

# EXHIBIT 1

March 10, 2011

**BY COURIER**

SRI International  
333 Ravenswood Avenue  
Menlo Park, CA 94025

Re: Personalized User Model LLP v. Google Inc., C.A. No. 09-00525-LPS

To Whom it May Concern:

On July 16, 2010, my client Personalized User Model, LLP brought a civil action against Google, Inc. for patent infringement in the United States District Court for the District of Delaware. You are being contacted because SRI International is likely to have documents and other information relevant to the case arising from its association and dealings with Google, Inc. Please see the attached subpoena and exhibits for instructions on how to respond.

Kind regards,

/s/ Jennifer D. Bennett

Jennifer D. Bennett

Enclosure

UNITED STATES DISTRICT COURT
for the
Northern District of California

Personalized User Model, LLP
Plaintiff
v.
Google, Inc.
Defendant
Civil Action No. 1:09-cv-525 (LPS)
District of Delaware

SUBPOENA TO TESTIFY AT A DEPOSITION
OR TO PRODUCE DOCUMENTS IN A CIVIL ACTION

To: SRI International
333 Ravenswood Avenue, Menlo Park, CA 94025

Testimony: YOU ARE COMMANDED to appear at the time, date, and place set forth below to testify at a deposition to be taken in this civil action. If you are an organization that is not a party in this case, you must designate one or more officers, directors, or managing agents, or designate other persons who consent to testify on your behalf about the following matters, or those set forth in an attachment:

Table with 2 columns: Place and Date and Time. Place: SNR Denton US LLP, 1530 Page Mill Road, Suite 200, Palo Alto, CA 94304. Date and Time: 03/21/2011 09:00

The deposition will be recorded by this method: Stenographic and video

Production: You, or your representatives, must also bring with you to the deposition the following documents, electronically stored information, or objects, and permit their inspection, copying, testing, or sampling of the material:

The provisions of Fed. R. Civ. P. 45(c), relating to your protection as a person subject to a subpoena, and Rule 45 (d) and (e), relating to your duty to respond to this subpoena and the potential consequences of not doing so, are attached.

Date: 03/10/2011

CLERK OF COURT

OR

/s/ Jennifer Bennett

Signature of Clerk or Deputy Clerk

Attorney's signature

The name, address, e-mail, and telephone number of the attorney representing (name of party)

Personalized User Model, LLP, who issues or requests this subpoena, are:

Jennifer Bennett
SNR Denton US LLP
1530 Page Mill Road, Suite 200, Palo Alto, CA 94304; T: 650.798.0300; Email: jennifer.bennett@snrdenton.com

Civil Action No. 1:09-cv-525 (LPS)

**PROOF OF SERVICE**

*(This section should not be filed with the court unless required by Fed. R. Civ. P. 45.)*

This subpoena for *(name of individual and title, if any)* \_\_\_\_\_  
was received by me on *(date)* \_\_\_\_\_.

I personally served the subpoena on the individual at *(place)* \_\_\_\_\_  
\_\_\_\_\_ on *(date)* \_\_\_\_\_; or

I left the subpoena at the individual's residence or usual place of abode with *(name)* \_\_\_\_\_  
\_\_\_\_\_, a person of suitable age and discretion who resides there,  
on *(date)* \_\_\_\_\_, and mailed a copy to the individual's last known address; or

I served the subpoena on *(name of individual)* \_\_\_\_\_, who is  
designated by law to accept service of process on behalf of *(name of organization)* \_\_\_\_\_  
\_\_\_\_\_ on *(date)* \_\_\_\_\_; or

I returned the subpoena unexecuted because \_\_\_\_\_; or

Other *(specify):* \_\_\_\_\_.

Unless the subpoena was issued on behalf of the United States, or one of its officers or agents, I have also  
tendered to the witness fees for one day's attendance, and the mileage allowed by law, in the amount of  
\$ \_\_\_\_\_.

My fees are \$ \_\_\_\_\_ for travel and \$ \_\_\_\_\_ for services, for a total of \$ \_\_\_\_\_ 0.00 \_\_\_\_\_.

I declare under penalty of perjury that this information is true.

Date: \_\_\_\_\_  
\_\_\_\_\_  
*Server's signature*

\_\_\_\_\_  
*Printed name and title*

\_\_\_\_\_  
*Server's address*

Additional information regarding attempted service, etc:

## Federal Rule of Civil Procedure 45 (c), (d), and (e) (Effective 12/1/07)

### (c) Protecting a Person Subject to a Subpoena.

**(1) Avoiding Undue Burden or Expense; Sanctions.** A party or attorney responsible for issuing and serving a subpoena must take reasonable steps to avoid imposing undue burden or expense on a person subject to the subpoena. The issuing court must enforce this duty and impose an appropriate sanction — which may include lost earnings and reasonable attorney’s fees — on a party or attorney who fails to comply.

#### **(2) Command to Produce Materials or Permit Inspection.**

**(A) Appearance Not Required.** A person commanded to produce documents, electronically stored information, or tangible things, or to permit the inspection of premises, need not appear in person at the place of production or inspection unless also commanded to appear for a deposition, hearing, or trial.

**(B) Objections.** A person commanded to produce documents or tangible things or to permit inspection may serve on the party or attorney designated in the subpoena a written objection to inspecting, copying, testing or sampling any or all of the materials or to inspecting the premises — or to producing electronically stored information in the form or forms requested. The objection must be served before the earlier of the time specified for compliance or 14 days after the subpoena is served. If an objection is made, the following rules apply:

**(i)** At any time, on notice to the commanded person, the serving party may move the issuing court for an order compelling production or inspection.

**(ii)** These acts may be required only as directed in the order, and the order must protect a person who is neither a party nor a party’s officer from significant expense resulting from compliance.

#### **(3) Quashing or Modifying a Subpoena.**

**(A) When Required.** On timely motion, the issuing court must quash or modify a subpoena that:

**(i)** fails to allow a reasonable time to comply;

**(ii)** requires a person who is neither a party nor a party’s officer to travel more than 100 miles from where that person resides, is employed, or regularly transacts business in person — except that, subject to Rule 45(c)(3)(B)(iii), the person may be commanded to attend a trial by traveling from any such place within the state where the trial is held;

**(iii)** requires disclosure of privileged or other protected matter, if no exception or waiver applies; or

**(iv)** subjects a person to undue burden.

**(B) When Permitted.** To protect a person subject to or affected by a subpoena, the issuing court may, on motion, quash or modify the subpoena if it requires:

**(i)** disclosing a trade secret or other confidential research, development, or commercial information;

**(ii)** disclosing an unretained expert’s opinion or information that does not describe specific occurrences in dispute and results from the expert’s study that was not requested by a party; or

**(iii)** a person who is neither a party nor a party’s officer to incur substantial expense to travel more than 100 miles to attend trial.

**(C) Specifying Conditions as an Alternative.** In the circumstances described in Rule 45(c)(3)(B), the court may, instead of quashing or modifying a subpoena, order appearance or production under specified conditions if the serving party:

**(i)** shows a substantial need for the testimony or material that cannot be otherwise met without undue hardship; and

**(ii)** ensures that the subpoenaed person will be reasonably compensated.

### (d) Duties in Responding to a Subpoena.

**(1) Producing Documents or Electronically Stored Information.** These procedures apply to producing documents or electronically stored information:

**(A) Documents.** A person responding to a subpoena to produce documents must produce them as they are kept in the ordinary course of business or must organize and label them to correspond to the categories in the demand.

**(B) Form for Producing Electronically Stored Information Not Specified.** If a subpoena does not specify a form for producing electronically stored information, the person responding must produce it in a form or forms in which it is ordinarily maintained or in a reasonably usable form or forms.

**(C) Electronically Stored Information Produced in Only One Form.** The person responding need not produce the same electronically stored information in more than one form.

**(D) Inaccessible Electronically Stored Information.** The person responding need not provide discovery of electronically stored information from sources that the person identifies as not reasonably accessible because of undue burden or cost. On motion to compel discovery or for a protective order, the person responding must show that the information is not reasonably accessible because of undue burden or cost. If that showing is made, the court may nonetheless order discovery from such sources if the requesting party shows good cause, considering the limitations of Rule 26(b)(2)(C). The court may specify conditions for the discovery.

#### **(2) Claiming Privilege or Protection.**

**(A) Information Withheld.** A person withholding subpoenaed information under a claim that it is privileged or subject to protection as trial-preparation material must:

**(i)** expressly make the claim; and

**(ii)** describe the nature of the withheld documents, communications, or tangible things in a manner that, without revealing information itself privileged or protected, will enable the parties to assess the claim.

**(B) Information Produced.** If information produced in response to a subpoena is subject to a claim of privilege or of protection as trial-preparation material, the person making the claim may notify any party that received the information of the claim and the basis for it. After being notified, a party must promptly return, sequester, or destroy the specified information and any copies it has; must not use or disclose the information until the claim is resolved; must take reasonable steps to retrieve the information if the party disclosed it before being notified; and may promptly present the information to the court under seal for a determination of the claim. The person who produced the information must preserve the information until the claim is resolved.

**(e) Contempt.** The issuing court may hold in contempt a person who, having been served, fails without adequate excuse to obey the subpoena. A nonparty’s failure to obey must be excused if the subpoena purports to require the nonparty to attend or produce at a place outside the limits of Rule 45(c)(3)(A)(ii).

**IN THE UNITED STATES DISTRICT COURT  
FOR THE DISTRICT OF DELAWARE**

**PERSONALIZED USER MODEL, )  
L.L.P., )**

Plaintiff, )

v. )

**GOOGLE, INC., )**

Defendant. )

C.A. No. 09-525 (LPS)

**NOTICE OF RULE 30(b)(6) DEPOSITION OF SRI INTERNATIONAL**

PLEASE TAKE NOTICE that, pursuant to Rules 26 and 30 of the Federal Rules of Civil Procedure, Plaintiff Personalized User Model, L.L.P. (“P.U.M.”) will take the deposition of Third Party SRI International (“SRI”) concerning the topics identified in Exhibit A, beginning at 9:00 a.m. on March 21, 2011, or at an otherwise mutually agreeable date, and will be held at the offices of SNR Denton US LLP, 1530 Page Mill Road, CA 94304, or at an otherwise mutually agreeable location. If the deposition is not completed on the date set out above, the taking of the deposition will continue day to day thereafter or pursuant to the parties’ agreement. The deposition will be recorded by stenographic, videographic, and/or audiographic means.

Pursuant to Rule 30(b)(6) of the Federal Rules of Civil Procedure, SRI is directed to designate one or more officers, directors, or managing agents, or other persons who will testify on its behalf, who are most knowledgeable regarding the matters identified in the attached Exhibit A. SRI is requested to provide a written designation of the names and positions of the officers, directors, managing agents, or other persons designated to testify

concerning the matters identified in the attached Exhibit and, for each person, identify the matters on which he or she will testify.

P.U.M. reserves the right to serve additional 30(b)(6) notices.

<p>Dated: March 10, 2011</p>	<p>By: <u>/s/ Jennifer D. Bennett</u> Jennifer D. Bennett (California State Bar No. 235196) SNR Denton US LLP 1530 Page Mill Road, Suite 200 Palo Alto, CA 94304 Telephone: (650) 798-0300 Facsimile: (650) 798-0310 E-Mail: jennifer.bennett@snrdenton.com</p> <p>Marc S. Friedman SNR Denton US LLP 1221 Avenue of the Americas New York, NY 10020-1089 Telephone: (212) 768-6700 Facsimile: (212) 768.6800 E-Mail: marc.friedman@snrdenton.com</p> <p>Attorneys for Plaintiff PERSONALIZED USER MODEL, L.L.P.</p>
------------------------------	--

CERTIFICATE OF SERVICE

I hereby certify that on March 10, 2011, copies of the foregoing were caused to be served by e-mail upon the following:

Richard L. Horwitz  
David E. Moore  
POTTER ANDERSON & CORROON LLP  
1313 N. Market St., 6<sup>th</sup> Floor  
Wilmington, DE 19801  
**rhowitz@potternanderson.com**  
**dmoore@potteranderson.com**

Brian C. Cannon  
QUINN EMANUEL URQUHART OLIVER  
& HEDGES, LLP  
**briancannon@quinnemanuel.com**

Charles K. Verhoeven  
QUINN EMANUEL URQUHART OLIVER  
& HEDGES, LLP  
**charlesverhoeven@quinnemanuel.com**

David A. Perlson  
QUINN EMANUEL URQUHART OLIVER  
& HEDGES, LLP  
**davidperlson@quinnemanuel.com**

Antonio R. Sistos  
QUINN EMANUEL URQUHART OLIVER  
& HEDGES, LLP  
**antoniosistos@quinnemanuel.com**

Eugene Novikov  
QUINN EMANUEL URQUHART OLIVER  
& HEDGES, LLP  
**eugenenovikov@quinnemanuel.com**

/s/ Jennifer D. Bennett  
Jennifer D. Bennett (Cal. Bar. No. 235196)  
SNR Denton US LLP  
1530 Page Mill Road, Suite 200  
Palo Alto, CA 94304-1125  
(650) 798-0300



## EXHIBIT A

### **I DEFINITIONS**

1. “SRI,” “YOU,” and “YOUR,” means SRI International, and its officers, directors, current and former employees, counsel, agents, consultants, representatives, and any other persons acting on behalf of any of the foregoing, and SRI International’s affiliates, parents, divisions, joint ventures, licensees, franchisees, assigns, predecessors and successors in interest, and any other legal entities, whether foreign or domestic, that are owned or controlled by SRI International, and all predecessors and successors in interest to such entities.

2. “Google” means Google, Inc. and its officers, directors, current and former employees, counsel, agents, consultants, representatives, attorneys, and any other persons acting on behalf of any of the foregoing, and Google’s affiliates, parents, divisions, joint ventures, licensees, franchisees, assigns, predecessors and successors in interest, and any other legal entities, whether foreign or domestic, that are owned or controlled by Google, and all predecessors and successors in interest to such entities.

3. “Lawsuit” means the case styled *Personalized User Model LLP v. Google, Inc.*, 1:09-cv-525, in the United States District Court for the District of Delaware.

4. “‘040 PATENT” means U.S. Patent No. 6,981,040, entitled “Automatic, Personalized Online Information and Product Services,” all underlying patent applications, all continuations, continuations-in-part, divisionals, reissues, and any other patent applications in the ‘040 patent family

5. “‘031 PATENT” means U.S. Patent No. 7,320,031, entitled “Automatic, Personalized Online Information and Product Services,” all underlying patent applications, all continuations, continuations-in-part, divisionals, reissues, and any other patent applications in the

‘031 patent family.

6. “‘276 PATENT” means U.S. Patent No. 7,685,276, entitled “Automatic, Personalized Online Information and Product Services,” all underlying patent applications, all continuations, continuations-in-part, divisionals, reissues, and any other patent applications in the ‘031 patent family.

7. “PATENTS-IN-SUIT” shall refer to the ‘040 PATENT, the ‘031 PATENT, and the ‘276 PATENT individually and collectively.

8. “DOCUMENT” shall mean all materials and information that are discoverable pursuant to Rule 34 of the Federal Rules of Civil Procedure. A draft or non-identical copy is a separate document within the meaning of this term.

9. “PUM” and “PLAINTIFF” shall mean Personalized User Model LLP., Plaintiff in the civil case captioned Personalized User Model, LLP v. Google Inc., Case No. 09-525 (JJF).

10. The term “PERSON” shall refer to any individual, corporation, proprietorship, association, joint venture, company, partnership or other business or legal entity, including governmental bodies and agencies.

11. “REFLECT,” “REFLECTING,” “RELATE TO,” “REFER TO,” “RELATING TO,” and “REFERRING TO” shall mean relating to, referring to, concerning, mentioning, reflecting, pertaining to, evidencing, involving, describing, discussing, commenting on, embodying, responding to, supporting, contradicting, or constituting (in whole or in part), as the context makes appropriate.

12. “Include” and “including” shall mean including without limitation.

13. Use of the singular also includes the plural and vice-versa.

14. The words “or” and “and” shall be read in the conjunctive and in the disjunctive

wherever they appear, and neither of these words shall be interpreted to limit the scope of these Interrogatories.

15. The use of a verb in any tense shall be construed as the use of the verb in all other tenses.

### **DEPOSITION TOPICS**

1. All facts and circumstances, including but not limited to all communications whether written, oral or otherwise, between Google and SRI, concerning all transactions, contracts, agreements and understandings, and payments between Google and SRI concerning the patents-in-suit or any invention(s) claimed therein, and/or Yochai Konig.
2. The work performed by Yochai Konig while at SRI.
3. Any and all documents or other evidence that Dr. Konig developed the inventions claimed in the patents-in-suit using SRI's equipment, supplies, facility, or trade secret information, or during the time of day when he was supposed to be working for SRI.
4. All documents provided by SRI to Google regarding Yochai Konig or work performed by him for SRI.
5. All invoices submitted by SRI to Google for work responding to discovery in connection with this lawsuit.
6. SRI's knowledge of Yochai Konig and/or Utopy's work after Dr. Konig left the employment of SRI.
7. Activities of the SRI Speech Technology and Research (STAR) Laboratory from 1996 through 1999.
8. All business relationships or contracts between SRI and Google, or subsidiary or affiliate of Google, including, but not limited to (a) all work performed by SRI for Google, or subsidiary or affiliate of Google, in the last 10 years; (b) all work performed by Google, or subsidiary or affiliate of Google, for SRI in the last 10 years, and (c) all sums of money received by SRI from

Google, or any subsidiary or affiliate of Google, or any officers or directors of these entities in the last 10 years.

9. All documents produced by SRI to PUM under the previously served subpoena, including, but not limited, to the authenticity of such documents and the manner in which they were created and kept.

10. All information received from third parties relating to any of the above subjects.

# EXHIBIT 2

**IN THE UNITED STATES DISTRICT COURT  
FOR THE DISTRICT OF DELAWARE**

PERSONALIZED USER MODEL, L.L.P.,	)	
	)	
Plaintiff,	)	C.A. No. 09-525-LPS
	)	
v.	)	
	)	
GOOGLE, INC.,	)	
	)	
Defendant.	)	
	)	

---

**DECLARATION OF ROY TWERSKY IN SUPPORT OF PLAINTIFF PERSONALIZED  
USER MODEL, L.L.P.’S BRIEF IN OPPOSITION TO DEFENDANT GOOGLE, INC.’S  
MOTION FOR LEAVE TO FILE ITS MOTION FOR SUMMARY JUDGMENT**

I, Roy Twersky, declare:

1. I make this declaration in support of P.U.M.’s Brief in Opposition to Defendant Google, Inc.’s Motion for Leave to File its Motion for Summary Judgment.
2. I received my undergraduate degree in economics and mathematics from Tel Aviv University in 1986. Five years later, I received my MBA from the Wharton School of Business. I also completed the All-But-Dissertation requirements for the PhD program at Stanford University's Graduate School of Business. While at both Tel Aviv University and Stanford, I took courses in computer science and mathematics.
3. Around early 1999, Yochai Konig and I started thinking generally about the problem of information overload on the Internet and whether there might be solutions to that problem. At that time, Mr. Konig and I also discussed starting a company to develop technology to address the problem. In order to form the company, I knew we would first need to generate capital. During the Spring and Summer of 1999 I reached out to potential investors to raise capital to fund the company. I approached investors with the general idea of personalized information services that we hoped to develop after receiving funding. The presentation that I

gave to investors in July 1999, attached hereto as Exhibit A, represents the general idea of the technology that might solve the problem.

REDACTED

REDACTED

REDACTED

REDACTED

I declare under penalty of perjury under the laws of the United States of America that the foregoing is true and correct.

Dated: March 10, 2011

/s/ Roy Twersky  
Roy Twersky



EXHIBIT A  
FULLY REDACTED

**EXHIBIT B**  
**FULLY REDACTED**

# EXHIBIT 3



speech recognition, statistical machine learning, stochastic processes, optimization theory and neural computation. Thus, by the beginning of 1996, I had been working with machine learning concepts for more than eight years.

4. On April 8, 1996, I began working as a research scientist at SRI in the Speech Technology and Research Laboratory ("Star Lab"). While employed by SRI, I researched over-the-telephone large vocabulary speech recognition. I later became the principal investigator of a research project funded by the Department of Defense, researching Nonlinear Discriminant Feature Extraction for Robust Text-Independent Speaker Recognition. My research at SRI was in no way related to the Internet, search, or the development of personalized on-line information services. Rather, the focus of my research was creating speaker recognition systems that are robust to telephone handset distortion by discriminative feature design. *See* published research, attached as Exhibits D-F.

5. Toward the conclusion of my employment at SRI, Mr. Twersky and I began to think generally about the problem of information overload on the Internet and concluded that perhaps we should team with each other and we started to discuss an approach to achieving a solution to this problem. Prior to my departure from SRI, we had not yet conceived of the inventions claimed in the patents-in-suit. For this reason, I did not think it was necessary to disclose our very general ideas to SRI pursuant to my Employment Agreement.

6. Also, while I was employed by SRI I used no equipment, supplies, facilities or trade secret information belonging to SRI to work on the ideas that eventually became the inventions disclosed in the patents-in-suit. And, any time that I did spend working on potential solutions to the information overload problem was entirely my personal time (*i.e.*, nights and weekends).

7. In sum, the inventions claimed in the patents-in-suit did not result from my work for SRI.

8. I ended my employment at SRI on August 5, 1999 and shortly thereafter joined Roy Twersky as a co-founder of Utopy, Inc. After terminating my employment at SRI, Mr.

Twersky and I had the time and resources to dedicate to conceiving the inventions claimed in the patents-in-suit.

9. In conceiving the inventions claimed in the patents-in-suit, I did not rely on my work at SRI's Star Lab involving discriminative feature design in speaker recognition systems to make them robust to telephone handset distortion. Instead, I drew upon my education and training that I had acquired well before joining SRI, including education and training in machine learning, speech recognition, stochastic modeling, neural networks and statistical pattern recognition. I applied my expertise in these areas to conceive of the inventions in the patents-in-suit.

REDACTED

I declare under penalty of perjury under the laws of the United States of America that the foregoing is true and correct.

Dated: March 10, 2011

/s/ Yochai Konig  
Yochai Konig

EXHIBIT A  
FULLY REDACTED

**EXHIBIT B**  
**FULLY REDACTED**



# EXHIBIT C

**REMAP: Recursive Estimation and Maximization of A  
Posteriori Probabilities in Transition-based Speech  
Recognition**

by

Yochai Konig

B.S. (Technion, Israel Institute of Technology) 1990

A dissertation submitted in partial satisfaction of the  
requirements for the degree of  
Doctor of Philosophy

in

Computer Science

in the

GRADUATE DIVISION

of the

UNIVERSITY of CALIFORNIA at BERKELEY

Committee in charge:

Professor Nelson Morgan, Chair  
Professor Jerome Feldman  
Professor Charles Stone

1996

**REMAP: Recursive Estimation and Maximization of A  
Posteriori Probabilities in Transition-based Speech  
Recognition**

Copyright 1996

by  
Yochai Konig

## Abstract

REMAP: Recursive Estimation and Maximization of A Posteriori Probabilities in  
Transition-based Speech Recognition

by

Yochai Konig

Doctor of Philosophy in Computer Science

University of California at Berkeley

Professor Nelson Morgan, Chair

In this thesis we present a framework for training and modeling continuous speech recognition systems based on the theoretically optimal Maximum a Posteriori (MAP) criterion. In contrast, most state-of-the-art systems are trained according to the Maximum Likelihood (ML) criterion. Specifically, we introduce a discriminant training algorithm (REMAP) for statistical sequence classification which, for any given sentence, monotonically increases the posterior probability of the correct sentence while reducing the probabilities of all rival models.

Based on the studies described here, which show that explicitly modeling transitions between speech units can improve recognition performance, REMAP is developed in the context of a transition-based model (although it is also applicable to non-transition-based models). Furthermore, the model uses local transition probabilities (i.e., the posterior probability of the current state given the current acoustic vector and the previous state) to estimate global posteriors of sentences. Thus, it is a true recognition model, i.e., it directly maps from acoustic sequences to sentences, unlike Hidden Markov Models (HMMs) which model the inverse relation (the likelihood of producing an acoustic sequence given an assumed state sequence).

Experimental results support the proposed framework. In comparison to a baseline system, the results show an increase in the estimates of posterior probabilities of the correct sentences after training, and a significant decrease in error rate.

To my parents, Pnina and Aron König

# Contents

<b>List of Figures</b>	<b>vi</b>
<b>1 Introduction and Goals</b>	<b>1</b>
1.1 Thesis Overview . . . . .	1
1.2 Thesis Goals . . . . .	2
1.2.1 A Recognition Model . . . . .	2
1.2.2 Transition-based Modeling . . . . .	3
1.2.3 A Discriminant Training Algorithm . . . . .	4
1.3 Thesis Structure . . . . .	5
<b>2 Problem Formulation and Existing Models</b>	<b>7</b>
2.1 Introduction . . . . .	7
2.1.1 Bayes Decision Theory . . . . .	7
2.1.2 Problem Formulation . . . . .	8
2.1.3 Definitions and Notation . . . . .	10
2.2 Hidden Markov Models (HMMs) . . . . .	12
2.2.1 Introduction . . . . .	12
2.2.2 Assumptions . . . . .	13
2.2.3 Definition of the Hidden Markov Model . . . . .	14
2.3 Hybrid Systems . . . . .	15
2.3.1 Multilayer Perceptrons (MLPs) . . . . .	15
2.3.2 Motivations . . . . .	16
2.3.3 MLPs as Statistical Estimators . . . . .	17
2.4 Non-stationary Modeling . . . . .	21
2.4.1 Segment-Based Approaches . . . . .	22
<b>3 Training Algorithms and Optimization Criteria</b>	<b>25</b>
3.1 Likelihood Estimation and Training . . . . .	25
3.1.1 Introduction . . . . .	25
3.1.2 The Relation between ML and MAP . . . . .	26
3.1.3 Implementation - The EM Algorithm . . . . .	27
3.1.4 Summary . . . . .	28
3.2 Discriminant Approaches . . . . .	29

3.2.1	Introduction . . . . .	29
3.2.2	Maximum Mutual Information (MMI) . . . . .	29
3.2.3	Generalized Probabilistic Descent (GPD) . . . . .	31
<b>4</b>	<b>Transition-based Modeling</b>	<b>34</b>
4.1	Introduction . . . . .	34
4.2	Perceptual and Physiological Point of View . . . . .	35
4.3	The Time-Index Model . . . . .	36
4.3.1	Introduction . . . . .	36
4.3.2	Deng’s Trended HMM . . . . .	36
4.3.3	An Introduction to the Time-Index Model . . . . .	37
4.3.4	An Example . . . . .	38
4.3.5	An Implementation of the Time-Index Model . . . . .	39
4.3.6	Experiments . . . . .	39
<b>5</b>	<b>Discriminant HMM (DHMM)</b>	<b>42</b>
5.1	Introduction . . . . .	42
5.2	Estimation of the Posterior Probability of Word Sequences . . . . .	43
5.2.1	Recognition . . . . .	45
5.3	MAP Constraints . . . . .	46
5.4	Early Experiments and Error Analysis . . . . .	47
5.4.1	Early Experiments . . . . .	47
5.4.2	Error Analysis . . . . .	49
<b>6</b>	<b>REMAP Training of HMM/MLP Hybrids</b>	<b>51</b>
6.1	Introduction . . . . .	51
6.1.1	Motivations . . . . .	51
6.1.2	Problem Formulation . . . . .	52
6.2	Solution - REMAP . . . . .	53
6.2.1	Introduction . . . . .	53
6.2.2	Definitions . . . . .	54
6.2.3	Theorem 1 . . . . .	54
6.2.4	Theorem 2 . . . . .	54
6.2.5	Theorem 3 . . . . .	55
6.2.6	Summary and Discussion . . . . .	58
6.3	REMAP Training . . . . .	59
6.3.1	Introduction . . . . .	59
6.3.2	Target Estimation . . . . .	59
6.3.3	Forward Recursion . . . . .	61
6.3.4	Estimating the Previous State Distribution . . . . .	62
6.3.5	REMAP Training Algorithm . . . . .	62
6.3.6	Complexity Issues . . . . .	63
6.4	The Role of the Language Model . . . . .	64

<b>7</b>	<b>Experimental Results with REMAP</b>	<b>67</b>
7.1	Isolated Speech Experiments . . . . .	67
7.2	Continuous Speech Experiments . . . . .	68
7.3	Analysis and Discussion . . . . .	70
7.3.1	Invalid State Sequences . . . . .	70
7.3.2	Hard Targets . . . . .	72
<b>8</b>	<b>Conclusions and Future Work</b>	<b>73</b>
8.1	Conclusions . . . . .	73
8.2	Future Work . . . . .	74
8.2.1	Incorporating Language Information . . . . .	74
8.2.2	Extensions . . . . .	76
8.3	Epilog . . . . .	77
<b>A</b>	<b>REMAP Convergence - Theorem Proof</b>	<b>79</b>
A.1	Theorem 1 . . . . .	79
A.2	Theorem 2 . . . . .	80
A.3	Theorem 3 . . . . .	82



## List of Figures

1.1	Goal - A Recognition Model . . . . .	3
1.2	Goal - A Discriminant Training Algorithm . . . . .	5
2.1	HMM - An Example . . . . .	14
2.2	MLP - An Example . . . . .	18
4.1	Depiction of the time index along with the current state . . . . .	37
4.2	The topology of the time-index model . . . . .	38
4.3	The time index net . . . . .	40
5.1	An example of a Discriminant HMM for the word “cat”. The variable $x$ refers to a specific acoustic observation $x_n$ at time $n$ . . . . .	43
5.2	An example MLP that estimates local conditional probabilities. . . . .	44
6.1	An illustration of REMAP . . . . .	63
7.1	The probability of a transition (changing state) for every frame in the utterance “one.” The Y-axis represents the probability of transition, and the X-axis the time step . . . . .	69

## Acknowledgements

First and foremost I would like to thank my advisor, Nelson Morgan for sharing with me his vision and enthusiasm for research and science. Morgan introduced me to the field of speech recognition, and has guided and collaborated with me along the long way that has culminated in this thesis. He set very high standards for everything starting with scientific work and ending with funny stories. I feel fortunate to be part of a research group led by Morgan, as he inspired spirit of cooperation, bonding, and fun, while doing ground-breaking research.

I want to thank my “second” advisor Hervé Bourlard for sharing his wisdom, knowledge, and Belgian sense of humor with me. Hervé is a role model for me in terms of his analytic and systematic approach to scientific research. My images of research, are discussing and arguing over equations written in green and blue on a whiteboard with Hervé and Morgan. Most of the research described in this thesis is a result of collaborating with Morgan and Hervé.

I would like to thank Hynek Hermansky for sharing his deep understanding of speech science with me, for stimulating discussions, and for his friendship. My thanks to Steve Greenberg for his patience, sound advice, and for making available to me his vast knowledge and wisdom. I am grateful to Jerry Feldman, for his insight in recommending that I work on speech recognition with Morgan, for serving as a member of my qualification and thesis committees, and for providing such a wonderful research environment at ICSI. My gratitude goes to John Wawrzynek and John Ohala for serving as members of my qualification committee, and to Charles Stone for being a member of my thesis committee.

I wish to thank my past and present peers at the realization group, here at ICSI. I feel really lucky to know, work, collaborate, and socialize with such a great group of people. In particular I would like to thank Chuck Wooters, Su-Lin Wu, Krstè Asanovic, Jeff Bilmes, Eric Fosler, Dan Jurafsky, David Johnson, Nikki Mirghafori, Brian Kingsbury, Mike Shire, Warner Warren, Andreas Stolcke, Gary Tajchman, Steve Renals, Chris Bregler, and Kristine Ma.

My life in Berkeley would not be the same without my friends. They made Berkeley not just a great school, but also a fun and exciting place to live. To mention a few, I would like to thank Amir Guttman, Kim Zetter, Shlomo Zilberstein, Gil Sudai, Eitan Reuveny, Adriana Helmer, Tina Perry, Chih-Po Wen, Roy Twersky, Michele Gill, Yizhaq Makovsky,

and Benjamin Fuchs.

Last but not least, I want to thank my family, my sister Mickey, and my brother Dudy, for their encouragement, support, and visits. This thesis is dedicated to my parents, Pnina and Aron. Their values, education, warmth, support, and love, made it all possible.

# Chapter 1

## Introduction and Goals

### 1.1 Thesis Overview

Many pattern recognition problems that are of crucial importance today are inherently sequential in nature. Some examples include recognizing an utterance given a sequence of samples from a speech signal, or deciphering a hand written sentence given a digitized pen trajectory. Theoretically the optimal way to classify an input sequence is to choose the class with the highest posterior probability given this sequence (Duda & Hart 1973). Therefore, at training one wants to maximize the posterior probability of the correct model (sentence) given the evidence (sequence of acoustic vectors). An optimization criterion for parameter estimation that achieves this goal during training is the Maximum A Posteriori (MAP) criterion. Most state-of-the-art systems, however, are trained according to other criteria such as Maximum Likelihood (ML), which achieve this goal only under strong assumptions as discussed in Chapter 3. In this thesis we present a framework for training and modeling continuous speech recognition systems based on the MAP criterion. Specifically, we introduce a discriminant training algorithm for statistical sequence classification that monotonically increases the posterior probability of the correct sentence while reducing the probabilities of all rival models (sentences). Roughly speaking, instead of modeling the distribution of each class observation, the boundaries between classes are modeled. Thus, the correctness of the model is not assumed and during training one can optimize the overall goal of minimizing recognition errors.

REMAP can be used in a new form of hybrid Hidden Markov Models (HMM)/ Artificial Neural Network (ANN) which, in addition to the advantages of standard HMM/ANN

hybrids, uses “full” posterior probabilities for training and recognition. Furthermore, in the new HMM/ANN hybrid, the ANN targets and weights are iteratively re-estimated, a process that guarantees an increase of the posterior probability of the correct model, hence reducing the error rate.

Based on the studies described here, which show that explicitly modeling transitions between speech units can improve recognition performance, a training algorithm is developed in the context of a transition-based model (the algorithm is general, and also applicable to other models). Our interest in transition-based models was motivated by perceptual and physiological evidence, e.g., (Lindblom & Studdert-Kennedy 1967; Furui 1986b; Kiang 1984; Ruggero 1994; Smith & Zwislocki 1975), which show that spectral transition information is crucial for human perception. Furthermore, the Discriminant HMM (DHMM) model, uses local transition probabilities (the posterior probability of the current state given the current acoustic vector and the previous state) to estimate global sentence posteriors. It is a true recognition model. It directly maps from acoustic sequences to sentences, unlike HMMs that model the inverse relation (the likelihood of producing an acoustic sequence given an assumed state sequence) (Levinson *et al.* 1983b; Jelinek 1976; Jelinek & Mercer 1980; Baker 1975).

Experimental results support the proposed framework. In comparison to a baseline system, the results show an increase in the estimates of posterior probabilities of the correct sentences after training, and a significant decrease in error rate. Thus, a posterior based approach may be a viable alternative to current paradigms.

In the following sections we expand on the goals and the motivations for this study.

## 1.2 Thesis Goals

### 1.2.1 A Recognition Model

In automatic speech recognition the following two steps are usually performed:

1. The first step is to transform the speech signal into a finite sequence of numbers hopefully without loss of relevant information for recognition. This step, usually called *feature extraction*, is constrained on the one hand by the need to reduce the dimensionality of the data and on the other hand by the need to have good information for recognizing the utterance.

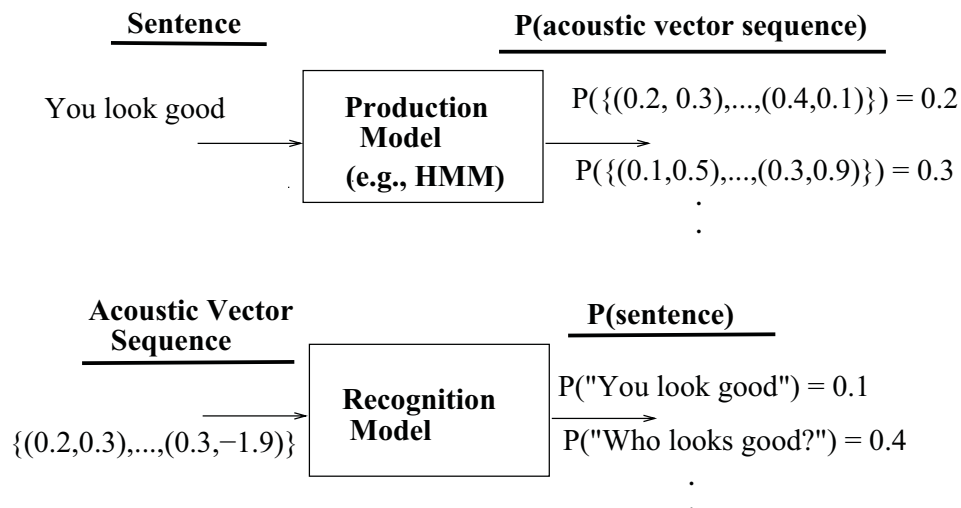


Figure 1.1: Goal - A Recognition Model

2. The second step is to construct a model that provides a mapping from the acoustic vector sequence into sentences, i.e., a *recognition model*.

In current state-of-the-art systems however, inverse modeling is usually used. Specifically, given a sentence, these systems model the likelihood of producing an acoustic sequence, a *production model*. This involves modeling all the possible ways that a given sentence can be spoken, taking into account differences between speakers, different speaking rates, styles, pronunciations, and variability due to environmental and channel noise.

An illustration of the difference between a production model and a recognition model is given in Figure 1.1. The illustration specifies the domain and range for production and recognition models. When recognizing an utterance, the input is a sequence of acoustic vectors and the output should be the recognized sentence, i.e., the domain and range are the same as in the recognition model. Hence, directly optimizing the recognition model parameters (during training) to increase the posterior probabilities of the correct models is more straightforward, compared with the production model with inverse range and domain.

### 1.2.2 Transition-based Modeling

One basic goal of this thesis is to explicitly model transitions between speech units<sup>1</sup>. Specifically, the goal is to estimate “soft boundaries” between speech units, i.e., to estimate

<sup>1</sup>The choice of the speech unit set determines the type of transitions that we model

the posterior probability of the current state given the current acoustic vector and the previous state, where *states* represent basic speech units such as phones. Roughly speaking, a “soft boundary” represents a probability of having a transition at a particular time point, i.e., values between zero and one. This is in contrast to a “hard boundary” that marks each time point as a transition, or non-transition point, i.e., values of ones and zeros only. This goal is motivated by our experimental results as described in Chapter 4, other related work (Goldenthal 1994; Deng 1992), and perceptual and physiological evidence (Lindblom & Studdert-Kennedy 1967; Furui 1986b; Kiang 1984; Ruggero 1994; Sachs *et al.* 1988; Smith & Zwislocki 1975; Seneff 1988) that show that spectral transition information is crucial for human perception.

The underlying reason for the “soft boundaries” is that human speech production is a continuous process, and modeling it as a sequence of discrete states is at best an approximation. A popular model of human speech production is that an utterance is organized as a succession of vocal-tract states, where each of these states represents a different configuration of the articulators. Further, at the level of articulatory performance, the gestures invoked to actualize these states are relatively slow. They merge spatially and temporally into a continuous process that usually only approximates the intended states and is highly sensitive to co-articulation (Deng & Sun 1994). Thus, there is not a single time point of transition between two states, but rather a “window” of transition which can be modeled with “soft boundaries.” Similarly, in human speech perception, the point at which a new phonetic element could be perceived is distributed in time.

### 1.2.3 A Discriminant Training Algorithm

Statistical pattern recognition theory tells us that the optimal recognition procedure (in terms of minimum classification error) is the one that chooses the sentence with the highest posterior probability given the acoustic vectors and all other available sources of knowledge (Duda & Hart 1973). Consequently, a training algorithm should estimate the parameters of a classifier to increase the posterior probability of the correct sentence (known during the training phase) while reducing the posteriors of all rival models (sentences). An illustration of a discriminant algorithm is given in Figure 1.2.

Roughly speaking, instead of modeling the distribution of the observations of each class, the boundaries between classes are modeled. Theoretically, both of these model-

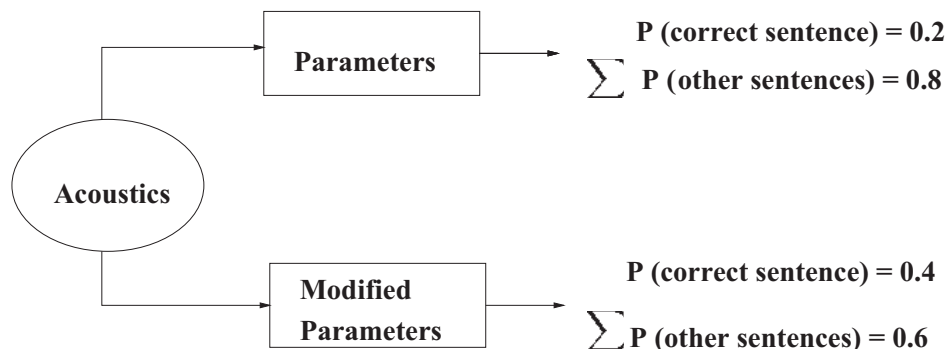


Figure 1.2: Goal - A Discriminant Training Algorithm

ing techniques can be used to achieve an optimal classifier. In practice, however, some assumptions have to be made in order to implement these approaches. In discriminant training model correctness is not assumed, and a training algorithm can directly minimize recognition errors. Generally, fewer free parameters are needed since the model describes boundaries, because there is no need to model the data everywhere (Lubensky *et al.* 1994; Renals *et al.* 1992).

### 1.3 Thesis Structure

Our starting point is to formulate the problem. This formulation serves both as a goal and an aid in evaluating existing solutions to this problem. It is followed by a description of the most popular models and algorithms used to solve this problem. These include Hidden Markov Models (HMMs) and hybrid systems of Artificial Neural Networks (ANNs) and HMMs, and their various training algorithms. In particular there is an emphasis on discriminant approaches such as Maximum Mutual Information (MMI) and Generalized Probabilistic Descent (GPD). The models are described in Chapter 2, and their training algorithms in Chapter 3.

In Chapter 4 the motivation for the use of transition-based models is given. Our early interest in modeling transitions was based on perceptual and physiological studies. Later, our own experiments reaffirmed the usefulness of it, and other recent related work strengthened this decision. In Chapter 5 we describe a particular model, Discriminant HMM (DHMM) that is a transition-based recognition model, and discusses its related mathematical theory. Experiments with this model motivated the need for a new training algorithm,



REMAP, which is described in Chapter 6. The training algorithm, REMAP, is an approach for Recursively Estimating and Maximizing A posteriori Probabilities of transition-based Hidden Markov Models given input sequences. It is a discriminant training algorithm; REMAP maximizes the posterior probability of the correct sentence while reducing the posteriors of all rival models. Experimental results are reported in Chapter 7. These experiments were both on isolated and continuous speech tasks. The results show an increase in the estimates of the posterior probabilities of the correct sentences after training, and a significant decrease in error rates in comparison to a baseline system. The thesis ends with conclusions and discussion of the remaining open problems, and future research directions in Chapter 8. The text proper is followed by an Appendix giving the proofs for theorems used in REMAP.

## Chapter 2

# Problem Formulation and Existing Models

### 2.1 Introduction

#### 2.1.1 Bayes Decision Theory

Bayes decision theory is a statistical approach to pattern classification. The approach is aimed at finding a *decision rule* that tells us which action to take for every possible observation. To simplify the discussion, let us assume a two-class classification problem wherein we have to choose between two classes  $M_1$  and  $M_2$ . We assume some a *priori probability*  $P(M_1)$  that the next observation belongs to  $M_1$ , and similarly for  $M_2$ . These a priori probabilities reflect our beliefs and prior knowledge about the models before seeing any evidence, e.g., if we needed to classify males and females based on their voices, we would assign equal prior probabilities<sup>1</sup> (Konig *et al.* 1993; Konig & Morgan 1993).

Each observation is a feature vector  $x$ , for instance (continuing the gender classification example) the fundamental frequency of the speaker and other speaker-dependent spectral information. The distribution of the feature vectors is class-dependent. Specifically, given a class say  $M_j$ , the distribution of the feature vectors is the likelihood function of the class,  $p(x|M_j)$ . For instance, given that the class is “male”, we give higher probability to feature vectors with fundamental frequency around 120 Hz than to feature vectors with

---

<sup>1</sup>We assume equal number of males and females in the training set.

fundamental frequency around 210 Hz. If the class is “female” it would be the other way around.

Bayes rule specifies how observing the value of  $x$  transforms the a priori probability  $P(M_j)$  into the a posteriori probability  $P(M_j|x)$  (for the 2-class example):

$$P(M_j|x) = \frac{P(x|M_j)P(M_j)}{P(x)} \quad (2.1)$$

where

$$P(x) = \sum_{j=1}^2 P(x|M_j)P(M_j) \quad (2.2)$$

Furthermore, according to Bayes decision theory, in order to minimize the average probability of error we should select the class with the highest a posterior probability  $P(M_j|x)$  (Duda & Hart 1973). Hence the general decision rule for  $I$  classes is:

$$\text{Choose } M_j \text{ if } P(M_j|x) > P(M_i|x) \quad \forall i \neq j \quad (2.3)$$

Note that this would be equivalent to

$$\text{Choose } M_j \text{ if } P(x|M_j)P(M_j) > P(x|M_i)P(M_i) \quad \forall i \neq j \quad (2.4)$$

as  $P(x)$  is fixed during recognition. Nonetheless, we are still left with the problem of choosing a model and estimating its parameters. The implication of Bayes decision rule is that we should optimize at training time the same measure that we use at recognition time, i.e., the a posteriori probability of the model given the observation.

### 2.1.2 Problem Formulation

#### Parameter Estimation

In statistical pattern classification as described above, it is known that a system leading to the minimum probability of error is the one that is trained to maximize the a posteriori probability of the correct class conditioned on the evidence (Duda & Hart 1973) and uses that same criterion during recognition. In real-life problems, however, we rarely have accurate knowledge about the structure of the probability functions. Given the typical high-dimensional input space and the limited training set, we can not evaluate the probability function directly from the training samples. Therefore, some parameterization of the probability function is needed, e.g., if we assume a normal density we only need to estimate its mean and covariance matrix.

In speech recognition two sources of knowledge are commonly used: “acoustic” and “language” knowledge. The acoustic knowledge provides the relation between the sound wave (or spectral patterns in the sound wave) and the linguistic identity of the utterance. The language information tells us about the phonemes and words, i.e., the pronunciation of each word in our vocabulary. In addition, the language model estimates the probability of a word given the hypothesized previous words.

Hence, these two sources of knowledge are used to parameterize the probability function that maps the acoustics of the utterance to the space of possible sentences. Specifically, if the speech signal is sampled at some fixed interval, and the acoustic vectors are concatenated, we can represent the input sequence to be classified with  $X = \{x_1, \dots, x_n, \dots, x_N\}$ . Additionally, we denote by  $M_i$  ( $i = 1, \dots, I$ ) one from the set of all possible sentences. Then we parameterize  $P(M_i|X)$  (the a posteriori probability of a sentence given a sequence of acoustic vectors) with  $L$  the parameter set which represents the language knowledge, (both a lexicon and a probabilistic grammar), and  $\Theta$ , the parameter set that represents the acoustic information. Thus, we want to estimate in training and use in recognition the following probability function  $P(M|X, L, \Theta)$ .

### Problem Formulation

Bayes Decision theory shows (as described above) that  $X$  will be optimally assigned to the sentence associated with model  $M_j$  if

$$M_j = \underset{M_i}{\operatorname{argmax}} P(M_i|X, L, \Theta), \quad i = 1, \dots, I \quad (2.5)$$

During training we should optimize the measure that we use in recognition. Thus, the ideal training algorithm should determine the set of parameters  $(\hat{\Theta}, \hat{L})$  that will maximize  $P(M_{w_j}|X_j, L, \Theta)$  for all training utterances<sup>2</sup>  $X_j$  ( $j = 1, \dots, J$ ), associated with  $M_{w_j}$ <sup>3</sup>, i.e.,

$$(\hat{\Theta}, \hat{L}) = \underset{(\Theta, L)}{\operatorname{argmax}} \prod_{j=1}^J P(M_{w_j}|X_j, L, \Theta) \quad (2.6)$$

---

<sup>2</sup>Note that  $J$  is the number of training sentences in the training set, while  $N$  is the number of acoustic vectors in a particular acoustic vector sequence.

<sup>3</sup> $M_{w_j}$  represents the correct model associated with the specific input sequence  $X_j$  that is known at training time. Strictly speaking,  $w_j$  is the index of the correct model for the acoustic vector sequence  $j$ .

with the following constraint:

$$\sum_{i=1}^I P(M_i|X_j, L, \Theta) = 1, \quad \forall j; \quad (2.7)$$

for every  $X$ , and where the sum over  $i$  represents the sum over all possible models. Note that this constraint makes the Maximum a Posteriori (MAP) criterion (2.6) discriminant. That is, when increasing the posterior probability of the correct model, the total probability mass assigned to all other models will automatically be reduced.

## Preview

The problem formulation above does not specify the model, i.e., how to estimate the posterior probability of a sentence given the parameters, acoustics, and the language model. In the following sections we describe several solutions to this problem. The discussion is divided into existing models and training algorithms. Specifically, we discuss models such as the Hidden Markov Model (HMM), hybrid systems of Artificial Neural Networks (ANNs) and HMMs, and segment-based models. The next chapter includes a discussion of several training algorithms and approaches for the models mentioned above. In particular there is an emphasis on discriminant approaches, such as Maximum Mutual Information (MMI) and Generalized Probabilistic Descent (GPD). Overall, we show that these solutions do not maximize the MAP criterion. They either maximize other criteria or approximate them.

### 2.1.3 Definitions and Notation

To facilitate the following discussion notation and basic terms are defined:

- A set of states  $\mathcal{Q} = \{q_1, \dots, q_K\}$ , that contains all the states from which phone and word models will be built. Each state class will be associated with a specific probability density function (PDF) or with specific statistical properties (see “conditional transition probabilities” in Section 5.2). For instance, if one wants to model the acoustic production in the beginning, middle, and end of a phone differently, each phone is modeled by three different states.
- $X = \{x_1, \dots, x_N\}$  is a sequence of acoustic vectors that is associated with a specific utterance. A sub-sequence of acoustic vectors that is local to the current vector,

extending  $c$  frames into the past and  $d$  frames into the future is expressed by  $X_{n-c}^{n+d} = \{x_{n-c}, \dots, x_n, \dots, x_{n+d}\}$ . Each acoustic vector is calculated from a short interval of speech, usually around 20 ms (shifted every 10 ms).

- $M_i$  is defined for  $i \in \mathcal{I} = \{1, 2, \dots, \mathcal{I}\}$ , the set of possible model indices;  $\mathcal{I}$  is the number of possible models (in the case of continuous speech,  $I$  indicates the number of possible sentences allowed by the grammar, although this is generally infinite). In continuous speech recognition the models represent all possible sentences.  $M_i$  is a directed graph with  $C_i$  states each of which belongs to  $Q$ .
- $M_{w_j}$ ,  $w_j \in \mathcal{I}$ , is the “correct” model associated with a specific training sequence  $X_j$ ,  $j = 1, \dots, J$ . Strictly speaking,  $w_j$  is the index of the correct model for the training sequence  $X_j$ .
- The parameter set describing all models is defined as  $\Theta = \{\lambda_1, \dots, \lambda_i, \dots, \lambda_I\}$ , in which  $\lambda_i$  represents only the parameters present in  $M_i$ . Of course, the different  $M_i$ 's, for  $i = 1, \dots, \mathcal{I}$ , can share some common parameters. In the hybrid systems discussed in this study, all models share the same set of parameters  $\Theta$  through a common neural network, which will be parameterized in terms of  $\Theta$ .
- The set of parameters that are only present in  $M_{w_j}$  will be denoted  $\Theta_{w_j}$ , which is a subset of  $\Theta$ .
- $q^n$  denotes the state at time  $n$ .
- $q_k^n$  means that state  $q_k$  has occurred at time  $n$ . Strictly speaking,  $q_k^n$  denotes the assignment of the value  $k$  to the random variable  $q^n$ .
- $\Gamma$  is a state sequence of length  $N$ . Each  $\Gamma$  is a realization of the stochastic state process, where the values assigned to the state process at each time step are taken from  $Q$ . Sometimes we write it explicitly as  $\Gamma_j = \{q_{\Gamma_j}^1, \dots, q_{\Gamma_j}^n, \dots, q_{\Gamma_j}^N\}$ . A state sequence which is legal in a given model  $M_i$  is also called a *path* in  $M_i$ .
- $P(\cdot)$  represents probabilities, while  $p(\cdot)$  denotes probability density functions (PDFs) and likelihoods.

## 2.2 Hidden Markov Models (HMMs)

### 2.2.1 Introduction

This section is a short review of the classical HMM approach to speech recognition. For a complete explanation, see (Huang *et al.* 1990; Levinson *et al.* 1983a; Rabiner 1989; Baker 1975; Jelinek 1976; Jelinek & Mercer 1980). Currently, this approach is very popular and a number of large-vocabulary, speaker-independent, continuous speech recognition state-of-the-art systems have been based on this approach. An HMM used for speech is a production model, and for each sentence it models the inherent statistical variations in speaking rate, pronunciation, and the differences between speakers (in speaker-independent models). The main idea is that we can approximate a continuous process as a sequence of short steady states. Roughly speaking we can model human speaking as moving between different states such that each state has fixed properties in terms of its repertoire of sounds.

In order to implement practical systems based on HMMs, a number of simplifying assumptions are typically made about the signal. For instance, although speech is a non-stationary process, HMMs model the sequence of feature vectors as a piecewise stationary process. That is, an utterance  $X = \{x_1, \dots, x_n, \dots, x_N\}$  is modeled as if it were produced by a succession  $L$  of discrete stationary states  $q_\ell \in \mathcal{Q}$ , with instantaneous transitions between these states. In this case, an HMM is defined (and represented) as a stochastic finite state automaton with a particular topology (usually strictly left-to-right, since speech is sequential). The approach defines two concurrent stochastic processes: the sequence of HMM states (modeling the temporal structure of speech), and a set of state output processes (modeling the [locally] stationary character of the speech signal). The HMM is called a “hidden” Markov model because there is an underlying stochastic process, the sequence of states, that is not observable but that affects the observed sequence of events. It is called “Markov” because the statistics of the current state are modeled as being dependent only on the current and the previous state (for the first-order Markov case).

Ideally, there should be an HMM for every possible utterance. However, this is clearly infeasible for all but extremely constrained tasks. Generally a hierarchical scheme must be adopted to reduce the number of possible models. First, a sentence is modeled as a sequence of words. To further reduce the number of parameters (and the required amount of training material) and to avoid the need of a new training each time a new word is added to the lexicon, sub-word units are usually preferred to word models. Although there are good

linguistic arguments for choosing units such as syllables or demi-syllables (Fujimura 1975; Segui *et al.* 1980; Levelt & Wheeldon 1994), the unit most commonly used is the phone (or context-dependent versions such as the triphone). This is the unit that we have generally used in our work, resulting in a selection of between 50 and 70 sub-word models. In this case, word models consist of concatenations of phone models (constrained by pronunciations from a lexicon), and sentence models consist of concatenations of word models (constrained by a grammar).

In the following section we describe the assumptions of HMM modeling.

### 2.2.2 Assumptions

Traditionally, in HMM modeling, the probability estimation process is divided into two parts: (1) *language modeling*, which does not depend on the acoustic data, and (2) *acoustic modeling*. The goal of the language model is to estimate prior probabilities of sentence models  $P(M_i|L)$ . The acoustic modeling role is to estimate the model-dependent probability densities  $p(X|M_i, \Theta)$ .

Additionally, several additional assumptions are usually required to make the estimation of  $p(X|M_i, \Theta_i)$  tractable:

- **Output-independence Assumption:** Acoustic vectors are not correlated (observational independence). The current acoustic vector  $x_n$  is assumed to be conditionally independent of the previous acoustic vectors (e.g.,  $X_1^{n-1}$ ). Furthermore, the assumption is that observations within the same speech segment (generated by the same state) are independent and identically distributed (i.i.d.), an unrealistic assumption given the non-stationary nature of speech.

To limit the impact of these assumptions, acoustic vectors at time  $n$  are usually complemented by their first and second time derivatives (Furui 1986b; Poritz & Richter 1986) computed over a span of a few frames, allowing very limited acoustical context modeling. Another solution to limit these assumptions is to consider a few adjacent frames (typically 3-5 frames in total) on which linear discriminant analysis is performed to reduce the dimension of the acoustic features (Haeb-Umbach & Ney 1992).

- **Markov Assumption:** Markov models are generally first-order Markov chains. Explicitly, the probability that the Markov chain is in state  $q_\ell$  at time  $n$  depends only



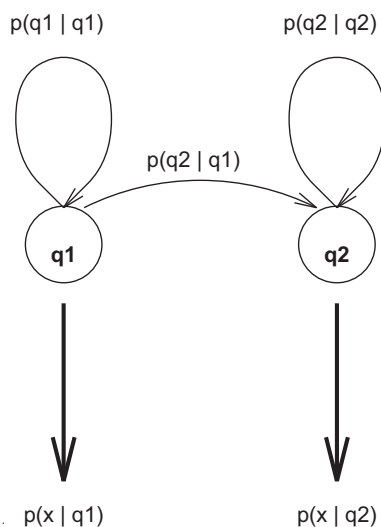


Figure 2.1: HMM - An Example

on the state of the Markov chain at time  $n - 1$ , and is conditionally independent of the past (both the past acoustic vector sequence and the states before the previous one).

### 2.2.3 Definition of the Hidden Markov Model

In this section we formally describe HMMs. An HMM can be defined by:

- $A = \{a_{ij} | a_{ij} = P(q^{n+1} = j | q^n = i)\}$ , a state transition probability distribution, where  $a_{ij}$  denotes the transition probability from state  $i$  to state  $j$ . Usually the assumption is that this probability distribution is the same for all time steps.
- $B = \{b_j(x_i) | b_j(x_i) = P(x_i | q = j)\}$ , for each state, there is a corresponding output probability (a discrete probability distribution in the discrete case and a continuous probability density function in the continuous case). It refers to the probability of generating some discrete symbol  $x_i$  in state  $q_j$ . Usually these probabilities are called emission probabilities.
- $\Pi = \{\pi_i | \pi_i = P(q^1 = i)\}$ , denotes the initial state distribution.

These definitions in addition to the relevant definitions in Section 2.1.3, represent an HMM,  $\lambda = (A, B, \Pi)$ . A very simplistic HMM is pictured in Figure 2.1.

## 2.3 Hybrid Systems

### 2.3.1 Multilayer Perceptrons (MLPs)

In this thesis, the discussion of neural networks for speech will be limited to the Multi-Layer Perceptron (MLP), a form of ANN that is commonly used for speech recognition. However, the analyses that follow are generally extensible to other kinds of ANN, e.g., a recurrent neural network (Robinson 1994), or a Time-Delay Neural Network (TDNN) (Waibel *et al.* 1989).

MLPs have a layered feed-forward architecture with an input layer, zero or more hidden layers, and an output layer. Each layer computes a set of linear discriminant functions (Duda & Hart 1973) (via a weight matrix) followed by a nonlinear function, which is often a sigmoid function

$$f(x) = \frac{1}{1 + \exp(-x)} \quad (2.8)$$

As discussed in (Bourlard & Morgan 1994), this nonlinear function performs a different role for the hidden and the output units. On the hidden units, it serves to generate high order moments of the input; this can be done effectively by many nonlinear functions, not only by sigmoids. On the output units, the nonlinearity can be viewed as a differentiable approximation to the decision threshold of a threshold logic unit or perceptron (Rumelhart *et al.* 1986), essentially to count errors. For this purpose, the output nonlinearity should be a sigmoid or sigmoid-like function. Alternatively, a function called *softmax* can be used. For an output layer of  $K$  units, this function is defined as

$$f(x_i) = \frac{\exp(x_i)}{\sum_{n=1}^K \exp(x_n)} \quad (2.9)$$

It can be proved that MLPs with enough hidden units can (in principle) provide arbitrary mappings  $g(x)$  between input and output. The MLP parameter set  $\Theta$  (the elements of the weight matrices) are trained to associate a “desired” output vector with an input vector. This is generally achieved via the Error Back-Propagation (EBP) algorithm (Rumelhart *et al.* 1986) that uses a steepest descent procedure to iteratively minimize a cost function in their parameter space. Since in our approach the HMMs will be described by the parameters of the neural network, we also denote the MLP parameter space by  $\Theta$ .

Popular cost functions are, among others, the Mean Square Error (MSE) criterion:

$$E = \sum_{n=1}^N \|g(x_n, \Theta) - d(x_n)\|^2 \quad (2.10)$$

or the relative entropy criterion<sup>4</sup>:

$$E_e = \sum_{n=1}^N \sum_{k=1}^K d_k(x_n) \ln \frac{d_k(x_n)}{g_k(x_n, \Theta)} \quad (2.11)$$

where  $g(x_n, \Theta) = (g_1(x_n, \Theta), \dots, g_k(x_n, \Theta), \dots, g_K(x_n, \Theta))^t$  represents the actual MLP output vector (depending on the current input vector  $x_n$  and the MLP parameters  $\Theta$ ),  $d(x_n) = (d_1(x_n), \dots, d_k(x_n), \dots, d_K(x_n))^t$  represents the desired output vector (as given by the labeled training data),  $K$  the total number of classes, and  $N$  the total number of training patterns.

MLPs, as well as other neurally-inspired architectures, have been used for many speech-related tasks. For instance, for some problems the entire temporal acoustic sequence is processed as a spatial pattern by the MLP. For isolated word recognition, for instance, each word can be associated with an output of the network. However, this approach has not been useful for continuous speech recognition and will not be discussed further here (Lippmann 1989).

### 2.3.2 Motivations

ANNs have several advantages that make them particularly attractive for Automatic Speech Recognition (ASR), e.g.:

- They can provide discriminant learning between speech units or HMM states that are represented by ANN output classes. That is, when trained for classification (using common cost functions such as MSE or relative entropy), the parameters of the ANN output classes are trained to minimize the error rate while maximizing the discrimination between the correct output class and the rival ones. In other words, ANNs not only can train and optimize the parameters of each class on the data belonging to that class, but also can attempt to reject data belonging to the other (rival) classes. This is in contrast to the likelihood criterion that does not lead to minimization of the error rate.

---

<sup>4</sup>In a number of references, including (Bourlard & Morgan 1994), this criterion is defined differently. In particular, the desired outputs are sometimes assumed to be independent, binary random variables and as a result this criterion gets a different form (which is sometimes called the cross entropy (Richard & Lippmann 1991)). However, viewing the network outputs as a posterior distribution over the values of one random variable (class conditioned on acoustic data), a discrete version of the classical definition of relative entropy may be used, as given here.

- Because ANNs can incorporate multiple constraints and find optimal combinations of constraints for classification, feature vectors do not need to be assumed independent. More generally, there is no need for strong assumptions about the statistical distributions of the input features (as is usually required in standard HMMs).
- They have a very flexible architecture that easily accommodates contextual inputs and feedback, and both binary and continuous inputs.
- ANNs are typically highly parallel and regular structures, which makes them especially amenable to high-performance architectures and hardware implementations.

A general formulation of statistical ASR can be summarized simply by a question: how can an input sequence (e.g., a sequence of spectral vectors) be explained in terms of an output sequence (e.g., a sequence of phones or words) when the two sequences are not synchronous (since there are multiple acoustic vectors associated with each pronounced word or phone)? It is true that neural networks are able to learn complex mappings between two vector variables. However, a purely connectionist formalism is not very well suited to solve the sequence-mapping problem. Most early applications of ANNs to speech recognition have depended on severe simplifying assumptions, e.g., small vocabulary, isolated words, known word or phone boundaries (Lippmann 1989). We shall see here that further structure (beyond a simple MLP) is required to perform well on continuous speech recognition, and that HMMs provide one solution to this problem. First, the relation between ANNs and HMMs must be explored.

### 2.3.3 MLPs as Statistical Estimators

MLPs can be used to classify speech classes such as words. However, MLPs by themselves classifying complete temporal sequences have not been successful for continuous speech recognition (Lippmann 1989). In fact, used as spatial pattern classifiers, they are not likely to work well for continuous speech since the number of possible word sequences in an utterance is generally infinite. On the other hand, HMMs provide a reasonable structure for representing sequences of speech sounds or words. One good application for MLPs is to provide the local distance measure for HMMs, while alleviating typical drawbacks such as lack of discrimination and assumptions of no correlation between acoustic vectors.

For statistical recognition systems, the role of the local estimator is to approximate

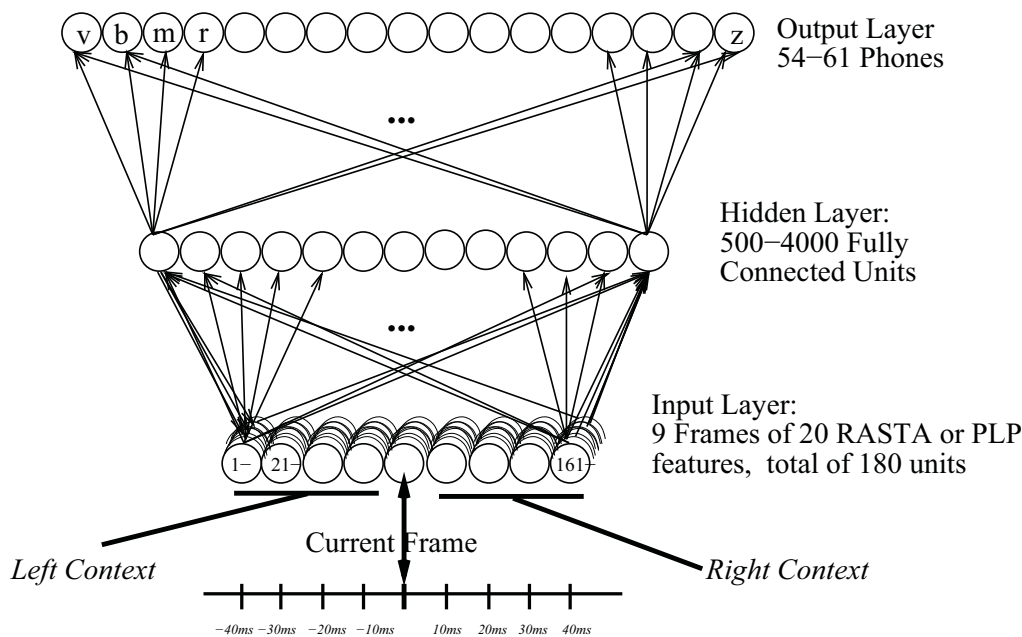


Figure 2.2: MLP - An Example

probabilities or probability density functions. In particular, given the basic HMM equations, we would like to estimate something like  $p(x_n|q_k)$ , which is the value of the probability density function (pdf) of the observed data vector given the hypothesized HMM state. The MLP can be trained to produce the *posterior* probability  $P(q_k|x_n)$  of the HMM state given the acoustic data. This can be converted to emission probability density function values using Bayes' rule.

Several authors (Bouclard & Wellekens 1989a; Bouclard & Morgan 1994; Gish 1990; Richard & Lippmann 1991) have shown that ANNs can be trained to estimate a *posteriori* probabilities of output classes conditioned on the input pattern. Recently, this property has been successfully used in HMM systems, referred to as *hybrid HMM/ANN* systems, in which ANNs are trained to estimate local probabilities  $P(q_k|x_n)$  of HMM states given the acoustic data (see, e.g., (Lubensky *et al.* 1994)).

Since MLPs require supervised training, all these systems have been used so far in the framework of Viterbi training, which provides the segmentation of the training sentences in terms of  $q_k$ 's and, hence, MLP training targets. A typical MLP is pictured in Figure 2.2. The principle of these systems is briefly recalled here.

Let  $c_k$ , with  $k = 1, \dots, K$ , be an output class of an MLP. Since we will use the

MLP for probability estimation associated with each HMM state  $q_k$  ( $k = 1, \dots, K$ ), there is a one-to-one equivalence between the  $q_k$ 's and the  $c_k$ 's that is associated with the discrete stationary states of  $Q$ . Also, we associate the parameter set  $\Theta$  as defined for HMMs with the MLP parameter set.

The output activation of the  $k$ -th MLP output class for a given set of parameters  $\Theta$  and an input  $x_n$  is denoted  $g_k(x_n, \Theta)$ . Since MLP training is supervised we will also assume the training set consists of a sequence of  $N$  acoustic vectors  $\{x_1, x_2, \dots, x_n, \dots, x_N\}$  labeled in terms of  $q_k$ 's. At time  $n$ , the input pattern of the MLP is acoustic vector  $x_n$ , and is associated with a state  $q_k^n$ .

For the popular MLP cost functions, it can be proved (as noted above) that the optimal MLP output values are estimates of the probability distribution over classes conditioned on the input  $\hat{P}(q_k|x_n)$ :

$$g_k(x_n, \Theta^{opt}) = \hat{P}(q_k|x_n) \quad (2.12)$$

if:

- the MLP contains “enough” parameters to be able to reasonably approximate the input/output mapping function,
- the network is not over-trained (which can be assured by stopping the training before the decline of generalization performance on an independent cross-validation set),
- the training does not get stuck at a local minimum.

In (2.12),  $\Theta^{opt}$  represents the parameter set minimizing (2.10) or (2.11).

It has been experimentally observed that, for systems trained on a large speech corpus, the outputs of a properly trained MLP do in fact approximate posterior probabilities (Bourlard & Morgan 1994).

This conclusion can easily be extended to other cases. For example, if we provide the MLP input not only with the acoustic vector  $x_n$  at time  $n$ , but also with some acoustic context  $X_{n-c}^{n+d} = \{x_{n-c}, \dots, x_n, \dots, x_{n+d}\}$ , the output values of the MLP will estimate

$$g_\ell(x_n, \Theta^{opt}) = \hat{P}(q_\ell^n | X_{n-c}^{n+d}), \quad \forall \ell = 1, \dots, K \quad (2.13)$$

This windowing over time has been used in the standard hybrid HMM/ANN system (briefly summarized later in this section) to account for correlation between acoustic vectors. If the

previous class is also provided to the input layer (leading to a quasi-recurrent network), the MLP output values will be estimates of

$$g_\ell(x_n, \Theta^{opt}) = \hat{P}(q_\ell^n | X_{n-c}^{n+d}, q_k^{n-1}), \quad \forall k, \ell = 1, \dots, K \quad (2.14)$$

It will be shown in Chapter 4 that this is a form of the local probability the hybrid HMM/MLP theory tells us to use. This will be referred to as the “conditional transition probability” and will be the major theme throughout this thesis.

In addition, this conclusion remains valid for other kinds of networks, given similar training conditions. For example, recurrent networks (Robinson 1994) and radial basis function networks (Renals *et al.* 1991) can also be used to estimate posterior probabilities.

There is another important generalization of this property that will be essential later in this thesis. If the ANNs are trained with an estimate of the posterior probabilities of the output states (as opposed to the “1-from-K” binary output targets used for a classification mode training), then (2.12) remains valid. In other words, if the targets come from some independent “expert”, the net will learn to produce posterior probabilities as well.<sup>5</sup> Although this property is mentioned in (Bourlard & Wellekens 1989a; Bourlard & Morgan 1994; Richard & Lippmann 1991), it has not previously used in hybrid HMM/MLP systems because of the lack of a full algorithm for the convergence to better probabilities. Such an algorithm has now been developed, and will be presented in this thesis.

### Estimating HMM Likelihoods with MLP

Since the network outputs approximate Bayesian probabilities,  $g_k(x_n, \Theta)$  is an estimate of

$$P(q_k | x_n) = \frac{p(x_n | q_k)P(q_k)}{p(x_n)} \quad (2.15)$$

which implicitly contains the a priori class probability  $P(q_k)$ . It is thus possible to vary the class priors during classification without retraining, since these probabilities occur only as multiplicative terms in producing the network outputs. As a result, class probabilities can be adjusted during use of a classifier to compensate for training data with class probabilities that are not representative of actual use or test conditions (Richard & Lippmann 1991).

---

<sup>5</sup>Actually, it is easy to prove that, for the popular MLP cost functions,  $g(x_n)$  will be an estimate of  $E\{d(x_n) | x_n\}$ , where  $E$  stands for the expected value.

Scaled likelihoods  $p(x_n|q_k)$  for use as emission probabilities in standard HMMs can be obtained by dividing the network outputs  $g_k(x_n)$  by the relative frequency of class  $q_k$  in the training set, which gives us an estimate of:

$$\frac{P(q_k|x_n)}{P(q_k)} = \frac{p(x_n|q_k)}{p(x_n)} \quad (2.16)$$

During recognition, the scaling factor  $p(x_n)$  is a constant for all classes and will not change the classification. It could be argued that, when dividing by the priors, we are using a scaled likelihood, which is no longer a discriminant criterion. However, this need not be true, since the discriminant training has affected the parametric optimization for the system that is used during recognition. Thus, this permits the use of the standard HMM formalism, while taking advantage of ANN characteristics.

## 2.4 Non-stationary Modeling

This section includes a brief a number of models that have been proposed to remedy some of the shortcomings of HMMs. A quick solution might be to represent each unit of speech by enough different states to approximate its non-stationary nature in a stepwise fashion. For instance, a vowel could be represented by ten different states. This solution has two major limitations:

- There are too many free and independent model parameters. This necessitates more training data, and also might be more prone to capturing irrelevant sources of variance in the data than a simpler model.
- Such a model does not capture the correlation and dependence between the different states. For states with a short duration, this would be even more pronounced, since the change between two states in a sequence would correspond to only a small movement of articulators for a given speaker.

Several extensions to the basic HMM have been proposed in order to overcome some of these deficiencies. For example, autoregressive HMMs condition the emission probability of a given state on previous observations (Juang & Rabiner 1985). However, none of these extensions have explicitly modeled the emission in a given phone as a non-stationary process. In general this is too difficult to handle with a practical number of parameters.



A number of HMM alternatives model the sequence of frames emitted in a given sub-word unit as correlated and dependent on each other (Digalakis 1992; Ostendorf & Roukos 1989; Ghitza & Sondhi 1993; Deng 1992) (see also Section 4.3.2). The models differ in their assumptions about the nature of the correlation between the frames in the sequence. For instance, some assume that only consecutive frames are correlated, while others assume that all the frames in the sequence are dependent on each other. In general these models do not require the HMM assumption of independent and identically distributed observations. In the following section we survey segment-based approaches that are in this family.

### 2.4.1 Segment-Based Approaches

#### Introduction

In segment-based models the basic unit is a sequence of acoustic vectors emitted in a given speech unit (a “segment”), as opposed to a single acoustic vector as used for HMMs. The production of the acoustic vectors in a segment may be described as a three step procedure (Digalakis 1992):

1. Select the length of the segment according to  $P(L|s_k)$ , where  $L$  is the random variable that denotes the length of the segment, and  $s_k$  is a particular speech unit.
2. Generate a fixed length segment  $M$  according to the distribution  $P(y_1, y_2, \dots, y_M|s_k)$ . The distribution models the trajectory of the sound in the feature vector space.  $M$  is chosen to be greater than all the possible values of  $L$ .  $Y = y_1, y_2, \dots, y_M$  is called the *hidden* sequence of acoustic vectors.
3. Down-sample  $Y$  using the time-warping transformation  $T_L$  and output the observed sequence of acoustic vectors  $X = \{x_1, x_2, \dots, x_L\}$ . This transformation can be either linear or non-linear depending on the specific segmental model.

#### Segmental Models

These models differ in the form of the distribution ( $Y|s_k$ ) and in the time-warping transformation  $T_L$ . Ostendorf and Roukos (Ostendorf & Roukos 1989) have used (among a number of methods) linear time sampling in their study, i.e., sampling  $Y$  in equal intervals along the time axis as their time warping procedure. Their specific implementation had ten 14-dimensional vectors of cepstral coefficients. They used a multivariate Gaussian to

represent the entire segment, which can require a 140 by 140 full covariance matrix for each phone (assuming that feature dependence is accounted for).

Ghitza and Sondhi developed a model (Ghitza & Sondhi 1993) that can also be viewed as a stochastic segment model with the following distinctions:

- Their warping procedure is a dynamic time warping technique, instead of the linear time warping method used by Ostendorf and Roukos.
- They used diphones as their sub-word units, as opposed to the phones in Ostendorf and Roukos' stochastic segment model (Ostendorf & Roukos 1989).
- They maintained the HMM framework and assumed a semi-hidden Markov chain, i.e., each state has an explicit duration distribution.

These stochastic segment models are not inherently subject to the constraints of the i.i.d. assumptions discussed earlier. However, there are some practical difficulties:

1. There are many free parameters that must be estimated reliably from the data, e.g., a large covariance matrix. As a result, independence assumptions are often made, leading to less powerful models.
2. These methods explicitly assume a particular parametric form for the observation distributions, e.g., multivariate Gaussian. This assumption is already faulty for standard HMMs, but may be even a worse approximation once observation interdependencies are taken into account. (Nonetheless, it is a sensible place to start.)
3. All the models assume a given segmentation, e.g., the knowledge of the boundaries between the basic speech units, which is difficult to obtain. One solution is to do an exhaustive search of all reasonable segmentations.
4. Warping the data to a fixed length segment may delete or obscure relevant information.

### **A Stochastic Dynamic System Approach**

This model assumes a discrete-time, linear, stochastic dynamic system, with a state process as the source for the observation process. To model an underlying dynamic system, some assumptions are required. For example, Digalakis has proposed two possible model constraints:

1. *Trajectory invariance*: It is assumed that the unobserved trajectory of state vectors in the state space is the same for each speech segment length and is the source for all possible realizations of the speech segment. Given that the state vector at each time step is a vector random variable, this translates into a fixed sequence of state transition matrices. The observed speech segment is then a down-sampled version of the trajectory of the feature vectors created by the system. For a long realization (a long observation sequence) the underlying trajectory is sampled at shorter intervals than a short observation sequence. Consequently, long observation sequences have higher correlation between successive frames than short observation sequences.
2. *Correlation invariance*: It is assumed that the underlying trajectory in state space varies with the realization length, and the sequence of state-transition matrices for a particular realization depends on the realization length. It is assumed that the correlation between two observations depends only on the relative location of the observations in the segment. The correlation is invariant under the time-warping transformation. Roughly speaking the trajectory length is chosen according to the length of the realization. Thus, the state change rate is slower than under the trajectory invariance assumption, making it somehow a more realistic approach.

In his study, Digalakis assumes that the observed segment of speech is the output of a piecewise time-invariant linear dynamical system. He uses up to five invariant regions for each model. The models based on the *correlation invariance* assumption outperformed the models based on the *trajectory invariance* assumption for the task of phone classification. For more details see (Digalakis 1992; Digalakis *et al.* 1993). The stochastic dynamic system approach appears to have more modeling power than an HMM, and can potentially capture the dynamics of acoustic vectors within a segment of speech. However, there are still open issues about the structure of the dynamic system, such as the arguable assumptions of linearity and several types of invariance.

## Chapter 3

# Training Algorithms and Optimization Criteria

### 3.1 Likelihood Estimation and Training

#### 3.1.1 Introduction

Theoretically the optimal way to classify an input sequence is to choose the class with the highest posterior probability given this sequence (Duda & Hart 1973). Therefore, at training we want to maximize the posterior probability of the correct model (sentence) given the evidence (sequence of acoustic vectors). An optimization criterion for parameter estimation that achieves this goal during training, is the Maximum a Posteriori (MAP) criterion. Most speech recognition systems however, are trained according to a maximum likelihood criterion that maximizes, in the parameter space, the likelihood of the data given some model. In this section, besides describing the ML criterion, we discuss under what assumptions it results in the best possible parameter set.

In HMMs, this likelihood can be represented as  $P(X|M, \Theta)$ . Training HMMs according to the ML criterion is aimed at finding the best set of parameters  $\hat{\Theta}$  such that:

$$\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} \prod_{j=1}^J P(X_j|M_{w_j}, \Theta) \quad (3.1)$$

Note that this criterion is different from the criterion that is specified in (2.6). The main difference is that ML is not discriminant; maximization of (3.1), the likelihood of the correct

sentence producing the utterance, does not necessarily decrease the likelihood of all other models, unlike the MAP criterion that “ties” all models by (2.7).

### 3.1.2 The Relation between ML and MAP

Classically, the relation between likelihood estimation and MAP estimation is obtained by applying Bayes’ rule:

$$P(M|X, L, \Theta) = \frac{P(X|M, L, \Theta)P(M|L, \Theta)}{P(X|L, \Theta)} \quad (3.2)$$

For practical reasons, it is assumed that:

1. The parameters  $\Theta$  of the acoustic model are independent of the parameters  $L$  of the language model, yielding

$$\begin{aligned} P(X|M, L, \Theta) &\approx P(X|M, \Theta) \\ P(M|L, \Theta) &\approx P(M|L) \end{aligned} \quad (3.3)$$

2. Despite the fact that  $\Theta$  and  $L$  vary during training,  $P(X|L, \Theta)$  is assumed to be constant.

Incorporating these assumptions for a particular model  $M_{w_j}$  associated with the training utterance  $x_j$  we get,

$$P(M_{w_j}|X_j, L, \Theta) = \frac{p(X_j|M_{w_j}, \Theta)P(M_{w_j}|L)}{p(X_j|L, \Theta)} \quad (3.4)$$

Since  $p(X_j|L, \Theta)$  is not constant during training

$$\begin{aligned} P(M_{w_j}|X_j, L, \Theta) &\approx; \\ &\frac{p(X_j|M_{w_j}, \Theta)P(M_{w_j}|L)}{\sum_{i \neq w_j} p(X_j|M_i, \Theta)P(M_i|L) + P(X_j|M_{w_j}, \Theta)P(M_{w_j}|L)} \end{aligned} \quad (3.5)$$

Maximum likelihood training maximizes  $p(X_j|M_{w_j}, \Theta)$  which appears both in the numerator and the denominator of (3.5). However, it does not necessarily maximize the left side, which is our goal. Furthermore, besides increasing the probability of the correct model we might also be increasing the probabilities of the incorrect models.

### 3.1.3 Implementation - The EM Algorithm

The most popular approach to iterative maximization (3.1) has been described in a number of classic papers (Baum & Petrie 1966; Baum *et al.* 1970; Baum 1972; Liporace 1982). Starting from initial guesses  $\Theta^0$ , the model parameters are iteratively updated according to the “Forward-Backward” algorithm [or, equivalently, the Expectation-Maximization (EM) algorithm (Dempster *et al.* 1977)] so that (3.1) is maximized at each iteration. This kind of training algorithm, often referred to as Baum-Welch training in the particular case of HMMs, can also be interpreted in terms of gradient techniques (Levinson *et al.* 1983a; Levinson 1985).

Below we describe the general idea behind the EM algorithm. Besides its usage in HMM training as described above, it has a role in the training algorithm that will be presented in Chapter 6.

#### The EM Algorithm

This section is a brief description of the EM algorithm. Its application to HMM parameter estimation is described in (Huang *et al.* 1990; Levinson *et al.* 1983b; Baker 1975; Jelinek 1976; Jelinek & Mercer 1980; Lee 1989). Roughly speaking, there is an optimization problem that would be greatly simplified by the knowledge of additional variables  $Y$ . For example, optimizing (3.1), the likelihood of producing the acoustics, would be easier if the state sequence that produced the data was known. So we *estimate* the missing data  $Y$  (state sequence) using the observed data (acoustic sequences) and the current set of parameters  $\Theta_t$ . We *maximize* the function using the estimated  $Y$  and we get a new set of parameters  $\Theta_{t+1}$  that lead to a new estimate of the missing data  $Y$ . The estimation and maximization steps are iterated until the guaranteed convergence to a local maximum (Dempster *et al.* 1977).

More precisely, following (Dempster *et al.* 1977), the goal is to maximize the following likelihood function  $L(X; \Theta)$ , where  $\Theta$  is the set of the parameters of the function and  $X$  the observations. The function estimates the log-likelihood of producing the observations given  $\Theta$ . The problem would be simplified by the knowledge of additional variables  $Y$ . In that case we will maximize  $L_c(X \cup Y; \Theta)$ , the log-likelihood of producing the complete data  $X \cup Y$  given  $\Theta$ , and usually called the complete likelihood. However, since  $Y$  is not observable, the EM algorithm relies on integrating over the distribution of  $Y$ , with the following

auxiliary function:

$$Q(\Theta, \Theta_t) = E_Y[L_c(X \cup Y; \Theta) | X, \Theta_t] \quad (3.6)$$

which is the expected value of the complete data likelihood, given the observed data  $X$  and the current estimate of the parameters  $\Theta_t$ .

The main difference between the likelihood function and the auxiliary function is that the auxiliary function is a deterministic function and we can optimize it, while the likelihood function is a stochastic one (it depends on the unobserved data  $Y$ ). The EM algorithm iterates between the following two steps:

- **Estimation:** Compute  $Q(\Theta, \Theta_t)$  based on  $\Theta_t$ .
- **Maximization:** Maximize the auxiliary function.

$$\Theta_{t+1} = \underset{\Theta}{\operatorname{argmax}} Q(\Theta, \Theta_t) \quad (3.7)$$

In the case that we are unable to maximize the auxiliary function as required by the EM algorithm  $Q(\Theta, \Theta_t)$ , but can only increase it, the algorithm is called Generalized EM (GEM), and is also guaranteed to converge to a local maximum.

### 3.1.4 Summary

The two main conceptual problems with the maximum likelihood approach are:

- It is implicitly assumed that the model (with all its assumptions relative to its topology and probability density functions) is accurate and reflects the structure of the data (although the data might not adhere to the constraints imposed by the HMMs). If we had enough training data, it might be preferable to infer all the parameters of the models (including topology and non-parametric probability density functions) directly from the data. This can be seen as implicitly using a Bayes or MAP criterion (maximizing  $P(M|X)$ ) during training instead of ML. Since the posterior probability includes the effects of prior information, the language model would also be inferred from the training data. However, it appears that this would require a prohibitive amount of training data. Nonetheless, the role of domain-specific knowledge is irreplaceable. It has the role of reducing the search space, i.e., the search is constrained to

all models that are feasible given our knowledge. In addition, the use of other knowledge sources might result in better generalization, as the dependency on a particular training set is weaker. Hence, a combination of data-driven and knowledge-based approaches is desired.

- By training with ML instead of MAP, we strongly reduce the discriminant properties of HMMs. Ideally, each HMM should be trained not only to generate high probabilities for its own class, but also to discriminate against rival models.

Both of these two points (but particularly the second one) are related to the discussion that follows on discriminant criteria for HMM training.

The ML criterion can lead to optimal recognition performance only if the model is an exact statistical model of the source and the amount of training data is infinite (Nadas 1983; Duda & Hart 1973). However, these conditions are rarely (if ever) satisfied in speech recognition.

## 3.2 Discriminant Approaches

### 3.2.1 Introduction

In this section we describe discriminant approaches that have been proposed for speech recognition in general, and for HMMs in particular. Roughly speaking, discriminant approaches try to model the boundaries between the classes, and not model the classes themselves. Thus, a priori knowledge of class models does not play the same role. With an inaccurate model, the best we can do is to optimize its ability to distinguish between the underlying classes, which is typically achieved by replacing the ML criterion by a discriminant one.

### 3.2.2 Maximum Mutual Information (MMI)

Initially introduced in (Bahl *et al.* 1986; Brown 1987), this method aims at maximizing the mutual information (Cover & Joy 1991) between a set of (sentence) models  $M_{w_j}$  and the associated sequences of acoustic vectors  $X_j$ . This mutual information is then defined as (Cover & Joy 1991):

$$I(M_{w_j}, X_j | L, \Theta) = \sum_{M_{w_j}, X_j} p(M_{w_j}, X_j | L, \Theta) \log \frac{p(M_{w_j}, X_j | L, \Theta)}{P(M_{w_j} | L, \Theta) p(X_j | L, \Theta)} \quad (3.8)$$



$$= E_{p(M_{w_j}, X_j|L, \Theta)} \left\{ \log \frac{p(M_{w_j}, X_j|L, \Theta)}{P(M_{w_j}|L, \Theta)p(X_j|L, \Theta)} \right\} \quad (3.9)$$

where  $\Theta$  is the whole parameter space (for all models) in which optimization is performed and the sum over  $(M_{w_j}, X_j)$  represents a sum over all training utterances.  $E_{p(M_{w_j}, X_j|L, \Theta)}$  stands for the expected value according to the mass function  $p(M_{w_j}, X_j|L, \Theta)$ . For one particular  $(M_{w_j}, X_j)$  set, we then have:

$$\begin{aligned} I(M_{w_j}, X_j|L, \Theta) &= \log \frac{p(M_{w_j}, X_j|L, \Theta)}{P(M_{w_j}|L, \Theta)p(X_j|L, \Theta)} \\ &= \log \frac{p(X_j|M_{w_j}, L, \Theta)}{\sum_{i=1}^I p(X_j|M_i, L, \Theta)P(M_i|L, \Theta)} \end{aligned} \quad (3.10)$$

As already mentioned in Section 2.2, the language model parameters  $L$  are often assumed independent of the acoustic parameters  $\Theta$  and are estimated independently from a (large) text copra. Furthermore, the likelihoods  $p(X_j|M_i, \Theta)$  depend only on the parameters  $\Theta_i$  present in  $M_i$ . As a consequence, (3.10) can be rewritten as

$$I(M_{w_j}, X_j|L, \Theta) = \log \frac{p(X_j|M_{w_j}, \Theta_{w_j})}{\sum_{i=1}^I p(X_j|M_i, \Theta_i)P(M_i|L)} \quad (3.11)$$

in which the contribution of each term in the denominator is weighted according to the prior probability of the associated sentence as given by the language model, independent of the acoustic training data.

It is obvious that (3.10) and (3.11) are discriminant criteria. In (Bahl *et al.* 1986) it is shown that it is possible to get some kind of re-estimation recursion of local probabilities but, unfortunately, there is no proof that the recursion converges and there is no guarantee that the new estimates of (e.g., transition) probabilities are positive. As a consequence, a local gradient ascent method is usually used for optimization and the standard (likelihood-based) forward-backward recurrences are used to estimate the gradient. This is similar to the Alpha-Nets presented in (Bridle 1990) in which the gradient of the mutual information criterion takes the form of the backward recurrence used in the Baum-Welch algorithm. In the framework of hybrid HMM/ANN systems, this MMI criterion has been used in (Bengio *et al.* 1992), in which the ANN generates the sequence of acoustic vectors for the HMM and is trained to optimize the (global) MMI. In that paper, it is shown that it is possible to compute the gradient of the HMM training criterion (MMI or ML) with respect to the parameters of the ANNs.

However, in addition to “theoretical” problems, this algorithm suffers from a “practical” problem for continuous speech recognition. Indeed, optimization of  $\Theta$  to maximize (3.10) requires not only a forward recurrence for the numerator, but also many forward recurrences for the denominator to estimate the contribution of all possible rival models.

Several solutions have been proposed to alleviate this problem, including:

1. If phoneme models are trained, the use of a “looped” phonetic model, i.e., a word model that allows any possible phoneme sequence (Merialdo 1988). This model may generate all possible phoneme sequences and, by running the forward recursion through it, may provide the summed probability in the denominator of (3.10).
2. Estimating the denominator in (3.10) by running an N-best algorithm (Schwartz & Chow 1990; Schwartz & Austin 1991), that provides the N-best (rival) sentences through which we run the forward recursion.

With the algorithm proposed in this thesis, in addition to all the advantages of “standard” HMM/MLP hybrids (i.e., local nonlinear discrimination, time correlation and no significant assumptions about probability density functions as discussed below), multiple forward recursions are not needed. Also, all probabilities will always be estimates of actual (local and global) posterior probabilities, will be positive, and will sum to one.

### 3.2.3 Generalized Probabilistic Descent (GPD)

Generalized Probabilistic Descent (GPD) is another discriminant approach that is sometimes used to train speech recognition systems. GPD is actually very close in spirit to MMI, although it permits generalization to different kinds of training criteria (Katagiri *et al.* 1991).

The idea of GPD is simple and can be summarized as follows. Given the whole set of parameters<sup>1</sup>  $\Theta$ , define a discriminant function associated with each (word or sentence) model  $M_i$  as  $g_i(X; \Theta)$ . This discriminant function can be any differentiable distance function or probability distribution. Several instances of this are discussed in (Katagiri *et al.* 1991), each of them leading to different interpretations (as is also the case for MMI and MAP training). However, often the discriminant function is defined as:

$$g_i(X; \Theta) = -\log p(X|M_i, \Theta) \tag{3.12}$$

---

<sup>1</sup>In the following discussion we omit the language model  $L$

Here again, (3.12) can be considered as the “full” (word or sentence) likelihood, the best-path (Viterbi) approximation (referred to as “segmental GPD training”) or any intermediate solution like a sum over the  $S$ -best matching path scores. Another solution could be to define  $g_i(X; \Theta)$  as the MMI in (3.11). However, since this will then be used in a discriminant measure (as defined below) taking all the classes into account, it can be easily shown that using MMI or full likelihoods as discriminant functions results in the same misclassification measure.

Classification (i.e., recognition) will then be based on that discriminant function according to the (usual) rule

$$X \in M_j \text{ if } j = \underset{i}{\operatorname{argmax}} g_i(X; \Theta) \quad (3.13)$$

Given this discriminant function, we can define a misclassification measure that will measure the “distance” between one specific class and all the others. Here again, several measures can be used, each of them leading to different interpretations. However, one of the most general ones given in (Katagiri *et al.* 1991) is:

$$d_j(X; \Theta) = g_j(X, \Theta) - \log \left[ \frac{1}{I-1} \sum_{i \neq j} \exp(\eta g_i(X; \Theta)) \right]^{1/\eta} \quad (3.14)$$

in which  $I$  represents the total number of possible reference models.

It is easy to see that if  $\eta = 1$ , (3.14) is then similar to (3.11), assuming (3.12), as a weighted ratio of the likelihood of the correct model by the likelihoods all models.

The error measure (3.14) could be used as the criterion for optimization by a gradient-like procedure, which would result in something very similar to MMI training. However, the goal of GPD is to minimize the actual misclassification rate, which can be achieved by passing  $d_j(X; \Theta)$  through a nonlinear, nondecreasing, differentiable function  $F$  (such as a sigmoid function) and then minimizing

$$E(\Theta) = \sum_j \sum_{X \in M_j} F(d_j(X; \Theta)) \quad (3.15)$$

Other functions can be used to approximate the error rate. For example, we can also assign zero cost when an input is correctly classified and a unit cost when it is not properly classified, which is then another formulation of the minimum Bayes risk.

As briefly shown above, this approach is certainly very general and includes several discriminant approaches as special cases. For some problems such as continuous speech

recognition, however, this approach has the same potential difficulty as MMI, i.e., the need to estimate “scores” (whatever they might be) of both the correct model and all possible rival models.

## Chapter 4

# Transition-based Modeling

### 4.1 Introduction

In this section we describe our motivation for studying a transition-based modeling approach to speech recognition, a particular model (the time-index model), and our experiments with the model. Our early interest in modeling transitions between speech units was motivated by perceptual and physiological studies. Therefore, we decided to study the following question: given accurate transition information between speech units, can we significantly improve recognition performance? Specifically, given boundary information between speech units, we explicitly condition the emission probability of a state on the time index, where time index is defined as the number of frames between entering a state and the current frame. Experiments with the time-index model as described below established the necessary condition that accurate transition information can significantly improve recognition performance. These results were corroborated by a recent study by Goldenthal (Goldenthal 1994). Goldenthal found a consistent improvement in phone recognition results when enhancing his segment-based models with explicit transition models. He used a set of 200 *canonical* transitions that were created by clustering all the transitions in the training set. Each canonical transition modeled the trajectory of a fixed number of frames centered about the transition boundary.

## 4.2 Perceptual and Physiological Point of View

A popular model of human speech production is that an utterance is organized as a succession of vocal-tract states, where each of these states represents a different configuration of the articulators. Further, at the level of articulatory performance, the gestures invoked to actualize these states are relatively slow. They merge spatially and temporally into a continuous process that usually only approximates the intended states and is highly sensitive to co-articulation (Deng & Sun 1994). For example, when vowels are co-articulated with consonants, the spectral pattern of the speech signal varies such that the acoustic targets found in isolated vowels are never fully realized in the changing spectrum. This is usually called “undershoot” (Lindblom & Studdert-Kennedy 1967). Obviously, in humans this problem is normally coped with.

The question whether human phonemic recognition is context-free or context-dependent has been addressed by Lindblom and Studdert-Kennedy among others (Lindblom & Studdert-Kennedy 1967). In their experiments, they tested the role of formant transitions in vowel recognition. Specifically they tested whether vowel recognition is done by the steady-state of the formants, i.e., zero rate of change, or by the short-term acoustic context, such as the direction and rate of adjacent formant transition. In the experiment, American English listeners were asked to identify monosyllabic nonsense speech. Each consonant-vowel-consonant (CVC) syllable was a sequence of three elements: transition + pattern at point of closest approach to target + transition. The vowel-formant patterns assigned to the points of closest approach to the target were selected from a continuum ranging from [I] to [u]. The rate and direction of the adjacent transitions were varied by the choice of two consonantal frames: [w-w] and [j-j]. Their results showed that the identity of the vowel stimulus is determined not only by the formant pattern at the point of closest approach to the target, but also by the direction of the adjacent formant transitions. For instance the same vowel pattern:  $F_1 = 350Hz, F_2 = 1578Hz, F_3 = 2604Hz$  was recognized by almost all the subjects as [I] in the context of [jVj] and was recognized by most of them as [U] in the context of [wVw]. In general, it was shown that the categorization of the continuum is adjusted in the different environments so as to compensate for an undershoot effect in the vowel stimuli.

In a more recent study Furui has suggested that sufficient information for both vowel and consonant identification is contained across the same initial part of each syllable

(Furui 1986a). This part includes the maximum spectral transition. He also verified that the steady-state portions of the formants are not the only key for either vowel or syllable perception.

Auditory physiologists have collected a vast amount of data describing the response of mammalian auditory-nerve fibers to simple signals (Kiang 1984; Ruggero 1994) as well as more complex signals such as synthetic speech (Sachs *et al.* 1988). From these data it is clear that some sort of frequency analysis is performed and the dynamics of the response to non-steady-state signals is an important aspect of the auditory processing. Specifically, during the initial 15 ms of acoustic stimulation, the discharge rate of auditory-nerve fibers is often significantly higher than during the steady-state level. The decrease in response rate is referred to as “adaptation.” Usually there is a very rapid initial decay in rate immediately after the onset, followed by a slower decay to a steady-state level, about 50 ms after the onset (Smith & Zwislocki 1975).

## 4.3 The Time-Index Model

### 4.3.1 Introduction

Motivated by the studies mentioned above we decided to consider the following question: given accurate transition information between speech units can we significantly improve recognition performance? The model that we chose to answer this question is the time-index model. The main idea behind the time-index model is that all the trajectories of a phone in the acoustic vector space share a stochastic dependency on the time elapsed since the beginning of the phone. This dependency can be modeled by a parametric distribution as in Deng’s model (Deng 1992) or by using an MLP in our model.

### 4.3.2 Deng’s Trended HMM

Deng described a model that explicitly conditioned the emission probability of a state on the time index, i.e., on the number of frames between the current frame and the previous state transition. For example, if the Markov chain has two states A and B, and we assume a specific realization that alternates between the states, the time index for a given state is depicted in Figure 4.1 (the figure does not show all the “machinery” of the HMM). Deng has coined his model the “trended HMM” (Deng 1992). In this model, a sequence of

**State Sequence : A A B B B A A**  
**Time Index:     1 2  1 2 3  1 2**

Figure 4.1: Depiction of the time index along with the current state

observation vectors generated in a given state is a combination of a stationary process and a deterministic function of time, as illustrated in the following equation for the multivariate normal distribution:

$$p(x_t|state, ti) = \frac{\exp(-(x_t - g_{state}(ti))^T(\Sigma)^{-1}(x_t - g_{state}(ti)))}{(2\pi)^{\frac{n}{2}}(\det\Sigma)^{0.5}} \quad (4.1)$$

where  $ti$  is the time index as defined above,  $g_{state}()$  is a deterministic function of the time index and has parameters that may differ from state to state. In this simplified example,  $g_{state}()$  shifts the mean vector of distribution as a function of the time index, while the stationary part is the variance-covariance matrix  $\Sigma$ . In principle this model explicitly conditions the emission probability on the time index, and a sequence of observations emitted from a given state is no longer assumed i.i.d. We do not know however the optimal form of  $g_{state}()$  for each unit of speech. For example, one would expect a different time index dependence for vowels and stops. In the following, the time-index idea has been incorporated into a connectionist context.

### 4.3.3 An Introduction to the Time-Index Model

We propose a time-index model that 1) differs from an HMM in that the observations emitted in a given phone are no longer identically distributed and that 2) differs from Deng's model and others by its use of posterior probabilities estimated by a connectionist network. In the time-index model, the realizations of the state process are no longer sequences of values taken from the phone set, but rather are chosen from a set of pairs consisting of a phone and a time index. The time index is defined as the number of frames between the current frame and the previous state transition. For this model, the probability of generating a sequence of observations  $X = \{x_t, x_{t+1}, \dots, x_{t+T}\}$  in a given phone  $q_j$  is:

$$P(X_t^{t+T}|q_j) = \prod_{i=t}^{t+T} P(x_i|(q_j, i - t + 1)) \quad (4.2)$$



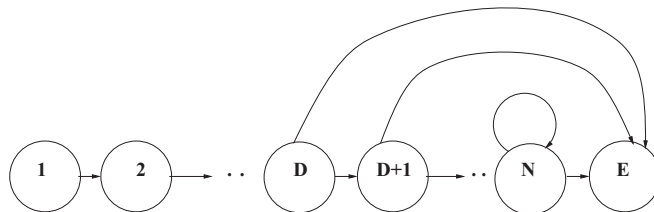


Figure 4.2: The topology of the time-index model

We can see that the  $q$ 's in the traditional HMM are replaced by a  $q$  and time index pair, as the state process is defined differently <sup>1</sup>. See also (Konig & Morgan 1994).

#### 4.3.4 An Example

Figure 4.2 shows the topology of a basic unit of speech. In this case phones are the basic speech units. Only the last state in the model has a self loop. For states with indices smaller than the minimum duration,  $D$ , for that phone, only a transition to the next state (corresponding to a time index increment of one) is permitted. For all other states, transitions are permitted either to the next state or to the exit state. This model differs from a traditional HMM (assuming a similar representation for duration) primarily in that now the emission probability for each state (for each time associated with a phone or subphone unit type) is not constrained to be equal. Specifically, the emission probability of a state in the Markov chain is  $P(x|(q_j, ti))$ , where  $ti$  is the time index. Similar equations could be used for multi-state HMMs that are also commonly used, in which the basic speech unit is smaller than a phone. While certainly one could define a standard HMM with the kind of model shown in Figure 4.2, and with a separate emission probability for each state, the basic problems are how to reliably estimate so many free parameters and to model the correlation between these these states. One possible solution is to share parameters between the estimates for the separate densities. Another solution would be to assume a parametric form for the trajectory, as was done by Deng. Reported here is a multi-layer perceptron (MLP) approach which, in our previous work at ICSI (as discussed in Section 2.3.3), has proved useful for such estimates (Bourlard & Morgan 1994).

<sup>1</sup>Formally the values of the state process are ordered pairs of the phone and the time-index.

### 4.3.5 An Implementation of the Time-Index Model

In our model we define the emission probability of a state as  $P(x|q_j, ti)$ . While such a quantity can always be defined, the important question is how to estimate it. Now consider the following decomposition according to Bayes' law:

$$\frac{P(x|q_j, ti)}{P(x)} = \frac{P(q_j|ti, x)P(ti|x)}{P(ti, q_j)} \quad (4.3)$$

Where  $ti$  is the value of the time index,  $x$  is the acoustic vector, and  $q_j$  is a specific phone. Alternatively, this can be decomposed as follows:

$$\frac{P(x|q_j, ti)}{P(x)} = \frac{P(ti|q_j, x)P(q_j|x)}{P(ti, q_j)} \quad (4.4)$$

Each of the terms conditioned on  $x$  can be estimated by an MLP with an acoustic vector (or a local neighborhood of acoustic vectors) as input, as well as any additional conditioning terms as input (for instance, an additional input representing time index  $ti$  in order to estimate  $P(q_j|ti, x)$ ). The targets correspond to a discrete binary coding of the class identity that is to the left of the condition bar (e.g.,  $q_j$  for estimating  $P(q_j|ti, x)$ , or  $ti$  for estimating  $P(ti|x)$ ). We have currently chosen to represent the  $ti$  inputs with a continuous-valued input as a smoother representation that requires fewer parameters. The first form of the equations given above requires the estimation of  $P(q_j|ti, x)$ , and this can be done with the MLP shown in Figure 4.3.

$P(ti, phone_j)$  can be estimated by counting the relative frequencies in the training set. Given that there is an accurate estimate of the boundaries between the phones we can calculate  $P(ti|x)$ ; otherwise an estimation of this probability is a difficult problem (Glass 1988).

### 4.3.6 Experiments

The Resource Management (RM) speaker independent task (Price *et al.* 1988) and the TIMIT database were used for initial experiments. In the RM experiments our training data consisted of 3990 read, continuous-speech sentences, and the 300-sentence Feb89 test set for network cross-validation and testing. The time-index net (as shown in Figure 4.3) had 1000 hidden units and 61 outputs (the size of phone set). There were 235 inputs to the net, including 234 that consisted of 9 frames of 26 features each (PLP12 + log gain + delta features for each of these 13) (Hermansky 1990), and a final time-index input. With

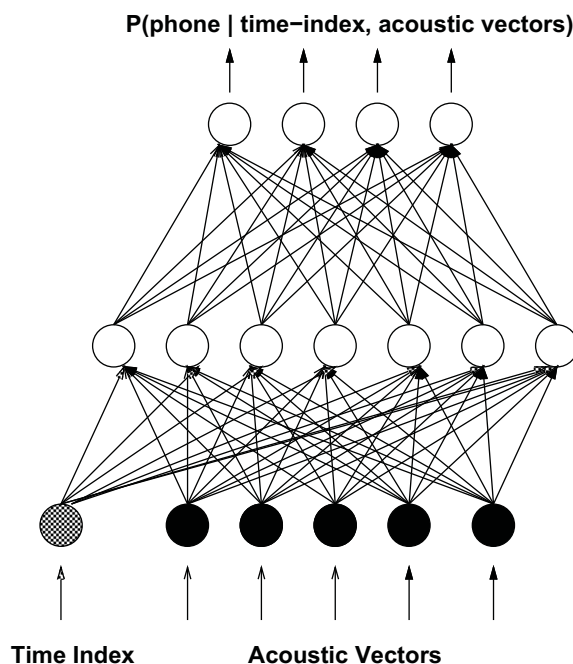


Figure 4.3: The time index net

the exception of this final input feature, this net was the same as the hybrid HMM/MLP system as described in (Boullard & Morgan 1994).

During training we can find the boundaries between phones by running a Viterbi alignment on the known word strings and from these boundaries to calculate the time index as the distance of the current frame from the beginning of the phone as marked by the previous boundary. For the preliminary tests the boundaries between the phones were determined by an automatic alignment (Viterbi) procedure on the known word string (Viterbi 1967), also during recognition. This side information about the word sequence was used only to generate boundaries and no explicit phonetic information was preserved. Obviously during realistic recognition experiments the identity of the spoken sentence is not known. Therefore, one can expect little improvement over the boundary detection found by the Viterbi procedure with a known word sequence. Hence, these initial time-index results serve as a lower bound on the error. The results are summarized in Table 4.3.6.

The TIMIT corpus was chosen for the second set of experiments because it is phonetically balanced, and in addition there are time-aligned phonetic transcriptions of all the sentences in the database. The goals were to verify the potential of the model on

<i>Task</i>	<i>Baseline</i>	<i>Time-Index + Boundary Information</i>
RM	4.8%	1.1%
TIMIT	36.4%	25.0%

Table 4.1: Time-Index Results with Known Boundaries

a different test set and also to answer a potential criticism that the reduction of error is due to restricting the recognizer to utterances with the same number of phones as in the answers. By using the boundaries from the Viterbi alignment on the known word strings, we restrict the potential answers to have the same number of boundaries as in the answers, i.e., the same number of phones.

The experiments were done on a 200-sentence development set that was selected from the official training set and which was not used for the training. The size of the nets and the features were the same as in the RM task experiments. We used 3300 sentences for training and 396 sentences for cross-validation (the 200 sentence development set is a subset of the cross-validation set). No language model was used in these experiments. All results are on the full 61 TIMIT phone set. The standard system had 36.4% phone errors on this task, while the incorporation of the knowledge of phoneme boundaries in the time-index network reduced the error to 25.0%. See also Table 4.3.6. When we restricted our standard system to sentences that have the same number of phones as in the known answers, the error rate was still 36.4%, but with a different mix of insertions, deletions, and substitutions. Hence, the crucial information was the timing of the boundaries and not their number.

These results suggest that, given good information about the phoneme boundary locations, recognition error rate can be greatly reduced. This was a necessary result for the transition-based approach to be ultimately useful, but it is certainly not sufficient. There is still the difficult problem of either specifically locating boundaries, or getting smooth estimates of them.

## Chapter 5

# Discriminant HMM (DHMM)

### 5.1 Introduction

Motivated by the experimental results that indicate the importance of the timing of transitions between basic speech units as described above, we chose to study the Discriminant HMM (DHMM) model, that uses local transition probabilities to estimate posterior probabilities of word sequences. Furthermore it is a true recognition model, i.e., it directly maps from acoustic sequences to sentences, unlike HMMs that model the inverse relation (the likelihood of producing an acoustic sequence).

The original theory of DHMM was described in (Bourlard & Wellekens 1989b). However, in the years following the original theoretical formulation, simplified systems to reduce the dependence on distributional assumptions for the observation space, and to make the probability estimates more discriminant. These simplified approaches did not make use of the full power of the initial scheme. Nonetheless, in controlled tests they displayed significant strength (Lubensky *et al.* 1994; Renals *et al.* 1994; Robinson *et al.* 1993). The basic scheme consisted of training neural networks to estimate probabilities of HMM states, and then using simple functions of these probabilities to label the training data using Viterbi decoding (dynamic programming). This procedure was repeated iteratively to train the system. During recognition the Viterbi procedure was used with probabilities produced by trained networks.

The remainder of this section will describe the original theory, but with the benefit of hindsight because of more recent developments. In addition, we describe experiments incorporating the original theory. An analysis of these experimental results suggests a

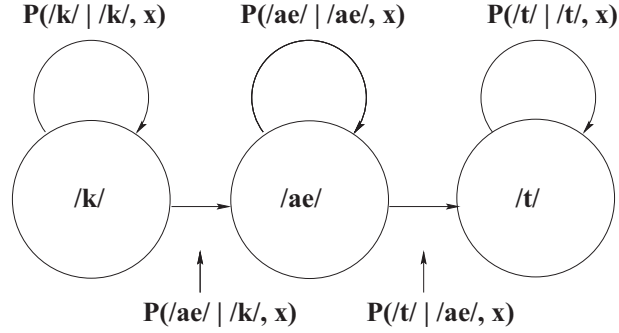


Figure 5.1: An example of a Discriminant HMM for the word “cat”. The variable  $x$  refers to a specific acoustic observation  $x_n$  at time  $n$ .

partial explanation why the original theory did not initially work. In the next chapter we describe a new training algorithm that overcomes some of the problems.

## 5.2 Estimation of the Posterior Probability of Word Sequences

In (Boullard & Morgan 1994) it was shown that it is possible to compute the global posterior probability  $P(M|X, L, \Theta)$  of (2.5) and (2.6) as:

$$P(M|X, L, \Theta) = \sum_{\Gamma_j} P(M, \Gamma_j|X, L, \Theta) \quad (5.1)$$

$$= \sum_{\Gamma_j} P(M|\Gamma_j, X, L, \Theta)P(\Gamma_j|X, L, \Theta) \quad (5.2)$$

in which “ $\Gamma_j$ ” represents a legal state sequence in  $M$ , see Section 2.1.3. Considering the second factor of (5.2) as the acoustic model and assuming that it is independent of the language model parameters, i.e.,:

$$P(\Gamma|X, L, \Theta) \approx P(\Gamma|X, \Theta) \quad (5.3)$$

and writing  $P(\Gamma|X, \Theta)$  explicitly as  $P(q^1, \dots, q^N|X, \Theta)$ , allows us to factor it as follows:

$$P(q^1, \dots, q^N|X, \Theta) = P(q^1|X, \Theta) P(q^2|X, q^1, \Theta) \dots \dots P(q^N|X, q^1, \dots, q^{N-1}, \Theta) \quad (5.4)$$

$$= \prod_{n=1}^N P(q^n|X, Q_1^{n-1}, \Theta) \quad (5.5)$$

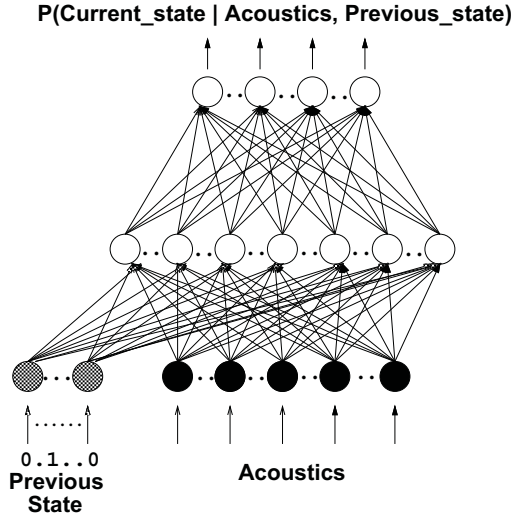


Figure 5.2: An example MLP that estimates local conditional probabilities.

where  $q^n$  represents the state observed at time  $n$  and  $Q_1^N$  the state sequence associated with  $X_1^N$ . Probabilities  $P(q^1, \dots, q^N | X, \Theta)$  can thus be calculated from “local” probabilities  $P(q^n | X, Q_1^{n-1}, \Theta)$ . These local probabilities may be simplified by relaxing the conditional constraints, by assuming dependency only on the previous state (first-order Markov model assumption) and on a temporal window  $X_{n-c}^{n+d}$  around the acoustic vector at time  $n$  (acoustic correlation limited to the contextual window). These local contributions can then be approximated by

$$P(q^n | X, q^1, \dots, q^{n-1}, \Theta) \approx P(q^n | X_{n-c}^{n+d}, q^{n-1}, \Theta) \quad (5.6)$$

where input contextual information is taken into account. These probabilities can be estimated at the outputs of an MLP with contextual input and output feedback, as described in Figure 5.2.

Thus, using Bayes’ law, we can then rewrite (5.4) for a particular path  $\Gamma_j$  as:

$$P(\Gamma_j | X, \Theta) \approx \prod_{n=1}^N P(q_{\Gamma_j}^n | q_{\Gamma_j}^{n-1}, X_{n-c}^{n+d}, \Theta) \quad (5.7)$$

A simple example of the model is given in Figure 5.2.

These new acoustic models, referred to as Discriminant HMMs (DHMM)<sup>1</sup>, are now described in terms of conditional transition probabilities  $P(q_n^\ell | q_{n-1}^k, x_n)$ , in which  $q_n^\ell$

<sup>1</sup>It could be argued that these models are no longer HMMs but more like “stochastic finite state acceptors”, a name suggested recently by John Bridle.

stands for the specific state  $q^\ell$  of  $Q$  hypothesized at time  $n$ . As with traditional hybrid HMM/ANN systems, conditional transition probabilities can be estimated by an ANN (in our case a multi-layer perceptron) with  $K$  output units and in which the acoustic input  $x_n$  is<sup>2</sup> complemented by a set of additional input units representing the state  $q^\ell$  hypothesized at the previous time step  $n - 1$ . The conditional transition probabilities are also functions of  $\Theta$ , the ANN parameter set, and can be written as  $P(q_n^\ell | q_{n-1}^k, x_n, \Theta)$ .

In estimating the first factor in (5.2)  $P(M|\Gamma_j, X, L, \Theta)$  given the state sequence and the language model we can assume no dependence on the acoustic sequence:

$$P(M|\Gamma_j, X, L, \Theta) \approx P(M|\Gamma_j, L, \Theta) \quad (5.8)$$

Thus, (5.2) can be rewritten as:

$$P(M|X, \Theta, L) \approx \sum_{\Gamma_j} P(\Gamma_j|X, \Theta) P(M|\Gamma_j, L, \Theta) \quad (5.9)$$

in which  $P(\Gamma_j|X, \Theta)$  is computed as described above.

### 5.2.1 Recognition

In the case that there are no homonyms, i.e., there is not a phone sequence that corresponds to more than one sentence, then the term  $P(M|\Gamma_j, L, \Theta)$  serves as a 0/1-valued filter. Equation (5.9) can then be rewritten as follows:

$$P(M_i|X, \Theta, L) \approx \sum_{\Gamma_j \in M_i} P(\Gamma_j|X, \Theta) \quad (5.10)$$

And it can be estimated with the following efficient forward recursion<sup>3</sup>.

#### Forward Recursion

We start with some definitions:

$$\alpha_n(k|M_j) = \sum_{\Gamma_n \in M_j} P(\Gamma_n, q_k^n | X_1^n, \Theta) \quad (5.11)$$

---

<sup>2</sup>As done with previous hybrid HMM/ANN systems,  $x_n$  will usually be replaced by  $X_{n-c}^{n+d} = \{x_{n-c}, \dots, x_n, \dots, x_{n+d}\}$  to take some acoustic context into account.

<sup>3</sup>Note this forward recursion will be different from the forward recursion that is used in training time and will be defined in the next section.



summing up all the paths of length  $n$  in model  $M_j$  that end in state  $q_k$  at time  $n$ . We initialize the recursion with:

$$\alpha_1(k|M_j) = P(q_k^1|x_1, \Theta) \quad (5.12)$$

for the first  $k$  valid states in model  $M_j$ . Now for the dynamic programming step:

$$\alpha_{n+1}(\ell|M_j) = \sum_{\Gamma_{n+1} \in M_j} P(\Gamma_{n+1}, q_\ell^{n+1} | X_1^{n+1}, \Theta) \quad (5.13)$$

$$= \sum_{\Gamma_n \in M_j} P(\Gamma_n | X_1^{n+1}, \Theta) P(q_\ell^{n+1} | X_1^{n+1}, \Gamma_n, \Theta)$$

assuming causality :  $P(\Gamma_n | X_1^{n+1}, \Theta) = P(\Gamma_n | X_1^n, \Theta)$ , denoting

the  $i$ th state of  $\Gamma_n$  by  $k$ , and assuming 1st order Markov model

$$= \sum_k \alpha_i(k|M_j) p(q_\ell^{n+1} | x_{n+1}, q_k^n, \Theta) \quad (5.14)$$

Thus, equation (5.10) can be computed using the forward recursion:

$$P(M_i | X, L, \Theta) \approx \sum_{q_f} \alpha_N(q_f | M_i) \quad (5.15)$$

where  $q_f$  is a legal final state for the model  $M_i$ . The use of the language model (in the case that we have homonyms) is discussed in Section 6.4.

### 5.3 MAP Constraints

In the previous section we described how to estimate the global posterior probability of a Markov model given the acoustic sequence. In this section we show that it can be done while meeting the constraint specified in Equation (2.7), i.e.,

$$\sum_{i=1}^I P(M_i | X, L, \Theta) = 1 \quad (5.16)$$

where the sum over  $i$  represents the sum over all possible Markov models (Bourlard *et al.* 1994). Here lies the difference between an ML and an MAP criterion. Any modification of the parameters of a model  $M_i$  must be complemented by a modification of all the parameters of the other models so as to preserve this constraint, hence making the MAP procedure discriminant.

It is also important to show that, in this case, if the “local” constraint:

$$\sum_{k=1}^K p(q_k^n | x_n, q_\ell^{n-1}, \Theta) = 1 \quad (5.17)$$

is met (which will be the case, at least approximately, with sigmoidal MLP outputs<sup>4</sup>), the constraint (2.7) on the global MAP probabilities is also met. Indeed, summing over the set of all possible paths  $\{q^1, \dots, q^n, \dots, q^N\}$  in all possible Markov models  $M_i$ , we have:

$$\begin{aligned}
\sum_i P(M_i|X, L, \Theta) &= \sum_{\Gamma_j} \sum_i P(M_i, \Gamma_j|X, L, \Theta) & (5.18) \\
&= \sum_{\Gamma_j} \sum_i P(M_i|X, \Gamma_j, L, \Theta) P(\Gamma_j|X, L, \Theta) \\
&\quad \text{Assuming (5.3) and (5.8)} \\
&= \sum_{\Gamma_j} P(\Gamma_j|X, \Theta) \sum_i P(M_i|\Gamma_j, L, \Theta) \\
&\quad \text{and, assuming } \sum_i P(M_i|\Gamma_j, L, \Theta) = 1, \forall \Gamma_j : \\
&= \sum_{\ell_1=1}^K P(q_{\ell_1}^1|x_1, \Theta) \dots \\
&\quad \dots \left( \sum_{\ell_2=1}^K P(q_{\ell_2}^2|x_2, q_{\ell_1}^1, \Theta) \dots \left( \sum_{\ell_N=1}^K P(q_{\ell_N}^N|x_N, q_{\ell_{N-1}}^{N-1}, \Theta) \right) \dots \right) \\
&= 1
\end{aligned}$$

It is however important to remember that this property is valid only if one considers all possible paths through the models.

Besides the advantage of forcing discrimination, numerical problems that plague the classical HMM are avoided when using discriminant models: namely, the lack of balance between the transition probability values (which only depend on the topology of the model) and the emission probability values (which decrease with the dimension of the input feature space) (Brown 1987).

## 5.4 Early Experiments and Error Analysis

### 5.4.1 Early Experiments

Given the theoretical properties of the Discriminant HMM/MLP model described earlier, we felt that empirical evaluations of this model would be a good first step in improving our understanding of transition-based systems. In particular, we began to empirically

---

<sup>4</sup>This constraint is precisely met in the case of a softmax output layer, since the outputs are normalized to sum to 1.

evaluate conditional transition probabilities as used in Discriminant HMM/MLP systems on phoneme classification and phonemic frame classification tasks.

As presented in the initial theory (Bourlard & Morgan 1994) the paradigm for training (and recognition) was to use the Viterbi approximation, i.e., to consider only the most probable state sequence in assigning phonetic labels to acoustic frames. The local discriminant probabilities (2.14) were estimated by an MLP as represented in Figure 5.2. In this case, the previous state is coded as additional binary inputs, one for each possible previous state. For every hypothesized previous state we set the corresponding input to one and the rest to zero. The set of possible previous states (or the set of possible successor states for a given  $q_k$  at the input) will be given by the topology of the HMMs (and by the currently hypothesized states of the matching process).

In Viterbi training (as used so far) we know the correct previous state (again considering only the most probable state sequence), either by having a hand-segmented database such as TIMIT, or by running an automatic alignment (Forced-Viterbi) on the training data. During recognition however, the MLP outputs will have to be hypothesized for every possible previous state (possibly constrained by a particular HMM topology or a language model).

The TIMIT corpus (Garofolo 1988) was chosen for the experiments because it is phonetically balanced and in addition there are time-aligned phonetic transcriptions of all the sentences in the database. The experiments were done on a 200-sentence development set that was selected from the official training set and was not used for the training. We used 3300 sentences for training and 396 sentences for cross-validation (where the 200 sentence development set is a subset of the cross-validation set). No language model was used in these experiments. All the results were on the full TIMIT 61 phone set. Phone models were simple one-state-per-phone models.

The net that estimated the local discriminant probabilities (as shown in Figure 5.2) had 1000 hidden units, 61 outputs (the size of the phone set). There were 295 inputs to the net, including 234 that consisted of 9 frames of 26 features each (PLP12 + log gain + delta features for each of these 13) (Hermansky 1990), and 61 binary inputs that represented the possible previous state. With the exception of these binary inputs, this net was the same as the hybrid HMM/MLP system as described in (Bourlard & Morgan 1994). The baseline HMM/MLP system (Bourlard & Morgan 1994) had 36.3% phone error on this task. When evaluating the Discriminant HMM on this task the error rate was 40.4%. This was

an intriguingly negative result; *increasing* the amount of input information led to a *decrease* in generalization performance. Why should this be so? Although it is difficult to draw firm conclusions from a negative result, it can at least inspire directions of inquiry. This result motivated the error analysis as described in the following section.

### 5.4.2 Error Analysis

As shown in the following, the error analysis suggested two potential reasons<sup>5</sup> for the observed performance loss: (1) poor transition detection, and (2) mismatch between the input space distribution of the MLP during training and recognition.

The first potential problem is missing transitions; implicitly the net is a transition detector because when it determines that the current state is different from the previous one it signals a transition, and transition detection between phonemes is known to be a hard problem (see (Glass 1988)). In order to test this assumption we compared the performance of the MLP described above on two kinds of acoustic frames: *transition frames* that start a new segment, i.e., their phonetic label is different from the previous frame, and all other frames, *self-loop frames*. *While presenting the correct previous state*, the frame level performance on the development set was:

1. Self-loop frames: 85.5% of correct phonemic classification.
2. Transition frames: 39.2% of correct transition detection and classification.

Transitions thus seems harder to detect and classify than “steady-state” frames. However, we suspect that this is not (only) due to the properties of transitions, but to two problems related to the training and testing procedure:

1. We have much less training data for transition frames than for “steady-state” frames (less than one-sixth). Thus the learning ability of the classifier will tend to focus on the steady-state phonetic classification.
2. Our training procedure assumes that a single frame is the transition and that its neighbors are not transitions<sup>6</sup>. This does not make sense in terms of the acoustic phonetics, since many spectral transitions are gradual. This makes a difficult classification function for a network to learn.

---

<sup>5</sup>Other than bugs.

<sup>6</sup>A transition frame is a frame whose target for the current state is different from the hypothesized previous state.

The second potential problem is the possible disparity between training and recognition input populations. During training we only present to the net “correct” pairs of acoustic vectors and the correct previous state, while in recognition we expect the net to generalize to all possible combinations of acoustic vectors and previous states. Some of these recognition inputs can be completely meaningless, e.g., like the combination of the acoustics of a middle of a vowel and a previous state that corresponds to a plosive. The net is not trained on anything close to these “impossible” pairs, but through the vagaries of interpolation these pairs could end up having the highest MLP outputs during recognition. This problem is often referred to as the “lack of negative training example” problem and sometimes can be partially overcome by presenting additional negative training examples to the net (Zavaliagos *et al.* 1994).

In order to test this hypothesis we computed the frame level performance of the net on the development set for the following two cases:

1. When the correct previous state is presented as input, the highest probability output was correct 79.4% of the time.
2. Presenting all possible previous states and taking as the winner the output with the highest activity, i.e., taking for every frame the maximum of 61 by 61 probabilities (61 outputs for each possible previous state), and checking if it was the correct pair of previous state and current state. In this case the result was 15.9% correct, which was the weighted average of 18.3% correct on self-loop frames and 0.4% correct on transition frames.

These results seem to suggest that, even for “steady-state” frames there is a problem of mismatch between the space of training and testing for hypothesized inputs. Of course, performance is also hurt by the problems mentioned earlier.

All the problems identified here motivated the REMAP training and recognition algorithm for HMM/MLP hybrids that is presented below. Specifically, the first experiment showed that “hard” transitions are difficult to detect. As we will see, the full MAP training will provide the nets with soft targets and soft decisions, i.e., with conditional probabilities of transitions. Furthermore, by considering all possible paths and transitions, we will reduce the mismatch between training and recognition. A formalism will be introduced that automatically considers negative training examples without the need for explicit enumeration of impossible input combinations.

## Chapter 6

# REMAP Training of HMM/MLP Hybrids

### 6.1 Introduction

#### 6.1.1 Motivations

The discriminant HMM/MLP theory as described above uses transition-based probabilities as the key building block for acoustic recognition. However, it is well known that estimating transitions accurately is a difficult problem (Glass 1988). Due to the inertia of the articulators, the boundaries between phones are blurred and overlapped in continuous speech (Deng & Sun 1994). It is also likely that some time variability exists in the human perception of the onset of a new phonetic region. Consequently, we would like to have a “window” of possible transitions instead of a single transition. Ideally the width of the transition window should depend on the specific bi-phone and on the speaker. Thus we need an automated way of estimating the transition windows to be used as targets in the MLP training.

In hybrid HMM/MLP systems, targets are typically obtained using an automatic alignment procedure incorporating a Viterbi approximation. However, this procedure yields rigid transition targets, and thus suffers from the problems mentioned above. Furthermore, our preliminary experiments with this procedure yielded poor transition detection performance.

Another related problem in our Viterbi-based MLP training procedure is a dispar-

ity between the training input space of the MLP and the input space used in recognition. Specifically, in training the network only processes inputs consisting of “correct” pairs of acoustic vectors and correct previous state, while in recognition we expect the net to generalize to all possible combinations of acoustic vectors and previous states.

One possible solution to these problems is to use a full MAP algorithm to find transition probabilities at each frame for all possible transitions with a forward-backward-like algorithm (Liporace 1982), which takes all possible paths into account. Furthermore a MAP algorithm would increase the a posteriori probability of the correct model and reduce the posterior probabilities of all other models. Thus, it might improve the approximation to an optimal Bayes classifier.

In the rest of this chapter, we describe a set of procedures that were derived in order to train and use the desired discriminant probability estimators in a full MAP framework.

### 6.1.2 Problem Formulation

Global MAP training of Discriminant HMMs should find the optimal acoustic parameter set  $\Theta$  maximizing (2.6). In the following derivation the dependency on the language model  $L$  is omitted, in order to concentrate on the acoustic model parameters (the role of the language model will be discussed in Section 6.4). Although, in principle, we could use a generalized back-propagation-like gradient descent in  $\Theta$  could maximize (2.6) (see, e.g., (Bengio *et al.* 1992)), an EM-like algorithm would have better convergence properties and would preserve the statistical interpretation of the ANN outputs. In this case, “full” MAP training of transition-based HMM/ANN hybrids requires a solution to the following problem: given a trained ANN at iteration  $t$  providing a parameter set  $\Theta^t$  and, consequently, estimates of  $P(q_n^\ell | x_n, q_{n-1}^k, \Theta^t)$ , how can we determine new ANN targets that:

1. Will be smooth estimates of conditional transition probabilities,  $\forall$  possible  $(k, \ell)$  state transition pairs in  $M$  and  $\forall n \in [1, N]$ .
2. When used in training the ANN for iteration  $t + 1$ , will lead to new estimates of  $\Theta^{t+1}$  and  $P(q_n^\ell | x_n, q_{n-1}^k, \Theta^{t+1})$  which are guaranteed to incrementally increase (2.6)?

## 6.2 Solution - REMAP

### 6.2.1 Introduction

In this section we describe a solution to the problem mentioned above. Specifically, we prove that for any training sentence  $X$ , an iteration consisting of

1. estimating new MLP training targets from a previously trained MLP via a backward recursion, and
2. training the MLP with the new targets

will increase the global MAP probability of the sentence model given the sequence of acoustic vectors<sup>1</sup>  $P(M|X)$ . It is easy to see that this proof can be generalized to several training sentences, since this is then simply equivalent to training on a long sentence built by concatenating all training sentences (with additional start and end point constraints). We describe below the main ideas and steps underlying the proof. Some of the technical details are discussed in Appendix A. The proof has three main steps:

1. Defining an auxiliary function such that maximizing this function is equivalent to maximizing the global posterior probability of the correct model and, since such probabilities must sum to 1 for the complete set of possible models, minimizing the posterior probabilities of the rival models.
2. Finding new targets for training the MLP that maximize the auxiliary function.
3. Showing that training the MLP with those new targets (using a weighted relative entropy error criterion) leads to an increase in the value of the auxiliary function.

Note that while this thesis has largely assumed the use of an MLP for the required probability estimation, other gradient-trained estimators (such as a recurrent network) could also be used.

---

<sup>1</sup>To simplify our notation, in all the following discussion only one training sentence  $X$  is associated with the Markov model  $M$ , but it is easy to see that the discussion remain valid in the case of several training sentences  $X_j$  associated with  $M_{w_j}$ , for  $j = 1, \dots, J$ .



### 6.2.2 Definitions

Let us define an auxiliary function<sup>2</sup>  $R(v_1, v_2)$  as a function on two probability sets  $v_1, v_2$ , where  $\Upsilon$  is defined below:

$$\Upsilon = \{P_z(q_\ell^n | x_n, q_k^{n-1}) : \forall k, \ell \in [1, \dots, K], \forall n \in [1, \dots, N], z \in \mathcal{Z}\} \quad (6.1)$$

Each set contains  $K^2(N - 1)$  possible conditional transition probabilities, where  $K$  is the number of states in the model, and  $N$  is the number of acoustic vectors, and  $\mathcal{Z}$  the number of possible legal sets. Note that the probability sets can be a function of a probability estimator. In our study, these probabilities are estimated by an MLP with parameter (weight) set  $\Theta$ , in which we denote the probabilities as  $P(q_\ell^n | x_n, q_k^{n-1}, \Theta)$ . In this case, the probability set  $v$  also becomes a function of  $\Theta$  and is then denoted by  $v(\Theta)$ .

The auxiliary function  $R(v_1, v_2)$  is defined as

$$R(v_1, v_2) = \frac{1}{P(M|X, v_1)} \sum_{\Gamma} P(M, \Gamma | X, v_1) \log P(M, \Gamma | X, v_2) \quad (6.2)$$

where  $\Gamma$  is a legal path (state sequence) in model  $M$  and  $P(M, \Gamma | X, v_i)$  represents the probability of a specific path  $\Gamma$  in  $M$  given a probability set  $v_i$ .

### 6.2.3 Theorem 1

**Theorem 1:** *IF  $R(v_1, v_2) \geq R(v_1, v_1)$*

*THEN  $P(M|X, v_2) \geq P(M|X, v_1)$ .*

In other words, if we can find a new set of probabilities  $v_2$  increasing  $R$ , the new set of probabilities will also increase the posterior probability of the model  $M$ . The proof is described in Appendix A.

### 6.2.4 Theorem 2

The question that arises from the first theorem is how to find a new set of probabilities  $v_2$  that increases the value of the auxiliary function  $R(\cdot)$  and, consequently, the posterior probability of the correct model (and therefore also minimizes the posterior probability of the rival models).

---

<sup>2</sup>This auxiliary function is usually denoted  $Q(\cdot)$  (Dempster *et al.* 1977). However, to avoid any possible ambiguity with an HMM state sequence, we denote it  $R(\cdot)$  in this thesis.

**Theorem 2:**

Given  $v_1$  a fixed set of probabilities that is estimated by an MLP with a fixed set of weights  $\Theta$ , we show that  $R(v_1, v_2)$  attains its maximum value when the conditional transitional probabilities  $P_{v_2}(q_\ell^n | x_n, q_k^{n-1}) \in v_2$  are defined as<sup>3</sup>

$$P_{v_2}(q_\ell^n | x_n, q_k^{n-1}) = P(q_\ell^n | X, q_k^{n-1}, v_1(\Theta), M) \quad (6.3)$$

Theorem 2 says the following: a trained MLP with a fixed set of parameters (MLP weights)  $\Theta$  provides us with estimates of conditional transition probabilities  $P(q_\ell^n | x_n, q_k^{n-1}, v_1(\Theta))$  (estimated on a given training data set  $X = \{x_1, \dots, x_n, \dots, x_N\} \forall n = 1, \dots, N$  and  $\forall k, \ell = 1, \dots, K$ ). Given these estimates, obtained at the output of the MLP, it is possible to compute re-estimates of the conditional transition probabilities  $P_{v_2}(q_\ell^n | x_n, q_k^{n-1}) = P(q_\ell^n | X, q_k^{n-1}, v_1(\Theta), M)$  by the backward recursion given in (6.14) to increase the global posterior probability of the correct model  $M$  over  $P(M|X, v_1)$ . The proof is described in Appendix A.

**6.2.5 Theorem 3**

As opposed to the “standard” EM algorithm (Baum *et al.* 1970; Baum 1972), Theorems 1 and 2 are not enough to prove convergence of the training process for two reasons:

1. The MLP training is usually minimizing a function (e.g., least mean square or relative entropy) that is different from the function optimized in Theorem 2. As a consequence, convergence must be proved through the same auxiliary function  $R(\cdot)$ .
2. Theorem 2 gives new (“optimal”) values (MLP targets) for the conditional transition probabilities which are going to be used to train the MLP. If the cost function can be trained to reach its optimal minimum, the MLP will just “learn” the targets and we will have

$$g_\ell(x_n, q_k^{n-1}, \Theta) = p(q_\ell^n | X, q_k^{n-1}, M, v_1) \quad (6.4)$$

which, by Theorem 2, is known to increase  $R(\cdot)$  and, consequently,  $P(M|X)$ . In this case, of course, we proved that MLP training is increasing  $P(M|X)$  and we do not

---

<sup>3</sup>Of course, all  $x_n$ 's in the following should be replaced by  $X_{n-c}^{n+d}$  if local contextual input is used, or  $X_1^n$  for a recurrent network.

need anything more. However, in general, the nets will not be trained to their optimal minimum because of

- “overlapping” of input patterns (e.g., two instances of the same pattern with two different classifications).
- use of cross-validation (early stopping) (Bourlard & Morgan 1994) to avoid overfitting and to get better estimates of actual probabilities.

Below we describe a training procedure for the MLP and a corresponding error criterion. We show that by minimizing this criterion we are maximizing the auxiliary function  $R(\cdot)$ . Thus by Theorem 1 the posterior probability of the correct sentence is increased. By this we show convergence (at least to a local minimum) on the training set.

Specifically, given a trained MLP with a set of weights  $\Theta^t$  which provides a set of conditional transition probabilities  $v_1(\Theta^t)$  and given a sequence of acoustic vectors  $X$  and a model  $M$ , the Discriminant HMM backward recursion can be used to compute a new set of probabilities

$$T = \{P_T(q_\ell^n | x_n, q_k^{n-1}) = P(q_\ell^n | X, q_k^{n-1}, v_1(\Theta^t), M); \forall k, l = 1, \dots, K; \forall n = 1, \dots, N\} \quad (6.5)$$

which will be used as targets to adapt further the MLP weights to a new set of weights  $\Theta^{t+1}$  and, consequently, a new set of conditional transition probabilities  $v_2(\Theta^{t+1})$ .

In the following we prove this property in the case of a weighted relative entropy  $E_e$ , similar to a common cost function for MLP training.<sup>4</sup> In this case, given a sequence of acoustic vectors  $X$  and a model  $M$  and the current set of parameters  $\Theta^t$ , the parameters  $\Theta^{t+1}$  of the MLP are trained to minimize

$$E_e(\Theta^{t+1}) = \mathcal{E}_{P(x_n, q_k^{n-1} | M, X, \Theta^t)} \sum_{\ell=1}^K P_T(q_\ell^n | x_n, q_k^{n-1}) \log \frac{P_T(q_\ell^n | x_n, q_k^{n-1})}{g_\ell(x_n, q_k^{n-1}, \Theta^{t+1})} \quad (6.6)$$

where  $g_\ell(x_n, q_k^{n-1}, \Theta^{t+1})$  is the  $\ell$ -th output of the MLP given weight set  $\Theta^{t+1}$  and inputs  $(x_n, q_k)$ . Note that the expected value  $\mathcal{E}_{P(x_n, q_k^{n-1} | M, X, \Theta^t)}$  is taken according to the distribution of the input variables that in the case of the Discriminant HMM are the concatenation

---

<sup>4</sup>Relative entropy is a particularly common error criterion for classification and probability estimation tasks, and we have used it for all of the speech training systems that we have developed over the last few years. The new criterion will actually only differ in that the expectation leading to its formulation will be taken with respect to the entire network input space, which includes a choice for the previous state.

of the acoustic input and the previous state. In this case (6.6) can also be expressed as:

$$\begin{aligned}
E_e(\Theta^{t+1}) &= \sum_{x_n, q_k^{n-1}} P(x_n, q_k^{n-1} | M, X, \Theta^t) \sum_{\ell=1}^K P_T(q_\ell^n | x_n, q_k^{n-1}) \log \frac{P_T(q_\ell^n | x_n, q_k^{n-1})}{g_\ell(x_n, q_k^{n-1}, \Theta^{t+1})} \\
&= \sum_{n=1}^N \sum_{k=1}^K P(q_k^{n-1} | M, X, x_n, \Theta^t) P(x_n | X, M, \Theta^t) \\
&\quad \left[ \sum_{\ell=1}^K P_T(q_\ell^n | x_n, q_k^{n-1}) \log \frac{P_T(q_\ell^n | x_n, q_k^{n-1})}{g_\ell(x_n, q_k^{n-1}, \Theta^{t+1})} \right] \quad (6.7)
\end{aligned}$$

Given that  $P(q_k^{n-1} | M, X, x_n, \Theta^t) = P(q_k^{n-1} | M, X, \Theta^t)$  and using the fact that  $P(x_n | X, M, \Theta^t) = 1$ ,

$$\begin{aligned}
E_e(\Theta^{t+1}) &= \sum_{n=1}^N \sum_{k=1}^K P(q_k^{n-1} | M, X, \Theta^t) \\
&\quad \left[ \sum_{\ell=1}^K P_T(q_\ell^n | x_n, q_k^{n-1}) \log \frac{P_T(q_\ell^n | x_n, q_k^{n-1})}{g_\ell(x_n, q_k^{n-1}, \Theta^{t+1})} \right] \\
&\quad \text{replacing } P_T(\cdot) \text{ with its definition (6.5) we get :} \\
&= \sum_{n=1}^N \sum_{k=1}^K P(q_k^{n-1} | M, X, \Theta^t) \\
&\quad \left[ \sum_{\ell=1}^K P(q_\ell^n | X, q_k^{n-1}, v_1(\Theta^t), M) \log \frac{P(q_\ell^n | X, q_k^{n-1}, v_1(\Theta^t), M)}{g_\ell(x_n, q_k^{n-1}, \Theta^{t+1})} \right] \quad (6.8)
\end{aligned}$$

It is easy to see that the above criterion will reach its global minimum when the outputs of the MLP are equal to the targets<sup>5</sup>. Note that the relative entropy between two probability mass functions is always greater or equal to zero (Cover & Joy 1991). Given that the targets are posterior probabilities, a network trained to the global minimum of error criterion (6.8) will estimate the posteriors.

An important point is that previous state  $q_k^{n-1}$  is not one of the features that are extracted from the speech waveform. Thus the scaling factor  $P(q_k^{n-1} | X, M, v_1(\Theta^t))$  is needed to compute the expected value over the “extended” input space. There are several ways to implement this scaling, one being to choose the previous state uniformly and to scale the error signal that is back propagated by this factor. An alternative way in stochastic gradient descent training (online training) is to implement this criterion by first choosing

---

<sup>5</sup>In the case that both the target probability and the net output are zero, this still holds given  $\lim_{\epsilon \rightarrow 0} \epsilon \log \frac{\epsilon}{\epsilon} = 0$

the acoustic frame  $x_n$  at random from the training test, and then choosing the previous state according to  $P(q_k^{n-1}|X, M, v_1(\Theta^t))$ .

**Theorem 3:**

Minimizing the weighted relative entropy criterion (6.6 ) with the target set  $T$  (which is calculated from a probability set  $v_1$ ) maximizes the auxiliary function  $R(\cdot)$ . Specifically, the new set of probabilities  $v_2$ , implemented by the trained MLP, satisfies the following:

$$E_e(\Theta^{t+1}) \leq E_e(\Theta^t) \implies R(v_1(\Theta^t), v_2(\Theta^{t+1})) \geq R(v_1(\Theta^t), v_1(\Theta^t)) \quad (6.9)$$

By proving Theorem 3 as described in Appendix A, we show that minimizing error criterion (6.8) is equivalent (within a scaling factor) to maximizing the auxiliary function. In combination with the previous Theorems, this confirms that a network trained using error criterion (6.8) and targets defined by Theorem 2 will increase the auxiliary function. This in turn means that the global probability of the correct model will be increased. In practice the change of the error measure on a cross-validation set is used to guide the training schedule of the MLP, e.g., for deciding the learning rate and the stopping point.

### 6.2.6 Summary and Discussion

Like the EM algorithm, REMAP training consists of two major steps: estimation (which in this case is estimating new targets for the MLP), and maximization (which here consists of adapting the MLP weights to maximize performance on the new set of targets). Here we have proved three theorems that together show that each iteration of REMAP training increases the posterior probability of the training sentence. It is assumed that the training set is a good sample of the overall input space, and cross-validation techniques are used to check that we have not over-fit to the training data, e.g., by computing the change of the posterior probability on an independent set after every iteration of the REMAP algorithm. In principle, REMAP should ultimately provide improved recognition accuracy for practical systems. However, as with all other gradient-based optimization techniques, we will be vulnerable to potential difficulties with local minima.

## 6.3 REMAP Training

### 6.3.1 Introduction

Since it is now well-known (Bourlard & Morgan 1994) how to train an MLP to lead to good estimates of posterior probabilities (whether the MLP targets are “1-from-K” binary vector or themselves estimates of posterior probabilities), the remaining problem is to find an efficient algorithm to express  $P(q_\ell^n | X, q_k^{n-1}, M)$  in terms of  $P(q_\ell^n | X_{n-c}^{n+d}, q_k^{n-1})$  so that the next iteration’s targets can be found, and also how to estimate  $P(q_k^{n-1} | X, M)$  in order to select the previous state in an adequate way. Note that from now on, to simplify notation, and only when there is no risk of confusion, we will drop the indices of  $M$  and  $X$ , keeping in mind that  $M$  will represent either a specific model associated with a specific training sentence  $X$  or one of the many possible hypothesized models during recognition of  $X$ .

### 6.3.2 Target Estimation

In this section we describe dynamic programming procedures for efficient estimation of the MLP targets. By simple statistical rules the desired targets may be expressed as:

$$P(q_\ell^n | X, q_k^{n-1}, M) = \frac{p(X, q_k^{n-1}, q_\ell^n, M)}{p(X, q_k^{n-1}, M)} \quad (6.10)$$

$$= \frac{p(X, q_k^{n-1}, q_\ell^n, M)}{\sum_\ell p(X, q_k^{n-1}, q_\ell^n, M)} \quad (6.11)$$

in which

$$p(X, q_k^{n-1}, q_\ell^n, M) = \sum_i p(X, \Gamma_i, q_k^{n-1}, q_\ell^n, M) \quad (6.12)$$

is equal to the sum of the probabilities of all possible paths  $\Gamma_i$  in a particular  $M$  visiting  $q_k$  at time  $n-1$  and  $q_\ell$  at time  $n$ . In a similar way, the denominator represents the sum of the probabilities of all possible paths in  $M$  visiting  $q_k$  at time  $n-1$ . Since we cannot afford to compute and memorize all possible paths in (6.12), to compute  $P(q_\ell^n | X, q_k^{n-1}, M)$  we need to find efficient recursions in terms of local probabilities  $P(q_\ell^n | X_{n-c}^{n+d}, q_k^{n-1})$  generated by the MLP. Denoting

- $\beta_n(k, \ell | M) = p(X_n^N, q_\ell^n | X_1^{n-1}, q_k^{n-1}, M)$  as the “backward” probability, defined as the probability of observing the rest of the sequence and starting from state  $q_\ell$  at time

$n$  given the previously observed  $X_1^{n-1}$ , and given state occupancy of class  $q_k$  at time  $n - 1$ .

$\beta_n(k, \ell|M)$  is explicitly written with a conditional on  $M$  in the parentheses to remind us that the backward recurrences will be run through a specific Markov model  $M$  with a specific topology.

### Backward Recursion

By using simple statistical rules, we have:

$$\begin{aligned}
\beta_n(j, k|M) &= p(X_n^N, q_k^n | X_1^{n-1}, q_j^{n-1}, M) \\
&= \sum_{\ell} p(x_n, X_{n+1}^N, q_k^n, q_{\ell}^{n+1} | X_1^{n-1}, q_j^{n-1}, M) \\
&= \sum_{\ell} p(X_{n+1}^N, q_{\ell}^{n+1} | X_1^n, q_j^{n-1}, q_k^n, M) p(x_n, q_k^n | X_1^{n-1}, q_j^{n-1}, M) \\
&\quad \text{assuming 1st order HMMs :} \\
&\simeq \sum_{\ell} \beta_{n+1}(k, \ell|M) P(q_k^n | X_1^n, q_j^{n-1}, M) p(x_n | X_1^{n-1}, q_j^{n-1}, M) \quad (6.13)
\end{aligned}$$

As with standard HMMs, the conditional of “local” probabilities  $P(q_{\ell}^{n+1} | X_1^{n+1}, q_k^n, M)$  will be restricted to a contextual acoustic window  $X_{n-c}^{n+d}$  and the local posterior (as generated by the neural network) will be assumed independent of  $M$ , leading to the following backward recursion:

$$\beta_n(j, k|M) = p(q_k^n | X_{n-c}^{n+d}, q_j^{n-1}) c_n(j) \sum_{\ell} \beta_{n+1}(k, \ell|M) \quad (6.14)$$

where  $c_n(j)$  denotes  $p(x_n | X_1^{n-1}, q_j^{n-1}, M)$ .

Note that the last simplification, i.e.,  $P(q_k^n | X_1^n, q_j^{n-1}) \simeq p(q_k^n | X_{n-c}^{n+d}, q_j^{n-1})$  is not necessary in the case of a recurrent net for estimation of the conditional transition probabilities. However, in practice it might not make a difference if an MLP is used to estimate the conditional transition probabilities, i.e., to do the last simplification. Initialization of this backward recursion can be done according to

$$\beta_{N+1}(L, F|M) = p(X_{N+1}^N, q_F^{N+1} | X_1^N, q_L^N, M) = 1 \quad (6.15)$$

in which  $q_F$  is the (non-emitting) final state and  $q_L$  represents any last emitting HMM state (associated with  $M$ , as imposed by the conditional).

### MLP Output Targets Update

Equation (6.10) may now be expressed in terms of the backward recursion:

$$\begin{aligned}
P(q_k^n | X, q_j^{n-1}, M) &= \frac{p(X, q_j^{n-1}, q_k^n | M)}{p(X, q_j^{n-1} | M)} \\
&\text{dividing numerator and denominator by } P(X_1^{n-1}, q_j^{n-1} | M) \\
&= \frac{\beta_n(j, k | M)}{\sum_{\ell} \beta_n(j, \ell | M)} \\
&\text{following (6.13)} \\
&= \frac{\sum_{\ell} \beta_{n+1}(k, \ell | M) P(q_k^n | X_{n-c}^{n+d}, q_j^{n-1}) c_n(j | M)}{\sum_h \sum_{\ell} \beta_{n+1}(h, \ell | M) P(q_h^n | X_{n-c}^{n+d}, q_j^{n-1}) c_n(j | M)} \\
&= \frac{\sum_{\ell} \beta_{n+1}(k, \ell | M) P(q_k^n | X_{n-c}^{n+d}, q_j^{n-1})}{\sum_h \sum_{\ell} \beta_{n+1}(h, \ell | M) P(q_h^n | X_{n-c}^{n+d}, q_j^{n-1})} \tag{6.16}
\end{aligned}$$

This final form of the equation shows that the probabilities required to determine MLP targets can be obtained from the previous MLP outputs and the beta recursions alone.

### 6.3.3 Forward Recursion

In order to efficiently calculate the posterior probability  $P(q_k^{n-1} | X, M)$  we define a forward recursion<sup>6</sup>. We start with some definitions:

$$\alpha_n(k | M_j) = P(q_k^n | X_1^n, M_j, \Theta) \tag{6.17}$$

i.e., the probability of being at  $q_k$  at time  $n$  given the acoustic vectors from the beginning of the utterance until time  $n$ . Also

$$\alpha_1(k | M_j) = P(q_k^1 | x_1, M, \Theta) \tag{6.18}$$

Now for the recursion step:

$$\begin{aligned}
\alpha_{n+1}(\ell | M_j) &= \sum_{\Gamma_{n+1}} P(\Gamma_{n+1}, q_{\ell}^{n+1} | X_1^{n+1}, M, \Theta) \\
&= \sum_{\Gamma_n} P(\Gamma_n | X_1^{n+1}, M, \Theta) P(q_{\ell}^{n+1} | X_1^{n+1}, \Gamma_n, M, \Theta) \\
&\text{assuming causality : } P(\Gamma_n | X_1^{n+1}, M, \Theta) = P(\Gamma_n | X_1^n, M, \Theta)), \text{ denoting} \\
&\text{the } i\text{th state of } \Gamma_n \text{ by } k, \text{ and assuming 1st order Markov model} \\
&= \sum_k \alpha_i(k | M_j) P(q_{\ell}^{n+1} | x_{n+1}, q_k^n, M, \Theta)
\end{aligned} \tag{6.19}$$

---

<sup>6</sup>The forward recursion as described here is used at training time, i.e., when we know the correct model.



where the transition from  $q_k$  to  $q_\ell$  is legal in model  $M_j$ . When the local probabilities are “tied” across models the recursion becomes:

$$\alpha_{n+1}(\ell|M_j) = \sum_k \alpha_i(k|M_j)P(q_\ell^{n+1}|x_{n+1}, q_k^n, \Theta) \quad (6.20)$$

### 6.3.4 Estimating the Previous State Distribution

We can compute  $P(q_k^{n-1}|X, \Theta^t)$ , the posterior probability of being in class  $q_k$  at time  $n - 1$  used in the training of the MLP as specified below, according to the following:

$$P(q_k^{n-1}|X, M, \Theta^t) = \frac{P(q_k^{n-1}, X|M, \Theta^t)}{P(X|M, \Theta^t)} \quad (6.21)$$

$$\begin{aligned} & \text{dividing numerator and denominator by } P(X_1^{n-1}|M, \Theta) \\ &= \frac{\alpha_{n-1}(k|M) \sum_\ell \beta_n(k, \ell|M)}{\sum_k [\alpha_{n-1}(k|M) \sum_\ell \beta_n(k, \ell|M)]} \end{aligned} \quad (6.22)$$

### 6.3.5 REMAP Training Algorithm

The general scheme of the MAP Forward-Backward training of hybrid HMM/MLP systems can be summarized as follow:

1. Start from some initial net providing  $P(q_\ell^n|X_{n-c}^{n+d}, q_k^{n-1}, \Theta^t)$ ,  $t = 0$ , for all possible  $(k, \ell)$ -pairs<sup>7</sup>.
2. Run backward recurrences to compute MLP targets  $P(q_\ell^n|X_j, M_j, q_k^{n-1}, \Theta^t)$  according to (6.16), for all training sentences  $X_j$  associated with HMM  $M_j$ , for all possible  $(k, \ell)$  state transition pairs in  $M_j$  and for all  $x_n$ ,  $n = 1, \dots, N$  in  $X_j$  (see next point). Also as part of the forward and backward recurrences we compute  $P(q_k^{n-1}|X, \Theta^t)$  (6.21), the posterior probability of being in class  $q_k$  at time  $n - 1$ , conditioned on the acoustic vector sequence  $X$ , to be used in the training of the MLP as specified below.
3. For every  $x_n$  (or  $X_{n-c}^{n+d}$ ) in the training set choose  $q_k^{n-1}$  according to  $P(q_k^{n-1}|X, \Theta^t)$ , train the MLP to minimize the relative entropy between the network outputs and targets which equal to  $P(q_\ell^n|X, q_k^{n-1}, \Theta^t)$ . (See Appendix A for a more theoretical explanation.) This provides us with a new set of parameters  $\Theta^t$ , for  $t = t + 1$ .
4. Iterate from 2 until convergence.

---

<sup>7</sup>This can be done, for instance, by training up such a net from a hand-labeled database like TIMIT or from some initial forward-backward estimator of equivalent local probabilities (usually referred to as “gamma” probabilities in the Baum-Welch procedure).

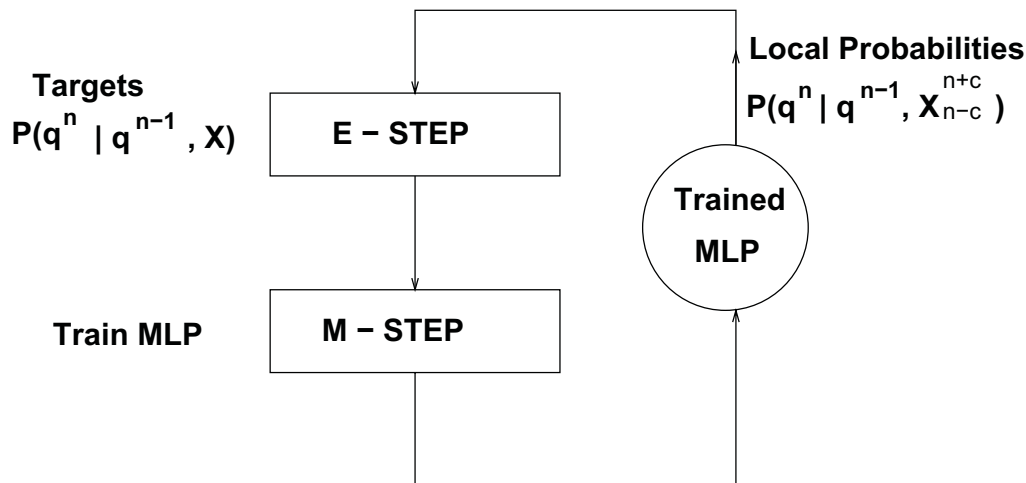


Figure 6.1: An illustration of REMAP

This procedure is composed of two steps: an Estimation (E) step, corresponding to step 2 above, and a Maximization (M) step, corresponding to step 3 above. In this regards, it is reminiscent of the Estimation-Maximization (EM) algorithm as discussed in (Dempster *et al.* 1977). The algorithm is illustrated in Figure 6.1. However, in the standard EM algorithm, the M step involves the actual maximization of the likelihood function. In a related approach, usually referred to as the Generalized EM (GEM) algorithm, the M step does not actually maximize the likelihood, but simply increases it (by using, e.g., a gradient procedure). Similarly, REMAP increases the global posterior function during the M step (in the direction of targets that actually maximize that global function), rather than actually maximizing it. Recently, a similar approach was suggested for mapping input sequences to output sequences (Bengio & Frasconi 1995).

### 6.3.6 Complexity Issues

It is important the computational cost of REMAP, particularly in comparison to the standard approach of training hybrid HMM/MLP systems (Bourlard & Morgan 1994). The computation of the MLP targets and the distribution of the previous state is done by running the forward and backward recurrences as described above. Denoting the number of states in a particular model  $M_i$  by  $K$  and the number of frames in a particular training sentence associated with  $M_i$  by  $N$ , the complexity of running the recurrences is  $O(KN)$ . This is explained by the fact that the number of operations per state per frame is a constant,

assuming a bounded branching factor of the directed graph that represents the topology of<sup>8</sup>  $M_i$ . This is the same computational cost as running a forced alignment (Viterbi) to get targets for the MLP training as done in the standard HMM/MLP systems. The difference, though, is that to obtain the conditional transition probabilities that are the input for the recurrences, we need to perform more computations in the REMAP case. Specifically, we have to present every acoustic frame once for every possible previous state<sup>9</sup>. Hence, we have to do  $K$  feed-forward runs for every acoustic frame. This is in contrast to the standard HMM/MLP system training, which requires presenting each acoustic frame only once.

During training, for each acoustic frame, the previous state is selected by random according to the previous state distribution as computed by the forward and backward recurrences (in the current implementation). Thus, due to the extra input which represents the previous state, there are up to  $K$  times more training patterns. In practice most of the previous states have very low probability ( $< 0.05$ ) and can be ignored, resulting in only two to three times more patterns. Therefore, the overall training time is two to three times longer than the standard HMM/MLP training, which has the standard back-propagation complexity of  $O(TW)$ , where  $T$  is the number of acoustic frames, and  $W$  is the number of the MLP weights.

## 6.4 The Role of the Language Model

So far, in the description of the suggested training algorithm, we have omitted the language model. In this section we describe some initial ideas about incorporating the language model into the REMAP framework. Current techniques of speech modeling assume a separation between the language model parameters and the acoustic model parameters as described in Section 3.1. Usually the language model estimates prior probabilities of sentences by computing n-grams (the probability of a word given the previous (n-1) words), by counting the relative frequencies in the training set or in a large text corpus. Acoustic models use a separate set of parameters, which are estimated independently from the language model parameters. In the posterior-based framework proposed here, the posterior probability of the sentence should be estimated given *both* the acoustics and high-level knowledge such as a language model (see Chapter 2). Furthermore, at training time we

---

<sup>8</sup>This is a reasonable assumption given that the transitions are strictly from left to right in our models.

<sup>9</sup>This computation cost could be reduced by pruning based on the probability of each previous state

want to optimize the measure that we use in recognition. Therefore, it would be desirable to maximize the posterior probability of the correct model given both the acoustics and the language model. However, the training algorithm that is proposed in this study, REMAP, increases the posterior probability of the correct model based only on acoustic information.

An alternative to the separation assumption is a unified model with one parameter set, e.g., a neural network for isolated word recognition, with one output unit for each possible word. It is not clear however, how to avoid the separation assumption in modeling continuous speech recognition, given the unavoidable hierarchical modeling. In hierarchical modeling the larger units (sentences) are built from smaller units (words and phone-like units) that are shared across the larger representation. Thus, there is usually one model for the relation between phones and acoustic vectors (e.g., emission and transition probabilities), one for the relation between phones and words (e.g., pronunciation models), and one for the relation between words and sentences (e.g., language model). Consequently, instead of directly estimating  $P(M|X, L, \Theta)$ , we decompose it into terms which can be estimated separately. Currently, the model  $M$  represents a sequence of words, i.e., there is one-to-one mapping between sentences and models<sup>10</sup>. Thus, going one level down in the modeling hierarchy, the model may be described in terms of state sequences:

$$P(M|X, L, \Theta) = \sum_{\Gamma_i} P(\Gamma_i|X, L, \Theta)P(M|\Gamma_i, X, L, \Theta) \quad (6.23)$$

Although acoustic information such as fundamental frequency, silences, and energy levels, might help to estimate the word sequence from the state sequence, our algorithm currently does not utilize it. Thus, assuming that given the state sequence and the language model, the dependence on the acoustic sequence can be dropped, we get:

$$P(M|X, L, \Theta) \approx \sum_{\Gamma_i} P(\Gamma_i|X, L, \Theta)P(M|\Gamma_i, L, \Theta) \quad (6.24)$$

Using Bayes' rule, we also have:

$$P(M|X, \Theta, L) \approx \sum_{\Gamma_i} P(\Gamma_i|X, \Theta, L) \frac{P(\Gamma_i|M, L, \Theta)P(M|L, \Theta)}{P(\Gamma_i|L, \Theta)}$$

Approximating that the effect of the language model can be ignored in the acoustic term<sup>11</sup>,

---

<sup>10</sup>In contrast to *speech understanding* tasks where the model may represent an action (e.g., dialing a particular phone number) that can be expressed by several sentences.

<sup>11</sup>This approximation is especially inaccurate at word boundaries.

i.e.,  $P(\Gamma_i|X, L, \Theta) \approx P(\Gamma_i|X, \Theta)$ , yields:

$$P(M|X, L, \Theta) \approx \sum_{\Gamma_i} P(\Gamma_i|X, \Theta) \frac{P(\Gamma_i|M, L, \Theta)P(M|L, \Theta)}{P(\Gamma_i|L, \Theta)} \quad (6.25)$$

in which  $P(\Gamma_i|X, \Theta)$  is computed as described in Chapter 5.  $P(M|L, \Theta)$  can be assumed independent of the acoustic model parameters and can be estimated using standard language modeling techniques. In principle  $P(\Gamma_i|M, L, \Theta)$  and  $P(\Gamma_i|L, \Theta)$  can be estimated during training by dynamic programming techniques similar to our  $\alpha$  and  $\beta$  recurrences (Bourlard *et al.* 1994), and the ratio of these two terms represents the additional state transition information that is gained by knowing the specific word sequence.

## Chapter 7

# Experimental Results with REMAP

### 7.1 Isolated Speech Experiments

In this section we report on experiments with isolated speech, where recognition was based on acoustic information. The initial task was isolated speech recognition task on the Digits+ corpus in use at ICSI, which is a subset of a larger database recorded over a clean telephone line at Bellcore. It is composed of 200 speakers saying the words “zero” through “nine”, “oh”, “no”, and “yes”. The additive noise in these experiments was automotive sound that was recorded over a cellular telephone. Noise was randomly selected from this source and then added to the clean speech waveforms (10 dB S/N ratio). In order to better utilize the data we used a jackknife procedure (Efron 1982). For each of four experiments, three fourths of the data was used for training and cross-validation, and one fourth was used for testing. Specifically, in each experiment there were 1720 utterances for training, 230 for cross-validation and 650 (from 50 speakers) for testing. All nets had 214 inputs: 153 inputs for the acoustic features, and 61 to represent the previous state (one unit for every possible previous state). The acoustic features were combined from 9 frames with 17 features each (RASTA-PLP8 + delta features + delta log gain) computed with an analysis window of 25 ms computed every 12.5 ms (overlapping windows) and the sampling rate was 8 Khz. The nets had 200 hidden units and 61 outputs. All the initial weights for the nets were from a net that was trained on the NTIMIT database. The combined

System	cut 1	cut 2	cut 3	cut 4	Overall Error	Avg. Posterior
Baseline Hybrid	3.5%	2.3%	4.4%	3.4%	3.4%	-
DHMM, pre-REMAP	2.5%	2.6%	2.9%	2.5%	2.7%	0.1269
1 REMAP iteration	2.8%	2.0%	2.5%	2.8%	2.5%	0.1731
2 REMAP iterations	2.9%	2.0%	2.5%	2.6%	2.5%	0.1773

Table 7.1: Training and testing on noisy isolated digits.

results for all four cuts are summarized in Table 7.1. Note that the row entitled “Baseline Hybrid” refers to an ANN trained on targets of 1’s and 0’s that were obtained from a forced Viterbi procedure by our standard HMM/ANN system as described in (Boumlard & Morgan 1994); the row entitled “DHMM, pre-REMAP” means a Discriminant HMM using the same training approach, with hard targets determined by the first system, and additional inputs to represent the previous state. The rightmost column gives the average probability of the correct model over all test words as determined during recognition. The recognition rate after the first and second iterations of REMAP is significantly better (at  $p < 0.05$  level) than the baseline hybrid system. Although the contribution of the REMAP step is small for this task, The overall improvement including the effect of using the transition-based, posterior framework as done in the Discriminant HMM, is significant. In figure 7.1 we illustrate the effect of REMAP iterations of the probability of transition, i.e., changing state for every frame in the utterance “one.” Specifically, we calculated the probability that the current state is different from the previous state. As a result of REMAP iterations we get smoother transition probabilities as desired.

## 7.2 Continuous Speech Experiments

The next step was to test whether this improved performance can also be obtained with continuous speech. For this purpose we chose the Numbers’93 corpus. It is a continuous-speech database collected by CSLU at the Oregon Graduate Institute. It consists of numbers spoken naturally over telephone lines on the public-switched network (Cole *et al.* 1994). The Numbers’93 database consists of 2167 speech files of spoken numbers produced by 1132 callers. We used 877 of these utterances for training and 657 for cross-validation and testing (200 for cross-validation). There are 36 words in the vocabulary, namely *zero, oh, 1,*

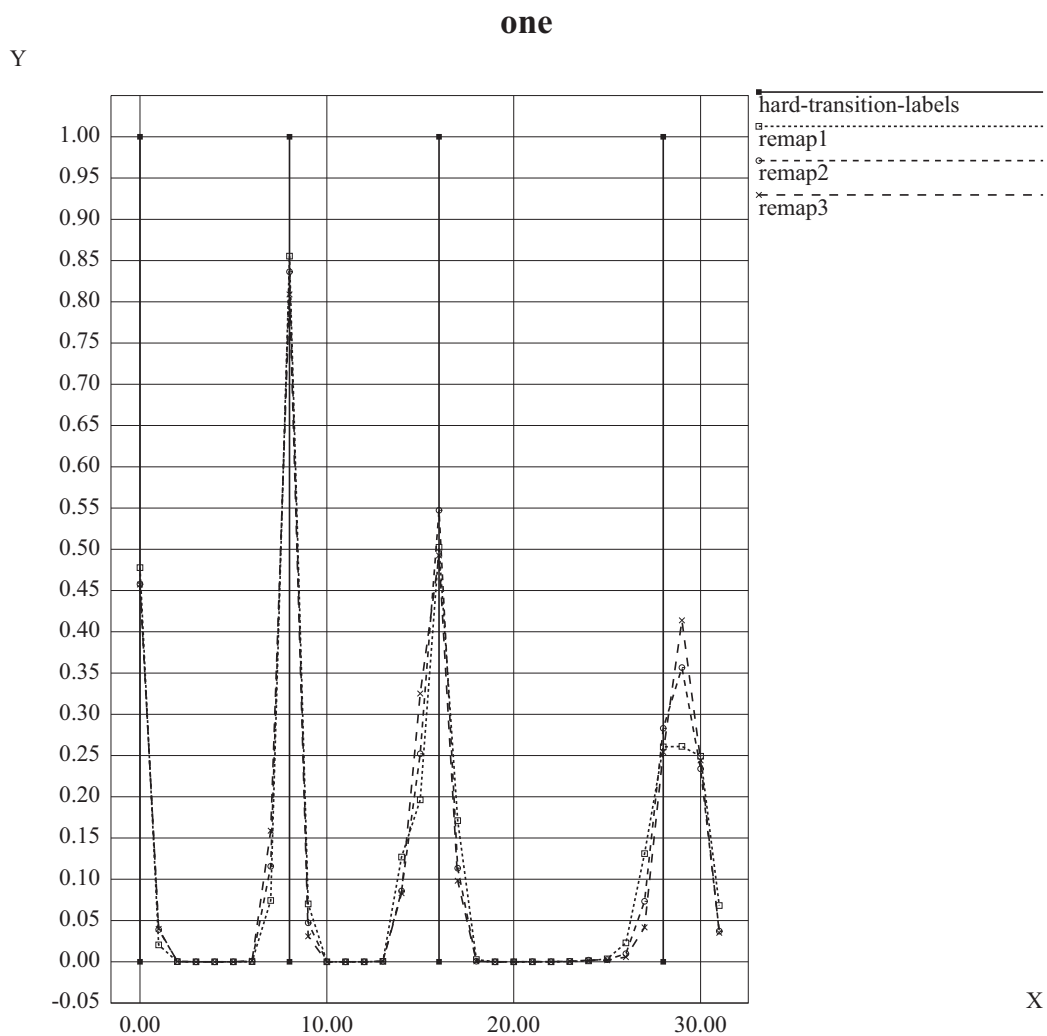


Figure 7.1: The probability of a transition (changing state) for every frame in the utterance “one.” The Y-axis represents the probability of transition, and the X-axis the time step



System	Error Rate
DHMM, pre-REMAP	14.9%
1 REMAP iteration	13.6%
2 REMAP iterations	13.2%

Table 7.2: Training and testing on continuous numbers, no syntax, no durational models.

*2, 3,...,20, 30, 40, 50,...,100, 1000, a, and, dash, hyphen, and double.* As before, all the nets have 214 inputs: 153 inputs for the acoustic features, and 61 to represent the previous state (one unit for every possible previous state). All the initial weights for the nets were from a net that was trained on the NTIMIT database. The acoustic features were combined from 9 frames with 17 features each (RASTA-PLP8 + delta features + delta log gain) computed with an analysis window of 25 ms computed every 12.5 ms (overlapping windows) and the sampling rate was 8 Khz. The nets have 200 hidden units and 61 outputs. Our results are summarized in Table 7.2. Note that the row entitled “DHMM, pre-REMAP” means a Discriminant HMM using our standard training approach, with hard targets determined by the first system, and additional inputs to represent the previous state.

The improvement in the recognition rate as a result of REMAP iterations was significant at  $p < 0.05$ . However, all the experiments were done using acoustic information alone. Using our (baseline) hybrid system under equal conditions, i.e., no duration information and no language information, yielded 31.6% word error; adding the duration information back reduced the error rate to 12.4%.

## 7.3 Analysis and Discussion

### 7.3.1 Invalid State Sequences

The improvement in the recognition performance as a result of REMAP training was significant for the continuous speech experiment and not significant for the isolated speech experiment. Most of the improvement in the isolated speech task came from using the posterior-based model, DHMM. This is an intriguing result, especially given the significant increase in the posterior probabilities of the correct models in the DIGITS+ experiment. One explanation for this might be related to phone sequences that do not represent any

valid sentence.

Essentially, REMAP discriminates between phone sequences (through the state sequences). In Section 5.3 we show that the MAP constraint is met by the DHMM model, i.e., the sum of all posterior probabilities of all state sequences is 1. However, some of these state sequences<sup>1</sup> are not valid in any particular model. They can be viewed as part of a big “garbage model” that includes all invalid state sequences. Further, the sum of the posterior probability of this “garbage model” and the posterior probabilities of all valid models is one. The question is how big is the share of the “garbage model” ? This share could be evaluated for the DIGITS+ case (given that we only have 13 legal models). As it turns out, the average posterior probability of the “garbage model” before the REMAP iterations is 0.87, and REMAP reduced it to 0.82. Therefore, it might be the case that most of the increase in the average posterior probability of the correct models was at the expense of the “garbage model” and not at the expense of other valid rival models. We tested this hypothesis by measuring the arithmetic average of the ratio of the posterior probability of the correct model and the sum of the posterior probabilities of all valid models, before and after REMAP. Before REMAP this ratio was 0.971, and REMAP increased it to 0.975. Explicitly, the average posterior probability of the correct models before REMAP was 0.1269 and the posterior probability of all other incorrect yet feasible models before REMAP sum to 0.00379, and REMAP increased them to 0.1773 and 0.00455 respectively. Therefore, REMAP only slightly increased the posterior probability of the correct model relative to the other valid models. In other words, for this task the model that had lower probabilities were primarily “illegal” models. This might be more useful in a word spotting task, where illegal sequences may be input.

The same measures could not be calculated for the NUMBERS’93 task given the very large number of valid models (continuous speech). However, the share of the probability mass dedicated to the “garbage model” may not be very big. The reason is that, given more valid models, more state sequences are encompassed. Therefore, there are fewer state sequences left for the “garbage model.” This might explain the significant improvement in the continuous speech experiment. Hopefully, for large tasks (in terms of vocabulary and perplexity) the valid models will cover more possible sound sequences, resulting in a less influential “garbage model.”

---

<sup>1</sup>Each state sequence represents one phone sequence, as each state belongs to a particular phone.

### 7.3.2 Hard Targets

#### Number of Transitions

In chapter 5 we observed that for a phone classification task training DHMM with “hard” targets resulted in worse performance than for the baseline system. In contrast, on both the DIGITS+ and the NUMBERS’93 experiments, we got better results with DHMM trained on “hard” targets than with the baseline system under equivalent training and testing conditions. One possible explanation is related to the number of phone transitions that the network has to learn. In the DIGITS+ corpus and the NUMBERS’93 task the number of possible transitions is small (less than a hundred) , in contrast to more than 1000 possible transitions in the TIMIT case<sup>2</sup>. Therefore, it might be possible to learn a small number of transitions with “hard” targets, but this approach may not be scalable to a task with a large number of transitions, or at the very least might require significantly more data for training.

#### Initialization

In order to start the REMAP iteration there is a need for a trained MLP to supply the initial conditional transition probabilities for the forward and backward recursions. In these experiments, these probabilities were estimated using a network trained on “hard targets.” It might be that this initialization limits the training to networks that are too similar to the initial network. One way try to overcome this potential limitation might be to add noise (analogous to simulated annealing techniques) to the initial MLP weights. Alternatively we can blur the “hard” targets, e.g., by adding a Gaussian centered around each target. Both of these approaches might result in initial transition probabilities that are smoother than the ones that are estimated by the MLP trained on “hard targets.”

---

<sup>2</sup>The exact number of possible transitions are function of the phone set in use

## Chapter 8

# Conclusions and Future Work

### 8.1 Conclusions

This thesis introduced a framework for training and modeling continuous speech recognition systems based on the theoretically optimal MAP criterion. In contrast, most current state-of-the-art systems are trained according to other criteria such as Maximum Likelihood (ML). Our proposed HMM/ANN paradigm is based on re-estimating ANN targets and weights to guarantee an increase in the posterior probability of the correct model (sentence). We have described the theory behind the new framework, included a convergence proof for the training algorithm, and reported on experimental results that support the proposed paradigm.

Our studies have shown that explicitly modeling transitions between speech units can improve recognition performance. Specifically, we have shown that accurate boundary information between phones can improve recognition performance significantly. We have studied a transition-based model that uses local transition probabilities (i.e., the posterior probability of the current state given the current acoustic vector and the previous state) to estimate global posterior of sentences. Therefore, it is a true recognition model, and it directly maps from acoustic sequences to sentences, unlike HMMs that model the inverse relation (the likelihood of producing an acoustic sequence)

Early experiments with the proposed model and the Viterbi-based training procedure showed the need for a new training algorithm. Specifically, the Viterbi-based training (when applied to the transition-based model) suffered from poor transition detection and a mismatch between the input space distribution during training and recognition. These

reasons, combined with the goal of training HMM/ANN systems based on the theoretically optimal MAP criterion, led to a new discriminant training algorithm, REMAP.

REMAP is a method to estimate parameters and train systems that classify sequences according to the MAP criterion. This can be used in a new form of hybrid HMM/MLP that, in addition to the advantages of standard HMM/MLP hybrids, uses “full” posterior probabilities for training and recognition. We still use neural nets (in our case MLPs, though recurrent nets or TDNNs could be used) to estimate local posterior probabilities (conditional transition probabilities), but our nets are trained with probabilistic targets that are themselves estimates of local posterior probabilities (conditioned on the acoustic data and the previous state).

There is now a way, similar in spirit to the forward-backward recursions of the Baum-Welch algorithm, of estimating these optimal targets given a previously trained neural network. Additionally, we have a convergence proof that guarantees iterative increase of the global posterior probability. This method is valid for any hybrid HMM/ANN system but, in this thesis, was developed in the framework of “Discriminant HMMs” using conditional transition probabilities.

While our gains with this approach are still small, in many ways the experimental results support the proposed framework. Experimenting with an isolated speech recognition task resulted in an increase in the estimates of the posterior probabilities of the correct sentences after training. Extending these experiments to a continuous speech recognition task achieved a significant decrease in error rate in comparison to a baseline system. While learning a larger number of transitions may be difficult, it still may be the case that larger tasks will be better for REMAP, since less probability mass will be associated with “illegal” state sequences.

## 8.2 Future Work

### 8.2.1 Incorporating Language Information

This work and other current approaches to speech recognition modeling assume independent models for language and acoustic information. Usually the language model estimates prior probabilities of sentences by computing n-grams (the probability of a word given the previous (n-1) words), by counting the relative frequencies in the training test

or in a large text corpus<sup>1</sup>. Acoustic models use a separate set of parameters, which are estimated independently from the language model parameters.

In the posterior-based framework proposed here, we should estimate the posterior probability of the sentence given *both* the acoustics and high-level knowledge such as a language model (see Chapter 2). Furthermore, during training we want to optimize the measure that we use in recognition. Therefore, we want to maximize the posterior probability of the correct model given both the acoustics and the language model. However, the training algorithm that is proposed in this study, REMAP, increases the posterior probability of the correct model based only on acoustic information.

One possible solution is described in Section 6.4. Specifically, Equation (6.25) describes a way to incorporate the language model into the estimation of the posterior probability of a particular model. So far, we have not not experimented with this approach. Another possible solution is to estimate “local” probabilities, e.g., local conditional posterior probabilities in this work, based on both the local acoustic and on language information. Specifically, we could take into account whether the transition between two phones occurs in the middle of a word or between two words, e.g., one phone is the last phone of the current word and the second is the first in the following word. A possible implementation of this idea is to add an additional input to the MLP described in Figure 5.2 to represent the words of the source and target phones. In the case of large vocabularies, the words can be clustered into classes as has been done in class grammar techniques. Unfortunately, this increases the dimensionality of the input space, resulting in a need for more training data and more computation.

Several researchers have worked on learning pronunciation models using both the linguistic and the acoustic training data (Stolcke & Omohundro 1993b; Stolcke & Omohundro 1993a; Wooters 1993). Extending this line of work towards learning probabilistic grammars based on both the linguistic and the acoustic training data seems like a direction worth pursuing. Eventually we might learn the models themselves (not just the parameters of the models) from the data. Complementing this data-driven direction as much domain-specific knowledge should be incorporated as possible. For instance, new experiments are now being done at ICSI using enforced minimum duration constraints with REMAP.

---

<sup>1</sup>Generally, smoothing techniques are used to get a better estimate of rare or non-occurring word transitions, e.g., averaging the bigram probabilities with the product of the corresponding unigram probabilities.

## 8.2.2 Extensions

### A Smaller Phone Set

In section 7.3 we describe a potential problem regarding invalid state sequences. We showed that in a small isolated speech experiment, around 80% of the probability mass was dedicated to a “garbage model”, which contains all phone sequences that are not part of any valid model. Furthermore, our speculation was that in the continuous speech task, this problem is not so severe, given the large number of legal models which encompass more phone sequences. Additionally, a potential partial solution is to use a smaller phone set, such that the phone sequences “covered” by the legal models will have a larger portion from all phone sequences.

### Recurrent Nets

In all the mathematical development through this thesis we rely on the assumption that to estimate conditional transition probabilities, we can represent both the acoustic past and the previous state sequence by the most recent previous state and by a temporal window of several acoustic frames into the past. This simplifying assumption can be relaxed by replacing the MLP used in this study by a recurrent neural network (Robinson 1994). A recurrent neural network encodes the past in its state units. Hence, the estimation of the transition probabilities is conditioned on the captured history of the utterance.

### Negative Training

In section 5.4.2 we describe a potential problem regarding a mismatch between the input space distribution during training and recognition. During training the net is only presented with “correct” pairs of acoustic vectors and the correct previous state, while in recognition the net is expected to generalize to all possible combinations of acoustic vectors and previous states. A possible solution might be to do explicit negative training. Specifically, we would train the network with all possible combinations of acoustic vectors and previous states, but for the “incorrect” previous states we would present uniform targets, i.e., all targets would have the same value. We can estimate the posterior probability of each previous state, by running the alpha and beta recursions in a fully-connected phone model. According to this estimated distribution of the previous state we would select the

pairs of acoustic vectors and previous states as described in Section 6.3. In other work, Lyon and Yaeger recently showed that explicit training on negative examples can improve recognition for on-line hand-recognition with ANNs (Lyon & Yaeger 1996).

### Modeling Extensions

In this study although we could increase the complexity we have used a simple 1st order Markov model. The extension to an M-th order model is described in (Bourlard *et al.* 1994). However, the increased complexity of the model requires more training data to reliably estimate the model parameters. Further, a more complex model has no guarantee to work better. Other modeling extensions might be to encode speaker-dependent information such as gender, speaking rate, and accent as extra inputs to the Discriminant HMM. Another possible extension is to apply REMAP to perceptually-motivated models such as SPAM (Morgan *et al.* 1994; Morgan *et al.* 1995).

## 8.3 Epilog

In experimental science such as speech recognition, a division between *data description* and *data modeling* can be drawn (Cohen 1995). Data description is taking an existing model and increasing its complexity, adding more exceptions to better fit the observed data; for instance, adding triphones to HMMs. In contrast, data modeling refers to finding models that are a better match to the underlying phenomena, such as a model based on a parameter for vocal tract length. Furthermore, these new models should be a simpler explanation of the observed data, with better generalization capabilities, and not just a more complex version of the existing models.

Currently, in speech recognition, the most popular models are HMMs, with a number of labs and companies using hybrid HMM/ANN in different forms. In this thesis we suggest a different framework than HMMs for speech recognition. The main difference between HMMs and our suggested framework is in the modeling objective. HMMs model a production system, i.e., all the different ways that a given sentence can be realized. Their input is a sentence and their output its acoustic realizations. Thus, a mismatch with observed data could be solved through a more detailed model of the production process, e.g., having different models for males and females. The framework proposed here models a recognition process, i.e., the input is an acoustic utterance and the output is a sentence.



We estimate a fixed number of classes (phones in our implementation) given the context. The context can represent the previous state (as done in this study), the speaking rate of the speaker, the gender of the speaker, the noise level, etc.

The long term usefulness of the proposed framework may be determined on the performance on larger and more challenging speech tasks (such as the Switchboard corpus (NIST 1992)) than tested here and on application of it to other domains that are also sequential in nature such as hand-writing recognition. For REMAP to become a popular approach for speech recognition, two main objectives must be achieved. The first objective is to incorporate more speech specific knowledge into this framework, such as duration constraints, speaking rate, etc. Initial experiments with enforcing minimum duration constraints are beginning to produce promising results. The second objective is to integrate the language information into the framework. We do believe that the framework presented here is a good start.

## Appendix A

# REMAP Convergence - Theorem Proof

### A.1 Theorem 1

**Theorem 1:**

*IF*  $R(v_1, v_2) \geq R(v_1, v_1)$

*THEN*  $P(M|X, v_2) \geq P(M|X, v_1)$ .

In other words, if we can find a new set of probabilities  $v_2$  increasing  $R$ , the new set of probabilities will also increase the posterior probability of the model  $M$ .

**Proof:**

$$\begin{aligned}
 & \log \frac{P(M|X, v_2)}{P(M|X, v_1)} \\
 &= \log \left[ \sum_{\Gamma} \frac{P(M, \Gamma|X, v_1)}{P(M|X, v_1)} \frac{P(M, \Gamma|X, v_2)}{P(M, \Gamma|X, v_1)} \right] \\
 &\geq \sum_{\Gamma} \frac{P(M, \Gamma|X, v_1)}{P(M|X, v_1)} \log \left[ \frac{P(M, \Gamma|X, v_2)}{P(M, \Gamma|X, v_1)} \right] \\
 &\quad \text{(because of Jensen's inequality and concavity of log function)} \\
 &= R(v_1, v_2) - R(v_1, v_1)
 \end{aligned}$$

(Note that the random variable used for the Jensen's inequality is a  $\frac{P(M, \Gamma|X, v_2)}{P(M, \Gamma|X, v_1)}$  which is a

function of the random variable  $\Gamma$ ). As a consequence, we have:

$$\log \frac{P(M|X, v_2)}{P(M|X, v_1)} \geq R(v_1, v_2) - R(v_1, v_1) \quad (\text{A.1})$$

which proves the theorem. If a new set of probabilities  $v_2$  that makes the right-hand side of (A.1) positive can be found, then the model re-estimation algorithm can be guaranteed to increase the posterior probability of the model to  $P(M|X, v_2)$ .

## A.2 Theorem 2

### Theorem 2:

Given  $v_1$ , a fixed set of probabilities that is estimated by an MLP with a fixed set of weights  $\Theta$ , we show that  $R(v_1, v_2)$  attains its local maximum value when the conditional transitional probabilities  $P_{v_2}(q_\ell^n | x_n, q_k^{n-1})$  are defined as<sup>1</sup>

$$P_{v_2}(q_\ell^n | x_n, q_k^{n-1}) = P(q_\ell^n | X, q_k^{n-1}, v_1(\Theta), M) \quad (\text{A.2})$$

### Proof:

We now treat the conditional transitional probabilities  $P_{v_2}(\cdot)$  as the variables for the optimization. Thus, we need to maximize  $R(\cdot)$  in the space of transition-probability sets  $\Upsilon$ , under the  $KN$  constraints

$$\sum_{j=1}^K P_{v_2}(q_j^n | x_n, q_k^{n-1}) = 1, \quad \forall k = 1, \dots, K; \forall n = 1, \dots, N \quad (\text{A.3})$$

Using Lagrange multipliers  $\Lambda = (\lambda_{1,1}, \dots, \lambda_{1,N}, \dots, \lambda_{K,1}, \dots, \lambda_{K,N})^t$ , maximization of  $R(\cdot)$  as defined in (6.2) under the constraints specified in (A.3) is then equivalent to maximization of

$$R^*(v_1, v_2, \Lambda) = R(v_1, v_2) + \sum_{k,n} \lambda_{k,n} \left( 1 - \sum_{j=1}^K P_{v_2}(q_j^n | x_n, q_k^{n-1}) \right) \quad (\text{A.4})$$

So there are  $K^2N$  unknowns that are the conditional transition probabilities in  $v_2$  and  $KN$  unknowns that are the Lagrange multipliers. Fortunately, there are the same number of equations as we compute the partial derivative of  $R^*(\cdot)$  relative to each unknown and

---

<sup>1</sup>Of course, all  $x_n$ 's in the following should be replaced by  $X_{n-c}^{n+d}$  if local contextual input is used, or  $X_1^n$  for a recurrent network.

equalize it to zero. Furthermore, it turns out that we can solve each of the  $K^2N$  equations described above independently and find solutions that satisfy the  $KN$  constraints.

Considering a specific transition  $(q_k^{n-1}, q_\ell^n)$ , we then have:

$$\frac{\partial R^*(\cdot)}{\partial \lambda_{k,n}} = 0 \quad (\text{A.5})$$

which returns the constraint (A.3). For the partial derivative of  $R^*(\cdot)$  with respect to  $P_{v_2}(\cdot)$ , we first use the following decomposition:

$$P(M, \Gamma|X, v_2) = P(\Gamma|X, v_2)P(M|\Gamma, X, v_2) \quad (\text{A.6})$$

According to (5.7), the first factor in (A.6) can be expressed as

$$P(\Gamma|X, v_2) = \prod_{n=1}^N P_{v_2}(q_\ell^n|x_n, q_k^{n-1}) \quad (\text{A.7})$$

Also, the second factor in (A.6) can be assumed independent of the conditional transition probabilities (i.e., given a state sequence, the probability of the model does not depend on the transition probabilities), in which case we have:

$$P(M|\Gamma, X, v_2) = P(M|\Gamma, X) \quad (\text{A.8})$$

Taking partial derivatives, then, the second term in (A.6) has no effect, since it can be assumed to have no dependence on  $P_{v_2}(q_\ell^n|x_n, q_k^{n-1})$  and since it only appears as an additive term once the logarithmic function has been applied.

We then get

$$\begin{aligned} & \frac{\partial R^*(\cdot)}{\partial P_{v_2}(q_\ell^n|x_n, q_k^{n-1})} \\ &= \frac{1}{P(M|X, v_1)} \sum_{\Gamma_{k,\ell,n}} \left[ P(M, \Gamma_{k,\ell,n}|X, v_1) \frac{1}{P_{v_2}(q_\ell^n|x_n, q_k^{n-1})} \right] - \lambda_{k,n} \\ &= 0 \end{aligned} \quad (\text{A.9})$$

where  $\Gamma_{k,\ell,n}$  stands for those paths containing the transition  $(q_k^{n-1}, q_\ell^n)$ . Solving (A.9) gives:

$$P_{v_2}(q_\ell^n|x_n, q_k^{n-1}) = \frac{1}{\lambda_{k,n}} \frac{\sum_{k,\ell,n} P(M, \Gamma_{k,\ell,n}|X, v_1)}{P(M|X, v_1)} \quad (\text{A.10})$$

We now have to find the value (or ‘‘a’’ value) of  $\lambda_{k,n}$  that guarantees that the new estimates of  $P_{v_2}(q_\ell^n|x_n, q_k^{n-1})$  will meet the constraint. It is possible to find it without

directly solving the set of equations. It is indeed easy to show that:

$$\begin{aligned}
\frac{\sum_{k,\ell,n} P(M, \Gamma_{k,\ell,n} | X, v_1)}{P(M|X, v_1)} &= \frac{\sum_{\forall \Gamma} P(M, \Gamma, q_\ell^n, q_k^{n-1} | X, v_1)}{P(M|X, v_1)} \\
&= \frac{P(M, q_\ell^n, q_k^{n-1} | X, v_1)}{P(M|X, v_1)} \\
&= P(q_\ell^n, q_k^{n-1} | X, M, v_1) \\
&= P(q_\ell^n | q_k^{n-1}, X, M, v_1) P(q_k^{n-1} | X, M, v_1) \quad (\text{A.11})
\end{aligned}$$

Since the second factor in (A.11) is a function of  $k$  and  $n$  only, we can set  $\lambda_{k,n}$  to  $P(q_k^{n-1} | X, M, v_1)$  which then gives us, according to (A.10):

$$P_{v_2}(q_\ell^n | x_n, q_k^{n-1}) = P(q_\ell^n | X, q_k^{n-1}, M, v_1) \quad (\text{A.12})$$

This is a valid solution since the constraint

$$\sum_{\ell=1}^K P(q_\ell^n | X, q_k^{n-1}, M, v_1) = 1, \forall k \text{ and } \forall n$$

given in (A.3) is automatically met.<sup>2</sup> In order to verify that we have a local maximum point we have to compute the Hessian matrix. It is easy to see by looking at (A.9) that all the non diagonal elements are zero. In computing the diagonal elements we get

$$\begin{aligned}
&\frac{\partial^2 R^*(\cdot)}{\partial P_{v_2}^2(q_\ell^n | x_n, q_k^{n-1})} \\
&= -\frac{1}{P(M|X, v_1)} \sum_{\Gamma_{k,\ell,n}} \left[ P(M, \Gamma_{k,\ell,n} | X, v_1) \frac{1}{P_{v_2}(q_\ell^n | x_n, q_k^{n-1})} \right] \quad (\text{A.13})
\end{aligned}$$

and it is obvious that for probabilities (i.e., positive numbers) we get negative diagonal elements. Thus, this is a local maximum point. This proves Theorem 2.

### A.3 Theorem 3

#### Theorem 3:

When minimizing the weighted relative entropy criterion (20) with the target set  $T$  (which is calculated from a probability set  $v_1$ ), the auxiliary function  $R(\cdot)$  is maximized. Specifically

---

<sup>2</sup>We cannot prove that this is a unique solution since most of the equations are nonlinear, but we know this is at least one valid solution.

the new set of probabilities  $v_2$ , implemented by the trained MLP satisfies the following:

$$E_e(\Theta^{t+1}) \leq E_e(\Theta^t) \implies R(v_1(\Theta^t), v_2(\Theta^{t+1})) \geq R(v_1(\Theta^t), v_1(\Theta^t)) \quad (\text{A.14})$$

**Proof:**

$$\begin{aligned} E_e(\Theta^{t+1}) - E_e(\Theta^t) &= \sum_{n=1}^N \sum_{k=1}^K P(q_k^{n-1}|X, M, v_1(\Theta^t)) \\ &\quad \sum_{\ell=1}^K P(q_\ell^n|X, q_k^{n-1}, M, v_1(\Theta^t)) \log \frac{g_\ell(x_n, q_k^{n-1}, \Theta^t)}{g_\ell(x_n, q_k^{n-1}, \Theta^{t+1})} \\ &= \sum_{n=1}^N \sum_{k=1}^K \sum_{\ell=1}^K P(q_\ell^n, q_k^{n-1}|X, M, v_1(\Theta^t)) \log \frac{g_\ell(x_n, q_k^{n-1}, \Theta^t)}{g_\ell(x_n, q_k^{n-1}, \Theta^{t+1})} \end{aligned} \quad (\text{A.15})$$

Below we show that the change in the auxiliary function  $R(\cdot)$  has the same magnitude then  $E_e$ , but was with the opposite sign.

$$\begin{aligned} R(v_1(\Theta^t), v_2(\Theta^{t+1})) - R(v_1(\Theta^t), v_1(\Theta^t)) &= \sum_{\Gamma} \frac{P(M, \Gamma|X, v_1)}{P(M|X, v_1)} \log \left[ \frac{P(M, \Gamma|X, v_2)}{P(M, \Gamma|X, v_1)} \right] \\ &= \sum_{\Gamma} \frac{P(M, \Gamma|X, v_1)}{P(M|X, v_1)} \log \left[ \frac{P(M|\Gamma, X, v_2)P(\Gamma|X, v_2)}{P(M|\Gamma, X, v_1)P(\Gamma|X, v_1)} \right] \end{aligned} \quad (\text{A.16})$$

Given a state sequence, the probability of the model does not depend on the transition probabilities ( $v_1$  and  $v_2$ ). As stated in (A.8), we then get:

$$\begin{aligned} R(v_1(\Theta^t), v_2(\Theta^{t+1})) - R(v_1(\Theta^t), v_1(\Theta^t)) &= \sum_{\Gamma} \frac{P(M, \Gamma|X, v_1)}{P(M|X, v_1)} \log \left[ \frac{P(\Gamma|X, v_2)}{P(\Gamma|X, v_1)} \right] \\ &= \sum_{\Gamma} \frac{P(M, \Gamma|X, v_1)}{P(M|X, v_1)} \log \left[ \frac{\prod_{n=1}^N P(q_\ell^n|x_n, q_k^{n-1}, \Theta_{t+1})}{\prod_{n=1}^N P(q_\ell^n|x_n, q_k^{n-1}, \Theta_t)} \right], \quad [\text{as in (A.7)}] \\ &= \sum_{\Gamma} P(\Gamma|M, X, v_1) \log \left[ \frac{\prod_{n=1}^N P(q_\ell^n|x_n, q_k^{n-1}, \Theta_{t+1})}{\prod_{n=1}^N P(q_\ell^n|x_n, q_k^{n-1}, \Theta_t)} \right] \\ &\quad \text{rearranging the terms in the summation :} \\ &= \sum_{n=1}^N \sum_{k=1}^K \sum_{\ell=1}^K \sum_{\Gamma_{k,\ell,n}} P(\Gamma_{k,\ell,n}|M, X, v_1) \log \frac{g_\ell(x_n, q_k^{n-1}, \Theta^{t+1})}{g_\ell(x_n, q_k^{n-1}, \Theta^t)} \end{aligned}$$

$$\begin{aligned}
&= \sum_{n=1}^N \sum_{k=1}^K \sum_{\ell=1}^K \sum_{\Gamma} P(\Gamma, q_k^{n-1}, q_\ell^n | M, X, v_1(\Theta_t)) \log \frac{g_\ell(x_n, q_k^{n-1}, \Theta^{t+1})}{g_\ell(x_n, q_k^{n-1}, \Theta^t)} \\
&= \sum_{n=1}^N \sum_{k=1}^K \sum_{\ell=1}^K P(q_k^{n-1}, q_\ell^n | M, X, v_1(\Theta_t)) \log \frac{g_\ell(x_n, q_k^{n-1}, \Theta^{t+1})}{g_\ell(x_n, q_k^{n-1}, \Theta^t)}
\end{aligned} \tag{A.17}$$

in which  $\Gamma_{k,\ell,n}$  stands for those paths containing the transition  $(q_k^{n-1}, q_\ell^n)$ .

A closer look at the last equation shows the term that we got for the difference in the auxiliary function  $R(\cdot)$  is with an opposite sign and the same magnitude as the difference in  $E_e$  in (A.15). Thus, minimizing the cost function  $E_e$  (as part of the MLP training) is equivalent to maximizing the auxiliary function  $R(\cdot)$ . Hence, we have proved Theorem 3, and in fact showed that minimizing the error criterion (6.8) is equivalent (within a scaling factor) to maximizing the auxiliary function.

# Bibliography

- BAHL, L. R., P. F. BROWN, P. V DE SOUZA, & R. L. MERCER. 1986. Maximum mutual information estimation of hidden Markov model parameters. In *Proc. Int. Conf. on Acoust., Speech, and Signal Processing*, 49–52, Tokyo, Japan.
- BAKER, J. K., 1975. *Stochastic Modeling as a Means of Automatic Speech Recognition*. Carnegie Mellon University dissertation.
- BAUM, L. E. 1972. An inequality and associated maximization technique in statistical estimation for probabilistic functions of Markov processes. *Inequalities* 3.1–8.
- , & T. PETRIE. 1966. Statistical inference for probabilistic functions of finite state Markov chains. *The Annals of Mathematical Statistics* 36.1554–1563.
- , T. PITRIE, G. SOULES, & N. WEISS. 1970. A maximization technique occurring in the statistical analysis of probabilistic functions in Markov chains. *The Annals of Mathematical Statistics* 41.164–171.
- BENGIO, Y., & P. FRASCONI. 1995. An input output HMM architecture. In *Advances in Neural Information Processing Systems*, ed. by G. Tesauro, D. Touretzky, & T. Leen, volume 7. Cambridge: MIT press.
- , R. DE MORI, G. FLAMMIA, & R. KOMPE. 1992. Global optimization of a neural network-hidden Markov model hybrid. *IEEE trans. on Neural Networks* 3.252–258.
- BOURLARD, H., Y. KONIG, & N. MORGAN. 1994. REMAP: Recursive estimation and maximization of a posteriori probabilities, application to transition-based connectionist speech recognition. Technical Report TR-94-064, International Computer Science Institute, Berkeley, CA.



- , & N. MORGAN. 1994. *Connectionist Speech Recognition - A Hybrid Approach*. Kluwer Academic Publishers.
- , & C. J. WELLEKENS. 1989a. Links between Markov models and multilayer perceptrons. In *Advances in Neural Information Processing Systems 1*, ed. by D.J. Touretzky, 502–510, San Mateo. Morgan Kaufmann.
- , & —. 1989b. Speech pattern discrimination and multilayer perceptrons. *Computer, Speech, and Language* 3.1–19.
- BRIDLE, J. S. 1990. Alpha-nets: A recurrent "neural" network architecture with a hidden Markov model interpretation. *Speech Communication* 9.83–92.
- BROWN, P. F., 1987. *The Acoustic-Modelling Problem in Automatic Speech Recognition*. Pittsburgh, PA: CMU dissertation.
- COHEN, J. R., 1995. Informal Communication.
- COLE, R. A., M. FANTY, & T. LANDER. 1994. Telephone speech corpus development at CSLU. In *Proceedings Int'l Conference on Spoken Language Processing*, Yokohama, Japan.
- COVER, T. M., & A. T. JOY. 1991. *Elements of Information Theory*. New York: John Wiley and Sons, Inc.
- DEMPSTER, A. P., N. M. LAIRD, & D. B. RUBIN. 1977. Maximum likelihood from incomplete data via the *EM* algorithm. *Journal of the Royal Statistical Society, Series B* 34.1–38.
- DENG, L. 1992. A generalized hidden Markov model with state-conditioned trend functions of time for the speech signal. *Signal Processing* 27.65–78.
- , & D.X. SUN. 1994. A statistical approach to automatic speech recognition using the atomic speech units constructed from overlapping articulatory features. *J. Acoustical Society of America* 95.2702–2719.
- DIGALAKIS, V. V., 1992. *Segment-Based Stochastic Models of Spectral Dynamics for Continuous Speech Recognition*. Boston University dissertation.

- , J. R. ROHLICEK, & M. OSTENDORF. 1993. Segment-based stochastic models of spectral dynamics for continuous speech recognition. *IEEE trans. on Speech and Audio Processing* 1.431–442.
- DUDA, R. O., & P. E. HART. 1973. *Pattern Classification and Scene Analysis*. John Wiley and Sons, Inc.
- EFRON, B. 1982. The jackknife, the bootstrap, and other resampling plans. Society for Industrial and Applied Mathematics, Philadelphia, Pa.
- FUJIMURA, O. 1975. Syllable as a unit of speech recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing* ASSP-23.82–87.
- FURUI, S. 1986a. On the role of spectral transition for speech perception. *J. Acoustical Society of America* 80.1016–1025.
- 1986b. Speaker independent isolated word recognizer using dynamic features of speech spectrum. *IEEE Trans. on Acoustics, Speech, and Signal Processing* 34.52–59.
- GAROFOLO, J. S., 1988. *Getting Started with the DARPA TIMIT CD-ROM: an Acoustic Phonetic Continuous Speech Database*. National Institute of Standards and Technology (NIST), Gaithersburgh, Maryland.
- GHITZA, O., & M. M. SONDEHI. 1993. Hidden Markov models with templates as non-stationary states: an application to speech recognition. *Computer Speech and Language* 2.101–119.
- GISH, H. 1990. A probabilistic approach to the understanding and training of neural network classifiers. In *Proc. Int. Conf. on Acoust., Speech, and Signal Processing*, 1361–1364, Albuquerque, NM.
- GLASS, J. R., 1988. *Finding Acoustic Regularities in Speech Applications to Phonetic Recognition*. M.I.T dissertation.
- GOLDENTHAL, W. D., 1994. *Statistical Trajectory Models for Phonetic Recognition*. M.I.T dissertation.

- HAEB-UMBACH, R., & H. NEY. 1992. Linear discriminant analysis for improved large vocabulary continuous speech recognition. In *Proc. Int. Conf. on Acoust., Speech, and Signal Processing*, volume 1, 13–16, San Francisco, USA.
- HERMANSKY, H. 1990. Perceptual linear predictive (PLP) analysis of speech. *J. Acoustical Society of America* 87.
- HUANG, X. D., Y. ARIKI, & M. A. JACK. 1990. *Hidden Markov Models For Speech Recognition*. Edinburgh University Press.
- JELINEK, F. 1976. Continuous speech recognition by statistical methods. *Proceedings of the IEEE* 64.532–556.
- , & R. L. MERCER. 1980. Interpolated estimation of Markov source parameters from sparse data. In *Pattern Recognition in Practice*, ed. by E. S. Gelsema & L. N. Kanal, 381–397. Amsterdam, The Netherlands: North-Holland Publishing Company.
- JUANG, B. H., & L. R. RABINER. 1985. Mixture autoregressive hidden Markov models for speech signals. *IEEE ASSP Magazine* 6.1404–1413.
- KATAGIRI, S., C. H. LEE, & JUANG B. H. 1991. New discriminative training algorithms based on the generalized probabilistic decent method. In *Proc. of the IEEE Workshop on Neural Networks for Signal Processing*, ed. by B.H. Juang, S.Y. Kung, & C.A. Kamm, 299–308.
- KIANG, N. Y. S. 1984. Peripheral neural processing of auditory information. In *In Handbook of physiology - the nervous system, volume III, sensory processes, part 2*, ed. by J. M. Brookhart & V. Mountcastle, 639–674. Bethesda: American Physiological Society.
- KONIG, Y., & N. MORGAN. 1993. Supervised and unsupervised clustering of the speaker space for connectionist speech recognition. In *Proc. Int. Conf. on Acoust., Speech, and Signal Processing*, Minneapolis, Minnesota. IEEE.
- , & —. 1994. Modeling dynamics in connectionist speech recognition - the time index model. In *Proceedings Int'l Conference on Spoken Language Processing*, 1523–1526, Yokohama, Japan.

- , —, C. WOOTERS, V. ABRASH, M. COHEN, & H. FRANCO. 1993. Modeling consistency in a speaker independent continuous speech recognition system. In *Advances in Neural Information Processing Systems 5*, ed. by J.S. Hanson & J.D. Cowan and C.L. Giles, volume 5, San Mateo. Morgan Kaufman.
- LEE, K. F. 1989. *Automatic Speech Recognition: The Development of the SPHINX System*. Kluwer Academic Publishers.
- LEVELT, W. J. M., & L. WHEELDON. 1994. Do speakers have access to a mental syllabary? *Cognition* 50.239–269.
- LEVINSON, S. E. 1985. Structural methods in automatic speech recognition. *Proceedings of the IEEE* 73.1625–1650.
- , L. R. RABINER, & M. M. SONDHI. 1983a. An introduction to the application of the theory of probabilistic functions on a Markov process to automatic speech recognition. *Bell System Technical Journal* 62.243–272.
- , LAWRENCE R. RABINER, & MAN MOHAN SONDHI. 1983b. An introduction to the application of the theory of probabilistic functions of a Markov process to automatic speech recognition. *Bell Syst. Tech. J.* 62.1035–1074.
- LINDBLOM, B. E. F., & M. STUDDERT-KENNEDY. 1967. On the role of formant transitions in vowel recognition. *J. Acoustical Society of America* 42.830–843.
- LIPORACE, L. A. 1982. Maximum likelihood estimation for multivariate observations of markov sources. *IEEE Trans. on Information Theory* IT-28.729–734.
- LIPPMANN, R. P. 1989. Review of neural networks for speech recognition. *Neural Computation* 1.1–38.
- LUBENSKY, D. M., A. O. ASADI, & J. M. NAIK. 1994. Connected digit recognition using connectionist probability estimators and mixture-gaussian densities. In *Proceedings Int'l Conference on Spoken Language Processing*, 295–298, Yokohama, Japan.
- LYON, R. F., & L. S. YAEGER. 1996. On-line hand-printing recognition with neural networks. In *MicroNeuro '96 – Fifth International Conference on Microelectronics for Neural Networks and Fuzzy Systems*, 201–212, Lausanne, Switzerland. IEEE.

- MERIALDO, B. 1988. Phonetic recognition using hidden Markov models and maximum mutual information training. In *Proc. Int. Conf. on Acoust., Speech, and Signal Processing*, 111–114, New York.
- MORGAN, N., H. BOURLARD, S. GREENBERG, & H. HERMANSKY. 1994. Stochastic perceptual auditory-event-based models for speech recognition. In *Proceedings Int'l Conference on Spoken Language Processing*, 1943–1946, Yokohama, Japan.
- , S. WU, & H. BOURLARD. 1995. Digit recognition with stochastic perceptual speech models. In *Proceedings of the Sixth European Conference on Speech Communication and Technology*, Madrid, Spain.
- NADAS, A. 1983. A decision-theoretic formulation of a training problem in speech recognition and a comparison of training by unconditional versus conditional maximum likelihood. *IEEE Transactions on Acoustics, Speech, and Signal Processing* 31.814–817.
- NIST, 1992. Switchboard Corpus: Recorded Telephone Conversations. National Institute of Standards and Technology Speech Disc 9-1.1 to 9-29.1, Collected by Texas Instruments.
- OSTENDORF, M., & S. ROUKOS. 1989. A stochastic segment model for phoneme-based continuous speech recognition. *IEEE ASSP trans.* 37.1857–1869.
- PORITZ, A. B., & A. L. RICHTER. 1986. On hidden Markov models in isolated word recognition. In *Proc. Int. Conf. on Acoust., Speech, and Signal Processing*, volume 1 of *Tokyo, Japan*, 705–708.
- PRICE, P., W. FISHER, J. BERNSTEIN, & D. PALLET. 1988. The darpa 1000-word resource management database for continuous speech recognition. In *Proceedings IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, volume 1, 651–654, New York. IEEE.
- RABINER, L.R. 1989. A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE* 77.257–285.
- RENALS, S., N. MORGAN, & H. BOURLARD. 1991. Probability estimation by feed-forward networks in continuous speech recognition. Technical Report TR-91-030, International Computer Science Institute, 1947 Center St., Suite 600, Berkeley, CA 94704.

- , —, —, M. COHEN, & H. FRANCO. 1994. Connectionist probability estimators in HMM speech recognition. *IEEE Trans. on Speech and Audio Processing* 2.161–174.
- , N. MORGAN, M. COHEN, H. FRANCO, & H. BOURLARD. 1992. Connectionist probability estimation in the DECIPHER speech recognition system. In *Proc. Int. Conf. on Acoust., Speech, and Signal Processing*, 601–604, San Francisco, California. IEEE.
- RICHARD, M. D., & R. P. LIPPMANN. 1991. Neural network classifiers estimate bayesian a posteriori probabilities. *Neural Computation* 3.461–483.
- ROBINSON, A. J. 1994. An application of recurrent nets to phone probability estimation. *IEEE transactions on Neural Networks* 5.298–305.
- , L. ALMEIDA, J.-M. BOITE, H. BOURLARD, F. FALLSIDE, M. HOCHBERG, D. KERSHAW, P. KOHN, Y. KONIG, N. MORGAN, J. P. NETO, S. RENALS, M. SAERENS, & C. WOOTERS. 1993. A neural network based, speaker independent, large vocabulary, continuous speech recognition system: The WERNICKE project. In *Proceedings European Conf. on Speech Communication and Technology. (EUROSPEECH)*, 1941–1944, Berlin, Germany.
- RUGGERO, M. 1994. Physiology and coding of sound in the auditory nerve. In *The Mammalian Auditory Pathway: Neurophysiology*, ed. by A. Popper & R. Fay, 34–93. New York, USA: Springer.
- RUMELHART, D. E., G. E. HINTON, & R. J. WILLIAMS. 1986. Learning internal representations by error propagation. In *Parallel Distributed Processing. Explorations of the Microstructure of Cognition*, ed. by D. E. Rumelhart & J. L. McClelland, volume 1: Foundations. MIT Press.
- SACHS, M. B., R. L. WINSLOW, & C. C. BLACKBURN. 1988. Representation of speech in the auditory periphery. In *Auditory Function: Neurobiological Bases of Hearing*, ed. by G. M. Edelman, W. E. Gall, & W. M. Cowan, 747–774. New York, USA: Wiley.
- SCHWARTZ, R., & S. AUSTIN. 1991. A comparison of several approximate algorithms for finding multiple (N-Best) sentence hypotheses. In *Proc. Int. Conf. on Acoust., Speech, and Signal Processing*, Toronto, Canada.

- , & Y. CHOW. 1990. The N-Best algorithm: An efficient and exact procedure for finding the n most likely sentence hypotheses. In *Proc. Int. Conf. on Acoust., Speech, and Signal Processing*, New Mexico, USA.
- SEGUI, J., E. DUPOUX, & J. MEHLER. 1980. The role of the syllable in speech segmentation, phoneme identification, and lexical access. In *Cognitive models of speech processing: Psycholinguistic and computational perspectives. ACL-MIT Press series in natural language processing*, ed. by G. T. M. Altmann, 263–280. Cambridge, MA, USA: MIT Press.
- SENEFF, S. 1988. A joint synchrony/mean-rate model of auditory processing. *Journal of Phonetics* 16.55–76.
- SMITH, R., & J. J. ZWISLOCKI. 1975. Short-term adaptation and incremental responses of single auditory-nerve fibers. *Biological Cybernetics* 17.169–182.
- STOLCKE, A., & S. OMOHUNDRO. 1993a. Best-first model merging for Hidden Markov Model induction. Technical report, International Computer Science Institute, 1947 Center St. Suite 600, Berkeley, CA.
- , & —. 1993b. Hidden Markov Model induction by Bayesian model merging. In *Advances in Neural Information Processing Systems 5*, 11–18. San Mateo, Ca.: Morgan Kaufman.
- VITERBI, A. J. 1967. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE Trans. on Information Theory* 13.260–269.
- WAIBEL, A., HANAZAWA T., HINTON G., S. KIYOHRO, & LANG K. J. 1989. Phoneme recognition using time-delay neural networks. *IEEE Trans. on Acoust., Speech, Signal Processing* 37.328–339.
- WOOTERS, C., 1993. *Lexical Modeling in a Speaker Independent Speech Understanding System*. Berkeley, CA: University of California dissertation.
- ZAVALIAGKOS, G., Y. ZHAO, R. SCHWARTZ, & J. MAKHOUL. 1994. A hybrid segmental neural net/hidden Markov model system for continuous speech recognition. *IEEE Transactions on Speech and Audio Processing* 2.151–160.

# EXHIBIT D



# NONLINEAR DISCRIMINANT FEATURE EXTRACTION FOR ROBUST TEXT-INDEPENDENT SPEAKER RECOGNITION

*Yochai Konig, Larry Heck, Mitch Weintraub, and Kemal Sonmez*

Speech Technology and Research Laboratory  
SRI International  
Menlo Park, CA 94025

## RÉSUMÉ

Cet article propose une méthode basée sur l'analyse discriminative non-linéaire pour extraire et sélectionner un ensemble de vecteurs acoustiques utilisés pour l'identification de locuteurs. L'approche consiste à mesurer et grouper un grand nombre de mesures acoustiques (correspondant à plusieurs trames de données consécutives), et à réduire la dimensionalité du vecteur résultant au moyen d'un réseau de neurones artificielles. Le critère utilisé pour optimiser les poids du réseau consiste à maximiser une mesure de la séparation entre les locuteurs d'une base de données d'apprentissage. L'architecture du réseau est telle que l'une de ses couches intermédiaires représente la projection des vecteurs acoustiques d'entrée sur un espace de dimensionalité inférieure. Après la phase d'apprentissage, cette partie du réseau peut être isolée et utilisée pour projeter les vecteurs acoustiques d'une base de données de test. Les vecteurs acoustiques projetés peuvent alors être classifiés. Combiné à un classificateur cepstral, le classificateur utilisant ces nouveaux vecteurs acoustiques réduit de 15% le taux d'erreur de classification de la base de données définie par NIST en 1997 pour l'évaluation des systèmes de reconnaissance du locuteur.

## ABSTRACT

We study a nonlinear discriminant analysis (NLDA) technique that extracts a speaker-discriminant feature set. Our approach is to train a multilayer perceptron (MLP) to maximize the separation between speakers by nonlinearly projecting a large set of acoustic features (e.g., several frames) to a lower-dimensional feature set. The extracted features are optimized to discriminate between speakers and to be robust to mismatched training and testing conditions. We train the MLP on a development set and apply it to the training and testing utterances. Our results show that by combining the NLDA-based system with a state of the art cepstrum-based system we improve the speaker verification performance on the 1997 NIST Speaker Recognition Evaluation set by 15% in average compared with our cepstrum-only system.

## 1. INTRODUCTION

Our goal is to extract and select features that are more invariant to non-speaker-related conditions such as handset type, sentence content, and channel effects. Such features will be robust to mismatched training and testing conditions of speaker verification systems. With current feature sets (e.g., cepstrum) there is a big performance gap between matched and mismatched tests [8] even after applying standard channel compensation techniques [4]. In order to find these features, the feature extraction step should be directly optimized to increase discrimination between speakers, and to filter out the non-relevant information.

Our proposed solution is to train a multilayer perceptron (MLP) to nonlinearly project a large set of acoustic features to a lower-dimensional feature set, such that it maximizes speaker separation. We train the MLP on a development set that includes several realizations of the same speakers under different conditions. We then apply the learned transformation (MLP in feed-forward mode) to the training and testing utterances. Finally, we use the resulting features for training the speaker recognition system, e.g., Bayesian adapted Gaussian mixture system [9].

We begin by reviewing related studies in Section 2. We describe the proposed feature extraction technique in Section 3. The Development database is described in Section 4. In Section 5, we report the experimental results on the 1997 NIST evaluation set. We continue with analysis of the results in Section 6. Finally, we conclude and describe directions for future work in Section 7.

## 2. RELATED STUDIES

The related studies to the NLDA technique can be divided into two main categories: robust speaker verification systems, and data-driven feature extraction techniques. Previously proposed approaches to increase robustness to mismatched training and testing conditions, especially to handset variations, include handset-dependent background

models [3], and a handset-dependent score normalization procedure known as Hnorm [9]. Data-driven feature extraction techniques were mainly suggested for speech recognition tasks. Rahim, Bengio and LeCun suggested optimizing a set of parallel class specific (e.g., phones) networks performing feature transformation based on minimum classification (MCE) criterion [7]. Fontaine, Ris and Boite used 2-hidden layer MLP to perform NLDA for isolated word, large vocabulary speech recognition task [2]. The training criterion for the MLPs was phonetic classification. Bengio and his colleagues suggested a global optimization of a neural network-hidden Markov (HMM) hybrid, where the outputs of the neural network constitute the observation sequence for the HMM [1].

### 3. NONLINEAR DISCRIMINANT ANALYSIS (NLDA)

We explore a nonlinear discriminant analysis (NLDA) technique that finds a nonlinear projection of the original feature space into a lower dimensional space that maximizes speaker recognition performance. This maximization problem can be expressed as

$$A^* = \underset{A}{\operatorname{argmax}} J\{A(X)\} \quad (1)$$

Where  $A(X)$  is a nonlinear projection of the original feature space  $X$  onto a lower dimensional space, and  $J\{\}$  is a closed-set speaker identification performance measure. To find the best  $A$  we train a 5 layer multilayer perceptron (MLP) to discriminate between speakers in a carefully selected development set (as described below). The MLP is constructed from a large input layer, a first large nonlinear hidden unit, a small (“bottleneck”) second linear hidden layer, a large third nonlinear hidden layer, and a softmax output layer (see Figure 1). The idea is that  $A$  is the projection of the input features speaker onto the “bottleneck” layer. After training the 5-layer MLP (denoted ‘MLP5’) we can remove the last hidden layer and the output layer, and use the remaining 3-layer MLP to project the target speaker data. Then, we use the transformed features to train the speaker verification system, for example, a Bayesian adapted GMM system (see Figure 2). The underlying assumption is that the transformation as found in the development set will be invariant across different speaker populations.

### 4. DEVELOPMENT DATABASE

To train the 5-layer MLP, we chose 855 Switchboard sentences (about 2 hours) from 31 speakers with a balanced mix of carbon and electret handsets, and balanced across gender. The input consists of 17 cepstral coefficients

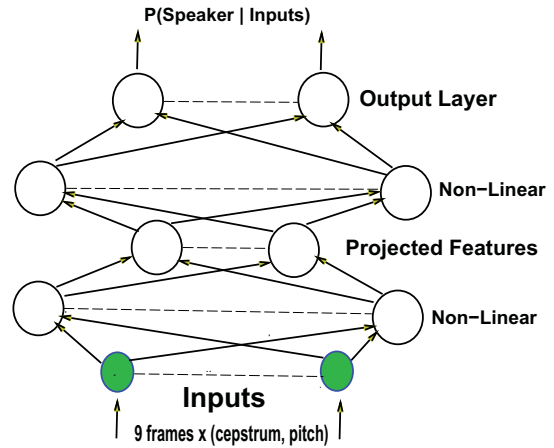


Figure 1: MLP5 for Speaker Recognition

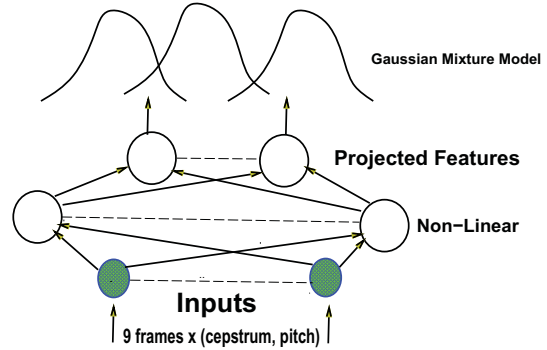


Figure 2: MLP3 for Feature Transformation

and an estimate of the pitch for the current frame, four past frames and four future frames, resulting in a 162-dimension vector. The first hidden layer has 500 sigmoidal units, the bottleneck layer has 34 linear units, the third hidden layer has 500 sigmoidal units, and a softmax output layer has 31 outputs (one for each speaker in the development set). After training the MLP5, we chopped the upper two layers. The resulting MLP (‘MLP3’) has one hidden layer and was used to transform the data of the target and impostor speakers in a test set as described above.

### 5. EXPERIMENTAL RESULTS

We used the 1997 NIST Speaker Recognition Evaluation corpus for testing. We report results for three different systems: (1) our best cepstrum system, which is our implementation of the state of the art in text independent speaker verification systems [6]) with 33 input features comprised of 10 cepstral coefficients, energy term, and first and second time derivatives (2) the NLDA based system described in this paper, (3) a combination of the cepstrum and the

Test	Cepstrum	NLDA	Combined
female 3	18.4%	23.0%	16.7%
female 10	12.1%	14.6%	10.8%
female 30	10.5%	12.4%	9.0%
male 3	14.9%	19.4%	14.4%
male 10	13.2%	12.9%	11.1%
male 30	7.9%	11.0%	7.1%

Table 1: Equal Error Rate (EER) Results of the 1997 NIST Eval., 1h condition

Test	Cepstrum	NLDA	Combined
female 10	13.5%	17.0%	12.5%
male 10	11.3%	14.4%	10.5%

Table 2: Equal Error Rate (EER) Results of the 1997 NIST Eval., 1s condition

NLDA systems. The third system is a linear combination of the normalized scores with weights of 0.7 for the cepstrum system scores and 0.3 for the NLDA system scores (except for the 3 second cases, where we used 0.6 for the cepstrum system and 0.4 for the NLDA system). We use the equal error rate (EER) between misses and false alarms as a performance measure for reporting results. In Table 1, we summarize the results for the 1h condition in the NIST evaluation. In this condition the training consists of 2 phone calls from the same handset, each 1 minute in duration. There are three different test lengths: 3, 10, and 30 seconds. We report the results for each gender separately, by pooling all the test data together (matched and mismatched telephone number).

The results show a consistent win for the combined system over our state of the art cepstrum system. We observe the same consistent win for another condition, 1s, in the 1997 NIST Speaker Recognition Evaluation as demonstrated for the 10 second case in Table 2, and across all regions of the DET (false alarm probability versus miss probability) curves as illustrated in Figure 3 for the male, 10 seconds (1h condition) for the cepstrum only system and the combined system. These results are consistent with our initial results for the 1998 Evaluation corpus.

## 6. RESULT ANALYSIS

In this section, we examine our “black box” approach, provide insight to its success and give directions for potential improvements. In order to examine the importance of the pitch input, the 9 frame temporal window, and the degradation loss as a result of the dimension reduction

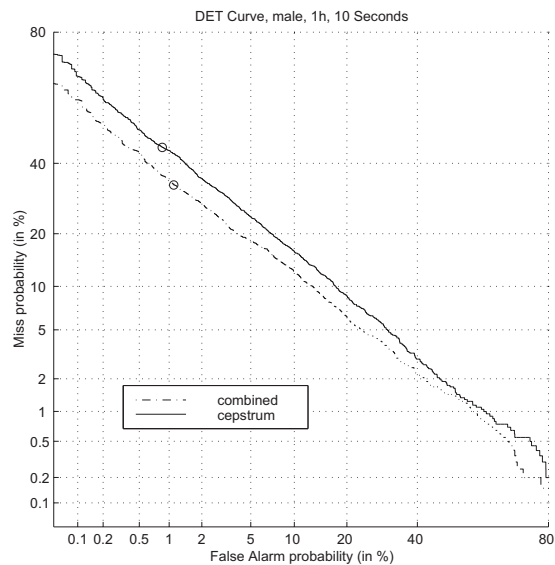


Figure 3: DET Curve for male, 1h, 10 seconds

Inputs	Name	Frame Correct
9 frames + pitch	MLP3	37.2%
9 frames + pitch	MLP5-34	28.9%
9 frames + pitch	MLP5-50	29.0%
9 frames, no pitch	MLP5-NO	25.9%
1 frame + pitch	MLP5-1fr	18.6%

Table 3: Frame-level results on the cross-validation set

from 162 inputs to 34 hidden units in the bottleneck layer, we trained several MLPs and tested their cross-validation, frame-level performance on a close set speaker recognition (our development set as described above). In the development phase we found a strong correlation between these frame-level results and the “full cycle” results of the speaker verification system. The results are summarized in Table 3.

We trained two types of MLPs: a 5-layer MLP, and a “vanilla” MLP with three layers including one hidden layer (denoted ‘MLP3’). As mentioned above there were 31 speakers in our development set, 687156 frames for training and 77904 for cross-validation. Our baseline MLP is the MLP5 described above with 162 inputs and 3 hidden layers with 500, 34, and 500 units (named ‘MLP5-34’). The output layer of all our nets has 31 outputs, one output for each speaker in our development set. The MLP5 named ‘MLP5-NO’ is the same as the baseline but without pitch information (only 153 inputs). The MLP5 named ‘MLP5-1’ is the same as the baseline but with only one input frame (as compared to the 9 frames used in the other systems)

Training a 5-layer MLP is difficult given the complex

nonlinear error surface and requires a lot of training data preferably a ratio of at least 10 between frames than free parameters. In these experiments the ratio was around 4.7 (700k frames to 150k parameters). This might explain the disparity in performance between the MLP3 to the MLP5. This is not due to the bottleneck size as shown by the result of the MLP5 named ‘MLP5-50’ (the same as ‘MLP5-34’ but with 50 hidden units in the bottleneck layer). In our speech recognition experiments [5] with NLDA, with the right ratio between frames to free parameters, we did not observe any performance loss because of the dimension reduction at the bottleneck layer. Thus, we plan to increase the size of the development set and hopefully improve the performance of the MLP5 and the overall technique. Additionally comparing the second row to the fourth and fifth rows in Table 3, we observe from these results that that we get a 3% absolute gain from the pitch information, and 10.3% absolute gain from the temporal window.

Another set of interesting results is the correlation between the cepstrum and the NLDA scores on 1997 Eval. set, 1h condition, as summarized in Table 4. From these results, we observe that the NLDA technique contribute a significant amount of new information, especially for the shorter test lengths. This is consistent with the results previously shown in Table 1.

Test Length	Male	Female
3	0.61	0.47
10	0.68	0.71
30	0.76	0.77

Table 4: Correlation Coefficients between NLDA and Cepstrum systems on 1997 Eval. set, 1h condition

## 7. CONCLUSIONS AND FUTURE WORK

We presented a nonlinear discriminant analysis (NLDA) technique that extracts a speaker-discriminant feature set. Our results on the 1997 NIST evaluation show a consistent (across 12 different tests) and significant (around 15% in relative error) improvement when combining the system trained with the NLDA features with cepstrum based system. Our initial results on 1998 NIST evaluation are consistent with 1997 results. Furthermore, our analysis suggests that there is a potential for performance improvement given more development data. We also plan to experiment with other types of input data such as speech over cellular phones and speaker-phone speech. In addition, we plan to extend this study by using a wider range of input representations and resolutions such as first and second derivatives of cepstrum, filter-bank energy levels, and

different analysis windows. Finally we want to note that although the training of the MLP with 5 layers is computationally expensive (25 x real time), the application of the MLP3 in a feed forward mode is very fast (less than 0.4 real-time), thus the NLDA approach is feasible in realistic settings.

## 8. REFERENCES

- [1] Y. Bengio, R. De Mori, G. Flammia, and R. Kompe. Global optimization of a neural network-hidden Markov model hybrid. *IEEE trans. on Neural Networks*, 3(2):252–258, March 1992.
- [2] V. Fontaine, C. Ris, and J. M. Boite. Nonlinear discriminant analysis for improved speech recognition. In *Proceedings European Conf. on Speech Communication and Technology. (EUROSPEECH)*, Rhodes, Greece, 1997.
- [3] L. P. Heck and M. Weintraub. Handset-dependent background models for robust text-independent speaker recognition. In *Proceedings International Conference on Acoustics Speech and Signal Processing (ICASSP)*, Munich, Germany, 1997.
- [4] H. Hermansky, N. Morgan, A. Bayya, and P. Kohn. Compensation for the effect of the communication channel in auditory-like analysis of speech (RASTAPLP). *Proceedings European Conf. on Speech Communication and Technology. (EUROSPEECH)*, pages 1367–1370, 1991.
- [5] Y. Konig and M. Weintraub. Acoustic modeling session - SRI site presentation. In *NISTLVCSR Workshop*, Linthicum Heights, MD, October 1996.
- [6] NIST. Result summary. In *Speaker Recognition Workshop Notes*, Linthicum Heights, Maryland, 1997.
- [7] M. Rahim, Y. Bengio, and Y. LeCun. Discriminative feature and model design for automatic speech recognition. In *Proceedings European Conf. on Speech Communication and Technology. (EUROSPEECH)*, Rhodes, Greece, 1997.
- [8] D. A. Reynolds. The effects of handset variability on speaker recognition performance experiments on the switchboard corpus. In *Proceedings International Conference on Acoustics Speech and Signal Processing (ICASSP)*, Atlanta, GA, 1996.
- [9] D. A. Reynolds. Comparison of background normalization methods for text-independent speaker verification. In *Proceedings European Conf. on Speech Communication and Technology. (EUROSPEECH)*, Rhodes, Greece, 1997.

# EXHIBIT E



# EXPLICIT WORD ERROR MINIMIZATION IN N-BEST LIST RESCORING

Andreas Stolcke

Yochai König

Mitchel Weintraub

Speech Technology and Research Laboratory

SRI International, Menlo Park, CA, U.S.A.

<http://www.speech.sri.com/>

{stolcke, konig, mw}@speech.sri.com

## ABSTRACT

We show that the standard hypothesis scoring paradigm used in maximum-likelihood-based speech recognition systems is not optimal with regard to minimizing the word error rate, the commonly used performance metric in speech recognition. This can lead to sub-optimal performance, especially in high-error-rate environments where word error and sentence error are not necessarily monotonically related. To address this discrepancy, we developed a new algorithm that explicitly minimizes expected word error for recognition hypotheses. First, we approximate the posterior hypothesis probabilities using N-best lists. We then compute the expected word error for each hypothesis with respect to the posterior distribution, and choose the hypothesis with the lowest error. Experiments show improved recognition rates on two spontaneous speech corpora.

## 1. INTRODUCTION

The standard selection criterion for speech recognition hypotheses aims at maximizing the posterior probability of a hypothesis  $W$  given the acoustic evidence  $X$  [1]:

$$\begin{aligned} W^* &= \operatorname{argmax}_W P(W|X) \\ &= \operatorname{argmax}_W \frac{P(W)P(X|W)}{P(X)} \end{aligned} \quad (1)$$

$$= \operatorname{argmax}_W P(W)P(X|W) \quad (2)$$

Here  $P(W)$  is the prior probability of a word sequence according to a *language model*, and  $P(X|W)$  is given by the acoustic model. Equation (1) is Bayes' Rule, while (2) is due to the fact that  $P(X)$  does not depend on  $W$  and can therefore be ignored during maximization. Bayes decision theory (see, e.g., [2]) tells us that this criterion (assuming accurate language and acoustic models) maximizes the probability of picking the correct  $W$ ; i.e., it minimizes *sentence* error rate. However, speech recognizers are usually evaluated primarily for their *word* error rates.

Empirically, sentence and word error rates are highly correlated, so that minimizing one tends to minimize the other. Still, if only for theoretical interest, two questions arise:

- (A) Are there cases where optimizing expected word error and expected sentence error produce different results?
- (B) Is there an effective algorithm to optimize expected word error explicitly?

Note that question (A) is not about the difference between word and sentence error in a particular instance of  $X$  and its correct transcription, since obviously the two error criteria would likely pick different best hypotheses in any given instance. Instead, we are concerned with the *expected* errors, as they would be obtained by averaging over many instances of the same acoustic evidence with varying true word sequences, i.e., if we sampled from the true posterior distribution  $P(W|X)$ .

We will answer question (A) first by way of a constructed example, showing that indeed the two error metrics can diverge in their choice of the best hypothesis. Regarding question (B), we develop a new N-best rescoring algorithm that explicitly estimates and minimizes word error. We then verify that the algorithm produces lower word error on two benchmark test sets, thus demonstrating that question (A) can be answered in the affirmative even for practical purposes.

## 2. AN EXAMPLE

The following is a hypothetical list of recognition outputs with attached (true) posterior probabilities.

$w_1$	$w_2$	$P(w_1 w_2   X)$	$P(w_1   X)$	$P(w_2   X)$	$E[\text{correct}]$
a	d	.0	.44	.4	.84
a	e	.24	.44	.34	.78
a	f	.2	.44	.26	.7
b	d	.2	.26	.4	.66
b	e	.05	.26	.34	.6
b	f	.01	.26	.26	.52
c	d	.2	.3	.4	.7
c	e	.05	.3	.34	.64
c	f	.05	.3	.26	.56

For simplicity we assume that all hypotheses consist of exactly two words,  $w_1$  and  $w_2$ , shown in the first two columns. The third column shows the assumed joint posterior probabilities  $P(w_1 w_2 | X)$  for these hypotheses. Columns 4 and 5 give the posterior probabilities  $P(w_1 | X)$  and  $P(w_2 | X)$  for individual words. These posterior word probabilities follow from the joint posteriors but summing over all hypotheses that share a word in a given position. For example, the posterior  $P(w_1 = a | X)$  is obtained by summing

$P(w_1 w_2 | X)$  of all hypotheses such that  $w_1 = a$ . Column 6 shows the expected number of correct words  $E[\text{correct}]$  in each hypothesis, under the assumed posterior distribution. This is simply the sum of  $P(w_1 | X)$  and  $P(w_2 | X)$ , since

$$\begin{aligned} E[\text{words correct}(w_1 w_2) | X] \\ &= E[\text{correct}(w_1) | X] + E[\text{correct}(w_2) | X] \\ &= P(w_1 | X) + P(w_2 | X) \end{aligned}$$

As can be seen, although the first hypothesis (“a d”) has posterior 0, it has the highest expected number of words correct, i.e., the minimum expected word error. Thus, we have shown by construction that optimizing overall posterior probability (sentence error) does not always minimize expected word error. Of course the example was constructed such that two words that each have high posterior probability happen to have low (i.e., zero) probability when combined. Note that this is not unrealistic: for example, the language model could all but “prohibit” certain word combinations.

Furthermore, we can expect the discrepancy between word and sentence error to occur more at high error rates. When error rates are low, i.e., when there are at most one of two word errors per sentence, each word error corresponds to a sentence error and vice-versa. Thus, if we had an algorithm to optimize the expected word error directly, we would expect to see its benefits mostly at high error rates.

### 3. THE ALGORITHM

We now give an algorithm that minimizes the expected word error rate (WER) in the N-best rescoring paradigm [5]. The algorithm has two components: (1) approximating the posterior distribution over hypotheses and (2) computing the expected WER for N-best hypotheses (and picking the one with lowest expected WER).

#### 3.1. Approximating posterior probabilities

An estimate of the posterior probability  $P(W | X)$  of a hypothesis  $W$  can be derived from Equation (1), with modifications to account for practical limitations:

- The true distributions  $P(W)$  and  $P(X | M)$  are replaced by their imperfect counterparts, the language model probability  $P_{\text{LM}}(W)$  and the acoustic model likelihood  $P_{\text{AC}}(X | W)$ .
- The dynamic range of the acoustic model, due to unwarranted independence assumptions, needs to be attenuated by an exponent  $1/\lambda$  ( $\lambda$  is the language model weight commonly used in speech recognizers, and optimized empirically).
- The normalization term

$$P(X) = \sum_W P(W) P(X | W)$$

is replaced by a finite sum over all the hypotheses in the N-best list. This is not strictly necessary for the algorithm since it is invariant to constant factors on the posterior estimates, but it conveniently makes these estimates sum to 1.

Let  $W_i$  be the  $i$ th hypothesis in the  $N$ -best list; the posterior estimate is thus

$$P(W_i | X) \approx \frac{P_{\text{LM}}(W_i) P_{\text{AC}}(W_i | X)^{\frac{1}{\lambda}}}{\sum_{k=1}^N P_{\text{LM}}(W_k) P_{\text{AC}}(W_k | X)^{\frac{1}{\lambda}}}$$

This N-best approximation to the posterior has previously been used, e.g., in the computation of posterior word probabilities for keyword spotting [7].

#### 3.2. Computing expected WER

Given a list of N-best hypotheses and their posterior probability estimates, we approximate the expected WER as the weighted average word error relative to all the hypotheses in the N-best list. That is, we consider each of the  $N$  hypotheses in turn as the “truth” and weight the word error counts from them with the corresponding posterior probability:

$$E[\text{WE}(W) | X] \approx \sum_{i=1}^N P(W_i | X) \text{WE}(W | W_i) \quad (3)$$

where  $\text{WE}(W | W_i)$  denotes the word error of  $W$  using  $W_i$  as the reference string (computed in the standard way using dynamic programming string alignment).

#### 3.3. Computational Complexity

Rescoring  $N$  hypotheses requires  $N^2$  word error computations, which can become quite expensive for N-best lists of 1000 or more hypotheses. We found empirically that the algorithm very rarely picks a hypothesis that is not within the top 10 according to posterior probability. This suggests a shortcut version of the algorithm that only computes expected word error for the top  $K$  hypotheses, where  $K \ll N$ . Note that we still need to consider all  $N$  hypotheses to compute the expected word error as in Equation (3), otherwise these estimates become very poor and affect the final result noticeably. The practical version of our algorithm thus has complexity  $O(KN)$ .

#### 3.4. Other knowledge sources and weight optimization

Often other knowledge sources are added to the standard language model and acoustic scores to improve recognition, such as word transition penalties or scores expressing syntactic or semantic well-formedness (e.g., [4]). Even though these additional scores cannot always be interpreted as probabilities, they can still be combined with exponential weights; the weights are then optimized on a held-out set to minimize WER [5].

This weight optimization should not be confused with the word error minimization discussed here; instead, the two methods complement each other. The additional knowledge sources can be used to yield improved posterior probability estimates, based on which the algorithm described here can be applied. In this scheme, one should first optimize the language model and other knowledge source weights to achieve the best posterior probability estimates (e.g., by minimizing empirical *sentence* error).

	WER	SER
<b>Switchboard</b>		
Standard rescoring	52.7	84.0
WER minimization	52.2	84.4
<b>CallHome Spanish</b>		
Standard rescoring	68.4	80.9
WER minimization	67.8	81.2

Table 1. Word (WER) and Sentence error rates (SER) of standard and word-error-minimizing rescoring methods

So far, we have not implemented combined weight and word error optimization. The experiments reported below used standard language model weights and word transition penalties that had previously been determined as near-optimal in the standard recognition paradigm.

#### 4. EXPERIMENTS

We tested the new rescoring algorithm on 2000-best lists for two test sets taken from spontaneous speech corpora. Test set 1 consisted of 25 conversations from the Switchboard corpus [3]. Test set 2 were 25 conversations from the Spanish CallHome corpus collected by the Linguistic Data Consortium. Due to the properties of spontaneous speech, error rates are relative high on these data, making word error minimization more promising, as discussed earlier.

The results for both standard rescoring and WER minimization are shown in Table 1. On both test sets the WER was reduced by about 0.5% (absolute) using the word error minimization method. A per-sentence analysis of the differences in word error show that the improvement is highly significant in both cases (Sign test  $p < 0.0005$ ). Note that, as expected, the sentence error rate (SER) increased slightly, since we no longer were trying to optimize that criterion.

For comparison, we also applied our algorithm to the 1995 ARPA Hub3 development test set. This data yields much lower word error rates, between 10% and 30%. In this case the algorithm invariably picked the hypothesis with the highest posterior probability estimate, confirming our earlier reasoning that word error minimization was less likely to make a difference at lower error rates.

#### 5. DISCUSSION AND CONCLUSION

We have shown a discrepancy between the classical hypothesis selection method for speech recognizers and the goal of minimizing word error. A new N-best rescoring algorithm has been proposed that corrects this discrepancy by explicitly minimizing expected word error (as opposed to sentence error) according to the posterior distribution of hypotheses. Experiments show that the new algorithm results in small, but consistent (and statistically significant) reductions in word error under high error rate conditions.

In our experiments so far, the improvement in WER is small. However, the experiments confirm that the theoretical possibility of suboptimal WER using the standard

rescoring approach is manifest in practice. An important aspect of the WER minimization algorithm is that it can use other, more sophisticated posterior probability estimators, with the potential for larger improvements. Our experiments so far have been based on the commonly used acoustic and language model scores, but we are already experimenting with more complex posterior estimator methods based on neural network models [6].

#### REFERENCES

- [1] L. R. Bahl, F. Jelinek, and R. L. Mercer. A maximum likelihood approach to continuous speech recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 5(2):179–190, 1983.
- [2] R. O. Duda and P. E. Hart. *Pattern Classification and Scene Analysis*. John Wiley, New York, 1973.
- [3] J. J. Godfrey, E. C. Holliman, and J. McDaniel. SWITCHBOARD: Telephone speech corpus for research and development. In *Proceedings of the IEEE Conference on Acoustics, Speech and Signal Processing*, vol. I, pp. 517–520, San Francisco, 1992.
- [4] R. Moore, D. Appelt, J. Dowding, J. M. Gawron, and D. Moran. Combining linguistic and statistical knowledge sources in natural language processing for ATIS. In *Proceedings ARPA Spoken Language Systems Technology Workshop*, pp. 261–264, Austin, Texas, 1995.
- [5] M. Ostendorf, A. Kannan, S. Austin, O. Kimball, R. Schwartz, and J. R. Rohlicek. Integration of diverse recognition methodologies through reevaluation of N-best sentence hypotheses. In *Proceedings DARPA Speech and Natural Language Processing Workshop*, pp. 83–87, Pacific Grove, CA, 1991. Defense Advanced Research Projects Agency, Information Science and Technology Office.
- [6] M. Weintraub, F. Beaufays, Z. Rivlin, Y. Konig, and A. Stolcke. Neural-network based measures of confidence for word recognition. In *Proceedings of the IEEE Conference on Acoustics, Speech and Signal Processing*, vol. II, pp. 887–890, Munich, 1987.
- [7] M. Weintraub. LVCSR log-likelihood ratio rescoring for keyword spotting. In *Proceedings of the IEEE Conference on Acoustics, Speech and Signal Processing*, vol. I, pp. 297–300, Detroit, 1995.



# EXHIBIT F

# DYNAMO: An Algorithm for Dynamic Acoustic Modeling

*Françoise Beaufays, Mitch Weintraub, Yochai Konig*

Speech Technology and Research Laboratory  
SRI International  
Menlo Park, CA 94025

## ABSTRACT

This paper summarizes part of SRI's effort to improve acoustic modeling in the context of the Large Vocabulary Continuous Speech Recognition (LVCSR) project. It concentrates on two problems that are believed to contribute to the large error rates observed with LVCSR databases: (1) the lack of discriminative power of the speech models in the acoustic space, and (2) the discrepancy between the criterion used to train the models (typically frame-level maximum likelihood) and the task expected from the models (word-level recognition).

We address the first issue by searching for features that help in narrowing the model distributions, and by proposing a neural-network-based architecture to combine these features. The neural networks (NNET) are used in association with a set of large Gaussian mixture models (GMM) whose mixture weights are dynamically estimated by the neural networks, for each frame of incoming data. We call the resulting algorithm DYNAMO, for dynamic acoustic modeling. To address the second problem, we propose two discriminative training criteria, both defined at the sentence level. We report preliminary results with the Spanish Callhome database.

## 1. Introduction

Many factors contribute to the relatively low performance of state-of-the-art speech recognizers operating on spontaneous, telephone speech. A few of these factors are: the diversity of speakers and speaking styles, the typically relaxed articulation, the multitude of pronunciation variants, the presence of extraneous noises, the superposition of more than one voice in some segments, and the distortion due to the communication channel. Whereas some of these factors can be efficiently dealt with by explicit modeling (*e.g.* vocal tract normalization (*e.g.* [AKC94]), pronunciation modeling (*e.g.* [Slo95, FW97])), many others are left for the acoustic models's multi-modal distributions to model implicitly. This, however, has the well-known result of broad overlapping distributions which often lead to recognition errors.

In this context, identifying features that act as discriminants in the acoustic space would be useful to narrow the acoustic distributions. If such features can be found, the problem becomes how to use them, and how to ensure that sufficient data sharing is allowed for the model parameters to be reliably estimated. These are the main issues that motivated this work.

In the past decade, contextual linguistic features have been widely used in conjunction with decision tree models, and have significantly improved recognition performance (*e.g.* [BdSG<sup>+</sup>91, YOW94]). Decision trees, however, make data sharing among different states difficult, and are not well suited to the use of features that are continuous

in nature, as opposed to binary. For these reasons, we chose instead to base our models on neural networks.

More recently, Ostendorf *et al.* [OBB<sup>+</sup>97] showed that a combination of acoustic and prosodic features could greatly help identifying speech segments that were erroneously recognized (32% predictability improvement for a 10-hour training subset of Switchboard). Similar results were reported by various researchers working on confidence measures for word recognition (*e.g.* [WBR<sup>+</sup>97]). Presumably, some of these features, which include various measures of speaking rate, SNR, energy, fundamental frequency, stress pattern, and syllable position, could be directly used to disambiguate large acoustic distributions.

In the field of speaker recognition, the use of handset detectors has dramatically decreased recognition error rates by sorting out carbon button from electret handsets [Rey96, HW97]. The handset type could also be used as an input to the acoustic modeling algorithms.

Another important issue in acoustic modeling is how to capture the dynamics of the speech signal. Much research has recently been devoted to relaxing the independence assumption imposed by most hidden Markov modeling approaches (HMM) and to modeling the correlation between successive frames of data, leading to the family of so-called segment models [ODK96]. Without embarking in this level of complexity, and following a feature-based approach, we propose to include in the acoustic models time features similar to the time index proposed in [GN93, DASW94] and [KM94]. These features don't model correlation but they do alleviate the independence assumption.

Our goal here is to explore the usefulness of such knowledge sources as acoustic discriminants, and to propose an efficient and robust architecture to incorporate them in the acoustic models. Clearly, the richness of the acoustic space representation will have a strong influence on how far this approach can be pushed, but the success of the experiments cited above (handset classification, feature-based error prediction, etc.) indicate that the cepstrum-based representation that most systems use offers enough flexibility for the acoustic models to be significantly improved.

As mentioned before, the architecture we propose relies on neural networks. An important issue related to this choice is the selection of a training criterion to optimize the weights of the networks. The desirable properties for this criterion are (1) to be discriminative, (2) to be closely related to the metric used to evaluate the performance of the recognizer (typically the word error rate (WER)), and (3) to be differentiable with respect to the weights of the neural networks.

Not all the above issues will be discussed in the paper since this

System	Eval '95	Eval '96
baseline	71.00	65.22
+ DT	67.77	64.37
+ CI (size: 1/16 DTs)	68.77	65.22
+ CI (size: 1/8 DTs)	68.27	65.22
+ CI (size: 1/4 DTs)	68.34	65.10
+ CI (size: 1/2 DTs)	67.98	64.49
+ CI (size: 1/1 DTs)	67.69	64.31
N-best error rate	52.54	/

Table 1: N-best list rescoring with decision tree models and context-independent phone models of different sizes: WER in %.

work is still in an early stage. Our first goals were to validate the architecture we propose and to investigate different discriminative training criteria. These two points will be addressed. Feature selection, however, will be the object of future work: for our preliminary experiments, we used a set generic knowledge sources including linguistic features and time indices.

## 2. Baseline System and Databases

The baseline system for this work is a speaker-independent continuous speech recognition system trained with 75 conversations of Callhome Spanish data and 80 conversations from Callfriend Spanish. It is based on continuous-density, genonic HMMs [DMM96], and uses a multipass recognition strategy [MBDW93] with a vocabulary of 8K words, non-cross-word acoustic models, and a bigram language model. N-best lists are generated, and rescored with the original acoustic models, a trigram language model, and additional acoustic models such as decision-tree-based cross-word models (DT) or large context-independent phone GMMs (CI).

## 3. Recognition with Large Context-Independent Models

Using the Spanish Callhome database, we conducted a series of N-best list rescoring experiments with decision tree models and with large context-independent GMMs. The numbers of Gaussians in the GMMs were chosen to be fractions of the numbers of Gaussians used in the corresponding decision tree models. The smallest models had 16 times fewer Gaussians than the decision tree models, and the largest models had exactly the same size. Recognition experiments were performed with two sets of 200 sentences selected at random from the male evaluation test sets of 1995 and 1996. The results, reported in Table 1, show that, for this database, context-independent models perform as well as or slightly better than decision tree models, provided that the numbers of parameters are equal.

## 4. The DYNAMO Algorithm

The architecture we propose is based on a hybrid system combining feedforward neural networks and context-independent phone models. Each phone is modeled with a large GMM whose mixture weights are dynamically estimated by a neural network (see Fig. 1), hence the name of the algorithm, DYNAMO. The means and variances of the GMMs are held constant. The inputs to the neural network are the knowledge sources discussed in the introduction. For each data frame, the knowledge sources for each phone are evaluated and input into the corresponding NNET. Each NNET outputs

a set of mixture weights, and the likelihood of the observed data is computed from the corresponding phone GMM.

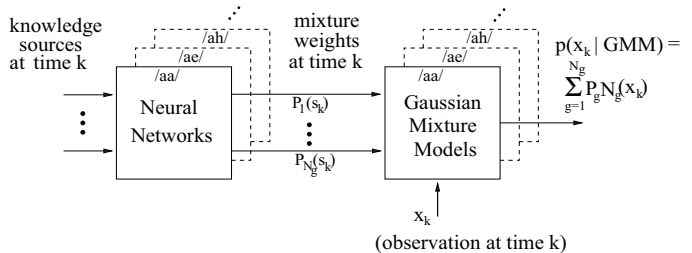


Figure 1: A hybrid NNET-GMM model for dynamic acoustic modeling.

Specifically, the likelihood of an observation,  $\mathbf{x}_k$ , with respect to phone  $\varphi$  is given by

$$p(\mathbf{x}_k | \mathcal{N}^\varphi, \mathcal{G}^\varphi) = \sum_{g=1}^{N_g} P_g^\varphi(\mathbf{s}_k^\varphi) N_g^\varphi(\mathbf{x}_k), \quad (1)$$

where  $\mathcal{N}^\varphi$  and  $\mathcal{G}^\varphi$  denote, respectively, the NNET and the GMM associated to phone  $\varphi$ ,  $N_g$  is the number of Gaussians in  $\mathcal{G}^\varphi$ ,  $N_g^\varphi(\cdot)$  and  $P_g^\varphi(\cdot)$  are, respectively, the  $g^{th}$  mixture component and the  $g^{th}$  mixture weight in  $\mathcal{G}^\varphi$ , and  $\mathbf{s}_k^\varphi$  represents the vector of knowledge sources for phone  $\varphi$ , at time  $k$ .

Because the mixture weights for each phone must sum to one, the training of the neural networks is a constrained optimization problem. To simplify the training procedure, we chose to hard-wire this constraint in the architecture of the neural networks by using a “softmax” output layer [Bri90]:

$$P_g(\mathbf{s}) = \frac{e^{y_g(\mathbf{s})}}{\sum_j e^{y_j(\mathbf{s})}}, \quad (2)$$

where  $y_g(\cdot)$  is the  $g^{th}$  output of the neural network, before the softmax layer.

The Gaussians in each phone model can be interpreted as a set of basis functions. A multimodal probability density function is then estimated for each observation by taking a linear combination of the basis functions, the weights of which are computed dynamically by the neural network. The discriminative emphasis of certain portions of the acoustic space at each instant has the effect of narrowing the distributions around the acoustic areas where the data are expected to lie.

This architecture thus outputs the likelihoods of the observations. This is in contrast with NNET-HMM hybrids trained for state classification [BM90], where the outputs are state posterior probabilities that need to be converted into likelihoods, and with approaches such as REMAP [BKM95, KBM96] that estimate global posterior probabilities of word sequences.

### 4.1. Training of the DYNAMO Models

The DYNAMO models are trained in two phases. First, the context-independent phone GMMs are trained with the expectation-maximization (EM) algorithm to maximize the log-likelihood of the

training data. The means and variances of these models are retained; the mixture weights are discarded. Then, the adaptive parameters of the neural networks are trained with the stochastic steepest descent algorithm to optimize some criterion  $\xi$ . The neural network weights are thus updated according to

$$\Theta_{n+1}^\varphi = \Theta_n^\varphi + \Delta\Theta_n^\varphi \quad (3)$$

$$\Delta\Theta_n^\varphi = \mu \hat{\nabla}_{\Theta_n^\varphi} \xi_\varphi, \quad (4)$$

where  $\Theta_n^\varphi$  denotes the set of neural network weights for phone  $\varphi$  at iteration  $n$ ,  $\hat{\nabla}_{\Theta_n^\varphi} \xi_\varphi$  is the instantaneous gradient of the optimization criterion  $\xi_\varphi$  for phone  $\varphi$ , and  $\mu$  is a constant that controls the learning rate.

Note that the optimization criterion  $\xi_\varphi$  does not need to be identical to the criterion used to train the GMMs (ML). Indeed, we argue in the next sections that discriminative training is better suited to this task. For now, however, we will assume for simplicity that  $\xi_\varphi$  is the average log-likelihood of the data,

$$\xi_\varphi = \sum_k \log p(\mathbf{x}_k | \mathcal{N}^\varphi, \mathcal{G}^\varphi), \quad (5)$$

where the sum is taken over all the observations  $\mathbf{x}_k$  aligned to phone  $\varphi$ .

Applying the chain rule to the derivatives of Eq. 5, and taking Eq. 2 into account, we find

$$\hat{\nabla}_{\Theta_n^\varphi} \xi_\varphi = \sum_j \frac{\partial \xi_\varphi}{\partial y_j} \frac{\partial y_j}{\partial \Theta_n^\varphi}, \quad (6)$$

where

$$\delta_j \triangleq \frac{\partial \xi_\varphi}{\partial y_j} = \frac{P_j^\varphi N_j^\varphi}{\sum_g P_g^\varphi N_g^\varphi} - P_j^\varphi \quad (7)$$

can be backpropagated through the neural network, as in the traditional backpropagation algorithm [RMT86].

Intuitively, the backpropagation term,  $\delta$ , for Gaussian  $j$  is large in absolute value if the posterior probability of the Gaussian is very different from its prior probability  $P_j$ , with both probabilities being functions of the knowledge sources for the current data frame.

To hasten the convergence of the neural networks and steer them away from uninteresting local minima, we initially set their weights so that the network outputs are equal to the mixture weights estimated with the EM algorithm.

## 5. Recognition Experiments with ML-trained Dynamo Models

We performed a set of rescoreing experiments with ML-trained DYNAMO models, using linguistic questions and, in some experiments, time features. We chose the linguistic features to be identical to those selected by the decision trees in previous DT-rescoreing experiments (Table 1). The time features for a hypothesized phone aligned to  $T$  frames of data were the phone duration,  $T$ , and the relative time index  $t/T$ , where  $t = 0 \dots T - 1$ .

Results are given in Table 2, where the baseline obtained by rescoreing the N-best lists with the GMMs is given for comparison. These

GMM size	Experiment	WER
$\times 1/16$	no NNETs – baseline	68.77
$\times 1/16$	NNETs w/ ling. feat. & time feat.	69.20
$\times 1/16$	NNETs w/ ling. feat. only	68.92
$\times 1/8$	no NNETs – baseline	68.27
$\times 1/8$	NNETs w/ ling. feat. & time feat.	69.35

Table 2: Rescoreing experiments with ML-trained DYNAMO models: WER in %.

numbers show that the introduction of the ML-trained networks increased the overall WER. Further analysis of the results revealed that the likelihood of the test data had increased as a result of training but that the posterior probabilities of the correct models had decreased. This indicated that competing models scored higher than the correct model, which confirmed that discriminative training should be used instead.

## 6. Discriminative Training Criteria

Discriminative training of speech models was first introduced by Bahl *et al.* under the form of Maximum Mutual Information (MMI) estimation [BBdSM86]. In this framework, the speech models are trained to maximize the mutual information between the observation sequence  $\mathcal{X} = [\mathbf{x}_1, \dots, \mathbf{x}_k, \dots, \mathbf{x}_N]$  and the correct word sequence  $W_c$ :

$$\Theta^* = \arg \max_{\Theta} I_{\Theta}(W_c, \mathcal{X}), \quad (8)$$

with

$$I_{\Theta}(W_c, \mathcal{X}) = \frac{p(\mathcal{X}, W_c)}{p(\mathcal{X})p(W_c)} = \frac{p(\mathcal{X}|W_c)}{\sum_W p(\mathcal{X}|W)p(W)}, \quad (9)$$

where the sum in the denominator is taken over all possible word sequences,  $W$ .

Practical implementations of Eq. 9 for continuous speech recognition include the estimation of the denominator with a phone loop model [Mer88], and its approximation by a sum over the hypotheses in an N-best list [Cho90].

The first optimization criterion we propose is similar to the N-best list implementation of MMI, but differs in that we augment the N-best list with the correct word sequence,  $W_c$ . We then maximize the posterior probability of the correct word sequence,

$$P(W_c | \mathcal{X}) = \frac{p(\mathcal{X}|W_c)p(W_c)}{p(\mathcal{X}|W_c)p(W_c) + \sum_{h=1}^{N_h} p(\mathcal{X}|W_h)p(W_h)}, \quad (10)$$

where  $N_h$  is the N-best list depth. The inclusion of the joint probability of the observation and the correct word sequence in the denominator makes the criterion depart from the original MMI but has a useful property in terms neural network training, as we will show.

Another family of discriminative criteria stems from the motivation of directly optimizing the metric used to evaluate the recognizer, *i.e.* the word error rate. Bahl *et al.* proposed the heuristic ‘‘corrective training’’ procedure in [BBdSM88]. Katagiri *et al.* developed the Generalized Probabilistic Descent method that extends the idea of Bayes optimum classification by introducing smooth classification

error functions, and generalizes this framework to the classification of patterns of variable lengths [KLJ91].

The second criterion we propose consists in minimizing the average number of errors over the N-best list,

$$\text{ANER}(\mathcal{X}) = \frac{1}{N_h} \sum_{h=1}^{N_h} \text{NER}(W_h) P(W_h|\mathcal{X}), \quad (11)$$

where  $\text{NER}(W_h)$  denotes the number of errors in the  $h^{\text{th}}$  hypothesis, and  $P(W_h|\mathcal{X})$  is the posterior probability of the  $h^{\text{th}}$  hypothesis in the (non-augmented) N-best list.

Both criteria are optimized in a stochastic optimization framework, as we will discuss shortly. In both cases, the training procedure requires N-best lists for all the training data. This is typically quite costly but not infeasible, especially if the N-best list depth is limited to a small number of hypotheses (5 or 10).

### 6.1. Maximizing the posterior probability of the correct sentence

Let  $p(i)$  denote the joint probability of a word sequence  $i$  (reference or hypothesis) and of the corresponding acoustic sequence,

$$p(i) = p_{LM}(i) p_{AM}(i)^{1/\lambda}, \quad (12)$$

where  $p_{LM}(i)$  and  $p_{AM}(i)$  are shorthands for the language model and acoustic model probabilities,  $p(W_i)$  and  $p(\mathcal{X}|W_i)$ , respectively, and where  $\lambda$  is the language model weight.

With this notation, we can rewrite the posterior probability of the correct word sequence in Eq. 10 as

$$P(c) = \frac{p(c)}{p(c) + \sum_h p(h)}. \quad (13)$$

Likewise,

$$P(h) = \frac{p(h)}{p(c) + \sum_{h'} p(h')}, \quad (14)$$

denotes the posterior probability of the  $h^{\text{th}}$  hypothesis in the augmented N-best list. (All posteriors and likelihoods are conditioned upon the set of acoustic models  $\{\mathcal{N}_\varphi, \mathcal{G}_\varphi\}$  for  $\varphi = 1 \dots N_\varphi$ .)

The first training criterion can be expressed as

$$\xi = \frac{1}{N_s} \sum_s \log P_s(c) \quad (15)$$

where  $N_s$  is the number of sentences in the training set, and  $P_s(c)$  represents the posterior probability of the correct transcription of sentence  $s$ .

Adapting the neural network weights according to this criterion amounts to adjusting them after the presentation of each training sentence by an amount proportional to (stochastic gradient update)

$$\nabla \log P_s(c) = \sum_h P_s(h) \left[ \nabla \log p_{AM}(c) - \nabla \log p_{AM}(h) \right], \quad (16)$$

where we made use of the property

$$P_s(c) + \sum_{h=1}^{N_h} P_s(h) = 1. \quad (17)$$

Since the acoustic log-likelihoods can be expanded into sums over the observations,  $\mathbf{x}_k$ , in the sentence, the above weight update formula modifies the neural network weights only for those frames where the reference and the hypothesis strings do not coincide. In that case, positive training is given to the correct model (c) and negative training is given to the erroneously hypothesized model (h). The log-likelihood gradients  $\nabla \log p(\cdot)$  are calculated according to Eqs. 6 and 7. This property results from the fact that the N-best list was augmented with the correct transcription (Eq. 10).

Another desirable feature of this training criterion is that more training is given to hypotheses with high posterior probabilities (the multiplicative term,  $P(h)$ ).

A potential disadvantage is that the correct hypothesis is often not in the N-best list for databases with high error rates. Improving the posterior of the correct sentence may thus result in decreasing the probability of the best (although erroneous) hypothesis in the N-best list.

### 6.2. Minimizing the average number of errors in the N-best list

The second training criterion we propose is given by

$$\xi = \frac{1}{N_s} \sum_s \text{ANER}_s, \quad (18)$$

where the average number of errors  $\text{ANER}_s$  in a sentence was defined in Eq. 11.

Note that here the posterior probability of a hypothesis is computed only with respect to the other hypotheses in the N-best list (*i.e.* without taking the reference into account):

$$P_s(h) = \frac{p(h)}{\sum_{h'} p(h')}. \quad (19)$$

Intuitively, minimizing  $\text{ANER}_s$  “redistributes” the posterior probability mass to favor hypotheses with few errors and penalize hypotheses with more errors.

Again, the weight update formula can be derived by taking the instantaneous gradient of  $\xi$  with respect to the weights of the neural networks. The weight update for each sentence is therefore proportional to

$$-\nabla \text{ANER}_s = \sum_h P_s(h) \nabla \log p_{AM}(h) \left[ \text{ANER}_s - \text{NER}_s(h) \right]. \quad (20)$$

The characteristics of this weight update formula are quite different from those of the previous criterion. Negative training is given to hypotheses that have a number of errors above average, and positive training is given to hypotheses with a number of errors below



average. Of course, this average,  $ANER_s$ , evolves with the training of the models. If the learning process progresses correctly,  $ANER_s$  decreases with time, thereby progressively decreasing the number of hypotheses that receive positive training. In the limit, all the posteriors  $P(h)$  converge to zero except the one that corresponds to the hypothesis with the lowest number of errors,  $h^*$ , and  $ANER_s$  converges to  $NER_s(h^*)$ , thereby bringing the training process to an end.

The main disadvantage of this criterion is that positive training is given to all the frames in the best hypothesis, including those associated with incorrectly recognized words. This criterion, however, is closer to the WER metric that we ultimately wish to optimize.

## 7. Recognition Experiments with Discriminatively Trained Dynamo Models

These experiments were limited to the training of small models (NNETs associated to GMMs  $\times 1/16$ ), with linguistic and time features only. Fig.2 shows the results of a self-test experiment (*i.e.* the test data is identical to the training data) with the 627 male sentences of the Eval'96 test set of the Spanish Callhome database. The N-best list depth was limited to 10 hypotheses.

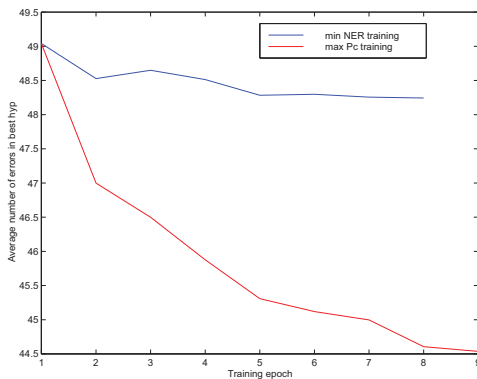


Figure 2: Average number of errors as a function of the training epoch, for both training criteria.

The N-best error rate for this set of sentences was 41.49%. The learning curves show that for the self-test experiment the ANER criterion shows more promise. This, however, is not a fair experiment, and the generalization properties of the max-posterior criterion may be superior. N-best rescoring of 200 randomly selected male sentences of the Eval'96 test set with the neural networks trained to minimize the ANER gave a significant WER improvement (see Table 3).

models	WER
GMMs $\times 1/16$ -- baseline	65.22
min ANER NNETs	63.89

Table 3: N-best rescoring with ANER NNETs, self-test experiment: WER in %.

A fair experiment was conducted with the max-posterior criterion. A set of neural networks was trained from linguistic and time features

to output mixture weights for the same small phone models (GMMs  $\times 1/16$ ). The training data consisted of all 15K male sentences in the training set, of which 10 % was held as a cross-validation set. The models were tested on the same subset of Eval'96 as in the previous experiments. The N-best list depth was limited to 5 hypotheses. The error rate is given in Table 4. The WER improvement is modest but since the phone GMMs in this experiments were small and hence not very detailed, little margin for improvement was left to the NNETs.

models	WER
GMMs $\times 1/16$ -- baseline	65.22
max log-post NNETs	64.79

Table 4: N-best rescoring with log-posterior NNETs, fair experiment: WER in %.

## 8. Conclusions

We described an algorithm to incorporate new knowledge sources in a set of acoustic models, with the objective of dynamically increasing or decreasing the likelihoods of the different modes of the models, thereby narrowing their distributions. The algorithm makes use of feedforward neural networks to dynamically estimate the mixture weights of the speech models, given the knowledge sources for the current data frame.

We argued that the neural networks need to be discriminatively trained, and we proposed two training criteria: maximizing the log-posterior probability of the correct transcription and minimizing the average number of errors in the N-best list. Preliminary experiments showed a modest but encouraging improvement in WER. We are currently experimenting with larger phone models and increased N-best list depths.

## References

- [AKC94] A. Andreou, T. Kamm, and J. Cohen. Experiments in vocal tract normalization. In *Proc. the CAIP Workshop: Frontiers in Speech Recognition II*, 1994.
- [BBdSM86] L. R. Bahl, P. F. Brown, P. V. de Souza, and R. L. Mercer. Maximum mutual information estimation of hidden markov model parameters for speech recognition. In *Proc. IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, 1986.
- [BBdSM88] L. R. Bahl, P. F. Brown, P. V. de Souza, and R. L. Mercer. A new algorithm for the estimation of hidden markov model parameters. In *Proc. IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, New York, NY, April 1988.
- [BdSG<sup>+</sup>91] L. R. Bahl, P. V. de Souza, P. S. Gopalakrishnan, D. Nahamoo, and M.A. Picheny. Context dependent modeling of phones in continuous speech using decision trees. In *DARPA Proc. Speech and Natural Language Workshop*, Pacific Grove, CA, February 1991.
- [BKM95] H. Bourlard, Y. Konig, and N. Morgan. REMAP: Recursive Estimation and Maximization of A Posteriori Probabilities, applications to transition-based connectionist speech recognition. Technical Report TR-94-064, ICSI, Berkeley, CA, March 1995.
- [BM90] H. Bourlard and N. Morgan. A continuous speech recognition system embedding MLP into HMM. In

- D. S. Touretzky, editor, *Advances in Neural Information Processing Systems*, volume 2. Morgan Kaufmann, 1990.
- [Bri90] J. S. Bridle. Training stochastic model recognition algorithms as networks can lead to maximum mutual information estimation of parameters. In D. S. Touretzky, editor, *Advances in Neural Information Processing Systems*, volume 2. Morgan Kaufmann, 1990.
- [Cho90] Y. L. Chow. Maximum mutual information estimation of HMM parameters for continuous speech recognition using the N-best algorithm. In *Proc. IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, Albuquerque, NM, 1990.
- [DASW94] L. Deng, M. Aksmanovic, D. Sun, and J. Wu. Speech recognition using hidden markov models with polynomial regression functions as nonstationary states. *IEEE Trans. Speech, Audio Processing*, 2(4), 1994.
- [DMM96] V. V. Digalakis, P. Monaco, and H. Murveit. Genones: Generalized mixture tying in continuous hidden markov model-based speech recognizers. *IEEE Trans. Speech, Audio Processing*, 4(4), July 1996.
- [FW97] M. Finke and A. Waibel. Speaking mode dependent pronunciation modeling in large vocabulary conversational speech recognition. In *Proc. Eurospeech*, Rhodes, Greece, September 1997.
- [GN93] H. Gish and K. Ng. A segmental speech model with applications to word spotting. In *Proc. IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, volume II, 1993.
- [HW97] L. P. Heck and M. Weintraub. Handset-dependent background models for robust text-independent speaker recognition. In *Proc. IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, Munich, Germany, April 1997.
- [KBM96] Y. Konig, H. Bourlard, and N. Morgan. REMAP: Experiments with speech recognition. In *Proc. IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, Atlanta, GA, May 1996.
- [KLJ91] S. Katagiri, C.-H. Lee, and B.-H. Juang. New discriminative training algorithms based on the generalized probabilistic descent method. In *Proc. Workshop on Neural Networks for Signal Processing*, 1991.
- [KM94] Y. Konig and N. Morgan. Modeling dynamics in connectionist speech recognition - the time index model. In *Proc. Intl. Conf. on Speech and Language Processing*, 1994.
- [MBDW93] H. Murveit, J. Butzberger, V. V. Digalakis, and M. Weintraub. Large-vocabulary dictation using SRI's DECIPHER(TM) speech recognition system: Progressive-search techniques. In *Proc. IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, pages II-319:II-322, April 1993.
- [Mer88] B. Merialdo. Phonetic recognition using hidden markov models and maximum mutual information training. In *Proc. IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, New York, NY, April 1988.
- [OBB<sup>+</sup>97] M. Ostendorf, W. Byrne, M. Bacchiani, M. Finke, A. Gunawardana, K. Ross, S. Roweis, E. Shriberg, D. Talkin, A. Waibel, B. Wheatley, and T. Zeppenfeld. Modeling systematic variations in pronunciation via a language-dependent hidden speaking mode. Technical Report LVCSR Summer Research Workshop, Johns Hopkins U., 1997.
- [ODK96] M. Ostendorf, V. V. Digalakis, and O. A. Kimball. From HMM's to segment models: A unified view of stochastic modeling for speech recognition. *IEEE Trans. Speech, Audio Processing*, 4(5), 1996.
- [Rey96] D.A. Reynolds. Mit lincoln laboratory site presentation. In *NIST Speaker Recognition Workshop*, Linthicum Heights, MD, March 1996.
- [RMT86] D.E. Rumelhart, J.L. McClelland, and The PDP Group, editors. *Parallel Distributed Processing*, volume 1. The MIT Press, Cambridge, MA, 1986.
- [Slo95] T. Sloboba. Dictionary learning: Performance through consistency. In *Proc. IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, 1995.
- [WBR<sup>+</sup>97] M. Weintraub, F. Beaufays, Z. Rivlin, Y. Konig, and A. Stolcke. Neural - network based measures of confidence for word recognition. In *Proc. IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing*, Munich, Germany, April 1997.
- [YOW94] S. J. Young, J. J. Odell, and P. C. Woodland. Tree-based state tying for high accuracy acoustic modelling. In *Proc. Human Language Technology Workshop*, pages 307-312, Plainsboro, NJ, March 1994.

EXHIBIT 4  
FULLY REDACTED



EXHIBIT 5  
FULLY REDACTED

EXHIBIT 6  
FULLY REDACTED