



# **Discover Users' Specific Geo Intent in Web Search**

**Xing Yi<sup>1</sup>, Hema Raghavan<sup>2</sup> and Chris Leggetter<sup>2</sup>**

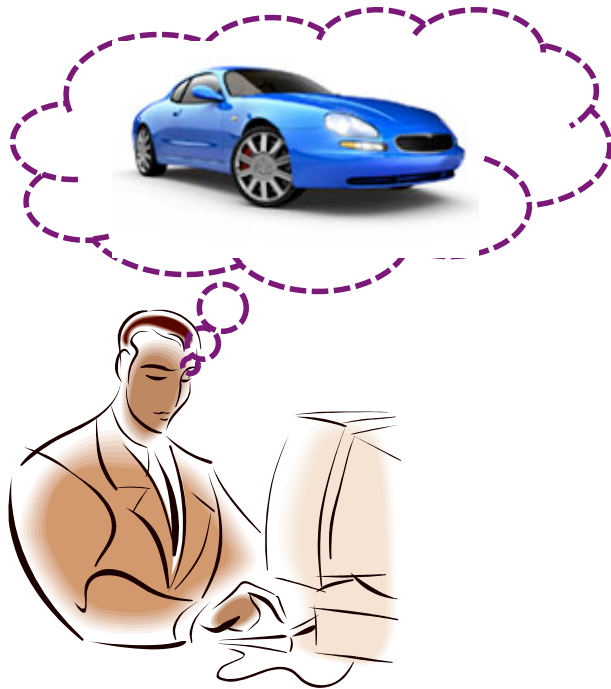
<sup>1</sup> Center for Intelligent Information Retrieval, University of Massachusetts Amherst

<sup>2</sup> Yahoo! Labs, 4401 Great America Pky, Santa Clara, CA

**YAHOO!**



# Queries with Geo Intent



Web | Images | Video | Local | Shopping | more ▾

car dealer sunnyvale  Options ▾ Custom

Also try: [toyota car dealer sunnyvale](#), [More...](#)

**Sunnyvale Car Dealers** SF Bay Area, CA

Deals on award-winning models. Find a Toyota dealer near you.  
[www.BuyAToyota.com](http://www.BuyAToyota.com)

**Sunnyvale Car Dealer** SF Bay Area, CA

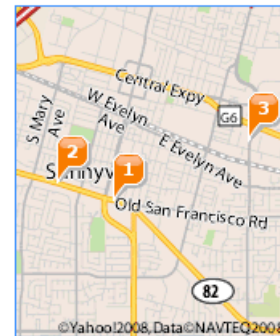
Find Great Offers on Chrysler® Vehicles at a Local Dealer Online.  
[www.ChryslerCalifornia.com](http://www.ChryslerCalifornia.com)

**California Dodge Dealer** SF Bay Area, CA

Find New Vehicle Deals at Your Local California Dodge Dealer.  
[www.DodgeCalifornia.com](http://www.DodgeCalifornia.com)

**California Jeep Dealer** SF Bay Area, CA

Learn About New Deals at Your Local California Jeep Dealer.  
[www.JeepCalifornia.com](http://www.JeepCalifornia.com)



© Yahoo! 2008, Data © NAVTEQ 2008  
[Yahoo! Shortcut](#) - [About](#)

**Car Dealer near Sunnyvale**

[local.yahoo.com](http://local.yahoo.com)

1. **Car Concepts**  
(408) 733-1000 - 310 W El Camino Real, Sunnyvale, CA  
[Get Directions](#) | [Official site](#)
2. **Toyota at Sunnyvale** - ★★★★★ (45)  
(408) 245-6640 - 898 W El Camino Real, Sunnyvale, CA  
[Get Directions](#) | [Reviews](#) | [Official site](#)
3. **Dannicks Auto Care** - ★★★★★ (3)  
(408) 732-4222 - 135 N Wolfe Rd, #40, Sunnyvale, CA  
[Get Directions](#) | [Reviews](#) | [Official site](#)

[More Results...](#)



# Geo Intent

**Geo intent** – a user's information need has some kind of entity which has a geographic (**geo**) location associated with it:

- **explicit:** “one bedroom apartment new york city”, “madrid guided tour”

- An **explicit** geo query has two portions:

e.g. “**car dealer** in **sunnyvale**”,

  
*non location part*

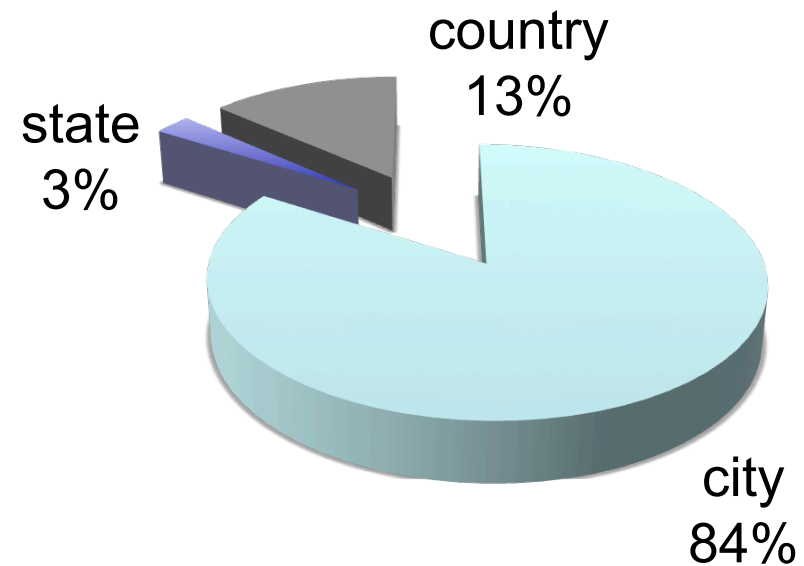
  
*location part*

- **implicit:** “pizza delivery”, “dental care”, “day care”, “rockefeller center”



## Observations about Web Geo queries

- many web queries contain geo info
  - About 13-14% queries have a place name (Jones *et al.*, Intl. J. of G.I. Science 2008, Sanderson & Kohler, SIGIR GIR workshop 2004)
  - About 30% queries may have geo intent; **only about half** of them have explicit geo info. ( Welch & Cho, SIGIR 2008)





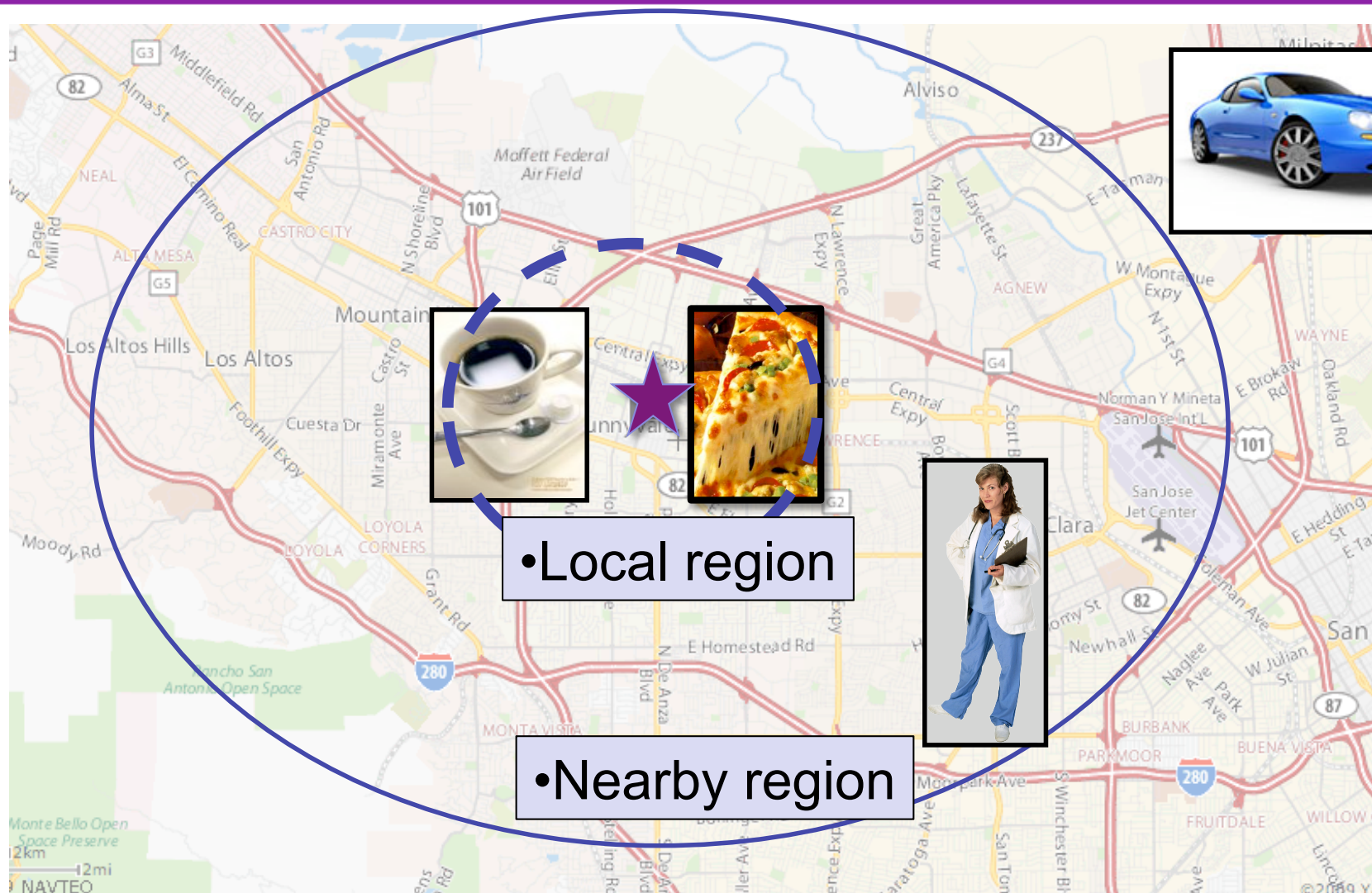
## Research questions/tasks

---

- ① Given a set of queries with no mentions of any location (town/city/state), can we predict which of these have **implicit** geo intent?



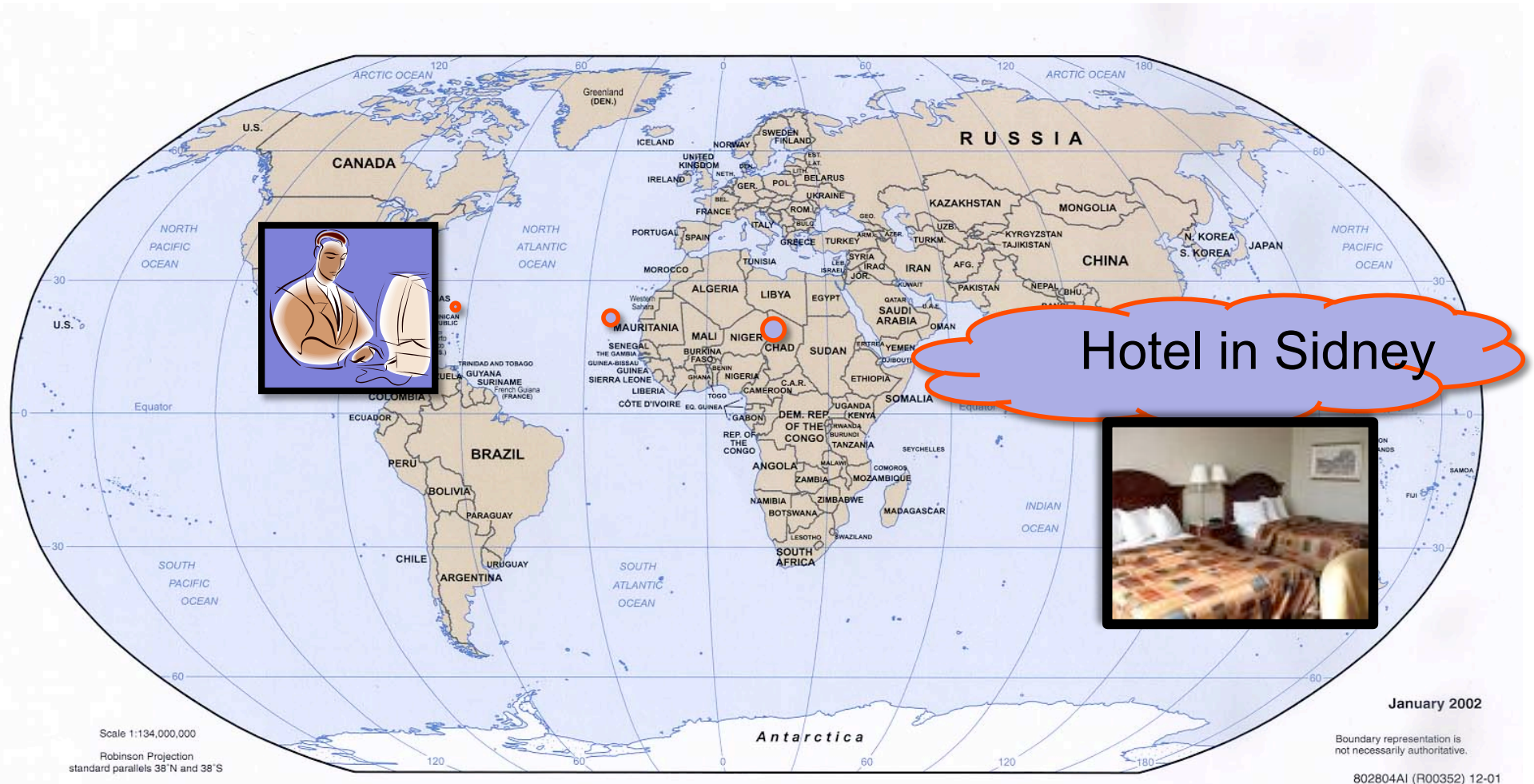
# Localization capability of a geo-intent query







# Queries with Geo intent but not localized





## Research questions/tasks

---

- ① Given a set of queries with no mentions of any location (town/city/state), can we predict which of these have **implicit** geo intent?
- ② What is the **localization capability** of a geo-intent query?
- ③ What is the **city** corresponding to the geo-intent?





# Applications

---

Benefits for finding users' geo intent:

- Personalizing web search results
- Better sponsored online advertisement matching

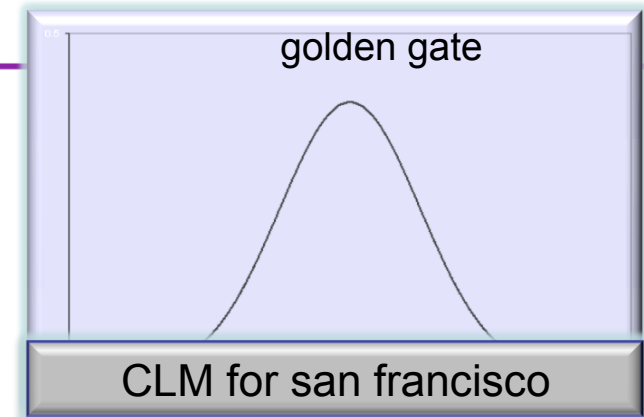
More benefits for finding users' specific geo intent at a fine-grained **city**/location level:

- Delivering more local goods and services
- Finding local news and events



# Outline of remainder of the talk

- Feature extraction
  - City Language Models
    - Entity Language Models (Raghavan *et al.* ACM LinkKDD 2004)
- Experiments for each of the 3 tasks:
  - Label generation : **millions of training samples from click data.**
  - Evaluation
- Conclusions and Future Work





# City Language Models

Feature Extraction

Use Internal tool

san francisco

query  $Q_{nc}$

freq

pizza

200

cheap hotel

150

49ers

125

zoo

100

golden gate

75

- Bigram language model
- Smoothed
- Details in the paper

$$P(q | C_k) = \prod_{i=1}^n P(w_i | w_1^{i-1}, C_k) \approx \prod_{i=1}^n P(w_i | w_{i-1}, C_k)$$



- Calculate the posteriors:

$$P(C_i | q) \propto P(C_i)P(q | C_i)$$

- These posteriors are used for predicting the locations for *location-specific* queries
- Top-10 posteriors are used as features for classifications



## Some examples

Feature Extraction

“Disney world ticket”		“Harvard University”	
City Name	$P(C_i   q)$	City Name	$P(C_i   q)$
Orlando	0.98011	Cambridge	0.63545
Kissimmee	0.01386	Princeton	0.05360
Anaheim	0.00240	Longwood	0.05334
New Castle	0.00135	Boston	0.01979
San Antonio	0.00044	Tuskegee	0.01719
...	...	...	...





# Geo Information Units (GIU)

Feature Extraction

san francisco	$Q_{nc}$	freq	Global information units
	pizza	200	
	cheap hotel	150	
	49ers	125	
	zoo	100	
	golden gate		

GIUs like ***pizza*** co-occur with many different city names



# Features based on GIUs

Feature Extraction

Examples:

- Probability  $P_g(w_i^{i+n-1})$  of a GIU appearing in geo queries
- Probability  $P(w_i^{i+n-1})$  of this GIU appearing in all the queries
- The pair-wise mutual information (PMI) between the  $w_i^{i+n-1}$  and each location

**Aggregate features and individual GIUs as features**



# Experiments

---

## Overall Data Description

Three learning tasks:

- Classifier I: Detecting implicit geo queries
- Classifier II: Discriminating different localization capabilities of geo queries: local geo intent, neighbor region geo intent, *etc.*
- City language models: Predicting geo entities related to a query

## Evaluations and Results



Slice of traffic:

- Training data
  - 1.44b queries in May 08
  - 96.2m are explicit geo queries (**training geo subset**)
- Testing data
  - 1.42b queries in June 08
  - 96.7m are explicit geo queries (**testing geo subset**)

Weakly supervised automatic label generation



# Generating labeled data for Classifier I

Task 1

training geo subset

san francisco

pizza

san francisco

cheap hotel

san francisco

49ers

san francisco

zoo

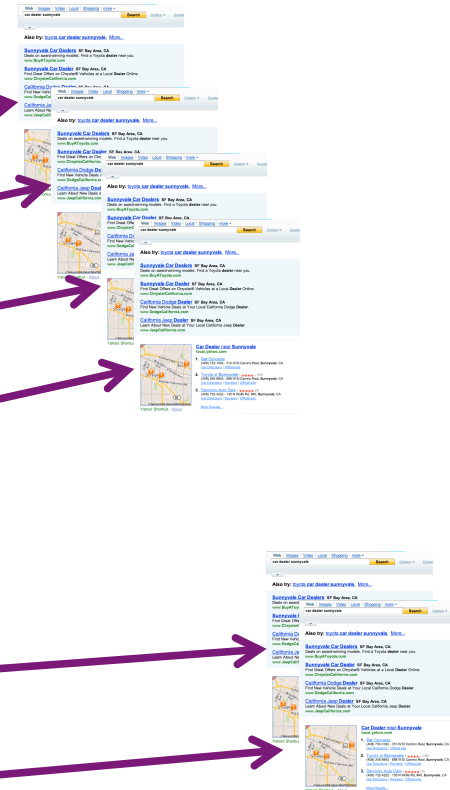
.....

new york

pizza

new york

zoo



Logs of what was shown to the  
users

Step 1: get the clicked url for each query (domain name)





# Generating labeled data for Classifier I

Task 1

**Step 2:**  $DN1$  is set of top 100 clicked domains from Step 1.

**Step 3:**  $DN2$  is set of top 100 clicked domains from queries in *training set* and not in *training geo subset*.

**Step 4:**

- $DN+ = DN1 \setminus DN2$        $DN- = DN2 \setminus DN1$
- If a query in training geo subset has clicked domain in  $DN+$   $\rightarrow$  positive sample
- *non-location* parts of positive samples as the final implicit geo intent queries.
- randomly sample 20,000 implicit geo queries and 20,000 non-geo queries to train classifiers



## Some examples in DN+ and DN-

Task 1

DN+ – DN as Positive label	DN- – DN as Negative label
<u><a href="http://www.citysearch.com">www.citysearch.com</a></u> <u><a href="http://www.yellowpages.com">www.yellowpages.com</a></u> <u><a href="http://local.yahoo.com">local.yahoo.com</a></u> <u><a href="http://www.local.com">www.local.com</a></u> <u><a href="http://travel.yahoo.com">travel.yahoo.com</a></u> <u><a href="http://www.tripadvisor.com">www.tripadvisor.com</a></u> <u><a href="http://www.yellowbook.com">www.yellowbook.com</a></u> <u><a href="http://www.city-data.com">www.city-data.com</a></u>	<u><a href="http://en.wikipedia.org">en.wikipedia.org</a></u> <u><a href="http://answers.yahoo.com">answers.yahoo.com</a></u> <u><a href="http://search-desc.ebay.com">search-desc.ebay.com</a></u> <u><a href="http://www.youtube.com">www.youtube.com</a></u> <u><a href="http://www.amazon.com">www.amazon.com</a></u> <u><a href="http://www.myspace.com">www.myspace.com</a></u> <u><a href="http://www.nextag.com">www.nextag.com</a></u>



# Generating labels: Test data

Task 1

Two testing subsets from testing data

## Testing data I:

- Same as training data process, but on testing data.
- 40,000 implicit geo queries + 40,000 negative queries

## Testing data II:

- Extract all queries that have DNs in DN+ or DN-.
- Remove all possible location information using WOE
- Sample 40,000 implicit geo queries + 40,000 negative queries
- May have queries that had implicit geo intent to begin with.



## Three classifiers

- Support Vector Machines (linear kernel and RBF gaussian kernel)
- Gradient boosting decision trees (Treenet)
- Multinomial Logistic Regression

## 5-fold cross-validation



## Results using CLM features + aggregated GIU features

Task 1

	P	R	A
<b>Testing Set I</b>			
SVM-Linear	<b>91.7%</b>	82.6%	87.6%
SVM-RBF	91.4%	86.0%	<b>89.0%</b>
Treenet	89.4%	<b>87.4%</b>	88.5%
Logistic-R	91.3%	83.5%	87.8%
<b>Testing Set II</b>			
SVM-Linear	<b>80.9%</b>	35.7%	63.7%
SVM-RBF	80.4%	36.2%	63.7%
Treenet	78.1%	<b>40.9%</b>	<b>64.7%</b>
Logistic-R	80.2%	36.4%	63.7%





## Using aggregate GIU stats as well as GIUs as individual features.

Task 1

	P	R	A
<b>Testing Set I</b>			
SVM-Linear	99.9%	66.0%	83.0%
SVM-RBF	98.5%	62.8%	80.9%
<b>Testing Set II</b>			
SVM-Linear	99.9%	48.8%	74.4%
SVM-RBF	97.8%	48.0%	73.5%

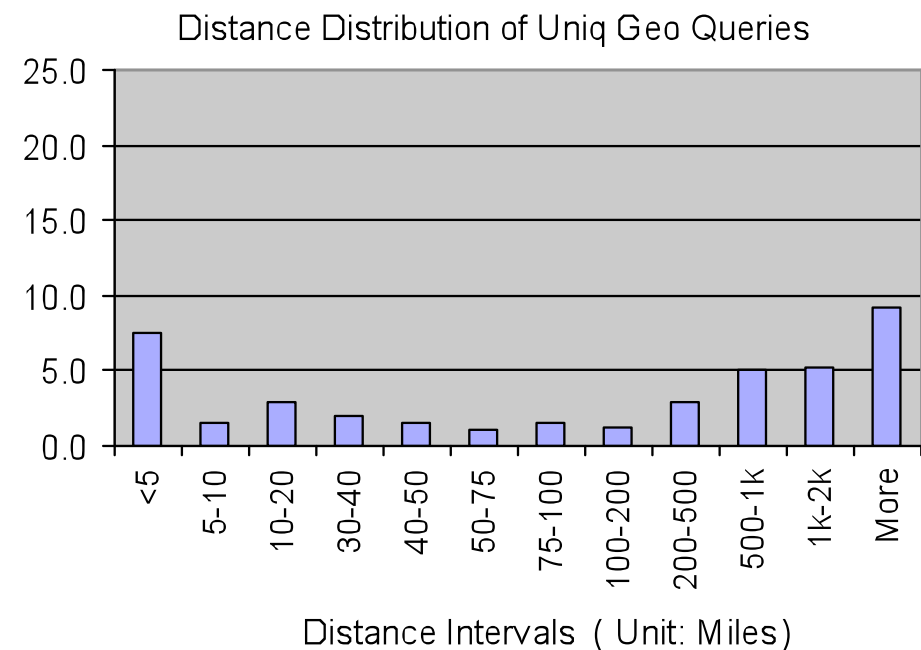
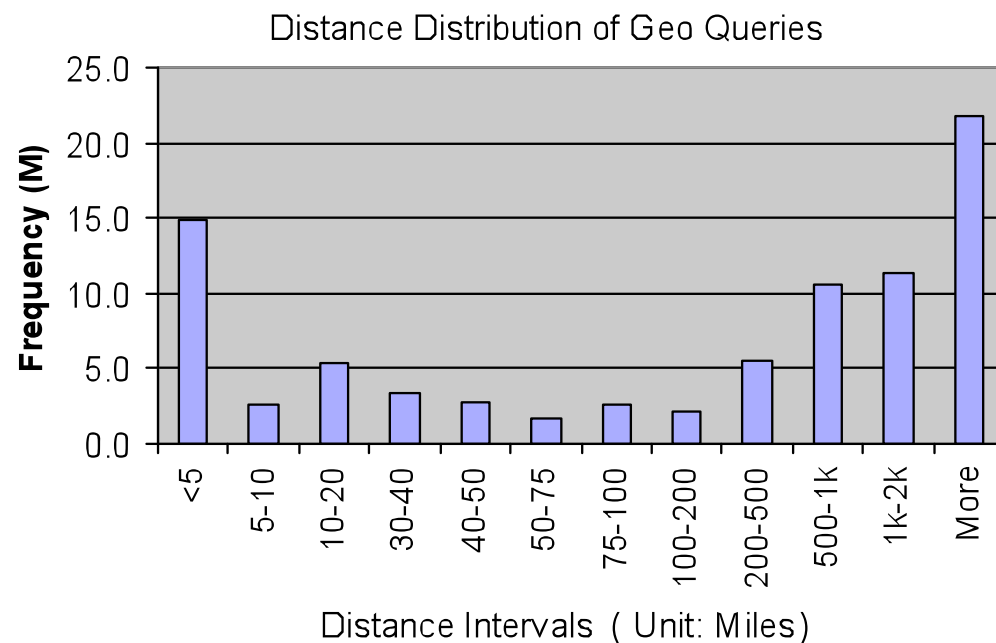


## Classifier II: Localization capability of a query

### Task 2

$L(q, C)$  = distance between the city  $C$  ( $Q_c$ ) in a query  $q$  ( $Q_{nc}$ ) and the user IP

$L^m(q)$  = median of all  $L(q, C) > 0$  for all cities  $C$  associated with  $q$ .





## Data generation for Classifier II

Task 2

- 3 classes:
  - $L^m(q) \leq 50$  miles  $\rightarrow$  q is a local geo query (LG)
  - $50 < L^m(q) < 100$  miles  $\rightarrow$  neighbor region query (NG)
  - other geo queries (OG)



# Results

## Task 2

	A (LG/OG)	B (LG/NG)	C (NG/OG)	D (ALL)
<b>Case I</b>	aggregate GIU features			
SVM-Linear	61.3%	53.5%	61.0%	42.6%
SVM-RBF	62.0%	53.9%	<b>61.8%</b>	43.2%
Treenet	<b>62.8%</b>	<b>54.2%</b>	60.8%	<b>44.1%</b>
Logistic-R	61.2%	53.4%	61.0%	42.6%
<b>Case II</b>	high dimensional features			
SVM-Linear	99.6%	97.2%	96.9%	87.0%
SVM-RBF	<b>99.6%</b>	<b>98.0%</b>	<b>98.0%</b>	<b>96.6%</b>



## Predict Locations for Location-Specific Queries

Task 3

- Queries with mentions of an entity that is directly associated with a location: eg., hotels, local tv and radio channels, local newspapers, universities etc.
  - “airport check metro airport” → detroit
  - “woodfield mall jobs” → schauburg
- Data is generated using several rules (refer paper).
- Top cities using the City Language Models ( $P(C_i | q)$ ) were taken as predictions.





# Human Evaluation: key points

Task 3

669 randomly sampled *location-specific* queries and their predicted related locations

Request annotators to answer two questions with `yes/no/?`:

- whether the selected query was a *location-specific* query (84.5 % inter annotator agreement)
- Whether the predicted location was correct (73% agreement)

Of queries marked location specific, **accuracy** of predicting a location was **84.5%**.



## Concluding Remarks

---

- Methods for
  - identifying users' implicit city-level geo intent.
  - discriminating different localization capabilities of geo queries.
  - predicting the city corresponding to the geo-intent in a *location-specific* query.
- The models are learned from large amounts of click-through data and involve little supervision.
- Future Work:
  - Incorporate our CLM into retrieval models.
  - Use geo intent analysis results for helping search engines provide better query suggestions.
  - Exploit other data sources

# Questions?



## THANK YOU !

YAHOO!