

$\hat{W}$  with highest forward probability is then chosen as the recognized word. When subword models are used, word HMMs can be easily constructed by concatenating corresponding phoneme HMMs or other types of subword HMMs according to the procedure described in Chapter 9.

#### 12.2.4. Continuous Speech Recognition

Search in continuous speech recognition is rather complicated, even for a small vocabulary, since the search algorithm has to consider the possibility of each word starting at any arbitrary time frame. Some of the earliest speech recognition systems took a two-stage approach towards continuous speech recognition, first hypothesizing the possible word boundaries and then using pattern matching techniques for recognizing the segmented patterns. However, due to significant cross-word co-articulation, there is no reliable segmentation algorithm for detecting word boundaries other than doing recognition itself.

Let's illustrate how you can extend the isolated-word search technique to continuous speech recognition by a simple example, as shown in Figure 12.10. This system contains only two words,  $w_1$  and  $w_2$ . We assume the language model used here is a uniform unigram ( $P(w_1) = P(w_2) = 1/2$ ).

It is important to represent the language structures in the same HMM framework. In Figure 12.10, we add one starting state  $S$  and one collector state  $C$ . The starting state has a null transition to the initial state of each word HMM with corresponding language model probability ( $1/2$  in this case). The final state of each word HMM has a null transition to the collector state. The collector state then has a null transition back to the starting state in order to allow recursion. Similar to the case of embedding the phoneme (subword) HMMs into the word HMM for isolated speech recognition, we can embed the word HMMs for  $w_1$  and  $w_2$  into a new HMM corresponding to structure in Figure 12.10. Thus, the continuous speech search problem can be solved by the standard HMM formulations.

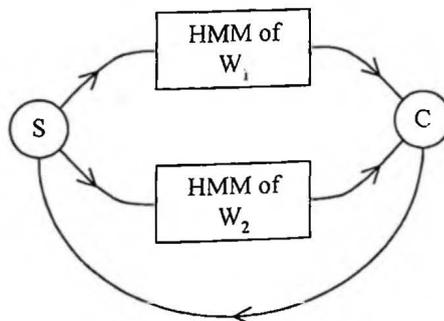


Figure 12.10 A simple example of continuous speech recognition task with two words  $w_1$  and  $w_2$ . A uniform unigram language model is assumed for these words. State  $S$  is the starting state while state  $C$  is a collector state to save fully expanded links between every word pair.

The composite HMMs shown in Figure 12.10 can be viewed as a stochastic finite state network with transition probabilities and output distributions. The search algorithm is essentially producing a match between the acoustic observation  $\mathbf{X}$  and a path<sup>4</sup> in the stochastic finite state network. Unlike isolated word recognition, continuous speech recognition needs to find the optimal word sequence  $\hat{\mathbf{W}}$ . The Viterbi algorithm is clearly a natural choice for this task since the optimal state sequence  $\hat{\mathbf{S}}$  corresponds to the optimal word sequence  $\hat{\mathbf{W}}$ . Figure 12.11 shows the HMM Viterbi trellis computation for the two-word continuous speech recognition example in Figure 12.10. There is a cell for each state in the stochastic finite state network and each time frame  $t$  in the trellis. Each cell  $C_{s,t}$  in the trellis can be connected to a cell corresponding to time  $t$  or  $t+1$  and to states in the stochastic finite state network that can be reached from  $s$ . To make a word transition, there is a null transition to connect the final state of each word HMM to the initial state of the next word HMM that can be followed. The trellis computation is done *time-synchronously* from left to right, i.e., each cell for time  $t$  is completely computed before proceeding to time  $t+1$ .

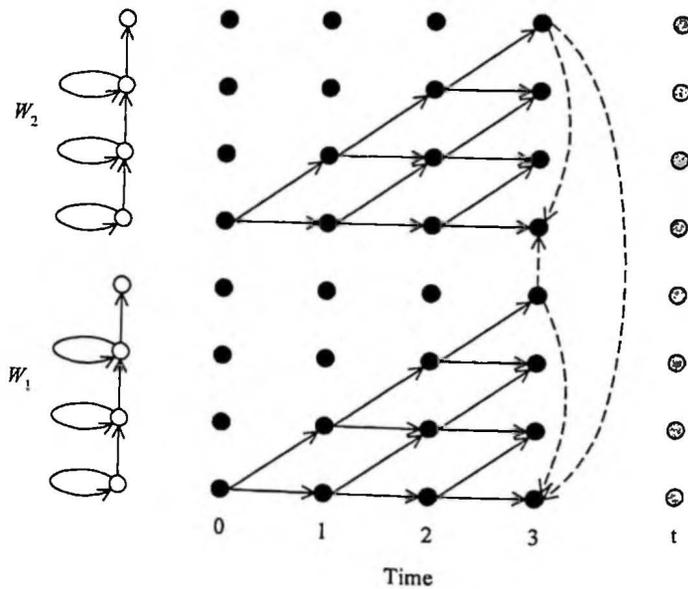


Figure 12.11 HMM trellis for continuous speech recognition example in Figure 12.10. When the final state of the word HMM is reached, a null arc (indicated by a dashed line) is linked from it to the initial state of the following word.

<sup>4</sup> A path here means a sequence of states and transitions.

## 12.3. LANGUAGE MODEL STATES

The state-space is a good indicator of search complexity. Since the HMM representation for each word in the lexicon is fixed, the state-space is determined by the language models. According to Chapter 11, every language model (grammar) is associated with a state machine (automata). Such a state machine is expanded to form the state-space for the recognizer. The states in such a state machine are referred to as *language model states*. For simplicity, we will use the concepts of state-space and language model states interchangeably. The expansion of language model states to HMM states will be done implicitly. The language model states for isolated word recognition are trivial. They are just the union of the HMM states of each word. In this section we look at the language model states for various grammars for continuous speech recognition.

### 12.3.1. Search Space with FSM and CFG

As described in Chapter 8, the complexity for the Viterbi algorithm is  $O(N^2T)$ , where  $N$  is the total number of states in the composite HMM and  $T$  is the length of input observation. A full time-synchronous Viterbi search is quite efficient for moderate tasks (vocabulary  $\leq 500$ ). We have already demonstrated in Figure 12.11 how to search for a two-word continuous speech recognition task with a uniform language model. The uniform language model, which allows all words in the vocabulary to follow every word with the same probability, is suitable for connected-digit task. In fact, most small vocabulary tasks in speech recognition applications usually use a finite state grammar (FSG).

Figure 12.12 shows a simple example of an FSM. Similar to the process described in Sections 12.2.3 and 12.2.4, each of the word arcs in an FSG can be expanded as a network of phoneme (subword) HMMs. The word HMMs are connected with null transitions with the grammar state. A large finite state HMM network that encodes all the legal sentences can be constructed based on the expansion procedure. The decoding process is achieved by performing a time-synchronous Viterbi search on this composite finite state HMM.

In practice, FSGs are sufficient for simple tasks. However, when an FSG is made to satisfy the constraints of sharing of different sub-grammars for compactness and support for dynamic modifications, the resulting non-deterministic FSG is very similar to context-free grammar (CFG) in terms of implementation. The CFG grammar consists of a set of productions or rules, which expand nonterminals into a sequence of terminals and nonterminals. Nonterminals in the grammar tend to refer to high-level task-specific concepts such as dates, names, and commands. The terminals are words in the vocabulary. A grammar also has a non-terminal designated as its start state.

Although efficient parsing algorithms, like chart parsing (described in Chapter 11), are available for CFG, they are not suitable for speech recognition, which requires left-to-right processing. A context-free grammar can be formulated with a recursive transition network (RTN). RTNs are more powerful and complicated than the finite state machines described in

Chapter 11 because they allow arc labels to refer to other networks as well as words. We use Figure 12.13 to illustrate how to embed HMMs into a recursive transition network.

Figure 12.13 is an RTN representation of the following CFG:

```

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

```

There are three types of arcs in an RTN, as shown in Figure 12.13:  $CAT(x)$ ,  $PUSH(x)$ , and  $POP(x)$ . The  $CAT(x)$  arc indicates that  $x$  is a terminal node (which is equivalent to a word arc). Therefore, all the  $CAT(x)$  arcs can be expanded by the HMM network for  $x$ . The word HMM can again be a composite HMM built from phoneme (or subword) HMMs. Similar to the finite state grammar case in Figure 12.12, each grammar state acts as a state with incoming and outgoing null transitions to connect word HMMs in the CFG.

During decoding, the search pursues several paths through the CFG at the same time. Associated with each of the paths is a grammar state that describes completely how the path can be extended further. When the decoder hypothesizes the end of the current word of a path, it asks the CFG module to extend the path further by one word. There may be several alternative successor words for the given path. The decoder considers all the successor word possibilities. This may cause the path to be extended to generate several more paths to be considered, each with its own grammar state.

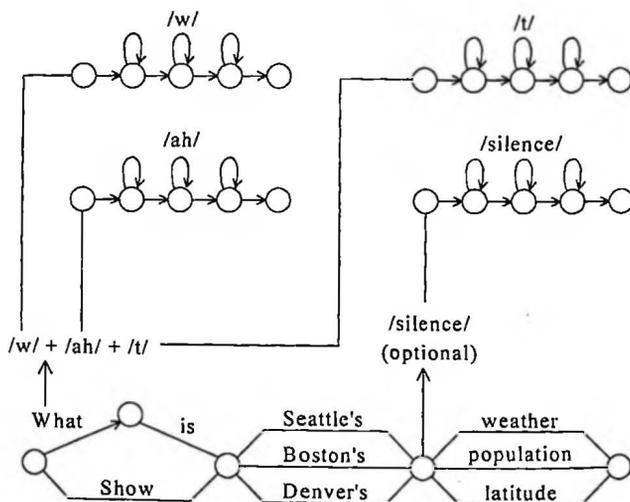


Figure 12.12 An illustration of how to compile a speech recognition task with finite state grammar into a composite HMM.

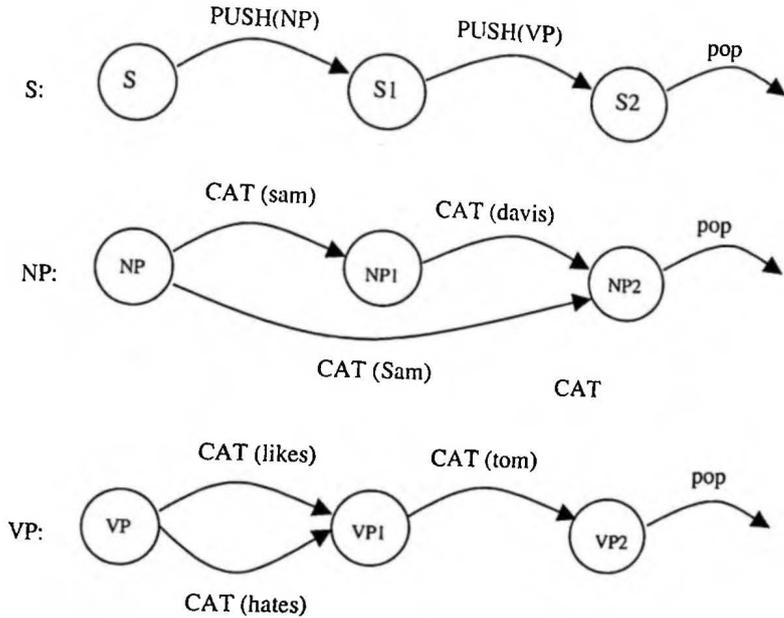


Figure 12.13 A simple RTN example with three types of arcs: CAT(x), PUSH(x), POP.

Readers should note that the same word might be under consideration by the decoder in the context of different paths and grammar states at the same time. For example, there are two word arcs CAT (Sam) in Figure 12.13. Their HMM states should be considered as distinct states in the trellis because they are in completely different grammar states. Two different states in the trellis also means that different paths going into these two states cannot be merged. Since these two partial paths will lead to different successive paths, the search decision needs to be postponed until the end of search. Therefore, when embedding HMMs into word arcs in the grammar network, the HMM state will be assigned a new state identity, although the HMM parameters (transition probabilities and output distributions) can still be shared across different grammar arcs.

Each path consists of a stack of production rules. Each element of the stack also contains the position within the production rule of the symbol that is currently being explored. The search graph (trellis) started from the initial state of CFG (state S). When the path needs to be extended, we look at the next arc (symbol in CFG) in the production. When the search enters a CAT(x) arc (terminal), the path gets extended with the terminal, and the HMM trellis computation is performed on the CAT(x) arc to match the model x against the acoustic data. When the final state of the HMM for x is reached, the search moves on via the null

transition to the destination of the  $CAT(x)$  arc. When the search enters a  $PUSH(x)$  arc, it indicates a nonterminal symbol  $x$  is encountered. In effect, the search is about to enter a sub-network of  $x$ ; the destination of the  $PUSH(x)$  arc is stored in a last-in first-out (LIFO) stack. When the search reaches a  $POP$  arc that signals the end of the current network, the control should jump back to the calling network. In other words, the search returns to the state extracted from the top of the LIFO stack. Finally, when we reach the end of the production rule at the very bottom of the stack, we have reached an accepting state in which we have seen a complete grammatical sentence. For our decoding purpose, that is the state we want to pick as the best score at the end of time frame  $T$  to get the search result.

The problem of connected word recognition by finite state or context-free grammars is that the number of states increases enormously when it is applied to more complex grammars. Moreover it remains a challenge to generate such FSGs or CFGs from a large corpus, either manually or automatically. As mentioned in Chapter 11, it is questionable whether FSG or CFG is adequate to describe natural languages or unconstrained spontaneous languages. Instead,  $n$ -gram language models are often used for natural languages or unconstrained spontaneous languages. In the next section we investigate how to integrate various  $n$ -grams into continuous speech recognition.

### 12.3.2. Search Space with the Unigram

The simplest  $n$ -gram is the unigram that is memory-less and depends only on the current word.

$$P(\mathbf{W}) = \prod_{i=1}^n P(w_i) \quad (12.11)$$

Figure 12.14 shows such a unigram grammar network. The final state of each word HMM is connected to the collector state by a null transition, with probability 1.0. The collector state is then connected to the starting state by another null transition, with transition probability equal to 1.0. For word expansion, the starting state is connected to the initial state of each word HMM by a null transition, with transition probability equal to the corresponding unigram probability. Using the collector state and starting state for word expansion allows efficient expansion because it first merges all the word-ending paths<sup>5</sup> (only the best one survives) before expansion. It can cut the total cross-word expansion from  $N^2$  to  $N$ .

<sup>5</sup> In graph search, a partial path still under consideration is also referred to as a theory, although we will use paths instead of theories in this book.

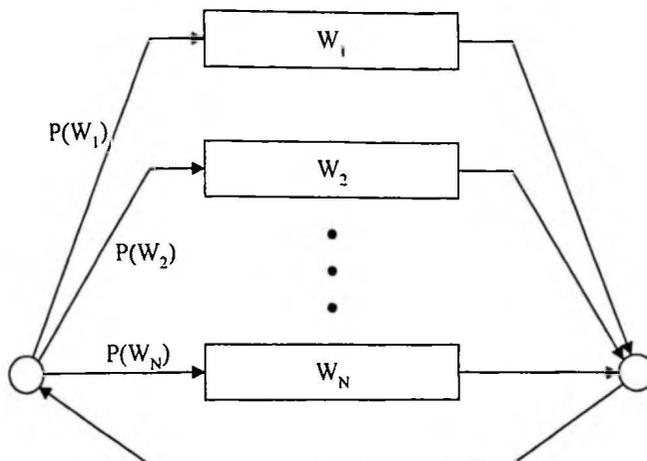


Figure 12.14 A unigram grammar network where the unigram probability is attached as the transition probability from starting state  $S$  to the first state of each word HMM.

### 12.3.3. Search Space with Bigrams

When the bigram is used, the probability of a word depends only on the immediately preceding word. Thus, the language model score is:

$$P(\mathbf{W}) = P(w_1 | \langle s \rangle) \prod_{i=2}^n P(w_i | w_{i-1}) \quad (12.12)$$

where  $\langle s \rangle$  represents the symbol of starting of a sentence.

Figure 12.15 shows a grammar network using a bigram language model. Because of the bigram constraint, the merge-and-expand framework for unigram search no longer applies here. Instead, the bigram search needs to perform expand-and-merge. Thus, bigram expansion is more expensive than unigram expansion. For a vocabulary size  $N$ , the bigram would need  $N^2$  word-to-word transitions in comparison to  $N$  for the unigram. Each word transition has a transition probability equal to the corresponding bigram probability. Fortunately, the total number of states for bigram search is still proportional to the vocabulary size  $N$ .

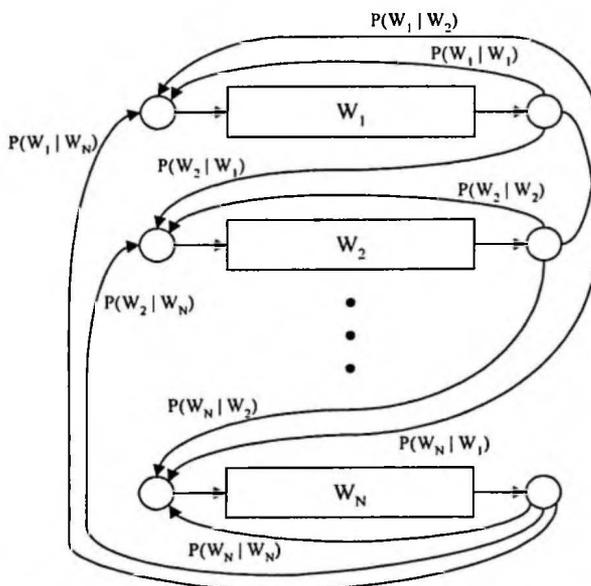


Figure 12.15 A bigram grammar network where the bigram probability  $P(w_j | w_i)$  is attached as the transition probability from word  $w_i$  to  $w_j$  [19].

Because the search space for bigram is kept manageable, bigram search can be implemented very efficiently. Bigram search is a good compromise between efficient search and effective language models. Therefore, bigram search is arguably the most widely used search technique for unconstrained large-vocabulary continuous speech recognition. Particularly for the multiple-pass search techniques described in Chapter 13, a bigram search is often used in the first pass search.

### 12.3.3.1. Backoff Paths

When the vocabulary size  $N$  is large, the total bigram expansion  $N^2$  can become computationally prohibitive. As described in Chapter 11, only a limited number of bigrams are observable in any practical corpora for a large vocabulary size. Suppose the probabilities for unseen bigrams are obtained through Katz's backoff mechanism. That is, for unseen bigram  $P(w_j | w_i)$ ,

$$P(w_j | w_i) = \alpha(w_i)P(w_j) \quad (12.13)$$

where  $\alpha(w_i)$  is the backoff weight for word  $w_i$ .

Using the backoff mechanism for unseen bigrams, the bigram expansion can be significantly reduced [12]. Figure 12.16 shows the new word expansion scheme. Instead of full bigram expansion, only observed bigrams are connected by direct word transitions with correspondent bigram probabilities. For backoff 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 transition probability equal to its corresponding unigram probability  $P(w_j)$ . Readers should note that there are now two paths from  $w_i$  to  $w_j$  for an observed bigram  $P(w_j | w_i)$ . One is the direct link representing the observable bigram  $P(w_j | w_i)$ , and the other is the two-link backoff path representing  $\alpha(w_i)P(w_j)$ . For a word pair whose bigram exists, the two-link backoff path is likely to be ignored since the backoff unigram probability is almost always smaller than the observed bigram  $P(w_j | w_i)$ . Suppose there are only  $N_b$  different observable bigrams, this scheme requires  $N_b + 2N$  instead of  $N^2$  word transitions. Since under normal circumstance  $N_b \ll N^2$ , this backoff scheme significantly reduces the cost of word expansion.

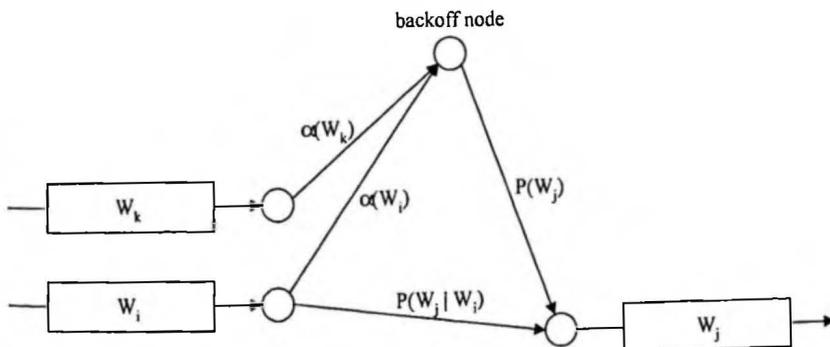


Figure 12.16 Reducing bigram expansion in a search by using the backoff node. In addition to normal bigram expansion arcs for all observed bigrams, the last state of word  $w_i$  is first connected to a central backoff node with transition probability equal to backoff weight  $\alpha(w_i)$ . The backoff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

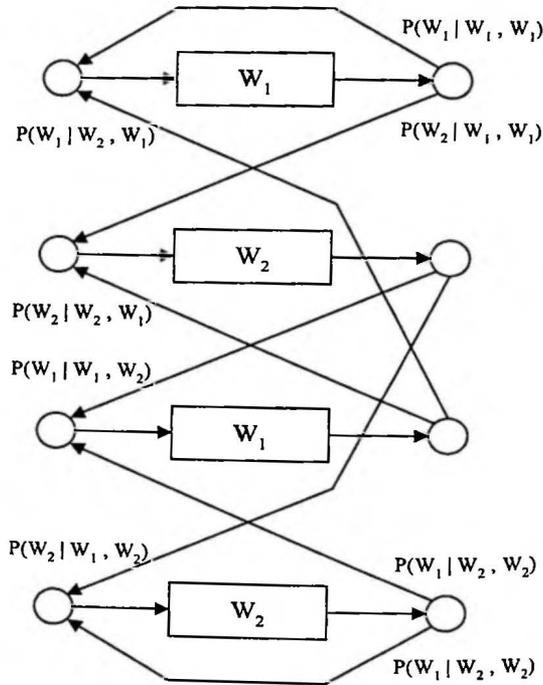
### 12.3.4. Search Space with Trigrams

For a trigram language model, the language model score is:

$$P(W) = P(w_1 | \langle s \rangle) P(w_2 | \langle s \rangle, w_1) \prod_{i=3}^n P(w_i | w_{i-2}, w_{i-1}) \tag{12.14}$$

The search space is considerably more complex, as shown in Figure 12.17. Since the equivalence grammar class is the previous two words  $w_i$  and  $w_j$ , the total number of grammar states is  $N^2$ . From each of these grammar states, there is a transition to the next word [19].

Obviously, it is very expensive to implement large-vocabulary trigram search given the complexity of the search space. It becomes necessary to dynamically generate the trigram search graph (trellis) via a graph search algorithm. The other alternative is to perform a multiple-pass search strategy, in which the first-pass search uses less detailed language models, like bigrams, to generate an  $n$ -best list or word lattice, and then a second-pass detailed search can use trigrams on a much smaller search space. Multiple-pass search strategy is discussed in Chapter 13.



**Figure 12.17** A trigram grammar network where the trigram probability  $P(w_k | w_i, w_j)$  is attached to transition from grammar state word  $w_i, w_j$  to the next word  $w_k$ . Illustrated here is a two-word vocabulary, so there are four grammar states in the trigram network [19].

### 12.3.5. How to Handle Silences Between Words

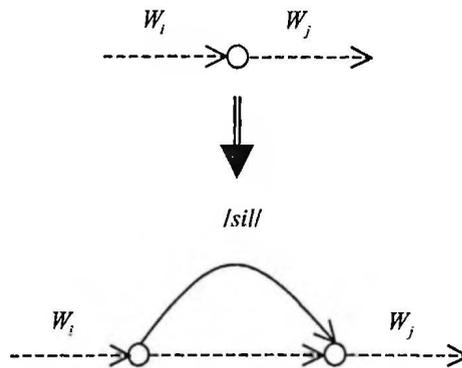
In continuous speech recognition, there are unavoidable pauses (silences) between words or sentences. The pause is often referred to as silence or a non-speech event in continuous speech recognition. Acoustically, the pause is modeled by a silence model<sup>6</sup> that models background acoustic phenomena. The silence model is usually modeled with a topology flexible enough to accommodate a wide range of lengths, since the duration of a pause is arbitrary.

It can be argued that silences are actually linguistically distinguishable events, which contribute to prosodic and meaning representation. For example, people are likely to pause more often in phrasal boundaries. However, these patterns are so far not well understood for unconstrained natural speech (particularly for spontaneous speech). Therefore, the design of almost all automatic speech recognition systems today allows silences occurring just about anywhere between two lexical tokens or between sentences. It is relatively safe to assume that people pause a little bit between sentences to catch breath, so the silences between sentences are assumed mandatory while silences between words are optional. In most systems, silence is often modeled as a special lexicon entry with special language model probability. This special language model probability is also referred to as silence insertion penalty that is set to adjust the likelihood of inserting such an optional silence between words.

It is relatively straightforward to handle the optional silence between words. We need only to replace all the grammar states connecting words with a small network like the one shown in Figure 12.18. This arrangement is similar to that of the optional silence in training continuous speech, described in Chapter 9. The small network contains two parallel paths. One is the original null transition acting as the direct transition from one word to another, while the other path will need to go through a silence model with the silence insertion penalty attached in the transition probability before going to the next word.

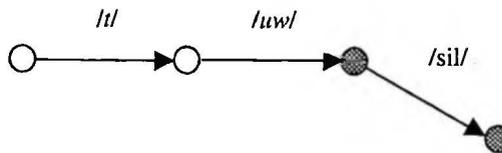
One thing to clarify in the implementation of Figure 12.18 is that this silence expansion needs to be done for every grammar state connecting words. In the unigram grammar network of Figure 12.14, since there is only one collector node to connect words, the silence expansion is required only for this collector node. On the other hand, in the bigram grammar network of Figure 12.15, there is a collector node for every word before expanding to the next word. In this case, the silence expansion is required for every collector node. For a vocabulary size  $|V|$ , this means there are  $|V|$  numbers of silence networks in the grammar search network. This requirement lies in the fact that in bigram search we cannot merge paths before expanding into the next word. Optional silence can then be regarded as part of the search effort for the previous word, so the word expansion needs to be done after finishing the optional silence. Therefore, we treat each word as having two possible pronunciations, one with the silence at the end and one without. This viewpoint integrates silence in the word pronunciation network like the example shown in Figure 12.19.

<sup>6</sup> Some researchers extend the context-dependent modeling to silence models. In that case, there are several silence models based on surrounding contexts.



**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.

For efficiency reasons, a single silence is sometimes used for large-vocabulary continuous speech recognition using higher order  $n$ -gram language model. Theoretically, this could be a source of pruning errors.<sup>7</sup> However, the error could turn out to be so small as to be negligible because there are, in general, very few pauses between word for continuous speech. On the other hand, the overhead of using multiple silences should be very minimal because it is less likely to visit those silence models at the end of words due to pruning.



**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.

## 12.4. TIME-SYNCHRONOUS VITERBI BEAM SEARCH

When HMMs are used for acoustic models, the acoustic model score (likelihood) used in search is by definition the forward probability. That is, all possible state sequences must be considered. Thus,

<sup>7</sup> Speech recognition errors due to sub-optimal search or heuristic pruning are referred to as *pruning errors*, which will be described in detail in Chapter 13.

$$P(\mathbf{X} | \mathbf{W}) = \sum_{\text{all possible } s_0^T} P(\mathbf{X}, s_0^T | \mathbf{W}) \quad (12.15)$$

where the summation is to be taken over all possible state sequences  $\mathbf{S}$  with the word sequence  $\mathbf{W}$  under consideration. However, under the trellis framework (as in Figure 12.11), more bookkeeping must be performed since we cannot add scores with different word sequence history. Since the goal of decoding is to uncover the best word sequence, we could approximate the summation with the maximum to find the best state sequence instead. The Bayes' decision rule, Eq. (12.5), becomes

$$\hat{\mathbf{W}} = \arg \max_{\mathbf{W}} P(\mathbf{W})P(\mathbf{X} | \mathbf{W}) \cong \arg \max_{\mathbf{W}} \left\{ P(\mathbf{W}) \max_{s_0^T} P(\mathbf{X}, s_0^T | \mathbf{W}) \right\} \quad (12.16)$$

Equation (12.16) is often referred to as the *Viterbi approximation*. It can be literally translated to "the *most likely word sequence* is approximated by the *most likely state sequence*." Viterbi search is then sub-optimal. Although the search results by using forward probability and Viterbi probability could, in principle, be different, in practice this is rarely the case. We use this approximation for the rest of this chapter.

The Viterbi search has already been discussed as a solution to one of the three fundamental HMM problems in Chapter 8. It can be executed very efficiently via the same trellis framework. To briefly reiterate, the Viterbi search is a time-synchronous search algorithm that completely processes time  $t$  before going on to time  $t+1$ . For time  $t$ , each state is updated by the best score (instead of the sum of all incoming paths) from all states in at time  $t-1$ . This is why it is often called *time-synchronous Viterbi search*. When one update occurs, it also records the backtracking pointer to remember the most probable incoming state. At the end of search, the most probable state sequence can be recovered by tracing back these backtracking pointers. The Viterbi algorithm provides an optimal solution for handling nonlinear time warping between hidden Markov models and acoustic observation, word boundary detection and word identification in continuous speech recognition. This unified Viterbi search algorithm serves as the basis for all search algorithms as described in the rest of the chapter.

It is necessary to clarify the backtracking pointer for time-synchronous Viterbi search for continuous word recognition. We are generally not interested in the optimal state sequence for speech recognition.<sup>8</sup> Instead, we are only interested in the optimal word sequence indicated by Eq. (12.16). Therefore, we use the backtrack pointer just to remember the word history for the current path, so the optimal word sequence can be recovered at the end of search. To be more specific, when we reach the final state of a word, we create a history node containing the word identity and current time index and append this history node to the existing backtrack pointer. This backtrack pointer is then passed onto the successor node if it

<sup>8</sup> While we are not interested in optimal state sequences for ASR, they are very useful in deriving phonetic segmentation, which could provide important information for developing ASR systems.

is the optimal path leading to the successor node for both intra-word and inter-word transition. The side benefit of keeping this backtrack pointer is that we no longer need to keep the entire trellis during the search. Instead, we only need space to keep two successive time slices (columns) in the trellis computation (the previous time slice and the current time slice) because all the backtracking information is now kept in the backtrack pointer. This simplification is a significant benefit in the implementation of a time-synchronous Viterbi search.

Time-synchronous Viterbi search can be considered as a *breadth-first search* with dynamic programming. Instead of performing a tree search algorithm, the dynamic programming principle helps create a search graph where multiple paths leading to the same search state are merged by keeping the best path (with minimum cost). The Viterbi trellis is a representation of the search graph. Therefore, all the efficient techniques for graph search algorithms can be applied to time-synchronous Viterbi search. Although so far we have described the trellis in an explicit fashion—the whole search space needs to be explored before the optimal path can be found—it is not necessary to do so. When the search space contains an enormous number of states, it becomes impractical to pre-compile the composite HMM entirely and store it in the memory. It is preferable to dynamically build and allocate portions of the search space sufficient to search the promising paths. By using the graph search algorithm described in Section 12.1.1, only part of the entire Viterbi trellis is generated explicitly. By constructing the search space dynamically, the computation cost of the search is proportional only to the number of active hypotheses, independent of the overall size of the potential search space. Therefore, dynamically generated trellises are key to heuristic Viterbi search for efficient large-vocabulary continuous speech recognition, as described in Chapter 13.

### 12.4.1. The Use of Beam

Based on Chapter 8, the search space for Viterbi search is  $O(NT)$  and the complexity is  $O(N^2T)$ , where  $N$  is the total number of HMM states and  $T$  is the length of the utterance. For large-vocabulary tasks these numbers are astronomically large even with the help of dynamic programming. In order to avoid examining the overwhelming number of possible cells in the HMM trellis, a heuristic search is clearly needed. Different heuristics generate or explore portions of the trellis in different ways.

A simple way to prune the search space for breadth-first search is the beam search described in Section 12.1.3.2. Instead of retaining all candidates (cells) at every time frame, a threshold  $T$  is used to keep only a subset of promising candidates. The state at time  $t$  with the lowest cost  $D_{\min}$  is first identified. Then each state at time  $t$  with cost  $> D_{\min} + T$  is discarded from further consideration before moving on to the next time frame  $t+1$ . The use of the beam alleviates the need to process all the cells. In practice, it can lead to substantial savings in computation with little or no loss of accuracy.

Although beam search is a simple idea, the combination of time-synchronous Viterbi and beam search algorithms produces the most powerful search strategy for large-vocabulary speech recognition. Comparing paths with equal length under a time-

synchronous search framework makes beam search possible. That is, for two different word sequences  $W_1$  and  $W_2$ , the posterior probabilities  $P(W_1 | x'_0)$  and  $P(W_2 | x'_0)$  are always compared based on the same partial acoustic observation  $x'_0$ . This makes the comparison straightforward because the denominator  $P(x'_0)$  in Eq. (12.5) is the same for both terms and can be ignored. Since the score comparison for each time frame is fair, the only assumption of beam search is that an optimal path should have a good enough partial-path score for each time frame to survive under beam pruning.

The time-synchronous framework is one of the aspects of Viterbi beam search that is critical to its success. Unlike the time-synchronous framework, time-asynchronous search algorithms such as stack decoding require the normalization of likelihood scores over feature streams of different time lengths. This, as we will see in Section 12.5, is the Achilles' heel of that approach.

The straightforward time-synchronous Viterbi beam search is ineffective in dealing with the gigantic search space of high perplexity tasks. However, with a better understanding of the linguistic search space and the advent of techniques for obtaining  $n$ -best lists from time-synchronous Viterbi search, described in Chapter 13, time-synchronous Viterbi beam search has turned out to be surprisingly successful in handling tasks of all sizes and all different types of grammars, including FSG, CFG, and  $n$ -gram [2, 14, 18, 28, 34, 38, 44]. Therefore, it has become the predominant search strategy for continuous speech recognition.

## 12.4.2. Viterbi Beam Search

To explain the time-synchronous Viterbi beam search in a formal way [31], we first define some quantities:

$D(t; s; w) \equiv$  total cost of the best path up to time  $t$  that ends in state  $s$ , of grammar word state  $w$ .

$h(t; s; w) \equiv$  backtrack pointer for the best path up to time  $t$  that ends in state  $s$ , of grammar word state  $w$ .

Readers should be aware that  $w$  in the two quantities above represents a grammar word state in the search space. It is different from just the word identity since the same word could occur in many different language model states, as in the trigram search space shown in Figure 12.17.

There are two types of dynamic programming (DP) transition rules [30], namely intra-word and inter-word transition. The intra-word transition is just like the Viterbi rule for HMMs and can be expressed as follows:

$$D(t; s; w) = \min_{s_{t-1}} \{d(\mathbf{x}_t, s_t | s_{t-1}; w) + D(t-1; s_{t-1}; w)\} \quad (12.17)$$

$$h(t; s; w) = h(t-1, b_{\min}(t; s; w); w) \quad (12.18)$$

where  $d(\mathbf{x}_t, s_t | s_{t-1}; w)$  is the cost associated with taking the transition from state  $s_{t-1}$  to state  $s_t$  while generating output observation  $\mathbf{x}_t$ , and  $b_{\min}(t; s_t; w)$  is the optimal predecessor state of cell  $D(t; s_t; w)$ . To be specific, they can be expressed as follows:

$$d(\mathbf{x}_t, s_t | s_{t-1}; w) = -\log P(s_t | s_{t-1}; w) - \log P(\mathbf{x}_t | s_t; w) \quad (12.19)$$

$$b_{\min}(t; s_t; w) = \arg \min_{s_{t-1}} \{d(\mathbf{x}_t, s_t | s_{t-1}; w) + D(t-1; s_{t-1}; w)\} \quad (12.20)$$

The inter-word transition is basically a null transition without consuming any observation. However, it needs to deal with creating a new history node for the backtracking pointer. Let's define  $F(w)$  as the final state of word HMM  $w$  and  $I(w)$  as the initial state of word HMM  $w$ . Moreover, state  $\eta$  is denoted as the pseudo initial state. The inter-word transition can then be expressed as follows:

$$D(t; \eta; w) = \min_v \{\log P(w | v) + D(t; F(v); v)\} \quad (12.21)$$

$$h(t; \eta; w) = \langle v_{\min}, t \rangle :: h(t, F(v_{\min}); v_{\min}) \quad (12.22)$$

where  $v_{\min} = \arg \min_v \{\log P(w | v) + D(t; F(v); v)\}$  and  $::$  is a link appending operator.

The time-synchronous Viterbi beam search algorithm assumes that all the intra-word transitions are evaluated before inter-word null transitions take place. The same time index is used intentionally for inter-word transition since the null language model state transition does not consume an observation vector. Since the initial state  $I(w)$  for word HMM  $w$  could have a self-transition, the cell  $D(t; I(w); w)$  might already have an active path. Therefore, we need to perform the following check to advance the inter-word transitions.

$$\begin{aligned} \text{if } D(t; \eta; w) < D(t; I(w); w) \\ D(t; I(w); w) = D(t; \eta; w) \text{ and } h(t; I(w); w) = h(t; \eta; w) \end{aligned} \quad (12.23)$$

The time-synchronous Viterbi beam search can be summarized as in Algorithm 12.6. For large-vocabulary speech recognition, the experimental results show that only a small percentage of the entire search space (the beam) needs to be kept for each time interval  $t$  without increasing error rates. Empirically, the beam size has typically been found to be between 5% and 10% of the entire search space. In Chapter 13 we describe strategies of using different level of beams for more effectively pruning.

## 12.5. STACK DECODING (A\* SEARCH)

If some reliable heuristics are available to guide the decoding, the search can be done in a depth-first fashion around the best path early on, instead of wasting efforts on unpromising paths via the time-synchronous beam search. Stack decoding represents the best attempt to

**ALGORITHM 12.6: TIME-SYNCHRONOUS VITERBI BEAM SEARCH**

**Step 1: Initialization:** For all the grammar word states  $w$  which can start a sentence,

$$D(0; I(w); w) = 0$$

$$h(0; I(w); w) = null$$

**Step 2: Induction:** For time  $t = 1$  to  $T$  do

For all active states do

Intra-word transitions according to Eq. (12.17) and (12.18)

$$D(t; s_t; w) = \min_{s_{t-1}} \{d(x_t, s_t | s_{t-1}; w) + D(t-1; s_{t-1}; w)\}$$

$$h(t; s_t; w) = h(t-1, b_{\min}(t; s_t; w); w)$$

For all active word-final states do

Inter-word transitions according to Eq. (12.21), (12.22) and (12.23)

$$D(t; \eta; w) = \min_v \{\log P(w | v) + D(t; F(v); v)\}$$

$$h(t; \eta; w) = \langle v_{\min}, t \rangle :: h(t, F(v_{\min}); v_{\min})$$

$$\text{if } D(t; \eta; w) < D(t; I(w); w)$$

$$D(t; I(w); w) = D(t; \eta; w) \text{ and } h(t; I(w); w) = h(t; \eta; w)$$

Pruning: Find the cost for the best path and decide the beam threshold

Prune unpromising hypotheses

**Step 3: Termination:** Pick the best path among all the possible final states of grammar at time  $T$ . Obtain the optimal word sequence according to the backtracking pointer  $h(t; \eta; w)$

use A\* search instead of time-synchronous beam search for continuous speech recognition. Unfortunately, as we will discover in this section, such a heuristic function  $h(\bullet)$  (defined in Section 12.1.3) is very difficult to attain in continuous speech recognition, so search algorithms based on A\* search are in general less efficient than time-synchronous beam search.

*Stack decoding* is a variant of the heuristic A\* search based on the forward algorithm, where the evaluation function is based on the forward probability. It is a tree search algorithm, which takes a slightly different viewpoint than the time-synchronous Viterbi search. Time-synchronous beam search is basically a breadth-first search, so it is crucial to control the number of all possible language model states as described in Section 12.3. In a typical large-vocabulary Viterbi search with  $n$ -gram language models, this number is determined by the equivalent classes of language model histories. On the other hand, stack decoding as a tree search algorithm treats the search as a task for finding a path in a tree whose branches correspond to words in the vocabulary  $V$ , non-terminal nodes correspond to incomplete sentences, and terminal nodes correspond to complete sentences. The search tree has a constant branching factor of  $|V|$ , if we allow every word to be followed by every word. Figure 12.20 illustrates such a search tree for a vocabulary with three words [19].

An important advantage of stack decoding is its consistency with the forward-backward training algorithm. Viterbi search is a graph search, and paths cannot be easily summed because they may have different word histories. In general, the Viterbi search finds the optimal state sequence instead of optimal word sequence. Therefore, Viterbi approximation is necessary to make the Viterbi search feasible, as described in Section 12.4. Stack decoding is a tree search, so each node has a unique history, and the forward algorithm can be used within word model evaluation. Moreover, all possible beginning and ending times (shaded areas in Figure 12.21) are considered [24]. With stack decoding, it is possible to use an objective function that searches for the optimal word string, rather than the optimal state sequence. Furthermore, it is in principle natural for stack decoding to accommodate long-range language models if the heuristics can guide the search to avoid exploring the overwhelmingly large unpromising grammar states.

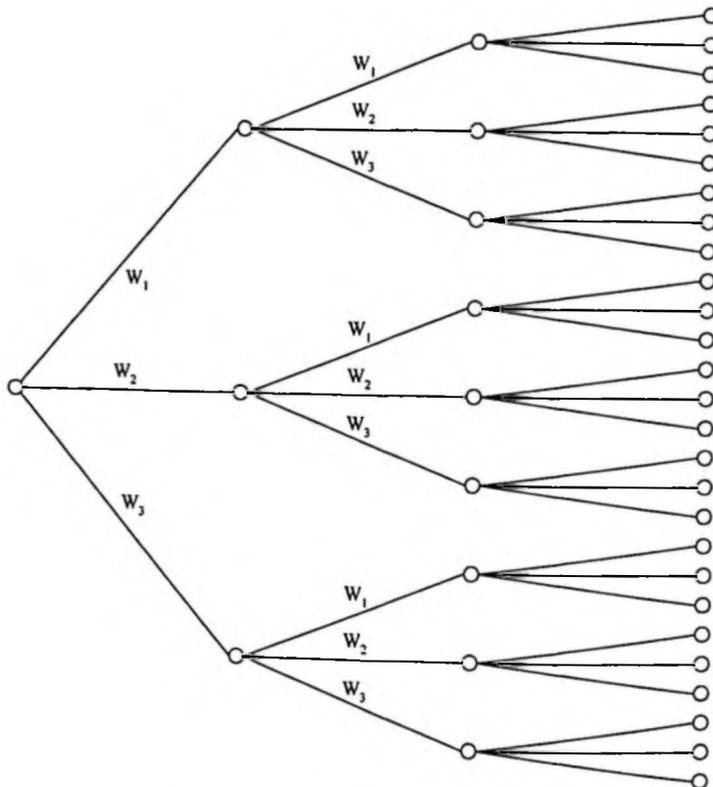


Figure 12.20 A stack decoding search tree for a vocabulary size of three [19].

By formulating stack decoding in a tree search framework, the graph search algorithms described in Section 12.1 can be directly applied to stack decoding. Obviously, blind-search methods, like depth-first and breadth-first search, that do not take advantage of the goodness measurement of how close we are getting to the goal, are usually computationally infeasible in practical speech recognition systems. A\* search is clearly attractive for speech recognition, given the hope of a sufficient heuristic function to guide the tree search in a favorable direction without exploring too many unpromising branches and nodes. In contrast to the Viterbi search, it is not time-synchronous and extends paths of different lengths.

The search begins by adding all possible one-word hypotheses to the *OPEN* list. Then the best path is removed from the *OPEN* list, and all paths from it are extended, evaluated, and placed back in the *OPEN* list. This search continues until a complete path that is guaranteed to be better than all paths in the *OPEN* list has been found.

Unlike Viterbi search, where the acoustic probabilities being compared are always based on the same partial input, it is necessary to compare the goodness of partial paths of different lengths to direct the A\* tree search. Moreover, since stack decoding is done asynchronously, we need an effective mechanism to determine when to end a phone/word evaluation and move on to the next phone/word. Therefore, the heart and soul of the stack decoding are clearly in

1. Finding an effective and efficient heuristic function for estimating the future remaining input feature stream and
2. Determining when to extend the search to the next word/phone.

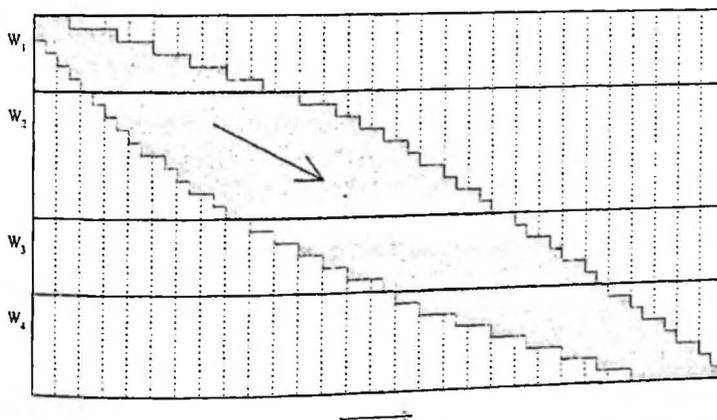


Figure 12.21 The forward trellis space for stack decoding. Each grid point corresponds to a trellis cell in the forward computation. The shaded area represents the values contributing to the computation of the forward score for the optimal word sequence  $w_1, w_2, w_3, \dots$  [24].

In the following section we describe these two critical components. Readers will note that the solutions to these two issues are virtually the same—using a normalization scheme to compare paths of different lengths.

### 12.5.1. Admissible Heuristics for Remaining Path

The key issue in heuristic search is the selection of an evaluation function. As described in Section 12.1.3, the heuristic function of the path  $H_N$  going through node  $N$  includes the cost up to the node and the estimate of the cost to the target node from node  $N$ . Suppose path  $H_N$  is going through node  $N$  at time  $t$ ; then the evaluation for path  $H_N$  can be expressed as follows:

$$f(H_N^t) = g(H_N^t) + h(H_N^{t,T}) \quad (12.24)$$

where  $g(H_N^t)$  is the evaluation function for the partial path of  $H_N$  up to time  $t$ , and  $h(H_N^{t,T})$  is the heuristic function of the remaining path from  $t+1$  to  $T$  for path  $H_N$ . The challenge for stack decoders is to devise an admissible function for  $h(\bullet)$ .

According to Section 12.1.3.1, an admissible heuristic function is one that always underestimates the true cost of the remaining path from  $t+1$  to  $T$  for path  $H_N$ . A trivially admissible function is the zero function. In this case, it results in a very large OPEN list. In addition, since the longer paths tend to have higher cost because of the gradually accumulated cost, the search is likely to be conducted in a breadth-first fashion, which functions very much like a plain Viterbi search. The evaluation function  $g(\bullet)$  can be obtained easily by using the HMM forward score as the true cost up to current time  $t$ . However, how can we find an admissible heuristic function  $h(\bullet)$ ? We present the basic concept here [19, 35].

The goal of  $h(\bullet)$  is to find the expected cost for the remaining path. If we can obtain the expected cost per frame  $\psi$  for the remaining path, the total expected cost,  $(T-t)*\psi$ , is simply the product of  $\psi$  and the length of the remaining path. One way to find such expected cost per frame is to gather statistics empirically from training data.

1. After the final training iteration, perform Viterbi forced alignment<sup>9</sup> with each training utterance to get an optimal time alignment for each word.
2. Randomly select an interval to cover the number of words ranging from two to ten. Denote this interval as  $[i \dots j]$ .
3. Compute the average acoustic cost per frame within this selected interval according to the following formula and save the value in a set  $\Lambda$ :

<sup>9</sup> Viterbi forced alignment means that the Viterbi is performed on the HMM model constructed from the known word transcription. The term “forced” is used because the Viterbi alignment is forced to be performed on the correct model. Viterbi forced alignment is a very useful tool in spoken language processing as it can provide the optimal state-time alignment with the utterances. This detailed alignment can then be used for different purposes, including discriminant training, concatenated speech synthesis, etc.

$$\frac{-1}{j-i} \log P(\mathbf{x}'_i | \mathbf{w}_{i\dots j}) \quad (12.25)$$

where  $\mathbf{w}_{i\dots j}$  is the word string corresponding to interval  $[i\dots j]$ .

4. Repeat Steps 2 and 3 for the entire training set.
5. Define  $\psi_{\min}$  and  $\psi_{\text{avg}}$  as the minimum and average value found in set  $\Lambda$ .

Clearly,  $\psi_{\min}$  should be a good under-estimate of the expected cost per frame for the future unknown path. Therefore, the heuristic function  $h(H_N^{t,T})$  can be derived as:

$$h(H_N^{t,T}) = (T-t)\psi_{\min} \quad (12.26)$$

Although  $\psi_{\min}$  is obtained empirically, stack decoding based on Eq. (12.26) will generally find the optimal solution. However, the search using  $\psi_{\min}$  usually runs very slowly, since Eq. (12.26) always under-estimates the true cost for any portion of speech. In practice, a heuristic function like  $\psi_{\text{avg}}$  that may over-estimate has to be used to prune more hypotheses. This speeds up the search at the expense of possible search errors, because  $\psi_{\text{avg}}$  should represent the average cost per frame for any portion of speech. In fact, there is an argument that one might be able to use a heuristic function even more than  $\psi_{\text{avg}}$ . The argument is that  $\psi_{\text{avg}}$  is derived from the correct path (training data) and the average cost per frame for all paths during search should be more than  $\psi_{\text{avg}}$  because the paths undoubtedly include correct and incorrect ones.

### 12.5.2. When to Extend New Words

Since stack decoding is executed asynchronously, it becomes necessary to detect when a phone/word ends, so that the search can extend to the next phone/word. If we have a cost measure that indicates how well an input feature vector of any length matches the evaluated model sequence, this cost measure should drop slowly for the correct phone/word and rise sharply for an incorrect phone/word. In order to do so, it implies we must be able to compare hypotheses of different lengths.

The first thing that comes to mind for this cost measure is simply the forward cost  $-\log P(\mathbf{x}'_1, s, | w_1^k)$ , which represents the likelihood of producing acoustic observation  $\mathbf{x}'_1$  based on word sequence  $w_1^k$  and ending at state  $s$ . However, it is definitely not suitable because it is deemed to be smaller for a shorter acoustic input vector. This causes the search to almost always prefer short phones/words, resulting in many insertion errors. Therefore, we must derive some normalized score that satisfies the desired property described above. The normalized cost  $\hat{C}(\mathbf{x}'_1, s, | w_1^k)$  can be represented as follows [6, 24]:

$$\hat{C}(\mathbf{x}'_1, s, | w_1^k) = -\log \left[ \frac{P(\mathbf{x}'_1, s, | w_1^k)}{\gamma^t} \right] = -\log [P(\mathbf{x}'_1, s, | w_1^k)] + t \log \gamma \quad (12.27)$$

where  $\gamma$  ( $0 < \gamma < 1$ ) is a constant normalization factor.

Suppose the search is now evaluating a particular word  $w_k$ ; we can define  $\hat{C}_{\min}(t)$  as the minimum cost for  $\hat{C}(\mathbf{x}'_1, s_t | w_k^t)$  for all the states of  $w_k$ , and  $\alpha_{\max}(t)$  as the maximum forward probability for  $P(\mathbf{x}'_1, s_t | w_k^t)$  for all the states of  $w_k$ . That is,

$$\hat{C}_{\min}(t) = \min_{s_t \in \mathcal{H}_t} [\hat{C}(\mathbf{x}'_1, s_t | w_k^t)] \quad (12.28)$$

$$\alpha_{\max}(t) = \max_{s_t \in \mathcal{H}_t} [P(\mathbf{x}'_1 | w_k^t, s_t)] \quad (12.29)$$

We want  $\hat{C}_{\min}(t)$  to be near 0 just as long as the phone/word we are evaluating is the correct one and we have not gone beyond its end. On the other hand, if the phone/word we are evaluating is the incorrect one or we have already passed its end, we want the  $\hat{C}_{\min}(t)$  to be rising sharply. Similar to the procedure of finding the admissible heuristic function, we can set the normalized factor  $\gamma$  empirically during training so that  $\hat{C}_{\min}(T) = 0$  when we know the correct word sequence  $\mathbf{W}$  that produces acoustic observation sequence  $\mathbf{x}_1^T$ . Based on Eq. (12.27),  $\gamma$  should be set to:

$$\gamma = \sqrt[T]{\alpha_{\max}(T)} \quad (12.30)$$

Figure 12.22 shows a plot of  $\hat{C}_{\min}(t)$  as a function of time for correct match. In addition, the cost for the final state  $FS(w_k)$  of word  $w_k$ ,  $\hat{C}(\mathbf{x}'_1, s_t = FS(w_k) | w_k^t)$ , which is the score for  $w_k$ -ending path, is also plotted. There should be a valley centered around 0 for  $\hat{C}(\mathbf{x}'_1, s_t = FS(w_k) | w_k^t)$ , which indicates the region of possible ending time for the correct phone/word. Sometimes a stretch of acoustic observations may match better than the average cost, pushing the curve below 0. Similarly, a stretch of acoustic observations may match worse than the average cost, pushing the curve above 0.

There is an interesting connection between the normalized factor  $\gamma$  and the heuristic estimate of the expected cost per frame,  $\psi$ , defined in Eq. (12.25). Since the cost is simply the logarithm on the inverse posterior probability, we get the following equation:

$$\psi = \frac{-1}{T} \log P(\mathbf{x}_1^T | \hat{\mathbf{W}}) = -\log [\alpha_{\max}(T)^{1/T}] = -\log \gamma \quad (12.31)$$

Equation (12.31) reveals that these two quantities are basically the same estimate. In fact, if we subtract the heuristic function  $f(H'_N)$  defined in Eq. (12.24) by the constant  $\log(\gamma^T)$ , we get exactly the same quantity as the one defined in Eq. (12.27). Decisions on which path to extend first based on the heuristic function and when to extend the search to the next word/phone are basically centered on comparing partial theories with different lengths. Therefore, the normalized cost  $\hat{C}(\mathbf{x}'_1, s_t | w_k^t)$  can be used for both purposes.

Based on the connection we have established, the heuristic function,  $f(H'_N)$ , which estimates the goodness of a path is simply replaced by the normalized evaluation function  $\hat{C}(\mathbf{x}'_1, s_t | w_k^t)$ . If we plot the un-normalized cost  $C(\mathbf{x}'_1, s_t | w_k^t)$  for the optimal path and other

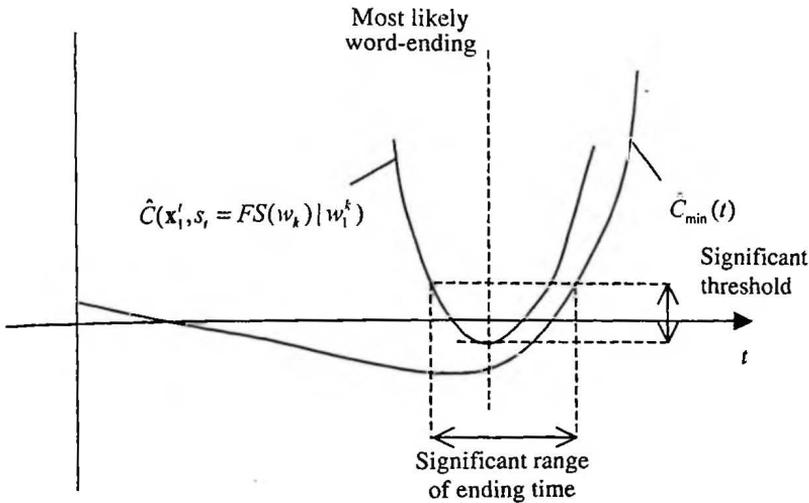


Figure 12.22  $\hat{C}_{\min}(t)$  and  $\hat{C}(x'_1, s_t = FS(w_k) | w_1^k)$  as functions of time  $t$ . The valley region represents possible ending times for the correct phone/word.

competing paths as the function time  $t$ , the cost values increase as paths get longer (illustrated in Figure 12.23) because every frame adds some non-negative cost to the overall cost. It is clear that using un-normalized cost function  $C(x'_1, s_t | w_1^k)$  generally results in a breadth-first search. What we want is an evaluation that decreases slightly along the optimal path, and hopefully increases along other competing paths. Clearly, the normalized cost function  $\hat{C}(x'_1, s_t | w_1^k)$  fulfills this role, as shown in Figure 12.24.

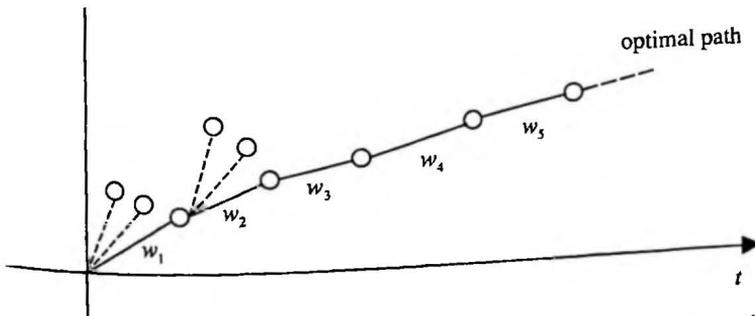


Figure 12.23 Unnormalized cost  $C(x'_1, s_t | w_1^k)$  for optimal path and other competing paths as a function of time.

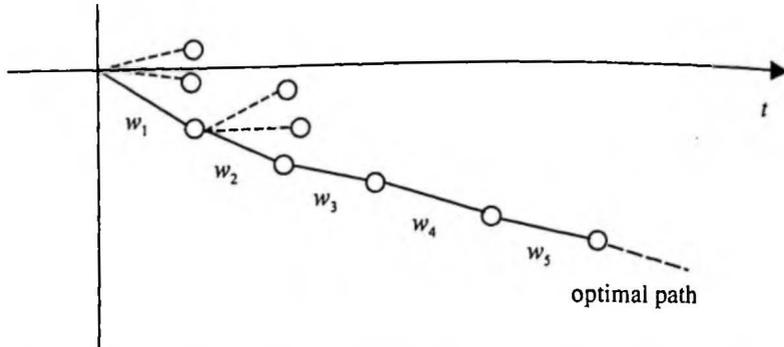


Figure 12.24 Normalized cost  $\hat{C}(x'_i, s_i | w_i^t)$  for the optimal path and other competing paths as a function of time.

Equation (12.30) is a context-less estimation of the normalized factor, which is also referred to as zero-order estimate. To improve the accuracy of the estimate, you can use context-dependent higher-order estimates like [24]:

$\gamma_i = \gamma(x_i)$	first-order estimate
$\gamma_i = \gamma(x_i, x_{i-1})$	second-order estimate
$\gamma_i = \gamma(x_i, x_{i-1}, \dots, x_{i-N+1})$	$n$ -order estimate

Since the normalized factor  $\gamma$  is estimated from the training data that is also used to train the parameters of the HMMs, the normalized factor  $\gamma_i$  tends to be an overestimate. As a result,  $\alpha_{\max}(t)$  might rise slowly for test data even when the correct phone/word model is evaluated. This problem is alleviated by introducing some other scaling factor  $\delta < 1$  so that  $\alpha_{\max}(t)$  falls slowly for test data for when evaluating the correct phone/word model. The best solution for this problem is to use an independent data set other than the training data to derive the normalized factor  $\gamma_i$ .

### 12.5.3. Fast Match

Even with an efficient heuristic function and mechanism to determine the ending time for a phone/word, stack decoding could still be too slow for large-vocabulary speech recognition tasks. As described in Section 12.5.1, an effective underestimated heuristic function for the remaining portion of speech is very difficult to derive. On the other hand, a heuristic estimate for the immediate short segment that usually corresponds to the next phone or word may be feasible to attain. In this section, we describe the fast-match mechanism that reduces phone/word candidates for detailed match (expansion).

In asynchronous stack decoding, the most expensive step is to extend the best subpath. For a large-vocabulary search, it implies the calculation of  $P(x_i^{t+k} | w)$  over the entire vocabulary size  $|V|$ . It is desirable to have a fast computation to quickly reduce the possible

words starting at a given time  $t$  to reduce the search space. This process is often referred to as *fast match* [15, 35]. In fact, fast match is crucial to stack decoding, of which it becomes an integral part. Fast match is a method for the rapid computation of a list of candidates that constrain successive search phases. The expensive *detailed match* can then be performed after fast match. In this sense, fast match can be regarded as an additional pruning threshold to meet before new word/phone can be started.

Fast match, by definition, needs to use only a small amount of computation. However, it should also be accurate enough not to prune away any word/phone candidates that participate in the best path eventually. Fast match is, in general, characterized by the approximations that are made in the acoustic/language models in order to reduce computation. There is an obvious trade-off between these two objectives. Fortunately, many systems [15] have demonstrated that one needs to sacrifice very little accuracy in order to speed up the computation considerably.

Similar to *admissibility* in A\* search, there is also an *admissibility* property in fast match. A fast match method is called admissible if it never prunes away the word/phone candidates that participate in the optimal path. In other words, a fast match is admissible if the recognition errors that appear in a system using the fast match followed by a detailed match are those that would appear if the detailed match were carried out for all words/phones in the vocabulary. Since fast match can be applied to either word or phone level, as we describe in the next section, we explain the admissibility for the case of word-level fast match for simplicity. The same principle can be easily extended to phone-level fast match.

Let  $V$  be the vocabulary and  $C(\mathbf{X} | w)$  be the cost of a detailed match between input  $\mathbf{X}$  and word  $w$ . Now  $F(\mathbf{X} | w)$  is an estimator of  $C(\mathbf{X} | w)$  that is accurate enough and fast to compute. A word list selected by fast match estimator can be attained by first computing  $F(\mathbf{X} | w)$  for each word  $w$  of the vocabulary. Suppose  $w_b$  is the word for which the fast match has a minimum cost value:

$$w_b = \arg \min_{w \in V} F(\mathbf{X} | w) \quad (12.32)$$

After computing  $C(\mathbf{X} | w_b)$ , the detailed match cost for  $w_b$ , we form the fast match word list,  $\Lambda$ , from the word  $w$  in the vocabulary such that  $F(\mathbf{X} | w)$  is no greater than  $C(\mathbf{X} | w_b)$ . In other words,

$$\Lambda = \{w \in V | F(\mathbf{X} | w) \leq C(\mathbf{X} | w_b)\} \quad (12.33)$$

Similar to the admissibility condition for A\* search [3, 33], the fast match estimator  $F(\bullet)$  conducted in the way described above is admissible if and only if  $F(\mathbf{X} | w)$  is always an under-estimator (lower bound) of detailed match  $C(\mathbf{X} | w)$ . That is,

$$F(\mathbf{X} | w) \leq C(\mathbf{X} | w) \quad \forall \mathbf{X}, w \quad (12.34)$$

The proof is straightforward. If the word  $w_c$  has a lower detailed match cost  $C(\mathbf{X} | w_c)$ , you can prove that it must be included in the fast match list  $\Lambda$  because

$$C(\mathbf{X} | w_c) \leq C(\mathbf{X} | w_b) \text{ and } F(\mathbf{X} | w_c) \leq C(\mathbf{X} | w_c) \Rightarrow F(\mathbf{X} | w_c) \leq C(\mathbf{X} | w_b)$$

Therefore, based on the definition of  $\Lambda$ ,  $w_c \in \Lambda$ .

Now the task is to find an admissible fast match estimator. Bahl et al. [6] proposed one fast match approximation for discrete HMMs. As we will see later, this fast match approximation is indeed equivalent to a simplification of the HMM structure. Given the HMM for word  $w$  and an input sequence  $x_1^T$  of codebook symbols describing the input signal, the probability that the HMM  $w$  produces the VQ sequence  $x_1^T$  is given by (according to Chapter 8):

$$P(x_1^T | w) = \sum_{s_1, s_2, \dots, s_T} \left[ P_w(s_1, s_2, \dots, s_T) \prod_{i=1}^T P_w(x_i | s_i) \right] \quad (12.35)$$

Since we often use Viterbi approximation instead of the forward probability, the equation above can be approximated by:

$$P(x_1^T | w) \cong \max_{s_1, s_2, \dots, s_T} \left[ P_w(s_1, s_2, \dots, s_T) \prod_{i=1}^T P_w(x_i | s_i) \right] \quad (12.36)$$

The detailed match cost  $C(\mathbf{X} | w)$  can now be represented as:

$$C(\mathbf{X} | w) = \min_{s_1, s_2, \dots, s_T} \left\{ -\log \left[ P_w(s_1, s_2, \dots, s_T) \prod_{i=1}^T P_w(x_i | s_i) \right] \right\} \quad (12.37)$$

Since the codebook size is finite, it is possible to compute, for each model  $w$ , the highest output probability for every VQ label  $c$  among all states  $s_k$  in HMM  $w$ . Let's define  $m_w(c)$  to be the following:

$$m_w(c) = \max_{s_k \in w} P_w(c | s_k) = \max_{s_k \in w} b_k(c) \quad (12.38)$$

We can further define the  $q_{\max}(w)$  as the maximum state sequence with respect to  $T$ , i.e., the maximum probability of any complete path in HMM  $w$ .

$$q_{\max}(w) = \max_T [P_w(s_1, s_2, \dots, s_T)] \quad (12.39)$$

Now let's define the fast match estimator  $F(\mathbf{A} | w)$  as the following:

$$F(\mathbf{X} | w) = -\log \left[ q_{\max}(w) \prod_{i=1}^T m_w(x_i) \right] \quad (12.40)$$

It is easy to show the fast match estimator  $F(\mathbf{X} | w) \leq C(\mathbf{X} | w)$  is admissible based on Eq. (12.38) to Eq. (12.40).

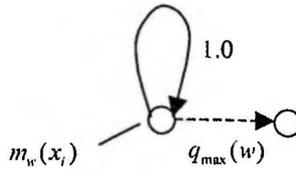


Figure 12.25 The equivalent one-state HMM corresponding to fast match computation defined in Eq. (12.40) [15].

The fast match estimator defined in Eq. (12.40) requires  $T+1$  additions for a vector sequence of length  $T$ . The operation can be viewed as equivalent to the forward computation with a one-state HMM of the form shown in Figure 12.25. This correspondence can be interpreted as a simplification of the original multiple-state HMM into such a one-state HMM. It thus explains why fast match can be computed much faster than detailed match. Readers should note that this HMM is not actually a true HMM by strict definition, because the output probability distribution  $m_w(c)$  and the transition probability distribution do not add up to one.

The fast match computation defined in Eq. (12.40) discards the sequence information with the model unit since the computation is independent of the order of input vectors. Therefore, one needs to decide the acoustic unit for fast match. In general, the longer the unit, the faster the computation is. and, therefore, the smaller the under-estimated cost  $F(\mathbf{X}|w)$  is. It thus becomes a trade-off between accuracy and speed.

Now let's analyze the real speedup by using fast match to reduce the vocabulary  $V$  to the list  $\Lambda$ , followed by the detailed match. Let  $|V|$  and  $|\Lambda|$  be the sizes for the vocabulary  $V$  and the fast match short list  $\Lambda$ . Suppose  $t_f$  and  $t_d$  are the times required to compute one fast match score and one detailed match score for one word, respectively. Then, the total time required for the fast match followed by the detailed match is  $t_f|V| + t_d|\Lambda|$ , whereas the time required in doing the detailed match alone for the entire vocabulary is  $t_d|V|$ . The speed-up ratio is then given as follows:

$$\frac{1}{\left( \frac{t_f}{t_d} + \frac{|\Lambda|}{|V|} \right)} \quad (12.41)$$

We need  $t_f$  to be much smaller than  $t_d$  and  $|\Lambda|$  to be much smaller than  $|V|$  to have a significant speed-up using fast match. Using our admissible fast match estimator in Eq. (12.40), the time complexity of the computation for  $F(\mathbf{X}|w)$  is  $T$  instead of  $N^2T$  for  $C(\mathbf{X}|w)$ , where  $N$  is the number of states in the detailed acoustic model. Therefore, the  $t_f/t_d$  saving is about  $N^2$ .

In general, in order to make  $|\Lambda|$  much smaller than  $|V|$ , one needs a very accurate fast match estimator that could result in  $t_f \approx t_d$ . This is why we often relax the constraint of admissibility, although it is a nice principle to adhere to. In practice, most real-time speech recognition systems don't necessarily obey the admissibility principle with the fast match. For example, Bahl et al. [10], Laface et al., [22] and Roe et al., [36] used several techniques

to construct off-line groups of acoustically similar words. Armed with this grouping, they can use an aggressive fast match to select only a very short list of words, and words acoustically similar to the words in this list are added to form the short word list  $\Lambda$  for further detailed match processing. By doing so, they are able to report a very efficient fast match method that misses the correct word only 2% of the time. When non-admissible fast match is used, one needs to minimize the additional search error introduced by fast match empirically.

Bahl et al. [6] use a one-state HMM as their fast match units and a tree-structure lexicon similar to the lexical tree structures introduced in Chapter 13 to construct the short word list  $\Lambda$  for next-word expansion in stack decoding. Since the fast match tree search is also done in an asynchronous way, the ending time of each phone is detected using normalized scores similar to those described in Section 12.5.2. It is based on the same idea that this normalized score rises slowly for the correct phone, while it drops rapidly once the end of phone is encountered (so the model is starting to go toward the incorrect phones). During the asynchronous lexical tree search, the unpromising hypotheses are also pruned away by a pruning threshold that is constantly changing once a complete hypothesis (a leaf node) is obtained. On a 20,000-word dictation task, such a fast match scheme was about 100 times faster than detailed match and achieved real-time performance on a commercial workstation with only 0.34% increase in the word error rate being introduced by the fast match process.

#### 12.5.4. Stack Pruning

Even with efficient heuristic functions, the mechanism to determine the ending time for phone/word, and fast match, stack decoding might still be too slow for large-vocabulary speech recognition tasks. A beam within the stack, which saves only a small number of promising hypotheses in the *OPEN* list, is often used to reduce search effort. This *stack pruning* is very similar to beam search. A predetermined threshold  $\epsilon$  is used to eliminate hypotheses whose cost value is much worse than the best path so far.

Both fast match and stack pruning could introduce search errors where the eventual optimal path is thrown away prematurely. However, the impact could be reduced to a minimum by empirically adjusting the thresholds in both methods.

The implementation of stack decoding is, in general, more complicated, particularly when some inevitable pruning strategies are incorporated to make the search more efficient. The difficulty of devising both an effectively admissible heuristic function for  $h(\bullet)$  and an effective estimation of normalization factors for boundary determination has limited the advantage that stack decoders have over Viterbi decoders. Unlike stack decoding, time-synchronous Viterbi beam search can use an easy comparison of same-length path without heuristic determination of word boundaries. As described in the earlier sections, these simple and unified features of Viterbi beam search allow researchers to incorporate various sound techniques to improve the efficiency of search. Therefore, time-synchronous Viterbi Beam search enjoys a much broader popularity in the speech community. However, the principle of stack decoding is essential particularly for n-best and lattice search. As we describe in Chapter 13, stack decoding plays a very crucial part in multiple-pass search strate-

gies for  $n$ -best and lattice search because the early pass is able to establish a near-perfect estimate of the remaining path.

### 12.5.5. Multistack Search

Even with the help of normalized factor  $\gamma$  or heuristic function  $h(\bullet)$ , it is still more effective to compare hypotheses of the same length than those of different lengths, because hypotheses with the same length are compared based on the true forward matching score. Inspired by the time-synchronous principle in Viterbi beam search, researchers [8, 35] propose a variant stack decoding based on multiple stacks.

Multistack search is equivalent to a best-first search algorithm running on multiple stacks time-synchronously. Basically, the search maintains a separate stack for each time frame  $t$ , so it never needs to compare hypotheses of different lengths. The search runs time-synchronously from left to right just like time-synchronous Viterbi search. For each time frame  $t$ , multistack search extracts the best path out of the  $t$ -stack, computes one-word extensions, and places all the new paths into the corresponding stacks. When the search finishes, the top path in the last stack is our optimal path. Algorithm 12.7 illustrates the multistack search algorithm.

This time-synchronous multistack search is designed based on the fact that by the time the  $t^{\text{th}}$  stack is extended, it already contains the best paths that could ever be placed into it. This phenomenon is virtually a variant of the dynamic programming principle introduced in Chapter 8. To make multistack more efficient, some heuristic pruning can be applied to reduce the computation. For example, when the top path of each stack is extended for one more word, we could only consider extensions between minimum and maximum duration. On the other hand, when some heuristic pruning is integrated into the multistack search, one might need to use a small beam in Step 2 of Algorithm 12.7 to extend more than just the best path to guarantee the admissibility.

#### ALGORITHM 12.7: MULTISTACK SEARCH

**Step 1: Initialization:** for each word  $v$  in vocabulary  $V$   
for  $t = 1, 2, \dots, T$

    Compute  $C(x'_t | v)$  and insert it to  $t^{\text{th}}$  stack

**Step 2: Iteration:** for  $t = 1, 2, \dots, T - 1$

    Sort the  $t^{\text{th}}$  stack and pop the top path  $C(x'_t | w_t^k)$  out of the stack

        for each word  $v$  in vocabulary  $V$

            for  $\tau = t + 1, t + 2, \dots, T$

                Extend the path  $C(x'_t | w_t^k)$  by word  $v$  to get  $C(x'_\tau | w_\tau^{k+1})$

                where  $w_\tau^{k+1} = w_t^k || v$  and  $||$  means string concatenation

                Place  $C(x'_\tau | w_\tau^{k+1})$  in  $\tau^{\text{th}}$  stack

**Step 3: Termination:** Sort the  $T^{\text{th}}$  stack and the top path is the optimal word sequence

## 12.6. HISTORICAL PERSPECTIVE AND FURTHER READING

Search has been one of the most important topics in artificial intelligence since the origins of the field. It plays the central role in general problem solving [29] and computer games. [43], Nilsson's *Principles of Artificial Intelligence* [32] and Barr and Feigenbaum's *The Handbook of Artificial Intelligence* [11] contain a comprehensive introduction to state-space search algorithms. A\* search was first proposed by Hart et al. [17]. A\* was thought to be derived from Dijkstra's algorithm [13] and Moore's algorithm [27]. A\* search is similar to the *branch-and-bound* algorithm [23, 39], widely used in *operations research*. The proof of admissibility of A\* search can be found in [32].

The application of *beam search* in speech recognition was first introduced by the HARP system [26]. It wasn't widely popular until BBN used it for their BYBLOS system [37]. There are some excellent papers with detailed description of the use of time-synchronous Viterbi beam search for continuous speech recognition [24, 31]. Over the years, many efficient implementations and improvements have been introduced for time-synchronous Viterbi beam search, so real-time large-vocabulary continuous speech recognition can be realized on a general-purpose personal computer.

On the other hand, stack decoding was first developed by IBM [9]. It is successfully used in IBM's large-vocabulary continuous speech recognition systems [3, 16]. Lacking a time-synchronous framework, comparing theories of different lengths and extending theories are more complex as described in this chapter. Because of the complexity of stack decoding, far fewer publications and systems are based on it than on Viterbi beam search [16, 19, 20, 35]. With the introduction of multistack search [8], stack decoding in essence has actually come very close to time-synchronous Viterbi beam search.

Stack decoding is typically integrated with fast match methods to improve its efficiency. Fast match was first implemented for isolated word recognition to obtain a list of potential word candidates [5, 7]. The paper by Gopalakrishnan et al. [15] contains a comprehensive description of fast match techniques to reduce the word expansion for stack decoding. Besides the fast match techniques described in this chapter, there are a number of alternative approaches [5, 21, 41]. Waast's fast match [41], for example, is based on a binary classification tree built automatically from data that comprise both phonetic transcription and acoustic sequence.

### REFERENCES

- [1] Aho, A., J. Hopcroft, and J. Ullman, *The Design and Analysis of Computer Algorithms*, 1974, Addison-Wesley Publishing Company.
- [2] Alleva, F., X. Huang, and M. Hwang, "An Improved Search Algorithm for Continuous Speech Recognition," *Int. Conf. on Acoustics, Speech and Signal Processing*, 1993, Minneapolis, MN, pp. 307-310.
- [3] Bahl, L.R. and et. al., "Large Vocabulary Natural Language Continuous Speech Recognition," *Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1989, Glasgow, Scotland, pp. 465-467.

- [4] Bahl, L.R., et al., "Language-Model/Acoustic Channel Balance Mechanism," *IBM Technical Disclosure Bulletin*, 1980, 23(7B), pp. 3464-3465.
- [5] Bahl, L.R., et al., "Obtaining Candidate Words by Polling in a Large Vocabulary Speech Recognition System," *Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1988, pp. 489-492.
- [6] Bahl, L.R., et al., "A Fast Approximate Acoustic Match for Large Vocabulary Speech Recognition," *IEEE Trans. on Speech and Audio Processing*, 1993(1), pp. 59-67.
- [7] Bahl, L.R., et al., "Matrix Fast Match: a Fast Method for Identifying a Short List of Candidate Words for Decoding," *Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1989, Glasgow, Scotland, pp. 345-347.
- [8] Bahl, L.R., P.S. Gopalakrishnan, and R.L. Mercer, "Search Issues in Large Vocabulary Speech Recognition," *Proc. of the 1993 IEEE Workshop on Automatic Speech Recognition*, 1993, Snowbird, UT.
- [9] Bahl, L.R., F. Jelinek, and R. Mercer, "A Maximum Likelihood Approach to Continuous Speech Recognition," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 1983(2), pp. 179-190.
- [10] Bahl, L.R., et al., "Constructing Candidate Word Lists Using Acoustically Similar Word Groups," *IEEE Trans. on Signal Processing*, 1992(1), pp. 2814-2816.
- [11] Barr, A. and E. Feigenbaum, *The Handbook of Artificial Intelligence: Volume I*, 1981, Addison-Wesley.
- [12] Cettolo, M., R. Gretter, and R.D. Mori, "Knowledge Integration" in *Spoken Dialogues with Computers*, R.D. Mori, Editor 1998, London, Academic Press, pp. 231-256.
- [13] Dijkstra, E.W., "A Note on Two Problems in Connection with Graphs," *Numerische Mathematik*, 1959, 1, pp. 269-271.
- [14] Gauvain, J.L., et al., "The LIMSI Speech Dictation System: Evaluation on the ARPA Wall Street Journal Corpus," *Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1994, Adelaide, Australia, pp. 129-132.
- [15] Gopalakrishnan, P.S. 2and L.R. Bahl, "Fast Match Techniques," in *Automatic Speech and Speaker Recognition*, C.H. Lee, F.K. Soong, and K.K. Paliwal, eds., 1996, Norwell, MA, Kluwer Academic Publishers, pp. 413-428.
- [16] Gopalakrishnan, P.S., L.R. Bahl, and R.L. Mercer, "A Tree Search Strategy for Large-Vocabulary Continuous Speech Recognition," *Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1995, Detroit, MI, pp. 572-575.
- [17] Hart, P.E., N.J. Nilsson, and B. Raphael, "A Formal Basis for the Heuristic Determination of Minimum Cost Paths," *IEEE Trans. on Systems Science and Cybernetics*, 1968, 4(2), pp. 100-107.
- [18] Huang, X., et al., "Microsoft Windows Highly Intelligent Speech Recognizer: Whisper," *IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1995, pp. 93-96.
- [19] Jelinek, F., *Statistical Methods for Speech Recognition*, 1998, Cambridge, MA, MIT Press.

- [20] Kenny, P., *et al.*, "A\*—Admissible Heuristics for Rapid Lexical Access," *IEEE Trans. on Speech and Audio Processing*, 1993, 1, pp. 49-58.
- [21] Kenny, P., *et al.*, "A New Fast Match for Very Large Vocabulary Continuous Speech Recognition," *IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1993, Minneapolis, MN, pp. 656-659.
- [22] Laface, P., L. Fissore, and F. Ravera, "Automatic Generation of Words toward Flexible Vocabulary Isolated Word Recognition," *Proc. of the Int. Conf. on Spoken Language Processing*, 1994, Yokohama, Japan, pp. 2215-2218.
- [23] Lawler, E.W. and D.E. Wood, "Branch-and-Bound Methods: A Survey," *Operations Research*, 1966(14), pp. 699-719.
- [24] Lee, K.F. and F.A. Alleva, "Continuous Speech Recognition" in *Recent Progress in Speech Signal Processing*, S. Furui and M. Sondhi, eds., 1990, Marcel Dekker, Inc.
- [25] Lee, K.F., H.W. Hon, and R. Reddy, "An Overview of the SPHINX Speech Recognition System," *IEEE Trans. on Acoustics, Speech and Signal Processing*, 1990, 38(1), pp. 35-45.
- [26] Lowerre, B.T., *The HARPY Speech Recognition System*, PhD Thesis in Computer Science Department, 1976, Carnegie Mellon University.
- [27] Moore, E.F., "The Shortest Path Through a Maze," *Int. Symp. on the Theory of Switching*, 1959, Cambridge, MA, Harvard University Press, pp. 285-292.
- [28] Murveit, H., *et al.*, "Large Vocabulary Dictation Using SRI's DECIPHER Speech Recognition System: Progressive Search Techniques," *Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1993, Minneapolis, MN, pp. 319-322.
- [29] Newell, A. and H.A. Simon, *Human Problem Solving*, 1972, Englewood Cliffs, NJ, Prentice Hall.
- [30] Ney, H. and X. Aubert, "Dynamic Programming Search: From Digit Strings to Large Vocabulary Word Graphs," in *Automatic Speech and Speaker Recognition*, C.H. Lee, F. Soong and K.K. Paliwal, eds., 1996, Boston, Kluwer Academic Publishers, pp. 385-412.
- [31] Ney, H. and S. Ortmanns, *Dynamic Programming Search for Continuous Speech Recognition*, in *IEEE Signal Processing Magazine*, 1999, pp. 64-83.
- [32] Nilsson, N.J., *Principles of Artificial Intelligence*, 1982, Berlin, Germany, Springer Verlag.
- [33] Nilsson, N.J., *Artificial Intelligence: A New Synthesis*, 1998, Academic Press/Morgan Kaufmann.
- [34] Normandin, Y., R. Cardin, and R.D. Mori, "High-Performance Connected Digit Recognition Using Maximum Mutual Information Estimation," *IEEE Trans. on Speech and Audio Processing*, 1994, 2(2), pp. 299-311.
- [35] Paul, D.B., "An Efficient A\* Stack Decoder Algorithm for Continuous Speech Recognition with a Stochastic Language Model," *Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1992, San Francisco, California, pp. 25-28.

- [36] Roe, D.B. and M.D. Riley, "Prediction of Word Confusabilities for Speech Recognition," *Proc. of the Int. Conf. on Spoken Language Processing*, 1994, Yokohama, Japan, pp. 227-230.
- [37] Schwartz, R., *et al.*, "Context-Dependent Modeling for Acoustic-Phonetic Recognition of Speech Signals," *Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1985, Tampa, FLA, pp. 1205-1208.
- [38] Steinbiss, V., *et al.*, "The Philips Research System for Large-Vocabulary Continuous-Speech Recognition," *Proc. of the European Conf. on Speech Communication and Technology*, 1993, Berlin, Germany, pp. 2125-2128.
- [39] Taha, H.A., *Operations Research: An Introduction*, 6th ed, 1996, Prentice Hall.
- [40] Tanimoto, S.L., *The Elements of Artificial Intelligence : An Introduction Using Lisp*, 1987, Computer Science Press, Inc.
- [41] Waast, C. and L.R. Bahl, "Fast Match Based on Decision Tree," *Proc. of the European Conf. on Speech Communication and Technology*, 1995, Madrid, Spain, pp. 909-912.
- [42] Winston, P.H., *Artificial Intelligence*, 1984, Reading, MA, Addison-Wesley.
- [43] Winston, P.H., *Artificial Intelligence*, 3rd ed, 1992, Reading, MA, Addison-Wesley.
- [44] Woodland, P.C., *et al.*, "Large Vocabulary Continuous Speech Recognition Using HTK," *Proc. of the IEEE Int. Conf. on Acoustics, Speech and Signal Processing*, 1994, Adelaide, Australia, pp. 125-128.



---

# CHAPTER 13

---

## Large-Vocabulary Search Algorithms

Chapter 12 discussed the basic search techniques for speech recognition. However, the search complexity for large-vocabulary speech recognition with high-order language models is still difficult to handle. In this chapter we describe efficient search techniques in the context of time-synchronous Viterbi beam search, which becomes the choice for most speech recognition systems because it is very efficient. We use Microsoft Whisper as our case study to illustrate the effectiveness of various search techniques. Most of the techniques discussed here can also be applied to stack decoding.

With the help of beam search, it is unnecessary to explore the entire search space or the entire trellis. Instead, only the promising search state-space needs to be explored. Please keep in mind the distinction between the implicit search graph specified by the grammar network and the explicit partial search graph that is actually constructed by the Viterbi beam search algorithm.

In this chapter we first introduce the most critical search organization for large-vocabulary speech recognition—tree lexicons. Tree lexicons significantly reduce potential search space, although they introduce many practical problems. In particular, we need to

address problems such as reentrant lexical trees, factored language model probabilities, subtree optimization, and subtree polymorphism.

Various other efficient techniques also are introduced. Most of these techniques aim for clever pruning with the hope of sparing the correct paths. For more effective pruning, different layers of beams are usually used. While fast match techniques described in Chapter 12 are typically required for stack decoding, similar concepts and techniques can be applied to Viterbi beam search. In practice, the look-ahead strategy is equally effective for Viterbi beam search.

Although it is always desirable to use all the knowledge sources (KSs) in the search algorithm, some are difficult to integrate into the left-to-right time-synchronous search framework. One alternative strategy is to first produce an ordered list of sentence hypotheses (a.k.a. *n-best list*), or a lattice of word hypotheses (a.k.a. *word lattice*) using relatively inexpensive KSs. More expensive KSs can be used to rescore the *n-best list* or the word lattice to obtain the refined result. Such a multipass strategy has been explored in many large-vocabulary speech recognition systems. Various algorithms to generate sufficient *n-best lists* or the word lattices are described in the section on multipass search strategies.

Most of the techniques described in this chapter rely on nonadmissible heuristics. Thus, it is critical to derive a framework to evaluate different search strategies and pruning parameters.

## 13.1. EFFICIENT MANIPULATION OF A TREE LEXICON

The lexicon entry is the most critical component for large-vocabulary speech recognition, since the search space grows linearly along with increased linear vocabulary. Thus an efficient framework for handling large vocabulary undoubtedly becomes the most critical issue for efficient search performance.

### 13.1.1. Lexical Tree

The search space for *n*-gram discussed in Chapter 12 is organized based on a straightforward linear lexicon, i.e., each word is represented as a linear sequence of phonemes, independent of other words. For example, the phonetic similarity between the words *task* and *tasks* is not leveraged. In a large-vocabulary system, many words may share the same beginning phonemes. A tree structure is a natural representation for a large-vocabulary lexicon, as many phonemes can be shared to eliminate redundant acoustic evaluations. The lexical tree-based search is thus essential for building a real-time<sup>1</sup> large-vocabulary speech recognizer.

<sup>1</sup> The term *real-time* means the decoding process takes no longer than the duration of the speech. Since the decoding process can take place as soon as the speech starts, such a real-time decoder can provide real instantaneous responses after speakers finish talking.

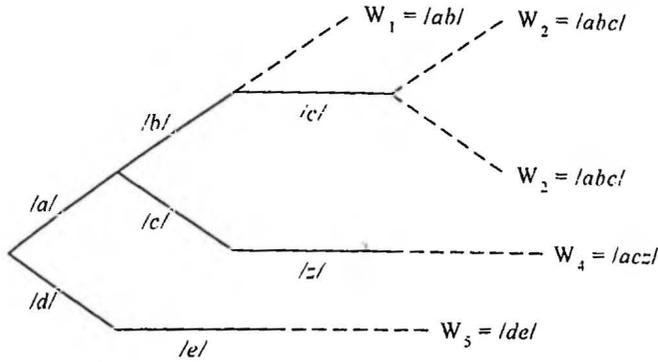


Figure 13.1 An example of a lexical tree, where each branch corresponds to a shared phoneme and the leaf corresponds to a word.

Figure 13.1 shows an example of such a lexical tree, where common beginning phonemes are shared. Each leaf corresponds to a word in the vocabulary. Please note that an extra null arc is used to form the leaf node for each word. This null arc has the following two functions:

1. When the pronunciation transcription of a word is a prefix of other ones, the null arc can function as one branch to end the word.
2. When there are homophones in the lexicon, the null arcs can function as linguistic branches to represent different words such as *two* and *to*.

The advantage of using such a lexical tree representation is obvious: it can effectively reduce the state search space of the trellis. Ney et al. [32] reported that a lexical tree representation of a 12,306-word lexicon with only 43,000 phoneme arcs had a saving of a factor of 2.5 over the linear lexicon with 100,800 phoneme arcs. Lexical trees are also referred to as *prefix trees*, since they are efficient representations of lexicons with sharing among lexical entries that have a common prefix. Table 13.1 shows the distribution of phoneme arcs for this 12,306-word lexical tree. As one can see, even in the fifth level the number of phoneme arcs is only about one-third of the total number of words in the lexicon.

Table 13.1 Distribution of the tree phoneme arcs and active tree phoneme arc for a 12,306-word lexicon using a lexical tree representation [32].

Level	1	2	3	4	5	6	≥7
Phoneme arcs	28	331	1511	3116	4380	4950	29,200
Average active arcs	23	233	485	470	329	178	206

The saving by using a lexical tree is substantial, because it not only results in considerable memory saving for representing state-search space but also saves tremendous time by searching far fewer potential paths. Ney et al. [32] report that a tree organization of the lexicon reduces the total search effort by a factor of 7 over the linear lexicon organization. This is because the lion's share of hypotheses during a typical large-vocabulary search is on the first and second phonemes of a word. Haeb-Umbach et al. [23] report that for a 12,306-word dictation task, 79% and 16% of the state hypotheses are in the first and second phonemes, when analyzing the distribution of the state hypotheses over the state position within a word. Obviously, the effect is caused by the ambiguities at the word boundaries. The lexical tree representation reduces that effort by evaluating common phonetic prefixes only once. Table 13.1 also shows the average number of active phoneme arcs in the layers of the lexical tree [32]. Based on this table, you can expect that the overall search cost is far less than the size of the vocabulary. This is the key reason why lexical tree search is widely used for large-vocabulary continuous speech recognition systems.

The lexical tree search requires a sophisticated implementation because of a fundamental deficiency—a branch in a lexical tree representation does not correspond to a single word with the exception of branches ending in a leaf. This deficiency translates to the fact that a unique word identity is not determined until a leaf of the tree is reached. This means that any decision about the word identity needs to be delayed until the leaf node is reached, which results in the following complexities.

- Unlike a linear lexicon, where the language model score can be applied when starting the acoustic search of a new word, the lexical tree representation has to delay the application of the language model probability until the leaf is reached. This may result in an increased search effort, because the pruning needs to be done on a less reliable measure, unless a factored language model is used, as discussed in Section 13.1.3.
- Because of the delay of language model contribution by one word, we need to keep a separate copy of an entire lexical tree for each unique language model history.

### 13.1.2. Multiple Copies of Pronunciation Trees

A simple lexical tree is sufficient if no language model or a unigram is used. This is because the decision at time  $t$  depends on the current word only. However, for higher-order  $n$ -gram models, the linguistic state cannot be determined locally. A tree copy is required for each language model state. For bigrams, a tree copy is required for each predecessor word. This may seem to be astonishing, because the potential search space is increased by the vocabulary size. Fortunately, experimental results show only a small number of tree copies are required, because efficient pruning can eliminate most of the unneeded ones. Ney et al. [32] report that the search effort using bigrams is increased by only a factor of 2 over the unigram

case. In general, when more detailed (better) acoustic and/or language models are used, the effect of a potentially increased search space is often compensated by a more focused beam search from the use of more accurate models. In other words, although the static search space might increase significantly by using more accurate models, the dynamic search space can be under control (sometimes even smaller), thanks to improved evaluation functions.

To deal with tree copies [19, 23, 37], you can create redundant subtrees. When copies of lexical trees are used to disambiguate active linguistic contexts, many of the active state hypotheses correspond to the same redundant unigram state. due to the postponed application of language models. To apply the language model sooner, and to eliminate redundant unigram state computations, a successor tree,  $T_i$ , can be created for each linguistic context  $i$ .  $T_i$  encodes the nonzero  $n$ -grams of the linguistic context  $i$  as an isomorphic subgraph of the unigram tree,  $T_u$ . Figure 13.2 shows the organization of such successor trees and unigram tree for bigram search. For each word  $w$  a successor tree,  $T_w$  is created with the set of successor words that have nonzero bigram probabilities. Suppose  $u$  is a successor of  $w$ ; the bigram probability  $P(u|w)$  is attached to the transition connecting the leaf corresponding to  $u$  in the successor tree  $T_w$ , with the root of the successor tree  $T_w$ . The unigram tree is a full-size lexical tree and is shared by all words as the back-off lexical tree. Each leaf of the unigram tree corresponds to one of  $|V|$  words in the vocabulary and is linked to the root of its bigram successor tree ( $T_u$ ) by an arc with the corresponding unigram probability  $P(u)$ . The backoff weight,  $\alpha(u)$ , of predecessor  $u$  is attached to the arc which links the root of successor tree  $T_u$  to the root of the unigram tree.

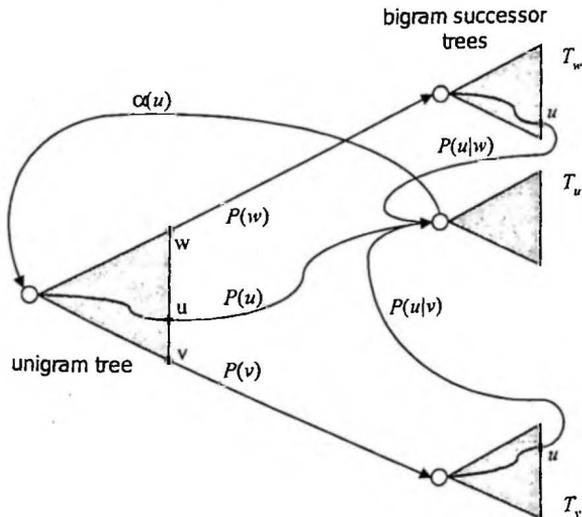


Figure 13.2 Successor trees and unigram trees for bigram search [13].

A careful search organization is required to avoid computational overhead and to guarantee a linear time complexity for exploring state hypotheses. In the following sections we describe techniques to achieve efficient lexical tree recognizers. These techniques include factorization of language model probabilities, tree optimization, and exploiting subtree dominance.

### 13.1.3. Factored Language Probabilities

As mentioned in Section 13.1.2, search is more efficient if a detailed knowledge source can be applied at an early stage. The idea of *factoring* the language model probabilities across the tree is one such example [4, 19]. When more than one word shares a phoneme arc, the upper bound of their probability can be associated to that arc.<sup>2</sup> The factorization can be applied to both the full lexical tree (unigram) and successor trees (bigram or other higher-order language models).

An unfactored tree only has language model probabilities attached to the leaf nodes, and all the internal nodes have probability 1.0. The procedure for factoring the probabilities across the tree computes the maximum of each node  $n$  in the tree according to Eq. (13.1). The tree can then be factored according to Eq. (13.2) so when you traverse the tree you can multiply  $F^*(n)$  along the path to get the needed language probability.

$$P^*(n) = \max_{x \in \text{child}(n)} P(x) \quad (13.1)$$

$$F^*(n) = \frac{P^*(n)}{P^*(\text{parent}(n))} \quad (13.2)$$

An illustration of the factored probabilities is shown in Table 13.2. Using this lexicon, we create the tree depicted in Figure 13.3(a). In this figure the unlabeled internal nodes have a probability of 1.0. We distribute the probabilities according to Eq. (13.1) in Figure 13.3(b), which is factored according to Eq. (13.2), resulting in Figure 13.3(c).

**Table 13.2** Sample probabilities  $P(w_j)$  and their pseudoword pronunciations [4].

$w_j$	Pronunciation	$P(w_j)$
$w_0$	/a b c/	0.1
$w_1$	/a b c/	0.4
$w_2$	/a c z/	0.3
$w_3$	/d e/	0.2

<sup>2</sup> The choice of upper bound is because it is an admissible estimate of the path no matter which word will be chosen later.

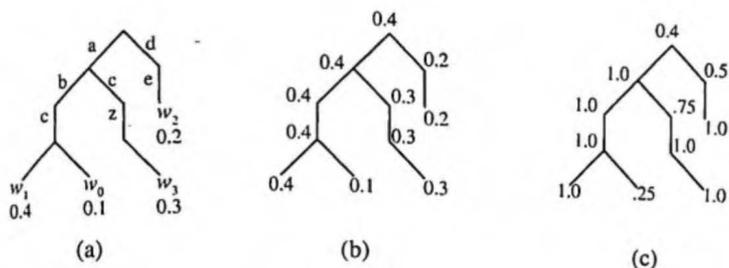


Figure 13.3 (a) Unfactored lexical tree; (b) distributed probabilities with computed  $P^*(n)$ ; (c) factored tree  $F^*(n)$  [4].

Using the upper bounds in the factoring algorithm is not an approximation, since the correct language model probabilities are calculated by the product of values traversed along each path from the root to the leaves. However, you should note that the probabilities of all the branches of a node do not sum to one. This can be solved by replacing the upper-bound (max) function in Eq. (13.1) with the sum.

$$P^*(n) = \sum_{x \in \text{child}(n)} P(x) \tag{13.3}$$

To guarantee that all the branches sum to one, Eq. (13.2) should also be replaced by the following equation:

$$F^*(n) = \frac{P^*(n)}{\sum_{x \in \text{child}(\text{parent}(n))} P^*(x)} \tag{13.4}$$

A new illustration of the distribution of LM probabilities by using sum instead of upper bound is shown in Figure 13.4. Experimental results have shown that the factoring method with either sum or upper bound has comparable search performance.

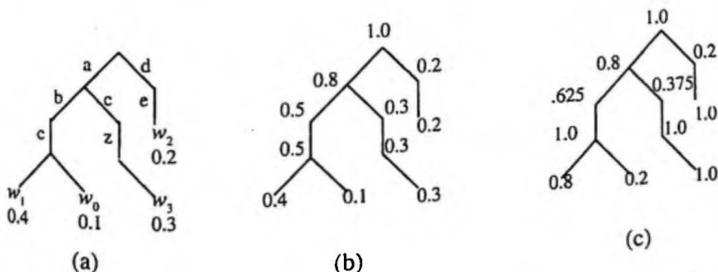


Figure 13.4 Using sum instead of upper bound when factoring tree, the corresponding (a) unfactored lexical tree; (b) distributed probabilities with computed  $P^*(n)$ ; (c) factored tree with computed  $F^*(n)$  [4].

One interesting observation is that the language model score can be regarded as a heuristic function to estimate the linguistic expectation of the current word to be searched. In a linear representation of the pronunciation lexicon, application of the linguistic expectation was straightforward, since each state is associated with a unique word. Therefore, given the context defined by the hypothesis under consideration, the expectation for the first phone of word  $w_i$  is just  $P(w_i | w_i^{-1})$ . After the first phone, the expectation for the rest of the phones becomes 1.0, since there is only one possible phone sequence when searching the word  $w_i$ . However, for the tree lexicon, it is necessary to compute  $E(p_j | p_i^{j-1}, w_i^{-1})$ , the expectation of phone  $p_j$  given the phonetic prefix  $p_i^{j-1}$  and the linguistic context  $w_i^{-1}$ . Let  $\phi(j, w_i)$  denote the phonetic prefix of length  $j$  for  $w_i$ . Based on Eqs. (13.1) and (13.2), we can compute the expectation as:

$$E(p_j | p_i^{j-1}, w_i^{-1}) = \frac{P(w_c | w_i^{-1})}{P(w_p | w_i^{-1})} \quad (13.5)$$

where  $c = \arg \max_k (w_k | w_i^{-1}, \phi(j, w_k) = p_i^j)$  and  $p = \arg \max_k (w_k | w_i^{-1}, \phi(j-1, w_k) = p_i^{j-1})$ . Based on Eq. (13.5), an arbitrary  $n$ -gram model or even a stochastic context-free grammar can be factored accordingly.

### 13.1.3.1. Efficient Memory Organization of Factored Lexical Trees

A major drawback to the use of successor trees is the large memory overhead required to store the additional information that encodes the structure of the tree and the factored linguistic probabilities. For example, the 5.02 million bigrams in the 1994 NABN (North American Business News) model require 18.2 million nodes. Given a compact binary tree representation that uses 4 bytes of memory per node, 72.8 million bytes are required to store the predecessor-dependent lexical trees. Furthermore, this tree representation is not as amenable to data compression techniques as the linear bigram representation.

The factored probability of successor trees can be encoded as efficiently as the  $n$ -gram model based on Algorithm 13.1, i.e., one  $n$ -gram record results in one constant-sized record. Step 3 is illustrated in Figure 13.5(b), where the heavy line ends at the most recently visited node that is not a direct ancestor. The encoding result is shown in Table 13.3.

#### ALGORITHM 13.1: ENCODING THE LEXICAL SUCCESSOR TREES (LST)

For each linguistic context:

**Step 1:** Distribute the probabilities according to Eq. (13.1).

**Step 2:** Factor the probabilities according to Eq. (13.2).

**Step 3:** Perform a depth-first traversal of the LST and encode each leaf record,  
 (a) the depth of the most recently visited node that is not a direct ancestor,  
 (b) the probability of the direct ancestor at the depth in (a),  
 (c) the word identity.

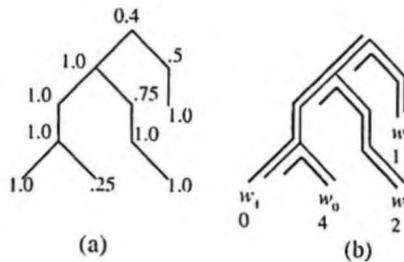


Figure 13.5 (a) Factored tree; (b) tree with common prefix-length annotation.

Clearly the new data structure meets the requirements set forth, and, in fact, it only requires additional  $\log(n)$  bits per record ( $n$  is the depth of the tree). These bits encode the common prefix length for each word. Naturally this requires some modification to the decoding procedure. In particular, the decoder must scan a portion of the  $n$ -gram successor list in order to determine which tree nodes should be activated. Depending on the structure of the tree (which is determined by the acoustic model, the lexicon, and language model), the tree structure can be interpreted at runtime or cached for rapid access if memory is available.

Table 13.3 Encoded successor lexical tree; each record corresponds to one augmented factored  $n$ -gram.

$w_j$	Depth	$F^*(w_j)$
$w_1$	0	0.4
$w_0$	4	0.25
$w_3$	2	0.75
$w_2$	1	0.5

### 13.1.4. Optimization of Lexical Trees

We now investigate ways to handle the huge search network formed by the multiple copies of lexical trees in different linguistic contexts. The factorization of lexical trees actually makes it easier to search. First, after the factorization of the language model, the intertree transitions shown in Figure 13.2 no longer have the language model scores attached because they are already applied completely before leaving the leaves. Moreover, as illustrated in Figure 13.3, many transitions toward the end of a single-word path now have an associated transition probability that is equal to 1. This observation implies that there could be many duplicated subtrees in the network. Those duplicated subtrees can then be merged to save both space and computation by eliminating redundant (unnecessary) state evaluation. Unlike pruning, this saving is based on the dynamic programming principle, without introducing any potential error.

### 13.1.4.1. Optimization of Finite State Network

One way to compress the lexical tree network is to use a similar algorithm for optimizing the number of states in a deterministic finite state automaton. The optimization algorithm is based on the *indistinguishable* property of states in a finite state automaton. Suppose that  $s_1$  and  $s_2$  are the initial states for automata  $T_1$  and  $T_2$ , then  $s_1$  and  $s_2$  are said to be *indistinguishable* if the languages accepted by automata  $T_1$  and  $T_2$  are exactly the same. If we consider our lexical tree network as a finite state automaton, the symbol emitted from the transition arc includes not only the phoneme identity, but also the factorized language model probability.

The general set-partitioning algorithm [1] can be used for the reduction of finite state automata. The algorithm starts with an initial partition of the automaton states and iteratively refines the partition so that two states  $s_1$  and  $s_2$  are put in the same block  $B_i$  if and only if  $f(s_1)$  and  $f(s_2)$  are both in the same block  $B_j$ . For our purpose,  $f(s_1)$  and  $f(s_2)$  can be defined as the destination state given a phone symbol (in the factored trees, the pair  $\langle \text{phone}, \text{LM-probability} \rangle$  can be used). Each time a block is partitioned, the smaller subblock is used for further partitioning. The algorithm stops when all the states that transit to some state in a particular block with arcs labeled with the same symbol are in the same block. When the algorithm halts, each block of the resulting partition is composed of *indistinguishable* states, and those states within each block can then be merged. The algorithm is guaranteed to find the automaton with the minimum number of states. The algorithm has a time complexity of  $O(MN \log N)$ , where  $M$  is the maximum number of branching (fan-out) factors in the lexical tree and  $N$  is the number of states in the original tree network.

Although the above algorithm can give optimal finite state networks in terms of number of states, such an optimized network may be difficult to maintain, because the original lexical tree structure could be destroyed and it may be troublesome to add any new word into the tree network [1].

### 13.1.4.2. Subtree Isomorphism

The finite state optimization algorithm described above does not take advantage of the tree structure of the finite state network, though it generates a network with a minimum number of states. Since our finite state network is a network of trees, the indistinguishability property is actually the same as the definition of subtree isomorphism. Two subtrees are said to be *isomorphic* to each other if they can be made equivalent by permuting the successors. It should be straightforward to prove that two states are indistinguishable, if and only if their subtrees are isomorphic.

There are efficient algorithms [1] to detect whether two subtrees are isomorphic. For all possible pairs of states  $u$  and  $v$ , if the subtrees starting at  $u$  and  $v$ ,  $ST(u)$  and  $ST(v)$ , are isomorphic,  $v$  is merged into  $u$  and  $ST(v)$  can be eliminated. Note that only internal nodes need to be considered for subtree isomorphism check. The time complexity for this algorithm is  $O(N^2)$  [1].

### 13.1.4.3. Sharing Tails

A *linear tail* in a lexical tree is defined as a subpath ending in a leaf and going through states with a unique successor. It is often referred as a *single-word subpath*. It can be proved that such a linear tail has unit probability attached to its arcs according to Eqs. (13.1) and (13.2). This is because LM probability factorization *pushes forward* the LM probability attached to the last arc of the linear tail, leaving arcs with unit probability. Since all the tails corresponding to the same word  $w$  in different successor trees are linked to the root of successor tree  $T_w$ , the subtree starting from the first state of each linear tail is isomorphic to the subtree starting from one of the states forming the longest linear tail of  $w$ . A simple algorithm to take advantage of this share-tail topology can be employed to reduce the lexical tree network.

Figure 13.6 and Figure 13.7 show a lexical tree network before and after shared-tail optimization. For each word, only the longest linear tail is kept. All other tails can be removed by linking them to an appropriate state in the longest tail, as shown in Figure 13.7.

Shared-tail optimization is not global optimization, because it considers only some special topology optimization. However, there are some advantages associated with shared-tail optimization. First, in practice, duplicated linear tails account for most of the redundancy in lexical tree networks [12]. Moreover, shared-tail optimization has a nice property of maintaining the basic lexical tree structure for the optimized tree network.

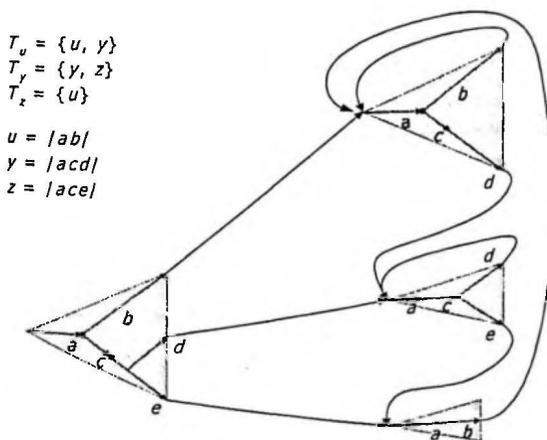


Figure 13.6 An example of a lexical tree network without shared-tail optimization [12]. The vocabulary includes three words,  $u$ ,  $y$ , and  $z$ .  $T_u$ ,  $T_y$ , and  $T_z$  are the successor trees for  $u$ ,  $y$ , and  $z$  respectively [13].

<sup>1</sup>We assume bigram is used in the discussion of "sharing tails."

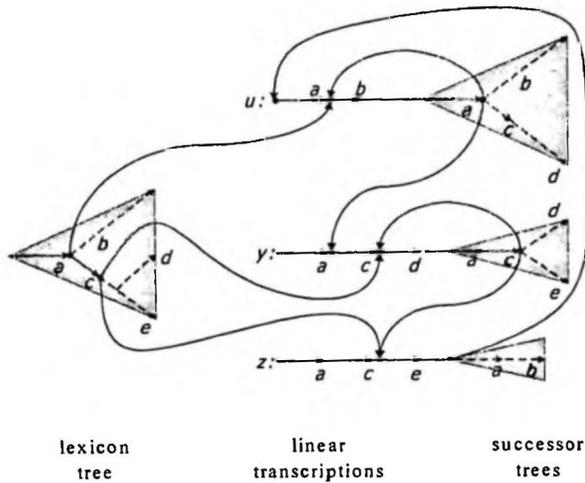


Figure 13.7 The lexical tree network in Figure 13.6 after shared-tail optimization [12].

### 13.1.5. Exploiting Subtree Polymorphism

The techniques of optimizing the network of successor lexical trees can only eliminate identical subtrees in the network. However, there are still many subtrees that have the same nodes and topology but with different language model scores attached to the arcs. The acoustic evaluation for those subtrees is unnecessarily duplicated. In this section we exploit *subtree dominance* for additional saving.

A subtree instance is *dominated* when the best outcome in that subtree is not better than the worst outcome in another instance of that subtree. The evaluation becomes redundant for the dominated subtree instance. Subtree isomorphism and shared-tail are cases of subtree dominance, but they require prearrangement of the lexical tree network as described in the previous section.

If we need to implement lexical tree search dynamically, the network optimization algorithms are not suitable. Although subtree dominance can be computed using minimax search [35] during runtime, this requires that information regarding subtree isomorphism be available for all corresponding pairs of states for each successor tree  $T_w$ . Unfortunately, it is not practical in terms of either computation or space.

In place of computing strict subtree dominance, a *polymorphic* linguistic context assignment to reduce redundancy is employed by estimating subtree dominance based on local information and ignoring the subgraph isomorphism problem. Polymorphic context assignment involves keeping a single copy of the lexical tree and allowing each state to assume the linguistic context of the most promising history. The advantage of this approach is that it employs maximum sharing of data structures and information, so each node in the tree is

evaluated, at most, once. However, the use of local knowledge to determine the dominant context could introduce significant errors because of premature pruning. Whisper [4] reports a 65.7% increase in error rate when only the dominant context is kept, based on local knowledge.

To recover the errors created by using local linguistic information to estimate subtree dominance, you need to delay the decision regarding which linguistic context is most promising. This can be done by keeping a heap of contexts at each node in the tree. The heap maintains all contexts (linguistic paths) whose probabilities are within a constant threshold  $\epsilon$ , of that of the best global path. The effect of the  $\epsilon$ -heap is that more contexts are retained for high-probability states in the lexical tree. The pseudocode fragment in Algorithm 13.2 [3] illustrates a transition from state  $s_n$  in context  $c$  to state  $s_m$ . The terminology used in Algorithm 13.2 is listed as follows:

- $(-\log P(s_m | s_n, c))$  is the cost associated with applying acoustic model matching and language model probability of state  $s_m$  transited from  $s_n$  in context  $c$ .
- $InHeap(s_m, c)$  is true if context  $c$  is in the heap corresponding to state  $s_m$ .
- $Cost(s_m, c)$  is the cost for context  $c$  in state  $s_m$ .
- $StateInfo(s_m, c)$  is the auxiliary state information associated with context  $c$  in state  $s_m$ .
- $Add(s_m, c)$  adds context  $c$  to the state  $s_m$  heap.
- $Delete(s_m, c)$  deletes context  $c$  from state  $s_m$  heap.
- $WorstContext(s_m)$  retrieves the worst context from the heap of state  $s_m$ .

**ALGORITHM 13.2: HANDLING MULTIPLE LINGUISTIC CONTEXTS  
IN A LEXICAL TREE**

```

1.  $d = Cost(s_n, c) + (-\log P(s_m | s_n, c))$ 
2. if  $InHeap(s_m, c)$  then
    if  $d < Cost(s_m, c)$  then
         $Cost(s_m, c) = d$ 
         $StateInfo(s_m, c) = StateInfo(s_n, c)$ 
    else if  $d < BestCost(s_m) + \epsilon$  then
         $Add(s_m, c)$ ;  $StateInfo(s_m, c) = StateInfo(s_n, c)$ 
         $Cost(s_m, c) = d$ 
    else
         $w = WorstContext(s_m)$ 
        if  $d < Cost(s_m, w)$  then
             $Delete(s_m, w)$ 
             $Add(s_m, c)$ ;  $StateInfo(s_m, c) = StateInfo(s_n, c)$ 
             $Cost(s_m, c) = d$ 

```

When higher-order  $n$ -gram is used for lexical tree search, the potential heap size for lexical tree nodes (some also refer to *prefix nodes*) could be unmanageable. With decent acoustic models and efficient pruning, as illustrated in Algorithm 13.2, the average heap size for active nodes in the lexical tree is actually very modest. For example, Whisper's average heap size for active nodes in the 20,000-word WSJ lexical tree decoder is only about 1.6 [3].

### 13.1.6. Context-Dependent Units and Inter-Word Triphones

So far, we have implicitly assumed that context-independent models are used in the lexical tree search. When context-dependent phonetic or subphonetic models, as discussed in Chapter 9, are used for better acoustic models, the construction and use of a lexical tree become more complicated.

Since senones represent both subphonetic and context-dependent acoustic models, this presents additional difficulty for use in lexical trees. Let's assume that a three-state context-dependent HMM is formed from three senones, one for each state. Each senone is context-dependent and can be shared by different allophones. If we use allophones as the units for lexical tree, the sharing may be poor and fan-out unmanageable. Fortunately, each HMM is uniquely identified by the sequence of senones used to form the HMM. In this way, different context-dependent allophones that share the same *senone sequence* can be treated as the same. This is especially important for lexical tree search, since it reduces the order of the fan-out in the tree.

Interword triphones that require significant fan-ins for the first phone of a word and fan-outs for the last phones usually present an implementation challenge for large-vocabulary speech recognition. A common approach is to delay full interword modeling until a subsequent rescoring phase.<sup>4</sup> Given a sufficiently rich lattice or word graph, this is a reasonable approach, because the static state space in the successive search has been reduced significantly. However, as pointed out in Section 13.1.2, the size of the dynamic state space can remain under control when detailed models are used to allow effective pruning. In addition, a multipass search requires an augmented set of acoustic models to effectively model the biphone contexts used at word boundaries for the first pass. Therefore, it might be desirable to use genuine interword acoustic models in the single-pass search.

Instead of expanding all the fan-ins and fan-outs for inter-word context-dependent phone units in the lexical tree, three *metaunits* are created.

1. The first metaunit, which has a known right context corresponding to the second phone in the word, but uses open left context for the first phone of a word (sometimes referred to as the *word-initial unit*). In this way, the fan-in is represented as a subgraph shared by all words with the same initial left-context-dependent phone.

---

<sup>4</sup> Multipass search strategy is described in Section 13.3.5.

2. Another metaunit, which has a known left context corresponding to the second-to-last phone of the word, but uses open right context for the last phone of a word (sometimes referred to as the *word-final unit*). Again, the fan-out is represented as a subgraph shared by all words with the same final right-context-dependent phone.
3. The third metaunit, which has both open left and right contexts, and is used for single-phone word unit.

By using these metaunits we can keep the states for the lexical trees under control, because the fan-in and fan-out are now represented as a single node.

During recognition, different left or right contexts within the same metaunit are handled using Algorithm 13.2, where the different acoustic contexts are treated similarly as different linguistic contexts. The open left-context metaunit (fan-ins) can be dealt with in a straightforward way using Algorithm 13.2, because the left context is always known (the last phone of the previous word) when it is initiated. On the other hand, the open right-context metaunit (fan-out) needs to explore all possible right contexts because the next word is not known yet. To reduce unnecessary computation, fast match algorithms (described in Section 13.2.3) can be used to provide both expected acoustic and language scores for different context-dependent units to result in early pruning of unpromising contexts.

## 13.2. OTHER EFFICIENT SEARCH TECHNIQUES

Tree structured lexicon represents an efficient framework of manipulation of search space. In this section we present some additional implementation techniques, which can be used to further improve the efficiency of search algorithms. Most of these techniques can be applied to both Viterbi beam search and stack decoding. They are essential ingredients for a practical large-vocabulary continuous speech recognizer.

### 13.2.1. Using Entire HMM as a State in Search

The state in state-search space based on HMM-trellis computation is, by definition, a Markov state. Phonetic HMM models are the basic unit in most speech recognizers. Even though subphonetic HMMs, like senones, might be used for such a system, the search is often based on phonetic HMMs.

Treating the entire phonetic HMM as a state in state-search has many advantages. The first obvious advantage is that the number of states the search program needs to deal with is smaller. Note that using the entire phonetic HMM does not in effect reduce the number of states in the search. The entire search space is unchanged. All the states within a phonetic HMM are now bundled together. This means that all of them are either kept in the beam, if the phonetic HMM is regarded as promising, or all of them are pruned away. For any given time, the minimum cost among all the states within the phonetic HMM is used as the cost for the phonetic HMM. For pruning purposes, this cost is used to determine the promising

degree of this phonetic HMM, i.e., the fate of all the states within this phonetic HMM. Although this does not actually reduce the beam beyond normal pruning, it has the effect of processing fewer candidates in the beam. In programming, this means less checking and bookkeeping, so some computation savings can be expected.

You might wonder if this organization might be ineffective for beam search, since it forces you to keep or prune all the states within a phonetic HMM. In theory, it is possible that only one or two states in the phonetic HMM need to be kept, while other states can be pruned due to high cost score. However, this is, in reality, very rare, since a phone is a small unit and all the states within a phonetic HMM should be relatively promising when the search is near the acoustic region corresponding to the phone.

During the trellis computation, all the phonetic HMM states need to advance one time step when processing one input vector. By performing HMM computation for all states together, the new organization can reduce memory accesses and improve cache locality, since the output and transition probabilities are held in common by all states. Combining this organization strategy with lexical tree search further enhances the efficiency. In lexical tree search, each hypothesis in the beam is associated with a particular node in the lexical tree. These hypotheses are linked together in the heap structure described in Algorithm 13.2 for the purposes of efficient evaluation and heuristic pruning. Since the node corresponds to a phonetic HMM, the HMM evaluation is guaranteed to execute once for each hypothesis sharing this node.

In summary, treating the entire phonetic HMM as a state in state-search space allows you to explore the effective data structure for better sharing and improved memory locality.

### 13.2.2. Different Layers of Beams

Because of the complexity of search, it often requires pruning of various levels of search to make search feasible. Most systems thus employ different pruning thresholds to control what states participate. The most frequently used thresholds are listed below:

- $\tau_s$  controls what states (either phone states or senone states) to retain. This is the most fundamental beam threshold.
- $\tau_p$  controls whether the next phone is extended. Although this might not be necessary for both stack decoding and linear Viterbi beam search, it is crucial for lexical tree search, because pruning unpromising phonetic prefixes in the lexical trees could improve search efficiency significantly.
- $\tau_w$  controls whether hypotheses are extended for the next word. Since the branching factor for word boundaries is very large, we need this threshold to limit search to only the promising ones.
- $\tau_c$  controls where a linguistic context is created in a lexical tree search using higher-order language models. This is also known as  $\epsilon$ -heap in Algorithm 13.2.

Pruning can introduce search errors if a state is pruned that would have been on the globally best path. The principle applied here is that the more constraints you have available, the more aggressively you decide whether this path will participate in the globally best path. In this case, at the state level, you have the least constraints. At the phonetic level there are more, and there are the most at the word level. In general, the number of word hypotheses tends to drop significantly at word boundaries. Different thresholds for different levels allow the search designer to fine-tune those thresholds for their tasks to achieve best search performance without significant increase in error rates.

### 13.2.3. Fast Match

As described in Chapter 12, fast match is a crucial part of stack decoding, which mainly reduces the number of possible word expansions for each path. Similarly, fast match can be applied to the most expensive part—extending the phone HMM fan-outs within or between lexical trees. Fast match is a method for rapidly deriving a list of candidates that constrain successive search phases in which a computationally expensive *detailed match* is performed. In this sense, fast match can be regarded as an additional pruning threshold to meet before a new word/phone can be started.

Fast match is typically characterized by the approximations that are made in the acoustic/language models to reduce computation. The factorization of language model scores among tree branches in lexical trees described in Section 13.1.3 can be viewed as fast match using a language model. The factorized method is also an admissible estimate of the language model scores for the future word. In this section we focus on acoustic model fast match.

#### 13.2.3.1. Look-Ahead Strategy

Fast match, when applied in time-synchronous search, is also called *look-ahead* strategy. since it basically searches ahead of the time-synchronous search by a few frames to determine which words or phones are likely to extend. Typically the look-ahead frames are fixed, and the fast match is also done in time-synchronous fashion with another specialized beam for efficient pruning. You can also use simplified models, like the one-state HMMs or context-independent models [4, 32]. Some systems [21, 22] have tried to simplify the level of details in the input feature vectors by aggregating information from several frames into one. A straightforward way for compressing the feature stream is to skip every other frame of speech for fast match. This allows a longer-range look-ahead, while keeping computation under control. The approach of simplifying the input feature stream instead of simplifying the acoustic models can reuse the fast match results for detailed match.

Whisper [4] uses phoneme look-ahead fast match in lexical tree search, in which pruning is applied based on the estimation of the score of possible phone fan-outs that may follow a given phone. A context-independent phone-net is searched synchronously with the

search process but offset  $N$  frames into the future. In practice, significant savings can be obtained in search efforts without increase in error rates.

The performance of word and phoneme look-ahead clearly depends on the length of the look-ahead frames. In general, the larger the look-ahead window, the longer is the computation and the shorter the word/phone  $\Lambda$  list. Empirically, the window is a few tens of milliseconds for phone look-ahead and a few hundreds of milliseconds for word look-ahead.

### 13.2.3.2. The Rich-Get-Richer Strategy

For systems employing continuous-density HMMs, tens of mixtures of Gaussians are often used for the output probability distribution for each state. The computation of the mixtures is one of the bottlenecks when many context-dependent models are used. For example, Whisper uses about 120,000 Gaussians. In addition to using various beam pruning thresholds in the search, there could be significant savings if we have a strategy to limit the number of Gaussians to be computed.

The *Rich-Get-Richer* (RGR) strategy enables us to focus on most promising paths and treat them with detailed acoustic evaluations and relaxed path-pruning thresholds. On the contrary, the less promising paths are extended with less expensive acoustic evaluations and less forgiving path-pruning thresholds. In this way, locally optimal candidates continue to receive the maximum attention while less optimal candidates are retained but evaluated using less precise (computationally expensive) acoustic and/or linguistic models. The RGR strategy gives us finer control in the creation of new paths that has potential to grow exponentially.

RGR is used to control the level of acoustic details in the search. The goal is to reduce the number of context-dependent senone probability (Gaussian) computations required. The context-dependent senones associated with a phone instance  $p$  would be evaluated according to the following condition:

$$\begin{aligned} \text{Min}[ci(p)] * \alpha + \text{LookAhead}[ci(p)] &< \text{threshold} \\ \text{where } \text{Min}[ci(p)] &= \min_s \{ \text{cost}(s) \mid s \in ci\_phone(p) \} \\ \text{and } \text{LookAhead}[ci(p)] &= \text{look-ahead estimate of } ci(p) \end{aligned} \quad (13.6)$$

These conditions state that the context-dependent senones associated with  $p$  should be evaluated if there exists a state  $s$  corresponding to  $p$ , whose cost in linear combination with a look-ahead cost score corresponding to  $p$  falls within a threshold. In the event that  $p$  does not fall within the threshold, the senone scores corresponding to  $p$  are estimated using the context-independent senones corresponding to  $p$ . This means the context-dependent senones are evaluated only if the corresponding context-independent senones and the look-ahead start showing promise. RGR strategy should save significant senone computation for clearly unpromising paths. Whisper [26] reports that 80% of senone computation can be avoided without introducing significant errors for a 20,000-word WSJ dictation task.

### 13.3. N-BEST AND MULTIPASS SEARCH STRATEGIES

Ideally, a search algorithm should consider all possible hypotheses based on a unified probabilistic framework that integrates all *knowledge sources* (KSs).<sup>5</sup> These KSs, such as acoustic models, language models, and lexical pronunciation models, can be integrated in an HMM state search framework. It is desirable to use the most detailed models, such as context-dependent models, interword context-dependent models, and high-order  $n$ -grams, in the search as early as possible. When the explored search space becomes unmanageable, due to the increasing size of vocabulary or highly sophisticated KSs, search might be infeasible to implement.

As we develop more powerful techniques, the complexity of models tends to increase dramatically. For example, language understanding models in Chapter 17 require long-distance relationships. In addition, many of these techniques are not operating in a left-to-right manner. A possible alternative is to perform a multipass search and apply several KSs at different stages, in the proper order to constrain the search progressively. In the initial pass, the most discriminant and computationally affordable KSs are used to reduce the number of hypotheses. In subsequent passes, progressively reduced sets of hypotheses are examined, and more powerful and expensive KSs are then used until the optimal solution is found.

The early passes of multipass search can be considered fast match that eliminates those unlikely hypotheses. Multipass search is, in general, not admissible because the optimal word sequence could be wrongly pruned prematurely, due to the fact that not all KSs are used in the earlier passes. However, for complicated tasks, the benefits of computation complexity reduction usually outweigh the nonadmissibility. In practice, multipass search strategy using progressive KSs could generate better results than a search algorithm forced to use less powerful models due to computation and memory constraints.

The most straightforward multipass search strategy is the so-called  $n$ -best search paradigm. The idea is to use affordable KSs to first produce a list of  $n$  most probable word sequences in a reasonable time. Then these  $n$  hypotheses are rescored using more detailed models to obtain the most likely word sequence. The idea of the  $n$ -best list can be further extended to create a more compact hypotheses representation—namely word lattice or graph. A word lattice is a more efficient way to represent alternative hypotheses.  $N$ -best or lattice search is used for many large-vocabulary continuous speech recognition systems [20, 30, 44].

In this section we describe the representation of the  $n$ -best list and word lattice. Several algorithms to generate such an  $n$ -best-list or word lattice are discussed.

---

<sup>5</sup> In the field of artificial intelligence, the process of performing search through an integrated network of various knowledge sources is called *constraint satisfaction*.

### 13.3.1. *N*-best Lists and Word Lattices

Table 13.4 shows an example *n*-best (10-best) list generated for a North American Business (NAB) sentence. *N*-best search framework is effective only for *n* of the order of tens or hundreds. If the short *n*-best list that is generated by using less optimal models does not include the correct word sequence, the successive rescoring phases have no chance to generate the correct answer. Moreover, in a typical *n*-best list like the one shown in Table 13.4, many of the different word sequences are just one-word variations of each other. This is not surprising, since similar word sequences should achieve similar scores. In general, the number of *n*-best hypotheses might grow exponentially with the length of the utterance. Word lattices and word graphs are thus introduced to replace *n*-best list with a more compact representation of alternative hypotheses.

*Word lattices* are composed by word hypotheses. Each word hypothesis is associated with a score and an explicit time interval. Figure 13.8 shows an example of a word lattice corresponding to the *n*-best list example in Table 13.4. It is clear that a word lattice is more efficient representation. For example, suppose the spoken utterance contains 10 words and there are 2 different word hypotheses for each word position. The *n*-best list would need to have  $2^{10} = 1024$  different sentences to include all the possible permutations, whereas the word lattice requires only 20 different word hypotheses.

*Word graphs*, on the other hand, resemble finite state automata, in which arcs are labeled with words. Temporal constraints between words are implicitly embedded in the topology. Figure 13.9 shows a word graph corresponding to the *n*-best list example in Table 13.4. Word graphs in general have an explicit specification of word connections that don't allow overlaps or gaps along the time axis. Nonetheless, word lattices and graphs are similar, and we often use these terms interchangeably.<sup>6</sup> Since an *n*-best list can be treated as a simple word lattice, word lattices are a more general representation of alternative hypotheses. *N*-best lists or word lattices are generally evaluated on the following two parameters:

**Table 13.4** An example 10-best list for a North American Business sentence.

1.	I will tell you would I think in my office
2.	I will tell you what I think in my office
3.	I will tell you when I think in my office
4.	I would sell you would I think in my office
5.	I would sell you what I think in my office
6.	I would sell you when I think in my office
7.	I will tell you would I think in my office
8.	I will tell you why I think in my office
9.	I will tell you what I think on my office
10.	I Wilson you I think on my office

<sup>6</sup> We will use the term *word lattice* in the rest of this chapter..

- *Density*: In the  $n$ -best case, it is measured by how many alternative word sequences are kept in the  $n$ -best list. In the word lattice case, it is measured by the number of word hypotheses or word arcs per uttered word. Obviously, we want the density to be as small as possible for successive rescoring modules, provided the correct word sequence is included in the  $n$ -best list or word lattice.
- *The lower bound word error rate*: It is the lowest word error rate for any word sequence in the  $n$ -best list or the word lattice.

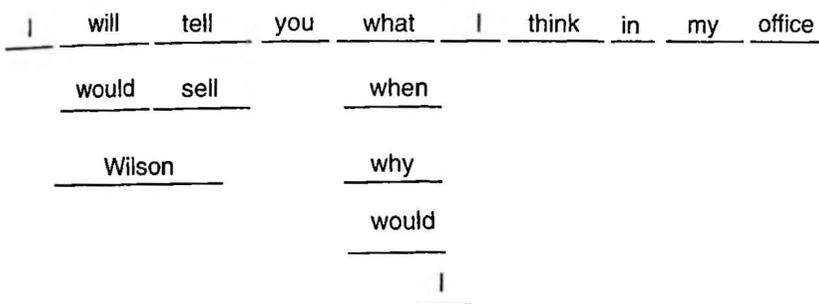


Figure 13.8 A word lattice example. Each word has an explicit time interval associated with it.

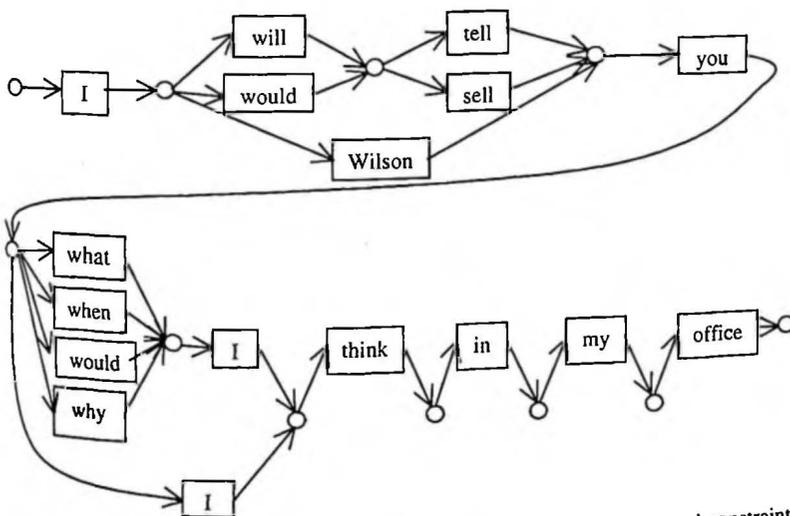
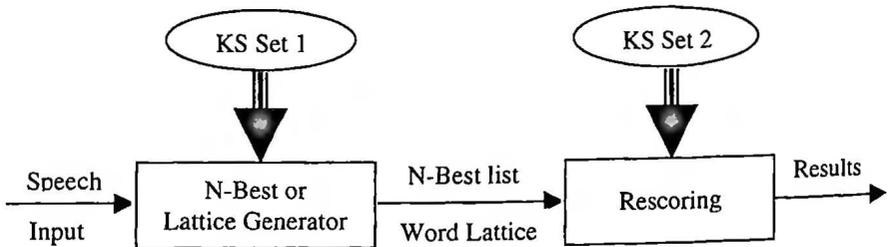


Figure 13.9 A word graph example for the  $n$ -best list in Table 13.4. Temporal constraints are implicit in the topology.

Rescoring with highly similar  $n$ -best alternatives duplicates computation on common parts. The compact representation of word lattices allows both data structure and computation sharing of the common parts among similar alternative hypotheses, so it is generally computationally less expensive to rescore the word lattice.

Figure 13.10 illustrates the general  $n$ -best/lattice search framework. Those KSs providing most constraints, at a lesser cost, are used first to generate the  $n$ -best list or word lattice. The  $n$ -best list or word lattice is then passed to the rescoring module, which uses the remaining KSs to select the optimal path. You should note that the  $n$ -best and word-lattice generators sometimes involve several phases of search mechanisms to generate the  $n$ -best list or word lattice. Therefore, the whole search framework in Figure 13.10 could involve several ( $> 2$ ) phases of search mechanism.

Does the compact  $n$ -best or word-lattice representation impose constraints on the complexity of the acoustic and language models applied during successive rescoring modules? The word lattice can be expanded for higher-order language models and detailed context-dependent models, like inter-word triphone models. For example, to use higher-order language models for word lattice entails copying each word in the appropriate context of preceding words (in the trigram case, the two immediately preceding words). To use inter-word triphone models entails replacing the triphones for the beginning and ending phone of each word with appropriate interword triphones. The expanded lattice can then be used with detailed acoustic and language models. For example, Murveit et al. [30] report this can achieve trigram search without exploring the enormous trigram search space.



**Figure 13.10**  $N$ -best/lattice search framework. The most discriminant and inexpensive knowledge sources (KSs 1) are used first to generate the  $n$ -best/lattice. The remaining knowledge sources (KSs 2, usually expensive to apply) are used in the rescoring phase to pick up the optimal solution [40].

### 13.3.2. The Exact $N$ -best Algorithm

Stack decoding is the choice of generating  $n$ -best candidates because of its best-first principle. We can keep it generating results until it finds  $n$  complete paths; these  $n$  complete sentences form the  $n$ -best list. However, this algorithm usually cannot generate the  $n$  best candidates efficiently. The efficient  $n$ -best algorithm for time-synchronous Viterbi search was first introduced by Schwartz and Chow [39]. It is a simple extension of time-synchronous Viterbi search. The fundamental idea is to maintain separate records for paths

with distinct histories. The history is defined as the whole word sequence up to the current time  $t$  and word  $w$ . This exact  $n$ -best algorithm is also called *sentence-dependent  $n$ -best algorithm*. When two or more paths come to the same state at the same time, paths having the same history are merged and their probabilities are summed together; otherwise, only the  $n$ -best paths are retained for each state. As commonly used in speech recognition, a typical HMM state has 2 or 3 predecessor states within the word HMM. Thus, for each time frame and each state, the  $n$ -best search algorithm needs to compare and merge 2 or 3 sets of  $n$  paths into  $n$  new paths. At the end of the search, the  $n$  paths in the final state of the trellis are simply re-ordered to obtain the  $n$ -best word sequences.

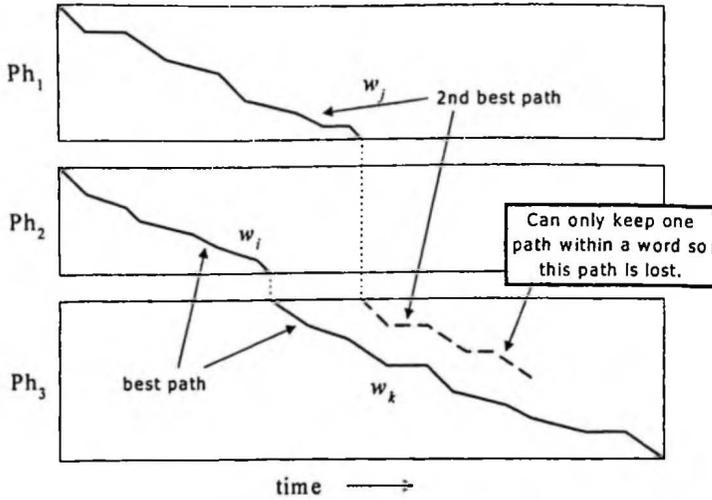
This straightforward  $n$ -best algorithm can be proved to be admissible<sup>7</sup> in normal circumstances [40]. The complexity of the algorithm is proportional to  $O(n)$ , where  $n$  is the number of paths kept at each state. This is often too slow for practical systems.

### 13.3.3. Word-Dependent $N$ -best and Word-Lattice Algorithm

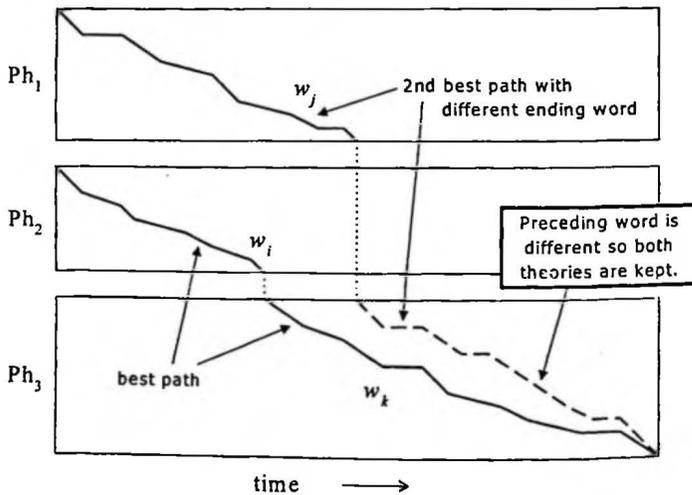
Since many of the different entries in the  $n$ -best list are just one-word variations of each other, as shown in Table 13.4, one efficient algorithm can be derived from the normal 1-best Viterbi algorithm to generate the  $n$ -best hypotheses. The algorithm runs just like the normal time-synchronous Viterbi algorithm for all within-word transitions. However for each time frame  $t$ , and each word-ending state, the algorithm stores all the different words that can end at current time  $t$  and their corresponding scores in a *traceback* list. At the same time, the score of the best hypothesis at each grammar state is passed forward, as in the normal time-synchronous Viterbi search. This obviously requires almost no extra computation above the normal time-synchronous Viterbi search. At the end of search, you can simply search through the stored traceback list to get all the permutations of word sequences with their corresponding scores. If you use a simple threshold, the traceback can be implemented very efficiently to only uncover the word sequences with accumulated cost scores below the threshold. This algorithm is often referred as *traceback-based  $n$ -best algorithm* [29, 42] because of the use of the traceback list in the algorithm.

However, there is a serious problem associated with this algorithm. It could easily miss some low-cost hypotheses. Figure 13.11 illustrates an example in which word  $w_k$  can be preceded by two different words  $w_i$  and  $w_j$  in different time frames. Assuming path  $w_i - w_k$  has a lower cost than path  $w_j - w_k$  when both paths meet during the trellis search of  $w_k$ , the path  $w_j - w_k$  will be pruned away. During traceback for finding the  $n$ -best word sequences, there is only one best starting time for word  $w_k$ , determined by the best boundary between the best preceding word  $w_i$  and it. Even though path  $w_j - w_k$  might have a very low cost (let's say only marginally higher than that of  $w_i - w_k$ ), it could be completely overlooked, since the path has a different starting time for word  $w_k$ .

<sup>7</sup> Although one can show, in the worst case, when paths with different histories have near identical scores for each state, the search actually needs to keep all paths ( $> N$ ) in order to guarantee absolute admissibility. Under this worst case, the admissible algorithm is clearly exponential in the number of words for the utterance, since all permutations of word sequences for the whole sentence need to be kept.



**Figure 13.11** Deficiency in traceback-based  $n$ -best algorithm. The best subpath,  $w_j - w_k$ , will prune away subpath  $w_j - w_i$  while searching the word  $w_k$ ; the second-best subpath cannot be recovered [40].



**Figure 13.12** Word-dependent  $n$ -best algorithm. Both subpaths  $w_i - w_k$  and  $w_j - w_k$  are kept under the word-dependent assumption [40].

The *word-dependent*  $n$ -best algorithm [38] can alleviate the deficiency of the traceback-based  $n$ -best algorithm, in which only one starting time is kept for each word, so the starting time is independent of the preceding words. On the other hand, in the sentence-dependent  $n$ -best algorithm, the starting time for a word depends on all the preceding words, since different histories are kept separately. A good compromise is the so-called word-dependent assumption: *The starting time of a word depends only on the immediate preceding word.* That is, given a word pair and its ending time, the boundary between these two words is independent of further predecessor words.

In the word-dependent assumption, the history to be considered for a different path is no longer the entire word sequence; instead, it is only the immediately preceding word. This allows you to keep  $k$  ( $\ll n$ ) different records for each state and each time frame in Viterbi search. Differing slightly from the exact  $n$ -best algorithm, a traceback must be performed to find the  $n$ -best list at the end of search. The algorithm is illustrated in Figure 13.12. A word-dependent  $n$ -best algorithm has a time complexity proportional to  $k$ . However, it is no longer admissible because of the word-dependent approximation. In general, this approximation is quite reasonable if the preceding word is long. The loss it entails is insignificant [6].

### 13.3.3.1. One-Pass $N$ -best and Word-Lattice Algorithm

As presented in Section 13.1, one-pass Viterbi beam search can be implemented very efficiently using a tree lexicon. Section 13.1.2 states that multiple copies of lexical trees are necessary for incorporating language models other than the unigram. When bigram is used in lexical tree search, the successor lexical tree is predecessor-dependent. This predecessor-dependent property immediately translates into the word-dependent property,<sup>8</sup> as defined in Section 13.3.3, because the starting time of a word clearly depends on the immediately preceding word. This means that different word-dependent partial paths are automatically saved under the framework of predecessor-dependent successor trees. Therefore, one-pass predecessor-dependent lexical tree search can be modified slightly to output  $n$ -best lists or word graphs.

Ney et al. [31] used a word graph builder with a one-pass predecessor-dependent lexical tree search. The idea is to exploit the word-dependent property inherited from the predecessor-dependent lexical tree search. During predecessor-dependent lexical tree search, two additional quantities are saved whenever a word ending state is processed.

$\tau(t; w_i, w_j)$ —Representing the optimal word boundary between word  $w_i$  and  $w_j$ , given word  $w_j$  ending at time  $t$ .

$h(w_j; \tau(t; w_i, w_j), t)$ —Representing the cumulative cost that word  $w_j$  produces acoustic vector  $\mathbf{x}_\tau, \mathbf{x}_{\tau+1}, \dots, \mathbf{x}_t$ .

<sup>8</sup> When higher order  $n$ -gram models are used, the boundary dependence will be even more significant. For example, when trigrams are used, the boundary for a word juncture depends on the previous two words. Since we generally want a fast method of generating word lattices/graphs, bigram is often used instead of higher order  $n$ -gram to generate word lattices/graphs.

At the end of the utterance, the word lattice or  $n$ -best list is constructed by tracing back all the permutations of word pairs recorded during the search. The algorithm is summarized in Algorithm 13.3.

**ALGORITHM 13.3: ONE-PASS PREDECESSOR-DEPENDENT LEXICAL TREE SEARCH FOR  $N$ -BEST OR WORD-LATTICE CONSTRUCTION**

**Step 1:** For  $t = 1..T$ ,

1-best predecessor-dependent lexical tree search;

$\forall (w_i, w_j)$  ending at  $t$

record word-dependent crossing time  $\tau(t; w_i, w_j)$ ;

record cumulative word score  $h(w_j; \tau(t; w_i, w_j), t)$ ;

**Step 2:** Output 1-best result;

**Step 3:** Construct  $n$ -best or word-lattice by tracing back the word-pair records ( $\tau$  and  $h$ ).

### 13.3.4. The Forward-Backward Search Algorithm

As described Chapter 12, the ability to predict how well the search fares in the future for the remaining portion of the speech helps to reduce the search effort significantly. The one-pass search strategy, in general, has very little chance of predicting the cost for the portion that it has not seen. This difficulty can be alleviated by multipass search strategies. In successive phases the search should be able to provide good estimates for the remaining paths, since the entire utterance has been examined by the earlier passes. In this section we investigate a special type of multipass search strategy—forward-backward search.

The idea is to first perform a forward search, during which partial forward scores  $\alpha$  for each state can be stored. Then perform a second pass search backward—that is, the second pass starts by taking the final frame of speech and searches its way back until it reaches the start of the speech. During the backward search, the partial forward scores  $\alpha$  can be used as an accurate estimate of the heuristic function or the fast match score for the remaining path. Even though different KSs might be used in forward and backward phases, this estimate is usually close to perfect, so the search effort for the backward phase can be significantly reduced.

The forward search must be very fast and is generally a time-synchronous Viterbi search. As in the multipass search strategy, simplified acoustic and language models are often used in forward search. For backward search, either time-synchronous search or time-asynchronous A\* search can be employed to find the  $n$ -best word sequences or word lattice.

### 13.3.4.1. Forward-Backward Search

Stack decoding, as described in Chapter 12, is based on the admissible A\* search, so the first complete hypothesis found with a cost below that of all the hypotheses in the stack is guaranteed to be the best word sequence. It is straightforward to extend stack decoding to produce the  $n$ -best hypotheses by continuing to extend the partial hypotheses according to the same A\* criterion until  $n$  different hypotheses are found. These  $n$  different hypotheses are destined to be the  $n$ -best hypotheses under a proof similar to that presented in Chapter 12. Therefore, stack decoding is a natural choice for producing the  $n$ -best hypotheses.

However, as described in Chapter 12, the difficulty of finding a good heuristic function that can accurately under-estimate the remaining path has limited the use of stack decoding. Fortunately, this difficulty can be alleviated by *tree-trellis forward-backward search* algorithms [41]. First, the search performs a time-synchronous forward search. At each time frame  $t$ , it records the score of the final state of each word ending. The set of words whose final states are active (surviving in the beam) at time  $t$  is denoted as  $\Delta_t$ . The score of the final state of each word  $w$  in  $\Delta_t$  is denoted as  $\alpha_t(w)$ , which represents the sum of the cost of matching the utterance up to time  $t$  given the most likely word sequence ending with word  $w$  and the cost of the language model score for that word sequence. At the end of the forward search, the best cost is obtained and denoted as  $\alpha^T$ .

After the forward pass is completed, the second search is run in reverse (backward), i.e., considering the last frame  $T$  as the beginning one and the first frame as the final one. Both the acoustic models and language models need to be reversed. The backward search is based on A\* search. At each time frame  $t$ , the best path is removed from the stack and a list of possible one-word extensions for that path is generated. Suppose this best path at time  $t$  is  $ph_{w_i}$ , where  $w_i$  is the first word of this partial path (the last expanded during backward A\* search). The exit score of path  $ph_{w_i}$  at time  $t$ , which now corresponds to the score of the initial state of the word HMM  $w_j$ , is denoted as  $\beta_t(ph_{w_i})$ .

Let us now assume we are concerned about the one-word extension of word  $w_i$  for path  $ph_{w_i}$ . Remember that there are two fundamental issues for the implementation of A\* search algorithm—(1) finding an effective and efficient heuristic function for estimating the future remaining input feature stream and (2) finding the best crossing time between  $w_i$  and  $w_j$ .

The stored forward score  $\alpha$  can be used for solving both issues effectively and efficiently. For each time  $t$ , the sum  $\alpha_t(w_i) + \beta_t(ph_{w_i})$  represents the cost score of the best complete path including word  $w_i$  and partial path  $ph_{w_i}$ .  $\alpha_t(w_i)$  clearly represents a very good heuristic estimate of the remaining path from the start of the utterance until the end of the word  $w_i$ , because it is indeed the best score computed in the forward path for the same quantity. Moreover, the optimal crossing time  $t^*$  between  $w_i$  and  $w_j$  can be easily computed by the following equation:

$$t^* = \arg \min_i [\alpha_t(w_i) + \beta_t(ph_{w_i})] \quad (13.7)$$

Finally, the new path  $ph'$ , including the one-word ( $w_i$ ) extension, is inserted into the stack, ordered by the cost score  $\alpha_i(w_i) + \beta_i(ph_w)$ . The heuristic function (forward scores  $\alpha$ ) allows the backward A\* search to concentrate search on extending only a few truly promising paths.

As a matter of fact, if the same acoustic and language models are used in both the forward and backward search, this heuristic estimate (forward scores  $\alpha$ ) is indeed a perfect estimate of the best score the extended path will achieve. The first complete hypothesis generated by backward A\* search coincides with the best one found in the time-synchronous forward search and is truly the best hypothesis. Subsequent complete hypotheses correspond sequentially to the  $n$ -best list, as they are generated in increasing order of cost. Under this condition, the size of the stack in the backward A\* search need only be  $N$ . Since the estimate of future is exact, the  $(N+1)^{\text{th}}$  path in the stack has no chance to become part of the  $n$ -best list. Therefore, the backward search is executed very efficiently to obtain the  $n$ -best hypotheses without exploring many unpromising branches. Of course, tree-trellis forward-backward search can also be used like most other multipass search strategies—inexpensive KSs are used in the forward search to get an estimate of  $\alpha$ , and more expensive KSs are used in the backward A\* search to generate the  $n$ -best list.

The same idea of using forward score  $\alpha$  can be applied to time-synchronous Viterbi search in the backward search instead of backward A\* search [7, 34]. For large-vocabulary tasks, the backward search can run 2 to 3 orders of magnitude faster than a normal Viterbi beam search. To obtain the  $n$ -best list from time-synchronous forward-backward search, the backward search can also be implemented in a similar way as a time-synchronous word-dependent  $n$ -best search.

#### 13.3.4.2. Word-Lattice Generation

The forward-backward  $n$ -best search algorithm can be easily modified to generate word lattices instead of  $n$ -best lists. A forward time-synchronous Viterbi search is performed first to compute  $\alpha_i(w)$ , the score of each word  $w$  ending at time  $t$ . At the end of the search, this best score  $\alpha^T$  is also recorded to establish the global pruning threshold. Then, a backward time-synchronous Viterbi search is performed to compute  $\beta_i(w)$ , the score of each word  $w$  beginning at time  $t$ . To decide whether to include word juncture  $w_i - w_j$  in the word lattice/graph at time  $t$ , we can check whether the forward-backward score is below a global pruning threshold. Specifically, supposed bigram probability  $P(w_j | w_i)$  is used, if

$$\alpha_i(w_i) + \beta_i(w_j) + [-\log P(w_j | w_i)] < \alpha^T + \theta \quad (13.8)$$

where  $\theta$  is the pruning threshold, we will include  $w_i - w_j$  in the word lattice/graph at time  $t$ . Once word juncture  $w_i - w_j$  is kept, the search continues looking for the next word-pair, where the first word  $w_i$  will be the second word of the next word-pair.