# Optimized GPU-Based Matrix Inversion

## Through The Use of Thread-Data Remapping

Samuel J Monson

Seattle Univerisity

2024-05-31

# INTRODUCTION

- cpu-inverse
  - ▶ Basic algorithm implementation to prove validity.

- **cpu-inverse**
  - ▶ Basic algorithm implementation to prove validity.
- **inverse**
  - ▶ GPU implementation written in CUDA.
  - ▶ Based on the works of Sharma et al [1] and DasGupta [2].

- **cpu-inverse**
  - Basic algorithm implementation to prove validity.
- **inverse**
  - GPU implementation written in CUDA.
  - Based on the works of Sharma et al [1] and DasGupta [2].
- **tdr-inverse**
  - Utilizes Thread-Data Remapping (TDR) to more efficiently use the GPU.
  - Actual TDR implementation is Dr. Cuneo's Harmonize library [3].

# Introduction to Inverses

## Definition (Inverse)

The <u>inverse</u> of $a$ is some value $a^{-1}$ such that $a \star a^{-1} = i$ where $i$ is the identity of $\star$.

## Definition (Inverse)

The <u>inverse</u> of $a$ is some value $a^{-1}$ such that $a \star a^{-1} = \boldsymbol{i}$ where $\boldsymbol{i}$ is the identity of $\star$.

## Definition (Identity)

The <u>identity</u> of an operation $\star$ and set $G$ is some value $\boldsymbol{i} \in G$ where for all $a \in G$, $\boldsymbol{i} \star a = a \star \boldsymbol{i} = a$.

## Definition (Inverse)

The <u>inverse</u> of $a$ is some value $a^{-1}$ such that $a \star a^{-1} = i$ where $i$ is the identity of $\star$.

## Definition (Identity)

The <u>identity</u> of an operation $\star$ and set $G$ is some value $i \in G$ where for all $a \in G$, $i \star a = a \star i = a$.

- For example, the set and operation $(\mathbb{R}, \times)$ has the identity $i = 1$ since $1 \times x = x \times 1 = x$ for all $x \in \mathbb{R}$.
- Thus the inverse of $a$ is $\frac{1}{a}$ since $a \times a^{-1} = 1 \rightarrow a = \frac{1}{a}$.
  - Note that this is only true because $a \times b = b \times a$ for all $a, b \in \mathbb{R}$.

$$3x + 2y = 2$$
$$-7x - 5y = 4$$

$$\begin{bmatrix} 3 & 2 \\ -7 & -5 \end{bmatrix}\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

$$\begin{bmatrix} 3 & 2 \\ -7 & -5 \end{bmatrix}^{-1}\begin{bmatrix} 3 & 2 \\ -7 & -5 \end{bmatrix}\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 3 & 2 \\ -7 & -5 \end{bmatrix}^{-1}\begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

$$I_2\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 5 & 2 \\ -7 & -3 \end{bmatrix}\begin{bmatrix} 2 \\ 4 \end{bmatrix}$$

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 18 \\ -26 \end{bmatrix}$$

# Methods of Matrix Inversion

### Definition (Matrix Inverse)

Let $A$ be an $n \times n$ matrix,

$$AA^{-1} = I = A^{-1}A$$

$$A^{-1} = \frac{1}{\det(A)}X$$

## Definition (Matrix Inverse)

Let $A$ be an $n \times n$ matrix,
$$AA^{-1} = I = A^{-1}A$$

$$A^{-1} = \frac{1}{\det(A)}X$$

## Determinate: $\det(\mathbb{R}^{n \times n}) \to \mathbb{R}$

$$\det\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}\right) = ad - bc$$

## Definition (Matrix Inverse)

Let $A$ be an $n \times n$ matrix,
$$AA^{-1} = I = A^{-1}A$$

$$A^{-1} = \frac{1}{\det(A)}X$$

## Determinate: $\det(\mathbb{R}^{n \times n}) \to \mathbb{R}$

$$\det\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}\right) = ad - bc$$

## Value of X

$$X = \det(A) \cdot A^{-1}$$

# General Definition

## Definition (Matrix Inverse)

Let $A$ be an $n \times n$ matrix,
$$AA^{-1} = I = A^{-1}A$$

$$A^{-1} = \frac{1}{\det(A)}X$$

## Determinate: $\det(\mathbb{R}^{n \times n}) \to \mathbb{R}$

$$\det\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}\right) = ad - bc$$

## Value of X

$$X = \det(A) \cdot A^{-1}$$

$$A \cdot X = \det(A) \cdot I$$

# General Definition

## Definition (Matrix Inverse)

Let $A$ be an $n \times n$ matrix,
$$AA^{-1} = I = A^{-1}A$$

$$A^{-1} = \frac{1}{\det(A)}X$$

## Determinate: $\det(\mathbb{R}^{n \times n}) \to \mathbb{R}$

$$\det\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}\right) = ad - bc$$

## Value of X

$$X = \det(A) \cdot A^{-1}$$

$$A \cdot X = \det(A) \cdot I$$

$$= \begin{bmatrix} \det(A) & 0 \\ 0 & \det(A) \end{bmatrix}$$

# General Definition

## Definition (Matrix Inverse)

Let $A$ be an $n \times n$ matrix,
$$AA^{-1} = I = A^{-1}A$$

$$A^{-1} = \frac{1}{\det(A)}X$$

## Determinate: $\det(\mathbb{R}^{n \times n}) \to \mathbb{R}$

$$\det\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}\right) = ad - bc$$

## Value of X

$$X = \det(A) \cdot A^{-1}$$

$$A \cdot X = \det(A) \cdot I$$

$$= \begin{bmatrix} \det(A) & 0 \\ 0 & \det(A) \end{bmatrix}$$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}\begin{bmatrix} x_1 & x_2 \\ x_3 & x_4 \end{bmatrix} = \begin{bmatrix} ad - bc & 0 \\ 0 & ad - bc \end{bmatrix}$$

## Definition (Matrix Inverse)

Let $A$ be an $n \times n$ matrix,

$$AA^{-1} = I = A^{-1}A$$

$$A^{-1} = \frac{1}{\det(A)}X$$

## Determinate: $\det(\mathbb{R}^{n\times n}) \to \mathbb{R}$

$$\det\left(\begin{bmatrix} a & b \\ c & d \end{bmatrix}\right) = ad - bc$$

## Value of X

$$X = \det(A) \cdot A^{-1}$$

$$A \cdot X = \det(A) \cdot I$$

$$= \begin{bmatrix} \det(A) & 0 \\ 0 & \det(A) \end{bmatrix}$$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}\begin{bmatrix} x_1 & x_2 \\ x_3 & x_4 \end{bmatrix} = \begin{bmatrix} ad - bc & 0 \\ 0 & ad - bc \end{bmatrix}$$

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}\begin{bmatrix} d & -b \\ -c & a \end{bmatrix} = \begin{bmatrix} ad - bc & 0 \\ 0 & ad - bc \end{bmatrix}$$

## Augmented Matrix

$$\text{aug}\,(A) = A|I = \begin{bmatrix} a & b & 1 & 0 \\ c & d & 0 & 1 \end{bmatrix}$$

# Gauss-Jordan Method

## Augmented Matrix

$$\text{aug}(A) = A|\boldsymbol{I} = \begin{bmatrix} a & b & 1 & 0 \\ c & d & 0 & 1 \end{bmatrix}$$

## Operations

1. Swap any two rows: $swap(R_i, R_j)$
2. Multiply any row by a non-zero value: $c \times R_i$
3. Add to any row a multiple of another row: $R_i + c \times R_j$

## Augmented Matrix

$$\operatorname{aug}(A) = A|I = \begin{bmatrix} a & b & | & 1 & 0 \\ c & d & | & 0 & 1 \end{bmatrix}$$

## Operations

1. Swap any two rows: $swap(R_i, R_j)$
2. Multiply any row by a non-zero value: $c \times R_i$
3. Add to any row a multiple of another row: $R_i + c \times R_j$

$$\begin{bmatrix} 3 & 2 & | & 1 & 0 \\ -7 & -5 & | & 0 & 1 \end{bmatrix} \xrightarrow{R_1 \leftrightarrow R_2} \begin{bmatrix} -7 & -5 & | & 0 & 1 \\ 3 & 2 & | & 1 & 0 \end{bmatrix}$$

# Gauss-Jordan Method

## Augmented Matrix

$$\operatorname{aug}(A) = A|\boldsymbol{I} = \begin{bmatrix} a & b & 1 & 0 \\ c & d & 0 & 1 \end{bmatrix}$$

## Operations

1. Swap any two rows: $swap(R_i, R_j)$
2. Multiply any row by a non-zero value: $c \times R_i$
3. Add to any row a multiple of another row: $R_i + c \times R_j$

$$\begin{bmatrix} 3 & 2 & 1 & 0 \\ -7 & -5 & 0 & 1 \end{bmatrix} \xrightarrow{R_1 \leftrightarrow R_2} \begin{bmatrix} -7 & -5 & 0 & 1 \\ 3 & 2 & 1 & 0 \end{bmatrix} \xrightarrow{R_2 \times 2} \begin{bmatrix} -7 & -5 & 0 & 1 \\ 6 & 4 & 2 & 0 \end{bmatrix}$$

# GAUSS-JORDAN METHOD

## Augmented Matrix

$$\text{aug}(A) = A|\boldsymbol{I} = \begin{bmatrix} a & b & | & 1 & 0 \\ c & d & | & 0 & 1 \end{bmatrix}$$

## Operations

1. Swap any two rows: $swap(R_i, R_j)$
2. Multiply any row by a non-zero value: $c \times R_i$
3. Add to any row a multiple of another row: $R_i + c \times R_j$

$$\begin{bmatrix} 3 & 2 & | & 1 & 0 \\ -7 & -5 & | & 0 & 1 \end{bmatrix} \xrightarrow{R_1 \leftrightarrow R_2} \begin{bmatrix} -7 & -5 & | & 0 & 1 \\ 3 & 2 & | & 1 & 0 \end{bmatrix} \xrightarrow{R_2 \times 2} \begin{bmatrix} -7 & -5 & | & 0 & 1 \\ 6 & 4 & | & 2 & 0 \end{bmatrix} \xrightarrow{R_1 + ^1/_2 R_2} \begin{bmatrix} -4 & -3 & | & -1 & 1 \\ 6 & 4 & | & 2 & 0 \end{bmatrix}$$

**Data:** $M$ is a matrix with $N$ rows
**foreach** *row $M_i$* **do**
    $M_i \leftarrow M_i / M_{ii}$;
    **foreach** *row $M_j$ in M where $j \neq i$* **do**
        $M_j \leftarrow M_j - M_{ji} \times M_i$
    **end**
**end**

**Data:** $M$ is a matrix with $N$ rows
**foreach** *row $M_i$* **do**
 $\quad M_i \leftarrow M_i / M_{ii}$;
 $\quad$ **foreach** *row $M_j$ in M where $j \neq i$* **do**
 $\quad\quad M_j \leftarrow M_j - M_{ji} \times M_i$
 $\quad$ **end**
**end**

$$\begin{bmatrix} 3 & 2 & | & 1 & 0 \\ -7 & -5 & | & 0 & 1 \end{bmatrix}$$

**Data:** $M$ is a matrix with $N$ rows
**foreach** *row $M_i$* **do**
    $M_i \leftarrow M_i / M_{ii}$;
    **foreach** *row $M_j$ in M where $j \neq i$* **do**
        $M_j \leftarrow M_j - M_{ji} \times M_i$
    **end**
**end**

$$\begin{bmatrix} 3 & 2 & 1 & 0 \\ -7 & -5 & 0 & 1 \end{bmatrix} \xrightarrow{R_1/3} \begin{bmatrix} \mathbf{1} & {}^{2}\!/_{3} & {}^{1}\!/_{3} & \mathbf{0} \\ -7 & -5 & 0 & 1 \end{bmatrix}$$

**Data:** $M$ is a matrix with $N$ rows
**foreach** *row $M_i$* **do**
    $M_i \leftarrow M_i / M_{ii}$;
    **foreach** *row $M_j$ in M where $j \neq i$* **do**
        $M_j \leftarrow M_j - M_{ji} \times M_i$
    **end**
**end**

$$\begin{bmatrix} 3 & 2 & | & 1 & 0 \\ -7 & -5 & | & 0 & 1 \end{bmatrix} \xrightarrow{R_1/3} \begin{bmatrix} \mathbf{1} & \mathbf{2/3} & | & \mathbf{1/3} & \mathbf{0} \\ -7 & -5 & | & 0 & 1 \end{bmatrix} \xrightarrow{R_2-(-7)R_1} \begin{bmatrix} 1 & 2/3 & | & 1/3 & 0 \\ \mathbf{0} & -1/3 & | & 7/3 & \mathbf{1} \end{bmatrix}$$

**Data:** $M$ is a matrix with $N$ rows
**foreach** *row $M_i$* **do**
$\quad$ $M_i \leftarrow M_i / M_{ii}$;
$\quad$ **foreach** *row $M_j$ in M where $j \neq i$* **do**
$\quad\quad$ $M_j \leftarrow M_j - M_{ji} \times M_i$
$\quad$ **end**
**end**

$$\begin{bmatrix} 3 & 2 & \bigm| & 1 & 0 \\ -7 & -5 & \bigm| & 0 & 1 \end{bmatrix} \xrightarrow{R_1/3} \begin{bmatrix} \mathbf{1} & \mathbf{2/3} & \bigm| & \mathbf{1/3} & \mathbf{0} \\ -7 & -5 & \bigm| & 0 & 1 \end{bmatrix} \xrightarrow{R_2-(-7)R_1} \begin{bmatrix} 1 & 2/3 & \bigm| & 1/3 & 0 \\ \mathbf{0} & \mathbf{-1/3} & \bigm| & \mathbf{7/3} & \mathbf{1} \end{bmatrix}$$

$$\xrightarrow{R_2/-1/3} \begin{bmatrix} 1 & 2/3 & \bigm| & 1/3 & 0 \\ \mathbf{0} & \mathbf{1} & \bigm| & \mathbf{-7} & \mathbf{-3} \end{bmatrix}$$

# Gauss-Jordan Algorithm

**Data:** $M$ is a matrix with $N$ rows
**foreach** *row $M_i$* **do**

$\quad M_i \leftarrow M_i / M_{ii}$;

$\quad$ **foreach** *row $M_j$ in M where $j \neq i$* **do**

$\quad\quad M_j \leftarrow M_j - M_{ji} \times M_i$

$\quad$ **end**

**end**

$$\begin{bmatrix} 3 & 2 & | & 1 & 0 \\ -7 & -5 & | & 0 & 1 \end{bmatrix} \xrightarrow{R_1/3} \begin{bmatrix} \mathbf{1} & \mathbf{2/3} & | & \mathbf{1/3} & \mathbf{0} \\ -7 & -5 & | & 0 & 1 \end{bmatrix} \xrightarrow{R_2-(-7)R_1} \begin{bmatrix} 1 & 2/3 & | & 1/3 & 0 \\ \mathbf{0} & \mathbf{-1/3} & | & \mathbf{7/3} & \mathbf{1} \end{bmatrix}$$

$$\xrightarrow{R_2/-1/3} \begin{bmatrix} 1 & 2/3 & | & 1/3 & 0 \\ \mathbf{0} & \mathbf{1} & | & \mathbf{-7} & \mathbf{-3} \end{bmatrix} \xrightarrow{R_1-2/3R_2} \begin{bmatrix} \mathbf{1} & \mathbf{0} & | & \mathbf{5} & \mathbf{2} \\ 0 & 1 & | & -7 & -3 \end{bmatrix}$$

# Class Activity

$$M = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$\left[\begin{array}{ccc|ccc} 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 2 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 \end{array}\right]$$

$$\begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix} \xrightarrow{R_1/1} \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{1} & | & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix} \xrightarrow{R_1/1} \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{1} & | & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix} \xrightarrow[R_3-1R_1]{R_2-0R_1} \begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ \mathbf{0} & \mathbf{2} & \mathbf{1} & | & \mathbf{0} & \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \mathbf{0} & | & \mathbf{-1} & \mathbf{0} & \mathbf{1} \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix} \xrightarrow{R_1/1} \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{1} & | & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix} \xrightarrow[R_3-1R_1]{R_2-0R_1} \begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ \mathbf{0} & \mathbf{2} & \mathbf{1} & | & \mathbf{0} & \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \mathbf{0} & | & \mathbf{-1} & \mathbf{0} & \mathbf{1} \end{bmatrix}$$

$$\xrightarrow{R_2/2} \begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ \mathbf{0} & \mathbf{1} & \mathbf{1/2} & | & \mathbf{0} & \mathbf{1/2} & \mathbf{0} \\ 0 & 1 & 0 & | & -1 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix} \xrightarrow{R_1/1} \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{1} & | & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix} \xrightarrow[R_3-1R_1]{R_2-0R_1} \begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ \mathbf{0} & \mathbf{2} & \mathbf{1} & | & \mathbf{0} & \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \mathbf{0} & | & \mathbf{-1} & \mathbf{0} & \mathbf{1} \end{bmatrix}$$

$$\xrightarrow{R_2/2} \begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ \mathbf{0} & \mathbf{1} & \mathbf{1/2} & | & \mathbf{0} & \mathbf{1/2} & \mathbf{0} \\ 0 & 1 & 0 & | & -1 & 0 & 1 \end{bmatrix} \xrightarrow[R_3-1R_2]{R_1-0R_2} \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{1} & | & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ 0 & 1 & 1/2 & | & 0 & 1/2 & 0 \\ \mathbf{0} & \mathbf{0} & \mathbf{-1/2} & | & \mathbf{-1} & \mathbf{-1/2} & \mathbf{1} \end{bmatrix}$$

# SOLUTION

$$\begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix} \xrightarrow{R_1/1} \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{1} & | & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix} \xrightarrow[R_3-1R_1]{R_2-0R_1} \begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ \mathbf{0} & \mathbf{2} & \mathbf{1} & | & \mathbf{0} & \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \mathbf{0} & | & \mathbf{-1} & \mathbf{0} & \mathbf{1} \end{bmatrix}$$

$$\xrightarrow{R_2/2} \begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ \mathbf{0} & \mathbf{1} & \mathbf{^1/_2} & | & \mathbf{0} & \mathbf{^1/_2} & \mathbf{0} \\ 0 & 1 & 0 & | & -1 & 0 & 1 \end{bmatrix} \xrightarrow[R_3-1R_2]{R_1-0R_2} \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{1} & | & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ 0 & 1 & ^1/_2 & | & 0 & ^1/_2 & 0 \\ \mathbf{0} & \mathbf{0} & \mathbf{-^1/_2} & | & \mathbf{-1} & \mathbf{-^1/_2} & \mathbf{1} \end{bmatrix}$$

$$\xrightarrow{R_3/-^1/_2} \begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ 0 & 1 & ^1/_2 & | & 0 & ^1/_2 & 0 \\ \mathbf{0} & \mathbf{0} & \mathbf{1} & | & \mathbf{2} & \mathbf{1} & \mathbf{-2} \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix} \xrightarrow{R_1/1} \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{1} & | & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ 0 & 2 & 1 & | & 0 & 1 & 0 \\ 1 & 1 & 1 & | & 0 & 0 & 1 \end{bmatrix} \xrightarrow[R_3-1R_1]{R_2-0R_1} \begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ \mathbf{0} & \mathbf{2} & \mathbf{1} & | & \mathbf{0} & \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \mathbf{0} & | & \mathbf{-1} & \mathbf{0} & \mathbf{1} \end{bmatrix}$$

$$\xrightarrow{R_2/2} \begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ \mathbf{0} & \mathbf{1} & \mathbf{^1/_2} & | & \mathbf{0} & \mathbf{^1/_2} & \mathbf{0} \\ 0 & 1 & 0 & | & -1 & 0 & 1 \end{bmatrix} \xrightarrow[R_3-1R_2]{R_1-0R_2} \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{1} & | & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ 0 & 1 & ^1/_2 & | & 0 & ^1/_2 & 0 \\ \mathbf{0} & \mathbf{0} & \mathbf{-^1/_2} & | & \mathbf{-1} & \mathbf{-^1/_2} & \mathbf{1} \end{bmatrix}$$

$$\xrightarrow{R_3/-^1/_2} \begin{bmatrix} 1 & 0 & 1 & | & 1 & 0 & 0 \\ 0 & 1 & ^1/_2 & | & 0 & ^1/_2 & 0 \\ \mathbf{0} & \mathbf{0} & \mathbf{1} & | & \mathbf{2} & \mathbf{1} & \mathbf{-2} \end{bmatrix} \xrightarrow[R_2-^1/_2R_3]{R_1-1R_3} \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{0} & | & \mathbf{-1} & \mathbf{-1} & \mathbf{2} \\ \mathbf{0} & \mathbf{1} & \mathbf{0} & | & \mathbf{-1} & \mathbf{0} & \mathbf{1} \\ 0 & 0 & 1 & | & 2 & 1 & -2 \end{bmatrix}$$
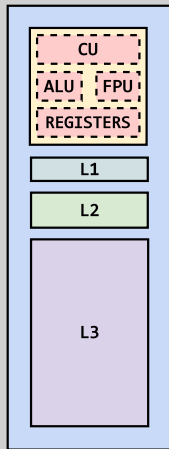
# GPU Programming

- Core
  - ▶ Processing unit
- L1 L2 L3
  - ▶ Caches for storing work

- Registers
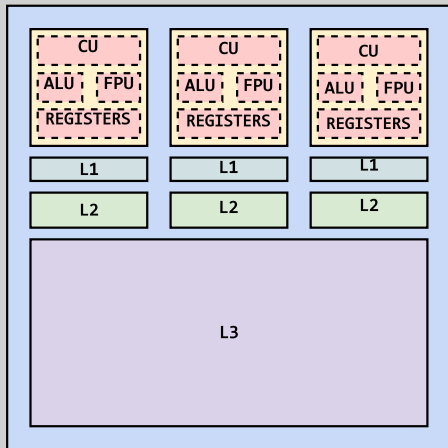  - ▶ Store localized data including the fetched instruction
- Control Unit (CU)
  - ▶ Decodes instruction and sends to appropriate logic unit
- Arithmetic Logic Unit (ALU) / Floating Point Unit (FPU)
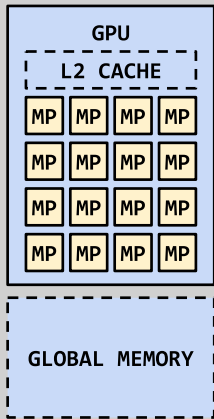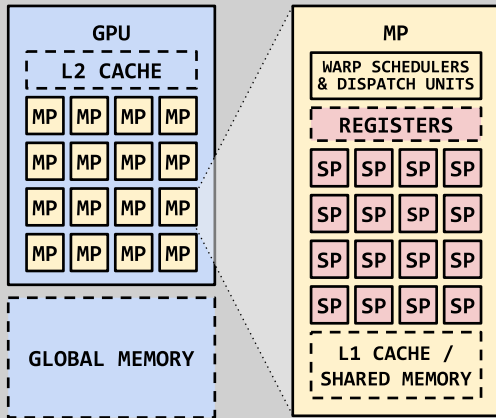  - ▶ Executes given instruction

- Registers
  - Store localized data including the fetched instruction
- Control Unit (CU)
  - Decodes instruction and sends to appropriate logic unit
- Arithmetic Logic Unit (ALU) / Floating Point Unit (FPU)
  - Executes given instruction
- Each core has its own registers and control logic

**GPU**

**L2 CACHE**

| MP | MP | MP | MP |
| MP | MP | MP | MP |
| MP | MP | MP | MP |
| MP | MP | MP | MP |

**GLOBAL MEMORY**

- Multiprocessor (MP)
  - A "core" that handles simultaneous execution of a vector of tasks

- Multiprocessor (MP)
  - A "core" that handles simultaneous execution of a vector of tasks
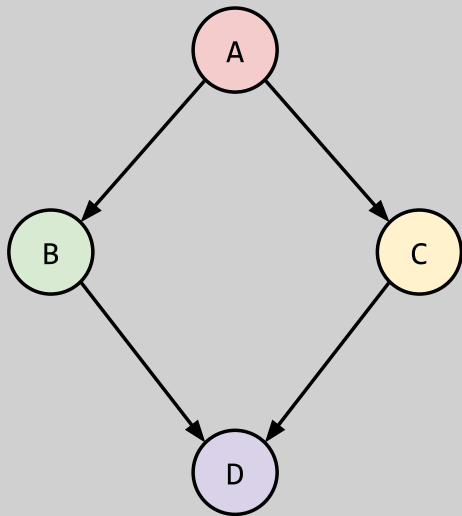- Scalar processor (SP)
  - Executes a single scalar component

- Warp
  - Group of SPs in a MP that execute in lockstep
  - Warps share registers; including the program counter
- Kernel
  - Group of warps that operate on the same method

```
A(x)

while(B1(x)){


        B2(x)




}

C(x);
```

```
1   // A special method invoked by the
2   // CPU to launch a GPU kernel
3   __global__ print(string message) {
4       int idx = threadIdx.x;
5       int jdx = blockIdx.x;
6       printf("%s from (%d, %d)\n",
7               message, jdx, idx);
8   }
9
10  // Standard C main
11  int main() {
12      // Call kernel launcher
13      print<<<2, 4>>>("Hello World");
14      // Wait for GPU to finish
15      cudaDeviceSynchronize();
16  }
```

### Output

Hello World from (0, 0)
Hello World from (0, 3)
Hello World from (1, 2)
Hello World from (1, 0)
Hello World from (0, 2)
Hello World from (1, 1)
Hello World from (1, 3)
Hello World from (0, 1)

```
1  for (size_t j = 0; j < rows; j++) {
2    fixRow<<<1, cols>>>(data_gpu, cols, j);
3    auto_throw(cudaDeviceSynchronize());
4
5    fixColumn<<<rows, cols>>>(data_gpu, cols, j);
6    auto_throw(cudaDeviceSynchronize());
7  }
```

```
1   __global__ void fixRow(
2       float *matrix, int size, int rowId) {
3     // the ith row of the matrix
4     __shared__ float Ri[MAX_BLOCK_SIZE];
5     // The diagonal element for ith row
6     __shared__ float Aii;
7     int colId = threadIdx.x;
8     Ri[colId] = matrix[size * rowId + colId];
9     Aii = matrix[size * rowId + rowId];
10
11    __syncthreads();
12    // Divide the whole row by the diagonal
13    Ri[colId] = Ri[colId] / Aii;
14    matrix[size * rowId + colId] = Ri[colId];
15  }
```

Example

$$\begin{bmatrix} 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 2 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 \end{bmatrix}$$

```
1  __global__ void fixColumn(
2      float *matrix, int size, int colId) {
3    int i = threadIdx.x, j = blockIdx.x;
4    // The colId column
5    __shared__ float col[MAX_BLOCK_SIZE];
6    // The jth element of the colId row
7    __shared__ float AColIdj;
8    // The jth column
9    __shared__ float colj[MAX_BLOCK_SIZE];
10   col[i] = matrix[i * size + colId];
11   __syncthreads();
12   colj[i] = matrix[i * size + j];
13   AColIdj = matrix[colId * size + j];
14   if (i != colId) {
15     colj[i] = colj[i] - AColIdj * col[i];
16   }
17   matrix[i * size + j] = colj[i];
18 }
```

### Example

$$\begin{bmatrix} 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 2 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 1 \end{bmatrix}$$

16

$$\left[\begin{array}{ccc|ccc} 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1/2 & 0 & 1/2 & 0 \\ 0 & 0 & -1/2 & -1 & -1/2 & 1 \end{array}\right]$$

$$\left[\begin{array}{ccc|ccc} 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1/2 & 0 & 1/2 & 0 \\ 0 & 0 & -1/2 & -1 & -1/2 & 1 \end{array}\right]$$

$$\left[\begin{array}{ccc|ccc} 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1/2 & 0 & 1/2 & 0 \\ 0 & 0 & -1/2 & -1 & -1/2 & 1 \end{array}\right] \rightarrow \left[\begin{array}{ccc|c} 1 & 0 & 0 & 1 \\ 0 & 1/2 & 0 & 1/2 \\ -1 & -1/2 & 1 & -1/2 \end{array}\right]$$
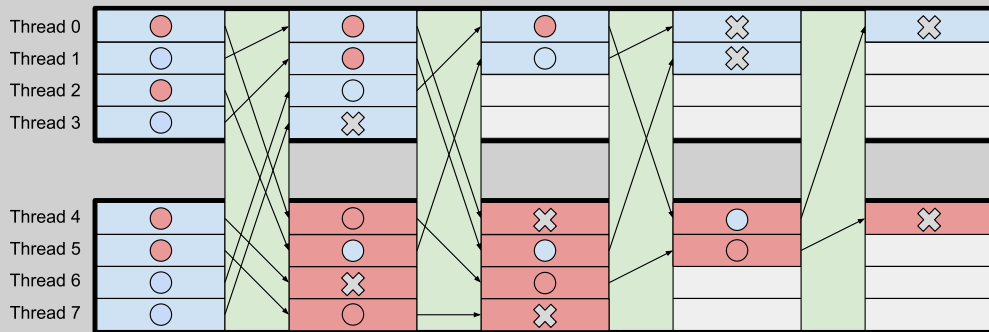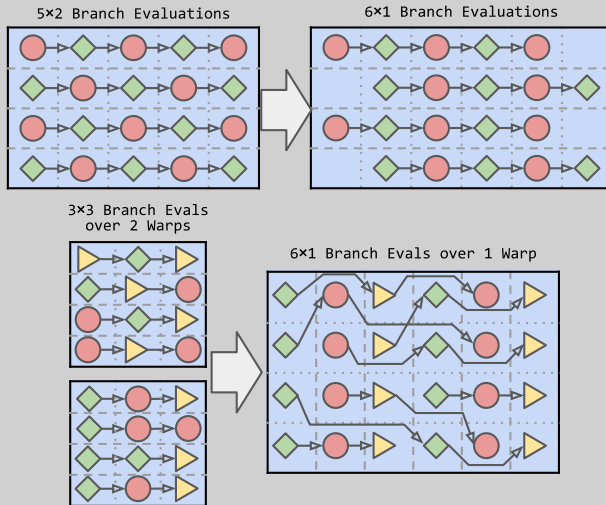
5×2 Branch Evaluations

6×1 Branch Evaluations

3×3 Branch Evals over 2 Warps

6×1 Branch Evals over 1 Warp

```cpp
struct FixRow {
  using Type = void(*)(index_t rowId, size_t colId);

  template<typename PROGRAM>
  __device__ static void eval(PROGRAM prog, index_t rowId, size_t colId) {
    index_t size = prog.device.size.row;
    matrix_t Ri  = prog.device.matrix[size*rowId + colId];
    matrix_t Aii = prog.device.Aij[rowId];

    Ri /= Aii;
    prog.device.matrix[size*rowId + colId] = Ri;

    if( Ri != 0.0 ) {
      prog.template async<SplitCol>(rowId, 0, prog.device.size.col-1, colId);
    }
  }
};
```

```
1   struct FixCol {
2     using Type = void(*)(index_t colId, index_t i_start,
3                          index_t i_end, index_t j);
4     template<typename PROGRAM>
5     __device__ static void eval(
6         PROGRAM prog, index_t colId,
7         index_t i_start, index_t i_end, index_t j) {
8       index_t size = prog.device.size.row;
9       for (index_t i = i_start; i <= i_end; i++) {
10        matrix_t col = prog.device.Aij[i];
11
12        if (col != 0) {
13          matrix_t colj   = prog.device.matrix[    i*size + j];
14          matrix_t AColIdj = prog.device.matrix[colId*size + j];
15          if (i != colId) {
16            colj -= AColIdj * col;
17            prog.device.matrix[i*size + j] = colj;
18  }}}};
```

# Results

- cpu-inverse

- cpu-inverse
- inverse

- cpu-inverse
- inverse
- tdr-inverse

Execution time for an n x n matrix, best of 3 runs

Execution time for an n x n matrix, best of 3 runs

- The use of TDR can offer fairly significant gains when applied to sparse matrix inversion.

# Conclusion / Closing Thoughts

- The use of TDR can offer fairly significant gains when applied to sparse matrix inversion.
- Possible next step is to implement Strassen Inversion [4].

[1]    G. SHARMA, A. AGARWALA, AND B. BHATTACHARYA, "A FAST PARALLEL GAUSS JORDAN ALGORITHM FOR MATRIX INVERSION USING CUDA," *Computers & Structures*, vol. 128, pp. 31–37, 2013, ISSN: 0045-7949. DOI: `10.1016/j.compstruc.2013.06.015`. [Online]. Available: `https://www.sciencedirect.com/science/article/pii/S0045794913002095` (cit. on pp. 3–5).

[2]    D. DASGUPTA, "IN-PLACE MATRIX INVERSION BY MODIFIED GAUSS-JORDAN ALGORITHM," *Applied Mathematics*, vol. 4, no. 10, pp. 1392–1396, Sep. 2013. DOI: `10.4236/am.2013.410188` (cit. on pp. 3–5).

[3] B. Cuneo and M. Bailey, **"Divergence reduction in monte carlo neutron transport with on-gpu asynchronous scheduling,"** *ACM Trans. Model. Comput. Simul.*, vol. 34, no. 1, Jan. 2024, ISSN: 1049-3301. DOI: 10.1145/3626957. [Online]. Available: https://dl.acm.org/doi/10.1145/3626957 (cit. on pp. 3–5).

[4] V. Strassen, **"Gaussian elimination is not optimal,"** *Numerische Mathematik*, vol. 13, pp. 354–356, Aug. 1969. DOI: 10.1007/BF02165411 (cit. on pp. 65, 66).

[5] J. R. Gardner, G. Pleiss, D. Bindel, K. Q. Weinberger, and A. G. Wilson, **"Gpytorch: Blackbox matrix-matrix gaussian process inference with gpu acceleration,"** in *Advances in Neural Information Processing Systems*, 2018. DOI: 10.48550/arXiv.1809.11165. [Online]. Available: https://arxiv.org/abs/1809.11165.

[6] D. BAILEY AND H. FERGUSON, **"A STRASSEN-NEWTON ALGORITHM FOR HIGH-SPEED PARALLELIZABLE MATRIX INVERSION,"** in *Supercomputing '88:Proceedings of the 1988 ACM/IEEE Conference on Supercomputing, Vol. I*, 1988, pp. 419–424. DOI: `10.1109/SUPERC.1988.44680`.

[7] J. XUEBIN, C. YEWANG, F. WENTAO, Z. YONG, AND D. JIXIANG, **"FAST ALGORITHM FOR PARALLEL SOLVING INVERSION OF LARGE SCALE SMALL MATRICES BASED ON GPU,"** *The Journal of Supercomputing*, vol. 79, no. 16, pp. 18 313–18 339, May 2023, ISSN: 1573-0484. DOI: `10.1007/s11227-023-05336-7`. [Online]. Available: `https://link.springer.com/article/10.1007/s11227-023-05336-7`.

[8]     S. C. ALTHOEN AND R. MCLAUGHLIN, "GAUSS-JORDAN REDUCTION: A BRIEF HISTORY," *American Mathematical Monthly*, vol. 94, no. 2, pp. 130–142, Feb. 1987, ISSN: 0002-9890. DOI: `10.2307/2322413`. [Online]. Available: `https://www.jstor.org/stable/2322413`.

[9]     H. COHN, R. KLEINBERG, B. SZEGEDY, AND C. UMANS, "GROUP-THEORETIC ALGORITHMS FOR MATRIX MULTIPLICATION," in *46th Annual IEEE Symposium on Foundations of Computer Science (FOCS'05)*, 2005, pp. 379–388. DOI: `10.1109/SFCS.2005.39`.

[10]    D. COPPERSMITH AND S. WINOGRAD, "ON THE ASYMPTOTIC COMPLEXITY OF MATRIX MULTIPLICATION," in *22nd Annual Symposium on Foundations of Computer Science (sfcs 1981)*, 1981, pp. 82–90. DOI: `10.1109/SFCS.1981.27`.

# References V

[11] I. DIMOV, T. DIMOV, AND T. GUROV, **"A NEW ITERATIVE MONTE CARLO APPROACH FOR INVERSE MATRIX PROBLEM,"** *Journal of Computational and Applied Mathematics*, vol. 92, no. 1, pp. 15–35, 1998, ISSN: 0377-0427. DOI: `10.1016/S0377-0427(98)00043-0`. [Online]. Available: `https://www.sciencedirect.com/science/article/pii/S0377042798000430`.

[12] C. MISRA, S. HALDAR, S. BHATTACHARYA, AND S. K. GHOSH, **"SPIN: A FAST AND SCALABLE MATRIX INVERSION METHOD IN APACHE SPARK,"** in *Proceedings of the 19th International Conference on Distributed Computing and Networking*, ser. ICDCN '18, Varanasi, India: Association for Computing Machinery, 2018, ISBN: 9781450363723. DOI: `10.1145/3154273.3154300`.

[13]     E. AGULLO, H. BOUWMEESTER, J. DONGARRA, J. KURZAK, J. LANGOU, AND L. ROSENBERG, **"TOWARDS AN EFFICIENT TILE MATRIX INVERSION OF SYMMETRIC POSITIVE DEFINITE MATRICES ON MULTICORE ARCHITECTURES,"** in *High Performance Computing for Computational Science – VECPAR 2010*, J. M. L. M. Palma, M. Daydé, O. Marques, and J. C. Lopes, Eds., Berlin, Heidelberg: Springer Berlin Heidelberg, 2011, pp. 129–138, ISBN: 978-3-642-19328-6. DOI: 10.1007/978-3-642-19328-6_14.

[14]     W. H. PRESS, S. A. TEUKOLSKY, W. T. VETTERLING, AND B. P. FLANNERY, **"SOLUTION OF LINEAR ALGEBRAIC EQUATIONS,"** in *Numerical Recipes: The Art of Scientific Computing*, 3rd ed. 32 Avenue of the Americas, New York, NY 10013-2473, USA: Cambridge University Press, 2007, ch. 2, pp. 37–109, ISBN: 978-0-521-88068-8.

[15]   H. ANTON AND C. ROSSES, **"APPLICATIONS OF LINEAR ALGEBRA,"** in *Elementary Linear Algebra*, 11th ed. 111 River Street, Hoboken, NJ 07030-5774, USA: John Wiley & Sons, Inc., 2014, ch. 10, pp. 527–713, ISBN: 978-1-118-43441-3.

[16]   H. NGUYEN, *GPU GEMS 3*, FIRST. ADDISON-WESLEY PROFESSIONAL, 2007, ISBN: 9780321545428.

https://github.com/scrufulufugus/tdr-inverse-materials