



# Wortdekodierung

Vorlesungsunterlagen  
Speech Communication 2, SS 2004

Franz Pernkopf/Erhard Rank

Institute of Signal Processing and Speech Communication

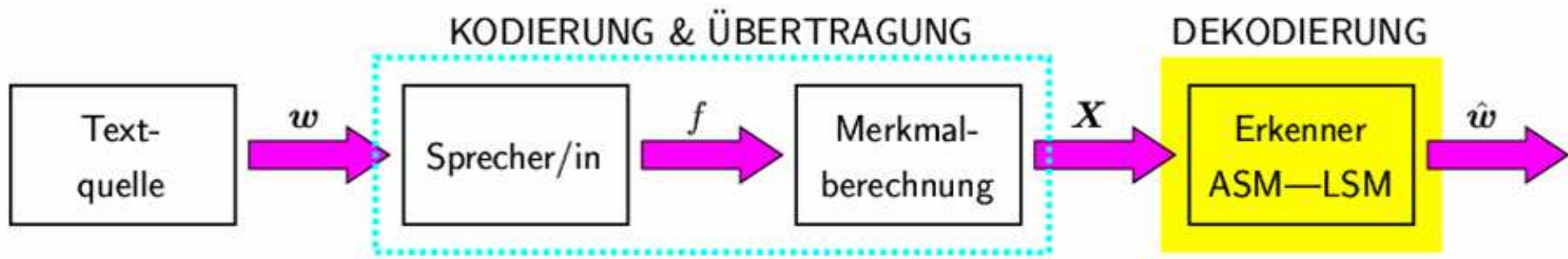
University of Technology Graz

Inffeldgasse 16c , 8010 Graz, Austria

Tel: +43 316 873 4436

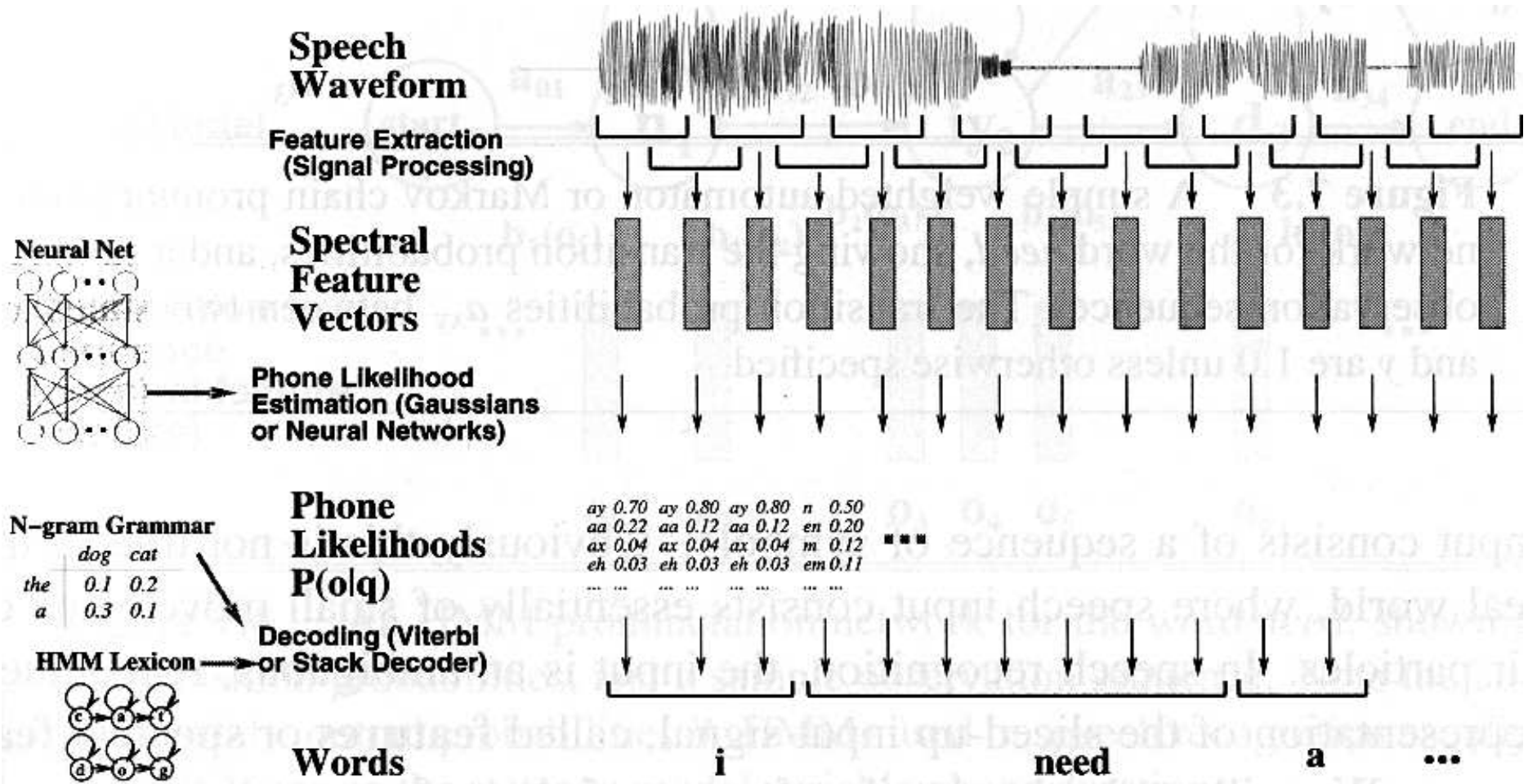
E-Mail: [pernkopf/rank@inw.tugraz.at](mailto:pernkopf/rank@inw.tugraz.at)

## Kommunikationstheoretisches Modell der Spracherzeugung und -erkennung:



DEKODIERUNG = Maximierungsaufgabe

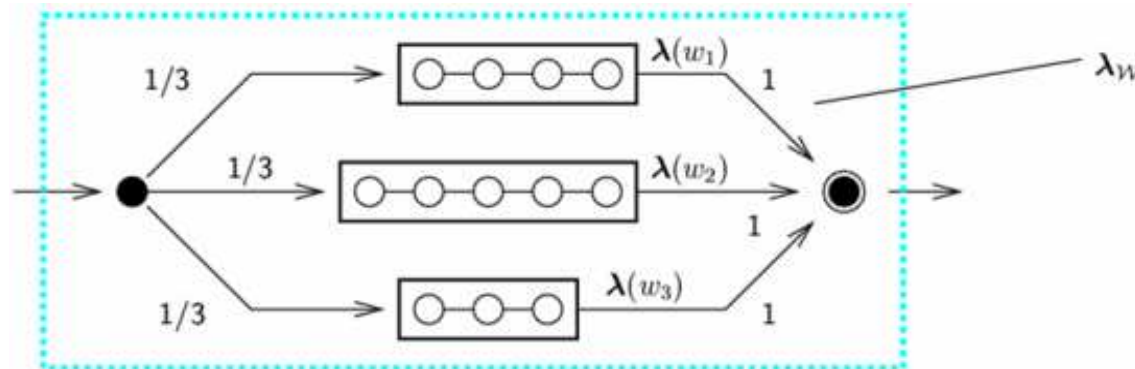
$$\hat{w} = \operatorname{argmax}_v P(\mathbf{X} | \mathbf{v}) \cdot P(\mathbf{v})$$



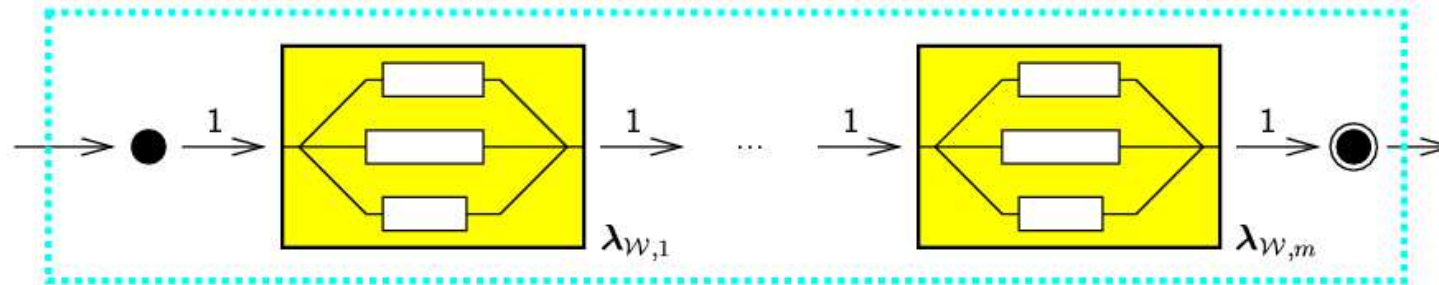
## Kompilierte HMM-Netzwerke

- Voraussetzung: einfache Wortübergangsgrammatik (*Bigramme*)
- Voraussetzung: jedes  $\lambda_w$  besitzt je einen A/E-Zustand
- Vernetzung der Wort-HMMs “im Sinne der Grammatik”
- Dekodierung: Viterbi-Algorithmus auf dem HMM-Netzwerk
- Optimale Zustandsfolge  $\Rightarrow$  Wahrscheinlichste Wortfolge

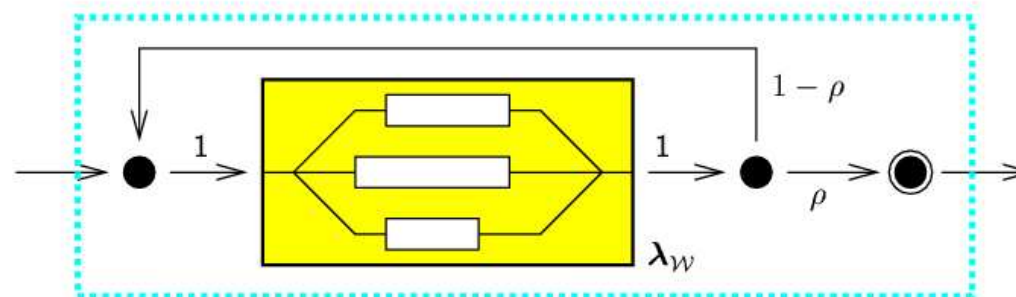
## HMM-Netzwerk zur Einzelworterkennung:



## Verbundwörtererkennung bei vorgegebener Satzlänge $m$



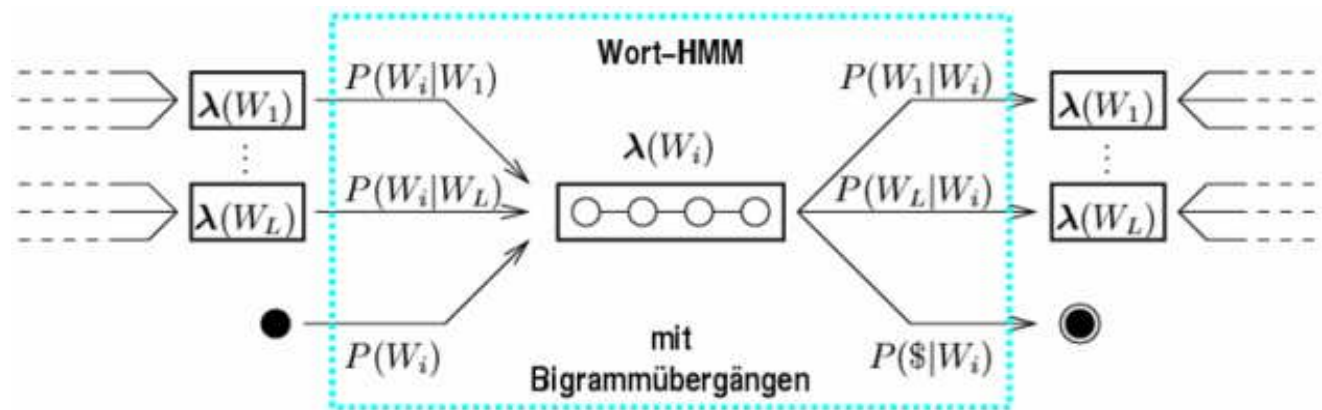
## Verbundwörtererkennung bei unbekannter Satzlänge



- **Fluchtwahrscheinlichkeit  $\rho$** 
  - Verlassen der Wiederholungsschleife
- Vermeiden der “Kantenexplosion” durch *konfluente* Zustände

## Verbundworterkennung

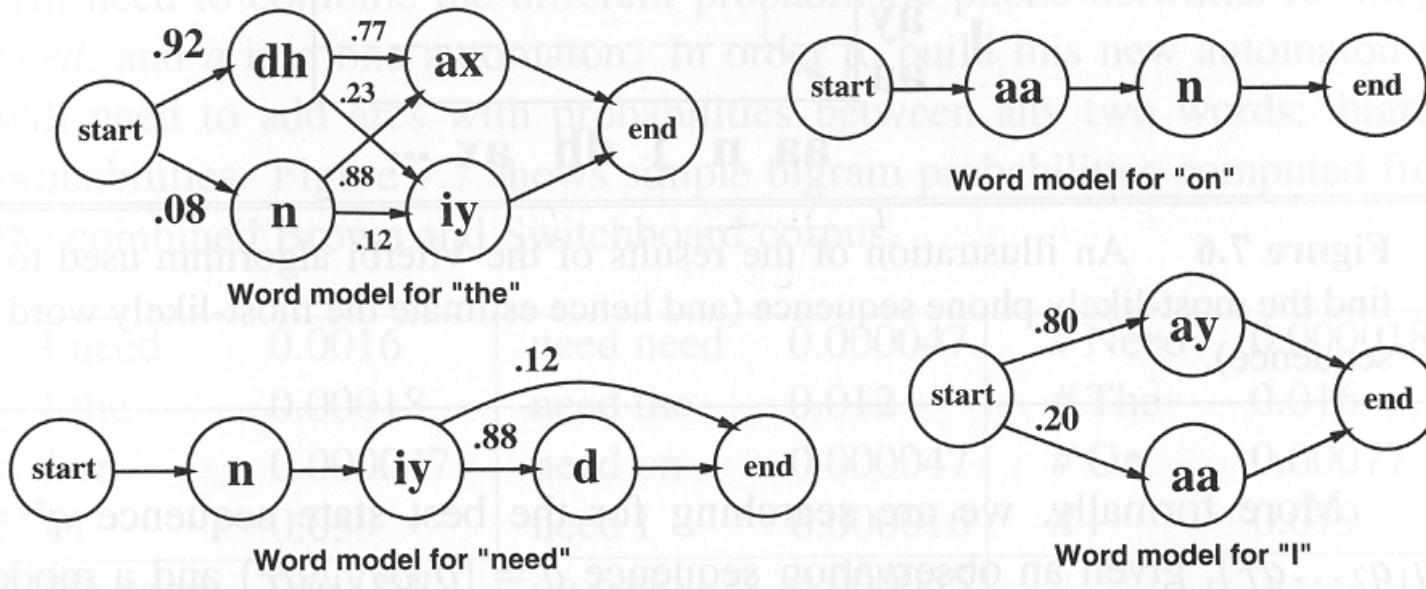
mit wortbezogener Bigrammgrammatik



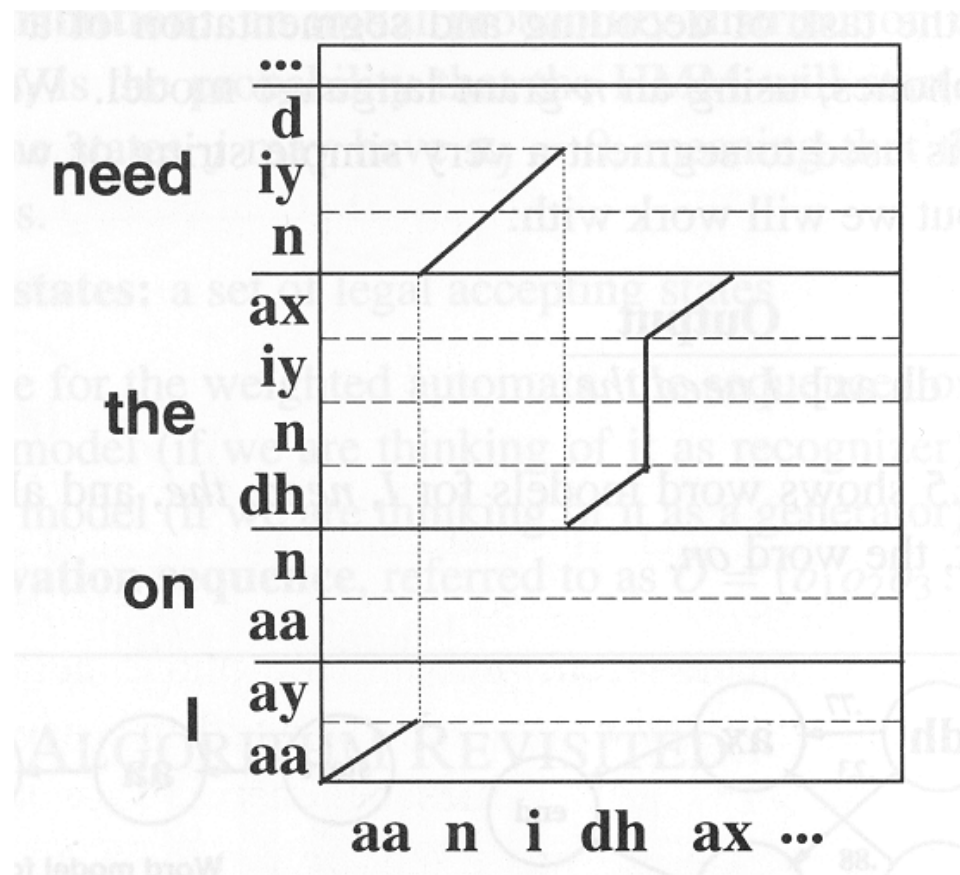
- $L$  Wortmodelle  
(wie ein Einzel- oder grammatikfreier Verbundworterkenner)
- $L^2$  Wort-HMM-Übergangskanten mit Bigrammwahrscheinlichkeiten

Example:      Input:      [aa n iy dh ax]  
                  Output:      *I need the*

Word models for *I*, *need*, *the*, and *on*:



Viterbi algorithm used to find the most-likely phone sequence within word models:



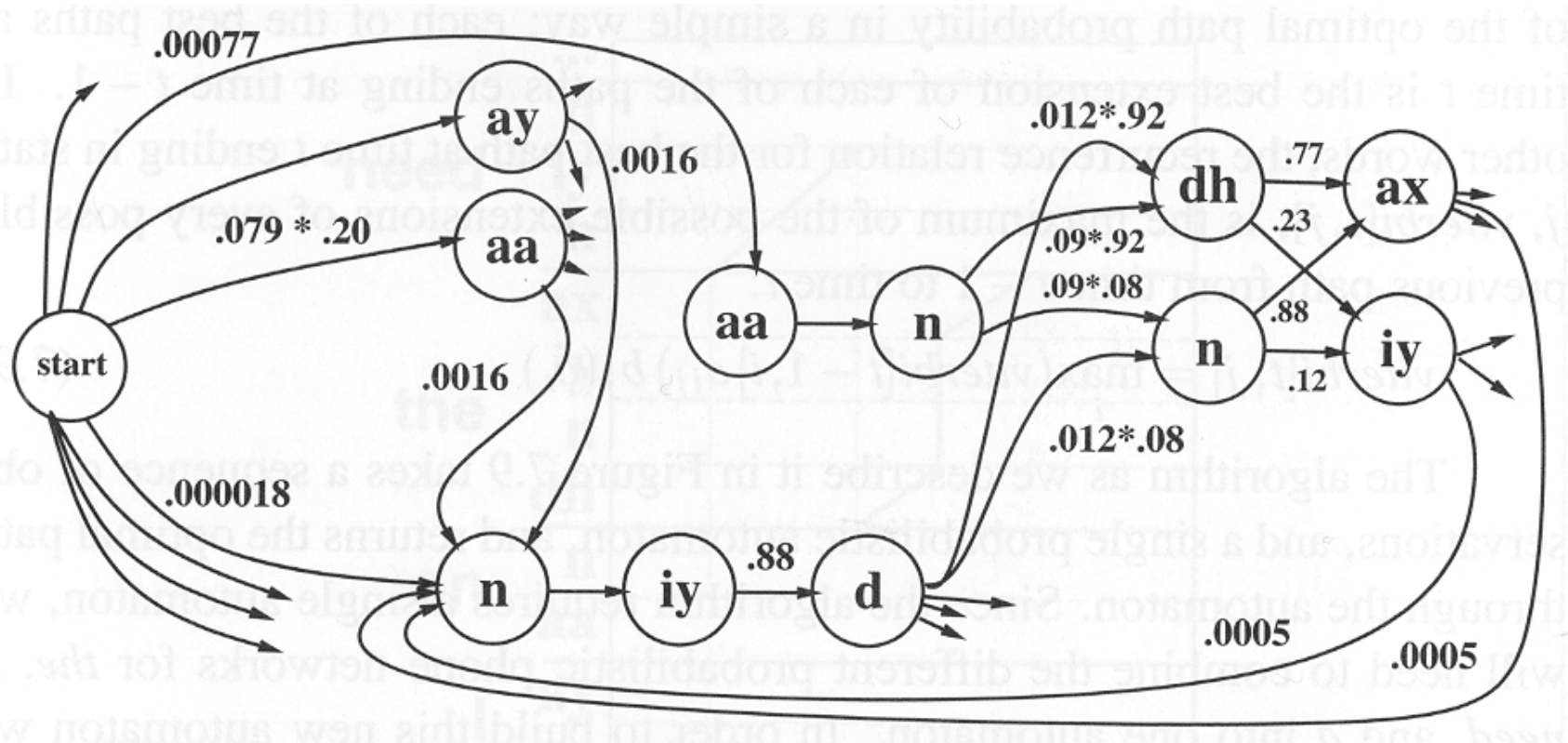


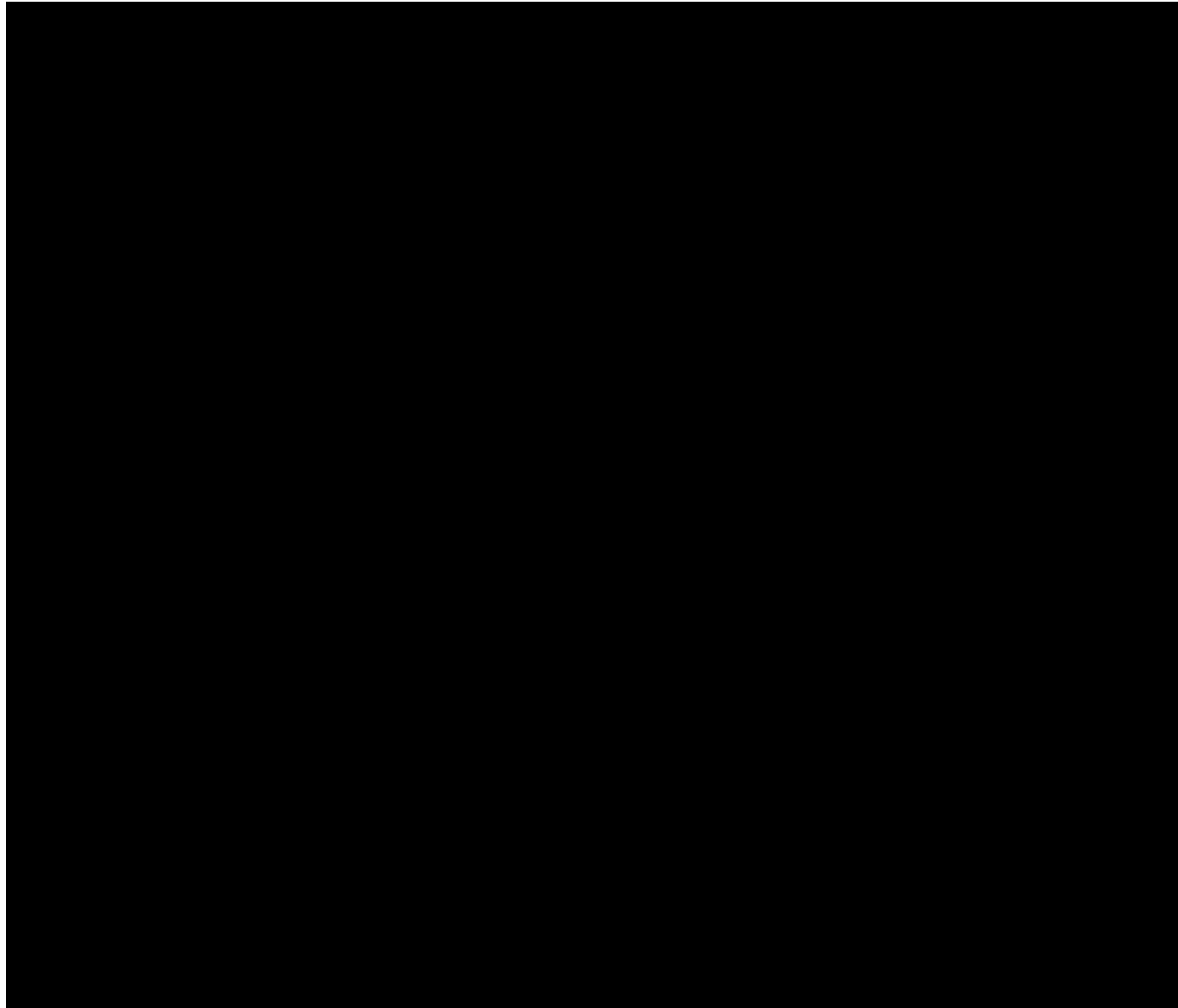


Example bigram probabilities for the words *the*, *on*, *need*, and *I* following each other, and starting a sentence (computed from Brown and Switchboard corpus).

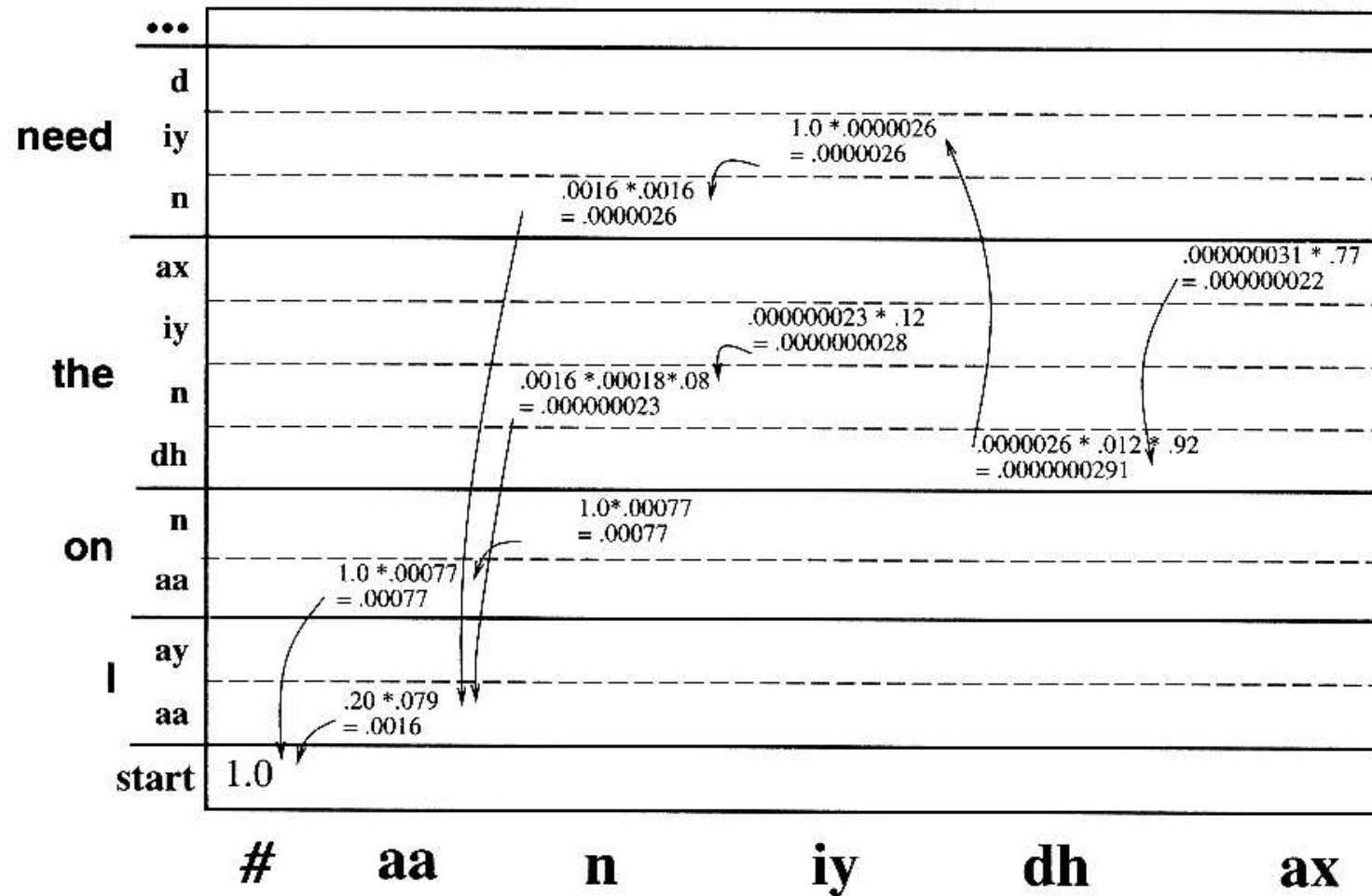
<i>I need</i>	0.0016	<i>need need</i>	0.000047	<i># Need</i>	0.000018
<i>I the</i>	0.00018	<i>need the</i>	0.012	<i># The</i>	0.016
<i>I on</i>	0.000047	<i>need on</i>	0.000047	<i># On</i>	0.00077
<i>II</i>	0.039	<i>need I</i>	0.000016	<i># I</i>	0.079
<i>the need</i>	0.00051	<i>on need</i>	0.000055		
<i>the the</i>	0.0099	<i>on the</i>	0.094		
<i>the on</i>	0.00022	<i>on on</i>	0.0031		
<i>the I</i>	0.00051	<i>on I</i>	0.00085		

Finite state machine for the words *I*, *need*, *on*, and *the*:



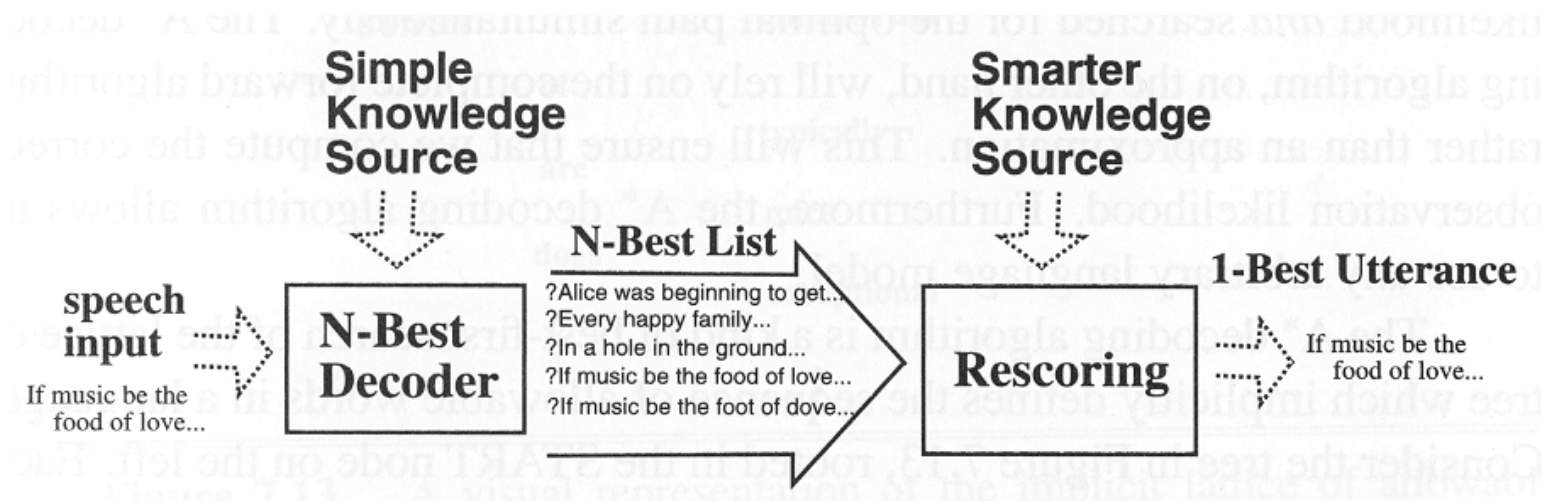


Viterbi decoding of [aa n iy dh ax]:



- Problems
  - Computes most likely *state-sequence* (e.g. phoneme sequence) not *word-sequence*
  - Only for *bigram* grammar models (in our example)
- Solutions
  - Viterbi algorithm as first decoding stage generates
    - ▶ *N-best word sequence hypotheses*, or
    - ▶ *word lattice* for further decoding stage(s)
  - Tree search: Stack decoder, A\* decoder
    - ▶ Based on the *forward algorithm*
    - ▶ Best-first *search in lattice/tree* of sequence of allowable words (in a language)

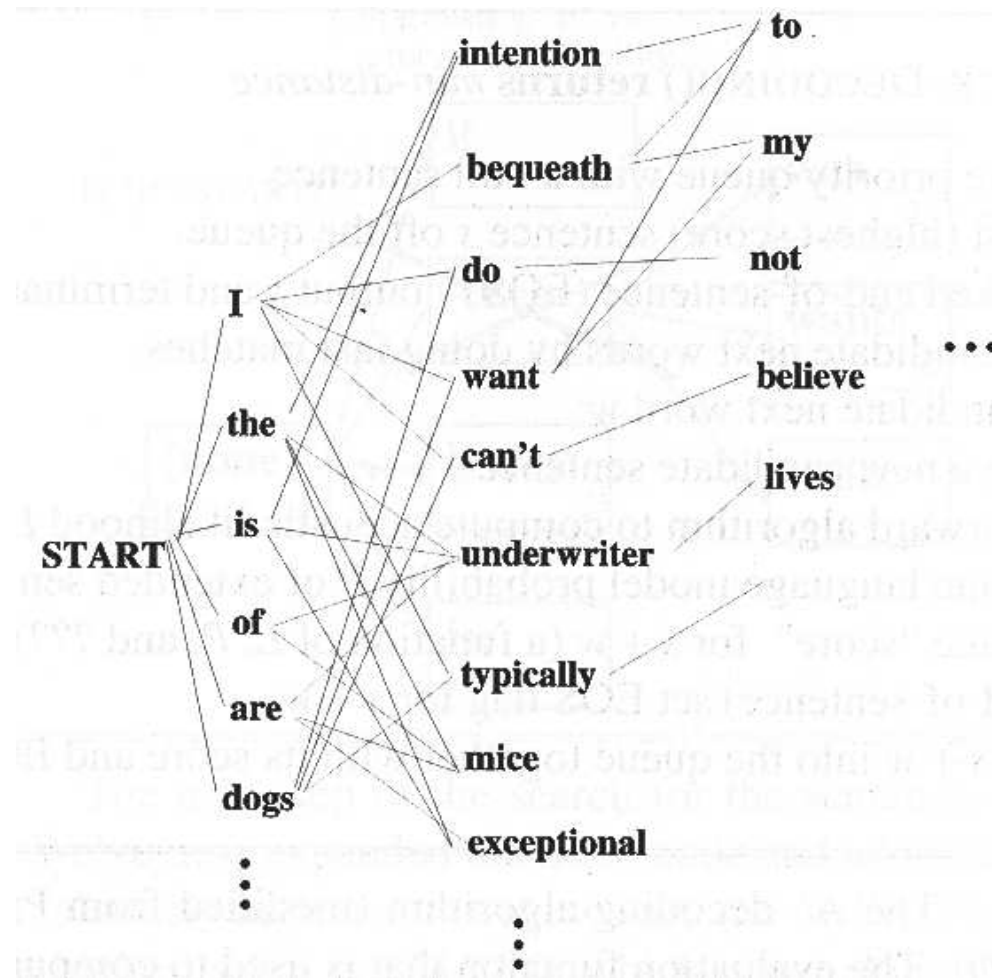
N-best decoding, e.g., as part of a two stage decoder:



N-Best Decoder:

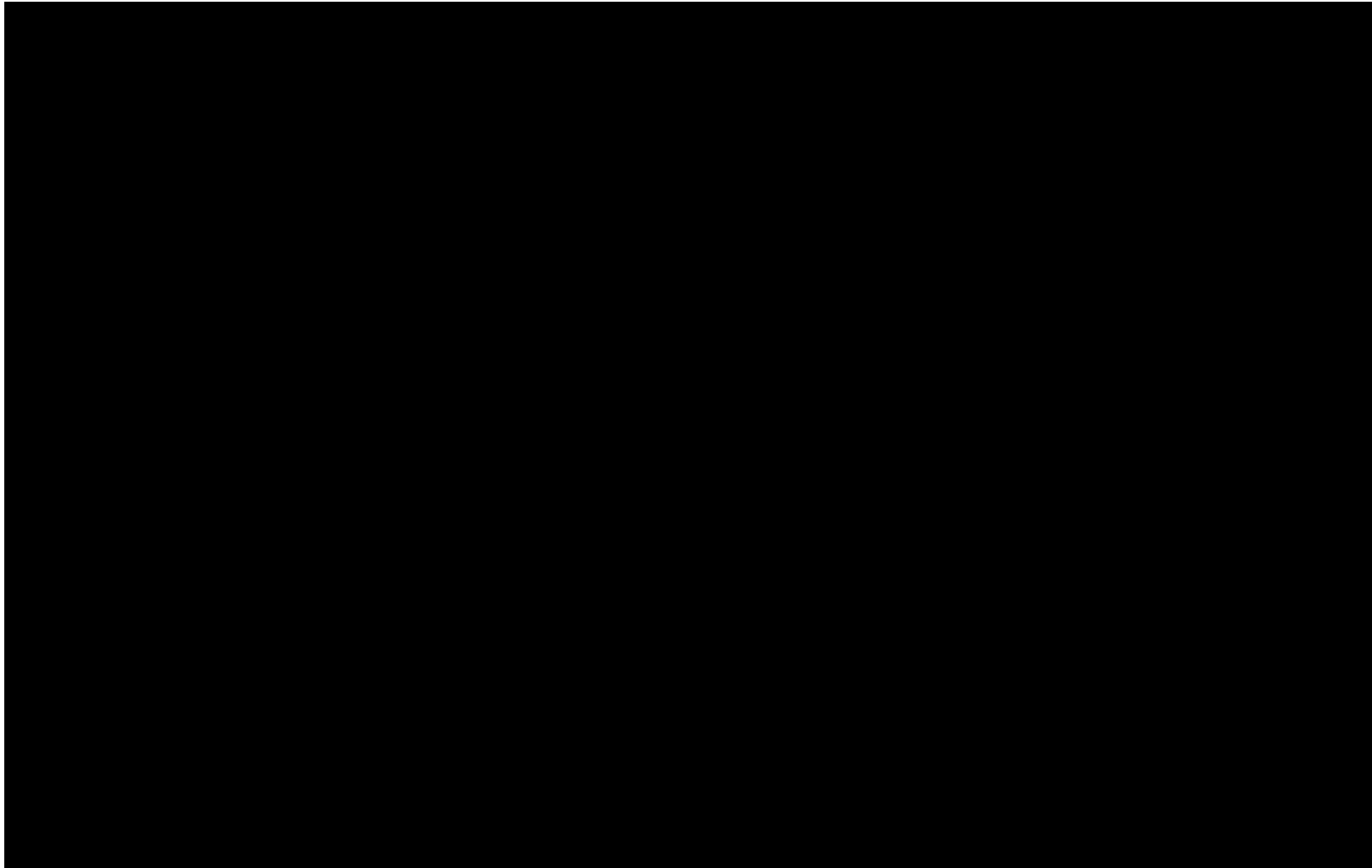
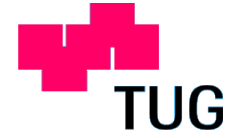
- Simple, fast
- Reduces search space for second stage

Lattice of allowable word sequences:





# Tree Search (Stack Decoding)

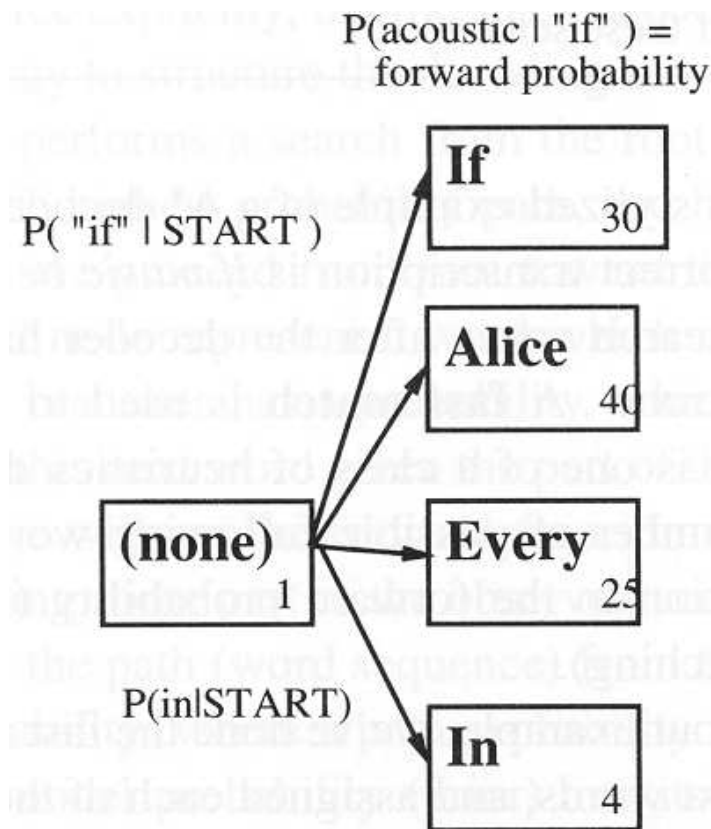






# Stack Dekoder (A\* Dekoder)

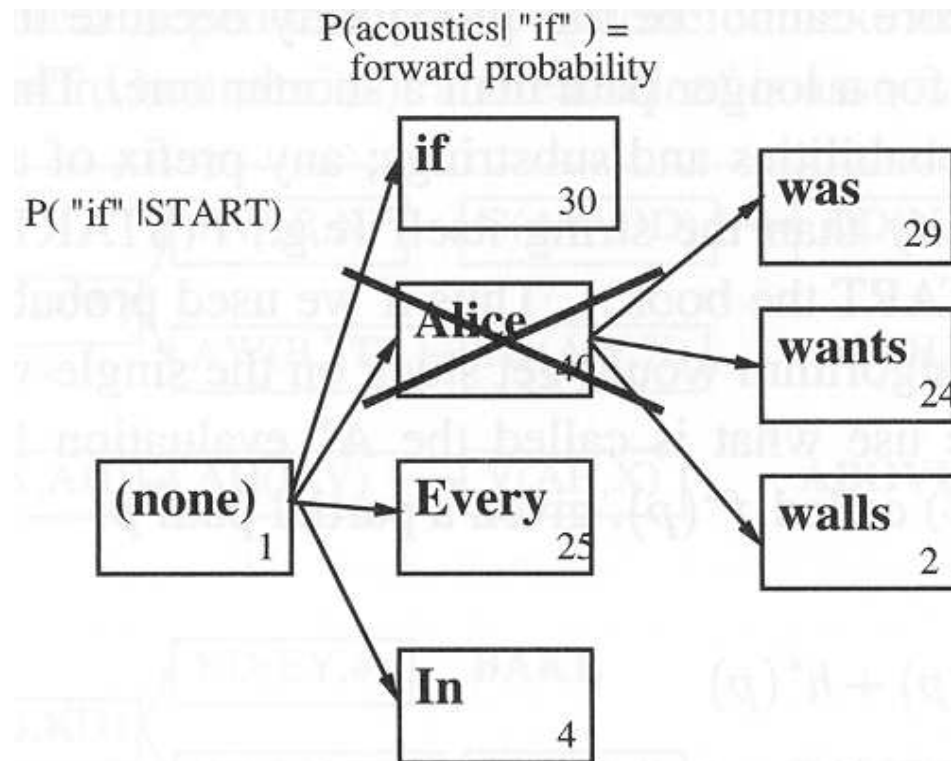
Begin of search for the sentence *If music be the food of love*.  
At this stage *Alice* is the most likely hypothesis.



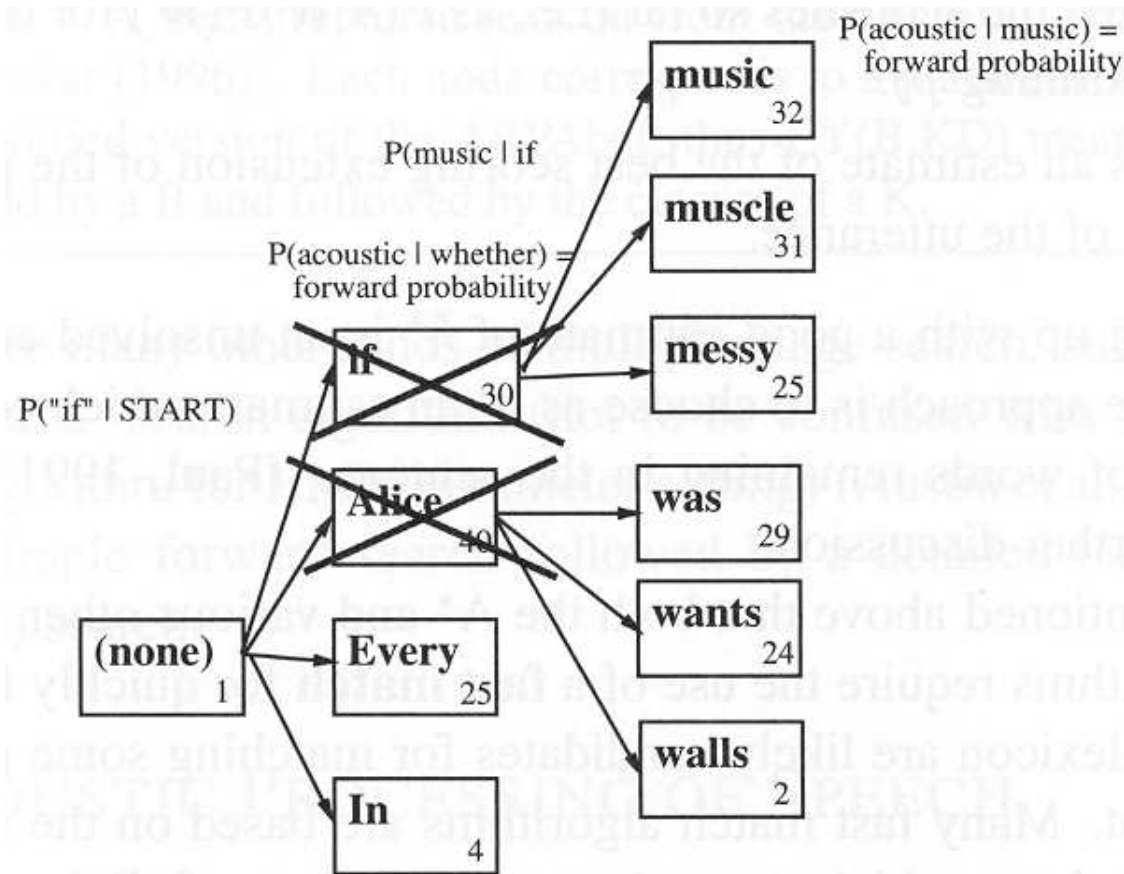
Next step: expand *Alice* node

⇒ three extensions with relatively high score

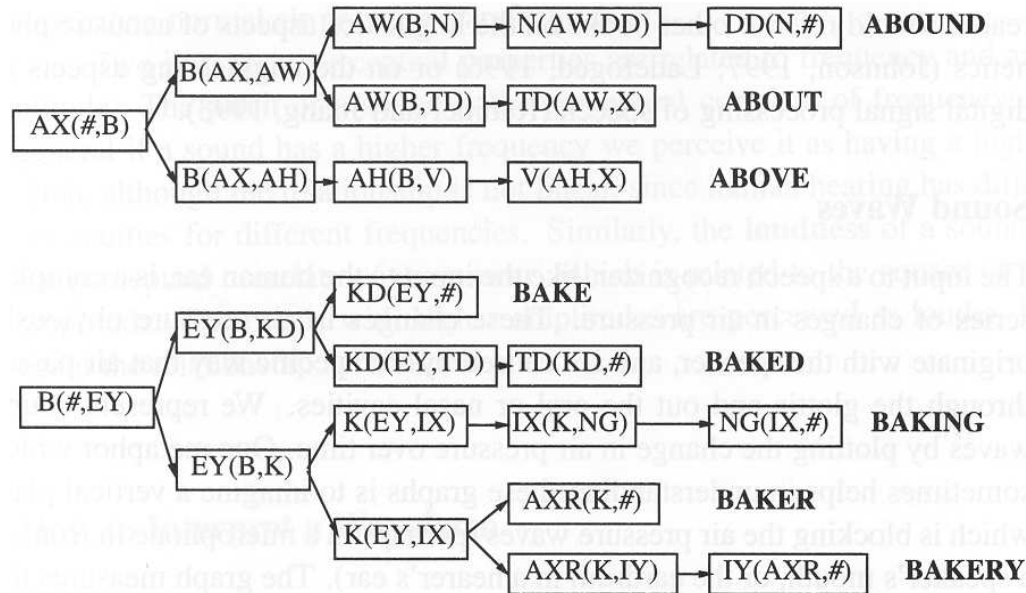
But: highest score for *START if*



Expand *if* node: highest score for *START if music*



- Score according to probability,  $P(X|W) P(W)$  ?
  - Gets smaller for longer word sequences
  - Score has to be extended with estimate of probability for the rest of an utterance
- Use **fast match** to find extensions for a word, e.g. based on tree structured lexicon:



- Kompiliertes HMM-Netzwerk
  - Viterbi Decoder
  - Vereint akustisches und linguistisches Modell
  - Synchrone Suche
- Tree Search
  - A\* Decoder
  - Heuristische Schätzung der Restwahrscheinlichkeiten
  - Asynchrone Suche
- Mehrstufige Dekodierung:
  - N-best Wortliste
  - Wortgitter (lattice)

- E.G. Schukat-Talamazzini, *Automatische Spracherkennung*, Vieweg-Verlag, 1995.
- D. Jurafsky and J.H. Martin, *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Prentice-Hall, 2000.
- F. Jelinek, *Statistical Methods for Speech Recognition*. MIT Press, 1997.