

# Platform Efficiency for large-scale workloads

Martin Hilgeman  
HPC Performance Lead  
Distinguished Member of Technical Staff

[martin.hilgeman@dell.com](mailto:martin.hilgeman@dell.com)

@MartinHilgeman

**DELL** Technologies

# Dell Technologies HPC market leadership

11

generations of technologies  
in HPC clusters



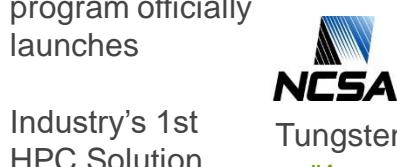
Cornell University  
Center for Advanced Computing

1st HPC cluster

1999

Industry's 1st  
HPC Solution  
Bundle

HPC solutions  
program officially  
launches



2004

Tungsten

#4



Thunderbird

#6

2005

DCS formed

2008

Dell EMC merger  
Isilon joins HPC  
portfolio

2015

2012

C-series joins  
PowerEdge



Pitzer



Pitzer

First systems with  
DCLC (OSC) and HDR

Dell EMC AI solutions announced

2017

2016

Zenith System  
launched at HPC &  
AI Innovation Lab

2021

2020

2019

2018



# Dell Technologies HPC Strategy



DEMOCRATIZING



OPTIMIZING



ADVANCING

# Dell Technologies HPC Strategy



## DEMOCRATIZING

Enable everyone to  
accelerate science,  
engineering and analytics



## OPTIMIZING



## ADVANCING

# Dell Technologies HPC Strategy



DEMOCRATIZING



OPTIMIZING

Offering a robust portfolio,  
optimized for HPC, providing  
performance and efficiency



ADVANCING

# Dell Technologies HPC strategy



DEMOCRATIZING



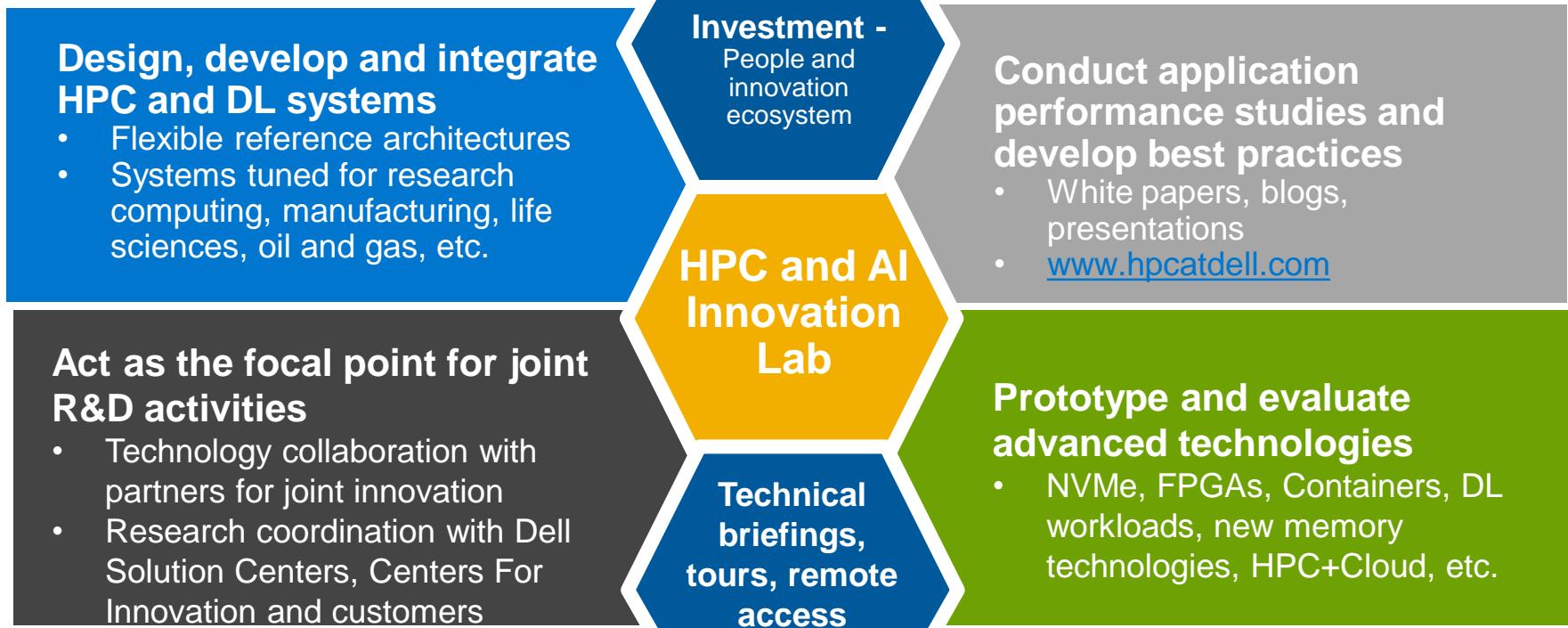
OPTIMIZING



ADVANCING

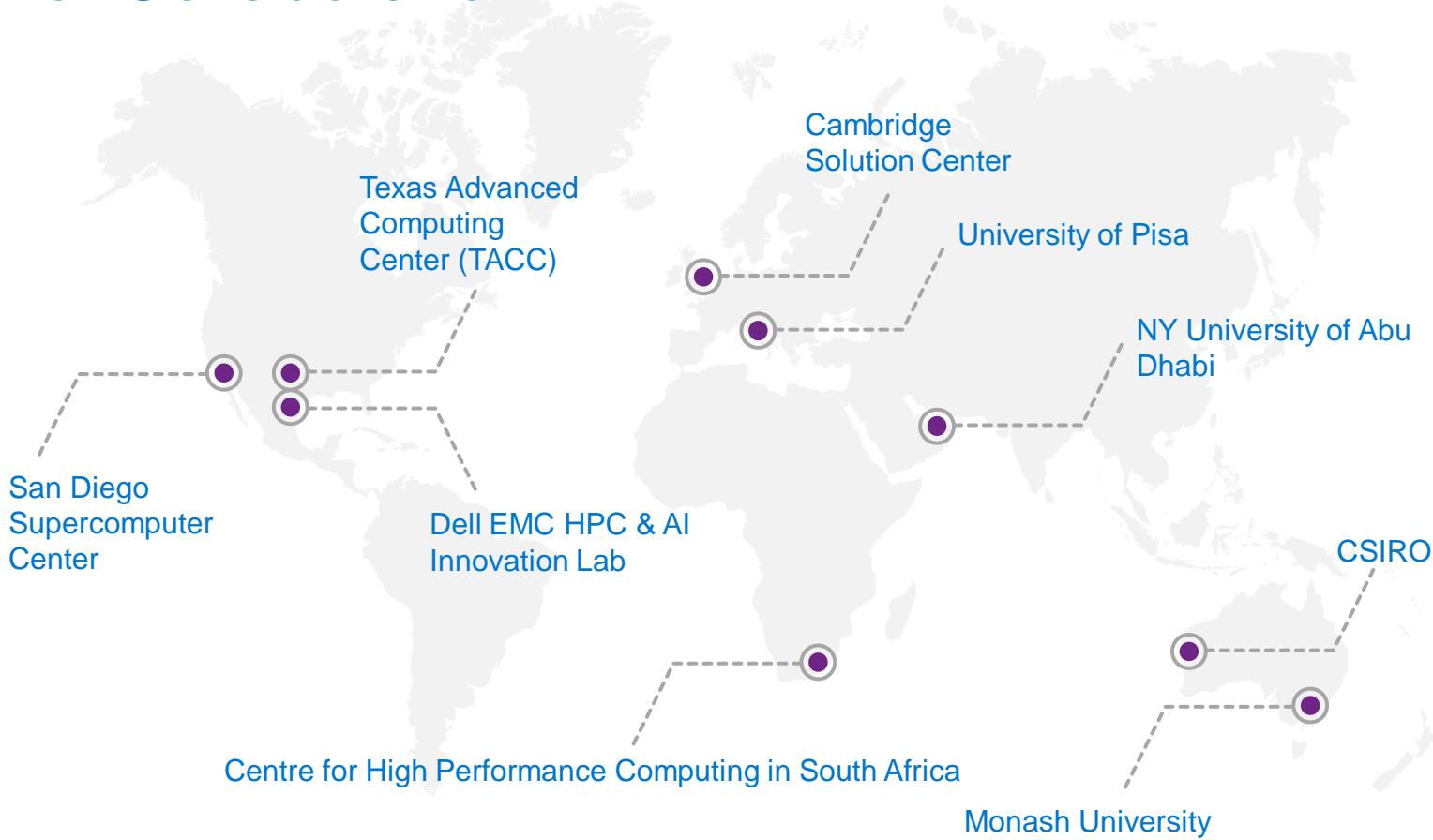
Providing expertise,  
innovation and partnerships  
— from workgroup to the  
TOP500

# Dell Technologies HPC and Deep Learning Team Charter



[www.dellemc.com/innovationlab](http://www.dellemc.com/innovationlab)

# The Spirit of Collaboration



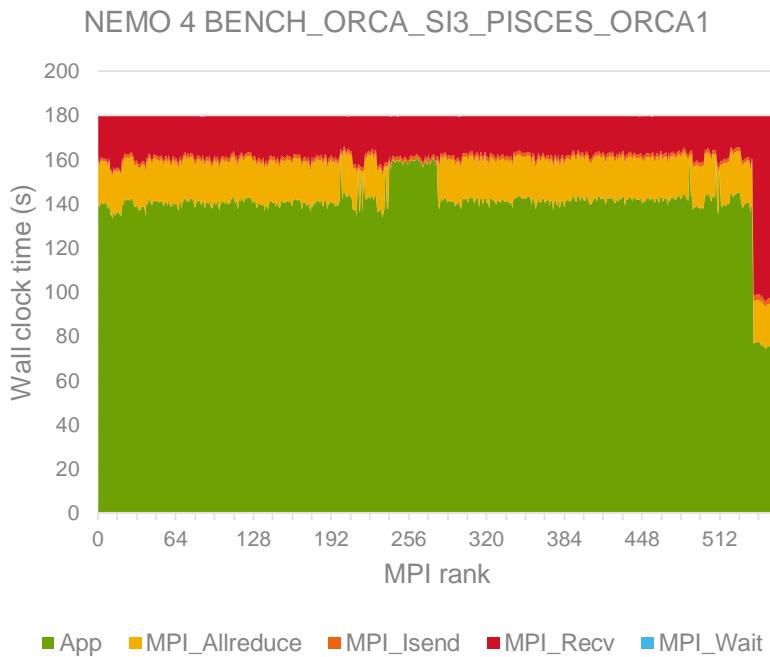


## Examples of MPI studies

# Collective MPI function rationale

- Collective functions are called by all ranks participating in the communicator used in the function
- A collective function implies a barrier
- Most MPI libraries have implemented various kernels that are tuned for a certain topology or message size:
  - Intel MPI: `I_MPI_ADJUST_ALLREDUCE=<val>` (1-Recursive doubling, 2-Rabenseifner's, 3-Reduce + Bcast, 4-Topology aware Reduce + Bcast, 5-Binomial gather + scatter, 6-Topology aware binomial gather + scatter, 7-Shumilin's ring, 8-Ring, 9-Knomial, 10-Topology aware SHM-based flat, 11-Topology aware SHM-based Knomial, 12-Topology aware SHM-based Knary)
  - Open MPI: `mpirun --mca coll_tuned_allreduce_algorithm <val>` (0-ignore, 1-basic\_linear, 2-nonoverlapping, 3-recursive\_doubling, 4-ring, 5-segmented\_ring, 6-rabenseifner)
  - MVAPICH2: algorithms are used based on thresholds set by environment variables like `MV2_SHMEM_ALLREDUCE_MSG` and `MV2_ALLREDUCE_2LEVEL_MSG`

# Collective MPI function rationale



- Collective functions spend a lot of time “fixing” load imbalance
- MPI ranks spend different amounts of time in MPI\_Allreduce

# MPI-3 feature: non-blocking collectives

- MPI-3 has added support for non-blocking collectives
  - They combine the potential benefits of nonblocking point-to-point operations, to exploit overlap and to avoid synchronization, with the optimized implementation and message scheduling provided by collective operations
  - The call initiates the operation, which indicates that the system may start to copy data out of the send buffer and into the receive buffer
  - Workings are similar to non-blocking point to point calls (Isend/Irecv)
- **Blocking:** `int MPI_Allreduce( const void *sendbuf, void *recvbuf,  
int count, MPI_Datatype datatype, MPI_Op op,  
MPI_Comm comm)`
- **Non-blocking:** `int MPI_Iallreduce(const void *sendbuf, void *recvbuf,  
int count, MPI_Datatype datatype, MPI_Op op,  
MPI_Comm comm, MPI_Request *request)`

# The simple collective example – non blocking

```
#include <stdio.h>
#include <stdlib.h>
#include <mpi.h>

int main(int argc, char* argv[])
{
    int size, my_rank, result = 0;
    MPI_Request request;

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &size);
    MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);

    MPI_Iallreduce(&my_rank, &result, 1, MPI_INT,
MPI_SUM, MPI_COMM_WORLD, &request);

    printf("[MPI Process %d] Doing something else
here.\n", my_rank);

    MPI_Wait(&request, MPI_STATUS_IGNORE);

    printf("[MPI Process %d] The sum of all ranks is
%d.\n", my_rank, result);

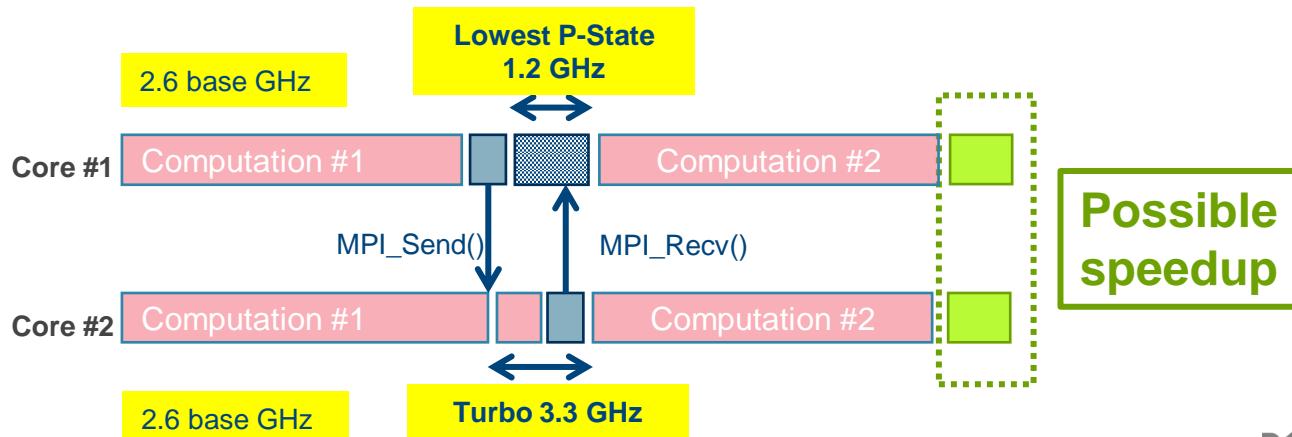
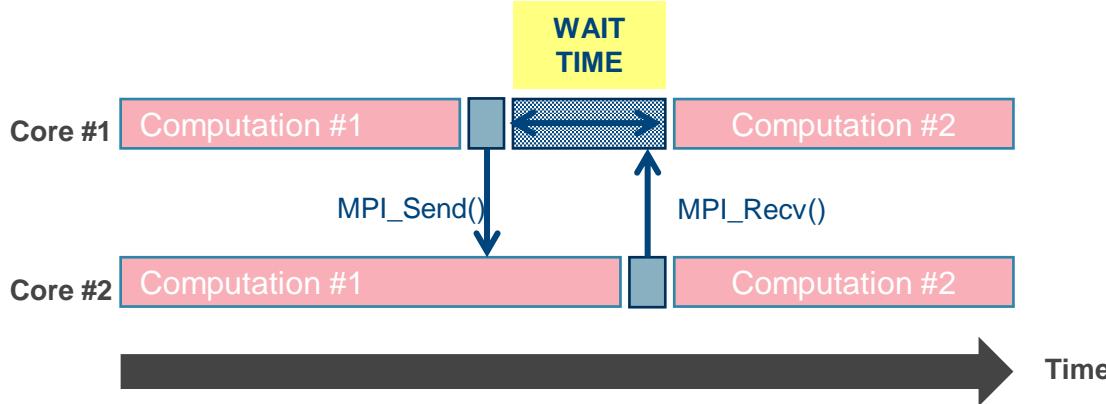
    MPI_Finalize();

    return EXIT_SUCCESS;
}
```

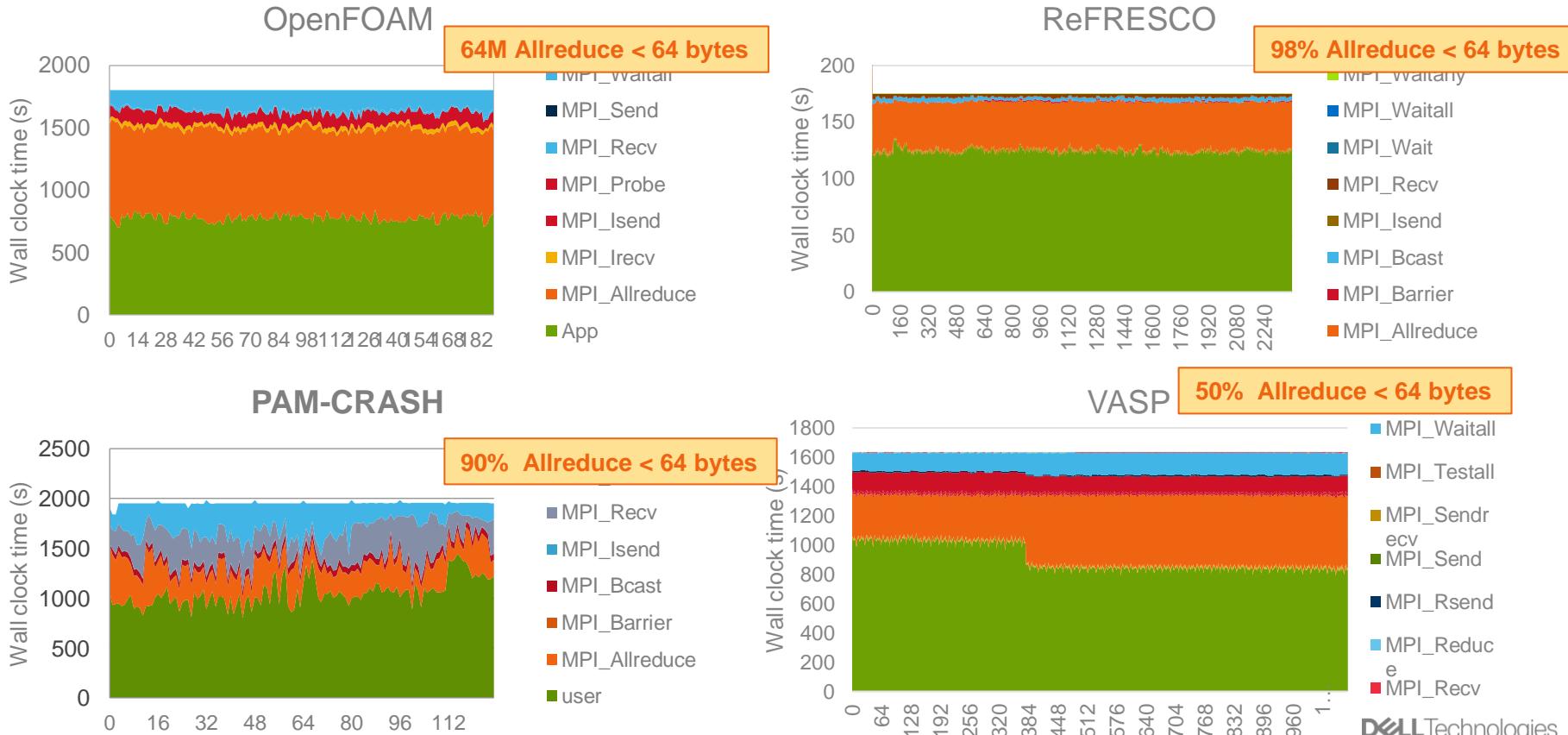
- Similar calling sequence as MPI\_Allreduce
- MPI\_Request parameter is extra, just like in Isend/Irecv

```
% mpirun -np 4 iallreduce.exe
[MPI Process 3] Doing something else here.
[MPI Process 0] Doing something else here.
[MPI Process 1] Doing something else here.
[MPI Process 1] The sum of all ranks is 6.
[MPI Process 2] Doing something else here.
[MPI Process 2] The sum of all ranks is 6.
[MPI Process 3] The sum of all ranks is 6.
[MPI Process 0] The sum of all ranks is 6.
```

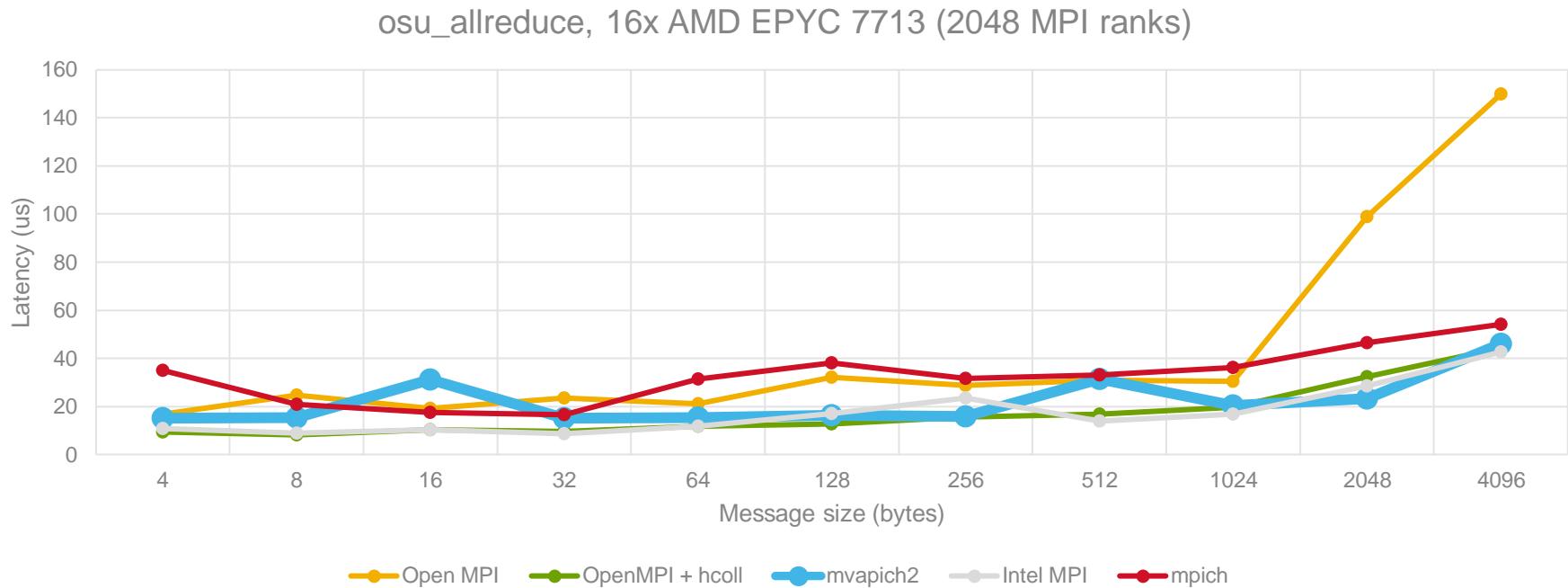
# The effect of turbo mode could help load imbalance



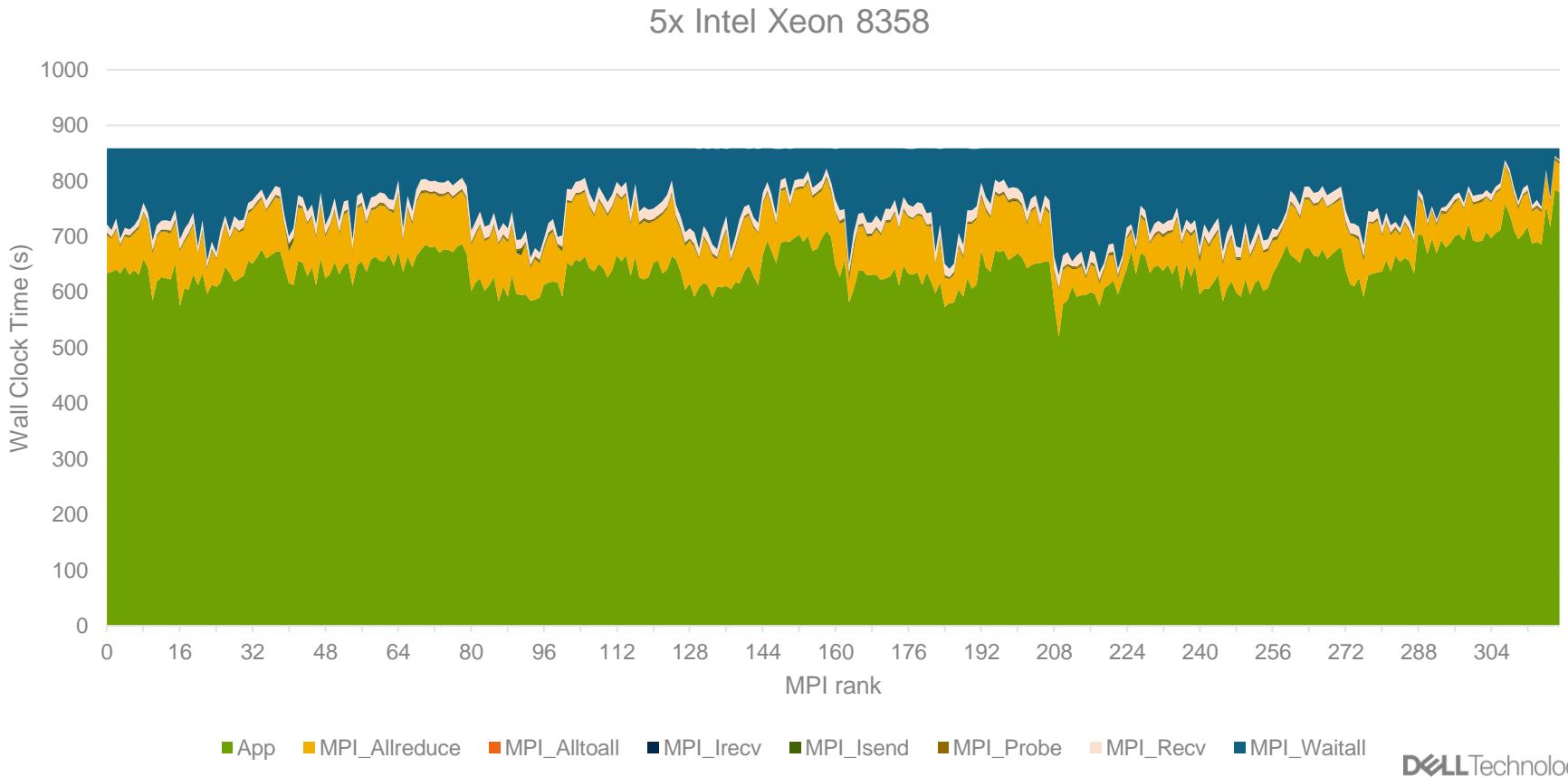
# Small MPI\_Allreduce is killing MPI performance



# Allreduce for small message sizes



# OpenFOAM v2106 simpleBenchmark2006



# OpenFOAM v2106 simpleBenchmark2006

5x Intel Xeon 8358



# Looking for real communication overhead

- We have seen a load imbalance in the OpenFOAM benchmark, what about real communication overhead?
- Profiling library simulates collective operations as follows:
  - **PMPI\_Barrier(comm) is doing an explicit synchronization to rule out load imbalance**
  - **The “real” MPI\_Allreduce call is now on perfectly synchronized data**

```
int MPI_Allreduce(const void *sendbuf, void *recvbuf, int count,
MPI_Datatype datatype, MPI_Op op, MPI_Comm comm)
{
    int ret, len = 0, msgsize = 0;
    double t0, t1, t2;

    t0 = my_time();
    if (datatype != MPI_DATATYPE_NULL ) {
        PMPI_Type_size(datatype, &msgsize);
        len = count * msgsize;
    }

    if (libtoolbox.coll_stats) {
        PMPI_Barrier(comm);
        t2 = my_time() - t0;
        mpi.coll_time[ALLREDUCE] += t2;
    }

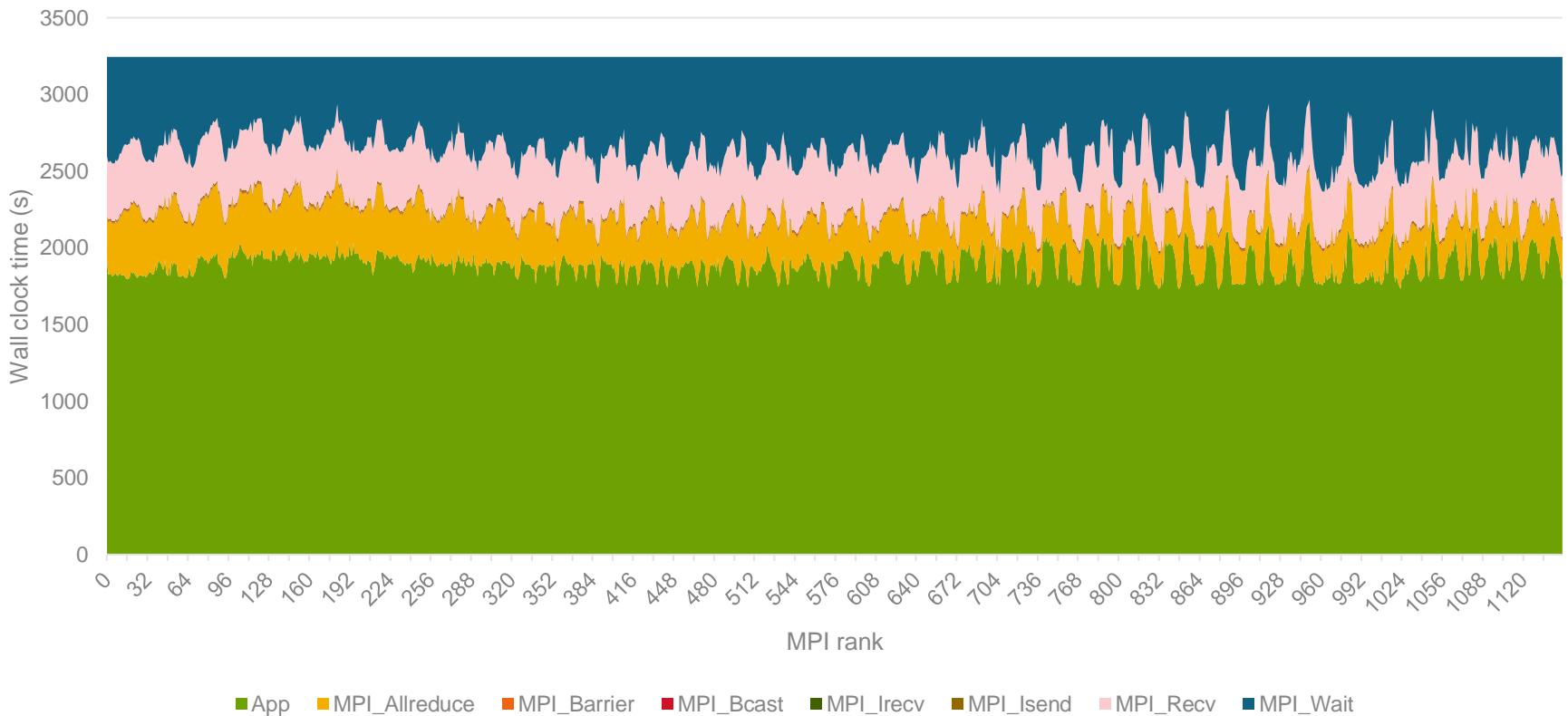
    ret = PMPI_Allreduce(sendbuf, recvbuf, count, datatype, op,
comm);

    if (libtoolbox.stats) {
        t1 = my_time() - t0;
        mpi.time[MPI_ALLREDUCE] += t1;
        mpi.len[MPI_ALLREDUCE] += len;
        update_bin(MPI_ALLREDUCE, len);
    }

    return ret;
}
```

# Allreduce example: MOM5 benchmark

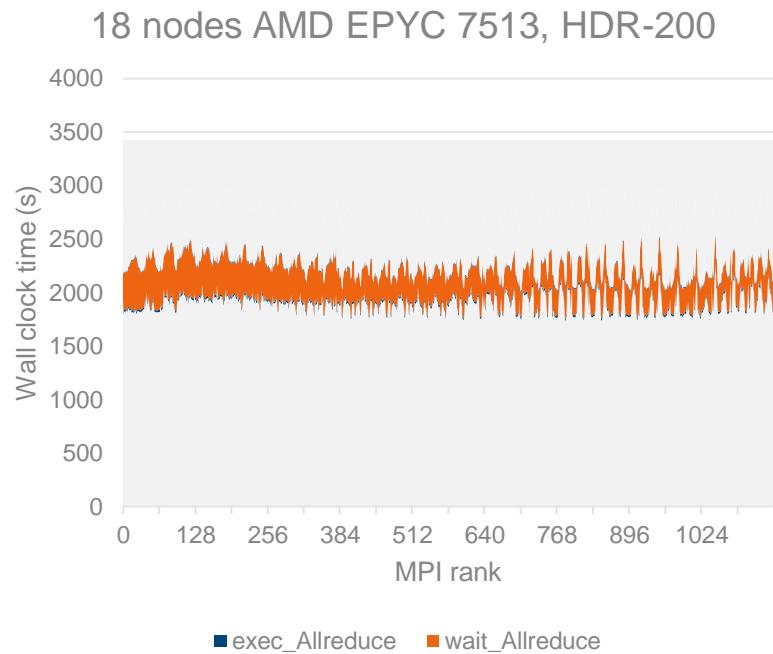
18 nodes, AMD EPYC 7513, HDR-200



# Result of perfect interconnect simulation (zoomed in)

- Wait\_\*: load imbalance
- Exec\_\*: actual communication
- Communication overhead is negligible
- Load imbalance is the cause of parallel overhead and potential scaling issue

123k calls, 119k < 64 bytes  
Minimum time: 191 seconds  
Maximum time: 463 seconds  
Transfer time: 18 seconds



# Result of perfect interconnect simulation (zoomed in)

- Wait\_\*: load imbalance
- Exec\_\*: actual communication
- Communication overhead is negligible
- Load imbalance is the cause of parallel overhead and potential scaling issue

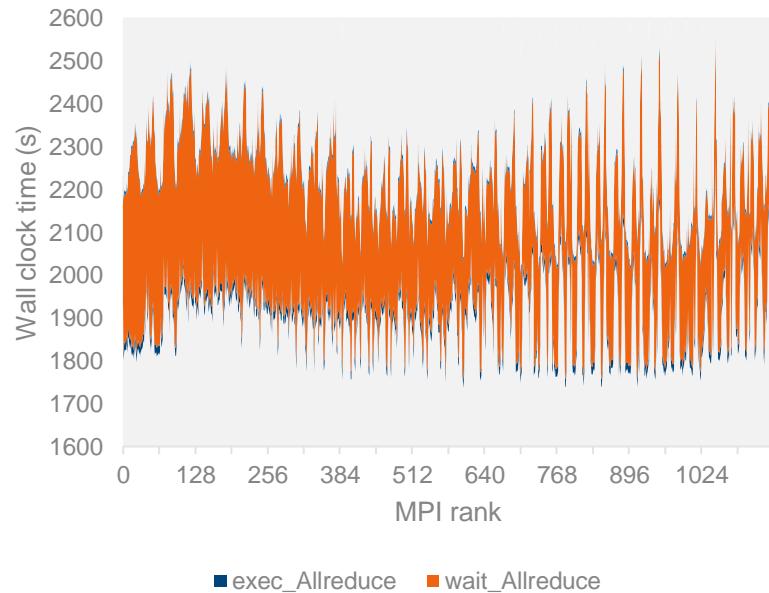
123k calls, 119k < 64 bytes

Minimum time: 191 seconds

Maximum time: 463 seconds

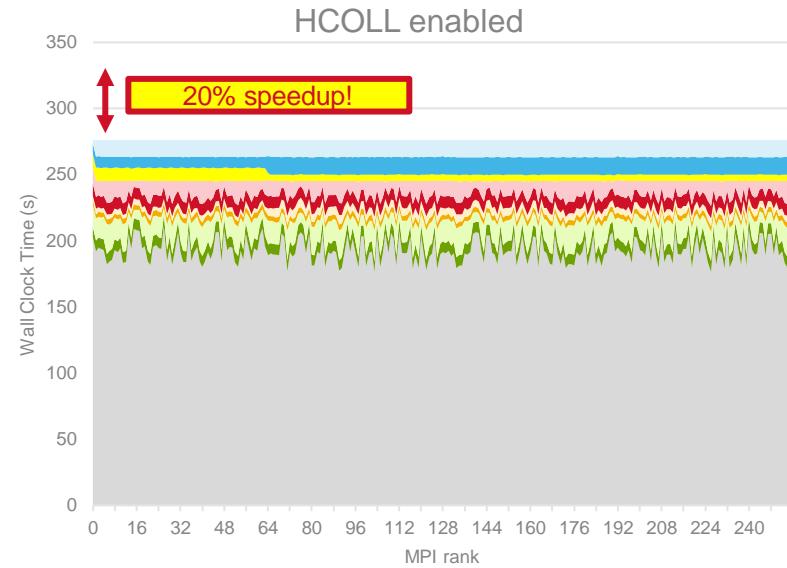
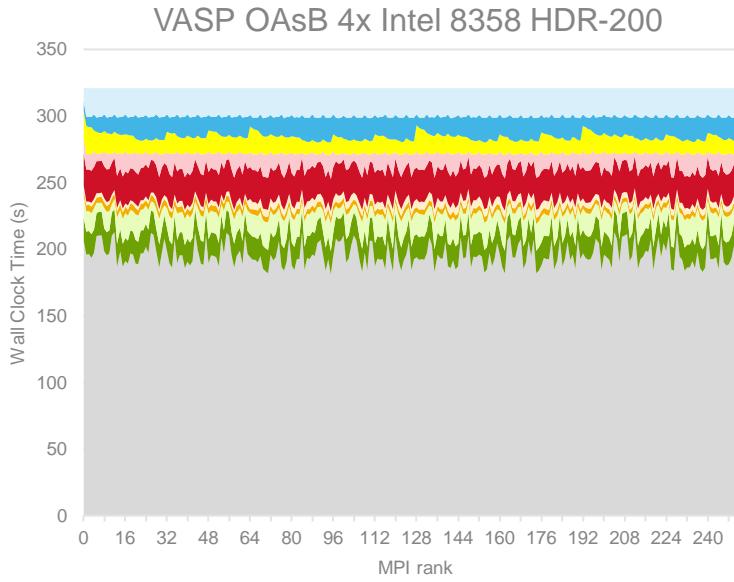
Transfer time: 18 seconds

18 nodes AMD EPYC 7513, HDR-200



# Use hardware offload of the IB card for collectives

```
% mpirun -genv HCOLL_ENABLE 1           -np 512 <app> # MVPAPICH2  
% mpirun --mca coll_hcoll_enable 1      -np 512 <app> # Open MPI  
% mpirun -genv I_MPI_COLL_EXTERNAL hcoll -np 512 <app> # Intel MPI
```

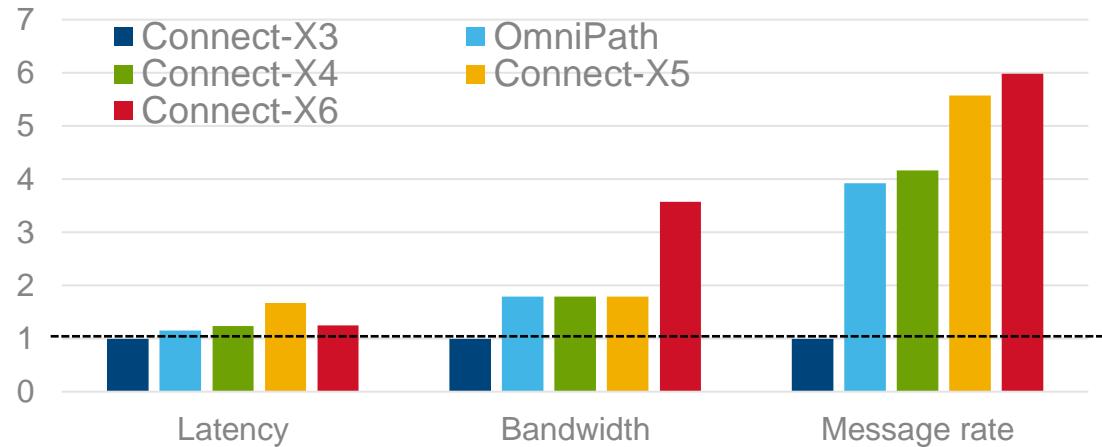
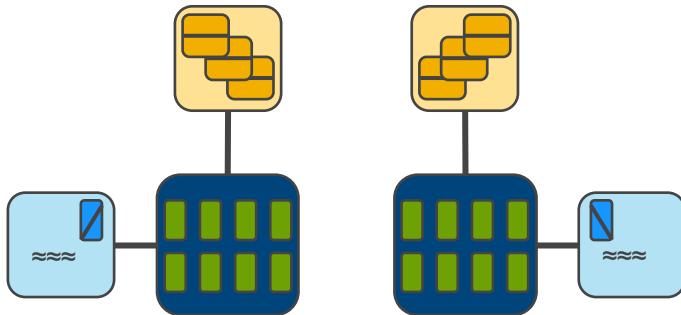


# Message rate matters!

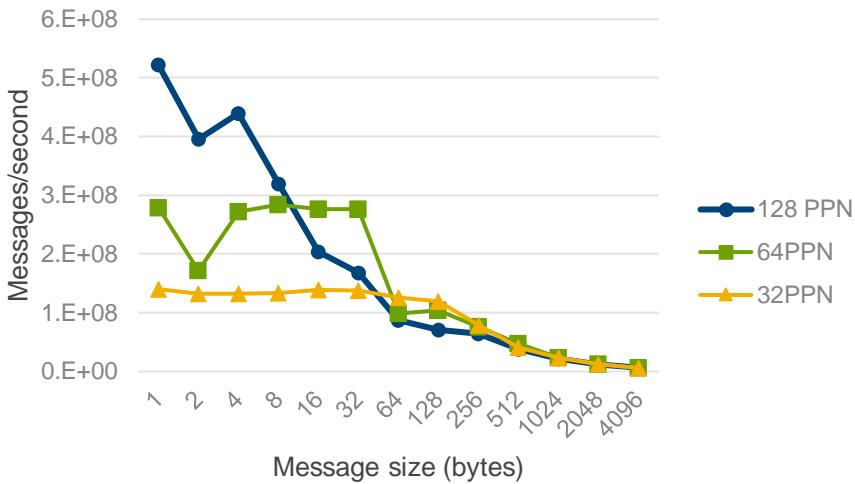
**MPI\_Allreduce()** between two nodes

**2012:** 16 cores per node, one FDR  
IB HCA:  $2 \times 16 = 32$  transfers  
(messages)

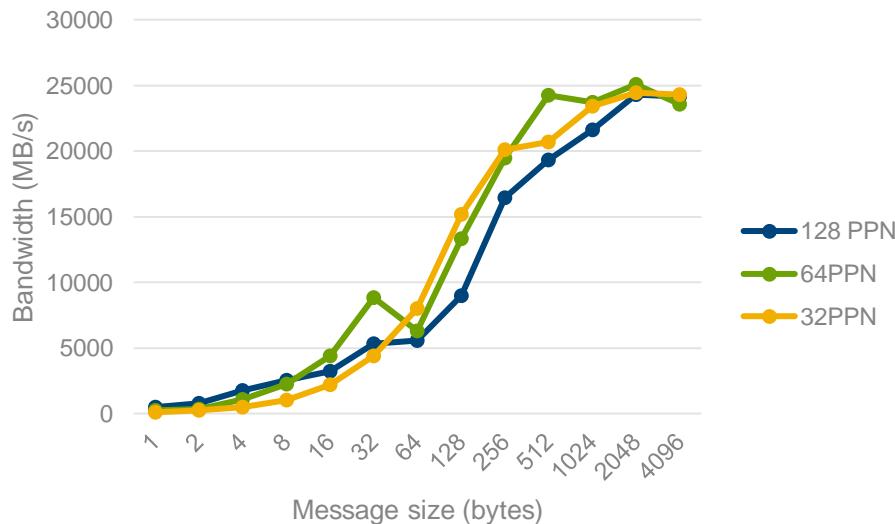
**2022:** 128 cores per node, one HDR  
IB HCA:  $2 \times 128 = 256$  transfers  
(messages)



# AMD EPYC 7713 message rate



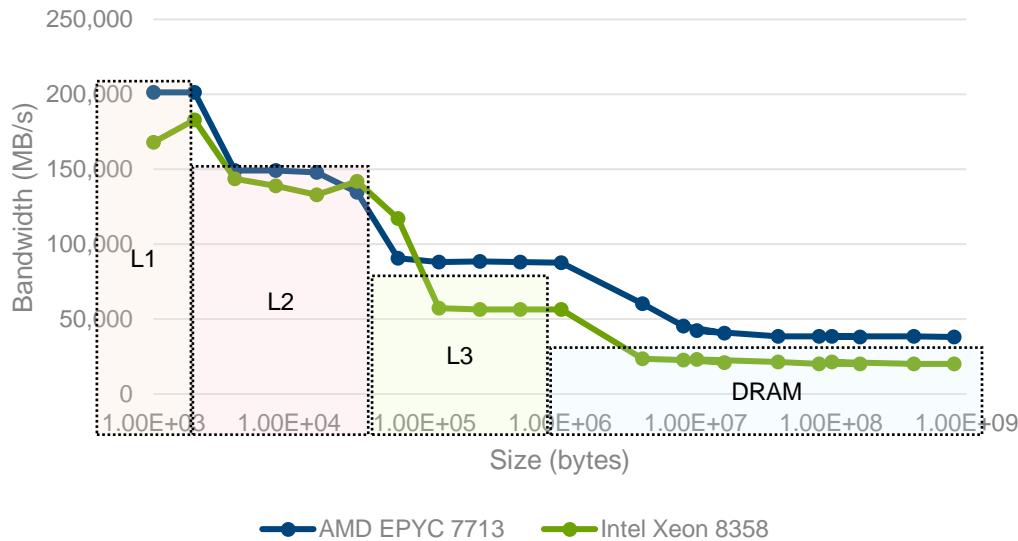
- But it doesn't add bandwidth



- Adding more cores helps message rate for small messages

# Measuring single core memory performance

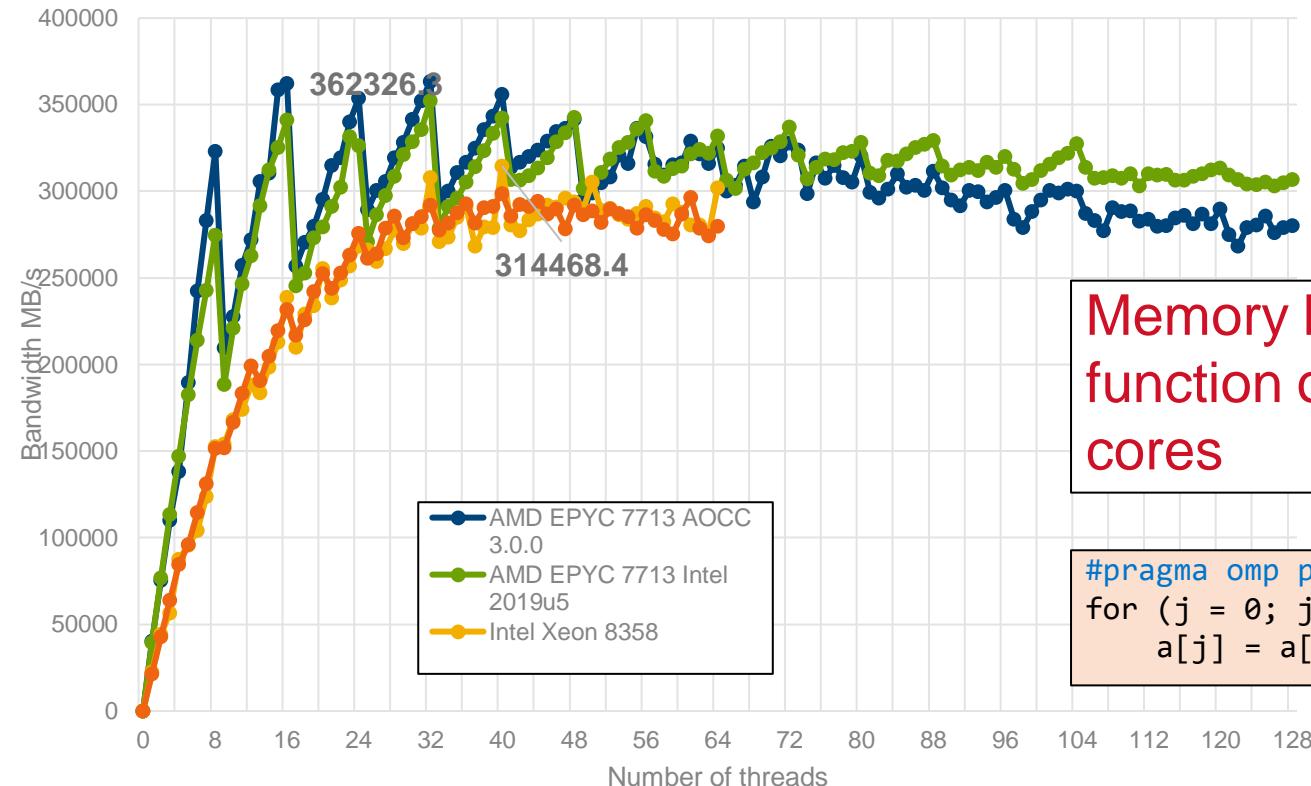
Single core DAXPY bandwidth



## DAXPY kernel

```
#pragma ivdep
for (j = 0; j < SIZE; j++)
    a[j] = a[j] + scalar * c[j];
```

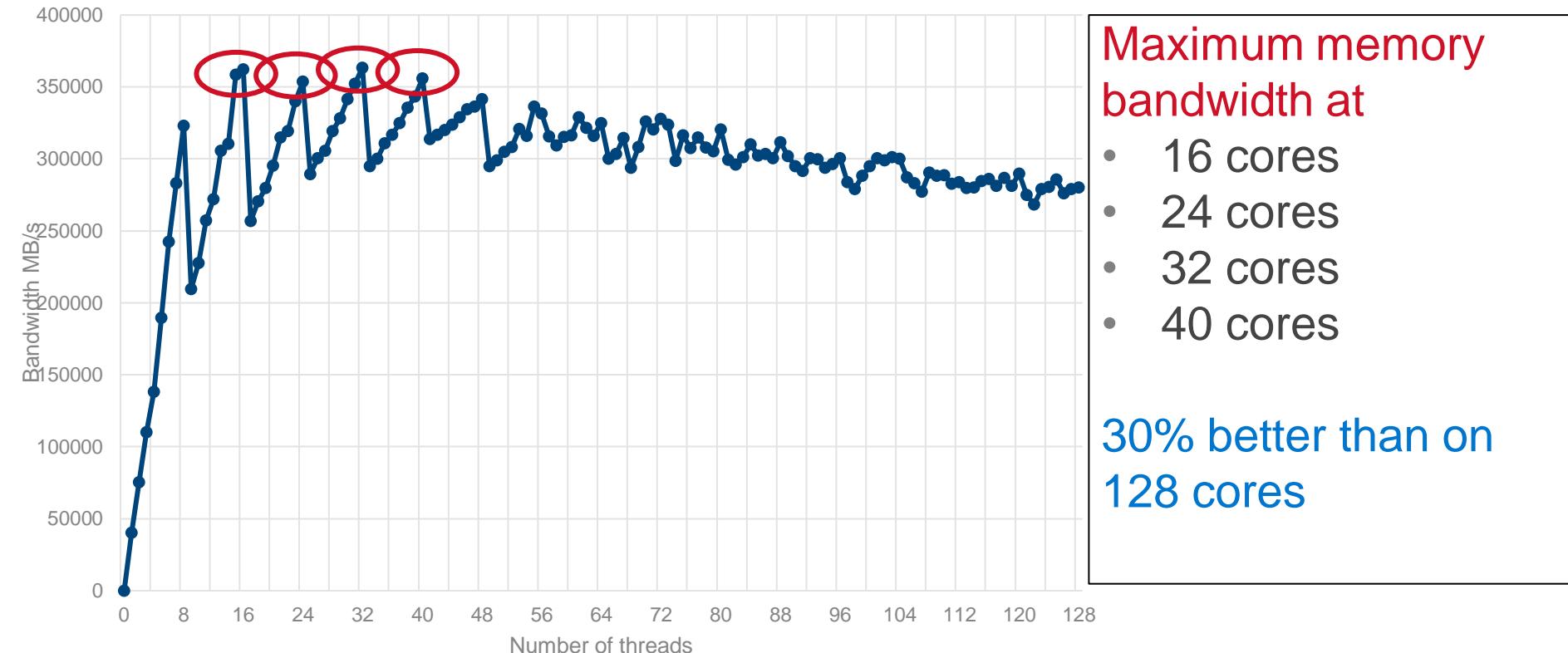
# Measuring node memory performance



Memory bandwidth as a function of the number of cores

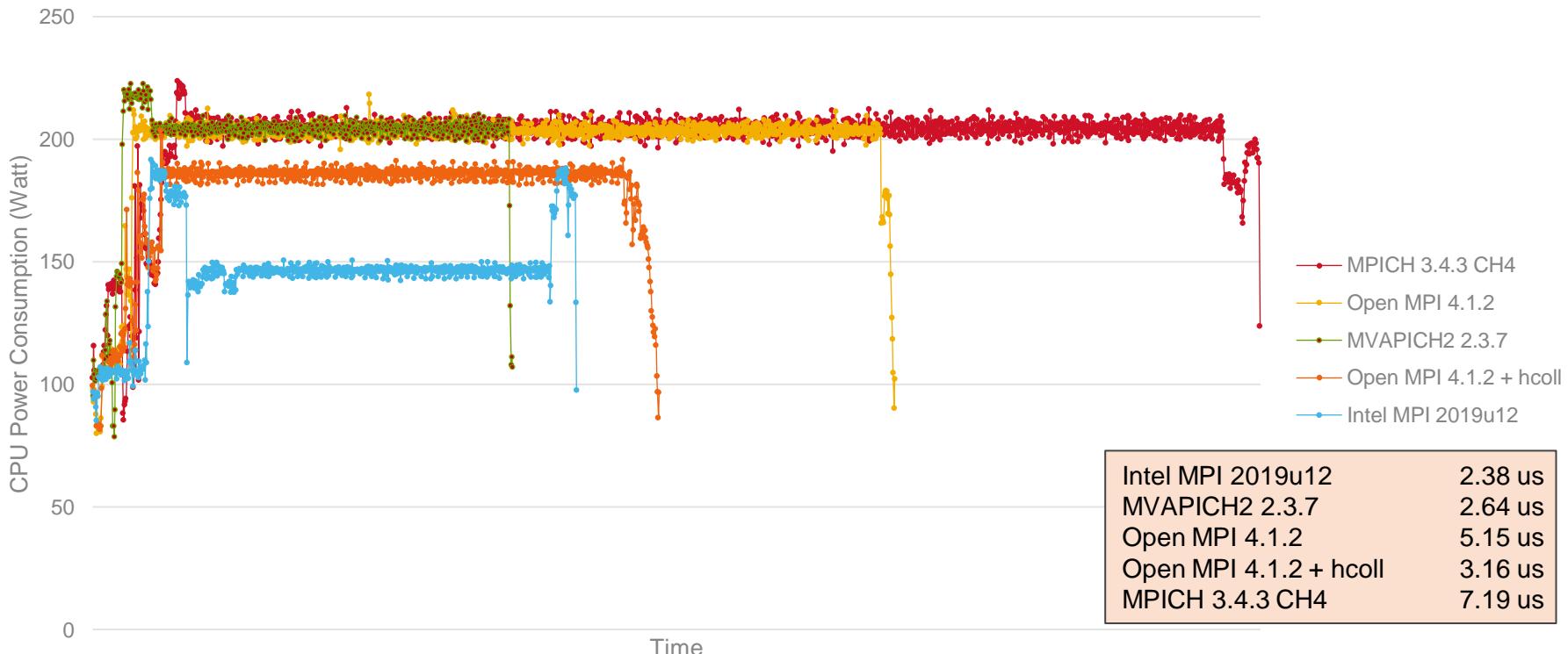
```
#pragma omp parallel for  
for (j = 0; j < STREAM_ARRAY_SIZE; j++)  
    a[j] = a[j] + scalar * c[j];
```

# Why is this important?

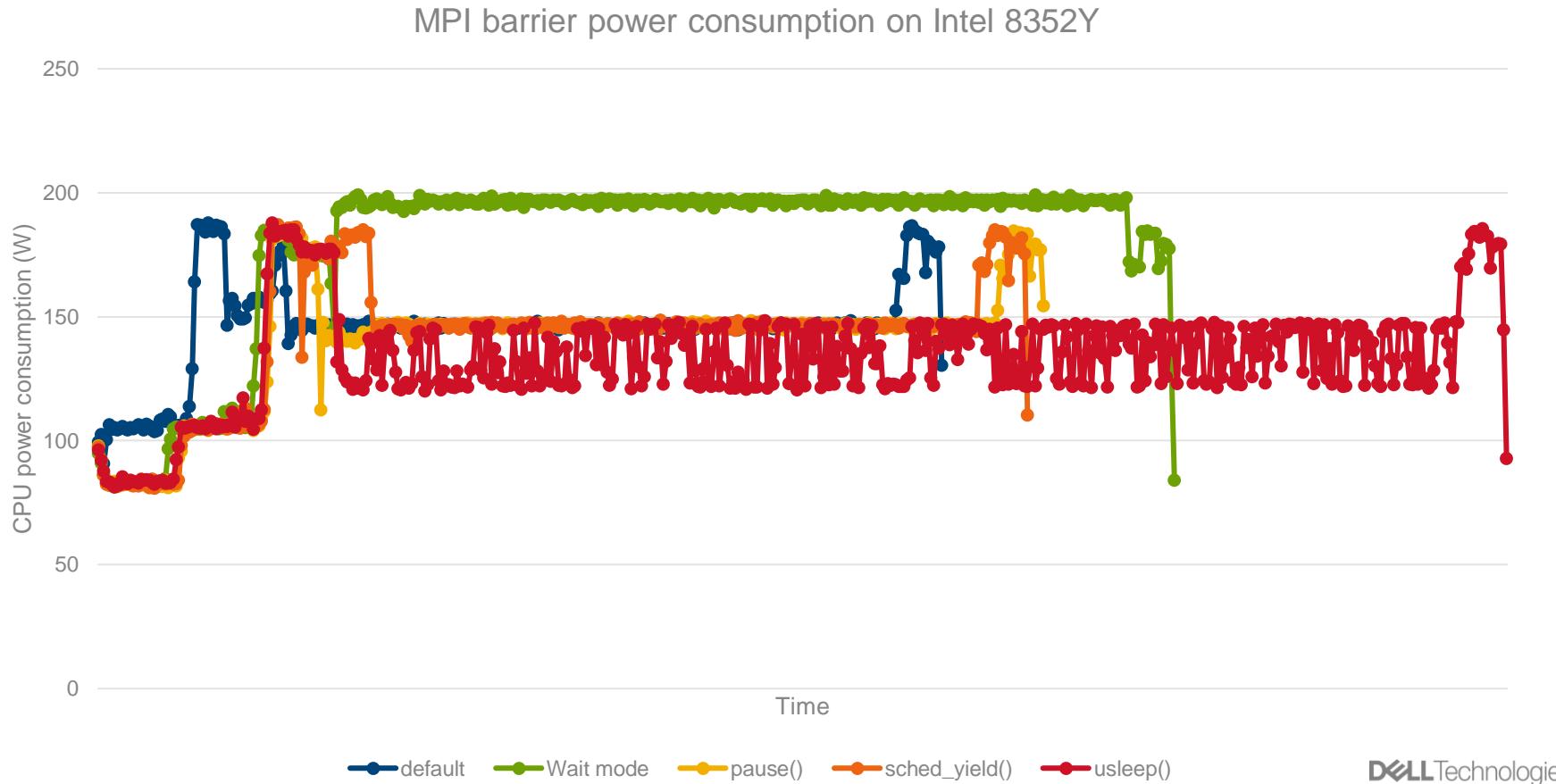


# MPI\_Barrier() power consumption

2 node Intel Xeon 8352Y 2.0 GHz 205W TDP

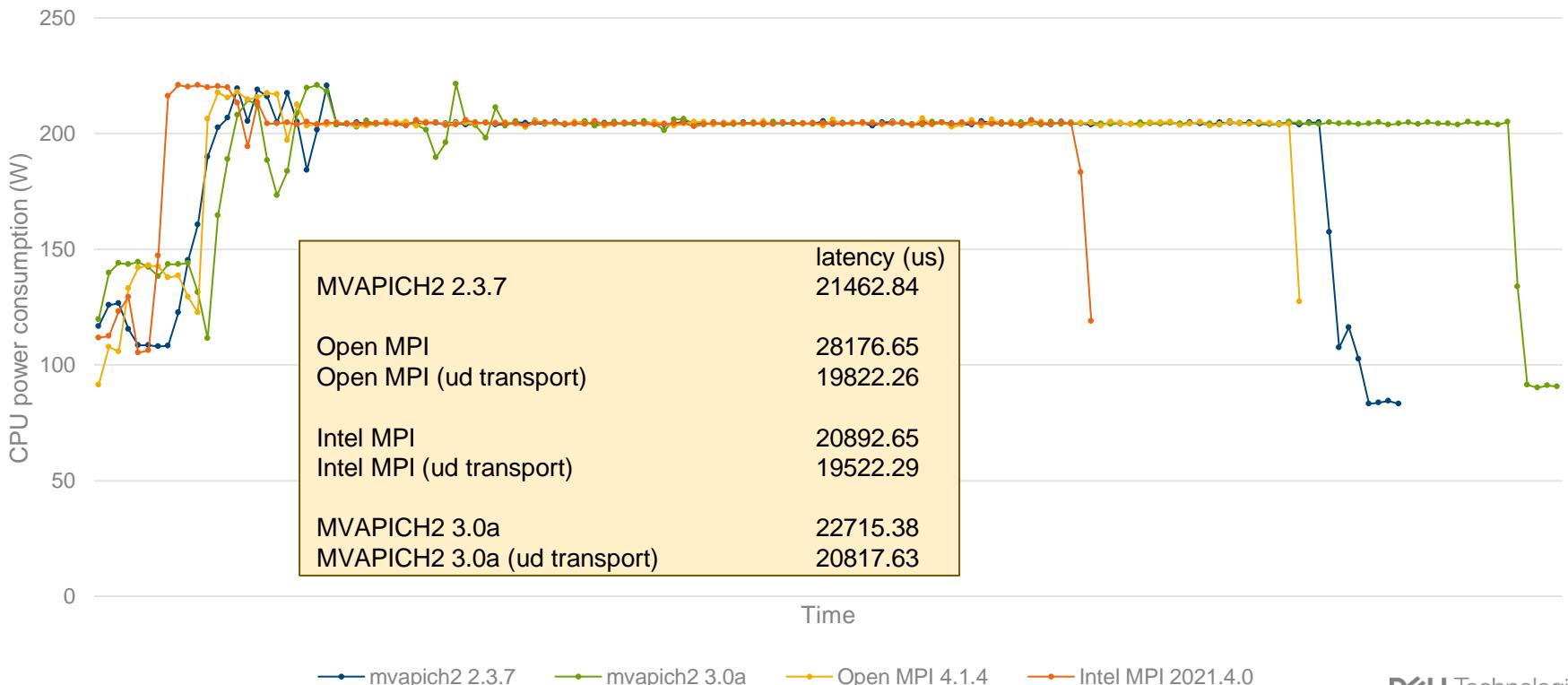


# MPI\_Barrier() power consumption



# MPI\_Alltoall() power consumption

16 node Intel Xeon 8352Y 2.0 GHz, 205 W TDP  
1024 MPI ranks, 4 KB message size





THANK YOU FOR  
BEING PART OF  
THE COMMUNITY

- Dell Technologies HPC Community: [dellhpc.org](https://dellhpc.org)
- HPC & AI Engineering: [hpcatdell.com](https://hpcatdell.com) and <https://infohub.delltechnologies.com/t/high-performance-computing/>
- HPC & AI Innovation Lab: [delltechnologies.com/innovationlab](https://delltechnologies.com/innovationlab)
- HPC/AI Centers of Excellence: [delltechnologies.com/coe](https://delltechnologies.com/coe)
- [delltechnologies.com/hpc](https://delltechnologies.com/hpc)
- Stay tuned to <https://insidehpc.com/dell>