

Computer Science 222: Succincter

Tuesday, December 11, 2012

Professor Mitzenmacher

Dan Bradley (dbradley@college), Saagar Deshpande (sdeshpande@college)

Final Project

1 Abstract

Using Mihai Patrascu's 2008 paper "Succincter" [1], we implement a way to store trits (ternary values) within 1.01% of the ideal space of $n * \log_2(3)$ while having lookup in $O(t)$ time, where t is the depth of our data structure. We use recursion to improve redundancy in the data structure, which allows us to accurately compute single decodes quickly without needing to decode the entire compressed block first. We find that this is both a fast and space efficient robust data structure for encoding and decoding operations with room for extension past simply storing trits.

2 Introduction

There are few effective methods for storing trits. In this paper, we focus on three of them: the naive method, arithmetic coding, and our succincter implementation.

The naive method is simple. As a trit has more information than a bit, but less than two bits, store each trit using two bits. Obviously this is not the most space-efficient method, you could encode $\frac{4}{3}$ as much information in the same space, so it is clearly wasteful. It is however, very fast, both on encode and decode. Encode is strictly linear, with very small constants. The encoder we used took 3.9 milliseconds to encode 500000 entries, and 13.7 milliseconds to decode 500000 entries.

Arithmetic coding is much more effective at reducing size. By reducing the entries to one number, we can encode in exactly $\lceil \frac{n * \log_2(3)}{8} \rceil$ bytes. Decoding, however, ideally takes linear time, as to reach any original entry, you must go through every entry prior, there is no way to randomly access any element. Furthermore, the nature of arithmetic coding lends itself to manipulation of large numbers, which is time consuming.

The succincter data structure is a recursive approach to the problem, and attempts to approach the ideal space usage without sacrificing on decoding time, especially for individual entries. Essentially, succincter takes an array of entries, and breaks them into blocks. These blocks are each treated as one large number, base K (in the first level of the structure when storing trits K is 3, but K changes by level). These large number are then divided by 2^M where M is a user-defined variable. The remainder is stored for transmission, and the quotient gets passed up to the next level of the data structure as a new array. The quotient will be some number $\{0, 1, 2, \dots, K - 1\}$, and so we can repeat the process on the new array by simply finding new blocks and treating them as numbers base K .

Mihai Patrascu outlined the idea of succincter data structures in [1], but we took some liberties with our implementation, as his approaches to finding block size and K proved less useful in practice than they might have been in theory.

3 Implementation

3.1 Encoding

To encode, we start with an array with fixed size of trits with pseudorandomly assigned values. For the sake of simplicity, and the fact that the values of the trits do not particularly matter for our inquiry, we used the *c* standard library `rand()` function seeded by the time of running.

We also decided that all remainders would be up to 32 bits long, which for the sake of discussion we will refer to as M .

With the size of the array and the base 3, we construct the necessary headers to be used by the decoder, specifically the size of the array at that level, the base at that level, and the block size at that level.

Given the size and base, the header of all levels are constructable. We are trying to fill an array of size `LEVELS` (an arbitrary constant we the user decide) + 1 with the size for that level, the block size for that level, and the base for that level. In the context of this paper, we will call this array headers. We find an n such that our new $K = \lceil \frac{3^n}{2^M} \rceil > 2$. n is the block size for the level, the base for the level is the old K , and the size for the level is the old size. The variables then update, with old K being set to new K , size decreasing by a factor `chunksize`, and we are back where we started, just a level higher up the headers.

Once the headers are found, we can begin to encode the actual data. We encode recursively, starting from the bottom level = 0 and continuing until level equals `LEVELS`. To encode, we are given an array to encode, we take the first `headers[level].chunksize` entries, and combine them into one large number. We then divide by 2^M , store the remainder, and pass the quotient to a new array we're building for the next level. We do this for each block, and when we reach the end of our entries, we call decode again, on the new array, with level increased.

Eventually level will equal `LEVELS`, and when it does, we naively code the array: we figure out how many bits it takes to completely store a number base K , and store each number in the array using that many bits.

We put in an output file the `SIZE`, `LEVELS`, and `M`, write all the headers to the file, then write out all the bits for the last level of storage.

```
Encode(array, level):
if (level = LEVELS OR headers[level].size = 1)
    encode_last(array, headers[level]);
counter = 0;
index = 0;
new_array[];
stringtotal;
while(counter < headers[level].size)
    for(int i = 0; i < headers[level].chunksize; i++)
        total[i] = array[counter];
        counter++;
    number = num_from_string(total, base h.K);
    number = number/2M;
    new_array[index] = number/2M;
    m = number%2M;
    fwrite(output, m);
    index++;
return Encode(new_array, level + 1);
```

```

Encode_last(array, header):
  counter = 0;
  while(header.K > 0)
    header.K = header.K/2;
    counter ++;
  cycler = 1;
  byte = 0;
  for(i = 0; i < header.size; i++)
    k = 1;
    for(j = 0; j < counter; j++)
      if(k & array[i])
        byte = byte|cycler;
      k = 2 * k;
      cycler = 2 * cycler;
      if(cyclor = 28)
        fwrite(output, byte);
        byte = 0;
        cyclor = 1;
  fwrite(output, byte);
  return;

```

3.2 Decoding

For both decoding schemes, we first need to read and interpret the input file, which is done essentially in reverse of encoding, taking SIZE, LEVELS, and M to be global variables, as well as the headers array, an array of all the remainders, and the last array that was encoded.

3.2.1 Decode All

We are given an array (to start with, the last array), and that we are on the highest level, and with that and the global variables, we can decode every original entry in the right order. To do this, we go through the array with index, take the remainder corresponding to the index in and find remainders[current_level - 1][index] + $2^M \cdot \text{array}[\text{index}]$. This is the same number we created by amalgamating each block in the encoder. To get the array for the next level down, we simply divide the number by the base of the next level down, and place each remainder in a new array. If we go through the entire array like this, we can traverse the data structure level by level passing new array and the current level, and eventually decode the entire structure.

```

Decode_all(array, level):
  if(level = 0)
    print(array);
  index = 0;
  for(i = 0; i < h.size; i++)
    number = remainders[level - 1][i];
    number = number + 2M * array[i];
    min = min(headers[level - 1].chunksize, headers[level - 1].size - index);
    for(j = 0; j < min; j++)
      new_array[index + min - j - 1] = number%headers[level - 1].K;
      number = number/headers[level - 1].K;
    index += min;
  return Decode_all(new_array, level - 1);

```

3.2.2 Decode One

Decoding one value alone is a more complicated process. We need to determine the path of remainders and array values to follow from the highest level to the lowest. To do this, we can simply find where in the completely decoded array the item would be, and then by using the headers, we find where in each block the number would be, and where in the highest level's array the object to be decoded is.

Once we have the path, it becomes simply a matter of following, instead of building arrays, you only keep the object identified by the path, and decode through there alone. That is why it only takes linear time with respect to the depth of the data structure.

```

Make_path(int n):
for(i = 0; i < levels + 1; i ++ )
    paths[i].where = n;
    paths[i].path = n % headers[i].chunksize;
    n = n / headers[i].chunksize;

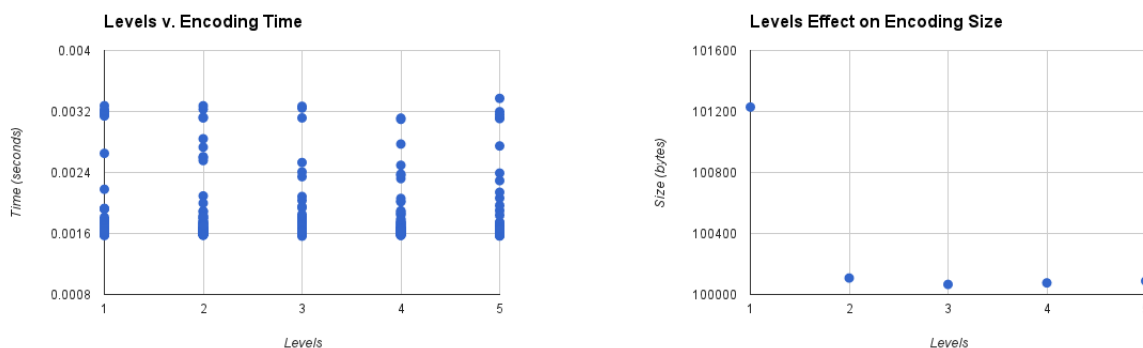
Decode_one():
num = final[paths[levels].path];
for(i = levels - 1; i >= 0; i -- )
    number = remainders[i][paths[i + 1].where];
    number = number + 2M * num;
    for(j = headers[i].chunksize - 1; j > paths[i].path; j -- )
        number = number / headers[i].K;
    number = number / headers[i].K;
    num = number;
return num;

```

The full encoding and decoding files are in the Appendix.

4 Results and Analysis

We ran our implementation of Succincter on 500000 trits except in the case of individual decoding, where we only use 200000 trits. We compared this implementation against a simple C implementation of arithmetic encoding as well as a naive C implementation of a trit to bit converter, which stores each trit as its corresponding bit.

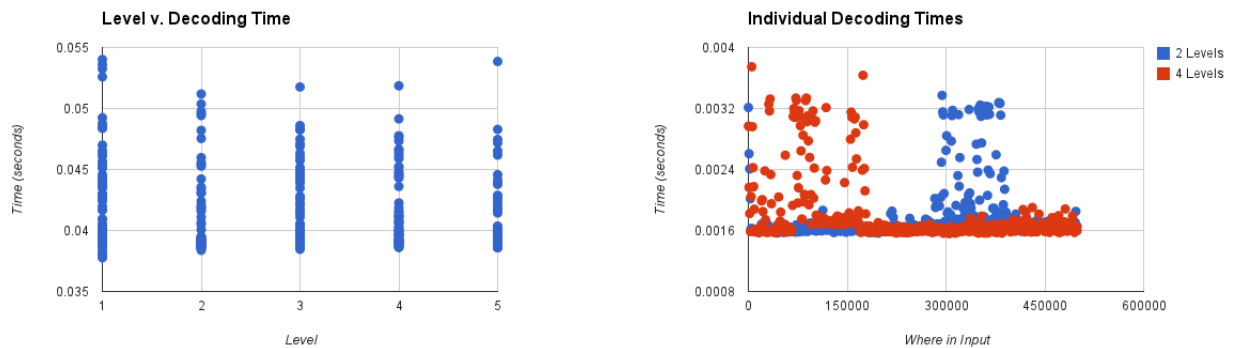


Initially, we compared the speed of encoding and decoding using multiple depths for the data structure as well as the space required for the encoding process. We noticed that the encoding time is relatively similar for each of the depths we run Succincter on, with level 4 edging the others by an insignificant amount. This indicates that the time taken is relatively similar for any depth of structure. One likely reason is the fact

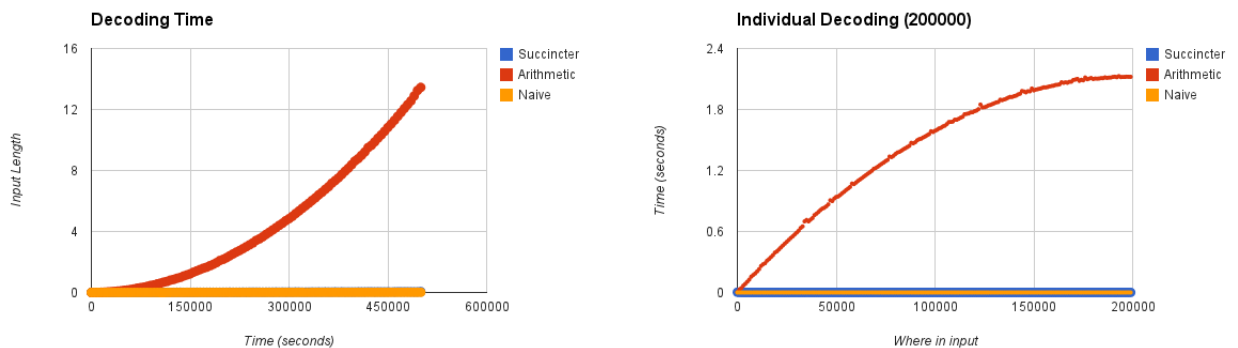
that the maximum useful depth is $\log(n)$, as we divide the size down each level by the block size, and an array with a size of 1 should just be encoded wasting the one bit. The reason that we see the high levels of encoding to be slightly faster is because the last layer of values need to be specially encoded from their current base to bits. For the lower levels, this is a more expensive process as the number of values left to be specially encoded is much larger. However, this time is still largely insignificant for large sets of values.

This reasoning for the encoding time also extends to the average encoding sizes we see in our experimentation. We see that as the level increases, we have an immediate decrease in size from level 1 to level 2, and after that the size remains roughly the same. This likely indicates that the size of the encoded data will remain close to the same for any level past level 1 given this size of input, though it is very possible that with larger inputs higher level caps become more effective.

Level	Avg. Encode Time (sec)	Avg. Decode Time (sec)	Avg. Encode Size (bytes)
1	0.00178906	0.043107825396825	101229
2	0.00177046	0.041240047619048	100108
3	0.00173798	0.041672920634921	100066
4	0.00173344	0.041674126984127	100076
5	0.00175082	0.040755396825397	100088

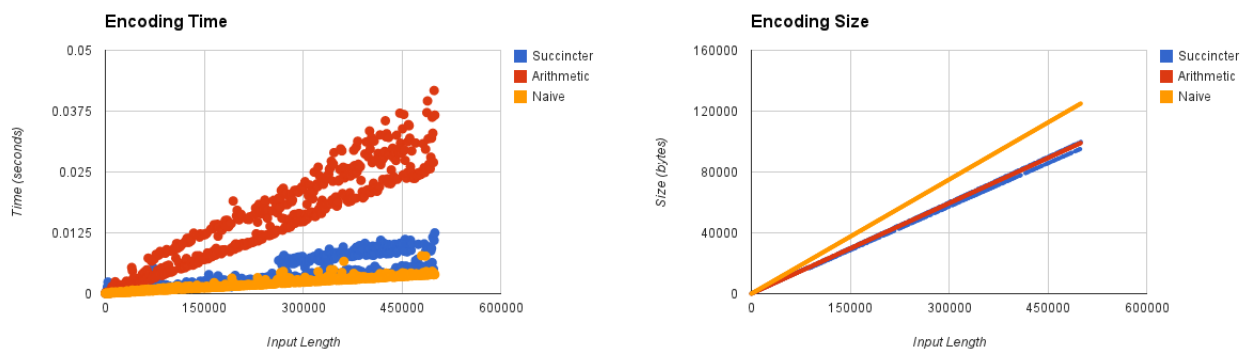


Looking at decoding across each level, we notice that the average decoding times are very similar, with a slight decrease as the level cap increases. The change is too small to be of much statistically interest, leading us to believe that the decode time is essentially the same for each depth on average. We can get a better sense of this by looking at the individual decoding time scatter plot. Here, it appears that in the majority of cases, both level 2 and level 4 decoding are roughly the same time. The red scatter mass on the left and the blue scatter mass on the right seem to be noise in the system. There does seem to be much longer times for level 2 for larger values, and longer times for level 4 on smaller individual values, but a good explanation for why that would be the case is elusive.



In general, it appears that decoding is not impacted by the level of encoding used on the input string, and the only factor we will have to worry about when decoding is the length of the input string. Of some more importance is the time required to decode. We compared the decoding time for the three implementations mentioned earlier, the naive trit to bit translation, the simple arithmetic encoder, and our Succincter model. We noticed that the naive method is definitely the fastest of the three, which intuitively makes sense because each trit is stored in its full bit form. Succincter is within a factor of 3 in speed to the naive decoder; which, while not great, can be seen as a necessary tradeoff for the improved efficiency. This means that we can encode information and store it in a more compact form than the naive method while still having similar decoding times across the entire input. Succincter handles the arithmetic encoder time with ease; it appears that as the length of the input increases, the arithmetic decoder grows in polynomial time, while the Succincter implementation barely grows at all. Because we use similar, or even less, space than the arithmetic encoding scheme, Succincter seems to be a viable replacement for arithmetic encoding if properly implemented. It is also possible however, that our arithmetic encoder was not quite of the same quality, and we could find (or make) one that would be better competition with Succincter. Especially since theoretically, the decode time should be on the same order.

The area that Succincter should be asymptotically faster than arithmetic coding is in individual decoding time. Arithmetic coding should run in linear time, while the naive approach and Succincter should be constant time algorithms. Once again, we see that the naive method is the fastest one; for any individual trits, we require no calculation and can simply convert the number of bits we see into trits. Succincter takes 300 times as long, to decode each individual, but still takes less than a millisecond to run. Succincter is able to stay within the constant arithmetic bounds for individual decoding time because of the properties of the spillover algorithm and universe as defined in the Patrascu paper. Arithmetic decoding of individual trits requires the entire string to be decoded first, which seems to end up slowing down the arithmetic decoder a lot, leading to the logarithmic curve seen in the graph above.



We also compared the encoding times and encoding sizes of the three approaches. The naive takes the least amount of time to encode, which is obvious as there is no real computation to be done for the encoding process. Succincter is only marginally slower than the naive progress, and it appears that both naive and Succincter have a much smaller slope of growth when compared to the arithmetic encoder. All three algorithms grow linearly with the size of the input string; this is true because when a new character is presented, we act on it when it appears, so we only need to read the input string one time for each implementation in order to encode it. Based on encoding time, Succincter beats arithmetic encoding by more than a factor of two on average, and naive encoding beats Succincter by more than a factor of two on average.

When we look at the encoding size, it is immediately apparent that the naive approach requires the most space. Each trit in the string is stored in its entirety as bits, so there is no change to encode and compress the information. When we compare Succincter and arithmetic coding, both are better compressed than naive, with Succincter just slightly performing better at longer string lengths.

Looking just at encoding time and size, we see that Succincter is an effective algorithm to use. In terms of encoding time, Succincter is comparable, though worse to the faster naive implementation, while

in terms of encoding size, Succincter is comparable to the arithmetic encoder implementation.

5 Conclusion

Succincter seems like an effective algorithm to use for while naive implementation is much faster on the whole, there are obviously data sets where Succincter could prove useful. Datasets that have plenty of space and need to be fast would be better served by using the naive approach, though the space limited user should use Succincter. Succincter allows us to have both speed and compression without major losses of one or the other. Furthermore, looking at decode time, Succincter has relatively quick decoding times for both the entire string and for individual decoding. The individual decoding aspect makes Succincter especially powerful because it is now possible to look up any single trit in the string without having to spend the total time to decode the entire encoded string. Arithmetic decoders require the entire encoded string to be decoded before an answer can be found; this is not true with Succincter.

Additional steps to be taken are clearly test to our structure against a better arithmetic encoder, and see how the times compare. Spacewise, even a perfect arithmetic encoder could not improve much on our structure, as the ideal $n * \log_2(3)$ space for $n = 500000$ is 99060 bytes; Succincter used 100066 bytes. Our time comparisons however, do not accurately represent a fully functioning arithmetic encoder, and so to understand the true utility of the structure, it would be useful to have an upper time bound.

Further (assuming this structure does stand the test of time) it can be applied to any time an arithmetic coder could be used. Assuming the arithmetic encoder uses rational numbers for its probabilities, the lowest common denominator of the probabilities could be used as the base for the first level of the data structure. Take the example where the probabilities are ('a' = .4, 'b' = .1, 'c' = .5). Multiplying by ten, it's clear these numbers could all be treated as numbers base ten, ('a' = 0,1,2,3, 'b' = 4, 'c' = 5,6,7,8,9). So to encode, take the integer values of the string as the initial array, and the size to be the length of the string, and with that you can encode and decode just as effectively as with arithmetic encoding.

6 Acknowledgements

We would like to thank Micheal Mitzenmacher for inspiring this project.

7 References

Patrascu, Mihai. Succincter. 49th IEEE Symposium on Foundations of Computer Science. 2008.

8 Appendix

8.1 Succincter Code

Listing 1: Succincter Encoder

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <stdint.h>
4 #include <time.h>
5 #include <gmp.h>
6 #include <math.h>
7 #include <sys/stat.h>
```

```

8
9 #define SIZE 500000
10 #define LEVELS 2
11 #define M 32
12 FILE* output;
13 int mainarray[SIZE];
14 int arrindex = 0;
15 mpz_t number, m, k, divisor, try;
16
17 typedef struct
18 {
19     unsigned int chunksize;
20     unsigned int K;
21     unsigned int size;
22 } header;
23
24 header headers[LEVELS + 1];
25 int* final;
26
27 void make_headers(void)
28 {
29     mpz_t newk, oldk, thing, div;
30     mpz_init(newk);
31     mpz_init(thing);
32     mpz_ui_pow_ui(thing, 2, M);
33     mpz_init_set(div, thing);
34     mpz_init_set_ui(oldk, 3);
35     int size = SIZE;
36     for(int i = 0; i <= LEVELS; i++)
37     {
38         mpz_set_ui(newk, 0);
39         int r = 1;
40         int n;
41         while(mpz_get_ui(newk) < 3)
42         {
43             r++;
44             double bottom = mpz_get_d(oldk);
45             n = floor(((M+1)*log(2.0) + log((double) r))/log(bottom));
46             mpz_pow_ui(thing, oldk, n);
47             mpz_cdiv_q(newk, thing, div);
48         }
49         headers[i].chunksize = n;
50         headers[i].K = mpz_get_ui(oldk);
51         mpz_set(oldk, newk);
52         headers[i].size = size;
53         size = ceil((double) size/n);
54     }
55     fwrite(headers, sizeof(header), LEVELS + 1, output);
56     mpz_clear(oldk);
57     mpz_clear(div);
58     mpz_clear(thing);
59     mpz_clear(newk);
60 }
61
62
63
64 void rand_trit(void)
65 {
66     FILE* trits = fopen("trits.txt", "w");
67     for(int i = 0; i < SIZE; i++)
68     {
69         mainarray[i] = rand() % 3;
70         char c = mainarray[i] + '0';
71         fwrite(&c, sizeof(char), 1, trits);
72     }
73     fclose(trits);
74 }
75

```



```

76 void encode_last(int* array, header h)
77 {
78     double counter = 0;
79     while (h.K > 0)
80     {
81         h.K = h.K / 2;
82         counter++;
83     }
84     int bytes = ceil(h.size*counter/8);
85     int cycler = 1;
86     uint8_t writer = 0;
87     for(int i = 0; i < h.size; i++)
88     {
89         int k = 1;
90         for(int j = 0; j < counter; j++)
91         {
92             if(k & array[i])
93             {
94                 writer = writer | cycler;
95             }
96             k = k << 1;
97             cycler = cycler << 1;
98             if (cycler == 256)
99             {
100                 fwrite(&writer, sizeof(uint8_t), 1, output);
101                 writer = 0;
102                 cycler = 1;
103             }
104         }
105     }
106     fwrite(&writer, sizeof(uint8_t), 1, output);
107     return;
108 }
109
110 void encode(int* array, header h, int level)
111 {
112     int allzeroes = 1;
113     if (level == LEVELS || h.size == 1)
114     {
115         return encode_last(array, h);
116     }
117     int counter = 0;
118     int index = 0;
119     int new_array[headers[level+1].size];
120     char* total = malloc(sizeof(char)*h.size);
121     if (total == NULL)
122     {
123         printf("ERROR!\n");
124         exit(1);
125     }
126     char* start = total;
127     while(counter < h.size)
128     {
129         if(h.size - counter < h.chunksize)
130         {
131             for(int i = 0; i < h.chunksize - h.size + counter; i++)
132                 total++;
133             int tot = h.size - counter;
134             for (int i = 0; i < tot; i++)
135             {
136                 total[i] = array[counter] + '0';
137                 counter++;
138             }
139         }
140         else
141         {
142             for (int i = 0; i < h.chunksize; i++)
143             {

```

```

144         total[i] = array[counter] + '0';
145         counter++;
146     }
147 }
148 mpz_set_str(number, total, h.K);
149 mpz_fdiv_qr(k,m,number,divisor);
150 if(mpz_sgn(k))
151     allzeroes = 0;
152 new_array[index] = mpz_get_ui(k);
153 unsigned int remains = mpz_get_ui(m);
154 fwrite(&remains, sizeof(unsigned int), 1, output);
155 index++;
156 }
157 free(start);
158 if (allzeroes)
159     return;
160 return encode(new_array, headers[level+1], level + 1);
161 }
162
163
164 int main(void)
165 {
166     srand(time(NULL));
167     rand_trit();
168
169     output = fopen("encoded.bin", "w");
170     int l[3] = {SIZE, LEVELS, M};
171     fwrite(&l, sizeof(int), 3, output);
172     make_headers();
173     mpz_init(number);
174     mpz_init(m);
175     mpz_init(k);
176     mpz_init(try);
177     mpz_ui_pow_ui(try, 2, M);
178     mpz_init_set(divisor, try);
179     mpz_clear(try);
180     encode(mainarray, headers[0], 0);
181     fclose(output);
182     mpz_clear(number);
183     mpz_clear(m);
184     mpz_clear(k);
185     mpz_clear(divisor);
186     struct stat st;
187     stat("encoded.bin", &st);
188     int size = st.st_size;
189     printf("%d\n", size);
190 }

```

Listing 2: Succincter Decoder

```

1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <stdint.h>
4 #include <time.h>
5 #include <gmp.h>
6 #include <math.h>
7
8 typedef struct
9 {
10     int chunksize;
11     int K;
12     int size;
13 } header;
14
15 typedef struct
16 {
17     int path;

```

```

18     int where;
19 } path;
20
21 FILE* input;
22 int M;
23 int levels;
24 int size;
25 header* headers;
26 mpz_t** remainders;
27 int* final;
28 path* paths;
29
30 void finish(int* array)
31 {
32     FILE* output = fopen("decoded.txt", "w");
33     char c;
34     for(int i = 0; i < headers[0].size; i++)
35     {
36         c = array[i] + '0';
37         fwrite(&c, sizeof(char), 1, output);
38     }
39     fclose(output);
40     return;
41 }
42
43 void make_path(int n)
44 {
45     paths = malloc(sizeof(int)*(levels + 1));
46     for(int i = 0; i < levels + 1; i++)
47     {
48         paths[i].where = n;
49         paths[i].path = n % headers[i].chunksize;
50         n = n / headers[i].chunksize;
51     }
52 }
53
54 void decode_all(int* array, header h, int level)
55 {
56     if(!level)
57     {
58         return finish(array);
59     }
60
61     int* new_array = malloc(sizeof(int)*headers[level-1].size);
62     int index = 0;
63     mpz_t number, spill;
64     mpz_init(number);
65     mpz_init(spill);
66     for(int i = 0; i < h.size; i++)
67     {
68         mpz_ui_pow_ui(spill, 2, M);
69         mpz_set(number, remainders[level-1][i]);
70         mpz_addmul_ui(number, spill, array[i]);
71         int min = headers[level-1].chunksize > (headers[level-1].size - index)?(headers[level-1].size - index):headers[level-1].chunksize;
72
73         for(int j = 0; j < min; j++)
74         {
75             new_array[index + min - j - 1] = mpz_fdiv_q_ui(number, number, headers[level-1].K);
76         }
77         index += min;
78     }
79     return decode_all(new_array, headers[level-1], level-1);
80 }
81
82 int decode_one(void)
83 {
84     int num = final[paths[levels].path];

```

```

85     mpz_t number, spill;
86     mpz_init(number);
87     mpz_init(spill);
88     mpz_ui_pow_ui(spill, 2, M);
89     for(int i = levels - 1; i >= 0; i--)
90     {
91         mpz_set(number, remainders[i][paths[i + 1].where]);
92         mpz_addmul_ui(number, spill, num);
93         for(int j = headers[i].chunksize - 1; j > paths[i].path; j--)
94         {
95             mpz_fdiv_q_ui(number, number, headers[i].K);
96         }
97         mpz_fdiv_r_ui(number, number, headers[i].K);
98         num = mpz_get_ui(number);
99     }
100     return num;
101 }
102
103 void extract_last(header h)
104 {
105     final = malloc(sizeof(int)*h.size);
106     int counter = 0;
107     while (h.K > 0)
108     {
109         h.K = h.K / 2;
110         counter++;
111     }
112     int index = -1;
113     int k;
114     int count = 0;
115     int cycler = 1;
116     uint8_t reader;
117     while(fread(&reader, sizeof(uint8_t), 1, input))
118     {
119         while(cycler != 256)
120         {
121             if(!(count % counter))
122             {
123                 index++;
124                 k = 1;
125             }
126             if(reader & cycler)
127             {
128                 final[index] = final[index] | k;
129             }
130             cycler = cycler << 1;
131             k = k << 1;
132             count++;
133         }
134         cycler = 1;
135     }
136 }
137
138 void extract_data(void)
139 {
140     input = fopen("encoded.bin", "r");
141     fread(&size, sizeof(int), 1, input);
142     fread(&levels, sizeof(int), 1, input);
143     fread(&M, sizeof(int), 1, input);
144     headers = malloc(sizeof(header) * (levels+1));
145     fread(headers, sizeof(header), levels + 1, input);
146     remainders = malloc(sizeof(mpz_t)*levels);
147     for(int i = 0; i < levels; i++)
148     {
149         remainders[i] = malloc(sizeof(mpz_t)*headers[i + 1].size);
150         for(int j = 0; j < headers[i + 1].size; j++)
151         {
152             unsigned int remains;

```

```
153         fread(&remains, sizeof(unsigned int), 1, input);
154         mpz_init_set_ui(remainders[i][j], remains);
155     }
156 }
157 extract_last(headers[levels]);
158 }
159
160
161 int
162 main(int argc, char*argv[])
163 {
164     extract_data();
165     if(argc == 2)
166     {
167         int num = atoi(argv[1]);
168         if (num < size)
169         {
170             make_path(num);
171             int p = decode_one();
172             printf("%d\n", p);
173         }
174         else
175             printf("error!\n");
176     }
177     else
178         decode_all(final, headers[2], 2);
179 }
```

8.2 Full Size Images

