# An Evaluation Dataset for Intent Classification and Out-of-Scope Prediction

Petr Marek

# Outline

Problem statement

New Dataset

Metrics

Baseline Models

My Initial Work

# Can a robot spice up the retail banking experience? HSBC's 'Pepper' is now on the job at Seattle branch

BY **KURT SCHLOSSER** on March 12, 2019 at 12:11 pm

Post a Comment      ✉ Email

Pepper the humanoid robot is seen in action at HSBC Bank in Beverly Hills, Calif. (Amanda C. Edwards Photo / @acefotopro)

Bank customers who have already embraced technology as a way of handling their transaction needs could be lured offline and back into a retail branch if one cute robot has anything to do with it. Seattle, don't pass the Pepper, come in and say hi.

What is my balance?

Can you give a mortgage?

I would like to open an account.

# Intent Classification

| Balance | Mortgage | Account |
|---|---|---|
| What is my balance? | Can I have a mortgage? | I want a new account! |
| Can you tell me how much money do I have left on my account? | I need a mortgage! | Can you open me an account? |
| How much money do I have? | Can you assist me with a mortgage? | I would like to become your customer. |
| Are there any money left on my account? | How can I get a mortgage? | I would like to open an account. |
| What is the balance of my account? | I would like a mortgage! | New account, I need it now! |

- Mortgage
- Account
- Balance

What is my balance? **Balance**

Can you give a mortgage? **Mortgage**

I would like to open an account. **Account**

Do you like movies?

How are my sport teams doing?

How to get from Paris to London?

# Other

- 🟢 Mortgage
- 🔵 Account
- 🔵 Balance

# Other

Movies

Mortgage
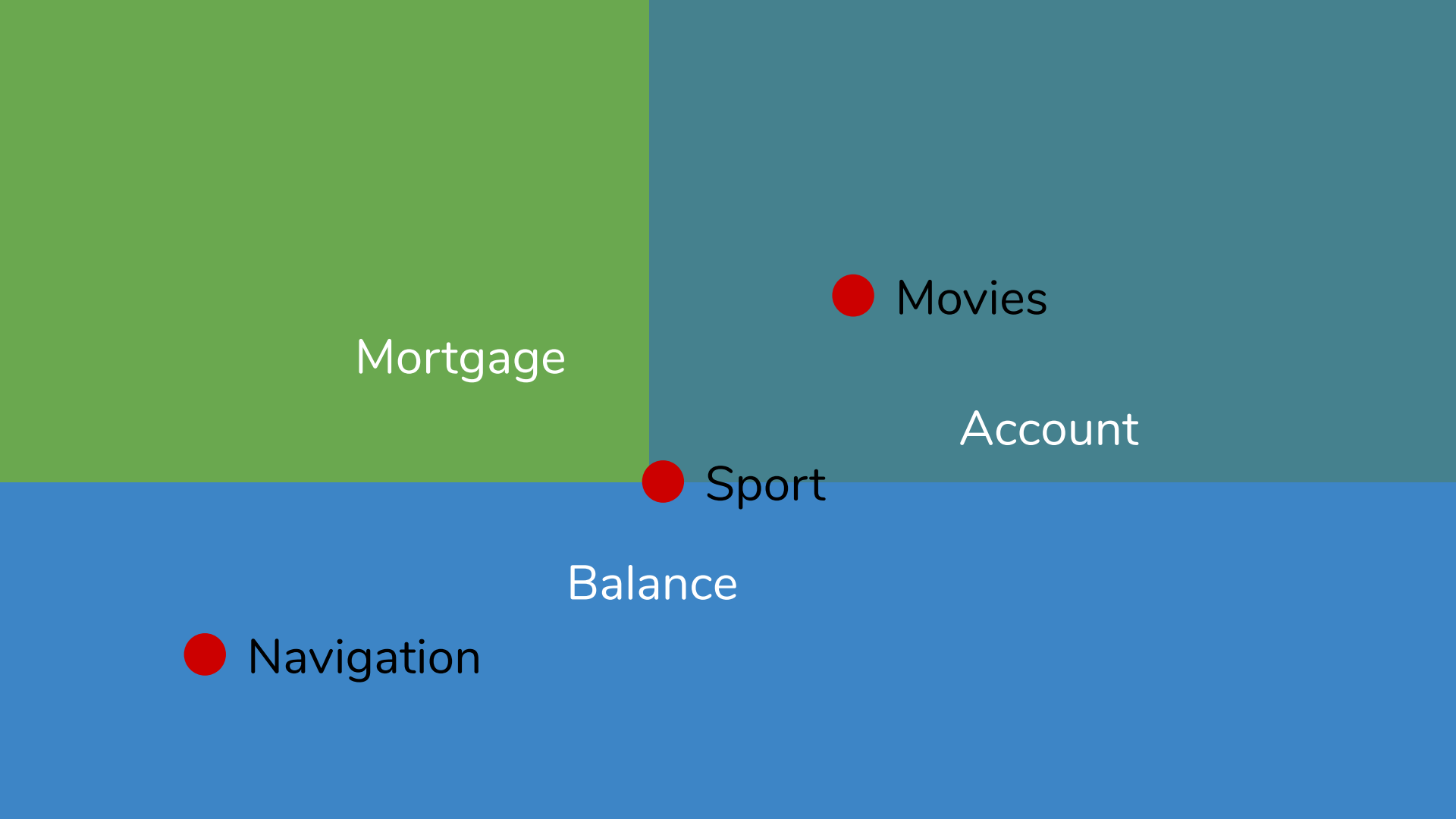
Account

Sport

Balance

Navigation

$$y^* = \text{argmin}_{\hat{y}} \sum_{y \in Y} L(y, \hat{y}) \, \text{Pr}[y|x]$$

Zero-one Loss

$$y^* = \text{argmax}_{y} \, \text{Pr}[y|x]$$

# Other

# What is the most similar intent?

# What is the most similar intent?

# AND

# Isn't it out of scope?

# An Evaluation Dataset for Intent Classification and Out-of-Scope Prediction

Stefan Larson    Anish Mahendran    Joseph J. Peper    Christopher Clarke
Andrew Lee    Parker Hill    Jonathan K. Kummerfeld    Kevin Leach
Michael A. Laurenzano    Lingjia Tang    Jason Mars
Clinc, Inc.
Ann Arbor, MI, USA
stefan@clinc.com

## Abstract

Task-oriented dialog systems need to know when a query falls outside their range of supported intents, but current text classification corpora only define label sets that cover every example. We introduce a new dataset that includes queries that are out-of-scope—i.e., queries that do not fall into any of the system's supported intents. This poses a new challenge because models cannot assume that every query at inference time belongs to a system-supported intent class. Our dataset also covers 150 intent classes over 10 domains, capturing the breadth that a production task-oriented agent must handle. We evaluate a range of benchmark classifiers on our dataset along with several different out-of-scope identification schemes. We find that while the classifiers perform well on in-scope intent classification, they struggle to identify out-of-scope queries. Our dataset and evaluation fill an important gap in the field, offering a way of more rigorously and realistically benchmarking text classification in task-driven dialog systems.

## 1 Introduction

Task-oriented dialog systems have become ubiquitous, providing a means for billions of people to interact with computers using natural language. Moreover, the recent influx of platforms and tools such as Google's DialogFlow or Amazon's Lex for building and deploying such systems makes them even more accessible to various industries and demographics across the globe.

Tools for developing such systems start by guiding developers to collect training data for intent classification: the task of identifying which of a fixed set of actions the user wishes to take based on their query. Relatively few public datasets exist for evaluating performance on this task, and those that do exist typically cover only a very small number of intents (e.g. Coucke et al. (2018), which has 7



Figure 1: Example exchanges between a user (blue, right side) and a task-driven dialog system for personal finance (grey, left side). The system correctly identifies the user's query in ①, but in ② the user's query is mis-identified as in-scope, and the system gives an unrelated response. In ③ the user's query is correctly identified as out-of-scope and the system gives a fallback response.
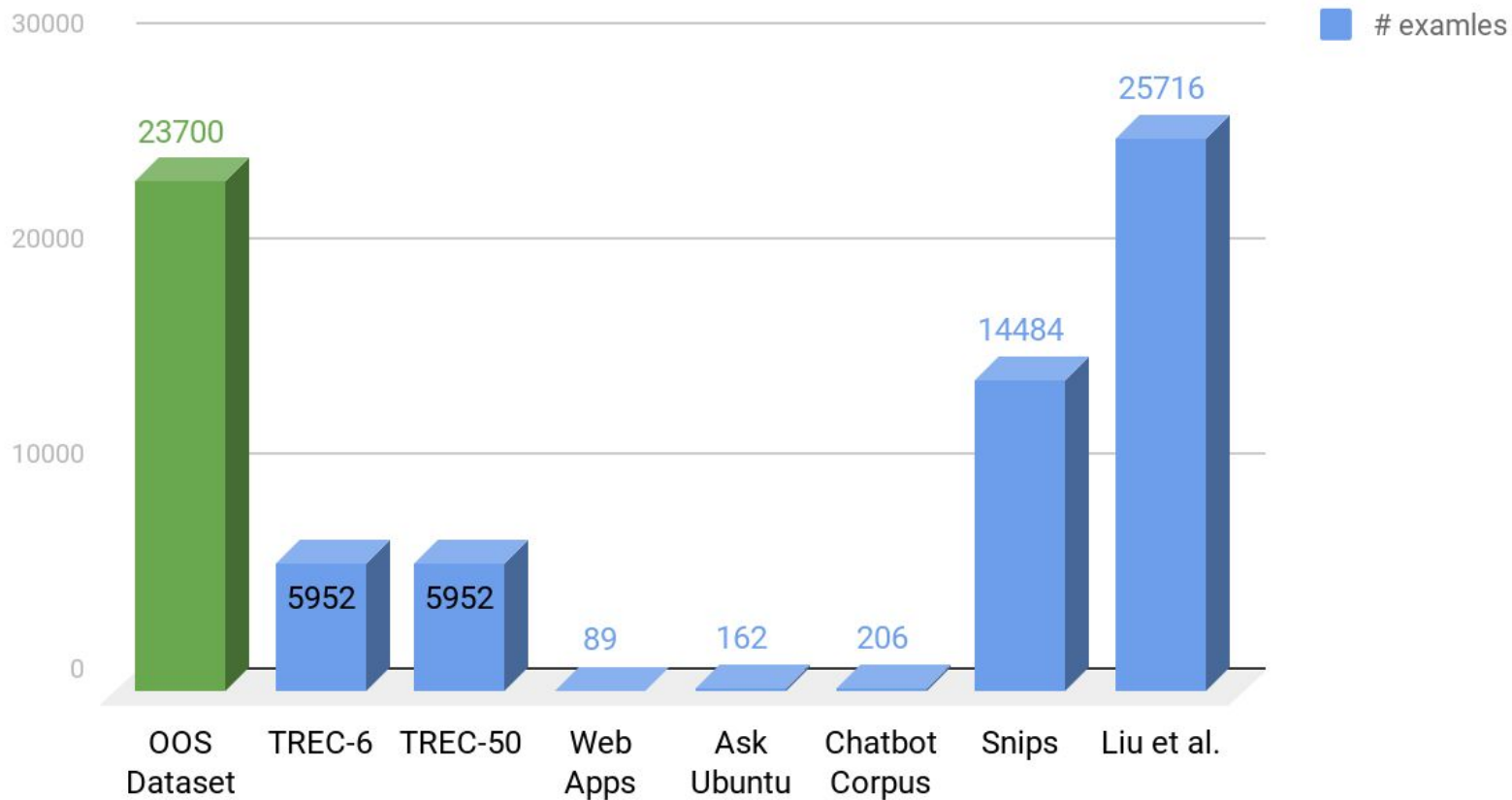
intents). Furthermore, such resources do not facilitate analysis of *out-of-scope* queries: queries that users may reasonably make, but fall outside of the scope of the system-supported intents.

Figure 1 shows example query-response exchanges between a user and a task-driven dialog system for personal finance. In the first user-system exchange, the system correctly identifies the user's intent as an in-scope BALANCE query. In the second and third exchanges, the user queries with out-of-scope inputs. In the second exchange, the system incorrectly identifies the query as in-scope and yields an unrelated response. In the third exchange, the system correctly classifies the user's query as out-of-scope, and yields a fallback response.
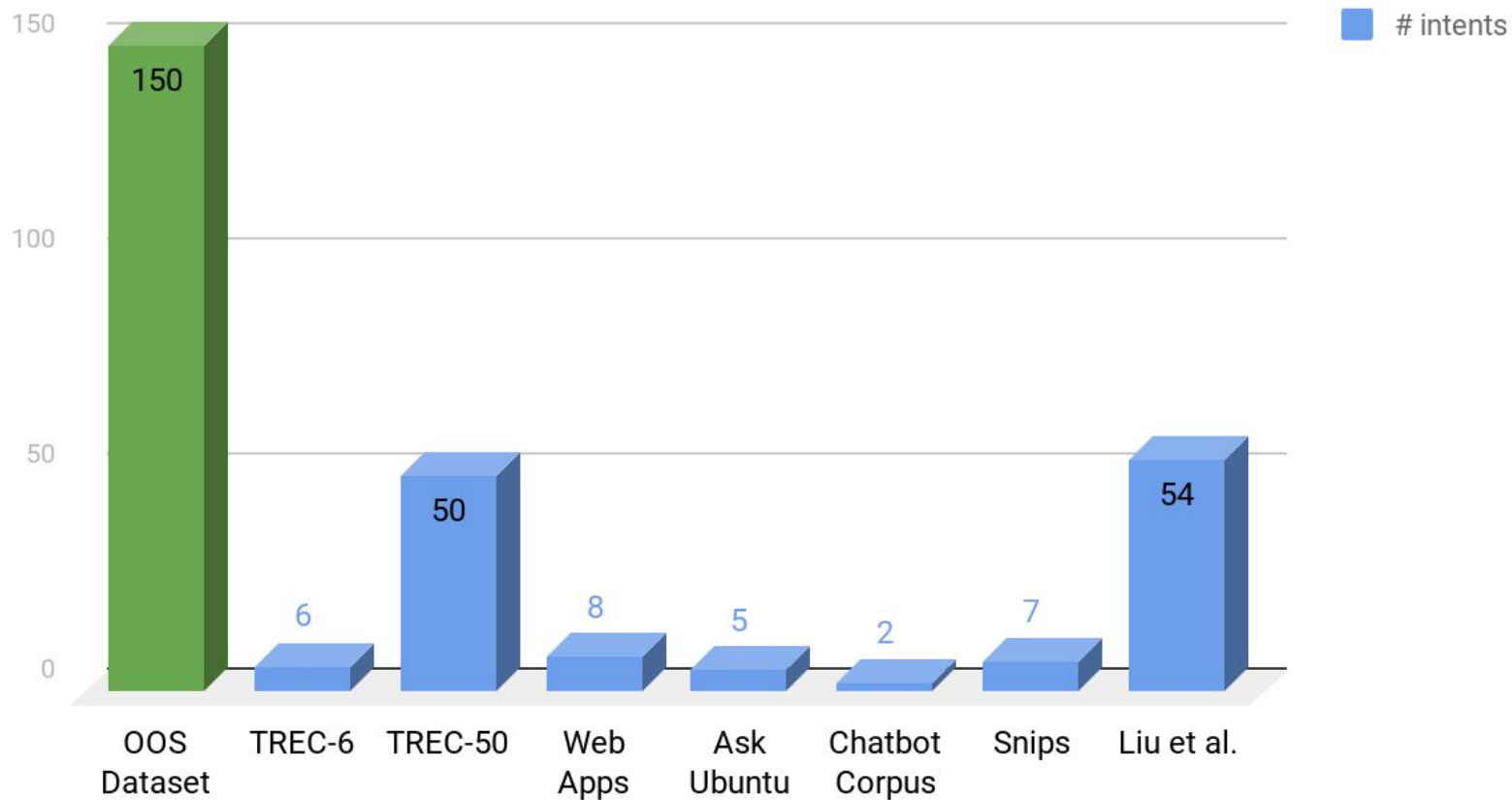
Out-of-scope queries are inevitable for a task-oriented dialog system, as most users will not be fully cognizant of the system's capabilities, which are limited by the fixed number of intent classes.

| Domain | Intent | Query |
|---|---|---|
| BANKING | TRANSFER | *move 100 dollars from my savings to my checking* |
| WORK | PTO REQUEST | *let me know how to make a vacation request* |
| META | CHANGE LANGUAGE | *switch the language setting over to german* |
| AUTO & COMMUTE | DISTANCE | *tell the miles it will take to get to las vegas from san diego* |
| TRAVEL | TRAVEL SUGGESTION | *what sites are there to see when in evans* |
| HOME | TODO LIST UPDATE | *nuke all items on my todo list* |
| UTILITY | TEXT | *send a text to mom saying i'm on my way* |
| KITCHEN & DINING | FOOD EXPIRATION | *is rice ok after 3 days in the refrigerator* |
| SMALL TALK | TELL JOKE | *can you tell me a joke about politicians* |
| CREDIT CARDS | REWARDS BALANCE | *how high are the rewards on my discover card* |
| OUT-OF-SCOPE | OUT-OF-SCOPE | *how are my sports teams doing* |
| OUT-OF-SCOPE | OUT-OF-SCOPE | *create a contact labeled mom* |
| OUT-OF-SCOPE | OUT-OF-SCOPE | *what's the extended zipcode for my address* |

# Dataset Size



| Dataset | # examples |
| --- | --- |
| OOS Dataset | 23700 |
| TREC-6 | 5952 |
| TREC-50 | 5952 |
| Web Apps | 89 |
| Ask Ubuntu | 162 |
| Chatbot Corpus | 206 |
| Snips | 14484 |
| Liu et al. | 25716 |

# Number of Intents

| Dataset | Out-of-Scope Utterances |
|---|---|
| OOS Dataset | ✔️ |
| TREC-6 | ❌ |
| TREC-50 | ❌ |
| Web Apps | ❌ |
| Ask Ubuntu | ❌ |
| Chatbot Corpus | ❌ |
| Snips | ❌ |
| Liu et al. | ❌ |

## OOS Train

Train Out-of-Scope as 151st intent

## OOS Threshold
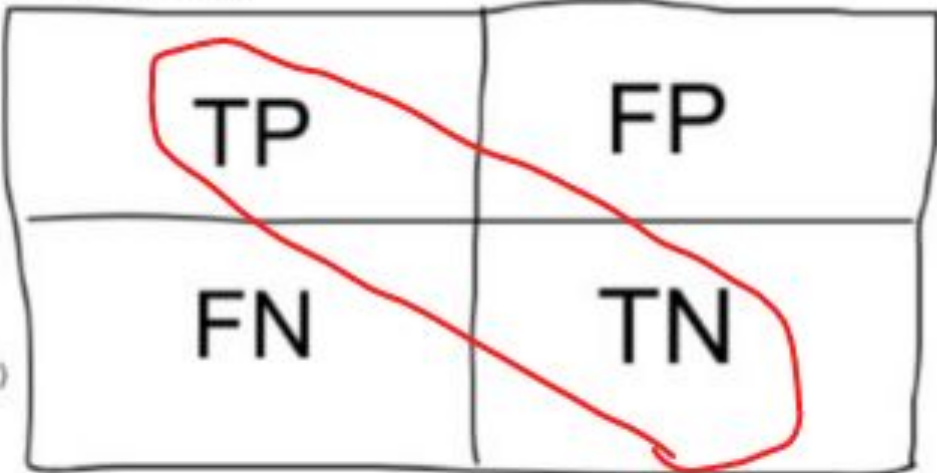
A threshold on the classifier's probability estimate

## OOS Binary

Two-stage process where we

first classify a query
as in- or
out-of-scope,

then classify it into one of
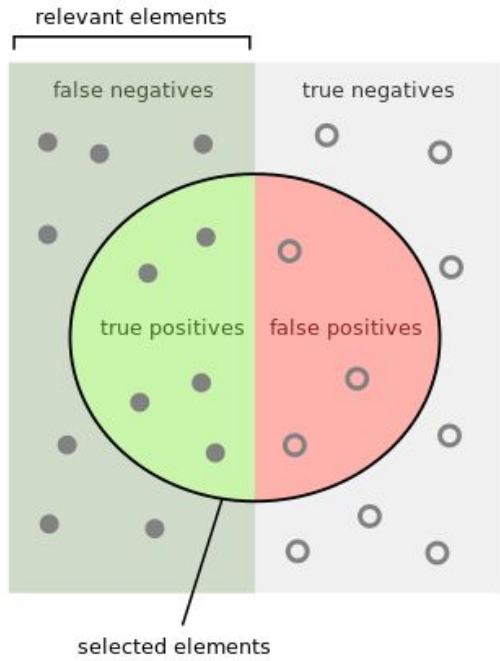the 150 intents if classified as
in-scope

Accuracy $= \dfrac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}}$

Actual

Positives(1)  Negatives(0)

Predicted

Positives(1)

TP    FP

Negatives(0)

FN    TN

$$\text{Recall} = \frac{TP}{TP + FN}$$

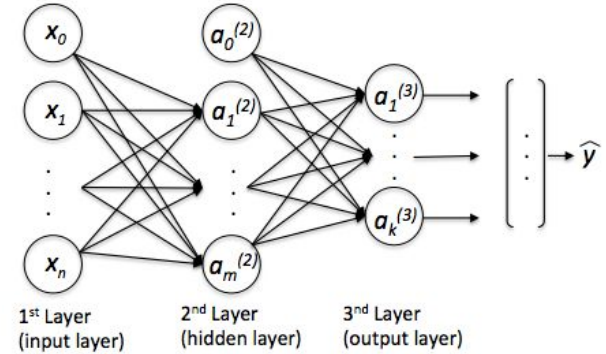| Baseline Models |
| --- |
| Support Vector Machine |
| Multilayer Perceptron |
| FastText |
| Convolutional Neural Network |
| Bert |
| Google's DialogFlow |
| Rasa NLU |

# Support Vector Machine

# Multilayer Perceptron

Universal Sentence Encoding

"How old are you?"
"What is your age?"
"My phone is good."
...

Embed

[0.3, 0.2, …]
[0.2, 0.1, …]
[0.9, 0.6, …]
...

$x_0$
$x_1$
$x_n$

$a_0^{(2)}$
$a_1^{(2)}$
$a_m^{(2)}$

$a_1^{(3)}$
$a_k^{(3)}$

$\widehat{y}$

1st Layer
(input layer)

2nd Layer
(hidden layer)

3nd Layer
(output layer)

# FastText



Document classes

Softmax

Document vector

Averaging

Word vector $x_1$

Word vector $x_2$

Word vector $x_N$

Document

# Convolutional Neural Network



wait
for
the
video
and
do
n't
rent
it

$n \times k$ representation of sentence with static and non-static channels

Convolutional layer with multiple filter widths and feature maps

Max-over-time pooling

Fully connected layer with dropout and softmax output

# BERT

# Google's DialogFlow

# Rasa NLU

# Baseline Results

OOS-Train

| | In-Scope Accuracy | Out-Of-Scope Recall |
|---|---|---|
| SVM | 91 | 14,5 |
| MLP | 93,5 | 47,4 |
| FastText | 89 | 9,7 |
| CNN | 91,2 | 18,9 |
| BERT | 96,9 | 40,3 |
| DialogFlow | 91,7 | 14 |
| Rasa | 91,5 | 45,3 |

**OOS-Threshold**
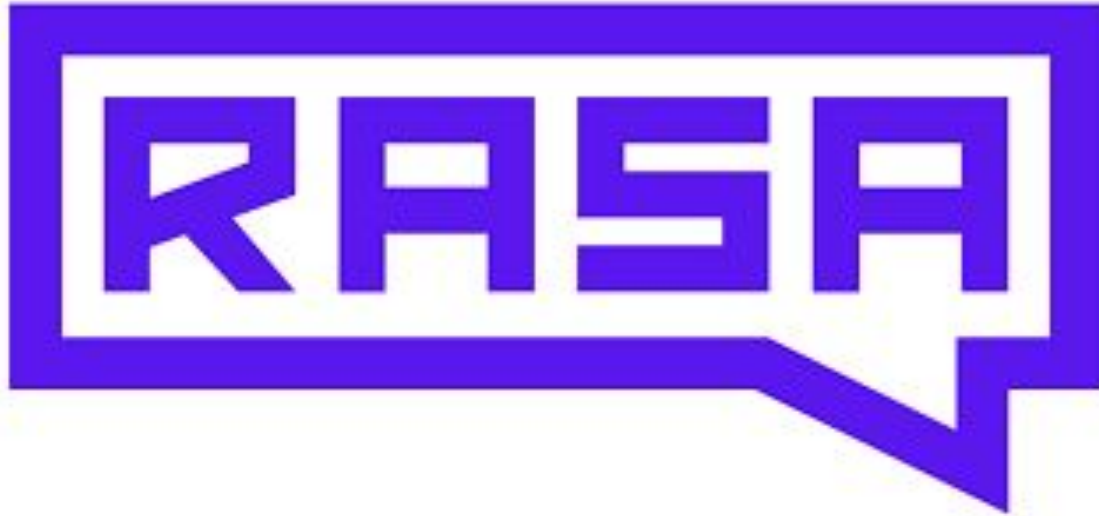
In-Scope Accuracy ▪ Out-Of-Scope Recall

| Model | In-Scope Accuracy | Out-Of-Scope Recall |
|---|---|---|
| SVM | 88,2 | 18 |
| MLP | 93,4 | 49,1 |
| FastText | 88,6 | 28,3 |
| CNN | 90,9 | 30,9 |
| BERT | 96,2 | 52,3 |
| DialogFlow | 90,8 | 26,7 |
| Rasa | 90,9 | 31,2 |

OOS-Binary

| | In-Scope Accuracy | Out-Of-Scope Recall |
|---|---|---|
| SVM | 88,4 | 32,2 |
| MLP | 90,1 | 52,8 |
| FastText | 88,1 | 22,7 |
| CNN | 89,8 | 25,6 |
| BERT | 94,4 | 46,5 |
| DialogFlow | 84,7 | 37,3 |
| Rasa | 87,5 | 37,7 |

Comparison of Approaches

In-Scope Accuracy
Out-Of-Scope Recall

| Approach | In-Scope Accuracy | Out-Of-Scope Recall |
|---|---|---|
| OOS-Train (MLP) | 93,5 | 47,4 |
| OOS-Threshold (BERT) | 96,2 | 52,3 |
| OOS-Binary (MLP) | 90,1 | 52,8 |

# My Initial Work

How
are
you

Dogs
are
fast

fastText

fastText

fastText

fastText

fastText

fastText

Fixed Dimension

Fixed Dimension

Cosine Similarity → Similarity Score

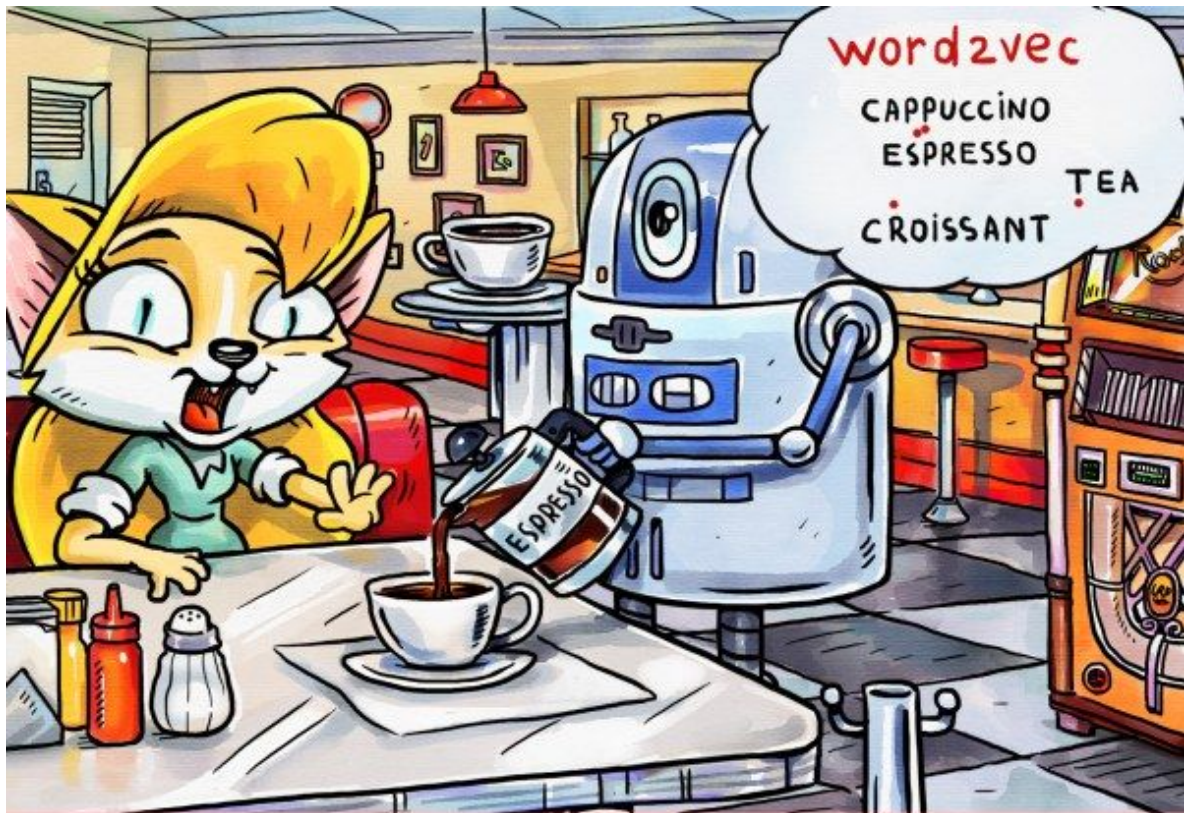$$similarity(A,B) = \frac{A \cdot B}{\|A\| \times \|B\|} = \frac{\sum_{i=1}^{n} A_i \times B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \times \sqrt{\sum_{i=1}^{n} B_i^2}}$$

$$sim(A,B) = \cos(\theta) = \frac{A \cdot B}{\|A\|\|B\|}$$

Burger

Sandwich

$= 0.6$

10

5

5

10

15

I would like to open account

Cos sim.

I need new account — Account 0.7

How to open account — Account 0.8

I would like a mortgage — Mortgage 0.5

Can I have mortgage — Mortgage 0.4

What's my balance — Balance 0.4

Tell me about sports — OOS 0.3

OOS-Train

In-Scope Accuracy
Out-Of-Scope Recall

| | Cos Sim. | SVM | MLP | FastText | CNN | BERT | DialogFlow | Rasa |
|---|---|---|---|---|---|---|---|---|
| In-Scope Accuracy | 67 | 91 | 93,5 | 89 | 91,2 | 96,9 | 91,7 | 91,5 |
| Out-Of-Scope Recall | 8 | 14,5 | 47,4 | 9,7 | 18,9 | 40,3 | 14 | 45,3 |

# Out of scope



| | | |
|---|---|---|
| ● translate | Y change_speed | ■ meeting_schedule |
| ● transfer | Y tire_pressure | ■ ingredients_list |
| ● timer | Y no | ■ report_fraud |
| ● definition | Y apr | ✱ measurement_conversion |
| ● meaning_of_life | Y nutrition_info | ■ smart_home |
| ● insurance_change | Y calendar | ■ book_hotel |
| ● find_phone | ■ uber | ■ current_location |
| ● travel_alert | ■ calculator | ■ weather |
| ● pto_request | ■ date | ✱ taxes |
| ● improve_credit_score | ■ carry_on | ■ min_payment |
| ● fun_fact | ■ pto_used | ■ whisper_mode |
| ● change_language | ■ schedule_maintenance | ■ cancel |
| ● payday | ■ travel_notification | ■ international_visa |
| ● replacement_card_duration | ■ sync_device | ★ vaccines |
| ● time | ■ thank_you | ★ pto_balance |
| ● application_status | ■ roll_dice | ★ directions |
| ● flight_status | ■ food_last | ★ spelling |
| ● flip_coin | ■ cook_time | ★ greeting |
| ● change_user_name | ■ reminder_update | ★ reset_settings |
| ▼ where_are_you_from | ■ report_lost_card | ★ what_is_your_name |
| ▼ shopping_list_update | ■ ingredient_substitution | ★ direct_deposit |
| ▼ what_can_i_ask_you | ■ make_call | ★ interest_rate |
| ▼ maybe | ■ alarm | ★ credit_limit_change |
| ▼ oil_change_how | ■ todo_list | ★ what_are_your_hobbies |
| ▼ restaurant_reservation | ■ change_accent | ★ book_flight |
| ▼ balance | ⬣ w2 | ★ shopping_list |
| ▼ confirm_reservation | ⬣ bill_due | ★ text |
| ▼ freeze_account | ⬣ calories | ★ bill_balance |
| ▼ rollover_401k | ⬣ damaged_card | ★ share_location |
| ▼ who_made_you | ⬣ restaurant_reviews | ★ redeem_rewards |
| ▼ distance | ⬣ routing | ★ play_music |
| ▼ user_name | ⬣ do_you_have_pets | ★ calendar_update |
| ▼ timezone | ⬣ schedule_meeting | ⬣ are_you_a_bot |
| ▼ next_song | ⬣ gas_type | ⬣ gas |
| ▼ transactions | ⬣ plug_type | ⬣ expiration_date |
| ▼ restaurant_suggestion | ⬣ tire_change | ⬣ update_playlist |
| ▼ rewards_balance | ⬣ exchange_rate | ⬣ cancel_reservation |
| ▼ pay_bill | ⬣ next_holiday | ⬣ tell_joke |
| Y spending_history | ⬣ change_volume | ⬣ change_ai_name |
| Y pto_request_status | ⬣ who_do_you_work_for | ⬣ how_old_are_you |
| Y credit_score | ⬣ credit_limit | ⬣ car_rental |
| Y new_card | ⬣ how_busy | ⬣ jump_start |
| Y lost_luggage | ⬣ accept_reservations | ⬣ meal_suggestion |
| Y repeat | ⬣ order_status | ⬣ recipe |
| Y mpg | ■ pin_change | ⬣ income |
| Y oil_change_when | ■ goodbye | ⬣ order |
| Y yes | ■ account_blocked | ⬣ traffic |
| Y travel_suggestion | ■ what_song | ⬣ order_checks |
| Y insurance | ■ international_fees | ⬣ card_declined |
| Y todo_list_update | ■ last_maintenance | ✕ oos |
| Y reminder | | |

Average