



# Recuperación de Información Multimedia

## Procesamiento de Videos (Codecs, Shots, Keyframes)

**CC5213 – Recuperación de Información Multimedia**

Departamento de Ciencias de la Computación

Universidad de Chile

Juan Manuel Barrios – <https://juan.cl/mir/> – 2019



# Videos

- CODEC “compressor-decompressor”
  - Compresión de frames
    - h.261, h.264, mpeg-2, Xvid, DV, ... etc.
  - Compresión de audio
    - MP3, speex, vorbis, aac, ... etc.
- Container
  - Guarda los frames comprimidos.
  - Guarda la información de las pistas (video, audio).
  - Guarda los metadatos del video.
  - Ej: avi, mov, mpg, mkv, ogv, etc.



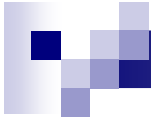
# Estándares

- MPEG (Moving Pictures Experts Group), grupo de trabajo ISO/IEC, creado en 1998.
  - MPEG-1, MPEG-2, MPEG-4, MPEG-7, MPEG-21.
- VCEG (Video Coding Experts Group or Visual Coding Experts Group ), grupo de trabajo de ITU-T, creado 1994.
  - h.261, h.262, h.263, h.264.
- Foco:
  - Intentar hacer un decoder simple.
  - El encoder es más complejo que el decoder.
  - Los estándares se enfocan en normar el decoder.



# Estándares MPEG

- MPEG-1: Estándar inicial de compresión en 5 partes. 1993.
  - Parte 1: Definición de archivo container (mpg).
  - Parte 2: Codificación de video. Basado en h.261.
  - Parte 3: Codificación de audio, 3 formatos posibles:
    - Layer I (mp1), Layer II (mp2), Layer III (mp3).
- MPEG-2: 11 partes. 1995.
  - Parte 1: Definición de archivo container incluyendo streams.
  - Parte 2: Extensión de MPEG-1, incluye calidad DVD y HDTV = h.262.
  - Parte 3: Extensión de MPEG-1, incluye multi-channel (5.1).
  - Parte 7: Advanced Audio Coding (AAC).
- MPEG-4: 1998, 28 partes, aún en desarrollo.
  - Parte 1, 12, 14, 15: Protocolos de transmisión y formato del container (mp4).
  - Parte 2: Compresión de video = h.263.
  - Parte 10: Advanced Video Coding (AVC) = h.264.
  - Parte 3: Inclusión de audio lossless y otros encoders.



# Otros estándares MPEG

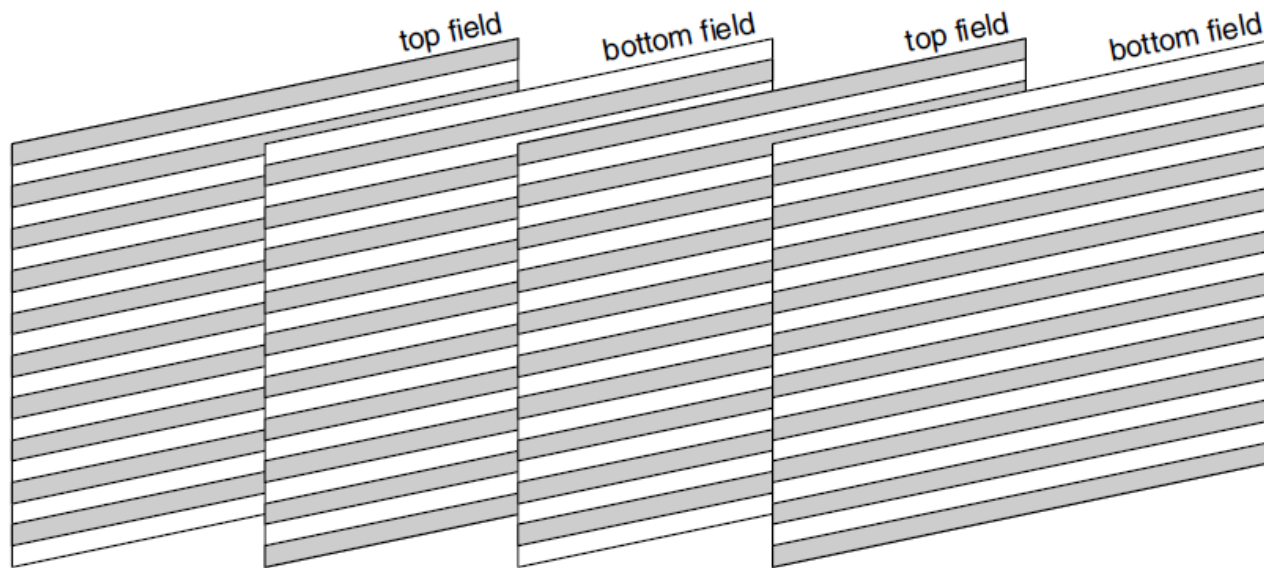
- MPEG-7:

- ☐ Estándar para incluir metadatos en contenido multimedia.
- ☐ Metadata de alto nivel y bajo nivel (descriptores visuales).

- MPEG-21:

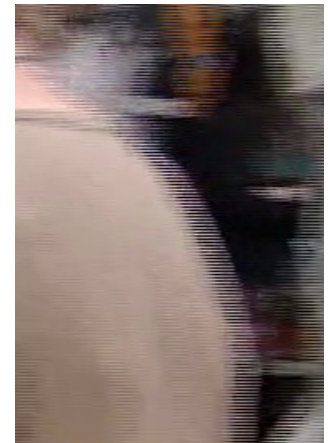
- ☐ Marco para intercambio de contenido multimedia.
- ☐ Mercado digital, restringiendo derechos de autor.

# Interlaced Videos



# Interlaced Videos

- Aparecen tramas horizontales y bordes dobles cuando hay movimiento:



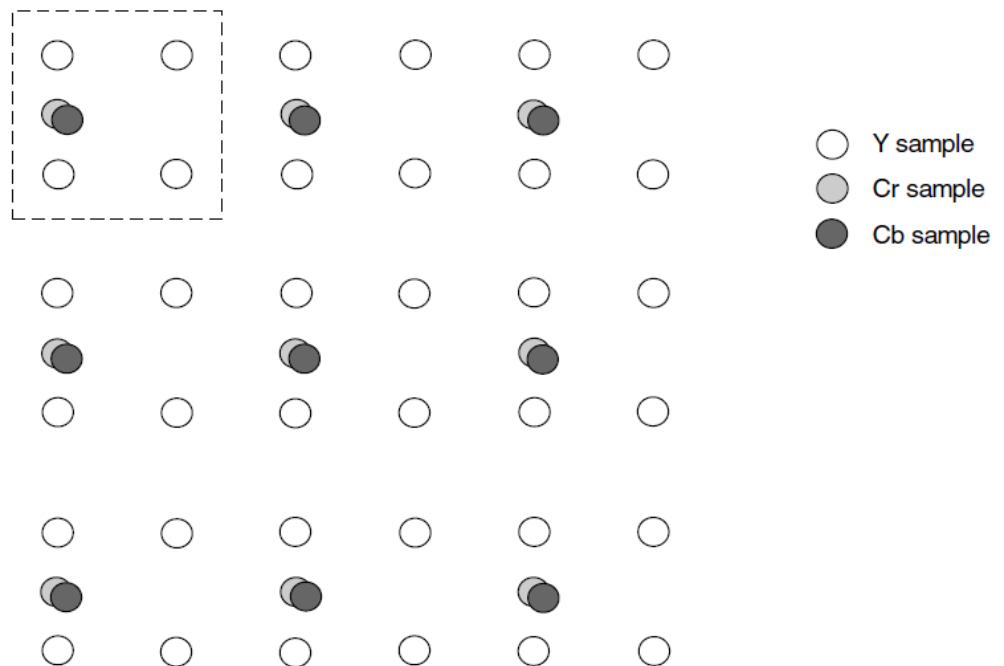


# **Codificación de videos MPEG-1**

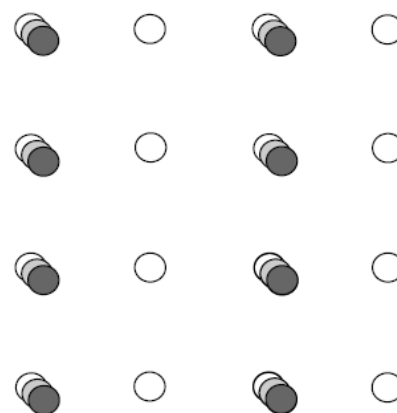


# Colores

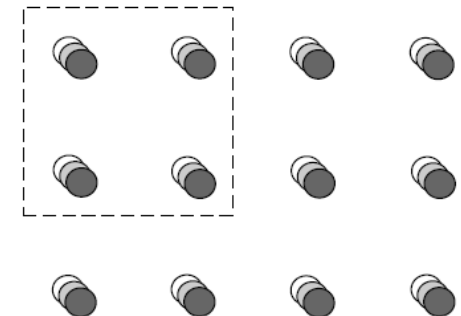
- YUV
- Reducción de los canales de color



4:2:0 sampling

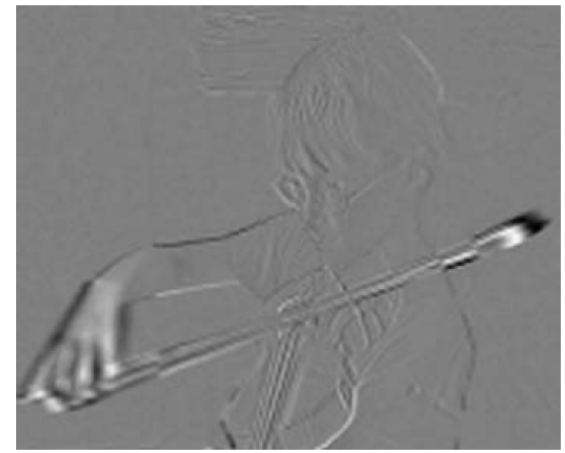


4:2:2 sampling



4:4:4 sampling

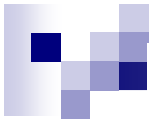
# Optical Flow





# Codificación de frames

- Usualmente dos frames consecutivos son muy similares.
- En vez de guardar dos frames consecutivos en forma independiente basta con guardar sólo la diferencia.
  - Si tiene muchos ceros tendrá mejor compresión.
  - Imagen Residual = Frame 2 – Frame 1



Frame 1



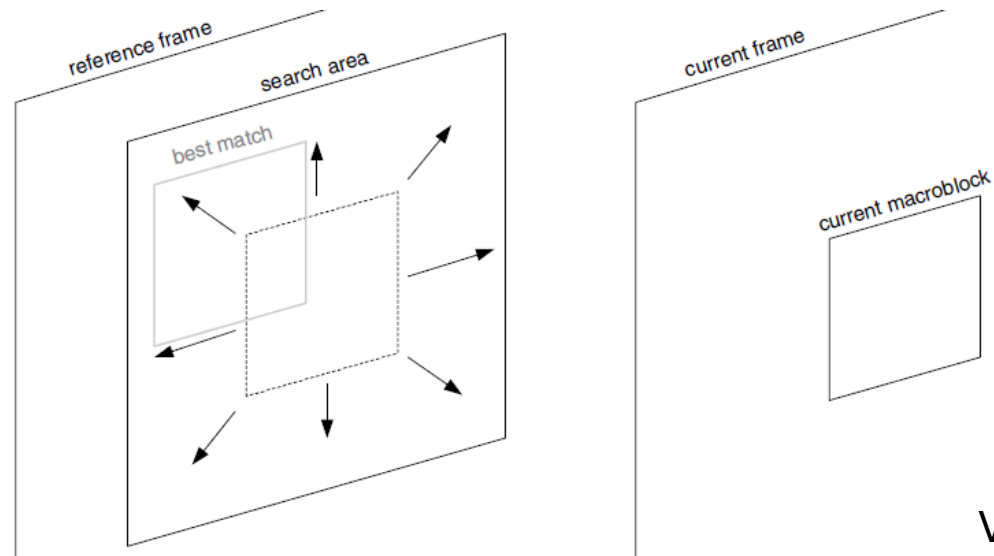
Frame 2



Imagen  
Residual

# Estimación de movimiento por bloques

- Se hace una estimación del movimiento usando bloques:
  - Se divide un frame en “macrobloques” de  $N \times N$ .
  - Cada macrobloque del frame 2 se forma con algún bloque de  $N \times N$  del frame 1 más la imagen residual.
  - El vector que apunta al lugar de donde obtener el bloque base desde el frame 1 se denomina “Motion Vector”.





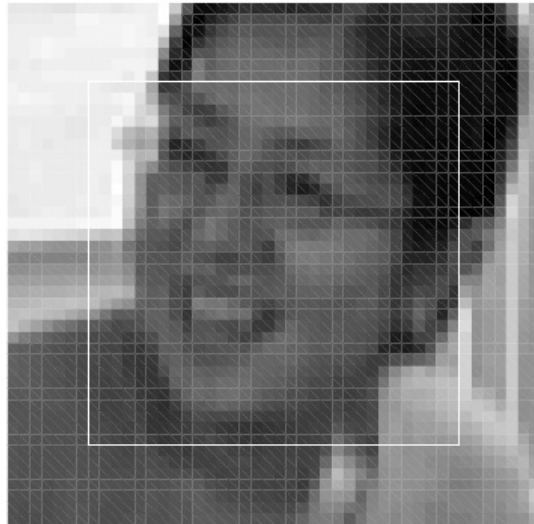
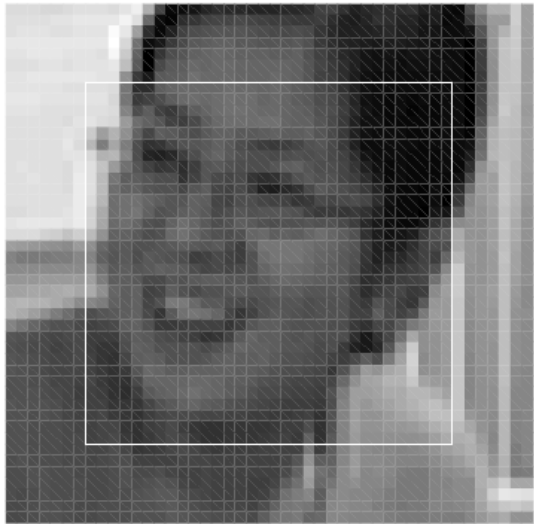
# Estimación del movimiento

- Dado dos frames, buscar la posición del macrobloque actual dentro de la imagen previa que minimiza el error:

Mean Squared Error: 
$$MSE = \frac{1}{N^2} \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} (C_{ij} - R_{ij})^2$$

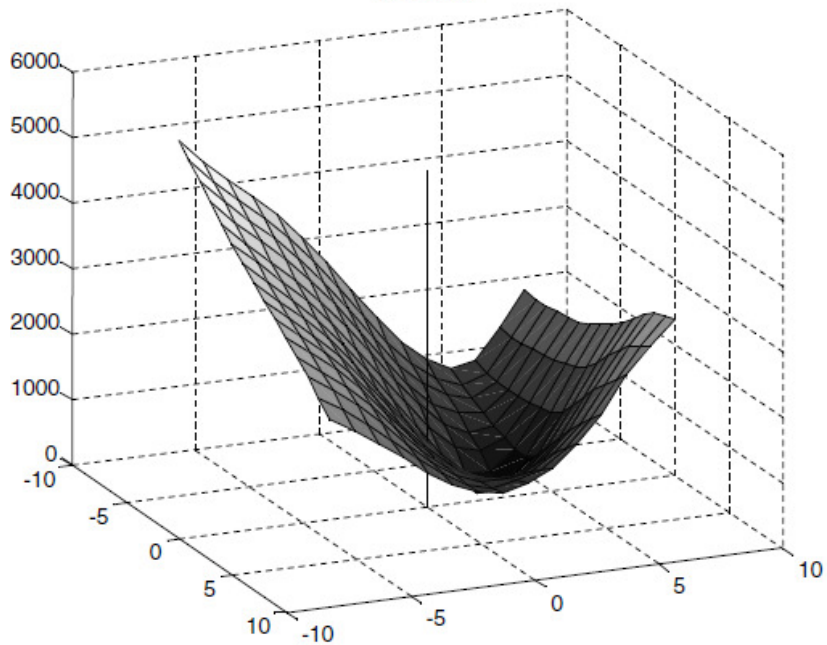
Mean Absolute Error: 
$$MAE = \frac{1}{N^2} \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} |C_{ij} - R_{ij}|$$

Sum of Absolute Errors: 
$$SAE = \sum_{i=0}^{N-1} \sum_{j=0}^{N-1} |C_{ij} - R_{ij}|$$

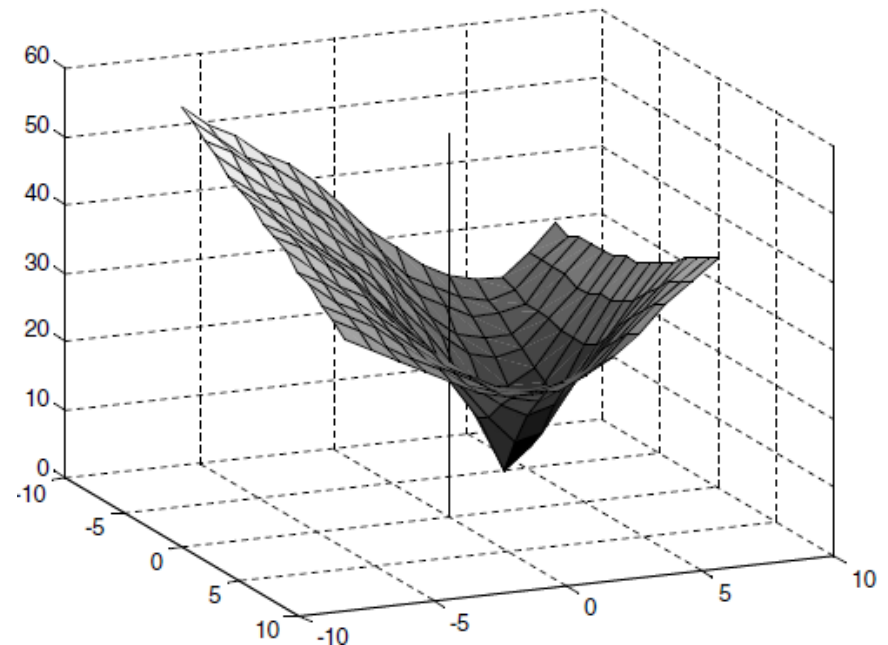


Mínimo error para  
el desplazamiento  
(2,0)

MSE map



MAE map



Ver Richardson, cap 7

# Motion Vectors para macrobloques 4x4

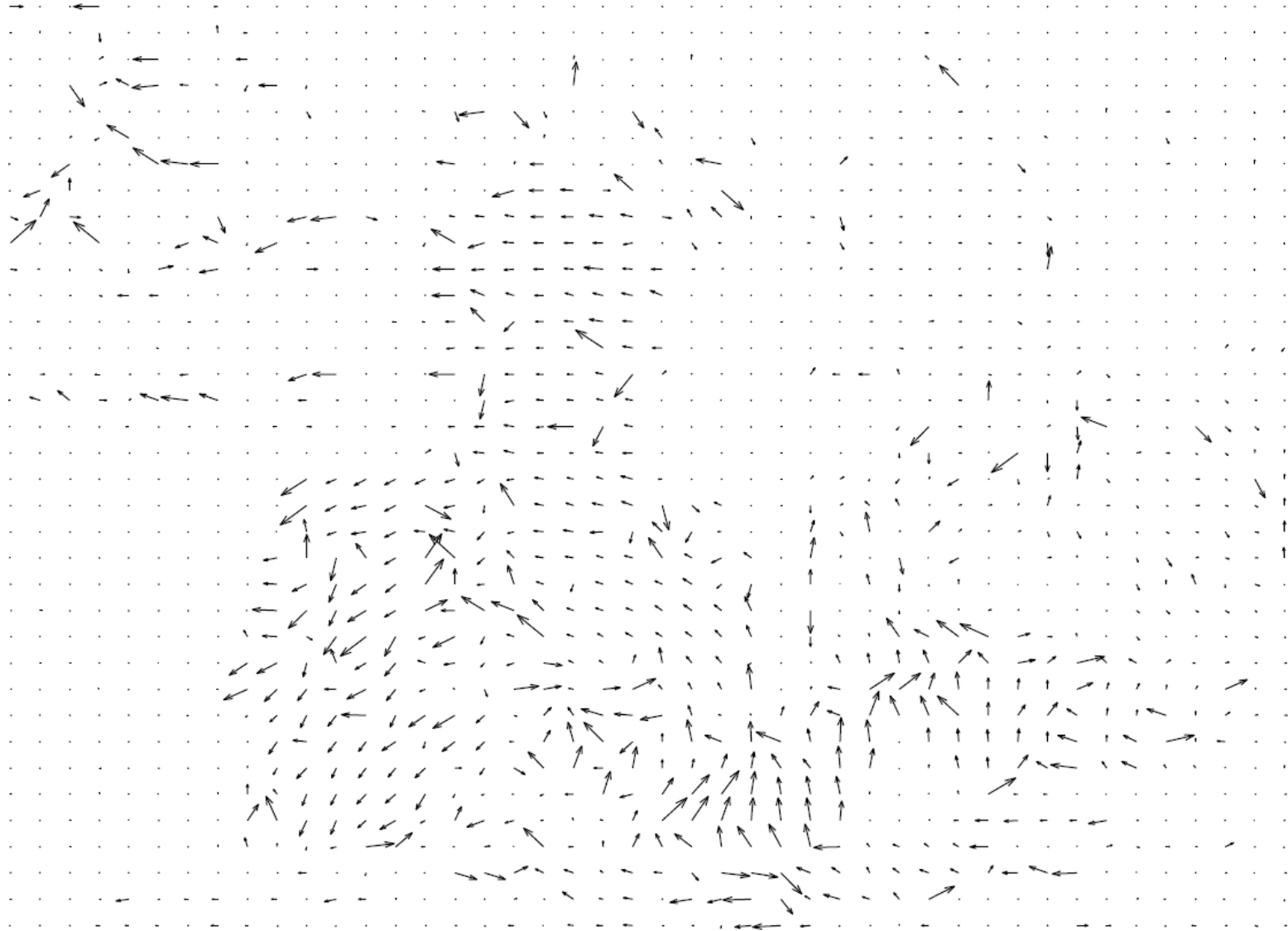




Imagen residual luego de ajustar frame 2 con motion vectors para macrobloques de 4x4

Frame 2 = Frame 1 + motion vectors + imagen residual



Imagen residual luego de ajustar frame 2 con motion vectors para macrobloques de 16x16

(aumenta la información en la imagen residual)



Ver Richardson, cap 3



# Codificación con MPEG-1

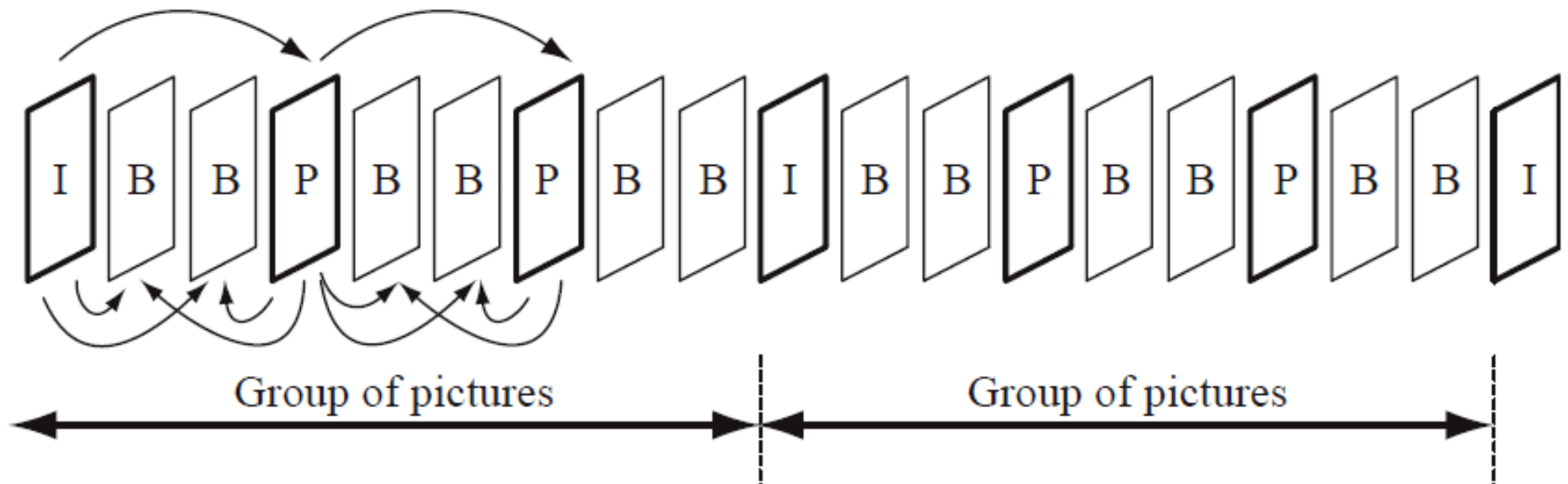
- Dado un video, se quieren comprimir los frames para reducir el tamaño del archivo.
- La compresión se basa en:
  - Estimación del movimiento usando motion vectors.
  - Transformación de la imagen residual con DCT, como en jpg.
  - Codificación de la entropía (codificación huffman o codificación aritmética), como en jpg.



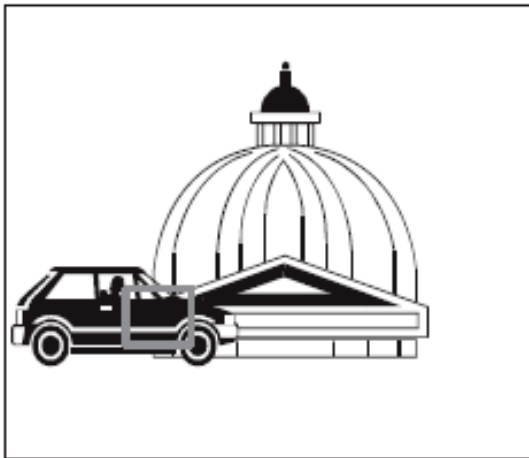
# Tipos de frames

- Los frames del video se clasifican en 3 tipos:
  - Intra-coded (frames I), se comprime como una imagen estática.
  - Predictive coded (frames P), se comprime usando motion vectors con un frame I o P previo.
  - Bidirectional predicted (frames B), se comprime usando como referencia frames I o P previos y posteriores.

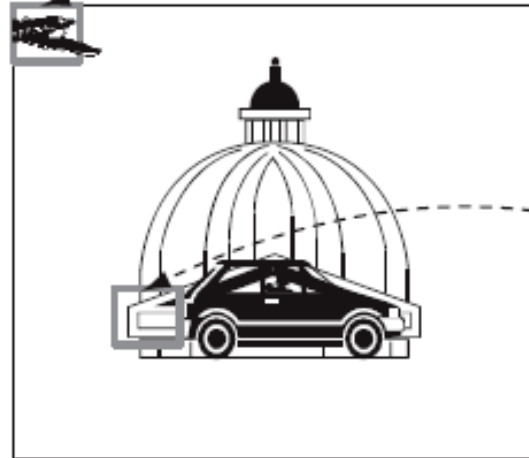
# Tipos de frames



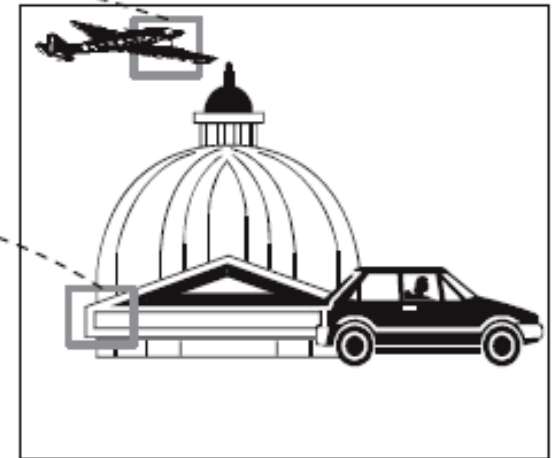
# Frames B



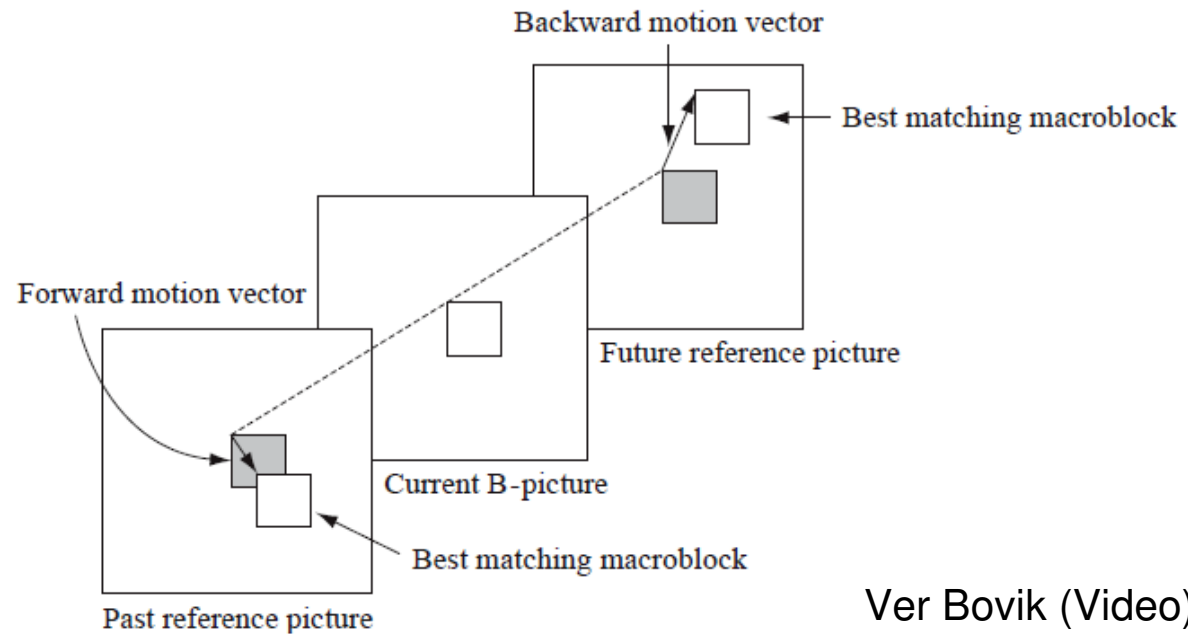
Frame  $N - 1$



Frame  $N$



Frame  $N + 1$





# Tipos de frames

## ■ Frames I

- ☐ No dependen de otro frame.
- ☐ Su compresión es relativamente baja.
- ☐ Finaliza la propagación de errores previos.

## ■ Frames P

- ☐ Dependen del frame I o P previo.
- ☐ Propagan los errores que pueden existir en frames I o P previos

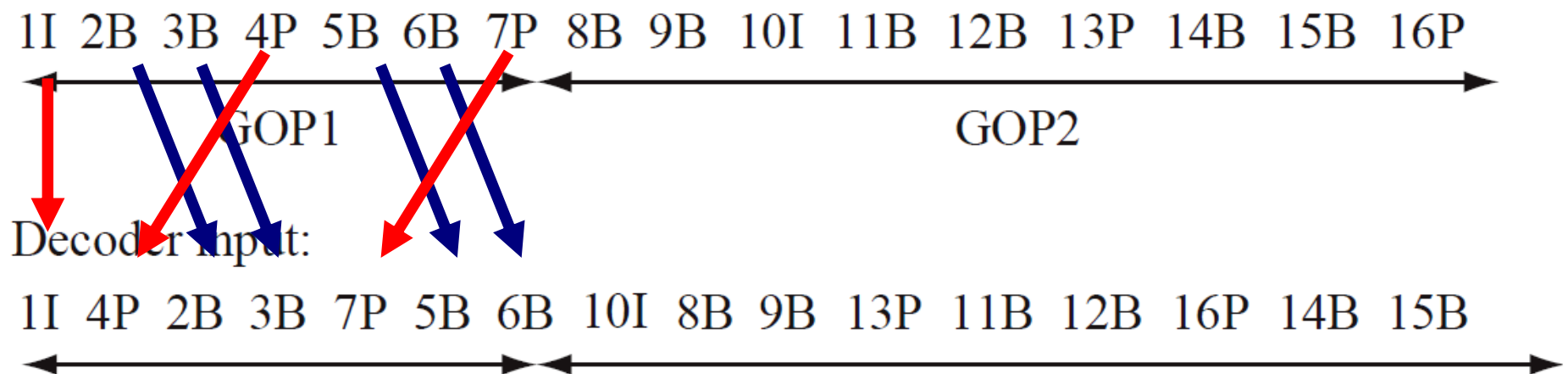
## ■ Frames B

- ☐ Dependen de frames previos y posteriores
- ☐ Mayor compresión.
- ☐ Para poder decodificarlos se debe decodificar frame posterior.
- ☐ No hay dependencias sobre frames B por lo que un error en un frame B no se propaga (pero propaga los errores de los frames en los que depende).

# Reordenamiento de frames

- Los frames no se guardan en orden correlativo, si no que se deben guardar primero los frames I o P y luego los B.
- Se requiere un buffer de decodificación.

Encoder input:







# **Detección de Shots**



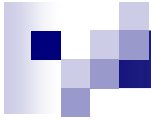
# Videos

- Frame: unidad mínima
- Shot: secuencia continua de frames procedente de una cámara que representa una acción continua en el tiempo y espacio
- Scene: conjunto de shots en una misma ubicación



# Detección de límites de shots

- Los frames pertenecientes a shots distintos presentan un cambio en su contenido
  - Detectar discontinuidades en el flujo del contenido de los frames
- En general, extraer un descriptor global al frame  $i$  y al frame  $i+1$ , calcular la distancia  $d(i, i+1)$ 
  - Si es mayor a un umbral entonces hay un cambio de shot



# Detección de límites de shots

- Diferencia de frames:
  - Hay un cambio de shot si la distancia  $L1$  entre frames consecutivos es mayor a un umbral
- Cantidad de pixeles cambiados:
  - Se define que un pixel cambia cuando la diferencia de intensidad entre dos frames supera un umbral2
  - Hay un cambio de shot cuando el número de pixeles que cambian sea mayor a un umbral1
- Reducir la imagen o usar filtro gaussiano para reducir ruido



# Detección de límites de shots

## ■ Diferencias estadísticas

- Dividir cada frame en zonas y conocer la media y varianza del canal Y para cada zona en el video
- Cuando las zonas se alejan de la media hay un límite de shot

## ■ Histogramas

- Cambio cuando la distancia entre histogramas consecutivos supera un umbral
- 4x4 zonas, histograma por zona, eliminar las 8 zonas con más cambios, hay cambio cuando la suma de las 8 menores supera un umbral



# Detección de límites de shots

- Falso positivo:
  - Fotografías con flash.
  - Comparar 2 pares:
    - Diferencia es  $\min\{ d(i, i+1), d(i-1, i+2) \}$
- Falso negativo:
  - Transiciones suaves entre shots.
  - Detectores específicos para transiciones (fade-in, fade-out)



# Detección de límites de shots

- TRECVID durante 2001-2007 evaluó la detección de shots:
  - Resultados en:  
<http://www-nlpir.nist.gov/projects/tvpubs/tv.pubs.org.html#2007>
- Algunos Papers:
  - J.S. Boreczky and L.A. Rowe. “Comparison of video shot boundary detection techniques”. 1996.
  - S.Eickeler and S.Müller. “Content-Based Video Indexing Of Tv Broadcast News Using Hidden Markov Models”. 1999.



# **Selección de Keyframes**





# Selección de Keyframes

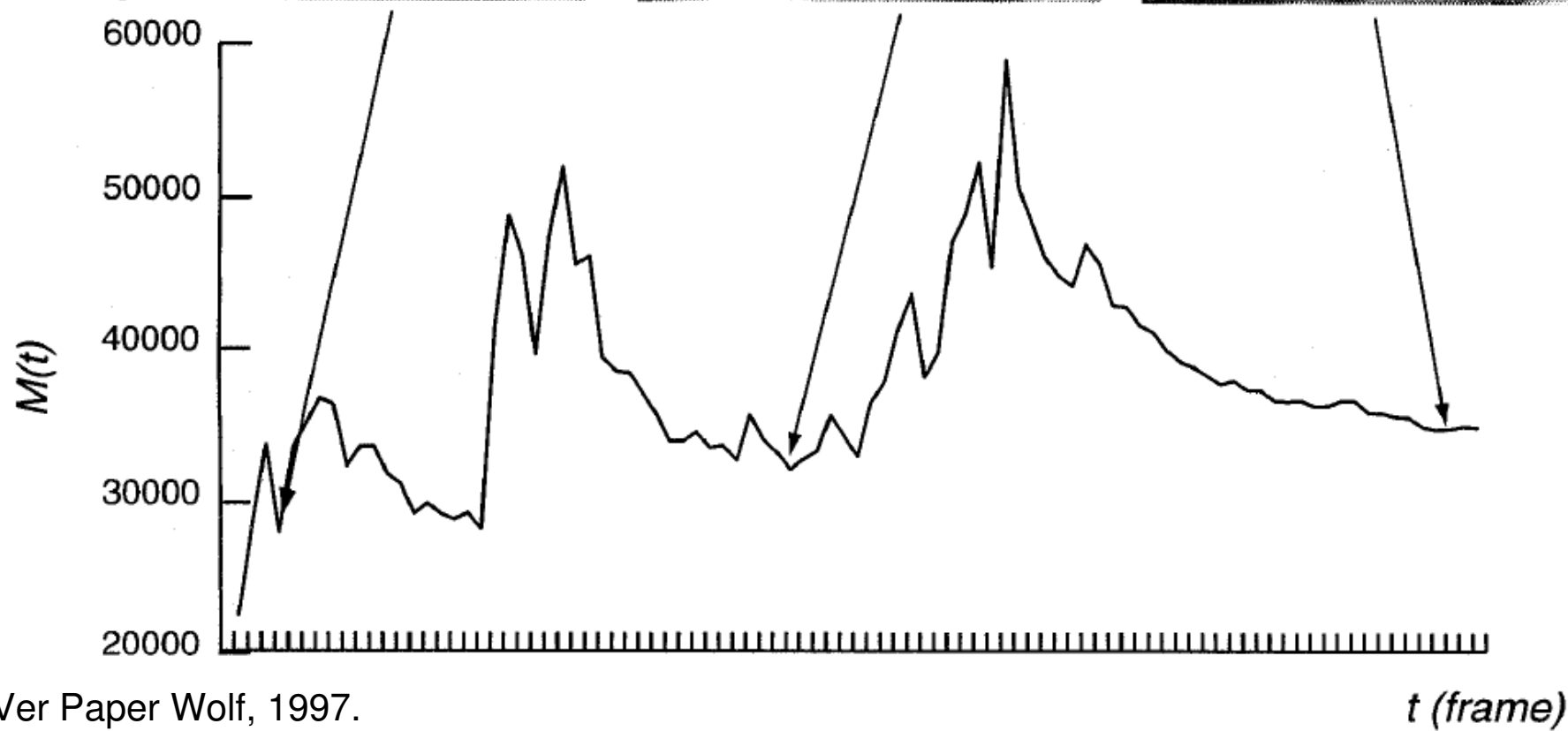
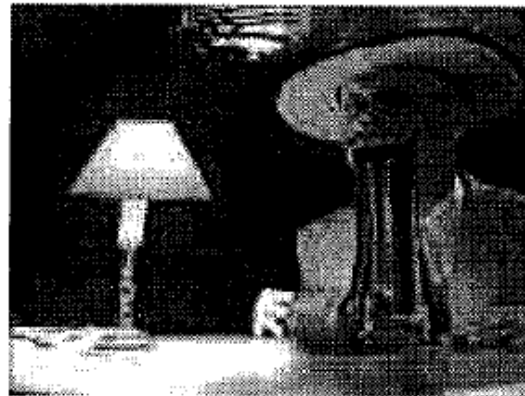
- Selección constante, por ej.:
  - ☐ 1 frame por segundo
  - ☐ 5 frames por segundo
  - ☐ 1 frame cada 3 segundos
- Calcular un descriptor global para todos los frames, clusterizar, y seleccionar los frames más cercano a los centroides



# Selección de Keyframes

- Dividir en shots, para cada shot tomar los frames estables
  - Menor diferencia con el anterior según un descriptor global
  - Menor movimiento según el optical flow

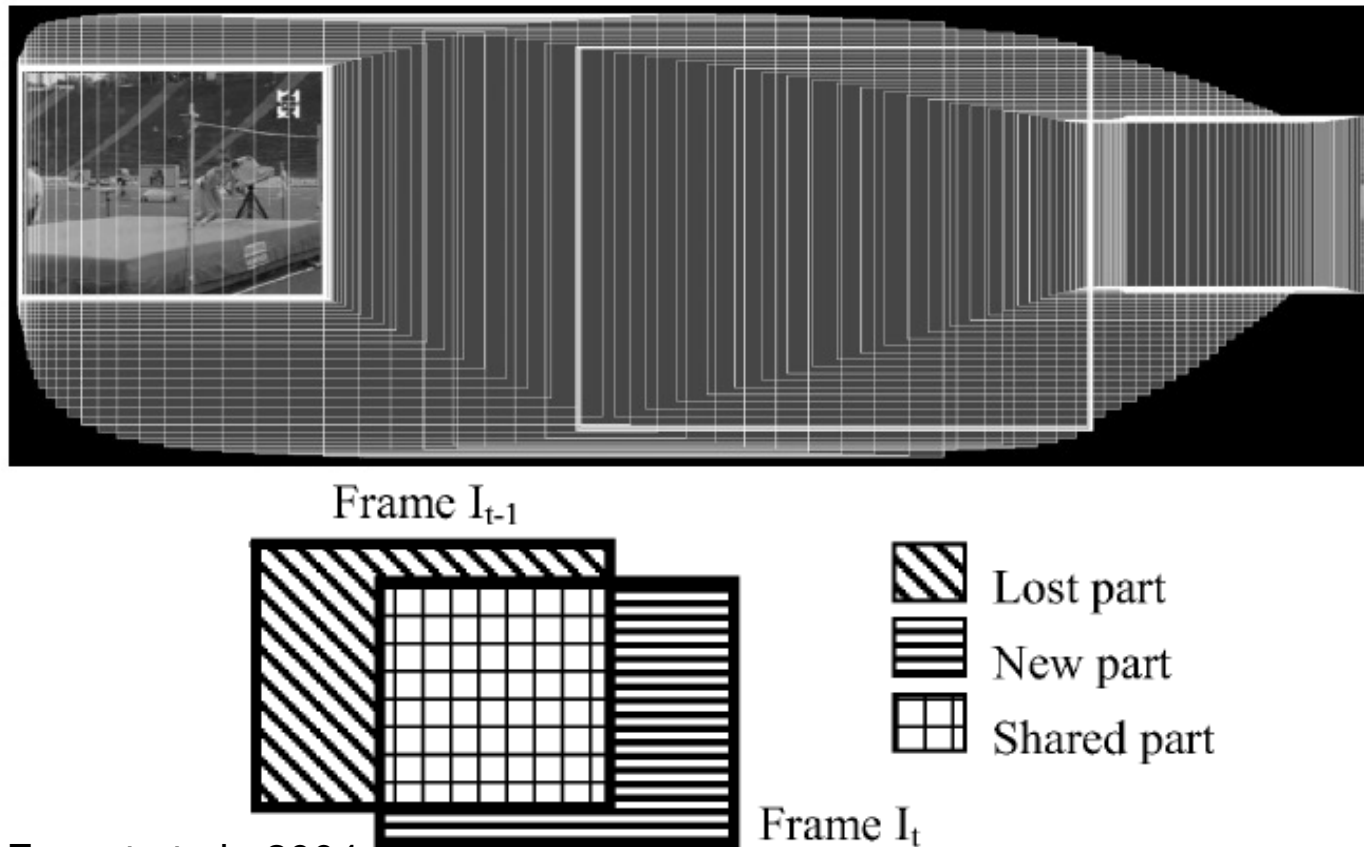
$$M(t) = \sum_i \sum_j |o_x(i, j, t)| + |o_y(i, j, t)|$$



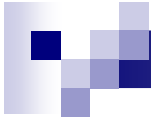
Ver Paper Wolf, 1997.

# Selección de Keyframes

- Por estimación del movimiento:



Ver Paper Fauvet et al., 2004.



# **Imagen + Audio**



# Combinación de descriptores

- En un video, para un segmento se puede obtener un descriptor visual y un descriptor de audio. Por ej:
  - Visual: Promedio de los valores de un descriptor de bordes (como Edge Histogram)
  - Visual: Promedio de los histogramas de color de todos los frames dentro del segmento
  - Audio: Promedio del vector de coeficientes de la escala Mel o del cepstrum
- ¿Como usar toda esta información?



# Late Fusion

- Utilizar cada uno por separado y luego combinar los resultados
- Hacer búsquedas por cada descriptor en forma independiente y combinar los resultados
- Simple, pero requiere que se puedan obtener buenos resultados con cada uno por separado.



# Early Fusion

- Crear una distancia combinada. Por ejemplo:

$$\delta_{av}(q, r) = \frac{w_1}{\tau_1} * L_1\text{-Eh}(q, r) + \frac{w_2}{\tau_2} * L_1\text{-Rgb}(q, r) + \frac{w_3}{\tau_3} * L_1\text{-Aud}(q, r)$$

- ¿Cómo escalar las distancias para que queden en rangos comparables?
  - Ver capítulo de multimétricas





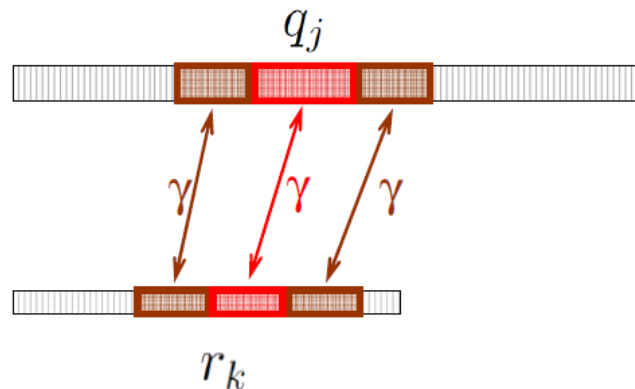
# Early Fusion

- Se realiza una sola búsqueda con la distancia combinada
- Puede localizar elementos que requieran de ambas modalidades a la vez
- Problema: La distancia es afectada por descriptores ruidosos
  - Se pueden descartar descriptores ruidosos en forma dinámica
  - Ver capítulo de multimétricas

# Distancia Temporal

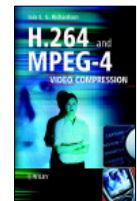
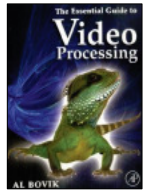
- Se puede aumentar la robustez de la distancia si se incluyen los segmentos anteriores y posteriores. Para ventana  $W$ :

$$\delta(q_j, r_k) = \frac{1}{W} \sum_{w=-\lfloor W/2 \rfloor}^{\lfloor W/2 \rfloor} \gamma(q_{j+w}, r_{k+w})$$



# Bibliografía

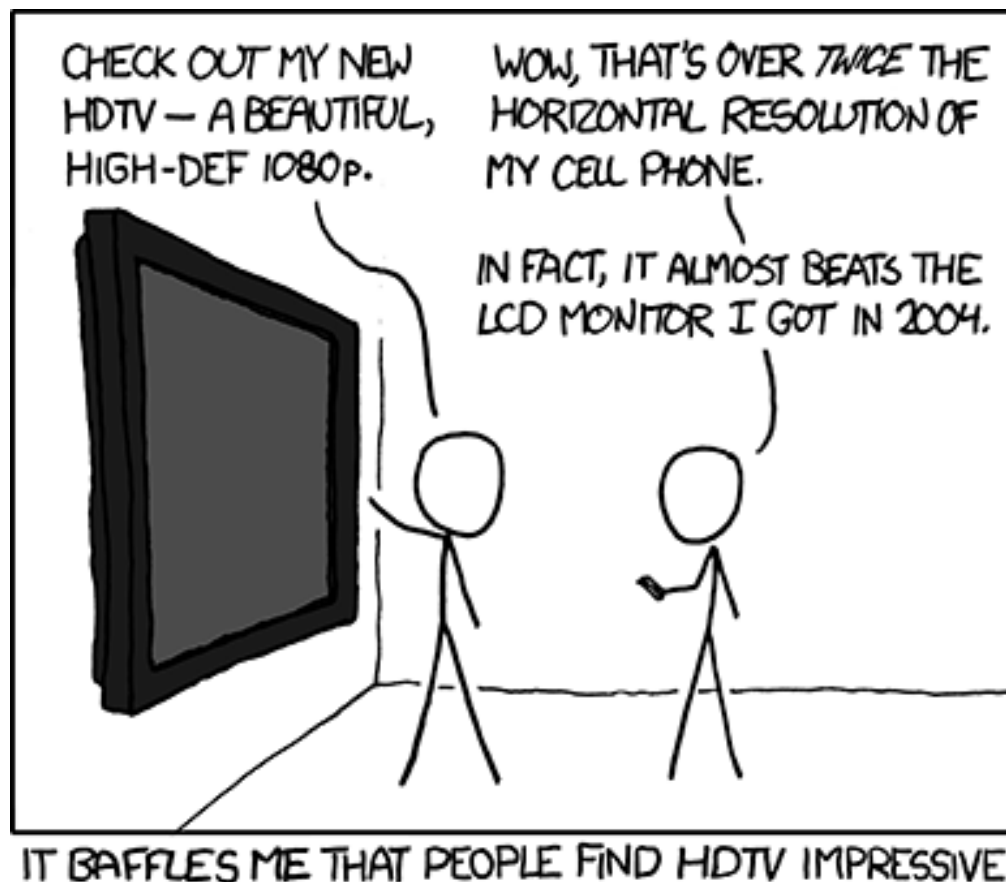
- **The Essential Guide To Video Processing.** Bovik. 2009.
  - Cap 9
- **H.264 and MPEG-4 Video Compression.** Richardson. 2003.
  - Cap 3 y 7





# Papers

- Wolf. “Keyframe selection by motion analysis”. 1996.
- Fauvet et al. “A Geometrical Key-Frame Selection Method Exploiting Dominant Motion Estimation in Video”. 2004.
- Sun et al. “Content-based representative frame extraction for digital video”. 1998.
- Zhuang et al. “Adaptive Keyframe Extraction Using Unsupervised Clustering”. 1998.



*"We're also stuck with blurry, juddery, slow-panning 24fps movies forever because (thanks to 60fps home video) people associate high framerates with camcorders and cheap sitcoms, and thus think good framerates look 'fake'."*

<http://xkcd.com/732/>