

ÉCOLE NORMALE SUPÉRIEURE DE PARIS
DÉPARTEMENT D'INFORMATIQUE

MASTER 1 INTERNSHIP REPORT

Year 2017–2018

Faster CakeML compilation with a verified linear scan register allocator

Théophile WALLEZ

Supervised by Magnus O. MYREEN

At Chalmers University, Göteborg, Sweden

Acknowledgments

I would like to express my gratitude to Magnus O. MYREEN for supervising me during this internship, for always being available when I had questions, and for allowing me to work on such an interesting project.

I would like to thank Oskar ABRAHAMSSON, Andreas LÖÖW, Johannes ÅMAN POHJOLA and Alejandro GÓMEZ-LONDOÑO for their answers to my countless questions about HOL4.

I would also like to thank Yong Kiam TAN for his precious help on the integration of my project in the current codebase.

Thanks to David REBOULLET for the numerous coffee-breaks which allowed me to practice rubber-duck debugging on a real human.

Last but not least, thanks to Jean-Christophe FILLIÂTRE for his amazing compilation course, and for telling me about the existence of the CakeML compiler.

Contents

1	Introduction	1
1.1	The HOL4 theorem prover	1
1.2	The CakeML verified compiler	1
1.3	The internship goal	1
2	First work on CakeML: improving the constant folding	1
3	Preliminaries	2
3.1	What is register allocation	2
3.2	The linear scan register allocation algorithm	3
3.3	HOL4's standard library	4
3.4	CakeML's current register allocator	5
4	Implementation and correctness of the algorithm in HOL4	6
4.1	Step 1: remove cutsets	6
4.2	Step 2: get the liveness intervals	8
4.2.1	The <code>get_intervals</code> function	8
4.2.2	The <code>get_intervals_withlive</code> function	10
4.2.3	Proving that the two functions compute the same thing	11
4.3	Step 3: do the allocation	12
4.3.1	Modifications to the original algorithm	12
4.3.2	The state and invariants used in the linear scan algorithm	14
4.3.3	Implementation of the linear scan algorithm	15
5	Evaluation	18
6	Future work	19
6.1	Split the SSA-form and the calling conventions	19
6.2	Optimise the constant of the algorithm	19
6.3	Remove the <code>fix_domination</code> function	19
6.4	Use the fact that we work on a SSA-like AST	19
7	Conclusion	19

1 Introduction

1.1 The HOL4 theorem prover

HOL4 is an interactive theorem prover for Higher Order Logic (HOL). It is a descendant of LCF, and is similar to Isabelle. Internally, all theorems in HOL4 are proved using basic inference rules [9]. The absence of dependent types in HOL4 make convenient proof automation available (such as rewriting and first-order provers).

1.2 The CakeML verified compiler

CakeML is a language based on a subset of Standard ML, with a compiler is written in HOL4. The formal semantics of CakeML is defined in a functional big-step style [7].

The CakeML compiler's backend [10] transforms a source AST (which lacks type annotations) to machine code for one of the five architectures: ARMv6, ARMv8, x86-64, MIPS-64 and RISC-V. The backend has been proved to produce code that has the same behavior as the source program.

The compiler has two frontends: the first one is a traditional parser from CakeML source [5] with a type inferencer [11], both of which have been proved to be sound and complete. The second frontend is a translator from HOL functions to CakeML AST [6, 3], which produces a proof that the generated AST has the same behavior as the HOL function. It can produce code that is stateful and performs I/O.

The compiler can bootstrap itself [5]: the first frontend combined with the backend is a HOL function which transforms a CakeML source code to machine code. This function is given to the second frontend combined with the backend to become machine code which transforms CakeML source code to machine code. Every step in this process is verified end-to-end with HOL4, therefore making CakeML a fully verified compiler.

1.3 The internship goal

The algorithm used in the register allocator of CakeML is the iterated register coalescing algorithm [1], which can be slow: its complexity is at least $\Omega(n^2)$ where n is the number of registers in the program. In fact, most of the compilation time is spent in the register allocation: it is the slowest part of the compiler.

For several applications (JIT, REPL, running on a verified processor on a FPGA) it would be very useful to have a short compilation time, and implementing the linear scan register allocator [8], would be useful.

The goal of my internship was to implement and verify a linear scan register allocator, and integrate it in the CakeML code base.

2 First work on CakeML: improving the constant folding

CakeML is a big and complex project, and starting by doing the linear scan register allocator first seemed a bit too ambitious. Therefore, I started by improving the code for constant propagation. I noticed that the algorithm was able to fold expressions like $x + (1 + 2)$ to $x + 3$, but was failing to fold expressions like $1 + (x + 2)$ or even $(x + 1) + 2$.

The constant folding is done by using smart constructors. A smart constructor `SmartOp op exp1 exp2` is in general equal to `Op op exp1 exp2` but can be smarter in special cases: for example `SmartOp Add exp1 (Number 0)` can be equal to `exp1`.

The smart constructor was previously implemented with code similar to the one below:

```
SmartOp op exp1 exp2 =
  case (exp1,exp2) of
  | (Number n1,Number n2) =>
    if op = Add then Number (n1 + n2)
    else if op = Sub then Number (n1 - n2)
    else if op = Mul then Number (n1 × n2)
    else Op op exp1 exp2
  | _ => Op op exp1 exp2
```

We can see that this code can't fold expressions like $(x + 1) + 2$.

When implementing the smart constructor, we need to be careful not to change the order of evaluation, since expressions can have side-effects executed. The smart constructor does not need to be careful of overflows since CakeML supports big integers by default.

The algorithm applies this smart-constructor in a bottom-up way: the arguments of the smart-constructor were also constructed using the smart-constructor hence the algorithm can fold expressions like $x + (1 + (2 + (3 + 4)))$ into $x + 10$.

Previously, `SmartOp op exp1 exp2` only checked if both `exp1` and `exp2` were `Numbers`: I extended it to check if `exp1` and `exp2` have the form `Number n`, `Op op (Number n) exp` or `Op op exp (Number n)`.

It is possible to remove the case `Op op exp (Number n)` by noticing that $exp + n = n + exp$, $exp \times n = n \times exp$ and $exp - n = -n + exp$. After doing these rewrites it is sufficient to check if `exp1` or `exp2` have the form `Number n` or `Op op (Number n) exp`.

Then, by enumerating all the cases it is possible to see that every combination of `Add` and `Sub` are foldable without changing the order of evaluation (for example, $(n_1 - x_1) - (n_2 - x_2) = (n_1 - n_2) - (x_1 - x_2)$)

I implemented this improvement in the intermediate language BVL, a functional language without closures, and I updated the proofs.

The pull request is at the following url: <https://github.com/CakeML/cakeml/pull/485>

3 Preliminaries

3.1 What is register allocation

During optimisation passes, compilers generally use a model with an infinite number of virtual registers. However, computers have a small fixed number of physical registers, hence we must map the virtual registers to use this small number of physical registers.

Since it might not be possible to fit all virtual registers on a small number of physical registers, we have a mechanism to "spill" some virtual registers to the stack. Therefore, we can still use an infinite number of physical registers, but we should avoid spilling virtual registers as much as possible since it degrades the performance of the code produced.

Two virtual registers must not be allocated to the same physical register when their value might be useful at the same time in the program. To avoid overlap, we compute the positions in the program where each virtual register might hold a useful value for the rest of the program execution.

At points the program where a virtual register might hold a useful value for the rest of the program, we say that the virtual register is live. Notice the word "might": without this word, computing the set of live registers at the position of a program is undecidable, so we compute a superset of the real live registers.

To compute the set of live registers at each position of the program, we could use Kildall's algorithm [4] which is the classical algorithm used to perform liveness analysis. However, in CakeML, the control-flow graphs have no cycles because loops are implemented as tail-recursive functions. In fact, we will see in

Section 3.4 that control-flow graphs in CakeML have a very simple shape: we don't need the complexity of Kildall's algorithm and the liveness analysis is more simple.

3.2 The linear scan register allocation algorithm

In this section, we present the linear scan register allocation algorithm [8], which is orders of magnitude faster than the iterated register coalescing algorithm [1] and produces machine code of only slightly lower quality.

First, the liveness analysis is performed, and the control-flow graph is flattened. After that, holes in the liveness intervals are removed: this is the key approximation made by the algorithm to achieve its speed.

For example, in the following example code:

```

0 /* Live = {}                                     */
  a ← ...
1 /* Live = {a}                                   */
  b ← ...
2 /* Live = {a,b}                                 */
  if ... :
3   /* Live = {b}                                 */
   c ← b
4   /* Live = {c}                                 */
  else:
5   /* Live = {a}                                 */
   c ← a
6   /* Live = {c}                                 */
7 /* Live = {c}                                   */
  print(c)
8 /* Live = {}                                     */

```

The liveness intervals are:

$\text{Live}(a) = \{1, 2, 5\} \subset [1, 5]$

$\text{Live}(b) = \{2, 3\} \subset [2, 3]$

$\text{Live}(c) = \{4, 6, 7\} \subset [4, 7]$

Each register has a liveness interval, and two registers interfere when their liveness intervals intersect.

The algorithm iterates over all registers, sorted by increasing beginning of liveness intervals. During the algorithm, two pieces of information are kept track of:

- the list of active intervals at the current point of the program: intervals that contain the start point of the interval we are currently processing
- a color pool, which contains the physical registers that are not used by registers in the active list

The colors are assigned greedily. When there are no colors available, the algorithm spills to the stack the register with the maximum end of liveness interval.

The algorithm produces a correct coloration, because if two intervals intersect, then one contains the beginning of the other, and they cannot get assigned the same color because when processing the second interval, the first one will be in the active list hence its color won't be in the colorpool.

More precisely, the algorithm works as follow:

```

Function LinearScanRegisterAllocation()
| active = {}
| foreach live interval i, in order of increasing startpoint[i] :
|   ExpireOldIntervals(i)
|   if colorpool is not empty :
|     col ← color removed from colorpool
|     color[i] ← col
|     add i to active
|   else:
|     SpillInterval(i)
end
Function ExpireOldIntervals(i)
| foreach live interval j in active such that endpoint[j] < startpoint[i] :
|   remove j from active
|   add color[j] to the colorpool
end
Function SpillInterval(i)
| tospill ← argmax(endpoint[j] for j ∈ active)
| if endpoint[tospill] > endpoint[i] :
|   color[i] ← color[tospill]
|   color[tospill] ← new stack location
|   remove tospill from active
|   add i to active
| else:
|   color[i] ← new stack location
end

```

3.3 HOL4's standard library

In this report, we will use several functions and datatypes defined in HOL4's standard library.

`unit` is the type with one element named `()`, `num` represents a natural number (including zero), `int` is an integer. They both are infinite, there is no maximal value.

`α list` is a list of objects with type `α`. `mem x l` is true iff `x` is in the list `l`. `every P l` is true iff the predicate `P` is true on every element of `l`. `el n l` is the `n`th element of the list `l`. `map` and `filter` behave like the usual map and filter functions in other functional languages. `x::xs` is the list `xs` prepended with `x`. `l1 ++ l2` is the concatenation of `l1` and `l2`. `all_distinct l` is true iff every element of `l` is unique. `length l` is the length of `l`.

`α × β` represents a pair. The first and second element can be accessed with the `fst` and `snd` functions.

`α option` is either `None` or `Some v`. The value can be extracted by the partial function `THE` or the total function `the`: `THE (Some v) = v`, `the d (Some v) = v` and `the d None = d`.

`α num_map` represents a function from a subset of `num` to `α`, or equivalently a total function from `num` to `α option`. `lookup k s` returns the value associated with `k` in the `α num_map s`. The type of `lookup k s` is `α option`, and `lookup k s = None` means that there is no element associated with `k` in `s`. `insert k v s` add the new mapping `k ↦ v` in `s`. `α num_map` is implemented using a trie on the binary form of numbers, therefore the complexity of `lookup k s` or `insert k v s` is $O(\log_2(k))$. `domain s` is the domain of definition of `s`: `domain s = { k | lookup k s ≠ None }`.

`num_set` represents a subset of `num`. It is actually a type alias for `unit num_map`: `x` is in the subset iff `lookup x s = Some ()`. The empty set is represented by `LN`. `union s1 s2` computes the union of `s1` and

s_2 , difference $s_1 \setminus s_2$ is the set of elements in s_1 but not in s_2 . `domain s` is the set represented by s .

There is some formalisation of sets, they are purely proofs objects and don't result in any computation. `image f set` is the image of `set` by the function f . `injective f set` checks if the function f is injective on `set`.

3.4 CakeML's current register allocator

The CakeML compiler has a lot of intermediate languages, and the register allocation happens close to the end of the compilation, in an intermediate language called `WORDLANG`.

A register allocation algorithm produces a coloration such that two interfering registers don't have the same color. This is checked by the predicate `colouring_ok`.

There is a pre-existing theorem stating that if `colouring_ok` says that a coloration is good, then applying the coloration doesn't change the semantics of the program.

Only the reads and writes of the `WORDLANG` AST are relevant to the register allocation. This information is summarized in the `clash_tree` datatype retrieved with the `get_clash_tree` function.

The function `check_clash_tree` checks if a coloration is compatible with the clash tree. Its prototype is: `check_clash_tree f clashtree live flive` where f is the coloration function, $live$ is the set of live variables after $clashtree$. Two invariants are maintained: $flive$ is f applied to the set $live$, and f is injective on $live$. When the coloration is correct, the function returns `Some (livein, flivein)` where $livein$ is the set of live variables before $clashtree$. When the coloration is not correct, the function returns `None`.

The following theorem says that if `check_clash_tree` says that a coloration is good then `colouring_ok` is always true. `wf_cutsets prog` is a technical hypothesis without any deep meaning (it says that the cutsets of $prog$ are well-formed). `LN` represents the empty set of the `num_set` datastructure.

$$\begin{aligned} &\vdash \text{wf_cutsets } prog \wedge \\ &\quad \text{check_clash_tree } f \text{ (get_clash_tree } prog) \text{ LN LN} = \text{Some } (livein, flivein) \Rightarrow \\ &\quad \text{colouring_ok } f \text{ prog LN} \end{aligned}$$

The `clash_tree` datatype is defined like this:

```
clash_tree =
  Delta (num list) (num list)
  | Set num_set
  | Branch (num_set option) clash_tree clash_tree
  | Seq clash_tree clash_tree
```

`Delta writes reads` represents an instruction that writes to a list of registers and reads from a list of registers.

`Set cutset` represents a cutset, the set of local variables which must be preserved past subroutine calls. It is useful for the garbage collector.

`Branch optcutset ct_1 ct_2` represents an if condition: the two programs are the if program and the else program. `optcutset` is an optional cutset.

`Seq ct_1 ct_2` represent the concatenation of two programs

The correctness theorem of the current register allocator looks like this:

$$\begin{aligned}
&\vdash \text{every } (\lambda(x,y). \text{in_clash_tree } ct \ x \wedge \text{in_clash_tree } ct \ y) \text{ forced} \Rightarrow \\
&\quad \exists \text{spcol livein flivein.} \\
&\quad \text{reg_alloc alg sc k moves ct forced} = \text{Success spcol} \wedge \\
&\quad \text{check_clash_tree (sp_default spcol) ct LN LN} = \text{Some (livein, flivein)} \wedge \\
&\quad (\forall x. \\
&\quad \quad \text{in_clash_tree } ct \ x \Rightarrow \\
&\quad \quad \quad x \in \text{domain spcol} \wedge \text{if is_phy_var } x \text{ then sp_default spcol } x = x \text{ div } 2 \\
&\quad \quad \quad \text{else if is_stack_var } x \text{ then } k \leq \text{sp_default spcol } x \text{ else T}) \wedge \\
&\quad (\forall x. x \in \text{domain spcol} \Rightarrow \text{in_clash_tree } ct \ x) \wedge \\
&\quad \text{every } (\lambda(x,y). \text{sp_default spcol } x = \text{sp_default spcol } y \Rightarrow x = y) \text{ forced}
\end{aligned}$$

Some registers represent physical registers, and shouldn't be allocated to any other register. They are recognised using the predicate `is_phy_var`. It is useful for two things: calling conventions, and some machine instructions which requires specific registers as input (e.g. x86_64's `idivq` uses only `%rax` and `%rdx`).

Some registers should be allocated on the stack: they are recognised using the predicate `is_stack_var`

The list *forced* is a list of pair of registers that should not be allocated to the same physical register. It is a constraint coming from some machine instructions in MIPS, RISC-V and ARMv8, where the input should not be equal to the output.

The first part of the correctness theorem is a function call to `check_clash_tree` which says that the coloration produced is a correct coloration. The second part says that registers that should be allocated on the stack are on the stack, and that physical registers are allocated to themselves. The last part says that the *forced* requirements are satisfied.

To ensure that the new allocator fits nicely in the CakeML codebase, it should try to have a correctness theorem as similar as the above pre-existing one.

4 Implementation and correctness of the algorithm in HOL4

Everything presented in this section is my own work.

The code is available here: https://github.com/CakeML/cakeml/blob/master/compiler/backend/reg_alloc/linear_scanScript.sml

The proofs are available here: https://github.com/CakeML/cakeml/blob/master/compiler/backend/reg_alloc/proofs/linear_scanProofScript.sml

The first step of the algorithm is to remove cutsets, because the linear scan algorithm can't do anything smart with cutsets. Doing this also allows to have a simpler color-checking function, hence simplify the proofs. The second step is to compute the liveness intervals. The third and last step is the actual linear scan algorithm.

4.1 Step 1: remove cutsets

The set of live registers is computed backwards. The reason for this is that when a register is written to, its value might not be useful after (e.g. it might never be read). However, when a register is read, we know that its value is useful earlier in the program.

Using the `clash_tree`, the set of live registers is computed by induction like this:

```

get_live_backward_ct (Delta writes reads) live =
  union (difference live (list_to_numset writes)) (list_to_numset reads)
get_live_backward_ct (Set cutset) live = cutset
get_live_backward_ct (Branch (Some cutset) ct1 ct2) live = cutset
get_live_backward_ct (Branch None ct1 ct2) live =
  union (get_live_backward_ct ct1 live) (get_live_backward_ct ct2 live)
get_live_backward_ct (Seq ct1 ct2) live =
  get_live_backward_ct ct1 (get_live_backward_ct ct2 live)

```

Cutsets provide liveness information which is too precise for the rough liveness intervals that the linear scan algorithm uses. It also makes proofs harder.

Therefore, we simplify `clash_tree` into `live_tree` which does not contain any cutsets, and `Delta` is split in two: `Writes` and `Reads`.

```

live_tree =
  Writes (num list)
  | Reads (num list)
  | Branch live_tree live_tree
  | Seq live_tree live_tree

```

Because this transformation loses information, it means that we can't compute sets of live variables as accurately as before. However, this is not a problem, since we can compute an over-approximation of the set of live variables.

Instead of treating cutsets as a fresh set of live registers, they are added to the set of live registers: we transform cutsets into `Reads`.

In fact, we don't lose much information by doing this because cutsets are already live registers: adding cutsets to the set of live registers is the same as ignoring cutsets.

Therefore, we can transform a `clash_tree` into a `live_tree` like this:

```

get_live_tree (Delta wr rd) =
  Seq (Reads rd) (Writes wr)
get_live_tree (Set cutset) =
  Reads (numset_to_list cutset)
get_live_tree (Branch None ct1 ct2) =
  Branch (get_live_tree ct1) (get_live_tree ct2)
get_live_tree (Branch (Some cutset) ct1 ct2) =
  Seq (Reads (numset_to_list cutset))
    (Branch (get_live_tree ct1) (get_live_tree ct2))
get_live_tree (Seq ct1 ct2) =
  Seq (get_live_tree ct1) (get_live_tree ct2)

```

The function to compute live variables is now very simple:

```

get_live_backward (Writes l) live =
  difference live (list_to_numset l)
get_live_backward (Reads l) live =
  union live (list_to_numset l)
get_live_backward (Branch ct1 ct2) live =
  union (get_live_backward ct1 live) (get_live_backward ct2 live)
get_live_backward (Seq ct1 ct2) live =
  get_live_backward ct1 (get_live_backward ct2 live)

```

The function `check_live_tree` is like `check_clash_tree`, but for `live_tree`.
 We can prove the following correctness theorem:

$$\begin{aligned} &\vdash \text{check_live_tree } col \text{ (get_live_tree } ct) \text{ LN LN} = \text{Some } (livein, flivein) \Rightarrow \\ &\quad \exists livein' flivein'. \text{check_clash_tree } col \text{ } ct \text{ LN LN} = \text{Some } (livein', flivein') \end{aligned}$$

This theorem is proved using the following lemma, which is proved by induction on the clash-tree `ct`.

$$\begin{aligned} &\vdash \text{image } col \text{ (domain } live) = \text{domain } flive \wedge \text{image } col \text{ (domain } live') = \text{domain } flive' \wedge \\ &\quad \text{injective } col \text{ (domain } live') \wedge \\ &\quad \text{domain } live \subseteq \text{domain } live' \wedge \\ &\quad \text{check_live_tree } col \text{ (get_live_tree } ct) \text{ } live' \text{ } flive' = \text{Some } (livein', flivein') \Rightarrow \\ &\quad \exists livein \ flivein. \\ &\quad \text{check_clash_tree } col \text{ } ct \text{ } live \text{ } flive = \text{Some } (livein, flivein) \wedge \\ &\quad \text{domain } livein \subseteq \text{domain } livein' \end{aligned}$$

The hypothesis $\text{image } col \text{ (domain } live) = \text{domain } flive$ and $\text{image } col \text{ (domain } live') = \text{domain } flive'$ are invariants maintained by `check_clash_tree` and `check_live_tree`. The hypothesis $\text{injective } col \text{ (domain } live')$ is also an invariant maintained by `check_clash_tree` and `check_live_tree`.

We also have $livein' = \text{get_live_backward } (get_live_tree \ ct) \ live'$ and $livein = \text{get_live_backward_ct } ct \ live$ (which are properties satisfied by `check_clash_tree` and `check_live_tree`). We also saw before that by construction, $\text{domain } (get_live_backward_ct \ ct \ live) \subseteq \text{domain } (get_live_backward \ (get_live_tree \ ct) \ live)$. This is why the hypothesis $\text{domain } live \subseteq \text{domain } live'$ is here, for the theorem to compose well when dealing with the induction case $ct = \text{Seq } ct_1 \ ct_2$.

4.2 Step 2: get the liveness intervals

4.2.1 The get_intervals function

We will now compute the liveness intervals as described in Section 3.2

The function `get_intervals` takes a `live_tree` and returns two `int num_map`: the beginning and the end of liveness interval for each register.

This function is very simple: when it treats a `Writes` the beginning of intervals are extended, when it treats a `Reads` the end of intervals are extended.

The function traverses the `live_tree` backwards, and maintains an integer `n` which represents the position in the program. We write it like this:

$$\begin{aligned} &\text{get_intervals (Writes } l) \ n \ int_beg \ int_end = \\ &\quad (n - 1, \text{numset_list_add_if_lt } l \ n \ int_beg, \text{numset_list_add_if_gt } l \ n \ int_end) \\ &\text{get_intervals (Reads } l) \ n \ int_beg \ int_end = \\ &\quad (n - 1, \ int_beg, \text{numset_list_add_if_gt } l \ n \ int_end) \\ &\text{get_intervals (Branch } lt_1 \ lt_2) \ n \ int_beg \ int_end = \\ &\quad \text{let } (n_2, \ int_beg_2, \ int_end_2) = \text{get_intervals } lt_2 \ n \ int_beg \ int_end \\ &\quad \text{in} \\ &\quad \text{get_intervals } lt_1 \ n_2 \ int_beg_2 \ int_end_2 \\ &\text{get_intervals (Seq } lt_1 \ lt_2) \ n \ int_beg \ int_end = \\ &\quad \text{let } (n_2, \ int_beg_2, \ int_end_2) = \text{get_intervals } lt_2 \ n \ int_beg \ int_end \\ &\quad \text{in} \\ &\quad \text{get_intervals } lt_1 \ n_2 \ int_beg_2 \ int_end_2 \end{aligned}$$

The `numset_list_add_if_lt` function verifies the following property:

$$\begin{aligned} &\vdash \text{lookup } r \text{ (numset_list_add_if_lt } l \ v \ s) = \\ &\quad \text{if mem } r \ l \ \text{then} \\ &\quad \quad \text{case lookup } r \ s \ \text{of} \\ &\quad \quad \quad \text{None} \Rightarrow \text{Some } v \\ &\quad \quad \quad | \text{Some } vr \Rightarrow \text{if } v \leq vr \ \text{then Some } v \ \text{else Some } vr \\ &\quad \text{else lookup } r \ s \end{aligned}$$

The `numset_list_add_if_gt` function verifies a similar property, with the opposite inequality.

Now, the function which checks a coloration with the liveness intervals is quite simple:

$$\begin{aligned} &\text{check_intervals } f \ \text{int_beg} \ \text{int_end} \iff \\ &\quad \forall r_1 \ r_2. \\ &\quad \quad r_1 \in \text{domain } \text{int_beg} \wedge r_2 \in \text{domain } \text{int_beg} \wedge \\ &\quad \quad \text{interval_intersect (THE (lookup } r_1 \ \text{int_beg}), \text{THE (lookup } r_1 \ \text{int_end}))} \\ &\quad \quad (\text{THE (lookup } r_2 \ \text{int_beg}), \text{THE (lookup } r_2 \ \text{int_end})) \wedge f \ r_1 = f \ r_2 \Rightarrow \\ &\quad \quad r_1 = r_2 \end{aligned}$$

And we have the following correctness theorem:

$$\begin{aligned} &\vdash (n_out, beg_out, end_out) = \text{get_intervals (fix_domination } lt) \ 0 \ \text{LN} \ \text{LN} \wedge \\ &\quad \text{check_intervals } f \ beg_out \ end_out \Rightarrow \\ &\quad \quad \exists \text{liveout } fliveout. \text{check_live_tree } f \ (\text{fix_domination } lt) \ \text{LN} \ \text{LN} = \text{Some (liveout, fliveout)} \end{aligned}$$

(the `fix_domination` function will be explained in Section 4.2.2)

To prove this theorem, we can prove that at every position where a register is live, this position is in the register's liveness interval. We can start by proving that the position is less than the end of the liveness interval:

$$\begin{aligned} &\vdash (n_out, beg_out, end_out) = \text{get_intervals } lt \ n_in \ beg_in \ end_in \wedge \\ &\quad (\forall r. r \in \text{domain } \text{live_in} \Rightarrow \exists v. \text{lookup } r \ end_in = \text{Some } v \wedge n_in \leq v) \Rightarrow \\ &\quad \quad \text{check_number_property} \\ &\quad \quad (\lambda n \ \text{live}. \forall r. r \in \text{domain } \text{live} \Rightarrow \exists v. \text{lookup } r \ end_out = \text{Some } v \wedge n + 1 \leq v) \ lt \ n_in \\ &\quad \quad \text{live_in} \end{aligned}$$

The function `check_number_property` takes a predicate on positions and live variables, and checks if it is true at every position of a `live_tree`.

We prove this theorem by induction on `live_tree`, using a few simple lemmas and the following monotonicity theorem:

$$\begin{aligned} &\vdash (\forall n' \ \text{live}'. n - \text{size_of_live_tree } lt \leq n' \wedge P \ n' \ \text{live}' \Rightarrow Q \ n' \ \text{live}') \wedge \\ &\quad \text{check_number_property } P \ lt \ n \ \text{live} \Rightarrow \\ &\quad \quad \text{check_number_property } Q \ lt \ n \ \text{live} \end{aligned}$$

We would like to prove a similar property about the beginning of intervals, but it turns out that the induction does not work well. The reason why the induction worked well for the end of intervals is that the algorithm processes the `live_tree` in a backward manner, so that the position of end of interval of a register is changed before this register becomes live. It is not true for beginning of intervals: beginning of intervals are changed after the registers are live. It means that the property we want to prove is not true locally.

4.2.2 The `get_intervals_withlive` function

The following function is a modification of `get_intervals`, with better invariants:

$$\begin{aligned}
&\vdash (\forall l \ n \ int_beg \ int_end \ live. \\
&\quad \text{get_intervals_withlive (Writes } l) \ n \ int_beg \ int_end \ live = \\
&\quad (n - 1, \text{numset_list_add_if_lt } l \ n \ int_beg, \text{numset_list_add_if_gt } l \ n \ int_end)) \wedge \\
&\quad (\forall l \ n \ int_beg \ int_end \ live. \\
&\quad \text{get_intervals_withlive (Reads } l) \ n \ int_beg \ int_end \ live = \\
&\quad (n - 1, \text{numset_list_delete } l \ int_beg, \text{numset_list_add_if_gt } l \ n \ int_end)) \wedge \\
&\quad (\forall lt_1 \ lt_2 \ n \ int_beg \ int_end \ live. \\
&\quad \text{get_intervals_withlive (Branch } lt_1 \ lt_2) \ n \ int_beg \ int_end \ live = \\
&\quad \text{let } (n_2, int_beg_2, int_end_2) = \text{get_intervals_withlive } lt_2 \ n \ int_beg \ int_end \ live; \\
&\quad (n_1, int_beg_1, int_end_1) = \\
&\quad \quad \text{get_intervals_withlive } lt_1 \ n_2 \ (\text{difference } int_beg_2 \ live) \ int_end_2 \ live \\
&\quad \text{in} \\
&\quad (n_1, \\
&\quad \quad \text{difference } int_beg_1 \\
&\quad \quad (\text{union (get_live_backward } lt_1 \ live) \ (\text{get_live_backward } lt_2 \ live)), int_end_1)) \wedge \\
&\quad \forall lt_1 \ lt_2 \ n \ int_beg \ int_end \ live. \\
&\quad \text{get_intervals_withlive (Seq } lt_1 \ lt_2) \ n \ int_beg \ int_end \ live = \\
&\quad \text{let } (n_2, int_beg_2, int_end_2) = \text{get_intervals_withlive } lt_2 \ n \ int_beg \ int_end \ live; \\
&\quad (n_1, int_beg_1, int_end_1) = \\
&\quad \quad \text{get_intervals_withlive } lt_1 \ n_2 \ int_beg_2 \ int_end_2 \ (\text{get_live_backward } lt_2 \ live) \\
&\quad \text{in} \\
&\quad (n_1, int_beg_1, int_end_1)
\end{aligned}$$

Here, we removed the set of live variables to `int_beg`, to have the invariant that `domain int_beg` and `domain live` are disjoint. This is the theorem stating that `get_intervals_withlive` preserves this invariant:

$$\begin{aligned}
&\vdash (n_out, beg_out, end_out) = \text{get_intervals_withlive } lt \ n_in \ beg_in \ end_in \ live \wedge \\
&\quad (\forall r. r \in \text{domain } live \Rightarrow r \notin \text{domain } beg_in) \Rightarrow \\
&\quad \forall r. r \in \text{domain } (\text{get_live_backward } lt \ live) \Rightarrow r \notin \text{domain } beg_out
\end{aligned}$$

Intuitively, we don't lose information by doing this because in a well-formed program, when an instruction reads the value of a register, the register has a value, so it was written to before in the program, independently of the path taken in the control-flow graph. We can also formulate this by saying that in a well-formed program, no variable lives at the beginning.

It means that during the execution of `get_intervals`, when a register `r` is live, the beginning of its interval will be overwritten, hence this value is not important and we can delete it.

We can think of a value missing in `int_beg` meaning that the beginning of interval at the end of the execution of the algorithm will be less than the number associated with the current position in the program.

This modification allows us to prove this theorem by induction:

$$\begin{aligned}
&\vdash (n_out, beg_out, end_out) = \text{get_intervals_withlive } lt \ n_in \ beg_in \ end_in \ live_in \wedge \\
&\quad (\forall r \ v. \text{lookup } r \ beg_in = \text{Some } v \Rightarrow n_in \leq v) \wedge (\forall r. r \in \text{domain } live_in \Rightarrow r \notin \text{domain } beg_in) \Rightarrow \\
&\quad \text{check_number_property} \\
&\quad (\lambda n \ live. \forall r. r \in \text{domain } live \Rightarrow (\text{case lookup } r \ beg_out \text{ of None } \Rightarrow n_out \mid \text{Some } x \Rightarrow x) \leq n) \\
&\quad lt \ n_in \ live_in
\end{aligned}$$

Based on the previous intuition, we would like to prove that on a well-formed program, those two functions compute the same thing.

But before that: do we actually have a proof that programs are well-formed, meaning that there are no live variables at the beginning of the program?

I found that the easiest way to ensure that this property is true is to force it, using the following function:

```

    ⊢ fix_domination lt =
      let live = get_live_backward lt LN
      in
        if live = LN then lt
        else Seq (Writes (map fst (toAList live))) lt
  
```

4.2.3 Proving that the two functions compute the same thing

Let's name *int_end* the end of intervals (which are computed the same way in both functions), *int_beg* the beginning of intervals produced by *get_intervals*, and *int_beg_withlive* the one produced by *get_intervals_withlive*.

It is easy to prove by induction that when *int_beg* and *int_beg_withlive* share a common key, then the values associated with this key are equal:

```

    ⊢ (n1, beg1, end1) = get_intervals lt n beg end ∧
      (n2, beg2, end2) = get_intervals_withlive lt n beg' end live ∧
      (∀ r v. lookup r beg = Some v ⇒ n ≤ v) ∧ (∀ r v. lookup r beg' = Some v ⇒ n ≤ v) ∧
      (∀ r v1 v2. lookup r beg = Some v1 ∧ lookup r beg' = Some v2 ⇒ v1 = v2) ⇒
      ∀ r v1 v2. lookup r beg1 = Some v1 ∧ lookup r beg2 = Some v2 ⇒ v1 = v2
  
```

Now we only have to prove that $\text{domain } int_beg = \text{domain } int_beg_withlive$

We define the set of registers in a *live_tree* like this:

```

    live_tree_registers (Writes l) = set l
    live_tree_registers (Reads l) = set l
    live_tree_registers (Branch lt1 lt2) = live_tree_registers lt1 ∪ live_tree_registers lt2
    live_tree_registers (Seq lt1 lt2) = live_tree_registers lt1 ∪ live_tree_registers lt2
  
```

We will prove the inclusions:

$\text{live_tree_registers } lt \subset \text{domain } int_end \subset \text{domain } int_beg_withlive \subset \text{domain } int_beg \subset \text{live_tree_registers } lt$

Proving $\text{domain } int_beg \subset \text{live_tree_registers } lt$ is easy by induction with the following theorem:

```

    ⊢ (nout, begout, endout) = get_intervals lt nin begin endin ⇒
      domain begout ⊆ domain begin ∪ live_tree_registers lt
  
```

Proving $\text{domain } int_beg_withlive \subset \text{domain } int_beg$ is easy by induction with the following theorem:

```

    ⊢ (n1, begout1, end1) = get_intervals_withlive lt n begin1 end live ∧
      (n2, begout2, end2) = get_intervals lt n begin2 end ∧
      domain begin1 ⊆ domain begin2 ⇒
      domain begout1 ⊆ domain begout2
  
```

Proving $\text{live_tree_registers } lt \subset \text{domain } int_end$ is easy with the following theorem:

```

    ⊢ (nout, begout, endout) = get_intervals_withlive lt nin begin endin livein ⇒
      domain endin ∪ live_tree_registers lt ⊆ domain endout
  
```

The only proof left is $\text{live_tree_registers } lt \subset \text{domain } int_beg_withlive$. It turns out it was quite difficult.

Before proving this inclusion, we prove the following lemma which says that `get_intervals_withlive` is Lipschitz-continuous¹

$$\begin{aligned} &\vdash (nout_1, begout_1, endout_1) = \text{get_intervals_withlive } lt \ n_1 \ beg_1 \ end_1 \ live \ \wedge \\ &\quad (nout_2, begout_2, endout_2) = \text{get_intervals_withlive } lt \ n_2 \ beg_2 \ end_2 \ live \ \wedge \\ &\quad \text{domain } beg_2 \subseteq \text{domain } beg_1 \cup \text{domain } s \Rightarrow \\ &\quad \text{domain } begout_2 \subseteq \text{domain } begout_1 \cup \text{domain } s \end{aligned}$$

It allows us to prove that $\text{domain } int_end \subset \text{domain } int_beg_withlive$:

$$\begin{aligned} &\vdash \text{domain } end_in \subseteq \text{domain } beg_in \cup \text{domain } live_in \ \wedge \\ &\quad (n_out, beg_out, end_out) = \text{get_intervals_withlive } lt \ n_in \ beg_in \ end_in \ live_in \Rightarrow \\ &\quad \text{domain } end_out \subseteq \text{domain } beg_out \cup \text{domain } (\text{get_live_backward } lt \ live_in) \end{aligned}$$

Proof sketch it is proved by induction on lt . The only hard case is $lt = \text{Branch } lt_1 \ lt_2$. Let's write $(n_2, beg_2, end_2) = \text{get_intervals_withlive } lt_2 \ n_in \ beg_in \ end_in \ live_in$ and

$$(n_1, beg_1, end_1) = \text{get_intervals_withlive } lt_1 \ n_2 \ (\text{difference } beg_2 \ live_in) \ end_2 \ live_in.$$

By the induction hypothesis, we have

$$\text{domain } end_2 \subseteq \text{domain } beg_2 \cup \text{domain } (\text{get_live_backward } lt_2 \ live_in).$$

Unfortunately, we can't use the induction hypothesis for lt_1 since we don't have the hypothesis

$$\text{domain } end_2 \subseteq \text{domain } (\text{difference } beg_2 \ live_in) \cup \text{domain } live_in.$$

We can force the hypothesis: let's write

$$\begin{aligned} &(n_1, beg'_1, end_1) = \\ &\text{get_intervals_withlive } lt_1 \ n_2 \ (\text{difference } (\text{union } beg_2 \ (\text{get_live_backward } lt_2 \ live_in)) \ live_in) \ end_2 \ live_in. \end{aligned}$$

We can see that

$$\begin{aligned} &\text{domain } (\text{difference } (\text{union } beg_2 \ (\text{get_live_backward } lt_2 \ live_in)) \ live_in) \cup \text{domain } live_in = \\ &\text{domain } beg_2 \cup \text{domain } (\text{get_live_backward } lt_2 \ live_in) \cup \text{domain } live_in \end{aligned}$$

which contains end_2 as we saw at the beginning of the proof.

From this, we can deduce using the induction hypothesis that

$$\text{domain } end_1 \subseteq \text{domain } beg'_1 \cup \text{domain } (\text{get_live_backward } lt_1 \ live_in)$$

and by the Lipschitz-continuity theorem, we have

$$\text{domain } beg'_1 \subseteq \text{domain } beg_1 \cup \text{domain } (\text{get_live_backward } lt_2 \ live_in).$$

Therefore,

$$\begin{aligned} &\text{domain } end_1 \subseteq \\ &\text{domain } beg_1 \cup \text{domain } (\text{get_live_backward } lt_1 \ live_in) \cup \text{domain } (\text{get_live_backward } lt_2 \ live_in) \end{aligned}$$

which proves the theorem. \square

Using all the previous lemmas, we can prove the following theorem:

$$\begin{aligned} &\vdash (n, beg, end) = \text{get_intervals_withlive } (\text{fix_domination } lt) \ 0 \ \text{LN} \ \text{LN} \ \text{LN} \ \wedge \\ &\quad (n', beg', end') = \text{get_intervals } (\text{fix_domination } lt) \ 0 \ \text{LN} \ \text{LN} \ \Rightarrow \\ &\quad \forall r. \text{lookup } r \ beg = \text{lookup } r \ beg' \end{aligned}$$

4.3 Step 3: do the allocation

4.3.1 Modifications to the original algorithm

In section 3.4, we saw that CakeML's register allocation needs a few features not covered by the original linear scan register allocation algorithm [8].

¹When f is an increasing function, the Lipschitz-continuity with a constant equal to 1 translates to $\forall x \forall y > 0, f(x+y) \leq f(x) + y$. Here the theorem says something like $f(T \cup S) \subset f(T) \cup S$

These features are:

- physical registers should be allocated themselves (e.g. for calling conventions)
- stack registers should be on the stack
- the *moves* list is a list of pair of registers that should be allocated to the same register if possible (to remove useless "move" instructions)
- the *forced* list is a list of pair of registers that can't be allocated to the same register
- registers that are spilled on the stack should be allocated in a "smart" way to minimise the stack-frame size

Small stack-frame size We used the classical solution to optimize the stack-frame size, which is to do the allocation in two passes. The first pass is the allocation as we saw previously, with a new fresh register on the stack each time we spill a register. During the second pass, we run the same algorithm only on the registers spilled on the stack to reallocate them in a smarter way.

The *moves* list The requirement that some pair of registers should be allocated to the same register if possible is easy to satisfy: before choosing a color from the color pool, check if a color of one of the registers in the *moves* list is available.

Stack registers The requirement that stack registers should be allocated on the stack is also easy to satisfy: when we encounter a such register, we can simply spill it.

The *forced* list The requirement that some pair of registers should not be allocated to the same register can be satisfied by adding some sanity check when we choose a color from the color-pool.

Physical registers The requirement that physical registers should be allocated to themselves is harder to satisfy. When we allocate the physical register *reg* which should have the color *col*, and that *col* is used by a register *reg₂* in the active list, what should we do? It's not possible to give a smart color to *reg₂* to free the color *col*: the only option is to spill *reg₂*. But that's a problem since it would produce very bad allocation.

I found a good solution to this, namely, to forget about this requirement, and fix the coloration afterward to fit this requirement. How do we fix the coloration? Afterward, we find a permutation of colors to ensure that each physical register is allocated to itself. A necessary and sufficient condition for a such permutation to exist, is that the physical registers must have different colors.

In summary: to ensure that physical registers are allocated to themselves, we first do the allocation and ensure that physical registers all have distinct colors, and we do a color permutation afterwards.

4.3.2 The state and invariants used in the linear scan algorithm

The internal state of the algorithm is represented by the two following records:

```

linear_scan_state = ⟨
  active : (int × num) list;
  colorpool : num list;
  physcols : num_set;
  colornum : num;
  colormax : num;
  stacknum : num
⟩
linear_scan_hidden_state = ⟨
  colors : num list
⟩

```

`active` is the list of active registers: the second component of the pair is the active register, the first component is the register's end of liveness interval. This list is sorted by increasing end of liveness interval. `colorpool` is the pool of available colors, along with the colors c such that $\text{colornum} \leq c < \text{colormax}$. `physcols` is the set of colors used by the physical registers (to ensure the colors are all distinct). `stacknum` is the color of a fresh register to spill on the stack.

`colors` is an array such that its regth element is the color of reg . In the state it is represented as a functional list, but when the HOL function is translated into a CakeML AST during the bootstrap process, it is represented as a static array with $O(1)$ element lookup. This is why `colors` is in a separate record. This record is said "hidden" because it is used in a state monad.

This state has a lot of invariants that are preserved during the algorithm. The invariants depend on several parameters: `int_beg` and `int_end` the beginning and end of liveness intervals, `st` and `sth` the state and the hidden state, the `forced` list, a list l of register that were processed by the algorithm, `pos` representing the position of the beginning of liveness interval of the last processed register.

The most important invariants are the one directly linked to the final correctness theorem:

- two different registers in l with intersecting liveness intervals must not have the same color
- two different registers in l and in the `forced` list must not have the same color

Some invariants describe the content of the `active` list:

- the `active` list is sorted by its first component
- all registers in the `active` list are also present in the l list
- the color of registers present in the `active` list are strictly lower than `colornum`
- for each element of `active`, the first component is equal to the end of liveness interval of the second component
- if a register of l is not spilled to the stack and its end of liveness interval is greater than `pos` then it is in the `active` list
- if a register r is in `active` then `pos` is less than its end of liveness interval incremented by one²

Some invariant describe the relationship between the colors manipulated by the algorithm:

- the colors of `colorpool` concatenated with the colors of the registers of `active` are all distinct³
- the color of physical registers are all distinct

²the increment is necessary to prove the correctness of `ExprireOldIntervals` by induction, in the case we are removing several intervals with the same endpoint

³this invariant also says that the colors in `colorpool` and the colors of registers in `active` are disjoint

And then there are miscellaneous invariants:

- the color of every register in l is strictly less than `stacknum`
- every color in `colorpool` is strictly less than `colornum`
- `colornum` is less than `colormax` which is less than `stacknum`
- if the color of a register of l is less than `colormax` then it is less than `colornum`
- the beginning of liveness interval of every register in l is less than pos
- `phycols` is the set of colors of registers reg in l such that reg is a physical register and reg is not spilled to the stack
- every register of l is strictly less than the length of `colors`

The invariants have an additional parameter `mincol` which is the minimal color used in the algorithm. It is useful to prove that the second pass that reallocates spilled registers keeps registers on the stack. It follows the following invariants:

- `mincol` is less than `colornum`
- `mincol` is less than every color in `colorpool`
- `mincol` is less than the color of every register in `active`

The invariants are checked by the predicate `good_linear_scan_state int_beg int_end st sth l pos forced mincol`

4.3.3 Implementation of the linear scan algorithm

In this section, we present the different functions used to implement the linear scan algorithm, along with the correctness theorems. The correctness theorems usually says that the function doesn't fail (i.e. there is no array out-of-bounds), that the invariants are preserved, that `length colors` and `colormax` are not changed, and express which colors are changed.

The `remove_inactive_intervals beg st sth` function corresponds to the `ExpireOldIntervals` function presented in Section 3.2: it removes active intervals whose endpoint is before beg .

It has the following correctness theorem (notice that pos is replaced by beg in the conclusion):

$$\begin{aligned} &\vdash \text{good_linear_scan_state } int_beg \ int_end \ st \ sth \ l \ pos \ forced \ mincol \ \wedge \\ &\quad pos \leq beg \Rightarrow \\ &\quad \exists stout. \\ &\quad (\text{Success } stout, sth) = \text{remove_inactive_intervals } beg \ st \ sth \ \wedge \\ &\quad \text{good_linear_scan_state } int_beg \ int_end \ stout \ sth \ l \ beg \ forced \ mincol \ \wedge \\ &\quad stout.colormax = st.colormax \end{aligned}$$

`find_color st forbidden` tries to find a color in the `colorpool` which is not in the set `forbidden`. When it succeeds, it returns $(stout, \text{Some } col)$ where $stout$ is the output state. The output color is removed from the `colorpool` in the output state.

$$\begin{aligned} &\vdash \text{good_linear_scan_state } int_beg \ int_end \ st \ sth \ l \ pos \ forced \ mincol \ \wedge \\ &\quad \text{domain } forbidden \subseteq \{ \text{el } r \ sth.colors \mid \text{mem } r \ l \} \ \wedge \\ &\quad \text{find_color } st \ forbidden = (stout, \text{Some } col) \Rightarrow \\ &\quad \text{good_linear_scan_state } int_beg \ int_end \\ &\quad (\text{stout with } colorpool := col :: stout.colorpool) \ sth \ l \ pos \ forced \ mincol \ \wedge \\ &\quad col < stout.colornum \ \wedge \ col \notin \text{domain } forbidden \ \wedge \\ &\quad st = stout \text{ with } \langle colorpool := st.colorpool; colornum := st.colornum \rangle \end{aligned}$$

`color_register st reg col rend sth` assign the color `col` to the register `reg`. It updates the active list with the end of interval `rend`. It has the following correctness theorem:

$$\begin{aligned} &\vdash \text{good_linear_scan_state } int_beg \ int_end \ st \ sth \ l \ pos \ forced \ mincol \ \wedge \\ &\text{forbidden_is_from_map_color_forced } forced \ l \ sth.colors \ reg \ forbidden \ \wedge \ col \notin \text{domain } forbidden \ \wedge \\ &(\text{is_phy_var } reg \Rightarrow \text{domain } st.phycols \subseteq \text{domain } forbidden) \ \wedge \\ &\neg \text{mem } col \ (st.colorpool \# \text{map } (\lambda (e,r). \text{el } r \ sth.colors) \ st.active) \ \wedge \\ &\text{the } 0 \ (\text{lookup } reg \ int_beg) = pos \ \wedge \ \text{the } 0 \ (\text{lookup } reg \ int_beg) \leq \text{the } 0 \ (\text{lookup } reg \ int_end) \ \wedge \\ &col < st.colornum \ \wedge \ mincol \leq col \ \wedge \ reg < \text{length } sth.colors \ \wedge \ \neg \text{mem } reg \ l \Rightarrow \\ &\exists stout \ sthout. \\ &(\text{Success } stout, sthout) = \text{color_register } st \ reg \ col \ (\text{the } 0 \ (\text{lookup } reg \ int_end)) \ sth \ \wedge \\ &\text{good_linear_scan_state } int_beg \ int_end \ stout \ sthout \ (reg::l) \ pos \ forced \ mincol \ \wedge \\ &\text{length } sthout.colors = \text{length } sth.colors \ \wedge \\ &(\forall r. r \neq reg \Rightarrow \text{el } r \ sth.colors = \text{el } r \ sthout.colors) \ \wedge \ stout.colormax = st.colormax \end{aligned}$$

The hypothesis of this theorem mostly says that `reg` and `col` are compatible with the invariants. An interesting part is the `forbidden_is_from_map_color_forced` which says that the `forbidden` set contains the color of registers that conflicts with `reg` in `forced`.

`spill_register st reg sth` spills the register `reg` on a fresh stack location. It has the following correctness theorem:

$$\begin{aligned} &\vdash (\neg \text{is_phy_var } reg \vee \neg \text{mem } reg \ l) \ \wedge \\ &\text{good_linear_scan_state } int_beg \ int_end \ st \ sth \ l \ pos \ forced \ mincol \ \wedge \\ ® < \text{length } sth.colors \ \wedge \ \text{the } 0 \ (\text{lookup } reg \ int_beg) \leq pos \Rightarrow \\ &\exists stout \ sthout. \\ &(\text{Success } stout, sthout) = \\ &\text{spill_register } (st \ \text{with } active := \text{filter } (\lambda (e,r). r \neq reg) \ st.active) \ reg \ sth \ \wedge \\ &\text{good_linear_scan_state } int_beg \ int_end \ stout \ sthout \ (reg::l) \ pos \ forced \ mincol \ \wedge \\ &\text{length } sthout.colors = \text{length } sth.colors \ \wedge \\ &(\forall r. r \neq reg \Rightarrow \text{el } r \ sth.colors = \text{el } r \ sthout.colors) \ \wedge \\ &stout.colormax = st.colormax \ \wedge \ st.colormax \leq \text{el } reg \ sthout.colors \end{aligned}$$

This theorem is a bit more complex than the other functions, because every function except this one can have the hypothesis $\neg \text{mem } reg \ l$. However as we see in the `SpillInterval` function in Section 3.2 we can spill a register that is in the active list (and therefore in `l`).

If $\neg \text{mem } reg \ l$ then `st.active = filter ($\lambda (e,r). r \neq reg$) st.active` and then this theorem is similar to the usual ones. However, if `mem reg l` we need the requirement `¬is_phy_var reg` because this function does not update `phycols`⁴. When we have `mem reg l` then `reg` needs to be removed to the active list: using the filter function is a brutal way to do this⁵.

The `find_spill st forbidden reg rend force sth` function corresponds to the `SpillInterval` function in Section 3.2. `force` is a boolean which says that the function should not spill `reg` to the stack if possible (it is best-effort). If the function "steals" the color of an active register, the function makes sure that it is not in the `forbidden` set. It has the following correctness theorem, which is similar to `color_register`'s correctness theorem.

⁴it could update `phycols`, but it is not done because in practice the hypothesis `mem reg l` \Rightarrow `¬is_phy_var reg` is satisfied

⁵hopefully, the filter function only appears in the proof

$$\begin{aligned}
&\vdash \neg \text{mem } reg \ l \wedge \\
&\text{good_linear_scan_state } int_beg \ int_end \ st \ sth \ l \ (\text{the } 0 \ (\text{lookup } reg \ int_beg)) \ \text{forced } mincol \wedge \\
® < \text{length } sth.colors \wedge \text{forbidden_is_from_map_color_forced } forced \ l \ sth.colors \ reg \ \text{forbidden} \wedge \\
&(\text{is_phy_var } reg \Rightarrow \text{domain } st.phycols \subseteq \text{domain } forbidden) \wedge \\
&\text{the } 0 \ (\text{lookup } reg \ int_beg) \leq \text{the } 0 \ (\text{lookup } reg \ int_end) \Rightarrow \\
&\exists stout \ sthout. \\
&\quad (\text{Success } stout, sthout) = \text{find_spill } st \ forbidden \ reg \ (\text{the } 0 \ (\text{lookup } reg \ int_end)) \ \text{force } sth \wedge \\
&\quad \text{good_linear_scan_state } int_beg \ int_end \ stout \ sthout \ (reg::l) \ (\text{the } 0 \ (\text{lookup } reg \ int_beg)) \\
&\quad \text{forced } mincol \wedge \text{length } sthout.colors = \text{length } sth.colors \wedge \\
&\quad (\forall r. \neg \text{mem } r \ (reg::l) \Rightarrow \text{el } r \ sthout.colors = \text{el } r \ sth.colors) \wedge \\
&\quad (\forall r. \text{mem } r \ l \wedge \text{is_phy_var } r \Rightarrow \text{el } r \ sthout.colors = \text{el } r \ sth.colors) \wedge \\
&\quad stout.colormax = st.colormax
\end{aligned}$$

`apply_reg_exchange l sth` takes a list l of physical registers, and do the exchange described in Section 4.3.1. The following correctness theorem says that the registers in l are assigned to the right registers, and that the exchange preserves the fact that two registers have the same color.

$$\begin{aligned}
&\vdash \text{all_distinct } (\text{map } (\lambda r. \text{el } r \ sth.colors) \ l) \wedge (\forall r. \text{mem } r \ l \Rightarrow \text{is_phy_var } r) \wedge \\
&(\forall r. \text{mem } r \ l \Rightarrow r < \text{length } sth.colors) \Rightarrow \\
&\exists sthout. \\
&\quad (\text{Success } (), sthout) = \text{apply_reg_exchange } l \ sth \wedge \text{length } sthout.colors = \text{length } sth.colors \wedge \\
&\quad (\forall r_1 \ r_2. \\
&\quad \quad r_1 < \text{length } sth.colors \wedge r_2 < \text{length } sth.colors \Rightarrow \\
&\quad \quad \text{el } r_1 \ sthout.colors = \text{el } r_2 \ sthout.colors \Rightarrow \text{el } r_1 \ sth.colors = \text{el } r_2 \ sth.colors) \wedge \\
&\quad \forall r. \text{mem } r \ l \Rightarrow \text{el } r \ sthout.colors = r \ \text{div } 2
\end{aligned}$$

In reality, the theorem has a more complex conclusion which gives information about how it preserves the fact that some registers are on the stack, but it was omitted here for simplicity.

`linear_reg_alloc_intervals int_beg int_end k forced moves reglist_unsorted sth` is the full algorithm. It does a first pass which corresponds to the original linear scan algorithm [8], then do the color exchange for physical registers that are not on the stack. Then, it runs a second pass for registers that were spilled to the stack and stack registers, then do the exchange for physical registers that are on the stack. This function has the following correctness theorem, which looks a lot like the correctness theorem in the original register allocator as seen in Section 3.4:

$$\begin{aligned}
&\vdash \text{every } (\lambda (r_1, r_2). \text{mem } r_1 \ reglist \wedge \text{mem } r_2 \ reglist) \ \text{forced} \wedge \\
&\text{every } (\lambda (r_1, r_2). r_1 < \text{length } sth.colors \wedge r_2 < \text{length } sth.colors) \ (\text{map } \text{snd } moves) \wedge \\
&\text{every } (\lambda r. r < \text{length } sth.colors) \ reglist \wedge \\
&\text{every } (\lambda r. \text{the } 0 \ (\text{lookup } r \ int_beg) \leq \text{the } 0 \ (\text{lookup } r \ int_end)) \ reglist \wedge \\
&\text{all_distinct } reglist \wedge \text{set } reglist = \text{domain } int_beg \wedge \text{domain } int_beg = \text{domain } int_end \Rightarrow \\
&\exists sthout. \\
&\quad (\text{Success } (), sthout) = \\
&\quad \text{linear_reg_alloc_intervals } int_beg \ int_end \ k \ \text{forced } moves \ reglist \ sth \wedge \\
&\quad \text{check_intervals } (\lambda r. \text{el } r \ sthout.colors) \ int_beg \ int_end \wedge \\
&\quad \text{every} \\
&\quad \quad (\lambda r. \\
&\quad \quad \quad \text{if } \text{is_phy_var } r \ \text{then } \text{el } r \ sthout.colors = r \ \text{div } 2 \\
&\quad \quad \quad \text{else if } \text{is_stack_var } r \ \text{then } k \leq \text{el } r \ sthout.colors \\
&\quad \quad \quad \text{else } \top) \ reglist \wedge \\
&\quad \text{every } (\lambda (r_1, r_2). \text{el } r_1 \ sthout.colors = \text{el } r_2 \ sthout.colors \Rightarrow r_1 = r_2) \ \text{forced} \wedge \\
&\quad \text{length } sthout.colors = \text{length } sth.colors
\end{aligned}$$

`linear_scan_reg_alloc k moves ct forced` is a function that combines the previous steps. It has exactly the same correctness theorem as the other register allocator seen in Section 3.4

5 Evaluation

The linear scan algorithm is compared against the iterated register coalescing (IRC) algorithm [1], and a simpler allocator which corresponds to the IRC algorithm without coalescing.

It is compared in term of compilation time, and in term of quality of the produced code on the standard benchmarks used for CakeML.

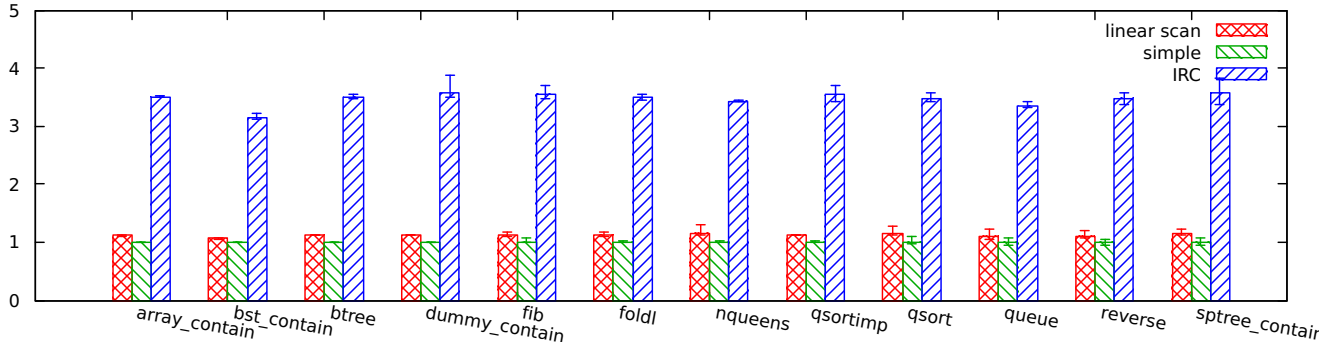


Figure 1: Compilation speed. The time is normalized on the time taken by the simple algorithm

We can see here that the linear scan algorithm is a bit slower than the simple allocator, but about three times faster than the IRC algorithm.

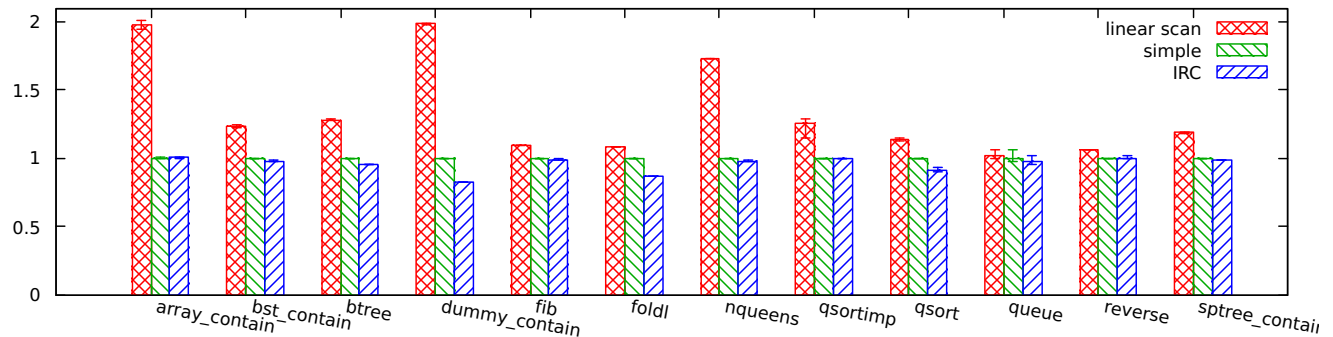


Figure 2: Produced code speed. The time is normalized on the time took by the simple algorithm

We can see here that the linear scan algorithm produces code of bad quality, even compared to the simple algorithm.

Why do we have such results, when the original paper announces that the code produced is about 10% slower than the code generated by the IRC algorithm [8, sec 5.3.2]? We found that the culprit were physical registers, more precisely calling conventions. The calling convention uses the first physical registers as the arguments of the function called. If in a program a function is called at the beginning, and a other function is called at the end, the liveness interval of the first physical register might be for example $\{1, 2, 3, 1001, 1002, 1003\} \subset [1, 1003]$. This means that the physical register might live for a very short time, but its liveness interval will be huge, and this physical register can't be used by another register in the middle of the program.

The solution to this problem is to do the allocation after the SSA pass, but before the calling convention is enforced. For this, we have to split a function in two parts since the SSA-form and the calling convention are enforced by the same function.

6 Future work

6.1 Split the SSA-form and the calling conventions

As we saw in Section 5, in order to produce code of decent quality we must split in two parts the function that transforms the program in SSA form and enforces calling conventions

6.2 Optimise the constant of the algorithm

There are a few things that can be done faster in the whole algorithm:

- we can compute liveness intervals directly on a `clash_tree` (the `live_tree` was useful for the proofs)
- we can use static arrays to store the endpoints of the liveness intervals (instead of a `int num_map`)
- the sorting function used is a quicksort on functional lists: we could write a faster version using static arrays

6.3 Remove the `fix_domination` function

Since proving the property that the set of live variables is empty at the beginning of the program is not easy, we force it using the function `fix_domination`. Proving the domination property would allow to remove this function call, which is in practice useless.

6.4 Use the fact that we work on a SSA-like AST

The WORDLANG AST in CakeML almost has an SSA form⁶, and the live variables have nice properties in SSA-form.

It is possible to prove that in a program in SSA form, when two registers interfere, they interfere at the definition of one of the two registers [2, lemmas 11 & 12].

The reason we filled the liveness holes to have one big interval was to have the following property: two intervals intersects iff one of them contains the starting point of the other. It turns out that for programs in SSA form, this property is true even if the lifetime has holes.

Using this property, it is possible to produce better allocation [12]

7 Conclusion

Before this internship, I had a bit of experience with Coq, but I never really went beyond some toy examples.

During this internship, I learned to use HOL4 and I saw how verification works on a large program. I implemented and verified end-to-end a new register allocation algorithm, which will enable CakeML to have shorter compilation time.

I realised that often, the hard part was not to do the actual proofs, but to find the theorems I needed to prove: for most theorems, once I found the correct invariants, the proofs were mostly straightforward.

It was a very interesting and rewarding experience. It changed the way I view maths, because I had to be absolutely rigorous since the computer doesn't accept any hand-waving arguments.

⁶"almost" meaning that ϕ -functions are already resolved: some variable might be assigned twice : at the end of both branches of a condition

References

- [1] Lal George and Andrew W. Appel. Iterated register coalescing. In *POPL*, pages 208–218. ACM Press, 1996.
- [2] Sebastian Hack and Gerhard Goos. Optimal register allocation for SSA-form programs in polynomial time. *Inf. Process. Lett.*, 98(4):150–155, 2006.
- [3] Son Ho, Oskar Abrahamsson, Ramana Kumar, Magnus O. Myreen, Yong Kiam Tan, and Michael Norrish. Proof-producing synthesis of CakeML with I/O and local state from monadic HOL functions. In *IJCAR*, volume 10900 of *Lecture Notes in Computer Science*, pages 646–662. Springer, 2018.
- [4] Gary A. Kildall. A unified approach to global program optimization. In *POPL*. ACM Press, 1973.
- [5] Ramana Kumar, Magnus O. Myreen, Michael Norrish, and Scott Owens. CakeML: a verified implementation of ML. In *POPL*, pages 179–192. ACM, 2014.
- [6] Magnus O. Myreen and Scott Owens. Proof-producing translation of higher-order logic into pure and stateful ML. *J. Funct. Program.*, 24(2-3):284–315, 2014.
- [7] Scott Owens, Magnus O. Myreen, Ramana Kumar, and Yong Kiam Tan. Functional big-step semantics. In *ESOP*, volume 9632 of *LNCS*. Springer, 2016.
- [8] Massimiliano Poletto and Vivek Sarkar. Linear scan register allocation. *ACM Trans. Program. Lang. Syst.*, 21(5):895–913, 1999.
- [9] Konrad Slind and Michael Norrish. A brief overview of HOL4. In *TPHOLs*, volume 5170 of *Lecture Notes in Computer Science*, pages 28–32. Springer, 2008.
- [10] Yong Kiam Tan, Magnus O. Myreen, Ramana Kumar, Anthony C. J. Fox, Scott Owens, and Michael Norrish. A new verified compiler backend for CakeML. In *ICFP*, pages 60–73. ACM, 2016.
- [11] Yong Kiam Tan, Scott Owens, and Ramana Kumar. A verified type system for CakeML. In *IFL*, pages 7:1–7:12. ACM, 2015.
- [12] Christian Wimmer and Hanspeter Mössenböck. Optimized interval splitting in a linear scan register allocator. In *VEE*, pages 132–141. ACM, 2005.