

# A Tale of Two Unions

Ashwin Menon

December 15, 2024

Bjarne Stroustrup, in *The C++ Programming Language*, recommends avoiding the use of unions. This is of course good advice, mainly because unions are error prone due to the mixing of types that they entail. If you really do need to use unions (for performance or space reasons) then he advises wrapping them in classes. But does that increase the size of the assembly code generated? Let's investigate.

One of my projects was a *BASIC* interpreter, that needed to store variables values. Each value was a union that could hold an integer, a floating point type, and a string (more precisely, a pointer to a string descriptor). This project was written in *C*. Let's see if implementing this union using *C++* classes introduces code bloat. For this experiment, I will simplify the union by omitting the string descriptors, i.e. each variable value can be a float or an int.

First, the *C* code:

```
1 enum ValueType {INT, REAL};
2
3 struct Value {
4     enum ValueType type;
5     union {
6         int integer;
7         float real;
8     } val;
9 };
10
11 #define CONST_INT(C, VAL) \
12     const struct Value C = { .type = INT, .val = { .integer = VAL }}
13 #define CONST_REAL(C, VAL) \
14     const struct Value C = { .type = REAL, .val = { .real = VAL }}
15
16 CONST_INT(a, 42);
17 CONST_REAL(b, 1.5);
18
19 struct Value add_vars(const struct Value *a, const struct Value *b)
20 {
21     struct Value ret;
22     ret.type = b->type;
23     switch(b->type) {
24         case INT:
25             ret.val.integer = a->val.integer + b->val.integer;
26             break;
27         default:
```

```

27     ret.val.real = a->val.real + b->val.real;
28     break;
29 }
30 return ret;
31 }
32
33 void add_a(struct Value *v) {
34     v->val.integer += a.val.integer;
35 }
36
37 void add_b(struct Value *v) {
38     v->val.real += b.val.real;
39 }

```

The union itself is a standard tagged union. But we want to see what happens when we actually use the union. Since variables can be added in *BASIC*, I've added an `add_vars` function that does just that. To keep this simple, it assumes that the values passed in are of the same type (i.e. both `float`'s or both `int`'s.).

*BASIC* has some built in constants (`PI`, for example). To simulate this, I've declared two constant values, an `int` and a `float`. I've also added two functions that take in a variable of the appropriate type and add one of these constants to it.

So what does this generate? I've used Matt Godbolt's excellent online tool to test the results. The code was compiled with the `-O2` optimization turned on. I've liberally commented the code.

A word about the x86\_64 ABI. The compiler here is gcc, so the SystemV AMD64 ABI is used here, which is discussed in this Wikipedia article. In a nutshell, the first six integer parameters to a function are passed in *RDI*, *RSI*, *RDX*, *RCX*, *R8*, and *R9*, while *XMM0*, *XMM1*, *XMM2*, *XMM3*, *XMM4*, *XMM5*, *XMM6* and *XMM7* are used for the first floating point arguments. Additional arguments are passed on the stack. Integer return values upto 64 bits are returned in *RAX*.

```

1 ; The add_vars function
2 ; The parameter 'a' is passed in rdi
3 ; The parameter 'b' is passed in rsi
4 add_vars:
5     ; Store parameter b's type in eax
6     mov     eax, DWORD PTR [rsi]
7     ; Store parameter a's value in edx
8     mov     edx, DWORD PTR [rdi+4]
9     ; Store parameter b's value in ecx
10    mov     ecx, DWORD PTR [rsi+4]
11    ; Check if b's type is INT
12    test    eax, eax
13    jne     .L2
14    ; If b is an INT, then add b's value to a's value using
integer
15    ; addition, and store the result in edx
16    add     edx, ecx
17    ; Move edx into the upper 32 bits of rdx, clearing the
lower 32

```

```

18     ; bits of rdx
19     sal     rdx, 32
20     ; Store b's type in the lower 32 bits of rdx
21     or      rax, rdx
22     ; At this point, rax contains the entire contents of the
    union.
23     ; The upper 32 bits contain the value, and the lower 32
    bits
24     ; contain the type.
25     ; Return the result in rax
26     ret
27 ; If b is a REAL, then add b's value to a's value using floating
    point
28 ; addition
29 .L2:
30     ; Store a's value in xmm0
31     movd    xmm0, edx
32     ; Store b's value in xmm1
33     movd    xmm1, ecx
34     ; Add a's value to b's value, storing the result in xmm0
35     addss   xmm0, xmm1
36     ; Move the result into edx
37     movd    edx, xmm0
38     ; Move edx into the upper 32 bits of rdx, clearing the
    lower 32
39     ; bits of rdx
40     sal     rdx, 32
41     ; Store b's type in the lower 32 bits of rdx
42     or      rax, rdx
43     ; At this point, rax contains the entire contents of the
    union.
44     ; The upper 32 bits contain the value, and the lower 32
    bits
45     ; contain the type.
46     ; Return the result in rax
47     ret
48
49 ; The add_a function
50 add_a:
51     ; The constant 'a' is 42, so just add 42 to 'a' and return
52     add     DWORD PTR [rdi+4], 42
53     ret
54
55 ; The add_b function
56 add_b:
57     ; Load the value of the constant 'b' into xmm0
58     movss   xmm0, DWORD PTR .LC0[rip]
59     ; Add the value of 'a' to 'b'
60     ; Remember, the union 'a' is passed in rdi
61     addss   xmm0, DWORD PTR [rdi+4]
62     ; Store the result in the back into 'a' and return
63     movss   DWORD PTR [rdi+4], xmm0
64     ret
65
66 ; The constant 'b'
67 b:
68     .long   1

```

```

69         .long    1069547520
70
71 ; The constant 'a'
72 a:
73     .long    0
74     .long    42
75
76 ; The constant 'b', again!
77 .LC0:
78     .long    1069547520

```

No surprises here, except for the redundant location *LC0* storing the value of *b*; *addl b* could just have taken the value of *b* from the location *b*.

And now, the *C++* code:

```

1  class Value {
2      enum class ValueType {Integer, Real};
3
4      enum ValueType type;
5      union {
6          int integer;
7          float real;
8      } val;
9
10     public:
11     // Constructors
12     constexpr Value(int i): type{ValueType::Integer}, val{.integer{i}} { }
13     constexpr Value(float f): type{ValueType::Real}, val{.real{f}} { }
14
15     // Getting the value
16     int integer() const { return val.integer; }
17     float real() const { return val.real; }
18
19     // + operator
20     Value operator +(const Value& v) const {
21         switch(v.type) {
22             case ValueType::Integer:
23                 return Value(val.integer + v.integer());
24             default:
25                 return Value(val.real + v.real());
26         }
27     }
28
29     // += operator
30     Value& operator +=(const Value& v) {
31         switch(v.type) {
32             case ValueType::Integer: val.integer += v.integer(); break;
33             default: val.real += v.real(); break;
34         }
35         return *this;
36     }
37 };
38
39 constexpr Value a{42};
40 constexpr Value b{1.5f};

```

```

41 Value add_vars(const Value& a, const Value& b) {
42     return a+b;
43 }
44
45
46 void add_a(Value& v) {
47     v += a;
48 }
49
50 void add_b(Value& v) {
51     v += b;
52 }

```

What does this *C++* code generate? Take a look:

```

1 add_vars(Value const&, Value const&):
2     mov     eax, DWORD PTR [rsi]
3     mov     ecx, DWORD PTR [rdi+4]
4     mov     edx, DWORD PTR [rsi+4]
5     test    eax, eax
6     jne     .L2
7     add     edx, ecx
8     sal     rdx, 32
9     or      rax, rdx
10    ret
11 .L2:
12    movd     xmm0, edx
13    movd     xmm1, ecx
14    mov     eax, 1
15    addss    xmm0, xmm1
16    movd     edx, xmm0
17    sal     rdx, 32
18    or      rax, rdx
19    ret
20 add_a(Value&):
21     add     DWORD PTR [rdi+4], 42
22     ret
23 add_b(Value&):
24     movss    xmm0, DWORD PTR .LC0[rip]
25     addss    xmm0, DWORD PTR [rdi+4]
26     movss    DWORD PTR [rdi+4], xmm0
27     ret
28 .LC0:
29     .long    1069547520

```

The assembly generated from the *C++* source is virtually identical! So much for “*C++* bloat”, eh?

Ok, so if the *C++* code generates identical assembly to *C* code, then what advantages has implementing the union in *C++* given me? Let’s see.

Take a look at how constants are defined in the two regimes. In *C*, we do this:

```

1 CONST_INT(a, 42);
2 CONST_REAL(b, 1.5);

```

That is just ugly; opaque `#define`’s getting in the way of code clarity. In contrast, in *C++*, we do this:

```

1 constexpr Value a{42};
2 constexpr Value b{1.5f};

```

Pretty transparent, I think. We don't even have to specify the type of the constant being declared, as we had to in *C*; the compiler deduces the type, thanks to the two constructors we've put in for the *Value* class.

Next, adding two values. In *C*'s *add* function, we had to do the right thing depending on the union's type tag:

```

1 struct Value add(const struct Value *a, const struct Value *b) {
2     struct Value ret;
3     ret.type = b->type;
4     switch(b->type) {
5         case INT:
6             ret.val.integer = a->val.integer + b->val.integer;
7             break;
8         default:
9             ret.val.real = a->val.real + b->val.real;
10            break;
11    }
12    return ret;
13 }

```

This is a lot of noise. Look at how we add two values in the *C++* version:

```

1 Value add(const Value& a, const Value& b) {
2     return a+b;
3 }

```

Add means Add, and nothing else! All the noise is pushed into the implementation of the *Value* class, so that the user of the class doesn't have to bother with it. You could argue that the *add* function in *C* is doing the same thing, i.e. hiding implementation behind an API. But in *C++*, I don't even need the *add* function; I only added it here to show you the assembly generated by adding two values. In *C++*, I can simply use the *+* operator to add two values.

So in conclusion, a modern compiler will give you a tight implementation of your code, and by using *C++*, you can achieve all the code clarity that using domain specific classes (in this case, the *Value* class) gives you.