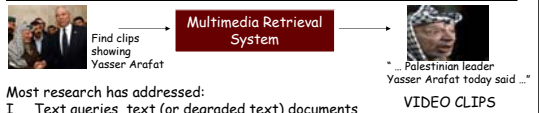


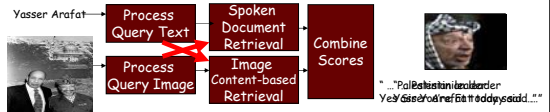
# Joint Visual-Text Modeling for Multimedia Retrieval

Résumé du JHU Workshop 2004  
Pour le Master USTV

## Big Picture: Multimedia Retrieval Task

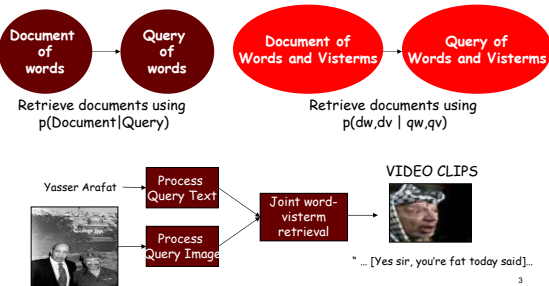


Most research has addressed:  
I. Text queries, text (or degraded text) documents  
II. Image queries, image data

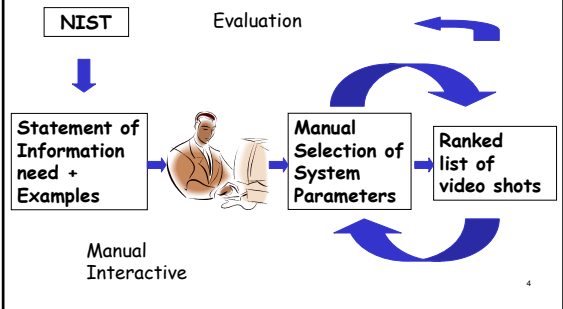


### Joint-Visual Text Models!

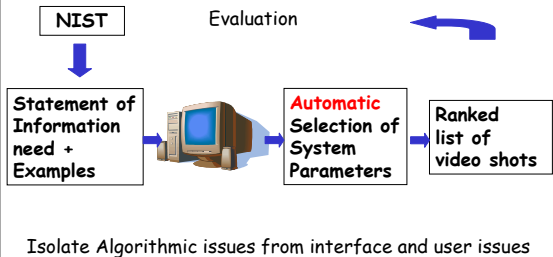
## Joint Visual-Text Modeling



## TRECVID Search task definition

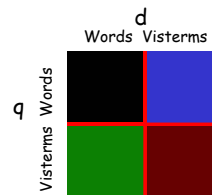


## Our search task definition



## Language Model based Retrieval

Rank documents with  $p(qw,qv|dw,dv)$



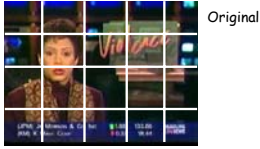
Baseline model

Relating document visterns to query words (MT, Relevance Model, HMMs)

Relating document words to query images (Text Classification experiments)

Visual-only retrieval models

## Experimental Setup: Visual Features



7

## Interest Point Neighborhoods (Harris detector)



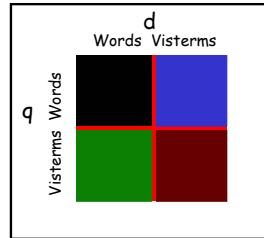
8

## Experimental Setup: Visual Feature list

- Regular partition
  - L\*a\*b Moments (COLOR)
  - Smoothed Edge Orientation Histogram (EDGE)
  - Grey-level Co-occurrence matrix (TEXTURE)
- Interest Point neighborhood
  - COLOR, EDGE, TEXTURE

9

## Presentation Outline



Translation (MT) models (Paola),

Relevance Models (Shao Lei, Desislava),

Graphical Models (Pavel, Brock)

Text classification models (Matt)

Integration & Summary (Dietrich)

10

## Inspiration from Machine Translation



the beautiful sun  
le soleil beau



the beautiful sun  
le soleil beau



$$p(f|e) = \sum_a p(f, a|e)$$

$$p(c|v) = \sum_a p(c, a|v)$$

Direct translation model 12

## Discrete Representation of Image Regions (visterms) to create analogy to MT

In Machine Translation → discrete tokens

In our task



sun sky waves sea  
concepts ✓



V10 V22 V35 V43

C5 C1 C38 C71



V20 V21 V50 V10

C15 C21 C83



V78 V28 V1 V1

C21 C19 C1 C56 C38



Solution : Vector quantization → visterms ✓

→ {f<sub>1</sub>, f<sub>2</sub>, ..., f<sub>nm</sub>} → v<sub>k</sub>



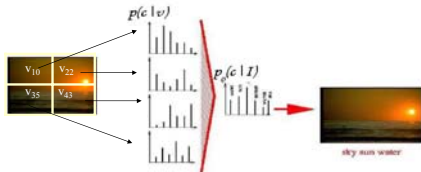
water harbor sky clouds sea

13

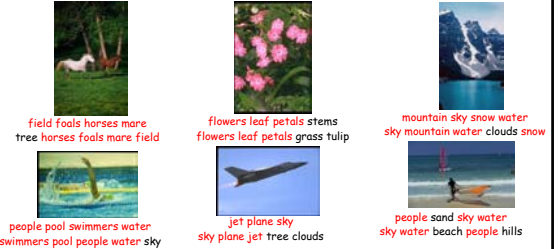
## Image annotation using translation probabilities

$p(c|v)$  : Probabilities obtained from direct translation  $p(\text{sun} | \text{sun})$

$$P_0(c | d_v) = \frac{1}{|d_v|} \sum_{v \in d_v} P(c | v)$$



## Annotation Results (Corel set)



Top: manual annotations, bottom: predicted words (top 5 words with the highest probability)  
Red: correct matches

15

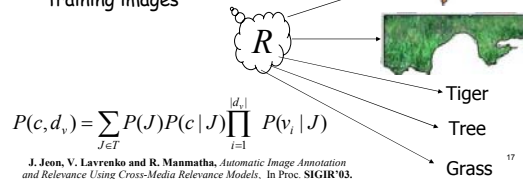
## Intuition

- Images are defined by spatial context.
  - Isolated pixels have no meaning.
  - Context simplifies recognition/retrieval.
  - E.g. Tiger is associated with grass, tree, water forest.
    - Less likely to be associated with computers.



## Cross Media Relevance Models (CMRM)

- Two parallel vocabularies: Words and Visterms
- Analogous to Cross-lingual relevance models
- Estimate the joint probabilities of words and visterms from training images



17

## Annotation Examples (Corel set)



Sky train railroad  
locomotive water



Cat tiger bengal  
tree forest



Snow fox arctic  
tails water



Tree plane zebra  
herd water



Birds leaf nest water  
sky



Mountain plane  
jet water sky

18

## Proposal: Using Dynamic Information for Video Retrieval

19

## Motivation

- Current models based on single frames in each shot.
- But video is dynamic
  - Has motion information.
- Use dynamic (motion) information
  - Better image representations (segmentations)
  - Model events/actions

20

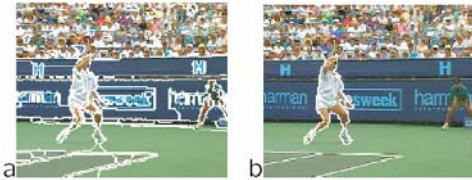
## Why Dynamic Information

- Model actions/events
  - Many Trecvid 2003 queries require motion information. E.g.
    - find shots of an airplane *taking off*.
    - find shots of a person *diving into* water.
  - Motion is an important cue for retrieving actions/events.
    - But using the optical flow over the entire image doesn't help.
    - Use motion features from objects.
- Better Image Representations
  - Much easier to segment moving objects from background than to segment static images.



21

## Segmentation Comparison



a: Segmentation using only still image information

b: Segmentation using only motion information

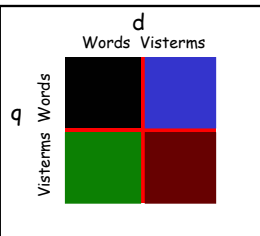
Pictures from Patrick Bouthemy's Website, INRIA

22

## Hidden Markov Models for Image Annotations

23

## Presentation Outline



Translation (MT) models (Paola).

Relevance Models (Shao Lei, Desislava).

Graphical Models (Pavel, Brock)

Text classification models (Matt)

Integration & Summary (Dietrich)

24

## Inadequacy of the annotations

- Corel database
  - Annotators often mark only interesting objects
- TRECVID database
  - Annotation concepts capture mostly semantics of the image and they are not very suitable for describing visual properties



beach  
palm  
people  
tree



man-made object



car  
transportation  
vehicle  
outdoors  
non-studio setting  
nature-non-vegetation  
snow

25

## Alignment problems

- There is **no** notion of **order** in the annotation words
  - Difficulties with automatic alignment between words and image regions



26

## Gradual Training Results

Forced alignment – flat-start training



jet, plane, sky

Forced alignment – gradual training



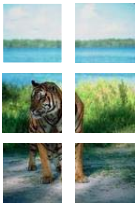
jet, plane, sky

### Results:

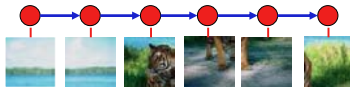
- Improved alignment of training images
- Annotation performance on test images did not change significantly

27

## Joint Segmentation and Labeling

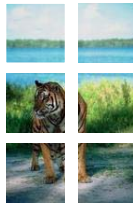


tiger, grass, sky

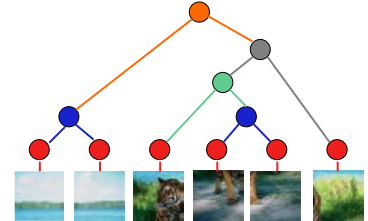


28

## Joint Segmentation and Labeling



tiger, grass, sky

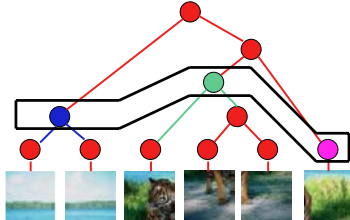


29

## Joint Segmentation and Labeling

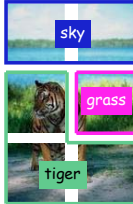


tiger, grass, sky

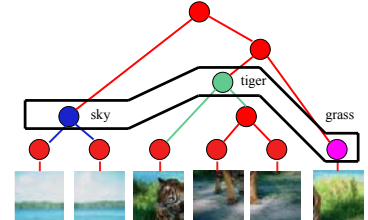


30

## Joint Segmentation and Labeling



tiger, grass, sky

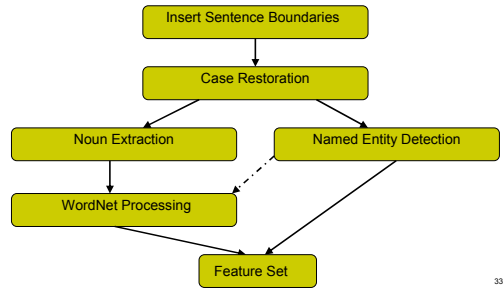


31

# Predicting Visual Concepts From Text

32

## Building Features



33

## ASR → Features Example

STEVE FOSSETT AND HIS BALLOON SOLO SPIRIT ARSENIDE OVER THE BLACK SEA DRIFTING SLOWLY TOWARDS THE COAST OF THE CAUCUSES HIS TEAM PLANS IF NECESSARY TO BRING HIM DOWN AFTER DAYLIGHT TOMORROW YOU THE CHECHEN CAPITAL OF GROZNY

34

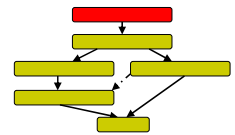
## ASR → Features Example

STEVE FOSSETT AND HIS BALLOON SOLO SPIRIT ARSENIDE.

OVER THE BLACK SEA DRIFTING SLOWLY TOWARDS THE COAST OF THE CAUCUSES.

HIS TEAM PLANS IF NECESSARY TO BRING HIM DOWN AFTER DAYLIGHT TOMORROW.

YOU THE CHECHEN CAPITAL OF GROZNY



35

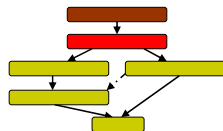
## ASR → Features Example

Steve Fossett and his balloon Solo Spirit arsenide.

Over the Black Sea drifting slowly towards the coast of the caucuses.

His team plans if necessary to bring him down after daylight tomorrow.

you the Chechan capital of Grozny....



36

## ASR → Features Example

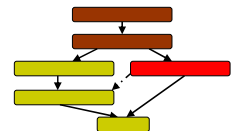
Steve Fossett and his balloon Solo Spirit arsenide.

- Named Entities
  - Male Person, Location (Region)

Over the **Black Sea** drifting slowly towards the coast of the caucuses.

His team plans if necessary to bring him down after daylight tomorrow.

you the Chechan capital of Grozny.



37

## ASR → Features Example

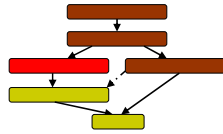
Steve Fossett and his balloon  
Solo Spirit arsenide.

Over the Black Sea drifting  
slowly towards the coast of the  
caucuses.

His team plans if necessary to  
bring him down after daylight  
tomorrow.

you the Chechan capital of  
Grozny.

- Named Entities
  - Male Person, Location (Region)



38

## ASR → Features Example

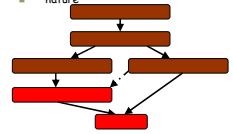
Steve Fossett and his balloon  
Solo Spirit arsenide.

Over the Black Sea drifting  
slowly towards the coast of the  
caucuses.

His team plans if necessary to  
bring him down after daylight  
tomorrow.

you the Chechan capital of  
Grozny.

- Named Entities
  - Male Person, Location (Region)
- Nouns
  - balloon, solo, spirit, coast, caucas, team, daylight, Chechan, capital, Grozny
- WordNet
  - nature



39

## Will this help for retrieval?

- "Find shots of a person diving into some water."
  - person, water\_body, non-studio\_setting, nature\_non-vegetation, person\_action, indoors
- "Find shots of the front of the White House in the daytime with the fountain running."
  - building, outdoors, sky, water\_body, cityscape, house, nature\_vegetation
- "Find shots of Congressman Mark Souder."
  - person, face, indoors, briefing\_room\_setting, text\_overlay

40

## Joint Visual-Text Video OCR

41

## Motivation

- "Find shots of Congressman Mark Souder"



42

## Motivation

- "Find shots of a graphic of Dow Jones Industrial Average showing a rise for one day. The number of points risen that day must be visible."



43

## Motivation

- Find shots of the Tomb of the Unknown Soldier in Arlington National Cemetery.



## Motivation



WEIF11 I1 NFWdJ TNNIF H

## Why use video OCR?



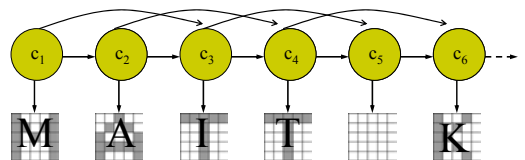
## Why use video OCR?

- Text overlays provide high precision information about query-relevant concepts in the *current* image.

## Image Processing

- Preprocessing
  - Normalize the text region's height
- Feature extraction
  - Color
  - Edge Strength and Orientation

## Proposal: HMM-based recognizer





## CONCLUSION : Resumé des requêtes multi-modales

		Document	
		Words $d_w$	Visterms $d_v$
Query	Words $q_w$	$p(q_w   d_w)$	$p(q_w   d_v)$
	Visterms $q_v$	$p(q_v   d_w)$	$p(q_v   d_v)$

50

## Méthodes

		Document	
		Words $d_w$	Visterms $d_v$
Query	Words $q_w$	$p(q_w   d_w)$	$p(q_w   d_v)$ •MT •Relevance Models •HMM
	Visterms $q_v$	$p(q_v   d_w)$ •Naive Bayes •Max. Ent •LM •SVM, Ada Boost, ...	$p(q_v   d_v)$

51