

Summer University

Speech Recognition

01 July – 30 August 2002



Lecture & Computer Experiments

Large Vocabulary

Speech Recognition

Andreas Wendemuth



Schedule:

Large Vocabulary Speech Recognition

- Refresh
- Gaussian densities
- What about LVSR?
- Integrating the Language Model
- Search techniques
- Lattice techniques

Modelle mit beschränktem Kontext

$$P(w_t | w_1 \dots w_{t-1}) \approx P(w_t | \underbrace{w_{t-n+1} \dots w_{t-1}}_{(n-1)\text{-Gramm-Kontext}})$$

Unigramm-, Bigramm- und Trigramm-Modelle

$$P_{1G}(\mathbf{w}) = \prod_{t=1}^T P(w_t)$$

$$P_{2G}(\mathbf{w}) = P(w_1) \cdot \prod_{t=2}^T P(w_t | w_{t-1})$$

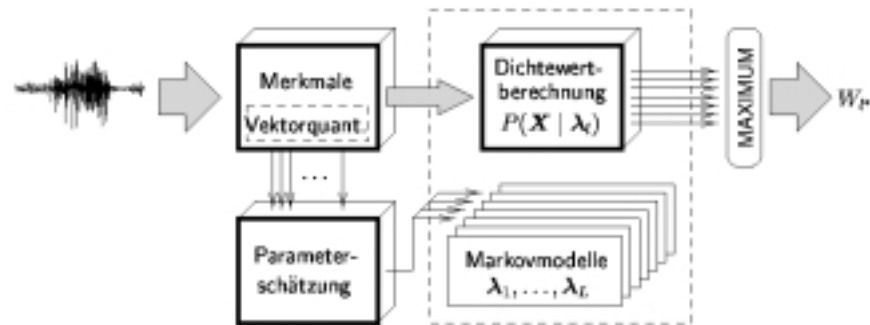
$$P_{3G}(\mathbf{w}) = P(w_1) \cdot P(w_2 | w_1) \cdot \prod_{t=3}^T P(w_t | w_{t-2}w_{t-1})$$

n -Gramm-Sprachmodelle

$$P_{nG}(\mathbf{w}) = \prod_{t=1}^T P(w_t | w_{t-n+1} \dots w_{t-1})$$

mit $w_0, w_{-1}, w_{-2}, \dots := \$$ (Satzanfangsmarkierungen)

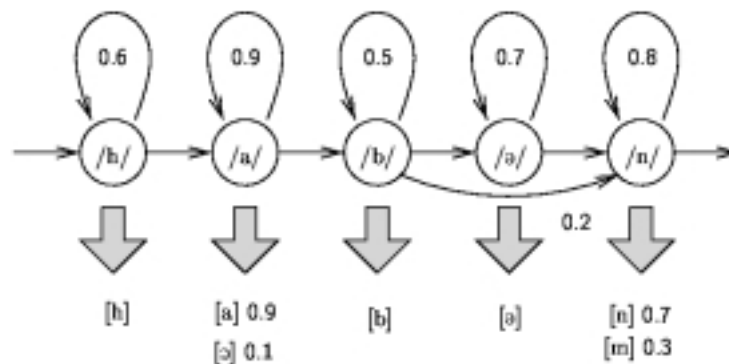
Statistische Wortdekodierung ...



$$l^* = \operatorname{argmax}_l P(W_l | \mathbf{X}) = \operatorname{argmax}_l \frac{P(\mathbf{X} | W_l) \cdot P(W_l)}{P(\mathbf{X})}$$

... mit wortbezog. „Markovmodellen“

Hidden Markov Model (HMM)



HIDDEN MARKOVMODELL mit DISKRETER Ausgabeverteilung

$$\lambda = (\pi, \mathbf{A}, \mathbf{B})$$

Normierungsbedingungen ($\pi_i \geq 0, a_{ij} \geq 0, b_{ik} \geq 0$)

$$\sum_{i=1}^N \pi_i = 1$$

$$\sum_{j=1}^N a_{ij} = 1 \quad (i = 1, \dots, N)$$

$$\sum_{k=1}^K b_{ik} = 1 \quad (i = 1, \dots, N)$$

3 offene Fragen zum Thema HMM:

⇨ Berechnung der **Produktionswahrscheinlichkeit**

$$P(\mathbf{X}|\lambda) = ?$$

⇨ Bestimmung der **wahrscheinlichsten Zustandsfolge**

$$P(\mathbf{q}, \mathbf{X} | \lambda) \xrightarrow{\mathbf{q}} \text{MAX}$$

⇨ Schätzung der bestpassenden **Modellparameter**

$$\hat{\lambda} = (\hat{\pi}, \hat{\mathbf{A}}, \hat{\mathbf{B}}) \text{ mit } P(\mathbf{X}|\hat{\lambda}) = \max_{\lambda} P(\mathbf{X}|\lambda)$$

... und die gesuchte

PRODUKTIONSWAHRSCHEINLICHKEIT

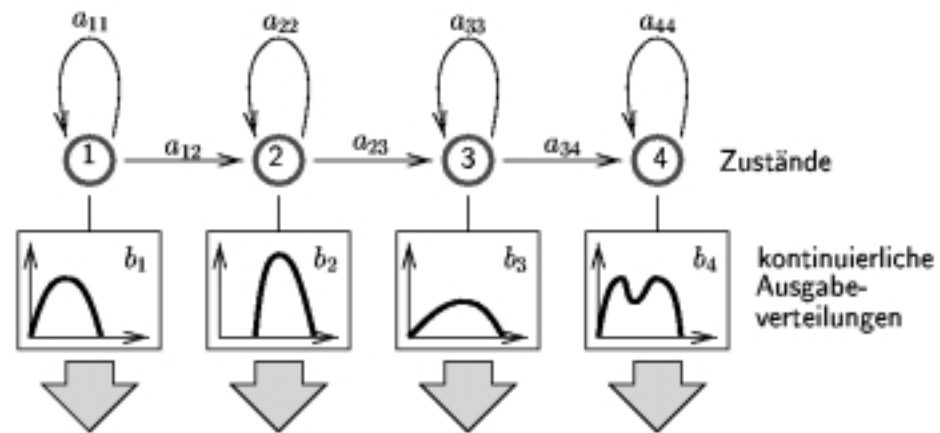
$$P(\mathcal{O} | \lambda) = \sum_{q \in \mathcal{Q}^T} P(\mathcal{O}, q | \lambda) = \sum_{q \in \mathcal{Q}^T} \pi_{q_1} b_{q_1}(O_1) \cdot \prod_{t=2}^T a_{q_{t-1}q_t} b_{q_t}(O_t)$$

⇔ ca. $2T \cdot N^T$ Multiplikationen (exponentielle Komplexität!)

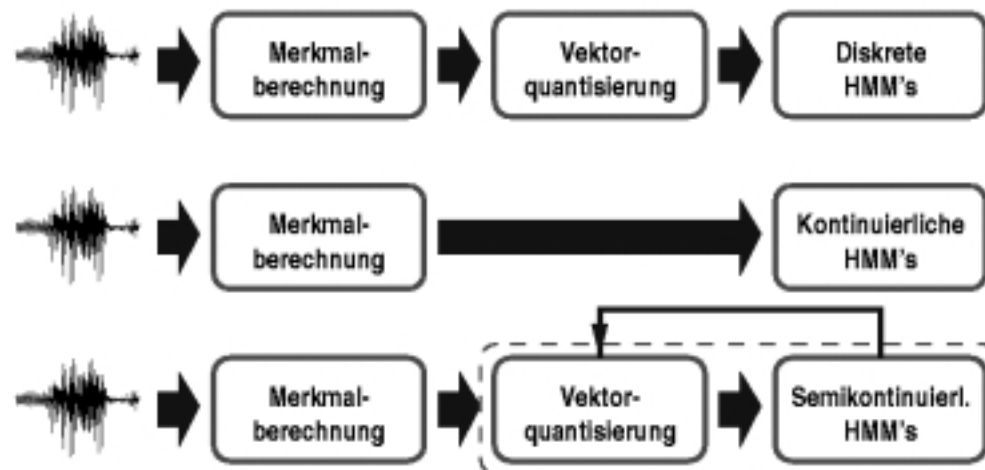
- From HMMs, we know that transitions and production probabilities give the total production prob. lined out above.
- These can be trained with Baum-Welch algorithm, using forward- and backward probabilities. This is a variant of Estimation-Maximization (EM) Techniques.
- Recognition (decoding) can be done with Viterbi algorithm.
- Now we are going to apply GAUSSIAN probability density functions:

Gaussian densities

HMM's mit kontinuierlichen Ausgabeverteilungen



Architektur von HMM-Spracherkennern:



Baum-Welch-Algorithmus

Berechnung „neuer“ Parameter $\hat{\lambda}$ aus den „alten“ λ mit

$$\mathcal{L}_{\text{HMM}}(\hat{\lambda}) \geq \mathcal{L}_{\text{HMM}}(\lambda)$$

a posteriori Übergangswahrscheinlichkeiten

für $s_t \rightarrow s_j$ zum Zeitpunkt t

$$\xi_t(i, j) := P(q_t = i, q_{t+1} = j \mid \mathbf{O}, \lambda) = \frac{\alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{i=1}^N \alpha_t(i) \beta_t(i)}$$

a posteriori Zustandswahrscheinlichkeiten

für s_t zum Zeitpunkt t

$$\gamma_t(i) := P(q_t = i \mid \mathbf{O}, \lambda) = \frac{\alpha_t(i) \beta_t(i)}{\sum_{j=1}^N \alpha_t(j) \beta_t(j)} = \sum_{j=1}^N \xi_t(i, j)$$

Baum-Welch-Formeln für diskretwertige HMMs

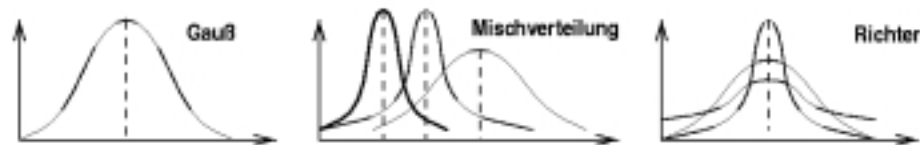
$$\begin{aligned} \hat{\pi}_i &= \gamma_1(i) = \frac{\alpha_1(i) \beta_1(i)}{\sum_{i=1}^N \alpha_1(i) \beta_1(i)} \\ \hat{a}_{ij} &= \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \gamma_t(i)} = \frac{\sum_{t=1}^{T-1} \alpha_t(i) a_{ij} b_j(O_{t+1}) \beta_{t+1}(j)}{\sum_{t=1}^{T-1} \alpha_t(i) \beta_t(i)} \\ \hat{b}_{jk} &= \frac{\sum_{t=1}^T \gamma_t(j) \mathbf{x}_{[O_t=v_k]}}{\sum_{t=1}^T \gamma_t(j)} = \frac{\sum_{t=1}^T \alpha_t(j) \beta_t(j) \mathbf{x}_{[O_t=v_k]}}{\sum_{t=1}^T \alpha_t(j) \beta_t(j)} \end{aligned}$$

Es sei $\mathbf{x}_{[c]} = 1$ für zutreffende Aussagen und 0 sonst.

HMM's mit normalverteilter Ausgabe

Ausgabeverhalten des Zustandes s_j

$$b_j(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$



Die Dekomposition der **Kullback-Leibler-Statistik** $Q(\cdot, \cdot)$ ergibt

$$Q_{\hat{b}_j}(\boldsymbol{\lambda}, \hat{\boldsymbol{\mu}}_j, \hat{\boldsymbol{\Sigma}}_j) = \sum_{t=1}^T \gamma_t(j) \log \hat{b}_j(\mathbf{x}_t) = \sum_{t=1}^T \gamma_t(j) \log \mathcal{N}(\mathbf{x} \mid \hat{\boldsymbol{\mu}}_j, \hat{\boldsymbol{\Sigma}}_j)$$

Schätzformeln für die NV-Dichteparameter:

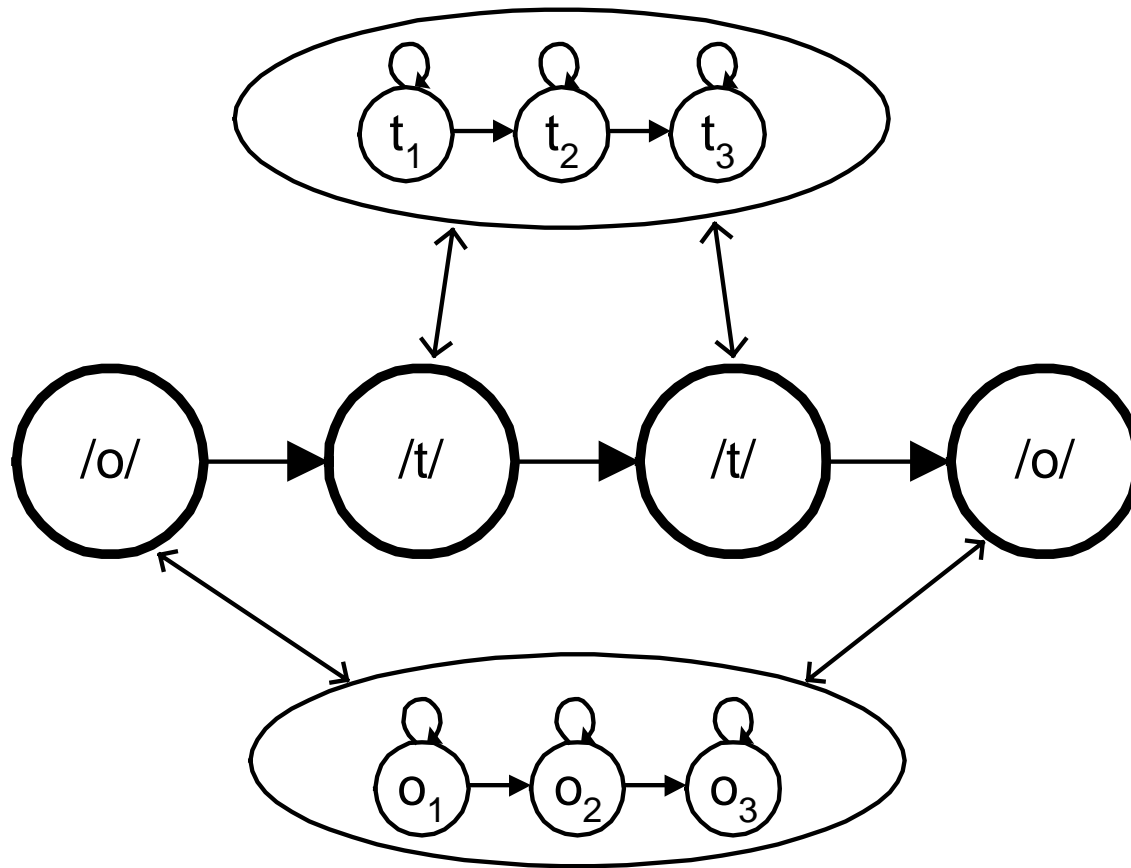
$$\begin{aligned} \hat{\boldsymbol{\mu}}_j &= \frac{1}{\sum_t \gamma_t(j)} \sum_{t=1}^T \gamma_t(j) \mathbf{x}_t \\ \hat{\boldsymbol{\Sigma}}_j &= \frac{1}{\sum_t \gamma_t(j)} \sum_{t=1}^T \gamma_t(j) (\mathbf{x}_t - \hat{\boldsymbol{\mu}}_j)(\mathbf{x}_t - \hat{\boldsymbol{\mu}}_j)^\top \\ &= \frac{1}{\sum_t \gamma_t(j)} \sum_{t=1}^T \gamma_t(j) \mathbf{x}_t \mathbf{x}_t^\top - \hat{\boldsymbol{\mu}}_j \hat{\boldsymbol{\mu}}_j^\top \end{aligned}$$

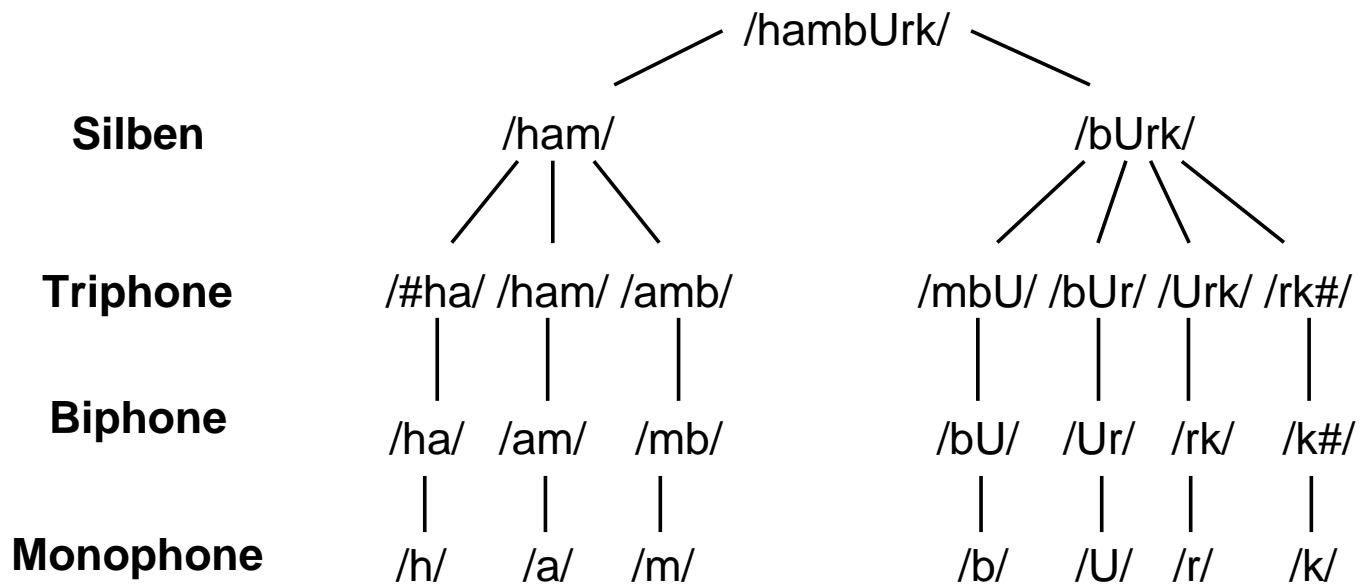
Why not word models?

- **Statistische Unzuverlässigkeit:** Für die Schätzung von Wortmodellen wäre eine Unmenge an Trainingsdaten erforderlich um Ganzwortmodelle zuverlässig zu schätzen. Bei wachsendem Erkennungsvokabular ist beispielsweise für ein Spracherkennungsprogramm bezüglich des Trainingsaufwandes die Grenze des Zumutbaren schon bei kleinerem Vokabular erreicht.
- **Rechen- und Speicherbedarf:** Für großes Vokabular steigt die Speicher und Rechenaufwand enorm an.
- **Geringe Flexibilität:** Bei der Verwendung von Wortmodellen ist eine Erkennung von Wörtern außerhalb des Erkennungsvokabulars nicht möglich.
- **Geringe Modularität:** Ganzwortmodelle berücksichtigen nicht die von der Satzkonstruktion abhängigen Verschleifungen am Wortrand.

Phoneme models!

- **Phonemes are the smallest entities of language which distinguish meaning**
- **Example: [f]ather, [r]ather**
- **There are 50 phonemes /k/,
4000 triphones [/a/ /th/ /r/]
20,000 syllables „rath“,
200,000 wordforms „fathers“**
- **Therefore phonemes are for acoustic speech what
letters are for written speech.**

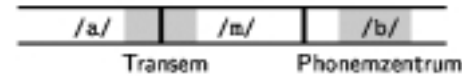
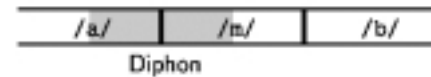
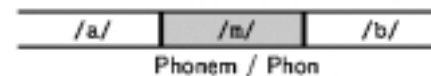
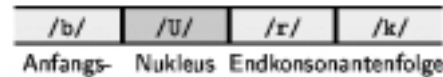
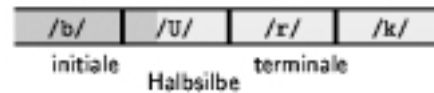




Word / Pronunciation Lexicon

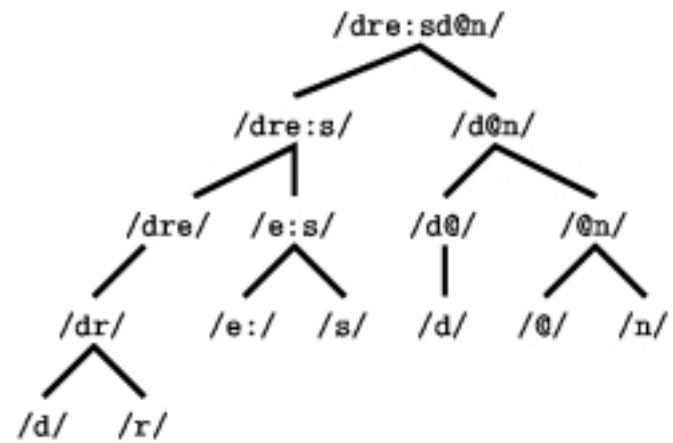
- Zero Z IH R OW
- One W AH N
- Two T UW
- Three TH R IY
- Four F OW R
- Five F AY V
- Six S IH K S
- Seven S EH V AX N
- Eight EY T
- Nine N AY N
- Oh OW

Phonologisch orientierte Wortuntereinheiten

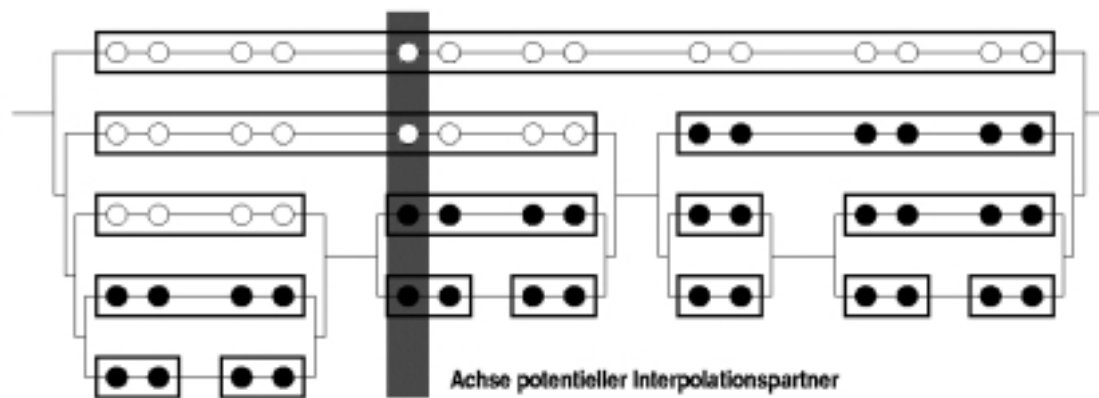


Phone	je nach Differenzierungsgrad 40–200 universelle Einheiten (+modular, –präzis)
Phoneme	je nach Sprache 20–60 (+modular, –präzis)
Silben	20 000 im Englischen, 100 im Japanischen; Koartikulation primär innerhalb (+präzis, –modular)
Halbsilben	800/2560 initiale/terminale im Deutschen (+trennscharf, ±modular)
Sylparts	47 AKF, 20 Nuklei, 159 EKF im Deutschen (guter Kompromiß)
Diphone	1000–1500 Einheiten (engl./ital.), ungünstige Nahtstellen (besser: Transeme)
Doppelhalbsilbe	Silbenkern–Silbenkern, 2 Mill. im Deutschen (++)trennscharf, --modular, --robust)

Hierarchische Wortrepräsentation



Hierarchische Modellstruktur



Schreibweise für **Triphone**

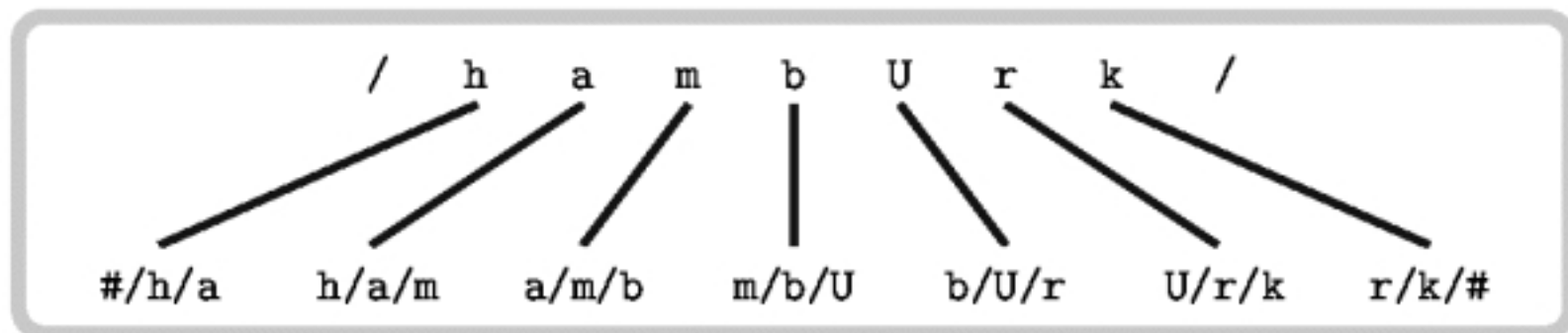
r in hambUrK \Rightarrow U/r/k

und für **rechte / linke Biphone** und **Monophone**

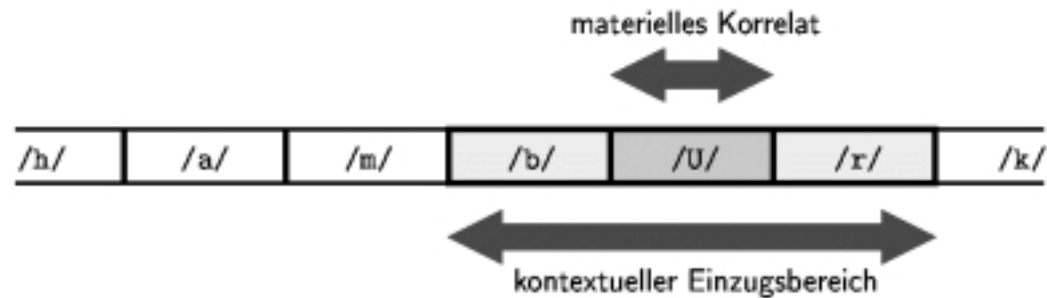
r \rightarrow /r/k

r \rightarrow U/r/

r \rightarrow /r/



Triphonmodelle



Training von Triphon-HMM's

1. Initialisierung / Training gewöhnlicher **Monophonmodelle**
2. Training der **Biphonmodelle** (⇔ Monophonparameter)
3. Training der **Triphonmodelle** (⇔ Biphonparameter)
4. Konstruktion der Wortmodelle des Erkennungsvokabulars:

◦ **Rückgriffstrategie:**

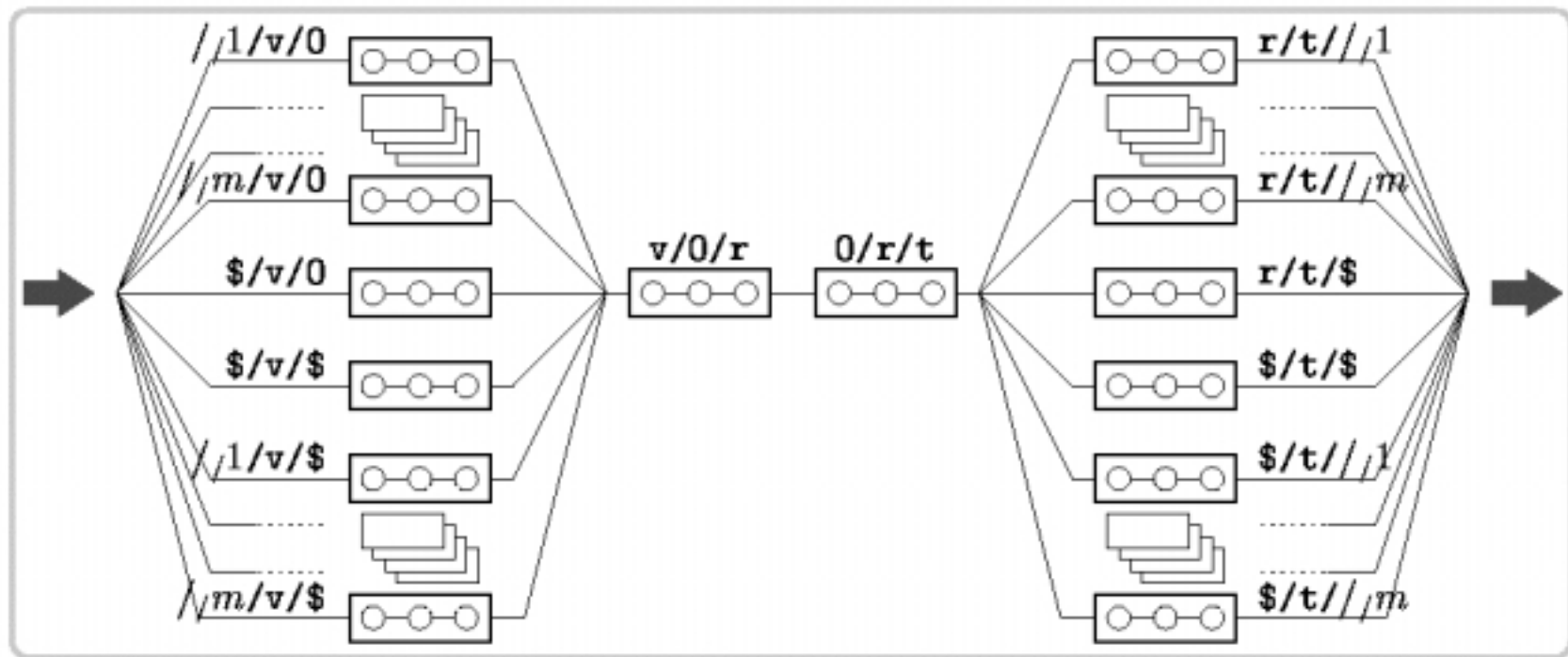
$$\lambda(b/ʊ/r) \rightsquigarrow \lambda(/ʊ/r) \rightsquigarrow \lambda(b/ʊ/) \rightsquigarrow \lambda(/ʊ/)$$

◦ **Interpolation:** Gewichte *uniform, heuristisch*, nach *Kopfzahl* oder EM-optimiert

⇨ **wortabhängige** Phonmodelle für häufige *Funktionswörter*
(stellen 3% des Vokabulars, 30% der Vorkommen und 60-75% der Fehler)

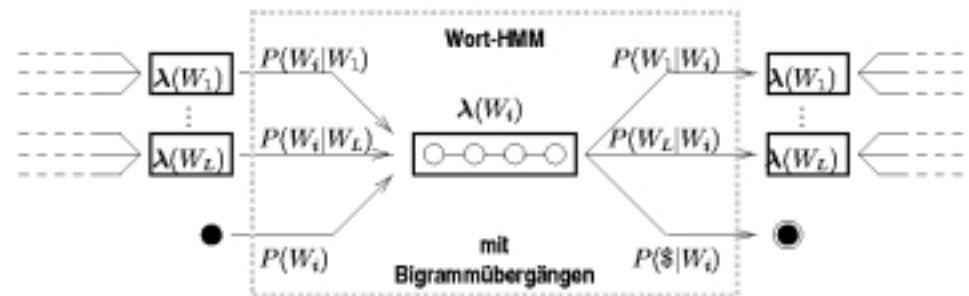
⇨ **prosodische Faktoren**

Akzentmarken; Silben-, Wort- und Morphemgrenzen



Kontextbedingt verästelted initiales/terminales Triphonmodell
des Wortes **v0rt** („Wort“)

Verbundworterkennung mit wortbezogener Bigrammgrammatik



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

Wann „realisiert“ ein HMM-Netz eine Grammatik?

$$P(\mathbf{X}, \mathbf{w} \mid \lambda) = P(\mathbf{X} \mid \lambda(\mathbf{w})) \cdot P(\mathbf{w})$$

mit

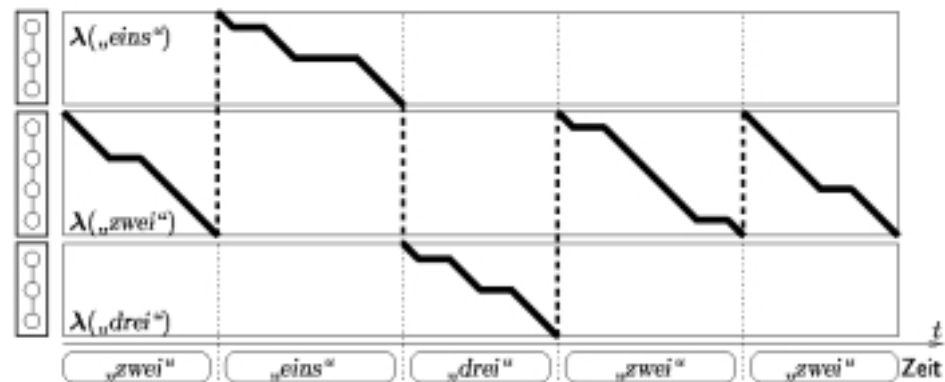
$$P(\mathbf{X}, \mathbf{w} \mid \lambda) := \sum_{q \in Q^T} P(\mathbf{X}, q \mid \lambda)$$

(Summation über alle q mit $w(q) = \mathbf{w}$)

Viterbi-Dekodierung

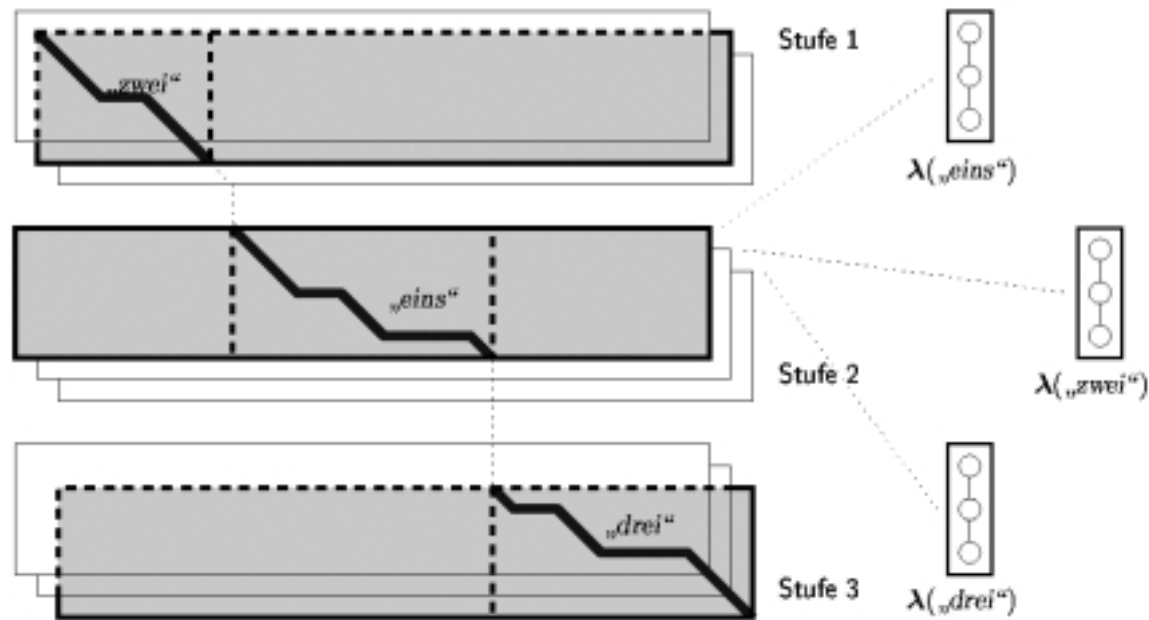
- Viterbi-Algorithmus auf dem HMM-Netzwerk
- Wahrscheinlichste Zustandsfolge \mathbf{q}^*
- Resultat ist $\mathbf{w}^* = w(\mathbf{q}^*)$
- ⇨ \mathbf{w}^* ist optimale Segmentierung bzgl. $P^*(\mathbf{X}|\lambda(\mathbf{w}))$

Einstufige Verbundwortdekodierung



- *Viterbi-Algorithmus auf dem Rückkopplungsmodell*
- „kurze“ Sprünge innerhalb eines Wortmodells
- „weitere“ Sprünge zwischen zwei Wortmodellen
- Rechenaufwand $T \cdot L \cdot \sum_t N_t$

Mehrstufige Verbundwortdekodierung



- *Viterbi-Algorithmus auf dem m -Wörter-Modell*
- Keine echten Zyklen im HMM-Netz
- ⇒ *synchrone* wie auch *asynchrone* $\vartheta_t(i)$ -Berechnung möglich
- Rechenaufwand $m \cdot T \cdot L \cdot \sum_l N_l$

Vorwärtsdekodierung

Näherungsweise Wortsegmentierung bzgl. $P(\mathbf{X}|\boldsymbol{\lambda}(w))$

- α -Berechnung (Summe) *innerhalb* der Wortmodelle
- ϑ -Berechnung (Max) *zwischen* den Wortmodellen

- **Initialisierung:** Setze für alle $j = 1, \dots, N$

$$\vartheta_1(j) = \pi_j b_j(\mathbf{x}_1) \quad \text{und} \quad \psi_1(j) = 0$$

- **Rekursion:** Für alle $j = 1, \dots, N$ setze

$$\psi_t(j) = \operatorname{argmax}_i \vartheta_{t-1}(i) a_{ij}$$

sowie

$$\vartheta_t(j) = \begin{cases} \max_i (\vartheta_{t-1}(i) a_{ij}) \cdot b_j(\mathbf{x}_t) & \text{falls } s_j \text{ Wortanfangszustand ist} \\ \sum_i (\vartheta_{t-1}(i) a_{ij}) \cdot b_j(\mathbf{x}_t) & \text{für alle sonstigen } s_j \end{cases}$$

- **Terminierung:** Setze

$$P^*(\mathbf{X} | \boldsymbol{\lambda}) = \vartheta_T(N) \quad \text{und} \quad q_T^* = \psi_T(N)$$

- **Rückverfolgung:** Für $t = T-1, \dots, 1$ setze

$$q_t^* = \psi_{t+1}(q_{t+1}^*)$$

- **Lösungswortkette:** Setze $w^* = w(q^*)$.

Language Model Recombination:

- **Scores: $\log P(A,W) = \log P(A|W) + \log P(W)$**
- **But the word chain probability $P(W)$ is only known after closure of the history**
- **Therefore, keep copies of (acoustic) history edges until new word transition**
- **Then, compute partial best score and close all sub-optimal histories which are the most in the past**
- **Example: current word: „is“, LM: trigram, 4 edges so far:**
The tree is, A tree is, The three is, <start> Three is
- **Best score: The tree is -> close „The“**
- **Next word investigated: „green“. New hypotheses are:**
Tree is green, three is green

Strahlsuche („beam search“)

- Verfolgt kleine Schar *wahrscheinlichster* aktueller Zustände
- Anzahl der Kandidaten (Strahlbreite) ist *adaptiv*
- Reformulierung des Viterbi-Algorithmus notwendig

Operationen des VA beim Fortschreiten $t \rightarrow t + 1$:

$$\triangleright \text{ für alle } j: \quad \vartheta_{t+1}(j) \leftarrow 0$$

$$\triangleright \text{ für alle } i, j: \quad \vartheta_{t+1}(j) \leftarrow \max \left\{ \begin{array}{l} \vartheta_{t+1}(j) \\ \vartheta_t(i) \cdot a_{ij} \cdot b_j(\mathbf{x}_{t+1}) \end{array} \right\}$$

Jede der N^2 Maximumop's erübrigt sich, wenn einer der 3 Faktoren = 0 wird!

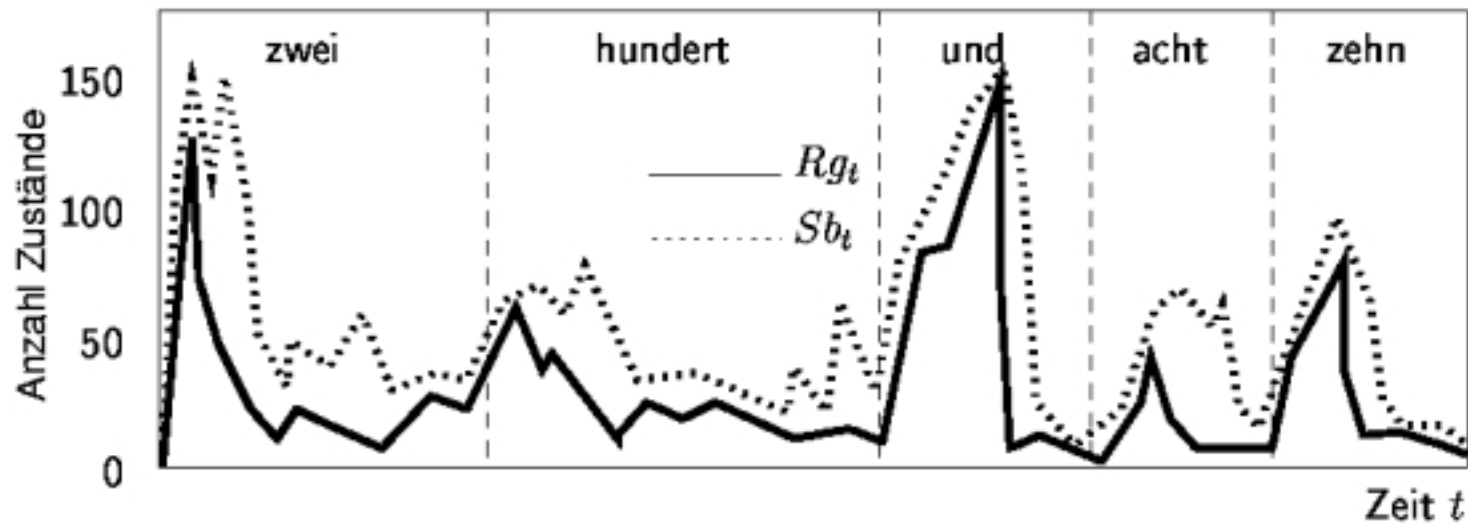
Menge der aktiven Zustände

$$\mathcal{O}_t = \{i \mid \vartheta_t(i) \neq 0\}$$

Beschneidungsstrategie ($B_0 > 0$ sehr klein):

$$\mathcal{O}_t = \{i \mid \vartheta_t(i) \geq B_0 \cdot \Lambda_t\} \quad \text{mit} \quad \Lambda_t = \max_j \vartheta_t(j)$$

Strahlbreite & Rang der wahrscheinl. Zustandsfolge



Die n besten Wortketten

(Näherungsweise) Berechnung mit n -Best-Viterbi-Algorithmen (NVA)

- **Zustandsbezogener NVA**

hält in jedem Gitterpunkt (t, j) die n besten Kandidaten in bewertungssortierter Liste $\mathcal{D}_t(j)$ und berechnet

$$\vartheta_t^{(k)}(j) = \max^{(k)} \underbrace{\left\{ \vartheta_{t-1}^{(l)}(i) \cdot a_{ij} \cdot b_j(\mathbf{x}_t) \mid 1 \leq i \leq N, 1 \leq l \leq n \right\}}_{\mathcal{D}_t(j)}$$

- **Satzbezogener NVA**

rekombiniert konkurrierende Kandidaten für gleiche Wortfolgen

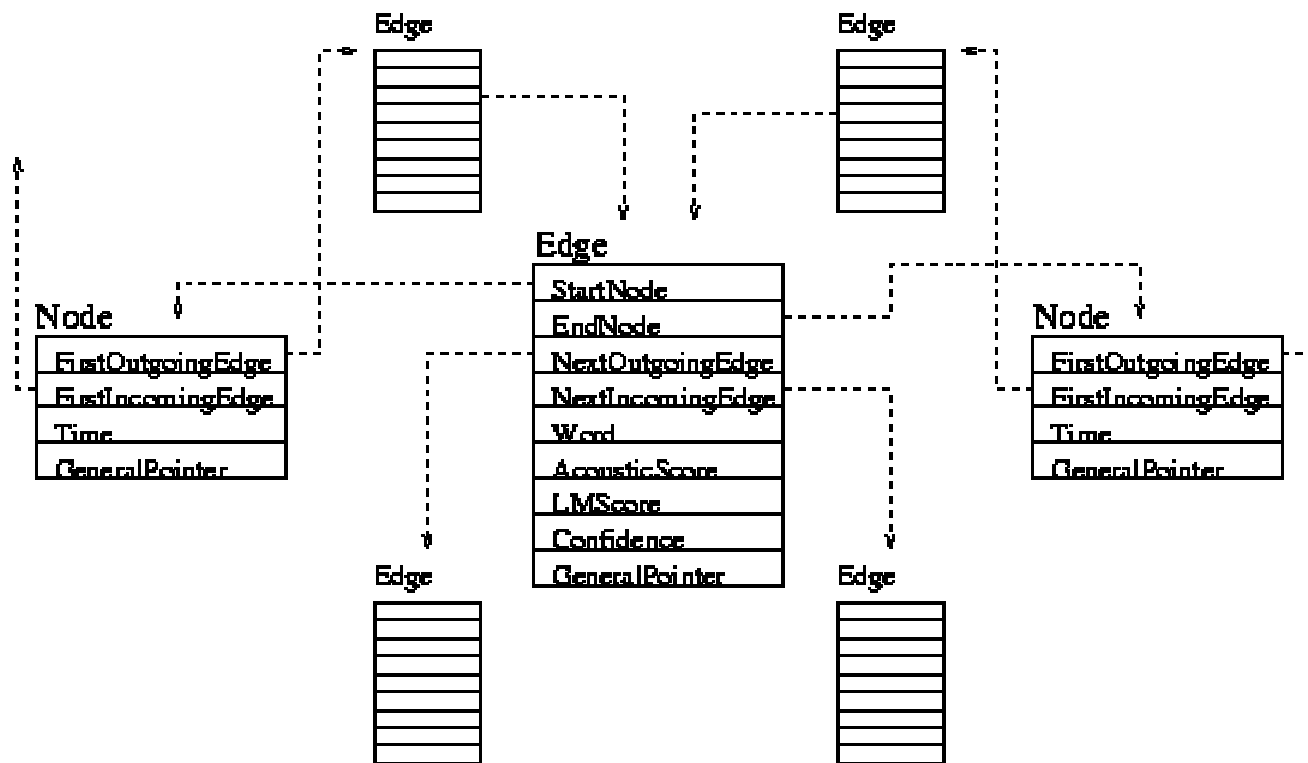
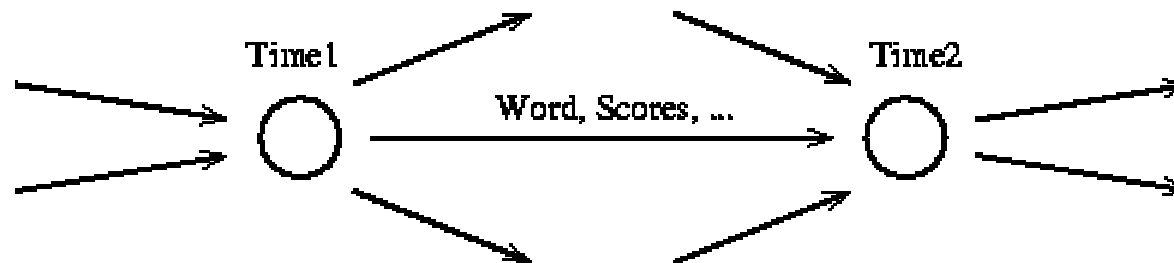
- **Gitterbezogener NVA**

keine Listen im Wortinneren, nur das *dichte Wortgitter*

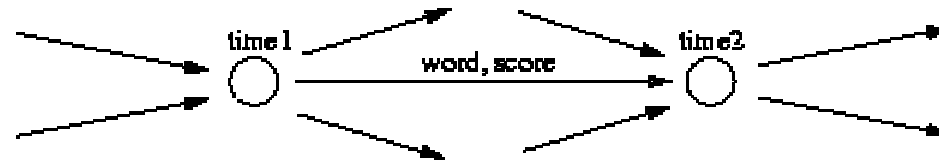
$$\{ \langle w, \vartheta_t(w), \tau_t(w) \rangle \mid t = 1, \dots, T, w \in \mathcal{W} \}$$

Lattices:

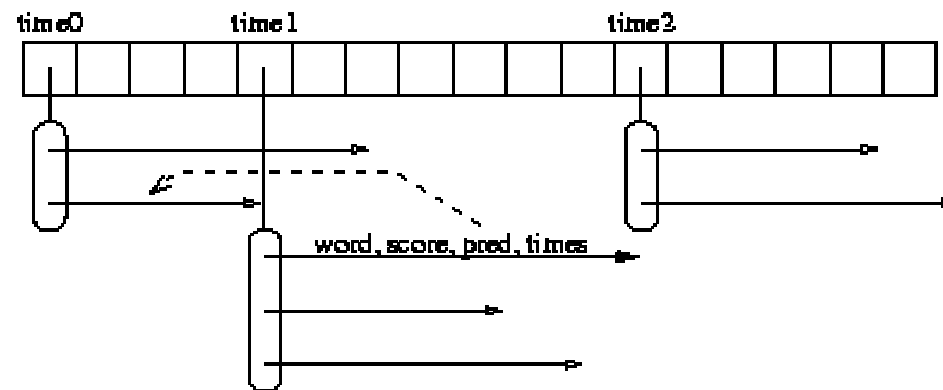
- **Use rather general LMs (Shakespeare bigram) and keep many edges (dense lattice)**
- **Write all edges with acoustic and Lm Score into „lattice“ = mostly time organised graph.**
- **On dense lattice, do more specific „rescoring“: examples:**
 - **More specific LMs (Hamlet trigram)**
 - **Silence penalties**
 - **Word transition penalties**
 - **LM penalties**



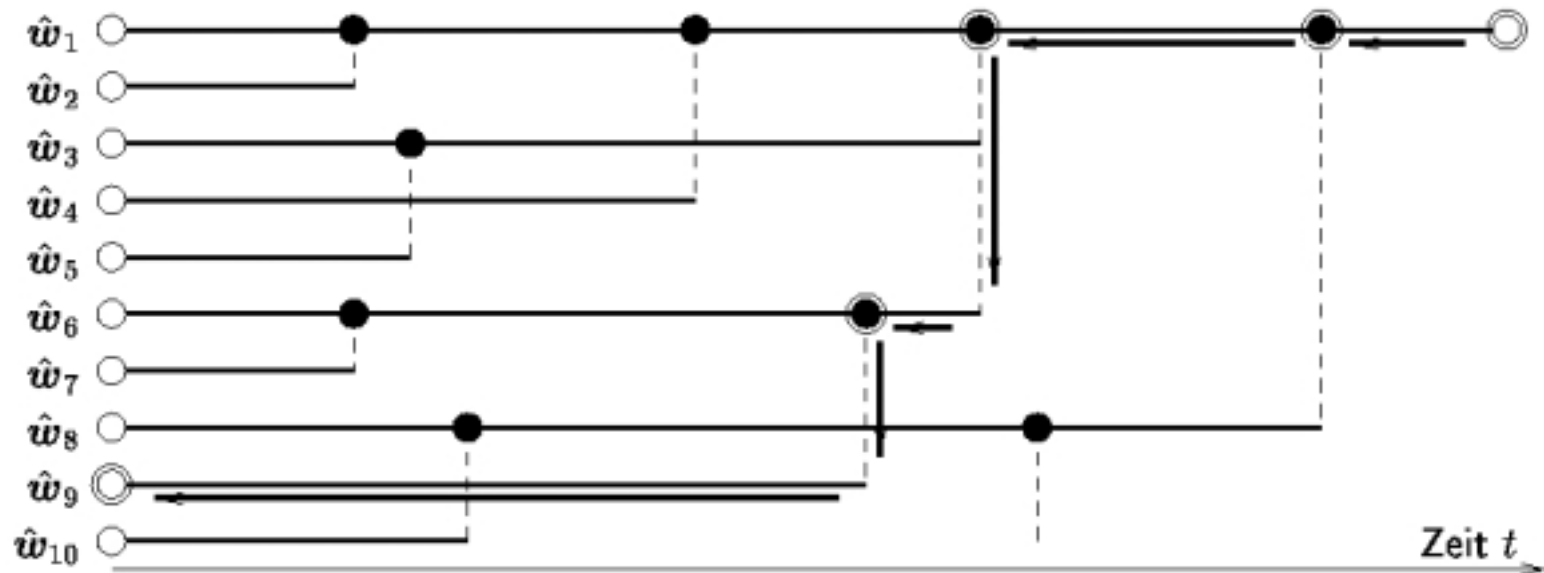
Example of new graph oriented lattice:



Example of old edge oriented lattice:



Rekonstruktion der 10 besten Ketten aus dem dichten Wortgitter



Lattice:

Header

N=24 L=39

Node definitions

I=0 t=0.00

I=1 t=0.25

I=2 t=0.26

I=3 t=0.61

Link definitions.

J=0	S=0	E=1	W=!ENTER	v=0	a=-1432.27	l=0.00
J=1	S=0	E=2	W=!ENTER	v=0	a=-1500.93	l=0.00
J=2	S=0	E=3	W=!ENTER	v=0	a=-3759.32	l=0.00
J=3	S=0	E=4	W=!ENTER	v=0	a=-3829.60	l=0.00

How to:

Start	End	cumulatexd	probability	
0	1	-1432,27-0	=	-1432,27
0	2	-1500,93-0	=	-1500,93
1	5	(-2434,05-87,29)-1432,27	=	-3953,61
2	5	(-2431,55-87,29) -1500,93	=	

Result

4	!ENTER
18	IT

Total score: -20218.25

Today:

Large Vocabulary Speech Recognition

- Refresh
- Gaussian densities
- What about LVSR?
- Integrating the Language Model
- Search techniques
- Lattice techniques