



University
of Basel

OpenMP Loop Scheduling Revisited: Making a Case for More Schedules

Florina M. Ciorba¹
hpc.dmi.unibas.ch

ADAC6 Workshop Zurich, 20.06.2018

Work conducted with Christian Iwainsky² and Patrick Buder¹, to appear at iWomp18

¹ University of Basel, Switzerland

² Technische Universität Darmstadt, Germany



SWISS NATIONAL SCIENCE FOUNDATION

HESSEN



Parsing the Title

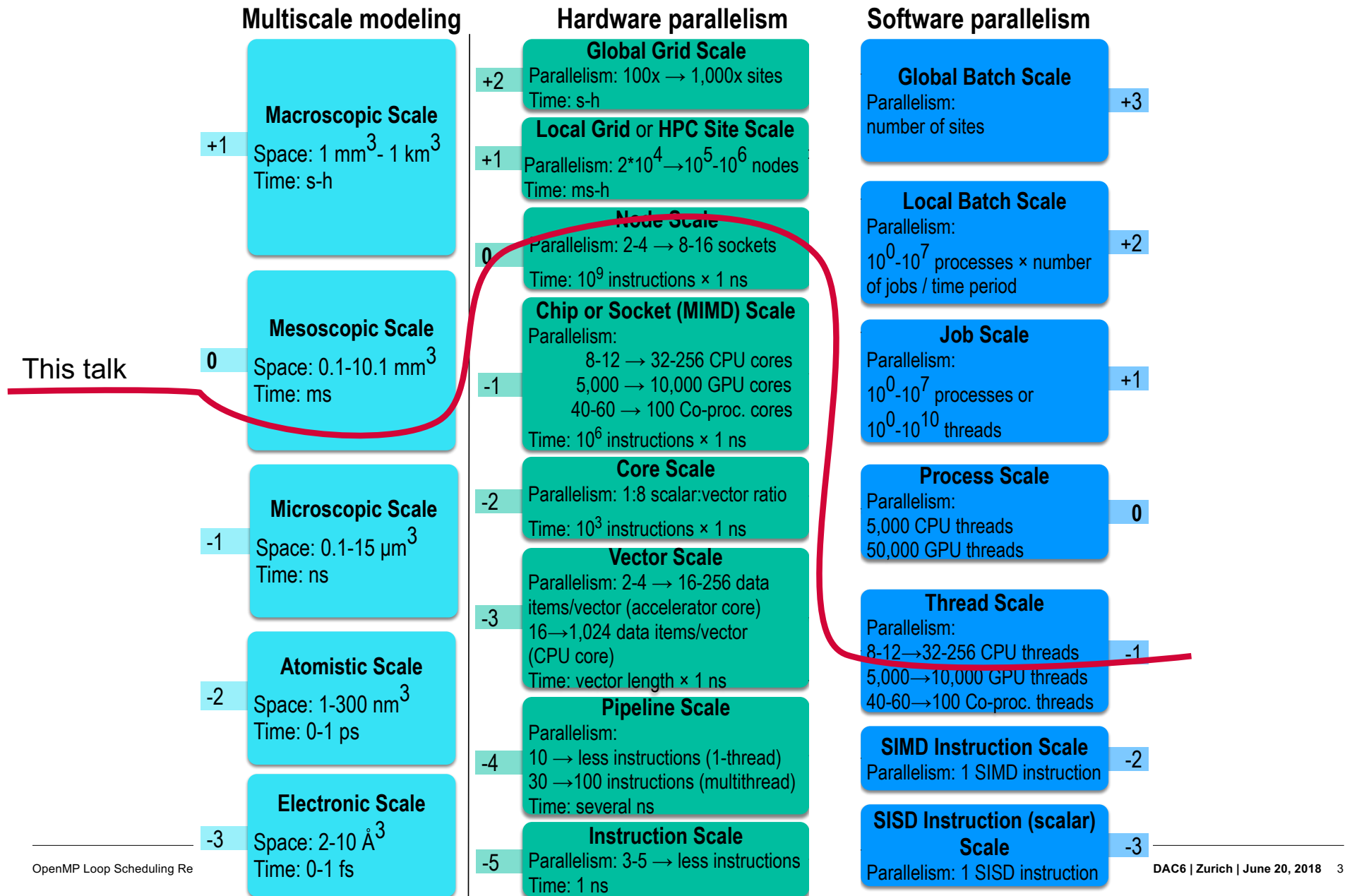
“Scheduling”

- ✧ Scheduling is a vital part of any successful effort of **coordinating** and **managing** parallelism in high performance computers*
 - ✧ Remains a challenge, at several levels, for Exscale computing**
 - ✧ For compute-intensive applications with irregular (nested) parallelism
- ✧ Multiple **types**, **levels**, and **forms** of parallelism
 - ✧ Focus of the SNSF project *Multilevel Scheduling in Large Scale High Performance Computers (2017-2020)*, p3.snf.ch/project-169123

* ETP4HPC SRA2: 5.2 System software (kernel and run-time), 5.3 Prog. env., 5.7 Math. and algo. for extreme scale HPC systems

** IESP 2.0: Runtimes, compilers, applications, algorithms, performance optimization

... To Multiple Types, Levels, and Forms of Parallelism in Parallel Computing



Increasing Hardware Parallelism

- ✧ Through increased node count, CPU core count (multi- and manycore), and accelerator core count

	Piz Daint @ CSCS	SUMMIT @ ORNL	TSUBAME 3.0 @ TokyoTech	
This talk {	CPU cores/node	1×12 (Xeon XC50) 2×18 (Xeon XC40)	2×22 (Power9)	2×14 (Xeon)
	GPU cores/node	1×3,584 (CUDA P100) (XC50)	6×640 (Tensor V100) 6×5,120 (CUDA V100)	4×3,584 (CUDA P100)
Nodes	5,320 (XC50) 1,813 (XC40)	4,608	540	

- ✧ Intel Xeon Phi x200 Knights Landing ≤ 72 CPU cores, 4 hardware threads/core

Parsing the Title

“OpenMP Loop Scheduling”

- ✧ **Loops** typically come to mind in the context of **shared memory** systems
- ✧ **Application** and **underlying system** characteristics determine the best schedule
 - ✧ No “one-size-fits-all” loop scheduling technique can address **all**
 - ✧ Sources of **load imbalance** for
 - ✧ Types of **scientific applications** on
 - ✧ Types of **computing platforms**
- ✧ **OpenMP**: 20+ years industry standard for shared-memory parallel programming
 - ✧ Widely used to parallel program a **broad variety of applications**
 - ✧ Supported by a **growing** number of hardware and software vendors
 - ✧ Several **benchmark suites** for performance evaluation (SPECComp, NAS)
- ✧ **Scheduling**: performance critical aspect of **loops** and important part of most OpenMP programs
 - ✧ Not overshadowed by the introduction of explicit tasks in OpenMP
 - ✧ Nor by the accelerated computing APIs
- ✧ The impact of **system-induced variability** is often neglected in loop scheduling research, particularly by OpenMP schedules

Parsing the Title

“OpenMP Loop Scheduling Revisited”

Are these schedules good enough to efficiently exploit HW parallelism in 2018+?

OpenMP standard **schedule()**

- ✧ **static, chunk**: predetermined allocation order offset by thread ID
- ✧ **dynamic, 1**: pure **self-scheduling SS** [Lusk, Overbeek '83]
- ✧ **dynamic, chunk**: chunk **self-scheduling CSS** [Kruskal, Weiss '85]
- ✧ **guided**: **guided self-scheduling GSS** [Polychronopoulos, Kuck '87]
- ✧ **guided, chunk**: **GSS** with **minimum chunk** size
- ✧ **auto**: implementation determines schedule; no “**chunk**” support

Are there any other schedules not yet in OpenMP?

Are these schedules sufficient for all apps and systems?

YES

Shared-memory self-schedules not in standard

- ✧ **tss**: **trapezoid self-scheduling TSS** [Tzen, Ni '93]
- ✧ **fac2**: practical **factoring FAC** [Flynn Hummel et al. '90-92]
- ✧ **wf2**: practical **weighted factoring WF** [Flynn Hummel et al. '96]
- ✧ **rand**: **random self-scheduling RAND**
- ✧ **taper**: tapering strategy
- ✧ **bold**: bold strategy

Parsing the Title

“Loop Scheduling Revisited”

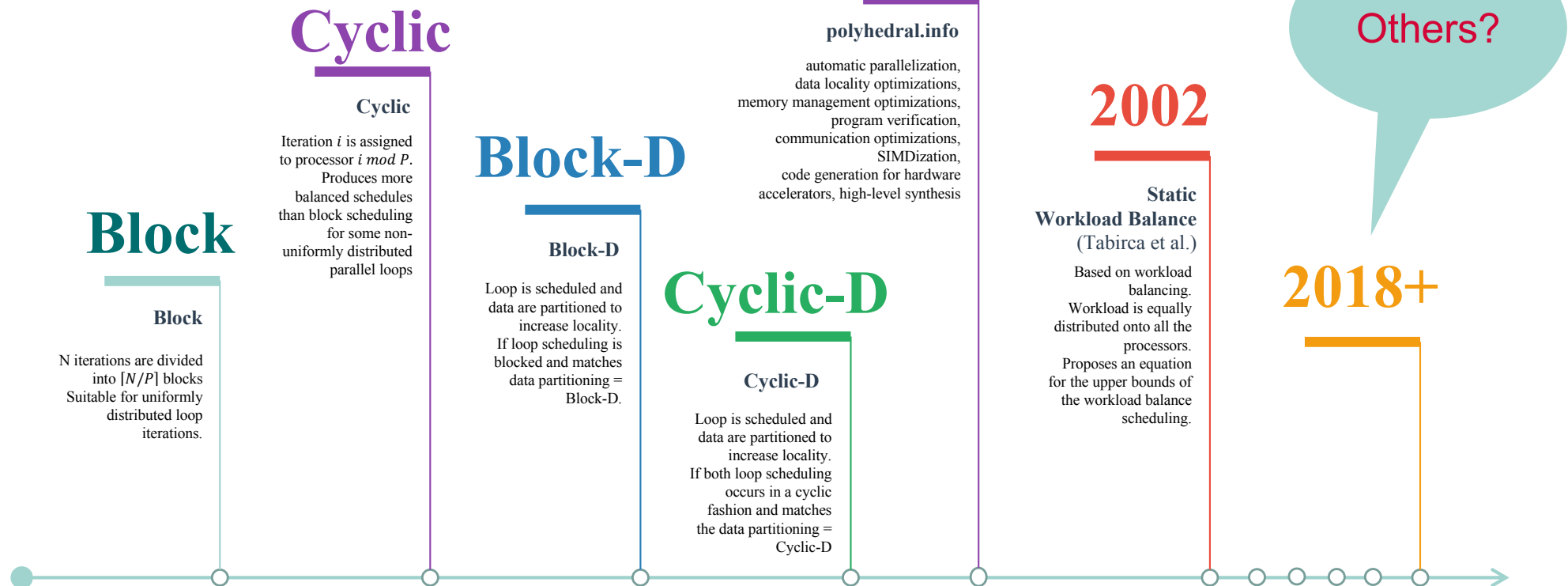
	Scheduling				Work Queue	Data	Optimization Goal	
	Partitioning	Assignment	Load Balancing				Explicit	Implicit
			Ordering	Timing				
Fully static (pre-scheduling)	compilation	compilation	compilation	compilation	central	central distributed	1/2 locality 1/2 scheduling overhead	load imbalance**
Work sharing (static allocation)	compilation execution	compilation	compilation	execution	central	central replicated distributed	1/2 locality 1/2 scheduling overhead	load imbalance**
Affinity & Work stealing	compilation execution	execution	compilation	execution	distributed	central distributed	1/2 locality 1/2 load imbalance**	scheduling overhead
Fully dynamic (self-scheduling)	execution	execution	[compilation] execution	execution	central	central replicated distributed	1/2 scheduling overhead 1/2 load imbalance***	locality

load imbalance** induced by problem and algorithm
load imbalance*** induced by problem, algorithm, and **system**

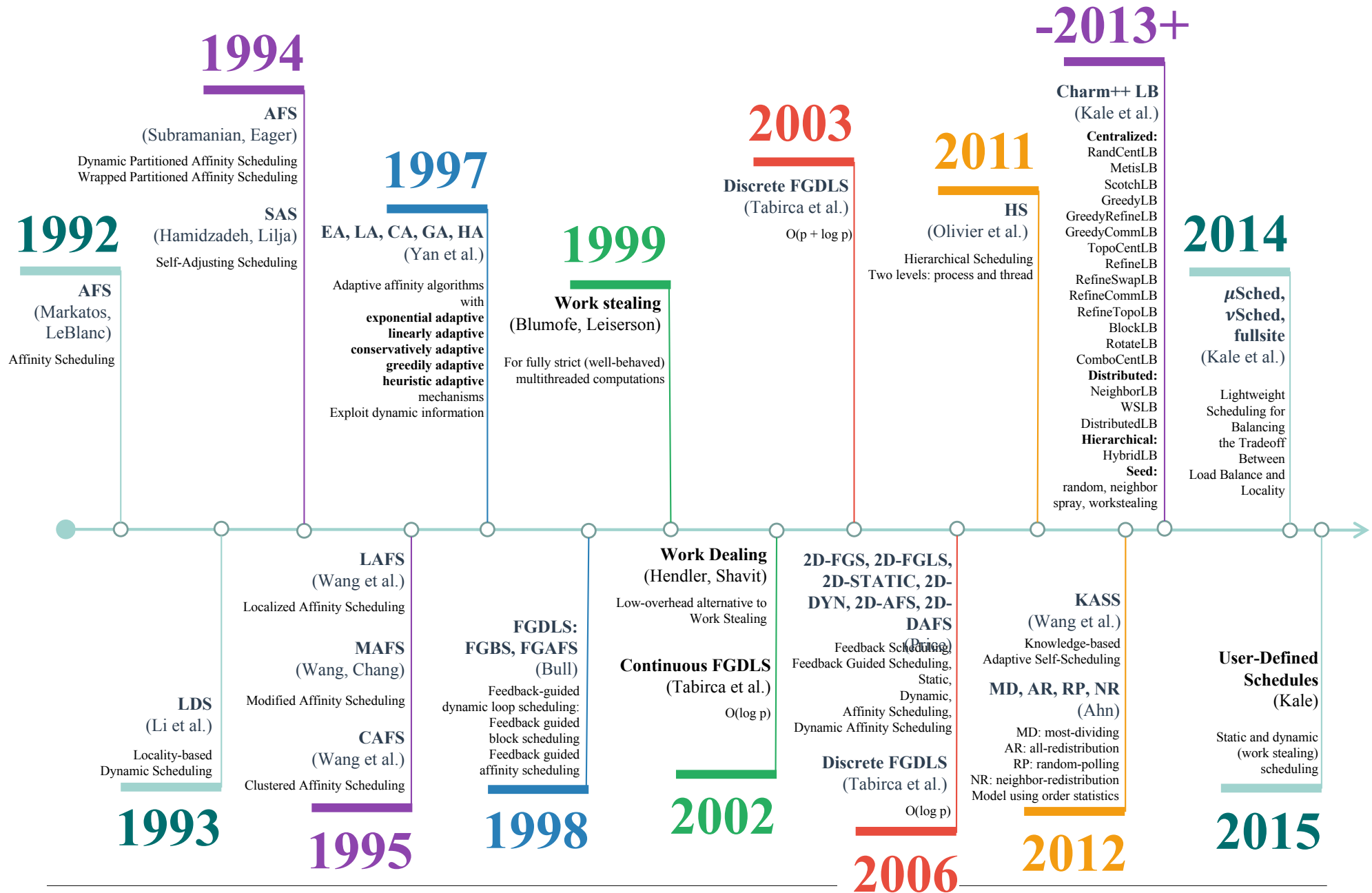
1/2 one goal vs. **explicit trade-off** between
 1/2 another goal two optimization goals

Static Scheduling and Work Sharing

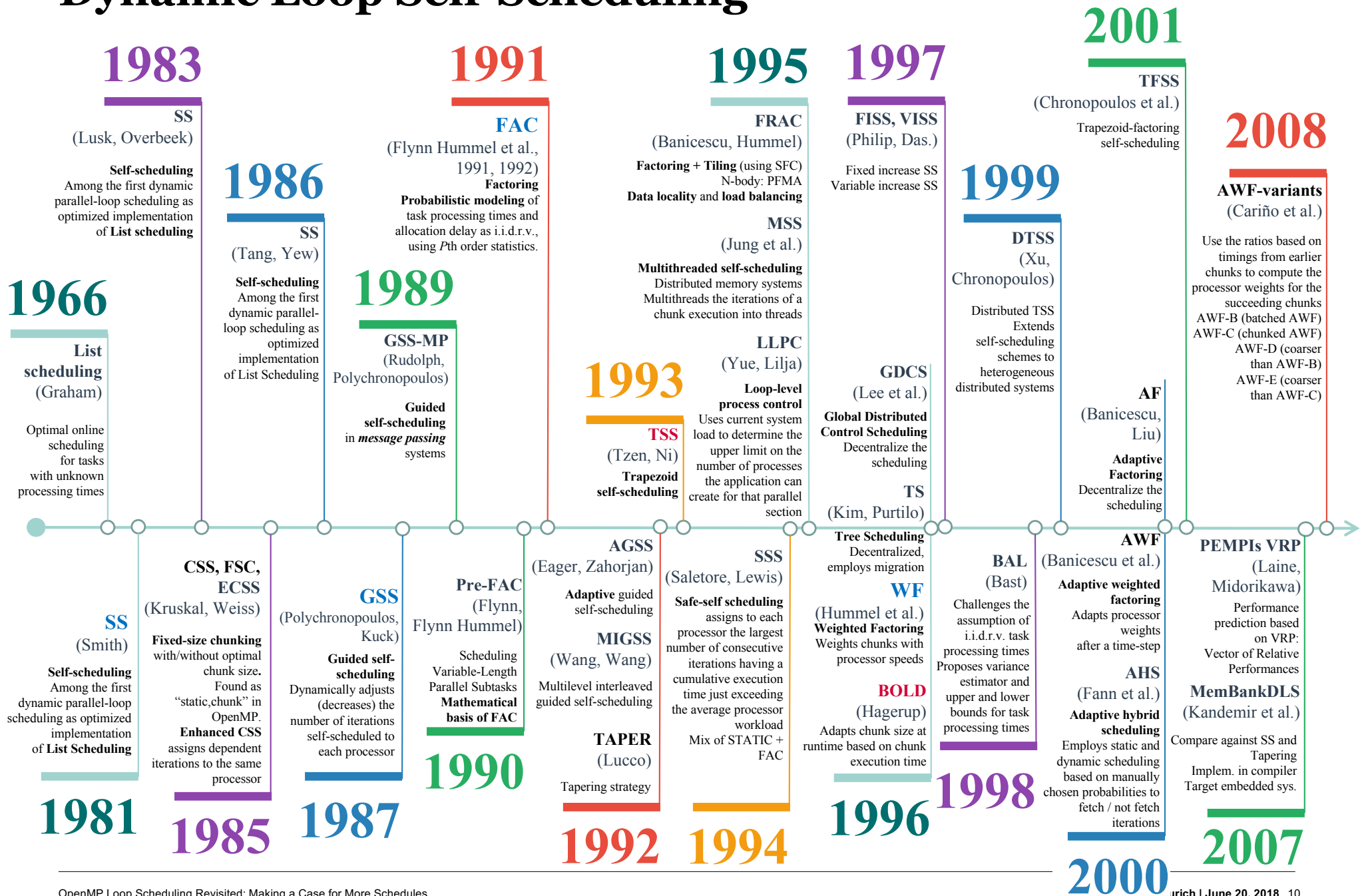
Polyhedral Compilation



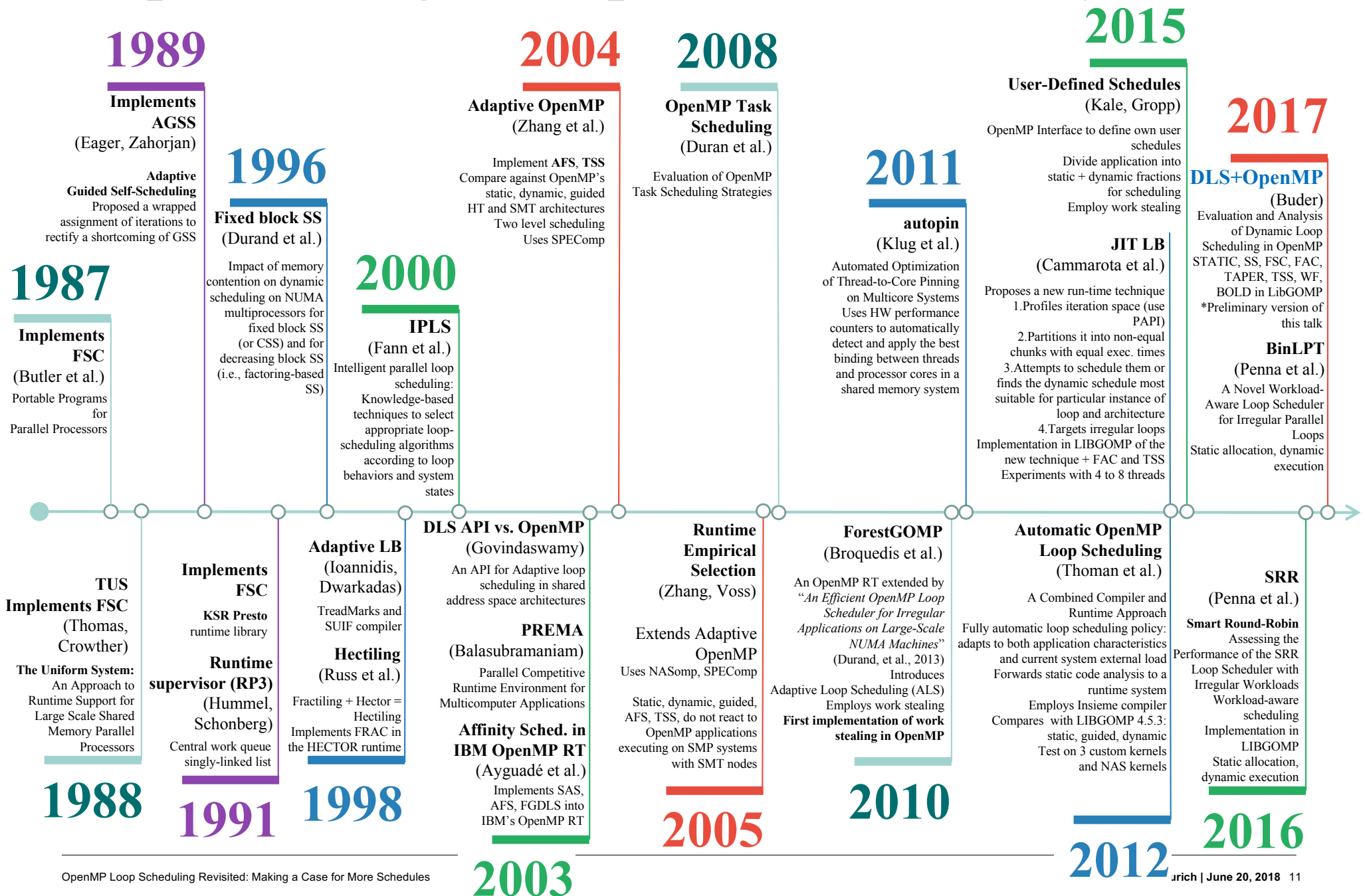
Affinity Scheduling and Work Stealing



Dynamic Loop Self-Scheduling



Loop Scheduling in Compilers and Runtime Systems



Parsing the Title

“More Schedules”

- ✧ OpenMP has not yet adopted **state of the art scheduling** (beyond SS and GSS)
- ✧ Why **more self-scheduling**?
 - ✧ Risk of unexploited parallelism due to increased core counts
 - ✧ **Load imbalance**: problem, application, system (e.g., OS preemption, migration; NUMA effects due to smaller caches / core)
 - ✧ **Central work queue**
 - ✧ Facilitates a dynamic, even distribution of load among processors
 - ✧ Ensures no processor remains idle while there is work to be done
 - ✧ Scalability through hierarchies and distribution
- ✧ **Self-scheduling** places the scheduling responsibility on the **runtime system** rather than on the operating system or the programmer
 - ✧ The **runtime**: optimized for a specific programming model and semantics
 - ✧ The operating system kernel primitives must be general enough to accommodate a variety of programming models and languages
 - ✧ The programmer: not (always) a scheduling expert

Parsing the Title

“Making the Case”

- ✧ One in-house linear algebra kernel
- ✧ Four molecular dynamics codes from various OpenMP benchmark suites
- ✧ Non-uniformly distributed loops \Rightarrow Problem and algorithmic variance

Benchmark suite: code	#LOC	#Parallel Loops	#Iterations	C.O.V. Iterations Exec. Time	Execution time on 1 thread	OpenMP Fraction	Not OpenMP Fraction
Adjoint convolution decreasing task: ac	235	1	10^6	57 %	591.39 s	99.99 %	0.01 %
OpenMP SCR: c_md	384	4	16×10^3	57 %	865.31 s	100 %	0.00 %
RODINIA: lava.md	430	1	13×10^4	14 %	5168.60 s	99.98 %	0.02 %
SPEC OpenMP2012: 350.md	3,701	10	27×10^3	8700 %	98.57 s	97.19 %	2.81 %
NAS OpenMP: MG Class C	1,466	13	10^1 - 10^3	0-1 %	55.70 s	89.04 %	10.96 %

Parsing the Title

“Making the Case”

- ✧ Newly added self-scheduling techniques
 - ✧ **tss**: trapezoid self-scheduling **TSS** ['93]
 - ✧ Collapses to `static`, chunk when first and last chunk equal $\#iterations/\#cores$
 - ✧ **fac2**: practical factoring **FAC** ['90-92]
 - ✧ Unknown mean and stdev of iteration execution times
 - ✧ **wf2**: practical weighted factoring **WF** ['96]
 - ✧ Unknown mean and stdev of iteration execution times
 - ✧ **rand**: random self-scheduling **RAND**
 - ✧ Random chunk $\in [\#iterations/100 \times \#cores, \#iterations/2 \times \#cores]$, $min \geq 1$, $max \geq min+1$
- ✧ Usage via `schedule(runtime)`
- ✧ Implementation into open source OpenMP runtime LaPeSD-libGOMP <https://github.com/lapesd/libgomp>

STATIC	GSS	TSS	FAC2	WF2	RAND	SS	THREADS
25	25	25	13	12	11	1	T0
25	19	22	13	10	2	1	T1
25	14	18	13	15	6	1	T2
25	11	14	13	14	9	1	T3
8	11	7	6	12	1		
6	7	7	5	6	1		
5	3	7	7	12	1		
3		7	6	6	1		
3		4	3	5	1		
2		4	4	4	1		
1		4	3	9	1		
1		4	4	3	1		
1		2	2	2	1		
1		2	2	10	1		
			2	3	1		
			2		1		
			1		1		
			1		1		
			1		1		
						...	
						1	

100 iterations, 4 threads

Parsing the Title

“Making the Case”

- ✧ miniHPC: Fully-controlled 22-node system used for research and teaching

miniHPC node	Characteristic
Sockets	2
Processor	Intel Xeon CPU E5-2650 v4
Clock speed	2.40GHz
Architecture	x86_64
L1D cache	32KB
L1I cache	32KB
L2 cache	256KB
L3 cache	25600KB
RAM	64GB
Physical CPU cores	20
HT CPU cores	40

Parsing the Title

“Making the Case”

✧ Executed five benchmarks with their OpenMP loops scheduled using **20 threads** with

✧ **STATIC, SS, GSS**

✧ **TSS, FAC2, WF, RAND**

✧ Originally: no schedule

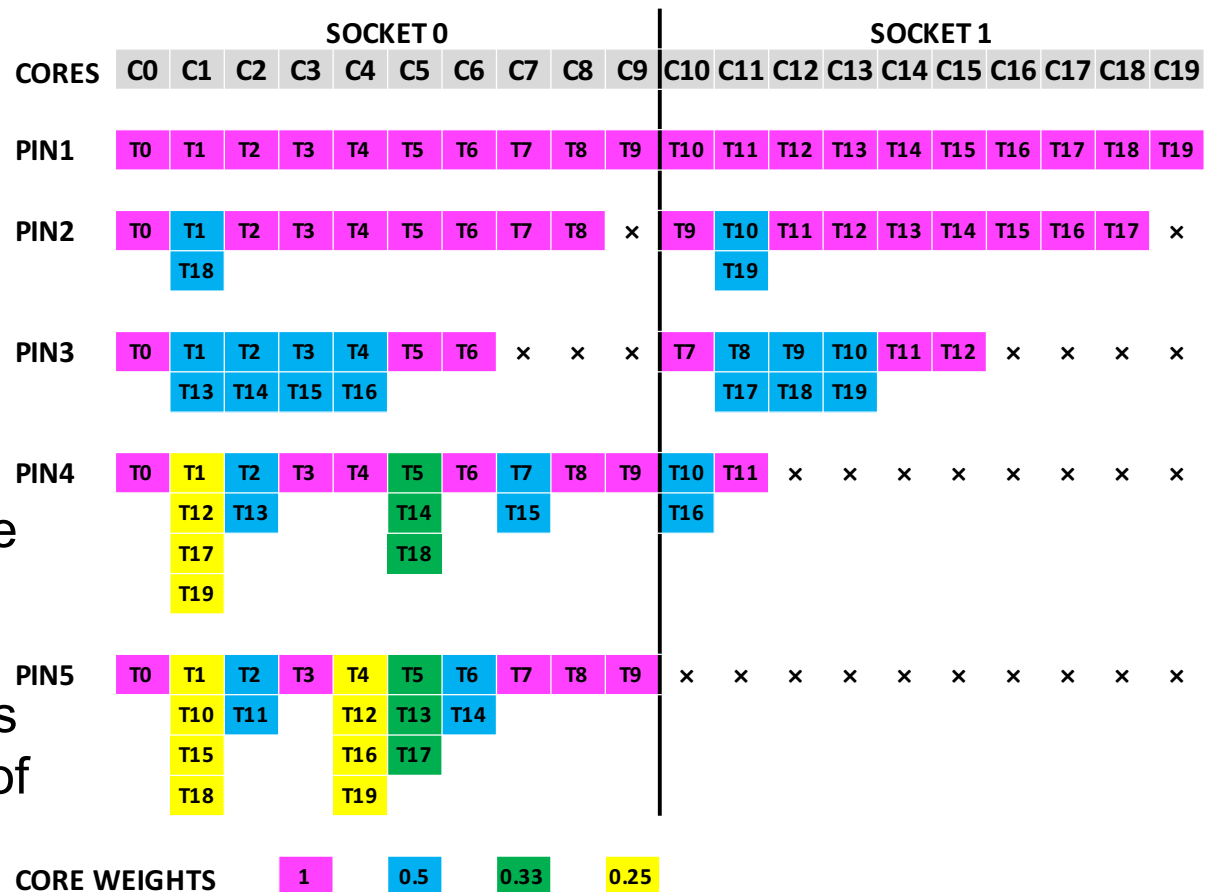
✧ System-induced load imbalance

✧ Five pinning strategies

✧ Parallel execution time statistics (median and stdev) of 20 runs of each experiment

✧ Does a schedule benefit a **parallel loop?**

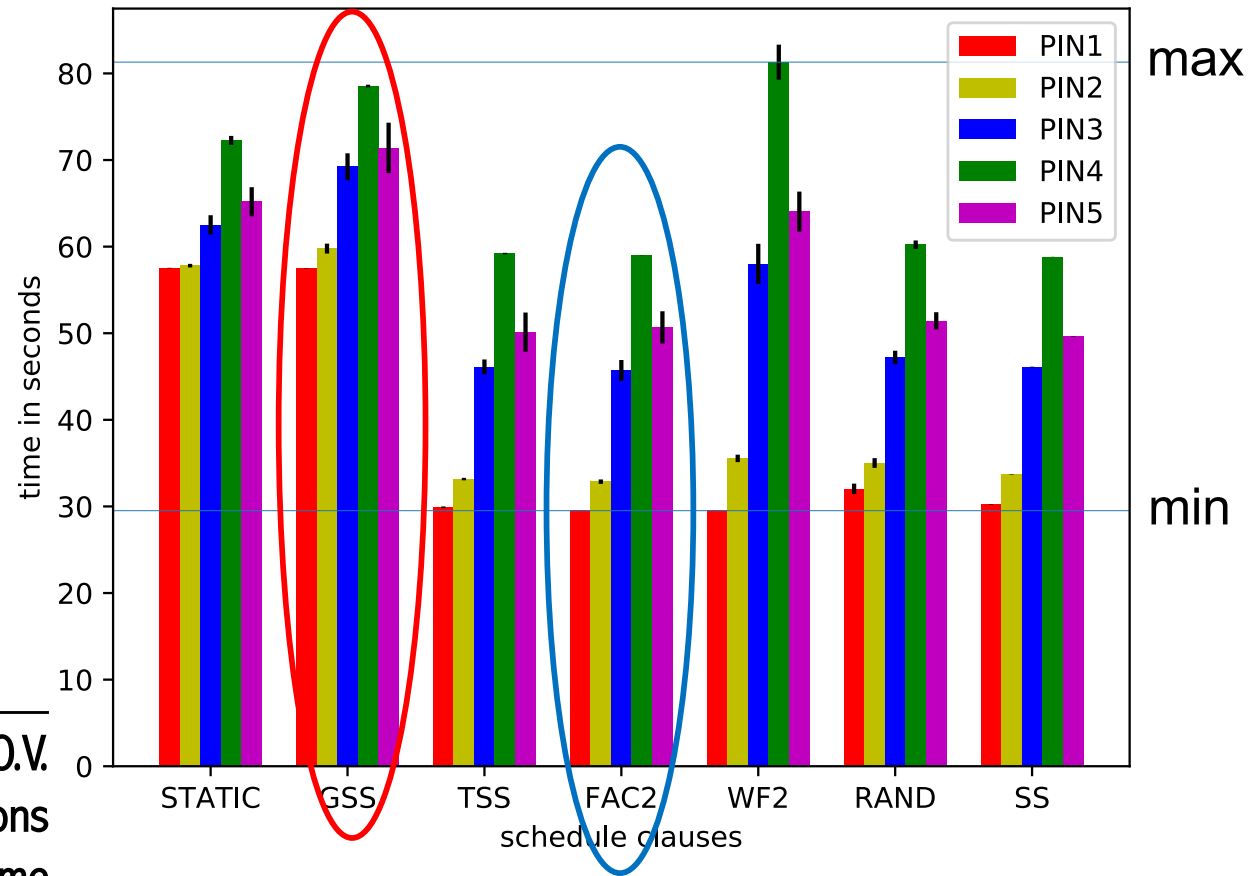
✧ Can it handle **HW heterogeneity?**



Parsing the Title

“Making the Case”

ac

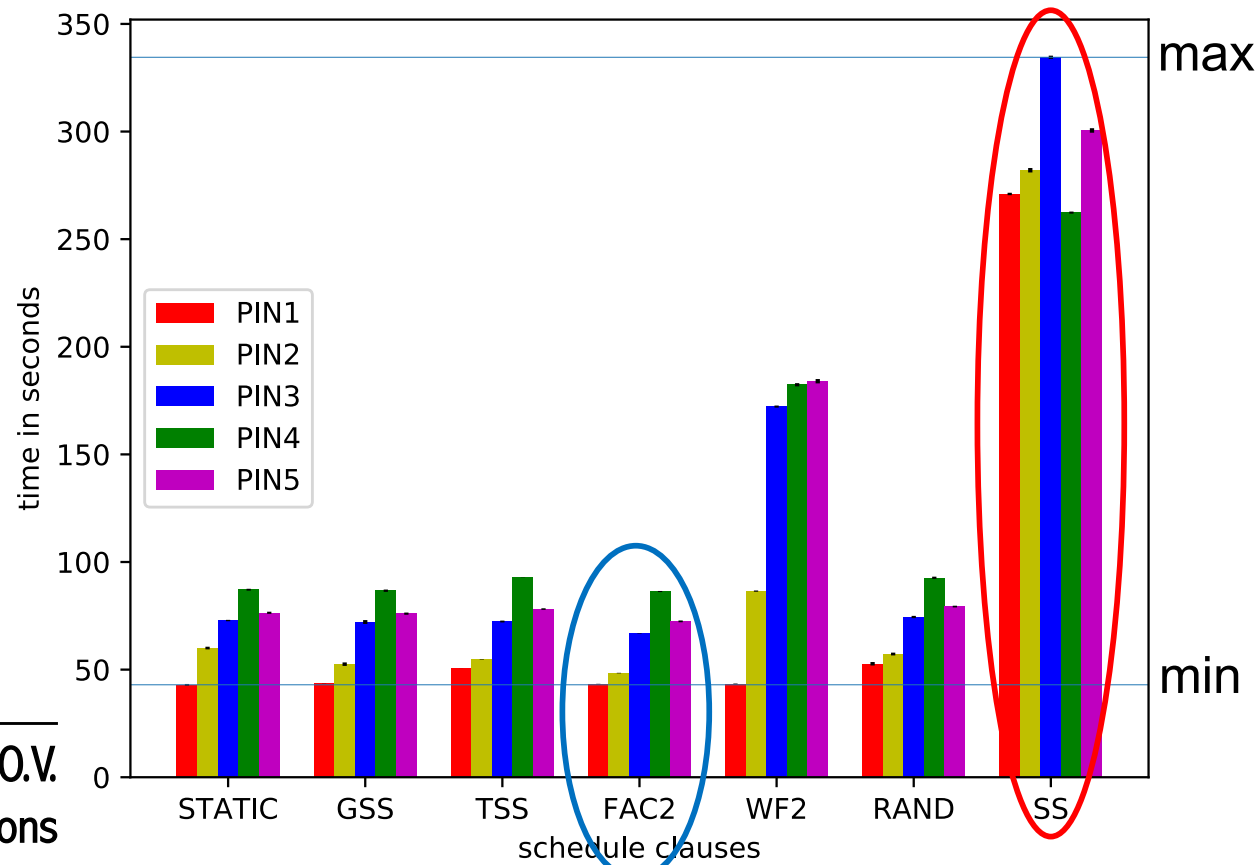


#Iterations	C.O.V.
10 ⁶	57 %
Iterations	
Exec. Time	

Parsing the Title

“Making the Case”

c_md

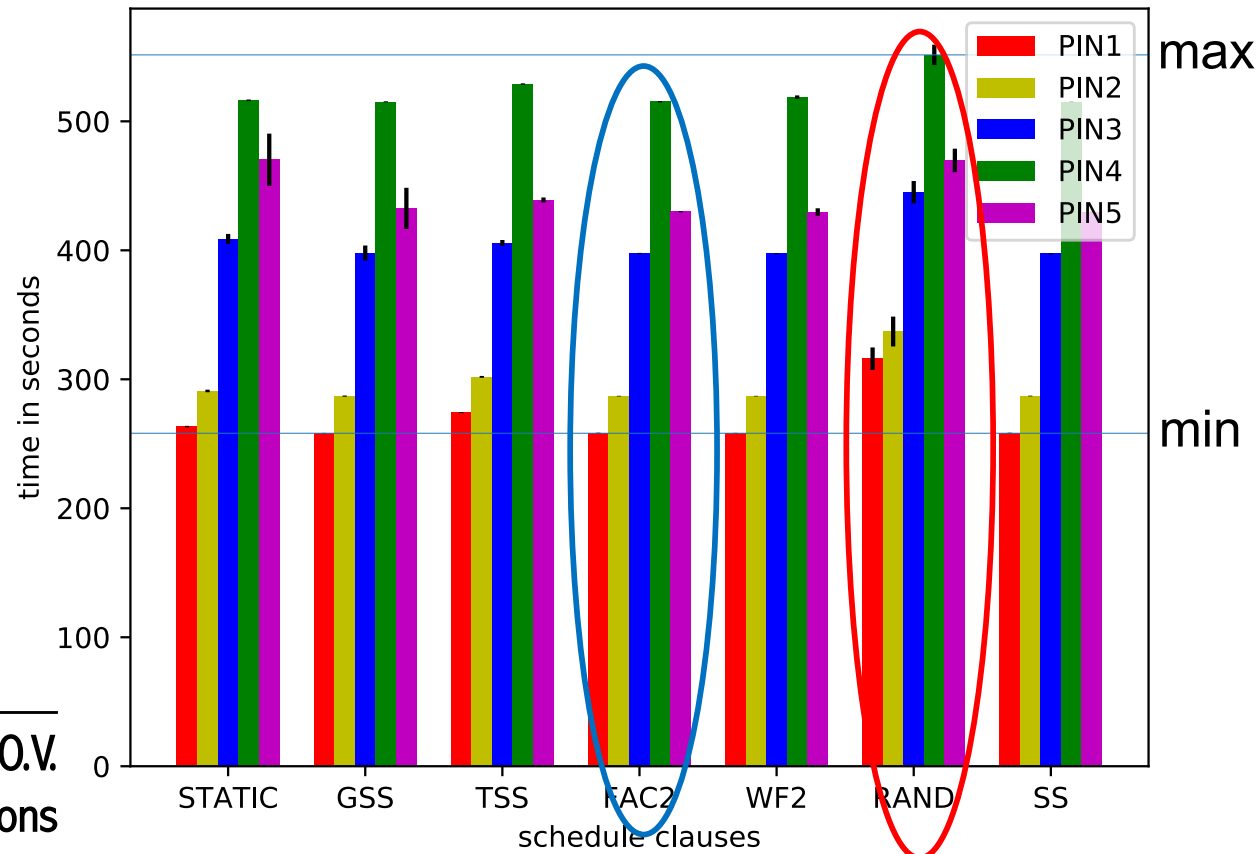


#Iterations	C.O.V. Iterations Exec. Time
16×10^3	57 %

Parsing the Title

“Making the Case”

lava.md

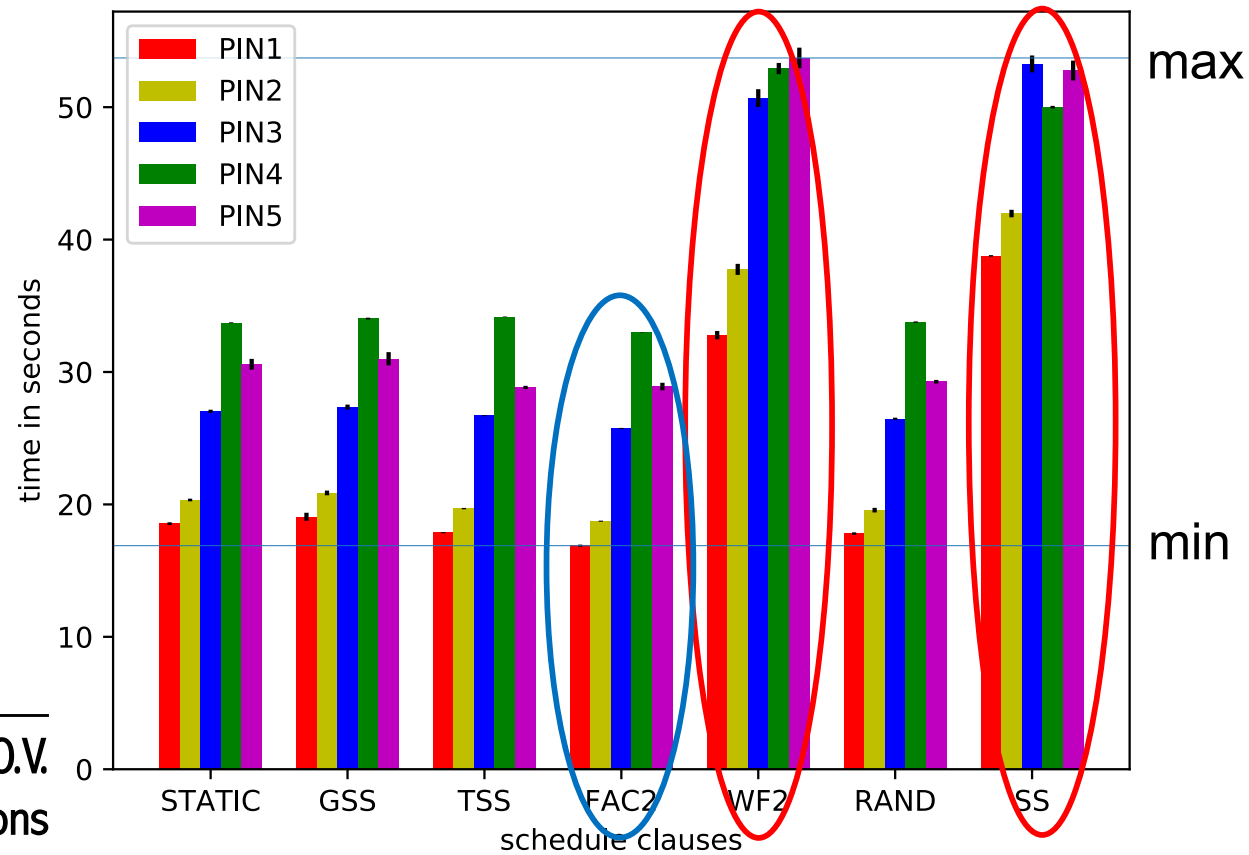


#Iterations	C.O.V. Iterations Exec. Time
13×10^4	14 %

Parsing the Title

“Making the Case”

350.md

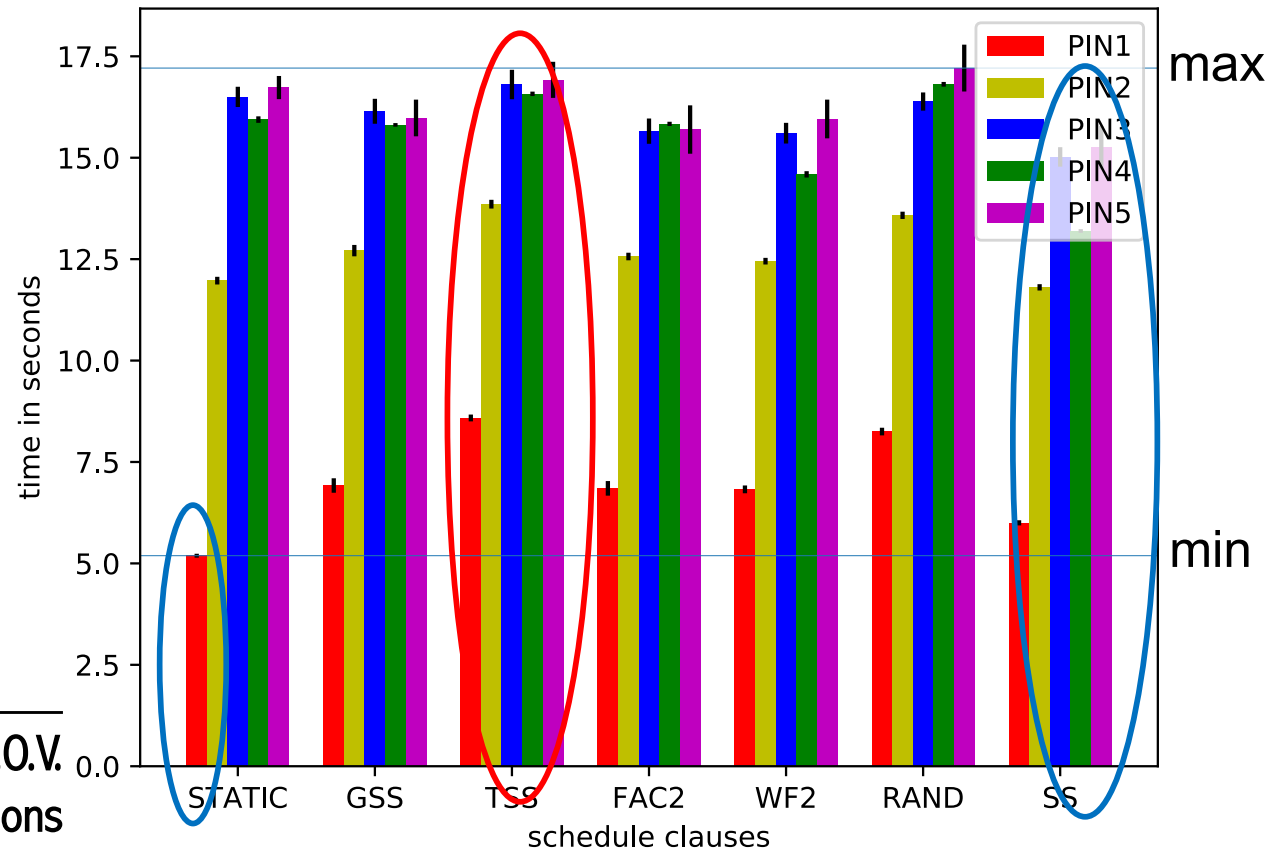


#Iterations	C.O.V.
Iterations	Exec. Time
27×10^3	8700 %

Parsing the Title

“Making the Case”

MG



C.O.V.

#Iterations Iterations
(Class C) Exec. Time

2-10³ 0-1 %

To Use or Not to Use Dynamic Loop Self-Scheduling?

No “One-Size-Fits-All”, Wide Gap Between Best and Worst

- ✧ Additional schedules provide benefit over existing schedules
- ✧ When application and system parallelism is **regular**, **STATIC is sufficient**
- ✧ When the **#iterations** is **too small** to generate enough work and when the **#threads** is **large**, then **STATIC is sufficient**
- ✧ When the cost of allocating loop iterations to a thread is **larger** than the cost to execute the loop iterations then **dynamic loop scheduling is not beneficial**
 - ✧ Static and affinity-based methods can be used instead
- ✧ In the other cases
 - ✧ High compute intensity
 - ✧ Nested and irregular parallelism
 - ✧ System-induced variabilities (e.g., OS, NUMA)**dynamic loop scheduling is needed** and **self-scheduling offers benefits** over affinity- and work stealing-based methods

So What?

Advantage

- ✧ The newly implemented DLS are **immediately usable** by existing programs using our non-standard prototype implementation via `schedule(runtime)`
 - ✧ Numerous OpenMP production codes in active use
 - ✧ Numerous multi/manycore platforms available
 - ✧ <https://bitbucket.org/PatrickABuder/libgomp/src>

Usefulness

- ✧ On heterogeneous platforms
 - ✧ Multi/manycore CPUs
 - ✧ Fat cores or faster connected cores self-schedule more frequently
 - ✧ Thin cores or slower connected cores self-schedule more rarely
 - ✧ Multi/manycore CPUs and accelerator cores
 - ✧ `schedule(runtime)` for the CPU threads
 - ✧ `dist_schedule(static, chunk)` with `schedule(runtime[, chunk])` for target teams and their threads on accelerator cores



Thank you for
your attention!

Now What?

- ✧ Advocate for the inclusion of **more self-scheduling techniques** into the OpenMP standard or as an interface for user-defined schedulers
 - ✧ To address **all sources of load imbalance** (problem, algorithmic, systemic) during execution
 - ✧ Runtime should exploit user expert knowledge about the application
 - ✧ Global OpenMP task scheduling still unaddressed
- ✧ Implement further **state-of-the-art loop self-scheduling** techniques (with feedback loops) into LLVM/Clang
- ✧ Extend the proof of concept beyond benchmarks into **real applications**
 - ✧ Combine with self-scheduling in MPI layer
 - ✧ PASC project *SPH-EXA*, www.pasc-ch.org/projects/2017-2020/sph-exa/
- ✧ Implement an **intelligent selection mechanism** among the many available options, based on previous work [Boulmier *et al.* 2017; Banicescu *et al.* 2013]