

Dynamic Probabilistic Multimedia Retrieval Model

Tzvetanka I. Ianeva

Introduction

Analysis of a video as a single entity by modeling one-second video sequence around the keyframe

Gaussian Mixture Model in DCT-space-time domain

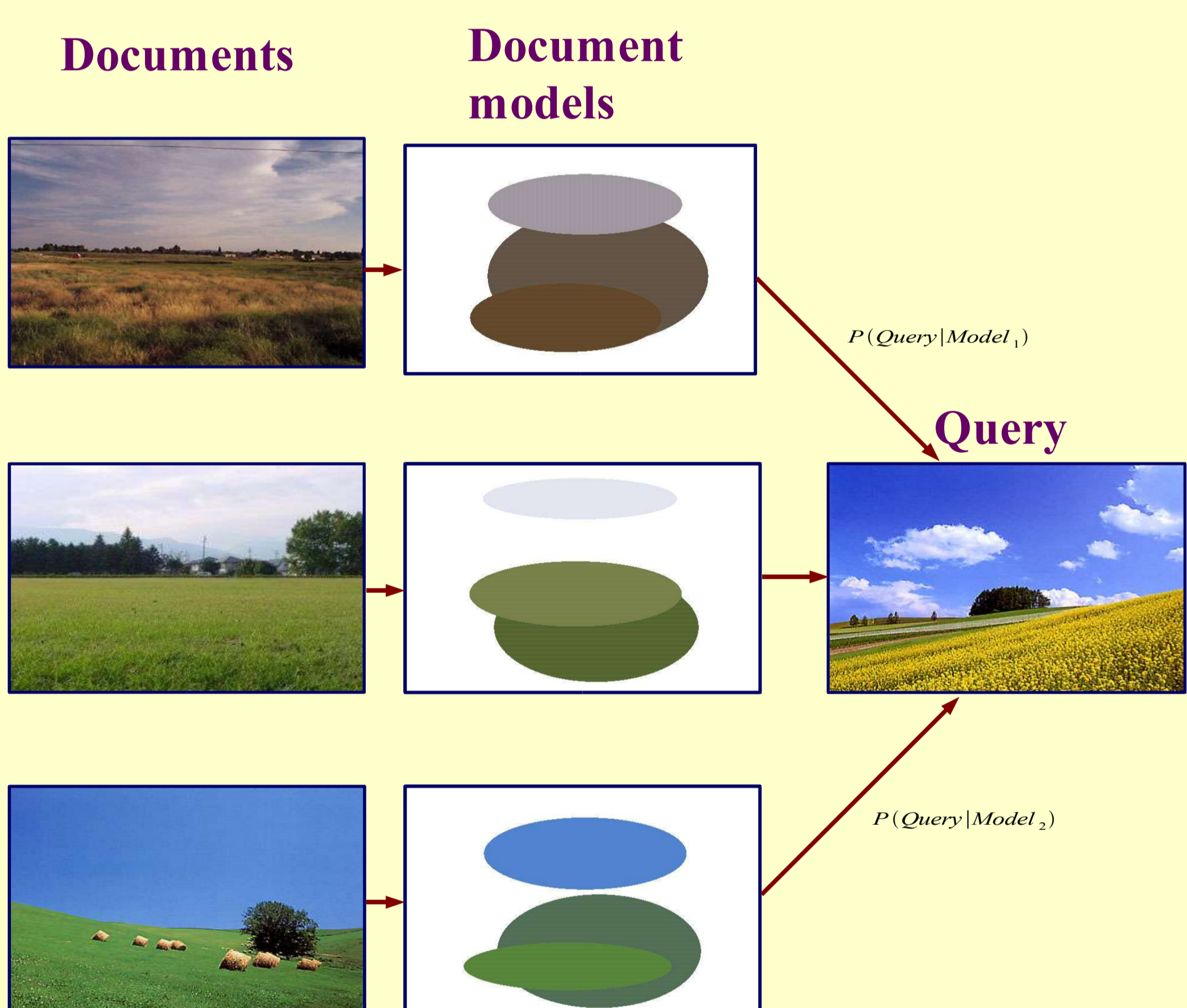
Extension of the Static Probabilistic Retrieval model
 Thijs Westerveld, Arjen P. de Vries: Experimental result analysis for a generative probabilistic image retrieval model. SIGIR2003: 135-142

Pursue *video* retrieval instead of keyframe retrieval

Static Retrieval Model

Given query image x consisting of N samples ($x = x_1, x_2, \dots, x_N$), collection image (document) w_i is compared to x by computing its ability (*RSV*) to explain the samples x :

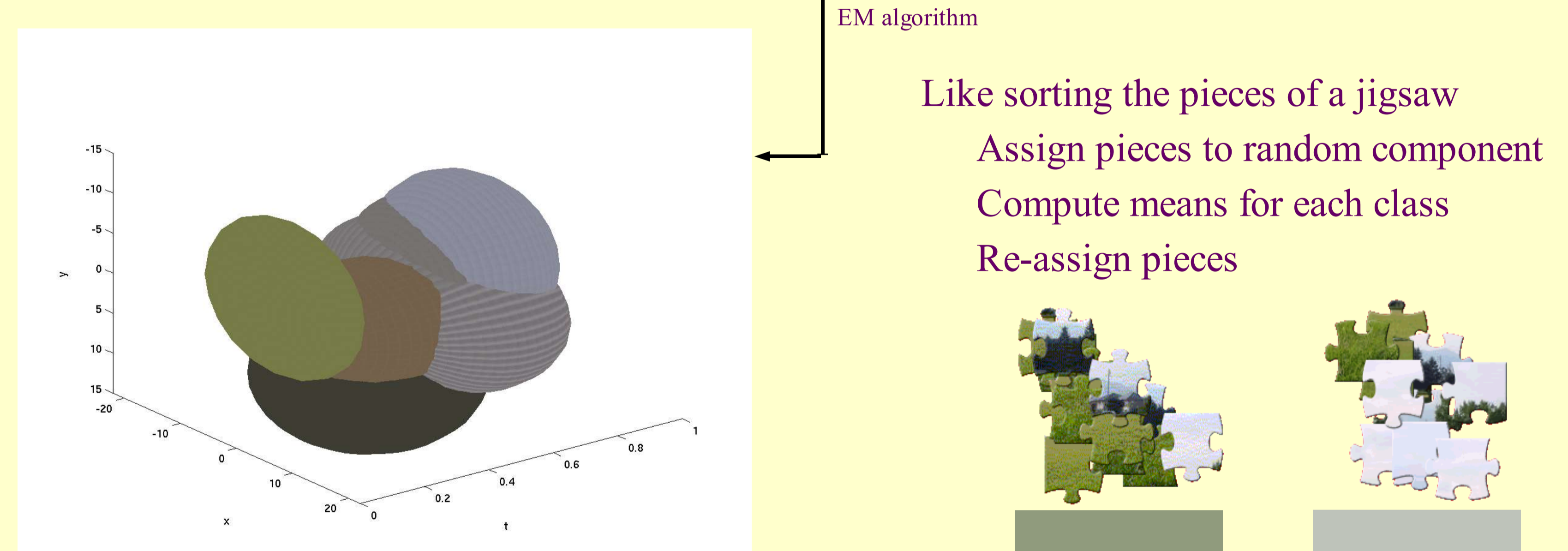
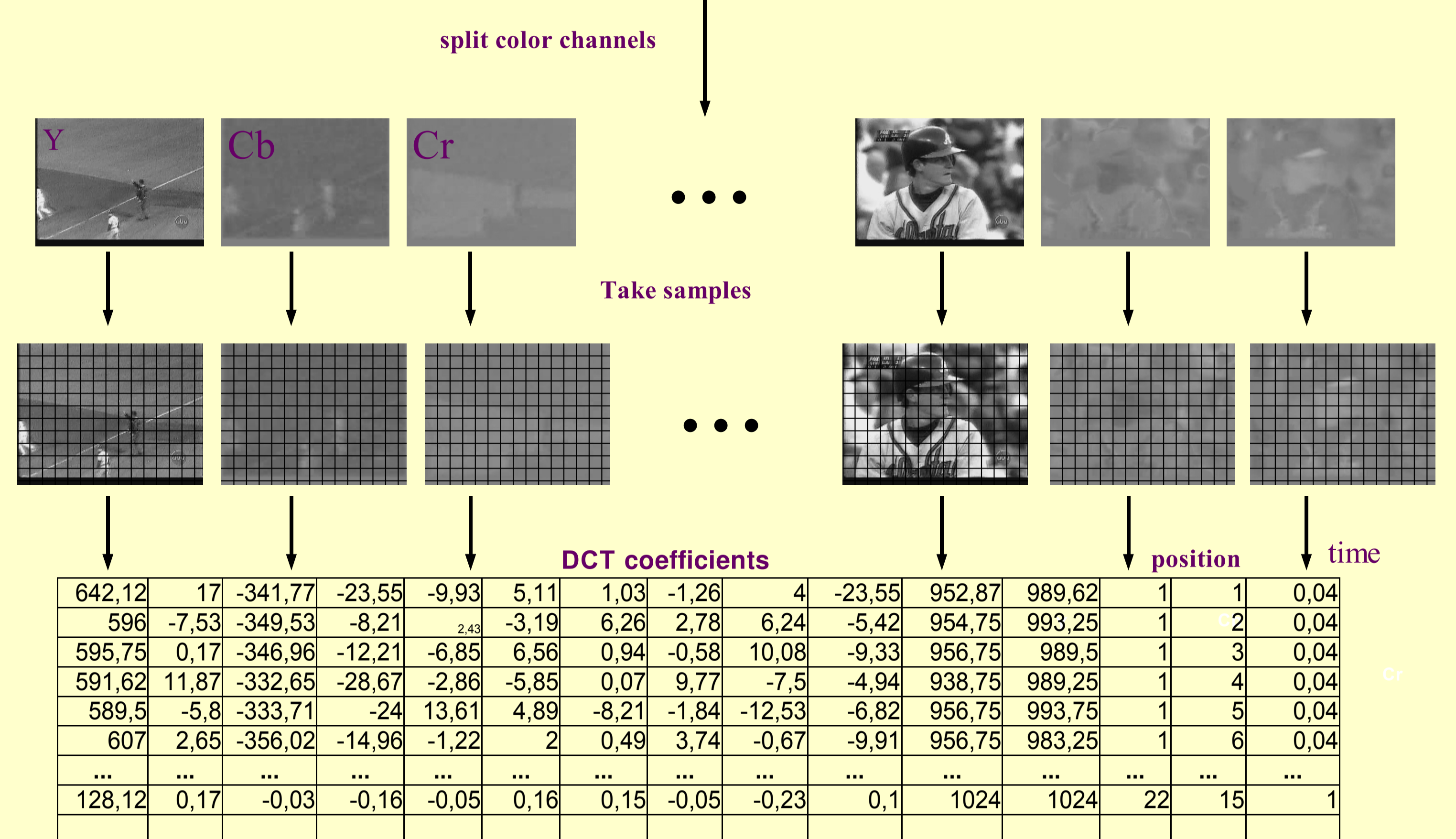
$$RSV(w_i) = \frac{1}{N} \sum_{j=1}^N \log[kP(x_j|w_i) + (1-k)P(x_j)]$$



Building Dynamic GMMs

The Dynamic GMM captures the disappearing of the grass in the video sequence. The corresponding component (green dynamic space-time blob) disappears at about $t=0.4$

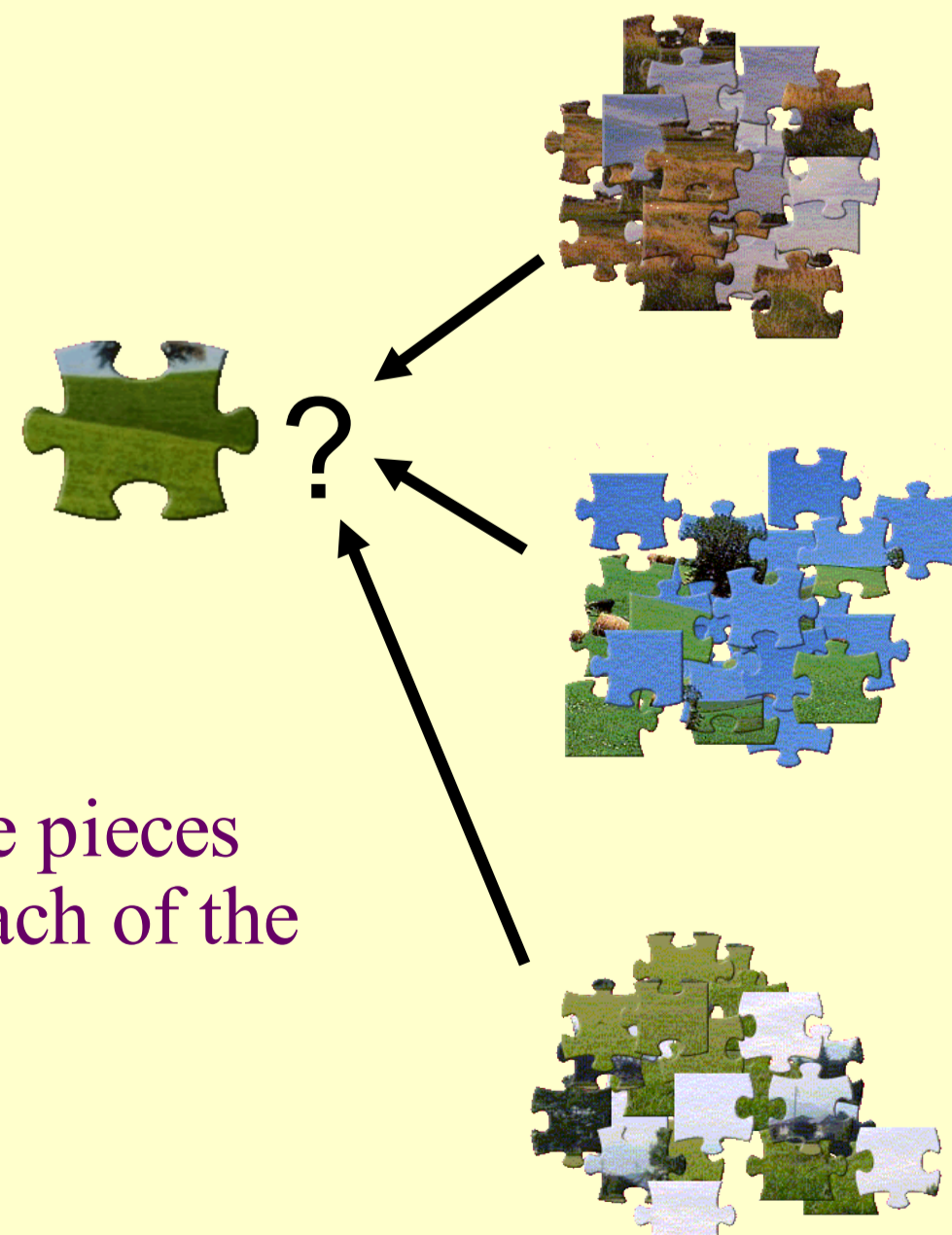
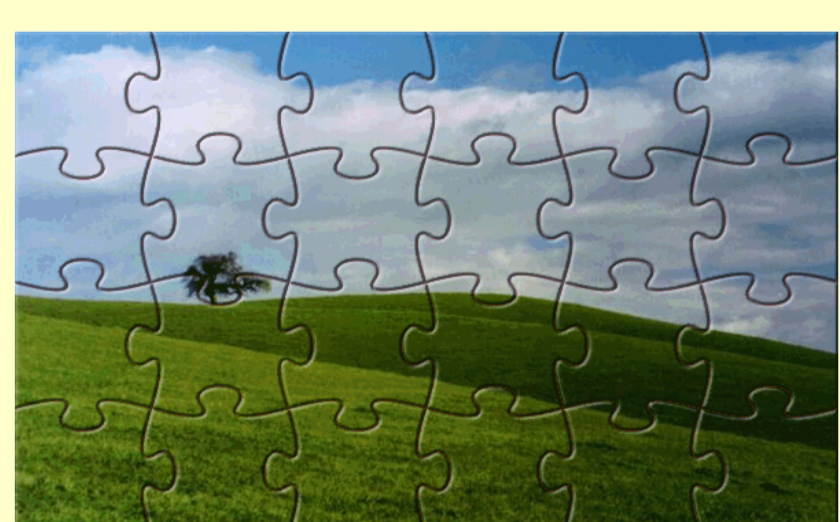
This effect is impossible to be captured from the static model where time is not taking into account



Probabilistic Image Retrieval

Example

Query:



What is the probability that the pieces of this query are drawn from each of the documents models?

Model [Vasconcelos 2000]

Gaussian Mixture Model (GMM)

each sample from an image is generated by 1 of C components

$$P(x|w_i) = \sum_{c=1}^{N_c} P(C_{i,c}) G(x, \mu_{i,c}, \Sigma_{i,c})$$

$$G(x, \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$

Conclusions/Future Work

Conclusions

- > appropriate keyframe less critical
- > can model (dis)appearance of objects
- > modeling sequence around the keyframe better than taking not connected frames from the full shot
- > combining textual and dynamic visual models better than ASR only
- > best results with multiple examples round-robin merged, topical filters and non-visual modalities

Future Work

- > more data needs more computation effort – optimizations ?
- > integration of audio
- > modeling of motion and texture over time