SCompiler 项目架构与代码解析

整体架构与运行流程

SCompiler 项目将源代码依次经过多个阶段处理,最终生成字节码并由自定义虚拟机执行。**编译器前端**包括词法分析、语法分析和语义分析;接着将AST转换为中间表示IR,再由**编译器后端**生成虚拟机可执行的指令序列,最后交由**虚拟机引擎**解释运行。具体流程如下:

- 1. **词法分析**: LexerEngine 读取源码文本,输出有序的Token序列。
- 2. **语法分析**: ParserEngine 将Token序列解析为抽象语法树(AST)节点列表。本语言采用类似Python的模块/函数/缩进语法,由顶层的 ModuleParser 解析模块, FunctionParser 解析函数,内部使用运算符优先解析表达式。
- 3. **语义分析**: SemanticAnalyzerRunner 遍历AST,检查标识符引用、类型匹配等语义规则,将符号登记到符号表并报告语义错误。
- 4. **IR生成**: IRProgramBuilder 遍历AST构建中间表示IR(Intermediate Representation),每个函数生成对应的 IRFunction ,以三地址码形式表示计算过程。
- 5. **目标代码生成**:对每个IR函数,先进行寄存器分配,将IR的虚拟寄存器映射到虚拟机的局部槽位,再由 VMCodeGenerator 将IR指令翻译为等价的虚拟机指令序列。例如,对一个加法表达式产生形如 LOAD 、 ADD 、 STORE 的指令。函数调用指令在生成时会填充或留存被调函数地址以供回填。
- 6. **虚拟机执行**: VirtualMachineEngine 加载指令列表,初始化运行栈,然后开始取指解释执行。每条指令由命令处理器解析并操作模拟的操作数栈、局部变量表和调用栈,直至程序结束。下面将对各模块进行详细分析。

词法分析模块

词法分析器(Lexer) 将源码文本转换成Token序列。 LexerEngine 是核心入口类,其构造时会根据定义的规则初始化一组 TokenScanner 子扫描器,并立即对输入源码进行扫描。各扫描器按优先级排列,包括:

- WhitespaceTokenScanner 跳过空白字符(空格、制表符等)
- NewlineTokenScanner 识别换行符,生成换行类型Token
- CommentTokenScanner 识别单行或多行注释文本
- NumberTokenScanner 识别整数和浮点数字面量
- IdentifierTokenScanner 识别标识符和关键字(如 module, if, return 等)
- StringTokenScanner 识别字符串字面量
- OperatorTokenScanner 识别运算符符号(+、-、*、/、== 等)
- SymbolTokenScanner 识别分隔符和符号(如逗号、括号、冒号等)
- UnknownTokenScanner 捕获无法识别的字符并标记为错误Token

扫描过程采用 **多通道扫描**: Lexer按字符顺序读取输入,对每个字符依次尝试上述扫描器,哪个扫描器的 canHandle() 方法返回true就交由 其 handle() 处理。每个扫描器从 LexerContext 获取当前字符流状态,消费相应字符序列并通过 TokenFactory 创建Token加入结果列表。例如, IdentifierScanner读取字母序列后由TokenFactory判断是标识符还是关键字。Lexer会跳过空白和注释,不生成多余Token。扫描循环持续直到输入结尾,然后显式追加一个EOF(TokenType.EOF)作为结束标记。整个词法分析的输出是有序的Token列表,可供语法分析器消费。

健壮性:当前Lexer对未知字符不会抛异常,而是生成类型为UNKNOWN的Token并继续。这在后续解析中会导致错误。可改进之处包括:当出现非法字符时Lexer直接报错或记录错误信息,而不是依赖Parser再处理。此外,Lexer已考虑换行和缩进控制,但未显式引入INDENT/DEDENT Token:目前通过冒号:和显式 end 标记处理块结构(见下文),没有真正依据缩进层级判断作用域。这种处理简化了实现,但使语法与缩进规则不完全一致,是设计上的折衷。

语法分析模块

语法分析器(Parser) 将Token序列按照语言文法规则还原为抽象语法树(AST)。SCompiler语言的文法大致为: *模块*由 module 模块名: 及缩进的多个函数定义组成,模块以 end module 结束;*函数定义*以 function 函数名(参数列表): 开始,内部缩进块包含若干语句,以 end function 结束。语句包括变量声明、赋值、条件 if/else、循环 loop、返回 return 以及表达式调用等。

Parser采用**递归下降**和**运算符优先解析**相结合的方法:顶层结构(模块、导入、函数)由不同解析器处理,表达式部分使用Pratt算法。主要组件如下:

- ParserEngine 语法分析驱动器,从Token流中识别顶层构造。它不断查看下一个Token,如果是模块关键字则调用模块解析器,直到遇到 EOF。目前顶层只注册了 module 关键字对应的解析器。若遇到未注册的顶级标记则抛异常。ParserEngine还会跳过空行以提高容错性。(注意:源码中 while (ts.isAtEnd()) 用来循环读取顶层构造,但 isAtEnd() 的实现返回值逻辑相反,实际应为 !isAtEnd() ,这一细节在实现中有所疏漏)。
- ModuleParser 模块解析器,实现了接口 TopLevelParser 。它期望当前Token序列形如 module: 模块名 NEWLINE ... end module 。描述了模块解析流程: 首先匹配关键字 module 及冒号,然后读取模块名称(IDENTIFIER),并期望接着一个换行开始模块体。显示了模块体内部的解析逻辑: 循环读取内部语句,跳过空行;遇到 import 开头则调用ImportParser解析导入语句列表,遇到 function 则调用FunctionParser解析函数定义,直到遇到 end 关键字标志模块结束。模块结尾强制要求 end module 成对出现。每个模块解析完成后,构造一个AST的ModuleNode节点,包含模块名、导入列表和函数列表。
- ImportParser 导入语句解析器,处理 import 模块名形式,将被导入模块名加入AST的ImportNode列表。项目中ImportParser会一次解析可能的多个连续导入语句,将结果返回给ModuleParser。需要注意的是,目前Import仅解析但并未触发跨模块代码加载(这在语义分析阶段验证)。
- FunctionParser 函数定义解析器,实现了 TopLevelParser 接口(模块内也会用到)。FunctionParser利用了一个自定义的**区块解析工** 具 FlexibleSectionParser ,将函数定义的不同部分(参数列表、返回类型、函数体)作为可选区块处理。函数解析流程:匹配 function: 标识并函数名,读入函数名后期望换行;然后通过 getSectionDefinitions() 预先注册三个区块解析器: parameter 参数列表、 return_type 返回类型和 body 函数体。接着调用 FlexibleSectionParser.parse() 循环解析这几个部分,直到遇到 end 关键字。其中:
 - 。 **参数列表**: Expect关键字 parameter: 后换行,然后对每行形如 declare 名称: 类型 的声明解析为ParameterNode节点。Parser会忽略参数声明行尾的注释。所有参数节点收集在列表中。
 - 。 返回类型: Expect return_type: 后读取类型名,如果省略则返回类型为null表示无返回。
 - 。 **函数体**: Expect body: 后逐行解析函数内部语句,直到遇到 end body 结束块。函数体解析使用了 StatementParserFactory 按当前行首关 键字选择对应的语句解析器。所有解析出的StatementNode语句节点加入函数体列表。函数末尾期望 end function 结束。最后组装 FunctionNode节点,包含函数名、参数列表、返回类型和函数体AST节点列表。
- StatementParsers 语句解析器集合,用于解析函数体内各类语句。根据语言设计,项目实现了若干种语句解析器并在 StatementParserFactory 中静态注册:包括变量声明(DeclarationStatementParser 解析 declare x: int = 表达式)、赋值语句(AssignmentParser 解析 x = 表达式 ,实现在Semantic阶段处理,语法上赋值被当作表达式或通过特殊形式解析)、条件语句(IfStatementParser 解析 if ... then ... else ... end 块)、循环语句(LoopStatementParser 解析 loop ... end 块)、返回语句(ReturnStatementParser 解析 return 表达式?)等。此外注册了默认的 ExpressionStatementParser 用于处理不匹配上述关键字的行,解释为一般表达式语句。StatementParserFactory在获取解析器时,会根据当前行第一个Token的字面量选择相应解析器,没有匹配的则使用默认表达式语句解析器。例如,函数体中的一行如果以 if 开头则交给IfStatementParser,否则如果是标识符开头(且不在其他分类)就作为表达式语句解析(常用于函数调用语句等)。
- 表达式解析 采用 Pratt 算法实现,位于 parser.expression 包下。由 PrattExpressionParser 根据运算符优先级解析中缀表达式,结合 PrefixParselet 和 InfixParselet 子类处理不同运算符。支持的表达式类型包括: 二元运算(加减乘除、比较等,由 BinaryOperatorParselet按优先级处理)、字面量(NumberLiteralParselet, StringLiteralParselet)、标识符引用(IdentifierParselet)、函数调用 (CallParselet识别 ()调用)、成员访问(MemberParselet识别 . 运算符)、括号分组(GroupingParselet)等。表达式解析结果构建对应的 ExpressionNode AST节点,如BinaryExpressionNode、CallExpressionNode等。优先级和结合性在Precedence枚举中定义。

健壮性:语法分析目前采用严格的期望匹配(expect()方法)和显式的 end 结束符,能够检测出缩进层级错误或缺失 end 等问题并报错。但由于没有维护一个专门的缩进栈,缩进不正确主要通过 end 不匹配来发现。设计上原本计划基于缩进来确定作用域(类似Python),但实际实现中通过:和成对的 end 显式标记块级范围,这种混合风格稍显冗余。改进方向可以是统一语法风格,例如完全采用缩进+换行而无 end 关键字,或采用显式花括号/ end 而无冒号。

另一个问题是Parser对顶层非模块内容的支持不完善:当前TopLevelParserFactory只注册了模块解析。如果源码未以 module 开头(比如仅有独立函数定义或语句),ParserEngine将无法识别顶层结构而报错。这在设计上可能考虑过支持**脚本模式**(无模块包裹的代码),从IR生成器看也有相应处理逻辑(将顶层Statement封装进_start函数),但由于顶层解析未实现直接处理Statement,实际使用中需要至少有一个模块声明。为提高

灵活性,可扩展Parser支持隐式模块包装或允许顶层函数定义。最后,Parser目前对错误恢复支持有限,一旦遇到语法错误通常终止分析;将来可考虑在捕获错误后跳过一定Token继续分析,收集多个错误再统一报告。

语义分析模块

语义分析器(Semantic Analyzer) 以AST为输入,在不改变结构的前提下检查和补充程序的意义信息。主要任务包括:标识符解析、作用域和生命周期管理、类型检查(如果有静态类型)、控制流合法性检查等。SCompiler的语义分析由 SemanticAnalyzerRunner 统一调度,内部调用 SemanticAnalyzer 执行具体步骤:

- 1. **模块和符号表初始化**: SemanticAnalyzer首先收集所有ModuleNode,注册模块到全局模块表(ModuleRegistry)。每个模块对应一个 ModuleInfo对象存储模块名、包含的函数签名等信息。接着调用 SignatureRegistrar 登记函数签名。SignatureRegistrar遍历每个模块,验证 Import的模块是否存在于ModuleRegistry中,并将该模块内每个函数的名称、参数类型、返回类型登记到ModuleInfo中。如果参数或返回类型名称无法识别(例如不存在的类型),则记录语义错误但尽可能继续分析。这样,在正式检查函数体前就建立了跨模块的函数签名表,便于 检测函数调用是否合法。
- 2. 函数体语义检查: 随后SemanticAnalyzer对每个模块内的函数逐一检查,由 FunctionChecker 执行具体逻辑。FunctionChecker会为每个函数 创建局部符号表(SymbolTable),并将函数形参作为变量首先注册到符号表中。然后遍历函数体中的每条语句节点,利用预先注册的语义分析器(StatementAnalyzer/ExpressionAnalyzer)对不同类型的语句和表达式进行检查。例如,对变量赋值语句,AssignmentAnalyzer会检查 变量是否已在当前或上层作用域声明、赋值类型是否兼容;对函数调用,CallExpressionAnalyzer会检索被调用函数在ModuleRegistry中是否 存在以及参数数量和类型是否匹配被调函数签名等;对控制结构,如IfAnalyzer检查条件表达式类型(预期为布尔)以及 then/else 分支的返回路径一致性,LoopAnalyzer检查循环初始化和更新部分等。每个Analyzer通过Context上下文对象获取所需信息并记录错误。符号表在进入 函数时创建,离开函数时销毁(当前实现中无嵌套块作用域,每个函数体使用一个SymbolTable,参数和局部变量都登记其中)。对于未声明的变量或不支持的语法,分析器将添加 SemanticError 错误对象。SemanticAnalyzer汇总所有错误,由 SemanticAnalysisReporter 统一输出错误列表;如果存在语义错误则终止后续编译流程。

项目提供了一些基础的语义检查实现,如:变量重复定义检测、变量未声明就使用、函数重复声明、函数调用不存在、参数不匹配等。类型系统方面,SCompiler目前是静态类型雏形:定义了BuiltinTypeRegistry含基本类型(int,string,float,bool,double等)枚举和值对象,但**尚未严格执行类型检查**。例如,变量声明了类型但赋值时类型不符的情况,Analyzer可能未完全实现检查。在SignatureRegistrar会对未知类型记录错误,但像算术运算左右项类型兼容性、条件表达式类型等并未全面检查。这是因为当前语言实现对类型要求不严格(更像动态类型处理)。**改进**:后续可加强类型检查逻辑,在ExpressionAnalyzer中加入对操作数类型的校验,对不匹配情况给出错误。

语义分析结束后,若无错误,即可确保AST上的每个标识符引用都有定义,每个函数调用可解析到目标函数,每条语句在语义上是合理的。这为后续代码生成提供了可靠基础。语义阶段还可以在AST节点上附加类型注解,便于代码生成选择合适的指令。目前Context的parseType会将类型名字符串转换为Type对象(如BuiltinType.INT等)供符号记录,但AST节点未存储类型属性,后端主要依赖变量的Type信息选择指令。整个SemanticAnalyzer设计采用了**注册表 + 分派**模式:通过AnalyzerRegistrar将具体AST节点类型与相应SemanticAnalyzer关联,FunctionChecker运行时按节点类型获取分析器。这种设计方便扩展新的语义检查规则,只需增加对应节点的Analyzer并注册即可。

健壮性: 语义分析目前存在一些可改进点: 一是错误收集不完善,目前遇到语义错误就记录并在最后可能退出,但对于不同模块或函数的多个错误能否一次收集输出没有详述。可以改进为**非终止检查**,在尽可能继续分析后面的代码前提下收集所有错误。二是模块导入解析还不完整,Import只验证模块是否存在于当前编译单元内的ModuleRegistry; 若要支持跨文件模块,需要扩展编译器能够根据import加载其它源文件或已编译中间文件(当前FileIOUtils等类可能为此做准备)。三是作用域管理仅支持函数级,缺少块级作用域支持(如在if/loop内部声明的变量不单独成域);如需支持,更细粒度的SymbolTable嵌套和Analyzer处理需实现。四是类型系统如上所述比较原始,没有检查表达式的类型正确性(例如算术运算混用int和string),未来可引入类型推断或强制类型检查,使语言更加健壮。

中间代码(IR)生成

在通过语法和语义检查后,编译器将AST转换为中间表示IR(Intermediate Representation)。IR是一种适合进一步翻译优化的**抽象指令序列**,通常比AST更贴近目标机但又独立于具体机器。本项目IR设计为**基于虚拟寄存器的三地址码**形式,每条IR指令类似于 dest = op(arg1, arg2) 的结构。主要实现位于 org.jcnc.snow.compiler.ir 包下,包括IR指令类、值类和IR构建器。

IRProgramBuilder负责遍历AST构建完整的IR程序。其 buildProgram(List<Node> roots) 方法对传入的AST根节点列表依次处理: 若节点是 ModuleNode,则对该模块下每个FunctionNode构建IRFunction加入IRProgram; 若节点本身是FunctionNode(顶层函数),直接构建加入; 若是 StatementNode(顶层语句),则会被包裹成一个伪函数"_start"再构建。这样支持了脚本模式下没有显式函数的语句执行。IRProgram相当于整个程序的IR表示,内部维护一个IRFunction列表及全局常量池等(当前实现主要管理函数列表)。

IRFunction对应源代码中的一个函数,包含该函数的IR指令序列、虚拟寄存器集合等信息。构建IRFunction由 FunctionBuilder 完成:它以AST的 FunctionNode为输入,创建一个空的IRFunction,然后:

- **参数处理**:为FunctionNode的每个参数创建一个新的虚拟寄存器(IRVirtualRegister)。IRFunction维护一个按序参数寄存器列表,参数寄存器 也算作函数体可用的虚拟寄存器。IRContext 会将参数名绑定到对应寄存器上,以便函数体后续IR生成时能找到参数。如所示:对每个参数 调用 irFunction.newRegister()生成新虚拟寄存器,然后通过 irContext.getScope().declare(参数名,寄存器)建立名字到寄存器的映射,并记录到IRFunction参数列表中。
- 函数体处理: FunctionBuilder创建 StatementBuilder 用于生成函数体内部各语句的IR。随后循环遍历AST函数体的每个StatementNode,由 StatementBuilder逐个处理。StatementBuilder根据语句类型构建不同的IR序列(具体见下文)。它持有同一个IRContext,保证多个语句间能 共享和更新作用域信息。所有IR指令通过 irContext.addInstruction() 添加到IRFunction的指令列表中。处理完所有语句后,一个IRFunction 就构造完成并返回给IRProgramBuilder,后者加入IRProgram。中各分支确保所有顶层代码最终都转换为IRFunction(若有顶层Statement,会包装到 _start 函数,从而统一处理方式)。

IR指令由抽象类 IRInstruction 及其子类表示,每种中间操作对应一个类。例如:IRAddInstruction表示整数加法,有属性dest(结果寄存器)、Ihs和rhs(操作数,可为寄存器或常量);IRReturnInstruction表示函数返回;IRJumpInstruction表示无条件跳转;CallInstruction表示函数调用等。每个IRInstruction实现方法 dest() 返回目标寄存器(若有),operands() 返回操作数列表,以及 op() 返回操作码枚举IROpCode。IROpCode是内部定义的枚举,列举了IR支持的操作码类型,如ADD_I32、SUB_I32、CALL等。当前IR实现针对32位整数运算进行了主要支持,例如IRAddInstruction的 op() 固定返回ADD_I32。虽然设计上考虑了不同数据类型(IROpCode中可能有ADD_F32等),并在ExpressionBuilder中有解析数值字面量后缀以决定操作数类型的逻辑(如识别 100L 为long型),但实际IR指令类并未针对其他类型定制,不同类型常量目前均以NumberLiteralNode和IRConstant统一处理,算术IR指令也都以I32后缀为主。这意味着当前IR阶段未充分区分多种数据类型,浮点和长整型等运算可能被错误地按I32处理或不支持。这属于功能上的不完整,未来可扩充更多IR指令种类或参数化IR指令以支持不同类型运算。

StatementBuilder实现将AST语句节点转换为一系列IR指令。它内部通过 exprBuilder = new ExpressionBuilder(irContext) 来处理表达式部分,从而实现递归构建。例如列出了不同语句的处理要点:

- ExpressionStatementNode(表达式语句):直接调用ExpressionBuilder构建表达式,将结果存入一个寄存器,但由于表达式语句的结果不需要保存,生成指令后忽略结果寄存器。比如函数调用语句,会生成CallInstruction得到返回值寄存器,但不进一步使用。
- AssignmentNode(赋值语句): 先构建右侧表达式获取结果寄存器,然后检查赋值的变量名在当前作用域是否已声明。如果未声明,说明是第一次赋值,则通过 irContext.getScope().declare(名称,寄存器) 将变量绑定到该寄存器(相当于变量声明);如果已存在,则用 scope.put() 更新变量名绑定的寄存器为新值寄存器。这种做法意味着每次赋值都会产生新的虚拟寄存器,旧值的寄存器不再作为该变量的当前值(形成一种Static Single Assignment风格)。这样在IR层实现了变量重新赋值,而无需单独的Store指令。
- **DeclarationNode**(声明语句):如果带初始化器,则构建初始化表达式的IR,将结果寄存器声明绑定到变量名;如果没有初始化,则仅调用 scope.declare(name) 分配一个新寄存器给该变量但不赋值。声明语句本身不产生独立的IR指令(除非有初始化表达式,需要先计算出结果)。
- **ReturnNode**(返回语句):若有返回值表达式,则先构建表达式IR得到值寄存器,再调用InstructionFactory.ret生成ReturnInstruction,将该寄存器标记为返回值;如果无返回值,则生成RetVoid指令表示返回。Return指令在IR级别会作为函数结尾指令,用于后端生成适当的RET/HALT等。
- 控制流语句(If/Loop等):目前StatementBuilder中没有针对IfNode或LoopNode的处理分支。实际上项目实现了IfNode、LoopNode的AST和相应SemanticAnalyzer,但IR生成尚未支持它们。这意味着遇到条件或循环语句,StatementBuilder会落入默认分支抛出异常"Unsupported statement"。因此当前版本编译器无法直接编译含有if/loop的函数(除非采取特殊措施绕过IR阶段)。这是实现上的一大不足,后续需补充相应IR构造逻辑。例如,可为IfNode生成条件判断和跳转IR(类似"三地址码"中的条件跳转),为LoopNode生成循环初始化、条件检查、跳转等IR指令(如IRJumpInstruction)。目前项目中虽有IRJumpInstruction类等,但未被使用。这属于功能不完整的部分,下文将建议改进。

IR生成完毕后, IRProgram 中汇集了所有函数的IR表示。调用 IRProgram.toString() 可以打印出IR的简要文本表示,例如变量寄存器通常显示 为 %0, %1 等格式,指令形如 %2 = %0 + %1。在SnowCompiler主程序中,会打印IR用于调试。IR作为中间产物,一方面可以进行优化(当前未实现显式优化流程,仅留下IROptimizer接口作为扩展点【17†output】),另一方面便于进行下一阶段目标码生成。

目标代码生成与虚拟机指令集

代码生成(Code Generation) 负责将IR翻译为虚拟机能够执行的指令序列。SCompiler的目标指令集是定制的Snow虚拟机字节码(文本形式表示),设计思想接近JVM:采用操作数栈+局部变量槽模型,并区分不同数据类型的操作。代码生成主要分两步:**寄存器分配和指令序列生成**。

• **寄存器分配(Register Allocation)**: IR中的虚拟寄存器需要映射到虚拟机的局部变量槽 (slot)。 RegisterAllocator.allocate(IRFunction fn) 完成这一任务。策略是简单的线性分配: 首先按函数参数列表顺序为每个参数分配连续 槽位0,1,...; 然后扫描函数体每条IR指令,若指令有dest结果寄存器且尚未分配槽,则分配新槽;再检查指令的每个源操作数,凡是

IRVirtualRegister且未分配槽的也立即分配。这种方式确保一个函数内每个出现过的虚拟寄存器都唯一定义一个槽号,寄存器生命周期不重叠也不会复用槽(未做活跃分析优化)。RegisterAllocator返回一个不可变的映射表 Map (IRVirtualRegister, Integer) 供后续使用。例如,若IRFunction有参数%0、%1,则它们映射slot0和1;函数体第一次出现的新寄存器%2映射slot2,以此类推。这样一来,变量和中间计算在虚拟机中都有了各自的"存储地址"。

- **虚拟机指令生成**: vMCodeGenerator 负责将单个IRFunction翻译为具体的指令文本。构造VMCodeGenerator时传入上一步计算的slotMap和一个 vMProgramBuilder 输出器。生成过程通过遍历IRFunction的指令列表实现: 对不同IR指令类型调用不同的生成函数:
 - 。 **常量加载**: IR的LoadConstInstruction表示将立即数载入目标寄存器。生成两条虚拟机指令: 首先将常量值压入操作数栈(例如使用 I_PUSH <value> 将整数常量压栈),然后从栈弹出该值存入目标槽位(如 I_STORE slot 将值存到本地变量表指定槽)。引用示例: IRConstant 5 -> I_PUSH 5; I_STORE 0(将5存入slot0)。
 - 。 **算术/二元运算**: IR的BinaryOperationInstruction涵盖加减乘除比较等操作。生成流程是:将两个操作数从本地槽加载到栈(如 I_LOAD slota; I_LOAD slota),执行对应运算指令,然后将结果存回目标槽。例如IRAdd(dest=%3, lhs=%1, rhs=%2)假设%1映射 slot1、%2映射slot2,则输出指令: I_LOAD 1; I_LOAD 2; I_ADD; I_STORE 3。这里 I_ADD 是Snow VM定义的"32位整数加法"操作码。实际上,VMCodeGenerator通过 IROpCodeMapper.toVMOp()将IR指令的IROpCode转换为相应的VM操作码助记符。目前由于IR基本都为I32类型操作码,映射多是直观对应的,如ADD_I32 -> I_ADD。若后续扩展不同类型运算,IROpCodeMapper会相应扩充映射,比如ADD F32映射 F ADD等。
 - 。 **一元运算**:如IR的UnaryOperationInstruction(目前可能用于数值取负等),生成类似过程:加载操作数,执行一元运算指令,将结果存 回。
 - 。 函数调用: IR的CallInstruction包含目标函数名和实参列表(IRValue列表)以及返回值寄存器。VMCodeGenerator对每个实参依次生成 I_LOAD argSlot 指令,将参数值压入栈。然后调用 VMProgramBuilder.emitCall(函数名,参数个数) 生成一条调用指令 CALL ⟨addr⟩ ⟨nArgs⟩。这里是被调函数入口地址: 如果该函数已处理过(已在labelAddr映射中),则直接填入其起始地址;如果尚未出现,则暂时填入占位符 -1 并记录到unresolved列表,待后面真正定义该函数时回补地址。这种单遍代码生成+回填机制使函数可以在定义前被调用。最后,如果调用有返回值,则紧接生成一条 I_STORE destSlot 将栈顶返回值存入调用方的目标槽。例如,调用CommonTasks.test(a,b)返回值赋给%5(slot5),且test函数在最终代码中位于地址10且有2参数,则调用处指令可能是: I LOAD ⟨slot a⟩; I LOAD ⟨slot b⟩; CALL 10 2; I STORE 5。
 - 。 **返回**: IRReturnInstruction翻译为虚拟机的返回或终止指令。若有返回值寄存器,则先加载该值到栈(如 I_Load slotx);然后需要根据返回的是主程序(入口函数)还是一般函数决定操作码: Snow VM约定主函数返回时应终止整个程序,而普通函数返回时跳回调用点继续执行。因此VMCodeGenerator做了判断: 如果当前函数名是"main",则输出 HALT 指令终止VM;否则输出一般的 RET 指令表示函数返回。在Snow虚拟机中,HALT会使主循环退出,而RET则弹出当前栈帧并跳转到先前保存的返回地址继续执行。这种特殊处理方式意味着**入口函数必须命名为"main"**才能在结束时正确停止虚拟机。如果用户主程序模块函数不用main命名,当前实现会错误地用RET返回导致VM尝试跳转到调用者(但main没有调用者)而可能引发异常或无效果。这是不够健壮之处,后文将讨论改进。

代码生成完成后,所有函数的指令已添加到 vMProgramBuilder 中。VMProgramBuilder的职责类似于汇编器,它维护一个 code 列表按顺序存放最终指令字符串,并处理函数标签和调用修补。当 beginFunction(name) 被调用时,会记录当前指令地址pc作为该函数入口地址并回填之前所有针对该函数的未解析调用。SnowCompiler主程序对每个IRFunction顺序调用VMCodeGenerator.generate后,最后调用 builder.build() 得到完整指令列表。在未出现未解析调用目标的前提下(否则build()会抛异常)即可进入执行阶段。

Snow虚拟机指令集支持整数、长整数、短整数、浮点、双浮点、字节等多种类型操作,对应不同前缀的操作码(I_, L_, S_, F_, D_, B_)。指令可分为: 栈操作(PUSH, POP, DUP, SWAP等)、算术运算(ADD, SUB, MUL, DIV, MOD及一元NEG,按类型区分如I_ADD, D_ADD等)、位操作(AND, OR, XOR等)、类型转换(例如I2F, F2D将栈顶值类型转换)、内存读写(LOAD/STORE将局部槽与栈之间传值,如I_LOAD n, I_STORE n)、函数调用与返回(CALL addr nArgs, RET)、条件跳转(比较指令如ICG="if greater"等通常会消耗栈顶两个操作数做比较并跳转)以及虚拟机控制(HALT停机)。在本项目生成代码中,主要用到了I_LOAD/STORE, I_PUSH, I_ADD等I_*系列指令以及CALL/RET/HALT等指令;针对其他类型的指令类虽然实现了许多(如BAddCommand, DAddCommand等),但由于编译器尚未生成这些指令,因此在执行中未实际用到。

虚拟机执行逻辑

编译生成的指令序列交由**Snow虚拟机**执行。虚拟机由 org.jcnc.snow.vm 包的一系列类实现,包括引擎、指令定义、运行时栈和辅助工具。其设计采用**栈式虚拟机模型**:每条指令从操作数栈弹出操作数,计算后再将结果压回栈,或在需要时与方法的本地变量区交互。主要组件如下:

• VirtualMachineEngine: 虚拟机核心引擎,提供 execute(List<String> program) 方法执行指令列表。Engine内部维护:操作数栈(OperandStack)、全局本地变量存储(LocalVariableStore)、调用栈(CallStack)、程序计数器(programCounter)等状态。执行前,会调用 ensureRootFrame() 创建一个"根调用帧"(对应主模块运行环境)并压入调用栈。根帧的返回地址通常无效,用于占位。然后进入主循环,每次根据 programCounter 取出对应指令字符串,解析其操作码和操作数。Engine通过 parseOpCode() 将助记符转为内部整数opcode(通常来自VMOpCode枚举)。随后调用 commandExecutionHandler.handle(opCode, parts, currentPC) 执行该指令。handler返回下一个指令的地址 nextPC,Engine更新 programCounter 进入下一循环。循环持续,直到遇到终止条件:如果nextPC返回特殊值 PROGRAM_END 或 -1 (例如

HALT指令或运行中出错)则跳出循环。执行完成后,引擎提供方法打印最终操作数栈和局部变量表供调试。整个执行过程即典型的**取指-译** 码-执行循环。

命令执行器(CommandExecutionHandler): 封装具体指令执行逻辑的调度。Engine在每次循环取到opcode后,将实际执行委托给

- CommandExecutionHandler。Handler持有对OperandStack、LocalVariableStore、CallStack的引用。 其 handle(int opCode, String[] parts, int currentPC) 首先通过 CommandFactory.getInstruction(opCode) 获取对应的Command对象实例,如果找不到则抛异常。然后调用该Command的 execute(String[] parts, int pc, OperandStack, LocalVariableStore, CallStack) 方法,执行指令所需操作,并返回下一条指令地址。所有具体指令的实现都遵循 Command 接口,在其execute方法中编写针对OperandStack等的数据
 - 如果找不到则抛异常。然后调用该Command的 execute(String[] parts, int pc, OperandStack, LocalVariableStore, CallStack) 万法,执行指令所需操作,并返回下一条指令地址。所有具体指令的实现都遵循 Command 接口,在其execute方法中编写针对OperandStack等的数据操作和状态变更。例如,IAddCommand的execute会从OperandStack弹出两个int值相加,将结果压回栈,并返回 pc+1 指向下一指令;CALL指令的execute会读取操作数(函数地址、参数个数),创建新的StackFrame并调整调用栈,设置programCounter跳转到目标地址开始执行被调函数;RET指令的execute会弹出当前栈帧并设置返回地址为nextPC;HALT指令则返回-1以通知引擎停止执行。
- **运行时栈与帧**:调用栈(CallStack)管理着一系列活动的栈帧(StackFrame)。每个StackFrame封装一次函数调用的执行环境,包括:返回地址(returnAddress)、局部变量存储引用和在其中的起始索引、以及一个方法上下文(MethodContext,保存函数名等调试信息)。LocalVariableStore在Snow VM中被设计为所有帧共享的内存(类似JVM将所有栈帧的本地变量保存在同一个连续内存,由每个帧记录基址和大小)。例如,根帧初始化时base=0长度=...,调用新的函数时可能在LocalVariableStore中分配新段。LocalVariableStore提供根据帧base和索引存取局部变量的方法,并支持可视化(LocalVariableStoreSwing GUI类用于调试查看局部变量)。OperandStack用于算术运算和临时值传递,是所有指令操作的主要对象。指令集通过严格的入栈/出栈顺序实现运算次序,比如前述 I_LOAD n 指令就是将本地槽n的值复制压入栈顶; I_ADD 则弹出两个栈顶数相加,再将结果压栈。
- 指令类: org.jcnc.snow.vm.commands 包按功能分类定义了大量Command实现类,每个类对应一种字节码指令。例如,IAddCommand实现将OperandStack弹出两个int相加;IStoreCommand实现从栈弹出一个值存入指定局部槽;CallCommand实现创建新帧并跳转;RetCommand实现返回上一帧。这些类大多遵循命名约定,例如 IAddCommand 处理int加法, LAddCommand 处理long加法等。指令执行中的类型转换、算术溢出等未做特别检查,直接利用Java对应运算,因此可能存在如整数溢出不检测的问题,属一般情况。CommandFactory内部预建了一个opcode到Command对象的映射表或工厂方法,用于快速获取指令实例(具体实现可能通过静态初始化将各Command注册到一个数组)。通过这种工厂方式,每次执行无需反射创建新对象,而是重用已有Command实例,减少开销。优化考虑:当前Engine每步都字符串拆分指令然后查表获取Command,理论上可以在加载指令时就将文本解析为Opcode和参数并存储,这样执行时直接使用预解析结果,可提升速度。但由于实现简洁优先,目前在每次循环中进行解析,对一般规模程序影响不大。

综上,Snow VM以相对经典的栈机方式运行字节码指令,支持函数调用栈和局部变量区,实现了基本的流程控制和运算功能。配合前端编译器输出的正确指令序列,即可完成源程序的实际运行。例如,对于前述示例函数 CommonTasks.test (将两个整数相加),编译器可能输出指令: I_LOAD 0; I_LOAD 1; I_ADD; I_STORE 2; I_LOAD 2; RET (假设参数num1/num2槽为0/1,结果槽2)。虚拟机执行这些指令,最终将正确的计算结果留在操作数栈或局部变量,从而实现预期的功能。

健壮性:目前虚拟机实现基本完备,但仍有改进空间。例如,主函数返回采用HALT特殊判断不够优雅,如用户误用导致未HALT可能悬挂进程。更好的方式是引入程序入口概念,由VM配置哪个函数作为入口并在其返回时自动HALT。又如,错误处理方面,CommandExecutionHandler捕获了执行过程中的异常并终止程序。对于常见运行时错误(如除零、空栈等)可考虑提前检查并提供明确错误信息。目前指令集的冗余也值得注意——大量类似的指令类(如加法指令就有6种类型)导致代码重复率高,可以通过泛型或模板减少重复。性能方面,如果追求更高执行效率,可考虑直接将字节码编译为主机机器码(JIT)或采用更紧凑的二进制格式并优化取指循环。但这些在当前1.0版本中都不是首要目标。

源码结构与关键模块概览

SCompiler项目源码按功能划分为编译器和虚拟机两大部分,各自下设子模块。下面以表格形式列出主要源码文件/类及其功能:

模块/文件	功能简介
org.jcnc.snow.compiler.lexer	间法分析模块:将源码文本读入并切分为Token序列。 LexerEngine — 核心词法分析器,初始化各种Token扫描器并驱动扫描过程。调用 getAllTokens()可获取完整Token列表。 LexerContext — 提供字符流读取、回退、位置跟踪等功能。 TokenScanner 接口 — 定义扫描器规范,每种类型Token对应一个实现类,如WhitespaceTokenScanner、NumberTokenScanner等,由LexerEngine按顺序调用其 canHandle/handle 进行识别。 Token / TokenType — 封装词法单元信息,类型枚举包括IDENTIFIER,KEYWORD,NUMBER等。Token记录词法类型、原始文本及行列位置等。 其他: TokenFactory 根据扫描到的字面串判断关键字或类型; TokenPrinter 用于调试打印Token序列。

模块/文件	功能简介
	语法分析模块:将Token列表解析为AST。 ParserEngine — 语法分析主控,循环读取TokenStream,根据下一个Token选择顶层解析器处理。跳过空行,并确保遇到未知结构抛错。 TokenStream — 封装Token列表并维护当前位置,提供 peek,next,expect 等方法辅助解析。注意:isAtEnd()逻辑有误,返回条件与命名相反,需更正。
org.jcnc.snow.compiler.parser	 语句解析: StatementParserFactory 根据语句起始关键字分派相应解析器。已实现解析器: Declaration、Assignment、If、Loop、Return,以及默认Expression解析器等。 每个解析器将构造对应的AST节点,如IfParser构造IfNode (含条件ExpressionNode和then/else子节点列表)。 表达式解析:采用Pratt算法,PrattExpressionParser协调多个Parselet完成。 包内定义了各种PrefixParselet(如解析字面量、标识符)和InfixParselet(解析运算符优先级)
	子类。ExpressionNode 派生类如BinaryExpressionNode、CallExpressionNode等用以表示表达式AST。
org.jcnc.snow.compiler.semantic	语义分析模块:检查AST的正确性并准备后端信息。
	检查结束后SymbolTable随函数退出作用域销毁。 • 类型系统: BuiltinType 枚举了基本类型(int, float, string等), Type 类是类型抽象。

模块/文件	功能简介
	目前变量和函数定义可指定类型名,但表达式类型推断不完善,多数Analyzer未严格验证类型一致性,这部分留待后续增强。 • 错误报告: SemanticAnalysisReporter 收集Context中的SemanticError列表,输出错误详情并在有致命错误时终止编译。当前错误恢复策略有限,一旦遇错通常中止后续阶段。
org.jcnc.snow.compiler.ir	输出错误详情并在有效命错误时终止编译。当前错误恢复策略有限,一旦遇错通常中止后续阶段。 中间表示(R)模块:将AST翻译为IR代码,并准备后端生成。 • IRProgram — 表示整个程序的IR容器,包含多个IRFunction。提供 functions() 遍历函数等方法,支持打印IR内容。 • IRFrogram — 表示整个函数的IR,包含函数名、指令列表(List)、参数寄存器列表、虚拟寄存器计数器等。通过 newRegister() 生成新IRVirtualRegister (其编号自动递增)。 addInstruction() 添加IR指令到函数尾部。 • IRInstruction 抽象类 — 所有IR指令的基类,定义了 dest() 目标、operands() 操作数列表、op() 操作码等抽象方法,以及 tostring() 用于调试打印。具体IR指令子类: BinaryOperationInstruction/UnaryOperationInstruction (抽象类,扩展IRInstruction以表示有两个/一个操作数的算术指令),IRAddInstruction、IRSubInstruction(将常量载入新寄存器)、IRJumpInstruction(语数调用,保存函数名及参数列表)、LoadConstinstruction(将常量载入新寄存器)、IRJumpInstruction(招限eturninstruction(控制流指令)等。IRValue接回统一了IR操作数,可以是IRVirtualRegister(虚拟寄存器)或IRConstant(常量值)或IRLabel(姚转目标标签)。 • IRVirtualRegister(虚拟寄存器,成形仅有一个int编号标识。例如 26.31 这样的表示方法。 • IR内建: IRProgramBuilder 遍历AST构建IRProgram。遇到ModuleNode则处理其内部函数,FunctionNode直接处理,顶层Statemen则包装入临时FunctionNode再处理。实际构建由 FunctionBuilder.build(FunctionNode) 完成:它创建IRFunction并用RContext管理当前作用域。然后对函数参数依次生成虚拟寄存器并由clare绑定。接着用 StatementBuilder 构建函数体。 • StatementBuilder — 核心IR生成器,根据不同AST语句节点构建IR。列学了几种处理:表达式语句 > 直接build表达式,丢弃结果;赋值语句 > build右侧表达式,得到新寄存器,然后变量若首次出现则declare(相当于变相声明),否则更新绑定;声时语句 > bidd该应值表达式(若有),然后生成Ret或RetVold指令;IfLoop等控制语句 > 目前来实现,将导致不支持异常。 • ExpressionBuilder — 递归构建表达式。对常量、变量引用返回IRConstant或已有寄存器;对算术或比较表达式,会build子表达式得到IRVirtualRegister,再生成对应IIHStruction;对函数调用表达式,先构建每个实验获得寄存器列表,再生成对CallInstruction;对函数调用表达式,先构建每个实验获得寄存器列表,再生成对CallInstruction;目标函数名字符单根据是否有模块的磁频接定,再生成对应用各类。 • IRVirtualRegister — 建归构建。

| org.jcnc.snow.compiler.backend | 后端代码生成模块: 负责将IR转换为Snow VM指令序列。

- RegisterAllocator 寄存器分配器,将IR虚拟寄存器映射到局部变量槽索引。先分配参数槽位,再扫描指令分配剩余寄存器槽。算法简单未考虑生存期重用,未来可优化。
- VMProgramBuilder 虚拟机指令构建器。维护内部List存放指令文本,以及函数名->地址映射表labelAddr和调用待回填列表。提供 beginFunction(name) 标记新函数入口、 emit(String line) 添加一条指令并递增PC、 emitCall(targetFn, nArgs) 生成CALL指令并处理未解析目标、 build() 返回最终代码列表。通过这些接口,后端可以顺序生成指令且无需关心函数地址回填细节。
- VMCodeGenerator 核心代码生成器,将单个IRFunction翻译为VM指令。初始化时保存slotMap和VMProgramBuilder。 generate(fn) 按顺序遍历IR指令,针对不同类型调用 genLoadConst 、 genBinOp 、 genUnary 、 genCall 、 genRet 等方法生成指令。每个方法利用 VMProgramBuilder.emit封装具体输出,例如genBinOp会输出形如 I_LOAD a; I_LOAD b; <OP>; I_STORE c 指令序列。这里 <OP> 通过 IROpCodeMapper.toVMOp() 将IR操作码映射到VM操作码助记符(如ADD_I32->I_ADD)。项目实现了常见指令的映射表 IROpCodeMapper。

slotMap用于将IR寄存器转为具体槽号,如 slot(vr) 函数得到某IRVirtualRegister对应整数索引,用于LOAD/STORE指令参数。

• 其他: IROpCodeMapper 定义静态映射,从IR操作码枚举常量到VMOpCode枚举常量或助记符字符串。例如将IROpCode.ADD_l32映射为"ADD"或对应VMOpCode对象,供VMCodeGenerator使用。 vMMode 定义虚拟机运行模式(RUN/STEP等,当前基本用RUN)。

(注: backend模块的命名有些混乱,理论上VMCodeGenerator和VMProgramBuilder关系更近似于"汇编器",而非生成可独立的目标代码格式。但 这里将Snow VM字节码视作目标机器码进行生成。) |

| org.jcnc.snow.vm | 虚拟机模块:解释执行编译后的指令序列,实现程序运行。

- VirtualMachineEngine 虚拟机核心类,包含主循环执行逻辑。初始化时根据模式构造 OperandStack 、 LocalVariableStore 和 CallStack 等,并创建CommandExecutionHandler。execute()方法按PC循环取指、解析、执行,直到遇到HALT等终止。提供printStack/printLocalVariables用于调试输出当前栈内容。
- OperandStack 操作数栈类,内部用列表或数组模拟栈操作,支持push/pop/peek等。
- LocalVariableStore 局部变量存储,模拟内存来存放各帧的局部变量值。实现上可用一个Object数组+每帧基址管理,也可按帧链表各自维护。提供按索引存取的方法。LocalVariableStoreSwing是其GUI表示,用于在窗口中监视各槽位值变化。
- CallStack 调用栈,内部维护一个栈结构的StackFrame列表。提供pushFrame(StackFrame)和popFrame()实现函数调用和返回。还有printCallStack用于调试输出调用链。
- StackFrame 栈帧类,保存单次函数调用的状态:包含returnAddress(调用者下一执行地址)、对LocalVariableStore的引用及当前帧局部变量 起始偏移、方法上下文(MethodContext)等。MethodContext包括函数名、所属模块等元信息,辅助日志和调试。
- •命令集合: vm.commands 子包按照指令类型分目录,包含数十个Command类。如commands.arithmetic下有IAddCommand/ISubCommand等实现整数算术,LAddCommand等实现长整运算;control子包下ICGCommand等实现条件跳转(Compare Greater)、JumpCommand无条件跳转;stack子包有Push/Pop/Dup/Swap等栈操作;memory子包有各类型Load/Store命令;function包有CallCommand和RetCommand;vm包有HaltCommand终止执行。每个Command类实现接口 Command 的execute方法,利用OperandStack/LocalVariableStore完成操作并返回下一个PC。示例:IAddCommand.execute会pop出两个int,加法运算后push结果,并返回currentPC+1。CallCommand.execute则会根据参数创建新StackFrame,将当前PC+1作为returnAddress压栈,并返回被调函数入口地址,实现跳转调用。RetCommand.execute弹出当前帧,取得先前保存的returnAddress返回。展示了CommandFactory获取并执行Command的过程。
- CommandFactory 命令工厂,提供静态方法 getInstruction(int opcode) 返回Optional。中通过opcode查找对应Command实例。如果opcode无效则抛IllegalArgumentException。CommandFactory很可能在静态初始化时就创建好所有Command实例映射,以避免执行时频繁创建对象。
- vMopCode 虚拟机操作码枚举,每个枚举值代表一种指令,并携带一个整数编码和助记符字符串等属性。提供从字符串助记符解析opcode的方法 valueOf(),Engine使用它在parseOpCode时将文本opcode转为内部code。(具体实现可能直接用Enum.name()匹配或维护Map映射)
- LoggingUtils / VMStateLogger 日志工具类,用于根据需要打印调试信息,比如每执行一步的栈和变量状态等。可配置虚拟机在不同VMMode 下输出不同粒度log。当前RUN模式下日志默认关闭。 |

存在的问题与可改进之处

综合分析当前实现,SCompiler项目在功能和设计上还有一些不足与可以改进的地方:

- 词法/语法错误处理不够健壮: Lexer遇到非法字符仅生成UNKNOWN Token继续,使后续解析难以定位问题。不如在Lexer阶段就抛出带位置的LexicalException或收集错误。Parser对错误恢复也未实现,一旦遇错整个编译停止,这在编译大型源码时不友好。改进建议:实现Lexer错误异常并中止解析,或者Lexer将错误存入列表供最终输出;Parser则可在某些错误后尝试跳到下一语句继续解析,以发现更多错误,最后集中报告。
- **顶层结构限制**:目前必须有 module 包装源码,缺少直接编译脚本或REPL段代码的能力。如果需要支持,可在ParserEngine中增加对无模块情况的处理,例如隐式创建一个默认ModuleNode包含所有顶层函数/语句。事实上IR生成器已有对应思路(识别顶层Statement包装_start函数),但Parser未跟进。通过在TopLevelParserFactory注册例如"function"关键字解析器来允许顶层函数,或对无module时由编译器注入一个默认模块,都可以增强灵活性。
- 状态标志命名逻辑错误: Parser部分的 TokenStream.isAtEnd() 实现与常规含义相反,导致代码中使用它的逻辑也颠倒。当前 isAtEnd() 返回 true 表示还**没有**到流末尾(peek()!= EOF),这既不直观也易埋下bug。事实上已经导致 ParserEngine 用 while (ts.isAtEnd()) 循环解析,表面看逻辑正确但实则依赖了错误实现(应为while not at end)。建议修正 isAtEnd() 定义为 return peek().getType() == EOF ,并相应修改调用处逻辑为while(!isAtEnd())等。保持命名和语义一致可减少维护困扰。
- 缩进和块结构处理:设计上DSL类似Python用缩进表示块,但实现上引入了 end 关键字辅助结束。这混合风格虽然工作正常,但存在冗余。比如函数定义用了 function name:和 end function包围,同时还需正确缩进代码块。当前编译器实际上并未根据缩进深度产生INDENT/DEDENT Token,也就是说缩进仅在视觉上要求,在解析时可有任意空格而不影响结果(Parser跳过多余空行和空格)。这可能导致一些不一致:如果用户缩进层级不对但写了正确的 end ,编译器仍接受;反之缩进对齐却少写 end 会报错。因此改进方向是统一块结构语法:要么完全采用缩进严格决定层次(像Python,不需要 end),要么采用显式块标记(如花括号或end)而忽略缩进。选择其一并修改Lexer/Parser,相应地大幅简化规则或提高代码健壮性。目前这种混用方式主要是为了易实现,但从语言设计角度应在后续版本调整。

- 控制流IR/代码生成未完成: 目前编译器尚不能处理 if/else 和 loop 语句的代码生成。虽然AST和Semantic对它们已有表示和检查,但IR阶段 StatementBuilder 直接将它们标记为不支持。这意味着源代码中如果出现条件或循环,将在编译阶段报错终止。这显然是功能缺失的部分。应尽快完善IfNode和LoopNode的IR翻译: 例如为IfNode生成条件比较(IR可能需要一个IRCompare指令)及两条有条件跳转IR指令(跳过then块或跳过else块),还需要在IR中引入标签IRLabel表示跳转目的地等。LoopNode可转换为类似 initializer; loopStartLabel: if (!condition) jump loopEnd; (loop body IR); jump loopStart; loopEndLabel: 的IR序列。当前IRInstruction已有IRJump等类,说明计划中考虑了这些操作,只是未集成。完成这部分将极大拓展语言功能。相应地,VMCodeGenerator也需增加对条件跳转和循环的指令生成,例如实现比较后条件跳转(Snow VM已有ICG/ICL等指令)。总之,**补全控制流编译**应是下一步重点。
- 类型支持不完整: SCompiler词法上支持如 123L 后缀表示long、 3.14 表示double等,TokenFactory也能将 int、string、bool、float、double 识别为类型Token。但是编译流程并未真正区分这些类型的操作。IR层所有算术指令都假定32位整数(ADD_I32等),即使Literal有后缀 d 或 f ,ExpressionBuilder的resolveSuffix尝试识别类型却未被用于选择不同IR操作码(IRAddInstruction的op()始终返回ADD_I32)。虚拟机虽然实现了不同类型指令(I_ADD, L_ADD, F_ADD等),但编译器从未生成过后者。结果就是:写 float a = 1.2; float b = 3.4; float c = a + b; 这样的代码,编译器仍当作int相加生成I_ADD而不是F_ADD。显然这是不正确的。改进措施:完善类型系统,在AST节点附加类型信息,在ExpressionAnalyzer中推导表达式类型,并在IR生成时根据节点类型选择相应IROpCode(如浮点加法应产生ADD_F32等)。目前预留的resolveSuffix等逻辑可配合类型规则使用。一旦IR准确区分类型,VMCodeGenerator的IROpCodeMapper需同步更新映射,以输出正确的指令助记符(如ADD_F32->F_ADD)。这样才能真正支持多类型运算。另外,诸如字符串拼接 + 、不同类型间运算的转换,也需要在语义阶段明确和在代码生成阶段处理(比如插入I2F等类型转换指令)。总的来说,当前类型检查与对应代码生成是不完整的,需要加强以避免隐患和拓展语言能力。
- 冗余的VM指令类与优化: Snow VM为不同数据类型的每个操作都定义了独立的Command类(如IAddCommand, LAddCommand逻辑几乎相同,区别只是操作数类型不同)。这样做直观但导致代码大量重复,不利于维护和扩展。如果以后增加新类型(比如字节、短整型已经增加了BAdd, SAdd等),每种操作都要添加新类。可以考虑优化设计:采用模板泛型或继承来减少重复,实现一个参数化的AddCommand,根据类型参数决定执行int还是long运算;或者在Command.execute内部判断操作数类型并执行不同逻辑(但这会牺牲执行效率)。另一种方式是在指令集中引入统一操作码,比如把IAdd/ LAdd/ FAdd等的公共部分抽象为一个指令,加一个类型码参数。不过这改变较大。性能层面,当前VM每步都解析字符串opcode并查找Command,解析开销相对执行计算可忽略,但仍有提升空间:可以在加载字节码时预把指令字符串转换为更紧凑的中间形式(如指令对象列表或(opcode, args)对),这样execute循环无需split字符串,每步直接取对象执行。Java实现里,也可考虑将频繁的小方法(如各Command.execute)改为内联或减少虚函数调用,但这由JIT自动优化即可,不必手动干预。若追求极致,可将VM用Java代码重写为一个大的switch执行(类似解释器线性代码),省去很多间接调用,不过会失去面向对象的清晰结构。总之,在功能完善后,可视情况对VM执行部分进行优化。
- 主函数返回和程序结束: 当前约定以函数名"main"识别主程序,并在其返回时用HALT终止。这种做法隐含假设每个编译单元的主模块函数叫main,否则可能出现"进程未HALT"的情况。如前所述,若用户没有使用main而用其他名字当作入口函数,程序运行完RET后不会退出主循环(除非调用者frame为空时Engine检测到programCounter越界而退出,但这依赖隐含行为)。这需要改进。可以在编译阶段明确标记入口函数(例如第一个Module的第一个Function默认为入口),在代码生成时对该函数无论名称为何都使用HALT结尾。或者在Engine中,当返回到根帧时自动HALT。目前StackFrame.root的returnAddress通常是0,RetCommand弹出根帧后programCounter=0,会导致主循环重新从0执行,从而重复运行程序(可能Engine的PROGRAM_END常量用于避免此情况,但代码未展示完整逻辑)。实际测试可能因为callStack为空跳出,但不明显可靠。建议:增加对入口点的配置,或者发射特定指令标记程序结束。一种方案是在指令序列末尾无条件附加HALT确保结束。总之,需要让编译器输出的代码在任何情况下都能正常结束,而不是依赖名称约定。
- 模块与导入的扩展:目前Import只是检查被导入模块是否存在于同一编译单元。没有支持真正多文件编译:例如import一个外部模块并自动加载编译。目前用户如果想引用另一模块的函数,需将两个模块源文件一起提供给编译器(或者手动把被导入模块代码也合并)。理想改进是在编译时检测到Import未解析,可以根据模块名去文件系统查找对应源码或已编译文件并编译/链接。这需要编译器具备分步编译和符号链接功能。目前pom.xml等未显示具体实现,因此猜测暂未支持。未来可引入例如编译输出中间二进制格式(类似JVM的class文件),在导入时加载这些中间码,从而实现真正的模块化编译。此外,对于重复导入、循环导入等情况,也需要完善检测避免无限递归加载等问题。
- 调试和工具:项目虽然有LoggingUtils、VMStateLogger等,但目前命令行接口简单,缺少调试器、单步执行等功能。可以考虑实现**断点调试**模式:比如利用VMMode.STEP逐条执行指令并输出状态,从而手动调试程序。也可以开发AST/IR可视化工具,将编译过程的各阶段结果呈现给开发者,有助于理解和优化。另一个改进方向是提供**更加友好的错误信息**:目前错误通常只有简短描述和Token文本,最好能包含行号列号(Token已记录位置信息,但SemanticError等输出未利用这些【17†output】)。完善错误提示有助于使用者调试源代码。
- 性能优化:对于复杂程序,编译性能和运行性能都可以优化。编译性能方面,可优化语法分析算法如引入预测分析减少回溯(当前输入规模小影响不大)。运行性能方面,可尝试JIT:将热点字节码翻译为宿主机器指令执行,或者直接目标生成原生代码而非解释,这样运行速度会有数量级提升。不过这将显著增加实现复杂度,不一定是项目目标。另外,考虑利用Java本身的性能:例如把Snow源码直接编译为Java字节码(利用ASM等库动态生成类),然后在JVM中加载执行,可省去自己写的解释器开销。这相当于把Snow作为一个DSL直接构建到JVM上,也是可探索的优化方向。
- 代码结构与工程: 随着功能扩展,目前代码组织上也可做一些重构提高可维护性。例如,虚拟机指令类众多,可以通过自动代码生成来避免手写重复代码(比如用模板文件生成各类型算术指令类)。又如,为每个模块内子包都写了个doc/README.md记录设计,很有益处,但如果文档内容与代码实现不符要及时更新。目前readme.md提到的一些设计(如未实现的IR优化阶段、基于缩进的语法等)和最终代码存在出入,需注意保持文档同步,以免误导开发者。项目pom.xml简单但可以考虑拆分编译器和VM为独立模块,便于单独测试和重用。最后,更多

的测试用例和示例也是必要的。当前仓库只附带了一两个示例模块(CommonTasks等),应增加单元测试,覆盖各语法和功能点,确保改进不引入回归。

综上,SCompiler在实现了基本的编译运行流程的同时,也暴露出一些典型问题。幸运的是,其架构清晰模块分明,大部分改进可以各自独立进行。例如,可先完善控制流IR/指令这一大功能,再逐步强化类型系统,最后针对性能和架构优化。通过循序迭代,项目有望从教学原型走向一个功能更健全、性能更高的自定义语言编译器。

后续优化方向

结合以上分析,SCompiler项目后续可从以下几个方面着手优化和扩展:

- 丰富语言特性:实现条件分支和循环等控制结构的完整支持,增加如 break/continue 控制,支持函数递归和更复杂的表达式组合。同时引入更多数据类型(char、数组、结构体等)以提高语言表达能力。可以考虑增加字符串拼接、逻辑运算符(&& , ||)等语法糖,以及简单的输入输出API,提升Snow语言实用性。
- 加强类型系统: 完善静态类型检查,支持基本类型的自动转换规则(例如int可用在float运算中,通过插入I2F指令)和类型兼容性检测。引入 布尔类型并强制条件判断必须为布尔表达式,而不是用数字非零即真。后续甚至可增加面向对象特性,引入类和对象、方法调度等,但这需 要更大的架构调整。
- 中间代码优化:实现IR层的优化模块,例如常量折叠(编译期计算常量表达式)、死代码消除(移除永不执行的代码)、公共子表达式消除等,以减少生成指令数量。还可引入简单的寄存器分配优化:目前RegisterAllocator映射后很多槽可能未被充分利用,可通过分析寄存器活跃区间来复用槽位,降低局部变量空间占用,这对减少栈桢大小有益。此外,可考虑在IR级别进行函数内联、循环展开等更高级优化,不过需权衡编译时间和收益。
- 改进编译性能:对于大型源码,解析和语义分析阶段可以通过更好的算法提速。例如使用预生成的预测分析器(如ANTLR之类的工具)替代手写Parser,可自动处理错误同步和大部分繁琐逻辑。不过手写解析器在这个项目规模下尚可接受。编译性能另一点在于IO和中间过程打印,尽量减少不必要的调试输出(当前默认打印IR和字节码,可在非调试模式下关闭)。如果需要批量编译多个文件,可考虑实现增量编译或多线程并行编译不同模块。
- 增强虚拟机性能:如前所述,引入字节码预解析或直接JIT。可以尝试实现一个字节码到Java字节码的转换器,把Snow程序动态转为等价的 Java类加载执行,借助JVM优化获得性能提升。此外,若保持解释执行,可优化指令调度,比如使用**直跳表**(switch-case)替代Command 对象调用,以降低每步的分派开销。对于热点函数,可以探讨**动态编译**:在运行时收集函数执行次数,对高频函数生成优化机器码替换解释 执行。虽然对这个教学编译器而言可能有些超前,但这些都是虚拟机优化的经典方向。
- 工具链和调试支持:开发交互式REPL,允许用户逐句输入Snow代码立即编译执行,便于调试和教学演示。增强错误信息,包含源代码片段、高亮出错位置等。构建更完善的测试集,涵盖各语法元素组合和边界情况,以防止将来修改引入回归。可以编写**性能基准测试**,评估不同规模程序编译和执行耗时,为优化提供依据。考虑编写**文档和教程**,包括语言规范、使用指南等,降低他人上手门槛。
- 模块化和生态:实现真正的多文件模块编译机制,支持将模块编译为中间文件并在别的程序中链接使用。可以设计自己的二进制格式(比如 .snb 文件,包含字节码和元数据),由虚拟机直接加载执行或静态链接多个模块。随着功能增强,可以编写一些标准库模块,提供常用功能(字符串处理、数学运算等),以丰富Snow语言生态。
- 代码质量和维护:在重构方面,可考虑精简重复代码,例如利用元编程自动生成各类型指令类的模板代码,或为相似逻辑的解析器/分析器编写共用辅助函数。引入静态代码检查工具和格式化工具,保证编码规范统一。由于项目采用了Java 17的新特性(如record,模式匹配等),继续关注这些语言特性在本项目中的应用,可以让代码更简洁。例如,可以使用Java的 sealed 类来限制IRInstruction子类,明确哪几种指令类型,利于优化switch分支。

总之,SCompiler作为一个**编译器**,已经具备了完整的前后端链路和一个可运行的虚拟机,架构清晰,模块划分合理。在后续迭代中,应优先补足当前未完成的功能(控制流、类型检查),然后着力提升生成代码质量和执行性能。在确保正确性的基础上,引入必要的优化,使其不仅能"编译运行",还能"高效运行"。随着功能完善,还可以尝试支持更复杂的语言特性,甚至借鉴成熟语言设计,使Snow语言更加实用有趣。通过以上改进,SCompiler项目将从一个教学原型逐步走向一个有实用价值的定制语言编译器,在实践中检验并展示编译器构造的奥妙。