

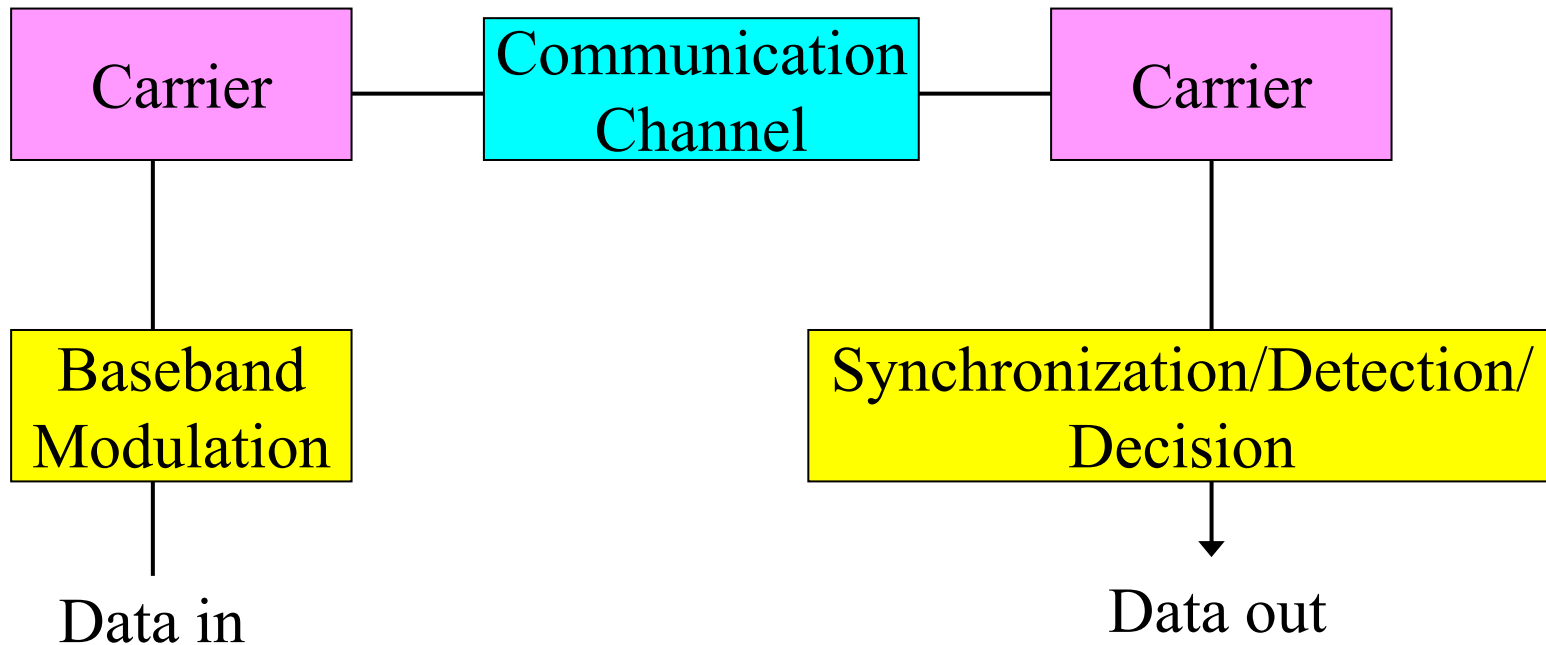
Part 3. Digital Modulation

What's Modulation & Demodulation?

➤ Digital modulation and demodulation:

- Modulation (demodulation) maps (retrieves) the digital information into (from) an analog waveform *appropriate for transmission over the channel*.
- Generally involve translating (recovering) the baseband digital information to (from) a bandpass analog signal at a carrier frequency that is very high compared to the baseband frequency.
- Examples: ASK, FSK, QPSK, 16QAM

Modulation & Demodulation



Why Carrier?

- *Effective radiation of electromagnetic waves* requires antenna dimensions comparable with the signal's wavelength:
 - Antenna for 3 kHz carrier would be ~100 km long
 - Antenna for 3 GHz carrier is 10 cm long
- *Frequency division multiplexing*
 - Shifting the baseband signals to different carrier frequencies
 - Sharing the communication channel resources

Geometric Representation (1)

- *Digital modulation involves choosing a particular analog signal waveform $s_i(t)$ from a finite set S of possible signal waveforms based on the information bits applied to the modulator.*
- For binary modulation schemes, a binary information bit is mapped directly to a signal and S contains only 2 signals, representing 0 and 1.
- For M -ary modulations, S contains more than 2 signals and each represents more than a single bit of information. With a signal set of size M , it is possible to transmit up to $\log_2 M$ bits per signal.

Geometric Representation (2)

- Any element of set S , $S = \{s_1(t), s_2(t), \dots, s_M(t)\}$, can be represented as a point in a vector space whose coordinates are basis signals $\phi_j(t)$, $j=1, 2, \dots, N$, such that

$$\int_{-\infty}^{\infty} \phi_i(t) \phi_j(t) dt = 0, i \neq j; (\rightarrow \text{orthogonal})$$

$$E = \int_{-\infty}^{\infty} [\phi_i(t)]^2 dt = 1; (\rightarrow \text{normalization})$$

$s_i(t)$ can be represented as a linear combination of the basis signals.

$$s_i(t) = \sum_{j=1}^N s_{ij} \phi_j(t), \quad i = 1, 2, \dots, M$$

Example: BPSK Geometric Representation

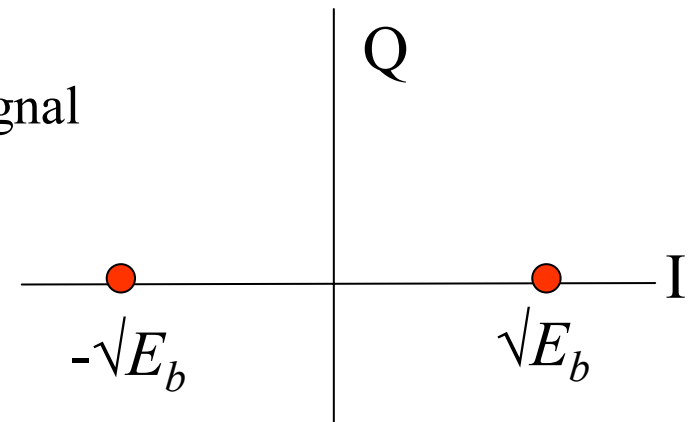
$$S_{BPSK} = \left\{ \begin{array}{l} \left[s_1(t) = \sqrt{\frac{2E_b}{T_b}} \cos(2\pi f_c t + \theta_c) \right], \\ \left[s_2(t) = -\sqrt{\frac{2E_b}{T_b}} \cos(2\pi f_c t + \theta_c) \right]; \end{array} \right. ; 0 \leq t \leq T_b$$

E_b = energy per bit; T_b = bit period

For this signal set, there is a single basis signal

$$\phi_1(t) = \sqrt{\frac{2}{T_b}} \cos(2\pi f_c t + \theta_c); 0 \leq t \leq T_b$$

$$S_{BPSK} = \left\{ \left[\sqrt{E_b} \phi_1(t) \right], \left[-\sqrt{E_b} \phi_1(t) \right] \right\}$$



Constellation diagram

Constellation Diagram

- *A graphical representation of the complex envelope of each possible signal*
- The x-axis represents the in-phase component and the y-axis represents the quadrature component of the complex envelope
- The distance between signals on a constellation diagram relates to *how different the modulation waveforms are* and *how well a receiver can differentiate between them* when random noise is present.

Performance Measures (1)

- *Two key performance measures of a modulation scheme are power efficiency and bandwidth efficiency*
- **Power efficiency** is a measure of how favorably the tradeoff between fidelity and signal power is made, and is expressed as the ratio of the signal energy per bit (E_b) to the noise PSD (N_0) required to achieve a given probability of error (say 10^{-5}):

$$\eta_p = \frac{E_b}{N_0} \quad \text{Small } \eta_p \text{ is preferred}$$

Performance Measures (2)

- **Bandwidth efficiency** describes the ability of a modulation scheme to accommodate data within a limited bandwidth, In general, it is defined as the ratio of the data bit rate R to the required RF bandwidth B :

$$\eta_B = \frac{R}{B} \text{ (bps/Hz)} \quad \text{Large } \eta_B \text{ is preferred}$$

- **Channel capacity** gives an upper bound of achievable bandwidth efficiency:

$$\eta_{B \max} = \frac{C}{B} = \log_2 \left(1 + \frac{S}{N} \right)$$

Modulation Schemes Classification

- **Linear modulation:** the amplitude of the transmitted signal, $s(t)$, varies linearly with the modulating digital signal, $m(t)$.
 - Bandwidth efficient but power inefficient
 - Example: ASK, QPSK
- **Nonlinear modulation:** the amplitude of the transmitted signal, $s(t)$, does not vary linearly with the modulating digital signal
 - Power efficient but bandwidth inefficient
 - Example: FSK, constant envelope modulation

Demodulation (1)

- **Coherent demodulation:** requires a replica carrier wave of the same frequency and phase at the receiver
 - The received signal and replica carrier are cross-correlated
 - Also known as synchronous demodulation
 - Carrier recovery methods
 - Using PLL to recover the carrier phase and frequency from the transmitted pilot carrier signal,
 - Recovering the carrier from the received signals using costas loop
 - Applicable to: PSK, FSK, ASK, etc.

Demodulation (2)

➤ *Example of BPSK coherent demodulator:*

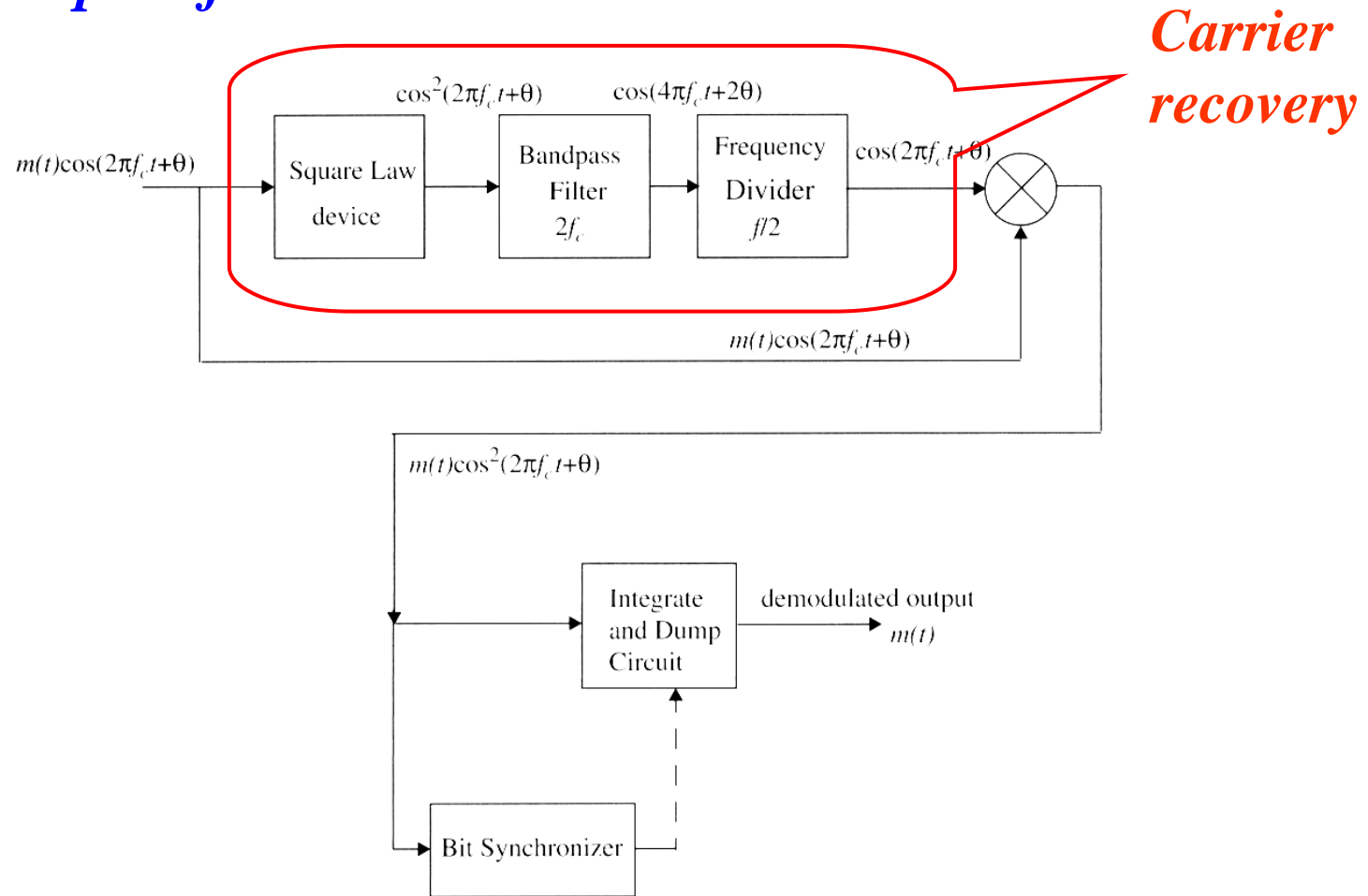


Figure 6.23 BPSK receiver with carrier recovery circuits.

Demodulation (3)

- **Non-coherent demodulation:** does not require a reference carrier wave
 - It is less complex than coherent demodulation (easier to implement), but has worse performance
 - Applicable to: DPSK, FSK, etc.
 - Example: FSK non-coherent demodulator

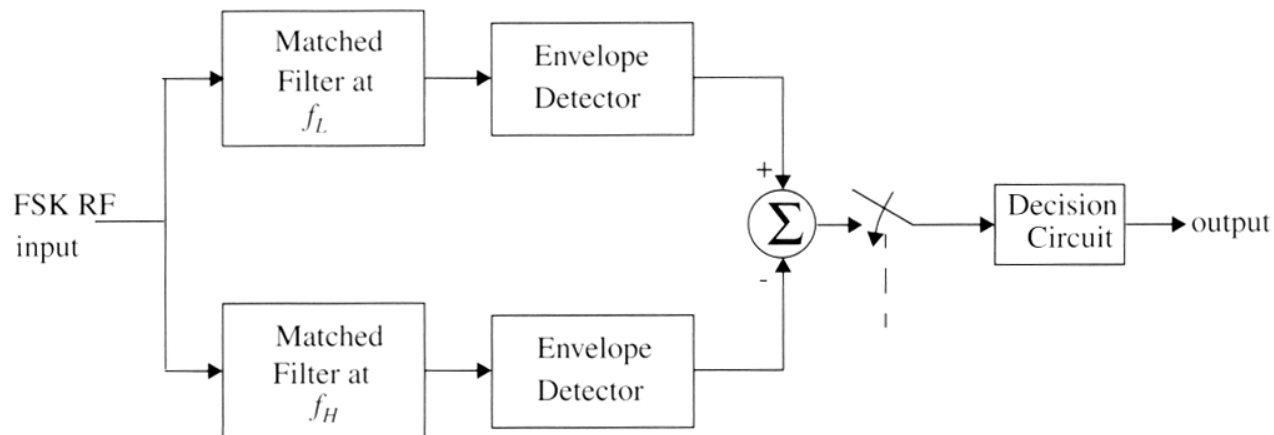


Figure 6.37 Block diagram of noncoherent FSK receiver.

Part 3.1 Basic Modulation

Modulation Process

$$f = f(a_1, a_2, a_3, \dots, a_n, t) \quad (\rightarrow \text{carrier})$$

$$a_1, a_2, a_3, \dots, a_n \quad (\rightarrow \text{modulation parameters})$$

$$t \quad (\rightarrow \text{time})$$

- Modulation implies varying one or more characteristics (modulation parameters a_1, a_2, \dots, a_n) of a carrier f in accordance with the information-bearing (modulating) baseband signal.
- Sinusoidal waves, pulse train, square wave, etc. can be used as carriers

Continuous Carrier

Carrier: $A \cos[\omega t + \varphi]$

- $A = \text{const}$
- $\omega = \text{const}$
- $\varphi = \text{const}$

□ **Amplitude modulation (AM)**

- $A = A(t)$ – carries information
- $\omega = \text{const}$
- $\varphi = \text{const}$

□ **Frequency modulation (FM)**

- $A = \text{const}$
- $\omega = \omega(t)$ – carries information
- $\varphi = \text{const}$

□ **Phase modulation (PM)**

- $A = \text{const}$
- $\omega = \text{const}$
- $\varphi = \varphi(t)$ – carries information

Modulation methods: using amplitude, phase or frequency of the carrier.

Basic Modulation

- Modulation involves operations on one or more of the three characteristics of a carrier signal: amplitude, frequency and phase.
- The three basic modulation methods are:
 - *Amplitude Shift Keying (ASK)*
 - *Phase Shift Keying (PSK)*
 - *Frequency Shift Keying (FSK)*
- These could be applied to binary or M-ary signals.
- There are other variants as well.

Amplitude Shift Keying (ASK) (1)

➤ The modulation signal set is

$$s_i(t) = A_i g(t) \cos[2\pi f_c t + \theta_c], \quad \begin{array}{l} i = 1, 2, \dots, M \\ 0 \leq t \leq T_s \end{array}$$

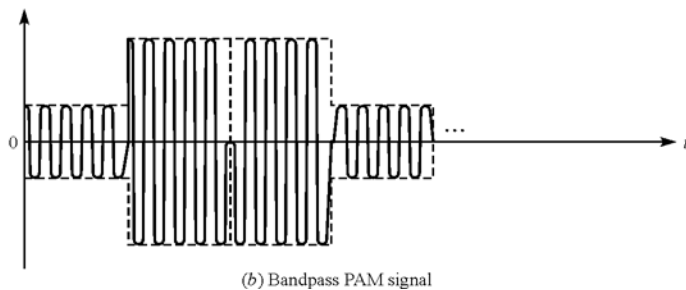
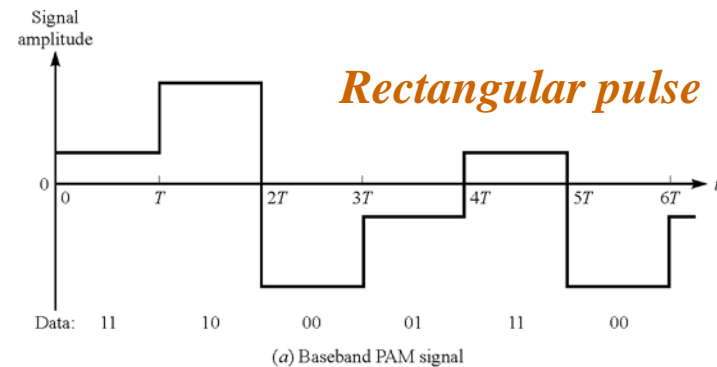
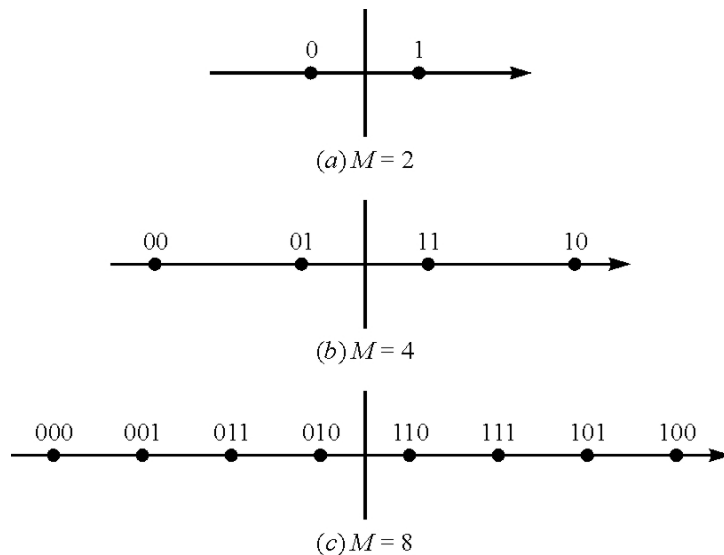
- T_s is the symbol period
- f_c is the carrier frequency, θ_c is the carrier initial phase
- $g(t)$ is a real-value signal pulse whose shape influences the spectrum of the transmitted signal; → **Pulse shaping**
 - Used to simultaneously reduce the intersymbol effects and the spectral width of a modulated digital signal
 - Example: rectangular pulse, Nyquist pulse shaping, raised cosine pulse shaping, Gaussian pulse shaping, etc.
- $A_i = (2i-1-M)d$, each symbol represents $\log_2 M$ bits

Amplitude Shift Keying (ASK) (2)

➤ The single basis signal is $\phi_1(t) = \sqrt{\frac{2}{\epsilon_g}} g(t) \cos[2\pi f_c t + \theta_c]$

➤ The modulated signal:

$$s_i(t) = s_i \phi_1(t), \quad s_i = A_i \sqrt{\epsilon_g / 2}$$



• *ASK demonstrates poor performance, as it is heavily affected by noise, fading, and interference. It is rarely used on its own.*

Phase Shift Keying (PSK)

➤ The modulation signal set is

$$\begin{aligned} s_i(t) &= A_c g(t) \cos[2\pi f_c t + \theta_c + \varphi_i], & i = 1, 2, \dots, M \\ \varphi_i &= \frac{2\pi}{M} (i - 1) & 0 \leq t \leq T_s \end{aligned}$$

- A_c is the carrier amplitude,
- φ_i carries information, each symbol represents $\log_2 M$ bits

Binary Phase Shift Keying (BPSK)

➤ $M=2$: minimum phase separation: 180°

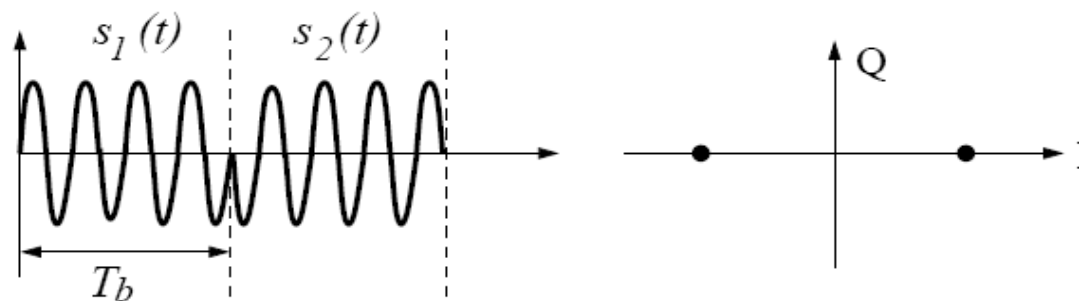
$$s_i(t) = A_c g(t) \cos[2\pi f_c t + \theta_c + (i-1)\pi], \quad i = 1, 2, \quad 0 \leq t \leq T_b$$

– $s_1(t)$ and $s_2(t)$ represent bit 0 and bit 1, respectively

– The single basis: $\phi_1(t) = \sqrt{\frac{2}{\mathcal{E}_g}} g(t) \cos[2\pi f_c t + \theta_c]$

– The set :

$$S = \left\{ A_c \sqrt{\frac{\mathcal{E}_g}{2}} \phi_1(t), -A_c \sqrt{\frac{\mathcal{E}_g}{2}} \phi_1(t) \right\}$$



Quadrature Phase Shift Keying (QPSK)

➤ $M=4$: symbol period $T_s=2T_b$, minimum phase separation: 90°

$$s_i(t) = A_c g(t) \cos \left[2\pi f_c t + \theta_c + (i-1) \frac{\pi}{2} \right], \quad i = 1, 2, 3, 4, \quad 0 \leq t \leq T_s$$

– The basis signals:

$$\phi_1(t) = \sqrt{\frac{2}{\epsilon_g}} g(t) \cos[2\pi f_c t + \theta_c], \quad \phi_2(t) = -\sqrt{\frac{2}{\epsilon_g}} g(t) \sin[2\pi f_c t + \theta_c]$$

– Constellation diagram:

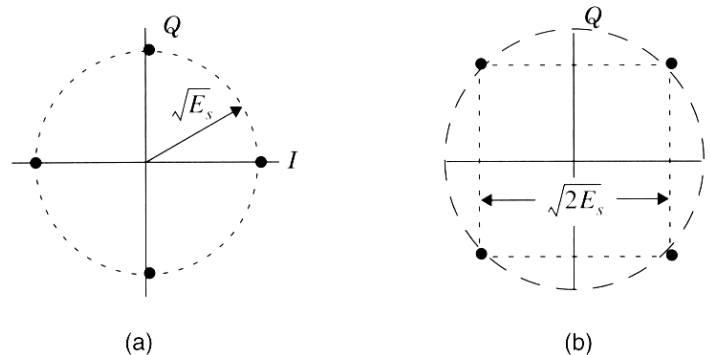


Figure 6.26 (a) QPSK constellation where the carrier phases are $0, \pi/2, \pi, 3\pi/2$; (b) QPSK constellation where the carrier phases are $\pi/4, 3\pi/4, 5\pi/4, 7\pi/4$.

PSK: Bandwidth vs. Power Efficiency

The system using ideal Nyquist pulse shaping is operated in AWGN channel.

Table 6.4 Bandwidth and Power Efficiency of M-ary PSK Signals

M	2	4	8	16	32	64
$\eta_B = R_b/B^*$	0.5	1	1.5	2	2.5	3
E_b/N_o for BER= 10^{-6}	10.5	10.5	14	18.5	23.4	28.5

* B : First null bandwidth of M-ary PSK signals

QPSK can be interpreted as two independent BPSK systems (one on the I-channel and the other on Q-channel), and thus *the same performance but twice the bandwidth efficiency.*

Quadrature Amplitude Modulation (QAM) (1)

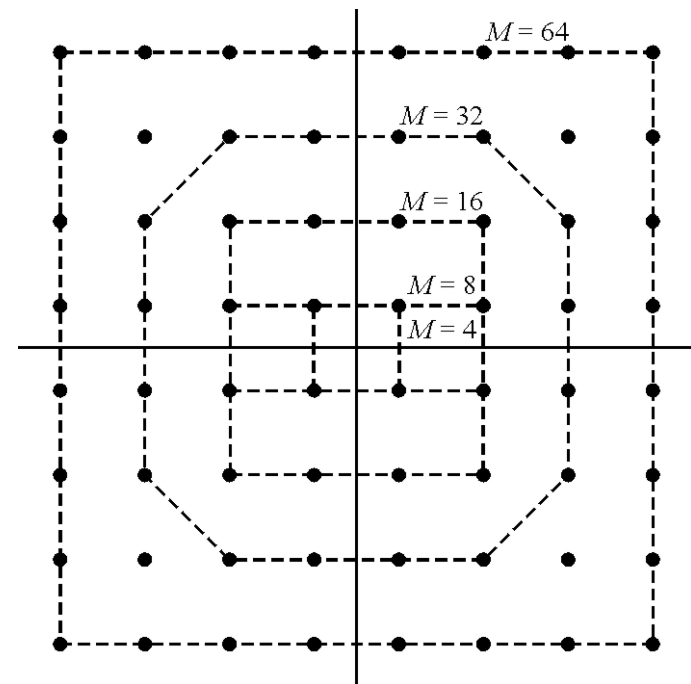
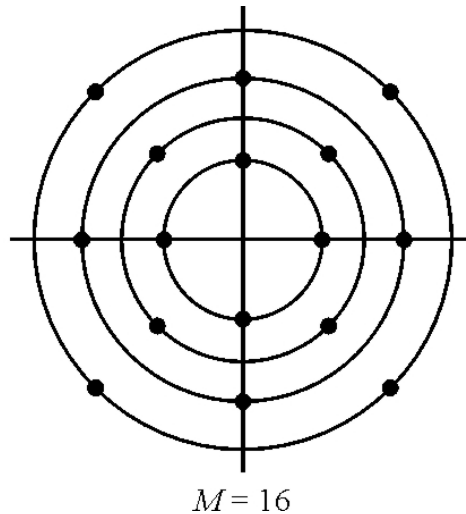
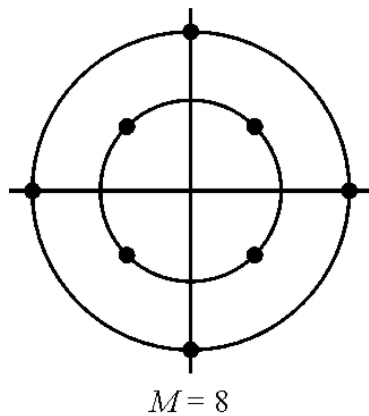
➤ *Combined amplitude/phase shift keying*

$$s_i(t) = A_i g(t) \cos \left[2\pi f_c t + \theta_c + \varphi_j \right],$$
$$i = 1, 2, \dots, M_1, j = 1, 2, \dots, M_2, 0 \leq t \leq T_s$$

- As both amplitude and phase are used to carry symbol information, it is very bandwidth efficient
- Signal set size $M=M_1M_2$: $2^1 \times 2^1 = 4$, $2^2 \times 2^2 = 16$, $2^3 \times 2^3 = 64$, etc 4QAM, 16QAM, 64QAM
- The larger M is, *the better bandwidth efficiency* but *lower robustness against noise and fading*

Quadrature Amplitude Modulation (QAM) (2)

➤ *Examples of constellation:*



QAM: Bandwidth vs. Power Efficiency

The system using optimum raised cosine pulse shaping is operated in AWGN channel.

Table 6.5 Bandwidth and Power Efficiency of QAM [Zie92]

M	4	16	64	256	1024	4096
η_B	1	2	3	4	5	6
E_b/N_o for BER = 10^{-6}	10.5	15	18.5	24	28	33.5

In terms of power efficiency, QAM is superior to M-ary PSK.

Frequency Shift Keying (FSK)

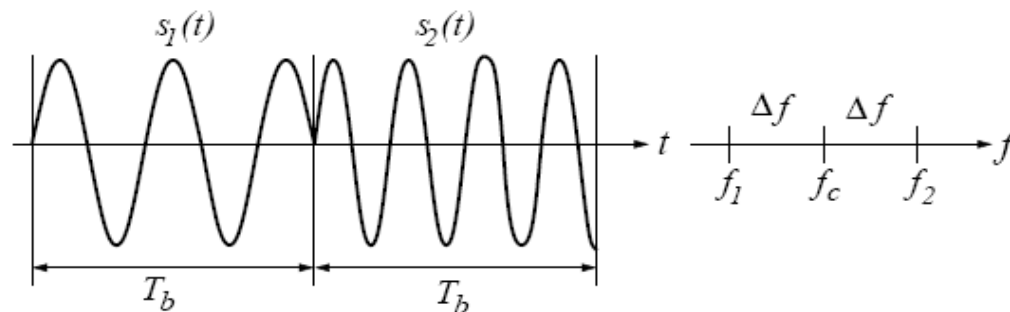
➤ The modulation signal set is

$$s_i(t) = A_c \cos[2\pi f_i t + \theta_c], \quad i = 1, 2, \dots, M, 0 \leq t \leq T_s$$

➤ **BFSK**: $M=2$

– Bit 0: $f_1 = f_c - \Delta f$, $s_1(t) = \sqrt{\frac{2\varepsilon_b}{T_b}} \cos(2\pi(f_c - \Delta f)t + \theta_c)$

– Bit 1: $f_2 = f_c + \Delta f$, $s_2(t) = \sqrt{\frac{2\varepsilon_b}{T_b}} \cos(2\pi(f_c + \Delta f)t + \theta_c)$



FSK: Bandwidth vs. Power Efficiency

Table 6.6 Bandwidth and Power Efficiency of Coherent M-ary FSK [Zie92]

M	2	4	8	16	32	64
η_B	0.4	0.57	0.55	0.42	0.29	0.18
E_b/N_o for BER = 10^{-6}	13.5	10.8	9.3	8.2	7.5	6.9

•Nonlinear modulation: *bandwidth inefficient but power efficient, no need for expensive linear amplifiers*

Other Modulations

➤ *Differential phase shift keying (DPSK)*

- The input binary sequence is differentially encoded before BPSK modulation ($d_k = m_k \oplus d_{k-1}$)
- Avoids the need for a coherent reference signal at the receiver

➤ *Offset QPSK*

- The phase transitions are limited to 90^0 , the transitions on the I and Q channels are staggered.

➤ */4 QPSK*

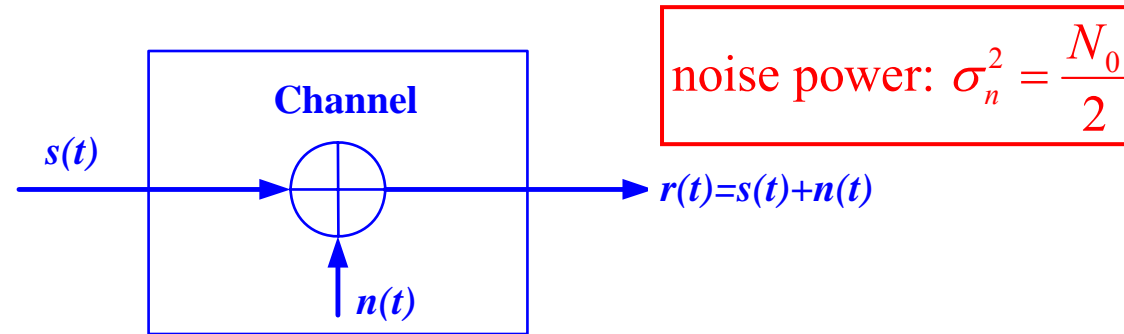
- The phase transitions are limited to 135^0

➤ *Continuous-phase FSK (CPFSK)*

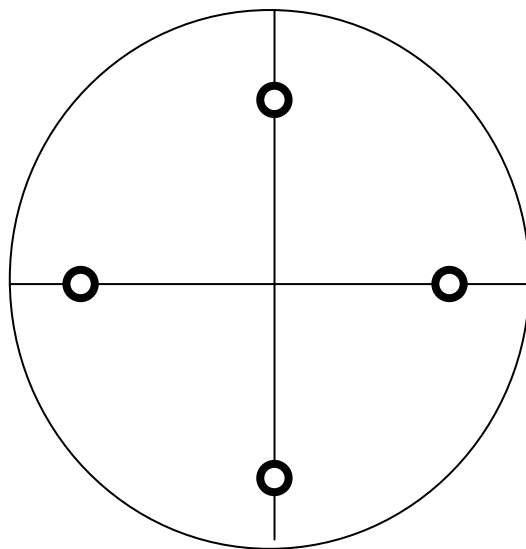
- Avoids sudden change in the signal frequency, i.e., large spectral side lobes outside of the main spectral band
- Minimum shift keying (MSK), Gaussian MSK(GMSK)

Performance in AWGN Channel (1)

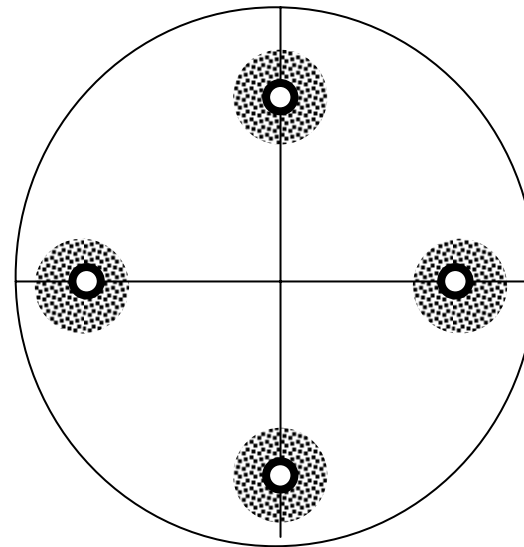
- The channel is assumed to corrupt the signal by the additive white Gaussian noise.



- *Distortion*



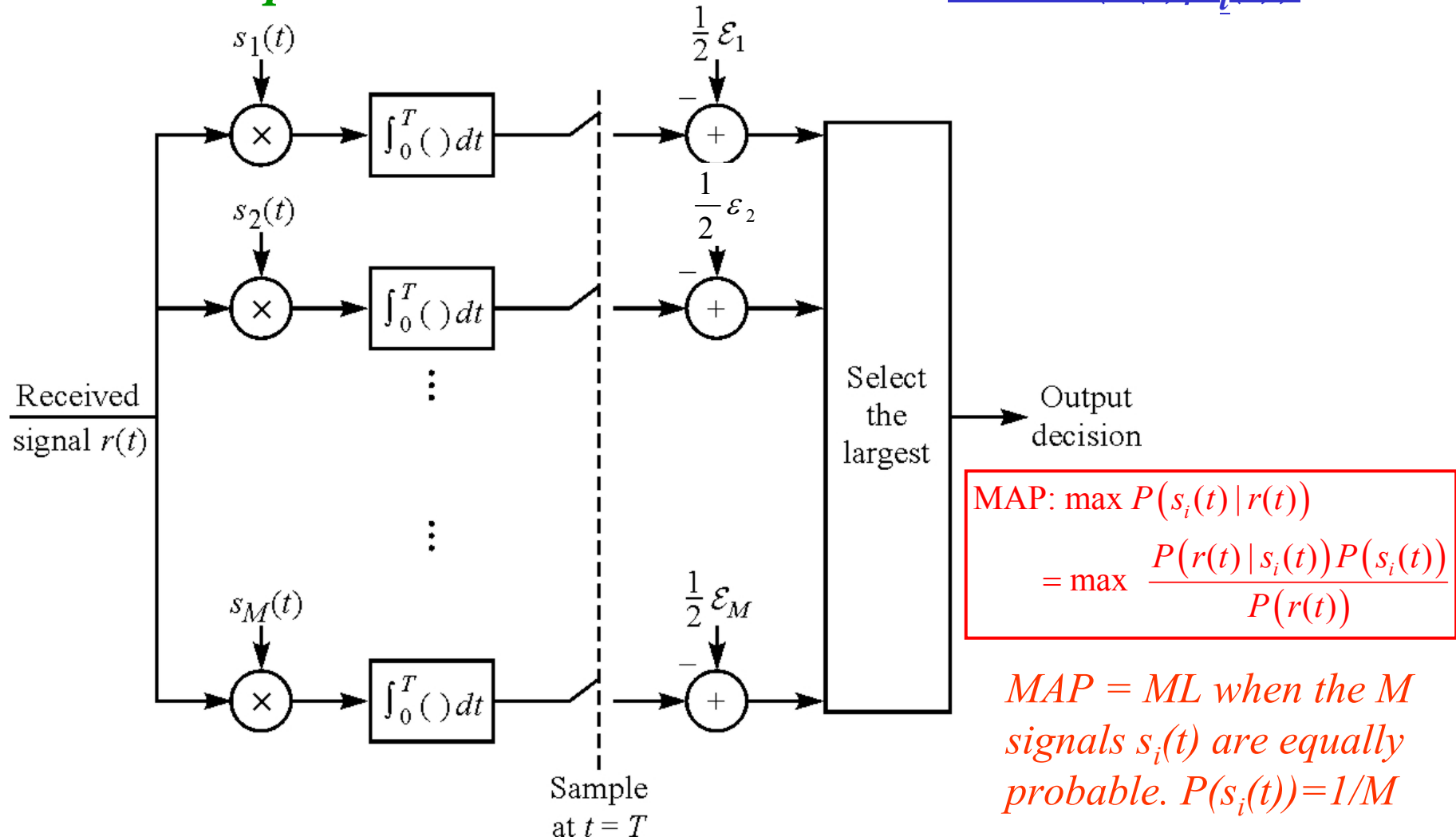
Perfect channel



White noise

Performance in AWGN Channel (2)

➤ *The optimum ML AWGN receiver: $\max P(r(t)/s_i(t))$*

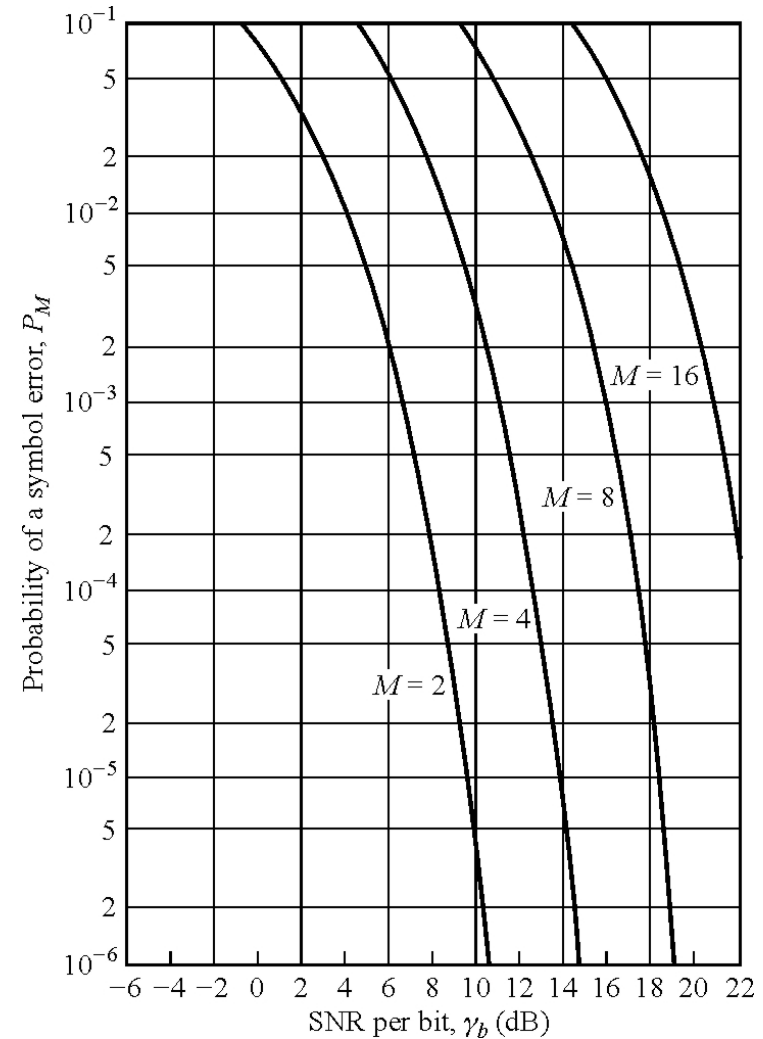


Performance in AWGN Channel (3)

➤ **ASK:** symbol error probability

$$P_M = \frac{2(M-1)}{M} Q \left(\sqrt{\frac{(6 \log_2 M) \mathcal{E}_{b,av}}{(M^2-1) N_0}} \right)$$

$\mathcal{E}_{b,av}$ is the average bit energy



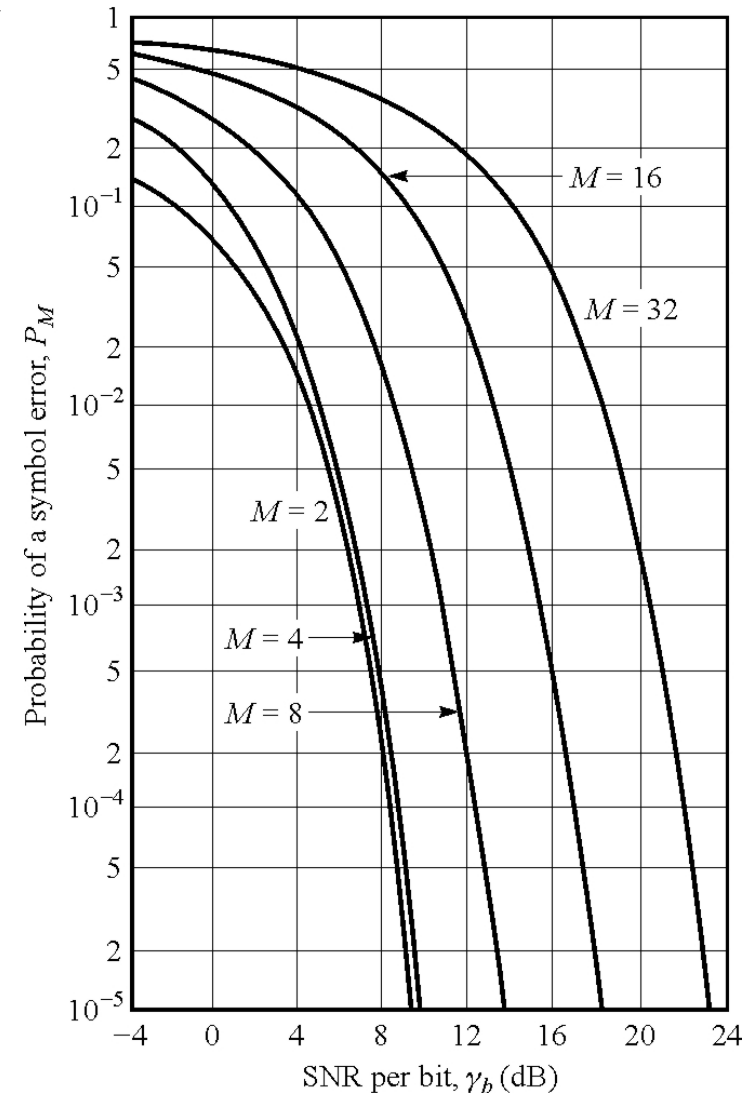
Performance in AWGN Channel (4)

➤ **PSK**: symbol error probability

$$P_M \approx 2Q\left(\sqrt{\frac{2\varepsilon_s}{N_0}} \sin\left(\frac{\pi}{M}\right)\right), \quad M > 4$$

$\varepsilon_s = (\log_2 M) \varepsilon_b$ is the symbol energy

$$\text{BPSK \& QPSK: } P_M = Q\left(\sqrt{\frac{2\varepsilon_b}{N_0}}\right)$$



Performance in AWGN Channel (5)

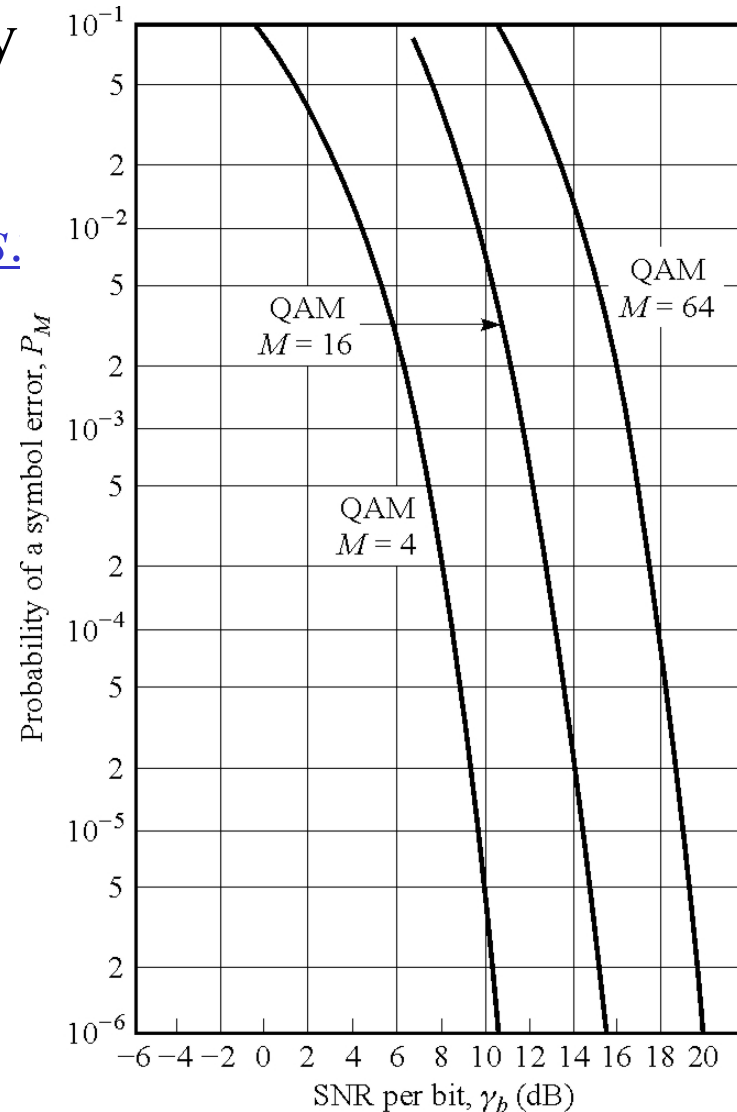
➤ **QAM**: symbol error probability

For rectangular QAM constellations.

$$P_M \leq 4Q\left(\sqrt{\frac{3\varepsilon_{s,av}}{(M-1)N_0}}\right),$$

The average symbol energy:

$$\varepsilon_{s,av} = (\log_2 M) \varepsilon_{b,av}$$

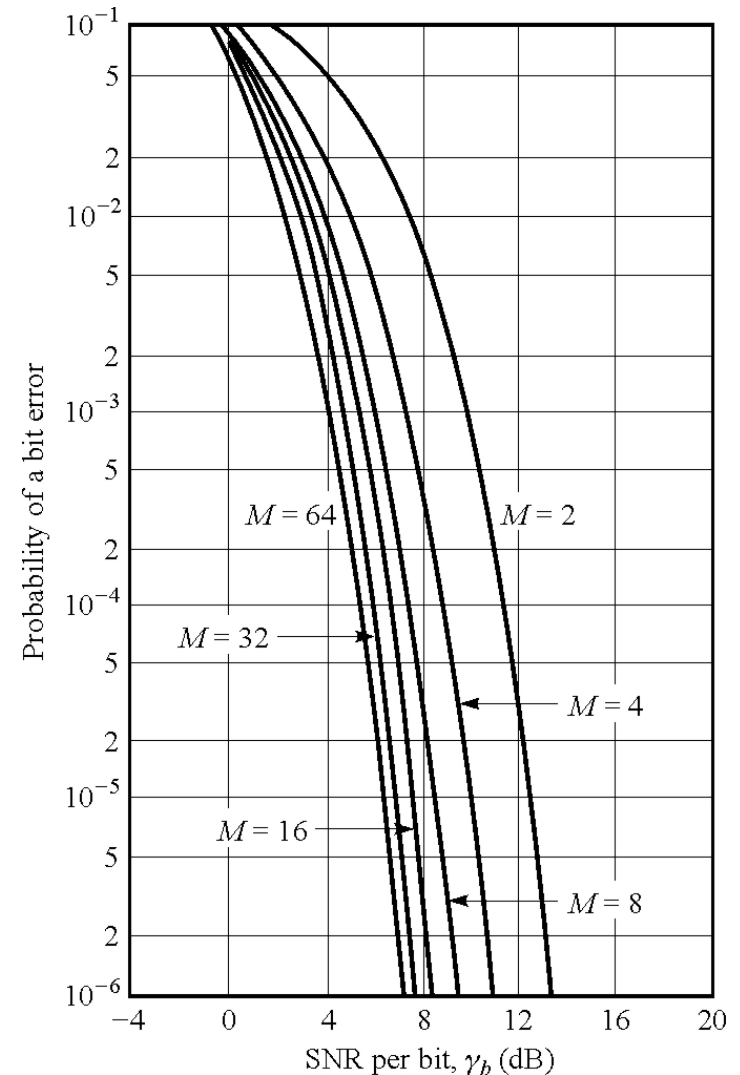


Performance in AWGN Channel (6)

➤ **FSK**: symbol error probability

$$P_M \leq (M-1)Q\left(\sqrt{\frac{\mathcal{E}_s}{N_0}}\right),$$

$$\mathcal{E}_s = (\log_2 M) \mathcal{E}_b$$



Part 3.2 Trellis Coded Modulation

Overview of TCM (1)

➤ Conventional coding

- Separate from modulation, performed at the digital level before modulation
- The insertion of redundant bits
 - Given the same information transmission rate, the symbol rate must be (n/k) times that of the uncoded system.
 - The redundancy provides coding gain, however, requires extra bandwidth.
- In a band-limited channel, the required additional bandwidth is unavailable.

Overview of TCM (2)

➤ **Solution:** Trellis coded modulation (TCM)

- The combination of coding and modulation
- Coding gain without expanding bandwidth
 - Using a constellation with more points than that required without coding
 - Typically, the number of points is doubled
 - The symbol rate is unchanged and the bandwidth remains unchanged.

Overview of TCM (3)

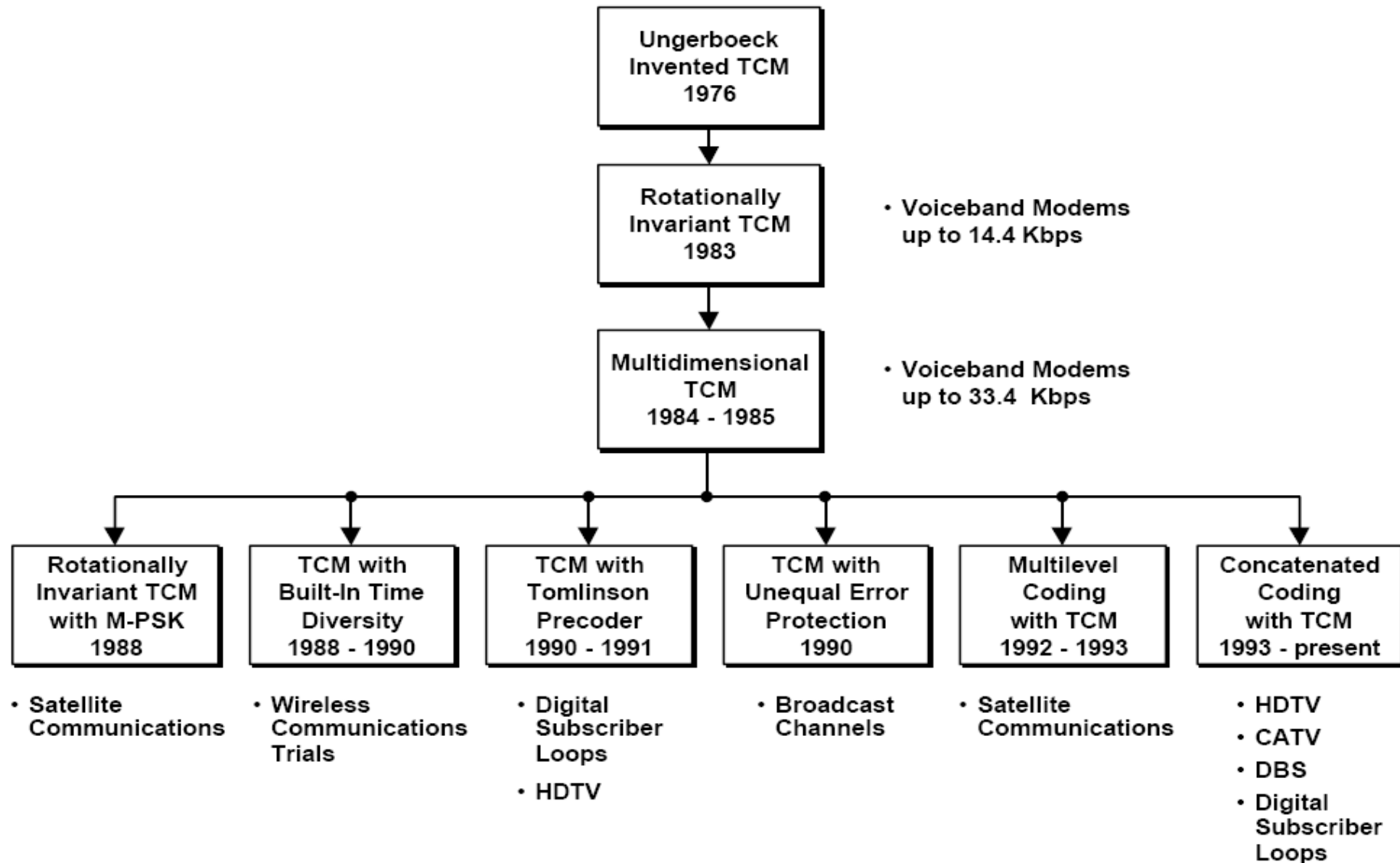
➤ *How to achieve the coding gain by TCM?*

- Introducing dependancy between every successive symbols
 - Only certain sequences of successive constellation points are allowed
- Maximizing the Euclidean distance between possible sequences of transmitted symbols
 - Minimum distance between the possible **sequences** of transmitted symbols in signal space (d_{min}) determines the performance:

$$P_e \sim e^{-d_{min}^2 / \sigma_n^2}$$

- It actually decreases the error probability for a given SNR, thus achieving coding gain

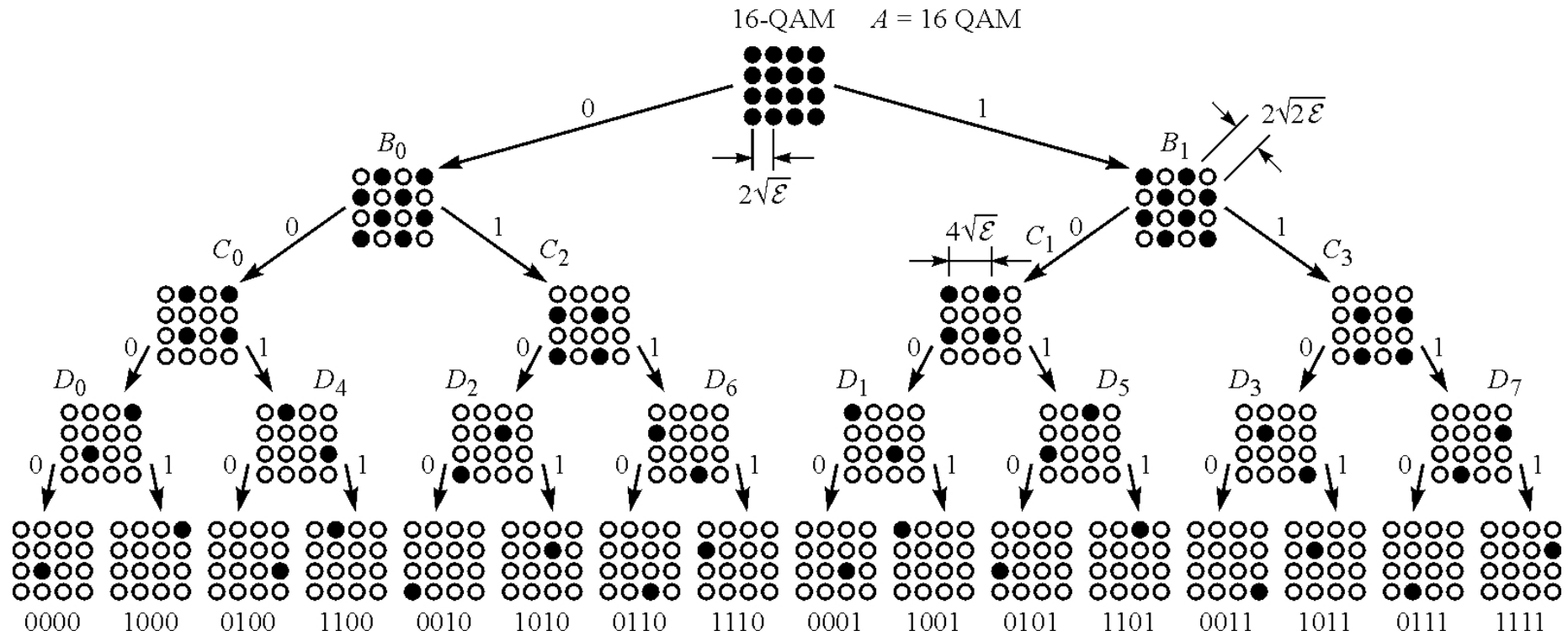
History of TCM



Basic Principles of TCM (1)

- TCM is to devise an effective method for mapping the coded bits into signal points such that the minimum Euclidean distance is maximized.
- Ungerboeck idea: mapping by set partitioning
 - The signal constellation is partitioned in a systematic manner to form a series of smaller subsets.
 - The resulting subsets have a larger minimum distance than their “parent”.
 - The goal of partitioning: each partition should produce subsets with increased minimum distance.

Example of Set Partitioning



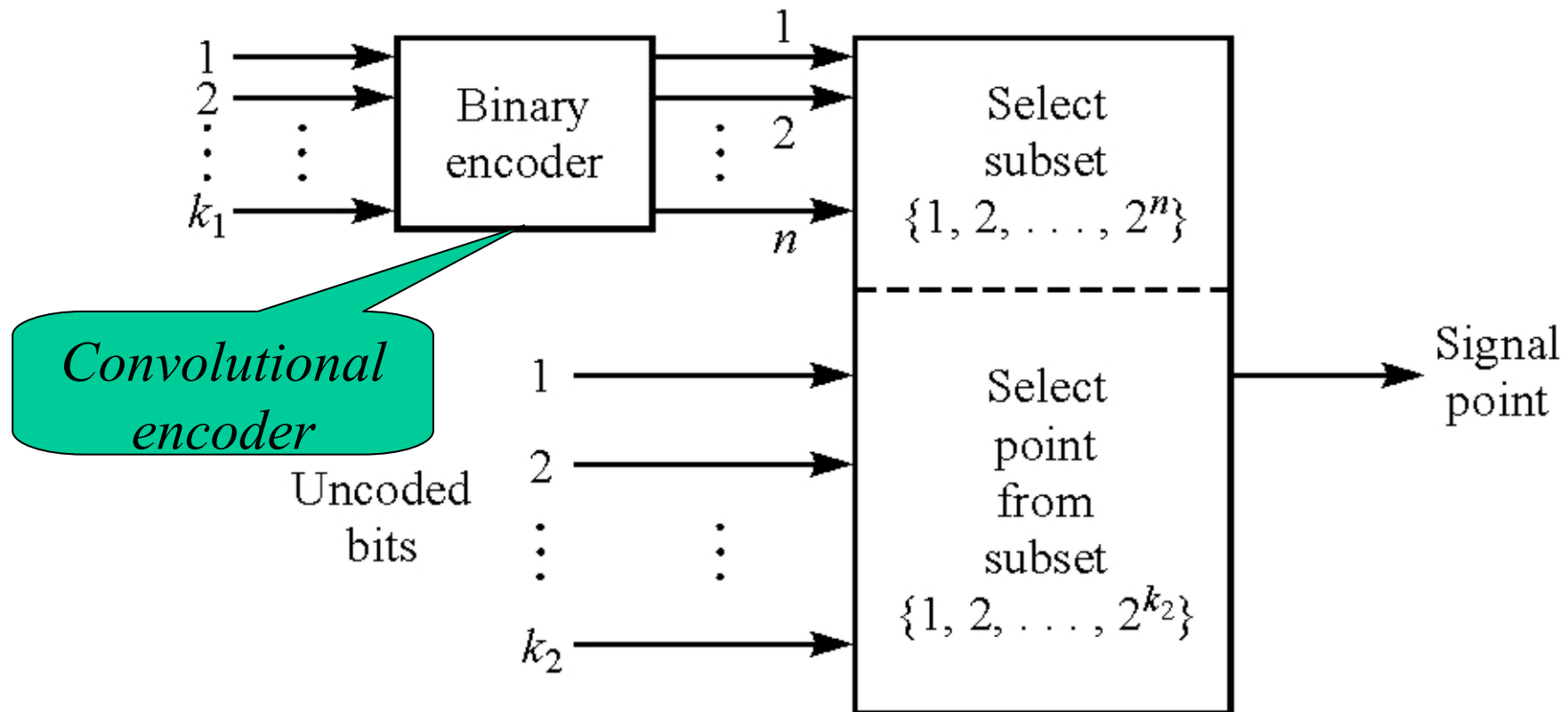
Basic Principles of TCM (2)

➤ *In general, the encoding is performed as follows:*

- A block of m information bits is separated into two groups of length k_1 and k_2 , respectively.
 - The k_1 bits are encoded into n bits, while the k_2 bits are left uncoded.
 - The n bits from the encoder are used to select one of the possible subsets in the partitioned signal set, while the k_2 bits are used to select one of 2^{k_2} signal points in each subset.
- The coder need not code all the incoming bits. When $k_2=0$, all m information bits are encoded.
- There are many ways to map the coded bits into symbols. The choice of mapping will drastically affect the performance of the code.

Basic Principles of TCM (3)

➤ General structure of encoder:



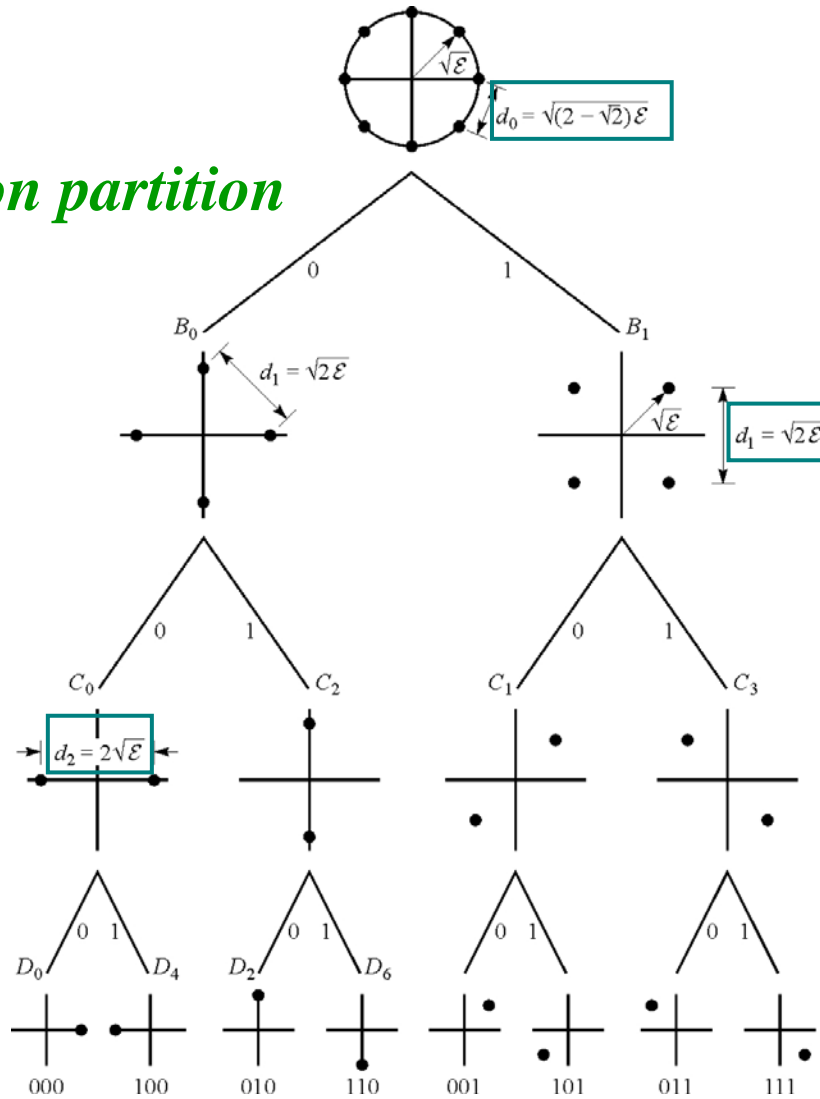
Basic Principles of TCM (4)

➤ *The basic rules for the assignment of signal subsets to state transitions in the trellis*

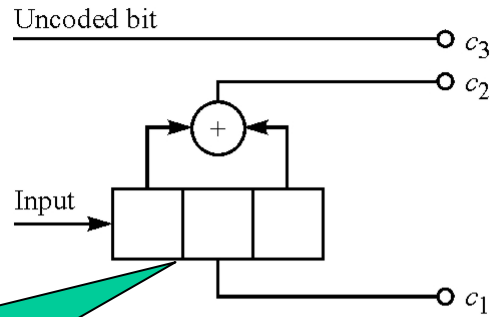
- Use all subsets with equal frequency in the trellis
- Transitions originating from the same state or merging into the same state in the trellis are assigned subsets that are separated by the largest Euclidean distance
- Parallel state transitions (when they occur) are assigned signal points separated by the largest Euclidean distance.
 - Parallel transitions in the trellis are characteristic of TCM that contains one or more uncoded information bits.

Examples of TCM (1)

8-PSK constellation partition

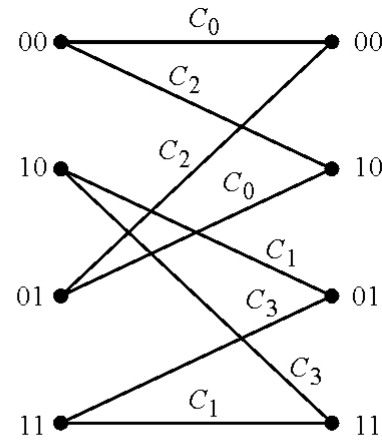


Examples of TCM (2)

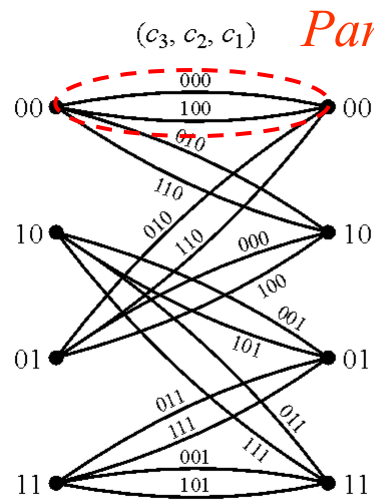
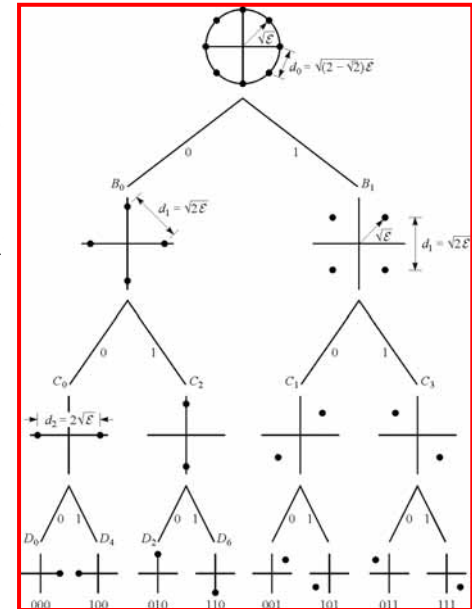


(a) Encoder

1/2 convolutional encoder with 4 states

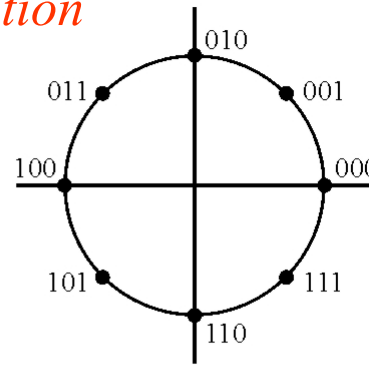


(b) Four-state trellis



(c) Mapping of bits to state transitions

Parallel transition

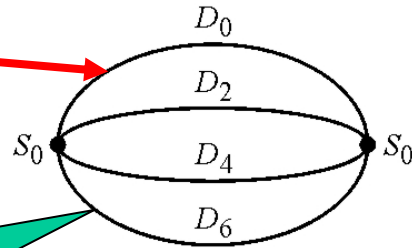


(d) Mapping of bits (c_3, c_2, c_1) to signal points corresponding to partition in Fig. 8.3-1 (note nonuniqueness of this mapping)

Examples of TCM (3)

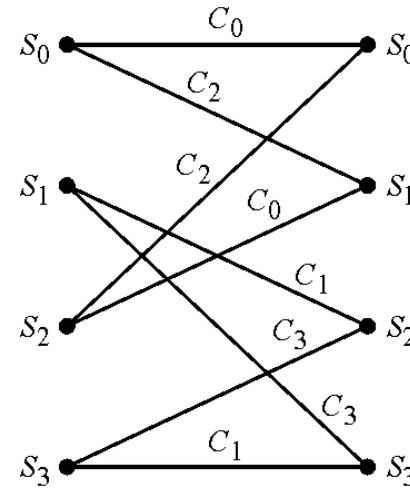
Minimum Euclidean distance:

$$d_{\min, \text{uncoded}} = \sqrt{2\varepsilon}$$



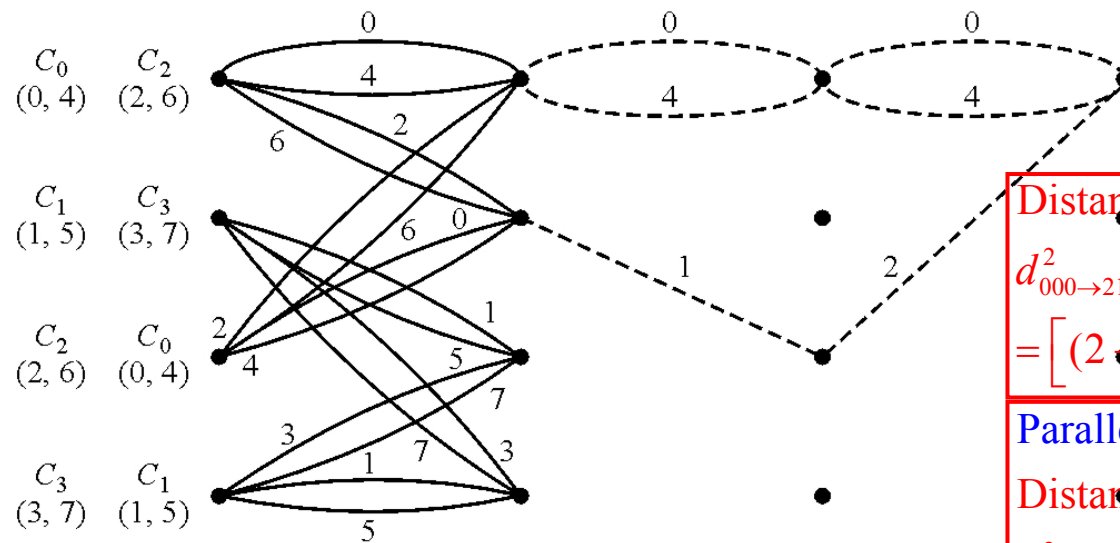
(a) One-state trellis

Uncoded QPSK



(b) Four-state trellis

Trellis coded 8PSK modulation

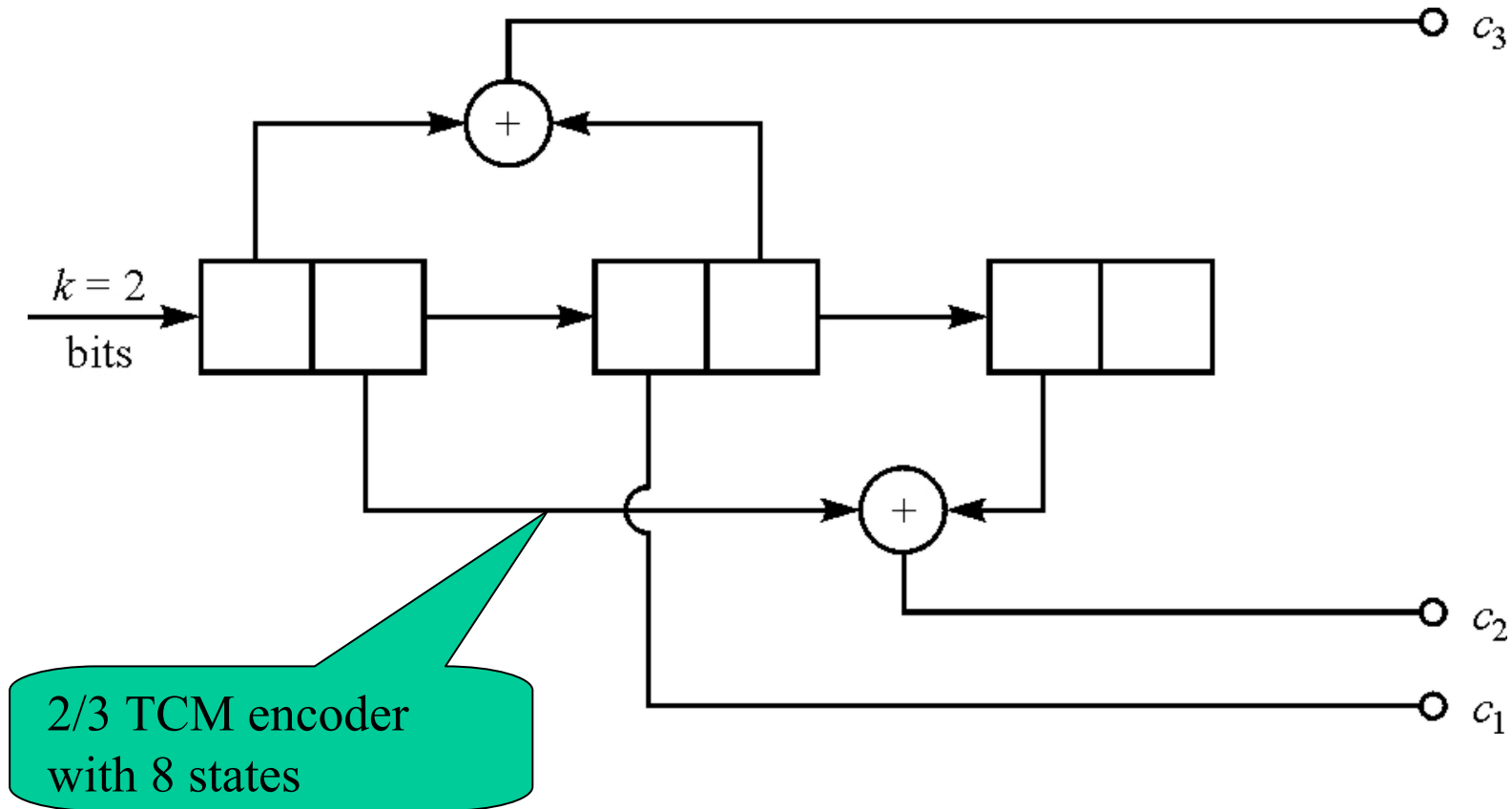


(c) Four-state trellis

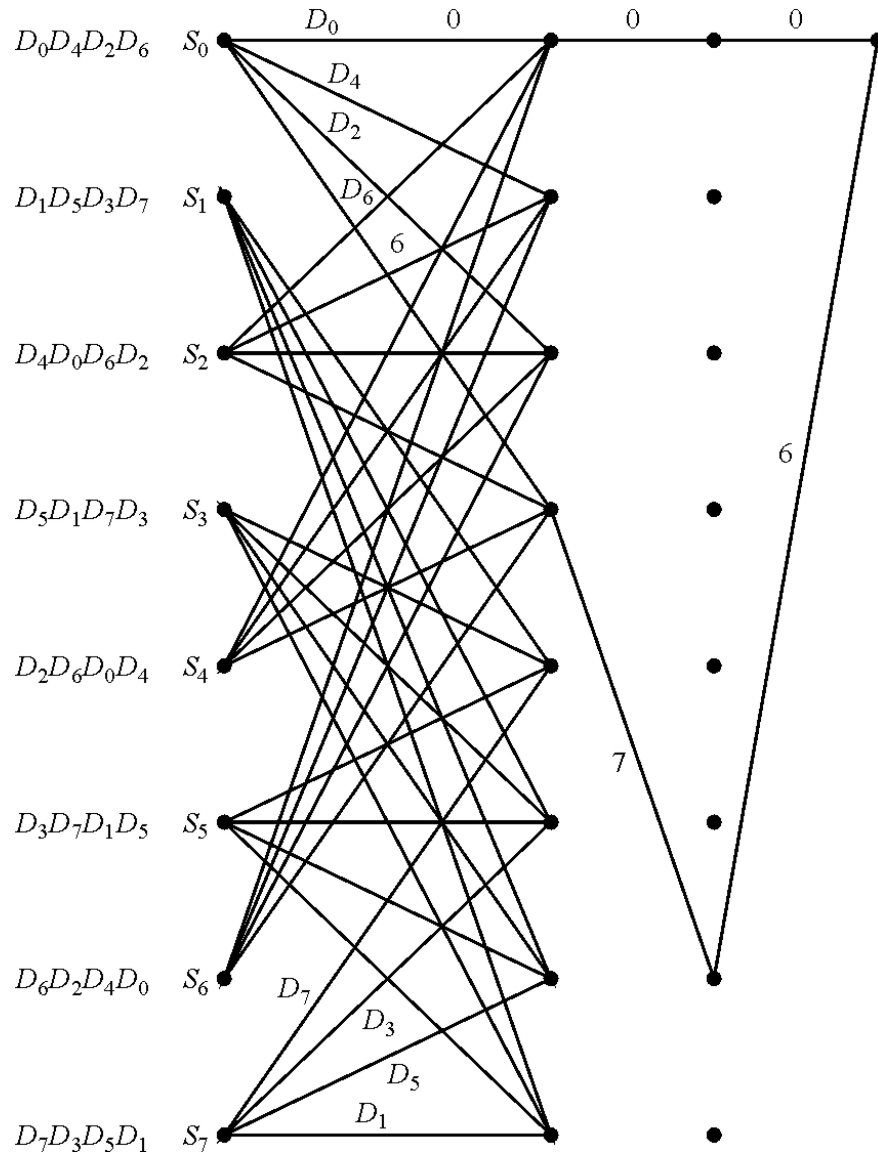
Distance : $(0, 0, 0) \rightarrow (2, 1, 2)$
 $d_{000 \rightarrow 212}^2 = d_0^2 + 2d_1^2$
 $= [(2 \cdot \sqrt{2})\varepsilon + 4\varepsilon] = 4.585\varepsilon$

Parallel transition:
 Distance : $(0, 0, 0) \rightarrow (4, 0, 0)$
 $d_{000 \rightarrow 400}^2 = d_2^2 = 4\varepsilon$

Examples of TCM (4)



Examples of TCM (5)

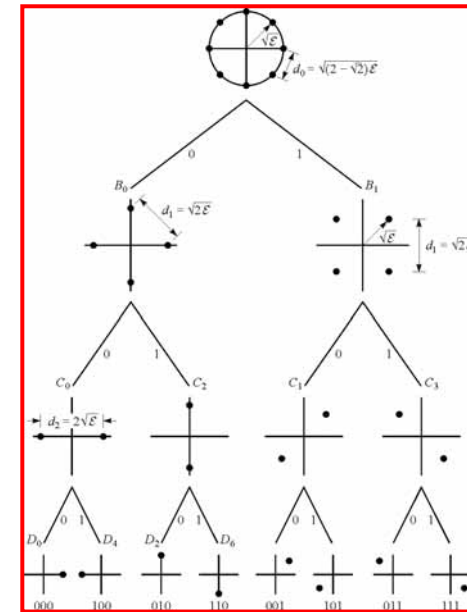


No parallel transition

Distance : $(0, 0, 0) \rightarrow (6, 7, 6)$

$$d_{000 \rightarrow 676}^2 = d_0^2 + 2d_1^2$$

$$= \left[(2 - \sqrt{2})\varepsilon + 4\varepsilon \right] = 4.585\varepsilon$$



Coding Gain (1)

- The minimum Euclidean distance between paths that diverge from any state and remerge at the same state in the trellis code is called free Euclidean distance D_{fed}
- **Asymptotic coding gain:**

$$\gamma = \left(\frac{E_{\text{uncoded}}}{d_{\text{min,uncoded}}^2} \right) / \left(\frac{E_{\text{coded}}}{D_{\text{fed,coded}}^2} \right)$$

where E is the normalized average received energy

$$\text{when } E_{\text{uncoded}} = E_{\text{coded}}, \quad \gamma = \frac{D_{\text{fed,coded}}^2}{d_{\text{min,uncoded}}^2}$$

In the 4-state example, $D_{fed} = 2\sqrt{\varepsilon}$, $d_{\text{min,uncoded}} = \sqrt{2\varepsilon}$

$\gamma=2 \Rightarrow 3\text{dB coding gain}$

Coding Gain (2)

- *Asymptotic coding gain can be increased by increasing the number of states and the rate of the convolutional encoder.*

In the 8-state example,

$$D_{fed} = \sqrt{4.585\varepsilon}, d_{\min, \text{uncoded}} = \sqrt{2\varepsilon}$$

$$\gamma = 2.2925 \Rightarrow 3.6\text{dB coding gain}$$

Coding Gain (3)

CODING GAINS FOR TRELLIS-CODED 16-PSK MODULATION

Number of states	k_1	Code rate $\frac{k_1}{k_1 + 1}$	$m = 3$ coding gain (dB) of 16-PSK versus uncoded 8-PSK	$m \rightarrow \infty$ N_{fed}
4	1	1/2	3.54	4
8	1	1/2	4.01	4
16	1	1/2	4.44	8
32	1	1/2	5.13	8
64	1	1/2	5.33	2
128	1	1/2	5.33	2
256	2	2/3	5.51	8

Source: Ungeboeck (1987).

Viterbi Decoding (1)

➤ *Two steps:*

□ *Step 1:* At each branch in the trellis,

- Compare the received signal to each of the signals allowed for that branch.
- Save the signal closest to the received signal
- Label the branch with a metric proportional to the Euclidean distance between the two signals.
- *Branch metric calculation*

Determining the best signal within each subset, i.e., the signal with the smallest distance to the received signal → *subset decoding*

Viterbi Decoding (2)

➤ *Two steps:*

□ *Step 2:*

- Apply the Viterbi algorithm to the trellis, with surviving partial paths corresponding to partial signal sequences that are closest to the received sequences.
- Select the ML path (the complete signal sequence closest in Euclidean distance to the received sequence) at the end of the trellis.
- *Path metric calculation*
- *Trellis update*

Error Rate Performance

- An error event happens when an erroneous path is selected at the decoder
- Error-event probability in AWGN channel:

$$P_e \approx N_{fed} Q \left(\sqrt{\frac{D_{fed}^2}{2N_0}} \right) \quad \text{under high SNR}$$

N_{fed} → the number of signal sequences with distance D_{fed} that diverge at any state and remerge at that state after one or more transitions

Announcement

- **Dr. Yuk's lecture notes could be downloaded from the department course material web page under "ELEC7073(tiyuk)".**

<https://www.eee.hku.hk/courses.msc/elec7073a/>