

ISTANBUL TECHNICAL UNIVERSITY
ELECTRICAL-ELECTRONICS FACULTY

**SUPPORTING CUSTOM INSTRUCTIONS WITH THE LLVM COMPILER FOR
RISC-V PROCESSOR**

SENIOR DESIGN PROJECT
INTERIM REPORT

Mehmet Eymen ÜNAY
Bora İNAN
Emrecañ YİĞİT

**ELECTRONICS AND COMMUNICATION ENGINEERING
DEPARTMENT**

January 2023

ISTANBUL TECHNICAL UNIVERSITY
ELECTRICAL-ELECTRONICS FACULTY

**SUPPORTING CUSTOM INSTRUCTIONS WITH THE LLVM COMPILER FOR
RISC-V PROCESSOR**

SENIOR DESIGN PROJECT
INTERIM REPORT

Mehmet Eymen ÜNAY
(040190218)

Bora İNAN
(040190205)

Emrecañ YiĞİT
(040190205)

**ELECTRONICS AND COMMUNICATION ENGINEERING
DEPARTMENT**

Project Advisor:Dr. Tankut AKGÜL

January 2023

We are submitting the Senior Design Project Interim Report entitled as “SUPPORTING CUSTOM INSTRUCTIONS WITH THE LLVM COMPILER FOR RISC-V PROCESSOR”. The Senior Design Project Interim Report has been prepared as to fulfill the relevant regulations of the Electronics and Communication Engineering Department of Istanbul Technical University. We hereby confirm that we have realized all stages of the Senior Design Project Interim Report by ourselves, and we have abided by the ethical rules with respect to academic and professional integrity .

**Mehmet Eymen ÜNAY
(040190218)**

**Bora İNAN
(040190205)**

**Emrecañ YİĞİT
(040190205)**

TABLE OF CONTENTS

	<u>Page</u>
FRONT COVER	i
INSIDE COVER	ii
ETHICAL CONFIRMATION.....	iii
TABLE OF CONTENTS.....	ii
LIST OF TABLES	iv
LIST OF FIGURES	v
1. INTRODUCTION	1
1.1 Project Work Plan and Possible Updates.....	1
1.2 Purpose of Project.....	1
2. BASICS OF A COMPILER	4
2.1 Front-End.....	4
2.1.1 Lexical Analysis.....	4
2.1.2 Syntax Analysis	4
2.1.3 Semantic Analysis.....	4
2.1.4 IR Code Generation	4
2.2 Back-End	4
2.2.1 Optimization	4
2.2.2 Target Code Generation	5
2.2.3 Instruction Selection	5
2.2.4 Register Allocation	5
2.2.5 Instruction Scheduling	5
3. THE LLVM COMPILER	6
3.1 Parts of the Clang Frontend.....	6
3.1.1 Clang Lex Library	6
3.1.2 Clang Parse Library	6
3.1.3 Clang Sema Library	7
3.1.4 Clang CodeGen Library.....	7
3.2 LLVM IR Optimizer	7
3.3 Stages of the LLVM RISC-V Back-end	7
3.3.1 Instruction Selection	7
3.3.2 Scheduling and Formation	7
3.3.3 SSA-based Machine Code Optimizations.....	7
3.3.4 Register Allocation	8
3.3.5 Prolog/Epilog Code Insertion	8
3.3.6 Code Emission	8
3.3.7 Linking.....	8
4. PATH OF AN INSTRUCTION	9
4.1 Clang AST	9
4.2 LLVM IR	9
4.3 SelectionDAG	9
4.3.1 First Optimization Pass	11
4.3.2 Instruction Legalization	11
4.3.3 Second Optimization Pass.....	13

4.3.4 Instruction Selection	13
4.3.5 Instruction Scheduling	15
4.3.6 Machine Instruction in SSA Form	15
4.4 Machine Code Instruction	17
5. ADDING CUSTOM INSTRUCTIONS	18
5.1 TableGen Reference	18
5.2 RISC-V TableGen Classes.....	18
5.3 Adding a New Instruction Using TableGen.....	19
6. UPCOMING WORK.....	21
6.1 Instructions to be Added.....	21
6.1.1 ROT/ROTI Instruction	21
6.1.2 S-BOX Instruction	21
6.2 Hardware Modifications	21
7. RISC-V	22
7.1 RISC-V ISA.....	22
7.2 RISC-V Base Instructions	22
7.3 RISC-V Standard Extensions	22
8. NEW INSTRUCTIONS	29
8.1 SHLXOR Instruction.....	29
8.2 RORI Instruction	30
9. TESTING	36
9.1 MC Test	36
9.2 Writing the MC Test.....	37
9.3 LLC/CodeGen Test.....	37
9.4 Running Tests	39
REFERENCES.....	40
APPENDICES.....	42
APPENDIX A.1	43
1.1 Installation of Software	43

LIST OF TABLES

	<u>Page</u>
Table 1.1 : Updated Work Plan	3

LIST OF FIGURES

	<u>Page</u>
Figure 2.1 : Compiler Stages	5
Figure 3.1 : Front-end and Back-end libraries connected by LLVM	6
Figure 4.1 : AST generated by Clang	10
Figure 4.2 : AST of MLA operation.....	10
Figure 4.3 : LLVM IR file generated at the output of Clang	10
Figure 4.4 : DAG before first optimization pass.....	11
Figure 4.5 : DAG before Legalization	12
Figure 4.6 : DAG before the second optimization.....	13
Figure 4.7 : DAG before Instruction Selection.....	14
Figure 4.8 : DAG before Instruction Scheduling.....	15
Figure 4.9 : Scheduling Dependency Graph.....	16
Figure 4.10 : Machine Instruction before Register Allocation.....	16
Figure 6.1 : S-Box Layer of ASCON	21
Figure 7.1 : Level of abstraction diagram [1]	23
Figure 7.2 : RISC-V registers [2]	23
Figure 7.3 : RISC-V base instruction formats [3]	24
Figure 7.4 : RV32I base instruction set [3].....	25
Figure 7.5 : List of standard extension sets [4]	26
Figure 7.6 : Bit manipulation extension groupings [5].....	26
Figure 7.7 : Bit RV32/RV64 compatibilites and groups [6]	27
Figure 7.8 : Cryptography extension subgroups [7]	28
Figure 8.1 : C code for SHLXOR.....	30
Figure 8.2 : SHLXOR assembly output	31
Figure 8.3 : C code with builtin function for RORI	32
Figure 8.4 : roti.ll output for roti.c.....	32
Figure 8.5 : ROTI assembly output	33
Figure 8.6 : C code for RORI	34
Figure 8.7 : RORI assembly output	35
Figure 9.1 : Result of first "RUN" command	36
Figure 9.2 : Result of second "RUN" command.....	37
Figure 9.3 : Use of pipe operator.....	37
Figure 9.4 : Output for successfully passed test.....	37
Figure 9.5 : Command for using utility script	38
Figure 9.6 : LLC test file	38
Figure 9.7 : Output for passed .ll test	39

1. INTRODUCTION

1.1 Project Work Plan and Possible Updates

The table we presented in Form 3 is demonstrated in Table 1.1 with updates. We completed the mac and SHLXOR [8] instructions earlier than expected and we consulted Prof. Örs Yalçın about our progress. We presented our work and since the instructions are not based on RISC-V, she advised us to change our instructions. She mentioned the Ibex processor which is a more stable core with RISC-V architecture and was extended in a past graduation project to support the ASCON encryption algorithm. We changed our instructions with compatible ones for the extended Ibex core. Convolution, deg, and AMB are replaced with ROT/ROTI and S-BOX instructions. Prof. Örs Yalçın also mentioned that the Ibex core is not able to implement the immediate rotation instruction and asked us to add the hardware. So our next tasks are implementing the S-BOX, ROT, and ROTI instructions on the LLVM compiler and adding the capability of immediate rotation to the Ibex processor.

1.2 Purpose of Project

Application Specific Instruction Set Processors (ASIPs) are becoming more popular with the development of embedded systems. The specialization of the core causes a tradeoff between flexibility and performance. For special purposes, using ASIPs increases efficiency, however, we can program a custom ASIP only by using assembly instructions that we defined. Programming custom processors with assembly language is not a preferred way of coding. We are also not able to use high-level languages because compiling tools are designed for common architectures with certain instructions. The ability to add custom instructions to compilers will enable us to make more use of custom hardware designs.

ASIPs are feasible for all application-specific embedded systems like consumer, industrial, automotive, home appliances, cryptology, medical, telecommunication, commercial, aerospace, and military applications. The custom backend that we will design under the supervision of Dr. Tankut Akgül, is going to serve the processor designed by Prof. Dr. Sıddika Berna Örs Yalçın's research team. When the project is completed, Prof. Yalçın is going to be able to produce the assembly codes that are compatible with the processor's extended instruction set in addition to RISC-V.

There are two ways to avoid designing a specific compiler for embedded microprocessors. The first one is using a common processor that already has compiling tools. The advantage of it is reducing the costs for both hardware and software designs. However, it reduces efficiency dramatically because the processor is not designed for a specific task and hardware cost increases while the speed is decreasing. Another way is using an application specific instruction set processor and programming with the assembly instructions that the hardware designer defined. Theoretically, all efficiency benefits can be achieved but programming with assembly instructions increases coding difficulty excessively.

Prof. Yalçın and her team are studying on designing application specific instruction set processors. The purpose of this project is to create a compiler backend for a processor that supports custom instructions on top of RISC-V instructions. This compiler is going to help to program the custom processor by using high-level languages. Existing RISC-V compilers are not able to produce efficient assembly code for ASIPs. Therefore a need arose for a compiler

backend. The main reason for choosing this project is that we wanted to meet an actual need for a critical existing problem. The project has the potential to be the bridge between hardware and software of custom hardware projects in research, enabling them to be candidates for critical applications. Our team has skills and experience in low-level programming and digital design. Our project advisor Dr. Tankut Akgül's lecture on microprocessors led us to work with him on the embedded systems. Our team member Mehmet Eymen Ünay is a double major student with computer engineering department. Bora İnan has experience in software development in the defense industry and Emrecaan Yiğit is working on low-level robotics programming. Emrecaan and Bora are taking digital system Design and application course from Prof. Yalçın to learn Verilog and get to know about hardware design. All of us have assembly, C, and FPGA programming experience. We also have worked with several open source frameworks on different Linux distros. We consider that our skills and experiences match the project we will be doing.

#	Task	Responsible Group Member (s)	WEEKS																											
			1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28
1	Researching the LLVM structure	Bora Inan Emrehan Yiğit Mehmet Eymen Unay	✓	✓																										
2	Learning Tablegen	Bora Inan Emrehan Yiğit Mehmet Eymen Unay				✓	✓																							
3	Implementing the mac instruction	Bora Inan Emrehan Yiğit Mehmet Eymen Unay					✓	✓																						
4	Implementing the SHLXOR instruction	Bora Inan Emrehan Yiğit Mehmet Eymen Unay							✓	✓																				
5	Implementing the ROT/ROTI instruction	Emrehan Yiğit									✓	✓	✓	✓	✓	✓	✓	✓	✓	✓										
6	Adding ROTI instruction to hardware	Bora Inan									✓	✓	✓	✓	✓	✓	✓	✓	✓	✓										
7	Implementing the s-box instruction	Mehmet Eymen Unay									✓	✓	✓	✓	✓	✓	✓	✓	✓	✓										
8	Thesis	Bora Inan Emrehan Yiğit Mehmet Eymen Unay																			✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Table 1.1 : Updated Work Plan

2. BASICS OF A COMPILER

A compiler is a software that converts source code written in a high-level programming language into machine code appropriate for a particular computer architecture. There are different stages of a compiler but they can be grouped into two main parts such as “Front-End” and “Back-End”. These parts of the compiler are also called the analysis and synthesis parts of the compiler. The analysis stage separates the source program into its individual components and applies a grammatical structure to them. The source code is then represented in an intermediate stage using this structure. The synthesis phase creates the final target program by using the intermediate representation. We can think of the compilation process as a series of phases, each of which takes the source program and transforms it into another representation [9]. These phases can be seen in Figure 2.1.

2.1 Front-End

2.1.1 Lexical Analysis

The compiler breaks down the source code into smaller units called lexemes, which are pieces of code that correspond to specific patterns in the code. These lexemes are then converted into tokens that can be used for syntax and semantic analyses.

2.1.2 Syntax Analysis

The compiler checks that the code follows the proper syntax for the programming language it is written in. This process is also called parsing. As part of this step, the compiler often creates abstract syntax trees in order to represent the logical structure of different parts of the code.

2.1.3 Semantic Analysis

The compiler checks that the code makes logical sense, not just that it follows the syntax of the programming language. This step goes beyond syntax analysis by ensuring that the code is correct. For example, the compiler might check that variables have been declared correctly and given the appropriate data types. This process is known as semantic analysis.

2.1.4 IR Code Generation

After the source code has been analyzed for lexemes, syntax, and semantics, the compiler creates an intermediate representation (IR) of the code. This intermediate code is going to be converted to machine code in the last two phases. These two phases are platform-dependent, meaning they are specific to a particular hardware architecture but the previous phases were not. Therefore, to create a new compiler, it is not necessary to start from scratch. Instead, it is possible to use the intermediate code from an existing compiler and build the final stages of the process for a specific platform. Because of that, we are interested in the back-end part for our project.

2.2 Back-End

2.2.1 Optimization

The intermediate code is prepared for the final code generation step. This process does not change the meaning or functionality of the code, but it can make the program run faster and

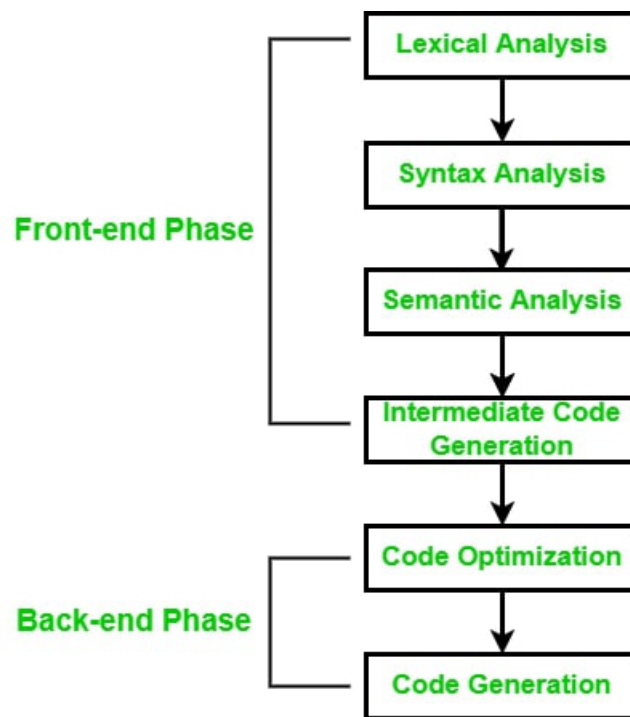


Figure 2.1 : Compiler Stages

more efficiently. A directed acyclic graph (DAG) used in the compiler design process might represent the dependencies between different instructions in the IR code, such as the order in which those instructions need to be executed or the data dependencies between them. The DAG can be used by the compiler to identify opportunities for optimization, such as removing unnecessary instructions, combining some of them or rearranging the order of execution to reduce the number of resources required by the code.

2.2.2 Target Code Generation

The target code generator is the final stage of the compilation process, and its main function is to convert the optimized code into a form that the machine can understand. The optimized code is turned into a relocatable machine code. The relocatable machine code is the input to the linker and loader, which are responsible for combining the code with other necessary resources and preparing it for execution. Target code generation can be divided into different parts:

2.2.3 Instruction Selection

IR is the input of the code generation step, and it maps the IR into the target machine's instruction set. There may be multiple ways for converting one representation, so the code generator tries to select the most suitable instructions.

2.2.4 Register Allocation

There may be many different variables/values in a program. The code generator decides which registers to use to keep these values.

2.2.5 Instruction Scheduling

The code generator determines the sequence in which instructions will be executed and creates schedules for the execution of those instructions.

3. THE LLVM COMPILER

LLVM is a collection of modular and flexible compiler and toolchain software that can be used to build a wide variety of compilers and other tools. LLVM compilers are organized into a set of libraries that implement the parts of a compiler. There are different front-end libraries for every language and different back-end libraries for every architecture. There is only one common intermediate representation optimizer that connects specific front-end and back-end.

LLVM IR is the common representation of languages. Various LLVM front-ends translate related languages into IR. Related back-end compiles IR into assembly according to the target hardware. This structure helps to increase flexibility between front-ends and back-ends. With this structure, we are able to have compilers for every combination of M source codes and N targets with M front-end and N back-end instead of $M \times N$ compilers. In our case, we do not have to struggle with the frontend also our backend will be working for every source code of LLVM because of this benefit. The front-end we use in the developing process will be clang which is the LLVM C/C++ front-end [10].

3.1 Parts of the Clang Frontend

3.1.1 Clang Lex Library

Clang Lex Library is a typical lexer implemented as finite state machines that read source code one character at a time and transition between different states based on the characters read. The Clang lexer, which is a front-end compiler for the C, C++, and Objective-C programming languages, uses this approach to filter out comments and white space, recognize and tokenize language elements such as keywords, identifiers, and operators, and handle escape sequences and string literals. The implementation files of the Clang lexer can be found in the `llvm-project/clang/lib/lex` directory within the LLVM infrastructure.

3.1.2 Clang Parse Library

Clang Parse Library is the parser that takes the tokens produced by the lexer and constructs an abstract syntax tree (AST) to represent the structure and meaning of the source code. The Clang parser checks the source code for proper syntax and resolves symbols and identifiers. It also performs type-checking to ensure the source code follows the rules of the programming

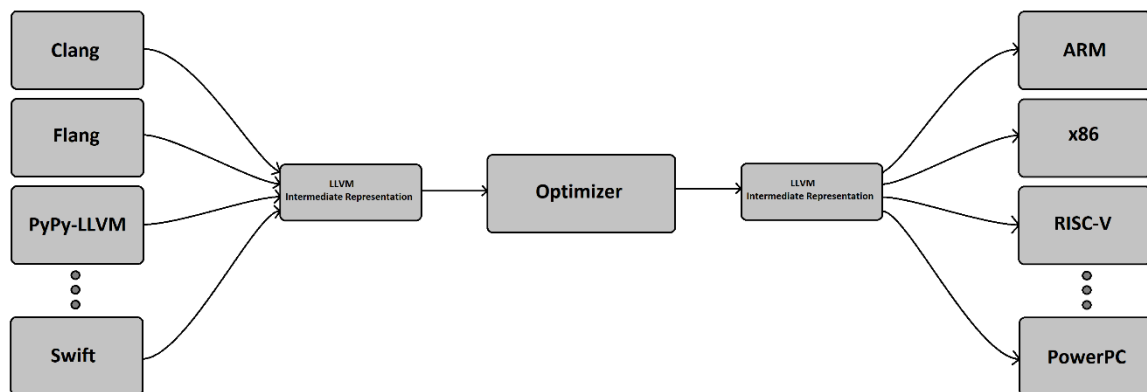


Figure 3.1 : Front-end and Back-end libraries connected by LLVM

language. It creates the AST, a tree-like structure, that represents the source code in a way that is easily processed by the compiler. The implementation files of the Clang lexer can be found in the `llvm-project/clang/lib/parse` directory within the LLVM infrastructure.

3.1.3 Clang Sema Library

Clang Sema Library is a semantic analyzer that involves examining the meaning and context of the source code in a program. In Clang, semantic analysis is a phase in the compilation process that analyzes the abstract syntax tree (AST) generated by the parser to verify that the source code conforms to the rules of the programming language and is properly constructed. Semantic analysis performs various checks and transformations on the AST to ensure the source code is correct. The implementation files of the Clang semantic analyzer can be found in the `llvm-project/clang/lib/sema` directory within the LLVM infrastructure.

3.1.4 Clang CodeGen Library

Clang CodeGen is the code generation library that takes the abstract syntax tree which is generated by the parser and corrected by the semantic analyzer as input. It generates the intermediate representation code and produces .ll file which will be used in the backend. The implementation files of the Clang lexer can be found in the `llvm-project/clang/lib/CodeGen` directory within the LLVM infrastructure.

3.2 LLVM IR Optimizer

LLVM Optimizer is a common optimization tool used for every possible source-target combination of a compiler. It takes the output file of CodeGen as input and runs the specified optimizations. The output of the optimizer becomes the input for the RISC-V back-end in our case. LLVM optimizer can be performed with `opt [options] [filename]` command [11].

3.3 Stages of the LLVM RISC-V Back-end

LLVM RISC-V Backend is responsible for compiling optimized IR down to RISC-V assembly or object code. LLVM Backend consists of libraries for the code generation steps [12].

3.3.1 Instruction Selection

SelectionDAG is the instruction selector of LLVM backends which is responsible for selecting the appropriate RISC-V instructions for a given intermediate representation instruction. It takes the target-independent LLVM code as input and generates the target-dependent DAG of instructions. SelectionDAG is the most important part of the compiler for us since we will be dealing with adding new instructions to the RISC-V Backend.

3.3.2 Scheduling and Formation

Scheduling is the phase of assigning an order to the DAG form of RISC-V instructions. The formation phase is responsible for converting the DAG into a list of machine instructions.

3.3.3 SSA-based Machine Code Optimizations

LLVM uses static single assignment (SSA) based optimizations before register allocation. SSA optimizations ensure that each variable is assigned and defined only once before it is used.

3.3.4 Register Allocation

The register allocator of RISC-V back-end is responsible for assigning physical registers to virtual registers in the IR. Each target has specific register count and order. Register allocator maps the registers by taking the RISC-V architecture registers into account. It uses the relevant `TargetRegisterInfo`, and `MachineOperand` classes. Register allocation will play a key role in our project while we manipulating the implementation of the selected instructions.

3.3.5 Prolog/Epilog Code Insertion

Prolog and epilog code insertion is another optimization phase which is responsible for frame-pointer elimination and stack packing.

3.3.6 Code Emission

Code emission stage is responsible for lowering the code generator abstractions down to the MC layer abstractions. It takes the assembly as input and creates the final RISC-V machine codes.

3.3.7 Linking

LLD is the LLVM linker library that is responsible for combining multiple object files into a single executable file. LLD is invoked after the code emission and generates a file by resolving symbol references, adjusting addresses, and performing other tasks as necessary.

4. PATH OF AN INSTRUCTION

In this chapter, the path of an instruction will be demonstrated and the corresponding DAG input of the most critical phases of SelectionDAG will be shown. We selected the input program as a function that performs multiplication and addition. This was our litmus test code used while adding MLA (Multiply and Add) instruction to the LLVM backend with TableGen. We explained our addition of MLA instruction thoroughly in Section 5.3.

```
1 int a,b,c;
2 void maddFunc() {
3     a = 3;
4     b = 103;
5
6     c = 127;
7     a = a * b +c;
8 }
```

Listing 4.1 : madd.c program

4.1 Clang AST

You can see the Abstract Syntax Tree (AST) produced by Clang in Figure 4.1. The AST consists of an expression tree with three levels. At the highest level there is an expression tree of multiplication between variables 'a' and 'b'. This expression tree's result becomes an argument for another expression tree with the addition operator. The second argument at this addition subtree is the variable 'c'. The expression tree at the root has assignment as an operator. The first argument to this tree is 'a' and the second argument is the result of multiplication and addition. Figure 4.2 shows the entire expression tree with the operations multiply and add.

4.2 LLVM IR

Clang CodeGen produces LLVM IR with the AST as the input. Figure 4.3 shows the produced LLVM IR. The optimized LLVM IR is the input to SelectionDAG to generate target-specific instructions.

4.3 SelectionDAG

Input DAGs to SelectionDAG's passes will be demonstrated so on. The following phases will be demonstrated:

1. First Optimization
2. Legalization
3. Second Optimization
4. Instruction Selection
5. Instruction Scheduling
6. Register Allocation


```

FunctionDecl 0x18ae318 <line:3:1, line:9:1> line:3:6 maddFunc 'void ()'
CompoundStmt 0x18ae600 <col:17, line:9:1>
  BinaryOperator 0x18ae3f8 <line:4:2, col:6> 'int' '='
    DeclRefExpr 0x18ae3b8 <col:2> 'int' lvalue Var 0x18ae100 'a' 'int'
    IntegerLiteral 0x18ae3d8 <col:6> 'int' 3
  BinaryOperator 0x18ae458 <line:5:2, col:6> 'int' '='
    DeclRefExpr 0x18ae418 <col:2> 'int' lvalue Var 0x18ae1c8 'b' 'int'
    IntegerLiteral 0x18ae438 <col:6> 'int' 103
  BinaryOperator 0x18ae4b8 <line:7:2, col:6> 'int' '='
    DeclRefExpr 0x18ae478 <col:2> 'int' lvalue Var 0x18ae248 'c' 'int'
    IntegerLiteral 0x18ae498 <col:6> 'int' 127
  BinaryOperator 0x18ae5e0 <line:8:2, col:13> 'int' '='
    DeclRefExpr 0x18ae4d8 <col:2> 'int' lvalue Var 0x18ae100 'a' 'int'
    BinaryOperator 0x18ae5c0 <col:6, col:13> 'int' '+'
      BinaryOperator 0x18ae568 <col:6, col:10> 'int' '*'
        ImplicitCastExpr 0x18ae538 <col:6> 'int' <LValueToRValue>
          DeclRefExpr 0x18ae4f8 <col:6> 'int' lvalue Var 0x18ae100 'a' 'int'
        ImplicitCastExpr 0x18ae550 <col:10> 'int' <LValueToRValue>
          DeclRefExpr 0x18ae518 <col:10> 'int' lvalue Var 0x18ae1c8 'b' 'int'
      ImplicitCastExpr 0x18ae5a8 <col:13> 'int' <LValueToRValue>
        DeclRefExpr 0x18ae588 <col:13> 'int' lvalue Var 0x18ae248 'c' 'int'

```

Figure 4.1 : AST generated by Clang

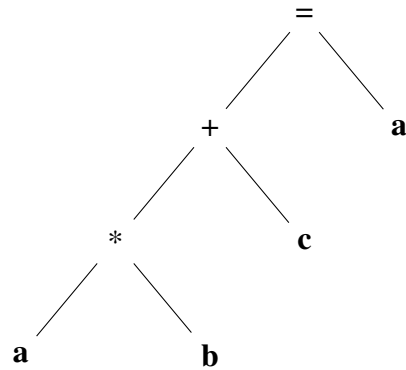


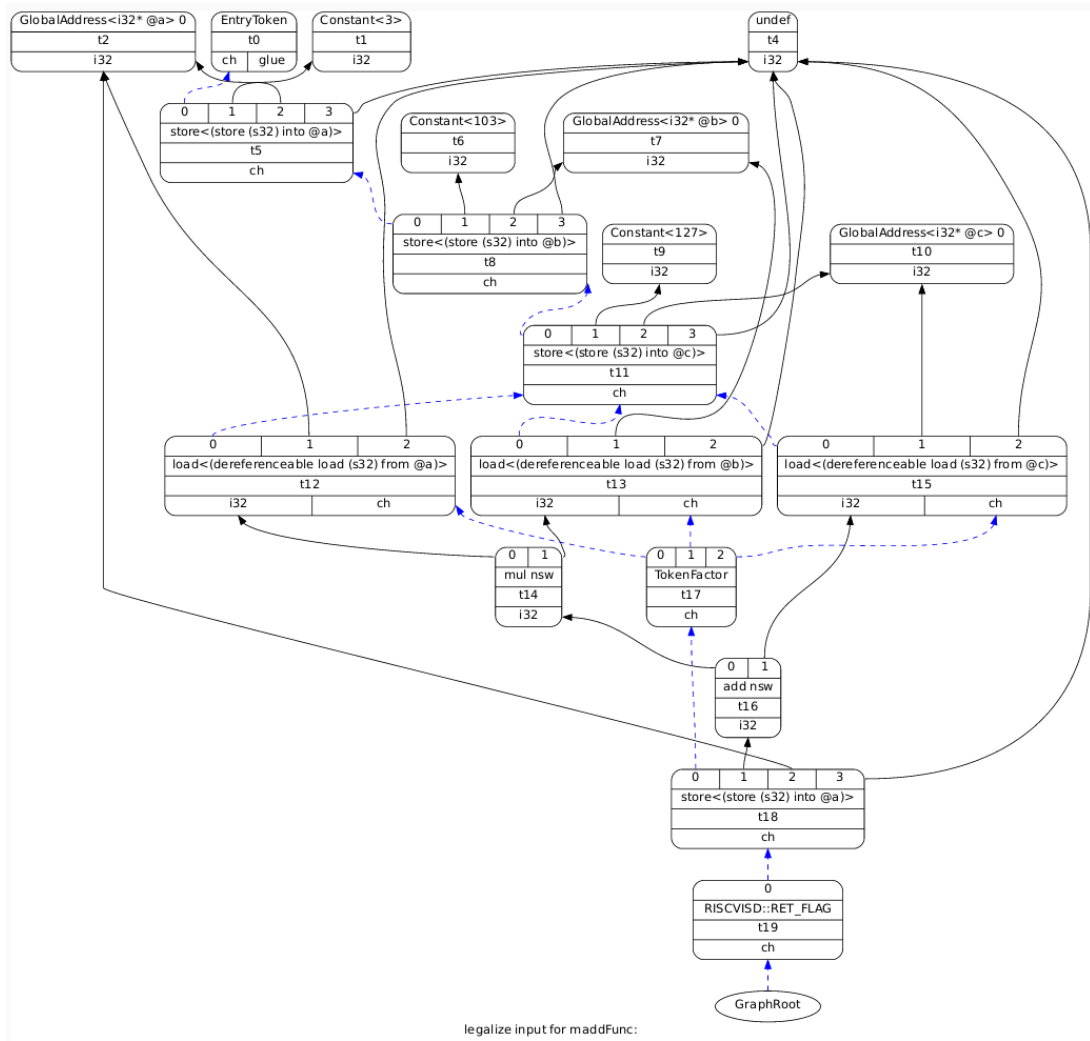
Figure 4.2 : AST of MLA operation

```

define dso_local void @maddFunc() #0 !dbg !19 {
    store i32 3, i32* @a, align 4, !dbg !23
    store i32 103, i32* @b, align 4, !dbg !24
    store i32 127, i32* @c, align 4, !dbg !25
    %1 = load i32, i32* @a, align 4, !dbg !26
    %2 = load i32, i32* @b, align 4, !dbg !27
    %3 = mul nsw i32 %1, %2, !dbg !28
    %4 = load i32, i32* @c, align 4, !dbg !29
    %5 = add nsw i32 %3, %4, !dbg !30
    store i32 %5, i32* @a, align 4, !dbg !31
    ret void, !dbg !32
}

```

Figure 4.3 : LLVM IR file generated at the output of Clang



legalize input for maddFunc:

Figure 4.5 : DAG before Legalization

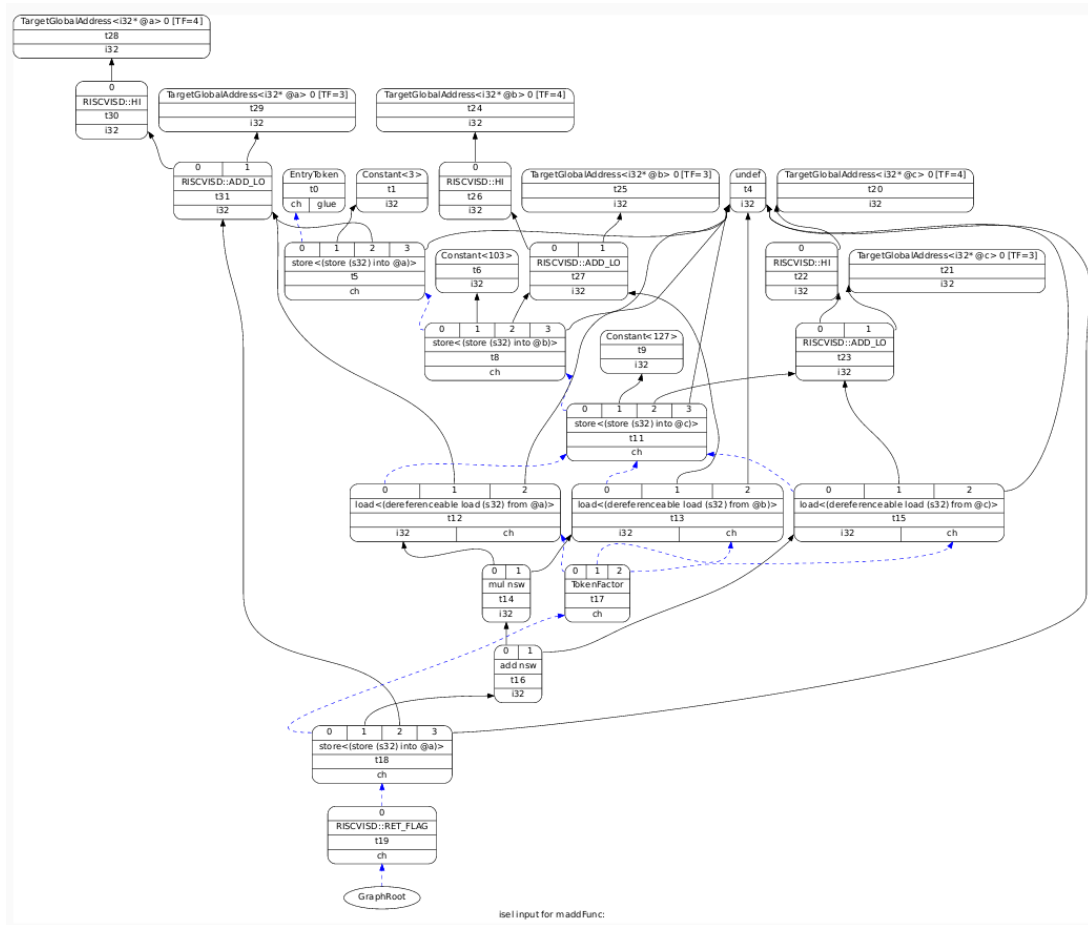


Figure 4.7 : DAG before Instruction Selection

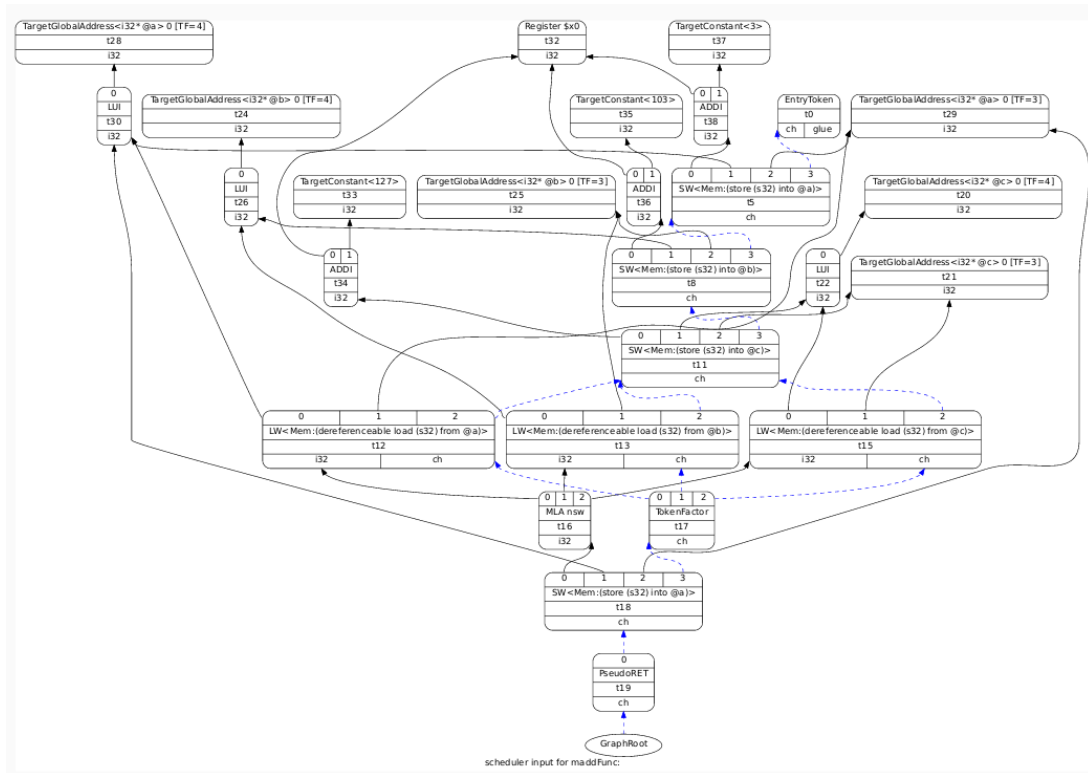


Figure 4.8 : DAG before Instruction Scheduling

We defined MLA instruction's DAG pattern as in the subgraph. The instruction selection phase took it as a reference, detected the pattern inside the global DAG and used it to place the MLA node. The addition process is explained more thoroughly in Section ??.

4.3.5 Instruction Scheduling

The DAG is transformed into a target-specific DAG with the result of legalization and selection phases. However, to generate a linear byte sequence, the DAG must be flattened. The instruction scheduling phase gets the DAG and linearises it according to the dependency graph of nodes. The scheduling dependency can be seen in Figure 4.9.

4.3.6 Machine Instruction in SSA Form

The generated Machine Instruction as a result of scheduling is shown in Figure 4.10. Because register allocation is not yet performed, the instructions are in SSA (Static Single Assignment) form. In SSA form, virtual registers are considered to be infinite unless some specific registers have to be used. In this case, '\$x0' is mentioned with ADDI instructions as they are hardwired zero in RISC-V.

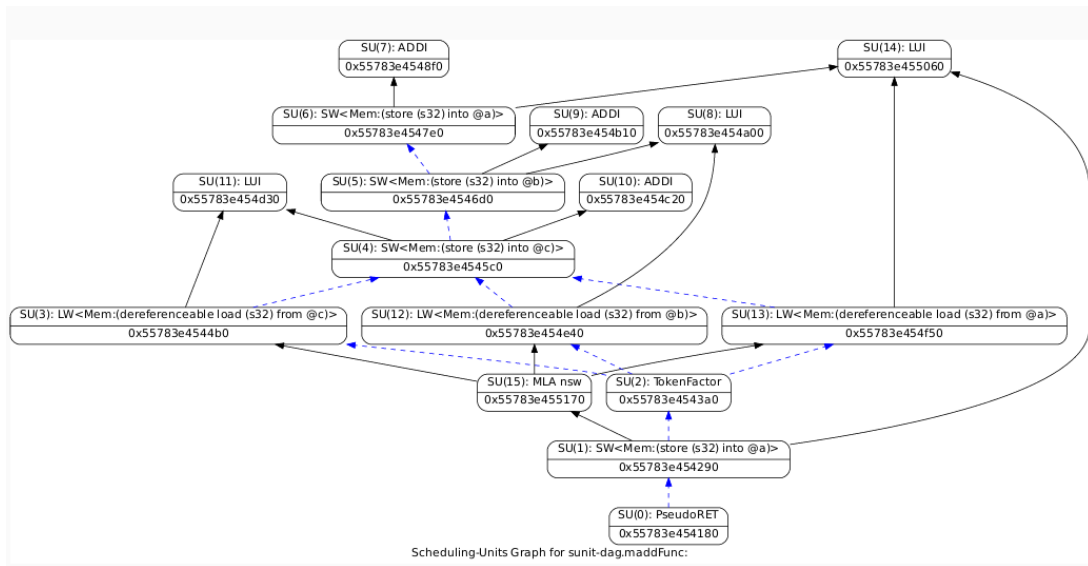


Figure 4.9 : Scheduling Dependency Graph

```
# After Instruction Selection:
# Machine code for function maddFunc: IsSSA, TracksLiveness

bb.0 (%ir-block.0):
  %0:gpr = LUI target-flags(riscv-hi) @a, debug-location !23; madd.c:4:4
  %1:gpr = ADDI $x0, 3
  SW killed %1:gpr, %0:gpr, target-flags(riscv-lo) @a, debug-location !23 :: (store (s32) into @a); madd.c:4:4
  %2:gpr = LUI target-flags(riscv-hi) @b, debug-location !24; madd.c:5:4
  %3:gpr = ADDI $x0, 103
  SW killed %3:gpr, %2:gpr, target-flags(riscv-lo) @b, debug-location !24 :: (store (s32) into @b); madd.c:5:4
  %4:gpr = LUI target-flags(riscv-hi) @c, debug-location !25; madd.c:7:4
  %5:gpr = ADDI $x0, 127
  SW killed %5:gpr, %4:gpr, target-flags(riscv-lo) @c, debug-location !25 :: (store (s32) into @c); madd.c:7:4
  %6:gpr = LW %0:gpr, target-flags(riscv-lo) @a, debug-location !26 :: (dereferenceable load (s32) from @a); madd.c:8:6
  %7:gpr = LW %2:gpr, target-flags(riscv-lo) @b, debug-location !27 :: (dereferenceable load (s32) from @b); madd.c:8:10
  %8:gpr = LW %4:gpr, target-flags(riscv-lo) @c, debug-location !29 :: (dereferenceable load (s32) from @c); madd.c:8:13
  %9:gpr = nsw MLA killed %6:gpr, killed %7:gpr, killed %8:gpr, debug-location !30; madd.c:8:12
  SW killed %9:gpr, %0:gpr, target-flags(riscv-lo) @a, debug-location !31 :: (store (s32) into @a); madd.c:8:4
  PseudoRET debug-location !32; madd.c:9:1

# End machine code for function maddFunc.
```

Figure 4.10 : Machine Instruction before Register Allocation

4.4 Machine Code Instruction

After register allocation, Machine Code Instruction (MCInst) representation of the code is created. MCInst can be thought of as an intermediate representation of the lower-level code. It can be used to produce both an object file and an Assembly file. The generated Assembly is presented below:

```
1      maddFunc:                                # @maddFunc
2 # %bb.0:
3      addi  sp, sp, -16
4 .Ltmp0:
5      sw   ra, 12(sp)                          # 4-byte Folded Spill
6      sw   s0, 8(sp)                           # 4-byte Folded Spill
7      addi  s0, sp, 16
8      lui  a0, %hi(a)
9      li   a1, 3
10     sw   a1, %lo(a)(a0)
11     lui  a1, %hi(b)
12     li   a2, 103
13     sw   a2, %lo(b)(a1)
14     lui  a2, %hi(c)
15     li   a3, 127
16     sw   a3, %lo(c)(a2)
17     lw   a3, %lo(a)(a0)
18     lw   a1, %lo(b)(a1)
19     lw   a2, %lo(c)(a2)
20     mla  a1, a3, a1, a2
21     sw   a1, %lo(a)(a0)
22     lw   ra, 12(sp)                          # 4-byte Folded Reload
23     lw   s0, 8(sp)                           # 4-byte Folded Reload
24     addi  sp, sp, 16
25     ret
```

Listing 4.2 : madd.s Assembly Output

5. ADDING CUSTOM INSTRUCTIONS

The instruction selection system we focused on at the back-end of the LLVM compiler is SelectionDAG among FastISel and GlobalISel. SelectionDAG is the most mature Instruction Selection framework with more target support. However, in the near future, it is worth considering GlobalISel as it is developed recently as an alternative to SelectionDAG. The reasons to replace it are to make it faster, smaller, and more open to low-level optimizations.

5.1 TableGen Reference

TableGen is a domain-specific language used in LLVM backend side to generate CPP header files. The purpose it serves is that it reduces redundancy of instruction declaration code which can be common to numerous architectures with minor differences. To maintain and scale the framework the minor differences are implemented level by level at a series of inheritance operations between TableGen classes.

LLVM Static Compiler, LLC, is responsible for converting LLVM IR to Assembly code. To add new instructions LLC is recompiled and its internals change. While the compilation operation of the LLC program, TableGen records are created which declare every instruction's encoding and describe its features. Referring to the records, directed acyclic graphs (DAG) are used in the process of instruction selection. DAG is a graph structure which has no cycles and has directions on the edges. Operations or functions are represented as nodes in the DAG. They are critical parts of declaring the logic or pattern of the new instruction.

The operations represented on the DAG can be LLVM intrinsics as well as instructions. LLVM instructions resemble conventional assembly instructions, in contrast, LLVM intrinsics have higher level abstraction depending on their functionality. Their instruction generation may vary depending on the target hardware. It is possible to define a new complicated instruction either by combining simple LLVM instructions and higher level intrinsics in DAG level or by creating a new LLVM intrinsic which gets created at the Intermediate Level of the compilation process.

5.2 RISC-V TableGen Classes

The most general instruction class used for every target architecture is the "InstructionEncoding" TableGen class defined in `llvm/include/llvm/Target/Target.td`. This class holds the decoder method and size of instruction in addition to minor variables. It gets inherited to the generic "Instruction" class which is defined in the same class. This class holds input and output DAGs and information which is useful to the compiler and is generalizable to all architectures.

The general class gets inherited to every target-specific class. In RISC-V's case, the next stop of the instruction is the "RVInst" class which inherits from the general "Instruction" class and it resides in `llvm/lib/Target/RISCV/RISCVInstrFormats.td` TableGen file. It defines the general bit patterns of RISC-V instructions. For example, the opcode being the first 7 bits. It defines additional information like the assembly string pattern. This general class is inherited by every type of instruction of R, I, S, B, U, and J types. As a simple example, XOR instruction can be traced. As XOR is an R type, a register-register instruction, it continues its inheritance journey from "RVInstR". It is common to R type instructions to have `funct7`, `rs2`, `rs1`, `funct3`, and `rd` format ordered from MSB to LSB. These variables are assigned corresponding bit fields in the class.

The RISC-V formats mentioned are included in the `llvm/lib/Target/RISCV/RISCVInstrInfo.td` file which is in the same directory as the `RISCVInstrFormats.td` file. After inclusion, the “RVInstR” class gets inherited by the “ALU_rr” class. The “ALU_rr” class adds the commutability feature which means swapping source 1 and source 2 does not create a different result like in addition but not in subtraction. In the end, XOR’s record is defined by putting `funct7`, `funct3` and assembly string manually in a single line with scheduling information added.

5.3 Adding a New Instruction Using TableGen

We created a tutorial to use as a reference while adding more complex instructions. Create a new tablegen file for custom additions and include it at the end of the `RISCVInstrInfo.td` file. We named it `RISCVInstrInfoCrypt.td` as it is going to be cryptography related.

```
1 include "RISCVInstrInfoCrypt.td"
```

Listing 5.1 : Include file

Add the specifications of the instruction to the `RISCVInstrInfoCrypt.td` file. Here we created a new class of instruction is defined named `ALU_rrr`. For the MLA instruction, the specifications are:

```
1 let hasSideEffects = 0, mayLoad = 0, mayStore = 0 in
2 class ALU\_rrr<bits<2> funct2, bits<3> funct3, string opcodestr,
3     bit Commutable = 0>
4     : RVInstR4<funct2, funct3, OPC_OP,
5     (outs GPR:$rd), (ins GPR:$rs1, GPR:$rs2, GPR:$rs3),
6     opcodestr, "$rd, $rs1, $rs2 , $rs3"> {
7   let isCommutable = Commutable;
8 }
```

Explanation of `ALU_rrr`:

```
1 let hasSideEffects = 0, mayLoad = 0, mayStore = 0 in
```

If the instruction has no side effect, `hasSideEffects` will be 0 If there is no need or possibility to load data from memory, `mayLoad` will be 0 If there is no need or possibility to store data from memory, `mayStore` will be 0

```
1 class ALU\_rrr<bits<2> funct2, bits<3> funct3, string opcodestr,
2     bit Commutable = 0>
```

Class is defined with `ALU_rrr` name. Variables are defined. `funct2` is two bits of binary number as `RVInstR4` is used which reserves 5 bits of `funct7` for another register. `funct3` is three bits of binary number. `opcodestr` is the string that will be shown in the assembly file. `Commutable` is one bit of zero which determines the importance of the order of the inputs.

```
1 : RVInstR4<funct2, funct3, OPC_OP, (outs GPR:$rd), (ins GPR:$rs1, GPR:$rs2, GPR:$rs3),
```

`RVInstR4` instruction type is called from `RISCVInstrFormats.td` file. `funct2`, `funct3`, `opcode`, `output`, and `inputs` are entered in function orderly.

```
1 opcodestr, "$rd, $rs1, $rs2 , $rs3"> {
2   let isCommutable = Commutable;
3 }
```

opcode string and activating commutability option.

Create another file in which we are going to add the schedules. In our case, the name of this file is “RISCVInstrInfoCrypt.td”. The definition of the instruction should be added to this file by using ALU_rrr class defined above and inputting the correct scheduling variables. For the MLA instruction, it is:

```
1 def MLA : ALU_rrr<0b10, 0b100, "mla">,
2 Sched<[WriteIMul, ReadIMul, ReadIMul]>;
```

MLA instruction is defined and ALU_rrr instruction type is used. funct2,funct3, and opcode strings are entered. Schedules are determined.

Add the instruction’s pattern defining its logic to the RISCVInstrInfoCrypt.td file. For the MLA instruction, it is:

```
1 def : Pat< (add (mul GPR:$src1, GPR:$src2), GPR:$src3),
2 (MLA GPR:$src1, GPR:$src2, GPR:$src3)>;
```

MLA instruction is called from RISCVInstrInfoCrypt.td file and inputs and outputs are determined.

6. UPCOMING WORK

6.1 Instructions to be Added

After experimenting with MLA instruction, we were given a new hardware target which lacked support in LLVM [14]. The hardware is based on the RISC-V instruction set extended by rotation (ROT) and s-box which is specific to the ASCON encryption algorithm the hardware aims to accelerate.

6.1.1 ROT/ROTI Instruction

Due to complications in LLVM at representing circular shift instead of linear shift at IR level, a new intrinsic function funnel shift (`llvm.fshl`) was introduced to the codebase. To see if it is generated at the assembly level, it is a better practice to declare the function in the `.ll` file that is fed into LLVM program. As ROT is a more popular instruction to be extended LLVM may have support already through for example the bit manipulation extension and moreover the ways to enable it will be investigated.

6.1.2 S-BOX Instruction

S-Box instruction used in the ASCON algorithm is a series of XORs and ORs between bits of 5 words which can be seen in Figure 6.1. The instruction is hard to generalize and the hardware implementation is implementing the instruction as from memory to memory which poses another set of challenges for us to overcome.

6.2 Hardware Modifications

The ROT instruction is implemented on hardware however the immediate version of it, ROTI, is not. After implementing the ROT instruction at LLVM we will try to modify the hardware without harming its integrity.

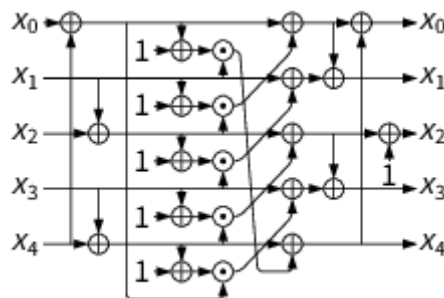


Figure 6.1 : S-Box Layer of ASCON

7. RISC-V

In this project, our target is a 32 bit RISC-V core. RISC stands for reduced instruction set computer and RISC-V is an open standard Instruction Set Architecture. [15] It is structured as a small base instruction set architecture and it has different additional extensions. The base instruction set architecture (ISA) is straightforward, rendering RISC-V appropriate for academic and learning purposes, yet extensive enough to function as a cost-effective and energy-efficient ISA for embedded systems. [16] Being open source and royalty free is another significant advantage and is an important reason why RISC-V is being commonly used. RISC-V was developed by Prof. Krste Asanović and his students Andrew Waterman and Yunsup Lee. They started working on this project in 2010 as a part of as part of the Parallel Computing Laboratory which was in UC Berkeley. Par Lab was sponsored by several companies and worked on advancing parallel computing.

7.1 RISC-V ISA

The Instruction Set Architecture (ISA) constitutes a part of a computer's abstract design that defines how the CPU is managed by the software. It serves as a bridge between the software and hardware, defining the processor's abilities and the methods by which it performs tasks. Its level in the system can be seen in figure 7.1.

There are different base integer variants of RISC-V such as RV32I, RV64I, and RV128I. These have address spaces of 32, 64, and 128 bits respectively. [17] In our project, we are interested in 32 bits. RISC-V has 32 general purpose registers. Their application binary interface names and purposes can be seen in figure 7.2. Also in the figure, we can see a different set of registers. These registers are used for floating point operations. Their ABI and purposes are also given.

7.2 RISC-V Base Instructions

There are four basic instruction formats in the base RV32I ISA. These are named R, I, S, U and all of these are 32 bits in length. There are two more additional variants named B and J as well. [3] These formats are given in figure 7.3.

Rs1 and Rs2 are the source registers and Rd is the destination register. An immediate value can also be used in some of the formats. The base instructions of the RV32I are given in figure 7.4. By inspecting their formats, we can see which type the instructions belong to. For example, the ADDI instruction is an I-type instruction and XOR is an R-type instruction.

7.3 RISC-V Standard Extensions

We had mentioned the extensions previously. Abbreviations for these extensions and what they are for are given in figure 7.5.

Thanks to these instruction extensions, more specific tasks can be implemented since we are not limited by the base instructions. Among these, the bit manipulation (B) standard extension is quite important for our project. This extension contains many instructions that operate on the bits. These extensions are also divided into several groups according to common properties. These subgroups and their purposes can be seen in figure 7.6.

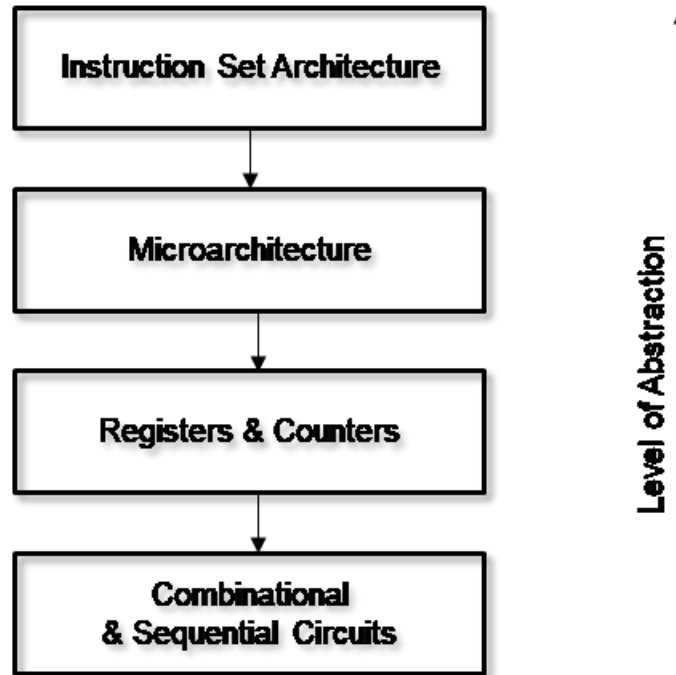


Figure 7.1 : Level of abstraction diagram [1]

Register	ABI Name	Description	Saver
x0	zero	Hard-wired zero	—
x1	ra	Return address	Caller
x2	sp	Stack pointer	Callee
x3	gp	Global pointer	—
x4	tp	Thread pointer	—
x5	t0	Temporary/alternate link register	Caller
x6–7	t1–2	Temporaries	Caller
x8	s0/fp	Saved register/frame pointer	Callee
x9	s1	Saved register	Callee
x10–11	a0–1	Function arguments/return values	Caller
x12–17	a2–7	Function arguments	Caller
x18–27	s2–11	Saved registers	Callee
x28–31	t3–6	Temporaries	Caller
f0–7	ft0–7	FP temporaries	Caller
f8–9	fs0–1	FP saved registers	Callee
f10–11	fa0–1	FP arguments/return values	Caller
f12–17	fa2–7	FP arguments	Caller
f18–27	fs2–11	FP saved registers	Callee
f28–31	ft8–11	FP temporaries	Caller

Figure 7.2 : RISC-V registers [2]

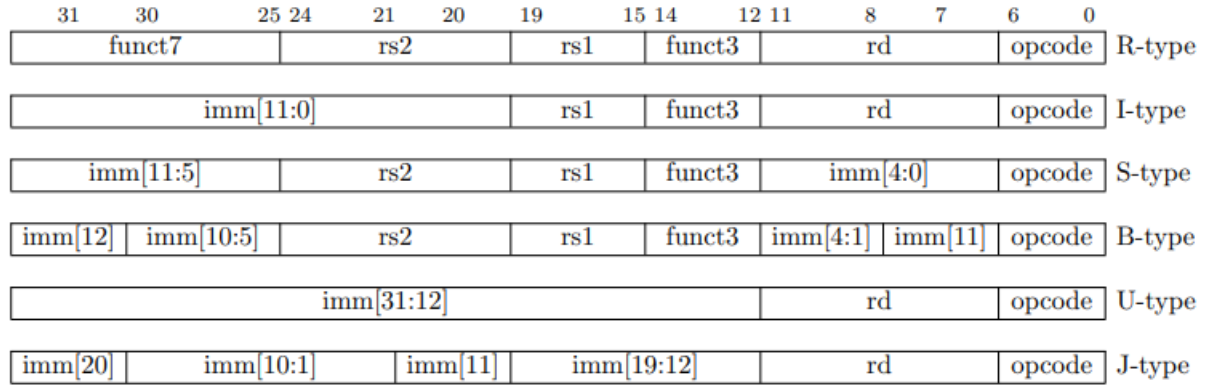


Figure 7.3 : RISC-V base instruction formats [3]

Grouping these instructions according to how commonly they are used and similarity of the operations that they perform makes them more organized and easier to work on. Some of these extensions are compatible with RV64 only. The compatibilities and the groups the instructions belong to are given in figure 7.7.

To give a clearer image of what bit manipulation (B) instructions do, let's explain few of them. For example, "CLZ" is an instruction for counting the leading zeros. Its purpose is to find out how many zeros there are before encountering a 1, starting from the most significant bit. Another example is "ORN" instruction. It negates the second operand and performs bitwise or with the first one. There is also a scalar cryptography instruction set extension for RISC-V. The RISC-V Scalar Cryptography extensions allow cryptographic tasks to be completed more quickly. Furthermore, these extensions significantly reduce the difficulty of implementing fast and secure cryptography in embedded devices and IoT. [18] This instruction set extension is also divided into subgroups according to the purpose and similarity of the instructions. The groups are given in figure 7.8.

These instructions are not a part of the base instruction set. Therefore, one cannot perform these operations with a single instruction if they use the base instructions only. Instruction extensions prove quite useful in these situations. One important thing we need to be careful about is that we should check these standard instruction extensions before we try to implement non-standard extensions by ourselves. These are comprehensive extensions and may already contain the extensions that we want to implement. After making sure that the extension we want is not present, we may try to implement non-standard extensions.

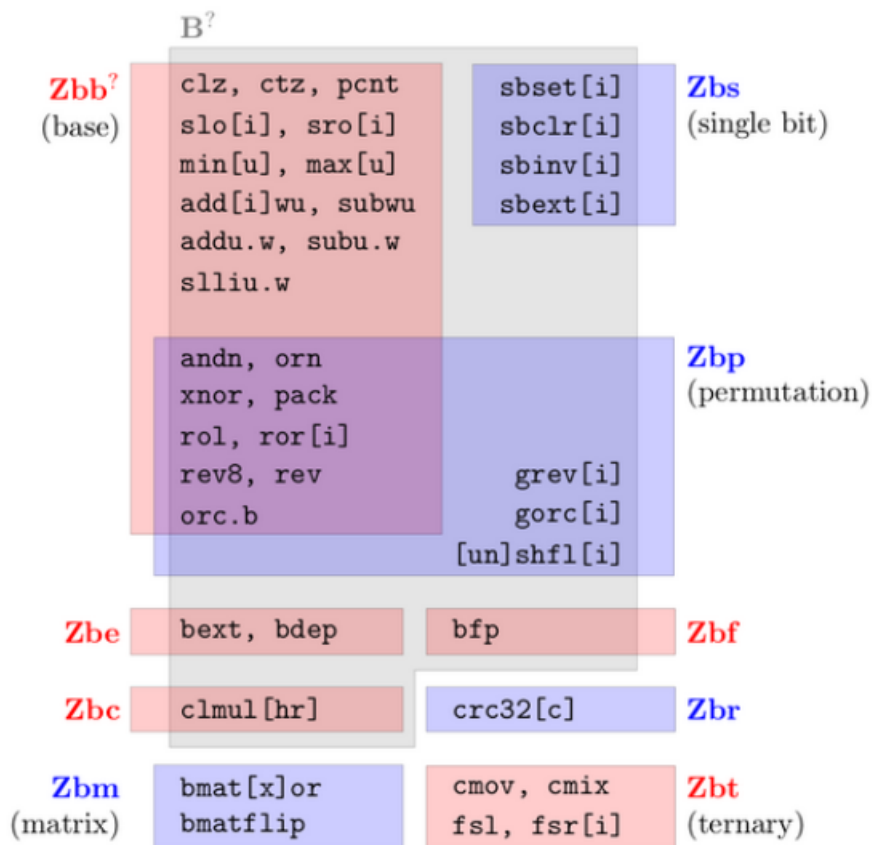
RV32I Base Instruction Set

imm[31:12]					rd	0110111	LUI
imm[31:12]					rd	0010111	AUIPC
imm[20:10:11:19:12]					rd	1101111	JAL
imm[11:0]			rs1	000	rd	1100111	JALR
imm[12:10:5]		rs2	rs1	000	imm[4:1:11]	1100011	BEQ
imm[12:10:5]		rs2	rs1	001	imm[4:1:11]	1100011	BNE
imm[12:10:5]		rs2	rs1	100	imm[4:1:11]	1100011	BLT
imm[12:10:5]		rs2	rs1	101	imm[4:1:11]	1100011	BGE
imm[12:10:5]		rs2	rs1	110	imm[4:1:11]	1100011	BLTU
imm[12:10:5]		rs2	rs1	111	imm[4:1:11]	1100011	BGEU
imm[11:0]			rs1	000	rd	0000011	LB
imm[11:0]			rs1	001	rd	0000011	LH
imm[11:0]			rs1	010	rd	0000011	LW
imm[11:0]			rs1	100	rd	0000011	LBU
imm[11:0]			rs1	101	rd	0000011	LHU
imm[11:5]		rs2	rs1	000	imm[4:0]	0100011	SB
imm[11:5]		rs2	rs1	001	imm[4:0]	0100011	SH
imm[11:5]		rs2	rs1	010	imm[4:0]	0100011	SW
imm[11:0]			rs1	000	rd	0010011	ADDI
imm[11:0]			rs1	010	rd	0010011	SLTI
imm[11:0]			rs1	011	rd	0010011	SLTIU
imm[11:0]			rs1	100	rd	0010011	XORI
imm[11:0]			rs1	110	rd	0010011	ORI
imm[11:0]			rs1	111	rd	0010011	ANDI
0000000		shamt	rs1	001	rd	0010011	SLLI
0000000		shamt	rs1	101	rd	0010011	SRLI
0100000		shamt	rs1	101	rd	0010011	SRAI
0000000		rs2	rs1	000	rd	0110011	ADD
0100000		rs2	rs1	000	rd	0110011	SUB
0000000		rs2	rs1	001	rd	0110011	SLL
0000000		rs2	rs1	010	rd	0110011	SLT
0000000		rs2	rs1	011	rd	0110011	SLTU
0000000		rs2	rs1	100	rd	0110011	XOR
0000000		rs2	rs1	101	rd	0110011	SRL
0100000		rs2	rs1	101	rd	0110011	SRA
0000000		rs2	rs1	110	rd	0110011	OR
0000000		rs2	rs1	111	rd	0110011	AND
fm	pred	succ	rs1	000	rd	0001111	FENCE
0000000000000			00000	000	00000	1110011	ECALL
0000000000001			00000	000	00000	1110011	EBREAK

Figure 7.4 : RV32I base instruction set [3]

A	Atomic instructions
M	Integer multiplication and division
F	Single precision floating point
D	Double precision floating point
Q	Quad precision floating point
C	Compressed instructions
L	Decimal floating point
B	Bit manipulation
P	Packed SIMD instructions
V	Vector operations
N	User level interrupts
J	Dynamically translated languages
T	Transactional memory

Figure 7.5 : List of standard extension sets [4]



Zbb (base): the operations of most common use.
Zbc (carry-less): carry-less multiplication.
Zbe (extract/deposit): extract/deposit a mask of multiple bits in a value.
Zbf (bit-field): placement of compact bit fields.
Zbp (permutation): large scale bit permutations (e.g.: rotations, reversals, shuffling...).
Zbm (matrix): matrix operations.
Zbr (redundancy): cyclic redundancy check operations (crc).
Zbs (single bit): single bit operations: set, clear, invert, extract.
Zbt (ternary): three operand operations.

Figure 7.6 : Bit manipulation extension groupings [5]

Extension	RV32/RV64	RV64 only
Zbb (*)	clz, ctz, cpop min, minu, max, maxu sext.b, sext.h, zext.h andn, orn, xnor rol, ror, rori rev8, orc.b	clzw, ctzw, cpow rolw, rorw, roriw
Zbp	andn, orn, xnor pack, packu, packh rol, ror, rori grev, grevi gorc, gorci shfl, shfli unshfl, unshfli xperm.n, xperm.b, xperm.h	packw, packuw rolw, rorw, roriw grevw, greviw gorcw, gorciw shflw unshflw xperm.w
Zbs	bset, bseti bclr, bclri binv, binvi bext, bexti	
Zba (*)	sh1add sh2add sh3add	sh1add.uw sh2add.uw sh3add.uw add.uw, slli.uw
Zbe	bcompress, bdecompress pack, packh	bcompressw, bdecompressw packw
Zbf	bfp pack, packh	bfpw packw
Zbc (*)	clmul, clmulh, clmulr	
Zbm		bmator, bmatxor, bmatflip unzip16, unzip8 pack, packu
Zbr	crc32.b, crc32c.b crc32.h, crc32c.h crc32.w, crc32c.w	 crc32.d, crc32c.d
Zbt	cmov, cmix fsl, fsr, fsri	fslw, fsrw, fsriw
B	All of the above except Zbr and Zbt	

Notes:

- * means the extensions are expected to be unchanged in the official version.

Figure 7.7 : Bit RV32/RV64 compatibilites and groups [6]

Scalar Cryptography Instruction Set Extension – Group Names

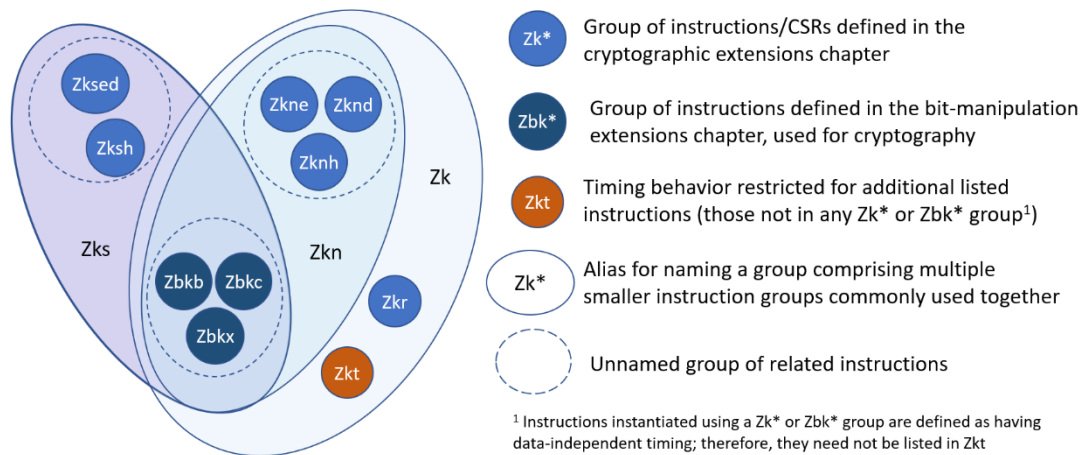


Figure 7.8 : Cryptography extension subgroups [7]

8. NEW INSTRUCTIONS

8.1 SHLXOR Instruction

We extended the instruction set by adding the SHLXOR instruction. The purpose of this instruction is to shift the first source operand one bit to the left and then XOR it with the second source operand. Then, the obtained result is stored in the destination register. By adding this instruction, we can perform this operation with a single instruction instead of using shift left and XOR instructions separately, making it more efficient. Let's give an example to make it clearer what this instruction does. RS1: 0x0101 RS2: 0xFFFF RD: 0xFDFD 0x0101 is shifted left by one and then XOR'ed with 0xFFFF, giving the result 0xFDFD. As we mentioned, this new instruction requires two source registers and one destination register. Therefore, unlike the MLA instruction, we don't need to create a new class to support it. There is already a class named ALU_rr in the RISCVINstrInfo.td file which has two source and one destination register. Therefore, the new SHLXOR instruction is going to belong to the ALU_rr class. The ALU_rr class is given below.

```
1 let hasSideEffects = 0, mayLoad = 0, mayStore = 0 in
2 class ALU_rr<bits<7> funct7, bits<3> funct3, string opcodestr,
3 bit Commutable = 0>
4 : RVInstr<funct7, funct3, OPC_OP, (outs GPR:\$rd), (ins GPR:\$rs1, GPR:\$rs2),
5 opcodestr, "\$rd, \$rs1, \$rs2"> {
6 let isCommutable = Commutable;
7 }
```

As we can see, encoding of this type of instructions consists of funct7, funct3, opcode, source registers, and the destination register. The encoding format and the other properties are described in the class. The source registers are described as inputs and the destination register is described as the output. In the RISCVINstrInfoCrypt.td file, we add the definition of the SHLXOR instruction by using the ALU_rr class. In this part, we define funct7, funct3, and the mnemonic of the new instruction as well as the schedulings.

```
1 def SHLXOR : ALU_rr<0b0011000, 0b111, "shlxor">,
2 Sched<[WriteIALU, ReadIALU, ReadIALU]>;
```

Also in the same file, we define the instruction's pattern. When we examine this definition, we can clearly see what the instruction performs and its pattern. In the inner parentheses, we can see the shifting of the first source operand by one bit. Then, the result of this shifting operation is used as an input for the XOR operation alongside the second source operand.

```
1 def : Pat< (xor (shl GPR:\$src1, (i32 1)), GPR:\$src2),
2 (SHLXOR GPR:\$src1, GPR:\$src2)>;
```

After doing these, we can try it with a simple C code given in 8.1. Let's get an assembly output from the C code by running the following commands:

```
1 clang -S -target riscv32-linux-gnu -emit-llvm -g shlxor.c
2 $~/CustomInstrLLVM/build/bin/llc -mtriple=riscv32 shlxor.ll$
```

Here, CustomInstrLLVM is the folder that we built llvm in.

```

1 int a,b;
2
3 void shlxor() {
4     a = 3;
5     b = 5;
6
7     a = b^(a<<1);
8 }

```

Figure 8.1 : C code for SHLXOR

This C code basically shifts the variable a by one bit and XOR's it with b variable. Then the result is stored in a. We can see that, the register that stores the value a is both the first source register and the destination register. The assembly output is given in figure 8.2.

In the assembly output, we can see the SHLXOR instruction. a1 and a2 are the sources and a1 is also the destination as we can see.

8.2 RORI Instruction

One of the Instructions that we worked on is RORI instruction. The purpose of this instruction is to take the operand and rotate it to the right by the amount of the immediate value. This is different from shifting right using an immediate. When shifting a number to the right, the rightmost (least significant) bits are deleted and the leftmost (most significant) bits are filled in with zeros. On the other hand, when a number is rotated right, the least significant bits that are pushed out are not deleted but written into the most significant bits. We want to add an instruction that performs this operation.

First of all, we tried pattern matching and added the definition and pattern of our new instruction to the InstrInfoCrypt.td file. However, that didn't work and we couldn't see the RORI instruction when we checked our assembly output created by using a simple C code that implements the rotation operation. This is because the pattern can have different combinations and may not match with what we expect. Therefore, the instruction cannot be recognized and we can't see it in the assembly output.

Realizing that, we looked for other options and tried to make use of intrinsics and builtin functions. We added the definitions into the InstrInfoCrypt.td file.

```

1 def ROTI : ALU_ri<0b101, "roti">;

```

It is ALU_ri type because one of the operands is an immediate value and we only need one register for this instruction.

```

1 def : Pat<(rotr GPR:\$rs1, simm12:\$imm12),
2 (ROTI GPR:\$rs1, simm12:\$imm12)>;

```

We used rotr here because we wanted to make use of the builtin function `__builtin_rotateright32`. rotr is defined in the RISVIselLowering.cpp file.

```

1 if (Subtarget.hasStdExtZbb() || Subtarget.hasStdExtZbkb()) {
2 if (Subtarget.is64Bit())
3 setOperationAction({ISD::ROTL, ISD::ROTR}, MVT::i32, Custom);

```

```

 9 shlxor:                                     # @shlxor
10 .Lfunc_begin0:
11     .loc      0 4 0                          # shlx.c:4:0
12     .cfi_sections .debug_frame
13     .cfl_startproc
14 # %bb.0:
15     addi      sp, sp, -16
16     .cfl_def_cfa_offset 16
17 .Ltmp0:
18     .loc      0 5 4 prologue_end              # shlx.c:5:4
19     sw        ra, 12(sp)                      # 4-byte Folded Spill
20     sw        s0, 8(sp)                      # 4-byte Folded Spill
21     .cfl_offset ra, -4
22     .cfl_offset s0, -8
23     addi      s0, sp, 16
24     .cfl_def_cfa s0, 0
25     lui      a0, %hi(a)
26     li       a1, 3
27     sw        a1, %lo(a)(a0)
28     .loc      0 6 4                          # shlx.c:6:4
29     lui      a1, %hi(b)
30     li       a2, 5
31     sw        a2, %lo(b)(a1)
32     .loc      0 9 6                          # shlx.c:9:6
33     lw        a1, %lo(b)(a1)
34     .loc      0 9 9 is_stmt 0                 # shlx.c:9:9
35     lw        a2, %lo(a)(a0)
36     .loc      0 9 7                          # shlx.c:9:7
37     shlxor    a1, a2, a1
38     .loc      0 9 4                          # shlx.c:9:4
39     sw        a1, %lo(a)(a0)
40     .loc      0 10 1 is_stmt 1                # shlx.c:10:1
41     lw        ra, 12(sp)                     # 4-byte Folded Reload
42     lw        s0, 8(sp)                      # 4-byte Folded Reload
43     addi      sp, sp, 16
44     ret

```

Figure 8.2 : SHLXOR assembly output

```

1 int a;
2
3 void ROT() {
4     a = 15;
5
6     a = __builtin_rotateright32(a,2);
7 }

```

Figure 8.3 : C code with builtin function for RORI

```

1 ; ModuleID = 'roti.c'
2 source_filename = "roti.c"
3 target datalayout = "e-m:e-p:32:32-i64:64-n32-S128"
4 target triple = "riscv32-unknown-linux-gnu"
5
6 @a = dso_local global i32 @0, align 4, !dbg !0
7
8 ; Function Attrs: noinline nounwind optnone
9 define dso_local void @ROT() #0 !dbg !15 {
10     store i32 15, i32* @a, align 4, !dbg !19
11     %1 = load i32, i32* @a, align 4, !dbg !20
12     %2 = call i32 @llvm.fshr.i32(i32 %1, i32 %1, i32 2), !dbg !21
13     store i32 %2, i32* @a, align 4, !dbg !22
14     ret void, !dbg !23
15 }

```

Figure 8.4 : roti.ll output for roti.c

```

4 } else {
5 setOperationAction({ISD::ROTL, ISD::ROTR}, XLenVT, Expand);
6 }

```

However, we need to change the “Expand” to “Legal” in here otherwise we will not see the ROTI instruction in the assembly output. This way, we are legalizing the action. If we don’t do this, we will see the fshr (funnel shift right) intrinsic in the .ll file but ROTI won’t make it into the assembly output. After doing these, we can try it with a simple C code. As mentioned before, we used a built in rotate function in the C code given in figure 8.3.

This code rotates 15 to the right by 2 bits. We can obtain the .ll file given in figure 8.4 by running the following command.

```
1 clang -S -target riscv32-linux-gnu -emit-llvm -g roti.c
```

The fshr intrinsic is visible in the twelfth line. After that, the assembly output can be obtained by running the following command. The assembly output is given in figure 8.5.

```
1 $~/CustomInstrLLVM/build/bin/llc -mtriple=riscv32 roti.ll$
```

```

1      .text
2      .attribute      4, 16
3      .attribute      5, "rv32i2p0"
4      .file      "roti.c"
5      .file      0 "/home/bora/Desktop" "roti.c" md5 0x6443649ef4e7311dcaf1e4ca685e9f
6      .globl      ROT                                     # -- Begin function ROT
7      .p2align      1
8      .type      ROT,@function
9 ROT:
10     .Lfunc_begin0:                                     # @ROT
11     .loc      0 3 0                                     # roti.c:3:0
12     .cfi_sections .debug_frame
13     .cfi_startproc
14 # %bb.0:
15     addi      sp, sp, -16
16     .cfi_def_cfa_offset 16
17 .Ltmp0:
18     .loc      0 4 4 prologue_end                         # roti.c:4:4
19     sw        ra, 12(sp)                                  # 4-byte Folded Spill
20     sw        s0, 8(sp)                                   # 4-byte Folded Spill
21     .cfi_offset ra, -4
22     .cfi_offset s0, -8
23     addi      s0, sp, 16
24     .cfi_def_cfa s0, 0
25     lui      a0, %hi(a)
26     li       a1, 15
27     sw        a1, %lo(a)(a0)
28     .loc      0 11 30                                     # roti.c:11:30
29     lw        a1, %lo(a)(a0)
30     .loc      0 11 6 is_stmt 0                           # roti.c:11:6
31     roti      a1, a1, 2
32     .loc      0 11 4                                     # roti.c:11:4
33     sw        a1, %lo(a)(a0)
34     .loc      0 14 1 is_stmt 1                           # roti.c:14:1
35     lw        ra, 12(sp)                                  # 4-byte Folded Reload
36     lw        s0, 8(sp)                                   # 4-byte Folded Reload
37     addi      sp, sp, 16
38     ret

```

Figure 8.5 : ROTI assembly output


```

1 #define XLEN 32
2 #include <stdint.h>
3 #define uint_xlen_t uint32_t
4
5 uint_xlen_t rotimm(uint_xlen_t rs1)
6 {
7     uint_xlen_t a = 0;
8     a = ((rs1>>2) | (rs1<<(XLEN-2)));
9     return a;
10 }

```

Figure 8.6 : C code for RORI

The ROTI instruction can be seen in the assembly output. As expected, it uses one register as both the destination and the source alongside an immediate value.

This way, we successfully obtained the ROTI instruction. However, using the builtin function while writing the C code wasn't the prettiest solution. Therefore, we looked further into how we can implement this. After that, we realized that there is a bit manipulation extension for RISC-V and what we need was in the RISC-VInstrInfoZb.td file. This file contains the instruction extensions for bit manipulations. These instructions operate on the bits of the data and RORI is one of those instructions. However, in order to utilize this extension, we need to add some flags to the command while running clang in order to get an assembly output from the C code we write. The simple C code for RORI is given in figure 8.6.

Here, we use 32 as the length because our target is 32 bit and the rotation is implemented in the eighth line. After that, we run the following command with additional flags as mentioned before.

```

1 $clang --target=riscv32 -O -S rori.c -march=rv32imaczbbs$

```

Here, -O defines the level of optimization. -S is used for getting an assembly file as an output. rori.c is the name of our simple C code. -march=rv32imaczbbs designates that we want to utilize the zbb subgroup of the bit manipulation extension. The assembly output is given in figure 8.7.

This way, we managed to successfully obtain RORI instruction in the assembly output.

```

1      .text
2      .attribute      4, 16
3      .attribute      5, "rv32i2p0_m2p0_a2p0_c2p0_zbb1p0"
4      .file   "rori.c"
5      .globl  rotimm
6      .p2align      1
7      .type   rotimm,@function
8 rotimm:
9      rori   a0, a0, 2
10     ret
11 .Lfunc_end0:
12     .size   rotimm, .Lfunc_end0-rotimm
13
14     .ident  "Ubuntu clang version 14.0.0-1ubuntu1"
15     .section      ".note.GNU-stack","",@progbits
16     .addrsig|

```

Figure 8.7 : RORI assembly output

9. TESTING

LLVM-lit coordinates the testing procedure. The comment lines that start with “RUN” call other programs via LLVM-lit. LLVM-lit also gives the output of a program to another program as an input.

FileCheck, as the name implies, controls the checking process. It basically compares the file and the corresponding lines of the output. It is used with CHECK-* command.

Regression testing is a core part of LLVM because of its size and active development. To make sure newly added features don't break the already present functionality it is a must to both build functionality and its corresponding tests.

9.1 MC Test

LLVM-MC is an abstracted assembler integrated with the Compiler. [19]

Before adding a new instruction, check its byte sequence. LLVM-MC is going to be used to get the output.

Here is a working test file:

```
1 \# RUN: llvm-mc %s -triple=riscv32 -riscv-no-aliases -show-encoding \  
2 \# RUN: | FileCheck -check-prefixes=CHECK-ASM,CHECK-ASM-AND-OBJ %s \  
3 \# RUN: llvm-mc -filetype=obj -triple=riscv32 < %s \  
4 \# RUN: | llvm-objdump -M no-aliases -d -r - \  
5 \# RUN: | FileCheck -check-prefixes=CHECK-ASM-AND-OBJ %s \  
6 \  
7 \# CHECK-ASM-AND-OBJ: shlxor s2, s2, s8 \  
8 \# CHECK-ASM: encoding: [0x33,0x59,0x89,0x81] \  
9 shlxor s2, s2, s8
```

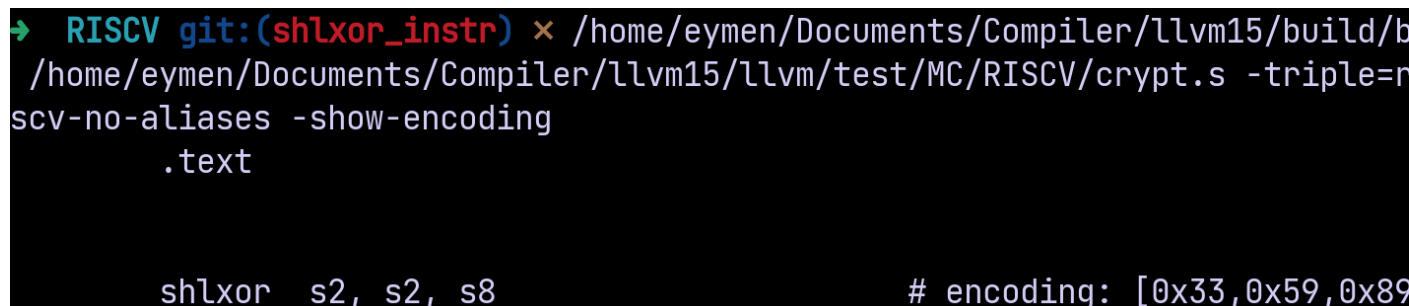
The first RUN command sequence results in the output given in figure 9.1.

Here, FileCheck is checking both the Assembly encoding and the string as it was provided with the prefixes: CHECK-ASM, CHECK-ASM-AND-OBJ.

The second RUN sequence with only the llvm-mc part spits out an ELF object which itself isn't useful. The console output can be seen in figure 9.2.

By using the pipe 'l' operator like in Shell, we can tell LLVM-lit to feed another program with the output of a program. An example is given in figure 9.3.

Since we observed the outputs of the commands, we can make sense of what FileCheck is checking. Here in the second RUN sequence, FileCheck is provided only with



```
→ RISCV git:(shlxor_instr) x /home/eymen/Documents/Compiler/llvm15/build/b  
/home/eymen/Documents/Compiler/llvm15/llvm/test/MC/RISCV/crypt.s -triple=r  
scv-no-aliases -show-encoding  
.text  
  
shlxor s2, s2, s8 # encoding: [0x33,0x59,0x89,0x81]
```

Figure 9.1 : Result of first "RUN" command

```
→ RISCVC git:(shlxor_instr) × /home/eymen/Documents/Compiler/llvm15/build/bin/llvm-mc /home/eymen/Documents/Compiler/llvm15/llvm/test/MC/RISCV/crypt.s -filetype=obj -triple=riscv32
ELF`4(3Y.text.strtab.symtabH48%
```

Figure 9.2 : Result of second "RUN" command

```
→ RISCVC git:(shlxor_instr) × /home/eymen/Documents/Compiler/llvm15/build/bin/llvm-mc -filetype=obj -triple=riscv32 < /home/eymen/Documents/Compiler/llvm15/llvm/test/MC/RISCV/crypt.s | /home/eymen/Documents/Compiler/llvm15/build/bin/llvm-objdump -M no-aliases -d -r -

<stdin>:          file format elf32-littleriscv

Disassembly of section .text:

00000000 <.text>:
      0: 33 59 89 81    shlxor  s2, s2, s8
```

Figure 9.3 : Use of pipe operator

CHECK-ASM-AND-OBJ and therefore it does not check the CHECK-ASM line. It seems that the object dump resulted the correct string.

In figure 9.4 we can observe that the test is passed:

Another complexity LLVM-lit handles is the path of our compiled binaries. The test case is free of the paths and %s placeholders are populated by LLVM-lit.

9.2 Writing the MC Test

LLVM-MC can be run before writing the test to see the encoding or hand-coding is also possible. It would be great if it was automatic like the llc test cases.

The header of the MC test can be used, if prefixes other than the present one are needed, add them to the --check-prefixes= side and make sure to use them in the file.

9.3 LLC/CodeGen Test

```
→ RISCVC git:(shlxor_instr) × ../../../../build/bin/llvm-lit crypt.s
-v
-- Testing: 1 tests, 1 workers --
PASS: LLVM :: MC/RISCV/crypt.s (1 of 1)

Testing Time: 0.07s
Passed: 1
```

Figure 9.4 : Output for successfully passed test

```
→ RISCv git:(shl_xor_instr) × ../../../../utils/update_llc_test_checks.py --llc-binary ../../../../build/bin/llc shl_xor.ll
```

Figure 9.5 : Command for using utility script

```
; NOTE: Assertions have been autogenerated by utils/update_llc_test_checks.py
; RUN: llc -mtriple=riscv32 -verify-machineinstrs < %s \
; RUN: | FileCheck %s -check-prefix=RV32R

define i32 @shl_xor(i32 %a, i32 %b) nounwind {
; RV32R-LABEL: shl_xor:
; RV32R:          # %bb.0:
; RV32R-NEXT:     shl_xor a0, a0, a1
; RV32R-NEXT:     ret
    %1 = shl i32 %a, 1
    %2 = xor i32 %1, %b
    ret i32 %2
}
```

Figure 9.6 : LLC test file

llc tests are more familiar since they check how an .ll produces Assembly-like (MCInst) strings. Instead of the manual checking, this tool can be used and batch-testing can be done.

Here is an example of the LLVM-IR code implemented for the SHLXOR instruction that we want. As it can be seen from its signature, it is a function taking two 32 bit integers and returning one. It takes one input, shifts it left by one and assigns it to a variable %1. %1 is then XOR'ed with the second input and the result is returned.

```
1 ; RUN: llc -mtriple=riscv32 -verify-machineinstrs < %s \
2 ; RUN: | FileCheck %s -check-prefix=RV32R
3
4
5 define i32 @shl_xor(i32 %a, i32 %b) nounwind {
6 %1 = shl i32 %a, 1
7 %2 = xor i32 %1, %b
8 ret i32 %2
9 }
```

In contrast to MC tests, we can use a utility script as shown in figure 9.5 to generate the expected result for us.

The script is located at llvm/utils/

The llc test file is populated with FileCheck lines as seen in figure 9.6.

We can observe that our newly added SHLXOR instruction is recognized and placed in the check lines. When LLVM-lit is run with the test file, the test is passed as the new instruction is implemented. Output for a successful test is given in figure 9.7.

9.4 Running Tests

```
→ RISCVC git:(shl_xor_instr) × ../../../../build/bin/llvm-lit shl_xor.ll -v

-- Testing: 1 tests, 1 workers --
PASS: LLVM :: CodeGen/RISCV/shl_xor.ll (1 of 1)

Testing Time: 0.08s
Passed: 1
```

Figure 9.7 : Output for passed .ll test

For running the tests, we should run the following command.

```
build/bin/llvm-lit llvm/test/PATH-OF-FILE -v
```

REFERENCES

- [1] (2021), Learn the basics of instruction set architecture - EDN Asia, <https://www.ednasia.com/learn-the-basics-of-instruction-set-architecture>, [Online; accessed 15. Apr. 2023].
- [2] (2020), RISC-V Assembly Language, <https://web.eecs.utk.edu/~smarz1/courses/ece356/notes/assembly>, [Online; accessed 15. Apr. 2023].
- [3] **Waterman, A., A.K.**, (2019), The RISC-V Instruction Set Manual Volume I: Unprivileged ISA.
- [4] **Gholizadehazari, E.**, (2021), An FPGA Implementation of a RISC-V Based SOC System with Custom Instruction Set for Image Processing Applications.
- [5] (2023), LLVM support for the draft Bit Manipulation Extension for RISC-V – Embecosm, <https://www.embecosm.com/2019/10/22/llvm-risc-v-bit-manipulation-extension>, [Online; accessed 15. Apr. 2023].
- [6] **Wolf, C.**, (2021), RISC-V Bitmanip Extension.
- [7] (2023), Scalar Cryptography Instruction Set Extension Group Names Diagram - Home - RISC-V International, <https://wiki.riscv.org/display/HOME/Scalar+Cryptography+Instruction+Set+Extension+Group+Names+Diagram>, [Online; accessed 15. Apr. 2023].
- [8] **Şairoğlu, M.**, (2020), Low-Power General Purpose Processor Design and Instruction Set Extension for AES, https://web.itu.edu.tr/~orssi/thesis/2020/MuhammedSairoglu_tez.pdf.
- [9] **Aho, A.** (2007). *Compilers: Principles, Techniques, and Tools*, Addison-Wesley.
- [10] (2022), “Clang” CFE Internals Manual — Clang 16.0.0git documentation, <https://clang.llvm.org/docs/InternalsManual.html>, [Online; accessed 5. Jan. 2023].
- [11] (2022), opt - LLVM optimizer — LLVM 16.0.0git documentation, <https://llvm.org/docs/CommandGuide/opt.html>, [Online; accessed 5. Jan. 2023].
- [12] (2022), The LLVM Target-Independent Code Generator — LLVM 16.0.0git documentation, <https://llvm.org/docs/CodeGenerator.html>, [Online; accessed 5. Jan. 2023].
- [13] (2022), LLVM: lib/Target/RISCV/RISCVISelLowering.h Source File, https://llvm.org/doxygen/RISCVISelLowering_8h_source.html, [Online; accessed 5. Jan. 2023].
- [14] **Eryilmaz, Y.E.**, (2022), Extending the Instrucion Set of RISC-V Processor for ASCON Algorithm, https://web.itu.edu.tr/~orssi/thesis/2022/YunusEryilmaz_bit.pdf.

- [15] (2023), About RISC-V – RISC-V International, <https://riscv.org/about>, [Online; accessed 15. Apr. 2023].
- [16] **Waterman, A.**, (2016), Design of the RISC-V instruction set architecture.
- [17] **Altınay, O.**, (2021), Instruction Extension of RV32I and GCC Back End for ASCON Lightweight Cryptography Algorithm.
- [18] (2023), RISC-V Cryptography Extensions Task Group Announces Public Review of the Scalar Cryptography Extensions – RISC-V International, <https://riscv.org/blog/2021/09/risc-v-cryptography-extensions-task-group-announces-public-review> [Online; accessed 15. Apr. 2023].
- [19] **Lattner, C.**, (2010), Intro to the LLVM MC Project, <https://blog.llvm.org/2010/04/intro-to-llvm-mc-project.html>, [Online; accessed 15. Apr. 2023].

APPENDICES

APPENDIX A.1 : Installation of Software

APPENDIX A.1

1.1 Installation of Software

As the LLVM codebase is large and has many options while building from source, finding the right options that our computers can handle easily was both essential to get started and critical as it decides the time it takes to see a change in code to get compiled. For this purpose, we accumulated the commands and created a tutorial that we can use in the future.

```
1 git clone https://github.com/llvm/llvm-project
2 cd llvm-project
3 mkdir build
4 cd build
5 sudo apt install cmake, ninja-build, clang, lld
```

Listing A.1 : Clone Repository and Install Necessary Packages

```
1 cmake -S ../llvm . -G Ninja -DCMAKE_BUILD_TYPE="Debug" \
2 -DBUILD_SHARED_LIBS=True -DLLVM_USE_SPLIT_DWARF=True \
3 -DLLVM_BUILD_TESTS=True -DCMAKE_C_COMPILER=clang \
4 -DCMAKE_CXX_COMPILER=clang++ -DLLVM_TARGETS_TO_BUILD="all" \
5 -DLLVM_EXPERIMENTAL_TARGETS_TO_BUILD="RISCV" -DLLVM_ENABLE_LLD=ON
```

Listing A.2 : CMake Configuration We Used

With this command, we are choosing the type as debug. Shared_libs=TRUE causes all libraries to be built shared instead of static libraries. ..SPLIT_DWARF is set to True to minimize memory usage at link time. We want to use clang as the C compiler. Therefore, it is specified in the command as DCMAKE_C_COMPILER=clang. In addition to that, we want to use lld as the linker instead of gold, so we specify that as well. This configuration is the most efficient in terms of memory and disk usage among our previous attempts at building LLVM from source.

```
1 Ninja
```

Listing A.3 : To build from
scratch or to
rebuild files
with change,
automatically

```
1 ninja llc
```

Listing A.4 : To build llc only
which is the binary
we modify

While running ninja, CPU and ram usage significantly increases. All available cores are used capacity. This may prevent doing other tasks while running ninja. In order to prevent this one may opt to use the following command instead. It allows you to choose how many cores are going to be utilized.

```
1 ninja llc -j<number of cores to use>
```