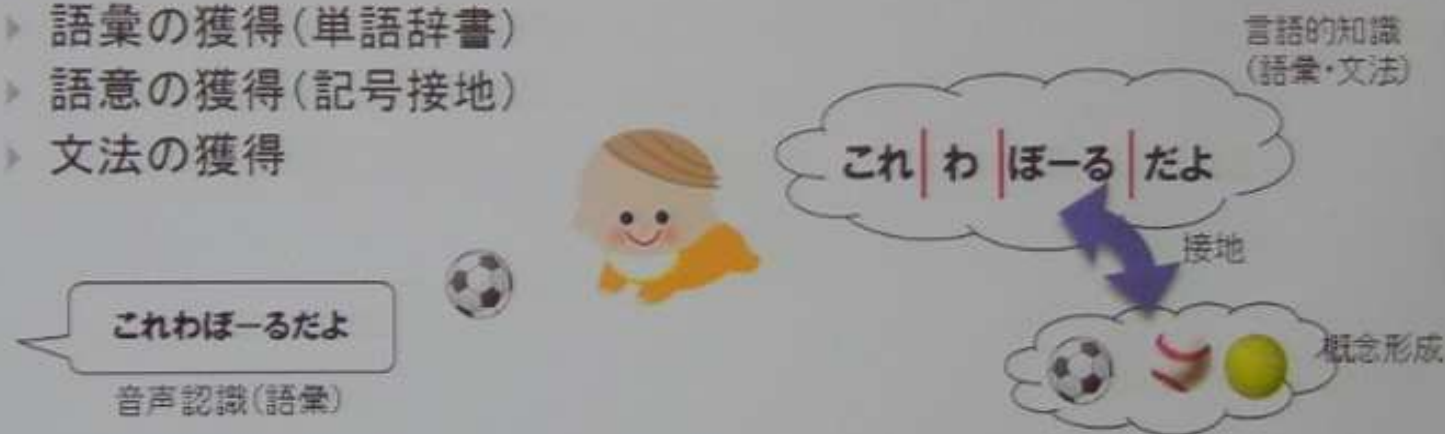


マルチモーダルカテゴリゼーション  
～階層ベイズモデルに基づくロボットによる概念・言語獲得～

電気通信大学 中村 友昭

## 概念・言語学習

- ▶ 人のように言語を獲得するロボットの実現
  - ▶ 人や環境とのインタラクションにより自律的に獲得
    - ▶ 概念形成
    - ▶ 語彙の獲得(単語辞書)
    - ▶ 語意の獲得(記号接地)
    - ▶ 文法の獲得

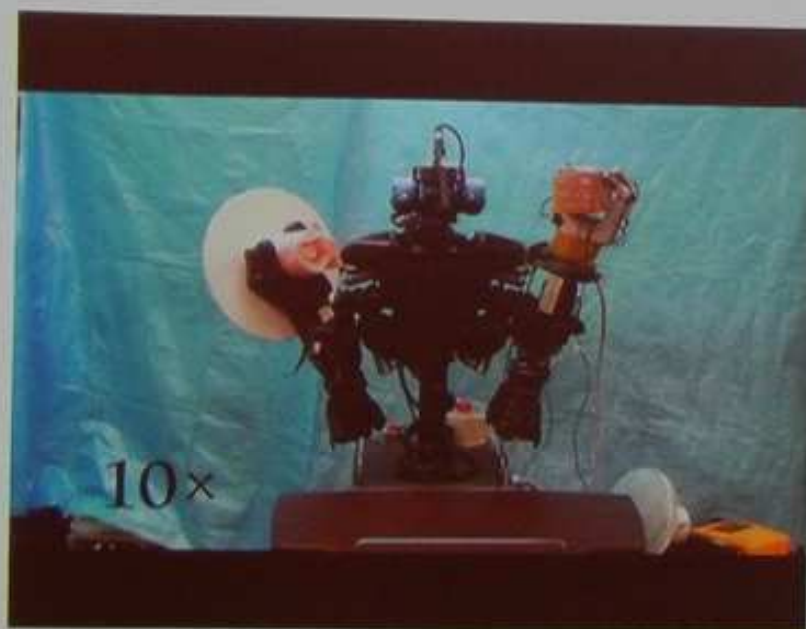


- ▶ 言語獲得アルゴリズムを確率モデルを用いて実現
  - ▶ 人間のような知能の実現
  - ▶ 人間の言語獲得過程の解明

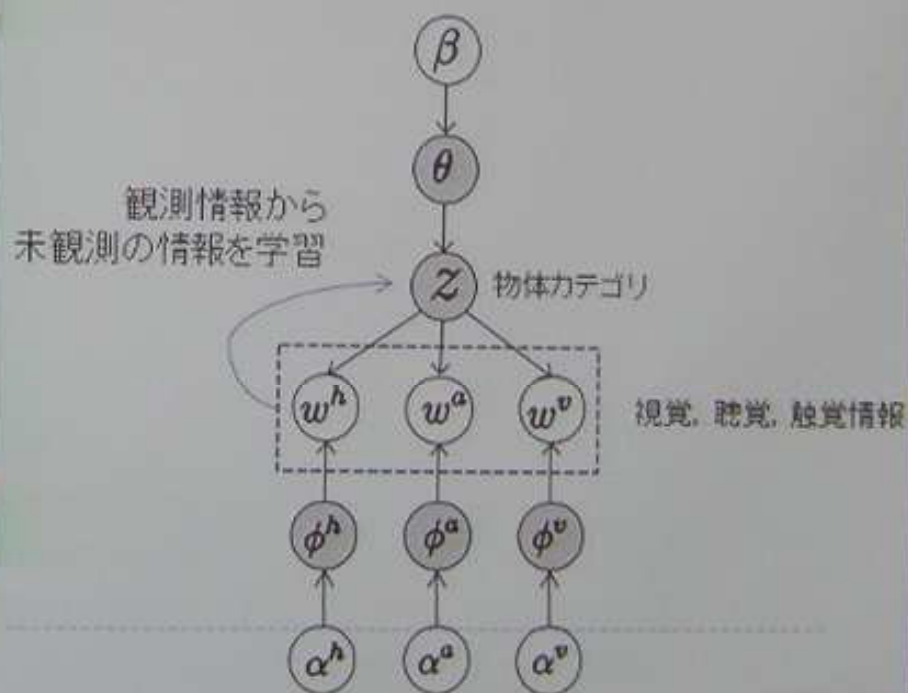
# 概念形成

## ▶ 概念の工学的定義

- ▶ 概念 = **知覚情報** のクラスタリングによって形成された **カテゴリ**
- ▶ 知覚情報: ロボットの視覚, 聴覚, 触覚情報
- ▶ これらの情報を確率モデルにより **教師なし** で分類

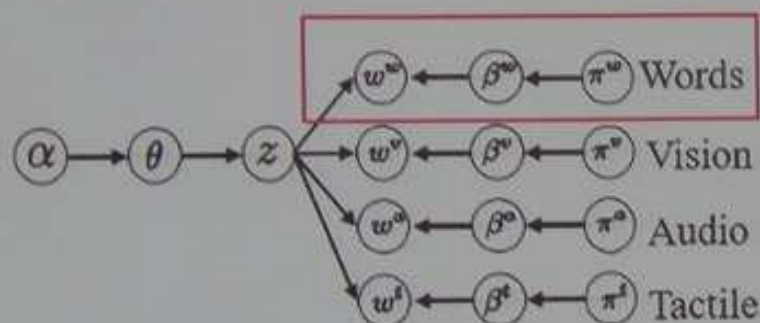


物体に見て・触れることで物体の概念を形成



# 語意の獲得

- ▶ マルチモーダル情報として単語を追加
- ▶ 単語も含めたマルチモーダルな概念を形成



それは  
ペットボトルだよ



- ▶ 単語からマルチモーダル情報を確率的に予測可能



$$P(\underline{w^v}, \underline{w^a}, \underline{w^t} | \underline{w^w})$$

マルチモーダル情報      単語

単語を“感覚的”に理解可能  
⇒ 単語の意味の理解



# 語彙の獲得

- ▶ 前の研究では語彙を持っていることが前提
  - ▶ 語彙＝音声認識・単語分割で使用する単語辞書(言語モデル)
- ▶ 語彙を獲得する際の問題
  - ▶ 語彙を持たないため  
音声が正しく認識できない
    - ▶ ぼーる？ ぼーう？ ごーる？
  - ▶ 単語の切れ目が分からない
    - ▶ ぼーる？ わぼー？ ぼーるだ？ るだよ？
- ▶ これらの問題を言語的知識と概念を相互学習することで解決
  - ▶ 言語のパターンに基づく単語分割
  - ▶ 同じ概念に含まれる物体には同じ単語が教示される可能性が高い

これわぼーるだよ



ぼーるがあるよ

これわぼーるだよ

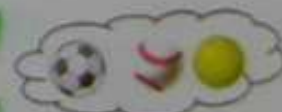


これわ | ぼーる | だよ

言語的知識



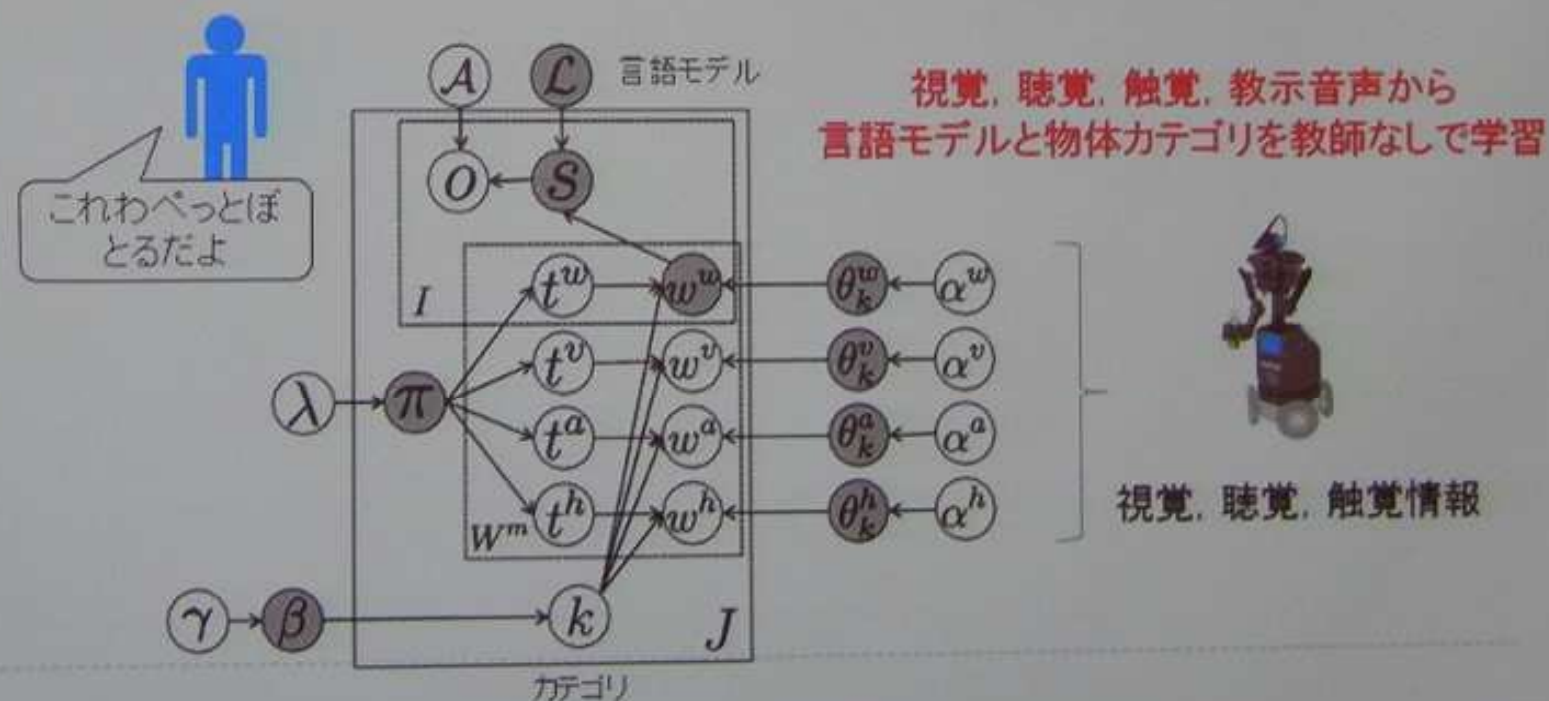
相互に学習



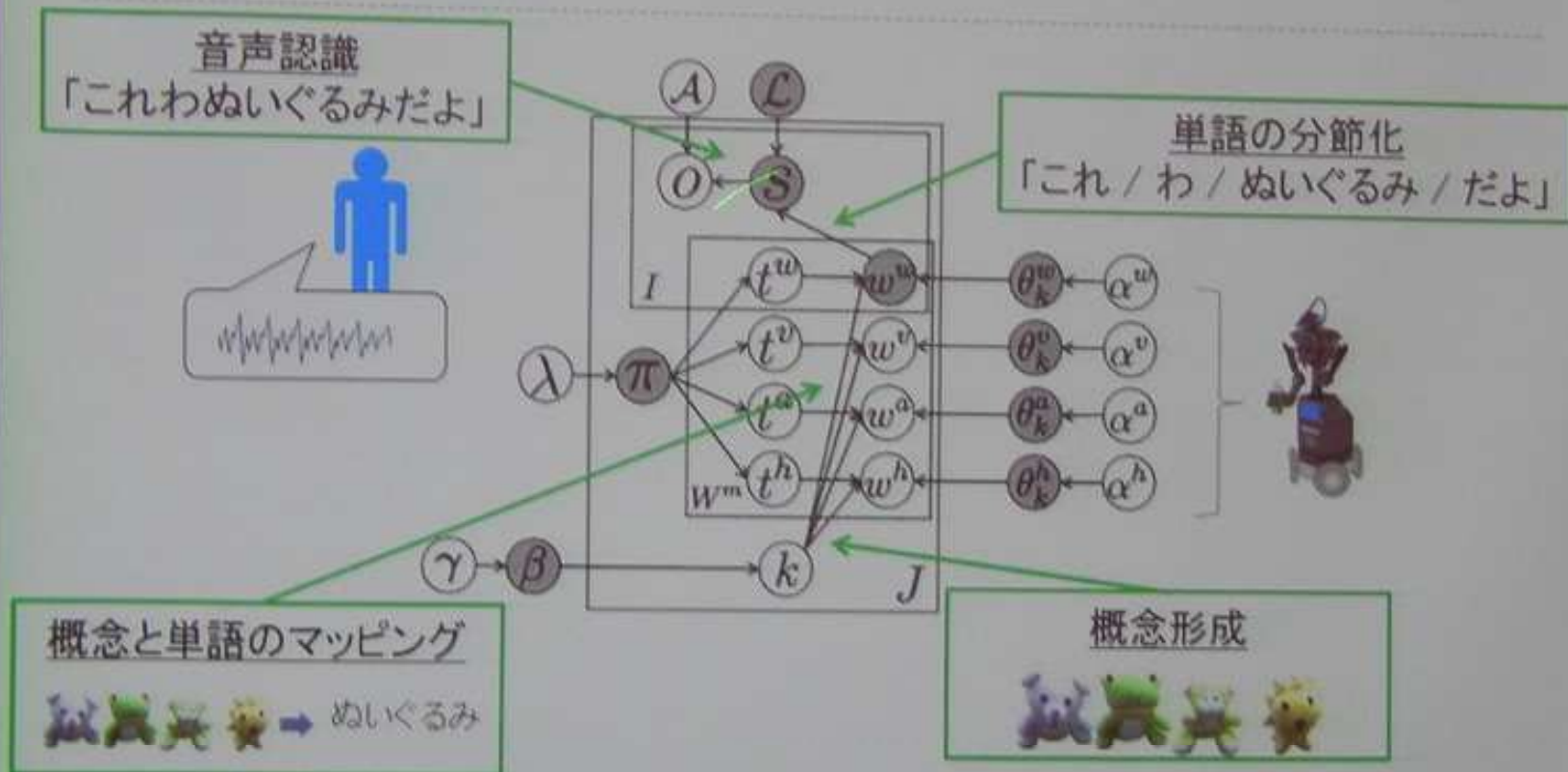
概念形成

# 提案モデル

- ▶ 教示音声・マルチモーダル情報の生成モデル
  - ▶ 概念の形成と同時に言語モデル(語彙)の獲得が可能
  - ▶ 物体カテゴリ $k$ と言語モデル $\mathcal{L}$ が結びついたモデル  $\rightarrow$  相互に影響
    - ▶ 同じ物体には同じ単語が発話される可能性が高い
    - ▶ 同じ単語が与えられた物体は同じカテゴリの物体である可能性が高い



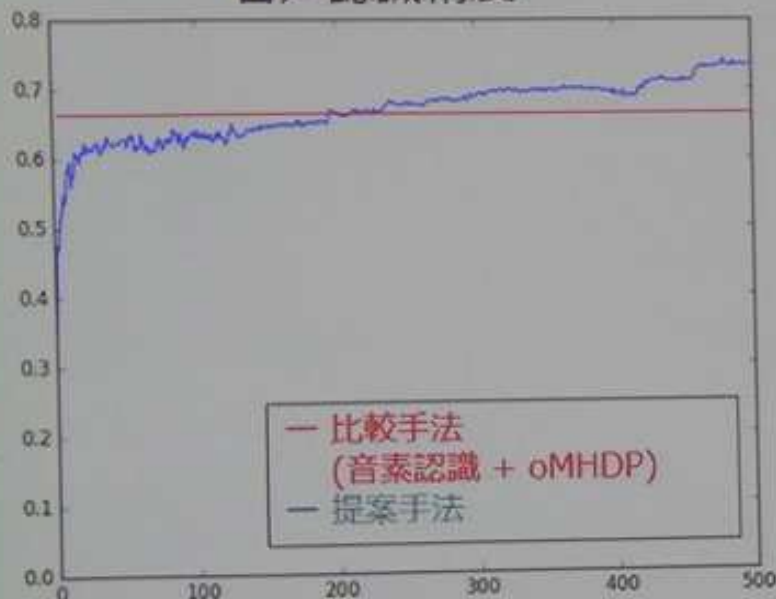
# モデルの構造



- ▶ 音声認識, 単語分節化, 概念形成, 概念と単語の結びつきを全て**教師なし**で学習

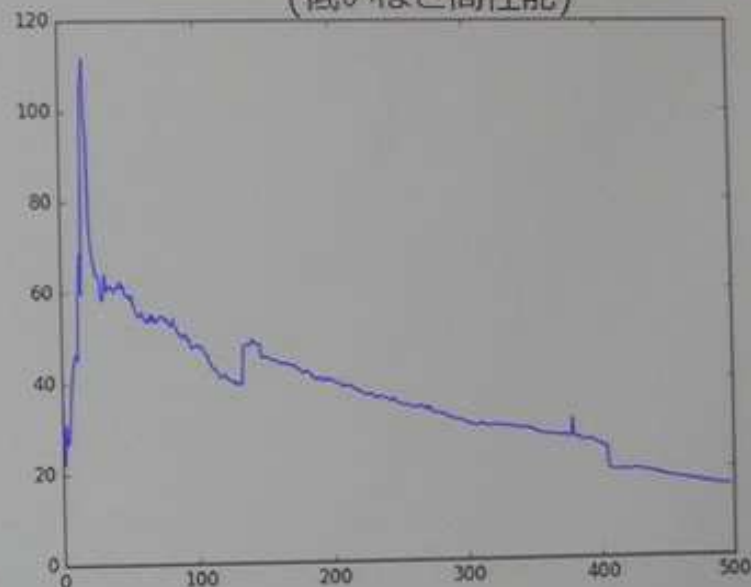
# 音声認識・言語モデル(語彙)の評価

音声認識精度



Number of learned objects

言語モデルのPerplexity  
(低いほど高性能)



Number of learned objects

- ▶ 学習により正しい言語モデル(語彙)が獲得された  
    ➡ 音声認識性能が向上





# 多様な概念と文法の学習

- ▶ 概念には物体だけでなく様々な概念が存在
  - ➡ 人の行動シーンから物体, 人, 場所, 動き概念を形成
- ▶ 概念クラスの順序規則を文法として学習
  - ➡ 文法を学習することで観測シーンの言語表現が可能



ロボットの観測シーン

物体情報

動き情報

場所情報

言語情報



ジュース  
ペットボトル

物体概念

単語



飲む  
食べる

動き概念

単語



ソファ  
リビング

場所概念

単語



女の子  
女性

人概念

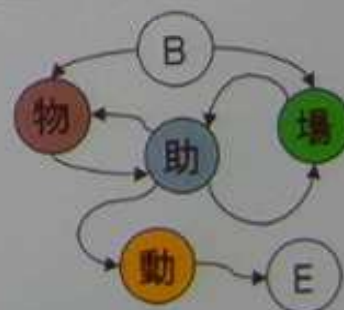
単語

概念形成・語意獲得

概念と単語の結び付きを獲得

ソファ で ジュース を 飲む  
↓ ↓ ↓ ↓ ↓  
場所 助詞 物体 助詞 動き

概念の遷移を文法として学習



# 複数概念モデルと文法モデル

## ▶ 複数概念モデル (mMLDA)

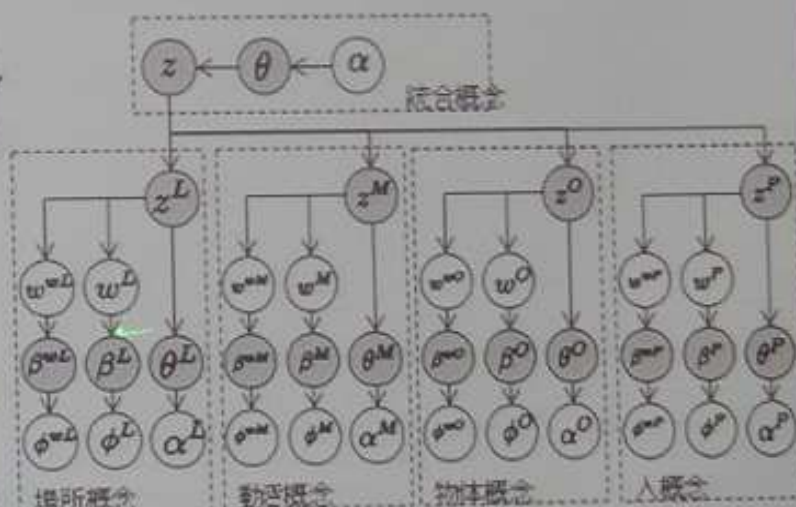
- ▶ 場所, 動き, 物体, 人概念の生成モデル
- ▶ ロボットが観測した情報から概念を教師なしで学習可能
- ▶ 単語と概念の結び付きも学習

## ▶ 文法モデル (HSMM)

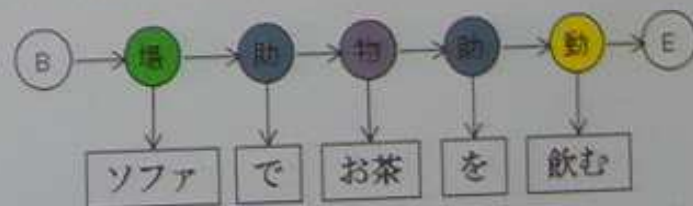
- ▶ 概念の結び付きから概念クラスの遷移順を文法として学習
- ▶ 概念モデルの結果を初期値として教師なし学習

## ▶ 文法と概念を相互に学習

- ▶ 言語的制約と知覚情報との共起性に基づく単語の接地

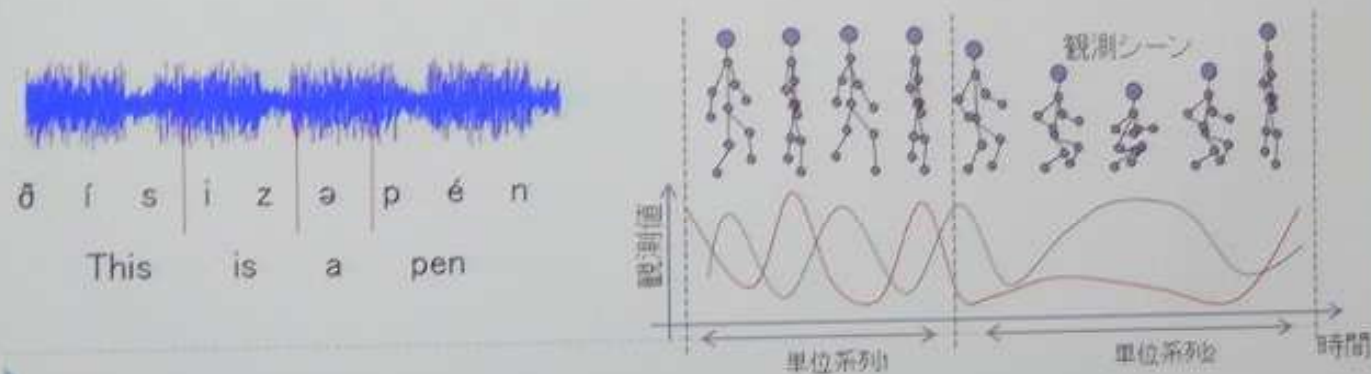


言語的な制約 ↑ ↓ 知覚情報との共起性



## 連続的なセンサ情報の分節化

- ▶ これまでの研究では**分節化**は考えていなかった
  - ▶ 動作概念なども単位動作毎にデータを人手で区切っていた
- ▶ 実際のセンサ情報は連続であり教師なしで分節・分類する必要がある
- ▶ ロボットが連続的なセンサ情報を**教師なしで分節・分類**し概念を獲得

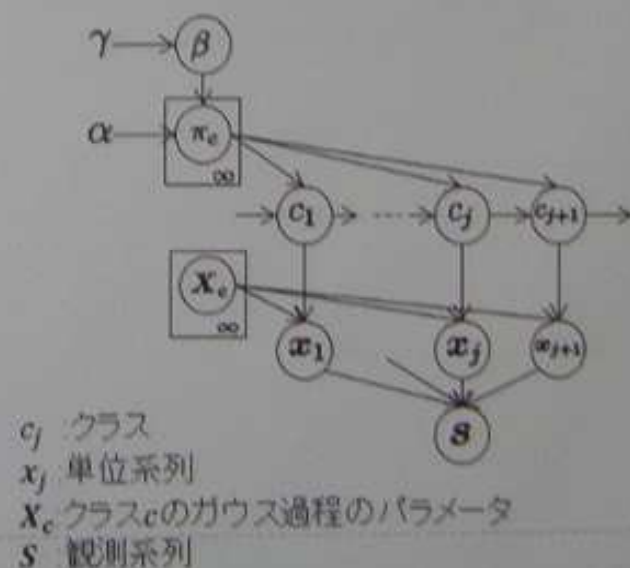
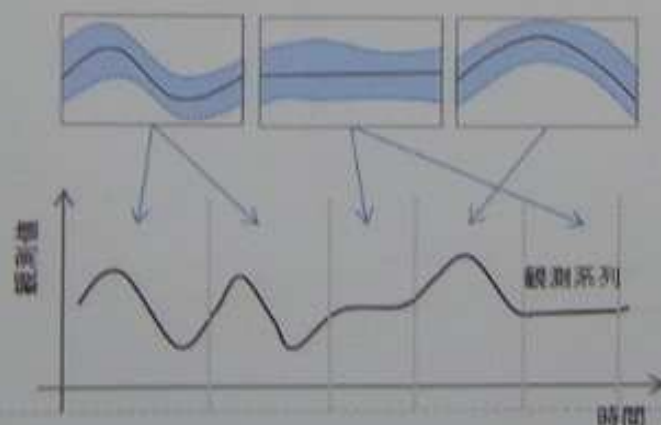




## 提案モデル

- ▶ 時系列データはガウス過程(GP)を出力分布とする隠れセミマルコフモデル(HSMM)によって生成されると仮定
  - ▶ ガウス過程: 分節化された定型パターンを表現
  - ▶ HSMM: 定型パターンの長さも隠れ変数としたHMM
- ▶ HSMMとGPのパラメータ推定することで, 連続的な情報の分節・分類が可能

定型パターン(ガウス過程)





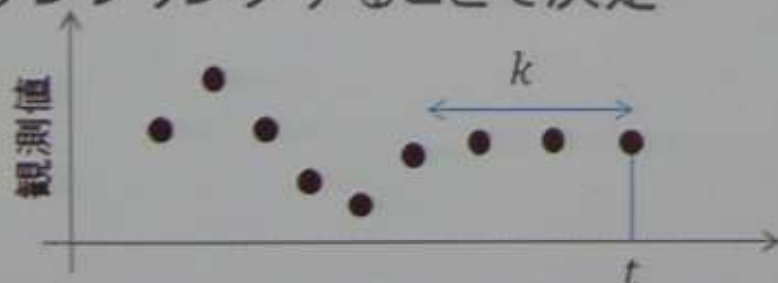
# GP-HSMMのパラメータ推定

- ▶ HSMMでは1つの状態に分類される  
系列長は状態によって異なる

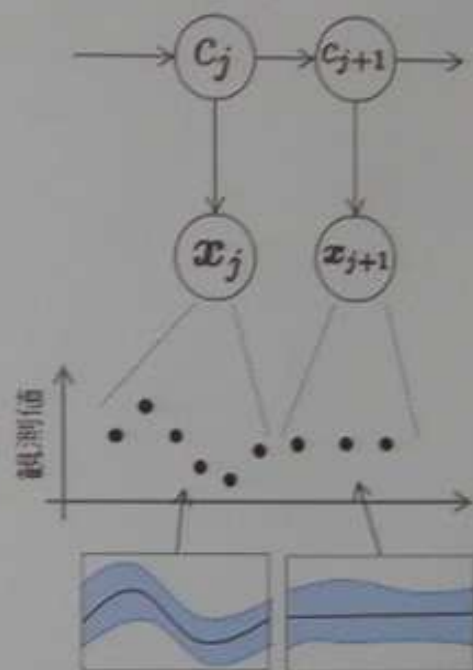
➡ 系列長も推定する必要がある

- ▶ パラメータ推定

- ▶ 時刻 $t$ のデータ点を終点とした長さ $k$ の  
系列のクラスが $c$ である確率から  
サンプリングすることで決定

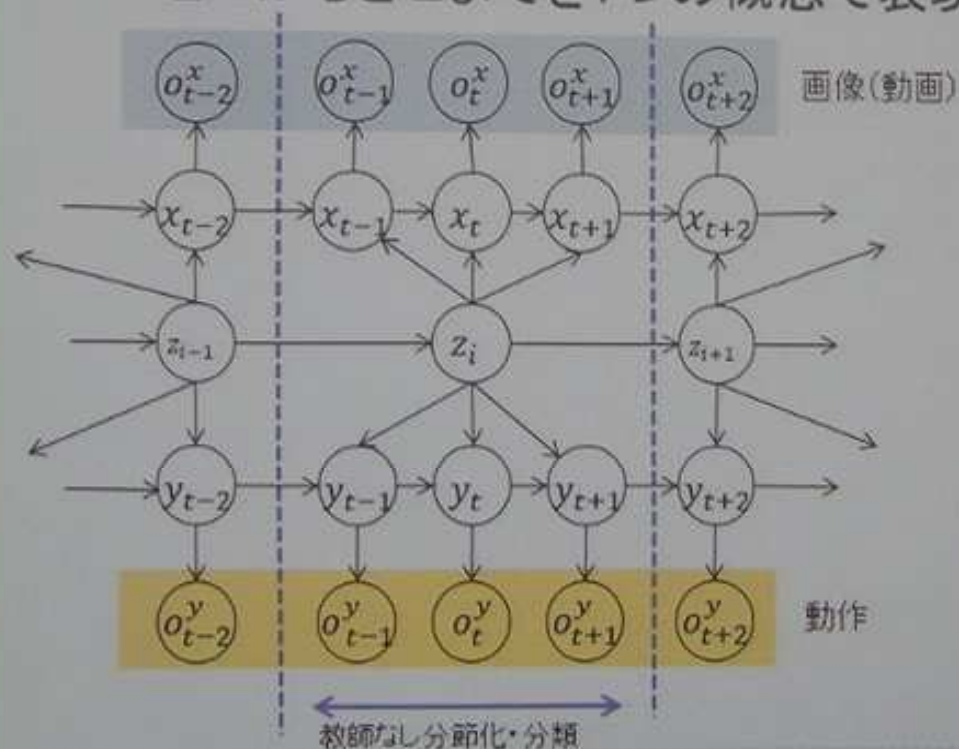


- ▶ あらゆる $k$ と $c$ の組み合わせの確率を計算する必要がある  
＝動的計画法(Forward filtering-Backward sampling)を利用



## マルチモーダル情報の分節・分類

- ▶ mMLDAの時間展開に基づく時系列マルチモーダル情報の分節・分類に基づく概念形成
- ▶ どこからどこまでを1つの概念で表現可能化を教師なしで推定

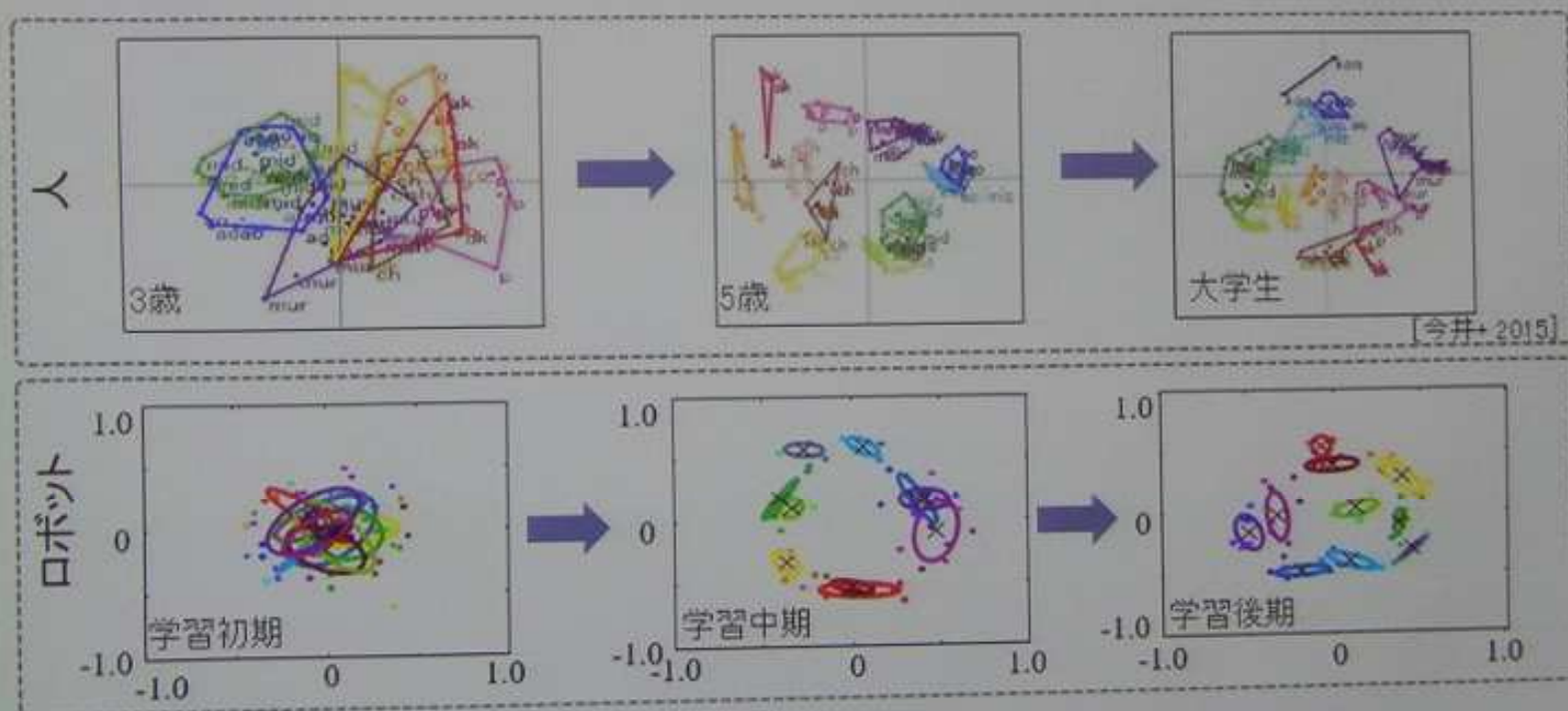


ロボットから取得可能な  
連続センサ情報からの  
柔軟な概念形成の実現

# 人の概念獲得との比較

## ▶ 提案モデルによって概念獲得過程が再現可能か検証

- ▶ 概念を二次元空間にプロット→似たような変化を確認
- ▶ 教示発話の違いによる概念学習の変化の検証

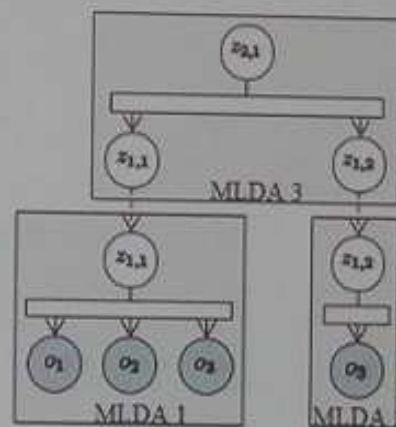
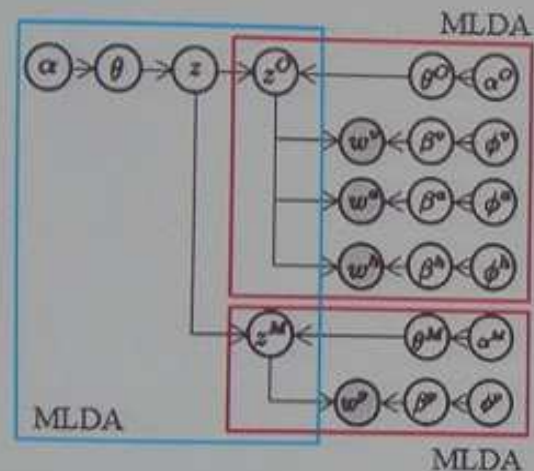


- 船田 他, "マルチモーダル概念形成における概念と言語の相互作用の解析", 人工知能学会全国大会, 2016
- Funada et al, "Analysis of the Effect of Infant-Directed Speech on Mutual Learning of Concepts and Language Based on MLDA and Unsupervised Word Segmentation", IROS2017: ML-HLOR, 2017

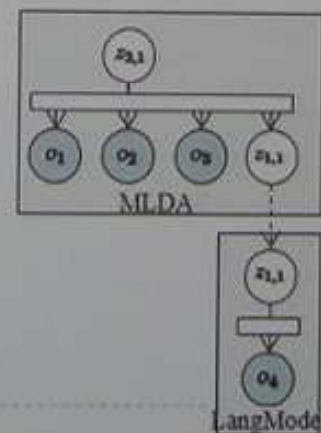
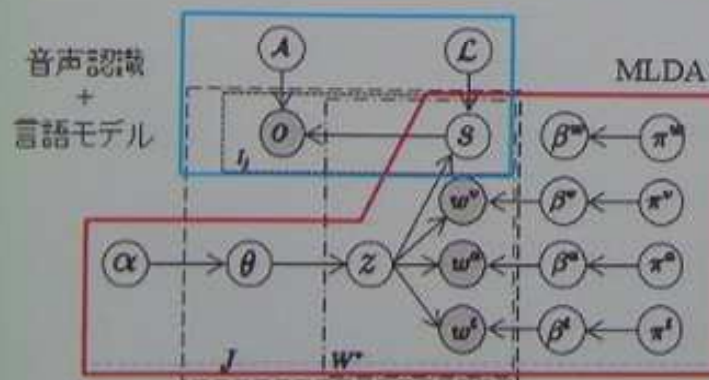


# 統計モデルの統合による大規模なモデルの構築

- ▶ 紹介したモデルは基礎的なモデルの統合によって構築



確率変数間の依存関係を定義することで、大規模なモデルを容易に実装可能なフレームワーク



より様々な能力を学習可能なロボット(汎用人工知能)の実現



## まとめ

- ▶ 本発表ではロボットによる概念・語彙・語意・文法獲得に関する研究を紹介
  - ▶ ロボットが取得したマルチモーダルな情報から教師なしでボトムアップに概念を形成
  - ▶ 人の教示発話から語彙・文法を学習
    - ▶ 単語と概念を結びつけることで単語の意味の学習
    - ▶ 単語と結びついた概念の遷移規則を文法として学習
  - ▶ 獲得した言語の情報がトップダウンに作用することでより人の感覚に近い概念の形成が可能