

Markus Geveler

@ PACO17 Tegernsee



Markus @ TUDo

www.mathematik.tu-dortmund.de/lsiii/_Markus_Geveler

markus.geveler@math.tu-dortmund.de

Energy-efficient simulations, Hardware-oriented numerics, Multilevel domain decomposition solvers



I.C.A.R.U.S.

Ministerium für Innovation, Wissenschaft und Forschung des Landes Nordrhein-Westfalen

www.icarus-green-hpc.org

@ TU Dortmund:

Insular Compute center Applied Sciences Renewables Unconventional System Integration Markus Geveler, Stefan Turek, Dirk Ribbrock, Hannes Ruelmann, Daniel Donner

@ MPI Magdeburg

Peter Benner, Jens Saak, Martin Köhler, Gerry Truschkewitz



Where everything is leading: simulation of technical flow

- Multiphysics: increasing complexity of problems and thus methods
- Models are versatile: applications in
- **Discrete continuum mechanics:** determine physical quantities for a huge number of points in space and repeat
- **Ressource-hungry:** memory, time, *energy*
- Methods: DD Newton-Krylov-Multigrid



Unconventional hardware for scientific computing

ARM-based SoCs with mobile GPUs plus Photovoltaic in a (big) box





History: FEAT and low power SoCs: embedded v commodity

- exploration of energy savings / performance tradeoffs using ARM processors on Tegra 2 SoCs (Cortex Ag only)
- results were very promising
- 2013: ICARUS accepted, Tegra K1 announced (GPU!)
- 2013 2014: preliminary experiments with Tegra 3, Tegra 4, Tegra K1 for ICARUS
- 2015: small cluster with 4 nodes @ PACO15
- **2015 march 2016:** ICARUS construction

Applications on ARM Cortex A15





6

technische universität dortmund fakultät für **m**







Computers and energy consumption: a (pessimistic) forecast



66

Conventional approaches are running into physical limits. Reducing the 'energy cost' of managing data on-chip requires coordinated research in new materials, devices, and architectures.

- the Semiconductor Industry Association (SIA), 2015



Technology revolutions we can use in scientific computing

	1980ies	1990ies	2000s	2010s	2020ies		
	Digital computing	/ signal processing /	embedded systems				
computing	In	ternet					
		Ν	lobile Internet				
			Mobile	Computing			
power supply	3. Indus	3. Industryrevolution (Automation, IT, FabLab)					
		Ene	rgy (renewables, bet	ter nets and consur	ners)		
				Traffic (elec	ctro)		
				AI			
					?		



Single devices, clusters and energy efficiency

- mobile processors: less power, less performance than commodity (x86, ..., desktop GPUs)
- same task, more devices of same type: more power, more performance
- Can we scale to same performance and spend less?
- scaling and energy: less execution time, same energy (because E is P integrated in time)

For two computer architectures A and B must hold: If A is more energy-efficient than B in executing a task then the powerdown from A to B executing this task must be larger than the respective speeddown.



performance



A performance model for a low energy cluster (dropping infrastructure)

- expected scaling penalty to power (switches)
- expected parallelisation penalty to performance (communication, numerics)

$$E = \left(n_{\text{nodes}} P_{\text{node}} + \left\lceil \frac{n_{\text{nodes}} - 1}{n_{\text{ports}}} \right\rceil \right) \left(\frac{t_{n_{\text{nodes}} = 1}}{n_{\text{nodes}}} + t_{\text{overhead}}(n_{\text{nodes}}) \right)$$

$$E_{A,B} = \frac{\Delta P_{A,B}}{\Delta t_{A,B}}$$

Hardware-oriented numerics and energy efficiency

How hardware and numerics determine energy to solution in (FEM-based) technical flow simulation





Hardware efficiency and performance engineering for technical flow simulation

- most of our codes are memory bandwidth-bound
- proper exploitation of SIMD is key to single core performance. Often: optimised SpMV.
- the memory interface is saturated with a (small) amount of cores
- GPGPU usually gives us a speedup of 5 10 through larger on-chip memory bandwidth. GPUs can also saturate that bandwidth.
- mixed precision provides another x1.5 max sustainable (double-single). More possible (double-half)
- low precision: with some methods
- baseline power of all devices has to be amortized via hybrid computation with careful load balancing



BSTREAM X 2

Arithmetic intensity

Numerical efficiency and performance engineering for technical flow simulation

- clever smoother construction example: **SPAI-types**
- in theory: SPAI with same structure as A gives convergence rates like GS (SPAI-1)
- works very well as MG smoother
- **construction phase:** different ways, we make progress (next)
- application phase: SpMV

$$\mathbf{x}^{k+1} \leftarrow \mathbf{x}^{k} + \omega M(\mathbf{b} - A\mathbf{x}^{k})$$

$$I - MA \parallel_{F}^{2} = \sum_{k=1}^{n} \parallel e_{k}^{\mathrm{T}} - m_{k}^{\mathrm{T}}A \parallel_{2}^{2} = \sum_{k=1}^{n} \parallel A^{\mathrm{T}}m_{k} - e_{k} \parallel_{2}^{2}$$

$$min_{m_{k}} \parallel A^{\mathrm{T}}m_{k} - e_{k} \parallel_{2}, \quad k = 1, \dots n.$$



Numerical efficiency: recent SPAI preconditioner results from EXA-DUNE project

• SPAI is exceptionally adaptable

- allows for good balancing of effort/energy to effectivity of preconditioner/smoother
- high reuse potential of once created approximate inverse
 - many screws to adapt to hardware (assembly stage)

 predefined sparsity pattern (SPAI-1)
 refinement of sparsity pattern
 refinement of coefficients
 rough inverses often good enough
 (half precision, use Machine Learning / Interpolation in knowledgebase)





Energy efficiency and performance engineering for technical flow simulation

Influences on incore performance of technical flow simulation (hardware efficiency):

FEM space(s)
 mesh adjacencies (fully unstructured)
 DOF numbering
 matrix storage (SELL)
 accuracy (mixed, low)
 assembly of matrices (SPAI methods)

Influences on incore performance of technical flow simulation (numerical efficiency):

assembly of matrices (SPAI- order)
 solver scheme
 preconditioners / smoothers

Building an insular compute center with a low energy cluster

And actually running it





Tegra K1 SoC

History of embedded/mobile processors is different than commodity archs'

32 Bit architecture, 4 x Cortex-A15 CPU, programmable Kepler GPU, LPDDR3, high SP Performance, most energy-efficient SoC of its time nowadays: Tegra X1, Tegra X2, ...

$\frac{\operatorname{Perf}_{\operatorname{peak}}^{\operatorname{TegraK1}}}{P_{\operatorname{peak}}^{\operatorname{TegraK1}}} =$	$=rac{35 \mathrm{GFlop}/s}{W}$
$\frac{\operatorname{Perf}_{\operatorname{peak}}^{\operatorname{Ref}}}{P_{\operatorname{peak}}^{\operatorname{Ref}}} =$	$\frac{11 \mathrm{GFlop}/s}{W}$
Ref [.] Xeo	n + K10

Jetson TK1 carrier board

Everything we need in a compute node on a single carrier board

Tegra K1 Carrier Board, P = 10 - 15W incl. fan (overkill), GiBit Ethernet, much I/O: serial, USB, SATA Linux, CUDA, nowadays: : Jetson TX1, Jetson TX2





Energy supply & -storage

Challenges: area, weather, operation at night

7.5 kWp Photovoltaik, freeland
2.5m x 16m,
8kWh Li-Ion battery,
2+1 inverter, charge
management
Aim: Full operation at daytime,
mild operation at night (3
seasons), just stay alive (winter)



Housing, cooling, heating

Challenges: isolation, area, ventilation, (cooling, heating)

Custom design modified high cube cargo container gomm isolation, fireproof, steel safety doors, heatable ventilation with separate power supply,





Cluster, networking, data storage

Challenges: space usage, avoid heat nests

60 x NVIDIA Jetson TK1, 240 ARM Cortex-A15 cores, 60 Kepler GPUs, 120 GB RAM (LPDDR3), 3 PDUs, sensors, 3 + 1 low power switches (Netgear) power dissipation (peak): ~ 1kW theoretical peak performance: ~ 20 TFlop/s mixed precision, 1 rack skeleton, many standoffs, cables, 3D-printed custom parts

Energy savings in the data storage system

Fully portable, BananaPi-based 10 x 1 TByte SSDs, redundant, e.g.: 5 TByte usable



MAX PLANCK INSTITUTE FOR DYNAMICS OF COMPLEX TECHNICAL SYSTEMS MAGDEBURG

Peak sustainable speed:

90 MByte/s (for comparison: 140 with commodity) -> **Speeddown = 1,6**

Average power: 30W (for comparison: 500W with commodity) -> Powerdown = 16,8

Benchmarks: basic kernels, single node



Benchmarks: basic kernels, single node



Applications on ARM Cortex A15





25

technische universität dortmund fakultät für **m**



Benchmarks: applications, full cluster, GPU, CPU



 \frown commodity(2015) \rightarrow commodity(2012 desktop) - Jetson TK1(2014) \rightarrow commodity(2012 compute)

Photovoltaic subsystem results

Uptime

ICARUS

Since spring 2016: Even at high stress (weather) no long downtimes in spring/summer autumn, mild slumber (16/60 cores online) in winter

Climate in the container On hot days: Max 33 °C ambient, Max 68 °C Jetson boards, Humidity 50% Aux power: 3 x 100W

Total cost 84k €

Devel time 3 years Battery temperature On cold days: hard to hold battery warm

May be reduced significantly without all the trial-and-error

Cost

Effort

Hardware market dynamics ...are fast! (next)

Lessons learned

Unconventional Preparedness of institutions?

HBSU 291532 Tegra KI small memory size, bandwidth, controller, (much better in later versions)

power

performance

Next?

- 'Mobile/embedded' is becoming more 'multi-purpose'.
- All compute hardware is becoming more energy-efficient.
- When will the employment of a newer architecture pay off?
- How much energy can we save during that time?
- Who will win the race? Will there be convergence?
- For which kind of codes does this pay of?

Next?

- autotuning for parameter settings (Exa-DUNE)
- better SPAI-eps and SAINV (Exa-DUNE)
- exploit Machine Learning more
- spread the lore:

Taking control of energy-consumption can make a huge difference

Thanks!

This work has been supported in part by the German Research Foundation (DFG) through the Priority Program 1648 'Software for Exascale Computing'.

ICARUS hardware is financed by MIWF NRW under the lead of MERCUR.

(Super-)Computers and power dissipation

(when we thought ICARUS in 2014)

Green500 rank	architecture	Top500 rank	Performance / P [GFlop/s / W]	P [MW]	 no Top500 topscorer (Top500#1:P=18)
1	Xeon + FirePro	168	5.2	0.057	Xeon CPU:
2	Xeon + PEZY	369	4.9	0.037	commodity
3	Xeon + K20x	392	4.4	0.035	GPUs (Fire,Kxx): also commodity
4	Xeon + K40m	361	4.0	0.044	 PEZY CPU: unconventional

www.top500.org/green500/list/2014/11

(Super-)Computers and power dissipation

(2016)

Green500 rank	architecture	Top500 rank	Performance / P [GFlop/s / W]	P [MW]	• Perf./P: x1.5!
1	Xeon + P100	28	9.5	0.350	Sunway: also
2	Xeon + P100	8	7.5	1.310	unconventional and Top500#1!
3	Xeon + PEZY	116	6.7	0.150	 Powerdown Top500#1: x0.8
4	Sunway	1	6.0	15.370	

www.top500.org/green500/list/2016/11

(Super-)Computers and power dissipation

(2017)

Green500 rank	architecture	Top500 rank	Performance / P [GFlop/s / W]	P [MW]	
1	Xeon + P100	61	14	0.14	Perf./P.: x1.5
2	Xeon + P100	465	14	0.03	 PEZY, Sunway systems: now
3	Xeon + P100	148	12	0.08	ranking 7, 15, 17
4	Xeon + P100	305	10	0.06	

www.top500.org/green500/list/2017/07