

## IV.2 Video Transcoding

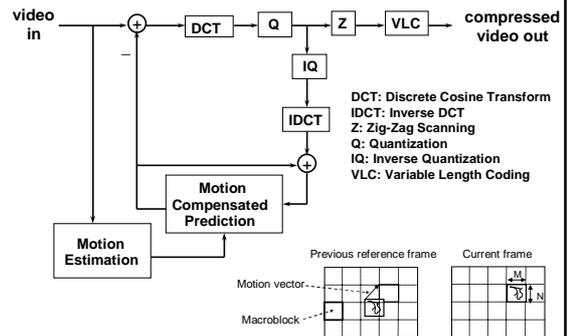
Prof. Chia-Wen Lin (林嘉文)  
 Department of CSIE  
 National Chung Cheng University  
 886-5-272-0411 ext. 33120  
<http://www.cs.ccu.edu.tw/~cwlin/>  
[cwlin@cs.ccu.edu.tw](mailto:cwlin@cs.ccu.edu.tw)

## Presentation Outline

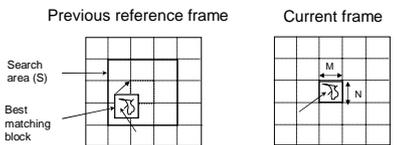
- Background
- Overview of Video Transcoding
  - Application Examples
  - Pixel-Domain Transcoder
  - DCT-Domain Transcoder
  - Fast Video Transcoding Architectures
- CCL's Strength in Video Transcoding
- Discussions

## Background

## Standard Video Encoder

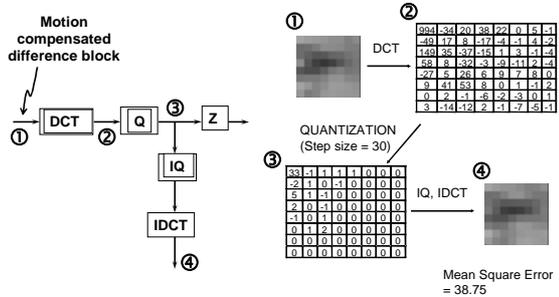


## Motion Estimation

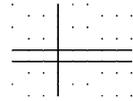


- Usually, based on Sum of Absolute Difference (SAD)

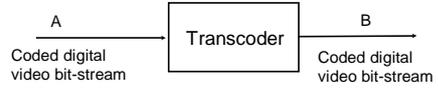
## DCT + Quantization



# Overview of Video Transcoding



## Introduction: Video Transcoding

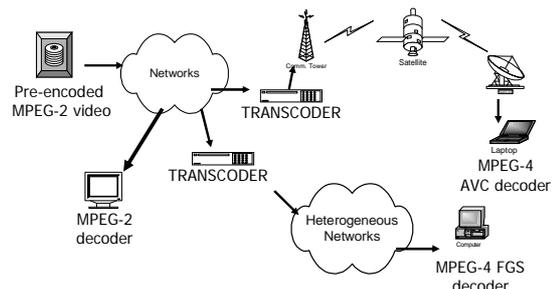


- Bit-rate/frame-rate/resolution conversion
- Format conversion
- Content manipulation

## Transcoder Applications

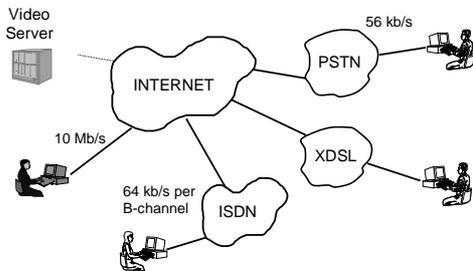
- Bit-rate/frame-rate/resolution adaptation
  - Video transport over heterogeneous networks
  - Multi-point video Conferencing
  - Statistical multiplexing ...
- Format conversion
  - Multi-standard terminal, HDTV to SDTV, MPEG-2 to H.264, MPEG-2 to MPEG-4 FGS, ...
- Content manipulation
  - Watermark/Logo insertion
  - Error resilience feature insertion
  - Editing/splicing ...

## Application Example: Video Delivery over Heterogeneous Networks



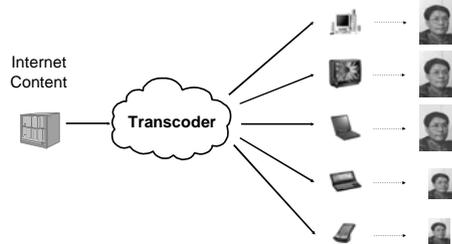
## Application Example: Video Delivery over Heterogeneous Networks

To deliver multimedia data over heterogeneous networks



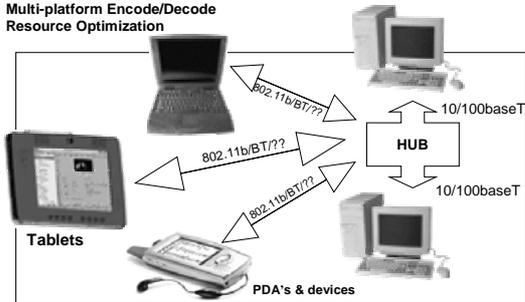
## Application Example: Video Delivery to Heterogeneous Clients

To deliver multimedia data to diverse devices with different capabilities (Universal Multimedia Access)



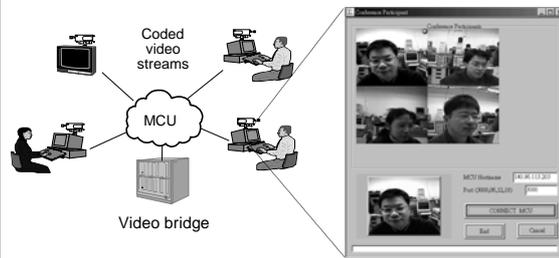
### Application Example: Seamless Home Networking

- Transcoding for bandwidth management
- Robust Wireless Streaming/transmission
- Multi-platform Encode/Decode
- Resource Optimization

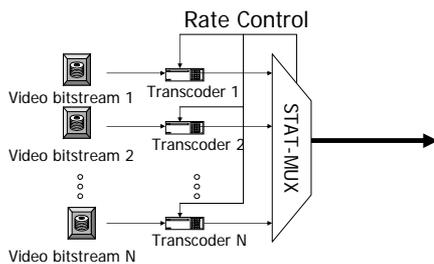


### Application Example: Dynamic Rate Control for Multipoint Video Conferencing

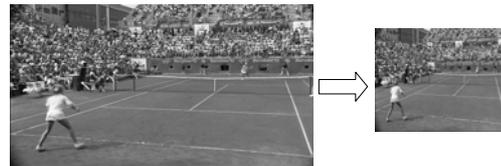
Example:  
 Input: 4 QCIF(176X144) videos @ 128 Kbps  
 Output: 1 CIF (352x288) videos @ 128 Kbps



### Application Example: Statistical Multiplexing



### Application Example: HDTV => SDTV Transcoding



### Application Example: Watemark/Logo Insertion



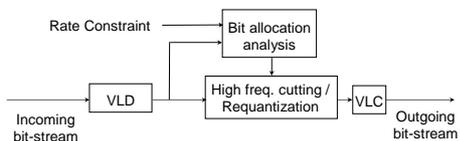
### Transcoding Research

- Issues
  - Complexity
  - Video quality
  - Flexibility
- Key
  - How to utilize available information from input compressed video – motion vectors, coding statistics, etc. for best video quality, lowest complexity, and highest flexibility?
  - A kind of two-pass encoding

# Transcoding Architectures

## Open-Loop Coded-Domain Transcoder

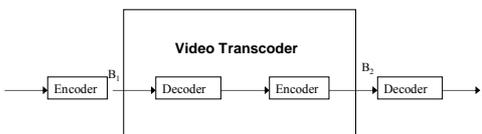
Truncation of high frequency or requantization of DCT coefficients



- No decoding and encoding, low complexity
- Significant drift between the front encoder and the end decoder

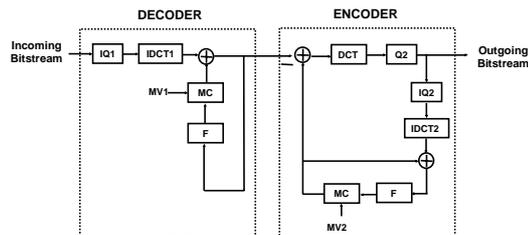
## Cascaded Pixel-Domain Transcoder

Direct implementation of Single Point Video Transcoder



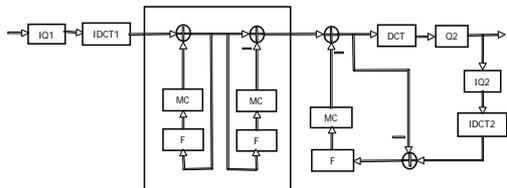
- Purpose: bit-rate/format/standard conversion
- In video telephony application complexity and quality are the most important factors.
- Information reused to speed up the re-encoding process
  - Header information re-use
  - Motion vectors re-use

## Cascaded Pixel-Domain Transcoder

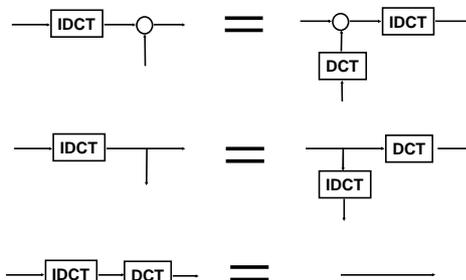


## Simplification of Transcoder (1)

• Remove one frame memory

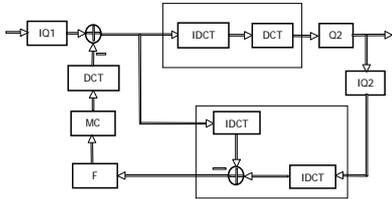


## Simplification of Transcoder (2)

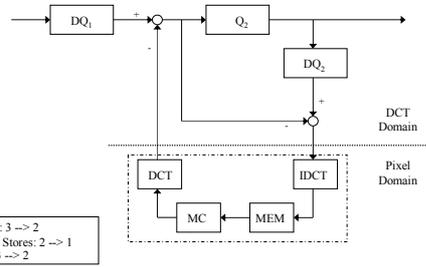


### Simplification of Transcoder (3)

- Re-arrange DCT/IDCT blocks

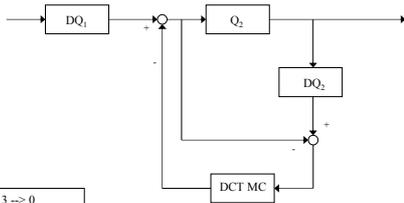


### Simplified Pixel-Domain Transcoder

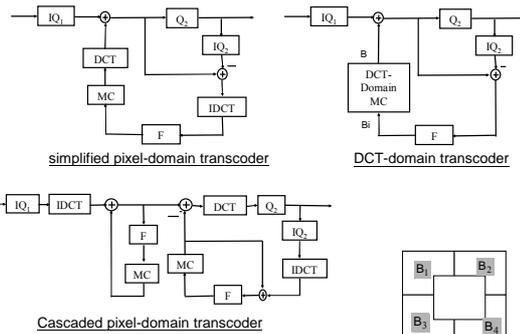


[G. Keesman et al]

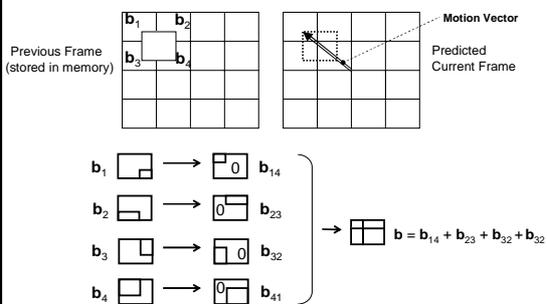
### Simplified DCT-Domain Transcoder



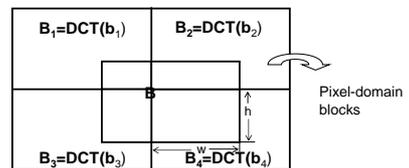
### Fast Transcoder Architectures



### Motion Compensation in the Pixel Domain



### DCT-Domain Motion Compensation (1)



- How to compute the DCT coefficients B from  $B_1, B_2, B_3, B_4$ ?

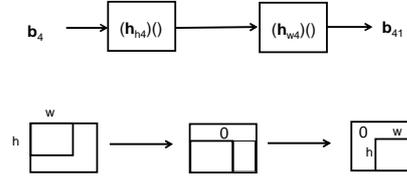
### DCT-Domain Motion Compensation (2)

- $\mathbf{b} = \mathbf{b}_{14} + \mathbf{b}_{23} + \mathbf{b}_{32} + \mathbf{b}_{41}$
- $\text{DCT}(\mathbf{b}) = \text{DCT}(\mathbf{b}_{14}) + \text{DCT}(\mathbf{b}_{23}) + \text{DCT}(\mathbf{b}_{32}) + \text{DCT}(\mathbf{b}_{41})$
- If we find the relationship between elements of  $(B_{14} - B_{41})$  and the elements  $(B_1 - B_4)$ , then we can calculate the  $B = \text{DCT}(\mathbf{b})$  directly from the elements  $(B_1 - B_4)$ .
- Consider the contribution from the original block  $B_4$

$$\mathbf{b}_{41} = \mathbf{h}_{h_4} \mathbf{b}_4 \mathbf{h}_{w_4}, \text{ where } \mathbf{h}_{h_4} = \begin{bmatrix} 0 & 0 \\ I_{h_4} & 0 \end{bmatrix}, \mathbf{h}_{w_4} = \begin{bmatrix} 0 & I_{w_4} \\ 0 & 0 \end{bmatrix}$$

- $I_{h_4}$  and  $I_{w_4}$  are identity matrices of size  $h_4 \times h_4$  and  $w_4 \times w_4$  respectively.

### DCT-Domain Motion Compensation (3)



### DCT-Domain Motion Compensation (3)

- All unitary orthogonal transforms such as DCT are distributive to matrix multiplication, i.e.,  $\text{DCT}(AB) = \text{DCT}(A)\text{DCT}(B)$
- The elements of  $B_{41}$  can be computed directly from the elements of  $B_4$  through geometric transforms

$$B_{41} = H_{h_4} B_4 H_{w_4}$$

- The DCT coefficients of block  $B$  can be computed using

$$B = \sum_{i=1}^4 H_{h_i} B_i H_{w_i}$$

- and can be pre-computed and stored in memory.

### DCT-Domain Motion Compensation (4)

| Neighboring blocks | Position    | $\mathbf{h}_{h_i}$                                   | $\mathbf{h}_{w_i}$                                   |
|--------------------|-------------|--|--|
| $\mathbf{b}_1$     | Lower right | $\begin{bmatrix} 0 & I_{h-h} \\ 0 & 0 \end{bmatrix}$ | $\begin{bmatrix} 0 & 0 \\ I_{w-w} & 0 \end{bmatrix}$ |
| $\mathbf{b}_2$     | Lower left  | $\begin{bmatrix} 0 & I_{h-h} \\ 0 & 0 \end{bmatrix}$ | $\begin{bmatrix} 0 & I_w \\ 0 & 0 \end{bmatrix}$     |
| $\mathbf{b}_3$     | Upper right | $\begin{bmatrix} 0 & 0 \\ I_h & 0 \end{bmatrix}$     | $\begin{bmatrix} 0 & 0 \\ I_{w-w} & 0 \end{bmatrix}$ |
| $\mathbf{b}_4$     | Upper left  | $\begin{bmatrix} 0 & 0 \\ I_h & 0 \end{bmatrix}$     | $\begin{bmatrix} 0 & I_w \\ 0 & 0 \end{bmatrix}$     |

### DCT-Domain MC Example

- Compute pixel values from neighboring blocks:



- Example:

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| 1 | 9 | 3 | 9 | 8 | 7 |
| 6 | 4 | 5 | 3 | 9 | 6 |
| 4 | 3 | 8 | 7 | 1 | 2 |
| 5 | 4 | 3 | 2 | 1 | 0 |
| 9 | 8 | 7 | 6 | 5 | 4 |
| 3 | 2 | 1 | 5 | 9 | 1 |

|   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 2 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |   |   |
| 1 | 0 | 0 | 6 | 5 | 4 | 0 | 0 | 1 | 0 | 0 | 2 | 1 |   |
| 0 | 1 | 0 | 5 | 9 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 6 | 5 |
| 0 | 0 | 1 | 9 | 8 | 7 | 0 | 1 | 0 | 1 | 0 | 0 | 7 | 1 |
| 0 | 0 | 0 | 3 | 9 | 6 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 7 | 1 | 2 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 | 9 | 3 | 0 | 0 | 0 | 0 | 8 | 0 | 0 | 0 |
| 0 | 0 | 0 | 6 | 4 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 4 | 3 | 8 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |

[S.F.Chang et al]

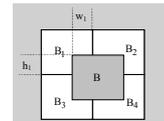
### Simplified DCT-Domain Transcoder

$$B = \sum_{i=1}^4 H_{h_i} B_i H_{w_i}$$

$H_{h_i}, H_{w_i}$  : geometric transform matrices

$w_i, h_i \in \{1, 2, \dots, 7\}$

$H_{h_1} = H_{h_2}, H_{h_3} = H_{h_4}, H_{w_1} = H_{w_3}, H_{w_2} = H_{w_4}$





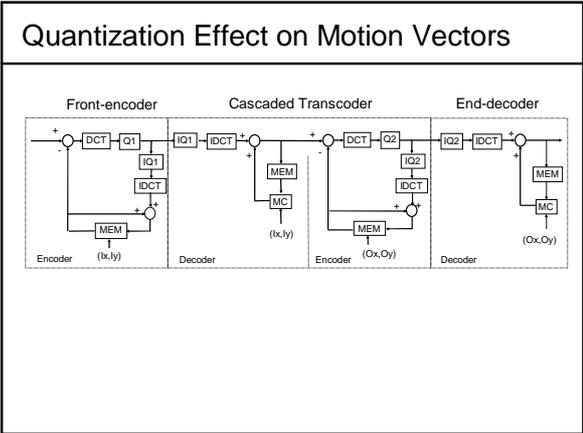
## Research Examples

- ### Research Examples
- Adaptive Motion Vector Refinement for High-Performance Video Transcoding
  - Dynamic Region-of-Interest Transcoding for Multipoint Video Conferencing
  - Error Resilience Transcoding for Video Streaming over WLAN
  - Key-frame Extraction/Transcoding for Channel-Aware Video Streaming

### Motion Vector Refinement for High-Performance Transcoding

IEEE Trans. Multimedia 1999

Prof. Chia-Wen Lin  
Department of CS  
National Chung Cheng University  
886-5-272-0411 ext. 33120  
cwlin@cs.ccu.edu.tw



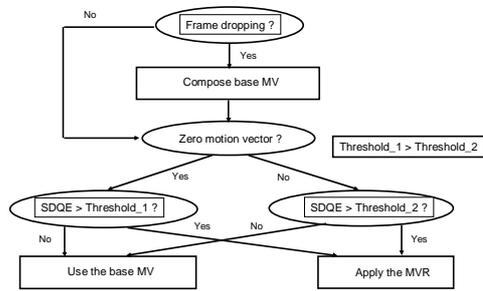
- ### Motion Vector Refinement
- Reduce the computation while keeping the high performance
  - Using base and delta motion vector
- 
- New search area      can be much smaller than the original search area (e.g.      vs.      pixels)

### Adaptive Motion Vector Refinement (1)

- **Sum of Differential Quantization Error (SDQE)**

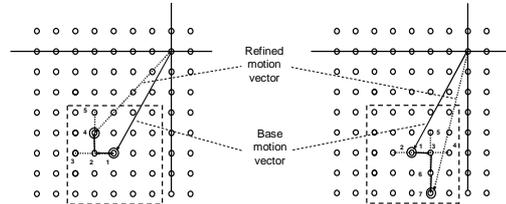
where  $q_1$  = quan. step-size extracted from incoming current frame  
 $q_2$  = quan. step-size used in the previous frame in the transcoder

## Adaptive Motion Vector Refinement (2)



## Fast Search Algorithms

### Horizontal And Vertical Search (HAV)



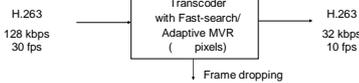
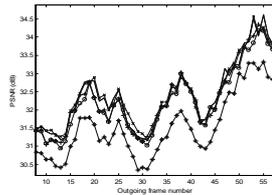
Best case (5 checking points)

Worst case (7 checking points)

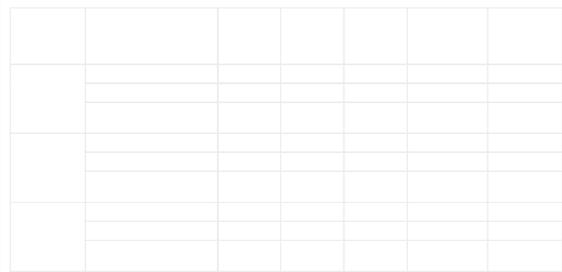
## Performance Comparison (1)

### "Carphone" Sequence

- \* : ONLY BASE MVs
- + : BBGDS+ADAP+MVR
- o : HAVS+ADAP+MVR
- x : FULL-MVR



## Performance Comparison (2)



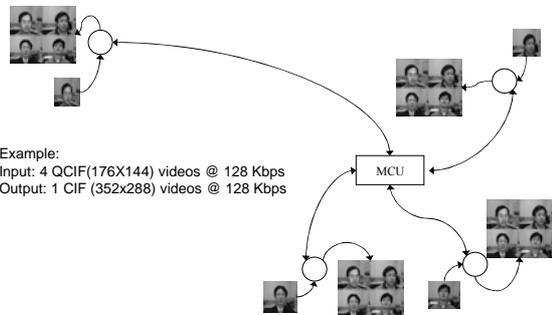
(Carphone Sequence, 128 kbps to 32 kbps)

## Dynamic Region-of-Interest Transcoding for Multipoint Video Conferencing

Presented at ISCAS 2000  
IEEE Trans. CSVT 2003

Prof. Chia-Wen Lin  
Department of CS  
National Chung Cheng University  
886-5-272-0411 ext. 33120  
cwlin@cs.ccu.edu.tw

## Multipoint Video Conferencing



## Motivation

- Observation
  - In a multipoint video conference, in general, only one or two conferees are active at one time.
- **Considerations in this work:**
  - Video quality enhancement
    - Sub-windows of high activity requires higher frame rate to maintain acceptable temporal resolution (motion smoothness), while the frame rate for the inactive sub-windows can be lower.
    - The bit-rate saved from skipping the inactive sub-windows can be used to enhance the quality of the active ones.
  - Computation reduction
    - Significant computation reduction can be achieved by sub-window skipping.
- Question: What frames should or should not be skipped?

## Dynamic Sub-Window Skipping

- Sub-window skipping criteria

$$\text{if } (S_m^{MV} < TH_{SAD}) \&\& \left( \frac{SAD_m - SAD_m^{ref}}{SAD_m^{ref}} < TH_{SAD} \right)$$

then

Skip Current Sub-window

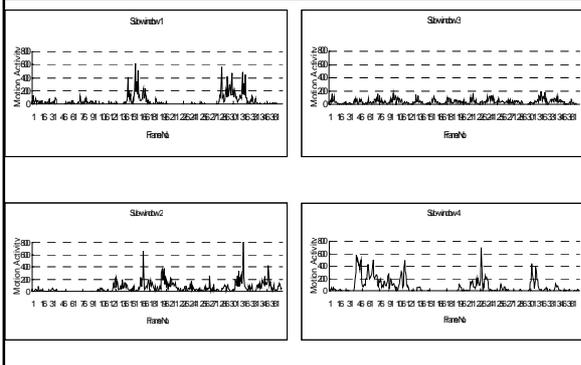
else

Encode Current Sub-window

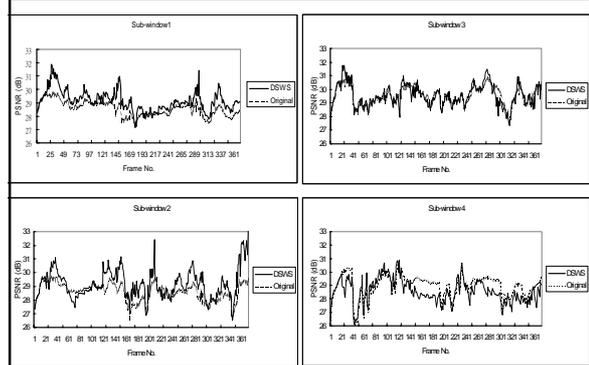
$$S_m^{MV} = \sum_{n=1}^N |MV_{m,n}^x| + |MV_{m,n}^y| \quad SAD_{m,n} = \sum_{n=1}^N \sum_{x,y \in MB_{m,n}} |f_m(x,y) - f_m^{ref}(x+MV_{m,n}^x, y+MV_{m,n}^y)|$$

- The static sub-windows are skipped (distortion is relatively invisible) and the saved bits are used to enhance other sub-windows
- Computation complexity:
  - The MV's are decoded from the incoming bitstreams, thus only the summation operation is required.
  - The SAD computation cost is low
  - Sub-window skipping can achieve significant computation saving

## Sub-window Activities



## Performance Comparison (1)



## Performance Comparison (1)

|                | Skipped frame No. | Average PSNR of all frames (dB) |       | Average PSNR of non-skipped frames (dB) |       |       |       |
|----------------|-------------------|---------------------------------|-------|---|-------|-------|-------|
|                |                   | Original                        | DSWS  | Original                                | DSWS  |       |       |
| Sub-sequence 1 | 151               | 28.54                           | 29.14 | +0.60                                   | 28.55 | 29.31 | +0.76 |
| Sub-sequence 2 | 75                | 28.59                           | 29.16 | +0.57                                   | 28.52 | 29.25 | +0.73 |
| Sub-sequence 3 | 54                | 29.54                           | 29.56 | +0.02                                   | 29.48 | 29.59 | +0.11 |
| Sub-sequence 4 | 139               | 28.99                           | 28.59 | -0.40                                   | 28.73 | 28.68 | -0.05 |
| Average        | 104.75            | 28.91                           | 29.11 | +0.20                                   | 28.82 | 29.21 | +0.39 |

Skipped frame number (total :400 frames) and average PSNR of four sub-windows ( $TH_1=5; TH_2=0.1$ )

**About 25% computation reduction is achieved**

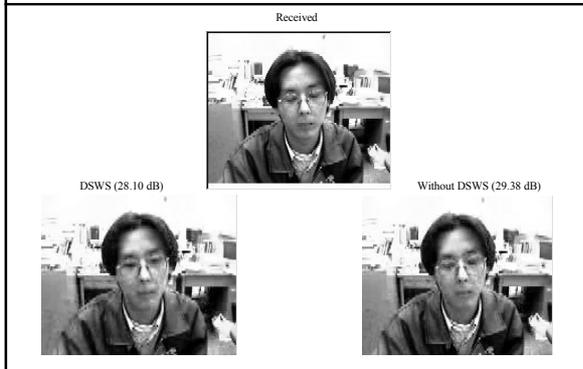
## Performance Comparison (2)



Without DSWS 29.45 dB

With DSWS 30.87 dB

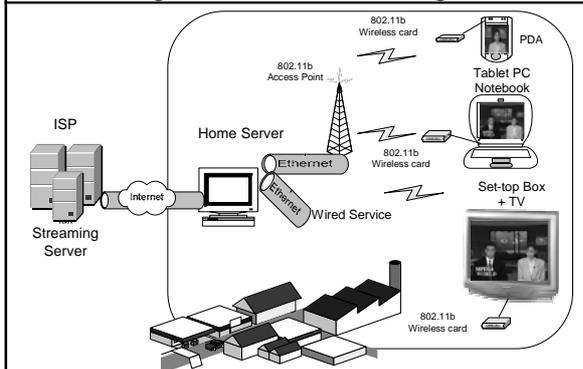
### Performance Comparison (3)



### Video Transcoding for Three-Tier Video streaming

Prof. Chia-Wen Lin  
 Department of CS  
 National Chung Cheng University  
 886-5-272-0411 ext. 33120  
 cwlin@cs.ccu.edu.tw

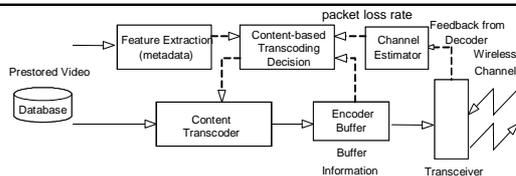
### Application Scenario: Three-Tier Video Streaming for Home Networking



### Video Adaptation for Three-Tier Video streaming

- WLANs are packet erasure channels
  - Burst errors due to slow-fading channels
  - Bit-level error-resilient coding tools and Bit/byte-level FEC are usually not useful
  - MAC-layer retransmission may degrade the quality seriously
- Transcoding is particularly useful for content adaptation in such application scenario
  - The number of home clients are usually small => reasonable complexity
  - The home server can easily capture the channel statistics through feedback information
  - Allow two-pass processing

### Content-Aware Video Adaptation

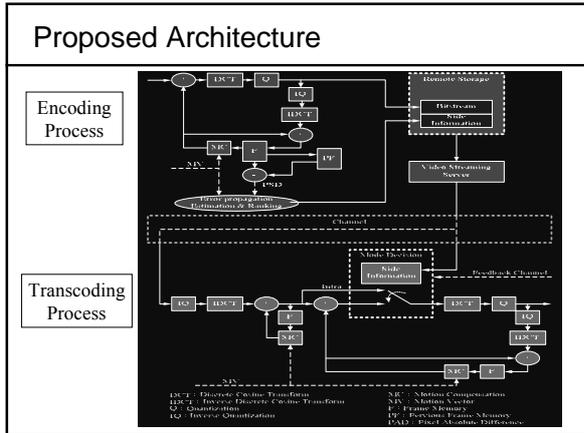


- Our Approach
  - The server offline extracts the features from the pre-encoded bit-stream
  - The features are then stored as auxiliary data to guide the real-time transcoding
  - Intra-refresh & ARQ are adopted as the tools for content adaptation because they doesn't need to change the decoder

### Error Resilience Transcoding for Video Streaming over WLAN

Presented at ISCAS 2004  
 J. Visual Commun. & Image Rep. 2005

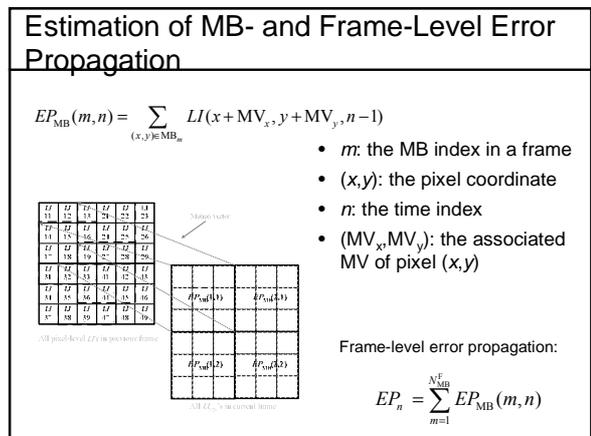
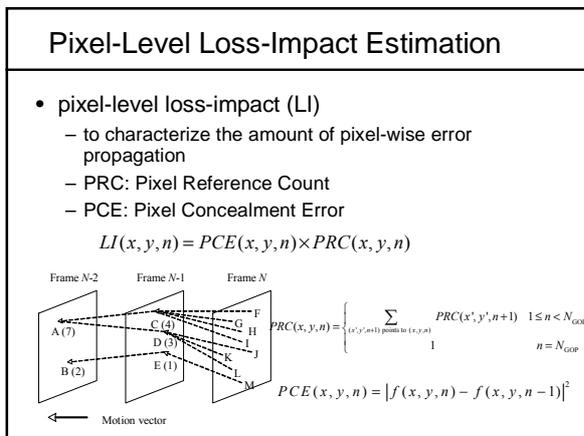
Prof. Chia-Wen Lin  
 Department of CS  
 National Chung Cheng University  
 886-5-272-0411 ext. 33120  
 cwlin@cs.ccu.edu.tw



- ### Two-Pass Error-Resilience Transcoding
- First pass – Encoding process
    - exploits the motion vectors and concealment distortion to evaluate the rank of the loss impact of each MB in a frame, each frame in a GOP
    - the rank serves as side information or transcoding hints to guide the second-pass (transcoding pass) intra-refresh decisions

- ### Two-Pass Error-Resilience Transcoding
- Second-pass – Transcoding process
    - collect channel statistics from the client terminals to estimate the channel packet loss-rate
    - use side information and feedback channel condition to determine the intra-refresh rate
    - perform intra-refresh on the highest-priority MBs (with top-ranked loss impact)

- ### Loss-Impact Estimation in First-Pass Encoding
- Involves three steps
    - Calculating the Pixel-level Loss Impact
    - Calculating the MB-level Error-Propagation
    - Calculating the Frame-level Error Propagation



### Transcoding with Adaptive Intra-Refresh: GOP-Level Allocation

- The number of intra-refreshed MBs in a GOP

$$N_{intra}^{GOP} = \frac{1}{N_{GOP}} \sum_{n=1}^{N_{GOP}} EP_n \times PLR$$

$$TH_{intra}$$

- $N_{GOP}$ : number of frames in a GOP
- $PLR$ : packet-loss rate
- $TH_{intra}$ : scaling parameter to characterize the relationship between  $N_{intra}^{GOP}$  and the error propagation effect in a GOP

### Transcoding with Adaptive Intra-Refresh: Frame-Level Allocation

- if  $n = 1$

$$N_{intra}^n = \frac{EP_n}{\sum_{i=1}^{N_{GOP}} EP_i} \times N_{intra}^{GOP}$$

- else if  $2 \leq n \leq N_{GOP}$

$$N_{intra}^n = \frac{EP_n}{\sum_{i=n}^{N_{GOP}} EP_i} \times \left( N_{intra}^{GOP} - \sum_{i=1}^{n-1} N_{intra}^i \right)$$

- endif

- if  $N_{intra}^n > k_{MB} N_{MB}^F$  then  $N_{intra}^n = k_{MB} N_{MB}^F$

Notation:

$n$ : frame index in a GOP

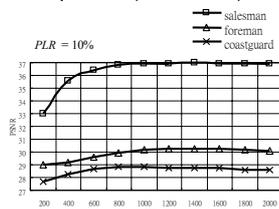
$N_{intra}^n$ : number of intra-coded MBs in frame  $n$

$N_{MB}^f$ : number of MBs in a frame

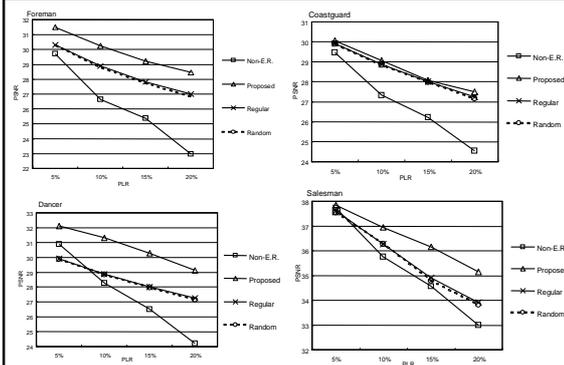
$k_{MB}$ : to constrain the number of intra-coded blocks in a frame not to exceed an upper limit ( $0 \leq k_{MB} \leq 1$ )

### Threshold Selection

- The parameter  $TH_{intra}$  is obtained empirically
  - the following figure shows the frame-by-frame PSNR with different  $TH_{intra}$ 's
  - the best PSNR performance is achieved at  $TH_{intra} \cong 1200$  for the three sequences ( $PLR = 10\%$ )



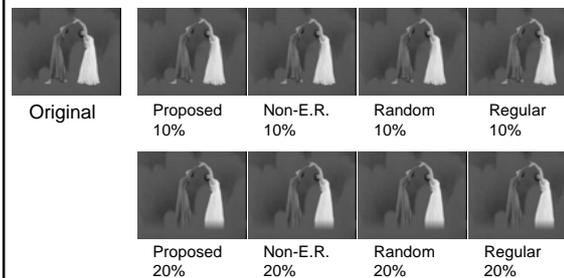
### Performance Comparison



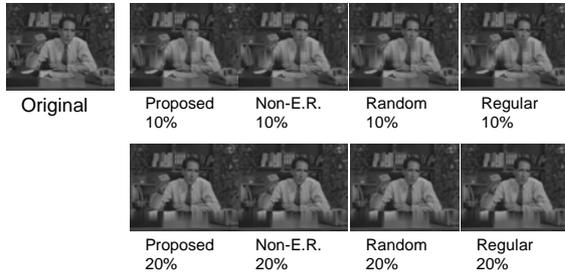
### Demo: Foreman (PLR = 10% & 20%)



### Demo: Dancer (PLR = 10% & 20%)



## Demo: Salesman ( $PLR = 10\% \& 20\%$ )



## Summary

- We proposed a novel two-pass prioritized intra-refresh strategy for error-resilience transcoding
- The extra computational complexity required for the transcoding process is almost negligible, thereby making it suitable for real-time applications
- The proposed algorithm can effectively mitigate the error propagation due to packet-loss
- The degree of error resilience can be dynamically adjusted to adapt to a channel with time-varying error characteristics

## Key-frame Extraction/Transcoding for Channel-Aware Video Streaming

Presented at ICIP 2004

Prof. Chia-Wen Lin  
Department of CS  
National Chung Cheng University  
886-5-272-0411 ext. 33120  
cwin@cs.ccu.edu.tw

## Motivation

- Key-frames form a compact representation of a video
  - Video summarization
  - Efficient video streaming and video retrieval
- Issues about generating a summary video:
  - How many key-frames are enough?
  - Which frames should be selected?
- We present an adaptive rate-constrained key-frame selection scheme for channel-aware realtime video streaming applications.

## Channel Rate Adaptation

- The features and the estimated channel conditions are used to determine the output bit-rate and the key-frame selection policy.
- After determining the output bit-rate, the target number of key-frames,  $N_{KF}$  with rate constraint  $R_{out}$  is obtained as

$$N_{KF} = \left\lfloor k_{adj} \frac{N_{total} R_{out}}{R_{in}} \right\rfloor$$

$R_{in}$ : bit-rate of the incoming video

$R_{out}$ : estimated bit-rate of the output video summary

$N_{total}$ : length of the video clip

## Key-Frame Extraction for Rate Adaptation

- R-D Optimized Shot-Level Key-Frame Allocation
  - The problem of key-frame allocation is to distribute  $N_{KF}$  key-frames into all shots.
  - Since different shots may have different characteristics (e.g., motion and texture), simply uniformly distributing the key-frames will result in non-uniform representation quality.
  - We propose an R-D optimized shot-level key-frame allocation scheme.

### Compressed-Domain Feature Extraction & Distance Metrics

- Given two finite point sets  $A = \{a_1, \dots, a_p\}$  and  $B = \{b_1, \dots, b_q\}$ , the Hausdorff distance is

$$H(A, B) = \max(h(A, B), h(B, A))$$

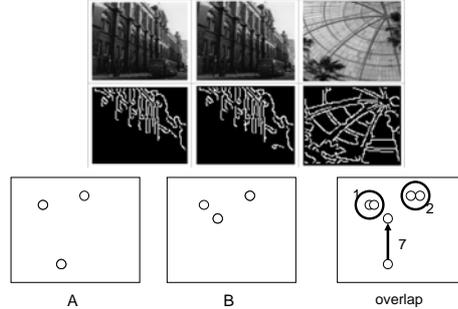
where  $h(A, B) = \max_{a \in A} \min_{b \in B} \|a - b\|$

- Because the Hausdorff distance is sensitive to noise, a modified Hausdorff distance metric is proposed

$$h_K(A, B) = K^{\text{th}} \min_{a \in A} \min_{b \in B} \|a - b\|$$

### Compressed-Domain Feature Extraction & Distance Metrics (Cont.)

Hausdorff distance of Canny edge images:



### Compressed-Domain Feature Extraction & Distance Metrics (Cont.)

- Extract the DC image and then apply the Canny edge operator to extract the Canny edgy images.
- The spatial distance between the  $i$ th and  $j$ th edge images :

$$d^S(f_i, f_j) = \sum_{n=i}^{j-1} h_K(f_n, f_{n+1})$$

- To take into account the object and camera motions, we define the temporal distance as the cumulative sum of motion vector magnitudes.

$$d^T(f_i, f_j) = \sum_{n=i}^{j-1} \sum_{m=1}^{N_{MB}} |MVX_m^n| + |MVY_m^n|$$

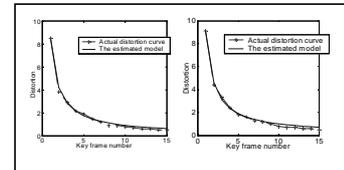
- Spatiotemporal distance

$$d(f_i, f_j) = k_1 d^T(f_i, f_j) + k_2 d^S(f_i, f_j)$$

### Shot-Level Rate-Distortion Modeling

- Assume that a video clip has  $N$ shot shots, and the  $i$ -th shot has its rate-distortion function  $D_i(N_i)$  to characterize the relationship between the average representation distortion ( $D_i$ ) and the number of key-frames ( $N_i$ ).
- We found the following first-order model is a good approximation of  $D_i(N_i)$ :

$$D_i(N_i) = \frac{a_i}{N_i} + b_i$$



### Optimal Shot-Level Key-frame Allocation

- The optimal key-frame allocation problem can be formulated as the following constrained optimization problem:

$$\min \sum_{i=1}^{N_{\text{shot}}} L_i D_i(N_i) \quad \text{subject to} \quad \sum_{i=1}^{N_{\text{shot}}} N_i = N_{\text{KF}}$$

- Then use the Lagrange multiplier to convert the above Eq. to the following unconstrained optimization problem

$$f = \sum_{i=1}^{N_{\text{shot}}} L_i \left( \frac{a_i}{N_i} + b_i \right) + \lambda \left( \sum_{i=1}^{N_{\text{shot}}} N_i - N_{\text{KF}} \right)$$

- By setting partial derivatives to zero (i.e.,  $\partial f / \partial N_i = 0$  and  $\partial f / \partial \lambda = 0$ ),

$$\lambda = \frac{\sum_{i=1}^{N_{\text{shot}}} \sqrt{L_i a_i}}{N_{\text{KF}}^2} \quad \text{and} \quad N_i = \sqrt{\frac{L_i a_i}{\lambda}}$$

### Experimental Results

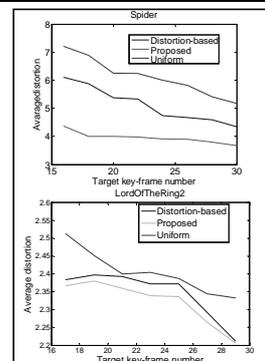
- Two movie clips consisting of 1132 and 1123 frames respectively are encoded at 30 fps and 1Mbps using an MPEG-4 encoder.
- Two other allocation schemes are compared

- uniform allocation

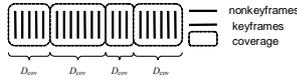
$$N_i = N_{\text{KF}} \times L_i / \sum_{i=1}^{N_{\text{shot}}} L_i$$

- distortion-weighted allocation

$$N_i = N_{\text{KF}} \times L_i D_i / \sum_{i=1}^{N_{\text{shot}}} L_i D_i$$



## Optimal Shot-Level Key-Frame Allocation



– Key-frame selection within a shot

- **Step 1.** Compute the average coverage range of each key-frame

$$D_{cov} = \frac{1}{N} \sum_{n=1}^{M-1} d(f_n, f_{n+1})$$

- Divide the shot into  $N$  intervals with approximately equal average distortion

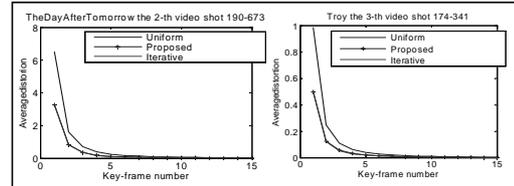
$$\sum_{n=i}^{i+1} d(f_n, f_{n+1}) \equiv D_{cov} \quad j = 1, 2, \dots, N$$

- **Step 2.** For each interval  $T_i$ , the key-frame is selected so as to minimize the distortion  $D(T_i, k_i)$ :

$$k_i = \arg \min_{k_i \in S_i} D(T_i, k_i), \quad i = 1, 2, \dots, N$$

## Result of Key-frame Extraction from a Shot

- The average distortion between frames and the corresponding key-frames is used for performance evaluation
- The proposed non-iterative approach can achieve close quality with significantly lower complexity as compared to the iterative approach



## Demo: Key-frame Extraction

Movie,  
1132 frames,  
70 key  
frames



Movie,  
1123 frames,  
55 key  
frames



## Demo: Key-frame Extraction

1132 frames, 20 target key-frame

Proposed

Uniform

Distortion-weighted



## Summary

- We proposed an adaptive sequential key-frame selection method for channel-aware video streaming applications
- R-D models are developed to characterize the key-frame representations of video contents
  - Based on offline extracted metadata
  - Non-iterative R-D constrained key-frame allocation & selection
  - Achieve lowest representation distortion
- The key-frames selection can be dynamically adjusted to adapt to a channel with time-varying characteristics