

Calebe Miranda Ferreira Braga de Castro

**Refinamento e Integração do Modelo LLM  
Compiler na Otimização de Código Alvo para  
Microcontroladores AVR na Robotics Language**

Jataí-GO

2025

Calebe Miranda Ferreira Braga de Castro

# **Refinamento e Integração do Modelo LLM Compiler na Otimização de Código Alvo para Microcontroladores AVR na Robotics Language**

Monografia apresentada ao curso de Ciência da Computação do Instituto de Ciências Exatas e Tecnológicas da Universidade Federal de Jataí (UFJ), como requisito para obtenção do título de Bacharel em Ciência da Computação.

Universidade Federal de Jataí

Orientador: Prof. Dr. Thiago Borges de Oliveira

Jataí-GO

2025

*Dedico este trabalho a todos que me apoiaram e acreditaram em mim durante esta jornada. Ao meu irmão e à minha família, que sempre me motivaram e incentivaram com amor e paciência, durante todo o processo. E também, ao professor e orientador Thiago Borges de Oliveira, que me apoiou e forneceu inestimáveis instruções.*

# Agradecimentos

*Gostaria de agradecer aos meus amigos e familiares que estiveram comigo durante todo o processo de realização deste trabalho, em especial ao estudante Alisson Ferreira Teles, que realizou um trabalho com tema semelhante e foi um grande parceiro durante todos os momentos difíceis desta caminhada. Agradeço à minha mãe, Lara Michelle Ferreira Braga, pelo amor incondicional e pelas palavras de motivação, assim como agradeço a todos os meus professores e avaliadores pela confiança para a realização deste trabalho. Muito Obrigado !*

*A ciência da computação não é sobre computadores. É sobre humanidade.*

*- Donald Knuth*

# Resumo

As bibliotecas padrão são um conjunto de bibliotecas predefinidas que acompanham a maioria das linguagens de programação, fornecendo funcionalidades essenciais para facilitar o desenvolvimento de software. Elas oferecem implementações de funções e rotinas que atendem a necessidades comuns, evitando que os desenvolvedores precisem reescrever código básico e complexo, o que aumenta a produtividade e a confiabilidade dos programas. Dentro das bibliotecas padrão, uma das mais utilizadas é a biblioteca matemática, que contém funções essenciais para realizar cálculos, além de oferecer operações como trigonometria, exponenciação, logaritmos, entre outras, frequentemente utilizadas em diversas áreas da ciência. Neste trabalho, investigam-se alternativas para a implementação de funções matemáticas na Robotics Language, desenvolvida na Universidade Federal de Jataí. Foram implementadas e avaliadas seis funções matemáticas: seno, cosseno, tangente, exponencial, logaritmo e raiz quadrada, utilizando algoritmos alternativos, como o CORDIC para as operações trigonométricas. Avaliou-se o tempo de execução e o tamanho do programa objeto para as plataformas AVR e STM32. Os resultados iniciais indicaram que os algoritmos são funcionalmente adequados e precisos, mas não competitivos em termos de tempo de execução e tamanho do programa objeto gerado, considerando o estágio atual de desenvolvimento do Robcmp e dos backends AVR e ARM do LLVM.

**Palavras-chaves:** Linguagem de programação; CORDIC; Bibliotecas padrão; Robcmp.

# Abstract

As bibliotecas padrão são um conjunto de bibliotecas predefinidas que acompanham a maioria das linguagens de programação, fornecendo funcionalidades essenciais para facilitar o desenvolvimento de software. Elas oferecem implementações de funções e rotinas que atendem a necessidades comuns, evitando que os desenvolvedores precisem reescrever código básico e complexo, o que aumenta a produtividade e a confiabilidade dos programas. Dentro das bibliotecas padrão, uma das mais utilizadas é a biblioteca matemática, que contém funções essenciais para realizar cálculos, além de oferecer operações como trigonometria, exponenciação, logaritmos, entre outras, frequentemente utilizadas em diversas áreas da ciência. Neste trabalho, investigam-se alternativas para a implementação de funções matemáticas na Robotics Language, desenvolvida na Universidade Federal de Jataí. Foram implementadas e avaliadas seis funções matemáticas: seno, cosseno, tangente, exponencial, logaritmo e raiz quadrada, utilizando algoritmos alternativos, como o CORDIC para as operações trigonométricas. Avaliou-se o tempo de execução e o tamanho do programa objeto para as plataformas AVR e STM32. Os resultados iniciais indicaram que os algoritmos são funcionalmente adequados e precisos, mas não competitivos em termos de tempo de execução e tamanho do programa objeto gerado, considerando o estágio atual de desenvolvimento do Robcmp e dos backends AVR e ARM do LLVM.

**Palavras-chaves:** Linguagem de programação; CORDIC; Bibliotecas padrão; Robcmp.

## Lista de ilustrações



# Lista de abreviaturas e siglas

IA	Inteligência Artificial
LLM	<i>Large Language Models</i> /Grandes modelos de linguagem
MCU	<i>Microcontroller Unit</i> /Microcontroladores
CPU	<i>Central Processing Unit</i>
RAM	<i>Random Access Memory</i>
ROM	<i>Read-Only Memory</i>
ROBL	<i>The Robotics Language</i>
IoT	<i>Internet of Things</i>
LLVM	<i>Low Level Virtual Machine</i>
LLVM-IR	<i>Low Level Virtual Machine-Intermediate Representation</i>
YaCoS	<i>Yet another Compiler Optimization Search</i>

# Sumário

# 1 Introdução

## 1.1 Motivação

Modelos de Inteligência Artificial (IA) têm sido amplamente utilizados em diversas tarefas no contexto da ciência da computação. Nos últimos anos, os avanços em Modelos de Linguagem de Grande Escala ou LLMs (*Large Language Models*) tornaram possível aplicar IA em atividades antes consideradas exclusivamente humanas, como geração de código, edição de projetos, sugestão e análise de codificação e, mais recentemente, otimização de código para compiladores.

Seria possível um modelo de linguagem treinado com exemplos de código de fases intermediárias do compilador substituir ou melhorar as estratégias de otimização empregadas por compiladores tradicionais? ??) partiram desta pergunta fundamental e treinaram um modelo baseado na arquitetura LLaMA 2<sup>1</sup> com 7 bilhões de parâmetros, alimentado por uma base de dados de funções LLVM-IR<sup>2</sup>, ou seja, códigos de linguagem intermediária gerados normalmente após a análise semântica de um compilador, não otimizados, associados a boas sequências de passes de otimização obtidas por busca semi-exaustiva (*autotuning*), além dos respectivos códigos otimizados. O principal objetivo era avaliar a possibilidade de um modelo de IA gerar a sequência correta de passes de otimização que diminuiria o tamanho do código binário resultante. O modelo treinado, chamado LLM-Compiler, foi capaz de gerar listas de passes de otimização que superaram as obtidas por estratégias como -Oz do LLVM, com uma média de melhoria de 3% a 5% na contagem de instruções, sem a necessidade de realizar múltiplas compilações como ocorre nos métodos tradicionais de *autotuning*. Além disso, ao ser perguntado sobre qual seria o código otimizado resultante da aplicação dos passes de otimização, o modelo demonstrou capacidade de raciocínio sobre código, atingindo mais de 90% de sucesso na geração de código compilável e uma taxa de equivalência exata de 70% com o código gerado pelo compilador usando os mesmos passes.

Os resultados de ??) se mostram promissores, visto que a *flag* -Oz é uma *flag* de otimização utilizada em diversos compiladores atuais como, por exemplo, Clang (que usa o *backend* LLVM), com o objetivo específico de gerar o menor tamanho de código binário possível durante a compilação. Existem outras *flags* comumente utilizadas como -Os, -O1, -O2 e -O3, sendo o -Os uma *flag* de otimização que tenta reduzir o tamanho do código sem comprometer muito a performance, e as *flags* de otimização de execução -O1, -O2 e -O3, que otimizam para uma velocidade boa, intermediária e rápida de execução de código,

<sup>1</sup> Site da Meta Llama 2: <<https://www.llama.com/llama2/>>

<sup>2</sup> Site do manual de referência da linguagem: <<https://llvm.org/docs/LangRef.html>>

respectivamente, porém, também aumentam o tamanho do código binário produzido e tornam o processo de compilação proporcionalmente mais lento. Como mostrado por ??), as *flags* de otimização tradicionais apresentam limitações. Os autores destacam que essas configurações genéricas muitas vezes não oferecem o melhor desempenho possível para todos os programas ou arquiteturas devido ao fato de aplicarem um conjunto fixo de transformações, sem considerar as características específicas do código-fonte ou do hardware de destino. Segundo ??), técnicas de aprendizado de máquina permitirão a seleção e ordenação dinâmica de passes de otimização com base em características específicas do código fonte e do ambiente de execução, personalizando o processo de compilação e maximizando ainda mais o desempenho. Sendo assim, é possível perceber que essa tecnologia é promissora e útil.

Atualmente, grande parte dos equipamentos de tecnologia em nossas residências e empresas compartilham uma característica em comum: em alguma instância, fazem uso de sistemas embarcados. Presentes em automóveis, telefones, impressoras, sensores, equipamentos médicos e uma infinidade de outros dispositivos eletrônicos essenciais para a indústria, comunicação, pesquisa e funcionamento de uma nação, sistemas embarcados são pequenos sistemas computacionais desenvolvidos para realizar uma tarefa específica, desde ligar um ar condicionado até auxiliar o controle de um drone (?). O cerne de um sistema embarcado são os microcontroladores ou MCUs (*Microcontroller Units*). Microcontroladores são circuitos integrados que compõem o “cérebro” de sistemas embarcados. Uma MCU possui uma CPU (*Central Processing Unit*), uma memória RAM (*Random Access Memory*), alguma forma de memória de longo prazo ou ROM (*Read-Only Memory*) e métodos de entrada e saída (?).

Segundo ??), microcontroladores desempenham um papel fundamental no funcionamento da sociedade atual e o mundo de hoje não funcionaria sem os *softwares* embarcados em nossos aparelhos. Logo, qualquer processo de otimização no funcionamento de MCUs e, consequentemente, de sistemas embarcados, resultaria em um enorme ganho para toda a sociedade. Porém, devido à sua natureza de atuação específica, as MCUs são extremamente restritas em relação ao poder de processamento e armazenamento de dados e necessitam de soluções eficientes de otimização de código (?).

O LLM-Compiler, desenvolvido por (?), embora seja uma solução promissora para atender os requisitos de otimização de código em geral, foi projetado e treinado no domínio de sistemas computacionais comuns, ou seja, para otimizar código para CPUs convencionais de 64 bits, X86\_64 e AArch64. Tais CPUs possuem um nível de capacidade computacional elevado comparado com microcontroladores. Por exemplo, a CPU Intel Core i9-10900K que implementa uma arquitetura de 64 bits, com diversas extensões ao conjunto de instruções, como Intel SSE4.1, Intel SSE4.2 e Intel AVX2, possui uma memória cache de 20 MB e suporta memória RAM DDR4 de até 128 GB, com uma

frequência base de *clock* de 3.70 GHz (??). Em comparação, arquiteturas como a AVR RISC de 8 bits, foco de atuação deste trabalho, presente em microcontroladores como o ATmega328P, possuem características extremamente limitadas, como 32 KB de memória *flash* programável, 2 KB de SRAM interna, e um *clock* máximo de 16 MHz (??). Sendo assim, levando em consideração a utilidade de sistemas embarcados, que utilizam MCUs, e também os resultados otimistas produzidos pelo trabalho de ??), surge a necessidade de se adaptar o LLM-Compiler para suprir a demanda por otimização de *software* em microcontroladores de sistemas embarcados.

Considerando estas demandas, o projeto *Robotics Language: Uma Linguagem de Programação de Propósito Específico para Microcontroladores* (PI05974-2024), desenvolvido na Universidade Federal de Jataí, assim como seu antecessor, Especificação e Construção de Protótipos Funcionais de Kits Robóticos de Baixo Custo para uso em Processos de Ensino-Aprendizagem (PI02361-2018), tem como objetivo desenvolver uma linguagem de programação, chamada *The Robotics Language* (ROBL), e seu compilador, denominado Robcmp, específica para robótica e microcontroladores.

Para o exercício deste trabalho, a capacidade do Robcmp de abstrair aspectos específicos do *hardware* e permitir um desenvolvimento compatível com diversos microcontroladores de forma dinâmica e sua integração com o *backend* AVR do LLVM (??), o torna a linguagem ideal para servir como ferramenta de conversão e compilação do conjunto de códigos de entrada e das respostas na saída do modelo.

## 1.2 Objetivo do Trabalho

Diante desse cenário, o presente trabalho teve como objetivo elaborar um conjunto de ferramentas, códigos e materiais necessários para realizar o refinamento do modelo LLM-Compiler (??), e possibilitar a plicação futura de um treinamento adicional com programas e otimizações no domínio de microcontroladores, especificamente a plataforma alvo AVR do LLVM.

Nossa proposta envolveu a formação de uma base de dados específica, composta por códigos de representação intermediária (IR), extraídos durante o processo de compilação na linguagem ROBL. Estes códigos, previamente convertidos da linguagem C, juntamente com suas respectivas listas de passes ótimos de compilação, obtidas através de estratégias semi-exaustivas, são a base necessária para a elaboração de um *dataset* dedicado ao refinamento supervisionado do modelo de (??). Nossa hipótese é que o refinamento permita que o modelo de IA seja capaz de compreender a estrutura e as restrições dos códigos voltados para MCUs AVR e, a partir disso, conseguir sugerir listas personalizadas de passes de compilação.

Os objetivos específicos deste trabalho foram:

1. Construir uma base de dados contendo códigos em C e posteriormente convertê-los para ROBL;
2. Compilar essa base de dados utilizando Robcmp para códigos IR específicos para Atmega328p;
3. Encontrar os passes ótimos de instruções para cada código IR;
4. Elaborar sequências de *prompt/label* utilizando o IR não otimizado e os passes ótimos encontrados, suficientemente para formar um dataset expressivo;

## 1.3 Contribuição do Trabalho

Este trabalho contribuirá com a criação, elaboração e disponibilização de um pipeline completo para o treinamento e refinamento do LLM-Compiler no domínio de microcontroladores AVR. A contribuição inclui, de forma organizada e reproduzível:

- O desenvolvimento e documentação de um fluxo (pipeline) automatizado para geração, conversão e preparação de datasets compatíveis com o modelo (incluindo o conversor `c2rob` e os scripts de processamento);
- A construção e curadoria de datasets prontos, com códigos convertidos para ROBL, pares *prompt/label* adequados para treinamento supervisionado e metadados que permitam reproduzir experimentos;
- A disponibilização das ferramentas e artefatos necessários ao processo de refinamento (scripts de autotuning, integração com YaCoS quando aplicável, mecanismos de validação e configuração para execução do fine-tuning em plataformas como Hugging Face).

A contribuição dessa pesquisa foi definida como o desenvolvimento do processo completo e reproduzível que permite o refinamento: ou seja, todo o conjunto de métodos, datasets e ferramentas necessários para, posteriormente, executar o fine-tuning do LLM-Compiler no domínio AVR, bem como avaliar seus resultados. Assim, o mérito científico e prático deste trabalho reside na preparação e disponibilização de um pipeline robusto e pronto para ser usado no refinamento do modelo, reduzindo significativamente a barreira de entrada para futuras execuções experimentais e permitindo que trabalhos subsequentes concluam o fine-tuning de maneira reproduzível e eficiente.

## 2 Referencial Teórico

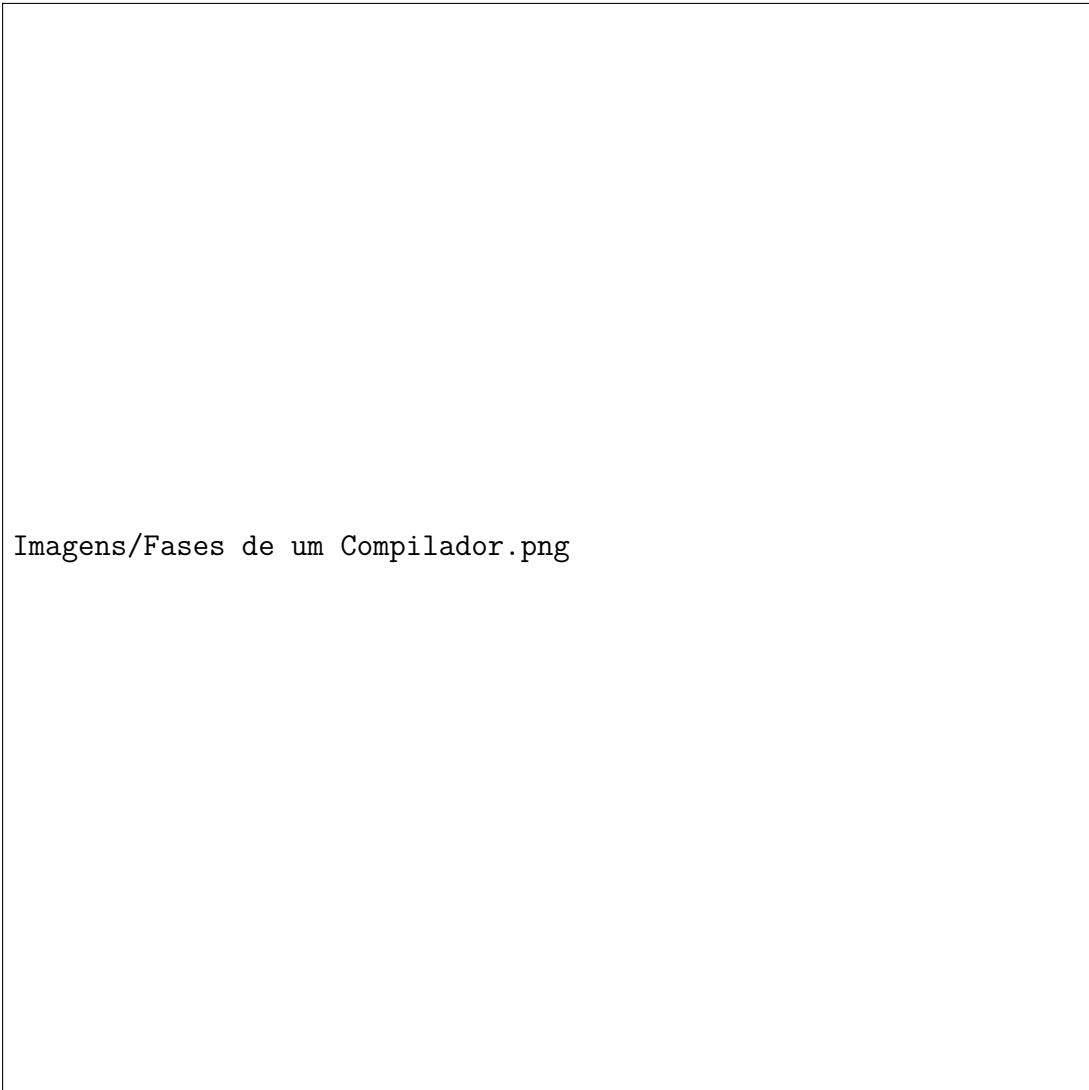
Este capítulo define os conceitos utilizados para a construção e estudo deste projeto. De início, conceitua-se compiladores e microcontroladores, o que são, como funcionam, onde atuam, suas limitações e capacidades. Posteriormente, o texto aborda Inteligência Artificial, detalhes de como funcionam, aprendizagem de máquina, em específico aprendizado supervisionado, LLMs e aplicações no mundo moderno. Por fim, este capítulo discute as possíveis estratégias de otimização de compiladores, com o uso de Inteligência Artificial.

### 2.1 Compilador

De acordo com ??), um compilador é um programa de computador que recebe outro programa em uma determinada linguagem de alto nível (linguagens de programação que se aproximam de linguagens humanas) e o traduz para uma linguagem de baixo nível (linguagem binária capaz de ser executada pelo processador do computador).

A importância de compiladores está atrelada ao desenvolvimento da tecnologia e à popularidade de sistemas de software em diversas aplicações. Com o advento da IoT (*Internet of Things*), conceito social que atrela a conexão de diversos dispositivos à rede e a computação na nuvem, a necessidade de compiladores otimizados se mostra maior do que nunca, visto que compiladores foram responsáveis pelo crescimento acelerado da tecnologia como conhecemos hoje (??).

A estrutura básica de um compilador convencional pode ser definida em três etapas: *Front-end*, *Middle-end* e *Back-end* que seguem um fluxo de funcionamento decrescente desde o código de entrada mais alto nível até o código da máquina alvo, assim como ilustrado na Figura ??. O *Front-end* é responsável pela análise do código-fonte, convertendo-o em uma árvore sintática. Nessa etapa, ocorrem processos como a análise léxica, sintática e semântica, que verificam a estrutura e o significado do programa de acordo com as regras da linguagem de programação utilizada.



Imagens/Fases de um Compilador.png

Figura 1 – Fases de um Compilador. Fonte: (??)

O *Middle-end* realiza otimizações independentes da arquitetura de hardware, transcrevendo o código original do programa de *input*(entrada) para uma versão intermediária otimizada de fácil discernimento, mantendo a semântica original do programa. O *Middle-end* pode ser representado como duas etapas que dependem da árvore sintática gerada pelo *front-end*, o gerador de código intermediário e o otimizador de código dependente da máquina alvo, ambas etapas geram ao final uma representação intermediária.

Por fim, o *Back-end* é encarregado de traduzir o código intermediário otimizado em código de máquina específico para a arquitetura alvo, com algumas alterações específicas para a melhor execução.

### 2.1.1 *Front-end*

O *front-end* é a etapa inicial do compilador. Ele é responsável por analisar o código-fonte e garantir que está bem estruturado para, posteriormente, poder ser traduzido



em uma representação intermediária do mesmo programa. O *front-end* age de forma independente ao contexto da máquina aplicada, tendo uma estrutura de execução quase universal, aplicável para diferentes contextos de diferentes linguagens de alto nível (??).

Sua abordagem consiste em analisar de forma léxica (o que é relativo ao vocabulário, à linguagem), depois analisar de forma sintática (o significado das expressões quanto à estrutura) e, por fim, analisar de forma semântica (o sentido da expressão no que tange ao contexto em que está inserida). Dessa forma, o *front-end* consegue confirmar que o código-fonte é bem formado, escrito corretamente e possui sentido (??).

### 2.1.2 Análise Léxica

O analisador léxico de um compilador é a primeira etapa percorrida do *front-end*. Essa parte, também conhecida como *Scanner*, tem como função transformar a sequência de caracteres do código-fonte em uma sequência de *tokens*, cada *token* representa uma unidade absoluta do código-fonte, como pontos, operadores, palavras-chave, números, símbolos, dentre outros. O *Scanner* é a primeira de três etapas que o compilador utiliza para entender o código-fonte, atuando diretamente sobre todo o código-fonte, fazendo com que essa etapa tenha um *input* maior do que as outras (??).

O analisador léxico é construído com base em expressões regulares que definem os padrões léxicos da linguagem (??) (??). Essas expressões são transformadas em autômatos finitos que processam os caracteres um a um e identificam onde cada *token* começa e termina. Quando um *token* é reconhecido, ele é emitido com uma etiqueta indicando seu tipo e, em alguns casos, com um valor associado (como o nome de uma variável ou o valor numérico de uma constante).

Além de identificar *tokens*, o analisador léxico pode também eliminar elementos que não são importantes para a execução do código, como espaços em branco, formatações e comentários. Essa limpeza torna o processo mais eficiente e simplifica o trabalho do analisador sintático. Também é comum que o *scanner* registre a posição dos *tokens* (linha e coluna) para auxiliar na identificação do código segundo o formalismo das expressões dos autômatos finitos (??).

### 2.1.3 Análise Sintática

O analisador sintático, também chamado de *parser*, é a segunda etapa do *front-end* de um compilador. Seu objetivo é verificar se a sequência de *tokens* produzida pelo analisador léxico está organizada de forma que respeite as regras gramaticais da linguagem de programação, geralmente descritas por uma gramática livre de contexto. Quando essa verificação tem sucesso, o parser constrói uma árvore sintática ou uma árvore de sintaxe abstrata (*Abstract Syntax Tree/AST*), que representa a estrutura hierárquica do programa

de acordo com sua sintaxe. Conforme ??), o parser é responsável por modelar a estrutura do programa de maneira que reflita sua lógica de construção, seu sucesso depende da definição cuidadosa da gramática da linguagem e da elaboração de seus algoritmos.

De acordo com ??), a análise sintática possui dois grandes grupos de técnicas, são eles: *parsing top-down* e *parsing bottom-up* (*parsing* de cima para baixo e de baixo para cima, respectivamente), cada um possui técnicas diferentes (LL e LR), que devem ser escolhidas de acordo com a linguagem e com a complexidade. *Parsers* LL são mais simples e geralmente implementados manualmente, enquanto *Parsers* LR são mais poderosos e geralmente aplicados por outros softwares. Em qualquer técnica utilizada o resultado esperado é uma AST que representa de forma precisa o código fonte. A AST atua como uma espécie de conexão entre o *front-end* e o *back-end*, visto que servirá como a base para todas as demais análises e transformações que seguirão.

A Figura ?? ilustra uma AST para um trecho de código simples. Nela, cada nó representa uma construção sintática do programa, como uma declaração de variável, um identificador, um tipo, um operador, ou um inteiro literal. Por exemplo, a raiz da árvore pode ser uma função ou um bloco de código, que se desdobra em declarações e expressões. A estrutura hierárquica da AST reflete a precedência e a associação das operações no código original, por exemplo, ao representar a expressão  $4 + 2 * 10 + 3 * (5 + 1)$ , a AST mantém a ordem e a lógica original das funções do programa e remove detalhes irrelevantes da sintaxe (como parênteses ou pontos e vírgulas), focando apenas nos elementos essenciais para a semântica.

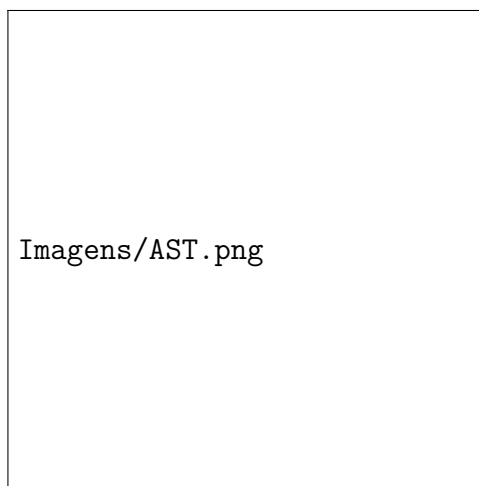


Figura 2 – Exemplo de uma *Abstract Syntax Tree* para  $4 + 2 * 10 + 3 * (5 + 1)$ . Fonte: <https://keleshev.com/abstract-syntax-tree-an-example-in-c/>

#### 2.1.4 Análise Semântica

O analisador semântico é a terceira etapa do *front-end* de um compilador, responsável por garantir que o programa, além de estar correto do ponto de vista sintático,

seja também semanticamente coerente. Ou seja, ele verifica se o programa faz sentido no contexto da linguagem, respeitando regras como declaração e uso de variáveis, estrutura de dados e outros aspectos que não podem ser descritos apenas por uma gramática livre de contexto. Essa análise é feita a partir da AST gerada pelo *parser*, geralmente com o auxílio de uma tabela de símbolos (??).

Segundo ??), a análise semântica percorre a AST executando diversas verificações contextuais que não podem ser descritas apenas por uma gramática livre de contexto. Essas verificações são realizadas em etapas bem definidas, que compõem o fluxo típico dessa fase no compilador:

1. **Construção e manutenção da tabela de símbolos:** à medida que a AST é percorrida, o compilador coleta informações sobre declarações de variáveis, funções, tipos e escopos, armazenando-as na tabela de símbolos para posterior consulta e validação (??).
2. **Verificação de tipos:** o analisador semântico assegura que as expressões e operações do programa estejam de acordo com as regras de tipagem da linguagem. Por exemplo, operadores aritméticos devem atuar sobre operandos numéricos compatíveis, e chamadas de função devem passar argumentos do tipo e quantidade corretos (??).
3. **Verificação de conversão de tipos (*casting*):** o compilador analisa se há necessidade de conversão entre tipos, automática ou explícita, e verifica se essas conversões são válidas no contexto da linguagem. Erros ou alertas são gerados quando conversões perigosas não são tratadas corretamente (??).
4. **Verificação de escopo:** a análise semântica valida se as variáveis e funções estão sendo acessadas dentro de seus escopos apropriados. Isso evita, por exemplo, o uso de variáveis declaradas fora do bloco de código atual.
5. **Verificação de fluxo de controle:** comandos como **break**, **continue** e **return** são analisados para garantir que aparecem em contextos válidos, como dentro de laços ou funções.
6. **Verificação de unicidade:** o compilador verifica se nomes de variáveis, membros de estruturas, funções ou rótulos não estão sendo definidos mais de uma vez dentro do mesmo escopo, o que seria ilegal.

Cada uma dessas etapas contribui para assegurar que o programa respeite as regras contextuais da linguagem de programação. Ao final dessa fase, a AST encontra-se enriquecida com informações de tipo e escopo, e pronta para ser traduzida para uma representação intermediária, iniciando a próxima etapa do processo de compilação (??).

### 2.1.5 Geração de código intermediário

Uma das últimas coisas que o *front-end* de um compilador faz, é a geração de código intermediário (??). Essa representação intermediária (IR, do inglês *Intermediate Representations*), geralmente utilizada para otimizações e para facilitar a geração do código de máquina (??), serve como entrada para o gerador de código, que, juntamente com as informações da tabela de símbolos, produz o código objeto equivalente. As IRs são a forma de como um compilador consegue representar os estados e as informações do código que ele compila, compiladores podem utilizar uma ou mais representações intermediárias, variando de acordo com o processo executado para a tradução até a linguagem alvo (??).

Existem diversos tipos de Representações Intermediárias (IRs) utilizadas em compiladores, cada uma com suas características e finalidades específicas (??). Alguns compiladores geram uma representação intermediária de baixo nível chamada, código de três endereços, nessa IR, cada instrução contém no máximo uma operação, e devido a sua simplicidade é mais simples aplicar otimização de código nesta representação, essa representação intermediária deve ter duas propriedades importantes: ser facilmente produzida e ser facilmente traduzida para a máquina alvo (??).

Como o LLVM-Compiler é voltado para arquiteturas de computadores que utilizam a infraestrutura LLVM, o foco deste trabalho será no LLVM-IR, a representação intermediária de compiladores a base de LLVM, projetada para ser uma “IR universal”, capaz de representar diversas linguagens de alto nível (??). A LLVM-IR pode ser dividida entre: uma representação em memória para o compilador, uma representação em *bitcode* para carregamento rápido e uma representação em linguagem *assembly* legível por humanos. Essa abrangência da LLVM-IR permite processos de otimização eficiente para transformações do compilador, além de detectar possíveis erros ou problemas na tradução.

A Figura ?? é uma representação de código intermediário LLVM, suas instruções como ‘zext’, ‘trunc’, ‘br’, ‘call’ e ‘ret’ exemplificam operações típicas da linguagem intermediária, como extensão de zero, truncamento de inteiros, desvios condicionais e incondicionais, chamadas de função e retornos, sua linguagem compreensível para humanos permite fácil compreensão de código e facilita a análise, otimização e posterior uso. Essa clareza é especialmente útil durante o desenvolvimento e treinamento de modelos baseados em aprendizado de máquina, como o LLVM-Compiler de ??), que dependem da extração de padrões sintáticos e semânticos do LLVM-IR para sugerir ou otimizar sequências de passes de compilação.

### 2.1.6 *Middle-end*

O *Middle-end* de um compilador atua como uma ponte entre a análise (*front-end*) e a geração de código final (*back-end*). Seu principal objetivo é transformar a árvore de

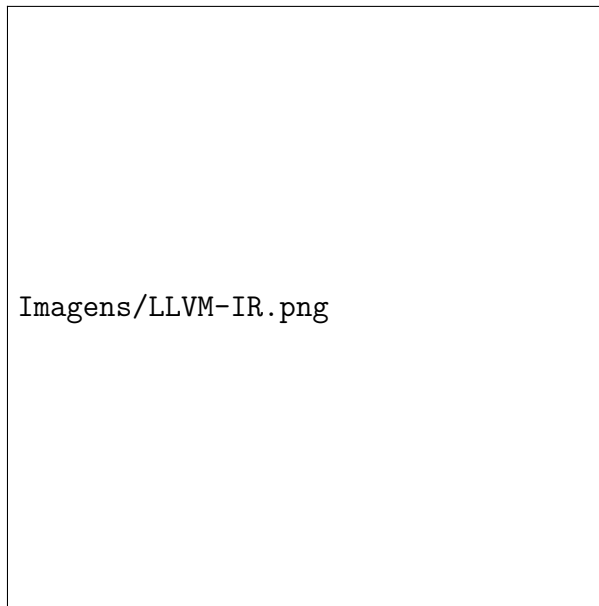


Figura 3 – Exemplo de um código em LLVM-IR. Fonte: <<https://github.com/colejcummins/llvm-syntax-highlighting>>

sintaxe abstrata (AST), produzida e validada pelas fases anteriores, juntamente com a tabela de símbolos, em uma representação intermediária (*Intermediate Representation/IR*) de forma que seja mais fácil para otimizar e depois gerar o código em linguagem baixo nível (??).

Uma boa representação intermediária abstrai os detalhes da linguagem fonte e da máquina alvo, permitindo que o compilador aplique uma série de análises e otimizações, como propagação de constantes, eliminação de código morto, análise de alcance de variáveis e redução de laços (??). Essas transformações visam melhorar o desempenho, reduzir o consumo de memória e aumentar a eficiência do código gerado, sem alterar o comportamento do programa original.

Além disso, o uso de uma IR padronizada, como o SSA (*Static Single Assignment*), proporciona uma base sólida para a construção de compiladores capazes de gerar código para diferentes arquiteturas a partir do mesmo núcleo de análise (??). Um bom exemplo é a LLVM (*Low Level Virtual Machine*) que adota o formato SSA como estrutura central em sua representação intermediária, o LLVM-IR. Essa escolha permite que múltiplas ferramentas de otimização e geração de código operem de forma coesa e eficiente sobre uma representação comum. Essa separação entre as fases do compilador, favorece a reutilização de componentes, a portabilidade e a manutenção do sistema de compilação (??).

É também no *Middle-end* que ocorre o processamento dos passes de otimização de código independentes de plataforma alvo, os quais serão alvo de escolha e ordenação pelo modelo de IA a ser refinado neste trabalho.

### 2.1.7 *Back-end*

O *back-end* de um compilador, segundo ??) e ??), é responsável por transformar a representação intermediária (IR) produzida pelo *front-end* em código de máquina ou código objeto específico para a arquitetura alvo. Essa fase ocorre após todas as análises e validações semânticas terem sido realizadas, e o programa já estar representado de forma estruturada e otimizada.

De acordo com ??), o *back-end* executa três funções principais: seleção de instruções, alocação de registradores e agendamento de instruções. A seleção de instruções envolve mapear operações da IR para instruções da linguagem de máquina. Já a alocação de registradores trata da distribuição eficiente das variáveis temporárias nos registradores físicos disponíveis, dado o número limitado desses recursos. O agendamento de instruções busca reorganizar o código para melhorar a performance, por exemplo, evitando dependências e aumentando o paralelismo entre instruções.

??) destacam que o objetivo do *back-end* é gerar código eficiente e correto para a máquina alvo, respeitando todas as restrições impostas pela arquitetura. Ele também ressalta a importância da modularidade: ao usar uma representação intermediária bem projetada, torna-se possível desenvolver diferentes *back-end* para várias arquiteturas, reutilizando o *front-end* e o *middle-end*. Isso é especialmente útil para compiladores que precisam ser portáteis ou suportar múltiplas plataformas.

## 2.2 Microcontroladores

Microcontroladores são componentes fundamentais no desenvolvimento de sistemas embarcados. De forma geral, um microcontrolador pode ser definido como um processador que possui recursos integrados, tais como memória RAM, espaço para código e interfaces periféricas, como linhas de entrada e saída (??). Essa integração de funcionalidades permite que os microcontroladores operem de forma independente em aplicações específicas, sem a necessidade de sistemas operacionais completos, diferentemente dos computadores de uso geral (??).

O uso de microcontroladores está amplamente difundido em dispositivos do cotidiano, como eletrodomésticos, brinquedos eletrônicos, sistemas automotivos, equipamentos médicos, sensores industriais, entre outros. Por serem aplicados em contextos de objetivo restrito, esses dispositivos normalmente enfrentam limitações significativas em termos de recursos computacionais, consumo de energia, e capacidade de armazenamento (??). Um bom exemplo disso é o modelo de MCU ATmega328p mostrado na Figura ??, pertencente a família de microcontroladores AVR, alvo deste trabalho. Este chip implementa uma arquitetura Advanced RISC de 8 bits e possui somente 256KB de memória *flash* programável (??).

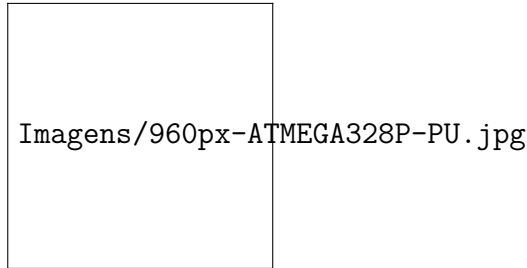


Figura 4 – Foto de exemplificação do microcontrolador estudado, Atmega328p Fonte: <<https://www.conectabit.com.br/MLB-4181008568-chip-atmega328p-atmega328p-u-atmega-328-pu-dip-32k-20mhz-JM>>

Segundo ??), a programação para microcontroladores é feita diretamente sobre o *bare metal*, ou seja, sem a presença de um sistema operacional intermediário. Isso significa que, ao escrever um comando como "acender um LED", o software se comunica diretamente com o hardware, conferindo maior controle e eficiência, características essenciais em sistemas com restrições severas de tempo real e consumo.

Além disso, microcontroladores geralmente estão acoplados a diversos sensores, atuadores e demais periféricos, sendo necessário lidar com a integração entre hardware e software de forma coordenada e robusta. Para enfrentar tais desafios, boas práticas de arquitetura de software, como modularidade, encapsulamento e uso de padrões de projeto, são essenciais para garantir flexibilidade, manutenção e reusabilidade do sistema ao longo de seu ciclo de vida (??).

### 2.2.1 *Domain Specific Languages*(DSL)

As *Domain-Specific Languages* (DSLs) são linguagens de programação ou notação projetadas para atender a um domínio específico de aplicação, oferecendo maior expressividade e facilidade de uso nesse contexto (??). Ao contrário das linguagens de programação de uso geral, que tentam atender a uma ampla variedade de problemas, as DSLs focam em um nicho restrito, como bancos de dados, planilhas eletrônicas, modelagem de *hardware* ou geração de relatórios.

Justamente pelo fato de serem desenvolvidas com um domínio específico em mente, essas linguagens tornam o trabalho mais simples e direto. Elas permitem, por exemplo, que profissionais da área possam interagir com sistemas complexos de forma mais intuitiva, sem necessitar de um alto grau de qualificação técnica (??). Essa especialização pode reduzir os custos de desenvolvimento e manutenção, ao mesmo tempo que aumenta a reutilização de *software* e a precisão dos sistemas (??).

Além disso, DSLs podem ser implementadas de diferentes formas, como interpretadores, compiladores, pré-processadores ou linguagens embutidas. Essa flexibilidade de

implementação e a capacidade de abstrair detalhes de baixo nível e expressar intenções de forma mais clara e concisa é particularmente vantajosa no contexto de sistemas embarcados, que operam sob severas restrições de memória, tempo de execução, consumo de energia e capacidade de processamento. De acordo com ??), práticas como modularidade, encapsulamento e reutilização de componentes são fundamentais para enfrentar os desafios comuns no desenvolvimento de *software* embarcado. DSLs bem projetadas podem fortalecer esses princípios ao oferecer construções específicas que se alinham diretamente com o hardware e com as tarefas que o software precisa desempenhar.

### 2.2.2 Robotics Language

A *Robotics Language* (ROBL) é uma linguagem de programação e compilador desenvolvidos com foco em aplicações de microcontroladores voltadas para robótica e Internet das Coisas (IoT). Seu principal diferencial está na abstração das particularidades de hardware diretamente na linguagem e em sua biblioteca padrão, permitindo que o mesmo código-fonte seja compilado para diferentes plataformas sem necessidade de ajustes manuais ou uso de diretivas condicionais (??).

Além de promover portabilidade entre arquiteturas, o ROBL realiza uma análise semântica aprofundada durante a compilação, prevenindo diversos tipos de erros que normalmente só seriam identificados em tempo de execução. Essa característica contribui para um desenvolvimento mais seguro e confiável, especialmente em sistemas embarcados, onde falhas podem ser críticas.

A implementação do compilador foi construída utilizando ferramentas clássicas no desenvolvimento de compiladores, como o *Flex* (versão 2.6.4) para análise léxica e o *Bison* (versão 3.8.2) para análise sintática e geração de código intermediário. O *backend* do compilador se apoia no *LLVM*, um *framework* moderno e modular que possibilita tanto otimizações quanto a geração de código para diferentes arquiteturas de microcontroladores.

O ecossistema do Robcmp também inclui suporte a depuração via simulador e integração com o Visual Studio Code (editor de código), por meio da extensão *RobCmpSyntax*, que fornece realce de sintaxe e facilita o desenvolvimento. Esses recursos tornam o Robcmp uma ferramenta especialmente atrativa em ambientes educacionais, oferecendo uma alternativa mais segura e acessível ao tradicional uso da linguagem C em projetos embarcados.

## 2.3 Inteligência Artificial

De acordo com ??) a Inteligência Artificial (IA) pode ser compreendida como o estudo de agentes inteligentes, entidades que percebem seu ambiente e tomam ações que maximizam suas chances de sucesso em atingir objetivos. Para que algo seja considerado



uma IA, ele deve ser capaz de realizar funções geralmente feitas por seres humanos, como perceber o ambiente, raciocinar e tomar decisões, aprender com dados e agir no mundo físico (????). Seja através da robótica, do aprendizado de máquina, de sistemas probabilísticos ou de mapeamento 3D em tempo real, softwares de Inteligência Artificial possuem a capacidade de agir de forma racional ou até imitar a capacidade humana até certo ponto (??).

Recentemente a IA tem-se demonstrado uma ferramenta poderosa para trabalho, logística e lazer dentro de diversos setores da sociedade atual (??). Tecnologias como reconhecimento de voz, tradução automática, veículos autônomos, sistemas de recomendação e diagnósticos médicos automatizados são somente alguns exemplos de como algoritmos e softwares de Inteligência Artificial são necessários e fazem parte do mundo moderno (????). A característica singular de adaptação de forma racional a ambientes complexos e incertos, antes presente somente em humanos, evidencia a importância desses sistemas inteligentes no desenvolvimento de novas tecnologias e na superação de desafios anteriormente insuperáveis (??).

No nível mais geral, uma IA funciona como um agente racional, que percebe seu ambiente por meio de sensores e age sobre ele por meio de atuadores. Esse agente processa sequências de percepções para decidir qual ação tomar, com base em alguma função desejada. O agente pode ser simples, reagindo diretamente aos estímulos, ou complexo, utilizando modelos internos, planejamento, raciocínio e aprendizado para adaptar seu comportamento a longo prazo (????).

### 2.3.1 *Machine Learning*

??) explicam que *machine learning* é uma subárea da Inteligência Artificial que estuda algoritmos e modelos que permitem a sistemas computacionais aprender com base em experiências, grupos de dados e diretrizes, para melhorar sua habilidade e execução de alguma tarefa. Em outras palavras é campo da IA que estuda como agentes podem melhorar automaticamente seu desempenho por meio da experiência (????).

De uma forma didática, *machine learning* pode ser explicado pela seguinte lógica: diferentemente de outros sistemas e softwares com execução linear, onde se tem o *input*(entrada), o método de processamento e se deseja saber a saída, em algoritmos de *machine learning* se possui o *input* e o *output*(saída) e se deseja saber ou fazer com que a máquina entenda, o meio termo, ou seja, o caminho até a saída. Na prática, em vez de programar explicitamente todas as regras para uma tarefa, fornece-se ao sistema uma grande quantidade de dados para que ele interprete padrões e regras por conta própria, podendo replicá-los futuramente (????).

Atualmente, qualquer treinamento de IA envolve três elementos principais: os

exemplos (dados de entrada e saída esperada), o modelo de representação e a função de avaliação e otimização. Os exemplos são o que alimenta a aprendizagem, são os dados iniciais que o modelo de representação recebe durante sua formação como IA, podem ser diversos tipos de arquivos ou outras formas de dados. O modelo de representação é a forma matemática ou computacional usada para expressar o conhecimento aprendido, pode ser uma árvore de decisão binária ou um sistema de nós que manipula informações (redes neurais). Por final a função de otimização avalia os resultados obtidos pelo modelo durante o treinamento e corrige seus parâmetros internos para melhorar os resultados com base em critérios de erro ou recompensa.

Dessa forma, ??) define que o tipo de *machine learning* aplicado a um modelo depende da natureza da tarefa e dos dados disponíveis para treiná-lo. Podendo serem divididos entre:

1. Aprendizado supervisionado, dados rotulados que ajudam o modelo a classificar semelhanças entre os dados em grupos dos rótulos facilitando o reconhecimento de padrões. Ex: Fotos de pássaro com o rótulo papagaio e fotos de árvores com o rótulo pinheiro, facilitariam um modelo de reconhecimento de imagens.
2. Aprendizado não supervisionado, dados não rotulados, permitindo ao modelo descobrir padrões não específicos nos exemplos, facilitando conexões não aparentes. Ex: Agrupamento de arquivos de diferentes tipos com base no conteúdo.
3. Aprendizado por reforço, um agente recebe recompensas conforme os resultados que suas ações geram no meio inserido. Ex: Ações de um boneco ao chegar mais longe em uma fase de um jogo eletrônico.

### 2.3.2 Aprendizado Supervisionado

Com base no que ??) o aprendizado supervisionado é uma das formas mais fundamentais e amplamente utilizadas de *machine learning*. Nesse tipo de aprendizado, o agente recebe um conjunto de exemplos rotulados, ou seja, entradas acompanhadas de suas respectivas saídas corretas. O objetivo do sistema é aprender uma função que generalize esses exemplos, ou seja, que seja capaz de prever corretamente a saída para novas entradas nunca vistas antes. O conjunto de exemplos é composto por pares de dados (*data,label*), esses pares serão analisados por uma árvore de decisão que compõe o sistema de aprendizado, permitindo que o modelo identifique padrões e classifique de forma precisa durante a execução ??).

Durante o treino cada exemplo de informa a resposta certa, e o algoritmo tenta ajustar seu modelo interno para minimizar os erros entre suas previsões e os valores reais. Com o tempo e com mais exemplos, o modelo vai melhorando sua capacidade de

prever resultados corretos mesmo quando confrontado com dados novos (??). A Figura ?? demonstra esse processo. O modelo de IA é alimentado por fotos de corujas, raposas e esquilos e também pelos rótulos respectivos de cada foto, compondo assim os pares (*data,label*), então o modelo é capaz de aprender regras e funções relativas a esse conjunto de dados que permite que ele classifique corretamente as imagens após terminado o treinamento.



Figura 5 – Representação didática do funcionamento de aprendizagem supervisionado. Fonte: <[https://datamapu.com/posts/ml\\_concepts/supervised\\_unsupervised](https://datamapu.com/posts/ml_concepts/supervised_unsupervised)>

Existem diversos exemplos de aprendizado supervisionado, os mais famosos são tarefas como classificação e regressão. Na classificação, o objetivo é prever uma categoria discreta, como determinar se uma foto é um papagaio ou uma árvore. Já na regressão, a saída é um valor contínuo, como prever o preço de uma casa com base em suas características (tamanho, localização, etc.).

Modelos comumente usados para aprendizado supervisionado incluem árvores de decisão, redes neurais, modelos lineares, entre outros. Cada um deles tem suas vantagens e é mais adequado para certos tipos de dados ou problemas. Independentemente do modelo escolhido, o desempenho é geralmente avaliado usando conjuntos de dados separados de teste, medindo métricas como acurácia, precisão, recall ou erro quadrático médio, dependendo da tarefa.

### 2.3.3 Large Language Models (LLM)

Um LLM, ou um grande modelo de linguagem, é um tipo de modelo estatístico treinado para prever a próxima palavra em uma sequência de texto, dado o contexto anterior. Essa tarefa, chamada de modelagem de linguagem, é um problema de aprendizado supervisionado (ou auto-supervisionado) onde, para cada entrada (um fragmento de texto), o modelo deve prever a saída correta (a próxima palavra ou *token*). Durante o treinamento, ele processa bilhões de palavras, ajustando seus bilhões de parâmetros para capturar padrões linguísticos, sintaxe, semântica e até aspectos pragmáticos do uso da linguagem (??).

O funcionamento básico de um LLM envolve três fases: o pré-treinamento, onde o modelo aprende representações gerais da linguagem com base em enormes volumes de texto (livros, sites, artigos); a refinamento (*fine-tuning*), em que ele é ajustado para tarefas específicas com dados mais controlados; e, em alguns casos, a aplicação de técnicas como aprendizado por reforço com *feedback* humano (RLHF), que serve para alinhar as respostas do modelo a valores humanos, como segurança, utilidade e cortesia (??).

LLMs são construídos sobre arquiteturas de redes neurais profundas chamadas *transformers*, que permitem o processamento eficiente de sequências longas de texto, capturando relações contextuais entre palavras, mesmo que estejam distantes no texto. Essa tecnologia presente na arquitetura de modelos de linguagem é fundamental para o funcionamento de paralelismo durante o treinamento do modelo, e também para permitir maior escalabilidade para maiores conjuntos de dados (????).

O treinamento de LLMs como ChatGPT<sup>1</sup>, Gemini<sup>2</sup> e Llama 2<sup>3</sup> se dá graças a um pré-treinamento e refinamento exaustivos e complexos, envolvendo um processo iterativo e em larga escala que se estende por várias etapas. Inicialmente, esses modelos passam por um pré-treinamento extensivo com um volume enorme de dados textuais, devido a isso o treinamento dessas IAs é extremamente custoso e demorado.

## 2.4 Otimização em Compiladores Modernos

A etapa de otimização em compiladores visa transformar o código intermediário (IR – *Intermediate Representation*) em uma versão mais eficiente, sem alterar sua semântica. Essa eficiência pode se referir à melhoria no desempenho de execução, economia de memória, redução de tamanho do binário, ou mesmo economia de energia, dependendo dos objetivos da aplicação alvo (????).

Tradicionalmente, compiladores como o LLVM aplicam conjuntos padronizados de

<sup>1</sup> Site dos modelos do ChatGPT: <<https://platform.openai.com/docs/models>>

<sup>2</sup> Site dos modelos do Gemini: <<https://ai.google.dev/gemini-api/docs/models?hl=pt-br>>

<sup>3</sup> Site da Meta para Llama2: <<https://www.llama.com/llama2/>>

passos de otimização, agrupados em níveis como -O1, -O2 e -O3, com sequências específicas voltadas à performance, e -Os ou -Oz, voltadas à compactação do código. Esses conjuntos são compostos por dezenas ou até centenas de transformações e análises estáticas, incluindo eliminação de código morto, propagação de constantes, *inlining* e ordenação de instruções (??).

Apesar da eficácia dessas abordagens, elas não são adaptáveis ao perfil específico de cada programa. Por isso, nos últimos anos, pesquisadores têm investigado o uso de aprendizado de máquina para guiar decisões de otimização de forma mais personalizada. A ideia central é que modelos possam aprender, a partir de exemplos anteriores, quais sequências de otimização produzem melhores resultados para diferentes classes de programas (??).

Um dos desafios centrais dessa abordagem é representar programas de forma mensurável. A extração de *features* — características numéricas derivadas da estrutura e comportamento do código — é essencial para alimentar modelos preditivos. Em trabalhos recentes, propõe-se representar programas como grafos, como forma de capturar dependências de controle e dados, permitindo que modelos de aprendizado relacional compreendam melhor o contexto das instruções (??).

Recentemente, modelos de linguagem de larga escala (LLMs) têm sido aplicados ao problema de otimização. Esses modelos são capazes de inferir boas sequências de passos de forma contextual, generalizando a partir de grandes volumes de dados de código e otimizações anteriores. O trabalho de ??) demonstra que LLMs treinados com representações intermediárias, como o LLVM-IR, podem alcançar desempenho competitivo com métodos heurísticos clássicos, muitas vezes sem a necessidade de compilações repetidas ou ajustes manuais extensivos.

Portanto, a otimização de compiladores pode, no futuro, passar por um momento de transição de heurísticas generalistas para abordagens baseadas em dados, aprendizado e raciocínio automatizado. Essa mudança combina os princípios tradicionais descritos por ??) e ??) com os avanços em IA e aprendizado de máquina discutidos por ??), abrindo caminho para compiladores mais inteligentes e adaptáveis.

## 3 Trabalhos Relacionados

### 3.1 Introdução

Este capítulo apresenta uma revisão acerca dos principais trabalhos relacionados ao tema desta pesquisa, com foco nas áreas de otimização de tamanho de código, busca de sequências de passes de compilação e, em especial, técnicas envolvendo modelos de linguagem de larga escala. O objetivo é contextualizar o refinamento do modelo LLM-Compiler para o domínio de microcontroladores AVR dentro do escopo da Robotics Language (ROBL), destacando contribuições, limitações e relações diretas com este estudo.

Inicialmente, são apresentados os critérios adotados para seleção dos trabalhos. Em seguida, são descritas as obras consideradas mais relevantes, discutindo suas metodologias, resultados e pertinência para o desenvolvimento desta pesquisa.

### 3.2 Metodologia de Análise

A análise dos trabalhos relacionados seguiu uma abordagem de Revisão Narrativa de Literatura (RNL), com busca realizada para quaisquer trabalhos que de forma direta ou indiretamente abordaram otimização de tamanho de código, associados ao ambiente LLVM. A seleção dos materiais também considerou referências citadas no próprio trabalho original do LLM-Compiler, em especial pesquisas que envolvem:

- **Modelos de IA aplicados à otimização de compiladores;**
- **Métodos de busca e seleção de passes de otimização** para redução de tamanho de código;
- **Ferramentas tradicionais de autotuning** como YaCoS;
- **Redução sistemática de sequências de passes;**
- **Construção de datasets** baseados em representações intermediárias;

### 3.3 Trabalhos Analisados

Com base nos critérios definidos, três trabalhos foram considerados fundamentais para o desenvolvimento desta pesquisa, por suas contribuições conceituais e práticas no problema de otimização de código e seleção de passes de compilação. As subseções a seguir apresentam um resumo estruturado de cada obra, destacando suas contribuições e sua

relação direta com a estrutura desta pesquisa e seu impacto em um futuro refinamento do modelo LLM-Compiler para o domínio AVR.

### 3.3.1 LLM-Compiler

O trabalho de (??) introduz o LLM-Compiler, um modelo de linguagem treinado para inferir sequências de otimização para o compilador LLVM, bem como gerar código otimizado equivalente ao produzido por essas sequências. O modelo, baseado em um LLaMA com 7 bilhões de parâmetros, é alimentado por um dataset composto por códigos LLVM-IR não otimizados e suas respectivas sequências de passes ótimos oriundos de métodos semi-exaustivos.

Os autores demonstram que o modelo é capaz de superar a flag `-Oz` em média de 3% a 5% na contagem de instruções, além de possuir alta capacidade de gerar código compilável e semanticamente equivalente. Esse trabalho constitui o ponto de partida desta pesquisa, uma vez que o presente estudo visa adaptar o LLM-Compiler para o domínio de microcontroladores AVR, expandindo sua capacidade de otimização para arquiteturas com restrições severas de memória e instruções.

### 3.3.2 Good Optimization Sequences Covering Program Space

O trabalho de (??) propõe um método para identificar sequências de otimização que cobrem eficientemente o espaço de programas, reduzindo a necessidade de testes exaustivos em todos os códigos de entrada. A abordagem utiliza remoção iterativa e análise de impacto dos passes, permitindo identificar sequências curtas, porém eficazes, na redução de tamanho e melhoria de desempenho.

Esta pesquisa é fundamental para o presente trabalho, pois seu método de redução de sequências é utilizado diretamente no processo de filtragem das listas geradas pelo autotuning descrito por Faustino et al. (2021). Assim, sua contribuição teórica fornece a base para garantir que somente os passes essenciais sejam mantidos nas listas que farão parte do dataset final.

### 3.3.3 Optimization Sequences for Code-Size Reduction

O trabalho de (??) apresenta um método sistemático para encontrar sequências de otimização voltadas especificamente à redução de tamanho de código no LLVM. Utilizando algoritmos genéticos combinados com a ferramenta YaCoS, os autores geram um grande conjunto inicial de possíveis sequências, filtram as melhores e, em seguida, aplicam o método de Purini e Jain (2013) para remover passes redundantes.

O presente trabalho utiliza diretamente essa metodologia para gerar o dataset de sequências ótimas para códigos ROBL compilados para LLVM-IR, servindo como fonte

principal para as listas de passes usadas no treinamento supervisionado do modelo refinado. Dessa forma, o trabalho de Faustino et al. constitui a principal referência prática para geração das sequências ótimas utilizadas.



## 4 Metodologia do desenvolvimento

### 4.1 Introdução

Neste capítulo será abordado o método aplicado no escopo desta pesquisa. Será detalhado a abordagem empregada durante as diversas etapas do processo de elaboração deste trabalho, especificando pontos chaves de cada etapa para garantir a replicabilidade desta pesquisa. Neste capítulo, também será exemplificado e descrito os códigos e ferramentas, utilizadas durante todo o processo, que podem ser encontrados [neste repositório do Github](#).

### 4.2 Elaboração do conjunto de códigos C

A primeira etapa do pipeline consiste na geração do conjunto inicial de programas em linguagem C, a partir de duas fontes distintas: códigos sintéticos produzidos automaticamente e códigos reais obtidos de repositórios públicos. O objetivo desta etapa é formar uma base ampla, variada e representativa o suficiente para sustentar o processo posterior de conversão, compilação e autotuning.

#### 4.2.0.1 Códigos sintéticos gerados com Csmith

Para a geração de códigos sintéticos, foi utilizado o software Csmith (??), disponibilizado publicamente em <https://github.com/csmith-project/csmith>. Após a instalação da ferramenta conforme instruções do repositório oficial, iniciou-se um processo exploratório para determinar parâmetros adequados de geração. Foram realizados diversos testes com diferentes combinações de opções do Csmith, a fim de identificar uma configuração capaz de produzir códigos: Não excessivamente extensos, evitando que os programas ultrapassassem a janela de contexto do modelo de linguagem a ser futuramente treinado e não demasiadamente simples, de modo a preservar diversidade estrutural e evitar um conjunto de dados pobre em variabilidade, que potencialmente conduzissem a *overfitting*.

Uma vez definidos os parâmetros ideais, foi desenvolvido um script em *shell*, denominado `csmithrun.sh`, responsável por automatizar a geração de grandes volumes de programas C. Esse script executa o Csmith repetidamente utilizando opções como `-max-funcs 2`, `-no-global-variables`, entre outras configurações que restringem a complexidade e asseguram compatibilidade posterior com o processo de conversão para ROBL. O código do script encontra-se disponível no repositório público: <https://github.com/Alisson-Teles/Pipeline-training/blob/main/AVR/CSMITH/csmithrun.sh>.

Com esse procedimento, foram gerados aproximadamente 80 mil códigos sintéticos em C. Embora a meta inicial fosse atingir 100 mil programas, estimativas de tempo de treinamento do modelo indicaram que 100k exemplos tornariam o processo excessivamente longo. Assim, adotou-se o conjunto de 80k códigos sintéticos como solução de compromisso entre volume, diversidade e viabilidade computacional.

#### 4.2.0.2 Códigos reais obtidos do AnghaBench

A segunda fonte de dados consiste em códigos reais provenientes do repositório AnghaBench (??), disponível em <https://github.com/brenocfg/AnghaBench/tree/master>. O AnghaBench contém aproximadamente um milhão de funções em C, extraídas automaticamente de diversos projetos hospedados no GitHub, apresentando grande variedade de estilos, padrões de escrita e estruturas de código típicas de aplicações reais.

Como os programas do AnghaBench já são curtos e subdivididos em funções isoladas, não houve necessidade de tratamento ou pré-processamento adicional. Para compor o subconjunto utilizado nesta pesquisa, foram selecionados 10 mil códigos aleatórios a partir dos diretórios do AnghaBench, utilizando o comando `find` no terminal para amostragem simples dos arquivos disponíveis.

Ao final desta etapa, o conjunto inicial de programas em C — combinando cerca de 80k códigos sintéticos e 10k códigos reais — foi consolidado e preparado para as etapas seguintes do pipeline.

### 4.3 Conversão de C para RobL

A segunda etapa do pipeline consiste exclusivamente na conversão dos programas escritos em C para a linguagem ROBL. Essa conversão é essencial para permitir que os códigos sejam compiláveis pelo Robcmp, e possam ser considerados aplicáveis em uma MCU AVR.

Para realizar a conversão, utilizou-se o conversor `c2rob`, disponível no repositório oficial do projeto Robotics Language: <https://github.com/thborges/c2rob>. O repositório foi clonado, e modificações pontuais foram aplicadas ao `Makefile` e a arquivos auxiliares, a fim de assegurar que o código convertido estivesse configurado especificamente para a arquitetura alvo adotada neste trabalho, o ATmega328P. Essas alterações garantiram conformidade com as características do backend AVR associado ao Robcmp, que compilava os códigos assim que já estivessem traduzidos, garantindo códigos RobL fiéis e compiláveis.

Como a linguagem ROBL não implementa todos os recursos da linguagem C, tornou-se necessária a aplicação de um mecanismo de filtragem automática dos programas obtidos na etapa anterior. A filtragem foi realizada diretamente a partir da saída do

**c2rob**: sempre que o conversor emitia um erro de tradução o código correspondente era automaticamente excluído do fluxo principal. Assim, o erro de sintaxe emitido pelo **c2rob** funcionou como um critério objetivo de incompatibilidade.

A escolha de gerar aproximadamente 80 mil programas C com o Csmith contribuiu significativamente para esse processo, garantindo ampla variedade estrutural e reduzindo a necessidade de intervenção manual. Os programas que não puderam ser convertidos para ROBL, mas ainda eram válidos em C, foram compilados com o Clang utilizando a flag **-O1**, garantindo ao menos um nível básico de tratamento e permitindo sua eventual utilização para fins experimentais e comparativos em etapas posteriores.

Ao final desta fase, obtém-se exclusivamente o conjunto de programas que foram efetivamente convertidos para ROBL e considerados aptos para seguir para as etapas seguintes do pipeline.

## 4.4 Método exaustivo de busca da melhor sequência de passes de otimização.

A terceira etapa do pipeline consiste na identificação da melhor sequência de passes de otimização para cada código convertido, com o objetivo de minimizar o tamanho do código objeto gerado para a plataforma AVR. Embora a metodologia completa descrita por (??) não tenha sido aplicada integralmente, elementos essenciais do seu trabalho foram incorporados ao presente estudo.

A base utilizada nesta etapa foi o conjunto de sequências de otimização conhecido como *Optimization Cache*, disponibilizado pelos autores em seu repositório oficial. Trata-se de um conjunto pré-computado de sequências particularmente eficazes para redução de tamanho de código no LLVM, servindo como uma alternativa robusta às estratégias tradicionais, como a **-Oz**. Contudo, essas sequências foram originalmente elaboradas para a versão LLVM-10, enquanto esta pesquisa adotou a versão mais recente, LLVM-20.

Durante os testes iniciais observou-se que apenas cerca de 30% das sequências podiam ser aplicadas diretamente no LLVM-20. As demais falhavam devido a modificações internas do compilador e no ambiente LLVM. Para contornar esse problema, realizou-se uma conversão manual dos passes, mapeando os equivalentes entre as versões 10 e 20 do LLVM. Após esse trabalho de ajuste, tornou-se possível recuperar 100% da lista original, composta por **1289 sequências únicas**, agora compatíveis com a infraestrutura moderna do LLVM-20.

Com a lista completa e compatível, foi construído um processo de compilação exaustiva utilizando o script `rodar_seq_uniq_AVR.sh`. Esse script recebe como entrada o arquivo `sequencias_unicas.txt`, contendo todas as 1289 sequências convertidas, e

executa o seguinte procedimento para cada arquivo IR:

1. Aplicação dos passes selecionados ao arquivo LLVM-IR do código de entrada;
2. Compilação para assembly AVR;
3. Geração do arquivo objeto correspondente;
4. Extração e armazenamento do DEC (tamanho em bytes do código objeto);
5. Armazenamento dos artefatos referentes à melhor sequência parcial (arquivos `.s`, `.o` e `.ll`).

Ao término da execução do script para cada código, é gerado um arquivo `melhor.txt` contendo a sequência ótima responsável pelo menor tamanho de código objeto produzido.

Como o volume de códigos do dataset é significativo, a execução completa do processo em modo sequencial seria impraticável. Para lidar com essa limitação, desenvolveu-se o script `parallel-run.py`, responsável por paralelizar a execução do `rodar_seq_uniq_AVR.sh` utilizando nove *threads* da CPU. Esse mecanismo distribui simultaneamente diferentes códigos entre os núcleos disponíveis, reduzindo o tempo total de execução.

Por fim, empregou-se o script `coleta_melhores.py` para consolidar os resultados. Esse script percorre todas as pastas geradas durante a fase de otimização e extrai o arquivo `melhor.txt` de cada código analisado, compilando todos esses resultados em um único arquivo `melhores.csv`. Este arquivo constitui o conjunto final de sequências ótimas associado ao dataset e será utilizado nas etapas posteriores de formação dos pares *prompt/label* e preparação do conjunto de treinamento do modelo. Todos os códigos e arquivos gerados estão disponíveis no repositório: <<https://github.com/Alisson-Teles/Pipeline-training/tree/main/AVR>>.

## 4.5 Geração dos pares *prompt/label* e elaboração do *dataset* de treinamento

Com as etapas anteriores concluídas, a conversão para ROBL, a geração do LLVM-IR não otimizado e a identificação da melhor sequência de passes para cada programa, torna-se possível estruturar o conjunto de dados que será utilizado para o treinamento supervisionado do modelo de linguagem. Esta fase consiste em transformar cada código processado em um par composto por *prompt* e *label*, seguindo o padrão metodológico proposto no trabalho original do LLM-Compiler (??).

Após a execução do script `rodar_seq_uniq_AVR.sh`, cada programa possui três informações fundamentais: (i) o arquivo LLVM-IR resultante da compilação sem otimizações

expressivas, (ii) a melhor sequência de passes encontrada dentre as 1289 sequências testadas e (iii) o arquivo de assembly otimizado correspondente a essa sequência. Com esses elementos, torna-se possível construir pares de entrada e saída adequados ao processo de aprendizado supervisionado (*Supervised Fine-Tuning*).

Em cada instância do conjunto de dados, o *prompt* consiste em uma instrução textual em inglês seguida do código LLVM-IR não otimizado. A instrução foi definida conforme o padrão empregado por (??), mantendo a estrutura funcional do modelo original. A mensagem inicial é:

*“Tell me what passes to run on the following LLVM-IR to reduce the object file size for avr assembly.”*

Em seguida, o código LLVM-IR correspondente ao programa é anexado diretamente após essa instrução, compondo o conteúdo completo do *prompt*. Esse formato padronizado orienta o modelo a identificar a tarefa de seleção de passes como um problema de recomendação condicional baseado no IR de entrada.

O *label*, por sua vez, contém a resposta correta esperada para o exemplo. Ele segue o formato descritivo adotado no LLM-Compiler, sendo estruturado da seguinte forma:

*“Run the following passes <lista\_de\_passes> to reduce the object file size to <tamanho\_em\_bytes>.”*

Após essa mensagem inicial, inclui-se o código assembly gerado a partir da aplicação da sequência ótima de passes. O *label* engloba, tanto a justificativa textual quanto o resultado concreto da aplicação dos passes, replicando a estrutura esperada pelo treinamento do modelo original.

A geração automática desses pares foi realizada por meio do script `promptLabel.py`, disponível no repositório público do pipeline (<[https://github.com/Alisson-Teles/Pipeline-training/blob/main/AVR/PROMPT\\_LABEL/promptLabel.py](https://github.com/Alisson-Teles/Pipeline-training/blob/main/AVR/PROMPT_LABEL/promptLabel.py)>). O script percorre todos os diretórios contendo os arquivos LLVM-IR, as sequências ótimas (`melhor.txt`) e os arquivos de assembly otimizados, combinando essas informações e produzindo um arquivo final no formato JSONL. Cada linha do arquivo está mapeada por uma coluna "messages", seguido de "content", assim como esperado pelo serviço de *autotrain* do *hugging face*, e a linha representa um par de dados, contendo os campos:

- "prompt": instrução textual seguida do LLVM-IR não otimizado;
- "label": sequência ótima de passes, redução de tamanho e assembly correspondente;

Esse procedimento resulta em um dataset estruturado, padronizado e diretamente compatível com a etapa de *fine-tuning* do modelo LLM-Compiler, garantindo consistência com a metodologia original e adequação ao domínio específico da arquitetura AVR.

## 4.6 Configuração de treinamento do modelo

A etapa final do pipeline consiste na configuração do ambiente necessário para o refinamento do modelo LLM-Compiler (??), de modo a adaptá-lo ao domínio específico de otimização de código para microcontroladores AVR. Diferentemente do trabalho original, cujo treinamento foi conduzido em infraestrutura própria e com forte customização de hiperparâmetros, o presente trabalho utiliza recursos em nuvem da plataforma Hugging Face, adequados às limitações de escopo e hardware disponíveis.

Para esse fim, foi criado um ambiente dedicado na plataforma, acessível por meio do espaço público <<https://huggingface.co/spaces/Cal-mfbc5446/STF-PFC2>>. Esse ambiente foi configurado para possibilitar o *Supervised Fine-Tuning* (SFT) do modelo, utilizando uma GPU Nvidia A10G na configuração *large*, contendo 12 vCPUs, 46 GB de RAM e 24 GB de VRAM. A opção pela execução em nuvem se justifica pela necessidade de memória de vídeo compatível com janelas de contexto amplas, inviáveis em hardware local.

O processo de treinamento foi conduzido utilizando a ferramenta AutoTrain Advanced (??), que possibilita a realização de ajuste fino supervisionado sem a necessidade de desenvolvimento de código adicional, operando por meio de interface gráfica. Embora essa abordagem reduza a flexibilidade na configuração de hiperparâmetros, permite um fluxo de experimentação mais rápido e simplificado, adequado ao escopo deste trabalho.

Os parâmetros de treinamento empregados foram adaptados a partir das configurações originalmente utilizadas por (??). Entretanto, apenas parte das estratégias apresentadas no trabalho original pôde ser reproduzida neste ambiente. Em particular, foi possível utilizar:

- o *scheduler* do tipo Cosine, também presente no LLM-Compiler;
- o otimizador AdamW, sem a possibilidade de especificar manualmente os valores de  $\beta_1$  e  $\beta_2$ ;
- *weight decay* e taxa de aprendizado configuráveis dentro dos limites da ferramenta.

Além disso, devido às restrições de memória da GPU disponibilizada na plataforma, o parâmetro `model_max_length` teve de ser ajustado para 8 192 tokens, metade do utilizado no trabalho original (16 384 tokens). A redução foi necessária para garantir que o modelo pudesse ser carregado e treinado sem exceder a capacidade de VRAM.

Todos os hiperparâmetros utilizados no treinamento foram registrados no arquivo `parametros-de-treinamento-IA.txt`, disponível no repositório oficial do pipeline no GitHub: <<https://github.com/Alisson-Teles/Pipeline-training/tree/main/AVR>>. O arquivo documenta as configurações efetivamente aplicadas, incluindo divisões do conjunto de dados, *learning rate*, tamanho do lote e número de épocas planejado.

Para a composição do conjunto de treinamento, seguiu-se a divisão tradicionalmente adotada em tarefas de *fine-tuning* supervisionado: 80% dos pares *prompt/label* foram destinados ao treinamento, 10% ao conjunto de validação e 10% ao conjunto de teste. A validação permite monitorar sobreajuste, enquanto o conjunto de teste é reservado para a avaliação final do modelo refinado durante a etapa de resultados.

Embora o treinamento completo não tenha sido executado dentro do prazo deste trabalho, toda a infraestrutura necessária para sua realização foi preparada, incluindo: (i) configuração do ambiente computacional; (ii) definição dos hiperparâmetros; (iii) organização do dataset final; e (iv) preparação do pipeline para execução direta do SFT.

## 5 Resultados

A partir do pipeline implementado neste trabalho foi possível construir um conjunto de aproximadamente **1 000 pares *prompt/label*** destinado ao treinamento supervisionado do LLM-Compiler especializado para AVR. Cada *prompt* corresponde ao código em ROBL compilado para LLVM-IR (versão não otimizada, compilada com a flag `-O1`) e o respectivo *label* contém a melhor sequência de passes de otimização identificada para a plataforma alvo (ATmega328P) e o código assembly gerado após a aplicação desses passes. O conjunto foi produzido a partir de programas originalmente escritos em C, gerados por Csmith e coletados no AnghaBench, convertidos para ROBL, compilados com diferentes sequências de passes e avaliados quanto ao tamanho do objeto gerado para a plataforma AVR.

A seleção das melhores sequências para cada programa foi realizada utilizando a lista de sequências derivada do trabalho de (??), denominada de *Optimization Cache* em conjunto com a comparação contra flags tradicionais do LLVM (`-Oz`, `-Os`, `-O1`, `-O2`, `-O3`). Essa etapa foi automatizada pelos scripts desenvolvidos no repositório do pipeline (por exemplo, `rodar_seq_uniq_AVR.sh`, `parallel-run.py` e `coleta_melhores.py`), que percorrem o conjunto de 1.289 sequências únicas consideradas e registram, para cada programa, a sequência que produz o menor tamanho de objeto e o respectivo `.s` / `.o` / `.11` resultante. A comparação detalhada entre as sequências do Faustino e as flags padrão encontra-se documentada na planilha *Comparação Otimização Faustino vs Padrão* disponível no link: [do repositório](#).

Dessa forma, os resultados práticos deste trabalho incluem os scripts:

- O script `rodar_seq_uniq_AVR.sh` responsável por aplicar, para cada programa em LLVM-IR, todas as 1289

Além dos scripts, este trabalho também contribuiu com dados prontos para consumo, que é o caso dos arquivos:

- **dec\_otimizacoes\_GLOBAL.csv**, contendo o registro global dos resultados das compilações realizadas pelas flags padrões. Esse arquivo sintetiza, para cada programa e para cada flag, o tamanho do objeto gerado e permite validação e comparação de desempenho entre diferentes técnicas de otimização.
- **melhores.csv**, que consolida somente as melhores sequências descobertas para cada programa após a execução da metodologia exaustiva. Esse arquivo é fundamental para a etapa de construção dos pares *prompt/label*, pois fornece a sequência de passes selecionada como rótulo supervisionado do modelo.



- **pairs.jsonl**, que contém o dataset final que poderá ser utilizado no treinamento supervisionado. Cada linha contém um objeto JSON estruturado com o *prompt* (texto inicial + código LLVM-IR não otimizado) e o *label* (sequência ótima de passes + assembly otimizado). Esse arquivo segue o formato recomendado para pipelines do Hugging Face e serve como dataset de refinamento do LLM-COMPILER para AVR.
- **requirements.txt**, que especifica o ambiente mínimo necessário para executar os scripts do pipeline e manipular os dados gerados, garantindo reprodutibilidade e facilitando a replicação do experimento em outras máquinas ou em ambientes.

Adicionalmente

Com o conjunto final de pares pronto e a infraestrutura de treinamento configurada (ambiente no Hugging Face, espaço `Cal-mfbc5446/STF-PFC2` e parâmetros de treinamento registrados), foram conduzidos testes exploratórios de inferência utilizando o modelo refinado no domínio STM32 como referência de método (modelo refinado para outra MCU) e o pipeline AVR como alvo principal desta pesquisa. Embora o modelo originalmente refinado para STM32 ofereça indícios de que a técnica de fine-tuning é viável, a transferência direta entre domínios distintos (STM32  $\rightarrow$  AVR) mostra-se limitada: diferenças de ABI, conjunto de instruções e restrições de memória tornam a generalização direta impraticável sem adaptação específica do dataset. Assim, o resultado prático obtido nesta fase foi a validação do pipeline (geração, conversão, compilação, seleção e empacotamento dos pares) e a preparação de um dataset de tamanho suficiente para iniciar experimentos de fine-tuning direcionados ao ATmega328P.

Foram realizados testes exploratórios de inferência para avaliar qualitativamente a capacidade do pipeline e do processo de ajuste fino (quando aplicado). Em inferências preliminares — conduzidas sobre um subconjunto de programas selecionados dos benchmarks Csmith e AnghaBench — observou-se que o modelo apresenta capacidade de sugerir sequências estruturalmente compatíveis com aquelas observadas no conjunto de treinamento em casos pontuais, porém com generalização limitada em razão do número reduzido de exemplos (1 000 pares) e das diferenças intrínsecas entre programas. Em suma, não se reivindica aqui um resultado conclusivo de ganho consistente sobre -0z para AVR; o que se demonstra é a maturidade do pipeline e sua reprodutibilidade para executar experimentos de fine-tuning e avaliação quantitativa futura.

Portanto, o principal resultado deste trabalho é a implementação de um **pipeline reprodutível e público** que cobre toda a cadeia: (i) geração e seleção de programas em C; (ii) conversão automatizada para ROBL; (iii) compilação com listas exaustivas de passes (1.289 sequências únicas) e comparação com flags padrão; (iv) coleta das melhores sequências e montagem dos pares *prompt/label*; e (v) preparação do ambiente de fine-tuning. A análise comparativa entre as sequências do Faustino e as flags padrão (arquivo

de suporte) corrobora a relevância de usar conjuntos de sequências especializados como ponto de partida para a construção de pares de treinamento e é disponibilizada como material de apoio ao presente trabalho. :contentReference[oaicite:1]index=1

## 6 Discussão

Durante o desenvolvimento deste trabalho, observou-se uma série de limitações que impactaram diretamente o escopo e os resultados alcançados. Os principais desafios enfrentados foram decorrentes de restrições de tempo, orçamento e gastos computacionais. Estes problemas acarretaram em uma série de mudanças no planejamento inicial, que influenciaram diretamente o escopo deste trabalho. Diversas mudanças foram aplicadas no decorrer desta pesquisa, a mais significativa delas foi que o refinamento, e consequentemente a inferência e validação, do modelo LLM-Compiler para AVR não poderia ser concluído dentro do prazo estipulado. Desta forma, a contribuição principal deste trabalho é o desenvolvimnto prático e teórico de um pipeline capaz de fornecer toda a base necessária para o refinamento futuro do modelo.

Outros problemas relevantes que causaram uma mudança no escopo deste projeto foi o custo computacional associado à execução do pipeline completo, em especial a etapa de autotuning. A execução do método completo, de elaborado por (??) demandaria tempo e adicionaria complexidade a um projeto já limitado por prazos e recursos. Em função dessas limitações, optou-se por restringir o autotuning a uma lista predefinida de sequências a *Optimization Cache* derivada de (??) e por processar um subconjunto controlado de programas priorizando a reprodutibilidade do pipeline, (da geração à produção dos pares *prompt/label*) em detrimento da exploração exaustiva do espaço combinatório de passes.

Outra limitação relevante refere-se ao tamanho dos exemplos em termos de tokens e à janela de contexto disponível para o fine-tuning. Em razão da GPU utilizada (Nvidia A10G alugada via Hugging Face), o parâmetro `model_max_length` foi ajustado para 8192 tokens. Vários exemplos gerados pelo pipeline excederam esse limite, exigindo truncamento para viabilizar o treinamento. O truncamento potencialmente remove contexto semântico importante para a escolha adequada de passes, prejudicando a qualidade do aprendizado. Porém, foi escolhido preservar os prompts completos no repositório, pensando em infraestruturas futuras que permitam janelas de contexto maiores.

Outros problemas técnicos superados, foi a conversão e mapeamento manual dos passes da lista *Optimization Cache*, originalmente definidos para LLVM-10, para torná-los utilizáveis no LLVM-20. Essa etapa de adaptação introduziu um esforço adicional para garantir que todas as sequências pudessem ser aplicadas corretamente, evitando falhas de compilação e garantindo a integridade do processo de autotuning.

Por fim, apesar das limitações supracitadas, o principal mérito prático do trabalho é a entrega de um pipeline reprodutível e público que cobre toda a cadeia de experimentação (geração de código, conversão para ROBL, aplicação e avaliação de 1 289 sequências únicas,

coleta das melhores sequências, montagem dos pares *prompt/label* e configuração de um ambiente e parâmetros de refinamento de IA), necessários para o fine-tuning futuro do LLM-Compiler no domínio AVR.

## 7 Conclusão e Trabalhos Futuros