

# OpenMP<sup>®</sup>

## SC24 Booth Talk Series



### NUMA:

## Stay Close to Home with OpenMP

Ruud van der Pas

Senior Principal Software Engineer, Oracle Linux Engineering

# What is NUMA?

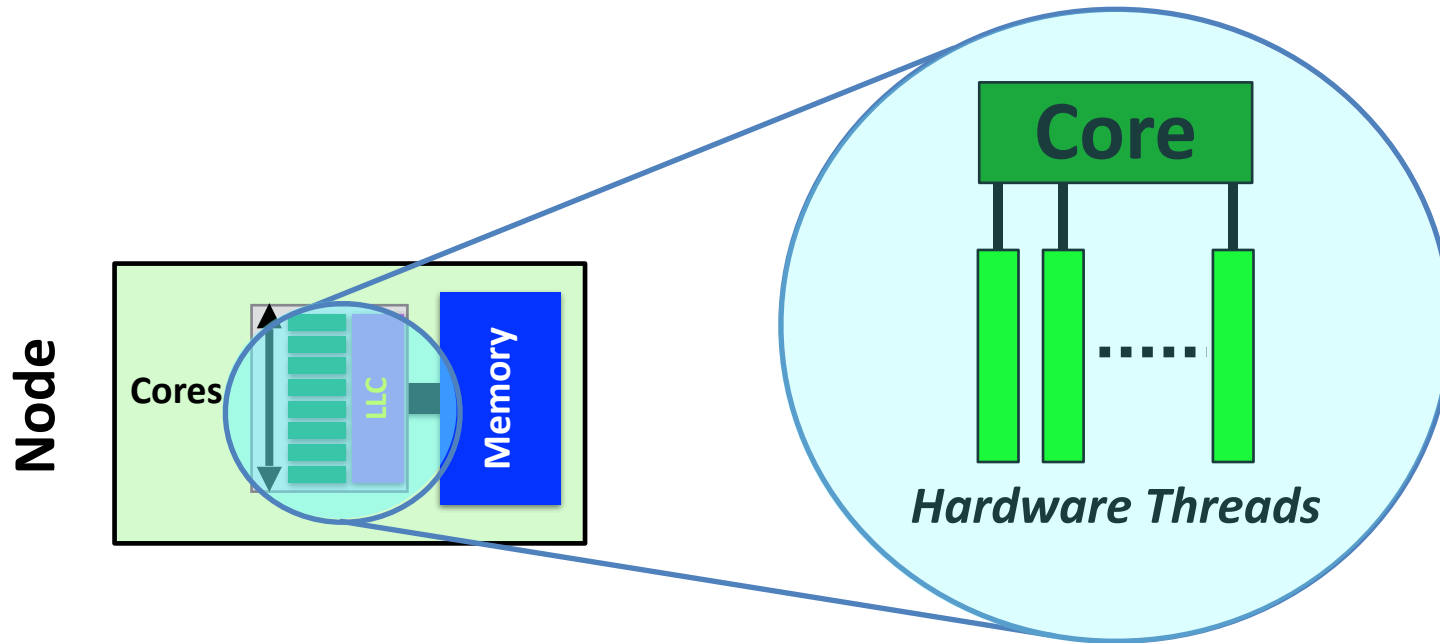


***Non-Uniform Memory Access (NUMA) used to be the realm of large servers only***

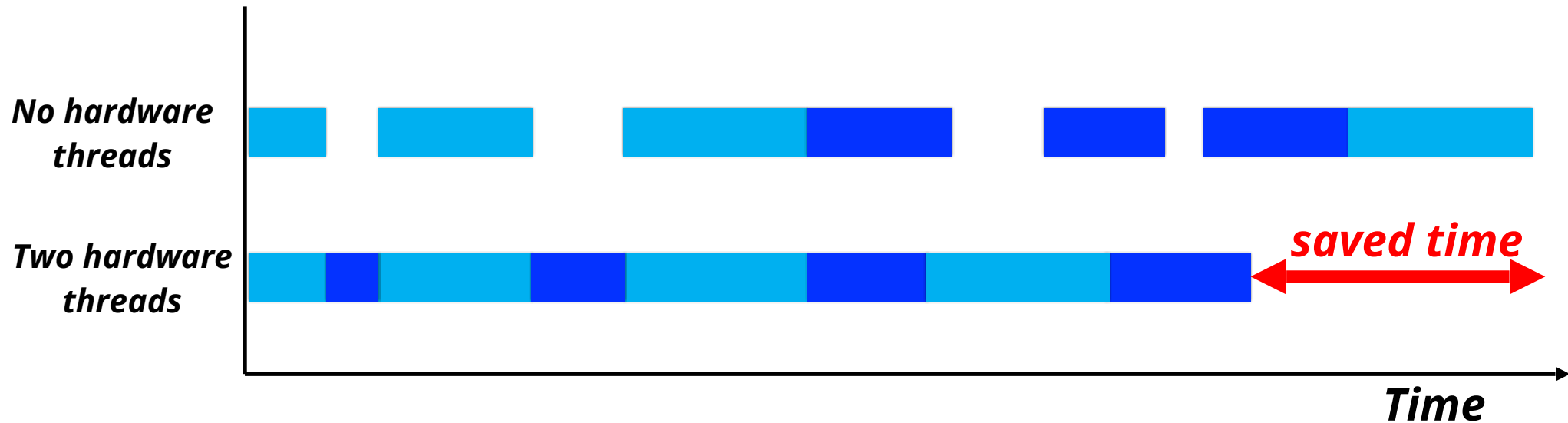
***This is no longer true and therefore **a concern to all*****

***The tricky thing is that “things just work”***

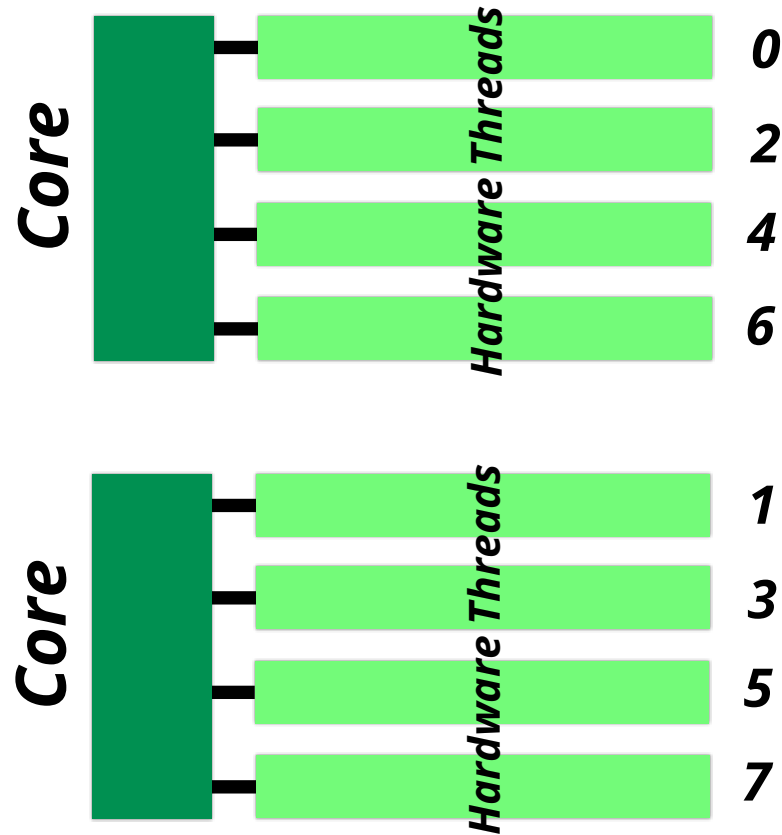
***But do you know how efficiently your code performs?***



# Intermezzo - How Hardware Threads Work

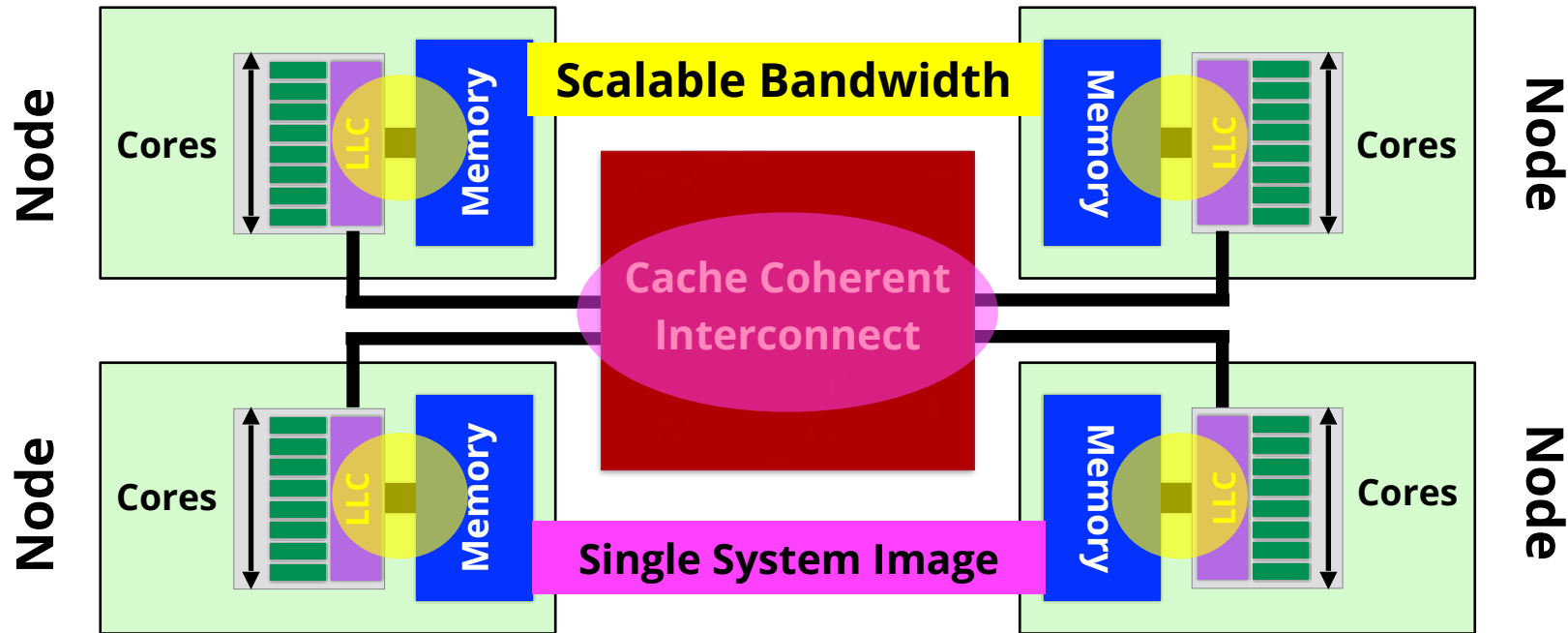


# Intermezzo - Hardware Thread IDs



# NUMA - The System Most of Us Use Today

*A Generic, but very Common and Contemporary NUMA System*



***Memory is physically distributed, but logically shared***

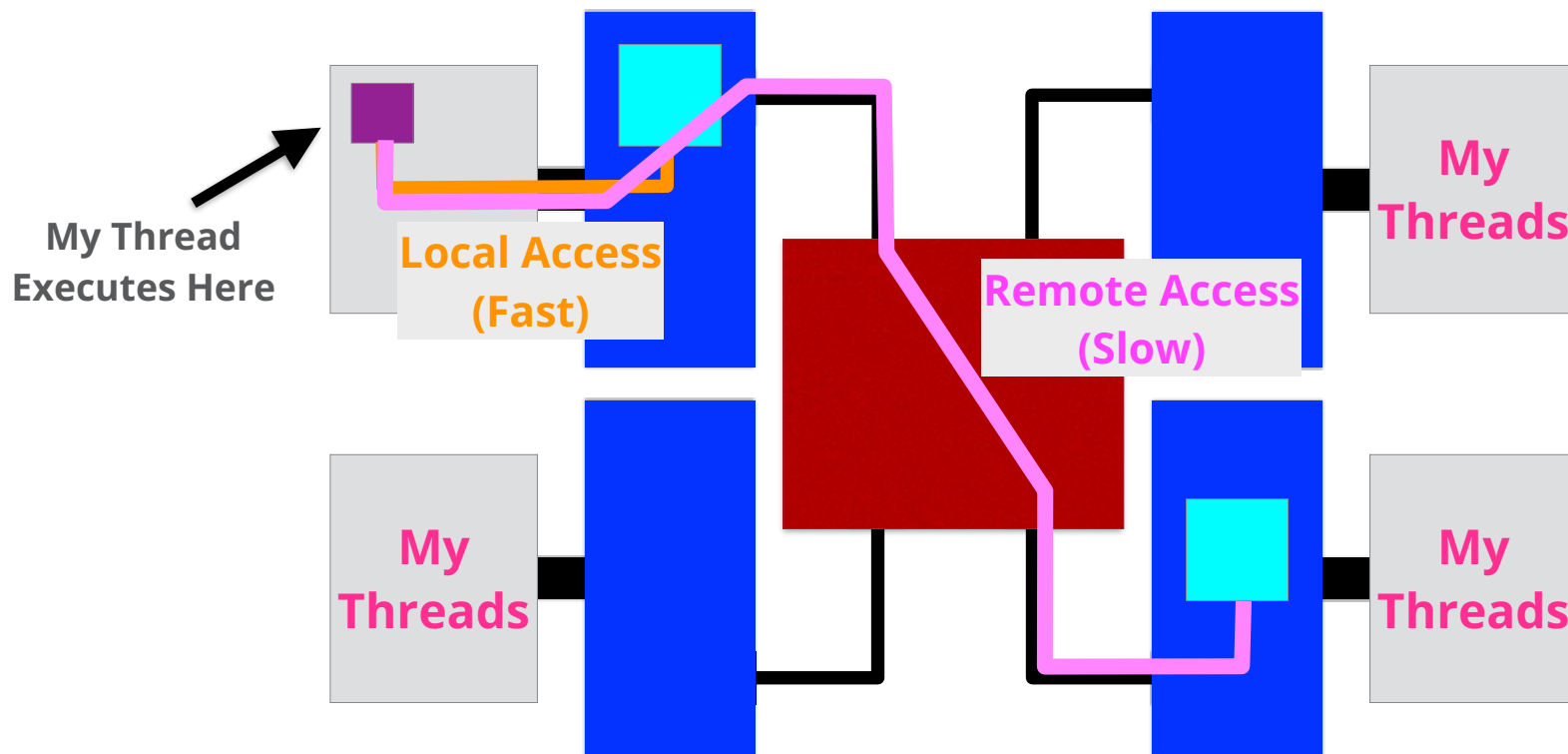
***Shared data is accessible to all threads***

***You don't know where the data is and it doesn't matter***

***Unless you care about performance ...***



# Local Versus Remote Access Times



# Understanding Your System



# The NUMA Information for a System

\$ **lscpu**

**8 cores/node**

**8 NUMA Nodes**

```
.....  
NUMA node0 CPU(s) : 0-7 , 64-71  
NUMA node1 CPU(s) : 8-15 , 72-79  
NUMA node2 CPU(s) : 16-23 , 80-87  
NUMA node3 CPU(s) : 24-31 , 88-95  
NUMA node4 CPU(s) : 32-39 , 96-103  
NUMA node5 CPU(s) : 40-47 , 104-111  
NUMA node6 CPU(s) : 48-55 , 112-119  
NUMA node7 CPU(s) : 56-63 , 120-127  
.....
```

\$ **numactl -H**

**node distances:**

node	0	1	2	3	4	5	6	7
0:	10	16	16	16	32	32	32	32
1:	16	10	16	16	32	32	32	32
2:	16	16	10	16	32	32	32	32
3:	16	16	16	10	32	32	32	32
4:	32	32	32	32	10	16	16	16
5:	32	32	32	32	16	10	16	16
6:	32	32	32	32	16	16	10	16
7:	32	32	32	32	16	16	16	10

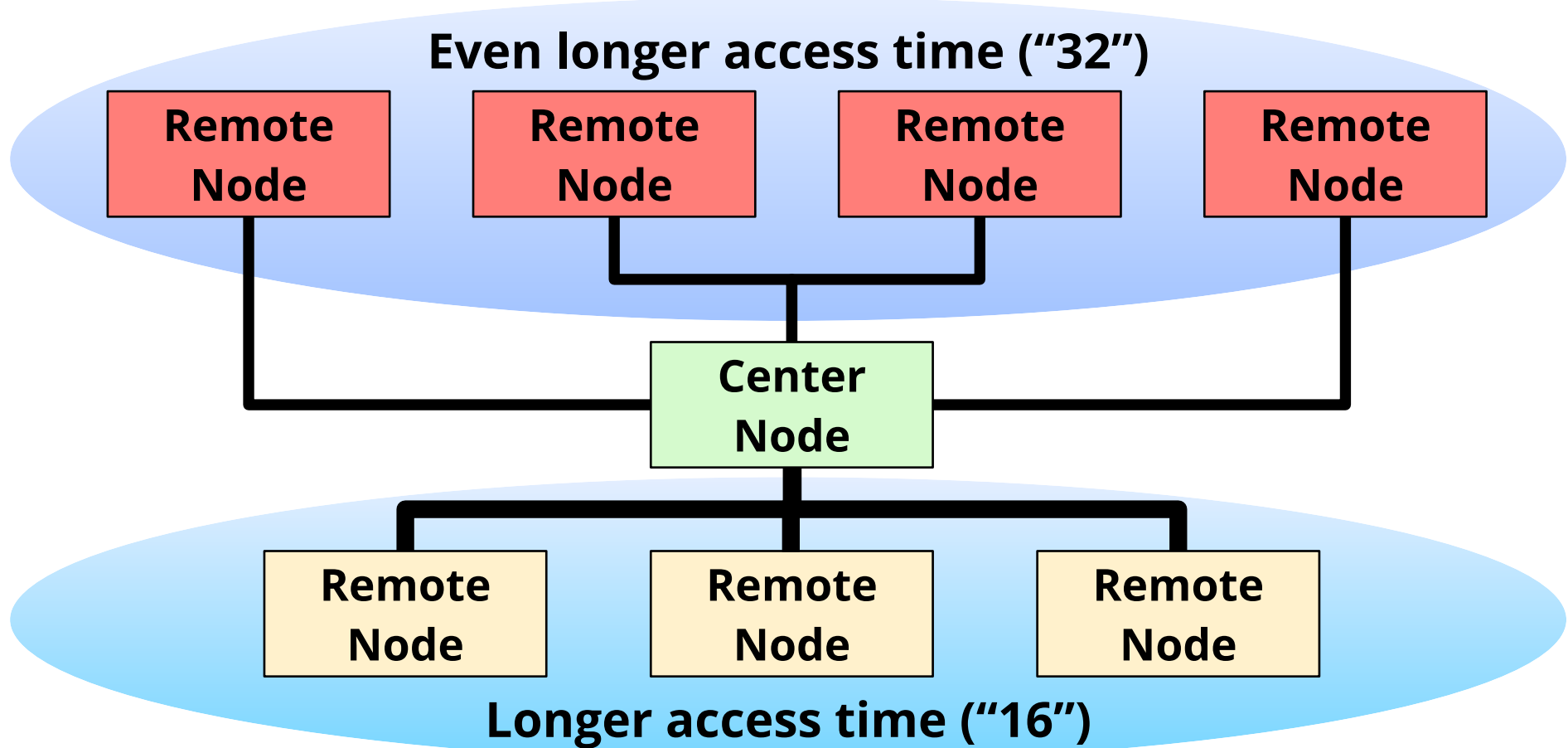
**2 columns => 2 hardware threads/core**

# The NUMA Structure of the System

<code>lscpu</code>	<b>There are 8 NUMA nodes</b>
<code>lscpu</code>	<b>There are 8 cores per node</b>
<code>lscpu</code>	<b>Each core has 2 hardware threads</b>
<code>numactl -H</code>	<b>Two levels of NUMA ("16" and "32")</b>

***This NUMA system has 64 cores and 128 hardware threads***

# The Abstract System Topology (numactl -H)



*The Goal: keep threads and their data close*

*The data location is fixed, a thread may be moved closer to it*

*Not the other way round, because that is more expensive*

*The **affinity constructs** in OpenMP control where threads run*

*This is a powerful feature, but it is up to you to get it right  
(in this context, "right" is not about correctness, but about the performance)*

# NUMA and Data Placement



*Question: where does data get allocated then?*

*The **First Touch Placement** policy allocates the data page in the memory closest to the thread accessing this page for the first time*

*This defines the fixed **Home Node** for the particular page*



*The goal of NUMA tuning*  
*Keep threads close to their Home Node*

***This policy is the default on Linux and other OSes***

***It is the right thing to do for a sequential application***

***But this may not work so well in a parallel application***

***What if a single thread initializes most, or all of the data?***

***Then, all the data ends up in the memory of a single node***

***This increases memory access times for certain threads  
(and may also cause congestion on the network)***

***Luckily, the solution is (often) surprisingly simple***

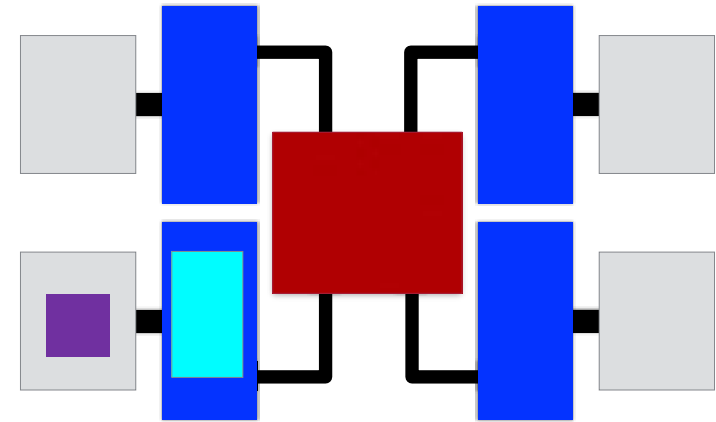
# A Sequential Initialization

```
for (int64_t i=0; i<n; i++)  
  a[i] = 0;
```

*One thread executes this loop*



*All of "a" is in a single node*



 = Thread  
 = Data

*Note: The allocation is on a virtual memory page basis*

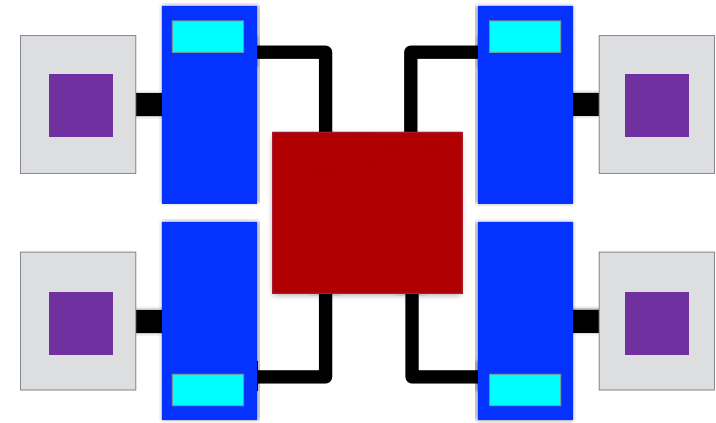
# Leverage the First Touch Placement Policy

```
#pragma omp parallel for schedule(static)  
for (int64_t i=0; i<n; i++)  
    a[i] = 0;
```

*Four threads execute this loop*



*The data is spread out*



 = Thread  
 = Data

*Note: The allocation is on a virtual memory page basis*

# OpenMP Support for NUMA Systems



*A place defines a set where a thread is allowed to run*

*The place list consists of a **set of places***

*The place list consists of either symbolic names, or a set of numbers*

- *An example of a symbolic name: **cores***
- *An example of a set of numbers: **{1, 3, 5, 7}** or **{1}, {3}, {5}, {7}***

  
**1 place**

  
**4 places**

## ***OMP\_PLACES***

***Defines the places where threads may run***

## ***OMP\_PROC\_BIND***

***Defines how threads map onto the OpenMP places  
(relevant if there are more places than threads)***



# Placement Targets Supported by OMP\_PLACES

<i>Keyword</i>	<i>Place definition</i>
<b><i>threads</i></b>	<b><i>A hardware thread</i></b>
<b><i>cores</i></b>	<b><i>A core</i></b>
<b><i>ll_caches</i></b>	<b><i>A set of cores that share the last level cache</i></b>
<b><i>numa_domains</i></b>	<b><i>A set of cores that share a memory with the same distance to that memory</i></b>
<b><i>sockets</i></b>	<b><i>A single socket</i></b>

***Note: The number of places may be restricted - For example: cores(4)***

*The **OMP\_PLACES** variable also supports hardware thread IDs*

*Places can be defined using any sequence of valid numbers*

*A compact set notation is supported as well*

*Notation: **{start:total:increment}***

*For example: **{0:4:2}** expands to **{0,2,4,6}***

*Threads are scheduled on the NUMA domains in the system:*

```
$ export OMP_PLACES=numa_domains
```

*Use 4 places using hardware thread IDs 0, 8, 16, and 24:*

```
$ export OMP_PLACES="{0},{8},{16},{24}"
```

```
$ export OMP_PLACES={0}:4:8
```

*Use 1 place using hardware thread IDs 0, 8, 16, and 24:*

```
$ export OMP_PLACES="{0,8,16,24}"
```

```
$ export OMP_PLACES={0:4:8}
```

*Use variable **OMP\_PROC\_BIND** to map threads onto places*

*The settings define the mapping of threads onto places*

*The following settings are supported:  
**true, false, primary, close, or spread***

*The definitions of close and spread are in terms of the place list*

*Threads are scheduled on the cores in the system:*

```
$ export OMP_PLACES=cores
```

*And they should be placed on cores as far away from each other as possible:*

```
$ export OMP_PROC_BIND=spread
```

# Remember this Example?

```
#pragma omp parallel for schedule(static)
for (int64_t i=0; i<n; i++)
  a[i] = 0;
```

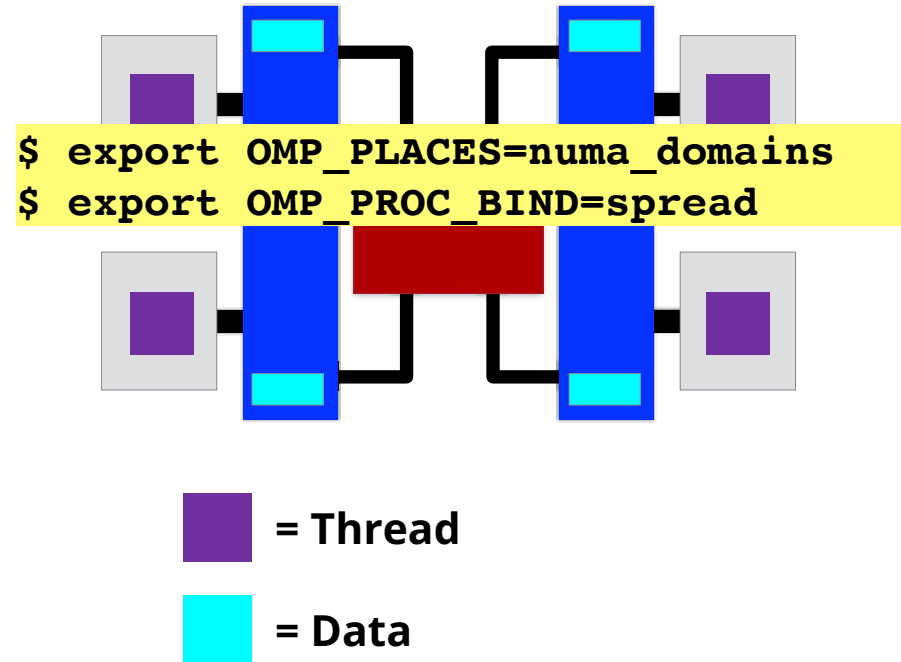
*Four threads execute this loop*



*Wishful Thinking*

*Data placement depends on  
where threads execute*

*Use the NUMA Controls*



# Important: Use the NUMA Diagnostics!

*It is very easy to make a mistake with the NUMA setup*

*Two very simple, but yet powerful features to assist you*

*Variable **OMP\_DISPLAY\_ENV** echoes the initial settings*

*Variable **OMP\_DISPLAY\_AFFINITY** prints information at run time*

*Highly recommended to use these diagnostic features!*

***Thank You And ... Stay Tuned!***

***Bad OpenMP  
Does Not Scale***

**Ruud van der Pas  
SC24 OpenMP Booth Talk**



The logo for OpenMP, featuring the word "Open" in a white sans-serif font with a horizontal line underneath, and "MP" in a larger, bold, white sans-serif font with a registered trademark symbol (®) to its upper right. The background is a dark teal and blue pixelated pattern.

# OpenMP<sup>®</sup>

## SC24 Booth Talk Series

**[openmp.org](https://openmp.org)**

OpenMP API specs, forum, reference guides, and more

**[link.openmp.org/sc24talks](https://link.openmp.org/sc24talks)**

OpenMP SC24 booth talk videos and presentations