

CSE - Conception des systèmes embarqués – 8

Multimedia : Audio / Video

Michel Starkier



CSE

ENCODAGE ET TRANSPORT DE L'AUDIO NUMÉRIQUE

25/05/2014

Conception systèmes embarqués / MSR



2

Bases d'acoustique

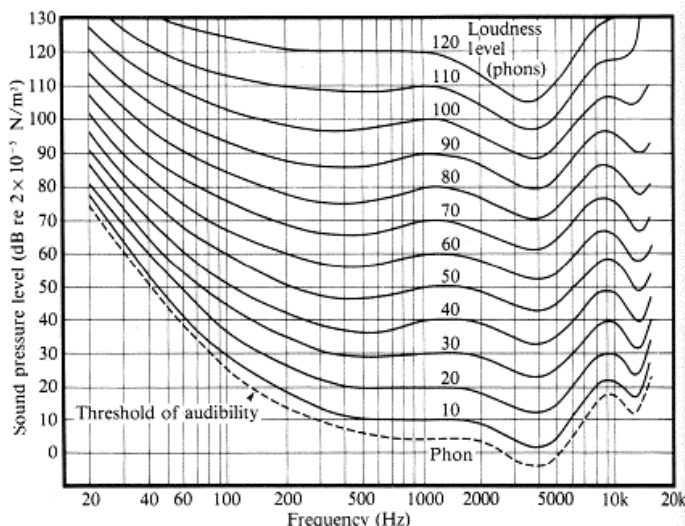
- Le son provient d'un objet qui entre en vibration mécanique et cause une variation de la pression atmosphérique dans l'air.
- Vitesse du son dans l'air : 340 m/s
- Niveau sonore mesuré en dB SPL
 - $S_{dB} = 20 \log_{10}(P_{eff}/P_{ref}) \text{ dB}$ => rapport de pression, P_{ref} est la pression de référence de 20 μPa (micro Pascal)
- Dynamique: 0dB (seuil d'audition) à 120dB (seuil de la douleur)
- Fréquence : 20Hz à 20KHz

25/05/2014

Conception systèmes embarqués / MSR

3

La perception du niveau sonore dépend de la fréquence. La courbe de Fletcher et Munson indique les lignes ou les niveaux sonores perçus sont identiques dans la plage des fréquences audibles.



Audio numérique

● Fréquence d'échantillonnage (non compressé) :

- Radio : 22/ 32 KHz / 16 bits
- CD : 44.1 KHz / 16 bits
- Audio pro : 48 KHz /16 ou 24 bits
- Audio HQ : 96 KHz / 24 bits
- Rapport signal /bruit > 90 dB

● Théorème de Shannon

- fréquence d'échantillonnage d'un signal doit être égale ou supérieure au double de la fréquence maximale

25/05/2014

Conception systèmes embarqués / MSR

5

Une sinusoïde de fréquence F_0 est échantillonnée à une fréquence F_s , la période d'échantillonnage est : $T = \frac{1}{F_s}$

$x[n] = \sin(2\pi F_0 t_s)$ avec $t_s = nT$ (n est un entier)

Nous pouvons écrire :

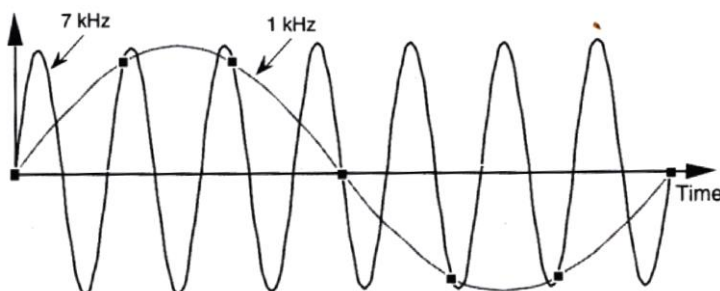
$x[n] = \sin\left(2\pi F_0 nT + 2\pi k n \frac{T}{T}\right)$, avec k entier, car le deuxième terme prends les valeurs $2\pi, 4\pi, 6\pi \dots$

d'où :

$x[n] = \sin(2\pi F_0 nT + 2\pi k n F_s T) = \sin(2\pi (F_0 + k F_s) nT)$

Avec un sampling rate de F_s samples/s, et k un entier (positif ou négatif), nous ne pouvons pas distinguer les valeurs échantillonnées d'une sinusoïde de fréquence F_0 Hz de celles d'une sinusoïde de fréquence $(F_0 + F_s)$

$F_s = 6 \text{ kHz}$



Compression lossless (sans perte)

- **Le plus simple : codage différentiel**
 - **Plus compliqué : prédiction linéaire**
 - **Run lenght (encodage des silences ...)**
 - **Compression 50 %**
-
- **Free lossless audio codec**
 - **Shorten**
 - **ALE (Apple)**



Le codage différentiel s'appuie sur le fait que les échantillons successifs sont fortement corrélés. Le codage différentiel, le plus simple consiste à encoder seulement la différence entre deux échantillons successifs. La prédiction linéaire, plus compliquée, consiste à prédire l'échantillon suivant et de ne coder que l'erreur entre l'échantillon prédit et l'échantillon réel suivant.

Liste non-exhaustive de codec lossless (Wikipedia)

Apple Lossless – ALAC (Apple Lossless Audio Codec)
 apt-X Lossless
 ATRAC Advanced Lossless
 Direct Stream Transfer – DST
 Dolby TrueHD
 DTS-HD Master Audio
 Free Lossless Audio Codec – FLAC
 Meridian Lossless Packing – MLP
 Monkey's Audio – Monkey's Audio APE
 La famille des MPEG
 OptimFROG
 RealPlayer – RealAudio Lossless
 Shorten – SHN
 TTA – True Audio Lossless
 WavPack – WavPack lossless
 WMA Lossless – Windows Media Lossless

Le streaming audio

Standards réseau

- **RTP Real-time Transport Protocol**
 - Protocole de transport
 - Voix sur IP, vidéo conférence
- **RTSP Real Time Streaming Protocol**
 - Couche application

Autres protocoles de streaming

- **WMA Windows Media Audio**
- **RTMP Real Time Messaging Protocol (Adobe)**
- **MP3**
- **Ogg Vorbis**
- **RealAudio**

Streaming Media With Linux by [Dave Phillips](#)

http://linuxdevcenter.com/pub/a/linux/2001/03/23/streaming_media.html



SILICON LABORATORIES
USB AUDIO CLASS TUTORIAL

<http://www.silabs.com/Support%20Documents/TechnicalDocs/AN295.pdf>

OpenSL ES - The Standard for Embedded Audio Acceleration

<http://www.khronos.org/opensles/>



CSE

ENCODAGE ET TRANSPORT DE LA VIDÉO NUMÉRIQUE

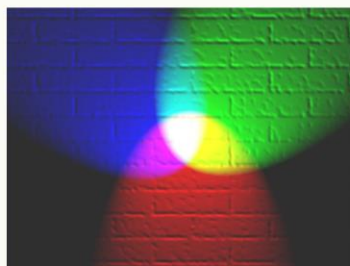
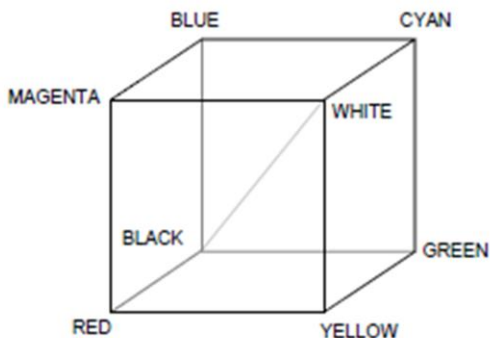
25/05/2014

Conception systèmes embarqués / MSR

13

Vidéo composante RVB

- Espace des couleurs RVB
- Addition des couleurs
- Codage 8 à 16 bits par couleur
- 24, 30, ou 48 bits par pixel



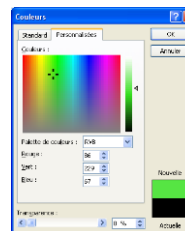
s embarqués / MSR

Les systèmes vidéo utilisent une synthèse additive des couleurs à partir de 3 couleurs primaires : rouge, vert et bleu.

- jaune = rouge + vert,
- violet (magenta) = bleu + rouge,
- bleu turquoise (cyan) = bleu + vert .

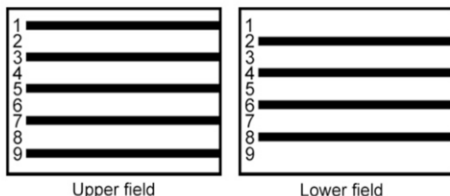
L'ensemble des couleurs que l'on peut obtenir par synthèse est représenté sur 3 axes par un cube, l'espace des couleurs.. Notez que rouge + vert + bleu = blanc. La quantification des couleurs se fait sur 8 à 16 bits, soit un mot de 24 à 48 bits par pixel (16 millions de couleurs en 24 bits).

Essayez de vous déplacer dans l'espace de couleur RVB de la fenêtre "Autres couleurs de.../Personnalisées" d'un programme Office. Observez les valeurs RVB correspondant à votre position.

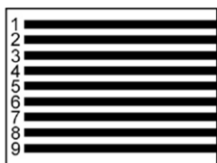


Vidéo entrelacée et progressive

● Entrelacement => 2 champs



● Progressif



25/05/2014

Conception systèmes embarqués / MSR

Les systèmes vidéo utilisaient un mécanisme appelé entrelacement pour limiter le nombre d'images par seconde. Les images diffusées par tube cathodique étaient balayées, c'est-à-dire diffusées point après point. L'entrelacement et le balayage tendent à disparaître avec l'utilisation des écrans LCD. Le progressif s'impose actuellement.

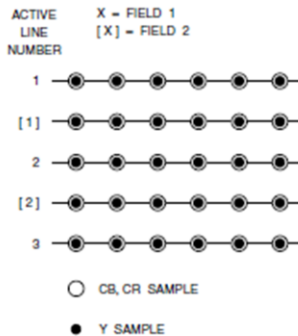
Entrelacement

Les lignes impaires de l'image sont d'abord balayées (trame impaire), puis les lignes paires sont balayées (trame paire). Ce mécanisme évite le scintillement en donnant l'impression que le nombre d'image est doublé, par exemple 25 images/s correspond à 50 trames/s.

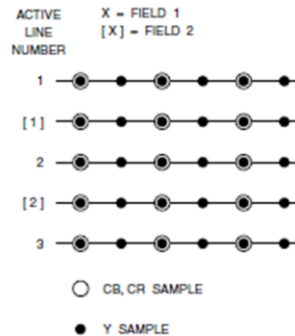
Progressif

L'image est diffusée en une seule fois. Le nombre d'images par seconde doit être au moins 50 i/s pour un confort visuel acceptable.

Echantillonnage



4:4:4



4:2:2

également 4:1:1 et 4:1:0

Chaque pixel de l'image est défini comme un échantillon du signal vidéo. Les valeurs d'échantillon sont transmises ligne par ligne de haut gauche vers bas droit de l'image.

4:4:4 : La fréquence d'échantillonnage de la luminance et celle de la chrominance sont identiques. Autant de pixels chrominance que de pixels luminance. Utilisé pour le RVB

4:2:2 : La fréquence d'échantillonnage de la luminance est deux fois plus élevée que la fréquence d'échantillonnage horizontal de la chrominance (un pixel couleur sur deux). Utilisé pour le YUV.

4:1:1 : La fréquence d'échantillonnage de la luminance est quatre fois plus élevée que la fréquence d'échantillonnage horizontal de la chrominance (un pixel couleur pour quatre pixels luminance)

4:2:0 : La fréquence d'échantillonnage de la luminance est deux fois plus élevée que la fréquence d'échantillonnage horizontal et que la fréquence d'échantillonnage vertical de la chrominance (un pixel couleur sur deux, une ligne couleur sur deux)

En 4:2:2 par exemple, le signal vidéo est encodé avec 4 mots (8 ou 19 bits) par pixel : Cb /Y0/Cr/Y1 => chrominance (bleu)/ luminance/chrominance/luminance

Compression vidéo

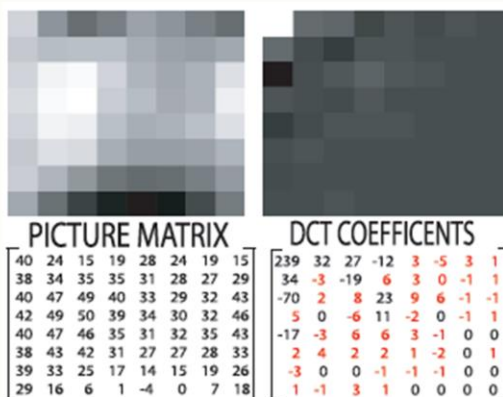
- **Standards courants** : MJPEG, MPEG-1, MPEG-2, MPEG-4
- **Méthodes de compression**:
 1. **Découpage de l'image en blocs (8x8 ou 16x16)**
 2. **Transformation**
 - DCT (TDC) Discrete Cosinus Transform
 - Ondelettes (moins utilisé)
 3. **Quantification => compression**
 4. **Prédiction (MPEG)**
 5. **Codage des données**
 - Codage arithmétique
 - Codage de Huffman

Méthodes de compression (Wikipedia)

Year	Standard	Publisher	Popular Implementations
1984	H.120	ITU-T	
1988	H.261	ITU-T	Videoconferencing, Videotelephony
1993	MPEG-1 Part 2	ISO, IEC	Video-CD
1995	H.262/MPEG-2 Part 2	ISO, IEC, ITU-T	DVD Video, Blu-ray, Digital Video Broadcasting, SVCD
1996	H.263	ITU-T	Videoconferencing, Videotelephony, Video on Mobile Phones (3GP)
1999	MPEG-4 Part 2	ISO, IEC	Video on Internet (DivX, Xvid ...)
2003	H.264/MPEG-4 AVC	Sony, Panasonic, Samsung, ISO, IEC, ITU-T	Blu-ray, HD DVD Digital Video Broadcasting, iPod Video, Apple TV,
2009	VC-2 (Dirac)	SMPTE	Video on Internet, HDTV broadcast, UHDTV
2013	H.265	ISO, IEC, ITU-T	

DCT, Discrete Cosine transform

- DCT => composantes fréquentielles du bloc
- Conserver les basse fréquences (haut gauche)



25/05/2014

Conception systèmes embarqués / MSR

21

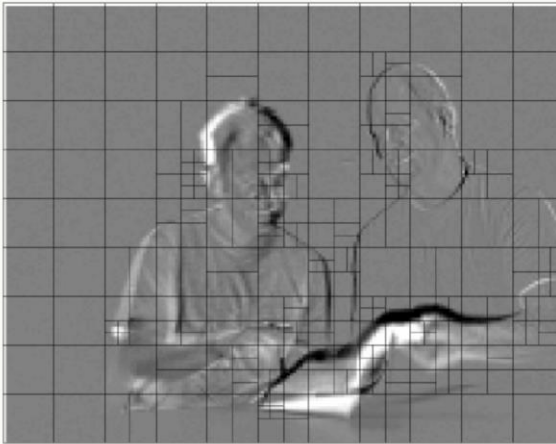
La transformée en cosinus discrète (DCT ou Discrete Cosine Transform) est une transformation proche de la transformée de Fourier. La DCT s'applique à une matrice carrée. Le résultat est représenté dans une matrice de même dimension. Il s'agit de représenter un bloc de 8x8 pixels par une matrice de même dimension mais dans le domaine fréquentiel. Au décodage, le signal est ramené dans le domaine temporel par une DCT inverse.

Les basses fréquences se trouvent en haut à gauche de la matrice, et les hautes fréquences en bas à droite.

$$X_k = \sum_{n=0}^{N-1} x_n \cos \left[\frac{\pi}{N} \left(n + \frac{1}{2} \right) k \right]$$

Division des images en blocs

- Les images sont divisées en **macroblocks** : 16 x 16 pixels, divisés en **blocs** plus petits - 8x8, 4x8, 8x4, 4x4



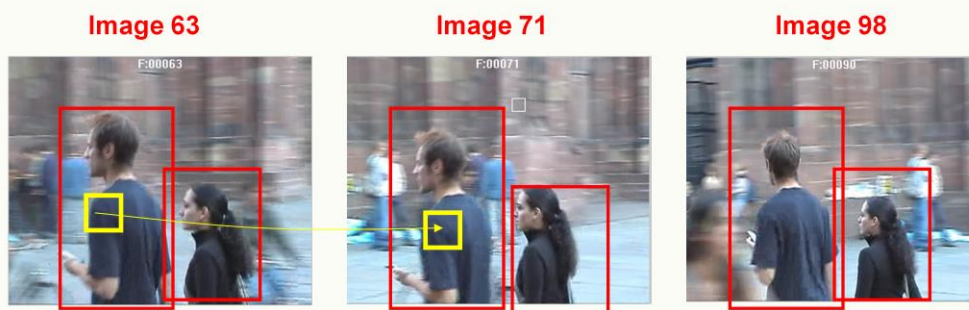
Choix optimum ?

Les systèmes de compression moins sophistiqués que le MPEG-4 AVC utilisent des blocs de taille fixe (en général 8x8). En fait, les blocs de petite taille sont plus adaptés à un niveau de détail important.

L'algorithme choisit la taille de bloc la plus adapté en fonction du niveau de détail. Plus le niveau de détail est important, plus le bloc contient de composantes haute fréquence.

Prédiction **inter** d'un bloc

- C'est **calculer** les valeurs des pixels de ce bloc, à partir des pixels provenant d'un bloc d'une **autre image** de la séquence vidéo



= > Réduire les **redondances temporelles** (les répétitions de valeurs semblables de pixels entre images successives)

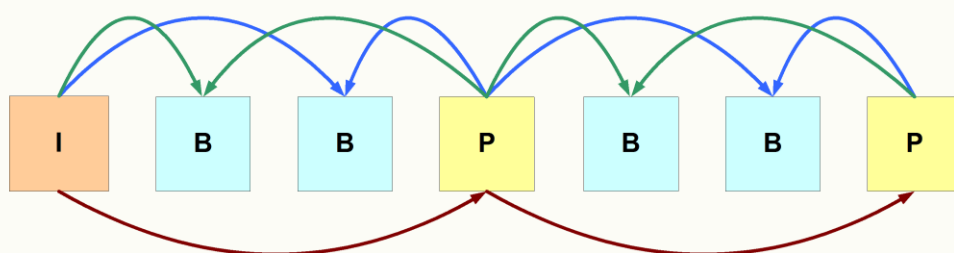
25

Toujours pour réduire la quantité d'information à transmettre, on déduit un bloc d'un autre bloc dans une autre image. Ceci est possible quand une partie d'une image se retrouve dans les images suivantes. C'est le cas de nombreuses images : soit le fond reste fixe (personnages dans une pièce par exemple), soit le fond change mais certains éléments sont toujours présents (voiture en mouvement par exemple).

Dans l'exemple ci-dessus, des promeneurs sont filmés. Les blocs composant l'image de ces promeneurs sont quasi identiques d'une image à l'autre. Comme précédemment, il est inutile de les encoder tous. Il suffit de coder les (faibles) différences entre ces blocs et la position des blocs (qui peut avoir changé d'une image à l'autre). Cette technique s'appelle : **prédiction inter**.

Séquence vidéo MPEG

- **I-Frames** : Images de référence , prédiction intra
- **P-Frames** : Images prédites à partir d'images I ou P
- **B-Frames** : Images prédites à partir de plusieurs images



Dans une séquence vidéo MPEG-4 AVC, les images ne sont pas toutes du même type:

Les **I-Frames** sont les images de référence. Le contenu d'une I-Frame ne dépend pas d'une autre image. Seule la prédiction intra est utilisée dans ces blocs. Ces images sont peu compressées (la compression intra n'est pas toujours possible) par quantification des blocs après DCT. La quantité de donnée est importante et la perte d'une I-Frame provoque un freeze de l'image (plusieurs images dépendant de l'I-Frame sont perdues).

Le **P-Frames** sont prédites à partir des I-Frames (ou parfois de P-Frames) précédentes. Ce sont des images de référence intermédiaires.

Les **B-Frames** peuvent être prédites à partir de I-Frames et de P-frames précédentes ou suivantes. Les B-Frames sont les images codées avec le minimum d'information.

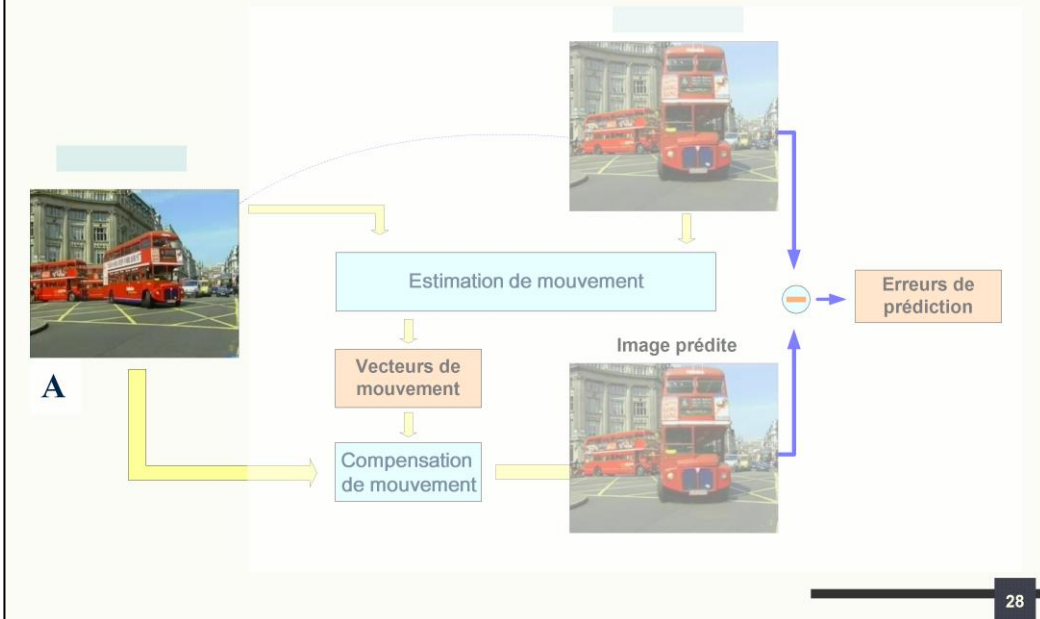
Exemple



- **Séquence vidéo : CIF 352x288 pixels -15 i/s - 1,5 Mb/s**
 - **Compression MPEG-4 AVC / H264**
 - **Un plan fixe avec des surfaces homogènes**
 - **Véhicules en mouvement**

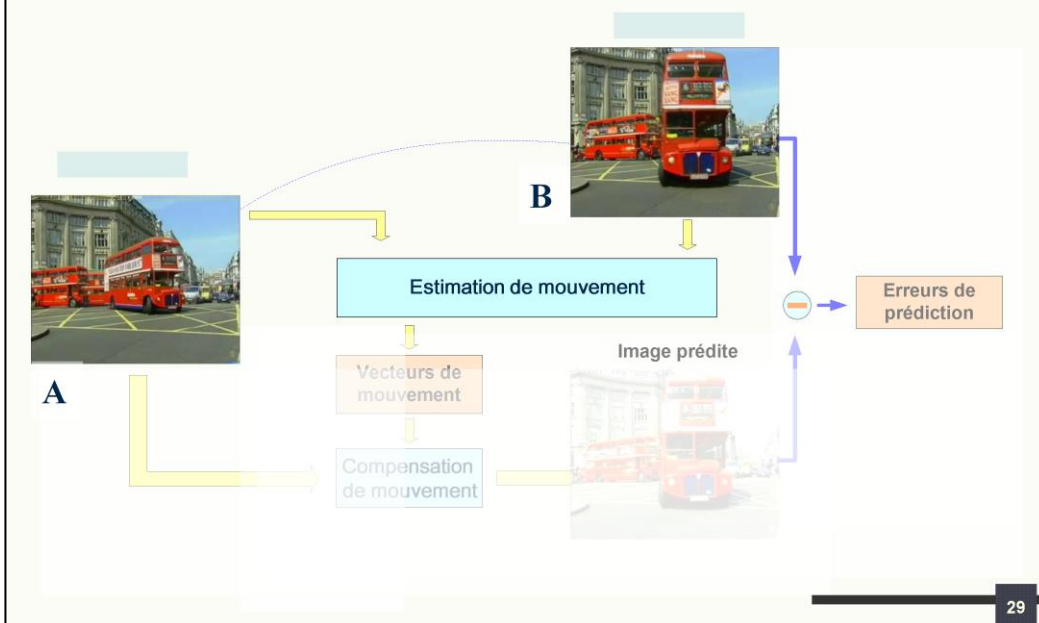
Dans l'exemple suivant, on utilise un film de test très particulier : Des bus rouges sont filmés en mouvement devant un fond fixe avec des grandes surfaces homogènes.

Prédiction Inter (codage)



1- L'image de référence est l'image A.

Prédiction Inter (codage)

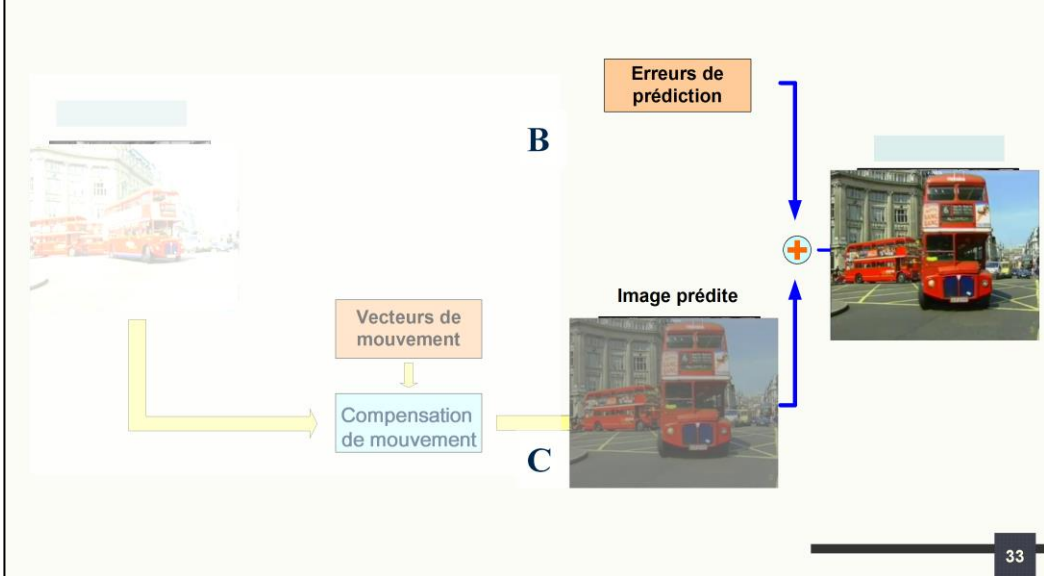


2- L'image qui doit être prédite à partir de l'image de référence A est l'image B située un peu plus loin dans la séquence vidéo.(image avec le bus de face).

L'algorithme de compression effectue une estimation de mouvement : A chaque bloc de l'image B, l'algorithme recherche un bloc plus ou moins identiques dans l'image A.

Il fait correspondre un (ou plusieurs blocs) de l'image B en calculant des vecteurs de mouvement. Les vecteurs de mouvement indiquent la direction et la distance de déplacement à appliquer à un bloc de l'image A pour obtenir le bloc correspondant de l'image B.

Prédiction Inter (décodage)



L'image B est ensuite calculée en ajoutant les erreurs de prédiction à l'image C.

Le streaming vidéo

- Taux de compression de 50 à 300
- Exemple de codec MPEG4 AVC
 - Résolution : 320x240 / 30 i/s => 76 800 pixels par image
 - Sans compression : 36 Mbits/s => Avec compression: 768 Kbits/s
- «Vrai» streaming
 - lecture quasi-instantanée à la volée d'un flux avec un buffer (FIFO)
 - nécessite un serveur de streaming
 - utilise les protocoles de diffusion RTP/RTCP sur UDP
- « Pseudo » streaming
 - téléchargement progressif d'un fichier
 - pas de serveur de streaming
 - utilise le protocole de diffusion HTTP sur TCP

Technologies de streaming :

- RealNetworks (RealMedia)
- Apple (QuickTimeMedia)
- Microsoft (Windows Media)
- Standard (ISMA)
- Cisco IP/TV
- VideoLAN (GNU Open Source)

Très bonne présentation :

Quelques mots sur la technologie de streaming

Nicolas.Meneceur

http://www.rap.prd.fr/pdf/technologie_streaming.pdf
