manuel Pardo Muñoz, y que fue financiado por la empresa Natural Vox. Duración: 1998-1999. Puesto: **Responsable Técnico / Investigador**.
2. **Proyecto INVOCA: Interfaces vocales para control de tráfico aéreo**, cuyo Investigador Responsable es Dr. José Manuel Pardo Muñoz, y que es financiado por AENA. Duración: 2001-2003. Puesto: **Investigador**.
3. **Proyecto SERVIVOX, Sistema para la automatización de servicios telefónicos**, cuyo Investigador Responsable fue Dr. José Manuel Pardo Muñoz, y que fue financiado por la empresa Hewlett Packard Española. Duración: 1997-1998. Puesto: **Responsable Técnico / Investigador**.
4. **Proyecto VOZ: Sistema automático de información telefónica**, cuyo Investigador Responsable fue Dr. José Manuel Pardo Muñoz, y que fue financiado por la UPM. Duración: 1995-1998. Puesto: **Responsable Técnico / Investigador**