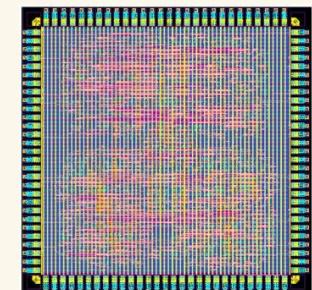




# Accelerator for Genetic Local Sequence Alignment

## Frontend Design Review

**Students:** Niv Bar-Tov, Dor Lerman  
**Supervisor:** Goel Samuel



# Outline

- Sequence Alignment
- Smith-Waterman Algorithm
- Project's Requirements and Specifications
- High-Level Overview
- Architectural Design
- Next Steps

# Sequence Alignment

- Method of arranging sequences of DNA, RNA or proteins
- Aims to identify regions of similarity that may indicate on a relationship between the sequences
- Can be used also for non-biological sequences



# Sequence Alignment

TACGCTTG  
CTACCTAG



TAC - CT  
TACGCT

# Global Alignment VS. Local Alignment

- Align the entire sequence (all letters)
- Suitable for closely related sequences
- Needleman-Wunsch Algorithm
- Align regions (substrings) having highest similarities
- Suitable for more divergent sequences
- Smith-Waterman Algorithm

# Smith-Waterman Algorithm

- Dynamic Programming algorithm used for local sequence alignment
- Compares segments of all possible lengths and optimizes the similarity measure
- Time complexity:  $O(n^2)$
- Space complexity:  $O(n^2)$

# Smith-Waterman Main Stages

1. Initialization: Construct a scoring matrix and initialize its first row and column with zeros
2. Scoring: Calculate score for each cell based on the scoring scheme
3. Traceback: Determine the optimal alignment by trace back from the highest-scoring cell in the matrix to a cell with a score of zero

# Smith-Waterman Algorithm

## Initialization

$$H_{i,j} = \max \begin{cases} H_{i-1,j-1} + MATCH/MISMTACH \\ H_{i-1,j} + GAP\_PENALTY \\ H_{i,j-1} + GAP\_PENALTY \\ 0 \end{cases}$$

Match	+1
Mismatch	-1
Gap Penalty	-2

S	T	A	C	G	C	T	T	T	G
0	0	0	0	0	0	0	0	0	0
C	0								
T	0								
A	0								
C	0								
C	0								
T	0								
A	0								
G	0								

# Smith-Waterman Algorithm

## Scoring

$$H_{i,j} = \max \begin{cases} H_{i-1,j-1} + MATCH/MISMTACH \\ H_{i-1,j} + GAP\_PENALTY \\ H_{i,j-1} + GAP\_PENALTY \\ 0 \end{cases}$$

Match	+1
Mismatch	-1
Gap Penalty	-2

S	T	A	C	G	C	T	T	G
0	0	0	0	0	0	0	0	0
C	0	0						
T	0							
A	0							
C	0							
C	0							
T	0							
A	0							
G	0							

# Smith-Waterman Algorithm

## Scoring

$$H_{i,j} = \max \begin{cases} H_{i-1,j-1} + MATCH/MISMTACH \\ H_{i-1,j} + GAP\_PENALTY \\ H_{i,j-1} + GAP\_PENALTY \\ 0 \end{cases}$$

Match	+1
Mismatch	-1
Gap Penalty	-2

S	T	A	C	G	C	T	T	T	G
C	0	0	0	0	0	0	0	0	0
T	0	1							
A	0								
C	0								
C	0								
T	0								
A	0								
G	0								

# Smith-Waterman Algorithm

## Scoring

$$H_{i,j} = \max \begin{cases} H_{i-1,j-1} + MATCH/MISMTACH \\ H_{i-1,j} + GAP\_PENALTY \\ H_{i,j-1} + GAP\_PENALTY \\ 0 \end{cases}$$

Match	+1
Mismatch	-1
Gap Penalty	-2

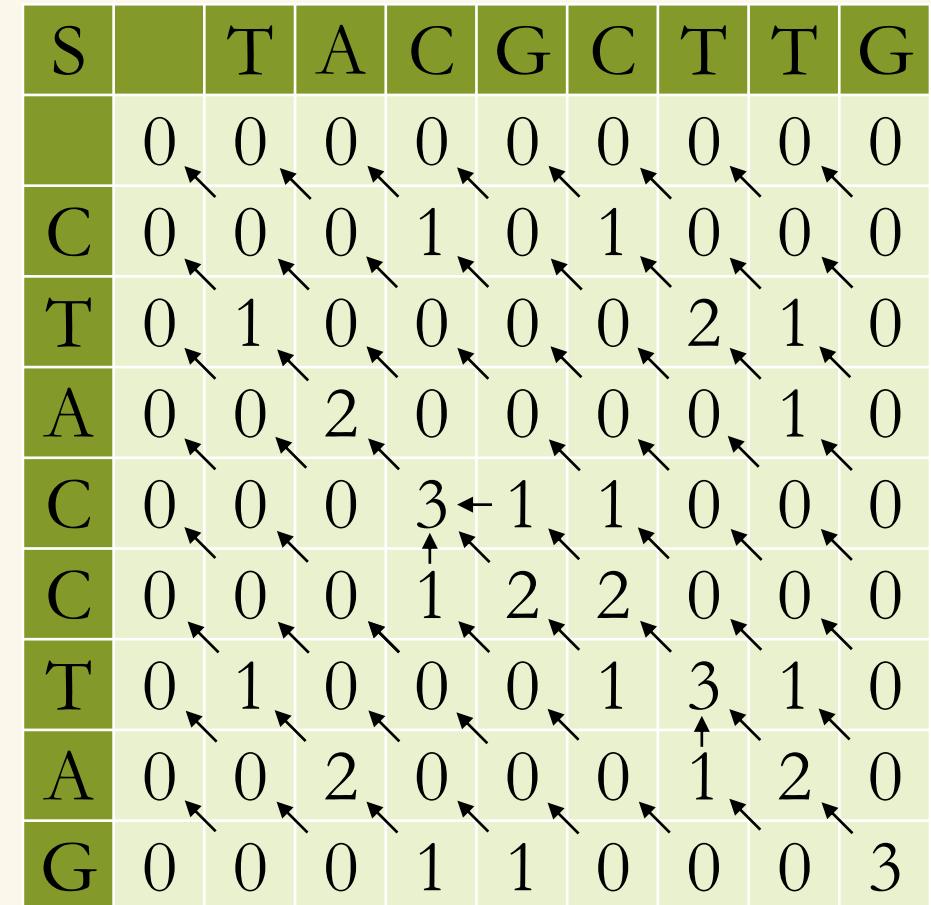
S	T	A	C	G	C	T	T	T	G
0	0	0	0	0	0	0	0	0	0
C	0	0	0	1					
T	0	1	0						
A	0	0							
C	0								
C	0								
T	0								
A	0								
G	0								

# Smith-Waterman Algorithm

## Scoring

$$H_{i,j} = \max \begin{cases} H_{i-1,j-1} + MATCH/MISMTACH \\ H_{i-1,j} + GAP\_PENALTY \\ H_{i,j-1} + GAP\_PENALTY \\ 0 \end{cases}$$

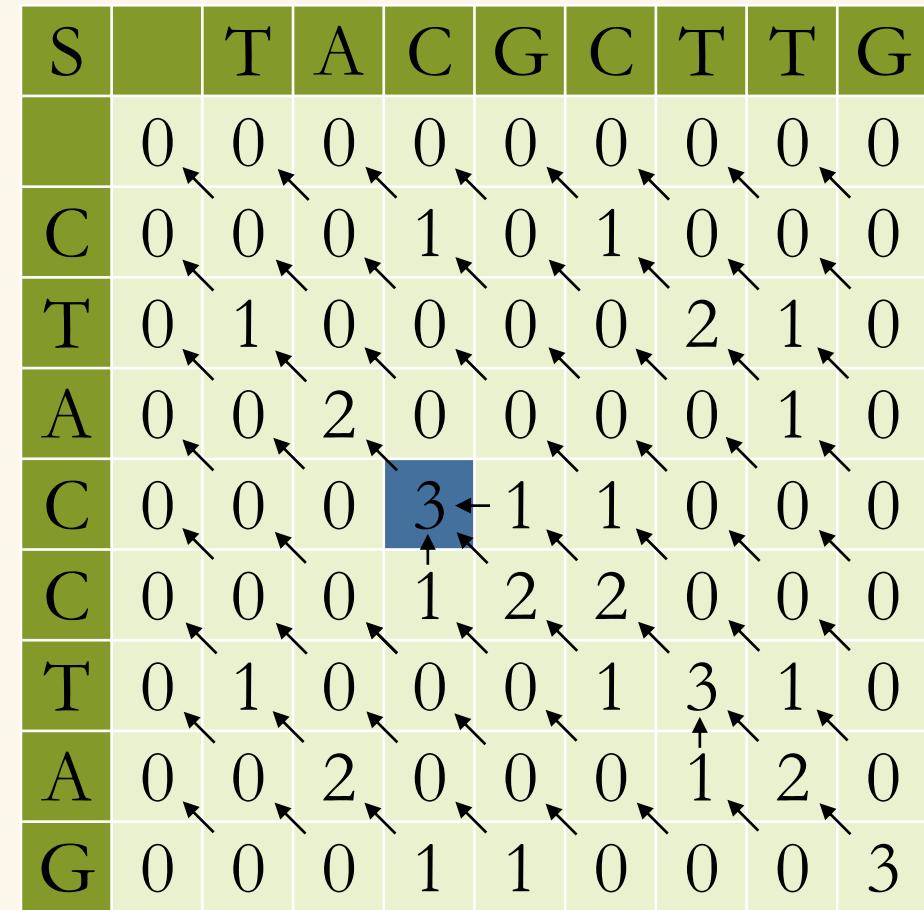
Match	+1
Mismatch	-1
Gap Penalty	-2



# Smith-Waterman Algorithm

## Traceback

Match	+1
Mismatch	-1
Gap Penalty	-2

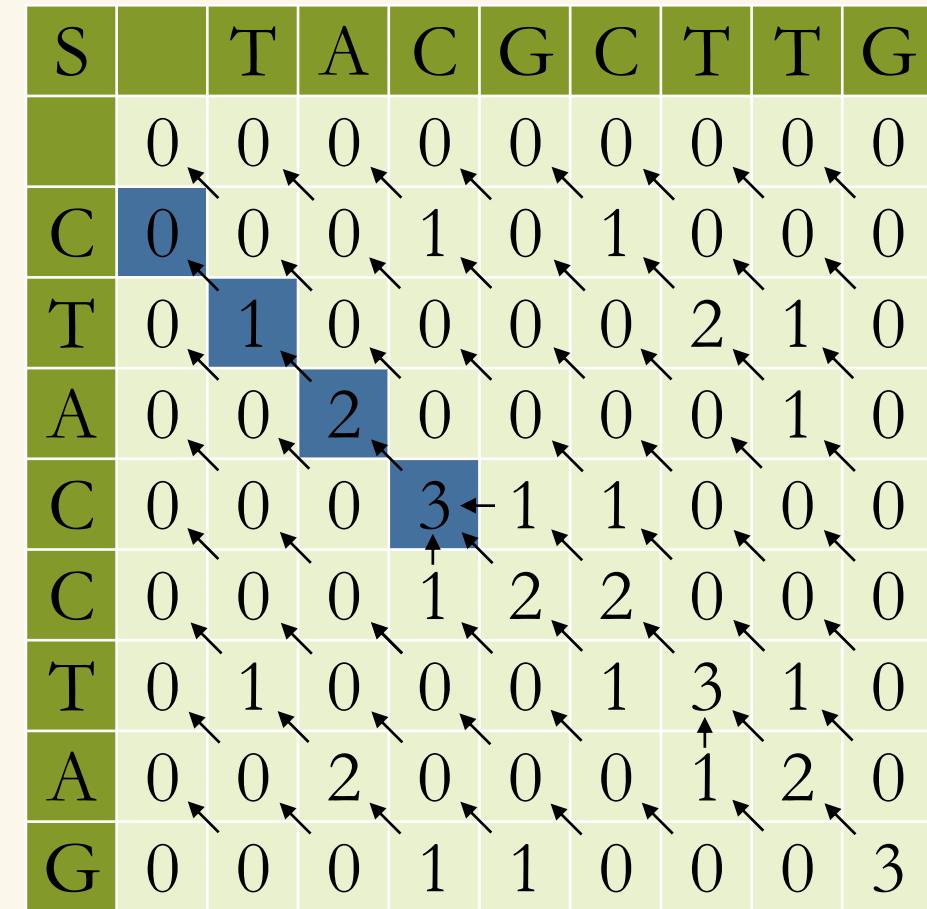


# Smith-Waterman Algorithm

## Traceback

TAC  
TAC

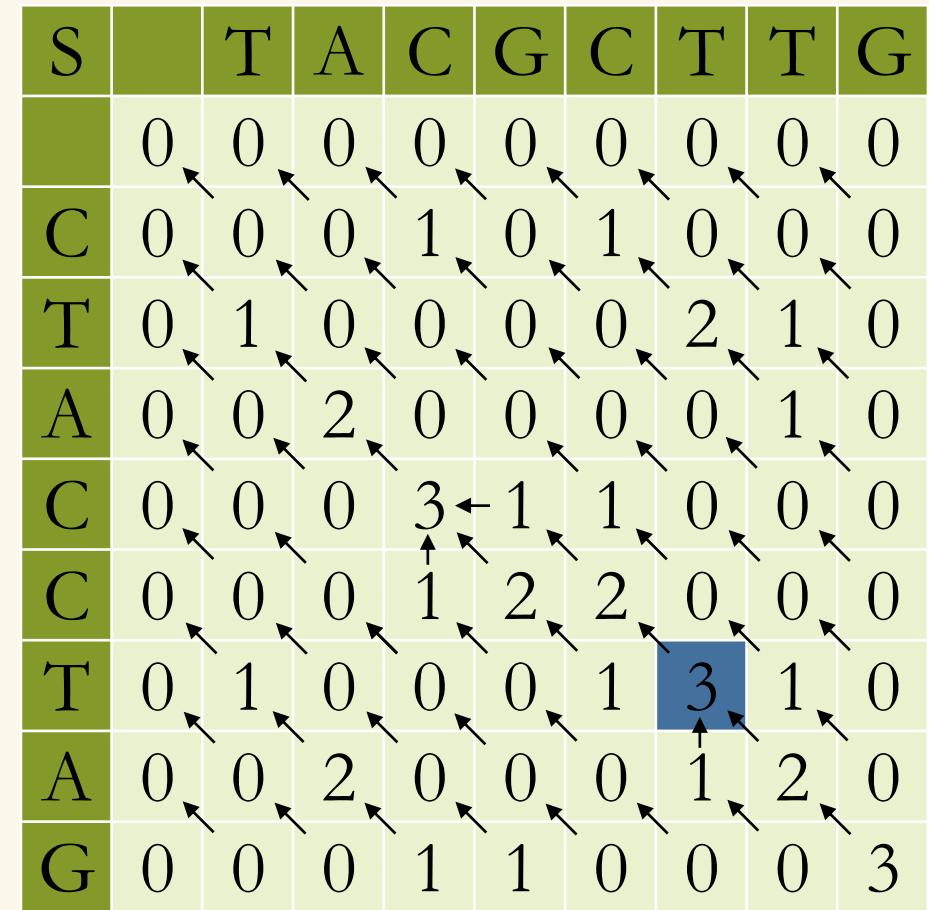
Match	+1
Mismatch	-1
Gap Penalty	-2



# Smith-Waterman Algorithm

## Traceback

Match	+1
Mismatch	-1
Gap Penalty	-2

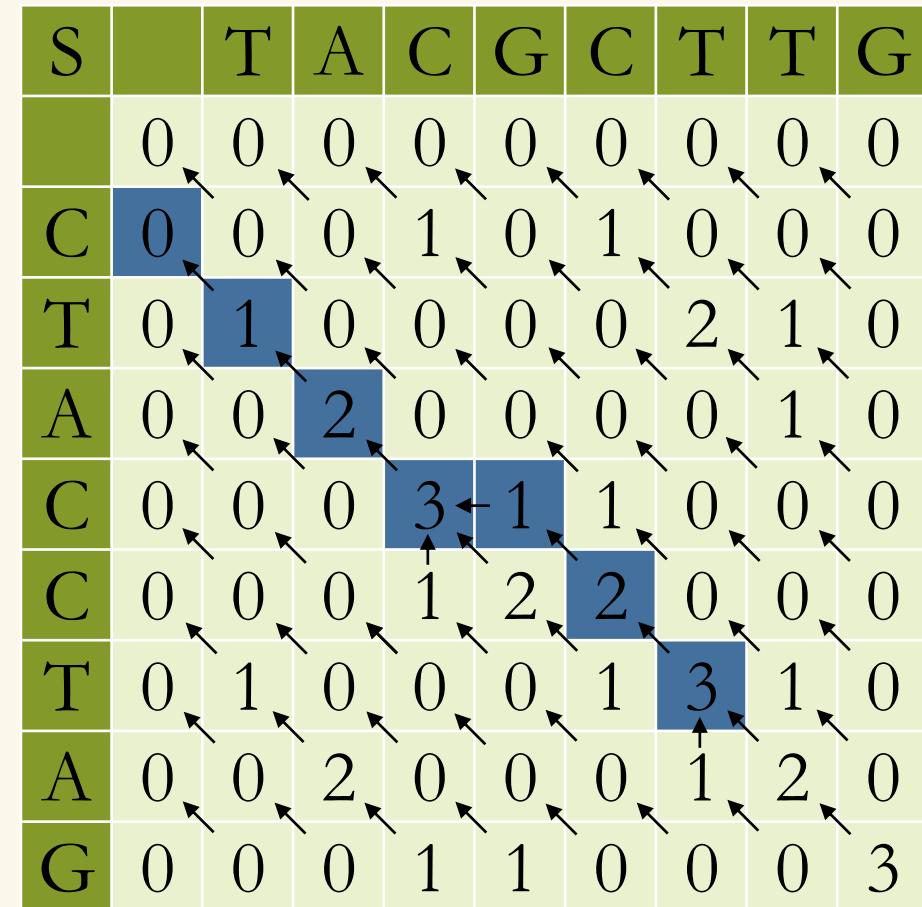


# Smith-Waterman Algorithm

## Traceback

TAC-CT  
TACGCT

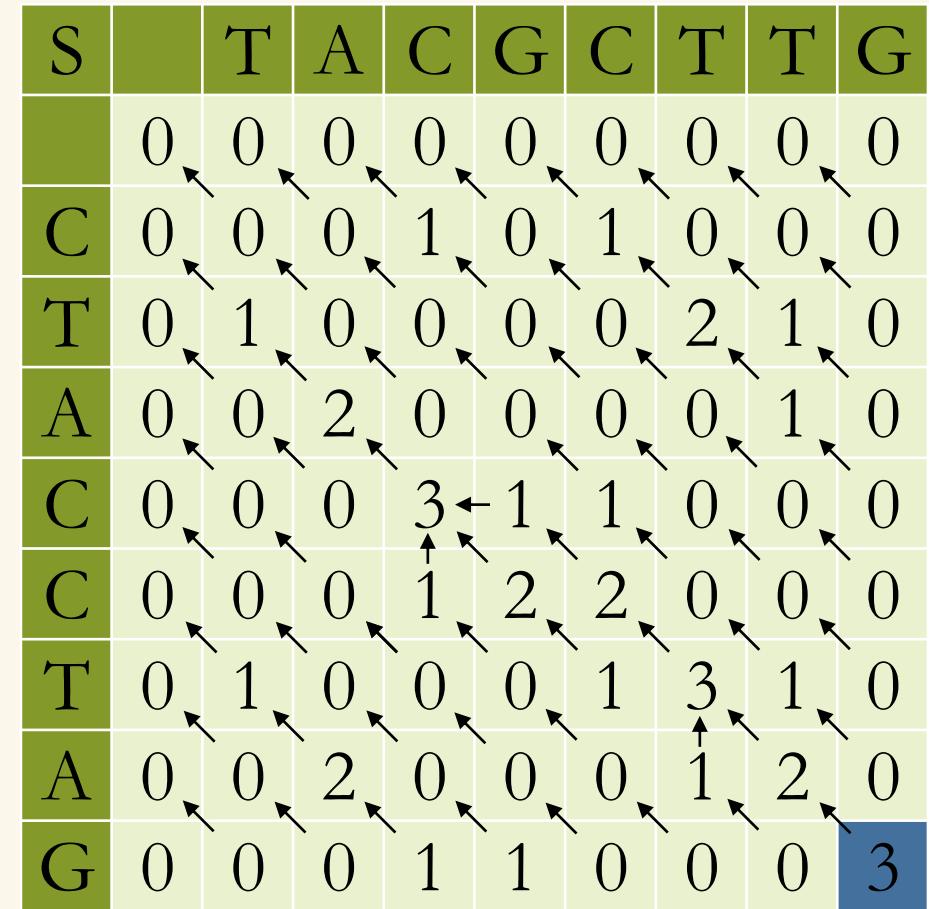
Match	+1
Mismatch	-1
Gap Penalty	-2



# Smith-Waterman Algorithm

## Traceback

Match	+1
Mismatch	-1
Gap Penalty	-2

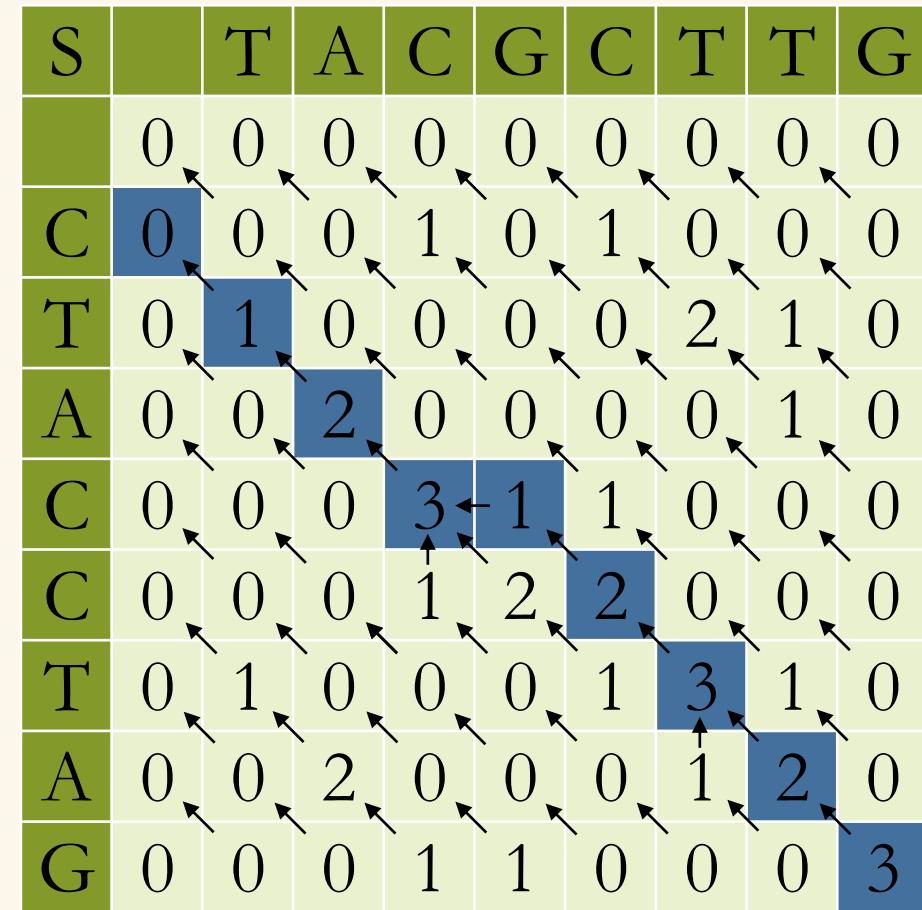


# Smith-Waterman Algorithm

## Traceback

TAC-CTAG  
TACGCTTG

Match	+1
Mismatch	-1
Gap Penalty	-2



# High-Level Overview

---

# Requirements & Specifications

- Technology process: TSMC 65 nm
- Sequences length: 32 residues
- Maximum area efficiency (1.0x1.0 sq. mm)
- Maximum of 40 I/O pins
- Clock frequency: up to 300 MHz
- Minimal on-chip memory
- Performance efficiency

# Sequences Representation

Input

A	00
G	01
T	10
C	11

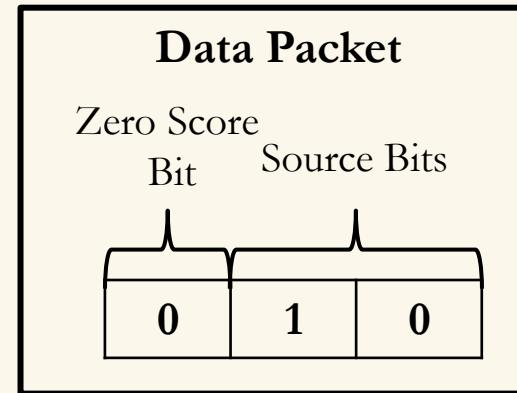
Output

A	000
G	001
T	010
C	011
—	100
Don't Care	101
Don't Care	110
Start/End Signal	111

# Data Packet

We need to keep in each cell:

1. Zero Score Bit – Indicates if the score of the cell equals zero
2. Source Bits – 2 bits indicate the neighbor cell from where the maximum score has been determined (left, top or diagonally located)

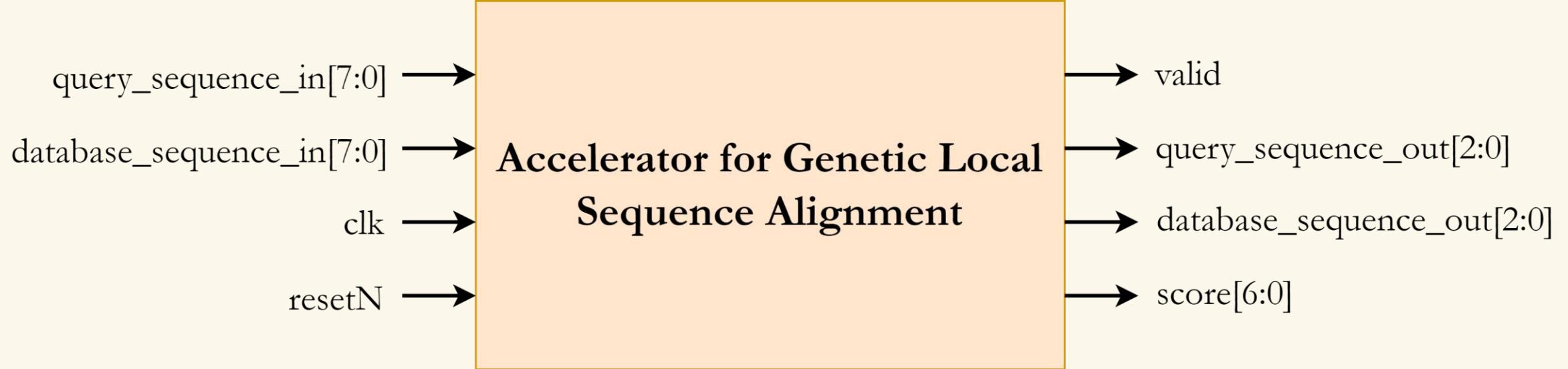


# High-Level Architecture Stages

Consists of 3 main stages:

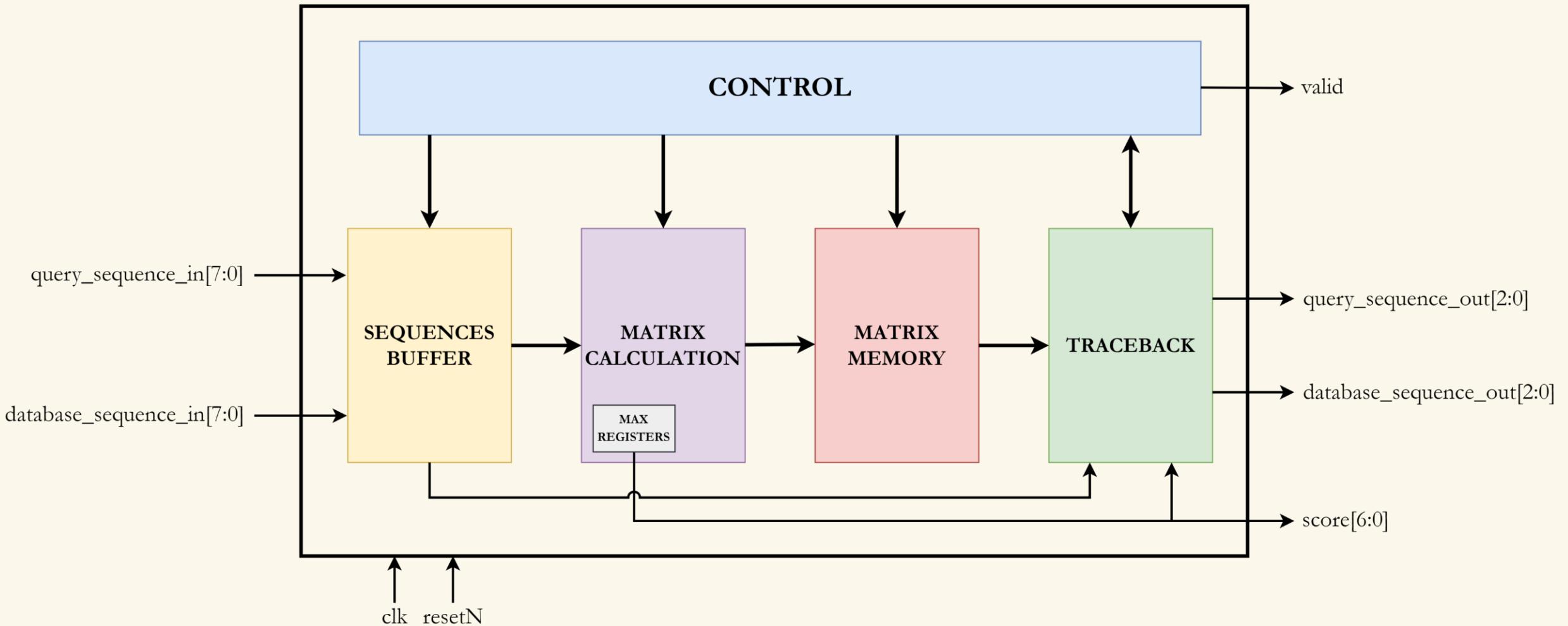
- 1. Load Sequences:** Load the two input sequences into the Sequences Buffer unit (4 residues from each sequence in parallel in each cycle).
- 2. Score Matrix Computation:** Initializing the Matrix Memory unit with zeros and then iteratively filling it according to the Matrix Calculation unit, while updating the Matrix Memory and Max Score Register.
- 3. Traceback:** Starting at the highest-scoring cell in the matrix and following a path of highest-scoring cells until zero cell score is reached.

# I/O Diagram



18 input pins, 14 output pins  
(+8 pins for  $V_{DD}$  and  $V_{SS}$ )

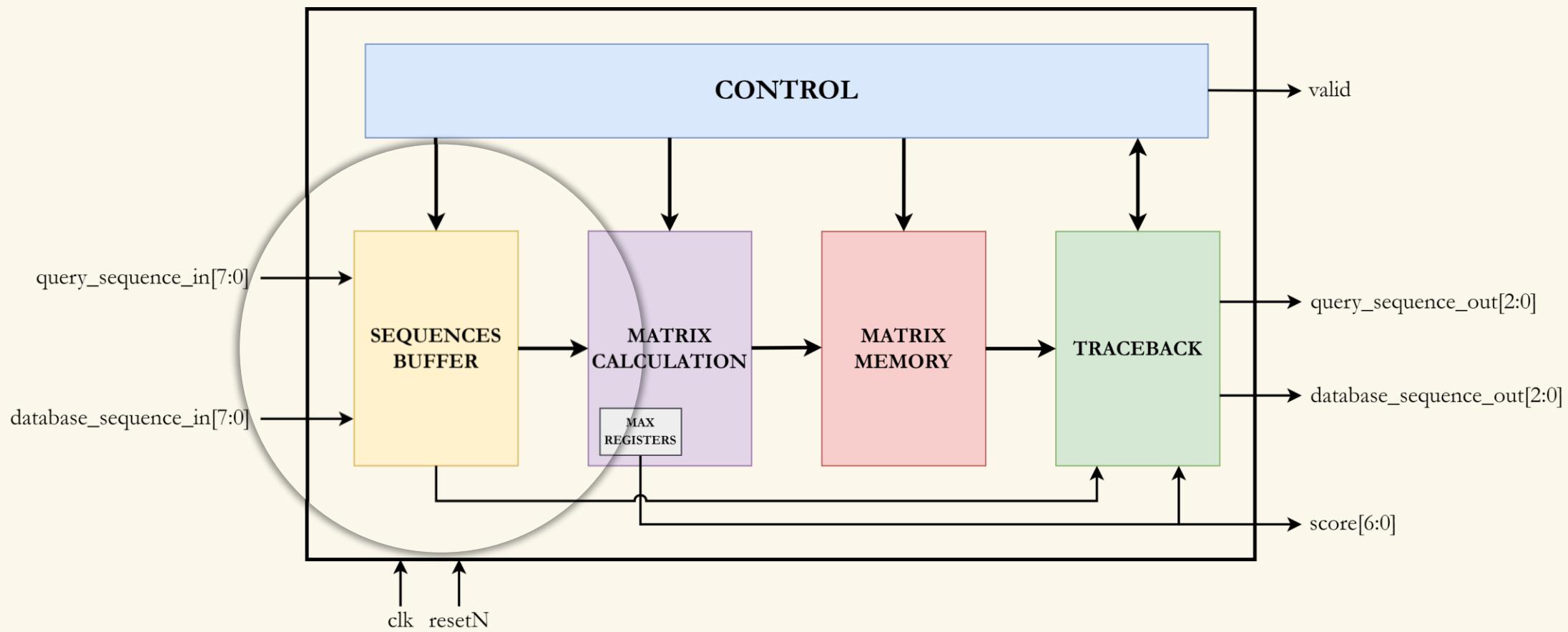
# Top Level Architecture



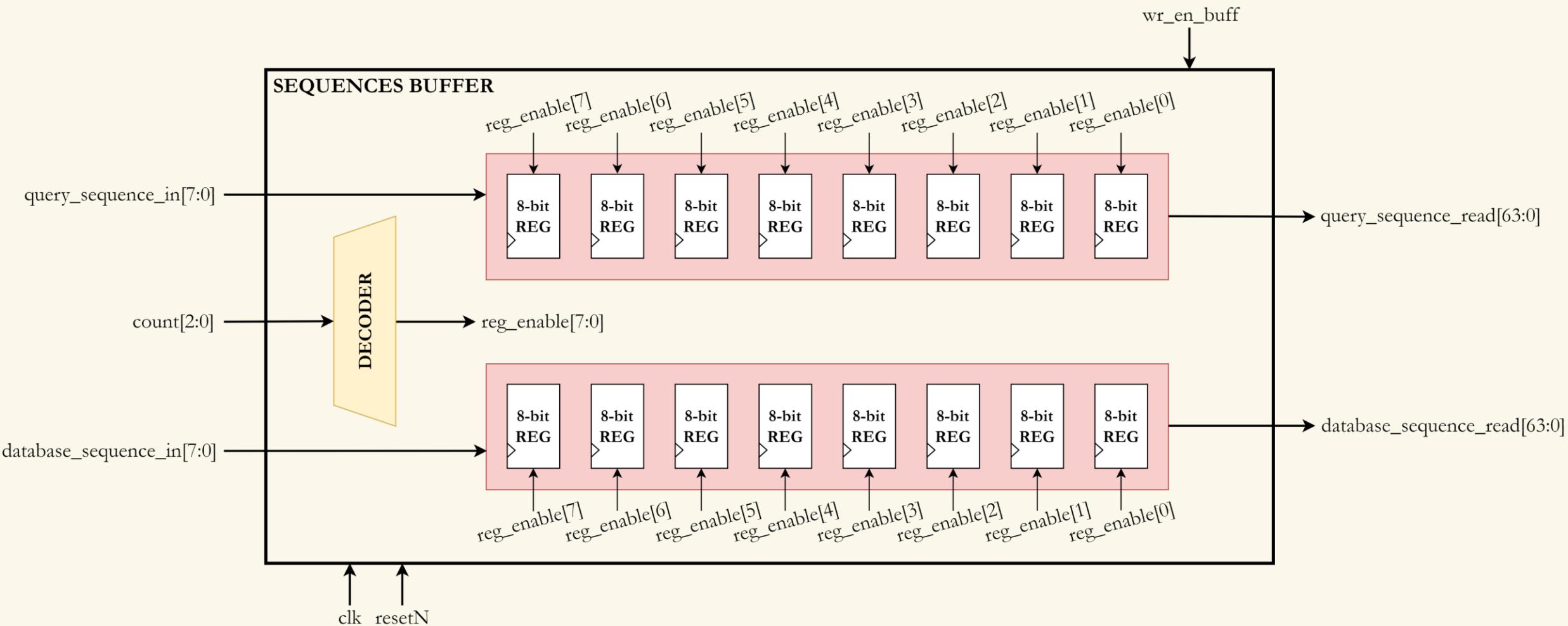
# Architectural Design

---

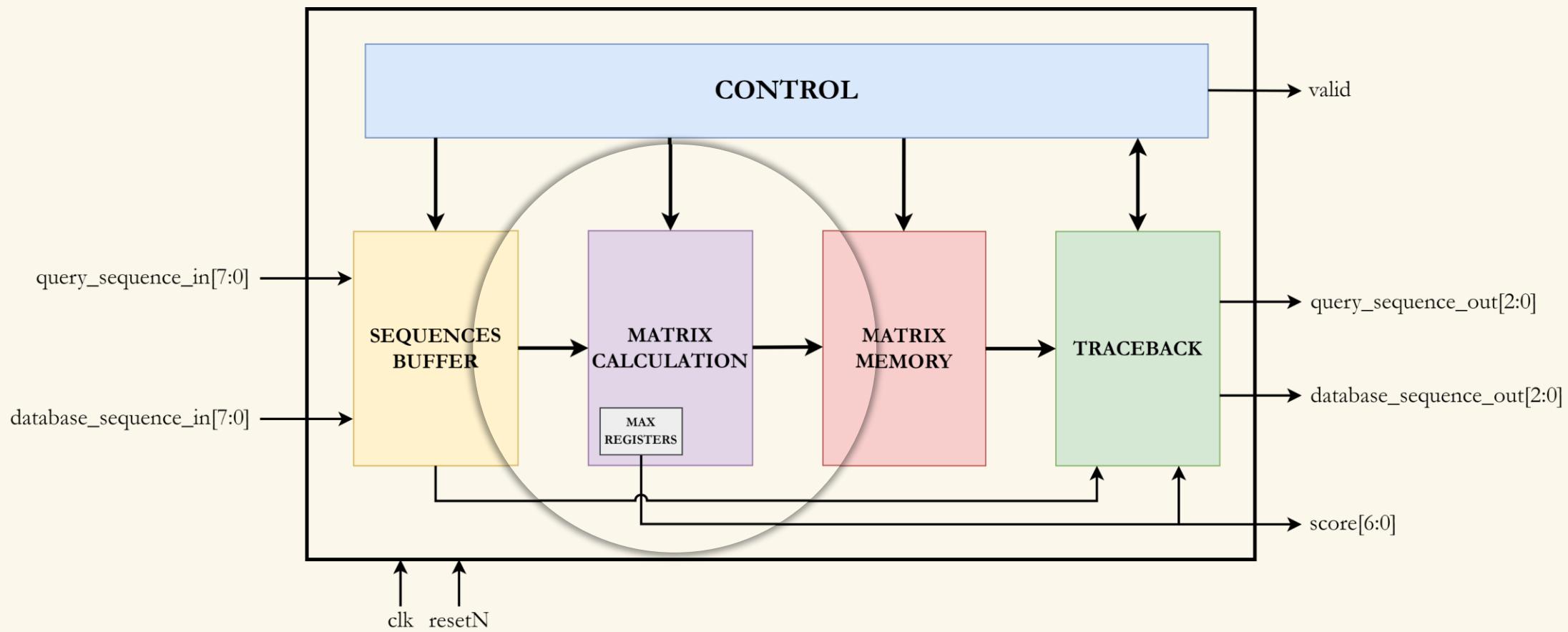
# Sequences Buffer



# Sequences Buffer



# Matrix Calculation





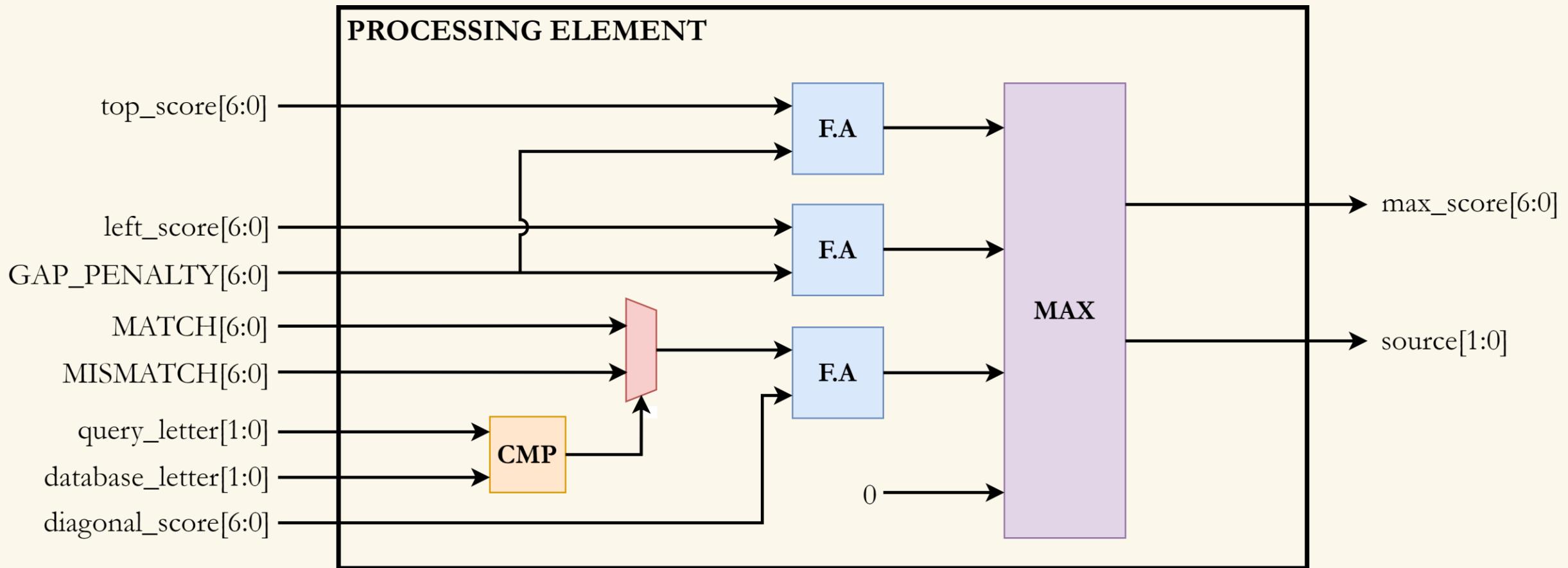




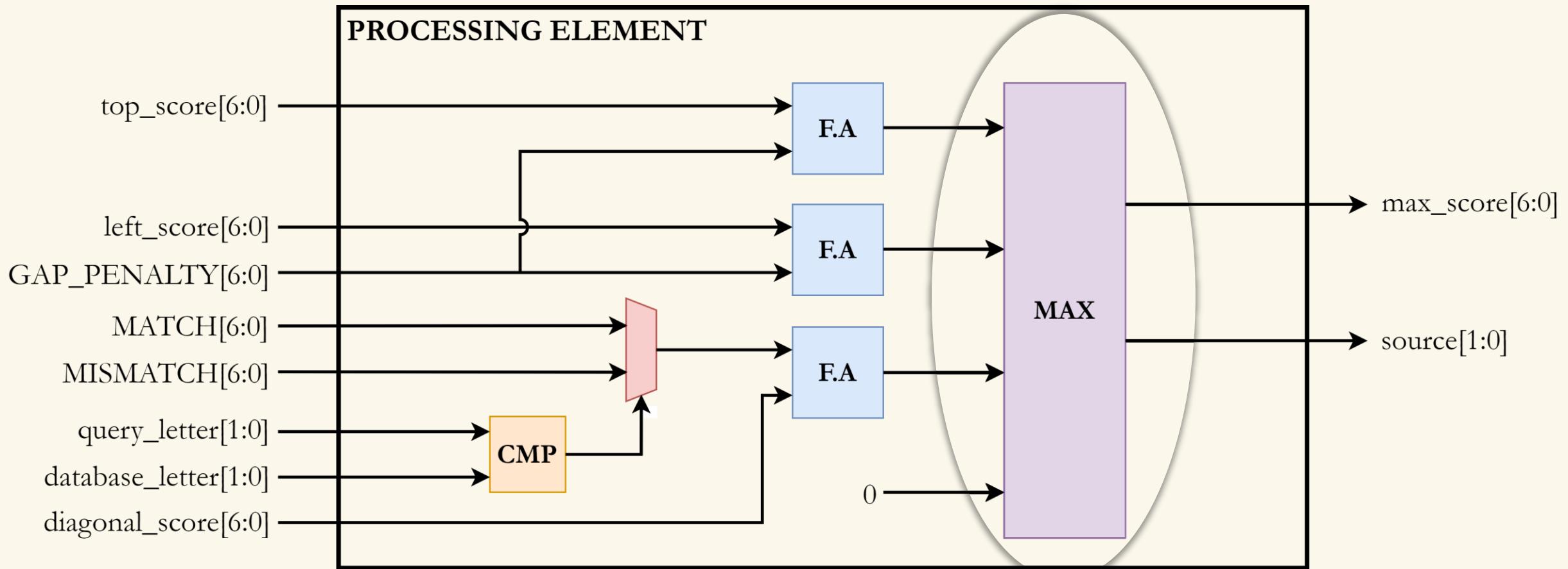




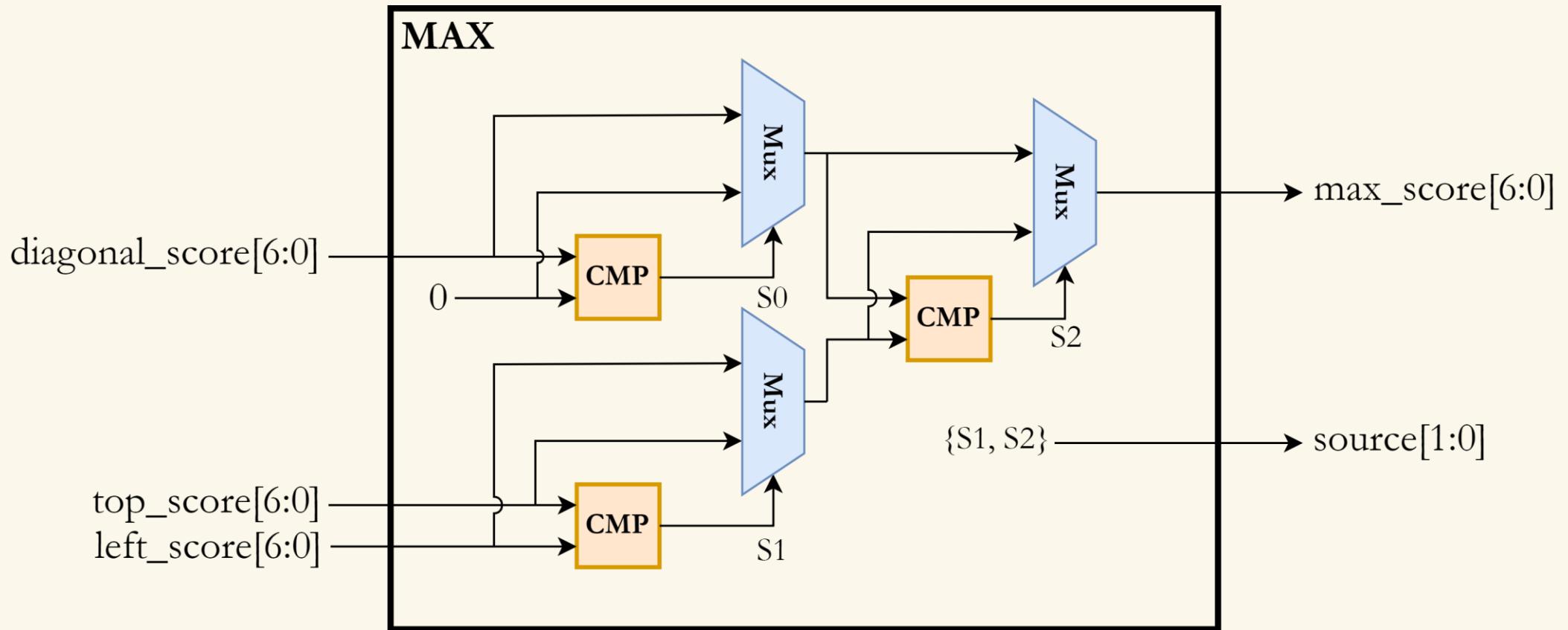
# Processing Element



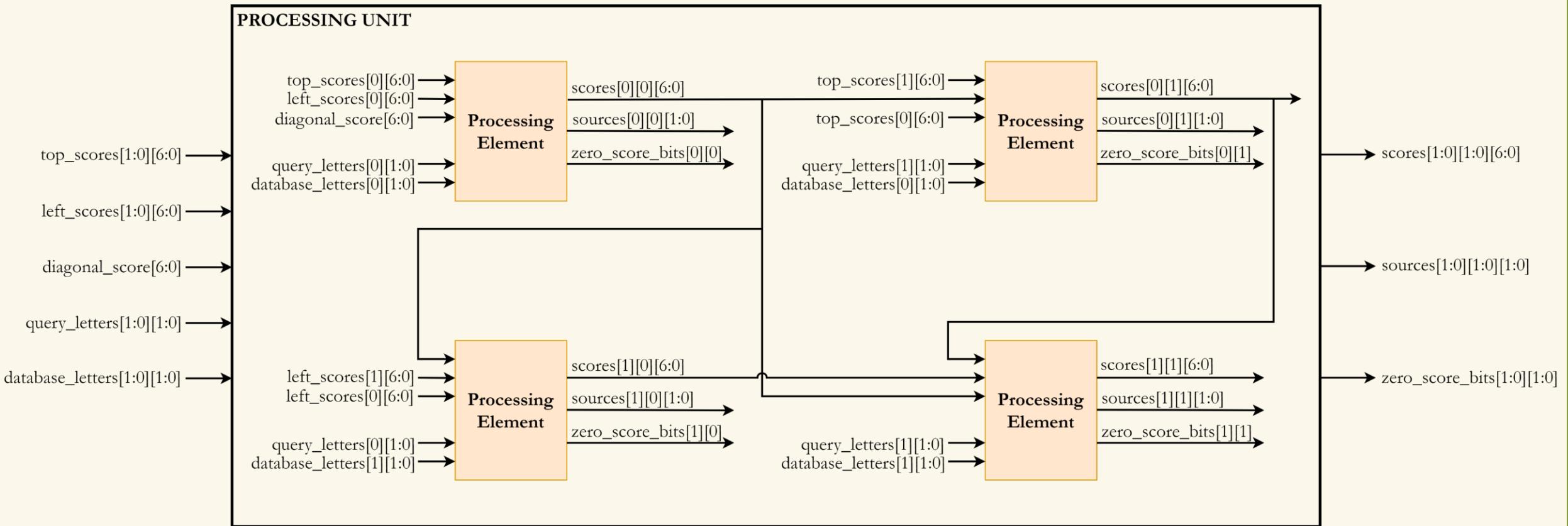
# Processing Element



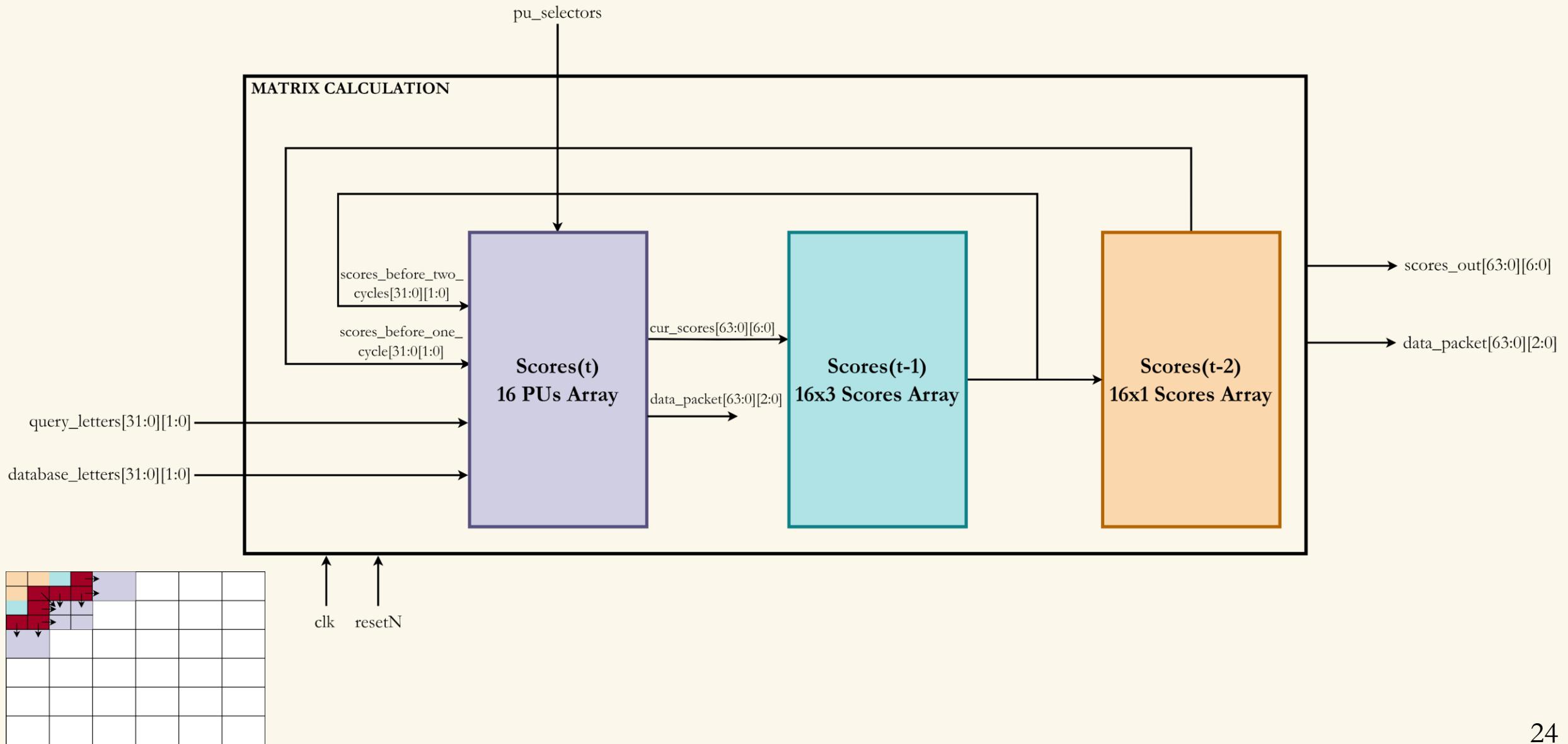
# Max Unit



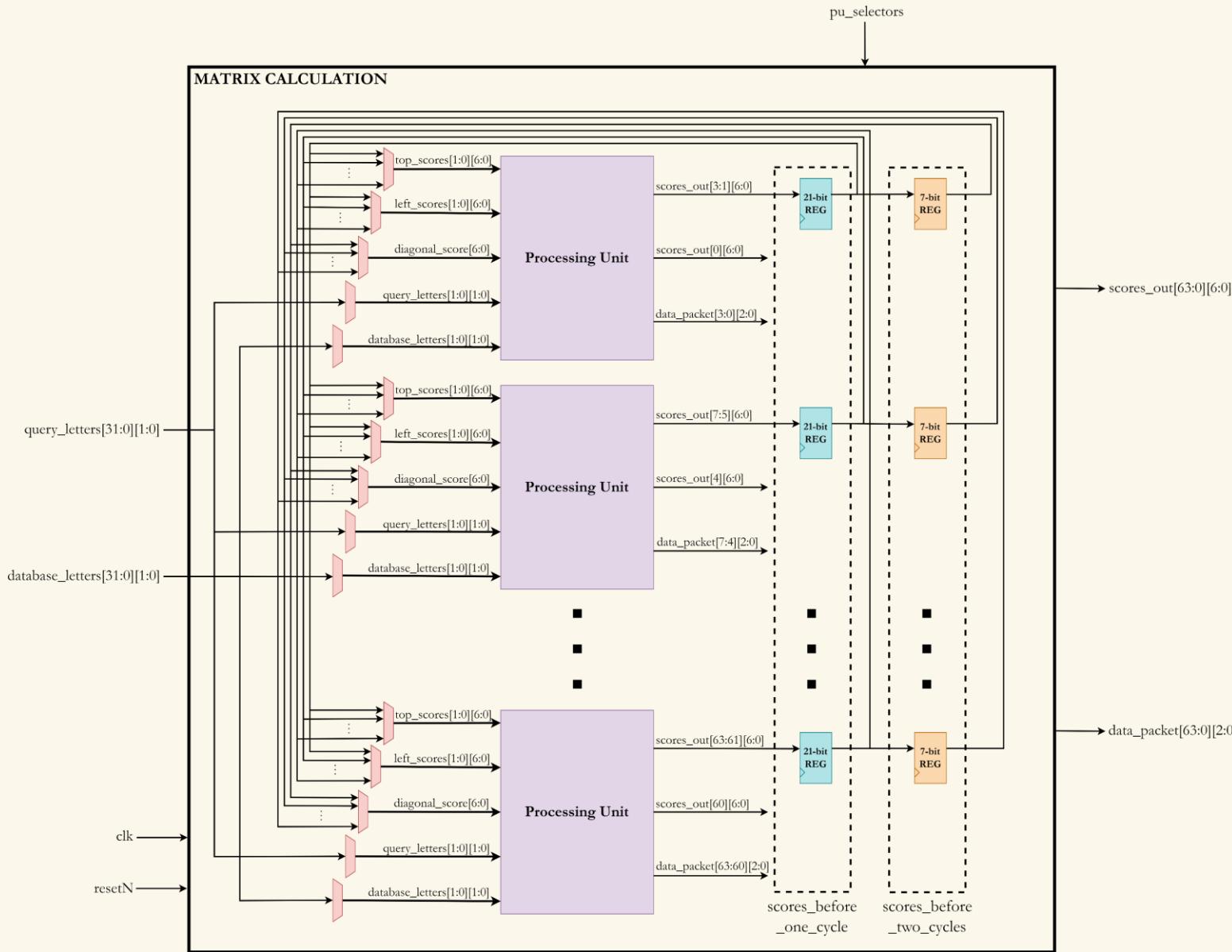
# Processing Unit



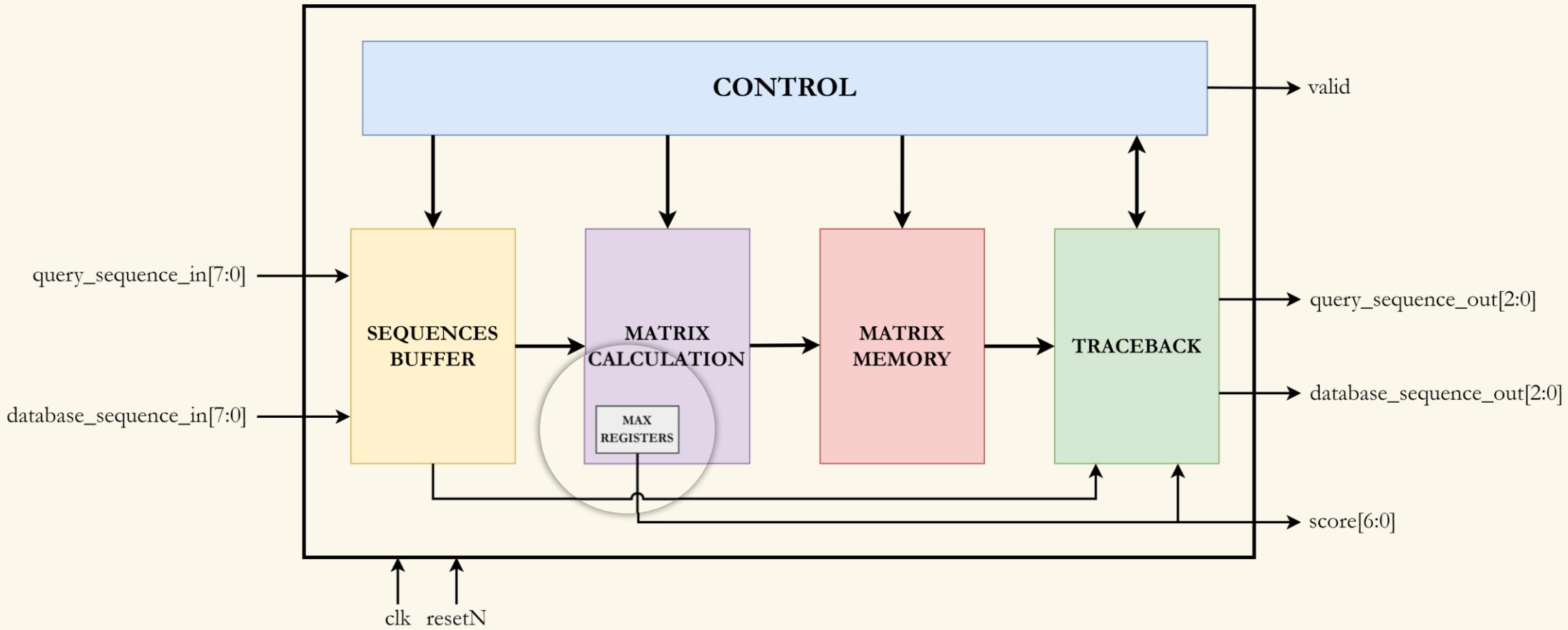
# Score Matrix Computation



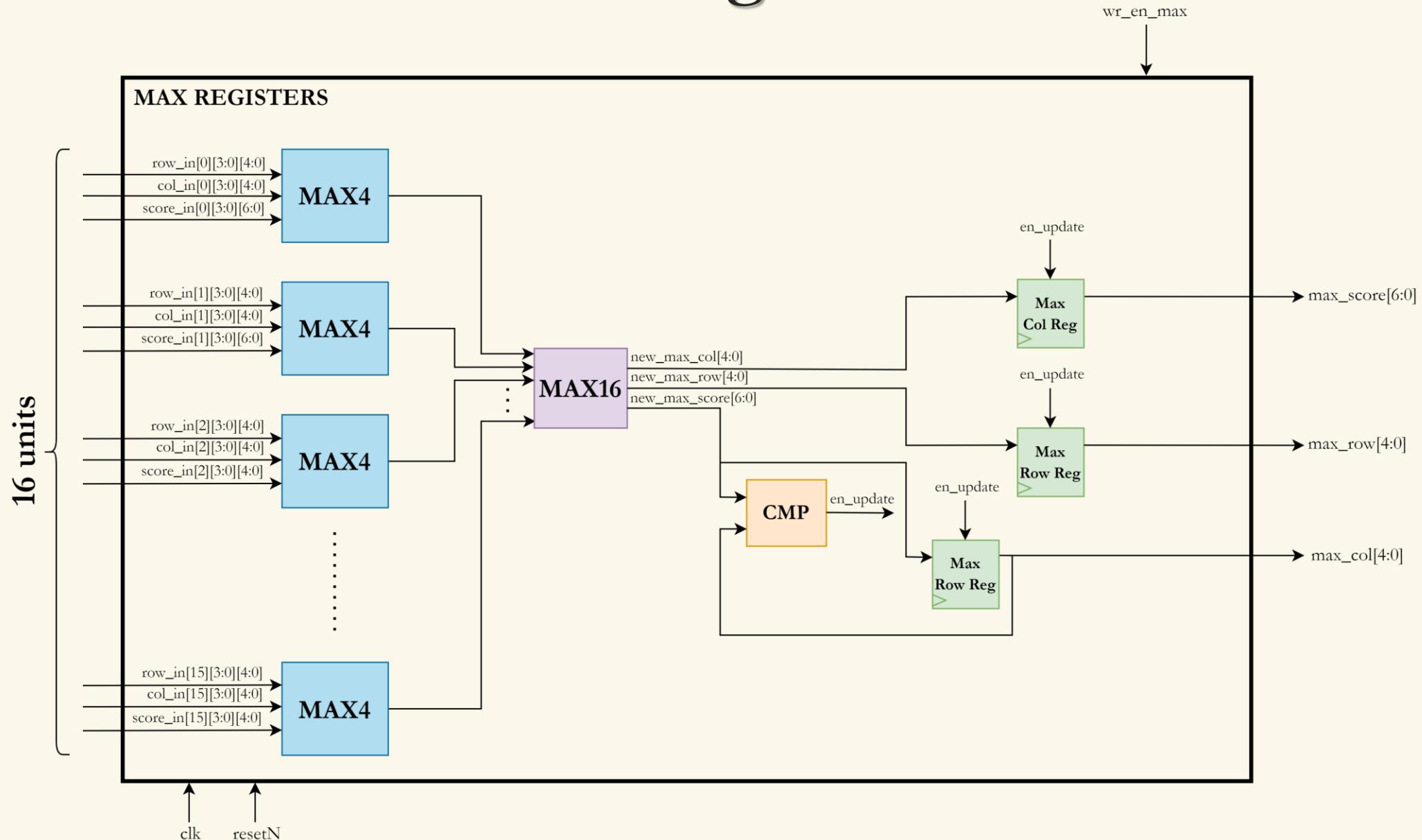
# Score Matrix Computation



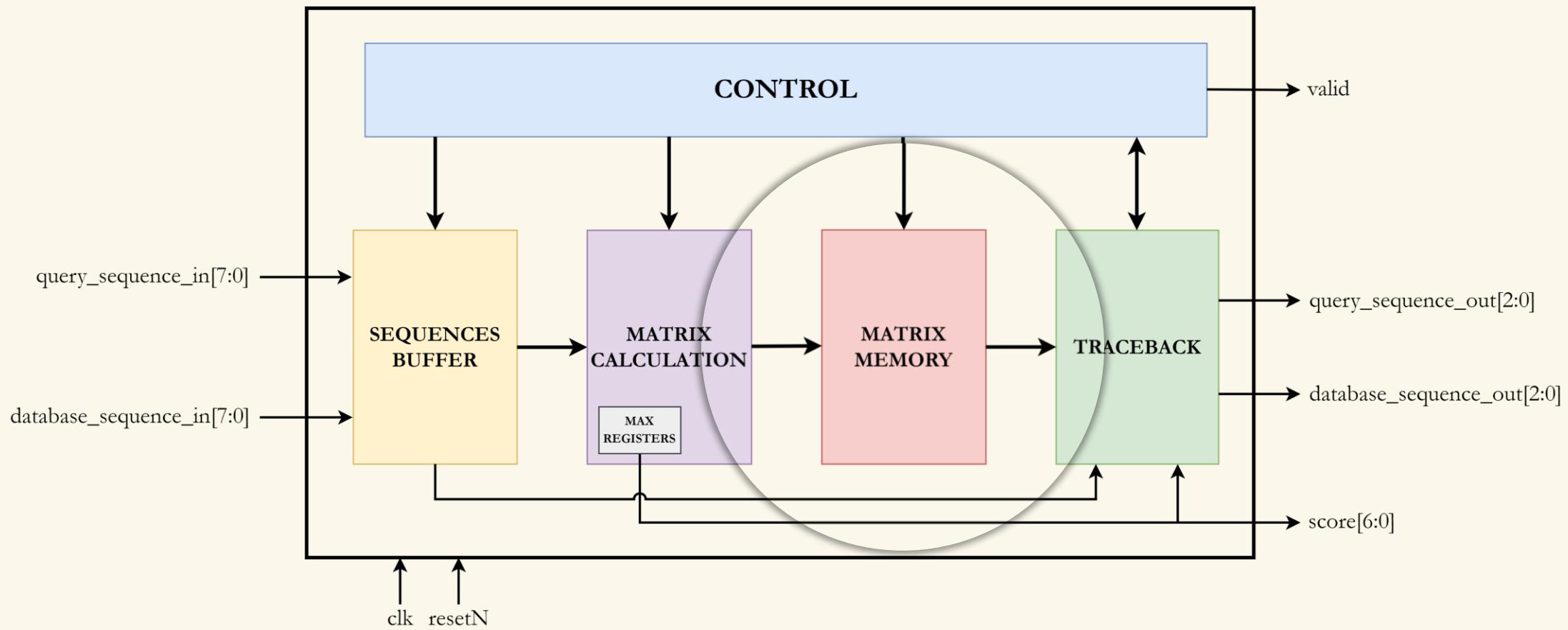
# Max Registers



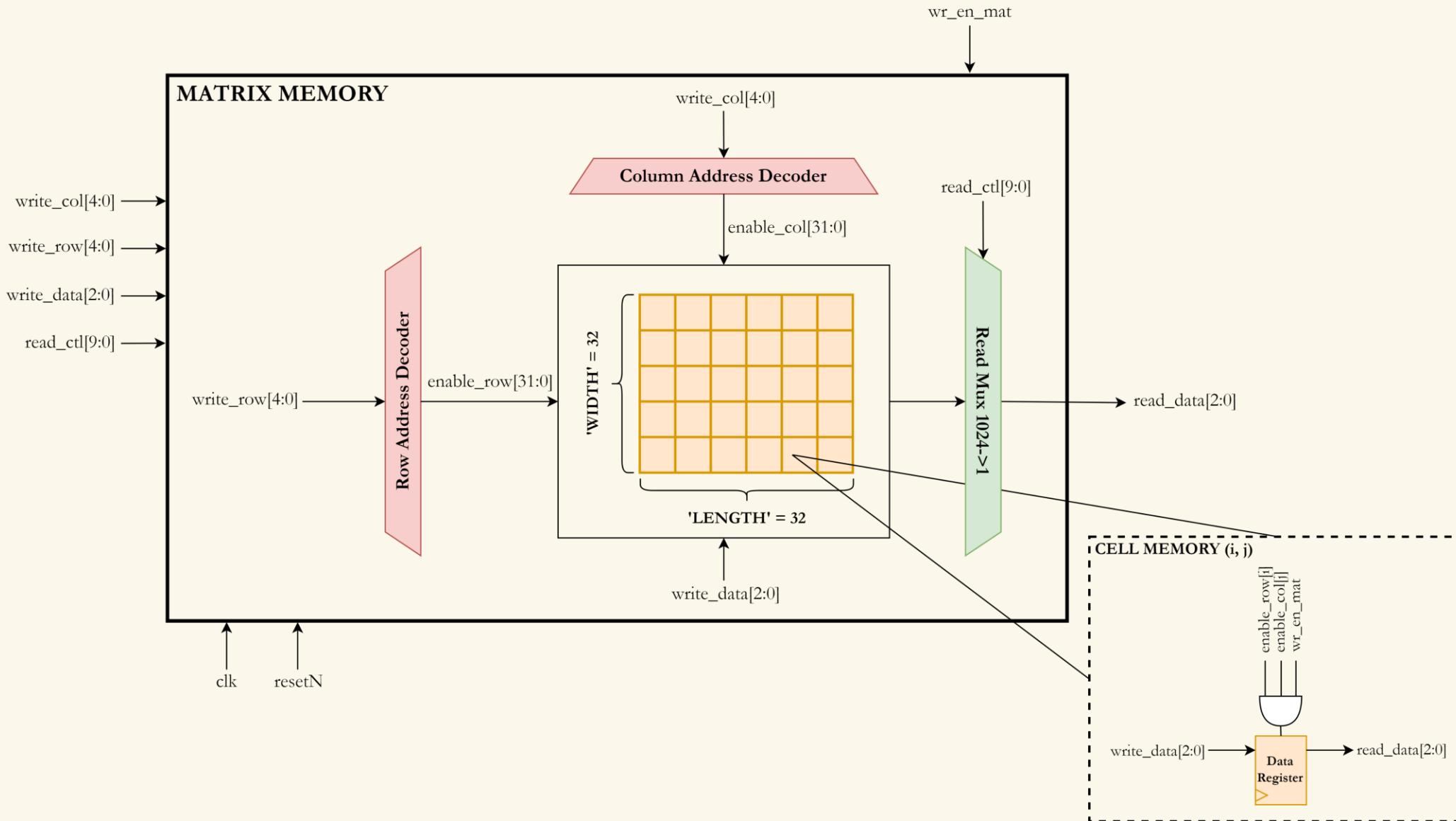
# Max Registers



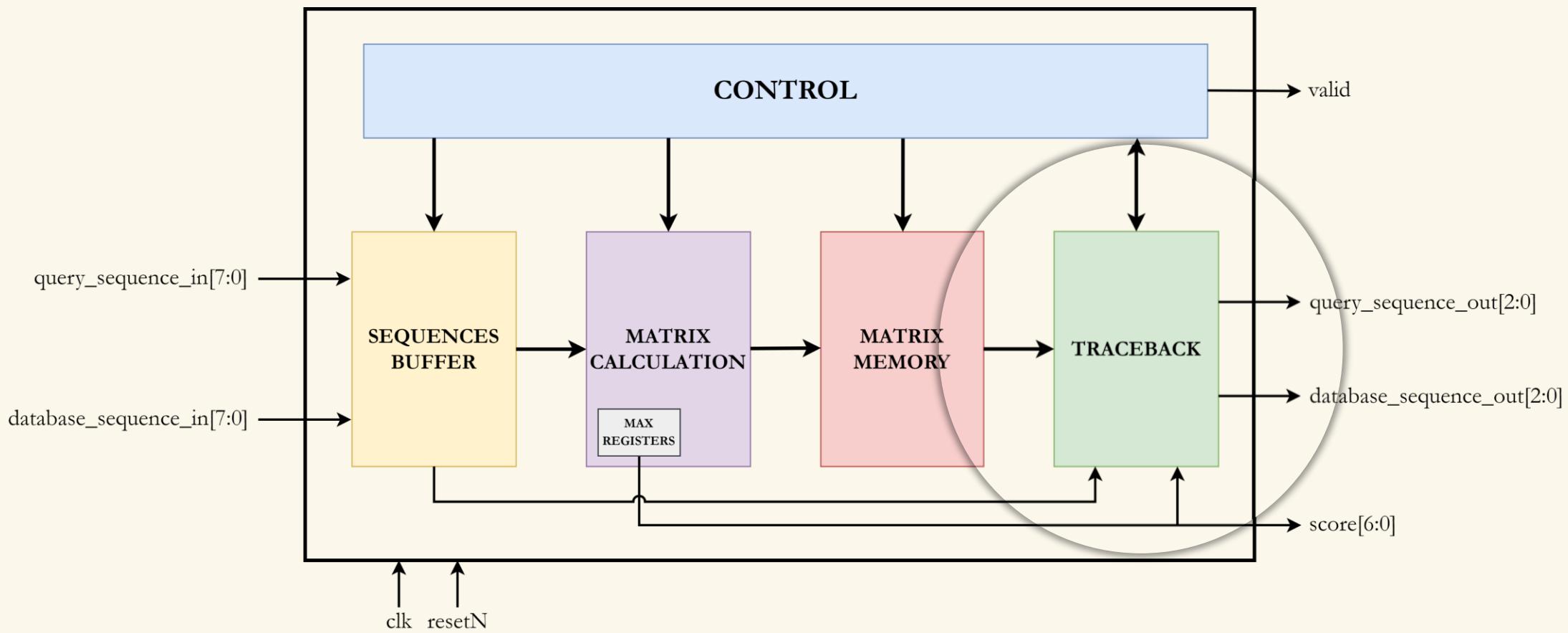
# Matrix Memory



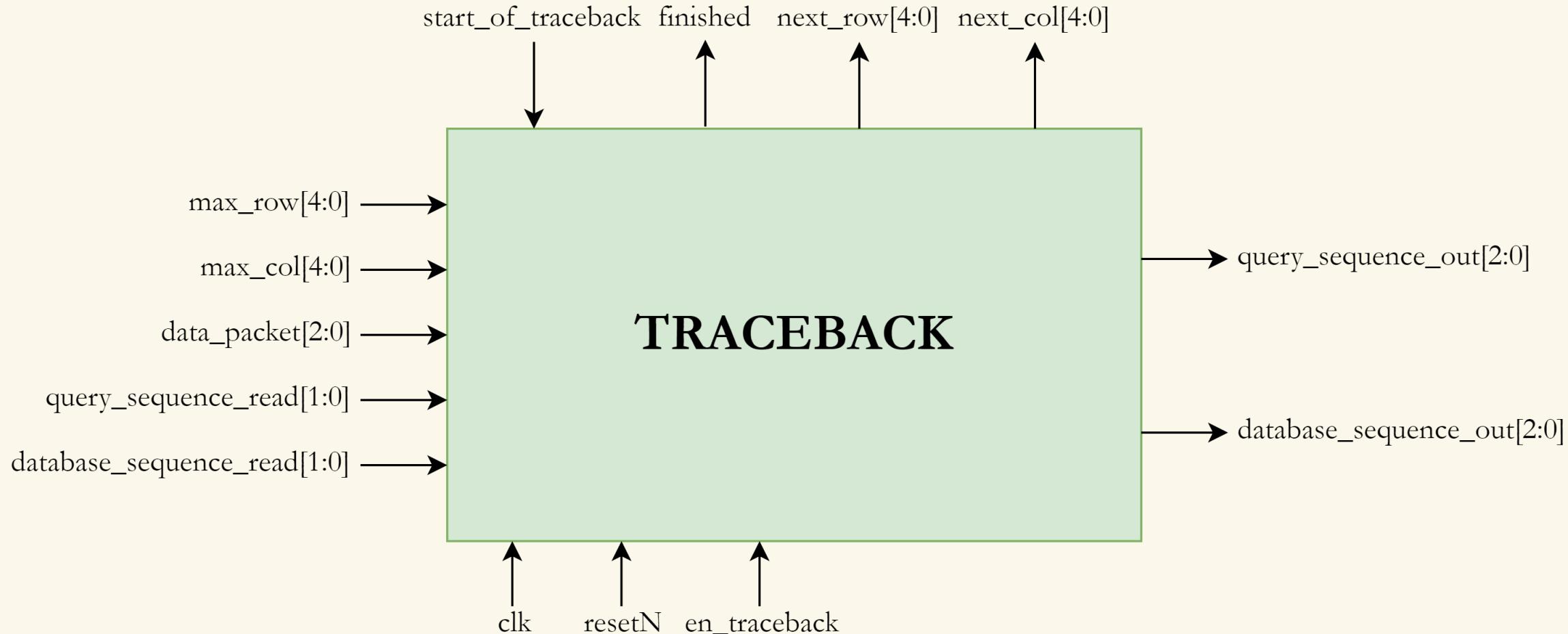
# Matrix Memory



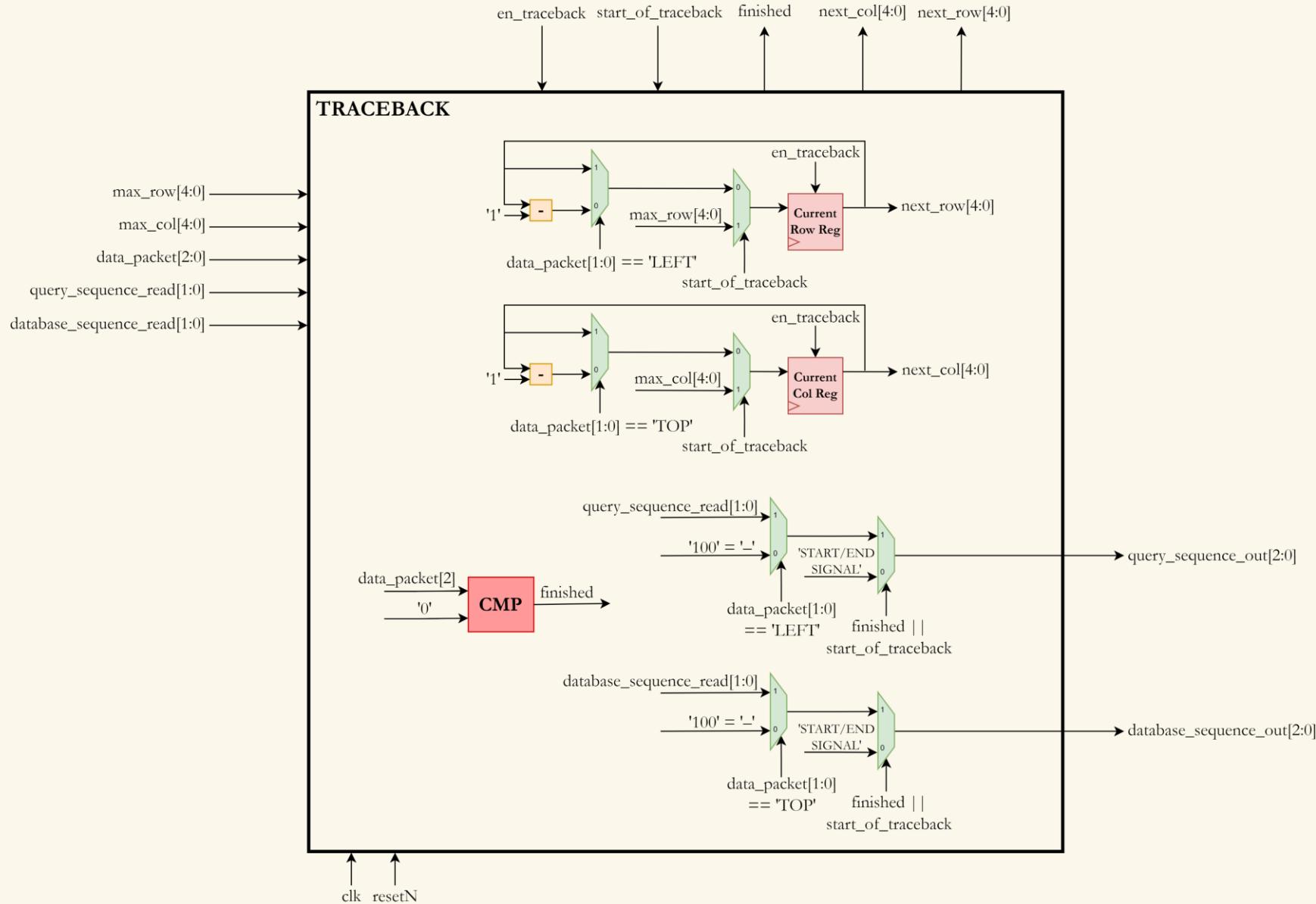
# Traceback



# Traceback – Interface



# Traceback



# Next Steps

- Finalize full system RTL simulations
- Synthesis
- Physical design
- Tape out (by end of December 2023)
- Post-silicon Testing



# THANK YOU!