

Introducing a new 150-intent dataset for evaluating *out-of-scope* prediction performance of intent classification systems

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
stefan@clinc.com ◦ www.clinc.com ◦ github.com/clinc/oos-eval



Task-driven dialog systems need to handle user queries that land within the set of supported intents, as well as queries that fall out-of-scope.

Out-of-scope queries are those that a user may (un)reasonably make given limited knowledge of a system, but do not fall into any of the system’s supported intents.

Consider the exchanges below, which are from a task-driven dialog system for banking:

① correct in-scope classification

What is my balance?

You have \$1,847.51 across your 3 accounts.



② out-of-scope misclassified as in-scope

Which team plays at the Barclays Center?

Your last payday was on the 1st of November.



③ out-of-scope identified correctly

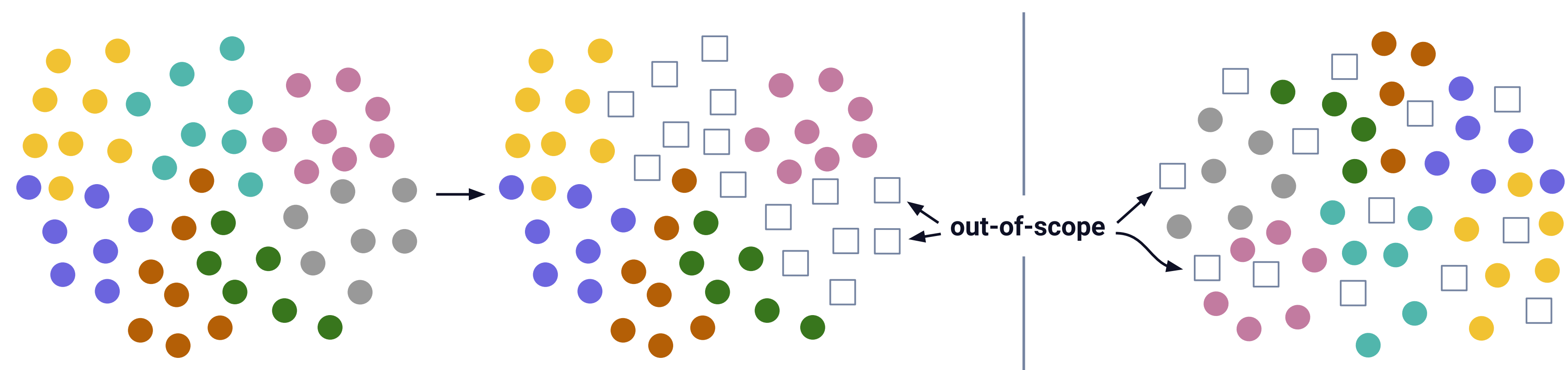
Who has the best record in the NBA?

Sorry, I can only answer questions about banking.



Our new dataset targets the task of intent classification in the presence of realistic, natural out-of-scope queries.

Other datasets and analyses either do not consider out-of-scope queries, or use held-out intents for out-of-scope.



Prior work creates out-of-scope data by holding out in-scope data, creating a single coherent out-of-scope intent.

Our work bootstraps a new dataset where the out-of-scope class is distributed, and not clearly defined.

Dataset Examples

Domain	Intent	Example
Banking	Balance	<i>how much do i have in my savings account</i>
Home	Todo List	<i>is vacuuming on my list of things to do</i>
Travel	Book Flight	<i>book a flight to los angeles from las vegas on american airlines</i>
Utility	Text	<i>i need leslie to be texted saying have a good day</i>
Work	PTO Used	<i>i want to know how many vacation days i have used</i>
Small Talk	Tell Joke	<i>can you tell me a joke about politicians</i>
Out-of-Scope	Out-of-Scope	<i>how many prime numbers are there between 0 and 100</i>
Out-of-Scope	Out-of-Scope	<i>how many students attend ucsb</i>
Out-of-Scope	Out-of-Scope	<i>change my seat to a window seat for my flight on monday</i>
Out-of-Scope	Out-of-Scope	<i>set a warning for when my bank account starts running low</i>

Benchmark

Classifier	In-Scope	OOS Recall
FastText	89.0	9.7
SVM	91.0	14.5
CNN	91.2	18.9
Rasa	91.5	45.3
DialogFlow	91.7	14.0
MLP + USE	93.5	47.4
BERT	96.9	40.3



github.com/clinc/oos-eval