# Submodularity-Inspired Data Selection for Goal-Oriented Chatbot Training Based on Sentence Embeddings

Mladen Dimovski<sup>1</sup>, Claudiu Musat<sup>2</sup>, Vladimir Ilievski<sup>1</sup>, Andreea Hossmann<sup>2</sup>, Michael Baeriswyl<sup>2</sup> <sup>1</sup>IC School, EPFL, Switzerland <sup>2</sup>Artificial Intelligence Group, Swisscom

BAD

CHOICE

GOOD

CHOICE



Slot Filling	Data Availability	Contributions
Are there any French restaurants in downtown Toronto ?	<ul> <li>Existing BiLSTM models perform reasonably well if given enough training data</li> <li>What if we can afford to label only small amount of data ?</li> </ul>	We show that the space of raw, unlabeled sentences contains information that we can use to choose the sentences to label
BiLSTM model	MAINIDEA	We create a submodularity-inspired ranking

Can you propose a good restaurant serving beef in the city center?

What is the best rated chinese restaurant in Lausanne?



0	0	0	<b>B-Cuisine</b>
Are	there	any	French
0	0	<b>B-Location</b>	I-Location
restaurants	in	downtown	Toronto

according to their usefulness and select only the best ones for labeling

![](_page_0_Picture_8.jpeg)

I want to eat sushi!

I want to eat pizza!

selecting the most useful sentences to label

3

We apply this data selection method to the problem of slot filling and prove that the model's performance can be considerably better with training samples chosen in an intelligent way

## Sentence embeddings and sentence similarity

![](_page_0_Figure_14.jpeg)

- Use a recently developed technique, sent2vec, that produces continuous vector representations of sentences
- **Define the similarity between two sentences** *X* and *Y* as:  $sim(x, y) = exp(-\beta ||e(x) - e(y)||_2)$
- where e(x) is the embedding of the sentence X and  $\beta$  is the inverse of the average distance between all pairs of embeddings
- Hypothesis : closeness in the embedding space is in line with the human perception of sentence similarity

#### What was the movie that featured Over the Rainbow?

Find me the movie with the song Over the Rainbow
What movie was the song Somewhere Out There featured in ?

 We experiment with 3 different publicly available datasets, each containing few thousands of sentences

![](_page_0_Picture_22.jpeg)

![](_page_0_Picture_23.jpeg)

![](_page_0_Picture_24.jpeg)

MIT Restaurant Dataset

t MIT Movie Dataset ATIS Dataset

Using some selection criteria, we select only a few dozen sentences, we reveal their labels and use them to train our model. This simulates the behavior of a system that needs to be trained for a newly available domain and we refer to it as the LOW-DATA regime.

![](_page_0_Picture_29.jpeg)

• What movie features the song Hakuna Matata?

Performance metric: F1 score on a separate test set

## Data selection methods

Use a well-chosen submodular function to evaluate the usefulness of each subset of sentences X

A function  $F : 2^V \to \mathbb{R}$  is called submodular if the value of an element diminishes as the context in which it is considered grows. Formally, F is submodular if  $F(s|Y) \leq F(s|X)$  for every  $X \subseteq Y$  and every  $s \in V \setminus Y$ .

 $- p_{\Theta}(\hat{y}_{w_i}|x))$ 

#### Baselines

- Random Data Selection
- Classic Active Learning
- Randomized Active Learning

Uncertainty of the model for a sentence X of k words

Coverage Score  

$$C(X) = \sum_{x \in X} \sum_{y \in V} sim(x, y)$$
  $C(s|X) = \sum_{y \in V} sim(s, y)$ 

inear-Penalty Marginal Gain  
$$D(s|X) = \sum_{y \in V} sim(s, y) - \alpha \sum_{x \in X} sim(s, x)$$

#### Ratio-Penalty Marginal Gain

$$F(s|X) = rac{\sum_{y \in V} sim(s, y)}{1 + \sum_{x \in X} sim(s, x)}$$

![](_page_0_Figure_44.jpeg)

#### Top 3 Ratio-Penalty Sentences (MIT Restaurant Dataset)

- I need to find somewhere to eat something close by, I'm really hungry. Can you see if there's any buffet style restaurants within 5 miles of my location?
- Find zenna noodle bar restaurant with kids friendly amenity around this place
- Where is the nearest 5 star restaurant that serves Italian food?

### Results

![](_page_0_Figure_50.jpeg)

For more information, see: M. Dimovski et al., Submodularity-Inspired Data Selection for Goal-Oriented Chatbot Training Based on Sentence Embeddings, IJCAI 2018