

# Hamming、Hadamard、Golay 纠错码实验

陈硕 chenshuo@chenshuo.com

最新版及讨论: <https://github.com/chenshuo/nuedc> [2026/04/03]

这篇个人笔记整理记录了我对几种早期纠错码的理解, 以及编解码软硬件实验, 并回顾了它们在历史上的实际用处。它们用到的数学知识非常简单, 有基本的线性代数概念 (矩阵、向量、线性空间等) 就能理解。本文介绍的 Hamming、Hadamard、Golay 等纠错码都诞生于 1954 年之前, 目前已基本淘汰。1960 年发明的 BCH 码、Reed-Solomon 码、LDPC 码至今仍有广泛的应用。我录制过 Python 实现 Reed-Solomon 纠错码的系列视频, 内容更深入一些。<sup>①</sup>

图 1 是数字通信系统的简化模型, 中间是理想的二元对称信道 (BSC), 每个 bit 发生翻转的概率是  $f$ , 并且 bit 翻转是独立事件。透过 BSC 传输  $n$  bits 时恰好翻转  $t$  bits 的概率满足二项分布  $\binom{n}{t} f^t (1-f)^{n-t}$ 。纠错码的思路就是编码器通过发送一些冗余的信息 ( $k$  bits  $\rightarrow$   $n$  bits), 即便这  $n$  bits 数据在传输过程中发生了少许错误, 设法让译码器能从接收到的  $n$  bits 数据中还原出原来发送的  $k$ -bit 消息。

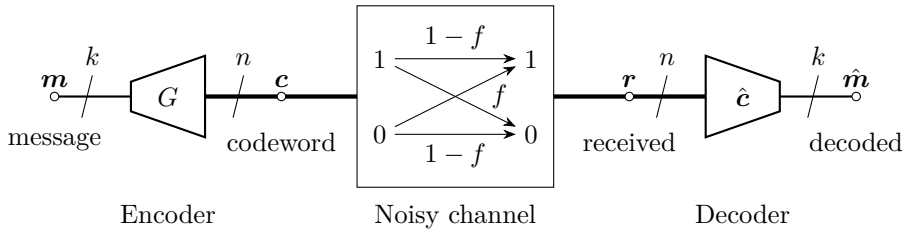


图 1: 数字通信系统的简化模型 (toy model)

**重复码——重要的事情说三遍** 将 1 bit 重复发送  $n$  次是最简单的纠错码。对  $n$  次重复码, 当  $n$  是奇数时用简单多数判决能容忍  $\frac{n-1}{2}$  个错, 那么译码出错的概率是  $p_b = 1 - \sum_{t=0}^{(n-1)/2} \binom{n}{t} f^t (1-f)^{n-t}$ 。<sup>②</sup>

例如对于  $n = 3$  的重复码, 若  $f = 0.1$  则  $p_b = 1 - (1-f)^3 - 3f(1-f)^2 = 1 - 0.9^3 - 3 \times 0.1 \times 0.9^2 = 0.028$ ; 如果  $f = 0.01$ ,  $p_b = 2.98 \times 10^{-4}$ 。换句话说, “重要的事情说三遍” 确实能降低错误概率。如果  $f < 0.5$ , 只要增加重复数, 可以把  $p_b$  降到任意值, 但是相应的码率  $R = 1/n$  也会趋于 0。例如在  $f = 0.1$  时, 如果要求  $p_b < 10^{-5}$ , 则需要重复  $n \geq 19$  次, 这是非常大的开销。

在 Claude E. Shannon 之前, 人们普遍认为要求  $p_b \rightarrow 0$  就意味着  $R \rightarrow 0$ , Shannon 在 1948 年证明<sup>[3]</sup>了  $p_b \rightarrow 0$  时可以有  $R = C > 0$ , 这里  $C$  是信道容量, 对于 BSC 来说  $C = 1 - H_b(f) = 1 + f \log f + (1-f) \log(1-f)$ 。若  $f = 0.1$ , 有  $C_{\text{BSC}} \approx 0.531$ , 这是一个惊人的结论, 等于说对于  $f = 0.1$  的 BSC, 理论上只需传输原来数据量的  $1/0.531 \approx 2$  倍就能让  $p_b$  达到任意小。

在 Shannon 发表他的结论的 60 多年后, 我们离他预言的极限已经很近了。[1]

<sup>①</sup>位于 <https://www.youtube.com/playlist?list=PL3wVcVGXqdnZtmggNEd48yPbElkPbj20A> 和 <https://space.bilibili.com/1356949475/lists/1980190>

<sup>②</sup>计算这个求和不用写循环代码, 直接调用 `scipy.special.bdnr` 函数即可。

# 1 Hamming 码

Hamming 码<sup>①</sup> 是 Richard W. Hamming 在 1940s 年代末发明的最早的纠错码，能够纠正单比特错误。他在 1950 年 4 月发表的论文<sup>[4]</sup> 是纠错码领域的开山之作，定义了 Hamming 距离等基本术语。这篇论文非常浅显易懂，没有任何高等数学，我估计高中生只要懂一点二进制就能大致看懂。

Richard Hamming<sup>②</sup> (1915–1998) 在 1968 获得了图灵奖，发明纠错码是获奖理由之一：For his work on numerical methods, automatic coding systems, and error-detecting and error-correcting codes.

本节先介绍 Hamming(7, 4) 码的编码与译码，再谈一谈线性分组码的一般理论。这是大学《通信原理》课程中关于信道编码的入门知识，国内教材一般翻译为“(7, 4) 汉明码”或“海明码”。

所谓“(n, k) 码”指的是编码前后的 bit 长度：把 k-bit 的消息 (message, *m*) 编码 (映射) 为 n-bit 的码字 (codeword, *c*)，即  $\{0, 1\}^k \Rightarrow \{0, 1\}^n$ 。Hamming(7, 4) 码就是把 4-bit 的消息 *m* 编码为 7-bit 的码字 *c*，4-bit 的 *m* 有 16 种取值，范围是 0 ~ 15，每一个  $m_i$  对应一个  $c_i$ 。表 1 列出了一种编码方式，取自一本优秀的本科教材<sup>[2]</sup>。这个表中的全部 16 个码字称为 Hamming(7, 4) 的码书 (codebook,  $\mathcal{C}$ )，也叫码表。对于一般的 (n, k) 码来说，码书有  $2^k$  个码字，每个码字为 n-bit，总的容量是  $n \cdot 2^k$  bits。

<i>m</i>	$m_3 m_2 m_1 m_0$	码字 <i>c</i>	十六进制	<i>m</i>	$m_3 m_2 m_1 m_0$	码字 <i>c</i>	十六进制
0	0000	0000 000	00	8	1000	1000 101	45
1	0001	0001 011	0B	9	1001	1001 110	4E
2	0010	0010 111	17	10	1010	1010 010	52
3	0011	0011 100	1C	11	1011	1011 001	59
4	0100	0100 110	26	12	1100	1100 011	63
5	0101	0101 101	2D	13	1101	1101 000	68
6	0110	0110 001	31	14	1110	1110 100	74
7	0111	0111 010	3A	15	1111	1111 111	7F

表 1: Hamming(7, 4) 编码表

注意到，码书 (表 1) 中 16 个码字 *c* 的高 4 bit 正好是对应的 *m* 的二进制表示  $m_3 m_2 m_1 m_0$ ，低 3 bit 是校验位 (parity)。这种编码型制叫系统码 (systematic code)，也就是码字 *c* 中把消息 *m* 和校验位前后分开，如图 2 所示。

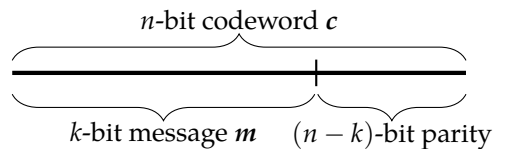


图 2: (n, k) 系统码的码字 *c*

为什么把 4-bit 的 *m* 编码为 7-bit 的 *c* 就能容错呢？因为码书  $\mathcal{C}$  中任意两个不同的码字  $c_i$  和  $c_j$  至少有 3 bit 不同，如果传输中码字 *c* 发生了单 bit 翻转，从接收到的 7 bit 数据 *r* 可以很容易推测其原本发送的消息。例如  $m = 1$ ，发送  $c_1 = 0001011$ ，传输时发生了最低 bit 翻转，收到的是  $r = 000101\textcircled{0}$ ，通过对比码书，我们发现 *r* 最有可能是从  $c_1$  变来，因为 *r* 与  $c_1$  只相差 1 bit，而与其他码字至少相差 2 bit (例如  $c_0, c_7, c_9$ )。

<sup>①</sup>[https://en.wikipedia.org/wiki/Hamming\\_code](https://en.wikipedia.org/wiki/Hamming_code)

<sup>②</sup>[https://en.wikipedia.org/wiki/Richard\\_Hamming](https://en.wikipedia.org/wiki/Richard_Hamming)

从而推断发送的是  $\hat{c} = c_1$ ，继而解码得到  $m = 1$ 。当然这种“查书”译码方法一般来说效率不算高，后面会介绍更高效的做法。

有兴趣的读者可以自行编程验证：从码书  $C$  中任取一个码字  $c$ ，任意翻转其中一个 bit，得到一个新的 7-bit 字符串  $r$ ，那么  $r$  与  $C$  中的其他 15 个码字至少有 2 bit 不同。由前面的构造方法可知， $r$  与  $c$  只相差 1 bit，对接收方来说， $r$  最有可能是从  $c$  变来，把  $r$  解码回  $c$  是最合理的。这就直接验证了 Hamming(7,4) 码总能正确纠正单 bit 错误。这种验证方法需要处理  $16 \times 7 \times 15 = 1680$  种情况，对于更长的码（例如 (72, 64) 码），这种暴力穷举式的验证方法会变得不可行，我们要用别的方法来验证其纠错能力，即后面讲的最小 Hamming 距离  $d_{\min}$ 。

而对于双 bit 或多 bit 错误，Hamming 码就无能为力了。特别是，如果翻转  $c$  中的两个 bits，那么  $r$  总是与另外的码字更像，即  $\hat{c} \neq c$ ，“纠错”反而会引入第 3 个 bit error。

### 1.1 Hamming(7, 4) 码的编码与译码实现

表 1 给出了 Hamming(7, 4) 码的码书，在现在内存很便宜的情况下，我们可以用简单的查表法实现编码与译码。编码用到一个长度为  $2^4 = 16$  的数组 `encode[16]`，每个元素是 7-bit 码字  $c$ ，总容量仅需 112 bits。译码要稍微复杂一些，用  $2^7 = 128$  的数组 `decode[128]`，每个元素是下标  $r$  对应的 4-bit 消息  $m$ ，总容量 512 bits。

Hamming 码发明的时候是 1940s 年代末，当时 Richard Hamming 用的是一台继电器组成的机电式计算机 (relay computer)<sup>①</sup>，他设计的纠错码必须能用极少的硬件来实现，实际上 Hamming 码可以用简单的组合逻辑来实现编解码。

Claude E. Shannon 在 1948 年发表的划时代论文<sup>[3]</sup>中就提到了 Hamming(7, 4) 码的编解码算法 (图 3)，他俩在 Bell Labs 曾共享同一间办公室<sup>②</sup>。这个原版的 Hamming 码只用异或逻辑就能直接定位单 bit 错误，图 3 的最后一句话“二进制数  $\alpha\beta\gamma$  给出了出错 bit 的下标”。

An efficient code, allowing complete correction of errors and transmitting at the rate  $C$ , is the following (found by a method due to R. Hamming):

Let a block of seven symbols be  $X_1, X_2, \dots, X_7$ . Of these  $X_3, X_5, X_6$  and  $X_7$  are message symbols and chosen arbitrarily by the source. The other three are redundant and calculated as follows:

$$\begin{aligned} X_4 \text{ is chosen to make } \alpha &= X_1 + X_5 + X_6 + X_7 \text{ even} \\ X_2 \text{ “ “ “ “ } \beta &= X_2 + X_3 + X_6 + X_7 \text{ “} \\ X_1 \text{ “ “ “ “ } \gamma &= X_1 + X_3 + X_5 + X_7 \text{ “} \end{aligned}$$

When a block of seven is received  $\alpha, \beta$  and  $\gamma$  are calculated and if even called zero, if odd called one. The binary number  $\alpha\beta\gamma$  then gives the subscript of the  $X_i$  that is incorrect (if 0 there was no error).

图 3: Shannon 1948 论文中提到了 Hamming(7, 4) 码

<sup>①</sup>[https://en.wikipedia.org/wiki/Model\\_V](https://en.wikipedia.org/wiki/Model_V)

<sup>②</sup>Hamming 在其著名的演讲 “You and Your Research” 中回忆到: [W]hen I came to Bell Labs, I shared an office for a while with Shannon. At the same time he was doing information theory, I was doing coding theory.

Shannon 论文提到的 Hamming(7, 4) 和现行教材里通常讲的版本略有不同，前面我们介绍的是系统码，其实只是码字  $c$  中 bits 的排列顺序不一样，没有本质区别。图 4 展示了这两种编码的转换。

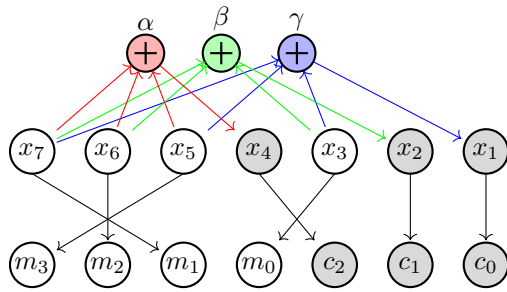


图 4: 原版 Hamming(7,4) 码字转换为“系统码”

Richard Hamming 在 1950 年 1 月申请了专利<sup>①</sup>之后再发表论文<sup>[4]</sup>，Hamming 论文的发表时间距 Hamming 码的发明晚了大约 18 个月。

编码实现不必查表，输入 4-bit 消息  $m = (m_3, m_2, m_1, m_0)$ ，输出 7-bit 码字  $c = (c_6, c_5, c_4, c_3, c_2, c_1, c_0)$ ，其对应关系如下，其中  $\oplus$  是异或运算。

- $(c_6, c_5, c_4, c_3) = (m_3, m_2, m_1, m_0)$
- $c_2 = m_1 \oplus m_2 \oplus m_3$
- $c_1 = m_0 \oplus m_1 \oplus m_2$
- $c_0 = m_0 \oplus m_1 \oplus m_3$

以上逻辑很容易转换为程序代码，也可以用数字电路实现（图 5）。

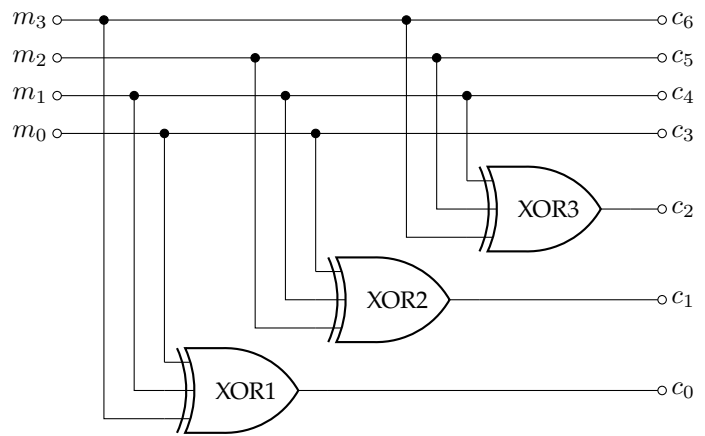


图 5: Hamming(7,4) 码的编码电路

对于纠错码（信道编码），一般来说译码比编码要复杂一些，信源编码（数据压缩）正好相反。Hamming 码算是译码比较简单的了，其他纠错能力更强的码的译码方法更复杂。

Hamming(7,4) 译码的总体框图见图 6。解码的基本流程分三步：

1. 从收到的 7-bit 向量  $r$  计算出 3-bit 向量  $s$ ，称为伴随式 (syndrome)，或校验子、校正子。 $s$  有 3 bit，能表示 8 个数，0 表示无错，1 ~ 7 分别指示  $r$  中的某一个 bit 出错， $s$  是错误的“症状”。
2. 从  $s$  映射为 7-bit 向量  $e$ ，称为错误图样 (error pattern)。 $s \neq 0$  时， $e$  只有某一个 1 bit 是 1。
3. 从  $r$  中扣除  $e$  得到译码结果  $\hat{c} = r \oplus e$

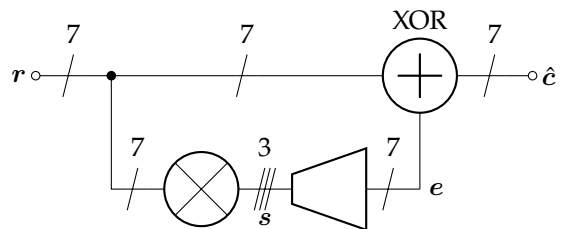


图 6: Hamming(7,4) 码的译码框图

<sup>①</sup>Richard W Hamming, Bernard D Holbrook. “Error-detecting and correcting system.” US Patent 2,552,629, filed Jan. 11, 1950, and issued May 15, 1951. <https://patents.google.com/patent/US2552629A>

以下 Python 代码实现了图 6 和 7 所示的译码过程。

```
def to_bits(x, N):
    bits = [0] * N
    for i in range(N):
        if (x & (1<<i)):
            bits[i] = 1
    return bits

def from_bits(bits: list[int]):
    x = 0
    for i, b in enumerate(bits):
        x |= (b << i)
    return x

def syndrome(r : list[int]):
    s0 = r[0] ^ r[3] ^ r[4] ^ r[6]
    s1 = r[1] ^ r[3] ^ r[4] ^ r[5]
    s2 = r[2] ^ r[4] ^ r[5] ^ r[6]
    return [s0, s1, s2]

def error_pattern(s : list[int]):
    [s0, s1, s2] = s
    e = [0] * 7
    e[0] = s0 & ~s1 & ~s2
    e[1] = ~s0 & s1 & ~s2
    e[2] = ~s0 & ~s1 & s2
    e[3] = s0 & s1 & ~s2
    e[4] = s0 & s1 & s2
    e[5] = ~s0 & s1 & s2
    e[6] = s0 & ~s1 & s2
    return from_bits(e)

def decode(r : int):
    s = syndrome(to_bits(r, 7))
    e = error_pattern(s)
    c = r ^ e
    return c >> 3
```

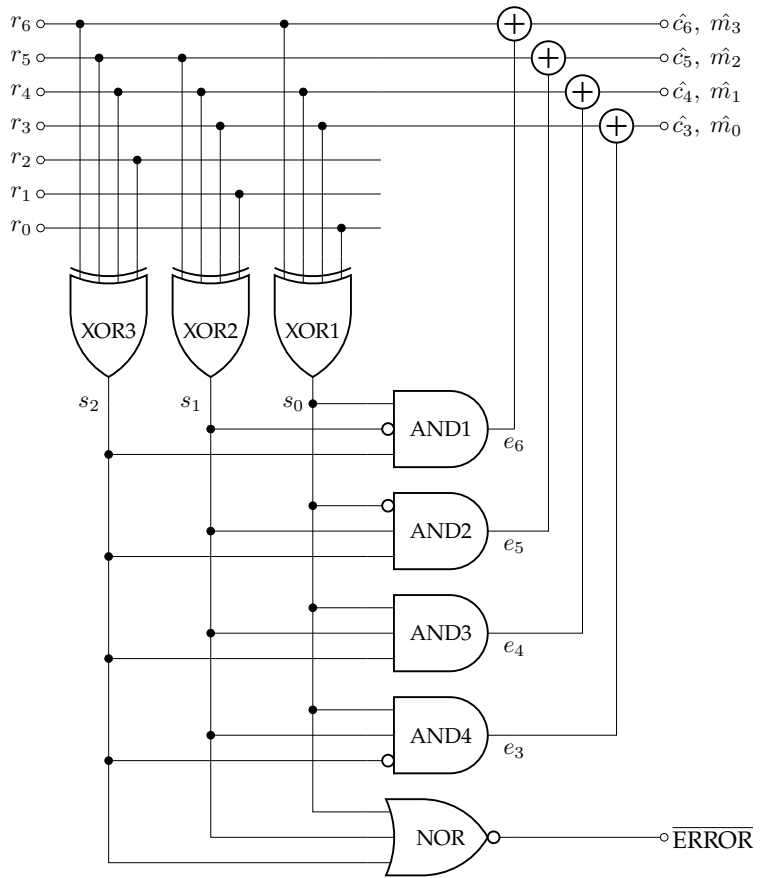


图 7: Hamming(7,4) 码的译码电路

图 7 是实际的译码电路，由于我们只关心最终译码的消息  $\hat{m}$ ，因此只用生成与  $\hat{c}_3, \hat{c}_4, \hat{c}_5, \hat{c}_6$  有关的  $e$ ，即只需解出  $e$  的高 4 bit。读者可以试试拿  $r = 0000011$  喂给以上实现，看看能否正确解出  $\hat{m} = 0001$ 。

## 1.2 Hamming 距离

因其纠错能力有限，Hamming 码现在已经没有多少实际用处了（见第 19 页的例外），但是 Hamming 引入的“Hamming 距离”概念在学习其他纠错码时仍然是必要的基础。

一个“码”之所以能够纠错，靠的是码字与码字不相像。两个汉字如果相像（形近），就容易认错，例如手写的“巳巳巳”、“人入”、“千千千”、“外处”、“拔拔”、“日曰”、“戎戎”等等，当然根据上下文我们一般不会误解。在手写票据时，需要防止涂改金额或日期，例如“一二三七十千”须写成“壹贰叁柒拾仟”，这样就不容易被人添一笔把“一”改为“二”或“十”。

码字的“像”与“不像”以 Hamming 距离论，所谓 Hamming 距离 (Hamming distance)<sup>①</sup>，本意

<sup>①</sup>[https://en.wikipedia.org/wiki/Hamming\\_distance](https://en.wikipedia.org/wiki/Hamming_distance)

是两个等长的向量有多少位置上的符号不同。对于纠错码而言，两个二进制向量  $\mathbf{x}$  和  $\mathbf{y}$  的 Hamming 距离就是它们有多少个 bit 不同，一般用  $d_H(\mathbf{x}, \mathbf{y})$  表示。易知  $d_H('000', '111') = 3$ 。

从数学上说，Hamming 距离  $d_H$  是一个度量 (metric)<sup>①</sup>，满足我们对于“距离”的直观认识：

1.  $d_H(x, x) = 0$ ，换言之  $d_H(x, y) = 0$  当且仅当  $x = y$
2.  $d_H(x, y) > 0$  if  $x \neq y$
3.  $d_H(x, y) = d_H(y, x)$
4.  $d_H(x, y) + d_H(y, z) \geq d_H(x, z)$

其实第 2 条不是独立的，可以通过另外三条推出来：

$$d_H(x, y) + d_H(x, y) = d_H(x, y) + d_H(y, x) \geq d_H(x, x) = 0$$

上式左起第一个等号用了性质 3（对称性），中间的  $\geq$  用了性质 4（三角不等式，令  $z = x$ ），最后一个等号用了性质 1，最终推出  $d_H(x, y) \geq 0$ 。

Hamming 距离其实就是二进制向量空间的  $L^1$  距离。<sup>②</sup>

在计算的时候，两个非负整数的 Hamming 距离即是这两个数异或 (xor) 之后得到的数的二进制表示里有多少个 bit 是 1。一个整数的二进制表示中有多少个“1”称为“Hamming 权重/Hamming 重量/Hamming weight”<sup>③</sup>，在现代的 CPU 上有专门的指令做这件事 (x86 是 POPCNT)，因此 C/C++ 中实现 Hamming 距离可以用以下简单的函数 (适用于 GCC/Clang 编译器)：

```
// Returns Hamming distance of two unsigned integers
uint32_t hamming_distance(uint32_t x, uint32_t y)
{
    return __builtin_popcount(x ^ y); // Or std::popcount() after C++20
}
```

对于 Hamming(7, 4) 码来说，每个码字长度为 7 bit，可以看成是 7 维空间中的一个点，因此很难在二维平面上画出全部  $2^k = 16$  个码字之间的距离关系。图 8 示意了  $c_0$  和  $c_1$  两个码字的距离  $d_H(c_0, c_1) = 3$ 。

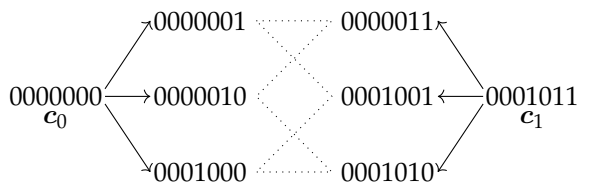


图 8: Hamming(7, 4) 其中两个码字的距离

**最小 Hamming 距离** 一个码  $C$  的检错与纠错能

力取决于其最小 Hamming 距离  $d_{\min}$ （简称“最小码距”），定义为任意两个不同码字  $(c_i, c_j)$  之间 Hamming 距离的最小值

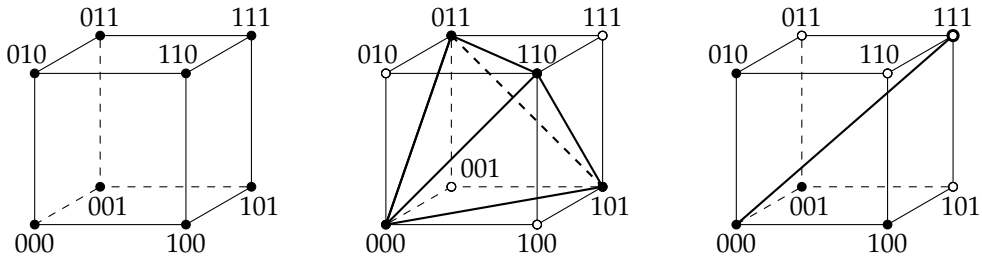
$$d_{\min}(C) = \min_{i \neq j} d_H(c_i, c_j)$$

这里我们先从图形的角度直观体验一下不同  $d_{\min}$  的纠错能力。图 9 以立方体的形式绘制了  $n = 3$  的几种码，便于观察码距。 $n = 3$  维 Hamming 空间有 8 个点，正好对于立方体的 8 个顶点。

<sup>①</sup>[https://en.wikipedia.org/wiki/Metric\\_space](https://en.wikipedia.org/wiki/Metric_space)

<sup>②</sup>[https://en.wikipedia.org/wiki/Taxicab\\_geometry](https://en.wikipedia.org/wiki/Taxicab_geometry)

<sup>③</sup>[https://en.wikipedia.org/wiki/Hamming\\_weight](https://en.wikipedia.org/wiki/Hamming_weight)

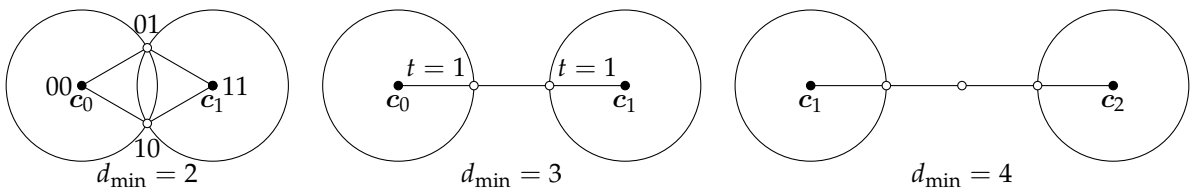


(a) ( $n = 3, k = 3, d_{\min} = 1$ ) (b) ( $n = 3, k = 2, d_{\min} = 2$ ) (c) ( $n = 3, k = 1, d_{\min} = 3$ )

图 9:  $n = 3$  的几种码

- 图 9(a) 中立方体的每个顶点都是码字,  $d_{\min}(C) = 1$ 。这个码没有纠错能力, 因为任何一个码字  $c$  翻转其中任一 bit 都是另一个合法的码字。
- 图 9(b) 中立方体的 4 个实心顶点是码字,  $C = \{000, 011, 101, 110\}$ ,  $d_{\min}(C) = 2$ 。这个码可以检测单 bit 错误, 因为任何一个码字, 如果翻转 1 bit, 会落到空心顶点, 表明传输有误。这个码是奇偶校验码 (Single parity-check, SPC), 可以看出码字的第 3 bit 其实就是前两个 bit 的异或。
- 图 9(c) 中立方体对角的两个顶点是码字,  $C = \{000, 111\}$ ,  $d_{\min}(C) = 3$ 。实心顶点表示距离左下角 (0,0,0) 较近, 空心顶点表示距离右上角 (1,1,1) 较近。这是简单的 3-重复码, 可以纠正单 bit 错误。

从空间的角度, 以码字为球心, 半径为  $t$  的 Hamming 球<sup>①</sup> 如果彼此不相交, 那么这个码就能纠正最多  $t$  个错误。下图中最左侧是  $C = \{00, 11\}$  的简单重复码, 其两个  $t = 1$  的球是相交的, 因此  $d_{\min} = 2$  的码只能检错不能纠错。如果发生了单 bit 错误, 接收方收到 “01” 这个向量, 并不能判断发送的是 “00” 还是 “11”, 因为概率一样。



上图中间  $d_{\min} = 3$  的一种具体情况可以参考图 9(c), 也就是  $c_0 = 000, c_1 = 111$ , 每个球连同球心由 4 个点构成。

从这些 Hamming 球的图形不难归纳出: 码  $C$  能纠正  $t = \lfloor \frac{d_{\min}(C) - 1}{2} \rfloor$  个错误, 这里的  $\lfloor x \rfloor$  表示向下取整, 例如  $\lfloor 2.5 \rfloor = 2, \lfloor 3 \rfloor = 3$ 。

$d_{\min}$  决定了一个码的纠错能力, 因此常缩写为  $d$ , 并与  $n, k$  一起作为码的重要参数标出来, 即 “[ $n, k, d$ ] 记法”。按此记法, Hamming(7, 4) 码是 [7, 4, 3] 码, 表明其  $n = 7, k = 4, d = 3$ , 能纠正  $t = 1$  个错误。表 2 列举了  $d_{\min}$  对应的检错与纠错能力。

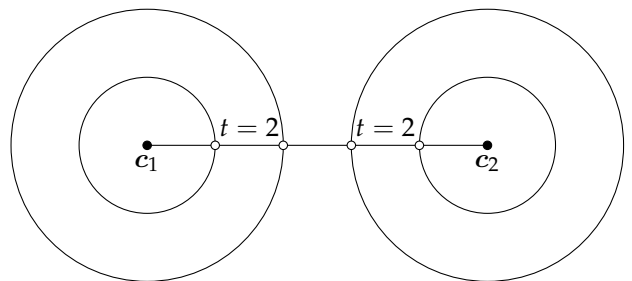


图 10:  $d_{\min} = 5$  情况下的 Hamming 球

<sup>①</sup>[https://en.wikipedia.org/wiki/Hamming\\_ball](https://en.wikipedia.org/wiki/Hamming_ball) 以  $c$  为球心, 满足  $d_H(x, c) \leq t$  的向量  $x$  构成

$d_{\min}$	$t = \lfloor \frac{d_{\min}-1}{2} \rfloor$	名称	纠错能力
1	0	无编码	不能检错或纠错
2	0	奇偶校验码, SPC	检测单 bit 错误
3	1	Hamming 码	纠正单 bit 错误
4	1	扩展 Hamming 码	检测两个错误或纠正单个错误 SEC-DED
5	2		纠正两个 bit 错误

表 2: 最小 Hamming 距离  $d_{\min}$  与纠错能力  $t$  的关系举例

Richard Hamming 本人在他的里程碑论文<sup>[4]</sup>中解决了  $d_{\min} \leq 4$  的情况, 因此习惯上把能纠正单 bit 错的码统称为 Hamming 码。

理论上存在能充分发挥纠错潜力的译码方法, 称为最小距离译码 minimum distance decoding。<sup>①</sup> 假如接收到一个向量  $r$ , 其中可能包含 bit 错, 我们在码书  $C$  中找一个与之距离最小的码字  $c$  作为译码结果  $\hat{c}$ :

$$\hat{c} = \arg \min_{c \in C} d_H(r, c)$$

在实践中, 因计算量太大, 只有很少的几种码能用这种方式译码。

回到本节的主题: Hamming(7, 4) 码, 表 3 列出了此码各个码字之间的距离, 从此表可看出  $d_{\min} = 3$ , 因此它能纠正  $t = \lfloor \frac{3-1}{2} \rfloor = 1$  个错误。这种验证方法只需要考察  $2^k(2^k - 1) = 16 \times 15 = 240$  种情况, 比前面的做法效率高一些, 但计算复杂度仍然高达  $O(2^{2k})$ , 我们还会找到更高效的验证方法。<sup>②</sup>

$m$	码字 $c$	Hex	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	000 0000	00	-	3	4	3	3	4	3	4	3	4	3	4	4	3	4	7
1	000 1011	0b	3	-	3	4	4	3	4	3	4	3	4	3	3	4	7	4
2	001 0111	17	4	3	-	3	3	4	3	4	3	4	3	4	4	7	4	3
3	001 1100	1c	3	4	3	-	4	3	4	3	4	3	4	3	7	4	3	4
4	010 0110	26	3	4	3	4	-	3	4	3	4	3	4	7	3	4	3	4
5	010 1101	2d	4	3	4	3	3	-	3	4	3	4	7	4	4	3	4	3
6	011 0001	31	3	4	3	4	4	3	-	3	4	7	4	3	3	4	3	4
7	011 1010	3a	4	3	4	3	3	4	3	-	7	4	3	4	4	3	4	3
8	100 0101	45	3	4	3	4	4	3	4	7	-	3	4	3	3	4	3	4
9	100 1110	4e	4	3	4	3	3	4	7	4	3	-	3	4	4	3	4	3
10	101 0010	52	3	4	3	4	4	7	4	3	4	3	-	3	3	4	3	4
11	101 1001	59	4	3	4	3	7	4	3	4	3	4	3	-	4	3	4	3
12	110 0011	63	4	3	4	7	3	4	3	4	3	4	3	4	-	3	4	3
13	110 1000	68	3	4	7	4	4	3	4	3	4	3	4	3	3	-	3	4
14	111 0100	74	4	7	4	3	3	4	3	4	3	4	3	4	4	3	-	3
15	111 1111	7f	7	4	3	4	4	3	4	3	4	3	4	3	3	4	3	-

表 3: Hamming(7, 4) 的码距

<sup>①</sup>[https://en.wikipedia.org/wiki/Decoding\\_methods#Minimum\\_distance\\_decoding](https://en.wikipedia.org/wiki/Decoding_methods#Minimum_distance_decoding) 也叫“最近邻居译码 nearest neighbour decoding”, 在离散无记忆信道 (DMC) 下等同于最大似然译码 (MLD)。

<sup>②</sup><https://cseweb.ucsd.edu/~daniele/Research/CodeComp.html> 总结来说, 找到线性码的最小码距是 NP-hard, 即  $O(2^k)$ 。

### 1.3 线性分组码

本节介绍线性分组码的一些最基本知识, 假设读者熟悉基础的线性代数运算(矩阵乘以向量), 本节的运算都是二进制算数, 又叫 mod 2 算数, 与常规四则运算惟一的区别是  $1 + 1 = 0$  和  $-1 = 1$ 。前文第 1.1 节我们介绍了 Hamming(7, 4) 码的编码与译码算法(主要以位运算实现), 验证了它确实能纠正单 bit 错误, 还画出对应的数字电路实现(图 5 和 7), 并在表 1 中给出了完整的码书。从实验和实用的角度, 似乎已经完成了任务。据 Richard Hamming 本人后来回忆, 他当时注意到了码字  $c \in \mathcal{C}$  有一定的代数结构, 例如在异或运算下 Hamming(7, 4) 的 16 个码字构成了一个群(group, 准确的说是阿贝尔群), 即任意两个码字的异或仍然是合法的码字( $c_1 \oplus c_2 = c_3, c_3 \oplus c_4 = c_7$  等等), 但他没有把这些数学内容写入论文中。

在 Hamming 发表这篇论文之后的十年中, 纠错码理论(coding theory)领域有了重大的进展。1956 年 David S. Slepian<sup>①</sup>用“子群 subgroup”、线性组合等数学语言来描述一大类分组码(block code)<sup>[5]</sup>, 也就是本节将要介绍的线性分组码。1961 年 W. Wesley Peterson<sup>②</sup>出版了第一本纠错码专著<sup>[6]</sup>, 定义了这一领域现在通用的话语体系(例如  $G$  矩阵与  $H$  矩阵)。

所谓线性码<sup>③</sup>, 指的是任意两个码字的线性组合仍然是合法的码字(若  $c_i, c_j \in \mathcal{C}$ , 则  $c_i + c_j \in \mathcal{C}$ ), 即码书  $\mathcal{C}$  是一个线性空间(又称向量空间<sup>④</sup>)。对于二进制码(binary code)来说, “线性组合”其实就是异或  $\oplus$ ,  $c_i + c_j = c_i - c_j = c_i \oplus c_j$ 。现在实用的纠错码似乎都是线性码, 包括 5G 里用到的 LDPC 码和 Polar 码都是线性分组码。

利用线性分组码的定义条件  $c_i \oplus c_j \in \mathcal{C}$  可以得到几条有用的性质:

1. 全 0 码字  $\mathbf{0} \in \mathcal{C}$ ,  $\mathbf{0}$  表示适当维度的全 0 向量, 因为  $c_1 \oplus c_1 = \mathbf{0} \in \mathcal{C}$
2. 根据  $d_H(c_i, c_j) = d_H(c_i - c_j, \mathbf{0})$ , 而  $c_i - c_j \in \mathcal{C}$ , 因此  $d_{\min}(\mathcal{C})$  是非  $\mathbf{0}$  码字的 Hamming 权重的最小值

$$d_{\min}(\mathcal{C}) = \min_{i \neq j} d_H(c_i, c_j) = \min_{c \neq \mathbf{0}} d_H(c, \mathbf{0})$$

判断一个码的纠错能力就不再需要像表 3 那样检查所有的码字对(pair), 只需要遍历所有码字即可。

$[n, k]$  线性分组码(linear block code, LBC)是指把  $k$ -bit 的消息  $m$  通过线性变换映射为  $n$ -bit 码字  $c$ , 线性变换其实就是用矩阵乘以向量  $m$ 。这样得到的码书  $\mathcal{C}$  是二元域  $\mathbb{F}_2$  上的  $n$  维线性空间  $\{0, 1\}^n$  中的一个  $k$  维线性子空间<sup>⑤</sup>。以上定义似乎有些抽象, 看几个例子就能明白。

线性分组码的一大好处是可以用生成矩阵(generator matrix)来方便地描述其编码过程。对于  $[n, k]$  线性分组码, 生成矩阵  $G$  是  $k$  行  $n$  列矩阵, 即  $k \times n$  阶矩阵。有了生成矩阵  $G$ , 线性码的编码过程就可以非常紧凑地写成  $c = m \cdot G$  或  $c = G^T \cdot m$ , 表示码字  $c$  是  $G$  行向量的线性组合。为了能唯一译码,  $G$  的  $k$  个行向量必须是线性无关的。线性分组码  $\mathcal{C}$  是生成矩阵  $G$  的行空间<sup>⑥</sup>,  $\mathcal{C} = \text{Row}(G)$ 。用  $G$  矩阵定

<sup>①</sup>[https://en.wikipedia.org/wiki/David\\_Slepian](https://en.wikipedia.org/wiki/David_Slepian) 是 Shannon 和 Hamming 在 Bell 实验的同事, 获得了 1974 年 Shannon 奖。

<sup>②</sup>[https://en.wikipedia.org/wiki/W.\\_Wesley\\_Peterson](https://en.wikipedia.org/wiki/W._Wesley_Peterson) 获得 1981 年 Shannon 奖, 发明了循环冗余校验(CRC)和最早的 BCH 码译码算法(PGZ 算法, Peterson-Gorenstein-Zierler)

<sup>③</sup>[https://en.wikipedia.org/wiki/Linear\\_code](https://en.wikipedia.org/wiki/Linear_code)

<sup>④</sup>[https://en.wikipedia.org/wiki/Vector\\_space](https://en.wikipedia.org/wiki/Vector_space)

<sup>⑤</sup>[https://en.wikipedia.org/wiki/Linear\\_subspace](https://en.wikipedia.org/wiki/Linear_subspace)

<sup>⑥</sup>[https://en.wikipedia.org/wiki/Row\\_and\\_column\\_spaces](https://en.wikipedia.org/wiki/Row_and_column_spaces)

一个线性分组码只需要  $n \times k$  bits，而不需要列出全部  $2^k$  个码字；对于系统码， $G$  的左边是  $k$  阶单位阵  $I_k$ ，那只需  $(n - k) \times k$  bits 定义剩余部分即可。

表 1 所示 Hamming(7, 4) 码对应的生成矩阵  $G$  大小是  $4 \times 7$ ：

$$G = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 & 1 \end{bmatrix}_{4 \times 7} = \begin{bmatrix} -g_3- \\ -g_2- \\ -g_1- \\ -g_0- \end{bmatrix}$$

其中  $g_i$  是矩阵  $G$  的行向量。 $g_0, g_1, g_2, g_3$  这几个向量其实是空间  $C$  的一组基 (basis)。

对于 4-bit 消息  $m = [m_3 m_2 m_1 m_0]$  而言，码字  $c$  是 7-bit 基向量  $g_0, g_1, g_2, g_3$  的线性组合  $c = m_0 g_0 + m_1 g_1 + m_2 g_2 + m_3 g_3$ ，当然，这里的加法是“模 2 加”，也就是异或运算  $\oplus$ 。

$$c = m \cdot G = [m_3 \ m_2 \ m_1 \ m_0] \cdot \begin{bmatrix} -g_3- \\ -g_2- \\ -g_1- \\ -g_0- \end{bmatrix} = m_3 g_3 + m_2 g_2 + m_1 g_1 + m_0 g_0$$

例如对于  $m = [0 \ 1 \ 0 \ 1]$ ，对应的码字  $c_5 = [0 \ 1 \ 0 \ 1] \cdot G = g_0 \oplus g_2 = 0101101$ 。

可能是受最早 Peterson 教材<sup>[6]</sup>的影响，纠错码领域习惯用行向量来表示  $m$  和  $c$ ，因此  $c = m \cdot G$  符合矩阵运算的规则，表示  $c$  是  $G$  行向量的线性组合。而在其他领域， $y = Ax$  是更常见的写法，意为将矩阵  $A$  所表示的线性变换作用于列向量  $x$ ，所得的列向量  $y$  是矩阵  $A$  的列向量的线性组合。这两种记法没有本质区别，只是不同领域的习惯不同罢了。

本文为了简单起见，在不至于引起歧义的情况下，对于矩阵与向量相乘，不仔细区分行向量和列向量。例如对于  $4 \times 7$  的  $G$ ，如果  $m$  是 4-bit 向量，那么  $m \cdot G$  和  $G^T \cdot m$  都表示对  $G$  的 7-bit 行向量做线性组合，运算结果是 7-bit 向量  $c$ 。这也符合 Python 3.5 引入的  $@$  运算符的规则。<sup>①</sup>

用 Python 实现 Hamming(7, 4) 编码（生成矩阵法）：

```
def encode(m : int) -> int:
    g3 = 0b1000101
    g2 = 0b0100110
    g1 = 0b0010111
    g0 = 0b0001011

    c = 0b0000000
    if m & 1:
        c ^= g0
    if m & 2:
        c ^= g1
    if m & 4:
        c ^= g2
    if m & 8:
        c ^= g3
    return c
```

<sup>①</sup><https://peps.python.org/pep-0465/>

```
# 以上代码等同于
[m0, m1, m2, m3] = to_bits(m, 4)
c = m0 * g0 ^ m1 * g1 ^ m2 * g2 ^ m3 * g3
return c
```

生成矩阵  $G$  对编码很有用，对于线性分组码的译码，我们要引入监督矩阵  $H$ ，又称一致校验矩阵 (parity-check matrix)，定义为“ $Hc = \mathbf{0}$  当且仅当  $c \in \mathcal{C}$ ”。也就是说， $\mathcal{C}$  是矩阵  $H$  的零空间<sup>①</sup>  $\mathcal{C} = \text{Nul}(H)$ ，或者说  $\mathcal{C}$  是齐次线性方程组  $Hc = \mathbf{0}$  的解空间。

对于  $[n, k]$  线性分组系统码，码字有  $n$  bit，其中  $k$  bit 是信息位，其余  $n - k$  bit 是监督位。 $H$  矩阵有  $n - k$  行  $n$  列，每一行对应一个监督位的校验规则，每一列对应码字中的一个 bit。 $Hr = \mathbf{0}$  表明  $r$  中每个监督位都符合校验规则，意味着  $r$  是合法码字。

假设码字  $c$  在传输时发生了一些 bit 翻转，错误图样  $e$  是  $n$ -bit 向量，接收到的向量  $r = c + e$ ，那么用  $H$  可以求出  $(n - k)$ -bit 伴随式 (syndrome)  $s$ ，这是译码的关键

$$s = Hr = H(c + e) = Hc + He = \mathbf{0} + He$$

表明伴随式  $s$  只取决于错误图样  $e$ ，和发送的码字  $c$  无关。 $s$  等于发生错误的位置对应的  $H$  列向量之和。

如果  $s = \mathbf{0}$ ，说明  $e = \mathbf{0}$  或者发生了无法检测的错误（例如  $e \in \mathcal{C}$ ）。

如果  $s \neq \mathbf{0}$ ，就需要设法从方程组  $He = s \pmod{2}$  解出可能性最大的  $e$ ，然后完成译码  $\hat{c} = r - e = c + e - e = c$ 。方程组  $He = s$  有  $n - k$  个方程，对应  $n - k$  个监督位，而  $e$  包含  $n$  个未知数，因此这是个不定方程，解空间有  $k$  维，一共  $2^k$  个解。求出一个特解不难，但是求得 Hamming 权重最小的解一般情况下是 NP-hard。<sup>[7]</sup>

对于系统码，从  $G$  求出  $H$  是很容易的：

$$G = [I_k \mid P], \quad H = [-P^\top \mid I_{n-k}]$$

这里  $I_k$  表示  $k \times k$  的单位矩阵， $P$  是  $k \times (n - k)$  矩阵，表示校验位和消息  $m$  的计算关系。我们先用这个方法找出 Hamming(7, 4) 码的监督矩阵  $H$ ，已知

$$G = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 0 & 1 & 1 \end{bmatrix} = [I_4 \mid P], \quad P = \begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 1 \end{bmatrix}_{4 \times 3}, \quad P^\top = \begin{bmatrix} 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \end{bmatrix}$$

对于二进制码，有  $-1 = 1$ ，从而得到

$$H = [-P^\top \mid I_3] = \begin{bmatrix} 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 \end{bmatrix}$$

如果  $e = [0, 1, 0, 0, 0, 0, 0]$ ，有  $s = He = [1, 1, 0]$ ，刚好是  $H$  的左起第 2 列，对应了发生错误的位置。

<sup>①</sup>[https://en.wikipedia.org/wiki/Kernel\\_\(linear\\_algebra\)](https://en.wikipedia.org/wiki/Kernel_(linear_algebra))

这个监督矩阵  $H$  表示了 7-bit 码字  $\mathbf{c}$  中的 3-bit 监督位  $c_2, c_1, c_0$  和信息位  $c_6, c_5, c_4, c_3$  的校验规则：

$$H = \begin{matrix} & c_6 & c_5 & c_4 & c_3 & c_2 & c_1 & c_0 \\ \begin{bmatrix} 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 \end{bmatrix} & c_2 = c_4 + c_5 + c_6 \\ & & & & & & & c_1 = c_3 + c_4 + c_5 \\ & & & & & & & c_0 = c_3 + c_4 + c_6 \end{matrix}$$

对于合法码字  $\mathbf{c} \in \mathcal{C}$ ，这三条校验规则均成立，则必有  $H\mathbf{c} = \mathbf{0}$ 。

$H$  矩阵取自 Hamming 的首字母，他本人在论文中提到：

This is equivalent to the earlier definition. To show this we form a matrix whose  $i$ -th row has 1's in the positions of the  $i$ -th parity check and 0's elsewhere. By assumption 1 the matrix is fixed and does not change from code symbol to code symbol. From 2 the rank of the matrix is  $k$ . This in turn

图 11: Hamming 论文谈及 parity-check matrix 的段落

这里关键的是这句：[A] matrix whose  $i$ -th row has 1's in the positions of the  $i$ -th parity check and 0's elsewhere. 这正是监督矩阵的意义。注意 Hamming 原文中的  $k$  和现在的意义不同，他用  $k$  表示有几个监督位，即现在的  $n - k$ 。

这里简述一下前面从  $G$  构造  $H$  的方法为何成立，已知  $G = [I_k | P]$ ,  $H = [-P^\top | I_{n-k}]$ ，按分块矩阵的乘法

$$HG^\top = [-P^\top | I_{n-k}] \begin{bmatrix} I_k \\ P^\top \end{bmatrix} = -P^\top I_k + I_{n-k} P^\top = -P^\top + P^\top = \mathbf{0}_{(n-k) \times k}$$

因此对于任何  $k$ -bit 消息  $\mathbf{m}$ ， $HG^\top \mathbf{m} = H(G^\top \mathbf{m}) = H\mathbf{c} = (HG^\top) \mathbf{m} = \mathbf{0} \mathbf{m} = \mathbf{0}$ 。另外，传统的纠错码书籍一般会用到  $GH^\top = \mathbf{0}$ ，除了用前面的方法证，还可以  $GH^\top = (HG^\top)^\top = (\mathbf{0}_{(n-k) \times k})^\top = \mathbf{0}_{k \times (n-k)}$ 。

从  $H$  矩阵的  $n$  列之间的线性相关性可以判断一个码的纠错能力，即  $d_{\min}$ 。假设一个  $H$  有  $a, b, c, d$  四列是线性相关的，即这些列按位异或得到全 0 向量。既然  $H$  矩阵的每一列对应码字中的一个 bit，那么一个只有第  $a, b, c, d$  bit 为 1，其余 bit 为 0 的码字  $\mathbf{c}$ ，它一定是合法码字  $H\mathbf{c} = \mathbf{0}$ ，这个码字与全 0 码字的距离为  $d_H(\mathbf{c}, \mathbf{0}) = 4$ 。

关于  $H$  有一个重要定理：如果  $H$  的任何  $d - 1$  或更少的列都线性无关，则  $d_{\min}$  至少是  $d$ 。此外， $d_{\min} = d$  的充要条件是  $H$  有  $d$  列线性相关，并且少于  $d$  列都线性无关。这个定理的详细证明在教材上都有，这里就不重复了。它有几个简单的推论：

- 如果  $H$  包含全 0 列，则  $d_{\min} = 1$ ，即这个码不能检错或纠错。理由如下： $H$  的每一列对应码字的一个 bit，如果某一列是全 0，则码字  $\mathbf{c}$  中对应的这一 bit 是 0 或 1 都不影响  $H\mathbf{c} = \mathbf{0}$ ，因此两个合法码字只相差 1 bit。
- 如果  $H$  不包含全 0 列，但包含两个相同的列，则  $d_{\min} = 2$ ，即这个码最多只能检测单 bit 错，但不能纠错。理由是，假设  $H$  中第  $i$  列和第  $j$  列相同且均不为  $\mathbf{0}$ ，则码字中对应的第  $i$  位和第  $j$  位可以互换 ( $10 \leftrightarrow 01$ ) 或同时翻转 ( $00 \leftrightarrow 11$ ) 而不影响  $H\mathbf{c} = \mathbf{0}$ ，因此存在两个合法码字只相差 2 bit。

- 如果  $H$  不包含全 0 列，并且每一列都不同，意味着任何两列都线性无关，则  $d_{\min} \geq 3$ 。Hamming 码正是如此，例如 Hamming(7, 4) 码的  $H$  矩阵有 7 列，包含全部 7 个非 0 的 3-bit 向量。
- $H$  矩阵是  $(n - k) \times n$  阶，而且  $\text{rank}(H) = n - k$ ，那么最多有  $n - k$  列线性无关， $n - k + 1$  列必然线性相关。这意味着  $d_{\min} \leq n - k + 1$ ，这是  $[n, k]$  码能达到的最大  $d_{\min}$ ，称为 Singleton 界<sup>①</sup>。能达到这个界的码叫 MDS 码 (maximum distance separable)，Reed-Solomon 码就是 MDS 码，可以说把  $n - k$  个校验位用到了极致。

交换  $H$  矩阵的两列不会影响码的纠错能力，这样我们可以获得一组等价的码。比方说 Wikipedia 上介绍的 Hamming(7, 4) 码是另一种构造<sup>②</sup>

$$G = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{bmatrix} \quad H = \begin{bmatrix} 1 & 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 \end{bmatrix}$$

这与我们前面给出的定义是等价的。不过尽管纠错能力一样，有时候不同的  $H$  对应的硬件电路实现的复杂度（成本）和运行速度（延迟）有所区别，这是实践中需要考虑的。例如如果再把 Hamming(7, 4) 码的监督矩阵的列重排成循环码的形式（如下），就可以用线性反馈移位寄存器 (LFSR) 电路实现编码。

$$H = \begin{bmatrix} 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 \end{bmatrix}$$

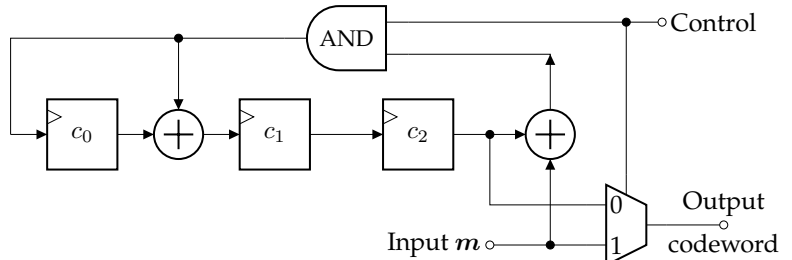


图 12: Hamming(7, 4) 循环码的移位寄存器编码电路

接下来，我们写一段 Python

程序从  $H$  矩阵求出  $d_{\min}$ 。

```

from itertools import combinations
import numpy as np

def find_d_min(H):
    n_rows, n_cols = H.shape
    d_min = -1
    for d in range(1, n_cols + 1): # d 从 1 到 n
        # 从全部 n 列中选 d 个不同的列，其下标位于数组 cols
        for cols in combinations(range(n_cols), d):
            s = np.zeros(n_rows, dtype=int)
            for i in cols: # 计算这些列之和
                s ^= H[:, i] # s ^= column[i]
            if s.sum() == 0: # 如果找到了线性相关的 d 列
                print(cols)
                d_min = d
                break
    if d_min > 0:
        break
    return d_min

```

<sup>①</sup>[https://en.wikipedia.org/wiki/Singleton\\_bound](https://en.wikipedia.org/wiki/Singleton_bound)

<sup>②</sup>[https://en.wikipedia.org/wiki/Hamming\\_code#Construction\\_of\\_G\\_and\\_H](https://en.wikipedia.org/wiki/Hamming_code#Construction_of_G_and_H)

```
print(find_d_min(np.array([[1, 1, 0, 0],
                          [1, 0, 1, 0],
                          [1, 0, 0, 1]],
                          dtype=int)))

# 输出:
# (0, 1, 2, 3)
# 4
```

以上代码中的  $3 \times 4$  示例  $H$  矩阵是四次重复码  $C = \{0000, 1111\}$  的监督矩阵，易知  $d_{\min} = 4$ 。

## 1.4 Hamming 码的构造原理

前面我们仔细研究了 Hamming(7, 4) 码的编码和解码，这是 Hamming 码的一个特例。Hamming 码是一类线性分组码，有  $m$  个监督位 ( $m \geq 2$ ) 的 Hamming 码是  $[n = 2^m - 1, k = n - m, d_{\min} = 3]$  码，其监督矩阵  $H$  有  $m$  行  $2^m - 1$  列，列向量包含全部  $2^m - 1$  个非 0 的  $m$  bit 向量。

监督位数 $m$	$n = 2^m - 1$	$k = 2^m - m - 1$	注释
2	3	1	三次重复码
3	7	4	教科书常用示例
4	15	11	
5	31	26	
...	...	...	
8	255	247	

表 4: Hamming 码的监督位数  $m$  与  $[n, k]$  的关系

Hamming(7, 4) 码因为监督位和信息位数量差不多，不容易体现其设计思路。看 Hamming(15, 11) 码就更清晰一些。Hamming 码是单 bit 纠错码，也就是说它的监督位要能定位码字中的单 bit 错误的位置。从信息论的角度，4 个监督位能定位  $2^4 = 16$  个位置，如果全 0 监督位表示无误，那么码字的最大长度是 15，剩下 11 bit 是信息位。Hamming 设计的巧妙之处在于，把出错位置用二进制表示，他让每个监督位负责这个位置信息中的一个 bit。以表 5 为例，把 15 个码字 bit 排成一行，选其中 2 的幂次的位置 (1, 2, 4, 8) 为监督位 (红字)，其余为信息位。第一行表示  $c_1 = c_3 + c_5 + \dots + c_{15}$ ，以此类推。假如第 10 个 bit 出错，那么只有  $p_1$  和  $p_3$  为 1，这就定位了这个错误。表 5 其实是非系统码的监督矩阵。

bits	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
$p_0$	1	0	1	0	1	0	1	0	1	0	1	0	1	0	1
$p_1$	0	1	1	0	0	1	1	0	0	1	1	0	0	1	1
$p_2$	0	0	0	1	1	1	0	0	0	0	0	1	1	1	1
$p_3$	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1

表 5: Hamming(15, 11) 码的 4 个监督位分别负责码字  $c$  中的 8 个 bit

由于每个监督位要负责检查大约一半的信息位，Hamming 码大概可以算作高密度奇偶校验码。

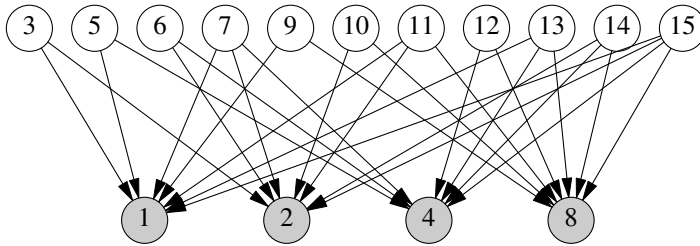


图 13: Hamming(15, 11) 码的 4 个监督位与 11 个信息位的关系

Hamming 码是完备码 (perfect code)<sup>①</sup>，意思是它的空间利用率高，不浪费。例如 Hamming(7, 4) 码有  $2^4 = 16$  个码字，每个码字是 7-bit。以码字  $c$  为球心，半径 (Hamming 距离) 为 1 的 Hamming 球有  $1 + 7 = 8$  个点。16 个码字对应的 16 个 Hamming 球刚好填满整个 7-bit 空间 ( $2^7 = 128$  个点)，因为  $8 \times 16 = 128$ 。这意味着，任何一个 7-bit 二进制向量，其要么是一个合法的码字；要么与某一个合法码字的距离为 1，而与其他码字的距离大于等于 2。

对于一般的有  $m$  个监督位的 Hamming 码，其  $n = 2^m - 1, k = n - m$ 。那么它共有  $2^k$  个码字，每个码字对应的  $t = 1$  的 Hamming 球有  $n + 1 = 2^m$  个点，这  $2^k$  个球一共有  $2^k \cdot 2^m = 2^{k+m} = 2^n$  个点，刚好填满整个  $n$ -bit 空间。就是说，对于任何接收到的  $n$ -bit 向量  $r$ ，刚好能无异议地译码成某个码字。

## 1.5 扩展 Hamming 码 (SEC-DED)

Hamming 码的  $d_{\min} = 3$ ，意味着它能纠正单 bit 错误。Hamming 在原始论文中指出，如果给码字添加 1 bit，就能让  $d_{\min} = 4$ ，即可以纠正单 bit 错，或者检测到双 bit 错 (Single Error Correction, Double Error Detection, 缩写 SEC-DED)，这种码叫扩展 Hamming 码 (extended Hamming code)，记为  $[n = 2^m, k = 2^m - m - 1, d_{\min} = 4]$ 。例如 Hamming(8, 4, 4) 码的生成矩阵  $G$  和监督矩阵  $H$  分别为

$$G = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & 0 & 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 1 & 1 \end{bmatrix} \quad H = \begin{bmatrix} 1 & 1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 0 & 1 \end{bmatrix}$$

与第 12 页定义的 Hamming(7, 4) 码相比，这里的生成矩阵多了最右列，监督矩阵多了一行一列。

用第 13 页的 `find_d_min` 函数可以方便地验证  $d_{\min} = 4$ 。

```
P = np.array([
    [1,1,1,0],
    [0,1,1,1],
    [1,0,1,1],
    [1,1,0,1]], dtype=int)
H = np.concatenate([P, np.eye(4, dtype=int)], axis=1)
print(H)
print(find_d_min(H))
```

<sup>①</sup>[https://en.wikipedia.org/wiki/Hamming\\_bound](https://en.wikipedia.org/wiki/Hamming_bound)

```

# 输出:
[[1 1 1 0 1 0 0 0]
 [0 1 1 1 0 1 0 0]
 [1 0 1 1 0 0 1 0]
 [1 1 0 1 0 0 0 1]]
(0, 1, 2, 4)
4

```

这个 Hamming(8, 4) 码的硬件编码电路很简单，如果用组合逻辑实现，电路类似图 5，增加一个三输入异或门即可。由于其校验位的规律性（每个校验位负责三个循环相邻的消息位），时序逻辑电路也容易用循环移位寄存器实现，并可复用三个三输入异或门，如图 14 所示。

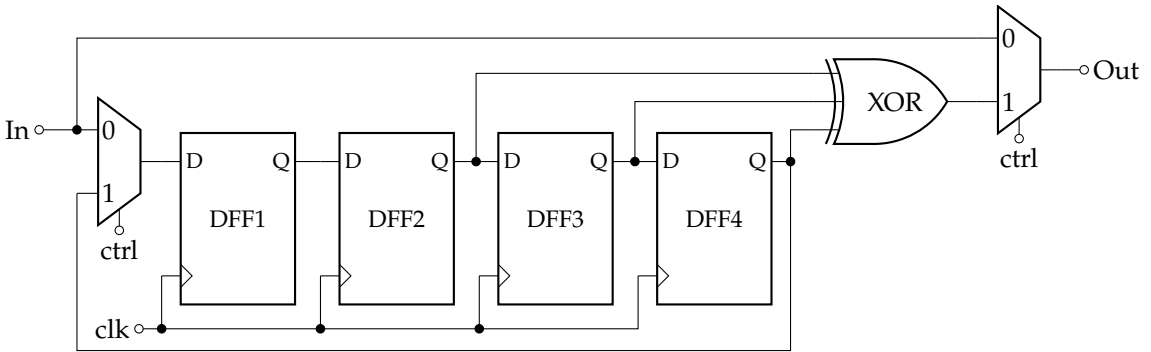


图 14: Hamming(8, 4) 码的一种编码电路

Hamming(8, 4) 码有 4 个监督位，伴随式 (syndrome)  $s$  也是 4-bit，一共有  $2^4 = 16$  种可能取值：全 0 表示无错；另有 8 个值表示 8-bit 码字中单 bit 错的错误位置，可以纠错；剩下 7 个值表示发生了双 bit 错，但无法确定具体位置。具体  $s$  的含义见表 6。从信息论的角度，4-bit 监督位不足以惟一表达全部  $\binom{8}{0} + \binom{8}{1} + \binom{8}{2} = 1 + 8 + 28 = 37$  种错误图样  $e$ 。

$s$	$e$
0000	00000000
1011	10000000
1101	01000000
1110	00100000
0111	00010000
1000	00001000
0100	00000100
0010	00000010
0001	00000001

$s$	$e$ 可能值
0011	00000011, 00010100, 01100000, 10001000
0101	00000101, 00010010, 01001000, 10100000
0110	00000110, 00010001, 00101000, 11000000
1001	00001001, 00110000, 01000100, 10000010
1010	00001010, 00100100, 01010000, 10000001
1100	00001100, 00100010, 01000001, 10010000
1111	00011000, 00100001, 01000010, 10000100

表 6: Hamming(8, 4) 码的伴随式  $s$  与错误图样  $e$  的关系

如果  $s = 0$ ，表示  $r = c + e = c + 0 = c$ ，意味着没有错误或者错误太多超出检测能力。否则， $s \neq 0$  表示发生了错误，其中再分为：

- 如果  $s$  有奇数个 1，表示发生了单 bit 错，可以纠正；
- 如果  $s$  有偶数个 1，表示发生了双 bit 错（或双数 bit 错），但无法纠正。

以上判断逻辑可以用图 15 所示的组合逻辑电路表达出来。

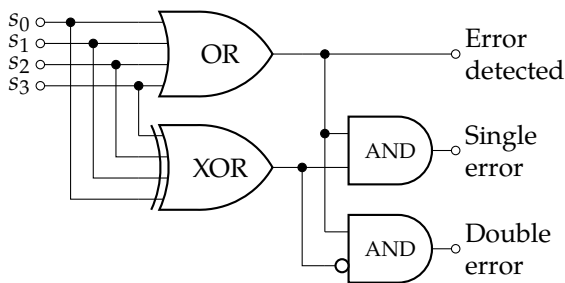


图 15: Hamming(8, 4) 码的部分解码电路

## 1.6 Hsiao 码与 ECC 内存

在实际的存储系统中，为了便于寻址，我们希望信息位数  $k$  是 2 的幂次，例如  $k = 16, k = 32, k = 64$  等等。Mu-Yue Hsiao<sup>①</sup> 为 1961 年面世的 IBM 7030 晶体管超级计算机的磁芯内存和硬盘设计了  $d_{\min} = 4$  的 SEC-DED 码<sup>②</sup>，一般称为 Hsiao 码 [9]。Hsiao 码有多种长度，例如 [22, 16]、[39, 32]、[72, 64] 等等。以下是 [22, 16] 和 [39, 32] 两种 Hsiao 码的监督矩阵，前者曾用于 IBM 7950 的磁带机<sup>③</sup>（22 条磁道，其中 16 条存信息位，6 条放校验位），后者曾用于 IBM 7030 的硬盘驱动器（39 个盘面，其中 32 个信息位，7 个校验位）。

$$H_{16} = \begin{bmatrix} 11111100 & 00100010 & 100000 \\ 11100011 & 11001000 & 010000 \\ 10001011 & 10000111 & 001000 \\ 00000110 & 01110111 & 000100 \\ 01010001 & 01111100 & 000010 \\ 00111100 & 10011001 & 000001 \end{bmatrix} \quad H_{32} = \begin{bmatrix} 11111111 & 00000010 & 00010010 & 10000011 & 1000000 \\ 00001001 & 11111111 & 00100100 & 10000100 & 0100000 \\ 00010000 & 00010000 & 11111111 & 00110110 & 0010000 \\ 00100010 & 00100101 & 10000000 & 11111111 & 0001000 \\ 01100101 & 01001001 & 00001111 & 01101000 & 0000100 \\ 10000110 & 10001110 & 11111000 & 00001000 & 0000010 \\ 11011000 & 11110000 & 01000001 & 01010001 & 0000001 \end{bmatrix}$$

Hsiao 码的监督矩阵的设计思路是

1. 没有全 0 列，保证  $d_{\min} \geq 2$
2. 没有相同的列，保证  $d_{\min} \geq 3$
3. 每一列有奇数个 1，保证  $d_{\min} \geq 4$
4. 整个矩阵的 1 数量尽量少
5. 每一行 1 的数量最好相等，或尽量接近

前三条是保证 Hsiao 码的纠错能力满足 SEC-DED，后两条是方便硬件实现，提高性能。第 4 条的目的是降低电路成本，因为每个 1 对应一个异或门的输入管脚，监督矩阵中的 1 越多，需要的电路门数

<sup>①</sup>网上资料不多，推测姓萧。他 1933 年 7 月出生在湖南长沙，在台北念了高中和大学本科。1958 年去美国，1960 年从 Illinois 大学获得了数学硕士学位，然后进入 IBM 工作了几年，1967 年从 Florida 大学拿了博士学位。此后长期在 IBM 纽约州 Poughkeepsie 市系统研发部门工作，1984 年当选 IBM Fellow。

<sup>②</sup>[http://www.bitsavers.org/pdf/ibm/7030/ADescriptionOfStretch\\_Dec59.pdf](http://www.bitsavers.org/pdf/ibm/7030/ADescriptionOfStretch_Dec59.pdf) 第 13 页和第 26 页

<sup>③</sup><http://www.brouhaha.com/~eric/retrocomputing/ibm/stretch/>

也越多。第5条的目的是降低电路延迟，提高运行速度，因为数字电路的运行速度（时钟频率）受限于最慢的路径（关键路径）的延迟。

例如 Hsiao(22, 16) 码有 6 个监督位，其监督矩阵  $H_{16}$  每行有 9 个 1，意味着每个监督位负责码字中的 9 bit，硬件实现可用两层共  $3 + 1 = 4$  个三输入异或门。其整个编码和解码电路可以复用 19 个三输入异或门，用于生成监督位或计算伴随式，电路的延迟是两级组合逻辑。

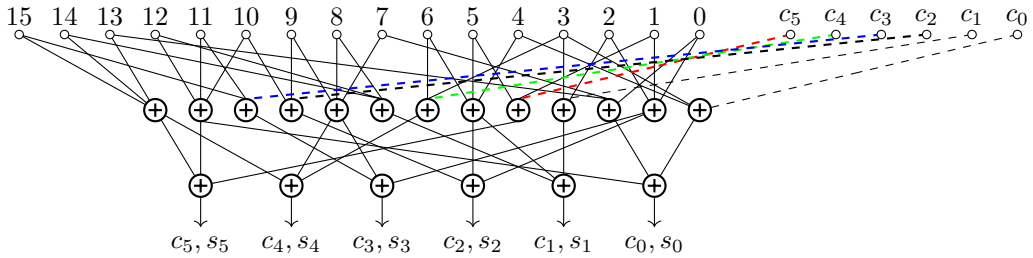


图 16: Hsiao(22, 16) 码监督矩阵  $H_{16}$  的硬件电路，共 19 个三输入异或门

图 16 所示的电路可同时用于编码和解码。在编码的时候，输入是 16-bit 消息向量  $m$ ，输出是 6 个监督位（码字  $c$  的低 6 bit），虚线部分接 0。在解码的时候，输入是 22-bit 码字（虚线部接的是监督位），输出是 6-bit 伴随式  $s$ 。

Hsiao(72, 64) 码曾广泛用于服务器或工作站上的 ECC 内存<sup>①</sup>，它以 8 个字节为一组，添加 1 个字节的监督位，就实现了 SEC-DED。表 7 是 Hsiao(72, 64) 码一种监督矩阵。

byte 0	byte 1	byte 2	byte 3	byte 4	byte 5	byte 6	byte 7	parity
11111111	00100110	01001001	10010000	00010000	00010011	00011100	11100000	10000000
11100000	11111111	00100110	01001001	10010000	00010000	00010011	00011100	01000000
00011100	11100000	11111111	00100110	01001001	10010000	00010000	00010011	00100000
00010011	00011100	11100000	11111111	00100110	01001001	10010000	00010000	00010000
00010000	00010011	00011100	11100000	11111111	00100110	01001001	10010000	00001000
10010000	00010000	00010011	00011100	11100000	11111111	00100110	01001001	00000100
01001001	10010000	00010000	00010011	00011100	11100000	11111111	00100110	00000010
00100110	01001001	10010000	00010000	00010011	00011100	11100000	11111111	00000001

表 7: Hsiao(72, 64) 码的一种监督矩阵，每行有 27 个 1， $d_{\min} = 4$

如果用三输入异或门，每个 syndrome bit 可以用三级共  $9 + 3 + 1 = 13$  个门电路来生成。在某个实际电路中，8 个监督位一共用 87 个三输入异或门。<sup>②</sup>

如果内存的每个 8-bit 字节单独配置一个奇偶校验 bit，只能检出本字节内的单 bit 错，但无法纠错，空间的额外开销是  $1/8 = 12.5\%$ 。但是如果把 8 个字节的 8 个校验位放到一起，配合适当的纠错码，在空间开销不变的情况下，却能纠正 64-bit 内的单 bit 错并检测双 bit 错 (SEC-DED)。见以下示意图。

<sup>①</sup>[https://en.wikipedia.org/wiki/ECC\\_memory](https://en.wikipedia.org/wiki/ECC_memory)

<sup>②</sup>Mu-Yue Hsiao, Eugene Kolankowsky. "Optimum apparatus and method for check bit generation and error detection, location and correction." US Patent 3,623,155, filed Dec. 24, 1969, and issued Nov. 23, 1971. <https://patents.google.com/patent/US3623155A>

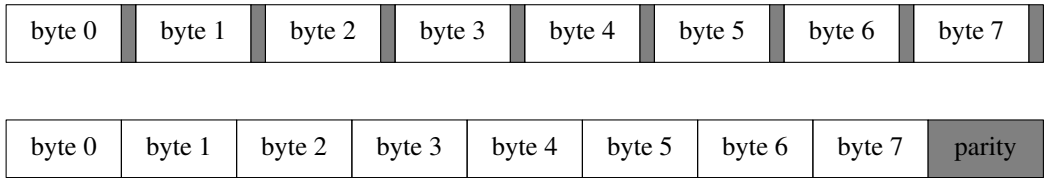


图 17: 简单奇偶校验 bit 与 ECC RAM 的监督位 (灰底) 对比

以前的家用台式机内存条一般有 8 块存储芯片 (或双面共 16 块), 而带 ECC 功能的服务器内存条有 9 块存储芯片 (或双面共 18 块), 这比较直观地体现了监督位所占的比例。目前 ECC 内存的纠错码已不止 [72, 64] 这一种<sup>①</sup>。此外, DDR5 内存芯片要求片内 “on-die ECC”<sup>②</sup>, 某些产品是用  $d_{\min} = 3$  的 Hamming(136, 128) 码实现的 [10]。

### 1.7 Hamming 码的误码率

Hamming [7, 4] 码能纠正 1 bit 错, 那么 block 出错的概率  $p_B = 1 - \sum_{t=0}^1 \binom{n}{t} f^t (1-f)^{n-t} = 1 - (1-f)^7 - 7 \times f(1-f)^6$ 。若  $f = 0.1, p_B = 1 - 0.9^7 - 7 \times 0.1 \times 0.9^6 \approx 0.15$ ,  $p_B$  看起来比  $f$  还高。不过 bit 错误的概率  $p_b \approx \frac{3}{7} p_B \approx 0.064$ , 这里系数  $\frac{3}{7}$  的来历可参考 [2] p.17 对于练习 1.6 的解答。<sup>③</sup>

$f$	$p_B$	$p_b$	$R_3$ 的 $p_b$
0.1	0.150	0.064	0.028
0.01	$2.03 \times 10^{-3}$	$8.70 \times 10^{-4}$	$2.98 \times 10^{-4}$
0.001	$2.09 \times 10^{-5}$	$8.97 \times 10^{-6}$	$2.998 \times 10^{-6}$

表 8: Hamming [7, 4] 码在不同  $f$  下的  $p_B$  与  $p_b$ , 与 3 次重复码  $R_3$  对比

既然简单的重复码就可以实现  $d_{\min} = n$ , 要纠正  $t$  个错只需要  $n \geq 2t + 1$ , 说明纠错能力  $d_{\min}$  并不是唯一的性能指标, 至少还需要考虑码率  $R = k/n$ 。我们怎么综合评价一个纠错码的好与坏呢? 对于线性分组码, 有一个比较简单的指标 nominal coding gain  $\gamma_c$ , 定义是  $\gamma_c = \frac{kd}{n}$ 。例如 Hamming [7,4,3]  $\gamma_c = \frac{kd}{n} = \frac{4 \times 3}{7} \approx 1.714$ , 扩展 Hamming [8, 4, 4]  $\gamma_c = \frac{kd}{n} = \frac{4 \times 4}{8} = 2$ 。在  $\gamma_c$  指标下, 扩展 Hamming 码多用 1 位监督位是很划算的, 它提高了标称编码增益。而重复码  $\gamma_c = \frac{kd}{n} = \frac{1 \times n}{n} = 1$ , 跟没有编码 ( $\gamma_c = \frac{kd}{n} = \frac{1 \times 1}{1} = 1$ ) 无区别, 因此是很糟糕的编码。不过,  $\gamma_c$  并没有反映译码的复杂度, 这在实践中也是需要考虑的。

<sup>①</sup><https://en.wikipedia.org/wiki/Chipkill>

<sup>②</sup>[https://en.wikipedia.org/wiki/DDR5\\_SDRAM#on-die\\_ECC](https://en.wikipedia.org/wiki/DDR5_SDRAM#on-die_ECC)

<sup>③</sup>在没有进一步资料之时, 对于一般的线性分组码通常可粗略估计  $p_b \approx \frac{1}{2} p_B$ 。即如果一个 block 出错, 那么大概有一半的 bit 会翻转, 因为 block 出错并不意味着全部 bits 出错, 否则我们把全部 bits 翻转不就全对了吗。

## 2 Hadamard 码与水手号

前面我们介绍的 Hamming 码、扩展 Hamming 码、Hsiao 码都只能最多纠正单 bit 错误，本节将介绍一种能纠正 7 bit 错误的 Hadamard [32, 6, 16] 线性分组码，这个码曾用于 1969 发射的 Mariner (水手号) 火星探索任务，是最早用于深空通信的信道编码之一。本节所说的 Hadamard 码准确地说是 augmented Hadamard 码<sup>①</sup>，也叫 Walsh 码、Walsh-Hadamard 码，其实是一阶 (first-order) Reed-Muller 码。

Jacques S. Hadamard 是法国数学家，一般译作“雅克·阿达马”。注意姓的第一个字母 h 在法语里不发音，不过有些英语国家的人念作“哈达马”。类似的情况还有讲英语的人把 Hermitian matrix (厄米特矩阵，由法国数学家 Charles Hermite 发明，这里 h 也不发音) 念成“赫密辛”，因为在英文中对应的单词 hermit 中的 h 是发音的，这倒是可以理解。

Reed-Muller 码<sup>②</sup>有两个参数  $r$  和  $m$ ，均为整数，满足  $0 \leq r \leq m$ ，记为  $RM(r, m)$  码，基本参数为  $[n = 2^m, k = \sum_{i=0}^r \binom{m}{i}, d = 2^{m-r}]$ 。 $r$  是阶数， $r = 0$  是简单重复码； $r = 1$  是本节所谈的 Hadamard 码  $RM(1, m)$ ，参数为  $[n = 2^m, k = m + 1, d = 2^{m-1}]$ 。<sup>③</sup> Mariner 1969 用的是  $RM(1, 5)$  码，码字长度 32-bit，信息位 6-bit， $d_{\min} = 16$ ，可以纠正  $t = \lfloor \frac{16-1}{2} \rfloor = 7$  bit 错误，后称 Hardmard [32, 6, 16] 码。

Hadamard 码有不只一种构造方法，我们这里介绍一种直接构造生成矩阵的基向量的办法 [11]。 $RM(1, m)$  码的生成矩阵  $G_m$  是  $(m + 1) \times 2^m$  阶，有  $k = m + 1$  行， $n = 2^m$  列。这个矩阵的行向量  $\mathbf{g}_0, \mathbf{g}_1, \dots, \mathbf{g}_m$  是  $k$  维码字空间的一组基。 $G_m$  矩阵可以递归构造：先把  $G_{m-1}$  矩阵的  $m$  个行向量重复两遍 (例如 0011  $\rightarrow$  00110011)，得到新的  $\mathbf{g}_0, \mathbf{g}_1, \dots, \mathbf{g}_{m-1}$ ；再增加一个新的行向量  $\mathbf{g}_m$ ，其构成是  $2^{m-1}$  个 0 接上  $2^{m-1}$  个 1，这样就得到了  $G_m$  的  $m + 1$  个行向量 (具体代码见第 22 页)。表 9 给出了  $G_1, G_2, \dots, G_5$  的构造，注意由于  $\mathbf{g}_0 = \mathbf{1}$ ，这种方法构造出来的是非系统码。

基	$G_1$	$G_2$	$G_3$	$G_4$	$G_5$
$\mathbf{g}_0$	11	1111	1111 1111	1111 1111 1111 1111	1111 1111 1111 1111 1111 1111 1111 1111
$\mathbf{g}_1$	01	0101	0101 0101	0101 0101 0101 0101	0101 0101 0101 0101 0101 0101 0101 0101
$\mathbf{g}_2$		0011	0011 0011	0011 0011 0011 0011	0011 0011 0011 0011 0011 0011 0011 0011
$\mathbf{g}_3$			0000 1111	0000 1111 0000 1111	0000 1111 0000 1111 0000 1111 0000 1111
$\mathbf{g}_4$				0000 0000 1111 1111	0000 0000 1111 1111 0000 0000 1111 1111
$\mathbf{g}_5$					0000 0000 0000 0000 1111 1111 1111 1111

表 9: Hadamard 码的基向量

有了生成矩阵  $G$ ，Hadamard 码的编码可以简单表示为  $\mathbf{c} = \mathbf{m} \cdot G$  或  $\mathbf{c} = G^T \cdot \mathbf{m}$ 。这种方法需要  $m \cdot 2^m$  bit 的存储空间 (考虑  $\mathbf{g}_0$  是全 1 向量，意味着  $m_0 = 1$  时只需对  $\mathbf{c}$  按位取反，因此不必存储  $\mathbf{g}_0$ )，Hadamard [32, 6, 16] 码需要  $5 \cdot 2^5 = 160$  bit (合 20 字节) 的存储空间。这点空间在现在看来不值一提，但是对于 1960s 设计的 Mariner 1969 飞船而言还是太大了，人们找到了更节省硬件的编码方法。

<sup>①</sup>[https://en.wikipedia.org/wiki/Hadamard\\_code](https://en.wikipedia.org/wiki/Hadamard_code)

<sup>②</sup>1954 年 9 月，David E. Muller 发表了这类线性分组码的构造方法 <https://ieeexplore.ieee.org/document/6499441>；同月，Irving S. Reed 在另一个杂志发表了第一个高效的译码算法 <https://ieeexplore.ieee.org/document/1057465>。Irving Reed 也是 Reed-Solomon 码的发明人，一个人有两个带自己名字的码是很难得的。

<sup>③</sup>原版 Hadamard 码的参数是  $[n = 2^m, k = m, d = 2^{m-1}]$ ，这里我们介绍的是双正交 (biorthogonal) 版， $k$  要多 1 bit。

以 RM(1, 4) 码 (Hadamard [16, 5, 8] 码) 为例,  $m = 4, k = m + 1 = 5$ , 生成矩阵是

$$G_4 = \begin{bmatrix} 1111 & 1111 & 1111 & 1111 \\ 0101 & 0101 & 0101 & 0101 \\ 0011 & 0011 & 0011 & 0011 \\ 0000 & 1111 & 0000 & 1111 \\ 0000 & 0000 & 1111 & 1111 \end{bmatrix}_{5 \times 16}$$

注意观察  $G_4$  的 16 个 5-bit 列向量, 发现除了第 0 bit 固定是 1, 其余 4 bit 刚好是二进制计数值 0000, 0001, 0010, 0011, ..., 1100, 1101, 1110, 1111。因此用一个 4-bit 二进制计数器配合 4 个与门就完成了对消息位高 4 bit 的运算。再把运算结果和  $m_0$  异或, 就得到了码字中的 1 bit, 这个思路如图 18 所示。

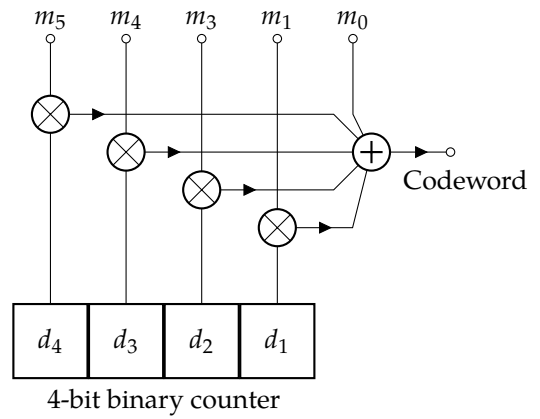


图 18: RM(1, 4) 码的硬件编码方案

这个编码电路可以用几片 TTL 74 数字逻辑芯片搭出来 (图 19)。电路上半部分是用四个 J-K 触发器 (两片 74LS73) 配合两个与门实现了 4-bit 同步计数器, 也可以直接用 74LS93 这种专门的计数电路。图 19 下半部分的 4 个与门可以用一片 74LS08, 五输入异或门可以用一片 74LS86, 这样一共三五块芯片在面包板上就能完成 Hadamard [16, 5, 8] 码 (RM(1, 4) 码) 的编码实验。

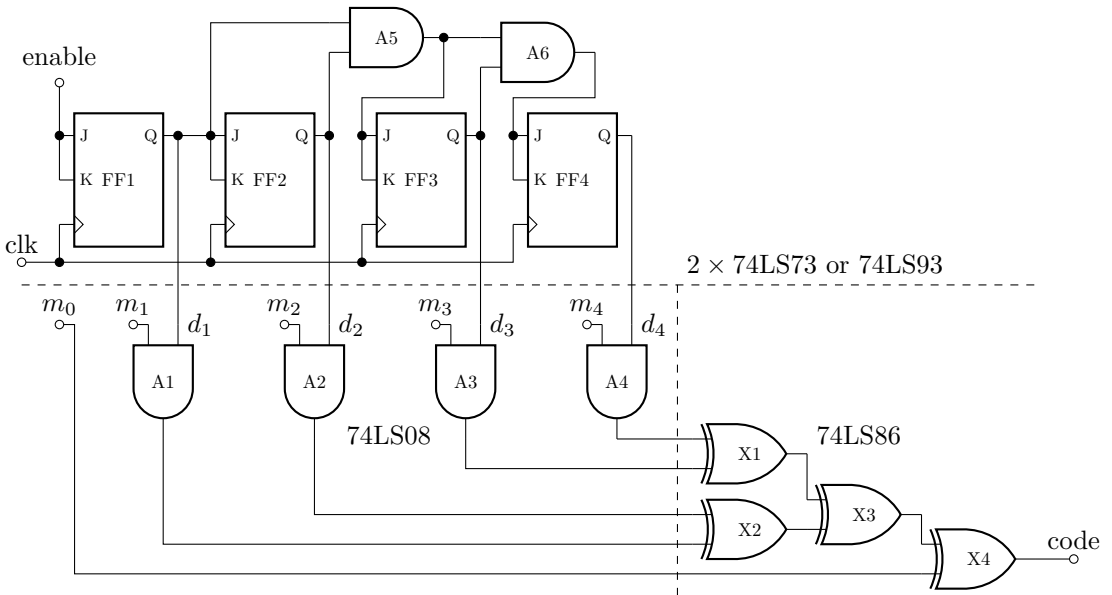


图 19: RM(1, 4) 码的一种编码电路

Mariner 1969 实际搭载的编码电路用的是古早的 DTL 逻辑芯片, 一共用了 24 片 Fairchild 9040 触发器 (R-S 模式和 J-K 模式可配置), 4 片 9046 两输入与非门, 1 片 9047 三输入与非门, 完整电路见文献 [12]。这个电路比图 19 复杂, 是因为它实现的是 comma-free<sup>①</sup> 的 Hadamard [32, 6, 16] 系统码 [13]。

<sup>①</sup>[https://en.wikipedia.org/wiki/Comma-free\\_code](https://en.wikipedia.org/wiki/Comma-free_code)



```

// Codebook for Hadamard[32, 6, 16]
uint32_t codebook[64];

static inline uint32_t popcnt(const uint32_t x)
{
    return __builtin_popcount(x);
}

// Minimum distance decoding
int32_t DecodeMDD(const uint32_t r)
{
    for (int m = 0; m < 64; ++m)
    {
        uint32_t d = popcnt(r ^ codebook[m]);
        if (d < 8)
        {
            return m;
        }
    }
    return -1;
}

```

用 clang 21.1 编译（针对 Intel® Rocket Lake 优化 `-O2 -march=rocketlake`）得到的 x86-64 汇编代码如下（节选），这段代码用到了 AVX-512 指令，其中 `vpopcntd` 是比较新的指令。<sup>①</sup>

```

# 256-byte codebook
codebook:
    .zero    256

.LCPI1_0:
    .long    8

# AVX2 + AVX-512 for Rocket Lake
DecodeMDD(unsigned int):
    vpbroadcastd    %edi, %ymm0
# 对比第 0~7 号码字
    vpxor           codebook(%rip), %ymm0, %ymm1
    vpopcntd        %ymm1, %ymm1
    vpcmpltd        .LCPI1_0(%rip){1to8}, %ymm1, %k0
    movb            $1, %dl
    kortestb        %k0, %k0
    je              .LBB1_2
    xorl            %ecx, %ecx
    jmp             .LBB1_9
# 对比第 8~15 号码字
.LBB1_2:
    vpxor           codebook+32(%rip), %ymm0, %ymm1
    vpopcntd        %ymm1, %ymm1
# 后略

```

我们略微解释一下这几条指令在干嘛，32-bit 寄存器 `edi` 是输入参数 `r`，`ymm0` 和 `ymm1` 都是 256-bit 寄存器（32 字节），可以装下 8 个 32-bit 整数，`ymm0` 保存 8 份相同的 `r`，`ymm1` 保存 `r` 每次与 8 个码字的对比结果，换句话说每条 SIMD 指令可以一次性处理 8 个码字。

<sup>①</sup><https://stackoverflow.com/questions/61880496/population-count-in-avx512>

以下简述每条指令的用途，图 20 配合展示了每条指令执行之后寄存器的内容。

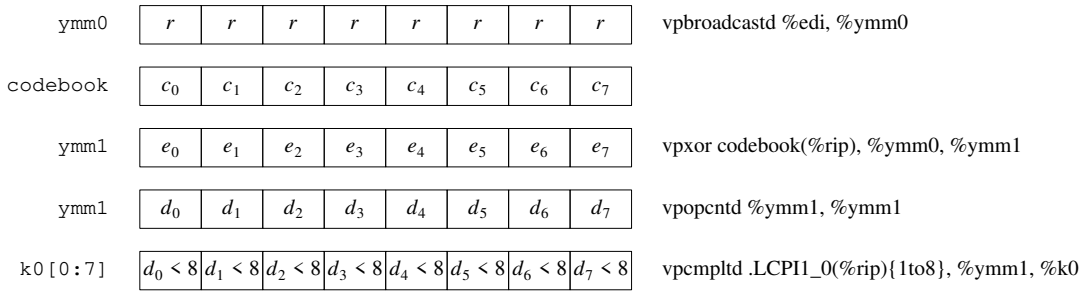


图 20: Hadamard 译码 SIMD 实现的一个步骤

1. `vpbroadcastd` 将 `esi` 中的  $r$  拷贝 8 份到寄存器 `ymm0` 中，`ymm0` 在后续执行中保持不变
2. `vpxor` 将内存 `codebook` 开始的 32 字节数据（前 8 个码字）与 `ymm0` 中的 8 份  $r$  作按位异或运算，结果存入寄存器 `ymm1`，这时 `ymm1` 保存了 8 个 32-bit 错误图样  $e_0, e_1, \dots, e_7$
3. `vpopcntd` 将寄存器 `ymm1` 中的 8 个 32-bit 整数分别数出有多少 bit 为 1，得到 8 个 Hamming 距离  $d_0, d_1, \dots, d_7$ ，结果存回 `ymm1`
4. `vpcmpltd` 将 `ymm1` 中的 8 个 32-bit 整数与 8 比较（`.LCPI1_0` 处存放的是立即数 8），结果存入掩码寄存器 `k0` 的低 8 bits。如果距离小于 8，那么 `k0` 中对应的 bit 是 1，否则是 0
5. `kortestb` 判断 `k0` 寄存器的低 8 bits 是否为全 0。如果是全 0，表明没有哪个距离小于 8，就跳转到对比后续 8 个码字（`je .LBB1_2`），否则说明  $r$  与某个码字的距离小于 8，则跳转到 `.LBB1_9` 找出对应的  $\hat{m}$
6. `.LBB1_2` 处开始处理码字 8 ~ 15，先将内存 `codebook+32` 开始的 32 字节数据（第 8 ~ 15 个码字）与 `ymm0` 作按位异或运算，后续的操作类似。

SIMD 指令可以一次性处理 8 个码字，一共 8 次 `vpxor` + `vpopcntd` 就能完成 64 个码字的对比。

如果用 `-O2 -march=znver4` 参数编译（针对 AMD Zen 4 优化），会生成使用 512-bit 寄存器 `zmm0`、`zmm1` 的代码（节选如下），思路与前面类似，但并行度更高，一次能对比 16 个码字：

```
# AVX-512 for Zen4
DecodeMDD(unsigned int):
    vpbroadcastd    %edi, %zmm0
# 对比第 0~15 号码字
    vpxord         codebook(%rip), %zmm0, %zmm1
    movb          $1, %dl
    vpopcntd      %zmm1, %zmm1
    vpcmpltd      .LCPI1_0(%rip){1to16}, %zmm1, %k0
    kortestw     %k0, %k0
    je            .LBB1_2
    xorl         %ecx, %ecx
    jmp          .LBB1_5
# 对比第 16~31 号码字
.LBB1_2:
    vpxord         codebook+64(%rip), %zmm0, %zmm1
    movl          $16, %ecx
    vpopcntd      %zmm1, %zmm1
```

这种暴力译码没啥技术含量，靠的是编译器优化和现代 CPU 的数据并行处理，在 1950s 年代，人们已经找到了更高效的译码算法。

## Majority logic 译码

Hadamard 码是一阶 Reed-Muller 码，Reed-Muller 码的标准译码方法叫“majority-logic 译码”，其实是一种投票表决机制，少数 (minority) 服从多数 (majority)。Majority-logic 译码最简单的例子是重复码，例如 5 次重复码的码字是 00000 和 11111，如果收到的 5-bit 向量  $r$  中 1 的数量超过半数 ( $\geq 3$ )，就认为  $\hat{m} = 1$ ，否则认为  $\hat{m} = 0$ 。

以 Hadamard[16, 5, 8] 码为例，其生成矩阵  $G_4$  是

$$G_4 = \begin{bmatrix} c_0 & c_1 & c_2 & c_3 & c_4 & c_5 & c_6 & c_7 & c_8 & c_9 & c_{10} & c_{11} & c_{12} & c_{13} & c_{14} & c_{15} \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} -g_0 - \\ -g_1 - \\ -g_2 - \\ -g_3 - \\ -g_4 - \end{bmatrix}$$

编码过程可记为  $c = m \cdot G_4 = [m_0 \ m_1 \ m_2 \ m_3 \ m_4] \cdot G_4 = m_0 g_0 \oplus m_1 g_1 \oplus m_2 g_2 \oplus m_3 g_3 \oplus m_4 g_4$ ，这样得到 16 bit 码字  $c = [c_0 \ c_1 \ c_2 \ \dots \ c_{15}]$ 。由  $G_4$  的列向量可知  $c$  与消息  $m$  的关系：

$$\begin{aligned} c_0 &= m_0 & , & \quad c_4 &= m_0 & \oplus m_3 \\ c_1 &= m_0 \oplus m_1 & , & \quad c_5 &= m_0 \oplus m_1 & \oplus m_3 \\ c_2 &= m_0 & \oplus m_2 & , & c_6 &= m_0 & \oplus m_2 \oplus m_3 \\ c_3 &= m_0 \oplus m_1 \oplus m_2 & , & \quad c_7 &= m_0 \oplus m_1 \oplus m_2 \oplus m_3 \\ &\vdots & & & \vdots & \end{aligned}$$

由异或运算的性质 ( $a \oplus a = 0$ ,  $a \oplus b \oplus a = b$ ) 可知：

$$\begin{aligned} c_0 \oplus c_1 &= m_1 & c_0 \oplus c_2 &= m_2 & c_0 \oplus c_4 &= m_3 & c_0 \oplus c_8 &= m_4 \\ c_2 \oplus c_3 &= m_1 & c_1 \oplus c_3 &= m_2 & c_1 \oplus c_5 &= m_3 & c_1 \oplus c_9 &= m_4 \\ c_4 \oplus c_5 &= m_1 & c_4 \oplus c_6 &= m_2 & c_2 \oplus c_6 &= m_3 & c_2 \oplus c_{10} &= m_4 \\ c_6 \oplus c_7 &= m_1 & c_5 \oplus c_7 &= m_2 & c_3 \oplus c_7 &= m_3 & c_3 \oplus c_{11} &= m_4 \\ &\dots & & \dots & & \dots & & \dots \end{aligned}$$

也就是说  $m_1 \ m_2 \ m_3 \ m_4$  中每个 bit 由码字中的 8 对 bits 提供校验，以上列出了部分情况。

如果传输过程中没有发生错误，那么  $r$  中的 bits 也满足上面的校验关系，如果发生了错误，会破坏一些校验关系，但是只要发生的错误不太多，我们可以用少数服从多数的原则纠正过来。这里举一个具体的例子，对  $m_{18} = [m_0 \ m_1 \ m_2 \ m_3 \ m_4] = [0 \ 1 \ 0 \ 0 \ 1]$  编码，得到  $c = m \cdot G_4 = g_1 \oplus g_4 = 0101 \ 0101 \ 1010 \ 1010$ 。假如传输过程中某些 bit 发生翻转，接收到  $r = \underline{1}101 \ \underline{1}001 \ 1010 \ 1010$ ，那么我们先把  $m_1 \ m_2 \ m_3 \ m_4$  的投票结果列在表 10 中。

	$\hat{m}_1$		$\hat{m}_2$		$\hat{m}_3$		$\hat{m}_4$
$r_0 \oplus r_1$	0	$r_0 \oplus r_2$	1	$r_0 \oplus r_4$	0	$r_0 \oplus r_8$	0
$r_2 \oplus r_3$	1	$r_1 \oplus r_3$	0	$r_1 \oplus r_5$	1	$r_1 \oplus r_9$	1
$r_4 \oplus r_5$	1	$r_4 \oplus r_6$	1	$r_2 \oplus r_6$	0	$r_2 \oplus r_{10}$	1
$r_6 \oplus r_7$	1	$r_5 \oplus r_7$	1	$r_3 \oplus r_7$	0	$r_3 \oplus r_{11}$	1
$r_8 \oplus r_9$	1	$r_8 \oplus r_{10}$	0	$r_8 \oplus r_{12}$	0	$r_4 \oplus r_{12}$	0
$r_{10} \oplus r_{11}$	1	$r_9 \oplus r_{11}$	0	$r_9 \oplus r_{13}$	0	$r_5 \oplus r_{13}$	0
$r_{12} \oplus r_{13}$	1	$r_{12} \oplus r_{14}$	0	$r_{10} \oplus r_{14}$	0	$r_6 \oplus r_{14}$	1
$r_{14} \oplus r_{15}$	1	$r_{13} \oplus r_{15}$	0	$r_{11} \oplus r_{15}$	0	$r_7 \oplus r_{15}$	1
	1		0		0		1

表 10: Hadamard[16, 5, 8] Majority-logic 译码示例

以上结果表明  $[\hat{m}_1 \hat{m}_2 \hat{m}_3 \hat{m}_4] = [1 0 0 1]$ 。接下来找出  $m_0$ ，办法是对  $\hat{\mathbf{m}} = [0 \hat{m}_1 \hat{m}_2 \hat{m}_3 \hat{m}_4]$  编码，得到  $\hat{\mathbf{c}} = \hat{\mathbf{m}} \cdot \mathbf{G}_4 = \mathbf{g}_1 \oplus \mathbf{g}_4 = 0101 0101 1010 1010$ ，再计算  $\hat{\mathbf{c}} \oplus \hat{\mathbf{r}} = 1000 1100 0000 0000$ ，这个向量中 0 的数量大于半数，于是  $\hat{m}_0 = 0$ 。最终结果  $\hat{\mathbf{m}} = [\hat{m}_0 \hat{m}_1 \hat{m}_2 \hat{m}_3 \hat{m}_4] = [0 1 0 0 1] = \mathbf{m}_{18}$ ，译码成功。

对于 Hadamard[32, 6, 16] 码，以下是 majority-logic 译码的 C 语言实现，可以纠正 7 个 bit 或以下的错误。这段代码没有处理译码失败的情况，例如刚好有 8 bits 翻转。

```
// Codebook for Hadamard[32, 6, 16]
uint32_t codebook[64];

// Majority-logic decoding
int DecodeMLD(const uint32_t r)
{
    uint8_t m1 = popcnt((r & 0xaaaaaaaa) ^ ((r & 0x55555555) << 1)) > 8;
    uint8_t m2 = popcnt((r & 0xcccccccc) ^ ((r & 0x33333333) << 2)) > 8;
    uint8_t m3 = popcnt((r & 0xf0f0f0f0) ^ ((r & 0x0f0f0f0f) << 4)) > 8;
    uint8_t m4 = popcnt((r & 0xffff0000) ^ ((r & 0x00ff00ff) << 8)) > 8;
    uint8_t m5 = popcnt((r & 0xffff0000) ^ ((r & 0x0000ffff) << 16)) > 8;

    uint8_t m = (m1 << 1) | (m2 << 2) | (m3 << 3) | (m4 << 4) | (m5 << 5);
    uint8_t m0 = popcnt(codebook[m] ^ r) > 16;

    return m | m0;
}
```

如果用数字电路实现译码，不必去数 bits 再比较大小，而可以直接用 majority gate<sup>①</sup> 求出多数 bit 是 0 还是 1。对于 Hadamard[32, 6, 16] 码，majority-logic 译码步骤如下：先算 5 次 16-bit 异或，用 majority gate 依次求出  $\hat{m}_1 \hat{m}_2 \hat{m}_3 \hat{m}_4 \hat{m}_5$ ，然后对  $\hat{\mathbf{m}}$  做一次编码得到  $\hat{\mathbf{c}}$ ，再对  $\mathbf{r}$  和  $\hat{\mathbf{c}}$  做一次 32-bit 异或和一次 32-bit majority gate，得到  $\hat{m}_0$ ，从而完成译码。以上步骤需要  $5 \times 16 + 32 = 112$  次异或（不计 majority gate 和编码  $\hat{\mathbf{m}}$  的开销<sup>②</sup>），远远小于最小距离译码需要的 2048 次异或。

不过 majority-logic 译码仍然是硬判决译码，与理想的 maximum likelihood decoding 还有相当的差距。在实际的高斯白噪声信道中，要想完全发挥 Hadamard 码的效力，需要用软判决译码。

<sup>①</sup>[https://en.wikipedia.org/wiki/Majority\\_function](https://en.wikipedia.org/wiki/Majority_function)

<sup>②</sup>如果按 Gray code 顺序求和，可以进一步缩小编码开销。<https://math.stackexchange.com/questions/2014142/>

## Hadamard 变换与软判决译码

在 Hamming 距离意义下的最小距离译码其实没有发挥 Hadamard 的全部纠错能力，因为实际的信道往往不是二元对称信道 (BSC)，实际的信道噪声也不是简单地按一定的概率随机地翻转一些 bits。这里以 Hadamard [8, 4, 4] 码举例，这个码  $d_{\min} = 4$ ，能纠正单 bit 错或检测双 bit 错，表 11 列出了完整的码书。

$m$	二进制	$c$	$m$	二进制	$c$
0	0000	00000000	8	1000	00001111
1	0001	11111111	9	1001	11110000
2	0010	01010101	10	1010	01011010
3	0011	10101010	11	1011	10100101
4	0100	00110011	12	1100	00111100
5	0101	11001100	13	1101	11000011
6	0110	01100110	14	1110	01101001
7	0111	10011001	15	1111	10010110

表 11: Hadamard [8, 4, 4] 码的码书

假如我们要发送消息  $m_{11}$ ，码字  $c = 10100101$ 。实际发送信号时，将码字  $c$  中的  $\{0, 1\}$  符号映射为  $\{+1, -1\}$  电平（或相位等模拟量），发送的信号是  $v = (-1, +1, -1, +1, +1, -1, +1, -1)$ 。信号  $v$  在传输过程中叠加了噪声<sup>①</sup>  $n = (-0.066, -1.038, 0.731, 0.043, -1.082, 0.142, -1.306, 0.621)$ ，注意有几个点的噪声幅度大过了信号本身  $|n_i| > 1$ 。

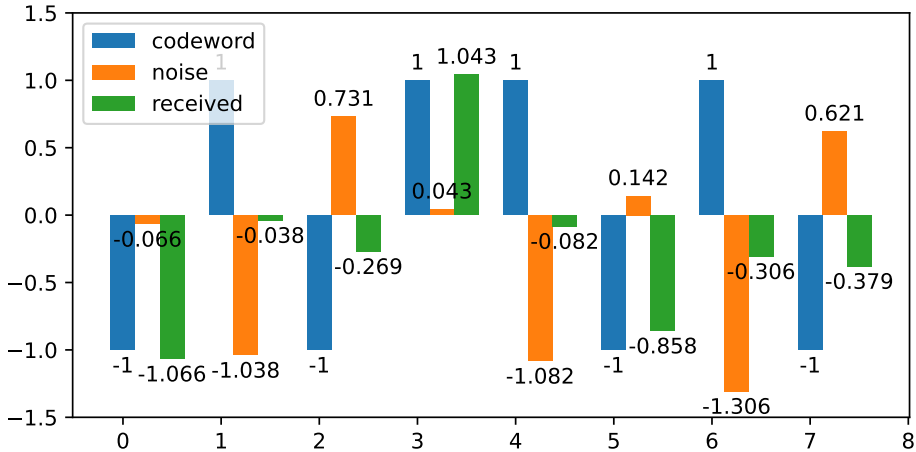


图 21: 叠加噪声的效果

收到的信号  $x = (-1.066, -0.038, -0.269, 1.043, -0.082, -0.858, -0.306, -0.379)$ ，即  $x_i = v_i + n_i$ 。这几个量的大小关系见图 21。如果以 0 为判决门限， $x_i \geq 0$  为 0， $x_i < 0$  为 1，则  $r = \underline{11101111}$ ，有 3 bits 发生了翻转，超过了 Hadamard [8, 4, 4] 码的纠错能力。如果用求最小 Hamming 距离的译码方法， $\hat{c} = 11111111$ ，推断  $\hat{m} = 1$ ，则译码结果错误。

<sup>①</sup>这个噪声样本是用 `np.random.normal(0, 1/2, 8)` 得到的。

注意到我们发送码字  $\mathbf{c}$  用了 8 份能量，而噪声  $\mathbf{n}$  只用了略大于 3 份能量就扰乱了译码结果。对于硬判决来说，翻转单个 bit 可谓“事半功倍”（如果信号是 0 和 1，判决门限是 0.5，那么只要把 0 扰乱到 0.51 就能更改判决结果为 1），这“颠倒黑白”的成本也太低了。

Hadamard 码的优点是可以利用 Hadamard 变换来实现软判决译码 (soft-decision decoding)。所谓 Hadamard 变换<sup>①</sup>，就是用 Hadamard 矩阵乘以输入向量  $\mathbf{x}$  得到向量  $\mathbf{y}$ ，注意这个运算是在实数域  $\mathbb{R}$  进行的，而不是  $\mathbb{F}_2$ 。Hadamard 矩阵  $H$  是  $n$  阶方阵，元素是 +1 或 -1，其  $n$  个行向量是彼此正交的，即  $HH^T = nI_n$ ，其中  $I_n$  是  $n$  阶单位阵。我们用到的 Hadamard 矩阵  $H_m$  是  $2^m$  阶方阵，可以递归构造：<sup>②</sup>（为了强调其运算遵循通常的实数运算法则，后用圆括号表示这个矩阵）

$$H_1 = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}, H_m = \begin{pmatrix} H_{m-1} & H_{m-1} \\ H_{m-1} & -H_{m-1} \end{pmatrix}, \text{ 以此方法递推构造 } H_2:$$

$$H_2 = \begin{pmatrix} H_1 & H_1 \\ H_1 & -H_1 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{pmatrix}, \text{ 类似的}$$

$$H_3 = \begin{pmatrix} H_2 & H_2 \\ H_2 & -H_2 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 & 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 & 1 & -1 & -1 & 1 \\ 1 & 1 & 1 & 1 & -1 & -1 & -1 & -1 \\ 1 & -1 & 1 & -1 & -1 & 1 & -1 & 1 \\ 1 & 1 & -1 & -1 & -1 & -1 & 1 & 1 \\ 1 & -1 & -1 & 1 & -1 & 1 & 1 & -1 \end{pmatrix}, H_4 \text{ 与 } H_5 \text{ 以此类推。}$$

Hadamard [8, 4, 4] 译码第 1 步，对收到的信号向量  $\mathbf{x}$  做 Hadamard 变换，得到实数向量  $\mathbf{y} = H_3\mathbf{x}$ 。

第 2 步，找出  $\mathbf{y}$  中绝对值最大的元素的下标，即使得  $|y_i|$  最大的下标  $z = \arg \max_i |y_i|$ 。

第 3 步，译码结果  $\hat{\mathbf{m}} = 2z + [y_z < 0]$ 。其中  $[y_z < 0]$  是一个布尔表达式<sup>③</sup>， $y_z < 0$  时取 1，否则取 0。例如，如果  $z = 1$ ， $y_1 \geq 0$  则  $\hat{\mathbf{m}} = 2$ ；如果  $z = 2$ ， $y_2 < 0$  则  $\hat{\mathbf{m}} = 5$ 。

我们先验证一下以上译码方法在无噪声的情况下是正确的。若要传递  $\mathbf{m} = 0$ ， $\mathbf{c} = 00000000$ ，发送  $\mathbf{v} = (+1, +1, +1, +1, +1, +1, +1, +1)$ ，如果没有噪声，收到  $\mathbf{x} = \mathbf{v}$ 。计算  $\mathbf{y} = H_3\mathbf{x} = (8, 0, 0, 0, 0, 0, 0, 0)$ ， $z = \arg \max_i |y_i| = 0$ ， $y_0 \geq 0$ ，所以  $\hat{\mathbf{m}} = 0$ 。要传递  $\mathbf{m} = 3$ ， $\mathbf{c} = 10101010$ ，发送  $\mathbf{v} = (-1, +1, -1, +1, -1, +1, -1, +1)$ ，收到  $\mathbf{x} = \mathbf{v}$ 。计算  $\mathbf{y} = H_3\mathbf{x} = (0, -8, 0, 0, 0, 0, 0, 0)$ ， $z = \arg \max_i |y_i| = 1$ ， $y_1 < 0$ ，所以  $\hat{\mathbf{m}} = 2z + 1 = 3$ 。

<sup>①</sup>[https://en.wikipedia.org/wiki/Hadamard\\_transform](https://en.wikipedia.org/wiki/Hadamard_transform)，可视为将输入向量  $\mathbf{x}$  分解为 Walsh 函数的叠加。

<sup>②</sup>[https://en.wikipedia.org/wiki/Hadamard\\_matrix#Sylvester's\\_construction](https://en.wikipedia.org/wiki/Hadamard_matrix#Sylvester's_construction) 信道编码常用的 Hadamard 矩阵是 James J. Sylvester 在 1867 年构造的，早于 Hadamard 本人在 1893 年发表成果。

<sup>③</sup>[https://en.wikipedia.org/wiki/Iverson\\_bracket](https://en.wikipedia.org/wiki/Iverson_bracket)

回到前面有噪声的例子，收到  $\mathbf{x} = (-1.066, -0.038, -0.269, 1.043, -0.082, -0.858, -0.306, -0.379)$ ，计算  $\mathbf{y} = H_3 \mathbf{x} = (-1.955, -1.491, -2.133, 0.987, 1.295, \underline{-3.189}, -1.623, -0.419)$ ， $z = \arg \max_i |y_i| = 5$ ， $y_5 < 0$ ，所以  $\hat{m} = 2z + 1 = 11$ ，译码正确。这个例子体现了软判决的纠错能力更强，因为硬判决损失了信号强度的部分信息。

我们现在的很多教材在讲纠错码的时候用的还是硬判决的思路，先对输入信号进行解调 (demodulation)，得到一串可能包含错误的比特流，再进行错误纠正 (error correction)，这种讲法不免有些过时了。早在 1960s 年代初，人们 (Robert M. Fano) 就已经认识到软判决的必要性，1960s 年代未发射的 Pioneer 和 Mariner 飞船的数字通信系统用的也是软判决译码。[20]

以下是用 NumPy/SciPy 实现的 Hadamard 变换译码算法：

```
def HadamardDecode(H, x):
    y = H @ x          # Hadamard 变换
    z = np.argmax(abs(y)) # 找出 y 中绝对值最大元素的下标 z
    m = int(z << 1)   # z 左移 1 位
    if y[z] < 0:
        m |= 1        # 必要时，最低位置 1
    return m
```

验证这段代码：

```
c = [1, 0, 1, 0, 0, 1, 0, 1] # 码字
v = (-1) ** np.array(c)      # 编码 {0, 1} -> {+1, -1}
n = np.random.normal(0, 0.5, 8) # 产生一些噪音
x = v + n                    # 叠加噪音
H3 = scipy.linalg.hadamard(2**3) # 生成 Hadamard 矩阵
print(HadamardDecode(H3, x))    # 译码，成功的话应得到 11
```

### 快速 Walsh-Hadamard 变换

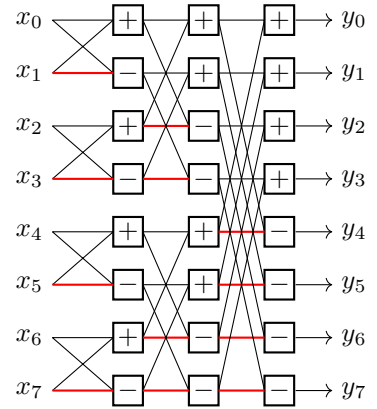
Hadamard 矩阵的元素是  $\pm 1$ ，所以 Hadamard 变换只用计算加减法。以下是 Hadamard [8, 4, 4] 译码的第一步计算：

$$\mathbf{y} = H_3 \mathbf{x} = \begin{pmatrix} x_0 + x_1 + x_2 + x_3 + x_4 + x_5 + x_6 + x_7 \\ x_0 - x_1 + x_2 - x_3 + x_4 - x_5 + x_6 - x_7 \\ x_0 + x_1 - x_2 - x_3 + x_4 + x_5 - x_6 - x_7 \\ x_0 - x_1 - x_2 + x_3 + x_4 - x_5 - x_6 + x_7 \\ x_0 + x_1 + x_2 + x_3 - x_4 - x_5 - x_6 - x_7 \\ x_0 - x_1 + x_2 - x_3 - x_4 + x_5 - x_6 + x_7 \\ x_0 + x_1 - x_2 - x_3 - x_4 - x_5 + x_6 + x_7 \\ x_0 - x_1 - x_2 + x_3 - x_4 + x_5 + x_6 - x_7 \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ y_6 \\ y_7 \end{pmatrix}$$

直接计算  $n$  阶 Hadamard 矩阵与向量的乘法需要  $O(n^2)$  次加减法，快速 Walsh-Hadamard 变换 (FWHT) 可以把运算量降到  $O(n \log n)$ 。表 22a 是用 8 点 FWHT 计算  $\mathbf{y} = H_3 \mathbf{x}$  的步骤，中间结果存入向量  $\mathbf{a}$  和  $\mathbf{b}$ ，一共  $8 \log 8 = 24$  次加减法。图 22b 是运算的流图，与快速 Fourier 变换 (FFT) 几乎一模一样。

$x$	$a$	$b$	$y$
$x_0$	$a_0 = x_0 + x_1$	$b_0 = a_0 + a_1$	$y_0 = b_0 + b_1$
$x_1$	$a_1 = x_2 + x_3$	$b_1 = a_2 + a_3$	$y_1 = b_2 + b_3$
$x_2$	$a_2 = x_4 + x_5$	$b_2 = a_4 + a_5$	$y_2 = b_4 + b_5$
$x_3$	$a_3 = x_6 + x_7$	$b_3 = a_6 + a_7$	$y_3 = b_6 + b_7$
$x_4$	$a_4 = x_0 - x_1$	$b_4 = a_0 - a_1$	$y_4 = b_0 - b_1$
$x_5$	$a_5 = x_2 - x_3$	$b_5 = a_2 - a_3$	$y_5 = b_2 - b_3$
$x_6$	$a_6 = x_4 - x_5$	$b_6 = a_4 - a_5$	$y_6 = b_4 - b_5$
$x_7$	$a_7 = x_6 - x_7$	$b_7 = a_6 - a_7$	$y_7 = b_6 - b_7$

(a) FWHT 计算  $y = H_3x$  的三轮步骤<sup>①</sup>



(b) 蝶形运算图，红线表示相减

图 22: 8 点快速 Walsh-Hadamard 变换

这里略微解释一下为什么 FWHT 和 Hadamard 变换是等价的，一种简单的证明思路是矩阵分解。用  $S$  矩阵表示表 22a 中的一列运算，容易验证  $H_3 = S^3$ ，那么  $y = H_3x = S^3x = S(S(Sx))$ 。由于  $S$  矩阵很稀疏，计算  $Sx$  只需要  $n$  次加减法，一共需要计算  $\log n$  轮，合起来是  $n \log n$  次加减法。

$$S = \begin{pmatrix} a_0 & a_1 & a_2 & a_3 & a_4 & a_5 & a_6 & a_7 \\ 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 1 & -1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & -1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 \end{pmatrix} \begin{matrix} b_0 = a_0 + a_1 \\ b_1 = a_2 + a_3 \\ b_2 = a_4 + a_5 \\ b_3 = a_6 + a_7 \\ b_4 = a_0 - a_1 \\ b_5 = a_2 - a_3 \\ b_6 = a_4 - a_5 \\ b_7 = a_6 - a_7 \end{matrix}$$

以下是 C 语言实现的 8 点 FWHT 计算，输入数据是 8-bit 有符号整数数组（实际上通常 3-bit 量化就足够发挥软判决的潜力），中间结果和输出是 16-bit 有符号整数数组。

```
// 8-point fast Walsh-Hadamard transform
void FWHT8(int8_t x[8], int16_t y[8])
{
    int16_t a[8], b[8];

    for (int i=0; i < 4; ++i) {
        a[i] = x[i*2] + x[i*2+1];
        a[i+4] = x[i*2] - x[i*2+1];
    }

    for (int i=0; i < 4; ++i) {
        b[i] = a[i*2] + a[i*2+1];
        b[i+4] = a[i*2] - a[i*2+1];
    }
}
```

<sup>①</sup>我为 Apache commons-math 的 FWHT 提交了一个文档补丁：<https://github.com/apache/commons-math/pull/295>

```

for (int i=0; i < 4; ++i) {
    y[i] = b[i*2] + b[i*2+1];
    y[i+4] = b[i*2] - b[i*2+1];
}
}

```

不出所料，现代编译器会用 SIMD 指令加速这段代码，以下是 clang 21.1 用 `-O2 -march=rocketlake` 参数编译得到的 AVX2 汇编代码（注释来自编译器本身），`xmm0`、`xmm1` 等寄存器是 128-bit，刚好可以装下 8 个 16-bit 整数。

```

FWHT8:
    vpmovsxbw    (%rdi), %xmm0
    vprold      $16, %xmm0, %xmm1
    vpaddw     %xmm0, %xmm1, %xmm2
    vpsubw     %xmm0, %xmm1, %xmm0
    vpblendw   $170, %xmm0, %xmm2, %xmm0 # xmm0 = xmm2[0],xmm2[1],xmm2[2],xmm2[3],
                                           #          xmm2[4],xmm2[5],xmm2[6],xmm2[7]

    vpsshufd   $177, %xmm0, %xmm1      # xmm1 = xmm0[1,0,3,2]
    vpaddw     %xmm1, %xmm0, %xmm2
    vpsubw     %xmm1, %xmm0, %xmm0
    vpblendw   $10, %xmm2, %xmm0, %xmm0 # xmm0 = xmm0[0],xmm2[1],xmm0[2],xmm2[3]

    vpsshufd   $78, %xmm0, %xmm1      # xmm1 = xmm0[2,3,0,1]
    vpaddw     %xmm1, %xmm0, %xmm2
    vpsubw     %xmm1, %xmm0, %xmm0
    vshufps    $27, %xmm0, %xmm2, %xmm0 # xmm0 = xmm2[3,2],xmm0[1,0]
    vmovups    %xmm0, (%rsi)
    retq

```

以下略微解释一下 SIMD 指令的运算逻辑，C 代码中的三个循环被向量化 (vectorized) 成三轮。第一轮：

1. `vpmovsxbw` 把输入数组  $x$  存入 `xmm0` 寄存器
2. `vprold` 将 `xmm0` 内的 4 个 32-bit 整数各循环左移 16 bit，结果存入 `xmm1`，这条是 AVX-512 指令
3. `vpaddw` 将 `xmm1` 和 `xmm0` 按 16-bit 整数数组相加，结果存入 `xmm2`，即  $xmm2 = xmm1 + xmm0$
4. `vpsubw` 将 `xmm1` 和 `xmm0` 按 16-bit 整数数组相减，结果存回 `xmm0`，即  $xmm0 = xmm1 - xmm0$
5. `vpblendw` 抽取 `xmm2` 和 `xmm0` 中的 16-bit 整数放入 `xmm0`，立即数 170 的二进制是 10101010，其中 1 的 bit 表示取 `xmm2` 对应位置的 16-bit 数，0 表示取 `xmm0` 对应位置的数。

下图配合展示了第一轮每条指令执行之后寄存器的内容。

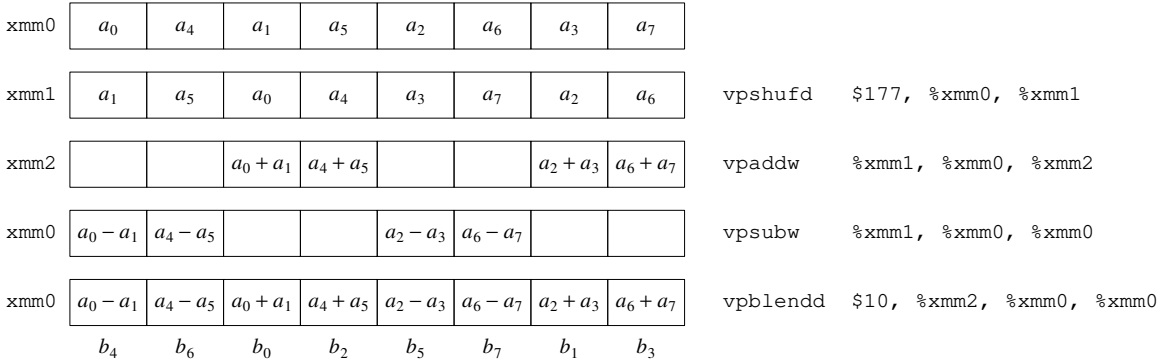
xmm0	$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	<code>vpmovsxbw (%rdi), %xmm0</code>
xmm1	$x_1$	$x_0$	$x_3$	$x_2$	$x_5$	$x_4$	$x_7$	$x_6$	<code>vprold \$16, %xmm0, %xmm1</code>
xmm2	$x_0 + x_1$		$x_2 + x_3$		$x_4 + x_5$		$x_6 + x_7$		<code>vpaddw %xmm0, %xmm1, %xmm2</code>
xmm0		$x_0 - x_1$		$x_2 - x_3$		$x_4 - x_5$		$x_6 - x_7$	<code>vpsubw %xmm0, %xmm1, %xmm0</code>
xmm0	$x_0 + x_1$	$x_0 - x_1$	$x_2 + x_3$	$x_2 - x_3$	$x_4 + x_5$	$x_4 - x_5$	$x_6 + x_7$	$x_6 - x_7$	<code>vpblendw \$170, %xmm0, %xmm2, %xmm0</code>
	$a_0$	$a_4$	$a_1$	$a_5$	$a_2$	$a_6$	$a_3$	$a_7$	

第一轮完成之后，`xmm0` 存放的是中间数组  $a$ ，注意下标不是顺序递增的。

第二轮：

1. `vpsufd` 将 `xmm0` 中的 4 个 32-bit 数按指定顺序存入 `xmm1`，立即数 177 的二进制是 10110001，表示了这 4 个 32-bit 数的排列顺序是 1032，即 `xmm1[0,1,2,3] = xmm0[1,0,3,2]`
2. `vpaddw` 将 `xmm0` 和 `xmm1` 按 16-bit 整数数组相加，结果存入 `xmm2`，即 `xmm2 = xmm0 + xmm1`
3. `vpsubw` 将 `xmm0` 和 `xmm1` 按 16-bit 整数数组相减，结果存回 `xmm0`，即 `xmm0 = xmm0 - xmm1`
4. `vpblendd` 抽取 `xmm0` 和 `xmm2` 中的 32-bit 整数放入 `xmm0`，立即数 10 的二进制是 1010，其中 1 的 bit 表示取 `xmm0` 对应位置的 32-bit 数，0 表示取 `xmm2` 对应位置的数。

下图配合展示了第二轮每条指令执行之后寄存器的内容。

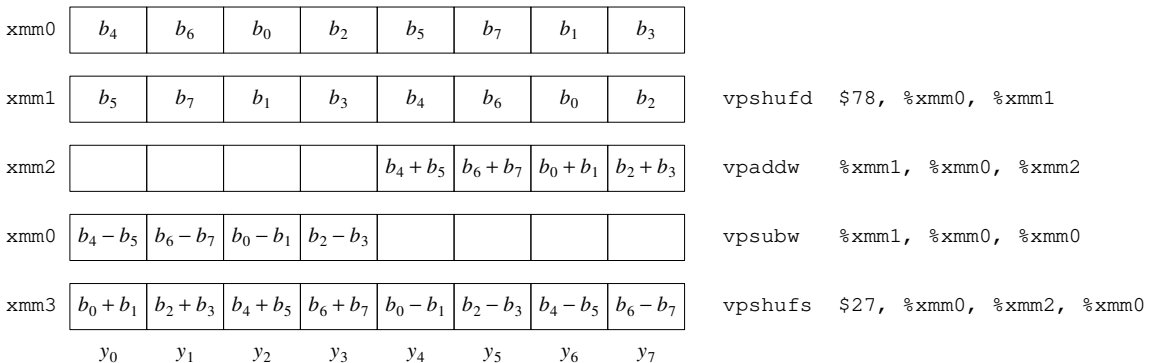


第二轮完成之后，`xmm0` 存放的是中间数组 `b`，注意下标不是顺序递增的。

第三轮：

1. `vpsufd` 将 `xmm0` 中的 4 个 32-bit 数按指定顺序存入 `xmm1`，立即数 78 的二进制是 01001110，表示了这 4 个 32-bit 数的排列顺序是 2301，即把 `xmm0` 的高 64-bit 存入 `xmm1` 的低 64-bit，且把 `xmm0` 的低 64-bit 存入 `xmm1` 的高 64-bit
2. `vpaddw` 将 `xmm0` 和 `xmm1` 按 16-bit 整数数组相加，结果存入 `xmm2`，即 `xmm2 = xmm0 + xmm1`
3. `vpsubw` 将 `xmm0` 和 `xmm1` 按 16-bit 整数数组相减，结果存回 `xmm0`，即 `xmm0 = xmm0 - xmm1`
4. `vsufps` 抽取 `xmm2` 和 `xmm0` 中的 32-bit 数放入 `xmm0`，立即数 27 的二进制是 00011011，表示抽取顺序是 3210，即 `xmm0[0] = xmm2[3]`，`xmm0[1] = xmm2[2]`，`xmm0[2] = xmm0[1]`，`xmm0[3] = xmm0[0]`

下图配合展示了第三轮每条指令执行之后寄存器的内容。



第三轮完成之后，`xmm0` 存放的是结果数组 `y`，注意下标刚好是顺序递增的。

以上快速算法很容易推广到更高阶的 Hadamard 变换，例如 32 点 FWHT 需要  $32 \log 32 = 32 \times 5 = 160$  次加减法，显著少于直接矩阵乘以向量所需的  $32 \times 32 = 1024$  次加减法。1969 年发射的 Mariner 火星探索飞船用的是 Hadamard [32, 6, 16] 码，地面的译码设备叫 Green 机 [14]，是快速 Walsh-Hadamard 变换的硬件实现。这台设备不是以颜色命名，而是以设计者 Richard Green 的姓氏命名。

### Hadamard 码的一些注脚

Hadamard 码是成功应用软判决译码的最早例子之一。硬判决译码其实是在 Hamming 空间找最近的邻居，而软判决译码的一种理解是在欧几里德空间（实内积空间）找最近邻居。Hamming 空间是一个向量空间，但不是内积空间（内积  $\langle \mathbf{x}, \mathbf{y} \rangle = \sum_i x_i y_i$ ），因为内积要求  $\mathbf{x} \neq \mathbf{0}$  时  $\langle \mathbf{x}, \mathbf{x} \rangle > 0$ 。对  $\mathbb{F}_2$  上的 Hamming 空间来说， $\langle \mathbf{x}, \mathbf{x} \rangle$  只有 0 和 1 两种取值，当  $\mathbf{x}$  有偶数个 1 时  $\langle \mathbf{x}, \mathbf{x} \rangle = 0$ ，这就不满足内积空间的要求。换言之，Hamming 空间里虽然可以定义两个向量的距离  $d_H(\mathbf{x}, \mathbf{y})$ ，但没法定义向量的长度  $|\mathbf{x}|$ ，也叫模长 (norm)。同理，也没法定义两个向量是否正交 (orthogonal)  $\iff \langle \mathbf{x}, \mathbf{y} \rangle = 0$ 。<sup>①</sup>

Hadamard 矩阵  $H_m$  是  $2^m$  阶方阵<sup>②</sup>，其元素是  $\pm 1$ ，因此各个行向量的模 (norm) 是相同的。 $H_m$  是对称阵  $H_m^T = H_m$ 。 $H_m$  的第 1 行和第 1 列的元素全是 +1，其余各行各列的 +1 和 -1 数量相等。 $H_m$  的  $2^m$  个行向量是彼此正交的，以  $H_m$  的行向量为码字的原始 Hadamard 码是正交码。而我们这里讲的 augmented Hadamard 码（一阶 Reed-Muller 码）是双正交码<sup>③</sup>，“双正交/ biorthogonal”这个词比较少见，意思是一组正交向量连同它们的“反向量”。例如 Hadamard [32, 6, 16] 码有  $2^6 = 64$  个码字，其中 32 个码字是 Hadamard 矩阵  $H_5$  的行向量，另外 32 个码字是矩阵  $-H_5$  的行向量。

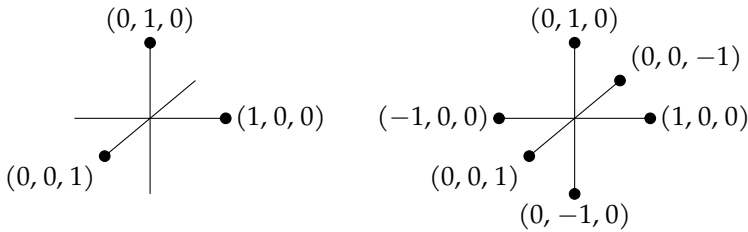


图 23:  $M = 3$  的 orthogonal 信号与  $M = 6$  的 biorthogonal 信号举例

## 2.2 Mariner 水手号行星探索计划

在 1960s 年代，NASA 开展了 Mariner 计划<sup>④</sup>，探索金星和火星，表 12 列出了这 10 次发射任务。Mariner 计划有很多载入史册的事迹，对普通人来说最直观的感受恐怕是它拍摄的火星表面照片。1962 年 8 月发射的 Mariner 2<sup>⑤</sup> 飞船是人类第一个成功接近其他行星并传回数据的空间探测器，可惜飞掠金星的 Mariner 2 没有安装拍照设备，因为当时人们认为金星的大气太厚，拍不到有价值的信息。

<sup>①</sup>参考 G. David Forney 在 MIT 6.451 课程中的讲解 <https://youtu.be/q4LsDyLkZcI?t=600>

<sup>②</sup>一般而言 Hadamard 矩阵的阶数不一定是  $2^m$ ，存在其他构造方法，见 <http://www.neilsloane.com/hadamard/>。

<sup>③</sup><https://errorcorrectionzoo.org/c/biorthogonal>

<sup>④</sup>[https://en.wikipedia.org/wiki/Mariner\\_program](https://en.wikipedia.org/wiki/Mariner_program)

<sup>⑤</sup>[https://en.wikipedia.org/wiki/Mariner\\_2](https://en.wikipedia.org/wiki/Mariner_2) 和 <https://www.space.com/18746-mariner-2.html>

太空飞船	发射日期	探测目标	成功与否	纠错码	灰度图像格式
Mariner 1	1962/07/22	飞掠金星	发射失败	无	200 × 200 × 6-bit
Mariner 2	1962/08/27	飞掠金星	成功	无	
Mariner 3	1964/11/05	飞掠火星	发射失败	无	
Mariner 4	1964/11/28	飞掠火星	成功	无	
Mariner 5	1967/12/04	飞掠金星	成功	无	
Mariner 6	1969/02/25	飞掠火星	成功	Hadamard 码	945 × 704 × 8-bit
Mariner 7	1969/03/27	飞掠火星	成功	Hadamard 码	945 × 704 × 8-bit
Mariner 8	1971/05/09	环绕火星	发射失败	Hadamard 码	832 × 700 × 9-bit
Mariner 9	1971/05/30	环绕火星	成功		
Mariner 10	1973/11/03	飞掠金星和水星	成功		

表 12: Mariner 行星探索计划, 从 1969 年的 Mariner 6/7 开始使用 Hadamard 纠错码。

苏联 1961 年发射的 Venera 1<sup>①</sup> 更早飞掠金星, 但其通信系统在发射之后不久就失灵了, 没能传回金星的观测数据。

Mariner 4<sup>②</sup> (水手 4 号) 在距今 60 年前 (1965 年 7 月) 飞掠 (flyby) 火星, 这是人类第一次成功飞掠火星; 它拍摄了第一幅近距离观察火星表面的照片, 这是人类第一次在外太空对其他行星拍照并传回地球。Mariner 4 一共拍摄了并传回了 22 张照片, 其中前 16 张有比较丰富的细节<sup>③</sup>。它拍摄的黑白照片分辨率是 200 × 200 像素, 每个像素是 6 bit 灰度, 每张照片数据量是 240 000 bits。Mariner 4 向地球通信的码率是 8 1/3 bps, 因此每张照片至少需要 8 小时才能发送完<sup>④</sup>。为了冗余, 照片数据传送了两遍, 一共花了约 20 天才传送完毕 (1965 年 7 月 15 日至 8 月 3 日)。Mariner 4 用不到 10W 的发射功率, 从距离地球 2 亿多公里的火星轨道上总共传回了超过 250 Mbits 数据, 这是通信技术史上的里程碑 [16, 17]。

[T]he spacecraft flew by Mars on July 14 and 15, 1965. Planetary science mode was turned on at 15:41:49 UT on 14 July. The camera sequence started at 00:18:36 UT on July 15 (7:18:49 p.m. EST on July 14) and 21 pictures plus 21 lines of a 22nd picture were taken. The images covered a discontinuous swath of Mars starting near 40 N, 170 E, down to about 35 S, 200 E, and then across to the terminator at 50 S, 255 E, representing about 1% of the planet's surface. The closest approach was 9,846 km from the Martian surface at 01:00:57 UT 15 July 1965 (8:00:57 p.m. EST 14 July). At the time of closest approach the spacecraft was 216 million km from Earth moving at a speed of approximately 7 km/sec relative to Mars (1.7 km/sec relative to Earth). The images taken during the flyby were stored in the onboard tape recorder. At 02:19:11 UT Mariner 4 passed behind Mars as seen from Earth and the radio signal ceased. The signal was reacquired at 03:13:04 UT when the spacecraft reappeared. Cruise mode was then re-established. Transmission of the taped images to Earth began about 8.5 hours after signal reacquisition and continued until 3 August. All images were transmitted twice to insure no data were missing or corrupt. <https://nssdc.gsfc.nasa.gov/nmc/spacecraft/display.action?id=1964-077A>

<sup>①</sup>[https://en.wikipedia.org/wiki/Venera\\_1](https://en.wikipedia.org/wiki/Venera_1)

<sup>②</sup>[https://en.wikipedia.org/wiki/Mariner\\_4](https://en.wikipedia.org/wiki/Mariner_4) 和 <https://www.directedplay.com/first-tv-image-of-mars>

<sup>③</sup>Mariner 4 pictures of Mars. <https://ntrs.nasa.gov/citations/19680006637>

<sup>④</sup>Mariner 4 通信的字长是 7 bit, 我估计每个 6-bit 像素会放到一个 7-bit word 里发送, 因此实际时间更长。

1969 年发射的 Mariner 6/7<sup>①</sup> 两艘飞船一共拍摄了 200 多张火星照片，这些黑白照片的分辨率是  $945 \times 704$ ，每个像素 8-bit。<sup>②</sup> Mariner 1969 每张照片的数据量约 5.3 Mbits，是 Mariner 4 的 22 倍多，好在 Mariner 1969 的最高通讯速率提高了近 2000 倍 (32.9dB)，达到 16.2 kbps<sup>[18]</sup>，发送一张照片只需要 5 分多钟时间。

为什么不先压缩再传送呢？我认为可能有软硬件两方面的原因。软件方面当时尚无合适的压缩算法，当时无损压缩只有 Huffman 编码和 run length encoding，而通用的 LZ77/LZ78 压缩算法还要等十多年才会发明。有损图像压缩的基础离散余弦变换 (DCT) 在 1972 年才提出，后来广为应用的 JPEG 压缩方案在 1992 年才完成标准化。硬件方面，1960 年集成电路才刚刚面世<sup>③</sup>，一度非常流行的 7400 系列 TTL 数字逻辑 IC 晚到 1967 年才上市，Mariner 1969 飞船上搭载的 Hadamard 编码电路尚且用的是 1964 年上市的 Fairchild DTL 芯片<sup>④</sup>。1971 年 Intel 才发布 4-bit 4004 微控制器。这都表明图像压缩技术在 Mariner 1969 上不现实。就算能进行图像压缩，压缩后的数据传输对误码率也有更高的要求，不一定划算。

TABLE I. COMPARISON OF TELEMETRY SYSTEMS

PARAMETER	MARINER IV	MARINER '69 HIGH-RATE	IMPROVEMENT
Transmitter Power	(8.9W) + 39.5 dbm W	(18.2W) + 42.60 dbm W	*+3.10 db
Modulation Loss	-5.3 db	-1.34 db	*+3.96 db
Transmitting Antenna Gain	+20.1 db	+20.21 db	*+.11 db
Space Loss	( $216 \times 10^6$ km) -266.2 db	( $97 \times 10^6$ km) -259.36 db	*+6.84 db
Receiving Antenna Gain	(85 ft.) +52.5 db	(210 ft.) +61.00 db	*+8.50 db
Receiver Loss	-3.4 db	-.44 db	*+2.96 db
Received Sideband Power	-162.8 dbm W	-137.33 db	+25.47 db
Bit Rate	(8-1/3 bit/sec) -9.2 db/s	(16.2 $\times 10^3$ bit/sec) -42.10 db/s	-32.90 db
Bit Energy	-172.0 dbm Ws	-179.43 dbm Ws	-7.43 db
Receiver Noise Temperature	(65°K) -180.5 dbm W/Hz	(25°K) -184.60 dbm W/Hz	*+4.10 db
Nominal Signal-to-Noise Ratio Per Bit	+8.5 db	+5.17 db	-3.33 db
Signal-to-Noise Ratio Required for Given Bit Error Probability	( $5 \times 10^{-3}$ ) +5.2 db	( $5 \times 10^{-3}$ ) +3.00 db	*+2.20 db
Margin	+3.3 db	+2.17 db	*-1.13 db
Total Gain in Bit Rate (Sum of starred items in last column)	-----	-----	+32.90 db

16

Combinatorial Structures In Planetary Reconnaissance

表 13: Mariner 1964 与 Mariner 1969 的远程通信系统对比 [15]

表 13 列出了这  $16200/8\frac{1}{3} = 1944$  倍 (32.9dB) 速度提升的各种因素。Mariner 4 的码率仅有  $8\frac{1}{3}$  bps，发射功率是 8.9W，也就是说每传送 1 bit 需要约 1.07 焦耳的能量。Mariner 1969 的发射功率翻倍

<sup>①</sup>[https://en.wikipedia.org/wiki/Mariner\\_6\\_and\\_7](https://en.wikipedia.org/wiki/Mariner_6_and_7)

<sup>②</sup><https://www.jpl.nasa.gov/news/mariner-television-cameras-experiment/>

<sup>③</sup><https://computerhistory.org/blog/who-invented-the-ic/>

<sup>④</sup>[https://en.wikipedia.org/wiki/Diode-transistor\\_logic](https://en.wikipedia.org/wiki/Diode-transistor_logic)

到 18.2W，同时传送每 bit 所需的能量降为原来的约 1/1000。一个重要因素是信号的传输距离缩短了一多半 ( $216 \times 10^6 \text{km} \rightarrow 97 \times 10^6 \text{km}$ )，根据平方反比律，信号强度相当于原来的  $(216/97)^2 \approx 4.96$  倍，这对应表中第 4 行 Space Loss。此外，地面接收天线的尺寸也成倍增大 (85 英尺  $\rightarrow$  210 英尺，合 64 米)，相应的接收面积增大为  $(210/85)^2 \approx 6.1$  倍，收到的信号能量也同比增加，这对应表中第 5 行。这几项物理因素占了全部 32.9dB 增益的一半。

根据距离和光速  $3 \times 10^5 \text{km/s}$  可以算出 Mariner 4 数据链路的单程带宽延迟积 (BDP) 是  $8\frac{1}{3} \times \frac{216 \times 10^6}{3 \times 10^5} = 6000 \text{ bits}$ ；Mariner 1969 数据链路的单程 BDP 是  $16200 \times \frac{97 \times 10^6}{3 \times 10^5} \approx 5.24 \text{ Mbits}$ 。后面这个数据和几毫秒延迟下的千兆以太网相当，因此在思考计算机网络延迟的时候可以想象成用 16 kbps 的 Modem 往光程 5 分钟远的火星发送数据，有时会有意想不到的观察。

在众多提速因素当中，Mariner 1969 信道编码仅提供了 2.2dB 的提升，这里我们略微深入分析一下这 2.2dB 是怎么得来的，主要是要计算 AWGN 噪声<sup>①</sup> 下的 BER (误比特率 / bit error rate)。Mariner 4 用的大概是 BPSK 调制，设计的 BER 是  $5 \times 10^{-3}$ ，根据 BPSK 的 BER 公式<sup>②</sup>  $P_b = Q(\sqrt{2E_b/N_0})$  可以算出所需的信噪比  $E_b/N_0 = 5.2\text{dB}$ ，这是表 13 倒数第 3 行。其中 Q 函数<sup>③</sup> 是标准正态分布的互补累积分布函数，定义及与误差函数的关系如下

$$Q(x) = P(X > x) = \frac{1}{\sqrt{2\pi}} \int_x^\infty e^{-x^2/2} dx, \quad Q(x) = 1 - \Phi(x) = \frac{1}{2} \operatorname{erfc}\left(\frac{x}{\sqrt{2}}\right) = \frac{1}{2} - \frac{1}{2} \operatorname{erf}\left(\frac{x}{\sqrt{2}}\right)$$

Q 函数及其反函数的 Python 实现：

```
from scipy import special

def qfunc(x):
    # ndtr: cumulative distribution of the standard normal distribution.
    return 1 - special.ndtr(x)

def qfuncinv(x):
    return special.ndtri(1 - x)
```

验算 Mariner 4 误码率：

```
def bpsk_ber(eb_n0):
    return qfunc(sqrt(2 * eb_n0))

def ber_to_bpsk_snr(ber):
    return qfuncinv(ber) ** 2 / 2

ber = 5e-3
snr = ber_to_bpsk_snr(ber)
db = 10 * log10(snr)
print("SNR %.3f, %.3f dB" % (snr, db))

# 输出
# SNR 3.317, 5.208 dB
```

<sup>①</sup>加性高斯白噪声 additive white Gaussian noise，这是深空通信的特性。

<sup>②</sup>简单说一句 BPSK 的 BER 公式是怎么来的，单 bit 信号能量  $E_b$ ，噪声功率谱密度  $N_0$ ，噪声  $X \sim \mathcal{N}(\mu = 0, \sigma^2 = \frac{1}{2} N_0)$ ，误码率  $P(X > \sqrt{E_b}) = Q(\sqrt{E_b}/\sigma) = Q(\sqrt{2E_b/N_0})$ 。<https://dsplog.com/2007/08/05/bit-error-probability-for-bpsk-modulation/>

<sup>③</sup><https://en.wikipedia.org/wiki/Q-function>

Mariner 1969 采用 Hadamard [32, 6, 16] 码, 信源以  $k = 6$  bit 为一组进行编码, 得到  $n = 32$  bit 码字。为了公平对比起见, 这个码字的能量  $E_s = k E_b$ 。

我们先给出 orthogonal signals 的 SER (symbol error rate) 公式, 再介绍 biorthogonal signals 的 SER 计算方法。 $M$ -ary orthogonal signals 是一组能量相等信号  $s_1, s_2, \dots, s_M$ , 信号一共有  $M$  个, 长度均为  $M$ , 而且这些信号是彼此正交的, 数字通信的教材上常以  $M$ -FSK 举例。 $M$  阶 Hadamard 矩阵的  $M$  个行向量也是符合条件的一组  $M$ -ary orthogonal signals。 $M$ -ary orthogonal signals 的 SER 公式:

$$\begin{aligned} P_s = 1 - P_c &= 1 - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-u^2/2} \left[ 1 - Q(u + \sqrt{2E_s/N_0}) \right]^{M-1} du \\ &= 1 - \int_{-\infty}^{\infty} \phi(u) \left[ \Phi(u + \sqrt{2E_s/N_0}) \right]^{M-1} du \end{aligned} \quad (1)$$

上式中  $\phi(\cdot)$  和  $\Phi(\cdot)$  分别是标准正态分布的 PDF 和 CDF。简单聊两句以上公式是怎么来的。

不失一般性, 假设发送  $s_1 = (\sqrt{E_s}, 0, \dots, 0)$ , 收到  $r = (\sqrt{E_s} + n_1, n_2, \dots, n_M)$ , 其中  $n_1, n_2, \dots, n_M$  是相互独立的高斯随机变量  $\mathcal{N}(\mu = 0, \sigma^2 = \frac{1}{2}N_0)$ 。正确判决的概率  $P_c = P(\sqrt{E_s} + n_1 > n_2, \sqrt{E_s} + n_1 > n_3, \dots, \sqrt{E_s} + n_1 > n_M)$ , 如果将判决门限  $\sqrt{E_s} + n_1$  视为已知, 那么  $P_c[n_1 \text{ 确定}] = P(n_2 < \sqrt{E_s} + n_1) \cdot P(n_3 < \sqrt{E_s} + n_1) \cdots P(n_M < \sqrt{E_s} + n_1) = \Phi\left(\frac{\sqrt{E_s} + n_1}{\sigma}\right) \cdot \Phi\left(\frac{\sqrt{E_s} + n_1}{\sigma}\right) \cdots \Phi\left(\frac{\sqrt{E_s} + n_1}{\sigma}\right) = \left[\Phi\left(\frac{\sqrt{E_s} + n_1}{\sigma}\right)\right]^{M-1}$ 。这里第一个等号成立是因为  $n_2, n_3, \dots, n_M$  是独立的随机变量  $P(A \cap B) = P(A) \cdot P(B)$ , 第二个等号成立是因为  $n_2, n_3, \dots, n_M$  具有相同的分布  $\mathcal{N}(\mu = 0, \sigma^2 = \frac{1}{2}N_0)$ 。

既然  $n_1$  本身也是一个高斯随机变量, 用连续全概率公式可得

$$P_c = \int_{-\infty}^{\infty} \left[ \Phi\left(\frac{\sqrt{E_s} + n_1}{\sqrt{N_0/2}}\right) \right]^{M-1} p(n_1) dn_1 \quad \text{其中 } p(n_1) = \frac{1}{\sqrt{\pi N_0}} e^{-n_1^2/N_0}$$

再令  $u = n_1/\sqrt{N_0/2}$  带入上式化解可得上面的公式 (1)。这个公式没有解析式, 一般用数值积分求值。在信噪比较高的时候也可以用 union bound 估算  $P_s < (M-1)Q(\sqrt{E_s/N_0})$ 。BER 与 SER 的关系  $P_b = \frac{M}{2M-2} P_s \approx \frac{1}{2} P_s$ 。用以上公式我们绘制了  $M = 32$  时 orthogonal signals 的 BER 与 BPSK 对比 (下页图 24), 图中也包括  $M = 64$  biorthogonal signals 的 BER。其实在  $M$  相等时, orthogonal 和 biorthogonal 的 BER 曲线几乎重叠, 只是 biorthogonal 占用的带宽少一半。

下面我们谈谈  $M$ -ary biorthogonal signals 的 SER 计算。 $M$ -ary biorthogonal signals  $\{s_1, s_2, \dots, s_M\}$  由两部分组成:  $\frac{M}{2}$  个 orthogonal signals 和这些信号的相反数, 每个信号的长度是  $\frac{M}{2}$ 。Hadamard [32, 6, 16] 码的  $2^6 = 64$  个码字是  $M = 64$  的 biorthogonal signals, 它包括 32 阶 Hadamard 矩阵的全部行向量以及负向量。 $M$ -ary biorthogonal signals 的 SER 公式:

$$\begin{aligned} P_s = 1 - P_c &= 1 - \frac{1}{\sqrt{2\pi}} \int_{-\sqrt{2E_s/N_0}}^{\infty} e^{-u^2/2} \left[ \operatorname{erf}\left(\frac{u + \sqrt{2E_s/N_0}}{\sqrt{2}}\right) \right]^{M/2-1} du \\ &= 1 - \frac{1}{\sqrt{2\pi}} \int_0^{\infty} e^{-(u - \sqrt{2E_s/N_0})^2/2} [1 - 2Q(u)]^{M/2-1} du \end{aligned} \quad (2)$$

据我查到的资料, 公式 (1) 和 (2) 最早是 Andrew J. Viterbi 在 1961 年发表 [19], 他时年 26 岁。

BERs of modulation methods

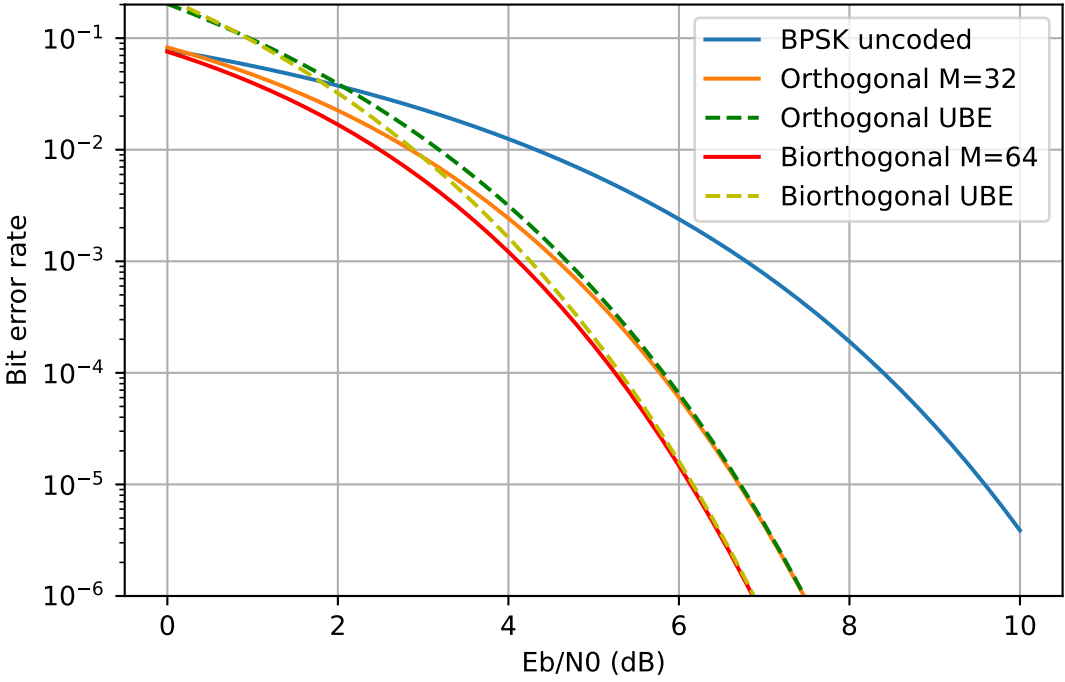


图 24: 几种调制方案的 BER 对比

为了便于推导以上 SER 公式 (2), 不妨取如下  $M$ -ary biorthogonal signals

$$\begin{aligned}
 s_1 &= (\sqrt{E_s}, 0, 0, \dots, 0) & s_2 &= (-\sqrt{E_s}, 0, 0, \dots, 0) = -s_1 \\
 s_3 &= (0, \sqrt{E_s}, 0, \dots, 0) & s_4 &= (0, -\sqrt{E_s}, 0, \dots, 0) = -s_3 \\
 &\vdots & &\vdots \\
 s_{M-1} &= (0, 0, 0, \dots, \sqrt{E_s}) & s_M &= (0, 0, 0, \dots, -\sqrt{E_s}) = -s_{M-1}
 \end{aligned}$$

假设发送的是  $s_1$ , 收到  $\mathbf{r} = (\sqrt{E_s} + n_1, n_2, \dots, n_{M/2})$ , 其中  $n_1, n_2, \dots, n_{M/2}$  是相互独立的高斯随机变量  $\mathcal{N}(\mu = 0, \sigma^2 = \frac{1}{2}N_0)$ 。最优的判决方法是先找到  $\mathbf{r}$  中幅度最大的分量  $r_i$ , 再根据  $r_i$  的正负号推断发送的是  $s_{2i-1}$  还是  $s_{2i}$ 。

对于发送  $s_1$  的情况, 判决得到正确结果的条件是  $r_1 = \sqrt{E_s} + n_1 > 0$  且  $r_1$  大于  $|n_2|, |n_3|, \dots, |n_{M/2}|$ 。如果将  $r_1 = \sqrt{E_s} + n_1$  视为已知, 那么  $P_c[r_1 \text{ 确定}] = P(|n_2| < r_1, |n_3| < r_1, |n_{M/2}| < r_1)$ 。考察其中一项  $P(|n_2| < r_1 \mid r_1 > 0) = \frac{1}{\sqrt{\pi N_0}} \int_{-r_1}^{r_1} e^{-x^2/N_0} dx = \text{erf}\left(\frac{r_1}{\sqrt{N_0}}\right)$ 。由于  $n_2, n_3, \dots, n_{M/2}$  是 i.i.d., 因此  $P_c[r_1 \text{ 确定}] = \left[\text{erf}\left(\frac{r_1}{\sqrt{N_0}}\right)\right]^{M/2-1}$ 。又因为  $r_1 \sim \mathcal{N}(\mu = \sqrt{E_s}, \sigma^2 = \frac{1}{2}N_0)$ , 有

$$P_c = \int_0^\infty \left[\text{erf}\left(\frac{r_1}{\sqrt{N_0}}\right)\right]^{M/2-1} p(r_1) dr_1 \quad \text{其中 } p(r_1) = \frac{1}{\sqrt{\pi N_0}} e^{-(r_1 - \sqrt{E_s})^2/N_0}$$

令  $u = r_1/\sqrt{N_0/2}$  带入上式，化解后就得到了公式 (2) 中的第二个式子。在信噪比较高时可以用 union bound 估算  $P_s < (M-2)Q(\sqrt{E_s/N_0}) + Q(\sqrt{2E_s/N_0})$ 。BER  $P_b \approx \frac{1}{2}P_s$ ， $M = 64$  的 BER 曲线见图 24。

为了验算 Mariner 1969 纠错码的 2.2dB 编码增益，我们计算达到相同 BER =  $5 \times 10^{-3}$  时 Hadamard [32, 6, 16] 码 ( $k = 6$ ,  $M = 2^k = 64$ ) 所需的  $E_b/N_0$ ，用 scipy 的数值积分和数值求根算出 SER =  $10^{-2}$  对应的  $E_b/N_0$ ：

```
def biorthogonal_ser(k, eb_n0):
    M = 2 ** k
    es_n0 = k * eb_n0
    Pc = integrate.quad(lambda u: exp(-u**2/2) * special.erf(u/sqrt(2)+sqrt(es_n0))**(M/2-1),
                        -sqrt(2*es_n0), np.inf)[0] / sqrt(2*math.pi)
    return (1-Pc)

desired_ser = 1e-2
root = optimize.root_scalar(lambda snr: biorthogonal_ser(6, snr) - desired_ser, x0 = 1)
if root.converged:
    snr = root.root
    print('SER %g, BER %g, Eb/N0 %.3f dB' % (desired_ser, biorthogonal_ser(6, snr)/2, decibel(snr)))
```

输出 SER 0.01, BER 0.005, Eb/N0 3.057 dB，因此与 Mariner 4 的 5.208 dB 相比，节约了 2.151 dB。

1971 发射的 Mariner 9 是第一个环绕其他行星运行的航天器，它传回了 7000 多张火星表面的照片，分辨率  $832 \times 700$ ，每个像素是 9-bit。<sup>①</sup> 1973 年发射的 Mariner 10 是第一个造访水星的航天器，也是第一艘飞掠多颗行星的航天器。Mariner 计划在科学史上创造了多个第一，拍摄照片和应用纠错码是其一小部分成就。

1968 年底发射的 Pioneer 9 号飞船是第一个进入太空的纠错码（卷积码），比 1969 年发射的 Mariner 6/7 早了几个月。[20, 21, 22]

### 2.3 Hadamard 码在 CDMA95 中的应用

CDMA（码分多址）是 2G 和 3G 移动通信的一种“制式”，Hadamard 码在其中扮演了重要角色，在 CDMA 文献中，称其为 Walsh 码或 Walsh 函数。这里以 2G 的 CDMA95 (IS-95) 为例，介绍 Hadamard 码的用处。在上行信道（手机到基站，以前的文献也叫反向链路），Hadamard 码是当作纠错码来用的。在下行信道（基站到手机，也叫前向链路），Hadamard 码的作用不是纠错，而是扩频，即 channelize（信道化）。

图 25 是 CDMA95 上行信道的信道编码简图。大致流程如下：语音信号按 8KHz 采样，然后每 20ms 分为一段 (segment)，每段 160 个样点 (samples) 做语音编码（信源编码），得到 172 bits 语音数据（一说 171 bits）。接着添加 12 bit 的检错码（如果接受方检测到错误，就把这 20ms 的语言静音）和 8 bit 零（用于终结紧接其后的卷积码编码器），得到 192 bits / 20ms 的码流，即 9 600 bps。然后经过码率为 1/3 的卷积码编码，得到  $9.6 \times 3 = 28.8\text{kbps}$  的码流，再经过一个  $32 \times 18$  的交织器，得到 28.8kbps 码流。接下来，用 Hadamard [64, 6] 编码，每 6 bit 为一个消息  $m$ ，编码得到 64-bit 的码字  $c$ ，这样得到  $28.8/6 \times 64 = 307.2\text{kbps}$  码流。然后与 1228.8 kbps 的长码（PN 伪随机码）叠加，最后的 1.2288Mbps

<sup>①</sup>Guide to the use of Mariner images. <https://ntrs.nasa.gov/citations/19760004316>

信号经过短 PN 码和 OQPSK 调制到 800MHz 载波上，通过天线发射出去。换句话说，上行信道用了级联码 (concatenated code)，外码是卷积码，内码是 Hadamard [64, 6] 分组码。相应的，基站的译码过程会先用 64 点的快速 Walsh-Hadamard 变换解 Hadamard 码，再用 Viterbi 算法解卷积码。

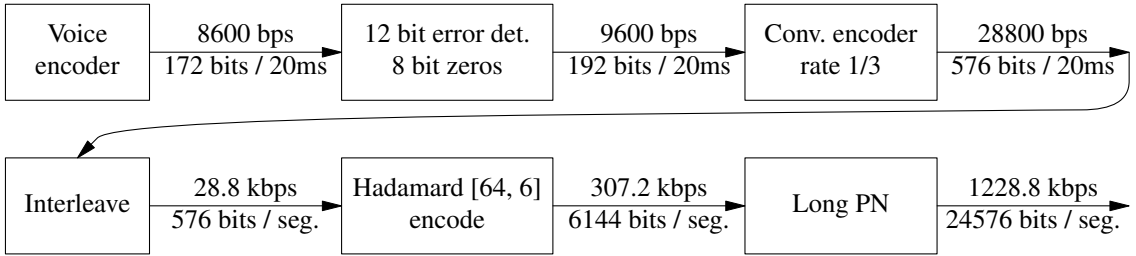


图 25: CDMA95 上行信道 (手机 → 基站) 发送方简图

下行信道也用到了 64 阶 Hadamard 矩阵，但其作用不是纠错，而是利用其行向量彼此正交的特性来实现扩频。图 26 是 CDMA95 下行信道的信道编码简图。大致流程如下：8 600bps 的有效载荷 (payload) 按 172bits / 20ms 分段，每段添加 12 bit 检错码和 8 bit 零 (作用和上行信道相同)，得到 192 bits / 20ms 的码流，即 9 600 bps。然后经过码率为 1/2 的卷积码编码，得到  $9.6 \times 2 = 19.2\text{kbps}$  的码流，再经过一个  $24 \times 16$  的交织器，并和长码叠加，得到 19.2kbps 码流。到目前为止，除了卷积码的码率不同，其他方面无甚区别，接下来轮到 Hadamard 码 / Walsh 码出场。

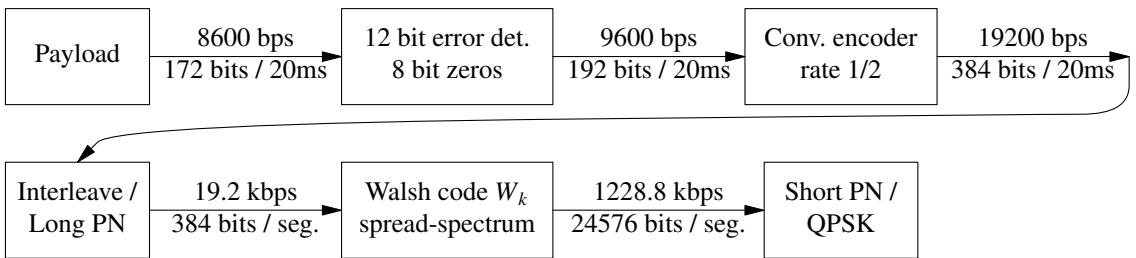


图 26: CDMA95 下行信道 (基站 → 手机) 发送方简图

为了区分不同的移动用户 (手机)，基站给每个手机分配 64 阶 Hadamard 矩阵其中的一行作为标识 (图中  $w_k$  码)，也就是说一个基站在一个频段上可以最多同时连接 64 台移动设备 (实际上有些行是保留的，如  $w_0, w_1, \dots, w_7, w_{32}$  等)。基站向第  $k$  号移动设备每发送 1 bit 时，如果是 1，就发送  $w_k$ ，是 0 则发送  $\bar{w}_k$ ，即  $w_k$  按位取反。这样一来，码率就提升到了  $19.2 \times 64 = 1228.8\text{kbps}$ ，再经过 QPSK 调制到载频上发射出去。由于 Hadamard 矩阵的各行是正交的，即两个用户的  $\langle w_i, w_j \rangle = 0$ ，对于接收方来说，只要用自己的  $w_k$  与收到的信号做一次相关运算 (correlation)，就能筛选出发给自己的信息，而发给其他用户的信息和噪声没有区别。换句话说，借助 Hadamard 码，CDMA95 在同一个 1.25MHz 的频段上分出了 64 个信道，这正是码分多址的意义。

在 3G 的 CDMA2000 中，用的是 256 阶的 Hadamard 矩阵。

### 3 Golay 码与旅行者号

Marcel J. E. Golay 在 1949 年用半页的篇幅介绍了二进制 [23, 12, 7] 码和三进制 [11, 6, 5]<sub>3</sub> 码 [23], 这里我们只谈谈二进制 Golay 码。Golay [23, 12, 7] 码是线性分组码, 其码字有 23 bit, 其中前 12 bit 是信息位, 后 11 bit 是监督位, 一共有  $2^{12} = 4096$  个合法码字。

Golay 论文的核心内容是给出了这个码的监督矩阵  $H_{23} = [P | I_{11}]$ , 矩阵  $H_{23}$  分左右两块, 右边是 11 阶方阵  $I_{11}$ , 左边是  $11 \times 12$  矩阵  $P$ , 如右图 27 所示。用 p.13 的 find\_d\_min() 函数可以验证其最小码距  $d_{\min} = 7$ , 因此它能纠正  $\lfloor \frac{7-1}{2} = 3 \rfloor$  bit 错, 这是第一个能纠正多 bit 错的码。

	$Y_1$	$Y_2$	$Y_3$	$Y_4$	$Y_5$	$Y_6$	$Y_7$	$Y_8$	$Y_9$	$Y_{10}$	$Y_{11}$	$Y_{12}$
$X_1$	1	0	0	1	1	1	0	0	0	1	1	1
$X_2$	1	0	1	0	1	1	0	1	1	0	0	1
$X_3$	1	0	1	1	0	1	1	0	1	0	1	0
$X_4$	1	0	1	1	1	0	1	1	0	1	0	0
$X_5$	1	1	0	0	1	1	1	0	1	1	0	0
$X_6$	1	1	0	1	0	1	1	1	0	0	0	1
$X_7$	1	1	0	1	1	0	0	1	1	0	1	0
$X_8$	1	1	1	0	0	1	0	1	0	1	1	0
$X_9$	1	1	1	0	1	0	1	0	0	0	1	1
$X_{10}$	1	1	1	1	0	0	0	0	1	1	0	1
$X_{11}$	0	1	1	1	1	1	1	1	1	1	1	1

图 27: Golay 码论文核心内容——监督矩阵 P

```
P = np.array([
    [1, 0, 0, 1, 1, 1, 0, 0, 0, 1, 1, 1],
    [1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 0, 1],
    [1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0],
    [1, 0, 1, 1, 1, 0, 1, 1, 0, 1, 0, 0],
    [1, 1, 0, 0, 1, 1, 1, 0, 1, 1, 0, 0],
    [1, 1, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1],
    [1, 1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0],
    [1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1, 0],
    [1, 1, 1, 0, 1, 0, 1, 0, 0, 0, 1, 1],
    [1, 1, 1, 1, 0, 0, 0, 0, 1, 1, 0, 1],
    [0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1]], dtype=int)

H23 = np.concatenate([P, np.eye(11, dtype=int)], axis=1)
print(find_d_min(H23))

# 输出:
# (0, 1, 2, 3, 4, 5, 22)
# 7
```

Golay 码是系统码, 用 p.11 的方法很容易求得其生成矩阵  $G_{23} = [I_{12} | P^T]$ , 左边是 12 阶方阵  $I_{12}$ , 右边是矩阵  $P$  的转置。编码过程可描述为  $c = G_{23}^T \cdot m$ , 实践中有更节约硬件资源的编码方法。<sup>①</sup>

跟 Hamming 码一样, Golay [23, 12, 7] 码也是完备码 (perfect code), 意味着 23 维的 Hamming 空间的全部  $2^{23}$  个点被完美地划分为  $2^{12}$  个互不相交的 Hamming 球, 每个球包含  $2^{11}$  个点,  $2^{12} \cdot 2^{11} = 2^{23}$ 。这里每个 Hamming 球的体积是, 以任意一个码字为中心, 与其相距 0, 1, 2, 3 bit 的点的个数:

$$\binom{23}{0} + \binom{23}{1} + \binom{23}{2} + \binom{23}{3} = 1 + 23 + \frac{23 \times 22}{2} + \frac{23 \times 22 \times 21}{2 \times 3} = 2048 = 2^{11}$$

完备码在数学上比较美<sup>②</sup>, 但其实对通信 (纠错) 来说似乎没有特别大的意义。

Golay [23, 12, 7] 码的译码出现过多种方法 [24, 25, 26, 28]。在如今的软硬件条件下, 我认为可以简单地用查表法, 对于接收到的 23-bit 向量  $r$ , 先计算 11-bit syndrome  $s = H_{23}r$ , 然后以  $s$  为下标查一个长度为 2048 的表找出错误图样  $e$ , 最后  $\hat{c} = r \oplus e$  完成译码。

<sup>①</sup><https://aqdi.com/articles/using-the-golay-error-detection-and-correction-code-3/>

<sup>②</sup><https://giam.southernct.edu/DecodingGolay/decoding.html>

给 Golay 码添加 1 bit 监督位，就得到了  $d_{\min} = 8$  的扩展 Golay [24, 12, 8] 码，这个码的生成矩阵可以在 Wikipedia 上看到<sup>①</sup>。扩展 Golay 码曾用在 1977 年发射的旅行者 (Voyager) 行星探索计划上，因为旅行者号拍摄的是彩色照片 (分辨率  $800 \times 800$ )，数据量比黑白照片大很多，为此采用了图像压缩技术，这对数据的码率 (最高达到 115,200 bits/s) 和误码率有更高的要求，纠错码是必须的。[29]

旅行者 1 号于 1977 年 9 月 5 日发射，它分别在 1979 年 3 月和 1980 年 11 月飞掠 (fly-by) 了木星 (Jupiter) 和土星 (Saturn)。在此之前，Pioneer 10 和 Pioneer 11 已造访过这两颗行星。<sup>②</sup>

旅行者 2 号的发射日期其实比旅行者 1 号早两周，利用百年难遇的太阳系行星排布，借助引力弹弓效应，它一口气造访了四颗行星，分别是木星 (1979 年 7 月)、土星 (1981 年 8 月)、天王星 (1986 年 1 月)、海王星 (1989 年 8 月)。图 28 展示了两艘旅行者号飞船的飞行轨迹。旅行者 2 号创下了多项记录，它是目前为止惟一造访过天王星 (Uranus) 和海王星 (Neptune) 的探测器。[30, 31]

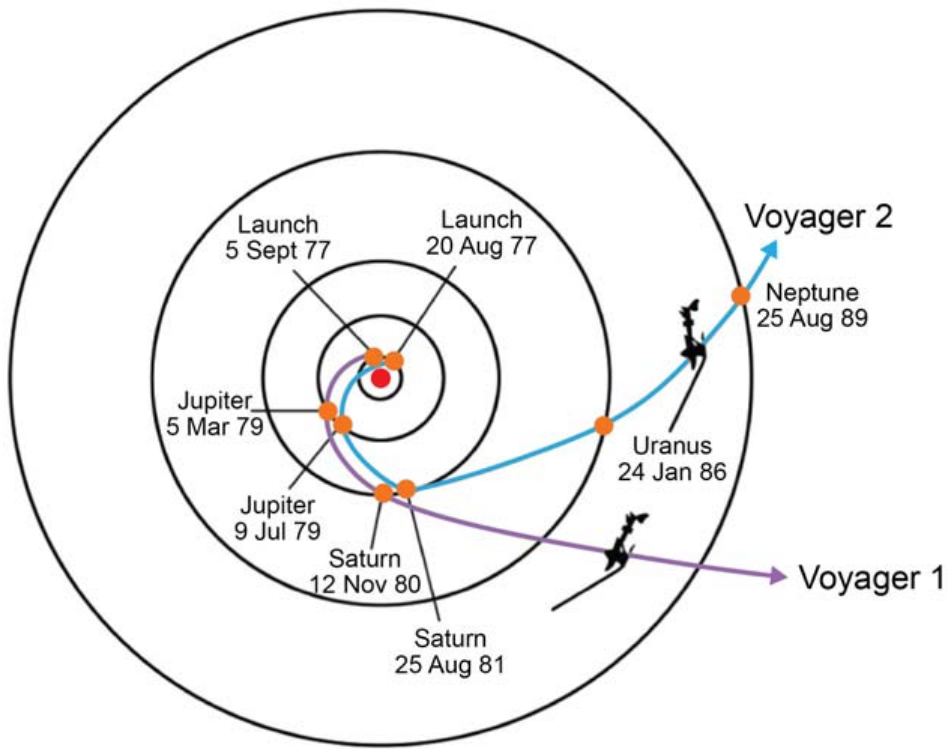


图 28: 旅行者计划的飞行轨迹

旅行者号在探索木星和土星时用的纠错码是码率  $R = k/n = 1/2$  的扩展 Golay [24, 12, 8] 码，在此之后换成了码率更高的 Reed-Solomon (255, 223) 码。<sup>[32, 33]</sup> 具体来说，旅行者号用的是级联码，内码是  $k = 7, R = 1/2$  的卷积码 (Viterbi 软判决译码)，外码是 Golay 码或 Reed-Solomon 码 (均为硬判决译码)，因此综合的码率是外码的一半。<sup>[22]</sup> Voyager 的级联码方案在 NASA 深空通信系统中用了很多年。

有人认为其实扩展 Golay 码的性能不如原始 Golay 码。[35]

<sup>①</sup>[https://en.wikipedia.org/wiki/Binary\\_Golay\\_code](https://en.wikipedia.org/wiki/Binary_Golay_code)

<sup>②</sup>[https://en.wikipedia.org/wiki/List\\_of\\_missions\\_to\\_the\\_outer\\_planets](https://en.wikipedia.org/wiki/List_of_missions_to_the_outer_planets)

## 参考文献

- [1] Daniel J. Costello and G. David Forney, "Channel coding: The road to channel capacity," *Proceedings of the IEEE*, vol. 95, issue 6, June 2007.
- [2] David J. C. MacKay, *Information Theory, Inference, and Learning Algorithms*, Cambridge University Press, 2003.
- [3] Claude E. Shannon, "A mathematical theory of communication," *Bell System Technical Journal*, vol. 27, pp. 379–423 and 623–656, 1948.
- [4] Richard W. Hamming, "Error detecting and error correcting codes," *Bell System Technical Journal*, vol. 29, pp. 147–160, 1950.
- [5] David S. Slepian, "A class of binary signalling alphabets," *Bell System Technical Journal*, vol. 35, pp. 203–234, 1956.
- [6] W. Wesley Peterson, *Error-Correcting Codes*, MIT Press, 1961.
- [7] Elwyn R. Berlekamp, Robert J. McEliece, and Henk C. A. van Tilborg, "On the inherent intractability of certain coding problems," *IEEE Transactions on Information Theory*, vol. 24, no. 3, pp. 384–386, 1978.
- [8] Richard C. Singleton, "Maximum distance  $q$ -nary codes," *IEEE Transactions on Information Theory*, vol. 10, no. 2, pp. 116–118, 1964.
- [9] Mu-Yue Hsiao, "A class of optimal minimum odd-weight-column SEC-DED codes," *IBM Journal of Research and Development*, vol. 14, no. 4, pp. 395–401, 1970.
- [10] Minesh Patel, Jeremie S. Kim, Hasan Hassan, and Onur Mutlu, "Understanding and Modeling On-Die Error Correction in Modern DRAM: An Experimental Study Using Real Devices," *Proceedings of the 49th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN 2019)*, Portland, OR, USA, June 2019.
- [11] David F. Brailsford, Building the  $H_5$  Reed-Muller Code used by Mariner 9. Reed-Muller Code (64 Shades of Grey pt2) - Computerphile <https://www.youtube.com/watch?v=Ct0CqKpti7s>
- [12] J. O. Duffy, "Detailed Design of a Reed-Muller (32,6) Block Encoder," *JPL Space Programs Summary 37-47* vol. III, pp. 263–267, 1967.
- [13] Leonard D. Baumert and Howard C. Rumsey, Jr., "Combinatorial Communications: The Index of Comma Freedom for the Mariner Mars 1969 High Data Rate Telemetry Code," *JPL Space Programs Summary 37-46* vol. IV, pp. 221–226, 1967.
- [14] Richard R. Green, "A serial orthogonal decoder," *JPL Space Programs Summary 37-39* vol. IV, pp. 247–253, 1966.
- [15] Edward C. Posner, "Combinatorial Structures in Planetary Reconnaissance," *Error Correcting Codes*, ed. Henry B. Mann, Wiley, NY 1968.
- [16] Richard P. Mathison, "Mariner Mars 1964 Telemetry and Command System," *JPL Technical Report 32-684*, 1965.

- [17] John A. Hunter, "Mariner Mars 1964 Telecommunication System," JPL Technical Report 32-836, 1965.
- [18] Robert C. Tausworthe, Mahlon F. Easterling, and A. J. Spear, "A High-Rate Telemetry System for the Mariner Mars 1969 Mission," JPL Technical Report 32-1354, 1969.
- [19] Andrew J. Viterbi, "On Coded Phase-Coherent Communications," *IRE Transactions on Space Electronics and Telemetry*, vol. SET-7, no. 1, pp. 3–14, March 1961,
- [20] James L. Massey, "Deep-Space Communications and Coding: A Marriage Made in Heaven," *Lecture Notes in Control and Information Sciences*, vol. 192, ed. Joachim Hagenauer, Springer 1992.
- [21] G. David Forney, "Coding and its application in space communications," *IEEE Spectrum*, June 1970.
- [22] Michael Rice, "A History of Channel Coding in Aeronautical Mobile Telemetry and Deep-Space Telemetry," *Entropy* 26(8):694, 2024.
- [23] Marcel J. E. Golay, "Notes on digital coding," *Proceedings of the I.R.E.*, vol. 37, pp. 657, 1949.  
Elwyn R. Berlekamp 将其誉为 "**best single published page**" in coding theory.
- [24] Elwyn R. Berlekamp, "Decoding the Golay Code," JPL Technical Report 32-1526, vol. XI, 1972.
- [25] Vera Pless, "Decoding the Golay Codes," *IEEE Transactions on Information Theory*, vol. 32, no. 4, July 1986.
- [26] T. K. Truong, J. K. Holmes, I. S. Reed, X. Yin, "A Simplified Procedure for Decoding the (23,12) and (24,12) Golay Codes," TDA Progress Report 42-96, 1988.
- [27] T. K. Truong, I. S. Reed, X. Yin, "Decoding of 1/2-Rate (24,12) Golay Codes," TDA Progress Report 42-97, 1989.
- [28] H. P. Lee, C. H. Chang, S. I. Chu, "High-Speed Decoding of the Binary Golay Code," *Journal of Applied Research and Technology*, vol. 11, June 2013.
- [29] Robert F. Rice, "Channel Coding and Data Compression System Considerations for Efficient Communication of Planetary Imaging Data," JPL Technical Memorandum 33-695, 1974.
- [30] E. C. Stone and E. D. Miner, "The Voyager 2 Encounter with the Uranian System," *Science* vol. 233, pp. 39–43, July 1986.
- [31] E. C. Stone and E. D. Miner, "The Voyager 2 Encounter with the Neptunian System," *Science* vol. 246, pp. 1417–21, Dec 1989.
- [32] Richard P. Larser, William I. McLaughlin, and Donna M. Wolff, "Engineering Voyager 2's encounter with Uranus," *Scientific American*, vol. 255, pp. 36–45, Nov 1986.
- [33] Stephen B. Wicker and Vijay K. Bhargava, "Reed-Solomon Codes and the Exploration of the Solar System," in *Reed-Solomon Codes and Their Applications*, Wiley-IEEE, pp. 25–40, 1994.
- [34] Roger Ludwig and Jim Taylor, "Voyager Telecommunications," 2002.
- [35] Jon Hamkins, "The Golay Code Outperforms the Extended Golay Code Under Hard-Decision Decoding," <https://arxiv.org/abs/1602.05620>, 2016.