



A PROBABILISTIC MODEL FOR FOOD RECOGNITION IN RESTAURANTS

Luis Herranz, Ruihan Xu, Shuqiang Jiang

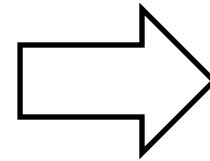
Institute of Computing Technology
Chinese Academy of Sciences

中国科学院计算技术研究所
Institute of Computing Technology, Chinese Academy of Sciences

Outline

- Introduction
- Proposed approach
- Experiments
- Conclusions

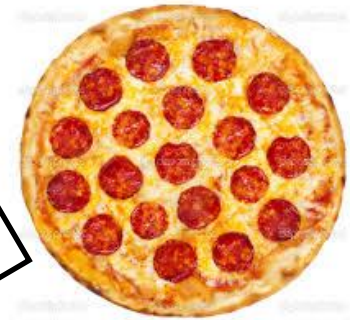
Food recognition



Pepperoni pizza

Why food recognition?

- Food images are everywhere
- Many potential applications



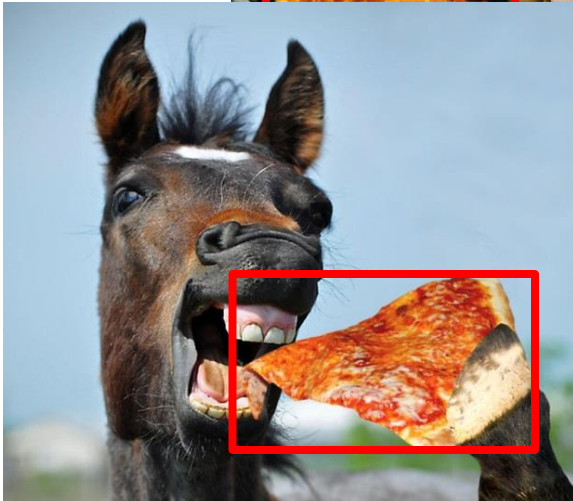
Retrieval

Pepperoni pizza

Nutrition Facts

Calories	298	(1246 kJ)
Calories from fat 109		
% Daily Value		
Total Fat	12.1g	19%
Sat. Fat	5.3g	26%
Cholesterol	29mg	10%
Sodium	683mg	28%
Total Carbs.	34g	11%
Dietary Fiber	1.6g	6%
Sugars	4.1g	
Protein	13.3g	
Calcium	168.5mg	
Potassium	185.8mg	

Annotation



Dietary/medical

Similarity and variability

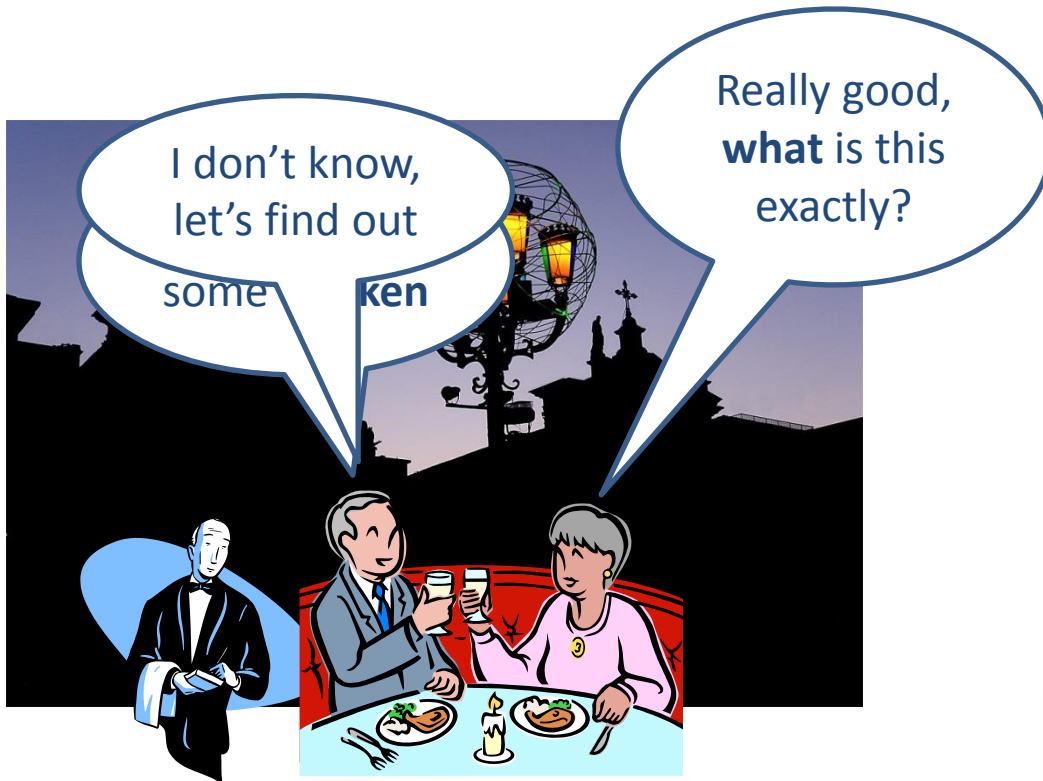
Beef vegetable stir-fry noodles



Spaghetti bolognese



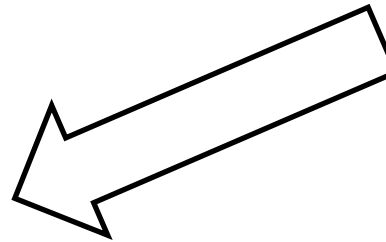
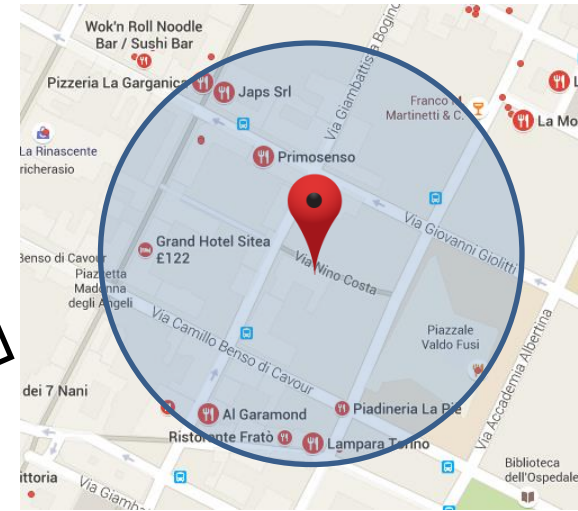
Dish recognition



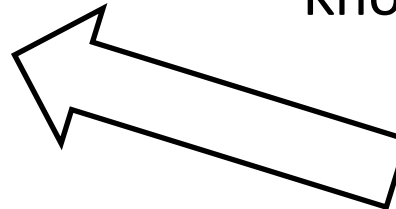
More accurate than user Requires expertise
Automatic (no interaction) Requires user interaction

Exploiting context and external data

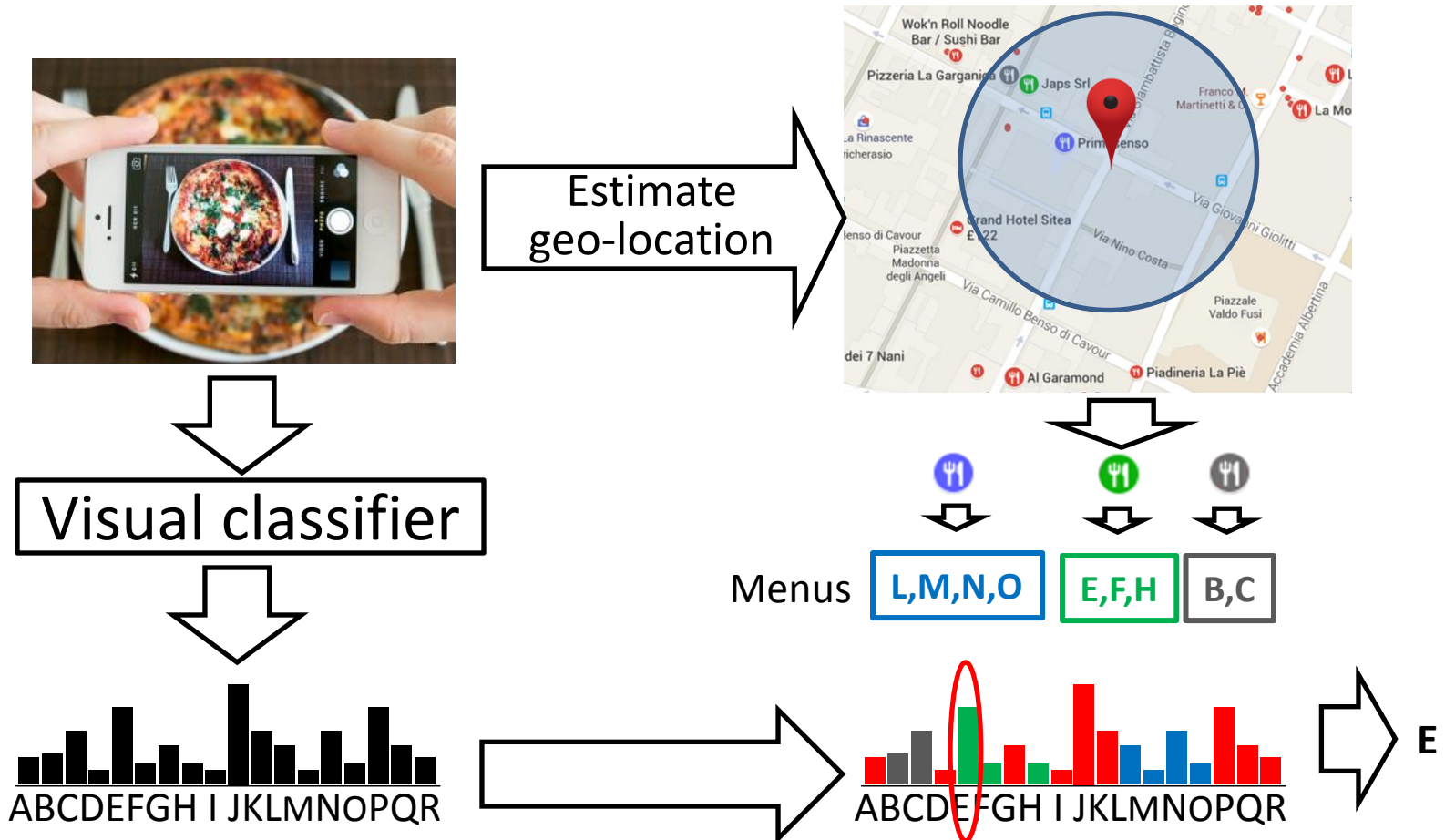
Geo-location



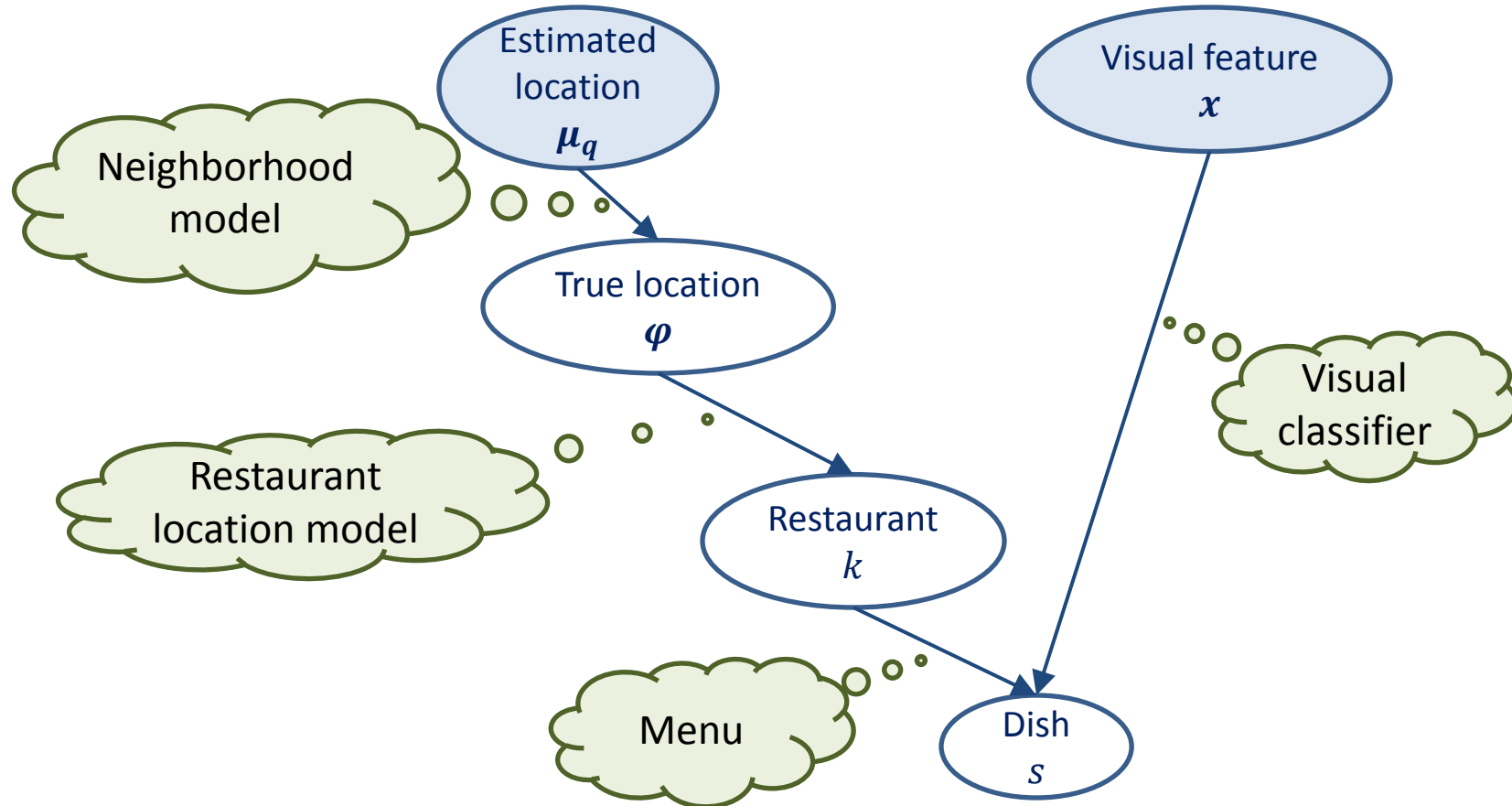
Knowledge about restaurants



Shortlist approach for dish recognition

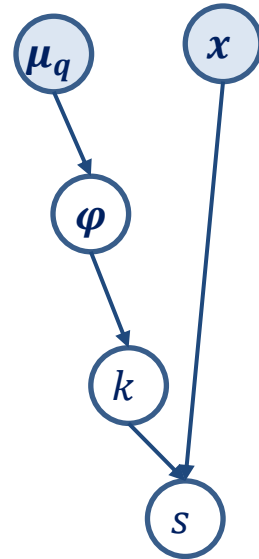


Probabilistic graphical model



Probabilistic graphical model

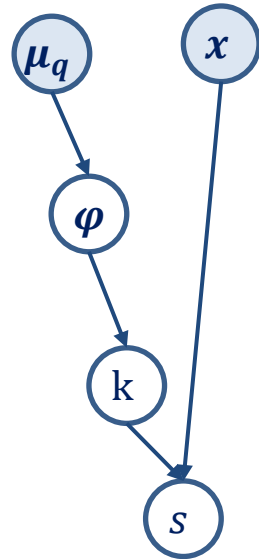
- If we model $p(s, k, \boldsymbol{\varphi} | \boldsymbol{\mu}_q, \boldsymbol{x})$ (joint distribution)
$$p(s, k, \boldsymbol{\varphi} | \boldsymbol{\mu}_q, \boldsymbol{x}) = \underbrace{p(\boldsymbol{\varphi} | \boldsymbol{\mu}_q)}_{\text{Location}} \underbrace{p(k | \boldsymbol{\varphi})}_{\text{Restaurant location}} \underbrace{p(s | k, \boldsymbol{x})}_{\text{Menu+visual}}$$
- By marginalizing we can infer either $s, k, \boldsymbol{\varphi}$
 - Dish $p(s | \boldsymbol{\mu}_q, \boldsymbol{x})$
 - Restaurant $p(k | \boldsymbol{\mu}_q, \boldsymbol{x})$
 - Location $p(\boldsymbol{\varphi} | \boldsymbol{\mu}_q, \boldsymbol{x})$



Predicting dish

- Marginalizing $p(s, k, \boldsymbol{\varphi} | \boldsymbol{\mu}_q, \mathbf{x})$ over k and $\boldsymbol{\varphi}$

$$\begin{aligned} p(s | \boldsymbol{\mu}_q, \mathbf{x}) &= \sum_{k=1}^K \int_{\boldsymbol{\varphi}} p(s, k, \boldsymbol{\varphi} | \boldsymbol{\mu}_q, \mathbf{x}) d\boldsymbol{\varphi} \\ &= \sum_{k=1}^K p(s | k, \mathbf{x}) \int_{\boldsymbol{\varphi}} p(\boldsymbol{\varphi} | \boldsymbol{\mu}_q) p(k | \boldsymbol{\varphi}) d\boldsymbol{\varphi} \end{aligned}$$



- And the predicted dish is

$$s^* = \operatorname{argmax}_s p(s | \boldsymbol{\mu}_q, \mathbf{x})$$

Revisiting shortlist

- Neighborhood model: hard **circular neighborhood**

$$p(\boldsymbol{\varphi}|\boldsymbol{\mu}_q) = [\|\boldsymbol{\varphi} - \boldsymbol{\mu}_q\| \leq \epsilon]$$

- Restaurant location model: **point**

$$p(k|\boldsymbol{\varphi}) = \delta(\|\boldsymbol{\varphi} - \boldsymbol{\mu}_k\|)$$

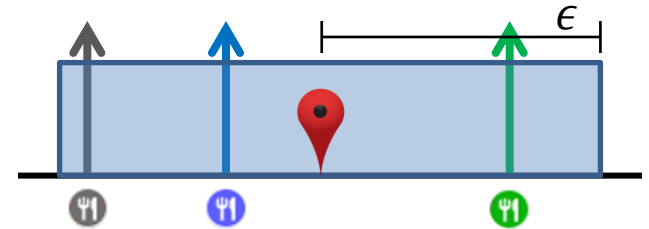
- Visual classifier: **global, filtered by menu**

$$p(s|k, \boldsymbol{x}) \propto p_G(s|\boldsymbol{x})[s \in M_k]$$

- Dish conditional probability

$$p(s|\boldsymbol{\mu}_q, \boldsymbol{x}) \propto p_G(s|\boldsymbol{x}) \left[s \in \bigcup_{k \in H_q} M_k \right]$$

with $H_q = \{k | \boldsymbol{\mu}_k \in B_q\}$ and $B_q = \{\boldsymbol{\varphi} | \|\boldsymbol{\varphi} - \boldsymbol{\mu}_q\| \leq \epsilon\}$



Better models (Gaussian)

- Neighborhood model: **Gaussian**

$$p(\boldsymbol{\varphi}|\boldsymbol{\mu}_q) = \mathcal{N}(\boldsymbol{\varphi}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q)$$

$$\boldsymbol{\Sigma}_q = \begin{pmatrix} \sigma_q & \mathbf{0} \\ \mathbf{0} & \sigma_q \end{pmatrix}$$

- Restaurant location model: **Gaussian**

$$p(k|\boldsymbol{\varphi}) = \mathcal{N}(\boldsymbol{\varphi}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

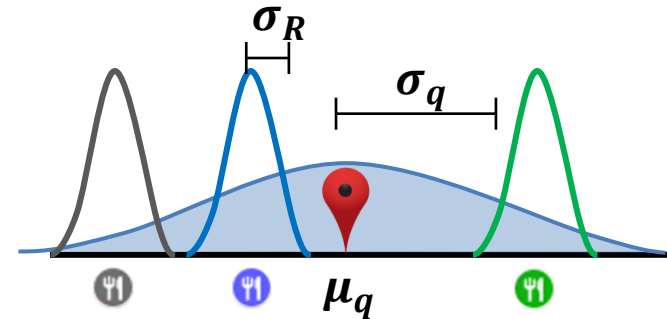
$$\boldsymbol{\Sigma}_k = \boldsymbol{\Sigma}_R = \begin{pmatrix} \sigma_R & \mathbf{0} \\ \mathbf{0} & \sigma_R \end{pmatrix}$$

- Visual classifier: global, filtered by menu

$$p(s|k, \mathbf{x}) \propto p_G(s|\mathbf{x})[s \in M_k]$$

- Dish conditional probability

$$p(s|\boldsymbol{\mu}_q, \mathbf{x}) \propto p_G(s|\mathbf{x}) \sum_{k=1}^K \frac{[s \in M_k]}{|M_k|} \mathcal{N}(\boldsymbol{\mu}_k|\boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q + \boldsymbol{\Sigma}_k)$$



Estimating restaurant and location

- Marginalizing $p(s, k, \boldsymbol{\varphi} | \boldsymbol{\mu}_q, \mathbf{x})$ over s and $\boldsymbol{\varphi}$

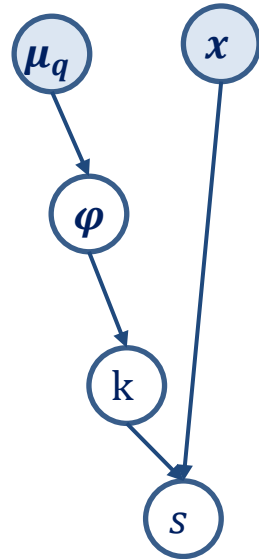
$$p(k | \boldsymbol{\mu}_q, \mathbf{x}) \propto \mathcal{N}(\boldsymbol{\mu}_k | \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q + \boldsymbol{\Sigma}_k) \frac{\sum_{s \in M_k} p(s | \mathbf{x})}{|M_k|}$$

$$k^* = \operatorname{argmax}_k p(k | \boldsymbol{\mu}_q, \mathbf{x})$$

- Marginalizing $p(s, k, \boldsymbol{\varphi} | \boldsymbol{\mu}_q, \mathbf{x})$ over s and k

$$p(\boldsymbol{\varphi} | \boldsymbol{\mu}_q, \mathbf{x})$$

- $\boldsymbol{\varphi}$ is continuous
- Maximum likelihood estimation
- Not closed form. Iterative procedure



Experiments: Dishes dataset

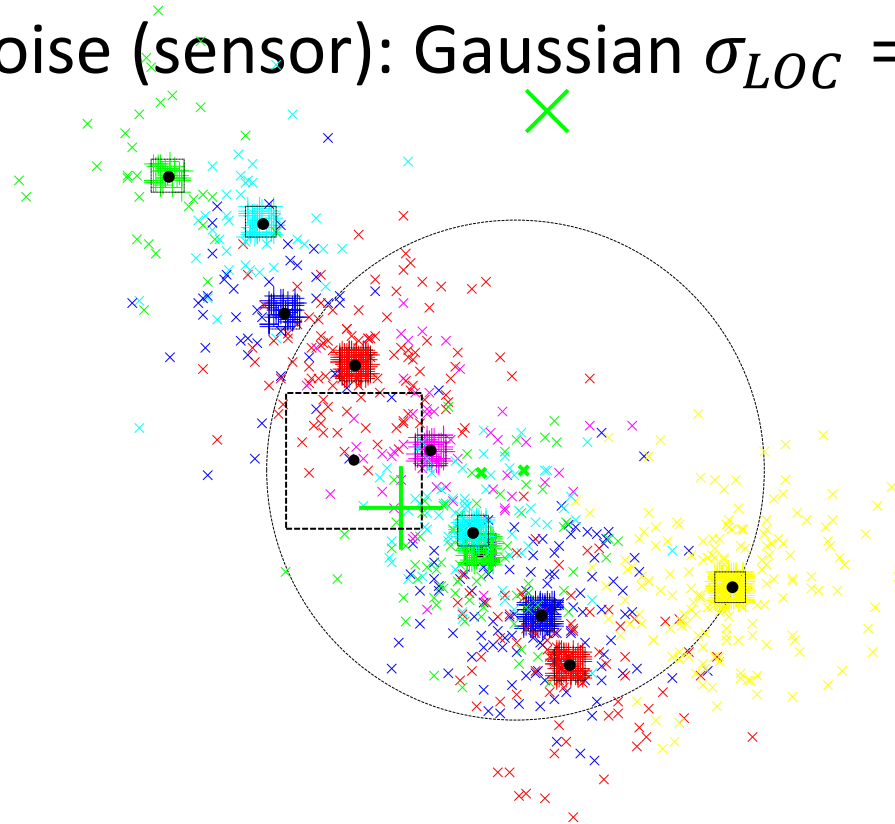
- Restaurants and dishes collected from dianping.com
 - Beijing city
 - Mostly Chinese dishes
- 187 restaurants
 - Menu (≥ 3 dishes/restaurant)
 - Dish images (≥ 15 images/dish)
 - Geo-location
- 701 unique dish categories (1173 in total)

Experiments: settings

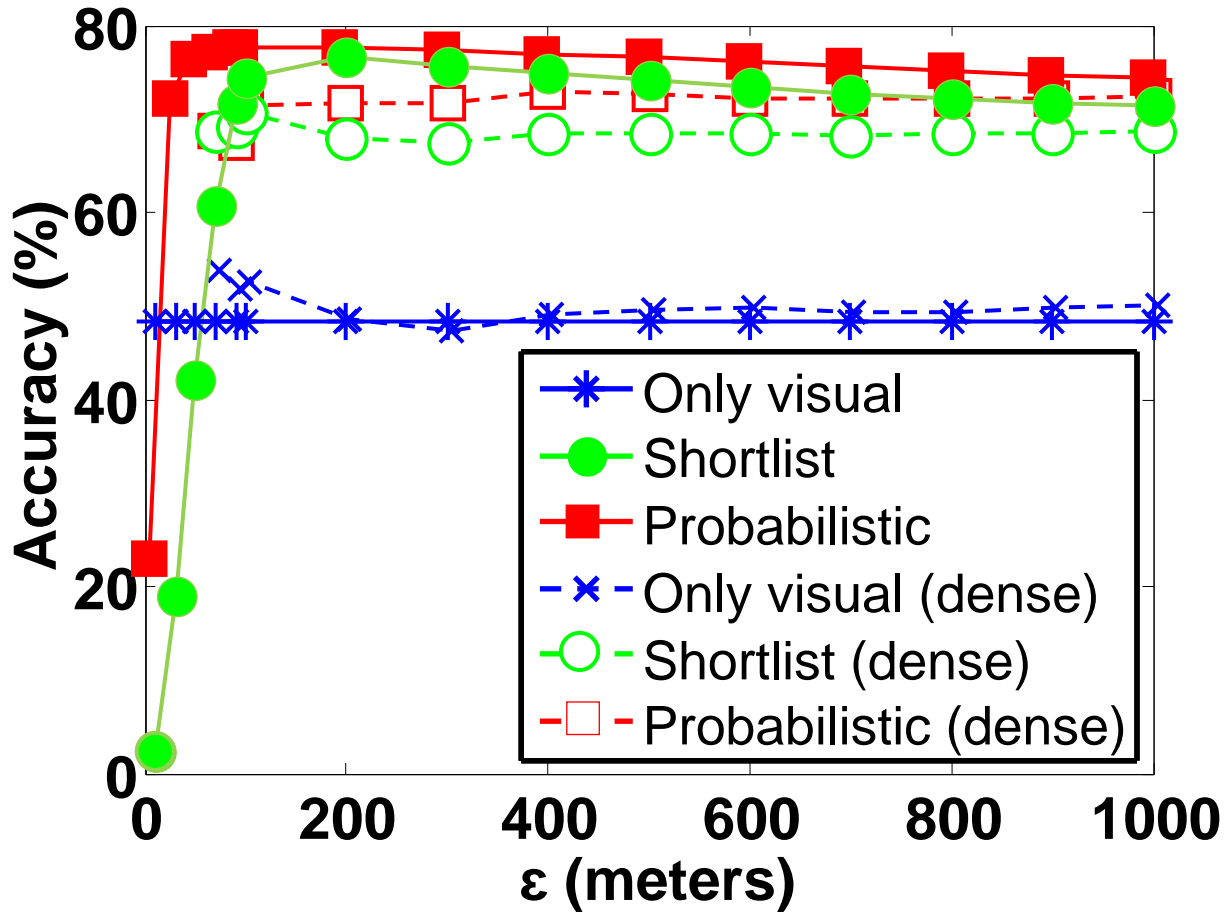
- Visual features: CNN FC7 (DeCAF)
- SVM: #training 10 images/dish (#test \geq 5 img/dish)
- Dense setting: queries with \geq 5 restaurants in neighborhood
- To compare shortlist/probabilistic: align parameters as $\epsilon = 3\sigma_q$

Simulating geo-locations for test

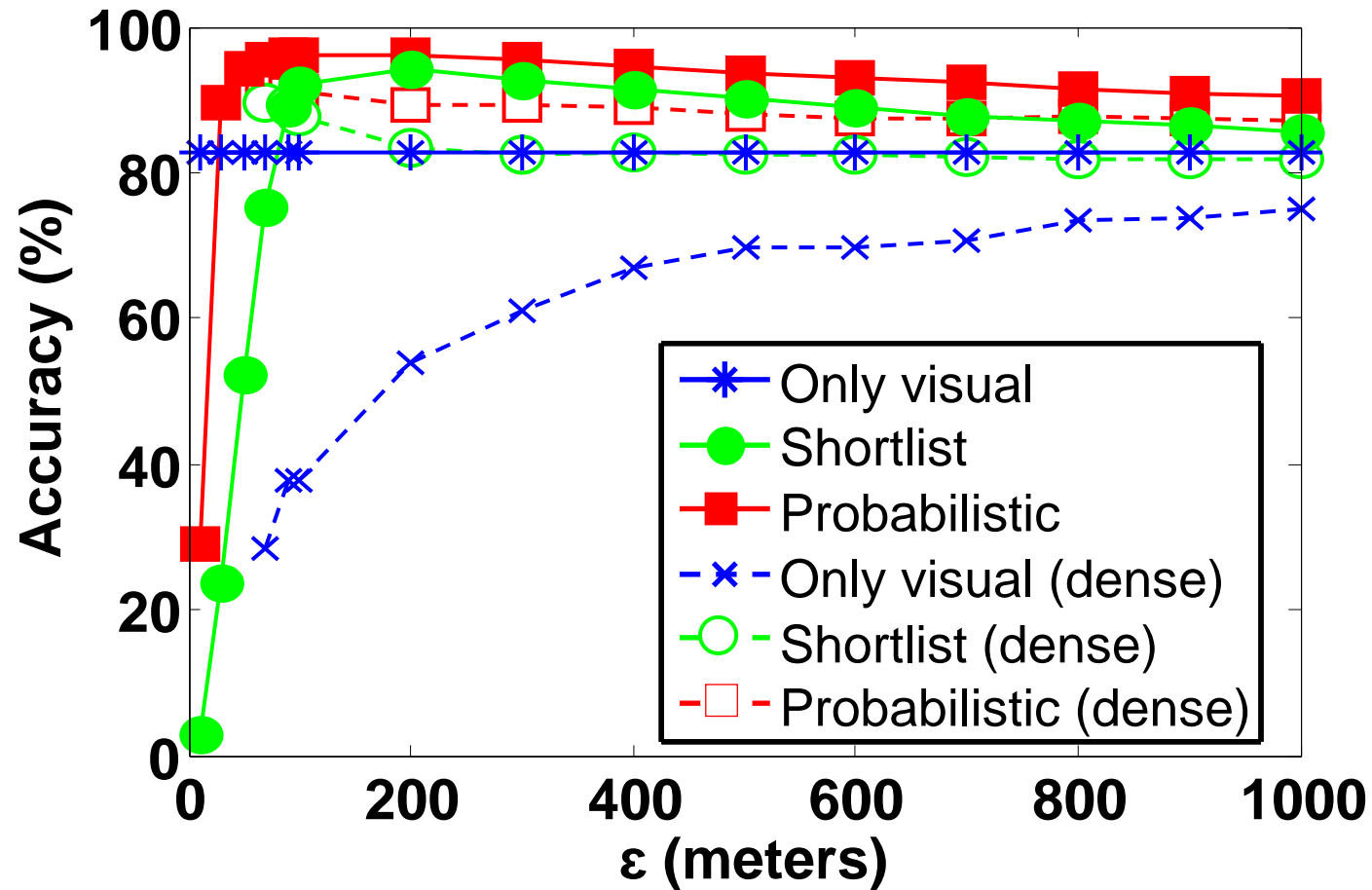
- Restaurants modeled as 25x25 m² squares
- Random location in the restaurant: uniform
- Random noise (sensor): Gaussian $\sigma_{LOC} = 40$ m



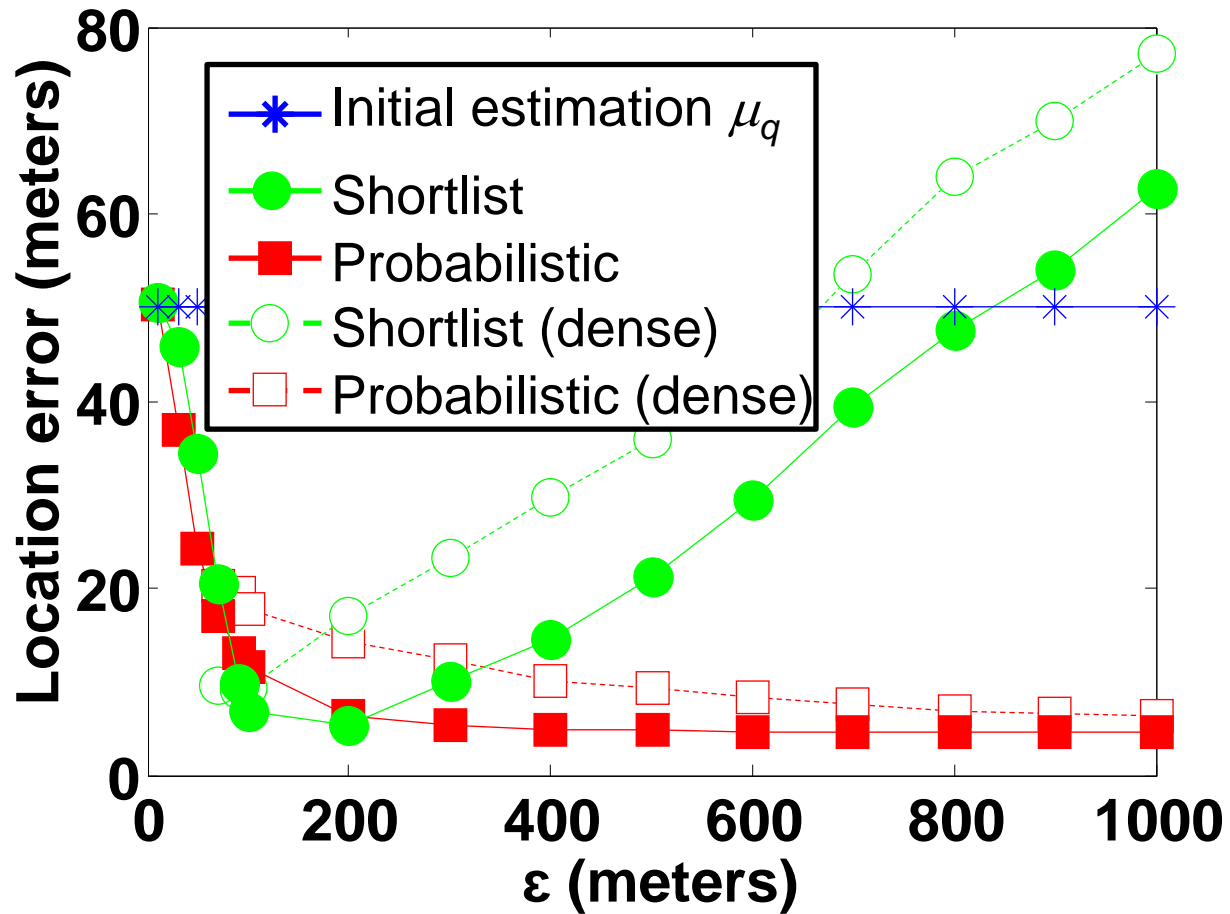
Results: dish recognition



Results: restaurant recognition



Results: geolocation refinement



Conclusions

- Dish recognition in restaurants
 - Food recognition+context+external knowledge
- Better using a probabilistic framework
 - Improves performance
 - More robust to the choice of parameters
- Restaurant recognition and geolocation refinement as byproducts

THANK YOU!

