

Similarity

What is a cube?

What is a cylinder?

How can you measure the similarity?



Computational Linguistics

Words and Images

Stefan Thater & Dietrich Klakow

FR 4.7 Allgemeine Linguistik (Computerlinguistik)

Universität des Saarlandes

Summer 2013

Example

Look at two examples

- Retrieving images using text queries
- Generating captions for images

Image Retrieval using Query Words



Tiger



Sign in

Search

About 761,000,000 results (0.17 seconds)

SafeSearch moderate



Web

Related searches: [ek tha tiger](#) [white tigers](#) [wild tigers](#) [indian tigers](#) [cute tigers](#) [tiger cubs](#)

Images

Maps

Videos

News

Shopping

More

Any time

Past 24 hours

Past week

Custom range...

All results

By subject

Any size

Large

Medium

Icon

Larger than...

Exactly...

Any color

Full color

Black and white





Tiger



Sign in

Search

About 761,000,000 results (0.17 seconds)

SafeSearch moderate



Web

Related searches: [ek tha tiger](#) [white tigers](#) [wild tigers](#) [indian tigers](#) [cute tigers](#) [tiger cubs](#)

Images

Maps

Videos

News

Shopping

More



220px-Tiger in Rantha
en.wikipedia.org
220 × 148 - Tiger
Similar More sizes



Any time

- Past 24 hours
- Past week
- Custom range...

All results

By subject

Any size

- Large
- Medium
- Icon
- Larger than...
- Exactly...

Any color

- Full color
- Black and white





Tiger



Sign in

Search

About 761,000,000 results (0.17 seconds)

SafeSearch moderate



Web

Related searches: [ek tha tiger](#) [white tigers](#) [wild tigers](#) [indian tigers](#) [cute tigers](#) [tiger cubs](#)

Images

Maps

Videos

News

Shopping

More

Any time

Past 24 hours

Past week

Custom range...

All results

By subject

Any size

Large

Medium

Icon

Larger than...

Exactly...

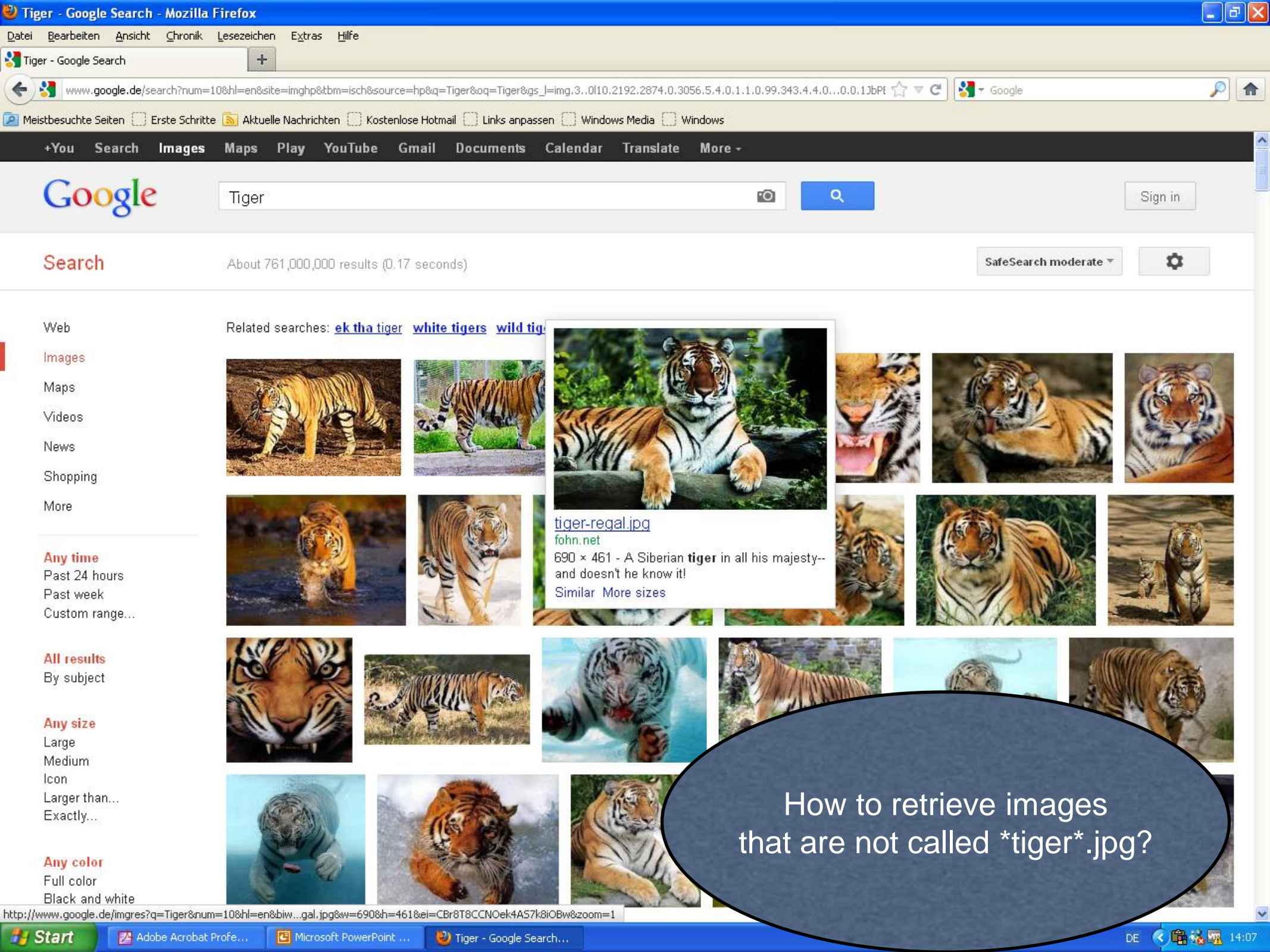
Any color

Full color

Black and white



220px-MarwellAmurTiger
de.wikipedia.org
220 × 147 - Tiger - Wikipedia
Similar More sizes



Tiger



Sign in

Search

About 761,000,000 results (0.17 seconds)

SafeSearch moderate



- Web
- Images
- Maps
- Videos
- News
- Shopping
- More

Related searches: [ek tha tiger](#) [white tigers](#) [wild tiger](#)



[tiger-regal.jpg](#)
fohn.net
690 x 461 - A Siberian **tiger** in all his majesty- and doesn't he know it!
Similar More sizes

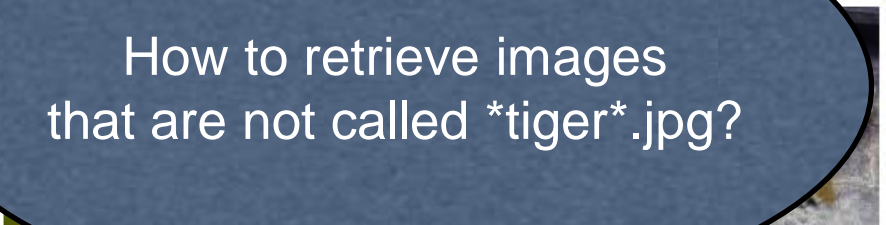


- Any time**
- Past 24 hours
 - Past week
 - Custom range...

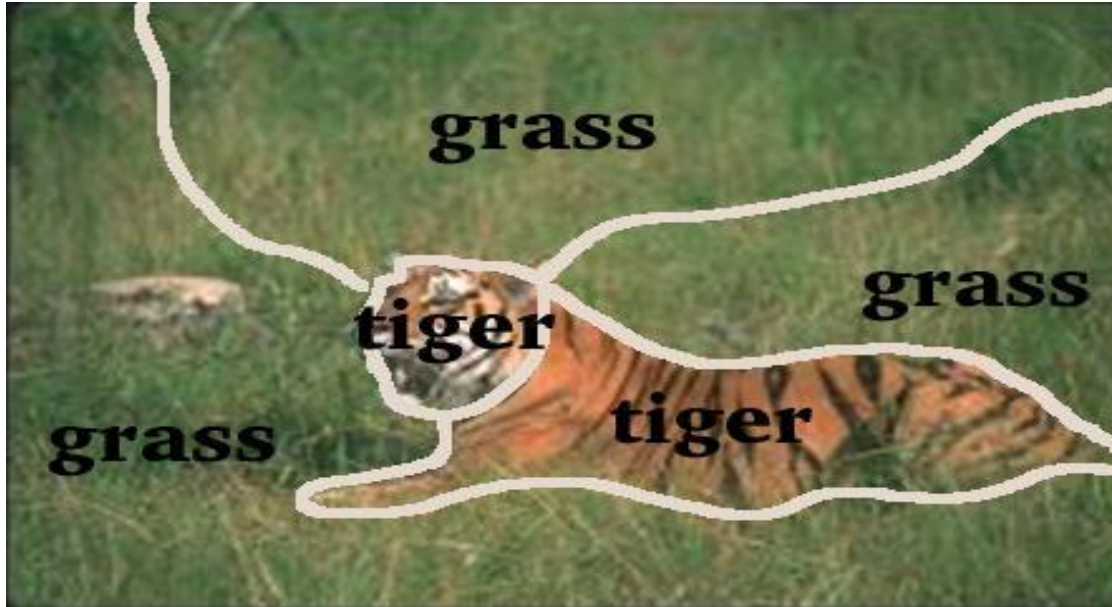
- All results**
- By subject

- Any size**
- Large
 - Medium
 - Icon
 - Larger than...
 - Exactly...

- Any color**
- Full color
 - Black and white



Are Tigers striped?



- Image from [Barnard2003]
- 5000 images plus captions

How important is it to detect “tiger color” and “tiger texture” together?

How important is it to detect “tiger color” and in the middle of the image?

Correlation concept/region is stronger than color/texture!

Overview

- Measuring regional mutual information
- Retrieval and annotation model:
 - Regional language models
- Experiments

The TRECVID 2003 Data

Training : 9413 shots

Validation: 4085 shots

Test : 4787 shots

Num concepts = 75

On average 5 concepts per shot

Features extracted from grids

$5 \times 7 = 35$ grids per image

Horse, outdoors, snow



- Feature extraction:

- Use existing tools (e.g. SIFT features using <http://www.vlfeat.org/~vedaldi/code/sift.html>)

- Features are quantized using k-Means (visterms)

Regional Mutual Information

$$MI(r, c) = \sum_v N(v, r, c) \log \left(\frac{N(v, r, c)}{N(v, r)N(c)} \right)$$

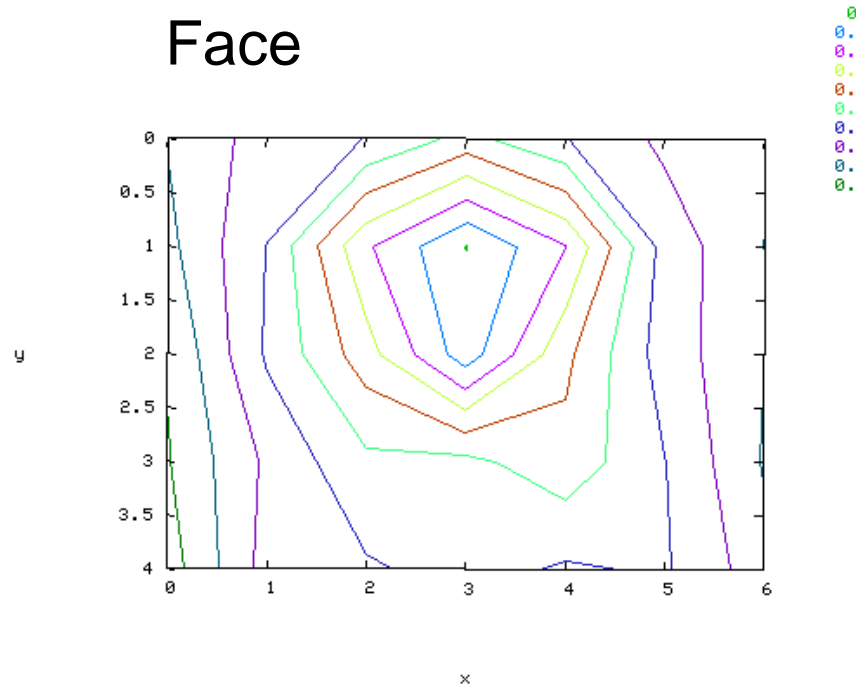
r : label of region (5x7 grid on the image)

c : concept label (e.g. face, vehicle, text_overlay)

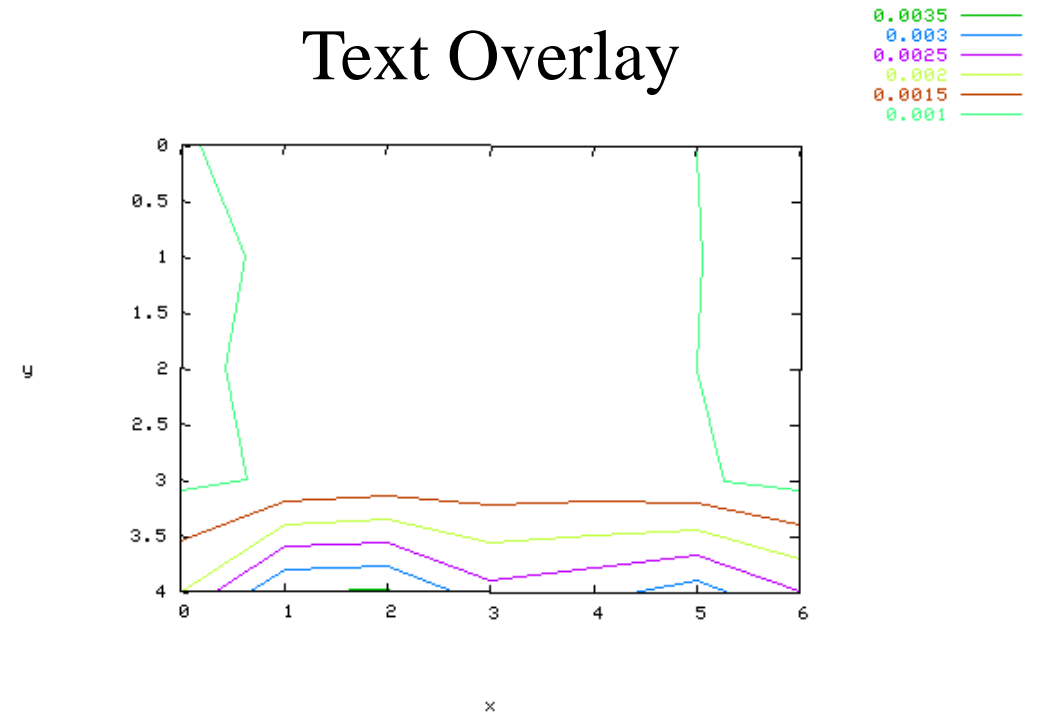
v : visterm (discrete representation of color, texture and edge information)

Regional Mutual Information for two Concepts

Face



Text Overlay



Same visterm in different regions can trigger different concepts

Regional Bayes Classifier

Bayes classifier $\bar{c} = \arg \max_c P(v_1 \dots v_R | c) P(c)$

Independence assumption (Naïve Bayes)

$$P(v_1 \dots v_R | c) = \prod_{r=1}^R P(v_r | c)$$

One model per region

$$P(v_1 \dots v_R | c) = \prod_{r=1}^R P_r(v_r | c)$$

Absolute Discounting

$$P_r(v | c) = \max\left(\frac{N(v, r, c) - d}{N(r, c)}, 0\right) + \frac{d B}{N(r, c)} P_{BG}(v | c)$$

Most widely used in speech recognition

d : discounting parameter

$P_{BG}(v|c)$: backing-off distribution

B : number of v for which $N(v,r,c) > d$

Smoothing Methods

Linear Interpolation

$$P_r(v | c) = (1 - \lambda) \frac{N(v, r, c)}{N(r, c)} + \lambda P_{BG}(v | c)$$

λ : interpolation weight

Also known as Jelinek-Mercer-Smoothing

Dirichlet Prior

$$P_r(v | c) = \frac{N(v, r, c) + \mu P_{BG}(v | c)}{N(r, c) + \mu}$$

μ : weight for prior

Background Distributions

Uniform Distribution (“zerogram”)

$$P_{BG}^{Zero}(v | c) = \frac{1}{|V|}$$

|V|: number of visterms)

Frequency of visterms (“unigram”)

$$P_{BG}^{Unigram}(v | c) = \frac{N(v)}{\sum_v N(v)}$$

Results

(mean average precision)

	Baseline	
Model	BG Zero	BG Unigram
Absolute Discounting	0.149	0.149
Dirichlet Prior	0.149	0.150
Linear Interpolation	0.145	0.148

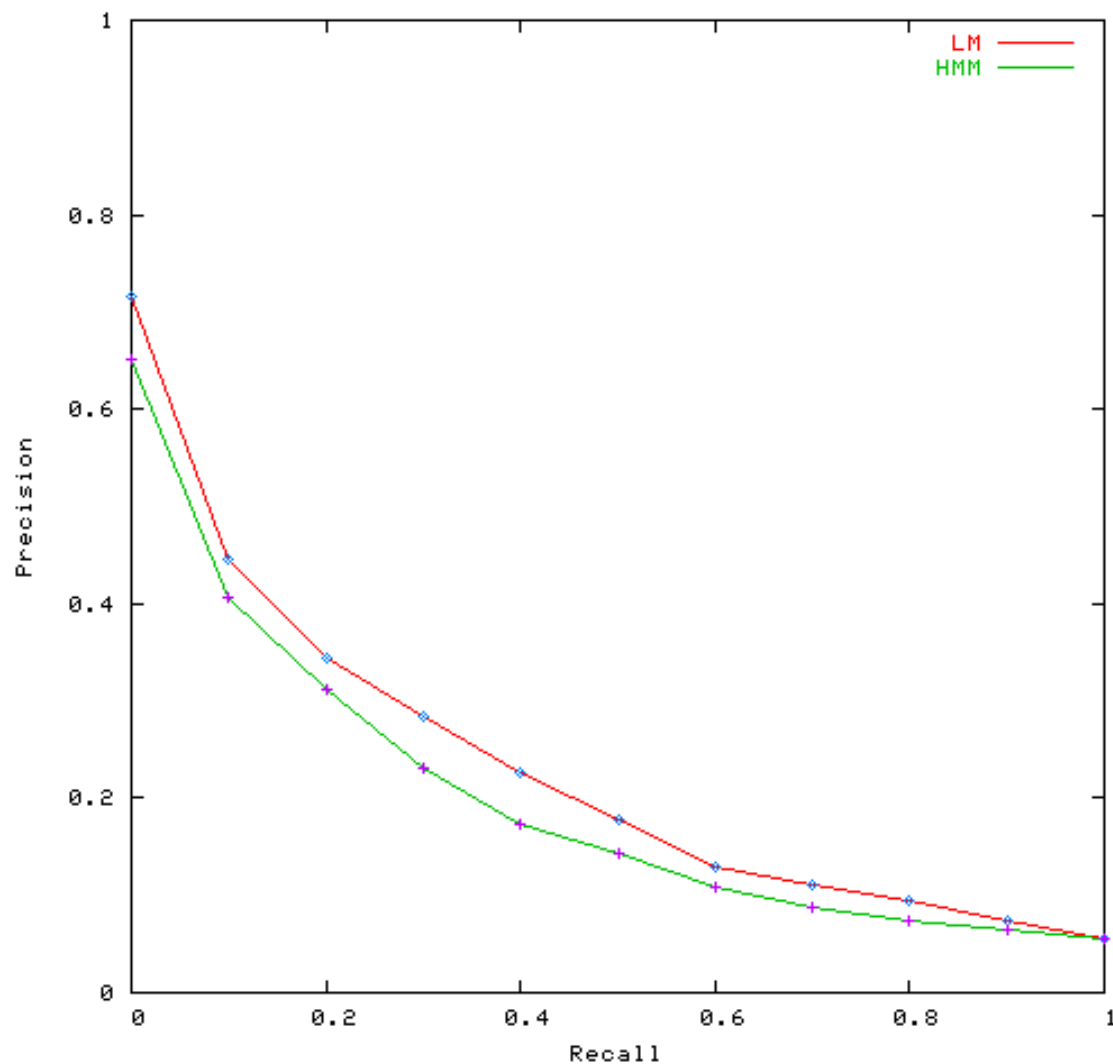
Results

(mean average precision)

	Baseline		+ regional modeling	
Model	BG Zero	BG Unigram	BG Zero	BG Unigram
Absolute Discounting	0.149	0.149	0.209	0.215
Dirichlet Prior	0.149	0.150	0.207	0.218
Linear Interpolation	0.145	0.148	0.215	0.221

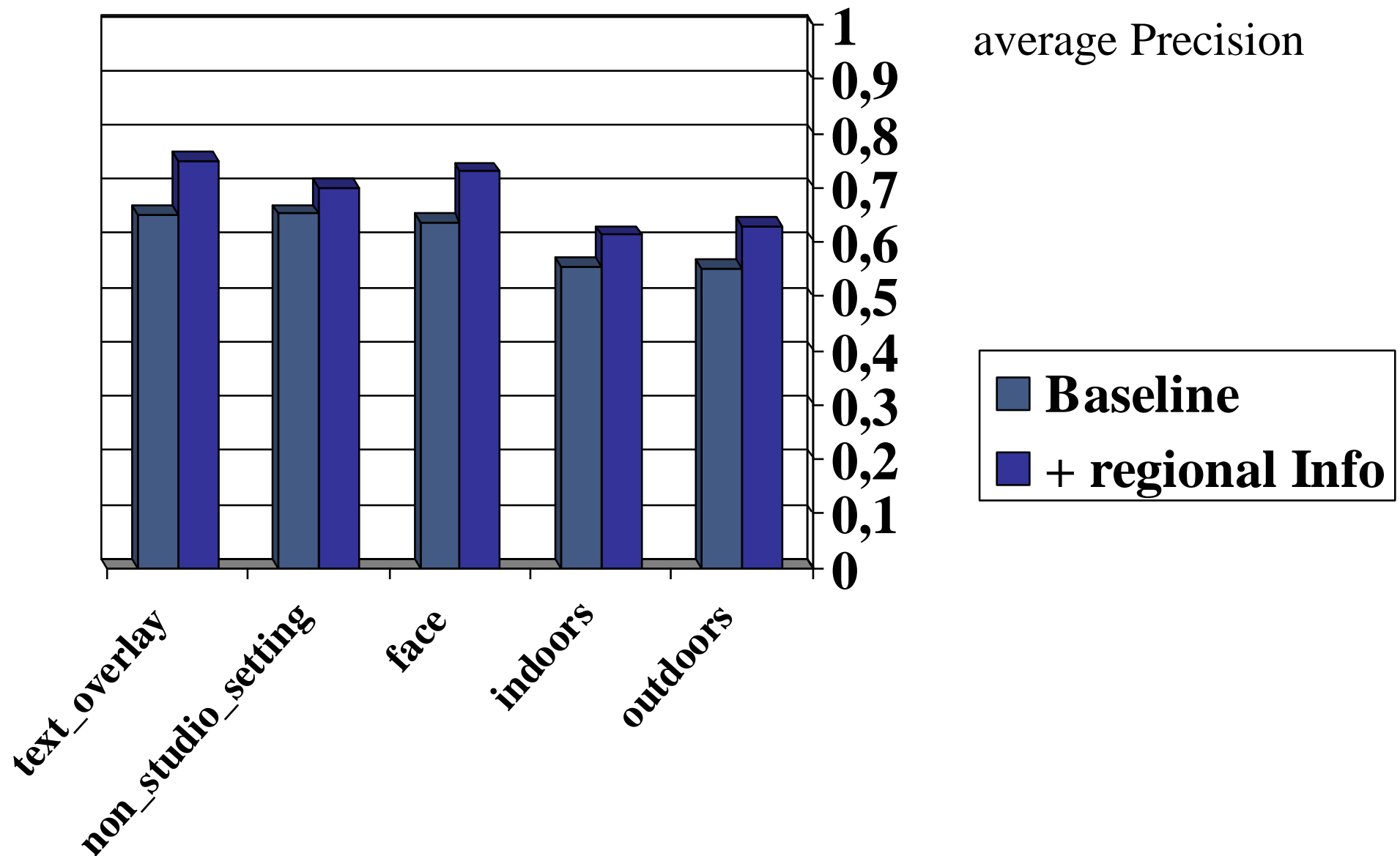
⇒ regional information gives significant improvement (32%)

Precision-Recall Graph



On TREC-VID 03
Consistent
improvement over
HMM model by
Ghoshal et al.
(SigIR 2005)

Improvement on most frequent concepts



Summary

- Analysis by regional mutual information
- Region specific modeling is important
- Best variant uses
 - Linear interpolation
 - Backing off to background unigram
- Extremely simple method

Generating Captions

**How Many Words Is a Picture Worth?
Automatic Caption Generation for News Images**

Yansong Feng and Mirella Lapata
School of Informatics, University of Edinburgh
10 Crichton Street, Edinburgh EH8 9AB, UK
Y.Feng-4@sms.ed.ac.uk, mlap@inf.ed.ac.uk

Task

Thousands of Tongans have attended the funeral of King Taufa'ahau Tupou IV, who died last week at the age of 88. Representatives from 30 foreign countries watched as the king's coffin was carried by 1,000 men to the official royal burial ground.



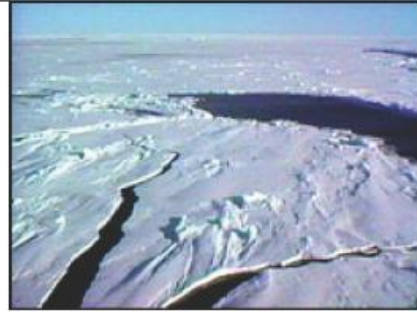
King Tupou, who was 88, died a week ago.

Contaminated Cadbury's chocolate was the most likely cause of an outbreak of salmonella poisoning, the Health Protection Agency has said. About 36 out of a total of 56 cases of the illness reported between March and July could be linked to the product.



Cadbury will increase its contamination testing levels.

A Nasa satellite has documented startling changes in Arctic sea ice cover between 2004 and 2005. The extent of "perennial" ice declined by 14%, losing an area the size of Pakistan or Turkey. The last few decades have seen ice cover shrink by about 0.7% per year.



Satellite instruments can distinguish "old" Arctic ice from "new".

A third of children in the UK use blogs and social network websites but two thirds of parents do not even know what they are, a survey suggests. The children's charity NCH said there was "an alarming gap" in technological knowledge between generations.



Children were found to be far more internet-wise than parents.

First Step

Annotate each image with key words

(e.g. using method from the previous section)

Second step: generate caption

Extractive approaches

Idea: pick a suitable sentence from the text to be the caption

Approach:

Measure similarity between sentence and key words describing image using

Word overlap

Cosine

KL-divergence

Second step: generate caption

abstractive approaches

Idea: pick the most likely word sequence given the key words

Estimate N-gram based on

Surrounding text

Key words

Combine the two using LM adaptation techniques

Optional: include a phrase based constraint

Results

<p>G: King Tupou, who was 88, died a week ago.</p> <p>KL: Last year, thousands of Tongans took part in unprecedented demonstrations to demand greater democracy and public ownership of key national assets.</p> <p>A_W: King Toupou IV died at the age of Tongans last week.</p> <p>A_P: King Toupou IV died at the age of 88 last week.</p>
<p>G: Cadbury will increase its contamination testing levels.</p> <p>KL: Contaminated Cadbury's chocolate was the most likely cause of an outbreak of salmonella poisoning, the Health Protection Agency has said.</p> <p>A_W: Purely dairy milk buttons Easter had agreed to work has caused.</p> <p>A_P: The 105g dairy milk buttons Easter egg affected by the recall.</p>
<p>G: Satellite instruments can distinguish "old" Arctic ice from "new".</p> <p>KL: So a planet with less ice warms faster, potentially turning the projected impacts of global warming into reality sooner than anticipated.</p> <p>A_W: Dr less winds through ice cover all over long time when.</p> <p>A_P: The area of the Arctic covered in Arctic sea ice cover.</p>
<p>G: Children were found to be far more internet-wise than parents.</p> <p>KL: That's where parents come in.</p> <p>A_W: The survey found a third of children are about mobile phones.</p> <p>A_P: The survey found a third of children in the driving seat.</p>

Summary

Image retrieval using text queries
Caption generation