

Chapter 5

Fundamental Concepts in Video

[5.1 Analog Video](#)

[5.2 Digital Video](#)

[5.3 Video Display Interfaces](#)

[5.4 3D Video and TV](#)

5.1 Analog Video

- An analog signal $f(t)$ samples a time-varying image. So-called “progressive” scanning traces through a complete picture (a frame) row-wise for each time interval.
- In TV, and in some monitors and multimedia standards as well, another system, called “interlaced” scanning is used:
 - a) The odd-numbered lines are traced first, and then the even-numbered lines are traced. This results in “odd” and “even” fields — two fields make up one frame.
 - b) In fact, the odd lines (starting from 1) end up at the middle of a line at the end of the odd field, and the even scan starts at a half-way point.

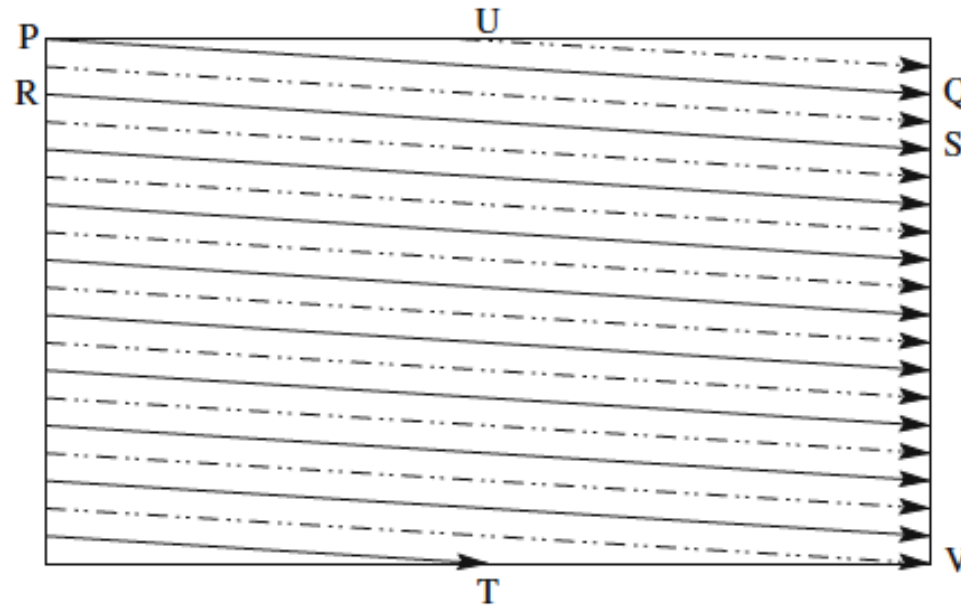


Fig. 5.1: Interlaced raster scan

- c) Figure 5.1 shows the scheme used. First the solid (odd) lines are traced, P to Q, then R to S, etc., ending at T; then the even field starts at U and ends at V.
- d) The jump from Q to R, etc. in Figure 5.1 is called the **horizontal** retrace, during which the electronic beam in the CRT is blanked. The jump from T to U or V to P is called the **vertical** retrace.

- Because of interlacing, the odd and even lines are displaced in time from each other — generally not noticeable except when very fast action is taking place on screen, when blurring may occur.
- For example, in the video in Fig. 5.2, the moving helicopter is blurred more than is the still background.



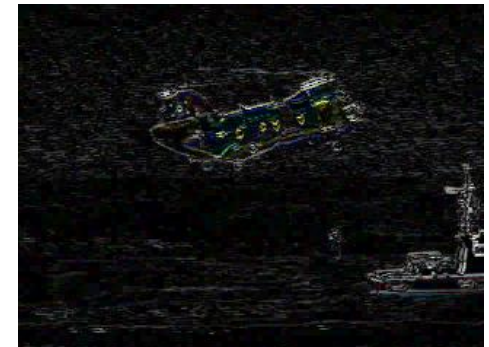
(a)



(b)



(c)



(d)

Fig. 5.2: Interlaced scan produces two fields for each frame. (a) The video frame, (b) Field 1, (c) Field 2, (d) Difference of Fields

- Since it is sometimes necessary to change the frame rate, resize, or even produce stills from an interlaced source video, various schemes are used to “de-interlace” it.
 - a) The simplest de-interlacing method consists of discarding one field and duplicating the scan lines of the other field. The information in one field is lost completely using this simple technique.
 - b) Other more complicated methods that retain information from both fields are also possible.
- Analog video use a small voltage offset from zero to indicate “black”, and another value such as zero to indicate the start of a line. For example, we could use a “blacker-than-black” zero signal to indicate the beginning of a line.

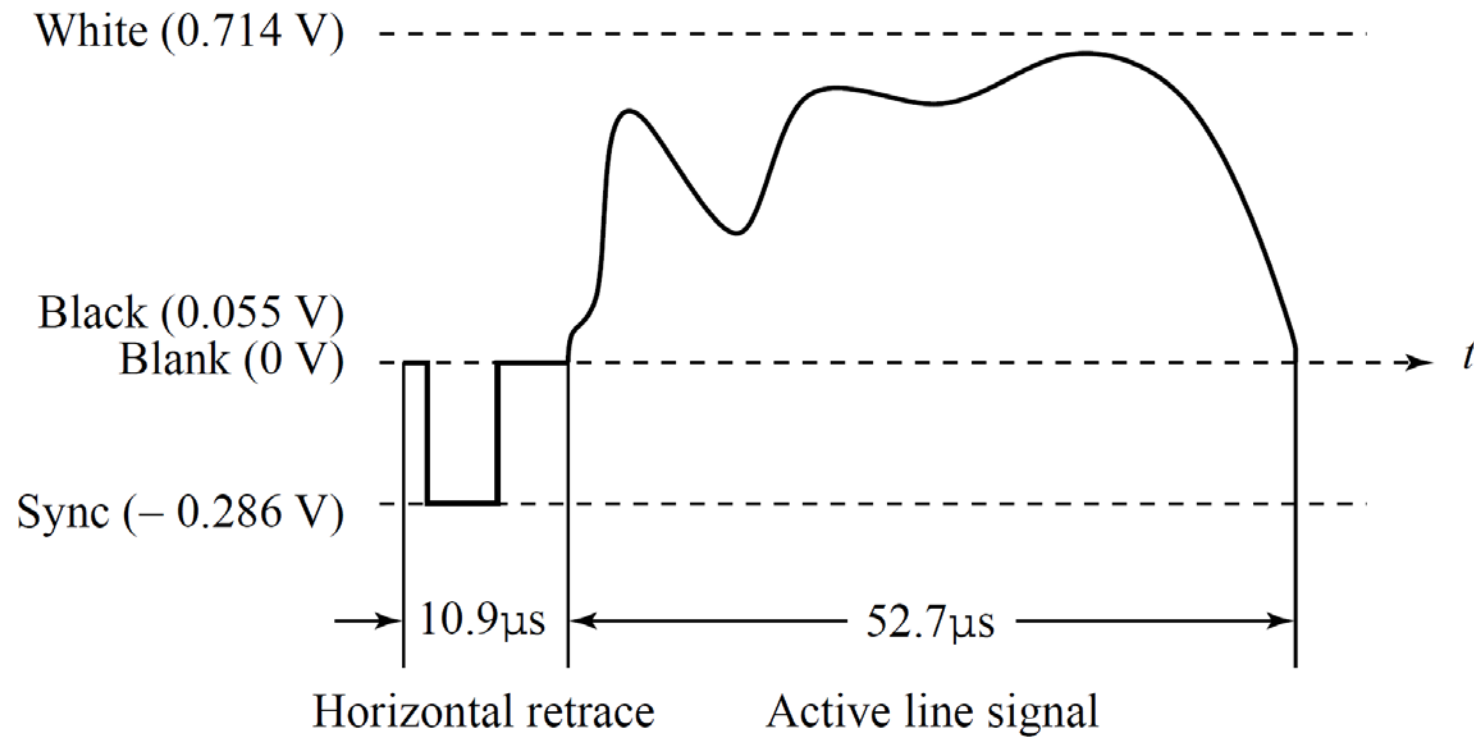


Fig. 5.3: Electronic signal for one NTSC scan line.

5.1.1 NTSC Video

- NTSC (National Television System Committee) TV standard is mostly used in North America and Japan. It uses the familiar 4:3 **aspect ratio** (i.e., the ratio of picture width to its height) and uses 525 scan lines per frame at 30 frames per second (fps).
 - a) NTSC follows the interlaced scanning system, and each frame is divided into two fields, with 262.5 lines/field.
 - b) Thus the horizontal sweep frequency is $525 \times 29.97 \approx 15,734$ lines/sec, so that each line is swept out in $1/15,734 \text{ sec} \approx 63.6 \mu\text{sec}$.
 - c) Since the horizontal retrace takes $10.9 \mu\text{sec}$, this leaves $52.7 \mu\text{sec}$ for the active line signal during which image data is displayed (see Fig.5.3).

- Fig. 5.4 shows the effect of “vertical retrace & sync” and “horizontal retrace & sync” on the NTSC video raster.

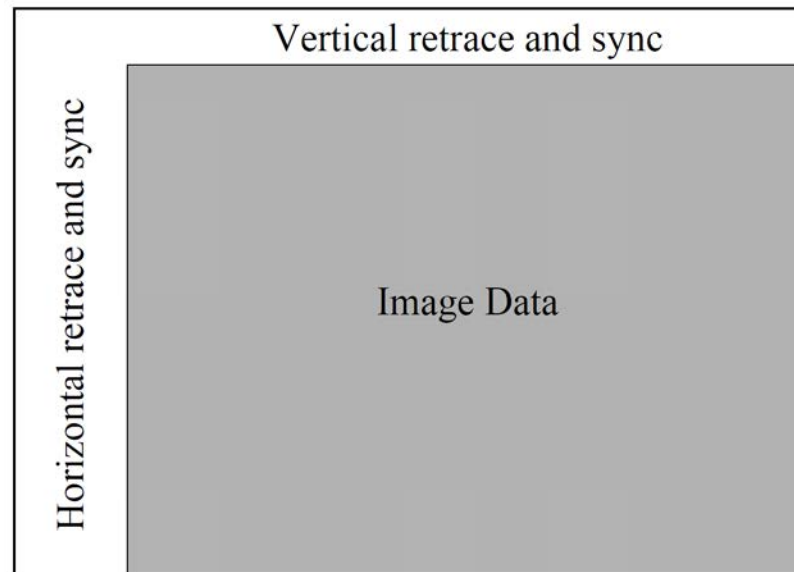


Fig. 5.4: Video raster, including retrace and sync data

- a) Vertical retrace takes place during 20 lines reserved for control information at the beginning of each field. Hence, the number of active *video lines* per frame is only 485.
- b) Similarly, almost 1/6 of the raster at the left side is blanked for horizontal retrace and sync. The non-blanking pixels are called *active pixels*.
- c) Since the horizontal retrace takes $10.9\ \mu\text{sec}$, this leaves $52.7\ \mu\text{sec}$ for the active line signal during which image data is displayed (see Fig.5.3).
- d) It is known that pixels often fall in-between the scan lines. Therefore, even with non-interlaced scan, NTSC TV is only capable of showing about 340 (visually distinct) lines, i.e., about 70% of the 485 specified active lines. With interlaced scan, this could be as low as 50%.

- NTSC video is an analog signal with no fixed horizontal resolution. Therefore one must decide how many times to sample the signal for display: each sample corresponds to one pixel output.
- A “pixel clock” is used to divide each horizontal line of video into samples. The higher the frequency of the pixel clock, the more samples per line there are.
- Different video formats provide different numbers of samples per line, as listed in Table 5.1.

Table 5.1: Samples per line for various video formats

Format	Samples per line
VHS	240
S-VHS	400-425
Betamax	500
Standard 8 m	300
Hi-8 mm	425

Color Model and Modulation of NTSC

- NTSC uses the YIQ color model, and the technique of **quadrature modulation** is employed to combine (the spectrally overlapped part of) I (in-phase) and Q (quadrature) signals into a single chroma signal C :

$$C = I \cos(F_{sc}t) + Q \sin(F_{sc}t) \quad (5.1)$$

- This modulated chroma signal is also known as the **color subcarrier**, whose magnitude is $\sqrt{I^2 + Q^2}$, and phase is $\tan^{-1}(Q/I)$. The frequency of C is $F_{sc} \approx 3.58$ MHz.
- The NTSC composite signal is a further composition of the luminance signal Y and the chroma signal as defined below:

$$\text{composite} = Y + C = Y + I \cos(F_{sc}t) + Q \sin(F_{sc}t) \quad (5.2)$$

- Fig. 5.5: NTSC assigns a bandwidth of 4.2 MHz to Y , and only 1.6 MHz to I and 0.6 MHz to Q due to human insensitivity to color details (high frequency color changes).

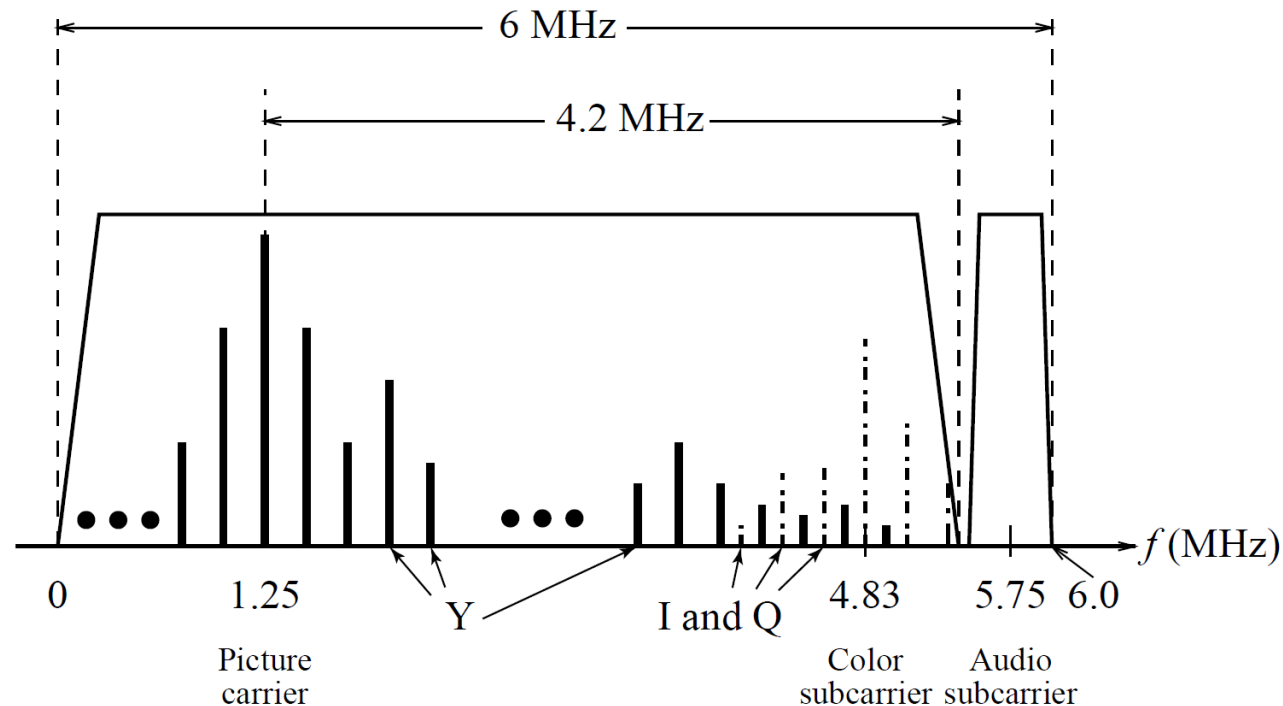


Fig. 5.5: Interleaving Y and C signals in the NTSC spectrum.

Decoding NTSC Signals

- The first step in decoding the composite signal at the receiver side is the separation of Y and C .
- After the separation of Y using a low-pass filter, the chroma signal C can be demodulated to extract the components I and Q separately. To extract I :
 1. Multiply the signal C by $2\cos(F_{sc}t)$, i.e.,

$$\begin{aligned}C \cdot 2\cos(F_{sc}t) &= I \cdot 2\cos^2(F_{sc}t) + Q \cdot 2\sin(F_{sc}t)\cos(F_{sc}t) \\&= I \cdot (1 + \cos(2F_{sc}t)) + Q \cdot 2\sin(F_{sc}t)\cos(F_{sc}t) \\&= I + I \cdot \cos(2F_{sc}t) + Q \cdot \sin(2F_{sc}t)\end{aligned}$$

- 2) Apply a low-pass filter to obtain I and discard the two higher frequency ($2F_{sc}$) terms.
- Similarly, Q can be extracted by first multiplying C by $2\sin(F_{sc}t)$ and then low-pass filtering.

- The NTSC bandwidth of 6 MHz is tight. Its audio subcarrier frequency is 4.5 MHz. The Picture carrier is at 1.25 MHz, which places the center of the audio band at $1.25 + 4.5 = 5.75$ MHz in the channel (Fig. 5.5). But notice that the color is placed at $1.25 + 3.58 = 4.83$ MHz.

- So the audio is a bit too close to the color subcarrier — a cause for potential interference between the audio and color signals. It was largely due to this reason that the NTSC color TV actually slowed down its frame rate to $30 \times 1,000 / 1,001 \approx 29.97$ fps.
- As a result, the adopted NTSC color subcarrier frequency is slightly lowered to

$$f_{sc} = 30 \times 1,000 / 1,001 \times 525 \times 227.5 \approx 3.579545 \text{ MHz},$$

** 227.5 is the # of color samples per scan line in NTSC broadcast TV.

5.1.2 PAL Video

- **PAL (Phase Alternating Line)** is a TV standard widely used in Western Europe, China, India, and many other parts of the world.
- PAL uses 625 scan lines per frame, at 25 frames/second, with a 4:3 aspect ratio and interlaced fields.
 - (a) PAL uses the YUV color model. It uses an 8 MHz channel and allocates a bandwidth of 5.5 MHz to Y, and 1.8 MHz each to U and V. The color subcarrier frequency is $f_{sc} \approx 4.43$ MHz.

- (a) In order to improve picture quality, chroma signals have alternate signs (e.g., +U and -U) in successive scan lines, hence the name “Phase Alternating Line”.
- (b) This facilitates the use of a (line rate) comb filter at the receiver — the signals in consecutive lines are averaged so as to cancel the chroma signals (that always carry opposite signs) for separating Y and C and obtaining high quality Y signals.

5.1.3 SECAM Video

- SECAM stands for *Système Electronique Couleur Avec Mémoire*, the third major broadcast TV standard.
- SECAM also uses 625 scan lines per frame, at 25 frames per second, with a 4:3 aspect ratio and interlaced fields.
- SECAM and PAL are very similar. They differ slightly in their color coding scheme:
 - (a) In SECAM, U and V signals are modulated using separate color subcarriers at 4.25 MHz and 4.41 MHz respectively.
 - (b) They are sent in alternate lines, i.e., only one of the U or V signals will be sent on each scan line.

- Table 5.2 gives a comparison of the three major analog broadcast TV systems.

Table 5.2: Comparison of Analog Broadcast TV Systems

TV System	Frame Rate (fps)	# of Scan Lines	Total Channel Width (MHz)	Bandwidth Allocation (MHz)		
				Y	I or U	Q or V
NTSC	29.97	525	6.0	4.2	1.6	0.6
PAL	25	625	8.0	5.5	1.8	1.8
SECAM	25	625	8.0	6.0	2.0	2.0

5.2 Digital Video

- The advantages of digital representation for video are many. For example:
 - (a) Video can be stored on digital devices or in memory, ready to be processed (noise removal, cut and paste, etc.), and integrated to various multimedia applications;
 - (b) Direct access is possible, which makes nonlinear video editing achievable as a simple, rather than a complex task;
 - (c) Repeated recording does not degrade image quality;
 - (d) Ease of encryption and better tolerance to channel noise.

5.2.1 Chroma Subsampling

- Since humans see color with much less spatial resolution than they see black and white, it makes sense to “decimate” the chrominance signal.
- Interesting (but not necessarily informative!) names have arisen to label the different schemes used.
- To begin with, numbers are given stating how many pixel values, per four original pixels, are actually sent:
 - a) The chroma subsampling scheme “4:4:4” indicates that no chroma subsampling is used: each pixel’s Y, Cb and Cr values are transmitted, 4 for each of Y, Cb, Cr.

- b) The scheme “4:2:2” indicates horizontal subsampling of the Cb, Cr signals by a factor of 2. That is, of four pixels horizontally labelled as 0 to 3, all four Ys are sent, and every two Cb’s and two Cr’s are sent, as (Cb0, Y0)(Cr0, Y1)(Cb2, Y2)(Cr2, Y3)(Cb4, Y4), and so on (or averaging is used).
 - c) The scheme “4:1:1” subsamples *horizontally* by a factor of 4.
 - d) The scheme “4:2:0” subsamples in both the *horizontal* and *vertical* dimensions by a factor of 2. Theoretically, an average chroma pixel is positioned between the rows and columns as shown Fig.5.6.
- Scheme 4:2:0 along with other schemes is commonly used in JPEG and MPEG (see later chapters in Part 2).

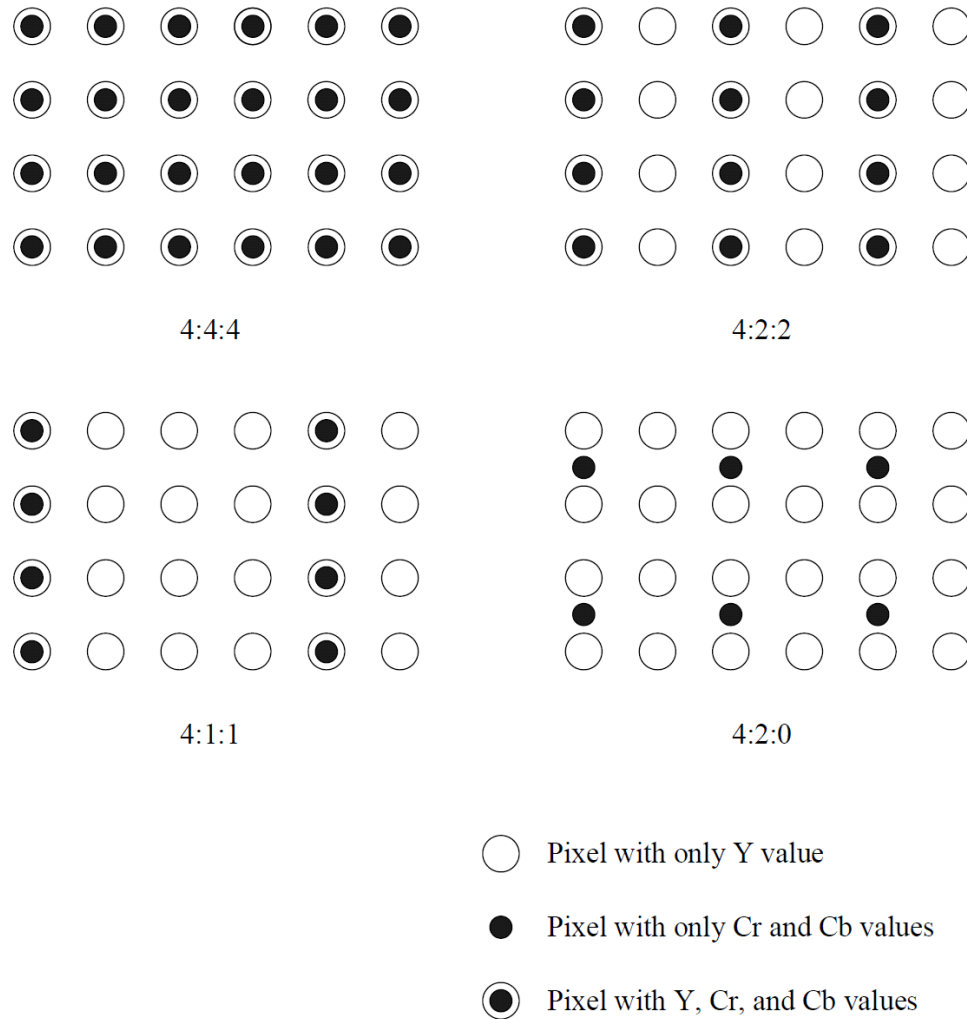


Fig. 5.6: Chroma subsampling

5.2.2 CCIR and ITU-R Standards for Digital Video

- CCIR is the Consultative Committee for International Radio, and one of the most important standards it has produced is CCIR-601, for component digital video.
 - This standard has since become standard ITU-R-601, an international standard for professional video applications
 - adopted by certain digital video formats including the popular DV video.
- Table 5.3 shows some of the digital video specifications, all with an aspect ratio of 4:3. The CCIR 601 standard uses an interlaced scan, so each field has only half as much vertical resolution (e.g., 240 lines in NTSC).

- CIF stands for Common Intermediate Format specified by the CCITT.
 - (a) The idea of CIF is to specify a format for lower bitrate.
 - (b) CIF is about the same as VHS quality. It uses a progressive (non-interlaced) scan.
 - (c) QCIF stands for “Quarter-CIF”. All the CIF/QCIF resolutions are evenly divisible by 8, and all except 88 are divisible by 16; this provides convenience for block-based video coding in H.261 and H.263, discussed later in Chapter 10.

- a) Note, CIF is a compromise of NTSC and PAL in that it adopts the NTSC frame rate and half of the number of active lines as in PAL.

Table 5.3: ITU-R digital video specifications

	CCIR 601 525/60 NTSC	CCIR 601 625/50 PAL/SECAM	CIF	QCIF
Luminance resolution	720 x 480	720 x 576	352 x 288	176 x 144
Chrominance resolution	360 x 480	360 x 576	176 x 144	88 x 72
Colour Subsampling	4:2:2	4:2:2	4:2:0	4:2:0
Fields/sec	60	50	30	30
Interlaced	Yes	Yes	No	No

5.2.3 High Definition TV (HDTV)

- The main thrust of HDTV (High Definition TV) is not to increase the “definition” in each unit area, but rather to increase the visual field especially in its width.
 - a) The first generation of HDTV was based on an analog technology developed by Sony and NHK in Japan in the late 1970s.
 - b) MUSE (MULTiple sub-Nyquist Sampling Encoding) was an improved NHK HDTV with hybrid analog/digital technologies that was put in use in the 1990s. It has 1,125 scan lines, interlaced (60 fields per second), and 16:9 aspect ratio.

- a) Since uncompressed HDTV will easily demand more than 20 MHz bandwidth, which will not fit in the current 6 MHz or 8 MHz channels, various compression techniques are being investigated.
- b) It is also anticipated that high quality HDTV signals will be transmitted using more than one channel even after compression.

- A brief history of HDTV evolution:
 - In 1987, the FCC decided that HDTV standards must be compatible with the existing NTSC standard and be confined to the existing VHF (Very High Frequency) and UHF (Ultra High Frequency) bands.
 - In 1990, the FCC announced a very different initiative, i.e., its preference for a full-resolution HDTV, and it was decided that HDTV would be simultaneously broadcast with the existing NTSC TV and eventually replace it.
 - Witnessing a boom of proposals for digital HDTV, the FCC made a key decision to go all-digital in 1993. A “grand alliance” was formed that included four main proposals, by General Instruments, MIT, Zenith, and AT&T, and by Thomson, Philips, Sarnoff and others.
 - This eventually led to the formation of the ATSC (Advanced Television Systems Committee) — responsible for the standard for TV broadcasting of HDTV.
 - In 1995 the U.S. FCC Advisory Committee on Advanced Television Service recommended that the ATSC Digital Television Standard be adopted.

- The standard supports video scanning formats shown in Table 5.4. In the table, “I” mean interlaced scan and “P” means progressive (non-interlaced) scan.

Table 5.4: Advanced Digital TV formats supported by ATSC

# of Active Pixels per line	# of Active Lines	Aspect Ratio	Picture Rate
1,920	1,080	16:9	60P 60I 30P 24P
1,280	720	16:9	60P 30P 24P
704	480	16:9 or 4:3	60P 60I 30P 24P
640	480	4:3	60P 60I 30P 24P

- For video, MPEG-2 is chosen as the compression standard. For audio, AC-3 is the standard. It supports the so-called 5.1 channel Dolby surround sound, i.e., five surround channels plus a subwoofer channel.
- The salient difference between conventional TV and HDTV:
 - a) HDTV has a much wider aspect ratio of 16:9 instead of 4:3.
 - b) HDTV moves toward progressive (non-interlaced) scan. The rationale is that interlacing introduces serrated edges to moving objects and flickers along horizontal edges.

- The services provided include:
 - SDTV (Standard Definition TV): the NTSC TV or higher.
 - EDTV (Enhanced Definition TV): 480 active lines or higher, i.e., the third and fourth rows in Table 5.4.
 - HDTV (High Definition TV): 720 active lines or higher. So far, the popular choices are:
 - 720P (1,280×720, progressive scan, 30 fps)
 - 1080I (1,920 × 1,080, interlaced, 30 fps)
 - 1080P (1,920 × 1,080, progressive scan, 30 or 60 fps).

5.2.4 Ultra High Definition TV (UHD TV)

- UHD TV is a new generation of HDTV. The standards announced in 2012 support 4K UHD TV: 2160P (3,840×2,160, progressive scan) and 8K UHD TV: 4320P (7,680×4,320, progressive scan).
- The aspect ratio is 16:9. The bit-depth can be up to 12 bits, and the chroma subsampling can be 4:2:0 or 4:2:2.
- The supported frame rate has been gradually increased to 120 fps.
- The UHD TV will provide superior picture quality, comparable to IMAX movies, but it will require a much higher bandwidth and bitrate.

5.3 Video Display Interfaces

5.3.1 Analog Display Interfaces

Analog video signals are often transmitted in one of three different interfaces: *Component video*, *Composite video*, and *S-video*. Figure 5.7 shows the typical connectors for them.



Fig. 5.7: Connectors for typical analog display interfaces. From left to right: Component video, Composite video, S-video, and VGA.

Component video

- **Component video:** Higher-end video systems make use of three separate video signals for the red, green, and blue image planes. Each color channel is sent as a separate video signal.
 - a) Most computer systems use Component Video, with separate signals for R, G, and B signals.
 - b) For any color separation scheme, Component Video gives the best color reproduction since there is no “crosstalk” between the three channels.
 - c) This is not the case for S-Video or Composite Video, discussed next. Component video, however, requires more bandwidth and good synchronization of the three components.

Composite Video

- **Composite video:** color (“chrominance”) and intensity (“luminance”) signals are mixed into a *single* carrier wave.
 - a) **Chrominance** is a composition of two color components (I and Q, or U and V).
 - b) In NTSC TV, e.g., I and Q are combined into a chroma signal, and a color subcarrier is then employed to put the chroma signal at the high-frequency end of the signal shared with the luminance signal.

- a) The chrominance and luminance components can be separated at the receiver end and then the two color components can be further recovered.
- b) When connecting to TVs or VCRs, Composite Video uses only one wire and video color signals are mixed, not sent separately. The audio and *sync* signals are additions to this one signal.
- Since color and intensity are wrapped into the same signal, some interference between the luminance and chrominance signals is inevitable.

S-Video

- **S-Video:** as a compromise, (separated video, or Super-video, e.g., in S-VHS) uses two wires, one for luminance and another for a composite chrominance signal.
- Less crosstalk between the color information and the crucial gray-scale information.
- The reason for placing luminance into its own part of the signal is that black-and-white information is most crucial for visual perception.
 - Humans are able to differentiate spatial resolution in grayscale images with a much higher acuity than for the color part of color images.
 - As a result, we can send less accurate color information than must be sent for intensity information — we can only see fairly large blobs of color, so it makes sense to send less color detail.

5.3.2 Digital Display Interfaces

Digital interfaces emerged in 1980s (e.g., Color Graphics Adapter (CGA)), and evolved rapidly. Today, the most widely used digital video interfaces include Digital Visual Interface (DVI), High-Definition Multimedia Interface (HDMI), and DisplayPort.



Fig. 5.8: Connectors of different digital display interfaces.
From left to right: DVI, HDMI, DisplayPort.

High-Definition Multimedia Interface (HDMI)

- HDMI is a newer digital audio/video interface developed to be backward compatible with DVI.
 1. HDMI doesn't carry analog signal and hence is not compatible with VGA.
 2. HDMI supports both RGB and YCbCr 4:4:4 or 4:2:2. [DVI is limited to the RGB color range (0-255).]
 3. HDMI supports digital audio, in addition to digital video.
- The maximum pixel clock rate for HDMI 1.0 is 165 MHz, which is sufficient to support 1080P (1,920 × 1,200) at 60 Hz.
- HDMI 2.0 was released in 2013, which supports 4K resolution at 60 frames per second.

5.4 3D Video and TV

- Three-dimensional (3D) pictures and movies have been in existence for decades.
- The rapid progress in R&D of 3D technology and the success of the 2009 film *Avatar* have pushed 3D video to its peak.
- Increasingly, it is in movie theaters, broadcast TV (e.g., sporting events), PCs, and various hand-held devices.
- The main advantage of the 3D video is that it enables the experience of immersion – be there, and really Be there!

5.4.1 Cues for 3D Percept

Monocular Cues:

- Shading, Perspective scaling, Relative size, Texture gradient, Blur gradient, Haze, Occlusion, Motion parallax.
- Among the above monocular cues, it is said that Occlusion and Motion parallax are more effective.

Binocular Cues:

- The human vision system uses binocular vision, i.e., *stereo vision*, aka. *Stereopsis*.
- Our left and right eyes are separated by a small distance, on average approximately two and half inches, or 65 mm. This is known as the *interocular distance*.
- The left and right eyes have slightly different views. The amount of the shift, or *disparity*, is dependent on the object's distance from the eyes, i.e., its *depth*, thus providing the binocular cue for the 3D percept.
- Current 3D video/TV systems are almost all based on stereopsis because it is believed to be the most effective cue.

5.4.2 3D Camera Models

Simple Stereo Camera Model:

- The left and right cameras are identical (same lens, same focal length, etc.). The cameras' optical axes are in parallel, pointing at the Z-direction, the scene depth.

$$d = fb/Z$$

where f is the focal length, b is the length of the baseline, $d = x_l - x_r$, is the *disparity* or *horizontal parallax*.

- Zero disparity for objects at the infinity.

Toed-in Stereo Camera Model:

- Similar to human vision system
- When humans focus on an object at a certain distance, our eyes rotate around a vertical axis in opposite directions in order to obtain (or maintain) single binocular vision — the so-called *vergence*.
- Disparity $d = 0$ at the object of focus, and at the locations that have the same distance from the observer as the object of focus.
- $d > 0$ for objects farther than the object of focus (the so-called *positive parallax*).
- $d < 0$ for nearer objects (*negative parallax*).

5.4.3 3D Movie and TV Based on Stereo Vision

3D Movie Using Colored Glasses: glasses tinted with complementary colors, usually red on the left and cyan on the right — *Anaglyph 3D*.

3D Movie Using Circularly Polarized Glasses: polarized glasses that the audience wears allows one of the two polarized pictures to pass through while blocking the other, e.g., in the RealD cinemas.

3D TV with Shutter Glasses: the liquid crystal layer on the glasses becomes opaque (behaving like a shutter) when some voltage is applied. The glasses are actively (e.g., via Infra-Red) synchronized with the TV set that alternately shows left and right images in a Time Sequential manner.

5.4.4 The Vergence-Accommodation Conflict

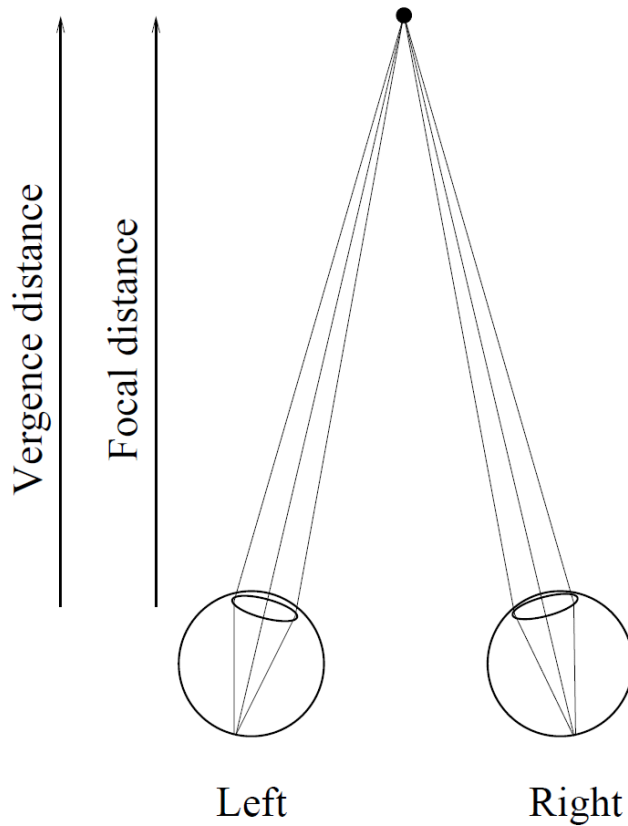
- *Accommodation* — to maintain a clear (focused) image on an object when its distance changes.
- As in Figure 5.9(a), in human vision, normally,

Focal distance = Vergence distance.

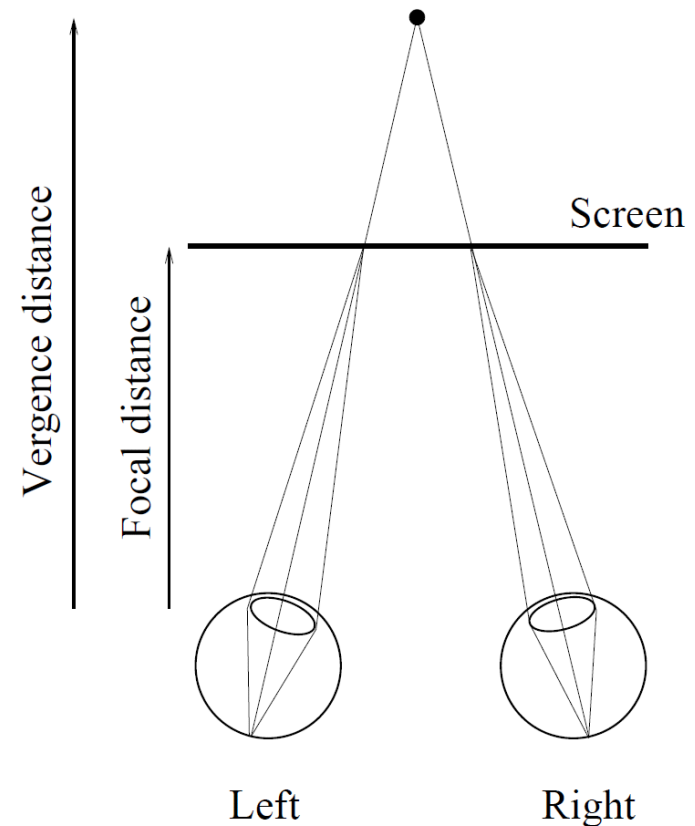
- As in Figure 5.9(b), most 3D video/movie/TV viewing requires

Focal distance \neq Vergence distance.

This creates the **Vergence-Accommodation Conflict** — one of the main reasons for eye fatigue and strain while viewing 3D video/TV/movie.



(a) Real World



(b) 3D Display

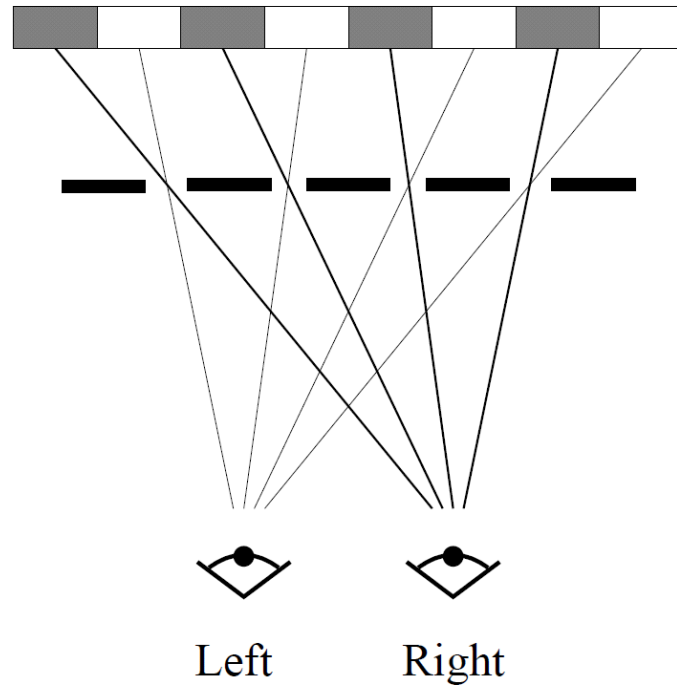
Fig. 5.9: The Vergence-Accommodation Conflict.

5.4.5 Autostereoscopic (Glasses-Free) Display Devices

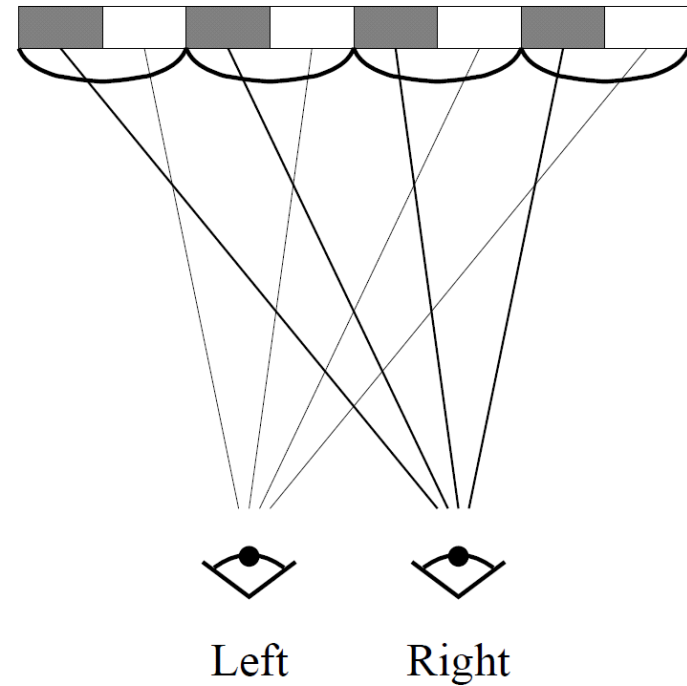
Parallax Barrier — a layer of opaque material with slits is placed in front of the normal display device, e.g., an LCD. See Fig. 5.10(a). Used in e.g., Nintendo 3DS, Fujifilm 3D camera, and Toshiba's glasses-free 3D TV.

Lenticular Lens — Instead of barriers, columns of magnifying lenses can be placed in front of the display to direct lights properly to the left and right eyes. See Fig. 5.10(b).

Integral Imaging — Instead of cylindrical lenses as shown above, an array of spherical convex microlenses can be used to generate a large number of distinct micro-images. It enables the rendering of multiple views from any directions.



(a) Parallax Barrier



(b) Lenticular Lens

Fig. 5.10: Autostereoscopic Display Devices.

5.4.6 Disparity Manipulation in 3D Content Creation

- **Disparity Range** — map (often suppress) the original disparities into the range that will fit in the *comfort zone* of most viewers.
- **Disparity Sensitivity** — human vision is more capable of discriminating different depths when they are nearby, so do nonlinear disparity mapping (suppress the far range).
- **Disparity Gradient** — human vision has a limit of disparity gradient in binocular fusion, so avoid it in depth editing.
- **Disparity Velocity** — we cannot rapidly process large accommodation and vergence changes (i.e., disparity changes), so slow it down!