# COLFAX
*Customized Solutions*

# PROGRAMMING AND OPTIMIZATION FOR INTEL® ARCHITECTURE

Hands-On Workshop (HOW) Series "Deep Dive"

Session 8

*Colfax International* — *colfaxresearch.com*

WELCOME

# DISCLAIMER

While best efforts have been used in preparing this training, Colfax International makes no representations or warranties of any kind and assumes no liabilities of any kind with respect to the accuracy or completeness of the contents and specifically disclaims any implied warranties of merchantability or fitness of use for a particular purpose. The publisher shall not be held liable or responsible to any person or entity with respect to any loss or incidental or consequential damages caused, or alleged to have been caused, directly or indirectly, by the information or programs contained herein. No warranty may be created or extended by sales representatives or written sales materials.

**WELCOME**

# COURSE ROADMAP

- ▷ Module I. Programming Models
  - 01. Intel Architecture and Modern Code
  - 02. Xeon Phi, Coprocessors, Omni-Path
- ▷ Module II. Expressing Parallelism
  - 03. Automatic vectorization
  - 04. Multi-threading with OpenMP
  - 05. Distributed Computing, MPI
- ▷ Module III. Performance Optimization
  - 06. Optimization Overview: N-body
  - 07. Scalar tuning, Vectorization
  - 08. Common Multi-threading Problems
  - 09. Multi-threading, Memory Aspect
  - 10. Access to Caches and Memory

COURSE ROADMAP

Course page:

colfaxresearch.com/how-series

- ▷ Slides
- ▷ Code
- ▷ Video
- ▷ Chat

More workshops:

colfaxresearch.com/training

colfaxresearch.com/how-series

# GET YOUR QUESTIONS ANSWERED: FORUMS

| | READ | WATCH | LEARN | FORUMS | CONNECT | JOIN |

**Forum**

## Colfax Cluster

Discussion of Colfax Cluster usage policies, troubleshooting.

## Developer Training, HOW Series

Questions about any of the Colfax trainings? Usage of training servers, experience with specific exercises, inquiries on what's inside, suggestions for future trainings - post them here.

## Performance Optimization and Parallelism

Discuss with Colfax Research and colleagues any topics related to computational science, parallel programming, performance optimization and code modernization.

## colfaxresearch.com/forum

# HANDS-ON EXERCISES AND REMOTE ACCESS

- ▷ All registrants receive an invitation from cluster@colfaxresearch.com
- ▷ Queue-based access to Intel Xeon E5, Intel Xeon Phi (KNC and KNL)
- ▷ Can access the cluster the entire 2 weeks of the workshop
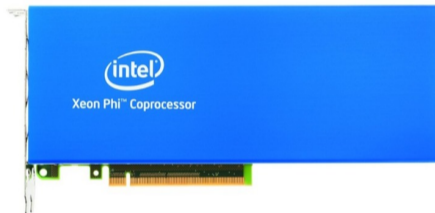
# §2. REFRESH

# PERFORMANCE OPTIMIZATION

# COMPUTING PLATFORMS

Intel Xeon
Processor

Intel Xeon Phi
Coprocessor, 1st generation

Intel Xeon Phi
Processor, 2nd generation*



Current: Broadwell
Upcoming: Skylake

* socket and coprocessor versions

Knights Corner (KNC)

Knights Landing (KNL)

Multi-Core Architecture

Intel Many Integrated Core (MIC) Architecture

**PERFORMANCE OPTIMIZATION**

# OPTIMIZATION AREAS



Scalar Tuning

Vectorization

Threading

Memory

Node-level

Cluster-level

Communication

PERFORMANCE OPTIMIZATION

## "HELLO WORLD" OPENMP PROGRAM

```cpp
#include <omp.h>
#include <cstdio>

int main(){
  // This code is executed by 1 thread
  const int nt=omp_get_max_threads();
  printf("OpenMP with %d threads\n", nt);

#pragma omp parallel
  { // This code is executed in parallel
    // by multiple threads
    printf("Hello World from thread %d\n",
                   omp_get_thread_num());
  }
}
```

▷ OpenMP = "Open Multi-Processing" = computing-oriented framework for shared-memory programming

▷ Threads – streams of instructions that share memory address space

▷ Distribute threads across CPU cores for parallel speedup

# §3. MULTI-THREADING: COMMON ISSUES
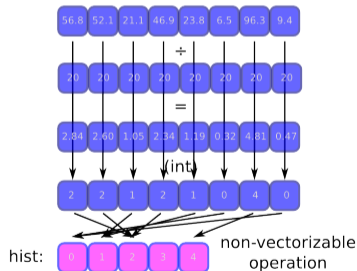
# TOO MUCH SYNCHRONIZATION

# EXAMPLE: BINNING PROBLEM

```cpp
void Histogram(
  // Ages, values from 0.0f to 100.0f:
  const float* age,
  // Size of array age, n=100000000:
  const int n,
  // Output: counts in groups:
  int* const hist,
  // Size of array hist, m=5:
  const int m,
  const float grpWidth) {
    for (int i = 0; i < n; i++) {
      const int j = int(age[i]/grpWidth);
      hist[j]++;
    }
}
```

▷ Vector dependence in
  `hist[j]++`

▷ Strip-mine or use
  conflict detection



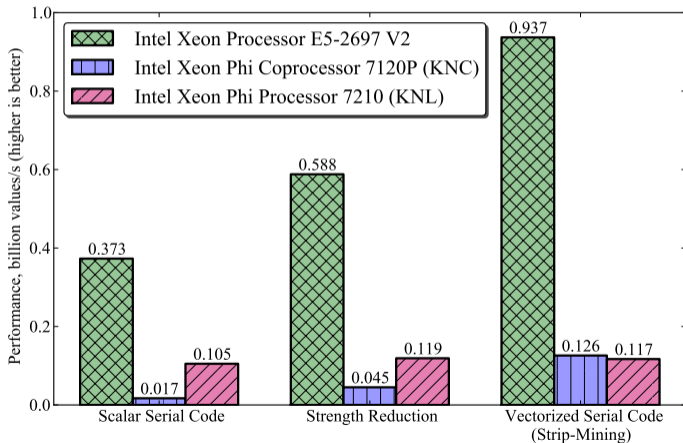**TOO MUCH SYNCHRONIZATION**

# THE SAME CALCULATION, STRIP-MINED, VECTORIZED

```
1   const float recGrpWidth = 1.0f/grpWidth; // precompute the reciprocal
2
3   for (int ii = 0; ii < n; ii += 16) { // strip-mining
4
5     int index[16]; // a block of indices
6     for (int i = ii; i < ii + 16; i++) // vectorizable
7       index[i-ii] = (int) ( age[i] * recGrpWidth ); // unit-stride access
8
9     for (int c = 0; c < 16; c++) // not vectorizable
10      hist[index[c]]++; // indirect access
11  }
```

## STRIP-MINING FOR VECTORIZATION

Vectorization improves performance.

More work is needed to take advantage of multiple cores.



**TOO MUCH SYNCHRONIZATION**

Incorrect solution: unprotected data races

```
#pragma omp parallel for schedule(guided)
for (int ii = 0; ii < n; ii += vecLen) {
  int index[vecLen] __attribute__((aligned(64)));
#pragma vector aligned
  for (int i = ii; i < ii + vecLen; i++)
    index[i-ii] = (int) ( age[i] * invGroupWidth );
  for (int c = 0; c < vecLen; c++)
    // Multiple threads will write into a single shared container
    // These data races lead to incorrect results!
    hist[index[c]]++;
}
```

**TOO MUCH SYNCHRONIZATION**

# HISTOGRAM CALCULATION EXAMPLE: ADDING THREAD PARALLELISM

Correct, but inefficient solution:

```
#pragma omp parallel for schedule(guided)
for (int ii = 0; ii < n; ii += vecLen) {
  int index[vecLen] __attribute__((aligned(64)));
#pragma vector aligned
  for (int i = ii; i < ii + vecLen; i++)
    index[i-ii] = (int) ( age[i] * invGroupWidth );
  for (int c = 0; c < vecLen; c++)
    // Protect the ++ operation with the atomic mutex (inefficient!)
#pragma omp critical
    { hist[index[c]]++; }
}
```

**TOO MUCH SYNCHRONIZATION**
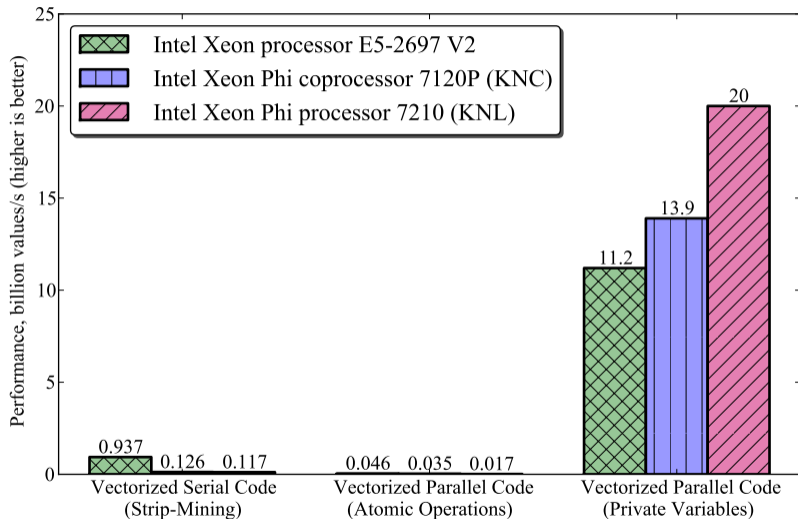
# HISTOGRAM CALCULATION EXAMPLE: ADDING THREAD PARALLELISM

Correct, but inefficient solution:

```
#pragma omp parallel for schedule(guided)
for (int ii = 0; ii < n; ii += vecLen) {
  int index[vecLen] __attribute__((aligned(64)));
#pragma vector aligned
  for (int i = ii; i < ii + vecLen; i++)
    index[i-ii] = (int) ( age[i] * invGroupWidth );
  for (int c = 0; c < vecLen; c++)
    // Protect the ++ operation with the atomic mutex (inefficient!)
#pragma omp atomic
    hist[index[c]]++;
}
```

**TOO MUCH SYNCHRONIZATION**

# CORRECT AND EFFICIENT SOLUTION WITH REDUCTION

```
1   #pragma omp parallel
2   {
3     int hist_priv[m]; // Better idea: thread-private storage
4     hist_priv[:] = 0;
5     int index[vecLen] __attribute__((aligned(64)));
6   #pragma omp for schedule(guided)
7     for (int ii = 0; ii < n; ii += vecLen) {
8   #pragma vector aligned
9       for (int i = ii; i < ii + vecLen; i++)
10        index[i-ii] = (int) ( age[i] * invGroupWidth );
11      for (int c = 0; c < vecLen; c++)
12        hist_priv[index[c]]++;
13    }
14    for (int c = 0; c < m; c++) {
15  #pragma omp atomic
16      hist[c] += hist_priv[c];
17  } }
```
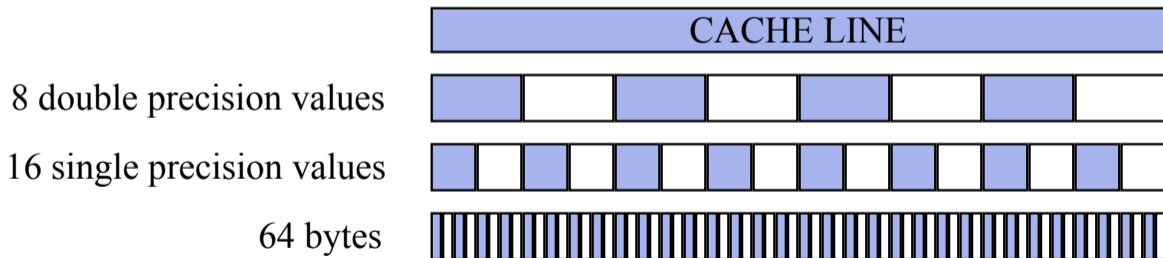
**TOO MUCH SYNCHRONIZATION**

# USING REDUCTION INSTEAD OF SYNCHRONIZATION

**FALSE SHARING**

# CACHE LINES
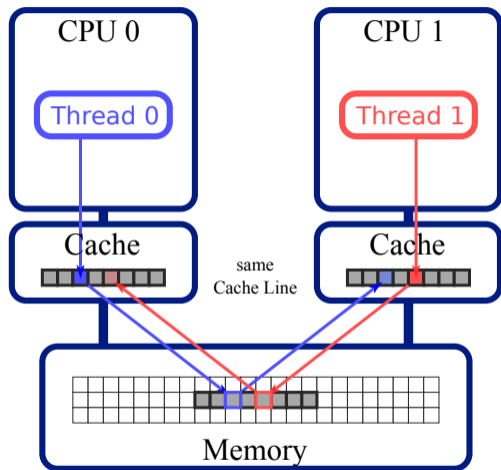
▷ Minimal block of data transferred between memory and cache

▷ 64 bytes long in Intel Architecture

▷ Aligned on 64-byte boundaries in memory



CACHE LINE

8 double precision values

16 single precision values

64 bytes

**FALSE SHARING**

# FALSE SHARING. DATA PADDING AND PRIVATE VARIABLES



CPU 0

Thread 0
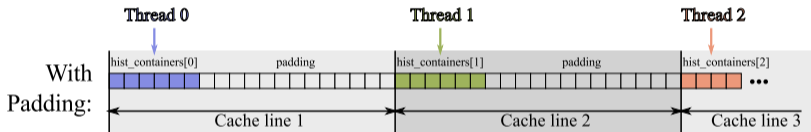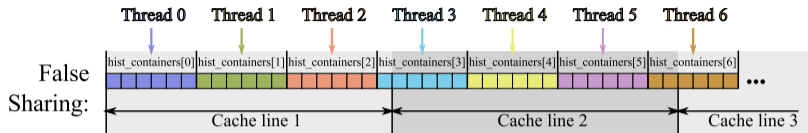
CPU 1

Thread 1

Cache

same
Cache Line

Cache

Memory

- ▷ Occurs when 2 or more threads acess the same cache line, and at least one of the accesses is for writing
- ▷ Caused by *coherent caches*
- ▷ Cache line is 64-byte wide (in modern Intel architectures)

FALSE SHARING

# FALSE SHARING. DATA PADDING AND PRIVATE VARIABLES

```
1  const int m = 5;
2  int hist_thr[nThreads][m];
3  #pragma omp parallel for
4  for (int ii = 0; ii < n; ii += vecLen) {
5    // ...
6    // False sharing occurs here
7    for (int c = 0; c < vecLen; c++)
8      hist_thr[iThread][index[c]]++;
9  }
10 // Reducing results from all threads to the common histogram hist
11 for (int iThread = 0; iThread < nThreads; iThread++)
12   hist[0:m] += hist_thr[iThread][0:m];
```

 ▷ The value of m=5 is small
 ▷ Array elements hist_thr[0][:] are within m*sizeof(int)=20
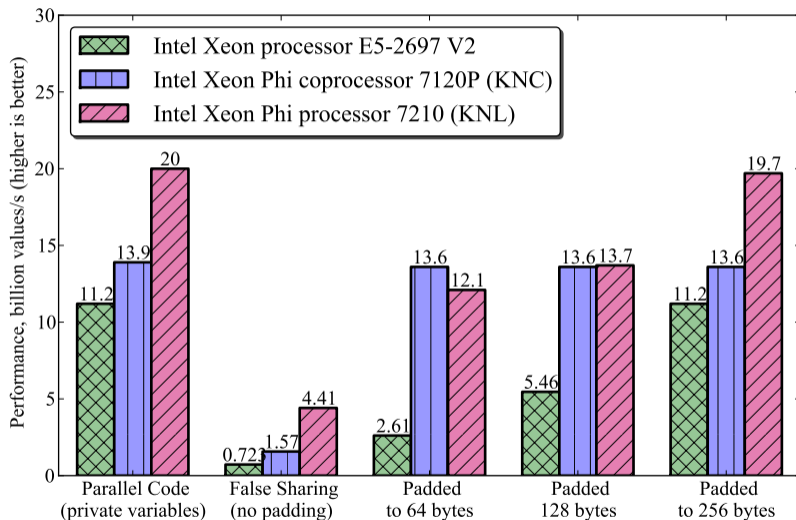   bytes of array elements hist_thr[1][:]

# PADDING TO AVOID FALSE SHARING



```
// Padding to avoid sharing a cache line between threads
const int paddingBytes = 64;
const int paddingElements = paddingBytes / sizeof(int);
const int mPadded = m + (paddingElements-m%paddingElements);
int hist_containers[nThreads][mPadded]; // New container
```

FALSE SHARING
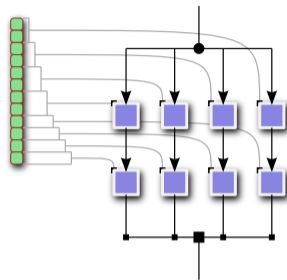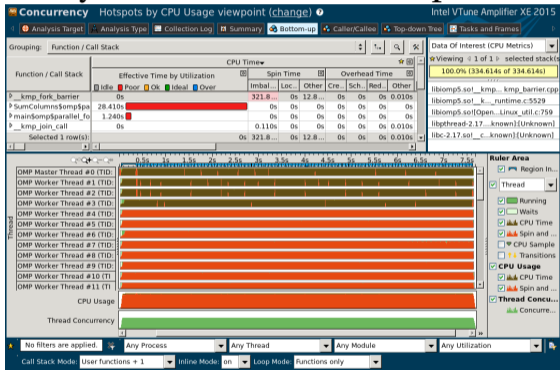
# PADDING TO AVOID FALSE SHARING

# INSUFFICIENT PARALLELISM

# INSUFFICIENT PARALLELISM

## Analysis in Intel VTune Amplifier XE





▷ Occurs when there are not enough iterations or parallel work-items exposed to the parallel loop in OpenMP.

# EXAMPLE: DEALING WITH INSUFFICIENT PARALLELISM

$$S_i = \sum_{j=0}^{n} M_{ij}, \ i = 0 \dots m. \tag{1}$$

▷ m=4 is small, smaller than the number of threads in the system

▷ $n \approx 10^8$ is large enough so that matrix does not fit into cache

```
1  void sum_unoptimized(const int m, const int n, long* M, long* s){
2  #pragma omp parallel for
3    for (int i=0; i<m; i++) { // m=4
4      long total=0;
5  #pragma vector aligned
6      for (int j=0; j<n; j++) // n=100000000
7        total+=M[i*n+j];
8    s[i]=total; }}
```

**INSUFFICIENT PARALLELISM**

# DOES NOT WORK: PARALLELIZING INNER LOOP

Inner loop has more iterations, parallelize there?

```
1  void SumParallelInnerLoop(const int m, const int n, long* M, long* s){
2    for (int i = 0; i < m; i++) { // m=4
3      long total = 0;
4  #pragma omp parallel for reduction(+: total)
5      for (int j = 0; j < n; j++) { // n=100000000
6        total += M[i*n + j];
7      }
8      s[i] = total;
9    }
10 }
```

Does not work well: code must spawn and stop threads many times;
OpenMP does not see the entire parallel region.

**INSUFFICIENT PARALLELISM**

# LOOP COLLAPSE: PRINCIPLE

Idea: combine iterations spaces of the inner loop and the outer loop.

```
1   #pragma omp parallel for collapse(2)
2     for (int i = 0; i < m; i++)
3       for (int j = 0; j < n; j++) {
4         // ...
5         // ...
6       }
```

```
1   #pragma omp parallel for
2     for (int c = 0; c < m*n; c++)  {
3       i = c / n;
4       j = c % n;
5       // ...
6     }
```

# DOES NOT WORK, BUT CORRECT DIRECTION: LOOP COLLAPSE
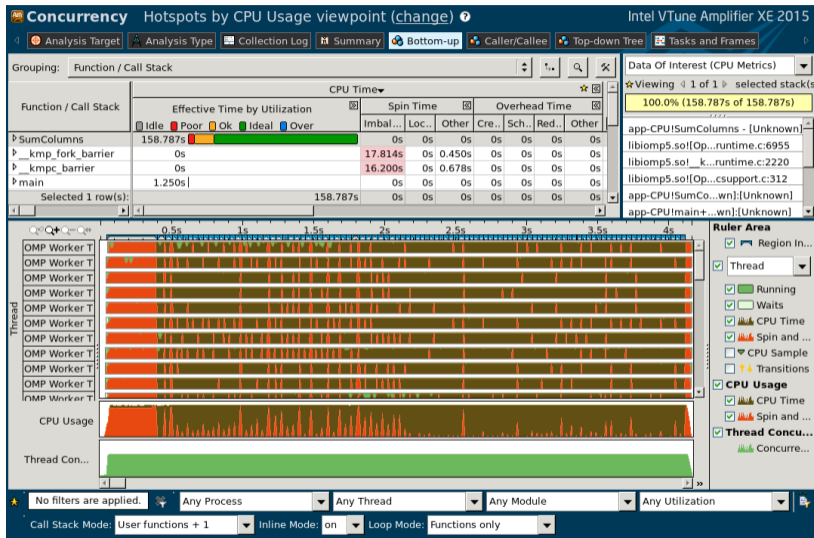
Loop collapse applied to the wide short matrix example:

```cpp
void SumCollapse(const int m, const int n, long* M, long* s){
  s[:]=0;
#pragma omp parallel
  { // Each thread will have a private container
    long total[m];  total[:] = 0;
#pragma omp for collapse(2)
    for (int i = 0; i < m; i++) // m=4
      for (int j = 0; j < n; j++) // n=100000000
        total[i] += M[i*n + j];
    for (int i = 0; i < m; i++)
#pragma omp atomic
      s[i]=total[i];
} }
```
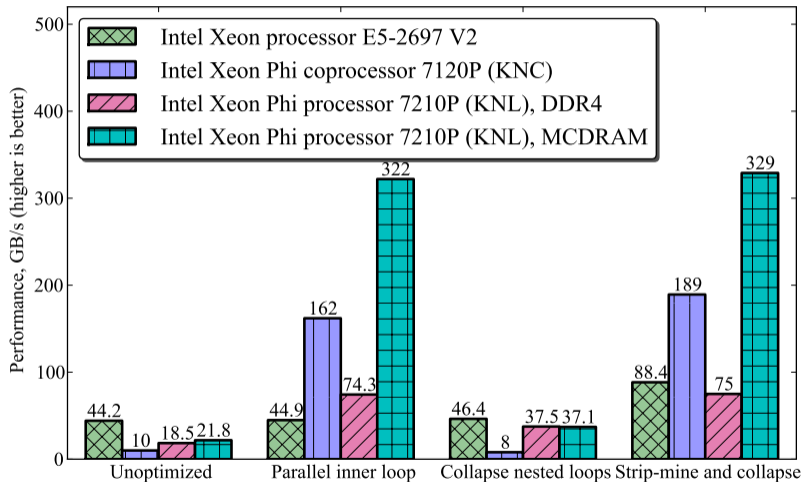
Does not work: automatic vectorization fails.

# EXPOSING PARALLELISM: STRIP-MINING AND LOOP COLLAPSE

```
1  void sum_stripmine(const int m, const int n, long* M, long* s){
2    const int STRIP=1024;
3    assert(n%STRIP==0);
4    s[0:m]=0;
5  #pragma omp parallel
6    {
7      long total[m];   total[0:m]=0;
8  #pragma omp for collapse(2) schedule(guided)
9      for (int i=0; i<m; i++)
10       for (int jj=0; jj<n; jj+=STRIP)
11 #pragma vector aligned
12         for (int j=jj; j<jj+STRIP; j++)
13           total[i]+=M[i*n+j];
14     for (int i=0; i<m; i++)          // Reduction
15 #pragma omp atomic
16         s[i]+=total[i];
17 }   }
```
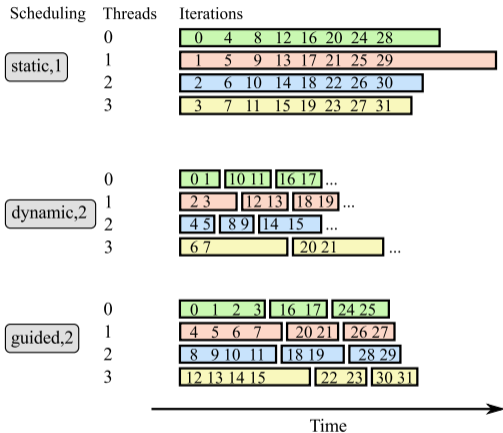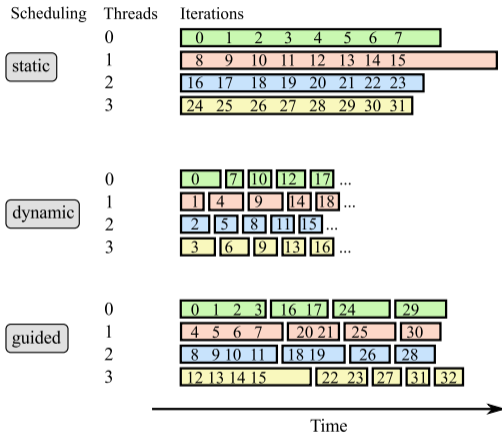
# EXPOSING PARALLELISM: STRIP-MINING AND LOOP COLLAPSE



**INSUFFICIENT PARALLELISM**

# DEALING WITH INSUFFICIENT PARALLELISM

# LOAD IMBALANCE

# LOOP SCHEDULING MODES IN OPENMP

**LOAD IMBALANCE**

# CONTROL OF SCHEDULING MODES

To set scheduling for a particular loop in code (example):

```
1  #pragma omp parallel for schedule(dynamic,4)
2    // ...
```

To set scheduling for the entire application at run time (example):

```
1  #pragma omp parallel for schedule(runtime)
2    // ...
```

```
vega@lyra% export OMP_SCHEDULE=dynamic,4
vega@lyra% ./run-my-app
```
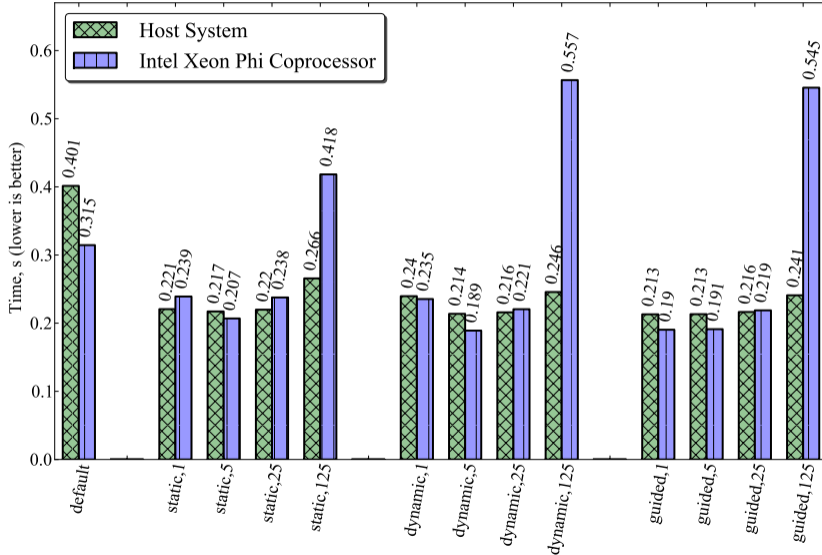
# ITERATIVE JACOBI SOLVER

```
int IterativeSolver(int n, double* M, double* b, double* x, double minAccuracy){
  double accuracy; int iters=0; double bTrial[n] __attribute__((aligned(64)));
  x[0:n] = 0.0; // Initial guess
  do { iters++; // The Jacobi method - iterate until convergence
      for (int i = 0; i < n; i++) {
        double c = 0.0;
#pragma vector aligned
        for (int j = 0; j < n; j++) c += M[i*n+j]*x[j]; // Iterate
        x[i] = x[i] + (b[i] - c)/M[i*n+i]; }
      bTrial[:] = 0.0; // Verification
      for (int i = 0; i < n; i++)
#pragma vector aligned
        for (int j = 0; j < n; j++) bTrial[i] += M[i*n+j]*x[j];
      accuracy = RelativeNormOfDifference(n, b, bTrial); // Check convergence
    } while (accuracy > minAccuracy); // Must achieve the requested accuracy
  return iters; }
```

# AN ITERATIVE JACOBI SOLVER

```
1   #pragma omp parallel for
2       for (int c = 0; c < nBVectors; c++)
3           IterativeSolver(n, M, &b[c*n], &x[c*n], accuracy[c]);
```

**LOAD IMBALANCE**

# AN ITERATIVE JACOBI SOLVER WITH DYNAMIC SCHEDULING

```
#pragma omp parallel for schedule(dynamic,4)
    for (int c = 0; c < nBVectors; c++)
        IterativeSolver(n, M, &b[c*n], &x[c*n], accuracy[c]);
```

LOAD IMBALANCE

# PERFORMANCE OF ITERATIVE JACOBI SOLVER

# §4. REVIEW AND WHAT'S NEXT

## SUMMARY

This session:

1. Synchronization is necessary to resolve data races
2. Mutexes must be moved out of innermost loops
3. False sharing can be resolved with padding
4. Loop collapse can help to expose parallelism
5. Strip-mining to make vectorization co-exist with threading
6. Trade-off between load balance and low scheduling overhead

Next session: optimization of thread affinity, NUMA locality, nested parallelism.

https://colfaxresearch.com/