

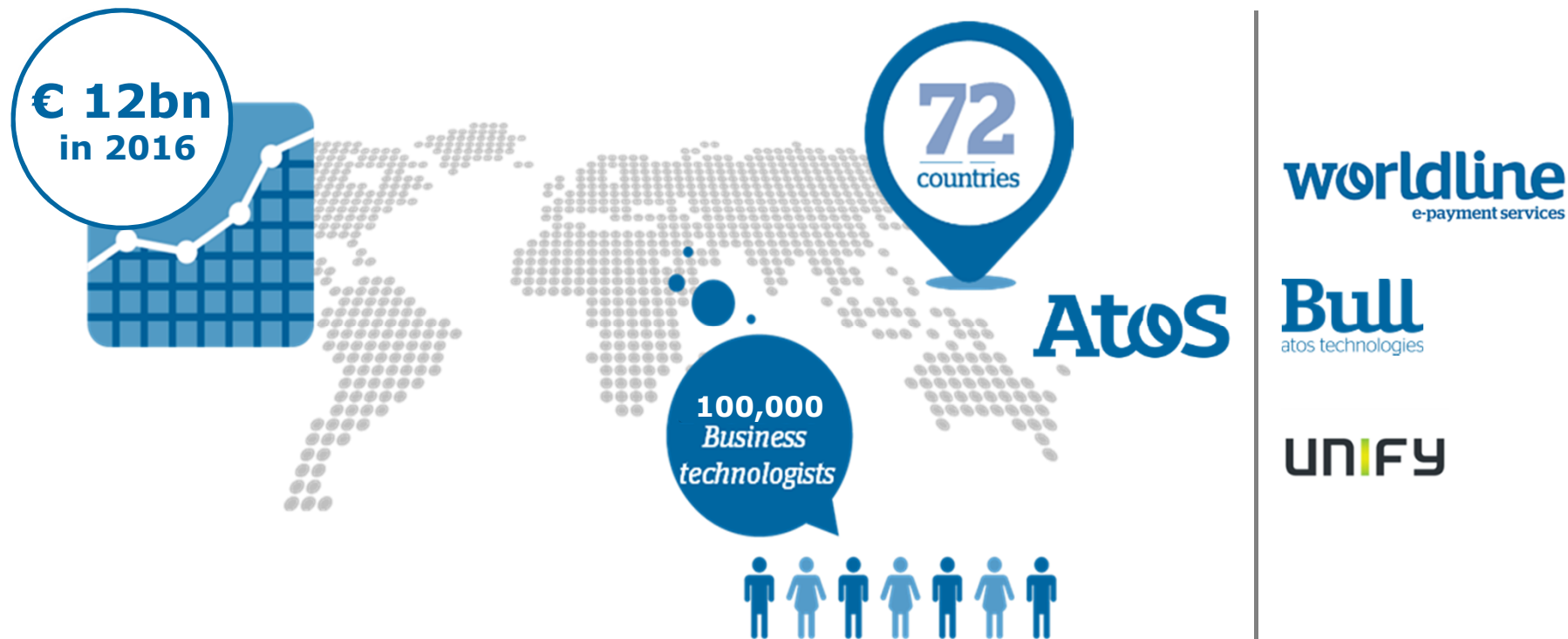
Extreme Computing

Strategic directions for the 2020s

04-07-2017

We are THE European IT Leader

and in the top 5 of worldwide Digital services players



Top500 trend

For the last 10 years

▶ 22 systems in top500 of Nov. 2016

▶ #1 in Europe among new systems deployed during the last 6 months

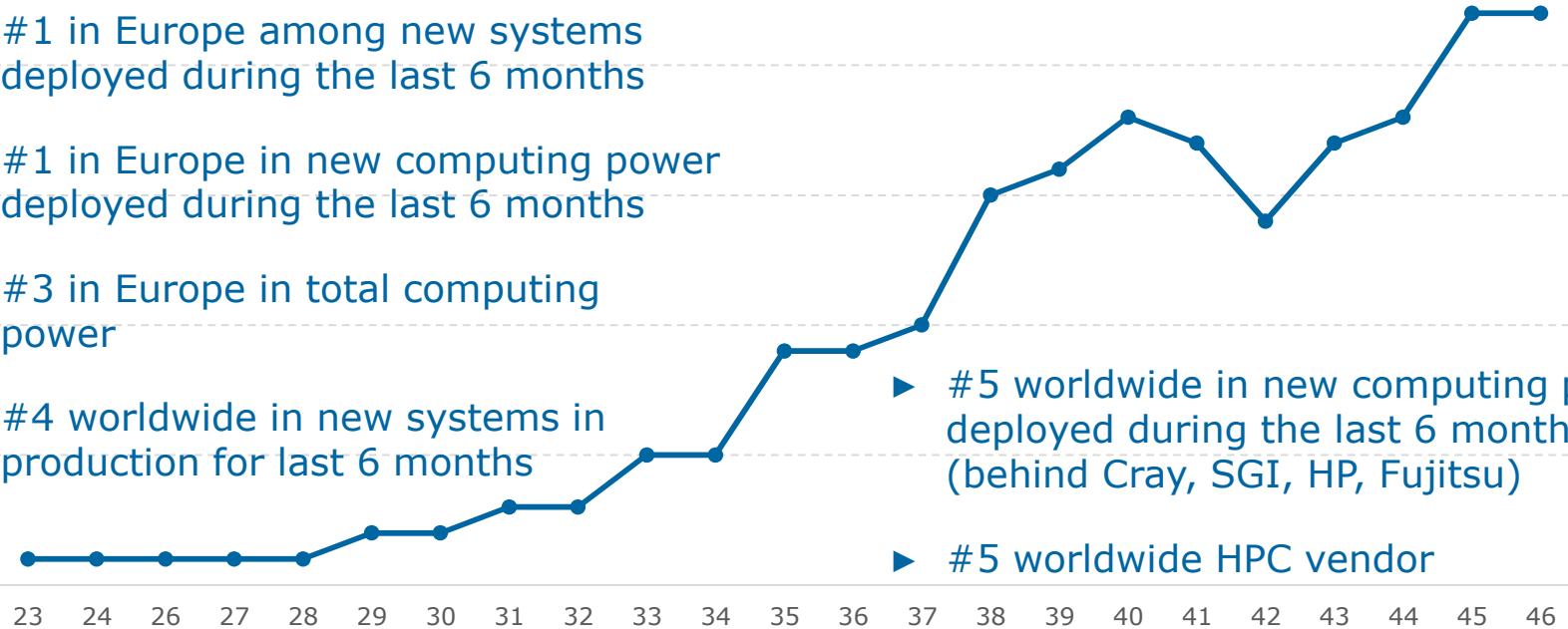
▶ #1 in Europe in new computing power deployed during the last 6 months

▶ #3 in Europe in total computing power

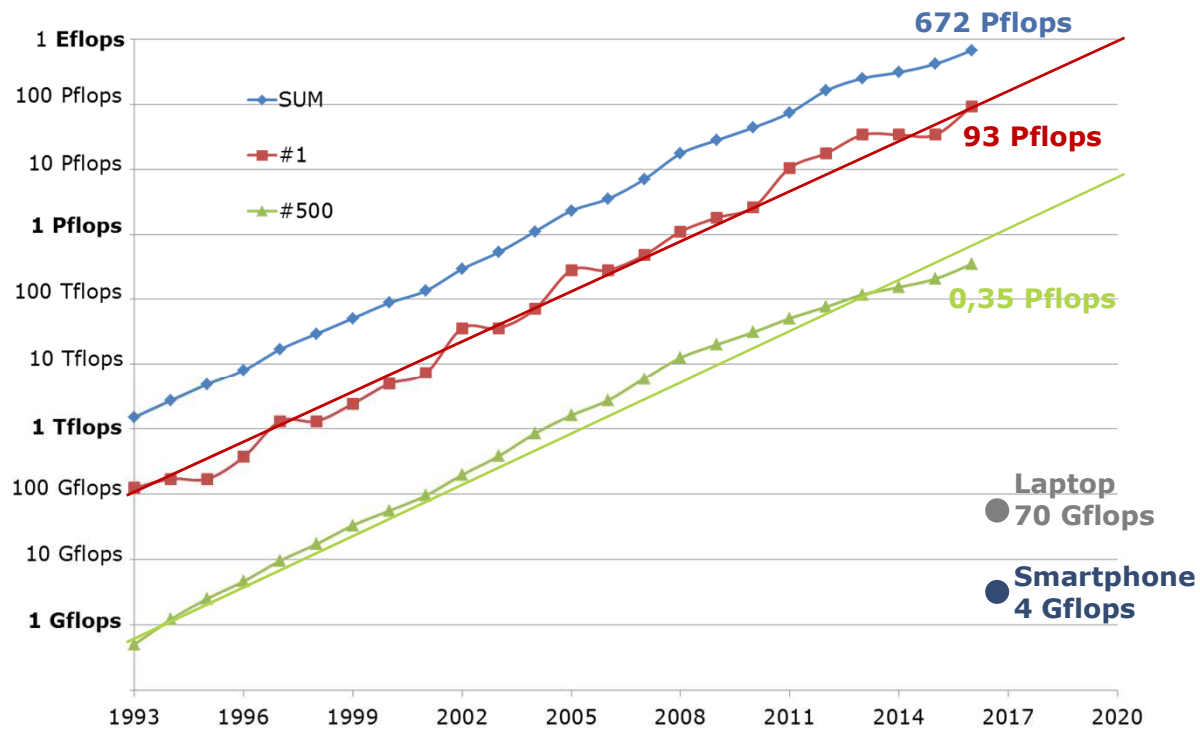
▶ #4 worldwide in new systems in production for last 6 months

▶ #5 worldwide in new computing power deployed during the last 6 months (behind Cray, SGI, HP, Fujitsu)

▶ #5 worldwide HPC vendor

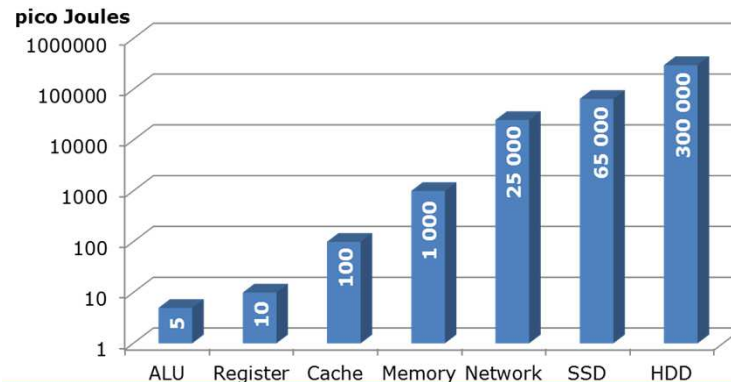
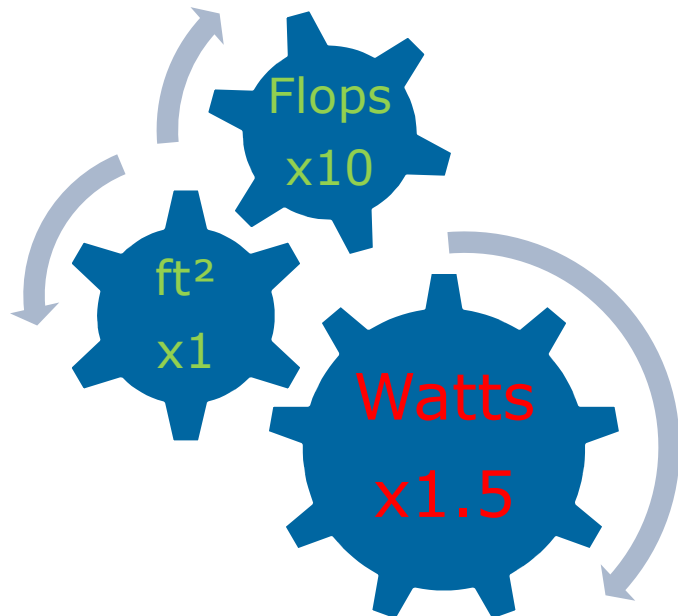


Are all lights green?



Exascale era still early 2020s

Some divergences



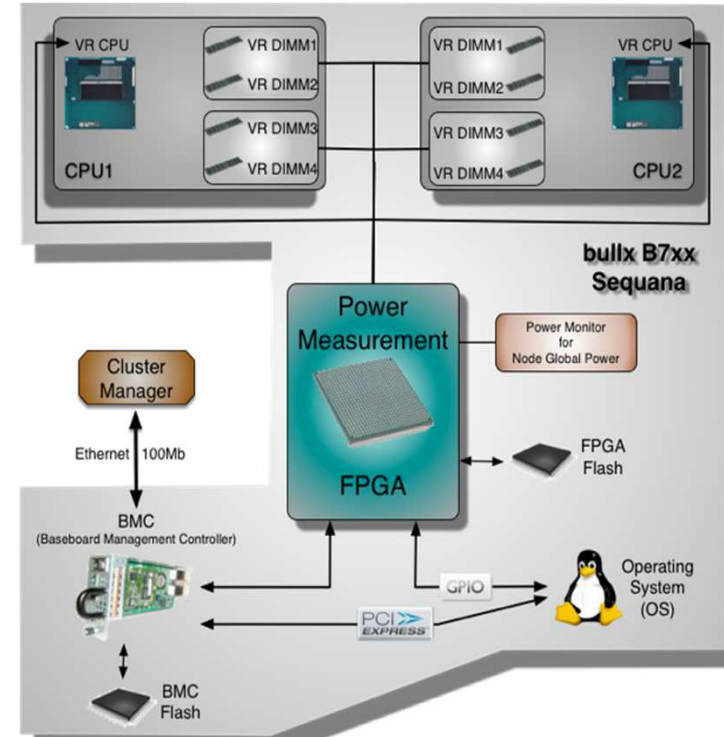
Ultra-efficient PUE

- ▶ Free cooling data centers
- ▶ 100% direct liquid cooling
- ▶ Up to 40°C for inlet water

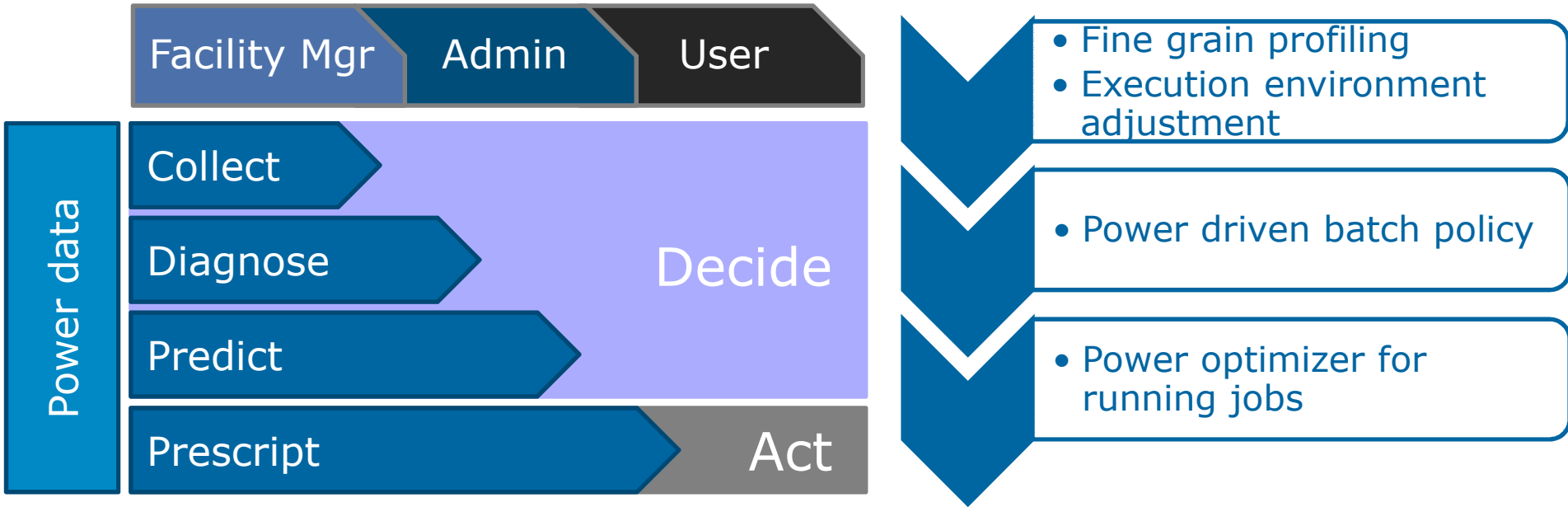


HDEEM: detailed power consumption

- ▶ 1,000 Hz Time stamped sampling
- ▶ CPUs
- ▶ DRAM
- ▶ Interconnect
- ▶ 8 hours ring buffer
- ▶ High accuracy



Prescriptive data analytics for power



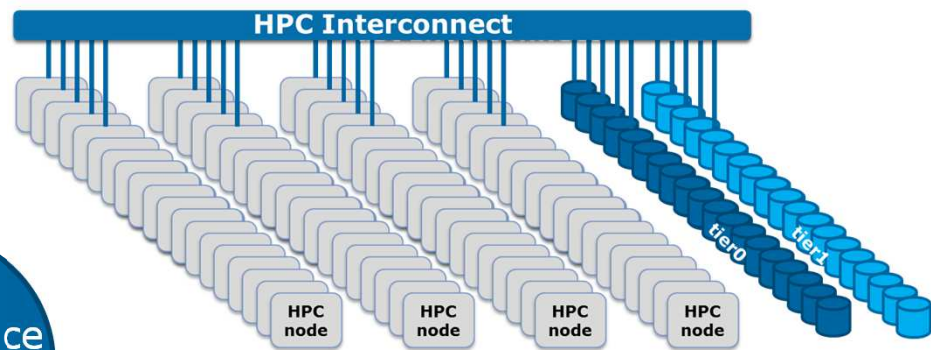
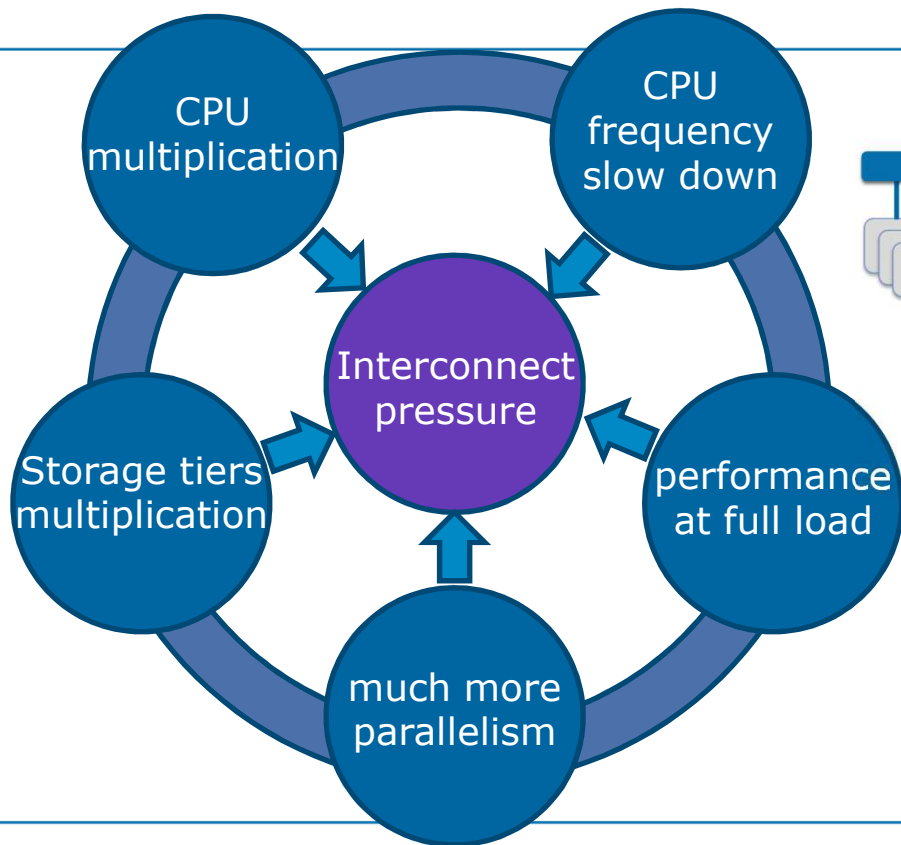
Make energy optimization happen

NEMO (**N**ucleus for **E**uropean **M**odelling of the **O**cean)



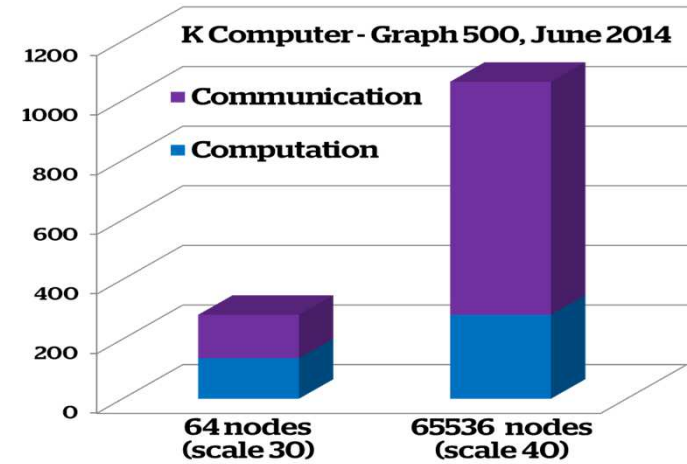
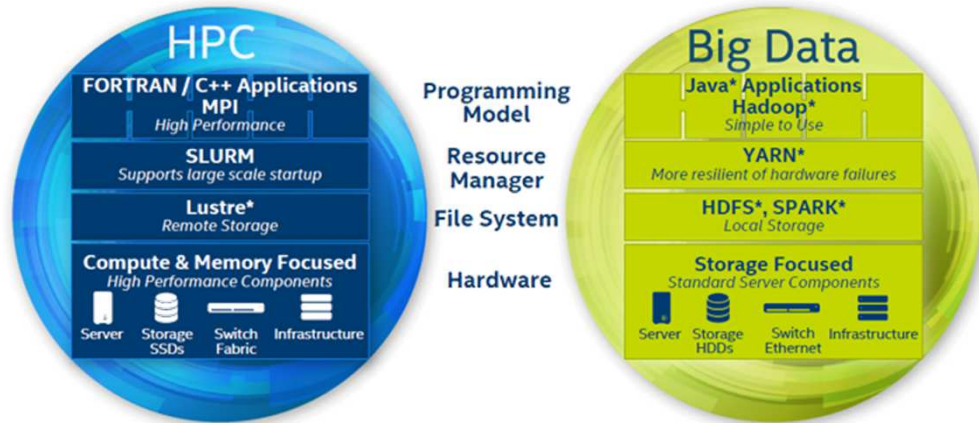
	Fixed frequency	Supervised frequency
Energy (J)	693 704	635 964
Execution time (s)	516	521
Energy saving (J)		8.3%
Execution penalty (s)		1.0%

Supercomputing trend



HPDA: which gap do we have to close?

- ▶ Rationalize the infrastructure
- ▶ Algorithms will be boosted by data capacity and bandwidth, not Flops anymore



Full acceleration in hardware for HPC applications

- HW coded Portals 4, a rich low level network API for message passing (MPI & PGAS)
- Ultra fast path inside the NIC for PGAS / MPI one-sided messaging
- HW offloaded collective operations
- Sustained performance under heavy load

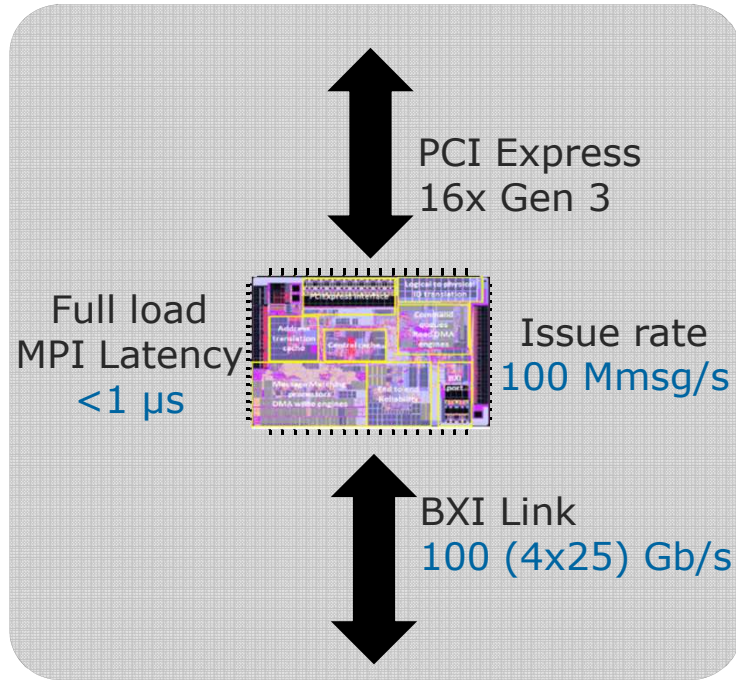
Highly scalable, efficient and reliable

- 64k nodes
- Small memory footprint
- Adaptive Routing
- Quality of Service
- End-to-end error checking + link level CRC + ASIC ECC

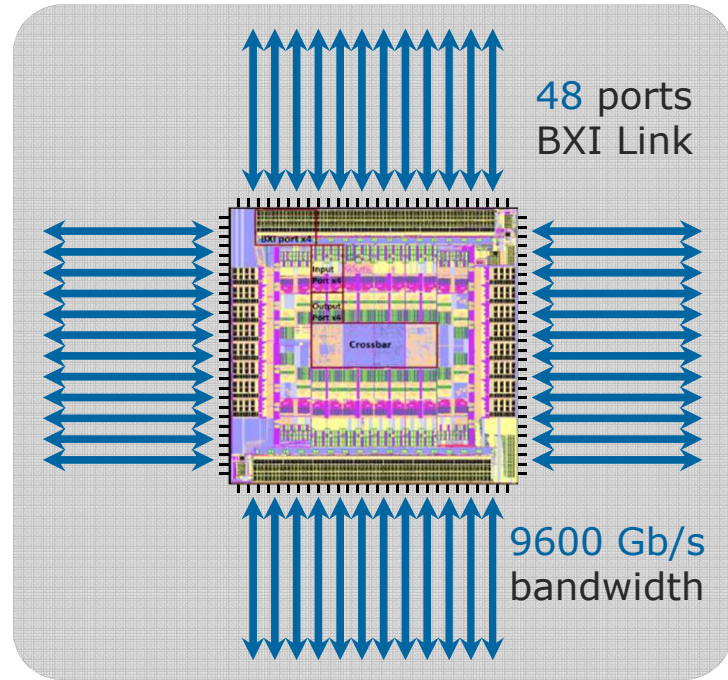
Designed with and for end users

- Co-design in collaboration with CEA

BXI Network is based on 2 in house ASICs

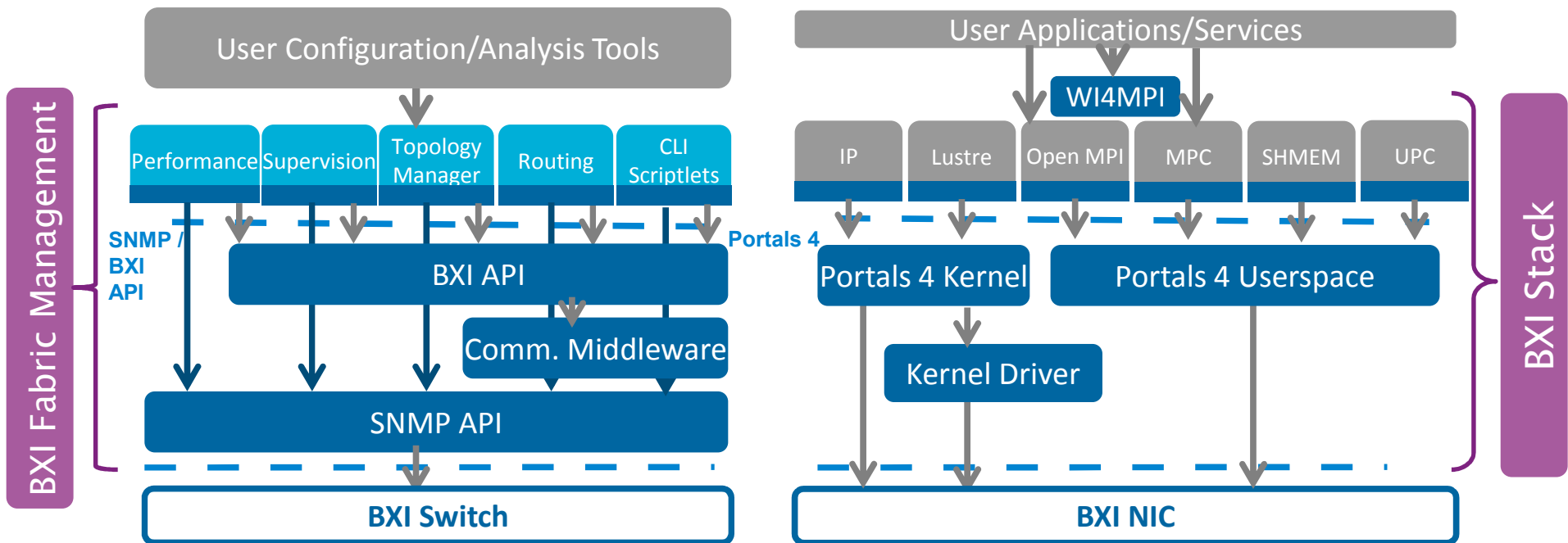


Lutetia



Divio

BXI Software Suite



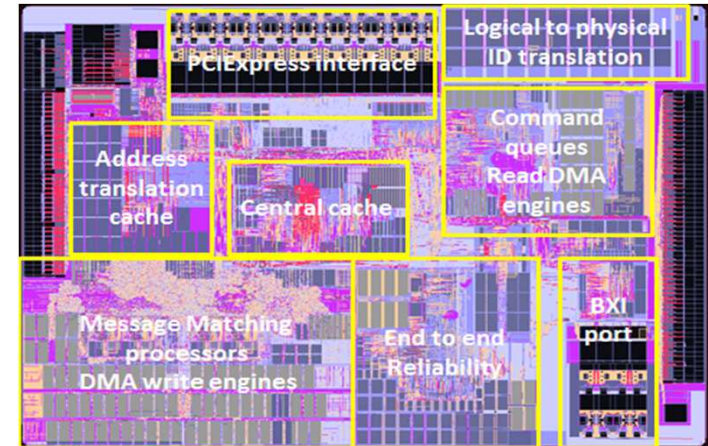
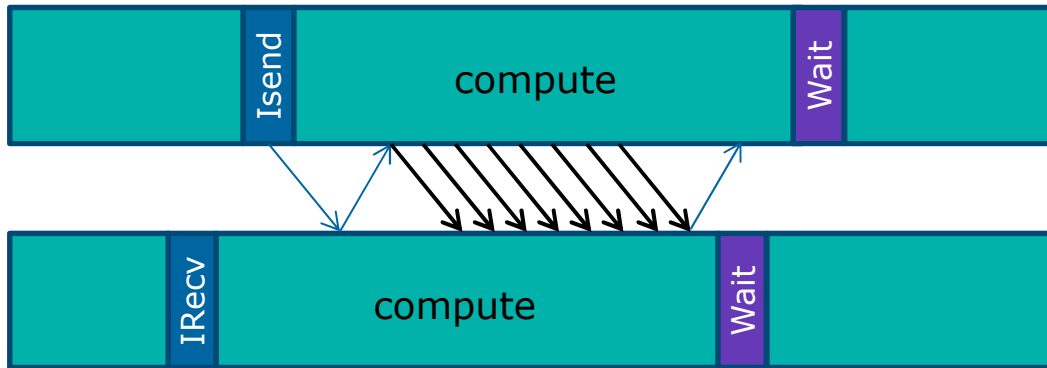
BXI: offloading MPI communication in HW

```

#include <mpi.h>

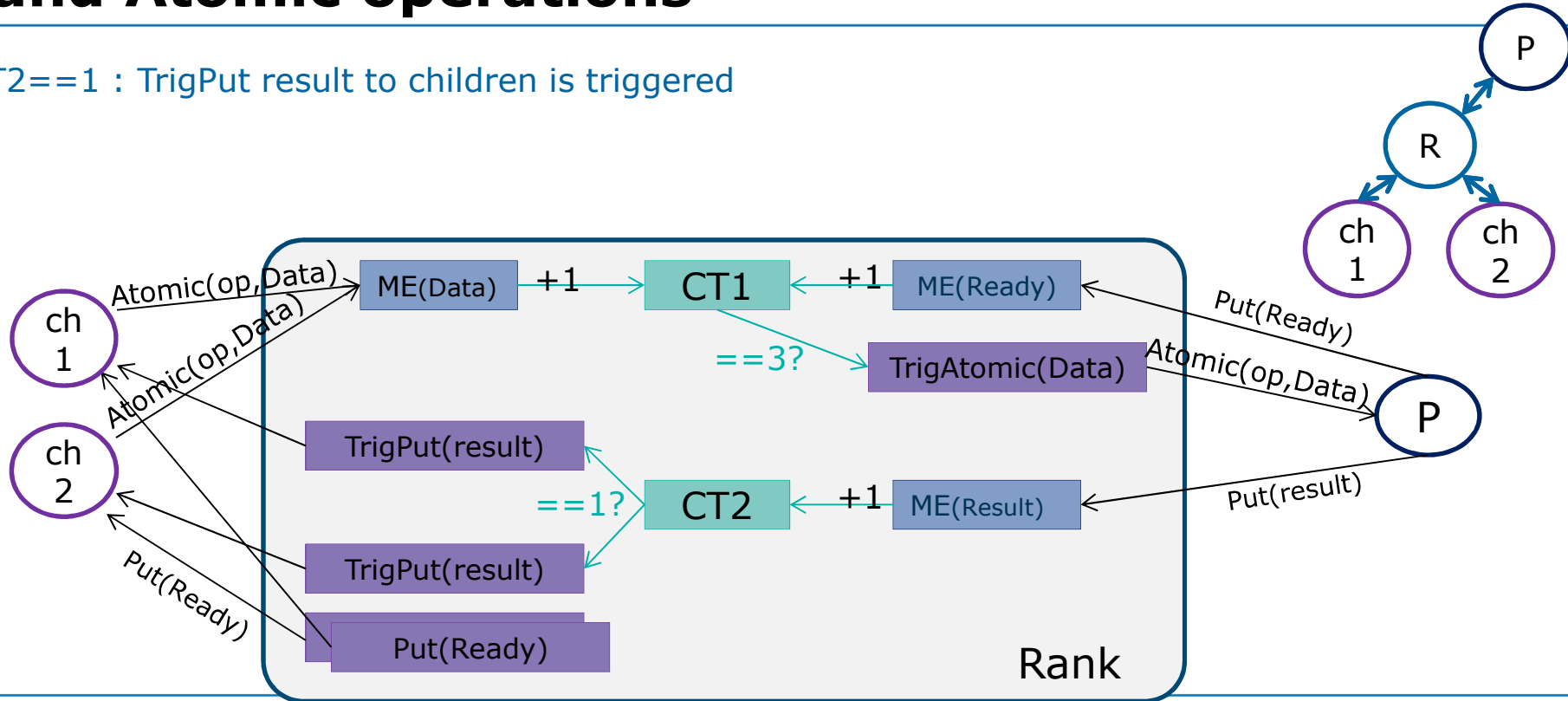
int MPI_Isend( const void *buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm, MPI_Request *request)
int MPI_IRecv(void *buf, int count, MPI_Datatype datatype, int source, int tag, MPI_Comm comm, MPI_Request *request)
int MPI_Wait(MPI_Request *request, MPI_Status *status)
  
```

address V2P ↑
 size ↑
 rank L2P ↑
 message order ↑



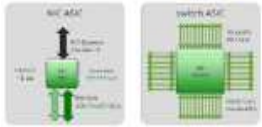
Allreduce example using Triggered and Atomic operations

CT2==1 : TrigPut result to children is triggered



On Track for Exascale era

2016	2017	2018	2019	2020	2021	2022
------	------	------	------	------	------	------



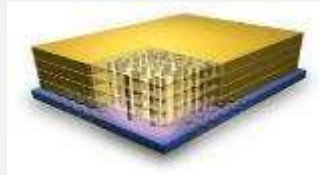
New Exascale
BXI
Interconnect



Sequana
New flexible
packaging



New integrated
computing architecture



Compute + Memory +
Network + Storage
integration



Extreme scale
packaging



Energy &
performance oriented

10¹⁸
Exascale



Thanks

For more information please contact:

T+ 33 4 76297270

F+ 33 4 76297607

M+ 33 6 80357914

eric.monchalin@atos.net

Atos, the Atos logo, Atos Codex, Atos Consulting, Atos Worldgrid, Worldline, BlueKiwi, Bull, Canopy the Open Cloud Company, Unify, Yunano, Zero Email, Zero Email Certified and The Zero Email Company are registered trademarks of the Atos group. August 2016. © 2016 Atos. Confidential information owned by Atos, to be used by the recipient only. This document, or any part of it, may not be reproduced, copied, circulated and/or distributed nor quoted without prior written approval from Atos.

Bull
atos technologies