

Edgar H. Sibley
Panel Editor

The state of the art in data compression is arithmetic coding, not the better-known Huffman method. Arithmetic coding gives greater compression, is faster for adaptive models, and clearly separates the model from the channel encoding.

ARITHMETIC CODING FOR DATA COMPRESSION

IAN H. WITTEN, RADFORD M. NEAL, and JOHN G. CLEARY

Arithmetic coding is superior in most respects to the better-known Huffman [10] method. It represents information at least as compactly—sometimes considerably more so. Its performance is optimal without the need for blocking of input data. It encourages a clear separation between the model for representing data and the encoding of information with respect to that model. It accommodates adaptive models easily and is computationally efficient. Yet many authors and practitioners seem unaware of the technique. Indeed there is a widespread belief that Huffman coding cannot be improved upon.

We aim to rectify this situation by presenting an accessible implementation of arithmetic coding and by detailing its performance characteristics. We start by briefly reviewing basic concepts of data compression and introducing the model-based approach that underlies most modern techniques. We then outline the idea of arithmetic coding using a simple example, before presenting programs for both encoding and decoding. In these programs the model occupies a separate module so that different models can easily be used. Next we discuss the construction of fixed and adaptive models and detail the compression efficiency and execution time of the programs, including the effect of different arithmetic word lengths on compression efficiency. Finally, we outline a few applications where arithmetic coding is appropriate.

Financial support for this work has been provided by the Natural Sciences and Engineering Research Council of Canada.

UNIX is a registered trademark of AT&T Bell Laboratories.

© 1987 ACM 0001-0782/87/0600-0520 75¢

DATA COMPRESSION

To many, data compression conjures up an assortment of ad hoc techniques such as conversion of spaces in text to tabs, creation of special codes for common words, or run-length coding of picture data (e.g., see [8]). This contrasts with the more modern model-based paradigm for coding, where, from an *input string* of symbols and a *model*, an *encoded string* is produced that is (usually) a compressed version of the input. The decoder, which must have access to the same model, regenerates the exact input string from the encoded string. Input symbols are drawn from some well-defined set such as the ASCII or binary alphabets; the encoded string is a plain sequence of bits. The model is a way of calculating, in any given context, the distribution of probabilities for the next input symbol. It must be possible for the decoder to produce exactly the same probability distribution in the same context. Compression is achieved by transmitting the more probable symbols in fewer bits than the less probable ones.

For example, the model may assign a predetermined probability to each symbol in the ASCII alphabet. No context is involved. These probabilities can be determined by counting frequencies in representative samples of text to be transmitted. Such a *fixed* model is communicated in advance to both encoder and decoder, after which it is used for many messages.

Alternatively, the probabilities that an *adaptive* model assigns may change as each symbol is transmitted, based on the symbol frequencies seen so far in the message. There is no need for a representative

sample of text, because each message is treated as an independent unit, starting from scratch. The encoder's model changes with each symbol transmitted, and the decoder's changes with each symbol received, in sympathy.

More complex models can provide more accurate probabilistic predictions and hence achieve greater compression. For example, several characters of previous context could condition the next-symbol probability. Such methods have enabled mixed-case English text to be encoded in around 2.2 bits/character with two quite different kinds of model [4, 6]. Techniques that do not separate modeling from coding so distinctly, like that of Ziv and Lempel [23], do not seem to show such great potential for compression, although they may be appropriate when the aim is raw speed rather than compression performance [22].

The effectiveness of any model can be measured by the entropy of the message with respect to it, usually expressed in bits/symbol. Shannon's fundamental theorem of coding states that, given messages randomly generated from a model, it is impossible to encode them into less bits (on average) than the entropy of that model [21].

A message can be coded with respect to a model using either Huffman or arithmetic coding. The former method is frequently advocated as the best possible technique for reducing the encoded data rate. It is not. Given that each symbol in the alphabet must translate into an integral number of bits in the encoding, Huffman coding indeed achieves "minimum redundancy." In other words, it performs optimally if all symbol probabilities are integral powers of $\frac{1}{2}$. But this is not normally the case in practice; indeed, Huffman coding can take up to one extra bit per symbol. The worst case is realized by a source in which one symbol has probability approaching unity. Symbols emanating from such a source convey negligible information on average, but require at least one bit to transmit [7]. Arithmetic coding dispenses with the restriction that each symbol must translate into an integral number of bits, thereby coding more efficiently. It actually achieves the theoretical entropy bound to compression efficiency for any source, including the one just mentioned.

In general, sophisticated models expose the deficiencies of Huffman coding more starkly than simple ones. This is because they more often predict symbols with probabilities close to one, the worst case for Huffman coding. For example, the techniques mentioned above that code English text in 2.2 bits/character both use arithmetic coding as the final step, and performance would be impacted severely

if Huffman coding were substituted. Nevertheless, since our topic is coding and not modeling, the illustrations in this article all employ simple models. Even so, as we shall see, Huffman coding is inferior to arithmetic coding.

The basic concept of arithmetic coding can be traced back to Elias in the early 1960s (see [1, pp. 61–62]). Practical techniques were first introduced by Rissanen [16] and Pasco [15], and developed further by Rissanen [17]. Details of the implementation presented here have not appeared in the literature before; Rubin [20] is closest to our approach. The reader interested in the broader class of arithmetic codes is referred to [18]; a tutorial is available in [13]. Despite these publications, the method is not widely known. A number of recent books and papers on data compression mention it only in passing, or not at all.

THE IDEA OF ARITHMETIC CODING

In arithmetic coding, a message is represented by an interval of real numbers between 0 and 1. As the message becomes longer, the interval needed to represent it becomes smaller, and the number of bits needed to specify that interval grows. Successive symbols of the message reduce the size of the interval in accordance with the symbol probabilities generated by the model. The more likely symbols reduce the range by less than the unlikely symbols and hence add fewer bits to the message.

Before anything is transmitted, the range for the message is the entire interval $[0, 1)$, denoting the half-open interval $0 \leq x < 1$. As each symbol is processed, the range is narrowed to that portion of it allocated to the symbol. For example, suppose the alphabet is $\{a, e, i, o, u, !\}$, and a fixed model is used with probabilities shown in Table I. Imagine trans-

TABLE I. Example Fixed Model for Alphabet $\{a, e, i, o, u, !\}$

Symbol	Probability	Range
<i>a</i>	.2	[0, 0.2)
<i>e</i>	.3	[0.2, 0.5)
<i>i</i>	.1	[0.5, 0.6)
<i>o</i>	.2	[0.6, 0.8)
<i>u</i>	.1	[0.8, 0.9)
!	.1	[0.9, 1.0)

mitting the message *eaii!*. Initially, both encoder and decoder know that the range is $[0, 1)$. After seeing the first symbol, *e*, the encoder narrows it to $[0.2, 0.5)$, the range the model allocates to this symbol. The second symbol, *a*, will narrow this new range to the first one-fifth of it, since *a* has been

allocated $[0, 0.2]$. This produces $[0.2, 0.26]$, since the previous range was 0.3 units long and one-fifth of that is 0.06. The next symbol, i , is allocated $[0.5, 0.6]$, which when applied to $[0.2, 0.26]$ gives the smaller range $[0.23, 0.236]$. Proceeding in this way, the encoded message builds up as follows:

Initially	$[0, 1]$
After seeing e	$[0.2, 0.5]$
a	$[0.2, 0.26]$
i	$[0.23, 0.236]$
i	$[0.233, 0.2336]$
!	$[0.23354, 0.2336]$

Figure 1 shows another representation of the encoding process. The vertical bars with ticks represent the symbol probabilities stipulated by the model. After the first symbol has been processed, the model is scaled into the range $[0.2, 0.5]$, as shown in

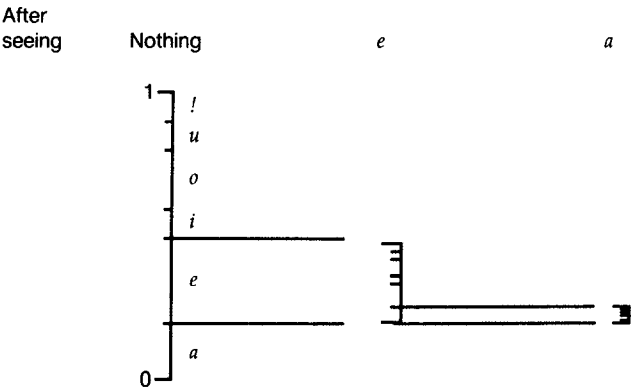


FIGURE 1a. Representation of the Arithmetic Coding Process

Figure 1a. The second symbol scales it again into the range $[0.2, 0.26]$. But the picture cannot be continued in this way without a magnifying glass! Consequently, Figure 1b shows the ranges expanded to full height at every stage and marked with a scale that gives the endpoints as numbers.

Suppose all the decoder knows about the message is the final range, $[0.23354, 0.2336]$. It can immediately deduce that the first character was e , since the range lies entirely within the space the model of Table I allocates for e . Now it can simulate the operation of the encoder:

Initially	$[0, 1]$
After seeing e	$[0.2, 0.5]$

This makes it clear that the second character is a , since this will produce the range

After seeing a $[0.2, 0.26]$,

which entirely encloses the given range $[0.23354, 0.2336]$. Proceeding like this, the decoder can identify the whole message.

It is not really necessary for the decoder to know both ends of the range produced by the encoder. Instead, a single number within the range—for example, 0.23355—will suffice. (Other numbers, like 0.23354, 0.23357, or even 0.23354321, would do just as well.) However, the decoder will face the problem of detecting the end of the message, to determine when to stop decoding. After all, the single number 0.0 could represent any of $a, aa, aaa, aaaa, \dots$. To resolve the ambiguity, we ensure that each message ends with a special EOF symbol known to both encoder and decoder. For the alphabet of Table I, ! will be used to terminate messages, and only to termi-

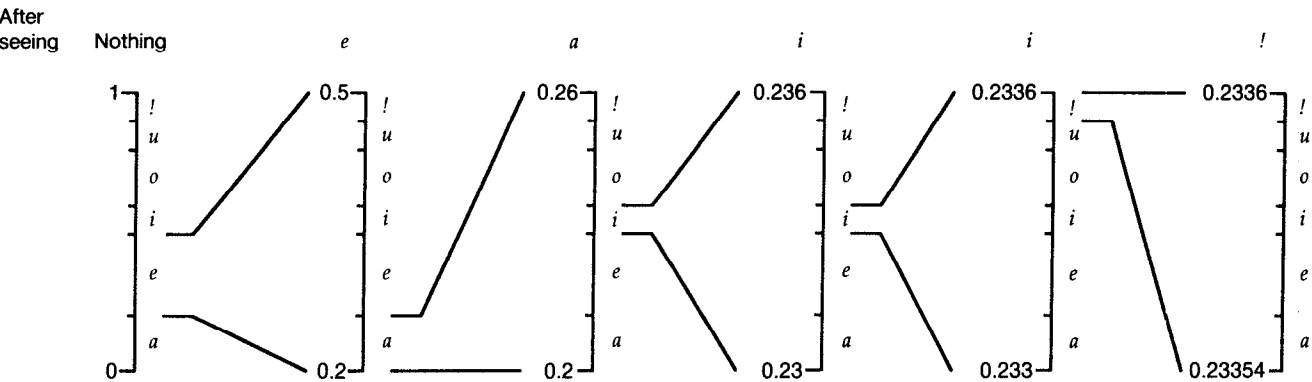


FIGURE 1b. Representation of the Arithmetic Coding Process with the Interval Scaled Up at Each Stage

```

/* ARITHMETIC ENCODING ALGORITHM. */

/* Call encode_symbol repeatedly for each symbol in the message.          */
/* Ensure that a distinguished "terminator" symbol is encoded last, then  */
/* transmit any value in the range [low, high).                            */

encode_symbol(symbol, cum_freq)
    range = high - low
    high = low + range*cum_freq[symbol-1]
    low = low + range*cum_freq[symbol]

/* ARITHMETIC DECODING ALGORITHM. */

/* "Value" is the number that has been received.                          */
/* Continue calling decode_symbol until the terminator symbol is returned. */

decode_symbol(cum_freq)
    find symbol such that
        cum_freq[symbol] <= (value-low)/(high-low) < cum_freq[symbol-1]
        /* This ensures that value lies within the new */
        /* [low, high) range that will be calculated by */
        /* the following lines of code.                  */

    range = high - low
    high = low + range*cum_freq[symbol-1]
    low = low + range*cum_freq[symbol]
    return symbol

```

FIGURE 2. Pseudocode for the Encoding and Decoding Procedures

nate messages. When the decoder sees this symbol, it stops decoding.

Relative to the fixed model of Table I, the entropy of the five-symbol message *eaii!* is

$$\begin{aligned}
 & -\log 0.3 - \log 0.2 - \log 0.1 - \log 0.1 - \log 0.1 \\
 & = -\log 0.00006 \approx 4.22
 \end{aligned}$$

(using base 10, since the above encoding was performed in decimal). This explains why it takes five decimal digits to encode the message. In fact, the size of the final range is $0.2336 - 0.23354 = 0.00006$, and the entropy is the negative logarithm of this figure. Of course, we normally work in binary, transmitting binary digits and measuring entropy in bits.

Five decimal digits seems a lot to encode a message comprising four vowels! It is perhaps unfortunate that our example ended up by expanding rather than compressing. Needless to say, however, different models will give different entropies. The best single-character model of the message *eaii!* is the set of symbol frequencies $\{e(0.2), a(0.2), i(0.4), !(0.2)\}$, which gives an entropy of 2.89 decimal digits. Using this model the encoding would be only three digits long. Moreover, as noted earlier, more sophisticated models give much better performance in general.

A PROGRAM FOR ARITHMETIC CODING

Figure 2 shows a pseudocode fragment that summarizes the encoding and decoding procedures developed in the last section. Symbols are numbered, 1, 2, 3, The frequency range for the i th symbol is from $\text{cum_freq}[i]$ to $\text{cum_freq}[i - 1]$. As i decreases, $\text{cum_freq}[i]$ increases, and $\text{cum_freq}[0] = 1$. (The reason for this “backwards” convention is that $\text{cum_freq}[0]$ will later contain a normalizing factor, and it will be convenient to have it begin the array.) The “current interval” is $[low, high)$, and for both encoding and decoding, this should be initialized to $[0, 1)$.

Unfortunately, Figure 2 is overly simplistic. In practice, there are several factors that complicate both encoding and decoding:

Incremental transmission and reception. The encode algorithm as described does not transmit anything until the entire message has been encoded; neither does the decode algorithm begin decoding until it has received the complete transmission. In most applications an incremental mode of operation is necessary.

The desire to use integer arithmetic. The precision required to represent the $[low, high)$ interval grows with the length of the message. Incremental operation will help overcome this, but the potential for

overflow and underflow must still be examined carefully.

Representing the model so that it can be consulted efficiently. The representation used for the model should minimize the time required for the decode algorithm to identify the next symbol. Moreover, an adaptive model should be organized to minimize the time-consuming task of maintaining cumulative frequencies.

Figure 3 shows working code, in C, for arithmetic encoding and decoding. It is considerably more detailed than the bare-bones sketch of Figure 2! Implementations of two different models are given in Figure 4; the Figure 3 code can use either one.

The remainder of this section examines the code of Figure 3 more closely, and includes a proof that decoding is still correct in the integer implementation and a review of constraints on word lengths in the program.

```

arithmetic_coding.h


---


1  /* DECLARATIONS USED FOR ARITHMETIC ENCODING AND DECODING */
2
3
4  /* SIZE OF ARITHMETIC CODE VALUES. */
5
6  #define Code_value_bits 16          /* Number of bits in a code value */
7  typedef long code_value;          /* Type of an arithmetic code value */
8
9  #define Top_value (((long)1<<Code_value_bits)-1) /* Largest code value */
10
11
12 /* HALF AND QUARTER POINTS IN THE CODE VALUE RANGE. */
13
14 #define First_qtr (Top_value/4+1) /* Point after first quarter */
15 #define Half (2*First_qtr) /* Point after first half */
16 #define Third_qtr (3*First_qtr) /* Point after third quarter */
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
25
```

```

encode.c


---


39  /* MAIN PROGRAM FOR ENCODING. */
40
41  #include <stdio.h>
42  #include "model.h"
43
44  main()
45  {
46      start_model();                /* Set up other modules. */
47      start_outputting_bits();
48      start_encoding();
49      for (;;) {                    /* Loop through characters. */
50          int ch; int symbol;
51          ch = getc(stdin);          /* Read the next character. */
52          if (ch==EOF) break;        /* Exit loop on end-of-file. */
53          symbol = char_to_index[ch]; /* Translate to an index. */
54          encode_symbol(symbol,cum_freq); /* Encode that symbol. */
55          update_model(symbol);      /* Update the model. */
56      }
57      encode_symbol(EOF_symbol,cum_freq); /* Encode the EOF symbol. */
58      done_encoding();              /* Send the last few bits. */
59      done_outputting_bits();
60      exit(0);
61  }

```



```

arithmetic_encode.c


---


61  /* ARITHMETIC ENCODING ALGORITHM. */
62
63  #include "arithmetic_coding.h"
64
65  static void bit_plus_follow(); /* Routine that follows */
66
67
68  /* CURRENT STATE OF THE ENCODING. */
69
70  static code_value low, high; /* Ends of the current code region */
71  static long bits_to_follow; /* Number of opposite bits to output after */
72                             /* the next bit. */
73
74
75  /* START ENCODING A STREAM OF SYMBOLS. */
76
77  start_encoding()
78  {
79      low = 0; /* Full code range. */
80      high = Top_value;
81      bits_to_follow = 0; /* No bits to follow next. */
82  }
83
84  /* ENCODE A SYMBOL. */
85
86  encode_symbol(symbol,cum_freq)
87  {
88      int symbol; /* Symbol to encode */
89      int cum_freq[]; /* Cumulative symbol frequencies */
90      long range; /* Size of the current code region */
91      range = (long) (high-low)+1;
92      high = low + (range*cum_freq[symbol-1])/cum_freq[0]-1; /* Narrow the code region */
93      low = low + (range*cum_freq[symbol])/cum_freq[0]; /* to that allotted to this */
94                                                         /* symbol. */

```

FIGURE 3. C Implementation of Arithmetic Encoding and Decoding (*continued*)

```

95     for (;;) {                                     /* Loop to output bits. */
96         if (high<Half) {
97             bit_plus_follow(0);                     /* Output 0 if in low half. */
98         }
99         else if (low>=Half) {                       /* Output 1 if in high half.*/
100             bit_plus_follow(1);
101             low -= Half;
102             high -= Half;                           /* Subtract offset to top. */
103         }
104         else if (low>=First_qtr                     /* Output an opposite bit */
105                 && high<Third_qtr) {               /* later if in middle half. */
106             bits_to_follow += 1;
107             low -= First_qtr;                       /* Subtract offset to middle*/
108             high -= First_qtr;
109         }
110         else break;                                /* Otherwise exit loop. */
111         low = 2*low;
112         high = 2*high+1;                            /* Scale up code range. */
113     }
114 }
115
116
117 /* FINISH ENCODING THE STREAM. */
118
119 done_encoding()
120 {
121     bits_to_follow += 1;                            /* Output two bits that */
122     if (low<First_qtr) bit_plus_follow(0);          /* select the quarter that */
123     else bit_plus_follow(1);                        /* the current code range */
124 }                                                    /* contains. */
125
126 /* OUTPUT BITS PLUS FOLLOWING OPPOSITE BITS. */
127
128 static void bit_plus_follow(bit)
129     int bit;
130 {
131     output_bit(bit);                                /* Output the bit. */
132     while (bits_to_follow>0) {
133         output_bit(!bit);                          /* Output bits_to_follow */
134         bits_to_follow -- 1;                        /* opposite bits. Set */
135     }                                                /* bits_to_follow to zero. */
136 }
137
138 decode.c
139
140
141 /* MAIN PROGRAM FOR DECODING. */
142
143 #include <stdio.h>
144 #include "model.h"
145
146 main()
147 {
148     start_model();                                  /* Set up other modules. */
149     start_inputting_bits();
150     start_decoding();
151     for (;;) {
152         int ch; int symbol;
153         symbol = decode_symbol(cum_freq);           /* Decode next symbol. */
154         if (symbol==EOF_symbol) break;              /* Exit loop if EOF symbol. */
155         ch = index_to_char[symbol];                /* Translate to a character.*/
156         putc(ch,stdout);                           /* Write that character. */
157         update_model(symbol);                       /* Update the model. */
158     }
159     exit(0);
160 }

```

FIGURE 3. C Implementation of Arithmetic Encoding and Decoding (continued)

```

arithmic_decode.c

155  /* ARITHMETIC DECODING ALGORITHM. */
156
157  #include "arithmic_coding.h"
158
159
160  /* CURRENT STATE OF THE DECODING. */
161
162  static code_value value;      /* Currently-seen code value      */
163  static code_value low, high;  /* Ends of current code region */
164
165
166  /* START DECODING A STREAM OF SYMBOLS. */
167
168  start_decoding()
169  {
170      int i;
171      value = 0;
172      for (i = 1; i <= Code_value_bits; i++) {
173          value = 2*value+input_bit();
174      }
175      low = 0;
176      high = Top_value;
177  }
178
179  /* DECODE THE NEXT SYMBOL. */
180
181  int decode_symbol(cum_freq)
182  {
183      int cum_freq[];
184      long range;
185      int cum;
186      int symbol;
187      range = (long) (high-low)+1;
188      cum = 0;
189      for (symbol = 1; cum_freq[symbol]>cum; symbol++) {
190          cum = cum + cum_freq[symbol-1];
191          high = low + (range*cum_freq[symbol-1])/cum_freq[0]-1;
192          low = low + (range*cum_freq[symbol])/cum_freq[0];
193      }
194      for (;;) {
195          if (high<Half) {
196              /* nothing */
197          }
198          else if (low>=Half) {
199              value -= Half;
200              low -= Half;
201              high -= Half;
202          }
203          else if (low>=First_qtr
204                  && high<Third_qtr) {
205              value -= First_qtr;
206              low -= First_qtr;
207              high -= First_qtr;
208          }
209          else break;
210          low = 2*low;
211          high = 2*high+1;
212          value = 2*value+input_bit();
213      }
214      return symbol;
215  }

```

FIGURE 3. C Implementation of Arithmetic Encoding and Decoding (continued)


```

bit_input.c


---


216  /* BIT INPUT ROUTINES. */
217
218  #include <stdio.h>
219  #include "arithmetic_coding.h"
220
221
222  /* THE BIT BUFFER. */
223
224  static int buffer;          /* Bits waiting to be input          */
225  static int bits_to_go;      /* Number of bits still in buffer    */
226  static int garbage_bits;    /* Number of bits past end-of-file   */
227
228
229  /* INITIALIZE BIT INPUT. */
230
231  start_inputting_bits()
232  {
233      bits_to_go = 0;          /* Buffer starts out with          */
234      garbage_bits = 0;        /* no bits in it.                  */
235  }
236
237  /* INPUT A BIT. */
238
239  int input_bit()
240  {
241      int t;
242      if (bits_to_go==0) {      /* Read the next byte if no      */
243          buffer = getc(stdin); /* bits are left in buffer.      */
244          if (buffer==EOF) {
245              garbage_bits += 1; /* Return arbitrary bits*/
246              if (garbage_bits>Code_value_bits-2) { /* after eof, but check */
247                  fprintf(stderr, "Bad input file\n"); /* for too many such. */
248                  exit(-1);
249              }
250              bits_to_go = 8;
251          }
252          t = buffer&1;          /* Return the next bit from      */
253          buffer >>= 1;          /* the bottom of the byte.      */
254          bits_to_go -- 1;
255          return t;
256  }

```

FIGURE 3. C Implementation of Arithmetic Encoding and Decoding (*continued*)

```

bit_output.c
-----
257  /* BIT OUTPUT ROUTINES. */
258
259  #include <stdio.h>
260
261
262  /* THE BIT BUFFER. */
263
264  static int buffer;          /* Bits buffered for output          */
265  static int bits_to_go;     /* Number of bits free in buffer */
266
267
268  /* INITIALIZE FOR BIT OUTPUT. */
269
270  start_outputting_bits()
271  {
272      buffer = 0;              /* Buffer is empty to start */
273      bits_to_go = 8;         /* with.                    */
274  }
275
276  /* OUTPUT A BIT. */
277
278  output_bit(bit)
279  {
280      int bit;
281      buffer >>= 1;           /* Put bit in top of buffer.*/
282      if (bit) buffer |= 0x80;
283      bits_to_go --;
284      if (bits_to_go == 0) {   /* Output buffer if it is   */
285          putchar(buffer, stdout); /* now full.                */
286          bits_to_go = 8;
287      }
288  }
289
290  /* FLUSH OUT THE LAST BITS. */
291
292  done_outputting_bits()
293  {
294      putchar(buffer >> bits_to_go, stdout);
295  }

```

FIGURE 3. C Implementation of Arithmetic Encoding and Decoding (*continued*)

```

fixed_model.c

1  /* THE FIXED SOURCE MODEL */
2
3  #include "model.h"
4
5  int freq[No_of_symbols+1] = {
6      0,
7      1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 124, 1, 1, 1, 1, 1,
8      1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
9
10     /*      !      "      #      $      %      &      '      (      )      *      +      ,      -      .      /      */
11     1236, 1, 21, 9, 3, 1, 25, 15, 2, 2, 2, 1, 79, 19, 60, 1,
12
13     /* 0      1      2      3      4      5      6      7      8      9      :      ;      <      =      >      ?      */
14     15, 15, 8, 5, 4, 7, 5, 4, 4, 6, 3, 2, 1, 1, 1, 1,
15
16     /* @      A      B      C      D      E      F      G      H      I      J      K      L      M      N      O      */
17     1, 24, 15, 22, 12, 15, 10, 9, 16, 16, 8, 6, 12, 23, 13, 11,
18
19     /* P      Q      R      S      T      U      V      W      X      Y      Z      [      \      ]      ^      _      */
20     14, 1, 14, 28, 29, 6, 3, 11, 1, 3, 1, 1, 1, 1, 1, 3,
21
22     /* `      a      b      c      d      e      f      g      h      i      j      k      l      m      n      o      */
23     1, 491, 85, 173, 232, 744, 127, 110, 293, 418, 6, 39, 250, 139, 429, 446,
24
25     /* p      q      r      s      t      u      v      w      x      y      z      {      |      }      ~      *      */
26     111, 5, 388, 375, 531, 152, 57, 97, 12, 101, 5, 2, 1, 2, 3, 1,
27
28     1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
29     1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
30     1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
31     1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
32     1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
33     1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
34     1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
35     1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
36     1
37 };
38
39
40 /* INITIALIZE THE MODEL. */
41
42 start_model()
43 {
44     int i;
45     for (i = 0; i < No_of_chars; i++) {
46         char_to_index[i] = i+1;
47         index_to_char[i+1] = i;
48     }
49     cum_freq[No_of_symbols] = 0;
50     for (i = No_of_symbols; i > 0; i--) {
51         cum_freq[i-1] = cum_freq[i] + freq[i];
52     }
53     if (cum_freq[0] > Max_frequency) abort();
54 }
55
56 /* UPDATE THE MODEL TO ACCOUNT FOR A NEW SYMBOL. */
57
58 update_model(symbol)
59     int symbol;
60 {
61     /* Do nothing. */

```

FIGURE 4. Fixed and Adaptive Models for Use with Figure 3

```

adaptive_model.c

1  /* THE ADAPTIVE SOURCE MODEL */
2
3  #include "model.h"
4
5  int freq[No_of_symbols+1];      /* Symbol frequencies          */
6
7
8  /* INITIALIZE THE MODEL. */
9
10 start_model()
11 {
12     int i;
13     for (i = 0; i < No_of_chars; i++) {
14         char_to_index[i] = i+1;
15         index_to_char[i+1] = i;
16     }
17     for (i = 0; i < No_of_symbols; i++) {
18         freq[i] = 1;
19         cum_freq[i] = No_of_symbols-i;
20     }
21     freq[0] = 0;
22 }
23
24 /* UPDATE THE MODEL TO ACCOUNT FOR A NEW SYMBOL. */
25
26 update_model(symbol)
27 {
28     int symbol;          /* Index of new symbol          */
29     int i;               /* New index for symbol         */
30     if (cum_freq[0] == Max_frequency) {
31         int cum;
32         cum = 0;
33         for (i = No_of_symbols; i >= 0; i--) {
34             freq[i] = (freq[i]+1)/2;
35             cum_freq[i] = cum;
36             cum += freq[i];
37         }
38     }
39     for (i = symbol; freq[i] == freq[i-1]; i--) ; /* Find symbol's new index. */
40     if (i < symbol) {
41         int ch_i, ch_symbol;
42         ch_i = index_to_char[i];
43         ch_symbol = index_to_char[symbol];
44         index_to_char[i] = ch_symbol;
45         index_to_char[symbol] = ch_i;
46         char_to_index[ch_i] = symbol;
47         char_to_index[ch_symbol] = i;
48     }
49     freq[i] += 1;
50     while (i > 0) {
51         i -= 1;
52         cum_freq[i] += 1;
53     }
54 }

```

FIGURE 4. Fixed and Adaptive Models for Use with Figure 3 (continued)

Representing the Model

Implementations of models are discussed in the next section; here we are concerned only with the interface to the model (lines 20–38). In C, a byte is represented as an integer between 0 and 255 (a *char*). Internally, we represent a byte as an integer between 1 and 257 inclusive (an *index*), EOF being treated as a 257th symbol. It is advantageous to sort the model into frequency order, so as to minimize the number of executions of the decoding loop (line 189). To permit such reordering, the *char/index* translation is implemented as a pair of tables, *index_to_char*[] and *char_to_index* []. In one of our models, these tables simply form the *index* by adding 1 to the *char*, but another implements a more complex translation that assigns small indexes to frequently used symbols.

The probabilities in the model are represented as integer frequency counts, and cumulative counts are stored in the array *cum_freq* []. As previously, this array is “backwards,” and the total frequency count, which is used to normalize all frequencies, appears in *cum_freq*[0]. Cumulative counts must not exceed a predetermined maximum, *Max_frequency*, and the model implementation must prevent overflow by scaling appropriately. It must also ensure that neighboring values in the *cum_freq* [] array differ by at least 1; otherwise the affected symbol cannot be transmitted.

Incremental Transmission and Reception

Unlike Figure 2 the program in Figure 3 represents *low* and *high* as integers. A special data type, *code_value*, is defined for these quantities, together with some useful constants: *Top_value*, representing the largest possible *code_value*, and *First_qtr*, *Half*, and *Third_qtr*, representing parts of the range (lines 6–16). Whereas in Figure 2 the current interval is represented by [*low*, *high*], in Figure 3 it is [*low*, *high*]; that is, the range now includes the value of *high*. Actually, it is more accurate (though more confusing) to say that, in the program in Figure 3, the interval represented is [*low*, *high* + 0.11111 ...]. This is because when the bounds are scaled up to increase the precision, zeros are shifted into the low-order bits of *low*, but ones are shifted into *high*. Although it is possible to write the program to use a different convention, this one has some advantages in simplifying the code.

As the code range narrows, the top bits of *low* and *high* become the same. Any bits that are the same can be transmitted immediately, since they cannot be affected by future narrowing. For encoding, since we know that $low \leq high$, this requires code like

```
for (;;) {
    if (high < Half) {
        output_bit(0);
        low = 2*low;
        high = 2*high+1;
    }
    else if (low ≥ Half) {
        output_bit(1);
        low = 2*(low-Half);
        high = 2*(high-Half)+1;
    }
    else break;
}
```

which ensures that, upon completion, $low < Half \leq high$. This can be found in lines 95–113 of *encode_symbol()*, although there are some extra complications caused by underflow possibilities (see the next subsection). Care is taken to shift ones in at the bottom when *high* is scaled, as noted above.

Incremental reception is done using a number called *value* as in Figure 2, in which processed bits flow out the top (high-significance) end and newly received ones flow in the bottom. Initially, *start_decoding()* (lines 168–176) fills *value* with received bits. Once *decode_symbol()* has identified the next input symbol, it shifts out now-useless high-order bits that are the same in *low* and *high*, shifting *value* by the same amount (and replacing lost bits by fresh input bits at the bottom end):

```
for (;;) {
    if (high < Half) {
        value = 2*value+input_bit();
        low = 2*low;
        high = 2*high+1;
    }
    else if (low > Half) {
        value = 2*(value-Half)+input_bit();
        low = 2*(low-Half);
        high = 2*(high-Half)+1;
    }
    else break;
}
```

(see lines 194–213, again complicated by precautions against underflow, as discussed below).

Proof of Decoding Correctness

At this point it is worth checking that identification of the next symbol by *decode_symbol()* works properly. Recall from Figure 2 that *decode_symbol()* must use *value* to find the symbol that, when encoded, reduces the range to one that still includes *value*. Lines 186–188 in *decode_symbol()* identify the symbol for which

$cum_freq[symbol]$

$$\leq \left\lfloor \frac{(value - low + 1) * cum_freq[0] - 1}{high - low + 1} \right\rfloor$$

$$< cum_freq[symbol - 1],$$

where $\lfloor \cdot \rfloor$ denotes the “integer part of” function that comes from integer division with truncation. It is shown in the Appendix that this implies

$$low + \left\lfloor \frac{(high - low + 1) * cum_freq[symbol]}{cum_freq[0]} \right\rfloor$$

$$\leq v \leq low$$

$$+ \left\lfloor \frac{(high - low + 1) * cum_freq[symbol - 1]}{cum_freq[0]} \right\rfloor - 1,$$

so that $value$ lies within the new interval that $decode_symbol()$ calculates in lines 190–193. This is sufficient to guarantee that the decoding operation identifies each symbol correctly.

Underflow

As Figure 1 shows, arithmetic coding works by scaling the cumulative probabilities given by the model into the interval $[low, high]$ for each character transmitted. Suppose low and $high$ are very close together—so close that this scaling operation maps some different symbols of the model onto the same integer in the $[low, high]$ interval. This would be disastrous, because if such a symbol actually occurred it would not be possible to continue encoding. Consequently, the encoder must guarantee that the interval $[low, high]$ is always large enough to prevent this. The simplest way to do this is to ensure that this interval is at least as large as $Max_frequency$, the maximum allowed cumulative frequency count (line 36).

How could this condition be violated? The bit-shifting operation explained above ensures that low and $high$ can only become close together when they straddle $Half$. Suppose in fact they become as close as

$$First_qtr \leq low < Half \leq high < Third_qtr.$$

Then the next two bits sent will have opposite polarity, either 01 or 10. For example, if the next bit turns out to be zero (i.e., $high$ descends below $Half$ and $[0, Half]$ is expanded to the full interval), the bit after that will be one, since the range has to be above the midpoint of the expanded interval. Conversely, if the next bit happens to be one, the one after that will be zero. Therefore the interval can safely be expanded right now, if only we remember

that, whatever bit actually comes next, its opposite must be transmitted afterwards as well. Thus lines 104–109 expand $[First_qtr, Third_qtr]$ into the whole interval, remembering in $bits_to_follow$ that the bit that is output next must be followed by an opposite bit. This explains why all output is done via $bit_plus_follow()$ (lines 128–135), instead of directly with $output_bit()$.

But what if, after this operation, it is *still* true that

$$First_qtr \leq low < Half \leq high < Third_qtr?$$

Figure 5 illustrates this situation, where the current $[low, high]$ range (shown as a thick line) has been expanded a total of three times. Suppose the next bit turns out to be zero, as indicated by the arrow in Figure 5a being below the halfway point. Then the next three bits will be ones, since the arrow is not only in the top half of the bottom half of the original range, but in the top quarter, and moreover the top eighth, of that half—this is why the expansion can occur three times. Similarly, as Figure 5b shows, if the next bit turns out to be a one, it will be followed by three zeros. Consequently, we need only count the number of expansions and follow the next bit by that number of opposites (lines 106 and 131–134).

Using this technique the encoder can guarantee that, after the shifting operations, either

$$low < First_qtr < Half \leq high \quad (1a)$$

or

$$low < Half < Third_qtr \leq high. \quad (1b)$$

Therefore, as long as the integer range spanned by the cumulative frequencies fits into a quarter of that provided by $code_values$, the underflow problem cannot occur. This corresponds to the condition

$$Max_frequency \leq \frac{Top_value + 1}{4} + 1,$$

which is satisfied by Figure 3, since $Max_frequency = 2^{14} - 1$ and $Top_value = 2^{16} - 1$ (lines 36, 9). More than 14 bits cannot be used to represent cumulative frequency counts without increasing the number of bits allocated to $code_values$.

We have discussed underflow in the encoder only. Since the decoder's job, once each symbol has been decoded, is to track the operation of the encoder, underflow will be avoided if it performs the same expansion operation under the same conditions.

Overflow

Now consider the possibility of overflow in the integer multiplications corresponding to those of Figure 2, which occur in lines 91–94 and 190–193

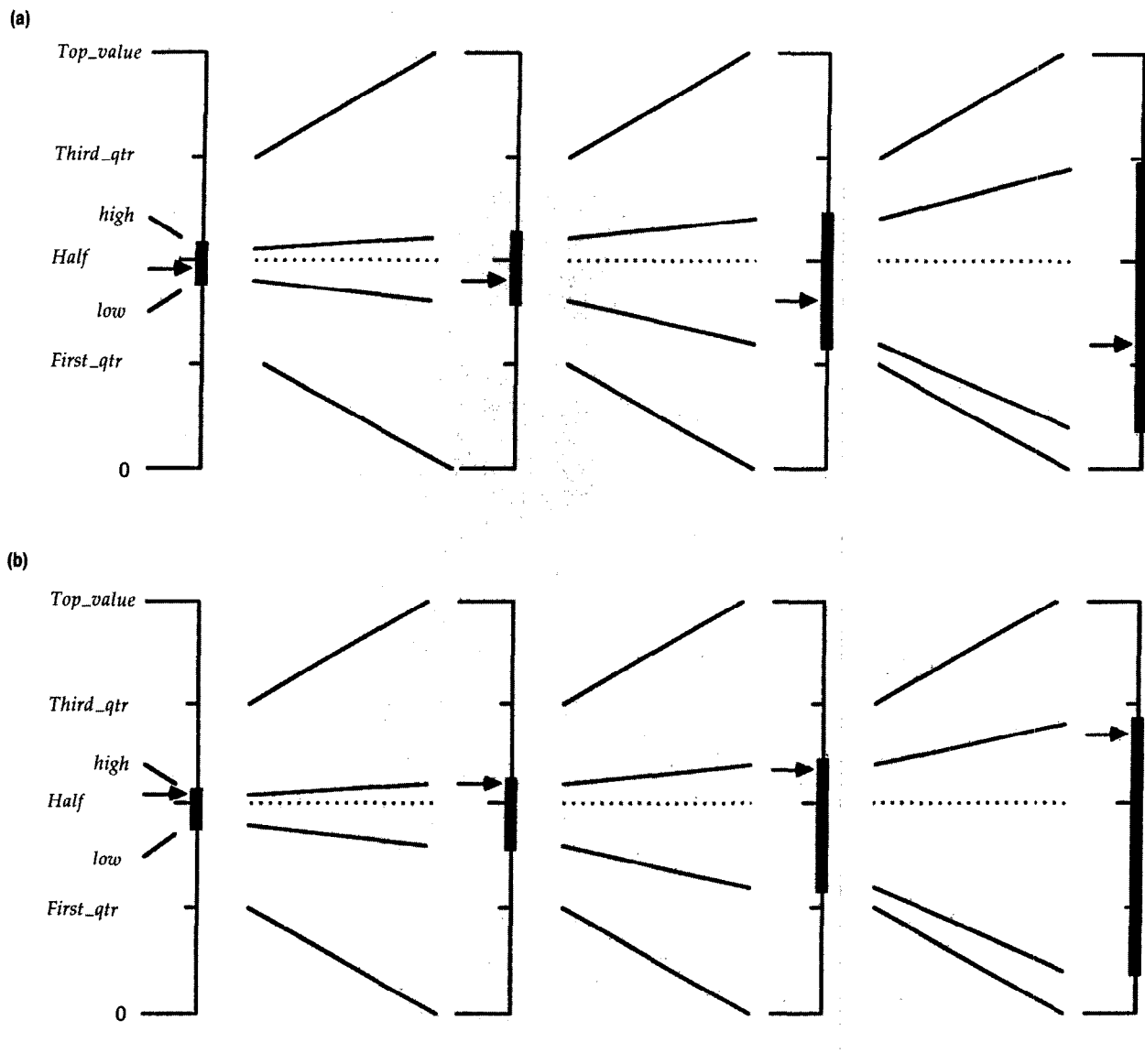


FIGURE 5. Scaling the Interval to Prevent Underflow

of Figure 3. Overflow cannot occur provided the product

$$\text{range} * \text{Max_frequency}$$

fits within the integer word length available, since cumulative frequencies cannot exceed *Max_frequency*. *Range* might be as large as *Top_value* + 1, so the largest possible product in Figure 3 is $2^{16}(2^{14} - 1)$, which is less than 2^{30} . *Long* declarations are used for *code_value* (line 7) and *range* (lines 89, 183) to ensure that arithmetic is done to 32-bit precision.

Constraints on the Implementation

The constraints on word length imposed by underflow and overflow can be simplified by assuming that frequency counts are represented in f bits, and *code_values* in c bits. The implementation will work correctly provided

$$\begin{aligned} f &\leq c - 2 \\ f + c &\leq p. \end{aligned}$$

the precision to which arithmetic is performed.

In most C implementations, $p = 31$ if *long* integers

are used, and $p = 32$ if they are *unsigned long*. In Figure 3, $f = 14$ and $c = 16$. With appropriately modified declarations, *unsigned long* arithmetic with $f = 15$ and $c = 17$ could be used. In assembly language, $c = 16$ is a natural choice because it expedites some comparisons and bit manipulations (e.g., those of lines 95–113 and 194–213).

If p is restricted to 16 bits, the best values possible are $c = 9$ and $f = 7$, making it impossible to encode a full alphabet of 256 symbols, as each symbol must have a count of at least one. A smaller alphabet (e.g., the 26 letters, or 4-bit nibbles) could still be handled.

Termination

To finish the transmission, it is necessary to send a unique terminating symbol (*EOF_symbol*, line 56) and then follow it by enough bits to ensure that the encoded string falls within the final range. Since *done_encoding()* (lines 119–123) can be sure that *low* and *high* are constrained by either Eq. (1a) or (1b) above, it need only transmit 01 in the first case or 10 in the second to remove the remaining ambiguity. It is convenient to do this using the *bit_plus_follow()* procedure discussed earlier. The *input_bit()* procedure will actually read a few more bits than were sent by *output_bit()*, as it needs to keep the low end of the buffer full. It does not matter what value these bits have, since EOF is uniquely determined by the last two bits actually transmitted.

MODELS FOR ARITHMETIC CODING

The program in Figure 3 must be used with a model that provides a pair of translation tables *index_to_char[]* and *char_to_index[]*, and a cumulative frequency array *cum_freq[]*. The requirements on the latter are that

- $cum_freq[i - 1] \geq cum_freq[i]$;
- an attempt is never made to encode a symbol i for which $cum_freq[i - 1] = cum_freq[i]$; and
- $cum_freq[0] \leq Max_frequency$.

Provided these conditions are satisfied, the values in the array need bear no relationship to the actual cumulative symbol frequencies in messages. Encoding and decoding will still work correctly, although encodings will occupy less space if the frequencies are accurate. (Recall our successfully encoding *eaii!* according to the model of Table I, which does not actually reflect the frequencies in the message.)

Fixed Models

The simplest kind of model is one in which symbol frequencies are fixed. The first model in Figure 4 has symbol frequencies that approximate those of English (taken from a part of the Brown Corpus [12]).

However, bytes that did not occur in that sample have been given frequency counts of one in case they do occur in messages to be encoded (so this model will still work for binary files in which all 256 bytes occur). Frequencies have been normalized to total 8000. The initialization procedure *start_model()* simply computes a cumulative version of these frequencies (lines 48–51), having first initialized the translation tables (lines 44–47). Execution speed would be improved if these tables were used to reorder symbols and frequencies so that the most frequent came first in the *cum_freq[]* array. Since the model is fixed, the procedure *update_model()*, which is called from both *encode.c* and *decode.c*, is null.

An *exact* model is one where the symbol frequencies in the message are exactly as prescribed by the model. For example, the fixed model of Figure 4 is close to an exact model for the particular excerpt of the Brown Corpus from which it was taken. To be truly exact, however, symbols that did not occur in the excerpt would be assigned counts of zero, rather than one (sacrificing the capability of transmitting messages containing those symbols). Moreover, the frequency counts would not be scaled to a predetermined cumulative frequency, as they have been in Figure 4. The exact model can be calculated and transmitted before the message is sent. It is shown by Cleary and Witten [3] that, under quite general conditions, this will *not* give better overall compression than adaptive coding (which is described next).

Adaptive Models

An adaptive model represents the changing symbol frequencies seen *so far* in a message. Initially all counts might be the same (reflecting no initial information), but they are updated, as each symbol is seen, to approximate the observed frequencies. Provided both encoder and decoder use the same initial values (e.g., equal counts) and the same updating algorithm, their models will remain in step. The encoder receives the next symbol, encodes it, and updates its model. The decoder identifies it according to its current model and then updates its model.

The second half of Figure 4 shows such an adaptive model. This is the type of model recommended for use with Figure 3, for in practice it will outperform a fixed model in terms of compression efficiency. Initialization is the same as for the fixed model, except that all frequencies are set to one. The procedure *update_model(symbol)* is called by both *encode_symbol()* and *decode_symbol()* (Figure 3, lines 54 and 151) after each symbol is processed.

Updating the model is quite expensive because of the need to maintain cumulative totals. In the code

of Figure 4, frequency counts, which must be maintained anyway, are used to optimize access by keeping the array in frequency order—an effective kind of self-organizing linear search [9]. *Update_model()* first checks to see if the new model will exceed the cumulative-frequency limit, and if so scales all frequencies down by a factor of two (taking care to ensure that no count scales to zero) and recomputes cumulative values (Figure 4, lines 29–37). Then, if necessary, *update_model()* reorders the symbols to place the current one in its correct rank in the frequency ordering, altering the translation tables to reflect the change. Finally, it increments the appropriate frequency count and adjusts cumulative frequencies accordingly.

PERFORMANCE

Now consider the performance of the algorithm of Figure 3, both in compression efficiency and execution time.

Compression Efficiency

In principle, when a message is coded using arithmetic coding, the number of bits in the encoded string is the same as the entropy of that message with respect to the model used for coding. Three factors cause performance to be worse than this in practice:

- (1) message termination overhead;
- (2) the use of fixed-length rather than infinite-precision arithmetic; and
- (3) scaling of counts so that their total is at most *Max_frequency*.

None of these effects is significant, as we now show. In order to isolate the effect of arithmetic coding, the model will be considered to be exact (as defined above).

Arithmetic coding must send extra bits at the end of each message, causing a message termination overhead. Two bits are needed, sent by *done_encoding()* (Figure 3, lines 119–123), in order to disambiguate the final symbol. In cases where a bit stream must be blocked into 8-bit characters before encoding, it will be necessary to round out to the end of a block. Combining these, an extra 9 bits may be required.

The overhead of using fixed-length arithmetic occurs because remainders are truncated on division. It can be assessed by comparing the algorithm's performance with the figure obtained from a theoretical entropy calculation that derives its frequencies from counts scaled exactly as for coding. It is completely negligible—on the order of 10^{-4} bits/symbol.

The penalty paid by scaling counts is somewhat larger, but still very small. For short messages (less than 2^{14} bytes), no scaling need be done. Even with messages of 10^5 – 10^6 bytes, the overhead was found experimentally to be less than 0.25 percent of the encoded string.

The adaptive model in Figure 4 scales down all counts whenever the total threatens to exceed *Max_frequency*. This has the effect of weighting recent events more heavily than events from earlier in the message. The statistics thus tend to track changes in the input sequence, which can be very beneficial. (We have encountered cases where limiting counts to 6 or 7 bits gives better results than working to higher precision.) Of course, this depends on the source being modeled. Bentley et al. [2] consider other, more explicit, ways of incorporating a recency effect.

Execution Time

The program in Figure 3 has been written for clarity rather than for execution speed. In fact, with the adaptive model in Figure 4, it takes about 420 μ s per input byte on a VAX-11/780 to encode a text file, and about the same for decoding. However, easily avoidable overheads such as procedure calls account for much of this, and some simple optimizations increase speed by a factor of two. The following alterations were made to the C version shown:

- (1) The procedures *input_bit()*, *output_bit()*, and *bit_plus_follow()* were converted to macros to eliminate procedure-call overhead.
- (2) Frequently used quantities were put in register variables.
- (3) Multiplies by two were replaced by additions (C “+=”).
- (4) Array indexing was replaced by pointer manipulation in the loops at line 189 in Figure 3 and lines 49–52 of the adaptive model in Figure 4.

This mildly optimized C implementation has an execution time of 214 μ s/252 μ s per input byte, for encoding/decoding 100,000 bytes of English text on a VAX-11/780, as shown in Table II. Also given are corresponding figures for the same program on an Apple Macintosh and a SUN-3/75. As can be seen, coding a C source program of the same length took slightly longer in all cases, and a binary object program longer still. The reason for this will be discussed shortly. Two artificial test files were included to allow readers to replicate the results. “Alphabet” consists of enough copies of the 26-letter alphabet to fill out 100,000 characters (ending with a partially completed alphabet). “Skew statistics” contains

TABLE II. Results for Encoding and Decoding 100,000-Byte Files

	VAX-11/780			Macintosh 512 K		SUN-3/75	
	Output (bytes)	Encode time (μ s)	Decode time (μ s)	Encode time (μ s)	Decode time (μ s)	Encode time (μ s)	Decode time (μ s)
Mildly optimized C implementation							
Text file	57,718	214	262	687	881	98	121
C program	62,991	230	288	729	950	105	131
VAX object program	73,501	313	406	950	1,334	145	190
Alphabet	59,292	223	277	719	942	105	130
Skew statistics	12,092	143	170	507	645	70	85
Carefully optimized assembly-language implementation							
Text file	57,718	104	135	194	243	46	58
C program	62,991	109	151	208	266	51	65
VAX object program	73,501	158	241	280	402	75	107
Alphabet	59,292	105	145	204	264	51	65
Skew statistics	12,092	63	81	126	160	28	36

Notes: Times are measured in microseconds per byte of uncompressed data.

The VAX-11/780 had a floating-point accelerator, which reduces integer multiply and divide times.

The Macintosh uses an 8-MHz MC68000 with some memory wait states.

The SUN-3/75 uses a 16.67-MHz MC68020.

All times exclude I/O and operating-system overhead in support of I/O. VAX and SUN figures give user time from the UNIX® *time* command; on the Macintosh, I/O was explicitly directed to an array.

The 4.2BSD C compiler was used for VAX and SUN; Aztec C 1.06g for Macintosh.

10,000 copies of the string *aaaabaaaac*; it demonstrates that files may be encoded into less than one bit per character (output size of 12,092 bytes = 96,736 bits). All results quoted used the adaptive model of Figure 4.

A further factor of two can be gained by reprogramming in assembly language. A carefully optimized version of Figures 3 and 4 (adaptive model) was written in both VAX and M68000 assembly languages. Full use was made of registers, and advantage taken of the 16-bit *code_value* to expedite some crucial comparisons and make subtractions of *Half* trivial. The performance of these implementations on the test files is also shown in Table II in order to give the reader some idea of typical execution speeds.

The VAX-11/780 assembly-language timings are broken down in Table III. These figures were obtained with the UNIX profile facility and are accurate only to within perhaps 10 percent. (This mechanism constructs a histogram of program counter values at real-time clock interrupts and suffers from statistical variation as well as some systematic errors.) "Bounds calculation" refers to the initial parts of *encode_symbol()* and *decode_symbol()* (Figure 3, lines 90–94 and 190–193), which contain multiply and divide operations. "Bit shifting" is the major loop in both the encode and decode routines (lines 95–113 and 194–213). The *cum* calculation in *decode_symbol()*, which requires a multiply/divide, and the following loop to identify the next symbol (lines 187–189), is "Symbol decode." Finally, "Model

TABLE III. Breakdown of Timings for the VAX-11/780 Assembly-Language Version

	Encode time (μ s)	Decode time (μ s)
Text file		
Bounds calculation	32	31
Bit shifting	39	30
Model update	29	29
Symbol decode	—	45
Other	4	0
	104	135
C program		
Bounds calculation	30	28
Bit shifting	42	35
Model update	33	36
Symbol decode	—	51
Other	4	1
	109	151
VAX object program		
Bounds calculation	34	31
Bit shifting	46	40
Model update	75	75
Symbol decode	—	94
Other	3	1
	158	241

update" refers to the adaptive *update_model()* procedure of Figure 4 (lines 26–53).

As expected, the bounds calculation and model update take the same time for both encoding and decoding, within experimental error. Bit shifting was quicker for the text file than for the C program and object file because compression performance was

better. The extra time for decoding over encoding is due entirely to the symbol decode step. This takes longer in the C program and object file tests because the loop of line 189 was executed more often (on average 9 times, 13 times, and 35 times, respectively). This also affects the model update time because it is the number of cumulative counts that must be incremented in Figure 4, lines 49–52. In the worst case, when the symbol frequencies are uniformly distributed, these loops are executed an average of 128 times. Worst-case performance would be improved by using a more complex tree representation for frequencies, but this would likely be slower for text files.

SOME APPLICATIONS

Applications of arithmetic coding are legion. By liberating *coding* with respect to a model from the *modeling* required for prediction, it encourages a whole new view of data compression [19]. This separation of function costs nothing in compression performance, since arithmetic coding is (practically) optimal with respect to the entropy of the model. Here we intend to do no more than suggest the scope of this view by briefly considering

- (1) adaptive text compression,
- (2) nonadaptive coding,
- (3) compressing black/white images, and
- (4) coding arbitrarily distributed integers.

Of course, as noted earlier, greater coding efficiencies could easily be achieved with more sophisticated models. Modeling, however, is an extensive topic in its own right and is beyond the scope of this article.

Adaptive text compression using single-character adaptive frequencies shows off arithmetic coding to good effect. The results obtained using the program in Figures 3 and 4 vary from 4.8–5.3 bits/char-

acter for short English text files (10^3 – 10^4 bytes) to 4.5–4.7 bits/character for long ones (10^5 – 10^6 bytes). Although adaptive Huffman techniques do exist (e.g., [5, 7]), they lack the conceptual simplicity of arithmetic coding. Although competitive in compression efficiency for many files, they are slower. For example, Table IV compares the performance of the mildly optimized C implementation of arithmetic coding with that of the UNIX *compact* program that implements adaptive Huffman coding using a similar model. (*Compact*'s model is essentially the same for long files, like those of Table IV, but is better for short files than the model used as an example in this article.) Casual examination of *compact* indicates that the care taken in optimization is roughly comparable for both systems, yet arithmetic coding halves execution time. Compression performance is somewhat better with arithmetic coding on all the example files. The difference would be accentuated with more sophisticated models that predict symbols with probabilities approaching one under certain circumstances (e.g., the letter *u* following *q*).

Nonadaptive coding can be performed arithmetically using fixed, prespecified models like that in the first part of Figure 4. Compression performance will be better than Huffman coding. In order to minimize execution time, the total frequency count, *cum_freq*[0], should be chosen as a power of two so the divisions in the bounds calculations (Figure 3, lines 91–94 and 190–193) can be done as shifts. Encode/decode times of around 60 μ s/90 μ s should then be possible for an assembly-language implementation on a VAX-11/780. A carefully written implementation of Huffman coding, using table lookup for encoding and decoding, would be a bit faster in this application.

Compressing black/white images using arithmetic coding has been investigated by Langdon and Rissanen [14], who achieved excellent results using

TABLE IV. Comparison of Arithmetic and Adaptive Huffman Coding

	Arithmetic coding			Adaptive Huffman coding		
	Output (bytes)	Encode time (μ s)	Decode time (μ s)	Output (bytes)	Encode time (μ s)	Decode time (μ s)
Text file	57,718	214	262	57,781	550	414
C program	62,991	230	288	63,731	596	441
VAX object program	73,546	313	406	76,950	822	606
Alphabet	59,292	223	277	60,127	598	411
Skew statistics	12,092	143	170	16,257	215	132

Notes: The mildly optimized C implementation was used for arithmetic coding.

UNIX *compact* was used for adaptive Huffman coding.

Times are for a VAX-11/780 and exclude I/O and operating-system overhead in support of I/O.

a model that conditioned the probability of a pixel's being black on a template of pixels surrounding it. The template contained a total of 10 pixels, selected from those above and to the left of the current one so that they precede it in the raster scan. This creates 1024 different possible contexts, and for each the probability of the pixel being black was estimated adaptively as the picture was transmitted. Each pixel's polarity was then coded arithmetically according to this probability. A 20–30 percent improvement in compression was attained over earlier methods. To increase coding speed, Langdon and Rissanen used an approximate method of arithmetic coding that avoided multiplication by representing probabilities as integer powers of $\frac{1}{2}$. Huffman coding cannot be directly used in this application, as it never compresses with a two-symbol alphabet. Run-length coding, a popular method for use with two-valued alphabets, provides another opportunity for arithmetic coding. The model reduces the data to a sequence of lengths of runs of the same symbol (e.g., for picture coding, run-lengths of black followed by white followed by black followed by white . . .). The sequence of lengths must be transmitted. The CCITT facsimile coding standard [11] bases a Huffman code on the frequencies with which black and white runs of different lengths occur in sample documents. A fixed arithmetic code using these same frequencies would give better performance; adapting the frequencies to each particular document would be better still.

Coding arbitrarily distributed integers is often called for in use with more sophisticated models of text, image, or other data. Consider, for instance, the locally adaptive data compression scheme of Bentley et al. [2], in which the encoder and decoder cache the last N different words seen. A word present in the cache is transmitted by sending the integer cache index. Words not in the cache are transmitted by sending a new-word marker followed by the characters of the word. This is an excellent model for text in which words are used frequently over short intervals and then fall into long periods of disuse. Their paper discusses several variable-length codings for the integers used as cache indexes. Arithmetic coding allows any probability distribution to be used as the basis for a variable-length encoding, including—among countless others—the ones implied by the particular codes discussed there. It also permits use of an adaptive model for cache indexes, which is desirable if the distribution of cache hits is difficult to predict in advance. Furthermore, with arithmetic coding, the code space allotted to the cache indexes can be scaled down to accommo-

date any desired probability for the new-word marker.

APPENDIX. Proof of Decoding Inequality

Using one-letter abbreviations for *cum_freq*, *symbol*, *low*, *high*, and *value*, suppose

$$c[s] \leq \left\lfloor \frac{(v-l+1) \times c[0] - 1}{h-l+1} \right\rfloor < c[s-1];$$

in other words,

$$c[s] \leq \frac{(v-l+1) \times c[0] - 1}{r} - \epsilon \leq c[s-1] - 1, \quad (1)$$

where

$$r = h - l + 1, \quad 0 \leq \epsilon \leq \frac{r-1}{r}.$$

(The last inequality of Eq. (1) derives from the fact that $c[s-1]$ must be an integer.) Then we wish to show that $l' \leq v \leq h'$, where l' and h' are the updated values for *low* and *high* as defined below.

$$\begin{aligned} \text{(a)} \quad l' &\equiv l + \left\lfloor \frac{r \times c[s]}{c[0]} \right\rfloor \\ &\leq l + \frac{r}{c[0]} \left\lfloor \frac{(v-l+1) \times c[0] - 1}{r} - \epsilon \right\rfloor \\ &\quad \text{from Eq. (1),} \\ &\leq v + 1 - \frac{1}{c[0]}, \end{aligned}$$

so $l' \leq v$ since both v and l' are integers and $c[0] > 0$.

$$\begin{aligned} \text{(b)} \quad h' &\equiv l + \left\lfloor \frac{r \times c[s-1]}{c[0]} \right\rfloor - 1 \\ &\geq l + \frac{r}{c[0]} \left\lfloor \frac{(v-l+1) \times c[0] - 1}{r} + 1 - \epsilon \right\rfloor - 1 \\ &\quad \text{from Eq. (1),} \\ &\geq v + \frac{r}{c[0]} \left[-\frac{1}{r} + 1 - \frac{r-1}{r} \right] = v. \end{aligned}$$

REFERENCES

1. Abramson, N. *Information Theory and Coding*. McGraw-Hill, New York, 1963. This textbook contains the first reference to what was to become the method of arithmetic coding (pp. 61–62).
2. Bentley, J.L., Sleator, D.D., Tarjan, R.E., and Wei, V.K. A locally adaptive data compression scheme. *Commun. ACM* 29, 4 (Apr. 1986), 320–330. Shows how recency effects can be incorporated explicitly into a text compression system.

3. Cleary, J.G., and Witten, I.H. A comparison of enumerative and adaptive codes. *IEEE Trans. Inf. Theory* IT-30, 2 (Mar. 1984), 306–315. Demonstrates under quite general conditions that adaptive coding outperforms the method of calculating and transmitting an exact model of the message first.
4. Cleary, J.G., and Witten, I.H. Data compression using adaptive coding and partial string matching. *IEEE Trans. Commun. COM-32*, 4 (Apr. 1984), 396–402. Presents an adaptive modeling method that reduces a large sample of mixed-case English text to around 2.2 bits/character when arithmetically coded.
5. Cormack, G.V., and Horspool, R.N. Algorithms for adaptive Huffman codes. *Inf. Process. Lett.* 18, 3 (Mar. 1984), 159–166. Describes how adaptive Huffman coding can be implemented efficiently.
6. Cormack, G.V., and Horspool, R.N. Data compression using dynamic Markov modeling. Res. Rep., Computer Science Dept., Univ. of Waterloo, Ontario, Apr. 1985. Also to be published in *Comput. J.* Presents an adaptive state-modeling technique that, in conjunction with arithmetic coding, produces results competitive with those of [4].
7. Gallager, R.G. Variations on a theme by Huffman. *IEEE Trans. Inf. Theory* IT-24, 6 (Nov. 1978), 668–674. Presents an adaptive Huffman coding algorithm, and derives new bounds on the redundancy of Huffman codes.
8. Held, G. *Data Compression: Techniques and Applications*. Wiley, New York, 1984. Explains a number of ad hoc techniques for compressing text.
9. Hester, J.H., and Hirschberg, D.S. Self-organizing linear search. *ACM Comput. Surv.* 17, 3 (Sept. 1985), 295–311. A general analysis of the technique used in the present article to expedite access to an array of dynamically changing frequency counts.
10. Huffman, D.A. A method for the construction of minimum-redundancy codes. *Proc. Inst. Electr. Radio Eng.* 40, 9 (Sept. 1952), 1098–1101. The classic paper in which Huffman introduced his famous coding method.
11. Hunter, R., and Robinson, A.H. International digital facsimile coding standards. *Proc. Inst. Electr. Electron. Eng.* 68, 7 (July 1980), 854–867. Describes the use of Huffman coding to compress run lengths in black/white images.
12. Kucera, H., and Francis, W.N. *Computational Analysis of Present-Day American English*. Brown University Press, Providence, R.I., 1967. This large corpus of English is often used for computer experiments with text, including some of the evaluation reported in the present article.
13. Langdon, G.G. An introduction to arithmetic coding. *IBM J. Res. Dev.* 28, 2 (Mar. 1984), 135–149. Introduction to arithmetic coding from the point of view of hardware implementation.
14. Langdon, G.G., and Rissanen, J. Compression of black-white images with arithmetic coding. *IEEE Trans. Commun. COM* 29, 6 (June 1981), 858–867. Uses a modeling method specially tailored to black/white pictures, in conjunction with arithmetic coding, to achieve excellent compression results.
15. Pasco, R. Source coding algorithms for fast data compression. Ph.D. thesis, Dept. of Electrical Engineering, Stanford Univ., Stanford, Calif., 1976. An early exposition of the idea of arithmetic coding, but lacking the idea of incremental operation.
16. Rissanen, J.J. Generalized Kraft inequality and arithmetic coding. *IBM J. Res. Dev.* 20 (May 1976), 198–203. Another early exposition of the idea of arithmetic coding.
17. Rissanen, J.J. Arithmetic codings as number representations. *Acta Polytech. Scand. Math.* 31 (Dec. 1979), 44–51. Further develops arithmetic coding as a practical technique for data representation.
18. Rissanen, J., and Langdon, G.G. Arithmetic coding. *IBM J. Res. Dev.* 23, 2 (Mar. 1979), 149–162. Describes a broad class of arithmetic codes.
19. Rissanen, J., and Langdon, G.G. Universal modeling and coding. *IEEE Trans. Inf. Theory* IT-27, 1 (Jan. 1981), 12–23. Shows how data compression can be separated into *modeling* for prediction and *coding* with respect to a model.
20. Rubin, F. Arithmetic stream coding using fixed precision registers. *IEEE Trans. Inf. Theory* IT-25, 6 (Nov. 1979), 672–675. One of the first papers to present all the essential elements of practical arithmetic coding, including fixed-point computation and incremental operation.
21. Shannon, C.E., and Weaver, W. *The Mathematical Theory of Communication*. University of Illinois Press, Urbana, Ill., 1949. A classic book that develops communication theory from the ground up.
22. Welch, T.A. A technique for high-performance data compression. *Computer* 17, 6 (June 1984), 8–19. A very fast coding technique based on the method of [23], but whose compression performance is poor by the standards of [4] and [6]. An improved implementation of this method is widely used in UNIX systems under the name *compress*.
23. Ziv, J., and Lempel, A. Compression of individual sequences via variable-rate coding. *IEEE Trans. Inf. Theory* IT-24, 5 (Sept. 1978), 530–536. Describes a method of text compression that works by replacing a substring with a pointer to an earlier occurrence of the same substring. Although it performs quite well, it does not provide a clear separation between modeling and coding.

CR Categories and Subject Descriptors: E.4 [Data]: Coding and Information Theory—*data compaction and compression*; H.1.1 [Models and Principles]: Systems and Information Theory—*information theory*
General Terms: Algorithms, Performance

Additional Key Words and Phrases: Adaptive modeling, arithmetic coding, Huffman coding

Received 8/86; accepted 12/86

Authors' Present Address: Ian H. Witten, Radford M. Neal, and John G. Cleary, Dept. of Computer Science, The University of Calgary, 2500 University Drive NW, Calgary, Canada T2N 1N4.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

In response to membership requests . . .

CURRICULA RECOMMENDATIONS FOR COMPUTING

- Volume I: Curricula Recommendations for Computer Science
- Volume II: Curricula Recommendations for Information Systems
- Volume III: Curricula Recommendations for Related Computer Science Programs in Vocational-Technical Schools, Community and Junior Colleges and Health Computing

Information available from Deborah Cotton—Single Copy Sales (212) 869-7440 ext. 309