

Chapter 5 The Hilbert transform and linear modulation theory

5.1 The Hilbert transform

In AC circuit theory we use complex exponentials to represent real sinusoidally varying quantities with the understanding that the real part of the complex exponential gives the physical quantity of interest. The **analytic signal** of a real function plays the same role for more general waveforms. Given a **real function** $f(t)$ with Fourier transform $F(\nu)$, the analytic signal $f_a(t)$ is defined by

$$f_a(t) = 2 \int_0^{\infty} F(\nu) \exp(j2\pi\nu t) d\nu \quad (5.1)$$

This is just the inverse Fourier transform of the positive frequency part of $F(\nu)$. The Fourier transform $F_a(\nu)$ of $f_a(t)$ is given by

$$F_a(\nu) = 2u(\nu)F(\nu) \quad (5.2)$$

Lemma: The real part of $f_a(t)$ is equal to $f(t)$.

Proof: The real part of $f_a(t)$ is

$$\frac{1}{2}[f_a(t) + f_a^*(t)] = \int_0^{\infty} F(\nu) \exp(j2\pi\nu t) d\nu + \int_0^{\infty} F^*(\nu) \exp(-j2\pi\nu t) d\nu \quad (5.3)$$

$$= \int_0^{\infty} F(\nu) \exp(j2\pi\nu t) d\nu + \int_0^{\infty} F(-\nu) \exp(-j2\pi\nu t) d\nu \quad (5.4)$$

since $F(\nu)$ is a Hermitian function. Substituting $-\nu$ for ν in the second integral and combining the two integrals yields the inverse Fourier transform of $F(\nu)$ which is just $f(t)$.

The **Hilbert transform** of $f(t)$ is defined to be the imaginary part of $f_a(t)$ and is denoted by $\hat{f}(t)$. Thus

$$f_a(t) = f(t) + j\hat{f}(t) \quad (5.5)$$

The Fourier transform of this relationship is

$$2u(\nu)F(\nu) = F(\nu) + j\hat{F}(\nu) \quad (5.6)$$

where \hat{F} is the Fourier transform of $\hat{f}(t)$. Solving this yields

$$\hat{F}(\nu) = -j \operatorname{sgn}(\nu)F(\nu) \quad (5.7)$$

The **envelope** of $f(t)$ is $|f_a(t)| = \sqrt{f(t)^2 + \hat{f}(t)^2}$.

Note: Unlike the Fourier transform, the Hilbert transform of a function of t is also a function of t .

Exercise: Show that the analytic signal of $R \cos(\omega t + \phi)$ is $A \exp(j\omega t)$ where $A = R \exp(j\phi)$.

Using the transform pair

$$\frac{1}{\pi t} \leftrightarrow -j \operatorname{sgn}(\nu) \quad (5.8)$$

we can take the inverse Fourier transform of (5.7) to obtain

$$\hat{f}(t) = f(t) * \frac{1}{\pi t} \quad (5.9)$$

$$= \int_{-\infty}^{\infty} \frac{f(\tau)}{\pi(t-\tau)} d\tau \quad (5.10)$$

where the integral is to be interpreted as a Cauchy principal value. This convolution can be interpreted as a filtering operation with a **quadrature filter** which shifts all sinusoidal components by a phase shift of $-\pi/2$.

Similarly, starting from

$$F(\nu) = j \operatorname{sgn}(\nu) \hat{F}(\nu) \quad (5.11)$$

we find that

$$f(t) = - \int_{-\infty}^{\infty} \frac{\hat{f}(\tau)}{\pi(t-\tau)} d\tau \quad (5.12)$$

Note that the amplitudes of the Fourier transforms of $f(t)$ and of $\hat{f}(t)$ are identical, they differ only in phase. The ear is largely insensitive to the phases of the Fourier components in a signal and so the Hilbert transform of a speech signal is still understandable.

Exercises:

1. Show from the integral (5.10) that the Hilbert transform of $\cos(\omega t)$ is $\sin(\omega t)$.
2. Find the Hilbert transform of $\Pi(t)$ and sketch the result.
3. If $m(t)$ is a nonnegative real-valued band-limited function and ν_c is larger than the bandwidth of $m(t)$, show that the envelope of $m(t) \cos(2\pi\nu_c t + \phi)$ is $m(t)$. Thus the definition of the envelope given above coincides with the intuitive concept.

5.2 The Fourier transform of a real causal function

Given a function $h(t)$ which is causal in the sense that $h(t) = 0$ for $t < 0$ and which has no singularities at the origin, we see that

$$h(t) = u(t)h(t) \quad (5.13)$$

Taking the Fourier transform of this relationship yields

$$H(\nu) = \left(\frac{1}{2} \delta(\nu) + \frac{1}{j2\pi\nu} \right) * H(\nu) \quad (5.14)$$

$$= \frac{1}{2} H(\nu) - \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{jH(\nu')}{\nu - \nu'} d\nu' \quad (5.15)$$

Writing $H(\nu)$ in terms of its real and imaginary parts leads to

$$\Im H(\nu) = -\frac{1}{\pi} \int_{-\infty}^{\infty} \frac{\Re H(\nu')}{\nu - \nu'} d\nu' \quad (5.16)$$

$$\Re H(\nu) = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{\Im H(\nu')}{\nu - \nu'} d\nu' \quad (5.17)$$

Thus if h is causal and regular at the origin, given the real part of the Fourier transform of h , it is possible to deduce the imaginary part and vice versa.

If in addition h is real so that H is Hermitian,

$$\Im H(\nu) = -\frac{2\nu}{\pi} \int_0^{\infty} \frac{\Re H(\nu')}{\nu^2 - \nu'^2} d\nu' \quad (5.18)$$

$$\Re H(\nu) = \frac{2}{\pi} \int_0^{\infty} \frac{\nu' \Im H(\nu')}{\nu^2 - \nu'^2} d\nu' \quad (5.19)$$

These relationships may be generalized to include cases in which h has singularities. They are often known as **dispersion relations** because similar equations were first derived in optics (by Kronig and Kramers) to relate the real and imaginary parts of the frequency-dependent refractive index of a material. The real part of the refractive index determines the speed of propagation of waves and the dispersion whereas the imaginary part determines the absorption. From general physical constraints such as causality, it is possible to derive relationships between these quantities. For example, a frequency range of anomalous dispersion (near resonances in the atomic levels) coincides with a range of large absorption.

Exercise: Show that if $h(t)$ is a real causal function and $\Im H(\nu) = -\nu/(a^2 + \nu^2)$ where a is real, then $\Re H(\nu) = a/(a^2 + \nu^2)$.

5.3 Transmission and reception of information

When sending information from one place to another, a suitable long-range physical process (such as electromagnetic radiation of the appropriate frequency) is used as a **carrier**. The message is used to modify some characteristic of the carrier in a process called **modulation**. At the receiver, the process of **demodulation** is used to recover the message.

Typically, the carrier is a sinusoidally varying signal at some frequency ν_c which is much higher than the frequencies contained in the message. For example, for medium-wave broadcast radio, ν_c is in the range 0.5–1.6 MHz. Short-wave radio extends to about 30 MHz and VHF television to about 300 MHz. For a pure sinusoidal signal, the spectrum consists of (arbitrarily narrow) delta functions at $\pm\nu_c$. As we shall see, the process of modulation gives the signal a nonzero bandwidth. However the bandwidth usually is much smaller than ν_c and so the signal may be regarded as a narrow band process. This is an important feature of modulation since the dispersion and differential phase-shift introduced by the medium may lead to distortion unless the fractional bandwidth is small.

As a general rule, when independent (non-cooperating) signals are sent together, they must occupy non-overlapping frequency bands in order not to interfere with each other. When considering different modulation techniques, several issues arise

1. How easy is it to do the modulation and demodulation?
2. What bandwidth is occupied by the transmitted signal?
3. How much power is required to send the transmitted signal?
4. How tolerant is the scheme to added noise which may corrupt the signal on its way to the receiver?

At this stage, we can consider the first three of these issues. We shall be considering a few examples to illustrate how the transform methods developed so far help us to understand the processes involved. Each of these examples are based on linear amplitude modulation. Other modulation techniques include frequency modulation and many forms of pulse modulation.

We shall use the following notation:

1. $c(t) = A \cos(2\pi\nu_c t)$ is the carrier waveform of amplitude $A > 0$ and frequency ν_c ,
2. $m(t)$ is the modulating waveform which contains the message. The Fourier transform of this is $M(\nu)$ and we shall assume that its bandwidth is ν_m , i.e., $M(\nu) = 0$ for $|\nu| > \nu_m$. Without loss of generality we shall assume that $|m(t)| \leq 1$ for all t .
3. $f(t)$ shall denote the transmitted signal and $F(\nu)$ is its spectrum.

5.4 Double sideband (DSB) modulation

The transmitted signal is the product of the modulating and carrier functions

$$f(t) = m(t)c(t) = Am(t) \cos(2\pi\nu_c t) \quad (5.20)$$

Taking the Fourier transform of this yields

$$F(\nu) = AM(\nu) * \frac{1}{2}[\delta(\nu - \nu_c) + \delta(\nu + \nu_c)] = \frac{1}{2}A[M(\nu - \nu_c) + M(\nu + \nu_c)] \quad (5.21)$$

The spectrum of DSB consists of two copies of the spectrum of the modulating signal, one shifted to ν_c and the other to $-\nu_c$. The frequencies on either side of the carrier are called the **sidebands**. The bandwidth of a DSB signal is thus twice the bandwidth of the modulating signal. If $\nu_c > \nu_m$, the analytic signal of a DSB signal is

$$f_a(t) = A \exp(j2\pi\nu_c t)m(t) \quad (5.22)$$

and so the envelope is $A|m(t)|$. This absolute value means that if $m(t)$ goes negative, the modulating signal cannot be recovered simply using envelope detection.

Exercise: Sketch the DSB signal, its envelope and spectrum when $m(t) = \cos(2\pi\nu_m t)$. Notice how there are phase reversals in $f(t)$ at the zeros of $m(t)$.

Demodulation of DSB involves multiplying $f(t)$ by a reconstructed carrier $\cos(2\pi\nu t)$ followed by a low-pass filter with cutoff frequency ν_m . This reconstructed carrier must be provided by the receiver.

$$f(t) \cos(2\pi\nu_c t) = Am(t) \cos^2(2\pi\nu_c t) \quad (5.23)$$

$$= \frac{1}{2}Am(t) + \frac{1}{2}Am(t) \cos(4\pi\nu_c t) \quad (5.24)$$

The low-pass filter removes the second term leaving a copy of the modulating signal. The demodulation process can also be understood in the frequency domain in terms of an additional convolution.

Exercise: What happens if the reconstructed carrier is not exactly correct so that it differs from the true carrier in phase or frequency?

The need to reconstruct an exact copy of the carrier is a major problem with DSB. One solution is add a small amount of the carrier (called a **pilot carrier**) to the transmitted signal so that it can be recovered at the receiver and used to demodulate the signal. This also alleviates the inconvenience that when $m(t) = 0$ the transmitted signal $f(t)$ also vanishes, making it impossible to tune to a transmission when there is no modulation.

Modulators and demodulators for DSB are essentially multipliers which are often called **balanced** modulators or demodulators for historical reasons. These are often based on the quarter-square multiplier principle in which the signals $c(t) + m(t)$ and $c(t) - m(t)$ are each passed through devices with matched quadratic nonlinearities. The DSB signal is obtained using

$$\frac{1}{4}[(c(t) + m(t))^2 - (c(t) - m(t))^2] = c(t)m(t) \quad (5.25)$$

They are moderately difficult to build using analogue electronics.

5.5 Amplitude modulation (AM)

This is the form of modulation most commonly used for commercial broadcasting on the medium-wave and short-wave radio bands. The signal is the same as DSB except that the carrier is injected and transmitted along with the signal

$$f(t) = (1 + m(t))c(t) = A(1 + m(t)) \cos(2\pi\nu_c t) \quad (5.26)$$

Since $|m(t)| \leq 1$ by assumption, the factor $A(1 + m(t))$ is always non-negative. It is easy to show that if $\nu_c > \nu_m$, $A(1 + m(t))$ is the envelope of $f(t)$ and so $m(t)$ can be recovered using a simple envelope detector. This ease of demodulation largely accounts for the popularity of AM for commercial broadcasting since receivers can be made very inexpensively.

In terms of the spectra, it is easy to show that

$$F(\nu) = \frac{1}{2}A[\delta(\nu - \nu_c) + \delta(\nu + \nu_c) + M(\nu - \nu_c) + M(\nu + \nu_c)] \quad (5.27)$$

which is just the same as DSB except for the addition of the delta functions representing the carrier. The bandwidth is the same as for DSB but additional power is required to transmit the carrier.

Exercise: Show that if $m(t) = B \cos(2\pi\nu_m t)$ where $|B| < 1$ the fraction of the total power in $f(t)$ that is contained within the carrier is $2/(2 + B^2)$. Note that this is always at least $\frac{2}{3}$.

Amplitude modulators are relatively easy to make and can work at the very high powers used for broadcasting. They are usually based on passing $c(t) + m(t)$ through some nonlinearity and then filtering the result (e.g. with a resonant tuned circuit) around ν_c .

5.6 Single sideband (SSB) modulation

Amplitude modulation is very inefficient since at least two-thirds of the total signal power is contained in the carrier component which does not contain any information. Furthermore both sidebands are transmitted, although one is essentially a mirror-image of the other (since $M(\nu)$ is

Hermitian) and so only one is strictly necessary. Single sideband avoids both of these problems at the cost of a more complicated modulator. Its high efficiency however has made it the most popular technique for voice-based communication (radio-telephones etc.) on the HF bands.

The SSB modulated signal is

$$f(t) = A\Re(m_a(t) \exp(j2\pi\nu_c t)) \quad (5.28)$$

$$= A[m(t) \cos(2\pi\nu_c t) - \hat{m}(t) \sin(2\pi\nu_c t)] \quad (5.29)$$

where m_a denotes the analytic signal and \hat{m} the Hilbert transform of m . Taking the Fourier transform of this result we see that

$$m_a(t) \exp(j2\pi\nu_c t) \leftrightarrow 2u(\nu - \nu_c)M(\nu - \nu_c) \quad (5.30)$$

$$f(t) = A\Re[m_a(t) \exp(j2\pi\nu_c t)] \leftrightarrow A[u(\nu - \nu_c)M(\nu - \nu_c) + u(-\nu - \nu_c)M^*(-\nu - \nu_c)] \quad (5.31)$$

Graphically, the positive frequency components of m are shifted by ν_c and the negative frequency components are shifted by $-\nu_c$. Comparing this with the spectrum of AM found earlier, we see that the spectrum of f only contains the upper sideband (frequencies $|\nu| > \nu_c$). Both the lower sideband and the carrier are not transmitted and so all the power in the transmitted signal sends information about m . Note that the components at negative frequencies $\nu < -\nu_c$ cannot be suppressed since we want $f(t)$ to be a **real** signal.

The above signal is more precisely called an USB (upper sideband) modulated signal. Similarly, we can define a LSB (lower sideband) modulated signal by

$$f(t) = A\Re(m_a(t) \exp(-j2\pi\nu_c t)) \quad (5.32)$$

$$= A[m(t) \cos(2\pi\nu_c t) + \hat{m}(t) \sin(2\pi\nu_c t)] \quad (5.33)$$

SSB signals can be demodulated by multiplying them with a reconstructed carrier $\cos(2\pi\nu_c t)$ and low-pass filtering the result, just as for DSB. Unlike DSB however, the precise phase and frequency of the reconstructed carrier is not so critical for telephony (voice) transmission.

Exercise: Show that if the reconstructed carrier is $\cos(2\pi\nu'_c t + \phi)$, the result of demodulating the USB signal (5.29) is

$$\frac{A}{2}\Re[m_a(t) \exp(-j[2\pi\Delta\nu t + \phi])] \quad (5.34)$$

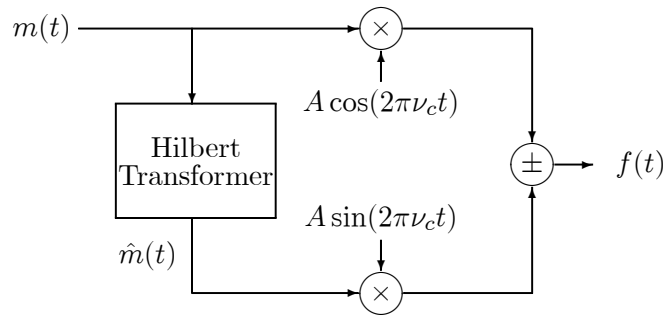
where $\Delta\nu = \nu'_c - \nu_c$ is the detuning. In particular, if $m(t) = \cos(2\pi\nu_m t)$ show that the demodulated signal is proportional to $\cos[2\pi(\nu_m - \Delta\nu)t - \phi]$.

From the exercise, we see that a detuning simply shifts the frequency of each sinusoidal component within $m(t)$ by a fixed amount. Speech that is distorted in this way is still readily understandable for detunings of a few Hz. In DSB however, the corresponding detuning causes the demodulated signal to fade in and out at a rate $\Delta\nu$.

(Application note: The ability to shift the frequency of a voice signal by a small amount is sometimes used in public address systems to reduce the likelihood of howl-around or oscillation when the volume is turned up.)

5.6.1 The phasing method for generating SSB

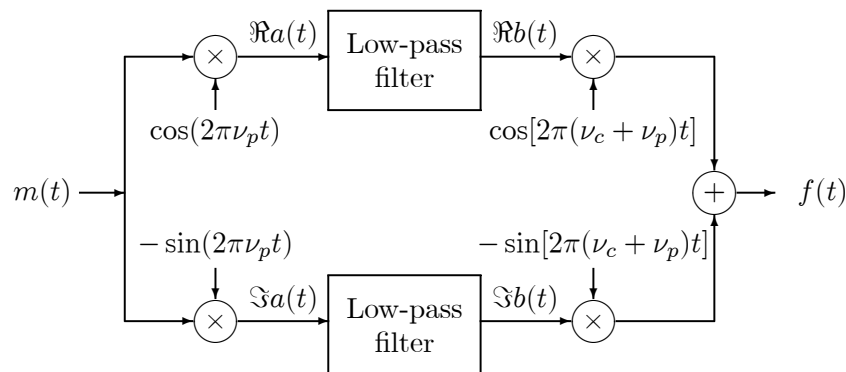
This is a direct implementation of (5.29) or (5.33). The modulating signal $m(t)$ is passed into a Hilbert transformer and $m(t)$ and $\hat{m}(t)$ are fed into two balanced modulators (multipliers) driven with $\cos(2\pi\nu_c t)$ and $\sin(2\pi\nu_c t)$. By adding or subtracting the outputs of the balanced modulators, either LSB or USB modulation can be performed.



The main technical difficulty with this approach is the implementation of the Hilbert transformer for $m(t)$ which must operate over a wide range of frequencies (about 300 Hz to 3 kHz for telephone-quality voice). One common solution is to build two RC networks with matched amplitude characteristics but such that one shifts all frequencies by 90° more than the other. If $m(t)$ is fed into both networks, an approximate Hilbert transform relationship holds between the outputs of the two networks and these are fed into the balanced modulators. Such networks require high-precision components with non-preferred values to give reasonable results.

5.6.2 Weaver's method

This requires a more complex modulator but avoids the use of the wide-band phase-shifter.



The operation of this modulator is best understood by thinking in terms of complex-valued signals and considering the spectra at different stages. Starting with $m(t)$, we define the complex signal

$$a(t) = m(t) \exp(-j2\pi\nu_p t) \quad (5.35)$$

where $\nu_p = \frac{1}{2}\nu_m$ and ν_m is the bandwidth of $m(t)$. In the frequency domain, $A(\nu) = M(\nu + \nu_p)$ which is a shifted version of $M(\nu)$. The positive frequency part of $M(\nu)$ now straddles the origin and falls in the band $[-\nu_p, \nu_p]$. If we now low-pass filter $a(t)$ using a filter with cutoff frequency ν_p , this will extract those components which originally formed the positive frequency part of $M(\nu)$. Let this filtered version of $a(t)$ be denoted $b(t)$. It is easy to see that

$$B(\nu) = u(\nu + \nu_p)M(\nu + \nu_p) \quad (5.36)$$

We now consider $b(t) \exp(j2\pi(\nu_c + \nu_p)t)$. This shifts the spectrum of B by $\nu_c + \nu_p$ resulting in

$$B(\nu - \nu_c - \nu_p) = u(\nu - \nu_c)M(\nu - \nu_c) \quad (5.37)$$

If we now consider the real part $\Re[b(t) \exp(j2\pi(\nu_c + \nu_p)t)]$, the spectrum is

$$\frac{1}{2}[B(\nu - \nu_c - \nu_p) + B^*(-\nu - \nu_c - \nu_p)] = \frac{1}{2}[u(\nu - \nu_c)M(\nu - \nu_c) + u(-\nu - \nu_c)M^*(-\nu - \nu_c)] \quad (5.38)$$

This is the spectrum of a USB signal, showing that the modulator works.

We now need to consider the physical implications of working with complex signals. A complex number may be regarded simply as a pair of real numbers (the real part and the imaginary part) together with appropriate rules for addition and multiplication. In the diagram, the upper half implements the real part of the system and the lower half the imaginary part. Thus for example, the real and imaginary parts of (5.35) are

$$\Re a(t) = m(t) \cos(2\pi\nu_p t) \quad \text{and} \quad \Im a(t) = -m(t) \sin(2\pi\nu_p t) \quad (5.39)$$

Low-pass filtering of the complex signal $a(t)$ corresponds to filtering the real and imaginary parts separately. In the final stage,

$$\Re[b(t) \exp(j2\pi(\nu_c + \nu_p)t)] = \Re[b(t)] \cos(2\pi(\nu_c + \nu_p)t) - \Im[b(t)] \sin(2\pi(\nu_c + \nu_p)t) \quad (5.40)$$

which is implemented using the last two balanced modulators and the summer.

5.7 Quadrature Modulation/Demodulation – The Complex Envelope

It is possible to send two real modulating signals using a single carrier by using quadrature modulation. If $m_1(t)$ and $m_2(t)$ are real-valued modulating signals, we define the complex modulating signal as

$$m(t) = m_1(t) + jm_2(t) \quad (5.41)$$

Quadrature modulation involves the transmission of the signal

$$f(t) = \Re[m(t) \exp(j2\pi\nu_c t)] \quad (5.42)$$

The spectrum of the quadrature modulated signal is

$$F(\nu) = \frac{1}{2}[M(\nu - \nu_c) + M^*(-\nu - \nu_c)] \quad (5.43)$$

This consists of sidebands on either side of the carrier ν_c , but unlike DSB, the upper and lower sidebands are not mirror images of each other. The output is a narrowband process around the carrier frequency ν_c .

The **complex envelope** of a narrowband process $g(t)$ with respect to frequency ν_0 is defined to be the signal

$$g_a(t) \exp(-j2\pi\nu_0 t) \quad (5.44)$$

where $g_a(t)$ is the analytic signal of $g(t)$. It is easy to see that the complex envelope of the quadrature modulated signal $f(t)$ with respect to the carrier frequency is $m(t)$.

Quadrature modulation is demodulated by multiplying the signal with $\exp(-j2\pi\nu_c t)$ and low-pass filtering the result. The real and imaginary parts of the low-pass filter output are $m_1(t)$ and $m_2(t)$ respectively. Just as in DSB demodulation, it is **essential** that the phase and frequency of this reconstructing carrier be the same as that at the transmitter.

Exercise: Splitting real and imaginary parts as discussed above, sketch a block diagram showing how to perform quadrature modulation and demodulation.

5.8 The Optical Homodyne Detector

Extracting the complex envelope of an optical signal presents special challenges because the carrier frequency ν_c is very high. This is of course also the main advantage of using light for communications since the fractional bandwidth occupied by the modulated signal is very small.

Typically we will wish to determine $m(t)$ from a low-intensity light field whose electric field is given by

$$f(t) = \Re[m(t) \exp(j2\pi\nu_c t)] \quad (5.45)$$

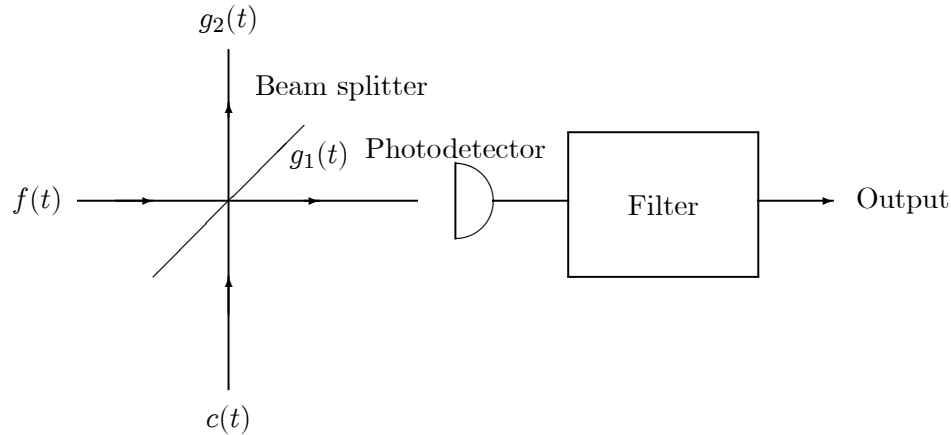
The **optical homodyne detector** consists of a beam splitter which linearly combines the input signal $f(t)$ and a strong coherent local oscillator for which the field is

$$c(t) = \Re[A \exp(j[2\pi\nu_c t - \phi])] \quad (5.46)$$

where A is real. If the amplitude transmission coefficient of the beam-splitter is T and the amplitude reflection coefficient is R , the conditions for a lossless beam-splitter are

$$|T|^2 + |R|^2 = 1 \quad (5.47)$$

$$TR^* + RT^* = 0 \quad (5.48)$$



There are two output ports of the beam splitter. The fields in each are

$$g_1(t) = \Re\{[RA \exp(-j\phi) + T m(t)] \exp(j2\pi\nu_c t)\} \quad (5.49)$$

$$g_2(t) = \Re\{[TA \exp(-j\phi) + R m(t)] \exp(j2\pi\nu_c t)\} \quad (5.50)$$

In an **unbalanced** homodyne detector, a photodetector which measures the **intensity** of the light is placed in one of the output ports and its output is low-pass filtered. If this is in the first port, we see that the intensity is proportional to

$$I_1(t) \propto |RA \exp(-j\phi) + T m(t)|^2 \quad (5.51)$$

$$= |R|^2 A^2 + 2A \Re\{R^* T m(t) \exp(j\phi)\} + |T|^2 |m(t)|^2 \quad (5.52)$$

The second term is the desired output. By adjusting ϕ appropriately we can recover $\Re[m(t)]$, $\Im[m(t)]$ or any intermediate quadrature. Notice that the homodyne detector has amplified the weak input signal by $A|RT|$ which in principle can be made as large as desired. The first term is a constant background due to the local oscillator and the third is usually small since $|m(t)| \ll A$.

In a **balanced** homodyne detector, a second photodetector is placed in the other output port of the beam splitter and the intensity difference is computed before filtering. The reflection and transmission coefficients are adjusted so that $|R|^2 = |T|^2 = 0.5$. We find that

$$I_1(t) - I_2(t) \propto 4A\Re(R^*Tm(t)\exp(j\phi)) \quad (5.53)$$

which consists of only the desired term.

The above discussion is a purely classical description of the homodyne detector. In the classical theory, it is possible to produce a light signal in which $m(t)$ is a constant (complex) quantity. This corresponds to a light wave of known amplitude and phase. In quantum mechanics however, it turns out that there is always some residual fluctuation in $m(t)$ no matter how well the light source is stabilized. The Heisenberg uncertainty principle shows that the **product** of the root-mean-square fluctuations in $\Re[m(t)]$ and $\Im[m(t)]$ cannot be less than a certain amount. For a coherent light source, the r.m.s. fluctuations in $\Re[m(t)]$ and $\Im[m(t)]$ are equal to each other whereas for a **squeezed** light source, these fluctuations are unequal (although they still satisfy the uncertainty principle.) Note that since the fluctuations in $m(t)$ have components at all frequencies, it is usually necessary to describe the behaviour of the fluctuations in different frequency bands corresponding to the placement of different filters after the photodetector. This gives rise to the concept of the **squeezing spectrum**.

5.9 The television video signal

In a television set, an electron beam scans across the screen and its intensity is altered in order to produce a picture (see Figure C.1). In $1/50$ s, the spot traces out a **field** of $312\frac{1}{2}$ lines starting at the top left of the screen and finishing at the bottom centre. On the next field, the spot starts at the top centre and ends at the bottom right, the lines being **interlaced** with the lines of the previous field. The two fields together make a **frame** of 625 lines in $1/25$ s. A single line is completed in $64\mu\text{s}$ which corresponds to a **horizontal scan frequency** of $f_h = 15625$ Hz. The **vertical field frequency** is $f_v = 50$ Hz and the **frame rate** is $f_f = 25$ Hz. The frame rate was chosen to give the illusion of continuous motion while the field rate makes the flicker imperceptible.

The intensity of the scene is encoded as a **luminance** signal $x_Y(t)$. Between adjacent lines and fields, synchronization information is added to allow the scan generators in the receiver to lock onto those in the transmitter. The resulting signal is called the **composite video** signal, and waveform of a single line is as shown in figure C.2.

In a monochrome television system, the composite video signal is amplitude modulated onto the VHF or UHF carrier. A vestigial sideband (VSB) system is used in which the carrier is transmitted but the lower sideband is partially suppressed to reduce the overall bandwidth of the signal. Since the composite video signal is purely real, this does not cause any loss of information.

5.9.1 The spectrum of the composite video signal

It is relatively easy to analyze the spectrum of a composite video signal when there is no motion in the scene. Instead of retraced scanning, we consider a periodically repeated image in two dimensions and a single unbroken scanning path as shown in figure C.3. Each new scan line corresponds to the scanning path entering the image in the next column while each new field corresponds to the

scanning path entering the image in the next row. The synchronization information may be thought of as occurring in the borders around each image.

At time t the coordinates of the scanning spot are

$$x = s_x t, \quad y = s_y t \quad (5.54)$$

where $s_x = f_h H$ and $s_y = f_v V$ so that the spot traverses a horizontal distance of $f_h H$ and a vertical distance of $f_v V$ each second.

Let $I(x, y)$ denote the distribution of intensity. Since this is a two-dimensional periodic image, it is expressible as a two-dimensional Fourier series

$$I(x, y) = \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} c_{kl} \exp \left[j2\pi \left(\frac{kx}{H} + \frac{ly}{V} \right) \right] \quad (5.55)$$

where the Fourier coefficients c_{kl} are given by

$$c_{kl} = \frac{1}{HV} \int_0^H dx \int_0^V dy I(x, y) \exp \left[-j2\pi \left(\frac{kx}{H} + \frac{ly}{V} \right) \right] \quad (5.56)$$

In general, we expect c_{kl} to fall to zero as k and l become large as these correspond to very rapid variations of intensity across the image (i.e., by the Riemann-Lebesgue lemma). The luminance signal is given by substituting (5.54) into (5.55)

$$x_Y(t) = I(s_x t, s_y t) = \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} c_{kl} \exp \left[j2\pi \left(\frac{kx}{H} + \frac{ly}{V} \right) \right] \quad (5.57)$$

$$= \sum_{k=-\infty}^{\infty} \sum_{l=-\infty}^{\infty} c_{kl} \exp [j2\pi (k f_h + l f_v) t] \quad (5.58)$$

Thus $x_Y(t)$ is a sum of complex exponentials at the discrete frequencies $k f_h + l f_v$ where k and l are integers. This leads to a spectrum consisting of clusters of lines spaced by f_v around multiples of f_h since $f_v \ll f_h$ (Figure C.4).

When we consider a moving picture, the lines broaden slightly, but not by much since the timescale of changes in the picture is much slower than the line and field rates. Note that if any spectral lines are present which are not of the form $k f_h + l f_v$, they will necessarily represent a non-stationary pattern on the screen. The significant spectral lines for a broadcast-quality television picture extend to about 5.5 MHz.

5.9.2 The NTSC colour system

This system is used in the USA, Japan, Canada and Mexico. It was developed in 1954 as a fully compatible colour TV signal which fits into the bandwidth of a monochrome signal. The PAL system used in Germany, Britain, Australia and New Zealand is an improvement of the NTSC system. We shall discuss the NTSC-4.43 system based on a 50 Hz field frequency rather than the NTSC-3.58 system based on the 60 Hz field frequency actually used in the US.

The human visual response to any colour can be simulated using a mixture of red, green and blue lights. A straightforward approach to colour television would be to send three separate video signals $x_R(t)$, $x_G(t)$ and $x_B(t)$, one for each primary colour. (Assume without loss of generality that each of these is of modulus no greater than one). This would occupy a much greater bandwidth than a monochrome signal and would not be compatible with preexisting monochrome equipment.

To simulate the luminance signal of a monochrome camera, the colour signals are combined according to

$$x_Y(t) = 0.30x_R(t) + 0.59x_G(t) + 0.11x_B(t) \quad (5.59)$$

the coefficients having been chosen to match the different sensitivities of the eye to the primary colours. The spectrum of this luminance signal is given by the analysis above. It consists of clusters of spectral lines around multiples of the line frequency separated by gaps in which the spectrum is almost zero. The key idea of the compatible colour systems is to fit the colour information in the gaps between the clusters in such a way that this does not cause an objectional patterning on a monochrome receiver. Subjectively, such additional frequency components do not produce objectionable effects if they are of sufficiently high frequency that the intensity fluctuations introduced are fine-grained and change rapidly between successive lines or fields so that their effects tend to average out. In the NTSC system, two additional signals x_I and x_Q are defined as

$$x_I(t) = 0.60x_R(t) - 0.28x_G(t) - 0.32x_B(t) \quad (5.60)$$

$$x_Q(t) = 0.21x_R(t) - 0.52x_G(t) + 0.31x_B(t) \quad (5.61)$$

These are called the **chrominance** signals. Given x_I , x_Q and x_Y it is easy to regenerate the signals x_R , x_G and x_B .

Another fact about the human visual system which is exploited in colour television is that the eye is much more sensitive to resolution in the luminance information than in the chrominance information. In fact, before colour photography was widely used, pseudo-colour photographs were made by hand-painting black and white photographs. If the original black and white photograph was clear and in-focus, the fact that the colours were not applied very carefully did not matter very much. This means that it is possible to restrict the bandwidth of the chrominance signals x_I and x_Q quite severely (to about 1–2 MHz) compared with the 5.5 MHz bandwidth of x_Y without significant perceived loss of picture quality. (Recall that restricting the bandwidth in frequency corresponds to convolving in time with the impulse response of a filter whose temporal width is inversely dependent on the bandwidth. The narrower the bandwidth in frequency, the more a sharp transition is smeared out by this convolution. Similarly for a picture, restricting the bandwidth corresponds to defocusing the camera and smearing out fine detail).

The signals x_I and x_Q are band-limited and quadrature modulated onto a **colour subcarrier** signal whose frequency f_s is chosen to lie exactly at 283.5 times the line frequency so that $f_s = 283.5 f_h = 4.4296875$ MHz. The sideband structure of this quadrature modulated signal also consists of clusters spaced by f_h around f_s (the same analysis as for x_Y applies except that the modulating signal $x_I + jx_Q$ is complex and so the sidebands are not Hermitian symmetric about f_s). By choosing f_s half way between two multiples of f_h , the chrominance sidebands interleave between the luminance sidebands with minimal interference.

Since quadrature modulation requires the receiver to have an accurate copy of the carrier in order to demodulate the two channels of chrominance information, a **colour burst** consisting of about ten cycles of the colour subcarrier is included in the transmitted signal immediately after the line synchronization pulse. The receiver has a phase-locked loop which is locked precisely to the subcarrier frequency and phase and this is used to reconstitute x_I and x_Q . The presence of the colour burst signal is used by colour receivers to distinguish between colour transmissions and monochrome transmissions so that compatibility is achieved.

The main problem with the NTSC system is that various propagation conditions can cause the phase of the decoded chrominance information to be shifted relative to the original. The net effect of this is to shift all the hues of the picture to incorrect values. NTSC receivers usually have a hue control to allow manual adjustment of the subcarrier phase so that the colours look natural.

5.9.3 The PAL (Phase Alternate Line) system

In the PAL system, the chrominance signals are defined as slightly different linear combinations of the primary colours (Table 5.1, Figure 5.2)

$$x_V(t) = 0.877(x_R(t) - x_Y(t)) \quad (5.62)$$

$$x_U(t) = 0.493(x_B(t) - x_Y(t)) \quad (5.63)$$

These signals are used to form a complex modulating signal and this is quadrature modulated onto a colour subcarrier just as in the NTSC system. Unlike the NTSC system however, the complex modulating signal is

$$m(t) = \begin{cases} x_U(t) + jx_V(t) & \text{if } t \text{ is within an odd numbered line} \\ x_U(t) - jx_V(t) & \text{if } t \text{ is within an even numbered line} \end{cases} \quad (5.64)$$

This phase-reversal in $x_V(t)$ on alternate lines gives the system its name. If propagation conditions introduce a phase shift, the phase-reversal on alternate lines cause opposite colour shifts which tend to cancel each other out (Figures 5.4, 5.5).

It is easy to determine the effect of the phase alternation on the spectrum of the chrominance signal. Using a similar construction to that used for finding the spectrum of the monochrome signal we consider a two-dimensional infinite periodic repetition of the image. The phase reversals make adjacent columns of images no longer identical and so the effective horizontal period is $2H$ rather than H . Repeating the analysis shows that instead of clusters of spectral lines around multiples of f_h , the clusters are separated by $\frac{1}{2}f_h$. It is no longer appropriate to place the subcarrier frequency half-way between two clusters of the luminance signal as in the NTSC system since the closer chrominance cluster frequency spacing would cause major chrominance sidebands to interfere with the luminance. Instead, the subcarrier is shifted to approximately $283.75 f_h$. A more detailed analysis of the sideband structure indicates that the subcarrier frequency should be of the form $(n \pm \frac{3}{8})f_v$ for some integer n . Applying this small correction leads to the PAL subcarrier frequency of $283.75 f_h + 0.5 f_v = 4.43361875 \text{ MHz}$.

As in the NTSC system, a colour burst is included to allow the receiver to reconstruct the subcarrier required for quadrature demodulation of the chrominance information. This colour burst is deliberately shifted by $\pm 45^\circ$ relative to the true subcarrier phase on alternate lines so that the receiver can distinguish between odd and even lines. The PAL colour burst is thus called a **swinging burst**.

References

Carlson, A.B., Communication Systems, (second edition) 1975, McGraw-Hill.

King, G.J., Beginner's Guide to Colour Television, (second edition) 1973, Newnes-Butterworths.