

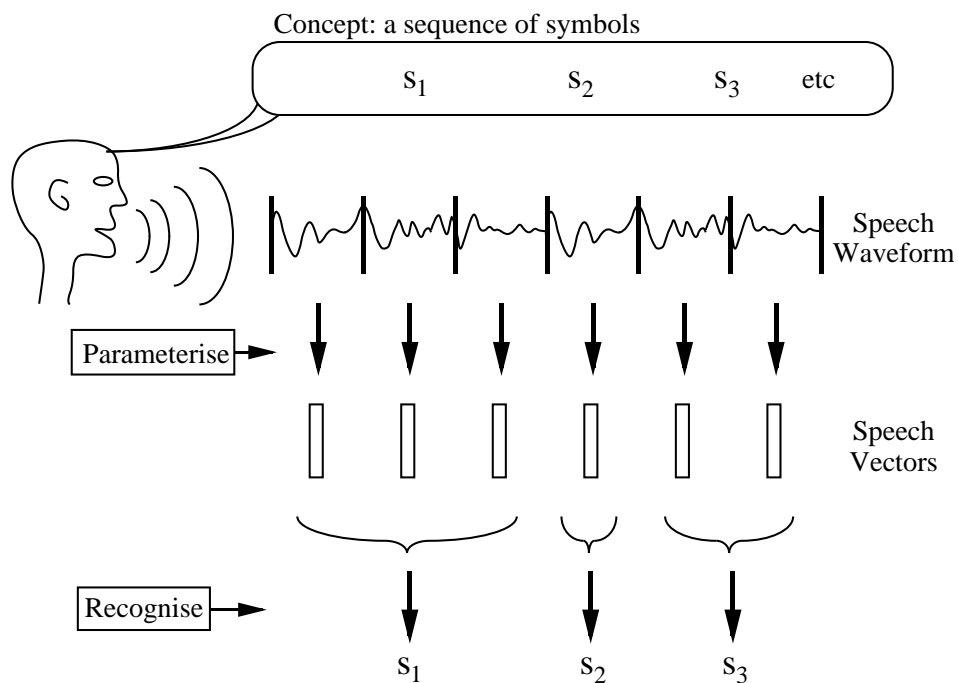
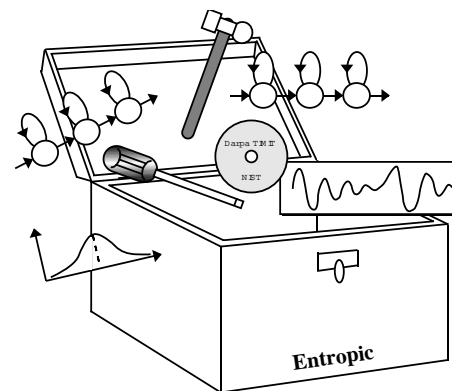
Grammars for speech recognition

Dr Philip Jackson

- Example application
 - Isolated digit recognition
 - Building the grammar
 - Training & testing
- Task grammars
 - Isolated word recognition
 - Connected word recognition
 - Null states and word networks
 - Continuous speech recognition

Isolated Word Recognition application

- Data preparation
- Training
- Testing
- Analysis



Message encoding and decoding.*

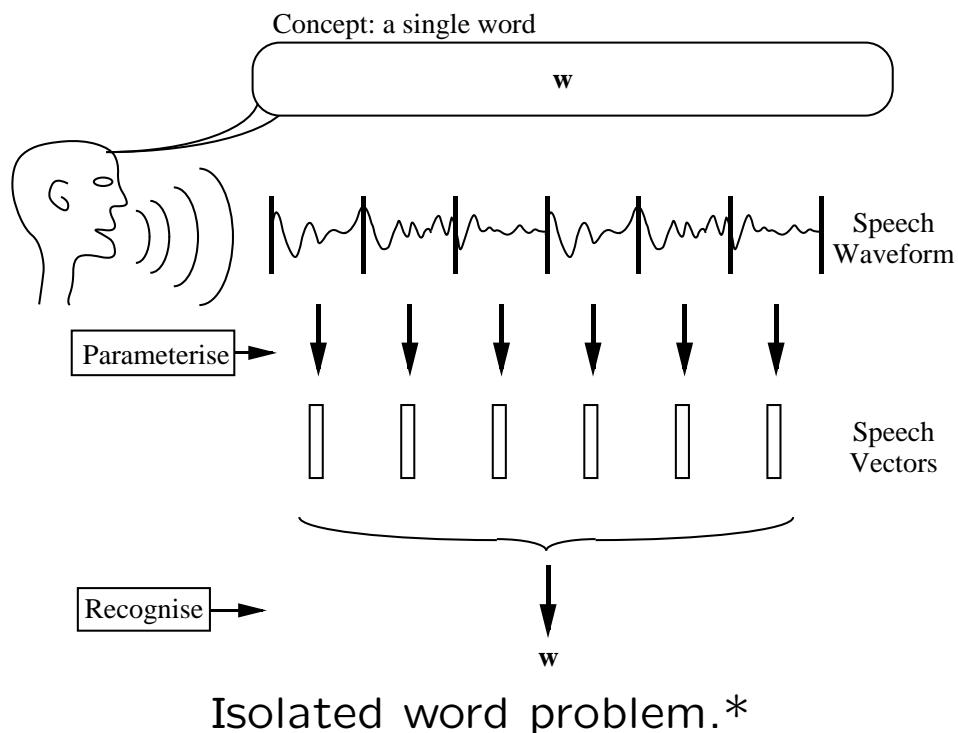
Isolated Word Recognition (IWR) task

The problem is to find

$$\hat{w} = \arg \max_i \{P(w_i | \mathcal{O})\} \quad (1)$$

where according to Bayes

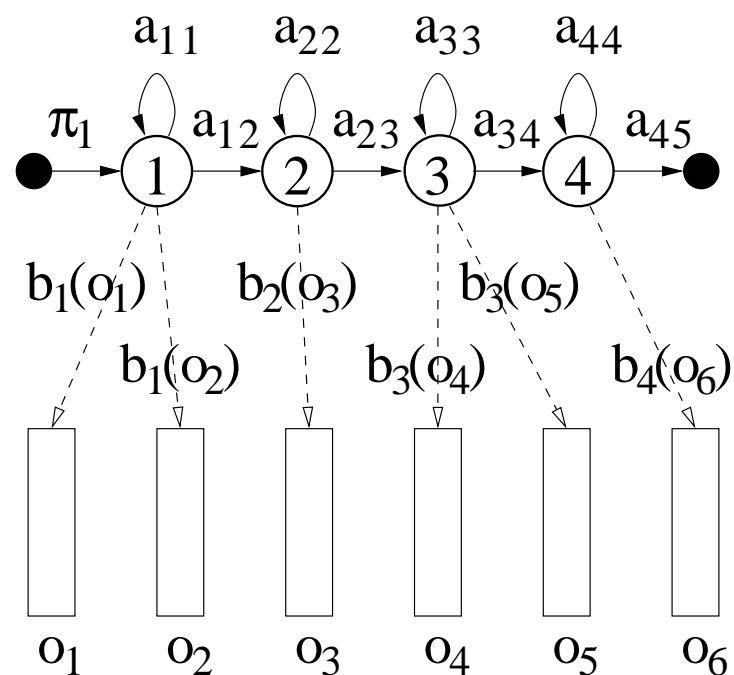
$$P(w_i | \mathcal{O}) = \frac{P(\mathcal{O} | w_i) P(w_i)}{P(\mathcal{O})} \quad (2)$$



The hidden Markov model

In this case, we assume

$$\begin{aligned} P(O|w_i) &\approx P(O|\lambda_i) \\ &\approx \max_X \left[\pi_{x_1} b_{x_1}(o_1) \left(\prod_{t=2}^T a_{x_{t-1}x_t} b_{x_t}(o_t) \right) \eta_{x_T} \right] \end{aligned} \quad (3)$$



The Markov generation model.

Null states in an HMM

Properties of non-emitting null states:

- determine beginning and end of a model
- do not generate any observations
- can be used to match a given model to an isolated test utterance
- useful for joining models together in a grammar
- transitions associated with null states can be modified to incorporate language model

IWR: Building the grammar

Example utterances:

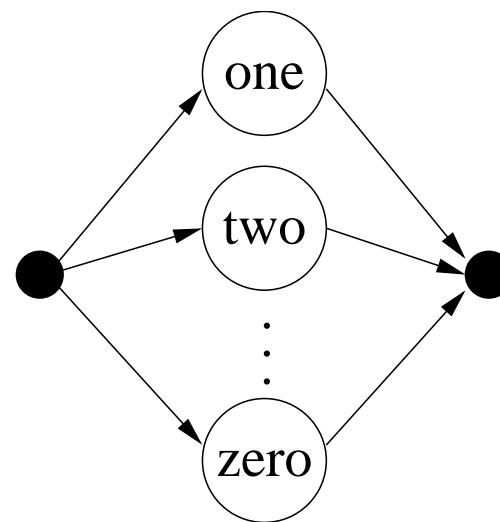
- eight
- oh
- six

Task grammar:

```
$digit = ONE | TWO | THREE |  
        FOUR | FIVE | SIX |  
        SEVEN | EIGHT | NINE |  
        OH | ZERO;  
( SENT-START $digit SENT-END )
```

Key:

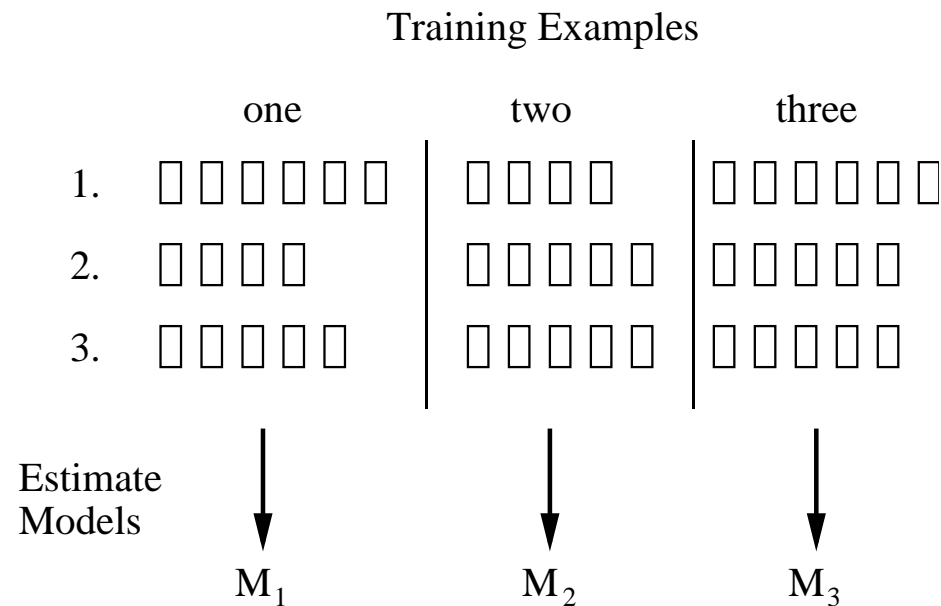
| alternatives



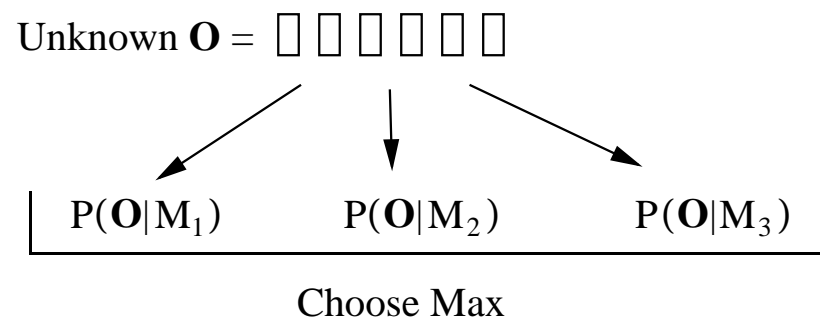
Grammar for isolated
digit recognition.

Training and test data

(a) Training



(b) Recognition



Using HMMs for isolated word recognition.*

Calculating recognition performance

Types of recognition error (e.g., for ground truth “A-B-C”):

- Substitution, S (e.g., “A-D-C”)
- Deletion, D (e.g., “A-C”)
- Insertion, I (e.g., “A-B-E-C”)

$$\% \text{ Correct} = 100 \times \frac{N - S - D}{N} \quad (4)$$

$$\% \text{ Accuracy} = 100 \times \frac{N - S - D - I}{N} \quad (5)$$

Task grammars

IWR: Binary word grammar

Example utterances:

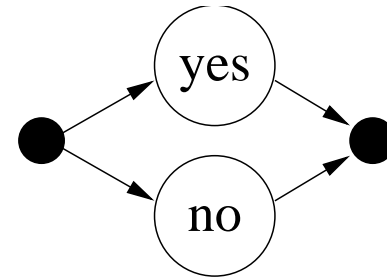
- yes
- no

Task grammar:

```
$answer = YES | NO;  
( SENT-START $answer SENT-END )
```

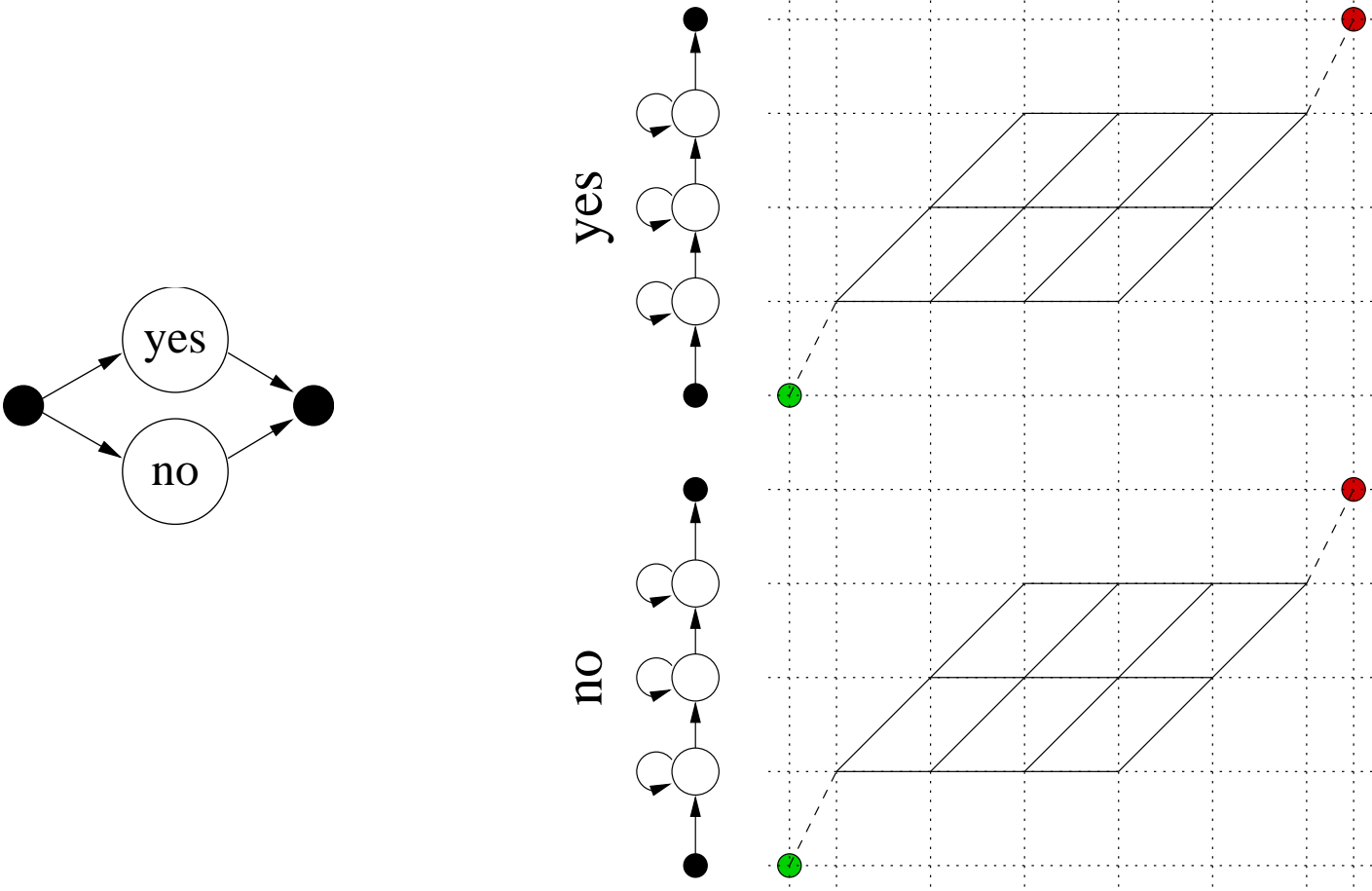
Key:

| alternatives



Grammar for isolated binary recognition.

IWR: binary word trellis



Grammar (left) and trellis (right) for a two-word IWR network.

IWR: Isolated digit grammar

Example utterances:

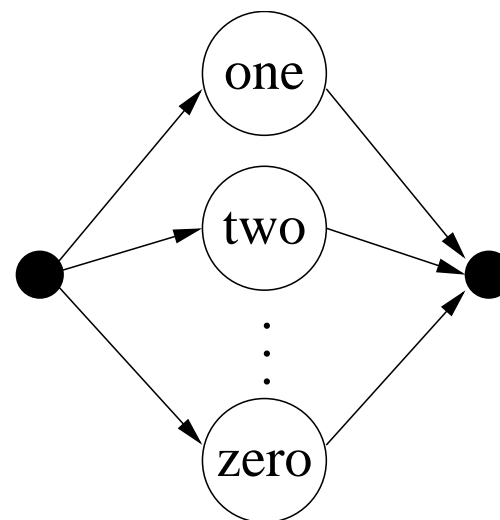
- eight
- oh
- six

Task grammar:

```
$digit = ONE | TWO | THREE |  
        FOUR | FIVE | SIX |  
        SEVEN | EIGHT | NINE |  
        OH | ZERO;  
( SENT-START $digit SENT-END )
```

Key:

| alternatives



Grammar for isolated
digit recognition.

CWR: Connected digit grammar

Example utterances:

- six eight six zero three one
- one oh one
- six oh four four

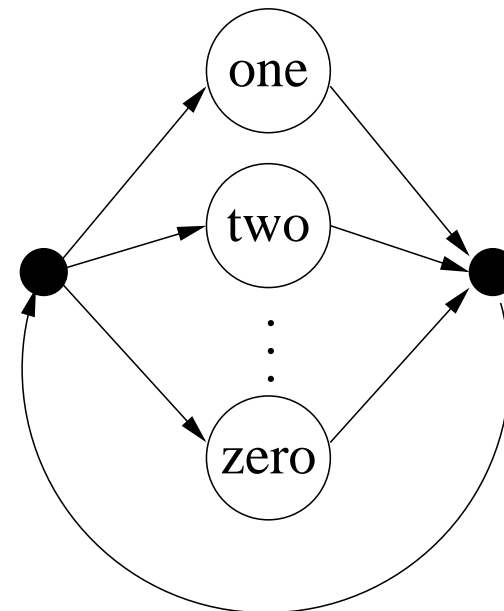
Task grammar:

```
$digit = ONE | TWO | THREE |  
        FOUR | FIVE | SIX |  
        SEVEN | EIGHT | NINE |  
        OH | ZERO;  
( SENT-START <$digit> SENT-END )
```

Key:

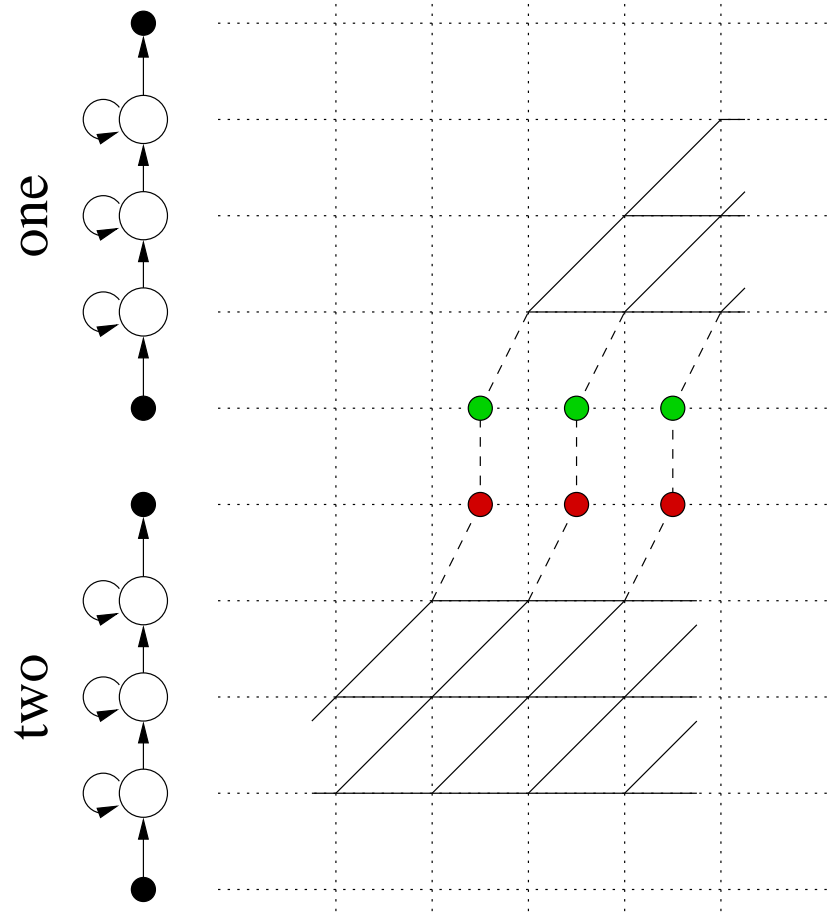
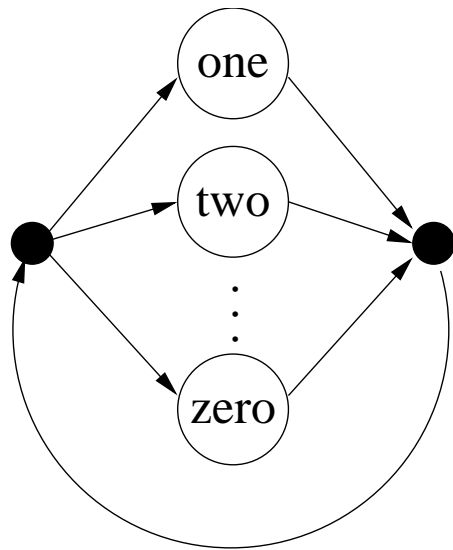
| alternatives

<.> one or more reps



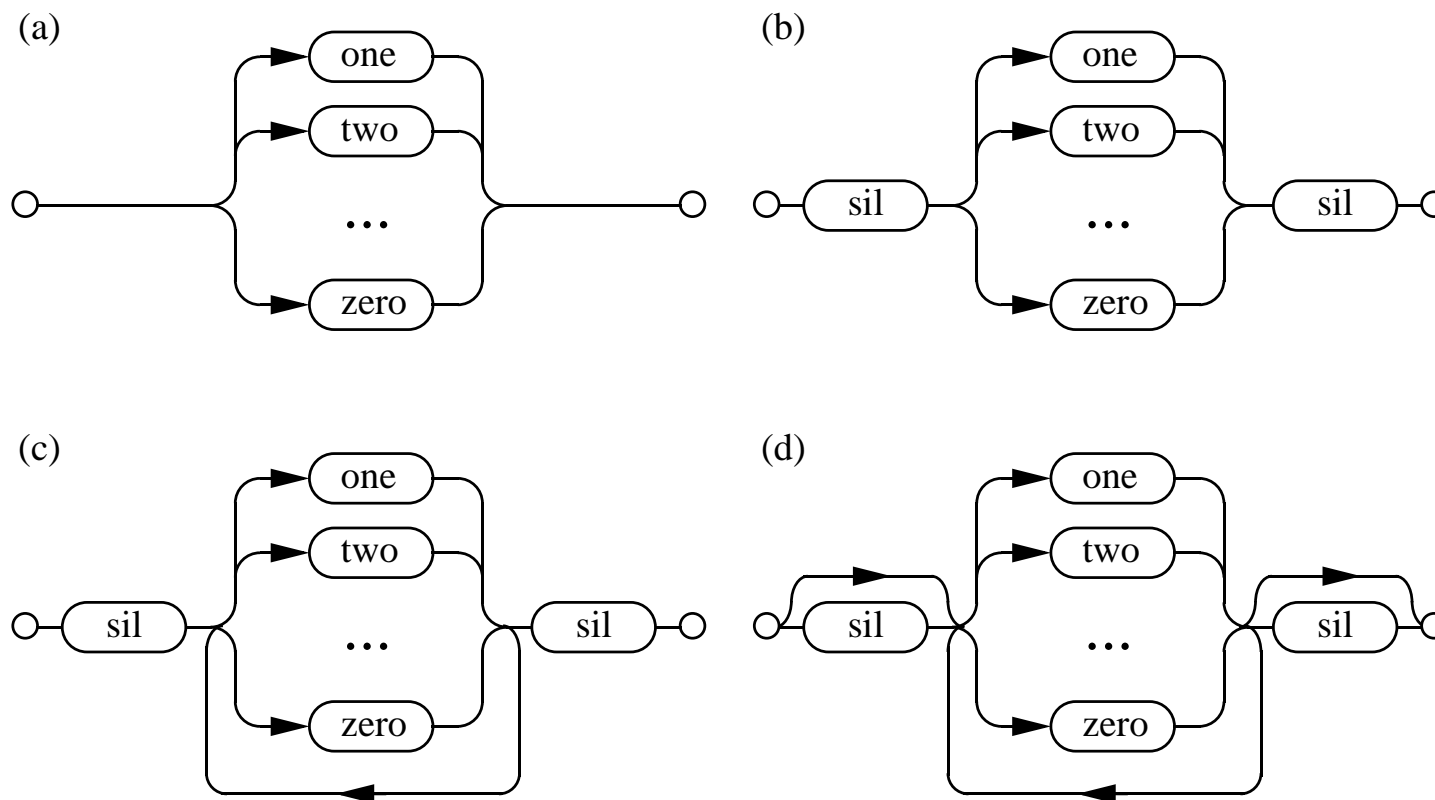
Grammar for connected digit recognition.

CWR: connected-digit trellis



Grammar (left) and trellis (right) for connected-digit recognition.

Isolated- & connected-digit grammars



Example networks for digit recognition tasks:* (a) IWR, (b) IWR with end-point adjustment using silence model, (c) CWR with silence model, (d) with optional silence.

CWR: Connected word grammar

Example utterances:

Dial three three two six five four

Phone Woodland

Call Steve Young

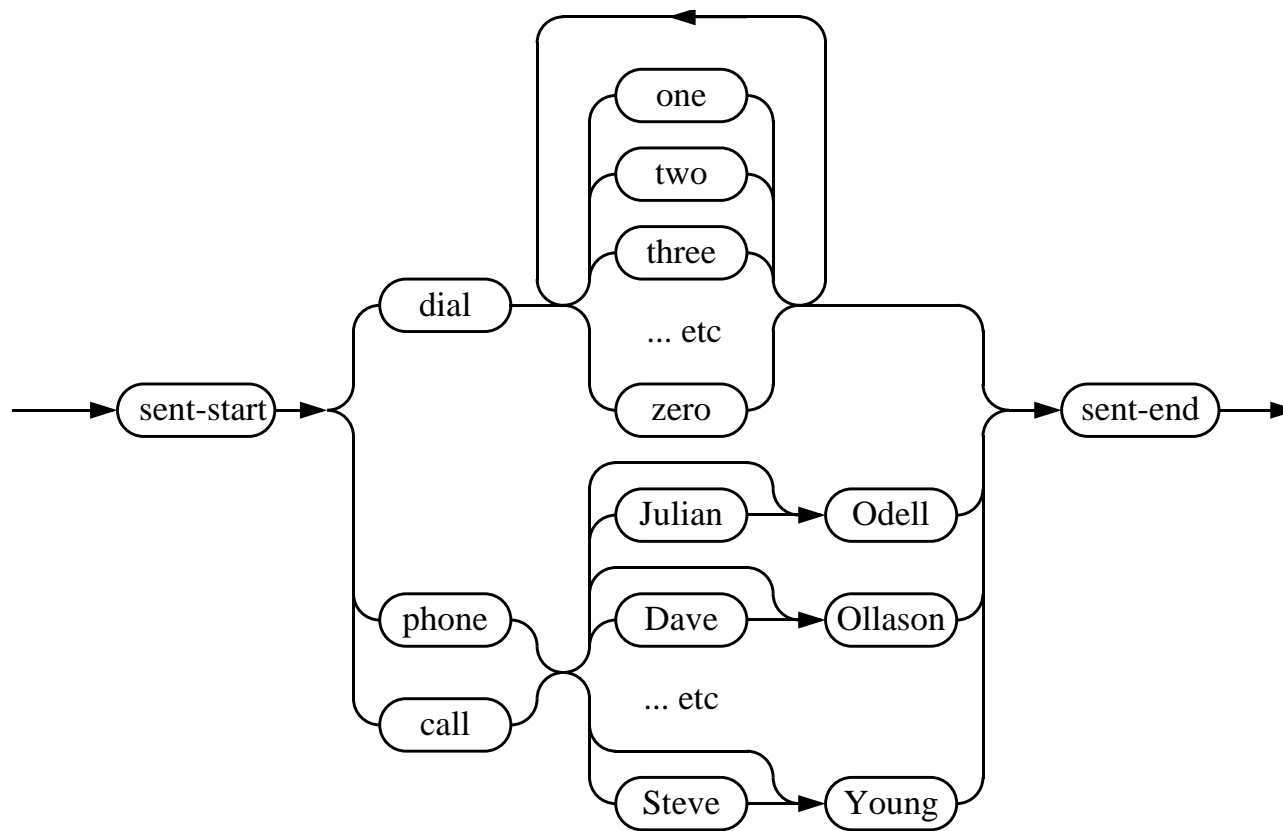
Task grammar:

```
$digit = ONE | TWO | THREE | FOUR | FIVE |  
        SIX | SEVEN | EIGHT | NINE | OH | ZERO;  
$name  = [ JULIAN ] ODELL |  
        [ DAVE ] OLLASON |  
        [ PHIL ] WOODLAND |  
        [ STEVE ] YOUNG;  
( SENT-START ( DIAL <$digit> | (PHONE|CALL) $name) SENT-END )
```

Key:

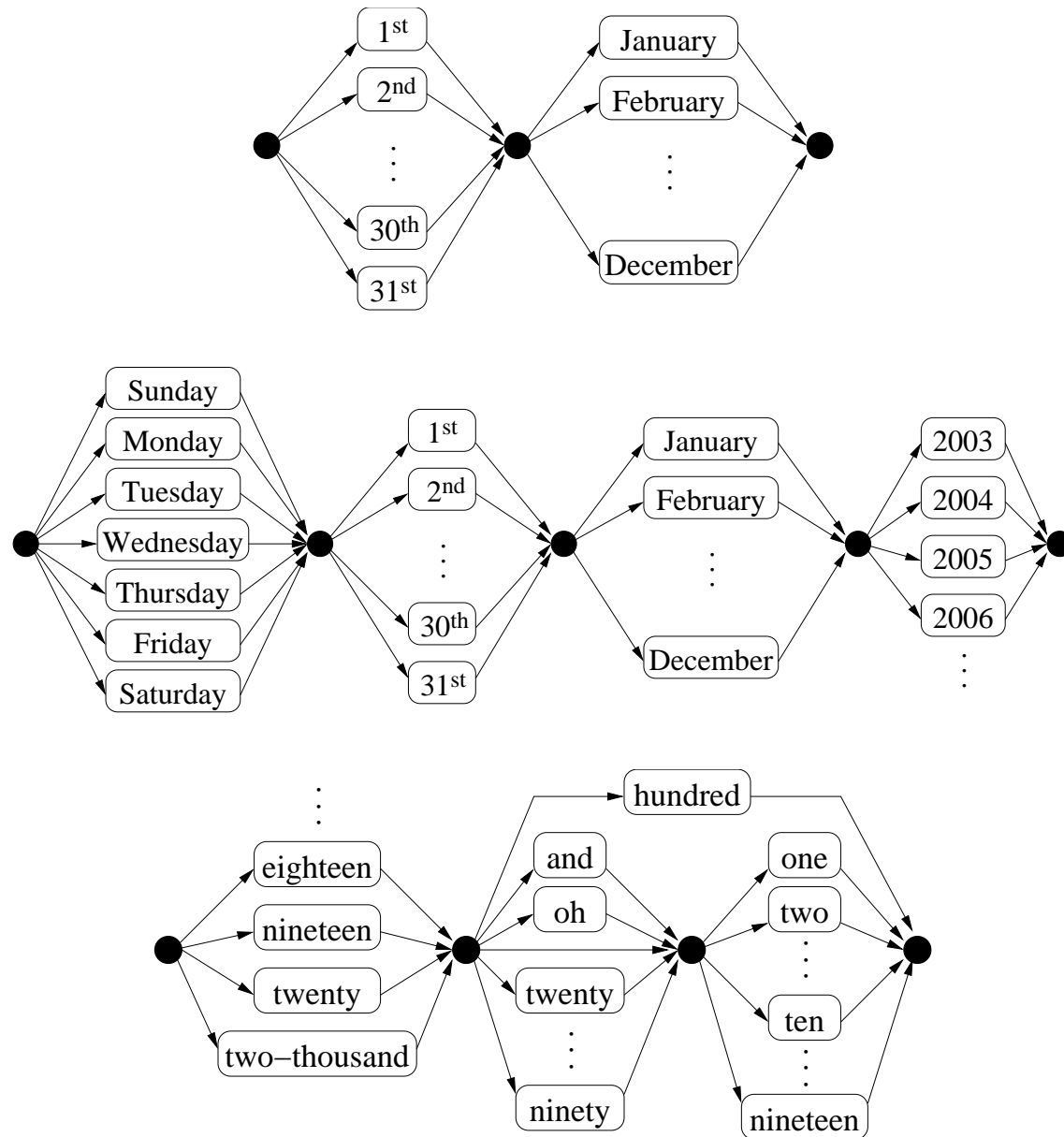
| alternatives, [.] optional, <.> one or more reps

CWR: Word network



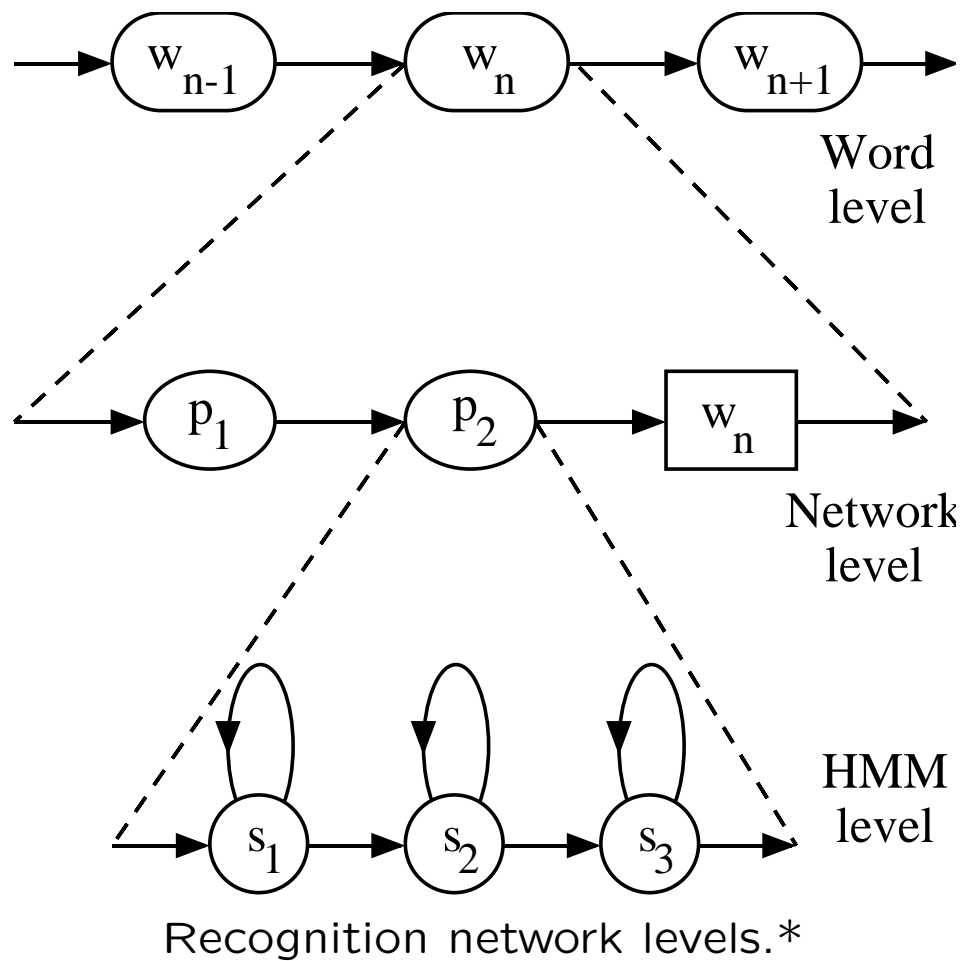
Grammar for voice dialling.*

CWR: date, day and year example grammars



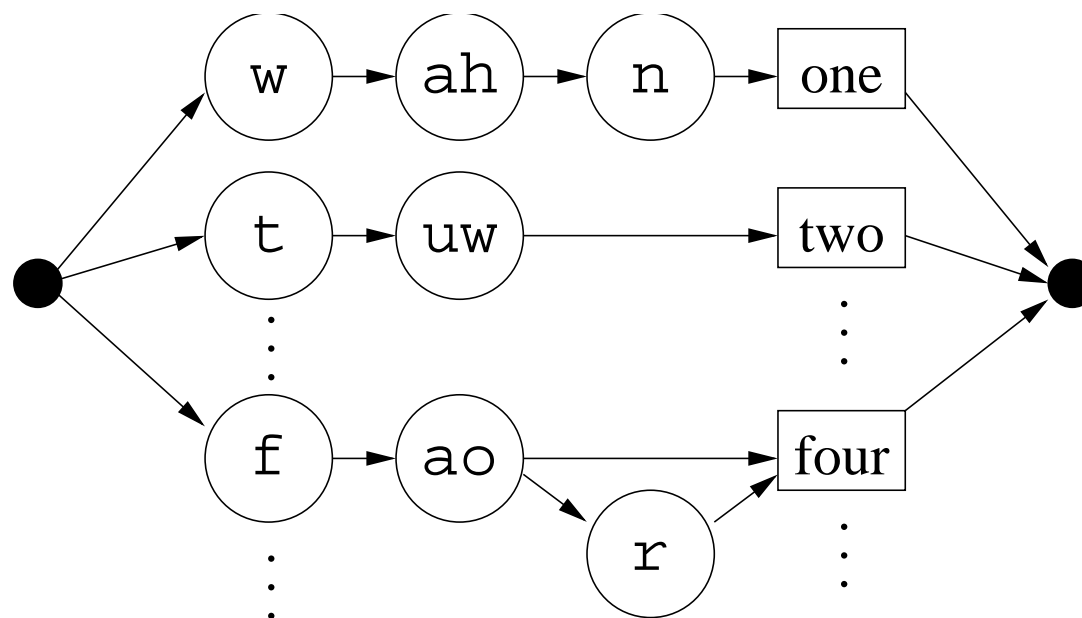
Example grammars (from top): date, day and year.

Hierarchy of recognition networks



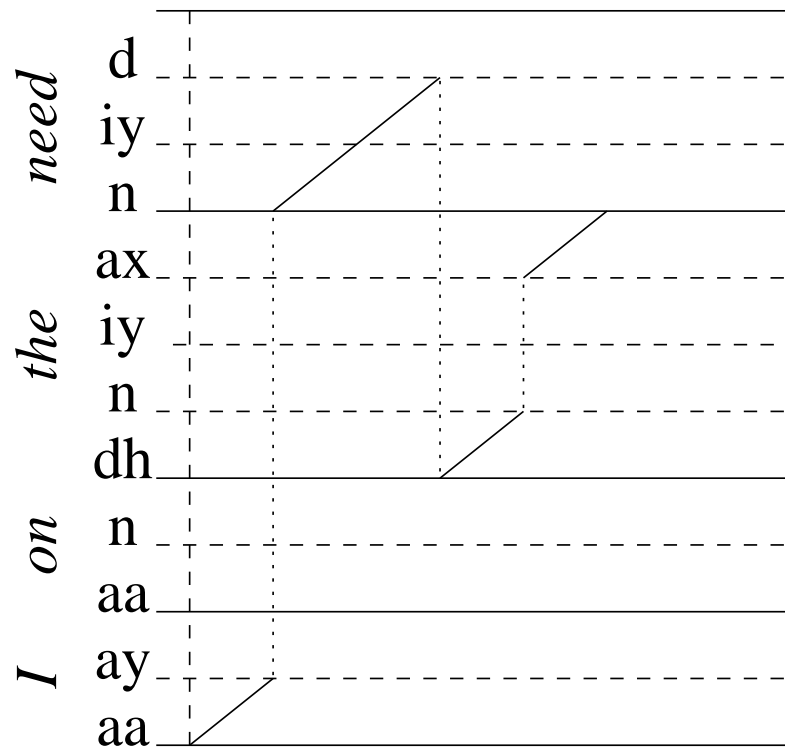
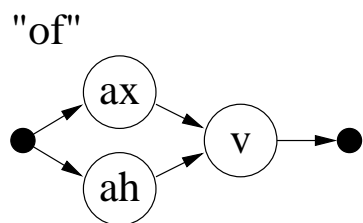
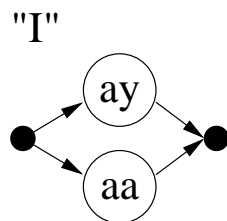
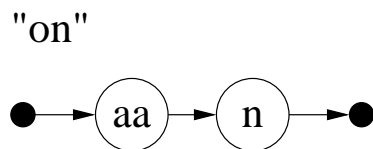
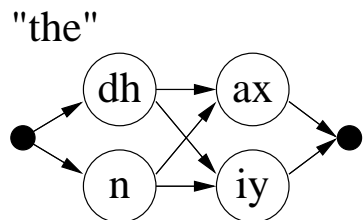
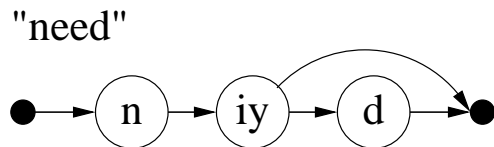
IWR: Phone-based digit dictionary

ONE	w ah n	SIX	s ih k s
TWO	t uw	SEVEN	s eh v n
THREE	th r iy	EIGHT	ey t
FOUR	f ao	NINE	n ay n
FOUR	f ao r	OH	ow
FIVE	f ay v	ZERO	z ia r ow



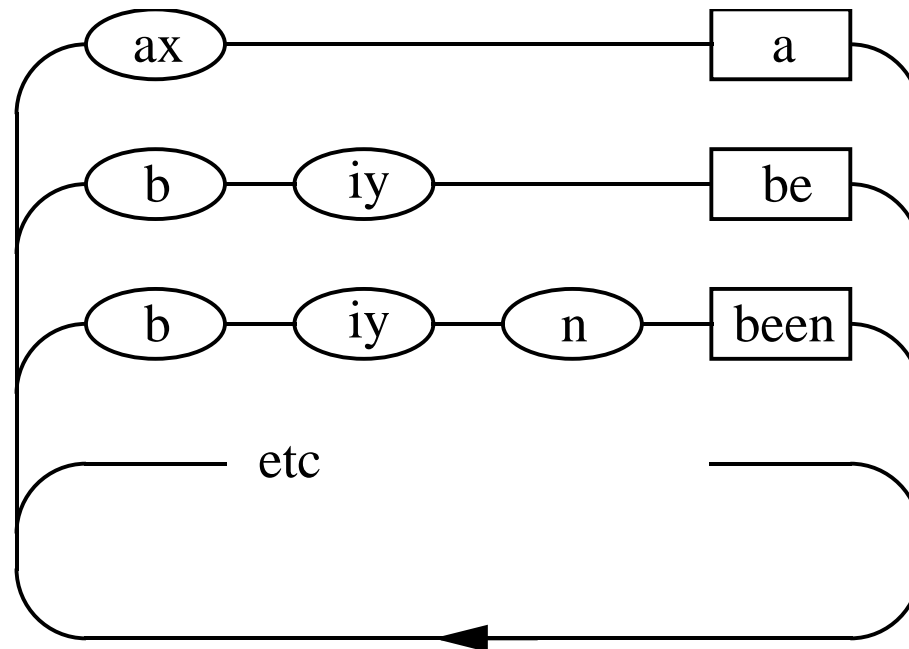
Grammar for phone-based isolated digit recognition.

CSR: continuous-speech grammar and trellis



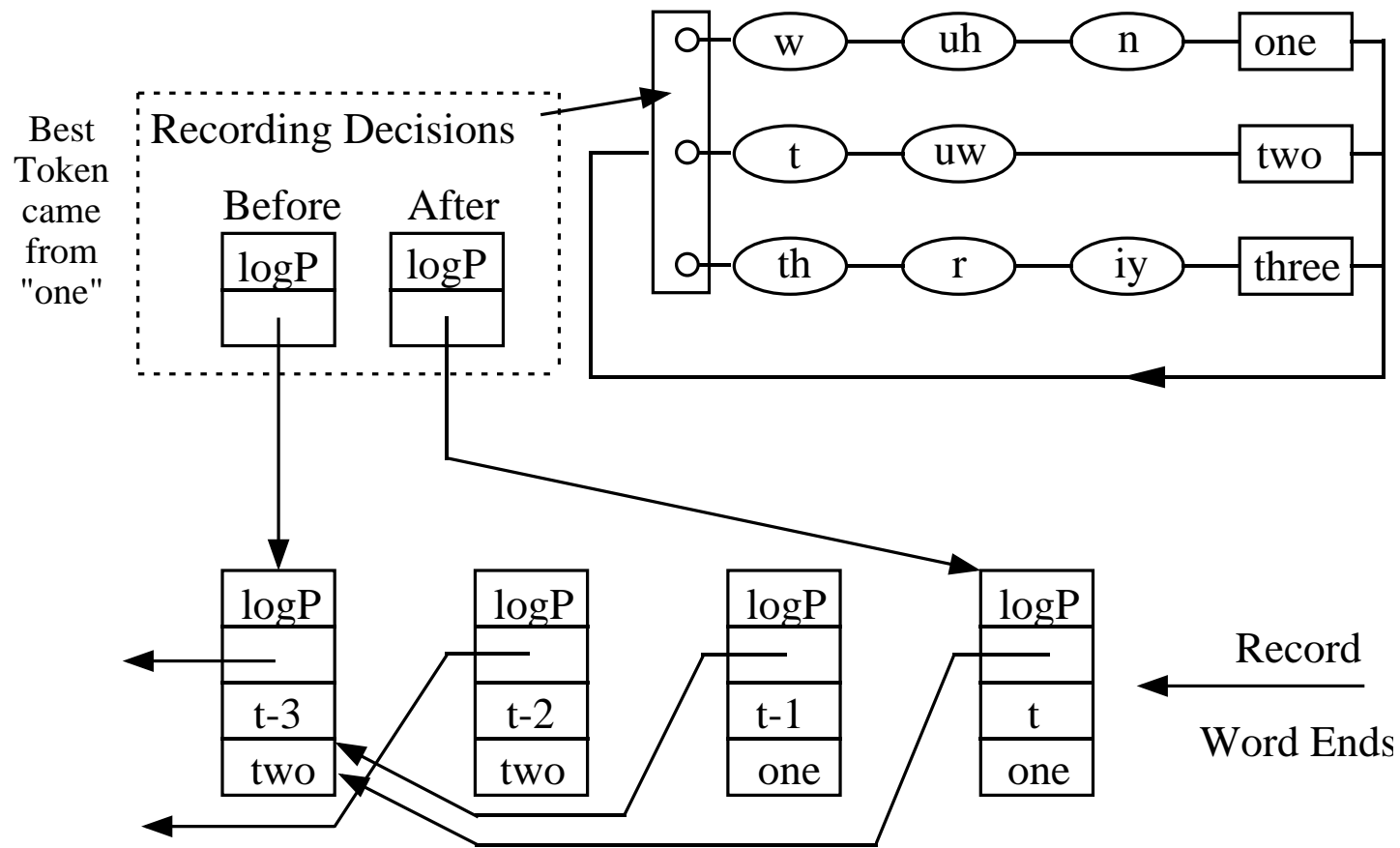
Networks (left) & trellis (right) for continuous-speech recognition.

CSR: word labels in the network



Word network for continuous speech recognition.*

CSR: recognition and traceback



Recording word-boundary decisions during continuous speech recognition.*

Grammar summary

- Isolated digit recognition
 - Null states, grammar & trellis diagrams
 - Training, testing & scoring
- Task grammars
 - Isolated word recognition (IWR)
 - Connected word recognition (CWR)
 - Context-free grammar
 - Hierarchy of HMMs
 - IWR using phone units
 - Continuous speech recognition (CSR)