

# UT Dialogue System at NTCIR-12 STC

---

Shoetsu Sato<sup>1</sup>, Shonosuke Ishiwatari<sup>1</sup>,  
Naoki Yoshinaga<sup>2</sup>, Masashi Toyoda<sup>2</sup>, Masaru Kitsuregawa<sup>2,3</sup>

The University of Tokyo, <sup>2</sup>IIS, the University of Tokyo, <sup>3</sup>NII, Japan



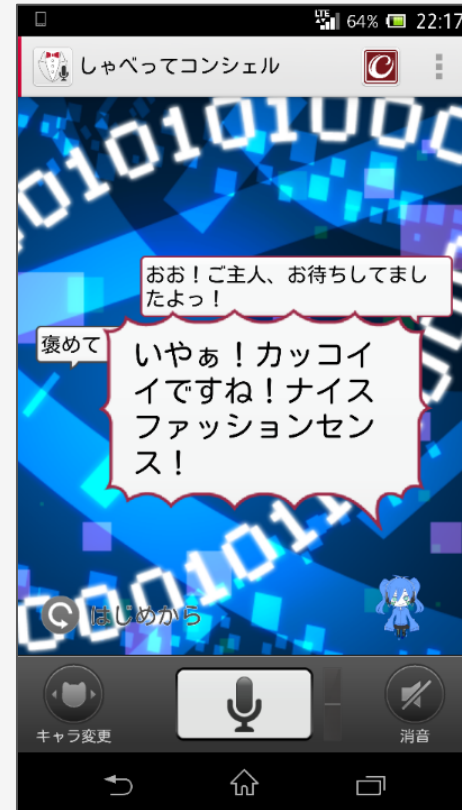
# A lot of dialogue systems that can chat have appeared



Siri (Apple)



Cortana (Microsoft)



しゃべってコンシェル (NTT Docomo)

<http://www.idownloadblog.com/>  
<http://techcrunch.com/2015/01/05/facebook-wit-ai/>  
<http://ameblo.jp/cos-120/entry-11748747974.html>



# Recent approaches for chatting dialogue systems

- Data-driven approaches using dialogue data from **social media** are promising [Ritter+, '10]

U: Utterance  
R: Response

U: また残業か . . .  
R: 生き残ろうな . . .

U: あの人どう思う？  
R: ああいう人間ほんと嫌い



U: 魚介嫌いでした？  
R: そんなこと無いですよ。



# Challenge we have tackled in STC task

- The diversity of **domains** (**topics**, **speaking styles**, etc...) makes it difficult to learn

U: Utterance  
R: Response

U: また**残業**か . . .  
R: 生き残ろうな . . .

U: あの**人**どう思う？  
R: ああいう人間ほんと嫌い



U: **魚介**嫌い**でした**？  
R: そんなこと無いですよ。



# OVERVIEW



# Goal: building a domain-aware dialogue model

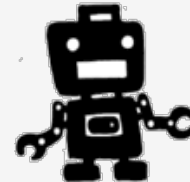
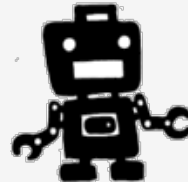
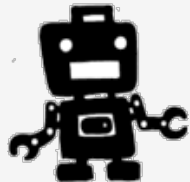
**Idea:** Divide conversation data into **domain-consistent subsets** to train multiple specific LSTM-based dialogue models

**Evaluation:** response selection from candidates

U: また**残業**か . . .  
R: 生き残ろうな . . .

U: あの**人**どう思う ?  
R: ああいう人間ほんと嫌い

U: **魚介**嫌い**でした** ?  
R: そんなこと無いですよ。



**Does domain consistence compensate for reduction of training data per a model?**

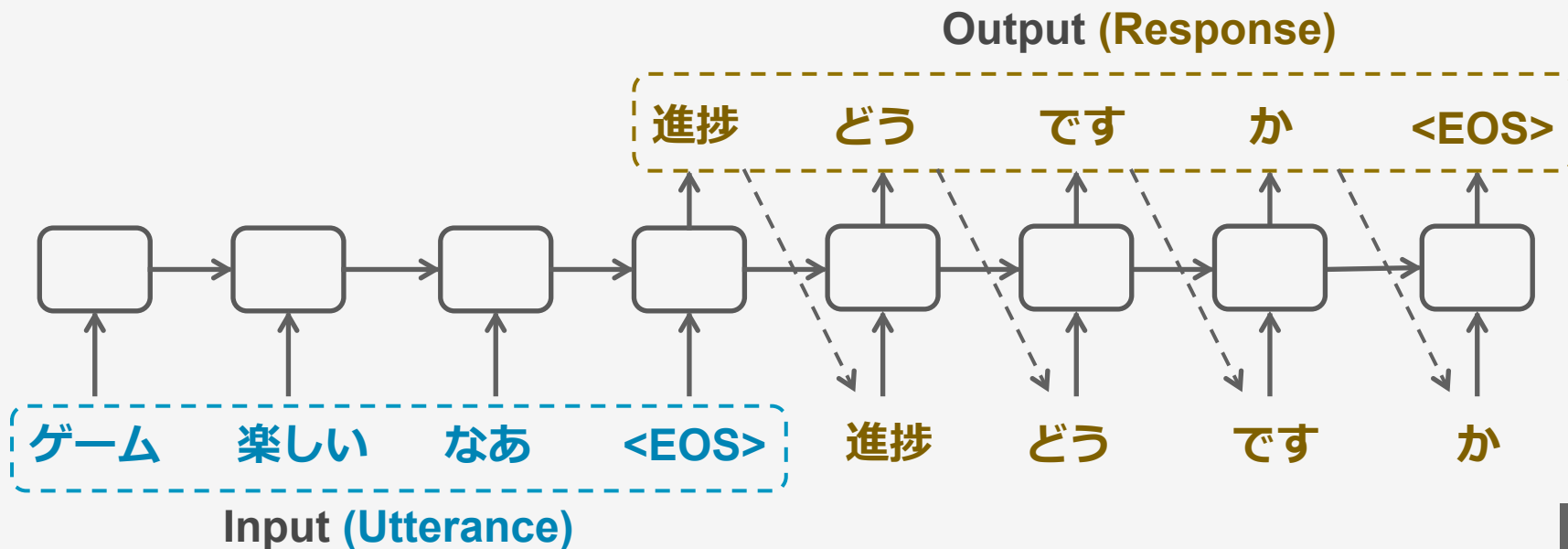


# RELATED WORK



# Recent promising approach to generate responses

- We employed recent promising Long-Short Term Memory based Recurrent Neural Network (LSTM-RNN) dialogue model [Vinyals+, '14]

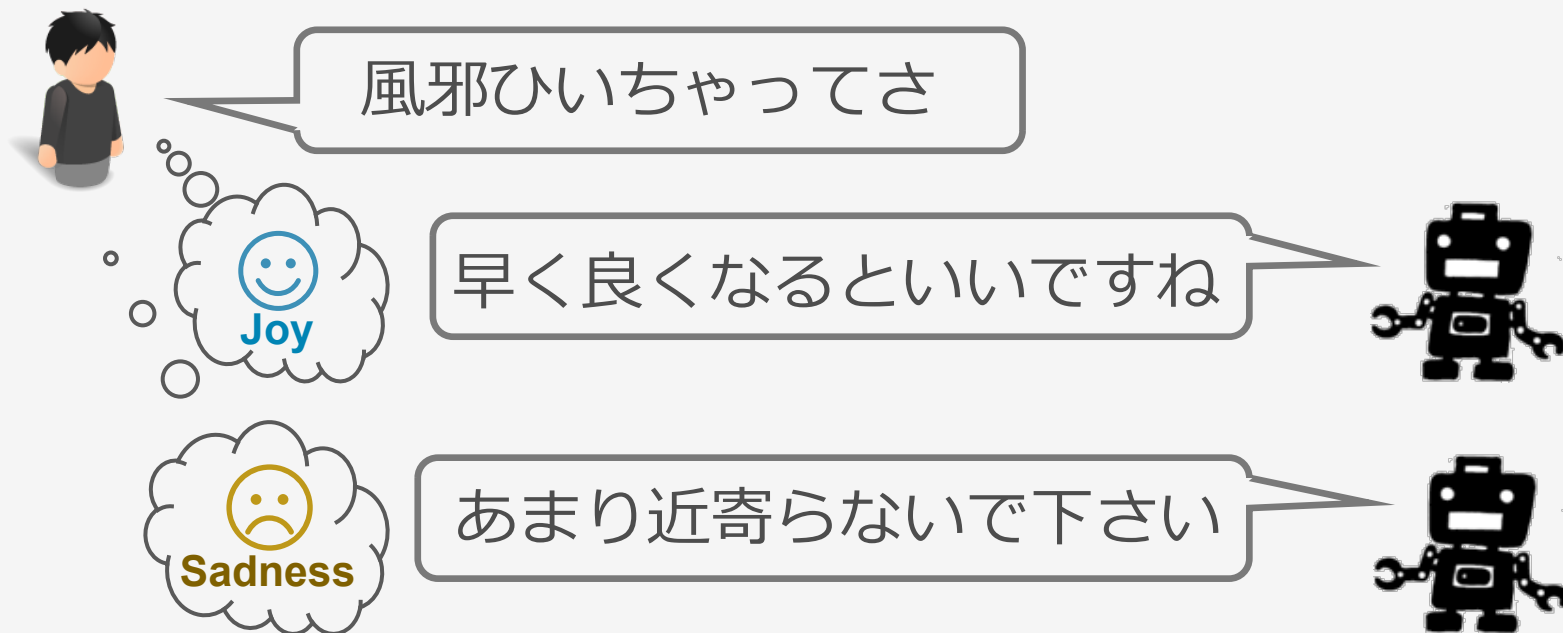






# Predicting and Eliciting Addressee's Emotion in Online Dialogue [Hasegawa+, '13]

- Generate a response that elicits a **specific emotion** in the addressee's mind





# Overview of the related work and target of our method

## Point

- General LSTM based methods employed a single model trained from all data
- It is impossible to enumerate all domains in human dialogues

## Purpose

- Capture the difference of domains **automatically** as clusters and train **multiple** models





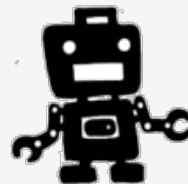
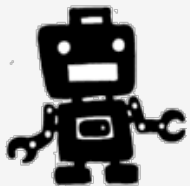
# Our approach: K-cluster model (1/2)

- Cluster the dialogues for each of the **unlabeled domain**, and train multiple models

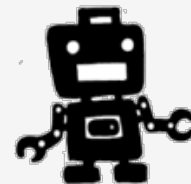
U: また**残業**か . . .  
R: 生き残ろうな . . .

U: あの**人**どう思う ?  
R: ああいう人間ほんと嫌い

U: **魚介**嫌い**でした** ?  
R: そんなこと無いですよ。



train LSTM





# Our approach: K-cluster model (2/2)



Utterance

就職したくない . . .



Find the nearest domain by human utterance and utterances in training subsets

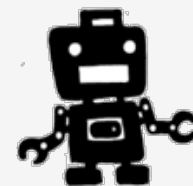
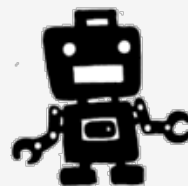
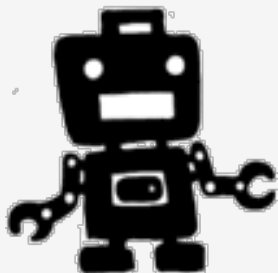
U: また残業か . . .

U: あの人どう思う ?

U: 魚介嫌いでした ?

Response

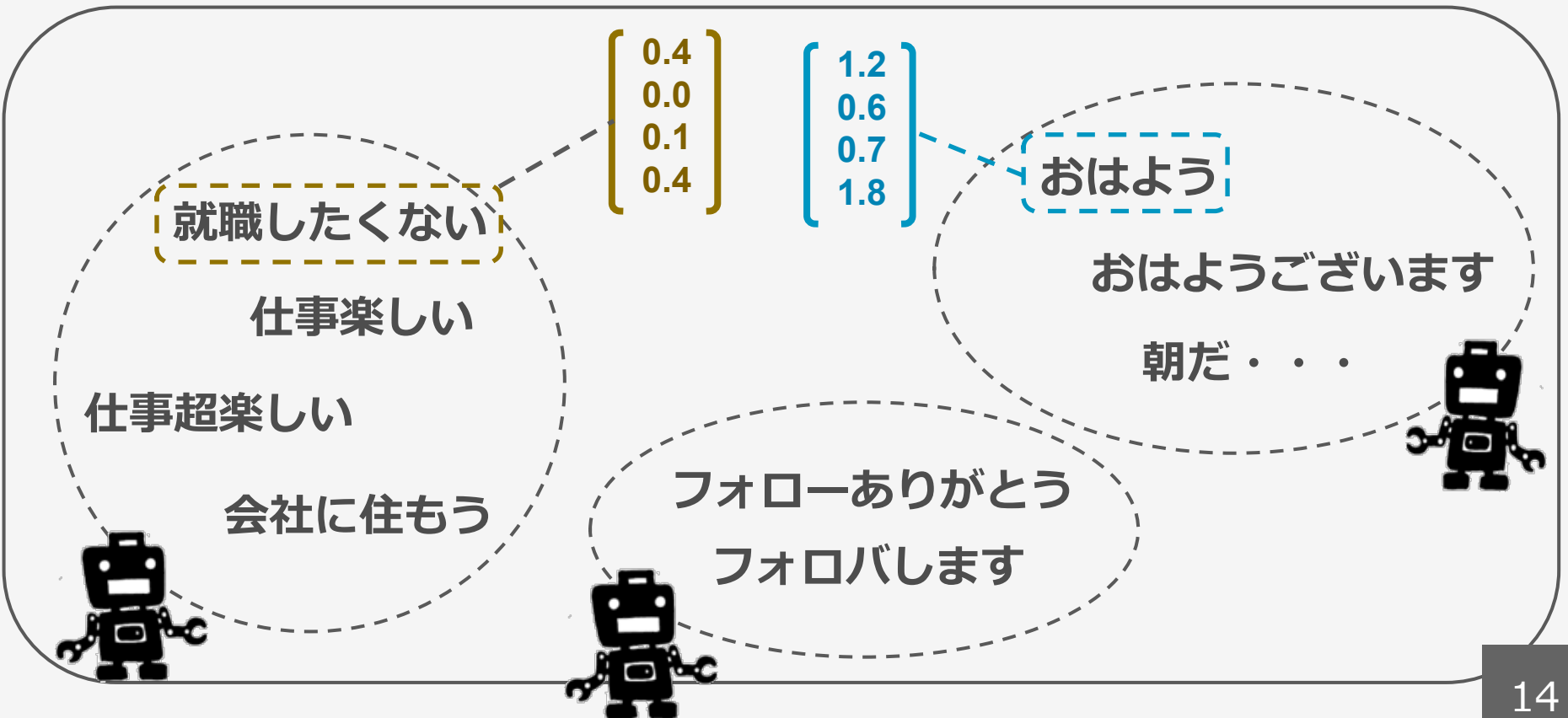
市民、労働は義務です。





# How to automatically handle the domains in each utterance (1/2)

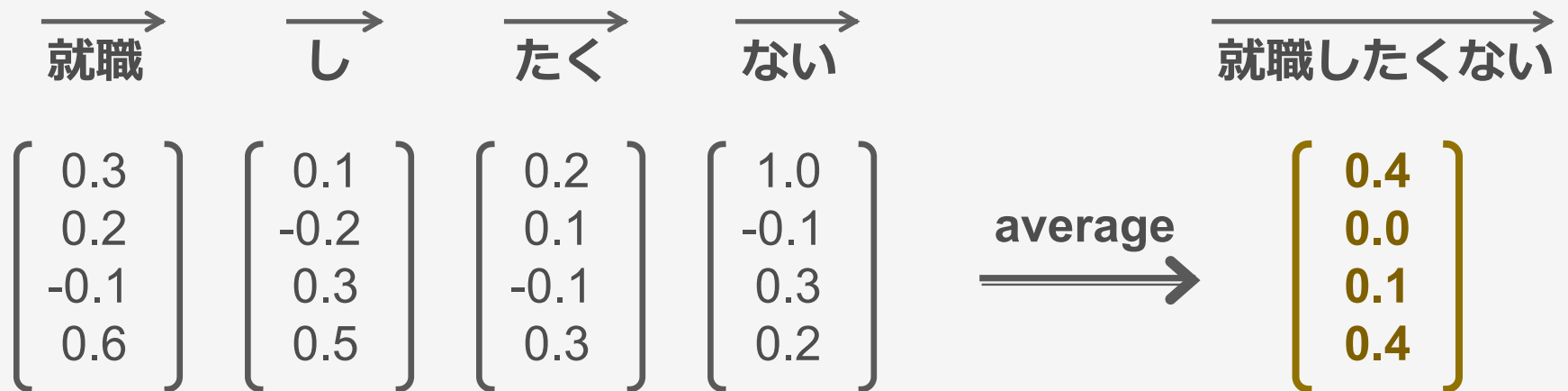
- Apply **k-means clustering** to the utterance vectors and regard clusters as subsets of the training data





# How to automatically handle the domains in each utterance (2/2)

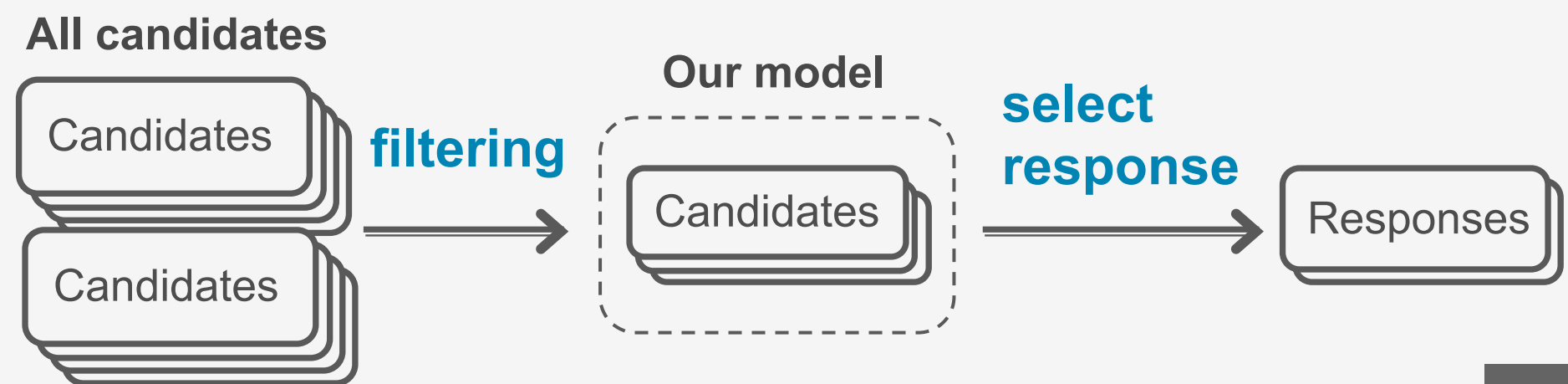
- Represent each utterance as a vector built from **word embeddings** [Mikolov+, '13]
- The density of word embeddings would solve sparseness problems in short texts compared with Bag-of-Words





# Response candidate filtering

- In response selection task from many candidates, our model's **high computational cost** causes a problem
- To reduce the number of candidates into **500**, we employed a fast SVM classifier [Yoshinaga+, '10]







# EXPERIMENTS



## 3 experiments we did

- **Experiment 1: Small response selection task**
  - Evaluate **how our method effects** in response selection
  - Select response from 20 candidates **without filtering**
- **Experiment 2: Filtering performance**
  - Evaluate **to what extent** our filter can select proper candidates
- **Experiment 3: NTCIR-12 formal run**
  - Evaluate **whole performance** of our system (clustered-LSTM, and filter)



# Experiment 1 : small response selection task

- **Dataset: Twitter**

Utterance-response (tweet-reply) pairs crawled from Twitter: **100K** for training, **1K** for test

- Provided for *NTCIR-12 Short Text Conversation Japanese Task* [Shang+, '16]

- **Evaluation: Response selection**

The proportion of test tweets where we succeeded to select the **correct (actually replied) response** from randomly chosen **20 candidates**



# Evaluation detail

## Utterance

発表つらいんだけど

Correct Response

## Ranked 20 Response Candidates

①

わかる

②

自分の研究を知ってもら  
う良い機会だと思うよ

③

今完全に鬱だよ

④

その店美味しいよね

⋮

②⑩

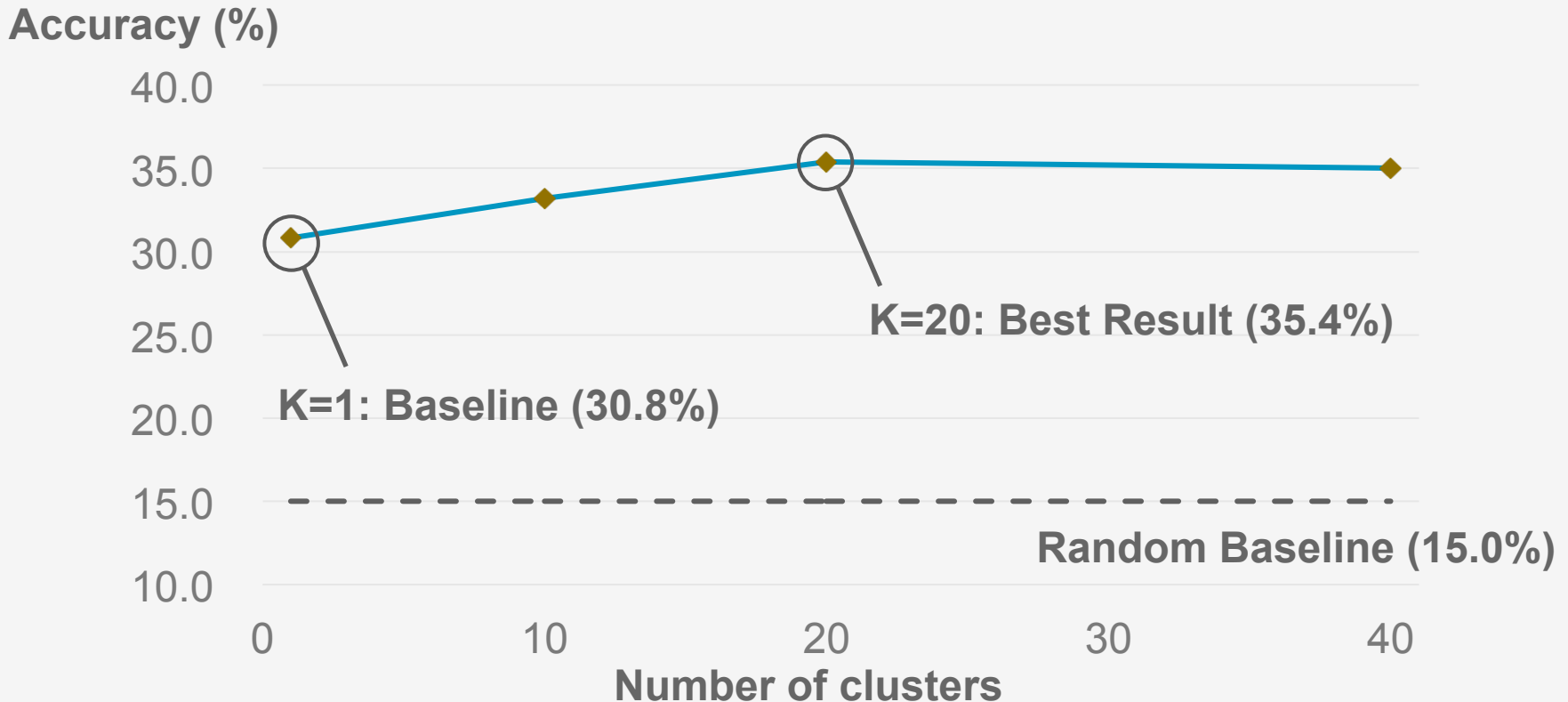
くぁwせdrftg

We defined it as success if the **top-3 responses** include the **correct response**



# Results of K-cluster model

- We compared 1, 10, 20, and 40 cluster models increasing number of clusters until the accuracy was saturated



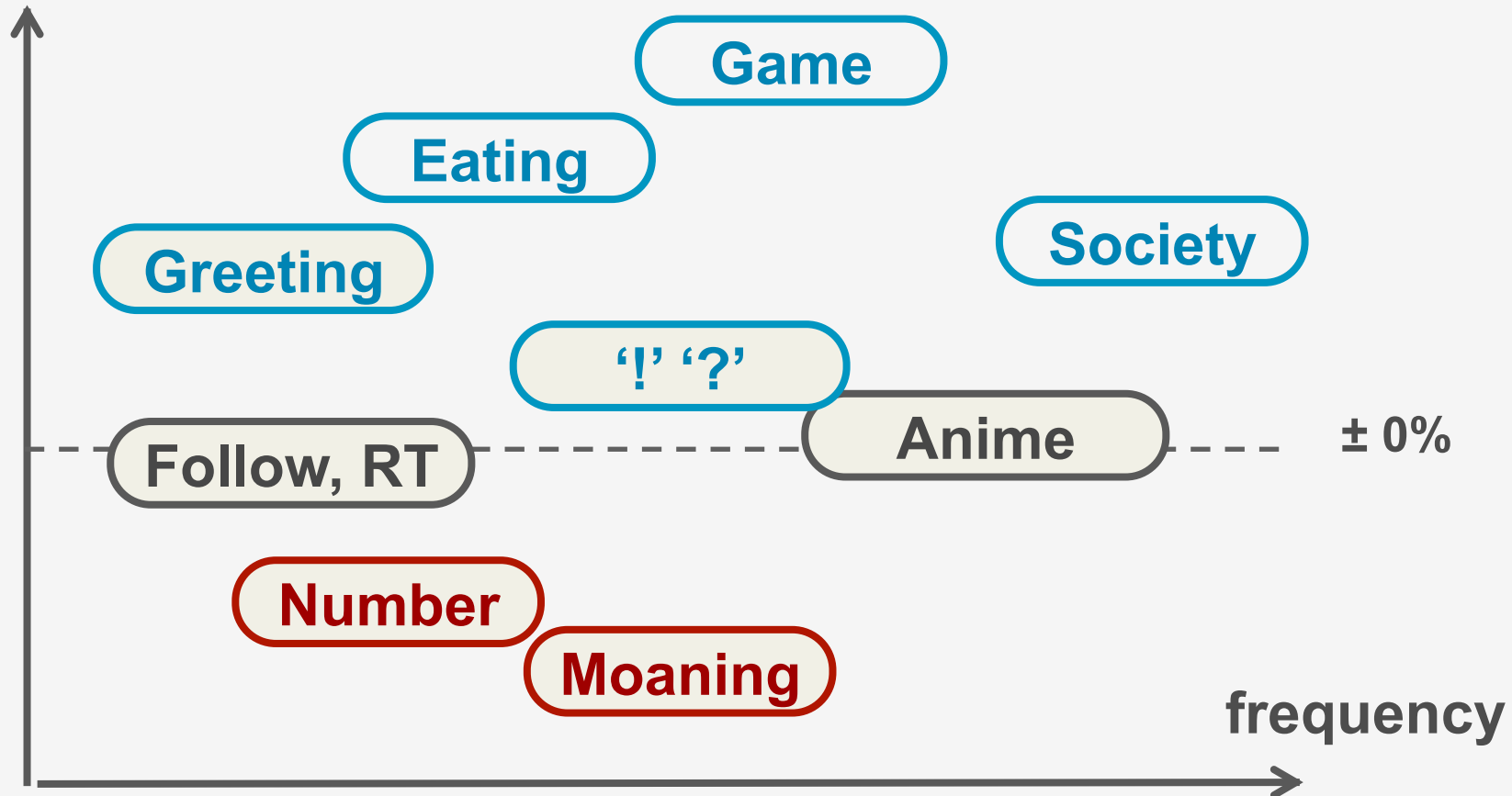
Our best 20-cluster model outperformed baseline at **4.6%**



# Examples of clusters we labeled by hand from their frequent utterances

- 13 models out of 20 cluster got improved

Improvement rate





# Examples of the selected responses

(1)

<b>Utterance</b>	あ、見るの忘れてた。おめでとう!
<b>Proposed</b>	ありー! 見直してくれてありがとう!
<b>Baseline</b>	今年は1年ありがとうございました

(2)

<b>Utterance</b>	カントリーマアムのドリンクのやつが 見つかりません。
<b>Proposed</b>	ローソン限定じゃなかったっけ?
<b>Baseline</b>	先輩、おはよーございます♪

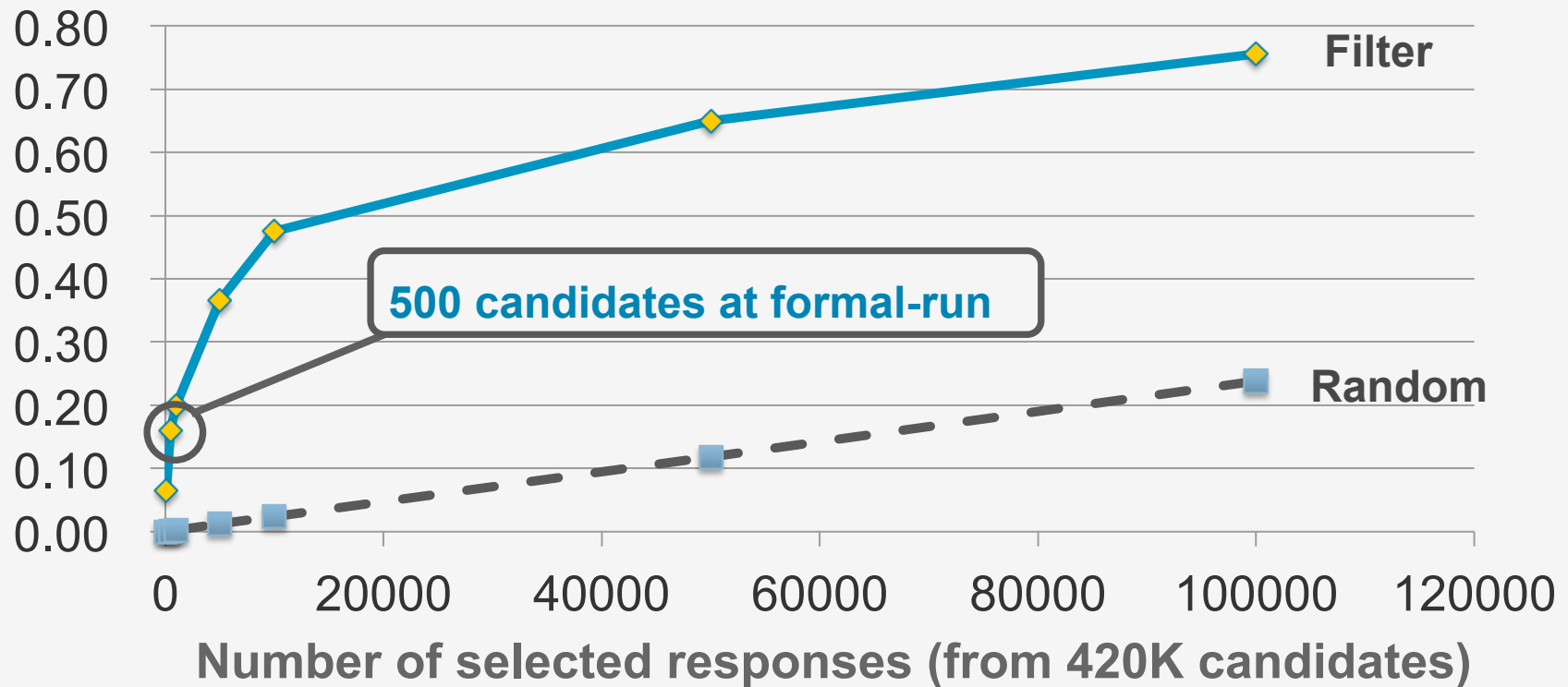
Our proposed model tends to stop selecting typical responses



## Experiment 2: Filtering performance

- Evaluate the filter by recall, whether top-N filtered candidates include the **correct response**

Recall



**Filtering effectively reduced the number of candidates**





## Experiment 3 : NTCIR-12 STC Japanese Task

- **Model: 20 cluster model**

The best one evaluated at **experiment 1**, 20 cluster model trained from **100k** utterance-response pairs

- **Evaluation:**

- For the **204** provided test utterances, select responses from **500k** candidates
- responses are assigned scores of **0 (inappropriate)**, **1 (appropriate in some context)**, **2 (appropriate)** by human annotators

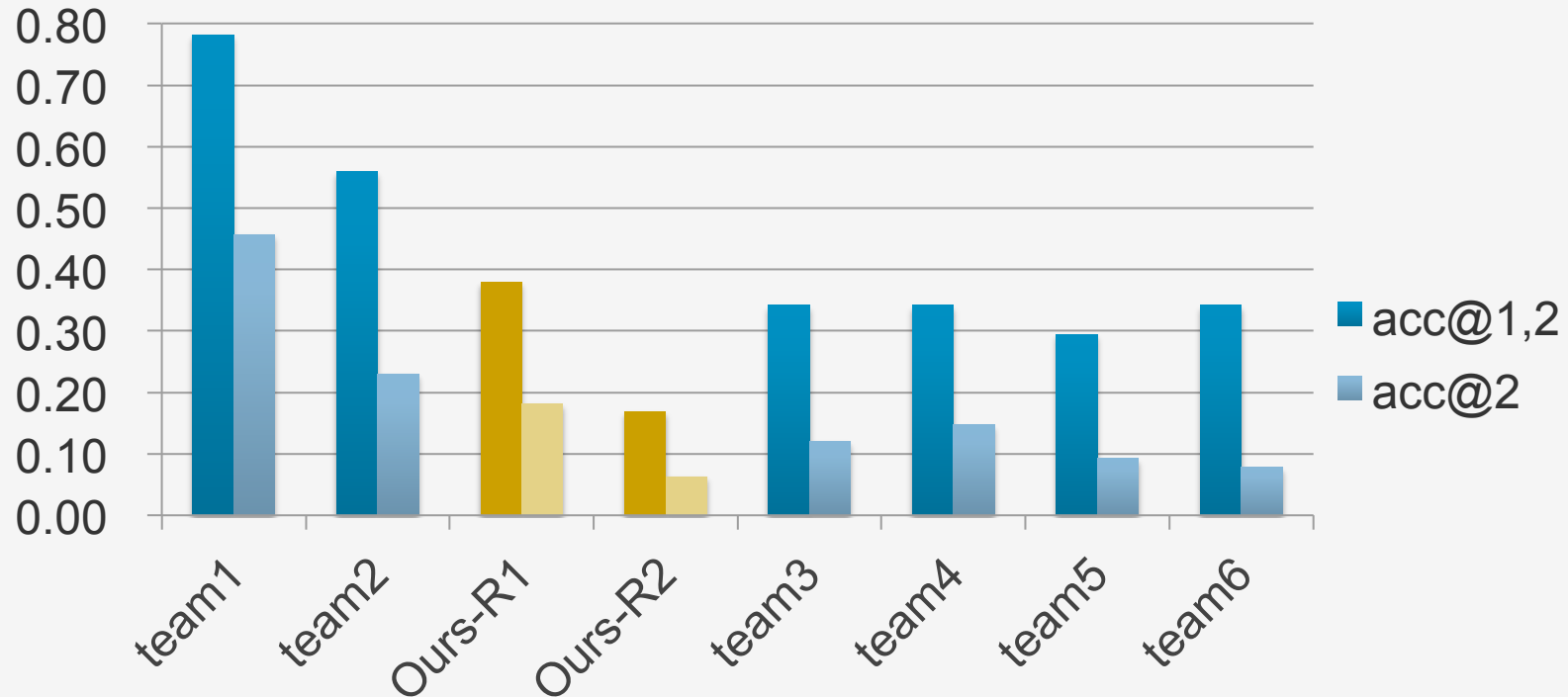


# Accuracies of selected top-1 responses at NTCIR-12 STC Japanese Task

Ours-R1: Filtering + 20-cluster LSTM model

Ours-R2: Filtering only

Accuracy @3



**Our system selected better responses from  
filtered candidates**



# Summary

By response selection test we confirmed the effect of cluster-based **domain-aware** dialogue model

- Domain-consistent training subsets made better results in spite of reduction of training data
- By filtering candidates, our system could effectively select responses



# RESULTS FOR EACH CLUSTER



# Results in each cluster (20-cluster model)

domain (topics, wording, writing style)	#elems		#corr		improvement $\frac{\Delta \# \text{corr}}{\# \text{elems (test)}}$
	train	test	ours	baseline	
-	11801	108	<b>38</b>	27	10.19%
-	11524	124	<b>37</b>	32	4.03%
politics, economics, social matters	10294	130	<b>48</b>	38	7.69%
-	9743	94	<b>32</b>	23	9.57%
animation, comics	6747	56	<b>11</b>	10	1.79%
-	6552	66	<b>24</b>	23	1.52%
game	5677	50	<b>13</b>	5	16.00%
-	5627	45	<b>14</b>	13	2.22%
end with '?' r '!	5190	63	<b>17</b>	15	3.17%
moaning (esp., sleepy, weary)	5064	52	17	<b>21</b>	-7.69%
-	4908	50	22	<b>24</b>	-4.00%
numbers	3803	31	5	<b>7</b>	-6.45%
eating	2630	16	<b>6</b>	4	12.50%
frank acknowledgment (follow, RT)	2252	33	29	<b>30</b>	-3.03%
end with '!!!'	1869	17	<b>8</b>	<b>8</b>	0.00%
polite acknowledgement (follow, RT)	1553	13	<b>12</b>	<b>12</b>	0.00%
greetings	1537	21	<b>7</b>	6	4.76%
end with '...'	1326	12	<b>3</b>	2	8.33%
polite morning greetings	1174	13	<b>9</b>	6	23.08%
shouting with word lengthing or repetition	729	6	<b>2</b>	<b>2</b>	0.00%
	100000	1000	354	308	4.60%