

UT Dialogue System at NTCIR-12 STC

Shoetsu Sato¹, Shonosuke Ishiwatari¹, Naoki Yoshinaga²,
Masashi Toyoda², and Masaru Kitsuregawa^{2,3}
{shoetsu, ishiwatari, ynaga, toyoda, kitsure}@tkl.iis.u-tokyo.ac.jp

¹ The University of Tokyo, ² IIS, the University of Tokyo, ³ NII, Japan

Background

In data-driven approaches for chat-dialogue modeling, the diversity of **domains** (topics, speaking styles, emotions..) makes it difficult to learn

U: フォローしました!
R: ありがとうございます!

U: Utterance
R: Response

U: また残業か・・・
R: 生き残ろうな・・・



U: ラーメン食べたい気分
R: 今日の夜行こうぜ

Related work

Classify training data by **several emotion types** each response elicits and train multiple models [Hasegawa+, '13]



But it is impossible to enumerate all domains in human dialogues by hand

Our approach

Cluster conversation data to **automatically** capture the difference of domains and train specific models



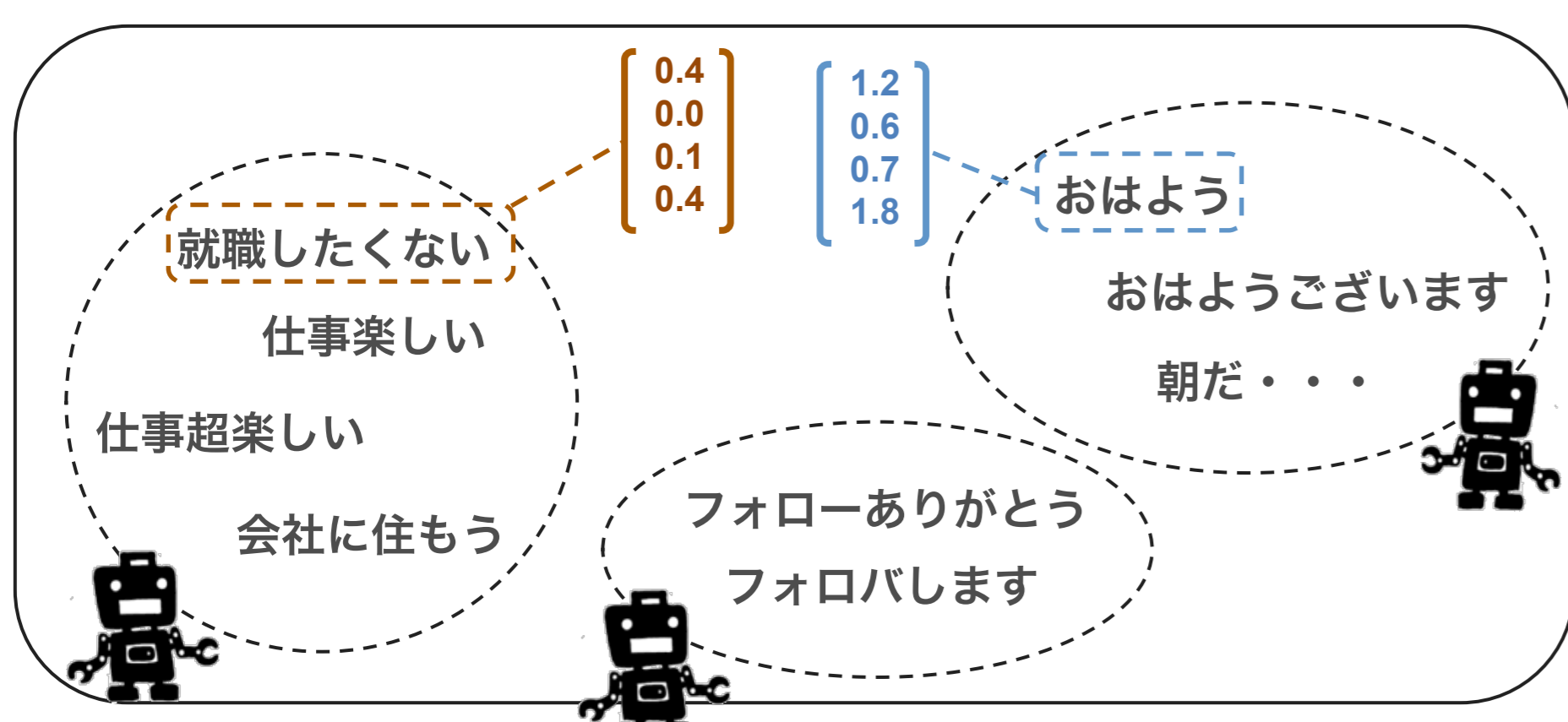
Domain-consistent responses



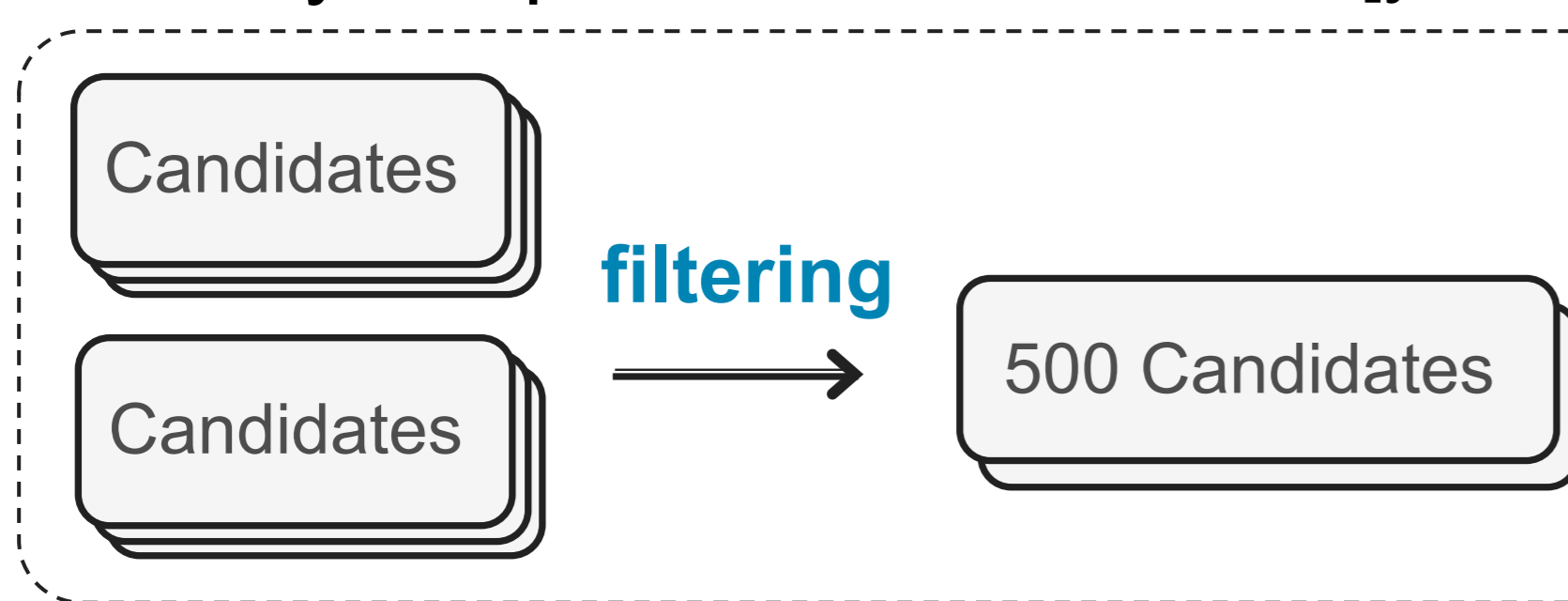
Smaller size of the training data per a model

Proposed Method

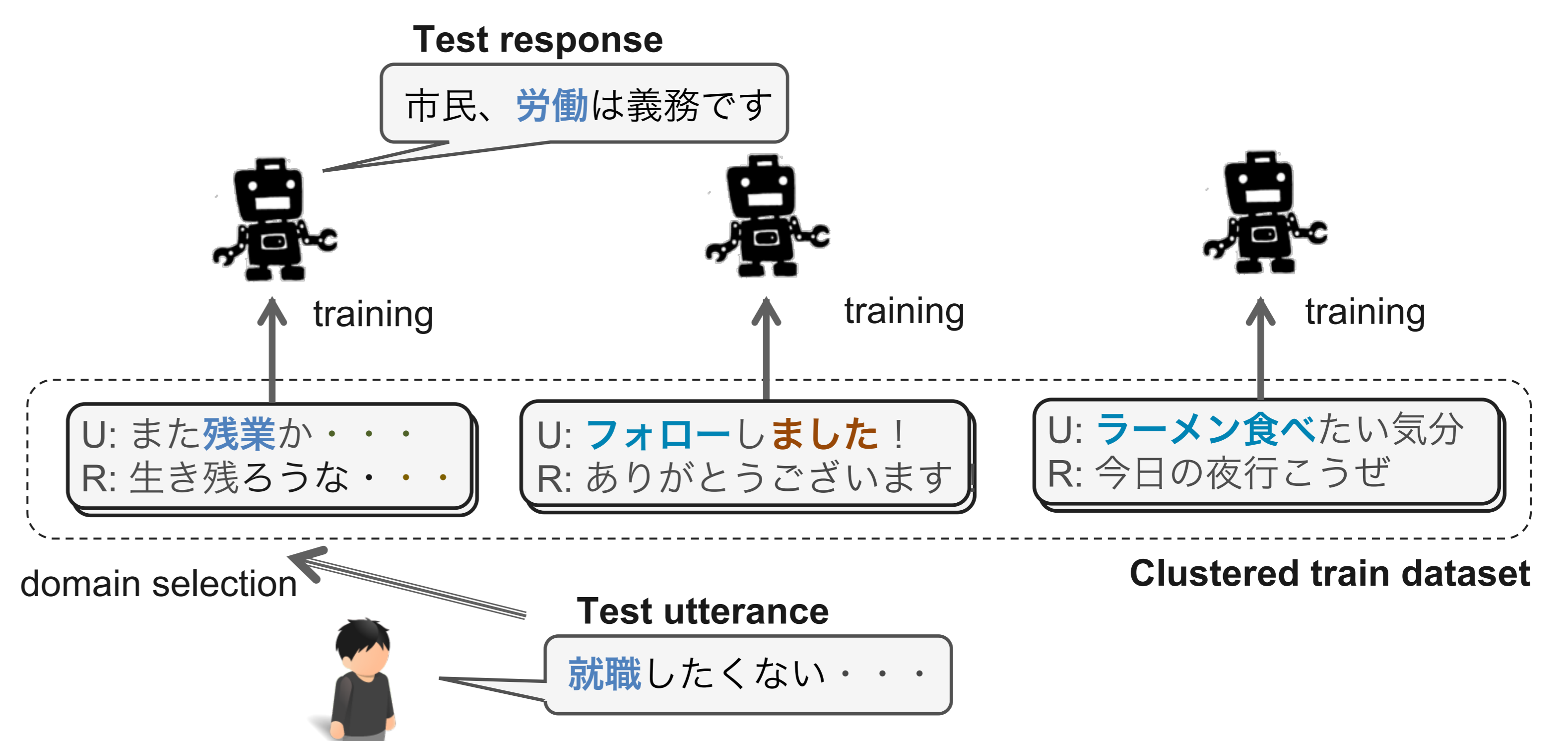
- Apply **k-means clustering** to the utterance vectors and regard clusters as subsets of the training data



- Narrow the number of the candidates to reduce computation by the pre-trained classifier [yoshinaga+, '10]



- Train multiple LSTM-based dialogue models by **each domain-specific training data subset**



- Select the model to respond from distance between **cluster's centroids and the utterance vector** and response from candidates

Experiments

Effectiveness of clustering

Data

100K (tweet-reply) pairs for train, 1K for test

Evaluation method

Utterance

発表づらいんだけど

Ranked 20 Candidates

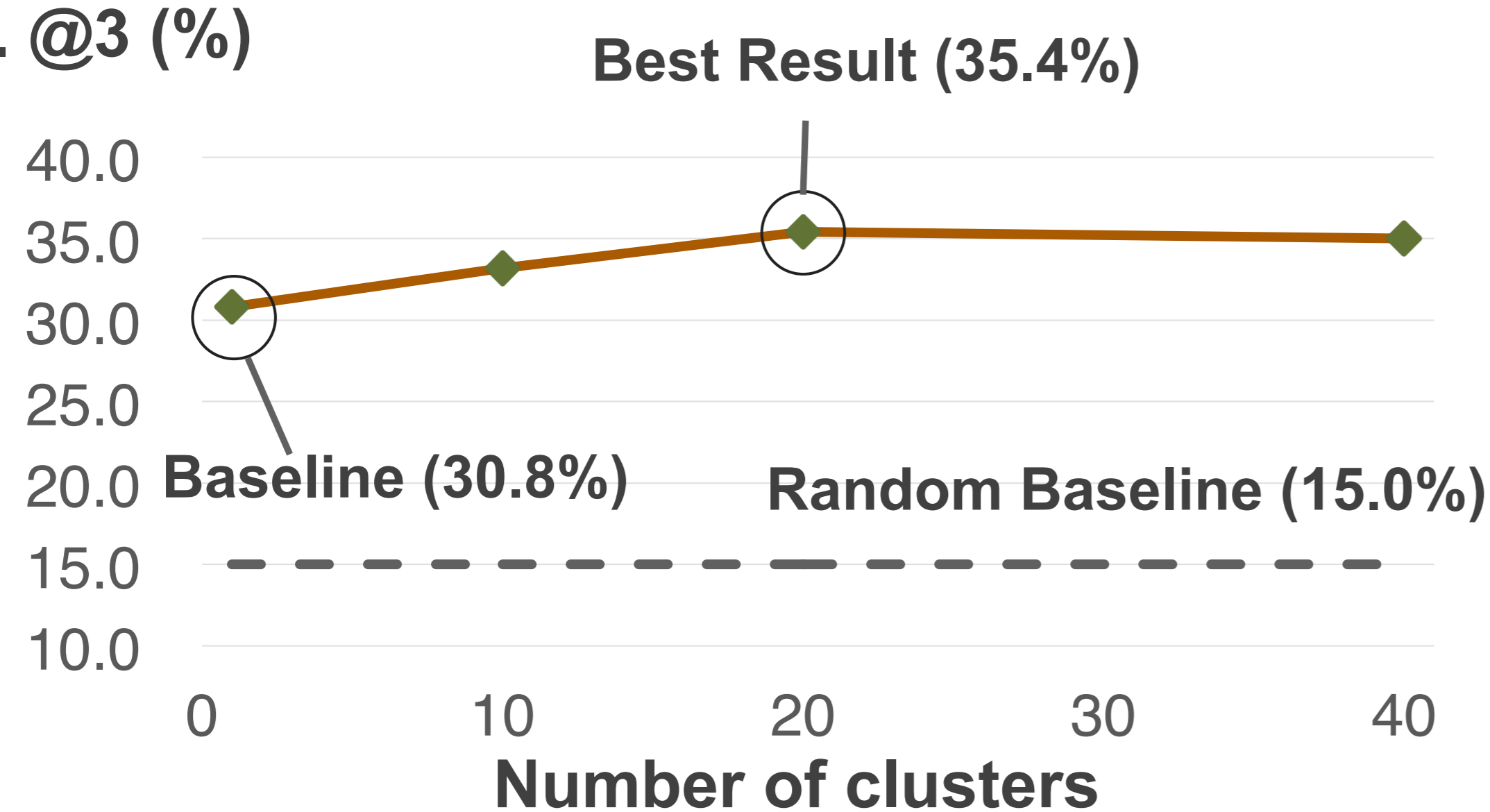
- わかる
- 自分の研究を知ってもらいたい機会だよ
- 今完全に鬱だよ
- その店美味いよね
- ...
- くあwせdrftg

Correct Response

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

Results

Acc. @3 (%)



Difference between baseline and proposed method

Utterance

あ、見るの忘れてた。おめでとう!

Baseline

今年は1年ありがとうございました

Proposed

ありー! 見なおしてくれてありがとう!

Our method less frequently select typical responses by extracting them as other domains

NTCIR-12 STC formal-run

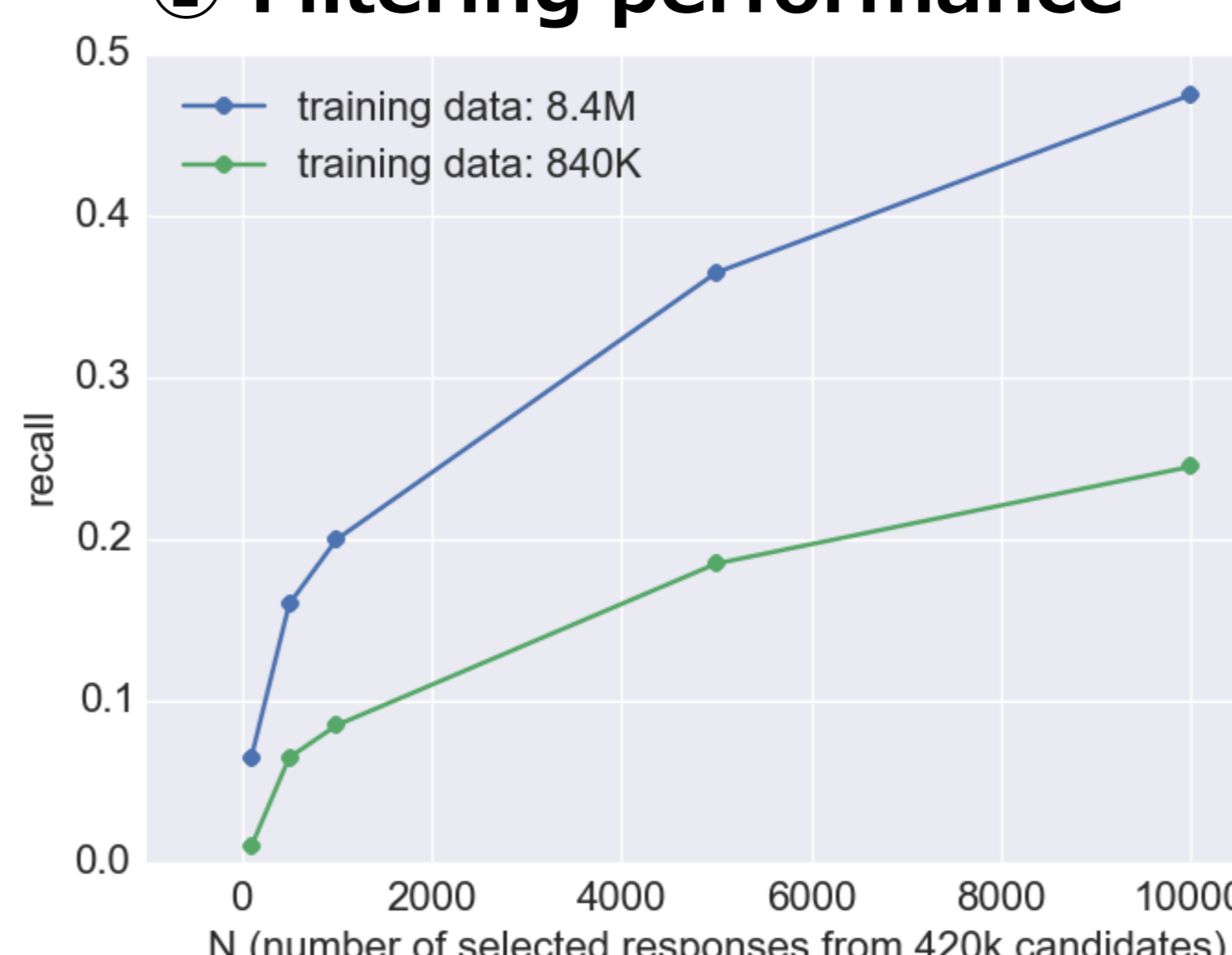
Evaluation

- Evaluate filters trained on different size of training data, by recall whether **top-N** filtered candidates including the correct response
- Selected responses are assigned score of **0 (inappropriate)**, **1 (appropriate in some context)**, and **2 (appropriate)** by human, and evaluated the proportion of **1 and 2**, or **only 2** for the **top-1** or **top-5** selected responses.

R1 : Responses selected by our system from filtered candidates

R2 : Responses only pre-filtered

① Filtering performance



② Accuracy on the NTCIR-12 STC task

