

Vector operations, tiled operations, distributed execution, task graphs, ...

What next?

Samuel Thibault, University of Bordeaux Inria STORM Team

- Frontier supercomputer
- #1 on top500.org LINPACK benchmark
- 1.102 EFlop/s
- 9472 nodes
 - 1 AMD Epyc CPU
 - 4 AMD Instinct 250X GPUs
 - 4TB flash memory
 - Infinity Fabric interconnect
- Slingshot 100GB/s network
- 21.1 MW

Ínría

• Also #1 Green 500



Classical parallel programming

threads



Classical parallel programming

https://starpu.gitlabpages.inria.fr/

threads

Ínría

• MPI+threads



Classical parallel programming

https://starpu.gitlabpages.inria.fr/

threads

- MPI+threads
- MPI+threads+GPU



Classical parallel programming

https://starpu.gitlabpages.inria.fr/

threads

- MPI+threads
- MPI+threads+GPU
- Data transfers



Classical parallel programming

https://starpu.gitlabpages.inria.fr/

threads

- MPI+threads
- MPI+threads+GPU
- Data transfers



Classical parallel programming

https://starpu.gitlabpages.inria.fr/

- threads
- MPI+threads
- MPI+threads+GPU
- Data transfers
- Disk storage



Classical parallel programming

- threads
- MPI+threads
- MPI+threads+GPU
- Data transfers
- Disk storage





Ínría

Classical parallel programming

- threads
- MPI+threads
- MPI+threads+GPU
- Data transfers
- Disk storage





Classical parallel programming

- threads
- **MPI+threads**
- MPI+threads+GPU
- Data transfers
- Disk storage ۲

Ínnía

Manage it all by hand?!



Classical parallel programming

- threads
- MPI+threads
- MPI+threads+GPU
- Data transfers
- Disk storage

Ínnía

Manage it all by hand?!

Need high-level abstractions



















Vector computers (end '70s, '80s)



Vector computers (end '70s, '80s)



Vector computers (end '70s, '80s)





Vector computers (end '70s, '80s)



Computers with cache (early '90s-today), blocked operations : BLAS



Computers with cache (early '90s-today), blocked operations : BLAS





Computers with cache (early '90s-today), blocked operations : BLAS



Computers with cache (early '90s-today), blocked operations : BLAS



Distributed Computers (late '90s-today), 2D-block-cyclic distribution



Distributed Computers (late '90s-today), 2D-block-cyclic distribution



Distributed Computers (late '90s-today), 2D-block-cyclic distribution



Distributed Computers (late '90s-today), 2D-block-cyclic distribution







Task graph (~'08-today)












Task graph (~'08-today)





Task graph (~'08-today)





Task graph (~'08-today)





Task graph (~'08-today), priorities





Task graph (~'08-today), priorities



Ínría https://starpu.gitlabpages.inria.fr/

Task graph (~'08-today), priorities





Task graphs in HPC

Runtime systems

- OmpSs, PARSEC, StarPU, SuperGlue/DuctTeip, XKaapi...
 Standards
- OpenMP4.0 introduced task dependencies

Applications

- Chameleon, DPLASMA, SLATE for dense linear algebra
- qr_mumps, PaStiX for sparse linear algebra
- ScalFMM for FMM
- ...

nría

Task-based support

- Task/data scheduling
 - Load balancing
 - Data prefetching
 - Pipelining
 - GPU memory limitation management
- Distributed execution through MPI
- Out-of-core: optimized swapping to disk
- High-level performance analysis
- Performance bounds

nnía

- Debugging sequential execution
- Reproducible performance simulation





- Small enough to get parallelism to feed all processing units
- Large enough to efficiently use the processing units



From PARSEC : « Hierarchical DAG Scheduling for Hybrid Distributed Systems », Wu, Bouteiller, Bosilca, Faverge, Dongarra

https://starpu.gitlabpages.inria.fr/

Inría

GPUs

• Efficient only with large tile sizes

CPUs

nría

- Need many tasks
- \rightarrow Hybrid task size

More generally, hierarchical task graphs

• Seen in CEA, OmpSs, PaRSEC, StarPU

Outline

- How big should a task be?
- Hierarchical task graphs
- Opportunities









- Lower bound due to runtime overhead ullet
 - Proposed by Martin Tillenius for DuctTeip (U. Uppsala)



https://starpu.gitlabpages.inria.fr/



From PARSEC : « Hierarchical DAG Scheduling for Hybrid Distributed Systems », Wu, Bouteiller, Bosilca, Faverge, Dongarra

https://starpu.gitlabpages.inria.fr/

Inría

GPUs

- Have thousands for cores to feed
- Newer generations require yet larger sizes
- Can run several kernels at the same time
 - Still limited



From PARSEC : « Hierarchical DAG Scheduling for Hybrid Distributed Systems », Wu, Bouteiller, Bosilca, Faverge, Dongarra

https://starpu.gitlabpages.inria.fr/

Inría

CPUs

- Have many independent cores
 - \rightarrow Need many tasks
- Can use parallel implementations (e.g. from MKL)
 - But better have subtasks to interleave them

How size does not fit all

Gwenolé Lucas experimented with Chameleon

- Cholesky inversion
- 2 NVIDIA V100
- 2 Xeon Gold 6142
- Different matrix tile sizes : 2880 / 960 / 320

Hybrid tile sizes



- Non-Hierarchical: 960
- Non-Hierarchical: 320

Innía

- Hierarchical: 2880 / 960
- Hierarchical: 2880 / 960 / 320

https://starpu.gitlabpages.inria.fr/

Tuning the tile size?

- B = "big" tile size, for GPUs
 - Large enough for GPUs
 - Not too large for parallelism
- b = "small" tile size, for CPUs
 - Large enough for CPUs
 - Not too large for parallelism

With Parsec, Wu, Bouteiller, Bosilca, Faverge and Dongarra experimented tuning



From PARSEC : « Hierarchical DAG Scheduling for Hybrid Distributed Systems », Wu, Bouteiller, Bosilca, Faverge, Dongarra

https://starpu.gitlabpages.inria.fr/

Ínría



From PARSEC : « Hierarchical DAG Scheduling for Hybrid Distributed Systems », Wu, Bouteiller, Bosilca, Faverge, Dongarra

https://starpu.gitlabpages.inria.fr/

Ínría



From CEA : « A hierarchical fast direct solver for distributed memory machines with manycore nodes », Augonnet, Goudin, Kuhn, Lacoste, Namyst, Ramet

https://starpu.gitlabpages.inria.fr/

Ínría

Multiple answers

- Depends on available platform parallelism
- Depends on available application parallelism
- Depends on application phases

Automatically adapt?







Applications themselves are recursive

• e.g. h-matrices



From Airbus Group



Applications themselves are recursive

• e.g. h-matrices





Applications themselves are recursive

• e.g. h-matrices

























```
• In OpenMP?
```

```
TRSM(A, B) {
```

#pragma omp task depend(in:A, inout:B) wait

{

```
#pragma omp task \
```

```
depend(in:A[0][0], inout:B[0][0])
TRSM(A[0][0], B[0][0])
```

```
#pragma omp task \
```

```
depend(in:A[0][0], inout:B[1][0])
TRSM(A[0][0], B[1][0])
```

```
#pragma omp task \
```

```
depend(in:A[1][0], in:B[0][0], inout:B[0][1])
GEMM(A[1][0], B[0][0], B[0][1])
```


Synchronization concerns

- Fork-join parallelism
- Hindered by synchronization induced by task graph



Inría

```
#pragma omp task \
```

depend(in:A[0][0], inout:B[0][0]) TRSM(A[0][0], B[0][0])

#pragma omp task \

```
depend(in:A[0][0], inout:B[1][0])
TRSM(A[0][0], B[1][0])
```

#pragma omp task \

depend(in:A[1][0], in:B[0][0], inout:B[0][1]) GEMM(A[1][0], B[0][0], B[0][1])





- Explicit data registration
- Recursive partitioning



- Explicit data registration
- Recursive partitioning





- Explicit data registration
- Recursive partitioning



- Explicit data registration
- Recursive partitioning



In StarPU

Ínría

- Explicit data registration
- Recursive partitioning



https://starpu.gitlabpages.inria.fr/

In StarPU

Innía

- Multi-level data coherency, among CPUs-GPUs-disk
- Enables seamless recursive tasks
- More details in Gwenolé Lucas' talk Thursday 14:40



https://starpu.gitlabpages.inria.fr/

Multi-level data coherency

- Synchronization pseudo-tasks
 - Only when needed
- Result is exactly as appropriate
- \rightarrow Just split at will!



Multi-level data coherency

- Synchronization pseudo-tasks
 - Only when needed
- Result is exactly as appropriate
- \rightarrow Just split at will!



No extra synchronization

 \rightarrow Can consider task graph subdivision as a tree



Recursion

No extra synchronization

 \rightarrow Can consider task graph subdivision as a tree



No extra synchronization

 \rightarrow Can consider task graph subdivision as a tree



No extra synchronization

- \rightarrow Can consider task graph subdivision as a tree
- \rightarrow Decide at will where and when to stop recursing



https://starpu.gitlabpages.inria.fr/

Innía





GPU / CPU efficiency management

- Stop recursing at desired tile size
- Care for latency of the task graph critical path



Parallel submission

Ínría

- Unroll the graph in parallel
- While one PU is computing the first task



Delay unrolling

- Observe behavior before unrolling the rest accordingly
- Decorrelate submission and execution



Ínría h

Delay unrolling

- Observe behavior before unrolling the rest accordingly
- Decorrelate submission and execution



https://starpu.gitlabpages.inria.fr/

Ínría

Scaling at large

Ínría

- Master unrolls higher recursion levels, schedules result
- Slaves unroll the rest
- Master still contention point



https://starpu.gitlabpages.inria.fr/

Scaling at large

- All nodes unroll higher recursion levels, determine task mapping •
- Nodes unroll their own remaining recursion ۲



Scaling at large

- All nodes unroll higher recursion levels, determine task mapping ٠
- Nodes unroll their own remaining recursion ۲



Scheduling at large

- High-level global task graph scheduling, e.g. mapping, critical path
- High-level local task graph scheduling, e.g. memory-based ordering
- O(N²) or even O(N³) not that costly



Inría

- Computation efficiency
- Scaling at large
- Scheduling at large

NumPEx PEPR project



NumPEx project

Building a national ecosystem for the future exaflops machine hosted in France

- Facilitate conception, deployment, and monitoring of composite applications at large scale
- Efficiently exploit a heterogeneous architecture with many accelerators
- Control the energy envelope

80 teams in France ● Maths/computation/data ★ Applications/demonstrators ▲ NumPEx contributors

https://starpu.gitlabpages.inria.fr/





NumPEx project

PC2-ExaSoft

Ínría

Holistic approach

- Contribute to a sound, consistent software stack
- Most components should fit together!
- Bridge the gap between existing languages/software/tools
- Integrate state-of-the-art research results
- Demonstrate relevance on representative applications
- WP3 : Runtime Systems at Exascale





https://starpu.gitlabpages.inria.fr/

Conclusion

- Programming large systems requires abstraction
- Recursive applications are there
- Hierarchical task graphs seems an opportunity for
 - Controlling task size
 - Controlling task submission
 - High-level task graph scheduling



The StarPU runtime system

Development context

- History
 - Started in 2008
 - PhD Thesis of Cédric Augonnet
 - StarPU main core ≈ 70k lines of code
 - Written in C
- Open Source
 - Released under LGPL
 - Sources freely available
 - git repository and nightly tarballs
 - See https://starpu.gitlabpages.inria.fr/
 - Open to external contributors
- [HPPC'08]
- [Europar'09] [CCPE'11],... >1500 citations

Features

- C, OpenMP, OpenCL APIs
- CPUs, GPUs, Phi
- Advanced task mapping & scheduling
- Optimized data transfers
- Cluster Support through MPI
- Out-Of-Core support
- Simulation support

Innía

• Performance analysis tools

https://starpu.gitlabpages.inria.fr/

The StarPU runtime system

Success stories

Task-based programming actually makes things easier!

- QR-Mumps (sparse linear algebra)
 - Non-task version: only 1D decomposition
 - Task version: 2D decomposition, flurry of parallelism
 - With seamless memory control
- H-Matrices (compressed linear algebra, AirBus)
 - Out-of-core support
 - Could run cases unachievable before
 - e.g. 1600 GB matrix with 256 GB memory
 - Shipped to AirBus customers
- Implemented CFD, FMM, CG, stencils, ...





The StarPU runtime system

Supported platforms

- Supported architectures
 - Multicore CPUs (x86, PPC, ...)
 - NVIDIA GPUs
 - OpenCL devices (eg. AMD cards)
 - HIP
 - FPGA (ongoing)
 - Old Intel Xeon Phi (MIC), SCC, Kalray MPPA, Cell (decommissioned)
- Supported Operating Systems
 - Linux

Inría

- Mac OS
- Windows

Applications on top of StarPU

Using CPUs, GPUs, distributed, out of core, ...

- Dense linear algebra
 - Cholesky, QR, LU, ... : Chameleon (based on Plasma/Magma)
- Sparse linear algebra
 - QR_MUMPS
 - PaStiX
- Compressed linear algebra
 - BLR, h-matrices
- Fast Multipole Method
 - ScalFMM
- Conjugate Gradient
- Other programming models : Data flow, skeletons
 - SignalPU, SkePU