



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación

OpenMP tasking:

Extensions and optimizations for performance, predictability and resilience

Sara Royuela
sara.royuela@bsc.es



PPC
Predictable
Parallel
Computing

OpenMP User's Monthly Telecon

October 28, 2022

Predictable Parallel Computing in OpenMP



Dr. Eduardo Quiñones
Team leader



Chenle Yu
PhD candidate



Adrian Munera
PhD candidate



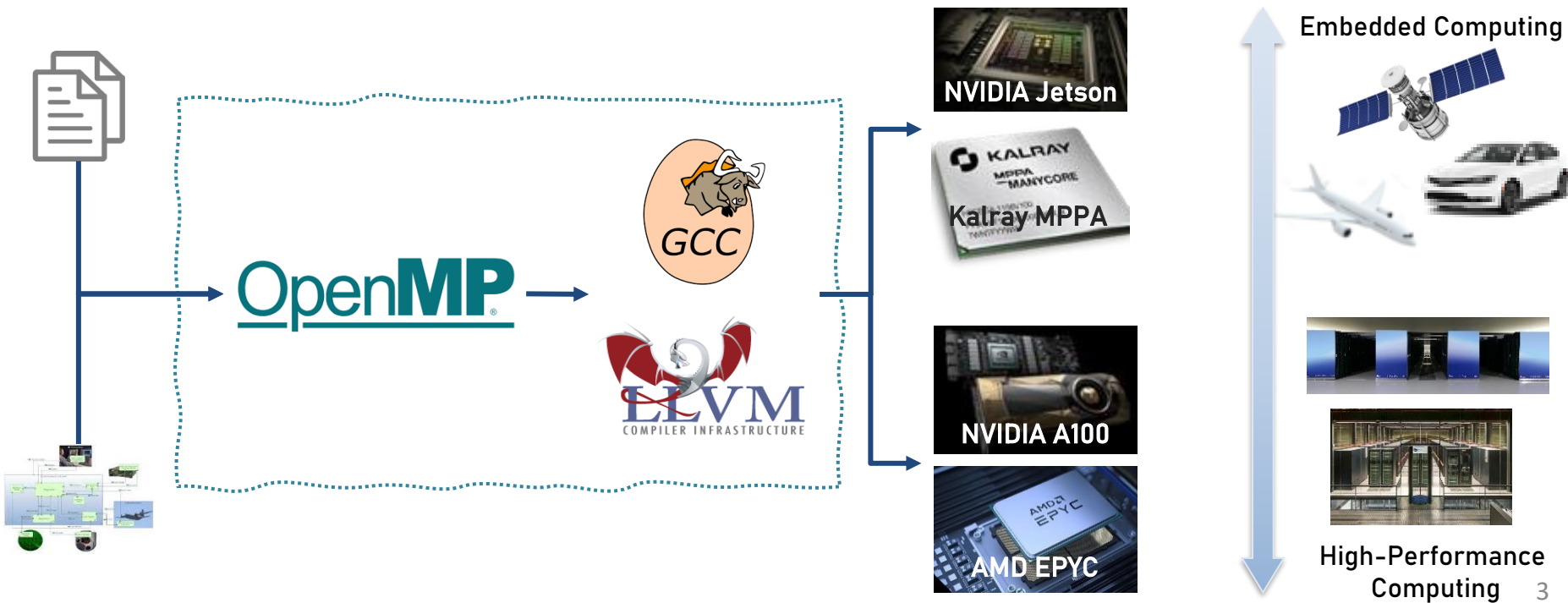
Dr. Sara Royuela
Senior Researcher



www.ampere-euproject.eu

PPC in the scope of OpenMP

Specification, compiler and runtime support in OpenMP
targeting **performance, predictability and resilience** in multiple domains



- The overhead of tasking
- The Task Dependency Graph
 - Performance
 - Memory consumption
 - Interoperability with CUDA graphs
- The RISING Stars and the AMPERE project

Tasking overheads

Podobas, A., Brorsson, M., and Faxén, K. F.

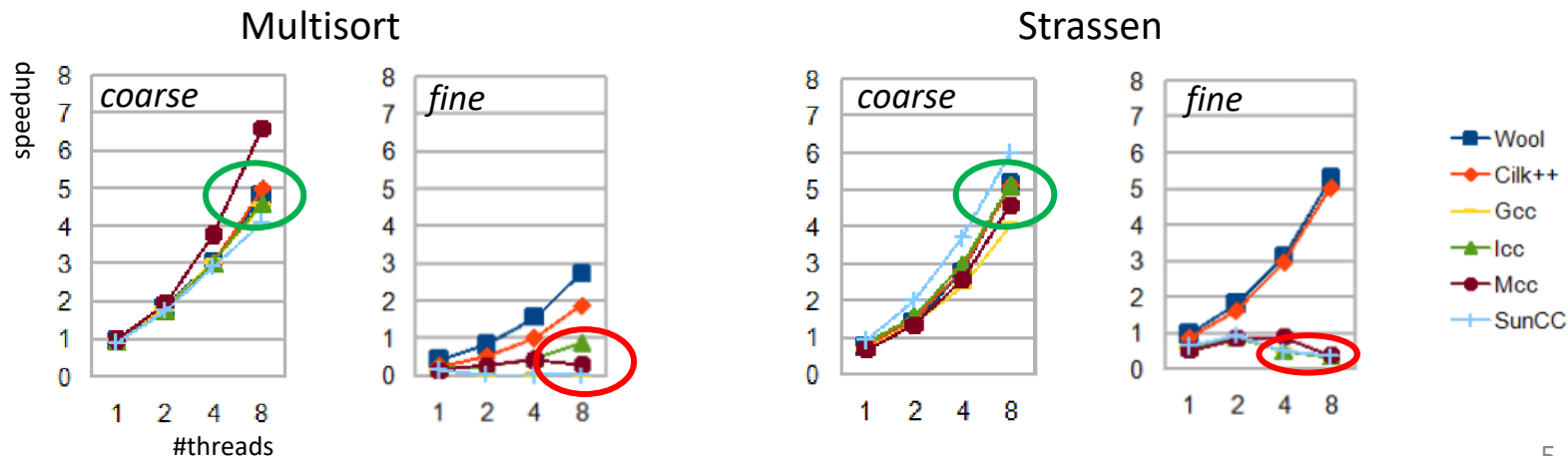
In 3rd workshop on programmability issues for multi-core computers.

A comparison of some recent task-based parallel programming models. 2010.

OpenMP 3.1

✓ Motivation for tasking: focus on the exposing parallelism rather than figuring out how to fit in a specific machine.

✗ Real limitations: fine granularities and deep cut-offs introduce too much overhead, reducing potential speed up.



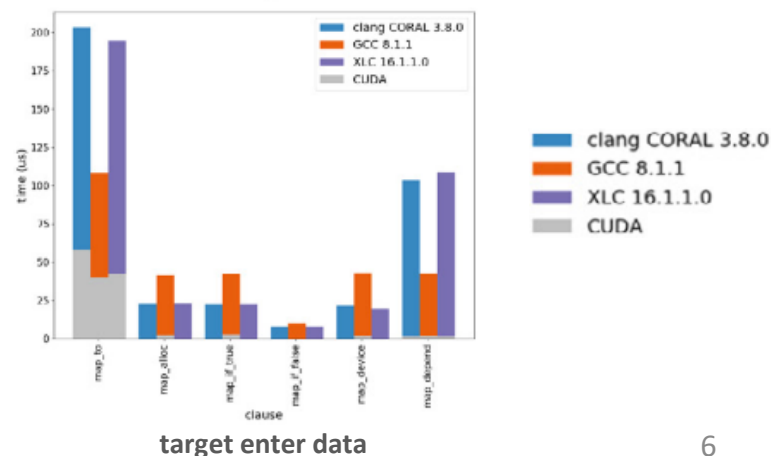
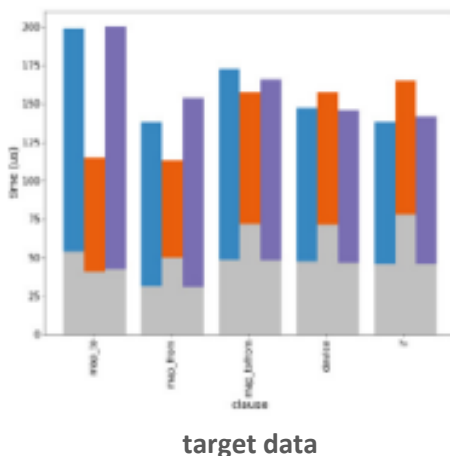
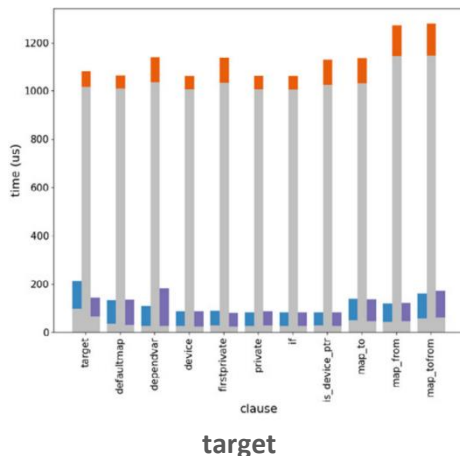
Tasking overheads

Diaz, J.M., Friedline, K., Pophale, S., Bernholdt, D.E., and Chandrasekaran, S.
Parallel Computing. **Analysis of OpenMP 4.5 offloading in implementations:
correctness and overhead.** 2019.

OpenMP 4.0

✓ Motivation for tasking: implement offloading capabilities to exploit accelerator devices.

✗ Real limitations: code transformation might not be optimized, e.g., memory allocations, synchronizations.



Tasking overheads

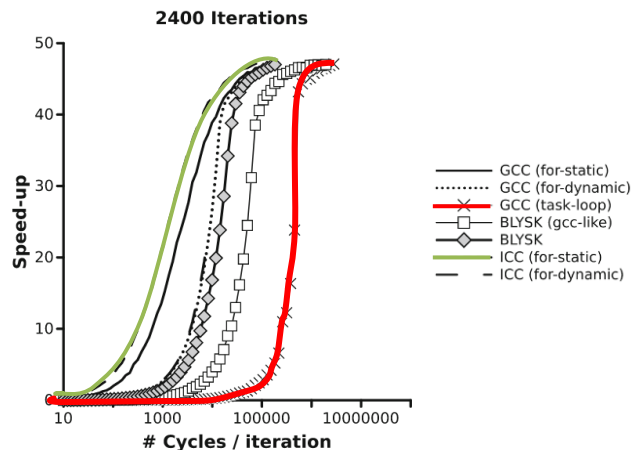
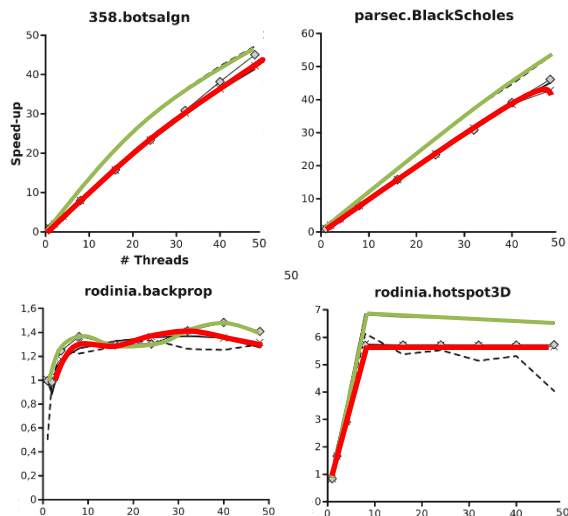
Podobas, A., and Karlsson, S.
In *International Workshop on OpenMP*.

Towards unifying OpenMP under the task-parallel paradigm. 2016.

OpenMP 4.5

✓ Motivation for tasking: taskloop could eliminate the need for thread-parallelism.

✗ Real limitations: the implementation is crucial and determining granularity becomes a challenge.



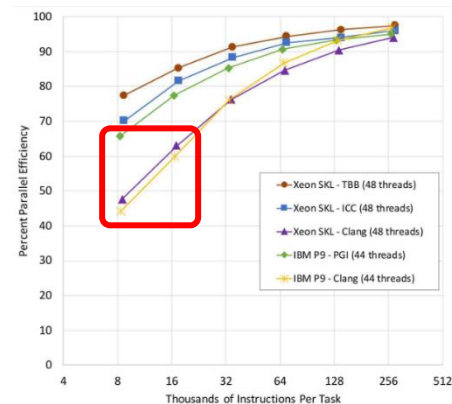
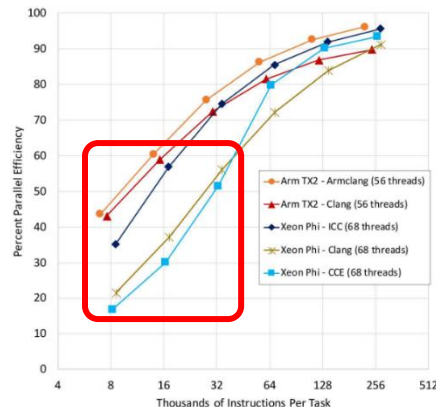
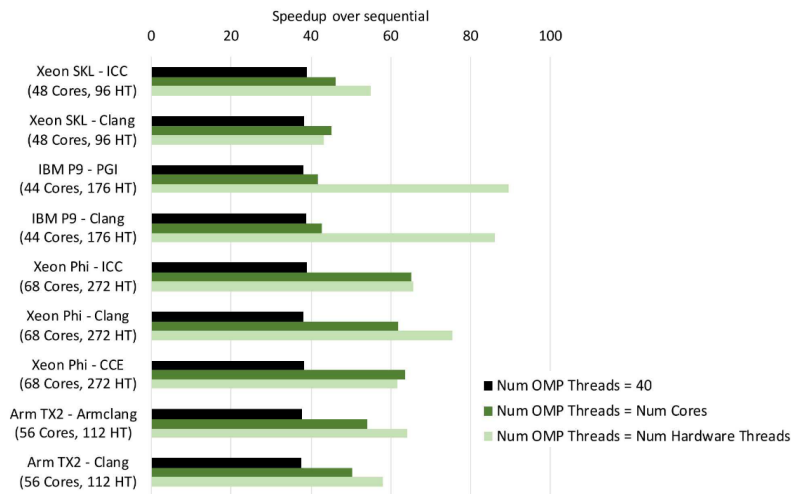
Tasking overheads

Olivier, S.L.. In *International Workshop on OpenMP*.
Evaluating the Efficiency of OpenMP Tasking for Unbalanced Computation on Diverse CPU Architectures. 2020.

OpenMP 5.0

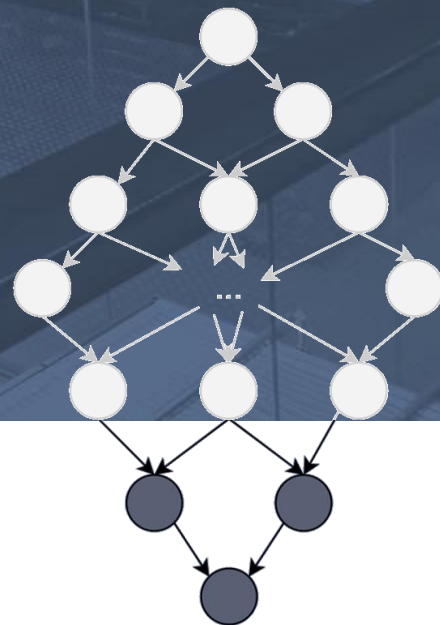
✓ Motivation for tasking: many different extensions that have changed the internals of the implementations.

✗ Real limitations: fine granularity provides poor efficiency 17% - 77%; acceptable granularity from tens of ms.



Where are we

- Tasking is convenient to expose parallelism
- Implementation overheads limit its use for:
 - Fine grained parallelism
 - Loop parallelism
 - Accelerator devices



Let's capitalize on the Task Dependency Graph!

Leveraging the Task Dependency Graph

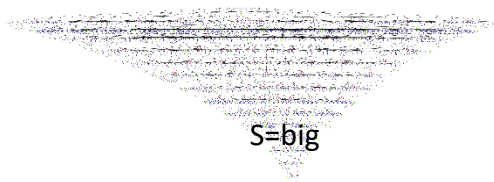
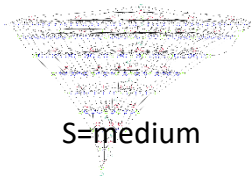
A region of code that can be fully represented as a TDG:

1. Taskified region:

- a) All computations are enclosed in tasks.
- b) Non-taskified code has no side effects on the tasks (e.g., induction variables in loops).

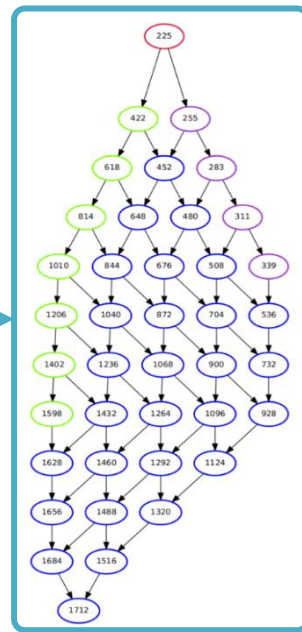
2. TDG shape:

- a) Shape does not change across TDG executions.
- b) Provide information for recomputing the TDG (e.g., a clause with the variables shaping the TDG).



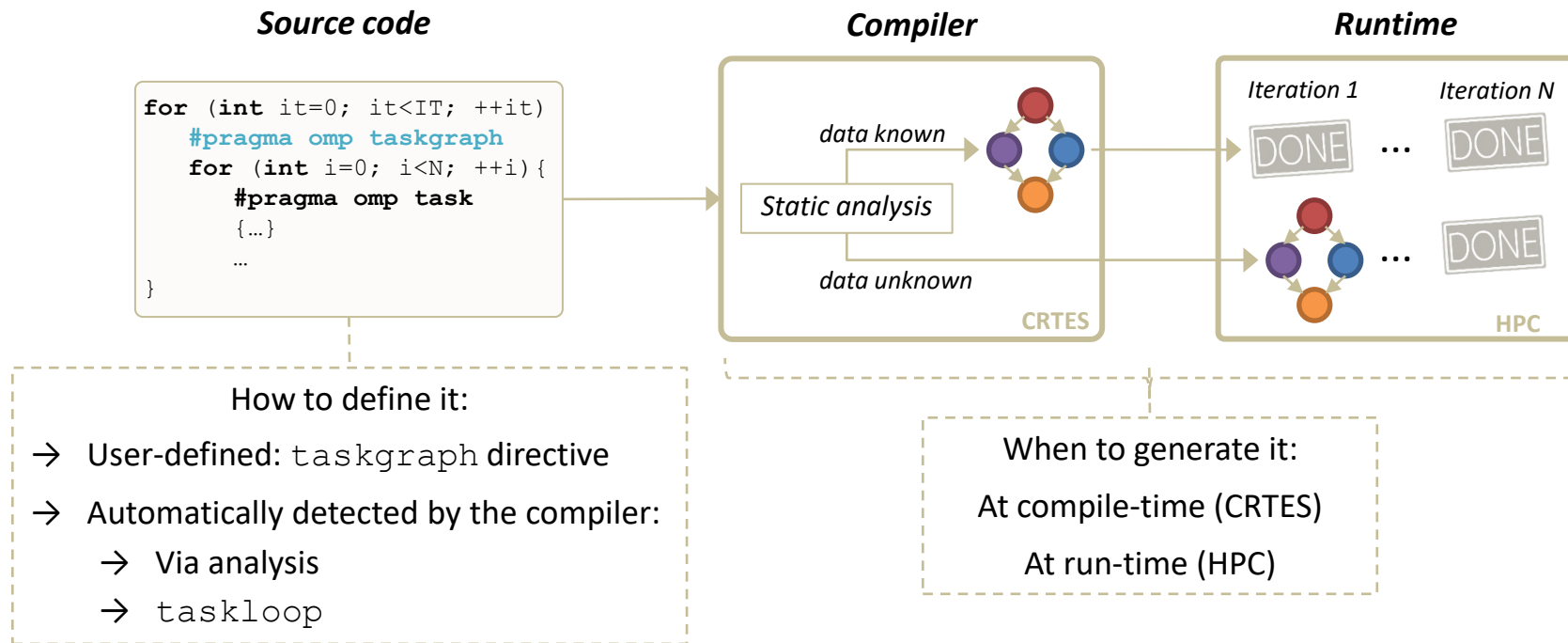
Pedestrian Detector

```
#pragma omp parallel
#pragma omp single
{
    for (i=0; i<N_ITER; ++i) {
        for (by=0; by < BY; by+=BS) {
            for (bx=0; bx < BX; bx+=BS) {
                if (bx==0 && by==0) {
                    #pragma omp task depend(...)
                    {...}
                }
                else if (by==0) {
                    #pragma omp task depend(...)
                    {...}
                }
                else if (bx==0) {
                    #pragma omp task depend(...)
                    {...}
                }
                else {
                    #pragma omp task depend(...)
                    {...}
                }
            }
        }
    }
}
```

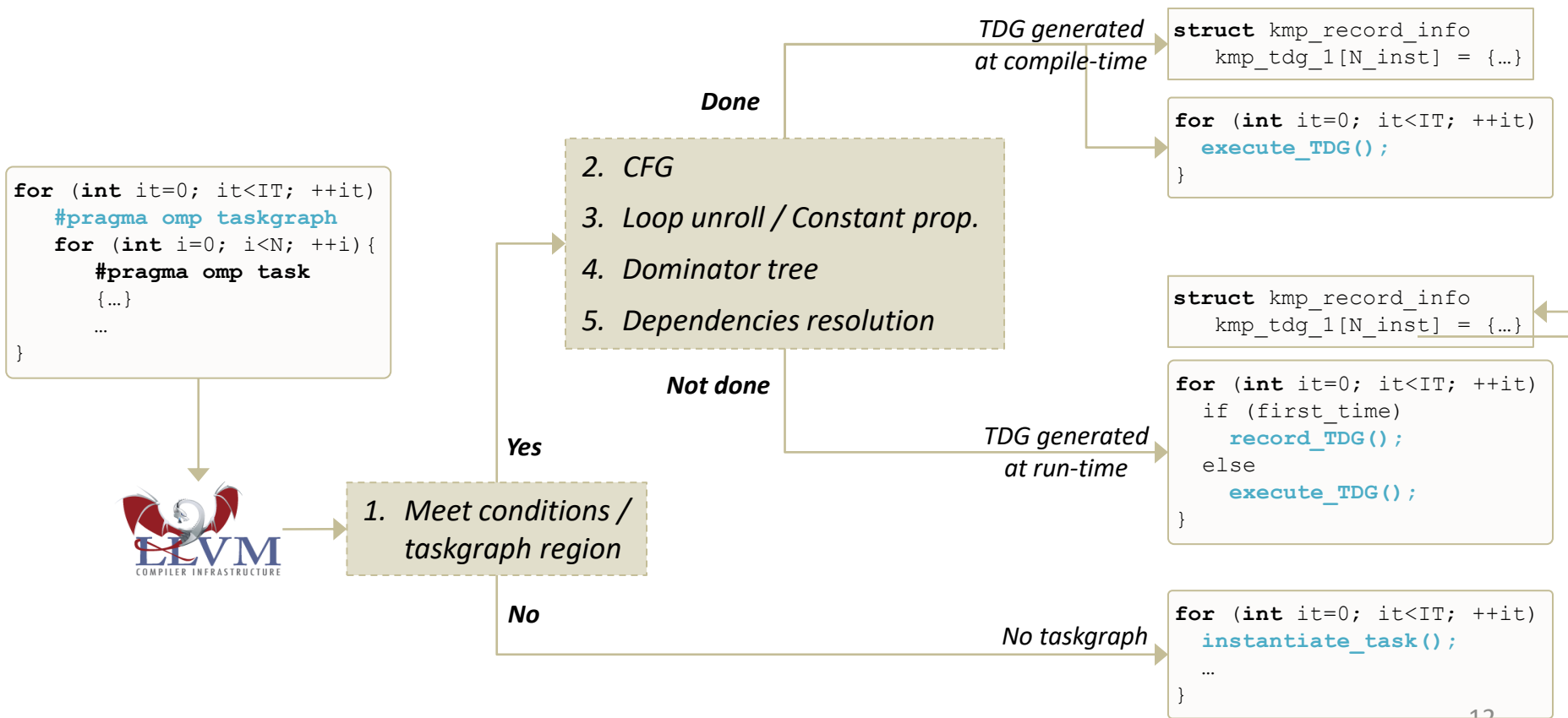


TDG-driven framework

- ❑ **Goals:** Reduce overhead due to task orchestration and dependency resolution
- ❑ **Methodology:** Eliminate the execution of user code to instantiate tasks



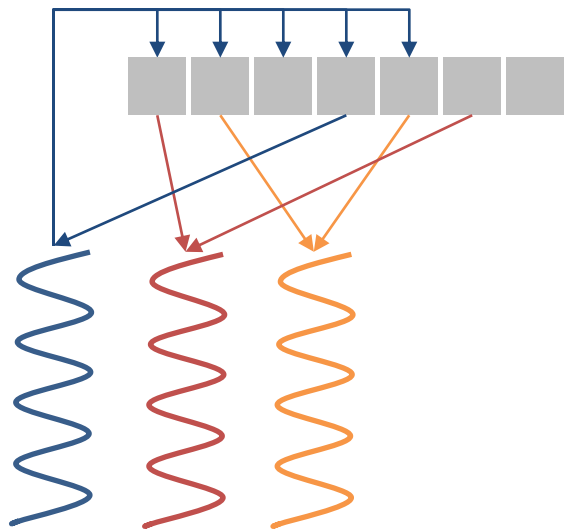
Compiler transformations



Runtime execution: LLVM vs. GCC



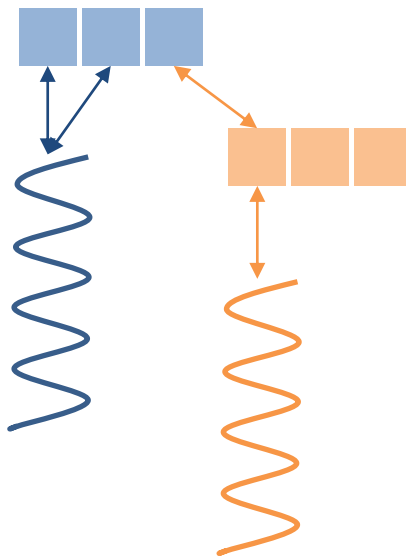
One thread pushes to a single queue, from which all threads pull.



Vanilla

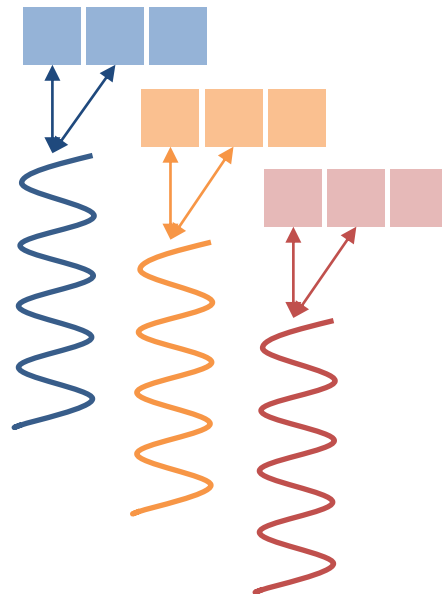


One thread pushes to its queue and the rest steal work from it.



TDG

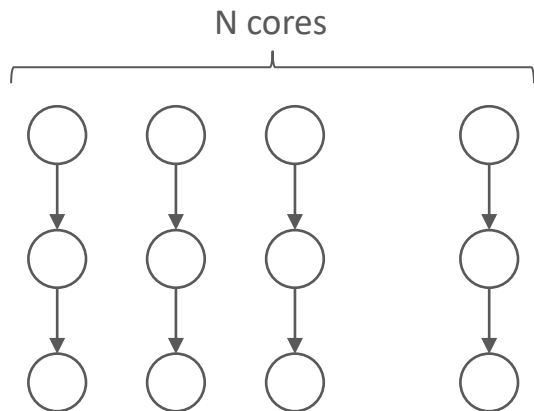
Each thread pushes and pulls from its own queue. Work-stealing is allowed.



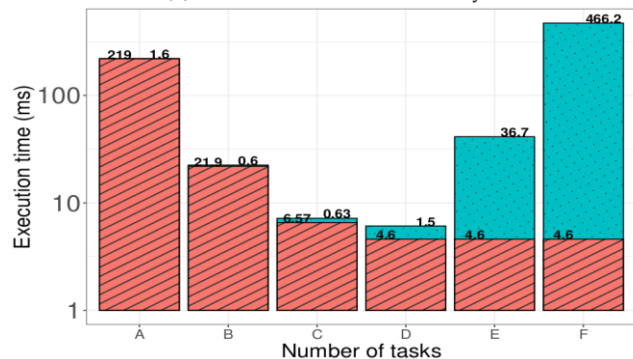
Can we reduce overhead?

Synthetic

```
#pragma omp parallel
#pragma omp single
{
    for (int i=0; i<N_Tasks; ++i) {
        int index = I % N_Cores;
        #pragma omp task depend(out:deps[index])
        fn();
    }
}
```



Task orchestration overhead



$$\text{Computation} = \frac{\text{serial_time}}{\text{\#threads}}$$

$$\text{Overhead} = \text{Total_time} - \text{Computation}$$

A = 1 task of 10⁹ inst., B = 10 tasks of 10⁸ inst., C = 10² tasks of 10⁷ inst.,
D = 10³ tasks of 10⁶ inst., E = 10⁴ tasks of 10⁵ inst., F = 10⁵ tasks of 10⁴ inst.

- ✓ Reduce #instructions needed to orchestrate tasks
- ✓ Alleviate contention

#tasks	10 ⁰	10 ¹	10 ²	10 ³	10 ⁴	10 ⁵
Vanilla	1.6	0.6	0.6	1.5	36.7	466.2
Taskgraph	0.2	0.1	0.2	0.3	9.9	132.2

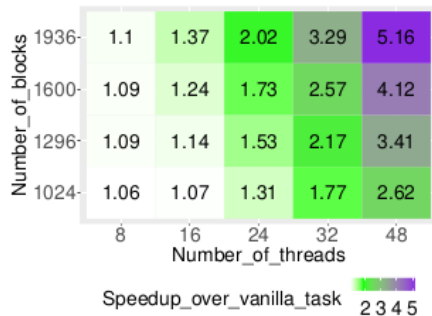
Scalability

Speed-up of *TDG-driven* execution compared to vanilla *task* and *taskloop* implementations using different number of threads and different task granularities

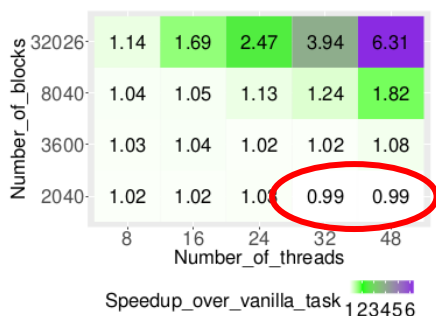


Explicit tasks

Heat propagation simulation

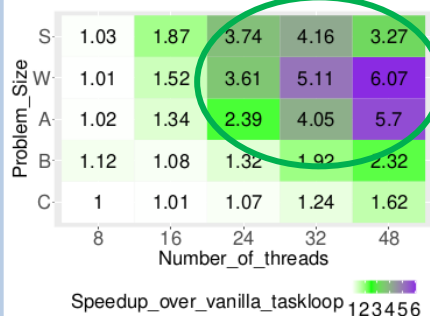


HoG Object detector

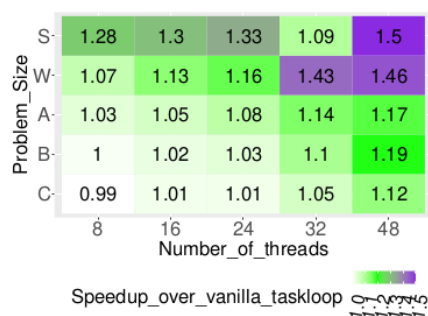


Taskloop

NAS_CG



NAS_LU



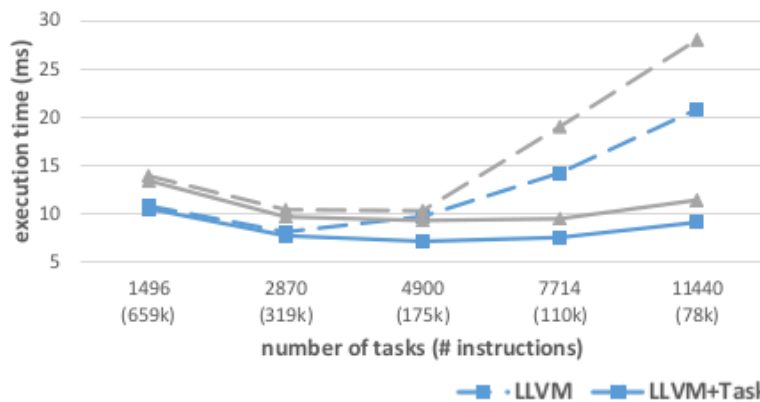
→ Negligible negative impact

→ Considerable positive impact for fine granularities and high thread contention

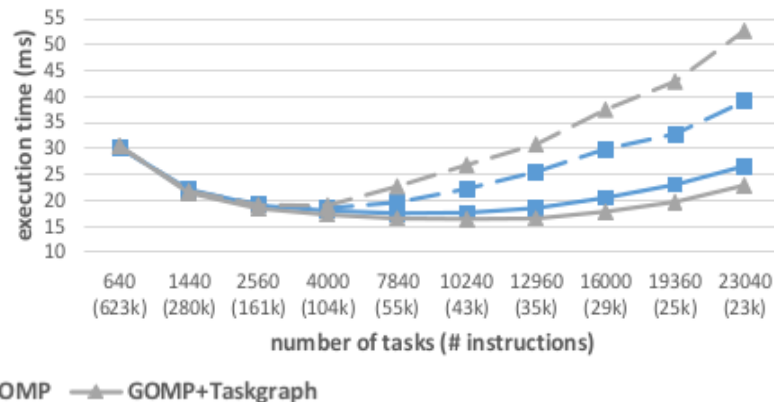
Portability

Speed-up of *TDG-driven* execution compared to vanilla GCC and LLVM implementations using different task granularities

SparseLU matrix decomposition



Heat diffusion simulator



- Coarse grained tasks provide comparable results
- Fine grained tasks show certain stability
- Benefits are portable across compilers/RTLs

Memory management

Goals:

1. Avoid dynamic allocation of task structures
2. Reduce and bound the memory requirements of the OpenMP RTL

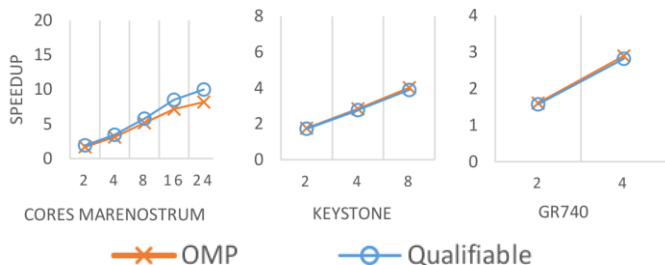
Methodology:

1. Compiler: static generation of the required task structures
2. Runtime: *lazy task creation* (task created when dependencies fulfilled)

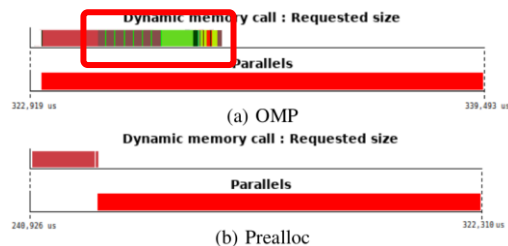


Space-Time Adaptive Processing

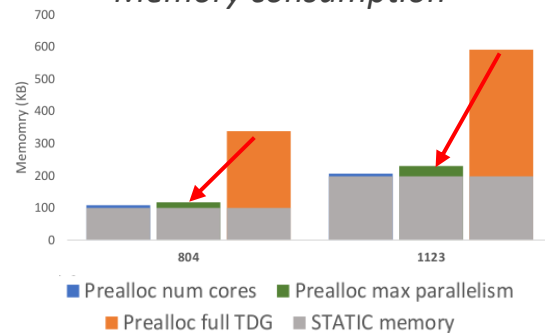
Performance speedup



Use of dynamic memory

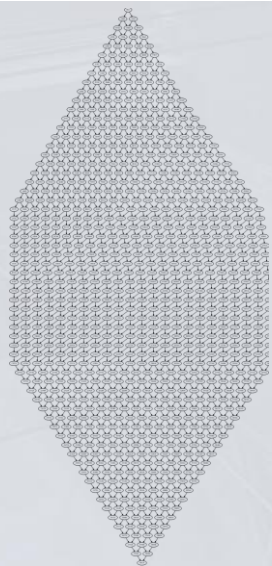


Memory consumption

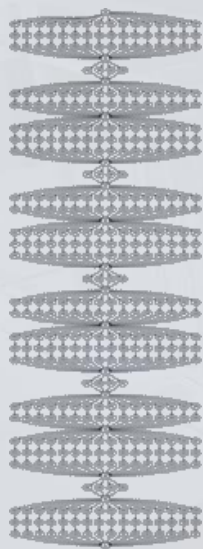


Applicability

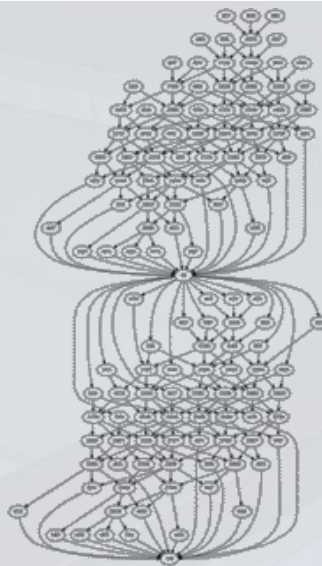
Pedestrian detector
(automotive)



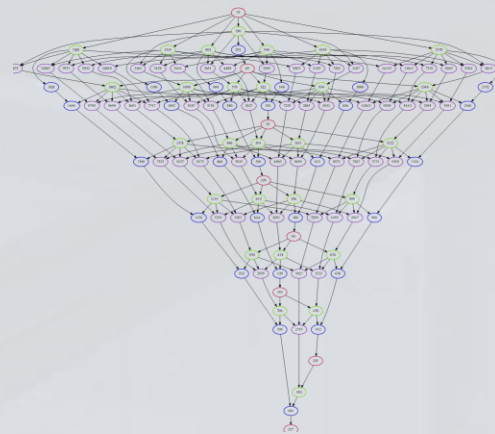
Infra-red sensor processing
(space)



3D Path Planning
(avionics)



Cholesky Factorization
(HPC)



Many applications can be represented as a TDG!

The RISING stars project

- ❑ Enable a versatile and efficient **data acquisition** providing interoperability between different programming models (OpenMP, CUDA)
- ❑ Expose **data acquisition/transfer** in the **programming model**
- ❑ Introduce **real-time** oriented features in the **programming model** to define periodicity, preemption, migration, and allocation.
- ❑ Use cases: Adaptive optics, adaptive beamforming, the Square Kilometer Array and Space Situational Awareness.



Rising
STARS

www.risingstars-project.eu

RISE International Network
for Solutions Technologies
and Applications of
Real-time Systems

CUDA Graphs

TDG

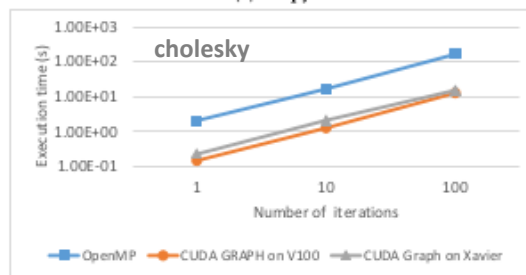


```

cudaGraphNode_t node_17 ;
cudaKernelNodeParams nodeArgs_17 = { 0 } ;
nodeArgs_17.func = (void *) potrf;
void * kernelArgs_17[3] = {&Ah[1][1], &ts, &ts};
nodeArgs_17.kernelParams = (void **) kernelArgs_17;
cudaGraphAddKernelNode(&node_17, graph[0], NULL, 0, &nodeArgs_17);

cudaGraphNode_t node_82 ;
cudaHostNodeParams nodeArgs_82 = {0} ;
nodeArgs_82.func = (void *) trsm;
void * hostArgs_82[4] = {&Ah[1][1], &Ah[1][1], &ts, &ts};
nodeArgs_82.kernelParams = (void **) hostArgs_82;
cudaGraphAddHostNode(&node_82, graph[0], &node_17, 1, &nodeArgs_82);

```



20

CUDA memory management strategies



Execution time in ms.

*Faster when there is no need
for unified memory*

App	Graph nodes	Unified memory + prefetch	Unified memory (non-prefetch)	Zero copy	cudaMemcpy	cudaMemcpyAsync
Vector addition	1024	739	865	687	448	435
	32	592	738	595	344	343
Saxpy	1024	658	916	538	441	441
	32	459	572	452	286	255
Nbody	1024	6080	6058	6464	6041	6094
	32	762	775	806	764	765

*Faster execution
when reducing the
number of nodes*

*Pre-fetching reduces
page faults*

The AMPERE project

- ❑ Use of **Domain Specific Modeling Languages** and **high-level synthesis methods** for building *correct-by-construction* systems.
- ❑ Use OpenMP to provide the **performance** needed to develop complex Cyber-Physical Systems:
 - Predictive Cruise Control (automotive)
 - Obstacle Detection and Avoidance System (railway)
- ❑ Provide mechanisms to guarantee **non-functional requirements**: time predictability, resilience and energy consumption.



A Model-driven development framework for highly Parallel and
EnErgy-Efficient computation supporting multi-criteria optimisation

www.ampere-euproject.eu

AMALTHEA DSML to OpenMP

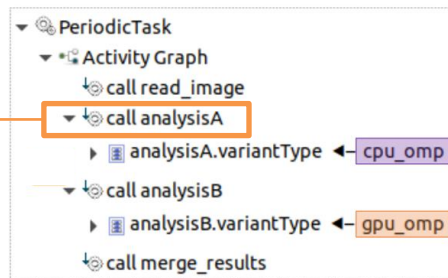


Automatic
code generation



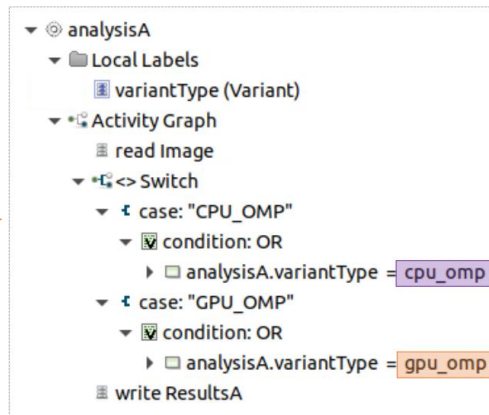
Goals:

1. Exploit parallelism within OS tasks with OpenMP (host and target) tasks
2. Exploit heterogeneity through specializations



OS
task

```
void PeriodicTask() {  
    #pragma omp parallel  
    #pragma omp single  
    {  
        #pragma omp task depend(out:Image)  
        { read_image(); }  
        #pragma omp task depend(in:Image) \  
            depend(out:ResultsA) cpu_omp  
        { analysisA(); }  
        #pragma omp target depend(in:Image) \  
            depend(out:ResultsA) \  
            map(to:Image) map(from:ResultsA) \  
            gpu_omp  
        { analysisB(); }  
        #pragma omp task depend(in:ResultsA, ResultsB)  
        { read_image(); }  
    }  
}
```

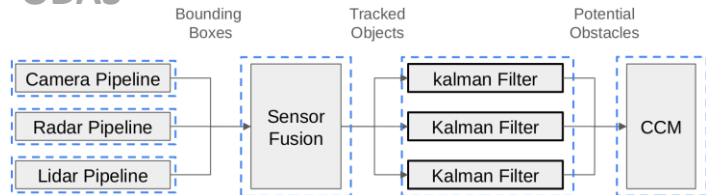


OpenMP
task

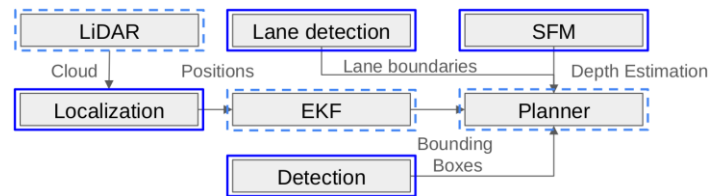
```
void analysisA_gpu() { ... }  
  
#pragma omp declare variant(analysisA_gpu) \  
    match(construct={target}) \  
    implementation={extension(gpu_omp) }  
void analysisA() { ... }
```

AMALTHEA DSML to OpenMP: Performance

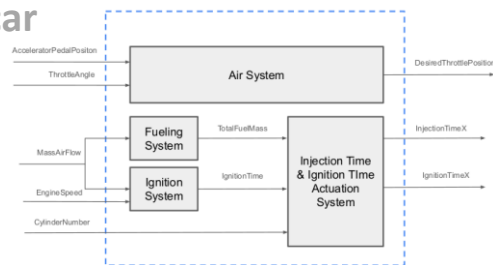
ODAS



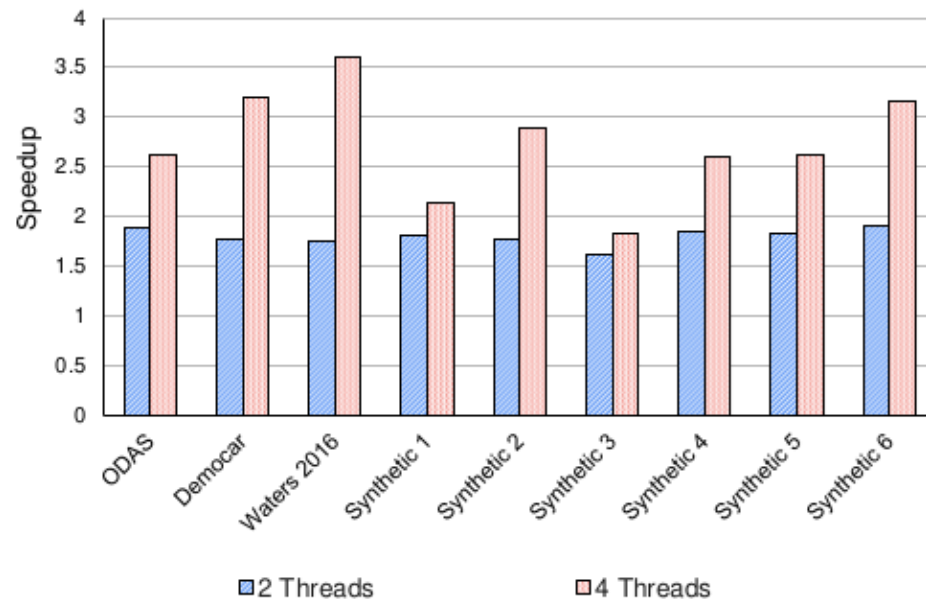
WATERS



Democar



NVIDIA Jetson TX2 board
with a GPU and a 4-core ARM CPU

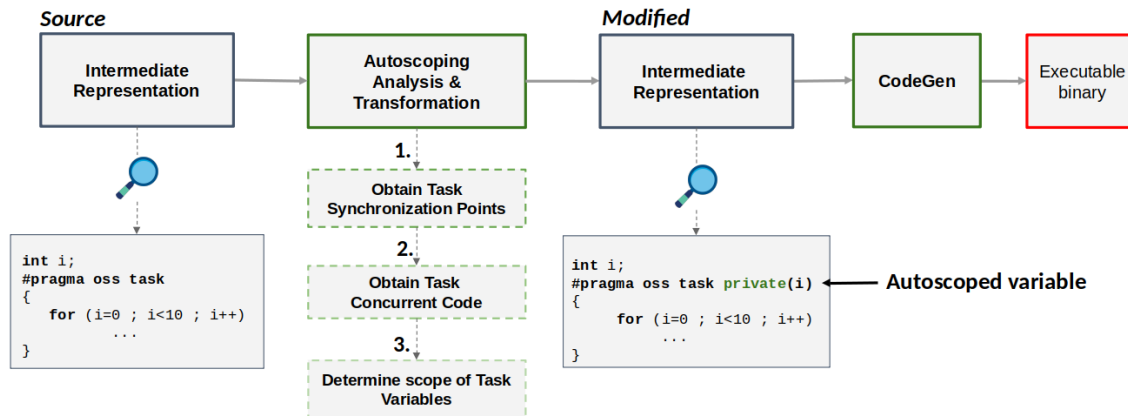


Correctness analysis for OpenMP

Goals:

1. Detect/resolve race conditions
2. Detect/correct wrong synchronizations (task dependencies, memory fences)
3. Detect inconsistencies in the data-sharing attributes

*LLVM Compiler pipeline to
define/correct data-
sharing attributes:*



- A. Munera, S Royuela, R Ferrer, R Peñacoba, E Quiñones. **Static analysis to enhance the programmability and performance in OmpSs-2**, In *ISC*. 2020.
- S. Royuela, R Ferrer, D Caballero, X Martorell. **Compiler analysis for OpenMP tasks correctness**, In *CF*. 2015.
- S. Royuela, A Duran, X Martorell. **Compiler automatic discovery of OmpSs task dependencies**, In *LCPC*. 2012.
- S. Royuela, A Duran, C Liao, DJ Quinlan. **Auto-scoping for OpenMP tasks**, In *IWOMP*. 2012.

GOAL: Assess predictability of OpenMP to allow schedulability analysis

Requirements:

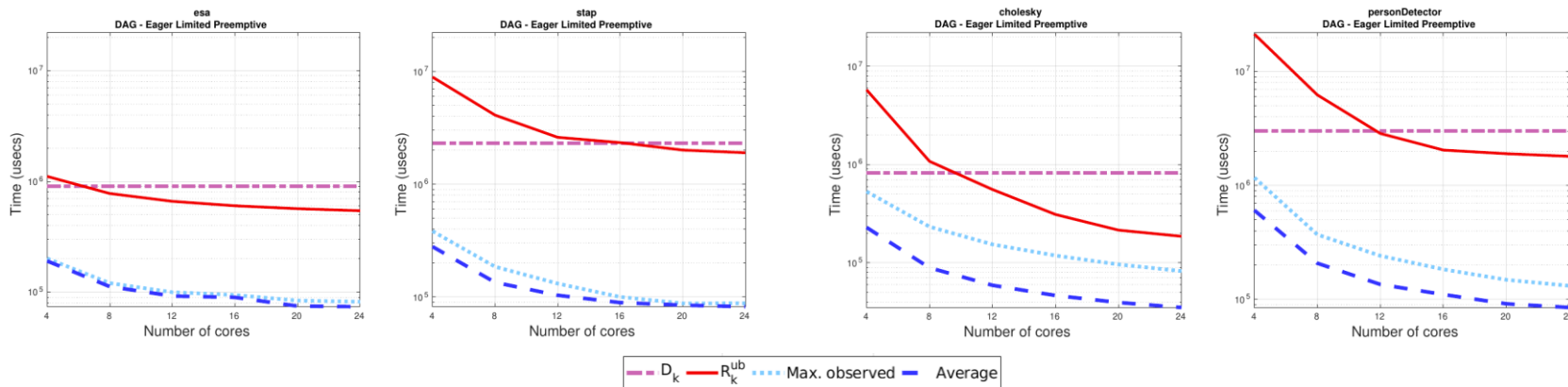
1. *Work-conserving* scheduler for non-pessimistic WCRT analysis
2. Prescriptive *task priorities* to support fixed priority schedulers
3. Prescriptive implementation of *Task Scheduling Points* to allow limited preemptive scheduling
 - *Taskyield*, to alleviate pessimism and enhance schedulability

Response time upper bound (R_k^{ub}):

$$R_k^{ub} \leftarrow \text{len}(G_k) + \frac{1}{m} (\text{vol}(G_k) - \text{len}(G_k)) + \frac{1}{m} (I_k^{hp} + I_k^{lp})$$

R_k^{ub} bounds the maximum observed execution time

Real-time system schedulable with at least 16 cores



Resilience with OpenMP

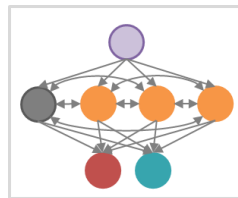
GOAL: Task-level replication for fault-detection

Augmented OpenMP code

```
#pragma omp task redundancy(  
    spatial|temporal|spatial_temporal,  
    3, /*n replicas*/  
    var1:func1,  
    var2:func2)  
{ /*code*/ }
```

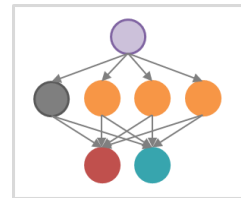
- Data capturing
- Replicated tasks
- Consolidation function / variable to check

TDG with replication



Temporal:

OpenMP mutexinoutset
between replicas

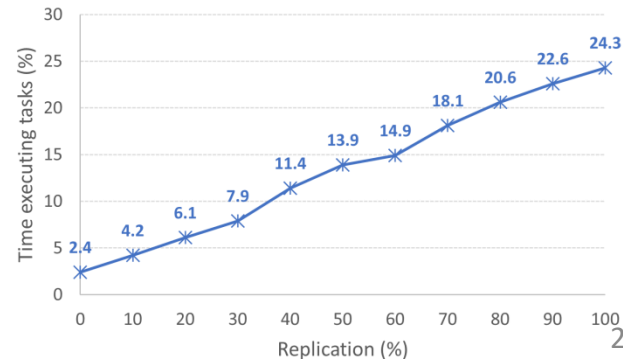
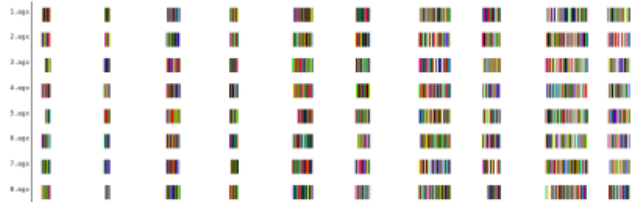


Spatial:

Each replica in a
different core

Engine control management system

- TDG1 (T10, T20, T50)
- TDG2 (T10, T50)
- TDG3 (T10, T20)
- TDG4 (T10)



Optimizations

1. C. Yu, S. Royuela, E. Quiñones, **Enhancing OpenMP tasking model: Performance and portability**, In IWOMP 2021.
2. A. Munera, S. Royuela, E. Quiñones, **Towards a qualifiable OpenMP framework for embedded systems**, In DATE 2020.
3. R. E. Vargas, S. Royuela, M. A. Serrano, X. Martorell, E. Quiñones, **A Lightweight OpenMP4 Run-time for Embedded Systems**, In ASP-DAC 2016.

Interoperability

4. C. Yu, S. Royuela, E. Quiñones, **OpenMP to CUDA graphs: A compiler-based transformation to enhance the programmability of NVIDIA devices**, In SCOPES 2020.
5. S. Royuela, L.M. Pinho, E. Quiñones, **Enabling Ada and OpenMP runtimes interoperability through template-based execution**, In JSA 2020.
6. S. Royuela, L.M. Pinho, E. Quinones, **Converging Safety and High-performance Domains: Integrating OpenMP into Ada**, In DATE 2018.

Functional safety

7. S. Royuela, A. Duran, M. A. Serrano, E. Quiñones, **A functional safety OpenMP for critical real-time embedded systems**, In IWOMP 2017.
8. S. Royuela, X. Martorell, E. Quinones, L. M. Pinho, **OpenMP Tasking Model for Ada: Safety and Correctness**, In AEiC 2017.
9. S. Royuela, R. Ferrer, D. Caballer, X. Martorell, **Compiler analysis for OpenMP tasks correctness**, In CF 2015.

Predictability and CRTES

10. M. A. Serrano, S. Royuela, E. Quiñones. **Towards an OpenMP Specification for Critical Real-time Systems**. In IWOMP 2018.
11. M. A. Serrano, A. Melani, S. Kehr, M. Bertogna, E. Quiñones, **An Analysis of Lazy and Eager Limited Preemption Approaches under DAG-based Global Fixed Priority Scheduling**, In ISORC 2017.
12. M. A. Serrano, A. Melani, M. Bertogna, E. Quiñones, **Response-Time Analysis of DAG Tasks under Fixed Priority Scheduling with Limited Preemptions**, In DATE 2016.
13. M. A. Serrano, A. Melani, R. Vargas, A. Marongiu, M. Bertogna, E. Quiñones, **Timing Characterization of OpenMP4 Tasking Model**, In CASES 2015.

Projects and collaborations



DENSO



BOSCH

THALES



**Barcelona
Supercomputing
Center**
Centro Nacional de Supercomputación

High-performance OpenMP tasking

Sara Royuela
sara.royuela@bsc.es



OpenMP User's Monthly Telecon

October 28, 2022