

DEPARTAMENTO DE INGENIERÍA ELECTRÓNICA

**ESCUELA TÉCNICA SUPERIOR DE INGENIEROS DE
TELECOMUNICACIÓN**



TESIS DOCTORAL

**ESTRATEGIAS DE INCORPORACIÓN DE
CONOCIMIENTO SINTÁCTICO Y SEMÁNTICO EN
SISTEMAS DE COMPRESIÓN DE HABLA CONTINUA
EN CASTELLANO**

José Colás Pasamontes

Ingeniero de Telecomunicación

Director de la Tesis

Dr. Ingeniero José Manuel Pardo Muñoz

1.999

Tesis Doctoral: **ESTRATEGIAS DE INCORPORACIÓN DE CONOCIMIENTO SINTÁCTICO Y SEMÁNTICO EN SISTEMAS DE COMPRESIÓN DE HABLA CONTINUA EN CASTELLANO**

Autor: **JOSÉ COLÁS PASAMONTES**

Director: **Dr. INGENIERO JOSÉ MANUEL PARDO MUÑOZ**

El tribunal nombrado para juzgar la tesis doctoral arriba citada, compuesto por los doctores:

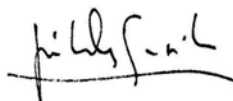
Presidente: **JAVIER FERREIROS**



Vocales: **NATIVIDAD PRIETO**
LUIS A. HERNANDEZ
ANTONIO RUBIO



Secretario: **JOSE CARLOS GONZALEZ**



acuerda otorgarle la calificación de **Sobresaliente cum Laude (5 votos)**

Madrid, a 1 de julio de 1.999

RESUMEN

Este trabajo de tesis doctoral se ha planteado analizar la problemática del proceso de **comprensión de habla**, no sólo desde un punto de vista científico sino también técnico, concluyendo con el **diseño, implementación y evaluación de un Sistema de Comprensión de Habla en castellano**.

En el marco de esta tesis se ha realizado una **revisión de distintas soluciones que han sido propuestas por diversos grupos de investigación internacionales para resolver el problema de la comprensión de habla**.

Se ha definido una **arquitectura no integrada novedosa para la comprensión del habla en castellano**, es decir, teniendo en cuenta características del castellano como lengua natural que no aparecen o aparecen menos acentuadas en otras lenguas. Esta arquitectura pretende ser la base de futuros trabajos en esta línea en el Grupo de Tecnología del Habla, de la Universidad Politécnica de Madrid. Con el objetivo de comprender habla perteneciente a un **dominio semántico restringido** (limitado por los conceptos que existen en el dominio de una aplicación concreta), esta arquitectura tiene como características principales:

- La **robustez**, es decir, la posibilidad de procesar frases que contienen errores (inserciones, borrados o sustituciones de palabras) producidos por el módulo de decodificación acústica (sistema de reconocimiento de habla continua), o que tienen rasgos de agramaticalidad producidos por la propia naturaleza del lenguaje hablado, o con problemas de cobertura a nivel léxico, sintáctico o semántico.
- La **modularidad**, que permite seguir mejorando sin necesidad de rediseñar e implementar el sistema completo.
- La **flexibilidad**, con el fin de independizar la arquitectura de una aplicación concreta dentro, lógicamente, de ciertas restricciones, impuestas por la naturaleza de los sistemas de información o sistemas de control automáticos.
- La **potencia**, definida como la posibilidad de procesar frases de un cierto nivel de complejidad lingüística.

Los diferentes módulos incorporan conocimiento lingüístico de distinta naturaleza, lo que ha permitido **estudiar la interacción de distintas fuentes de conocimiento lingüístico y un modo eficaz de integrarlas**, en el proceso de comprensión.

Se ha utilizado **información semántica en forma de rasgos**, que completan la ya modelada por las categorías semánticas del diccionario, a **gramáticas contextuales simplificadas** (definiendo **lenguajes específicos de reglas y algoritmos de análisis o ejecución de estas reglas**), que en forma de reglas solucionan principalmente problemas de ambigüedad semántica y elipsis, y una **gramática semántica de contexto libre** (utilizando el algoritmo de Earley con capacidad para procesar frases con ambigüedad) que pretende, basándose en una clasificación o taxonomía de los

conceptos del dominio que reduce en gran medida las reglas necesarias, obtener la información estructural de las mismas que ayuda al procesamiento de frases de una cierta complejidad manteniendo el proceso de **traducción a SQL, necesario en sistemas de información con acceso a bases de datos, dentro de unos límites de sencillez sorprendentes**, mediante el uso de **plantillas semánticas**.

Con el fin de evaluar el comportamiento del módulo de decodificación acústica se ha implementado un sistema de reconocimiento de habla continua modular, con capacidad para integrar conocimiento gramatical en base a cualquier gramática probabilística de tipo N-gram, de naturaleza morfo-sintáctica o semántica. Se han evaluado distintas gramáticas guiando el proceso de decodificación acústica. Con el fin de mantener la eficacia del este módulo a pesar de la incorporación de gramática en el proceso, **se ha estudiado con profundidad un mecanismo de reducción del espacio de búsqueda ampliamente utilizado conocido como “recorte de caminos” o “recorte del haz” (*beam-search*)**, presentando un método que se ha desarrollado en esta tesis que permite la determinación del *umbral de recorte* basado en la probabilidad (o distancia) del mejor estado del espacio de búsqueda en cada trama (estadístico) de antemano, utilizando los datos de entrenamiento y conociendo la influencia que tendrá dicho umbral en el proceso de reconocimiento. Además, se han evaluado dos variantes ya conocidas: el uso de *uno o dos umbrales de recorte*, uno basado en la probabilidad (o distancia) del mejor de los últimos estados de cada modelo en cada trama (estadístico del último estado) y otro en la probabilidad (o distancia) del mejor del resto de los estados distintos del último en cada trama (estadístico del resto de los estados), y se han aportado nuevas conclusiones al respecto. Todo ello ha permitido profundizar en el funcionamiento de esta técnica ya conocida pero no tan estudiada. Además, el decodificador acústico ha sido modificado para permitir la generación de varias hipótesis (frases) de salida (las N mejores), y **se ha estudiado la relación entre el valor de N (número de caminos o hipótesis) y la calidad del sistema de reconocimiento (mejora de la tasa de acierto de palabras o reducción del error del sistema)**, para aplicaciones como la que ha sido objeto en esta tesis. **Se ha comprobado que con un número de hipótesis reducido (N muy pequeño) se consigue que el módulo acústico se recupere de muchos errores que afectarían al proceso de comprensión de la frase hablada reconocida.**

ABSTRACT

This Ph.D. thesis work is aimed at analyzing the problems when facing **automatic speech understanding**, from both scientific and technical points of view, concluding with the **design, implementation and evaluation of a Castilian Spanish Speech understanding system**.

In this thesis, **some of the alternatives, that have been proposed by international research groups in order to solve the speech understanding problem, have been reviewed.**

A **novel non integrated architecture for speech understanding in Spanish** has been defined, taking into account the specific characteristics of Spanish as a natural language, not found or rarely found in other languages. This architecture intends to be the baseline of future work in this topic in the Speech Technology Group, in the Universidad Politécnica de Madrid.

To achieve the objective of understanding speech in **limited semantic domains** (limited by the concepts used in the domain of a specific application) this architecture has been designed with the following main characteristics:

- **Robustness**, that is, the possibility of processing sentences with errors (word insertions, deletions or substitutions) produced by the acoustic decoding module (a continuous speech recognition system); or non-grammatical constructions, due to the inherent characteristics of spoken language; or problems in lexical, syntactic or semantic coverage.
- **Modularity**, that permits improving the system without redesigning or implementing the whole system.
- **Flexibility**, in order to have an application-independent architecture, obviously under certain restrictions, imposed by the characteristics of both automatic information and control systems.
- **Power**, defined as the possibility of processing sentences with a certain degree of linguistic complexity.

Those modules incorporate linguistic knowledge of different kinds, and this has allowed us to **study the interaction of different linguistic knowledge sources and an efficient way of integrating them** in the understanding process.

Features to represent the semantic information have been used, completing the one already modelled by the dictionary semantic categories; simplified **contextual grammars** (defining **specific rules languages, and rule analysis or execution algorithms**), which mainly solve some of the semantic ambiguity and ellipsis problems; and a **semantic context free grammar** (using the Earley algorithm with its possibility of processing ambiguous sentences). The latter intends to obtain the structural information of the sentence, using a taxonomy of the domain concepts that heavily reduces the number of needed rules. Moreover, it helps the processing of complex

sentences, while keeping the SQL translation process surprisingly simple, by using **semantic templates**. This translation process is **needed in information systems accessing databases**.

In-order to evaluate the acoustic decoder module behaviour, a modular continuous speech recognition system has been implemented. It is able to integrated grammatical knowledge based on any stochastic morpho-syntactic or semantic N-gram. To keep the efficiency of this module, even when the grammar information is used, **a search space reduction mechanism (beam-search) has been deeply studied**. A new method developed in this Thesis allows to analyse and to determine, in advance, a *pruning threshold* based on the probability (or distance) of the best state in the search space for every frame (stochastic), making use of training data and knowing the impact this threshold will have in the recognition process. Besides, two well-known variants have been evaluated: the use of *one or two pruning thresholds*, one based on the probability (distance) of the best last states for every model in every frame (stochastic parameter of the last state) and the other one based on the probability of the best of the rest of the states in every frame (stochastic parameter of the rest of the states). New conclusions have been drawn from this study. All this allowed us to deepen in this well known but not so well understood technique. Moreover the acoustic decoder has been modified to allow the generation of several output hypothesis (N-best sentences), and **the relationship between the value N (number of paths or hypothesis) and the speech recognition system performance (improvement of the word error rate)**, for applications such as the one aimed in this Thesis. **We have checked that with a small number of hypothesis (very low N), the acoustic module is able to recover from a lot of errors that would severely affect the understanding process of the recognised spoken sentence.**

**ESTRATEGIAS DE INCORPORACIÓN DE CONOCIMIENTO SINTÁCTICO
Y SEMÁNTICO EN SISTEMAS DE COMPRESIÓN DE HABLA CONTINUA
EN CASTELLANO**

ÍNDICE

Agradecimientos	V
Resumen	VII
Abstract	IX
Indice	XI
Indice de Figuras	XVII
Indice de Tablas	XXI
CAPÍTULO 1. INTRODUCCIÓN	
1.1 Objetivos de la Tesis. Justificación	1-1
1.2 Contenido de la Tesis	1-3
CAPÍTULO 2. ENCUADRE CIENTÍFICO-TECNOLÓGICO	
2.1 Introducción	2-1
2.2 La Comprensión del Habla: Un Problema Abierto	2-3
2.3 Elección del Dominio	2-4
2.4 Definición de una Arquitectura para la Integración	2-5
2.4.1 Reconocimiento Automático de Habla	2-7
2.4.1.1 Definición	2-7
2.4.1.2 El Problema del Reconocimiento Automático del Habla	2-8
2.4.1.3 Clasificación del Problema de Reconocimiento	2-10
2.4.1.4 Técnicas más utilizadas aplicadas al Reconocimiento de Habla	2-10
2.4.1.5 Clasificación de los Sistemas de Reconocimiento de Habla según su Arquitectura	2-11
2.4.1.6 Incorporación de Conocimiento Lingüístico en los Sistemas de Reconocimiento de Habla	2-12
2.4.1.7 Interacción del Sistema de Reconocimiento y el Sistema de Comprensión	2-15
2.4.1.8 El Problema de la Eficiencia (Reducción del Espacio de Búsqueda)	2-16
2.4.1.9 Compilación de los Diccionarios. Su influencia en el Espacio de Búsqueda	2-16
2.4.1.10 Los Problemas de Robustez y de Cobertura del Modelo Gramatical ..	2-17
2.4.1.11 Generación de las N Mejores Hipótesis de Salida	2-17
2.4.1.12 El Módulo de Reconocimiento en esta Tesis	2-17
2.4.2 Interpretación Semántica	2-18
2.4.2.1 Estrategias de Interpretación Semántica	2-19
2.4.3 Interpretación Contextual	2-20
2.4.3.1 Estrategias de Interpretación Contextual	2-22

2.4.3.2 Problemas Característicos en la Representación del Conocimiento	2-23
2.4.4 Razonamiento de la Aplicación	2-23
2.4.5 Generación de una Respuesta Hablada	2-23
2.5 Descripción de la Solución Propuesta en esta Tesis	2-24
2.6 Habla Espontánea	2-25
2.7 Evaluación de los Sistemas	2-26
2.7.1 Evaluación Acústica de los Sistemas. Medidas de Calidad Acústica Objetivas	2-28
2.7.1.1 Algoritmo de Evaluación de la Calidad Acústica	2-30
2.7.1.2 Problemas de las Figuras de Mérito Descrietas	2-30
2.7.1.3 Una nueva medida: el ETP (Error Total Ponderado)	2-31
2.7.1.4 Alineamiento de Palabras vs. Alineamiento Fonológico	2-31
2.7.1.5 Alineamiento utilizando Marcas Temporales	2-32
2.7.1.6 Los Test de Significancia: Validación de los Sistemas	2-32
2.7.1.7 Bandas de Probabilidad	2-33
2.7.2 Evaluación de Sistemas de Comprensión de Habla	2-33
2.8 Evaluación del Sistema de Comprensión Desarrollado en la Tesis	2-37

CAPÍTULO 3. EL DOMINIO DE LA APLICACIÓN : SISTEMA DE INFORMACIÓN NAVAL CON ACCESO A BASES DE DATOS

3.1 Sistemas de Recuperación de Información en Lenguaje Natural	3-2
3.1.1 Definición	3-2
3.1.2 Clasificación	3-2
3.2 Descripción del Dominio Semántico Restringido de la Aplicación	3-2
3.3 Descripción de las Bases de Datos de Información Naval	3-3
3.4 Conceptos de la Aplicación (Ontología)	3-8
3.5 El Problema de la Ambigüedad Conceptual	3-13
3.6 Las Categorías Semánticas de la Aplicación	3-13
3.7 Bases de Datos de Habla y de Texto utilizadas en el Desarrollo	3-14
3.7.1 Base de Datos de Texto en Castellano	3-14
3.7.2 Base de Datos de Habla	3-15
3.7.3 Captura de Nuevos Datos Textuales	3-15

CAPÍTULO 4. DISEÑO E IMPLEMENTACIÓN DE UN SISTEMA DE COMPREENSIÓN DE HABLA

4.1 Descripción General de la Arquitectura Desarrollada	4-1
---	-----

CAPÍTULO 5. DECODIFICADOR ACÚSTICO BASADO EN HMM

5.1 Módulo Gramatical para N-gramas	5-4
5.1.1 Gramática Nula (GN)	5-8
5.1.2 Gramática Bigrama de 160 Macrocategorías (160MC)	5-9
5.1.3 Suavizado de la Gramática 160MC usando la técnica Back-Off	5-13
5.1.4 Gramática Semántica GRSEM-S	5-15
5.1.4.1 La Categoría y el Concepto BASURA: Robustez	5-16
5.1.4.2 Entrenamiento de la Gramática Semántica	5-17
5.1.4.3 Perplejidad y Cobertura de la Gramática GRSEM-S	5-18

5.2 Módulo Léxico para Diccionarios Lineales	5-19
5.2.1 Transiciones Fonológicas Entre Palabras en Habla Continua	5-20
5.2.2 Nombres Propios Compuestos como una Palabra	5-21
5.3 Módulo Acústico (Algoritmo de Un Paso)	5-22
5.3.1 Espacio de Búsqueda Estático (Algoritmo de Un Paso Básico)	5-24
5.3.1.1 El Problema del Silencio Acústico (Pausa entre Palabras)	5-27
5.3.2 Estrategia de Recorte de Caminos (Beam-Search): Eficiencia	5-27
5.3.2.1 Recorte basado en un Ancho de Haz Constante	5-28
5.3.2.2 Estudio del Espacio de Búsqueda: Distribución de las Distancias de los Estados del Espacio	5-31
5.3.2.3 Estimación de los Umbrales de Recorte: Método basado en un Histograma de Distancias de los Estados de los Caminos Óptimos de las Frases de Entrenamiento y en un Factor de Conservación	5-33
5.3.2.4 Influencia de la Gramática y del Modelado Acústico en el Umbral de Recorte	5-39
5.3.2.5 Umbrales basados en el <i>Parámetro</i> de la Trama Anterior	5-40
5.3.2.5 Algoritmo de Un Paso con Estrategia de Recorte	5-43
5.3.3 Generación de las N Mejores Soluciones (Pseudo N-Best)	5-44
5.3.3.1 Modificación del Algoritmo de Un Paso para la obtención de las N Mejores Hipótesis de Salida	5-46
5.3.3.2 Generando un Grafo con las N Mejores Soluciones	5-49
5.4 Evaluación del Módulo Acústico	5-50
5.4.1 Evaluación de la Calidad Acústica (Tasa de Error de Palabras) y Validación Estadística de los Resultados	5-50
5.4.1.1 Sin gramática	5-52
5.4.1.1.1 Modelos Discretos. 2 Codebooks	5-52
5.4.1.1.2 Modelos Semicontínuos con Pausado. 3 Codebooks	5-54
5.4.1.1.3 Comparando los Sistemas sin Gramática. Análisis de la Influencia de la Gramática Léxica o Diccionario	5-56
5.4.1.2 Con las gramáticas morfosintácticas GR160S-J y GR160S-S	5-58
5.4.1.3 Con la gramática semántica GRSEM-S	5-62
5.4.2 Evaluación de las N Mejores Salidas (Recuperación de Errores)	5-67
5.4.2.1 Modificación del Algoritmo de Evaluación	5-67
5.4.2.2 Resultados de la Evaluación de la Estrategia N-Caminos	5-67
5.4.2.2.1 Sin Gramática	5-68
5.4.3 Evaluación de la Eficiencia (Reducción del Espacio de Búsqueda)	5-71
5.4.3.1 Sin Gramática	5-72
5.4.3.2 Con la Gramática GR160S-J	5-75
5.4.4 Conclusiones Generales del Capítulo 5	5-81

CAPÍTULO 6. EL MÓDULO DE COMPRESIÓN

6.1 Decodificador Conceptual	6-1
6.1.1 Introducción	6-1
6.1.2 Autómatas Finitos Conceptuales. Categorías Semánticas	6-2
6.1.3 Generación del Autómata Conceptual del Dominio de Aplicación	6-4
6.1.4 El Problema de la Robustez. Categoría y Concepto Basura	6-5
6.1.5 El Categorizador Semántico	6-6
6.1.6 Ambigüedad Semántica. Grafo de Pares Palabra-Categoría	6-6

6.1.7 El Módulo de Segmentación Conceptual en el Sistema de Comprensión ..	6-7
6.1.8 Limitaciones del Segmentador Conceptual	6-8
6.2 Mapeador Conceptual	6-9
6.2.1 Introducción	6-9
6.2.2 Descripción General	6-10
6.2.3 Lenguaje de Reglas. Primitivas y Sintaxis	6-14
6.2.4 Reglas del Mapeador	6-15
6.2.5 El Módulo de Mapeado Conceptual en el Sistema de Comprensión	6-16
6.3 Analizador-Clasificador Estructural	6-17
6.3.1 Introducción. Justificación	6-17
6.3.2 Descripción General	6-19
6.3.3 Clasificación Conceptual (Taxonomía)	6-22
6.3.4 Clasificación Estructural de las Frases de Entrenamiento	6-23
6.3.5 Reglas Estructurales Libres de Contexto	6-25
6.3.6 El Módulo de Análisis Estructural en el Sistema de Comprensión	6-26
6.4 Transformador Estructural	6-27
6.4.1 Introducción	6-27
6.4.2 Descripción General	6-28
6.4.3 Lenguaje de Reglas. Primitivas y Sintaxis	6-31
6.4.4 Reglas de Transformación	6-33
6.4.5 El Módulo de Transformación Estructural en el Sistema de Comprensión.	6-33
6.5 Control	6-34
6.5.1 Introducción	6-34
6.5.2 Descripción General	6-36
6.6 Traductor a SQL	6-36
6.6.1 Introducción	6-36
6.6.2 Descripción General	6-37
6.6.3 Reglas de Traducción SQL	6-39
6.6.4 El Módulo Traductor a SQL en el Sistema de Comprensión	6-41
6.7 Procesador Funcional	6-43
6.7.1 Introducción	6-43
6.7.2 Descripción General	6-44
6.7.3 Funciones Monarias (Internas)	6-44
6.7.4 Funciones Binarias (Externas)	6-47
6.7.5 Método Alternativo	6-48
6.8 Conclusiones	6-50
6.8.1 Cobertura	6-50
6.8.2 Limitaciones	6-53
6.9 Evaluación del Sistema de Comprensión Desarrollado	6-55

CAPÍTULO 7. CONCLUSIONES Y LÍNEAS DE TRABAJO FUTURAS

7.1 Conclusiones	7-1
7.2 Líneas de Trabajo Futuras	7-5

APÉNDICES

Apéndice 5.1 Gramática 3-gram de Categorías	A-1
Apéndice 5.2 Tabla de Alófonos Independientes del Contexto	A-5

Apéndice 5.3 Lista de Nombres Propios Compuestos Compilados como una sola Palabra (Diccionarios D-C y D-J).....	A-7
Apéndice 5.4 Detalles del Diccionario Categorizado asociado a la Gramática GR160S-J	A-13
Apéndice 5.5 Categorías Semánticas utilizadas en la Gramática Semántica (Conceptual) GRSEM-S	A-15
Apéndice 5.6 Recorte de Caminos (Beam-Search): Umbrales y Resultados de Reconocimiento	A-21
Apéndice 6.1 Conceptos del Dominio de Aplicación del Sistema de Información Naval	A-33
Apéndice 6.2 Categorías Semánticas y Conceptos del Dominio de Aplicación (Sistema de Información Naval)	A-39
BIBLIOGRAFÍA	B-1