



# CRAFT DEVELOPMENT DIARY

Daniel McCarthy



# Introduction

This is a development diary of the “Craft” compiler also known previously by the name “Goblin”.

The “Craft” compiler is designed to compile “Craft” code a language which I have designed.

Unfortunately, I have only documented the project since 25<sup>th</sup> July 2016 so there is months of work undocumented. Anyhow it is titled day 1 and onwards even though it was not day 1 of the project it was day 1 of the diary.

Throughout this document you will see all the struggles and design choices I have made to get the Craft compiler completed.

# Information

Before continuing it is important that you understand the names used in this diary, I will help explain the ones an experienced reader may have trouble with.

For anything else such as trouble understanding what a Lexer, Parser, virtual method, or virtual pure method is. Please research on them to understand this diary.

## Code generator

The code generator is the base class for all code generators, it parses the abstract syntax tree from the root and invokes methods in its inheritor/child depending on the branch it finds in the tree. For example, if it should find an assignment such as “x = 50” then it would invoke a method on its inheritor for handling assignments.

## Goblin Bytecode generator

The Goblin bytecode generator extends the “CodeGenerator” class its methods get invoked by the code generator depending on the branch the code generator has found. Depending on the method that was invoked the goblin bytecode generator would write data to the stream or do further iteration of the AST (abstract syntax tree) until it finds what it is looking for.

## Craft bytecode generator

The Goblin bytecode generator was renamed to Craft bytecode generator so throughout this diary know that they are the same.

# Day 1 - 25<sup>th</sup> July 2016

This is the first day of documenting the Craft compiler. Unfortunately, I left it a bit late so I am actually writing this on the 26<sup>th</sup> July 2016 however I will write about the 25<sup>th</sup> as if it is the 25<sup>th</sup>.

Today I was working more on implementing arrays, at this point arrays were coming along well, I could only support 1 dimensional arrays although up to 3 dimensions can be parsed. I managed to make array access in expressions possible, although not perfect it worked, I was in the process of fixing a bug where array access would only work in expressions if it was on the left hand side, this was because the "A" register was used and if something else was on the left hand side of the expression, let's say a number it would occupy register "A" so the array access on the right hand side of the expression would generate byte code that would have overwritten the number at run time. I was in the process of fixing this or I did fix it I don't remember and then all of a sudden my computer crashed. When I booted back up I found my entire goblin bytecode generator NULLED. I searched my git repository for the most recent commit I could find and unfortunately the commit I found was on the 7<sup>th</sup> of July. The code that was present at that time would be completely obsolete a lot of development has been done since then so I chose to write the goblin code generator again. I rewrote a basic template for the goblin bytecode generator and I worked a little bit on the code generator improving the design.

# Day 2 - 26<sup>th</sup> July 2016

Today early hours in the morning I was focusing on improving the design of the code generator, as well as implementing more code on the goblin bytecode generator since the file was NULLED yesterday due to a system crash.

The work early hours in the morning today was continued from yesterday I continued into the 26<sup>th</sup> I usually do work late.

I decided that I should modify the code generator to use a “Scope” class instead of all the functionality defined in the code generator. This made sense as it is certainly possible to have child scopes which would have caused problems in my current design. So I created a “Scope” class and this class contains variables as well as methods for working with these variables such as creating new variables or getting variables by name. A lot of the related content was moved from the code generator into the “Scope” class and modified. The “Scope” class has the ability to have one parent and one child although this will be changed later to allow for multiple children.

I also removed a lot of obsolete methods from the code generator, the old design approach was to have a virtual pure method for every possible scenario, for example I would have a method for handling a number, and let the code generator invoke a virtual pure method every time a number was found. The child class would then be able to write the machine/byte code to the stream.

This design was not appropriate and my new design approach is to have a method for only more structured elements such as functions, and assignments.

At the moment only the assignment method is implemented and was implemented a few days ago, the function method does not exist yet, in its place is the “scope\_start” method. This however also needs to be removed as every scope can act differently. This will be replaced with a method such as “function\_start” and this would be for function declarations including the function scope.

I also fixed the “cleanBranch” method in the Parser class. Originally it was not working or was not working as expected. I made it clean the tree of branches who only have one child starting at the branch passed as the first argument. This meant that branches with only one child are replaced with its only child, leading to a cleaner tree with information that is no longer needed removed.

I thought about the parser rules today and how I perhaps should change them, for example my expressions will have three branches the first one being the left operand, the second being the operator and the third being the right operand, this could be changed allowing the parent branch to be the operator and then its two children being the left operand and the right operand. This is how most designs look when I search for “Abstract syntax tree” on google images. There will be no definite change yet as I need to outline the pros of this change and cons if there is any.

## Day 3 - 28<sup>th</sup> July 2016

Today I made the Exception class extend one of C++'s standard exceptions "logic\_error". My reason for this is that should an exception be thrown and not caught I or others will be able to see the output in the console window.

I also fixed the design problems caused by the rule changes and the parser update where I fixed the "cleanBranch" method. I also made the ability to have multiple children in the Scope class which is required because scopes can have multiple scopes inside of them.

I have considered changing certain parser rules to output different types of branches than it does at the moment. Google image searches have given me many results with expressions that look like this:

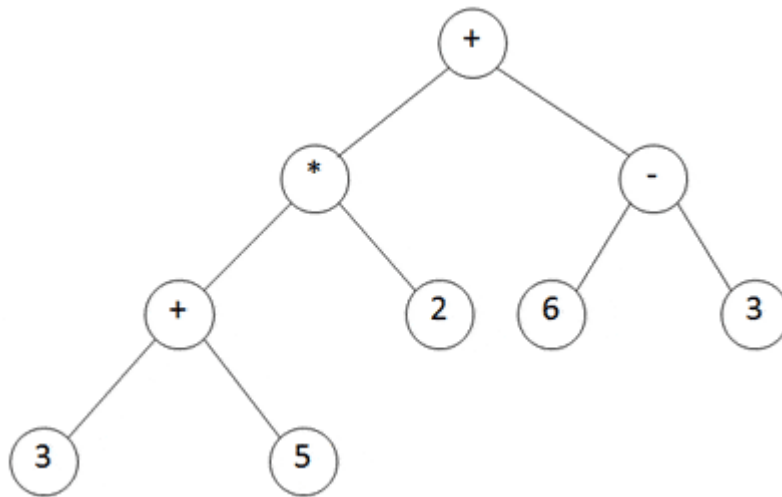


Fig 1.

My system does things a little differently and I need to consider whether or not the above image would be the best way to go about things. My tree is workable but I do believe the above image represents a better structure, the parser rules will need to be thought through carefully and not just the ones for expressions. I partially think that because I am not doing a similar structure as shown in Fig 1. May be partly the reason certain expressions cause parsing errors but I cannot be sure as of yet I am still in a bit of a grey area with this situation.

Another bug was discovered in function calls, a function call cannot have arguments that are also function calls, this is a parsing error.

In other news I attempted to update the parser to do a look ahead first before choosing a rule as the fact it does not do a look ahead has already caused me countless problems and requires rules to be placed in a pacific position to work correctly. Although today I have failed to implement the look ahead feature as the previous parser design would loop through all the rules and then apply it to every branch in the root, so in the first phase for example "3 + 3 + a" would become "E + E + E" and then as it goes to other rules would be broken down further. A look

ahead would not be possible with this design. I changed the design to apply the rules on a left to right basis without applying it to every branch at once, for example "3 + 3 + a" would become "E + 3 + a" now because there is no rule for "E + 3 + a" only "E + E" the parser will now fail to generate an abstract syntax tree, more thinking will need to be done but I plan to make the parser better before continuing the project as if I do not it may cause me even worse problems later on.

One possible solution that I would rather avoid because it might come across as a bad design perhaps, is the ability to state that a particular rule can be any of the following branches. E.g the "E" branch can be the "identifier" token or the "number" token.

## Day 4 - 29<sup>th</sup> July 2016

Today I made the parser a look ahead parser which will help prevent rule conflixtions, for example two rules may exist, "E:operator:E" and "E:operator:E:operator", The second rule is a rule that you probably would never implement but I will use it as an example as it's the perfect scenario. Now take the expression "E+E+" you would expect this expression to match the second rule but in the first parser this did not happen it would pick the first rule that it matched to so it would be "E:operator:E". Now that the parser is changed it will now look ahead and find the most appropriate rule so it would pick "E:operator:E:operator" which is the correct rule.

Upon implementing the new Parser, I removed some code that should be their this caused an issue and prevented branches from being excluded from the tree, this was a simple bug to fix and I fixed it quickly.

Since the parser is now more efficient I no longer need the hash tag to represent functions and function calls I will be considering soon weather to remove them or keep them the same just to make the language a little different.

In other news I am considering rewriting the parser rules all together to try and do a better job, or at least rewrite some of them, also another problem was found and that is with the cleaning system for the parser. This system does exactly what it is supposed to do but in certain situations this would cause problems.

For example, the fourth branch in the "FUNC" branch holds a bunch of statement branches, now because of this cleaning system if only one statement in the function body exists then the cleaning system will replace that "STMT" branch with its child branch, I do not like this at all as obviously the function body should have some sort of root branch for its self. Anyway more thought will need to be given before any change is made.

If I do make this change, then I would set it up a bit similar to the way you exclude branches in the parser rules. In the parser rules you exclude branches from the tree by using the quote character: " '" in the current rule.



# Day 5 - 30<sup>th</sup> July 2016

Today I was working on fixing a bug with input such as “c = (50 + 19) \* 8” or “c = 50 + 19 \* 8”. This bug is a parser bug where the parser cannot break down the branches correctly.

This issue occurred because of the parser rules and possibly the parser its self. Take these rules for example.

*E:identifier*

*E:number*

*E:E:operator:E*

*E:'symbol@(:E:'symbol@)*

*ASSIGN:E:'symbol@=:E*

The following happens:

*C = 50 + 19 \* 8*

*E = 50 + 19 \* 8*

*E = E + 19 \* 8*

*ASSIGN + 19 \* 8*

As you can see the parser is left with “+ 19 \* 8” unparsed, this is because the ASSIGN rule matched correctly while the other rules did not. A look ahead exists but it will not work with branches who have not yet become a child because of this the parser is not working correctly.

Possible solutions

1. Change the parser rules
2. Make the parser branch look what is ahead and then see if it matches a different rule

I would rather implement solution 2. Although I am not sure if it is the best solution yet.

Priority is now changing the parser to do just this.

## Day 6 - 31<sup>th</sup> July 2016

Today I decided to start rewriting the majority of the parser, I will not rewrite the rule system but the actual AST(Abstract Syntax Tree) generator. The problem with the old design was that it did not follow standard procedures to generating AST's and because of this the code was a bit ugly. I will now be using stacks instead of vectors, this did cross my mind in the past but I chose to go with the vector.

I hope to implement a shift reduce parser, properly this time.

## Day 7 - 9<sup>th</sup> August 2016

Today I worked more on rewriting the parser, the development was coming along well but unfortunately I have some segmentation faults that I need to fix. The parser code is not quite right but I feel that I am nearly there. Once complete this will be the 3<sup>rd</sup> time I have wrote the parser.

## Day 8 - 11<sup>th</sup> August 2016

Today not much was accomplished, I fixed a bug where I was creating a new instance of an array with a size of zero, this was because I did not check that functions existed before creating the array, e.g “new function[total\_functions]” if the “total\_functions” variable is zero then it will attempt to create a new function array of 0 elements, this causes issues. I believe that similar issues still exist in the system which I plan to clean out once the parser is complete.

As for the parser its self it is not coming along well at all and I am considering going back to the old working design, although not the best design it did work rather well.

I appear to be having many segmentation faults while implementing the new parser design and it is holding me back.

## Day 9 - 16<sup>th</sup> August 2016

Today I successfully created a shift reduce parser that uses a stack, at the moment it does not do a look ahead and is incomplete as it also is not recursive when it comes to reductions.

I hope to soon make the parser do a look ahead and do recursions on reductions until no more reductions can be done for the current stack state.

## Day 10 - 17<sup>th</sup> August 2016

Today I made the shift reduce parser recursive when it comes to reductions, this was simple to do, I simply just recall the same method when a reduction has been successful, this in turn will attempt to reduce the stack further.

I also attempted to implement the look ahead but was unsuccessful.

# Day 11 - 18<sup>h</sup> August 2016

Today I was attempting to implement a shift reduce parser and during my attempts I realised that I was popping from the back of the stack and not the front, I was curious to why the branches were getting created in order and it is because of that reason.

This setup worked fine for a stack that did not need a look ahead but in a real situation a look ahead is required.

I decided to have a cigarette and thought through the best approach I could do to get this working and the solution I came up with is to make the rule requirements also a stack, and make the start of the input get popped and then pushed to the parser stack so it's a left to right parser rather than a right to left. Then finally during checking for reductions I would pop from the back of the parser stack so its right to left.

Take the following example

Rule requirements: identifier:operator:number

Now take the following input

A + 9

Now what would happen is during every reduction attempt we pop from the right most symbol, and then we pop from the right most rule requirement. So as you can see the rule requirement is "number" and the right most input is also a "number".

This is why this sort of setup seems to work.

Now as for the look ahead symbol since we are now pushing data onto the parse stack left to right instead of the previous right to left, the look ahead will always be the next symbol to the right.

I started writing code that would check if there was a rule match with the parser stack plus the look ahead symbol and if it was it would shift the look ahead symbol to the parser stack and then reduce with the rule it matched with. However, if there was no match it would then try to match a rule with the parse stack alone.

There was a serious problem with this design. I fixed this problem by making it only choose the parser stack plus the look ahead symbol reduction should there be a match for the parser stack alone and also the parser stack plus the look ahead symbol. This solved the problem and my parser could work well with a look ahead symbol, however I then set the following rules:

*E:identifier*

*E:number*

*E:E:operator:E*

*T:E:operator:E:E*

and even tried

*T:E:operator:E:identifier*

I may have also tried

*T:identifier:operator:identifier:identifier*

Once I tried those rules everything went terribly wrong, it appears the reduction of “identifier” into branch “E” and “number” into branch “E” causes issues with the current parser logic.

I am also concerned about the cleanliness of the parser and the approaches I am taking.

It appears I am back to the drawing board.



## Day 12 - 6<sup>th</sup> September 2016

Today I worked on the parser and made the parser work with a look ahead properly I also made the parser reduce the correct way around rather than the previous opposite. Due to the way a stack works previously the branches were being written backwards this is fixed now. I also fixed a few bugs in the parser and partially implemented the ability of custom branches being set based on rules. E.g the newly reduced branch based on rule "V\_DEF" becomes the "V\_DEF" branch in memory.

I also changed the ParserRule class and made a method that returns the total requirements of its self.

# Day 13 - 7<sup>th</sup> September 2016

On the 7<sup>th</sup> of September I ran into problems trying to implement an infinite operator due to the design of the parser. The infinite operator would basically say this rule applies if there is one or infinite of a particular branch.

I then realised at that point that I am going to waste months more trying to get this parser to work so I decided to ditch the rule system all together and that's exactly what I did. I started rewriting the entire parser and this time the parser would generate branches through code rather than rules.

## Day 14 - 8<sup>th</sup> September 2016

On the 8<sup>th</sup> of September I attempted to implement expressions in the newly designed parser but sadly failed even after 3 or 4 attempts.

## Day 15 - 9<sup>th</sup> September 2016

Today I managed to implement expressions successfully. I also began work on “IF” statements but it did not work, I forgot to invoke the “process\_if\_stmt()” method and will do it tomorrow.

# Day 16 - 10<sup>th</sup> September 2016

Today I managed to partially implement “IF” statements in the parser, although “ELSE IF” and “ELSE” statements are currently not present.

For the parser I plan to soon make a few more methods for tasks that I believe are repeated quite frequently. This will improve design and partially increase the speed of code development.

## Day 17 - 11<sup>th</sup> September 2016

Today I implemented the “ELSE” statement but it does not feel right so I will rewrite it differently. Ideally the “ELSE” statement and “ELSE IF” should nest into other “IF” statement children. Currently because I implemented “ELSE” statements I did not consider nesting, nor would their need to be, but with an “ELSE IF” setup the final “ELSE IF” should have the “ELSE” branch.

# Day 18 - 13<sup>th</sup> September 2016

Today I worked on the parser and everything written for this date is parser related.

I made the “IF”, and “ELSE IF” statements, I also changed the “ELSE” statement. All these statements nest within each other.

I also fixed a bug in the parser where any keyword would be valid for a variable. You could literally type "else var\_name" and the parser would register a variable declaration branch. This is completely illegal and should not be allowed.

I also fixed a bug in the parser where expressions with no valid token caused a segmentation fault. This bug fix did not go as planned and introduced another bug which I then fixed. I do not remember what the bug was.

I created methods that are designed to increase readability and faster development time, I then changed some code in the parser to make it use these methods. I decided to do this because I was writing two statements where it could have been done in one function call.

I implemented the ability to set declared variables to a value while declaring them. This allowed variables to be defined and set at the same time. For example: “uint8 x = 50;”

I also added some code to check for a semicolon after a variable declaration in body statements, this is important as otherwise we cannot guarantee there is a semicolon and the programmer would receive syntax errors that may not make sense.

I also implemented structures in the parser.

I also implemented pointers in the parser although they still do not work in the global scope yet as I want to think about the cleanest way to go about it.

I finally implemented the ability to get addresses of entities using the “&” operator. However, this currently does not work in the global scope and currently does not work when defining and setting the structure variable for example: “struct test a = &b”

# Day 19 - 14<sup>th</sup> September 2016

On the 14<sup>th</sup> of September I implemented the “include” macro. This allows people to include other files in their source files. Currently it is possible to crash the program upon including the same file or including something that will lead to the same file being included again. I hope to fix this soon.

I also made the ability to load code generators from DLL files.



## Day 20 - 15<sup>th</sup> September 2016

On the 15<sup>th</sup> of September I changed the design a bit and made the Craft compiler its own DLL rather than an executable file. I had to do this due to problems linking other DLL files with the compiler. Since the compiler its self is now a DLL file I have an executable file interfacing with it.

Obviously these DLL files could be Linux equivalents in the future but for now it will work with Windows only.

I also was implementing expressions in the newly created code generator called “8086CodeGen” its target code is for the 8086 processor. These expressions did not work as I did not implement the order of operations the parser.

# Day 21 - 16<sup>th</sup> September 2016

On the 16<sup>th</sup> of September I made the parser prioritise multiplication and division before other forms of mathematics, this is known as the order of operations. Before this was implemented problems were present in the code generator when generating expressions, as no order of operations was acknowledged this lead to different results than you would get from a calculator.

## Day 22 - 17<sup>th</sup> September 2016

On the 17<sup>th</sup> of September I was writing assembly language in an 8086 emulator to try and find the best code design for functions, this was proven difficult as a lot of variables come into play. For example, a function may have many variables but they may get assigned to either scope variables, argument variables or global variables. The way to access each one of these is completely different from each other. I did not work out the best way to do this on the 17<sup>th</sup> as it passed midnight.

# Day 23 - 18<sup>th</sup> September 2016

Today early in the morning continuing from last night, I managed to successfully handle expressions, assignments and functions in the 8086 code generator. I spent a long time compiling code in other compilers and reading the assembly output to try and understand it. The “8cc” compiler assisted me as well as an online C, C++ compiler that pacifically shows you assembly output.

Eventually after studying these compilers I managed to implement a design of my own, although things are done a bit differently which I am not yet sure if it will be a problem in the future.

Should I want to access an argument variable of my own function I access these variables in 16 bit words relative to the base pointer. The offset is “+4” for argument one, “+6” for argument two and so on.

Should I wish to access a scope variable it’s the opposite way, I minus the base pointer starting at offset “-2”. So for example scope variable “a” our first variable would be located at “bp-2” and our second variable “b” would be at “bp-4”.

Global variable access is currently not permitted inside a scope. Functionality for this will be added shortly.

Here is example code as well as assembly offset of early this morning’s accomplishment

```
uint8 main(uint8 a, uint8 test)
{
    uint8 b;
    b = 10;

    uint8 c;
    b = (50 * b);
    return b;
}
```

## Assembly output

```
_main:
push bp
mov bp, sp
mov ax, 10
mov [bp-2], ax
mov ax, 50
mov bx, [bp-2]
mul bx
mov [bp-2], ax
mov ax, [bp-2]
pop bp
ret
```

In other news while writing code for the code generator I wrote code that could turn a positive to a minus and I forgot that this was legal in most programming languages, Take the expression

`-(50 + 10)`

The result of that expression should be “-60” the parser currently requires that all operators have a left and right operand, I will need to tweak this design as I would have had to anyway due to the Boolean operator “!” used to turn a false to a true or a true to a false. I will attempt to implement this in the parser and code generator around the same time I implement the Boolean operator “!”

## Day 24 - 19<sup>th</sup> September 2016

I do not remember exactly what I done on the 19<sup>th</sup> of September but I know it was either the 19<sup>th</sup> or 20<sup>th</sup> that I successfully implemented function calls.

I have decided that I am going to write diary entries as soon as possible for now on, to prevent forgetting important information.

## Day 25 - 20<sup>th</sup> September 2016

On the 19<sup>th</sup> or 20<sup>th</sup> of September 2016 I successfully implemented function calls in the code generator. Any function arguments are pushed to the stack before the function is called.

I also started working on making pointers and accessing memory with these pointers. I did not manage to finish this on the 20<sup>th</sup> as it passed midnight.

## Day 26 — 21<sup>st</sup> September 2016

On the 21st of September I was just still trying to implement pointers.



# Day 27 — 23<sup>rd</sup> September 2016

19:33:

I managed to implement pointer access in the parser this allows one to do things like “return \*a;”, “\*a = 50;” and more.

While implementing this pointer access I realised that the current design of the parser especially in the expressions must be improved and by the looks of it I do not believe it is compatible with structure access. I hope to improve this and implement structure access in expressions.

20:44:

I implemented structure access in the parser, it turns out the current parser design is actually compatible with structure access. The design of the parser is not too bad but improvements could be made in the future.

# Day 28 — 24<sup>th</sup> September 2016

03:18:

In the 8086 code generator I implemented the ability to access a value that a pointer is pointing to and I also implemented the ability to set the value of the memory a pointer is pointing to.

I have still not made the appropriate classes for some of the branches created for this process.

I am also considering soon making a method for setting registers and the 8086 code generator will store these registers in memory so should a register be set to the exact same thing again it will ignore it, although more thought needs to be taken as if I am not careful I may not reset it during appropriate instructions resulting in problems. It may be best to just finish the code generator and make these amendments later on.

Source code

```
uint8 test(uint8* a)
{
    uint8 c;
    *a = *a + 5;
    return *a;
}
```

```
uint8 main(uint8 a)
{
    uint8 b;
    b = test(&a) + 2;
    return b;
}
```

Assembly output:

```
_test:
    push bp
    mov bp, sp
    mov bx, [bp+4]
    push bx
    mov bx, [bp+4]
```

```
mov ax, [bx]
pop bx
mov cx, 5
add ax, cx
mov [bx], ax
push bx
mov bx, [bp+4]
mov ax, [bx]
pop bx
pop bp
ret
_main:
push bp
mov bp, sp
mov ax, bp
add ax, 4
push ax
call _test
add sp, 2
mov cx, 2
add ax, cx
mov [bp-4], ax
mov ax, [bp-4]
pop bp
ret
```

# Day 29 — 25<sup>th</sup> September 2016

23:09:

In the 8086 code generator I have been trying to implement compare expressions in the compiler as these are required before “if” statements will be possible.

I managed to make compare expressions possible such as “a == 5” this will set register “AX” to “1” if variable “a” is equal to five or “0” if variable “a” is not equal to five. “1” can be expressed as “true” and “0” can be expressed as “false”.

## Source code

```
uint8 main(uint8 a)
{
    a = 5;
    a = a == 5;
    return a;
}
```

## Assembly code

```
_main:
    push bp
    mov bp, sp
    mov ax, 5
    mov [bp+4], ax
    mov ax, [bp+4]
    mov cx, 5
    cmp ax, cx
    jne _0
    mov ax, 1
    jmp _1
_0:
    mov ax, 0
_1:
```

```
mov [bp+4], ax
```

```
mov ax, [bp+4]
```

```
pop bp
```

```
ret
```

Currently multiple compare expressions in sequence are not implemented, e.g “a == 5 && b == 6”.

There is a slight problem with the order of operations as it does not prioritise the order of operations for these compare operators. This leads expressions such as “a > 5 + 6” to check if variable “a” is above “5” and then add “6” this causes the compiler to generate pointless instructions that are overwritten at the end of an assignment.

I hope to allow multiple compare expressions in a sequence and fix the order of operations soon so that I can move onto developing the “if” statement.

Finally, a bug was found, the ability to declare a variable and set it at the same time is currently not possible: “uint8 b = 8;” I probably just forgot to add this functionality and will add it soon.

# Day 30 — 28<sup>th</sup> September 2016

20:48:

Today I have found a bug in the parser where parsing pointers causes the expressions that use the multiplication symbol to cause an infinite loop.

The solution may be to keep track of pointer variables in the parser and let that determine how the tree is generated.

20:56

I have decided to just use the “@” symbol to represent pointer access for now as this will fix the problems I am facing. I do intend to eventually change it back to the “\*” operator and fix this issue properly. I believe that similar issues will exist with the “&” operator as this can be used as a bitwise “&” therefore it is relevant for expressions so this too will be changed to a question mark symbol for now “?”. Again I do intend to fix both of these issues in the future.

21:27

Successfully implemented compare expressions in the parser.

22:32

I have successfully implemented compare expressions in the 8086 code generator, the only problem which has now been fixed appeared to be the use of the “>”, “<”, “>=” and the “<=” operations. It did not perform well when it was its opposite so “>” would check if it is below it did not like that at all. I have since made them and only them use the correct instructions for their operators although I believe this will cause problems in the future when multiple compare expressions are present, for example “a > 4 && b < 6”

# Day 31 — 29<sup>th</sup> September 2016

05:23:

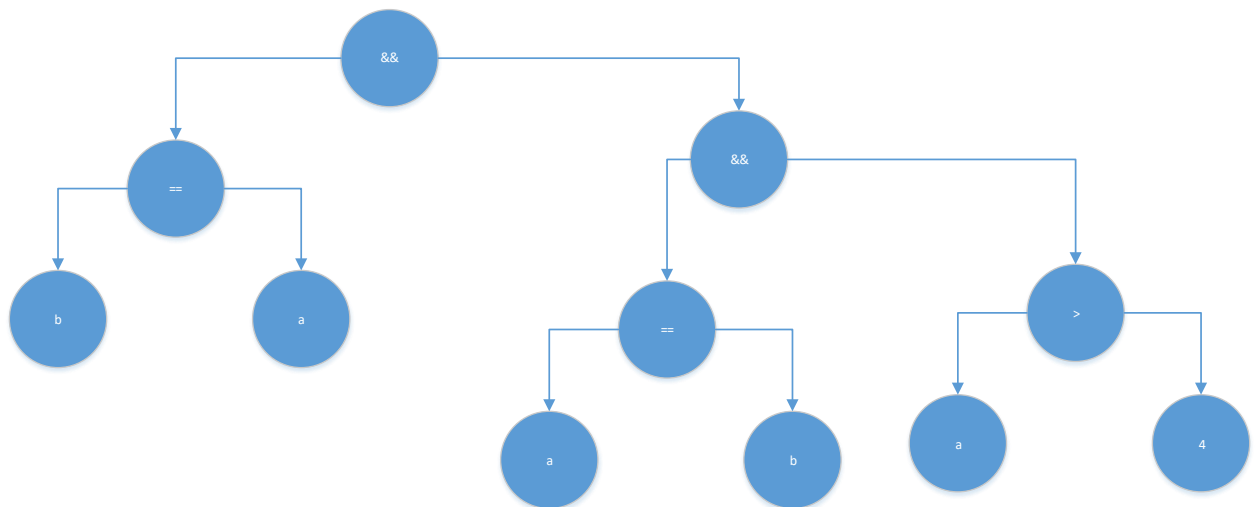
In the parser I have successfully implemented the ability for multiple conditions in an expression e.g “a == b && c > 1”.

I ran into many problems implementing this and it lead to me having to rename the “process\_expression” method to “process\_expression\_part”. I then created a new “process\_expression” method that would call “process\_expression\_part” if it needs to.

In the “process\_expression” method the system first calls the “process\_expression\_part” method after it has returned it will check if the next token is a logic operator if it is then it will then need to invoke the “process\_expression\_part” method again and then put it together creating the “E” branch with the value being that of the logical operator. This new “E” branch will contain the children of the left and right expressions.

Take the following expression “a == b && c > 1” the left expression is “a == b” and the right expression is “c > 1”.

**Graphical representation: b == a && a == b && a > 4**



# Day 32 — 1st October 2016

18:58:

Fixed a bug in the 8086 code generator where function calls did not push function parameters on backwards, before this patch function parameters would be accessed the wrong way around from within a function body.

The other possible solution would be to change how the function body accesses data passed to the function.

The reason pushing function parameters on backwards fixed the problem is because of how the 8086 stack works.

The “gcc” compiler also pushes function parameters on backwards to achieve the same thing.

19:31

In the 8086 code generator I implemented the ability to declare variables and assign them at the same time, previously this was not possible. You can now do “uint8 a = 50;”

21:23

In the 8086 code generator the ability for multiple conditions in an expression e.g “a == b && c == 1” is now possible. However, the expression “a == b || c == 1” is currently not possible.

In other news due to the fact that multiple expressions cause a “push” and a “pop” to happen should the statement be “false” the value will never get popped back off the stack which could be a problem but I don’t know yet.

23:41

It turns out the “push” and “pop” statements generated talked about at “21:23” has become a problem, functions will not return to the correct address in memory due to the old value not being popped off. I plan to fix this soon; at the moment I have just commented them out but they are required or another alternative will need to be factored.

23:48

Multiple compare expressions in the 8086 code generator with use of the “||” logical operator are now possible. This means both logical operators “&&” and “||” are implemented.

23:56

Fixed a bug in the 8086 code generator, the bug was where I forgot to turn the compare expression flag off once it was complete. I also changed a bit of code and tidied things up a bit.



# Day 33 — 2nd October 2016

00:25:

I have just fixed a bug in the 8086 code generator where multiple expressions would always default to the “add” instruction. This is wrong as an expression such as “(80 + 20) \* (20 + 40)” would act as if it was “(80 + 20) + (20 + 40)”.

00:34

I have fixed a bug that was introduced to the previous bug fix at “00:25”. Essentially the “make\_math\_instruction” method should not be called for logical operators.

00:42

I have fixed a bug where “push” and “pop” instructions were generated during compare expressions. This was explained in detail yesterday.

03:22

In the 8086 code generator I forgot to fully implement all of the possible compare expressions correctly this is fixed now.

03:46

I managed to implement “IF” statements in the 8086 code generator, currently “ELSE” and “ELSE IF” statements are not implemented.

# Day 34 — 4th October 2016

01:43:

Successfully implemented “ELSE IF” statements in both the parser and the 8086 code generator. The “ELSE IF” statements do not use a branch class of their own as they are no different from “IF” statements property wise.

02:01:

Successfully implemented “ELSE” statements in both the parser and the 8086 code generator. The “IF” statement, “ELSE” and the “ELSE IF” statements were fairly easy to implement as the hard part was compiling an expression, the expression is a part of these statements but does not define them.

02:10:

Currently while providing negative values in “IF” statements the Lexer does not notice them, they do become a “number” token but it is positive. I need to fix this and perhaps look into making expressions default to signed numbers.

# Day 35 — 7th October 2016

03:18:

Successfully implemented pointer assignments in the 8086 code generator, you can now do something such as `*ptr = 50;` to set the memory that the variable `ptr` is pointing to.

03:29:

I have realised a few weeks ago that in the parser accessing a pointer inside an expression can be a problem as how does the compiler know if you are multiplying or not, e.g `a = b *a;` Now imagine variable `a` being a pointer this should be a syntax error as `*a` is used to access the value that variable `a` is pointing at.

My point is the parser does not know if it is doing a multiplication or accessing a pointer value, the solution is to keep a variable or pointer list in the parser to describe the variables it is accessing. Should the variable it is accessing be a pointer then we know for sure we are accessing its value.

05:52

My previous diary entry at `03:29` was incorrect. I have successfully made the system work without needing to check if a variable is a pointer, I will remove the variable map and methods created with it shortly.

Pointers in expressions currently do not work as the assembly language generated will not get the value pointed to by an address stored in a particular register. I will get onto this soon and fix it.

# Day 36 — 9th October 2016

00:12:

On the 7<sup>th</sup> of October at 05:52 I stated that pointers in expressions are not working due to the assembly language generated and that it will not get the value pointed to by an address stored in a particular register.

I was wrong. The reason it is not working is because a variable is overwritten on the stack when calling a function. I obviously have made a bad design choice and I will hopefully have it fixed as soon as possible.

00:44:

I have identified many problems in the 8086 code generator, upon calling functions it is possible that the current scope variables will be overwritten, this is because I do not increment the stack pointer to accommodate for the scope of the current body.

00:54:

A bug in the 8086 code generator has been identified when accessing variables pushed to a function but it appears that there is not always a problem. The problem only appears to be valid while using function calls.

01:11:

The bug in the 8086 code generator described at 00:54 has been fixed. The problem was I was not clearing the function arguments stack upon processing them, this lead to old function arguments still being stored on the stack which lead to illegal offsets. Also the bug was not because of function calls but because of function definitions.

01:44:

The bug described at 00:44 has now been fixed, I have subtracted the stack pointer by the scope size at that current time and added it upon completion of the function call. This bug fix also fixed the bug described at 00:12 with pointers.

01:51:

Found a bug where scope variables are not cleared for each scope. This meant that after generating a function the next function generated shared the previous functions variable scope leading to illegal offsets.

01:56:

Fixed the bug described at 01:51.

02:07:

After all the bug fixing pointers in the 8086 code generator now appear to be working properly.

02:15:

Found a bug with the following “IF” statement expression:

```
“if (  
    (*a >= 97 && *a <= 122) ||  
    (*a >= 65 && *a <= 90)  
)”
```

Unfortunately, the “||” acts as “&&” I will look into this issue and fix it as soon as possible, obviously there is an issue with logical expressions.

# Day 37 — 10th October 2016

03:15:

Fixed a bug that was caused with assignments due to not checking if a compare expression was present. Before I checked this in the method that handles compare expressions but had to remove it.

# Day 38 — 13th October 2016

05:40:

Made the 8086 code generator use bytes instead of words where possible.  
For example, I will not use a word on the stack for something that should be only a byte.

21:48:

Fixed a bug in the parser where structure access did not work as expected.

# Day 39 — 14th October 2016

06:04:

Tidied up the 8086 code generator expression methods. They were quite poorly written.



# Day 40 — 17th October 2016

03:57:

Implemented the ability to get the address of a structure variable. I am a little unsure about the way I have gone about things with the structures.

04:22:

I have realised that you should not pass variable names to methods in a code generator but instead only pass the branches themselves. For example, I have a method for setting variables and it takes a string as its first argument, this string is the variable name you want to set. But now what if your setting a structure variable? E.g “a.b = 30;” you cannot just set variable “b” this would never work; you can’t set variable “a” either as variable “a” is the structure that holds the variable you want to set. If I allowed only branches to be passed to these methods, I would have never had this problem. I plan to fix this as soon as possible.

# Day 41 — 20th October 2016

02:36:

Yesterday and early this morning I was refactoring the 8086 code generator due to the fact variable strings were not decent enough as variables can also be accessed by structures and as an array.

The refactoring process took hours and I spent a good few hours debugging the code generator due to me forgetting to turn get a “STRUCT” branch based on the “STRUCT\_DEF” branch. Instead I just returned the “STRUCT\_DEF” branch and casted it to a “STRUCT” branch. I did not realise my mistake until intense debugging. The bug is fixed now and I can now access structure variables but currently not structures within structures.

I am pretty worried as things are getting more complicated and I am running out of time. I will have to implement pointer functionality for structures as well as array functionality. This will come in time. First I must fix this bug and get the system up to the standard I am happy with before continuing with more complicated tasks.

03:10:

I have identified the bug with accessing structures within structures the bug is located in the “getScopeVariable” method of the 8086 code generator. I plan to fix this as soon as I can.

03:47:

I have successfully fixed the bug relating to the “getScopeVariable” method explained at “03:10” today. Access to structure variables in expressions and assignments are now possible due to refactoring today.

05:12:

I have successfully implemented “FOR” loops in the 8086 code generator.

# Day 42 — 25th October 2016

04:41:

While attempting to implement the ability to parse arrays in the parser I realised that I made a really bad mistake by making pointers change the “ASSIGN” branch to “PTR\_ASSIGN” likewise I do the same for the “V\_DEF” branches they become “VDEF\_PTR”. The “PTR” represents a pointer.

The problem with this design is code reuse, it means I have to write the same code over and over in the parser I plan to tackle this issue shortly.

05:40:

In the parser I have now fixed the issues described at 04:41 today. However due to these fixes I have broken the ability for the parser to recognise if something is multiplication or pointer access. I shall fix this when I can.

## Day 43 — 2nd November 2016

02:20:

Fixed the bug described on Day 42 at 05:40. I just had to move a few things around and remove some code that is no longer needed that was written on Day 42, 25<sup>th</sup> of October 2016.

05:15:

I am currently attempting to implement compatibility with signed numbers in compare expressions. So far my efforts are starting to pay off, the intel 8086 instruction set has a rich variety of instructions for dealing with signed numbers which I am using to aid me.

# Day 44 — 4th November 2016

01:59:

Signed number support for compare expressions is now fully implemented.

03:59:

Upon attempting to implement arrays I realised that I made a grave mistake in design, it led me to two choices

1. Create an identifier branch where I can have methods to access its children
2. Create an “ACCESS” branch where the identifier will be as well as any other requirements.

I believe option 2. To be the correct way to go about this issue.

# Day 45 — 5th November 2016

04:16:

I have successfully improved the design of the compiler to the best of my ability. The problem was described yesterday on the 4<sup>th</sup> of November at 03:59. I chose option 2 but instead of naming the new branch “ACCESS” I named it “VAR\_IDENTIFIER”. I am still very unsure on the name of this branch but “VAR\_IDENTIFIER” will be good enough for now.

# Day 46 — 6th November 2016

16:46:

I have decided that the current design of the abstract syntax tree is still not perfect. I made an error in design where I assumed all pointer access would be to a variable. The problem is pointer access to direct memory should also be possible, because of this I will be changing the “PTR” branch to hold an expression branch rather than a variable branch. I will also may be making variable definitions use a Boolean to represent them as a pointer rather than being a child of the “PTR” branch. The “PTR” branch would still be used but only for pointer access. This sort of setup should be better with the one disadvantage and that is you would not be able to see that the variable definition is a pointer definition just from the tree alone unless I have a different branch for pointer variable definitions, although I would rather avoid that as this sort of setup may be an issue in the future.

17:44:

I have successfully completed refactoring of the “PTR” branch and also the “V\_DEF” branch and all its inheritor’s. The “PTR” branch now holds branches that will allow a code generator to permit access to a direct memory address specified by the pointer access. For example, “\*(300)” that would access memory address 300. Or more commonly you could do “\*b” which would get the value that the pointer called “b” is pointing to. Both situations are now supported where previously only pointer variable access was allowed. However, changes are currently only applied in the parser and the “8086 Code generator” must now be refactored to suit the changes. Pointers now no longer require variables and do not link to them in any way, pointers now hold expressions. The “V\_DEF” branch and all its inheritor’s are now also guaranteed to hold a “VAR\_IDENTIFIER” branch which also cleans up some previous design problems where I would need to check if it was a “PTR” branch and act accordingly if it is. If the variable defined is a pointer e.g “uint8\* ptr;” then a Boolean will be set. This is important for people to know as in no way does a branch signify that a “V\_DEF” branch is a pointer. It seemed the best solution to just mark it as a pointer and will stay that way unless a better solution arises.

# Day 47 — 7th November 2016

02:05:

After my PC crashing resulting in my compiler suite not working I wasted hours of time getting everything working again. However, once I finally got the C++ compiler to work again I continued rewriting the parser, I created something call strict mode which will prevent assignments in expressions, this is needed in certain parts of the program for a parse to be correct. I don't know whether or not to consider this a hacky approach but it doesn't look too bad and is straight forward to do.

All assignments now rely on the "process\_expression" method, this is a better setup as it now means assignments can happen in expressions which makes assignments such as "x = b = a = 5" possible.

All changes are for the parser only they will be implemented into the 8086 code generator shortly. A massive rewrite will likely have to take place since the 8086 code generator relies entirely on the old setup.



# Day 48 — 9th November 2016

21:33:

I have done a lot of changes today such as partially implemented pointers again as well as array access, I won't go into much details on them instead I want to talk about a problem I have just noticed.

The issue is that when you have a variable whose data type corresponds to one byte and then after it a variable whose data type corresponds to a word you end up with a problem. That problem is whenever I change the value of the data type whose size is a word it will overwrite the first variable whose data type is a byte in size.

This should only affect scope variables and its reason for it is that in my system scope variables are accessed offsetting from the base pointer towards the beginning of memory, the access grows downwards.

If you imagine a variable called "a" and a variable called "b" now "a" is a byte in size and "b" is a word in size. Let's pretend the "b" variable is directly under the "a" variable. Now let's pretend you set "a" to 10.

Below is our memory after setting "a" to 10.

9	8	7	6	5	4	3	2	1	0
a									

Now because when writing words to memory it will assign upwards not downwards the resulting memory becomes.

9	8	7	6	5	4	3	2	1	0
b	b								

As you can see this completely overwrites our "a" variables value.

The solution I can think of is that if it's a word to offset it one more into the negative zone avoiding this problem. I have tested it and don't believe this would cause future problems. It's the best I have for now so I will just go with it.

22:59:

I have realised a problem with accessing arrays, the problem is that accessing array indexes will grow downwards from the base pointer. This means that when pointing to an array and doing this  $(a + 5)$ . Instead of growing downwards it will go forwards causing it to access variables that are not part of the array. The same issue applies when referencing the array its self.

The solution is when accessing arrays to return the bottom of the array rather than the beginning for example “uint8 a[20];” we should return the end of the array at “a[20]” so when addition is done it will grow upwards towards the end of memory. This will fix these issues. This solution is for scope variables, argument variables should be the complete opposite.

23:26:

Both of the problems described today have been fixed.

# Day 49 — 16th November 2016

All code generator entries are related to the 8086 code generator for this diary entry.

18:41:

Many bug fixes have been done in the past few days, I have failed to record what has been going on, I plan to make further use of this diary when I can.

19:59:

Implemented pointers a little more, I forgot to implement them completely after my refactoring. There may also be a bug with multidimensional arrays but I can't be certain until more research is conducted although I am quite sure.

20:07:

Found a bug with the FOR loops, they fail to break out of the loop when the condition is no longer met. Another bug with the FOR loops is when you do not provide any statement for the "INIT" part. E.g "for(; a < 3; a = a + 1)".

20:18:

The first bug with the FOR loop has been fixed, it now breaks out of loops correctly. The problem was with the scope size calculating functionality. The bug was located in the "handle\_body" method.

20:39:

I am planning to soon fix the structure access in the compiler, I have currently just fixed it in the parser but now this must be applied to the code generator.

# Day 50 — 18th November 2016

21:40:

Fixed a bug in the 8086 code generator where array variables of a structure were not being included in the size of a variable.

22:42:

Implemented the ability for assembly instructions to have segments, this is important to separate data from code.

# Day 51 — 19th November 2016

02:31:

Successfully implemented global variables in the 8086 code generator, but global assignments are not yet supported.

03:37:

Upon implementing strings in the 8086 code generator, I realised that the way I access arrays is incorrect. I use the BX register for scope variables, by default the BX register will use the DS(Data Segment) register. This is a problem as I wanted to point to the stack for scope arrays. I will need to rewrite this when I get the chance.

# Day 52 — 22nd November 2016

18:46:

Yesterday I was working on fixing the problem with arrays in the 8086 code generator stated on the 19<sup>th</sup> of November.

I created a segmentation system allowing assembly to be written to other segments either yesterday or the day before. I modified the code for creation of arrays as well, my plan was to use the extra segment and just set it to the stack segment or the data segment, depending on what I need to access whether it be the stack segment or the data segment.

After writing code and it hit the end of my coding session, about 4 AM this morning. I was using the 16 bit “bcc” compiler to generate 8086 code and compare it to how I do it myself. I noticed they used the “BX” register and was also curious to why they offset so far off the “BP(base pointer)”. I then asked my questions on stack overflow,

<http://stackoverflow.com/questions/40733871/compilers-and-assembly> .

The questions I asked was “Why is it that compilers never start offsetting from index zero of the base pointer? Also how is the code generated above even safe to use as the BX register defaults to the data segment, which shouldn't be pointing at the stack anyway.”

The replies I got was quite interesting, the reason the “bcc” 16-bit compiler can use the “BX” register to access the stack because it follows the “tiny” model which means all the segment registers are equal to the same location, since the DS(Data Segment) and the SS(Stack Segment) would point to the same location in memory there would be no problem using the “BX” register to access stack variables. I was quite happy with this response.

As for the BP(Base Pointer) offsetting quite high they pointed out the extra pushing and popping that is being done which I completely didn't realise at the time, it makes sense why they would do that now.

This research has told me I should follow the “tiny” model, therefore arrays will now use the “BX” register again and will not use the ES(Extra Segment).

# Day 53 — 27<sup>th</sup> November 2016

01:22:

I have not written a diary entry for a few days now, it's a bad habit of mine. Any how a lot has happened recently with the 8086 code generator I have re-created structures and have created arrays, I have also programmed the ability to use them together for example if we have a structure that has two variables "a and b" and we have an array of that structure, you could access an element on that array like this. "array[5].a = 30". Currently you cannot set a structure to a structure I will add this soon.

I have also implemented the ability for pointers to be able to point to structures and also structure arrays and arrays in general.

I discovered a slight bug in the parser which I did not think about before, I don't even know whether to consider it a bug or not bug rather an incompatible design, currently pointers could be accessed like this "`*ptr+5`". This would return the value of the pointer which would be the address that it is pointing to, and then add five to it. Now what that line of code should really do is get the value from memory that the pointer is pointing to and then add five to it.

I will fix this soon and force people to do "`*(ptr+5)`" should they wish to access the memory address the pointer is pointing to + 5.

01:50:

Fixed the parser pointer bug explained at 01:22 today.

# Day 54 — 28<sup>th</sup> November 2016

01:24:

I have successfully created structure to structure assignments in the 8086 code generator, however pointer structure to pointer structure assignments are currently not supported, I may have to reconsider the design a little bit and make some adjustments as I feel the design may be slipping.



# Day 55 — 29<sup>th</sup> November 2016

01:26:

Fixed some bugs with pointers in the 8086 code generator.

# Day 56 — 1st December 2016

12:27:

Implemented inline assembly in the 8086 code generator, this allows more power to be given to the programmer, currently the inline assembly of the code generator only supports single variables so you can't pass an array element to it or a structure variable, however more support will be added in the future.

Example code:

```
uint8 print_c(uint8 c)
{
    __asm("mov ah, 0eh");
    __asm("mov al, [%i]", c);
    __asm("int 10h");
    return 0;
}

uint16 main()
{
    print_c(97);
    return 0;
}
```

12:35:

The 8086 code generator has generated some code that printed its first "Hello World" currently the second argument "size" of the "print" function causes issues with the "for" loop so I had to ignore it. The result can be seen in Figure 1. On the next page.

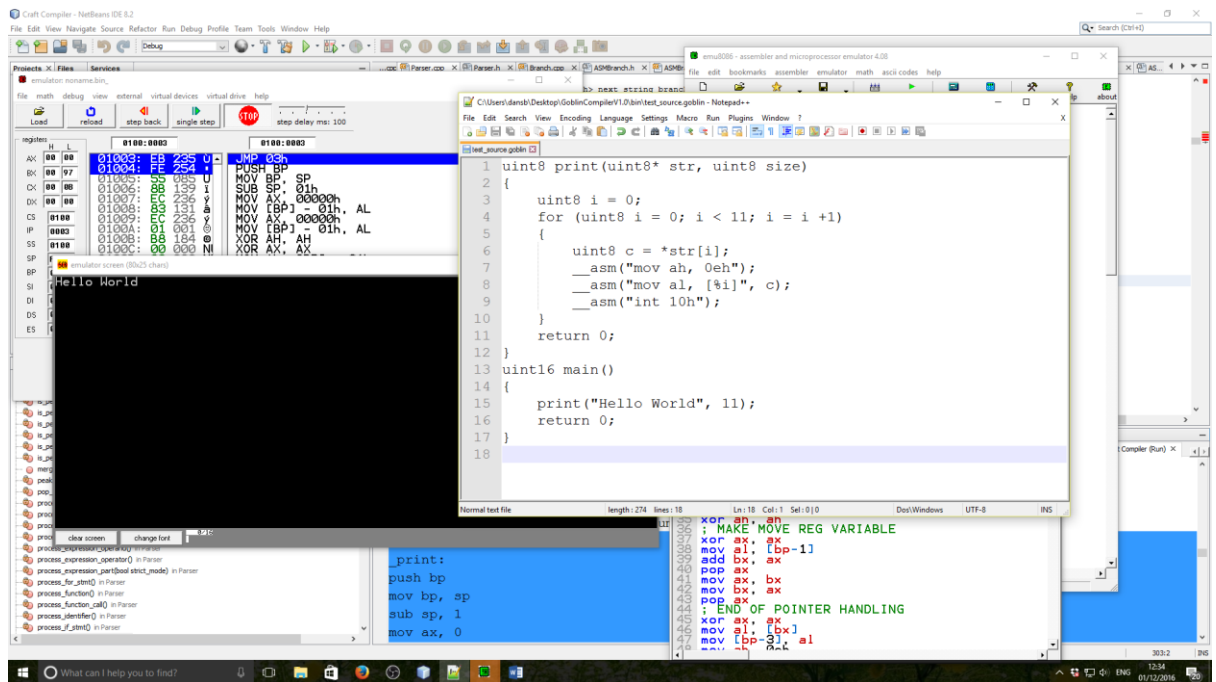


Figure 1- 8086 Code generator first "Hello World"

18:39:

In the compiler I created the ability to get the size of scope variables from the “BODYBranch” and the “FORBranch”,

# Day 57 — 3rd December 2016

11:15:

I have spent the last two days trying to implement the ability for the framework to give you the position of a variable relative to zero. I have just been successful and have even given the ability for the caller to adjust the alignment in real time. They do this by passing a lambda expression that will be invoked every time a “VDEFBranch” is found or an inheritor of the “VDEFBranch” is found.

13:44:

Removed the ability to adjust the alignment in real time for getting the position of a variable with the framework. Also added the ability to have the position start at the file size, this is useful for scopes that grow downwards.

14:52:

Fixed a bug with the “VDEFBranch” with the method for getting the variable position.

# Day 58 — 8th December 2016

12:15:

For the past two days I have been implementing a linker for the OMF(Object Module Format). I implemented a good few records correctly but then unfortunately when I got to the “PUBDEF” record I noticed that assemblers do not appear to be following the specification. I have researched this for many hours and have posted the question online:

<http://stackoverflow.com/questions/41022380/omfobject-module-format-length-field-appears-incorrect>

I have decided to put the linker on hold and may abandon the OMF linker entirely as I am really running out of time, the repository for the OMF linker I have been working on can be found here:

<https://github.com/nibblebits/MagicOMF/>

20:24:

In the framework I implemented the ability to get a variables VDEFBranch. This is really good as it takes a lot of work off the code generator, however upon doing this I realised another problem. Upon calling the “getPositionRelZero” method in the VDEFBranch it is impossible to work out the position of a structure variable, this is because only one instance of a structure is created, for this to be possible the structure contents would have to be in the STRUCTDEF branch its self, this way I could accurately work through the tree. However doing this may be inefficient I will try to think of the best solution to get around this problem.

One possible way is to make a copy of the VDEFBranch in question when accessing a structure. You would then be able to add extra information to help guide the compiler, however this may cause issues elsewhere such as if pointer values are being compared.

20:54:

I have thought of a solution for the problem described at 20:24, the solution is that the STRUCT branch’s body will be copied in memory for every STRUCT\_DEF branch that is defined. A pointer to this cloned STRUCT branch’s body will be placed in the STRUCT\_DEF branch. The new STRUCT branch’s body will never be placed on the tree. Only one STRUCT branch’s body will be displayed and that will be in the STRUCT branch was originally defined. However, in memory there can be an infinite number of clones of the STRUCT branch’s body.

# Day 59 — 9th December 2016

00:32:

I implemented the ability to clone branches ready for the task that is described at 20:54 yesterday.

# Day 60 — 10th December 2016

20:09:

Today I successfully finished upgrading the framework, any structure variable will now have their local scope point to a unique scope for the structure definition that defined it. You can also now easily get the position in memory of a variable in a structure relative to zero. Bare in mind this is not currently possible with “FOR” loops as this functionality is not yet implemented for them, I will do this soon.

# Day 61 — 13th December 2016

16:01:

For the past three days I have forgotten to write a diary entry, the system can now get the position of any variable in memory as well as the ability to get any “VDEFBranch” in the system. These framework changes have now allowed accessing structures and arrays to be a lot simpler, the code generator now generates a lot less code for accessing these things.



## Day 62 — 15th December 2016

17:53:

I have made the tree improver improve expression branches so expressions such as  $5 + 10 + 5$  could become only one branch of  $20$  and expressions such as  $5 + a + 5$  would become an expression containing a number of  $10$  and a variable identifier of  $a$ . This will take the weight off code generators and allow them to not generate instructions that they do not need to generate, as  $5 + 10 + 5$  will always be the same value there is no need to calculate this at run time. Bear in mind I am sure that there are bugs with this but do not have time to test every scenario. Perhaps with nested expressions bugs may be present.

23:05:

I have made global variables working in the 8086 code generator again, these include structure variables. There shouldn't be much more to go on the code generator now.

# Day 63 — 18th December 2016

23:15:

There was a bug with compare statements in the 8086 code generator, this bug did not allow for compare expressions such as “if(a)” and “if(a && b)”. The bug was resolved. Nested compare expressions have not yet been tested.

# Day 64 — 19th December 2016

00:11:

Made the compiler allow function definitions, also made the function definitions work inside the 8086 code generator. I also made the code generator mark all functions as global.