

How to compile pattern matching

Jules Jacobs

May 2, 2021

Abstract

This is a note about compiling ML-style pattern matching to decision trees. I try to keep it simple while still presenting the key idea of efficient pattern match compilation [Mar08].

1 Introduction

Our goal is to compile ML-style pattern matching:

```
match a with
| Add(Zero, Zero)     $\Rightarrow e_1$ 
| Mul(Zero, x)        $\Rightarrow e_2$ 
| Add(Succ(x), y)     $\Rightarrow e_3$ 
| Mul(x, Zero)        $\Rightarrow e_4$ 
| Mul(Add(x, y), z)   $\Rightarrow e_5$ 
| Add(x, Zero)        $\Rightarrow e_6$ 
| x                   $\Rightarrow e_7$ 
```

In order to keep things straight, I'll use green variables *a* for bound local variables that have a value and red variables *x* for variables that will be given a value by pattern matching.

The naive way to compile this is to test the patterns from top to bottom, and try to match *a* against each pattern. This is inefficient: if the outer constructor of *a* is an Add but the first pattern fails, then the second pattern will test if *a* is a Mul, even though it is already known to be an Add at this point. Then the third pattern will test once more whether *a* is an Add, which is redundant. We'd like to compile pattern matching to code that does **no unnecessary tests**.

This does mean that we *must* start by trying to match the first pattern: if the first pattern succeeds then any work trying to match other patterns is unnecessary. So our strategy will be to emulate the naive strategy that tries to match the first pattern completely, but when we fail in the middle of the matching process, we remember what we learnt to compile a specialised version of the remaining pattern matches that takes what we learnt into account.

2 Compiling the example

It will be useful to represent the pattern matching problem as a list of conditions to be tested:

```
| a is Add(Zero, Zero)     $\Rightarrow e_1$ 
| a is Mul(Zero, x)        $\Rightarrow e_2$ 
| a is Add(Succ(x), y)     $\Rightarrow e_3$ 
| a is Mul(x, Zero)        $\Rightarrow e_4$ 
| a is Mul(Add(x, y), z)   $\Rightarrow e_5$ 
| a is Add(x, Zero)        $\Rightarrow e_6$ 
| a is x                   $\Rightarrow e_7$ 
```

In general, each clause can have multiple tests and will be of the form

$$| a_1 \text{ is } pattern_1, \dots, a_k \text{ is } pattern_k \implies e$$

Where a_1, \dots, a_k are bound variables, and $pattern_1, \dots, pattern_k$ are patterns. This additional flexibility will turn out to be useful during pattern match compilation.

Our goal will be an algorithm that takes as input such a list of clauses, and outputs a tree of simple pattern matches that test against a single constructor:

$$\begin{array}{ll} \text{match } a \text{ with} & \\ | C(a_1, \dots, a_n) & \implies A \\ | _ & \implies B \end{array}$$

Let's see how to do this for the example. We start working on the first equation by testing a against the Add constructor with the following pattern match:

$$\begin{array}{ll} \text{match } a \text{ with} & \\ | \text{Add}(a_1, a_2) & \implies A \\ | _ & \implies B \end{array}$$

Then we have to solve the following sub problem for A:

$$\begin{array}{ll} | a_1 \text{ is Zero}, a_2 \text{ is Zero} & \implies e_1 \\ | a_1 \text{ is Succ}(x), a_2 \text{ is } y & \implies e_3 \\ | a_1 \text{ is } x, a_2 \text{ is Zero} & \implies e_6 \\ | a \text{ is } x & \implies e_7 \end{array}$$

Notice how we now have multiple equations per clause, corresponding to the a_1 and a_2 in the generated **match**, which become bound variables a_1 and a_2 in A. This is the sub problem for B:

$$\begin{array}{ll} | a \text{ is Mul}(\text{Zero}, x) & \implies e_2 \\ | a \text{ is Mul}(x, \text{Zero}) & \implies e_4 \\ | a \text{ is Mul}(\text{Add}(x, y), z) & \implies e_5 \\ | a \text{ is } x & \implies e_7 \end{array}$$

We can simplify these problems by pushing bare equations such as $a_2 \text{ is } y$ into the right hand sides. Then the sub problem for A becomes:

$$\begin{array}{ll} | a_1 \text{ is Zero}, a_2 \text{ is Zero} & \implies e_1 \\ | a_1 \text{ is Succ}(x) & \implies \text{let } y = a_2 \text{ in } e_3 \\ | a_2 \text{ is Zero} & \implies \text{let } x = a_1 \text{ in } e_6 \\ | & \implies \text{let } x = a \text{ in } e_7 \end{array}$$

(actually, we could have done this for the last equation already in the previous step)

We now continue matching by trying to match $a_1 = \text{Zero}$, by generating the pattern match:

$$\begin{array}{ll} \text{match } a_1 \text{ with} & \\ | \text{Zero} & \implies C \\ | _ & \implies D \end{array}$$

And continuing recursively for C and D.

3 The general algorithm

Let's generalise and see what's going on. Given a list of clauses to generate a pattern matching tree for, we use this algorithm:

1. Push bare equations a is y into the right hand sides, so that all the remaining equations test constructors.
2. Select one of the equations $a = C(\mathcal{P}_1, \dots, \mathcal{P}_n)$ in the first clause using some heuristic.
3. Generate this pattern match:

$$\begin{array}{ll} \text{match } a \text{ with} & \\ | C(a_1, \dots, a_n) & \Rightarrow A \\ | _ & \Rightarrow B \end{array}$$

4. Create the two sub problems A and B as follows by iterating over all the clauses. One of three cases can happen:
 - (a) The clause contains an equation a is $C(\mathcal{P}_1, \dots, \mathcal{P}_n), \dots \text{REST} \dots$ for a .
Add the expanded clause a_1 is $\mathcal{P}_1, \dots, a_n$ is $\mathcal{P}_n, \dots \text{REST} \dots$ to A .
Make sure that the fresh variable names a_1, \dots, a_n are used consistently in **match** a with $C(a_1, \dots, a_n)$ and in the equations a_1 is $\mathcal{P}_1, \dots, a_n$ is \mathcal{P}_n .
 - (b) The clause contains an equation a is $D(\mathcal{P}_1, \dots, \mathcal{P}_n), \dots \text{REST} \dots$ where $D \neq C$.
Add this clause to B unchanged.
 - (c) The clause contains no equation for a . Add this clause to both A and B .
(note that each clause can only have one equation for a)
5. Recursively generate code for A and B .

The recursion has two base cases:

- The list of clauses is empty: all patterns failed, so generate an error:
“Error: Non-exhaustive pattern match.”
- The first clause is empty (has zero equations): the first clause succeeded to match, so simply return the right hand side e_i .

4 Heuristic

How do we decide which equation to pick from the first clause to branch on? Any equation works, but we’d like to generate a compact pattern matching tree. Whenever we are in case (c) for one of the other clauses, we have to add that clause to both A and B . That can cause code explosion. To avoid it we want to select the equation that causes the least code explosion. We therefore select an equation that is present in the maximum number of other clauses.

Consider this example:

$$\begin{array}{ll} \text{match } a \text{ with} & \\ | \text{Add}(\text{Add}(x, y), \text{Zero}) & \Rightarrow e_1 \\ | \text{Add}(\text{Mul}(x, y), \text{Zero}) & \Rightarrow e_2 \\ | \text{Add}(x, \text{Mul}(y, z)) & \Rightarrow e_3 \\ | \text{Add}(x, \text{Add}(y, z)) & \Rightarrow e_4 \\ | \text{Add}(x, \text{Zero}) & \Rightarrow e_5 \end{array}$$

For the outer constructor we have no choice, but for the inner constructors we can either test $\text{Add}(x, y)$ first or Zero first. If we test $\text{Add}(x, y)$ first, then we end up duplicating the e_3 and e_4 clauses. Our heuristic chooses to match on Zero instead, which leads to optimal code without any

duplication:

```
match a with
| Add(a1, a2)  ⇒
  match a2 with
  | Zero  ⇒
    match a1 with
    | Add(x, y)  ⇒ e1
    | Mul(x, y)  ⇒ e2
    | _         ⇒ e5
    | Mul(y, z)  ⇒ e3
    | Add(y, z)  ⇒ e4
    | _         ⇒ ERR
  | _       ⇒ ERR
```

(I've merged subsequent tests of the same variable into a single match – see below)

5 Discussion

This algorithm is a bit different than the algorithms in the literature. It is based on the following observations and considerations.

Some of the literature spends effort on avoiding exponential code explosion, and opts to generate redundant tests instead [Aug85]. Exponential code explosion doesn't occur in practice [SR00]. Therefore the right approach seems to me to never generate redundant tests, and to try to avoid code duplication using a heuristic as e.g. [Mar08] does.

In fact, the literature shows that for real world code, different pattern matching algorithms generate almost identical code [SR00, Mar08]. Our approach here is to (1) always work on matching the first clause first, to avoid unnecessary tests and (2) greedily try to minimize duplication using the heuristic. This isn't exactly what is in the literature, but it combines some aspects of the best heuristics in [Mar08].

Generating binary tests instead of n-way tests that simultaneously branch on several possible constructors, may be necessary to avoid duplication:

```
match (a, b) with
| (A, _)      ⇒ e1
| (_, A)      ⇒ e2
| (_, B)      ⇒ e3
| (_, C)      ⇒ e4
| (_, D)      ⇒ e5
| (B, E)      ⇒ e6
| (C, F)      ⇒ e7
```

If we do a 3-way branch on $a = A, B, C$, then we have to duplicate the middle clauses for both B and C . On the other hand, if we simply test if $a = A$, and then test $b = A, B, C, D, E, F$, then we have no duplication.

Executing an n-way branch may be compiled more efficiently using a jump table than a series of 2-way branches. It is not difficult to detect a series of 2-way branches on the same variable, and convert those back into an n-way branch.

Some of the literature uses decision DAGs instead of decision trees. The easiest way to generate a decision DAG is simply to generate a decision tree and then compress it to a DAG using hash consing [Mar08]. A simpler though less optimal solution is to only do this for the right hand sides, and not for the internal nodes of the decision tree: we could generate a separate basic block for each unique right hand side, and replace the right hand sides by a jump to the corresponding basic block.

To take advantage of types, you can keep track of the remaining possible constructors of each variable. This way you can avoid generating the error cases when there are no remaining possible constructors. You can then do exhaustiveness checking by looking if an error case was generated or not. You can check for unreachable patterns by looking at whether one of the original right hand sides e_i doesn't appear in the decision tree.

6 Code

Scala code that implements the algorithm can be found at <https://julesjacobs.com/notes/patternmatching/pmatch.sc>

Yorick Peterse implemented [a much more advanced version in Rust](#). His version supports integer literals and ranges, or-patterns, guards, and exhaustiveness and redundancy checking. The repository also contains a good README explaining the algorithm and how to implement the advanced features. The syntax “*a is pat*” was invented by Yorick. I initially used “*a = pat*”, but I think Yorick's syntax is clearer.

References

- [Aug85] Lennart Augustsson. Compiling pattern matching. In *Proc. of a conference on Functional programming languages and computer architecture*, pages 368–381, Berlin, Heidelberg, January 1985. Springer-Verlag.
- [Mar92] Luc Maranget. Compiling lazy pattern matching. In *Proceedings of the 1992 ACM conference on LISP and functional programming*, LFP '92, pages 21–31, New York, NY, USA, January 1992. Association for Computing Machinery.
- [Mar08] Luc Maranget. Compiling pattern matching to good decision trees. In *Proceedings of the 2008 ACM SIGPLAN workshop on ML*, ML '08, pages 35–46, New York, NY, USA, September 2008. Association for Computing Machinery.
- [Pet92] Mikael Pettersson. A Term Pattern-Match Compiler Inspired by Finite Automata Theory. In *Proceedings of the 4th International Conference on Compiler Construction*, CC '92, pages 258–270, Berlin, Heidelberg, October 1992. Springer-Verlag.
- [Ses96] Peter Sestoft. ML pattern match compilation and partial evaluation. In Olivier Danvy, Robert Glück, and Peter Thiemann, editors, *Partial Evaluation*, pages 446–464, Berlin, Heidelberg, 1996. Springer Berlin Heidelberg.
- [SR00] Kevin Scott and Norman Ramsey. When Do Match-compilation Heuristics Matter? Technical Report, University of Virginia, USA, 2000.