

Doctoral Course in Speech Recognition

September - November 2003

Part 3

Kjell Elenius

Nov 28 2003

Speech recognition course

1

CHAPTER 12

BASIC SEARCH ALGORITHMS

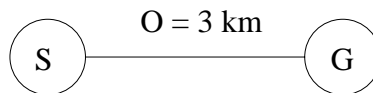
Nov 28 2003

Speech recognition course

2

State-based search paradigm

- Triplet S, O, G
 - S , set of initial states
 - O , set of operators applied on a state to generate a transition with corresponding cost to another state
 - G , set of goal states



Nov 28 2003

Speech recognition course

3

12.1.1 General Graph Searching Procedures

- Dynamic Programming is powerful but cannot handle all search problems, e.g. NP-hard problems
- NP hard problems
 - Definition: The complexity class of decision problems that are intrinsically harder than those that can be solved by a nondeterministic Turing machine in polynomial time.
- Examples
 - The traveling salesman problem
 - Visit all cities once, find shortest distance
 - The 8 Queen problem
 - Place 8 Queens on a chessboard so no-one can capture any of the other

Nov 28 2003

Speech recognition course

4

Simplified Salesman Problem

- Find shortest path from S to G

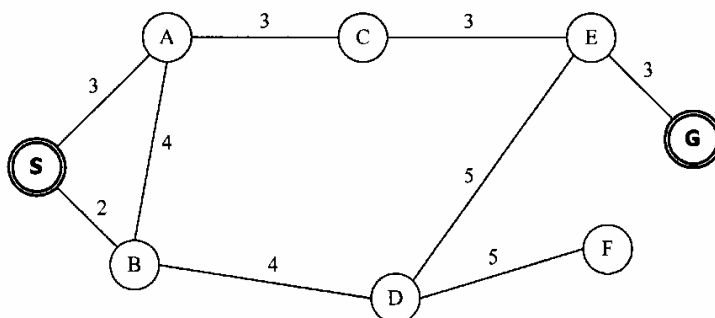


Figure 12.1 A highway distance map for cities S, A, B, C, D, E, F, and G. The salesman needs to find a path to travel from city S to city G [42].

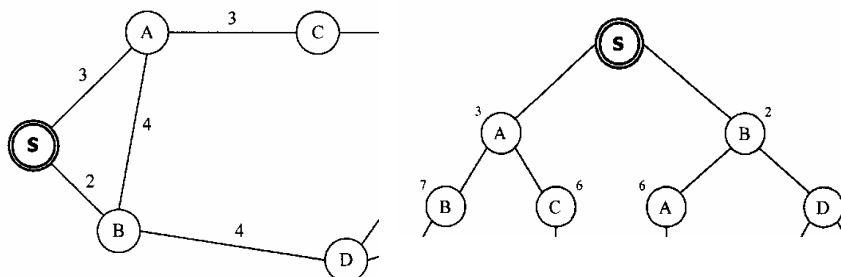
Nov 28 2003

Speech recognition course

5

Expand paths

- The successor (move) operator
 - generates all successors of a node and computes all costs associated with an arc
- Branching factor
 - average number of successors for each node



Nov 28 2003

Speech recognition course

6

Fully expanded Search Tree (Graph)

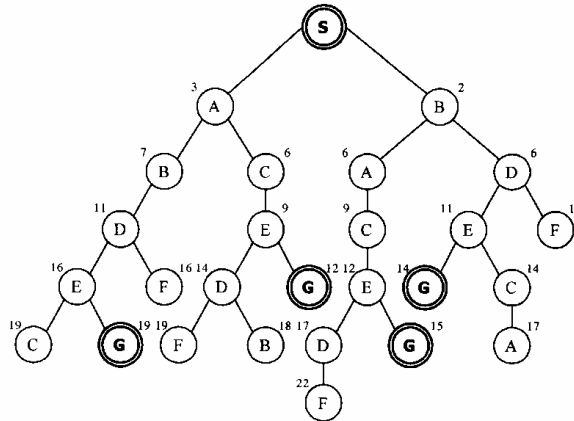


Figure 12.2 The search tree (graph) for the salesman problem illustrated in Figure 12.1. The number next to each node is the accumulated distance from start city to end city [42].

Nov 28 2003

Speech recognition course

7

Explicit search impractical for large problems

- Use Graph Search Algorithm
 - Dynamic Programming principle
 - Only keep the shortest path to a node
- Forward direction (reasoning) normal
- Backward reasoning
 - more initial states than goal states
 - backward branching factor smaller than the forward one
- Bi-directional search
 - start from both ends simultaneously

Nov 28 2003

Speech recognition course

8

A bad case for bi-directional search

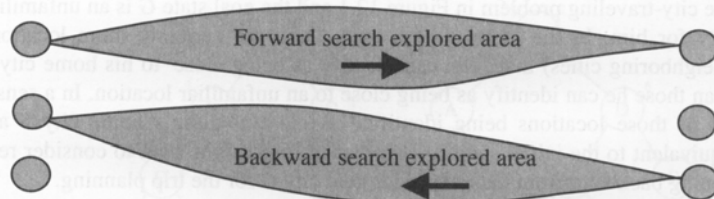


Figure 12.3 A bad case for bi-directional search, where the forward search and the backward search crossed each other [42].

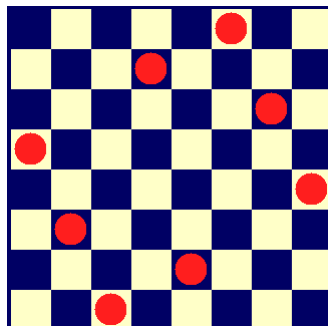
Nov 28 2003

Speech recognition course

9

The 8 Queen problem

- 1 of 12 solutions



Nov 28 2003

Speech recognition course

10

12.1.2 Blind Graph Search Algorithms

- Find an acceptable path - need not be the best one
- Blindly expand nodes without using domain knowledge
- Also called Uniform search or Exhaustive search

Nov 28 2003

Speech recognition course

11

Depth-First Search

- Deepest nodes are expanded first
- Nodes of equal depth are expanded arbitrarily
- Backtracking
 - If a dead-end is reached go back to last node and proceed with another one
- If Goal reached, exit

Nov 28 2003

Speech recognition course

12

Depth-First Search

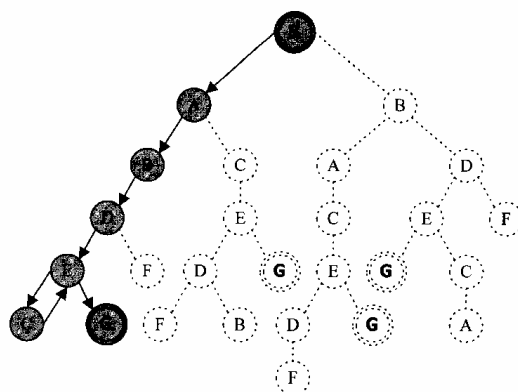


Figure 12.4 The node-expanding procedure of the depth-first search for the path search problem in Figure 12.1. When it fails to find the goal city in node *C*, it backtracks to the parent and continues the search until it finds the goal city. The gray nodes are those that are explored. The dotted nodes are not visited during the search [42].

Nov 28 2003

Speech recognition course

13

Breadth-First Search

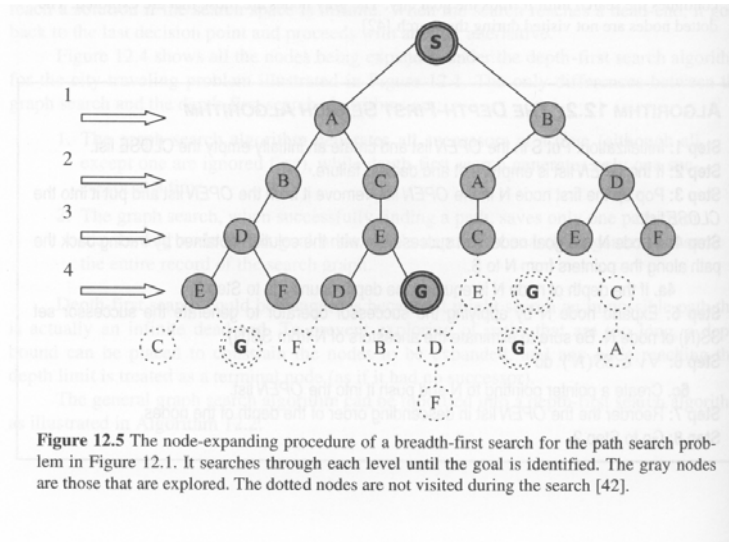
- Same level nodes are expanded before going to the next level
- Stop when goal is reached
- Guaranteed to find a solution if one exists
- Can find optimal solution after all solutions have been found -- brute-force search

Nov 28 2003

Speech recognition course

14

Breadth-First Search



Nov 28 2003

Speech recognition course

15

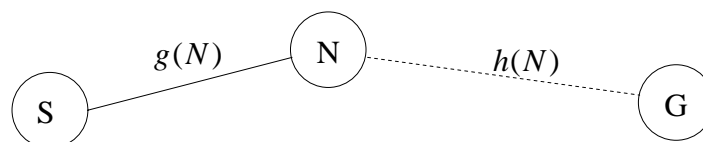
12.1.3 Heuristic Graph Search

- Use domain-specific (heuristic) knowledge to the guide search

$h(N)$ Heuristic estimate of remaining distance from node N to G

$g(N)$ The distance of the partial path from root S to node N

$f(N) = g(N) + h(N)$ Estimate of the total distance from S to N



Nov 28 2003

Speech recognition course

16

Best-First (A^* Search)

- A search is said to be admissible if it can guarantee to find an optimal solution if one exists
- If $h(n)$ is an underestimate of the remaining distance to G the best first search is admissible. This is called A^* search.

Nov 28 2003

Speech recognition course

17

City travel problem

- Use straight-line distance to goal as heuristic information (bold digits)

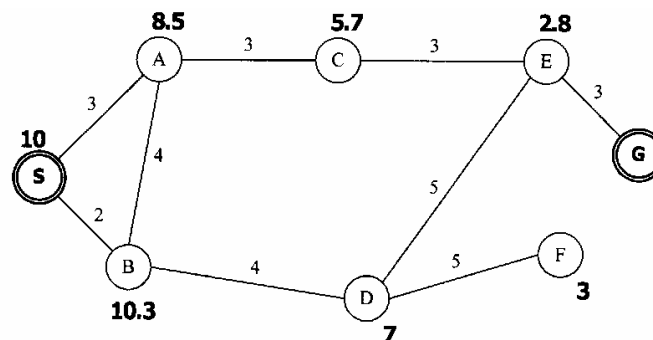


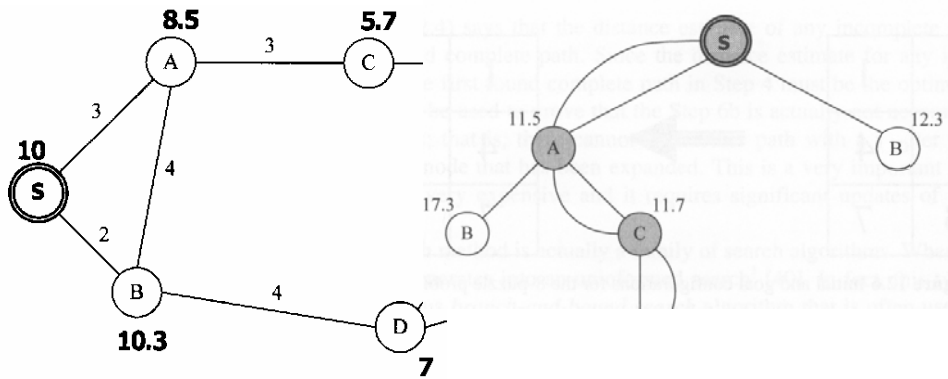
Figure 12.7 The city-travel problem augmented with heuristic information. The numbers beside each node indicate the straight-line distance to the goal node G [42].

Nov 28 2003

Speech recognition course

18

City travel problem with heuristics



Nov 28 2003

Speech recognition course

19

City travel problem with heuristics

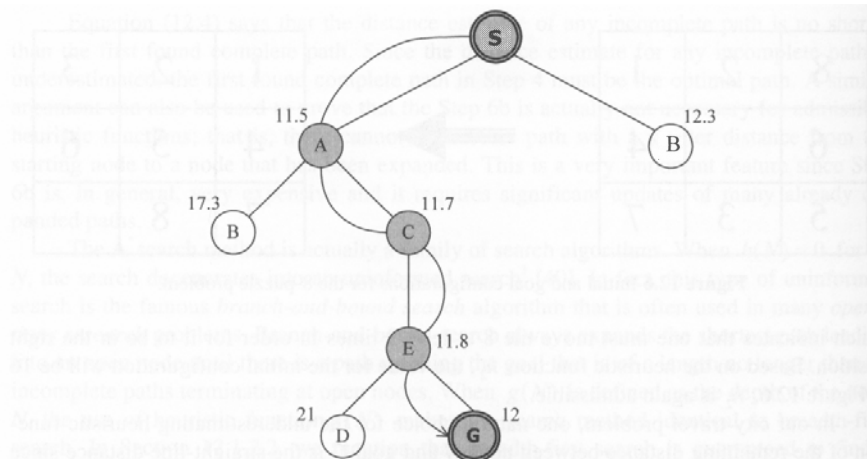


Figure 12.8 The search progress of applying A^* search for the city-travel problem. The search determines that path S-A-C-E-G is the optimal one. The number beside the node is f values on which the sorting of the *OPEN* list is based [42].

Nov 28 2003

Speech recognition course

20

Beam Search

- Breadth-first type of search but only expand paths likely to succeed at each level
- Only these nodes are kept in the beam and the rest are ignored, pruned
- In general a fixed number of paths, w , are kept at each level

Nov 28 2003

Speech recognition course

21

Beam Search

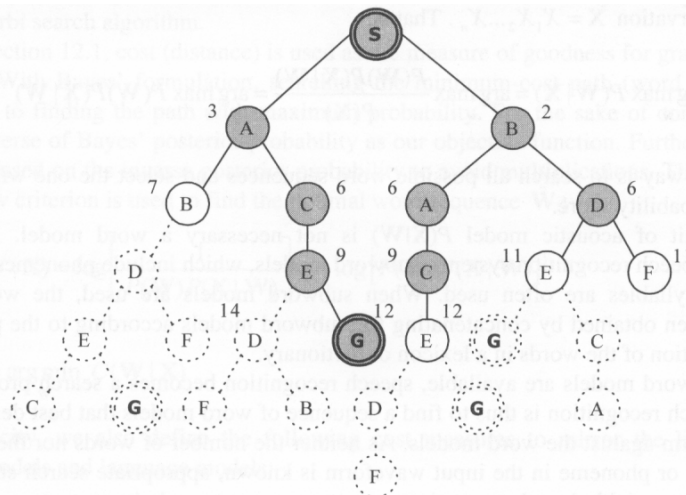


Figure 12.9 Beam search for the city-travel problem. The nodes with gray color are the ones kept in the beam. The transparent nodes were explored but pruned because of higher cost. The dotted nodes indicate all the savings because of pruning [42].

Nov 28 2003

Speech recognition course

22

12.2 Search Algorithms for Speech Recognition

- Goal of speech recognition
 - Find word sequence with maximum posterior probability
- Change to minimum criterion function C for consistency with search
 - use inverse of Bayes posterior probability and use log to avoid multiplications

$$\hat{W} = \arg \max_w P(W|X) = \arg \max_w \frac{P(W)P(X|W)}{P(X)} \propto \arg \max_w P(W)P(X|W)$$

$$C(W|X) = \log \left[\frac{1}{P(W)P(X|W)} \right] = -\log[P(W)P(X|W)]$$

$$\hat{W} = \arg \min_w C(W|X)$$

Nov 28 2003

Speech recognition course

23

12.2.2 Combining Acoustic and Language Models

- To balance the language model probability with the acoustic model probability a language model weight LW is introduced
 - thus we get the language model $P(W)^{LW}$
- Since generally $LW > 1$ every new word gets a penalty
 - if penalty is large the decoder prefers few long words else many short words
- If LW is used primarily to balance the acoustic model a special Insertion Penalty IP may be used
 - thus we get the language model $P(W)^{LW} IP^{N(W)}$

Nov 28 2003

Speech recognition course

24

12.2.3 Isolated Word Recognition

- Boundaries known
- Calculate $P(X|W)$ using forward algorithm or Viterbi
- Chose W with highest probability
- When subword models (monophones, triphones, ...) are used HMMs may be easily concatenated

Nov 28 2003

Speech recognition course

25

12.2.4 Continuous Speech Recognition

- Complicated
 - no reliable segmentation of words
 - each word can theoretically start at any time frame
 - the search space becomes huge for large vocabularies

Nov 28 2003

Speech recognition course

26

Simple cont. task with 2 words

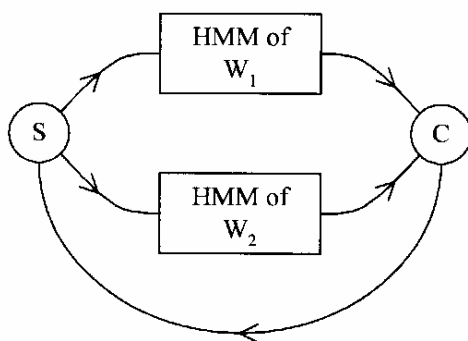


Figure 12.10 A simple example of continuous speech recognition task with two words w_1 and w_2 . A uniform unigram language model is assumed for these words. State S is the starting state while state C is a collector state to save fully expanded links between every word pair.

Nov 28 2003

Speech recognition course

27

HMM trellis for 2 word cont. rec.

- Viterbi search
 - stochastic finite state network with transition probabilities and output distributions

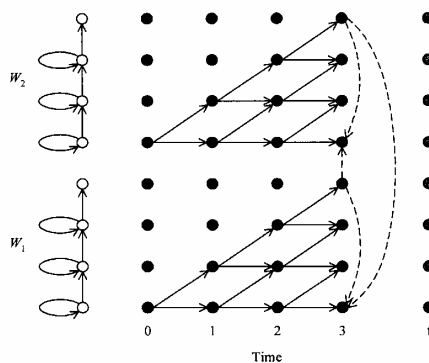


Figure 12.11 HMM trellis for continuous speech recognition example in Figure 12.10. When the final state of the word HMM is reached, a null arc (indicated by a dashed line) is linked from it to the initial state of the following word.

Nov 28 2003

Speech recognition course

28

HMM trellis for 2 word cont. rec.

- Viterbi search
 - the computations are done *time-synchronously* from left to right, i.e. each cell for time t is computed before proceeding to time $t+1$

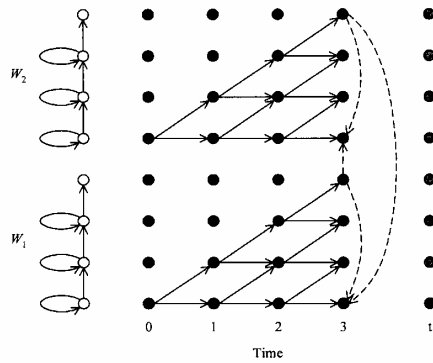


Figure 12.11 HMM trellis for continuous speech recognition example in Figure 12.10. When the final state of the word HMM is reached, a null arc (indicated by a dashed line) is linked from it to the initial state of the following word.

Nov 28 2003

Speech recognition course

29

12.3 Language Model States

- Search Space with FSM and CFG
- Search space with Unigram, Bigrams and Trigrams
- How to Handle Silence Between Words

Nov 28 2003

Speech recognition course

30

12.3.1 Search Space with FSM and CFG

- FSM
 - word network expanded into phoneme network (HMMs)
 - sufficient for simple tasks
 - very similar to CFG when using sub-grammars and word classes
- CFG
 - set of production rules expanding non-terminals into sequence of terminals (words) and non-terminals (e.g. dates, names)
 - Chart parsing not suitable for speech recognition which requires left-to-right processing
 - Formulated with Recursive Transition Network (RTN)

Nov 28 2003

Speech recognition course

31

Finite-State Machines (FSM)

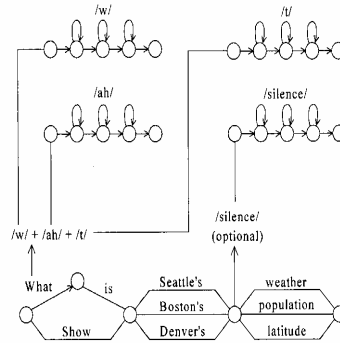
- word network expanded into phoneme network (HMMs)
- search using time-synchronous Viterbi
- sufficient for simple tasks
- similar to CFG when using sub-grammars and word classes

Nov 28 2003

Speech recognition course

32

FSM



Nov 28 2003

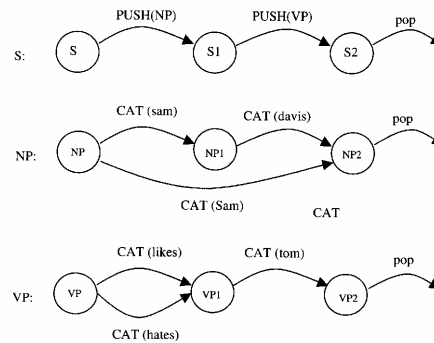
Speech recognition course

33

Search with CFG

- Example formulation with RTN

S → NP VP
NP → sam | sam davis
VP → VERB tom
VERB → likes | hates



CAT arcs can be expanded to HMMs and searched

Nov 28 2003

Speech recognition course

34

12.3.2 Search Space with Unigrams

- The simplest n-gram is the memoryless unigram

$$P(\mathbf{W}) = \prod_{i=1}^n P(w_i)$$

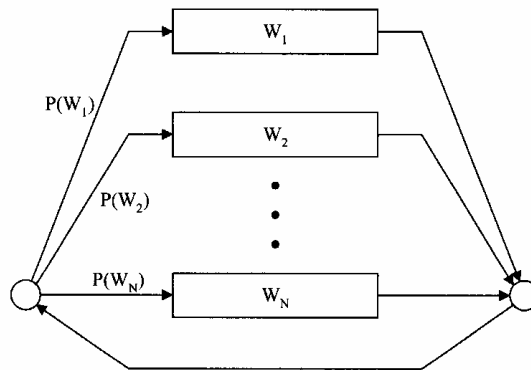


Figure 12.14 A unigram grammar network where the unigram probability is attached as the transition probability from starting state S to the first state of each word HMM.

Nov 28 2003

Speech recognition course

35

12.3.3 Search Space with Bigrams

N states, N^2 word transitions $P(\mathbf{W}) = P(w_i | < s >) \prod_{i=1}^n P(w_i | w_{i-1})$

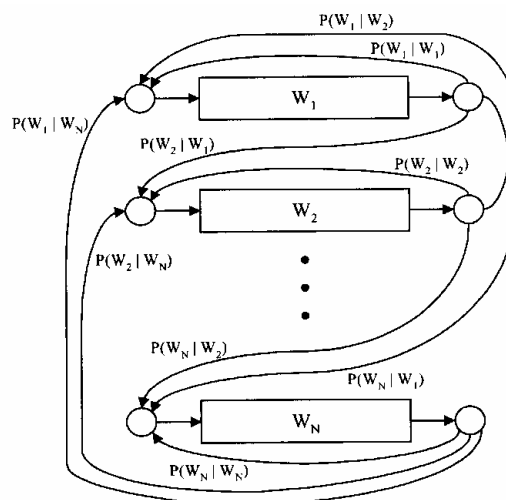


Figure 12.15 A bigram grammar network where the bigram probability $P(w_i | w_j)$ is attached as the transition probability from word w_j to w_i [19].

12.3.3.1 Backoff Paths

For an unseen bigram $P(w_j | w_i) = \alpha(w_i)P(w_j)$ where $\alpha(w_i)$ is the backoff weight for word w_i

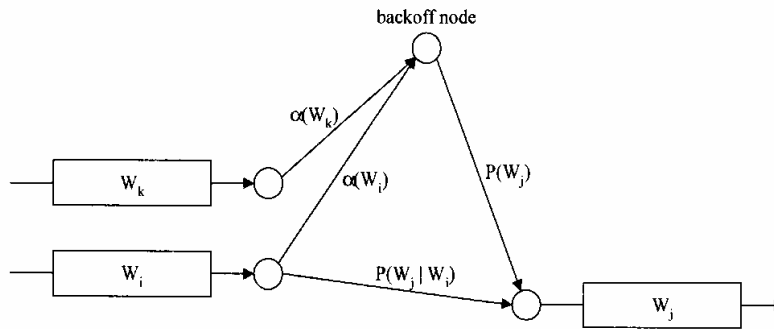
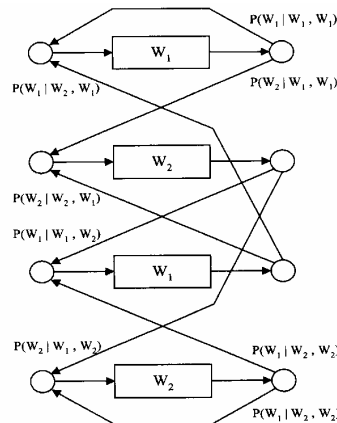


Figure 12.16 Reducing bigram expansion in a search by using the backoff node. In addition to normal bigram expansion arcs for all observed bigrams, the last state of word w_i is first connected to a central backoff node with transition probability equal to backoff weight $\alpha(w_i)$. The backoff node is then connected to the beginning of each word w_j with its corresponding unigram probability $P(w_j)$ [12].

12.3.4 Search Space with Trigrams

- The search space is considerably more complex
 - N^2 grammar states and from each of these there is a transition to the next word



12.3.5 How to Handle Silences Between Words

- Insert optional silence between words

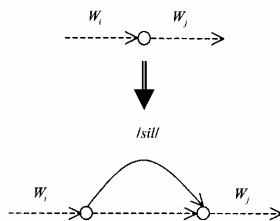


Figure 12.18 Incorporating optional silence (a non-speech event) in the grammar search network where the grammar state connecting different words is laced by two parallel paths. One is the original null transition directly from one word to the other, while the other first goes through the silence word to accommodate the optional silence.

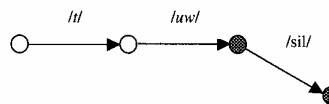


Figure 12.19 An example of treating silence as of the pronunciation network of word TWO. The shaded nodes represent possible word-ending nodes: one without silence and the other one with silence.

Nov 28 2003

12.4 Time-Synchronous Viterbi Beam search

- The Viterbi approximation
 - The *most likely word sequence* is approximated by the *most likely state sequence*
 - For time t each state is updated by the best score of time $t-1$
 - e.g. *time synchronous Viterbi search*
 - record backtracking pointer at each update
 - better to only record word history at end of each word
 - then we need only 2 successive time slices for the Viterbi computations

Nov 28 2003

Speech recognition course

40

12.4.1 The Use of Beam

- The search space for Viterbi search is $O(NT)$ and the complexity $O(N^2T)$ where
 - N is the total number of HMM states
 - T is the length of the utterance
- For large vocabulary tasks these numbers are astronomically large even with the help of dynamic programming
- Prune search space by beam search
- Calculate lowest cost D_{\min} at time t
- Discard all states with cost larger than $D_{\min} + T$ before moving on to the next time sample $t+1$

Nov 28 2003

Speech recognition course

41

12.4.2 Viterbi Beam Search

- Empirically a beam size of between 5% and 10% of the total search space is enough for large-vocabulary speech recognition.
- This means that 90% to 95% can be pruned off at each time t .
- The most powerful search strategy for large vocabulary speech recognition

Nov 28 2003

Speech recognition course

42

Forward trellis space for stack decoding

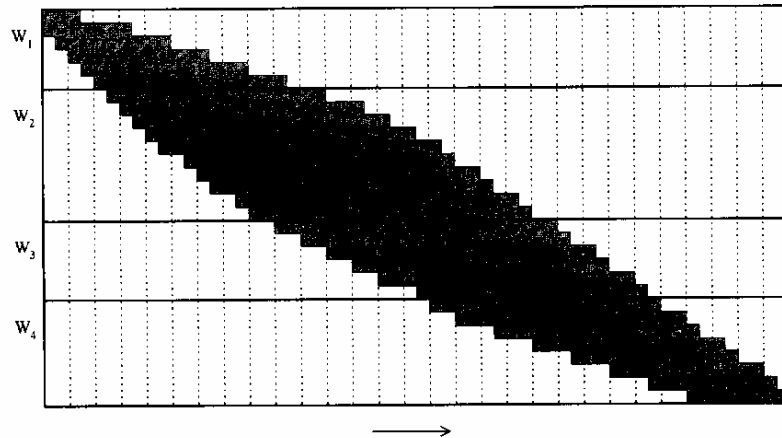


Figure 12.21 The forward trellis space for stack decoding. Each grid point corresponds to a trellis cell in the forward computation. The shaded area represents the values contributing to the computation of the forward score for the optimal word sequence w_1, w_2, w_3, \dots [24].

12.5.1 Admissible Heuristics for Remaining Path

- $f(t) = g(t) + h(T-t)$
- Calculate the expected cost per frame Ψ from the training set by using forced alignment
- $f(t) = g(t) + (T-t)\Psi$

Unnormalized cost

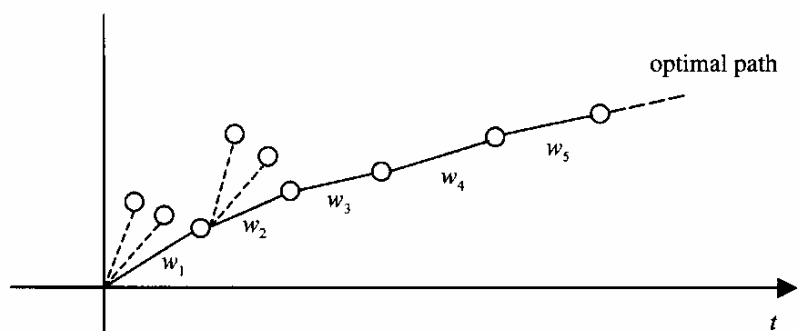


Figure 12.23 Unnormalized cost $C(\mathbf{x}'_1, s_t | w_t^k)$ for optimal path and other competing paths as a function of time.

Nov 28 2003

Speech recognition course

47

Normalized cost

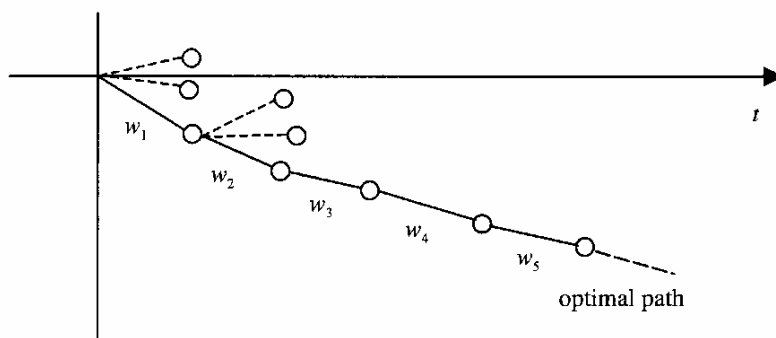


Figure 12.24 Normalized cost $\hat{C}(\mathbf{x}'_1, s_t | w_t^k)$ for the optimal path and other competing paths as a function of time.

Nov 28 2003

Speech recognition course

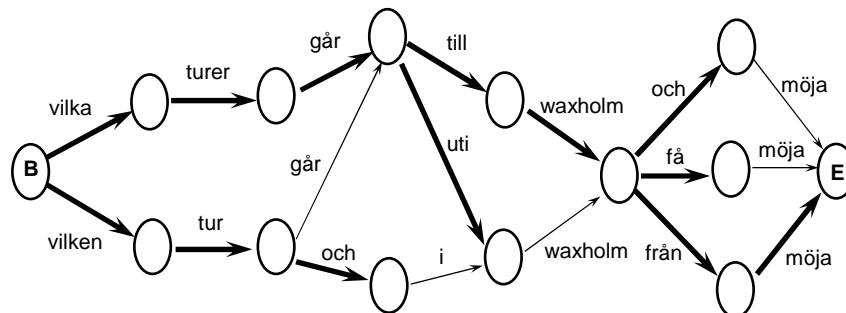
48

Dynamisk programmering - Viterbi

Exempel: finn bästa vägen genom ordgrafan

Processa noderna från vänster till höger.
För varje nod, hitta den "inkommande" båge som ger minst kostnad från B.
Markera noden för bakåtspårning och spara den minsta kostnaden.

Bästa vägen hittas genom bakåtspårning från E längs de markerade bågarna.



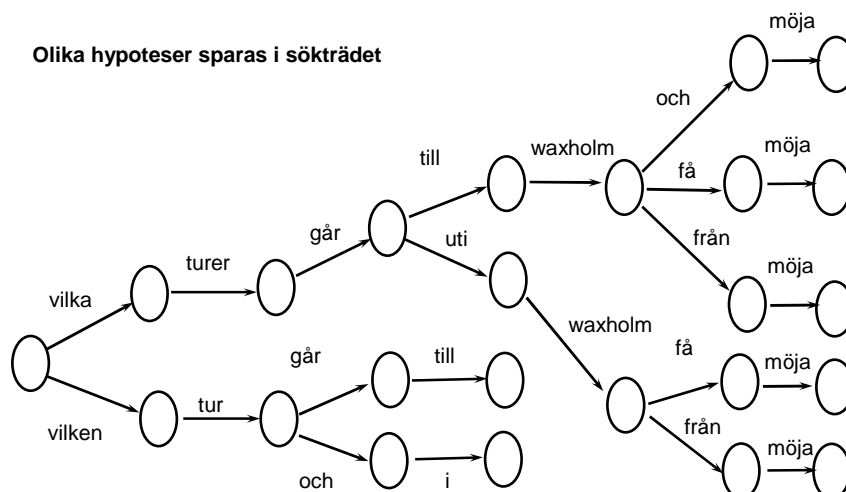
Nov 28 2003

Speech recognition course

49

N-bästa-sökning

Olika hypoteser sparas i sökträdet



Nov 28 2003

Speech recognition course

50

N-bästa sökning med A*

- Utför en "Best First" sökning med flera hypoteser
- Den första som täcker yttrandet är bäst, andra näst bäst osv.
- Problem: Monotont sjunkande sannolikhet vid ökande hypoteslängd => korta hypoteser gynnas -> långsamt
- Lösning: Lägg till varje hypotes' maximala sannolikhet (h^*) för den ej sökta delen av yttrandet -> alla hypoteser får samma längd.
- Metrik: $P(\text{hyp}, t) = g[0, t] + h^*[t+1, T]$ (mätt + uppskattad)
- h^* för alla tider och tillstånd erhålls med normal Viterbi-sökning i ett inledande pass.

Nov 28 2003

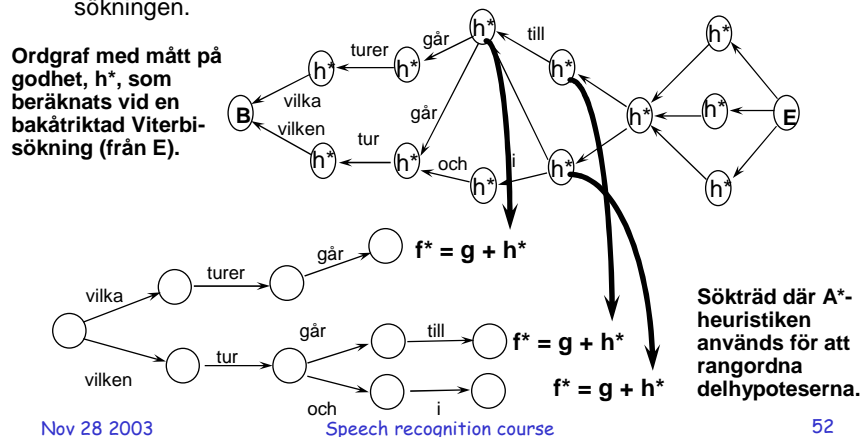
Speech recognition course

51

A*-sökning

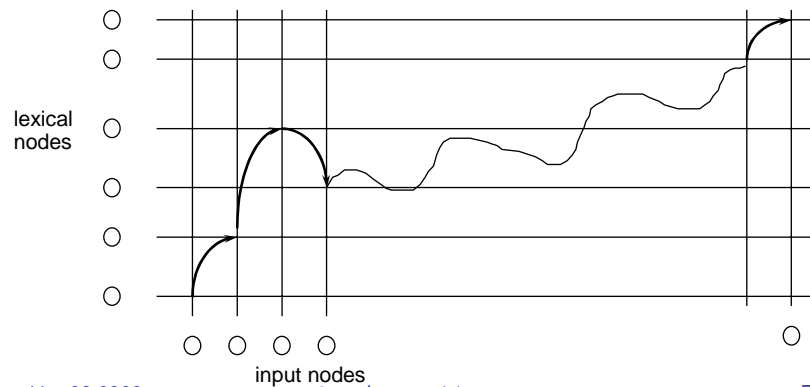
Delhypoteserna ordnas efter sin ackumulerade godhet, g , och h^* -heuristiken, dvs godheten för bästa vägen från den aktuella slutnoden till slutet av yttrandet, ($f^* = g + h^*$).

h^* -mättet beräknas i ett Viterbi pass från E d v s motsatt håll mot A*-sökningen.



Lexikal sökning

- Sökmatriisen består av alla par av inputnoder och lexikala noder
- En tillåten väg i sökmatriisen är en serie steg mellan olika input- och lexikala noder sådana att det finns bågar som binder ihop dem i respektive nät.

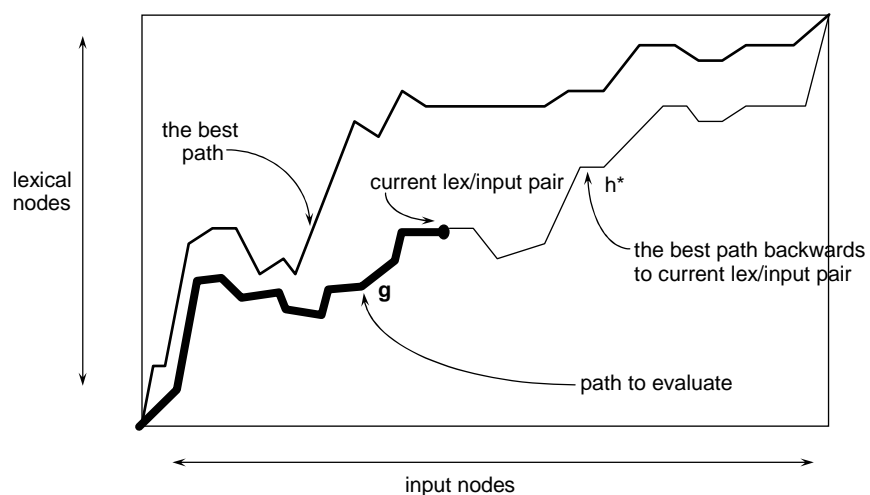


Nov 28 2003

Speech recognition course

53

The A* Heuristic in the Search Matrix

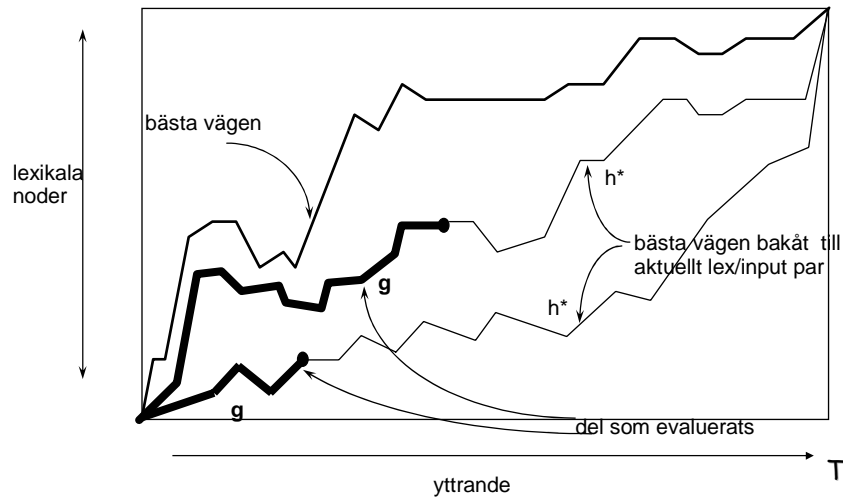


Nov 28 2003

Speech recognition course

54

A*-heuristik i sökmatriisen

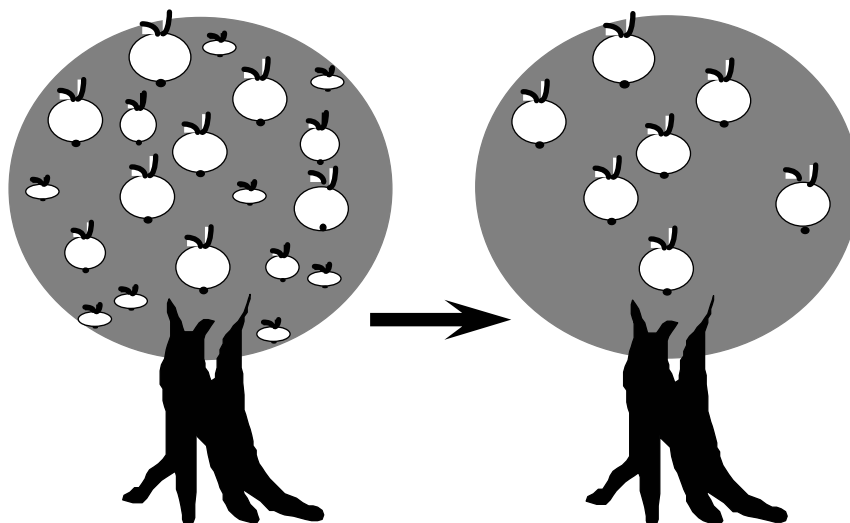


Nov 28 2003

Speech recognition course

55

Beskärning (pruning)

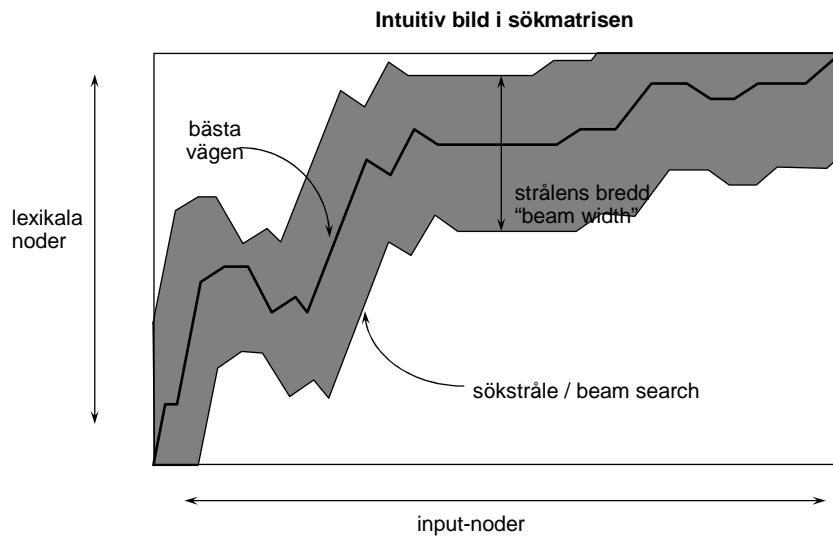


Nov 28 2003

Speech recognition course

56

Begränsad sökning - "Beam Search"



Nov 28 2003

Speech recognition course

57

N-bästa exempel de 10 akustiskt bästa meningarna

Nr	Akustisk poäng	Igenkänd mening
1	-716,914	hur åka till svartlöga
2	-717,444	du åka till svartlöga
3	-718,289	ju åka till svartlöga
4	-719,614	hur åker man till svartlöga
5	-719,730	hur åka av hamn till svartlöga
6	-720,260	du åka av hamn till svartlöga
7	-720,365	ut åka till svartlöga
8	-720,554	hur båtarna till svartlöga
9	-720,630	hur åker hamn till svartlöga
10	-720,699	nu åka till svartlöga

Nov 28 2003

Speech recognition course

58

SpeechDat

Philips Business Systems NL	TeleNor Research and Development, N
British Telecom UK	Digital Media Institute Tampere Technical University, SF
GEC-Marconi, Marconi Speech and Information Systems, UK	Dept. Speech, Music and Hearing Kung. Tekniska Hogskolan, S
GEC-Marconi Hirst Division, UK	Center for Personkommunikation University of Aalborg, DK
GPT Ltd. UK	TeleDanmark DK
Yocalis Ltd. UK	Speech Processing Expertise Centre, NL
Lemout & Hauspie B	Philips Forschungslaboratorien D
Matra Communication F	Department of Phonetics University of Munich, D
ELRA F	SIEMENS AG D
PTT Suissex CH	Institute of Electronics University of Maribor, SI
IDIA P CH	CSELT I
INESC PT	Universitat Politècnica de Catalunya, E
Portugal Telecom PT	Knowledge SA GR
	Wire Communications Laboratory University of Patras, GR

Partners

Nov 28 2003 Speech recognition course 59

EU-projektet

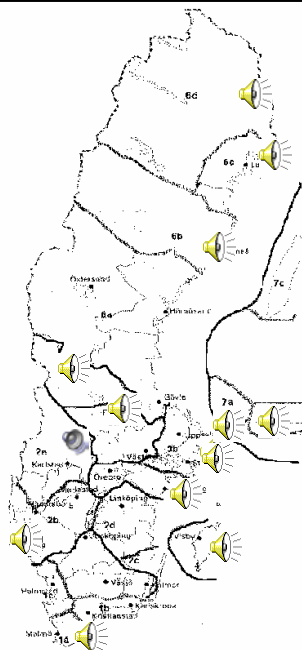
- Inspelat tal över telefonnätet för att träna och testa taligenkänningsystem
 - alla 11 officiella EU-språk samt varianter som finlandssvenska, schweizertyska, walesiska
 - totalt över 60 000 talare inspelade
- balansera talare enligt
 - dialekt, ålder och kön
- ca 50 yttranden per talare
 - siffror, datum, tider, penningbelopp, enkla kommandon, fonetiskt rika meningar och ord
- SpeechDat i Sverige
 - 5000 talare inspelade över vanlig telefon
 - 1000 talare inspelade över mobiltelefon

Nov 28 2003 Speech recognition course 60

Svenska dialekter

"Flyget, tåget och bilbranschen tävlar om lönsamhet och folkets gunst".

 Född i
USA
 ex-Jugoslavien



Störningar och annat

- Mobiltelefoni
 - bil, trottoar, restaurang
 -  • *Bengt Dennis ger inga avskedsintervjuer inför sin avgång vid årsskiftet*
 -  • *Det handlar bara om ett glapp på 18 månader*
- Dialektalt uttryckssätt
 -  • *Han försökte förgäves rädda sin hustru på övervåningen*
- Den mänskliga faktorn
 -  • *Kvinnan är mycket nära en total kollaps och gråter oupphörligt*

ARPA-projektet i USA

- Advanced Projects Research Agency
 - började 1984
 - deltagare
 - CMU, SRI, BBN, MIT, Lincoln Labs, Dragon Systems
 - "competitive evaluations" varje år
 - domäner
 - RM, Resource Management, ~1000 ord
 - ATIS, Air Traffic Information System, ~1000 ord
 - Flygbokning
 - WSJ, Wall Street Journal, 5000 - 60000 ord
 - Tidningstext, uppläst
 - SWITCHBOARD
 - samtal över telefon med okänd person om givet ämne
 - CALL HOME
 - samtal över telefon med närmaste familjekretsen
 - NAB, National Broadcast News
 - radiotal, olika talare, telefon ibland, musik

svårare

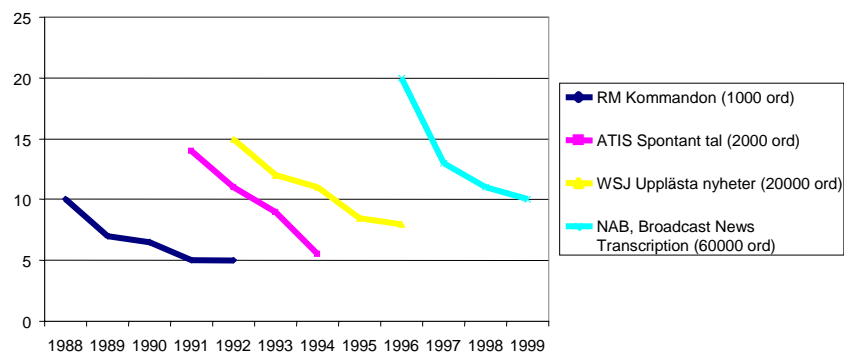
Nov 28 2003

Speech recognition course

63

Prestandautveckling DARPA-utvärdering 1988-1999

Ordfel (%)



Nov 28 2003

Speech recognition course

64

Igenkänningsresultat för olika ARPA-databaser

Databas	År	Inspelade timmar för inläring	Miljoner ord för inläring av språkmodellen	Ordfel
Wall Street Jour. kontor telefon	1996	60	200	7%
		60	200	37%
Switchboard telefon	1996	70	2	38%
Call Home telefon	1996	13	0,2	50%
Broadcast news kontor & telefon	2000	200	~ 1000	20%

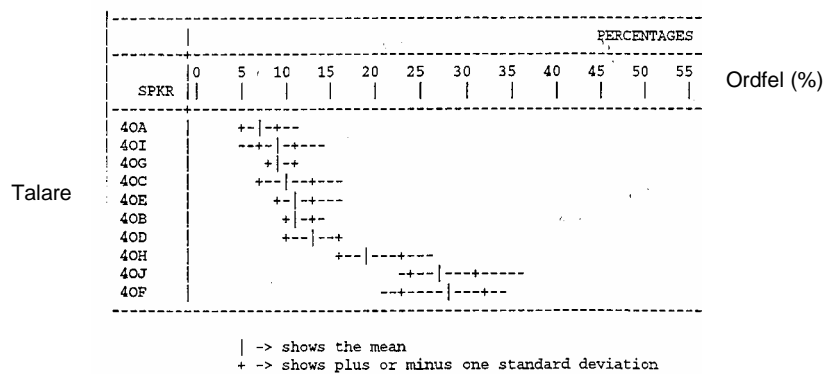
Nov 28 2003

Speech recognition course

65

Får och getter ("bra" och "dåliga" talare)

ARPA-utvärdering, Wall Street Journal
nov 1993 Hub1



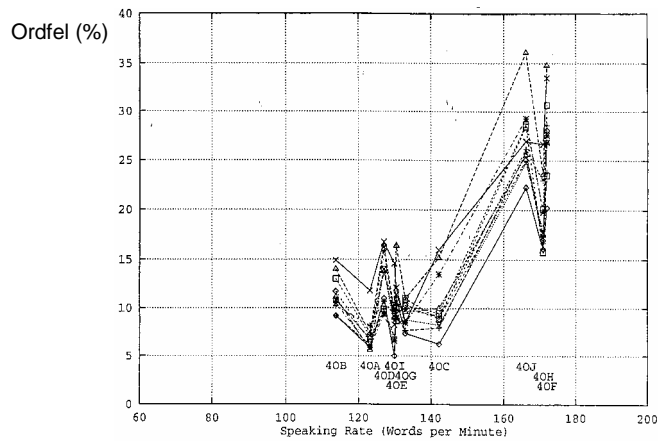
Nov 28 2003

Speech recognition course

66

Talshastighetens inverkan

ARPA-utvärdering, Wall Street Journal
nov 1993 Hub1

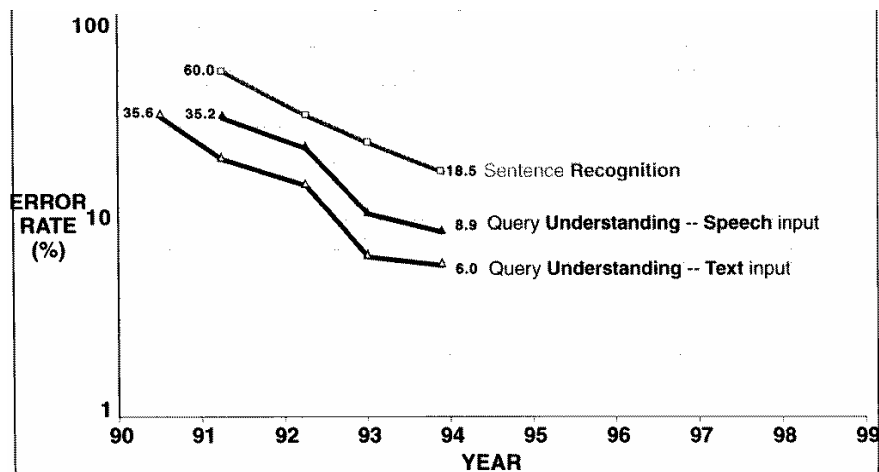


Nov 28 2003

Speech recognition course

67

ARPA-resultat från ATIS-projektet



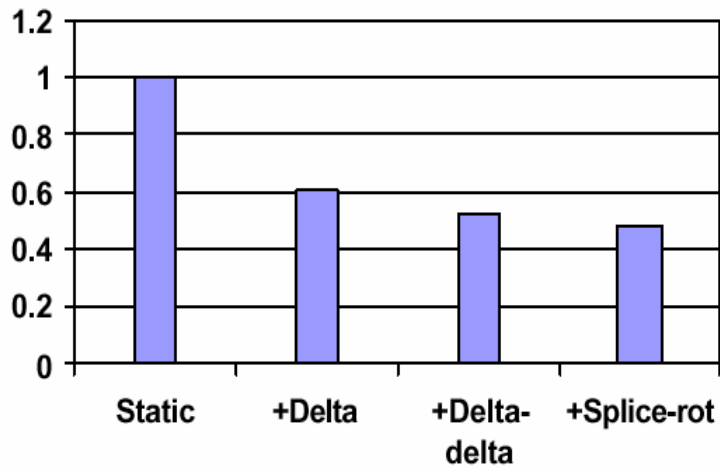
Från presentationen, *Human Languages Technology*, av George Doddington vid ARPA Workshop, New Jersey, 8 - 11 mars, 1994.

Nov 28 2003

Speech recognition course

68

Effekten av att addera parametrar (ordfelsprocent)

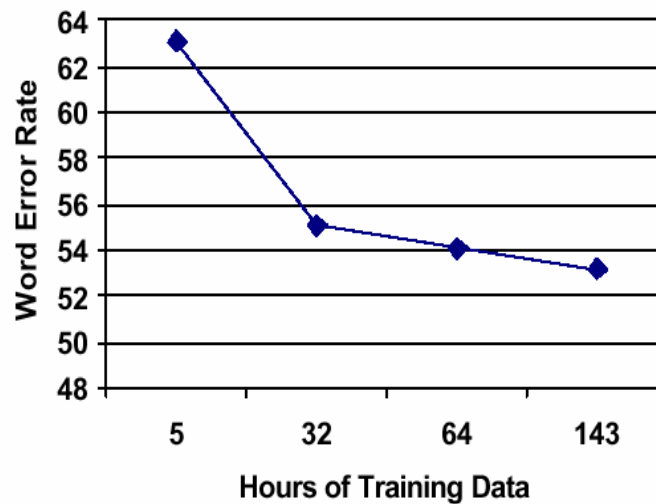


Nov 28 2003

Speech recognition course

69

Ordfel mot antal timmar tal för träning BBN, Switchboard data

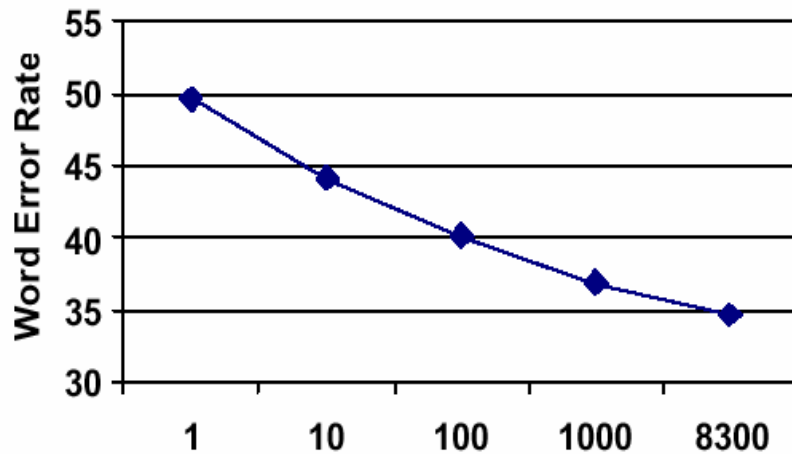


Nov 28 2003

Speech recognition course

70

Ordfel mot antal tusen meningar som använts för att träna en trigramspråkmodell



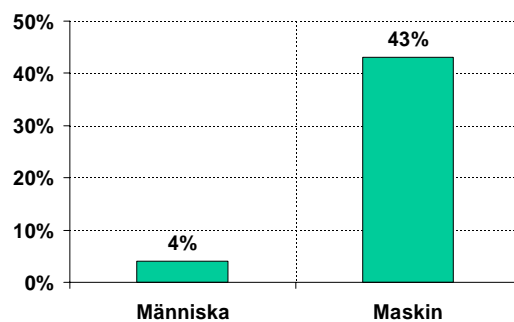
Nov 28 2003

Speech recognition course

71

Kreditkortsämnet i Switchboard

Resultat för mänsklig transkription av hela samtal jämfört med automatisk igenkänning av extraherade meningar från samma material.



R. P. Lippman: "Speech perception by humans and machines", Proc. of the workshop on the auditory basis of speech perception, Keele University, UK, 15-19 July, 1996, pp. 309-316.

Nov 28 2003

Speech recognition course

72

Emotional/ubiquitous computing - do we want it?
Early BBC vision - the conversational toaster



Nov 28 2003

Speech recognition course

73

SLUT

Nov 28 2003

Speech recognition course

74