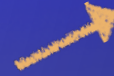ACM Multimedia 2018

# Knowledge-aware Multimodal Dialogue Systems

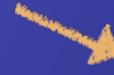Lizi Liao[1], Yunshan Ma[1], Xiangnan He[1], Richang Hong[2], Tat-Seng Chua[1]

[1]National University of Singapore, [2]Hefei University of Technology

24 October 2018
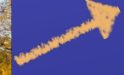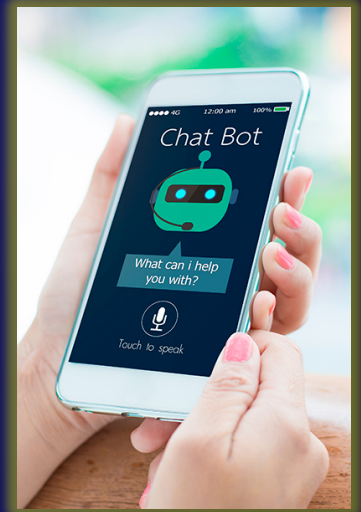
# Why Multimodal Dialogue?
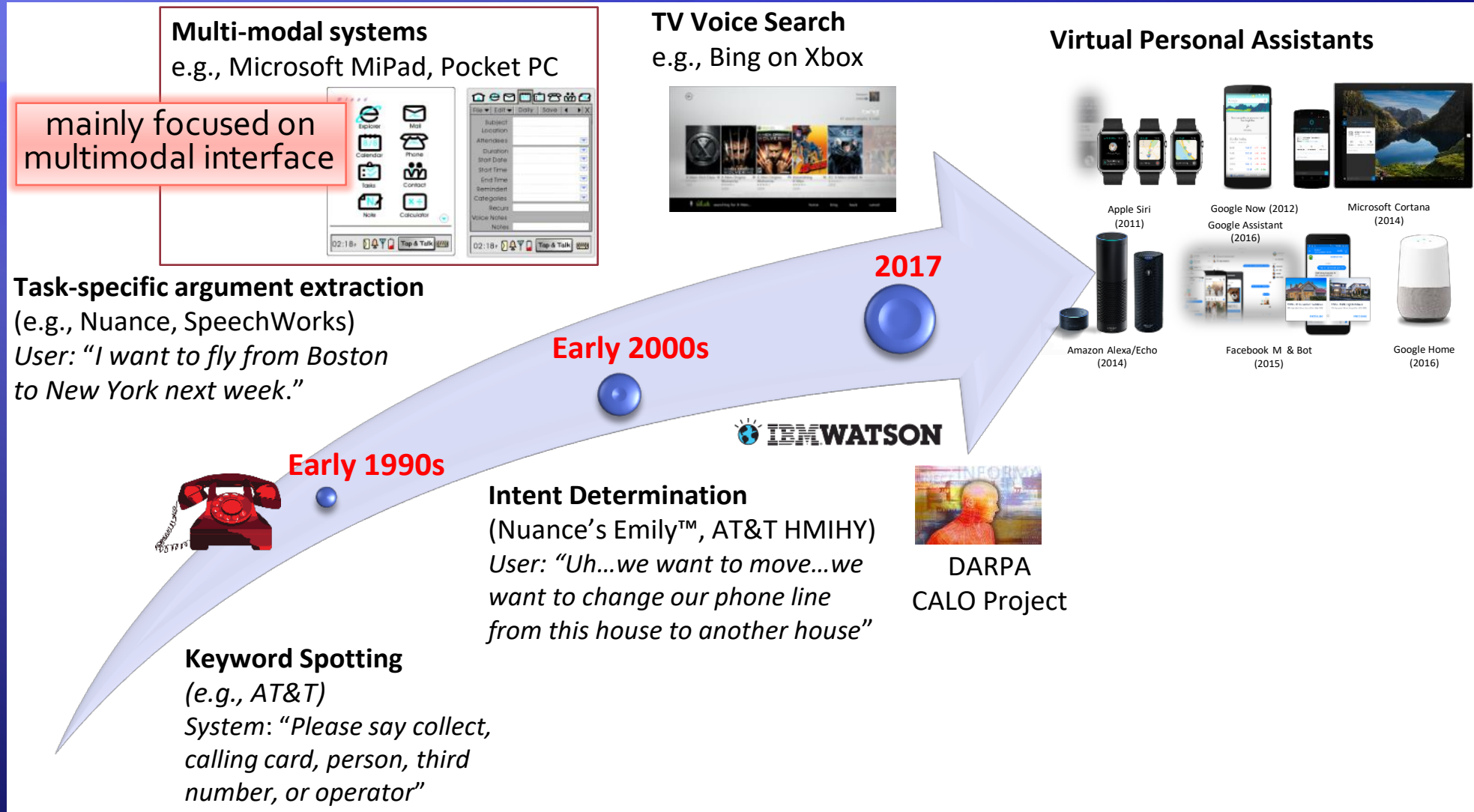


Any similar one in blue?

How to match with it?

Is there any such restaurant nearby?

Is there any shop selling this nearby?

# Evolution of Dialogue Systems

**Multi-modal systems**
e.g., Microsoft MiPad, Pocket PC

mainly focused on multimodal interface

**TV Voice Search**
e.g., Bing on Xbox

**Virtual Personal Assistants**

Apple Siri (2011)

Google Now (2012)
Google Assistant (2016)

Microsoft Cortana (2014)

Amazon Alexa/Echo (2014)

Facebook M & Bot (2015)

Google Home (2016)

**Task-specific argument extraction**
(e.g., Nuance, SpeechWorks)
*User: "I want to fly from Boston to New York next week."*

**2017**

**Early 2000s**

IBM WATSON

**Early 1990s**

**Intent Determination**
(Nuance's Emily™, AT&T HMIHY)
*User: "Uh…we want to move…we want to change our phone line from this house to another house"*

DARPA
CALO Project

**Keyword Spotting**
*(e.g., AT&T)*
*System: "Please say collect, calling card, person, third number, or operator"*

3

# Challenges



**1** **Understanding semantics from text and image**

**2**

**3**

# Challenges



**1** Understanding semantics from text and image

**2** Incorporating domain knowledge

**3**

# Challenges



1 **Understanding semantics from text and image**

2 **Incorporating domain knowledge**

3 **Improving Dialogue flow**

4

# System Overview

◆ Hierarchical RNN       **+ 3 core components**

# 1. Learning Taxonomy-based V

◆ Human perception of product organization and product similarity

- ◆ General to specific
- ◆ Exclusive and Independent relationships (EI)



6

◆ Map images and text into a joint visual semantic space

◆ Leverage EI tree taxonomy to guide fashion concepts learning

# 2. Incorporating Domain Knowledge

◆ Incorporate Knowledge by Multimodal Knowledge Memory Network

# 3. Training with Reinforcement Signals

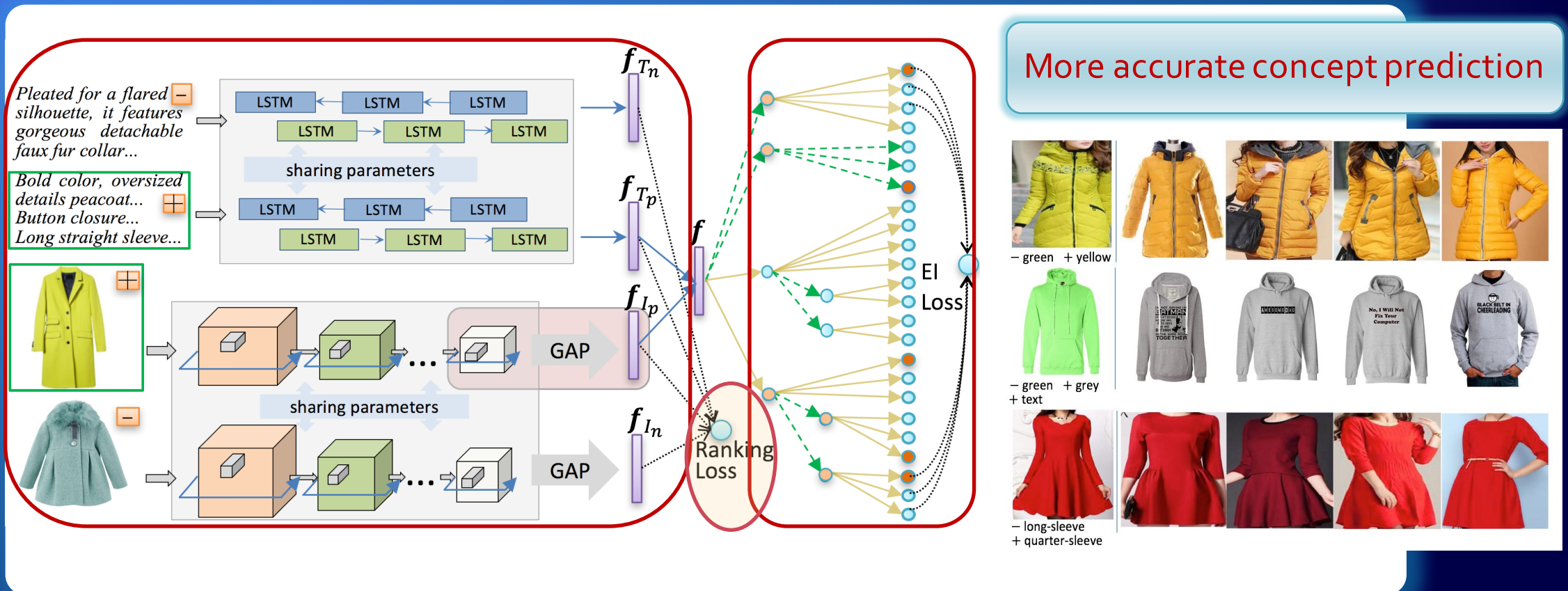◆ Improve dialogue flow via reinforcement signals in two stages training



**1** Predict a generated target utterance given the dialogue context in a **supervised** fashion

**2** Initialized the policy model using the model trained during the first stage, start **fine-tune**

# 3. Training with Reinforcement Signals

◆ Improve dialogue flow via reinforcement signals in two stages training

**rewards**

- **Text response**

$$R(h, r) = BLEU\ score$$

- **Image response**
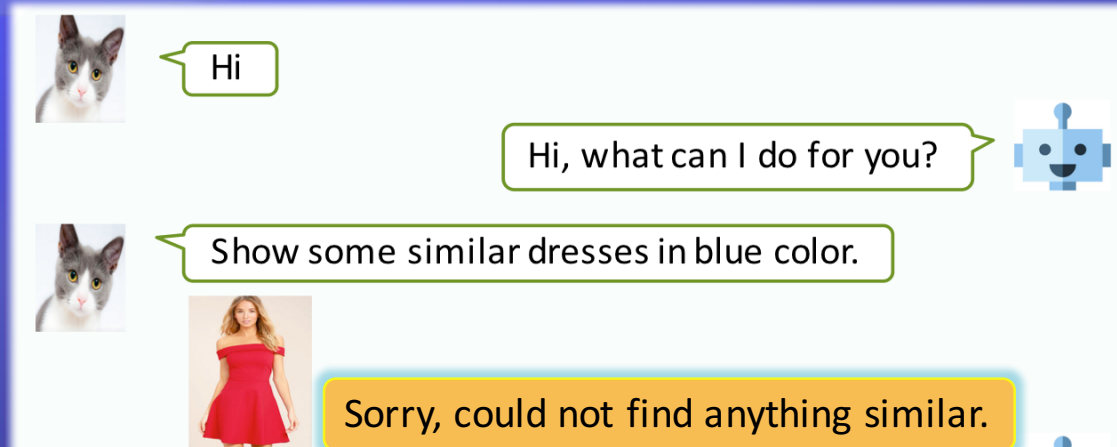
$$R(h, r) = sim(\mathbf{I}, \mathbf{I}^+) - sim(\mathbf{I}, \mathbf{I}^-)$$

**1** Predict a generated target utterance given the dialogue context in a **supervised** fashion

**2** Initialized the policy model using the model trained during the first stage, start **fine-tune**

# Experiments

◆ **Dataset:** 150 K conversation sessions, 1.05 M products, avg. 4 images each

　+ TK ◆ learns more informative representations for fashion products

　+ EK ◆ generates responses not only based on conversation context but also on domain knowledge

　+ RL ◆ fine-tunes the backbone network and optimize the BLEU score or image similarity as rewards

# Experiments

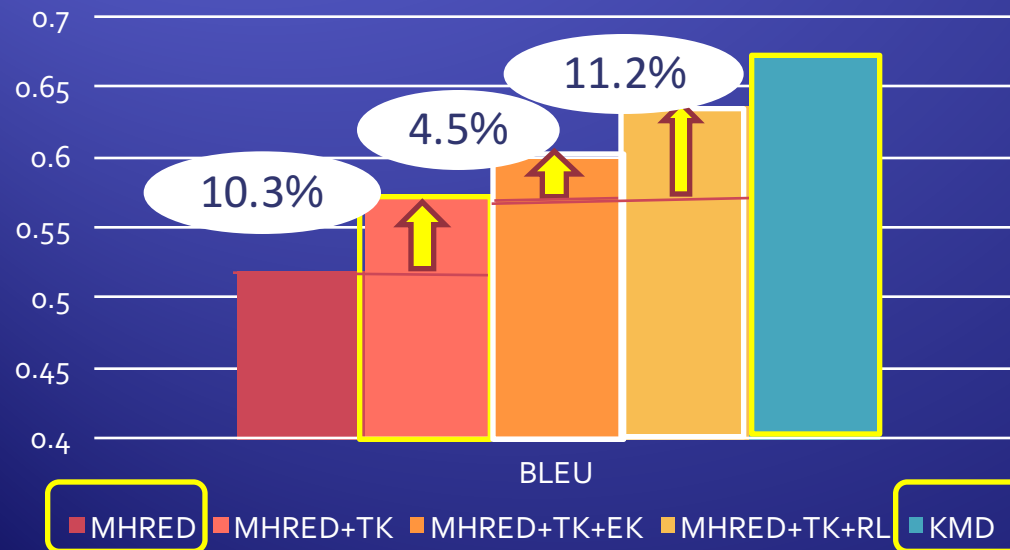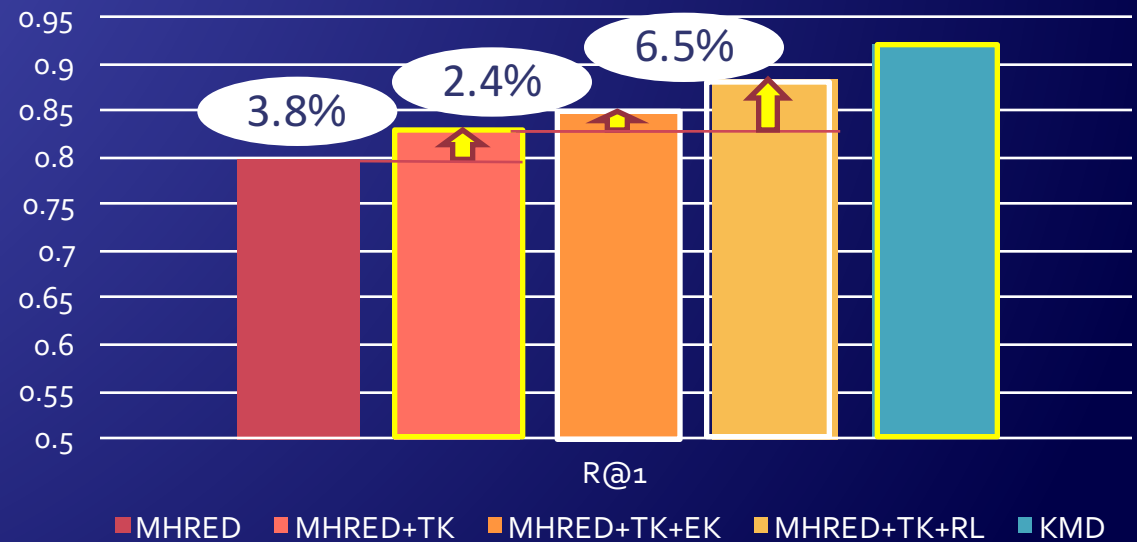- ◆ Sample responses

**Example 1**

**USER:** What is the style in the 1st and 2nd images?

Taxonomy-based semantic learning

**GT:** the style of the formal shoes is oxford in the 1st image; party in the 2nd image

**MHRED:** the style of the scarf is is in the 1st and image image image

**KMD:** the style of the formal shoes is oxford in the 1st image in the image

**Example 2**

**USER:** Which all will go with at least one of these results?

Domain knowledge incorporation

**GT:** it can go well with suede style , suede upper material , suede material running shoes

**MHRED:** it can go well with <unk> , , and and and

**KMD:** it can go well with suede, suede material,, and and shoes

11

# Conclusion and Future Work

- ◆ **Multimodal Dialogue Systems**
  - ✦ Offer an effective way for information seeking
  - ✦ Provide a general scheme for dialogue systems with in-depth visual understanding
  - ✦ Emphasize domain knowledge incorporation for enhancing bot intelligence

- ◆ **Future Work**
  - ✦ Maintain and update the domain knowledge base
  - ✦ Generalize to other domains such as travel, healthcare
  - ✦ Analyze dialogue acts to increase interpretability of dialogue flow control
  - ✦ Start procedural knowledge learning for performing tasks such as nudging customers

# Thank You
## Q & A