

Binary Code Compression Based on Decision

Trees

Tamás Gergely Ferenc Havasi

Tibor Gyimóthy

Department of Software Engineering,

University of Szeged,

Árpád tér 2.

6720 Szeged, Hungary

`{gertom,havasi,gyimothy}@inf.u-szeged.hu`

Abstract

More and more embedded systems are used in the world. These machines have limited resources, e. g. background storage size. It is always a task to increase the storage capacity of these machines. A possible solution for this is compressing the files before they are written on the storage device. This is usually done by a general compressor.

There are files containing text or binary data and in many cases binary program code. A general compressor compresses all of them with almost the same efficiency. But a compressor specialized for one type of input

would compress files of that type much better; thus using many special compressors would save space and increase the virtual capacity of the storage device. Using this idea, our goal was to create a compression algorithm for binary program code.

Most compression methods can be separated into two parts: model and coder. In this paper we introduce a decision-tree based modeller method. We combined this method with an arithmetic coder and applied it in a modified JFFS2 file system of a Linux distribution running on a PDA machine. The original file system used only one method for file compression (zlib), the modified file system uses many compressors including our model-based one. This PDA machine has an ARM processor thus our method were implemented for the compression of ARM instructions.

The results were very promising: depending on the parameters of the method at least 12.6% of the 13MB image created with only zlib compression were saved, and it costed only at most 3 times boot speed slowdown. The slowdown of loading a file cannot be detected otherwise (except some special cases).

Keywords: code compression, decision trees

1 Introduction

There are more and more embedded systems spreading all around the world (e. g. mobile phones, PDAs). The capacity of these small computers is limited in many ways. These limitations include processor speed, available energy, dynamic memory and background storage size.

In many cases background storage is some kind of static memory (e. g. flash), which keeps its state without consuming energy. The sizes of these memories are minimized in the embedded tools. But on the other hand, embedded systems contain more and more programs and data. To save space it is a good solution to compress programs/data.

This compression is usually done by a compressed file system. Such a file system logically stores fixed sized blocks, but physically stores them compressed in variable length blocks, logically increasing the capacity of the storage device. In most cases the used compression algorithm is a general one, which can compress programs and data with the same efficiency (e. g. gzip).

But the blocks to be compressed can be classified. There are blocks containing text data, binary data or program code. File systems compress all these blocks with the same compression method. However special compression methods could be used to compress different kinds of blocks. The size of a block compressed by its special compressor is usually smaller than the size of the same block compressed by the general compressor. Thus compressing each blocks with its special compressor makes the compressed blocks smaller and saves space.

To take advantage of this size difference, we created a compression algorithm for binary program code. We aimed this algorithm on the ARM instruction set, which is the most frequent code set in the embedded systems. The implementation is called *ARMlib*. However the method can be applied for other binary codes that satisfy some requirements (e. g. well separatable

instruction/parameter/condition bits).

There are lots of code compression methods introduced in the literature, many of them are overviewed in a survey from Árpád Beszéd et. al. [1]. These methods usually can be separated into a model and a coder part. Our method also has these parts: based on previous works of C. W. Fraser [2] and M. Garofalakis et. al. [3] we used special decision trees as the model and an arithmetic coder [4] as the coder.

The method was implemented and tested on real iPAQ machines with a specially modified file system. The original JFFS2 [5] file system was designed for flash devices and included zlib compression. This was modified to handle both ARMLib and zlib compressed blocks at the same time. This file system modification itself is not subject of this paper (it was presented on *CSCS2004* [6]), but it was necessary for trying our compression method in a real environment.

The compression efficiency of JFFS2 was improved with ARMLib: on an image with about 25MB uncompressed size, 1.7 - 2.6MB were saved based on the parameters of the algorithm. This means that the compressed image were 12.6 - 19.4% smaller when ARMLib was used than the original (zlib-only) compressed image, thus this amount of flash memory could have been saved.

The drawback of ARMLib is its speed. The boot time of the iPAQ when ARMLib was used were 2-4 times the original boot time (when only zlib was used). The slowdown is hard to be recognized by a regular user in other ways, except in some extreme cases.

The main contributions of this paper are: We introduce a decision tree

based code compression method which combines Fraser's idea of using decision trees as compression models [2] and the effective tree construction algorithm of Garofalakis et. al. [3], which were originally designed for building trees for classification. We applied this method in a real environment [6], namely the JFFS2 file system in the Familiar Linux distribution were modified to use both zlib and our method on an iPAQ machine. (The method is usable on other binary machine codes too. We tried it for Thumb where the method saved 3.5-8.5% of the uncompressed size, which means 9-20% save relative to zlib-compressed size. However these ratios are for pure thumb code compression and not for system-wide use due to our lack of a Thumb hardware unit.)

The rest of the paper is organized as follows: In Section 2 some terms of compression and decision trees are introduced, in Section 3 some existing code compression and decision tree building methods are described, in Section 4 our method is introduced. In Section 5 an existing implementation on iPAQ machines is introduced and in Section 6 the results are presented. In the last section the results are summarized and some issues of further improvement are mentioned.

2 Background

In this section we briefly overview what compression means and how it usually works, then we shortly describe decision trees.

2.1 Compression

The term “compression” can be defined as in [1]: “storing data in a format that requires less space than usual”. In other words: represent some data in a form that is smaller than the original representation. The term “decompression” means the inverse of compression, thus restoring the data in the original form.

2.1.1 Theoretical Background

The theory behind compression is based on results of *information theory*. In this section we will review some terms of information theory.

We define the input of a compression method as a sequence of input symbols. These symbols can be the bits or bytes of the input as well as more complex entities. These entities are usually called tokens. The input sequence may contain values from a fixed set of symbols (token values). The basic idea of most compression algorithm is to assign a *code* to each symbol in the way that the sequence of the codes will be shorter than the sequence of the symbols.

Each symbol in the input has a *probability* value. By giving a shorter code to a more frequent symbol and a longer code to a less frequent one, the overall size of the output sequence will be smaller than the size of the input sequence. Most compression methods utilizes this idea to produce smaller output sequence.

Formally: Let $A = \{a_1, a_2, \dots, a_N\}$ be a set of symbols. Let $X = x_1, x_2, \dots, x_m$ be a sequence with $x_i \in A, (i = 1..m)$. Now, all $a \in A$ has a $P_X(a) \geq 0$ probability in the sequence X , with $\sum_{a \in A} P_X(a) = 1$.

Now, we can compute the *information content* of X . The less information

is stored in X , the shorter encoded sequence can represent it. The information content can be measured by the entropy of X using the following formula:

$$H_A(X) = - \sum_{a \in A} P_X(a) \log_2 P_X(a)$$

This formula gives the minimum average number of bits required to encode one symbol in the sequence X . If $H_A(X)$ is multiplied by m , we get the theoretically minimal size (in bits) of the encoded sequence of X .

2.1.2 Compression model

As mentioned above, most of the compression methods are based on the information content of the input sequence, by assigning shorter codes to more frequent tokens. Thus the compression method can be separated into two parts: (1) gathering information about the input sequence and (2) assigning suitable codes to the tokens. The first component is called *modeller*, the second component is called *coder* and it uses the model provided by the modeller during the code assigning process. Most compression methods contain separate modeller and coder, although these are usually fine-tuned on each other. Many modellers can be used with the same coder and vice versa, thus these two can be treated as separate topics of research.

The goal of the model is to provide a good probability distribution on the next token in the input: for all token values tell what is the chance of the event that the next token has the said value. Modeling can mean almost anything: the simplest model is the probability distribution of the token values in the input,

but the probability distribution provided by the model may change from token to token. For example, the values provided by the model may depend on the value of the last token.

The coders also vary. There are coders that directly assign codeword to each individual input token value (e.g. Huffman) and there are coders that assign a final codeword to a sequence of tokens (e.g. arithmetic). The *decoder* implements the inverse of the coder (it restores the tokens from the codes), however it uses the same model that were used by the coder.

2.2 Code compression

Code compression covers the compression of almost any form of a program, including intermediate representation or binary program code but excluding source code (which is rather a special case text compression than code compression). Although all these forms are considered “code” there can be many differences between them, e.g. the IR code can be some kind of tree (e.g. AST - Abstract Syntax Tree) while binary code is a sequence of machine instruction.

2.3 Decision trees

Decision trees are used for storing information about a set of objects. More precisely, decision trees are trees that provide some information about the target-attribute of an object using some other attributes of it. A usual application of them is the classification of objects, where the object is put into one of the given classes based on some properties (attributes) of the object. [7, 8]

In Figure 1 the objects are given by their attributes and assigned with a

class (+ or -). The decision tree on this figure encodes this classification.

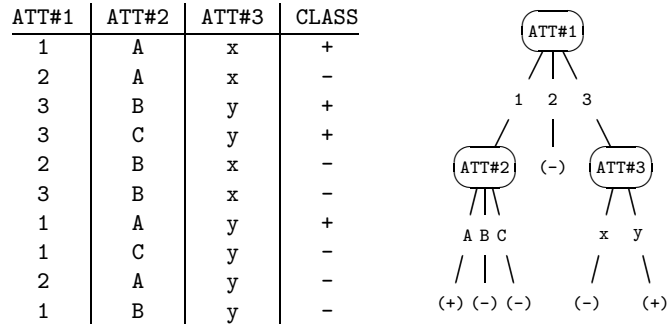


Figure 1: A simple decision tree

The tree contains an expression with attributes (*predictors*) in each of its internal nodes. Such an expression can be the attribute itself (as in the example) or the comparison of the attribute to one of its possible values or more complicated expressions too. An outgoing edge is assigned to each of the possible results of the expression in the node. These edges end in subtrees. In the leaves of the decision tree the information (target-attribute) is stored that corresponds to the decisions made on the route from the root to the actual leaf.

To extract the information about an object (which is given by its attribute values) from a decision tree, first evaluate the decision in the root using the attribute values of the given object, then check the end of the edge assigned with the result. If there is a subtree then repeat the process on that. When a leaf is reached it will contain the required information.

Decision trees are mainly used because they can be automatically built. Tree construction algorithms require only a great number of examples (where the target-attributes are also known) to automatically create a decision tree.

The most known decision tree building algorithms are ID3 [7] and C4.5 [8].

3 Previous works

In this section previous works on binary code compression and decision tree building are described, showing their advantages and disadvantages.

3.1 Code compression methods

Code compression is used for various reasons. The goal can be the reduction of energy consumption which can be result of the smaller stored size and/or the smaller data size send through the channels between hardware elements. Sometimes compression has no other reason than saving space thus improving the capacity of the storage device. For any reasons, code compression methods can be very different.

We concentrate on binary code compression although we borrow some concepts from other code compression as well.

Benini et al. created a *transparent* compression method [9], where the compressed code is decoded by a hardware unit between the memory and the CPU. The decoder unit is a simple table that contains the most frequent 255 codes, assigned to a byte-long codeword. The 256th value is an escape character to state an uncoded instruction. The most frequent codes come from statistics and the codewords are assigned using Minimal Hamming Distance.

Wolfe et al. described a *block-based* compression method usable on hardware with memory cache [10]. The cache pages are stored in the memory in com-

pressed form and decompressed into the cache when a cache miss occurs. The relation between the in-cache and real memory addresses are resolved using a Line Allocation Table (LAT). Many coders were tried, and a version of Huffman coding were proposed by the authors to be used.

Breternitz and Smith enhanced the previous method [11] and eliminated the LAT. The memory addresses in the code are modified to contain a cache page address and an offset of the target instruction. This solution released more problems, for example when the program runs through a cache page boundary without a jump. (This was later solved by automatically generated jump instructions.) The best compression ratio they achieved was 0.56.

Laketsas et. al. improved Wolfe's method [12]. They proposed the decomposition of RISC instructions into *streams*, thus the different parts of the instructions (e.g. operation code, target register) are encoded in different sequences. They also proposed arithmetic coding. Their models were Markov models and dictionaries. Their average compression ratios were 0.5 – 0.7.

The work of Lefturgy et. al. is similar to Wolfe's solution [13], but they assign codewords not only to single instructions but to instruction sequences too. The used model is a dictionary, the codewords have fixed size. The CodePack method of IBM [14] is similar to this but much more complicated. The 32 bits long instructions are divided into two 16 bits long parts and then these are encoded with a variable-length coding. The decoder is also a hardware unit, but the software-based decoding was also studied. The compression ratios for these methods were about 0.6.

The methods above are designed for hardware decompression. Their main task was to reduce the energy consumption of the embedded systems, thus the simplicity was more important than the greater compression ratio. Their models are usually based on some kind of a dictionary. In our case the energy consumption is secondary, the decompression can be done by software and compression ratio is more important, thus our model can be a more complex one.

In [2] Fraser creates a model for code compression using machine learning methods. The coder and decoder are not the subject of his paper; he focused on finding the relevant statistical information existing in the code to be compressed and automatically extracting it.

Fraser works on an intermediate representation (*IR*) not on binary code. The information he extracts from this representation are probability distributions usable by his coder. He stores these distributions in the “leaves” of this decision tree like model. To do this predictors are required that describe the context of the actual token. Fraser used the last 10-20 token values before the actual token (so-called “Markov” predictors) as predictors, and some computed predictors like stack depth. He used a simple tree building algorithm to infer his model automatically: it gets a large number of internal representation code and builds the tree for this data set. Fraser reduced the size of the model in the way he made a DAG (*directed acyclic graph*) from the tree by merging similar leaves. The presented results showed a 0.19 compression ratio on IR code. But these results did not count the size of the model, which can be very large.

3.2 Decision tree building

Decision trees can be automatically generated on large and representative *training data sets* and thus it is easy to use them. One of the best known tree building algorithms is ID3 [7]. It is based on entropy gain. Let X be a set of objects, T is the target-attribute, A is a non-target attribute. Let \mathbf{T} and \mathbf{A} denote the set of target and non-target attribute values. Now define the probability $P_X(t)$ and set $X_{A,a}$ for all $t \in \mathbf{T}$ and $a \in \mathbf{A}$ as:

$$P_X(t) = \frac{|\{x|x \in X, T(x) = t\}|}{|X|},$$

$$X_{A,a} = \{x|x \in X, A(x) = a\}.$$

Now the entropy gain of attribute A on the set of objects X using the formula $H_S(X)$ introduced in Section 2.1.1 is:

$$EG_A(X) = |X|H_{\mathbf{T}}(X) - \sum_{a \in \mathbf{A}} (|X_{A,a}| H_{\mathbf{T}}(X_{A,a})).$$

Tree building works as follows: a set of objects (X) whose attributes are known is assigned with the root; the *best attribute* (A) which produces the highest entropy gain ($EG_A(X)$) is then selected; the set of objects is split into subsets based on the best attribute ($X_{A,a}$) and each subset is assigned with a child of the root. The algorithm is then invoked for the children. If all attributes produce negative entropy gain, the node becomes a leaf encoding information of the target attribute values of the objects.

A problem can be the *overfitting*. If a tree trained on a noisy training set (that contains many errors) exactly fits the training set, then the answers the tree gives will also contain errors. Two kinds of methods can solve this problem.

The first is applied during tree building. When a new node is examined whether to be expanded or not, some *stopping criterions* are checked. These may depend on various properties of the node and the assigned training set. If any of these criterion becomes true, the node is not expanded.

The other method applied on the fully built tree. The subtrees are examined and replaced with leaves if necessary. This is called *pruning*. It gives more accurate result than the previous method, because all subtrees are well known during pruning, while in the previous method the subtrees are not known before the stopping decision. On the other hand, it may happen that a great subtree is built first and then dropped during pruning.

An enhancement of ID3 algorithm is C4.5 [8]. It contains a pruning algorithm that works on a rule-set derived from the tree. Each leaves are represented with the decisions made in the tree from the root to that leaf, then these rules are merged and sorted, and the resulted rule list is used.

In [3] Garofalakis et. al. introduced some methods for the efficient building of decision trees. Their trees are built for classifications, and used the MDL (Minimal Description Length) measure for pruning.

Their pruning method works on the tree. Let R be the root of the (sub)tree to be pruned, C_{leaf} is the MDL cost of a leaf that encodes the information of the training set assigned with R . The C_{node} cost of the (sub)tree can be recursively

computed. The pruning algorithm works as follows:

```

Procedure Prune( $R$ )

 $C_{leaf} :=$  cost of a leaf at node  $R$ 

 $C_{node} := \sum_{i \in \text{Children}(R)} \text{Prune}(R_i) +$ 
           cost of encoding internal node  $R$ 

If  $C_{leaf} \leq C_{node}$  Then

    Replace the subtree rooted at  $R$  with a leaf

    Return  $C_{leaf}$ 

Else

    Return  $C_{node}$ 

```

In [3] two enhanced versions of this methods were described too. In the first the size of the tree (number of nodes in it) can be maximized by replacing subtrees of the pruned tree with leaves using dynamic programming method. In the second the minimal precision of the tree can be set. This uses the previous method initially setting the bound on tree size to 1, then increasing the bound by 1 until the precision of the tree reaches the required value.

In the same paper a third method was also described that can use the bounds set for tree size during tree building phase (using *branch and bound* technique).

4 ARMLib - Code compression for ARM

In this section our compression method is described through *ARMLib*.

ARMLib is a function library that contains functions for ARM binary code compression. ARM binary code was selected because ARM instruction set [15]

is the most popular instruction set in the embedded systems.

Based on the work of C. W. Fraser [2], our model in ARMLib is a decision tree whose leaves contain probability distributions. For an application of ARMLib the decision tree is pre-built (*trained*) on a host machine using a great number of binary code blocks (*training set*).

The building of the decision tree is based on the works of C. W. Fraser [2] and M. Garofalakis et al. [3]. The building is automatic: it optimizes for image size, thus not only the compression ratio is minimized (because it would result in a very precise but huge model) but the size of the tree itself is also counted. These models provide best performance (i.e. greatest compression) on code sequences that are very similar to the code they were trained on. By setting the parameters of the training, the efficiency and speed of the compression can be varied.

An arithmetic coder based on [4] is used as the coder in ARMLib.

4.1 The parts of ARMLib

ARMLib consists of 3 main parts: model generator and compressor and decompressor modules. The model generator creates the model from a great number of example programs, the training data set. The compressor is to compress, the decompressor is to decompress the binary ARM code; both use the previously created model. The tokenizer and detokenizer algorithms are also the part of the ARMLib. Tokenizer transforms the raw ARM binary code into the sequence of tokens usable by the modeller and coder, and detokenizer transforms the se-

quence of tokens into ARM binary codes after decompression. Figure 2. shows how the parts of the ARMLib work together.

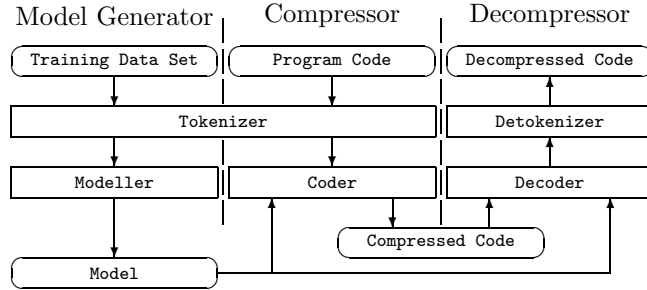


Figure 2: The structure of ARMLib

4.2 Pre- and post-processing

The input of model generation and compression and the output of decompression are ARM binary code. But the presented methods work on a sequence of tokens, so we need to transform the raw 32 bits long ARM codes into tokens and the tokens back to ARM codes. These transformations are done by the tokenizer/detokenizer modules.

ARM instructions are 32 bits long that would be too large for our method, so an instruction must be split into more and smaller tokens. In ARMLib the instructions are fragmented into 8 tokens which are 4 bits long each. So the first step of the model generation and the compression is to create a sequence of tokens from the sequence of ARM instructions, which is done by the tokenizer module. The last step of the decompression is to create the ARM instruction sequence from the token sequence. This is done by the detokenizer.

The 4 bits long tokens were chosen for two reasons: the first is that the ARM instructions functionally consist of parts whose length in bits are a multiply of 4.

The second reason is that 4 bits means 16 token values and this small number is very kind for the model. The order of the tokens in the token stream of an instruction were determined and fixed after some preliminary measurements had been made on test data. This was necessary because the value of token T_a of an instruction may imply some values of token T_b of the same instruction, but T_b has no such influence on T_a . In this case T_a must precede T_b .

The method can be easily applied on any architecture whose instructions have a fixed length. And if the bits of the instructions can be assigned into functional groups, the method might be as successful as on ARM.

4.3 The model

ARMLib uses special decision trees. The decision trees of ARMLib have the following traits: the objects are the tokens; the predictors are the last 16 tokens before the actual token and the type of the actual token (e.g. destination register); the information stored in the leaves are the probability distributions of the possible token values.

In [2] Fraser uses many kinds of *computed* and *reduced* predictors. In our case predictors that required complex computations did not improve the efficiency of the method so much that they have been worth being used. The "last n token" type predictors were combined with only one computed predictor, which states the type of the actual token. The last 8×2 tokens were also used as predictors, thus the last 2 values of each of the 8 types of tokens were stored as predictor. These are easy to compute, and the possibility of decomposing the input into

many streams still remains, however it is up to the tree building algorithm to decide if it is worth using the type predictor at the root of the tree or not. (Using it corresponds to the stream separation.)

In the internal nodes of the tree many types of decisions can be used. The simplest decision is a predictor itself (as in ID3 and C4.5 trees). A node with this type of decision has as many children as many values the predictor has (8 in case of the type predictor, 16 otherwise). Another type of decision is the comparison of a predictor with one of its possible values. This can be an equality or a less-than comparison. A node with this kind of decision has two children: one for the *true* and one for the *false* value.

The MDL measure used by Garofalakis et. al. is not enough to build an optimal tree in this case, because the size of the tree and the size of the compressed code have to be minimized together. To exactly compute this value, three things must be known: the storage method and size of the tree and its subtrees; the code to be compressed; and the coding method and its compression effect on the code. If these three things were known the optimal tree could be built. But there are some problems.

The way the tree is stored is known and the size of the tree or any of its subtrees can be exactly computed, as long as the tree itself is not compressed. But it may happen that the tree itself is also compressed in some way (which is not a crazy idea in a compressed file system). In this case the compressed sizes of the subtrees depend on each other and cannot be computed individually. The code to be compressed is known. Each subtree will compress those tokens that

are in the training set of one of the subtree's leaves. The coding method is also known, but to compute the exact size of the compressed code the coding or at least a fake coding (that computes only length information) must be done. This means too much computations, so we have to make a deal.

The size of a subtree is computed from the uncompressed size of the subtree (which is exactly computed) and a pre-defined *compression-multiplier*. This means a few additional computations and approximates the exact final size well (if the given multiplier is good). The size of the tokens compressed by a subtree is approximated with the sum of the entropies of the training sets associated with the leaves of the subtree. This is also a good approximation if the output of the coder is close to the entropy.

Thus the building of an ARMLib model for a given token sequence works as follows. A training set is made from the token sequence which is not more than assigning the predictor values with each tokens. The tree is built for this training set using the general tree building algorithm. When this is done, the pruning algorithm is invoked for the root with the previously described cost computations. Some kind of parallelisation of building and pruning is implemented in ARMLib. As soon as a subtree is built, the pruning method is invoked for its root immediately. This may reduce the maximal tree size in the memory during the tree building.

4.4 The coder

ARMLib uses arithmetic coder [4]. The concept of the arithmetic coder is that it does not assign codewords, but assigns an interval between 0 and 1 to the tokens. It also assigns an interval to the sequence of tokens (which interval is computed from the intervals of the tokens). The codeword assigned to a sequence of tokens is a random number in the final interval stored with high precision enough for correct decoding.

Arithmetic coder produces a codeword whose length is very close to the entropy. Technically the halving of the size of the actual interval means one more bit in the codeword. This allows the arithmetic coder to virtually assign a fraction of a bit to a token whose probability is greater than 50% (while Huffman coder assigns at least 1 bit to each tokens regardless its probability).

5 An ARMLib implementation

ARMLib were tested in a real environment, on *iPAQ* machines [16]. These machines have *StrongARM* processors, 32MB of memory and 16MB of flash. The operating system used was *Familiar Linux* with a modified *JFFS2* (Journaling Flash File System 2) file system.

The original JFFS2 file system stores the files logically in fixed size blocks. But it is able to compress these blocks in a transparent way when they are stored in the flash, which means that the user of the file system does not see the compression. In this way the capacity of the flash logically increased. The used compression methods are the methods of *zlib*. (See Figure 3a.)

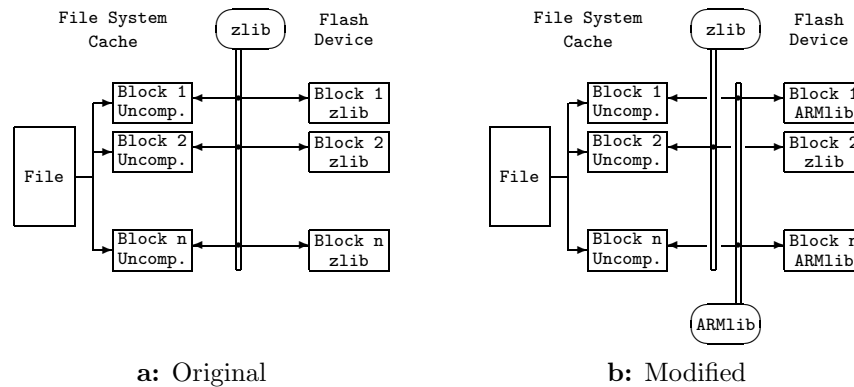


Figure 3: JFFS2

The compression methods of ARMLib were integrated into JFFS2 as can be seen in Figure 3b. The two compressors work side by side in the file system. Each block of the file system stores the information which method it was encoded by, so decompression can be done by the appropriate decompressor.

The model generation is done on a host machine. First the files and directory structure of iPAQ is created. Then the files (programs) that contain ARM code are selected, but some of them that have to be compressed by zlib for some technical issues (e. g. kernel files, ARMLib itself) are excluded from the selection. The selected programs are split into blocks (as in the file system), and the blocks that do not contain ARM code are also dropped. The training set is created from the remained blocks. It is possible to classify the blocks so creating many distinct training sets and build a separate tree for each such set.

After the model generation had been finished the model(s) is(are) added to the directory structure and the binary image is made. During this all blocks are compressed with both ARMLib and zlib, and the smaller variant is kept. The

block stores the identifier of the model it was compressed with, if necessary (if there were more than one model). Model file blocks can only be compressed with zlib.

After the image is ready, it is uploaded into the iPAQ. The bootloader (which knows zlib only) loads the kernel modules including ARMLib. At initialization the file system reserves some memory. The models listed in a configuration file are loaded into this area, and the rest of the reserved memory becomes the file cache. After this ARMLib is ready to work.

6 Results

Measurements were done in two kind of environments. In the test environment the modeller, coder and decoder were programs like gzip, and all were run on a host machine. The second environment was the JFFS2. Here modeller and coder ran on a host machine, but coder and decoder functions are used on the iPAQ. This environment were used to test the practical usability of the method. The uncompressed size of the image of the iPAQ were 25MB.

In both cases the compression were block-based. We used 4KB long blocks. Only those blocks were used for training that contained ARM binary code.

6.1 Size results

As can be seen on Figure 4. the larger the input is the better the compression ratio is. The results were done in the test environment by training a tree for each test object and compressing its blocks with the tree trained on it. The

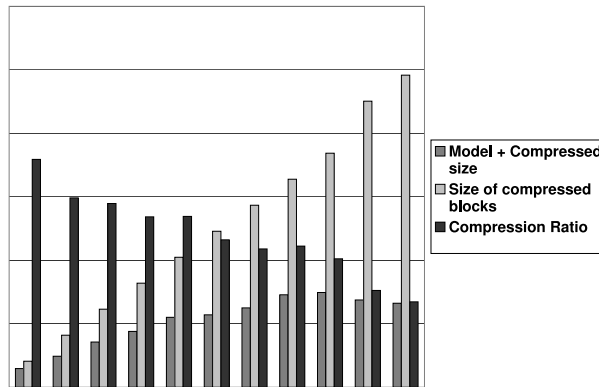


Figure 4: Compressed sizes of different inputs

test objects contained 1-12 blocks. (On real environment we did not make such measurement, but on a smaller image with about 18MB uncompressed size ARMLib produced about one percent worse results.)

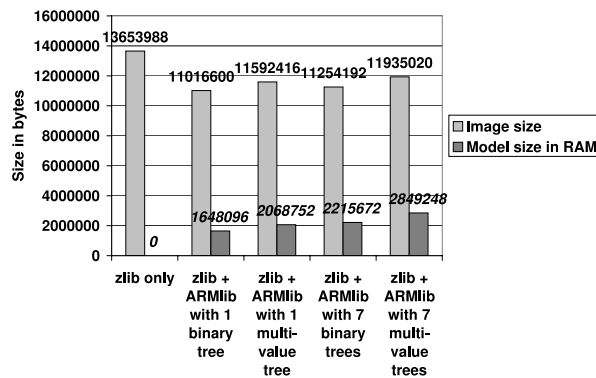


Figure 5: Compressed sizes

Different image sizes can be seen on Figure 5. The size reductions relative to the image compressed with zlib only are between 12.6% and 19.3%. Trees with binary nodes only are better than trees with multi-value decisions (because multi-value nodes can be substituted with binary nodes, but it is not true on the

contrary). A small image size enlargement is observed when the one model was replaced with 7 smaller models. (These were trained on 7 randomly separated distinct parts of the original training set.) The results showed that one tree is more effective in compression size (at least when the split of the data set is random) than the 7 smaller models together, but the (de)compression speed is better for the 7 models than for only one.

The same figure shows the sizes of the models in the memory. Binary models are smaller, and the 7 smaller models requires more memory together.

6.2 Speed results

The speed of the compression is worse than gzip's speed. Compression and decompression in ARMLib requires almost the same amount of time, but this makes decompression very slow relative to zlib (in zlib, decompression is 10 times faster than compression). Fortunately, Linux caches the file system. This means that not all accesses to the block provokes compression or decompression, usually the first and the last accesses only. This improves the average speed of the file system, and makes ARMLib slowness almost undetectable to a regular user.

The boot time (from turn on to the set up of the system) is 2.3-3 times of the original boot time (Figure 6), but considering the save in size (12.6-19.3%) this seems to be a good deal.

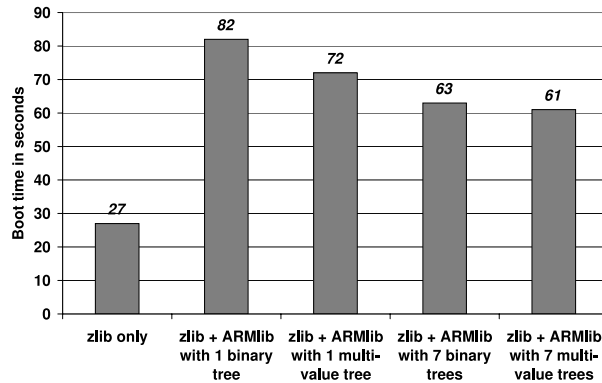


Figure 6: Boot times

7 Conclusion and future work

Our goal was to create a code compression method which compresses ARM binary codes better than current algorithms do. We combined a method that used decision trees as the compression models [2] and an efficient tree building algorithm [3]. In addition the tree building and pruning phases were also combined within the method. The method were aimed for ARM instruction set, but it can be applied for many instruction sets that satisfy some requirements.

The implementations of the functions of the method were collected in a library called *ARMLib*. *ARMLib* used an arithmetic coder [4]. *ARMLib* were used in the JFFS2 file system [5] on iPAQ machines. This modified file system [6] used either zlib or *ARMLib* compression on a block. The use of *ARMLib* reduced the image size by 12.6-19.3% depending on the parameters of tree building. The method can be used in situations when the size is more important than speed.

Possible improvement can be the automatic sorting of the input. The implemented method was measured with one tree built on the full image, and with

7 trees built on 7 separated parts of the image. But these parts were randomly separated. A good classification of the image blocks would improve not only the performance but the compression ratio of ARMLib too.

The measurement functions could also be modified to estimate the speed of the compression (perhaps by utilizing the fact that compression speed is proportional to the depth of the tree). The importance of compression ratio and compression speed then could also be scaled by setting a parameter.

References

- [1] Beszédés, Á., Ferenc, R., Gyimóthy, T., Dolenc, A., Karsisto, K.: Survey of Code-Size Reduction Methods. *ACM Computing Surveys*, 2003, **35**, 223–267
- [2] Fraser, C. W.: Automatic Inference of Models for Statistical Code Compression. In *Proc. PLDI*, 1999, 242-246.
- [3] Garofalakis, M., Hyun, D., Rastogi, R., Shim, K.: Efficient Algorithms for Constructing Decision Trees with Constraints. Bell Laboratories, 2000, <http://www.bell-labs.com/user/rastogi/sigkdd2000.ps>.
- [4] Witten, I. H., Neal, R. M., Cleary, J. G.: Arithmetic Coding for Data Compression. *Communications of the ACM*, 1987, **30**, 520–540.
- [5] Woodhause, D.: JFFS: The Journaling Flash File System. *Ottawa Linux Symposium*, 2001,

<http://www.linuxsymposium.org/2001/proceedings/pdf/jffs2.pdf>.

- [6] Havasi, F.: Increasing compression performance of block based file systems. Presentation on *CSCS2004*, 2004.
- [7] Quinlan, J. R.: Induction of decision trees. In *Machine Learning*, 1986, **1**, 81–106.
- [8] Quinlan, J. R.: *C4.5: Programs for Machine Learning*. Morgan Kaufmann, San Mateo, 1993.
- [9] Benini, L., Macii, A., Macii, E., Poncio, M.: Selective instruction compression for memory energy reduction in embedded systems. In *Proceedings of ISLPED '99*, 1999, 206–211.
- [10] Beneš, M., Nowick, S. M., Wolfe, A.: A fast asynchronous decompression circuit for embedded processors. In *Proceedings of ASYNC98*, 1998, 43–56.
- [11] Breternitz, M. Jr., Smith, R.: Enhanced compression techniques to simplify program decompression and execution. In *Proceedings of ICCD'97*, 1997, 170–176.
- [12] Laketsas, H., Henkel, J., Wolf, W.: Code compression for low power embedded system design. In *Proceedings of ACM/IEEE DAC*, 2000, 430–439.
- [13] Lefturgy, C. R., Piccininni, E., Mudge, T. N.: Reducing code size with run-time decompression. In *Proceedings of HPCA*, 2000, 218–227.
- [14] IBM. Codepack: PowerPC Code Compression Utility User's Manual Version 3.0. International Business Machines (IBM) Corporation, 1998.

- [15] ARM Instruction Set Quick Reference Card v2.1. Online documentation:
http://www.arm.com/pdfs/QRC0001H.rvct_v2.1.arm.pdf.
- [16] iPAQ H3000 Pocket PC Reference Guide. Online documentation:
[http://h200005.www2.hp.com/bc/docs/support/UCR/
SupportManual/TPM_177711-001/TPM_177711-001.pdf](http://h200005.www2.hp.com/bc/docs/support/UCR/SupportManual/TPM_177711-001/TPM_177711-001.pdf).