



STANFORD

Lecture 4

Capacity, Separation Thm, & C-OFDM

April 11, 2024

JOHN M. CIOFFI

Hitachi Professor Emeritus of Engineering

Instructor EE379B – Spring 2024

Announcements & Agenda

- Announcements
 - Problem Set #2 due Wednesday April 17 at 17:00
 - Sections 2.3-2.5
 - Need to leave off hours around 10:40 to catch flight today.

- Agenda
 - Capacity Continued
 - Chain Rule
 - Separation Theorem
 - Coded MultiTone



Capacity Continued

Sections 2.4 – 2.5

[See PS2.3 \(Prob 2.10\)](#)

General Capacity Theorem from L3

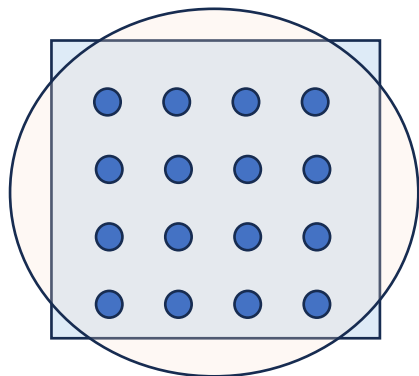
$$\frac{|A_N^\epsilon(\mathbf{x})|}{|A_N^\epsilon(\mathbf{x}/\mathbf{y})|} \rightarrow 2^{\mathcal{I}(\mathbf{x};\mathbf{y})} \quad \text{since } \mathcal{I}(\mathbf{x};\mathbf{y}) = \mathcal{H}\mathbf{x} - \mathcal{H}\mathbf{x}/\mathbf{y}$$

- Good codes will have only 1 codeword per conditional entropy subset.
- MAP detector decision region is then $\sim A_N^\epsilon(\mathbf{x}/\mathbf{y})$ - on average; but we can find it for one good code.
- If $A_N^\epsilon(\mathbf{x})$ were any larger, all codes (good or bad) will have at least one $A_N^\epsilon(\mathbf{x}/\mathbf{y})$ that contains 2+ codewords, which mean the MAP has to “flip a coin” – not good (high error prob).
- SHANNON’S CAPACITY THEOREM
 - Number of codewords is limited by mutual info $b \leq \mathcal{I}(\mathbf{x};\mathbf{y})$.
 - Which is per-subsymbol equivalent with random code $\tilde{b} \leq \mathcal{I}(\tilde{\mathbf{x}};\tilde{\mathbf{y}})$.
 - If maximized over input distributions $\tilde{b} < \tilde{c} \leq \max_{p_{\tilde{\mathbf{x}}}} \mathcal{I}(\tilde{\mathbf{x}};\tilde{\mathbf{y}}) \frac{\text{bits}}{\text{subsymbol}}$.



The uniform part is most important (from L3)

- The Gaussian distribution corresponds to marginal of uniform distribution over a hypersphere.
 - This uniform distributions marginals are asymptotically Gaussian.
 - This is a special case where uniform and Gaussian are basically the same.
 - Because all the Gaussian infinite-length vectors (codewords) have same energy (zero variance of the energy).
- All the points (really volume) are (is) at the surface.
- The Gaussian marginal dist'n is important only for shaping gain (< 1.53 dB).
- The (AEP) uniform spacing of points (no matter where the majority of them sit, surface or otherwise) remains for the fundamental gain.



The uniform spacing separates codewords in the union of the hypersquare (orthotope) and hypersphere.

Thus, good codes can be based on sequences from uniformly spaced PAM/QAM subsymbols.

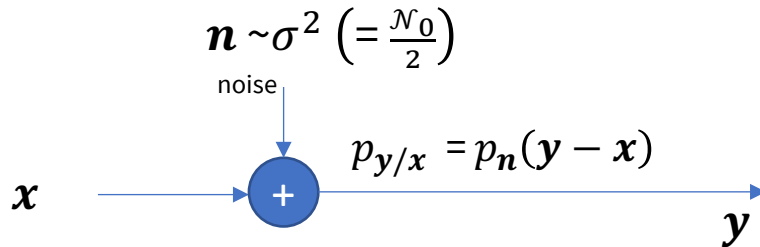
And the rest is MMSE estimation,

With a chain-rule twist in some situations

Vector Coding is always MAP, ML, & MMSE, a special case.



AWGN Capacity Review (379A)



$$\bar{C} = \frac{1}{2} \cdot \log_2 \left(1 + \underbrace{\frac{\bar{\mathcal{E}}_x}{\sigma^2}}_{SNR} \right)$$

- Often “gain” $\|h\|^2$ is absorbed into energy, really $g = \frac{\|h\|^2}{\sigma^2}$ so a “channel gain” $\bar{C} = \frac{1}{2} \cdot \log_2(1 + g \cdot \bar{\mathcal{E}}_x)$.
 - Note g here is per real dimension, but if complex- baseband channel, it would be $\tilde{C}_x = \log_2(1 + g \cdot \bar{\mathcal{E}}_x)$.
 - Know context and be consistent with numerator/denominator dimensionality.
- SNR=4.7 dB (3 and $g=1$), then $\bar{C} = 1$ bit/dimension.
- SNR=20 dB (100 and $g=1$), then 3.33 bits/dimension – and thus 6.67 bits/complex subsymbol.
- What SNR gives 7 bits per dimension? $10 \cdot \log_{10}(2^{14} - 1) = 14 \cdot 3 = 42$ dB.



Chain Rule

Subsection 2.3.2

Chain Rule

$$\mathbb{I}(\mathbf{x}; \mathbf{y}) = \sum_{n=1}^N \mathbb{I}(\tilde{\mathbf{x}}_n; \mathbf{y} / [\tilde{\mathbf{x}}_{n-1} \cdots \tilde{\mathbf{x}}_1])$$

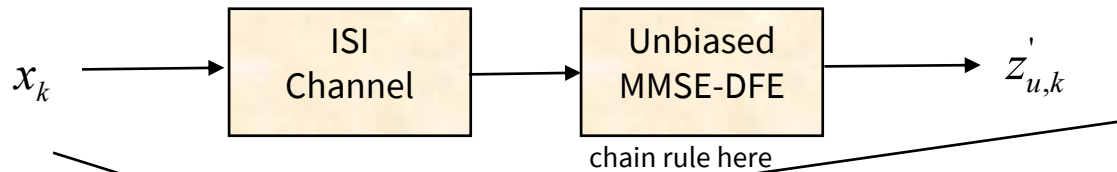
- If parallel independent channels (we know this by now!), the $[\tilde{\mathbf{x}}_{n-1} \cdots \tilde{\mathbf{x}}_1]$ provide no help;
 - just sum the individual channels $\mathbb{I}(\tilde{\mathbf{x}}_n; \tilde{\mathbf{y}}_n)$.
- But suppose not: each term represents a code (MMSE-related if Gaussian) problem with SNR, capacity, etc.

Matrix AWGN: GDFE (sometimes also called “successive decoding”)

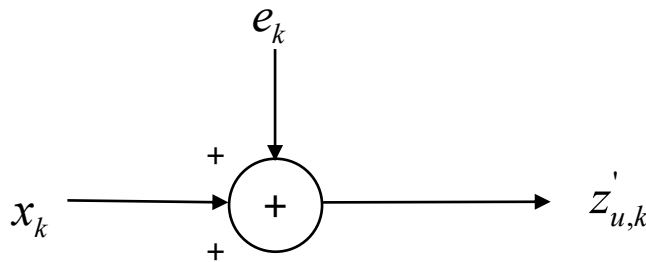
- Estimate (MMSE) and decode $[\tilde{\mathbf{x}}_{n-1} \cdots \tilde{\mathbf{x}}_1]$ first, then simpler single component problem.
 - So not just linear MMSE, linear MMSE + subtract “earlier” subsymbols’ effect (nonlinear).
- It’s parallel channels, but with a twist to make them independent step by step (“decision-feedback”).



CDEF Example EE379A, L18



equivalent to



$$\mathcal{I}(\tilde{x}; \tilde{y}) = \mathcal{H}_{\tilde{x}_k} - \underbrace{\mathcal{H}_{\tilde{x}_k / [\tilde{y} \quad \tilde{x}_{k-1} \dots -\infty]}}_{\text{MMSE-DFE}}$$

$$SNR = SNR_{mmse-dfe,u} = 2^{2\mathcal{I}(\tilde{x}; \tilde{y})} - 1$$

- The MMSE-DFE achieves the highest rate (with $\Gamma = 0$ dB) also:
 - $I = C$ if water-filling spectrum is at transmitter.
 - But, this spectra may be impossible with a single DFE, so can be several parallel DFES (see Section 3.12).
 - No error propagation (true if P_e is zero), and **canonical** (reliably achieves capacity).

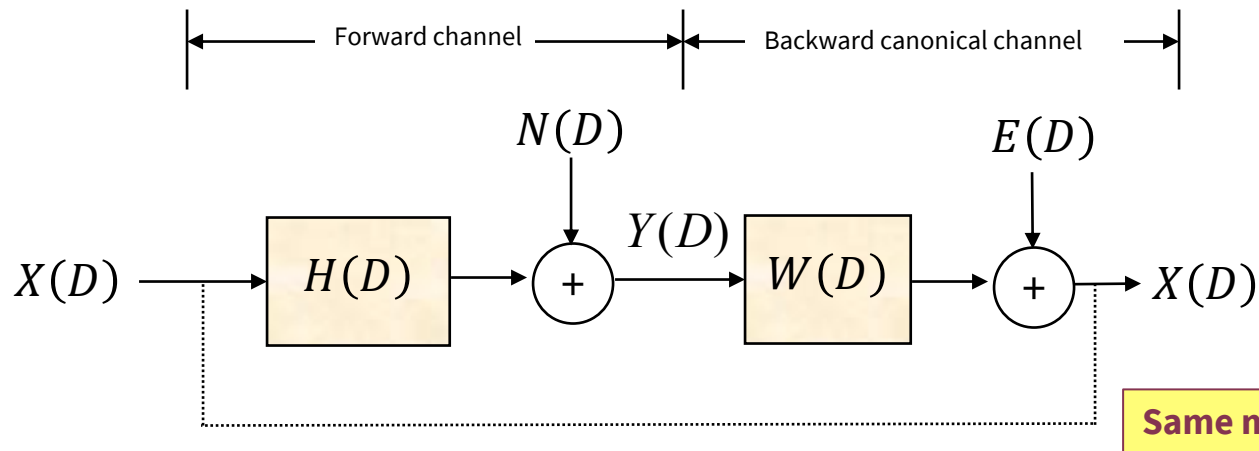


There are many chain rule orders

- $\mathcal{I}(\mathbf{x}; \mathbf{y}) = \sum_{n=1}^N \mathcal{I}(\tilde{\mathbf{x}}_{\pi(n)}; \mathbf{y} / [\tilde{\mathbf{x}}_{\pi(n-1)} \cdots \tilde{\mathbf{x}}_{\pi(1)}]) = \sum_{n=1}^N \mathcal{I}(\tilde{\mathbf{y}}_{\pi(n)}; \mathbf{x} / [\tilde{\mathbf{y}}_{\pi(n-1)} \cdots \tilde{\mathbf{y}}_{\pi(1)}])$
- $N!$ orders exist for each of $\mathcal{I}(\mathbf{x}; \mathbf{y}) = \mathcal{I}(\mathbf{y}; \mathbf{x})$.
- Every order corresponds to different set of parallel channels (some with feedback, a few without).
- But they all produce the same maximum data rate (achieved with good code that has zero gap).
- Thus, not only are there a lot of good codes – there are a lot of good MMSE-based modulation/demodulation designs also!



Forward and its Backward Canonical Models



$$r(t) = h_c(t) * h_c^*(-t) = \|h\|^2 \cdot q(t)$$

$$y(t) \rightarrow h_c^*(-t) \rightarrow \frac{1}{T} \rightarrow Y(D)$$

$$Y(D) = R(D) \cdot X(D) + N(D)$$

$$\underbrace{\quad}_{\frac{N_0}{2} \cdot R(D)}$$

Forward Canonical Model

$$X(D) = \underbrace{W(D)}_{MMSE-LE} \cdot Y(D) + \underbrace{E(D)}_{\frac{N_0}{2} \cdot W(D)}$$

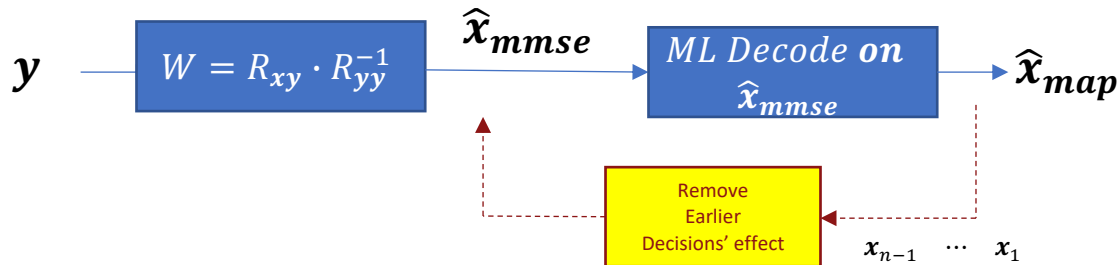
Backward Canonical Model
chain rule helps more here



For the filtered/matrix AWGN

- The MAP and MMSE determine the performance, and also the chain rule suggests a simpler decoder:

$$\begin{aligned}
 \mathcal{I}(\tilde{\mathbf{x}}; \tilde{\mathbf{y}}) &= \mathcal{H}_{\tilde{\mathbf{y}}} - \mathcal{H}_{\tilde{\mathbf{y}}/\tilde{\mathbf{x}}} \\
 &= \log_2 \left(\frac{|R_{\tilde{\mathbf{y}}\tilde{\mathbf{y}}}|}{|R_{\tilde{\mathbf{n}}\tilde{\mathbf{n}}}|} \right) \text{ bits/subsymbol} \\
 &= \mathcal{H}_{\tilde{\mathbf{x}}} - \mathcal{H}_{\tilde{\mathbf{x}}/\tilde{\mathbf{y}}} \\
 &= \log_2 \left(\frac{|R_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}}|}{|R_{ee}|} \right) \text{ bits/subsymbol} \\
 &= \log_2 |I - W \cdot H| \quad \text{Forward} \\
 &= \log_2 |I - H \cdot W| \quad \text{Backward}
 \end{aligned}
 \qquad = \log_2(SNR_{mmse})$$



Still MAP if “previous” decisions are correct sequentially decodes

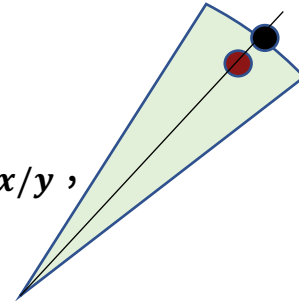


Optimal Detectors for Good Codes

- x codewords/subsymbols selected from Gaussian, $[x \ y]$ are jointly Gaussian (as is then n).
- **ML = MAP** since all good code's x codewords/symbols are equally likely (uniform, AEP):

$$\frac{MAP}{ML} \ni \min_{\{\tilde{x}_k\}} \sum_{k=-\infty}^{\infty} \|\tilde{y}_k - H \cdot \tilde{x}_k\|^2 \neq \sum_{k=-\infty}^{\infty} \|\tilde{n}_k\|^2.$$

Same as $\max_x p_{x/y}$,
where x has ∞ length



- **MMSE = MAP** The smallest sum will reduce $\{\tilde{x}_k\}$ magnitude slightly because it also shrinks noise (trade-off in sum):

$$MMSE \ni \min_{\{\tilde{x}_k\}} \left\{ \lim_{K \rightarrow \infty} \frac{1}{2K + 1} \sum_{K=-K}^K \|\tilde{x}_k - W \cdot \tilde{y}_k\|^2 \right\}$$

min over entire sum
 W is MMSE filter

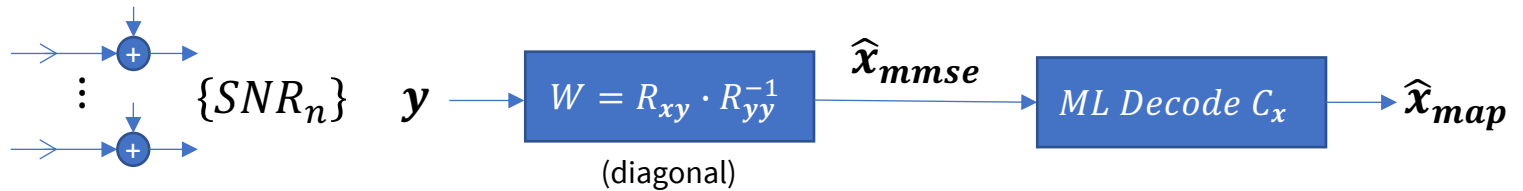
- By LLN, this sum is MMSE & MAP and has optimum $\hat{x} = E[\tilde{x}/\tilde{y}]$ on average over the random code set.
- The bias removal is unnecessary because of the hyper-conical decision regions (like QPSK where decision regions don't change) for a zero-gap AWGN code, but we now know MMSE is a "DFE-like" structure (chain rule)



Separation of Coding and Modulation

Subsection 4.4

The best (MAP) receiver

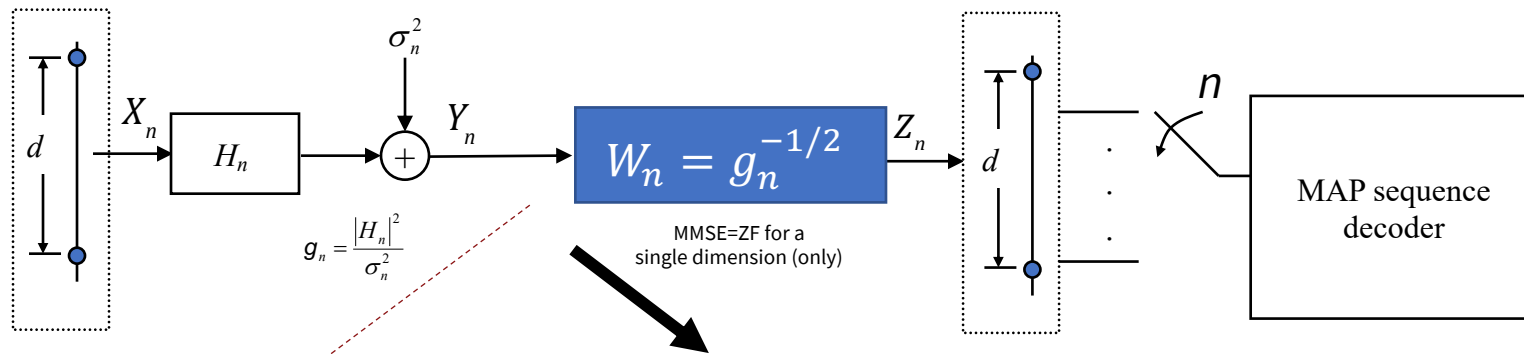


- Each parallel channel has $\mathcal{I}(x_n; y_n) = \log_2 SNR_{mmse,n}$, since they are independent
- Suppose each dimension is a dimension within the same code?
 - The dimensional subsymbols will remain uncorrelated (but not independent because of the same code).
 - On average over all (Gaussian) codes these dimensions are independent, not for specific code.
- W is a scalar multiply for each such uncorrelated dimension (so does not change signal to noise)
 - Does use of MMSE matter? (not for VC or DMT with $\tilde{b}_n = \log_2(1 + SNR_n)$)
- **BUT YES, IT DOES MATTER – IF, a constant** bits/subsymbol $b_n \equiv \tilde{b}$ (and/or SNR_{geo}) – **Coded OFDM**
 - because it impacts the weight of different dimensions before the final ML decoding; this (it turns out) is the same as the earlier bit-loading, in effect for same energy distribution/Rxx.

$$LL_n = -\frac{1}{2\sigma_n^2} \cdot \|y_n - \hat{x}_n\|^2$$



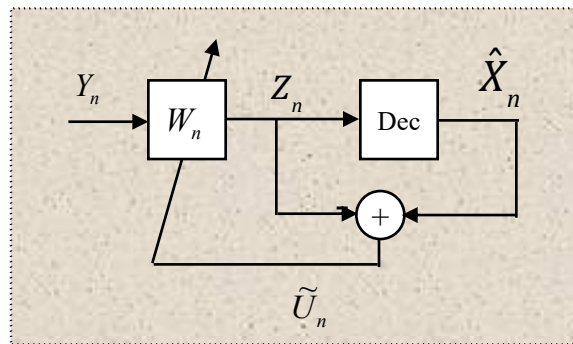
Dimensional Normalizer is MMSE



Must have each dimension scalar for ZF=MMSE

Zero-Forcing Algorithm

$$W_{n,k+1} = W_{n,k} + \mu \cdot \tilde{U}_n \cdot \hat{X}_n$$



see also L5:6, L4:30

- MAP Decoder is

$$\min_{\{X_{n,k}\}} \left\{ \sum_{k=0}^{\infty} \sum_{n=0}^{N-1} |X_{n,k} - W_n \cdot Y_{n,k}|^2 \right\}$$

dimensions with low gains g_n have greater contribution to minimization, higher soft/intrinsic info, and the decoder is consequently affected.



Separation Theorem

Theorem 4.4.1 [Separation of Coding and Modulation] *Given a set of independent partitioned AWGN-channel dimensions with energies/dimension $\bar{\mathcal{E}}_n$ and gains g_n with equivalent*

$$SNR_{geo} = \left[\prod_{n=1}^N (1 + \bar{\mathcal{E}}_n \cdot g_n) \right]^{1/N} - 1 ,$$

N repeated uses of a single good code with $\Gamma \rightarrow 0$ dB and

$$\bar{b} \leq \bar{\mathcal{I}} = \frac{1}{2} \cdot \log_2(1 + SNR_{geo})$$

and corresponding constant constellation $|C| = 2^{\bar{b} + \bar{p}}$ achieves the same performance as using N instances of that same good code with $\Gamma \rightarrow 0$ dB each with variable constellation $|C_n|$ and bits per tone

$$\bar{b}_n \leq \bar{\mathcal{I}}_n = \frac{1}{2} \cdot \log_2(1 + \mathcal{E}_n \cdot g_n) .$$

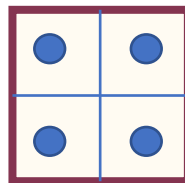
■ Critical are:

- the independent parallel dimensions (not code, the partitioned matrix/filtered AWGN), &
- the good code for which the input to the parallel dimensions comes from large constellation with subsymbol distribution approaching Gaussian \sim random coding,
- the code can apply “down the symbol” (over the subsymbols that all have same constellation).

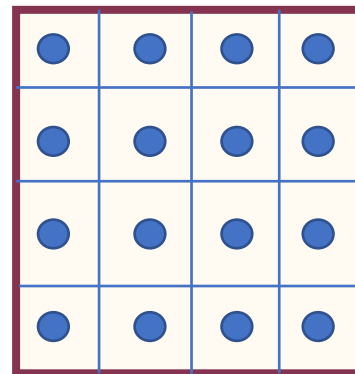


Simple Separation Thm Example

- The two tones' ave is $\tilde{I}_{ave} = 2$.
- The ST says use of a single constellation with $|C_{ave}| = 8$ is sufficient.
- Decoder must consider channel gains.

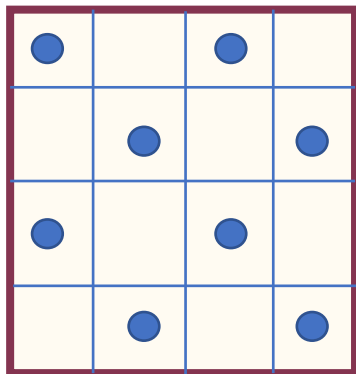


$$\tilde{I}_1 = 1$$
$$|C_1| = 4$$



$$\tilde{I}_2 = 3$$
$$|C_2| = 16$$

$$\tilde{I}_{ave} = 2$$
$$|C_{ave}| = 8$$



- Looks like 2 identical uses of a single AWGN with geometric-average SNR, equivalently $SNR = 2^{\tilde{I}_{ave}} - 1$



Separation T is widely applicable

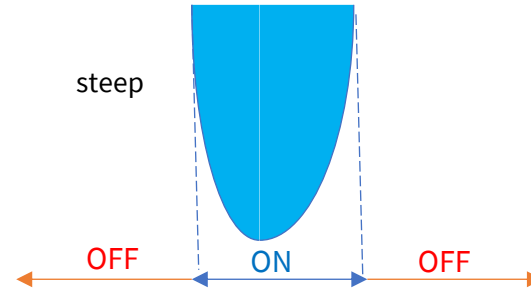
- This works for partitioning with
 - SVD,
 - Eigenvectors,
 - DFT/FFT (becomes Coded-OFDM here), &
 - other bases.
- The transmitter does not need to know individual \tilde{b}_n , just the sum for any symbol.
- It works for any \mathcal{E}_n and leads to highest rate for those energies $\mathcal{I}(\tilde{x}; \tilde{y})$.
 - Water-fill set gives highest data rate (highest mutual information).
- We've seen in our examples that water-fill is close to on/off.
 - So, if the designer guesses well the on/off, the system requires ALMOST no feedback of bit distribution to transmitter.
 - In practice, the constellation size and redundancy need specification, and thus on some indication of the value of $\mathcal{I}(\tilde{x}; \tilde{y})$ for the channel.
- **Example:** Wireless "MCS" (modulation coding scheme) specifies code rate and constellation size only in feedback to transmitter. The on/off distribution?
 - 5G/Wi-Fi ignore this for time-frequency and just use flat over the entire band.
 - 5G/Wi-Fi do excite spatial "streams" unequally in that some can carry data and while others are zeroed.

ONLY IF $\Gamma \rightarrow 0$



Caution on Water-filling and on/off

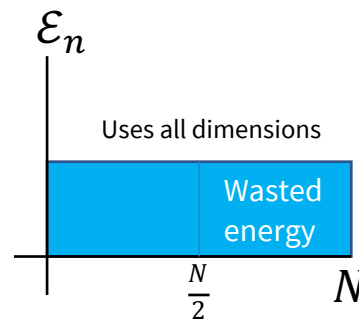
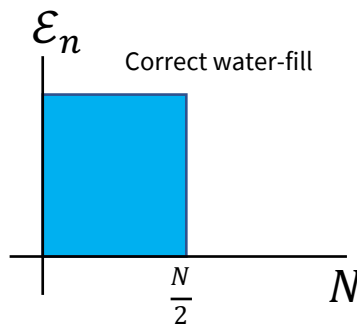
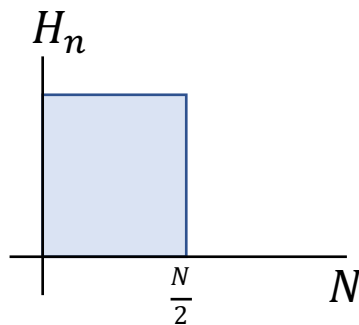
- Most water-fill will satisfy $\left(1 + \frac{SNR_n^*}{\Gamma}\right) \cong \frac{SNR_n^*}{\Gamma}$
 - IF dimension carries nonzero energy.
- The energy closely approximates flat:
 - RA: $K - \frac{\Gamma}{g_n} = \frac{\epsilon_x}{L^*} + \frac{\Gamma}{L^*} \cdot \sum_{l=1, l \neq n}^{L^*} 1/g_l$ is roughly the same (no one dimension dominates).
 - MA: $K - \frac{\Gamma}{g_n} = \Gamma \cdot \left(\frac{2\tilde{b}}{\prod_{l=1}^{L^*} g_l}\right)^{1/L^*} - \frac{\Gamma}{g_n} = \frac{\Gamma}{g_n} \cdot \left\{ \left(\frac{SNR_{geo}}{\prod_{l=1, l \neq n}^{L^*} g_l}\right)^{1/L^*} - 1 \right\}$ is roughly the same (again no one dim dominates).
- This is true on waterfill's ENERGIZED (“on”) dimensions, NOT for zeroed dimensions (“off”).



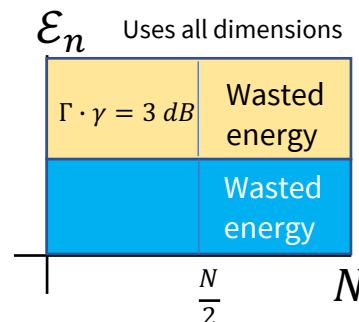
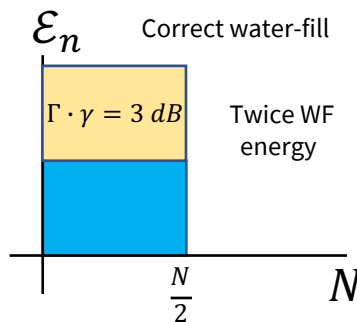
So, it is NOT true that wireless' C-OFDM is the same as DMT, UNLESS the used dimensions are close to the same and $\Gamma \rightarrow 0$!



Half-Band Example: Revisit 379A's L18:12-13



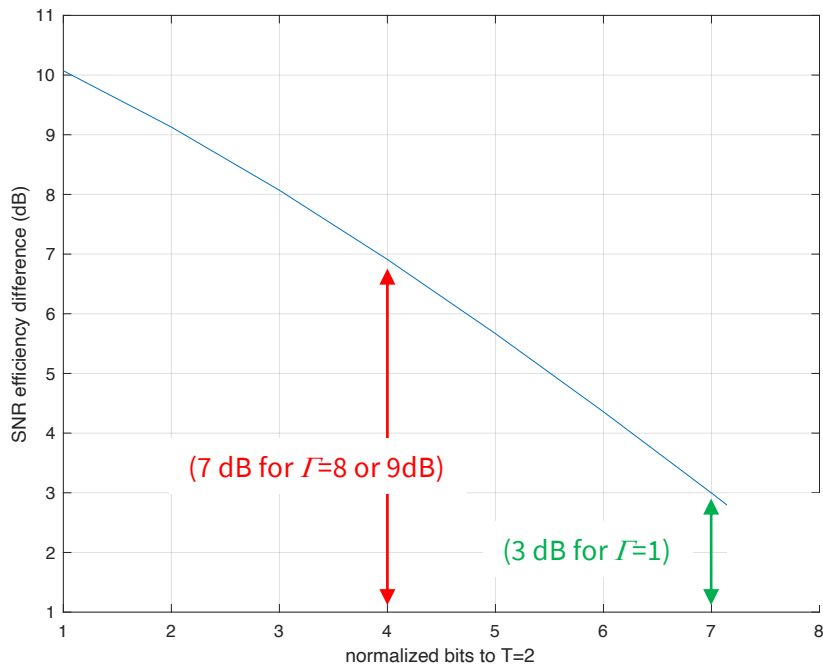
- The geometric SNR for water-fill is 3 dB higher if capacity-achieving codes are used
 - Or could run the water-fill system at same data rate at 3 dB less energy
- This amount is amplified below capacity by non-unity (not 0 dB) gap-margin product



4x WF energy!



margin difference for half-band optimum versus full band



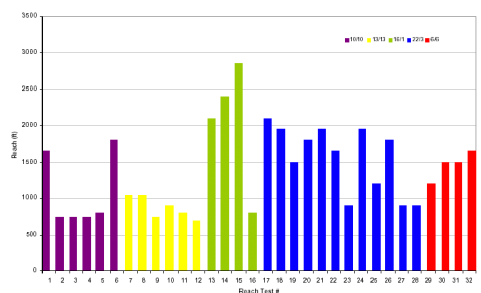
margin difference for half-band optimum versus full band

- Capacity of AWGN with WF is 8 bits/subsymbol (4 bits/dimension)

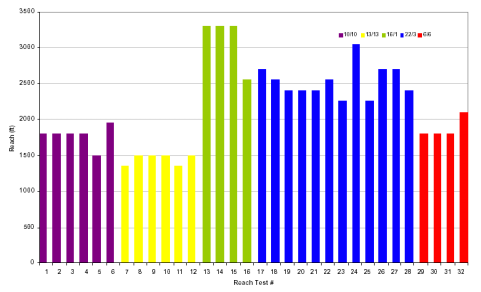
1993 ADSL Olympics – Bellcore
Margin differences at 1.6 Mbps, 4 miles, 11+dB
DMT 4x faster (6 Mbps) at 2 miles

2003 VDSL Olympics - Bellcore

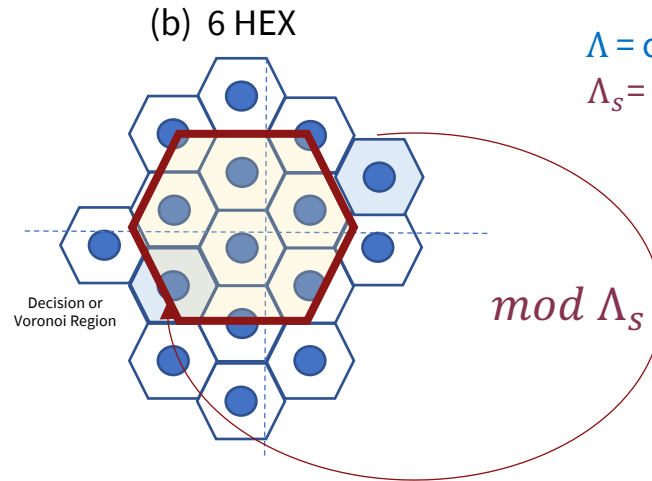
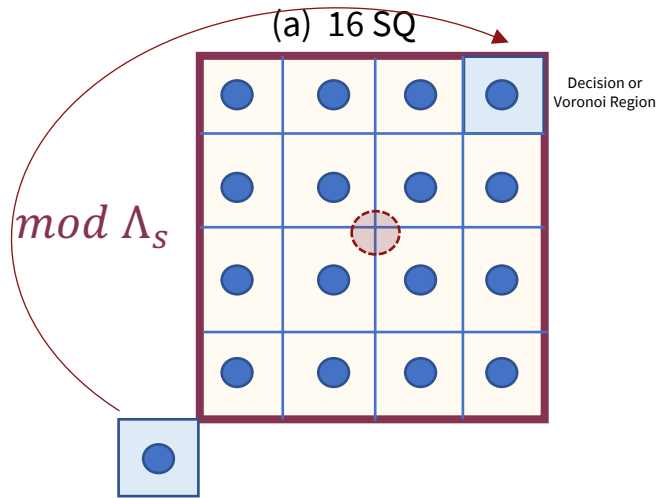
Variable f_c and 1/T single-carrier QAM results



DMT* results – exact same channels as QAM



Coding Gain Refresher



$\Lambda =$ coding lattice for d_{min}
 $\Lambda_s =$ shaping lattice for \mathcal{E}_x

$$\gamma \triangleq \frac{\left(\frac{d_{\min}^2(\mathbf{x})}{\bar{\mathcal{E}}_{\mathbf{x}}} \right)}{\left(\frac{d_{\min}^2(\check{\mathbf{x}})}{\bar{\mathcal{E}}_{\check{\mathbf{x}}}} \right)} = \underbrace{\left(\frac{\frac{d_{\min}^2(\mathbf{x})}{V^{2/N}(\Lambda)}}{\frac{d_{\min}^2(\check{\mathbf{x}})}{V^{2/N}(\check{\Lambda})}} \right)}_{\gamma_f \text{ fundamental gain}} \cdot \underbrace{\left(\frac{\frac{V^{2/N}(\Lambda)}{\bar{\mathcal{E}}_{\mathbf{x}}}}{\frac{V^{2/N}(\check{\Lambda})}{\bar{\mathcal{E}}_{\check{\mathbf{x}}}}} \right)}_{\gamma_s \text{ shaping gain}}$$

Basic principle extends $\bar{N} \rightarrow \infty$
 Hexagon \rightarrow hypersphere (Gaussian marginals)

**good codes can follow
 from $\Lambda_s / \Lambda = |C|$**



SQ constellations vs “Gaussian”

- There is always a loss for a non-hyper-spherical constellation boundary on the (any matrix/filtered) AWGN.
 - The max shaping gain, $\gamma_{s,max} < 1.53$ dB (when $\tilde{b} \geq 1$), relative to hypercube.
 - Hypercube is often the assumed reference system (so Λ for fundamental and scaled Λ_s for shaping).
- All of random coding/AEP can repeat with the input distribution being uniform in any dimension (instead of Gaussian) – hypercube-energy constraint.
- The MMSE Estimator can still be used with decoder, and it’s basically

$$\tilde{C} = \log_2(1 + SNR_{mmse,u}/\gamma_{s,max}).$$

- There is loss of 0.5 bit/complex dimension using SQ constellations.

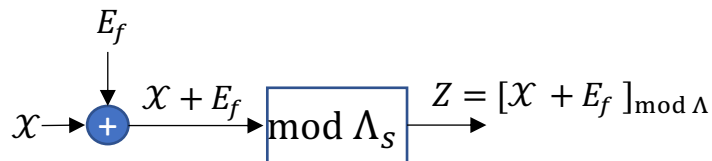
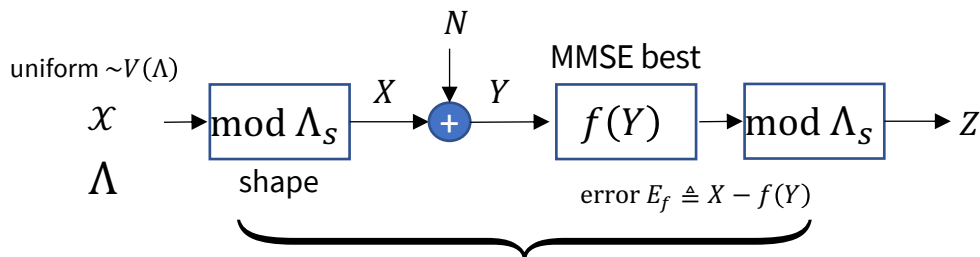


Forney's Crypto MMSE equivalence

1.53 dB (max) loss

$$\tilde{j} \geq \tilde{c} - \underbrace{\log_2 \left(\pi \cdot e \cdot \frac{\epsilon}{V(\Lambda_s)} \right) - \log_2 \left(\frac{\sigma_{E_f}^2}{\sigma_{mmse}^2} \right)}_0$$

$$\Lambda_s = \sqrt{\frac{|C|}{2}} \cdot Z^2$$



- See also Section 2.8 – there is a shaping loss with any Λ_s that is not a hypersphere (SQ is worst in practice) so various shaping methods can apply; however the separation theorem still applies to them all, with random coding used on uniform over Λ_s 's Voronoi region.



Coded OFDM/MT

Subsection 4.4.1

SQ constellations vs “Gaussian”

- Matrix/filtered-AWGN loss for “square” constellations

$$\gamma_s \leq 1.53 \text{ dB}$$

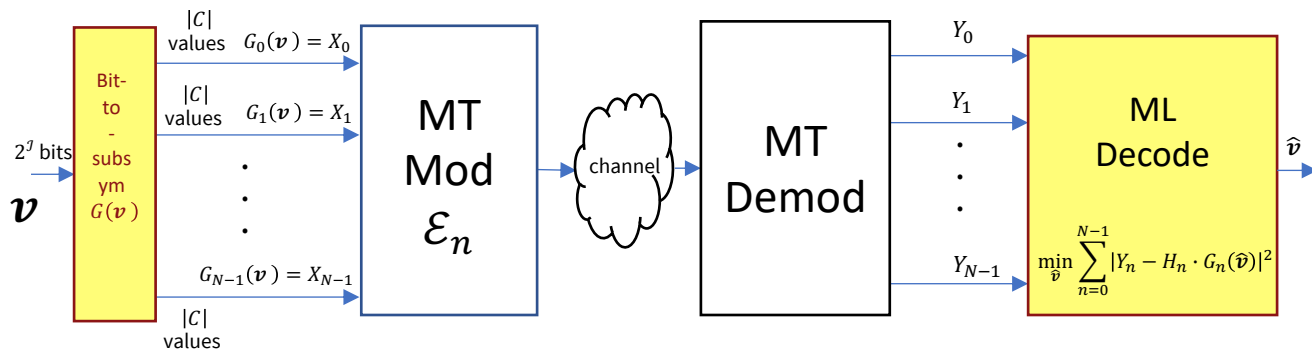
shaping gain

$$\tilde{C} = \log_2(1 + SNR_{geo}/\gamma_{s,max})$$

- When $\tilde{C} \leq 1$, $\gamma_{s,max} \cong 0$ dB.
 - There is tiny **low-SNR-AWGN** shaping loss for binary codes.
- AEP applies to hypercube (with shaping loss) boundary and random codes.
- MMSE estimator precedes MAP decoder for **original** code:
 - ISI/crosstalk is optimally handled linearly with parallel ind subchannels.
 - Nonlinear decision feedback needed when NOT parallel independent channels (chain rule).



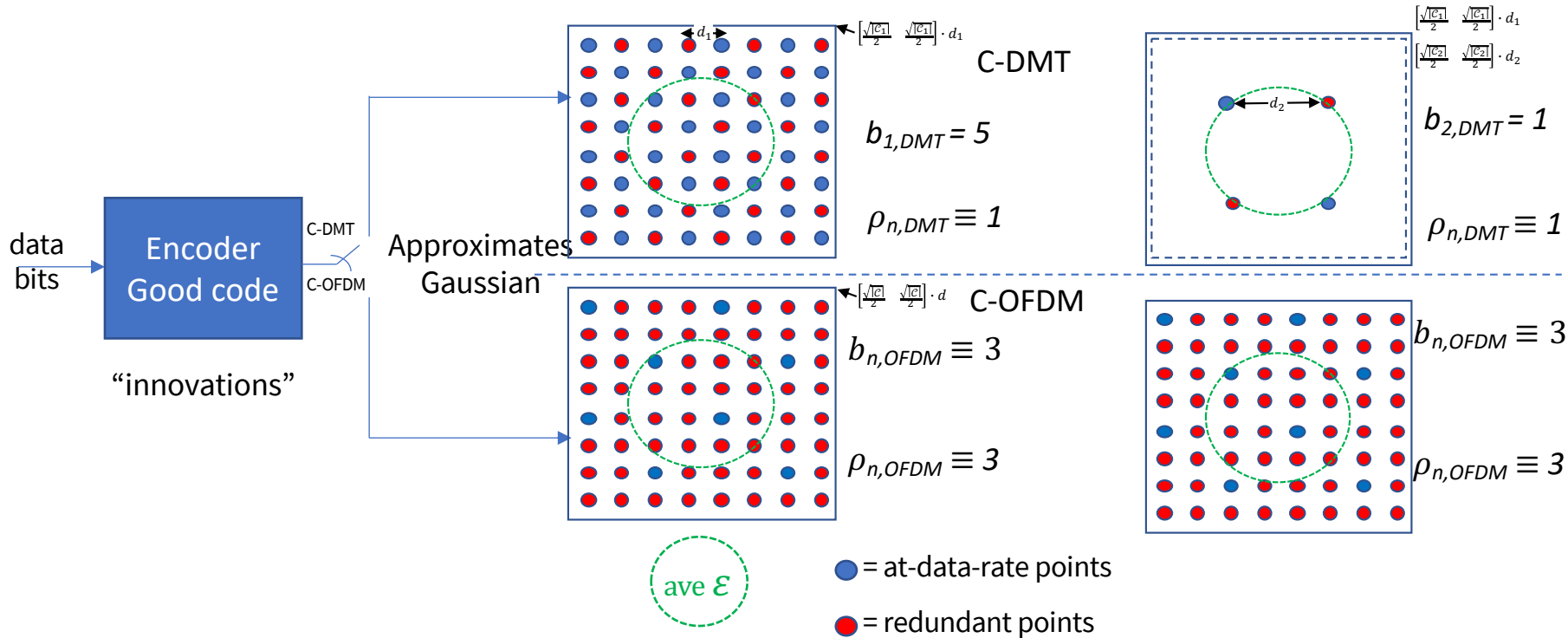
Coded-MT/OFDM



- Treats a pre-agreed known set of dimensions as repeated constant SNR_{geo} dimensions.
 - No transmitter bit loading, and energy is on/off on the pre-agreed set.
- The MT could be replaced by space-time MIMO, “Coded-Vector-Coding” – same basic principle.
- Usually wireless MIMO does allow “water-fill” over spatial dimensions (but not temporal).



Comparison of variable and fixed constellation



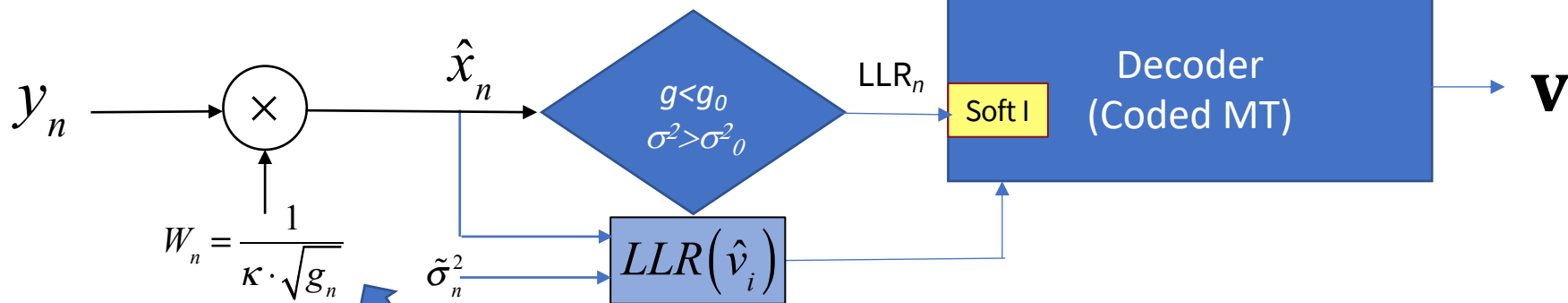
- These types of system are heavily used in practice



Full MAP Decoder – Coded MT

LLR = log likelihood ratio

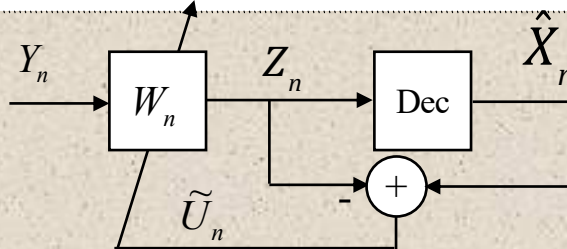
Computed from Gaussian noise dist'n & from input code constraints, each subsymbol and/or bit ($\tilde{\sigma}_n^2$)



$$W_{n,k+1} = W_{n,k} + \mu \cdot \tilde{U}_n \cdot \hat{X}_n$$

$$\tilde{\sigma}_{n,k+1}^2 = (1 - \mu') \cdot \tilde{\sigma}_{n,k}^2 + \mu' \cdot |\tilde{U}_{n,k}|^2$$

$$g_n = \frac{1}{\tilde{\sigma}_n^2}$$





End Lecture 4

BSC and BEC

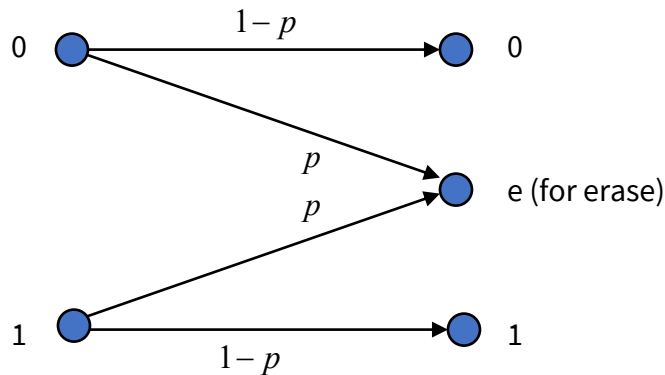
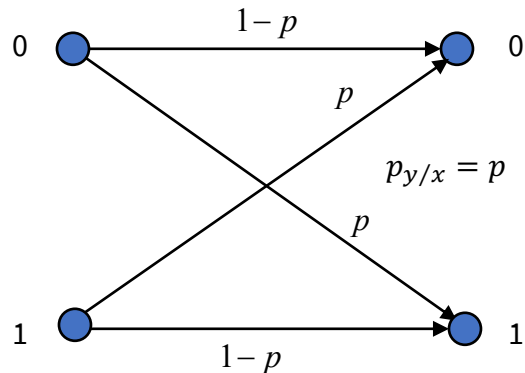
■ **BSC** has $\bar{C} = 1 - \mathcal{H}(p) = 1 - p \cdot \log_2 p - (1 - p) \cdot \log_2(1 - p)$.

- $p = 1/2 \rightarrow 0$ bits possible (makes sense).
- $p = 0 \rightarrow 1$ bit/dimension reliably (makes sense).
- $0 \leq \bar{C} \leq 1$.

■ **BEC** has $\bar{C} = 1 - p$.

- $p = 1/2 \rightarrow 1/2$ bits/dim reliable (no errors only erasures).
- $p = 0 \rightarrow 1$ bit/dimension reliably (makes sense).
- $0 \leq \bar{C} \leq 1$.

■ BEC is better than BSC (higher capacity) – decoders can use erasures with $N > 1$ to improve (reduce) P_e (soft info, 379A).



Symmetric DMC

- Generally, just a discrete probability transition matrix (Appendix A).
- q -ary (example 0,...,255 for a byte = subsymbol)

$$\mathcal{C} = b - p_s \cdot \log_2 \frac{2^b - 1}{p_s} + (1 - p_s) \cdot \log_2 (1 - p_s) \leq b \text{ bits.}$$

- $p_s = .01$
- $\mathcal{C} = 7.88$ bits/subsymbol.

