

Data Representation

Next, lets add support for

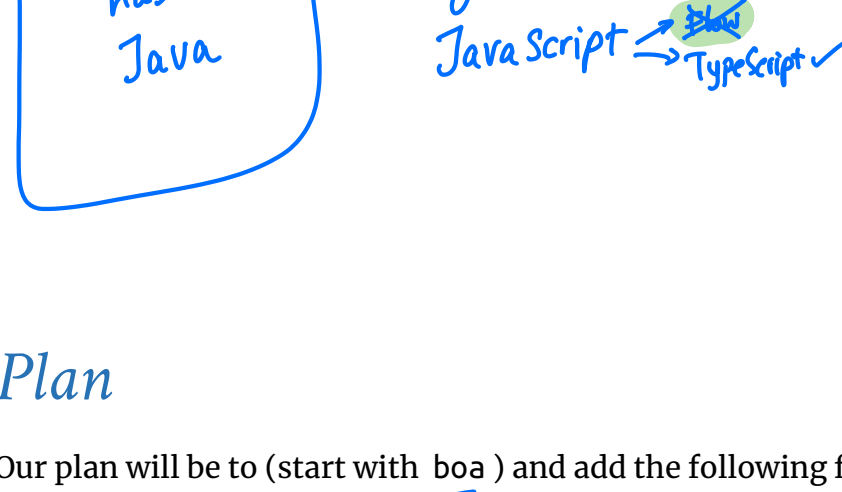
- Multiple datatypes (number and boolean)
- Calling external functions

01 - adder ✓ scores out

03 - cobra due 4/29.

In the process of doing so, we will learn about

- Tagged Representations
- Calling Conventions



Plan

Our plan will be to (start with `boa`) and add the following features:

- Representing boolean values (and numbers)
- Arithmetic Operations
- Arithmetic Comparisons $e_1 < e_2$, $e_1 \geq e_2$, $not\ P$
- Dynamic Checking (to ensure operators are well behaved)

$2 + true$ ✗ $3 < false$ ✗

1. Representation

Motivation: Why booleans?

In the year 2021, its a bit silly to use

- 0 for false and
- non-zero for true .

"avoid confusion", improve readability
support for `8b`, ..

`int` → `int + bool`

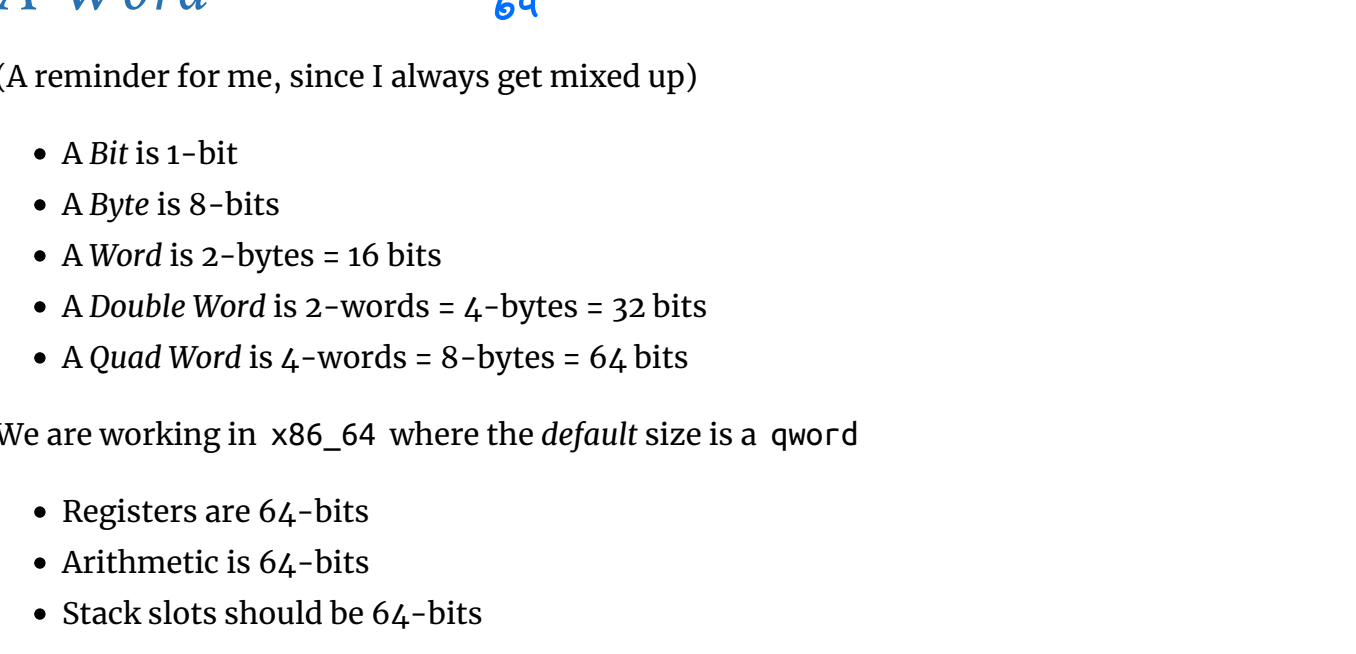
But really, `boolean` is a stepping stone to other data

- Pointers
- Tuples
- Structures
- Closures

The Key Issue

How to distinguish numbers from booleans?

- Need extra information to mark values as `number` or `bool` .



A Word

(A reminder for me, since I always get mixed up)

- A *Bit* is 1-bit
- A *Byte* is 8-bits
- A *Word* is 2-bytes = 16 bits
- A *Double Word* is 2-words = 4-bytes = 32 bits
- A *Quad Word* is 4-words = 8-bytes = 64 bits

We are working in `x86_64` where the *default* size is a *qword*

- Registers are 64-bits
- Arithmetic is 64-bits
- Stack slots should be 64-bits
- etc.

Option 1: Use Two (Quad-)Words

How to distinguish numbers from booleans?

Need extra information to mark values as `number` or `bool` .

First word is 0 means `bool` , is 1 means `number` , 2 means pointer etc.

Value	Representation (HEX)
0	[0x0-----0][0x0-----0]
3	[0x00000000][0x00000003]
5	[0x00000000][0x00000005]
12	[0x00000000][0x0000000c]
42	[0x00000000][0x0000002a]
false	[0x00000000][0x00000000]
true	[0x00000000][0x00000001]

Pros
- 64-bits

Cons
- use 64-bits which is *horrible*
- duplc mstr (but fancy Hw solutions)

Pros

- Can have *lots* of different types, but

Cons

- Takes up *double* memory,
- Operators `+`, `-` require *two* memory reads.

In short, rather wasteful! We don't need so *many* types.

`Int` `Integer`

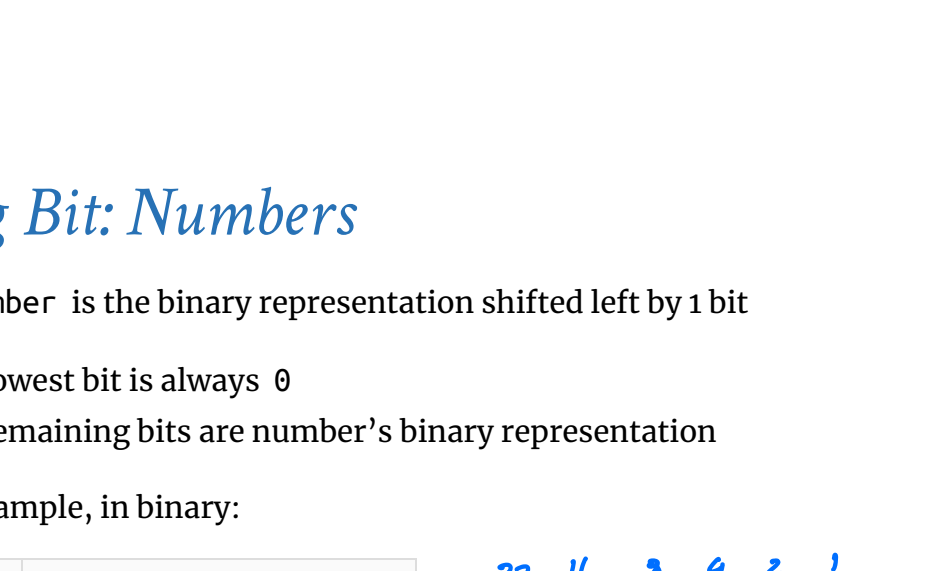
Option 2: Use a Tag Bit

Can distinguish two types with a *single* bit.

Least Significant Bit (LSB) is

- 0 for number
- 1 for boolean

Question: why not 0 for boolean and 1 for number ?



Tag Bit: Numbers

So `number` is the binary representation shifted left by 1 bit

- Lowest bit is always 0
- Remaining bits are number's binary representation

For example, in binary:

Value	Representation (Binary)
3	[0b00000110]
5	[0b00001010]
12	[0b00011000]
42	[0b01010100]

32 16 8 4 2 1
.
.
.
.
.
.
.
.

Or in hexadecimal:

Value	Representation (HEX)
3	[0x06]
5	[0x0a]
12	[0x18]
42	[0x54]

Tag Bit: Booleans

Most Significant Bit (MSB) is

- 1 for true
- 0 for false

For example

Value	Representation (Binary)
true	[0b1001]
false	[0b0001]

Or, in HEX

Value	Representation (HEX)
true	[0x80000001]
false	[0x00000001]

(eliding the 32/8 zeros in the "most-significant" DWORD)

Types

Lets extend our source types with `boolean` constants

```
data Expr a
= ...
| Boolean Bool a
```

Correspondingly, we extend our assembly `Arg` (values) with

```
data Arg
= ...
| HexConst Int
```

So, our examples become:

	Value	Representation (HEX)
Boolean False	HexConst	0x00000001
Boolean True	HexConst	0x80000001
Number 3	HexConst	0x00000006
Number 5	HexConst	0x0000000a
Number 12	HexConst	0x0000000c
Number 42	HexConst	0x0000002a

Transforms

Next, lets update our implementation

The `parse`, `anf` and `tag` stages are straightforward.



Compiler Pipeline

Lets focus on the `compile` function.

A TypeClass for Representing Constants

Its convenient to introduce a type class describing Haskell types that can be represented as x86 arguments:

```
class Repr a where
  repr :: a -> Arg
```

We can now define instances for Int and Bool as:

```
instance Repr Int where
  repr n = Const (Data.Bits.shift n 1) -- left-shift `n` by 1

instance Repr Bool where
  repr False = HexConst 0x00000001
  repr True  = HexConst 0x80000001
```

Immediate Values to Arguments

Boolean b is an immediate value (like Number n).

Lets extend immArg that transforms an immediate expression to an x86 argument.

```
immArg :: Env -> ImmTag -> Arg
immArg (Var x _) = ...
immArg (Number n _) = repr n
immArg (Boolean b _) = repr b
```

Compiling Constants

Finally, we can easily update the compile function as:

```
compileEnv :: Env -> AnfTagE -> Asm
compileEnv _ e@(Number _ _) = [IMov (Reg RAX) (immArg env e)]
compileEnv _ e@(Boolean _ _) = [IMov (Reg RAX) (immArg env e)]
```

(The other cases remain unchanged.)

Lets run some tests to double check.

QUIZ

What is the result of:

```
ghci> exec "15"
```

- A. Error
- B. 0
- C. 15
- D. 30

Output Representation

Say what?! Need to update our run-time printer in main.c

```
void print(int val){
  if (val == CONST_TRUE)
    printf("true");
  else if (val == CONST_FALSE)
    printf("false");
  else // should be a number!
    printf("%d", d >> 1); // shift right to remove tag bit.
}
```

and now we get:

```
ghci> exec "15"
15
```

Can you think of some other tests we should write?

QUIZ

What is the result of

```
ghci> exec "let x = 15 in x"
```

- A. Error
- B. 0
- C. 15
- D. 30

QUIZ

What is the result of

```
>>> exec "if 3: 12 else: 49"
>>> exec "if 0: 12 else: 49"
>>> exec "if true: 12 else: 49"
>>> exec "if false: 12 else: 49"
```

- A. Error
- B. 0
- C. 12
- D. 49

Lets go and fix the code so the above do the right thing!

2. Arithmetic Operations

Constants like 2, 29, false are only useful if we can perform computations with them.

First lets see what happens with our arithmetic operators.

QUIZ: Addition

What will be the result of:

```
ghci> exec "12 + 4"
```

- A. Does not compile
- B. Run-time error (e.g. segmentation fault)
- C. 16
- D. 32
- E. 0

Shifted Representation and Addition

We are representing a number n by shifting it left by 1

n has the machine representation 2*n

0 Number
1 Bool

Thus, our source values have the following _representations:

Source Value	Representation (DEC)
3	6
5	10
3 + 5 = 8	6 + 10 = 16
n1 + n2	2*n1 + 2*n2 = 2*(n1 + n2)

That is, addition (and similarly, subtraction) works as is with the shifted representation.

n1 2^63 2^62
n2 4 2
2^65
n1 * n2

QUIZ: Multiplication

What will be the result (using our code so far) of:

```
ghci> exec "12 * 4"
```

- A. Does not compile
- B. Run-time error (e.g. segmentation fault)
- C. 24
- D. 48
- E. 96

Shifted Representation and Multiplication

We are representing a number n by shifting it left by 1

n has the machine representation 2*n

Thus, our source values have the following _representations:

Source Value	Representation (DEC)
3	6
5	10
3 * 5 = 15	6 * 10 = 60
n1 * n2	2*n1 * 2*n2 = 4*(n1 + n2)

Thus, multiplication ends up accumulating the factor of 2.

- Result is two times the desired one.

Strategy

Thus, our strategy for compiling arithmetic operations is:

Addition and Subtraction “just work” – as shifting “cancels out”,

Multiplication result must be “adjusted” by dividing-by-two

- i.e. right shifting by 1

Types

The source language does not change at all, for the Asm lets add a “right shift” instruction (shr):

```
data Instruction
= ...
| IShr Arg Arg
```

Transforms

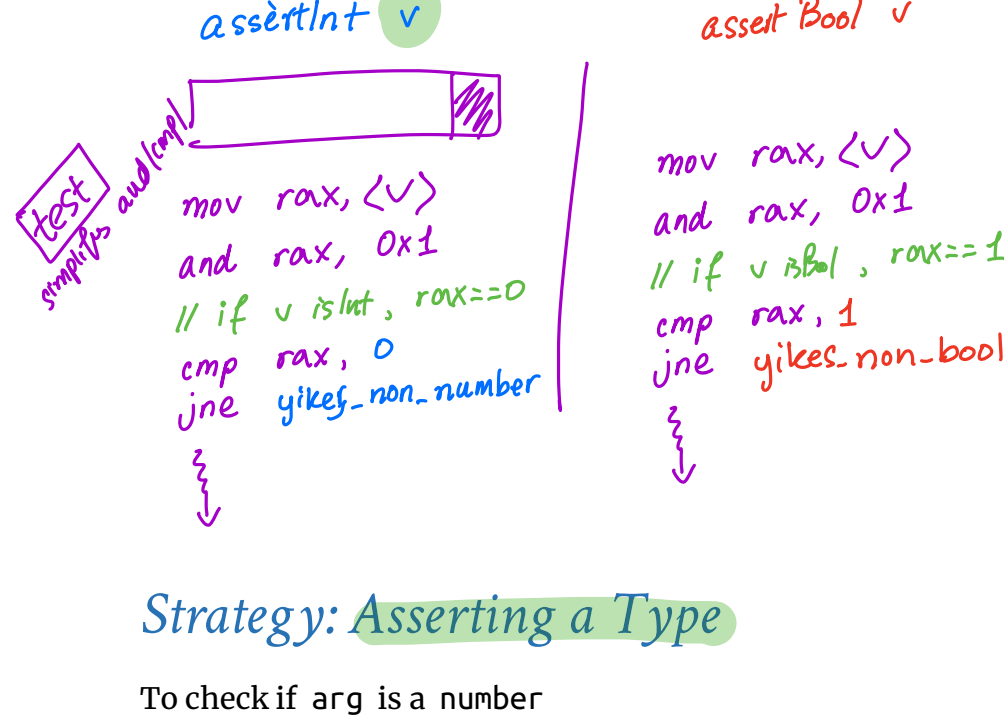
We need only modify compileEnv to account for the “fixing up”

```
compileEnv :: Env -> AnfTagE -> [Instruction]
compileEnv env (Prim2 o v1 v2 _) = compilePrim2 env o v1 v2
```

where the helper compilePrim2 works for Prim2 (binary) operators and immediate arguments:

!	bool	
if	bool	
=	int or bool	int or bool

handler



Strategy: Asserting a Type

To check if `arg` is a number

- Suffices to check that the LSB is 0
- If not, jump to special `error_non_int` label

For example

```

mov rax, arg
mov rbx, rax           ; copy into rbx register
and rbx, 0x00000001    ; extract lsb
cmp rbx, 0             ; check if lsb equals 0
jne error_non_number
...

```

at `error_non_number` we can call into a C function:

```

error_non_number:
    mov rdi, 0           ; pass error code
    mov rsi, rax         ; pass erroneous value
    call error           ; call run-time "error" function

```

Finally, the `error` function is part of the *run-time* and looks like:

```

void error(long code, long v){
    if (code == 0) {
        fprintf(stderr, "Error: expected a number but got %#010x\n",
v);
    }
    else if (code == 1) {
        // print out message for errorcode 1 ...
    }
    else if (code == 2) {
        // print out message for errorcode 2 ...
    }
    ...
    exit(1);
}

```

-3 * 4

Strategy By Example

Lets implement the above in a simple file `tests/output/int-check.s`

```

section .text
extern error
extern print
global our_code_starts_here
our_code_starts_here:
    mov rax, 1           ; not a valid number
    mov rbx, rax         ; copy into rbx register
    and rbx, 0x00000001  ; extract lsb
    cmp rbx, 0           ; check if lsb equals 0
    jne error_non_number
error_non_number:
    mov rdi, 0
    mov rsi, rax
    call error

```

Alas

```

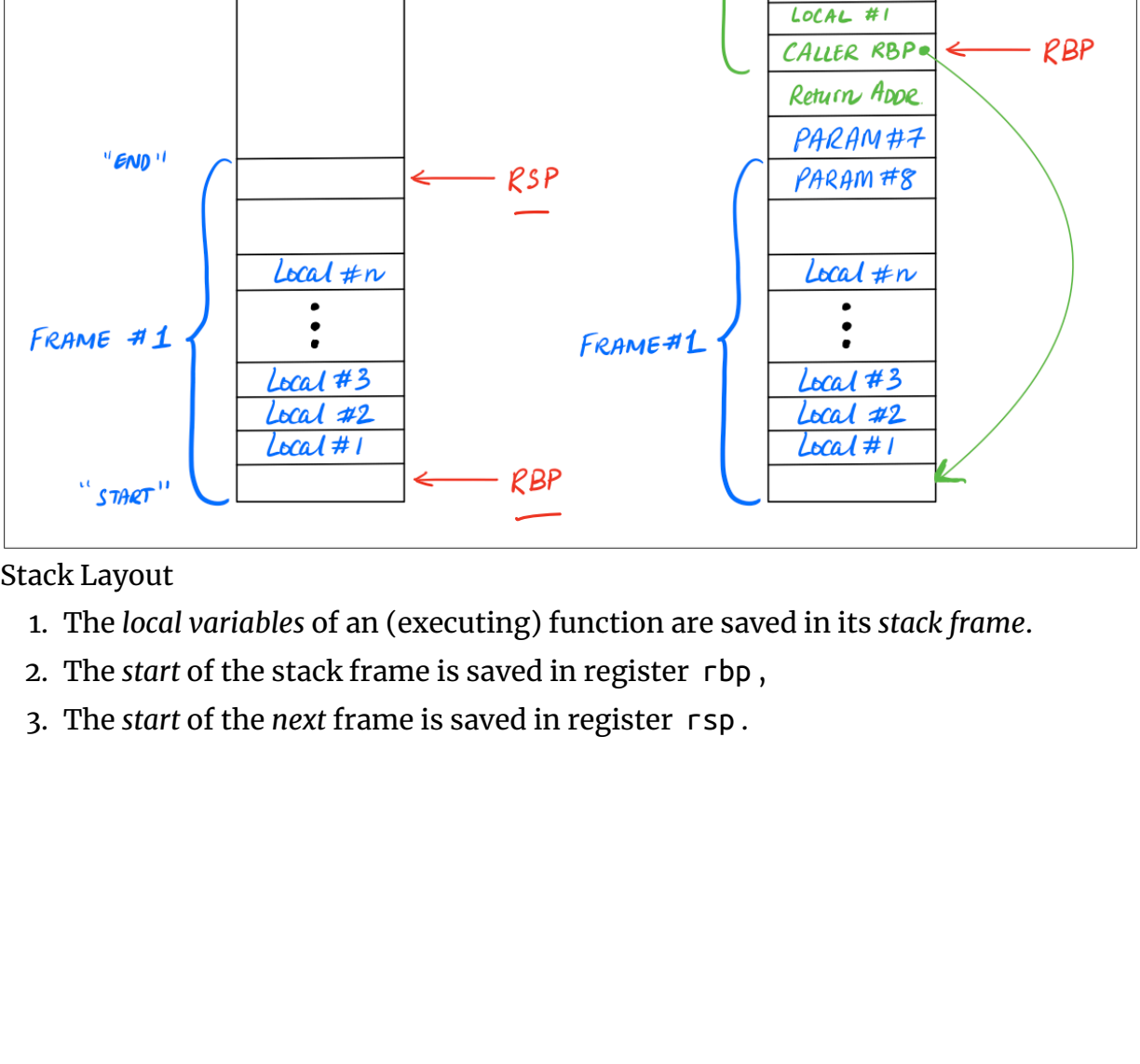
make tests/output/int-check.result
... segmentation fault ...

```

What happened ?

Managing the Call Stack

To properly call into C functions (like `error`), we must play by the rules of the [C calling convention](#)



Stack Layout

1. The *local variables* of an (executing) function are saved in its *stack frame*.
2. The *start* of the stack frame is saved in register `rbp`,
3. The *start* of the next frame is saved in register `rsp`.

Calling Convention

We must **preserve** the above invariant as follows:

In the Callee

At the start of the function

```

push rbp           ; SAVE (previous) caller's base-pointer on stack
mov rbp, rsp       ; set our base-pointer using the current stack-pointer
sub rsp, 8*N        ; ALLOCATE space for N local variables

```

At the end of the function

```

add rsp, 8*N0       ; FREE space for N local variables
pop rbp            ; RESTORE caller's base-pointer from stack
ret                ; return to caller

```

Fixed Strategy By Example

Lets implement the above in a simple file `tests/output/int-check.s`

```

section .text
extern error
extern print
global our_code_starts_here
our_code_starts_here:
    push rbp           ; save caller's base-pointer
    mov rbp, rsp       ; set our base-pointer
    sub rsp, 1600       ; alloc '100' vars

    mov rax, 1           ; not a valid number
    mov rbx, rax         ; copy into rbx register
    and rbx, 0x00000001  ; extract lsb
    cmp rbx, 0           ; check if lsb equals 0
    jne error_non_number

    add rsp, 1600        ; de-alloc '100' vars
    pop rbp            ; restore caller's base-pointer
    ret

error_non_number:
    mov rdi, 0
    mov rsi, rax
    call error

```

Aha, now the above works!

```

make tests/output/int-check.result
... expected number but got ...

```

Q: What NEW thing does our compiler need to compute?

Hint: Why do we `sub esp, 1600` above?

Types

Lets implement the above strategy.

To do so, we need a new data type for run-time types:

```

data Ty = TNumber | TBoolean

```

a new `Label` for the error

```

data Label
= ...
| TypeError Ty      -- Type Error Labels
| Builtin Text      -- Functions implemented in C

```

and thats it.

Transforms

The compiler must generate code to:

1. Perform dynamic type checks,
2. Exit by calling `error` if a failure occurs,
3. Manage the stack per the convention above.

1. Type Assertions

The key step in the implementation is to write a function

```

assertType :: Env -> IExp -> Ty -> [Instruction]
assertType env v ty
= [ IMov (Reg RAX) (immArg env v)
  , IMov (Reg RBX) (Reg RAX)
  , IAnd (Reg RBX) (HexConst 0x00000001)
  , ICmp (Reg RBX) (typeTag ty)
  , IJne (TypeError ty)
  ]

```

where `typeTag` is:

```

typeTag :: Ty -> Arg
typeTag TNumber = HexConst 0x00000000
typeTag TBoolean = HexConst 0x00000001

```

You can now splice `assertType` prior to doing the actual computations, e.g.

```

compilePrim2 :: Env -> Prim2 -> ImmE -> ImmE -> [Instruction]
compilePrim2 env Plus v1 v2 = assertType env v1 TNumber
                                ++ assertType env v2 TNumber
                                ++ [ IMov (Reg RAX) (immArg env v1)
                                    , IAdd (Reg RAX) (immArg env v2)
                                    ]

```

2. Errors

We must also add code at the `TypeError TNumber` and `TypeError TBoolean` labels.

```

errorHandler :: Ty -> Asn
errorHandler t =
[ ILabel (TypeError t)      -- the expected-number error
, IMov (Reg RDI) (ecode t)  -- set the first "code" param,
, IMov (Reg RSI) (Reg RAX) -- set the second "value" param fi
rst,
, ICall (Builtin "error")   -- call the run-time's "error" function.
]

ecode :: Ty -> Arg
ecode TNumber = Const 0
ecode TBoolean = Const 1

```

3. Stack Management

Maintaining `rsp` and `rbp`

We need to make sure that *all* our code respects the [C calling convention](#).

To do so, just *wrap* the generated code, with instructions to save and restore `rbp` and `rsp`

```

compileBody :: AnTagE -> Asn
compileBody e = entryCode e
                ++ compileEnv emptyEnv e
                ++ exitCode e

entryCode :: AnTagE -> Asn
entryCode e = [ IPush (Reg RBP)                                -- SAVE caller's RBP
               , IMov (Reg RBP) (Reg RSP)                      -- SET our RBP
               , ISub (Reg RSP) (Const (argBytes n))           -- ALLOC n local-vars
               ]

where
    n = countVars e

```



```

exitCode :: AnfTagE -> Asm
exitCode e [ IAdd (Reg RSP) (Const (argBytes n)) -- FREE n to
cal-vars
, IPop (Reg RBP) -- RESTORE c
aller's RBP
, IRet -- RETURN to
caller
]
where
n = countVars e

```

the `rsp` needs to be a multiple of 16 so:

```

argBytes :: Int -> Int
argBytes n = 8 * n'
where
n' = if even n then n else n + 1

```

Q: But how shall we compute `countVars`?

Here's a shady kludge:

```

countVars :: AnfTagE -> Int
countVars = 1000

```

Obviously a sleazy hack (why?), but lets use it to *test everything else*; then we can fix it.

4. Computing the Size of the Stack

Ok, now that everything (else) seems to work, lets work out:

`countVars :: AnfTagE -> Int`
StackSize
 Finding the exact answer is **undecidable** in general (CSE 105), i.e. is *impossible* to compute.

However, it is easy to find an *overapproximate* heuristic, i.e.

- a value guaranteed to be *larger* than the than the max size,
- and which is reasonable in practice.

As usual, lets see if we can work out a heuristic by example.

QUIZ

How many stack slots/vars are needed for the following program?

`1 + 2`

A. 0

B. 1

C. 2

$$V_1 \oplus V_2 \Rightarrow 0$$

QUIZ

How many stack slots/vars are needed for the following program?

```

let x = 1
, y = 2
, z = 3
in
x + y + z

```

ANF + compile ENV etc
 what is MAX i such
`MOV [RBP - 8 * i]`

A. 0

B. 1

C. 2

D. 3

E. 4

$$\text{let } x = e_1 \text{ in } e_2 \Rightarrow \max(\{e_1\}, 1 + \{e_2\})$$

QUIZ

How many stack slots/vars are needed for the following program?

```

if true:
let x = 1
, y = 2
, z = 3
in
x + y + z
else:
0

```

$$\text{if } e_0 : e_1 \Rightarrow \max(e_0, e_1, e_2)$$

A. 0

B. 1

C. 2

D. 3

E. 4

QUIZ

How many stack slots/vars are needed for the following program?

```

let x =
let y =
let z = 3
in
z + 1
in
y + 1
in
x + 1

```

A. 0

B. 1

C. 2

D. 3

E. 4

Strategy

Let `countVars e` be:

- The **maximum** number of let-binds in scope at any point *inside* `e`, i.e.
- The **maximum** size of the Env when compiling `e`

Lets work it out on a case-by-case basis:

- **Immediate values** like `Number` or `Var`
 - are compiled *without pushing* anything onto the Env
 - i.e. `countVars = 0`
- **Binary Operations** like `Prim2` o `v1` `v2` take immediate values,
 - are compiled *without pushing* anything onto the Env
 - i.e. `countVars = 0`
- **Branches** like `If v e1 e2` can go either way
 - can't tell at compile-time
 - i.e. worst-case is larger of `countVars e1` and `countVars e2`
- **Let-bindings** like `Let x e1 e2` require
 - evaluating `e1` and
 - *pushing* the result onto the stack and then evaluating `e2`
 - i.e. larger of `countVars e1` and `1 + countVars e2`

Implementation

We can implement the above a simple recursive function:

```

countVars :: AnfTagE -> Int
countVars (If v e1 e2) = max (countVars e1) (countVars e2)
countVars (Let x e1 e2) = max (countVars e1) (1 + countVars e2)
countVars _ = 0

```

Naive Heuristic is Naive

The above method is quite simplistic. For example, consider the expression:

```

let x = 1
, y = 2
, z = 3
in
0

```

`countVars` would tell us that we need to allocate 3 stack spaces but clearly *none* of the variables are actually used.

Will revisit this problem later, when looking at optimizations.

Recap

We just saw how to add support for

- **Multiple datatypes** (number and boolean)
- **Calling external functions**

and in doing so, learned about

- **Tagged Representations**
- **Calling Conventions**

To get some practice, in your assignment, you will add:

1. Dynamic Checks for Arithmetic Overflows (see the `jo` and `jno` operations)
2. A Primitive `print` operation implemented by a function in the `c` run-time.

And next, we'll see how to add **user-defined functions**.

