

## L' estrazione automatica dei metadati e loro affidabilita'

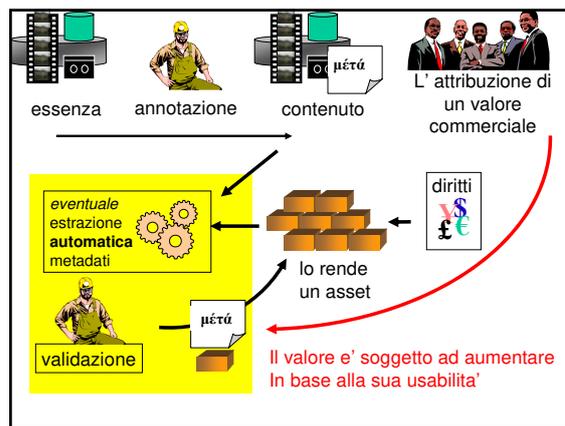
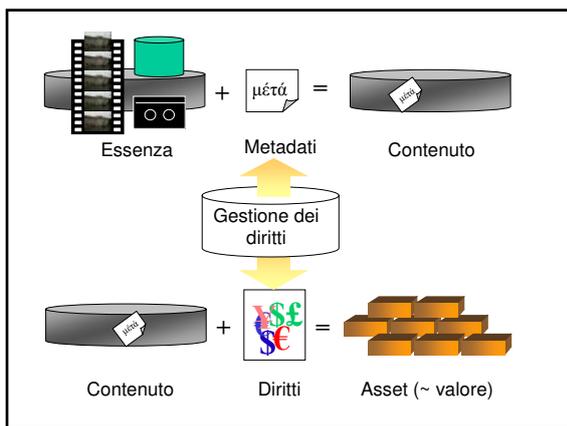
**Igino Manfre'**  
 Consulente in tecnologie televisive  
[igino.manfre@gmail.com](mailto:igino.manfre@gmail.com)



# μετά

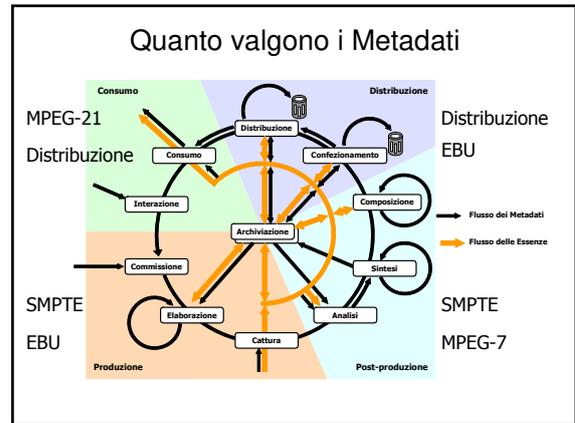
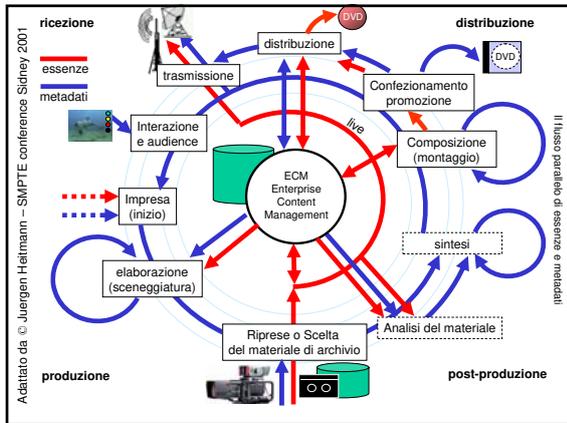
Oggi, ultimo incontro parliamo delle *finalita' del corso*

# sopra



## Il lexicon (molto vuoto) di SanJose

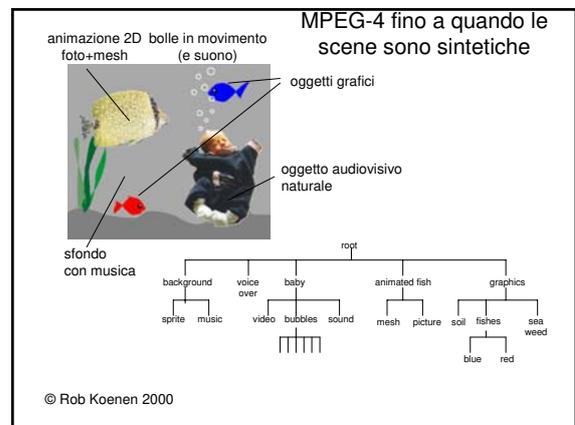
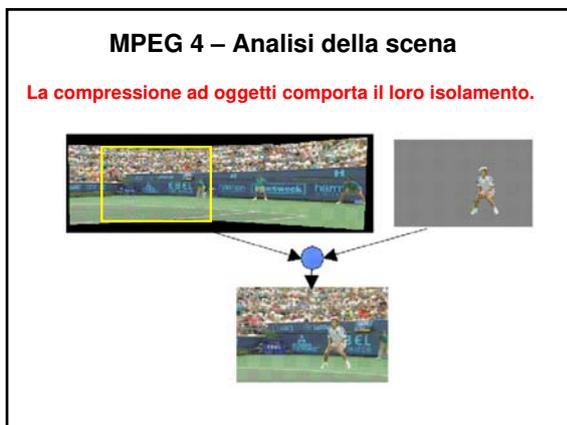
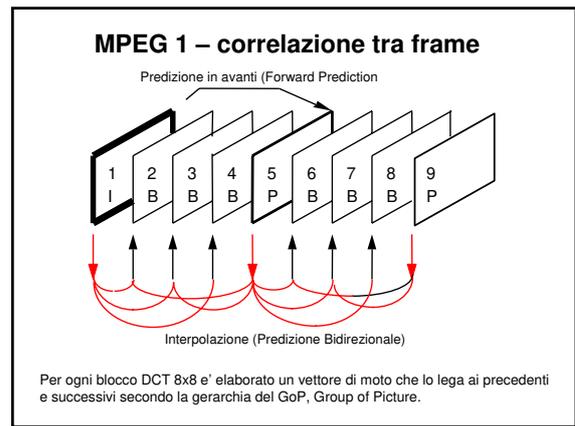
```
<?xml version="1.0" encoding="UTF-8" ?>
<!DOCTYPE complete SYSTEM "http://www.mpeg.org/2001/MPEG-7_Schema" xmlns="http://www.w3.org/2000/10/XMLSchema-instance"
  xmlns:schema="http://www.mpeg.org/2001/MPEG-7_Schema"
  <ContentDescription xsi:type="ContentEntityType">
    <MediaContent xsi:type="VideoType">
      <Video>
        <TemporalDecomposition>
          <VideoSegment>
            <TextAnnotation type="scene description" relevance="1" confidence="1" />
            <MediaTime>
              <MediaTimePoint time="00:00:00:00" />
              <MediaDuration time="PT1M30S" />
            </MediaTime>
            <SpatialTemporalDecomposition>
              <ISIRISegment>
                <ISIRISegment>
                  <TextAnnotation type="scene description" relevance="1" confidence="1" />
                  <MediaTime>
                    <MediaTimePoint time="00:00:04:20F30" />
                    <MediaDuration time="PT1M30S" />
                  </MediaTime>
                  <MediaTime>
                    <MediaTimePoint time="00:00:04:20F30" />
                    <MediaDuration time="PT1M30S" />
                  </MediaTime>
                </ISIRISegment>
              </ISIRISegment>
            </SpatialTemporalDecomposition>
          </VideoSegment>
        </TemporalDecomposition>
      </Video>
    </MediaContent>
  </ContentDescription>
</complete>
```

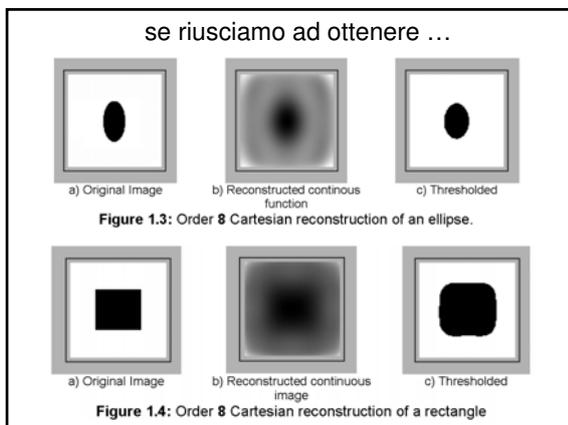
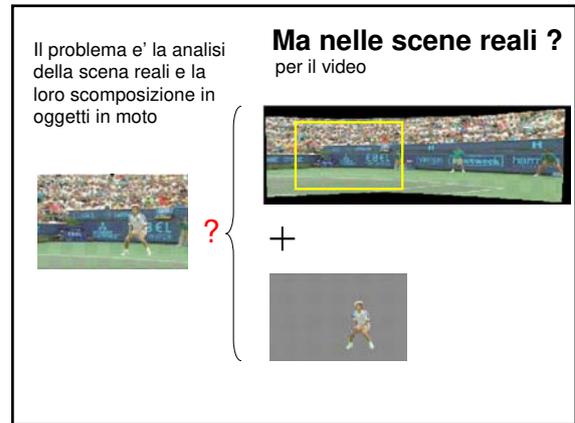


## AMWA [www.amwa.tv](http://www.amwa.tv)

Se tutte le imprese investono nella loro presenza in simili associazioni la ragione e' **solo** nei **diritti**

Arbitron, Inc., Autodesk, Automatic Duck, Inc., BPI Improve, BT Media & Broadcast, Digital Vision, E! Entertainment, eBus Limited, Marquis Broadcast, National Geographic, NBC Universal, Nielsen, Ninsight, OmniBus Systems, Panasonic, Perspective Media Group, Pro-Bel, SAIC, SGI Japan, VRT

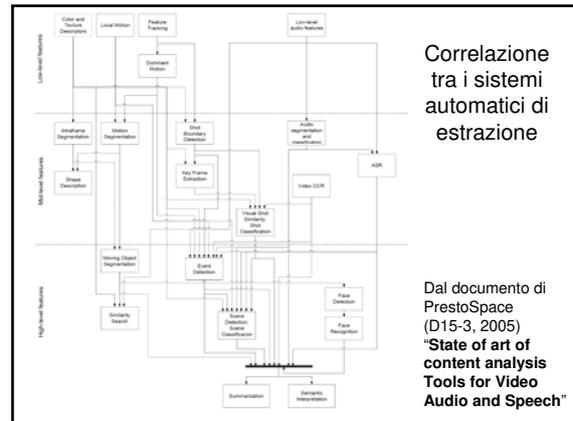




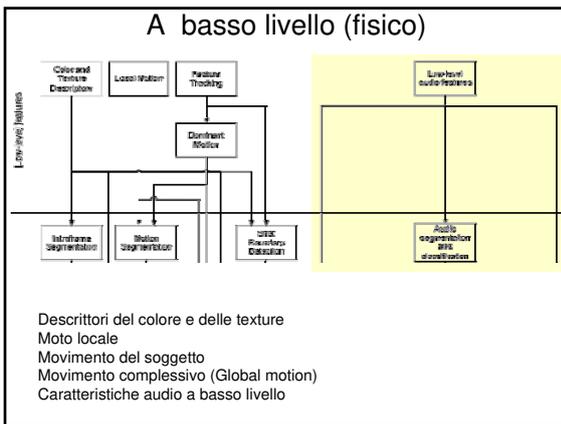
## Scopo della analisi del contenuto

Lo scopo dell'analisi del contenuto e' quello di descrivere il materiale in modo compatto ed efficiente scambiabile e in modo trasferibile.

Il fine e' fare in modo che queste descrizioni non sono ambigue e che siano creabili senza altri dati oltre quelli del materiale stesso.



## A basso livello (fisico)



## Analisi a basso livello

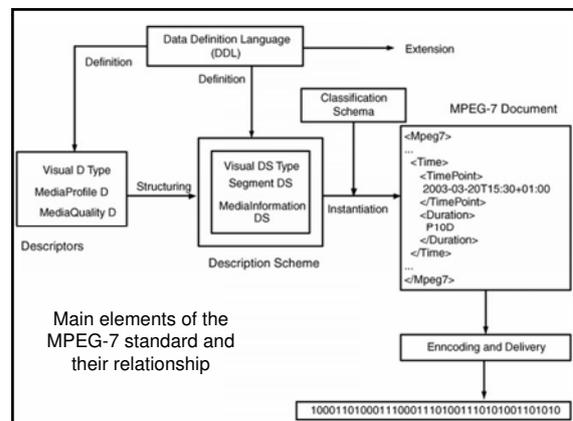
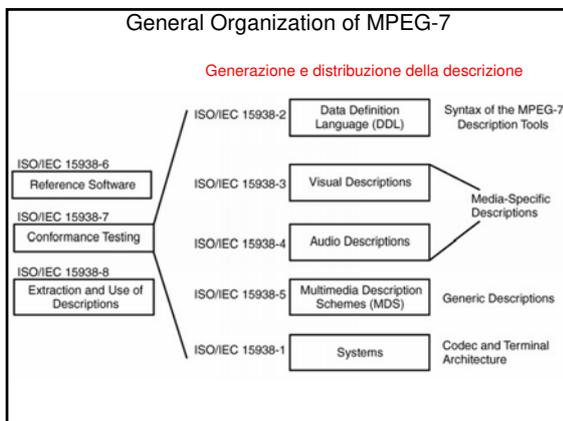
Le caratteristiche a basso livello sono singole proprieta' visuali legate al mezzo fisico come colore, texture, profilo, movimento. PrestoSpace li descrive come oggetto principale della nostro sistema visivo.

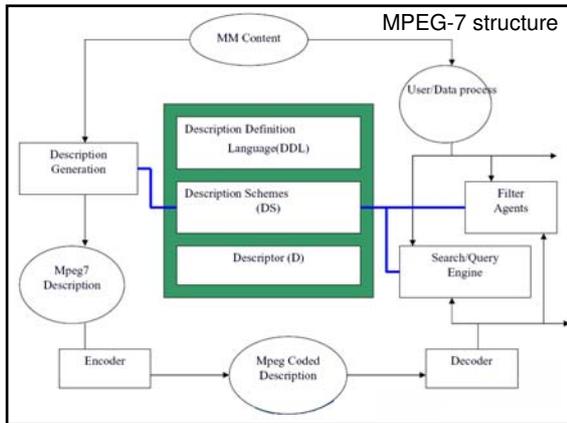
L'analisi a basso livello e' la base per un gran numero di algoritmi dei livelli superiori. Ogni descrittore consta di

- Rappresentazione: cosa e' contenuto nel descrittore e cosa significa.
- Estrazione (opzionale): come crearlo a partire da un dato visivo o sonoro
- Comparazione (opzionale): come determinare la collocazione nel relativo spazio, la distanza tra due descrittori dello stesso tipo.

L' MPEG-7 definisce un vasto set di descrittori visivi a basso livello.

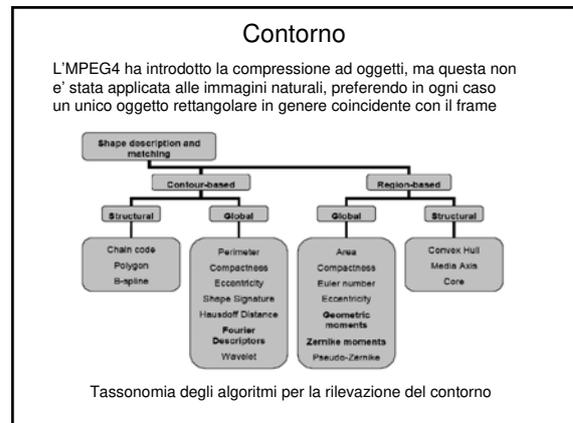
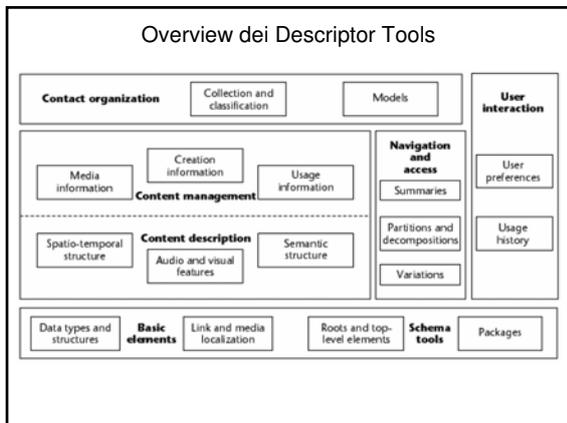
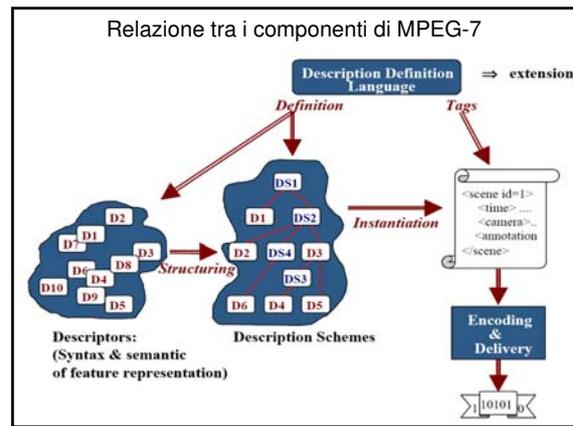
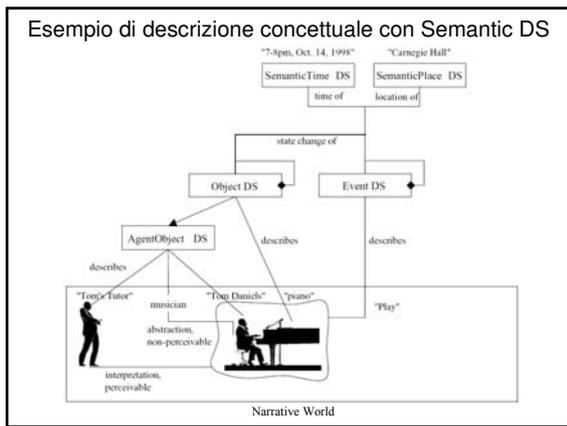
## General Organization of MPEG-7

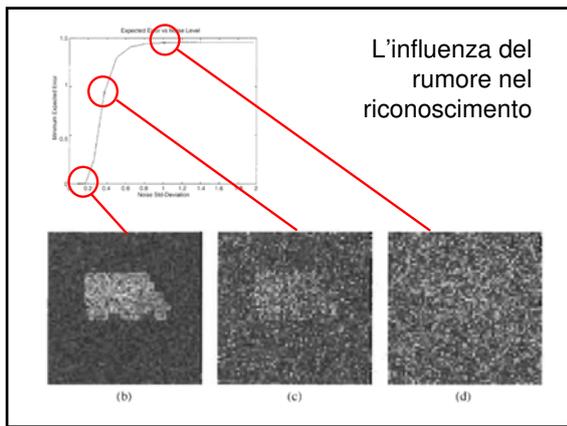
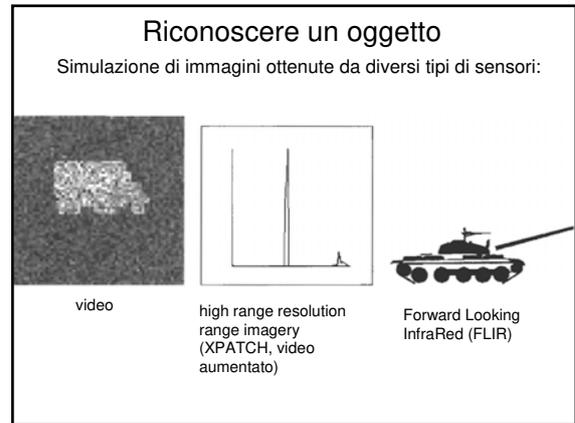
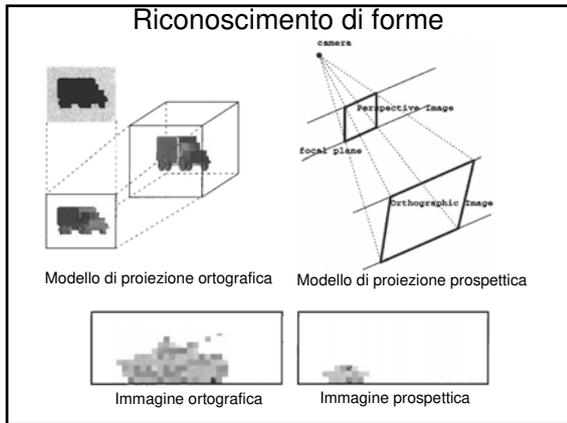




### Moduli dell' MPEG-7

MPEG-7 modules	Functionality	Features
Descriptors (Ds)	Low level audio or visual features	color ,texture, motion, audio energy and so on.
Description schemes (DSs)	Higher level audio or visual features	region, segments, objects, events and immutable metadata
Description definition language (DDL)	Based on XML scheme language	coding schemes has both text and binary format





### Movimento

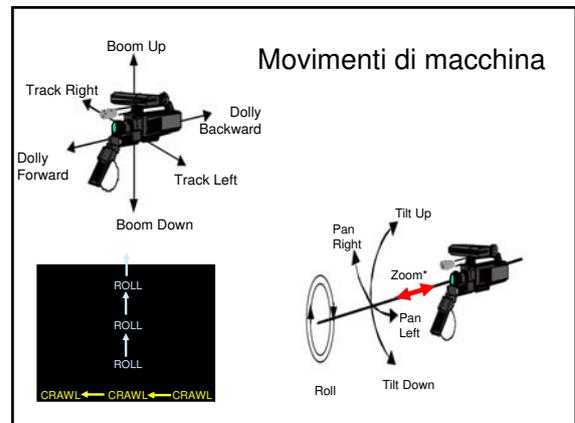
Si tratta di descrivere le caratteristiche di movimento degli oggetti nell'insieme della scena.

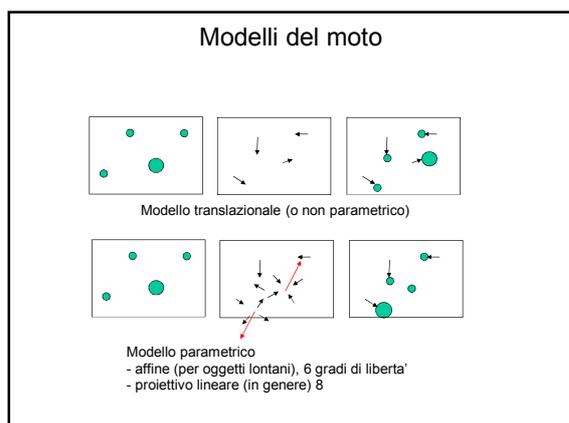
Il movimento possono essere relativo agli oggetti in primo piano (persone, edifici, etc) oppure al sistema di ripresa: cio' permette di descrivere lo sfondo in funzione dei movimenti di macchina (facilmente codificabili).

Tali relazioni possono essere estratte dalla analisi tridimensionale delle sequenze di immagini.

Vi sono sostanzialmente due tipi di moto:

- **Movimento globale** (global motion, g.m.) il movimento di tutti i pixel di una immagine (caso tipico i titoli di coda)
- oppure il **moto relativo** di un oggetto rispetto ad altri o allo sfondo.





- ### Stima e strategia di ricerca del moto
- Stima gerarchica su risoluzioni via via decrescenti (rende piu' facile l'analisi dei grandi movimenti, puo' essere applicato a qualsiasi modello)
  - Flusso Ottico estensione dirette dei modelli di moto, idoneo per grossi movimenti orientati (appunto: flussi)
  - Metodi basati sulla divisione in blocchi (tipici dell'MPEG1 e 2) che lavorano gerarchicamente alla ricerca di grossi movimenti di piccole aree.
  - Nel dominio della frequenza la correlazione di fase (utilizzata anche nell'ultima generazione di codificatori MPEG2) e metodi di stima Bayesiani. Una implementazione a bassa risoluzione di questi produce il caratteristico ringing della compressione MPEG4

### Analisi a basso livello audio

L'analisi a basso livello dell'audio e' molto sviluppata in quanto base per la i sistemi ad alto livello per il **ricoscimento automatico del parlato** (Automated Speech Recognition, ASR), funzione di analisi ad alto livello .

L'audio e' segmentato in piccole finestre (da 10 mSec in poi) ed e' poi processato sia nel dominio del tempo che in quello della frequenza o loro derivati.

Il dimensionamento di queste finestra puo' o meno rendere significativa l'elaborazione

Tra i domini piu' complessi c'e' il cepstrum, detto anche *spazio degli eventi*, sostanzialmente una trasformata di trasformata di fourier nel quale anche semplici operatori matematici (come il logaritmo) e' in grado di estrarre informazione particolarmente significativa.

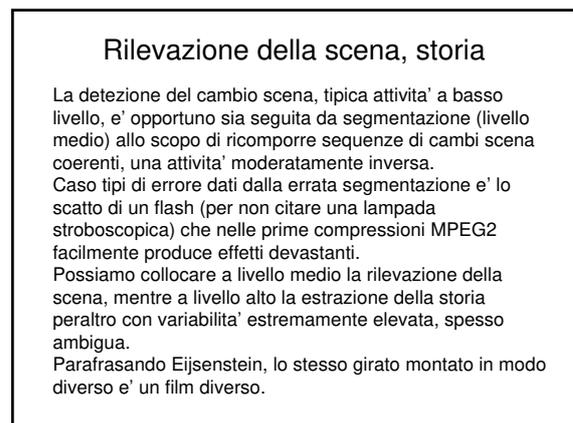
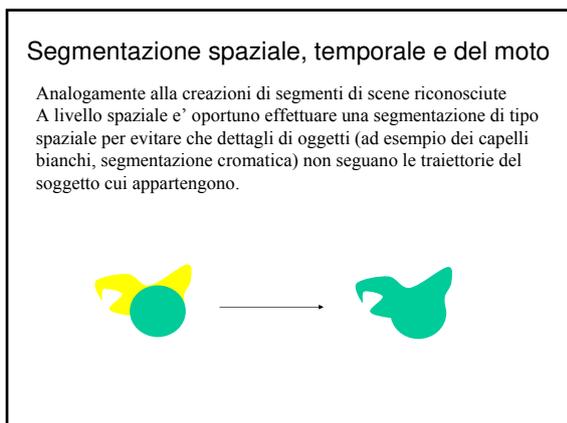
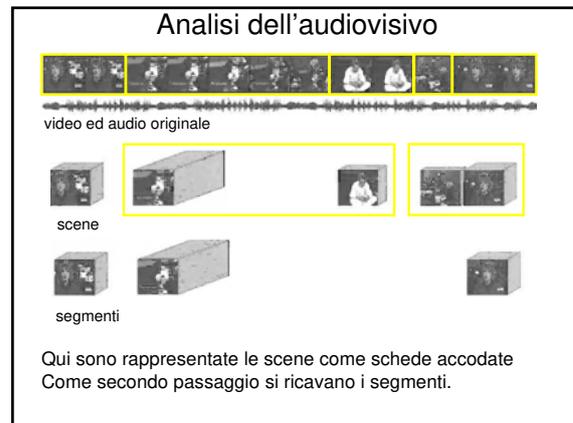
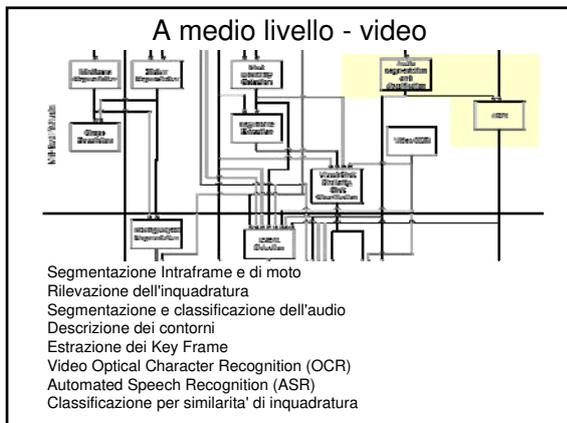
### Analisi a basso livello audio (2)

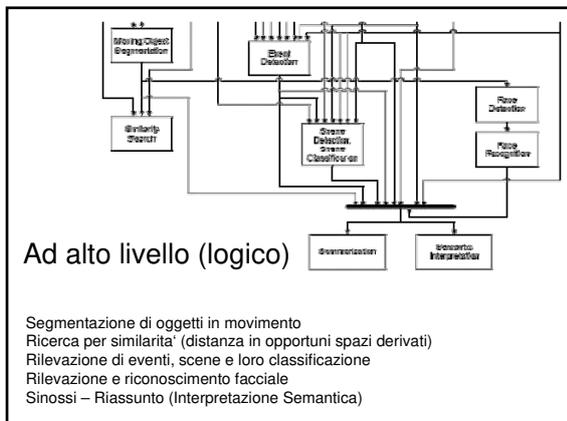
In particolare per quanto concerne gli stadi necessari alla conversione dei dialoghi (o del suono) al testo. La discriminazione tra musica e parlato, tra strumenti e oratori diversi, tra strumenti diversi, anche della stessa famiglia.

La complessita' degli algoritmi spesso comunque impedisce - al momento - la elaborazione in tempo reale.

La comprensione del linguaggio umano e' una collezione di analisi interdisciplinare ad alto livello, e costituisce a tutti gli effetti uno dei sistemi piu' affidabili per la generazione automatica di metadati.

Sono disponibili sistemi ASR per tutte le lingue dei paesi avanzati piu' l'urdu (per scopi militari)

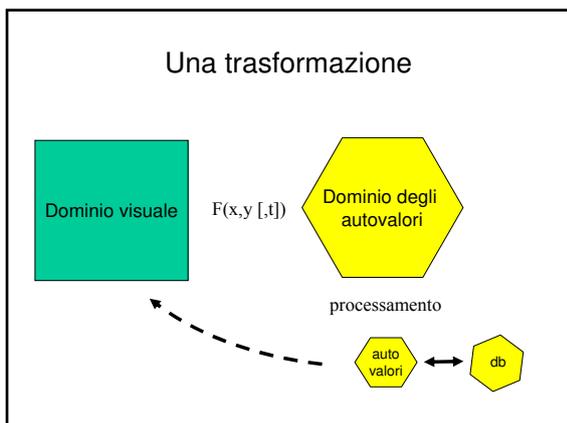




### Il riconoscimento del volto umano

Lo sviluppo di un modello computazionale per il riconoscimento degli oggetti naturali come il volto umano e' abbastanza difficile perche' il nostro riconoscimento e' basato su complessi stimoli significativi multidimensionali  
 Piu' che altro si puo' lavorare sulla densita' di probabilita' di apparizione (Probability Density Function, PDF) di un volto umano in due dimensioni.

Immagine da analizzare      Rilevazione del volto      Centatura del volto      Rilevazione della espressione facciale



### alla fine e' sempre una compressione

si tratta di ridurre l'informazione alla sola significativa

Immagine originale allineata      Immagine ricomposta dagli autovalori (85 byte)      Immagine compressa in JPEG (530 byte)

### Caratteri intrapersonali e extrapersonali

Gli autovalori possono far parte di due classi diverse

autovalori **intrapersonali** se le loro differenze' possono essere causate da sottili variazioni di espressione o illuminazione della ripresa

autovalori **extrapersonali** se tali diversita' fanno riferimento a differenze sostanziali (come il calvizie, presenza o meno di barba o occhiali etc).

### Non sempre un volto intero

Immagine in ingresso e a destra la probabilita'  $S(i,j)$  che il rettangolo comprenda un occhio sinistro

E' possibile **limitare l'esame a piccole aree della immagine** (gli occhi, la bocca, posizioni mutue con la punta del naso) che producono degli autovalori particolari, detti **eigenfeature** legati alle **espressioni**

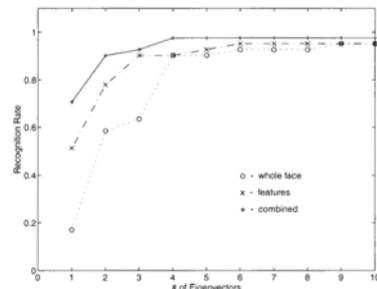
## Aspetto ed espressione

macchine da addestrare



A sinistra addestramento al riconoscimento delle **espressioni** effettuato limitando l'esame a piccole porzioni della immagine e a destra risultato del riconoscimento attraverso tale addestramento

## Piu' informazione, migliori risultati



Arece per il riconoscimento delle espressioni e rapporto di corretto Riconoscimento migliore in caso di uso di pochi autovalori misti

## autovalori diversi per volti e espressioni

Immagini in test



Volti riconosciuti



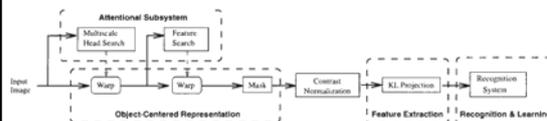
Espressioni riconosciute



## Workflow



Ricerca della testa in diverse inquadrature (multiscale)



## 92% di risultati positivi

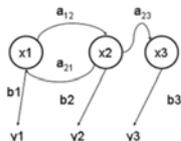


Alcune delle foto utilizzate per testare l'accuratezza del riconoscimento facciale. Il successo e' stato del 92%, l'errore del riconoscimento ha presentato una deviazione standards di 15 gradi

## Trascrizione dei dialoghi

Le applicazioni di riconoscimento automatico del parlato (ASR), hanno goduto di notevoli finanziamenti di provenienza militare, e solo negli ultimi 10 anni dalla unione europea (dove si parlano 15 lingue!). Per le applicazioni relative alla metadateazione **non** possono essere utilizzati gli algoritmi che fanno uso di raffronto con un vocabolario di fonemi, e **neanche** sistemi che necessitino di addestramento. I sistemi in uso si basano sui modelli markoviani nascosti (Hidden Markov Model, HMM), dei quali il riconoscimento vocale e' stato negli anni '70 una delle prime applicazioni.

## Andrey Markov (1856 – 1922)

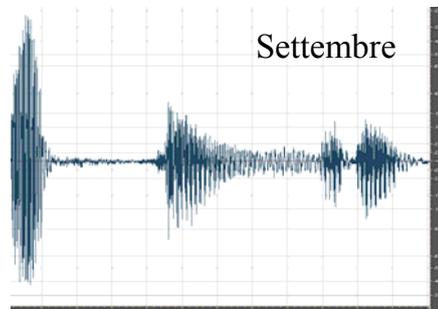


Un **processo di Markov** è un processo stocastico nel quale la probabilità di transizione che determina il passaggio ad uno stato di sistema dipende **unicamente dallo stato di sistema immediatamente precedente e non** dal come si è giunti a tale stato

Una **catena di Markov** è un processo di Markov a stati discreti, ovvero è un processo stocastico discreto per cui ad ogni istante  $t$  estraiamo dal processo una variabile casuale discreta

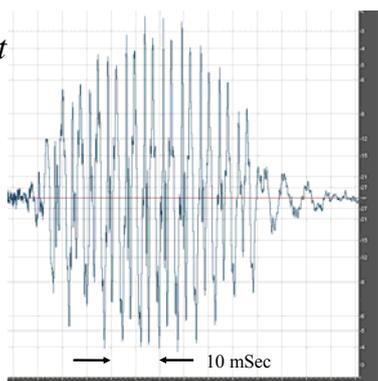
Un **Modello Nascosto di Markov** è una catena di Markov i cui stati non sono osservabili direttamente.

Una delle possibili ragioni dell'uso dei modelli nascosti di Markov e' che il discorso umano puo' essere considerato come una sequenza eventi stazionari di breve durata, ad es. 10 millisecondi.



Milena Gabanelli che pronuncia *settembre* in 600 msecondi

Set



## Si, ma come funziona?

Semplificando, si utilizza una FFT (Fast Fourier Transform) con una finestra da 10 mSec, e si produce un vettore di coefficienti cepstrali (correlazione DCT della trasformata di Fourier). L'applicazione di un modello nascosto di Markov elaborando i piu' significativi di questi valori (i primi, quelli piu' lontani dalla specificita' del singolo speaker) produce per ogni fonema una distribuzione di probabilita' dipendente (al piu') a fronte del precedente. Concatenando questa attivita' si puo' contare di riconoscere il discorso.



C'e' da considerare che poiche' il riconoscimento automatico del parlato ha una **grande importanza nel settore militare**, si e' raggiunta l'eccellenza in quasi tutte le lingue del primo mondo... ed in molte del terzo...

## Non perdere mai di vista l'obiettivo

### 3.1 Basic definitions

Let  $A \subset \mathbb{N}$ ,  $A = 0, \dots, H-1 \times 0, \dots, W-1$ , where  $H$  and  $W$  are the height and width of a frame, respectively.

**Definition 3.1 (Frame)** A frame  $f_t$  is a function from  $A$  to  $\mathbb{N}$  where for each spatial position  $(x, y)$  in  $A$ ,  $f_t(x, y)$  represents the value of the pixel  $(x, y)$  at time  $t$ .

**Definition 3.2 (Video)** A video  $V$  is a sequence of frames  $f_t$  with an accompanying audio track and can be described by  $V = (f_t)_{t \in [0, T-1]}$ , where  $T$  is the number of frames contained in the video. The number of frames is directly associated with the frequency and the duration of visualization.

Il formalismo matematico e' una necessita': qui non si tratta di contestarlo, ma per favore, una volta partorite definizioni come queste qua sopra, tornate alla realta', pensate alla pellicola, al fotogramma, alla scansione, alla cromaticita'.

## Database degli atteggiamenti



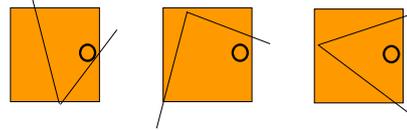
Esempio di ciascuna classe del database degli atteggiamenti elaborato dal progetto europeo MUSCLE (Multimedia Understanding through Semantics, Computation and Learning, progetto IST FP6-507752, da "Progress on Applications of Machine Learning Techniques" (Deliverable: D8.1c, 2005)

## ...e sua efficacia



Matrice di confusione per il database video degli atteggiamenti pesati uniformemente (a sinistra) e in base alla classe (a destra).  
 E' difficile sapere cosa e' vero in queste dichiarazioni e cosa e' artefatto.  
 Concediamoci il diritto di poter verificare in pratica qualora servissero.

## Inquadratura



Esempi di ciascuna delle tre classi del database "Basketball" di MUSCLE

## Fine

Spero abbiate potuto trovare degli spunti interessanti.

**Ricordate che la curiosità e' la molla della invenzione.**

Manderò una e-mail a tutti per informarvi quando sarà possibile scaricare il materiale (stato in cui rimarrà per un paio di settimane).

[Igino.manfre@tiscali.it](mailto:Igino.manfre@tiscali.it)

