

The image features the Cray logo at the top, consisting of the word "CRAY" in a bold, blue, sans-serif font with a registered trademark symbol. Below it, "EMEA RESEARCH LAB" is written in a smaller, blue, sans-serif font. The background is a complex digital graphic with a grid of white dots forming a curved surface, overlaid with blue lines and binary digits (0s and 1s).

CRAY[®]
EMEA RESEARCH LAB

**Explicit Data Movement Machinery on Cray systems:
Time-to-solution is not the only game in town**

Adrian Tate, Director, CERL

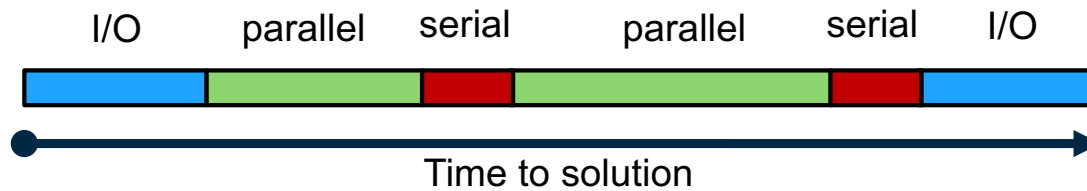
**Utz-uwe Haus, Alessandro Rigazzi, Tim Dykes, Edward Fauchon-Jones,
Karthee Sivalingam, Harvey Richardson, Nina Mujkanovic, Aniello Esposito, Clement Foyer**

On my travels I hear:

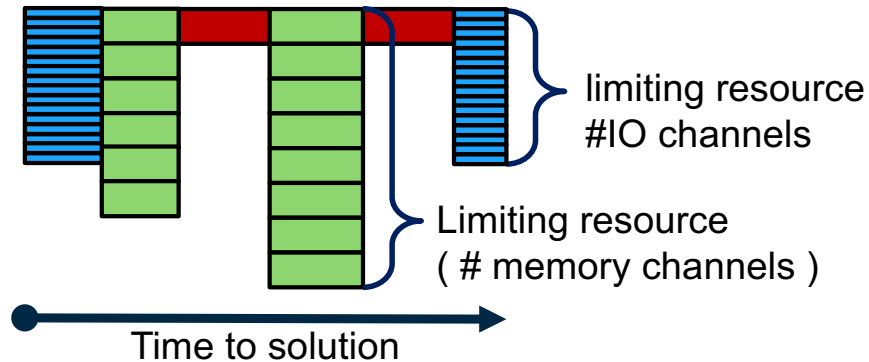


- Data Movement is key**
- Massive interest in flexibility and insulation/abstraction**
 - Cpu / accelerator
 - Memory
 - Interconnect
 - Abstraction layers to not tie into an architecture
- Technical Collaboration is a must**
 - Challenges are too hard!
 - Memory is 20x slower than 1980
 - The complex memory hierarchy isn't even here yet!
 - When it does come, how will we use it?
- HPC + AI : yes but in more than DL frameworks**
- Time-to-solution is not the only game in town**

Time to solution

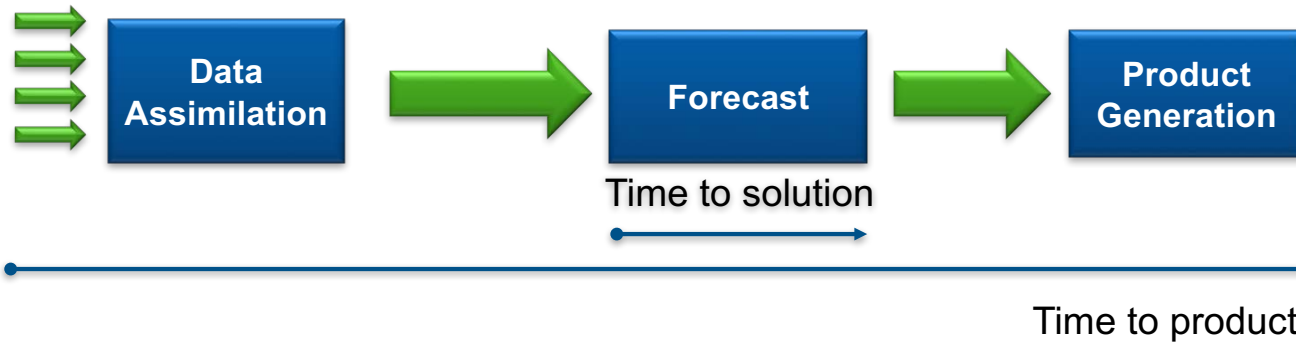


- Many forms of parallelism
- Algorithmic advances
- Code optimization (Compiler & hand)
- ISA features
- Programming Models
- Performance Abstraction
- Systems Software / Operational
- Network, memory increases



Time to scientific product / insight

NWP:

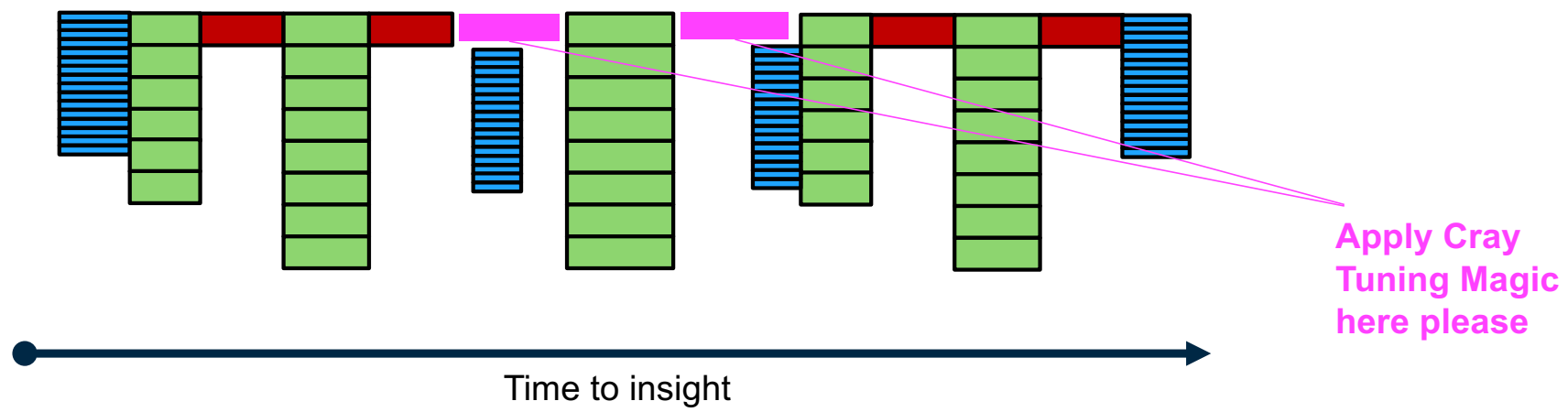
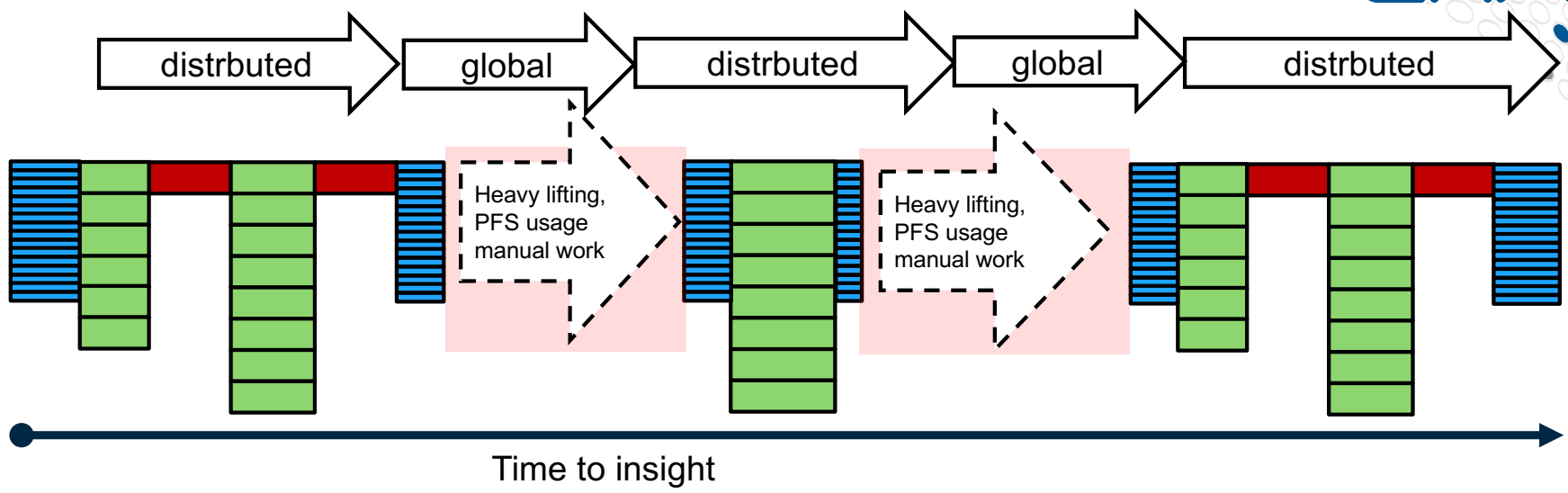


**Neuro-
science:**



What does  mean?

How do we optimise  to be smaller?



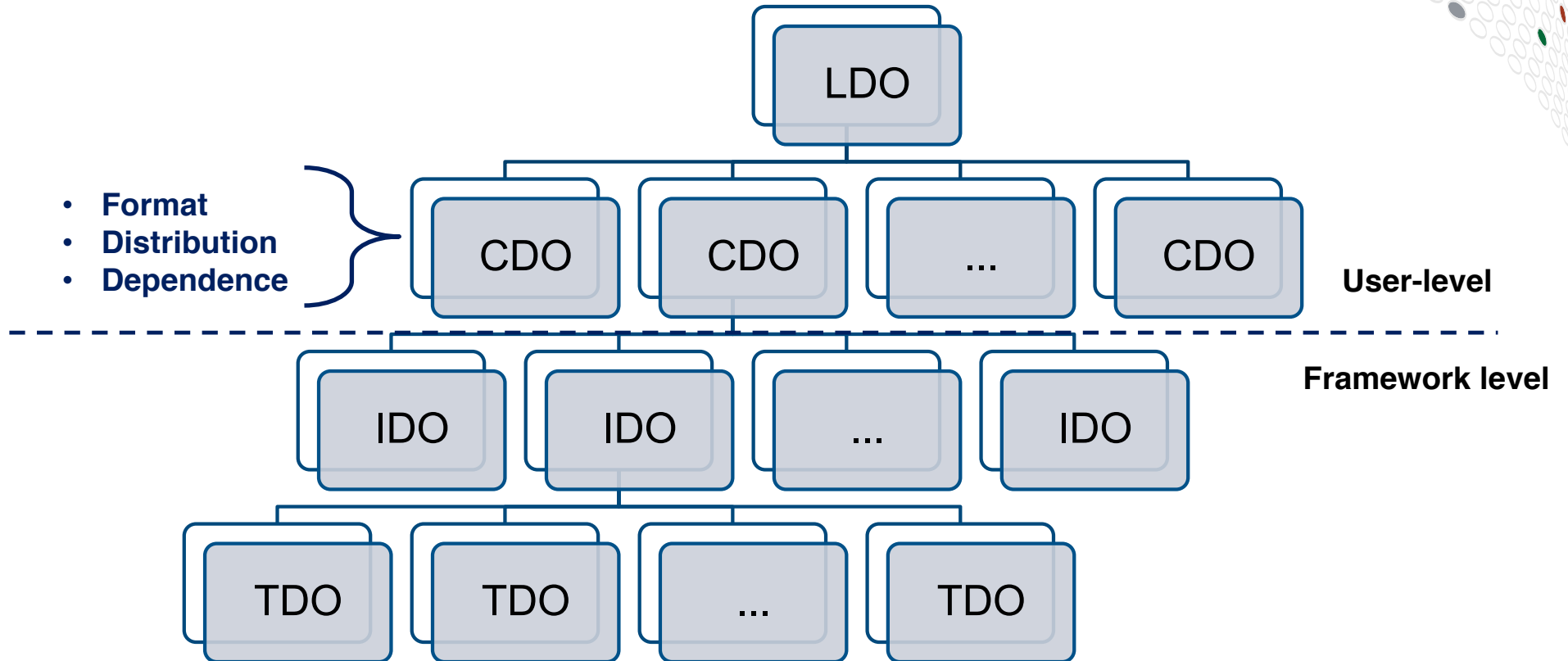
How to start looking at this problem



- Lacking in abstractions to tackle the problem**
 - No operators to minimise
 - MPI, PGAS express data movement semantics inside job
 - Distinct jobs need to export credentials
 - How to represent the data, its format, its distribution

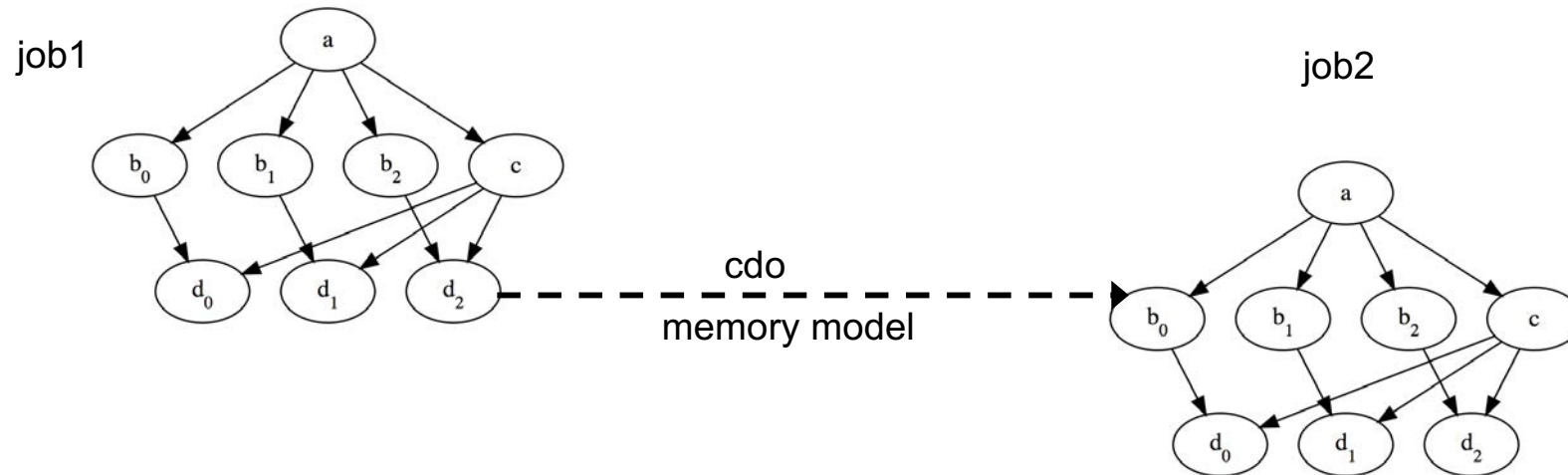
- Most importantly – parallel data needs to be redistributed**
 - Ignoring the fact that data is parallel is ignoring the general problem

Cray Data Object Model



“Octopus” Concept and Project

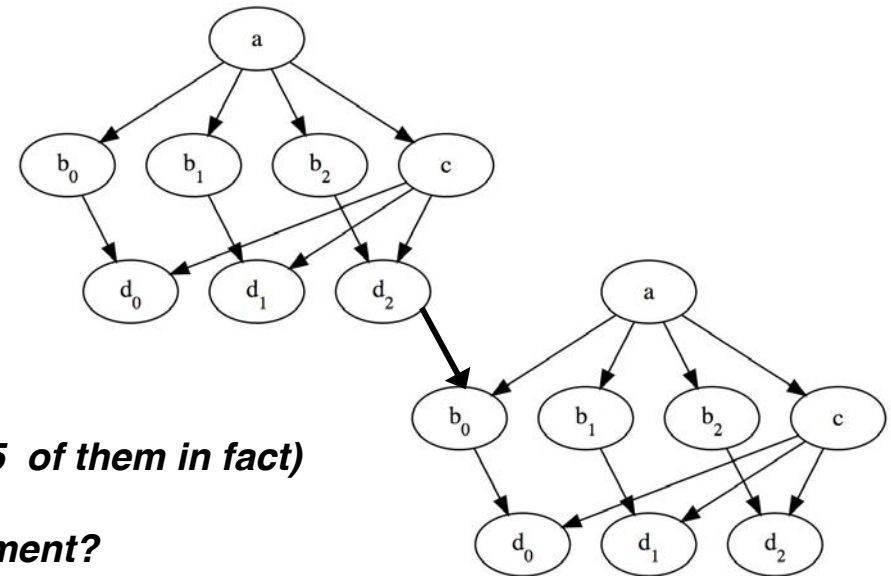
1. Use object hierarchy to express data dependencies (at a task level)
2. Build a useful model of memory hierarchy
3. Resolve task graph into execution graph
4. Execute graph with data in correct resource



Octopus Concept



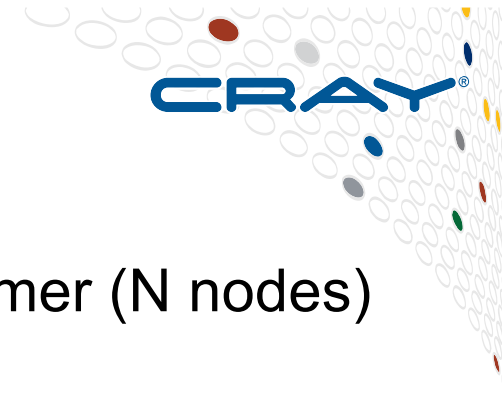
1. Use CDO to represent data dependencies
2. Build a useful model of memory hierarchy
3. Resolve task graph into execution graph
4. Execute graph with data in correct resource



There are no shortage of tools to execute task graph (185 of them in fact)

How many of those can manage the (parallel) data movement?

Universal Data Junction (UDJ)

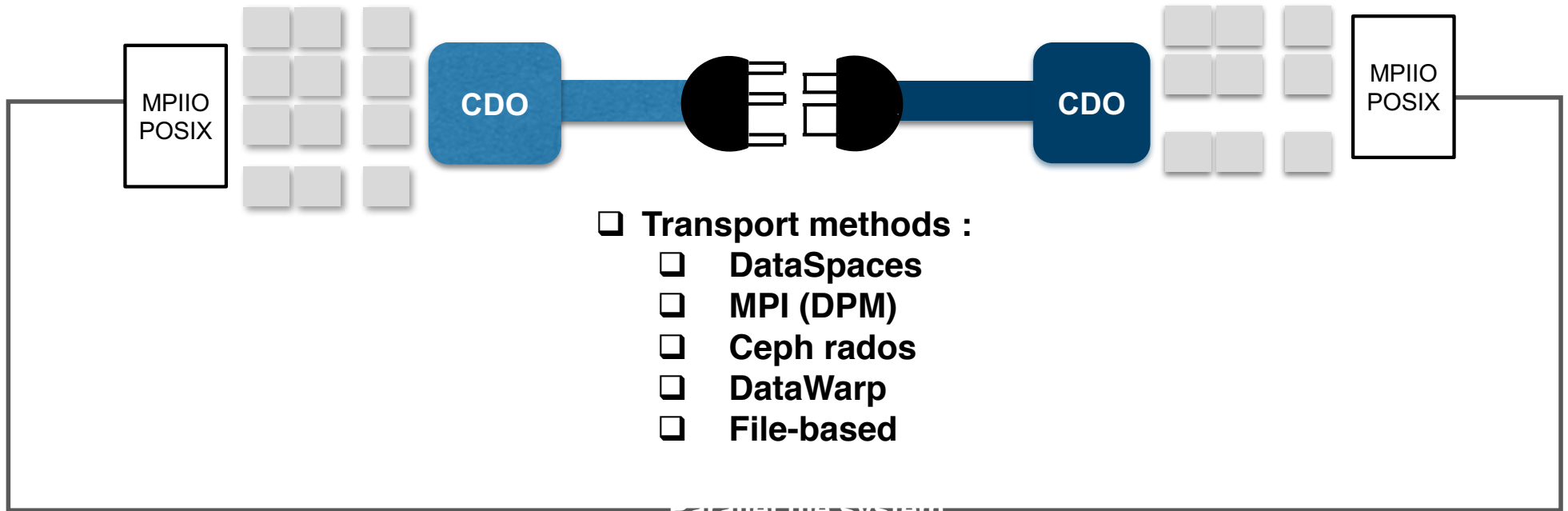


Producer (M nodes)

Consumer (N nodes)

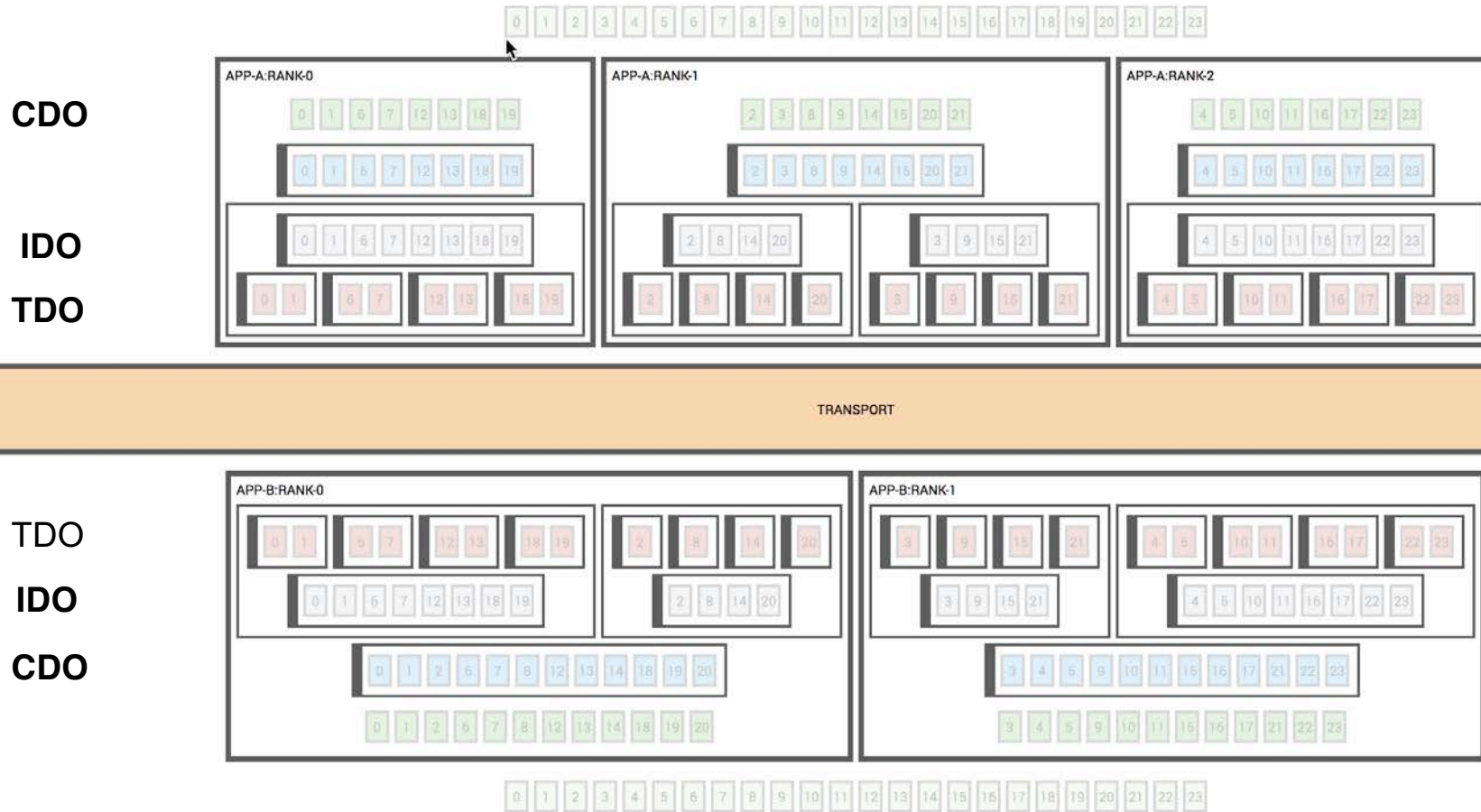
- Distribution (contig, none, cyclic)
- Format (array, HDF5, Conduit, text)

- Distribution
- Format



- Transport methods :
 - DataSpaces
 - MPI (DPM)
 - Ceph rados
 - DataWarp
 - File-based

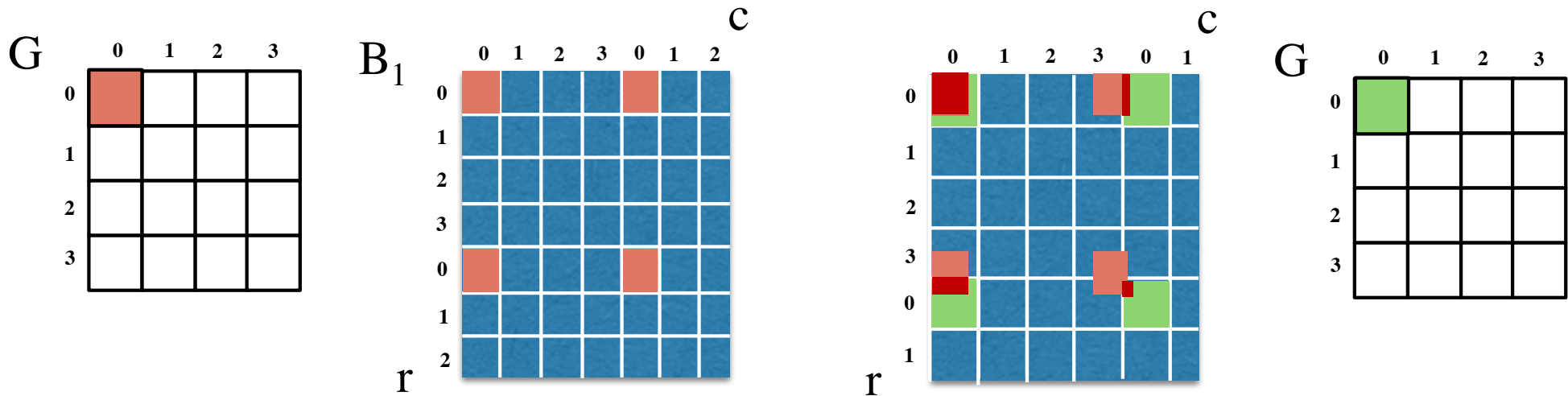
How the redistribution map onto Core Data Object Model



Non-triviality of Producer-Consumer Redistribution

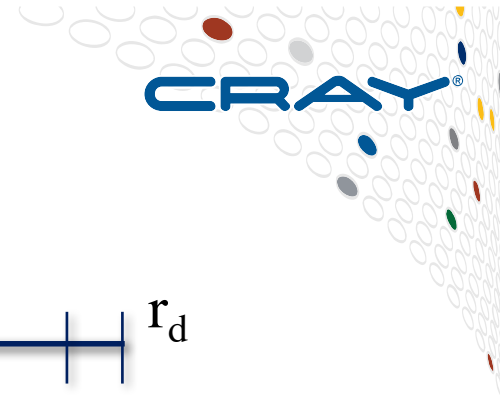


- 2d data set dim $r \times c$ in memory
- Distributed according to some distribution scheme $D_1 = (G, B_1)$



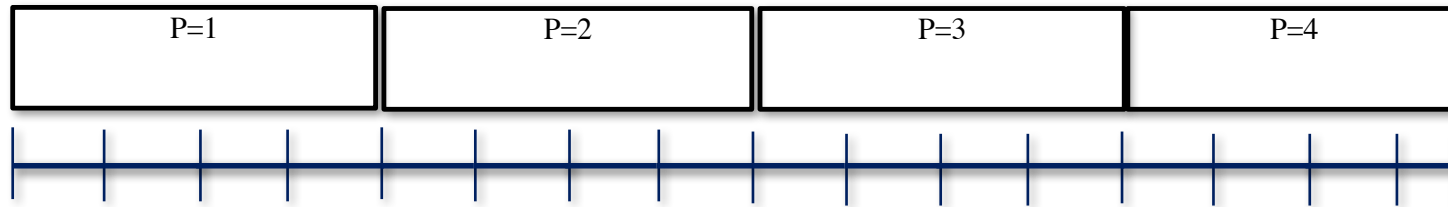
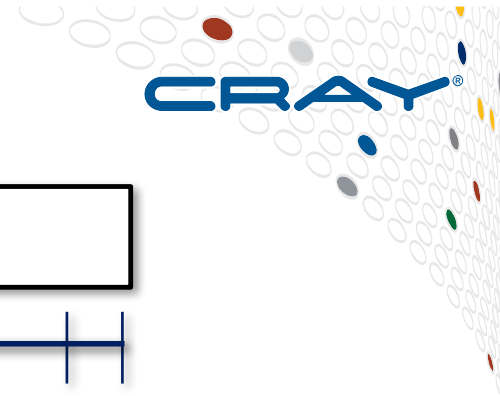
- Re-distributed according to new distribution scheme $D_2 = (G, B_2)$ on same grid G
- Must communicate the non-trivial intersection data (red) for every process pair

Classical Redistribution algorithm



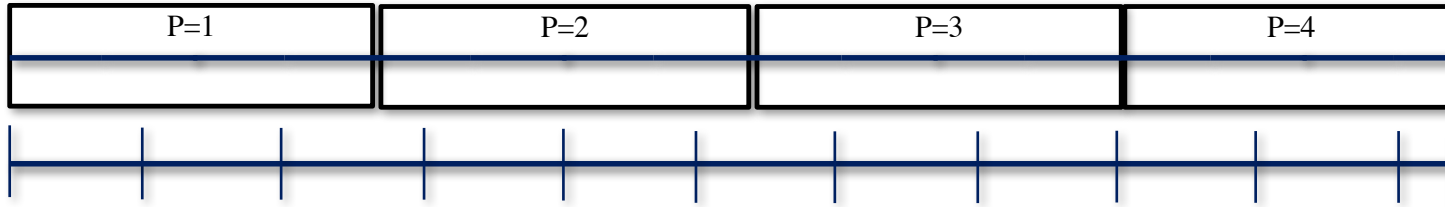
- For each d in #dimensions
- Consider the vector of length r_d
- Divide by blocks of length b_d

Classical Redistribution



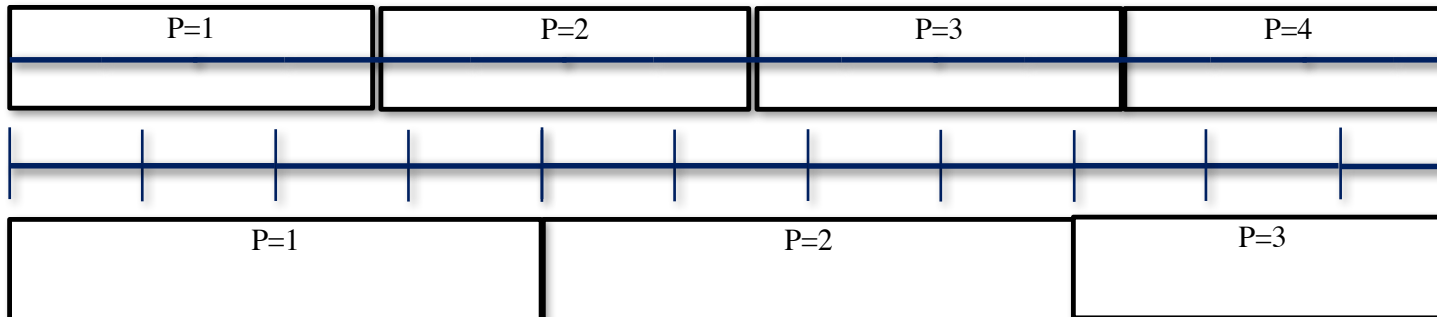
- For each d in #dimensions
- Consider the vector of length r_d
- Divide by blocks of length b_d
- Map the blocks to process rows/column

Classical Redistribution



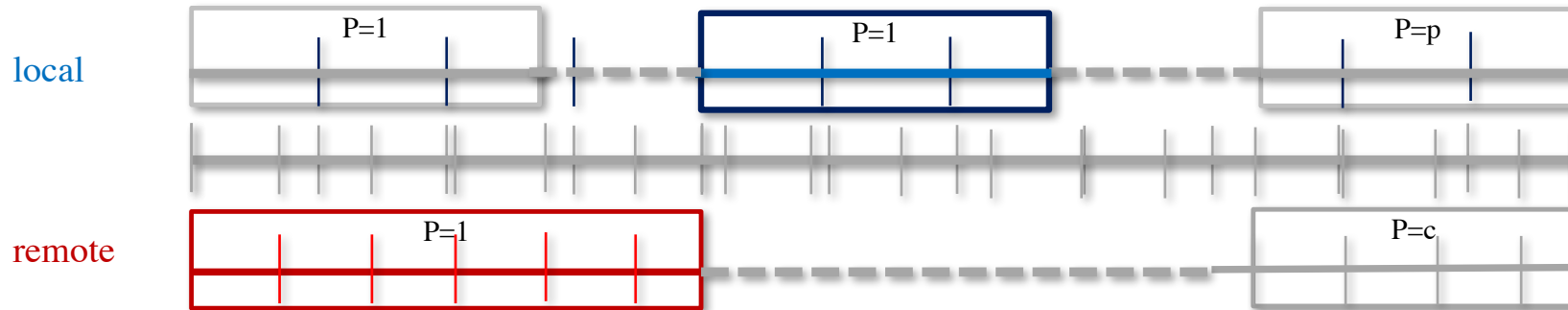
- For each d in #dimensions
- Consider the vector of length r_d
- Divide by blocks of length b^{prod}_d
- Map the blocks to process rows/column
- Divide consumer by blocks of length b^{cons}_d

Classical Redistribution



- For each d in #dimensions
- Consider the vector of length r_d
- Divide by blocks of length b^{prod}_d
- Map the blocks to process rows/column
- Divide consumer by blocks of length b^{cons}_d

Classical Redistribution



On each producer rank:

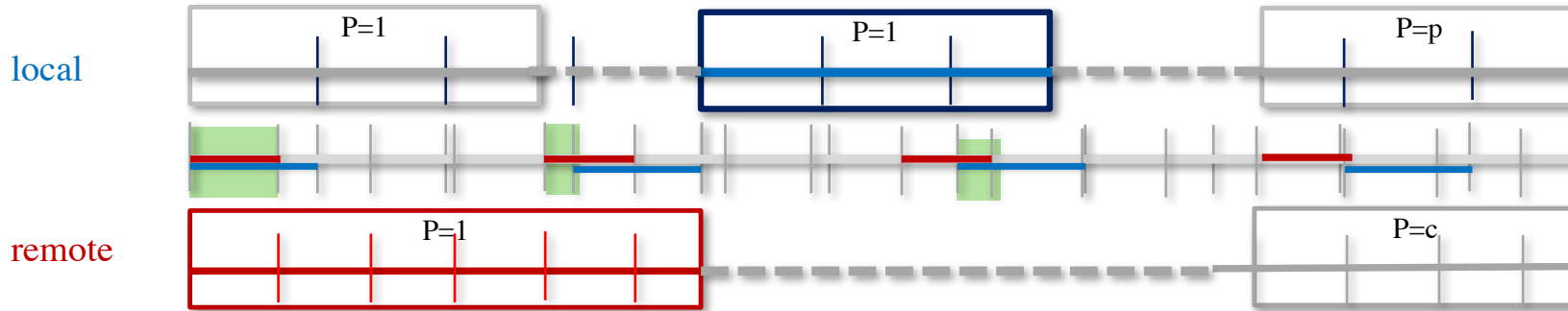
For each d in #Dimensions

For each p in length(consumer_grid(d))

For each loc in #NumLocalBlocks

For each rem in #NumRemoteBlocks

Classical Redistribution



On each producer rank:

For each d in #Dimensions

For each p in length(consumer_grid(d))

For each loc in #NumLocalBlocks

For each rem in #NumRemoteBlocks

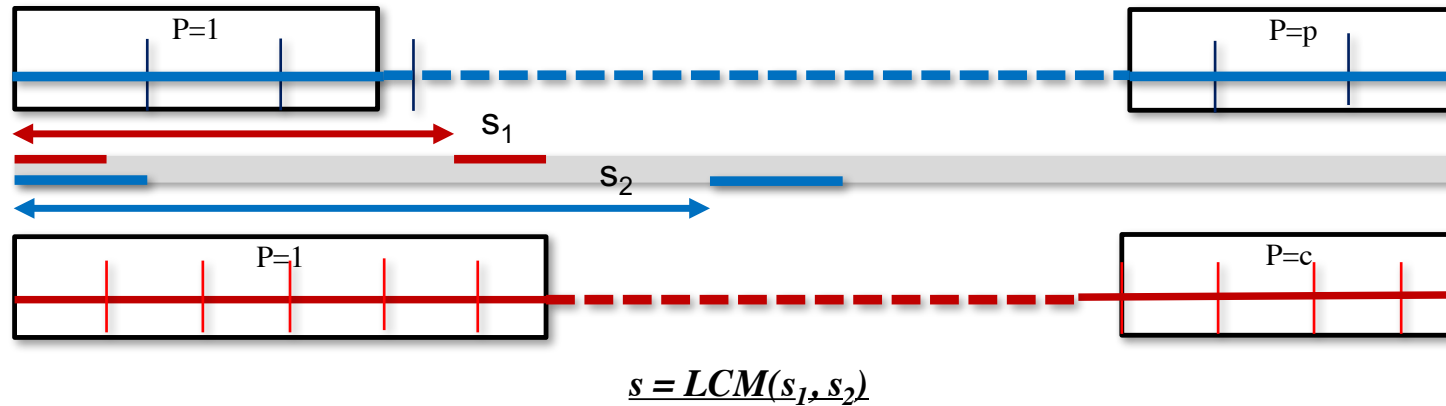
```
if MAX(loc2glob(loc), loc2glob(rem)) <
    MIN(loc2glob(loc+b1), loc2glob(rem+b2)) → Add to intersection
```

Intersection = $i_1 \times i_2 \dots i_d$

Complexity: $O(\#Dim \cdot P \cdot C \cdot n^{local} \cdot n^{remote})$

Ignores three types of periodicity!

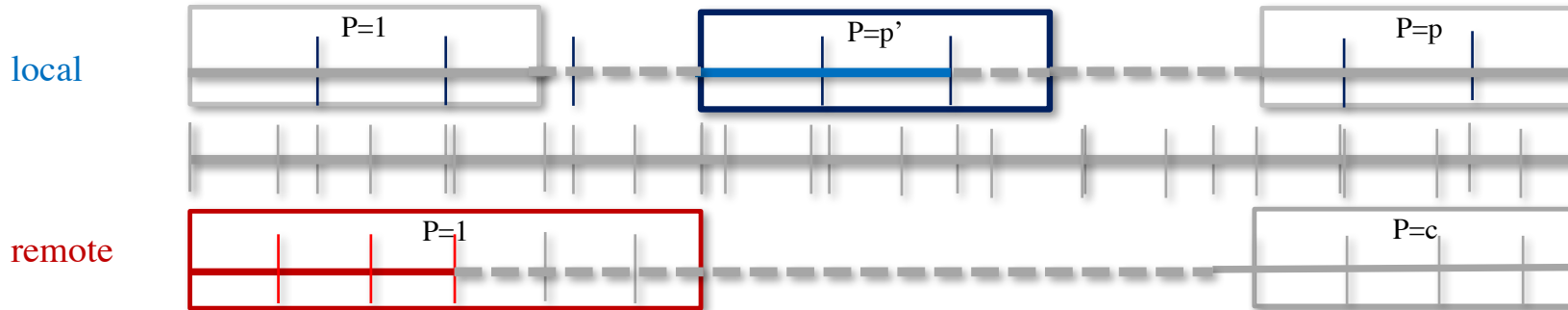
Exploit Periodicity of Intersection



Theorem: For a process pair (p_1, p_2) successive blocks of local data are separated globally by s_{1l} and successive blocks of remote data are separated by s_2 then for every element i in the intersection, element $i + s = LCM(s_1, s_2)$ will also be in the intersection

Source: Guo/Nakata "A Framework for Efficient Data Redistribution on Distributed Memory Multicomputers"

Guo-Nakata / FALLS Redistribution



On each producer rank:

For each d in #Dimensions

For each p in length(consumer_grid(d))

For each loc in #NumLocalBlocks

For each rem in #NumRemoteBlocks

Seems unnecessary

```
if MAX(loc2glob(loc), loc2glob(rem)) <
    MIN(loc2glob(loc+b1), loc2glob(rem+b2)) → Add to intersection
```

Intersection = $i_1 \times i_2 \dots i_d$

Complexity : $O(\#Dim . P . C . \hat{n}_{local} \hat{n}_{remote})$

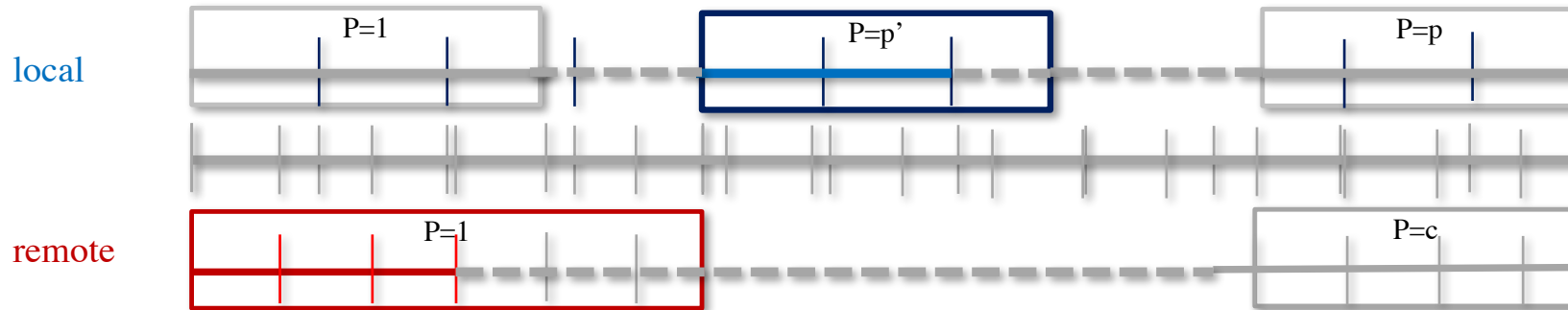
Still ignores 2 types of periodicity!

Exploit Periodicity of Remote Data



- Remote data is of a periodic form
- We can remove one loop if we construct periodic (modular) relations

Exploit Periodic Remote Data



On each producer rank:

For each d in #Dimensions

For each p in length(consumer_grid(d))

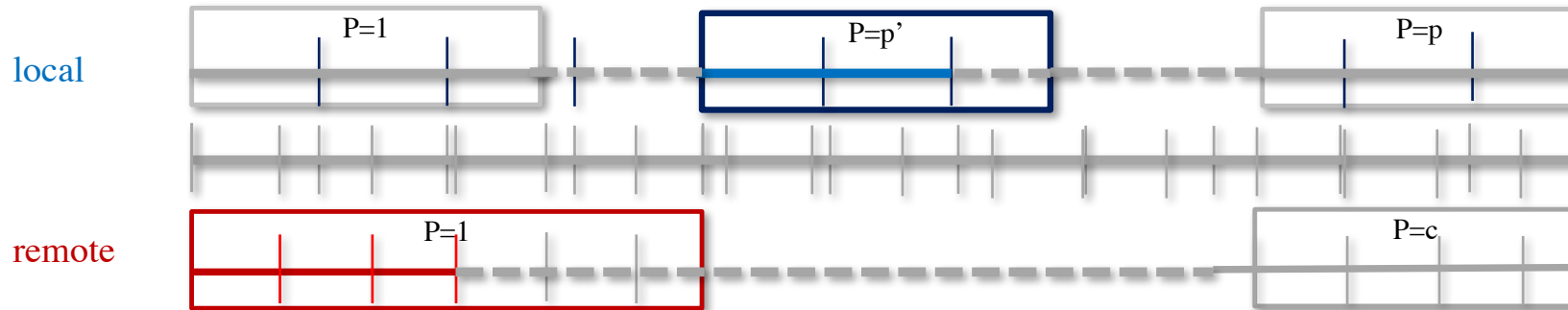
For each loc in #NumLocalBlocks

if ($loc2glob(loc) \% s2$) $\leq b2$ \rightarrow Add to intersection

Complexity: $O(\#Dim \cdot P \cdot C \cdot \hat{n}_{local})$

Intersection = $i_1 \times i_2 \dots i_d$

Exploit Periodic Remote Data



On each producer rank:

For each d in #Dimensions

For each p in length(consumer_grid(d))

Loop order prevents further optimisation

For each loc in #NumLocalBlocks

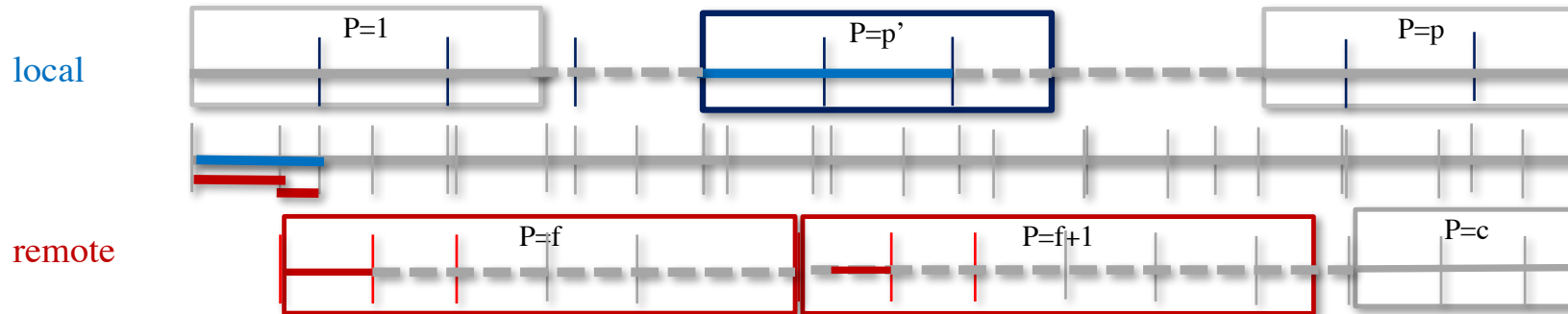
```
if (loc2glob(loc) % s2 ) <= b2 → Add to intersection
```

Intersection = $i_1 \times i_2 \dots i_d$

Complexity: $O(\#Dim . P . C . \hat{n}_{local})$

Exploit Periodic Nature of Grid Data

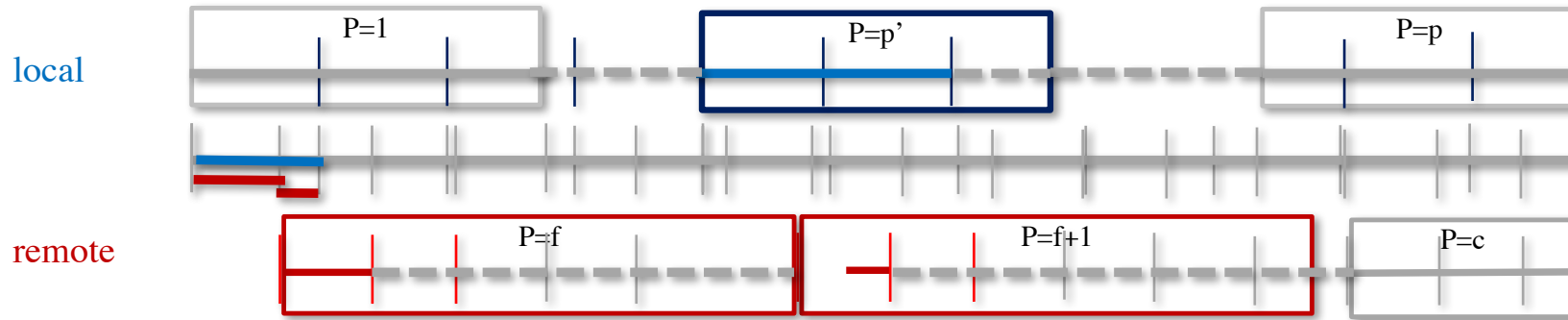
CRAY®



Theorem: *If local block l shares intersection with remote processor f , then if an adjacent block of length b_2 also exists, that block will share intersection with remote processor $(f+1)\%c$, where c is the length of the consumer grid*

Source: Foyer and Tate, "Efficient data redistribution for Producer-Consumer Grids" (xarchiv)

Exploit Periodic Nature of Grid Data



On each producer rank:

For each d in #Dimensions

For each p in length(consumer_grid(d))

if (loc2glob(loc) % s_2) <= b_2 → Add to intersection

For each sub in b_1 / b_2

→ Add sub to intersection

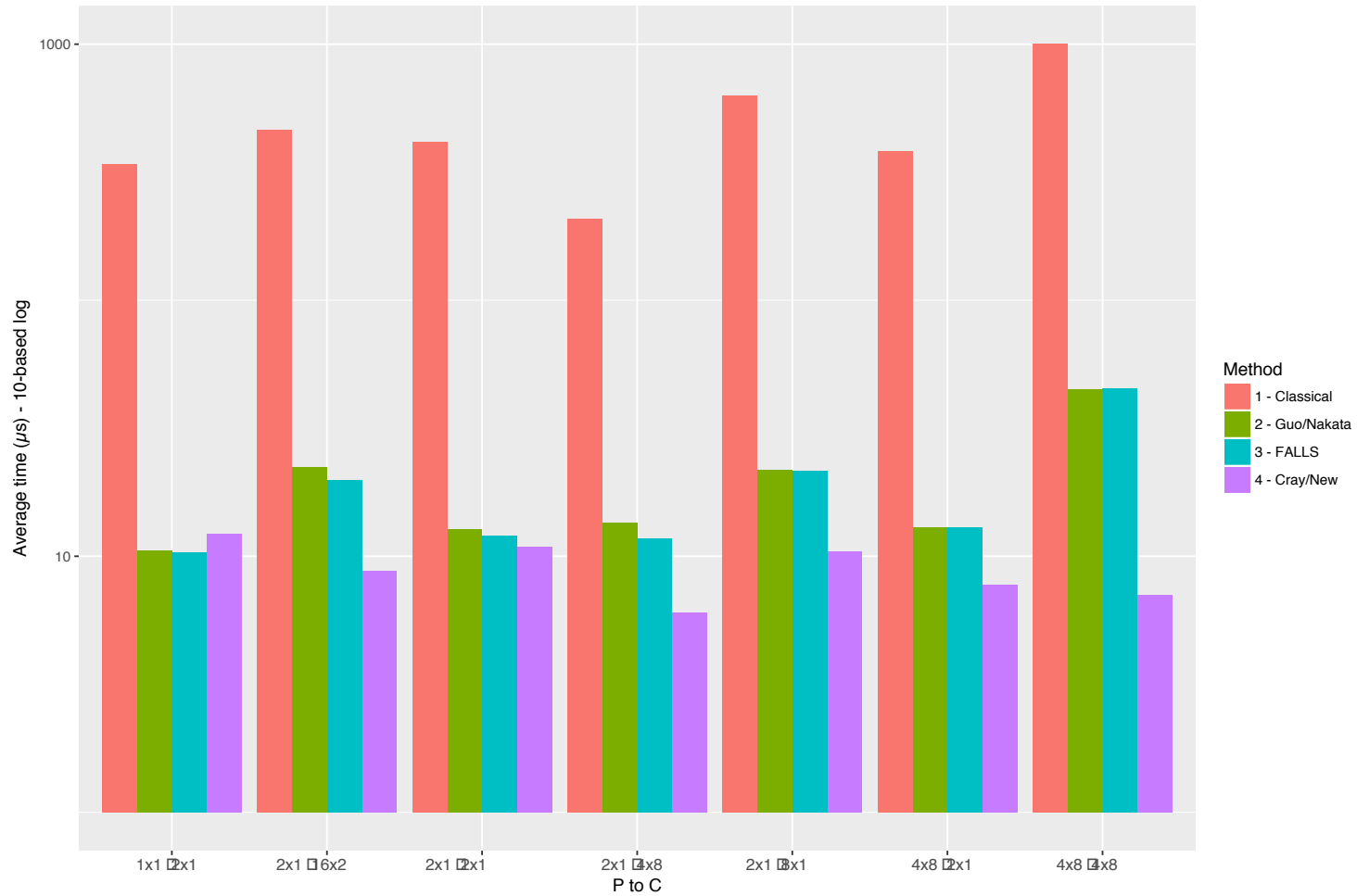
Intersection = $i_1 \times i_2 \dots i_d$

Complexity: $O(\#Dim \cdot P \cdot \hat{n}_{local} \cdot b_1/b_2)$

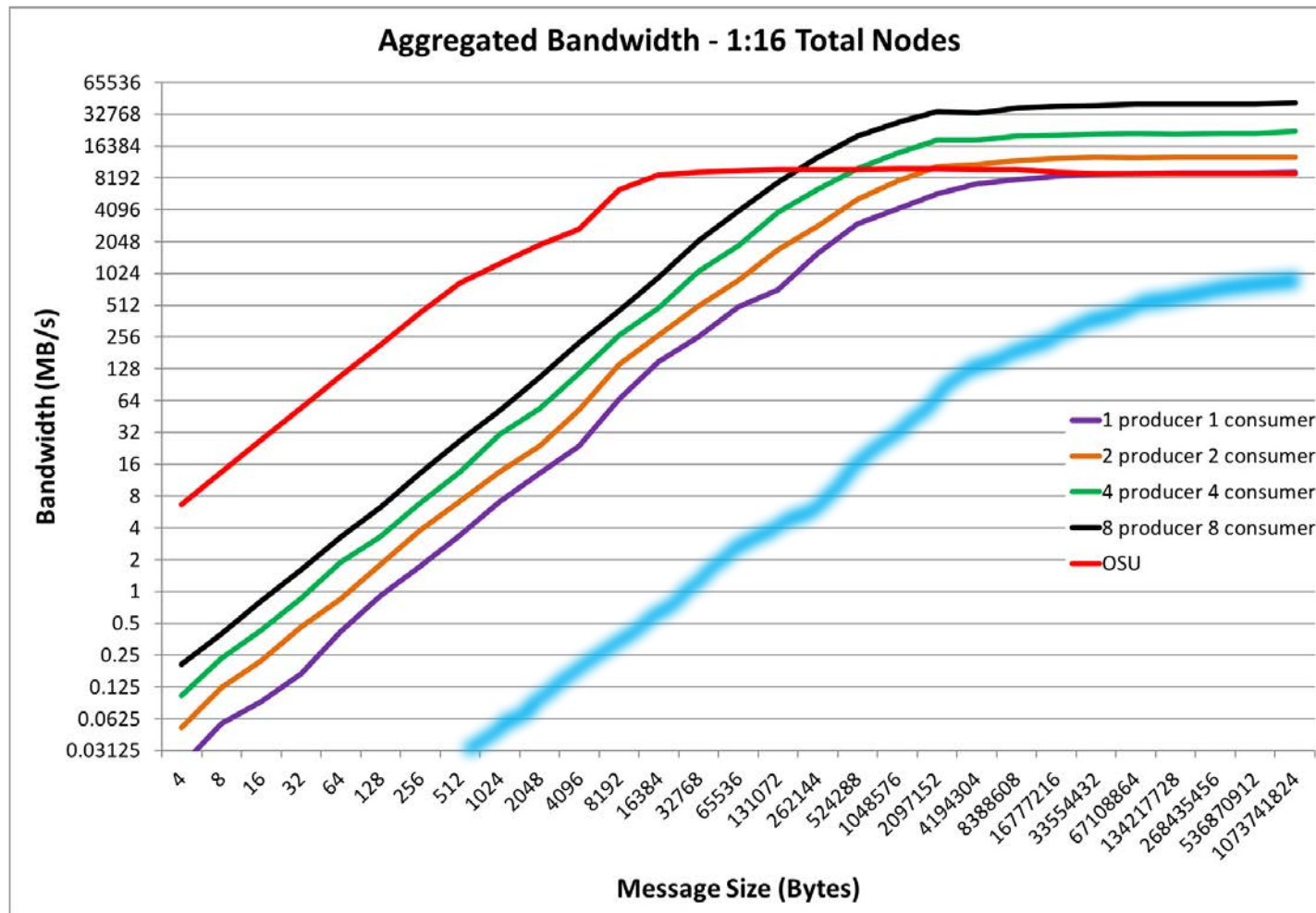
Results – M:N Node Redistribution

Average time for P to C redistribution (10-based log)

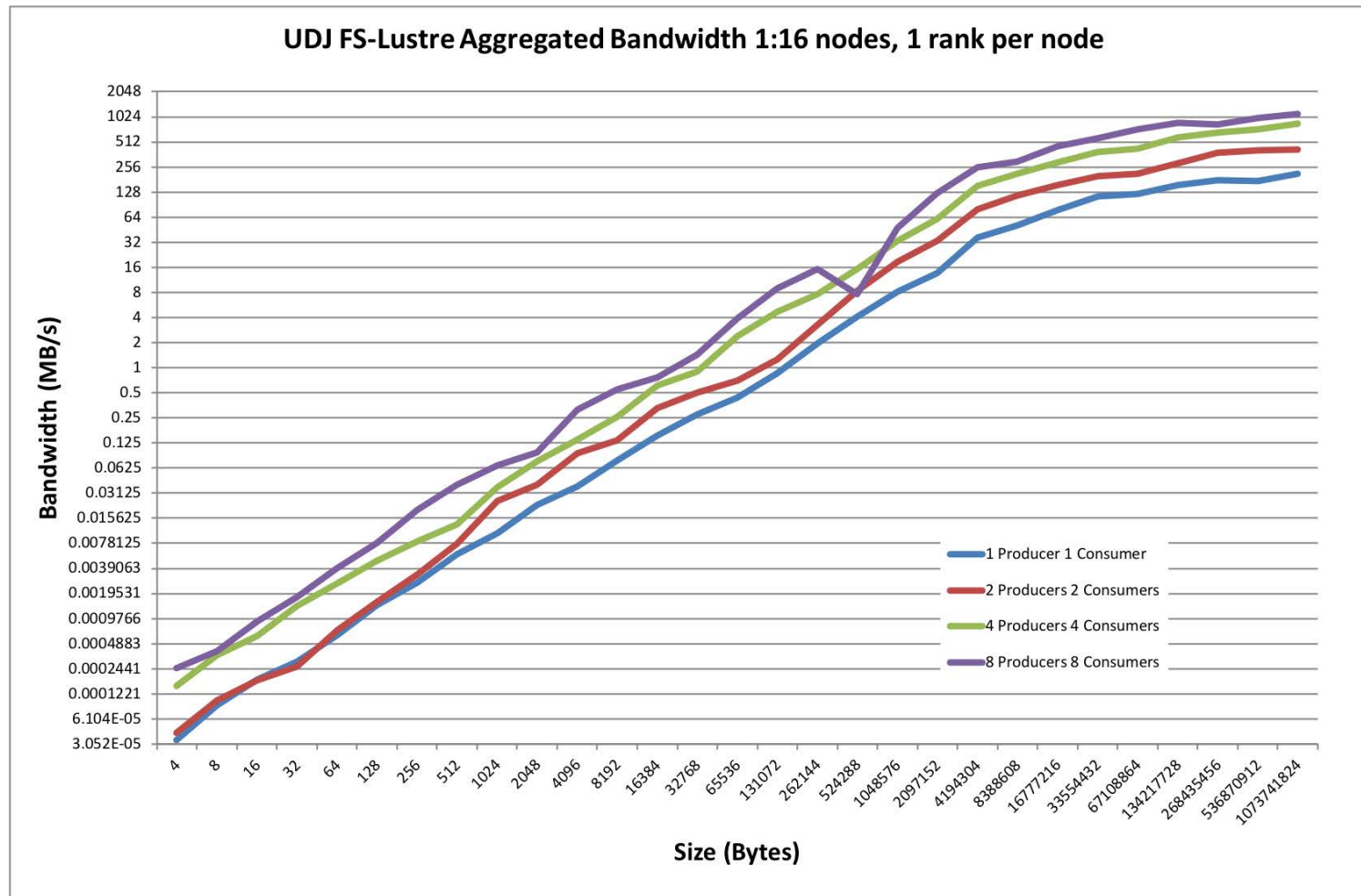
Producer block size: 256x256, Consumer block size: 1024x1024



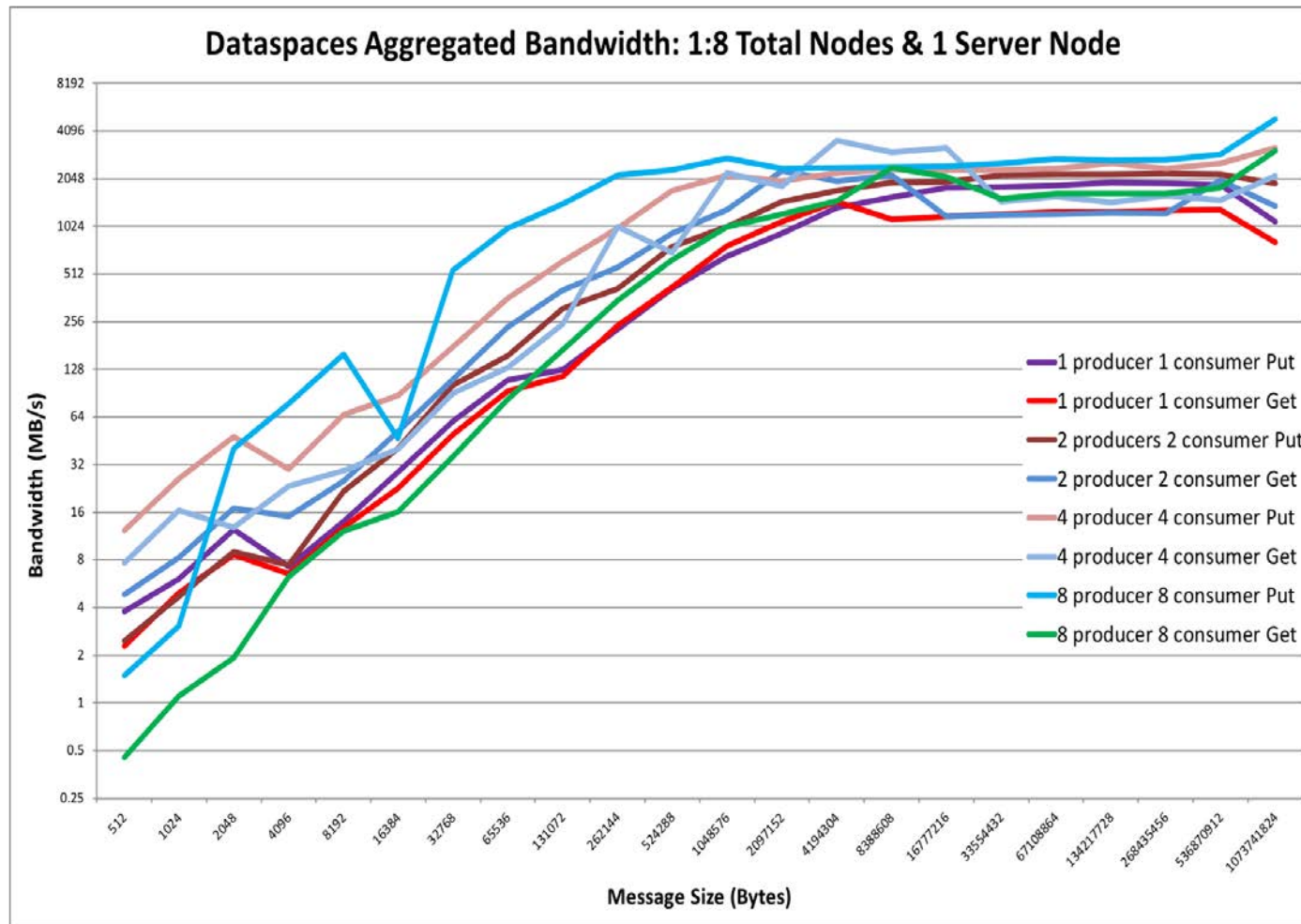
Results UDJ - MPI



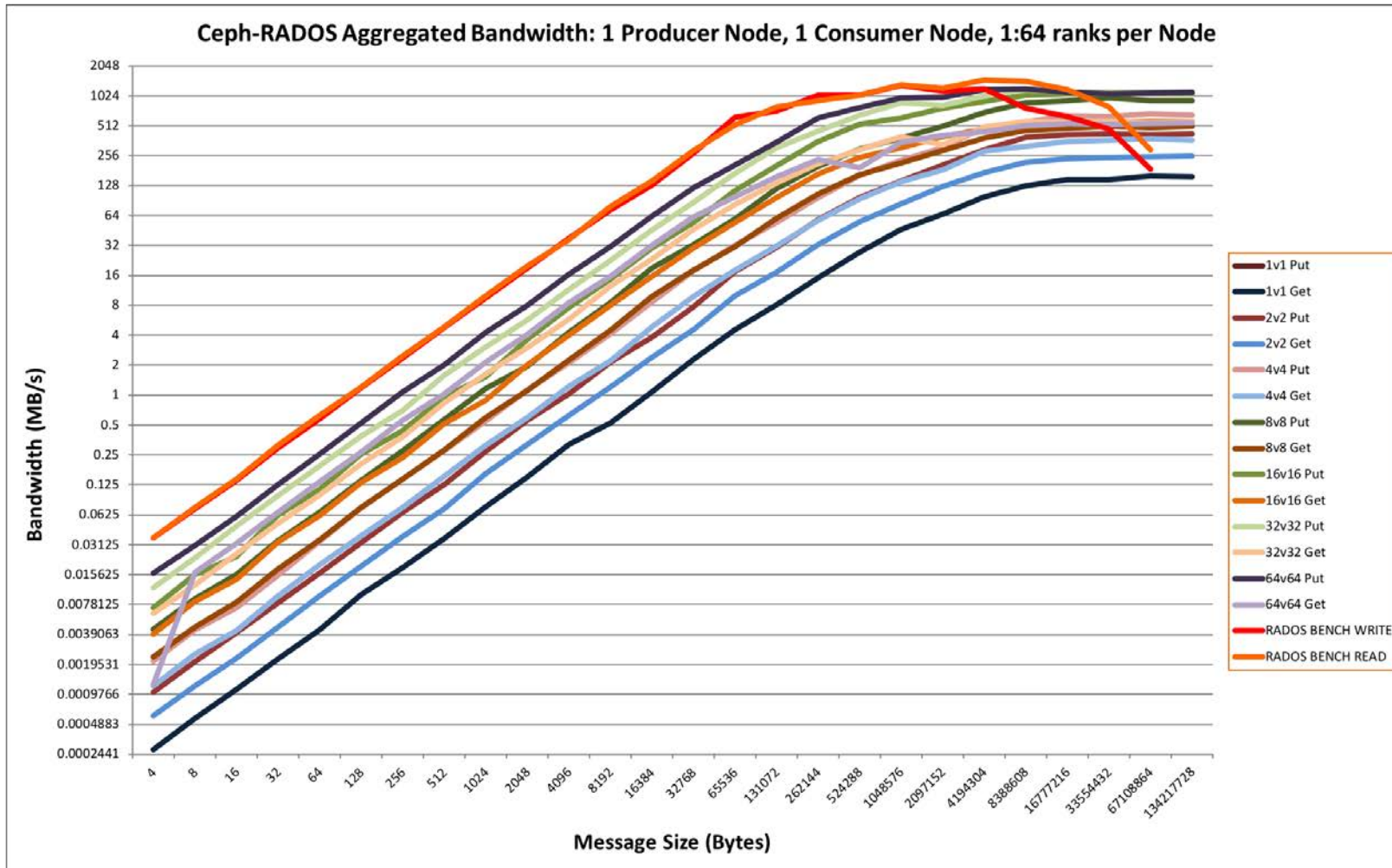
Results UDJ - Lustre



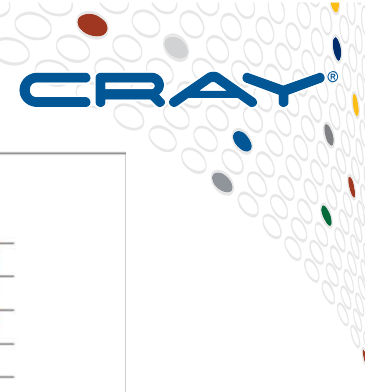
Results – UDJ DataSpaces



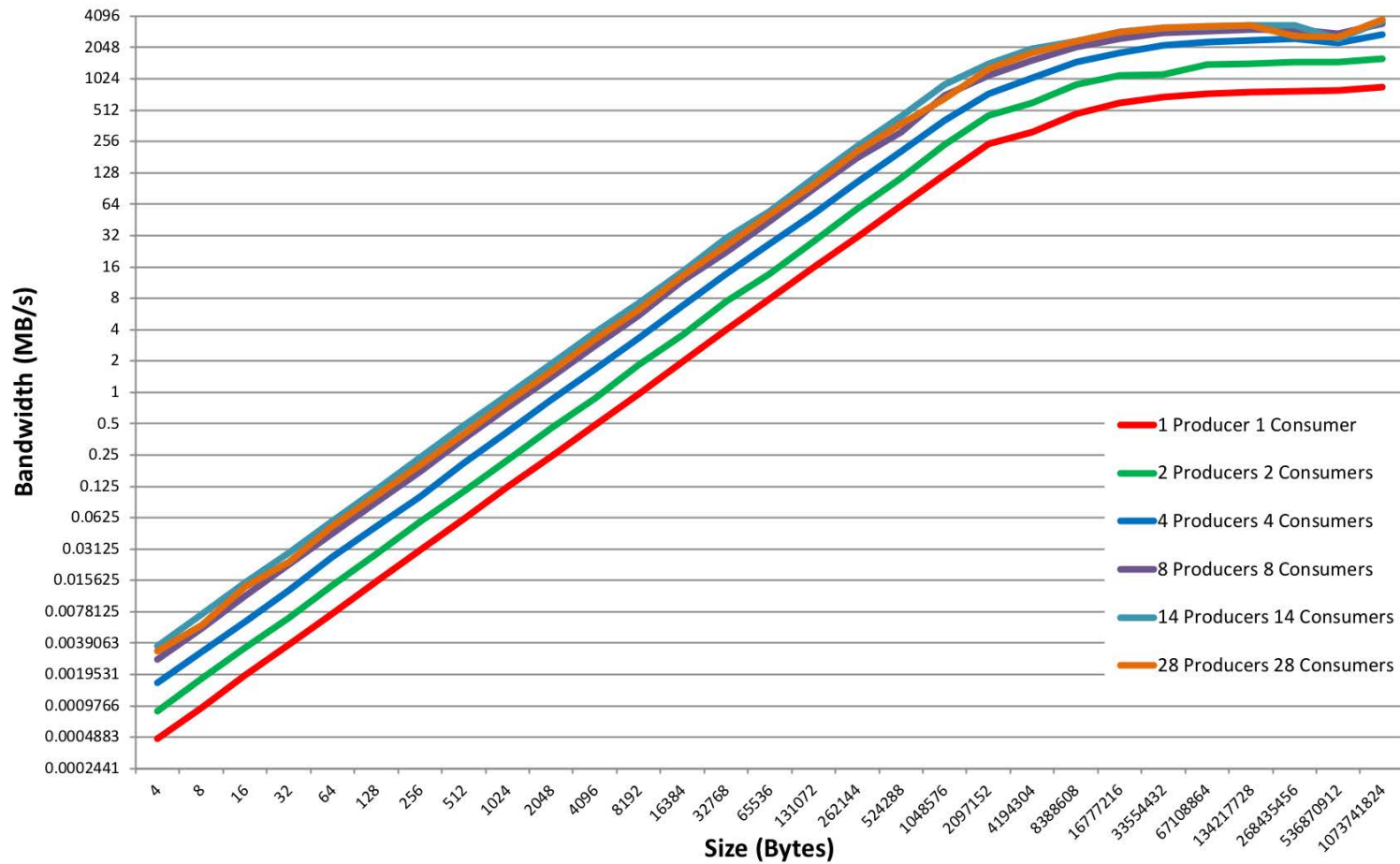
Results UDJ – Ceph RADOS



UDJ on DataWarp



FS Datawarp Transport : 1 Producer Node, 1 Consumer node, 1:28 ranks per node



Where Next?



- Memory-and data-Aware workflow middleware



- Memory Hierarchy Abstraction Programming



- Asynchronous Task models with data migration

KAUST Cray Centre of Excellence



- Implicit usage of NVRAM 

- Redistribution programming abstraction

- Customised workflow management on Shaheen?

Acknowledgement and Call



Special thanks to CERL team and collaborators for UDJ

Funding bodies

HBP-PCP : UDJ development



Human Brain Project

Plan4res EU project : Data Model

<https://www.plan4res.eu/>



MCSA-ITN EXPERTISE : data redistribution approaches

www.msca-expertise.eu/



Centre for Doctoral Training in Data-Intensive Science (U. Cardiff)

Contribute to UDJ and Octopus (BSD license)

Contact adrian@cray.com for access pre-release