



# Symbol Emergence in Robotics: Language Acquisition via Real-world Sensorimotor Information

Tadahiro Taniguchi

- 1) Professor, College of Information Science & Engineering,  
Ritsumeikan University
- 2) Visiting General Chief Scientist, AI solution center,  
Panasonic

Invited talk at Gatsby-Kakenhi Joint Workshop on AI and Neuroscience

12<sup>th</sup> May 2017



# Tadahiro Taniguchi @tanichu

- Professor, Emergent System Laboratory,  
College of Information Science and Engineering,  
Ritsumeikan University, Japan
  - 2003-2006: PhD student, Kyoto University
  - 2005-2008: JSPS research fellow, Kyoto University
  - 2008: Assistant professor, Ritsumeikan University
  - 2010: Associate professor, Ritsumeikan University
  - 2015-2016: Visiting Associate Professor,  
Imperial College London
  - 2017: Professor, Ritsumeikan University
  - 2017: Visiting General Chief Scientist,  
Panasonic Corporation  
AI solution center (20% C.A.)
- **Research Topics**
  - Machine learning, Intelligent Robotics,  
Symbol emergence in robotics, Language acquisition



# Contents

## 1. Introduction

## 2. Lexical acquisition tasks

- A) Direct phoneme and word discovery from speech signals
- B) Simultaneous acquisition of word units and multimodal categories
- C) Online spatial concept acquisition

## 3. Future challenges

# Computational Understanding of Mental Development

## From Behavioral Learning to Language Acquisition



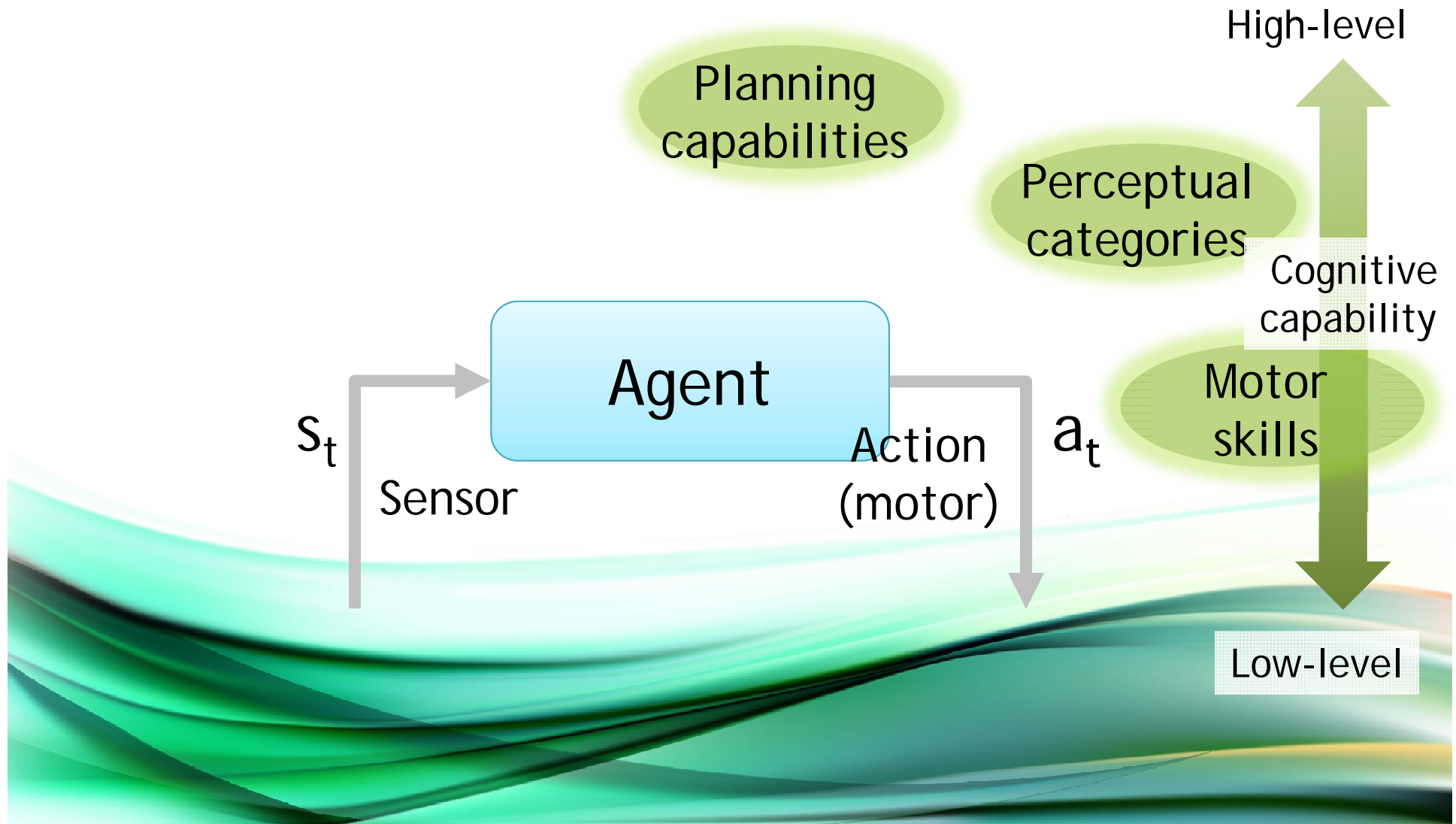
- ❑ A human child acquires many physical skills, concepts, and knowledge, including language, through physical and social interaction with his/her environment.
- ❑ How do we become able to communicate via symbols?
- ❑ We'd like to obtain an understanding of the **computational process** of mental development and language acquisition.

### Constructive approach

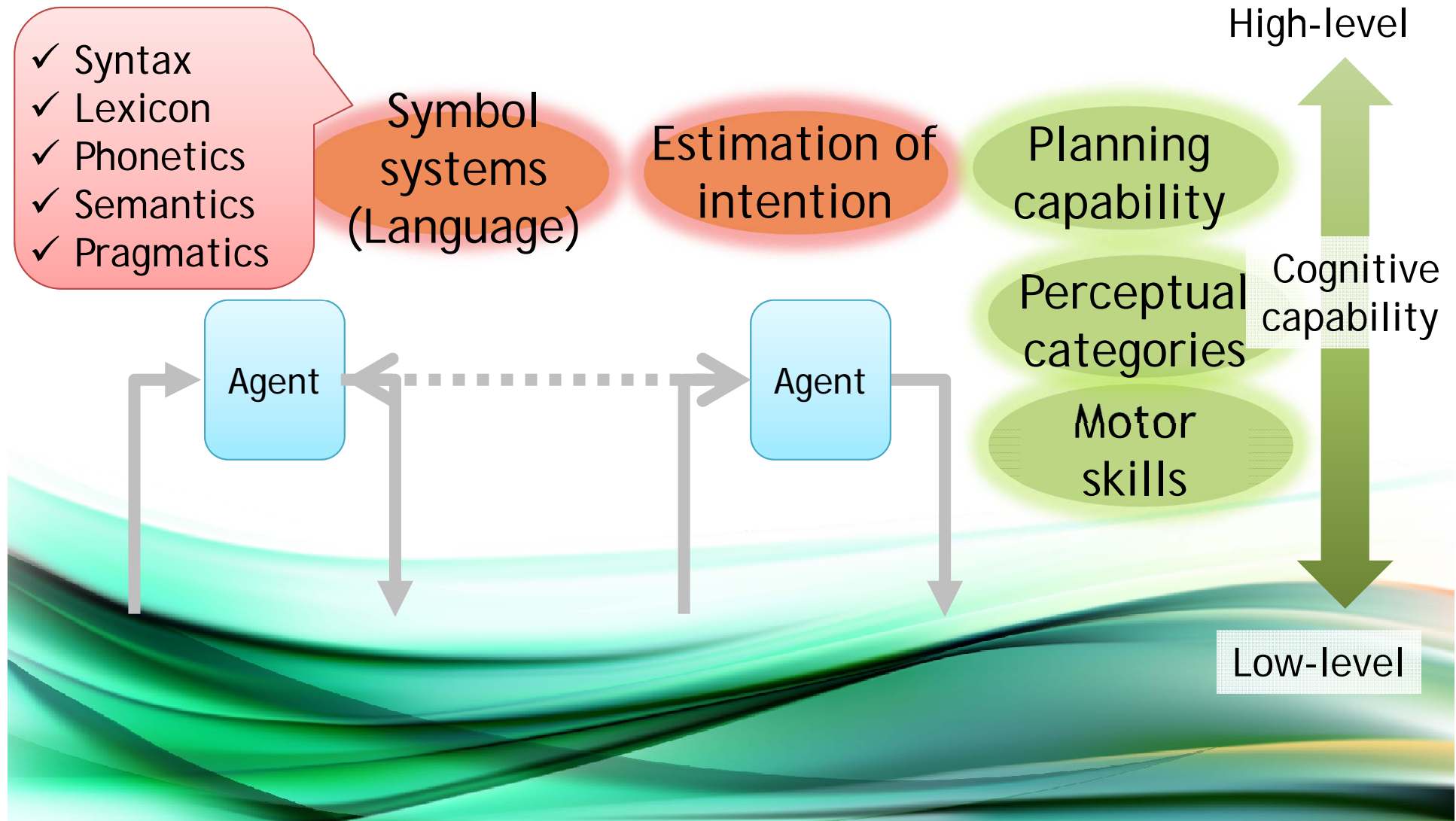
Develop robotic and computational models to better understand the original

## Symbol Emergence in Robotics

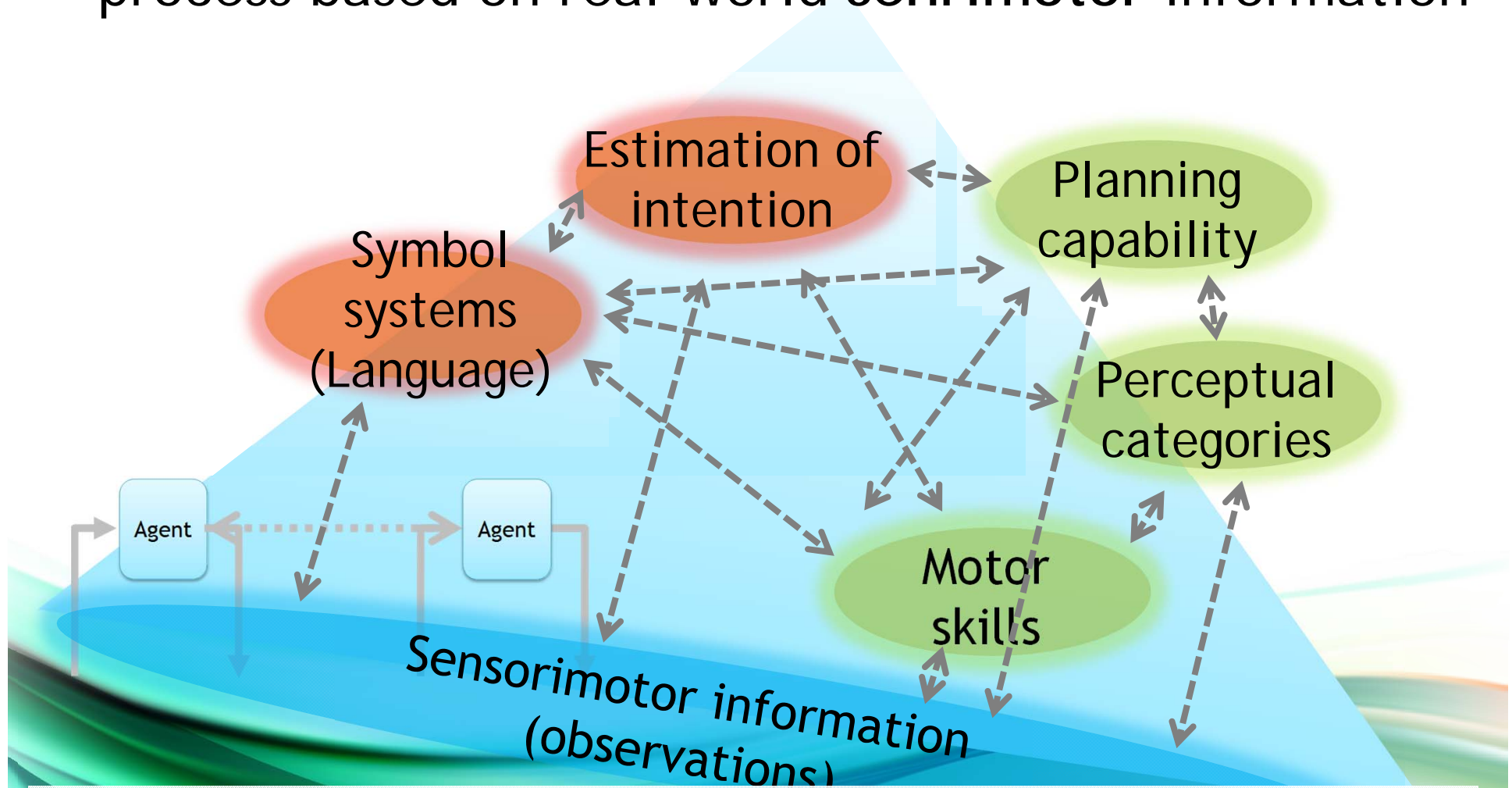
Development through a self-organizational learning process based on real-world sensorimotor information



# Development through a self-organizational learning process based on real-world sensorimotor information



Development through a self-organizational learning process based on real-world sensorimotor information

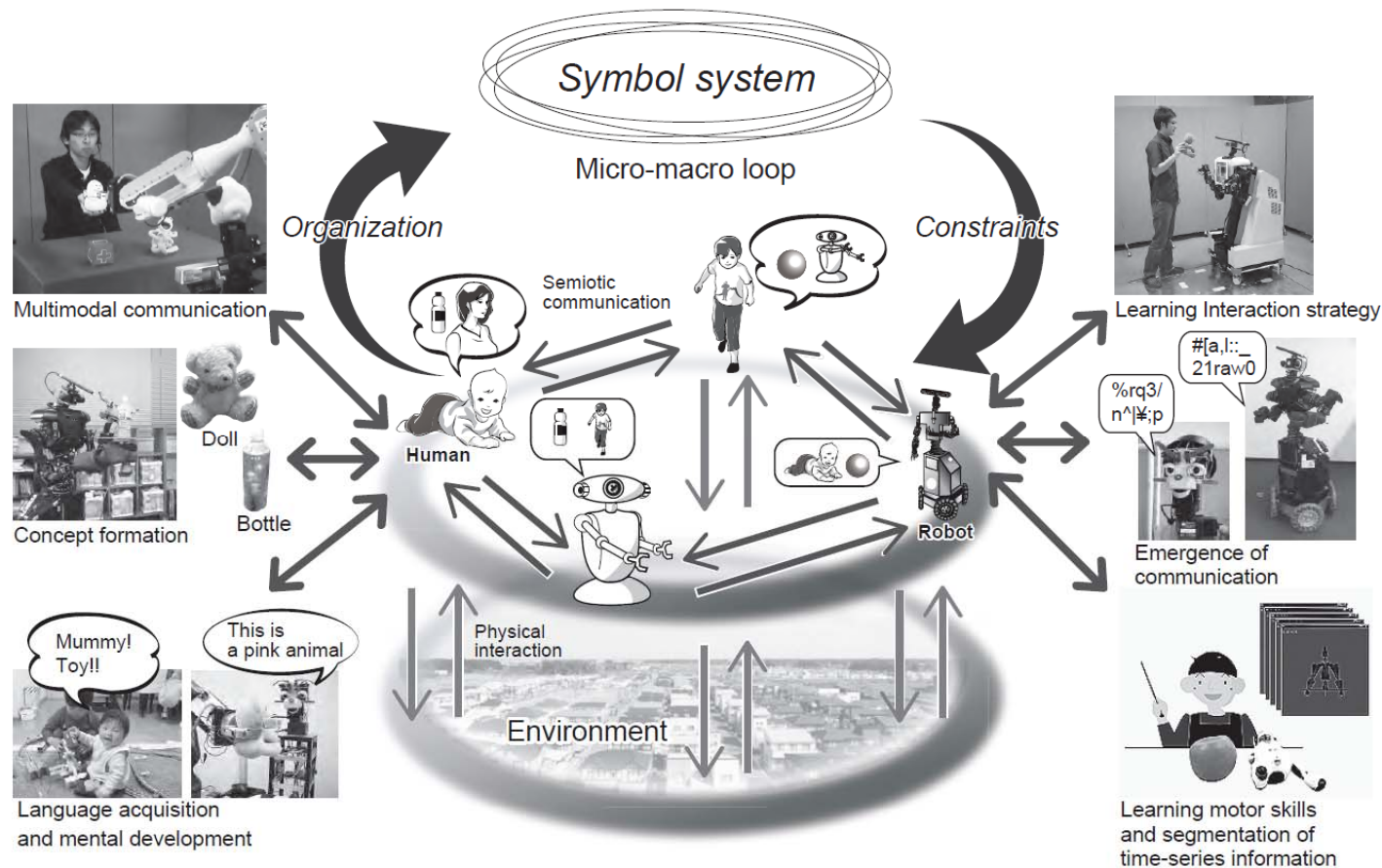


✓ Bottom-up organization of a variety of mutually dependent cognitive functions based on sensorimotor information

SURVEY PAPER

# Symbol emergence in robotics: a survey

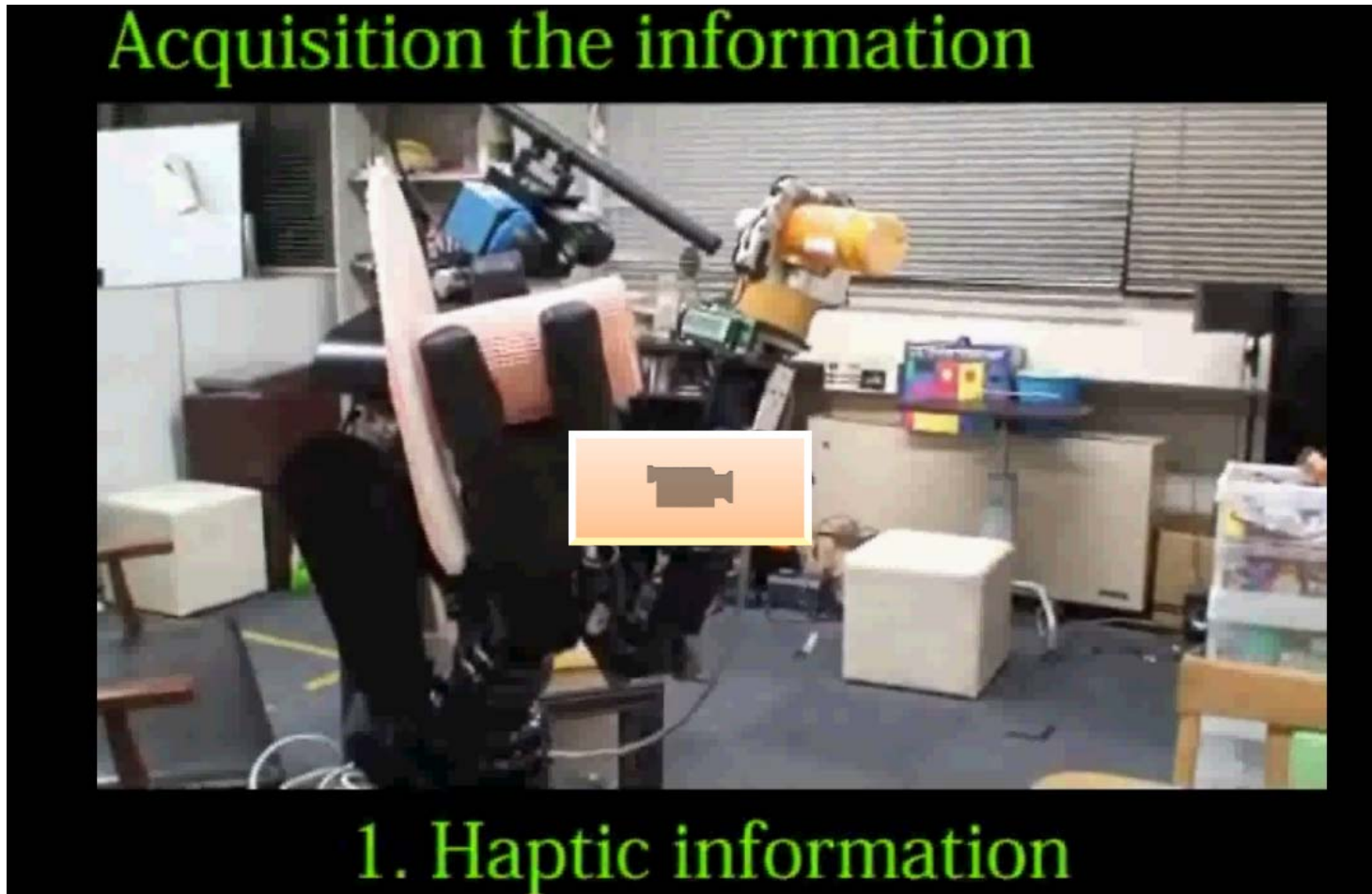
Tadahiro Taniguchi<sup>a</sup>, Takayuki Nagai<sup>b</sup>, Tomoaki Nakamura<sup>b</sup>, Naoto Iwahashi<sup>c</sup>, Tetsuya Ogata<sup>d</sup> and Hideki Asoh<sup>e</sup>



Tadahiro Taniguchi, Takayuki Nagai, Tomoaki Nakamura, Naoto Iwahashi, Tetsuya Ogata, and Hideki Asoh, Symbol Emergence in Robotics: A Survey  
Advanced Robotics, .(2016)DOI:10.1080/01691864.2016.1164622



# Multimodal Categorization and Lexical Acquisition by an Autonomous Robot [Nakamura+ 2009-]

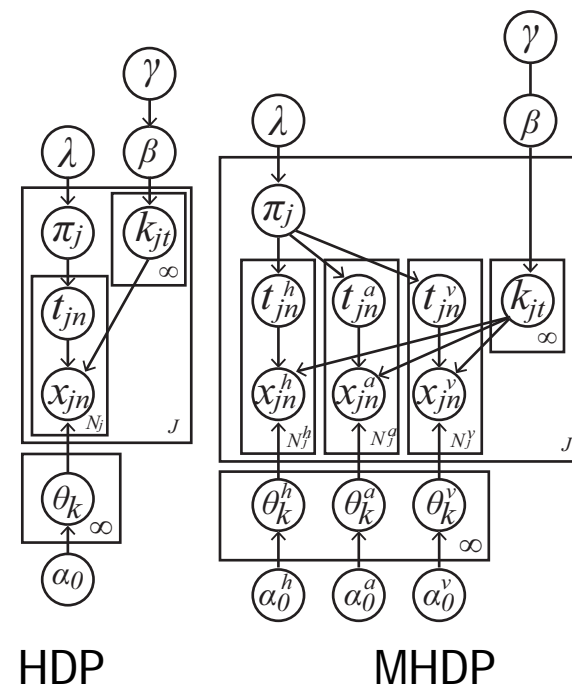


Takaya Araki, Tomoaki Nakamura, Takayuki Nagai, Shogo Nagasaka, Tadahiro Taniguchi, Naoto Iwahashi.  
Online Learning of Concepts and Words Using Multimodal LDA and Hierarchical Pitman-Yor Language Model.  
IEEE/RSJ International Conference on Intelligent Robots and Systems 2012 (IROS 2012), 1623-1630 .(2012)

# Multimodal latent Dirichlet allocation(MLDA) / Hierarchical Dirichlet Processes(MHDP)

[Nakamura+ 2009, 2011]

- The MLDA is a multimodal categorization method that is an extension of the LDA [Blei+ 2014].
- The MLDA was originally proposed for making a robot form **object categories** in an unsupervised manner.
- **Multimodal Hierarchical Dirichlet process (MHDP)** is a nonparametric Bayesian extension of MLDA [Nakamura+ 2011].



HDP [Teh+ 2006] MHDP [Nakamura+ 2011]

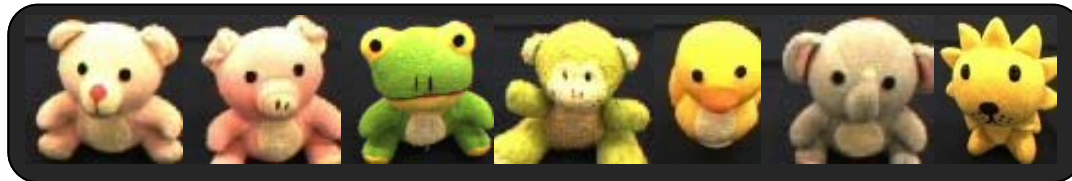
	Observations	Latent variable
LDA/HDP	Words in a document (i.e., Bag of words)	Topic
MLDA/MHDP	Multimodal (visual, auditory, and haptic) features obtained from an object (i.e., Bag of features)	Object category

[Teh+ 2006] Y.W. Teh, M.I. Jordan, M.J. Beal, and D.M. Blei. Hierarchical dirichlet processes. Journal of the American Statistical Association, 101(476):1566-1581, 2006.

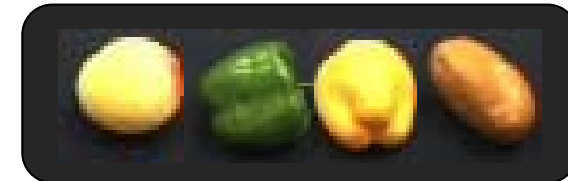
[Nakamura+ 2011] Tomoaki Nakamura, Takayuki Nagai, and Naoto Iwahashi. Multimodal categorization by hierarchical Dirichlet process. In IEEE/RSJ International Conference on Intelligent Robots and Systems, pages 1520-1525, 2011.

# Categorization result based on real-world multimodal sensorimotor information

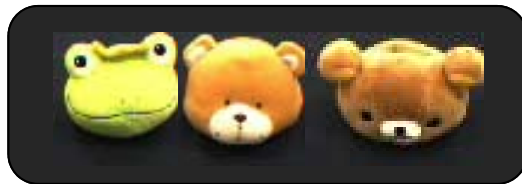
Stuffed animals



Toy vegetables



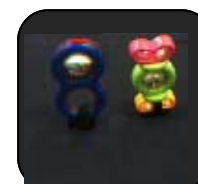
Stuffed animals with a bell



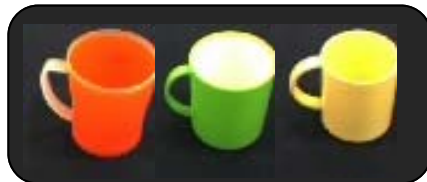
Maracas



Rattles



Cups



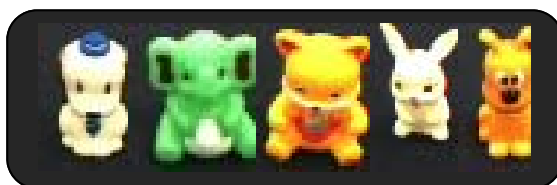
Blocks



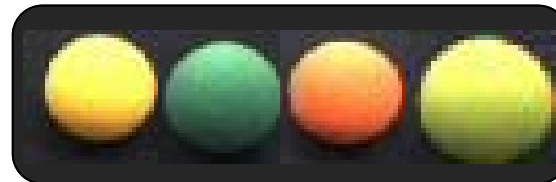
plastic bottles



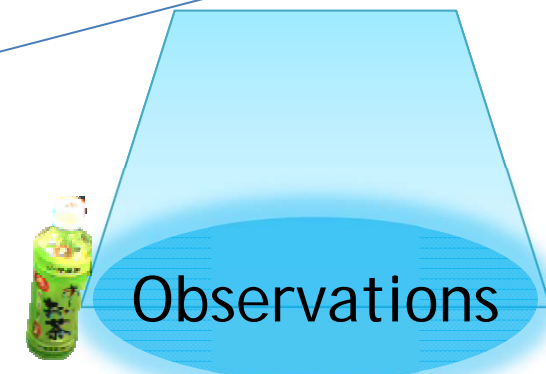
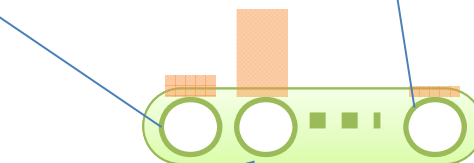
Rubber toys



Sponge balls



*Categorical distribution over topics*



By integrating multimodal information, the robot formed categories represented by latent variables that were similar to most of the human participants.

# Contents

## 1. Introduction

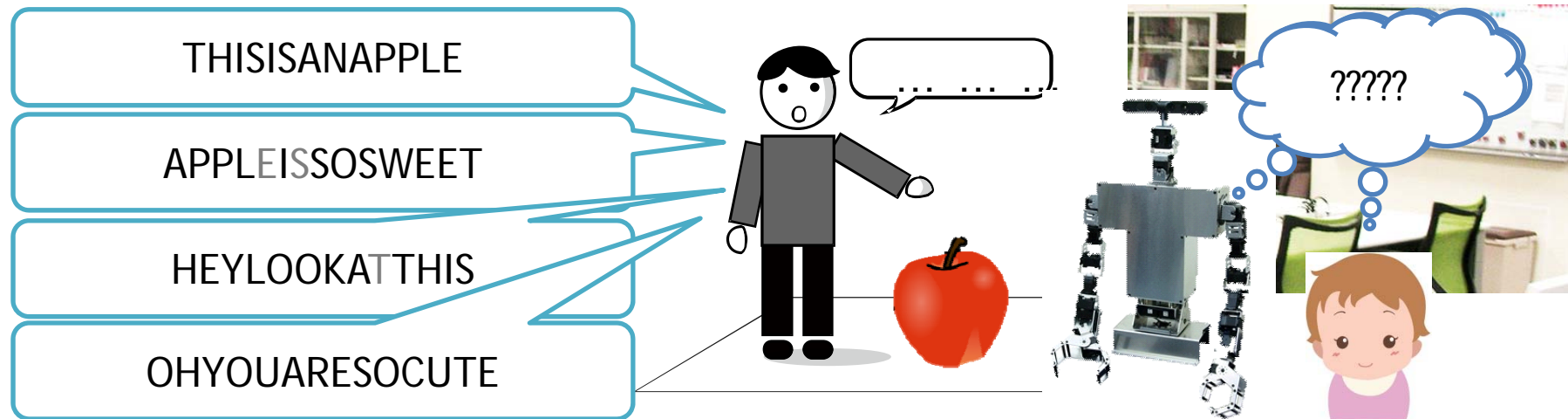
## 2. Lexical acquisition tasks

- A) Direct phoneme and word discovery from speech signals
- B) Simultaneous acquisition of word units and multimodal categories
- C) Online spatial concept acquisition

## 3. Future challenges

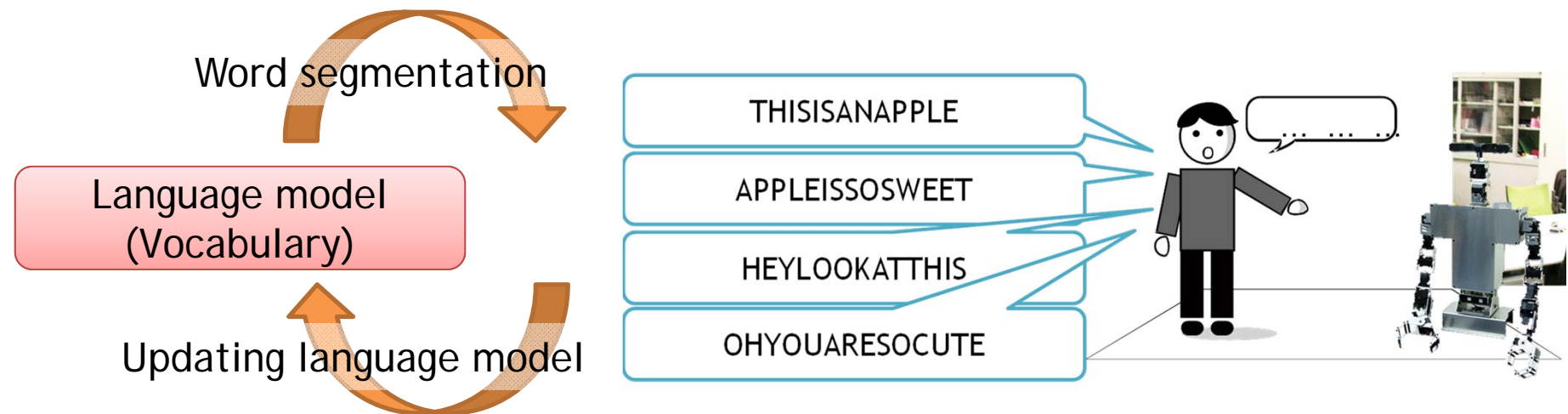
# Word discovery (segmentation) in language acquisition

- ❑ When parents speak to their children, they rarely use “isolated words,” but use continuous word sequences, i.e. sentences.
- ❑ Word and phoneme discovery (segmentation) is a primary task of language acquisition.
- ❑ The child has to perform word segmentation without pre-existing knowledge of vocabulary because children do not know lists of words before they learn.



# Unsupervised word segmentation

- Unsupervised word segmentation
  - No preexisting dictionaries are used.
  - Input data is test data.
  - A nonparametric Bayesian framework for word segmentation [Goldwater+ 09]
  - Unsupervised word segmentation method based on the Nested Pitman-Yor language model (NPYLM) [Mochihashi+ 09].



S. Goldwater, T. L. Griffiths, and M. Johnson, "A Bayesian framework for word segmentation: exploring the effects of context.," *Cognition*, vol. 112, no. 1, pp. 21-54, 2009.

Daichi Mochihashi, Takeshi Yamada, Naonori Ueda. "Bayesian Unsupervised Word Segmentation with Nested Pitman-Yor Language Model". ACL-IJCNLP 2009. pp.100-108. 2009.

# Analysis “Alice in Wonderland”の解析



first, she dreamed of little Alice herself, and once again the tiny hands were clasped up on her knee, and the bright eager eyes were looking up into hers -- she could hear the very tones of her voice, and see that queer little toss of her head to keep back the wandering hair that would always get into her eyes -- and still as she listened, or seemed to listen, the whole place around her became alive the strange creatures of her little sister's dream. the long grass rustled at her feet as the white rabbit hurried by -- the frightened mouse splashed his way through the neighbouring pool -- she could hear the rattle of the tea cups as the March hare and his friends shared their never-ending meal, and the shrill voice of the queen...



first, she dream ed of little Alice herself ,and once again the tiny hand s were clasped upon her knee ,and the bright eager eyes were looking up into hers -- she could hear the very tone s of her voice , and see that queer little toss of her head to keep back the wandering hair that would always get into hereyes -- and still as she listened , or seemed to listen , the whole place a round her became alive the strange creatures of her little sister 's dream. the long grass rustled at her feet as the white rabbit hurried by -- the frightened mouse splashed his way through the neighbour ing pool -- she could hear the rattle of the tea cups as the March hare and his friends shared their never -ending me a l ,and the ...

# Challenges in Real-world Unsupervised Word Discovery Tasks

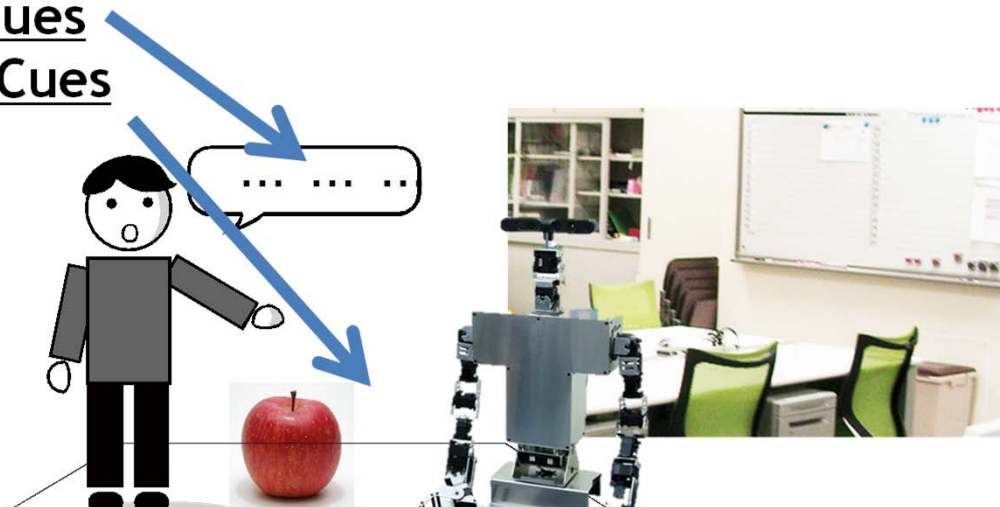
- NPYLM presumes that the target document (sentences) is transcribed without errors. If there are **phoneme recognition errors**, its performance becomes dramatically worse.

- A) Mitigating negative effects of phoneme recognition errors
- B) Learning a phoneme system (acoustic model).
- C) Grounding discovered words

Prosodic Cues

Distributional Cues

Co-occurrence Cues



[Saffuran 1996]



# Contents

1. Introduction

2. Lexical acquisition tasks

A) Direct phoneme and word discovery from speech signals

B) Simultaneous acquisition of word units and multimodal categories

C) Online spatial concept acquisition

3. Future challenges

# Simultaneous acquisition of phoneme and language models

- ✓ In most related studies about unsupervised word discovery, they used a pre-existing phoneme model and did not make a robot learn a phoneme system (an acoustic model).
- ✓ There are still few studies about unsupervised simultaneous learning of phoneme and language models from speech signals [Kamper+ 15, Lee+ 15].
- ✓ It has been suggested that the learning processes of acoustic and language models (phoneme and lexical systems) are mutually dependent.

Making full use of the **Distributional Cues** directly from speech signals

Prosodic Cues

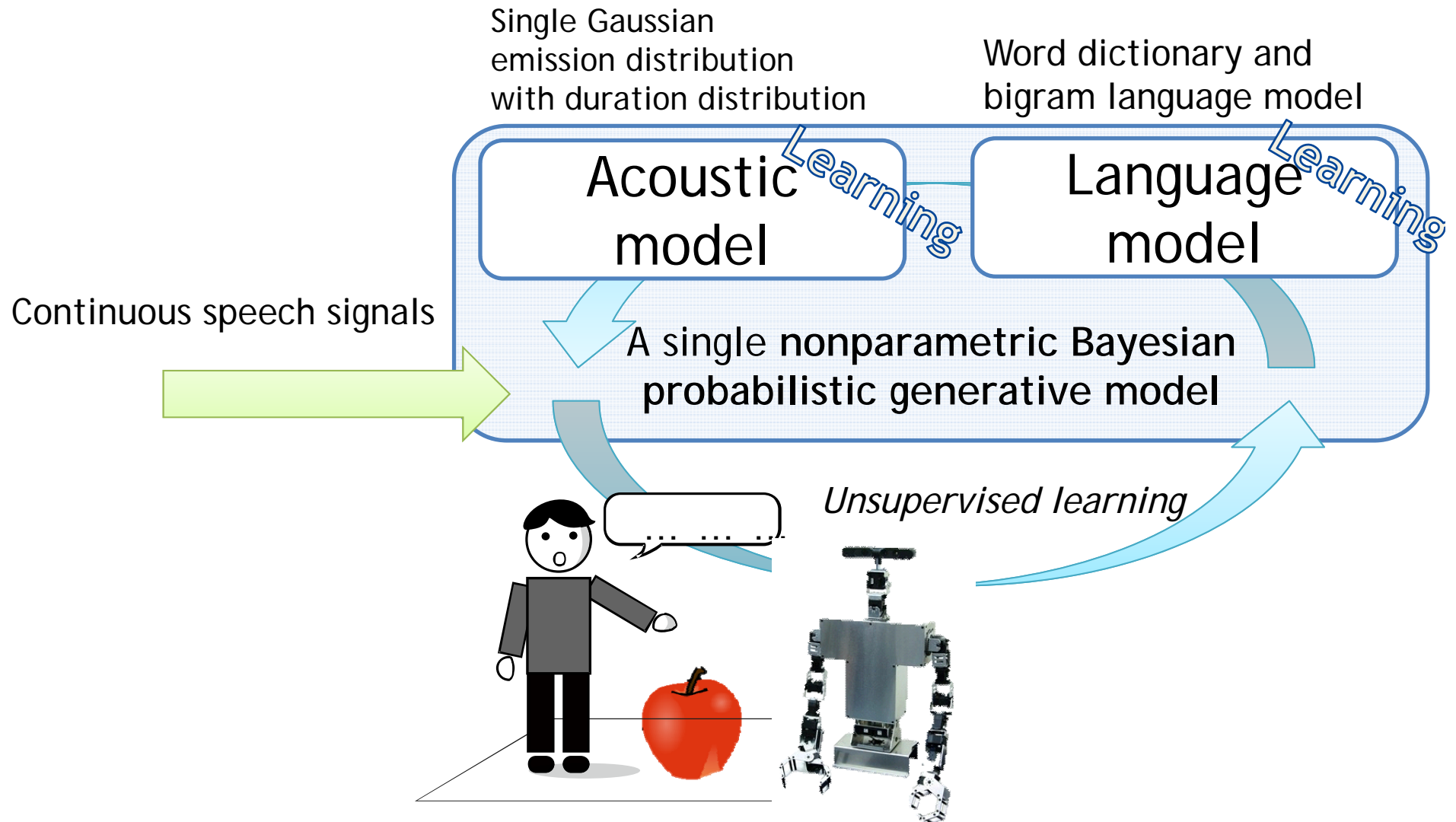
Co-occurrence Cues



H. Kamper, A. Jansen, and S. Goldwater, "Fully Unsupervised Small-Vocabulary Speech Recognition Using a Segmental Bayesian Model," in INTERSPEECH 2015, 2015.

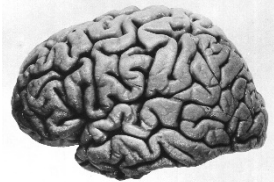
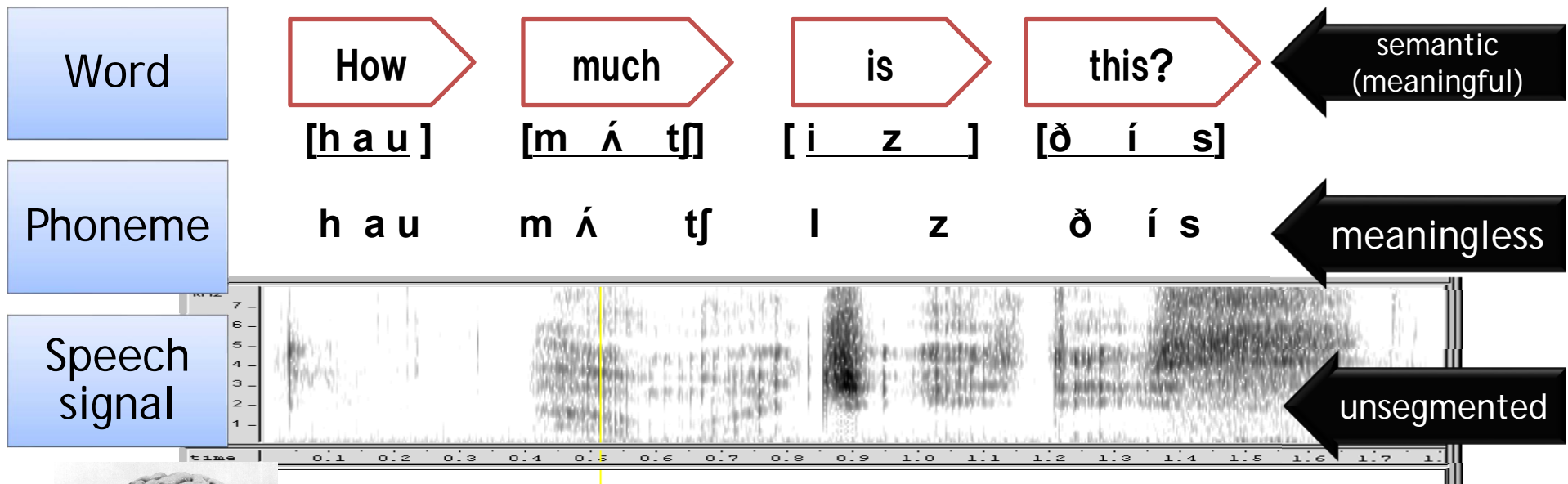
C.-y. Lee, T. J. O. Donnell, and J. Glass, "Unsupervised Lexicon Discovery from Acoustic Input," Transactions of the Association for Computational Linguistics, vol. 3, pp. 389-403, 2015.

# Direct phoneme and word discovery from speech signals



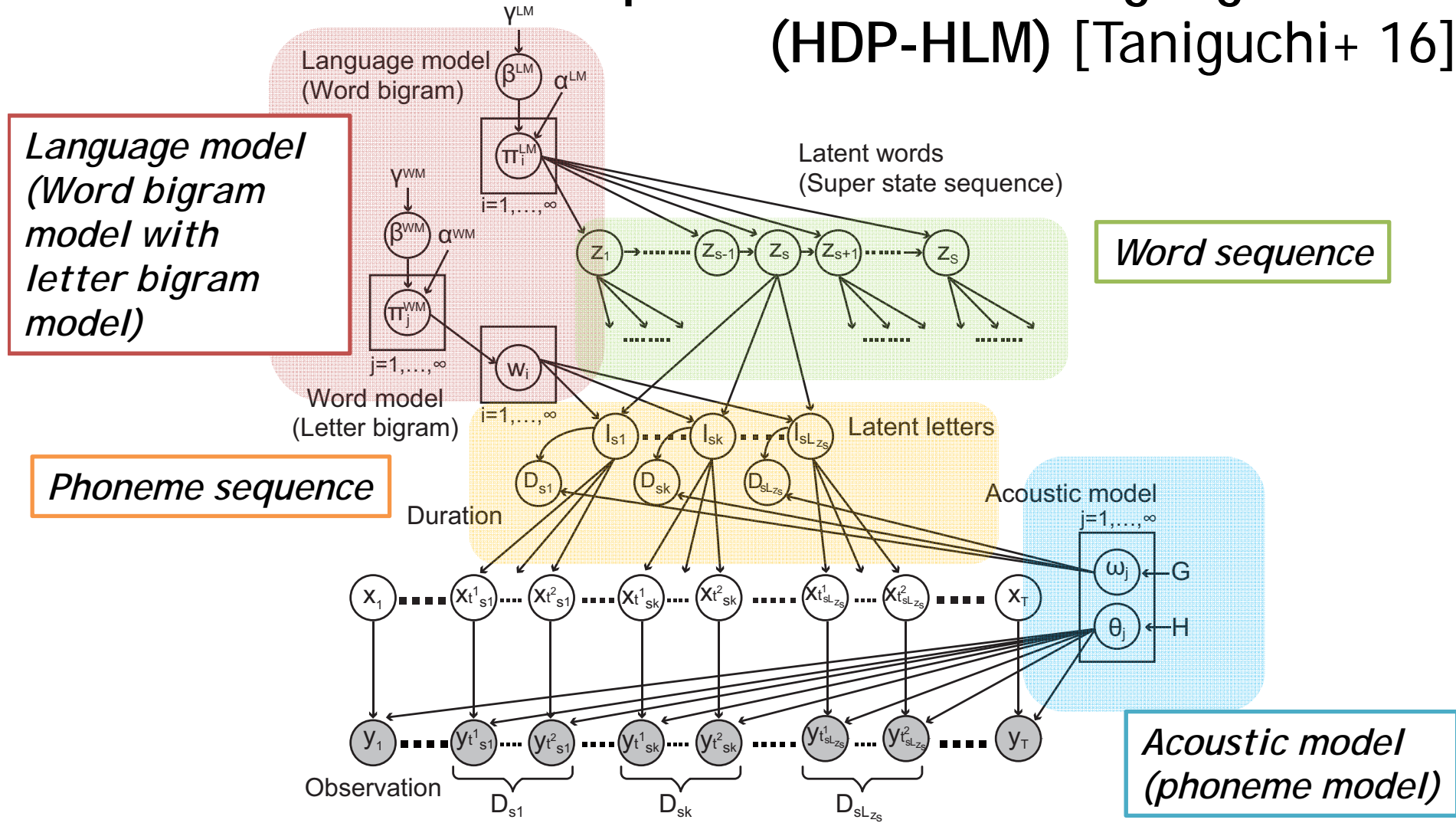
# Double articulation structure in semiotic data

- Semiotic time-series data often has double articulation
  - Speech signal is a continuous and high-dimensional time-series.
  - Spoken sentence is considered a sequence of **phonemes**.
  - Phonemes are grouped into words, and people give them meanings.



Does the human brain have a special capability to analyze double articulation structures embedded in time-series data?

# Hierarchical Dirichlet process hidden language model (HDP-HLM) [Taniguchi+ 16]

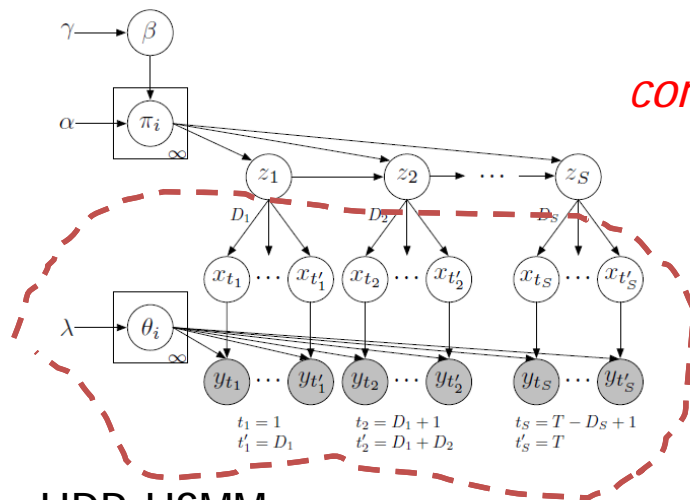


A probabilistic generative model for time-series data  
having double articulation structure

Tadahiro Taniguchi, Shogo Nagasaka, Ryo Nakashima, Nonparametric Bayesian Double Articulation Analyzer for Direct Language Acquisition from Continuous Speech Signals, IEEE Transactions on Cognitive and Developmental Systems. (2016)

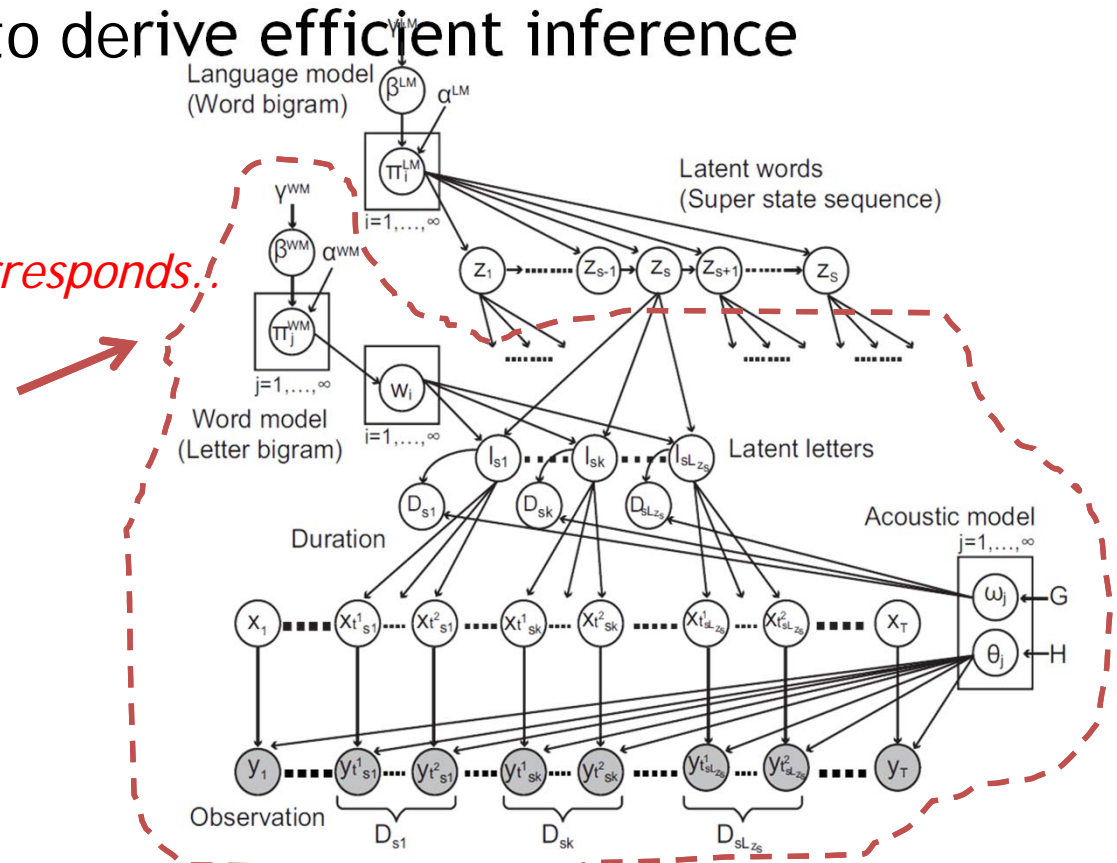
# HDP-HLM as an extension of HDP-HSMM

- ✓ HDP-HLM can be regarded as an extension of HDP-HSMM [Johnson'13]
- ✓ This property helps us to derive efficient inference procedure.



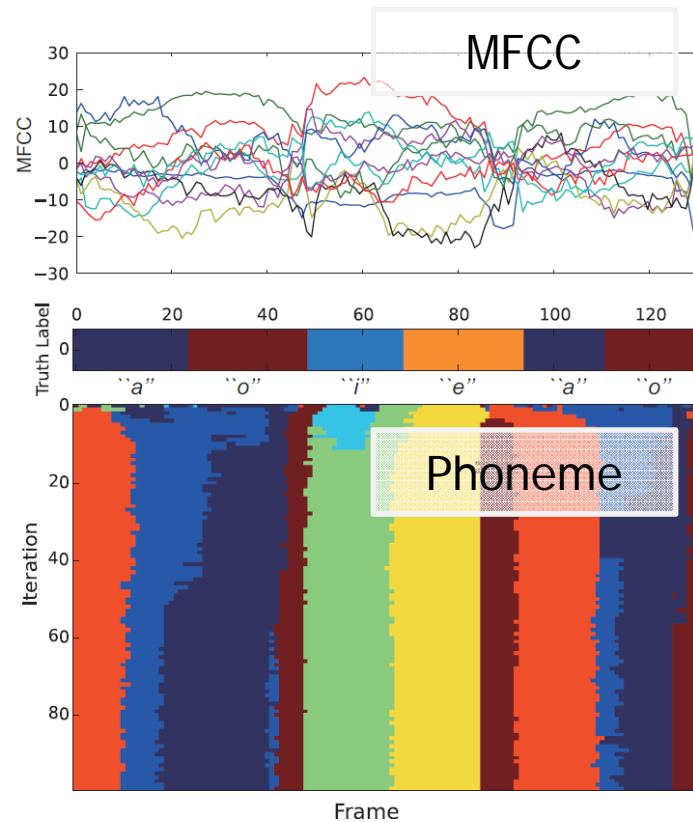
HDP-HSMM  
(hierarchical Dirichlet process hidden semi-Markov model)

corresponds!



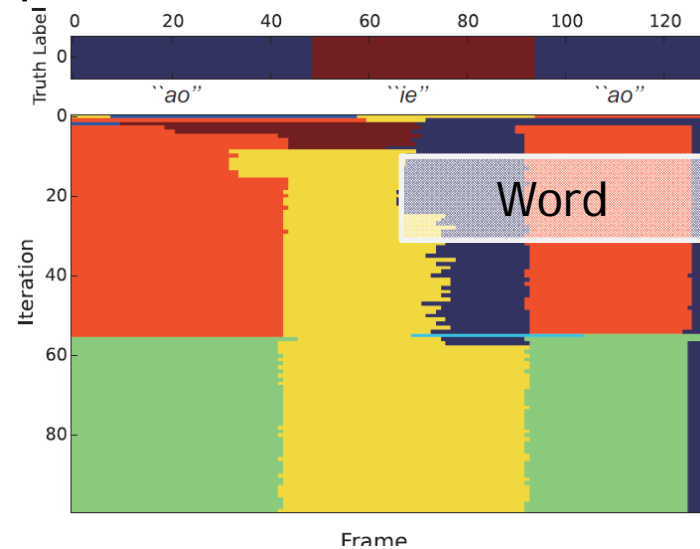
# Evaluation experiment using artificial 2 or 3 words sentences with Japanese five vowels

- ✓ Five artificial words {aioi, aue, ao, ie, uo} prepared by connecting five Japanese vowels.
- ✓ 30 sentences (25 two-word and 5 three-word sentences) are prepared and each sentence is recorded twice by four Japanese speakers.
- ✓ MFCC (frame size = 25ms, shift = 10ms, frame rate 100hz)
  - \* HDP-HLM are trained separately for each speaker.



✓ The inference procedure could gradually estimate the boundaries of words and phonemes.

ex) ao-ie-ao



Iteration of Gibbs sampler

# Experimental results

Table 1. ARI for estimated latent letters and words.

Method	Letter ARI	Word ARI	AM	LM
NPB-DAA with DSAE (MAP)	<b>0.589</b>	<b>0.705</b>		
NPB-DAA with DSAE	0.426	<b>0.398</b>		
NPB-DAA (MAP)	<b>0.612</b>	0.328		
NPB-DAA	0.551	0.359		
Conventional DAA	0.584	0.072		
Julius				
(GMM monophone + phoneme dictionary + NPYLM)	0.483	0.315	✓	
(GMM monophone + phoneme dictionary + lattice)	0.483	0.26	✓	
(GMM monophone + phoneme dictionary + lattice + LM)	0.483	0.30	✓	
(GMM triphone + phoneme dictionary + latticelm)	0.269	0.203	✓	
Julius				
(DNN triphone + phoneme dictionary + NPYLM)	0.634	0.333	✓	
Julius				
(GMM monophone + word dictionary)	0.565	0.548	✓	✓
Julius				
(GMM triphone + word dictionary)	0.516	0.636	✓	✓
Julius				
(DNN triphone + word dictionary)	0.675	0.779	✓	✓

The NPB-DAA with DSAE even outperformed MFCC-based off-the-shelf speech recognition system.



# Contents

1. Introduction

2. Lexical acquisition tasks

A) Direct phoneme and word discovery from speech signals


B) Simultaneous acquisition of word units and multimodal categories

C) Online spatial concept acquisition

3. Future challenges

# Making use of co-occurrence cue

- ✓ To detect the co-occurrence of an object and a phrase, the robot has to form the category of the object beforehand, or simultaneously.

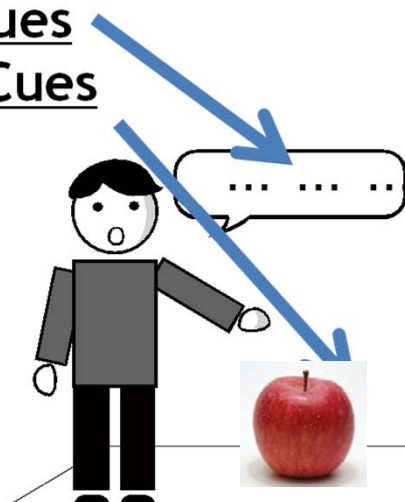
Example,  and, an *“apple”*

- ✓ How can a robot form “object categories” without knowing the names of objects?

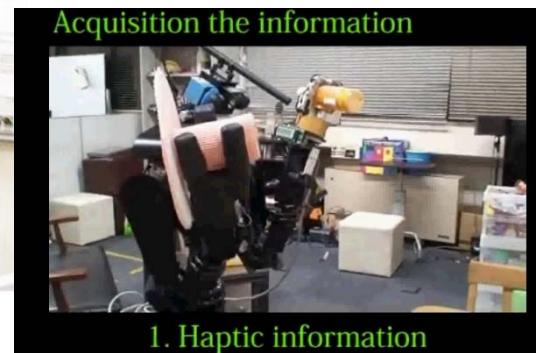
Prosodic Cues

Distributional Cues

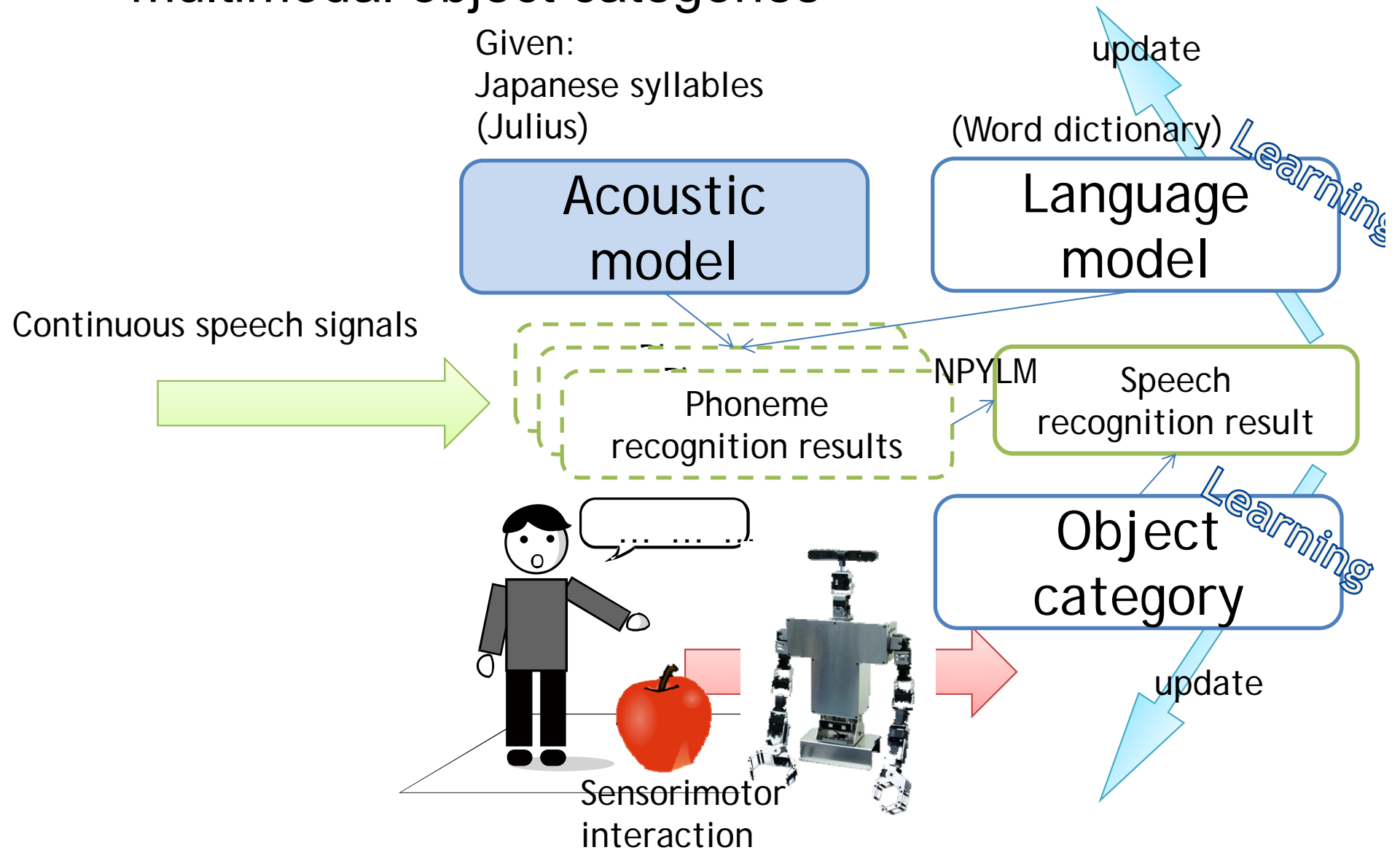
Co-occurrence Cues



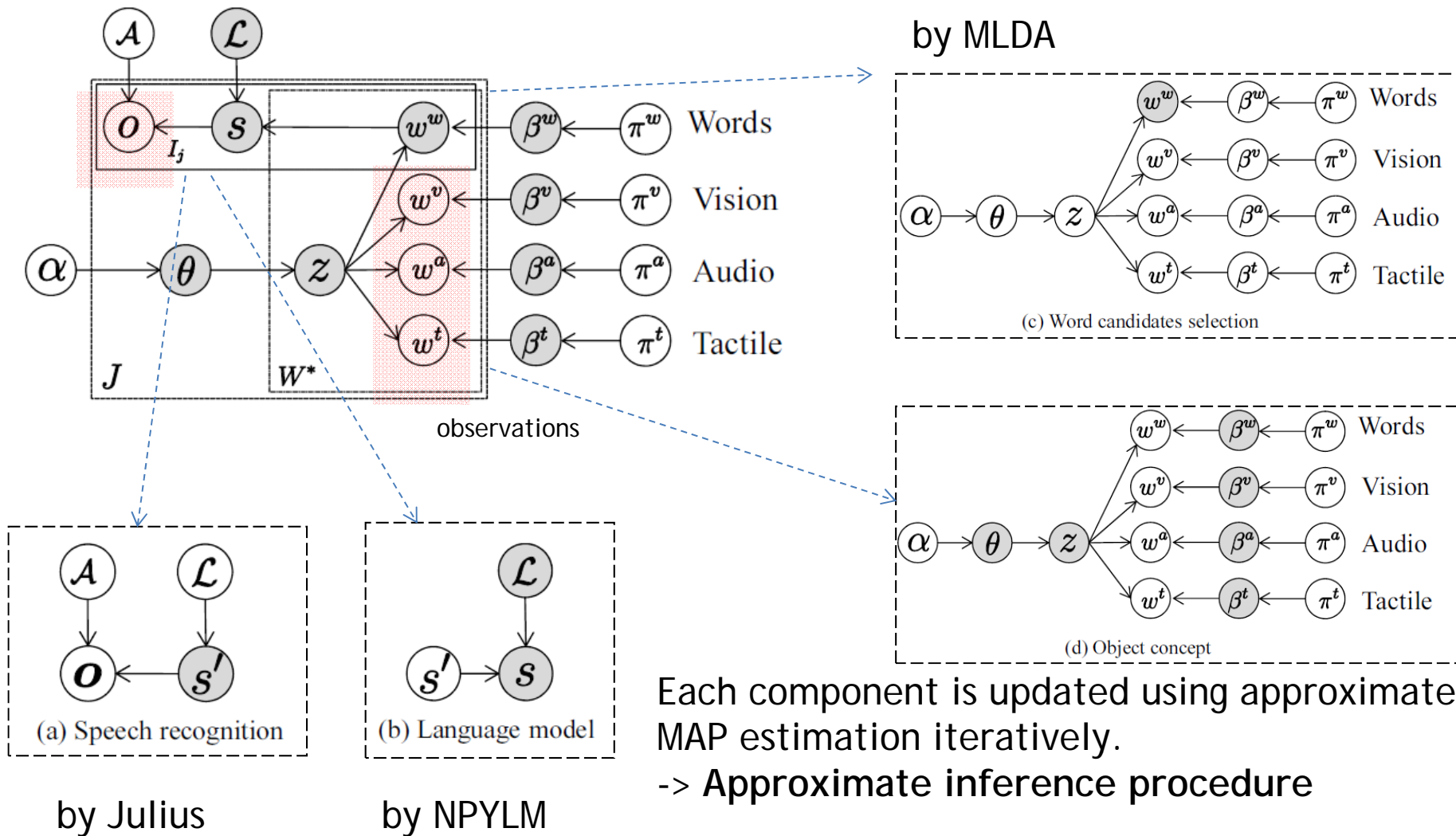
Multimodal object categorization



# Simultaneous acquisition of word units and multimodal object categories



# Probabilistic generative model for simultaneous acquisition of word units and multimodal object categories

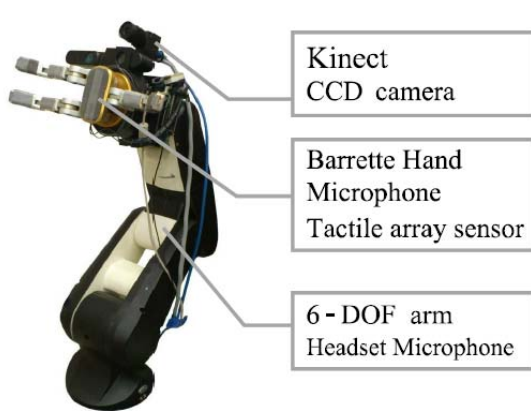


Each component is updated using approximate MAP estimation iteratively.

-> Approximate inference procedure

Tomoaki Nakamura, Takayuki Nagai, Kotaro Funakoshi, Shogo Nagasaka, Tadahiro Taniguchi, and Naoto Iwahashi, Mutual Learning of an Object Concept and Language Model Based on MLDA and NPYLM, 2014 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS'14), 600 - 607 .(2014)

# Overview of experiment and results



Robot



This is a red spray can. (ko re wa a ka i su pu re e ka N)  
 This makes a sound when shaken. (ko re wa o to ga shi ma su)  
 This is made of metal and is hard. (ko re wa ki N zo ku de de ki te i te ka ta i)



A green plushie of a frog. (mi do ri no ka e ru no nu i gu ru mi)  
 This is soft. (ko re wa ya wa ra ka i)  
 This is an animal. (ko re wa do u bu tsu)



A green plastic bottle. (mi do ri no pe tto bo to ru)  
 This is green tea. (ko re wa ryo ku cha)



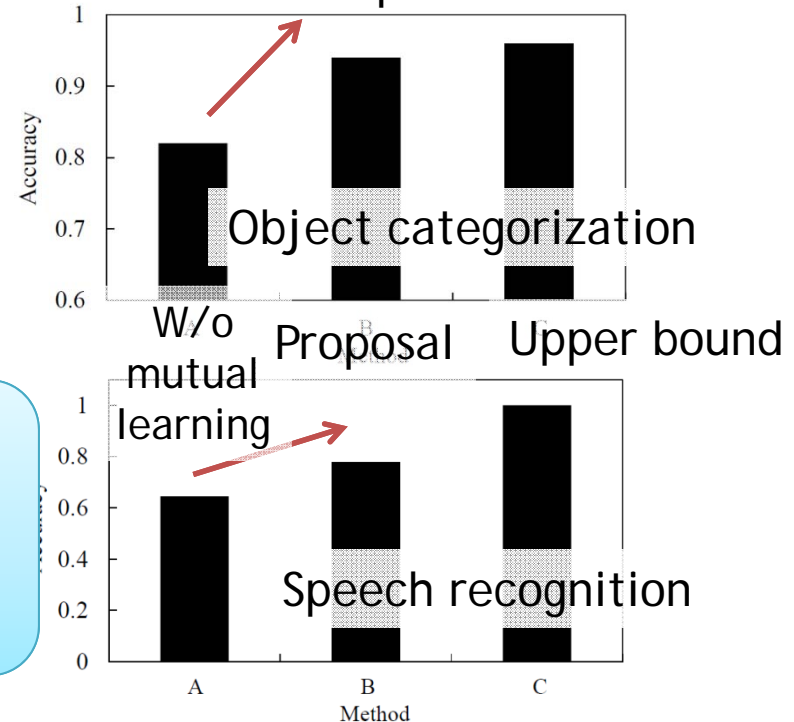
(a)

(b)

(c)

Obtaining multimodal sensory information

Example sentences used in the experiments



Both object categorization and speech recognition performances increased using co-occurrence cues.

# Contents

1. Introduction

2. Lexical acquisition tasks

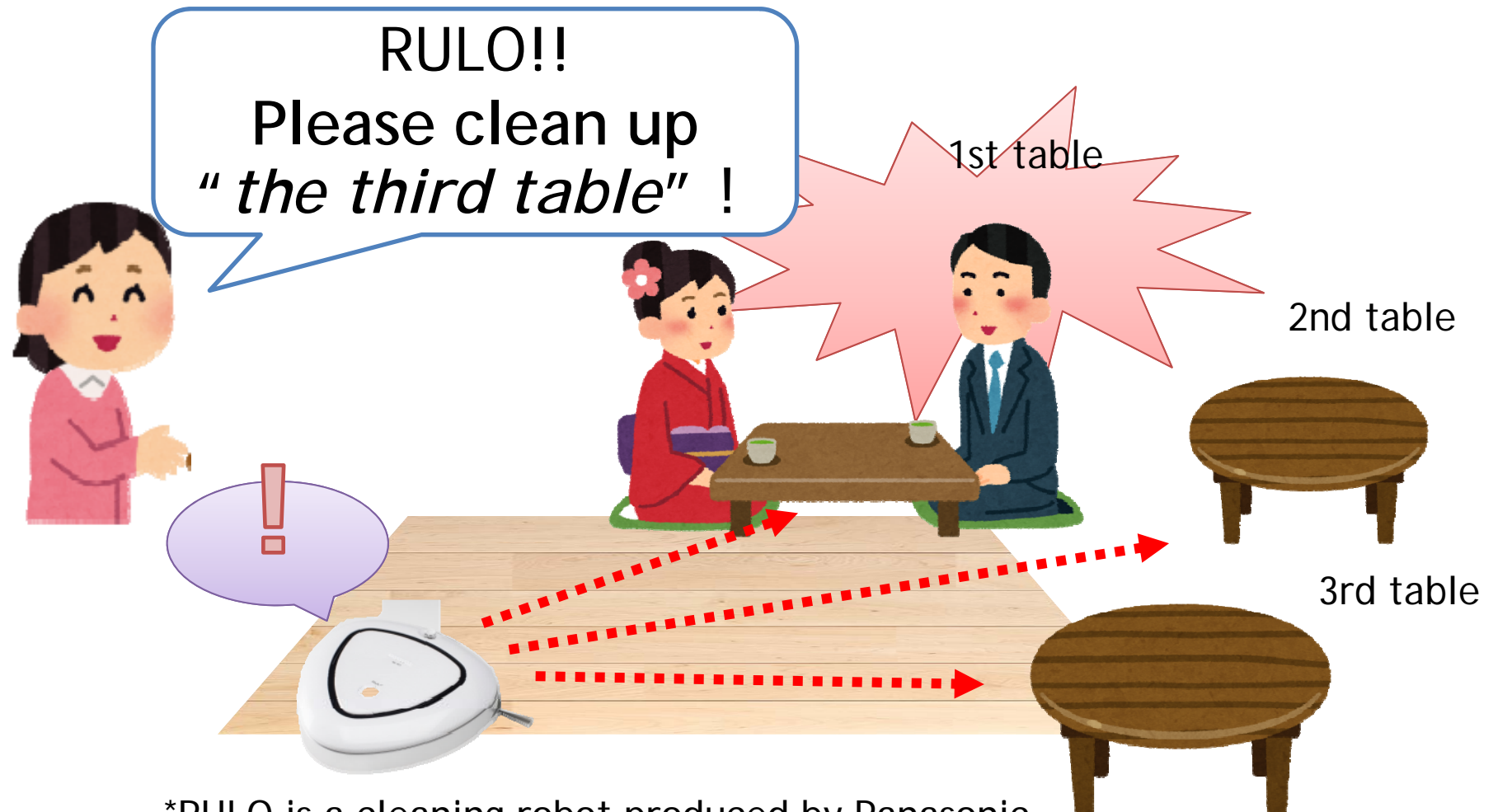
A) Direct phoneme and word discovery from speech signals

B) Simultaneous acquisition of word units and multimodal categories

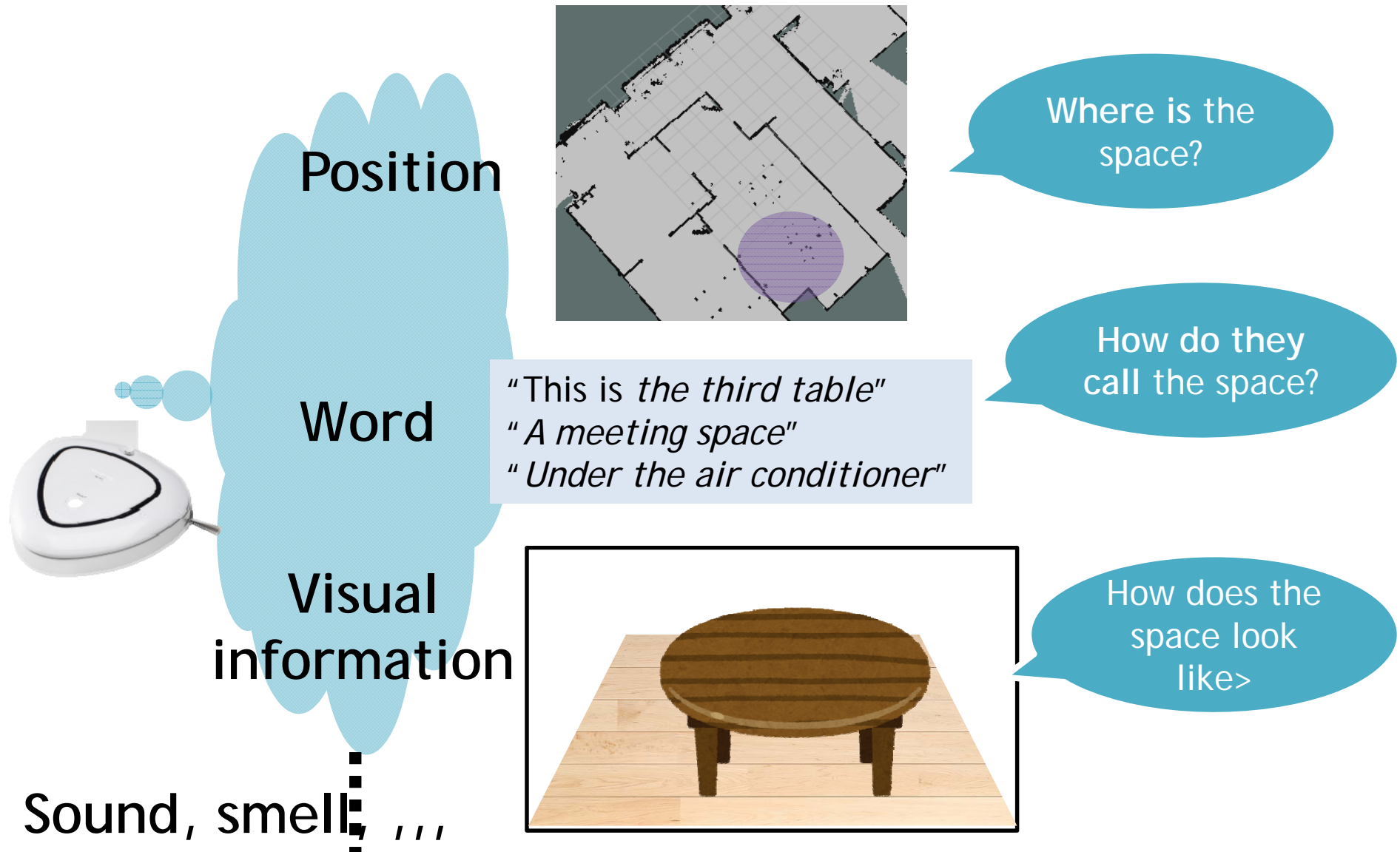
C) Online spatial concept acquisition

3. Future challenges

# Learning place name is very important for home service robot



# Spatial concept is multimodal

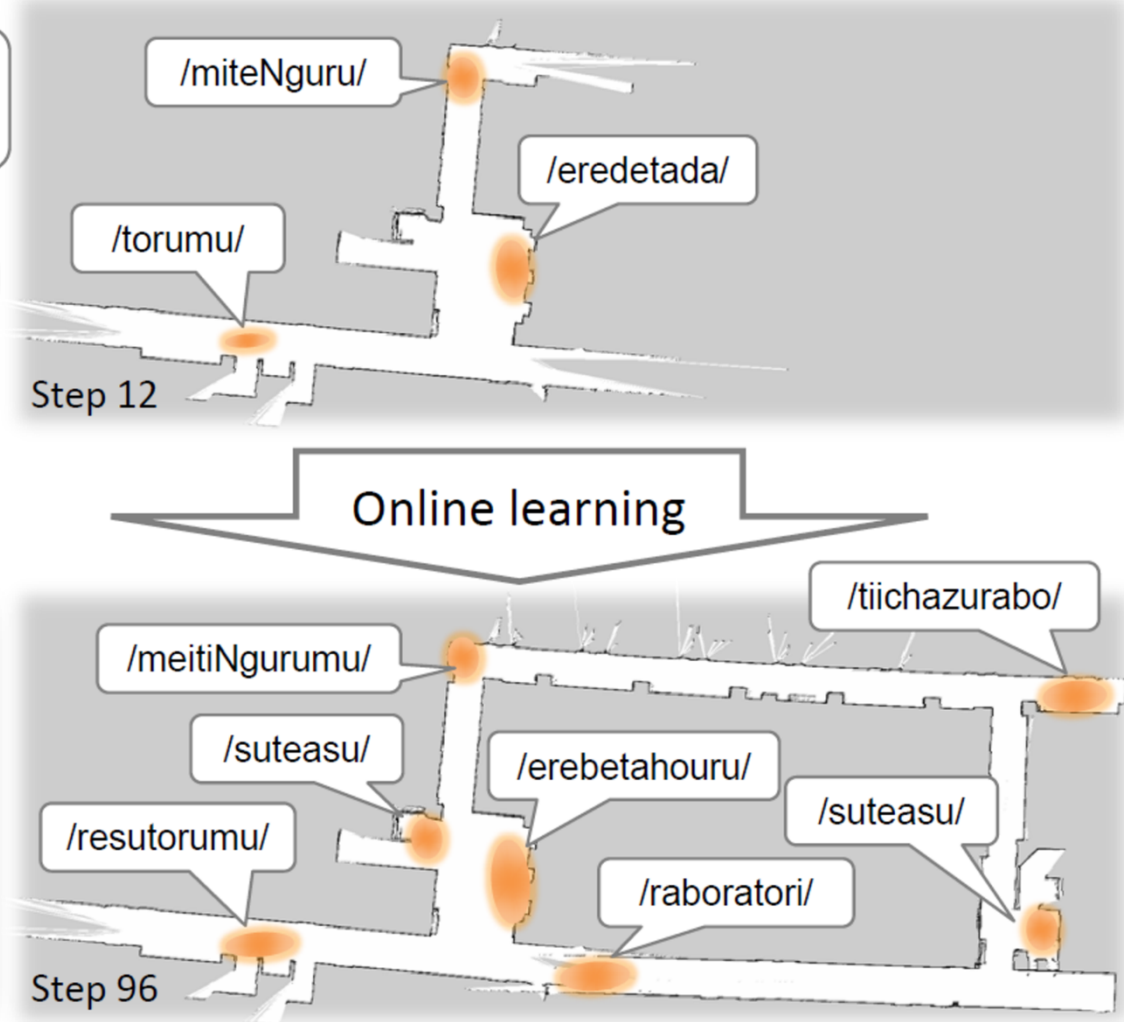
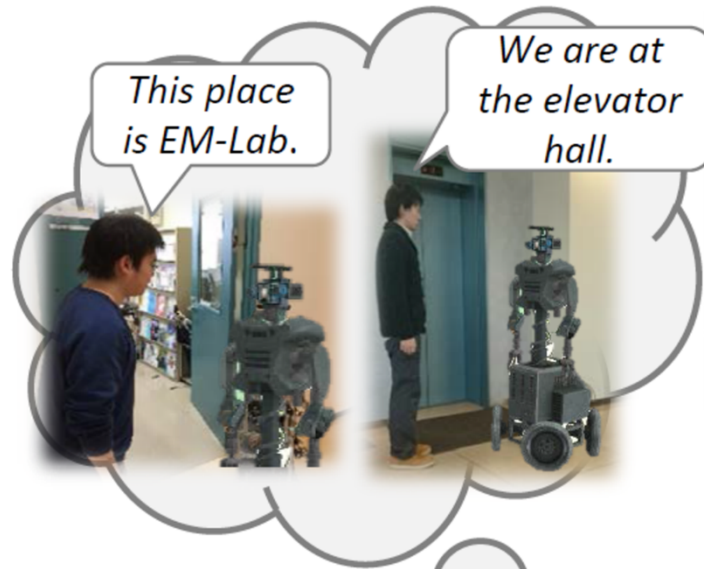




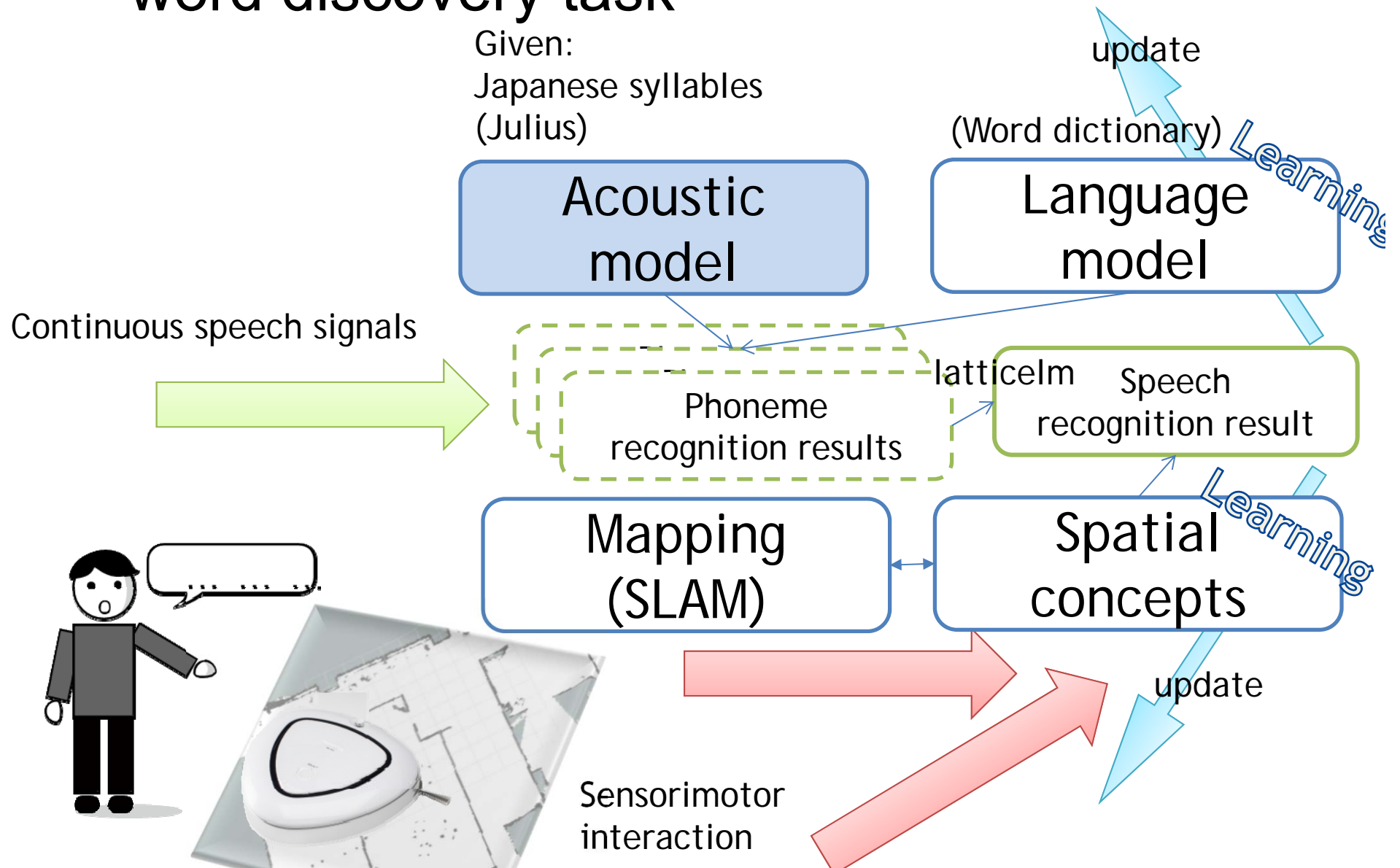
# Online spatial concept acquisition method

## SpCoSLAM [Taniguchi+ 2017]

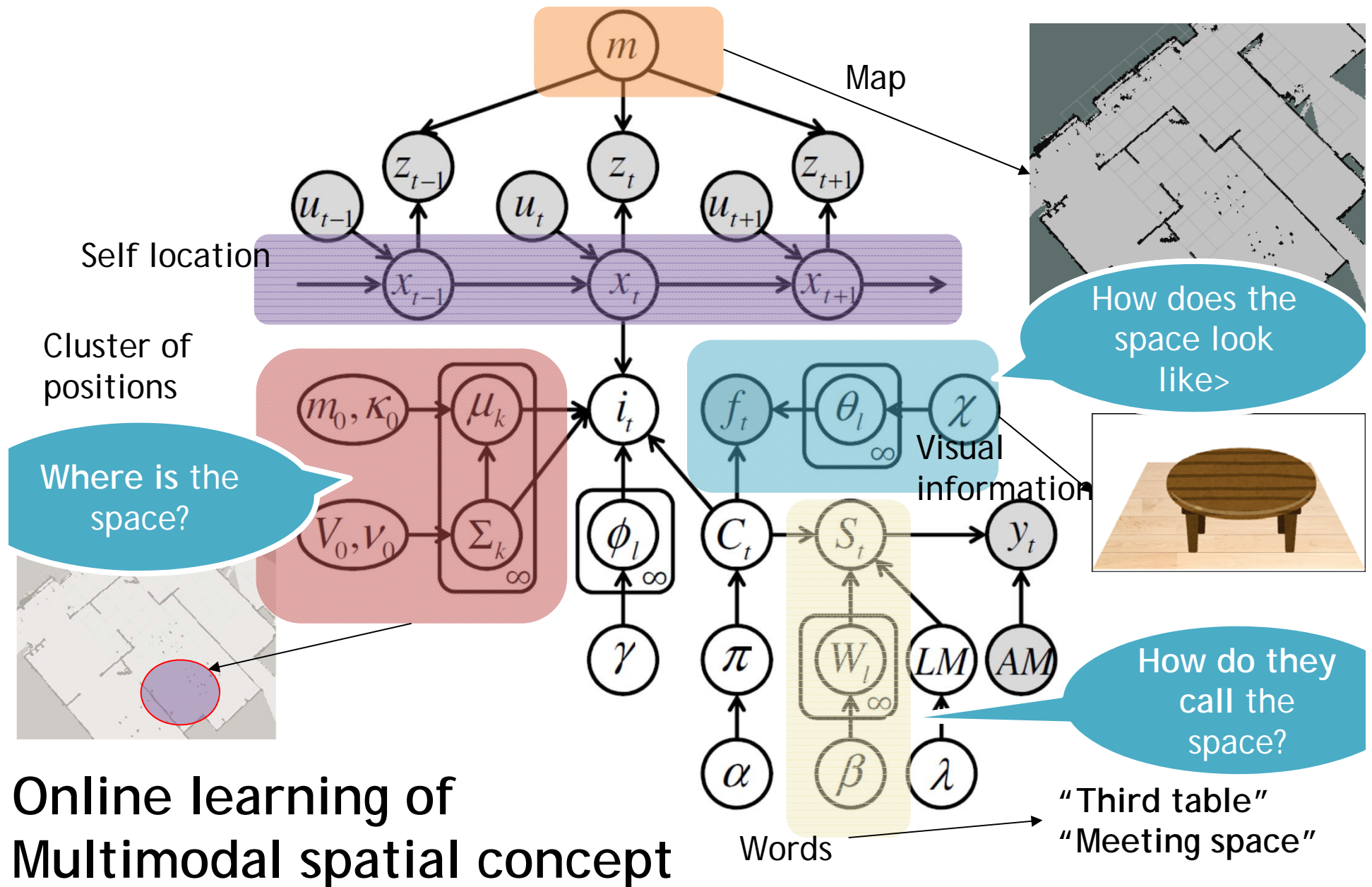
(including word discovery task)



# Online spatial concept acquisition with word discovery task



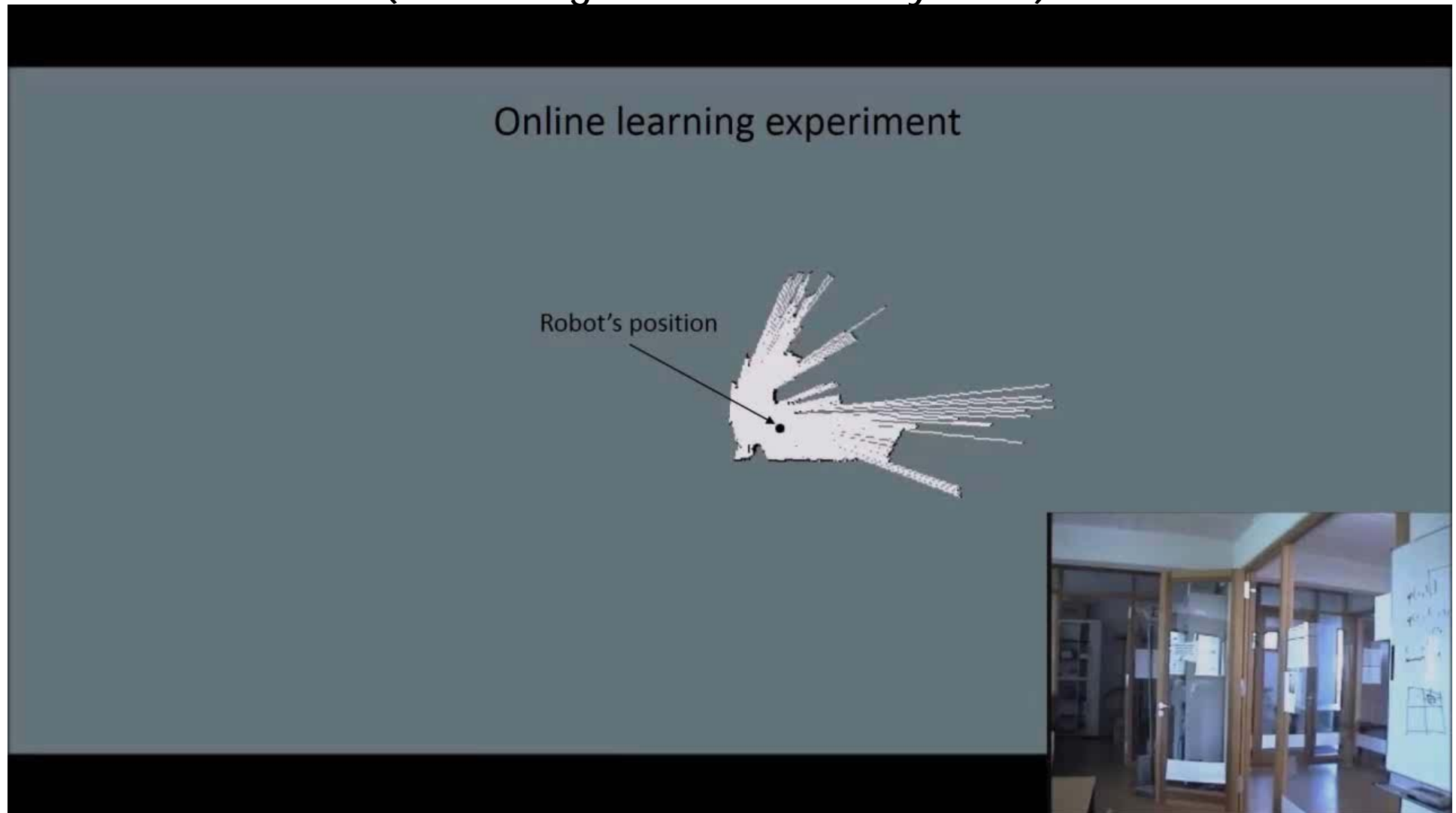
# Graphical model of SpCoSLAM



Online learning of Multimodal spatial concept

"Third table"  
"Meeting space"

# Online spatial concept acquisition method **SpCoSLAM** [Taniguchi+ 2017] (including word discovery task)



Akira Taniguchi, Yoshinobu Hagiwara, Tadahiro Taniguchi and Tetsunari Inamura, Online Spatial Concept and Lexical Acquisition with Simultaneous Localization and Mapping, IEEE IROS 2017 (submitted)

# Contents

1. Introduction

2. Lexical acquisition tasks

A) Direct phoneme and word discovery from speech signals

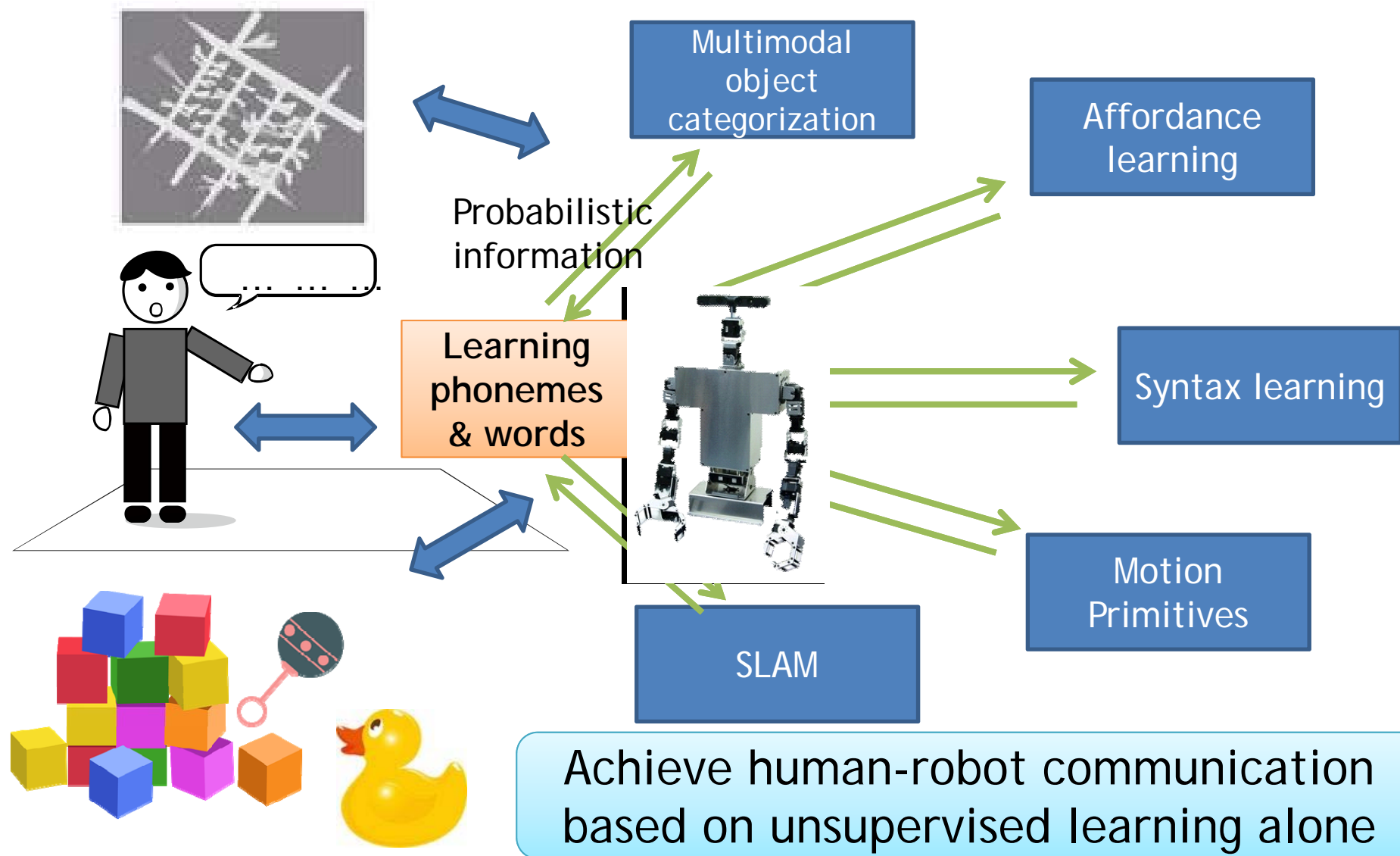
B) Simultaneous acquisition of word units and multimodal categories

C) Online spatial concept acquisition

3. Future challenges

# Current challenge

Unsupervised learning of lexicons grounded via a robot's sensorimotor information



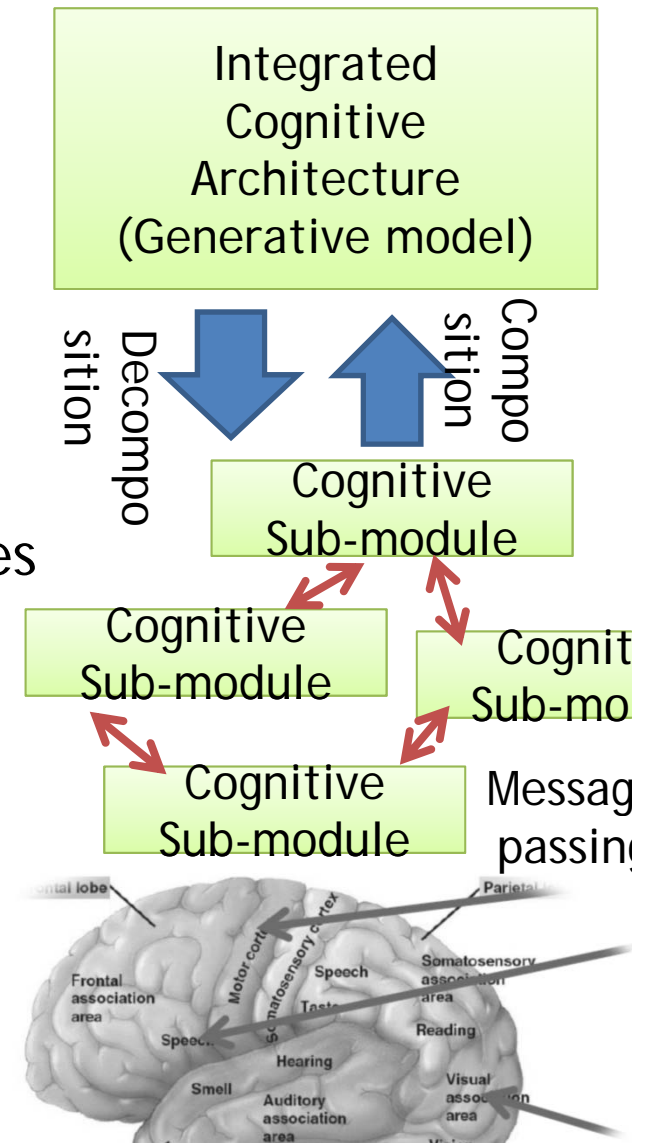
# Employing deep learning for unsupervised machine learning for symbol emergence in robotics

- ❑ **Complex emission distribution and automatic feature extraction**
  - ✓ Structured Variational Auto Encoding(e.g., GMM+VAE [Johnson 2016] )
- ❑ **Using deep learning in inference procedure of probabilistic generative models**
  - ✓ Amortized inference
- ❑ **Modeling language using recurrent neural network**
  - ✓ LSTM, GRU, and so on.

Bayesian deep learning

# Scaling up.....

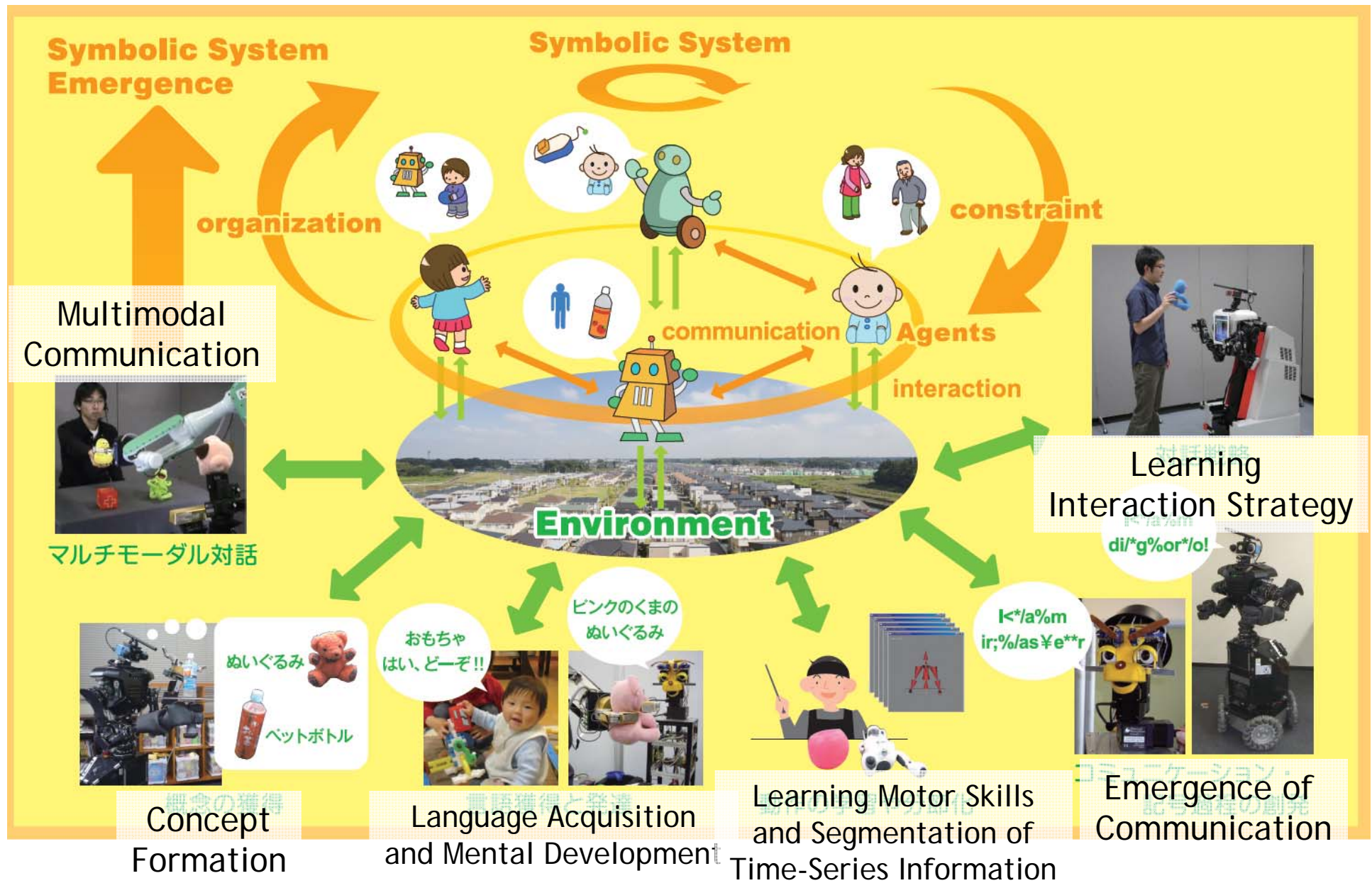
- ❑ Software engineering problems: to deal with complex (huge) graphical models representing mutually dependent cognitive modules.
  - ✓ Developing independent cognitive modules as probabilistic generative models and integrate them into an integrated cognitive module.
  - ✓ Using probabilistic programming. (e.g. Anglican, Venture, Edward)
  - ✓ Employing automatic gradient (e.g. Tensorflow, Autograd) for variational inference



Develop a developmental cognitive model by referring to human whole-brain architecture



# Symbol Emergence System



# Information

## Special Thanks

- Ritsumeikan University
  - R. Nakashima, S. Nagasaka, A. Taniguchi,
- UEC
  - T. Nagai, T. Nakamura, T. Araki, Y. Ando
- Okayama Pref. Univ.
  - N. Iwahashi



Email: [taniguchi@ci.ritumei.ac.jp](mailto:taniguchi@ci.ritumei.ac.jp)

Twitter: @tanichu

Facebook: Tadahiyo Taniguchi

HP: <http://www.tanichu.com>

Funded by



## The 2<sup>nd</sup> Workshop on Machine Learning Methods for High-Level Cognitive Capabilities in Robotics 2017 @IROS2017



<b>Workshop Home</b>	<b>Workshop Home</b>
Objectives	
Organizers	
Topics of interest	
▼ Accepted papers	
Instruction for presentations	
Call for papers	
Invited speakers	
Program	
Sitemap	

**Welcome to the 2nd Workshop on Machine Learning Methods for High-Level Cognitive Capabilities in Robotics**

Integrating high-level and low-level cognitive capabilities is essential for developing robotic systems that can adaptively act in our daily environment in active collaboration with humans.

The main objective of this workshop is to share knowledge about the state-of-the-art machine learning methods that contribute to modeling high-level cognitive capabilities in robotics and to exchange views among cutting-edge robotics researchers who are interested in adaptive high-level cognitive capabilities in robotics.