# OpenCL exercise 6: Prefix sum

Kaicong Sun

# Prefix sum

- Prefix sum = all prefix sums for an input vector
- For input values $x_0, x_1, x_2, \ldots$ compute:

$$
\begin{aligned}
y_0 &= x_0 \\
y_1 &= x_0 + x_1 \\
y_2 &= x_0 + x_1 + x_2
\end{aligned}
$$

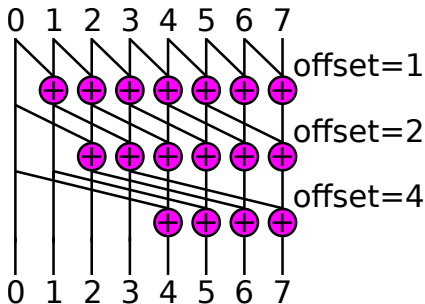- Also can be some other associative binary operation instead of $+$, e.g. min, max, ...

# Prefix sum

Host code:

```
1 cl_int sum = h_input[0];
2 h_output[0] = sum;
3 for (std::size_t i = 1; i < h_input.size (); i++) {
4     sum += h_input[i];
5     h_output[i] = sum;
6 }
```

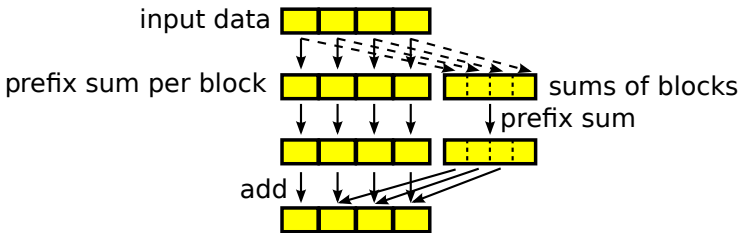# Parallel prefix sum

▶ Parallel prefix sum:

# GPU

- ► Task: Implement prefix sum on GPU
  - ► Plus usual code for performance measurements
- ► Kernel should:
  - ► Load input data to local memory
  - ► Loop over offsets
  - ► Write results to global memory
- ► Use one work item per value
- ► Do not forget to add `barrier` calls for synchronization

# GPU

- ▶ Problem: Can use only one work group
- ▶ Solution: Work with blocks



- ▶ Do prefix sum per block
    - ▶ Also write sum of block to another array, `temp1`
- ▶ Do prefix sum for `temp1` (recursivly, using `temp2` as temp array)
- ▶ For all blocks except the first: Add `temp2[blockIndex-1]` to all values in the current `temp1` block and then recursively add `temp1[blockIndex-1]` to all values in the current `d_output`.(write second kernel for this step)