

The background features a complex network of blue lines and arrows. Solid lines intersect at various angles, while dashed lines form loops and paths. Small circles, some solid and some hollow, are placed at various points along these lines, suggesting a network or data flow.

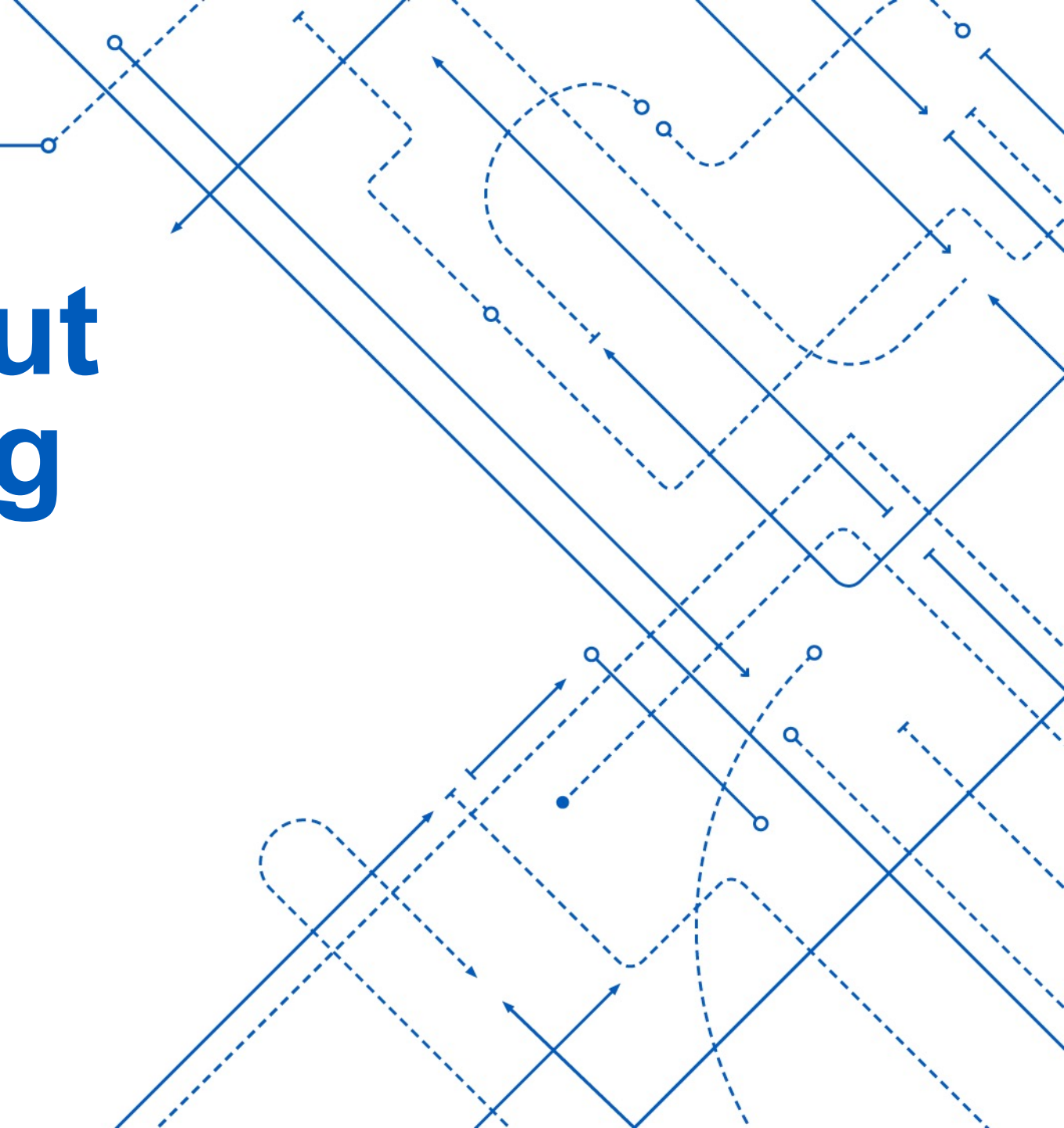
Implementation of Parallel Bitonic Sort using MPI and Intel TBB

Presented for CSE702 Fall 2021

Instructor: Dr. Russ Miller

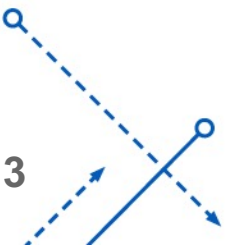
Presenter: Zainul Abideen Sayed

Why care about parallel sorting algorithms ?



Why parallel sorting?

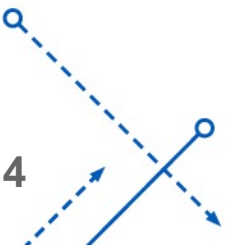
- Sorting is one of the most common operations in computation.
- The advancement in parallel hardware.
- Increasing nodes in a cluster and cores in a processor.
- Efficient utilization of resources.
- Therefore, good parallel sorting algorithms are needed.



Popular parallel sorting algorithms

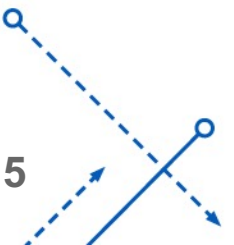
- Bitonic sort
- Sample sort
- Merge sort
- Quick sort
- Radix sort

Bitonic sorting algorithm is based on bitonic sorting network. The key operation is based on the sorting network which converts a given sequence into a bitonic sequence and finally bitonic merge can produce a monotonically increasing or decreasing sequence.

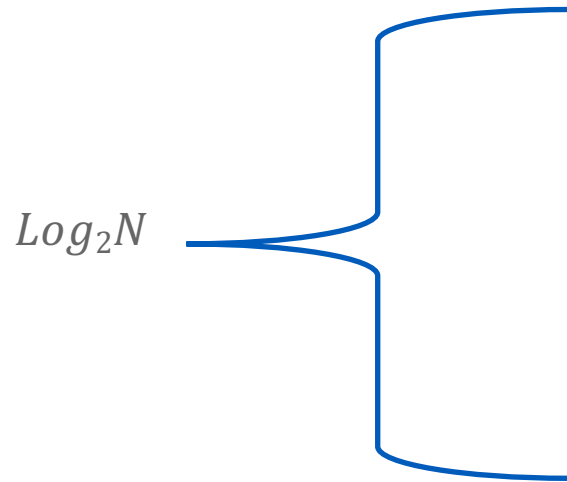


Bitonic Sort Principle

- Bitonic sequence $\langle 1, 2, 4, 7, 6, 0 \rangle$ $\langle 8, 9, 2, 1, 0, 4 \rangle$ $\langle 0, 4, 8, 9, 2, 1 \rangle$ $\langle 3, 4, 7, 8, 6, 5, 2, 1 \rangle$
- Let $s = \langle a_0, a_1, \dots, a_{n-1} \rangle$
- $s_1 = \{ \min(a_0, a_{n/2}), \min(a_1, a_{n/2+1}), \dots, \min(a_{n/2-1}, a_{n-1}) \}$
- $s_2 = \{ \max(a_0, a_{n/2}), \max(a_1, a_{n/2+1}), \dots, \max(a_{n/2-1}, a_{n-1}) \}$
- In sequence s_1 , there is an element $b_i = \min\{ a_i, a_{n/2+i} \}$ such that all the elements before b_i are from the increasing part of the original sequence and all the elements after b_i are from the decreasing part.
- Opposite case for $b_i = \max\{ a_i, a_{n/2+i} \}$



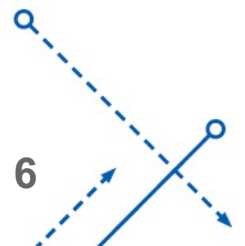
Example



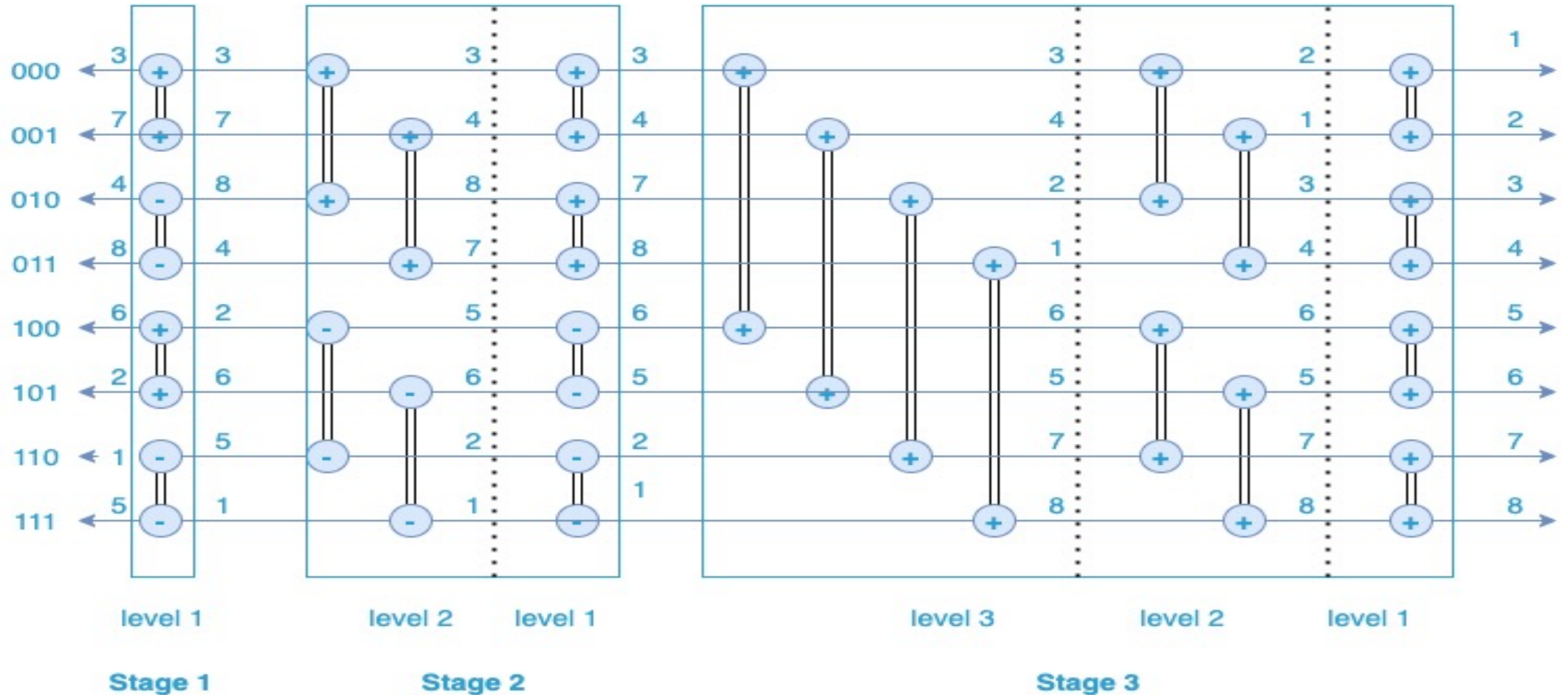
3 4 7 8 | 6 5 2 1
3 4 2 1 | 6 5 7 8
2 1 | 3 4 | 6 5 | 7 8
1 | 2 | 3 | 4 | 5 | 6 | 7 | 8



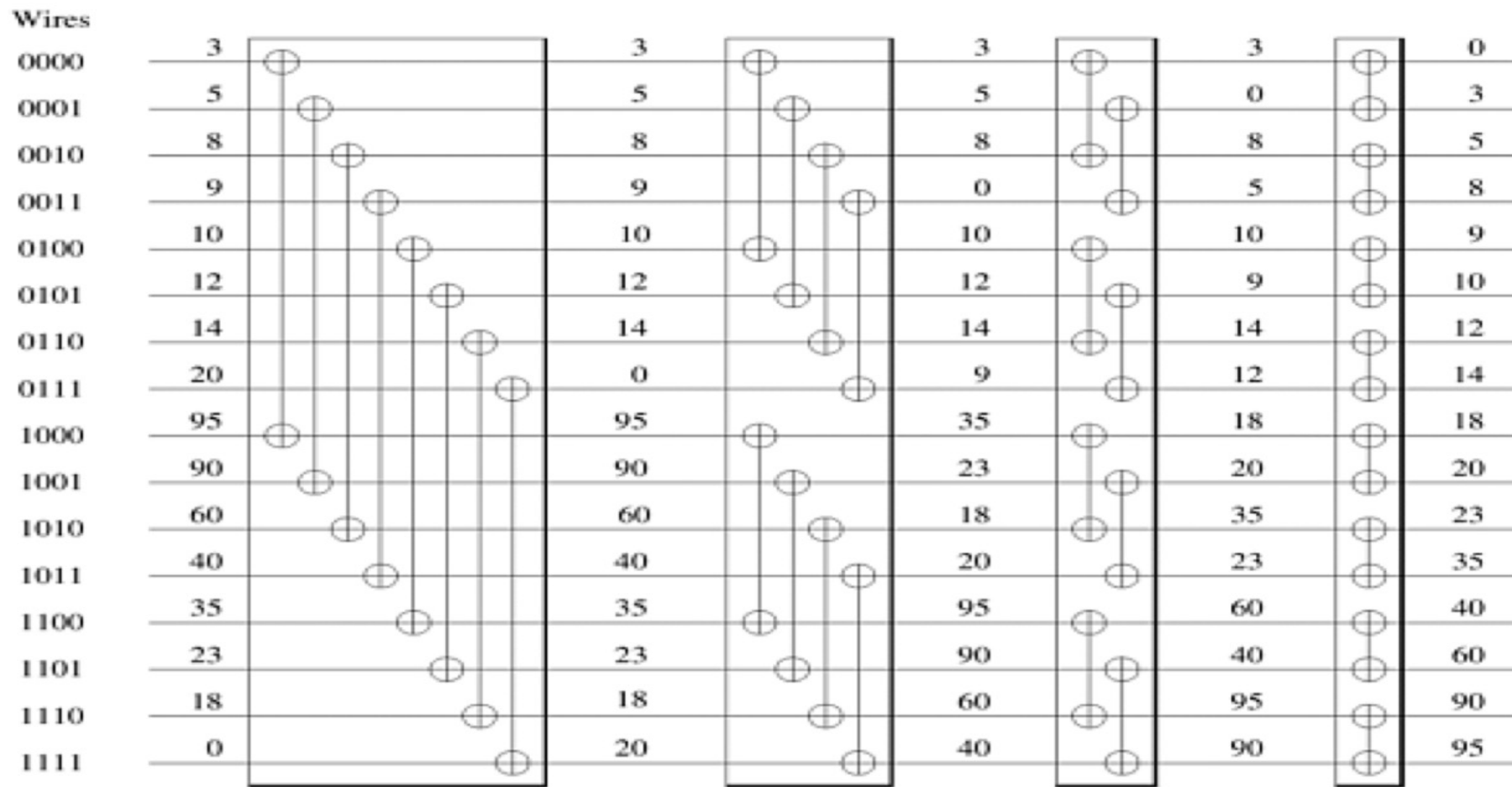
All elements in second sequence is greater than the first



Example: Sorting Network (Sort + Merge)



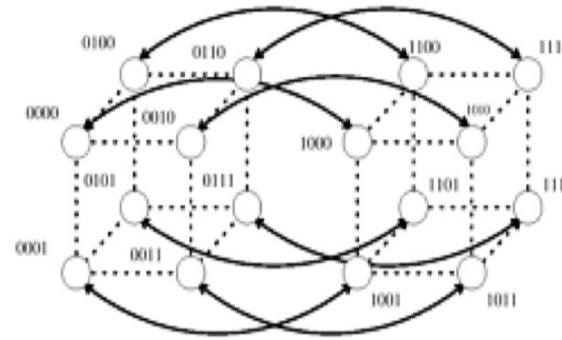
Example: 16 lines



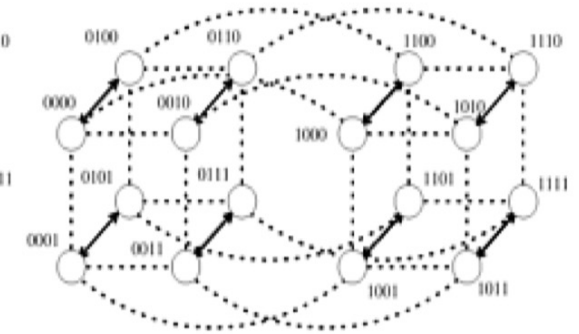
*Introduction to Parallel Computing 2nd Edition, Ananth Grama

Example: 16 lines

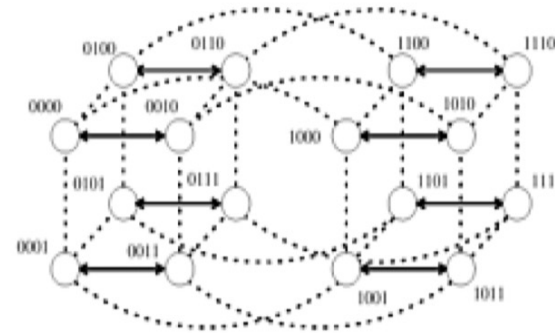
Hypercube



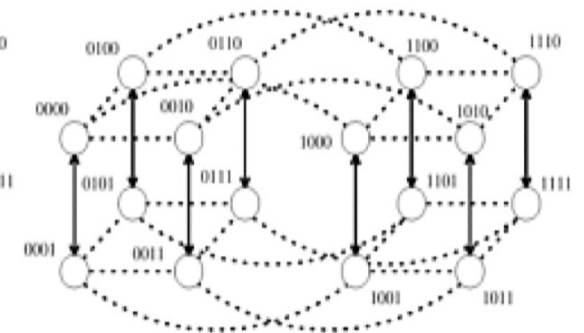
Step 1



Step 2



Step 3



Step 4

Algorithm

COMPARE SECTION

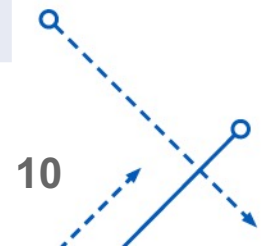
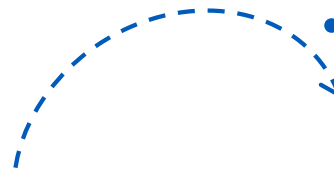
```

procedure BITONIC_SORT(label, d)
begin
  for i := 0 to d - 1 do
    for j := i downto 0 do
      if (i + 1) st bit of label != j th bit of label then
        comp_exchange_max(j);
      else
        comp_exchange_min(j);
    end
  end
end BITONIC_SORT
    
```

Complexity

- $(1 + \log n) (\log n) / 2$
- $T_p = \theta(n \log^2 n)$

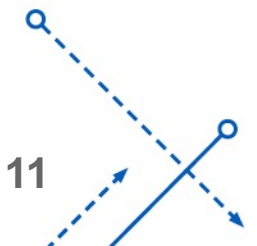
i	j	0	1	0	0
1	0	comp_exchange_min(0);			
2	1	comp_exchange_max(1);			
2	0	comp_exchange_max(0);			
3	2	comp_exchange_max(2);			
3	1	comp_exchange_min(1);			
3	0	comp_exchange_min(0);			



More on complexity...

- N/P Block of data per processor
- Fast sequential sort
 - Merge sort $\theta((n/p) \log(n/p))$
- Bitonic Merge
 - $\theta(\log^2 p)$

$$T_P = \overbrace{\Theta\left(\frac{n}{p} \log \frac{n}{p}\right)}^{\text{local sort}} + \overbrace{\Theta\left(\frac{n}{p} \log^2 p\right)}^{\text{comparisons}} + \overbrace{\Theta\left(\frac{n}{p} \log^2 p\right)}^{\text{communication}}.$$



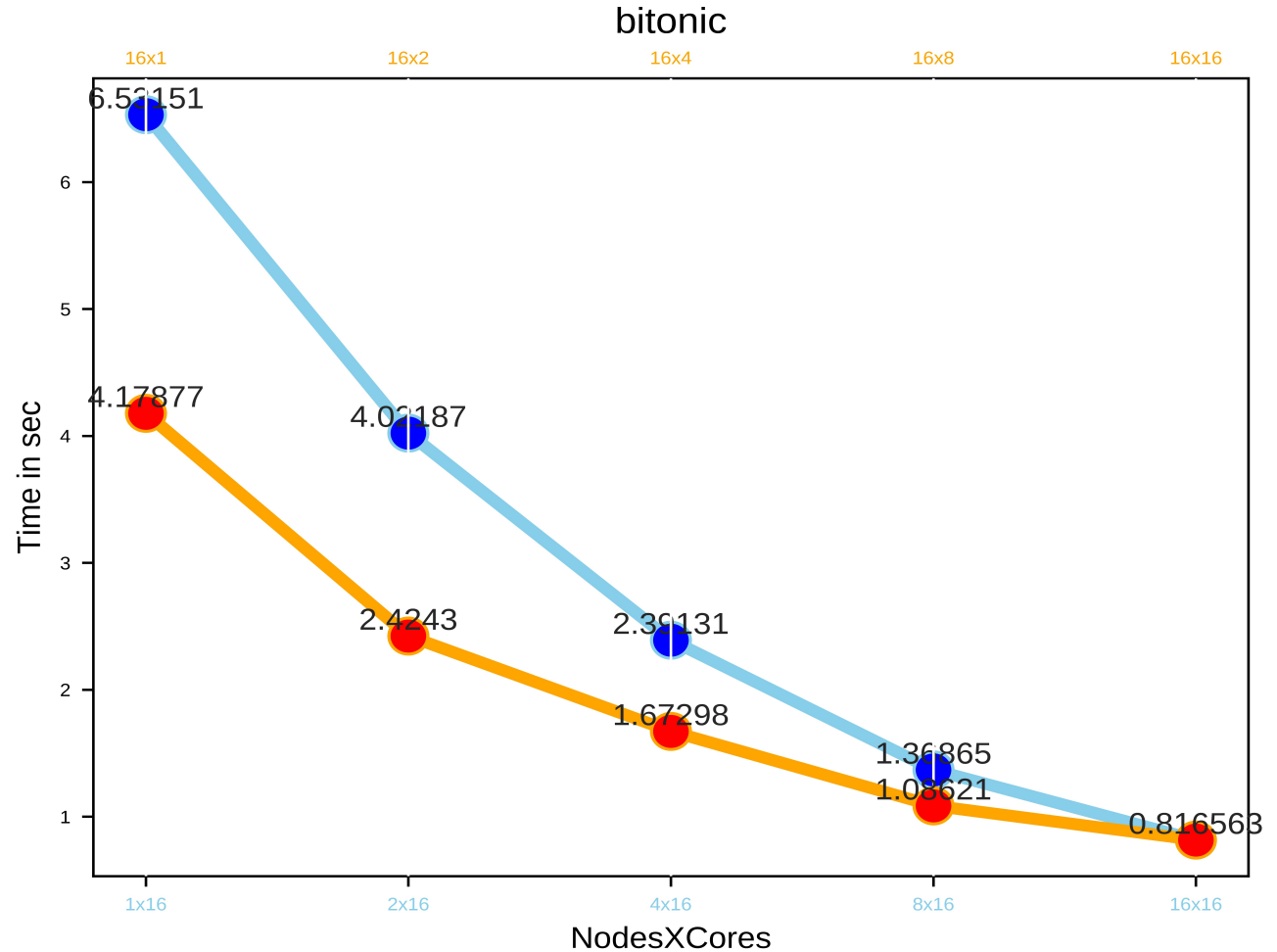
Results for 1 Billion keys

Key Observation: Inter-node parallelism give better performance than Intra-node parallelism.

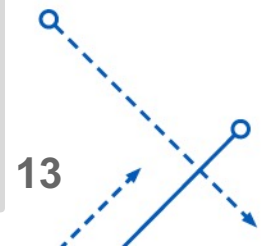
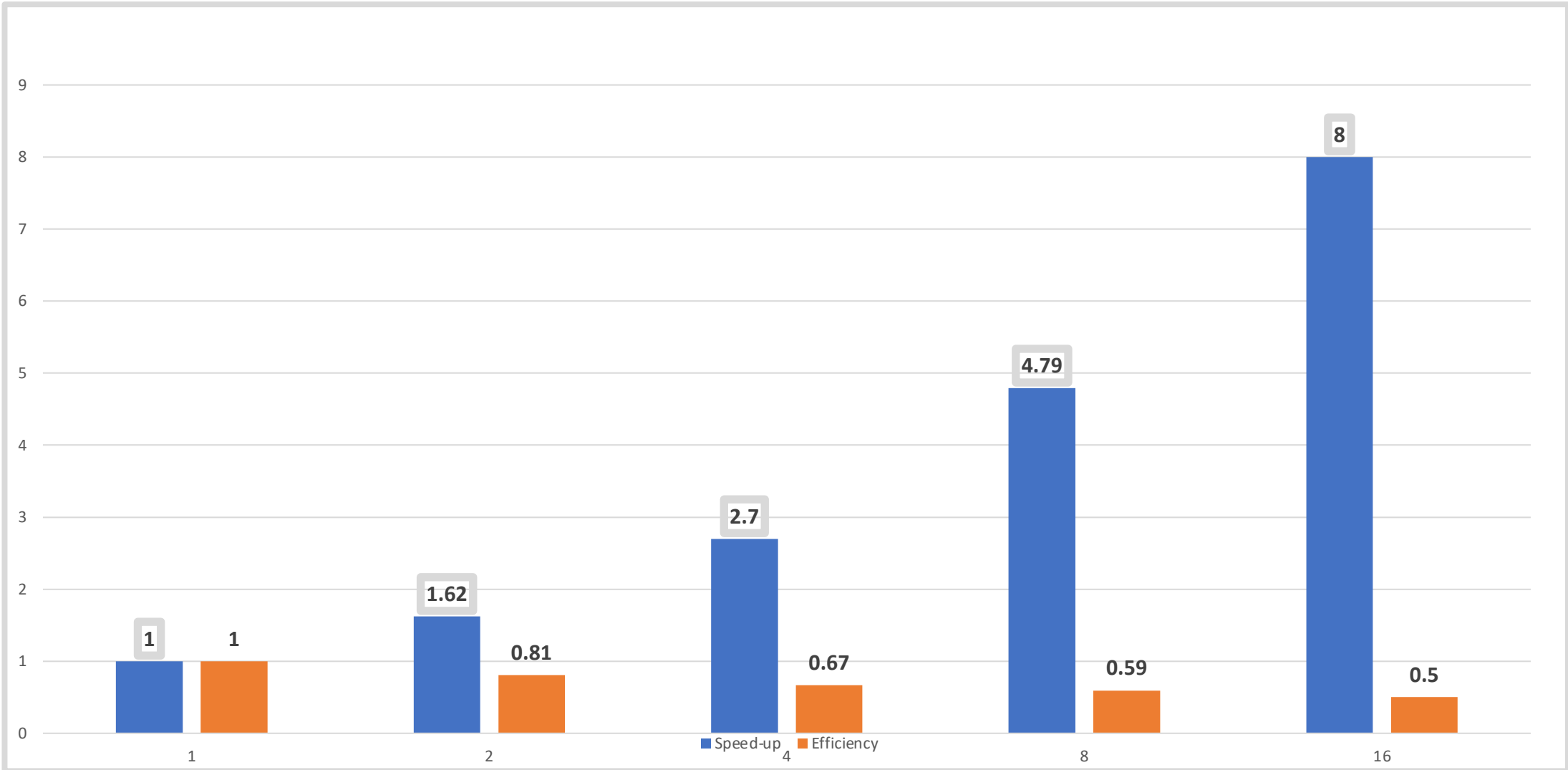
Relative speed-up:
 $T_1 = 6.53$; $T_{16} = 0.8$

$S_p = 8$
 $E_p = 0.5$

- Was able to sort 100 Billion keys on 512 ranks in 44.7524 sec
- Input Data size : 383G



Speed-up and Efficiency

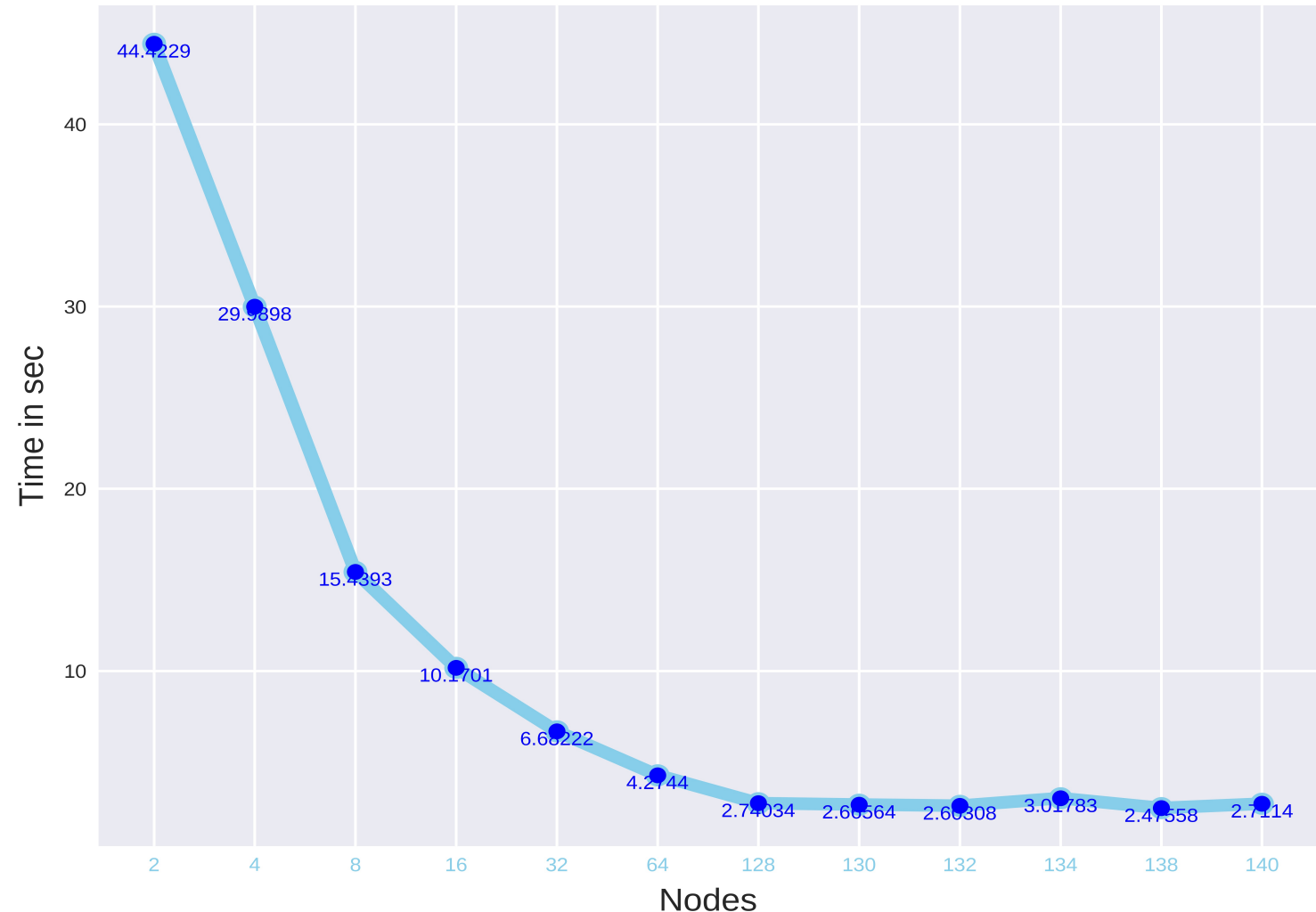


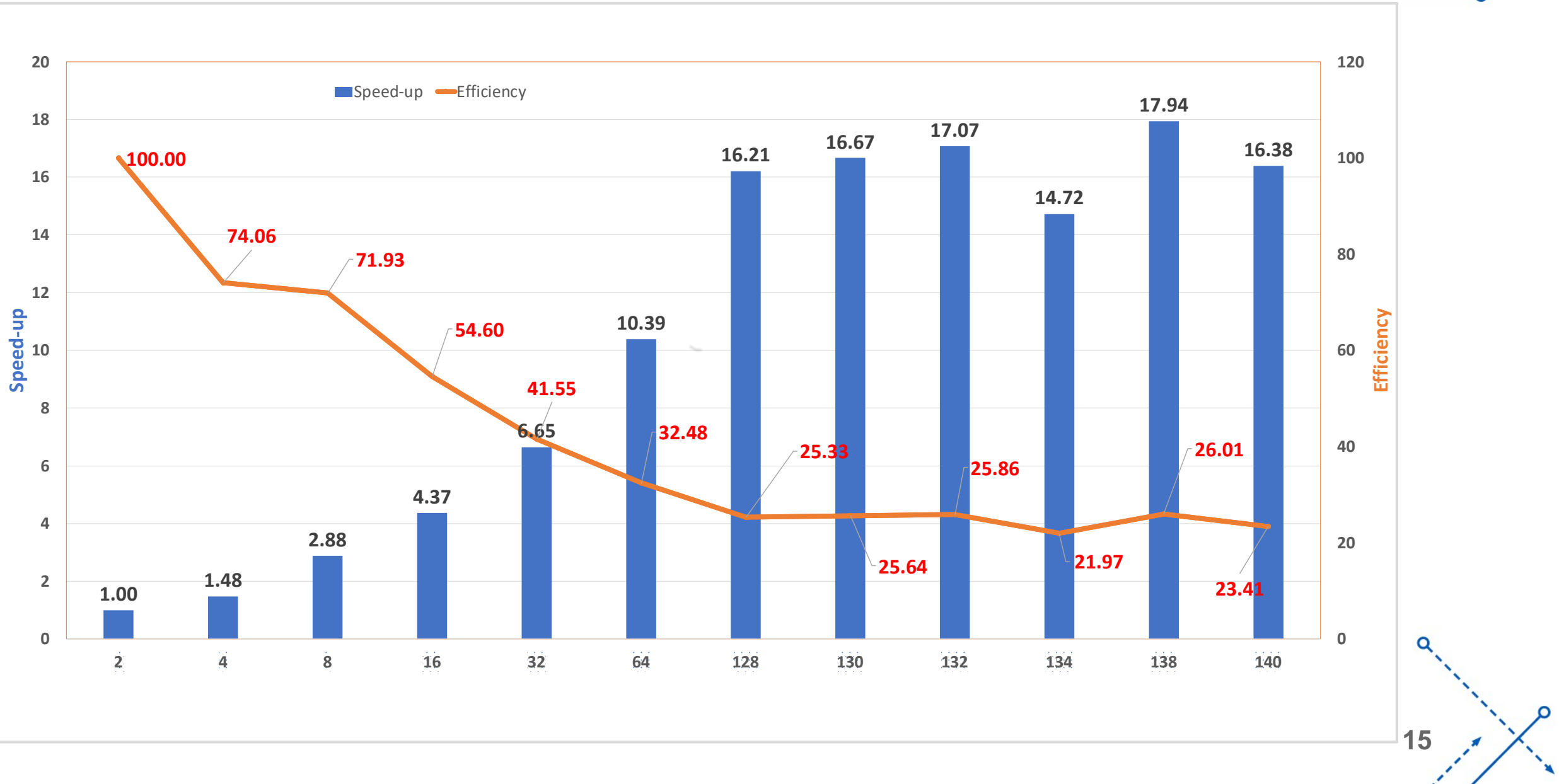
Results for 1 Billion keys

Observation:

- The code stops scaling after 128 Nodes.
- The execution time increases in some cases.
- However, It will be interesting to see the performance at 256 nodes.

Bitonic-Sort using Intel-MPI on Mellanox-Infiniband



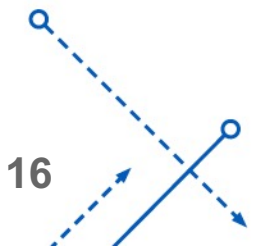


Amdahl's Law

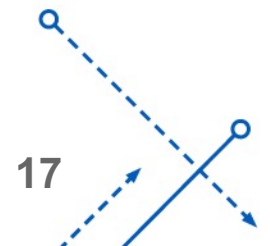
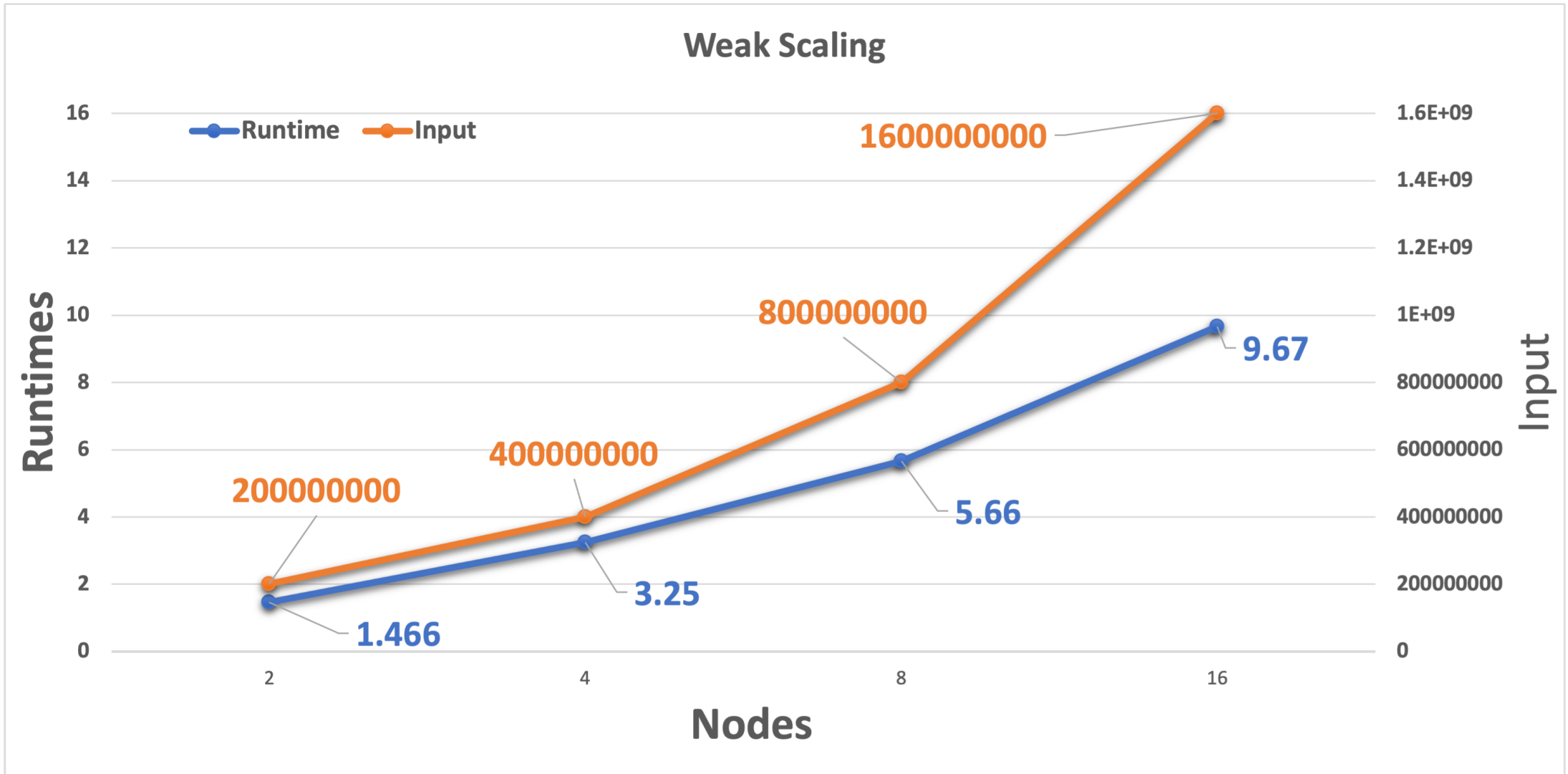
The Law focuses on strong scaling where input remains constant, and we increase the processors expecting the runtime to reduce in proportion to number of processors added maintaining reasonable efficiency.

$$S_p = \frac{1}{\beta + \frac{(1-\beta)}{p}} \quad \left\{ \begin{array}{l} \text{where } \beta \text{ is the serial part of the code which cannot be parallelized.} \end{array} \right.$$

- In the above experiment: ($S_p = 16, p = 140$) Therefore, β comes to = 5.5% .
- So according the Amdahl's Law 5.5% of code will never be parallelized.
- For Input size of 1 Billion keys we are able to strongly scale upto 32 nodes. After that the efficiency decrease dramatically.



Gustafson's law



Gustafson's law

As we increase the processor, we are able to solve bigger and bigger problems thus, achieving weak scaling

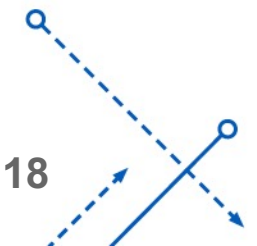
$$S_p = p - \alpha (p - 1)$$

$$\alpha = \frac{T_{seq}}{T_{seq} + T_{par}}$$

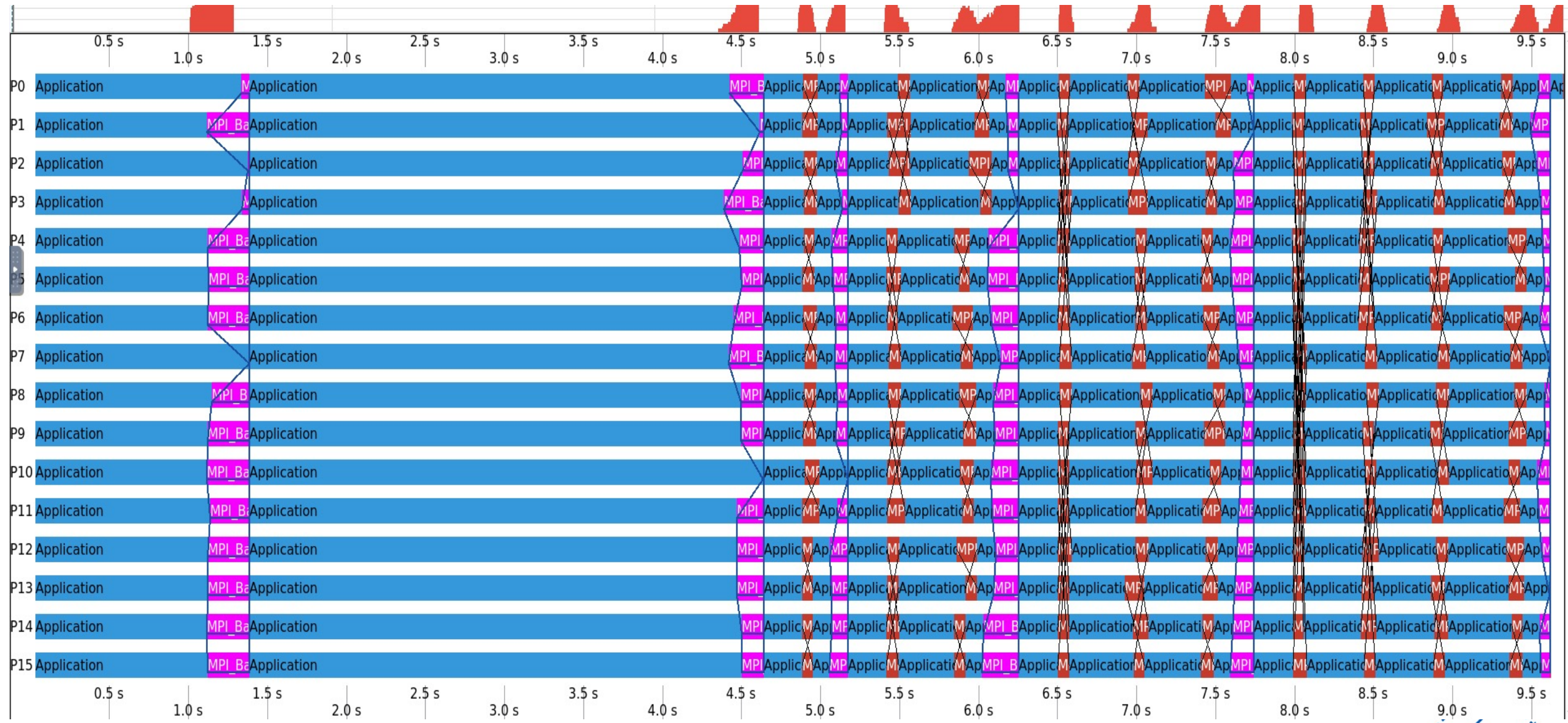
Substituting the values from above experiment : ($T_{seq} = 1.4$, $T_{par} = 9.6$, $\alpha = 0.127$)

$$S_p = 14.05$$

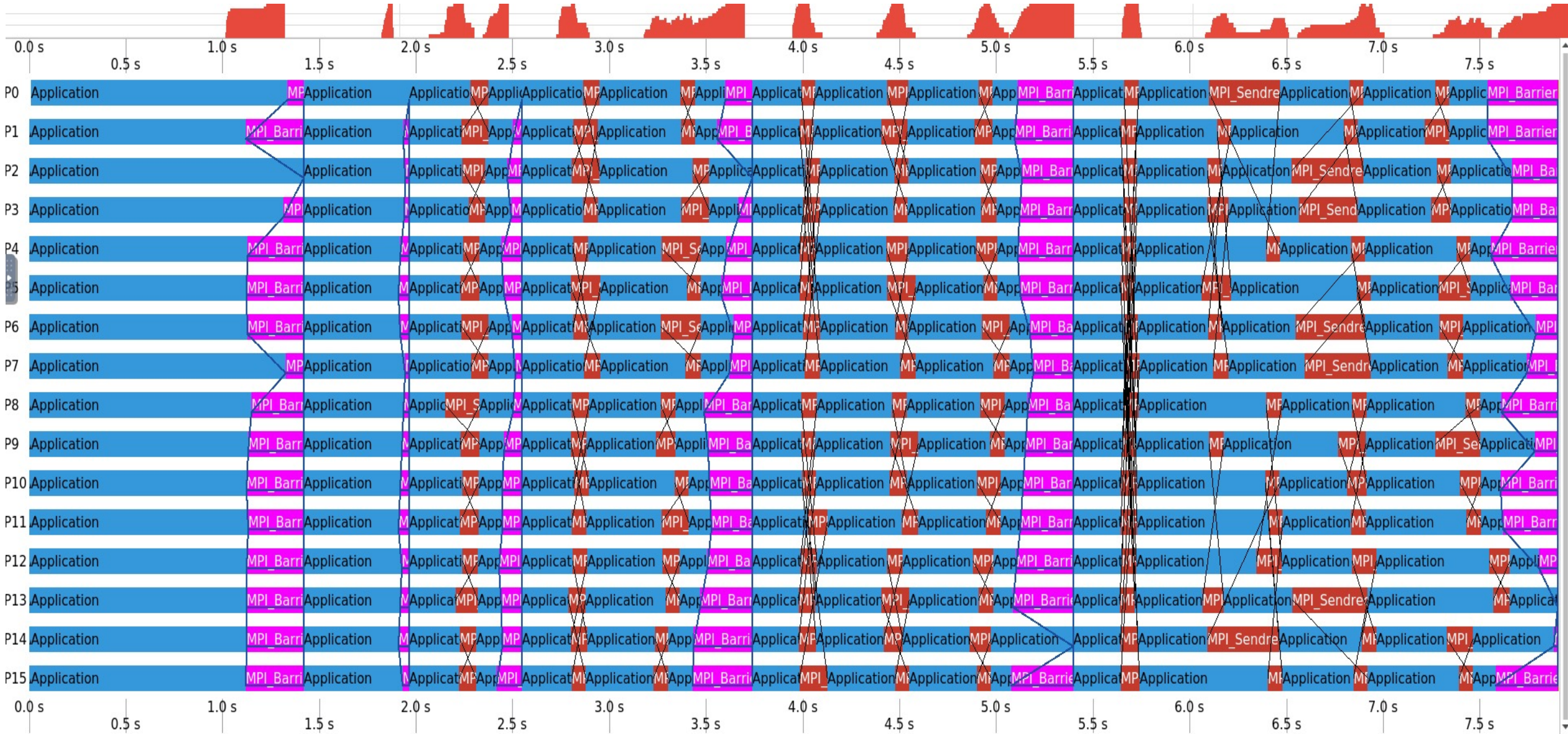
Therefore, the algorithm is weakly scalable at 16 nodes as N:P are in ratio.



Running without TBB (Inter MPI Tracer)



Running with TBB (Inter MPI Tracer)



Slurm.sh

```
1  #!/bin/bash
2
3
4  # SBATCH --mem=64000
5
6  # SBATCH --exclusive
7  # SBATCH --constraint=IB
8
9  # SBATCH --job-name="702"
10
11 # SBATCH --partition=general-compute
12 # SBATCH --qos=general-compute
13 # SBATCH --account=zsayed
14
15 # SBATCH --output=~/.panasas/logs/bitonic_%x_%j.stdout
16 # SBATCH --error=~/.panasas/logs/bitonic_%x_%j.stderr
17
18 # SBATCH --time=00:10:00
19
20 # SBATCH --nodes=16
21 # SBATCH --ntasks-per-node=1
22
23 module load intel-oneapi-2021.3
24 module load intel-oneapi-mpi/2021.3.0
25 export I_MPI_PMI_LIBRARY=/usr/lib64/libpmi.so
26
27 module load intel-tbb/2019.3
28 source /util/academic/intel/19.3/compilers_and_libraries/linux/tbb/bin/tbbvars
29
30
31 srun --mpi=pmi2 /user/zsayed/projects/bitonic-sort/bin/bitonic 1000000000
32
33 #mpirun -trace /user/zsayed/projects/bitonic-sort/bin/bitonic 1000000000
34
```

Readme

```
1  ### Bitonic sort
2  ```bash
3
4  ./build.sh
5
6
7  mpirun -np 8 ./bin/bitonic 1024
8
9  mpirun -np 8 <tau_exec> ./bin/bitonic 1024
10
11 mpirun -trace ./bin/bitonic 1024
12
13 sbatch slurm.sh
14
15 ```
16
17 - install tau pdt to instrument and profile
18
19 ```bash
20
21 export TAU_TRACE=1; export TAU_PROFILE=1; export TAU_COMM_MATRIX=1;
22 export TRACEDIR=./tau_trace; export PROFILEDIR=./tau_trace;
23
24 tau_treemerge.pl;
25
26 tau2slog2 tau.trc tau.edf -o tau.slog2;
27
28 jumpshot tau.slog2
29
30 ```
```

References

- [Introduction to Parallel Computing Solutions Manual on the Web](#)
Grama, Gupta, Karypis & Kumar
- R. Miller and L. Boxer, [*Algorithms Sequential and Parallel: A Unified Approach*](#), Third Edition, Cengage Learning, Boston, Mass., 2013.

Thank you!