

Assignment

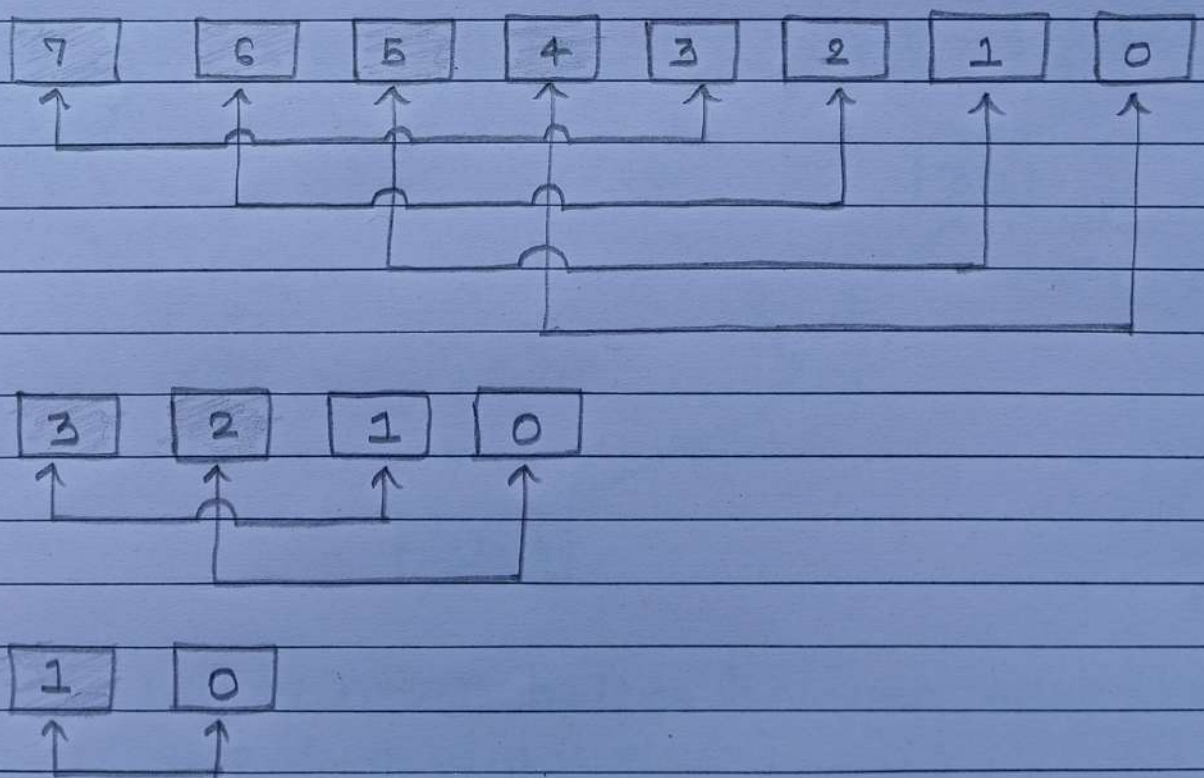
- Title :-
Implement Min, Max, Sum and Avg operations using parallel reduction.
- Objective:-
Understand about the Concept of Min, Max, Sum & Avg operations using parallel reduction.
- Problem Statement :-
Perform the Min, Max, Sum & Avg operations using parallel reduction.
- S/W & H/W Requirements:-
 - 64 bit open source linux & Windows & its derivatives.
 - CPU

• Theory :-

Parallel Reduction :-

- Parallel reduction works by using half number of threads of elements in dataset.
- Every thread calculates minimum of its own element and some other element. The resultant element is forwarded to next.

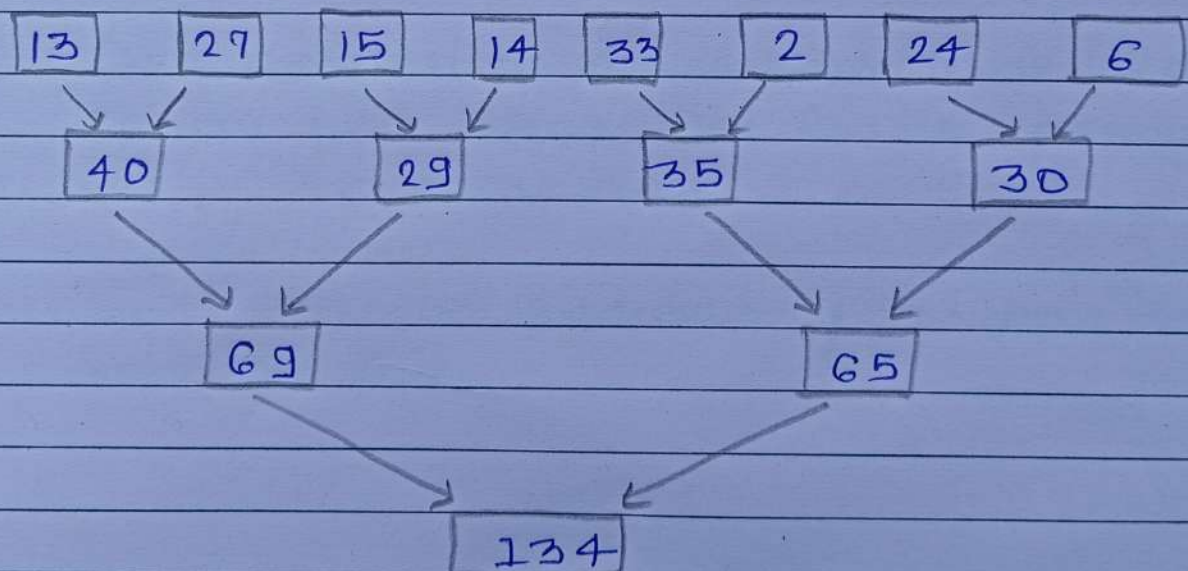
- The number of threads is then reduced by half and then process repeated, until there is just single element remaining which is result of a operation.
- With CUDA you must remember that execution unit for given SM is a Warp. Thus, any amount of threads less than one Warp is underutilization hardware.
- Following fig. See item being Compared With one from other half of dataset.



- Parallel reduction refers to algorithms which combine an array of elements producing a single value as a result.
- Problems eligible for their algorithm include those which involve operators that are associative and commutative in nature.

Parallel Sum

Adding values is an associative operation. So, we can try something like this,
 $((13+27) + (15+14)) + ((33+2) + (24+6))$



This way is much better because now we can execute it in parallel.

Parallel Programming

- Task Parallelism

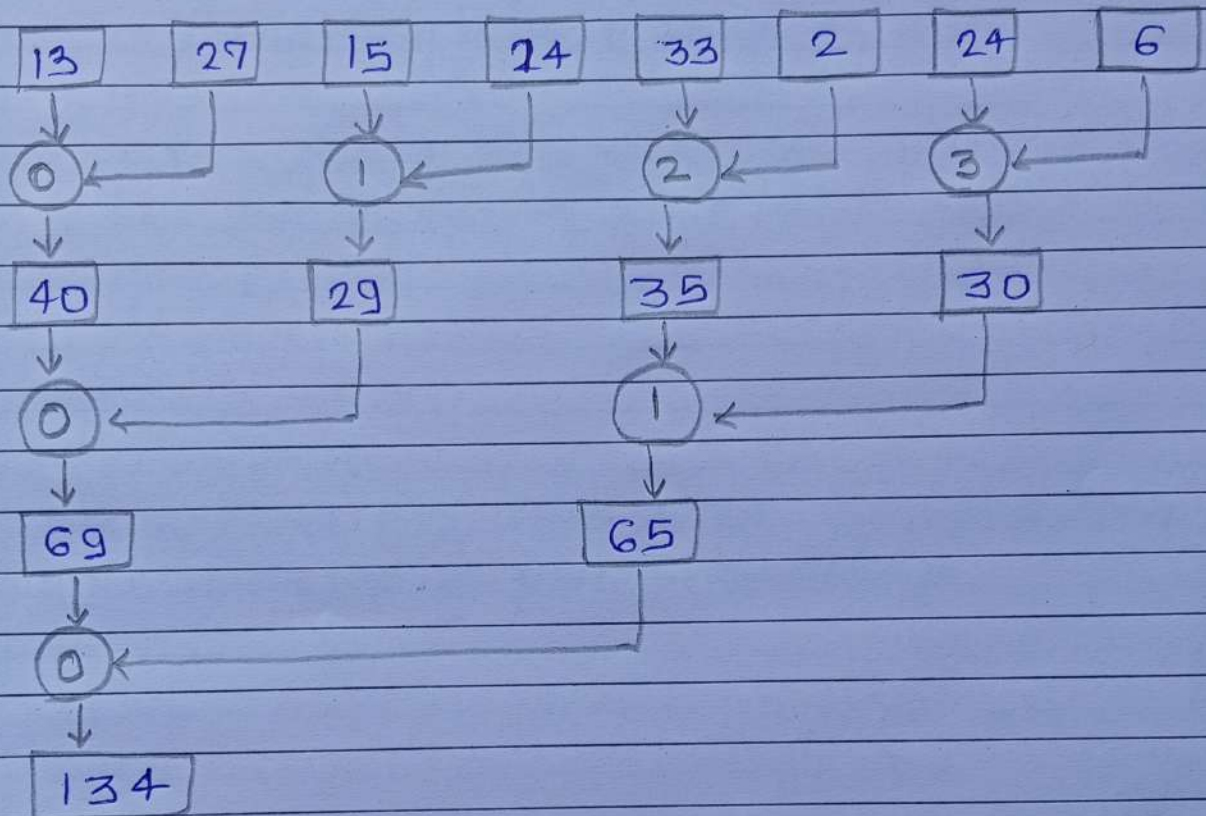
- Task parallelism arises when there are many tasks or functions that can be operated independantly and largely in parallel.
- It focuses on distributing functions across multiple cores.

- Data Parallelism

- Data parallelism arises when there are many data items that can be operated on at the same time.
- It focuses on distributing the data across multiple cores.

CUDA

Let's figure out using CUDA.



- Assume N as the no. of elements of the elements in an array, we start $N/2$ threads, one thread for every 2 elements.
- Each thread computes the sum of the corresponding 2 elements, storing the result at the position of first one.
- Iteratively each Step :

- The no. of threads halved
(for e.g. Starting With 4 then 2, 1)
 - Doubles the Step Size between Corresponding 2 elements
(Starting With 1, then 2, 4)
- After Some iterations, the reduction result will be stored in the first element of the array.

• Conclusion :-

We perform the Min, Max, Sum and Avg operations using parallel reduction.