

Ergodic-HMMを用いた確率付きネットワーク文法の自動獲得

・確率付きネットワーク文法



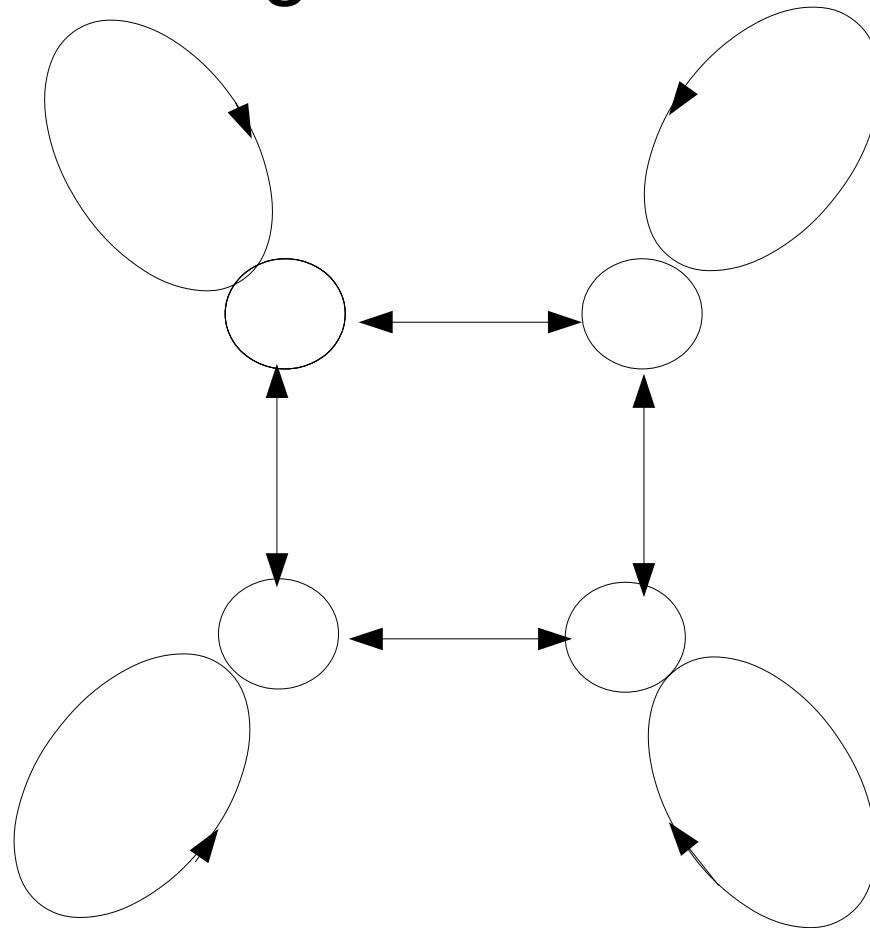
等値
(シンボル出力確率
= 単語の出力確率)

・Ergodic HMM

Ergodic HMM

言語データ

もしもし
はい そうです
えーっと ちょっと その
会議 の こと で ね
はい どうぞ
えーっと それ を ちょっ
と クレジットカード を
持っ て い な い 者 が
い る ん で す け れ ど も



Baum-Welch Algorithm

確率つきネットワーク文法

文の例

はい もしもし

えーっと そちら 第 1 回 の 通 訊 電 話 国 際 会 議 の 事 務 局 で
しょうか

はい そうです

えーっと ちょっと その 会 議 の こと でね

はい どうぞ

えーっと 今 手 元 に あの 登 録 用 紙 が ある ん です け れ ども

えーっと その 中 で ちょっと あの クレジットカード を ね

あ の クレジットカード の 名 前 と な ん か ナンバー を 書 く ところ
が ある ん です が

はい そうです

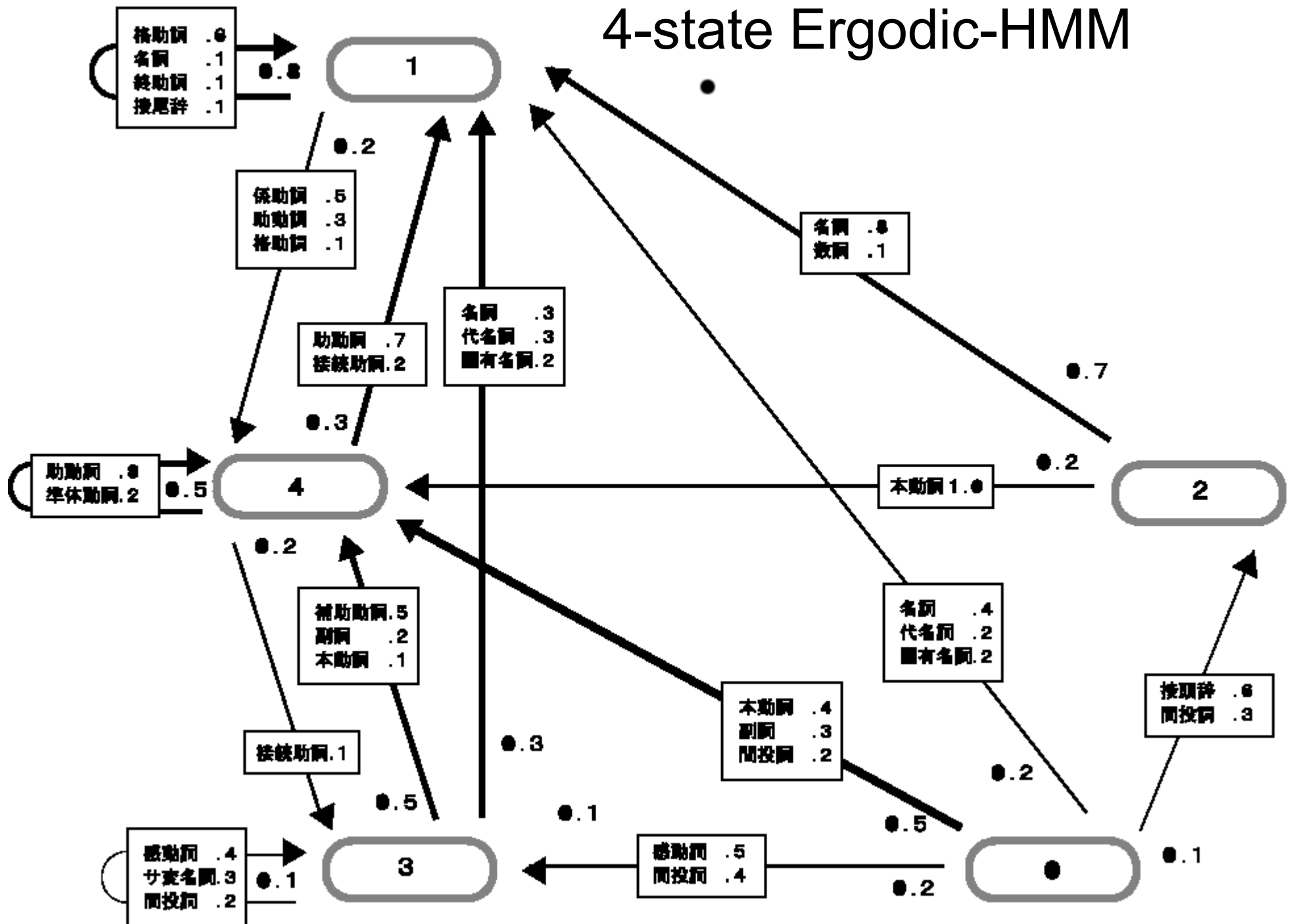
えーっと それ を ちょっと クレジットカード を 持 っ て い な い 者 が
い る ん です け れ ども

その 場 合 は どう な ん で し ょ う か

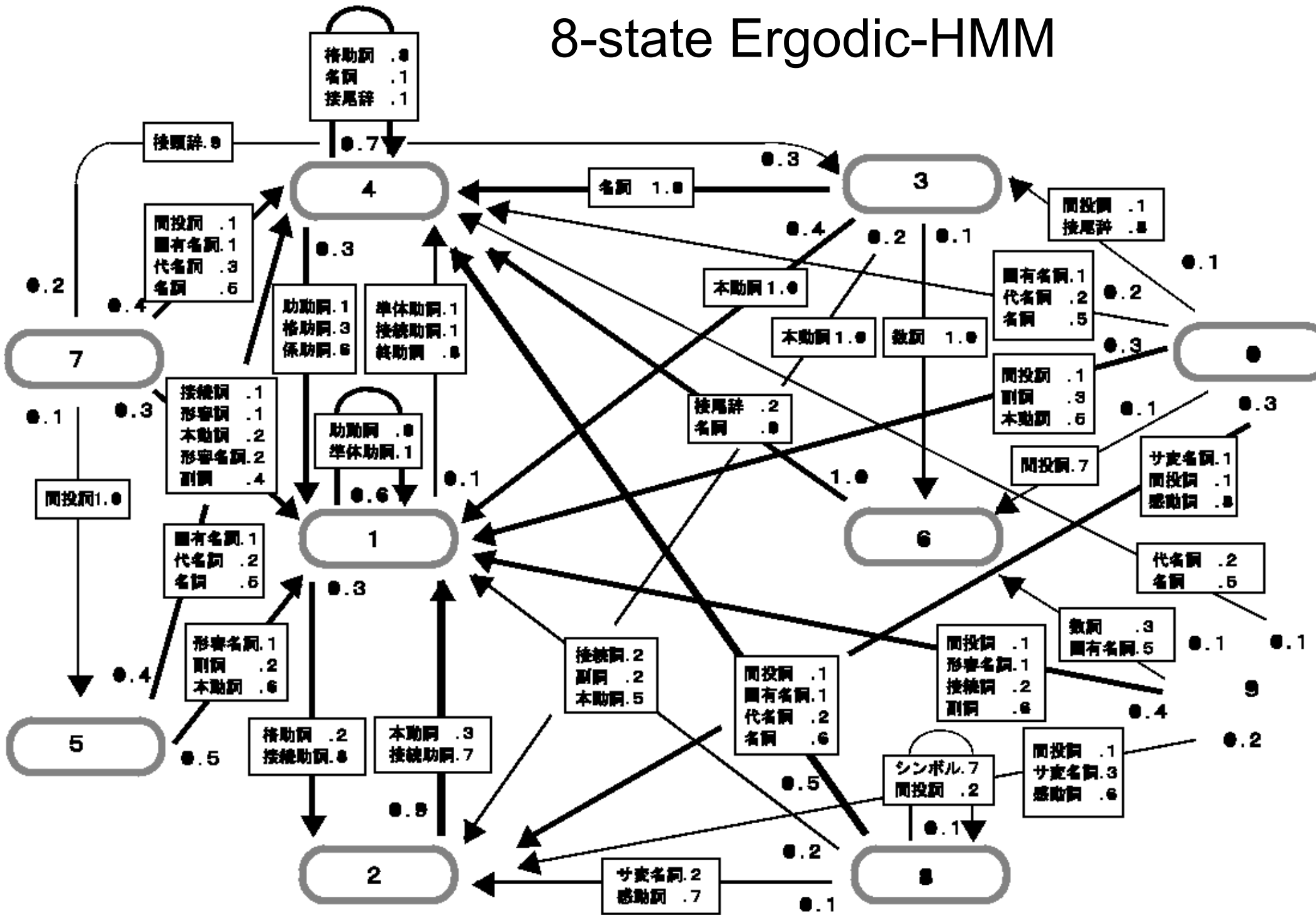
言語モデル生成実験の条件

HMMの構造	状態遷移出力型Ergodic HMM
HMMの状態数	4状態,8状態,16状態
HMMの出力シンボル	単語
開始 終了状態	任意
初期状態遷移確率	ランダム
初期シンボル出力確率	ランダム
初期状態確率	均等
語彙数	6418
学習データ set	odd4000,odd2000,odd1000
学習データ数	4000文,2000文,1000文
単語総数	57354単語,20730単語,13299単語
学習終了条件	尤度上昇率1%未滿

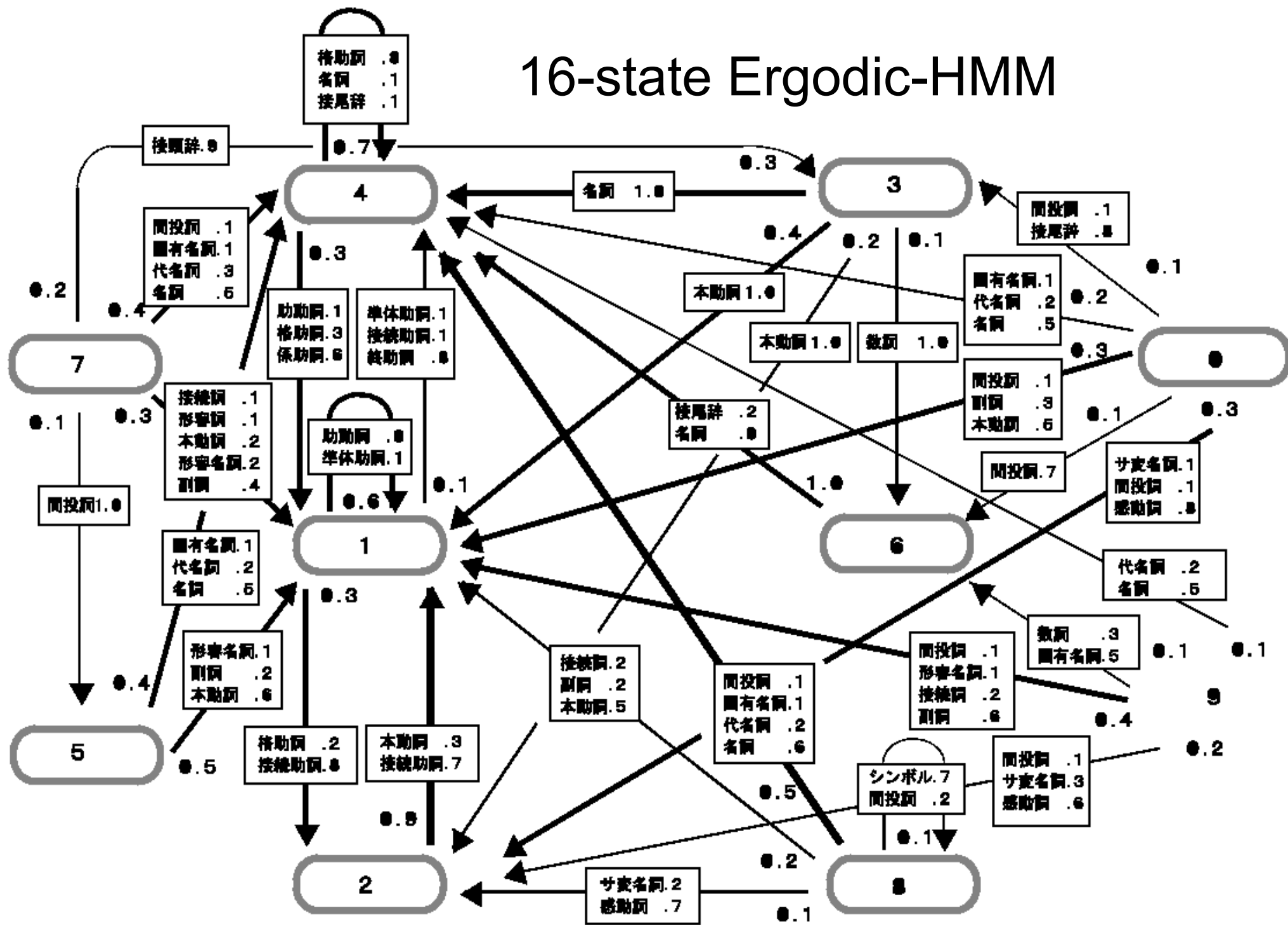
4-state Ergodic-HMM



8-state Ergodic-HMM



16-state Ergodic-HMM



まとめ

- ・4状態 体言と用言のグループ化
- ・8状態 活用系でのグループ化
- ・16状態 品詞でのグループ化
- ・確率付きネットワーク文法の自動獲得が可能

メモリ量および計算量を削減した
Baum - Welch アルゴリズムの提案と
言語モデルへの適用

1: 目的

状態数が大きいErgodic HMMの
Baum-Welch学習

1.2: 問題点

状態数が小さい

→ Perplexityが高い

状態数が多い

→ Baum-Welch 学習が不可能

1.3: 解決方法

メモリ量および計算量を削減した
Baum-Welchアルゴリズム

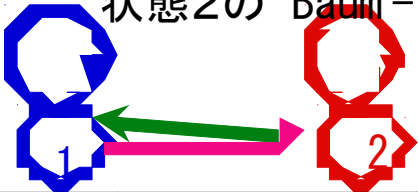
小さいシンボル出力確率の削除 (メモリ量, 計算量の削減)

シンボル出力確率が閾値 (10^{-300}) より
小さいとき 0.0

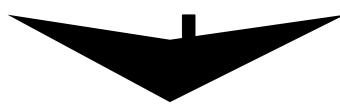
再推定およびメモリから削除。

状態数の逐次増加

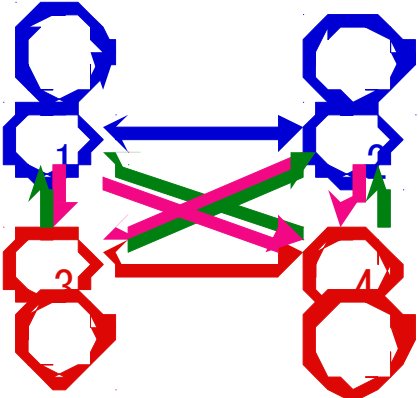
状態2の Baum-Welch 学習の終了後のパラメータ



$\pi_1 = 0.3$	$a_{11} = 0.4$
$\pi_2 = 0.7$	$a_{22} = 0.6$
	$a_{21} = 0.1$
	$a_{12} = 0.9$



状態4の初期モデルのパラメータ



$\pi_{41} = \pi_{21} = 0.3$	$a_{11} = a_{21} = 0.4$
$\pi_{42} = \pi_{22} = 0.3$	$a_{12} = a_{22} = 0.4$
$\pi_{43} = \pi_{23} = 0.7$	$a_{13} = a_{23} = 0.6$
$\pi_{44} = \pi_{24} = 0.7$	$a_{14} = a_{24} = 0.6$

ただしシンボル出力確率は乱数を使用

N状態の Ergodic HMM のパラメータ

$\pi_N(i)$; $i=1, \dots, N$; 初期状態確率
 $a_N(i, j)$; $i=1, \dots, N, j=1, \dots, N$; 状態遷移確率
 $b_N(i, j, w)$; $i=1, \dots, N, j=1, \dots, N, w=1, \dots, V$; シンボル出力確率
 V ; 語彙数

2N状態の Ergodic HMM の

初期状態確率および状態遷移確率の初期パラメータ

$\pi_{2N}(i) = 0.5 \times \pi_N(i/2)$; $i=1, \dots, 2N$
 $a_{2N}(i, j) = 0.5 \times a_N(i/2, j/2)$; $i=1, \dots, 2N, j=1, \dots, 2N$
 $b_{2N}(i, j, w) = b_N(i/2, j/2, w) \times \text{random}(i, j, w)$;
 $i = 1, \dots, 2N, j = 1, \dots, 2N, w = 1, \dots, V.$ ただし $\sum_w b_{2N}(i, j, w) = 1.0$
" / " 切り上げを意味

具体的なアルゴリズム

初期 Ergodic HMM
state number = 1

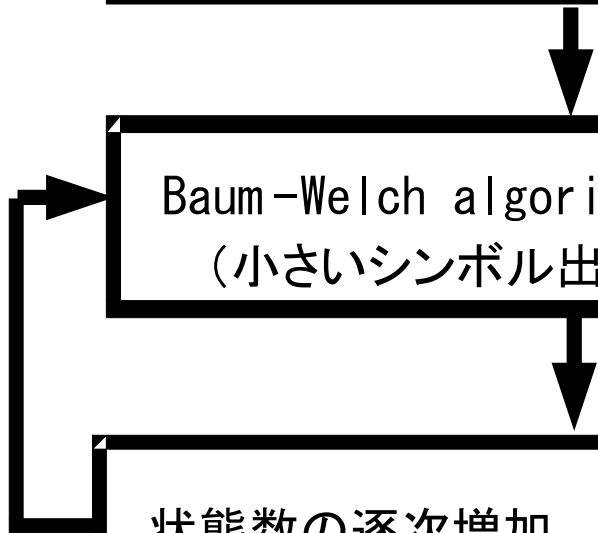
π, A, B ; random
state number = 1

Baum-Welch algorithm
(小さいシンボル出力確率の削除)

state number = N

状態数の逐次増加

state number = 2N

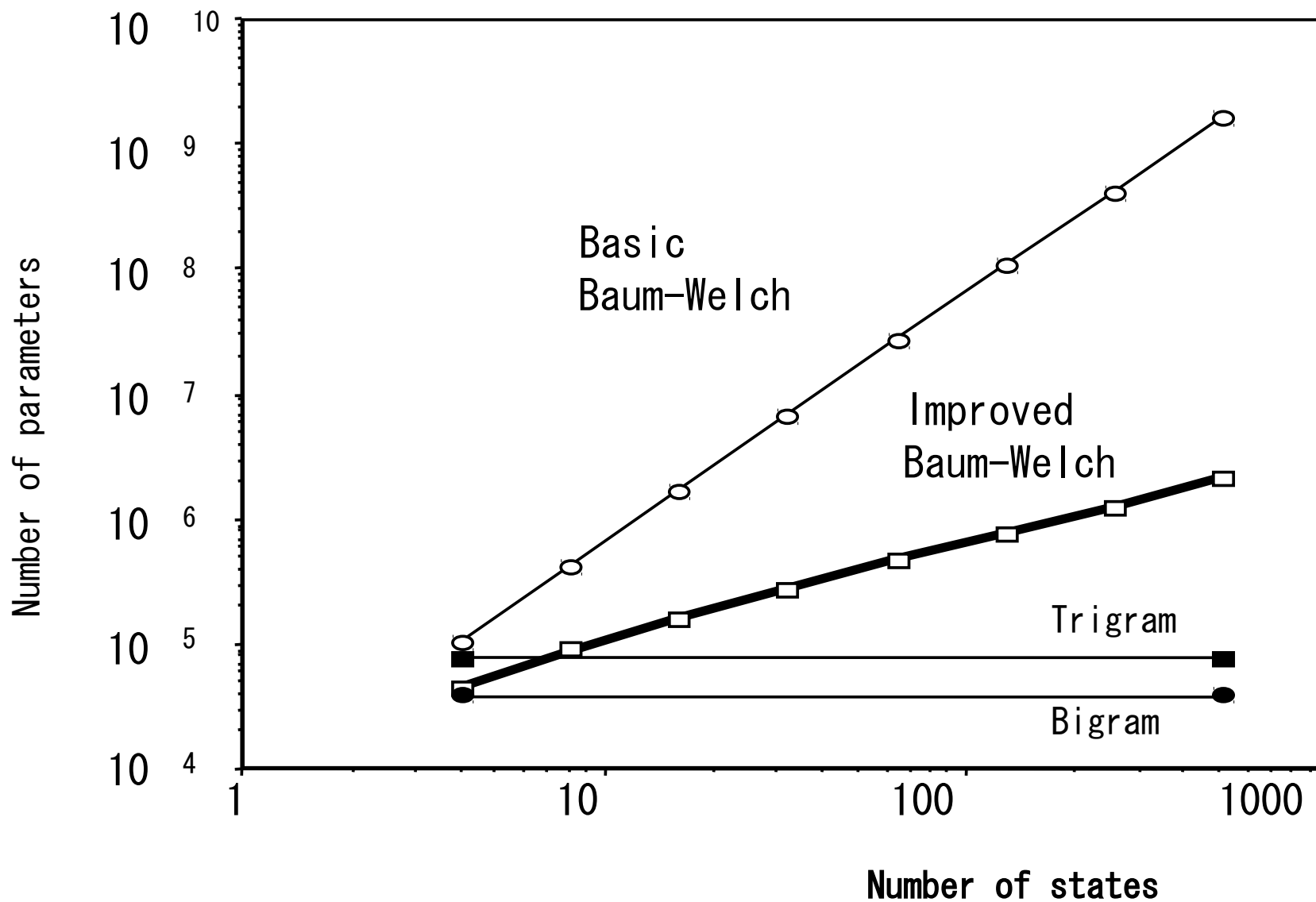


確率付きネットワーク文法の獲得の実験

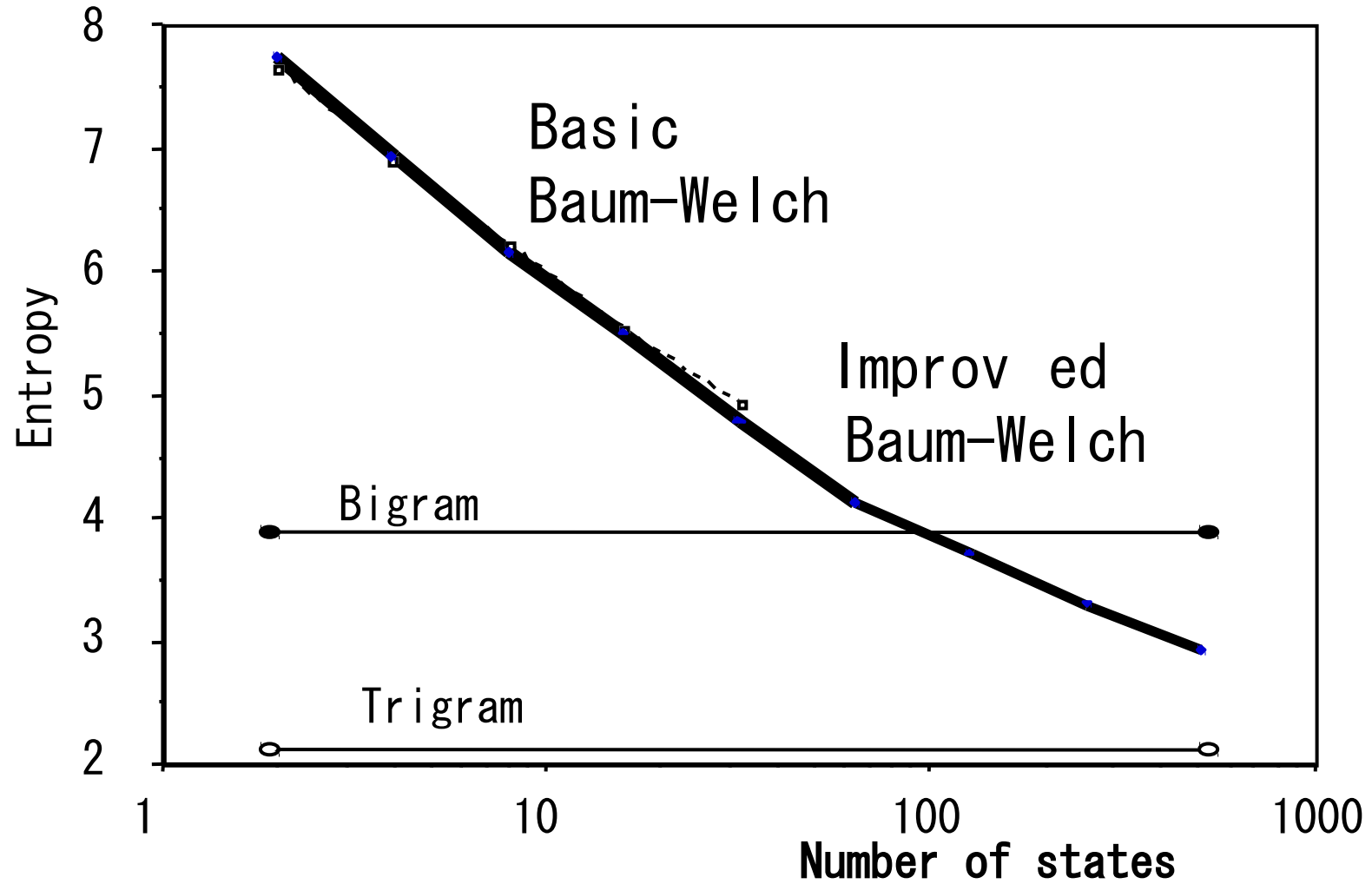
HMMの構造	状態遷移出力型
学習語彙数	6420 単語
学習データ数	8475
総単語数	57354
Baum-Welchアルゴリズムの終了条件	40回の繰り返し

言語モデル生成実験の条件

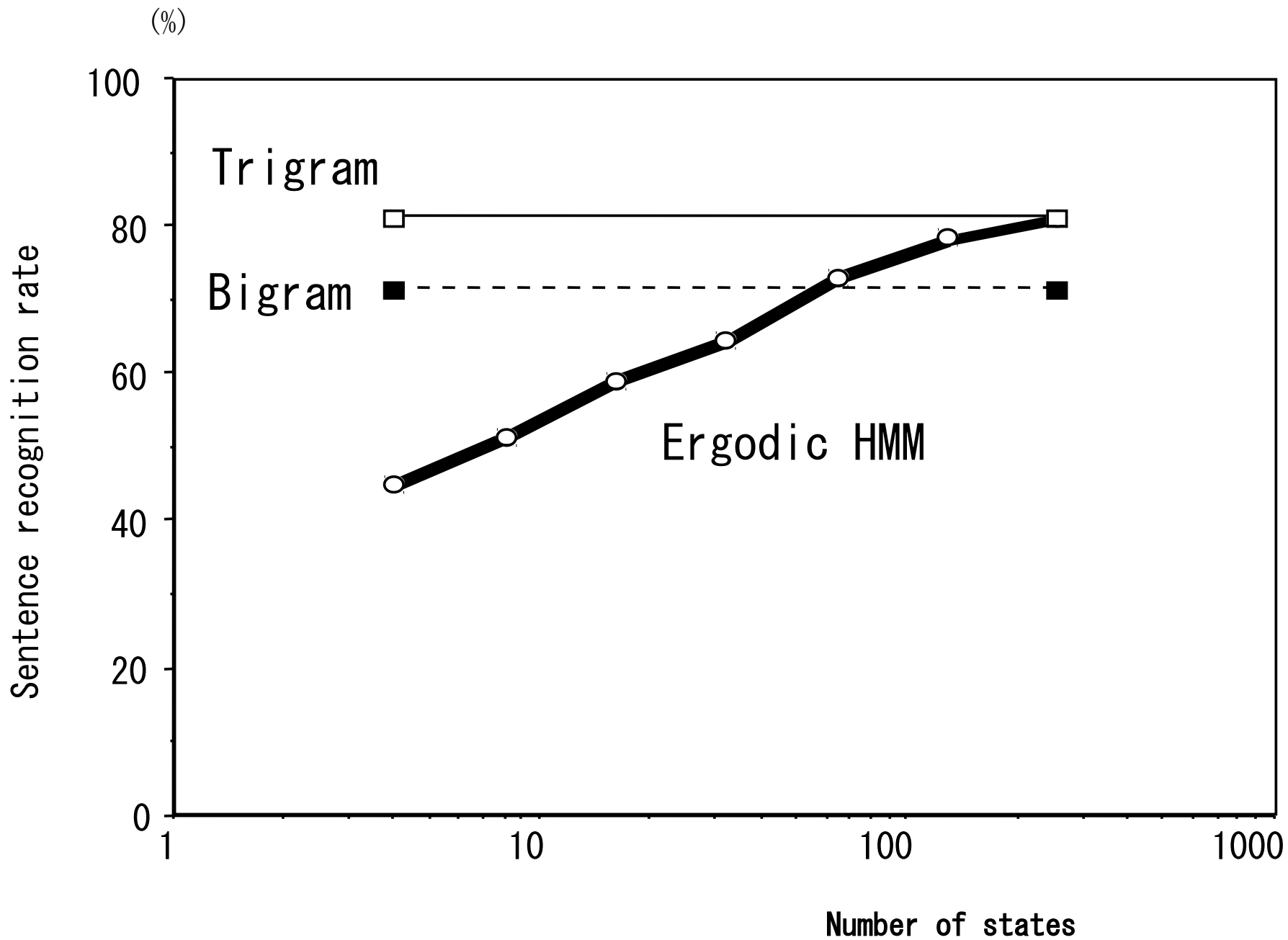
パラメータの数 vs 状態数



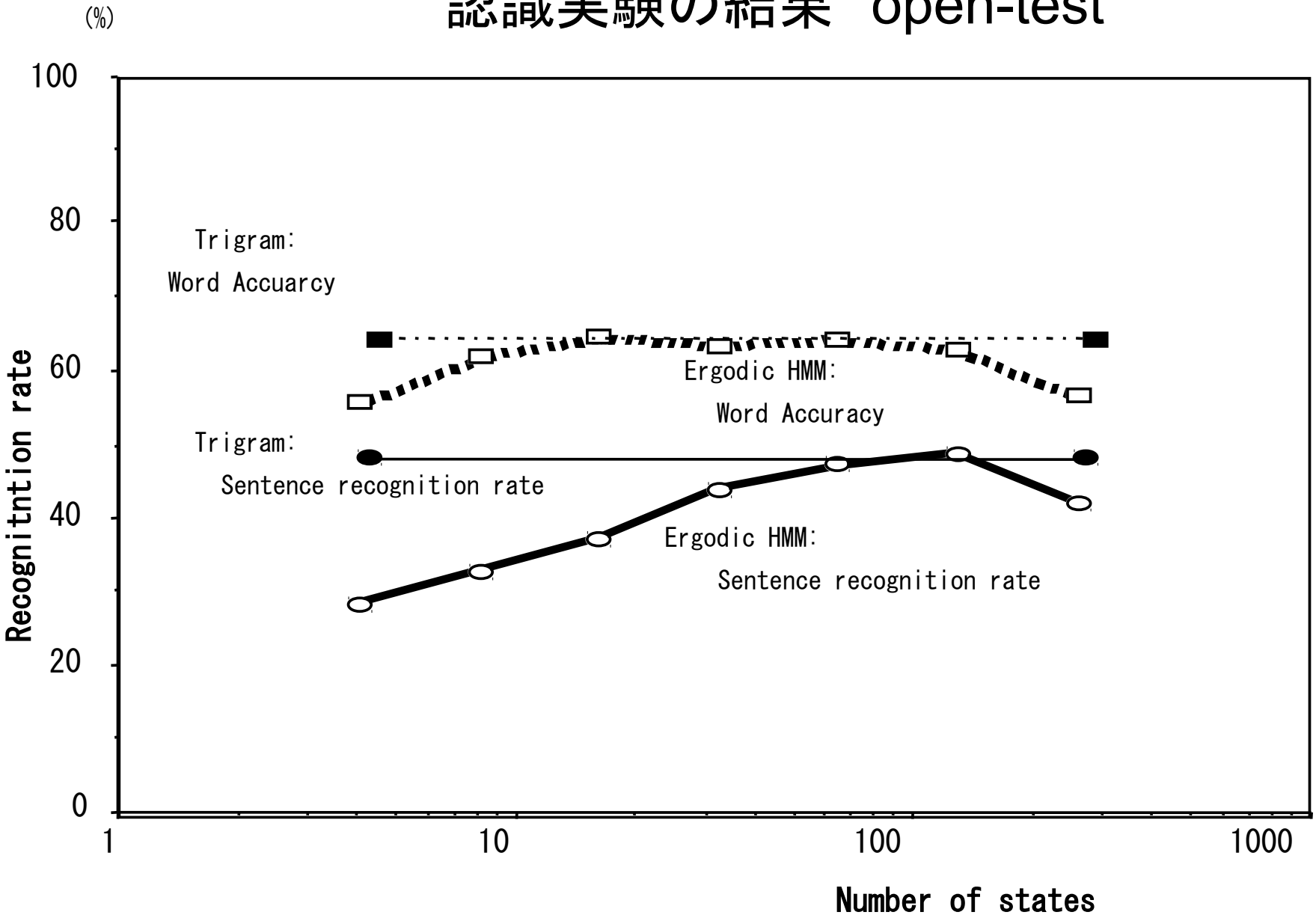
エントロピーの変化



認識実験の結果 closed-test



認識実験の結果 open-test



まとめ

メモリ量および計算量を削減した
Baum-Welch アルゴリズムを提案

確率つきネットワーク文法の獲得
perplexity:
bigram > Ergodic HMM > trigram

言語モデルとして
連続音声認識に利用

認識率

text-closed: bigram < Ergodic HMM < trigram
text-open : trigram < Ergodic HMM