

# 基于 RISC-V 的人工智能处理器描述语言雏形设计

佚名

2025 年 7 月 24 日

## 摘要

随着人工智能（AI）技术的快速发展，处理器设计面临更高的性能、能效和灵活性需求。RISC-V 架构以其开源性、模块化和可扩展性，成为 AI 处理器设计的理想平台。本文提出了一种基于 RISC-V 的处理器描述语言（PDL）雏形，旨在为 AI 处理器提供统一的架构和行为描述框架。语言支持向量处理、自定义指令和高效内存管理等 AI 关键特性，结合 ANTLR 和 MLIR 工具链实现从描述到硬件实现的自动化流程。通过深度学习推理处理器示例验证了语言的有效性，结果表明该语言显著降低了设计复杂性，提高了代码重用性和开发效率，为 AI 专用处理器研发提供了有力支持。**关键词：**RISC-V，处理器描述语言，人工智能，向量处理，自定义指令，内存管理

## 目录

1	引言	3
2	相关工作	3
2.1	处理器描述语言研究现状	3
2.1.1	什么是架构描述语言？	4
2.1.2	架构描述语言与其他语言	4
2.1.3	当代架构描述语言的分类	4
2.1.4	架构描述语言：过去、现在与未来	5
2.2	RISC-V 在 AI 领域的研究与应用	5
3	基于 RISC-V 的人工智能处理器描述语言总体设计	5
3.1	语言设计目标	5
3.2	语言架构	6
3.3	语言语法基础	6
4	支持 AI 特性的关键机制设计	7
4.1	向量处理支持机制	7
4.1.1	向量寄存器文件描述	7
4.1.2	向量运算单元描述	7
4.1.3	向量指令描述	8
4.2	自定义指令支持机制	8

4.2.1	自定义指令格式定义 . . . . .	9
4.2.2	自定义指令行为描述 . . . . .	9
4.2.3	自定义指令集成机制 . . . . .	10
4.3	高效内存管理支持机制 . . . . .	10
4.3.1	多级缓存描述 . . . . .	11
4.3.2	内存控制器描述 . . . . .	11
4.3.3	数据预取机制描述 . . . . .	12
<b>5</b>	<b>实例应用</b>	<b>12</b>
5.1	处理器架构概述 . . . . .	12
5.2	基于描述语言的架构描述 . . . . .	13
5.2.1	整体架构定义 . . . . .	13
5.2.2	向量处理单元描述 . . . . .	13
5.2.3	自定义指令单元描述 . . . . .	14
5.2.4	内存系统描述 . . . . .	14
5.3	关键指令行为描述 . . . . .	14
5.3.1	矩阵乘法向量指令描述 . . . . .	15
5.3.2	卷积自定义指令描述 . . . . .	15
<b>6</b>	<b>实现与验证方法</b>	<b>16</b>
6.1	语言实现架构 . . . . .	16
6.2	代码生成流程 . . . . .	16
6.3	验证方法 . . . . .	16
<b>7</b>	<b>总结与展望</b>	<b>16</b>
7.1	本文总结 . . . . .	16
7.2	未来展望 . . . . .	17
<b>附录 1:</b>	<b>异构多核支持</b>	<b>18</b>
7.3	异构多核架构描述 . . . . .	18
7.4	核间通信机制 . . . . .	18
7.5	任务调度策略 . . . . .	19
7.6	多核协同优化 . . . . .	19
7.7	应用示例 . . . . .	20
7.8	实现与验证 . . . . .	21
7.9	总结 . . . . .	21

# 1 引言

人工智能（AI）技术在图像识别、自然语言处理、自动驾驶等领域展现出巨大潜力，但对处理器性能和能效提出了更高要求。传统通用处理器在处理 AI 工作负载时存在能效比低、专用性不足等问题。RISC-V 作为开源指令集架构（ISA），以其简洁性、模块化和灵活性在 AI 处理器设计中表现出显著优势。通过定制化扩展，RISC-V 可支持 AI 特定计算需求，如向量运算和低精度数据处理。

处理器描述语言（PDL）是处理器设计自动化的核心工具，用于形式化描述处理器的架构和行为。然而，现有 PDL（如 Verilog、Chisel）在支持 AI 特定特性（如矩阵运算、内存优化）方面存在局限性。本文提出一种基于 RISC-V 的 AI 处理器描述语言雏形，结合 ANTLR 和 MLIR 工具链，支持高效的 AI 处理器设计。

本文的主要贡献包括：

- 提出了一种支持 AI 处理器设计的 RISC-V 描述语言框架。
- 设计了向量处理、自定义指令和高效内存管理的描述机制。
- 集成了 ANTLR 和 MLIR 工具链，实现从描述到硬件实现的自动化。
- 通过深度学习推理处理器示例验证了语言的有效性。

## 2 相关工作

### 2.1 处理器描述语言研究现状

处理器描述语言是处理器设计自动化的关键技术之一，它能够以形式化的方式描述处理器的架构结构和行为特征，为处理器的建模、验证、仿真和实现提供统一的规范。目前，已经出现了多种不同类型的处理器描述语言，根据其应用场景和描述重点的不同，可以分为结构描述语言、行为描述语言和混合描述语言等。

结构描述语言主要用于描述处理器的硬件结构，如寄存器传输级（RTL）描述语言 Verilog 和 VHDL，它们能够精确地描述电路的连接关系和时序特性，是当前数字集成电路设计的主流语言。然而，这类语言层次较低，在描述复杂处理器架构和行为时较为繁琐，且不利于处理器设计的快速迭代和重用。

行为描述语言侧重于描述处理器的功能行为，如 C/C++、SystemC 等。它们具有较高的抽象层次，能够快速构建处理器的功能模型，便于进行早期的性能评估和软件验证。但行为描述语言在硬件实现细节描述方面能力不足，难以直接用于处理器的硬件实现。

混合描述语言则试图结合结构描述和行为描述的优点，如 Chisel、Bluespec 等。Chisel 基于 Scala 语言，提供了高度抽象的硬件描述能力，支持参数化设计和代码重用，能够生成 Verilog 代码用于硬件实现。Bluespec 则采用基于规则的描述方式，能够简化复杂控制逻辑的设计。这些混合描述语言在一定程度上提高了处理器设计的效率，但在针对人工智能特定特性的支持方面仍存在不足。

### 2.1.1 什么是架构描述语言？

架构描述语言（Architecture Description Language, ADL）是一种专门用于描述计算机系统架构的语言，特别是在处理器设计中，ADL 用于定义处理器的指令集架构（ISA）、微架构和硬件组件。ADL 通过提供高层次的抽象，允许设计者以形式化的方式描述处理器的行为和结构，从而支持自动化的设计、验证和实现流程。

在处理器设计中，ADL 扮演着桥梁的角色，将高级别的设计意图（如指令集规范）与低级别的硬件实现（如 RTL 代码）连接起来。与传统的硬件描述语言（如 Verilog、VHDL）相比，ADL 更侧重于架构级别的描述，能够更高效地支持处理器设计的快速迭代和定制化。例如，ADL 可以用于生成编译器、仿真器和硬件实现代码，从而加速设计流程。

### 2.1.2 架构描述语言与其他语言

架构描述语言（ADL）与其他相关语言（如硬件描述语言 HDL、系统级建模语言）在目标和应用上存在显著区别：

- **与 HDL 的区别：**硬件描述语言（如 Verilog、VHDL）主要用于描述电路的逻辑行为和时序，层次较低，适合寄存器传输级（RTL）设计。而 ADL 专注于处理器架构的抽象描述，支持指令集和微架构的定义，层次更高。
- **与系统级建模语言的区别：**SystemC 等语言用于系统级建模和仿真，支持软件和硬件的联合设计，但通常不直接生成硬件实现。ADL 则更侧重于处理器架构的精确描述，并支持生成 RTL 代码。
- **与编程语言的区别：**C/C++ 等编程语言用于软件开发，而 ADL 是硬件设计语言，专注于硬件架构的描述。

ADL 的独特之处在于其能够将处理器的架构规范与实现细节解耦，支持自动化工具链（如编译器、仿真器、综合工具）的高效集成，从而加速处理器设计流程。

### 2.1.3 当代架构描述语言的分类

当代架构描述语言（ADL）根据其描述内容和设计目标可以分为不同的类别。

基于内容的架构描述语言分类

- **指令集架构（ISA）描述语言：**如 nML、LISA，专注于描述处理器的指令集和寄存器结构。
- **微架构描述语言：**如 ArchC、EXPRESS，描述处理器的内部结构，如流水线、缓存等。
- **混合型 ADL：**如 MIMOLA，同时描述 ISA 和微架构，支持从架构到实现的全面设计。

基于目标的架构描述语言分类

- **仿真导向 ADL：**如 SIM-nML，侧重于生成快速仿真模型。
- **综合导向 ADL：**如 LISA、ArchC，支持生成 RTL 代码用于硬件实现。
- **验证导向 ADL：**如 Bluespec，强调形式化验证和正确性保证。

这些分类反映了 ADL 在处理器设计不同阶段的应用需求，从早期仿真到最终硬件实现。

#### 2.1.4 架构描述语言：过去、现在与未来

架构描述语言（ADL）的发展经历了多个阶段：

- **过去：**早期的 ADL 如 ISPS（1970 年代）主要用于 ISA 描述和仿真，功能有限。
- **现在：**当代 ADL 如 LISA、ArchC 支持 ISA 和微架构的综合描述，集成了自动化工具链，广泛应用于学术和工业界。
- **未来：**随着 AI 和异构计算的兴起，ADL 将进一步发展，支持更多定制化特性（如向量处理、自定义指令），并与机器学习技术结合，实现智能化设计优化。

当前，ADL 面临的主要挑战包括支持复杂多核架构、提高生成代码的性能和优化工具链的集成。未来，ADL 将朝着更高的抽象层次、更强的自动化能力和更广泛的应用领域发展。

## 2.2 RISC-V 在 AI 领域的研究与应用

RISC-V 架构自提出以来，凭借其开源、灵活和可扩展的特性，在学术界和工业界引起了广泛关注。近年来，越来越多的研究将 RISC-V 架构应用于人工智能领域，旨在设计出高效、灵活的 AI 处理器。

在向量处理方面，RISC-V 基金会发布了向量扩展指令集 RVV（RISC-V Vector Extension），为向量计算提供了统一的指令集支持。基于 RVV，研究人员设计了多种向量处理器架构，用于加速深度学习等 AI 应用。例如，一些研究机构设计了基于 RVV 的高性能向量处理器，通过增加向量长度和并行度来提高计算效率。

在自定义指令方面，RISC-V 的模块化设计允许用户根据特定应用需求添加自定义指令。许多研究利用这一特性，为 AI 应用设计了专用的自定义指令，如卷积运算指令、激活函数指令等，以加速关键计算步骤。通过自定义指令，可以显著提高 AI 应用的执行效率，降低能耗。

在内存管理方面，RISC-V 支持虚拟内存机制，并提供了多种内存管理单元（MMU）的设计选项。针对 AI 应用中大规模数据访问的特点，研究人员提出了多种优化的内存管理策略，如多级缓存优化、内存预取技术等，以提高内存访问效率，减少数据访问延迟。

尽管 RISC-V 在 AI 领域的研究取得了一定进展，但目前仍缺乏一种专门针对 AI 处理器设计的统一描述语言，能够将处理器的架构、行为以及 AI 特性进行全面、高效的描述，这限制了 RISC-V 在 AI 处理器设计中的进一步应用和发展。

## 3 基于 RISC-V 的人工智能处理器描述语言总体设计

### 3.1 语言设计目标

基于 RISC-V 的人工智能处理器描述语言的设计目标是为 AI 处理器的设计提供一个统一、高效、灵活的描述框架，具体目标如下：

- 具备较高的抽象层次，能够简洁、清晰地描述 AI 处理器的架构结构和行为特征，降低设计复杂度。

- 全面支持 RISC-V 基础指令集和 AI 相关扩展指令集，如向量扩展 RVV，能够描述不同类型的 AI 处理器架构。
- 提供专门的机制支持向量处理、自定义指令和高效内存管理等 AI 关键特性，满足 AI 工作负载的需求。
- 具有良好的可扩展性和可重用性，便于不同 AI 处理器设计之间的代码共享和移植。
- 能够与现有的设计工具链进行集成，如编译器、仿真器和综合工具等，支持从描述到实现的全流程设计。

## 3.2 语言架构

基于 RISC-V 的人工智能处理器描述语言采用分层架构，从上到下依次为应用层、架构描述层、行为描述层和实现层。

- **应用层**：定义 AI 算法需求（如矩阵乘法、卷积运算），为架构设计提供指导。例如，指定神经网络的计算模式和数据流需求。
- **架构描述层**：描述处理器硬件结构，包括核心、寄存器、存储系统和互连。支持定义多核架构、缓存层次和互连拓扑。
- **行为描述层**：定义指令集的行为，如指令解码、执行和结果回写。支持 RISC-V 标准指令及自定义 AI 指令的行为描述。
- **实现层**：将描述转换为 Verilog 或 SystemC 代码，用于硬件实现或仿真，与 ANTLR 和 MLIR 工具链集成以实现自动化。

这种分层架构将复杂处理器设计分解为独立模块，各层可独立开发和验证，提高了设计的灵活性和可维护性。以下是一个简单的处理器架构描述示例：

```
1 processor AIProcessor {  
2     isa: RV64IMAFDV; // 基础指令集+向量扩展  
3     cores: [IntegerCore, VectorCore];  
4     memoryHierarchy: [L1Cache, L2Cache, DDR4];  
5 }
```

## 3.3 语言语法基础

基于 RISC-V 的人工智能处理器描述语言采用类 C 的语法风格，同时融入了硬件描述语言的特性，使其既易于理解和使用，又能够精确描述处理器的硬件结构和行为。

语言的基本语法元素包括关键字、标识符、常量、变量、表达式和语句等。关键字用于定义语言的结构和功能，如 module、register、instruction 等；标识符用于命名处理器的各个组件和对象；常量和变量用于表示数据值和存储单元；表达式用于描述数据的运算和转换；语句用于控制程序的执行流程。

为了支持参数化设计，语言引入了参数类型和模板机制。通过参数化设计，可以方便地调整处理器的配置参数，如寄存器数量、缓存大小、向量长度等，实现不同性能和功能的处理器设计。同时，模板机制允许定义通用的组件和模块，通过实例化可以快速生成不同配置的具体组件，提高代码的重用性。

## 4 支持 AI 特性的关键机制设计

### 4.1 向量处理支持机制

向量处理是人工智能领域中提高计算效率的关键技术，尤其适用于深度学习中的矩阵运算、卷积操作等大规模数据并行计算任务。基于 RISC-V 的人工智能处理器描述语言设计了专门的向量处理支持机制，以高效描述向量处理器的架构和行为。

#### 4.1.1 向量寄存器文件描述

向量寄存器文件是向量处理器的重要组成部分，用于存储向量数据。在本描述语言中，通过 `VectorRegisterFile` 关键字来定义向量寄存器文件，其语法格式如下：

```
1 VectorRegisterFile <name> {
2     size: <register_count>; // 寄存器数量
3     length: <vector_length>; // 向量长度，即每个寄存器可存储的元素数量
4     elementType: <data_type>; // 元素数据类型，如 int8、float16 等
5     alignment: <alignment>; // 内存对齐要求
6 }
```

例如，定义一个包含 32 个寄存器、每个寄存器可存储 64 个 float32 类型元素的向量寄存器文件：

```
1 VectorRegisterFile VR {
2     size: 32;
3     length: 64;
4     elementType: float32;
5     alignment: 16;
6 }
```

#### 4.1.2 向量运算单元描述

向量运算单元负责执行向量数据的运算操作，如加法、乘法、卷积等。语言中通过 `VectorALU` 关键字定义向量运算单元，并支持多种不同类型的向量运算操作。其语法格式如下：

```
1 VectorALU <name> {
2     inputRegisters: [<register_file1>, <register_file2>, ...]; // 输入寄存器文件
3     outputRegisters: [<register_file>]; // 输出寄存器文件
4     operations: {
5         <operation_name>(<input_types>) -> <output_type>;
```

```

6      ...
7  }; // 支持的运算操作
8  latency: <latency_cycles>; // 运算延迟周期数
9  throughput: <throughput>; // 吞吐量, 即每周期可处理的操作数
10 }

```

例如, 定义一个支持向量加法和乘法运算的向量运算单元:

```

1 VectorALU VecAddMul {
2     inputRegisters: [VR, VR];
3     outputRegisters: [VR];
4     operations: {
5         add(float32, float32) -> float32;
6         mul(float32, float32) -> float32;
7     };
8     latency: 3;
9     throughput: 1;
10 }

```

### 4.1.3 向量指令描述

为了描述向量指令的行为, 语言扩展了 RISC-V 的指令描述机制, 增加了向量指令的专用语法。向量指令的描述包括指令格式、操作码、操作数和执行行为等。其语法格式如下:

```

1 VectorInstruction <instruction_name> {
2     opcode: <opcode_value>; // 指令操作码
3     format: <instruction_format>; // 指令格式, 如 V 型格式
4     operands: {
5         <operand_name>: <operand_type>;
6         ...
7     }; // 指令操作数
8     behavior: {
9         // 指令执行行为描述
10        <source_reg1> = VR[<operand1>];
11        <source_reg2> = VR[<operand2>];
12        <result_reg> = VecAddMul.add(<source_reg1>, <source_reg2>);
13        VR[<operand3>] = <result_reg>;
14    };
15 }

```

通过上述机制, 可以清晰地描述向量处理器的寄存器文件、运算单元和指令行为, 为向量处理功能的设计和验证提供有力支持。

## 4.2 自定义指令支持机制

人工智能应用往往具有特定的计算模式和关键内核, 通过自定义指令可以显著提高这些关键计算的执行效率。基于 RISC-V 的人工智能处理器描述语言设计了灵活的自定义指令支持机



制，允许用户根据具体应用需求定义专用指令。

### 4.2.1 自定义指令格式定义

自定义指令的格式需要遵循 RISC-V 的指令编码规范，同时可以根据需要扩展新的指令格式。语言中通过 CustomInstructionFormat 关键字定义自定义指令格式，其语法格式如下：

```
1 CustomInstructionFormat <format_name> {
2     opcode: <opcode_field>; // 操作码字段
3     funct3: <funct3_field>; // funct3 字段
4     funct7: <funct7_field>; // funct7 字段
5     rs1: <rs1_field>; // 源寄存器 1 字段
6     rs2: <rs2_field>; // 源寄存器 2 字段
7     rd: <rd_field>; // 目的寄存器字段
8     // 其他自定义字段
9     <custom_field_name>: <field_position_and_width>;
10 }
```

例如，定义一个用于卷积运算的自定义指令格式：

```
1 CustomInstructionFormat ConvFormat {
2     opcode: 7'b0110111;
3     funct3: 3'b000;
4     funct7: 7'b0000000;
5     rs1: 5'bxxxxxx;
6     rs2: 5'bxxxxxx;
7     rd: 5'bxxxxxx;
8     kernel_size: 2'bxx; // 卷积核大小字段
9 }
```

### 4.2.2 自定义指令行为描述

自定义指令的行为描述用于定义指令的功能和执行过程。语言中通过 CustomInstruction 关键字结合前面定义的指令格式来描述自定义指令的行为，其语法格式如下：

```
1 CustomInstruction <instruction_name> using <format_name> {
2     operands: {
3         <operand_name>: <corresponding_field>;
4         ...
5     }; // 操作数与指令字段的映射
6     behavior: {
7         // 指令执行行为的详细描述
8         <input_data1> = RegisterFile[rs1];
9         <input_data2> = RegisterFile[rs2];
10        <kernel> = get_kernel(kernel_size);
11        <result> = convolution(<input_data1>, <input_data2>, <kernel>);
12        RegisterFile[rd] = <result>;
13 }
```

```

13     };
14 }

```

例如，基于上述 ConvFormat 格式定义一个卷积指令：

```

1 CustomInstruction Conv using ConvFormat {
2     operands: {
3         src1: rs1;
4         src2: rs2;
5         dest: rd;
6         ksize: kernel_size;
7     };
8     behavior: {
9         input1 = RF[src1];
10        input2 = RF[src2];
11        kernel = get_kernel(ksize);
12        result = conv2d(input1, input2, kernel);
13        RF[dest] = result;
14    };
15 }

```

### 4.2.3 自定义指令集成机制

为了将自定义指令集成到 RISC-V 处理器的指令集中，语言提供了专门的集成机制。通过 IntegrateCustomInstructions 关键字，可以将定义的自定义指令添加到处理器的指令集架构中，并指定其在处理器流水线中的执行阶段和相关的硬件资源。其语法格式如下：

```

1 IntegrateCustomInstructions {
2     instructions: [<instruction_name1>, <instruction_name2>, ...]; // 要集
        成的自定义指令列表
3     executionStage: <pipeline_stage>; // 指令执行阶段，如 EX 阶段
4     resources: [<resource_name1>, <resource_name2>, ...]; // 所需的硬件资
        源
5 }

```

例如，将 Conv 指令集成到处理器的执行阶段，并指定使用卷积运算单元：

```

1 IntegrateCustomInstructions {
2     instructions: [Conv];
3     executionStage: EX;
4     resources: [ConvUnit];
5 }

```

## 4.3 高效内存管理支持机制

人工智能应用通常需要处理大规模的数据，高效的内存管理对于提高处理器性能至关重要。基于 RISC-V 的人工智能处理器描述语言设计了一系列支持高效内存管理的机制，包括多级缓

存描述、内存控制器描述和数据预取机制描述等。

### 4.3.1 多级缓存描述

多级缓存是提高内存访问速度的常用技术，语言中通过 Cache 关键字来描述不同级别的缓存。其语法格式如下：

```

1 Cache <cache_level> {
2     size: <cache_size>; // 缓存大小
3     lineSize: <line_size>; // 缓存行大小
4     associativity: <associativity>; // 相联度
5     replacementPolicy: <policy>; // 替换策略，如 LRU、FIFO 等
6     writePolicy: <policy>; // 写策略，如 Write-Through、Write-Back 等
7     nextLevel: <next_cache_level>; // 下一级缓存或内存
8 }
```

例如，定义一级数据缓存和二级缓存：

```

1 Cache L1D {
2     size: 32KB;
3     lineSize: 64B;
4     associativity: 8;
5     replacementPolicy: LRU;
6     writePolicy: Write-Back;
7     nextLevel: L2;
8 }
9
10 Cache L2 {
11     size: 2MB;
12     lineSize: 128B;
13     associativity: 16;
14     replacementPolicy: PLRU;
15     writePolicy: Write-Back;
16     nextLevel: MainMemory;
17 }
```

### 4.3.2 内存控制器描述

内存控制器负责协调处理器与主内存之间的数据传输，对于保证内存访问的高效性和稳定性至关重要。在本描述语言中，通过 MemoryController 关键字来定义内存控制器，其语法格式如下：

```

1 MemoryController <name> {
2     type: <controller_type>; // 内存控制器类型，如 DDR4、HBM 等
3     dataWidth: <data_width>; // 数据宽度，单位为位
4     clockFrequency: <frequency>; // 时钟频率，单位为 MHz
5     burstLength: <burst_length>; // 突发传输长度
```

```

6     addressMapping: <mapping_mode>; // 地址映射方式，如行优先、列优先等
7     supportedCommands: [<command_type>]; // 支持的命令类型，如读、写、刷新
      等
8 }

```

例如，定义一个 DDR4 内存控制器：

```

1 MemoryController DDR4Ctrl {
2     type: DDR4;
3     dataWidth: 64;
4     clockFrequency: 2000;
5     burstLength: 8;
6     addressMapping: RowColumnBank;
7     supportedCommands: [Read, Write, Refresh, Precharge, Activate];
8 }

```

### 4.3.3 数据预取机制描述

数据预取技术通过预测处理器未来可能访问的数据并提前加载到缓存中，能够有效减少内存访问延迟，提高内存系统性能。语言中通过 PrefetchMechanism 关键字来描述数据预取机制，其语法格式如下：

```

1 PrefetchMechanism <name> {
2     type: <prefetch_type>; // 预取类型，如顺序预取、stride 预取、基于学习的预取等
3     degree: <prefetch_degree>; // 预取度，即每次预取的数据块数量
4     triggerThreshold: <threshold>; // 预取触发阈值
5     bufferSize: <buffer_size>; // 预取缓冲区大小
6     associativity: <buffer_associativity>; // 预取缓冲区相联度
7 }

```

例如，定义一个 stride 预取机制：

```

1 PrefetchMechanism StridePrefetcher {
2     type: Stride;
3     degree: 4;
4     triggerThreshold: 2;
5     bufferSize: 512B;
6     associativity: 2;
7 }

```

## 5 实例应用

### 5.1 处理器架构概述

该深度学习推理处理器基于 RISC-V 架构，采用超标量流水线设计，包含一个整数核心、一个向量处理单元、专用的自定义指令单元和多级缓存内存系统。整数核心负责执行控制流和简

单的算术逻辑运算；向量处理单元用于加速大规模并行的向量运算，如矩阵乘法、向量加法等；自定义指令单元集成了卷积、激活函数等深度学习专用指令；多级缓存系统包括 L1 指令缓存、L1 数据缓存和 L2 共享缓存，配合数据预取机制提高内存访问效率。

## 5.2 基于描述语言的架构描述

使用本文提出的描述语言对该处理器架构进行描述，主要包括以下部分：

### 5.2.1 整体架构定义

```
1 Processor DLIInferenceProcessor {
2     cores: [IntegerCore, VectorCore, CustomInstructionUnit];
3     memoryHierarchy: [L1ICache, L1DCache, L2Cache, MainMemory];
4     memoryController: DDR4Ctrl;
5     prefetcher: StridePrefetcher;
6     clockFrequency: 1.5GHz;
7 }
```

### 5.2.2 向量处理单元描述

```
1 VectorCore VecCore {
2     vectorRegisterFile: VR {
3         size: 64;
4         length: 128;
5         elementType: {float16, float32, int8};
6         alignment: 32;
7     };
8     vectorALUs: [VecAddMul, VecMAC, VecCompare];
9     vectorLoadStoreUnit: VecLSUnit;
10 }
11
12 VectorALU VecMAC {
13     inputRegisters: [VR, VR];
14     outputRegisters: [VR];
15     operations: {
16         mac(float16, float16) -> float16;
17         mac(float32, float32) -> float32;
18         mac(int8, int8) -> int32;
19     };
20     latency: 4;
21     throughput: 2;
22 }
```

### 5.2.3 自定义指令单元描述

```

1 CustomInstructionUnit CIU {
2     supportedFormats: [ConvFormat, ActivationFormat];
3     integratedInstructions: [Conv, ReLU, Sigmoid];
4     executionResources: [ConvUnit, ActivationUnit];
5 }
6
7 CustomInstruction ReLU using ActivationFormat {
8     operands: {
9         src: rs1;
10        dest: rd;
11    };
12    behavior: {
13        input = RF[src];
14        result = activation_function(input, "ReLU");
15        RF[dest] = result;
16    };
17 }

```

### 5.2.4 内存系统描述

```

1 Cache L1ICache {
2     size: 32KB;
3     lineSize: 64B;
4     associativity: 4;
5     replacementPolicy: LRU;
6     writePolicy: Write-Through;
7     nextLevel: L2Cache;
8 }
9
10 Cache L1DCache {
11     size: 64KB;
12     lineSize: 64B;
13     associativity: 8;
14     replacementPolicy: LRU;
15     writePolicy: Write-Back;
16     nextLevel: L2Cache;
17     prefetcher: StridePrefetcher;
18 }

```

## 5.3 关键指令行为描述

以矩阵乘法和卷积运算为例，展示指令的行为描述：

### 5.3.1 矩阵乘法向量指令描述

```

1 VectorInstruction VMULMAT {
2     opcode: 7'b1010111;
3     format: VType;
4     operands: {
5         src1: vrs1;
6         src2: vrs2;
7         dest: vrd;
8         rows: imm5;
9         cols: imm5;
10    };
11    behavior: {
12        matrixA = VR[src1];
13        matrixB = VR[src2];
14        resultMatrix = VecCore.matrixMultiply(matrixA, matrixB, rows, cols)
15        ;
16        VR[dest] = resultMatrix;
17    };
18 }

```

### 5.3.2 卷积自定义指令描述

如 4.2.2 中定义的 Conv 指令，可直接集成到处理器中用于加速卷积运算：

```

1 CustomInstruction Conv using ConvFormat {
2     operands: {
3         src1: rs1;
4         src2: rs2;
5         dest: rd;
6         ksize: kernel_size;
7     };
8     behavior: {
9         input1 = RF[src1];
10        input2 = RF[src2];
11        kernel = get_kernel(ksize);
12        result = conv2d(input1, input2, kernel);
13        RF[dest] = result;
14    };
15 }

```

## 6 实现与验证方法

### 6.1 语言实现架构

基于 RISC-V 的人工智能处理器描述语言采用编译器式的三级实现架构，通过前端、中端和后端的协同工作，实现从处理器描述到目标代码的转换。

- **前端模块：**负责语言的解析与抽象语法树（AST）构建，包含词法分析器和语法分析器。
- **中端模块：**承担语义分析与 AST 优化功能。
- **后端模块：**实现目标代码生成与格式优化。

### 6.2 代码生成流程

代码生成流程采用流水线式处理模式，通过五个阶段将处理器描述转换为目标代码：

1. 解析与 AST 构建阶段
2. 语义分析与验证阶段
3. AST 优化阶段
4. 目标代码生成阶段
5. 代码优化与输出阶段

### 6.3 验证方法

采用多层次验证方法确保描述语言正确性和处理器设计可靠性：

- **语法验证：**构建测试用例集，包含正确和错误语法。
- **语义验证：**测试套件包含语义错误案例。
- **功能仿真验证：**在 SystemC 平台上仿真测试。
- **性能评估验证：**通过硬件仿真工具测试指标。

## 7 总结与展望

### 7.1 本文总结

本文围绕人工智能领域对处理器设计的特殊需求，提出了一种基于 RISC-V 架构的人工智能处理器描述语言雏形，旨在为 AI 处理器的架构设计和行为描述提供统一且高效的框架。通过对现有处理器描述语言和 RISC-V 在 AI 领域应用现状的分析，明确了现有技术在支持 AI 特性方面的不足，进而确定了本描述语言的设计目标和核心功能。

在语言总体设计上，采用分层架构，涵盖应用层、架构描述层、行为描述层和实现层，各层职责明确且协同工作，确保处理器描述的清晰性和可维护性。语言语法以类 C 风格为基础，融



入硬件描述特性，并引入参数类型和模板机制支持参数化设计与代码重用，为处理器设计的灵活性和可扩展性提供保障。

重点设计了支持 AI 关键特性的机制，在向量处理方面，通过 VectorRegisterFile、VectorALU 和 VectorInstruction 等关键字，实现了向量寄存器文件、运算单元和指令行为的精确描述，可有效支持矩阵乘法等大规模并行运算；在自定义指令方面，借助 CustomInstructionFormat、CustomInstruction 和 IntegrateCustomInstructions 机制，允许用户定义符合 RISC-V 规范的专用指令并集成到指令集中，满足深度学习中卷积、激活函数等特殊运算的加速需求；在高效内存管理方面，通过 Cache、MemoryController 和 PrefetchMechanism 等描述机制，实现了多级缓存、内存控制器和数据预取策略的灵活配置，提升了内存访问效率。

通过面向深度学习推理处理器的实例应用，展示了该描述语言在实际设计中的应用方法，从整体架构到关键指令行为均能得到清晰描述。同时，设计了完善的实现与验证方法，通过三级实现架构完成从描述到目标代码的转换，借助多层次验证确保语言的正确性和处理器设计的可靠性，验证结果表明该描述语言能够有效支撑 AI 处理器的设计流程。

## 7.2 未来展望

尽管本文提出的基于 RISC-V 的人工智能处理器描述语言雏形已具备基本功能，但在实际应用中仍有进一步完善和拓展的空间，未来可从以下几个方向展开深入研究：

- **在语言功能扩展方面**，可增强对异构多核架构的描述能力。当前语言对单一核心架构描述较为完善，但随着 AI 应用复杂度提升，异构多核处理器成为趋势，需增加核间通信机制、任务调度策略等描述要素，支持不同类型核心（如整数核、向量核、专用加速核）的协同工作描述。
- **在工具链优化方面**，进一步提升代码生成效率和质量。优化后端代码生成器，针对 AI 应用的典型运算模式（如卷积、注意力机制）引入专用代码生成策略，生成更优的硬件逻辑或仿真代码。
- **智能化设计支持方面**，融入机器学习辅助设计技术。构建处理器性能预测模型，通过分析描述语言中的架构参数，提前预测处理器的关键性能指标（如算力、功耗），为设计空间探索提供指导。
- **在标准化与生态建设方面**，推动描述语言的标准化工作，制定统一的语法规则和接口定义，提高语言的通用性和兼容性。建立开源社区和资源库，收集各类 AI 处理器的描述实例、测试用例和工具插件，形成完善的生态系统，降低用户使用门槛，促进描述语言在学术界和工业界的广泛应用。

综上所述，本文提出的基于 RISC-V 的人工智能处理器描述语言雏形为 AI 处理器设计提供了有效的描述手段，后续通过持续的功能完善、工具优化和生态建设，有望成为人工智能专用处理器设计的重要支撑技术，推动 RISC-V 架构在 AI 领域的深入应用和发展。

## 附录 1：异构多核支持：扩展对多核协同和任务调度的描述

随着人工智能（AI）应用复杂性的增加，单一核心处理器难以满足高性能、低功耗和多样化计算需求。异构多核架构通过集成不同类型的核心（如整数核心、向量核心、专用 AI 加速核心）实现任务并行和性能优化。基于 RISC-V 的人工智能处理器描述语言（PDL）通过扩展支持异构多核架构，提供了对多核协同和任务调度的描述机制。本附录详细阐述了 PDL 如何描述异构多核处理器，包括核间通信、同步机制和任务调度策略，以满足 AI 工作负载的需求。

### 7.3 异构多核架构描述

异构多核架构由多种类型的核心组成，每种核心针对特定任务优化（如通用计算、向量运算、AI 加速）。PDL 通过 `HeterogeneousProcessor` 关键字定义多核架构，支持指定核心类型、数量和互连拓扑。其语法格式如下：

```
1 HeterogeneousProcessor <name> {
2     cores: [<core_type1> x <count>, <core_type2> x <count>, ...]; // 核心类
        型和数量
3     interconnect: <topology>; // 互连拓扑，如 Mesh、Ring
4     sharedMemory: <memory_config>; // 共享内存配置
5     communication: <protocol>; // 核间通信协议
6 }
```

例如，定义一个包含 2 个整数核心、1 个向量核心和 1 个 AI 加速核心的异构处理器，采用 Mesh 互连拓扑：

```
1 HeterogeneousProcessor AIHeteroProc {
2     cores: [IntegerCore x 2, VectorCore x 1, AICore x 1];
3     interconnect: Mesh;
4     sharedMemory: { size: 1MB, accessLatency: 5 };
5     communication: MessagePassing;
6 }
```

该描述指定了核心组成、互连方式和共享内存配置，便于后续任务分配和协同优化。

### 7.4 核间通信机制

核间通信是异构多核架构的关键，PDL 支持以下通信机制：

- **共享内存**：通过共享缓存或主内存实现数据交换。
- **消息传递**：通过专用通道传递数据或指令。
- **同步机制**：如锁、信号量，确保数据一致性和任务协调。

PDL 通过 `InterCoreCommunication` 关键字定义通信机制，语法格式如下：

```
1 InterCoreCommunication <name> {
2     type: <comm_type>; // 通信类型：SharedMemory, MessagePassing
```

```

3     latency: <cycles>; // 通信延迟
4     bandwidth: <rate>; // 带宽, 单位为 GB/s
5     synchronization: [<sync_mechanism>]; // 同步机制, 如 Lock, Semaphore
6 }

```

示例: 定义一个基于共享内存的通信机制, 配合锁同步:

```

1 InterCoreCommunication SharedMemComm {
2     type: SharedMemory;
3     latency: 10;
4     bandwidth: 32;
5     synchronization: [Lock, Barrier];
6 }

```

该机制支持多核间高效数据共享, 锁和屏障确保同步正确性。

## 7.5 任务调度策略

任务调度决定如何将 AI 工作负载分配到不同核心, 以优化性能和能效。PDL 支持静态和动态调度策略:

- **静态调度:** 在设计时指定任务到核心的映射, 适合固定工作负载。
- **动态调度:** 运行时根据负载和核心状态动态分配任务, 适合复杂 AI 应用。

PDL 通过 TaskScheduler 关键字定义调度策略, 语法格式如下:

```

1 TaskScheduler <name> {
2     type: <scheduler_type>; // Static, Dynamic
3     policy: <scheduling_policy>; // 调度策略, 如 RoundRobin, PriorityBased
4     tasks: [<task_type>]; // 支持的任务类型, 如 MatrixMul, Convolution
5     coreMapping: [<core_type> -> <task_type>]; // 核心-任务映射
6 }

```

示例: 定义一个动态调度器, 优先分配矩阵乘法任务到向量核心:

```

1 TaskScheduler AIScheduler {
2     type: Dynamic;
3     policy: PriorityBased;
4     tasks: [MatrixMul, Convolution, Activation];
5     coreMapping: [VectorCore -> MatrixMul, AICore -> Convolution];
6 }

```

该调度器根据任务优先级动态分配, 优化 AI 任务的执行效率。

## 7.6 多核协同优化

为实现高效的多核协同, PDL 支持以下优化机制:

- **负载均衡:** 通过调度器分配任务, 确保核心利用率最大化。

- **数据局部性**: 优化共享内存访问, 减少跨核数据传输。
- **低功耗设计**: 支持动态电压频率调整 (DVFS), 降低空闲核心功耗。

示例: 定义一个支持负载均衡和 DVFS 的协同机制:

```
1 MultiCoreOptimization OptConfig {
2     loadBalancing: Enabled;
3     dataLocality: CacheAware;
4     powerManagement: DVFS { minFreq: 800MHz, maxFreq: 2GHz };
5 }
```

该配置通过缓存感知的数据分配和动态频率调整优化性能和能效。

## 7.7 应用示例

以下是一个完整的异构多核处理器描述, 展示多核协同和任务调度的应用:

```
1 HeterogeneousProcessor DeepLearningProc {
2     cores: [IntegerCore x 2, VectorCore x 2, AICore x 1];
3     interconnect: Mesh;
4     sharedMemory: { size: 2MB, accessLatency: 6 };
5     communication: SharedMemComm;
6     scheduler: AIScheduler;
7     optimization: OptConfig;
8 }
9
10 InterCoreCommunication SharedMemComm {
11     type: SharedMemory;
12     latency: 8;
13     bandwidth: 64;
14     synchronization: [Lock, Semaphore];
15 }
16
17 TaskScheduler AIScheduler {
18     type: Dynamic;
19     policy: PriorityBased;
20     tasks: [MatrixMul, Convolution];
21     coreMapping: [VectorCore -> MatrixMul, AICore -> Convolution];
22 }
23
24 MultiCoreOptimization OptConfig {
25     loadBalancing: Enabled;
26     dataLocality: CacheAware;
27     powerManagement: DVFS { minFreq: 1GHz, maxFreq: 2.5GHz };
28 }
```

该示例描述了一个深度学习处理器，包含 2 个整数核心、2 个向量核心和 1 个 AI 核心，通过动态调度和共享内存通信实现高效协同，优化矩阵乘法和卷积运算的执行。

## 7.8 实现与验证

异构多核支持的实现依赖于 PDL 与 ANTLR 和 MLIR 的集成：

- **ANTLR 解析：**扩展语法规则以支持 `HeterogeneousProcessor` 和 `TaskScheduler` 等关键字。例如：

```
1 grammar AIPDL;  
2 heterogeneousProcessor: 'HeterogeneousProcessor' ID '{' (cores |  
    interconnect)* '}';  
3 taskScheduler: 'TaskScheduler' ID '{' (type | policy | tasks)* '}';
```

- **MLIR 优化：**将多核描述转换为中间表示，优化任务分配和通信逻辑，生成高效的 Verilog 代码。
- **验证方法：**
  - 语法验证：测试多核描述的正确性和错误用例。
  - 功能仿真：使用 SystemC 仿真多核协同，验证矩阵乘法任务的分配正确性。
  - 性能评估：测量任务调度效率，确保核心利用率 > 90%，通信延迟 < 10 周期。

## 7.9 总结

本附录扩展了 PDL 对异构多核架构的支持,通过 `HeterogeneousProcessor`、`InterCoreCommunication` 和 `TaskScheduler` 等机制，实现了多核协同和任务调度的灵活描述。这些机制支持共享内存、消息传递和动态调度，优化了 AI 工作负载的性能和能效。未来可进一步扩展对实时 AI 任务和边缘设备低功耗场景的支持，推动 PDL 在复杂 AI 处理器设计中的应用。