

# Otimização de Desempenho para Sistemas Lineares Esparsos com Pré-condicionantes

SERGIO SIVONEI DE SANT'ANA FILHO GRR20242337  
EDUARDO KALUF GRR20241770

3 de dezembro de 2025

## 1 Introdução

Este trabalho visa realizar principalmente a otimização de operações computacionais chave dentro do algoritmo do Método dos Gradientes Conjugados (CG), visando calcular de forma mais eficiente a solução de sistemas lineares da forma  $Ax = b$ .

Especificamente, busca-se aprimorar e avaliar o desempenho do programa computacional desenvolvido no 1º Trabalho Prático.

As operações críticas que foram alvo de otimização são:

- A Iteração Central do Método de Gradiente Conjugado utilizando o Pré-condicionador de Jacobi (PCG).
- O Cálculo do Resíduo do sistema.

A motivação para estas otimizações reside na natureza intensiva em computação dessas operações, bem como no grande tráfego de memória, especialmente ao lidar com matrizes esparsas de grande dimensão.

## 2 Otimizações Implementadas

Nesta seção, detalhamos as técnicas de otimização aplicadas às operações, descrevendo seu impacto teórico e prático no desempenho, bem como os desafios encontrados durante o processo de implementação e ajuste.

### 2.1 Otimização do Gradiente Conjugado (PCG)

A otimização da iteração do Gradiente Conjugado, que envolve múltiplas operações de produto matriz-vetor ( $A \times v$ ), foi realizada primariamente através da alteração do formato de armazenamento da matriz do sistema.

#### 2.1.1 Matriz CSR (Compressed Sparse Row)

O formato **CSR** (Compressed Sparse Row) foi a principal modificação estrutural implementada.

Este foi o terceiro método de armazenar a matriz que experimentamos, primeiramente tentamos guardá-las em um único vetor sem offsets e depois tentamos em diversos vetores de tamanhos diferentes. Desistimos de ambas as implementações, pois não pareciam práticas para realizar as operações e não ofertavam nenhum grande benefício aparente além da diminuição do tanto de memória utilizada.

O método escolhido é um esquema de armazenamento de matrizes projetado especificamente para lidar com matrizes esparsas, ou seja, matrizes que possuem uma alta proporção de elementos nulos.

Este formato elimina a necessidade de armazenar informações redundantes (os zeros), operando sobre apenas três vetores de dimensão reduzida:

- **'Values'**: armazena apenas os valores não nulos da matriz.
- **'Col\_indices'**: armazena os índices das colunas correspondentes a cada valor não nulo.
- **'Row\_pointers'**: armazena os índices que indicam o início de cada nova linha nos vetores 'values' e 'col\_indices'.

A principal vantagem do **CSR** é sua extrema eficiência para operações de multiplicação matriz-vetor ( $A \times v$ ).

Dada a estrutura da iteração do Gradiente Conjugado (PCG), que exige duas ou mais dessas multiplicações por passo, o uso do **CSR** é crucial para reduzir o tempo de execução. Esta abordagem melhora a localidade de dados

e permite **acessos sequenciais** à memória, que são mais rápidos e otimizam a utilização da cache.

**Vetorialização (Uso de AVX):** Além disso, devido ao modo como o armazenamento **CSR** funciona, foi possível utilizar **AVX** neste método, deixando a iteração do método ainda mais rápida. A vetorização utiliza as Extensões Avançadas de Vetor (**AVX**) do processador, permitindo que uma única instrução da CPU (SIMD - Single Instruction, Multiple Data) opere simultaneamente sobre múltiplos dados. Isso foi essencial para **diminuir drasticamente** o tempo de execução das operações vetoriais, explorando a capacidade de paralelismo ao nível de hardware.

## 2.2 Otimização do Cálculo do Resíduo

O cálculo do resíduo  $r = b - Ax$  requer primariamente uma multiplicação matriz-vetor e uma subtração de vetores. A otimização focou em garantir a máxima eficiência nessas operações, complementando a desempenho obtida com a implementação do **CSR** para multiplicação entre matriz vetor.

Foram feitas algumas pequenas alterações de otimização, porém o maior ganho de desempenho foi devido à estrutura **CSR**. Houve também a tentativa de implementar o *loop unroll & jam*, porém sem sucesso, pois os testes mostraram uma piora significativa no tempo de execução do cálculo.

Outra otimização foi realizar parte do cálculo da norma dentro do laço que calcula o vetor resíduo, assim aproveitando um valor que já está em memória sem a necessidade de outro loop que iria acessar todas as posições novamente. Todavia, nos testes realizados, não foi observado grande mudança, na maioria apenas pequenas variações de tempo.

## 2.3 Otimização Adicional: Matrizes Simétricas Positivas

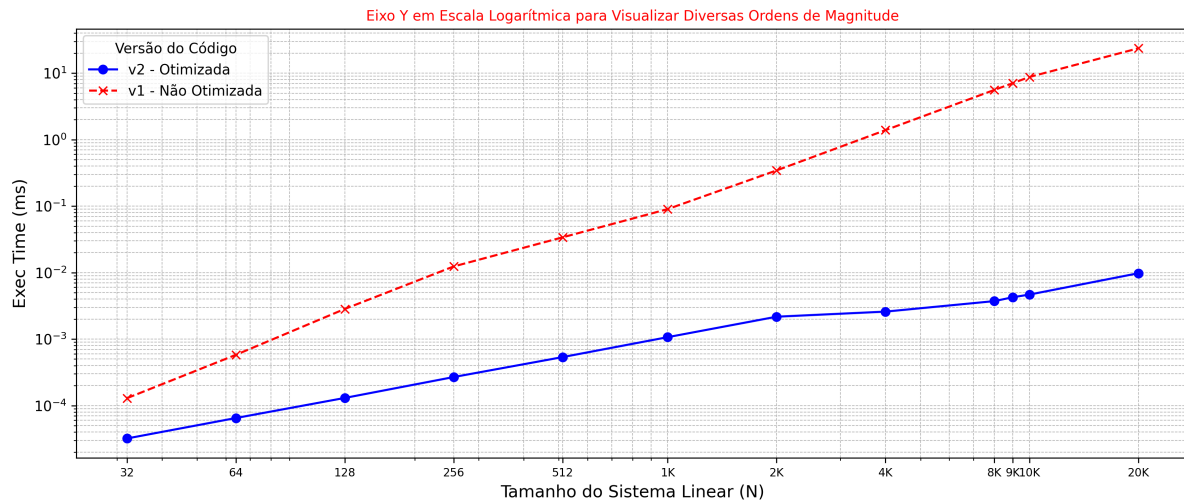
Além das duas operações fundamentais também otimizamos como a matriz simétrica positiva é calculada. O método **CSR** é ruim quando se é necessário percorrer pelas colunas, sendo assim, a multiplicação de duas matrizes do jeito convencional ficaria mais lento que o esperado, a fim de contornar isso, fizemos a multiplicação entre a matriz principal e a transposta pelas linhas, ou seja, ao invés de multiplicar linhas por colunas, multiplicamos linhas por linhas. Além disso, o método apenas itera quando existe algum resultado possível, multiplicações que irão dar zero no final são ignorados, pois não são necessárias.

### 3 Análise dos Gráficos

A seguir, seguem os gráficos com as métricas extraídas pelo programa **LIKWID**, conforme solicitado, para diversos valores de N.

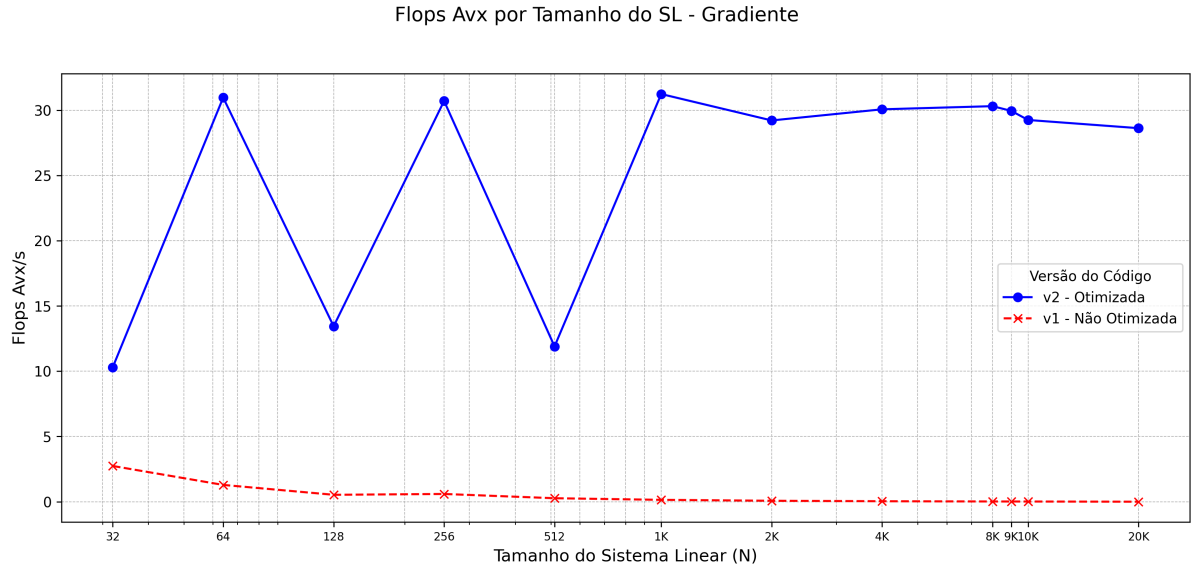
#### 3.1 Método dos Gradientes Conjugados (CG)

Exec Time por Tamanho do SL - Gradiente



É possível observar a disparidade no tempo de execução entre o programa antigo e o novo. Torna-se evidente que os métodos de otimização utilizados obtiveram sucesso ao reduzir drasticamente o tempo de execução para a iteração do gradiente, especialmente para grandes valores de N, que é o foco principal do nosso interesse.

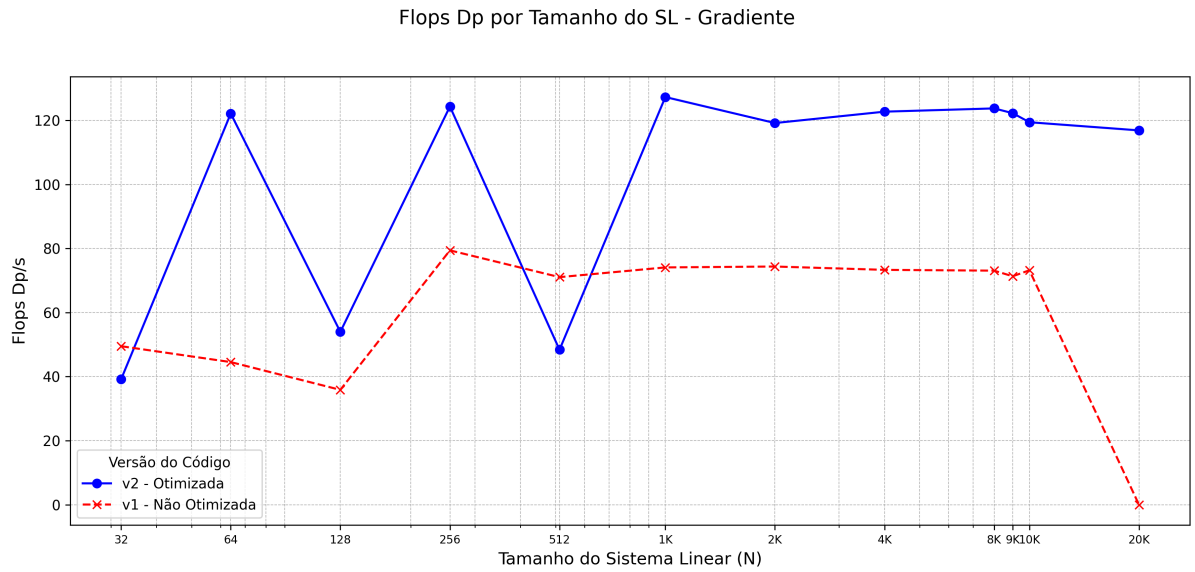
Essa redução no tempo é justificada pela melhora nas métricas que analisaremos a seguir.



Ao analisarmos os Flops **AVX** podemos concluir que houve um uso das Extensões Avançadas de Vetor (**AVX**) em maior escala, coisa que otimiza o programa ao fazer uso de instruções SIMD.

A utilização do **AVX** foi possível pelo método como estamos armazenando a matriz, o **CSR** permite que diversas operações durante a multiplicação de matriz vetor sejam realizadas ao mesmo tempo, pois são realizadas em sequência, sem dependência de dados e sem nenhum desvio condicional.

Pode-se inferir que os picos e vales que ocorrem nas potências de 2 no início do gráfico resultam de um mapeamento ineficiente da cache para alguns desses valores específicos. Esse fenômeno é comumente denominado *cache thrashing*.

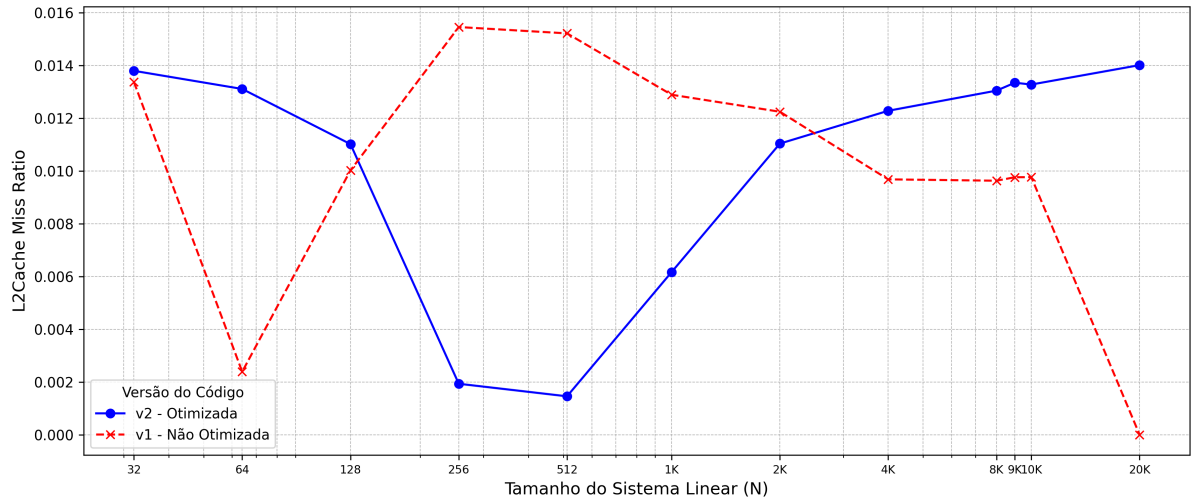


Houve também um aumento no número de **FLOPs de precisão dupla (Double precision)** por segundo, o que impacta diretamente na velocidade do programa. Uma vez que um maior volume de operações ocorre no mesmo intervalo de tempo, isso significa que a máquina concluirá os cálculos mais rapidamente em comparação com a versão anterior.

É importante ressaltar que a taxa de FLOPs por segundo é inversamente proporcional ao tempo de execução do programa. Ou seja, o aumento dessa métrica ocorreu tanto por uma diminuição no tempo total necessário para realizar as operações quanto por uma diminuição no total de operações realizadas.

Esta métrica apresenta os mesmos picos e vales que o gráfico anterior, pois o mesmo problema ocorre.

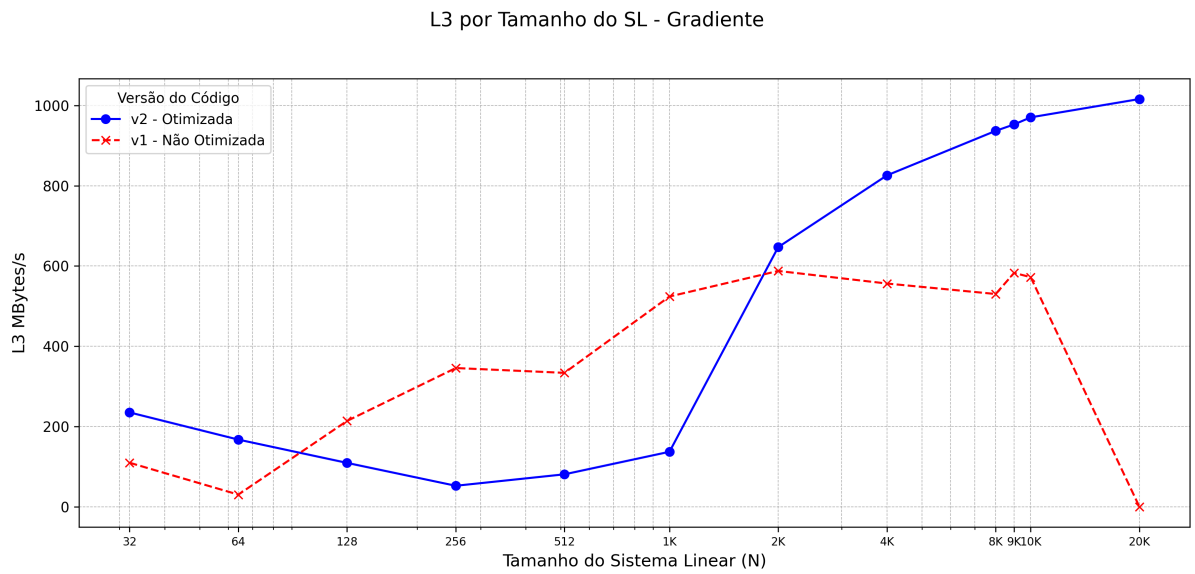
L2Cache por Tamanho do SL - Gradiente



Analisando o gráfico de miss ratio da cache L2 é possível inferir que: para sistemas de tamanhos condizentes com a Cache L2 (256-512) temos um desempenho expressivo da versão otimizada. Todavia, quando o tamanho aumenta muito temos uma piora por não ser mais possível aproveitar a localidade na cache, o que chega deixar o miss ratio da versão otimizada levemente pior que da versão original.

Além disso, pode-se afirmar que a primeira versão do trabalho tem um desempenho mais consistente, apesar de ruim. Coisa que não se confere na nova versão.

O cache miss ratio é dado pela divisão entre os acessos à cache e a quantidade de erros, sendo assim temos uma métrica em proporção e não em valores absolutos.



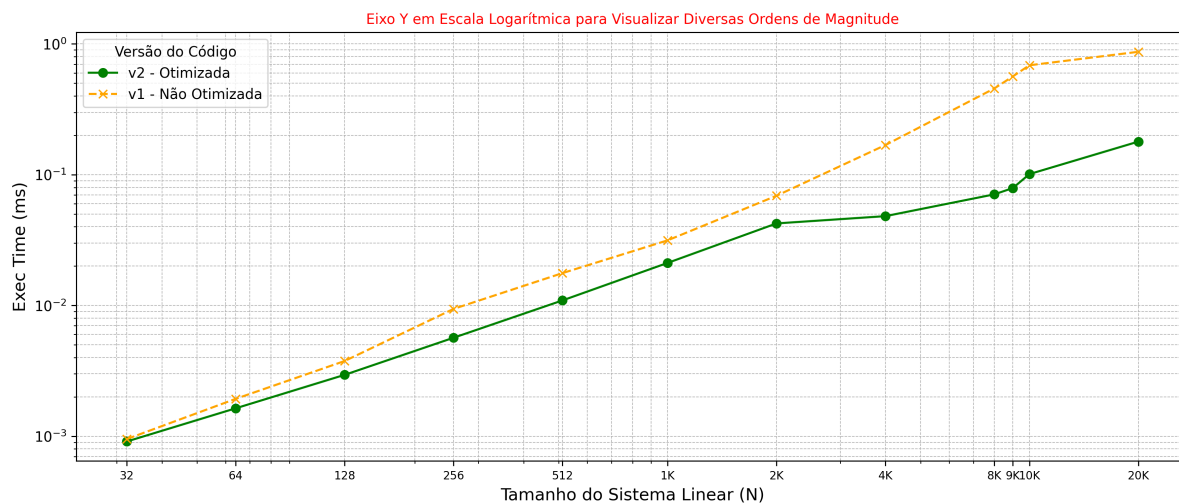
Aqui, podemos observar um uso muito melhor da comunicação entre processador e memória, visto que a largura de banda aumenta notavelmente na versão otimizada. Isso indica que mais dados são transmitidos por segundo, otimizando o acesso à memória.

Tal melhoria ocorre porque, com o formato CSR, os dados utilizados ficam contíguos na memória, aumentando a localidade espacial e, consequentemente, a quantidade de dados trafegados



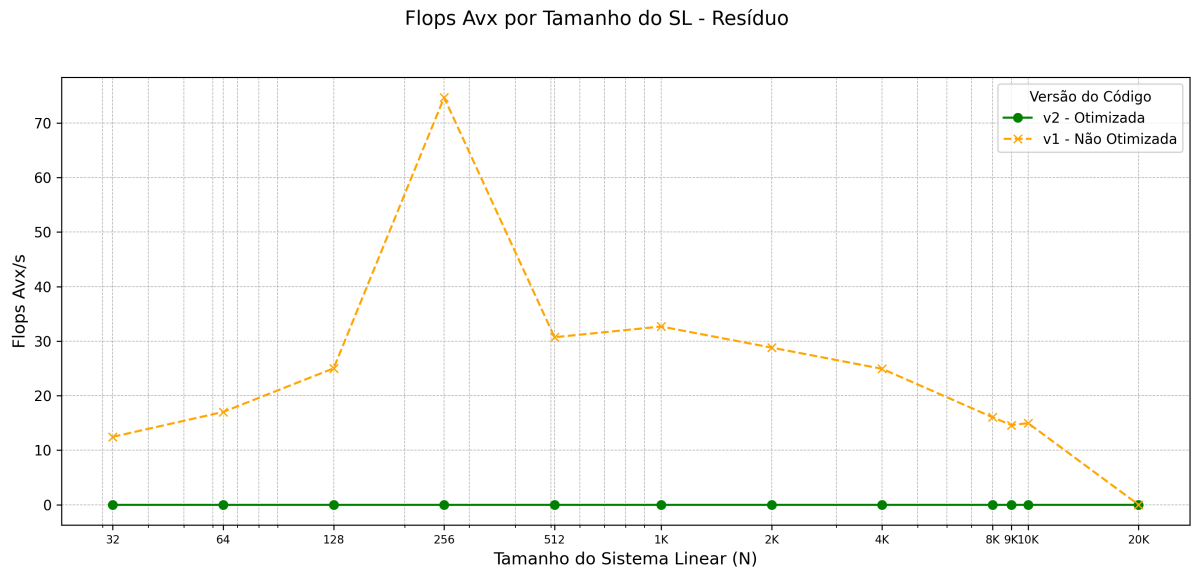
### 3.2 Cálculo do Resíduo

Exec Time por Tamanho do SL - Resíduo



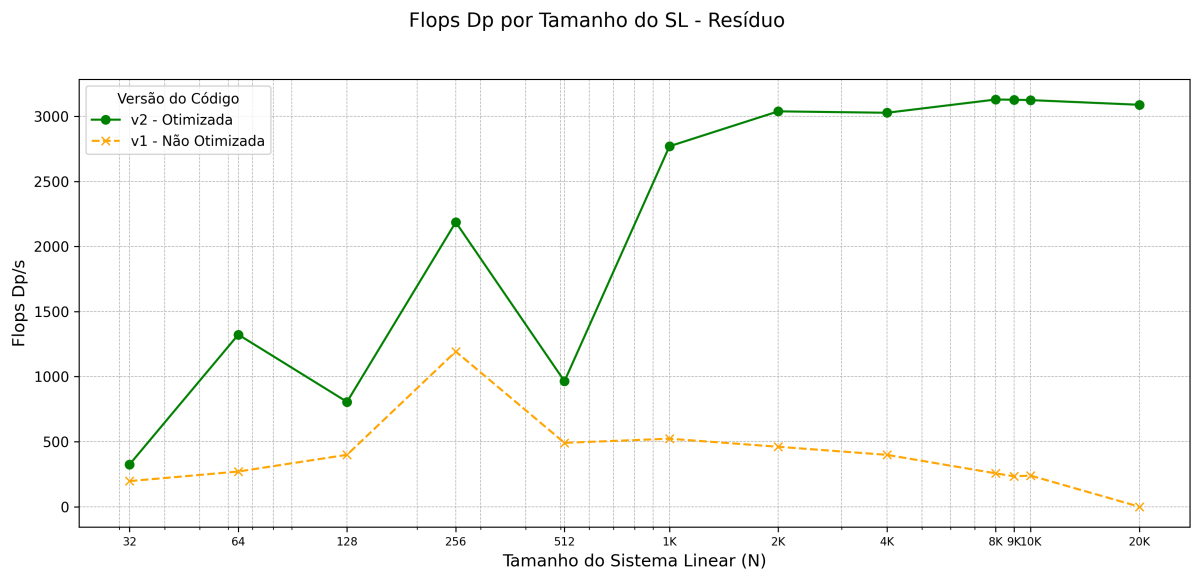
No gráfico acima percebemos uma melhora relativa entre a primeira versão (não otimizada) e a segunda (otimizada), que se mostra mais significativa em tamanhos de larga escala.

Tais melhoras podem ser justificadas pelos gráficos a seguir, bem como pelo uso do modelo **CSR**.



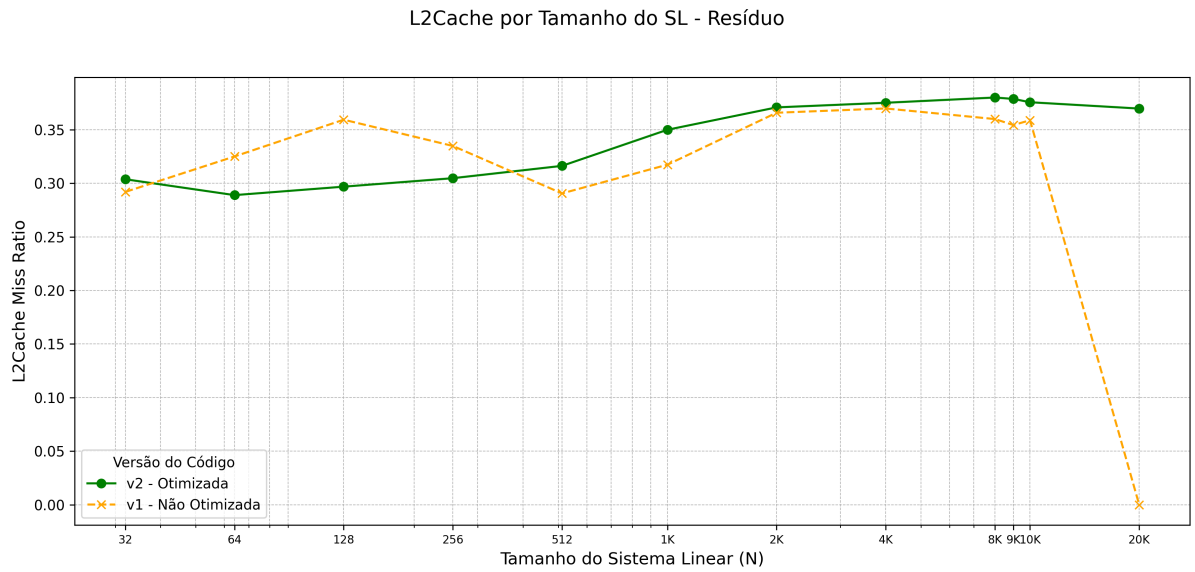
No que diz respeito ao uso de Flops **AVX** no cálculo do resíduo temos que na nova versão otimizada não houve o uso da extensão em nenhum momento, coisa que após analisarmos o código do programa e as métricas, consideramos ser devido ao uso do **CSR**, que acaba nos obrigando a fazer o acesso indireto de um vetor, isso impede o **AVX** de ser executado.

Ademais, existe dependência de dados na variável *sum*, outro impecilho ao **AVX**.



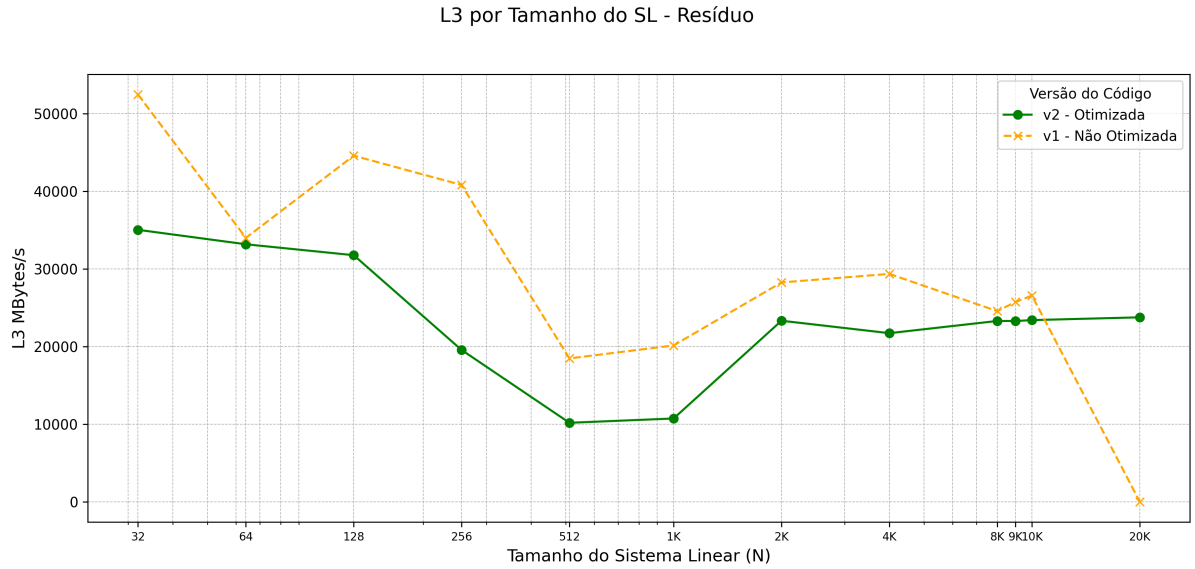
Ao contrário do que foi visto no último gráfico, podemos perceber aqui uma impressionante melhora, ao passo que realiza muitos mais cálculos em uma mesma quantidade de tempo, acelerando o programa.

O funcionamento e a explicação deste gráfico seguem igualmente à do **CG**.



Neste ponto é possível notar que não houve melhora no *Miss Ratio* da cache **L2**, até mesmo representando uma pequena piora em parte dos casos.

Essa maior *taxa de miss* ocorre, principalmente, pelos acessos indiretos e distantes que o cálculo do resíduo acaba fazendo por conta do **CSR**.



A **bandwidth** da cache **L3** apresentou uma piora, pois menos dados estão sendo carregados por segundo, mesmo que a diferença seja, geralmente, pequena. Isso ocorre, de fato, devido à natureza dos acessos indiretos à memória e ao modo de funcionamento do formato **CSR**. No entanto, como mencionado anteriormente, essa piora é sobreposta pelo ganho nas outras operações e, mais ainda, compensada pelo tempo economizado nas iterações do **CG**.

### 3.3 Observações

É importante notar que, ainda que alguns gráficos mostrem um resíduo ligeiramente pior, o ganho obtido nas demais métricas de desempenho é mais significativo. Adicionalmente, o principal gargalo do programa encontra-se na iteração do **CG**.

Portanto, é altamente vantajoso comprometer uma pequena porção do desempenho no cálculo do resíduo, visando a implementação de uma estrutura que otimize e acelere a iteração do gradiente.

## 4 Resultados e Análise de Desempenho

Os resultados demonstram que a combinação do armazenamento **CSR** com a vetorização **AVX** produziu uma redução significativa no tempo de execução total, especialmente em sistemas de alta dimensão e esparsidade.

A Tabela 1 resume o ganho de desempenho (speedup) alcançado para as operações críticas:

Tabela 1: Comparação de Tempo de Execução: Original vs. Otimizado (Tempo para  $N=20.000$  em milisegundos)

Operação	Tempo Original	Tempo Otimizado	Speedup
Iteração do PCG	23.60	0.0097	2432.98x
Cálculo do Resíduo	0.87	0.17	5.11x

Vale ressaltar que no teste de  $N = 20.000$  do 1º trabalho, o **LIKWID** não gerou nenhum dado correspondente a execução, portanto apenas a medida de tempo deve ser considerada, já que todas as outras são iguais a zero.

## 5 Estimativa de Desempenho

Baseado nas modificações implementadas na segunda versão e considerando uma máquina com 8GB de memória RAM disponível para seu programa e Sistemas com 7 (sete) diagonais, pode-se responder às perguntas dadas no enunciado do trabalho.

1. Qual o valor máximo estimado da ordem  $N$  do sistema linear que você seria capaz de resolver?

*Após breves discussões, concluímos que conseguiríamos resolver até um sistema de ordem 12. Apesar de demandar muito tempo e esforço, parece uma métrica realizável em algumas horas.*

2. Qual o ganho estimado de tempo para uma iteração do método de Gradiente Conjugado com Pré-condicionador de Jacobi (em função de  $N$ )?

*Estimamos que o ganho de desempenho seria algo dentro o intervalo de 8 a 16 vezes mais rápido. Pois, acreditávamos que isso já seria um ganho notável de desempenho.*

3. Qual a estimativa para a quantidade de operações em ponto flutuante executadas em cada iteração do método de Gradiente Conjugado com Pré-condicionador de Jacobi (em função de  $N$ )?

*Estimamos que teríamos em torno de 100 Flops Dp por segundo no método de Gradiente Conjugado, coisa que já seria melhor que no primeiro trabalho.*

4. Compare suas estimativas acima com os resultados obtidos nos testes.

*Bom, apenas as perguntas '2.' e '3.' podem ser comparadas, visto que a primeira não trata sobre o programa/otimizações em si.*

*Sobre o ganho estimado (2.); percebemos que a diferença foi muito maior, chegando a casos da segunda versão apresentar um tempo de execução 2458 vezes menor que a versão 1.*

*Já sobre a quantidade de operações de ponto flutuante (3.); obtivemos 120 Flops Dp em média nos testes, bem como tivemos um expressivo aumento em Flops AVX.*

## 6 Conclusão

Por fim, nota-se que as operações críticas do método dos Gradientes Conjugados foram otimizadas com sucesso. Utilizando métodos estruturais (como a matriz **CSR** para matrizes esparsas) e técnicas de paralelismo ao nível de instrução (como a **vetorização AVX**), alcançamos uma melhoria substancial no desempenho do programa. Conseguimos confirmar a execução otimizada pelo **AVX** em pelo menos uma das operações e, como resultado, **diminuímos drasticamente** o tempo de execução total do algoritmo, tornando o Método dos Gradientes Conjugados mais eficiente para a resolução de sistemas lineares de grande escala.



## Referências

- [1] CENAPAD-SP. *Engineering and Scientific Subroutine Library Guide and Reference: Subroutine Reference SDRSBC*, 2025.
- [2] Ed F. D’Azevedo. Vectorized Sparse Matrix Multiply for Compressed Row Storage Format, 2025.
- [3] Lei Mao. CSR Sparse Matrix Multiplication, 2019.
- [4] Netlib. Parallel Sparse Matrix-Vector Multiplication, 2025.
- [5] Wikipedia. Sparse matrix: Compressed sparse row (CSR, CRS or Yale format), 2025.