

Ontology Assisted Crowd Mining

Yael Amsterdamer

Susan B. Davidson*

Tova Milo

Slava Novgorodov

Amit Somech

Goals

- Crowd data sourcing is a powerful data procurement paradigm
- We would like harness it in order to allow users to:
 - specify their information needs in a declarative manner
 - efficiently mine the crowd for relevant data
 - obtain a concise list of answers that represent frequent, significant data patterns

A Motivating Example

Ann is planning a vacation in NYC with her family:

"I'm looking for activities to do at a child-friendly attraction in New York, and a good restaurant near by"

Answers should be relevant, represent popular recommendations, and may include additional advice from the crowd

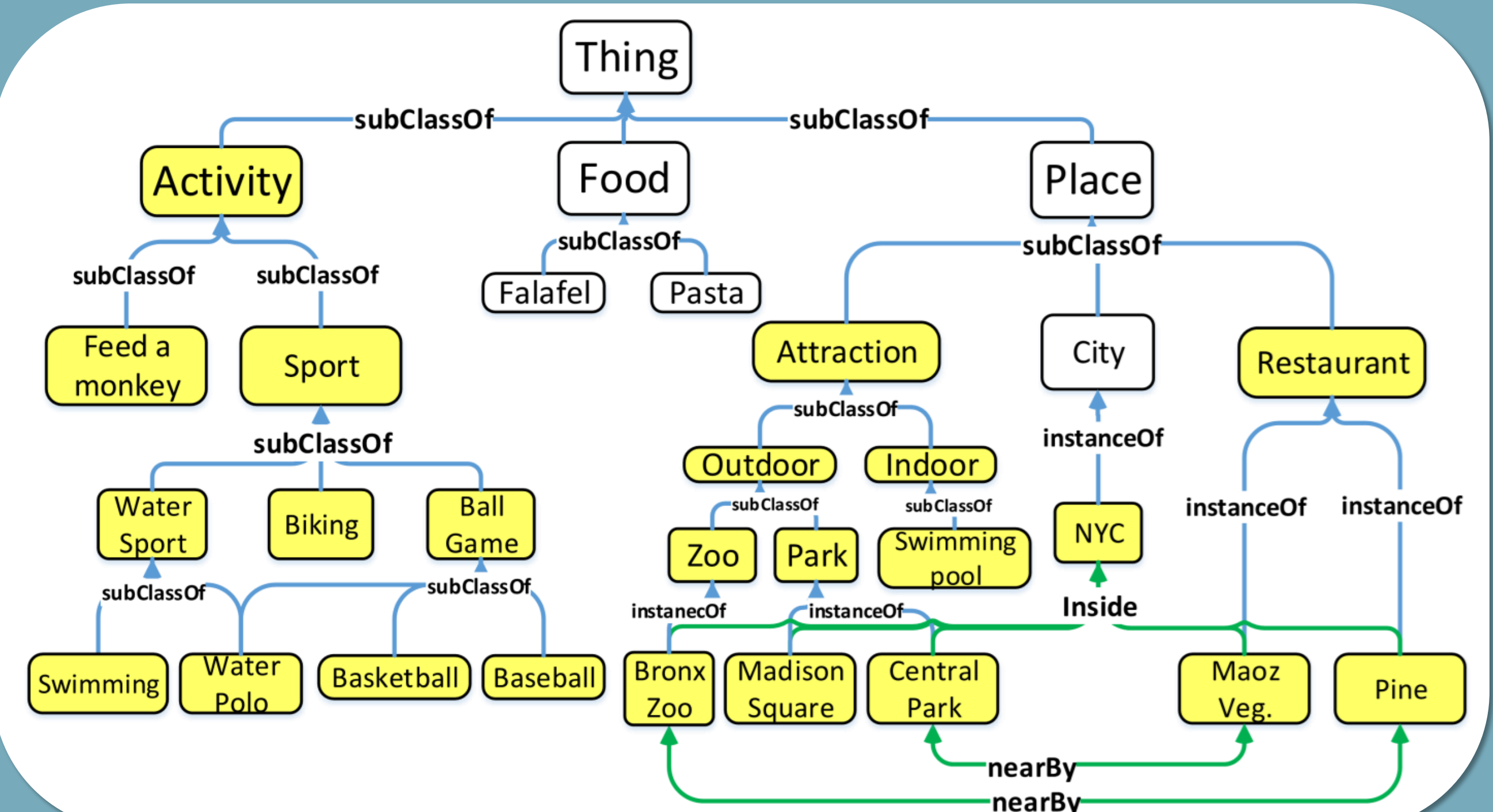
"You can go bike riding in Central Park and eat at Maoz Vegetarian.
Tips: Rent bikes at the boathouse"

"You can go visit the Bronx Zoo and eat at Pine Restaurant.
Tips: Order antipasti at Pine.
Skip dessert and go for ice cream across the street"

- Web search can only extract information from recorded data
 - It may be hard to extract combinations of restaurant and activity
- Forums can yield targeted answers,
 - which require reading, aggregating, identifying consensus, etc.

We propose a new, alternative approach using crowd mining

Formal Model Based on RDF



T1	I visited the Bronx Zoo and ate pasta at Pine on April 5th	[Visit doAt Bronx_Zoo]. [Pasta eatAt Pine]
T2	I played basketball in Central Park on April 13th	[Basketball playAt Central_Park]
T3	I played baseball in Central Park and ate falafel at Maoz Veg. on April 27th	[Baseball playAt Central_Park]. [Falafel eatAt Maoz_Veg]
...

OASSIS-QL: The Mining Language

```

1 SELECT VARIABLES
2 WHERE
3   {$w subClassOf* Attraction
4    $x instanceOf $w.
5    $x inside NYC.
6    $y subClassOf* Activity.
7    $z instanceOf Restaurant.
8    $z nearby $x}
9 SATISFYING
10  {$y+ doAt $x.
11   [] eatAt $z.
12   MORE}
13 WITH SUPPORT = 0.03
    
```

Evaluated over the ontology, to identify candidate data patterns

Retain the patterns that are significant for the crowd, and find additional advice

Efficient Query Evaluation Algorithm

- Lazily construct a semantic subsumption partial order
- Traverse it in a top-down manner
- Automatically generate crowd questions
- Prune insignificant parts

The goal is minimizing the # crowd questions

