#### **Tracing Performance Tools Back 25 Years**



TECHNISCHE UNIVERSITÄT DARMSTADT

Felix Wolf, Technical University of Darmstadt

25th Anniversary of APART



2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 1

Photo: Alex Becker / TU Darmstadt



**TECHNISCHE** 

UNIVERSITÄT DARMSTADT

#### Starting point – peak performance gap

2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 2

### Event traces - too large to analyze manually



- Needles-in-the-hay-stack problem
- Low abstraction level programmatic access needed



https://www.vampir.eu



### EARL - Event Analysis and Recognition Language

#### Extended Python interpreter

- Also Tcl, Perl interfaces
- Pattern specified in extended scripting language

#### Services

- Efficient random access to events
  - VAMPIR, ALOG, CLOG trace data formats
- Mapping of events to EARL event trace model

Abstractions defined in the EARL event trace model

• Permit use of simple pattern search algorithm





### EARL - Event Trace Model

**Basic Model** 

- Sequence of events:  $E = e_1, ..., e_n$
- Event types: enter, exit, send, recv
- Event attributes: time, loc, region, src, dest

#### Abstractions

- Pointer attributes connecting related events
  - enterptr : links event to entering of enclosing region instance
  - sendptr : links message receipt to message dispatch
- System states mapping events to execution context (set of events)
  - Region stack: set of enter events belonging to open region instances
  - Message queue: set of send events belonging to messages in transfer



### Example: Late Sender

MPI\_RECV is posted before corresponding MPI\_SEND

• Determine time during which receiver sits idle





#### Pattern Class: LateSender

```
class LateSender (Pattern) :
  [...]
 def recv callback(self, recv msg):
    enter recv = self.trace.event(recv msg[`enterptr'])
    send msg = self.trace.event(recv msg[`sendptr'])
    enter send = self.trace.event(send msg[`enterptr'])
    idle time = enter send[`time'] - enter recv[`time']
    if (idle time > 0 and
        enter send['region'] == "MPI SEND" and
        enter recv['region'] == "MPI RECV"):
      self.sum idle time = self.sum idle time + idle time
  [...]
 def confidence(self): return 1
 def severity(self): return self.sum idle time
```



### EXPERT - Extensible Performance Tool





### EXPERT Output

tatus	Pattern Region	
- File In	formation	
File: 0	cx 8x1.bpv	
Procs	: 8	
- Messa	ages	
MPT E	xmert 1 0 The Jul 25 12:30:54 2000	
Analu	ris of "cx 8x1 how" terminated successfully	
Searc	h focus: "WELD"	
Jearc		
Con	fidence: 1	
Res	erity: 0.035 ult:	
D	istribution : [0.17, 0.08, 0.01, 0.01, 0.01, 0.01, 0.05, 0.68]	



### Test Case: Czochralski Crystal Growth

Convection processes in a rotating cylindric crucible

- 3 dimensional cubical mesh
- 2 dimensional spatial decomposition

Most work is done in routine VELO

Calculating new velocity vectors

Patterns used for analysis

- Communication costs
- Late sender







### Idle times in VELO caused by Late Sender







#### Next generation GUI



### History

### 1998

- Start at the ZAM (Michael Gerndt and Bernd Mohr)
- Manifesto

M. Gerndt, B. Mohr, M. Pantano, F. Wolf: Automatic Performance Analysis for Cray T3E, Proceedings of the 7<sup>th</sup> Workshop on Compilers for Parallel Computers (CPC 98), University of Linköping, Sweden, June-July 1998

- First component: EARL trace-analysis toolkit (Diploma thesis, ZAM)
- 2000
  - First prototype of an automatic trace analyzer (EXPERT)
- 2001
  - Automatic OpenMP instrumentation (OPARI, POMP)
  - First prototype of the EPILOG tracing library
  - First prototype of the KOJAK GUI



13

# History (2)

### 2003

- KOJAK became collaboration between Forschungszentrum Jülich and the University of Tennessee
- First beta release of KOJAK
- 2004
  - Release 1.0 and 2.0b
  - New GUI component CUBE
- 6 Diploma theses
- 1 Ph.D. thesis
- 22 publications
  - 3 journal articles
  - 19 conference and workshop articles





### Automatic performance analysis

### Automatic performance analysis



- Take event traces of MPI/OpenMP applications
- Search for execution patterns
- Calculate high-level call path profile
  - Problem, call path, system location  $\Rightarrow$  time
- Display in performance browser





### **Typical performance problems**







### **Pattern hierarchy**





17



### **CUBE Performance algebra**

### Difference operator

- Compare different experiments
- Merge operator
  - Integrate performance data from multiple sources
- Mean operator
  - Summarize a series of experiments



- Obtain new CUBE instance as result
- Display it like ordinary CUBE instance





### (Re)moving waiting times

- Difference between before / after barrier removal
- Raised relief shows improvement
- Sunken relief shows degradation







# Virtual topology in SWEEP3D

- Three-dimensional domain (i,j,k)
- Two-dimensional domain decomposition (i,j)





- Wave-fronts from different directions
- Limited parallelism upon pipeline refill (late sender)

# Four new patterns

Refill from NW, NE, SE, SW

• Late sender combined with change of message direction

- 0 3

>

(7, 7)

- Topological knowledge needed to recognize direction change
- Caveat: message direction has two components
  - Special case: border processes



# Analysis process





# Trace size limits scalability

- Serially analyzing a single and potentially large global trace file does not scale to 1000s of processors
- Even if locality is exploited, main memory might be insufficient to store current working set
- Amount of trace data might not fit into single file













in der Helmholtz-Gemeinschaft



# Current prototype







# scalasca 🗖

- Scalable performance-analysis toolset for parallel codes
- Integrated performance analysis process
  - Performance overview on call-path level via runtime summarization
  - In-depth study of application behavior via event tracing
  - Switching between both options without recompilation or relinking
- Supported programming models
  - MPI-1, MPI-2 one-sided communication
  - OpenMP (basic features)
- Available under the New BSD open-source license
  - <u>http://www.scalasca.org/</u>

Joint project of



# **The team (2010)**





































### Scalability of parallel wait-state search (SWEEP3D)



sweep3d, jugene\_vn, scalasca-1.2

### Sweep3D Late sender

• Grid of 576 x 512 processes

>	Cube 3.3	Qt: epik_sw	eep3d_vn2	)4912_trace/trace.cube.gz 🥯	_ 🗆 🗡	٢
[	Peer distributi	on			•	•
	System tree	Topology 0	Topology 1			
				KAAX IIIIIII		2
					$\geq$	
					$\langle \rangle$	
	•	//		0.00	100.00	0
-	27.50			36.40 +/- 3.03	45.3	1

### **Redundant messages in XNS CFD code**



# **Delay analysis**



INTERNATIONAL CONFERENCE ON PARALLEL PROCESSING

- Delay counterpart of waiting time
- Distinction between direct and indirect waiting times
- Essentially scalable version of Meira Jr. et al.
- Analysis attributes costs of wait states to delay intervals
  - Requires backward replay



## **Origin of delay costs in Zeus-MP/2**



Computation





Delay costs

## **Delay analysis of code Illumination**

- Particle physics code (laser-plasma interaction)
- Delay analysis identified inefficient communication behavior as cause of wait states



Computation

Propagating wait states: Original vs. optimized code Costs of direct delay in optimized code



### Wait-state analysis of CICE





Indirect waiting time



Direct waiting time



#### Delay costs

Method: bi-directional replay of event traces

Illustrated propagation of wait states as a result of suboptimal load balancing

Motivated load-balancing simulator

#### The virtual institute in a...



• Partnership to develop advanced programming tools for complex simulation codes

VI-HPS

- Goals
  - Improve code quality
  - Speed up development
- Activities
  - Tool development and integration
  - Training
  - Support
  - Academic workshops
- www.vi-hps.org

### **Score-P** measurement system

Vampir	Sca	alasca		TAU			Periscope	
Interactive trace exploration	Perfo dyna wait	formance namics & nit states		Performance data base & data mining			Automatic online classification	
Tracing		Profiling			С	Online interface		
Score-P measurement infrastructure								
Application (MPI, OpenMP, accelerator, PGAS, hybrid)								











UNIVERSITY OF OREGON



2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 38

### **Performance model**



Formula that expresses a relevant performance metric as a function of one or more execution parameters



### **Empirical performance modeling**





2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 40

### Challenges





#### Run-to-run variation / noise





#### Cost of the required experiments

### How to deal with noisy data



- Introduce prior into learning process
  - Assumption about the probability distribution generating the data



### **Performance model normal form (PMNF)**





2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 43

#### 2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 44

#### Available at: https://github.com/extra-p/extrap



#### Extra-P 4.0



How many data points do we really need?





### Learning cost-effective sampling strategies [Ritter et al., IPDPS'20]



TECHNISCHE UNIVERSITÄT

DARMSTADT

2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 46

### **Case studies**





Applicatio	n	#Parameters	Extra points	Cost savings [%]	Prediction error [%]
FASTEST		2	0	70	2
Kripke		3	3	99	39
Relearn		2	0	85	11

2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 47

### **Optimized measurement point selection**





### **Optimized measurement point selection**

via Gaussian Process Regression (GPR)

TECHNISCHE UNIVERSITÄT DARMSTADT

Idea: Use covariance function



2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 49

### **Parameter selection**

[Copik et al, PPoPP'21]



- The more parameters the more experiments
- Modeling parameters without performance impact is harmful



### Case study – LULESH & MILC Influence of program parameters



TECHNISCHE UNIVERSITÄT DARMSTADT

LULESH	Total	р	size	regions	iters	balance		cost	p, size
Functions	349	2	40	15	1	1	2		40
Loops	275	2	78	29	1	1	2		78
MILC	Total	р	size	trajecs	warms steps	nrest. niter	mass, beta nfl.	u0	p, size
Functions	621	54	53	12	9	6	1	4	56
Loops	874	187	161	39	31	15	1	7	196

2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 51

#### **Relearn** [Rinke et. al, JPDC'18]





Scalable algorithm to simulate structural plasticity in the brain
Adaptation of Barnes-Hut algorithm (astrophysics)



2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 52

### **Complexity consideration**



- O(n \* log2n) = O(n \* log n \* log n):
  - Barnes–Hut at every level: O(n \* log n)
  - Tree depth: O(log n)
- Parallel complexity:
   O(n/p \* log<sup>2</sup> n + p)



#### Media coverage



#### TV - SAT.1 Regionalmagazin für Rheinland-Pfalz und Hessen



https://www.1730live.de/wissenschaftler-wollen-gehirn-nachbauen/





#6 of 2019

#### **Execution time** One connectivity update



TECHNISCHE UNIVERSITÄT DARMSTADT



2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 55

#### **Sharp complexity for scalable algorithm** [Czappa et al, JPDC'23]



- Depending on acceptance criterion Θ, subsequent Barnes–Hut steps have constant complexity
  - $O(n * (\log n + \log n)) = O(n * \log n)$
- Overall complexity for practical scenarios:
   O(n/p \* log n + p)

#### **FTIO: Frequency techniques for I/O** [Tarraf et al., IPDPS '24]

- Captures the period of I/O phases
- Operates in the **frequency domain**
- Quantifies the **confidence** in the results
- Online (prediction) w/ low overhead and offline (detection)





**FTIO** 



### Use case: I/O scheduling with IO-Sets



- Classify applications based on their period
  - Shared file-system access to different classes
  - Mutually exclusive access to individual jobs within the same class



### Thank you!





2/12/2024 | Technical University of Darmstadt, Germany | Felix Wolf | 59

### Literature



ΤοοΙ	Paper
KOJAK/ EXPERT	<ul> <li>Michael Gerndt, Bernd Mohr, Felix Wolf, Mario Pantano: Performance Analysis on Cray T3E. In <i>Proc.</i> of the 7th Euromicro Workshop on Parallel and Distributed Processing (PDP), Funchal, Madeira, Portugal, pages 241–248, IEEE, February 1999.</li> <li>Felix Wolf, Bernd Mohr: Automatic performance analysis of hybrid MPI/OpenMP applications. Journal of Systems Architecture, 49(10-11):421–439, November 2003.</li> </ul>
Scalasca	<ul> <li>Markus Geimer, Felix Wolf, Brian J. N. Wylie, Bernd Mohr: A scalable tool architecture for diagnosing wait states in massively parallel applications. Parallel Computing, 35(7):375–388, July 2009.</li> <li>Markus Geimer, Felix Wolf, Brian J. N. Wylie, Erika Ábrahám, Daniel Becker, Bernd Mohr: The Scalasca performance toolset architecture. Concurrency and Computation: Practice and Experience, 22(6):702–719, April 2010.</li> <li>David Böhme, Markus Geimer, Lukas Arnold, Felix Voigtländer, Felix Wolf: Identifying the root causes of wait states in large-scale parallel applications. ACM Transactions on Parallel Computing, 3(2):Article No. 11, 24 pages, July 2016.</li> </ul>
Extra-P	<ul> <li>Alexandru Calotoiu, Torsten Hoefler, Marius Poke, Felix Wolf: Using Automated Performance Modeling to Find Scalability Bugs in Complex Codes. In Proc. of the ACM/IEEE Conference on Supercomputing (SC13), Denver, CO, USA, pages 1–12, ACM, November 2013.</li> <li>Sergei Shudler, Yannick Berens, Alexandru Calotoiu, Torsten Hoefler, Alexandre Strube, Felix Wolf: Engineering Algorithms for Scalability through Continuous Validation of Performance Expectations. IEEE Transactions on Parallel and Distributed Systems, 30(8):1768–1785, August 2019.</li> </ul>
FTIO	<ul> <li>Ahmad Tarraf, Alexis Bandet, Francieli Boito, Guillaume Pallez, Felix Wolf: Capturing Periodic I/O Using Frequency Techniques. In Proc. of the 38th IEEE International Parallel and Distributed Processing Symposium (IPDPS), San Francisco, CA, USA, pages 1–14, IEEE, May 2024, (accepted)</li> </ul>

### Acknowledgment



Alexis Bandet Nikhil Battia Alexandru Calotoiu **Daniel Becker** Francieli Boito David Böhme **Dominic Eschweiler** Marcin Copik Jack Dongarra Christian Feld Wolfgang Frings Markus Geimer Alexander Geiß Michael Gerndt Marc-André Hermanns Torsten Hoefler Michael Knobloch

Daniel Lorenz Bernd Mohr Shirley Moore Guillaume Pallez Farzona Pulatova Sebastian Rinke Marcus Ritter Pavel Saviankou Martin Schulz **Christian Siebert** Sergei Shudler Fengguan Song Zoltán Szebenyi Ahmad Tarraf **Brian Wiley** [...]



#### HESSEN



Hessisches Ministerium für Wissenschaft und Kunst



Bundesministerium für Bildung und Forschung









Office of Science

Swiss National Science Foundation