**Imperial College London**

# Hybrid Parallelization of SU2
*A Comprehensive Introduction*

Pedro Gomes, Rafael Palacios

1st Annual SU2 Conference, 10-12 June 2020
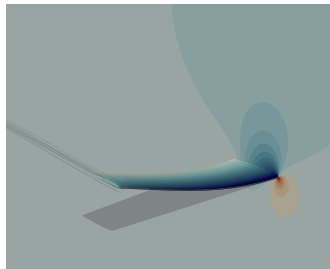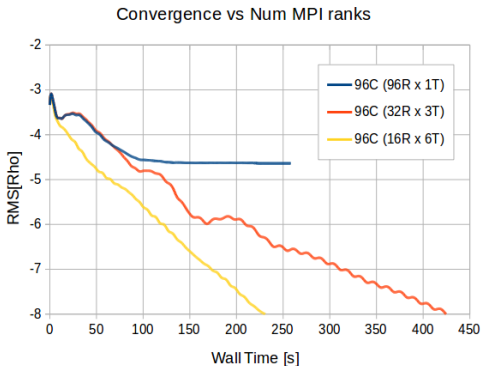
# Contents

- Motivation
- The hybrid parallel model
- OpenMP, an overview
- Challenges (and solutions)
- Implementation overview
- Concluding remarks

# Motivation

Faster and more robust code, that scales better.

▶ Algorithms work better;
▶ Dynamic load balancing;
▶ Reduced communication overhead;



Convergence vs Num MPI ranks

— 96C (96R x 1T)
— 96C (32R x 3T)
— 96C (16R x 6T)



Objective, fast medium scale optimizations.

# Motivation

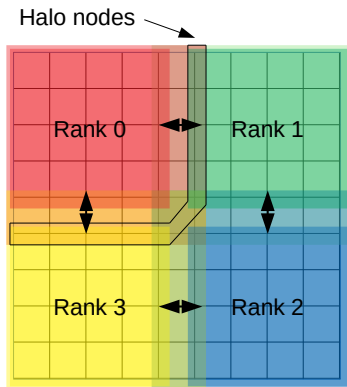Q: What is covered by the implementation?
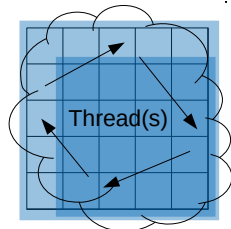A: Primal and forward AD compressible URANS FSI (and subsets).
Q: How do I use it?

```
1 ./meson.py ... -Dwith-omp=true ...
2
3 # auto number of threads/rank
4 SU2_CFD config.cfg
5 # 8 threads total
6 mpirun -n 2 --bind-to numa SU2_CFD -t 4 ...
7 # never --bind-to core
8 # mileage may vary, e.g. 2*4 != 4*2
9
10 # Useful environment variables:
11 # overrides default threads/rank
12 export OMP_NUM_THREADS=4
13 # better performance on some systems
14 export OMP_WAIT_POLICY=ACTIVE
```

# The hybrid parallel model

Domain decomposition for MPI (static) vs the (possibly) dynamic movement of threads within sub-domains.
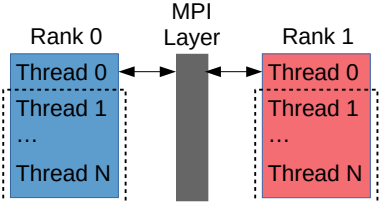


Halo nodes

Rank 0    Rank 1

Rank 3    Rank 2

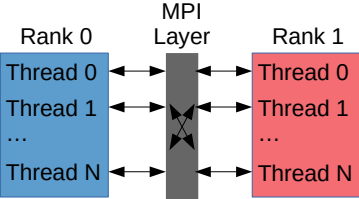Data from other ranks is obtained **indirectly** via messages.

Thread(s)

Threads within each partition can **directly** access any part of it.

# The hybrid parallel model

The threads can interact with MPI in different ways, currently communications are **funneled** (multiple is WIP).



**Funneled** communication
(only the main thread uses MPI)

**Multiple** communication
(any thread at any time can use MPI)

# OpenMP, an overview

An API that provides a simple and flexible interface (mostly in the form of pragmas) to develop portable parallel applications.

```cpp
const int N = 1024; // a shared variable
// start some threads
#pragma omp parallel
{
  int i; // a private variable
  // distribute loop indexes over threads
  #pragma omp for schedule(dynamic,32)
  for(i=0; i<N; ++i)
    myThreadSafeFunction(i);
}
```

We want to use this API to distribute the work (loops) in each MPI partition over its threads.

# OpenMP, an overview

How do threads "communicate" between themselves?

```cpp
// a function called by multiple threads
// x,y shared variables
void axpy(int N, double a, const double* x, double* y)
{
  // here there is no guarantee that the threads
  // have a consistent view of the arrays
  #pragma omp barrier
  // now there is
  #pragma omp for simd schedule(static,1024)
  for(int i=0; i<N; ++i) y[i] += a*x[i];
  // there are implicit barriers after most
  // worksharing directives
}
```

Identifying and counting threads:

```cpp
omp_get_num_threads() <=> "size"
omp_get_thread_num() <=> "rank"
```

# OpenMP, an overview

Shared vs private variables

```
1  #pragma omp parallel num_threads(4)
2  {
3    // each thread has its own stack -> private
4    double x[64] = {1.0};
5    // the heap is shared, but we made 4 y's...
6    vector<double> y(64,2.0);
7    // this will not do what we want...
8    axpy(64, 0.5, x, y.data());
9  }
```
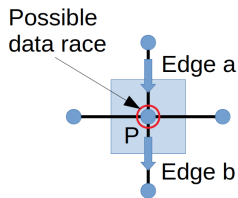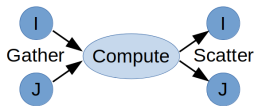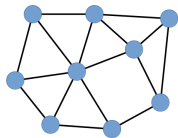
So do we need to declare/allocate everything outside parallel
regions? Yes, and no.

```
1  vector<double> y;
2  #pragma omp parallel num_threads(4)
3  {
4    #pragma omp master
5    y.resize(64,2.0); // only one thread allocates
6    #pragma omp barrier
7    ...
```

# Challenges

### Data races
When multiple threads simultaneously modify the same memory location (in an unregulated manner).



Builtin (OpenMP) solutions:

- Atomic operations;
- Critical directive;
- Locks;

Algorithmic solutions:

- Coloring / Partitioning;
- Scatter to Gather transformations;

# Challenges

Builtin solutions:

```cpp
// atomics are good for reduction operations
auto mySum = f(); // a private variable
#pragma omp atomic
ourSum += mySum; // safe update of shared variable

// critical for global resources
#pragma omp critical
cout << mySum << endl; // serialize output (unordered)

// locks for specific resources
const auto j = selectResource(omp_thread_num());
omp_set_lock(fileLocks[j]);
files[j] << mySum << endl;
omp_unset_lock(fileLocks[j]);
```
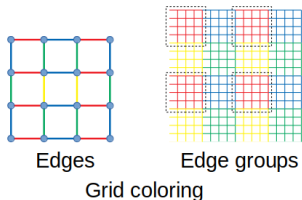
Pros: Small modifications to existing algorithms.
Cons: Overhead, poor scaling for resources used intensively.

# Challenges

**Partitioning**, re-partition the MPI sub-domains, not what we want.
**Coloring**, create non-intersecting sets of entities (data race free).



Edges        Edge groups
Grid coloring

```cpp
for(auto c : EdgeColoring) {
  SU2_OMP_FOR_DYN(groupSize)
  for(int k=0; k<c.size; ++k) {
    auto idx = c.indices[k];
    ...
```

Pros: Load balancing via dynamic scheduling.
Cons: Reduced locality, parallel inefficiency (not enough work chunks for all threads).

**This is our first choice for residual loops (use option EDGE_COLORING_GROUP_SIZE [512] to tune it).**

# Challenges

Scatter to Gather transformations:

```
1 for(edge : Edges) {
2   // gather
3   auto f = y[iPt]+y[jPt];
4   // scatter
5   x[iPt] += f;
6   x[jPt] += f;
7 }
```

```
1 for(iPt : Points) {
2   // gather
3   for(jPt : neighbors(iPt))
4     x[iPt]+=y[iPt]+y[jPt];
5 }
```

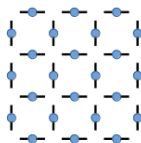Pros: Embarrassingly parallel code, if the FLOP/BYTE ratio is O(1) the code may perform better.

Cons: Needs adjacency matrix, 2x slower if FLOP/BYTE >> 1.

**This is what was done for preprocessing-type routines (gradients, limiters, sensors, etc.).**

# Challenges

Scatter to Gather transformations (two loop approach):

```
1 for(edge : Edges) {
2   // gather
3   f[edge] = y[iPt]+y[jPt];
4 }
5 for(iPt : Points) {
6   // gather
7   for(edge : Edges(iPt))
8     x[iPt] += f[edge];
9 }
```



Step 1- Fluxes    Step 2- Reduce

Reduction strategy

Pros: Approximately the same number of flops.
Cons: Extra storage needed for temporary variables, reduction loop has very low FLOP/BYTE ratio.

**When edge coloring fails, SU2 falls back to this approach. About 20% slower worst case.** A hybrid approach would probably be optimal.

# Implementation overview

All pragmas and functions used throughout the code are wrapped in omp_structure.hpp, this allows disabling everything when -Dwith-omp=false (default).

```
1 #define SU2_OMP_SIMD SU2_OMP(simd)
2 #define SU2_OMP_MASTER SU2_OMP(master)
3 #define SU2_OMP_ATOMIC SU2_OMP(atomic)
4 #define SU2_OMP_BARRIER SU2_OMP(barrier)
5 #define SU2_OMP_CRITICAL SU2_OMP(critical)
6 #define SU2_OMP_PARALLEL SU2_OMP(parallel)
7 ...
```

SU2_OMP_ $\approx$ SU2_MPI::

# Implementation overview

Threads are started once per iteration (and per integration) in
CIntegration (single or multi grid), output, to screen and file, is
not multi-threaded.

```
1 /*--- Start an OpenMP parallel region covering the
      entire MG iteration , if the solver supports it.
      ---*/
2 SU2_OMP_PARALLEL_(if(solver_container[iZone][iInst][
      MESH_0][Solver_Position]->GetHasHybridParallel()))
3 {
4 ...
```

All routines that are part of one iteration must be thread-safe (i.e.
no unguarded writes to colliding memory locations).

# Implementation overview

The "one numerics per thread" paradigm:
Numerics are shared objects (instantiated outside parallel regions)
with mutable state and thus cannot be used by multiple threads.

```
1 /*--- Pick one numerics object per thread. ---*/
2 CNumerics* numerics = numerics_container[CONV_TERM +
     omp_get_thread_num()*MAX_TERMS];
```

This kind of temporary variable must also be avoided:

```
1 class CSolver {
2   su2double *Solution,  /*!< \brief Auxiliary ... */
3   *Solution_i,          /*!< \brief Auxiliary ... */
4   ...
```

What about using one per thread too? Bad idea due to false
sharing.

# Implementation overview

Grid coloring or fallback strategies are setup in solver constructors, then in residual loops:

```
1   if (ReducerStrategy) {
2     EdgeFluxes.SetBlock(iEdge, residual);
3     Jacobian.SetBlocks(iEdge, ...
4   }
5   else {
6     LinSysRes.AddBlock(iPoint, residual);
7     LinSysRes.SubtractBlock(jPoint, residual);
8     Jacobian.UpdateBlocks(iEdge, iPoint, jPoint, ...
9   }
10  ...
11  } // end color loop
12
13  if (ReducerStrategy) {
14    SumEdgeFluxes();
15    Jacobian.SetDiagonalAsColumnSum();
16  }
```

# Implementation overview

## Other tricky areas

To go around a barrier, we need two barriers:

```
1  if(condition) {
2    SU2_OMP_BARRIER // wait for all threads to enter
3    SU2_OMP_MASTER {condition = f();} // before updating
4    SU2_OMP_BARRIER // or some might skip this barrier
5  }
```

Not so obvious deadlocks:

```
1  axpy(N,a,x,y); // this is fine, works in serial
2  SU2_OMP_PARALLEL {
3    axpy(N,a,x,y); // this is fine, works in parallel
4    SU2_OMP_MASTER {axpy(N,a,x,y);} // deadlock
5    // other threads missed the barrier inside axpy
6  }
```

## Implementation overview

AD-compatible funneled reductions (would be simpler with multiple communication):

```
1  su2double minElem(int N, const su2double* x) {
2    static su2double ourMin; // global var!!
3    SU2_OMP_BARRIER // consistent view of x
4    SU2_OMP_MASTER {ourMin = 1e30;} // init global
5    su2double myMin = 1e30; // init local
6    SU2_OMP_FOR_STAT(256)
7    for(int i=0; i<N; ++i) myMin = min(myMin, x[i]);
8    SU2_OMP_CRITICAL // serialize update of shared var
9    ourMin = min(ourMin, myMin);
10   SU2_OMP_BARRIER // wait for all updates
11   SU2_OMP_MASTER { // master communicates
12     myMin = ourMin;
13     SU2_MPI::Allreduce(&myMin, &ourMin,...
14   }
15   SU2_OMP_BARRIER // consistent view of ourMin
16   return ourMin; // same on all threads and ranks
17 }
```

# Concluding remarks

- ▶ Small set of OpenMP features used (also for eventual compatibility with reverse AD);
- ▶ That are still enough to improve scalability;
- ▶ Somethings require a bit more care, but essentially just be careful when writing to shared variables;
- ▶ Still lots of WIP, the solvers currently covered are a test bed for hybrid parallel strategies, it will take some time to cover everything.