

J. M. Pardo

**DEPARTAMENTO DE INGENIERÍA ELECTRÓNICA**  
**ESCUELA TÉCNICA SUPERIOR DE INGENIEROS**  
**DE TELECOMUNICACIÓN**



**TESIS DOCTORAL**

**SISTEMAS DE RECONOCIMIENTO**  
**DE HABLA CONTINUA Y AISLADA:**  
**COMPARACIÓN Y OPTIMIZACIÓN DE**  
**LOS SISTEMAS DE MODELADO**  
**Y PARAMETRIZACIÓN**

**Ricardo de Córdoba Herralde**  
**Ingeniero de Telecomunicación**

**Director de la tesis**  
**Dr. Ingeniero José Manuel Pardo Muñoz**

**1995**

## Resumen

El título de esta tesis puede parecer ambicioso por la amplitud de temas planteados, y puede que realmente lo sea, pues su objetivo es tratar de dominar toda una serie de técnicas con las que mejorar el modelado y parametrización de dos sistemas de reconocimiento de habla, que tienen las siguientes características:

- 1) Habla continua, vocabulario de 1000 palabras, dependiente del locutor, en inglés.
- 2) Habla aislada, vocabulario de dígitos, independiente del locutor, en castellano.

Los puntos básicos que se han tratado son:

1) *Parametrización*: se ha trabajado en la manera de caracterizar el espectro de la señal de voz, en cuanto a la forma de obtener la energía en distintos filtros en frecuencia, como paso previo a la obtención de los parámetros. Se han introducido distintas maneras de realizar un filtrado de dichos valores para conseguir que sean robustos frente al ruido telefónico.

2) *Tipo de modelado*: se ha trabajado con tres tipos de modelado: discreto, continuo y semicontinuo. En este último se ha introducido una técnica de preselección con la que reducir el tiempo de cálculo sin perder tasa de reconocimiento. En modelado continuo se ha tratado el equilibrio "número de unidades/número de mezclas por unidad", y la manera de "guiar" los experimentos hasta alcanzar un óptimo.

3) *Unidad de modelado*: se han estudiado distintas posibilidades en nuestros dos sistemas: palabra, fonema, trifenema generalizado y agrupamiento a nivel de estado. Se han introducido nuevas técnicas de agrupamiento de distribuciones con las que se consigue una compartición de parámetros óptima. También se ha

## *SISTEMAS DE RECONOCIMIENTO DE HABLA*

tratado el tema del suavizamiento de las unidades contextuales (mediante su interpolación con unidades no contextuales). Así mismo, se ha introducido una técnica con la que duplicar el número de modelos con el objeto de mejorar el reconocimiento suponiendo dos tipos de locutores.

### *4) Detección de principio y fin: (habla aislada)*

Se ha introducido una técnica basada en una red neuronal con la que se consigue una mejora significativa tanto en la precisión de las marcas obtenidas como en la tasa de reconocimiento que se alcanza a continuación cuando se utiliza.

Así mismo, utilizando modelos de ruido inicial y final comunes a todos los modelos, y como algoritmo de reconocimiento un *one-pass*, también se ha conseguido superar el efecto de una detección defectuosa.

## Abstract

The title of the thesis may sound too ambitious, because it covers a wide range of topics, and may be it is, as there are many techniques involved in our purpose of improving the modelling and parametrization of two speech recognition systems, with the following characteristics:

- 1) Continuous speech, 1000 words vocabulary, speaker dependent, in English.
- 2) Isolated speech, digits vocabulary, speaker independent, in Spanish.

The basic points covered are:

- 1) *Parametrization*: we have worked in the way to deal with the speech spectrum, to obtain the energy using different filters, as the previous step to obtain the parameters. We have introduced different ways to filter these values of energy to obtain robustness in the telephonic environment.

- 2) *Modelling technique*: we have worked with three types of modelling: discrete, continuous and semicontinuous. With semicontinuous, we have introduced a preselection technique to reduce computing time with no decrease in recognition rate. With continuous modelling we have dealt with the balance between the number of units and the number of mixtures in each unit, and the way to "guide" the experiments to reach an optimum.

- 3) *Modelling unit*: we have used different possibilities in our systems: word, phoneme, generalized triphone and state clustering. We have introduced new ways to cluster the distributions of the models in order to obtain an optimum parameter sharing. We have also dealt with the smoothing of contextual units (interpolating them with context independent units). We have introduced a technique to double the number of models in order to improve the recognition assuming there are two kinds of speakers.

## *SISTEMAS DE RECONOCIMIENTO DE HABLA*

### *4) Begin-end detection: (isolated speech)*

We introduced a technique based on a neural network that obtains a significant improvement both in accuracy of the end-pointing and in the recognition rate obtained after it.

By the other hand, using begin and end noise models shared by all the digit models, and one-pass as the recognition algorithm, we have override the effect of the unaccurate end-pointing.

# Índice

Resumen. . . . .	1
Abstract . . . . .	3
CAPÍTULO 1. INTRODUCCIÓN . . . . .	5
1. Enfoque general . . . . .	7
2. Técnicas más usadas en reconocimiento. . . . .	8
3. Resumen de objetivos . . . . .	10
3.1 Tipo de parametrización . . . . .	11
3.2 Tipo de modelado. . . . .	11
3.3 Unidades de modelado. . . . .	12
3.4 Detección de principio y final de palabra . . . . .	13
4. Características de los experimentos . . . . .	13
4.1 Sistema 1 . . . . .	13
4.2 Sistema 2 . . . . .	14
CAPÍTULO 2. ENCUADRE CIENTÍFICO-TECNOLÓGICO . . . . .	17
CAPÍTULO 3. PARAMETRIZACIÓN. . . . .	25
1. Parametrización original . . . . .	27
2. Sistema de habla continua . . . . .	28
2.1 Filtros triangulares . . . . .	29
2.2 Filtros sinusoidales . . . . .	30
2.3 Resultados . . . . .	33
3. Sistema de habla aislada . . . . .	33
3.1 Utilización de filtros perceptuales. . . . .	33
3.2 Filtrado temporal similar al RASTA . . . . .	37
3.3 Técnica J-RASTA . . . . .	44

## ÍNDICE

4. Número de parámetros y su agrupación en vectores . . . . .	50
CAPÍTULO 4. CUANTIFICACIÓN VECTORIAL . . . . .	53
1. Generación del <i>codebook</i> . . . . .	55
2. Variación en el número de centroides . . . . .	57
CAPÍTULO 5. ESTUDIO DE LOS TIPOS DE MODELADO . . . . .	59
1. Discreto . . . . .	61
2. Semicontinuo . . . . .	61
2.1 Sistema de habla continua . . . . .	62
2.1.1 Adaptación de los algoritmos. . . . .	62
2.1.2 Decisión de la inicialización del entrenamiento. . . . .	64
2.1.3 Decisión del número de mezclas a utilizar . . . . .	65
2.1.4 Resultados globales obtenidos . . . . .	67
2.2 Sistema de dígitos . . . . .	67
2.2.1 Resultados obtenidos . . . . .	67
2.2.2 Codificación del algoritmo . . . . .	68
2.2.3 Técnica de preselección . . . . .	68
2.2.4 Limitación de los estados considerados . . . . .	71
3. Continuo (sistema 1) . . . . .	72
3.1 Introducción teórica. . . . .	72
3.2 Base de datos . . . . .	73
3.3 Descripción del modelado . . . . .	74
3.4 Incremento del número de mezclas. . . . .	75
3.5 Número de mezclas variable. . . . .	77
CAPÍTULO 6. ELECCIÓN DE LA UNIDAD DE MODELADO . . . . .	79
1. Introducción . . . . .	81
2. Sistema 1 (habla continua) . . . . .	82
2.1 Trifonema generalizado . . . . .	83
2.1.1 Algoritmo de agrupamiento. . . . .	83

## ÍNDICE

2.1.2 Resultados obtenidos . . . . .	85
2.1.3 Suavizado de las unidades obtenidas . . . . .	87
2.1.4 Palabras función . . . . .	91
2.2 Distribuciones o agrupamiento a nivel de estados . . . . .	93
2.2.1 Algoritmo de agrupamiento. . . . .	94
2.2.2 Mejoras sobre el sistema básico. . . . .	95
2.2.3 Intercambios durante el agrupamiento . . . . .	102
2.2.4 Resultados con los 12 locutores. . . . .	112
2.2.5 Suavizamiento de las unidades obtenidas. . . . .	113
2.3 Agrupamiento de parámetros con modelos continuos. . . . .	124
2.3.1 Compartición de parámetros . . . . .	125
2.3.2 Agrupamiento a nivel de estados. . . . .	125
3. Sistema 2 (habla aislada) . . . . .	129
3.1 Elección del número de estados. . . . .	129
3.2 Doble conjunto de modelos . . . . .	131
CAPÍTULO 7. DETECCIÓN DE PRINCIPIO Y FIN . . . . .	137
1. Introducción . . . . .	139
2. Red neuronal . . . . .	139
2.1 Configuración . . . . .	140
2.2 Entrenamiento. . . . .	140
2.3 Ajuste de la salida obtenida en la fase de reconocimiento . . . . .	141
2.4 Experimentos . . . . .	142
3. Modelos de ruido inicial y final. . . . .	144
CAPÍTULO 8. CONCLUSIONES Y LÍNEAS DE TRABAJO FUTURAS. . . . .	149
1. Conclusiones. . . . .	151
2. Líneas de trabajo futuras. . . . .	154
BIBLIOGRAFÍA. . . . .	157