

第 1 章 绪论

解决数据压缩的问题通常可以从三步来分析：第一步是为什么要做，即数据压缩的必要性问题；第二步是为什么可以做，即分析信源数据的特性，并在此基础上进行数据压缩的可行性分析；第三步是在第二步分析的基础上，如何借助信号处理等技术实现对于数据的压缩。以下进行详细的分析。

一、数据压缩的必要性

在数据压缩处理的各种媒体中，视频和音频信息占据了主导地位。这是因为人们所接受的外部信息中，绝大部分是由视觉和听觉获得的。人们对视听享受的要求越来越高，视频从标清视频到高清再到超高清视频，音频从单声道到立体声、环绕立体声和 22.2 多声道等。下述表格中仅给出不同数字化视频格式中亮度信号的原始码率，包含色度信号的总码率在后续第二章媒体特性与表示中介绍彩色空间后给出。（视音频原始格式的分析也在第二章中给出）

表 1 数字化视频格式（仅计算亮度信号）

数字视频格式	每秒帧数	图像分辨率（宽×高，像素）	样本精度（bit）	原始码率（Mb/s）
CIF	30	352×288	8	24.33
ITU-R 601 PAL	25	720×576	8	82.944
HDTV 1080P	30	1920×1080	8	497.664
UHDTV1(4K)	60	3840×2160	10	4976.64
UHDTV2(8K)	120	7680×4320	10	39813.12

由上述数据可见：在人们常用的数据格式中，视频、图像、音频和语音等数据量是相当巨大的，尤其是视频图像的数据量。这仅是人们日常生活中经常接触到的数据类型。为探索更多人类未知的领域，诸如航空探测、基因分析、地质勘探等活动每天将产生以太字节和拍字节计的数据量（太： 10^{12} ，拍： 10^{15} ）。即使目前的存储介质成本迅速下降，但仍然赶不上数据量急剧增加的速度。从另一方面说，有些情况下的存储和传输能力也无法显著提高。例如，通过无线电波进行数据传输总是受到大气特性的限制。因此，高效数据压缩技术的研究和实现仍是人们坚持不懈的追求目标。

对原始信源的分析是非常重要的，不仅需要理解其数据构成的原理，还要熟悉其数据格式。例如：为什么视频压缩的对象是 YUV，而不是 RGB？音频数据是双极性信号，它的采样率跟人耳能听到的频率范围有关。这些都是基本常识，但往往被忽略，导致在系统设计的时候出现问题。

二、数据压缩的可行性和实现思路

一个典型的数据压缩系统通常包括两部分，即数据压缩和解压缩（也称重建，*reconstruction*），如图 1 所示。压缩算法将输入的原始数据 x 转换为更少码元表示的 y ，相应的解压缩（重建）算法对压缩数据流 y 进行还原，得到解压缩后的数据 \hat{x} 。数据压缩和解压缩是密不可分的，通常合称为 **codec**（编解码器）。

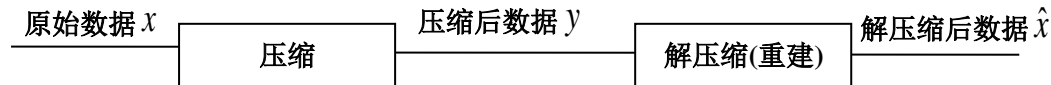


图 1-1 典型的数据压缩系统

之所以可以对信源数据进行压缩，是因为绝大多数信源数据中存在结构上的特点。这种结构上的特点便于我们利用其特征进行压缩。在现有的数据压缩技术中，主要是从三个方面着手进行压缩。

1. 利用统计结构特征进行压缩

什么叫信源的统计结构特征呢？通常我们是从随机过程的角度来对原始信源进行分析。对于信源符号间相互独立的无记忆信源来说，通常存在概率分布的不均匀性。此时采用统计匹配的方法进行码元分配，即大概率符号分配短码，小概率符号分配长码，可以降低平均码长，从而实现了数据压缩的目的。这种思路的典型案例就是摩尔斯电码，它是由摩尔斯在 19 世纪中期设计的。莫尔斯电码的设计是对电报发送的字符用点（·）和划（—）来表示。摩尔斯为出现频率高的字符分配了短码，如 e（·）和 a（·—）；为出现频率低的字符分配了长码，如 q（— — · —）和 j（· — — —），这样可以减少发送消息的平均时间。这也是后面统计编码所要讨论的基本思想。显然，对统计匹配编码（也称熵编码）来说，它最愿意接收的信源是：信源符号数不要太少，而且概率分布越不均匀越好。香农的信息论中有严格的证明，对信源进行无失真编码的平均码长下限就是信源熵。信源熵的通俗意义就是信源的不确定程度。对于某信源，概率分布越均匀（也就越不确定），熵越大，也越不好压缩；概率分布越不均匀，熵越小，越好压缩。（后面我们还会细致地从信息论的角度来分析这个问题，这是个很长的故事...）

像视音频这样的信源，符号之间是存在相关性的，通俗地说：就是相邻符号的取值很大概率上是颇为接近的。对这类信源符号间存在相关性的有记忆信源来说，其结构特征主要来自于相关性。相关性与原始信源的信号结构有关，可划分为一维相关（如文本、语音和音频信号）、二维相关（如图象和计算机图形）、三维相关（如视频信号）。

图 1-2 示出示例图像（Lena）水平相邻像素的二维联合直方图。 x 轴表示当前像素的取值， y 轴表示水平相邻像素的取值，对一幅图像中所有 (x,y) 值出现的次数进行统计（用 z 轴表示）。从图中可看出在该图像中水平相邻像素取值相近的概率较大。绝大多数图像都具有相似的特点，即空域存在较强的相关性。那么针对这种结构特征可以从哪里下手去进行压缩呢？

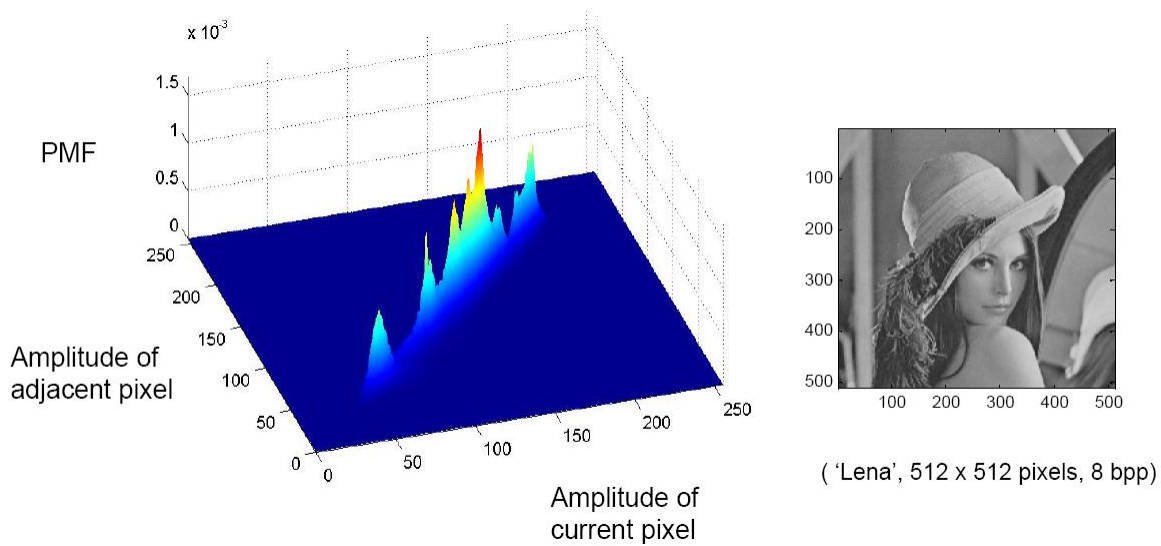


图 1-2 Lena 图像水平相邻像素的二维联合直方图

这里我们尝试采用预测的方法，即：假定我们先发送第一个值，然后发送后续值与其前一个值之间的差值（预测误差）。设某行像素值为： $s_n = 127, 128, 129, 130, 130, 132, 131, 131, 129, 128, 127$ ，使用这样的预测模型： $\hat{s}_1 = 0; \hat{s}_n = s_{n-1}$ for $n > 1$ (\hat{s}_n 表示第 n 个像素的预测值，对第一个像素值的预测值为 0，后面每个像素的预测值取前一个像素的幅值)。显然，这行像素的预测误差为 127, 1, 1, 1, 0, 2, -1, 0, -2, -1, -1。可以发现：（1）预测误差信号的幅度远远小于原始信号符号；（2）与原信号相比，预测误差的取值分布产生了很大的变化，其概率分布接近拉普拉斯分布。对图像采用其它的空间预测方法可产生相似的结果。如图 1-3 所示对某图像使用上方和左侧像素加权预测方法后，预测误差和原图像信号取值分布的对比。

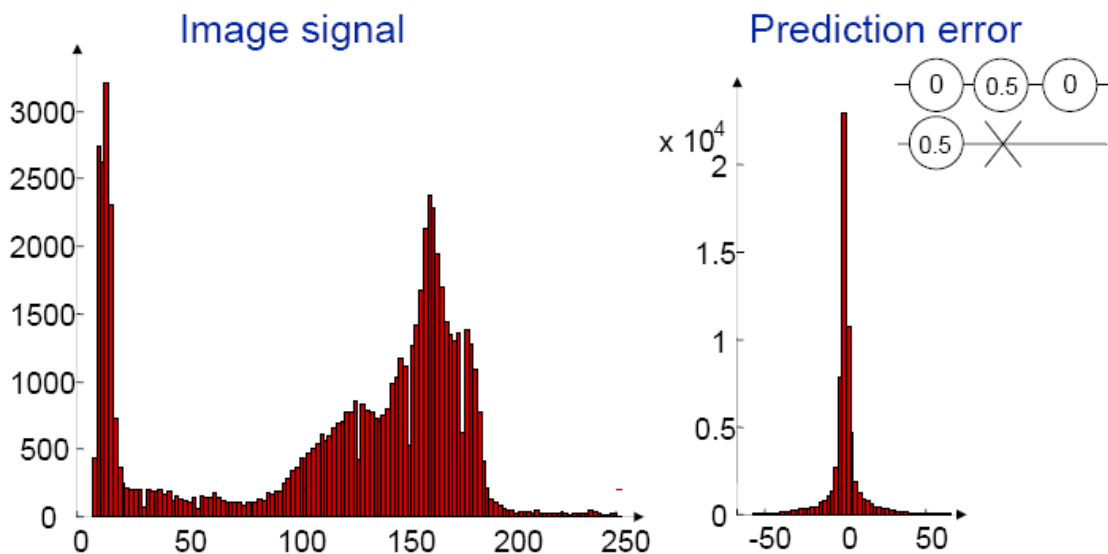


图 1-3 某图像原始图像信号与预测误差信号的概率分布对比

看到这里，我们会很自然地想到：与原始图像信号相比，预测误差信号更是熵编码的“菜”！由于预测误差信号的概率分布极不均匀，因此可以用较短的码表示预测误差为 0 附近的取值，用较长的码字表示其它预测误差较大的情况。显然，可以更有效地降低平均码长。到此，数据压缩中的杀手锏—算法组合的思想已经初露端倪了！通常我们会建立某个数学模型来重新表达规律性不那么明显（或是统计特性不太符合要求）的原始信源（在这个例子中使用了预测模型），对经过“改造”后的“新信源”符号再进行随后的处理（在这个例子中对预测误差信号使用了熵编码），将算法进行组合使用果然可以更好地压缩数据。而且很重要的是：在解码时可以完全无误地恢复出原始图像。

我们讲到的上述第一个算法组合就是**预测+熵编码**。理解其基本设计思想后，我们进一步分析其实现过程中数据及其特点。在经过预测后，预测误差信号的取值分布接近于拉普拉斯分布，但其取值范围发生了变化。以 8 比特的图像为例，原始图像信号的取值范围为 0 到 255，在进行预测处理后，预测误差信号的取值范围是-255 到 255（需要用 9 比特表示），反而增加了！但是数据压缩系统考虑的是最终的压缩性能。课程中后面的实验表明即使中间数据（预测误差信号）增加了数据量，但只要整体压缩性能达到要求，那就完全没问题。由此我们也看到数据压缩技术的整体性思想，对这个问题的理解还将延续到优化方法的设计上。如果将预测误差信号进行 8 比特量化（从 9 比特减少了 1 比特，更简洁地表达预测误差），然后再送入统计匹配编码器中进行编码，结果会是什么样？显然，量化带来了损失，而且这种损失是无法找回的。但我们后面的实验结果表明：量化带来的失真人眼并不能分辨，而且还有效地提高了压缩比。因此，量化是提高压缩效率的有效手段。

在数据压缩技术中，预测编码的思想几乎覆盖到所有信源，它出现在语音、音频、视频、图像等人们日常生活中使用的绝大多数信源的压缩方法中。只要原始信源存在着信号结构上的相关性，那么使用预测编码是理所当然的事。除了预测方法外，还有很多“改造”信源的技术（如变换等），在后面我们还会陆续学到。

2. 利用人的感知冗余进行压缩

利用信源的统计结构特征来进行压缩时只关注信源符号的统计特性，不管信源的性质（也即视频、音频、语音等等）。另一种数据压缩的思路是利用数据接收者的特性。毕竟任何多媒体信息经过压缩后毕竟是要送给用户体验（看和听）的，而人的感知能力是有限的，因此可以利用人类的感知局限，丢弃一些不必要的信息来实现数据压缩。利用这种思想来设计数据压缩系统的就是音频编码。感知音频编码技术完全围绕人的听觉特性设计，是围绕人的感知冗余来进行压缩的典型示例，所以前面一般会加“感知”两字。相对地，“感知视频编码”则不同，其核心技术仍然是以利用信源的统计结构特征为主，但在编码过程中考虑到人的视觉特性。

下面简要介绍感知音频编码技术的主要思路，详细的分析将在后续课程中给出。在算法设计时主要考虑了听觉系统对声音的三个感知特性，即响度、音高和掩蔽效应。如图 1-4 所示，人耳对响度的感知随频率变化而变化。例如：1 kHz 的 10dB（声压级）的声音和

200Hz 的 30dB（声压级）的声音，在人耳听起来具有相同的响度。人耳对不同频率的敏感程度差别很大，其中对 1kHz-4 kHz 范围的信号最为敏感，幅度很低的信号都能被人耳听到。而在低频区和高频区，能被人耳听到的信号幅度要高得多。更为重要的是人耳存在掩蔽效应这样一个事实。掩蔽效应是人听觉器官存在的缺陷。通常，频域中一个声压级较大的音可以掩蔽其它同时发生的声压级较小的音，低频纯音可以有效地掩蔽高频纯音，如图 1-5 所示。在掩蔽情况下，必须加大被掩蔽音的强度才能被人耳听到，此时的听阈称为掩蔽听阈。例如，在一个环境很喧闹的场合中说话，即使说话声音存在，但是其听阈由于掩蔽音（环境噪声）的存在而提高了。因此如果听不见这个说话声音的话，就完全没有必要去编码表示。

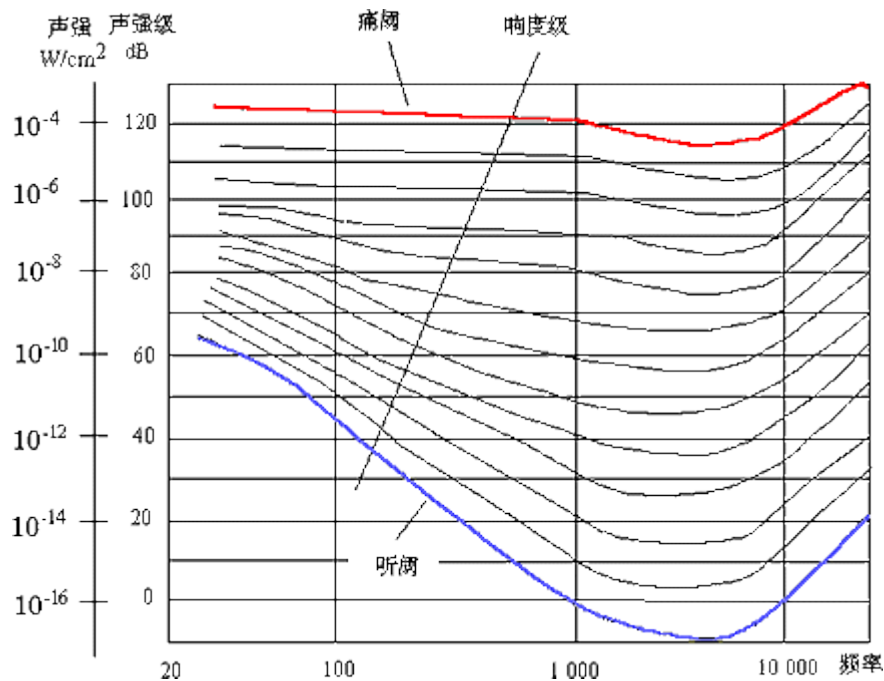


图1-4 人耳“听阈—频率”曲线

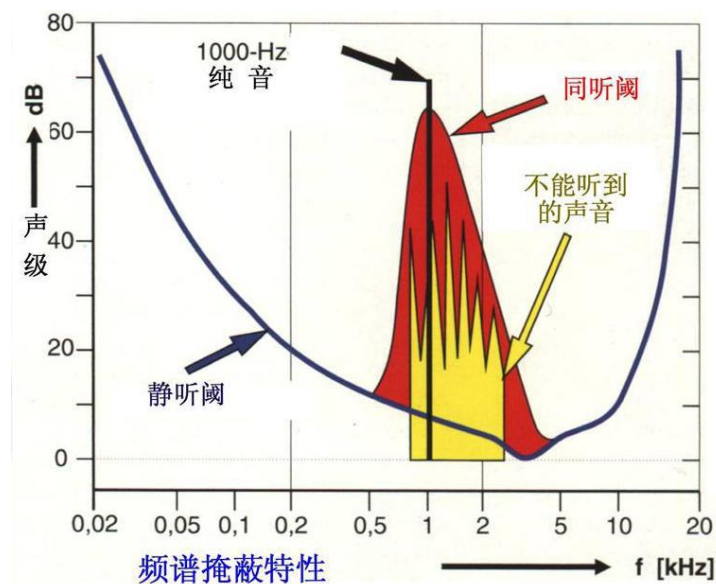


图 1-5 频域掩蔽效应

在理解上面的思路后，感知音频编码的大框架就显而易见了，如图 1-6 所示。将输入信号分割成若干个子频带信号，用心理声学模型来计算每个频带的掩蔽阈值。根据可听度来分配每个子频带系数所使用的量化比特数。重要的声音就分配多一些位数来确保可听的完整性，而对于轻言细语的编码位数就会少一些，听不到的声音就根本不进行编码，从而降低了比特速率。在解码时将这些量化后的子频带系数进行反量化，然后再进行频带合成，得到全频带信号。

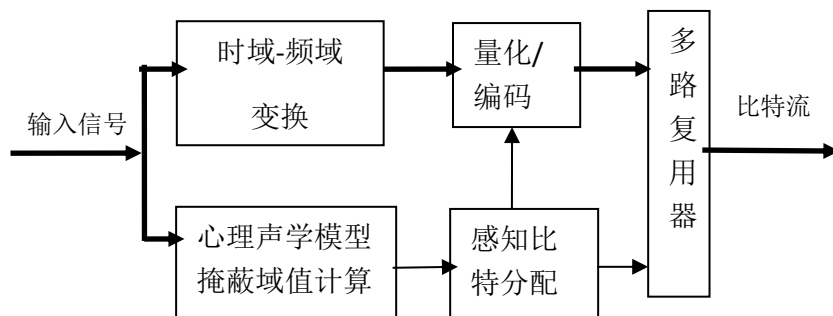


图 1-6 感知音频编码基本框架

3. 利用信源数据产生的物理过程进行压缩

第三类数据压缩的思路颇有些回归本源的意思。仔细观察视频、音频、语音、图像、动画、文字等信源，就会发现除了语音的产生以外，其余信源的产生都是没有规律可循的。唯有语音数据的物理产生过程对于所有人来说都是一致的，即：肺部的气流冲击喉咙的声带，通过舌头和口型变化发出不同的声音。如果能想办法从语音数据产生物理过程的建模入手，不直接传送语音信号本身，只传送其模型信息，那么就可能实现很大程度的数据压缩。在接收端只需要利用模型信息重建语音即可。

在这种思想的指导下，对语音信号进行编码的过程实质上是对语音信号进行分析和建模的过程。也即，将被分析的语音信号 $s(n)$ 视为某未知系统的输出，并假设它是由某序列激励未知系统（该系统的传输函数用 $V(z)$ 表示）而产生的，如图 1-7 所示。那这个未知系统到底是什么模样呢？

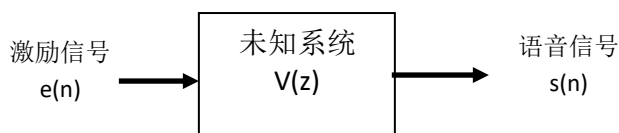


图 1-7 语音分析建模的过程

通常假设声道是一条具有时变的非均匀横截面积的声管，空气流中或声管壁上都没有热传导和粘滞损耗。声波是沿着声管的轴线传播的平面波，根据平面声波通过无损声管的模型进行理论和实验分析，可以推导出声道的等效传输函数。对于大多数语音来说是一个全极点传输函数。但对于鼻音和摩擦音则得到一个既有极点又有零点的传输函数。由于任何零点可以用多个极点来逼近，因此，常用具有时变全极点传输函数的数字滤波器来作为声道的模型。在语音编码中最著名的线性预测模型如图 1-8 所示。这是一个全极点滤波器，输入激励信号为准周期脉冲信号或随机噪声源。因此语音分析建模的过程就是在分析激励信号以及滤波器参数的过程，其本质也是用全极点模型估计语音信号的谱的谱估计问题。

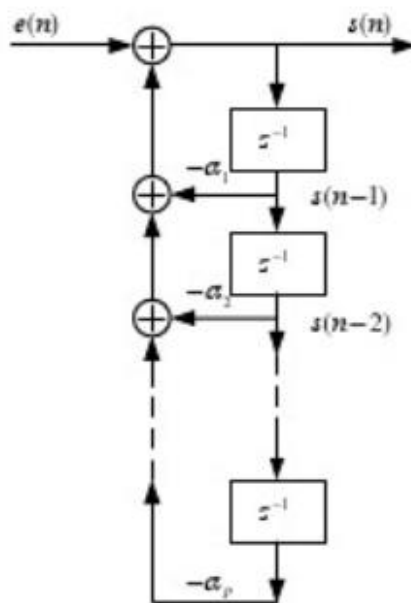


图 1-8 具有时变全极点传输参数的声道模型

上述语音编码的思路显然与前面两类方法不同，其主要区别在于传送的内容并非信源输出样本的表示，而是由发送端告诉接收端，如何重新生成这些输出。将这种分析/合成的思想扩展到图像和视频编码领域就是目前的基于模型的视频编码方法。某些特定类型

的视频（如正在讲话的头像）具有满足分析/合成方法所需要的条件，人们在讲话时，面部表情收到面部构造和物理运动规律的限制，因此就有可能基于三维模型进行面部图像序列的压缩。这是一个很有意思的研究领域。

综上所述，目前主要的数据压缩技术围绕上述三类角度（即信源数据本身的统计结构特征、信源数据使用者的感知特性、信源数据本身产生的模型）进行设计和实现。实际上，有许多种不同的方法可以描述信源数据的特征，不同的特征描述方式得到不同的压缩方案。在分析问题的过程中，我们发现数据压缩的具体方法和特定的应用有关。但从实现上来看，都经过了对信源进行**重新表达的过程（建模）**，只不过建模的方式和出发点不同。对建立模型之后的模型参数再进行量化、编码和表示。**建模—量化—熵编码**通常是数据压缩系统的一般性步骤。在具体的算法过程中，大量应用了信号处理技术。例如，在感知音频编码中首先要经过时域-频域变换，心理声学模型的计算中也需要将输入音频信号经过 **FFT** 后才能进行进一步的分析；语音编码的过程中需要进行语音信号的分析。这些都需要建立扎实的信号处理基础才能够理解得比较透彻。

对于数据压缩系统，理解其总体设计思想是非常重要的。只有站在一定的高度来观察整个系统，才可能比较深刻地理解每个（工具）模块的作用和整体优化方法，这是自顶向下去理解问题的方式。同样，也需要学习每种数据压缩工具（算法）的原理及实现方法，才有可能在面对具体的数据压缩问题时选择适当的工具并加以组合。数据压缩是一门科学，更是一门实验科学，要对一种算法有深入的理解，最好的方法就是实现（至少调试）该算法。这也是我们的课程中所重点强调的实验内容。

三、压缩系统的性能评价

之所以在这里提到压缩系统的性能评价问题，是因为在后面讨论每项技术和系统时都需要衡量其性能。目前，对于数据压缩系统的评价主要是从质量、比特率（压缩比）、编解码复杂度和通信时延四个方面来进行衡量。并非所有的应用都需要比较这四项指标，但是前两项一定是要合在一起的（通常以编码效率表示）。单独比较比特率或者质量毫无意义。对于数据压缩系统编码效率的评价可分为两方面：在规定的比特率下比较输出（解码重建）信号质量的优劣；或者在一定的输出信号质量下比较编码比特率的高低。质量包括主观质量和客观质量的评价。

对于数据压缩系统的复杂度分析通常用算法的运算量和需要的存储量来表示。系统的复杂度既包括编码器的复杂度，也包括解码器的复杂度。大多数数据压缩系统的特点是编码器复杂度大于解码器复杂度（高得多）。这是因为考虑到尽量简化数量庞大的终端（解码）设备。

由于编码算法自身的结构会导致产生时延，例如，视频编码中使用了帧间预测的方案，在使用双向预测时编码的帧序和显示的帧序不同。因此在对通信时延要求较高的应用中就需要考虑如何制定编码方案的问题。

对不同信源使用各种压缩技术及系统的性能评价将在后续课程中详细介绍。

四、课程整体内容的安排

欲先攻其事，必先利其器。数据压缩的主要对象包括文本、图形、图像、视频、音频和语音等信源和媒体数据。课程首先介绍其中重要的几种媒体特性及其常见的原始表达格式。**建模—量化—熵编码**是数据压缩系统算法组合的一般性步骤。其中，熵编码通常是最后一步，也是一个最为被动的步骤。熵编码的主要思想（统计匹配）来自于香农的信息论。因此我们先介绍与数据压缩有关的信息论基础，然后再介绍熵编码及其它的无失真信源编码方法。在绝大多数实际的数据压缩系统中，量化器的输入是信源建模后的输出，量化器的输出通常送入熵编码器。量化器是数据压缩的主要失真来源，值得好好研究和分析。熵编码和量化都分析完后，最后我们按照建模目的的不同，分别介绍预测、变换、子带、小波及相关的技术组合方案，并详细分析重要的数据压缩标准。