

# SYSTEMS PROGRAMMING

## *Hax ASM - a SIC/XE x86\_64 assembler*

```
Terminal - kandie@cornholio: ~/Workspace/Projects/Hax/assembler/bin
File Edit View Terminal Go Help

kandie@cornholio: ~/Workspace/Projects/Hax/... x kandie@cornholio: ~/Workspace/Projects/Hax/... x kandie@cornholio: ~/Workspace/Projects/Hax/... x kandie@cornholio: ~/Workspace/Projects/Hax/... x kandie@cornholio: ~/Workspace/Projects/Hax/... x

Info: skipping non-assemblable directive 'START'
t_record[1] => 0000 'FIRST' STL (0x14) RETADR 172053
t_record[1] => 0003 'LOOP' JSUB (0x48) WRREC 4B2021
t_record[1] => 0006 LDA (0x00) LENGTH 032050
t_record[1] => 0009 COMP (0x28) 0 290000
t_record[1] => 000C JEQ (0x30) ENDFIL 332006
t_record[1] => 000F JSUB (0x48) WRREC 4B2038
t_record[1] => 0012 J (0x3C) LOOP 3E2FEE
t_record[1] => 0015 'ENDFIL' LDA (0x00) =C' EOF' 032055
t_record[1] => 0018 STA (0x0C) BUFFER 0F2056
t_record[1] => 001B LDA (0x00) 3 010003
t_record[2] => 001E STA (0x0C) LENGTH 0F2048
t_record[2] => 0021 JSUB (0x48) WRREC 4B2029
t_record[2] => 0024 J (0x3C) RETADR 3E203F
Info: skipping non-assemblable directive 'USE'
Info: skipping non-assemblable directive 'RESW'
Info: skipping non-assemblable directive 'USE'
Info: skipping non-assemblable directive 'RESB'
Info: skipping non-assemblable directive 'EQU'
Info: skipping non-assemblable directive 'EQU'
Info: skipping non-assemblable directive 'USE'
t_record[3] => 0027 'RORC' CLEAR (0xB4) B410
t_record[3] => 0029 CLEAR (0xB4) B400
t_record[3] => 002B CLEAR (0xB4) B440
t_record[3] => 002D +LOT (0x74) MAXLEN 75101000
t_record[3] => 0031 'RLOOP' TD (0xE0) INPUT E32038
t_record[3] => 0034 JEQ (0x30) RLOOP 332FFA
t_record[3] => 0037 RD (0xD8) INPUT 0B2032
t_record[3] => 003A COMP (0x40) A004
t_record[3] => 003C JEQ (0x30) EXIT 332008
t_record[3] => 003F STCH (0x54) BUFFER 57A02F
t_record[3] => 0042 TLXR (0xB8) B850
t_record[4] => 0044 JLT (0x38) RLOOP 3B2FEA
t_record[4] => 0047 'EXIT' STX (0x10) LENGTH 13201F
t_record[4] => 004A RSUB (0x4C) 0 4F0000
Info: skipping non-assemblable directive 'USE'
t_record[5] => 005C 'INPUT' BYTE X'F1' F1
Info: skipping non-assemblable directive 'USE'
t_record[6] => 0040 'WRREC' CLEAR (0xB4) B410
t_record[6] => 004F LDI (0x74) LENGTH 772017
t_record[6] => 0052 'WLOOP' TD (0xE0) =X' 05' E3201B
t_record[6] => 0055 JEQ (0x30) WLOOP 332FFA
t_record[6] => 0058 LDCH (0x50) BUFFER 53A016
t_record[6] => 005B WD (0xD4) =X' 05' 0F2012
t_record[6] => 005E TLXR (0xB8) B850
t_record[6] => 0060 JLT (0x38) WLOOP 3B2FEF
t_record[6] => 0063 RSUB (0x4C) 0 4F0000
Info: skipping non-assemblable directive 'USE'
Info: skipping non-assemblable directive 'LORG'
t_record[7] => 0060 =C' EOF' (0x00) =C' EOF' 454F46
t_record[7] => 0070 =X' 05' (0x00) =X' 05' 000005
Info: skipping non-assemblable directive 'END'
dumping 7 text records
```

Author: **Ahmad Amireh**  
Supervisor: **Dr. Sami Al-Sarhan, Ph.D**

Computer Science dept., Fall 2011  
University of Jordan

## Preface

Material in this document is part of a series of assignments required by a Computer Science course, Systems Programming. The objective of said assignments is to design and program a primitive structural programming language, a compiler, an assembler, a loader, and finally a virtual machine to emulate it on the Intel x86\_64 architecture.

To fulfill this objective, I will be utilizing the knowledge gained in Discrete Mathematics, Compiler Design & Construction, and Systems Programming courses.

This part is concerned with the assembler of the above-noted pipeline, the Hax assembler.

<i>Hax</i>	<i>3</i>
<i>The Hax Language</i>	<i>3</i>
<i>The Hax Assembler</i>	<i>4</i>
<i>Overview</i>	<i>4</i>
<i>Specification</i>	<i>4</i>
<i>Technical Overview</i>	<i>4</i>
<i>Pass 1, parsing the input file</i>	<i>5</i>
<i>Pass 2, writing object program</i>	<i>5</i>
<i>Building</i>	<i>6</i>
<i>Usage</i>	<i>6</i>
<i>Limitations</i>	<i>7</i>
<i>Source Code</i>	<i>7</i>
<i>References</i>	<i>7</i>
<i>Credits</i>	<i>7</i>

# HAX

## The Hax Language

For the purposes of this document, the following assumptions are made:

1. Hax is a strongly-typed language that has built-in support for two basic datatypes: integers and characters
2. The Hax compiler produces SIC/XE compatible assembly programs
3. Target machine architecture will be emulated by a VM that represents the hypothetical SIC/XE architecture

# The Hax Assembler

## OVERVIEW

The Hax Assembler (HASM) translates SIC/XE assembly input (.hasm) file into a SIC/XE compatible object program (.obj) which can be loaded by the Hax loader.

## SPECIFICATION

HASM is a **two-pass** assembler. The assembler accepts and produce only relocatable programs, and as such, any starting addresses defined will be silently ignored. The assembler is expected to provide support for the following:

- SIC Format 1 and Format 2 instructions
- SIC/XE Format 3 and Format 4
- Literals and literal pools
- Program Blocks and Control Sections
- SIC/XE Register-to-Register instructions
- SIC/XE addressing modes:
  - Immediate addressing (in the extended Format 4 version)
  - Indirect addressing (Format 3-only)
  - Base-Relative indexed addressing
  - PC-relative addressing
- Expressions (and nested expressions)
- User-defined symbols

## TECHNICAL OVERVIEW

The assembler was written from the ground-up using C++. The decision of following an object-oriented approach to design the assembler as opposed to the classic structural one is an arbitrary one; while it does increase the complexity and has a slight impact on performance, it improves the readability and maintainability of the code.

The code makes heavy use of the STL, exceptions, and of the latest standards, C++0x.

## PASS 1, PARSING THE INPUT FILE

During the first pass, HASM scans the input file for entries, and begins preparing the objects required for processing by Pass 2. The assembler expects the input to begin with a definition of a control section, either in the form of a START operation, or a CSECT one. Once encountered, a control section with the given name (or a default-assigned one) is constructed with a default program block.

From that point on, every entry parsed is pre-processed to determine its type; from a top-view perspective, instructions could be either assembler directives, which are usually not assembled, or one of the supported formats. The type of the instruction is deduced from the opcode specified in the entry; as per the SIC/XE specification of operations, each one supports a particular subset of formats.

When an instruction is constructed, the parser proceeds to scanning for operands. Operands are created in a similar fashion to instructions; the type of operand is deduced from the content of the operand field in the entry. An operand could be a constant, a literal, an expression, or a symbol.

Once an instruction is built, attached to the current program block, its operand scanned and assigned, it is given the chance to “preprocess” - a routine in which they can usually calculate their length. Right after that, it is assigned an address based on the location counter of that block and the counter is stepped forward.

Reaching the end of the input file, the parser hands the control over to each control section, requesting them to *assemble* in turn, and the second pass is then engaged.

## PASS 2, WRITING OBJECT PROGRAM

The current control section begins its assembly routine by dumping the literals declared in the literal pool into the back of its instruction list. Afterwards, each instruction in the set is asked to assemble itself by calculating its object code.

An instruction's object code depends mainly on its operand's value, which could vary greatly between the types of operands. However, thanks to the power of abstraction, all the instruction knows is that it has an operand which has a value, so it asks the operand to evaluate and then stores the result. That value is then modified based on the type of the instruction; for example, an instruction of format 3 with the simple addressing mode might have to calculate whether it should use PC or base-relative mode.

Instructions are responsible for determining whether they require relocation, based on their operands. If they do, they track every symbol that requires relocation.

With the instructions assembled, their object code produced and ready, the control section calls the *serializer* to perform the actual object program writing. The serializer constructs

H,E,R,D,T, and M records based on the instructions contained in every program block of the current control section and writes them to the output file.

## **BUILDING**

The assembler was developed under a Linux environment using GCC 4.6.2 as a compiler. Since the code uses the C++0x standard, a compatible compiler is required to build the code. Any version of GCC later than 4.5 is compatible, however, as of the time of writing neither MSVC 2011 nor Apple's LLVM 4.2 yet fully support the latest standards and thus might not be able to compile the software.

CMake 2.8 is used to build the source, and the source is maintained in a git SCM repository. It is recommended to check out the code from the online git repository instead of the version found on the media provided with this document as it will be updated. To clone the repository, install git and run the following in a terminal:

```
git clone git@amireh.net:~/hax/hasm.git
```

Then from the same terminal session, we have to navigate to the source directory, run CMake, and build the project:

```
cd /path/to/hasm/  
mkdir build && cd build/  
cmake ..  
make -j3
```

If you would like to configure the build parameters, run `ccmake` instead of `cmake` in line 3 above. Hopefully, the above will be sufficient to build the code and you should receive no warnings or errors.

The Hax assembler was written and tested under Arch Linux on an Intel x86\_64 architecture.

## **USAGE**

A few test programs are provided with the source code, found under `hasm/text/fixture/`. These fixtures can be used for testing. The general usage format is:

```
./hasm input_file
```

This will process a file named `input_file` and write the resulting object program to `a.obj` if it succeeded. To test using a provided program, run in a terminal:

```
cd /path/to/hasm/bin  
./hasm -d -o out_program.obj ../test/fixture/program_87.asm
```

The last command will run the assembler with program\_87.asm as an input file, and it will write to out\_program.obj using the delimited output format. To view the arguments the assembler provides, run:

```
./hasm --help
```

## LIMITATIONS

1. There is no support as of yet for the ORG assembler directive
2. Expressions are not evaluated for being relative or absolute based on their terms, and thus the assembler will not recognize any such erroneous expressions and will instead ignorantly process them
3. Since HASM is a two-pass assembler, EQU operands need to consist of previously defined symbols. The same rule applies to expressions; all symbolic terms must be evaluated by the time pass 2 is engaged

## SOURCE CODE

You can view the source code online on Github using the following link:

<https://github.com/amireh/hasm>

To clone the repository, please refer to the Building section above. The code can also be found on the CD provided with the soft-copy of the document.

## References

1. System Software, An Introduction to Systems Programming 3rd Ed. (L. L. Beck)
2. SIC/XE Addressing Modes, link: [http://cis.csuohio.edu/~jackie/cis335/sicxe\\_address.txt](http://cis.csuohio.edu/~jackie/cis335/sicxe_address.txt)
3. SIC Simulator, link:  
<http://www.unf.edu/~cwinton/html/cop3601/supplements/sicsim.doc.html>

## Credits

Dr. Sami Al-Sarhan for being such an inspiration for his students, and his time spent in carrying us through core computer science theories.